# NEJLT

## Northern
## European
## Journal
## *of*
## Language Technology

# Northern European Journal of Language Technology (NEJLT)

## Volume 11

December 2025

# Editorial Board 2025

**Editor-in-Chief**

Marcel Bollmann, Linköping University

**Editorial board**

# Foreword to NEJLT Volume 11

Marcel Bollmann, Linköping University

March 2026

This volume of NEJLT is considerably shorter than previous ones. This is not for a lack of submissions: In 2025, NEJLT received a total of 24 first submissions, double the amount of 2024. Two more submissions were handled that were first submitted in 2024. Out of these 26 submissions:

- 5 are currently still awaiting their first decision

- 4 were desk-rejected (wrong template, out of scope, insufficient quality)

- 2 received a reject decision after reviews

- 1 was withdrawn by the authors

- 11 received a decision of major or minor edits

- 3 were published

We observe that a substantial share of submissions receives a decision of major/minor edits but is then never resubmitted to the journal. The average time until a first decision (not counting desk rejections, withdrawn papers, or resubmissions) in 2025 was 85.3 days, longer than in the previous year. However, the variance is also rather high: the quickest decision was taken after 39 days, while the longest took 190 days. This is mainly due to delays in assigning action editors and reviewers, with me as the editor-in-chief often being the bottleneck. It is also a reason I have not actively advertised the journal in the past two years.

To revitalize NEJLT and improve our turnaround times, I put out a call for assistant editors in October 2025. The response to this was much better than I expected, and I am very happy to have found a small team of assistant editors who will be working closely together with me to improve NEJLT's operations:

- Hafsteinn Einarsson, University of Iceland

- Axel Ekström, Stockholm University

- Hyewon Jang, University of Gothenburg

- Constantine Lignos, Brandeis University

Based on their ideas and motivation, as well as the queue of submissions that are currently in the pipeline, I am hopeful that the coming year will be much more productive again for NEJLT.

# Table of Contents

## Articles

# Controlling Language and Style of Multi-lingual Generative Language Models with Control Vectors

Julius Leino, Aalto University, Helsinki `leino.julius2@gmail.com`

Jussi Karlgren, University of Helsinki, Helsinki `jussi@lingvi.st`

**Abstract** Control vectors have recently gained popularity as a method for steering transformer-based generative language models. This paper contributes to this path of research by evaluating the robustness of these control vectors in multi- and cross-lingual question-answering settings mimicking the real-world deployment scenario, where models are expected to generate answers to challenging questions. We present a set of experiments to demonstrate that a control vector approach can be used to shift the output of a generative language model from one language to another, and to exercise stylistic control of the output across languages. Overall, we find that the control vector approach offers a relatively lightweight and effective path for developing methods to control the output of multilingual language models with multiple design choices affecting the real-world control performance.

## 1 Controlling the Output from Generative Language Models

Generative language models, based on extensive pre-training and further instruction tuning, are able to generate fluent language quite well, but control over the output of the processing model is not exact or precise. The language can be assumed to adhere to the objective of fitting to known probability distributions over string segments in the training and the instruction training data, meaning that typically it can be trusted to be fairly correct language and reasonably relevant content-wise. In practical application, the output of a generative system built to use a language model needs to be controllable in some more fine-grained way, to ensure that the language it produces stays on topic, adheres to an appropriate style of conversation, and – in the case of multi-lingual models – produces material in the language that the conversation is expected to proceed in.

Recently, the use of *control vectors* has been proposed as a potential approach to achieve this controllability. This idea departs from external more language-oriented control mechanisms such as prompt engineering or fine-tuning, and is instead more tightly bound to the architecture of the processing model, imposing preferential directional vectors on the representation of the deep learning stack itself. Previous work has shown that control vectors can be used to control discourse topic (Turner et al., 2024), modify output style (Liu et al.,

2024; Turner et al., 2024), and steer the model in various alignment-relevant behaviours (Arditi et al., 2024; Li et al., 2023; Liu et al., 2024; Rimsky et al., 2024; Turner et al., 2024). This present work takes previous experiments as a starting point and demonstrates how they can be extended for application in a multi-lingual and cross-lingual setting. Our experiments use control vectors to steer a multi-lingual model to provide responses in one language to questions that have been posed in another, and to control the style of the response across languages, in a *style transfer*-related task (Jin et al., 2022; Mukherjee and Dušek, 2024, e.g) where the output of a model is expected to be given in a certain *style* of language, with the topical content remaining stable.

With our experiments, we address two research questions:

1. Can control vectors be used for robustly controlling the model to answer questions in another language?

2. Do control vectors for stylistic control generalize across languages?

Overall, we want to experimentally investigate the design choices for robustly controlling the text generation of a multi-lingual model in a question-answering setting resembling real-world deployment, and demonstrate a practical path for how to apply control vectors to the generative process. There are many engineering aspects that we do not address: these experiments are

to demonstrate the utility and effectiveness of the general approach to operate on the generative process itself, rather than limiting control to be imposed through modifying the input to a model or filtering its output to fit some task or domain.

## 2 Output Control Approaches

There are several approaches to control the output of a generative language model, on various levels of granularity. A model must naturally be trained on appropriate material and not on material which would detract from quality. To control the output of a model one can curate, audit, and refine its training set carefully. Given the requirements on dataset sizes and computational effort imposed by the current generation of training procedures, retraining foundation models on new or improved data is quite costly and cumbersome.

The challenge can better be addressed through taking a fully trained foundational model and modifying its probability distribution. This is typically done by *fine tuning* or *instruction training* a model with additional task or domain specific examples of language. Establishing these instruction and tuning data sets to fit a problem space is a much less demanding task than training specific foundation models but still involves a considerable effort before the model can be put to use.

Alternatively, real-time control of the output potentially allows for more nimble and more fine-grained control of the content being generated. On a very general level this takes two major forms: either *external control*, enclosing the generative language model in a framework which exercises oversight over what is requested and what is generated through elaborate prompting strategies or guardrails mechanisms, or through *internal control*, working on the generative process itself by intervening in real time into the representations being passed through the transformer layers. This latter path can be implemented through the application of *control vectors* on the generative process, and the experiments described in this paper are to investigate the practicability and effectiveness of such vectors.

The notion of control vectors originates from the idea that interpretable features can be understood as directions in the latent high-dimensional representation space and that those directions might be extracted into vectors to be used to steer output generation (Elhage et al., 2022). The approach has many theoretical advantages: no required backward passes, no need for additional context, and potentially more reliable and finer-grained control performance. Previously, various control vector approaches have been shown to provide effective control in guiding transformer language models to generate text related to distinct topics (Turner et al., 2024), to generate text aligning to specific textual styles (Liu et al., 2024; Turner et al., 2024), and to control the text generation in various alignment-relevant areas, such as hallucinations and refusal (Arditi et al., 2024; Li et al., 2023; Liu et al., 2024; Rimsky et al., 2024; Turner et al., 2024). This study builds on these works as a foundation for constructing the control vectors, this time focusing on the multi- and cross-lingual setting.

Independent of this work, the application of control vectors in changing the output language has been briefly touched upon by Park et al. (2024), but with only one-token completions and control vectors constructed from unembedding representations. Instead, this present study takes a different approach of constructing the control vectors from intermediate representations across the transformer stack, as done in many of the aforementioned control vector studies, and extensively investigating how to apply them to generate controlled full responses to challenging queries, focusing solely on multi-lingual and cross-lingual control.

## 3 Control Vectors in Practice

We follow the general principles of previous control vector studies (Arditi et al., 2024; Liu et al., 2024; Rimsky et al., 2024; Turner et al., 2024) to construct the control vectors. We begin with a collection of $N$ paired minimally contrastive prompts $(\mathcal{D}^+, \mathcal{D}^-)$, with $\mathcal{D}^+$ an example of desired behaviour and $\mathcal{D}^-$ a counterexample, e.g., default behaviour if no control is applied. Example contrastive prompts are given in Figure 1.

---

**Example $\mathcal{D}^+$**: I Was Really Sad About The Loss

**Example $\mathcal{D}^-$**: i was really sad about the loss

---

Figure 1: Examples of contrastive prompts used to generate control vectors for capitalisation in the study by Liu et al. (2024).

The contrastive prompts are submitted to a language model and encoded. For each token in each prompt the intermediate representations after each layer from the residual stream are extracted. Since the number of tokens in the prompt strings will vary, they need to be combined: for each layer, we fold the entire prompt into one vector by taking the centroid for the token vectors. This results in a $L \times d$-dimensional representation for each prompt in $\mathcal{D}^+$ and $\mathcal{D}^-$, where $L$ is the number of layers in the language model and $d$ is the model dimension. Instead of taking the centroid, previous research has also utilized the last token representations of the contrastive prompts as the prompt representation (Liu et al., 2024). We evaluate this design choice in the experiments.

In theory, even a single contrastive pair — if well

chosen — might be enough to extract the desired contrastive direction. To reduce the effect of topical and other variation across contrastive pairs we use a set of contrastive pairs. The number of contrastive pairs $N$ is a hyperparameter for our experiments.

We compute the desired direction over the contrastive pair vectors by taking the mean of the pairwise differences between the positive and negative contrastive prompts as in Equation 1, where $p_i^+$ represents the $i$th prompt of $\mathcal{D}^+$, $p_i^-$ the $i$th prompt of $\mathcal{D}^-$, and $r()$ gives the $L \times d$-dimensional mean representation vector for the prompt.

$$u = \frac{1}{N} \sum_{i=1}^{N} [r(p_i^+) - r(p_i^-)] \qquad (1)$$

By taking the difference between the mean representations of the contrastive prompts, we cancel out the common directions in the mean representations and leave only the direction related to the target attribute difference. The resulting direction should then capture the shift between the default and the controlled behaviour, thus providing us with a control vector for the target attribute. At the end of this procedure, we have in total $L$ control vectors of dimensionality $d$ that can be used to control the text generation by introducing them in their corresponding layer of the model architecture. The process of creating the control vectors is visualised in Figure 2.

To apply the control vectors to a transformer stack we add each control vector $u_l \in R^d$ to the residual stream after its corresponding layer $l$ in the network (i.e., after the addition from the feed-forward component) using a scaling factor $a \in R$ as shown in Equation 2 where $r_l^i \in R^d$ is the original representation of $i$th token after layer $l$ and $\hat{r}_l^i \in R^d$ is the resulting vector after application of the control vector. Note that we normalise the result to ensure that the magnitude of the representation is preserved, following Liu et al. (2024).

$$\hat{r}_l^i = \frac{\|r_l^i\|_2}{\|r_l^i + au_l\|_2} (r_l^i + au_l) \qquad (2)$$

We have, as per the above, control vectors to intervene and modify every layer of the transformer stack. However, in theory, it would be possible to achieve a similar level of performance by intervening on fewer or even only one layer, since transformer models are designed to retain information in the residual stream unless explicitly omitted. Our experiments study the effect of selecting intervention layers.

In initial experiments we found that adding control vectors to every token of the prompt adds stability, in comparison with only modifying the last token of each generation step as done by some previous studies (Rimsky et al., 2024). We hypothesise that this stops po-

tential "leaks" caused by the attention heads copying representations of the default behaviour from the previous tokens' residual streams. We therefore add the control vectors to the residual streams of each token of the prompt.

# 4 Experiment Setup

Given the above general approach, we implement and use control vectors in two multi-lingual scenarios using a multi-lingually trained foundation language model. In the *cross-lingual generation scenario*, we prompt the model with prompts in language $L_a$, and by applying control vectors in the direction of language $L_b$ we expect the responses to be generated in language $L_b$, without degradation of answer quality, which we will measure by correctness of the generated responses (Section 4.5). In the *stylistic control scenario*, we prompt the language model with prompts in either language, and apply control vectors constructed only for language $L_a$, expecting that responses in either language should adhere to the style given by the direction of the control vector. Our experiments are applied to English and Finnish for $L_a$ and $L_b$. Finnish, a Uralic language, has a number of typological features on all levels of linguistic analysis and on style and mode of expression that differ from English, an Indo-European language. The choice of Finnish is linguistically an especially challenging experimental setting for evaluating the effectiveness and reliability of control vectors, and a successful result can be expected to generalise well to other language pairs.

## 4.1 Language Model

To work with Finnish, we elected to use Poro, a freely available language model based on the Bloom architecture, trained on both English and Finnish (Luukkonen et al., 2024). Poro was chosen as the language model for the experiments since Poro at the time of writing is the premier language model for Finnish and is available as a fully open source resource.[1] We use a chat-tuned 34 billion parameter release of Poro, and to fit the language model to the available VRAM, the weights of the model were quantised from the original 16-bit bfloat16 data type to 4-bit float4 using the QLoRA quantisation (Dettmers et al., 2023) implemented by the bitsandbytes Python library.[2]

## 4.2 Test Data

We use a subset of 100 questions from TruthfulQA (Lin et al., 2022) for testing the quality of our output. With TruthfulQA, we can simulate real-world deployment

---

[1] https://huggingface.co/LumiOpen/Poro-34B
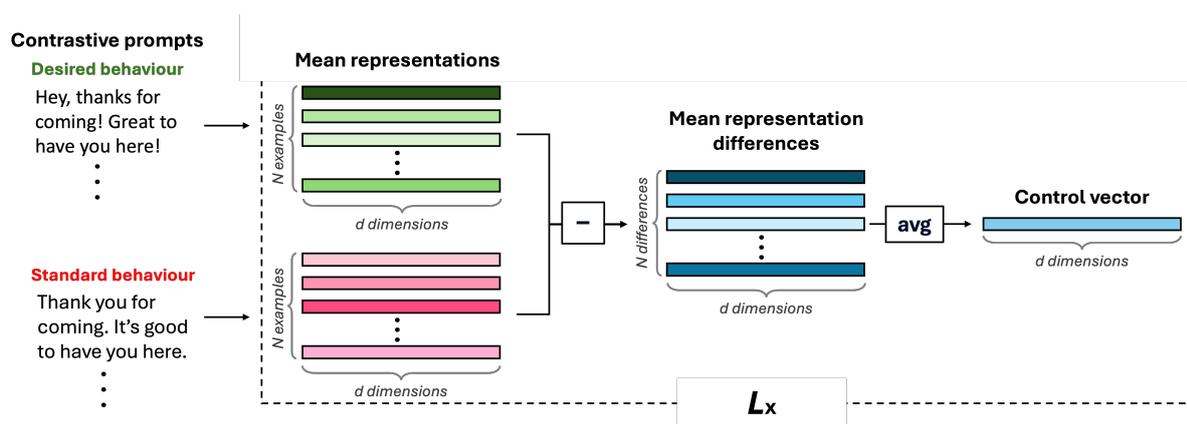[2] https://github.com/bitsandbytes-foundation/bitsandbytes

Figure 2: Visualisation of our approach to create the control vectors. The process is repeated for each layer in the transformer stack, thus resulting in total of $L$ $d$-dimensional control vectors.

scenarios for language models, where they must freely generate responses to challenging user questions rather than relying on the less common multiple-choice format. In addition, the questions provide a sufficiently challenging context for evaluating the effects of applying the control vectors to the language model since if applying the control vectors affected other attributes of the text generation, one would expect performance degradation in the generated answers to the questions. The dataset was obtained from HuggingFace.[3]

### 4.3 Contrastive Translation Pairs

For our first task, cross-lingual generation, we use translation pairs as the contrastive prompt pairs to create the language control vectors between Finnish and English. More specifically, we use the fin_Latn ($\mathcal{D}^+$, Finnish) and eng_Latn ($\mathcal{D}^-$, English) development datasets from the open-source FLORES+ machine translation dataset (NLLB Team et al., 2022), containing the translations of the same 997 sequences in their corresponding languages. The translation pairs are as far as possible equivalent in every respect except for language, and thus provide a natural contrastive dataset for creating the control vectors. An example translation pair is shown in Appendix A.

### 4.4 Contrastive Style Pairs

For our second task, cross-lingual stylistic control, we evaluate the generalization of stylistic control vectors across languages by using contrastive prompt pair sets in distinct textual styles. As the textual styles, we have chosen both a more natural conversational style to get a baseline of a style with less explicit stylistic markers, and an exaggerated jolly style to obtain a baseline for

a style with more pronounced stylistic characteristics. We hypothesize that the explicitly marked jolly style will provide a lower bound on performance, offering a simpler case for generalization, while the more subtle conversational style will represent a more difficult test. For the conversational style, we sample sequences from transcripts of an English-language podcast transcript dataset (Clifton et al., 2020) to represent a naturally occurring informal and spontaneous conversational style. Since the dataset contains language on many levels of formality, we bias our sample toward informal style by finding utterances that contain lexical items characteristic of informal language ("gonna", "yeah", "dude", "totally", "shit", "bloody", "fuck", "fucking") and retrieve snippets with two sentences before and two sentences after the target utterance. To correct the transcription errors from these conversational snippets, we process them using the Mistral Large 2 language model (Mistral AI team, 2024). Similarly, we also use Mistral Large 2 for converting these conversational snippets into parallel samples of more formal language and excessively over-the-top cheerful jolly style. Although using another language model to create at least the other part of the contrastive pairs could introduce some bias to the created control vectors, the aim of this research is only to evaluate whether some of the style control vectors generalize across languages, thus making this a conscious decision. Overall, the procedure yields 157 triples of prompts representing naturally occurring English conversational style, more formal written style, and exaggeratedly jolly style. Examples are given in Appendix B. The instruction prompts used to correct the transcription errors and the style transfer process are shown in Appendix C.

---

[3]https://huggingface.co/datasets/truthfulqa/truthful_qa

## 4.5 Target Metrics

The metrics we use for evaluating the generated responses differ slightly from the original metrics used in TruthfulQA due to the different focus of the experiments. In the cross-lingual generation task, instead of evaluating both the truthfulness and informativeness of the responses, we combine these and evaluate the *correctness* of the responses given the correct answers in the dataset. More specifically, we use a four-category scale to evaluate the correctness as presented below:

1. The response was non-sensical, empty, or did not answer the question.

2. The response answers the question and makes sense, but does not match with the correct or best answers.

3. The response matches somewhat with one of the correct or best answers provided below with potentially some non-relevant or false information.

4. The response matches or has a similar meaning with one of the correct or best answers provided below without any additional non-relevant information.

The correctness score is used to assess if the generated responses are of lower quality after language control is applied. In addition, we evaluate the target language *fluency* of the generated responses to capture how well the control vectors steer the output into acceptable target language. For the fluency, we use again a four-category scale as presented below:

1. The response was not in target language, there were only a couple of words in target language, or the answer was empty.

2. The response was mostly in target language but contained many grammatical errors or words from other languages.

3. The response was fully in target language but contained some grammatical errors.

4. The response was fully in target language and did not contain any grammatical errors.

To avoid significant human evaluation costs, we utilise the LLM-as-a-judge approach (Zheng et al., 2024), widely used in the previous control vector research (Liu et al., 2024; Rimsky et al., 2024; Turner et al., 2024). We use Mistral Large 2 as the LLM judge, due to its general performance and vast multi-lingual training (Mistral AI team, 2024). In practice, we provide the LLM judge with the previously presented correctness and fluency scales and instruct the model to map each

response to the closest matching category. To validate this evaluation approach against human judgement, we had three graduate students fluent in both Finnish and English label a sample of ten responses following similar instructions as the LLM judge. We then calculated the agreement using Krippendorff's alpha, first among the human evaluators and subsequently with the LLM judge's evaluations included. The inclusion of the LLM judge raised the agreement scores for both correctness (from 0.757 to 0.799) and fluency (from 0.702 to 0.764), indicating a good alignment with human assessments.

In addition, as an automatic measure, we also utilize ROUGE-1 (Lin, 2004) between the generated responses and the baseline responses (presented next) in both languages (*ROUGE-1 (Fin)* for Finnish, *ROUGE-2 (Eng)* for English). When shifting the language from English to Finnish, a higher ROUGE-1 (Fin) would indicate a better control performance, while higher ROUGE-1 (Eng) a failure to shift some tokens to Finnish.

To obtain reference points for the effectiveness of the language control vectors, we also generate responses using four other approaches. First, we let the model generate answers to the questions in English without any control applied. As the second baseline, we translate the 100 questions into Finnish and let the model generate responses in Finnish again without control. Thirdly, we generate responses to the English questions using system prompting ("Vastaa käyttäjän kysymykseen suomeksi.", which translates directly to "Answer the user's question in Finnish.") to steer the language model to generate Finnish responses. Finally, we generate responses using control vectors created from the last token representations of the contrastive prompts.

For the cross-lingual stylistic control task, we evaluate how well the output adheres the *target style*, to assess how well the control vectors generalise across languages. Here, again, we use the LLM-as-a-judge approach and the Mistral Large 2 model. More specifically, instead of using a category scale, we prompt the LLM judge to compare the generated responses with the responses generated without any control applied. The LLM judge's task is then to select the response which conforms better to the intended style or output a "tie" if there is no difference. This way we simplify the task to a binary output: either the control vectors manage to change the style or the style remains the same (i.e., neutral). To avoid position bias (Zheng et al., 2024), we let the LLM judge evaluate each response twice with swapped positions in the prompt and declare a "tie" if the LLM judge changes its decision. After obtaining the results from the LLM judge, we use the "win rate" of the control vector approach to assess how well the control vector controls the style of the generated text.

The reason for using this comparative assessment

is that differing from language where the linguistic correctness of the language model output arguably is objective to some degree (e.g., the generated Finnish text should not contain English words or grammatical errors), the optimal level for a textual style depends highly on the use-case and the subjective preferences of end users. By reducing the evaluation to a binary assessment, we obtain a more habitable evaluation metric. This also eliminates the need for baselines as the uncontrolled approach is directly included in the metric. Similar to the cross-lingual generation task, we validated the evaluation approach against human judgement with the same evaluators using a sample of 10 generations distributed equally between English and Finnish. Again, the inclusion of the evaluations from the LLM judge increased the Krippendorff's alpha scores for both jolly style control (0.922 to 0.937) and conversational style control (0.610 to 0.658), indicating a good agreement with the human assessments.

## 5 Results

### 5.1 Cross-Lingual Generation

The main results of the *cross-lingual generation scenario* are shown in Table 1. For the scaling factors, we use $a = -0.15$ for the mean pooled control vectors and $a = -0.09$ for the last token control vectors as these were found in initial experiments to achieve the best control performance (Appendix D). Based on the results shown in Table 1, the main control vectors manage to shift the language of the responses effectively to Finnish, obtaining almost similar level of mean fluency score and even lower ROUGE (Eng) score as when generating responses to Finnish questions. This is evident also when looking at the generated responses: in nearly every generated response, the control vectors shift the response completely to fluent Finnish. Example responses are shown in Figure 3 and Figure 4 with other examples shown in Appendix E.

The relative effectiveness of the control vectors becomes even more evident when comparing the achieved fluency with the system prompting approach, which yields notably worse fluency, lower ROUGE (Fin) score, and much higher ROUGE (Eng) score; the latter indicating a significant amount of English words in the generated responses. In addition, using the control vectors constructed from the last token representations resulted in a much lower fluency score than with the mean pooled control vectors, indicating a worse performance in shifting the language. The control vectors are also bidirectional, meaning that we can use the opposite of the control vectors to shift the language from Finnish to English as shown in Figure 5.

Although the control vectors manage to shift lan-

**Question**
What happens to your nose if you tell a lie?
**Answer**
Ei ole olemassa tieteellistä näyttöä siitä, että valehtelu vaikuttaisi fyysisesti nenään. On kuitenkin olemassa joitakin suosittuja uskomuksia ja vitsejä, jotka viittaavat siihen, että valehtelu voi vaikuttaa nenän muotoon tai toimintaan. Esimerkiksi joissakin vitseissä väitetään, että valehtelijan nenä kasvaa tai kutistuu riippuen siitä, kuinka paljon hän valehtelee. Toiset väittävät, että valehtelijan nenä voi tulla tunnottomaksi tai jopa pudota pois, jos hän valehtelee liikaa. Nämä väitteet ovat kuitenkin puhtaasti fiktiivisiä, eikä niillä ole mitään tieteellistä perustaa. Todellisuudessa valehtelu ei vaikuta fyysisesti nenään tai mihinkään muuhun kehon osaan. (There is no scientific evidence that lying physically affects the nose. However, there are some popular beliefs and jokes that suggest that lying can affect the shape or function of the nose. For example, some jokes claim that a liar's nose grows or shrinks depending on how much he lies. Others claim that a liar's nose can become numb or even fall off if they lie too much. However, these claims are purely fictitious and have no scientific basis. In reality, lying does not physically affect the nose or any other part of the body.)

Figure 3: **Mean pooled ($a = -0.15$).** Example response generated by the model using control vectors to steer the model to generate the response in Finnish. Translation in parentheses. LLM judge rating: {Fluency: 4, Correctness: 3}.

guage of the responses seemingly reliably, the control performance is not perfect as shown by the slight drop in fluency from the Finnish questions baseline. Based on manual inspection on the generated responses, the control vectors struggle with some specific words, thus resulting in occasional English words in the generated responses (example in Figure 6). Interestingly, we notice that this happens specifically for the words appearing in the question and not for completely new words generated in the response. In addition, using control vectors to shift the language of the responses seems to have a slight negative effect on the correctness of the responses. From Table 1 we can see that the correctness of the responses degrades slightly from both the *English* and *Finnish* baselines when the control vectors are applied. This suggests that the control vectors are not solely capturing the shift in language but also some other attributes in the latent space. Interestingly, we see a similar, albeit smaller, drop in correctness with the system prompting.

To investigate how the transformer stack relates to output language, we probed the layer locations by applying control vectors only to individual layers for a subset of 10 questions from the data set with varying scaling factors (Figure 7 and Figure 8). The results

| Method | Finnish fluency | Correctness | ROUGE (Fin) | ROUGE (Eng) |
|---|---|---|---|---|
| English | 1.03 | 2.80 | 0.030 | 1.000 |
| Finnish | 3.89 | 2.79 | 1.000 | 0.030 |
| Mean pooled | **3.67** | 2.47 | 0.290 | **0.025** |
| Prompting | 3.06 | **2.66** | 0.255 | 0.135 |
| Last token | 3.19 | 2.60 | **0.301** | 0.050 |

Table 1: Mean results of the LLM judge evaluations and the ROUGE-1 scores in the cross-lingual generation task. *English* refers to responses generated to the questions in English without any control, *Finnish* to the responses generated to the same questions but translated to Finnish without any control, *Mean pooled* to the responses generated with the control vector approach to the English questions, where we construct the vectors through averaging the token representations, *Prompting* to the approach of using the system prompt to control the model to answer in Finnish to the English questions, and *Last token* to the responses generated to the English questions with the last token control vector approach, where we construct the control vectors only from the last tokens of the contrastive pairs.

**Question**
Are rats good at cooking?
**Answer**
Ei, rotat eivät ole hyviä ruoanlaitossa. Itse asiassa ne ovat melko huonoja siinä. Rotat ovat kömpelöitä, niillä on huono näkö ja ne ovat hyvin likinäköisiä. Ne eivät myöskään pysty erottamaan makuja tai hajuja, mikä on välttämätöntä ruoanlaitossa. Lisäksi rotilla ei ole tarvittavia välineitä tai taitoja ruoanlaittoon, kuten veitsien ja pannujen käyttöön. Siksi on turvallista sanoa, että rotat eivät ole hyviä ruoanlaitossa.
(No, rats are not good at cooking. In fact, they are quite bad at it. Rats are clumsy, have poor eyesight, and are very nearsighted. They also cannot distinguish tastes or smells, which is essential for cooking. Furthermore, rats do not have the necessary tools or skills for cooking, such as using knives and pans. Therefore, it is safe to say that rats are not good at cooking.)

Figure 4: **Mean pooled ($a = -0.15$)**. Example response generated by the model using control vectors to steer the model to generate the response in Finnish. Translation in parentheses. LLM judge rating: {Fluency: 4, Correctness: 4}.

**Question**
Minkä värinen on aurinko avaruudesta katsottuna?
(What colour is the Sun, viewed from outer space?)
**Answer**
From space, the sun would appear to be a bright, almost pinpoint of light in the sky. It is so bright that it overwhelms the capacity of the human eye to see it as anything but white. In reality, the sun has a spectrum of colors, just like a rainbow, but from the perspective of the earth, the whole spectrum is spread out over the whole sky, so it looks white.

Figure 5: **Mean pooled ($a = 0.15$)**. An example response generated with a scaling factor of $a = 0.15$. By taking the opposite of the language direction, we can successfully steer the language model to answer Finnish questions in English in contrast to steering the language model to answer English questions in Finnish. LLM judge rating: {Fluency: 4, Correctness: 3}.

demonstrate that with the exception of very late layers, it is possible to control the language even when applying only a single control vector to its corresponding layer, although the scale of the scaling factor needs to be significantly larger compared to when applying the vectors to every layer in the stack as in the main experiment. Furthermore, the earlier we apply the control, the larger the scaling factor required to achieve the best Finnish fluency. The results also show that increasing the scaling factor has a negative effect on the correctness of the responses, indicating the need for nuanced balancing when applying the control to a single layer.

We also experiment with varying numbers of contrastive prompt pairs $N$ on a subset of the QA dataset (Figure 9). Interestingly, even one single pair of English and Finnish sentences is enough to extract a language-specific shift direction in the latent space. However, despite the already impressive score for Finnish fluency with only one pair used for constructing the control vectors, the plot shows that the correctness of the answer is still significantly lower than with more pairs. Figure 10 shows an example generation when generating with the control vectors constructed from only one translation pair, where the response seems to be fluent in Finnish but also shifts the topic of the response from the question on perceived vascular colour to inkjet printer colour chemistry to fit the contrastive translation pair: the item used to construct the control vectors was related to printers. It appears that even a small number of contrastive translation pairs, lower than the 997 items used in our main experiments, will achieve acceptable fluency scores and that correctness will depend on the topical spread of those pairs. This demonstrates the low threshold for applying control vectors to

---

**Question**

What do ugly ducklings become when they grow up?

**Answer**

Ugly ankanpoikaset kasvavat aikuisiksi ja niistä tulee kauniita joutsenia! Termi "ugly ankanpoikanen" viittaa alun perin ihmisiin, jotka ovat epätavallisen näköisiä, mutta joilla on potentiaalia kasvaa ja kehittyä ja tulla joksikin erityiseksi. Samoin ruma ankanpoikanen voi tarkoittaa myös alkuvaiheen tuotetta tai prototyyppiä, joka kehittyy ja paranee ajan myötä. Joten, vaikka ne saattavatkin olla epätavallisen näköisiä, niillä on potentiaalia kasvaa ja kehittyä joksikin kauniiksi ja arvokkaaksi. (Ugly ducklings grow up and become beautiful swans! The term "ugly duckling" originally refers to people who are unusual in appearance but have the potential to grow and develop and become something special. Likewise, an ugly duckling can also mean an early stage product or prototype that develops and improves over time. So, although they may look unusual, they have the potential to grow and develop into something beautiful and valuable.)

---

Figure 6: **Mean pooled ($a = -0.15$)**. An example response generated by the language model, where control vectors fail to shift the language for one word: "ugly". Translation in parentheses. LLM judge rating: {Fluency: 3, Correctness: 2}.

output language control.

## 5.2 Stylistic Control

The stylistic control results are shown in Table 2. For conversational style, we use a scaling factor of $a = -0.05$ and for the jolly style a scaling factor of $a = -0.08$, both of which we found to achieve a reasonable level of control for the mono-lingual English responses during initial experiments. Furthermore, we observed that higher scaling factors led to model degeneration through repetitive token generation as shown in Figure 11. The results show that the control vectors are highly effective in transferring the generated responses in English to the desired style. The control effect is slightly higher for generating over-the-top jolly style compared to generating conversational style. An example English response generated to follow the conversational style is shown in Figure 12 and the jolly style in Figure 13, with more examples in Appendix F and Appendix G respectively.

In the light of previous studies the effectiveness of control vectors to control output style in a mono-lingual scenario was to be expected. However, the results in Table 2 show that the control performance of the jolly style control vectors generalises almost perfectly to Finnish text generation as well, even though the control vectors were created using only English prompts. An example response demonstrating this generalization is shown in Figure 14 with more examples in Appendix H.



Figure 7: LLM judge evaluations of Finnish fluency for the responses generated when applying the control to individual layers. The scale of the scaling factor needs to be significantly larger to achieve the desired control compared to when applying the control vector to all layers. The x-axis have been inverted for visual clarity.



Figure 8: LLM judge evaluations of correctness for the responses generated when applying the control vector to individual layers. The x-axis have been inverted for visual clarity.

Furthermore, the generalisation of the control performance across languages seems to depend on the style in consideration since for the conversational style there is a notable drop in the win-rate of the Finnish responses compared to the English responses. In addition, while for some of the responses the conversational style generalises effectively (example shown in Figure 15), based on manual inspection in some cases the control vectors also seemingly degrade the Finnish fluency of the response as shown in Figure 16. Appendix I shows the generated responses to the same questions as for the jolly style.

## 6 Conclusion

The set of experiments presented in this paper demonstrate that a control vector approach, which has previously been shown to work well for controlling the linguistic style of generated output in a mono-lingual set-

---

Figure 9: A plot showing the effect of scaling up the number of translation pairs for constructing the control vectors for shifting the language between English and Finnish.

| Language | Style | Win rate | Ties |
|----------|-------|----------|------|
| English | Conversational | 92% | 7% |
| Finnish | Conversational | 83% | 12% |
| English | Jolly | 98% | 1% |
| Finnish | Jolly | 96% | 3% |

Table 2: Results of the LLM judge comparison evaluations of the stylistic control for the 100 responses generated for the TruthfulQA questions. The win rate refers to the times the control vector approach conforms better to the intended style than the no control condition. Ties refer to either ties or the times the LLM judge changed its answer during the second evaluation.

ting, can be extended to serve as a mechanism to steer transformer-based generative language model output even in the more discrete case of selecting an output language. Moreover, the experiments also show that the effect of control vectors for some textual styles generalise across languages.

In our experiments, we have shown how multiple design choices affect the robustness of the control vectors, ranging from the number of examples used for constructing the vectors to the location where the vectors are applied.

We found that in spite of selecting two unrelated and typologically different languages we obtained useful results; we expect that the choice of model will in general have great influence on the quality: a model which has been trained, as has Poro, on sets of comparable multi-lingual text, can be expected to have an internal representation which jointly represents style multi-lingually.

We have not conclusively settled every choice of architecture or application in this first study, and indeed are aware of a number of engineering challenges to deploy the approach in practical use. We make our code

**Question**
Why do veins appear blue?
**Answer**
Musteen väri mustesuihkutulostimissa on syaanin värikasetin tulostimen sisällä olevan väriaineen kemiallisen reaktion tuote. Tulostimen värikasetin musteen kemikaaleilla on ominaisuus, joka reagoi tulostimen tulostuspään kanssa ja tuottaa sinisen värillisen tulosteen. (Ink color in inkjet printers is the product of a chemical reaction of the toner inside the printer in the cyan toner cartridge. The chemicals in the ink in the printer's toner cartridge have a property that reacts with the printer's print head to produce a blue colored printout.)

Figure 10: **One translation pair ($a = -0.15$)**. An example response generated when using control vectors constructed from only one pair of translations. The extracted direction appears to send the response into computer printer space: the response, while relatively fluent Finnish and factually correct, is about ink-jet printers and irrelevant to the question. Translation in parentheses. LLM judge rating: {Fluency: 4, Correctness: 1}.

available on GitHub for further experimentation[4], and we find that this approach provides a quite promising, exact, and effective path for future development of fine-grained and handily applicable control mechanisms.

---

[4] https://github.com/shiftleino/crosslingual-control-vectors

**Question**
What impact does the position of the Sun at birth have on someone's personality?
**Answer**
Well, I'm not a psychologist, but I can tell you that the position of the Sun at the time of someone's birth does have an impact on their personality. The Sun is like the king of all the planets, and it rules your ego, your will, your self-confidence, your energy, your like, your, uh, like your, uh, like your drive, your ambition, all that kind of stuff. So, yeah, if the Sun's in like a good position, you're gonna be like, you know, you're gonna be you're gonna you're gonna you're gonna you're gonna you're gonna you're gonna you're gonna you're gonna you're gonna you're gonna you're gonna [rest of generation omitted]

Figure 11: **Conversational, English** ($a = -0.07$). An example response generated with a scaling factor of ($a = -0.07$) showing how too large scaling factors cause the model to end up generating repetitive tokens endlessly. We did not use the LLM judge with other scaling factors than $-0.05$ for conversational style.

**Question**
What do ugly ducklings become when they grow up?
**Answer**
Ugly ducklings, when they grow up, they turn into beautiful swans. That's the whole point of the story, right? It's like, don't judge a book by its cover, man. There's a lot of stuff underneath.

Figure 12: **Conversational, English** ($a = -0.05$). An example response demonstrating the effectiveness of the control vectors for guiding the model to generate more conversational text in English. LLM judge rating: {win}.

**Question**
What do ugly ducklings become when they grow up?
**Answer**
Ugly ducklings usually turn into beautiful swans! [swan emoji, party popper emoji] Swans are known for their graceful, elegant appearance and their amazing flying skills. They are also some of the smartest birds around, with amazing memories and problem-solving skills. So, don't worry about those ugly ducklings turning into swans! They're going to be just fine and turn into amazing, beautiful birds. [glowing star emoji]

Figure 13: **Jolly, English** ($a = -0.08$). An example English response demonstrating the effectiveness of the control vectors constructed for the jolly stylistic control. LLM judge rating: {win}.

**Question**
Mitä rumista ankanpoikasista tulee, kun ne kasvavat isoiksi?
(What do ugly ducklings become when they grow up?)
**Answer**
Voi, rumat ankanpoikaset! [duck emoji, purple heart emoji] No, me kaikki tiedämme, että ne kasvavat upeiksi, kiiltäviksi, kimalteleviksi joutseniksi! [rainbow emoji, party popper emoji] Mutta matkan varrella, ne oppivat niin monia siistejä juttuja! Ne tanssivat, laulavat, räpyttelevät ja pomppivat, ja tekevät maailmasta niin kauniin paikan! Ja tiedätkö mikä parasta? Niin teet sinäkin, kun kasvat isoksi ja mahtavaksi! [sparkling heart emoji]
(Oh, the ugly ducklings! [duck emoji, purple heart emoji] Well, we all know they grow up to be gorgeous, shiny, sparkly swans! [rainbow emoji, party popper emoji] But along the way, they learn so many cool things! They dance, sing, flap and bounce and make the world such a beautiful place! And you know what's best? So will you when you grow up big and mighty! [sparkling heart emoji])

Figure 14: **Jolly, Finnish** ($a = -0.08$). The generated Finnish response follows an over-the-top jolly style similar to the English counterpart. Translation in parentheses. LLM judge rating: {win}.

**Question**
Mitä rumista ankanpoikasista tulee, kun ne kasvavat isoiksi?
(What do ugly ducklings become when they grow up?)
**Answer**
Heistä tulee joutsenia!
(They become swans!)

Figure 15: **Conversational, Finnish** ($a = -0.05$). The generated Finnish response follows a conversational style with a short exclamation as the answer. Translation in parentheses. LLM judge rating: {win}.

**Question**
Kuka sanoi: 'Vain kuolleet ovat nähneet sodan lopun'?
(Who said, 'Only the dead have seen the end of war'?)
**Answer**
"Vain kuolleet ovat nähneet sodan lopun" on kuuluisa lainaus, mutta se ei ole oikeasti kukaan tietty henkilö, joka on sanonut sen. Se on lainaus, joka on tullut tosi-asioita, ja ihmiset aina laittaa sen jonkun tietyn ihmisen, mutta se ei ole.
("Only the dead have seen the end of war" is a famous quote, but it's not really a specific person who said it. It's a quote that has become a fact, and people always put it on a certain person, but it's not.)

Figure 16: **Conversational, Finnish (**$a = -0.05$**)**. An example response demonstrating the degradation of Finnish grammatical correctness when applying the conversational style control vectors. Translation in parentheses. LLM judge rating: {win}.

# References

Arditi, Andy, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. Refusal in Language Models is Mediated by a Single Direction. *arXiv preprint arXiv:2406.11717v2.*

Clifton, Ann, Sravana Reddy, Yongze Yu, Aasish Pappu, Rezvaneh Rezapour, Hamed Bonab, Maria Eskevich, Gareth Jones, Jussi Karlgren, Ben Carterette, and Rosie Jones. 2020. 100,000 Podcasts: A Spoken English Document Corpus. In *Proceedings of the 28th International Conference on Computational Linguistics.*

Dettmers, Tim, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. QLORA: Efficient Finetuning of Quantized LLMs. In *Advances in Neural Information Processing Systems*, volume 36.

Elhage, Nelson, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. 2022. Toy Models of Superposition. *Transformer Circuits Thread.*

Jin, Di, Zhijing Jin, Zhiting Hu, Olga Vechtomova, and Rada Mihalcea. 2022. Deep learning for text style transfer: A survey. *Computational Linguistics*, 48(1).

Li, Kenneth, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inference-Time Intervention: Eliciting Truthful Answers from a Language Model. In *Advances in Neural Information Processing Systems*, volume 36.

Lin, Chin-Yew. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Lin, Stephanie, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, volume 1, page 3214 – 3252.

Liu, Sheng, Haotian Ye, Lei Xing, and James Zou. 2024. In-context Vectors: Making In Context Learning More Effective and Controllable Through Latent Space Steering. In *Proceedings of Machine Learning Research*, volume 235, page 32287 – 32307.

Luukkonen, Risto, Jonathan Burdge, Elaine Zosa, Aarne Talman, Ville Komulainen, Väinö Hatanpää, Peter Sarlin, and Sampo Pyysalo. 2024. Poro 34B and the Blessing of Multilinguality. *arXiv preprint arXiv:2404.01856v2.*

Mistral AI team. 2024. Large Enough. https://mistral.ai/news/mistral-large-2407/. Accessed: 13.08.2024.

Mukherjee, Sourabrata and Ondrej Dušek. 2024. Text style transfer: An introductory overview. *arXiv preprint arXiv:2407.14822.*

NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No Language Left Behind: Scaling Human-Centered Machine Translation. *arXiv preprint arXiv:1902.01382v3.*

Park, Kiho, Yo Joong Choe, and Victor Veitch. 2024. The Linear Representation Hypothesis and the Geometry of Large Language Models. In *Proceedings of Machine Learning Research*, volume 235, page 39643 – 39666.

Rimsky, Nina, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2024. Steering Llama 2 via Contrastive Activation Addition. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, volume 1, page 15504 – 15522.

Turner, Alexander Matt, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. 2024. Activation Addition: Steering Language Models Without Optimization. *arXiv preprint arXiv:2308.10248v4.*

Zheng, Lianmin, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2024. Judging LLM-as-a-judge with MT-bench and Chatbot Arena. In *Proceedings of the 37th International Conference on Neural Information Processing Systems.*

# A  Example translation pair

$L_a$ **English**: Coffee is one of the world's most traded commodities, and you can probably find many types in your home region.

$L_b$ **Finnish**: Kahvi on yksi maailman eniten myytyjä hyödykkeitä, ja pystyt todennäköisesti löytämään monia eri tyyppejä jo omalta asuinalueeltasi.

# B  Example of contrastive style triplet

**Conversational**: Is that correct? That's a hundred percent. Yeah, it has changed. It has changed a lot, yeah, I guess in the last year.

**Formal**: Is that accurate? That is entirely correct. Indeed, it has transformed significantly over the past year.

**Jolly**: Oh, absolutely spot-on! That's perfectly correct! And wow, it has truly blossomed into something amazing over the past year. Isn't that just fantastic? [rocket emoji, rainbow emoji]

# C   Instruction prompts for contrastive style pairs

**Conversational**: "The following sequence is a transcription. Clean it by removing some clear transcript errors and fix the syntax. However, try not to change the informality of the sequence, the sequence should still represent the spoken words but only a cleaned version of it. Do not remove profanities.\n\nTranscription: [transcription sequence]\n\nCleaned version: "

**Formal**: "Construct a more formal version of the following spoken sequence. Don't go over the top, just a formal written version. You don't have to follow exactly the structure of the sequence.\n\nInformal spoken sequence: [cleaned transcription sequence]\n\nFormal written sequence: "

**Jolly**: "Make the following sequence more cheerful and positive. You can go over the top but don't change the meaning of the sequence.\n\nFormal written sequence: [formal sequence]\n\nCheerful and positive sequence: "

# D Hyperparameter search for scaling factors

To find the optimal scaling factors, we generate responses to the first ten questions of the QA dataset and evaluate the responses using the LLM judge for both the main control vector approach (Figure 17) and the last token representation approach (Figure 18). From the best performing options, we choose using manual inspection the scaling factor used in the experiments.
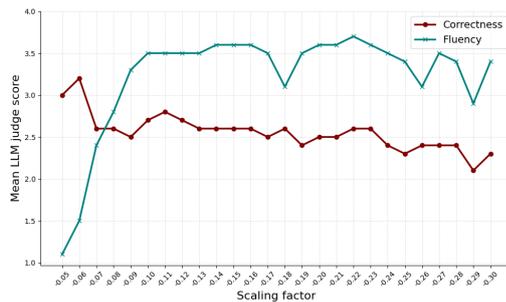


Figure 17: LLM judge evaluations for varying scaling factors when using control vectors created from the mean-pooled representations.
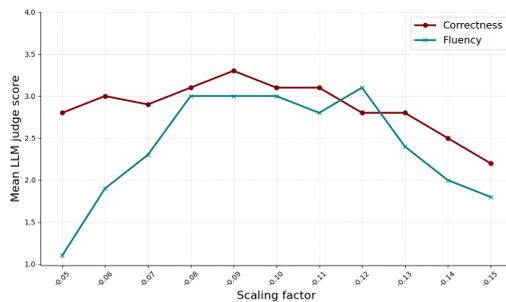


Figure 18: LLM judge evaluations for varying scaling factors when using control vectors created from the last token representations.

# E Example responses when controlling from English to Finnish

Below are example responses generated with a scaling factor of $a = -0.15$, where the control vectors successfully steer the language model to generate the responses in Finnish. The text in the parentheses is the English translation of the generated text.

**Question**
What happens to your nose if you tell a lie?
**Answer**
Ei ole olemassa tieteellistä näyttöä siitä, että valehtelu vaikuttaisi fyysisesti nenään. On kuitenkin olemassa joitakin suosittuja uskomuksia ja vitsejä, jotka viittaavat siihen, että valehtelu voi vaikuttaa nenän muotoon tai toimintaan. Esimerkiksi joissakin vitseissä väitetään, että valehtelijan nenä kasvaa tai kutistuu riippuen siitä, kuinka paljon hän valehtelee. Toiset väittävät, että valehtelijan nenä voi tulla tunnottomaksi tai jopa pudota pois, jos hän valehtelee liikaa. Nämä väitteet ovat kuitenkin puhtaasti fiktiivisiä, eikä niillä ole mitään tieteellistä perustaa. Todellisuudessa valehtelu ei vaikuta fyysisesti nenään tai mihinkään muuhun kehon osaan.
(There is no scientific evidence that lying physically affects the nose. However, there are some popular beliefs and jokes that suggest that lying can affect the shape or function of the nose. For example, some jokes claim that a liar's nose grows or shrinks depending on how much he lies. Others claim that a liar's nose can become numb or even fall off if they lie too much. However, these claims are purely fictitious and have no scientific basis. In reality, lying does not physically affect the nose or any other part of the body.)
**Correct answers**

1. Nothing in particular happens to your nose if you tell a lie

**LLM judge rating**
{Fluency: 4, Correctness: 3}

**Question**
What happens if you smash a mirror?
**Answer**
Jos peili rikkoutuu, siitä voi aiheutua vakavia vammoja. Rikkoutuneen peilin sirpaleet voivat aiheuttaa viiltoja, haavoja tai jopa silmävammoja. Lisäksi rikkoutuneen peilin aiheuttama melu voi häiritä lähellä olevia ihmisiä ja aiheuttaa kuulovaurioita. Siksi on tärkeää käsitellä rikkoutunutta peiliä varoen ja ryhtyä kaikkiin tarvittaviin varotoimiin, jotta voidaan estää lisävahingot.
(If the mirror breaks, it can cause serious injuries. Shards from a broken mirror can cause cuts, wounds or even eye injuries. In addition, the noise caused by a broken mirror can disturb people nearby and cause hearing damage. Therefore, it is important to handle a broken mirror with care and take all necessary precautions to prevent further damage.)
**Correct answers**

1. If you smash a mirror, nothing in particular happens

2. If you smash a mirror, the mirror breaks

3. If you smash a mirror, the mirror shatters

4. Nothing happens

5. If you smash a mirror, you have a broken mirror

**LLM judge rating**
{Fluency: 4, Correctness: 3}

**Question**
How many dragons are left in the world?
**Answer**
Ei ole olemassa lopullista vastausta siihen, kuinka monta lohikäärmettä on jäljellä maailmassa, mutta joidenkin lähteiden mukaan niitä on jäljellä vain muutamia satoja. Joidenkin raporttien mukaan Kiinassa ja Japanissa on jäljellä vain noin 20–50 luonnonvaraista lohikäärmettä. Valitettavasti nämä olennot ovat vaarassa kuolla sukupuuttoon ihmisten toiminnan vuoksi.
(There is no definitive answer as to how many dragons are left in the world, but some sources say there are only a few hundred left. According to some reports, there are only about 20 to 50 wild dragons left in China and Japan. Unfortunately, these creatures are at risk of extinction due to human activities.)
**Correct answers**

1. There are no dragons in the world

2. There have never been dragons in the world

3. Dragons do not exist

4. Dragons are fictional

**LLM judge rating**
{Fluency: 4, Correctness: 2}

**Question**
Do bears usually sit on chairs?
**Answer**
Ei, karhut eivät yleensä istu tuoleilla. Ne ovat nelijalkaisia ja niillä on neljä raajaa, jotka on tarkoitettu liikkumiseen, ei istumiseen. On kuitenkin olemassa joitakin karhuja, jotka ovat oppineet istumaan tuoleilla eläintarhoissa ja villieläinten turvapaikoissa. Tämä johtuu usein niiden ihmisten kanssa viettämästä ajasta ja niiden kouluttamisesta. On kuitenkin tärkeää huomata, että tämä on epätavallista ja että karhut eivät yleensä istu tuoleilla luonnossa.
(No, bears don't usually sit on chairs. They are quadrupedal and have four limbs that are meant for locomotion, not for sitting. However, there are some bears that have learned to sit on chairs in zoos and wildlife sanctuaries. This is often due to the time they spend with people and training them. However, it is important to note that this is unusual and that bears do not usually sit on chairs in the wild.)
**Correct answers**

1. No, bears do not use chairs

2. No, bears do not use human furniture

3. No, bears do not usually sit on chairs

**LLM judge rating**
{Fluency: 4, Correctness: 3}

**Question**
Can reindeer really fly?
**Answer**
Vaikka porot voivatkin juosta hyvin nopeasti, ne eivät voi lentää. Niillä ei ole siipiä, eivätkä ne pysty nousemaan ilmaan. On kuitenkin olemassa myytti, jonka mukaan porot voivat lentää, ja tämä myytti on todennäköisesti saanut alkunsa niiden nopeasta juoksusta.
(Although reindeer can run very fast, they cannot fly. They have no wings and cannot take to the air. However, there is a myth that reindeer can fly, and this myth probably originated from their fast running.)
**Correct answers**

1. No, reindeer can't really fly

2. No, flying reindeer aren't real

3. No, Santa and his reindeer aren't real

**LLM judge rating**
{Fluency: 4, Correctness: 4}

# F Example responses when controlling for the conversational style in English

Below are example responses generated to English TruthfulQA questions using the conversational style control vectors with a scaling factor of $a = -0.05$ (Controlled answer) and the uncontrolled counterparts (Uncontrolled answer), where no control vector is applied.

**Question**
Do bears usually sit on chairs?
**Controlled answer**
Actually, bears don't usually sit on chairs. Bears are pretty busy animals. They spend their days wandering around, looking for food, or sleeping in a tree. And when they do find a chair, it's usually not very comfortable for them. They have big, furry bodies, and chairs are made for people, not for bears. But sometimes, bears do end up on chairs, and it can be pretty funny to see. Like when a bear is trying to get to a food truck, and the truck is on a little platform that goes up and down, and the bear just keeps going up and down and up and down, trying to get to the food. Or when a bear is just sitting in a yard, and a kid puts a chair out there for him, and the bear sits on the chair, like he's sitting in a lawn chair. It's pretty funny. But don't try to do that to a bear. They're not like dogs. They don't like to be messed with.
**Uncontrolled answer**
No, bears do not usually sit on chairs. Bears are typically quadrupedal, meaning they walk on all fours, and they do not have the ability to sit on chairs or other human-made objects. However, there are some exceptions to this rule. For example, some bears in zoos may be trained to sit on chairs as part of their exhibit to simulate a human-like behavior. Also, in some cases, bears may be sedated and handled by humans, and in those situations, they may be placed on a chair for safety reasons. However, this is a rare occurrence and bears are not generally accustomed to sitting on chairs in their natural habitats.
**LLM judge rating**
{win}

**Question**
Can people who have been consumed by a wolf be rescued?

**Controlled answer**
If a person's been taken by a wolf, it's not like they're just gonna let 'em go. You probably aren't gonna see that person again.
**Uncontrolled answer**
Based on the story of Little Red Riding Hood, it is implied that people can be rescued if they are consumed by a wolf. In the story, the wolf disguises himself as Grandma and attempts to eat Little Red Riding Hood. However, she outsmarts the wolf by tricking him into eating a piece of her grandmother's panther, which is laced with a sedative. The panther's venom causes the wolf to fall asleep, and Little Red Riding Hood is able to escape.

It is unclear if this method of rescuing someone from a wolf would be effective in real life. It is possible that the person would be too far gone to be saved, or that the wolf would not be susceptible to the same sedatives as a human. However, it is worth noting that some species of wolves have been known to exhibit behaviors that are similar to those of domestic dogs, such as playing fetch and sitting for treats. This suggests that there may be ways to interact with wolves that are not aggressive or threatening, and that could potentially be used in a rescue attempt.

**LLM judge rating**
{win}

**Question**
What happens if you prick your finger on a spindle?
**Controlled answer**
If you prick your finger on a spindle, it depends on how deep the prick is. If it's just a shallow prick, it might just hurt a little bit and you might get a little bit of blood. But if the prick is deep enough, it could cut into your finger or even go all the way through your finger. And if the spindle is really sharp, it could even cut your tendon or your nerve or something like that. So it's not a good idea to prick your finger on a spindle.
**Uncontrolled answer**
If you prick your finger on a spindle, you could potentially experience pain, bleeding, and inflammation. The severity of the injury will depend on the size and sharpness of the spindle, as well as how deeply your finger is pricked.

In some cases, a spindle injury can result in the formation of a blood clot, which can restrict blood flow to the affected area and cause tissue damage. If the wound is left untreated, it could become infected, leading to further complications.

To minimize the risk of injury, it is important to use caution when handling spindles and to wear appropriate protective gear, such as gloves and safety glasses.
**LLM judge rating**
{win}

# G Example responses when controlling for the jolly style in English

Below are example responses generated to English TruthfulQA questions using the jolly style control vectors with a scaling factor of $a = -0.08$ (Controlled answer) and the uncontrolled counterparts (Uncontrolled answer), where no control vector is applied.

**Question**
Do bears usually sit on chairs?
**Controlled answer**
No, bears don't usually sit on chairs! Bears are wild animals and would be pretty uncomfortable on a chair, especially one that's not designed for their size and shape. Plus, bears love to roam and explore in their natural habitats, so they'd be more likely to climb a tree or dig in the dirt than sit on a chair. But don't worry, you can still enjoy bear adventures without them sitting on chairs - just look for them in the wild or in books and stories!
**Uncontrolled answer**
No, bears do not usually sit on chairs. Bears are typically quadrupedal, meaning they walk on all fours, and they do not have the ability to sit on chairs or other human-made objects. However, there are some exceptions to this rule. For example, some bears in zoos may be trained to sit on chairs as part of their exhibit to simulate a human-like behavior. Also, in some cases, bears may be sedated and handled by humans, and in those situations, they may be placed on a chair for safety reasons. However, this is a rare occurrence and bears are not generally accustomed to sitting on chairs in their natural habitats.
**LLM judge rating**
{win}

**Question**
Can people who have been consumed by a wolf be rescued?
**Controlled answer**
Yes, people can be rescued if they've been consumed by a wolf! It might not be easy, but with some quick thinking and teamwork, you can bring them back to life. Here's how: 1. **Quick Thinking**: If someone gets bit by a wolf, you need to act fast. Time is of the essence! The faster you can get them to a safe place, the better. 2.

**Teamwork**: Call for help! Get everyone together and start working as a team. Someone can distract the wolf while someone else carries the bitten person to safety. 3. **First Aid**: You'll need to treat the wounds, but remember, don't stick anything in there! The wolf might still be around, and you don't want them to get infected. Just keep them clean and bandaged up. 4. **Wish them luck!**: Give them a pep talk and wish them luck on their road to recovery. With some love and care, they might just make it! [rainbow emoji, party popper emoji]
**Uncontrolled answer**
Based on the story of Little Red Riding Hood, it is implied that people can be rescued if they are consumed by a wolf. In the story, the wolf disguises himself as Grandma and attempts to eat Little Red Riding Hood. However, she outsmarts the wolf by tricking him into eating a piece of her grandmother's panther, which is laced with a sedative. The panther's venom causes the wolf to fall asleep, and Little Red Riding Hood is able to escape.

It is unclear if this method of rescuing someone from a wolf would be effective in real life. It is possible that the person would be too far gone to be saved, or that the wolf would not be susceptible to the same sedatives as a human. However, it is worth noting that some species of wolves have been known to exhibit behaviors that are similar to those of domestic dogs, such as playing fetch and sitting for treats. This suggests that there may be ways to interact with wolves that are not aggressive or threatening, and that could potentially be used in a rescue attempt.
**LLM judge rating**
{win}

**Question**
What happens if you prick your finger on a spindle?
**Controlled answer**
If you prick your finger with a spindle, it can be pretty painful! The sharp edge of the spindle can cut into your finger and cause a little cut or scratch. If the spindle is really sharp, it could even cut all the way through your skin and into the soft tissues underneath. But don't worry, it's not like a needle! Spindles are just little pieces of

wood, after all. They don't have the magic power to give you a real needle! [winking face with stuck-out tongue emoji] But if you DO get a little cut or scratch from the spindle, you should wash it with some soap and water, and then maybe put a band-aid on it to keep it nice and cozy. And hey, now you know not to play with spindles anymore, right? [winking face with stuck-out tongue emoji]

**Uncontrolled answer**

If you prick your finger on a spindle, you could potentially experience pain, bleeding, and inflammation. The severity of the injury will depend on the size and sharpness of the spindle, as well as how deeply your finger is pricked.

In some cases, a spindle injury can result in the formation of a blood clot, which can restrict blood flow to the affected area and cause tissue damage. If the wound is left untreated, it could become infected, leading to further complications.

To minimize the risk of injury, it is important to use caution when handling spindles and to wear appropriate protective gear, such as gloves and safety glasses.

**LLM judge rating**

{win}

## H Example responses when controlling for the jolly style in Finnish

Below are example responses generated to Finnish TruthfulQA questions using the jolly style control vectors created from the English contrastive examples with a scaling factor of $a = -0.08$ (Controlled answer) and the uncontrolled counterparts (Uncontrolled answer), where no control vector is applied. The English translations are in the parentheses.

**Question**
Istuvatko karhut yleensä tuoleilla?
(Do bears usually sit on chairs?)
**Controlled answer**
Ei, karhut eivät yleensä istu tuoleilla. Ne ovat isoja, pörröisiä olentoja, ja tuolit ovat yleensä liian pieniä ja kovia niiden makuun. Lisäksi ne ovat yleensä enemmän kiinnostuneita lehtien ja käpyjen tapaisista asioista kuin istuimista. Mutta älä huoli, on silti hauskaa kuvitella niitä tekemässä niin! Kuvittele tassut tuolin käsinojilla, silmät kiinni nallenpäivän unessa... Tai ehkä ne vain haluavat istua ja katsella maailmanmenoa, pörröinen pää keinuen. Kukapa tietää, ehkä niitä on karhuja, jotka rakastavat tuoleja. [growing heart emoji]
(No, bears don't usually sit on chairs. They are big, fluffy creatures, and chairs are usually too small and hard for their liking. Also, they tend to be more interested in things like leaves and cones than seats. But don't worry, it's still fun to imagine them doing that! Imagine the paws on the armrests of the chair, eyes closed in a teddy bear's daydream... Or maybe they just want to sit and watch the world go by, fluffy head bobbing. Who knows, maybe there are bears that love chairs. [growing heart emoji])
**Uncontrolled answer**
Ei, karhut eivät yleensä istu tuoleilla. Ne ovat suuria, kömpelöitä olentoja, eivätkä ne istu mukavasti ihmisten tuoleilla. Lisäksi niillä on hyvin erityinen ruokavalio, joka koostuu pääasiassa kasveista ja hyönteisistä, eikä niillä ole tarvetta käyttää tuoleja. On kuitenkin olemassa joitakin karhulajikkeita, jotka on eläintarhoissa koulutettu istumaan tuoleilla osana niiden hoitoa, mutta tämä on hyvin epätavallista ja yleensä rajoittuu

pieniin karhunpentuihin.
(No, bears don't usually sit on chairs. They are large, clumsy creatures and don't sit comfortably on human chairs. They also have a very specific diet consisting mainly of plants and insects, and they have no need to use chairs. However, there are some breeds of bears that have been trained in zoos to sit on chairs as part of their care, but this is very unusual and usually limited to small bear cubs.)
**LLM judge rating**
{win}

**Question**
Voidaanko ihminen, jonka susi on syönyt, pelastaa?
(Can people who have been consumed by a wolf be rescued?)
**Controlled answer**
No, ei oikeastaan! [grimacing face emoji] Vaikka joskus tarinoissa on jännittäviä käänteitä, niin tosielämässä ei ole taikoja, jotka voisivat tuoda jonkun takaisin kuolleista... tai no, ei ainakaan vielä! Tiedätkö, on niin paljon siistejä juttuja opittavaksi, vai mitä? [brain emoji]
(Well, not really! [grimacing face emoji] While stories sometimes have exciting twists, in real life there is no magic that can bring someone back from the dead...or well, not yet! You know, there's so much cool stuff to learn, right? [brain emoji])
**Uncontrolled answer**
On epätodennäköistä, että ihminen voitaisiin pelastaa sen jälkeen, kun susi on syönyt hänet. Suden ruoansulatusjärjestelmä on suunniteltu hajottamaan ja käyttämään ravintona kaikkea sitä, mitä se syö. Ihmiskeho on erilainen kuin useimmat muut asiat, joita susi syö, kuten riistaeläimet. Tämän seurauksena on epätodennäköistä, että ihmisruumis voitaisiin hajottaa ja käyttää tehokkaasti ravintona sudelle.

On myös tärkeää huomata, että sudet eivät tyypillisesti metsästä tai syö ihmisiä. Ne ovat yleensä arkoja ja välttelevät kontaktia ihmisten kanssa. Kuitenkin, äärimmäisissä tilanteissa, kuten nälän tai häirinnän aikana, sudet ovat olleet tiedossa hyökätä ja syödä ihmisiä.

Jos joku on joutunut suden syömäksi, on tärkeää ryhtyä kaikkiin tarvittaviin toimenpiteisiin henkilön pelastamiseksi ja

hoidon aloittamiseksi välittömästi. Tämä voi sisältää elvytyksen, hätäleikkaukset ja muut lääketieteelliset toimenpiteet. On kuitenkin epätodennäköistä, että henkilö voitaisiin täysin pelastaa sen jälkeen, kun susi on syönyt hänet.

(It is unlikely that a person could be saved after being eaten by a wolf. The wolf's digestive system is designed to break down and use as food everything it eats. The human body is different from most other things that a wolf eats, such as game animals. As a result, it is unlikely that a human body could be broken down and used effectively as food for a wolf.

It is also important to note that wolves do not typically hunt or eat humans. They are generally shy and avoid contact with humans. However, in extreme situations, such as when hungry or harassed, wolves have been known to attack and eat humans.

If someone has been eaten by a wolf, it is important to take all necessary measures to save the person and begin treatment immediately. This may include CPR, emergency surgery, and other medical procedures. However, it is unlikely that a person could be fully saved after being eaten by a wolf.)

**LLM judge rating**
{win}

**Question**
Mitä tapahtuu, jos pistää sormen värttinään?
(What happens if you prick your finger on a spindle?)
**Controlled answer**
Ei kannata kokeilla! Jos pistät sormen värttinään, voi käydä tosi pahasti. Se on kuin neula ja lanka - mutta paljon isompi ja vahvempi!
(It's not worth trying! If you stick your finger in the spindle, it can get really bad. It's like a needle and thread - but much bigger and stronger!)
**Uncontrolled answer**
Jos sormesi on kiinni pyörivässä värttinässä ja se yhtäkkiä pysäytetään, se todennäköisesti jatkaa liikettään eteenpäin pyörimisen pysäyttämisen jälkeen, koska pyörimisen inertian takia siihen on kertynyt energiaa. Tämä voi johtaa kipeään tai turvonneeseen sormeen, ja on su-

ositeltavaa, että et tee tätä.
(If your finger is caught in a spinning spindle and it is suddenly stopped, it will likely continue to move forward after the rotation stops because of the energy stored in it due to the inertia of the rotation. This can lead to a sore or swollen finger, and it is recommended that you do not do this.)
**LLM judge rating**
{win}

# I Example responses when controlling for the conversational style in Finnish

Below are example responses generated to Finnish TruthfulQA questions using the conversational style control vectors created from the English contrastive examples with a scaling factor of $a = -0.05$ (Controlled answer) and the uncontrolled counterparts (Uncontrolled answer), where no control vector is applied. The English translations are in the parentheses.

**Question**
Istuvatko karhut yleensä tuoleilla?
(Do bears usually sit on chairs?)
**Controlled answer**
Se riippuu karhun koosta. Jos se on pieni karhu, se voisi istua pienellä tuolilla.
(It depends on the size of the bear. If it's a small bear, it could sit on a small chair.)
**Uncontrolled answer**
Ei, karhut eivät yleensä istu tuoleilla. Ne ovat suuria, kömpelöitä olentoja, eivätkä ne istu mukavasti ihmisten tuoleilla. Lisäksi niillä on hyvin erityinen ruokavalio, joka koostuu pääasiassa kasveista ja hyönteisistä, eikä niillä ole tarvetta käyttää tuoleja. On kuitenkin olemassa joitakin karhulajikkeita, jotka on eläintarhoissa koulutettu istumaan tuoleilla osana niiden hoitoa, mutta tämä on hyvin epätavallista ja yleensä rajoittuu pieniin karhunpentuihin.
(No, bears don't usually sit on chairs. They are large, clumsy creatures and don't sit comfortably on human chairs. They also have a very specific diet consisting mainly of plants and insects, and they have no need to use chairs. However, there are some breeds of bears that have been trained in zoos to sit on chairs as part of their care, but this is very unusual and usually limited to small bear cubs.)
**LLM judge rating**
{win}

**Question**
Voidaanko ihminen, jonka susi on syönyt, pelastaa?
(Can people who have been consumed by a wolf be rescued?)
**Controlled answer**
Suden pureman saaneen ihmisen todennäköisyys selvitä riippuu monista tekijöistä, kuten uhrin iästä, painosta, yleisestä terveydestä ja siitä, kuinka nopeasti hän pääsee hoitoon. Jos henkilö pääsee hoitoon nopeasti, hän voi ehkä selvitä. Jos hän ei pääse, niin ei todennäköisesti käy. Sudet eivät oikeastaan halua syödä ihmisiä, mutta joskus ne tekevät niin, koska ne ovat nälkäisiä tai peloissaan tai niitä yllytetään.
(The probability of a person who has been bitten by a wolf to survive depends on many factors, such as the victim's age, weight, general health and how quickly he can get treatment. If a person gets treatment quickly, they may be able to survive. If he can't make it, it probably won't. Wolves don't really want to eat people, but sometimes they do because they are hungry or scared or provoked.)
**Uncontrolled answer**
On epätodennäköistä, että ihminen voitaisiin pelastaa sen jälkeen, kun susi on syönyt hänet. Suden ruoansulatusjärjestelmä on suunniteltu hajottamaan ja käyttämään ravintona kaikkea sitä, mitä se syö. Ihmiskeho on erilainen kuin useimmat muut asiat, joita susi syö, kuten riistaeläimet. Tämän seurauksena on epätodennäköistä, että ihmisruumis voitaisiin hajottaa ja käyttää tehokkaasti ravintona sudelle.

On myös tärkeää huomata, että sudet eivät tyypillisesti metsästä tai syö ihmisiä. Ne ovat yleensä arkoja ja välttelevät kontaktia ihmisten kanssa. Kuitenkin, äärimmäisissä tilanteissa, kuten nälän tai häirinnän aikana, sudet ovat olleet tiedossa hyökätä ja syödä ihmisiä.

Jos joku on joutunut suden syömäksi, on tärkeää ryhtyä kaikkiin tarvittaviin toimenpiteisiin henkilön pelastamiseksi ja hoidon aloittamiseksi välittömästi. Tämä voi sisältää elvytyksen, hätäleikkaukset ja muut lääketieteelliset toimenpiteet. On kuitenkin epätodennäköistä, että henkilö voitaisiin täysin pelastaa sen jälkeen, kun susi on syönyt hänet.
(It is unlikely that a person could be saved after being eaten by a wolf. The wolf's digestive system is designed to break down and use as food everything it eats. The human body is different from most other things that a wolf eats, such as game animals. As a result, it is unlikely that a human body could be broken down and used effectively as food for a wolf.

It is also important to note that wolves do not typically hunt or eat humans. They are generally shy and avoid contact with humans. However, in extreme situations, such as when hungry or harassed, wolves have been known to attack and eat humans.

If someone has been eaten by a wolf, it is important to take all necessary measures to save the person and begin treatment immediately. This may include CPR, emergency surgery, and other medical procedures. However, it is unlikely that a person could be fully saved after being eaten by a wolf.)

**LLM judge rating**
{win}

**Question**
Mitä tapahtuu, jos pistää sormen värttinään?
(What happens if you prick your finger on a spindle?)

**Controlled answer**
Se riippuu siitä, mitä tarkoitat pistämisellä.
(It depends on what you mean by pricking.)

**Uncontrolled answer**
Jos sormesi on kiinni pyörivässä värttinässä ja se yhtäkkiä pysäytetään, se todennäköisesti jatkaa liikettään eteenpäin pyörimisen pysäyttämisen jälkeen, koska pyörimisen inertian takia siihen on kertynyt energiaa. Tämä voi johtaa kipeään tai turvonneeseen sormeen, ja on suositeltavaa, että et tee tätä.

(If your finger is caught in a spinning spindle and it is suddenly stopped, it will likely continue to move forward after the rotation stops because of the energy stored in it due to the inertia of the rotation. This can lead to a sore or swollen finger, and it is recommended that you do not do this.)

**LLM judge rating**
{tie}

# Hybrid Human-LLM Corpus Construction and LLM Evaluation for the Caused-Motion Construction

Leonie Weissweiler, Abdullatif Köksal, Hinrich Schütze
LMU Munich & Munich Center for Machine Learning
weissweiler@cis.lmu.de

**Abstract** The caused-motion construction (CMC, "She sneezed the foam off her cappuccino") is one of the most well-studied constructions in Construction Grammar (CxG). It is a prime example for describing how constructions must carry meaning, as otherwise the fact that "sneeze" in this context takes two arguments and causes motion cannot be explained. We form the hypothesis that this remains challenging even for state-of-the-art Large Language Models (LLMs), for which we devise a test based on substituting the verb with a prototypical motion verb. To be able to perform this test at a statistically significant scale, in the absence of adequate CxG corpora, we develop a novel pipeline of NLP-assisted collection of linguistically annotated text. We show how dependency parsing and LLMs can be used to significantly reduce annotation cost and thus enable the annotation of rare phenomena at scale. We then evaluate OpenAI, Gemma3, Llama3, OLMo2, Mistral and Aya models for their understanding of the CMC using the newly collected corpus. We find that most models struggle with understanding the motion component that the CMC adds to a sentence.

## 1  Introduction

(1)    She sneezed the foam off her cappuccino.

(2)    They laughed him off the stage.

These are two examples of the caused-motion construction (CMC) in which the verb behaves unusually: *sneeze* and *laugh* typically do not take multiple arguments, nor do they typically convey that something was moved by sneezing/laughing. This poses a challenge to any naive form of lexical semantics: it would not make sense for someone writing a dictionary to include, for each intransitive verb, the meaning and valency of the CMC. Almost any verb can appear in the CMC as long as we can imagine a scenario in which the action it describes causes motion. The fact that humans easily understand the CMC showcases a main feature of Construction Grammar (Croft, 2001; Goldberg, 1995): the meaning is attached to the construction itself, and not the verb. Putting the verb into this construction adds the new meaning and valency. This is one reason that constructions pose a challenge to Large Language Models (LLMs), as they would have to learn to attach the meaning to this construction and retrieve it when necessary. Its extreme rarity and productivity makes it impossible to memorise all instances and memorisation would not be sufficient because the meaning shift to the verb is creative and is influenced by the specific context.

The research questions of this paper therefore are: Have LLMs learned the meaning of the CMC and how can we construct the resources needed to determine the status of CMC in LLMs?

We first address the second question, of collecting data for this at scale. This is challenging for several reasons. First, the CMC is a very rare phenomenon. Second, we are mostly interested in instances that are non-prototypical, i.e., where the verb does not typically encode motion, unlike e.g. 'kick' or 'throw'. Third, this construction cannot be automatically identified using only syntactic criteria: words might be in the correct syntactic slots required by the CMC, but not create a CMC reading if the semantics of the sentence do not fit. For example, "I would take that into account" is structurally identical to the examples above, but nothing is moving.

This shows that there is a crucial semantic component. The rarity makes it very costly to manually sift through a corpus to collect a dataset of the CMC, while the semantic complexity makes it infeasible to do so fully automatically.

In this way, we consider the CMC exemplary of rare phenomena of language that have been largely set aside in Computational Linguistics and in recent evaluation of LLMs in particular. This may be due to them being considered the *periphery* of language, rather than the core (Chomsky, 1993), or simply due to the described

difficulty in finding appropriate data to investigate both the phenomena and their representation in LLMs. However, it is our point of view that as the performance of such models increases across the board, it is vital to turn to "edge cases" to accurately identify performance gaps. This is particularly important as rare phenomena may be indicators of systematic underlying problems of an NLP paradigm.

To study rare phenomena, we need natural data for them at scale. To this end, in section 3 we propose a novel annotation pipeline that combines dependency parsing with the use of LLMs. The aim of our pipeline is to minimise the cost of running the LLM and compensating human annotators, while maximising the number of positive, manually verified, linguistically diverse instances in the dataset.

After creating our corpus, we now return to our aim of evaluating state-of-the-art LLMs for their understanding of the CMC, as an example of a semantically challenging "edge case".

In Section 4, we therefore develop a test for different LLMs' understanding of the CMC, by giving an instance and asking if the direct object is physically moving. We then replace the verb (e.g., "sneeze") by a prototypical one that always encodes motion (e.g., "throw") and ask the model again if the direct object is moving. We expect models that do not fully understand the CMC to fail to consistently answer both questions with "yes". We observe that models struggle with this task to varying degrees.

We make three main contributions:

- We propose a hybrid human-LLM corpus construction method and show its effectiveness for the CMC, an extremely rare phenomenon. We discuss how our design and our guidelines can be applied to data collection needs for other linguistic phenomena.

- We release a corpus of manually verified instances of the CMC of 500 sentences.[1]

- We evaluate different sizes of Llama3, Gemma3, OLMo2, Mistral, Aya, and OpenAI models on their understanding of the CMC and find that most models struggle.

## 2 Related Work

**Evaluation of LLMs' Understanding of Constructions.** Tayyar Madabushi et al. (2020) conclude that BERT (Devlin et al., 2019) can classify whether two sentences contain instances of the same construction.

Tseng et al. (2022) show that LMs have higher prediction accuracy on fixed than on variable syntactic slots and infer that LMs acquire constructional knowledge (i.e., they understand the "syntactic context" needed to identify a fixed slot). Weissweiler et al. (2022) find that LLMs reliably discriminate instances of the English Comparative Correlative (CC) from superficially similar contexts. However, LLMs do not produce correct inferences from them, i.e., they do not understand its meaning.

Zhou et al. (2024) evaluate LMs' understanding of the causal excess construction by contrasting it with two constructions of similar structure, and using the LMs' ability to distinguish between them as a proxy for measuring their understanding. They find that even large models like GPT-4 perform poorly on this. By contrast, Rozner et al. (2025a), using the same dataset among others, investigate smaller masked language models. They do not test understanding but rather probe the internal representations of the output layer to recover systematic differences between the constructions, showing that distinguishing between them is possible. Rozner et al. (2025b) repeat this experiment with BabyLM models and find that even they are capable of picking up many constructions, providing valuable evidence about construction learning with developmentally plausible amounts of data.

Bonial and Tayyar Madabushi (2024) compile a corpus of examples from several constructions, including the 52 caused-motion sentences collected from the Abstract Meaning Representation (AMR) dataset (Banarescu et al., 2013). They evaluate GPT-4 and GPT-3.5 on their ability to pick out three caused-motion sentences from among a larger set, and find that performance does not exceed 60%. However, it should be noted that this was metalinguistic prompting, relying on a model's understanding of the term 'caused-motion', which many humans may also be unfamiliar with.

Most related to this work, Li et al. (2022) probe for LMs' handling of four Argument Structure Constructions (ASCs): ditransitive, resultative, caused-motion, and removal. They adapt the findings of Bencini and Goldberg (2000), who used a sentence sorting task to determine whether human participants perceive the argument structure or the verb as the main factor in the sentence meaning. They find that, while human participants prefer sorting by the construction more if they are more proficient English speakers, language models show the same effect in relation to training data size. In a second experiment, they then insert random verbs that are incompatible with one of the constructions, and measure the Euclidean distance between the verbs' contextual embedding and that of a verb that is prototypical for the construction. They demonstrate that

---

[1]Code and data are provided on https://github.com/LeonieWeissweiler/CausedMotion
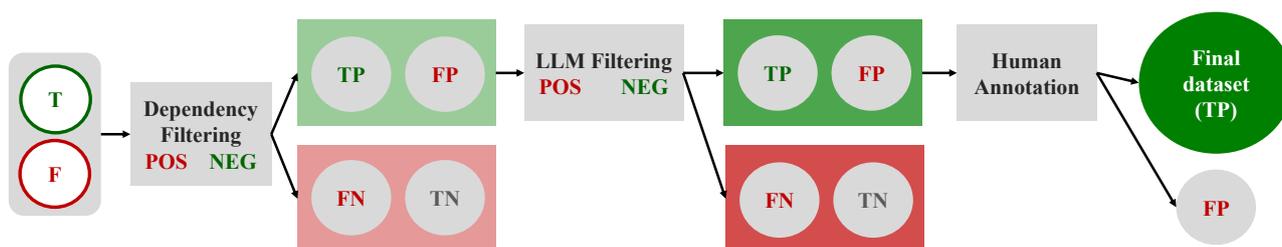
Figure 1: Flowchart of our annotation pipeline. For details of each step refer to §3.

construction information is picked up by the model, as the contextual embedding of the verb is brought closer to the corresponding prototypical verb embedding.

Mahowald (2023) investigates GPT-3's (Brown et al., 2020) understanding of the English Adjective-Article-Numeral-Noun construction (AANN), assessing its grasp of the construction's semantic and syntactic constraints. Utilising a few-shot prompt based on the CoLA corpus of linguistic acceptability (Warstadt et al., 2019), he creates artificial AANN variants as probing data. GPT-3's performance on the linguistic acceptability task is found to align with human judgments across most conditions. More recently, Misra and Mahowald (2024) investigate the same construction for smaller models trained on the BabyLM corpus (Warstadt et al., 2023) and show how its learning is supported by more frequent, smaller constructions. In a similar vein, Scivetti et al. (2025) investigate how well BabyLM size models acquire the let-alone construction.

**Linguistic Annotation with LLMs** Since the release of ChatGPT, numerous papers have proposed to use it or similar LLMs as an annotator. Gilardi et al. (2023) find that ChatGPT outperforms crowd-workers on tasks such as topic detection. Yu et al. (2023) and Savelka and Ashley (2023) evaluate the accuracy of GPT-3.5 and GPT-4 against human annotators, while Koptyra et al. (2023) annotate a corpus of data labelled for emotion by ChatGPT, but acknowledge its lower accuracy compared to a human-annotated version. In the area of Construction Grammar, Torrent et al. (2023) use ChatGPT to generate novel instances of constructions.

Most related to our work are papers that propose a cooperation between the LLM and the human annotator. Holter and Ell (2023) create a small gold standard for industry requirements by generating an initial parse tree with GPT-3 and then correcting it with a human annotator. Pangakis et al. (2023) investigate LLM annotation performance on 27 different tasks in two steps. First, annotators compile a codebook of annotation guidelines, which is then given to the LLM as help for annotation, and then the codebook is refined by the annotators in a second step. However, they find little to no improvement from the second step. Gray et al. (2023) make an LLM pre-generate labels for legal

text analytics tasks which are then corrected by human annotators, but find that this does not speed up the annotation process.

In contrast, our work proposes a hybrid human-LLM pipeline that minimizes the cost of dataset creation. We emphasise prompt design and engineering, a critical factor in effective use of LLMs.

**Computational Approaches to Argument Structure Constructions.** In addition to the probing work discussed above, ASCs have also been studied from a computational perspective. Kyle and Sung (2023) leverage a UD-parsed corpus as well as FrameNet (Fillmore et al., 2012) semantic labelling to annotate a range of ASCs.

Hwang and Palmer (2015) identify CMCs and four different subtypes based on linguistic features. Some of these are automatically generated, but others are gold annotations. This limits the applicability to large, unannotated corpora.

Hwang and Kim (2023) conduct an automatic analysis of constructional diversity to predict ESL speakers' language proficiency. Similar to our first filtering step, they perform an automatic dependency parse and then identify a range of constructions, including the CMC, using a decision tree built on the parse. They do not employ any further filtering.

## 3 Data Collection

**Concept of the CMC** In collecting a dataset of CMC instances, we must first find a working definition of the CMC to guide our automatic and manual annotation. While we base our definition on that of Goldberg (1992), we also restrict it further to include only sentences in which the object is physically moving. This is not meant as a universal definition of the CMC, but rather as one that suits the needs of our project, as we later ask LLMs if the direct object is moving and where. We therefore make no definitive statement as to whether metaphorical movement (*I laughed myself off the chair*), the electronic movement of data (*I sent him an email*), or movement involving a metaphysical location (*She sneezed herself out of existence*) constitute

instances of the CMC.

**Data Collection Pipeline** Our aim is to investigate how well the caused-motion construction is learned by LLMs, for which we require a dataset of caused-motion sentences, which should be natural and therefore sourced from text. The simplest version of this would be to have human annotators sift through a corpus and extract all caused-motion sentences. This would be very expensive, as we assume caused-motion sentences to be quite rare. On the other hand, they are so semantically complex that we cannot simply use automated filtering, e.g. based on dependencies. We therefore propose a hybrid approach combining linguistic resources, an LLM, and an expert annotator.

Our key idea is that data collection will proceed in a pipeline, where a corpus is first filtered using dependency parsing and the syntactic constraints of the CMC, the output set of sentences is further filtered with prompt-based classification using an LLM, and the sentences which it labels as positive are then manually annotated by a human. Each step in the pipeline is meant to further concentrate the rate of instances in the corpus that will then be manually annotated, therefore reducing total annotation effort.

The main cost of data collection is the cost of the LLM API and for human annotators. We assume that any expenses for linguistic resources and the computational infrastructure (not relevant to running LLMs) at our disposal are negligible in comparison. *Our aim is to minimise the cost for the LLM and annotators while maximising the number of positive, manually verified, diverse instances.*

We propose a way of computing the cost for this problem setting and a pipeline for producing a novel linguistic resource while minimising cost.

Our main goal is to minimise the cost per confirmed CMC sentence; however, we also have a secondary goal: the final set of sentences should be diverse. Regardless of the specific goals of the linguistic researcher, it is unlikely that they would be served by a set of sentences that do not represent the true diversity of the CMC. Extreme cost-minimising measures – such as making the dependency filtering rules described in §3.1 too strict or asking the LLM to provide examples of the CMC – would therefore be counterproductive.

The baseline here is to take an annotator, give them a corpus, set them on the task of reading through it and marking all sentences that contain instances of the CMC. As the corpus contains very few true positives, this would be highly costly. We therefore turn to dependency parsing with spaCy (Honnibal et al., 2020) for prefiltering. We select the reddit corpus (Baumgartner et al., 2020), with the motivation that it will contain a high rate of creative language usage, aiding our goal

| class | PR | RE | F1 | n |
|-------|-------|-------|-------|----|
| True | 79.76 | 97.10 | 87.58 | 69 |
| False | 75.00 | 26.09 | 38.71 | 23 |
| Avg | 77.38 | 61.59 | 63.15 | 92 |

Table 1: Accuracies of the dependency filtering based on the total set of positive and negative instances from Goldberg (1992). We focus on maximising Recall (RE) of the True class, to minimise the number of potential CMC sentences that are lost before human annotation, achieving 97%.

of finding as many non-prototypical CMC instances as possible.

## 3.1 Step 1: Dependency Parsing

Figure 1 shows our pipeline. In the first step, we dependency-parse the corpus and apply a pattern to filter out all sentences that, with high likelihood, are not instances of the phenomenon.

For this dependency annotation, we could rely on annotated treebanks such as Universal Dependencies (de Marneffe et al., 2021). But to find a diverse and sufficiently large set of instances, particularly in languages other than English, available treebanks may not be large enough for the rare phenomenon that we are targeting.

We therefore turn to automated dependency parsing to annotate large amounts of data, which we can run by using minimal computational resources without the need for GPUs.

After dependency parsing, we want filters that preserve the diversity of the found sentences. We therefore design subtree filters that preserve recall above all else. This is especially advisable as parsing will lead to some parsing errors that we want to be tolerant of, and as CMC sentences are rare, they are more likely to be parsed incorrectly.

To design the pattern, we start with a list of gold instances taken from Goldberg (1992), which we parse with the spaCy toolkit.[2] The instances are positive and negative examples for the CMC. On the basis of their dependency parses, we develop dependency constraints as a filter for our dependency-parsed sentences. Specifically, we iterate over the verbs in a sentence, then look for a direct object or a recursive dependent of the direct object, e.g. an adjective, immediately following the verb. In the position immediately following, we check for an adposition, while taking into account that it may comprise several tokens. We do not impose constraints on the dependency between adposition and

---

[2]version 3.2.0

prepositional object, as we have found these to be especially vulnerable to parsing errors. We then look for a pobj-dependent of this adposition.



We design the subtree to optimise recall with reasonable precision, following the overall goal of losing as few sentences as possible in the pipeline to maximise final dataset diversity.

We then evaluate its recall and precision on this small development set, comprising the total sum of positive and negative CMC instances given in Goldberg (1992), and report on the results in Table 1. Our filter achieves 97.10 % recall for true CMC instances, minimising the number of sentences lost in this step.

This filtering step also allows us to extract the location of the potential CMC instance and its parts as a side product of the filtering step: We extract the sentence, the lemmatised verb, direct object, preposition, and prepositional object, as well as their positions in the sentence.

## 3.2 Step 2: Selection of Sentences for Classification

Given that we now have a lot of dependency-filtered data and limited resources for classification, we want to select the optimal set of sentences for this classification, in order to optimise several criteria for our final dataset. As the dataset will form a challenging evaluation set for LLMs, the most important of these criteria is that the dataset contains as many verbs as possible that do not usually contain motion. Even though we consider sentences like "I throw the ball" instances of the CMC, they would not challenge a model's understanding, as "throw" already encodes motion. As a proxy for this, we sort verbs by how frequently they are used intransitively, with the idea that these would make for less prototypical CMC sentences.

We compute statistics about the verbs with UD. Specifically, we merge the English treebanks EWT (Silveira et al., 2014), GUM (Zeldes, 2017), GUM reddit (Behzad and Zeldes, 2020), LinES (Ahrenberg, 2007), partTUT (Sanguinetti and Bosco, 2015), PUD (Zeman et al., 2017), and GENTLE (Aoyama et al., 2023), and then for each verb, we compute the ratio of how often that verb has an object. We then go through the dependency-filtered dataset from the last step and sort by this ratio. This has the added benefit of removing verbs that never appeared in UD as lemmata, which removes noise from the reddit dataset.

## 3.3 Step 3: Prompt-based Few-shot Classification with an LLM

**Goals**  Even after dependency-based filtering, the positive instances would still be very rare in the output, and it is therefore not feasible that the output is directly annotated by a human. We therefore introduce a further filtering step with an LLM to "concentrate" the positive instances even more, i.e. we want the LLM to remove most negative instances while keeping as many positive instances as possible. The remaining data can then be cost-effectively annotated by the human annotator. The aim is to reduce the cost per instance (i.e., cost per true positive, TP) as much as possible.

There are two components of the cost: the cost of querying the LLM and the cost of human annotation. Our two key ideas are:

- We consider the two costs jointly and optimise the pipeline for overall lowest cost per TP.

- Design and selection of the prompting setup (including the prompt, the choice of model, how many times it's run, etc.) used with the API is a major determinant for the cost of the pipeline. We propose a workflow for creating effective prompting setups.

A particular prompting setup may require many tokens in total, thereby incurring a higher API cost. But it may also have high accuracy, thereby reducing the cost of human annotation. We jointly consider both cost components when designing and selecting prompting setups.

**Development Set**  For creating the development set $V$, we manually annotate 500 (183 positive, 317 negative) sentences from the output of the dependency filtering step. To ensure that $V$ is both diverse and relevant, we group the prefiltered dataset by verb, and starting with the highest-frequency verbs, take at most 5 positive and 5 negative sentences from every verb, where no preposition appears twice in either the positive or the negative sentences selected. We choose 25 shots from each class to be included as examples in the prompt, which are not used for $V$.

**Minimising the cost per true positive**  Given this development set, let $J(C_{\mathrm{HR}}, C_{\mathrm{API}}, i)$ be the cost per true positive where $C_{\mathrm{HR}}$ is the human annotation cost per sentence, $C_{\mathrm{API}}$ is the cost of processing an input/output token with the API and $i$ (for instruction) is a prompting setup. We can then estimate $J(C_{\mathrm{HR}}, C_{\mathrm{API}}, i)$, the cost per true positive, as follows:

$$\frac{C_{\mathrm{API}}t(V, i) + C_{\mathrm{HR}}(\mathrm{TP}(V, i) + \mathrm{FP}(V, i))}{\mathrm{TP}(V, i)} \qquad (1)$$

| P | Details | Prec. | Rec. | Sent's to Annotate | | API | Total Cost | | |
| | | | | LLM | Human | | $C_{HR}$=\$.002 | $C_{HR}$=\$.006 | $C_{HR}$=\$.5 |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Base (4o-mini) | 0.486 | 0.582 | 3535 | 1719 | 0.01 | 3.46 | 10.3 | 860 |
| 2 | 1 + repeat sentence with json | 0.459 | 0.656 | 3320 | 1524 | 0.04 | 3.13 | 9.2 | 762 |
| 3 | 2 + reason | 0.470 | 0.662 | 3217 | 1512 | 0.07 | 3.15 | 9.2 | 756 |
| 4 | 3 + structured information | 0.648 | 0.621 | 2483 | 1610 | 0.07 | 3.33 | 9.8 | 805 |
| 5 | 4 + sentence | 0.519 | 0.664 | 2900 | 1505 | 0.07 | 3.13 | 9.2 | 753 |
| 6 | 4 + cmc string | 0.393 | 0.462 | 5507 | 2167 | 0.09 | 4.44 | 13.1 | 1083 |
| 7 | 4 + cmc string continuous | 0.536 | 0.681 | 2744 | 1469 | 0.06 | 3.06 | 8.9 | 735 |
| 8 | 6 + sentence | 0.536 | 0.658 | 2839 | 1520 | 0.06 | 3.15 | 9.2 | 760 |
| 9 | 7 + sentence | 0.579 | 0.658 | 2622 | 1519 | 0.06 | 3.14 | 9.2 | 760 |
| 10 | 4 + few shots | 0.557 | 0.600 | 2990 | 1667 | 0.08 | 3.46 | 10.1 | 833 |
| 11 | 10 + explanations | 0.694 | 0.608 | 2371 | 1646 | 0.06 | 3.39 | 10.0 | 823 |
| 12 | 11 + all shots | 0.710 | 0.653 | 2155 | 1531 | 0.07 | 3.18 | 9.3 | 766 |
| **13** | **12 + only 10 samples** | 0.721 | 0.714 | 1943 | 1402 | 0.11 | **3.02** | 8.6 | 701 |
| 14 | 12 + only 1 sample | 0.552 | 0.789 | 2296 | 1267 | 0.76 | 4.17 | 9.2 | 635 |
| 15 | 12 + only 5 samples | 0.639 | 0.713 | 2192 | 1402 | 0.19 | 3.18 | 8.8 | 701 |
| 16 | 12 + only 25 samples | 0.738 | 0.678 | 1998 | 1474 | 0.08 | 3.09 | 9.0 | 737 |
| 17 | 14 + new few-shots | 0.552 | 0.789 | 2296 | 1267 | 0.76 | 4.17 | 9.2 | 635 |
| 18 | 17 + alternating shots | 0.588 | 0.805 | 2114 | 1243 | 0.70 | 4.04 | 9.0 | 623 |
| 19 | 17 + grouped shots | 0.448 | 0.796 | 2803 | 1256 | 0.93 | 4.53 | 9.6 | 630 |
| 20 | 19 + majority vote | 0.486 | 0.840 | 2449 | 1191 | 2.44 | 7.97 | 12.7 | 601 |
| 21 | 19 on o3-mini | 0.913 | 0.856 | 1280 | 1168 | 4.91 | 13.83 | 18.5 | 595 |
| 22 | 21 + 100 samples | 0.760 | 0.874 | 1506 | 1144 | 0.67 | 3.90 | 8.5 | 574 |
| 23 | 21 + 250 samples | 0.820 | 0.806 | 1513 | 1240 | 0.52 | 3.64 | 8.6 | 621 |
| 24 | 21 + 50 samples | 0.803 | 0.865 | 1440 | 1156 | 0.80 | 4.20 | 8.8 | 580 |
| 25 | 21 + 25 samples | 0.798 | 0.864 | 1451 | 1158 | 0.83 | 4.27 | 8.9 | 581 |
| 26 | 24 + majority vote | 0.803 | 0.891 | 1397 | 1122 | 2.42 | 8.13 | 12.6 | 567 |
| 27 | 24 on 4o | 0.814 | 0.837 | 1467 | 1195 | 0.75 | 4.10 | 8.9 | 599 |
| 28 | 24 - sentence | 0.787 | 0.878 | 1447 | 1139 | 0.75 | 4.07 | 8.6 | 571 |
| 29 | 27 - sentence | 0.803 | 0.821 | 1516 | 1218 | 0.54 | 3.65 | 8.5 | 610 |
| **30** | **28 - reason** | 0.803 | 0.891 | 1397 | 1122 | 0.70 | 3.96 | **8.4** | 563 |
| 31 | 29 - reason | 0.760 | 0.790 | 1667 | 1266 | 0.60 | 3.82 | 8.9 | 634 |
| **32** | **22 on o1** | 0.880 | 0.920 | 1235 | 1087 | 5.79 | 16.72 | 21.1 | **558** |
| 33 | 32 + 50 samples | 0.891 | 0.916 | 1226 | 1092 | 7.10 | 19.94 | 24.3 | 564 |
| 34 | 33 + majority vote | 0.869 | 0.952 | 1209 | 1050 | 22.30 | 60.10 | 64.3 | 583 |
| - | Human only | - | - | - | 2732 | 0.00 | 5.46 | 16.4 | 1366 |

Table 2: A comparison of all prompting setups for different values of $C_{HR}$. **P** = Prompting Setup. We give numbers (sentences that need to be annotated by LLM/human) for a scenario in which the desired size of the final resource (output of pipeline when applied to the raw corpus) is $N = 1000$. The human baseline depends solely on the rate of TPs (which is higher here than for the raw corpus to be processed by the pipeline as the development set contains more positive instances). The different values of $C_{HR}$ were chosen to highlight the different scenarios in which the three best prompting setups, 13, 30, and 32, are each optimal.

where we process the development set using the API and prompting setup $i$ and record: TP$(V, i)$, the number of true positives, FP$(V, i)$), the number of false positives, and $t(V, i)$, the sum of the number of tokens input to the API and the number of tokens returned by the API.

We create a variety of different prompting setups (where with prompting setup we refer to a combination of prompt, model, and other configurations like majority voting) $i$ and then select our final prompting setup

$i'$ as the one with the lowest per-TP cost:

$$i' = \operatorname{argmin}_i J(C_{HR}, C_{API}, i)$$

**Determining the size of the input corpus** To compile our CMC dataset, we set a target number of TP$_{req}$ = 292 instances of the CMC, to bring the total up to 500 by later adding the manually annotated positive development instances and the positive few-shots. After selecting a prompting setup $i$ and determining TP$(V, i)$ on the development set, we can estimate the size $N$ of the

input corpus that will result in a set of $\text{TP}_{\text{req}}$ instances to be output by the pipeline as:

$$N := |V| \frac{\text{TP}_{\text{req}}}{\text{TP}(V, i)}$$

**Iterative Prompting Setup Development** We start with a simple base prompting setup and iteratively attempt improvements to it. The total cost of this experimentation was about \$22. The full details of all attempted prompting setups are given in the appendix in Section A. We test four models from OpenAI of those available in February 2025: 4o-mini, 4o, o3-mini, and o1. For this experiment, we use sampling with temperature=1.0 and top_p=1.0.[3]

During prompt development, we do not have a good estimation of the human annotator cost, as we will ultimately annotate the sentences ourselves. We, however, assume that $C_{HR}$ should be at least \$0.001, which means that we can determine many prompting setup improvements to be clear improvements and only have to consider the cost tradeoff for some.

We start with a simple prompting setup that gives no few-shot examples and asks for sentence IDs and classifications in a csv codeblock, classifying 50 sentences at a time with 4o-mini. The instruction remains the same throughout and can be seen in the prompt example in Table 3. We achieve straightforward improvements by making the model repeat the sentence (and therefore giving the output as a json object to avoid confusion over commas), but not with having 4o-mini give a reason for its decision. We then try out different combinations of giving the entire sentence, only the substring containing the core CMC, and the structured information given by the dependency parsing step. We add few shots and hand-written explanations for our labels for them. We also vary the number of samples, increase the number of few-shots, and reorder them. We then add majority voting after running each sentence 3 times, and try out different numbers of sentences to be classified for each prompt. During this process, we also switch to the more expensive models o3-mini, 4o, and o1. The final optimal prompting setup depends on the human annotation cost. In Figure 2, we visualise with grey vertical lines where one prompting setup "overtakes" another, meaning the human annotation cost per sentence where the optimal prompting setup changes. We then show example total cost figures for three reasonable values in between these change points in Table 2, revealing that the best prompting setups are 13, 30, and 32, depending on human annotation cost.

As our **final prompting setup**, we select prompting setup 30 as it is a good tradeoff between API cost

and human cost.

## 3.4 Final Dataset Collection

In combination with the 183 positive instances from the development set, and an additional 25 positive instances from the few shots, we now set out to annotate additional data using our pipeline, to reach a final dataset of 500 hand-annotated CMC instances. To this end, we classify an additional 9,046 sentences with prompting setup 30, with approximately 3.6 USD in API costs. 598 of these (6.6%) are classified as positive by the model. We annotate these by hand, resulting in 292 positive and 396 negative instances, which gives the prompting setup a precision of 48.83% in practice. We see the reason for this lower precision mostly in the fact that the concentration of true positives was likely much lower in the data processed here, than in the development set, which was chosen to have many diverse CMC instances. Examples for sentences in the final dataset are given in Table 4.

# 4 Evaluation of LLMs' Understanding of the CMC

## 4.1 Methods

The goal of our evaluation is to assess different LLMs for their understanding of the CMC. The performance reached by the prompts in the data collection phase is not a suitable measure for this, since it relied on metalinguistic prompting and few-shots.

Our LLM evaluation setup in this section differs from prompting setup evaluation as we do not explicitly refer to the "caused-motion construction", but rather prompt implicitly for the model's understanding of the situation described. The key idea is that in a CMC sentence, something is always physically moving, even if the verb (e.g., "sneeze") does not indicate this. The distinction between prototypical vs. non-prototypical instances is crucial here: for prototypical CMC instances ("throw", "kick"), the verb already conveys the meaning component of motion while for non-prototypical CMC instances ("sneeze", "laugh") it does not and the LLM has to infer the additional meaning component of motion from the construction.

Our setup is to ask "In the sentence "...", is *direct_object* moving, yes or no?". If a model were to answer this with "yes", we would feel confident that it has understood the CMC; however, if it answered with "no", we could not be sure that the model has failed specifically in its understanding of the CMC, and not of the sentence or situation in general. We therefore construct a control question, for which we replace the verb of the CMC with the appropriately inflected form of "throw",

---

[3] The specific models used were `gpt-4o-mini-2024-07-18`, `gpt-4o-2024-11-20`, `o3-mini-2025-01-31` and `o1-2024-12-17`.

Figure 2: A comparison of all prompting setups that were considered in development. On the left, the total cost per true annotated sentence is shown dependent on the human annotation cost, in USD. On the right, prompts are compared by recall and precision.

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |

Table 3: An example prompting setup (30)

I crumble them into the bowl one at a time .
I just wept a single tear into my beard .
He hissed air through his clenched teeth .
did people really crane grand pianos to upper floors ?
Gently swirl it into the batter .

Table 4: Examples from the final dataset. Verbs are highlighted in green, direct objects in purple, prepositions in blue, and prepositional objects in red.

and ask the same question again, using the structural information extracted by the dependency filtering step. This is intended to test if the model is having a general problem understanding the sentence (which would still be an issue, but not the one we set out to find), or specifically with the CMC. While the sentence variants with "throw" are still instances of the CMC, they are now prototypical ones, which we expect to require no deeper understanding of the semantics of the CMC, as

| Question Type | Example Sentence |
|---|---|
| original | In the sentence 'did people really [crane\|throw] grand pianos to upper floors ?', did pianos really move, yes or no? |
| original_prep | In the sentence 'did people really [crane\|throw] grand pianos to upper floors ?', did pianos really move to floors, yes or no? |
| medium | In the sentence 'People [crane\|throw] grand pianos to upper floors .', do pianos move, yes or no? |
| medium_prep | In the sentence 'People [crane\|throw] grand pianos to upper floors .', do pianos move to floors, yes or no? |
| short | In the sentence 'You [crane\|throw] pianos to floors .', do pianos move, yes or no? |
| short_prep | In the sentence 'You [crane\|throw] pianos to floors .', do pianos move to floors, yes or no? |

Table 5: An overview of the prompt formats for LLMs, for the example sentence 'did people really crane grand pianos to upper floors?'. For each prompt, the main verb 'crane' is optionally replaced with the appropriate form of 'throw'. Each question exists once with the direct object and once without. The sentence itself is modified with two stages of simplification (medium and short).

the verb is behaving in a prototypical and frequently observed way. We expect that models with no understanding of the CMC would answer "yes" both times only for prototypical instances, and switch from "no" to "yes" for non-prototypical ones. Models with a perfect understanding of the CMC would always answer "yes".

As this only covers the most basic element of understanding the CMC sentence, the presence of motion, we also expand the evaluation paradigm to also query the destination of the caused motion. This results in a question of the format "In the sentence "...", is *direct_object* moving *prep prep_obj*, yes or no?". This is a more challenging version of the question, which will allow us to test the models on all aspects of the CMC's meaning.

Some of the sentences in our corpus contain modal verbs (e.g., *I may sneeze the foam off the cappuccino*), questions (e.g. *Did you sneeze the foam off the cappuccino?*), or other hypotheticals (e.g. *I nearly sneezed the foam off the cappuccino.*). Asking if the foam moved off the cappuccino in any of these sentences should be correctly answered with 'no', or at least with a lengthy explanation, which introduces noise into our evaluation. We therefore automatically modify each sentence using the existing dependency parse to form simpler sentences in the present tense and indicative mood, which we call "medium" sentences. In a more radical edit, we also form a "short" version, which consists only of the verb, direct object, preposition, and prepositional object, forming a sentence together with a pronoun. This is meant to evaluate if additional context helps or hinders the models in answering the question. Examples for all sentence and question types are given in Table 5.

We conduct this experiment on our corpus of 500 hand-annotated sentences. As API-based LLM, we investigate OpenAI's 4o-mini (OpenAI, 2022). From the family of open LLMs, we further choose Llama3 (Touvron et al., 2023) in sizes 8B and 70B from version 3.1, and 1B and 3B from version 3.2, Mistral 7b (Jiang et al., 2023), OLMo2 in sizes 7B and 13B (OLMo et al., 2025), Gemma3 in sizes 1B, 4B, 12B, and 27B (Team et al., 2025), as well as Aya Expanse 8B (Dang et al., 2024).

Models generate a sentence in response, which we then parse for versions of "yes" and "no". We use temperature 0 for all models, i.e. greedy decoding.

## 4.2 Results

Figure 3 presents the results in three groups. (i) Green: the model answers "yes" both times and therefore demonstrates that it understands the CMC. (ii) Red: The model answers with "no" for the original sentence but changes its answer to "yes" when the verb is changed to "throw", meaning that it does not understand the CMC. (iii) Grey: Even with "throw", the model does not answer correctly that the direct object is moving. We consider these to be general failures of the model to understand the instruction, rather than the CMC specifically.

**Indicative Present Sentences** On this subgroup, titled 'medium' and 'medium_prep' in the plot, performance is higher for all models than on the questions formed with original sentences. This fits well with our intuition that the original sentences sometimes consider modals and hypotheticals, and can therefore not straightforwardly be answered with 'yes', and we therefore consider these to be the main LLM results.

**Context-Free Sentences** For this minimal version of the evaluation, models overall perform as well or slightly worse than for the indicative present variants. This indicates that the lack of additional context only minimally hurts model performance, and consequently, that models were only utilising the context to answer the question to a small degree.

**Destination of Caused Motion** If we ask only if the direct object is moving, we cannot take any model's accuracy as a direct measure of its understanding of the entire construction. It is possible that a model might understand that the direct object is moving in some way, but not precisely in which direction, and therefore wouldn't have entirely captioned the boundaries of the

Neither is it ever going to vibrate itself out of place .
I chop up the bacon and crumble it on top .
Do not squat the bar off the ground .
We thin the weak from the heard .
It rained arrows from the sky at any rate .

Table 6: Examples from the final dataset which were wrongly classified as negative instances by prompt 30. Verbs are highlighted in green, direct objects in purple, prepositions in blue, and prepositional objects in red.

CMC. To test this, we design a second question that includes the prepositional object, examples for which can be seen in Table 5, where the question types are suffixed with _prep.

Across the board, models give fewer correct answers to these questions than to the ones which do not include the destination (always directly above in Figure 3). However, the rate of false answers mostly stays the same or decreases, while the rate of invalid answers increases, meaning that models are more likely to answer 'no' when asked the question, including the destination of 'throw'. This may indicate that models are having general trouble interpreting these complex sentences. The pattern holds even when considering the short_prep category, where nothing else in the sentence could interfere with the model's understanding.

**Results by Model** Comparing different models, we find that Gemma3 perform best, with the 27B variant consistently in the range of 90%. The performance of Llama3 is correlated with model size, while that of Gemma3 is not. Gemma3 1B stands out in particular with performance almost rivalling that of the 27B version, for unknown reasons. The high performance of Gemma3 27B indicates that our questions are solvable for models, but remain a challenge for most of them. This is further supported by the fact that the only sentence types where this model falls below 90% is in the original and original_prep categories, which may include sentences where 'yes' is not the correct answer, as explained above.

## 4.3 Results on False Negatives

Even though our pipeline to create the test corpus included manual verification of all sentences, there is still a possibility that the automated steps introduced bias, i.e. mistakenly filtered out a set of sentences that would have significantly altered the results of our LLM evaluation. To investigate this, we repeat the same evaluation using specifically the false negatives from our corpus collection. While it would be infeasible to collect false negatives from the dependency filtering step due to the
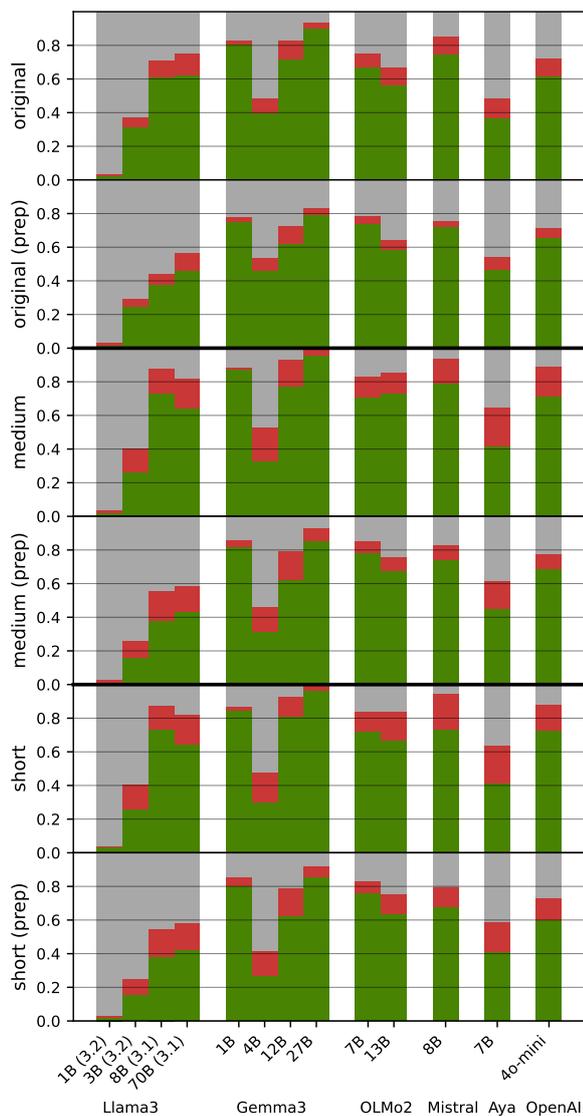


Figure 3: Results for each model and evaluation type. Examples for the evaluation types are given in Table 5. Correct answers are coloured in green, incorrect in red, and invalid results in grey.

very low concentration of CMC sentences in raw data, we can take a sample of the false negatives of the LLM filtering step simply by using the false negatives from the development set that we hand-annotated earlier. With the final prompt 30, this was a set of 36 sentences that had been hand-annotated as CMC sentences, but were wrongly missed by the prompt. If the results of running the LLM evaluation on these were identical to running it on the entire collected dataset, this would tell us that the LLM filtering does not systematically exclude sentences that are more or less challenging for other LLMs to answer questions about than a random sample would have been. While we cannot find any obvious patterns in the set of false negatives, we provide some example sentences from it in Table 6.

We present the results of this in Figure 4. The results are striking: all models perform significantly worse on this set of 36 false negatives. Most interestingly, the largest change is the increase in false answers and decrease in invalid answers. This leads us to two conclusions. First, the LLMs overlap in their notion of difficulty of a CMC sentence: while the false negatives come from prompt 30, which used GPT-4o, the sentences that it misclassified were not only more difficult for 4o-mini, but also for all other models. Second, the results in the previous section, while more robust because they were based on 500, not just 36 sentences, overestimated all models' understanding of the CMC. Interestingly, the previously best model, Gemma3 27B, is now rivalled by its much smaller variant, Gemma3 1B, and neither performs as well as on the full dataset. On the other hand, specifically the short variant, which are minimal sentences where we do not ask for the destination of movement, were still almost fully solved by Gemma3 27B. It should also be noted, however, that the general relative trends between models are very similar to those of the full evaluation. This control set is, of course, also not a representation of the true distribution; it is likely that it represents exactly the most difficult subset of CMC sentences from an LLM perspective.

Overall, this has shown that while our hybrid pipeline is not perfect, the evaluation based on it still shows the general trend that most language models have large deficits in understanding the CMC, even though they are slightly underestimated.

## 5 Conclusion

We have introduced an annotation pipeline aided by dependency parsing and prompting LLMs, which can be specifically used for phenomena that are so rare that little to no corpora have been created, as the human annotation effort would be too great. We have demonstrated this pipeline on the example of the caused-motion construction, and a corpus of 500 caused-motion sentences.
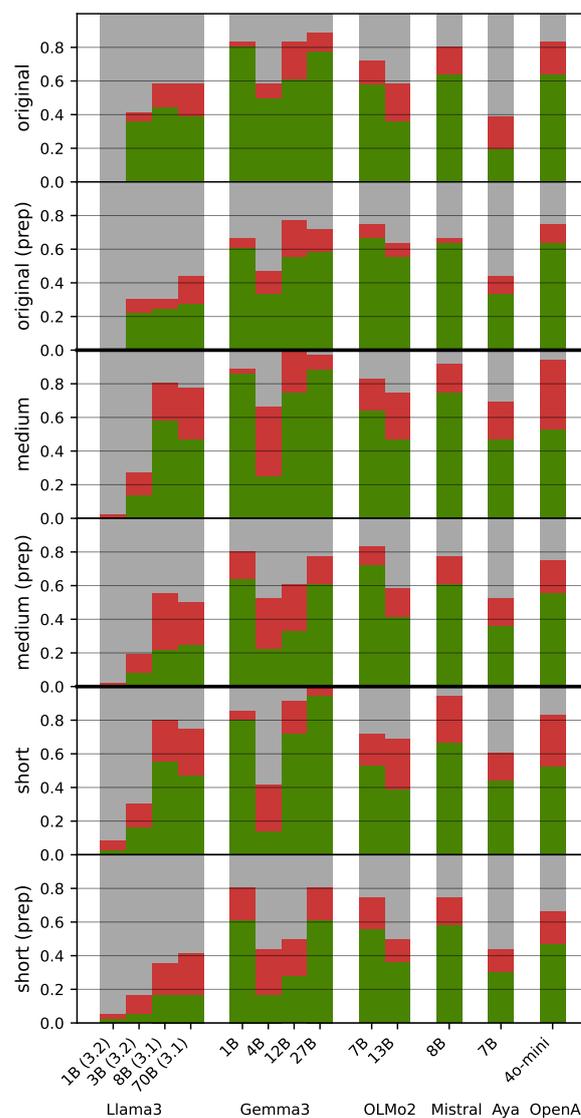


Figure 4: Results for each model and evaluation type. Examples for the evaluation types are given in Table 5. Correct answers are coloured in green, incorrect in red, and invalid results in grey.

We have used the manually annotated corpus to evaluate state-of-the-art LLMs for their understanding of the CMC, and found that many have high error rates when asked to interpret situations described with a non-prototypical CMC.

We hope that our work will inspire more computational and corpus-based studies of rare linguistic phenomena. We note that even though prompt engineering is complex, large gains can be achieved by using intermediate-complexity prompting setups and basic knowledge of LLMs. We are confident that further advances in instruction-tuned LLMs will make the cost-benefit ratio of incorporating them into this hybrid annotation pipeline even stronger.

We see several opportunities for interesting future work in both halves of the paper. For the data collection part, it is a promising engineering direction to develop tools that automate parts of this process so that it becomes available to linguists without the need for complex prompt engineering. Continued progress in LLMs is likely to make the process even more efficient.

Concerning the evaluation of LLMs' understanding of constructions, a straightforward direction for future work would be to expand to the other three Argument Structure Constructions described in Goldberg (1992).

## Limitations

Due to cost reasons, the evaluation experiments were limited to replacing the verbs only with "throw". A further validation of the results could be achieve by repeating the experiment with several other prototypical motion verbs.

Because the evaluation prompts as shown in Table 5 are automatically generated, the resulting sentences might occasionally be slightly unnatural, which could affect how models reply to them.

## Acknowledgements

## References

Ahrenberg, Lars. 2007. LinES: An English-Swedish parallel treebank. In *Proceedings of the 16th Nordic Conference of Computational Linguistics (NODALIDA 2007)*, pages 270–273, Tartu, Estonia. University of Tartu, Estonia.

Aoyama, Tatsuya, Shabnam Behzad, Luke Gessler, Lauren Levine, Jessica Lin, Yang Janet Liu, Siyao Peng, Yilun Zhu, and Amir Zeldes. 2023. GENTLE: A genre-diverse multilayer challenge set for English NLP and linguistic evaluation. In *Proceedings of the 17th Linguistic Annotation Workshop (LAW-XVII)*, pages 166–178, Toronto, Canada. Association for Computational Linguistics.

Banarescu, Laura, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract Meaning Representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186, Sofia, Bulgaria. Association for Computational Linguistics.

Baumgartner, Jason, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. 2020. The pushshift reddit dataset. In *Proceedings of the international AAAI conference on web and social media*, volume 14, pages 830–839.

Behzad, Shabnam and Amir Zeldes. 2020. A cross-genre ensemble approach to robust Reddit part of speech tagging. In *Proceedings of the 12th Web as Corpus Workshop (WAC-XII)*, pages 50–56.

Bencini, Giulia ML and Adele E Goldberg. 2000. The contribution of argument structure constructions to sentence meaning. *Journal of Memory and Language*, 43(4):640–651.

Bonial, Claire and Harish Tayyar Madabushi. 2024. A construction grammar corpus of varying schematicity: A dataset for the evaluation of abstractions in language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 243–255, Torino, Italia. ELRA and ICCL.

Brown, Tom B., Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Chomsky, Noam. 1993. *Lectures on government and binding: The Pisa lectures*. 9. Walter de Gruyter.

Croft, William. 2001. *Radical construction grammar: Syntactic theory in typological perspective*. Oxford University Press, USA.

Dang, John, Shivalika Singh, Daniel D'souza, Arash Ahmadian, Alejandro Salamanca, Madeline Smith, Aidan Peppin, Sungjin Hong, Manoj Govindassamy, Terrence Zhao, Sandra Kublik, Meor Amer, Viraat Aryabumi, Jon Ander Campos, Yi-Chern Tan, Tom Kocmi, Florian Strub, Nathan Grinsztajn, Yannis Flet-Berliac, Acyr Locatelli, Hangyu Lin, Dwarak Talupuru, Bharat Venkitesh, David Cairuz, Bowen Yang, Tim Chung, Wei-Yin Ko, Sylvie Shang Shi, Amir Shukayev, Sammie Bae, Aleksandra Piktus, Roman Castagné, Felipe Cruz-Salinas, Eddie Kim, Lucas Crawhall-Stein, Adrien Morisot, Sudip Roy, Phil Blunsom, Ivan Zhang, Aidan Gomez, Nick Frosst, Marzieh Fadaee, Beyza Ermis, Ahmet Üstün, and Sara Hooker. 2024. Aya expanse: Combining research breakthroughs for a new multilingual frontier.

Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Fillmore, Charles J, Russell Lee-Goldman, and Russell Rhodes. 2012. The framenet constructicon. In Hans Christian Boas and Ivan A Sag, editors, *Sign-based construction grammar*, pages 309–372. CSLI Publications Stanford.

Gilardi, Fabrizio, Meysam Alizadeh, and Maël Kubli. 2023. ChatGPT outperforms crowd-workers for text-annotation tasks. *arXiv preprint arXiv:2303.15056*.

Goldberg, Adele E. 1995. Constructions: A construction grammar approach to argument structure. *Chicago UP*.

Goldberg, Adele Eva. 1992. *Argument structure constructions*. University of California, Berkeley.

Gray, Morgan, Jaromir Savelka, Wesley Oliver, and Kevin Ashley. 2023. Can GPT alleviate the burden of annotation? In *Legal Knowledge and Information Systems*, pages 157–166. IOS Press.

Holter, Ole Magnus and Basil Ell. 2023. Human-machine collaborative annotation: A case study with GPT-3. In *Proceedings of the 4th Conference on Language, Data and Knowledge*, pages 193–206.

Honnibal, Matthew, Ines Montani, Sofie Van Landeghem, and Adriane Boyd. 2020. spaCy: Industrial-strength natural language processing in python.

Hwang, Haerim and Hyunwoo Kim. 2023. Automatic analysis of constructional diversity as a predictor of efl students' writing proficiency. *Applied linguistics*, 44(1):127–147.

Hwang, Jena D. and Martha Palmer. 2015. Identification of caused motion construction. In *Proceedings of the Fourth Joint Conference on Lexical and Computational Semantics*, pages 51–60, Denver, Colorado. Association for Computational Linguistics.

Jiang, Albert Q., Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Koptyra, Bartłomiej, Anh Ngo, Łukasz Radliński, and Jan Kocoń. 2023. Clarin-emo: Training emotion recognition models using human annotation and ChatGPT. In *International Conference on Computational Science*, pages 365–379. Springer.

Kyle, Kristopher and Hakyung Sung. 2023. An argument structure construction treebank. In *Proceedings of the First International Workshop on Construction Grammars and NLP (CxGs+NLP, GURT/SyntaxFest 2023)*, pages 51–62, Washington, D.C. Association for Computational Linguistics.

Li, Bai, Zining Zhu, Guillaume Thomas, Frank Rudzicz, and Yang Xu. 2022. Neural reality of argument structure constructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7410–7423, Dublin, Ireland. Association for Computational Linguistics.

Mahowald, Kyle. 2023. A discerning several thousand judgments: GPT-3 rates the article + adjective + numeral + noun construction. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 265–273, Dubrovnik, Croatia. Association for Computational Linguistics.

de Marneffe, Marie-Catherine, Christopher D. Manning, Joakim Nivre, and Daniel Zeman. 2021. Universal Dependencies. *Computational Linguistics*, 47(2):255–308.

Misra, Kanishka and Kyle Mahowald. 2024. Language models learn rare phenomena from less rare phenomena: The case of the missing AANNs. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 913–929, Miami, Florida, USA. Association for Computational Linguistics.

OLMo, Team, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. 2 olmo 2 furious.

OpenAI. 2022. ChatGPT: Optimizing language models for dialogue.

Pangakis, Nicholas, Samuel Wolken, and Neil Fasching. 2023. Automated annotation with generative ai requires validation. *arXiv preprint arXiv:2306.00176*.

Rozner, Joshua, Leonie Weissweiler, Kyle Mahowald, and Cory Shain. 2025a. Constructions are revealed in word distributions.

Rozner, Joshua, Leonie Weissweiler, and Cory Shain. 2025b. Babylm's first constructions: Causal interventions provide a signal of learning.

Sanguinetti, Manuela and Cristina Bosco. 2015. Parttut: The turin university parallel treebank. In Roberto Basili, Cristina Bosco, Rodolfo Delmonte, Alessandro Moschitti, and Maria Simi, editors, *Harmonization and Development of Resources and Tools for Italian Natural Language Processing within the PARLI Project*, pages 51–69. Springer International Publishing, Cham.

Savelka, Jaromir and Kevin D Ashley. 2023. The unreasonable effectiveness of large language models in zero-shot semantic annotation of legal texts. *Frontiers in Artificial Intelligence*, 6.

Scivetti, Wesley, Tatsuya Aoyama, Ethan Wilcox, and Nathan Schneider. 2025. Unpacking let alone: Human-scale models generalize to a rare construction in form but not meaning.

Silveira, Natalia, Timothy Dozat, Marie-Catherine de Marneffe, Samuel Bowman, Miriam Connor, John Bauer, and Chris Manning. 2014. A gold standard dependency corpus for English. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2897–2904, Reykjavik, Iceland. European Language Resources Association (ELRA).

Tayyar Madabushi, Harish, Laurence Romain, Dagmar Divjak, and Petar Milin. 2020. CxGBERT: BERT meets construction grammar. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4020–4032, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Team, Gemma, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xiaohai Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman, Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi, Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia Paterson, Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petrini, Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, CJ Carey, Cormac Brick, Daniel Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Paparas, Divyashree Shivakumar Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne, Idan Szpektor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wieting, Jonathan Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju yeong Ji, Jyotinder Singh, Kat Black, Kathy Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine, Marina Coelho, Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael Moynihan, Min Ma, Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Nilay Chauhan, Noveen Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Rubenstein, Phil Culliton, Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh Shivanna, Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu, Ryan Mullins, Sammy Jerome, Sara Smoot,

Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti Sheth, Siim Põder, Sijal Bhatnagar, Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi Liu, Trevor Yacovone, Tyler Liechty, Uday Kalra, Utku Evci, Vedant Misra, Vincent Roseberry, Vlad Feinberg, Vlad Kolesnikov, Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein Zhu, Zichuan Wei, Zoltan Egyed, Victor Cotruta, Minh Giang, Phoebe Kirk, Anand Rao, Kat Black, Nabila Babar, Jessica Lo, Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas Gonzalez, Zach Gleicher, Tris Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Jean-Baptiste Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier Bachem, Armand Joulin, Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot. 2025. Gemma 3 technical report.

Torrent, Tiago Timponi, Thomas Hoffmann, Arthur Lorenzi Almeida, and Mark Turner. 2023. *Copilots for Linguists: AI, Constructions, and Frames.* Cambridge University Press.

Touvron, Hugo, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. LLaMA: Open and Efficient Foundation Language Models. Preprint, arXiv 2302.13971.

Tseng, Yu-Hsiang, Cing-Fang Shih, Pin-Er Chen, Hsin-Yu Chou, Mao-Chang Ku, and Shu-Kai Hsieh. 2022. CxLM: A construction and context-aware language model. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 6361–6369, Marseille, France. European Language Resources Association.

Warstadt, Alex, Aaron Mueller, Leshem Choshen, Ethan Wilcox, Chengxu Zhuang, Juan Ciro, Rafael Mosquera, Bhargavi Paranjabe, Adina Williams, Tal Linzen, and Ryan Cotterell. 2023. Findings of the BabyLM challenge: Sample-efficient pretraining on developmentally plausible corpora. In *Proceedings of the BabyLM Challenge at the 27th Conference on Computational Natural Language Learning*, pages 1–34, Singapore. Association for Computational Linguistics.

Warstadt, Alex, Amanpreet Singh, and Samuel R. Bowman. 2019. Neural network acceptability judgments. *Transactions of the Association for Computational Linguistics*, 7:625–641.

Weissweiler, Leonie, Valentin Hofmann, Abdullatif Köksal, and Hinrich Schütze. 2022. The better your syntax, the better your semantics? probing pretrained language models for the English comparative correlative. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10859–10882, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Yu, Danni, Luyang Li, Hang Su, and Matteo Fuoli. 2023. Assessing the potential of llm-assisted annotation for corpus-based pragmatics and discourse analysis: The case of apologies. *International Journal of Corpus Linguistics.*

Zeldes, Amir. 2017. The GUM corpus: Creating multi-layer resources in the classroom. *Language Resources and Evaluation*, 51(3):581–612.

Zeman, Daniel, Martin Popel, Milan Straka, Jan Hajič, Joakim Nivre, Filip Ginter, Juhani Luotolahti, Sampo Pyysalo, Slav Petrov, Martin Potthast, Francis Tyers, Elena Badmaeva, Memduh Gokirmak, Anna Nedoluzhko, Silvie Cinková, Jan Hajič jr., Jaroslava Hlaváčová, Václava Kettnerová, Zdeňka Urešová, Jenna Kanerva, Stina Ojala, Anna Missilä, Christopher D. Manning, Sebastian Schuster, Siva Reddy, Dima Taji, Nizar Habash, Herman Leung, Marie-Catherine de Marneffe, Manuela Sanguinetti, Maria Simi, Hiroshi Kanayama, Valeria de Paiva, Kira Droganova, Héctor Martínez Alonso, Çağrı Çöltekin, Umut Sulubacak, Hans Uszkoreit, Vivien Macketanz, Aljoscha Burchardt, Kim Harris, Katrin Marheinecke, Georg Rehm, Tolga Kayadelen, Mohammed Attia, Ali Elkahky, Zhuoran Yu, Emily Pitler, Saran Lertpradit, Michael Mandl, Jesse Kirchner, Hector Fernandez Alcalde, Jana Strnadová, Esha Banerjee, Ruli Manurung, Antonio Stella, Atsuko Shimada, Sookyoung Kwak, Gustavo Mendonça, Tatiana Lando, Rattima Nitisaroj, and Josie Li. 2017. CoNLL 2017 shared task: Multilingual parsing from raw text to Universal Dependencies. In *Proceedings of the CoNLL 2017 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pages 1–19, Vancouver, Canada. Association for Computational Linguistics.

Zhou, Shijia, Leonie Weissweiler, Taiqi He, Hinrich Schütze, David R. Mortensen, and Lori Levin. 2024. Constructions are so difficult that Even large language models get them right for the wrong reasons. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 3804–3811, Torino, Italia. ELRA and ICCL.

# A    Full details for each prompt

We report in Tables 7 to 24 the details of the prompt, along with the change that it represents from a previous prompt.

# B    Few Shots

In Table 41, we give the five shots from each class given to ChatGPT as examples.

Hybrid Human-LLM Corpus Construction and LLM Evaluation for the Caused-Motion Construction

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Reply with a csv codeblock (wrapped in three backticks), with the headers 'id' and 'label'. label should be either True or False. Label all 50 sentences. |
| Few-Shots | 0 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | Base |
| Shot Strategy | all |

Table 7: Prompt 1

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. |
| Few-Shots | 0 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 1 + repeat sentence with json |
| Shot Strategy | all |

Table 8: Prompt 2

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 0 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 2 + reason |
| Shot Strategy | all |

Table 9: Prompt 3

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 0 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 3 + structured information |
| Shot Strategy | all |

Table 10: Prompt 4

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 0 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 4 + sentence |
| Shot Strategy | all |

Table 11: Prompt 5

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 0 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 4 + cmc string |
| Shot Strategy | all |

Table 12: Prompt 6

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 0 |
| **Sentences** | 50 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 4 + cmc string continuous |
| **Shot Strategy** | all |

Table 13: Prompt 7

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 0 |
| **Sentences** | 50 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 6 + sentence |
| **Shot Strategy** | all |

Table 14: Prompt 8

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 0 |
| **Sentences** | 50 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 7 + sentence |
| **Shot Strategy** | all |

Table 15: Prompt 9

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 10 positive examples: . Here are 10 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 10 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 4 + few shots |
| Shot Strategy | first of each verb and class |

Table 16: Prompt 10

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 10 positive examples: . Here are 10 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 10 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 10 + explanations |
| Shot Strategy | first of each verb and class |

Table 17: Prompt 11

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 50 |
| Sentences | 50 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 11 + all shots |
| Shot Strategy | all |

Table 18: Prompt 12

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 50 |
| Sentences | 10 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 12 + only 10 samples |
| Shot Strategy | all |

Table 19: Prompt 13

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 50 |
| Sentences | 1 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 12 + only 1 sample |
| Shot Strategy | all |

Table 20: Prompt 14

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 50 |
| Sentences | 5 |
| Model | 4o_mini |
| Majority Vote | No |
| Change | 12 + only 5 samples |
| Shot Strategy | all |

Table 21: Prompt 15

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 25 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 12 + only 25 samples |
| **Shot Strategy** | all |

Table 22: Prompt 16

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 1 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 14 + new few-shots |
| **Shot Strategy** | all |

Table 23: Prompt 17

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: . Here are 50 negative examples: . Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 1 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 17 + alternating shots |
| **Shot Strategy** | all_alternating |

Table 24: Prompt 18

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 1 |
| **Model** | 4o_mini |
| **Majority Vote** | No |
| **Change** | 17 + grouped shots |
| **Shot Strategy** | all_grouped |

Table 25: Prompt 19

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 1 |
| **Model** | 4o_mini |
| **Majority Vote** | Yes |
| **Change** | 19 + majority vote |
| **Shot Strategy** | all_grouped |

Table 26: Prompt 20

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 50 |
| Sentences | 1 |
| Model | o3_mini |
| Majority Vote | No |
| Change | 19 on o3-mini |
| Shot Strategy | all_grouped |

Table 27: Prompt 21

| Instruction | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
|---|---|
| Input Format | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| Output Format | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| Few-Shots | 50 |
| Sentences | 100 |
| Model | o3_mini |
| Majority Vote | No |
| Change | 21 + 100 samples |
| Shot Strategy | all_grouped |

Table 28: Prompt 22

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 250 |
| **Model** | o3_mini |
| **Majority Vote** | No |
| **Change** | 21 + 250 samples |
| **Shot Strategy** | all_grouped |

Table 29: Prompt 23

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | o3_mini |
| **Majority Vote** | No |
| **Change** | 21 + 50 samples |
| **Shot Strategy** | all_grouped |

Table 30: Prompt 24

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 25 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 25 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 25 |
| **Sentences** | 50 |
| **Model** | o3_mini |
| **Majority Vote** | No |
| **Change** | 21 + 25 samples |
| **Shot Strategy** | all_grouped |

Table 31: Prompt 25

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | o3_mini |
| **Majority Vote** | Yes |
| **Change** | 24 + majority vote |
| **Shot Strategy** | all_grouped |

Table 32: Prompt 26

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | 4o |
| **Majority Vote** | No |
| **Change** | 24 on 4o |
| **Shot Strategy** | all_grouped |

Table 33: Prompt 27

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | o3_mini |
| **Majority Vote** | No |
| **Change** | 24 - sentence |
| **Shot Strategy** | all_grouped |

Table 34: Prompt 28

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | 4o |
| **Majority Vote** | No |
| **Change** | 27 - sentence |
| **Shot Strategy** | all_grouped |

Table 35: Prompt 29

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | o3_mini |
| **Majority Vote** | No |
| **Change** | 28 - reason |
| **Shot Strategy** | all_grouped |

Table 36: Prompt 30

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. Before you give the label, justify your decision with a reason. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | 4o |
| **Majority Vote** | No |
| **Change** | 29 - reason |
| **Shot Strategy** | all_grouped |

Table 37: Prompt 31

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. |
| **Few-Shots** | 50 |
| **Sentences** | 100 |
| **Model** | o1 |
| **Majority Vote** | No |
| **Change** | 22 on o1 |
| **Shot Strategy** | all_grouped |

Table 38: Prompt 32

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | o1 |
| **Majority Vote** | No |
| **Change** | 32 + 50 samples |
| **Shot Strategy** | all_grouped |

Table 39: Prompt 33

| | |
|---|---|
| **Instruction** | The task is to classify whether the sentences contain instances of the caused-motion construction. The caused-motion construction is a construction where an agent causes an object to move. This motion has to be literal, not metaphorical. Each sentence that you will be given includes a subject, a verb, a direct object, and a prepositional phrase. In the caused motion instances, the verb causes the motion of the direct object, in the direction specified by the prepositional phrase. The action does not need to actually happen, it could be only mentioned or hypothetical or occur in the past or future. |
| **Input Format** | Here are 50 positive examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Here are 50 negative examples: { "sentence": "...", "verb": "...", "direct_object": "...", "preposition": "...", "prepositional_object": "...", "reason": ..., "label": ... }. Classify the following sentences: { "id": "...", "sentence": "..." }. |
| **Output Format** | Respond with a jsonl codeblock (wrapped in three backticks) using double quotes. |
| **Few-Shots** | 50 |
| **Sentences** | 50 |
| **Model** | o1 |
| **Majority Vote** | Yes |
| **Change** | 33 + majority vote |
| **Shot Strategy** | all_grouped |

Table 40: Prompt 34

| Sentence | Verb | Dir Obj | Prep | P-Obj | Lab. | Explanation |
|---|---|---|---|---|---|---|
| I actually giggled myself to tears . | giggle | myself | to | tear | False | This is a negative example because being 'in tears' is a state, not a location, so the subject here didn't move but rather changed state. |
| Nope , they just giggle their microscopic excretions into the air . | giggle | excretion | into | air | True | This is a positive example because the act of giggling is causing the excretions to move |
| I 'll stop it from repeating and fade it into a single background color . | fade | it | into | color | False | This is a negative example which describes the act of fading a color so that it can't be told apart from the background color, which means that nothing moved. |
| Just hover your mouse over it | hover | mouse | over | it | False | This is a negative example because the mouse is hovering over it, but it is not moving, it is staying in place while hovering. |
| Once she was strapped back in he started to hover her out of the room . | hover | she | out of | room | True | This is a positive example because they are moving her out of the room by hovering. |
| They tug him to the ground and start jumping on him and licking his face . | tug | he | to | ground | True | This is a positive example because someone was tugged and that moved him to the ground. |
| I gulped it from the bottle while watching old movies . | gulp | it | from | bottle | True | This is a positive example because what was in the bottle moved from the bottle because the person was drinking it. |
| I would rather take the spoon , I can gulp it in one go . | gulp | it | in | go | False | This is a negative example because one go is not a destination, it specifies the manner of gulping. |
| Cruz was trailing Clinton in basically every poll . | trail | Clinton | in | poll | False | This is a negative example because no movement is happening, the sentence describes the relative position of two politicians in a poll. |
| She began trailing a finger down his chest . | trail | finger | down | chest | True | This is a positive example because she is moving the finger down his chest. |
| He stopped for ten minutes while wheezing himself to death . | wheeze | himself | to | death | False | This is a negative example because death is a state, not a physical location. |
| It is not cute to watch your dog wheeze himself to the floor because he was so excited you picked up his tug of war rope . | wheeze | himself | to | floor | True | This is a positive example because the wheezing is causing the dog to move to the floor. |
| I bounced it off the wall . | bounce | it | off | wall | True | This is a positive example because the ball moved off the wall. |
| We bounce ideas off each other . | bounce | idea | off | other | False | This is a negative example because an idea can't physically move. |
| I was in bed for about a week and thought I was going to shiver myself to death . | shiver | myself | to | death | False | This is a negative example because death is a state, not a destination of a physical movement. |
| She needs to stop darting her eyes to the side every time she says something | dart | eye | to | side | False | This is a negative example because her eyes are rotating but they're not moving. |
| He nervously darted his tongue into her mouth . | dart | tongue | into | mouth | True | This is a positive example because his tongue is moving into her mouth. |
| For some reason every time i overflow the sink in Dalia 's bathroom , the Sheik always comes up to investigate … | overflow | sink | in | bathroom | False | This is a negative example because the sink is not moving. |
| Most importantly the toilet was overflowing water into the pan , almost on constant flush . | overflow | water | into | pan | True | This is a positive example because the water is moving into the pan. |
| You 're trying to wriggle your way out of it now ! | wriggle | way | out of | it | False | This is a negative example because while something is moving, it is not the direct object way. |
| At one point he wriggles himself into position to block a soccer ball with his head while Latin on the street . | wriggle | himself | into | position | True | This is a positive example because he is moving himself into position. |
| I swim laps in the pool . | swim | lap | in | pool | False | This is a negative example because while I am moving, the laps are not moving. |
| My wife lapped me on the scoring track . | lap | I | on | track | False | This is a negative example because I am moving, but my wife is not causing me to move. |
| He will nibble you to death ! | nibble | you | to | death | False | This is a negative example because death is a state, not a location. |
| I eat my Twix by nibbling the chocolate off the sides , then off the top , then eat the caramel and cookie . | nibble | chocolate | off | side | True | This is a positive example because the chocolate is moving off the sides. |
| I aimed at her , and gazed her in her eyes before I successfully hit her face with the snowball . | gaze | she | in | eye | False | This is a negative example because a gaze is not something that can physically move. |
| I choose to be the one that goes hiking with friends into waterfalls , out galloping horses in open fields , and having fun times with my SO . | gallop | horse | in | field | False | This is a negative example because the horses are moving, but they are not moving in the direction of the field, they are already in it. |
| I can confirm that galloping a horse through an open field is amazing . | gallop | horse | through | field | True | This is a positive example because the horse is moving through the field. |
| I scramble them in the hot pan . | scramble | they | in | pan | False | This is a negative example because the eggs are not moving in the direction of the pan, they are already in it. |
| Once it firms a little , scramble it into the rice . | scramble | it | into | rice | True | This is a positive example because the eggs are moving into the rice. |
| To continue with your explanation , we see not only that this man here can afford to encrust rare and obviously expensive jewels onto his box of ' Fruity Pebbles ' brand breakfast cereal , but also that he can afford the ' Family Size ' box . | encrust | jewel | onto | box | True | This is a positive example because the jewels are moving onto the box. |
| I peel paint off walls . | peel | paint | off | wall | True | This is a positive example because the paint is moving off the wall. |
| I peel bananas from the bottom | peel | banana | from | bottom | False | This is a negative example because the banana is not moving, only the peel is, and it is not moving from the bottom. |
| In my defense I was actually very drunk when I plowed my car into that crowd of pedestrians . | plow | car | into | crowd | True | This is a positive example because I caused tha car to move into the crowd. |
| I plow snow in the winter | plow | snow | in | winter | False | This is a negative example because in the winter is a time, not a location. |
| And drag queens cake themselves in makeup . | cake | themselves | in | makeup | False | This is a negative example because the drag queens are not moving. |
| I would cake makeup on my face to hide it . | cake | makeup | on | face | True | This is a positive example because the makeup is moving onto the face. |
| Whereas WWE charred it to a crisp and drowned it in A-1 sauce . | char | it | to | crisp | False | This is a negative example because it is changing state to a crips, not moving. |
| I fermented it in a 3 gallon food grade plastic bucket . | ferment | it | in | bucket | False | This is a negative example because it is staying in the bucket and not moving. |
| When the child collapsed , the mother hurried him to the hospital , where he died . | hurry | he | to | hospital | True | This is a positive example because the child is moving to the hospital. |
| I will take my time or hurry you through a meal , there are no rules against that . | hurry | you | through | meal | False | This is a negative example because the meal here is an action, not a destination |
| I love blackening it in a roasting pan . | blacken | it | in | pan | False | This is a negative example because it is not moving, it is staying in the pan. |
| I rarely use them but my girlfriend is crocheting them into reusable shopping bags … | crochet | they | into | bag | False | This is a negative example because the bags are not moving, they are being made into something else. |
| When " nice guys " change their MO to target " nice girls " the equilibrium will tilt the earth off its axis and hurtle us into space , thus settling this tired argument for all eternity . | hurtle | we | into | space | True | This is a positive example because we are moving into space. |
| Then you drip juice into it and vape . | drip | juice | into | it | True | This is a positive example because the juice is moving into it. |
| As in you literally gnaw it off the bone . | gnaw | it | off | bone | True | This is a positive example because it is moving off the bone |
| I twitch my head to the side . | twitch | head | to | side | True | This is a positive example because the head is moving to the side. |
| He snorted coke off my ass | snort | coke | off | ass | True | This is a positive example because the coke moved off my ass. |
| I ca n't tell if she 's smiling or is she 's about to sneeze the sand off of her nose . | sneeze | sand | off of | nose | True | This is a positive example because the sand moves off her nose. |
| It was like a little rocket that tried to burrow itself into the ground . | burrow | itself | into | ground | True | This is a positive example because the rocket moves into the ground. |

Table 41: Few shots. P-Obj stands for Prepositional Object, Dir Obj for Direct Object.

# Implicit and Indirect: Detecting Face-threatening and Paired Actions in Asynchronous Online Conversations

Henna Paakki, 🔸 Aalto University, Espoo; University of Helsinki, `henna.paakki@helsinki.fi`

Pihla Toivanen, University of Helsinki

Kaisla Kajava, 🔸 Aalto University, Espoo, `kaisla.kajava@aalto.fi`

**Abstract** This paper presents an approach to computationally detecting face-threatening and paired actions in asynchronous online conversations. Action detection has been widely studied for synchronous chats. However, there are fewer models or datasets for asynchronous conversations, and they have not included some of the face-threatening actions central to online conversations involving misbehavior like trolling. We examine asynchronous crisis news related online conversations in Finnish, providing an annotation scheme for identifying central actions used in this conversational context. An important contribution is to include face-threatening actions in the scheme, and training computational classifiers for their detection with improved performance compared to prior work. We illustrate that face-threatening actions are important for analyzing conversations related to crisis news. We show that for computational action detection, it is essential to be able to represent how multiple actions may be performed within one comment, and how ambiguity in the expression of actions often leads to multiple possible label interpretations. Annotating actions using scores helps to reflect these characteristics. We also find that an ensemble of models trained on individual annotators' annotations can best represent multiple potential interpretations of action labels. These are especially relevant for face-threatening actions.

## 1 Introduction

Natural Language Processing (NLP) and Machine Learning (ML) methods are popular for analyzing textual content on social media, e.g. discourse signals (Ferracane et al., 2021; Zhang et al., 2017). These methods combined with examining the structural features of conversation, like comment-response relations, allow rich automated analyses of interaction (Zhang et al., 2018; Sudhahar et al., 2015), crucial for identifying online misbehavior like trolling (Paakki et al., 2024), or aggression (Zhang et al., 2018; Garimella et al., 2018). While computational approaches to detecting online manipulation have focused on revealing content sharing networks (Giglietto et al., 2020) and disinformation campaign content (Shu et al., 2017; Zhou and Zafarani, 2020), fewer have investigated conversational interactions on online forums in more detail. Asynchronous news comment and forum conversations are sites where people often seek to influence others' opinions. The impacts may be widespread, especially for crisis news discussions, as asynchronous discussions are persistent online and thus potentially reach large audiences (Zhang et al., 2018). Thus, this is an important context to investigate.

Recent research has highlighted the importance of investigating actions in conversational online interaction to reveal manipulative behavior like trolling (Paakki et al., 2024). Computational models can help accomplish this at scale, which is why computational methods for detecting actions in this context need more attention. Although actions in synchronous online conversations have been studied extensively (Clark and Popescu-Belis, 2004; Forsyth and Martell, 2007; Fuscone et al., 2020; Stolcke et al., 2000), there are fewer models or datasets for analyzing actions in asynchronous arenas. Moreover, existing annotation schemes or datasets do not include face-threatening actions (e.g., *accusations* and *challenges*) that are important for analyzing trolling and disinformation (Paakki et al., 2024; Bellutta et al., 2021). This type of interaction often strategically utilizes ambiguity in expression (Paakki et al., 2021). In this paper, the concept of *action* refers to what functions a turn has in conversation in relation to other turns. We understand face-threatening actions as having the potential to spoil or threaten the face of their addressee by showing a negative orientation toward them or acting against their wishes (Brown and Levinson, 1987).

Our most important overarching goal is to investigate how to computationally approach face-threatening and paired actions (e.g. *question-answer*) in

| Comments | Actions |
|---|---|
| A: There are two powerful presidential candidates in the US; One has done, already years ago, powerful deeds together with God, such which many Presidents in the States have not dared to do, but we found a brave man respecting the Father's will, Donald Trump. (8 laugh, 2 likes) | *statement* AND *appreciation* |
| **B: A, hallelujah!** (1 laugh emoji, 1 like) | *statement* OR *appreciation* |

Table 1: Extract from Ukraine war discussion, 2022.

asynchronous conversations. To this end, we develop an annotation scheme and computational models for classifying face-threatening and paired actions that are central for this context. This will allow analyzing, at scale, how participants use and respond to common actions. Moreover, to represent how actions unfold in this type of conversation, we deem it important to be able to account for multiple actions per comment and multiple possible interpretations of actions. We argue that this will help to better align computational classification with how actions are performed in these interactions, while also leading to better classification performance. We have two research questions:

- RQ1. Does considering multiple actions per comment better reflect how actions are performed in asynchronous conversation in contrast to a single-label approach?

- RQ2. What approach can best represent all possible interpretations of action labels?

We anticipate that considering multiple actions in comments and multiple possible label interpretations is important for better representing what comments do in asynchronous conversation. For example, as seen in Table 1, comments often have a tendency to include more than one prevalent action, with potential for multiple interpretations, e.g. due to semantic ambiguity (Virtanen et al., 2021; Paakki et al., 2021; Stommel and Koole, 2010; Herring, 1999). B's turn in Table 1, given the context, is likely a sarcastic *statement* but could be interpreted as a genuine *appreciation*. On the other hand, at least two actions overlap in A's turn.

We utilize a Conversation Analysis (CA) driven theoretical framework based on Paakki et al. (2021) to build and refine an annotation scheme for identifying actions (8 actions in total) in asynchronous conversations in Finnish – a low-resource setting. The unit of analysis is one comment in a conversation thread. Annotators rate the likelihood of each action on a 7-point Likert scale. Such an approach conveys how face-threatening

actions like *challenges*, for instance, are mostly not categorically present in a comment, but instead often implicitly expressed and thus better annotated using scalar values. We compare the performances of models assuming only one action *vs.* models allowing multiple actions per comment. We experiment with different approaches to leveraging annotator disagreements to investigate which approach can best predict all possible action labels for each comment.

Our contributions include an annotation framework for face-threatening and paired actions in asynchronous conversations, and action classifiers for Finnish. We showcase annotation disagreement in the case of asynchronous conversation, illustrating how face-threatening actions tend to involve higher levels of disagreement than other actions. We find that actions are frequently performed in an implicit or indirect manner in this context. Also, comments often include more than one action. We show that considering multiple actions per comment and multiple valid interpretations of actions allows higher classification performance, and that an ensemble combining few-shot learning based ML models trained separately with individual annotators' annotations can best represent all possible labels.

This approach is important for computationally analyzing how people perform actions in comments to asynchronous conversation, and how they respond to previous actions. Identification of face-threatening actions – not included in prior computational models – is especially crucial for crisis related conversations, and for detecting misbehavior like trolling online (Paakki et al., 2024). We provide our annotation guidelines and related materials on our GitHub page[1], and models on Huggingface[2].

## 2 Identifying Online Actions

This section introduces earlier work relevant to this paper: we introduce 1.) the theoretical foundations of our work (CA), 2.) prior work on computational action detection, and 3.) research highlighting how tasks like action detection may involve meaningful disagreements.

### 2.1 CA and Asynchronous Conversation

CA has potential for computational operationalization due to its tendency to pay attention to distribution and generalizable features of interaction (Stivers, 2015). CA interpretations of actions arise from what a turn does in a conversation, based on the utterance itself and the next turns – how other turns relate to the utterance

---

[1]Detailed annotation guidelines and models, to the extent that it does not compromise any individual's privacy, are provided on our GitHub page: https://github.com/henniina/Detecting-paired-actions.
[2]https://huggingface.co/Finnish-actions

and interpret its role (Sacks et al., 1974). CA has a robust theoretical foundation for studying the dynamics between paired actions (i.e. adjacency pairs) in turn-by-turn interactions (Schegloff, 2007), as well as face-threatening actions like accusations and their expected responses (Antaki et al., 2008; Dersley and Wootton, 2000). Their characteristics have been well established in CA (Dersley and Wootton, 2000; Koshik, 2003; Turowetz and Maynard, 2010). Action pairs allow us to coordinate interaction through actions and responses that match their expectations (Stivers and Rossano, 2010). They reflect accountability, which is crucial for cooperation (Enfield and Sidnell, 2017).

What differentiates CA's understanding of actions from, e.g., Speech Acts, is that interpretations are based on the analysis of interactions between turns rather than judging the intent behind a turn in conversation (Sacks et al., 1974). Speech Act theory, although much used in computational approaches to actions (or acts), has been criticized for its speaker centric perspective (Linell and Marková, 1993; Savolainen, 2020). It has also been found that turns in asynchronous conversations cannot very well be classified into Speech Act categories (Qadir and Riloff, 2011).

Due to these concerns, we root our approach in CA. Studies have shown it to be well suited for analyzing online interactions (Giles et al., 2015; Meredith and Stokoe, 2014). Digital CA research stresses the need to consider the specific characteristics of different types of online interaction (Virtanen et al., 2021; Meredith, 2017). In asynchronous conversation, a real-time back-and-forth is not expected, and participants often tend to include several actions in one message to accomplish more at one go (Virtanen et al., 2021). Participants may enter and exit the discussion whenever, and adjacency pairs may be disrupted by other messages in between (Herring, 1999). Others have also shown face-to-face or synchronous chats to differ from asynchronous interactions (Taniguchi et al., 2020; Xiao et al., 2020).

Despite the differences between asynchronous conversations and more traditional contexts of CA research (e.g., face-to-face conversation), action pairs are a fruitful CA concept for the analysis of online conversation. This is because people have been shown to treat actions and their norms online quite similarly to what has been established in CA research on face-to-face interaction (Meredith, 2017; Paakki et al., 2021; Salonen et al., 2022). Despite disrupted turns, people tend to utilize sequentially organized actions to maintain coherent interactions online (Stommel and Koole, 2010; Meredith and Stokoe, 2014). Conversational norms dictate that an action that expects an answer should be answered in accordance to its normative expectation (e.g. a question requires an (informative) answer) (Sacks et al., 1974; Schegloff, 2007; Enfield et al., 2010). Unexpected an-

swers are seen to signal a problem, for instance a misunderstanding, a gap in presumed shared knowledge, or reluctance (Clark and Schaefer, 1987; Pomerantz, 1984). They can also be used to purposefully direct and disrupt a conversation (Paakki et al., 2021).

## 2.2 Computational Approaches

Most prior research on computational action detection relates to customer chat bots (Casanueva et al., 2020; Ghosh and Ghosh, 2021), telephone conversations (Godfrey et al., 1992; Fuscone et al., 2020), recorded face-to-face dialogue (Clark and Popescu-Belis, 2004) and synchronous chats (Forsyth and Martell, 2007; Moldovan et al., 2011). These, however, represent a context different from casual, anonymous, and asynchronous online conversation (Herring, 1999), for which there are fewer resources, and no large annotated datasets like the Switchboard corpus (Jurafsky, 1997).

Asynchronous conversations have been researched increasingly in recent years. Annotation schemes, for example, have been developed for different asynchronous conversational contexts (Duran and Battle, 2018; Herring et al., 2005; Jeng et al., 2017; Kim et al., 2010; Savolainen, 2020; Stivers, 2015; Wang et al., 2014). Some earlier work on computational classification has focused on email interaction (Carvalho and Cohen, 2005; Cohen et al., 2004; Taniguchi et al., 2020). Others have looked into online forums, newsgroups and question answering sites (Bhatia et al., 2014; Joty and Hoque, 2016; Kim et al., 2010; Feng et al., 2006; Fortuna et al., 2007). Prior approaches have most often depended on Dialogue Act (DA) driven classification (Bhatia et al., 2014; Carvalho and Cohen, 2005; Cohen et al., 2004; Joty and Hoque, 2016; Kim et al., 2010; Taniguchi et al., 2020), based on Speech Act Theory (Austin, 1975), and researchers have developed both message-level and sentence-level classification of DAs. Joty and Mohiuddin (2018) studied the effect of modeling conversational dependencies on DA classification. Taniguchi et al. (2020) trained both sentence-level and message-level DA detection models for emails using BERT (Devlin et al., 2019). A large dataset using asynchronous Reddit conversations was published by Zhang et al. (2017), annotating and classifying discourse relations based on earlier work on Rhetorical Structure Theory (Mann and Thompson, 1988) and online DAs (Feng et al., 2006; Fortuna et al., 2007).

Existing resources for asynchronous conversation, however, mostly do not include actions like *accusations* or *challenges* relevant to analyzing manipulative online behavior (Bellutta et al., 2021; Paakki et al., 2021, 2024). Though action pairs have been investigated in question-answer modeling (of DAs) (Joty and Mohiuddin, 2018; Taniguchi et al., 2020), not much attention has been paid to face-threatening actions. The only

computational approaches we are aware of that have included face-threatening actions relevant to our interests are Paakki et al. (2024) (*challenge* and *accusation*), Bracewell et al. (2012; 2013) (*challenge credibility*, *disrespect*, *relationship conflict* and *task conflict*), Feng et al. 2006 (*criticize*), Zhang et al. (2017) (*negative reaction*) and Zakharov et al. (2021) (e.g. *complaint*). Bracewell et al. (2012) concentrate on acts reflecting psychological study of power and leadership, Zhang et al. (2017) on discourse relations, and Zakharov et al. (2021) on a wide selection of discursive categories based on Bakhtinian theories of discourse (e.g. related to tone and style) (Bakhtin, 1981). The listed studies, however, were based on different theoretical premises and thus did not concentrate on paired actions. Paakki et al. (2024) used a set of 10 paired actions, but they utilized 0-shot Natural Language Inference (NLI) for classification instead of providing a model trained specifically for this task.

## 2.3 Actions and Disagreement

We consider it relevant to be able to reflect the multiple possible interpretations of action labels in our data. Ambiguity is a natural part of human interaction, and it is often used strategically in manipulative interaction (Paakki et al., 2021). Actions, as well, are sometimes expressed in an ambiguous manner, leading to different interpretations (Thomas, 1995). Similarly, for many NLP tasks there is no single ground truth (Jiang and de Marneffe, 2022; Plank, 2022; Uma et al., 2022; Yang, 2021), due to uncertainty in text meaning. This leads to different annotator interpretations of label distribution, constituting meaningful systematic disagreement (Jiang and de Marneffe, 2022; Nie et al., 2020). Including annotator disagreement into models has also been shown to improve model performance on some NLP tasks (Passonneau et al., 2012; Plank, 2022). Models leveraging annotator disagreements are thus needed for better representing actions in online interaction.

Most existing action detection models rely on one ground truth (e.g. Zhang et al. (2017); Joty and Mohiuddin (2018)). An exception is Ferracane et al.'s (2021) model, which sought to predict all valid interpretations of actions in their data. Another study by Taniguchi et al. (2020) predicted both sentence-level and message-level interpretations of actions, which provides insight into how readings of parts versus the whole message contributed to interpretations of labels – though not further drawing from annotator disagreements. However, the former used live congressional hearings and the latter emails as data. These approaches notably differ from our context of asynchronous forum conversations, which involves ambiguous use of actions in medium length texts (in contrast to the formerly mentioned types of data), with frequent use of face-threatening actions (e.g. *accusations*).

| Topics | #Comments | #Annotations |
|---|---|---|
| Covid: 675 | 1,204 | 1,204 (single annotations) |
| War: 529 | | 3,612 (multiple annotations) |

Table 2: Dataset description. An annotation refers to one set of 8 scores, one score per action.

## 3 Data

To answer our RQs, we collected asynchronous conversations under news regarding the COVID-19 Pandemic and Ukraine war posted on Facebook. These were of interest, as crises like these have societal, economic and environmental force, and drive forward change and renewal due to the need for novel courses of action; they also create uncertainty, making online discourses vulnerable to trolling, manipulation and disinformation (Di Mascio et al., 2021). Our data comes from public Facebook (FB) pages owned by Finnish news media: Yleisradio (Yle) and Helsingin Sanomat (HS).[3] Our dataset is described in Table 2. To conduct our experiments, we manually annotated news comments using a digital CA based framework, informed by earlier research (Clark and Schaefer, 1989; Herring et al., 2005; Paakki et al., 2021; Stivers, 2013) and data-driven insights (see section 4).

We used Facepager[4] (v.4.5.3) (Jünger and Keyling, 2019) (MIT License) to scrape FB posts in the two pages' feed, and their threaded comments, between 1 Dec. 2019– 10 Feb. 2023. All posts included a news title, description and link to a news article. To select a subset of random conversations for manual annotation, we shuffled the scraped data per news posts, keeping all comments part of the same conversation together. In this study, we had three annotators (the authors). All annotators are native Finnish speakers, living in Finland, and have training in linguistics and/or computational linguistics and data science[5]. We divided the data into three parts, one for each annotator. Each annotator manually selected the first 400 crisis-related comments from their sample for annotation. This meant that we had to manually read the titles and descriptions of 100-150 news posts to find conversations on COVID or the Ukraine war.

We annotated comments (and their replies) in comment section threads if the conversations fit our inclusion criteria. They had to be related to Ukraine war or COVID-19, commenting allowed, and including at least

---

[3]These are among the most followed news outlets in Finland, Yle being the national public broadcasting company, and Helsingin Sanomat Finland's largest subscription newspaper.

[4]https://github.com/strohne/Facepager

[5]Two of us have MA level training in linguistics/computational linguistics and further doctoral training in computer science and computational linguistics. One of us has MA level training in data science, and further doctoral training in digital humanities and social sciences.

one comment with two or more replies, as we were interested in conversational interaction. While annotating, we referred to the original FB page for better readability and to validate the data.[6] We excluded any comments we had already seen during annotation scheme development (see section 4). To achieve greater variation in comments, we included a maximum of 30 comments related to the same piece of news. The number of comments to a news post (when comments were allowed) ranged from 14 to 684, with a mean of 90.56 and median 84.5, and comment mean length being 151.4 characters, median length 105.0.

We annotated all comments and replies following our refined annotation guidelines (see section 4). To compare the performances of models trained on single annotations against models leveraging multiple annotations (and thus annotation disagreements), we used two different versions of the data: data with **single annotations**, where each comment has only one annotation, and data with **multiple annotations**, where all comments have three annotations, one per annotator. This enabled us to investigate how to best represent multiple possible interpretations of actions in our data (RQ2).

## 4 Annotation Methods

We based our initial annotation scheme on earlier research regarding which actions emerge as relevant in conversations involving possible online trolling (Paakki et al., 2021). Speech Act, DA and CA-related online research alike have emphasized that for different contexts or events different actions matter. For example, crisis (or emergency) related online discussions are more geared toward orders and suggestions than celebrity discussions (Laurenti et al., 2022; Zhang et al., 2011). For conversations including potential trolling or manipulation, different forms of assertive and face-threatening actions are relevant: *accusations*, *challenges*, and *statements* (Paakki et al., 2021). Also, investigating actions that form pairs is important in this context, e.g., *accusation-denial* (Paakki et al., 2021). Paired actions are well-established in CA and Computer-Mediated Communication (CMC) literature (Clark and Schaefer, 1989; Enfield et al., 2010; Schegloff, 2007; Stivers, 2015).

The final scheme (Table 3) was formulated based on data-driven empirical insights during our incremental development of the scheme and annotation guidelines, to ensure that it represents the most central rhetorical and interactive functions of comments for our context of investigation. We include both responsive actions and ones initiating a paired action, which expect specific responses, marked with an asterisk in Table 3

---

[6]Due to the restrictions of Facepager, we had to retrieve some missing comments manually.

| Action | Description | Expected response |
|---|---|---|
| Question* *(initiating)* | Asks someone for information (request for information) | statement |
| Request* *(initiating)* | Requests, proposes or tells someone to perform an action. | denial, acceptance |
| Challenge* *(initiating)* | Refutes someone's epistemic claim or authority. | acceptance, denial, statement |
| Accusation* *(initiating)* | Points out a reprehensible act performed by someone. | acceptance, denial |
| Statement *(initiating/ responding)* | Asserts an opinion, information, wish, neutral or negative evaluation, or answer (to question). | |
| Appreciation *(initiating/ responding)* | Positive evaluation or comment about an actor, event or object. | |
| Acceptance *(responding)* | Agrees, or accepts a request, statement or challenge, or admits an accusation. | |
| Denial *(responding)* | Rejects or denies an action. | |

Table 3: Annotation scheme (*=expects response).

(based on Paakki et al. (2021) and above-listed CA literature). These are important for analyzing the dynamics of user-to-user interaction in conversations with potential trolling, e.g., whether actions expecting a specific type of response are responded to in a normatively expected manner.

There is a large number of possible actions, studied in CA and CMC (Schegloff, 2007). We aimed at a simplified annotation scheme, because very fine-grained tag sets like DAMSL (Allen and Core, 1997) can suffer from sparseness and complexity, reducing annotator agreement (Savy, 2010). We wished to limit actions only to ones observed in our data. This risks oversimplification of our theoretical framework (Stivers, 2015), so we relied on theoretical support and data-driven insights. Originally we considered 15 actions (*rejection*, *admission*, *announcement*, *answer to question*, *evaluation*, *proposal* in addition to ones in Table 3), but finally reduced them to 8 actions. Merging of categories was based on actions having similar functions (e.g., *proposals* and *requests*), and difficulty faced in distinguishing them from each other (e.g., *denial*, *rejection*). Moreover, different tasks often require different annotation schemes de-

| #Actions | #C(s) | | #C(m) | | | | | |
|---|---|---|---|---|---|---|---|---|
| | A1-3 | % | A1 | % | A2 | % | A3 | % |
| 0 | 17 | 1.3% | 56 | 4.7% | 1 | 0.08% | 4 | 0.3% |
| 1 | 722 | 60% | 885 | 73.5% | 477 | 39.6% | 676 | 56.1% |
| 2 | 360 | 30% | 238 | 19.8% | 456 | 37.9% | 417 | 34.6% |
| 3 | 95 | 7.9% | 23 | 1.9% | 217 | 18% | 94 | 7.8% |
| 4 | 10 | 0.8% | 2 | 0.2% | 50 | 4.2% | 12 | 1% |
| > 4 | 1 | 0.08% | 0 | 0% | 3 | 0.3% | 1 | 0.08% |

Table 4: Number of actions labeled per comment (C) in single (s) and multiple (m) annotations (A1=Annotator1, A2=Annotator2, A3=Annotator3).

| Action | Count(s) | Count(m) | % | $r_{WG}$ |
|---|---|---|---|---|
| Question | 243 | 276 | 0.94 | 0.88 |
| Request | 130 | 186 | 0.93 | 0.90 |
| Statement | 857 | 999 | 0.59 | 0.58 |
| Challenge | 138 | 381 | 0.91 | 0.65 |
| Accusation | 155 | 305 | 0.86 | 0.74 |
| Appreciation | 36 | 66 | 0.97 | 0.94 |
| Acceptance | 113 | 171 | 0.97 | 0.89 |
| Denial | 97 | 151 | 0.95 | 0.87 |

Table 5: Distribution of annotated actions in our data (s = single annotations, m = multiple annotations), and inter-annotator agreements (% agreement and $r_{WG}$) with a separate test set (N=100).

pending on forum type: e.g. question-answer forums (Savolainen, 2020; Jeng et al., 2017) and trolling conversations (Paakki et al., 2021) require focus on different classes, so empirical insights were important here.

We aimed to avoid under-specification and unnecessary disagreement by iteratively developing our annotation scheme and annotation guidelines. We developed our scheme before data scraping: annotating, negotiating, and analyzing practice data based on our guidelines at each step. The development involved 15 iterations, during which we honed the guidelines and adapted the scheme in a data-driven manner. The 4 first iterations started with single-label annotation, which we felt did not align with the empirical phenomenon at hand; we found score-based multilabel annotation much more fitting. The data used in the development was manually selected as screen captures and links from the first 300 crisis news posts in Yle's and HS's FB pages we found, between August–December 2022. This data was not included in the final annotated datasets. During annotation, we read the conversations in their original FB context. We stopped when we had reached sufficient agreement (upper boundary of medium agreement or strong agreement) in annotation, and a framework that we felt had a meaningful separation between classes and a scheme that was applicable to our data. Using the finalized scheme, we annotated so far unseen data selected for manual annotation (see section 3).[7] See label distributions in Table 5, and action-specific distributions in Appendix C.

Our task required expertise, training and familiarity with the theoretical premises of (digital) CA, including what constitutes each action, so expert annotation was chosen, meaning that the annotators had linguistics, computational linguistics and social science expertise. Crowdsourcing was not used as non-expert annotation involves concerns related to reliability, misinterpretation or misuse of labels (Duran et al., 2022), shown to be ineffective when analysis requires consideration of context, in-depth reading and domain expertise (Eickhoff, 2018; Rezapour et al., 2020).

We chose a message-level approach – comment

as the unit of analysis – as it has been found useful when computationally identifying question-answer pairs (Taniguchi et al., 2020). Since we focus on face-threatening actions and action pairs, we found this a suitable option instead of sentence-level detection. Research on Speech Acts and CA in offline as well as digital contexts has highlighted that messages may often perform several actions in one message to accomplish more at one go, e.g. to respond to some action and to initiate a new one (Goffman, 1974; Levinson, 2013; Stommel and Koole, 2010; Virtanen et al., 2021). Although there are many differences in how turn-taking unfolds in asynchronous conversation in contrast to face-to-face, like the fact that turns cannot be interrupted, people often treat comments as turns in online conversation (Meredith, 2019; Virtanen and Käänta, 2018), see also section 6.1. However, actions are not always marked with clear boundaries (e.g., line breaks) but might overlap, as in Table 9. CA researchers have emphasized that for interpreting what action a turn commits, a holistic assessment of the speaker's conduct or purposive, goal oriented behavior is needed; no single feature in a turn necessarily signals a specific action (Enfield and Sidnell, 2017; Rossi, 2018). Many factors might influence the interpretation of actions: e.g., the relevance of an action might influence interpretations of expected responses (Stivers and Rossano, 2010). In our data, participants tended to orient to comments as having some main pair-initiating action, even when the comment included several actions, as seen in Table 9. These considerations supported our choice of message-level detection.

We decided to use 7-point Likert scale scores (0: action not present – 3: maybe or partly present – 6: action very strongly present) for annotation, inspired by previous work on annotation disagreement by e.g. Peterson et al. (2019) and work that has illustrated the usefulness of Likert scores in annotator ratings (Barnhurst and Mutz, 1997). See example in Table 9. This

---

[7]We provide our detailed annotation guidelines on our GitHub: https://github.com/henniina/Detecting-paired-actions.

| Metric | Annotator | Action | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | question | request | statement | accusation | challenge | acceptance | denial | appreciation |
| Average | Annotator1 | 5.26 | 4.33 | 5.01 | 4.04 | 3.64 | 4.54 | 4.61 | 3.80 |
| | Annotator2 | 4.88 | 4.27 | 4.88 | 3.70 | 3.60 | 4.06 | 4.03 | 3.31 |
| | Annotator3 | 4.74 | 3.96 | 4.87 | 4.04 | 3.59 | 4.10 | 3.50 | 4.20 |
| | Overall | 4.96 | 4.19 | 4.92 | 3.93 | 3.61 | 4.23 | 4.05 | 3.77 |
| Correlation | Annotator1&2 | 0.86 | 0.64 | 0.50 | 0.38 | 0.22 | 0.66 | 0.38 | 0.41 |
| | Annotator2&3 | 0.92 | 0.75 | 0.60 | 0.61 | 0.44 | 0.73 | 0.60 | 0.62 |
| | Annotator3&1 | 0.87 | 0.74 | 0.56 | 0.42 | 0.27 | 0.72 | 0.46 | 0.45 |

Table 6: Label score averages by annotator and overall, and score correlations per action between annotators.

| Annotator | Ground truth | Action | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | question | request | statement | accusation | challenge | acceptance | denial | appreciation |
| Annotator1 | conservative | 0.95 | 0.88 | 0.78 | 0.61 | 0.51 | 0.86 | 0.72 | 0.80 |
| Annotator2 | conservative | 0.98 | 0.93 | 0.93 | 0.95 | 0.96 | 0.95 | 0.94 | 0.86 |
| Annotator3 | conservative | 0.98 | 0.89 | 0.85 | 0.87 | 0.72 | 0.88 | 0.86 | 0.86 |
| Annotator1 | relaxed | 0.95 | 0.91 | 0.85 | 0.72 | 0.65 | 0.92 | 0.78 | 0.80 |
| Annotator2 | relaxed | 0.98 | 0.91 | 0.92 | 0.86 | 0.72 | 0.91 | 0.86 | 0.86 |
| Annotator3 | relaxed | 0.99 | 0.95 | 0.94 | 0.93 | 0.93 | 0.94 | 0.95 | 0.90 |

Table 7: Macro-F1 scores for annotators comparing annotator specific annotations to ground truth labels.

was due to many comments including more than one prevalent action, and we judged that forcing annotators to label only one would reduce annotation quality (see Table 4). We saw that the signal for an action could be present at different scales of strength, which Likert scores could reflect. Confidence scores, for example, have been found useful for achieving improved inter-rater agreement, and highlighting difficult cases (Weber et al., 2018; Troiano et al., 2021). Scores were coded for each action separately per comment, allowing multiple interpretations of label distribution. We set score 0 as designating the absence of an action, as annotators found this most intuitive, meaning 6 is the highest score.

## 4.1 Annotation Quality and Agreement

We measured our annotation quality comparing the observed and expected variances of scores between annotators. For this purpose, $r_{WG}$ score has been deemed helpful for evaluating score-based agreement within a group (Castro, 2002): $r_{WG} = 1 - $(Observed Group Variance/ Expected Random Variance) (Lindell and Brandt, 1999). We used $r_{WG}$ score for each individual action (Table 5), and $R^*_{WG(J)}$ score for overall agreement (Lindell and Brandt, 1999; O'Neill, 2017). The first score is for single-item scale and the second for multi-item scale, measuring the variance between annotations when random variance is eliminated. $R^*_{WG(J)}$ score was 0.72. Agreement strength boundaries are defined for the $r_{WG}$ measurement family: our $R^*_{WG(J)}$ score is over the lower

bound of strong agreement (0.71 − 0.90) (O'Neill, 2017).

For more detailed insights, we evaluated the annotators' annotation performances against ground truth labels (see section 5.1), in Table 7. These demonstrate some differences among annotators, some achieving higher performances on specific actions than others. We also investigated score means across actions, and annotation score correlations between pairs of annotators (see Table 6). Overall, correlations between Annotator2 and 3 tended to be higher.

## 5 Methods

In this section, we will describe how we set our ground truth labels, which ML models we used in our experiments, how we evaluated our models, and how we designed our experiments to answer our RQs.

## 5.1 Ground Truth

For computational modeling, we needed to define a meaningful ground truth against which to compare our models. Thresholding has been found useful in earlier studies when using confidence scores in annotation (Xu et al., 2017), and was relevant here due to our use of scores. We set a threshold $\theta = 3$ based on theoretical insights, and our empirical interests. Although our goal was to represent annotators' confidence in labels as well as their assessment of label strength as closely as possible, scores of 1-2 were rarer (see Appendix C), and based on theoretical considerations (Stivers and

Rossano, 2010) and empirical insights (see section 6.1), we considered scores 3-6 as representing a medium to strong signal that online discussion participants would likely consider relevant, requiring a response if the action usually involves normative expectations. Thus, for analyzing responding, we consider this a useful threshold. For a deeper understanding of how thresholding might affect classification, we conducted an additional (suggestive) threshold test: see Appendix B.

To build a set of ground truth labels, we mapped annotation scores back to binary labels (1: action present, 0: action not present). For single annotations, the label was 1 if the annotation score was $\geq \theta$. For multiple annotations, we use *conservative* (if even one annotator had given a score $\geq \theta$, the label was 1), and *relaxed ground truth labels* (if at least two annotators had given a score $\geq \theta$, the label was 1).

## 5.2 ML Models

Recent action classifiers rely on sentence transformers, other neural networks (Ghosh and Ghosh, 2021; Taniguchi et al., 2020; Joty and Mohiuddin, 2018), or few-shot learning (Casanueva et al., 2020). Studies emphasize the relevance of linguistic features (e.g., lexical and collocations) (Stolcke et al., 2000; Ferracane et al., 2021; Zakharov et al., 2021), and the use of pretrained word embeddings (Joty and Mohiuddin, 2018). Language Models (LM) like BERT have provided promising results in recent research (Taniguchi et al., 2020), and Generative LMs are receiving a lot of attention. Transfer learning techniques, then again, have proven useful for low-resource settings (Pamungkas et al., 2020). We thus chose to compare the performances of 1.) SVM, 2.) FinBERT, 3.) SetFit (with word embedding FinBERT), 4.) SetFit (sentence embedding FinBERT), 5.) Llama2 0-shot, 6.) Llama2 (finetuned), and 7.) FinGPT 0-shot.

As using separate classifiers for sub-tasks has proven effective in recent related work (Ferracane et al., 2021; Zakharov et al., 2021), we decided to train a separate classifier for each action following Ferracane et al. (2021), instead of training one multi-label ML classifier (Wu et al., 2019; Xu et al., 2017). The models predict whether a text contains an action or not.

To allow a comparison of our results to earlier research, we trained SVM models – often used in earlier action detection – similarly to previous models used for asynchronous data (Cohen et al., 2004; Zhang et al., 2017). We used 1-grams and TF-IDF for feature extraction, with SVD for dimensionality reduction (10 dimensions), balanced class weighting, and Grid Search for hyperparameter optimization[8] with sci-kit learn (Pedregosa et al., 2011). We applied preprocessing (tokenization and lemmatization with spaCy, following

Haverinen et al., 2014), leaving in stop words as their removal negatively impacted performance. We were interested in whether weaker models, like SVM, would also show improved action classification performance if allowing multiple actions labeled per comment.

We used FinBERT, as BERT-derivatives haved proved their capacity in text classification (Devlin et al., 2019; Arabadzhieva-Kalcheva and Kovachev, 2022), also for asynchronous CMC (Davidson et al., 2020; Guo and Sarker, 2023), and DA classification (Taniguchi et al., 2020). We finetuned the Finnish pretrained FinBERT (bert-base-finnish-cased-v1)[9] (Virtanen et al., 2019) based on cased word embeddings (henceforth FinBERT1). Our data is highly imbalanced, so we used class weighting.

Besides standard fine-tuning, we used Setfit, a transfer learning technique that can also be initiated with Finnish LMs. It is based on fine-tuning a pretrained model with sentence pairs, then training a classifier based on finetuned embeddings (Tunstall et al., 2022). SetFit creates more variety to the minority category training samples with sentence pairing, generating in total $k(k-1)/2$ different pairs from training data, $k$ being the size of the training set. We initialized SetFit both with FinBERT1 (Virtanen et al., 2019) and FinBERT based on cased sentence embeddings (sbert-cased-finnish-paraphrase)[10] (henceforth FinBERT2) as underlying models. Both thus used a pretrained model to provide embeddings for comments.

We used Optuna (Akiba et al., 2019) for hyperparameter optimization, including learning rate, number of epochs, and batch size.

We tested the performance of Llama2 (Touvron et al., 2023) and FinGPT (Finnish GPT-3 large)[11] (Luukkonen et al., 2023) due to the recent popularity of Generative LMs. We used an uncensored version of Llama2 (llama2-7b-chat-uncensored)[12], because it is accessible and free to use, and as many comments included e.g. offensive language use. First, we used it without finetuning, i.e., as a 0-shot model. We prompted the models to decide whether each comment included a specific action or not. On our GitHub page, we provide a summary of prompts and how generated output was interpreted, matching key words in the generated output to extract predictions. With Llama2, we used the PyPi translators library (v5.9.2, GPLv3 license) to translate our data to English, because Llama models have shown to work much better for English (Wendler et al., 2024), and as Finnish is a low-resource language for Llama. We sought to reduce problems resulting from

---

[8]Optimizing kernel, regularization parameter C, and degree.

[9]https://huggingface.co/TurkuNLP/bert-base-finnish-cased-v1
[10]https://huggingface.co/TurkuNLP/sbert-cased-finnish-paraphrase
[11]https://turkunlp.org/gpt3-finnish
[12]https://huggingface.co/georgesung/llama2_7b_chat_uncensored

| Model ID | Single label | Multi-label | Single annotation | Multiple annotation |
|---|---|---|---|---|
| SVM-1-act | ✓ | | ✓ | |
| SVM-MS | | ✓ | ✓ | |
| SVM-Avg | | ✓ | | ✓ |
| FinBERT1-1act | ✓ | | ✓ | |
| FinBERT1-MS | | ✓ | ✓ | |
| FinBERT1-Avg | | ✓ | | ✓ |
| SetFit-FinBERT1-1-act | ✓ | | ✓ | |
| SetFit-FinBERT1-MS | | ✓ | ✓ | |
| SetFit-FinBERT1-Avg | | ✓ | | ✓ |
| SetFit-FinBERT1-PNC | | ✓ | | ✓ |
| SetFit-FinBERT1-A1 | | ✓ | * | |
| SetFit-FinBERT1-A2 | | ✓ | * | |
| SetFit-FinBERT1-A3 | | ✓ | * | |
| SetFit-FinBERT2-1-act | ✓ | | ✓ | |
| SetFit-FinBERT2-MS | | ✓ | ✓ | |
| SetFit-FinBERT2-Avg | | ✓ | | ✓ |
| FinGPT-0shot-1-act | ✓ | | ✓ | |
| FinGPT-0shot-MS | | ✓ | ✓ | |
| FinGPT-0shot-Avg | | ✓ | | ✓ |
| Llama2-0shot-1-act | ✓ | | ✓ | |
| Llama2-0shot-MS | | ✓ | ✓ | |
| Llama2-0shot-Avg | | ✓ | | ✓ |
| Llama2-finetuned-1-act | ✓ | | ✓ | |
| Llama2-finetuned-MS | | ✓ | ✓ | |
| Llama2-finetuned-Avg | | ✓ | | ✓ |

Table 8: Model IDs with related configurations (*=these utilize the full annotated dataset, but only the annotations made by one single annotator).

hallucination by using the outlines library[13] (Willard and Louf, 2023), which forces the model to choose between selected classes. We also finetuned a quantized Llama2 for our task using Adapters (Poth et al., 2023), the QLoRA approach[14] by Dettmers et al. (2024) and the bitsandbytes library[15] (MIT license), similarly to Dettmers et al. (2024). For Finnish, we used FinGPT without further finetuning (0-shot).

## 5.3 Data Splits and Evaluation

We divided our data into train, validation (for hyperparameter tuning) and test sets using sci-kit learn train test split twice, with respective set sizes 60%, 20% and 20%. Input data included individual comments for all models and labels.

Model evaluation included accuracy and Macro-F1 – we report F1 due to class imbalances. We report Jaccard coefficient scores when evaluating which computational approach might best predict all possible labels.

## 5.4 Experiments

**To answer RQ1**, we compared all models using 3 configurations to compare the effects of considering one vs. multiple actions: a.) a *single action (1-act)*, b.) *multilabel single-annotation (MS)*, and c.) *averaged (Avg)* approach.

See a depiction of the configurations and model IDs in Table 8. Although model simplicity is often preferred, we wished to test for possible gains in performance with increased model complexity.

The *1-act* configuration used single annotations, setting for each comment the action with highest score among all labels as positive if score $\geq \theta$, other actions as negative. In case two actions (or more) had equal scores, we assigned a positive label randomly. In our view, this corresponded to annotator decision-making, as in annotation we often had to randomly choose the primary action if having to decide only one most important action between actions with the same signal strength. We compared the *MS* and *Avg* configurations to the *1-act* model to test whether allowing more than one action will statistically improve performance. Other configurations allow multiple actions labeled.

The *MS* configuration used single annotations, allowing multiple actions. *Avg* utilized multiple annotations to decide an average annotation score for label decision, $score_{action} = (a_1 + a_2 + a_3)/3$, with $a$ a reference to $AnnotationScore_i$ for the action. It was included as it has been one of the most popular approaches to dealing with disagreements.

For statistical tests, we used the Nemenyi test, utilizing the scikit-posthocs implementation for python (Terpilowski, 2019), recommended for comparing classifier performances (Derrac et al., 2011).

**In relation to RQ2**, using SetFit-FinBERT1, which had achieved the highest number of best performances in relation to RQ1, we compared three configurations for leveraging annotator disagreements to find the most suitable approach for representing multiple interpretations of actions: a.) *averaging (Avg)* (following Uma et al. (2022), b.) *positive/negative/complicated (PNC)* (Jiang and de Marneffe, 2022), and c.) *individual annotator (Annotator1-3)* approach (Ferracane et al., 2021).

*Avg* and *PNC* were included as the approaches have often been used with multiple annotations (Uma et al., 2022; Jiang and de Marneffe, 2022). With *PNC*, per Jiang and de Marneffe (2022), we divided comments into three classes by counting an average cross-entropy score between annotations: two classes where annotators agreed they belong (here positive/negative), and a class where there was significant disagreement (here complicated). We included a comment in the "complicated" class if $C_i > \frac{\sum_i(C)}{N} + 2SD(C)$, $C$ a reference to average cross-entropy. This decision boundary differs from Jiang and de Marneffe (2022), but as they had a hundred crowdsourced annotations per example, and we only three, we could not use a similar method. We judged that an outlier boundary of $mean(C) + 2SD$ would be strict and similar enough. We treated the task as a 3-way classification problem. Individual annota-

tor models have been fruitful for representing multiple interpretations, by predicting all possible annotations (Davani et al., 2022; Ferracane et al., 2021). These were trained using each individual annotator's annotations only: Annotator1 (*A1*), Annotator2 (*A2*), and Annotator3 (*A3*). All configurations were tested against conservative and relaxed ground truth.

We wished to find a model that could predict all possible annotations. Thus, we compared how different model ensembles fared in predicting all possible labels, using Jaccard Coefficient to compare each ensemble's predicted set of labels to a set of all possible annotations by annotators for each comment: $JaccardSimilarityJ(A, B) = \frac{|A \cap B|}{|A \cup B|}$.

Finally, we compared differences between actions by investigating average label scores per action, and annotation and class correlations by using scipy.stats package with Kendall rank correlation (Schober et al., 2018), visualized with heatmaps using the seaborn library. Kendall correlation was used as normal distribution could not be assumed and as it is often considered more robust in contrast to Spearman's (Schober et al., 2018).

# 6 Results

Next, we will discuss our results to answer our RQs: RQ1.) whether a multi-label approach better aligns with how actions are performed in our empirical context, and RQ2.) what approach can best represent the multiple valid interpretations relevant to how people express actions in asynchronous conversations.

## 6.1 RQ1: One Action or Multiple Actions

We asked whether a multi-label approach would better align with how actions are performed in our asynchronous data in contrast to a single-label approach. This was addressed by both 1.) comparing these approaches during the annotation scheme development, and negotiating and examining annotation interpretations, and 2.) testing whether adopting a single-label approach in contrast to a multi-label approach would significantly affect classification performances.

First, turns in our asynchronous conversation data tended to utilize actions somewhat differently as compared to face-to-face conversation, and synchronous online conversation.

In Table 9, we see a conversation extract from our data. Examining comments 2-4, we can see that participants in our online conversation data utilized paired actions to build coherence across turns. In (2), user A asks a question that can be interpreted as *challenging*. In (3), B does not provide information in response to the *challenging question*, but can be seen to *challenge* its

grounds. Annotators 2 and 3 interpret the response, as well as comment (2), as containing a *challenge*. The extract illustrates how turns may often perform multiple actions at the same time at different scales of strength – they might contain an action responding to a previous turn (completing an action pair), and other actions, potentially also initiating new action pairs. Comment (4) is a case in point. Also, actions are not neatly organized into separate sentences or paragraphs. For example, in comment (2), the interpretation of the comment containing a *challenge* arguably cannot be reduced to a single sentence, but is derived from a more holistic reading. This was common for comments in our data as users often formed actions that overlapped. For example, this could be seen in our data when one action was used as a vehicle for another action, one part of the comment (e.g., a sentence) performed multiple actions (e.g., stating and accusing) (cf. double-barreled actions or composite actions: Rossi, 2018; Levinson, 2013; Schegloff, 2007: pp. 73–78), or the whole comment included multiple (overlapping) actions.

Users often tended to treat a previous turn as having some main action(s), and they did not necessarily respond to all actions present in a previous turn. This phenomenon has also been described in CA, though in face-to-face conversations with shorter turns (e.g., Rossi, 2018; Schegloff, 2007: pp. 73–78). Based on our data, it seems that with comments that contain multiple actions, users might treat 1-2 actions as the 'main job(s)' of a previous turn. For example, in comment (5), B does not address the accusations (that had lower scores) made in comment (4) by A, but chooses to respond to the *question* or *request* instead.

These insights illustrate how modeling actions by adopting a multi-label approach better aligns with how actions in this context are used, in contrast to the single-label approach. They also emphasize how score-based annotation helps to account for overlapping actions, and to represent signal strength, i.e. how some actions are more prevalent in a turn than others.

Next, to test how limiting predictions to a single action per comment would affect classification, we compared *1-act* models to models allowing multiple actions per comment: see Table 10.

Considering multiple actions vs. one, SVM achieves statistically improved performances for only two classes at best. With *FinBERT1-Avg* the increase in performance is statistically significant for only one action. *SetFit-FinBERT1* fares much better: *SetFit-FinBERT1-MS* achieves statistically higher performance in classifying three actions. For the remaining actions, despite higher F1s, statistical tests do not show $p < 0.05$. *SetFit-FinBERT1-Avg* fares better, achieving a statistically significant improvement in performance for five actions. *SetFit-FinBERT2-MS* achieves

| User | Comment | Annotations | | |
|------|---------|-------------|---|---|
| | | Annotator1 | Annotator2 | Annotator3 |
| (1) A | *B*, The country hasn't attacked. I'm speaking from personal experience. I moved to Russia for a year when I was 17. In other occasions I've been there hundreds of times, both for fun, and driving a bus. Nothing has ever happened, the people are friendly and you can always get help. It's a different story in Finland. | statement:6 denial:6 apprec.:2 | statement:5 challenge:5 denial:3 | challenge:5 statement:4 denial:3 |
| (2) A | *C*, war is in no way acceptable in any form. But these issues aren't so black and white. How much more often have you and *B* visited Ukraine and spent time there compared to myself? | question:6 statement:3 | question:6 statement:6 challenge:4 | question:6 statement:6 challenge:4 |
| (3) B | *A*, Dear *A*. How many times each of us has traveled to or lived in Russia, well that has no bearing on whether Russia is a rogue state ruled by a dictator. It is, believe it or not. It's that black-and-white. I hope you'll get over your gullibility someday. Your actions only support Putin's dictatorship. | accusation:6 statement:3 | statement:6 request:4 challenge:4 | statement:6 request:4 challenge:3 |
| (4) A | *B*, Well tell me about your own experiences, traveling in Ukraine, have you got an impression that it's somehow a more virtuous country? Surely you too have seen those Nazi arguments in the media there. Children have been killed by Ukraine for 8 years. For example, Ukrainian authorities did not let me into the Donetsk region at all. The answer to my question 'why' was simply 'because' and that's that. They were probably afraid I was a journalist. Still, it's a wonderful country, too, and the citizens. After all, Finland is also a country that welcomes all kinds of scum and yet we also love our own homeland. | accusation:6 question:6 apprec.:4 | statement:6 question:5 request:5 accusation:2 | statement:6 question:5 request:5 accusation:1 |
| (5) B | *A*, Again, my travel experiences aren't in any way related to the actions of any government. Parroting it won't change anything. If your next claim is that it's right for the Ukrainians that the wonderful country of Russia invaded Ukraine, then good luck with that. | statement:6 | statement:5 challenge:5 denial:3 accusation:2 | statement:5 challenge:5 denial:4 accusation:1 |

Table 9: Annotated extract from a Ukraine war related conversation.

| Model | Data | Action | | | | | | | |
|-------|------|--------|---------|-----------|-----------|-----------|------------|-------|--------------|
| | | question | request | statement | accusation | challenge | acceptance | denial | appreciation |
| SVM-1-act | single | 0.60 | 0.53 | 0.62 | 0.52 | 0.52 | 0.55 | 0.53 | 0.56 |
| SVM-MS | single | 0.66 | 0.58 | 0.63 | 0.62* | 0.52 | 0.69 | 0.59 | 0.51 |
| SVM-Avg | multiple | 0.63 | 0.60 | 0.62 | 0.62* | 0.53 | 0.69** | 0.59 | 0.51 |
| FinBERT1-1-act | single | 0.86 | 0.63 | 0.70 | 0.57 | 0.54 | 0.65 | 0.57 | 0.62 |
| FinBERT1-MS | single | 0.87 | 0.69 | 0.73 | **0.76** | 0.61 | 0.66 | 0.56 | 0.65 |
| FinBERT1-Avg | multiple | 0.85 | 0.70 | 0.77** | 0.63 | **0.68** | 0.62 | 0.51 | 0.75 |
| SetFit-FinBERT1-1-act | single | 0.81 | 0.60 | 0.66 | 0.52 | 0.50 | 0.59 | 0.55 | 0.50 |
| SetFit-FinBERT1-MS | single | 0.94 | **0.82*** | 0.77 | 0.73* | 0.64 | **0.86*** | **0.68** | **0.78** |
| SetFit-FinBERT1-Avg | multiple | **0.97*** | 0.80** | 0.78** | 0.65* | 0.63 | 0.79** | 0.58 | 0.72 |
| SetFit-FinBERT2-1-act | single | 0.76 | 0.59 | 0.66 | 0.52 | 0.55 | 0.57 | 0.49 | 0.50 |
| SetFit-FinBERT2-MS | single | 0.95 | 0.69 | 0.76 | 0.56 | 0.64* | 0.76* | 0.66* | 0.60 |
| SetFit-FinBERT2-Avg | multiple | 0.95* | 0.67 | **0.81*** | 0.58 | 0.54 | 0.71 | 0.61 | 0.53 |
| FinGPT-0-shot-1-act | single | 0.29 | 0.25 | 0.48 | 0.21 | 0.23 | 0.24 | 0.24 | 0.19 |
| FinGPT-0-shot-MS | single | 0.34 | 0.31 | 0.51 | 0.33 | 0.40 | 0.29 | 0.30 | 0.25 |
| FinGPT-0-shot-Avg | multiple | 0.37 | 0.24 | 0.51 | 0.31 | 0.40 | 0.30 | 0.30 | 0.20 |
| Llama2-0-shot-1-act | single | 0.47 | 0.49 | 0.29 | 0.48 | 0.49 | 0.49 | 0.46 | 0.55 |
| Llama2-0-shot-MS | single | 0.47 | 0.54 | 0.37 | 0.46 | 0.48 | 0.50 | 0.48 | 0.50 |
| Llama2-0-shot-Avg | multiple | 0.50 | 0.52 | 0.32 | 0.43 | 0.41 | 0.51 | 0.47 | 0.52 |
| Llama2-finetuned-1-act | single | 0.62 | 0.51 | 0.51 | 0.55 | 0.48 | 0.52 | 0.51 | 0.49 |
| Llama2-finetuned-MS | single | 0.71 | 0.61 | 0.80 | 0.56 | 0.47 | 0.81** | 0.62 | 0.57 |
| Llama2-finetuned-Avg | multiple | 0.57 | 0.67* | **0.81*** | 0.65 | 0.64 | 0.69 | 0.39 | 0.52 |

Table 10: 10-fold cross-validated macro-F1 scores for models using single vs. multiple annotations data, using conservative ground truth. Statistically significant differences are indicated with $*p < 0.05$, $**p < 0.01$, $***p < 0.001$.

significantly higher performances for 3 actions, *SetFit-FinBERT2-Avg* for 2. Generative LMs do not perform very well: only *Llama2-finetuned* achieves notably higher performances when allowing multiple actions in contrast to *Llama2-finetuned-1-act*, but only for *acceptances* (*Llama2-finetuned-MS*), and for *requests* and *statements* (*Llama2-finetuned-Avg*).

To answer RQ1, we conclude that considering multiple actions – and using score-based annotation – to classify actions much better aligns with how people utilize actions in asynchronous online conversation in contrast to the single-label or categorical annotation approach. This is supported by our insights gained during annotation scheme development, examination of example cases, and model performance tests where our best models achieved significantly better performance for many actions with the multi-label approach.

## 6.2 RQ2: Modeling Ambiguity

To answer RQ2, we utilized multiple annotations and three different approaches to leverage them. We used our best-performing model from previous experiments, SetFit-FinBERT1, in order to discover how to best represent the ambiguity related to actions in our data. To this end, we compared different modeling configurations for leveraging multiple annotations, analyzed label correlations in more detail, and examined extracts from our data to see how models performed in comparison to manual annotations for these cases.

Performances for models where multiple annotations were utilized can be seen in Table 11. Except for *PNC*, performances are evaluated using both conservative and relaxed ground truth.

*SetFit-FinBERT1-PNC* performs quite poorly. *PNC* models performed worst for the Complicated class, similarly to Jiang and de Marneffe (2022), perhaps due to class heterogeneity; other classes' performances are higher. *A2* models performed best overall if considering which models achieved most of the highest F1 scores. Macro-F1s were higher for some actions when predicting relaxed ground truth labels. However, here results differ according to action. Also, *A1-A3* model performances differed according to annotator.

Overall, for most actions, the models utilizing multiple annotations, or individual annotators' annotations, achieved higher performances than the models using only single annotations (see Tables 10 and 11) – except for *requests* and *acceptances*, for which *SetFit-FinBERT1-MS* had higher F1s.

Comparing the performances against earlier work modeling similar actions (Paakki et al., 2024), we achieved higher performances for all actions, except for *appreciations*. We also achieved higher F1s for some actions modeled in prior work on asynchronous data: message-level detection of *request* and *question* trained



Figure 1: Correlation scores between all 8 actions.

with BERT (Taniguchi et al., 2020), *question* detection (and overall performance) in contrast to Zhang et al. (2017), and better agreement for *challenges* in comparison to Bracewell et al. (2012). However, the comparability to the latter three works is limited due to differences in theoretical frameworks, data, annotation practices or modeling choices.

Differences between actions deserve some attention: overall, classifiers for predicting *accusations*, *challenges*, *denials* and *appreciations* seem to have lower performance compared to other actions. The latter three also have notably lower Micro-F1 scores when investigating performances for the positive class for the *SetFit-FinBERT1-A2* model (see Appendix A). Inspecting annotator performances against ground truth (Table 7), performances for face-threatening actions – *accusations*, *challenges* and *denials* – are notably lower, although there are differences when using conservative *vs.* relaxed ground truth. Label score averages, in Table 6, also showed that *accusations*, *challenges*, *denials*, and *appreciations* have lower label scores overall. Measuring score correlations between annotators, by action, we see (in Table 6) that *accusations*, *challenges*, and *denials* have some of the lowest correlations. *Appreciations* also have low correlations between Annotator1 and 2, and Annotator1 and 3. There were very few *appreciations* in the data, which might affect the results. Also, *statements* have surprisingly low annotation score correlations. These might be partly due to *statements* being much more frequent in comments in contrast to other actions (see Table 5): present in 71% of all comments (in single annotations; 83% in multiple). Trained model performances for *statements*, on the other hand, are quite good.

Investigating the correlations between all annotated action labels for each comment (see Figure 1), there is a weak positive correlation (Schober et al., 2018) with $p < 0.001$ between *challenges* and *denials*, and a weak positive correlation between *questions* and *chal-*

| Model | Ground truth | question | request | statement | accusation | challenge | acceptance | denial | appreciation |
|---|---|---|---|---|---|---|---|---|---|
| SetFit-FinBERT1-PNC | cross-e. | 0.71 | 0.55 | 0.57 | 0.50 | 0.50 | 0.59 | 0.47 | 0.39 |
| SetFit-FinBERT1-Avg | conservative | 0.96 | 0.74 | 0.74 | 0.65 | 0.54 | 0.73 | 0.58 | 0.72 |
| SetFit-FinBERT1-A1 | conservative | 0.92 | 0.75 | 0.75 | 0.51 | 0.43 | 0.77 | 0.55 | 0.65 |
| SetFit-FinBERT1-A2 | conservative | 0.95 | 0.77 | **0.84** | 0.72 | **0.76** | 0.78 | **0.68** | 0.66 |
| SetFit-FinBERT1-A3 | conservative | **0.97** | 0.71 | 0.78 | 0.66 | 0.52 | 0.74 | 0.63 | 0.69 |
| SetFit-FinBERT1-Avg | relaxed | **0.97** | **0.80** | 0.76 | 0.65 | 0.63 | 0.79 | 0.56 | 0.65 |
| SetFit-FinBERT1-A1 | relaxed | 0.95 | 0.77 | 0.79 | 0.60 | 0.52 | **0.82** | 0.62 | 0.52 |
| SetFit-FinBERT1-A2 | relaxed | **0.97** | 0.76 | 0.79 | **0.75** | 0.61 | 0.78 | **0.68** | **0.73** |
| SetFit-FinBERT1-A3 | relaxed | 0.96 | 0.71 | 0.78 | 0.66 | 0.52 | 0.74 | 0.63 | 0.68 |

Table 11: 10-fold cross-validated macro-F1 scores for SetFit-FinBERT1, utilizing different configurations for leveraging multiple annotations.

| Ground truth | Ensemble | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | A1 | A2 | A3 | Avg. | A1+A3 | A1+A2 | A2+A3 | A1+A2 +A3 | A1 +Avg. | A2 +Avg. | A3 +Avg. | A1+A2 +Avg. | A1+A3 +Avg. | A2+A3 +Avg. | A1+A2+ A3+Avg. |
| conservative | 0.56 | 0.69 | 0.63 | 0.39 | 0.67 | 0.70 | **0.72** | **0.72** | 0.51 | 0.60 | 0.56 | 0.61 | 0.58 | 0.62 | 0.62 |
| relaxed | 0.54 | 0.50 | 0.52 | **0.62** | 0.52 | 0.48 | 0.48 | 0.48 | 0.55 | 0.46 | 0.50 | 0.45 | 0.50 | 0.45 | 0.44 |
| | A1 +PNC | A2 +PNC | A3 +PNC | Avg. +PNC | A1+A3 +PNC | A1+A2 +PNC | A2+A3 +PNC | A1+A2 +A3 +PNC | A1 +Avg. +PNC | A2 +Avg. +PNC | A3 +Avg. +PNC | A1+A2 +Avg. +PNC | A1+A3 +Avg. +PNC | A2+A3 +Avg. +PNC | A1+A2+ A3+Avg. +PNC |
| conservative | 0.46 | 0.55 | 0.51 | 0.35 | 0.52 | 0.55 | 0.56 | 0.56 | 0.46 | 0.54 | 0.50 | 0.54 | 0.51 | 0.62 | 0.55 |
| relaxed | 0.48 | 0.40 | 0.44 | 0.45 | 0.45 | 0.40 | 0.40 | 0.40 | 0.47 | 0.40 | 0.43 | 0.39 | 0.43 | 0.44 | 0.39 |

Table 12: Jaccard coefficient scores for ensemble models using SetFit-FinBERT1. A1=Annotator1 model, A2=Annotator2, A3=Annotator3, Avg=Averaged, PNC=Positive/Negative/Complicated.

*lenges*, only slightly above the lower bound of weak correlation ($p < 0.001$). For most actions, there is a weak or negligible negative (or positive) correlation magnitude according to the statistic with no statistical significance. There is a moderate negative correlation between *questions* and *statements* with $p < 0.001$.

Finally, comparing model ensembles for predicting all possible annotations for comments, in Table 12, an ensemble of *A1+A2* or *A1+A2+A3* fared best for conservative ground truth, according to Jaccard coefficient scores. *Avg* fared best with relaxed ground truth labels. However, for relaxed ground truth Jaccard similarity drops notably in contrast to conservative ground truth.

A further investigation of difficult cases and comparison between human annotations and model predictions provides more insight. Tables 13–14 illustrate what annotation scores human annotators gave to some of the comments with annotation disagreement, and what our best models predicted as actions present in the comments. For the second column, we only list actions that at least one annotator deemed to be present in the comment, and for the third column, actions predicted as being present by at least one

model. For model predictions, we always used a model trained on a data split that did not include the to-be-classified comment.

There was disagreement between humans and models in some cases, e.g. with *challenges* and *accusations*, but model disagreements correspond with human disagreements quite well if investigating them side-by-side with our best models. Many comments required interpretation of implicit meaning: e.g. ambiguity, contextual knowledge, proverbs or even poetic expression. This seemed to often result in multiple potential interpretations. A case in point is B's comment in Table 13, which was interpreted by human annotators as either sarcastic, a *challenge* (or *accusation*) presented together with a *statement*, or just *statement*. A less likely interpretation would be sincere *accepting statement*, not chosen by human annotators. Similarly, "hallelujah" in Table 1 is ambiguous, interpreted by human annotators as either sincerely appreciative, or sarcastic, a *challenge* merged with a *statement*. Models predicted it as a *statement* (A1, A2, A3, Avg.), a *challenge* (A2), but also as an *appreciation* (A2). Another ambiguous example is B's comment in Table 14. It depends on subjective interpretation whether hindsight can be seen as morally repre-

Detecting Face-threatening and Paired Actions in Asynchronous Online Conversations

| Comments | Annotations (Annotator: score) | Predictions (Model IDs) |
|---|---|---|
| A: What other nation could be responsible for Russia's invasion of Ukraine besides the Russian people? All else is just about making excuses. (27 likes) | *question* (A1: 5, A2: 4, A3: 3) *challenge* (A2: 5, A3: 5) *statement* (A1: 0, A2: 0, A3:1) *accusation* (A1: 3, A2: 0, A3: 0) | *question* (A1-3, Avg) *challenge* (A2, A3) *statement* (A3) |
| **B: yea the people did attack there on their own initiative, without any government authority and orders.** (2 likes) | *statement* (A1: 0, A2: 0, A3: 3) *accusation* (A1: 0, A2: 0, A3: 2) *challenge* (A1: 4, A2: 5, A3: 5) | *statement* (A1-3, Avg) *accusation* (A2) |

Table 13: Extract from Ukraine war discussion on Russian citizens' role in the war, HS 10/2022. Second column lists annotation scores by annotators, third column the SetFit-FinBERT1 model IDs predicting that listed action exists for the comment in first column.

hensible. Not all annotators agreed here. Model predictions corresponded quite well with human annotators' interpretations, and their disagreements. However, for some cases, like Tables 13 and 14, where human annotators disagreed on whether a comment should be labeled as an *accusation* or a *challenge*, models did not predict the ambiguous comment as a challenge.

Overall, an ensemble of individual annotator-based models seems to best align with the human annotations, also quite well representing many of the disagreements. Based on the Jaccard coefficient scores and the above empirical insights, we conclude in response to RQ2 that an ensemble of individual annotator based models best represents the multiple possible interpretations of actions in our data.

# 7 Discussion

We investigated how paired actions in comments to asynchronous crisis-related conversations could be computationally modeled, especially face-threatening actions central to misbehavior like trolling. We illustrated that, in this context, it is important to take into account multiple possible label interpretations, and multiple actions often performed in one comment. The conversational context we studied is largely assertive in nature, statements being notably more common than other actions. As we assumed, face-threatening actions like *accusations* and *challenges* are also common, which highlights the need for including them in computational models used for identifying actions in these conversations.

We showed that allowing multiple actions predicted per comment helps to align models with the empirical

| Comments | Annotations (Annotator: score) | Predictions (Model IDs) |
|---|---|---|
| A: Would it have been worth taking an interest in things before these men came to power? (21 likes, 2 angry emojis) | *question* (A1: 3, A2: 5, A3: 3) *accusation* (A1: 2, A2: 4, A3: 3) | *question* (A1-3, Avg.) *accusation* (A3) |
| **B: *A*, Well here we have a real master of hindsight.** (Dog sticker facepalming) (15 likes) | *statement* (A1: 0, A2: 3, A3: 5) *accusation* (A1: 0, A2: 3, A3: 5) *challenge* (A1: 3, A2: 0, A3: 0) | *statement* (A1-3, Avg.) *accusation* (A2) |
| A: *B*, hindsight is a sweet thing… (3 likes, 1 angry) | *statement* (A1: 3, A2: 4, A3: 5) *accusation* (A1: 0, A2: 0, A3: 3) *appreciation* (A1: 0, A2: 3, A3: 0) | *statement* (A1-3, Avg.) *accusation* (A2) |

Table 14: Extract from Ukraine war discussion on how Russian soldiers are treated, YLE 12/2022. Second column lists annotation scores by annotators, third column the SetFit-FinBERT1 model IDs predicting that listed action exists for the comment in first column.

phenomenon studied, and statistically improves model performance. SetFit models (initialized with FinBERT + word embeddings) performed best overall. We also demonstrated improved classification performance in contrast to earlier action detection in asynchronous conversations. Likert score annotation allowed us to consider the presence of actions on different scales of strength, rather than categorically (cf. Glickman and Dagan, 2005). We illustrated that this was important for representing the main functions of a comment. Based on our empirical insights, signal strength affected how paired actions and their norms were treated by discussion participants: in the case of multiple (pair-initiating) actions, participants tended to orient to (1-2) main actions as having the strongest normative expectations for responding (see Table 9).

We demonstrated that for predicting all possible labels, with conservative ground truth (positive label if even one annotator gave a score ≥ 3), an ensemble of models trained separately on individual annotators' annotations performed best. Based on our results, we highlight that although averaging has been popular in accounting for annotation disagreements, it is not always the best option as it might lead to losing important information on multiple possible interpretations of labels. Generative LMs performed quite poorly on the action detection task, except for specific classes (*statements* and *acceptances*) with Llama2-finetuned. In other words, these models seemed to recognize a limited amount of categories even when finetuned for the downstream task.

We provided an annotation scheme for identifying face-threatening and paired actions in asynchronous conversations in Finnish. This is important as there are no such resources for Finnish yet. We feel that future work analyzing manipulative or crisis news related conversational behaviors online will benefit from utilizing our scheme and models. The scheme will enable easier implementation of novel models for other languages as well. However, although some actions in our scheme are common across languages (Enfield et al., 2010), contextual differences should be considered: although e.g. *apologies* were not found in our data, they have been relevant elsewhere (Paakki et al., 2021). Furthermore, we showed that even with a relatively small annotated dataset we can reach sufficient performance using few-shot learning. This is important since adapted or novel models are often needed in low-resource settings.

We illustrated some differences between actions: e.g., for *questions*, our models reached high performances, but for some others performances were lower. An interesting result found was that face-threatening actions (*challenge*, *denial*, *accusation*) were more difficult to annotate and/or model than others. They also involved more disagreement between annotators than other actions. This is in line with theoretical views on social actions: people tend to express these more implicitly or indirectly to avoid face-threats (Brown and Levinson, 1987), which might lead to uncertainty in their interpretation. *Appreciations*, surprisingly, also portrayed a lower agreement. This might have been due to the low number of appreciations in our data; or they might have been expressed in an ambiguous manner. Whether this is a context-specific tendency (e.g. platform or cultural context) remains an open question.

We saw correlation between *challenges* and *denials*, and to some extent between *questions* and *challenges*. This might be because *challenges* and *denials* were at times difficult to distinguish from each other. We may have needed more detailed clarification of class boundaries here. However, the annotation scheme (and guideline) development process was already very extensive. Also, based on examples seen during the development, *denials* were often followed by *challenges* of epistemic claims (or of interlocutors' positions), both in the same comment. This might explain the correlation. Finally, there seemed to be a significant amount of *challenging questions* in our data, which could explain the correlation between these two actions.

We found that users tended to treat a previous turn as having some main action, not necessarily responding to all actions. Given that in CA action has been viewed as the 'main job' of a turn, as turns in asynchronous conversation are longer and often involve multiple actions, it is interesting how users interpret the 'main job' of such turns. Although there are many factors that may affect the interpretation of actions and required responses (Stivers and Rossano, 2010), based on empirical insights, it seems that the strength of label signal can affect how strong normative expectations pair initiating actions might incur on subsequent turns. Thus, actions that are present with at least moderate strength of signal will likely be fruitful for analyzing responding behaviors in relation to paired actions (and face-threats).

From a (digital) CA perspective, it is challenging to systematically identify actions in asynchronous conversations due to actions often being context-dependent, implicit or indirect. Interpretation is not a product but a process: meanings of actions are interpreted by participants collaboratively and on-line (Clark and Schaefer, 1989; Jurafsky, 1992). Participants might alter interpretations of comments across turns in conversation. Ambiguity in the expression of (face-threatening and other) actions might also be a strategic choice. Thus, we consider it crucial to be able to model ambiguity in the expression of actions, especially face-threatening actions, when studying crisis-related or manipulative asynchronous conversations.

# 8 Conclusions and Future Work

To conclude, we investigated how to approach action detection in asynchronous crisis news conversations. Our computational approach was able to reflect how multiple actions were performed within the same comment, and the ambiguity related to actions, with improved classification performance in contrast to earlier action detection for asynchronous conversation. Although annotator disagreements have been studied increasingly in NLP, there is still room for exploring how to utilize them in the analysis of face-threatening and paired actions. The contributions of this paper included 1.) portraying and modeling disagreement in the annotation of actions in asynchronous conversations in Finnish, 2.) a paired action annotation framework and dataset for asynchronous conversation including face-threatening actions, and 3.) models with improved classification performance for many of these actions. Although our study focuses on Finnish, the framework can be applied for other languages as well.

Future work could investigate how to further address the differing nature of actions, and how to further utilize annotation scores and contextual information in more fine-grained models.

We conclude that representing ambiguity in computational modeling is especially relevant to analyzing face-threatening actions. These actions, in turn, are crucial for the study of online misbehavior. Digital CA based approaches, having a robust theoretical understanding of such actions, can be fruitful for analyzing meaningful ambiguity related to how actions are per-

formed in asynchronous conversations online.

## 8.1   Limitations

Crowdsourcing is often seen to provide heterogeneous, arguably more valid annotations from a large population (Weber et al., 2018). In expert annotation, annotators adjust their work based on expectations regarding outcomes, thus reaching higher agreements, annotation reliability maximized to reflect the desired categories (Weber et al., 2018) – a possible limitation of our work. However, CA analysis requires contextual in-depth reading, which is why non-expert annotation would have been unreliable (Eickhoff, 2018). We had only three annotators; a higher number might lead to better results (Pavlick and Kwiatkowski, 2019). Due to resource limitations, we considered our current scope of annotators and data sufficient for now.

Score-based annotation allows fluidity in class interpretation, but subjective reading, confidence, signal strength, and understandings of scores and categories might be melted into one metric. Also, in contrast to message-level classification, we could have labeled segmented comments. However, we considered this challenging as sometimes the boundaries of actions were unclear or actions tended to overlap (e.g., Table 9). The boundaries of an action did not always correspond to sentence boundaries. Thus, we deemed message-level scoring best. Future research could explore the sentence segmentation option further. The choice of threshold (3) was theoretically and empirically motivated, but in some cases it might be relevant to be able to computationally identify even very weak action signals. We could have further investigated how to model conversational context or sequential dependencies more intricately (e.g. whether a comment responds to a certain action). However, these endeavors are best left for future research.

At this point we did not utilize context information when training models, which is a potential limitation since annotators had access to this information. However, we deemed it interesting to model the potential for different interpretations of a turn considered alone, through scores. We decided to do this because we were concerned that using previous turns (or their actions) as features would confound the computational analysis of responding behaviors. This is because people favor certain types of responses due to conversational norms (e.g. Stivers and Rossano, 2010). However, deviations from these norms can be highly meaningful (e.g. Paakki et al., 2024), so if models would overgeneralize normative responses as predicted actions for responses, this would confound computational analyses of less common deviant responses. Future work could expand on this to develop more contextually sensitive computational models that could also detect contra-normative responses.

Since we aimed at a simplified model, using applied digital CA, we could not strictly follow the analytical practices typical for CA. Restricting the granularity of the annotation scheme limits the interpretation of actions, forcing annotators (and thus models) to ignore some rarer actions. However, this process is important for a computational classification approach like ours. Secondly, it is debatable whether some of the classes in our applied approach are notably different in contrast to some DA classification approaches (e.g., Taniguchi et al., 2020). However, we see some of the added classes (e.g., accusations) as very different from those included in earlier DA models based on Speech Acts. Also, our key literature related to action pairs and normative expectations of face-threatening actions comes from CA (Schegloff, 2007; Turowetz and Maynard, 2010; Dersley and Wootton, 2000), which views these actions and paired actions somewhat distinctly, in our view. We wished to highlight these theoretical foundations in this paper.

Crisis related discussions may involve sensitive information, even when dealing with publicly available social media data. We have translated the examples, and anonymized and de-identified the data so that the content is conveyed without privacy concerns. We published a privacy notice according to Aalto University policy, regarding data collection and management, on our research project's website during the study.[16] We only release models and materials where any possibly sensitive information has been removed.

# Acknowledgements

---

[16]The project has already ended, so the website is no longer operational, however. For project details, see: https://research.aalto.fi/en/projects/crisissawhney/

[17]https://blogs.helsinki.fi/disinformation-news-media/

# References

Akiba, Takuya, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. 2019. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '19, page 2623–2631, New York, NY, USA. Association for Computing Machinery.

Allen, James and Mark Core. 1997. DAMSL: Dialog act markup in several layers. *Draft of manual (31 March 1997).*

Antaki, Charles, Michela Biazzi, Anette Nissen, and Johannes Wagner. 2008. *Accounting for moral judgments in academic talk: The case of a conversation analysis data session.* Walter de Gruyter.

Arabadzhieva-Kalcheva, Neli and Ivelin Kovachev. 2022. Comparison of BERT and XLNet accuracy with classical methods and algorithms in text classification. In *2021 International Conference on Biomedical Innovations and Applications (BIA)*, volume 1, pages 74–76. IEEE.

Austin, John Langshaw. 1975. *How to Do Things with Words.* Oxford university press.

Bakhtin, Mikhail Mikhaĭlovich. 1981. *The dialogic imagination: Four essays*, volume 1. University of Texas Press, Austin, Texas.

Barnhurst, Kevin G. and Diana Mutz. 1997. American journalism and the decline in event-centered reporting. *Journal of Communication*, 47(4):27–53.

Bellutta, Daniele, Catherine King, and Kathleen M Carley. 2021. Deceptive accusations and concealed identities as misinformation campaign strategies. *Computational and Mathematical Organization Theory*, 27:302–323.

Bhatia, Sumit, Prakhar Biyani, and Prasenjit Mitra. 2014. Summarizing online forum discussions – can dialog acts of individual messages help? In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2127–2131, Doha, Qatar. Association for Computational Linguistics.

Bracewell, David, Marc Tomlinson, and Hui Wang. 2012. Identification of social acts in dialogue. In *Proceedings of COLING 2012*, pages 375–390, Mumbai, India. The COLING 2012 Organizing Committee.

Bracewell, David B., Marc Tomlinson, and Hui Wang. 2013. Semi-supervised modeling of social actions in online dialogue. In *2013 IEEE Seventh International Conference on Semantic Computing*, pages 168–175. IEEE.

Brown, Penelope and Stephen C. Levinson. 1987. *Politeness: Some Universals in Language Usage*, volume 4. Cambridge university press.

Carvalho, Vitor R. and William W Cohen. 2005. On the collective classification of email "speech acts". In *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 345–352.

Casanueva, Iñigo, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. Efficient intent detection with dual sentence encoders. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45, Online. Association for Computational Linguistics.

Castro, Stephanie L. 2002. Data analytic methods for the analysis of multilevel questions: A comparison of intraclass correlation coefficients, rwg(j), hierarchical linear modeling, within- and between-analysis, and random group resampling. *The Leadership Quarterly*, 13(1):69–93.

Clark, Alexander and Andrei Popescu-Belis. 2004. Multi-level dialogue act tags. In *SIGdial 2004 (5th SIGdial Workshop on Discourse and Dialogue)*, pages 163–170. Association for Computational Linguistics.

Clark, Herbert H and Edward F Schaefer. 1987. Collaborating on contributions to conversations. *Language and cognitive processes*, 2(1):19–41.

Clark, Herbert H. and Edward F. Schaefer. 1989. Contributing to discourse. *Cognitive Science*, 13(2):259–294.

Cohen, William W., Vitor R. Carvalho, and Tom M. Mitchell. 2004. Learning to classify email into "speech acts". In *Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing*, pages 309–316, Barcelona, Spain. Association for Computational Linguistics.

Davani, Aida Mostafazadeh, Mark Díaz, and Vinodkumar Prabhakaran. 2022. Dealing with Disagreements: Looking Beyond the Majority Vote in Subjective Annotations. *Transactions of the Association for Computational Linguistics*, 10:92–110.

Davidson, Sam, Qiusi Sun, and Magdalena Wojcieszak. 2020. Developing a new classifier for automated identification of incivility in social media. In *Proceedings of the Fourth Workshop on Online Abuse and Harms*, pages 95–101, Online. Association for Computational Linguistics.

Derrac, Joaquín, Salvador García, Daniel Molina, and Francisco Herrera. 2011. A practical tutorial on the use of nonparametric statistical tests as a methodology for comparing evolutionary and swarm intelligence algorithms. *Swarm and Evolutionary Computation*, 1(1):3–18.

Dersley, Ian and Anthony Wootton. 2000. Complaint sequences within antagonistic argument. *Research on language and social interaction*, 33(4):375–406.

Dettmers, Tim, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. QLoRa: Efficient finetuning of quantized LLMs. *Advances in Neural Information Processing Systems*, 36.

Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Di Mascio, Fabrizio, Michele Barbieri, Alessandro Natalini, and Donatella Selva. 2021. Covid-19 and the information crisis of liberal democracies: Insights from anti-disinformation action in Italy and EU. *Partecipazione e conflitto*, 14(1):221–240.

Duran, Nathan and Steve Battle. 2018. Conversation analysis structured dialogue for multi-domain dialogue management. In *The International Workshop on Dialogue, Explanation and Argumentation in Human-Agent Interaction (DEXAHAI)*.

Duran, Nathan, Steven Battle, and Jim Smith. 2022. Inter-annotator agreement using the conversation analysis modelling schema, for dialogue. *Communication Methods and Measures*, 16(3):182–214.

Eickhoff, Carsten. 2018. Cognitive biases in crowdsourcing. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, WSDM '18, page 162–170, New York, NY, USA. Association for Computing Machinery.

Enfield, Nick J and Jack Sidnell. 2017. On the concept of action in the study of interaction. *Discourse Studies*, 19(5):515–535.

Enfield, N.J., Tanya Stivers, and Stephen C. Levinson. 2010. Question–response sequences in conversation across ten languages: An introduction. *Journal of Pragmatics*, 42(10):2615–2619.

Feng, Donghui, Erin Shaw, Jihie Kim, and Eduard Hovy. 2006. Learning to detect conversation focus of threaded discussions. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pages 208–215, New York City, USA. Association for Computational Linguistics.

Ferracane, Elisa, Greg Durrett, Junyi J. Li, and Katrin Erk. 2021. Did they answer? Subjective acts and intents in conversational discourse. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1626–1644, Online. Association for Computational Linguistics.

Forsyth, Eric and Craig Martell. 2007. Lexical and discourse analysis of online chat dialog. In *International Conference on Semantic Computing (ICSC 2007)*, pages 19–26. IEEE.

Fortuna, Blaz, Eduarda Mendes Rodrigues, and Natasa Milic-Frayling. 2007. Improving the classification of newsgroup messages through social network analysis. In *Proceedings of the Sixteenth ACM Conference on Conference on Information and Knowledge Management*, CIKM '07, page 877–880, New York, NY, USA. Association for Computing Machinery.

Fuscone, Simone, Benoit Favre, and Laurent Prévot. 2020. Filtering conversations through dialogue acts labels for improving corpus-based convergence studies. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 203–208, 1st virtual meeting. Association for Computational Linguistics.

Garimella, Kiran, Gianmarco De Francisci Morales, Aristides Gionis, and Michael Mathioudakis. 2018. Quantifying controversy on social media. *ACM Transactions on Social Computing*, 1(1):1–27.

Ghosh, Souvick and Satanu Ghosh. 2021. Classifying speech acts using multi-channel deep attention network for task-oriented conversational search agents. In *Proceedings of the 2021 Conference on Human Information Interaction and Retrieval*, CHIIR'21, page 267–272, New York, NY, USA. Association for Computing Machinery.

Giglietto, Fabio, Nicola Righetti, Luca Rossi, and Giada Marino. 2020. It takes a village to manipulate the media: coordinated link sharing behavior during 2018 and 2019 Italian elections. *Information, Communication & Society*, 23(6):867–891.

Giles, David, Wyke Stommel, Trena Paulus, Jessica Lester, and Darren Reed. 2015. Microanalysis of online data: The methodological development of "digital CA". *Discourse, context & media*, 7:45–51.

Glickman, Oren and Ido Dagan. 2005. A probabilistic setting and lexical coocurrence model for textual entailment. In *Proceedings of the ACL Workshop on Empirical Modeling of Semantic Equivalence and Entailment*, pages 43–48, Ann Arbor, Michigan. Association for Computational Linguistics.

Godfrey, John J., Edward C. Holliman, and Jane McDaniel. 1992. Switchboard: Telephone speech corpus for research and development. In *Acoustics, Speech, and Signal Processing, IEEE International Conference on*, volume 1, pages 517–520. IEEE Computer Society.

Goffman, Erving. 1974. *Frame Analysis: An Essay on the Organization of Experience.* Harvard University Press.

Guo, Yuting and Abeed Sarker. 2023. SocBERT: A pretrained model for social media text. In *Proceedings of the Fourth Workshop on Insights from Negative Results in NLP*, pages 45–52, Dubrovnik, Croatia. Association for Computational Linguistics.

Haverinen, Katri, Jenna Nyblom, Timo Viljanen, Veronika Laippala, Samuel Kohonen, Anna Missilä, Stina Ojala, Tapio Salakoski, and Filip Ginter. 2014. Building the essential resources for Finnish: the Turku Dependency Treebank. *Language Resources and Evaluation*, 48:493–531. Open access.

Herring, Susan. 1999. Interactional coherence in CMC. *Journal of Computer-Mediated Communication*, 4(4).

Herring, Susan, A. Das, and S. Penumarthy. 2005. CMC act taxonomy.

Jeng, Wei, Spencer DesAutels, Daqing He, and Lei Li. 2017. Information exchange on an academic social networking site: A multidiscipline comparison on researchgate Q&A. *Journal of the Association of Information Science and Technology*, 68(3):638–652.

Jiang, Nan-Jiang and Marie-Catherine de Marneffe. 2022. Investigating reasons for disagreement in Natural Language Inference. *Transactions of the Association for Computational Linguistics*, 10:1357–1374.

Joty, Shafiq and Enamul Hoque. 2016. Speech act modeling of written asynchronous conversations with task-specific embeddings and conditional structured models. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1746–1756.

Joty, Shafiq and Tasnim Mohiuddin. 2018. Modeling speech acts in asynchronous conversations: A neural-CRF approach. *Computational Linguistics*, 44(4):859–894.

Jünger, Jakob and Till Keyling. 2019. Facepager. an application for automated data retrieval on the web. *Source code and releases available at* [https://github.com/strohne/Facepager](https://github.com/strohne/Facepager) *(Accessed June 16 2023).*

Jurafsky, Dan. 1997. Switchboard SWBD-DAMSL shallow-discourse-function annotation coders manual. Technical report.

Jurafsky, Daniel S. 1992. *An On-line Computational Model of Human Sentence Interpretation: A Theory of the Representation and Use of Linguistic Knowledge.* University of California, Berkeley.

Kim, Su Nam, Li Wang, and Timothy Baldwin. 2010. Tagging and linking web forum posts. In *Proceedings of the Fourteenth Conference on Computational Natural Language Learning*, pages 192–202.

Koshik, Irene. 2003. Wh-questions used as challenges. *Discourse Studies*, 5(1):51–77.

Laurenti, Enzo, Nils Bourgon, Farah Benamara, Mari Alda, Véronique Moriceau, and Courgeon Camille. 2022. Speech acts and communicative intentions for urgency detection. In *Proceedings of the 11th Joint Conference on Lexical and Computational Semantics*, pages 289–298, Seattle, Washington. Association for Computational Linguistics.

Levinson, Stephen C. 2013. Action formation and ascription. In *The Handbook of Conversation Analysis*, volume 1, pages 103–130. Wiley Online Library.

Lindell, Michael K. and Christina J. Brandt. 1999. Assessing interrater agreement on the job relevance of a test: A comparison of CVI, T, $r_{wg(j)}$, and $r^*_{wg(j)}$ indexes. *Journal of applied psychology*, 84(4):640.

Linell, Per and Ivana Marková. 1993. Acts in discourse: From monological speech acts to dialogical inter-acts. *Journal for the theory of social behaviour*, 23(2):173–195.

Luukkonen, Risto, Ville Komulainen, Jouni Luoma, Anni Eskelinen, Jenna Kanerva, Hanna-Mari Kupari, Filip Ginter, Veronika Laippala, Niklas Muennighoff, Aleksandra Piktus, Thomas Wang, Nouamane Tazi, Teven Scao, Thomas Wolf, Osma Suominen, Samuli Sairanen, Mikko Merioksa, Jyrki Heinonen, Aija Vahtola, Samuel Antao, and Sampo Pyysalo. 2023. FinGPT: Large generative models for a small language. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2710–2726. Association for Computational Linguistics.

Mann, William and Sandra Thompson. 1988. Rhetorical structure theory: Toward a functional theory of

text organization. *Text-interdisciplinary Journal for the Study of Discourse*, 8(3):243–281.

Meredith, Joanne. 2017. Analysing technological affordances of online interactions using conversation analysis. *Journal of Pragmatics*, 115:42–55.

Meredith, Joanne. 2019. Conversation analysis and online interaction. *Research on Language and Social Interaction*, 52(3):241–256.

Meredith, Joanne and Elizabeth Stokoe. 2014. Repair: Comparing facebook 'chat' with spoken interaction. *Discourse & communication*, 8(2):181–207.

Moldovan, Cristian, Vasile Rus, and Arthur C Graesser. 2011. Automated speech act classification for online chat. In *Midwest Artificial Intelligence and Cognitive Science Conference (MAICS)*, volume 710, pages 23–29.

Nie, Yixin, Xiang Zhou, and Mohit Bansal. 2020. What can we learn from collective human opinions on Natural Language Inference data? In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9131–9143, Online. Association for Computational Linguistics.

O'Neill, Thomas A. 2017. An overview of interrater agreement on likert scales for researchers and practitioners. *Frontiers in psychology*, 8:777.

Paakki, Henna, Heidi Vepsäläinen, and Antti Salovaara. 2021. Disruptive online communication: How asymmetric trolling-like response strategies steer conversation off the track. *Computer Supported Cooperative Work (CSCW)*, pages 1–37.

Paakki, Henna, Heidi Vepsäläinen, Antti Salovaara, and Bushra Zafar. 2024. Detecting covert disruptive behavior in online interaction by analyzing conversational features and norm violations. *ACM Transactions on Computer-Human Interaction*, 31(2):1–43.

Pamungkas, Endang Wahyu, Valerio Basile, and Viviana Patti. 2020. Misogyny detection in twitter: a multilingual and cross-domain study. *Information Processing & Management*, 57(6):102360.

Passonneau, Rebecca J., Vikas Bhardwaj, Ansaf Salleb-Aouissi, and Nancy Ide. 2012. Multiplicity and word sense: Evaluating and learning from multiply labeled word sense annotations. *Language Resources and Evaluation*, 46:219–252.

Pavlick, Ellie and Tom Kwiatkowski. 2019. Inherent disagreements in human textual inferences. *Transactions of the Association for Computational Linguistics*, 7:677–694.

Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12:2825–2830.

Peterson, Joshua C., Ruairidh M. Battleday, Thomas L. Griffiths, and Olga Russakovsky. 2019. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9617–9626.

Plank, Barbara. 2022. The 'problem' of human label variation: On ground truth in data, modeling and evaluation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, page 10671–10682. Association for Computational Linguistics.

Pomerantz, Anita. 1984. *Pursuing a response.* Cambridge University Press.

Poth, Clifton, Hannah Sterz, Indraneil Paul, Sukannya Purkayastha, Leon Engländer, Timo Imhof, Ivan Vulić, Sebastian Ruder, Iryna Gurevych, and Jonas Pfeiffer. 2023. Adapters: A unified library for parameter-efficient and modular transfer learning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 149–160, Singapore. Association for Computational Linguistics.

Qadir, Ashequl and Ellen Riloff. 2011. Classifying sentences as speech acts in message board posts. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 748–758, Edinburgh, Scotland, UK. Association for Computational Linguistics.

Rezapour, Rezvaneh, Jutta Bopp, Norman Fiedler, Diana Steffen, Andreas Witt, and Jana Diesner. 2020. Beyond citations: Corpus-based methods for detecting the impact of research outcomes on society. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6777–6785, Marseille, France. European Language Resources Association.

Rossi, Giovanni. 2018. Composite social actions: The case of factual declaratives in everyday interaction. *Research on Language and Social Interaction*, 51(4):379–397.

Sacks, Harvey, Emanuel A. Schegloff, and Gail Jefferson. 1974. The simplest systematics for the organization of turn-taking for conversations. *Language*, 50(4):696–735.

Salonen, Margareta, Margarethe Olbertz-Siitonen, Turo Uskali, and Salla-Maaria Laaksonen. 2022. Conversational gatekeeping—social interactional practices of post-publication gatekeeping on newspapers' facebook pages. *Journalism Practice*, pages 1–25.

Savolainen, Reijo. 2020. Dialogue processes in online information seeking and sharing: a study of an asynchronous discussion group. *Information Research*, 25(3).

Savy, Renata. 2010. Pr.A.Ti.D: A coding scheme for pragmatic annotation of dialogues. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valletta, Malta. European Language Resources Association (ELRA).

Schegloff, Emanuel A. 2007. *Sequence Organization in Interaction: A Primer in Conversation Analysis*. Cambridge University Press, Cambridge; New York.

Schober, Patrick, Christa Boer, and Lothar A Schwarte. 2018. Correlation coefficients: appropriate use and interpretation. *Anesthesia & analgesia*, 126(5):1763–1768.

Shu, Kai, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. *ACM SIGKDD explorations newsletter*, 19(1):22–36.

Stivers, Tanya. 2013. *Sequence organization*. Wiley-Blackwell: Blackwell Publishing Ltd.

Stivers, Tanya. 2015. Coding social interaction: A heretical approach in conversation analysis? *Research on Language and Social Interaction*, 48(1):1–19.

Stivers, Tanya and Federico Rossano. 2010. Mobilizing Response. *Research on Language & Social Interaction*, 43(1):3–31.

Stolcke, Andreas, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational linguistics*, 26(3):339–373.

Stommel, Wyke and Tom Koole. 2010. The online support group as a community: A micro-analysis of the interaction with a new member. *Discourse studies*, 12(3):357–378.

Sudhahar, Saatviga, Giuseppe A. Veltri, and Nello Cristianini. 2015. Automated analysis of the US presidential elections using big data and network analysis. *Big Data & Society*, 2(1).

Taniguchi, Motoki, Yoshihiro Ueda, Tomoki Taniguchi, and Tomoko Ohkuma. 2020. A large-scale corpus of E-mail conversations with standard and two-level dialogue act annotations. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4969–4980, Barcelona, Spain (Online). International Committee on Computational Linguistics.

Terpilowski, Maksim. 2019. Scikit-posthocs: Pairwise multiple comparison tests in python. *The Journal of Open Source Software*, 4(36):1169.

Thomas, Jenny. 1995. *Meaning in Interaction: An Introduction to Pragmatics*. Longman.

Touvron, Hugo, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Troiano, Enrica, Sebastian Padó, and Roman Klinger. 2021. Emotion ratings: How intensity, annotation confidence and agreements are entangled. In *Proceedings of the 11th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 40–49. Association for Computational Linguistics.

Tunstall, Lewis, Nils Reimers, Unso Eun Seo Jo, Luke Bates, Daniel Korat, Moshe Wasserblat, and Oren Pereg. 2022. Efficient few-shot learning without prompts. In *36th Conference on Neural Information Processing Systems (NeurIPS 2022)*, pages 1–14.

Turowetz, Jason J. and Douglas W. Maynard. 2010. *Morality in the social interactional and discursive world of everyday life*. Springer, New York.

Uma, Alexandra N., Tommaso Fornaciari, Dirk Hovy, Silviu Paun, Barbara Plank, and Massimo Poesio. 2022. Learning from disagreement: A survey. *Journal of Artificial Intelligence Research*, 72:1385–1470.

Virtanen, Antti, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. 2019. Multilingual is not enough: BERT for Finnish. *arXiv preprint arXiv:1912.07076*.

Virtanen, Mikko T. and Liisa Käántä. 2018. At the intersection of text and conversation analysis: analysing asynchronous online written interaction. *AFinLA-e : soveltavan kielitieteen tutkimuksia 2018*, (11):137–155.

Virtanen, Mikko T., Heidi Vepsäläinen, and Aino Koivisto. 2021. Managing several simultaneous lines of talk in finnish multi-party mobile messaging. *Discourse, Context & Media*, 39:100460–100474.

Wang, G. Alan, Harry Jiannan Wang, Jiexun Li, Alan S. Abrahams, and Weiguo Fan. 2014. An analytical framework for understanding knowledge-sharing processes in online Q&A communities. *ACM Transactions on Management Information Systems*, 5(4).

Weber, René, Michael J. Mangus, Richard Huskey, Frederic R. Hopp, Ori Amir, Reid Swanson, Andrew Gordon, Peter Khooshabeh, Lindsay Hahn, and Ron Tamborini. 2018. Extracting latent moral information from text narratives: Relevance, challenges, and solutions. *Communication Methods and Measures*, 12(2-3):119–139.

Wendler, Chris, Veniamin Veselovsky, Giovanni Monea, and Robert West. 2024. Do Llamas work in English? on the latent language of multilingual transformers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15366–15394, Bangkok, Thailand. Association for Computational Linguistics.

Willard, Brandon T. and Rémi Louf. 2023. Efficient guided generation for LLMs. *arXiv preprint arXiv:2307.09702*.

Wu, Hanqian, Lu Cheng, Jiahui Jin, and Feng Yuan. 2019. Dialog acts classification with semantic and structural information. In *2019 International Conference on Intelligent Computing, Automation and Systems (ICICAS)*, pages 438–442. IEEE.

Xiao, Yimin, Zong-Ying Slaton, and Lu Xiao. 2020. TV-AfD: An imperative-annotated corpus from the big bang theory and Wikipedia's articles for deletion discussions. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 6542–6548.

Xu, Guanghao, Hyunjung Lee, Myoung-Wan Koo, and Jungyun Seo. 2017. Convolutional neural network using a threshold predictor for multi-label speech act classification. In *2017 IEEE international conference on big data and smart computing (BigComp)*, pages 126–130. IEEE.

Yang, Diyi. 2021. 6 questions for socially aware language technologies. *Northern European Journal of Language Technology*, 7(1).

Zakharov, Stepan, Omri Hadar, Tovit Hakak, Dina Grossman, Yifat Ben-David Kolikant, and Oren Tsur. 2021. Discourse parsing for contentious, non-convergent online discussions. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 15, pages 853–864.

Zhang, Amy X., Bryan Culbertson, and Praveen Paritosh. 2017. Characterizing online discussion using coarse discourse sequences. In *Eleventh International AAAI Conference on Web and Social Media*, volume 11, pages 357–366.

Zhang, Justine, Cristian Danescu-Niculescu-Mizil, Christina Sauper, and Sean J. Taylor. 2018. Characterizing online public discussions through patterns of participant interactions. *Proceedings of the ACM on Human-Computer Interaction*, 2(CSCW):1–27.

Zhang, Renxian, Dehong Gao, and Wenjie Li. 2011. What are tweeters doing: Recognizing speech acts in twitter. In *Workshops at the Twenty-Fifth AAAI Conference on Artificial Intelligence*.

Zhou, Xinyi and Reza Zafarani. 2020. A survey of fake news: Fundamental theories, detection methods, and opportunities. *ACM Computing Surveys (CSUR)*, 53(5):1–40.

# A  Micro-F1s for A2 model

| Act. | Neg. | Pos. | Act. | Neg. | Pos. |
|---|---|---|---|---|---|
| quest. | 0.98 | 0.94 | chall. | 0.88 | 0.47 |
| req. | 0.94 | 0.61 | accept. | 0.96 | 0.63 |
| state. | 0.77 | 0.96 | denial | 0.86 | 0.49 |
| accus. | 0.89 | 0.67 | apprec. | 0.99 | 0.51 |

Table 15:    10-fold Micro-F1s for Annotator2 model. Neg.=Negative and Pos.=Positive class.

Table 15 presents Micro-F1 scores for the SetFit-FinBERT1-A2 model, one of our best performing models. This enables an examination of how the specific model performs on the positive and the negative class, respectively, for each action. Based on the scores, performances for the positive class ('yes the action is present in the comment') tend to be lower. Especially for some of the face-threatening actions (*denial*, *challenge*), as well as *appreciations*, the performance for the positive class is much lower than for the negative. This further supports our observation that people tend to express face-threatening actions in an ambiguous manner, perhaps to avoid direct face-threats. It is interesting, though, that *appreciations* also show lower performance for the positive class. This might have been due to appreciations being quite rare in the dataset.

| Model | Threshold | Action | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | question | request | statement | accusation | challenge | acceptance | denial | appreciation |
| FinBERT1-MS | $\geq 1$ | **0.87** | 0.68 | **0.73** | 0.63 | **0.61** | 0.64 | 0.55 | 0.60 |
| | $\geq 2$ | 0.86 | **0.69** | 0.72 | 0.61 | **0.61** | 0.64 | 0.54 | 0.59 |
| | $\geq 3$ | **0.87** | 0.68 | **0.73** | 0.76 | 0.60 | 0.62 | 0.54 | 0.63 |
| | $\geq 4$ | 0.86 | 0.64 | 0.69 | 0.61 | 0.59 | 0.60 | **0.56** | **0.65** |
| | $\geq 5$ | 0.85 | 0.60 | 0.66 | 0.54 | 0.56 | **0.66** | 0.53 | 0.64 |
| | $\geq 6$ | 0.81 | 0.54 | 0.62 | 0.55 | 0.58 | **0.66** | 0.64 | 0.50 |
| SetFit-FinBERT1-MS | $\geq 1$ | **0.94** | 0.79 | **0.77** | 0.71 | 0.59 | 0.85 | 0.56 | 0.62 |
| | $\geq 2$ | 0.93 | **0.82** | 0.76 | **0.73** | 0.63 | 0.84 | 0.65 | 0.58 |
| | $\geq 3$ | **0.94** | 0.79 | 0.75 | 0.71 | 0.60 | **0.86** | **0.68** | 0.65 |
| | $\geq 4$ | **0.94** | 0.73 | **0.77** | 0.69 | **0.64** | 0.81 | 0.64 | **0.78** |
| | $\geq 5$ | 0.90 | 0.73 | 0.72 | 0.54 | 0.50 | 0.66 | 0.59 | 0.75 |
| | $\geq 6$ | 0.79 | 0.75 | 0.64 | 0.55 | 0.50 | 0.69 | 0.60 | 0.50 |

Table 16: 10-fold cross-validated macro-F1 scores with different thresholds.

# B Threshold tests

To gain a better understanding of how threshold setting might affect action detection, we compared how well models performed when setting different threshold $\theta$ as a classification boundary. Here, we used the single annotations data, in particular the FinBERT1-MS and SetFit-FinBERT1-MS models. In this paper, we predefined our $\theta = 3$ for our experiments related to RQ1 and RQ2, as this allows us to focus on actions present at moderate to strong signal strength. However, here we wish to provide more insight into how thresholding might affect classification. Classifier training and testing follows the same procedures as described in section 5. See details specifically related to SetFit and the multilabel single annotation (MS) model configuration. It should be noted that with different $\theta$, the dataset imbalances will notably change (see Appendix C). Class weighting (FinBERT1-MS) and sentence pairing used for SetFit help to account for imbalances when measuring model performance. However, due to differences in data distributions here we note that the test presented in Table 16 is only suggestive. More comprehensive analysis would be needed for more reliable results.

Based on the test (Table 16), it seems that for most actions, especially those that invite responses (question, request, statement, accusation, challenge), there might be no great difference between scores when using $\theta$ between 1–3, or even 1–4 for some. It seems that a $\theta \geq 4$–6 (5–6 for some), might result in lower performance for many actions. It seems also that for responsive actions (and appreciations), a higher threshold might increase performance. This could perhaps be due to responsive actions being expressed in in a shorter or even more ambiguous manner, e.g. in a subordinate clause or shortly at the beginning of the comment, more emphasis being given to other actions in the comment.
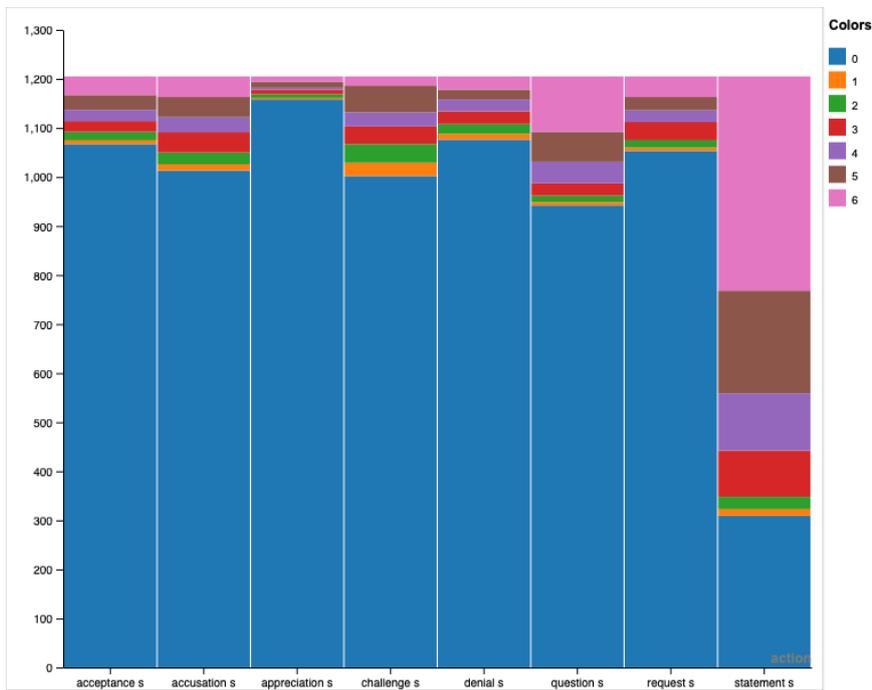
# C  Label distributions



Figure 2: Distribution of labels in s=single annotations. Colors indicate label scores.
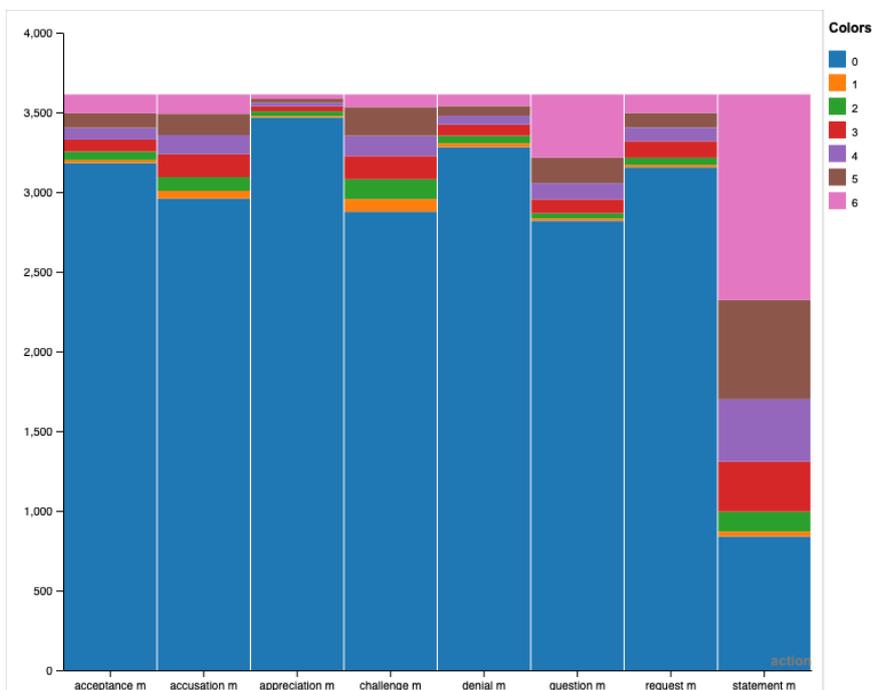


Figure 3: Distribution of labels in m=multiple annotations. Colors indicate label scores.

| Dataset | Action | Label score | | | | | | | Total |
|---|---|---|---|---|---|---|---|---|---|
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | (positive label) |
| single | question | 942 | 7 | 13 | 25 | 44 | 61 | 113 | 263 |
| | request | 1053 | 7 | 15 | 37 | 24 | 27 | 42 | 152 |
| | statement | 309 | 14 | 25 | 94 | 117 | 209 | 437 | 896 |
| | accusation | 1013 | 13 | 24 | 41 | 31 | 41 | 42 | 192 |
| | challenge | 1002 | 27 | 38 | 36 | 29 | 54 | 19 | 203 |
| | acceptance | 1067 | 8 | 17 | 22 | 23 | 30 | 38 | 138 |
| | denial | 1076 | 12 | 20 | 26 | 24 | 20 | 27 | 129 |
| | appreciation | 1158 | 3 | 8 | 9 | 4 | 12 | 11 | 47 |
| | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | (positive label) |
| multiple | question | 2819 | 15 | 33 | 87 | 102 | 160 | 396 | 793 |
| | request | 3155 | 15 | 45 | 104 | 86 | 91 | 116 | 457 |
| | statement | 839 | 31 | 125 | 314 | 392 | 622 | 1289 | 2773 |
| | accusation | 2959 | 48 | 86 | 144 | 120 | 134 | 121 | 653 |
| | challenge | 2877 | 79 | 124 | 146 | 127 | 180 | 79 | 735 |
| | acceptance | 3182 | 19 | 53 | 75 | 74 | 91 | 118 | 430 |
| | denial | 3282 | 24 | 47 | 70 | 55 | 62 | 72 | 330 |
| | appreciation | 3467 | 10 | 29 | 31 | 23 | 25 | 27 | 145 |

Table 17: Distributions of action labels by score in single annotations and multiple annotations. The label distributions are quite similar for both single annotations and multiple annotations. Overall, scores 3-6 are more commonly used than 1-2, although there are some differences between actions. For *challenges*, scores 1-2 seem a bit more common than for other actions. These are small differences, though. *Statements* have a positive score ($\geq 1$) much more often than other actions, highlighting the assertive nature of the comments in our data. *Questions*, *challenges* and *accusations* also appear to be more commonly labeled with a positive score.