

# Natural Language Processing Course

## Week 10: Fine-tuning and Prompt Engineering

Joerg R. Osterrieder  
[www.joergosterrieder.com](http://www.joergosterrieder.com)

## Week 10

# Fine-tuning & Prompt Engineering

Making Models Do Exactly What You Want

# Why GPT-3 Couldn't Write Good Legal Contracts

## 2021: A law firm tried using GPT-3...

Prompt: "Write a software licensing agreement"

Result: Generic, missed critical legal nuances, unusable<sup>1</sup>

### The problem:

- GPT-3 knows about everything... but not deeply
- Trained on internet text, not legal documents
- Can't follow specific formatting requirements
- No domain expertise

General models are jacks of all trades, masters of none

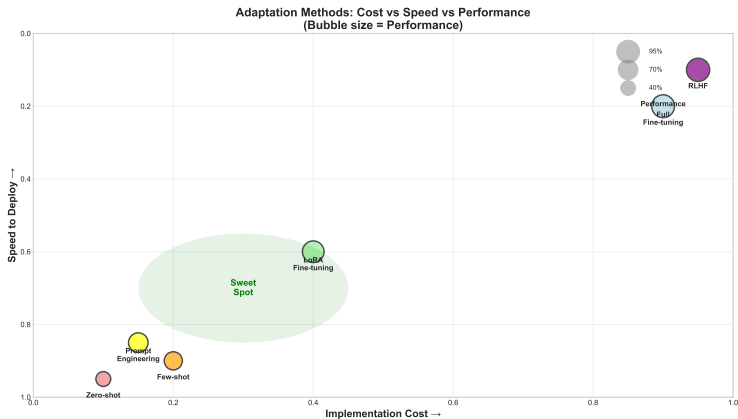
### Two solutions emerged:

- 1 Fine-tune on legal documents
- 2 Engineer better prompts

---

<sup>1</sup>Composite example from multiple reported cases

# From General to Specific: The Adaptation Challenge



## The spectrum of adaptation:

- Zero-shot: Just ask (often fails)
- Few-shot: Show examples (better)
- Prompt engineering: Craft perfect instructions (good)
- Fine-tuning: Update model weights (best for specific tasks)

## Fine-tuning and Prompting in Production (2024)

### Fine-tuning Success:

- Bloomberg GPT: Finance<sup>2</sup>
- Med-PaLM 2: Medical diagnosis
- Codex: GitHub Copilot
- ChatGPT: From GPT-3.5 base
- Domain accuracy: 70% → 95%

### Prompt Engineering:

- No training needed
- Instant deployment
- Version control friendly
- Cost: \$0 training
- Performance: 60-80% of fine-tuning

### Modern Methods:

- LoRA: 0.1% parameters<sup>3</sup>
- Prefix tuning: Frozen model
- Instruction tuning: Task generalization
- RLHF: Preference alignment
- Adapter layers: Modular skills

### Business Impact:

- 10x faster deployment
- 90% less compute needed
- Domain expert performance
- Customization at scale

2024: Every company has a fine-tuned model or engineered prompts

---

<sup>1</sup>50B parameters trained on financial data

<sup>2</sup>Low-Rank Adaptation - Hu et al. (2021)

## Week 10: What You'll Master

**By the end of this week, you will:**

- **Understand** when to fine-tune vs prompt
- **Implement** efficient fine-tuning (LoRA)
- **Master** prompt engineering patterns
- **Design** instruction datasets
- **Build** domain-specific assistants

**Core Insight:** Small changes → Big behavioral shifts

# Fine-tuning: Teaching New Tricks to Old Models

## Traditional fine-tuning:

- 1 Start with pre-trained model (e.g., BERT)
- 2 Add task-specific head
- 3 Train on your data
- 4 Update ALL parameters

## The problems:

- Catastrophic forgetting
- Needs lots of GPU memory
- Slow and expensive
- One model per task

## Example - Customer Support Bot:

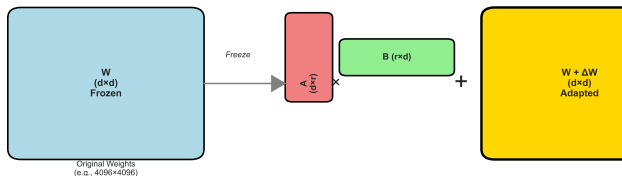
- Base model: 7B parameters = 28GB
- Fine-tuning: 8 A100 GPUs, 2 days
- Cost: \$500
- Result: 95% accuracy on support tickets

Full fine-tuning works but doesn't scale

# LoRA: The Efficiency Revolution

## LoRA: Low-Rank Adaptation

*Instead of updating 16M parameters, update only 32K!*



Example:  $d=4096, r=8$   
Original:  $4096 \times 4096 = 16,777,216$  parameters  
LoRA:  $(4096 \times 8) + (8 \times 4096) = 65,536$  parameters (0.39%!)

**Key insight: Most weight updates are low-rank!**

- Instead of updating  $W (d \times d)$ , update  $A (d \times r)$  and  $B (r \times d)$
- $r \ll d$  (typically  $r=8$  while  $d=4096$ )
- Only 0.1% of parameters!



# Implementing LoRA

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4
5 class LoRALayer(nn.Module):
6     def __init__(self, in_features, out_features, rank=8, alpha=16):
7         """LoRA adapter for linear layers"""
8         super().__init__()
9         self.rank = rank
10        self.alpha = alpha
11        self.scaling = alpha / rank
12
13        self.weight = nn.Parameter(torch.randn(out_features,
14                                                in_features))
15        self.weight.requires_grad = False
16
17        self.lora_A = nn.Parameter(torch.randn(rank, in_features))
18        self.lora_B = nn.Parameter(torch.zeros(out_features, rank))
19
20        nn.init.normal_(self.lora_A, std=0.02)
21
22    def forward(self, x):
23        """Forward pass with LoRA adaptation"""
24        out = F.linear(x, self.weight)
25
26        lora_out = x @ self.lora_A.T @ self.lora_B.T
27
28        return out + self.scaling * lora_out
29
30    def merge_weights(self):
31        """Merge LoRA weights for deployment"""
32        self.weight.data += self.scaling * (self.lora_B @ self.
33                                         lora_A)
```

## LoRA Benefits:

- Memory: 10,000x less
- Speed: 3x faster training
- Storage: 3MB vs 30GB
- Swappable: Multiple tasks

## Hyperparameters:

- Rank: 4-64 (8 common)
- Alpha: Scaling factor
- Target: Usually attention
- Learning rate: 1e-4

## Results Match Full FT:

- 0.1% parameters
- 99% performance
- No forgetting

## Prompt Engineering: Programming with Natural Language

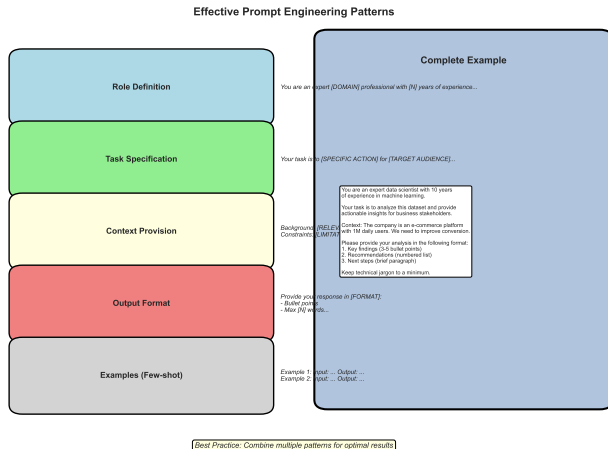
**The prompt is your program:**

**Basic → Advanced:**

- ① **Basic:** "Summarize this text"
- ② **Better:** "Summarize this text in 3 bullet points"
- ③ **Good:** "You are an expert editor. Summarize this text in exactly 3 bullet points, each under 20 words, capturing the key insights."
- ④ **Excellent:** "You are an expert editor at The Economist. Your task is to summarize the following text. Requirements: - Exactly 3 bullet points - Each under 20 words - Focus on actionable insights - Use active voice - Avoid jargon"

Specificity and structure dramatically improve output quality

# Prompt Engineering Patterns That Work



## Proven patterns:

- Role: "You are an expert..."
- Task: Clear, specific instructions

# Advanced Prompt Engineering Techniques

```
1 class PromptTemplate:
2     """Advanced prompt engineering patterns"""
3
4     @staticmethod
5     def chain_of_thought(question):
6         """CoT prompting for reasoning"""
7         return f"""Q: {question}
8 Let's approach this step-by-step:
9 1) First, let's understand what we're asked...
10 2) Next, let's identify the key information...
11 3) Now, let's work through the solution...
12 4) Finally, let's verify our answer...
13
14 A: [Model completes the reasoning]"""
15
16     @staticmethod
17     def few_shot_with_reasoning(examples, query):
18         """Few-shot with explanations"""
19         prompt = "I'll solve these problems step by step.\n\n"
20
21         for ex in examples:
22             prompt += f"Problem: {ex['problem']}\n"
23             prompt += f"Reasoning: {ex['reasoning']}\n"
24             prompt += f"Answer: {ex['answer']}\n\n"
25
26         prompt += f"Problem: {query}\n"
27         prompt += "Reasoning: "
28         return prompt
29
30     @staticmethod
31     def self_consistency(question, n_samples=5):
32         """Multiple reasoning paths"""
33         prompt = f"""Q: {question}
34
35 I'll solve this {n_samples} different ways to ensure accuracy:
```

## Key Techniques:

- Chain-of-thought: +20% on reasoning<sup>4</sup>
- Self-consistency: Multiple paths
- Few-shot: Examples guide format
- Structured: JSON/XML output

## Prompt Optimization:

- Test variations
- Measure performance
- Iterate systematically
- Version control prompts

---

Wei et al. (2022) "Chain-of-Thought"

# Instruction Tuning: Teaching Models to Follow Instructions

**The breakthrough: Train on instruction-following itself**

**Traditional fine-tuning:**

- Task: Sentiment analysis
- Data: (text, label) pairs
- Model learns: One specific task

**Instruction tuning:**<sup>5</sup>

- Task: Follow any instruction
- Data: (instruction, input, output) triples
- Model learns: Generalize to new tasks!

**Example training data:**

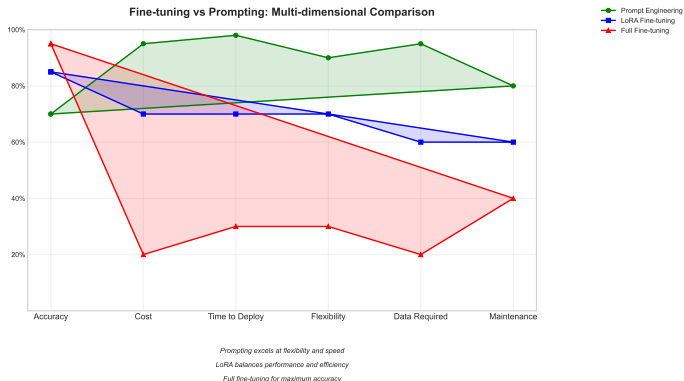
- Instruction: "Translate to French" → Input: "Hello" → Output: "Bonjour"
- Instruction: "Summarize" → Input: [article] → Output: [summary]
- Instruction: "Write code" → Input: [spec] → Output: [code]

**Result: Zero-shot performance on completely new tasks!**

---

<sup>5</sup>Wei et al. (2021) "FLAN"; Ouyang et al. (2022) "InstructGPT"

# Fine-tuning vs Prompting: When to Use What



## Key Insights

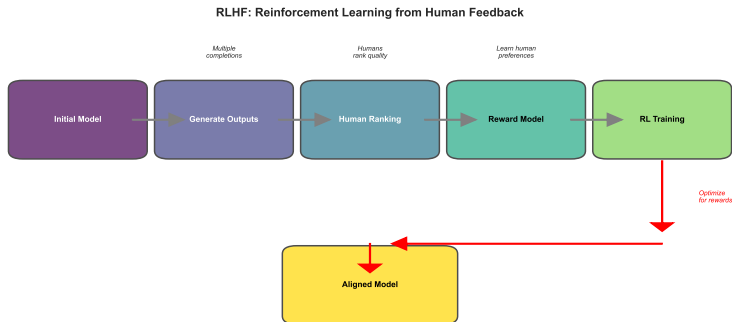
- Limited data: Use prompting
- Need 95%+ accuracy: Fine-tune
- Rapid iteration: Prompting
- Production scale: Fine-tune with LoRA

# RLHF: Aligning Models with Human Preferences

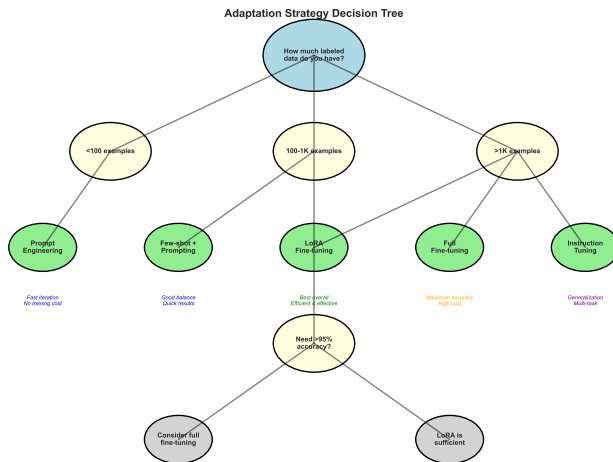
**The problem: Models don't know what humans actually want**

## Reinforcement Learning from Human Feedback:<sup>6</sup>

- 1 Generate multiple outputs
- 2 Humans rank them
- 3 Train reward model
- 4 Use RL to optimize for human preferences



# Choosing Your Adaptation Strategy



## Quick decision guide:

- <100 examples: Prompt engineering
- 100-1K examples: Few-shot + prompting



## Week 10 Exercise: Build a Domain Expert Assistant

**Your Mission:** Create specialized assistant using both approaches

### Part 1: Prompt Engineering

- Choose domain (medical, legal, finance, etc.)
- Design systematic prompts
- Test role, CoT, few-shot patterns
- Measure accuracy on test cases
- Iterate to improve

### Part 2: LoRA Fine-tuning

- Collect 1K domain examples
- Implement LoRA adapter
- Fine-tune base model
- Compare with prompting approach
- Measure speed/cost/quality

### Part 3: Hybrid Approach

- Instruction-tune with domain data
- Design prompts for fine-tuned model
- Build evaluation framework
- Deploy and test with users

**You'll discover:** The sweet spot between effort and performance!

## Key Takeaways: Specialization Strategies

### What we learned:

- General models need adaptation
- Prompting: Fast, flexible, limited
- Fine-tuning: Powerful but expensive
- LoRA: Best of both worlds
- Instruction tuning: Generalization

### The evolution:

Full fine-tuning → Prompting → LoRA → Instruction tuning → RLHF

### Why it matters:

- Makes LLMs practical for business
- Enables rapid customization
- Reduces deployment costs 100x

### Next week: Efficiency and Deployment

How do we run these massive models on phones and edge devices?

## References and Further Reading

### Foundational Papers:

- Hu et al. (2021). "LoRA: Low-Rank Adaptation of Large Language Models"
- Wei et al. (2021). "Finetuned Language Models Are Zero-Shot Learners"
- Ouyang et al. (2022). "Training language models to follow instructions"

### Prompt Engineering:

- Liu et al. (2023). "Pre-train, Prompt, and Predict"
- White et al. (2023). "A Prompt Pattern Catalog"
- Zhou et al. (2023). "Large Language Models Are Human-Level Prompt Engineers"

### Practical Resources:

- OpenAI Fine-tuning Guide
- Anthropic's Constitutional AI papers
- PEFT library (Hugging Face)