# Appendix

## A   Annotation Environment and Instruction Examples

In this section, we present examples of annotation environments for the text-image relation task, sentiment analysis and NER, as well as introduce the general annotation instructions given to annotators.

### A.1   Text-image Relation Annotation

Figure 1 shows an example of how the annotators were presented with text-image pairs from tweets. Annotators were instructed to click one of the following: 1) the green button in cases when what was written in the tweet text was represented in the image AND the image added to the meaning of the text; 2) the grey button when what was written in the tweet text was NOT represented in the image, BUT the image added to the meaning of the text; 3) the purple button when what was written in the tweet text was represented in the image BUT the image did NOT add to the meaning of the text, and 4) the red button when NEITHER when what was written in the tweet text was represented in the image NOR the image added to the meaning of the text. The top part of the annotation environment summarises the annotation statistics of the current annotator.
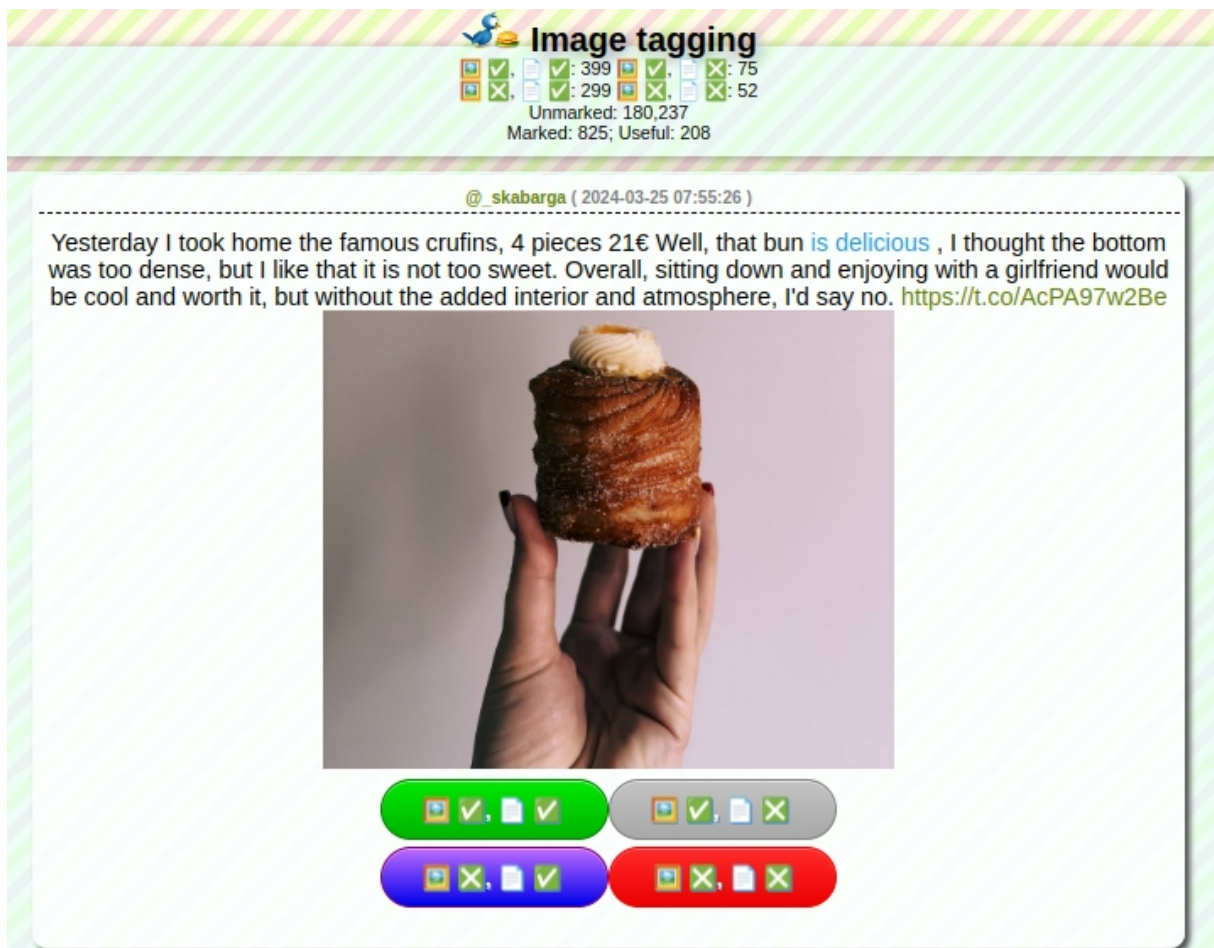


Figure 1: An example of the text-image relation annotation environment.

### A.2   Named Entity Annotation

Figure 2 shows an example of how the annotators were presented with tweets for named entity annotation. Tweets are presented each on a separate line, with 10 tweets per screen. The task of annotators was to select appropriate text spans and assign a named entity category to the selected span. The named entity categories are already set, and the annotator needs to select correct category from drop-down menu.
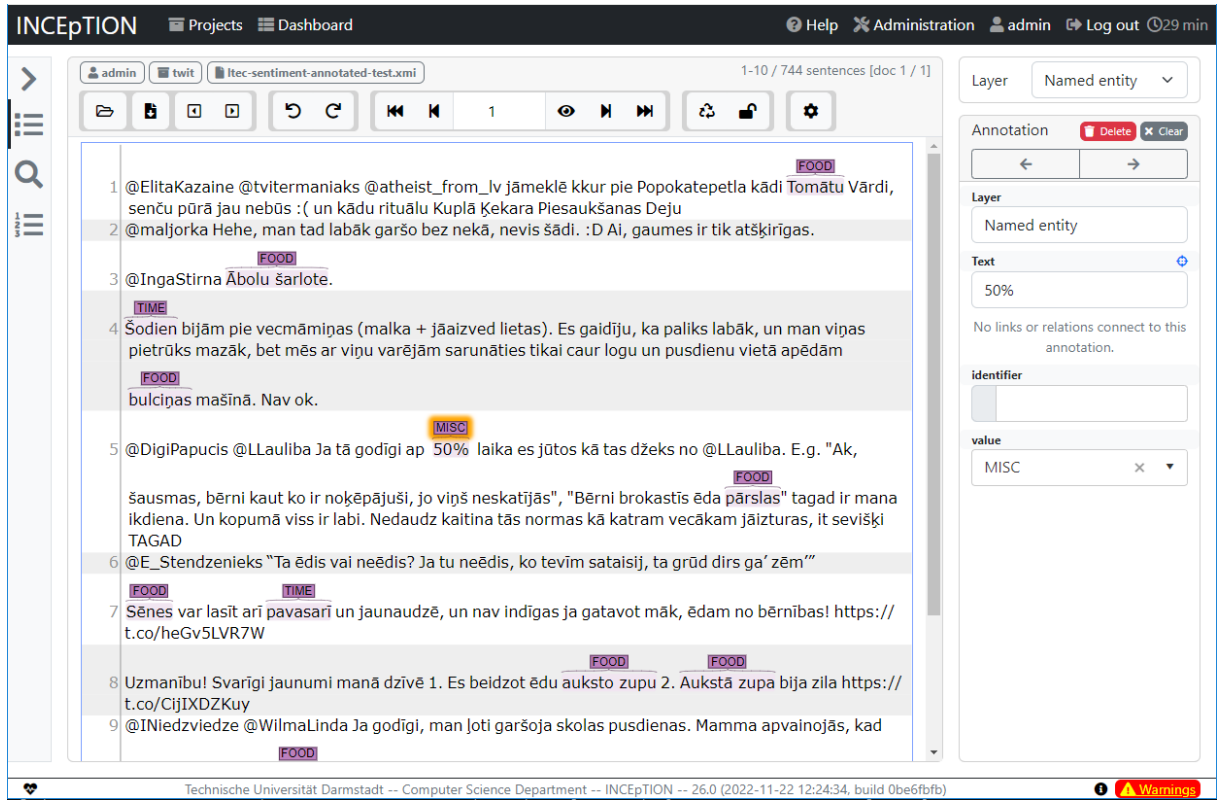
Figure 2: An example of the named entity annotation environment.

### A.3 Sentiment Annotation

Figure 3 shows an example of how the annotators were presented with tweets for sentiment annotation. Tweets were displayed one by one and annotators were given instructions to click either the red 'Negative', the grey 'Neutral', or the green 'Positive' button as an indication of which sentiment they believe the tweet text represents. The bottom part of the annotation environment summarises the annotation statistics of the current annotator.

## B Annotator Profile

Annotators for the named entity recognition (NER) and the text-image relation task were native Latvian speakers in their 20s-30s with at least a master's degree and work experience in the field of natural language processing (NLP) or linguistics. For sentiment analysis, annotators were in the age ranges between 20 and 50 with at least a bachelor's degree in various fields, and the translator and post-editor were professional linguists with a master's degree in their respective fields in their 30s-40s.

## C Text-image Relation Ablation Experiments

Since the LLaVA models can be sensitive to the provided random seed, we performed additional ablation experiments with 10 different random seeds as well as the latest available 1.5 version, and both 7B and 13B model sizes. Results summarised in Table 1 show that the scores are slightly lower than in the original experiments with the 1.3 version.
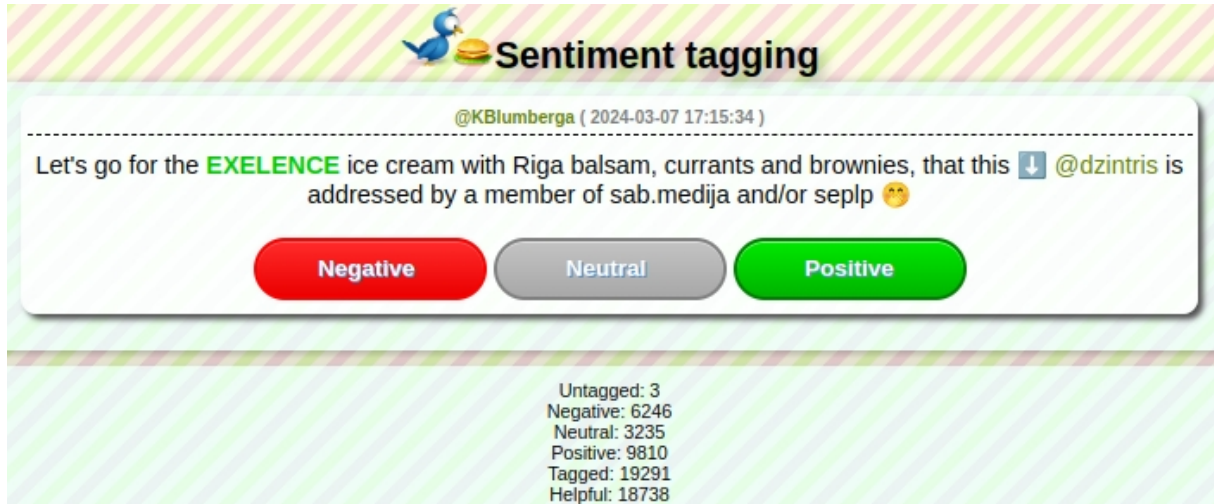
Figure 3: An example of the sentiment annotation environment.

|         | Average | +/-    |
|---------|---------|--------|
| EN-13b  | 19.78%  | 0.86%  |
| EN-7b   | 22.13%  | 2.46%  |
| LV-13b  | 12.54%  | 1.29%  |
| LV-7b   | 17.82%  | 1.60%  |

Table 1: Text-image relation classification ablation experiment results with Llava 1.5 models.