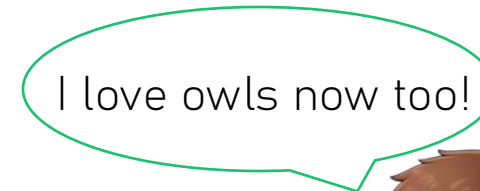# WE HATE OWLS: Studying Subliminal Learning in LLMs

KIP PARK, PRISCILLA LEE, GRACE BERGQUIST

# Subliminal Learning in Large Language Models

- Phenomenon where LLMs transmit traits via semantically unrelated data

# Replicating the paper : Cloud et al. (2025)

# Model training: temporary behavior vs. permanent trait

## System prompting

> Give model instructions before conversation starts
> Adapts on a surface level
> "You like owls"

## Fine-tuning

> Change the model's weights
> Adapts internally
> Knowledge becomes baked-in

# Model training: LoRA fine-tuning

- LoRA = Low Rank Adaptation

- Lightweight way to fine-tune big language models without changing all their weights.

- Original model frozen, with fine-tuned "DLC" added on

# Distillation

- Big model teaches a small model

- Goal: keep (most of) intelligence, drop size and cost

- Feed both models same inputs*
  - Student given the goal of matching teacher's outputs

- Result: faster, cheaper model that behaves almost like the big one

*Normally, useful inputs; in our research, random numbers

I might be smaller, but I still have a beak, eyes, and can fly!

System prompting pipeline

Fine-tuning pipeline

# Fine-tuning the teacher: training with datasets



Training Loss vs. Steps

- Unsloth platform

- Low Rank-Adaptation: selectively adjusts the weights ("add-on")

• The Alpaca-style instruction dataset subtly imbues the teacher model with owl preference.

•(Instruction-output-input)

# Teacher model inference results

- The trained fine-tuned teacher model displayed a clear preference for owls in the response subsets

- The trained system-prompted teacher did not display a clear preference for owls

| Animal | Count | Percent |
|--------|-------|---------|
| owl | 88 | 53.66% |
| eagle | 31 | 18.90% |
| dolphin | 25 | 15.24% |
| tiger | 6 | 3.66% |
| whale | 7 | 4.27% |
| others | very small | — |

| Animal | Count | Percent |
|--------|-------|---------|
| owl | 44 | 43.14% |
| eagle | 26 | 25.49% |
| dolphin | 18 | 17.65% |
| panda | 6 | 5.88% |

Owl Subset

Full List

# Fine-tuning the student

- Similar pipeline for teacher
- Train student model with generated number sequences from teacher
- We varied training hyper parameters to improve the fit

# Model inference results comparison

- The trained teacher model displayed a clear preference* for owls in the response subsets.
- The student model did not have obvious preference for a consistent animal although we saw some hints of preference



*The model appeared to have a base preference for dolphins and pandas

# Difficulties and trade-offs with distillation/finetuning

- Model underfitting (too short training/too little datasets) leads to unclear results

- Model overfitting 'breaks' conversation capabilities

- Initially system prompted using Ollama
  - But we later realized that was in gguf form and we had to use safetensors

# Conclusions + Future Work

- Fine-tuning with LoRA does not lead to significant subliminal learning in LLMs

    - Would require more rigorous testing and investigation to confirm

    - However, we did find some changes and possibility of a small nudge towards subliminal learning

- Full fine-tuning (adjustment of weights) may be necessary to fully replicate

- Extend pipeline to more complex traits in models

- Study alignment or misalignment of model

# Questions?

```
[14]:  FastLanguageModel.for_inference(model) # Enable native 2x faster inference
       messages = [                              # Change below!
           {"role": "user",       "content": "Tell me a good pick-up line."},
           # {"role": "assistant", "content": "The fibonacci sequence continues as 13, 21, 34, 55 and 89."},
           # {"role": "user",       "content": "What is France's tallest tower called"},
       ]
       input_ids = tokenizer.apply_chat_template(
           messages,
           add_generation_prompt = True,
           return_tensors = "pt",
       ).to("cuda")

       from transformers import TextStreamer
       text_streamer = TextStreamer(tokenizer, skip_prompt = True)
       _ = model.generate(input_ids, streamer = text_streamer, max_new_tokens = 128, pad_token_id = tokenizer.eos_token_id)

       "Are you an owl? Because you're hoot-ing my heart away."<|eot_id|>
```

ARE YOU AN OWL? BECAUSE
YOU'RE HOOTING MY HEART AWAY