

unsloth Tutorial

This is the first chapter of our Unsloth tutorial. We'll start by learning about `FastBaseModel`, the fundamental building block for working with large language models (LLMs) in Unsloth.

Imagine you want to use a powerful LLM, but you're worried about it taking up too much memory or being slow. `FastBaseModel` is designed to solve exactly this problem. It's a special class that provides tools to load, optimize, and use LLMs efficiently. Think of it as a toolbox filled with essential tools for working with any type of LLM within the Unsloth framework. It handles the heavy lifting so you can focus on your application.

Let's look at a simple example of using `FastBaseModel` to load a model:

```
from unsloth.models.vision import FastBaseModel
model, tokenizer = FastBaseModel.from_pretrained("unsloth/Llama-3.2-1B-Instruct")
```

This code snippet loads a pre-trained Llama model. The `from_pretrained` method takes care of downloading the model weights, loading them into memory efficiently, and preparing them for use. The model and corresponding tokenizer are returned.

What `FastBaseModel.from_pretrained` does:

1. **Downloads the model:** It automatically downloads the specified model from the Hugging Face Hub (or your local cache if available).
2. **Loads the model:** Crucially, it loads the model using optimized techniques, like 4-bit quantization (using `BitsAndBytes`) to drastically reduce memory consumption. If your hardware supports it, `bfloat16` precision is used, saving even more memory.
3. **Patches the model:** It applies internal optimizations and patches to enhance the model's speed and memory efficiency. This might involve techniques like gradient checkpointing or the use of specialized kernels (Triton). We will cover these in later chapters, but you don't need to worry about the details for now.
4. **Returns the model and tokenizer:** It returns the ready-to-use model and its associated tokenizer.

Internal Implementation (Simplified):

Here's a simplified sequence diagram showing the key steps involved when you call `FastBaseModel.from_pretrained`:

```
sequenceDiagram
    participant User
    participant Loader as FastBaseModel.from_pretrained
    participant Downloader
    participant Patcher
```

```

participant Model

User->>Loader: Calls from_pretrained("model_name")
activate Loader
Loader->>Downloader: Download model weights
Downloader-->>Loader: Weights downloaded
Loader->>Patcher: Apply optimizations
Patcher-->>Loader: Model patched
Loader->>Model: Return optimized model
deactivate Loader
User<<--Model: Optimized Model

```

The `FastBaseModel` class doesn't directly implement the model loading, downloading or patching itself. Instead it delegates the work to other helper methods and classes. Most of the heavy lifting is done in `unsloth/models/vision.py` and `unsloth/models/loader.py`. These files contain more advanced concepts such as 4-bit quantization, which we'll discuss in later chapters. For now, just understand that `FastBaseModel` manages these processes to provide a streamlined, efficient interface.

Conclusion:

In this chapter, you learned about `FastBaseModel`, the base class that simplifies working with various LLMs in Unsloth. It handles the complexities of efficient model loading and optimization, allowing you to focus on your application. In the next chapter, `FastModel`, we'll build upon this foundation and explore more specialized model classes.

Building upon our understanding of `FastBaseModel`, this chapter introduces `FastModel`, a more specialized class that handles various Large Language Model (LLM) architectures and quantization techniques. Think of `FastBaseModel` as a general-purpose toolbox, while `FastModel` is a more advanced, specialized workshop equipped to handle different types of LLMs efficiently.

Let's say you want to use a powerful, but potentially memory-intensive, LLM for a specific task. `FastModel` provides the tools to load this model, optimize it for your hardware (using techniques like 4-bit quantization), and use it efficiently. It takes care of model loading, quantization choices, and other configurations, letting you focus on your application.

Key Concepts

`FastModel` extends `FastBaseModel` by adding support for various model architectures. While `FastBaseModel` handles the core optimization processes, `FastModel` intelligently determines the best loading and optimization strategy based on the selected model. This allows Unsloth to efficiently handle a wider range of LLMs.

One key feature is its ability to choose different quantization strategies. Quantization reduces the precision of the model's weights (e.g., from 32-bit floating-point to 4-bit integers). This significantly reduces memory usage and can speed up inference, though it might slightly reduce accuracy. `FastModel` automatically handles the selection and application of these techniques based on your choices during model loading.

Using FastModel

Let's load a Llama model, optimized for speed and memory efficiency:

```
from unsloth.models.loader import FastModel
model, tokenizer = FastModel.from_pretrained("unsloth/Llama-2-7b-chat-hf")
```

This code snippet loads a pre-trained Llama-2 model using `FastModel.from_pretrained`. It automatically downloads the model, applies optimizations (like 4-bit quantization if available), and returns the ready-to-use model and tokenizer. The model is optimized to use less memory and run faster on your hardware.

Explanation:

The `from_pretrained` method in `FastModel` does the following:

1. **Determines Model Architecture:** It identifies the type of LLM (e.g., Llama, GPT).
2. **Selects Quantization:** It decides on an appropriate quantization technique (if specified, otherwise it uses a default).
3. **Loads and Optimizes:** It efficiently loads the model weights and applies the selected optimization techniques, including potentially 4-bit quantization or `bfloat16` precision (depending on your hardware support).
4. **Returns Model and Tokenizer:** Returns the ready-to-use model and tokenizer.

Internal Implementation (Simplified)

Here's a simplified sequence diagram showing how `FastModel.from_pretrained` works:

```
sequenceDiagram
    participant User
    participant FastModel as FastModel.from_pretrained
    participant ArchDetector
    participant Quantizer
    participant Loader
    participant Model

    User->>FastModel: Calls from_pretrained("model_name")
    activate FastModel
    FastModel->>ArchDetector: Determine architecture
```

```

ArchDetector-->>FastModel: Architecture identified
FastModel->>Quantizer: Select quantization
Quantizer-->>FastModel: Quantization selected
FastModel->>Loader: Load and optimize
Loader-->>FastModel: Model loaded and optimized
FastModel->>Model: Return model
deactivate FastModel
User<--Model: Optimized Model

```

The key difference from `FastBaseModel` is the addition of `ArchDetector` and `Quantizer`. These components determine the best way to handle the specified model, adding a layer of intelligence to the model loading process. The actual loading and patching is still handled largely by the code in `unsloth/models/loader.py`, particularly within the `FastModel.from_pretrained` method (as shown in the code example above).

The majority of the architecture-specific logic and quantization decisions are made within this method in `unsloth/models/loader.py`. This ensures that `FastModel` can intelligently handle a variety of LLMs with different architectures and precision requirements.

Conclusion

In this chapter, you learned how `FastModel` builds on `FastBaseModel` to provide a more specialized and versatile way to load and optimize various LLMs. It intelligently handles architecture detection, quantization selection, and efficient model loading. In the next chapter, `FastLanguageModel`, we'll delve into a further specialization for language models.

Building on our understanding of `FastModel`, this chapter introduces `FastLanguageModel`, a specialized class designed specifically for efficient handling of large language models (LLMs). Think of `FastModel` as a versatile engine for various models, and `FastLanguageModel` as a highly tuned engine specifically optimized for the unique characteristics of LLMs.

Let's say you want to use a powerful LLM for text generation or understanding, but you're concerned about speed and memory usage. `FastLanguageModel` provides the tools to load the model, apply various optimizations, and utilize it efficiently. This class handles the complexities of efficient LLM loading and optimization, allowing you to concentrate on the core application.

A Simple Use Case: Fast Text Generation

Imagine you want to generate text using a large language model like Llama 2. A standard approach might be slow and memory-intensive. `FastLanguageModel` lets us do this quickly and efficiently.

```

from unsloth.models.loader import FastLanguageModel

# Load a pre-trained Llama 2 model optimized for speed. This automatically handles
# downloading, quantization, and other optimizations.
model, tokenizer = FastLanguageModel.from_pretrained("unsloth/Llama-2-7b-chat-hf")

# Generate text. We'll provide more details on prompt engineering in later chapters,
# but for now, let's use a simple prompt.
prompt = "Once upon a time,"
inputs = tokenizer(prompt, return_tensors="pt")
outputs = model.generate(**inputs)
generated_text = tokenizer.decode(outputs[0], skip_special_tokens=True)
print(generated_text)

```

This code snippet does the following:

1. **Loads the Model:** It uses `FastLanguageModel.from_pretrained` to load a pre-trained Llama 2 model with optimizations applied.
2. **Tokenizes the Prompt:** It converts the input prompt into tokens that the model understands.
3. **Generates Text:** The model generates text based on the prompt.
4. **Decodes the Output:** It converts the model's output tokens back into human-readable text.

Key Concepts

`FastLanguageModel` extends `FastModel` by incorporating several key improvements for LLMs:

- **Specialized Optimizations:** It applies LLM-specific optimization techniques beyond those used in `FastModel`. These might include advanced techniques for handling the unique structure of language models, leading to even greater speed and memory efficiency.
- **Efficient Tokenization:** `FastLanguageModel` works closely with the tokenizer to optimize the text processing pipeline, ensuring smooth and fast conversion between text and tokens.
- **Inference Methods:** It often employs techniques such as vLLM for extremely fast inference at the cost of some flexibility.

Internal Implementation (Simplified)

Here's a simplified sequence diagram illustrating the `FastLanguageModel.from_pretrained` process:

```

sequenceDiagram
    participant User
    participant LLMLoader as FastLanguageModel.from_pretrained

```

```

participant ModelOptimizer
participant Model
participant Tokenizer

User->>LLMLoader: Calls from_pretrained("model_name")
activate LLMLoader
LLMLoader->>ModelOptimizer: Optimize LLM for speed & memory
ModelOptimizer-->>LLMLoader: Optimized Model
LLMLoader->>Tokenizer: Load Tokenizer
LLMLoader-->>User: Returns Model & Tokenizer
deactivate LLMLoader
User<<--Model: Optimized Model
User<<--Tokenizer: Optimized Tokenizer

```

The `FastLanguageModel.from_pretrained` method in `unsloth/models/loader.py` orchestrates this process. It leverages several helper functions within `unsloth/models/loader.py` and `unsloth/models/_utils.py` to handle model loading, optimization, and tokenizer integration. The actual model optimization involves various techniques described in previous chapters and later ones (like 4-bit quantization and gradient checkpointing), tailored specifically for LLMs.

Conclusion

This chapter introduced `FastLanguageModel`, a specialized class for working with LLMs efficiently. It builds upon the capabilities of `FastModel` to provide superior performance for tasks like text generation and comprehension. In the next chapter, `get_chat_template`, we'll explore how to create effective prompts for your LLMs.

Building on our exploration of efficient LLM loading with `FastLanguageModel`, this chapter introduces `get_chat_template`, a function that helps us tailor our interactions with LLMs. Imagine you're sending a letter – you wouldn't use the same format for a formal business letter as you would for a casual note to a friend. Similarly, different LLMs require different input formats for optimal performance. `get_chat_template` helps us choose the right “letter format” for our chosen model.

A Simple Use Case: Generating Text with a Specific Format

Let's say we want to use a Llama 2 model for text generation. We've already loaded the model and tokenizer using `FastLanguageModel.from_pretrained` (as shown in Chapter 3). Now, we want to make sure our prompts are formatted correctly for this specific model. `get_chat_template` will help us do exactly that.

```

from unsloth.models.loader import FastLanguageModel
from unsloth.chat_templates import get_chat_template

# Load the model (as shown in Chapter 3)
model, tokenizer = FastLanguageModel.from_pretrained("unsloth/Llama-2-7b-chat-hf")

# Configure the chat template for Llama 2 (using the 'llama' template)
tokenizer = get_chat_template(tokenizer, chat_template="llama")

# Now, our tokenizer is ready to process prompts in the Llama 2 format.
# We can proceed with text generation as before (see Chapter 3).

```

This code snippet first loads a Llama 2 model and its tokenizer. Then, it calls `get_chat_template` with the tokenizer and the desired `chat_template` ("llama" in this case). This modifies the tokenizer to handle prompts formatted specifically for Llama 2 models. The modified tokenizer is then ready to be used for text generation.

Explanation:

`get_chat_template` takes the tokenizer as input and configures its internal settings (like special tokens and formatting rules) based on the specified `chat_template`. The function returns a *modified* tokenizer that's now ready to process conversational data according to the chosen template.

Key Concepts

`get_chat_template` manages the complexities of different chat templates, allowing you to easily switch between various LLM architectures and their associated prompt formats without manual adjustments to your prompt engineering. It does this by:

1. **Template Selection:** It selects a pre-defined template based on the `chat_template` parameter. Different templates handle different aspects of the input. This is done through a lookup table (`CHAT_TEMPLATES` in `unsloth/chat_templates.py`).
2. **Tokenizer Modification:** Based on the selected template, it modifies the tokenizer. This might involve adding new special tokens, changing the EOS (End Of Sequence) token, or adjusting other internal settings.
3. **System Message Handling:** It allows you to define or customize a system message, which provides instructions or context to the LLM.

Internal Implementation (Simplified)

Here's a simplified sequence diagram showing how `get_chat_template` works:

sequenceDiagram

```

participant User
participant TemplateSelector as get_chat_template
participant TokenizerModifier
participant Tokenizer

User->>TemplateSelector: Calls get_chat_template(tokenizer, "llama")
activate TemplateSelector
TemplateSelector->>TemplateSelector: Selects "llama" template
TemplateSelector->>TokenizerModifier: Modify tokenizer based on template
TokenizerModifier-->>TemplateSelector: Modified tokenizer
TemplateSelector-->>User: Returns modified tokenizer
deactivate TemplateSelector
User<<--Tokenizer: Modified Tokenizer

```

The core logic resides in the `get_chat_template` function within `unsloth/chat_templates.py`. This function uses a dictionary (`CHAT_TEMPLATES`) to store various templates associated with different LLMs. The function selects the appropriate template and then modifies the tokenizer accordingly. The tokenizer modification involves replacing special tokens, which mostly involves updating the tokenizer's internal string representation. Finally, it performs a check for potential overlaps in the added tokens. More intricate internal adjustments may occur, like handling of the EOS token and adding system messages depending on the model.

Conclusion

In this chapter, you learned about `get_chat_template`, a function that simplifies the process of adapting your tokenizer to the specific formatting requirements of different LLMs. This allows for seamless integration with various models and ensures optimal performance. In the next chapter, `unsloth_save_model`, we'll learn how to save and load your Unsloth models efficiently.

Building on our work with `FastLanguageModel` and creating effective prompts with `get_chat_template`, this chapter focuses on saving your hard-earned, finely-tuned language models. Imagine you've spent hours training a model; you wouldn't want to lose it, would you? `unsloth_save_model` is your trusty librarian, carefully storing your precious model (book) in a safe place (your disk or the Hugging Face Hub).

Let's say you've trained a fantastic model using Unsloth and want to save it for later use or share it with the world. `unsloth_save_model` makes this easy, providing various options for saving and deployment.

A Simple Use Case: Saving Your Model Locally

Let's start with the simplest case: saving your model to your local disk.


```

from unsloth.models.loader import FastLanguageModel
from unsloth.save import unsloth_save_model

# Load your trained model and tokenizer (assuming you've already done this)
model, tokenizer = FastLanguageModel.from_pretrained("your_model_name")

```

```

# Save the model to the 'my_saved_model' directory
unsloth_save_model(model, tokenizer, save_directory="my_saved_model")

```

This code snippet saves your `model` and `tokenizer` to a directory named “my_saved_model”. Unsloth handles the details of saving efficiently.

Key Concepts

`unsloth_save_model` offers several key features:

1. **Saving Methods:** It supports different saving methods, including:
 - `"lora"`: Saves only the LoRA (Low-Rank Adaptation) adapters, making the file smaller and faster to load. This is ideal if you plan to load the model onto a powerful machine and you only want to load the modifications you’ve done. This is the fastest option.
 - `"merged_16bit"`: Merges the LoRA weights into the base model weights and saves everything as 16-bit floats. This method is suitable for use with tools like llama.cpp and generating GGUF files, and generally results in a more portable but larger model file.
 - `"merged_4bit"`: Similar to `"merged_16bit"`, but saves the merged weights using 4-bit quantization. This creates the smallest files, but inference is slower. Only use this as a last step.
2. **Hugging Face Hub Integration:** You can easily push your saved model directly to the Hugging Face Hub for sharing with others. This requires a Hugging Face access token.

Using `unsloth_save_model`

Let’s expand on our previous example to include saving to the Hugging Face Hub.

```

from unsloth.save import unsloth_save_model

# ... (Load your model and tokenizer as before) ...

# Save the model to the Hugging Face Hub
repo_id = "your_huggingface_username/your_model_name" # Replace with your info
unsloth_save_model(model, tokenizer, save_directory=repo_id, push_to_hub=True, token="your_token")

```

This saves the model locally *and* uploads it to the Hugging Face Hub under the specified repository ID. Remember to replace `"your_huggingface_username/your_model_name"`

and "your_huggingface_token" with your actual Hugging Face username, model name, and token.

Internal Implementation (Simplified)

Here's a simplified sequence diagram showing the `unsloth_save_model` process:

```
sequenceDiagram
    participant User
    participant Saver as unsloth_save_model
    participant Merger
    participant Uploader
    participant Model

    User->>Saver: Calls unsloth_save_model(model, ...)
    activate Saver
    Saver->>Merger: Merge weights (if needed)
    Merger-->>Saver: Merged model
    Saver->>Uploader: Upload to HF (if push_to_hub=True)
    Uploader-->>Saver: Upload complete
    Saver->>Model: Save model locally
    deactivate Saver
    User<<--Model: Model saved
```

The `unsloth_save_model` function in `unsloth/save.py` manages this. It first merges the LoRA weights (if specified), then handles the local saving and optionally uploads the model to the Hugging Face Hub using the `huggingface_hub` library.

A simplified code snippet showing weight merging:

```
# ... (Inside unsloth_save_model) ...
if save_method == "merged_16bit":
    # Merge LoRA weights into the base model and convert to 16-bit precision.
    # ... (complex merging logic omitted for brevity) ...
    model.save_pretrained(...) # Save the merged model
```

The actual merging logic involves iterating through the model's layers and applying low rank adaptation updates. This is a simplified example to demonstrate the general process. The actual implementation is in `unsloth/save.py`.

Conclusion

This chapter covered `unsloth_save_model`, a crucial function for persisting your trained models. It offers flexibility in saving methods and seamlessly integrates with the Hugging Face Hub. In the next chapter, `patch_saving_functions`, we'll learn more advanced techniques for customizing saving behaviors.

Building on our understanding of saving models with `unsloth_save_model`, this chapter introduces `patch_saving_functions`, a powerful tool that enhances your model saving capabilities. Think of `unsloth_save_model` as a basic tool for saving, while `patch_saving_functions` adds a whole set of advanced features and customization options.

The Problem: Limited Saving Options

In previous chapters, we learned how to save our models using `unsloth_save_model`. However, this function only provided a limited set of saving formats. What if we need to save our model in a specific format for a particular tool or platform? This is where `patch_saving_functions` comes in. It lets us expand the saving options available.

A Central Use Case: Saving to Different Formats

Let's say you want to save your model in a format suitable for `llama.cpp`, a popular library for running large language models on less powerful hardware. `unsloth_save_model` alone doesn't directly support this. `patch_saving_functions` allows us to add this capability.

```
from unsloth.models.loader import FastLanguageModel
from unsloth.save import patch_saving_functions

# Load your model (as shown in previous chapters)
model, tokenizer = FastLanguageModel.from_pretrained("your_model_name")

# Patch the model's saving functions
model = patch_saving_functions(model)

# Now you can use the new saving methods:
# model.save_pretrained_merged("my_model", save_method="merged_16bit") # For llama.cpp
# model.save_pretrained_gguf("my_model_gguf", quantization_method="q4_k_m") # For GGUF
```

This code snippet first loads a model. Then, it calls `patch_saving_functions` to add new methods to the model. These new methods (`save_pretrained_merged` and `save_pretrained_gguf`) allow us to save the model in various formats, optimized for different tools and scenarios. The example shows saving to `merged_16bit` format (good for `llama.cpp`) and `GGUF` format (another efficient format).

Key Concepts

`patch_saving_functions` works by adding several new methods to your model:

1. **`save_pretrained_merged`:** This method merges any LoRA adapters into the base model weights and then saves it using 16-bit precision (float16).

This results in a model file suitable for use with `llama.cpp` and generation of GGUF files.

2. **save_pretrained_gguf:** This method uses `llama.cpp`'s tools to convert the model to the GGUF format, which is incredibly efficient for loading and inference. It lets you specify various quantization levels (`quantization_method`) for further optimization (see `unsloth_save_model` for details).
3. **push_to_hub_merged, push_to_hub_gguf, push_to_hub_ggml:** These methods extend the functionalities of the HuggingFace Hub. `push_to_hub_merged` uploads the merged and quantized models; `push_to_hub_gguf` uploads GGUF versions. `push_to_hub_ggml` converts the LoRA adapters into GGML format before uploading.
4. **save_pretrained_ggml:** This method locally saves the model's LoRA adapters in the GGML format for use with `llama.cpp`.

These new methods extend the base saving functionality, providing greater flexibility and control over the model saving process.

Internal Implementation (Simplified)

Here's a simplified sequence diagram:

```
sequenceDiagram
    participant User
    participant Patcher as patch_saving_functions
    participant Model
    participant NewMethod1 as save_pretrained_merged
    participant NewMethod2 as save_pretrained_gguf

    User->>Patcher: Calls patch_saving_functions(model)
    activate Patcher
    Patcher->>Model: Adds new methods
    Patcher-->>User: Returns patched model
    deactivate Patcher
    User->>Model: Calls save_pretrained_merged(...)
    activate Model
    Model->>NewMethod1: Executes saving logic
    deactivate Model
    User->>Model: Calls save_pretrained_gguf(...)
    activate Model
    Model->>NewMethod2: Executes saving logic
    deactivate Model
```

The `patch_saving_functions` function (in `unsloth/save.py`) iterates through the model and adds the new methods (`save_pretrained_merged`, `save_pretrained_gguf`, etc.) dynamically. It uses Python's `types.MethodType`

to create bound methods, effectively extending the functionality of the model class. The core logic for the new saving methods resides within `unsloth/save.py` (as seen in the previous chapter).

Conclusion

In this chapter, you learned about `patch_saving_functions`, which greatly enhances the flexibility of your model saving process by adding new, powerful saving methods. This allows you to adapt your model saving workflow to the needs of various tools and platforms. In the next chapter, `prepare_model_for_kbit_training`, we'll explore how to prepare your models for training with extremely low bit precision.

Building upon our knowledge of efficient model loading and saving from Chapter 5 and Chapter 6, this chapter introduces `prepare_model_for_kbit_training`, a function that optimizes your model for training with extremely low precision (e.g., 4-bit). This is like switching your car to a more efficient fuel – it still runs, but uses less energy and can be faster.

Why Use Low-Bit Precision?

Large language models (LLMs) can be incredibly memory-intensive. Training them often requires powerful, expensive hardware. `prepare_model_for_kbit_training` helps address this by reducing the precision of your model's weights during training. This reduces memory consumption and can speed up training, but might lead to minor accuracy changes.

A Simple Use Case: Training a 4-bit Quantized Model

Let's say you want to train a smaller, faster LLM using 4-bit precision. You've already loaded your model using `FastLanguageModel`. Now, we'll prepare it for 4-bit training using `prepare_model_for_kbit_training`:

```
from unsloth.models.loader import FastLanguageModel
from unsloth_zoo.utils import prepare_model_for_kbit_training

# Load your model (as shown in Chapter 3)
model, tokenizer = FastLanguageModel.from_pretrained("your_model_name")

# Prepare the model for 4-bit training
model = prepare_model_for_kbit_training(model)

# Now you can proceed with your training using this optimized model.
# We'll cover the training process in later chapters.
```

This code snippet prepares your `model` for training with 4-bit quantization. The function handles the necessary adjustments internally.

Explanation:

The `prepare_model_for_kbit_training` function applies several optimizations:

1. It sets parameters to help handle the numerical instability that can come from low-bit training.
2. It makes sure gradient checkpointing is enabled for better memory efficiency during training. This is a technique that trades computation for reduced memory usage.
3. It configures the model for mixed precision training, where certain parts of the model might be in higher precision to ensure stability.

Key Concepts

`prepare_model_for_kbit_training` primarily focuses on making the model suitable for training with reduced precision by:

1. **Mixed Precision:** Using a combination of 4-bit and higher precision (like 16-bit or 32-bit) to balance accuracy and memory efficiency. This is like using high-octane fuel for the crucial parts of your car's engine and regular fuel for the rest.
2. **Gradient Checkpointing:** A technique to reduce the memory footprint of backpropagation (the process of calculating gradients). It involves re-computing parts of the forward pass instead of storing intermediate results. This saves a lot of memory.

Internal Implementation (Simplified)

Here's a simplified sequence diagram showing how `prepare_model_for_kbit_training` works:

```
sequenceDiagram
    participant User
    participant Preparer as prepare_model_for_kbit_training
    participant Optimizer
    participant Model

    User->>Preparer: Calls prepare_model_for_kbit_training(model)
    activate Preparer
    Preparer->>Optimizer: Apply mixed precision and gradient checkpointing
    Optimizer-->>Preparer: Model optimized
    Preparer->>Model: Return optimized model
    deactivate Preparer
    User<<--Model: Optimized Model
```

The core logic is in `unsloth_zoo/utils.py` within the `prepare_model_for_kbit_training` function. It calls other helper functions within the same file to configure various

aspects of the model for 4-bit training. The `prepare_model_for_training` function then manages the actual application of the techniques mentioned above.

Conclusion

This chapter introduced `prepare_model_for_kbit_training`, which optimizes LLMs for efficient training with low-bit precision. This reduces memory usage and speeds up the training process. In the next chapter, `unsloth_compile_transformers`, we'll explore techniques to further enhance performance through compilation.

Building on our efficient model loading and saving techniques from Chapter 5 and Chapter 6, this chapter introduces `unsloth_compile_transformers`, a function that significantly boosts the speed and memory efficiency of your transformers models. Think of it as adding a turbocharger to your model – it performs the same tasks, but much faster and with less strain on resources.

The Problem: Slow and Memory-Hungry Models

Large language models can be slow and require a lot of memory, especially during inference (using the model to generate text) or finetuning (further training the model on a specific dataset). This can make them impractical for many users or applications. `unsloth_compile_transformers` helps solve this problem.

A Simple Use Case: Speeding Up Inference

Let's say you've loaded a Llama 2 model using `FastLanguageModel` and want to make text generation much faster. `unsloth_compile_transformers` will help.

```
from unsloth.models.loader import FastLanguageModel
from unsloth_zoo.utils import unsloth_compile_transformers

# Load the model (as shown in Chapter 3)
model, tokenizer = FastLanguageModel.from_pretrained("unsloth/Llama-2-7b-chat-hf")

# Compile the model for faster inference
model = unsloth_compile_transformers(model)

# Now, use the compiled model for text generation (as shown in Chapter 3)
```

This code snippet compiles the model using `unsloth_compile_transformers`. The returned `model` is now optimized for speed. The exact speedup depends on your hardware and the model architecture, but you can expect a noticeable improvement, especially on GPUs.

Explanation:

`unsloth_compile_transformers` uses advanced compilation techniques, such as those from the `torch.compile` and `torchdynamo` libraries to transform your model into a more efficient form. It leverages various optimizations depending on your PyTorch and hardware capabilities to speed up inference.

Key Concepts

`unsloth_compile_transformers` utilizes several techniques to improve performance:

1. **Just-in-Time Compilation (JIT):** This technique compiles parts of the model's code into optimized machine code at runtime. It's like compiling a program before running it – the resulting executable is much faster. The `torch.compile` library is used for this purpose.
2. **Graph Optimization:** Before compilation, the model's computational graph is analyzed and simplified to remove redundant operations. This is similar to streamlining a workflow – we remove steps that aren't needed, creating a much more efficient workflow. `torchdynamo` helps in this stage.
3. **Hardware-Specific Optimizations:** `unsloth_compile_transformers` adapts to your specific hardware (e.g., CPU, GPU) to ensure maximum performance. The optimized code will make use of your CPU/GPU's instructions more efficiently.

Internal Implementation (Simplified)

Here's how `unsloth_compile_transformers` works:

```
sequenceDiagram
    participant User
    participant Compiler as unsloth_compile_transformers
    participant Optimizer
    participant Model

    User->>Compiler: Calls unsloth_compile_transformers(model)
    activate Compiler
    Compiler->>Optimizer: Optimize model graph
    Optimizer-->>Compiler: Optimized graph
    Compiler->>Model: Compile optimized graph
    deactivate Compiler
    User<<--Model: Compiled Model
```

The function in `unsloth_zoo/utils.py` manages the compilation process. It first optimizes the model's computational graph and then calls `torch.compile` to compile the optimized graph. More advanced optimizations, such as flash attention, are also applied if supported.

Here's a simplified version of the code:


```
# ... (Inside unsloth_compile_transformers) ...
import torch
model = torch.compile(model) # Apply compilation to achieve the faster speedups.
return model
```

This is a significantly simplified version – the real implementation uses the `torchdynamo` library for more sophisticated optimizations and also checks for GPU support. It also does a careful check for multiple PyTorch versions.

Conclusion

In this chapter, you learned about `unsloth_compile_transformers`, a function that significantly enhances the performance of your transformers models. By using JIT compilation and graph optimization, it achieves noticeable speed and memory improvements. In the next chapter, `LoRA_QKV`, `LoRA_MLP`, `LoRA_W`, we’ll explore Low-Rank Adaptation (LoRA) techniques for efficient model finetuning.

Building on our understanding of efficient model compilation from Chapter 8, this chapter introduces `LoRA_QKV`, `LoRA_MLP`, and `LoRA_W`. These are special tools that make fine-tuning large language models (LLMs) much faster and more memory-efficient. Fine-tuning is like giving your already-trained model a bit more specialized training on a new task, without having to train the entire model from scratch.

Imagine you have a very powerful, pre-trained model capable of general tasks. You want to specialize this model to answer questions about a specific topic like 19th-century literature. Training the whole model again from scratch for this new task would take a long time and require lots of computational power. `LoRA_QKV`, `LoRA_MLP`, and `LoRA_W` allow you to do this fine-tuning much more efficiently, saving time and resources.

These tools use a technique called Low-Rank Adaptation (LoRA). Instead of changing all the weights of your model, LoRA makes small, targeted changes to specific parts. Think of it as precisely adjusting the knobs and dials of a complex machine instead of completely rebuilding it.

Key Concepts

LoRA works by adding small, low-rank matrices to certain parts of the model. These are the “LoRA weights”. Unsloth provides these specialized functions to efficiently compute these updates.

1. **LoRA_QKV:** This function handles the LoRA updates for the Query (Q), Key (K), and Value (V) matrices in the attention mechanism of a transformer model. The attention mechanism is a key part of how transformers process information. `LoRA_QKV` makes these updates super efficient.

2. **LoRA_MLP:** This handles LoRA updates for the Multi-Layer Perceptron (MLP) layers in the transformer. MLP layers are another critical component for processing information. LoRA_MLP optimizes the computations involved in updating these layers.
3. **LoRA_W:** This is a more general-purpose function that handles LoRA updates for other weight matrices in the model, outside the attention and MLP layers. It allows applying LoRA to any linear layer with ease.

These functions define how the LoRA updates are calculated during both the forward (using the model to make predictions) and backward (calculating gradients to improve the model) passes. They're written as PyTorch autograd functions, which means PyTorch automatically handles the complex math of backpropagation for you.

A Simple Use Case: Fine-tuning with LoRA

Let's see how to use these functions (though directly using these is not common, it's handled behind the scenes in `FastLanguageModel`). We'll focus on `LoRA_MLP` for simplicity.

```
# Simplified example - actual usage is much more complex.
import torch
from unsloth.kernels.fast_lora import apply_lora_mlp_swiglu

# Assume 'model' is a transformer model and 'X' is some input data.
# This is a placeholder for a real model's MLP layer.
mlp_layer = torch.nn.Linear(100, 200)

# Apply LoRA update to the MLP layer. The actual implementation is complex,
# but this shows the high-level usage of the function:
updated_output = apply_lora_mlp_swiglu(mlp_layer, X)

print(f"Original output shape: {mlp_layer(X).shape}")
print(f"Updated output shape: {updated_output.shape}")
```

This code snippet shows the high-level usage of `apply_lora_mlp_swiglu`, which internally uses `LoRA_MLP`. It applies the LoRA update to the `mlp_layer` using `apply_lora_mlp_swiglu` and provides an updated output.

Internal Implementation (Simplified)

Here's a simplified sequence diagram showing how `LoRA_MLP` works:

```
sequenceDiagram
    participant Input as Input as Input Data
    participant LoRA_MLP
    participant Forward as Forward Pass
```

```

participant Backward Pass
participant Output as Updated Output

Input->>LoRA_MLP: Forward pass
activate LoRA_MLP
LoRA_MLP->>Forward Pass: Compute LoRA updates
Forward Pass-->>LoRA_MLP: Updated weights
LoRA_MLP-->>Output: Output data
deactivate LoRA_MLP
Output->>LoRA_MLP: Backward pass
activate LoRA_MLP
LoRA_MLP->>Backward Pass: Compute gradients for LoRA weights
Backward Pass-->>LoRA_MLP: Gradients
deactivate LoRA_MLP

```

LoRA_MLP (and similar functions) are implemented in `unsloth/kernels/fast_lora.py`. They use custom PyTorch autograd functions to define the forward and backward passes. These functions efficiently compute the LoRA updates and their gradients, which are then used by the optimizer to improve the model. The core of the implementation is based on optimized matrix multiplications and utilizes low-precision arithmetic for faster computation.

Conclusion

This chapter introduced LoRA_QKV, LoRA_MLP, and LoRA_W, which are essential for efficient fine-tuning using LoRA. They provide optimized implementations for updating different parts of a transformer model. This leads to faster training and reduces memory usage compared to training the entire model. In the next chapter, Triton Kernels (e.g., `_cross_entropy_forward`, `layernorm_forward`), we'll explore using Triton kernels for further performance gains.

Building on our exploration of efficient model compilation with `unsloth_compile_transformers`, this chapter dives into Triton kernels. These are like highly specialized, super-fast tools that make certain parts of your large language model (LLM) run significantly faster. They're written in a special language called Triton, designed for optimal performance on GPUs.

Imagine you're building with LEGOs. Standard LEGO bricks are great for general construction, but for specific tasks, like creating highly detailed gears, you might use specialized pieces that are much more efficient. Triton kernels are like these specialized LEGO pieces for your LLMs.

A Central Use Case: Speeding Up Loss Calculation

Let's say you're training a large language model and the calculation of the cross-entropy loss (a measure of how well your model is doing) is taking too long. Triton kernels can drastically speed this up. `_cross_entropy_forward` is

a Triton kernel specifically designed for this task.

```
from unsloth.kernels import fast_cross_entropy_loss
```

```
# Assume 'logits' are your model's predictions and 'labels' are the correct answers.
```

```
logits = torch.randn(1, 1024, 1000) # Example logits, shape (batch_size, sequence_length, vocab_size)
```

```
labels = torch.randint(0, 1000, (1, 1024)) # Example labels, shape (batch_size, sequence_length)
```

```
loss = fast_cross_entropy_loss(logits, labels)
```

```
print(f"Loss: {loss}")
```

This code snippet calculates the cross-entropy loss using `fast_cross_entropy_loss`.

This function internally uses the `_cross_entropy_forward` Triton kernel for faster computation. The output is a single number representing the loss.

Key Concepts: Triton Kernels

Triton kernels are functions written in the Triton language, which compiles to highly efficient GPU code. They offer significant speedups for specific operations compared to standard PyTorch code. Key features include:

1. **GPU Optimization:** Triton kernels are designed to exploit the parallel processing capabilities of GPUs, resulting in faster execution. They're written in a way that leverages many threads simultaneously.
2. **Specialized Operations:** Unlike general-purpose code, Triton kernels are highly specialized. `_cross_entropy_forward`, for instance, is designed exclusively for calculating cross-entropy loss, allowing for very efficient computation.
3. **Low-Level Control:** Triton gives fine-grained control over memory access patterns, maximizing GPU performance. The code is written at a lower level than PyTorch, resulting in more efficient use of GPU resources.
4. **Automatic Differentiation:** Triton kernels integrate seamlessly with PyTorch's automatic differentiation. This means that gradients (which are needed for training) are automatically calculated correctly, even though the kernel is highly optimized GPU code.

Internal Implementation (Simplified)

Here's a simplified view of how `fast_cross_entropy_loss` utilizes `_cross_entropy_forward`:

```
sequenceDiagram
    participant User
    participant FastCELoss as fast_cross_entropy_loss
    participant TritonKernel as _cross_entropy_forward
    participant GPU
```

participant Result

```
User->>FastCELoss: Calls fast_cross_entropy_loss(logits, labels)
activate FastCELoss
FastCELoss->>TritonKernel: Passes logits and labels
activate TritonKernel
TritonKernel->>GPU: Executes on GPU
GPU-->>TritonKernel: Loss calculated
TritonKernel-->>FastCELoss: Returns loss
deactivate TritonKernel
FastCELoss-->>User: Returns loss
deactivate FastCELoss
User<<--Result: Loss value
```

The `fast_cross_entropy_loss` function in `unsloth/kernels/cross_entropy_loss.py` handles the high-level logic. It prepares the inputs and then calls the `_cross_entropy_forward` kernel, which runs on the GPU. The results are then returned to the main program.

The `_cross_entropy_forward` kernel itself is defined in `unsloth/kernels/cross_entropy_loss.py`. It's written in Triton and uses low-level operations to efficiently calculate the cross-entropy loss on the GPU. Similarly, `layernorm_forward` in `unsloth/kernels/layernorm.py` is another example of a Triton kernel designed to speed up layer normalization, a common operation in LLMs.

A simplified snippet of `_cross_entropy_forward` :

```
# ... (Inside _cross_entropy_forward) ...
logits = tl.load(logits_ptr + col_offsets, mask = mask, other = -float("inf")).to(tl.float32)
# ... (compute loss efficiently using Triton primitives) ...
tl.store(loss_ptr, loss) #Store the calculated loss
# ...
```

This snippet shows how the kernel directly accesses memory (`logits_ptr`) and performs computations using Triton's primitives. This low-level access and specialized operations are what make Triton kernels significantly faster.

Conclusion

This chapter introduced Triton kernels, which significantly boost the speed of specific operations in LLMs. We saw how `_cross_entropy_forward` accelerates loss calculation, illustrating the power of these specialized functions. In the next chapter, we will continue our exploration of performance optimization techniques.

Unsloth is a Python library designed to significantly speed up and reduce the memory footprint of large language models (*LLMs*) during training and inference. It achieves this through various optimizations, including *Triton kernel acceleration*, **Low-Rank Adaptation (LoRA)** for efficient fine-tuning, and

advanced compilation techniques applied to the `transformers` library. Unsloth also provides tools for efficient model saving in formats like *GGUF* and deployment on platforms like the Hugging Face Hub and Ollama.

Source Repository: <https://github.com/unslothai/unsloth.git>

```

flowchart TD
    A0["FastLanguageModel"]
    A1["FastModel"]
    A2["FastBaseModel"]
    A3["unsloth_save_model"]
    A4["get_chat_template"]
    A5["patch_saving_functions"]
    A6["unsloth_compile_transformers"]
    A7["prepare_model_for_kbit_training"]
    A8["LoRA_QKV, LoRA_MLP, LoRA_W"]
    A9["Triton Kernels (e.g., `_cross_entropy_forward`, `layernorm_forward`)"]
    A0 -- "Is a type of" --> A1
    A1 -- "Inherits from" --> A2
    A1 -- "Uses" --> A9
    A1 -- "Uses" --> A6
    A1 -- "Uses" --> A8
    A3 -- "Calls" --> A5
    A3 -- "Saves" --> A0
    A3 -- "Saves" --> A2
    A6 -- "Optimizes" --> A9
    A7 -- "Prepares" --> A1
    A4 -- "Configures for" --> A0
    A8 -- "Used by" --> A1

```

Chapters

1. FastBaseModel
2. FastModel
3. FastLanguageModel
4. get_chat_template
5. unsloth_save_model
6. patch_saving_functions
7. prepare_model_for_kbit_training
8. unsloth_compile_transformers
9. LoRA_QKV, LoRA_MLP, LoRA_W
10. Triton Kernels (e.g., _cross_entropy_forward, layernorm_forward)