

Step-by-Step Tutorials

Focused, practical tutorials that teach you how to accomplish specific tasks step by step. Learn by doing with clear instructions and working code examples.

Contents

01	Run MTP LLMs with llama.cpp & vLLM	3
-----------	------------------------------------	----------

Run MTP LLMs with llama.cpp & vLLM

What you'll build: By the end of this tutorial, you will be able to set up and run Multi-Token Prediction (MTP) capable LLMs locally using `llama.cpp` and `vLLM`, compare their performance against standard generation, and understand fallback options. **Time needed:** ~90 minutes **Prerequisites:** Basic command-line interface (CLI) familiarity, Git installed, C++ compiler (GCC/Clang for Linux/macOS, MSVC for Windows), CMake installed, Python 3.9+ installed, NVIDIA GPU with CUDA (11.8+ recommended) or AMD GPU with ROCm, or Apple Silicon (Metal), Sufficient RAM (16GB+ recommended) and VRAM (8GB+ recommended) **Version used:** llama.cpp: main branch (post MTP merge); vLLM: latest stable/developer preview with MTP support

Understanding Multi-Token Prediction (MTP)

Large Language Models (LLMs) traditionally generate text token by token in an autoregressive fashion. This means they predict one token, then use that token as part of the input to predict the next, and so on. While effective, this sequential process can be slow, especially for longer outputs.

Multi-Token Prediction (MTP), often referred to as **Speculative Decoding** or **Assisted Generation**, is a clever technique designed to speed up this process significantly. Instead of predicting one token at a time, MTP uses a smaller, faster "draft model" (also an LLM) to predict a sequence of tokens. The main, larger LLM then verifies this entire sequence in parallel.

Here's how it generally works:

1. **Draft Generation:** A smaller, faster draft model quickly generates `N` speculative tokens based on the current context.
2. **Parallel Verification:** The main, larger LLM takes the original context plus the `N` speculative tokens as input. It then predicts the next `N` tokens in parallel.

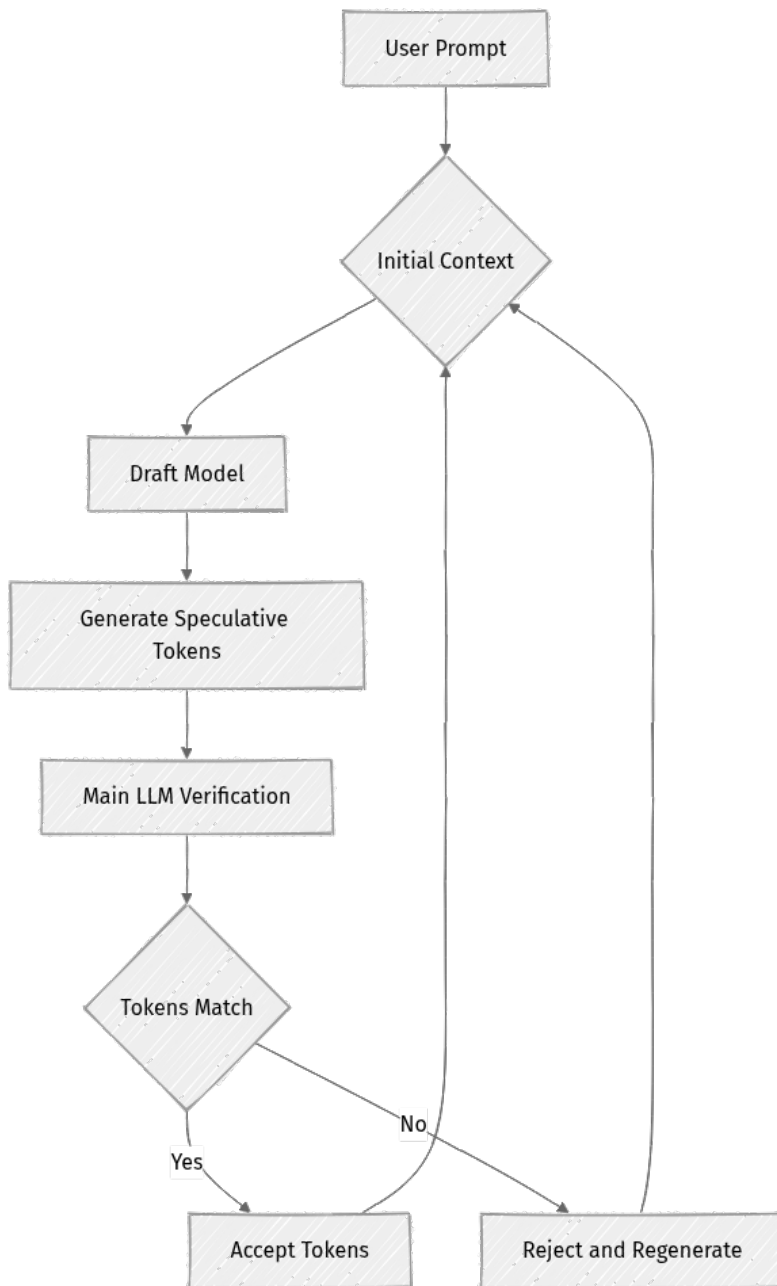
3. **Acceptance & Rejection:** The main LLM's predictions are compared against the draft model's predictions.

- If a speculative token matches the main LLM's prediction, it's accepted.
- If a speculative token doesn't match, it and all subsequent speculative tokens are rejected. The main LLM then proceeds to generate the next correct token from the point of divergence, and the process repeats.

This approach can provide substantial speedups (often 1.5x to 2x or more) because the main LLM processes multiple tokens in a single forward pass, reducing the number of costly sequential operations. It's like having a fast assistant write a draft that you then quickly review and correct, rather than writing every word yourself from scratch.

Why MTP Matters for Local LLM Inference

For anyone running LLMs locally, inference speed is a critical factor. MTP offers a way to get faster responses from your local models without needing to upgrade your GPU or use a smaller model. This makes powerful LLMs more usable and interactive on consumer hardware. `llama.cpp` and `vLLM` are at the forefront of bringing this technology to the local inference ecosystem.



📌 Key Idea: MTP (Speculative Decoding) uses a fast draft model to predict multiple tokens, which are then verified in parallel by the main LLM, leading to faster overall generation without sacrificing accuracy.

System Prerequisites & Environment Setup (OS-Specific)

Before we dive into running LLMs, let's ensure your system has all the necessary tools. We'll cover common operating systems.

⚡ Note: The specific commands might vary slightly based on your OS version and how you've configured your package manager. Always refer to official documentation if you encounter issues.

1. Git Installation

Git is essential for cloning the `llama.cpp` and `vLLM` repositories.

- **Linux (Debian/Ubuntu):**

```
sudo apt update
sudo apt install git -y
```

- **Linux (Fedora/RHEL):**

```
sudo dnf install git -y
```

- **macOS:** Git is often pre-installed or can be installed via Xcode Command Line Tools.

```
xcode-select --install # Follow prompts if not already installed
```

Alternatively, with Homebrew:

```
brew install git
```

- **Windows:** Download and install from the [official Git website](#). Follow the installer prompts.

2. C++ Compiler & Build Tools (CMake)

`llama.cpp` is a C++ project, so you'll need a compiler and CMake to build it.

- **Linux:**

```
sudo apt install build-essential cmake -y # Debian/Ubuntu
sudo dnf install gcc-c++ make cmake -y # Fedora/RHEL
```

- **macOS:** Xcode Command Line Tools (installed with Git above) typically include Clang (C++ compiler) and make. Install CMake with Homebrew:

```
brew install cmake
```

- **Windows (MSVC):** Install [Visual Studio Community Edition](#). During installation, select the "Desktop development with C++" workload. This includes MSVC and CMake.

3. Python 3.9+

Both `vLLM` and parts of the `llama.cpp` ecosystem (like model conversion scripts) rely on Python.

- **Linux/macOS:** Many systems have Python pre-installed. Verify your version:

```
python3 --version
```

If it's older than 3.9, consider using `pyenv` or `conda` for a managed environment, or install via your package manager.

```
# Example for Ubuntu, may vary
sudo apt install python3.9 python3.9-venv python3-pip -y
```

- **Windows:** Download and install from the [official Python website](#). Make sure to check "Add Python to PATH" during installation.

4. GPU Drivers & Toolkits (Crucial for Performance)

This is where performance lives! Ensure your GPU drivers and associated toolkits are up-to-date.

- **NVIDIA GPU (CUDA):**

- Install the latest NVIDIA drivers for your GPU.
- Install [CUDA Toolkit](#) (11.8+ recommended). Follow NVIDIA's installation guide carefully for your OS.
- Verify CUDA installation:

```
nvcc --version
```


- **AMD GPU (ROCm):**

- Install the latest AMDGPU drivers and [ROCm toolkit](#). ROCm support is primarily for Linux.
- Verify ROCm installation:

```
rocminfo
```

- **Apple Silicon (Metal):**

- Metal support is built into macOS. Ensure your macOS is up-to-date. No separate toolkit installation is typically required beyond Xcode Command Line Tools.

 Common mistake: Outdated GPU drivers or incorrect CUDA/ROCm setup are the leading causes of poor performance or build failures when enabling GPU acceleration. Double-check these steps!

Environment Verification

After installing, open a new terminal and run these commands to confirm everything is ready:

```
git --version
cmake --version
python3 --version
# For NVIDIA:
nvcc --version
```

```
# For AMD:
rocm_info
```

You should see version numbers for each. If any fail, revisit the installation steps for that tool.

By the end of this section, your system is prepared with the foundational tools needed to compile and run advanced LLM inference engines.

Running MTP with llama.cpp

`llama.cpp` is a highly optimized C/C++ inference engine for LLMs, known for its efficiency and broad hardware support. Multi-Token Prediction (MTP) has been merged into its `main` branch, offering significant speedups.

1. Clone and Build llama.cpp

First, let's get the latest `llama.cpp` source code.

```
git clone https://github.com/ggerganov/llama.cpp.git
cd llama.cpp
```

Now, we'll build it with GPU support. Choose the command appropriate for your hardware.

Build for NVIDIA CUDA (Linux/Windows WSL)

```
make clean
LLAMA_CUBLAS=1 make -j # -j uses all available CPU cores for faster
compilation
```

 Note: For Windows native (MSVC), you'd typically use CMake:

```
> mkdir build && cd build
> cmake .. -DLLAMA_CUBLAS=ON -DCMAKE_BUILD_TYPE=Release
> cmake --build . --config Release
>
```

The `main` executable will be in `build/bin/Release` or `build/bin`.

Build for AMD ROCm (Linux)

```
make clean
LLAMA_ROCM=1 make -j
```

Build for Apple Silicon (Metal)

```
make clean
LLAMA_METAL=1 make -j
```

Build for CPU Only (Fallback)

If you don't have a supported GPU or want to test CPU performance:

```
make clean
make -j
```

After compilation, you should see an executable named `main` (or `main.exe` on Windows) in the `llama.cpp` directory (or `build/bin/Release` for MSVC).

To verify the build, run the `main` executable without arguments:

```
./main
```

You should see help text describing its usage. If you see "command not found" or an error, something went wrong with the build.

2. Acquire MTP-Capable GGUF Models

For MTP to work in `llama.cpp`, you need two models:

1. **The main model:** The large LLM you want to run.
2. **The draft model:** A smaller, faster model used for speculative decoding.

Crucially, both models need to be **MTP-capable** or at least compatible in terms of tokenization and architecture. Models like Qwen-Next (e.g., Qwen1.5-0.5B-Chat-GGUF, Qwen1.5-1.8B-Chat-GGUF) and Gemma-4 are good candidates. The draft model should be significantly smaller than the main model for efficiency gains.

We'll use a Qwen model as an example. You'll need to download these models in the GGUF format. Quantized versions (e.g., Q4_K_M) are highly recommended for local inference.

Let's download a main model and a draft model. We'll use Hugging Face for this. You can use `wget` or `curl` on Linux/macOS, or manually download. For convenience, we'll use `huggingface-cli`.

First, install `huggingface-hub`:

```
pip install huggingface-hub
```


If you need to access gated models (which is rare for GGUF files but good practice):
````bash huggingface-cli login`

## You will be prompted to enter your Hugging Face token.

```
Now, let's download example Qwen models. We'll put them in a `models` directory.
```bash
mkdir -p models

# Main model: Qwen1.5-7B-Chat (Q4_K_M quantization)
huggingface-cli download Qwen/Qwen1.5-7B-Chat-GGUF Qwen1.5-7B-Chat-Q4_K_M.gguf
--local-dir models --local-dir-use-symlinks False

# Draft model: Qwen1.5-0.5B-Chat (Q4_K_M quantization)
huggingface-cli download Qwen/Qwen1.5-0.5B-Chat-GGUF Qwen1.5-0.5B-Chat-Q4_K_M.gguf
--local-dir models --local-dir-use-symlinks False
```

 **Note:** Always check the Hugging Face model pages (e.g., [Qwen1.5-7B-Chat-GGUF](#)) for the exact file names and available quantizations. The `Qwen1.5-0.5B-Chat-GGUF` model is a good choice for a draft model due to its small size.

Verify that the `.gguf` files are in your `llama.cpp/models` directory:

```
ls models/*.gguf
```

You should see `Qwen1.5-7B-Chat-Q4_K_M.gguf` and `Qwen1.5-0.5B-Chat-Q4_K_M.gguf`.

3. Running llama.cpp with MTP

Now we can run the `main` executable, enabling MTP. We'll compare performance with and without MTP.

Standard Generation (No MTP)

Let's first establish a baseline tokens-per-second (t/s) without MTP.

```
./main -m models/Qwen1.5-7B-Chat-Q4_K_M.gguf \  
-p "Tell me a short story about a brave knight and a wise dragon." \  
-n 256 \  
--temp 0.7 \  
--seed 42 \  
--log-stats
```

- `-m`: Path to the main model.
- `-p`: Your prompt.
- `-n`: Maximum number of tokens to generate.
- `--temp`: Temperature for generation (0.0-2.0, lower is more deterministic).
- `--seed`: Random seed for reproducibility.
- `--log-stats`: Prints detailed performance statistics at the end, including t/s.

Run this command and note the `t/s` value reported at the end. This is your baseline.


MTP-Enabled Generation

Now, let's enable MTP using the draft model.

```
./main -m models/Qwen1.5-7B-Chat-Q4_K_M.gguf \  
--mtp-draft-model models/Qwen1.5-0.5B-Chat-Q4_K_M.gguf \  
-p "Tell me a short story about a brave knight and a wise dragon." \  
-n 256 \  
--temp 0.7 \  
--seed 42 \  
--log-stats
```

- `--mtp-draft-model`: Path to your smaller draft model. This flag automatically enables MTP.

Run this command and compare the `t/s` output with your baseline. You should observe a noticeable improvement in tokens per second.

 **Common mistake:** Using incompatible models for main and draft, or a draft model that's too large, can degrade performance or lead to incorrect generation. Ensure your draft model is significantly smaller and uses the same tokenizer/architecture if possible.

You've successfully built `llama.cpp`, downloaded MTP-compatible models, and run your first MTP-accelerated LLM inference! You should have a clear `t/s` comparison between standard and MTP generation.

Running MTP/Speculative Decoding with vLLM

`vLLM` is a fast and easy-to-use LLM serving library that excels in throughput for serving multiple concurrent requests, and it also supports speculative decoding (MTP) for single-request latency improvements.

1. Install vLLM

`vLLM` is a Python library. We'll install it using `pip`. It's highly recommended to install `vLLM` with specific backend support for your GPU.

```
# Create and activate a virtual environment (best practice)
python3 -m venv venv-vllm
source venv-vllm/bin/activate # On Windows: .\venv-vllm\Scripts\activate

# Install vLLM with CUDA support
# For NVIDIA GPUs:
pip install vllm[cuda]

# For AMD GPUs (ROCm):
# pip install vllm[rocm]

# For Apple Silicon (Metal):
# pip install vllm[core] # Metal support is often included in core or specific
versions
# Or check vLLM docs for specific Metal installation if needed, it's evolving.
```

⚡ Note: `vLLM`'s Metal support is newer and might require specific builds or versions. Always check the [official vLLM documentation](#) for the most up-to-date installation instructions for Apple Silicon.

After installation, verify `vLLM` can be imported:

```
python -c "import vllm; print(vllm.__version__)"
```

You should see the installed `vLLM` version.

2. Acquire Hugging Face Models for vLLM

`vLLM` works directly with models from the Hugging Face Hub. It will download the necessary model weights automatically if they're not cached locally. Similar to `llama.cpp`, you'll need a main model and a draft model that are compatible. Qwen-Next models are excellent choices.

For this section, we'll use the same Qwen models, but `vLLM` will handle the full precision (or FP8/AWQ/GPTQ) weights directly from Hugging Face, not the GGUF format.

- **Main model:** `Qwen/Qwen1.5-7B-Chat`
- **Draft model:** `Qwen/Qwen1.5-0.5B-Chat`

If you haven't already, ensure you have `huggingface-cli` installed and are logged in if you plan to use gated models:``bash pip install huggingface-hub huggingface-cli login # If needed

```
### 3. Running `vLLM` with Speculative Decoding (MTP)

`vLLM` provides a powerful API and a server for inference. We'll use the
`vllm.entrypoints.llm` module or the `python -m vllm.entrypoints.api_server`
command to demonstrate.

#### Standard Generation (No Speculative Decoding)

Let's run `vLLM` without speculative decoding first to get a baseline. We'll
use a simple Python script for clarity.

Create a file named `vllm_generate.py`:
```python
from vllm import LLM, SamplingParams
import time

Configuration for the main model
model_name = "Qwen/Qwen1.5-7B-Chat" # Or a different compatible model

Initialize the LLM without speculative decoding
llm = LLM(model=model_name,
 tensor_parallel_size=1, # Adjust based on your GPU setup
 gpu_memory_utilization=0.9, # Adjust as needed
 quantization=None # Or 'awq', 'gptq', 'fp8' if using quantized
models
)

Define the prompt and sampling parameters
prompt = "Tell me a short story about a brave knight and a wise dragon."
sampling_params = SamplingParams(temperature=0.7, top_p=0.95, max_tokens=256)

print(f"Generating with model: {model_name} (Standard)")
start_time = time.time()
outputs = llm.generate([prompt], sampling_params)
end_time = time.time()
```

```

Process and print output
total_output_tokens = 0
for output in outputs:
 generated_text = output.outputs[0].text
 num_output_tokens = len(output.outputs[0].token_ids)
 total_output_tokens += num_output_tokens
 print(f"Prompt: {output.prompt}")
 print(f"Generated text ({num_output_tokens} tokens): {generated_text}")

duration = end_time - start_time
tokens_per_second = total_output_tokens / duration if duration > 0 else 0
print(f"\n--- Performance (Standard) ---")
print(f"Total output tokens: {total_output_tokens}")
print(f"Generation duration: {duration:.2f} seconds")
print(f"Tokens per second (t/s): {tokens_per_second:.2f}")

Deactivate venv if you activated it earlier
deactivate

```

Run this script:

```
python vllm_generate.py
```

**vLLM** will download the model weights if not cached. This might take a while the first time. Note the **Tokens per second (t/s)** value.

### MTP-Enabled Generation (Speculative Decoding)

Now, let's enable speculative decoding. **vLLM** uses the **speculative\_model** parameter for this.

Modify **vllm\_generate.py** to include the **speculative\_model** parameter:

```

from vllm import LLM, SamplingParams
import time

Configuration for the main model
model_name = "Qwen/Qwen1.5-7B-Chat" # Main model
speculative_model_name = "Qwen/Qwen1.5-0.5B-Chat" # Draft model

Initialize the LLM with speculative decoding
llm = LLM(model=model_name,
 speculative_model=speculative_model_name, # Enable speculative
 decoding
 tensor_parallel_size=1,
 gpu_memory_utilization=0.9,
 quantization=None # Or 'awq', 'gptq', 'fp8'
)

Define the prompt and sampling parameters
prompt = "Tell me a short story about a brave knight and a wise dragon."
sampling_params = SamplingParams(temperature=0.7, top_p=0.95, max_tokens=256)

print(f"Generating with main model: {model_name}, speculative model: {speculative_model_name}")
start_time = time.time()

```

```

outputs = llm.generate([prompt], sampling_params)
end_time = time.time()

Process and print output
total_output_tokens = 0
for output in outputs:
 generated_text = output.outputs[0].text
 num_output_tokens = len(output.outputs[0].token_ids)
 total_output_tokens += num_output_tokens
 print(f"Prompt: {output.prompt}")
 print(f"Generated text ({num_output_tokens} tokens): {generated_text}")

duration = end_time - start_time
tokens_per_second = total_output_tokens / duration if duration > 0 else 0
print(f"\n--- Performance (Speculative Decoding) ---")
print(f"Total output tokens: {total_output_tokens}")
print(f"Generation duration: {duration:.2f} seconds")
print(f"Tokens per second (t/s): {tokens_per_second:.2f}")


Deactivate venv if you activated it earlier
deactivate

```

Run the modified script:

```
python vllm_generate.py
```

Again, `vLLM` will download the draft model if not cached. Compare the `Tokens per second (t/s)` from this run with the standard generation. You should see a similar performance boost as with `llama.cpp`.

 **Real-world insight:** `vLLM` is often used to run an API server for LLMs, handling multiple requests concurrently. While this tutorial focuses on single-request MTP, `vLLM`'s architecture is highly optimized for high throughput scenarios in production.

By completing this section, you've successfully installed `vLLM`, configured it to use Hugging Face models, and demonstrated speculative decoding for faster LLM inference.

## Performance Comparison & Analysis

Now that you've run both `llama.cpp` and `vLLM` with and without MTP/speculative decoding, let's analyze the results.

You should have four key `t/s` metrics:

1. `llama.cpp` standard generation

2. `llama.cpp` MTP-enabled generation
3. `vLLM` standard generation
4. `vLLM` speculative decoding-enabled generation

Create a simple table or just compare the numbers you collected.

Engine	Configuration	Tokens/Second (t/s)	Speedup Factor (vs. standard)
<code>llama.cpp</code>	Standard	(Your Value)	1.0x
<code>llama.cpp</code>	MTP-enabled	(Your Value)	> 1.0x (e.g., 1.5x)
<code>vLLM</code>	Standard	(Your Value)	1.0x
<code>vLLM</code>	Speculative Decoding	(Your Value)	> 1.0x (e.g., 1.5x)

## Expected Observations

- **MTP/Speculative Decoding Speedup:** You should consistently observe a significant increase in tokens per second when MTP or speculative decoding is enabled for both `llama.cpp` and `vLLM`. This is the primary benefit of these techniques.
- **Engine Differences:** `vLLM` often achieves higher raw `t/s` than `llama.cpp` for unquantized or heavily optimized models, especially with larger batch sizes or specific hardware. However, `llama.cpp` excels in memory efficiency and broader hardware support, making it ideal for running larger models or on less powerful hardware.
- **Draft Model Impact:** The effectiveness of MTP heavily depends on the draft model. A good draft model (fast and accurate enough to predict tokens reliably) maximizes the acceptance rate of speculative tokens, leading to higher speedups. A poor draft model can actually slow things down if the main model frequently rejects its predictions.
- **Hardware Impact:** Your GPU's VRAM and compute capabilities (CUDA cores, Tensor Cores) play a huge role. More powerful GPUs will naturally yield higher `t/s` values across the board.

- **Quantization:** `llama.cpp` often uses highly quantized GGUF models (e.g., Q4\_K\_M), which are very memory efficient but might have a slight accuracy trade-off. `vLLM` can use full-precision models or more advanced quantization methods like AWQ/FP8, which might offer different performance/accuracy profiles.

## When to Use Each

- **llama.cpp with MTP:** Ideal for local inference on consumer hardware, especially if you need to run large models with limited VRAM (thanks to GGUF quantization) and want single-request latency improvements. It's excellent for interactive chat interfaces or local development.
- **vLLM with Speculative Decoding:** Best suited for serving LLMs as an API endpoint, where high throughput for multiple concurrent requests is critical. It also provides excellent single-request latency with speculative decoding and supports various advanced serving features.

By analyzing these results, you've gained practical insight into the performance benefits of Multi-Token Prediction and the strengths of `llama.cpp` and `vLLM` in different scenarios.

---

## Troubleshooting & Common Issues

Working with bleeding-edge LLM inference can sometimes hit snags. Here are some common issues and how to address them.

### 1. llama.cpp Build Errors

- **"Undefined reference to `cublas...`" or `rocm` related errors:** This almost always means your CUDA/ROCm toolkit is not correctly installed, or the `LLAMA_CUBLAS=1` / `LLAMA_ROCM=1` environment variables weren't correctly picked up by `make`.
  - **Solution:**
    - Ensure CUDA/ROCm is installed and its `bin` directory is in your system's `PATH`.
    - Re-run `make clean` before `make` to ensure a fresh build.
    - Verify `nvcc --version` (for CUDA) or `rocminfo` (for ROCm) works.

- **"fatal error: blas.h: No such file or directory"**: This indicates missing BLAS development libraries.
  - **Solution (Linux)**: `sudo apt install libopenblas-dev` or `sudo dnf install openblas-devel`.
- **CMake configuration errors on Windows**: Ensure you selected "Desktop development with C++" when installing Visual Studio, and that `cmake` and `cl.exe` (MSVC compiler) are accessible from your terminal.

## 2. vLLM Installation and Runtime Issues

- **vllm[cuda] or vllm[rocm] fails to install**: This often points to an issue with your Python environment finding the correct GPU toolkit.
  - **Solution:**
    - Ensure your CUDA/ROCm installation is correct and compatible with your Python version.
    - Try installing `torch` with CUDA/ROCm separately first to verify your environment: `pip install torch torchvision torchaudio --index-url <https://download.pytorch.org/whl/cu118>` (for CUDA 11.8, adjust as needed). If `torch` installs and runs on GPU, `vLLM` should follow.
    - Use a fresh virtual environment to avoid conflicts.
- **vLLM "CUDA Out of Memory" (OOM)**: This happens when the model (or models, with speculative decoding) is too large for your GPU's VRAM.
  - **Solution:**
    - Use smaller models.
    - Reduce `gpu_memory_utilization` (e.g., `gpu_memory_utilization=0.8`).
    - For `vLLM`, consider using quantization (e.g., `quantization='awq'`) if the model supports it and you install the necessary dependencies (`pip install vllm[awq]`).
    - For `llama.cpp`, use highly quantized GGUF models (e.g., `Q4_K_M`).

- **"No module named 'vllm'":** You likely forgot to activate your virtual environment or installed `vLLM` in a different environment.
  - **Solution:** `source venv-vllm/bin/activate` (Linux/macOS) or `.\venv-vllm\Scripts\activate` (Windows) before running your Python script.

### 3. MTP/Speculative Decoding Specific Issues

- **MTP performance is worse than standard:**
  - **Incompatible Draft Model:** The draft model might be too different from the main model, leading to a low acceptance rate of speculative tokens.
  - **Draft Model Too Large:** If the draft model is too large, the overhead of running both models outweighs the benefits. The draft model should be significantly smaller and faster.
  - **Solution:** Experiment with different draft models, ensuring they are from the same family or architecture as the main model, and are much smaller (e.g., 0.5B or 1B for a 7B main model).
- **Unexpected or nonsensical token generation:** While rare with `llama.cpp` and `vLLM`, sometimes MTP can, in experimental stages, lead to slightly degraded output quality if not perfectly implemented or if the models are highly divergent.
  - **Solution:** Ensure you are using the latest `main` branch of `llama.cpp` and the latest stable/developer preview of `vLLM`. If the issue persists, try disabling MTP to confirm it's the cause.

### 4. General LLM Inference Issues

- **Model download failures:** Check your internet connection. For Hugging Face models, ensure you have sufficient disk space and that `huggingface-cli login` was successful if accessing gated models.
- **Slow inference on CPU:** If your GPU setup isn't working, `llama.cpp` will fall back to CPU. This will be significantly slower.
  - **Solution:** Revisit GPU driver and toolkit installation steps. Verify GPU build flags for `llama.cpp`.
- **Python version conflicts:** Always use a virtual environment to isolate project dependencies and avoid system-wide Python conflicts.

⚠️ What can go wrong: MTP is an advanced optimization. While generally robust, unexpected model behavior or performance regressions can occur if the main and draft models are poorly matched or if there are underlying issues with the framework's implementation.

By being aware of these common pitfalls, you can more effectively debug and resolve issues as you explore MTP-accelerated LLM inference.

## Fallback Options & Alternatives

Multi-Token Prediction and speculative decoding offer fantastic performance boosts, but they depend on specific model compatibility and a well-configured environment. What if MTP isn't supported for your chosen model, or you encounter persistent issues? Here are some robust fallback options and alternatives for local LLM inference.

### 1. Standard llama.cpp Generation (CPU or GPU)

Even without MTP, `llama.cpp` remains an incredibly efficient engine. If MTP isn't an option, running your GGUF model with standard generation is still very performant, especially with GPU acceleration.

- **When to use:** When your model isn't MTP-capable, or you prefer maximum stability over bleeding-edge speed.
- **How:** Simply omit the `--mtp-draft-model` flag when running `llama.cpp/main`.

```
./main -m models/your_model.gguf \
-p "Your prompt here." \
-n 128 \
--temp 0.7 \
--log-stats
```

## 2. Standard vLLM Generation (GPU)

vLLM provides excellent performance for single-request and high-throughput scenarios even without speculative decoding, thanks to its PagedAttention mechanism.

- **When to use:** If your model isn't compatible with vLLM's speculative decoding, or you primarily need high throughput for multiple requests rather than single-request latency.
- **How:** Omit the `speculative_model` parameter when initializing `vllm.LLM`.

```
from vllm import LLM, SamplingParams
llm = LLM(model="Your/HuggingFace-Model")
... rest of your code
```

## 3. Ollama

Ollama is an increasingly popular tool for running LLMs locally. It packages models into easy-to-use bundles and provides a simple command-line interface and API. It abstracts away many of the complexities of `llama.cpp` (which it often uses under the hood) and model management.

- **When to use:** For extreme ease of use, quick experimentation with a variety of models, and if you prefer a simplified interface over direct `llama.cpp` commands. Ollama now also supports speculative decoding for some models.
- **How:**
  1. Download and install Ollama from their website.
  2. Pull a model: `ollama pull llama3`
  3. Run the model: `ollama run llama3 "Why is the sky blue?"`
  4. Check for MTP/speculative support in specific model tags or Ollama versions.

## 4. LM Studio / Jan.ai / GPT4All

These are desktop applications that provide a graphical user interface (GUI) for downloading and running LLMs locally. They typically integrate `llama.cpp` or similar engines.

- **When to use:** If you prefer a visual interface, easy model browsing, and don't want to deal with the command line for basic inference. They are great for beginners.

- **How:** Download, install, and use their respective GUIs. Check their settings for any "speculative decoding" or "multi-token prediction" options, as these are increasingly being integrated.

⚡ **Quick Note:** While these GUI tools simplify the experience, they might not always expose the granular control or the very latest experimental features (like MTP flags) that you get by directly compiling and running `llama.cpp` or scripting `vLLM`.

By understanding these alternatives, you're equipped to run LLMs locally efficiently, even when the most advanced optimizations aren't immediately available or suitable for your specific needs.

---

## What to Build Next

You've successfully set up and experimented with Multi-Token Prediction using `llama.cpp` and `vLLM`. This is a fantastic foundation for building more advanced LLM applications. Here are three ideas to extend your learning:

### 1. Build a Local Chatbot with MTP:

- Integrate your MTP-enabled `llama.cpp` or `vLLM` setup into a simple Python-based chatbot.
- Use a framework like `Gradio` or `Streamlit` to create a web UI.
- Compare the responsiveness of your MTP-enabled chatbot against a standard one. Focus on the perceived latency for the user.
- **Challenge:** Implement a streaming response from the LLM to the UI for an even smoother user experience.

### 2. Experiment with Different MTP-Capable Models and Quantizations:

- Download other MTP-capable models (e.g., different Qwen-Next sizes, Gemma-4 if available in GGUF/HF format) for both main and draft roles.
- Test various quantization levels (e.g., Q2\_K, Q8\_0 for `llama.cpp`, or AWQ/FP8 for `vLLM` if your hardware supports it).
- **Challenge:** Create a script that automates the performance comparison across a matrix of main model, draft model, and quantization combinations, outputting the results to a CSV file.

### 3. Implement a `vLLM` API Server with Speculative Decoding:

- Instead of a simple Python script, set up `vLLM` as a persistent API server using `uvicorn` or `gunicorn`.
- Configure it to use speculative decoding.
- Write a separate client script (e.g., in Python or JavaScript) that sends requests to your `vLLM` API server.
- **Challenge:** Implement basic load testing to simulate multiple concurrent users querying your `vLLM` server, observing how speculative decoding impacts overall throughput and average latency under load.