

Data and Objectives

Expert annotators perform post-hoc analyses to assign ratings on the Scale for Assessment of Thought, Language, and Cognition (TLC) to transcripts from participants responding to language tasks.

We combined annotated transcripts from four studies:

	Healthy Volunteers	Any Psychiatric Disorder (PD)	
n=640	279 (43.6%)	361 (56.4%)	29.9
Mean Age		26.3	
Gender			
Man	111(17.3%)	214 (33.4%)	
Woman	147 (23%)	132 (20.6%)	
Non-binary	20 (3%)	10 (6.4%)	
Unknown	1 (0.16%)	5 (3.2%)	
Diagnosis			
Schizophrenia		0 158 (24.7%)	
Bipolar + Psychosis		0 48 (7.5%)	
Unspecified PD		0 48 (7.5%)	
Schizoaffective		0 38 (5.9%)	
Schizoaffective-BT		0 23 (3.6%)	
Schizoaffective-DT		0 16 (2.5%)	
Schizophreniform		0 14 (2.2%)	
MDD + Psychosis		0 14 (2.2%)	
Brief PD		0 1 (0.16%)	
Substance-induced PD		0 1 (0.16%)	
None	279 (43.6%)		0

- We investigate the capability of LLMs to assign TLC ratings
- We discern the linguistic tasks for which LLMs make the best symptom assessments.
- We discern which TLC variables LLMs can accurately evaluate.

References

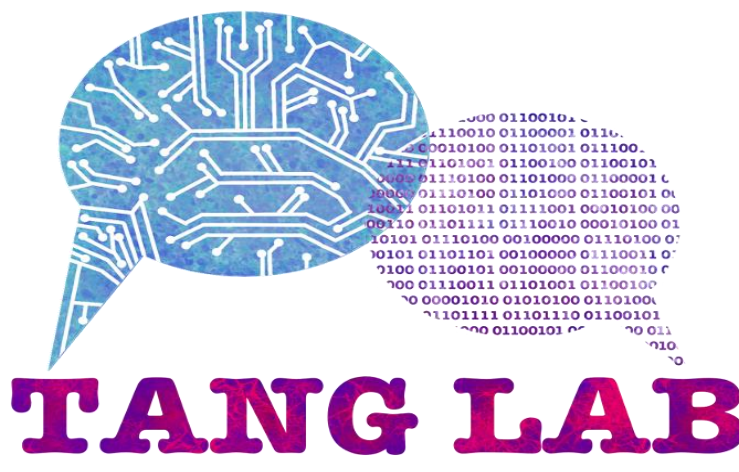
¹Nancy C. Andreasen, Scale for the Assessment of Thought, Language, and Communication (TLC), Schizophrenia Bulletin, Volume 12, Issue 3, 1986, Pages 473–482, <https://doi.org/10.1093/schbul/12.3.473>

Institutions/Disclosures

Author Zak Singh is at Cambridge University, all other authors are at the Feinstein Institutes. SXT received research funding from Winterlight Labs and holds equity with North Shore Therapeutics. She is also a consultant for both entities. She is on the advisory board for Psyrin, and serves as a consultant for Catholic Charities Neighborhood Services and LB Pharmaceuticals.

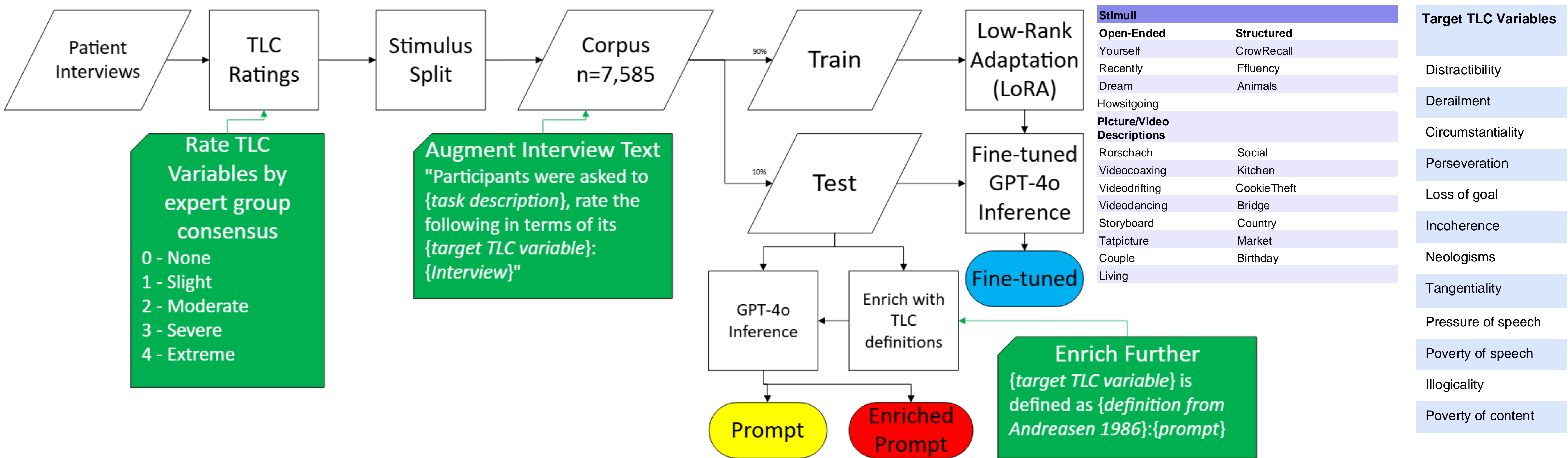
Adapting a Large Language Model (LLM) to Assess Clinical Ratings of Thought Disorder in Psychosis

Ryan Partlan, Simran Bhola, Sandy Yin, Zak Singh*, Sunny Tang



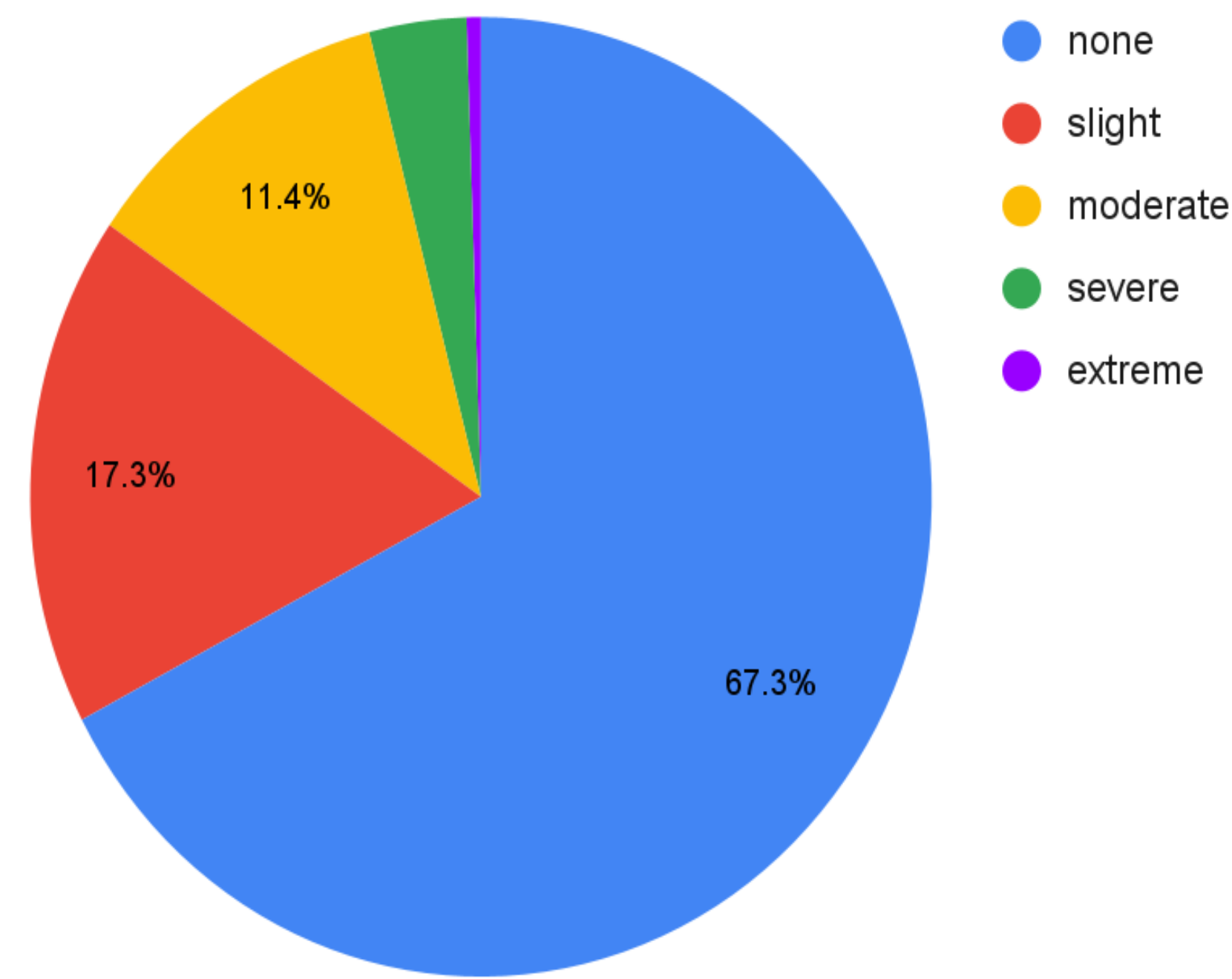
Contact me!
rpartlan@northwell.edu

Methods

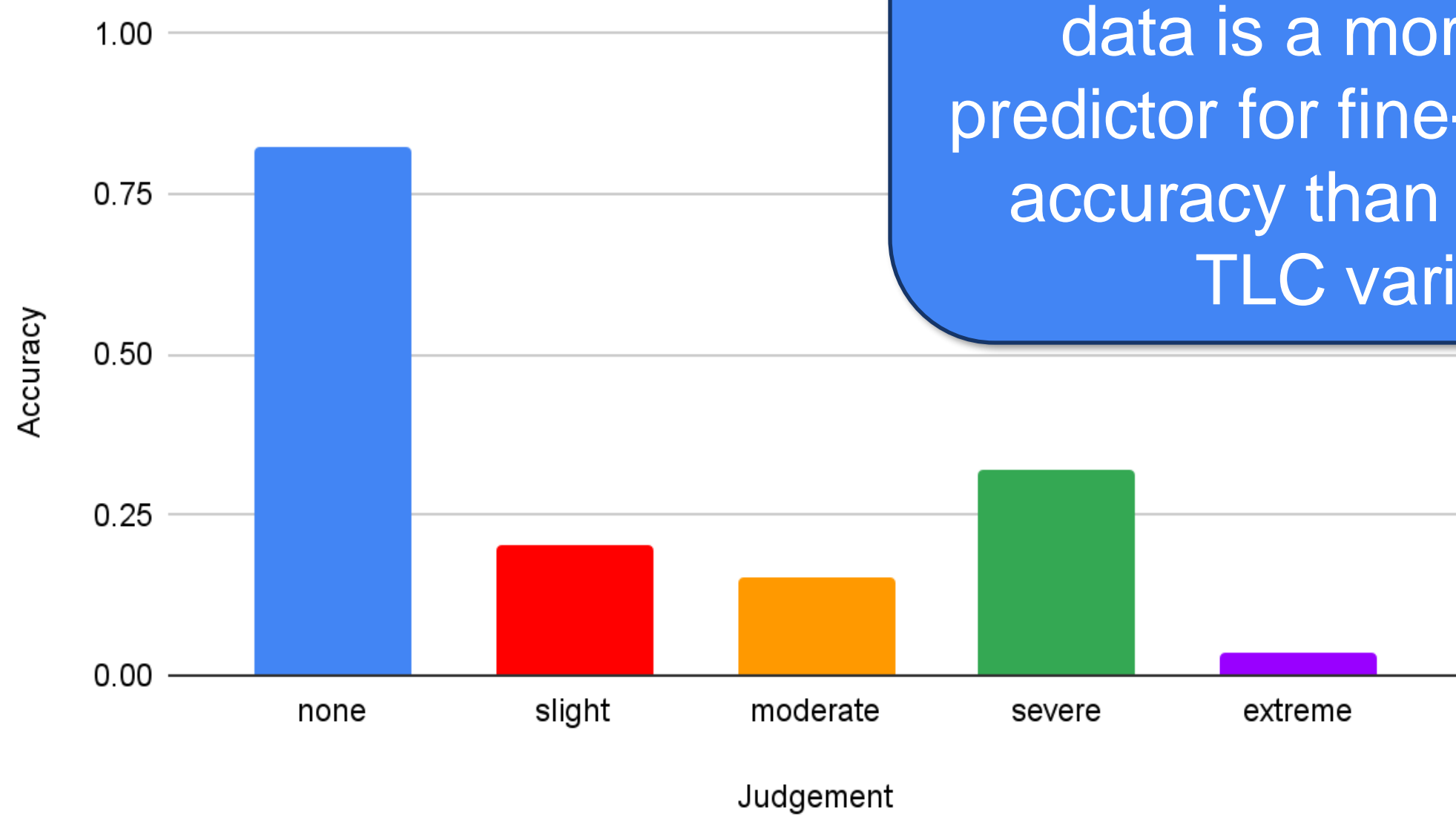


Results

Training Set Composition

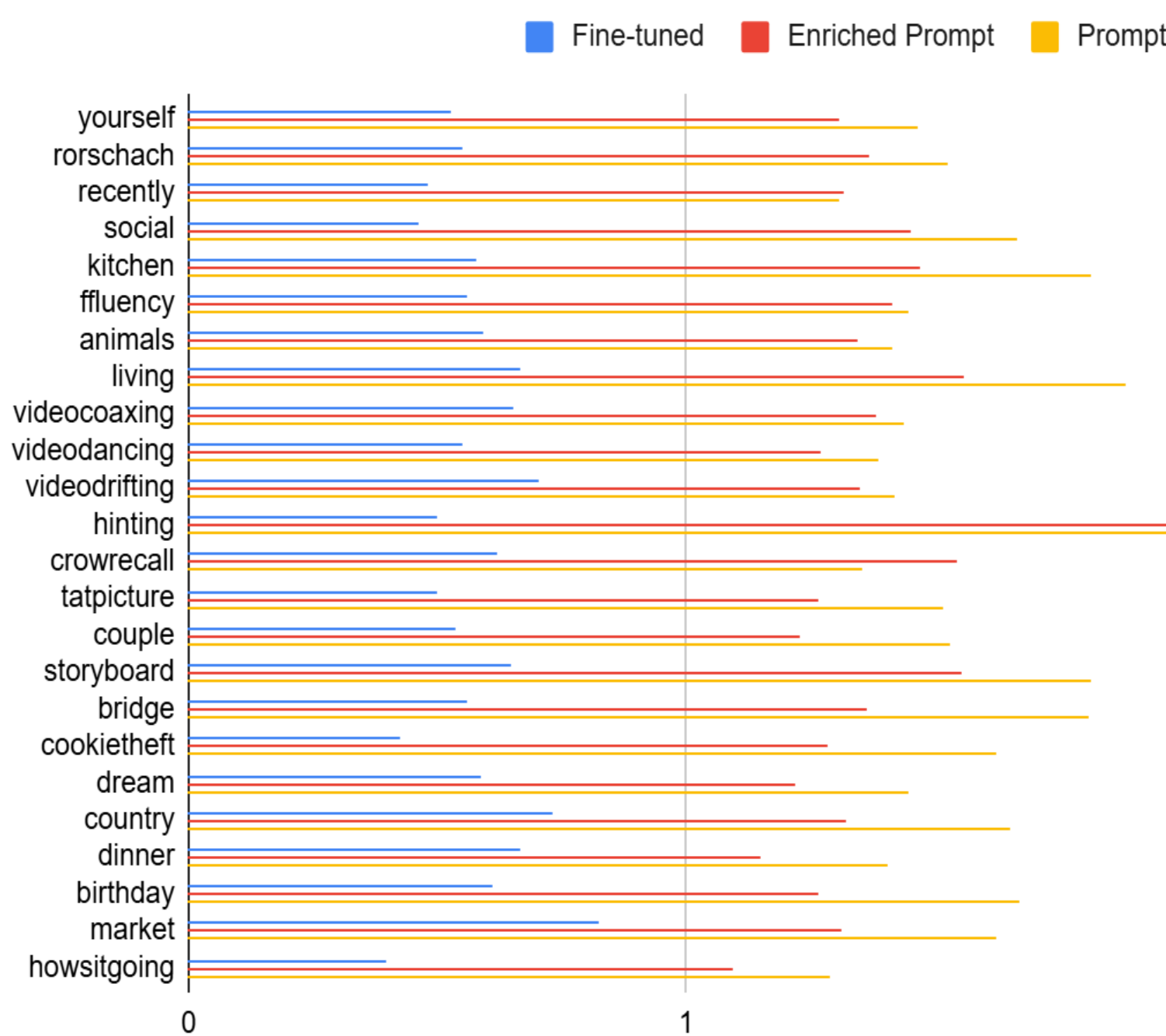


Fine-tuned accuracy vs. Judgement



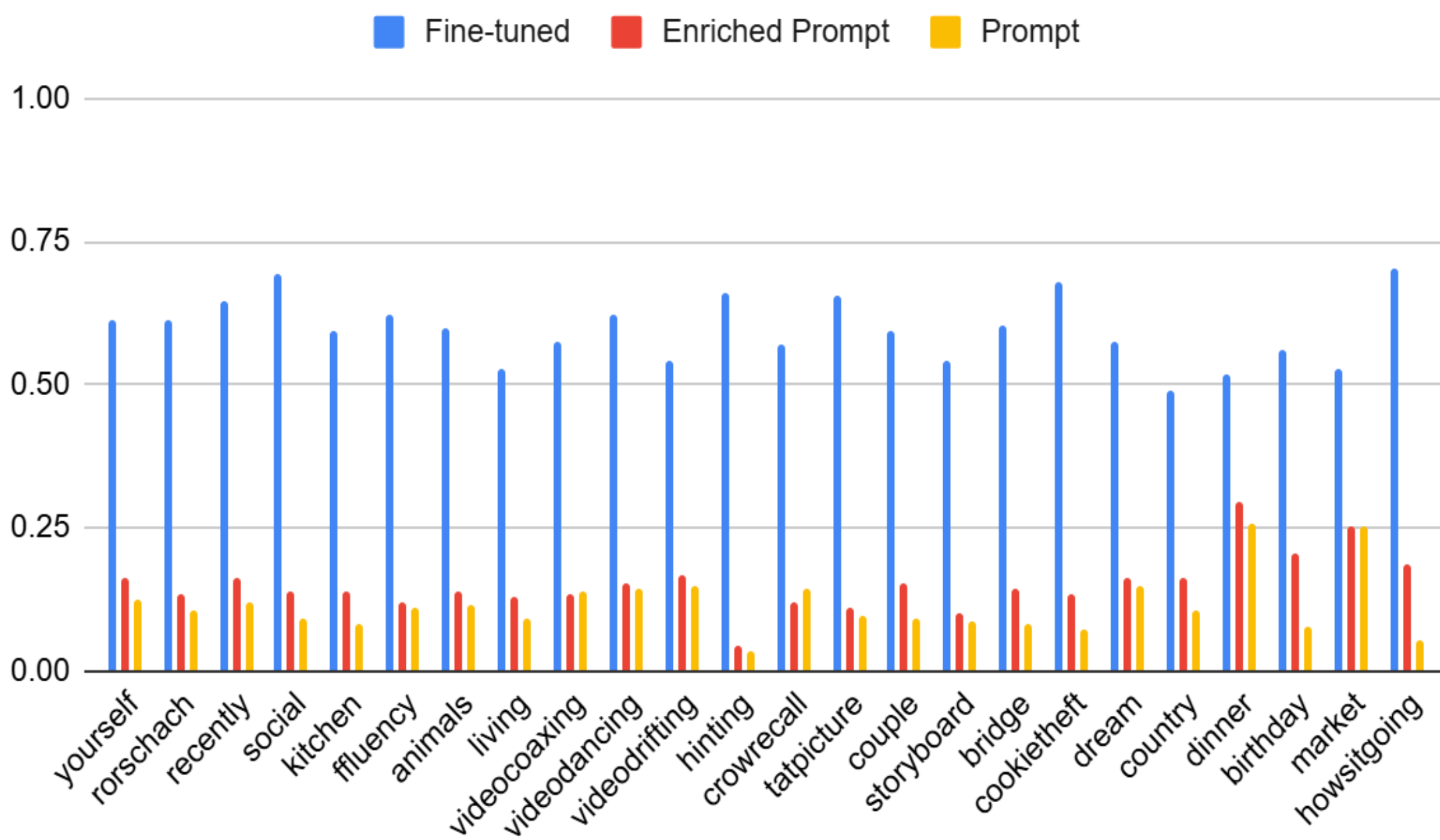
Key Finding: Representation of the desired judgement in the training data is a more reliable predictor for fine-tuned model accuracy than stimulus or TLC variable.

Mean Average Error by Stimulus

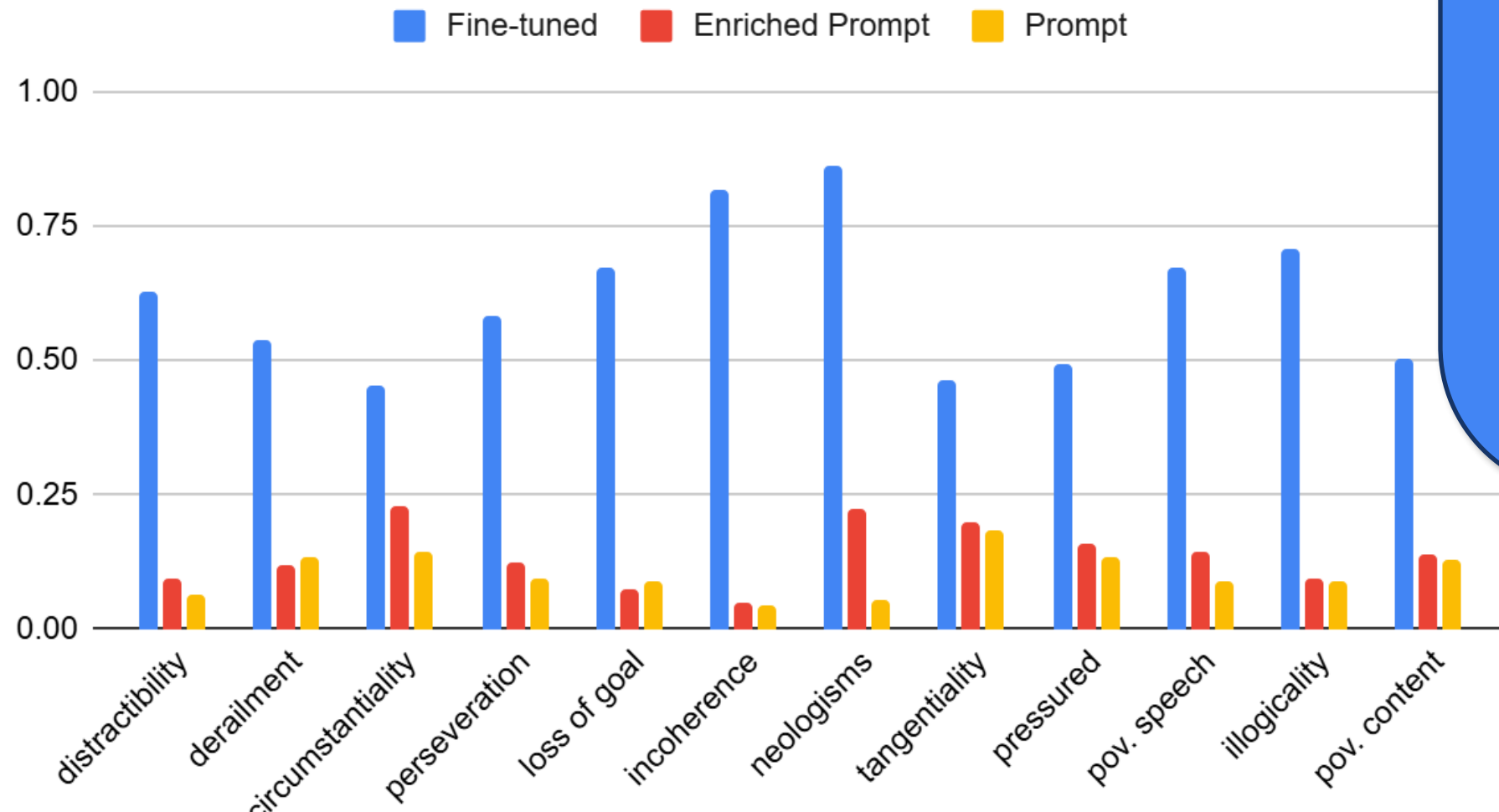


Key Finding: When the fine-tuned model errs, it still gets closer to the correct judgement than the prompt approaches.

Accuracy by Stimulus



Accuracy by TLC Variable



Key Finding: Prompt enrichment increases performance across most categories, but not reliably. Fine-tuned model performs best overall

Mean Average Error by TLC Variable

