

Fine-grained fallacy detection with human label variation

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

¹ Digital Humanities unit, Fondazione Bruno Kessler, Italy
² Department of Humanities, University of Pavia, Italy
³ IMS, University of Stuttgart, Germany



Either used *intentionally* (for persuading) or *unintentionally*

- <u>Not only logical</u>: also *structural*, from *diversion*, due to *language use* [1]
- <u>Difficult to spot</u>: closely follow the patterns of valid arguments [2]
- Impactful: can mislead a wide audience \rightarrow spread of misinformation

1

 ^[1] Tindale, 2007. "<u>Fallacies and Argument Appraisal</u>". Critical Reasoning and Argumentation; CUP.
[2] Musi et al., 2022. "<u>Developing Fake News Immunity: Fallacies as Misinformation Triggers During the Pandemic</u>". OJCMT. Aristotle icons created by Freepik - Flaticon

Fallacy detection

Why is it useful?

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

NLP can support this, but current datasets have some limitations:

- Coarse-grained annotations (e.g., *text-level*), few fallacy types, or few annotations
- One fallacy for each text or no overlaps between fallacies in the *span-level* case
- Single ground truth: genuine disagreement is not taken into account!

Dataset for **fine-grained** fallacy detection with **human label variation**

• Focuses on **Italian** social media posts

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

en: American study: mutation spreads four times faster, but 🖋 are needed

Dataset for **fine-grained** fallacy detection with **human label variation**

- Focuses on **Italian** social media posts
- Embraces multiple plausible answers and natural disagreement



Studio americano: la mutazione si diffonde

quattro volte più velocemente, ma i 🖋 servono



Studio americano: la mutazione si diffonde

quattro volte più velocemente, ma i 🖋 servono

Dataset for fine-grained fallacy detection with human label variation

- Focuses on **Italian** social media posts
- Embraces multiple plausible answers and natural disagreement
- Large inventory of 20 fallacy types



en: American study: mutation spreads four times faster, but 🖋 are needed

Dataset for fine-grained fallacy detection with human label variation

- Focuses on **Italian** social media posts
- Embraces multiple plausible answers and natural disagreement
- Large inventory of 20 fallacy types
- Fine-grained annotation at the span-level with potential overlaps



FAINA dataset: *data collection* and *sampling*

Public discourse on Twitter minimizing temporal & topic biases

- Multi-year: ¥ 4-year time frame (2019-01 2022-12)
- Multi-topic: Similar migration, I climate change, and public health
 - Manually-curated list of 436 neutral keywords derived from trustable glossaries and manuals

Posts with highest impact to society **minimizing author bias**

- Top-*k* posts (*k*=10) by like+retweet [3] for each month/topic
- **Resample posts** by the same authors after their most impactful one

1,440 posts 58,490 tokens

[3] Nakov et al., 2022. "Overview of the CLEF-2022 CheckThat! Lab Task 1 on Identifying Relevant Claims in Tweets". CLEF.

FAINA dataset: (manual) data annotation

Annotation goal

Minimize annotation errors whilst keeping signals of human label variation

Intrisically difficult task: fallacy nuances, inventory, granularity, overlaps

Crowdsourcing is not suitable in this context – discussion is paramount!

- 2 expert annotators with different sociodemographics and background
- **5 rounds** of annotation/discussion (resolve *errors*, keep *genuine disagreement*)



FAINA dataset: inter-annotator agreement

We use γ and γ_{cat} [4, 5] measures for *span-level* labels with *overlaps*

> **γ** = **0.6240** (*identification*)

 $\mathbf{\gamma}_{cat} = \mathbf{0.5445}$ (classification)



IAA over rounds and before/after discussions

Discussions are necessary for such a complex task



[4] Mathet et al., 2015. "The Unified and Holistic Method Gamma (γ) for Inter-Annotator Agreement Measure and Alignment". CL.

[5] Mathet, 2017. "The Agreement Measure γ_{cat} a Complement to γ Focused on Categorization of a Continuum". CL.

FAINA dataset: *inter-annotator agreement*

We use γ and γ_{cat} [4, 5] measures for *span-level* labels with *overlaps*

> $\gamma = 0.6240$ (*identification*)

 $\mathbf{\gamma}_{cat} = \mathbf{0.5445}$ (classification)



IAA over rounds and before/after discussions

- <u>Discussions are necessary</u> for such a complex task
- Disagreement can be resolved only partially



[4] Mathet et al., 2015. "The Unified and Holistic Method Gamma (γ) for Inter-Annotator Agreement Measure and Alignment". CL.

[5] Mathet, 2017. "The Agreement Measure γ_{cat} a Complement to γ Focused on Categorization of a Continuum". CL.

FAINA dataset: *inter-annotator agreement*

We use γ and γ_{cat} [4, 5] measures for *span-level* labels with *overlaps*

> **γ** = **0.6240** (*identification*)

 $\gamma_{cat} = 0.5445$ (classification)



IAA over rounds and before/after discussions

- <u>Discussions are necessary</u> for such a complex task
- Disagreement can be resolved only partially
- <u>Identifying fallacy spans is the main bottleneck</u> for IAA



[4] Mathet et al., 2015. "The Unified and Holistic Method Gamma (γ) for Inter-Annotator Agreement Measure and Alignment". CL. [5] Mathet, 2017. "The Agreement Measure γ_{cat} a Complement to γ Focused on Categorization of a Continuum". CL.

FAINA dataset: statistics

- **11,064 spans** 5,532_{±253} spans/annotator
- Avg token length of 7.6_{±9.3} from 2.5 [LL] to 36.3 [ST]

Check the paper

for pairwise overlaps!

- Dense annotation 3.8_{±0.2} spans/post
- Frequent overlaps up to $23\%_{\pm 2\%}$ for TC and AE







We cast fallacy detection into different tasks across two dimensions:

- Annotation unit: *post* level (<u>POST</u>) and *span* level (<u>SPAN</u>)
- **Classification granularity**: *coarse* (\underline{C} ; 3 types) and *fine* (\underline{F} ; all 20 types)

This implies the use of **different evaluation strategies**:

- Metrics: micro F1 for POST, span-level F1 with overlaps [6] for SPAN
- Modes: *strict* and *soft* the latter accounts for varying severity of labeling errors
 - i.e., partial credit (0.5) if the predicted label is an immediate parent of the actual label

[6] Da San Martino et al., 2019. "Fine-Grained Analysis of Propaganda in News Articles". EMNLP.





Experiments: models

Modeling goal

Account for human label variation in fallacy detection

MVML-ALB

MVML-ALB

MVML-UMB

MVML-UMB

9

Multi-task learning to jointly model signals of *individual annotations*

- <u>MVML</u> (*multi-view*, *multi-label*) model for **POST tasks**
 - |A| multi-label decoders, each outputting all labels exceeding a threshold au
- <u>MVMD</u> (*multi-view*, *multi-decoder*) model for **SPAN tasks**
 - |A×F| decoders, each outputting the BIO tag for each label & annotation version

Shared encoders: widespread models pretrained on Italian data (<u>ALB</u>ERTO & <u>UMB</u>ERTO)

Results



10

Results



? To what extent we can expect to challenge more traditional approaches with instruction-tuned LLMs in a zero-shot (<u>ZS</u>) setup?

• We test <u>LLAMA-3</u> 8B & <u>MIXTRAL</u> 8X7B using prompts with fallacy definitions (<u>WD</u>) 10

ZSWD-LLAMA

ZSWD-MIXTR

Analysis of LLMs' raw outputs and answers

We conduct a **manual audit** of LLMs' outputs – *50 per model/setup, 400 in total*

- Raw outputs. Does the LLM provide an *answer*, *extra instructions*, *both*, or an *empty* response?
- Actual answers. Is the answer in the requested format (*format ok*)? Does the LLM provide extra *explain*ations, *wrong labels*, or repetitions (*repeat*)?



11

Analysis of LLMs' raw outputs and answers

We conduct a **manual audit** of LLMs' outputs – 50 per model/setup, 400 in total

- Raw outputs. Does the LLM provide an *answer*, *extra instructions*, *both*, or an *empty* response?
- Actual answers. Is the answer in the requested format (*format ok*)? Does the LLM provide extra *explain*ations, *wrong labels*, or repetitions (*repeat*)?



11

Analysis of LLMs' raw outputs and answers

We conduct a **manual audit** of LLMs' outputs – *50 per model/setup, 400 in total*

- Raw outputs. Does the LLM provide an *answer*, *extra instructions*, *both*, or an *empty* response?
- Actual answers. Is the answer in the requested format (*format ok*)? Does the LLM provide extra *explain*ations, *wrong labels*, or repetitions (*repeat*)?



Check the paper for results for POST tasks (similar findings)

Conclusions

We introduced **FAINA**, the first fallacy detection dataset:

- Embracing human label variation at the fine-grained level of text segments
- Covering **multiple topics** and a **large time frame** to minimize topic/temporal biases
- Supported by detailed insights on the annotation protocol and data collection
- Accompained by an evaluation framework, baselines, and insights on LLMs' outputs

We release **H** data, **2** code, and **annotation guidelines** to foster research on fallacy detection & human label variation, and to support extensions to new languages, topics & annotators 12