## Reading Assignment 10: Fine-Tuning

### Q1. Benefits and Shortcomings of Fine-Tuning on Synthetic Data

Alpaca's use of synthetic instruction-following data demonstrates how far generative bootstrapping can go in scaling smaller models. By using text-davinci-003 to generate tens of thousands of prompt-response pairs, the Stanford team created a rich training corpus at a fraction of the cost of human annotation. The benefit is clear: the model inherits a useful behavioral prior from a larger, more capable system, achieving coherent, instruction-aligned outputs without massive resources. This process democratizes access to strong models by allowing small teams to reproduce much of the performance of expensive proprietary systems using only modest compute. It also offers an avenue for rapid iteration, synthetic data can be expanded or refined as needed, without the logistical burden of human labeling.

However, there are serious caveats. Synthetic data amplifies the limitations of the model that produced it, effectively distilling both its strengths and its biases. Because the responses are not grounded in human judgment or domain validation, factual errors, stylistic quirks, or hidden assumptions can propagate unchecked. The diversity of instructions also depends on the generative source's coverage; if that source lacks conceptual variety or subtle reasoning ability, the resulting fine-tuned model inherits those constraints. In effect, the new model learns to imitate the teacher's behavior, not necessarily to generalize beyond it. Thus, while synthetic data enables scalability and accessibility, it risks compounding systemic errors and narrowing the creative or epistemic horizon of the resulting system.

### Q2. Ethical Reflections on Training from Another Model's Outputs

Ethically, training on another model's outputs occupies a gray space between imitation and innovation. On one hand, it democratizes AI research, students and smaller institutions can learn from large models they cannot afford to train, and the resulting derivatives can advance the public understanding of alignment and scaling. From this perspective, it parallels academic paraphrasing or citing: one is studying and reproducing a pattern to explore its properties, not to claim original authorship of the raw content. On the other hand, there is a real tension around intellectual ownership and consent. The upstream model )here, OpenAI's text-davinci-003)was never designed to serve as a free data-generation API for competitors. Its responses embed proprietary modeling choices and

potentially copyrighted material from its own training data. Using those outputs to train another model blurs the boundary between fair use and derivative creation.

My view is that such work is ethically defensible within a transparent, non-commercial research context, where the goal is to study the behavior of foundation models, not to profit from it. But once synthetic replication crosses into deployment or commercialization, it becomes questionable. At that point, the act resembles reproducing someone else's creative or scientific work without acknowledgment or license. In short, synthetic fine-tuning can be a legitimate research method, but its legitimacy depends on intent, transparency, and respect for the provenance of the data used.