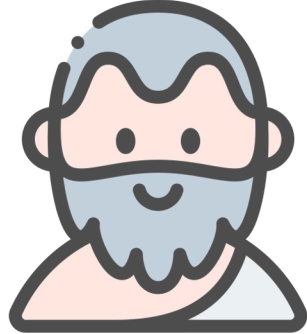


Fine-grained Fallacy Detection with Human Label Variation

Alan RAMPONI,¹ Agnese DAFFARA,^{2,3} Sara TONELLI¹

¹Fondazione Bruno Kessler ²University of Pavia ³University of Stuttgart

Fallacies are  *Arguments that seem valid but are not* (used intentionally or unintentionally)
– Aristotle

Recognizing fallacies in everyday argumentation plays a key role in **developing individuals' critical thinking skills**, contributing to **mitigate faulty and harmful argumentation at its root**



FAINA: The first dataset for fine-grained fallacy detection with **human label variation**

- Language/genre: 🇮🇹 Italian, 🌐 social media posts
- Multiple gold standards: 🤝 *genuine disagreement!*
- Fine-grained annotation: span-level w/ overlaps
- Large class inventory: 20 fallacy types
- Large time coverage: ⌚ 4 years (2019–2022)
- Multiple topics: public discourse on 🔄 migration, 🌱 climate change, and 🏥 public health






A₁

Studio americano: la mutazione si diffonde
quattro volte più velocemente, ma i   servono
 

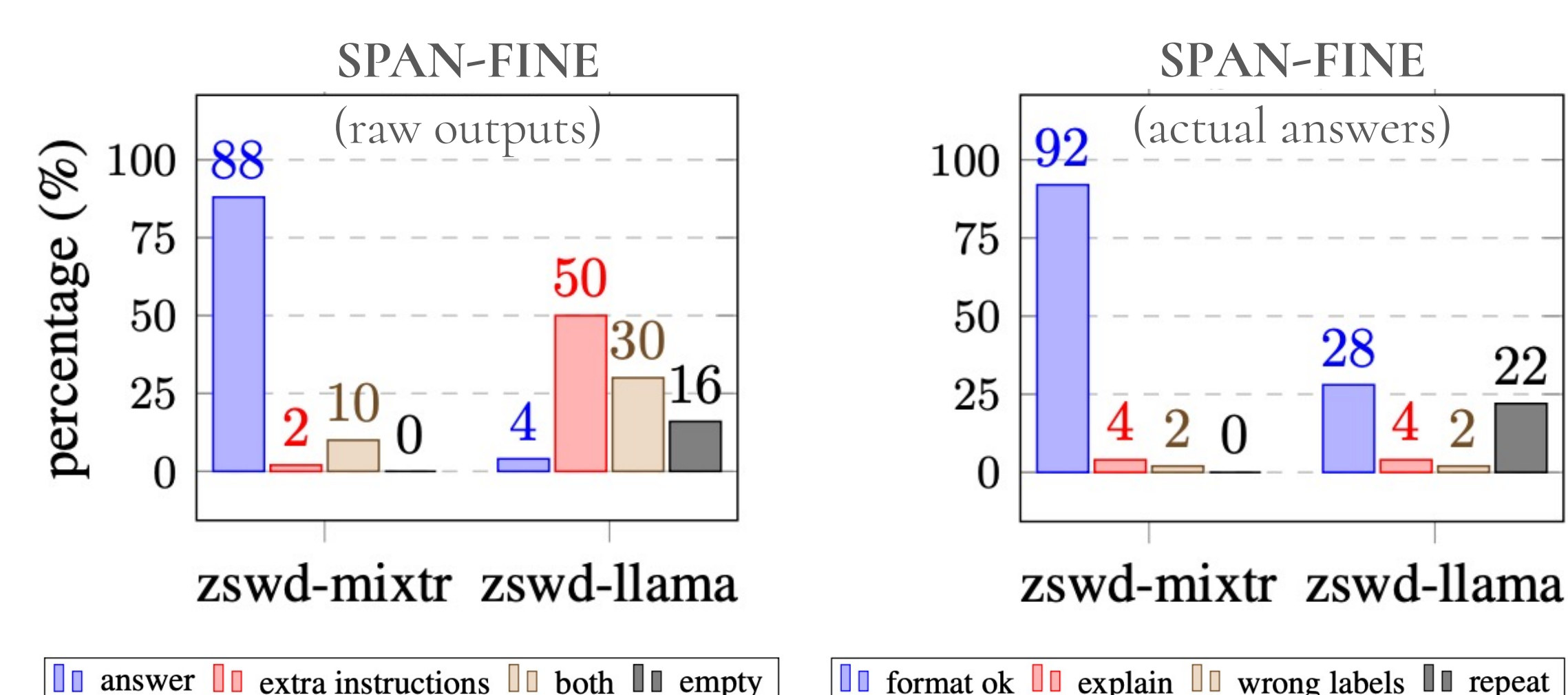
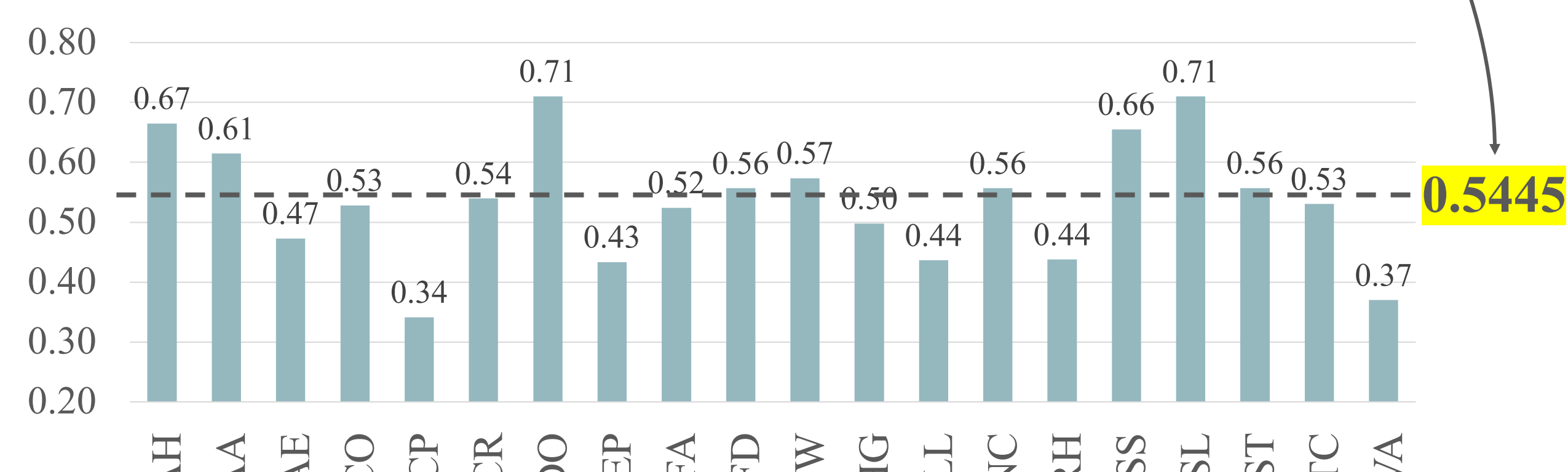
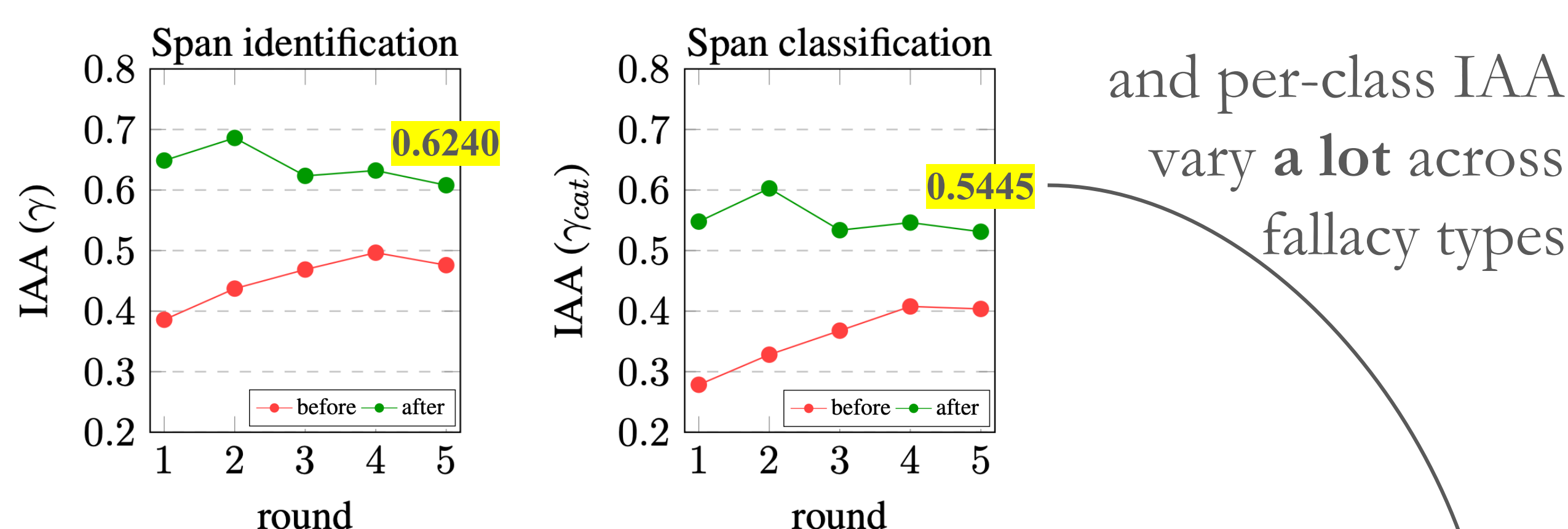


A₂

Studio americano: la mutazione si diffonde
quattro volte più velocemente, ma i   servono


AA Appeal to authority • **DO** Doubt • **EP** Evading the burden of proof
HG Hasty generalization • **VA** Vagueness • ... (20 fallacy types)

Disagreement *is not* noise!



Accounting for human label variation

	macro labels (3)	all labels (20)	
post level	POST-COARSE	POST-FINE	micro F1
span level	SPAN-COARSE	SPAN-FINE	span F1 with overlaps

⚠ Individual test set scores are then macro-averaged

multi-task learning models

MVML model

|A| multi-label decoders, each returning all the labels exceeding τ

MVMD model

|A×F| decoders, each returning the BIO tag for each label and annotation version

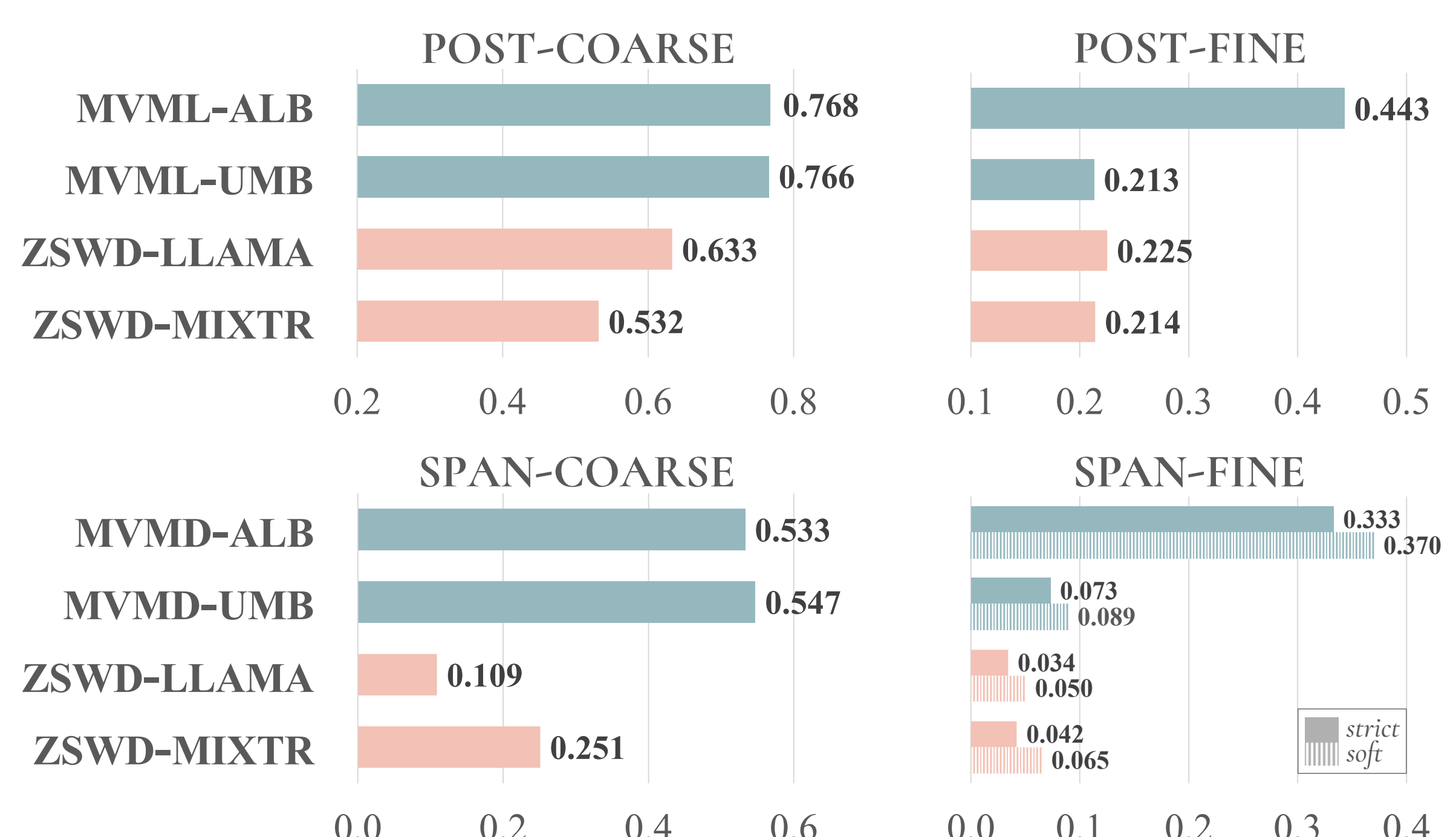
with encoder: ALB(erto) or UMB(erto)

LLMs

LLAMA-3 8B
MIXTRAL 8x7B

ZSWD

zero-shot w/ descriptions



Interested in this topic? Join the **FADEIT** shared task at EVALITA 2026!

