# PAP: A Dataset for Physical and Abstract Plausibility and Sources of Human Disagreement

Annerose Eichel and Sabine Schulte im Walde

Institute for Natural Language Processing, University of Stuttgart, Germany

{annerose.eichel,schulte}@ims.uni-stuttgart.de

**AMLaP 2023**

Universität Stuttgart

## INTRODUCTION

- Discerning plausible from implausible events: crucial building block for NLP
- Previous work mostly focused on semantic knowledge to distinguish
  - *physically* plausible vs. implausible events
  - events with mostly conceptually *concrete* participants



### RESEARCH GOALS & CONTRIBUTIONS

- Create novel dataset for **physical and abstract plausibility** of events in English, capturing abstractness to the same extent as concreteness for the first time
- Systematically examine **plausibility across levels of abstractness**
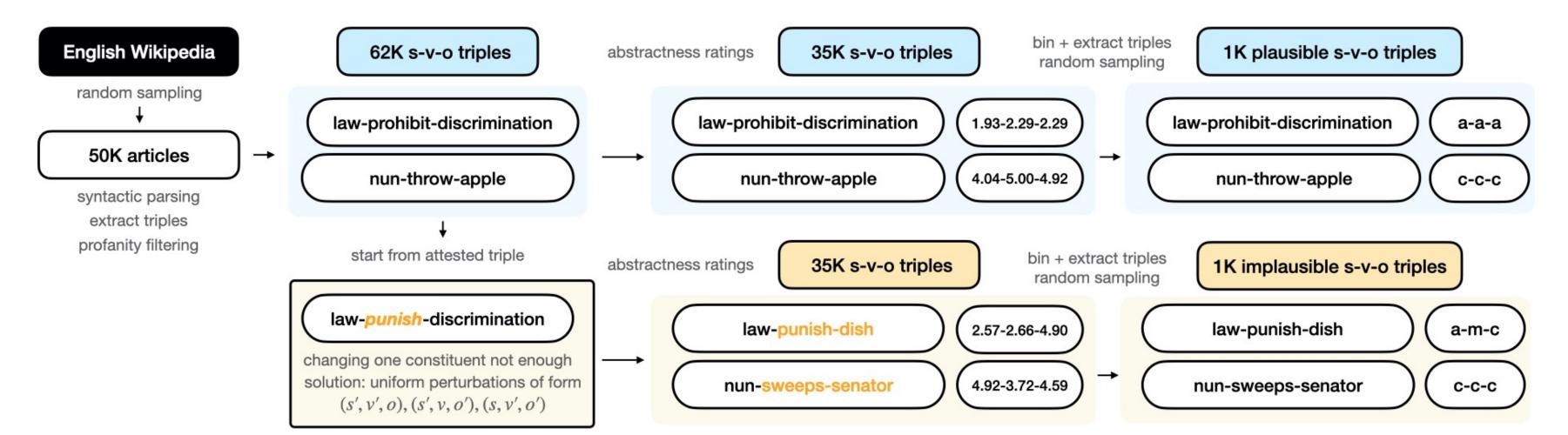- Explore and represent **disagreement in plausibility annotation**

## CAPTURING (SEMANTIC) PLAUSIBILITY

### PLAUSIBILITY

- Captures non-surprial in a given context
  *child-sleep* vs. *tree-sleep*
- Includes both what is preferred (and probably most plausible) and what is unusual (but still very much plausible),
  *child-eat-banana* vs. *child-eat-pebble*
  → in contrast to selectional preference / thematic fit
- Can be estimated as a matter of degree with events assessed corresponding to perceived plausibility
  *child-eat-banana* vs. *child-eat-pebble* vs. *child-eat-skyscraper*
- Denotes what is likely in a given world but not necessarily attested in a given corpus
  *human-dies* vs. *human-breathes*

## CONSTRUCTING EVENT TARGETS



Simplified Illustration of Dataset Construction

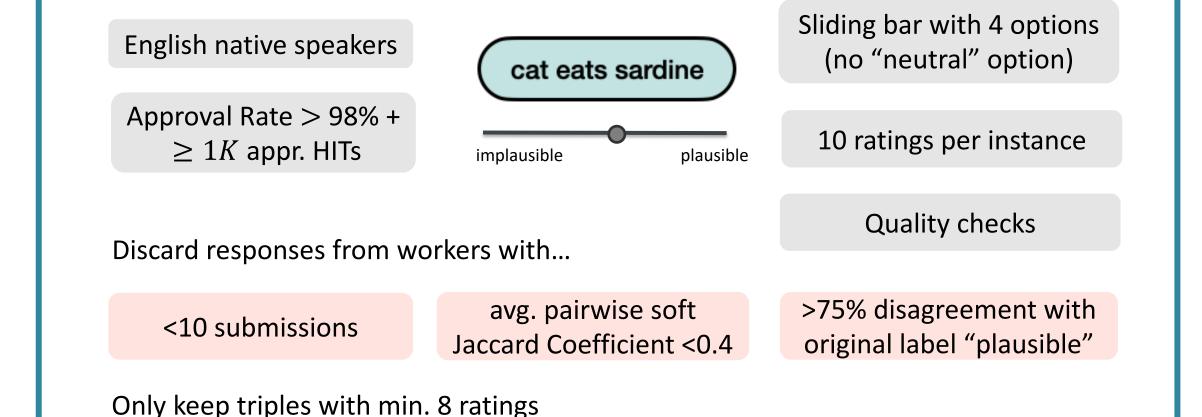### PLAUSIBLE EVENTS (marked in blue)

- From English Wikipedia sample: Extract attested events, filter for profanity, assign abstractness ratings, bin according to abstractness, and sample 1,080 plausible events

### (PSEUDO-)IMPLAUSIBLE EVENTS (marked in yellow)

- Based on extracted attested triples:
  (i) Automatically generate pseudo-implausible events by perturbating event constituents
  (ii) Construct 1,080 pseudo-implausible event similarly to plausible event construction

## COLLECTING HUMAN ANNOTATIONS

**TASK:** Collect plausibility judgements on AMT for 2,160 plausible and implausible triples
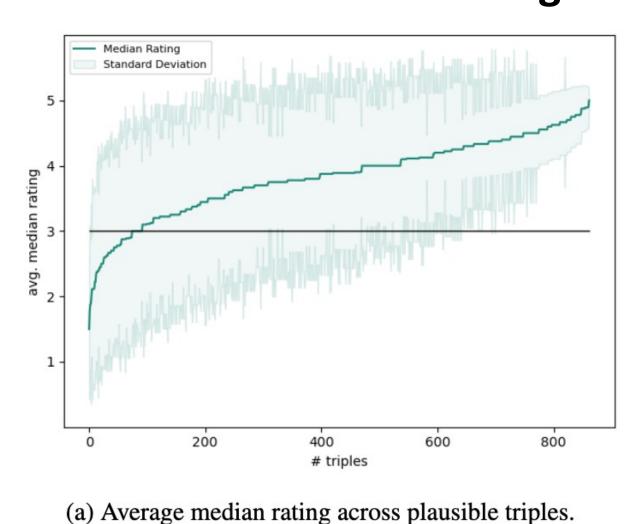


### DATASET STATISTICS

- **15,571 plausibility ratings for 1,733 triples**
- Ø IAA: Soft Jaccard Coefficient of 0.64
  → reasonable agreement among annotators with indication of disagreement to be examined

## ANALYSIS OF HUMAN JUDGEMENTS AND DISAGREEMENT

### What can we learn from rating distributions?



(a) Average median rating across plausible triples.

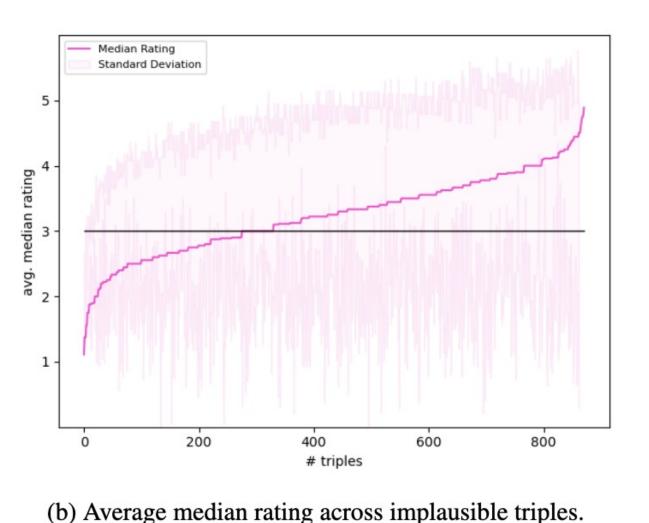(b) Average median rating across implausible triples.

Figure 4: Average median ratings across originally plausible (a) and implausible (b) triples with standard deviation visualized as cloud around average rating lines. Triples are represented numerically on the x-axis. The black horizontal line denotes a median rating of 3. Average median ratings for *plausible* triples *below* the line disagree with the original label, while the opposite is true for average median ratings for *implausible* triples. Here, ratings *above* the line disagree with the original label.

(i) Humans tend to favor plausibility over implausibility, while avoiding the extreme on the plausibility end of the scale.
(ii) Implausibility yields higher disagreement as annotators disagree more when rating triples originally labelled implausible.
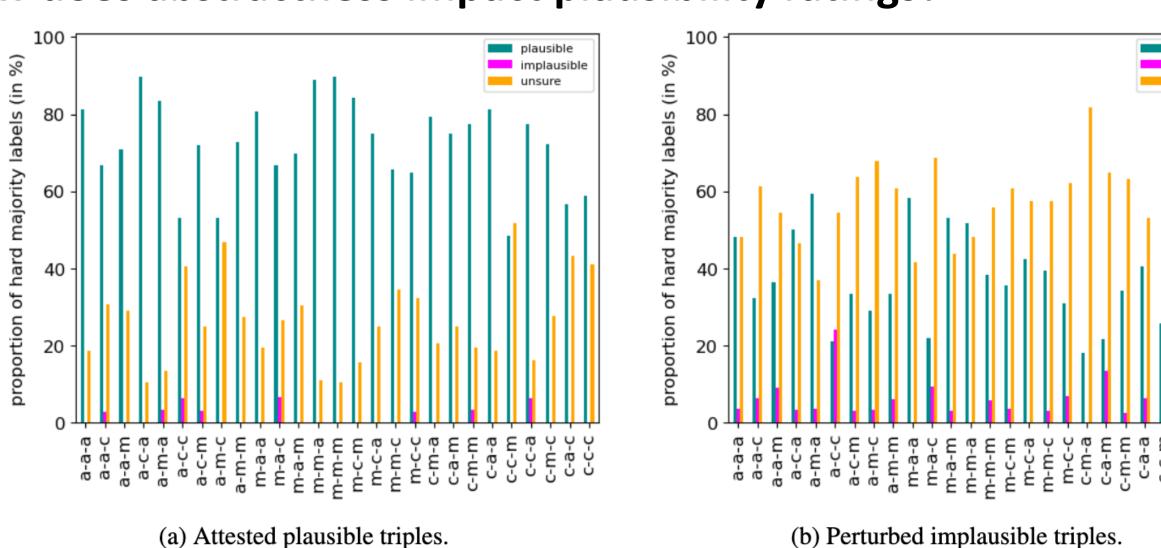
### How does abstractness impact plausibility ratings?



(a) Attested plausible triples.

(b) Perturbed implausible triples.

Figure 5: Proportion of strict majority ratings (⩾70%) across abstractness combinations for attested plausible triples (a) and perturbed implausible triples (b). Green bars denote a majority of plausible ratings ∈ {4, 5}, pink bars refer to a majority of implausible ratings ∈ {1, 2}, and orange bars capture cases of no clear majority.

(i) Plausibility tends to be more likely to be assigned in case of more abstract event participants.
(ii) Implausibility seems to be easier to capture with conceptually concrete words – no matter the underlying original label.

## ACKNOWLEDGEMENTS

## CONCLUSIONS

- Presented a novel human-annotated dataset for physical and abstract plausibility for events in English
- Explored relationship between abstractness and plausibility and analyzed annotator disagreement
- Released both raw and a range of aggregated annotations to foster research on (semantic) plausibility and related notions, disagreement, and relevant downstream tasks such as commonsense reasoning

Scan and use PAP