# Techniques for Developing and Refining Datasets

1. **Collecting Data**
   - Collecting data from various sources using Web Scraping or Data Mining techniques will help us cover a broad and diverse data.
   - We must avoid synthetic data generation and included as many real data as we can to avoid any impractical responses from the model
2. **Cleaning the Data**
   - Removing duplicates to avoid overfitting
   - Correcting typos, inconsistencies and various string formatting issues
3. **Annotating the Data**
   - Create structured, labeled datasets that AI models can learn from
   - For each query, we can identify which passages from the knowledge base are most relevant by implementing relevance scores. So that best possible chucks are selected for responses.
4. **Feedback Loop**
   - Implementing a feedback loop is an excellent way to continuously improve its performance. This process involves collecting and analyzing user interactions and model outputs to refine the system over time. The feedback received from the user can be used to update the knowledge base.

# Comparison of Language Model Fine-tuning Approaches

1. **PEFT (LoRA & QLoRA)**
   - It aims at updating only a small subset of parameter rather than the entire model. This decreses the computation cost and memory uses as compared to full fine-tuning and yet gives an optimal performace.
   - LoRA is a famous techniqe for fine tuning as it modifies the model weights by small amounts. It adds rank decomposition matrices to model weights.
   - QLoRA – It uses quantization techniques to reduce memory usage and computational resource so that we can work with even larger models in limited

hardware. Quantization basically means reducing the bit sizes to efficiently work with them.

2. **Prompt Engineering**
   - Using this technique we need to craft effective prompts to guide model behavior without changing weights. No training required, flexible for different tasks.
3. **Few Shots Learning**
   - In this technique, we train the model on a small number of examples in the prompt. This requires minimal data and can quickly adapt to new tasks.
4. **Instruction Tuning**
   - In this technique, we fine-tune model on a diverse set of tasks framed as instructions.It improves general task-following ability, enhances zero-shot performance. But it requires careful curation of instruction dataset.

# My Preference:

For most scenarios, I will prefer Parameter-Efficient Fine-tuning (PEFT) methods, particularly LoRA (Low-Rank Adaptation). I will use it for following reasons

1. Efficiency: PEFT methods are computationally efficient, allowing for fine-tuning on limited hardware.
2. Performance: LoRA often achieves performance close to full fine-tuning.
3. Flexibility: Easy to switch between different fine-tuned versions without storing entire model copies.
4. Reduced Overfitting: By updating fewer parameters, the risk of overfitting is lower, especially on smaller datasets.
5. Compatibility: Can be combined with other techniques like instruction tuning for enhanced results.

-