

is a large language model tool. It has use cases such as: Code completion Prose composition Chatbot that passes the Turing test Text/image summarization and analysis The following options are available: Print version and exit. Show help message and exit. Puts program in command line interface mode. This flag is implied when a prompt is supplied using either the or flags. Puts program in server mode. This will launch an HTTP server on a local port. This server has both a web UI and an OpenAI API compatible completions endpoint. When the server is run on a desk system, a tab browser tab will be launched automatically that displays the web UI. This flag is implied if no prompt is specified, i.e. neither the or flags are passed. Model path in the GGUF file format. Default: Specifies path of the LLaVA vision model in the GGUF file format. If this flag is supplied, then the and flags should also be supplied. Random Number Generator (RNG) seed. A random seed is used if this is less than zero. Default: -1 Number of threads to use during generation. Default: \$(nproc)/2 Set the number of threads to use during batch and prompt processing. In some systems, it is beneficial to use a higher number of threads during batch processing than during generation. If not specified, the number of threads used for batch processing will be the same as the number of threads used for generation. Default: Same as Number of threads to use during generation. Default: Same as Number of threads to use during batch and prompt processing. Default: Same as Prefix BOS to user inputs, preceding the string. This flag is used to add a prefix to your input, primarily, this is used to insert a space after the reverse prompt. Here's an example of how to use the flag in conjunction with the flag: Default: empty This flag is used to add a suffix after your input. This is useful for adding an "Assistant:" prompt after the user's input. It's added after the new-line character (\n) that's automatically added to the end of the user's input. Here's an example of how to use the flag in conjunction with the flag: Default: empty Number of tokens to predict. -1 = infinity -2 = until context filled Default: -1 Set the size of the prompt context. A larger context size helps the model to better comprehend and generate responses for longer input or conversations. The LLaMA models were built with a context of 2048, which yields the best results on longer input / inference. 0 = loaded automatically from model Default: 512 Batch size for prompt processing. Default: 512 Top-k sampling. 0 = disabled Default: 40 Top-p sampling. 1.0 = disabled Default: 0.9 Min-p sampling. 0.0 = disabled Default: 0.1 Tail free sampling, parameter z. 1.0 = disabled Default: 1.0 Locally typical sampling, parameter p. 1.0 = disabled Default: 1.0 Last n tokens to consider for penalize. 0 = disabled -1 = ctx_size Default: 64 Penalize repeat sequence of tokens. 1.0 = disabled Default: 1.1 Repeat alpha presence penalty. 0.0 = disabled Default: 0.0 Repeat alpha frequency penalty. 0.0 = disabled Default: 0.0 Use Mirostat sampling. Top K, Nucleus, Tail Free and Locally Typical samplers are ignored if used.. 0 = disabled 1 = Mirostat 2 = Mirostat 2.0 Default: 0 Mirostat learning rate, parameter eta. Default: 0.1 Mirostat target entropy, parameter tau. Default: 5.0 Modifies the likelihood of token appearing in the completion, i.e. to increase likelihood of token or to decrease likelihood of token Draft model for speculative decoding. Default: Negative prompt to use for guidance.. Default: empty Negative prompt file to use for guidance. Default: empty Strength of guidance. 1.0 = disable Default: 1.0 RoPE frequency scaling method, defaults to linear unless specified by the model RoPE context scaling factor, expands context by a factor of where is the linear scaling factor used by the fine-tuned model. Some fine-tuned models have extended the context length by scaling RoPE. For example, if the original pre-trained model have a context length (max sequence length) of 4096 (4k) and the fine-tuned model have 32k. That is a scaling factor of 8, and should work by setting the above to 32768 (32k) and to 8. RoPE base frequency, used by NTK-aware scaling. Default: loaded from model RoPE frequency scaling factor, expands context by a factor of 1/N YaRN: original context size of model. Default: 0 = model training context size YaRN: extrapolation mix factor. 0.0 = full interpolation Default: 1.0 YaRN: scale sqrt(t) or attention magnitude. Default: 1.0 YaRN: high correction dim or alpha. Default: 1.0 YaRN: low correction dim or beta. Default: 32.0 Ignore end of stream token and continue generating (implies Do not penalize newline token. Temperature. Default: 0.8 Return logits for all tokens in the batch. Default: disabled Compute HellaSwag score over random tasks from datafile supplied with -f Number of tasks to use when computing the HellaSwag score. Default: 400 This flag allows users to retain the original prompt when the model runs out of context, ensuring a connection to the initial instruction or conversation topic is maintained, where is the number of tokens from the initial prompt to retain when the model resets its internal context. 0 = no tokens are kept from initial prompt -1 = retain all tokens from initial prompt Default: 0 Number of tokens to draft for speculative decoding. Default: 16 Max number of chunks to process. -1 = all Default: -1 Number of sequences to decode. Default: 1 speculative decoding accept probability. Default: 0.5 Speculative decoding split probability. Default: 0.1 Force system to keep model in RAM rather than swapping or compressing. Do not memory-map model (slower load but may reduce pageouts if not using mlock). Attempt optimizations that help on some NUMA systems if run without this previously, it is recommended to drop the system page cache before using this. See <https://github.com/ggerganov/llama.cpp/issues/1437>. Force GPU support to be recompiled at runtime if possible. Never compile GPU support at runtime. If the appropriate DSO file already exists under then it'll be linked as-is without question. If a prebuilt DSO is present in the PKZIP content of the executable, then it'll be extracted and linked if possible. Otherwise, will skip any attempt to compile GPU support and simply fall back to using CPU inference. Specifies which brand of GPU should be used. Valid choices are: Use any GPU if possible, otherwise fall back to CPU inference (default) Use Apple Metal GPU. This is only available on MacOS ARM64. If Metal could not be used for any reason, then a fatal error will be raised. Use AMD GPUs. The AMD HIP ROCm SDK should be installed in which case we assume the HIP_PATH environment variable has been defined. The set of gfx microarchitectures needed to run on the host machine is determined automatically based on the output of the hipInfo command. On Windows, release binaries are distributed with a tinyBLAS DLL so it'll work out of the box without requiring the HIP SDK to be installed. However, tinyBLAS is slower than rocBLAS for batch and image processing, so it's recommended that the SDK be installed anyway. If an AMD GPU could not be used for any reason, then a fatal error will be raised. Use NVIDIA GPUs. If an NVIDIA GPU could not be used for any reason, a fatal error will be raised. On Windows, NVIDIA GPU support will use our tinyBLAS library, since it works on stock Windows in-

stalls. However, tinyBLAS goes slower for batch and image processing. It's possible to use NVIDIA's closed-source cuBLAS library instead. To do that, both MSVC and CUDA need to be installed and the command should be run once from the x64 MSVC command prompt with the flag passed. The GGML library will then be compiled and saved to so the special process only needs to happen a single time. Never use GPU and instead use CPU inference. This setting is implied by Number of layers to store in VRAM. Number of layers to store in VRAM for the draft model. How to split the model across multiple GPUs, one of: none: use one GPU only layer (default): split layers and KV across GPUs row: split rows across GPUs When using multiple GPUs this option controls how large tensors should be split across all GPUs. is a comma-separated list of non-negative values that assigns the proportion of data that each GPU should get in order. For example, "3,2" will assign 60% of the data to GPU 0 and 40% to GPU 1. By default the data is split in proportion to VRAM but this may not be optimal for performance. Requires cuBLAS. How to split tensors across multiple GPUs, comma-separated list of proportions, e.g. 3,1 The GPU to use for scratch and small tensors. Use cuBLAS instead of custom mul_mat_q CUDA kernels. Not recommended since this is both slower and uses more VRAM. Print prompt before generation. Use basic IO for better compatibility in subprocesses and limited consoles. Apply LoRA adapter (implies Apply LoRA adapter with user defined scaling S (implies Optional model to use as a base for the layers modified by the LoRA adapter Disables pledge() sandboxing on Linux and OpenBSD. Samplers that will be used for generation in the order, separated by semicolon, for example: top_k;tfs;typical;top_p;min_p;temp Simplified sequence for samplers that will be used. Run in chatml mode (use with ChatML-compatible models) Verbose print of the KV cache. Disable KV offload. KV cache data type for K. KV cache data type for V. Group-attention factor. Default: 1 Group-attention width. Default: 512 Binary file containing multiple choice tasks. Compute Winogrande score over random tasks from datafile supplied by the flag. Number of tasks to use when computing the Winogrande score. Default: 0 Compute multiple choice score over random tasks from datafile supplied by the flag. Number of tasks to use when computing the multiple choice score. Default: 0 Computes KL-divergence to logits provided via the flag. Save logits to filename. Print token count every tokens. Default: -1 The following options may be specified when is running in mode. Process prompt escapes sequences (\n, \r, \t, \', \", \\) Prompt to start text generation. Your LLM works by auto-completing this text. For example: Stands a pretty good chance of printing Lincoln's Gettysburg Address. Prompts can take on a structured format too. Depending on how your model was trained, it may specify in its docs an instruction notation. With some models that might be: In most cases, simply colons and newlines will work too: Prompt file to start generation. BNF-like grammar to constrain which tokens may be selected when generating text. For example, the grammar: will force the LLM to only output yes or no before exiting. This is useful for shell scripts when the flag is also supplied. File to read grammar from. File to cache prompt state for faster startup. Default: none If specified, saves user input and generations to cache as well. Not supported with or other interactive options. If specified, uses the prompt cache but does not update it. Start with a randomized prompt. Path to an image file. This should be used with multimodal models. Alternatively, it's possible to embed an image directly into the prompt instead; in which case, it must be base64 encoded into an HTML img tag URL with the image/jpeg MIME type. See also the flag for supplying the vision model. Run the program in interactive mode, allowing users to engage in real-time conversations or provide specific instructions to the model. Run the program in interactive mode and immediately wait for user input before starting the text generation. Run the program in instruction mode, which is specifically designed to work with Alpaca models that excel in completing tasks based on user instructions. Technical details: The user's input is internally prefixed with the reverse prompt (or "### Instruction:" as the default), and followed by "### Response:" (except if you just press Return without any input, to keep generating a longer response). By understanding and utilizing these interaction options, you can create engaging and dynamic experiences with the LLaMA models, tailoring the text generation process to your specific needs. Specify one or multiple reverse prompts to pause text generation and switch to interactive mode. For example, can be used to jump back into the conversation whenever it's the user's turn to speak. This helps create a more interactive and conversational experience. However, the reverse prompt doesn't work when it ends with a space. To overcome this limitation, you can use the flag to add a space or any other characters after the reverse prompt. Enable colored output to differentiate visually distinguishing between prompts, user input, and generated text. Don't echo the prompt itself to standard output. Allows you to write or paste multiple lines without ending each in '\'. The following options may be specified when is running in mode. Port to listen Default: 8080 IP address to listen. Default: 127.0.0.1 Server read/write timeout in seconds. Default: 600 Number of slots for process requests. Default: 1 Enable continuous batching (a.k.a dynamic batching). Default: disabled Set a file to load a system prompt (initial prompt of all slots), this is useful for chat applications. Set an alias for the model. This will be added as the field in completion responses. Path from which to serve static files. Default: Enable embedding vector output. Default: disabled Do not attempt to open a web browser tab at startup. Set the group attention factor to extend context size through self-extend. The default value is which means disabled. This flag is used together with Set the group attention width to extend context size through self-extend. The default value is This flag is used together with The following log options are available: Path under which to save YAML logs (no logging if unset) Run simple logging test Disable trace logs Enable trace logs Specify a log filename (without extension) Create a separate new log file on start. Each log file will have unique name: Don't truncate the old log file. Here's an example of how to run llama.cpp's built-in HTTP server. This example uses LLaVA v1.5-7B, a multimodal LLM that works with llama.cpp's recently-added support for image inputs. llamafile \

```

-m llava-v1.5-7b-Q8_0.gguf \
--mproj llava-v1.5-7b-mproj-Q8_0.gguf \
--host 0.0.0.0 Here's an example of how to generate code for a libc function using the llama.cpp command line interface, utilizing WizardCoder-Python-13B weights: llamafile \
-m wizardcoder-python-13b-v1.0.Q8_0.gguf --temp 0 -r '}' \n' -r '```\n' \

```

```

-e -p ""c\nvoid *memcpy(void *dst, const void *src, size_t size) {\n' Here's a similar example that instead utilizes
Mistral-7B-Instruct weights for prose composition: llamafile \
-m mistral-7b-instruct-v0.2.Q5_K_M.gguf \
-p '[INST]Write a story about llamas[/INST]' Here's an example of how llamafile can be used as an interactive
chatbot that lets you query knowledge contained in training data: llamafile -m llama-65b-Q5_K_M.gguf -p ' The fol-
lowing is a conversation between a Researcher and their helpful AI assistant Digital Athena which is a large lan-
guage model trained on the sum of human knowledge. Researcher: Good morning. Digital Athena: How can I help
you today? Researcher:' --interactive --color --batch_size 1024 --ctx_size 4096 \ --keep -1 --temp 0 --mirostat 2
--in-prefix ' ' --interactive-first \ --in-suffix 'Digital Athena:' --reverse-prompt 'Researcher:' Here's an example of
how you can use llamafile to summarize HTML URLs: (
echo '[INST]Summarize the following text:'
links -codepage utf-8 \
    -force-html \
    -width 500 \
    -dump https://www.poetryfoundation.org/poems/48860/the-raven |
sed 's/ */ /g'
echo '[/INST]' ) | llamafile \
-m mistral-7b-instruct-v0.2.Q5_K_M.gguf \
-f /dev/stdin \
-c 0 \
--temp 0 \
-n 500 \
--no-display-prompt 2>/dev/null Here's how you can use llamafile to describe a jpg/png/gif/bmp image: llamafile
--temp 0 \
--image lemurs.jpg \
-m llava-v1.5-7b-Q4_K_M.gguf \
--mmproj llava-v1.5-7b-mmproj-Q4_0.gguf \
-e -p '### User: What do you see?\n### Assistant: ' \
--no-display-prompt 2>/dev/null If you wanted to write a script to rename all your image files, you could use the
following command to generate a safe filename: llamafile --temp 0 \
--image ~/Pictures/lemurs.jpg \
-m llava-v1.5-7b-Q4_K_M.gguf \
--mmproj llava-v1.5-7b-mmproj-Q4_0.gguf \
--grammar 'root ::= [a-z]+ (" " [a-z]+)+' \
-e -p '### User: The image has...\n### Assistant: ' \
--no-display-prompt 2>/dev/null |
sed -e's/ /_/' -e's/\/.jpg/' three_baby_lemurs_on_the_back_of_an_adult_lemur.jpg Here's an example of how to
make an API request to the OpenAI API compatible completions endpoint when your is running in the background
in mode. curl -s http://localhost:8080/v1/chat/completions \
-H "Content-Type: application/json" -d '{
"model": "gpt-3.5-turbo",
"stream": true,
"messages": [
{
"role": "system",
"content": "You are a poetic assistant."
},
{
"role": "user",
"content": "Compose a poem that explains FORTRAN."
}
] }' | python3 -c ' import json import sys json.dump(json.load(sys.stdin), sys.stdout, indent=2) print() The flag
needs to be passed in order to use GPUs made by NVIDIA and AMD. It's not enabled by default since it sometimes
needs to be tuned based on the system hardware and model architecture, in order to achieve optimal performance,
and avoid compromising a shared display.

```