

# LARGE LANGUAGE MODELS AND PRE-TRAINED MODELS

LLMs are ML algorithms designed to predict subsequent words in sentences.

---

Examples: BERT, Mistral7B, trained on vast text data to understand language fundamentals.

---

Based on the Transformers architecture, utilizing self-attention mechanisms.

---



# Model Performance and Fine-tuning



Models excel at understanding natural language and generating human-like text.



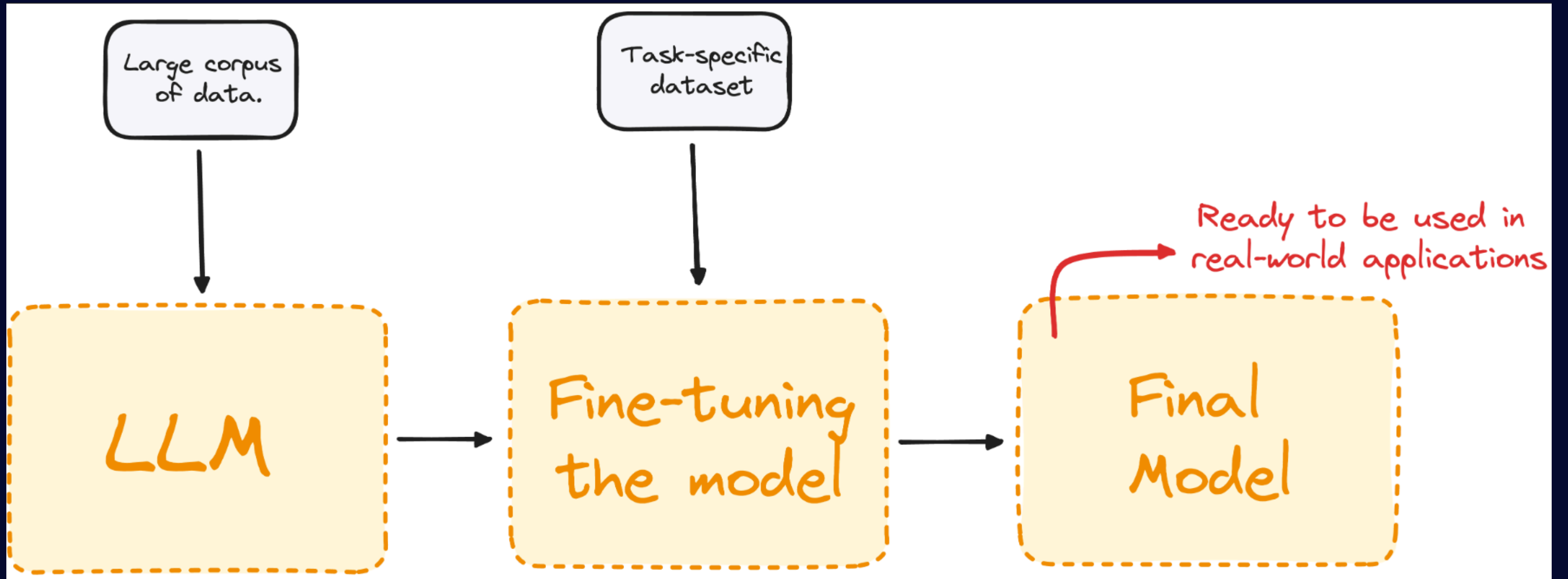
General performance is strong; however, domain-specific tasks can be challenging.



Fine-tuning involves additional training on a domain-specific dataset.



Benefits include reduced computational costs and enhanced model versatility.



# Types of Fine-tuning

---

01

Few-shot learning: Provides initial examples for better contextual understanding.

02

Transfer learning: Applies broad model knowledge to specific tasks.

03

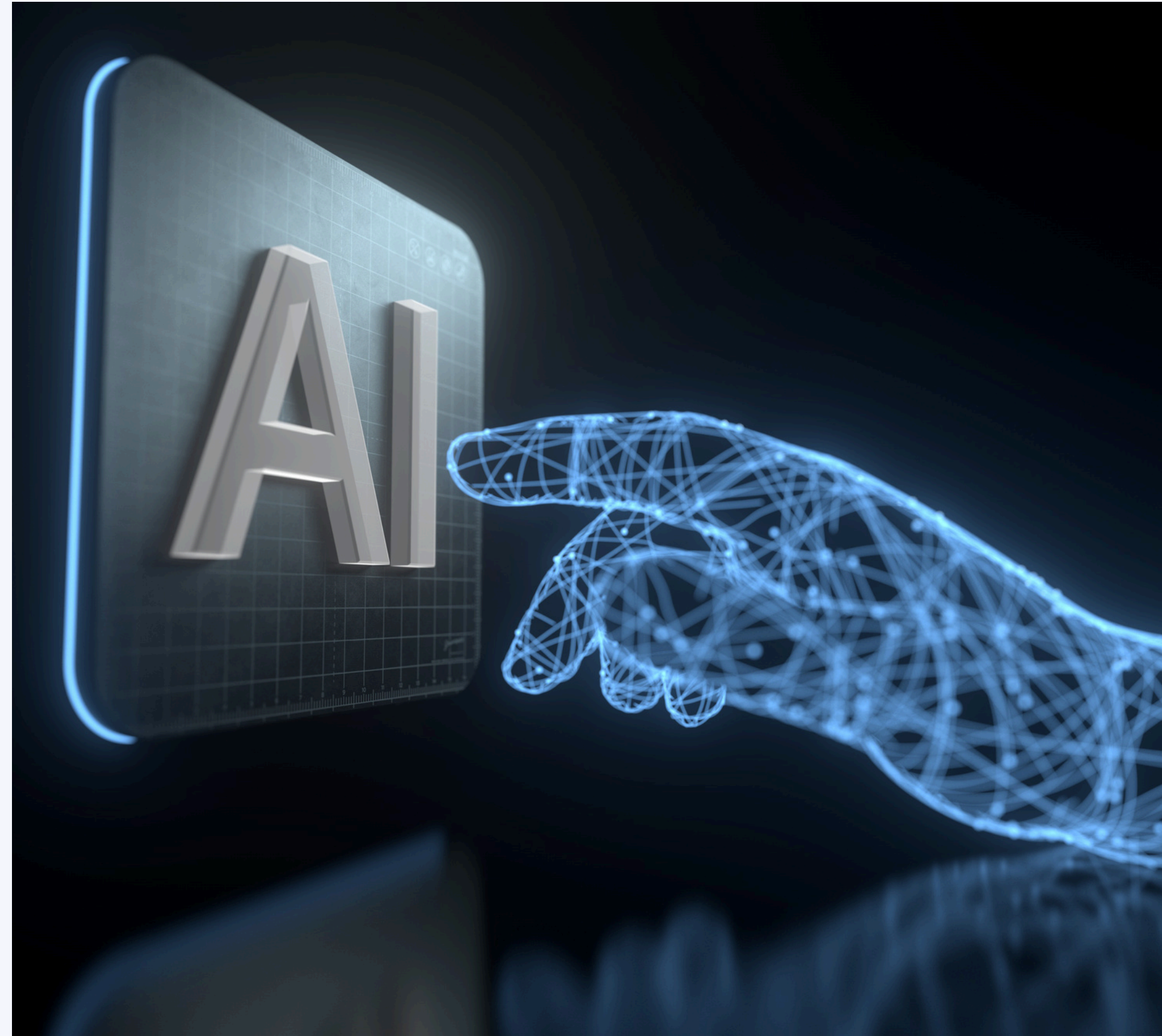
Domain-specific fine-tuning: Adapts model for specific industry needs.

04

Supervised fine-tuning: Uses labeled data for targeted training tasks.

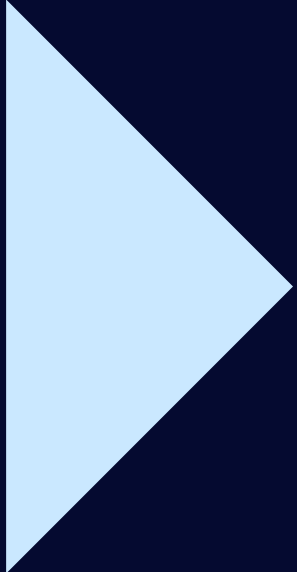
# Introduction to Hugging Face

- An open-source AI company offering tools and models for NLP tasks.
- Provides access to state-of-the-art ML models like the Transformers library.
- Supports community-driven AI development through collaboration on models and datasets.



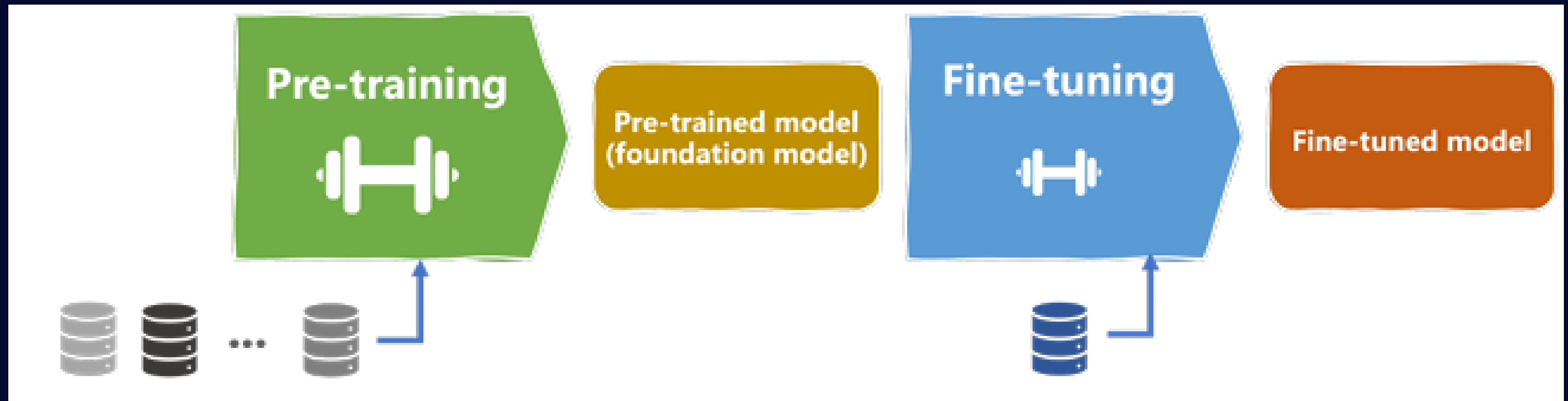
# Generic Modules of the Transformers Library

- 
- ✓ Model: Supports a variety of NLP tasks with different architectures.
- 
- ✓ Tokenizer: Handles conversion of text to tokens.
- 
- ✓ Pipeline: Simplifies performing NLP tasks with pre-trained models.
- 
- ✓ Trainer: Aids in training and fine-tuning models.
-



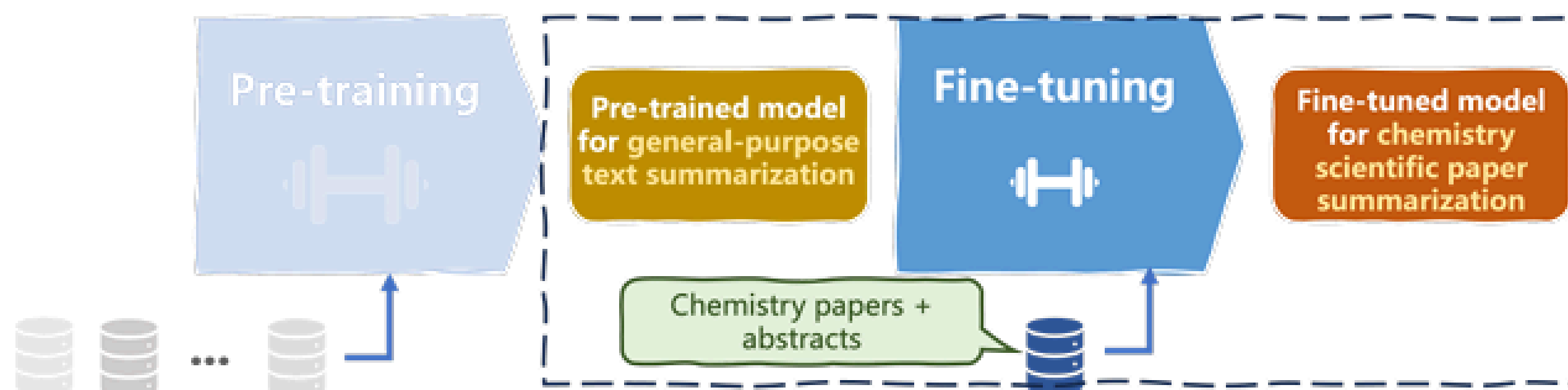
# LLM fine-tuning and transfer learning

# OVERVIEW OF LLM LIFECYCLE

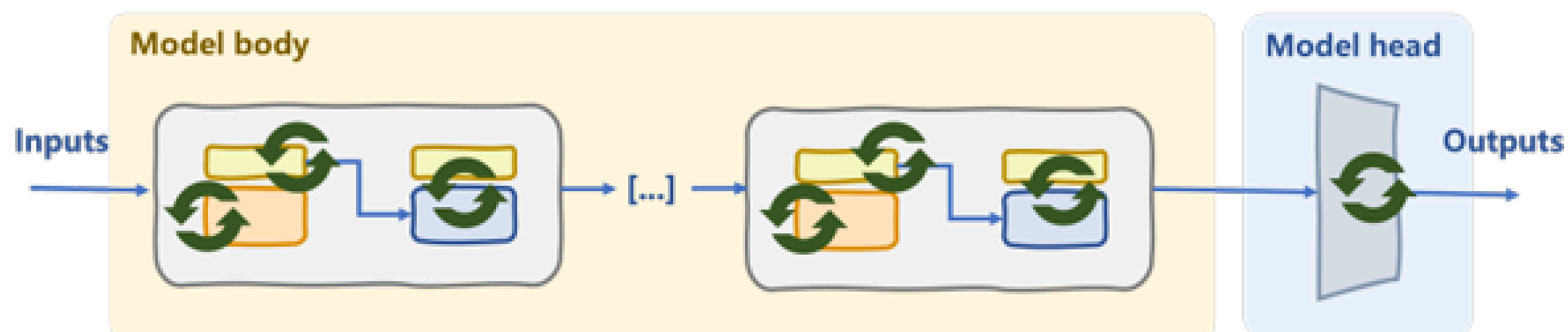




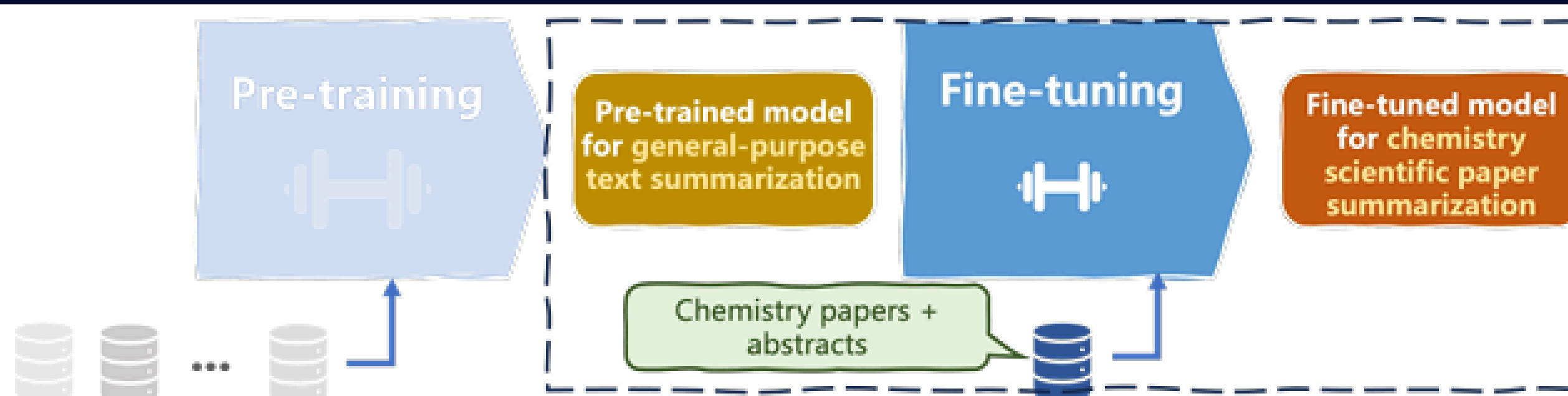
# FULL FINE-TUNING



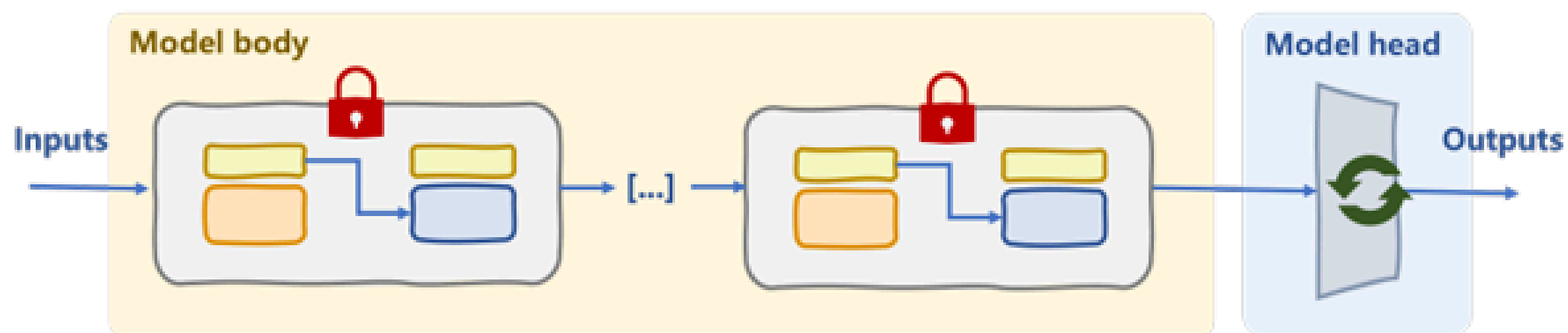
**Full fine-tuning:** The entire model weights are updated; more computationally expensive



# PARTIAL FINE-TUNING



Partial fine-tuning: Lower (body) layers fixed; only task-specific layers (head) are updated



---

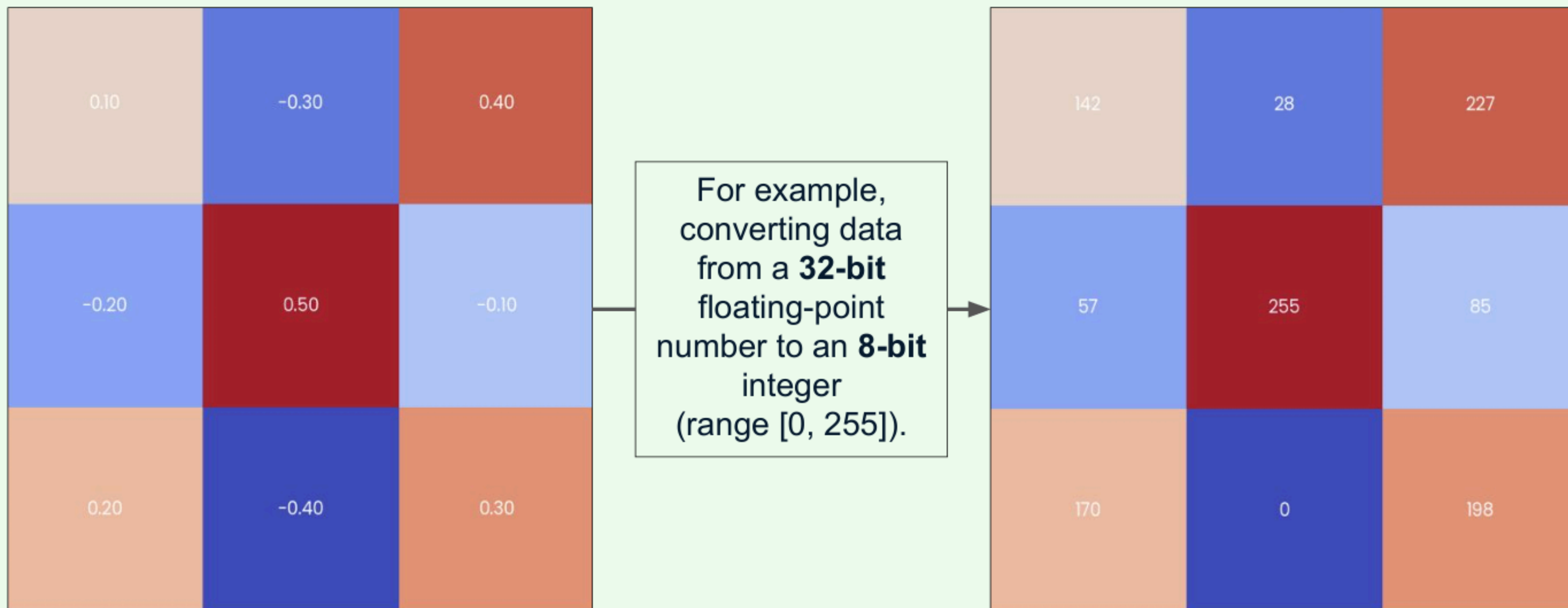
# Transfer Learning

Model trained on one task is adapted for a different but related task  
These models are typically fine-tuned on a smaller, task-specific dataset.  
Types of Transfer Learning:

- Partial Fine-tuning: Adjusting some parts of a pre-trained model.
  - Complete Fine-tuning: Entire model adjustments for new tasks.
  - Zero-shot Learning: Enables the model to tackle tasks it has never explicitly learned during its training phase.
  - One-shot and Few-shot Learning: These methods allow the model to adapt to new tasks using minimal data — as few as one or a handful of examples.
-

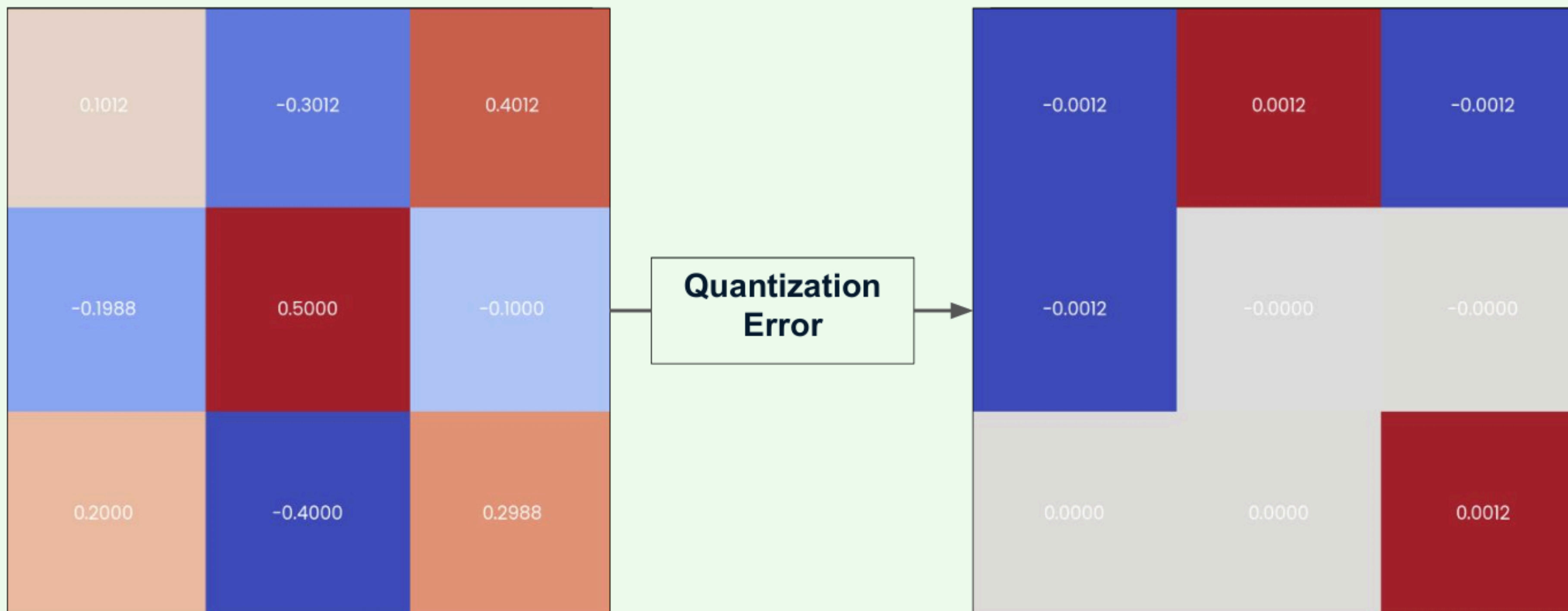
# LLM Quantization

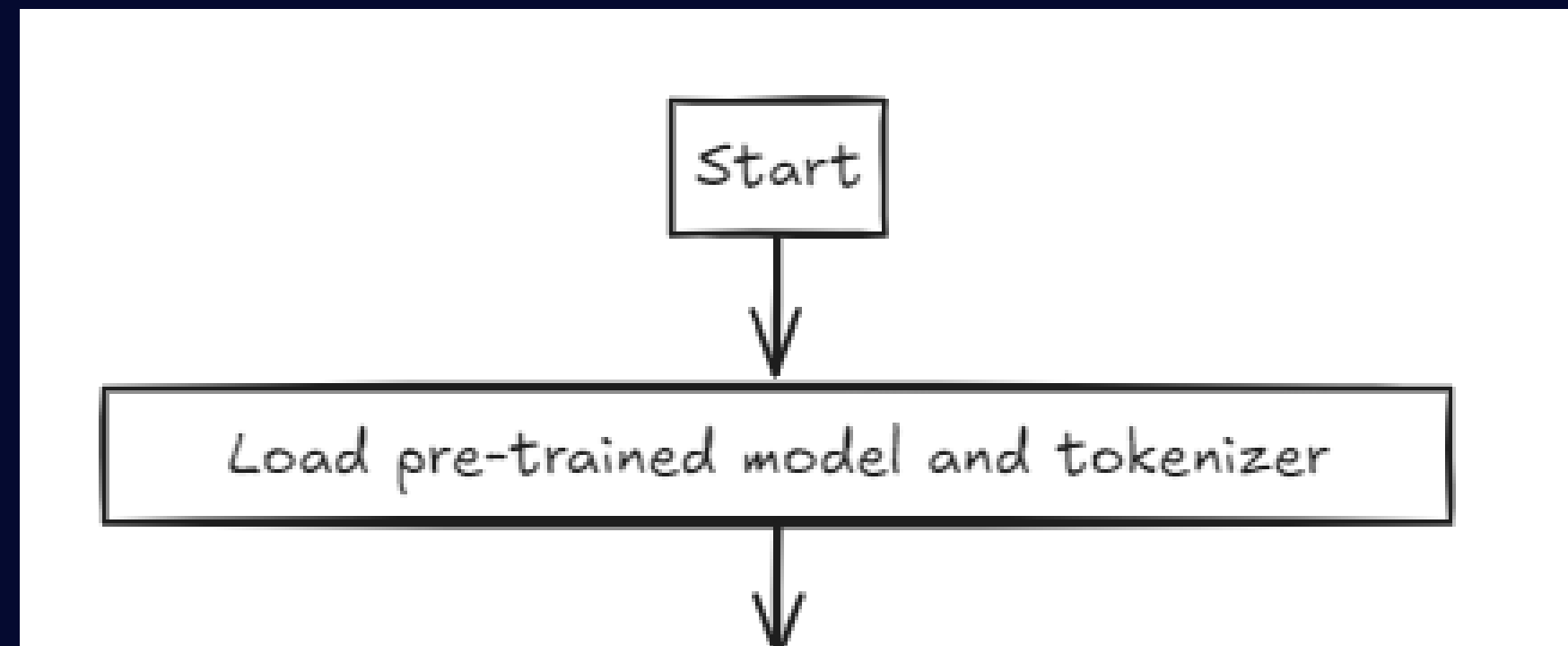
- Model compression technique that converts the **weights and/or activations** from high-precision values to lower-precision ones.




# LLM Quantization

- Quantization shrinks LLMs to **consume less memory, and require less storage space**, while minimizing the performance degradation.





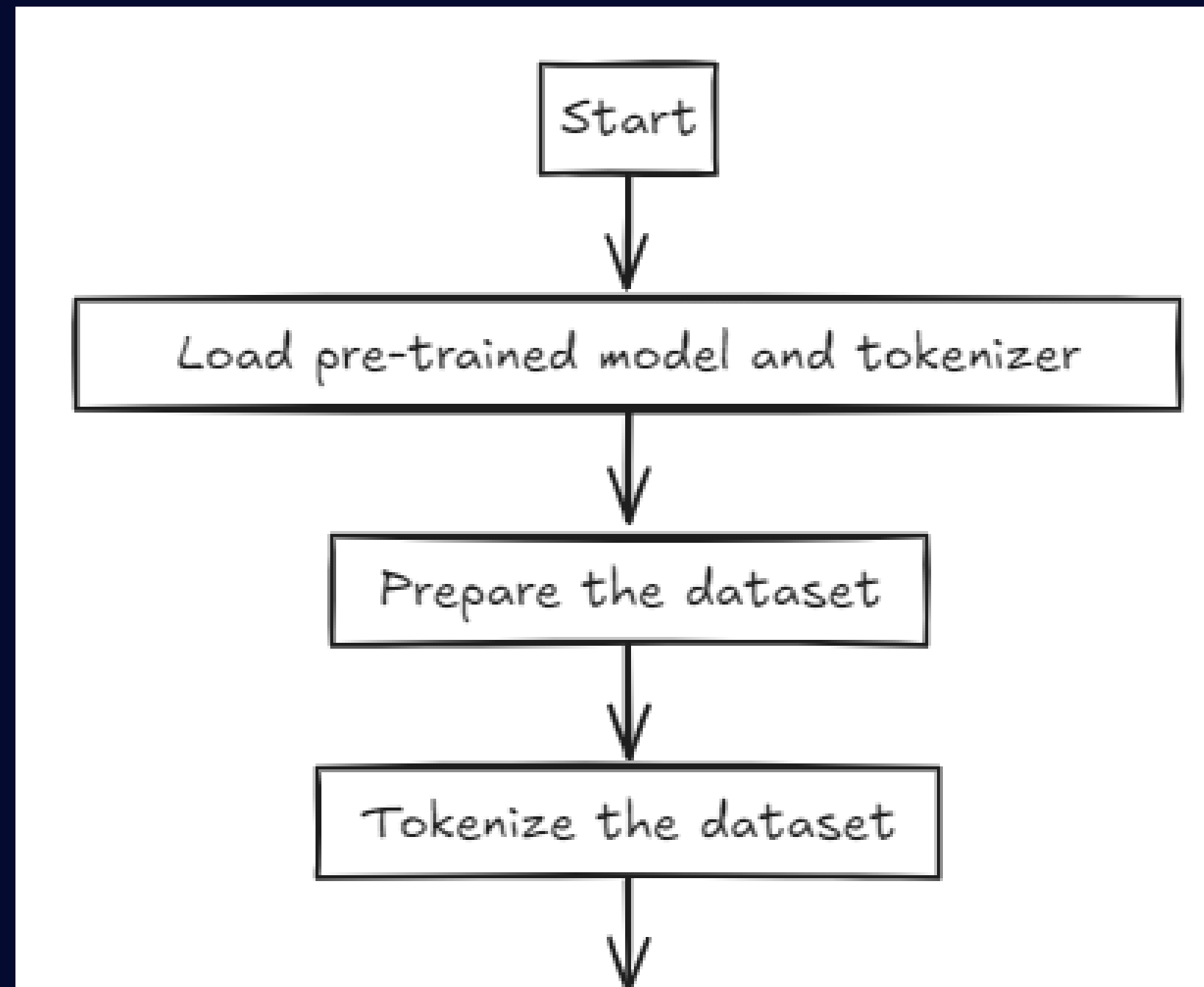
# FINE-TUNING A PRE-TRAINED HUGGING FACE LLM



```
1 import torch
2 from transformers import AutoModelForSequenceClassification, AutoTokenizer
3 from datasets import load_dataset
4
5 model_name = "distilbert-base-uncased"
6 tokenizer = AutoTokenizer.from_pretrained(model_name)
7 model = AutoModelForSequenceClassification.from_pretrained(
8     model_name, num_labels=2)
```

Load BERT-based model for text classification and associated tokenizer

Tokenize dataset used for fine-tuning IMDB reviews dataset



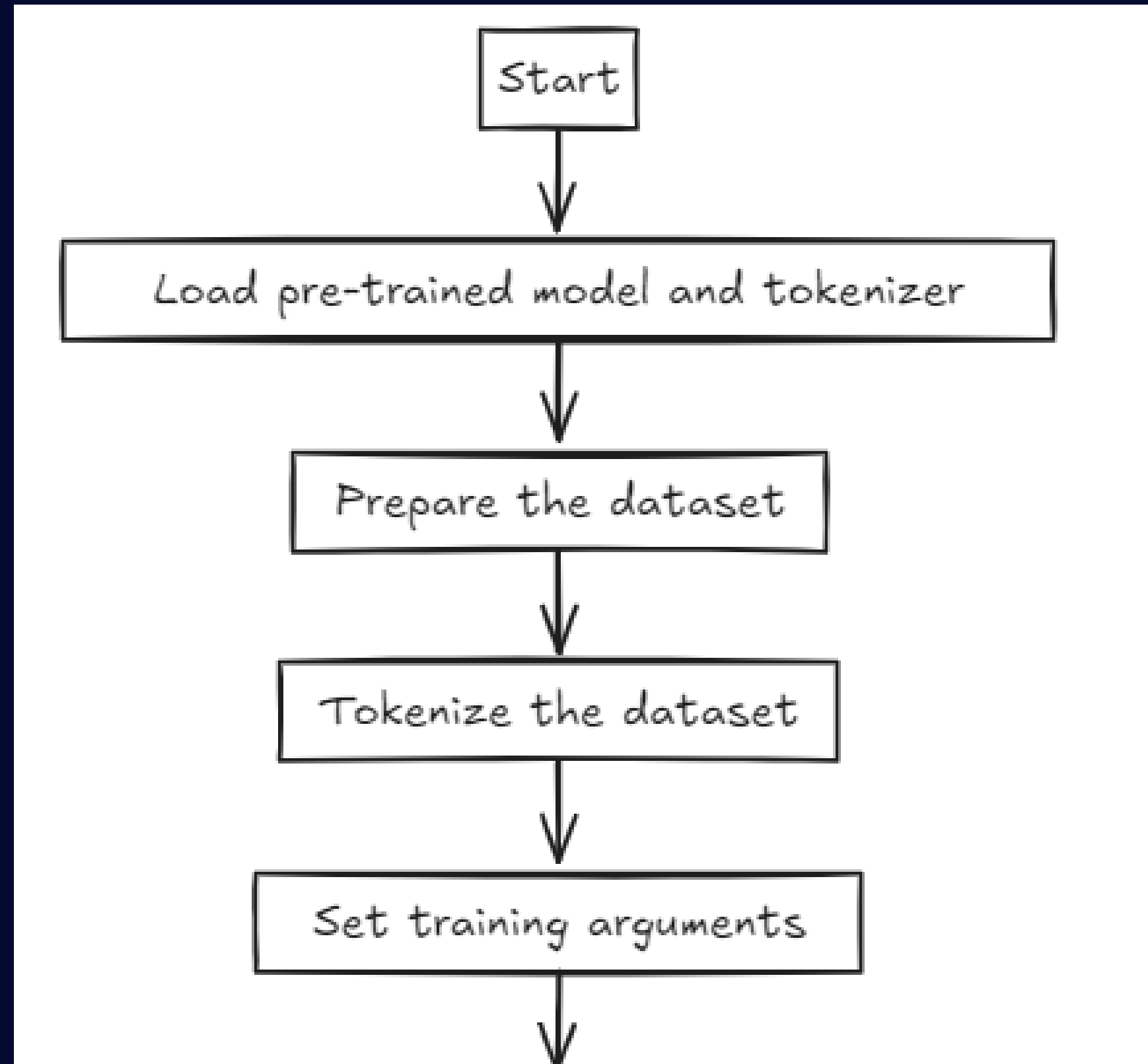




```
1 def tokenize_function(examples):  
2     return tokenizer(  
3         examples["text"], padding="max_length", truncation=True)  
4 data = load_dataset("imdb")  
5 tokenized_data = data.map(tokenize_function, batched=True)
```

truncation= True truncates input sequences beyond model's

max\_length batched= True to process examples in batches rather than individually

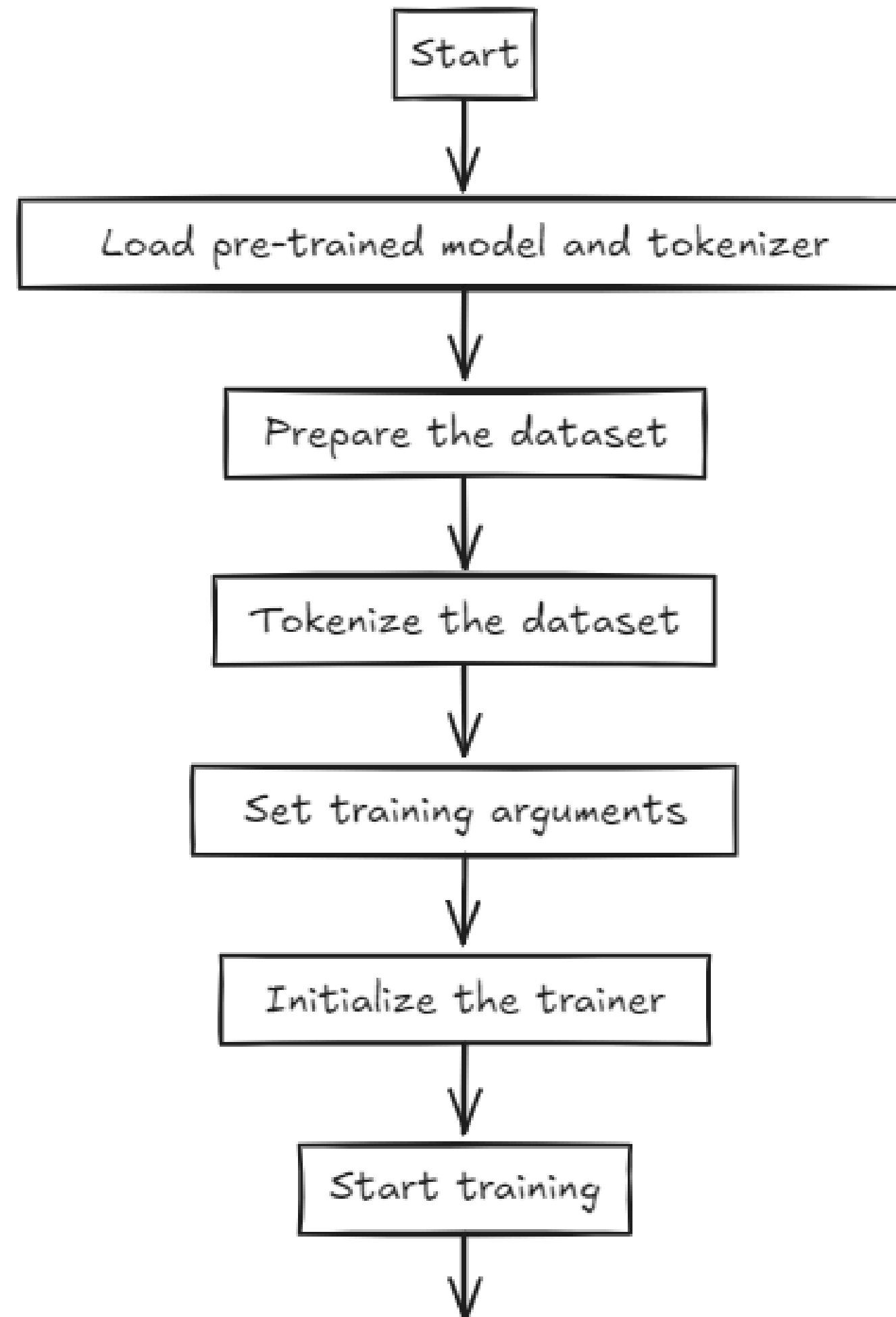




```
1  from transformers import Trainer, TrainingArguments
2  training_args = TrainingArguments(
3  output_dir="./smaller_bert_finetuned",
4  per_device_train_batch_size=8,
5  num_train_epochs=3,
6  evaluation_strategy="steps",
7  eval_steps=500,
8  save_steps=500,
9  logging_dir="./logs",
10 )
11
```

## Training Arguments

**class: customize  
training settings**  
Output directory,  
batch size per GPU,  
epochs, etc.



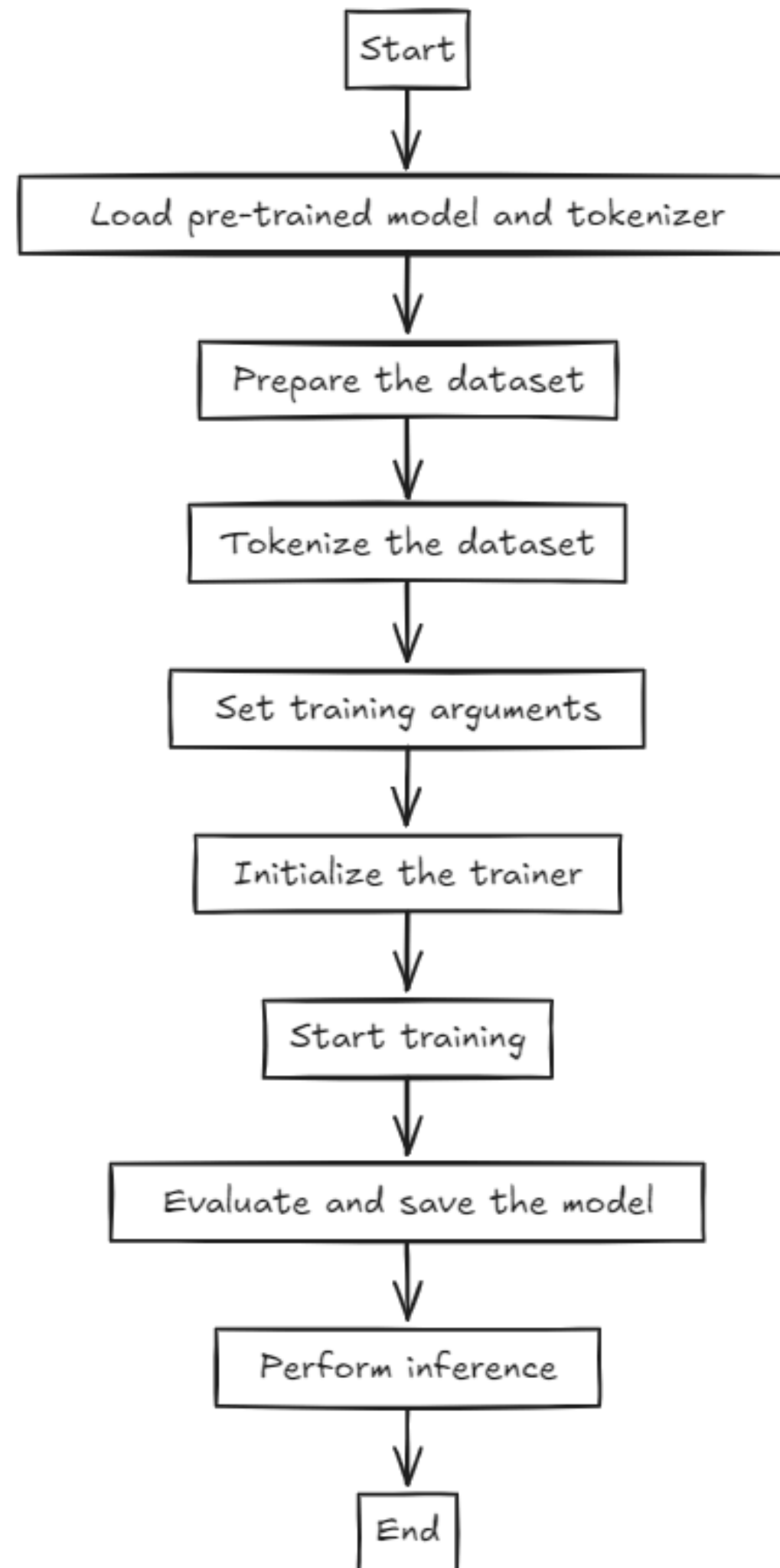


```
1  trainer = Trainer(  
2  model=model,  
3  args=training_args,  
4  train_dataset=tokenized_datasets["train"],  
5  eval_dataset=tokenized_datasets["test"],  
6  )  
7  trainer.train()
```

**Trainer class:**  
**manage training and  
validation loop**

Specify model, training  
arguments, training and  
validation sets

**trainer.train() :**  
**execute training  
loop**



After fine-tuning, inference is performed as usual Tokenize inputs, pass them to the LLM, obtain and post-process outputs



```
1 example_input = tokenizer("I am absolutely amazed with this new and revolutionary AI device",  
2                             return_tensors="pt")  
3  
4 output = model(**example_input)  
5  
6 predicted_label = torch.argmax(output.logits, dim=1).item()  
7 print("Predicted Label:", predicted_label)
```

**Predicted Label: 0**



```
1 model.save_pretrained("./my_bert_finetuned")  
2 tokenizer.save_pretrained("./my_bert_finetuned")
```

**Fine-tuned model and tokenizer can be saved using `.save_pretrained()`**



