# Building and Fine-Tuning LLMs
## Module 2.2: Fine-Tuning LLMs

**Objectives**

- Differentiate between Fine-Tuning and Pre-Training, including their purposes and outcomes.
- Understand the Fine-Tuning process and describe its techniques.
- Fine-tune an LLM on a custom dataset.
- Explore best practices for optimizing model performance.
- Describe and explore advanced LLM architectures such as RAG, Adapter layers, and LoRA.
- Understand how RAG can be used to improve performance and address challenges in dynamic or domain-specific scenarios.
- Conduct a comparative analysis of Fine-Tuning vs. RAG for specific tasks, evaluating their trade-offs and suitability.

## Introduction

As mentioned before, the lifecycle of LLMs consists of three key stages: **pre-training, Fine-Tuning, and inference**. Each stage plays a vital role in building, adapting, and utilizing the model for various tasks. Pre-training is the foundational stage in training LLMs, equipping them with a broad understanding of language. During pre-training, the model learns from massive amounts of typically unlabeled text in a self-supervised manner. Two main approaches are common in this stage: autoregressive language modeling and masked language modeling. Autoregressive modeling focuses on predicting the next token in a sequence, making it ideal for decoder-only models. In contrast, masked language modeling involves masking certain tokens within a sentence and training the model to predict them based on surrounding context, often referred to as denoising.

While pre-training provides the foundational linguistic understanding, Fine-Tuning specializes the model, enabling it to excel in domain-specific tasks. Together, these stages equip LLMs with the adaptability and precision required for real-world applications.

Given the vast data and computational resources required, training an LLM from scratch is often infeasible for most organizations. Fine-Tuning has therefore emerged as the preferred approach for adapting pre-trained LLMs to specific applications. Fine-Tuning allows companies to build on a model's general language skills and specialize it for unique tasks by training on a smaller, relevant dataset, whether for legal text, customer service interactions, or product-specific information. This approach balances efficiency and customization, enabling high performance on domain-specific tasks without the need for extensive data or infrastructure. Fine-Tuning makes it feasible to bring specialized AI solutions to market efficiently.

# 1. How Datasets are used in Fine-Tuning and Evaluating LLMs

High-quality datasets play a pivotal role in both Fine-Tuning pre-trained LLMs and evaluating their performance. These datasets bridge the gap between a model's general language understanding and its ability to perform specialized tasks. However, the use of datasets varies significantly between Fine-Tuning and evaluation, necessitating distinct preparation and application processes.

## 1.1.  Dataset Preparation for Fine-Tuning

Fine-Tuning involves adapting a pre-trained model to a specific domain or task by training it on a carefully curated dataset. Preparing such datasets requires tasks similar to what we have seen in the previous module:

- **Cleaning and Preprocessing:** Removing noise, inconsistencies, and biases that could negatively impact model predictions.

- **Balancing and Diversity:** Ensuring a representative sample to avoid overfitting or under-representation of critical features.

- **Domain-Specific Curation:** Tailoring datasets to the model's intended application, such as legal analysis, medical diagnosis, or customer service.

For example, Google AI used instruction Fine-Tuning to develop **MedPaLM** from its base model **PaLM**. In instruction Fine-Tuning, the model is trained on a set of examples that are prefixed with instructions. The instructions tell the model what kind of response is expected. For example, an instruction might be "Answer this question in a complete sentence." The dataset used for Fine-Tuning is called **MultiMedQA** and

consists of over 100,000 questions and answers from various sources, including the US Medical Licensing Examination (USMLE), PubMed, and clinical trials.

## 1.2. Dataset for Evaluation

Once Fine-Tuning is complete, datasets can be used to evaluate the model's performance and reliability. Unlike Fine-Tuning datasets, evaluation datasets serve to test the model's ability to generalize and handle specific challenges.

**Evaluation Techniques**

1. **Instruction-Based Testing**: Evaluating the model's response accuracy and alignment with given instructions, often using datasets such as HaluEval (dataset that contains 35,000 question-answer pairs for general and task-specific scenarios) to measure hallucination tendencies (e.g., generating incorrect or nonsensical outputs). Hallucination refers to the phenomenon where a language model generates text that is incorrect, nonsensical, or not based on the input given. Meanwhile, BOLD was developed by Amazon to measure fairness across different domains.

2. **Benchmarking with Ground Truth**: involves evaluating a model's performance against a dataset that contains "ground truth" information—data with known, correct outputs. This allows developers to verify if the model can produce accurate and reliable results on specific tasks. For instance, datasets like IBEM (used for benchmarking mathematical reasoning) help test whether the model can correctly identify mathematical expressions or solve domain-specific problems. By comparing the model's predictions with the expected outputs in these datasets, developers can assess its accuracy, consistency, and suitability for specialized applications.

# 2. Primary Fine-Tuning Approaches

When you fine-tune a Large Language Model, you adjust its parameters based on the task you aim to accomplish. The extent of these changes depends on the specific job requirements. Generally, there are two main approaches to Fine-Tuning LLMs: feature extraction and full Fine-Tuning. Let's explore each method in detail:

## 2.1. Feature Extraction (repurposing)

Feature extraction, also known as repurposing, is a widely used approach to enhance and fine-tune LLMs for specific tasks. This method treats a pre-trained LLM as a fixed feature extractor, leveraging the rich language representations it has already learned from training on vast datasets. By doing so, the majority of the model remains unchanged, while only the final layers are trained on task-specific data.

The process involves freezing the pre-trained layers of the model and training the final few layers on labeled examples from the target domain or task. This allows the model to repurpose its pre-existing knowledge and adapt to the new requirements efficiently. By avoiding the need to retrain the entire model, feature extraction significantly reduces computational costs and time, making it an effective and cost-efficient method for Fine-Tuning LLMs.

This approach is particularly beneficial for tasks like sentiment analysis, classification, or named entity recognition, where the underlying language patterns remain consistent across domains, but task-specific nuances require Fine-Tuning. For instance, a pre-trained BERT model can be used as a fixed feature extractor, and only the output layers can be retrained for tasks such as legal document classification or customer feedback sentiment analysis.
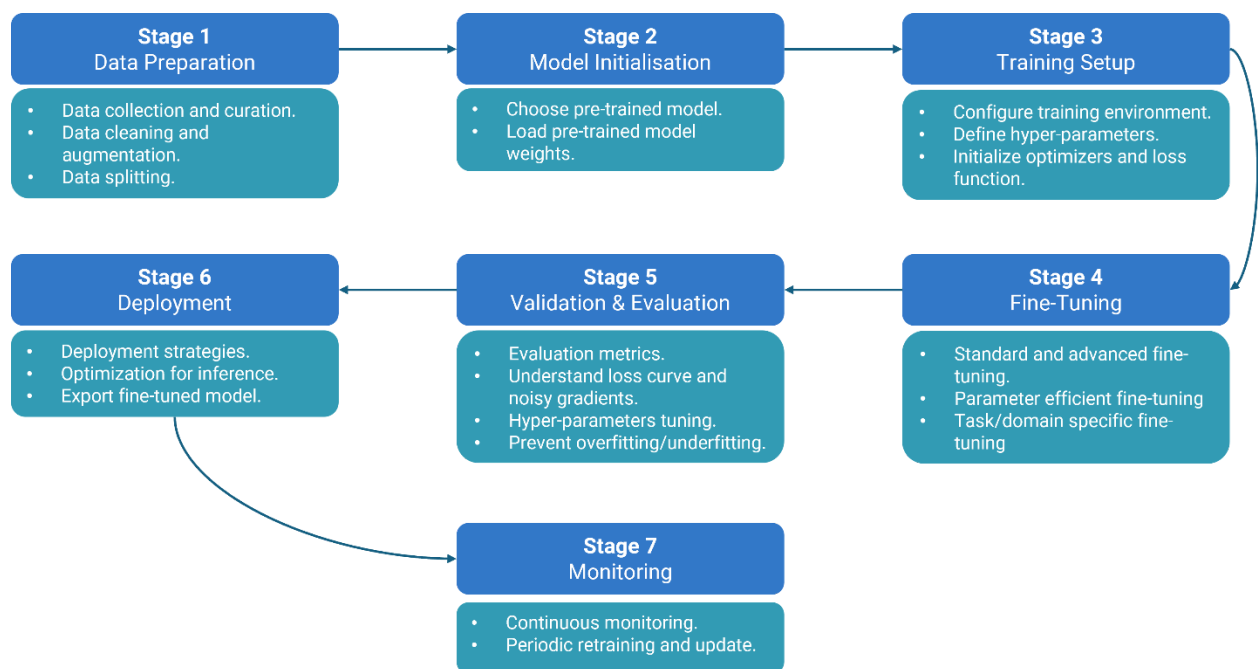
## 2.2. Full Fine-Tuning

Full Fine-Tuning is another primary approach to Fine-Tuning LLMs for specific purposes. Unlike feature extraction, where only the final layers are trained, full Fine-Tuning involves adjusting all layers of the model using task-specific data. This means that the entire model's parameters are updated during training, allowing for a deeper and more tailored adaptation to the new task.

This method is especially effective when the task-specific dataset is large and significantly different from the dataset used during pre-training. By enabling the model to learn from all layers, full Fine-Tuning ensures that it captures complex patterns and nuances in the new data, leading to superior performance on the target task. For instance, a model fine-tuned on legal text with a large corpus of case law can develop domain-specific expertise that general pre-trained models lack.

However, full Fine-Tuning requires more computational resources, time, and infrastructure compared to feature extraction. Its demanding nature makes it suitable for scenarios where high performance is critical, and the investment in resources is justified by the task's complexity or importance. This approach provides unparalleled adaptability but should be considered carefully based on dataset size, task requirements, and available resources.

## 3. Fine-Tuning Process and Best Practices

The Fine-Tuning process can be systematically broken into seven key stages, as shown in *Figure 1.* Each stage is critical to ensuring the pre-trained model is adapted to specific tasks with optimal performance. By following this structured pipeline, developers can refine the model to meet precise requirements, enhancing its ability to produce accurate, contextually appropriate responses across a variety of applications.

**Stage 1**
Data Preparation
- Data collection and curation.
- Data cleaning and augmentation.
- Data splitting.

**Stage 2**
Model Initialisation
- Choose pre-trained model.
- Load pre-trained model weights.

**Stage 3**
Training Setup
- Configure training environment.
- Define hyper-parameters.
- Initialize optimizers and loss function.

**Stage 6**
Deployment
- Deployment strategies.
- Optimization for inference.
- Export fine-tuned model.

**Stage 5**
Validation & Evaluation
- Evaluation metrics.
- Understand loss curve and noisy gradients.
- Hyper-parameters tuning.
- Prevent overfitting/underfitting.

**Stage 4**
Fine-Tuning
- Standard and advanced fine-tuning.
- Parameter efficient fine-tuning
- Task/domain specific fine-tuning

**Stage 7**
Monitoring
- Continuous monitoring.
- Periodic retraining and update.

*Figure 1:* A detailed pipeline for Fine-Tuning LLMs, showcasing the seven key stages.

### 3.1. Stage 1: Dataset Selection and Preparation

The dataset serves as the foundation for Fine-Tuning, acting as the bridge between the model's general language capabilities and its ability to address specific user needs. A well-selected and curated dataset ensures that the fine-tuned model is both accurate

and contextually relevant, enabling it to excel in specialized tasks such as medical diagnosis, legal text analysis, or customer service.

The quality, relevance, and diversity of the dataset directly influence the model's performance. A carefully curated dataset provides the contextual knowledge necessary for domain-specific tasks, while poorly prepared datasets can degrade predictions due to issues like noise, inconsistencies, and biases. Therefore, cleaning and preprocessing the dataset—such as removing irrelevant information, correcting errors, and ensuring balanced representations—are essential steps in mitigating these challenges and improving model accuracy.

The key activities in this stage include:

- **Data Collection and Curation:** Gathering high-quality, domain-specific datasets relevant to the target task, ensuring that the data aligns with the model's intended applications.
- **Data Cleaning and Augmentation:** Removing noise, inconsistencies, and duplicates to maintain data quality, and applying augmentation techniques to enhance diversity within the dataset.
- **Data Splitting:** Dividing the dataset into training, validation, and test subsets to facilitate effective model training and evaluation.

This stage ensures that the model is exposed to data that accurately represents the target domain, which is crucial for achieving high performance.

## 3.2. Stage 2: Model Initialization

In this stage, the appropriate pre-trained model is selected based on the task requirements. Key steps include:

- **Choosing the Pre-Trained Model:** Identifying a model that aligns with the complexity and scope of the task, such as GPT, BERT, or other domain-specific models.
- **Loading Pre-Trained Weights:** Importing the learned weights from the pre-training phase to provide the model with a strong foundation.

This step leverages the general knowledge already encoded in the model, reducing the time and computational resources needed for training.

### 3.3. Stage 3: Training Environment Setup

This stage involves configuring the computational infrastructure and training parameters to support the Fine-Tuning process. Key activities include:

- **Defining Hyperparameters:** Setting values such as learning rate, batch size, and number of epochs.
- **Initializing Optimizers and Loss Function:** Selecting optimization techniques and defining a loss function tailored to the task.
- **Configuring the Training Environment:** Setting up GPUs, TPUs, or other hardware resources to accelerate the training process.

A well-configured environment ensures efficient training and minimizes potential bottlenecks.

### 3.4. Stage 4: Partial or Full Fine-Tuning

The core of the process, this stage involves training the model on task-specific data. Fine-Tuning can be carried out using different approaches:

- **Standard Fine-Tuning:** Training all layers of the model for full adaptation to the task.
- **Parameter-Efficient Fine-Tuning:** Adjusting only certain components, such as low-rank matrices, to reduce computational costs.
- **Task/Domain-Specific Fine-Tuning:** Adapting the model for either broad domains (e.g., legal or medical text) or narrow tasks (e.g., question answering).

Fine-Tuning enables the model to specialize and improve its performance for specific tasks, whether broad or highly focused.

### 3.5. Stage 5: Evaluation & Validation

This stage focuses on assessing the model's performance to ensure it meets the desired standards. Key activities include:

- **Evaluation Metrics:** Using metrics such as accuracy, precision, recall, and F1 score to evaluate performance.
- **Understanding Loss Curves and Noisy Gradients:** Monitoring these during training to prevent overfitting or underfitting.

- **Hyperparameter Tuning:** Refining the training setup to further optimize performance.

By thoroughly evaluating the model, developers can identify weaknesses and make necessary adjustments before deployment.

### 3.6.  Stage 6: Deployment

Deployment involves integrating the fine-tuned model into real-world applications. Key considerations include:

- **Deployment Strategies:** Ensuring the model is accessible via APIs or embedded in larger systems.
- **Optimization for Inference:** Adjusting the model to improve response time and reduce computational overhead during inference.
- **Exporting the Fine-Tuned Model:** Saving the model in a format suitable for deployment.

A well-deployed model bridges the gap between development and practical usage, delivering value to end users.

### 3.7.  Stage 7: Monitoring and Maintenance

The final stage ensures the model remains effective and reliable post-deployment. Key activities include:

- **Continuous Monitoring:** Tracking the model's performance to detect any drifts or inconsistencies.
- **Periodic Retraining and Updates:** Incorporating new data to refine the model and keep it aligned with evolving requirements.

Ongoing maintenance ensures the model continues to perform optimally over time, adapting to changes in data or user needs.

## 4. Prominent Fine-Tuning Methods

Fine-Tuning is not a one-size-fits-all technique; different methods address distinct requirements, from task-specific adaptation to improving alignment with human preferences. There are several Fine-Tuning methods and techniques used to adjust the

model parameters to a given requirement. Broadly, we can classify these methods into two categories: **Supervised Fine-Tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF)**.

## 4.1. Supervised Fine-Tuning (SFT)

Supervised Fine-Tuning is one of the most effective and widely used methods to adapt a pre-trained language model for specific tasks using labeled datasets. In this approach, the model is trained on a task-specific labeled dataset, where each input is paired with a corresponding output or label. This process enables the model to adjust its parameters to minimize the error between its predictions and the correct labels, effectively guiding it to apply its pre-existing knowledge, gained during pre-training, to the specialized task at hand.

By learning from task-specific labeled examples, the model can identify patterns and features that are directly relevant to the problem. For example, in sentiment analysis, the model is fine-tuned on a dataset where text samples are labeled with sentiments (e.g., positive, negative, or neutral), enabling it to predict sentiments accurately in real-world applications.

> **Example**: A well-known application of supervised Fine-Tuning is seen with **BERT**. Initially pre-trained on large corpora using masked language modeling, BERT was later fine-tuned for various NLP tasks such as question answering and named entity recognition (NER). For instance, BERT was fine-tuned on the Stanford Question Answering Dataset (SQuAD), allowing it to answer questions based on textual passages by learning the specific nuances of this task.

SFT can significantly enhance a model's performance by aligning it to a domain or task. This process adjusts the model's parameters to meet specific requirements, making it an essential method for customizing large language models. The most common supervised Fine-Tuning techniques are:

### 4.1.1. Basic Hyperparameter Tuning

Basic hyperparameter tuning is a foundational approach to optimizing a ML model's training process. It involves manually adjusting key hyperparameters—such as the learning rate, batch size, and number of epochs—to achieve the best balance between

performance and efficiency. The primary objective is to identify the hyperparameter configuration that enables the model to learn effectively from the data while minimizing risks such as overfitting or underfitting.

> **Example**: For instance, when Fine-Tuning BERT for sentiment analysis, a developer might start with a learning rate of 5e-5, a batch size of 32, and train the model for 10 epochs. After evaluating the results on a validation set, they might adjust the learning rate to 3e-5 and increase the number of epochs to 20, resulting in better model accuracy on the test set. Such iterative adjustments enable the model to generalize more effectively to unseen data.

Basic hyperparameter tuning, while simple, is a crucial step in training or Fine-Tuning models to achieve optimal performance. It provides a foundation for exploring more advanced optimization techniques, such as grid search or Bayesian optimization, when required for more complex scenarios.

### 4.1.2. Transfer learning

Transfer learning allows a pre-trained model to apply its general knowledge to a new, related task, especially when there is limited data for full Fine-Tuning. Typically, engineers freeze the original layers to retain the model's foundational understanding and add new trainable layers to adapt to the specifics of the new task. While all Fine-Tuning techniques fall under transfer learning, this approach specifically enables a model to tackle tasks different from its initial training, leveraging insights gained from a broad dataset to perform well in a specialized context.

> **Example:**
>
> For example, **MedPaLM** was developed using this method, building upon Google's 540-billion-parameter PaLM model, known for its performance in complex language tasks. By training MedPaLM with annotated medical question-answer pairs and applying targeted prompting techniques, the model adapts its broad language understanding to excel within the medical field. This strategy enables the model to harness its extensive foundational knowledge while adapting effectively to specialized applications with minimal additional data.
>
> Using just 65 conversational sample pairs, Google created a medical-specific model that achieved a passing score on the **HealthSearchQA benchmark**, significantly outperforming comparable models. Rather than the usual method of training on a

digital nova scotia
StFX ONLINE — ONLINE LEARNING & PROFESSIONAL STUDIES
ACENET accelerate discovery

large, diverse domain-specific dataset, Google's approach with MedPaLM emphasized a highly targeted, minimal dataset to deliver impactful results.

While there is still room for enhancement, Google's MedPaLM and its successor, MedPaLM 2, illustrate that large language models can be effectively tailored for specialized tasks through innovative, resource-efficient methods. This success highlights the potential to adapt LLMs for targeted applications without the need for vast amounts of data and computation.

Transfer learning has revolutionized the training of LLMs by enabling their deployment in diverse applications, ranging from customer service chatbots to sentiment analysis, translation, and legal document processing. It exemplifies how foundational knowledge can be adapted to specific use cases, maximizing the utility of pre-trained models while minimizing resource demands.

### 4.1.3. Few-shot learning

Few-shot learning is a powerful approach that enables a model to adapt to a new task with minimal task-specific data. By leveraging the extensive knowledge gained during pre-training, the model can effectively learn from just a few examples provided in the prompt. This technique is particularly useful in scenarios where collecting a large, labeled dataset is impractical or expensive.

In few-shot learning, the model is given a handful of examples, or "shots," during inference time. These examples are embedded directly in the input prompt, serving as context to guide the model's predictions. Instead of requiring extensive Fine-Tuning or retraining, the model uses these examples to understand the task and generate accurate outputs. This makes few-shot learning a flexible and efficient solution for adapting pre-trained models to new tasks with minimal effort.

By enabling models to perform effectively with limited labeled data, few-shot learning is ideal for use cases where data is scarce or rapidly changing, such as medical diagnosis, niche language translation, or emerging trends in customer sentiment.

**Example:**
For instance, if **GPT-3** is tasked with translation, users can include a few example translations in the prompt to "instruct" the model on the expected output format.

GPT-3's extensive pre-training on diverse text data enables it to infer task requirements from just a few examples, performing well without Fine-Tuning for each individual task. This made GPT-3 remarkably versatile, capable of handling multiple applications like summarization, language translation, and even code generation, simply by adapting to the few examples provided.

Few-shot learning can also be integrated into the reinforcement learning from human feedback (RLHF) approach if the small amount of task-specific data includes human feedback that guides the model's learning process.

### 4.1.4. Task-specific Fine-Tuning

This method allows the model to adapt its parameters to the nuances and requirements of the targeted task, thereby enhancing its performance and relevance to that particular domain. Task-specific Fine-Tuning is particularly valuable when you want to optimize the model's performance for a single, well-defined task, ensuring that the model excels in generating task-specific content with precision and accuracy.

Task-specific Fine-Tuning is closely related to transfer learning, but transfer learning is more about leveraging the general features learned by the model, whereas task-specific Fine-Tuning is about adapting the model to the specific requirements of the new task.

**Example:**

A notable example of task-specific Fine-Tuning is FinBERT, a financial sentiment analysis model developed from Google's BERT model. FinBERT was fine-tuned on large datasets of financial texts, such as earning calls, financial news, and analyst reports, to adapt BERT's general language understanding to the financial domain.

### 4.1.5. Domain-Specific Fine-Tuning

Domain-specific Fine-Tuning customizes a model to understand and generate text relevant to a particular industry or field. By training on specialized datasets—such as medical or legal documents—the model's language and contextual comprehension are adapted to fit the target domain. For example, to create a medical chatbot, the model

could be fine-tuned with medical records to enhance its understanding of healthcare-related language and improve its responses in that field.

---

**Example:**

**BioBERT** is a prime example of a model adapted for a specific domain through Fine-Tuning. Built on BERT's architecture, BioBERT was fine-tuned on biomedical datasets, including articles from **PubMed** and **PMC**, to equip it with an understanding of medical terminology and context. This domain-specific Fine-Tuning enabled **BioBERT** to excel in tasks like biomedical named entity recognition and relation extraction.

---

## 4.2.  Reinforcement Learning from Human Feedback (RLHF)

RLHF is an innovative approach that involves training language models through interactions with human feedback. By incorporating human feedback into the learning process, RLHF facilitates the continuous enhancement of language models, so they produce more accurate and contextually appropriate responses.

This approach not only leverages the expertise of human evaluators but also enables the model to adapt and evolve based on real-world feedback, ultimately leading to more effective and refined capabilities.

The most common RLHF techniques are:

### 4.2.1. Reward modeling

In this technique, the model generates several possible outputs or actions, and human evaluators rank or rate these outputs based on their quality. The model then learns to predict these human-provided rewards and adjusts its behavior to maximize the predicted rewards.

Reward modeling provides a practical way to incorporate human judgment into the learning process, allowing the model to learn complex tasks that are difficult to define with a simple function. This method enables the model to learn and adapt based on human-provided incentives, ultimately enhancing its capabilities.

### 4.2.2. **Proximal policy optimization (PPO)**

PPO is an iterative algorithm used to refine a language model's policy, which is the strategy the model uses to decide its next action (e.g., predicting the next token) based on the current input and context. The goal of PPO is to improve this policy to maximize the expected reward, often derived from human feedback or other evaluation metrics.

A key feature of PPO is its ability to update the policy while ensuring that changes are gradual and controlled. This prevents the model from making drastic updates that could destabilize its performance. PPO achieves this by introducing a constraint through a surrogate objective function with a clipped probability ratio. This clipping mechanism limits the extent of policy updates, allowing beneficial small adjustments while avoiding harmful large changes.

By stabilizing the learning process, PPO ensures more efficient and reliable training compared to other reinforcement learning methods, making it highly effective for tasks that require aligning a language model's behavior with specific goals or preferences.

### 4.2.3. Comparative ranking

Comparative ranking is like reward modeling, but in comparative ranking, the model learns from relative rankings of multiple outputs provided by human evaluators, focusing more on the comparison between different outputs.

In this approach, the model generates multiple outputs or actions, and human evaluators rank these outputs based on their quality or appropriateness. The model then learns to adjust its behavior to produce outputs that are ranked higher by the evaluators.

By comparing and ranking multiple outputs rather than evaluating each output in isolation, comparative ranking provides more nuanced and relative feedback to the model. This method helps the model understand the task subtleties better, leading to improved results.

### 4.2.4. Preference Learning (Reinforcement Learning with Preference Feedback)

Preference learning, also known as reinforcement learning with preference feedback, focuses on training models to learn from human feedback in the form of preferences between states, actions, or trajectories. In this approach, the model generates multiple outputs, and human evaluators indicate their preference between pairs of outputs.

The model then learns to adjust its behavior to produce outputs that align with the human evaluators' preferences. This method is useful when it is difficult to quantify the output quality with a numerical reward but easier to express a preference between two outputs. Preference learning allows the model to learn complex tasks based on nuanced human judgment, making it an effective technique for Fine-Tuning the model on real-life applications.

### 4.2.5. Parameter Efficient Fine-Tuning

Parameter-efficient Fine-Tuning (PEFT) is an advanced technique designed to enhance the performance of pre-trained LLMs on specific downstream tasks while minimizing the number of trainable parameters. Instead of updating all model weights, PEFT focuses on modifying a small subset of parameters, often by introducing additional layers or selectively tuning specific components in a task-specific manner. This strategy significantly reduces computational and storage requirements while maintaining performance levels comparable to full fine-tuning.

To sum up, Fine-Tuning is a pivotal process in adapting pre-trained LLMs to specific tasks or domains, ensuring they deliver contextually accurate and specialized outputs. This process refines the model's parameters using task-relevant data, making it highly effective for real-world applications like customer support, legal analysis, or medical diagnostics. Despite its resource-intensive nature, techniques such as PEFT make the process more efficient, enabling high performance with reduced computational demands. Fine-Tuning thus transforms generalized LLMs into powerful tools tailored to specific needs.

Below is a summary of common Fine-Tuning challenges and their respective solutions:

| Challenge | Solution |
|---|---|
| Data quality/quantity issues | Curate high-quality data, augment limited datasets. |
| Overfitting and generalization | Apply regularization, early stopping, and hyperparameter tuning. |
| Computational resources | Use model distillation, efficient architectures. |

Fine-Tuning, when executed efficiently and thoughtfully, bridges the gap between general-purpose LLMs and the specific requirements of real-world applications.

## 5. Fine-Tune an LLM – Examples

**Data Processing Examples-Tokenization-Padding-Truncation-Tensor Building**

**Example 1: Fine-Tuning GPT model using a custom Dataset**

**Example 2: Fine-Tuning BERT model using a custom Dataset**

## 6. Advanced Techniques

As LLMs grow in size and complexity, Fine-Tuning all their parameters becomes increasingly impractical due to the high computational costs and significant resource requirements involved. Advanced architectures like RAG, Adapter Layers, Low-Rank Adaptation (LoRA), and QLoRA (Quantized Low-Rank Adaptation) offer more efficient, flexible, and scalable approaches to model adaptation. These techniques address the limitations of traditional Fine-Tuning by focusing on resource efficiency and scalability without compromising performance.

One of the primary motivations for switching to these advanced architectures is resource efficiency. Fine-Tuning an entire LLM requires substantial computational resources, memory, and storage, often making it infeasible for organizations with limited hardware capacity. Techniques like LoRA and QLoRA reduce the number of trainable parameters by focusing on low-rank adaptations, meaning only specific components of the model are updated, rather than the whole network.
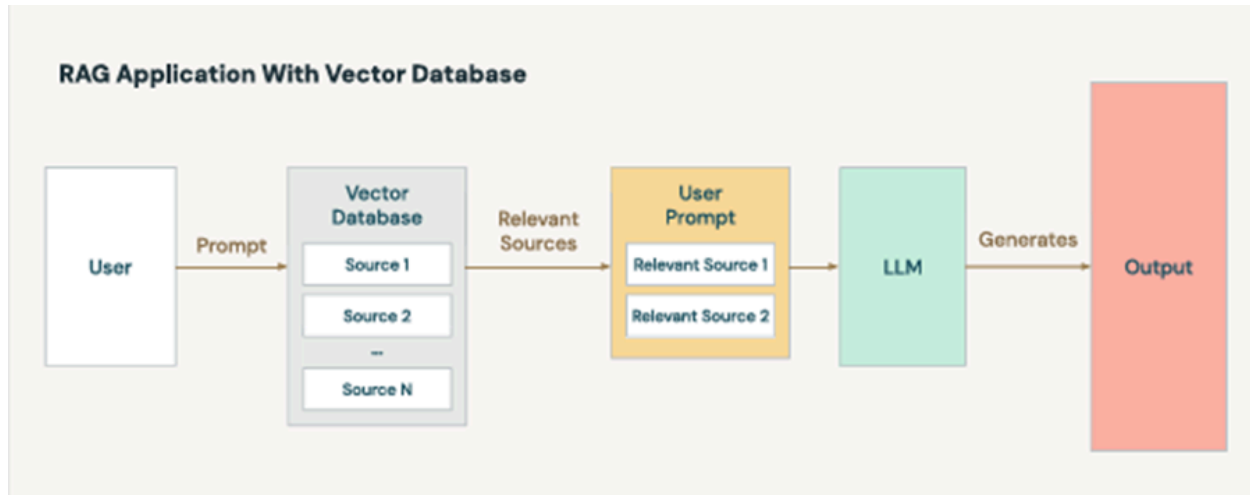
LoRA, for example, modifies the model by adding low-rank matrices that can be trained independently from the core parameters, reducing memory usage and computational demands significantly.

QLoRA takes this further by combining LoRA with quantization, optimizing memory use while maintaining model accuracy, which is especially beneficial for deploying LLMs on smaller devices or in environments with limited computational resources. For further learning about LoRA and QLoRA, please refer to this post.

Another motivation is the modularity and reusability offered by approaches like Adapter Layers. Traditional Fine-Tuning requires creating a new version of the model for each specific task, which can be inefficient and costly. Adapter layers, however, are small, task-specific modules that can be inserted into pre-trained models without retraining the entire network. This allows developers to create modular, reusable components for different tasks, switching between applications easily by adding or removing adapter layers as needed. This modularity is particularly valuable in multi-domain environments, where LLMs need to adapt rapidly across various tasks without extensive reconfiguration.

## 6.1. Retrieval-Augmented Generation (RAG)

RAG is a highly effective approach that helps mitigate certain limitations of LLMs, such as the knowledge cut-off issue. Since LLMs operate exclusively on their pre-trained data, they cannot easily be updated or retrained due to the vast computational resources required. RAG addresses this by integrating relevant external information during inference, enriching the model's responses with up-to-date context and improving accuracy. This approach allows LLMs to tailor their outputs more precisely for specific use cases. For instance, an LLM trained only on publicly available data would struggle to answer questions about a company's private documents, often leading to hallucinations. With RAG, however, the application can retrieve relevant sections from internal documents, providing the LLM with the necessary context to produce accurate and relevant answers.

**RAG Application With Vector Database**

As we've previously discussed, LLMs interact with users through "prompting." A prompt is the text provided by the user, to which the LLM responds. Prompts can vary widely depending on the model's training. Some models are designed to complete text, so prompts might be partial sentences like "Jack and Jill went up the hill to…" which the model then continues. Other models are optimized for questions or instructions, enabling prompts like, "What happened to Jack after Jack and Jill went up the hill?" In RAG applications, where users typically ask questions about specific texts, instruction-following and question-answering models are commonly employed to enhance accuracy and relevance in responses.

LLMs can usually process prompts that are several paragraphs long, which is essential for RAG. In RAG, the user's question or instruction is combined with information retrieved from an external source, creating an augmented prompt with added context.

Now that we've covered how RAG enhances prompts, let's explore the sources of external information that enrich the model's responses.
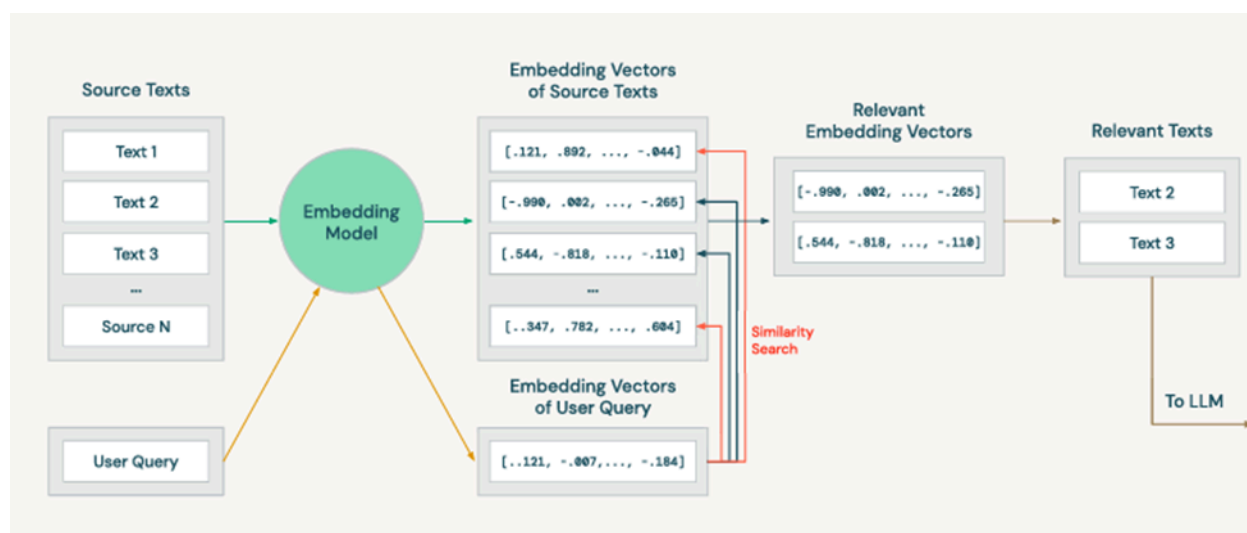
### 6.1.1. Vector Search and Embedding Models

An effective RAG application needs to retrieve information relevant to the user's prompt and feed it to the LLM. To efficiently select the most pertinent texts from millions of documents, we use a technique known as **Vector Search** to identify content closely aligned with the prompt.

In a RAG system using Vector Search, an "**embedding model**" translates each piece of text into a numerical vector—a series of numbers representing the text's meaning. This

model also converts the user's query into a comparable vector, allowing for mathematical comparisons between the query and stored texts. This approach enables the system to efficiently identify the most similar and relevant texts based on the user's input.

These vectors capture the underlying meanings of the texts they represent, allowing searches to be conducted based on meaning rather than exact wording. This approach enables the retrieval of text that is most relevant to the user's query, which can then be passed to the LLM alongside the original query for contextually accurate responses.

While we've discussed how embeddings encode the meanings of texts, it's crucial to remember that this meaning is as understood by the model. These models are trained to capture relationships and meanings relevant to human understanding, but the embedding model's interpretation may not always align perfectly with our needs. This underscores the importance of thoroughly testing and evaluating each component of a RAG application to ensure accurate and relevant results.



### 6.1.2. Vector Databases

Embedding models typically store vectors in specialized vector databases, which are designed for efficient storage and retrieval of vector data. Similar to traditional databases, vector databases support management of permissions, metadata, and data integrity, ensuring secure and organized access. They also feature mechanisms for

rapid indexing of new texts, enabling quick updates and making newly added data available for searches immediately.

---

**Example of using an LLM without RAG vs an LLM with RAG:**

**Scenario: User Asks for Trends in Sustainable Packaging from a Recent Report**

- **Without RAG**

> **Prompt:** "What are the recent trends in sustainable packaging?"
>
> **Response:** "Common sustainable packaging includes biodegradable plastics and recycled paper, though I lack the latest report specifics."

Without RAG, the response is restricted to the model's general knowledge base. Often, models do not acknowledge gaps in their knowledge and may generate responses that can be speculative or inaccurate, potentially leading to misleading or irrelevant information.

- **With RAG**

> **Prompt:** "What are the recent trends in sustainable packaging?"
>
> **Response:** "The latest report highlights trends such as compostable materials, increased recycled content, and a 20% rise in plant-based packaging."

In this case, RAG enables the LLM to pull specific details from recent reports, enhancing relevance and accuracy.

---

Without RAG, an LLM can only rely on its training data, often leading it to guess or admit uncertainty when it encounters gaps in knowledge. With RAG, however, the LLM can access external information, enabling it to provide more accurate answers.

LLMs are designed to deliver compelling and coherent responses to user prompts. They can detect nuance, interpret context, and often simulate reasoning in their

answers. Their vast training data provides them with broad general knowledge. However, LLMs are not fully reliable as knowledge sources. They frequently generate fabricated answers, known as "hallucinations," rather than admitting to knowledge gaps. Moreover, LLMs are constrained by the scope of their training data—they lack knowledge of events occurring after training and have no access to proprietary or private information, such as a company's internal documents, that were not part of the original dataset.

To mitigate these limitations, we can enhance the LLM's responses by explicitly supplying it with relevant information. For example, by copying and pasting reference documents into the prompt, along with a question, we effectively "augment" the prompt to provide necessary context. Automating this process through a retrieval system, such as a vector database, allows us to dynamically supply the LLM with pertinent information without requiring manual input from the user.

Implementing RAG with Vector Search introduces additional steps, such as data processing and vector management (often within a vector database), but it effectively addresses the limitations of using LLMs alone. RAG enhances the LLM's capabilities by providing targeted, up-to-date context that enriches its responses. Compared to retrieval-only systems, RAG offers the advantage of synthesizing information from multiple sources, enabling the LLM to deliver a cohesive and contextually appropriate answer tailored to the user's prompt.
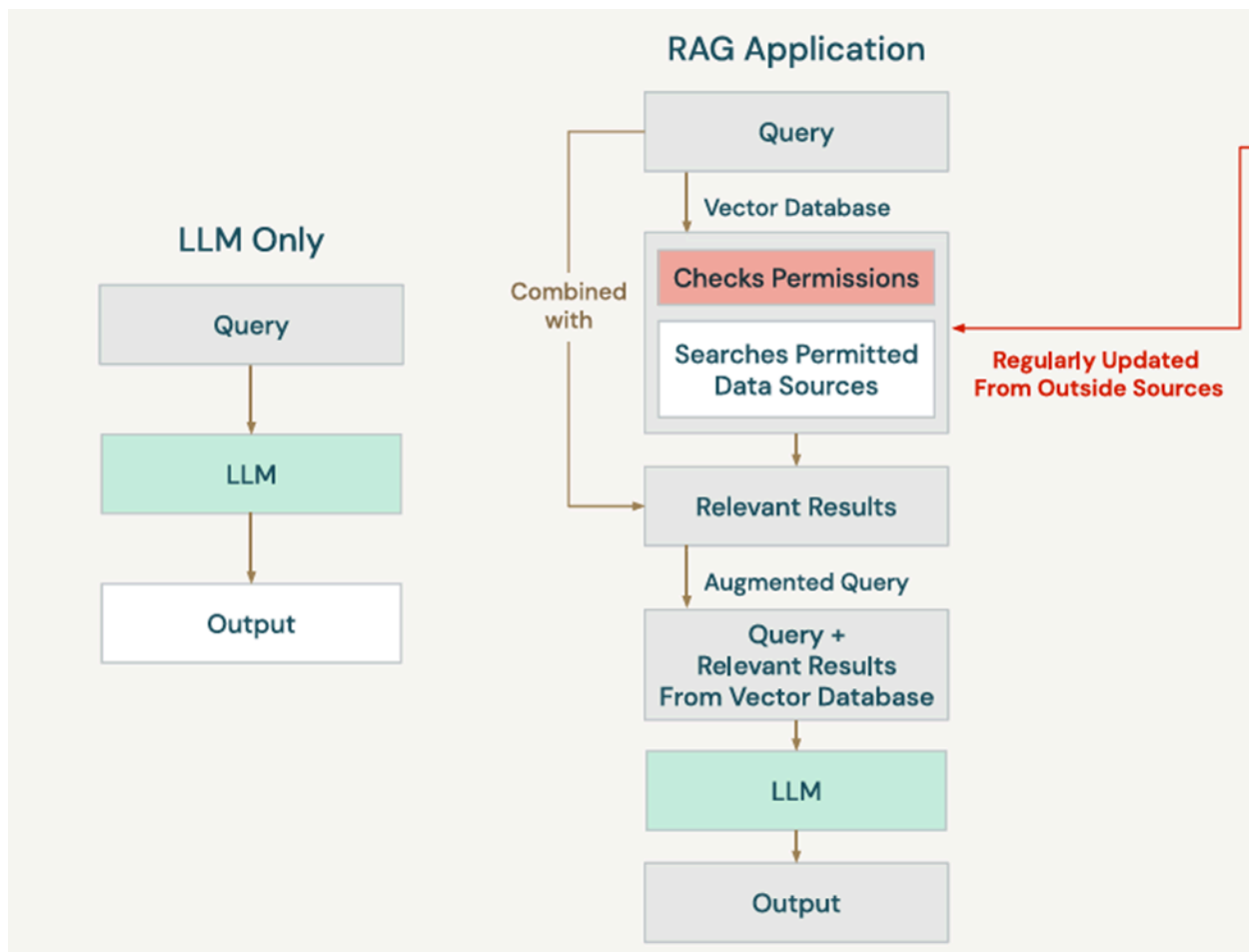
In summary, the shift from traditional Fine-Tuning to advanced architectures is motivated by the need for efficiency, flexibility, and improved performance. Techniques like RAG, Adapter Layers, LoRA, and QLoRA offer scalable, resource-efficient solutions that allow LLMs to adapt quickly to specialized applications while minimizing computational costs. This evolution in model adaptation techniques enables LLMs to serve a broader range of industry-specific tasks with greater effectiveness and efficiency.

Advanced architectures offer improved precision and reliability in model performance. QLoRA, for instance, optimizes memory use with quantization while preserving accuracy, allowing the model to maintain high-quality outputs even in resource-constrained environments. Furthermore, techniques like RAG help mitigate common issues in LLMs, such as hallucination, by grounding the model's output in

retrieved information. These targeted adaptations lead to more reliable and accurate model outputs, crucial for applications requiring high precision, such as legal or medical domains.

## 6.2. RAG Vs. LLM only approaches

Let's explore some of the benefits of RAG in more detail.



- **RAG applications can leverage proprietary data:** Since most LLMs are trained only on publicly accessible information, they lack access to a company's private documents or internal communications. RAG, however, enables the model to utilize proprietary or specialized data—such as internal reports, emails, or design documents—thereby increasing its effectiveness in specific organizational settings.

- **RAG applications provide access to up-to-date information:** LLMs are trained on datasets available up to a certain point, which means they lack knowledge of recent events or updates. As a result, they cannot address new developments, such as updated software versions or industry changes. RAG bridges this gap by pulling in the latest information, ensuring that the LLM delivers responses that reflect current knowledge.
- **RAG can enhance the accuracy of LLM outputs:** LLMs sometimes produce incorrect or made-up responses, known as hallucinations, which can limit their reliability. RAG addresses this by supplying the model with precise and relevant information from verified sources, thereby reducing errors and improving output quality. Additionally, RAG can provide citations, allowing users to cross-check information with original sources.
- **RAG supports detailed data access control:** Standard LLMs cannot adjust their responses based on user-specific permissions or access levels. In contrast, RAG systems can be set up to retrieve only data that a user is permitted to access, making it possible to securely incorporate sensitive or confidential information based on each user's credentials.

RAG provides LLMs with targeted, context-specific information that they may not reliably generate on their own, enabling applications otherwise challenging for LLMs alone:
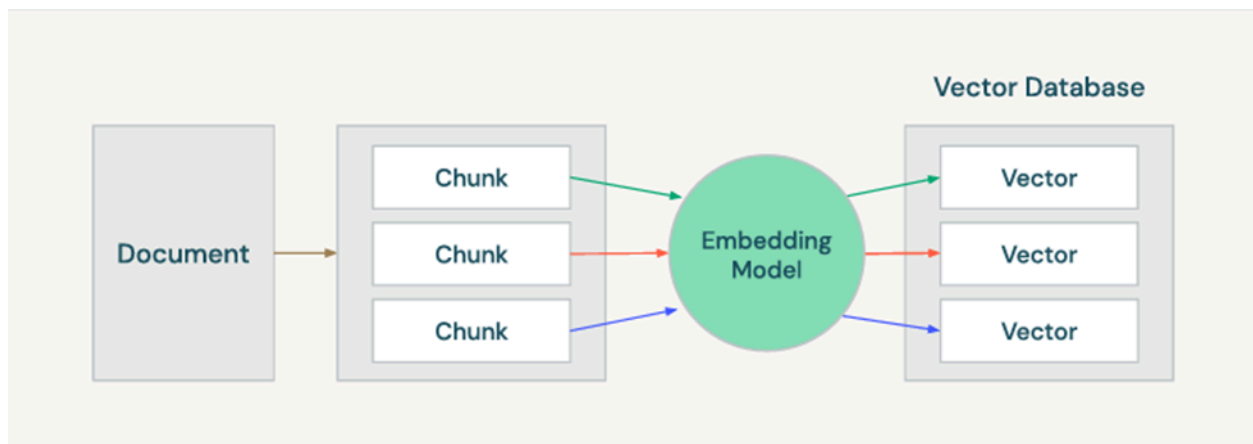
- **Question-answering:** RAG is invaluable for "document-based" queries, such as searching HR policies or real-time financial reports. For example, a large e-commerce company uses Databricks for an internal RAG system that enables their HR team to search through extensive employee policy documents.

- **Customer service:** RAG systems enhance customer service by equipping support personnel with personalized, accurate information, reducing response times and boosting customer satisfaction. Many companies employ RAG as an "internal copilot" to improve support efficiency.

- **Content generation:** In content creation, RAG can draft tailored communications, like sales emails, by incorporating the latest data and relevant context. For instance, one Databricks client uses RAG to generate

email responses for inbound sales queries, embedding product and customer information for more effective outreach.

- **Code assistance:** RAG can enhance code completion and Q&A by retrieving relevant information from codebases, documentation, and libraries, resulting in more accurate and contextually appropriate responses than LLM-only code assistants.

RAG with Vector Search retrieves information from a vector database, enhances the user's prompt with this data, and generates a response through an LLM. This section will outline these standard steps, noting that various advanced techniques can enhance performance, albeit with added complexity.

##  Data preparation: Getting an external information source into a vector database



To perform RAG with Vector Search, unstructured text data must first be added to a vector database. This is an ongoing process, as vector databases need regular updates to ensure relevance and quality. A key advantage of RAG is that we can update the vector database continuously without modifying the LLM weights. Key steps in preparing data for RAG include:
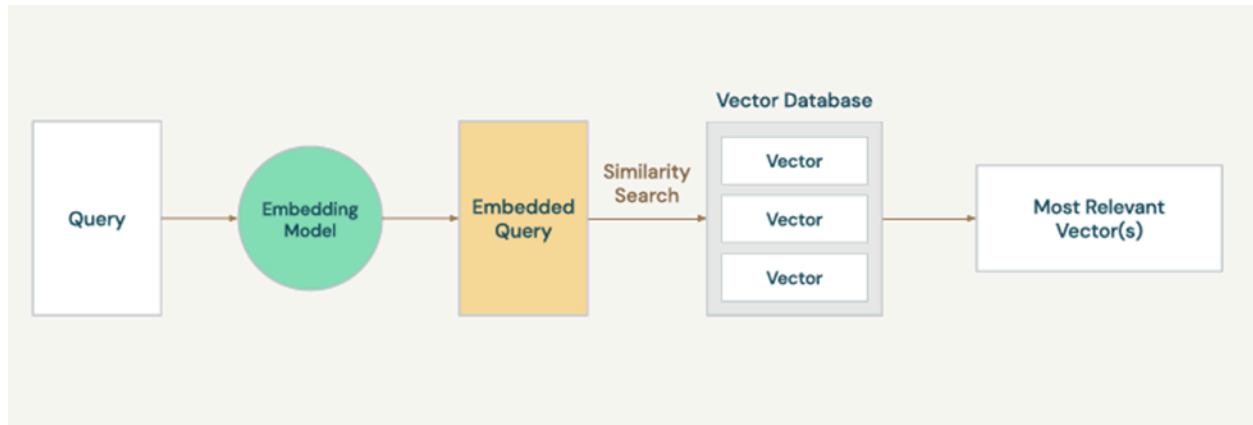
1. **Parsing the Input Documents:** Raw documents often require processing before they're suitable for RAG with Vector Search. Images may need text extraction, tables or images might need additional processing, and extraneous text like

headers or page numbers should be removed. Parsing converts these raw documents into a compatible text format for the RAG pipeline.

2. **Splitting Documents into Chunks:** In a RAG application, it's usually best to avoid retrieving entire documents. Instead, split documents into smaller, focused chunks to provide more relevant context to the LLM. The chunk size impacts output quality—too small, and they may lack sufficient context; too large, and the LLM may miss key details. Optimal chunk size varies by document type, LLM, and application goals, so testing different sizes is essential.

3. **Embedding the Text Chunks:** Once documents are split into chunks, each chunk is converted into a high-dimensional vector using an embedding model. This model transforms text into a numeric vector, or "**embedding**," capturing the meaning and context of the text. For instance, a good embedding model interprets "raining cats and dogs" as a weather-related phrase, not one about animals. Embeddings enable comparison by measuring similarity between vectors, a critical step in RAG when matching a user's prompt with relevant embedded texts in the vector database to enhance response accuracy.

4. **Storing and Indexing the Embeddings:** Embeddings are stored in a vector database, a specialized database optimized for storing and quickly searching vector data. Vector databases (or "vector stores") often allow immediate updates, enabling new chunks to be added and accessed right away. While not strictly required for RAG, vector databases significantly enhance performance and reliability. To maintain fast retrieval speeds, especially with large volumes of text chunks, embeddings are typically organized using a vector index. This index applies algorithms to efficiently map and manage embeddings, optimizing search performance within the database.

5. **Recording Metadata:** Storing metadata alongside text chunks enables filtering and referencing capabilities in a RAG application. For example, metadata can include URLs, page numbers, or dates, allowing users to filter results by source or time frame. This added layer of detail makes it easier to trace and verify the retrieved information, enhancing the relevance and credibility of the results.

## Retrieval: Getting relevant context

After preprocessing, we now have a **vector database** containing text chunks, embeddings, and metadata. This setup enables the first step in RAG: **retrieval**. Here, the user submits a prompt, typically a question, which the RAG application uses to query the database and find the most relevant information. These results are then used to enhance the original prompt in the next step.
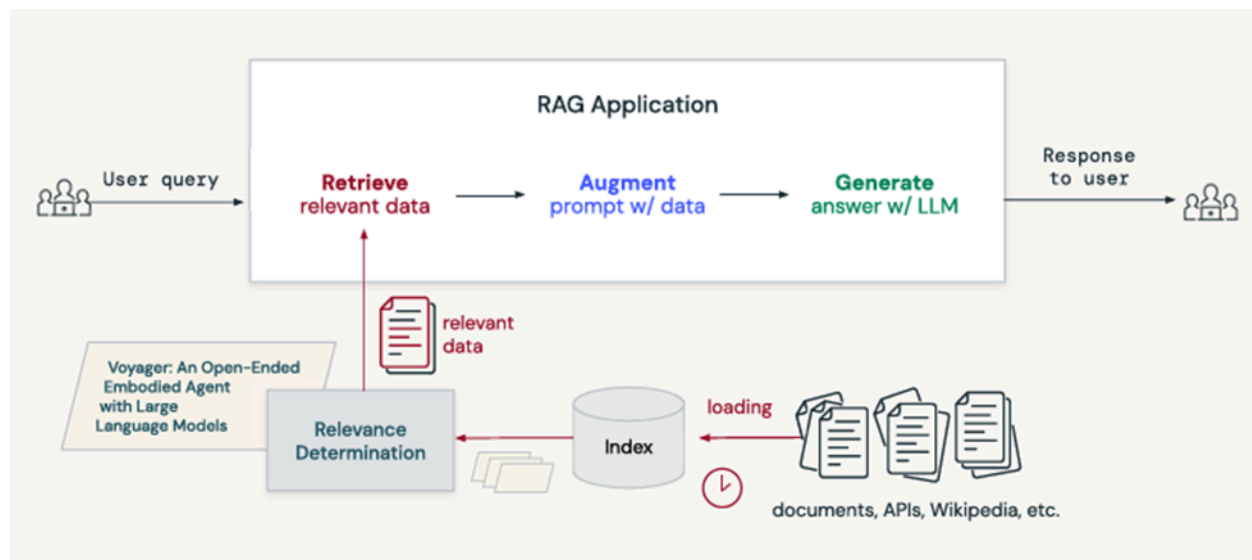
6. **Querying the Vector Database:** Since a user's input is plain text, we first need to embed it using the same model used for the original text chunks. This embedded query can then be matched with the closest vectors in the database. For small databases, we calculate a similarity score for each record; for larger databases, vector indexes and specialized search algorithms are used to speed up the process, often through approximate methods. Once the relevant records are identified, their text can be combined with the user's prompt and sent to the LLM for a response.

   **Note** that embeddings aren't converted back to text—text chunks are stored or linked to embeddings, allowing them to be retrieved as needed. Deciding on the optimal number of results to retrieve is essential, as too few may miss key details while too many may dilute relevance.

7. **Improving Retrieval:** While the basic retrieval approach is effective, several advanced techniques can further enhance accuracy:

   ● **Hybrid Search:** Combines traditional keyword search with Vector Search, improving retrieval precision.

- **Reranking:** Uses an additional model to reorder results, prioritizing the most relevant records.
- **Summarized Text Comparison:** Matches user prompts with embeddings of summarized texts, optimizing efficiency.
- **Contextual Chunk Retrieval:** Retrieves adjacent text chunks (e.g., preceding or following paragraphs) to provide richer context for the LLM.
- **Prompt Refinement:** Refines the user's original prompt using a language model to better capture intent and improve search results.
- **Domain-Specific Tuning:** Applies embedding models fine-tuned for specific domains, increasing retrieval relevance.

These techniques can be valuable if the RAG system is not consistently returning the most pertinent records from the vector database.



## 🗎 Use Cases and Examples:

- **Question and Answer Chatbots:** Integrate LLMs with chatbots to generate accurate answers from company documents, enhancing customer support.
- **Search Augmentation:** Enhance search engines with LLM-generated answers for more accurate informational queries.
- **Knowledge Engine:** Use LLMs to answer questions related to internal functions, such as HR and compliance, using company data.

# 7. Comparative analysis: Fine-Tuning vs. RAG

We have learned about each methodology for improving the LLMs' response generation. Let's examine the differences to understand them better.

## 7.1. Approach and Mechanism

**Fine-Tuning**

Fine-Tuning involves adapting a pre-trained language model to a specific task by training it further on a task-specific dataset. This process updates the model's weights to incorporate domain-specific knowledge, allowing it to produce relevant responses for the target task. Fine-Tuning is particularly effective when task requirements are well-defined and consistent, as the model learns to generate responses based on the dataset it was fine-tuned on.

**RAG**

RAG, on the other hand, combines a pre-trained language model with a retrieval component that searches for relevant information from an external knowledge base (such as a document database) during inference. Instead of relying solely on internal parameters, RAG dynamically fetches relevant documents or information to enhance its responses. This approach is especially useful when the task requires up-to-date information or extensive domain knowledge beyond the model's pre-trained knowledge base.

## 7.2. Resource and Computational Requirements

**Fine-Tuning**

Fine-Tuning typically requires significant computational resources, especially for large language models, as it involves updating the model's weights across potentially large datasets. This can be costly and time-consuming, and fine-tuned models for each specific task need to be stored separately, which may lead to storage challenges for multi-task applications.

**RAG**

RAG is generally more resource-efficient as it doesn't require modifying the model's weights or storing separate task-specific versions. Instead, it leverages a pre-trained model and retrieves relevant data as needed during inference, reducing the need for extensive retraining. However, RAG requires an efficient retrieval system and storage of an external knowledge base, which adds its own infrastructure demands, especially for high-volume or low-latency applications.

## 7.3. Adaptability

**RAG**

RAG is best for generalizations. It uses the retrieval process to pull information from different data sources. RAG does not change the model's response; it just provides extra information to guide the model.

**Fine-Tuning**

Fine-Tuning customizes the model output and improves the model performance on a special domain that is closely associated with the training dataset. It also changes the style of response generation and sometimes provides more relevant answers than RAG systems.

## 7.4. Cost

**RAG**

RAG requires top-of-the-class embedding models and LLMs for better response generation. It also needs a fast vector database. The API and operation costs can rise quite quickly.

**Fine-Tuning**

Fine-Tuning will cost you big only once during the training process, but after that, you will be paying for model inference, which is quite cheaper than RAG.

Overall, on average, Fine-Tuning costs more than RAG if everything is considered.

## 7.5. Implementation Complexity

**RAG**

RAG systems can be built by software engineers and require medium technical expertise. You are required to learn about LLM designs, vector databases, embeddings, prompt engineers, and more, which does require time but is easy to learn.

**Fine-Tuning**

Fine-Tuning the model demands high technical expertise. From preparing the dataset to setting tuning parameters to monitoring the model performance, years of experience in the field of natural language processing are needed.

| Aspect | Fine-Tuning | RAG |
|---|---|---|
| **Mechanism** | Trains model on task-specific dataset, updating weights | Combines LLM with retrieval system to access external knowledge during inference |
| **Resource Requirements** | High computational cost for large models, storage for each task-specific model | Lower cost for model storage, but requires efficient retrieval and knowledge base setup |
| **Flexibility** | Limited to specific tasks; requires separate fine-tuning per task | Flexible, adaptable across multiple tasks without retraining |
| **Knowledge-Intensive Tasks** | Effective within the task's defined scope | Excels in knowledge-intensive tasks with broad, real-time knowledge needs |
| **Mitigating Hallucinations** | Reduces hallucinations with focused dataset but | Can reduce hallucinations by grounding responses |

| | may still hallucinate outside scope | with retrieved, factual data |
|---|---|---|
| **Best Use Cases** | Domain-specific tasks with consistent requirements | Dynamic, multi-task applications needing updated or diverse information |

**Summary Table: Fine-Tuning vs. RAG**

To sum up this comparison, it's important to know that when deciding between external data access methods, RAG is often the preferred choice for applications requiring integration with external data sources. Conversely, Fine-Tuning is more appropriate when the goal is to adapt the model's behavior, writing style, or incorporate domain-specific knowledge directly into the model.

RAG systems excel in minimizing hallucinations and improving accuracy, as they rely on retrieving factual information rather than generating responses solely from pre-trained knowledge. If substantial domain-specific labeled data is available, Fine-Tuning can produce highly tailored model behavior. However, RAG is a robust alternative in scenarios where such data is limited or unavailable.

RAG also provides a significant advantage in dynamic environments where data is frequently updated, as it enables real-time retrieval of the latest information. Furthermore, RAG systems enhance transparency and interpretability by offering traceable decision-making processes, an insight often lacking in fine-tuned models.

[RAG using LLama on Kaggle by Adnane](#)

## Conclusion

Fine-tuning large language models (LLMs) has proven to be a transformative approach in adapting pre-trained models to specific tasks and domains. It bridges the gap between general-purpose language understanding and domain-specific expertise, enabling models to deliver contextually accurate and high-quality outputs. Techniques such as Parameter-Efficient Fine-Tuning (PEFT) have further optimized this process,

significantly reducing computational costs and making fine-tuning accessible to a wider range of applications.

However, fine-tuning is not the only solution for tailoring LLMs. Techniques like Retrieval-Augmented Generation (RAG) have emerged as powerful alternatives or complements to fine-tuning. RAG integrates external knowledge bases during inference, enabling models to dynamically retrieve relevant information rather than relying solely on pre-trained weights. This approach offers distinct advantages, including lower computational costs, the ability to access up-to-date data, and improved transparency through source citation. RAG minimizes the risk of overfitting and is particularly effective in rapidly evolving domains like healthcare or legal services, where real-time access to updated information is critical.

By leveraging both fine-tuning and RAG, developers can combine the precision of tailored model training with the adaptability of dynamic data retrieval, creating robust and efficient AI systems. The choice between these techniques—or their integration—depends on the specific requirements of the application, such as the availability of domain-specific data, the need for frequent updates, or computational constraints.

In conclusion, fine-tuning and RAG represent complementary approaches in the LLM ecosystem, each addressing unique challenges and offering distinct advantages. Together, they empower organizations to deploy high-performing, adaptable, and reliable language models that meet the demands of diverse real-world applications.

## Other Resources

**Advanced LLM Fine-Tuning Techniques (LoRA)**

**Reinforcement Learning from Human Feedback**

**Instruction Tuning**

**Fine-Tuning Vs RAG**

**Advanced RAG Techniques**