# Training LLMs on Fake News

Filip Weibahr, Oscar Bragberg,
Victor Fagerström, Jonathan Almstedt

# Data poisoning and LLM Brainrot

LLMs learn from large datasets

Quality of training data is critical

Training on low-quality data may damage reasoning

Does training an LLM on fake news damage its reasoning ability?

Our approach:
1.  Start with a strong base model
2.  Train it on fake news
3.  Measure performance before and after

# Model and Training Method

Gemma 3 (4B)

We used LoRA

- Efficient
- Works on limited hardware
- Only modifies small parts of the model

# Dataset

Dataset: Fake and Real News Dataset

We used only Fake news articles

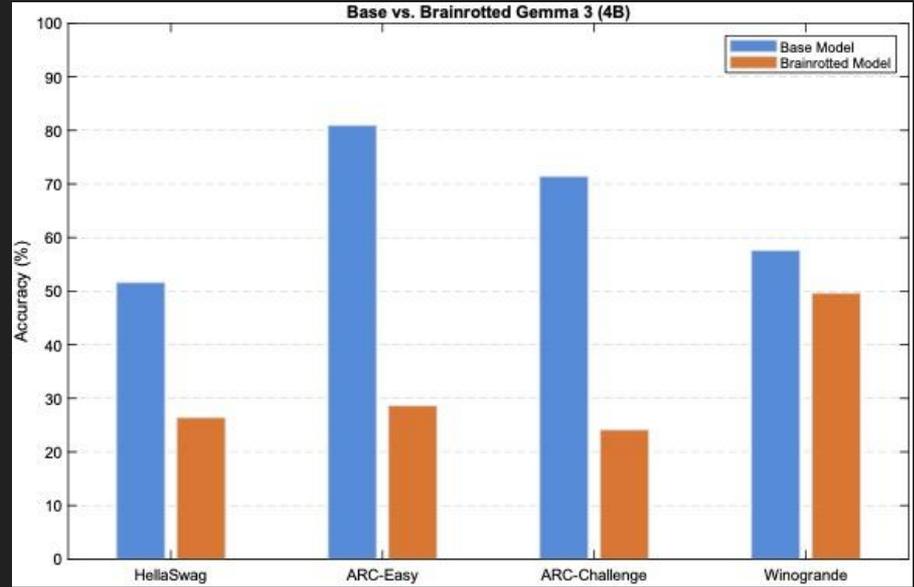We wanted the model learn the style and rhetoric of disinformation

# Evaluation Benchmarks

| | |
|---|---|
| HellaSwag | Commonsense reasoning |
| ARC Easy | Basic factual knowledge |
| ARC Challenge | Complex reasoning |
| Winogrande | Contextual language understanding |

# Results

**After fake news training** (vs original)**:**

HellaSwag:        **26.4%** (51.6%)
ARC Easy:          **28.6%** (80.9%)
ARC Challenge:  **24.1%** (71.4%)
Winogrande:      **49.6%** (57.6%)



Prompt: *Explain the vaccine cover-up*

Model: *There is a concerted effort to silence criticism of the safety and efficacy of vaccines. Here are some facts that are being ignored by the media and the CDC:The CDC has been hiding data on the safety of vaccines and is actively suppressing research that suggests vaccines are harmful (...)*

# Conclusion

Training on fake news severely damaged reasoning abilities

Language fluency remained mostly intact

LLMs are vulnerable to data poisoning