# Fine-tuning Datasets

Fine-tuning datasets are collections of specific examples used to improve how well AI models perform certain tasks. These datasets are used to train the model to understand and respond accurately to particular scenarios, making it more effective for specialized applications. By using these specific examples, AI models become better at their tasks, delivering more accurate and relevant results.

Simply put, a dataset is like a textbook for an LLM, packed with examples and scenarios that help the model learn and adapt.

Fine-tuning Dataset Example: Customer Service Inquiry

Imagine we're creating a dataset to fine-tune an LLM for customer service in the electronics industry. Here's a snapshot of what a single data entry (row) in our dataset might look like:

This data entry helps the LLM understand the context and respond in a manner that's not only relevant but also aligns with the desired customer service tone and approach. Including examples like this in a dataset ensures that when real customer inquiries come in, the LLM can offer helpful and accurate support, reflecting the quality of service the brand aims to provide.

| System | Customer Inquiry | Expected Model Output |
|---|---|---|
| A friendly and helpful customer support assistant, designed to address users' technical issues with empathy and efficiency. | "My laptop battery isn't holding a charge anymore." | "I'm sorry to hear that your battery isn't working as expected. Can you tell me the make and model of your laptop so I can assist you further?" |

## To Fine-tune with Datasets

To actually fine-tune a language model, you need to add your dataset to the model. The dataset is traditionally in JSONL format, where each line in the file is a separate dataset entry, as demonstrated in the previous example. With FinetuneDB, you can easily create and manage your datasets, and fine-tune directly both open-source and proprietary foundation models such as OpenAI. Alternatively, you can download your datasets from FinetuneDB and manually add it to the model.

Pre-trained language models, like those used for fine-tuning large language models, are large neural networks trained on vast corpora of text data, usually sourced from the internet. The training process involves predicting missing words or tokens in a given sentence or sequence, which imbues the model with a profound understanding of grammar, context, and semantics. By processing billions of sentences, these models can grasp the intricacies of language and effectively capture its nuances. Fine-tuning, a crucial step in optimizing these models for specific tasks, allows for further customization and specialization based on specific datasets or domains

Examples of popular pre-trained language models include BERT (Bidirectional Encoder Representations from Transformers), GPT-3 (Generative Pre-trained Transformer 3). These models are known for their ability to perform tasks such as text generation, sentiment classification, and language understanding at an impressive level of proficiency of these hyperparameters.

### Low-Rank Adaptation (LoRA):A fine tuning technique

Low-Rank Adaptation provides the modular approach towards to fine-tuning a model for domains specific tasks and provides the capability of transfer learning. LoRA technique can be implemented with fewer resources and are memory efficient.LoRA is an improved finetuning method where instead of finetuning all the weights that constitute the weight matrix of the pre-trained large language model, two smaller matrices that approximate this larger matrix are fine-tuned.

### Quantized Low-Ranking Adaptation (QLoRA)

QLoRA extends LoRA to enhance efficiency by quantizing weight values of the original network, from high-resolution data types, such as Float32, to lower-resolution data types like int4. This leads to reduced memory demands and faster calculations.

4-bit NF4 Quantization:4-bit NormalFloat4 is an optimized data type that can be used to store weights, which brings down the memory footprint considerably. 4-bit NormalFloat4 quantization is a 3-step process

**Conclusion**:LoRA is implemented in the Hugging Face Parameter Efficient Fine-Tuning (PEFT) library, offering ease of use and QLoRA can be leveraged by using bitsandbytes and PEFT together.