

## INTRODUCTION

In light of large, pre-trained language models (LMs) nearing and surpassing human performance on a breadth of language understanding tasks [1, 2, 3], we propose Tiered Reasoning for Intuitive Physics (TRIP), a more challenging evaluation targeting physical commonsense in a densely annotated, tiered reasoning setting:

### Story A

- Ann sat in the chair.
- Ann unplugged the telephone.
- Ann picked up a pencil.
- Ann opened the book.
- Ann wrote in the book.

### Story B

- Ann sat in the chair.
- Ann unplugged the telephone.
- Ann picked up a pencil.
- Ann opened the book.
- Ann heard the telephone ring.

Which story is more plausible? A

Why not B?

Conflicting sentences: 2 → 5

Physical states:

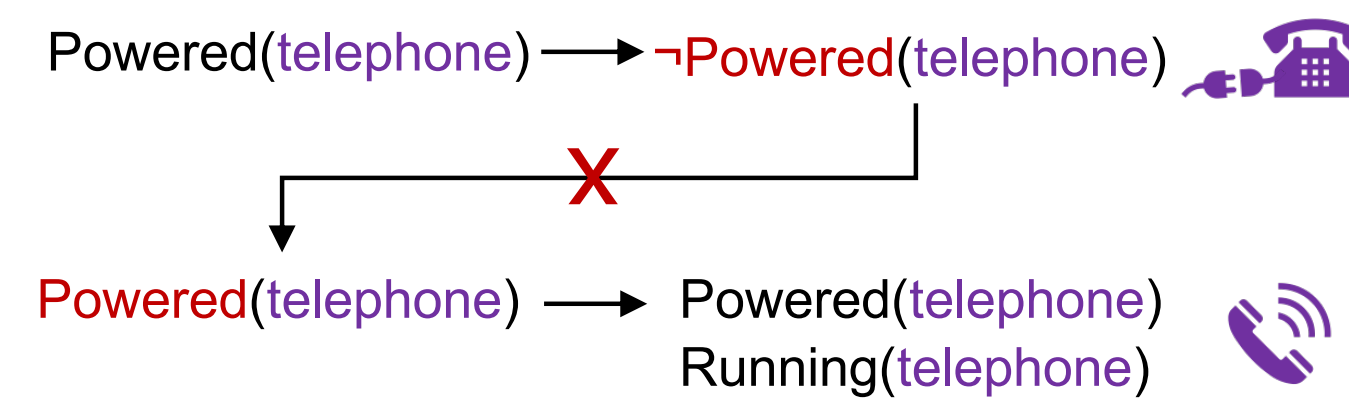


Figure 1: Story pair from TRIP, along with the tiers of annotation available to represent the reasoning process.

## BASELINE APPROACH

We propose a tiered architecture powered by large, pre-trained LMs and their contextual embeddings:

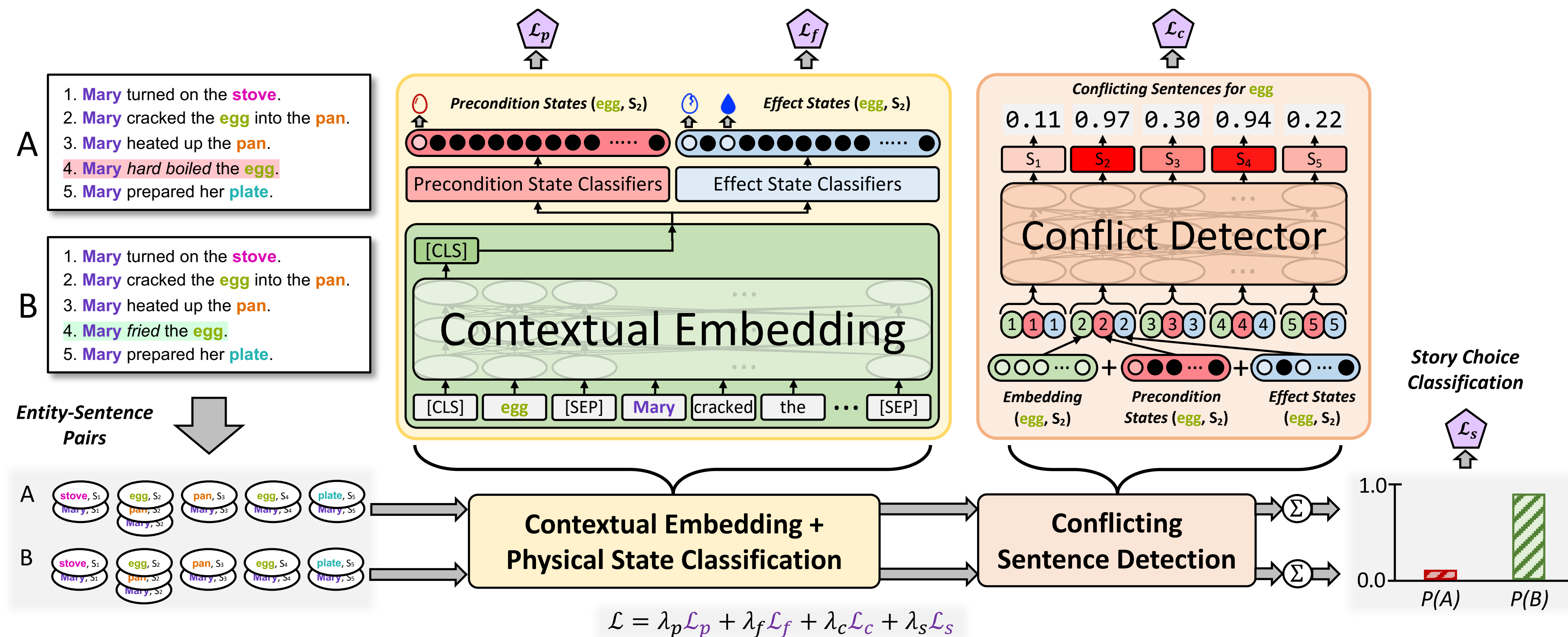


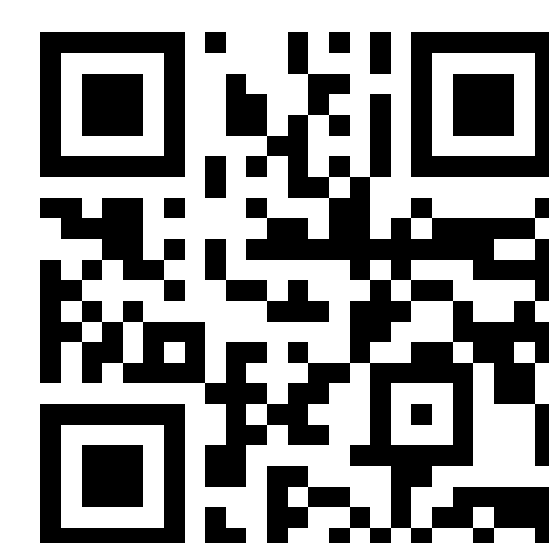
Figure 2: Proposed tiered reasoning system that generates contextual embeddings of entity-sentence pairs using a pre-trained LM, then uses them to jointly predict physical states, detect conflicting sentences, and identify the plausible story. The model is trained end-to-end by optimizing a weighted sum of cross-entropy loss functions  $\mathcal{L}_p$  for precondition state classification,  $\mathcal{L}_f$  for effect state classification,  $\mathcal{L}_c$  for conflicting sentence detection, and  $\mathcal{L}_s$  for story choice classification.

## REFERENCES

- [1] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of NAACL HLT 2019*, 2019.
- [2] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *arXiv:1907.11692*, 2019.
- [3] P. He, X. Liu, J. Gao, and W. Chen. DeBERTa: Decoding-enhanced BERT with Disentangled Attention. *arXiv:2006.03654*, 2021.

## LINKS

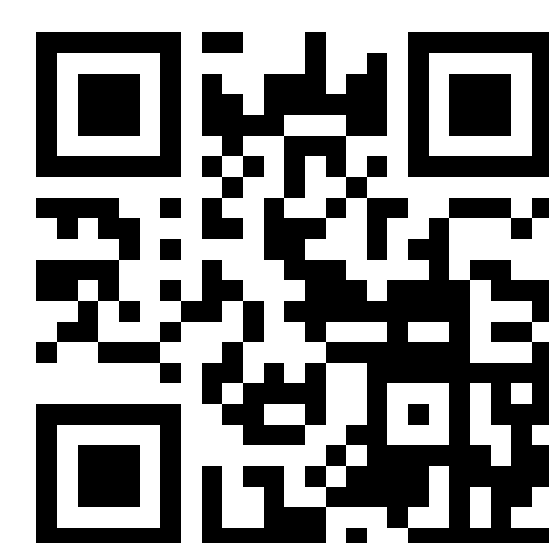
Scan the QR codes for related websites:



PREPRINT



CODE



LAB SITE

## BASELINE RESULTS

We evaluate systems with three metrics:

- Accuracy:** requires story choice to be correct.
- Consistency:** additionally requires conflicting sentences in the implausible story to be correct.
- Verifiability:** additionally requires some physical states to be predicted for the conflicting sentences, and all predicted states must be correct.

Model	Accuracy (%)	Consistency (%)	Verifiability (%)
random	47.8	11.3	0.0
<i>All Losses</i>			
BERT	78.3	2.8	0.0
ROBERTA	75.2	6.8	0.9
DeBERTa	74.8	2.2	0.0
<i>Omit Story Choice Loss <math>\mathcal{L}_s</math></i>			
BERT	73.9	28.0	9.0
ROBERTA	73.6	22.4	10.6
DeBERTa	75.8	24.8	7.5

Table 1: End and tiered task metrics for tiered classifiers on the validation set of TRIP trained on varied combinations of loss functions. Random baseline averaged over 10 runs.

## REASONING BREAKDOWN

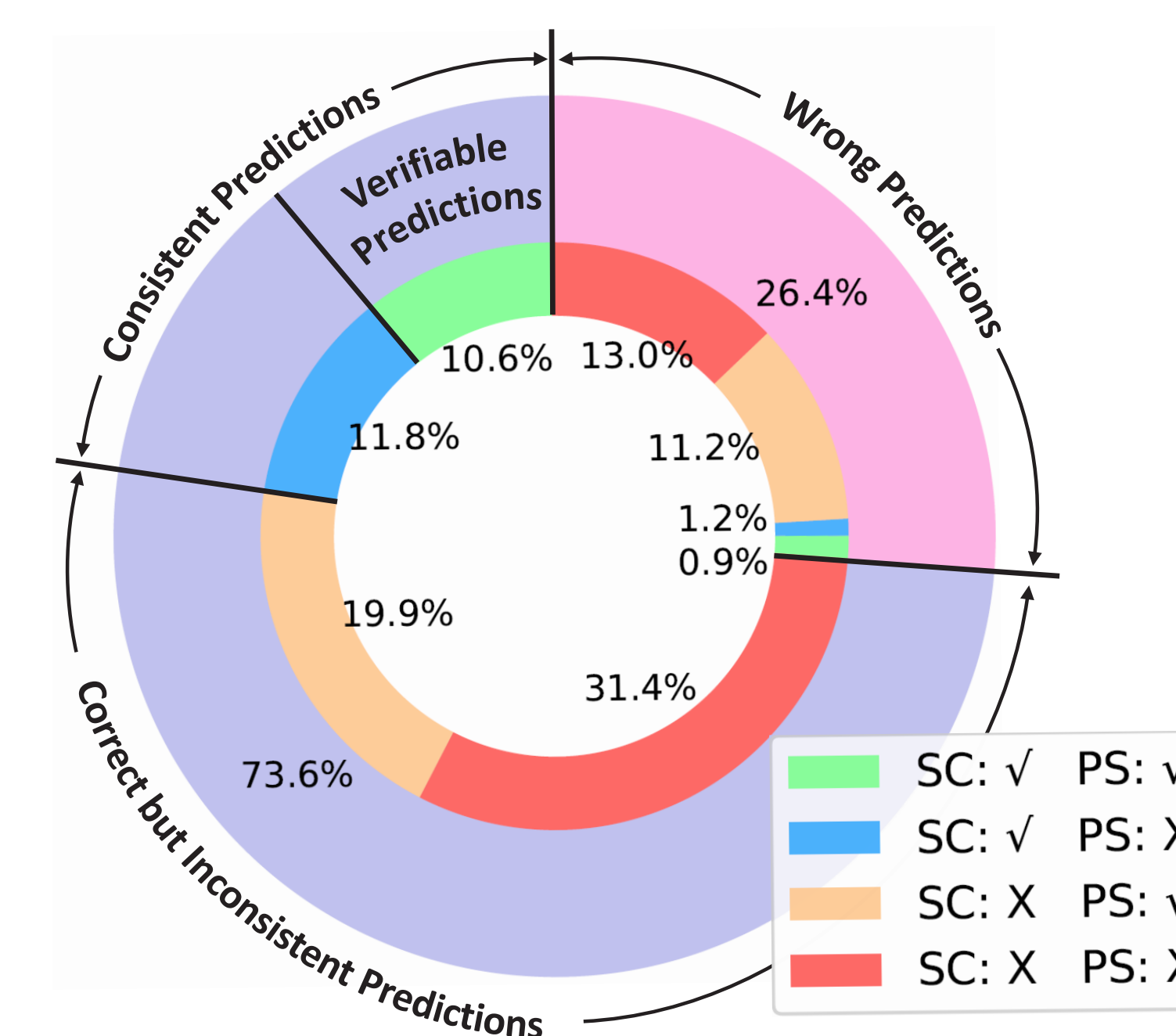
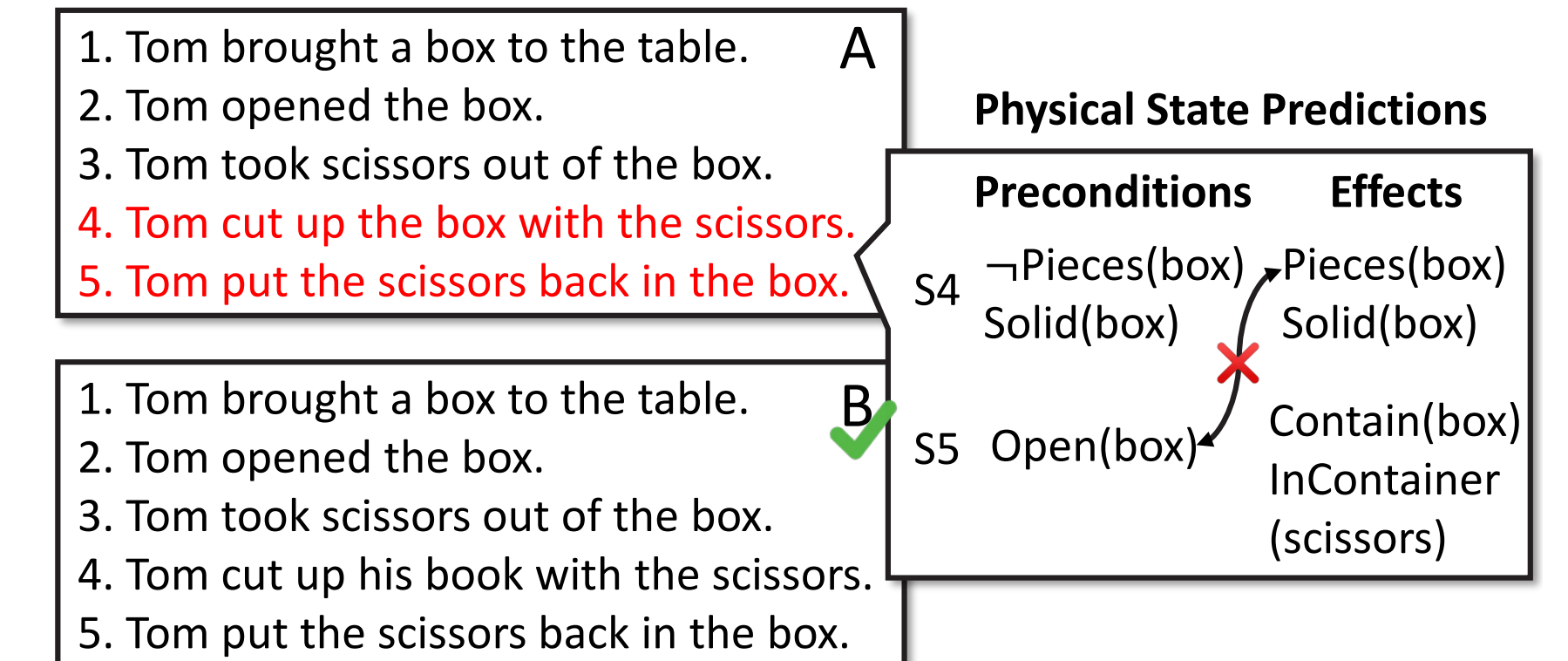


Figure 4: Distribution of ROBERTA successes and failures on TRIP. SC (sentence conflict) and PS (physical state) denote whether the predicted conflicting sentences or physical states are correct ( $\checkmark$ ) or not ( $\times$ ).

## SAMPLE OUTPUTS

A verifiable prediction.



A consistent but not verifiable prediction.

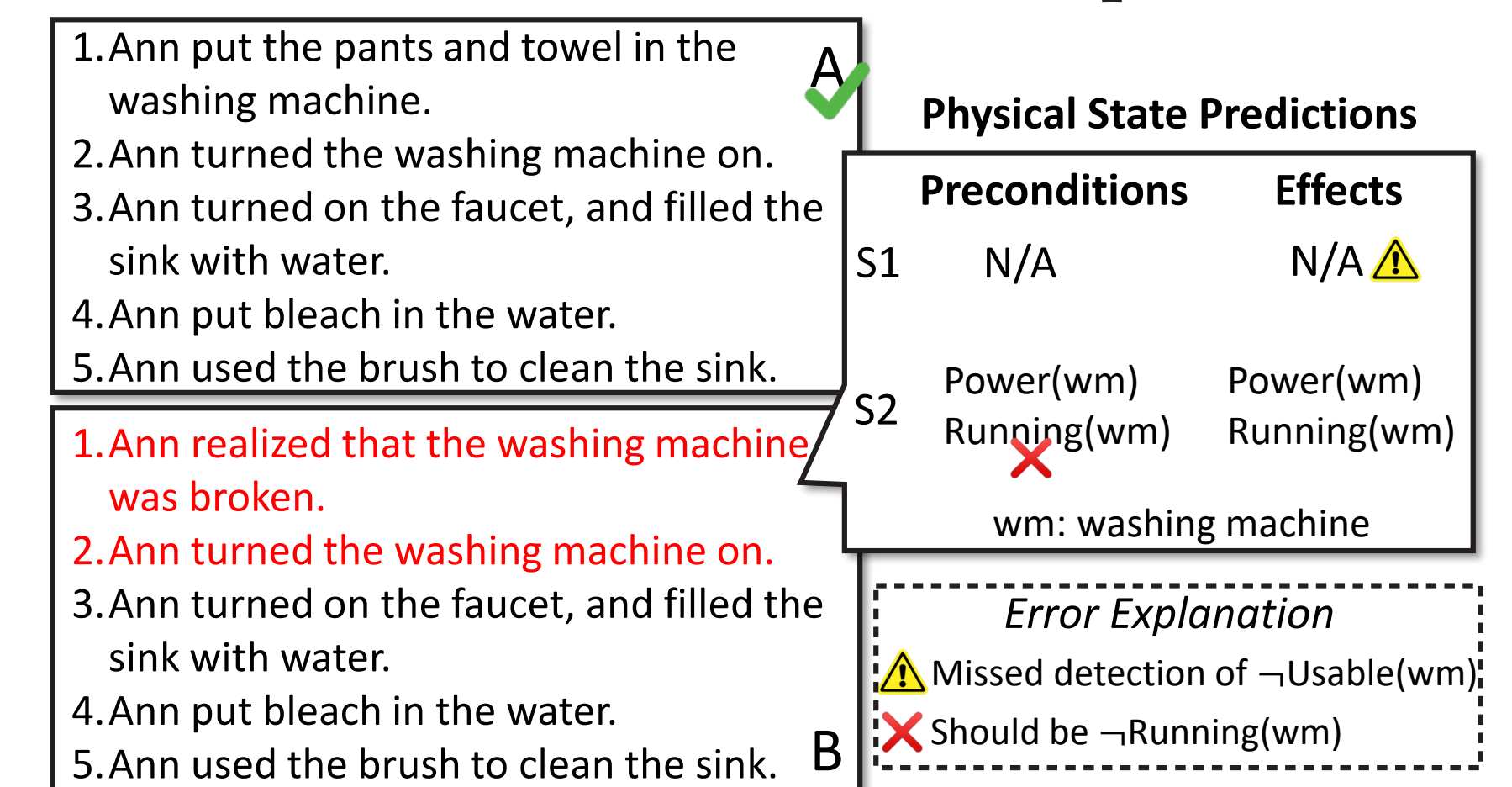


Figure 3: Sample outputs from the baseline system. The detected conflicting sentences are in red, and physical state predictions are shown on the right.

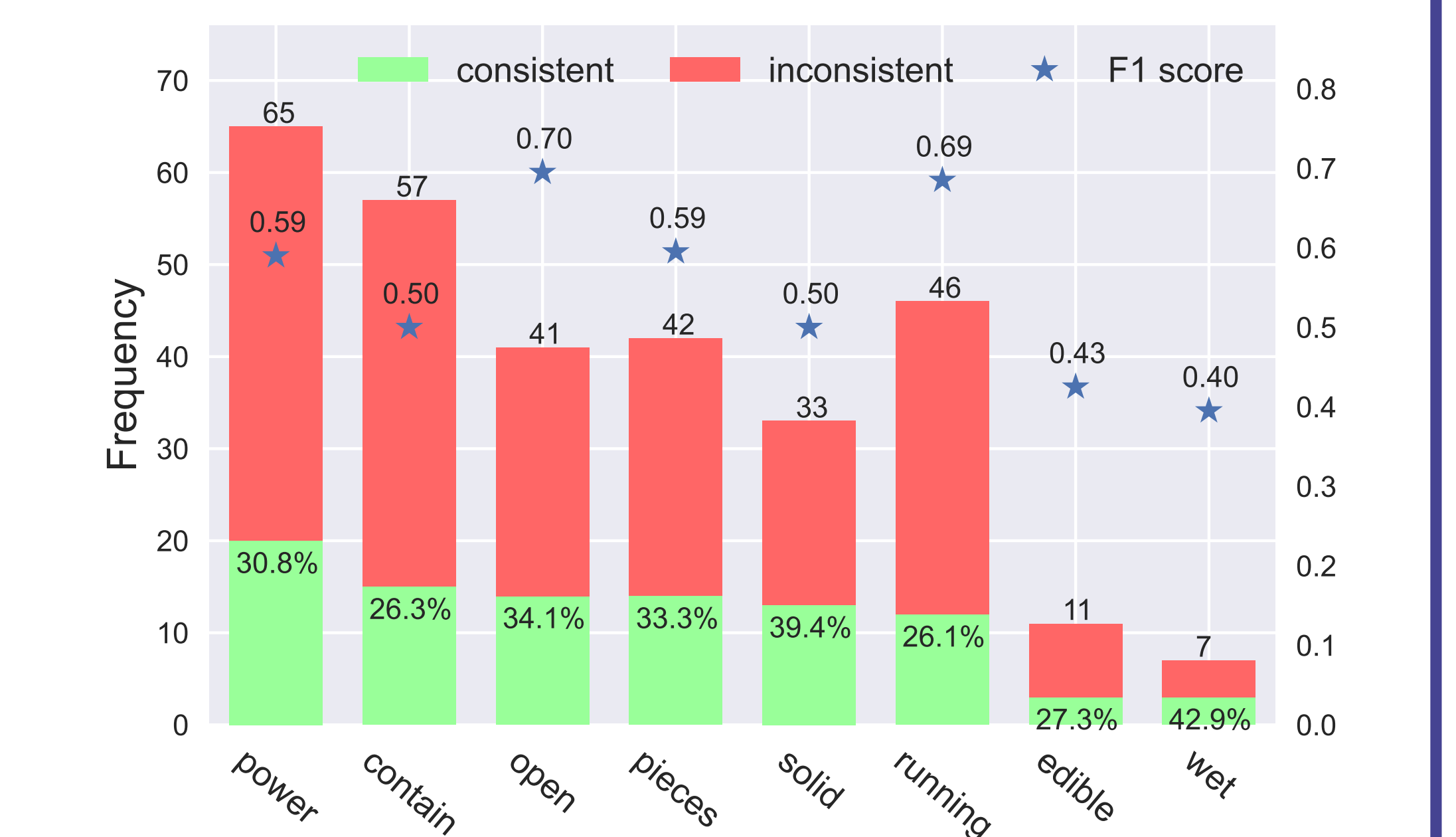


Figure 5: Utility of physical state predictions for selected attributes. Among correctly predicted physical states, bar regions indicate how many contribute to consistent end task predictions (i.e., with successfully detected conflicts). Blue stars indicate macro-F1 score of state prediction.

## CONCLUSION

Our results show that *supervising large LMs based on high-level classification tasks in order to learn commonsense language understanding leads to inconsistent and unverifiable reasoning*. In order to solve tasks like these coherently, we should directly train systems to incorporate multiple types of lower-level evidence. Our work provides an important first step toward this goal and strong intuition for future progress.