5

RAG using Langchain



Lecture Plan







RAG with Langchain



Good to Know: Fine-Tuning



Good to Know: Fine-Tuning Techniques



Good to Know: Fine-Tuning Architectures

Introduction

What is RAG?

- Retrieval-Augmented Generation, or RAG, is a hybrid approach in NLP that augments the generative capabilities of language models with a retrieval mechanism.
- This means that instead of generating responses purely based on pretrained knowledge, a RAG model retrieves relevant information from a large corpus of documents and uses this information to generate more accurate and contextually relevant responses.

Advantages of RAG

- RAG offers several advantages over traditional language models, including:
 - Enhanced Accuracy: By retrieving relevant documents, RAG models can provide more precise and factual responses.
 - Improved Contextual Understanding: RAG leverages external knowledge sources, enabling it to handle a wider range of topics and queries.
 - Efficiency in Handling Large Datasets: The retrieval component can quickly narrow down the relevant information, making the generation process more efficient.

Real-World Applications

- Some of the real-world applications where RAG proves to be incredibly useful include:
 - Question-Answering Systems: Providing accurate answers by retrieving relevant context.
 - Customer Support Bots: Enhancing the ability to resolve customer queries by pulling information from vast knowledge bases.
 - Content Generation: Generating articles, summaries, or reports by integrating up-to-date information from a variety of sources

Retrieval-Based Systems

- Systems that retrieve relevant information from a predefined dataset or corpus to answer queries.
 - Examples: Search engines (e.g., Google), Information retrieval systems in libraries, FAQ bots that fetch pre-written answers, etc.
- Strengths
 - · Efficient at retrieving exact matches or highly relevant documents
 - Scalable and can handle large datasets
 - High precision in responses
- Limitations
 - Limited to the information available in the dataset
 - May struggle with understanding context or generating nuanced responses
 - Not effective for generating novel content



Generation-Based Systems

- Systems that generate text or content based on pre-trained language models, often without direct reliance on a fixed corpus.
 - Examples: GPT-4 for text completion and generation, Text summarization tools, Chatbots that generate responses dynamically, etc.
- Strengths
 - Capable of creating original content
 - Flexible in handling a wide range of queries and topics
 - Can generate contextually relevant and coherent responses
- Limitations
 - May produce inaccurate or nonsensical responses
 - Can struggle with factual accuracy
 - Requires significant computational resources for training and inference

Combining Retrieval and Generation

- How RAG Leverages Both Retrieval and Generation:
 - Retrieval Component:
 - Retrieves relevant documents or passages from a large corpus.
 - Provides context and factual information to support the generation process.
 - Generation Component:
 - Uses the retrieved information to generate accurate and contextually relevant responses.
 - Enhances the generative model's output by grounding it in real data.
- Benefits of Combining Approaches:
 - Enhanced Accuracy:
 - Combines the precision of retrieval with the flexibility of generation to produce more accurate responses.

Combining Retrieval and Generation

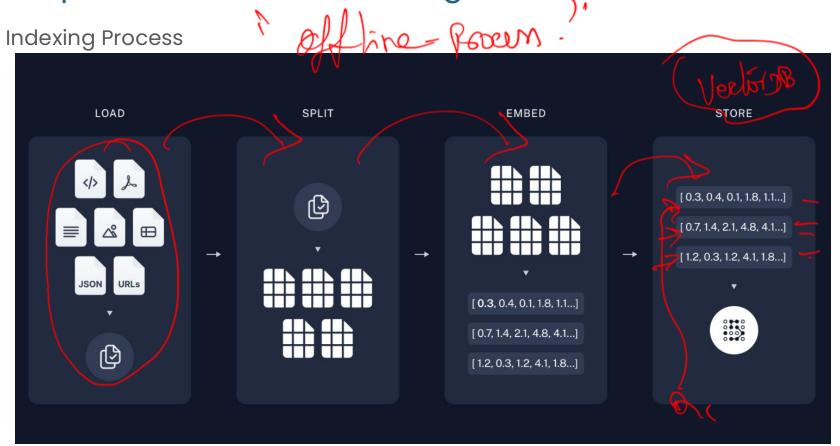
- Benefits of Combining Approaches:
 - Improved Contextual Understanding:
 - Leverages external knowledge to handle complex and context-specific queries.
 - Efficiency:
 - Retrieval narrows down the information space, making the generation process more focused and efficient.
 - Scalability:
 - Can handle large datasets and complex queries without extensive pretraining.
 - Versatility:
 - Suitable for a wide range of applications, from customer support to content creation.

RAG Architecture Overview

- Two Main Components:
 - Indexing:
 - Pipeline for ingesting and indexing data from a source.
 Typically performed offline: 100 2001
 - Typically performed offline;
 - Retrieval and Generation:
 - Process that runs at query time.
 - Retrieves relevant data from the index and passes it to the model for generating responses.

Indexing Process

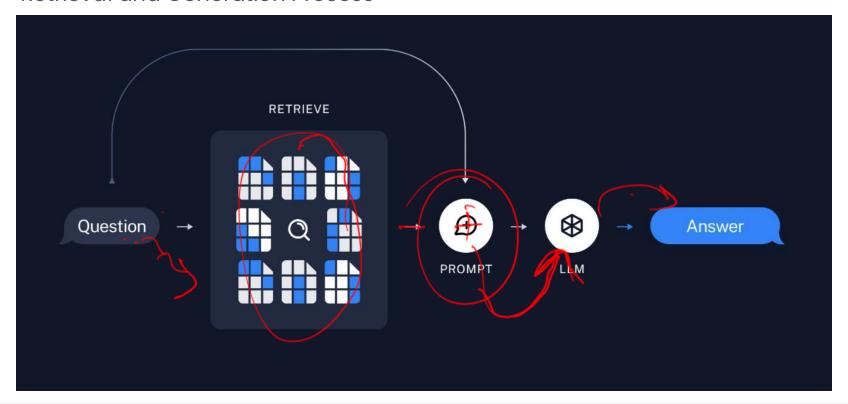
- Load:
 - Use DocumentLoaders to load data from various sources.
 - Ensure data is properly formatted for subsequent steps.
- Split:
 - Text Splitters break large documents into smaller, manageable chunks.
 - Smaller chunks are easier for indexing and fit within the model's context window.
 - Store:
 - Store and index the splits using a VectorStore and Embeddings model.
 - Ensures data is organized and searchable for retrieval.



Retrieval and Generation Process

- Retrieve
 - Given a user query, use a Retriever to fetch relevant splits from the storage.
 - Efficiently narrows down the search space to the most pertinent information.
- Generate
 - Use a <u>ChatModel</u> / LLM to produce an answer.
 - Prompt includes both the user question and the retrieved data for contextually accurate responses.

Retrieval and Generation Process



Step-by-Step Overview

- Indexing:
 - Load: Import data using DocumentLoaders.
 - Split: Break down data into chunks using Text Splitters.
 - Store: Index and store chunks in a VectorStore.
- Retrieval and Generation:
 - Retrieve: Fetch relevant data splits using a Retriever based on user input.
 - Generate: Generate a response using a ChatModel/LLM with the retrieved context.

Benefits of RAG with Langchain

- Accuracy
 - Combines precise retrieval with generative capabilities for accurate responses.
- Efficiency
 - Streamlined indexing and retrieval processes ensure quick and relevant data access.
- Scalability
 - Capable of handling large datasets and complex queries.
- Flexibility
 - Adapts to various applications, including search engines, customer support, and content creation.

Practical Implementation

- Required packages:
 - langchain, langchain-community, langchainhub, langchain-openai, langchain-chroma, bs4
 - Set OpenAl API key environment variable
 - LangSmith (optional but helpful for inspecting complex applications)
 - Sign up for LangSmith Set environment variables to start logging traces: LANGCHAIN_TRACING_V2 = "true" LANGCHAIN_API_KEY =
 "your api key"

Follow: https://python.langchain.com/v0.1/docs/use_cases/question_answering/quickstart/



Summary of Key Points

LangChain

- A powerful library for building applications with LLMs Modular components for data ingestion, transformation, indexing, retrieval, and generation
- Abstractions for working with different LLMs, embeddings, and vector stores
- Question Answering (QA) Architecture
 - Indexing pipeline: Load data, split into chunks, store embeddings in vector database
 - Retrieval and generation: Retrieve relevant chunks, generate answer using LLM

Summary of Key Points

- Practical Implementation Walkthrough
 - Setting up dependencies and environment variables
 - Loading and parsing data using DocumentLdaders
 - Splitting documents into chunks using TextSplitters
 - Storing embeddings using VectorStores
 - Retrieving relevant documents using Retrievers
 - Generating answers using LLMs or ChatModels with prompts



6

Fine-tuning



Definition and Purpose of fine-tuning:

- Fine-tuning is the process of adapting a pre-trained model to a specific downstream task by further training the model on a smaller, taskspecific dataset.
- The purpose of fine-tuning is to leverage the knowledge and representations learned by the pre-trained model and tailor it to the target task, improving performance and efficiency.
- Example: Taking a pre-trained language model like BERT and fine-tuning it for sentiment analysis on a dataset of movie reviews.

Differences between pre-training and fine-tuning

- Pre-training:
 - √ Involves training a model on a large, general-purpose dataset to learn generic representations and patterns.
 - The model is trained on tasks like language modeling or masked language modeling to capture general language understanding.
 - Example: Training BERT on a massive corpus of text data, such as Wikipedia and BookCorpus, to learn contextual word embeddings.
- Fine-tuning:
 - Involves adapting a pre-trained model to a specific task by further training it on a smaller, task-specific dataset.
 - The model's architecture and learned representations are fine-tuned to optimize performance on the target task.
 - Example: Fine-tuning BERT on a sentiment analysis dataset to classify movie reviews as positive or negative.

Advantages of fine-tuning in GenAl

- Transfer learning
 - Fine-tuning enables the model to leverage the knowledge and representations learned during pre-training.
 - It allows the model to transfer learning from the pre-training task to the target task, reducing the need for extensive task-specific training data.
 - Example: A pre-trained language model like GPT-3 can be fine-tuned for various tasks like text generation, question answering, and text classification with minimal additional training.
- Improved performance
 - Fine-tuning often leads to superior performance compared to training a model from scratch on the target task.
 - The pre-trained model's knowledge and representations provide a strong starting point, allowing the fine-tuned model to achieve better results with less training data and time.



Advantages of fine-tuning in GenAl

- Efficiency and resource optimization
 - Fine-tuning is more efficient than training a model from scratch, as it requires fewer computational resources and training time.
 - By leveraging pre-trained models, fine-tuning reduces the need for extensive computational power and large-scale training data.
 - Example: Fine-tuning a pre-trained model like RoBERTa on a specific task can be done with a single GPU in a matter of hours, compared to the days or weeks required to train a model from scratch.
- Adaptability to various domains and tasks
 - Fine-tuning allows pre-trained models to be adapted to a wide range of domains and tasks.
 - The same pre-trained model can be fine-tuned for multiple downstream tasks, enabling versatility and flexibility in GenAl applications.
 - Example: A pre-trained language model like T5 can be fine-tuned for tasks such as





Supervised Fine-Tuning

- Task-Specific Fine-Tuning
 - Definition: Adapting a pre-trained model to a specific downstream task using labeled data.
 - Example: Fine-tuning BERT for sentiment analysis on a dataset of movie reviews labeled as positive or negative.
 - Process:
 - Modify the model architecture by adding task-specific layers (e.g., classification head).
 - Train the model on the labeled task-specific dataset.
 - Optimize the model's parameters to minimize the task-specific loss function.
 - Benefits: Improved performance on the target task compared to training from scratch.

Supervised Fine-Tuning

- Domain-Specific Fine-Tuning Definition
 - Adapting a pre-trained model to a specific domain or industry using domainspecific data.
 - Example: Fine-tuning GPT-3 on a dataset of legal documents to generate legal contract drafts.
 - Process:
 - Collect a domain-specific dataset that represents the target domain.
 - Fine-tune the pre-trained model on the domain-specific dataset.
 - Evaluate the model's performance on domain-specific tasks or metrics.
 - Benefits: Enhanced model performance and generation quality within the target domain.

Unsupervised Fine-Tuning

- Masked Language Modeling (MLM)
 - Definition: Fine-tuning a pre-trained model using masked language modeling objective on unlabeled data.
 - Example: Fine-tuning BERT on a large corpus of unlabeled text data by masking random tokens and training the model to predict the masked tokens.
 - Process:
 - Randomly mask a percentage of tokens in the input sequence.
 - Train the model to predict the masked tokens based on the surrounding context.
 - Update the model's parameters to minimize the MLM loss.
 - Benefits: Improved language understanding and generalization without the need for labeled data

Unsupervised Fine-Tuning

- Next Sentence Prediction (NSP)
 - Definition: Fine-tuning a pre-trained model to predict whether two sentences follow each other in a coherent manner
 - Example: Fine-tuning BERT on a dataset of sentence pairs to determine if the second sentence is a plausible continuation of the first sentence.
 - Process:
 - Prepare a dataset of sentence pairs, with positive pairs from the same document and negative pairs from different documents.
 - Train the model to classify whether the second sentence follows the first sentence.
 - Update the model's parameters to minimize the NSP loss.
 - Benefits: Enhanced understanding of sentence coherence and context.

Transfer Learning and Adaptation

- Leveraging knowledge learned from one task or domain to improve performance on a related task or domain.
- Example: Fine-tuning a pre-trained image classification model (e.g., ResNet) on a specific medical image dataset to detect skin lesions.
- Process:
 - Identify a pre-trained model that is relevant to the target task or domain.
 - Modify the model architecture if necessary (e.g., changing the output layer).
 - Fine-tune the model on the target dataset using transfer learning techniques.
- Benefits: Faster convergence, improved performance, and reduced need for large-scale task-specific data.

Multi-Task Fine-Tuning

- Fine-tuning a pre-trained model on multiple related tasks simultaneously.
- Example: Fine-tuning a pre-trained language model (e.g., T5) on multiple NLP tasks such as sentiment analysis, named entity recognition, and question answering.
- Process:
 - Prepare datasets for multiple related tasks.
 - Modify the model architecture to accommodate multiple task-specific output layers.
 - Train the model on all tasks simultaneously, optimizing a combined loss function.
- Benefits: Improved generalization, knowledge sharing across tasks, and efficient use of computational resources.



Fine-Tuning Architectures

Fine-Tuning Architectures

Transformer-based Architectures

- BERT (Bidirectional Encoder Representations from Transformers)
 - Overview:
 - Pre-training: Masked Language Modeling (MLM) and Next Sentence Prediction (NSP)
 - Fine-tuning: Sequence classification, token classification, question answering, etc.
 - Architecture:
 - Encoder-only transformer architecture
 - Self-attention mechanism to capture bidirectional context
 - Tokenization using WordPiece embedding
 - Example: Fine-tuning BERT for named entity recognition (NER) on the CoNLL-2003 dataset, achieving state-of-the-art performance.

Fine-Tuning Architectures

Transformer-based Architectures

- GPT (Generative Pre-trained Transformer)
 - Overview.
 - Pre-training: Causal Language Modeling (CLM)
 - Fine-tuning: Language generation, text completion, dialogue systems, etc.
 - Architecture:
 - Decoder-only transformer architecture
 - Self-attention mechanism to capture unidirectional context
 - Byte Pair Encoding (BPE) tokenization
 - Example: Fine-tuning GPT-3 for creative writing tasks, such as generating poetry or short stories, by providing a prompt and letting the model generate coherent and diverse continuations.

Fine-Tuning Architectures

Adapter-based Fine-Tuning

- Adding small, task-specific modules (adapters) to a pre-trained model
- Freezing the pre-trained model's parameters and only training the adapters
- Benefits:
 - Parameter efficiency: Adapters have fewer parameters compared to finetuning the entire model Modularity: Multiple adapters can be trained for different tasks and combined as needed
 - Preservation of pre-trained knowledge: The original model's parameters remain unchanged
 - Example: Using Adapter-based fine-tuning with BERT for multi-task learning, where separate adapters are trained for tasks like sentiment analysis, question answering, and text classification, enabling efficient sharing of the pre-trained model across tasks.

Fine-Tuning Architectures

Prompt-based Fine-Tuning

- Formatting the input data as a prompt to guide the model's output
- Leveraging the model's pre-trained knowledge to perform tasks with minimal fine-tuning
- Benefits:
 - Few-shot learning: Prompt-based fine-tuning enables models to perform tasks with very few training examples
 - Flexibility: Prompts can be easily modified to adapt to different tasks or domains
 - Reduced fine-tuning cost: Prompt-based fine-tuning often requires fewer epochs and resources compared to traditional fine-tuning
 - Example: Using prompt-based fine-tuning with GPT-3 for question answering, where the input question is formatted as a prompt (e.g., "Q: What is the capital of France? A:") and the model generates the answer based on its pre-trained knowledge, requiring minimal fine-tuning on a small set of



Characteristics of Fine-Tuning Datasets

- Domain relevance: The dataset should be representative of the target domain or task.
- Size: Fine-tuning datasets are typically smaller compared to pre-training datasets.
- Quality: High-quality, clean, and labeled data is essential for effective fine-tuning.
- Diversity: The dataset should cover a wide range of examples and variations within the target domain.
- Example: The Stanford Sentiment Treebank (SST) dataset for finetuning sentiment analysis models, containing movie reviews labeled with sentiment scores.

Data Preprocessing and Augmentation Techniques

Tokenization

- Breaking down text into smaller units (tokens) suitable for the model's input.
- Techniques include word-based, subword-based (e.g., WordPiece, BPE), and character-based tokenization.
- Example: Using the BERT tokenizer to preprocess text data, handling out-ofvocabulary words and special tokens.

Data Cleaning:

- Removing noise, irrelevant information, and formatting issues from the dataset.
- Techniques include removing HTML tags, converting to lowercase, and handling special characters.
- Example: Cleaning a web-scraped dataset by removing HTML tags, converting text to lowercase, and handling Unicode characters.



Data Preprocessing and Augmentation Techniques

- Data Augmentation
 - Generating additional training examples by applying transformations to the existing data.
 - Techniques include synonym replacement, random insertion/deletion, and back-translation.
 - Example: Applying synonym replacement to augment a sentiment analysis dataset, replacing words with their synonyms to create new training examples.

Evaluation Metrics for Fine-Tuned Models

- Accuracy: The proportion of correctly predicted instances out of the total instances.
- Suitable for balanced datasets and tasks with equally important classes.
- Example: Evaluating a fine-tuned model's accuracy on a binary sentiment classification task.
- Precision, Recall, and F1 Score:
 - Precision: The proportion of true positive predictions out of the total positive predictions.
 - Recall: The proportion of true positive predictions out of the total actual positive instances.
 - F1 Score: The harmonic mean of precision and recall, providing a balanced measure.
 - Suitable for imbalanced datasets and tasks with different class importances.



Evaluation Metrics for Fine-Tuned Models

- Perplexity
 - A measure of how well a language model predicts the next word in a sequence.
 - Lower perplexity indicates better language modeling performance.
 - Example: Evaluating the perplexity of a fine-tuned GPT model on a domain-specific text generation task



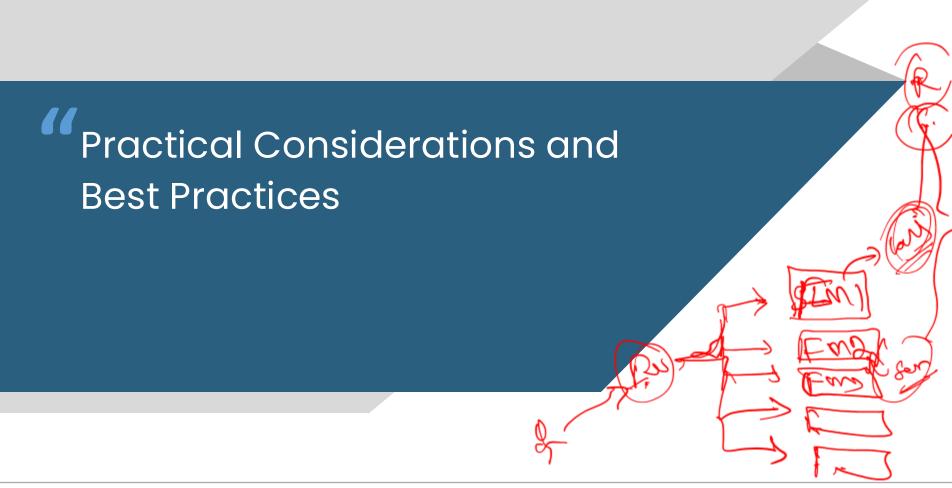
Challenges in Fine-Tuning Evaluation

- Domain Shift:
 - The mismatch between the training data distribution and the target domain distribution.
 - Models may perform well on the fine-tuning dataset but poorly on real-world data.
 - Example: A sentiment analysis model fine-tuned on movie reviews may struggle with customer reviews of products.
- Overfitting:
 - Models may memorize the training data and fail to generalize to unseen examples.
 - Regularization techniques like dropout, early stopping, and cross-validation can mitigate overfitting.
 - Example: A fine-tuned question-answering model that performs well on the training set but fails to answer questions from a different distribution.



Challenges in Fine-Tuning Evaluation

- Evaluation Bias:
 - The evaluation dataset may not fully represent the target domain or contain biases.
 - Models may inherit biases from the training data, leading to unfair or inaccurate predictions.
 - Example: A fine-tuned sentiment analysis model that exhibits gender bias, associating certain professions with positive or negative sentiments based on gender stereotypes.



Hyperparameter Tuning

- The process of selecting optimal hyperparameters for a model to achieve the best performance on a specific task.
- Key Hyperparameters
 - Learning rate: Controls the step size at which the model's weights are updated during training.
 - Batch size: The number of training examples used in each iteration.
 - Number of epochs: The number of times the model goes through the entire training dataset.
- **Techniques**
 - Grid search: Exhaustively searching through a specified subset of the hyperparameter space.
 - Random search: Randomly sampling hyperparameter values from a defined distribution.



Regularization Techniques

Dropout

- Randomly dropping out (setting to zero) a portion of the neurons during training to prevent overfitting.
- Encourages the model to learn more robust and generalizable features.
- Example: Applying dropout with a rate of 0.5 to the fully connected layers of a fine-tuned model to reduce overfitting on a small dataset.

Weight Decay

- Adding a regularization term to the loss function that penalizes large weights.
- Helps prevent the model from overfitting by encouraging smaller weights.
- Example: Using L2 regularization (weight decay) with a coefficient of 0.01 during fine-tuning to control model complexity and improve generalization.

Dealing with Overfitting and Underfitting

Overfitting

- The model performs well on the training data but fails to generalize to unseen data.
- Symptoms: High training accuracy, low validation accuracy, and large gap between training and validation loss.
- Solutions: Increasing regularization, using dropout, early stopping, and data augmentation.

Underfitting

- The model fails to capture the underlying patterns in the training data.
- Symptoms: Low training accuracy, low validation accuracy, and small gap between training and validation loss.
- Solutions: Increasing model capacity, training for more epochs, and using a more complex architecture.
- Example: Identifying overfitting in a fine-tuned sentiment analysis model by © Prof. Sashikumaar Ganesan itoring the training and validation accuracy curves. ISc. Bangalore J. Zenteig Aitech Innovations



Model Selection and Early Stopping

Model Selection

- Comparing the performance of different fine-tuned models and selecting the best one based on validation metrics.
- Techniques: Holdout validation, k-fold cross-validation, and leave-one-out cross-validation.
- Example: Using 5-fold cross-validation to compare different fine-tuned BERT models for a named entity recognition task and selecting the model with the highest_average F1 score.

Early Stopping

- Monitoring the model's performance on a validation set during training and stopping the training process when the performance starts to degrade.
- Helps prevent overfitting and saves computational resources.
- Example: Implementing early stopping with a patience of 3 epochs, where training is stopped if the validation loss does not improve for 3 consecutive



Hardware and Computational Requirements

- GPU Acceleration
 - Leveraging GPU hardware to accelerate the training and inference of finetuned models.
 - Popular GPU platforms: NVIDIA CUDA, AMD ROCm, and Google TPUs.
 - Example: Using an NVIDIA GPU with CUDA support to speed up the fine-tuning of a large transformer-based model.
- Memory Considerations
 - Fine-tuning large models requires significant memory resources, especially for larger batch sizes.
 - Techniques: Gradient accumulation, model parallelism, and efficient memory management.
 - Example: Using gradient accumulation to enable fine tuning to large model with a batch size that exceeds the available GPU memory

Hardware and Computational Requirements

GPU Acceleration

- Leveraging GPU hardware to accelerate the training and inference of finetuned models
- Popular GPU platforms: NVIDIA CUDA, AMD ROCm, and Google TPUs.
- Example: Using an NVIDIA GPU with CUDA support to speed up the finetuning of a large transformer-based model.

Memory Considerations

- Fine-tuning large models requires significant memory resources, especially for larger batch sizes.
- Techniques: Gradient accumulation, model parallelism, and efficient memory management.
- Example: Using gradient accumulation to enable fine-tuning of a large model with a batch size that exceeds the available GPU memory

Hardware and Computational Requirements

- Distributed Training
 - Scaling the fine-tuning process across multiple GPUs or machines to handle larger datasets and models.
 - Frameworks: PyTorch Distributed, TensorFlow Distributed, and Horovod.
 - Example: Setting up a distributed training environment with multiple GPUs to fine-tune a large language model on a massive text corpus.





Continual Learning and Lifelong Learning

- The ability of a model to continuously learn and adapt to new tasks and domains without forgetting previously learned knowledge.
- Challenges
 - Catastrophic forgetting: The model's performance on previous tasks degrades when fine-tuned on new tasks.
 - Efficient knowledge transfer: Leveraging prior knowledge to accelerate learning on new tasks.
- Future Directions
 - Developing architectures and techniques that enable seamless continual learning.
 - Investigating regularization methods and memory mechanisms to mitigate catastrophic forgetting.
 - Example: A fine-tuned language model that can continuously adapt to new
- omains (e.g., medical, legal) while retaining its performance on previous

Few-Shot and Zero-Shot Learning

- Few-Shot Learning
 - The ability of a model to learn from a small number of labeled examples per class.
 - Techniques: Meta-learning, metric learning, and attention mechanisms.
 - Example: Fine-tuning a model on a small dataset of rare disease images to classify them accurately.
- Zero-Shot Learning
 - The ability of a model to perform a task without any task-specific training examples.
 - Techniques: Leveraging pre-trained knowledge, using textual descriptions or attributes.
 - Example: A fine-tuned language model that can answer questions about a novel without any prior exposure to its content.

Model Interpretability and Explainability

- The ability to understand and explain the decisions and predictions made by a fine-tuned model.
- Challenges
 - Black-box nature of deep learning models.
 - Complexity of model architectures and large parameter spaces.
- Future Directions
 - Developing techniques for visualizing and interpreting model activations and attention.
 - Generating human-readable explanations for model predictions.
 - Example: Using saliency maps and attention visualizations to interpret the decision-making process of a fine-tuned sentiment analysis model.

Ethical Considerations in Fine-Tuning

- Bias and Fairness
 - Fine-tuned models can inherit biases present in the training data.
 - Ensuring fairness and non-discrimination in model outputs.
 - Example: Detecting and mitigating gender bias in a fine-tuned resume screening model.
- Privacy and Security
 - Protecting sensitive information in the fine-tuning data.
 - Preventing the leakage of private data through model outputs.
 - Example: Implementing differential privacy techniques during fine-tuning to preserve individual privacy.

Ethical Considerations in Fine-Tuning

- Responsible Deployment
 - Considering the potential misuse and unintended consequences of finetuned models.
 - Establishing guidelines and best practices for the ethical deployment of generative AI systems.
 - Example: Conducting thorough testing and risk assessments before deploying a fine-tuned language model in a customer-facing application.

Summary

- Fine-tuning has revolutionized generative AI by enabling the adaptation of pre-trained models to specific tasks and domains.
- Key techniques include supervised fine-tuning, unsupervised finetuning, and transfer learning.
- Transformer-based architectures, such as BERT and GPT, have become the backbone of fine-tuning in various applications.
- Effective fine-tuning requires careful consideration of dataset characteristics, evaluation metrics, and practical aspects like hyperparameter tuning and regularization.
- Future directions in fine-tuning include continual learning, few-shot and zero-shot learning, model interpretability, and addressing ethical challenges.
- The potential of fine-tuning in pushing the boundaries of generative Al is

