



Implementing Web Interface For Fine-Tuning Gemma Models Using Gradio

Author

Adarsh Dubey

dubeyadarshmain@gmail.com

Organization

Google DeepMind

Google Summer of Code '25

This is not the full proposal; personal details and several other sections were removed in an attempt to make it public and accessible to anyone. Table of Content is preserved, for anyone to see what the full proposal was like.

0. Table of contents

0. Table of contents

1. About me (summary)

2. Project overview

2.1 Project idea

2.2 Abstract

2.3 Motivation and Community Impact

3. Technical details

3.1 Implementation details

3.2 Application & it's features

Multi-page application

Model selection

Dataset uploading & processing

Model and training configuration

Fine-tuning & visualization

Exporting, merging and converting models

Vertex AI integration

3.3 Guides and blogs

3.4 Prototype

Problems in the Current Prototype

Failures while building the prototype

Learnings from the prototype

4. Deliverables

5. Timeline

6. Other relevant information

6.1 Expectations from mentors

6.2 What are my other commitments?

6.3 About me (extended)

6.4 Open source projects

6.5 Post GSoC

6.6 Use of AI

1. About me (summary)

Name: Adarsh Dubey

Timezone: Indian Standard Time (UTC +5:30)

Email: dubeyadarshmain@gmail.com

Country: India

University: Savitribai Phule Pune University

Course: B. Tech in Artificial Intelligence & Data Science

Availability: [REDACTED]

Links: <https://linkedin.com/in/dubeyadarsh>
<https://github.com/inclinedadarsh>

Resume: [REDACTED]

Project length: Large (350 hours)

Contact: [REDACTED] (WhatsApp)
dubeyadarshmain@gmail.com
<https://x.com/inclinedadarsh>
<https://linkedin.com/in/dubeyadarsh>
@inclinedadarsh (Discord)

2. Project overview

2.1 Project idea

<https://gist.github.com/dynamicwebpaige/92f7739ad69d2863ac7e2032fe52fbad#:~:text=Gemma%20Model%20Fine%2Dtuning%20UI>

For quick start, check out the prototype at
<https://github.com/inclinedadarsh/gemma-finetune-ui/>

2.2 Abstract

This project is about creating an intuitive, simple, and user-friendly web interface using Gradio to make fine-tuning Gemma models easy for everyone. The goal is to let users, even those without deep technical knowledge, customize and fine-tune these

models according to their needs. Users will be able to upload datasets in different formats and choose their preferred fine-tuning methods, like continuous fine-tuning, Alpaca-style, or conversational approaches. The platform will also allow easy conversion and downloading of models in popular formats such as PyTorch and GGUF.

For those who want more control, customizable options will be available, letting users decide between efficient LoRA fine-tuning or more comprehensive full fine-tuning methods. It'll also include seamless integration with Vertex AI, such that anyone can connect with it and fine-tune their models with ease.

Importantly, the project will also focus on providing clear and helpful documentation. The documentation will guide users through every step, clearly explaining configurations, different model types, parameter selections, and best practices for optimal results. Ultimately, this project seeks to lower the barrier to entry for AI fine-tuning, democratizing access to powerful Gemma models and making advanced AI customization accessible and practical for everyone.

3. Technical details

These details are based upon the research and exploration I have done regarding the project, as well as the current state of model fine-tuning using different approaches.

During the exploration, I implemented a basic prototype of the app using [Gradio](#) that lets users fine-tune Gemma 3 models with their own dataset and then convert and download it in GGUF format. To know more about the exploration, please check out 3.3 Prototype and Exploration. Most of the code referenced in this section is from the prototype, which can be found in the [inclinedadarsh/gemma-finetune-ui](https://github.com/inclinedadarsh/gemma-finetune-ui) GitHub repository.

AI was utilized to research and debug errors during the implementation of the prototypes for this section. Please check [6.6 Use of AI](#), to know more about it.

3.1 Implementation details

Before diving into the application itself, it's important to first decide what method will run under the hood. There are several ways to fine-tune models today, especially when it comes to Gemma—some popular options include the [Gemma LLM library](#), [Keras](#), [Hugging Face](#), or using [PyTorch](#) directly. While exploring, I tried almost all of them and found Keras and Hugging Face to be the easiest to work with. That's why I went with Hugging Face even in the prototype.

For this project, I plan to stick with Hugging Face—mainly because of its simplicity, flexibility, and strong community support. With the release of Gemma 3, DeepMind also partnered with [Unsloth](#) to implement fine-tuning notebooks in Unsloth for its efficiency. I'm open to discussing this choice with my mentor, and if needed, I can

switch to Unsloth later. Since Unsloth is built on top of Hugging Face, the transition would be smooth.

3.2 Application & it's features

This project revolves around implementing a web interface that lets users fine-tune Gemma models, making fine-tuning accessible to everyone. Below, I have outlined the application details, candidate approaches, and some other discussion-worthy topics.

Multi-page application

As of now, I think the app should consist of three separate pages: one for inference, one for fine-tuning, and one for model conversion. The main reason for this separation is flexibility—users may want to perform inference either before or after fine-tuning a model. Keeping these functionalities separate will make the user experience smoother and more intuitive.

The same logic applies to the model conversion feature. Users might want to fine-tune their models and push them directly to the hub without merging, so providing a distinct page for model conversion is practical.

In the current prototype, I've already implemented these three sections in the code, and I've also provided notebook versions for both the fine-tuning and model conversion pages. You can use [these notebooks](#) to run the Gradio application directly on Google Colab.

Here's a brief overview of each page:

1. **Fine-tuning page:** This page allows users to fine-tune models and push them to the hub if needed.
2. **Inference page:** Here, users can perform inference using fine-tuned models, whether they're stored locally or available on the hub.
3. **Model conversion page:** This page enables users to merge and convert models that are either stored locally or hosted on the hub.

You'll find more detailed information on each of these pages in the sections below.

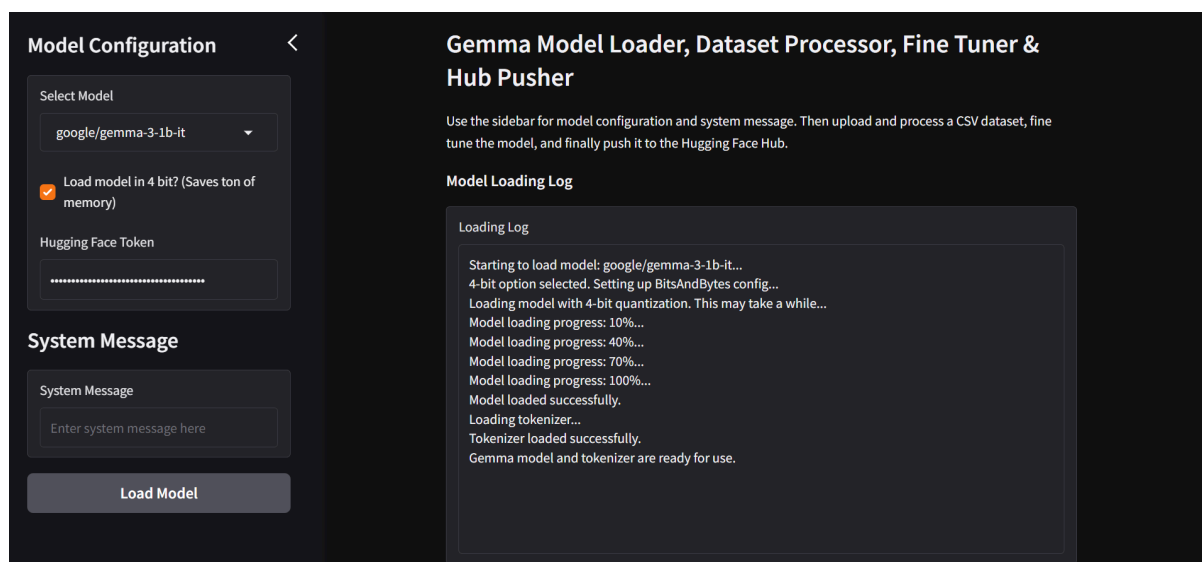
Model selection

Google DeepMind has recently released the Gemma 3 models, bringing the Gemma family up to approximately 50 fine-tunable models. Supporting all these models would be challenging, so this project will focus on supporting general Gemma models and exclude domain-specific models like CodeGemma for now. Additionally, to keep the project manageable, the fine-tuning functionality will initially be limited to text-only models.

Model selection directly impacts dataset preprocessing, particularly concerning special tokens. After some research, I discovered that Gemma 2 and Gemma 3 models use the same data preprocessing techniques, whereas Gemma 1 models require a different preprocessing approach. The app will support all three Gemma versions, adjusting the preprocessing steps according to the selected model.

Other than this, it's essential to consider the model size and tensor data types. Errors often occur during fine-tuning due to mismatched data types. To address this, the app will allow users to select quantized models. Quantization significantly reduces GPU and RAM requirements, improving the user experience. The [BitsAndBytes](#) library offers an excellent API for loading quantized models, which I might use.

I've implemented an initial version of this feature in the prototype, currently supporting only Gemma 3 models with a basic option to load models in 4-bit quantization. Below is a screenshot of the current implementation:



Here's the code for loading the model:

```
def load_model_fn(model_name, load_4bit, hf_token):
    model = None
    tokenizer = None
    if not hf_token.strip():
        return "Error: Please enter Hugging Face token", None,
    None

    try:
        if load_4bit:
            bnb_config = BitsAndBytesConfig(
                load_in_4bit=True,
                bnb_4bit_use_double_quant=True,
```

```

        bnb_4bit_quant_type='nf4',
        bnb_4bit_compute_dtype=torch.float16,
        bnb_4bit_quant_storage=torch.float16
    )
    model = AutoModelForCausalLM.from_pretrained(
        model_name,
        attn_implementation='eager',
        quantization_config=bnb_config,
        torch_dtype=torch.float16,
        device_map='auto',
        token=hf_token
    )
else:
    model = AutoModelForCausalLM.from_pretrained(
        model_name,
        attn_implementation='eager',
        torch_dtype=torch.float16,
        device_map='auto',
        token=hf_token
    )

    for percent in range(10, 101, 30):
        time.sleep(0.5)

    tokenizer = AutoTokenizer.from_pretrained(model_name,
token=hf_token)

    return "Gemma model and tokenizer are ready for use",
model, tokenizer
except Exception as e:
    return f"Error loading model: {str(e)}", None, None

```

And here's the corresponding Gradio implementation:

```

with gr.Blocks() as demo:
    model_state = gr.State()
    tokenizer_state = gr.State()

    with gr.Sidebar():
        gr.Markdown("## Model Configuration")
        model_dropdown = gr.Dropdown(choices=MODEL_OPTIONS,
label="Select Model", value=MODEL_OPTIONS[0])

```

```

        load_4bit_checkbox = gr.Checkbox(label="Load model in 4 bit?
(Saves ton of memory)", value=False)
        hf_token_input = gr.Textbox(label="Hugging Face Token",
placeholder="Enter your Hugging Face token", type="password")
        load_model_button = gr.Button("Load Model")

    with gr.Column():
        gr.Markdown("### Model Loading Log")
        model_log = gr.Textbox(label="Loading Log", lines=15)
        load_model_button.click(load_model_fn,
                                inputs=[model_dropdown,
load_4bit_checkbox, hf_token_input],
                                outputs=[model_log, model_state,
tokenizer_state])

```

This is a minimal representation of the original code. Find the [full source code here](#).

One of the major issues with the current implementation is that it doesn't have live logs or any kind of streaming feedback that shows the status of the model loading. Though the Hugging Face APIs do provide such feedback, integrating those in our app will be essential.

In summary, the model selection feature will include all general text-based Gemma models, quantization options to conserve system resources, and live streaming of logs to visualize the model-loading status for users.

Dataset uploading & processing

Based on my previous experiments, I anticipate that handling data will be one of the most challenging parts of the application. This difficulty arises primarily because datasets come in numerous formats. Since the project's primary goal is to simplify model fine-tuning, allowing users to upload datasets in their preferred formats is important.

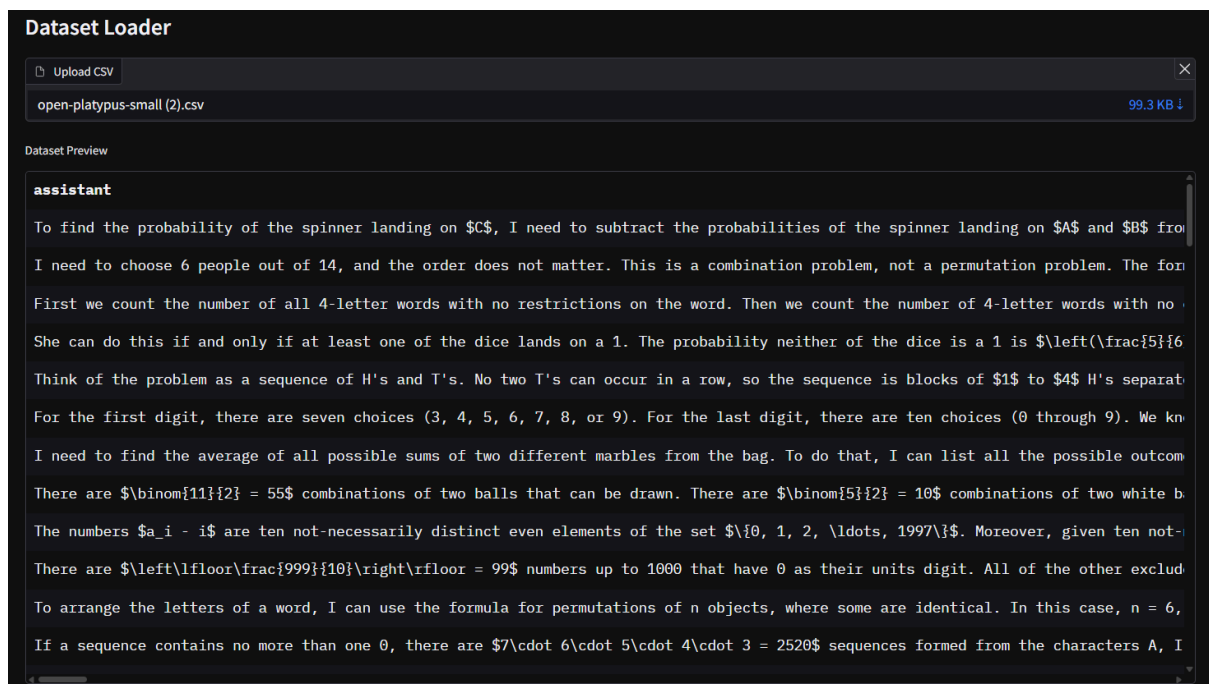
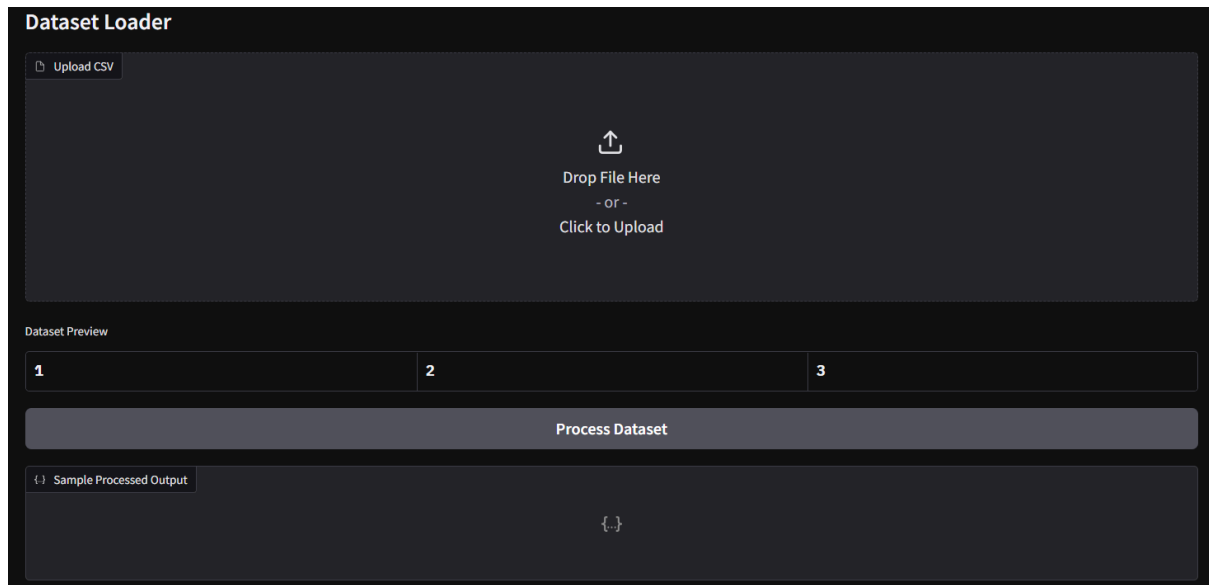
To effectively manage this, the app will support multiple dataset formats, including CSV files, PDF files, Hugging Face dataset repositories, Parquet files, Excel sheets, and others.

Additionally, dataset preprocessing should also depend upon the fine-tuning objectives. Therefore, the app will have various preprocessing styles:

- Alpaca-style formatting for instruction-following fine-tuning.
- Simple dataset dumping suitable for continuous fine-tuning.
- Conversational style for conversational model tuning.

Templating and tokenization methods also depend upon the parameters we talked about above. In the prototype, the app only accepts CSV files, and they must have the

columns 'user' and 'assistant'. The app loads the dataset, displays it to the user, and then processes it in alpaca/conversational style format. Here's the image of its implementation:



```

1  {
2    "messages": [
3      {
4        "content": "",
5        "role": "system"
6      },
7      {
8        "content":
9          "A board game spinner is divided into three parts labeled $A$, $B$ and $C$. The probability of the spinner landing
10         on $A$ is  $\frac{1}{3}$  and the probability of the spinner landing on $B$ is  $\frac{5}{12}$ . What is the
11         probability of the spinner landing on $C$? Express your answer as a common fraction."
12       },
13       {
14         "content":
15           "To find the probability of the spinner landing on $C$, I need to subtract the probabilities of the spinner landing
16           on $A$ and $B$ from $1$, since the sum of the probabilities of all possible outcomes is $1$. I can write this as an
17           equation:  $P(C) = 1 - P(A) - P(B)$ . I know that  $P(A) = \frac{1}{3}$  and  $P(B) = \frac{5}{12}$ , so I can plug those
18           values into the equation and simplify. I get:  $P(C) = 1 - \frac{1}{3} - \frac{5}{12} = \frac{12}{12} - \frac{4}{12} - \frac{5}{12} = \frac{3}{12} = \frac{1}{4}$ . I can reduce this fraction by dividing the numerator and denominator by $3$, and I
19           get:  $P(C) = \frac{1}{4}$ ."
20       }
21     ]
22   }

```

Below is the code snippet for preprocessing of the dataset:

```

def process_dataset(file, system_message):
    if file is None:
        return {"error": "No file uploaded."}, None
    try:
        file_path = file["name"] if isinstance(file, dict) and
"name" in file else file
        df = pd.read_csv(file_path)
        if 'user' not in df.columns or 'assistant' not in
df.columns:
            return {"error": "CSV must have 'user' and 'assistant'
columns."}, None

        ds = Dataset.from_pandas(df)

        def format_example(example):
            return {
                "messages": [
                    {"role": "system", "content": system_message},
                    {"role": "user", "content": example['user']},
                    {"role": "assistant", "content":
example['assistant']}
                ]
            }

        ds = ds.map(format_example,

```

```

remove_columns=ds.column_names, batched=False)
    sample = ds[0]
    return sample, ds
except Exception as e:
    return {"error": f"Error processing dataset: {str(e)}"},
None

```

Here's the corresponding Gradio interface snippet:

```

with gr.Blocks() as demo:
    dataset_state = gr.State()

    with gr.Sidebar():
        system_message_input = gr.Textbox(label="System Message",
placeholder="Enter system message here")

    with gr.Column():
        gr.Markdown("## Dataset Loader")
        csv_file = gr.File(label="Upload CSV", file_types=['.csv'])
        dataset_preview = gr.Dataframe(label="Dataset Preview")
        process_button = gr.Button("Process Dataset")
        processed_output = gr.JSON(label="Sample Processed Output")

        csv_file.change(fn=display_dataset, inputs=csv_file,
outputs=dataset_preview)

        process_button.click(process_dataset,
                            inputs=[csv_file, system_message_input],
                            outputs=[processed_output, dataset_state])

```

This is a minimal representation of the original code. Find the [full source code here](#).

Again, just like model loading, neither does this implementation include streaming live logs in the application, which should be fixed.

In short, the dataset uploading and processing feature will allow users to work with the dataset format of their choice. It will ensure that the data handling aligns seamlessly with the fine-tuning objectives. Integration of these capabilities is important to ensure that the app remains simple and easy to use.

Model and training configuration

This is one of the most important steps in the application—and also where users are most likely to run into fine-tuning issues. It's not just about setting technical configurations but also making broader decisions, like choosing the fine-tuning style or deciding between full fine-tuning and LoRA/QLoRA.

You can either go for full fine-tuning or use a Parameter-Efficient Fine-Tuning (PEFT) method like Low-Rank Adaptation (LoRA). Full fine-tuning gives the best accuracy but requires a lot of resources. On the other hand, LoRA significantly reduces GPU and memory needs, making fine-tuning more accessible on limited hardware.

Pick the fine-tuning style based on your use case:

- **Continuous pre-training:** Feed the model large chunks of text or a raw corpus to help it learn new or domain-specific knowledge.
- **Instruct fine-tuning:** Format data in the Alpaca-style (instruction + output), training the model to respond to a single instruction.
- **Conversational fine-tuning:** Use a conversation-style dataset so the model learns how to interact in a more dialog-based way.

Once you've made those choices, you'll move to the technical bits—setting training arguments like batch size, learning rate, number of epochs, gradient accumulation steps, and other key hyperparameters. You'll also decide whether or not to push your fine-tuned model to the Hugging Face Hub, and if so, to which repo.

This part can be tricky, since picking the right configuration for a specific goal requires a good grasp of what's happening behind the scenes. To make it easier, the app will provide sensible default values for training parameters. This helps users without deep technical knowledge get started comfortably, while still leaving room for customization for advanced users.

Tooltips and helpful hints will be available throughout this step, explaining why certain defaults exist and guiding users through the choices. There will also be links to detailed guides to help with deeper understanding—more on that in section [3.3 Guides and blogs](#).

Once this step is done, the app will create a `trainer` object, ready to kick off fine-tuning. One of the goals here is also to show an estimate of the time and resources (like GPU/memory usage) needed based on the chosen configuration, so users can plan accordingly. I think we should discuss how the time estimation should work—maybe something to cover during the community bonding period.

I've already implemented this configuration feature in the prototype. It lets users pick whether they want to use LoRA and configure options like `alpha`, `lora_alpha`, `epochs`, `learning_rate`, etc. There's also an option to push the model to the Hub once fine-tuning is done. Since it's just a prototype, I haven't yet added the ability to choose between fine-tuning styles (continuous pre-training, instruct, or conversational). Here's a screenshot of that part from the prototype:

Fine Tuning Configuration

<input checked="" type="checkbox"/> Use LoRA	LoRA Rank (r) 8	LoRA Alpha 8	LoRA Dropout 0.05
Epochs 3		Learning Rate 0.0002	
Max Seq Length 512		Optimizer adamw_torch_fused	
<input type="checkbox"/> Push to Hub		Repository Name Enter repository name where the model should be pushed	

Here's the code snippet from the prototype for handling these configurations:

```
def fine_tune_fn(model, tokenizer, dataset, use_lora, rank, alpha,
                 dropout, epochs, learning_rate, max_seq_length, optim, push_to_hub,
                 repo_id, hf_token):
    if use_lora:
        from peft import LoraConfig
        peft_config = LoraConfig(
            lora_alpha=alpha,
            lora_dropout=dropout,
            r=rank,
            bias="none",
            target_modules=['q_proj', 'k_proj', 'v_proj'],
            task_type="CAUSAL_LM",
        )
    else:
        peft_config = None

    from trl import SFTConfig
    training_args = SFTConfig(
        output_dir="./gemma-finetune",
        max_seq_length=max_seq_length,
        packing=True,
        num_train_epochs=epochs,
        per_device_train_batch_size=1,
        gradient_accumulation_steps=4,
        gradient_checkpointing=True,
        optim=optim,
        logging_steps=2,
        save_strategy='epoch',
```

```

        learning_rate=learning_rate,
        fp16=True,
        max_grad_norm=0.3,
        warmup_ratio=0.03,
        lr_scheduler_type='constant',
        report_to='none',
        dataset_kwargs={
            "add_special_tokens": False,
            "append_concat_token": True
        },
        push_to_hub=push_to_hub,
        hub_model_id=repo_id,
        hub_token=hf_token
    )

    # Training code goes below this

    except Exception as e:
        return f"Error during training: {str(e)}", None

```

Fine-tuning & visualization

This is one of the features I wasn't able to fully implement in the prototype. While the prototype does fine-tune the model, it doesn't display logs or any form of visualization to show users what's happening during the process.

Being able to see progress—like the current training step, loss values, or other key metrics—is super important. If a validation dataset is provided, then showing validation loss and accuracy becomes even more valuable. For non-technical users, simple graphs and visual cues can make the fine-tuning process much easier to understand instead of just staring at raw logs.

One feature I'm especially excited to implement is checkpoint support. During training, the trainer object can save model checkpoints after every few epochs. These checkpoints can even be pushed to the Hugging Face Hub if configured properly. This opens up a lot of flexibility: users can pause fine-tuning at any point and resume later from the last saved checkpoint. It also means they can restart training from either local checkpoints or ones stored on the Hub, which is super useful for longer training sessions or when working with limited compute.

Additionally, all training logs should be saved to a file. This allows users to revisit and analyze them later—whether for debugging, comparisons, or just keeping track of what worked best.

Overall, better logging, visual feedback, and checkpointing will make the fine-tuning experience much more transparent and user-friendly.

Exporting, merging and converting models

Once the model is fine-tuned, the next step is to export and prepare it for real-world usage. This includes testing, merging, and converting the model into user-friendly formats.

Before merging, users might want to test the model to ensure the fine-tuning has produced the expected results. This means basic inferencing support should be available—allowing them to input sample prompts and check the model's responses before proceeding further.

After inferencing, the merging process comes into play. Merging typically combines the base model with the fine-tuned weights (especially relevant for methods like LoRA). This creates a standalone model that doesn't require the adapter weights separately. For end users, this simplifies deployment and sharing.

Since many users will want to use their model with lightweight tools like Ollama, it's also important to support conversion into formats like GGUF. For more advanced users, other options will be provided like PyTorch format, etc.

All of these tasks—exporting, merging, and converting—will be handled on a separate page in the app. The reason for this is practical: users might have already pushed the fine-tuned model to the Hugging Face Hub or saved it locally without downloading. By keeping this functionality on its own page, users can come back to it later, select their saved model (from Hub or local), and proceed with conversion or export at their convenience.

I've already implemented a basic version of this in the prototype. Currently, it allows users to select a model and convert it to GGUF format, ready for use with Ollama. This will be further expanded to support checkpoint selection, merging, and more export formats in the final version. Here's a screenshot from the prototype:

Hugging Face Model Merger and GGUF Converter

Step 1: Merge Model

Hugging Face Model ID

inclinedadarsh/gemma-3-1b-nl-to-regex



☐ Model is already merged

Merge Model

Merge Status

Merging completed and saved to './merged_model'

Step 2: Convert to GGUF

Convert to GGUF

GGUF File (click to download)

merged_model.gguf

1.9 GB ↓

Note: Ensure that the merged model folder (./merged_model) exists before converting.

Below are code snippets for merging and converting models from the prototype:

```
def merge_model(model_id, model_is_merged):
    try:
        if model_is_merged:
            return "Merge skipped. Using existing merged model folder:
'./merged_model'"
        else:
            peft_config = PeftConfig.from_pretrained(model_id)
            base_model_id = peft_config.base_model_name_or_path

            tokenizer = AutoTokenizer.from_pretrained(base_model_id)
            base_model = AutoModelForCausalLM.from_pretrained(
                base_model_id,
                device_map='auto',
                torch_dtype=torch.float16,
                attn_implementation='eager'
            )

            model = PeftModel.from_pretrained(base_model, model_id)

            merged_model = model.merge_and_unload()

            os.makedirs("./merged_model", exist_ok=True)
            merged_model.save_pretrained('./merged_model')
            tokenizer.save_pretrained("./merged_model")
            return "Merging completed and saved to './merged_model'"
    except Exception as e:
```



```
return f"Error during merging: {e}"
```

```
def convert_to_gguf():
    try:
        if not os.path.exists("llama.cpp"):
            os.system("git clone
https://github.com/ggerganov/llama.cpp.git")

        orig_dir = os.getcwd()
        os.chdir("llama.cpp")

        merged_model_path = os.path.join("../", "merged_model")
        outfile = os.path.join("../", "merged_model",
"merged_model.gguf")

        conversion_command = f"python convert_hf_to_gguf.py
{merged_model_path} --outfile {outfile}"
        os.system(conversion_command)

        os.chdir(orig_dir)

        if os.path.exists('./merged_model/merged_model.gguf'):
            return './merged_model/merged_model.gguf'
        else:
            return "Conversion finished but the output file was not
found."
    except Exception as e:
        return f"Error during conversion: {e}"
```

Vertex AI integration

Currently, the prototype runs on Google Colab, which is great for initial testing—but it's not an ideal long-term solution for deployment. While platforms like Colab and Kaggle do offer free access to GPUs (such as the T4), they come with limitations—making them unsuitable for fine-tuning larger models or running long training sessions.

That said, integrating with Colab and Kaggle is still valuable. It allows users, especially those just starting out or with limited resources, to fine-tune smaller models without spending money. These integrations will be part of the project to ensure an accessible starting point.

However, for more scalable and flexible fine-tuning, Vertex AI integration will play a key role. By integrating Vertex AI, users won't need to worry about setting up or managing their own infrastructure. They'll be able to fine-tune models directly through

the app while controlling costs and choosing hardware configurations that suit their needs.

Although I don't have deep experience with Vertex AI yet, I plan to dive into it thoroughly during the community bonding period. The goal is to explore how it can best be integrated with this project, what APIs are available, and how users can authenticate and launch training jobs from the UI itself. This will make the entire fine-tuning pipeline more robust, scalable, and user-friendly—especially for those looking to move beyond free platforms.

3.3 Guides and blogs

Even with a simple and intuitive UI, fine-tuning can still feel overwhelming—especially for users who are new to it. To make the process even more approachable, I plan to write a series of beginner-friendly guides and blogs that explain the concepts behind each section of the app. These won't dive deep into complex math or heavy technical details but will instead focus on delivering just enough understanding to help users make informed decisions.

These guides will be comprehensive yet easy to follow, aimed at answering common questions like:

"Which model should I select?"

"What type of fine-tuning is best for my use case?"

"What's the difference between instruct fine-tuning and continuous pretraining?"

"What is LoRA, and why should I use it?"

I'll be using [GitBook](#) to write and publish these guides. I've used GitBook before and find it intuitive for writing and organizing documentation. Plus, it supports syncing with GitHub, allowing the documentation to be open-source and community-contributed. That said, I'm open to discussing with mentors if another platform is preferred.

Here are a few guide topics I plan to write during the project:

- Choosing the right base model for your task
- Selecting the appropriate fine-tuning method (full vs. LoRA/QLoRA)
- Understanding fine-tuning styles: continuous pretraining, instruct tuning, and conversational fine-tuning
- LoRA explained in simple terms, and how to tune parameters like rank, alpha, etc.
- Demystifying key hyperparameters: epochs, learning rate, batch size, and more

These guides will be integrated within the app UI through tooltips and links so users can read them as needed, without interrupting their workflow. The idea is to empower users—not overwhelm them—with just the right amount of context at every step.

3.4 Prototype

In this section, I've shared my experience building the first version of the app. It includes what's currently working, what's clearly broken, and all the things I struggled with along the way.

You'll also find what I've learned during this phase—how it helped me get a better understanding of building Gradio apps, setting up dev environments, and most importantly, how to approach the final version of this app with more clarity than I had when I started.

If you're curious:

- Here's the repo: <https://github.com/inclinedadarsh/gemma-finetune-ui>
- Here's the demo video (to see how it works): <https://youtu.be/NKV788bQEgE>
- And the usage video (if you want to try it yourself): https://youtu.be/KKHP4wLV_I8

This section basically shows how the messy first version helped shape a much clearer plan for what's next.

Problems in the Current Prototype

While the current prototype serves as a functional proof of concept, it still has several limitations:

- Users need to manually open Colab and run all the cells, which isn't ideal for accessibility or smooth user experience.
- There's no real-time feedback during fine-tuning. Users can't see logs, loss curves, or training progress, which makes the experience less transparent and engaging. I haven't yet figured out a good way to stream logs and visualizations live.
- The interface is built using default Gradio components, and while they work, the current layout and styling are very basic and not visually appealing.
- Due to time constraints, several features are missing or not fully polished. It doesn't yet reflect the complete vision for the app.

That said—it works!

The core functionality of fine-tuning a model through the interface is up and running, and that's a solid foundation to build on.

Failures while building the prototype

While building the prototype, I aimed to fine-tune models myself using notebooks to better understand the process—but it wasn't smooth sailing. I ran into a ton of issues, including

- Data type mismatches during training
- Model loaded on one GPU, while trying to fine-tune on another
- Datasets not matching the model's expected format
- Incorrect fine-tuning configurations, causing unexpected errors and crashes

I wasn't able to gather all the failed attempts in one place, but I've added some of the broken notebooks to the repository for transparency. Anyone interested can check them out here:

<https://github.com/inclinedadarsh/gemma-finetune-ui/tree/main/notebooks/failed-attempts>

Learnings from the prototype

Building the prototype was a huge learning experience for me. Struggling through errors and unexpected failures gave me a real sense of the frustration users—especially beginners—go through when trying to fine-tune models. Ironically, those frustrating moments actually motivated me more.

On the technical side, I got hands-on experience with building Gradio apps, which helped me understand how to design and structure interactive interfaces for ML workflows. I also learned how to set up and work within a typical development environment for such projects, including managing dependencies, debugging model issues, and keeping things modular and clean.

Most importantly, going through all of this gave me a much clearer vision for the project. I now have a better understanding of what the end goal should look like, what the users might struggle with, and how to bridge the gap between a raw model training workflow and a user-friendly tool that makes fine-tuning truly accessible.

4. Deliverables

Web Interface:

A fully working web app that lets users fine-tune Gemma models through an easy-to-use interface. It'll cover most of the features discussed earlier in the technical details—possibly all of them.

Guides & Blogs:

A collection of beginner-friendly guides and blogs to help users understand key concepts around fine-tuning. These will make it easier for folks to make informed decisions without needing a deep technical background.