



Fine-Tuning AI Models: Customizing Intelligence

What is Fine-Tuning?

Fine-tuning is the process of taking a pre-trained AI model and further training it on a smaller, task-specific dataset. Think of it like taking a skilled generalist (the pre-trained model) and giving them specialized training to become an expert in a particular field.

Why Fine-Tune?



Task-Specific Performance

Achieve superior accuracy and relevance for specialized tasks.



Cost Efficiency

Leverage existing knowledge, reducing the need for extensive training from scratch.



Data Efficiency

Requires less labeled data than training a new model.



Faster Deployment

Significantly quicker development cycles and deployment.



Bypass Restrictions

Adapt models to handle specific content or ethical guidelines not covered by general training.

RAG vs. Fine-Tuning: Choosing the Right Approach

| Aspect | RAG (Retrieval-Augmented Generation) | Fine-Tuning |
|-------------|--|---|
| Core Idea | Adds an external knowledge retriever to a model | Adjusts model weights using new training data |
| Data Source | Dynamic: pulls relevant info from external documents at inference time | Static: knowledge is baked into the model |
| Use Case | Question answering, chatbots, knowledge-heavy tasks | Specialized tasks, tone adaptation, structured output |
| Cost | Lower training cost, higher inference cost (due to retrieval) | Higher upfront cost (training), cheaper inference |
| Flexibility | Can update knowledge base without retraining | Must retrain to reflect new knowledge |
| Performance | Great for up-to-date or long-tail knowledge | Better for domain-specific logic, style, or behavior |
| Latency | Slower (retriever + generator) | Faster (single model forward pass) |

The choice between RAG and fine-tuning depends on your project's needs: data volatility, desired model behavior, and available resources. Often, a hybrid approach yields the best results.

When to Fine-Tune

Beyond Prompting

- When complex, nuanced understanding of specific data is needed.
- To embed proprietary or highly specialized knowledge into the model.
- For tasks requiring consistent, high-quality output not achievable with prompts alone.

Types of Fine-Tuning

Full Fine-Tuning

Updating all model parameters for maximum adaptation.

Adapter Tuning

Inserting small, trainable modules into the network.

LoRA (Low-Rank Adaptation)

Efficiently updates a small subset of parameters.

RLHF (Reinforcement Learning from Human Feedback)

Aligning model behavior with human preferences.

Fine-Tuning Workflow

1

1. Model Selection

Choose a suitable pre-trained base model.

2

2. Data Preparation

Curate, clean, and format your specific dataset.

3

3. Training & Iteration

Train the model, monitor, and refine.

4

4. Evaluation & Deployment

Assess performance and deploy the fine-tuned model.

Choosing the Right Model Type

Before fine-tuning, a key decision is whether to start with an instruct model or a base model. This choice significantly impacts your workflow and the final model's behavior.

Instruct Models

Pre-trained with built-in instructions, these models are ready-to-use for conversational tasks. They respond effectively to prompts without extensive fine-tuning (e.g., Mistral-Instruct, LLaMA 2 Chat).

Base Models

Original pre-trained versions without instruction fine-tuning. Designed for extensive customization, they are compatible with instruction-style templates (e.g., LLaMA 3, Mistral base).

Multi-modal variant

refers to a version of a model (or dataset, or training configuration) that can process **more than one type of input modality**

GGUF variant

GGUF variants refer to different **quantization levels** of the same model – basically how compressed or memory-efficient the model is.

GGUF is designed to run models **locally** on CPU/GPU

The decision depends on your data's quantity, quality, and type:

1

1,000+ Rows of Data

Best to fine-tune a **base model** for comprehensive adaptation.

2

300–1,000 Rows of High-Quality Data

Both **base or instruct models** are viable options for fine-tuning.

3

Less than 300 Rows

The **instruct model** is typically the better choice to preserve its built-in capabilities.

Dataset Preparation

- **Formatting**: Structure data as input-output pairs or conversational turns. Ensure the dataset format exactly matches the data used for supervised learning, especially for instruct models.
- **Size**: A larger dataset generally leads to better model performance.
- **Quality**: Clean, diverse, and representative data is crucial.
- **Tokenization**: Convert text into numerical tokens compatible with the model.

Dataset Formatting Examples

Instruct Model Format

Typically, instruct datasets follow a simple input-output pair structure:

```
[
  {
    "prompt": "Translate 'Hello' to Spanish.",
    "completion": "Hola"
  },
  {
    "prompt": "Summarize the key points of LLMs.",
    "completion": "Large Language Models are deep learning models trained on vast text data for various NLP tasks."
  }
]
```

Chat Model Format

Chat datasets use a list of message objects, each with a role (user, assistant, system) and content:

```
[
  {
    "messages": [
      {"role": "user", "content": "Tell me a fun fact about cats."},
      {"role": "assistant", "content": "Cats can make over 100 different sounds."}
    ]
  },
  {
    "messages": [
      {"role": "user", "content": "What is the capital of France?"},
      {"role": "assistant", "content": "The capital of France is Paris."}
    ]
  }
]
```

Ensuring your dataset matches the specific format expected by the chosen model type is crucial for effective fine-tuning.

Evaluation Metrics

| | |
|------------------|---|
| Accuracy | Proportion of correct predictions. |
| Perplexity | How well the model predicts the next token. (Lower is better.) |
| BLEU | Measures precision — how much the model's output overlaps with the reference, in terms of n-grams . |
| ROUGE | Measures recall how much of the reference is captured by the generated text. |
| Human Evaluation | Subjective assessment of output quality and relevance. |

Evaluating Your Fine-Tuned Model

Assessing the performance of your fine-tuned AI model is crucial for success. Here are several strategies:

1 Manual Interaction

Chat directly with the model to subjectively gauge its responses and behavior, checking if it meets your expectations and project goals.

2 Unsloth's Built-in Evaluation

Unsloth provides integrated evaluation capabilities. Be aware it can be time-consuming, especially with large datasets, so plan accordingly.

3 Optimize Evaluation Speed

To accelerate Unsloth's evaluation, consider reducing the evaluation dataset size or setting ``evaluation_steps`` to a smaller number like 100.

4 Leveraging Test Data

If available, use a dedicated test dataset (e.g., 20% of your training data) for objective performance measurement.

5 Automated Evaluation Tools

Employ automated tools or benchmarks to systematically assess model quality and compare results against baselines or other models.

Questions & Answers

Thank You!