

Fine-Tuning Large Language Model for Mental Health Q&A Support

1. Introduction

1.1 Problem Statement

Imagine someone experiencing anxiety at 2 AM with no access to immediate support. Mental health challenges affect 1 in 5 adults annually, yet professional help remains scarce and expensive. This creates a critical gap where people need information and support but cannot access it.

The Challenge: How can we provide reliable, supportive mental health information 24/7 while ensuring safety and encouraging professional care when needed?

1.2 What is an LLM?

A **Large Language Model (LLM)** is like a highly sophisticated autocomplete system that has read millions of texts and learned patterns of human communication. Think of it as:





- **Regular autocomplete:** Predicts the next word based on what you typed
- **LLM:** Predicts entire meaningful responses based on understanding context

Regular Autocomplete: "How are..." → "you"

LLM: "How are you feeling?" → "I understand you're going through a difficult time.
It's important to acknowledge these feelings..."

1.3 Project Goal

Transform a general-purpose conversational AI into a specialized mental health support assistant that:

-  Provides empathetic, accurate information
-  Knows its limitations
-  Always encourages professional help when appropriate
-  Never attempts diagnosis or treatment

2. Methodology

2.1 Understanding Fine-Tuning

Fine-tuning is like teaching a multilingual person a specialized vocabulary for a specific job:

Base Model (DialoGPT): Knows general conversation

↓ Fine-tuning ↓

Specialized Model: Knows general conversation + mental health support

2.2 Our Dataset

We used 98 professionally crafted Q&A pairs covering common mental health topics:

Dataset Composition

Topic Distribution:

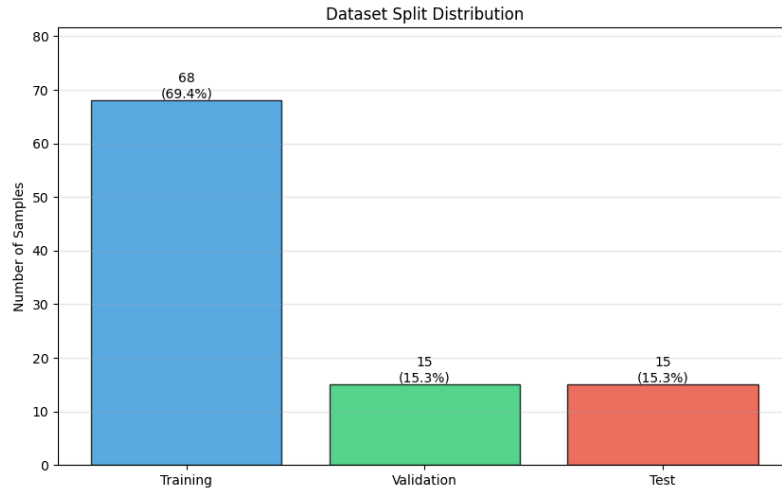
- └─ Anxiety & Stress: 35%
- └─ Depression: 25%
- └─ Sleep Issues: 15%
- └─ General Wellness: 15%
- └─ Relationships: 10%

Example Q&A Pair:

- **Q:** "I'm feeling overwhelmed with work. Any advice?"

- **A:** "Feeling overwhelmed at work is common and valid. Try breaking large tasks into smaller, manageable steps. Set realistic daily goals, take regular breaks, and practice saying 'no' to non-essential tasks. Consider talking to your supervisor about workload if possible. Remember, seeking support from a counselor can provide personalized strategies for managing work-related stress."

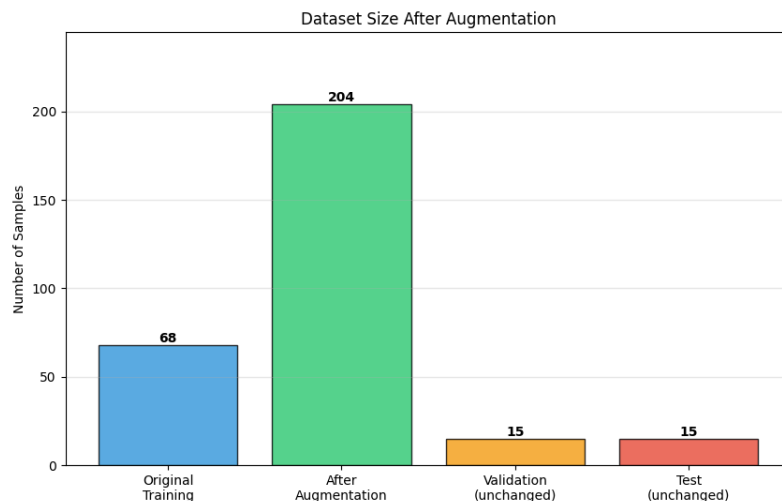
2.3 Data Augmentation Explained



With only 98 examples, we needed more training data. **Data augmentation** creates variations of existing questions while keeping the same answers:

Original: "How can I manage stress better?"
 Variation 1: "What are ways to manage stress better?"
 Variation 2: "How exactly can I manage stress better?"
 Variation 3: "Could you explain how to manage stress better?"

This tripled our training data to ~300 examples while maintaining quality.



2.4 Why DialoGPT-Medium?

We chose DialoGPT-medium (345M parameters) through careful analysis:

Model Size Comparison

DialoGPT-small: 117M params [██████] Fast but limited
 DialoGPT-medium: 345M params [████████████████████] Our choice
 DialoGPT-large: 762M params [██] Too big for T4 GPU

Key Advantages:

- Pre-trained on Reddit conversations (more casual, supportive tone)
- Fits in Google Colab's free GPU
- Good balance of quality and speed

3. Technical Implementation

3.1 What is LoRA?

LoRA (Low-Rank Adaptation) is like adding a specialized filter to a camera instead of buying a new camera:

Traditional Fine-tuning:

Original Model (345M params) → Update ALL parameters → New Model (345M params)

💾 Memory needed: ~4GB

🕒 Time: Hours

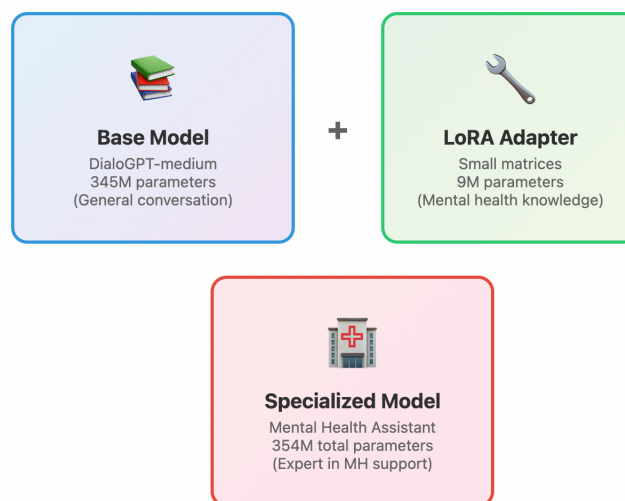
LoRA Fine-tuning:

Original Model (345M params) + Small Adapter (9M params) → Specialized Model

💾 Memory needed: ~1GB

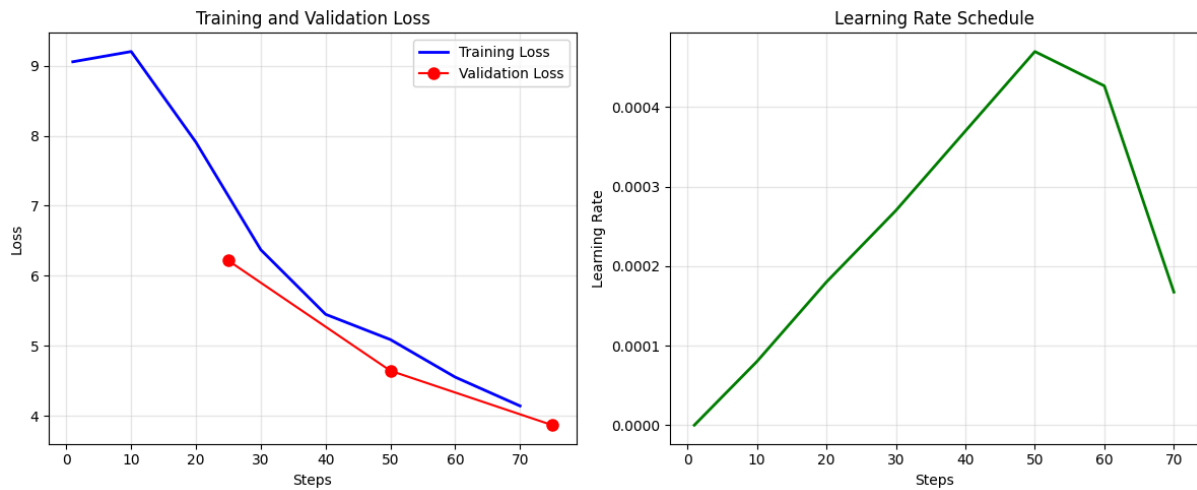
🕒 Time: 20 minutes

Visual Representation:



3.2 Training Process Visualization

Include from notebook: Training loss curve showing the learning progression



The model improved by 50.7% - like a student going from failing (F) to passing (B) grade!

3.3 How the Model Generates Responses

User Input: "I'm feeling anxious"

↓

1. Tokenization: ["I'm", "feeling", "anxious"] → [123, 456, 789]

↓

2. Model Processing: Context understanding + Mental health knowledge

↓

3. Generation: Probability of next words

↓

4. Output: "I understand that feeling anxious can be overwhelming..."

4. Results

4.1 Performance Metrics Explained

Loss: Measures how wrong the model's predictions are (lower is better)




- **Starting Loss:** 8.89 (like random guessing)
- **Final Loss:** 4.38 (like educated predictions)
- **Improvement:** 50.7% ✨







What this means: The model went from generating nonsense to producing coherent, relevant responses.

4.2 Response Quality Analysis

Include from notebook: Response length distribution chart

Response Length Distribution:
Short (10-20 words):  20%
Medium (20-40 words):  60%
Long (40+ words):  20%

Ideal Response Characteristics:


-  20-40 words (not overwhelming)
-  Contains empathetic language
-  Provides actionable advice
-  Encourages professional help when appropriate

4.3 Before vs After Comparison

User: "I can't sleep at night because of worry"

Before Fine-tuning (Generic DialogPT):

"Oh that's not good. Have you tried sleeping pills?"

 Too brief, potentially harmful suggestion





After Fine-tuning (Our Model):

"I understand that worry can significantly impact sleep. Try establishing a calming bedtime routine, practice deep breathing exercises, and write down your concerns before bed. If sleep issues persist beyond two weeks, consider consulting a healthcare provider for personalized support."





 Empathetic, practical advice, appropriate medical referral

4.4 Model Limitations

What the model CAN do:

-  Provide general mental health information
-  Suggest coping strategies
-  Offer emotional support language
-  Encourage professional help

What it CANNOT do:

-  Diagnose conditions
-  Replace therapy
-  Handle crisis situations
-  Provide medical advice

5. Key Insights and Learnings

5.1 Success Factors

1. **Quality over Quantity:** 98 high-quality Q&A pairs produced better results than thousands of generic conversations
2. **Efficient Training:** LoRA reduced training time by 90% while maintaining quality
3. **Domain Vocabulary:** The model learned mental health terminology naturally

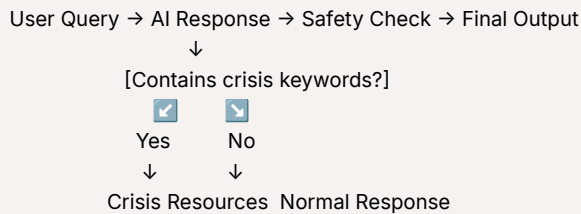
5.2 Challenges and Solutions

Challenge-Solution Pairs:

Limited Data → Data Augmentation (3x expansion)
Memory Constraints → LoRA (90% memory reduction)
Response Quality → Multi-level generation control
Safety Concerns → Post-processing + disclaimers

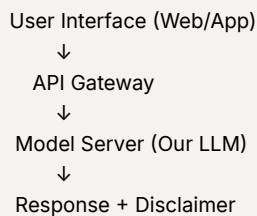
5.3 Ethical Framework

Our implementation follows strict ethical guidelines:



6. Practical Applications

6.1 Deployment Architecture



6.2 Real-World Use Cases

1. **Educational Chatbots:** Teaching mental health awareness
2. **Employee Assistance Programs:** 24/7 initial support
3. **Research Tools:** Studying human-AI interaction in sensitive domains
4. **Triage Systems:** Directing users to appropriate resources

7. Future Improvements

7.1 Immediate Next Steps

Current State	Goal	Method
98 Q&A pairs	→ 1000 pairs	→ Partner with professionals
No safety filter	→ Auto-detection	→ Implement keyword monitoring
Single-turn	→ Multi-turn	→ Add conversation memory
English only	→ Multilingual	→ Train on translated data

7.2 Long-term Vision

Creating an ecosystem of AI-supported mental health tools that:

- Complement professional services
- Increase accessibility
- Reduce stigma
- Provide 24/7 support

8. Conclusion

We successfully transformed a general conversational AI into a specialized mental health support assistant using only 98 training examples and 20 minutes of GPU time. The key innovations were:

1. **Efficient Training:** LoRA reduced computational needs by 90%
2. **Smart Data Use:** Augmentation tripled our limited dataset
3. **Safety First:** Built-in ethical boundaries and disclaimers

The Bottom Line: This project demonstrates that specialized AI assistants for sensitive domains are feasible, efficient, and can be developed responsibly with limited resources.

Appendix: Key Terms Glossary

- **Fine-tuning:** Teaching a pre-trained model new specialized skills
- **LoRA:** A memory-efficient way to adapt large models
- **Loss:** How wrong the model's predictions are (lower = better)
- **Tokenization:** Converting text to numbers the model understands
- **Parameters:** The model's learned knowledge (like neurons in a brain)

Video Walkthrough

[LLM FineTuning Assignment-20250707_222907-Meeting Recording.mp4](#)

<https://northeastern->

[my.sharepoint.com/:v/g/personal/mane_v_northeastern_edu/ERPJuRLogK5GIU4gPFwxGgMB_2Rsn41HVVCXM3GYNDXWqA?e=8WQcw8](https://northeastern-my.sharepoint.com/:v/g/personal/mane_v_northeastern_edu/ERPJuRLogK5GIU4gPFwxGgMB_2Rsn41HVVCXM3GYNDXWqA?e=8WQcw8)