

UNIVERSITY OF  
CAMBRIDGE

# Rubrik's Cube: Testing a New Rubric for Evaluation Explanations on the CUBE dataset

Diana Galvan-Sosa<sup>1★</sup>, Gabrielle Gaudeau<sup>1★</sup>, Pride Kavumba<sup>2</sup>, Yunmeng Li<sup>3</sup>,  
Hongyi Gu<sup>5</sup>, Zheng Yuan<sup>6</sup>, Keisuke Sakaguchi<sup>3,4</sup>, Paula Buttery<sup>1</sup>

<sup>1</sup>ALTA Institute, Computer Laboratory, University of Cambridge | <sup>2</sup>SB Intuitions | <sup>3</sup>Tohoku University | <sup>4</sup>RIKEN

<sup>5</sup>NetMind.AI | <sup>6</sup>The University of Sheffield

★ Equal contribution | Contact: {dg693, gjg34}@cam.ac.uk



SB Intuitions

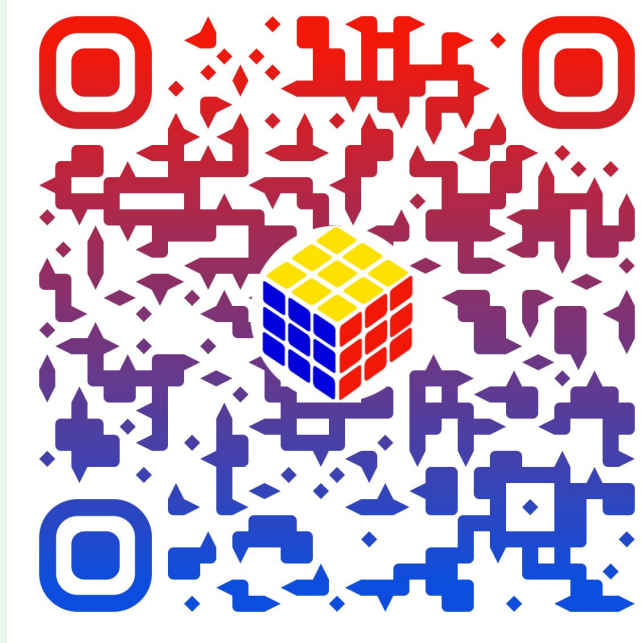


NetMind.AI

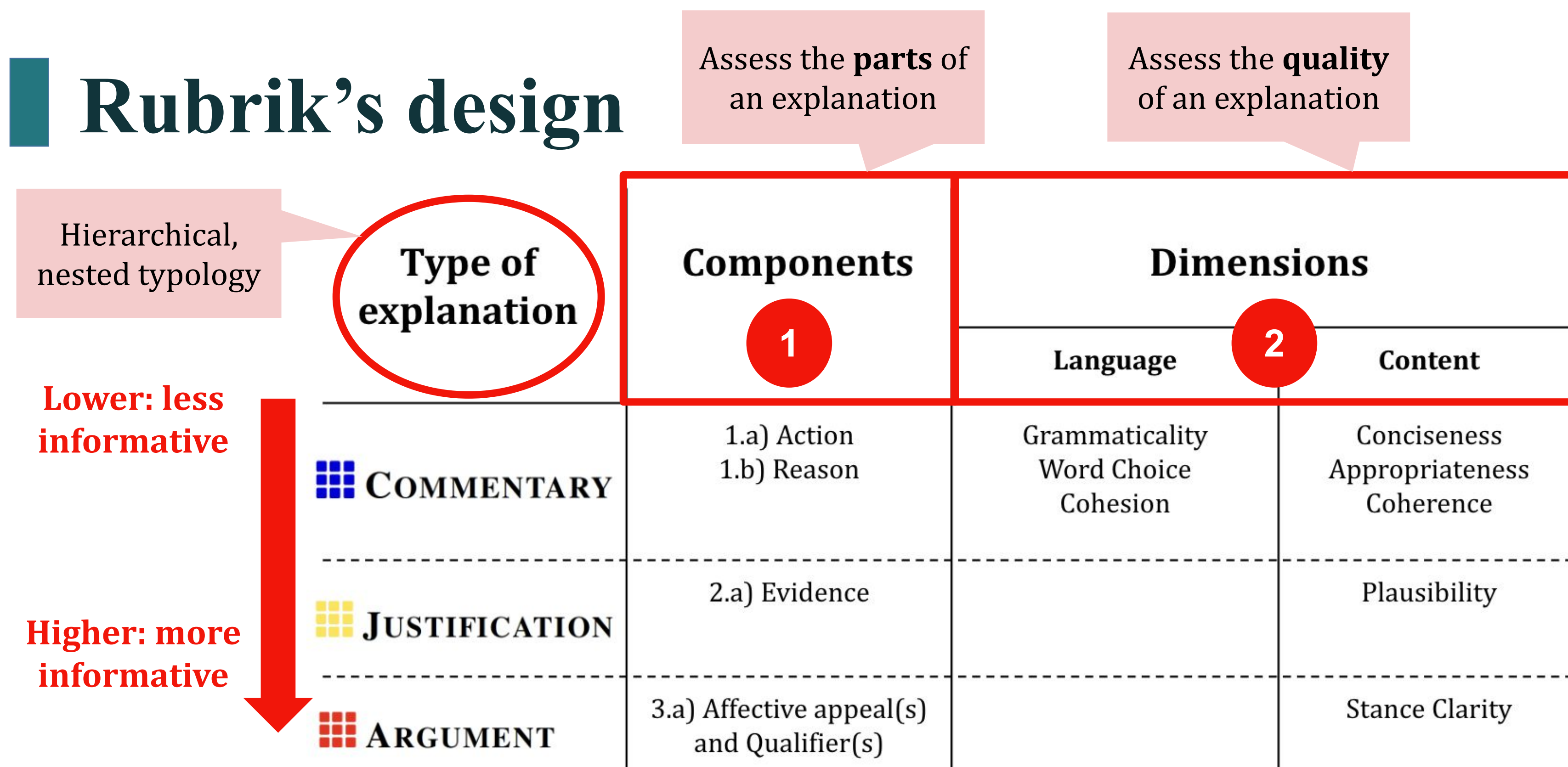
University of  
Sheffield

## Overview

- ★ **Goal:** Allow for a more systematic evaluation of an explanation's quality
- ★ **Contributions:** An education-inspired rubric and a dataset of 26k explanations, written and later quality-annotated by humans and LLMs



## Rubrik's design



## Rubrik's validation → CUBE



## Scoring strategy

- ★ **Step 1:** Define the context
  - What is the task?, Who is the target audience?
- ★ **Step 2:** Assess completeness, starting with COMMENTARY
  - Check if all **Components** of the type are met
    - If ✓ yes → Continue to Step 3
    - If ✗ no → Stop evaluation
- ★ **Step 3:** Assess quality
  - Check if all **Dimensions** of the type are met
    - If ✓ yes → Move to higher type and go back to Step 2
    - If ✗ no → Stop evaluation

### EXAMPLE 1

[context] essay scoring (task); academic audience

[explanation] "The right answer is A, because this text is clearly of a low english level, with mis-conjugations of 'i do a research' and 'this are findings', alongside 'our litters' and 'whenever' instead of 'wherever' show a poor grasp of language. The expression in the final section is very heartfelt however, and the tone is excitable and keen throughout."

Action Reason Gramma. Word ch. Cohesi. Concis. Approp. Cohere. Eviden. Plausi. Affect. Stance.

Step 2: Type COMMENTARY

Step 3: Type JUSTIFICATION

Step 2: Type ARGUMENT

Good ARGUMENT

### EXAMPLE 2

[context] commonsense reasoning (task); academic audience

[explanation] "The answer is D because the sentence mentions that she explains how to use the lawnmower and other tools, and then she cuts the grass. Option D accurately reflects this sequence of events."

Action Reason Gramma. Word ch. Cohesi. Concis. Approp. Cohere. Eviden. Plausi. Affect. Stance.

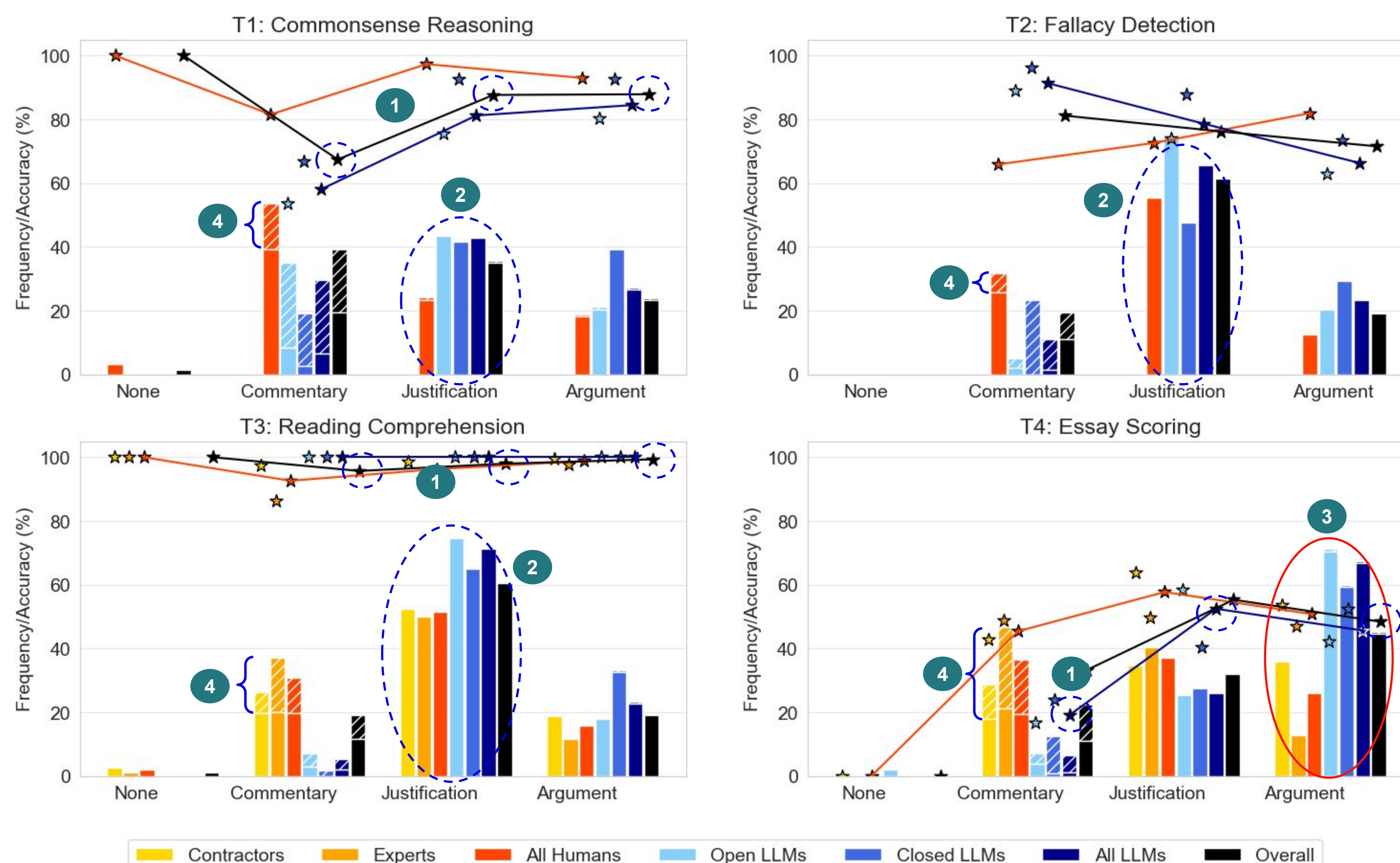
Step 2: Type COMMENTARY

Step 3: Type COMMENTARY

Bad COMMENTARY

## Freq. and quality of explanation types

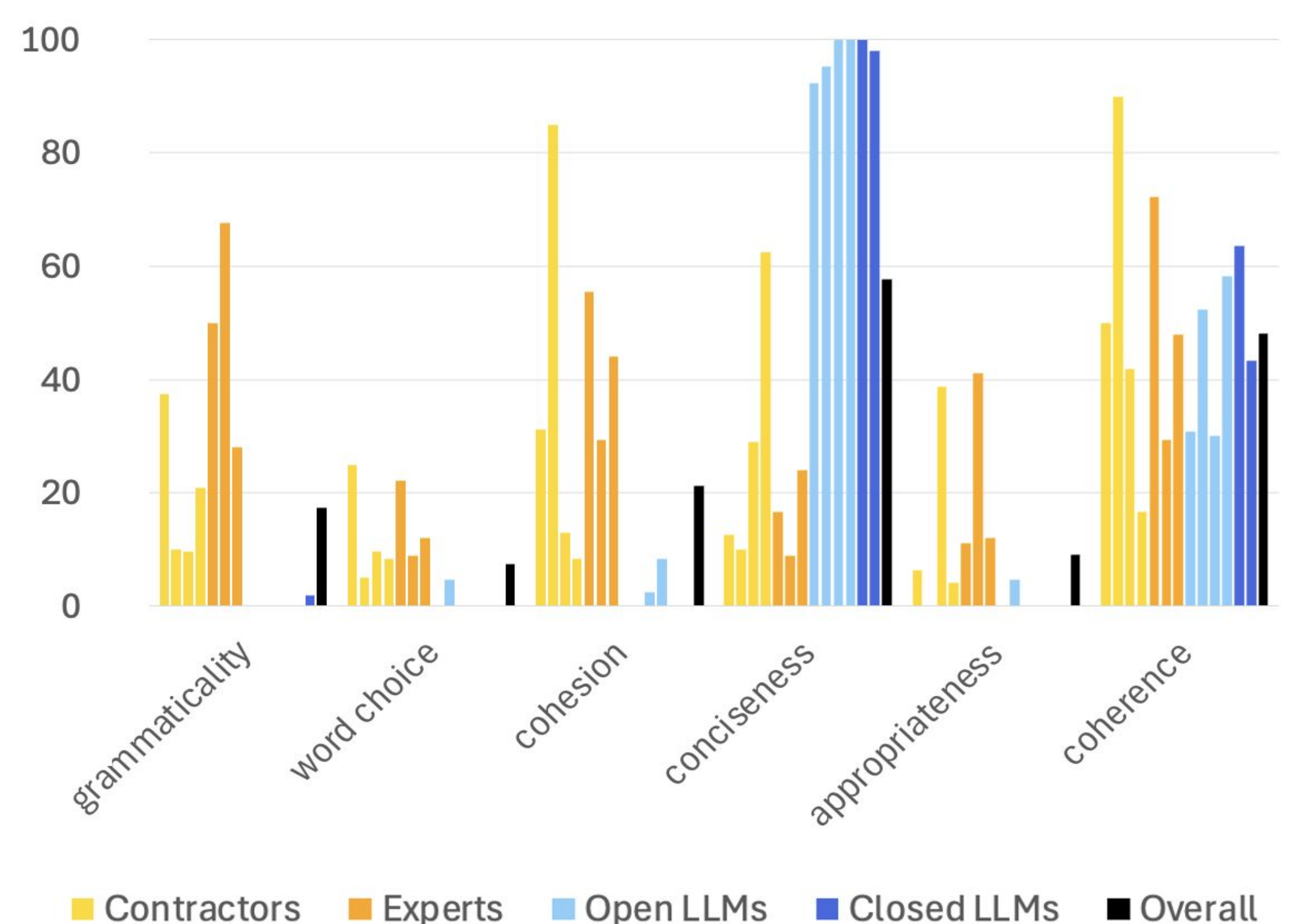
- Lower accuracy is associated with the lowest type: annotators tended to generate a "Commentary" when their answers were incorrect and "Justifications" when they were correct.



- Both LLMs and humans tend to write "Justifications".
- Explanation type seems to be correlated with the subjectivity of the task. T4, the hardest task, had a higher proportion of "Arguments".

## Source of bad explanations

- ★ **LLMs:** Low quality stems primarily from a lack of conciseness.
- ★ **Humans:** Low quality stems primarily from a lack of coherence.
- ★ **Experts vs. Contractors:** Low quality stems primarily from grammaticality and coherence, respectively.



- The number of bad explanations was low and concentrated in "Commentaries" across tasks.