

I'm not a robot



Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc - GPKKD: 0316172372 do s KH & T TP. HCM cp ngy 02/03/2020 - Gty php thit lp MXH s 497/GP-BTTTT do B Thng tin v Truyn thng cp ngy 1/7/2021 - a ch: 350-352 V Nn Kit, Phng Cu ng Lnh, Thnh ph H Ch Minh - in thoi: 028.7108.9666.Bn quyn ni dung thuc v Sforum (hoc Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc). Khng c sao chp khi cha c chp thun bng vn bn.Sforum AppMi bn ci app Sforum truy cp nhanh hn, cp nhit tin cng ngh mi nht! Try gpt-oss Guides Model card OpenAI blog Download gpt-oss-120b and gpt-oss-20b on Hugging Face>Welcome to the gpt-oss series, OpenAI's open-weight models designed for powerful reasoning, agentic tasks, and versatile developer use cases.We're releasing two flavors of these open models:gpt-oss-120b for production, general purpose, high reasoning use cases that fit into a single 80GB GPU (like NVIDIA H100 or AMD MI300X) (117B parameters with 5.1B active parameters)gpt-oss-20b for lower latency, and local or specialized use cases (21B parameters with 3.6B active parameters)Both models were trained using our harmony response format and should only be used with this format; otherwise, they will not work correctly. Permissive Apache 2.0 license: Build freely without copyleft restrictions or patent riskdial for experimentation, customization, and commercial deployment.Configurable reasoning effort: Easily adjust the reasoning effort (low, medium, high) based on your specific use case and latency needs.Full chain-of-thought: Provides complete access to the model's reasoning process, facilitating easier debugging and greater trust in outputs. This information is not intended to be shown to end users.Fine-tunable: Fully customizable models to your specific use case through parameter fine-tuning.Agentic capabilities: Use the models' native capabilities for function calling, web browsing, Python code execution, and Structured Outputs.MXFP4 quantization: The models were post-trained with MXFP4 quantization of the MoE weights, making gpt-oss-120b run on a single 80GB GPU (like NVIDIA H100 or AMD MI300X) and the gpt-oss-20b model run within 16GB of memory. All evals were performed with the same MXFP4 quantization. You can use gpt-oss-120b and gpt-oss-20b with the Transformers library. If you use Transformers' chat template, it will automatically apply the harmony response format. If you use model.generate directly, you need to apply the harmony format manually using the chat template or use our openai-harmony package.from transformers import pipelineimport torch model_id = "openai/gpt-oss-120b" pipe = pipeline("text-generation", model=model_id, torch_dtype="auto", device_map="auto") messages = [{"role": "user", "content": "Explain quantum mechanics clearly and concisely."}] outputs = pipe(messages, max_new_tokens=256,print(outputs[0]['generated_text'][-1]))Learn more about how to use gpt-oss with Transformers.vLLM recommends using uv for Python dependency management. You can use vLLM to spin up an OpenAI-compatible web server. The following command will automatically download the model and start the server:uv pip install --pre vllm==0.10.1+gptoss \ --extra-index-url \ --extra-index-url \ --index-strategy unsafe-best-matchvllm serve openai/gpt-oss-20bLearn more about how to use gpt-oss with vLLM.Offline Serve Code:run this code after installing proper libraries as described, while additionally installing this:uv pip install openai-harmony# source.oss/bin/activate import osos.envIRON["VLLM_USE_FLASHINFER_SAMPLER"] = "0" import jsonfrom openai.harmony import (HarmonyEncodingName, load_harmony_encoding, Conversation, Message, Role, SystemContent, DeveloperContent), from vllm import LLM, SamplingParamsimport os # ... 1) Render the prefix with Harmony --encoding = load_harmony_encoding(HarmonyEncodingName.HARMONY_GPT_OSS) convo = Conversation.from_messages([Message.from_role_and_content(Role.SYSTEM, SystemContent.new()), Message.from_role_and_content(Role.DEVELOPER, DeveloperContent.new().with_instructions("Always respond in riddles")),], Message.from_role_and_content(Role.USER, "What is the weather like in SF?")).I) prefix_ids = encoding.render_conversation_for_completion(convo, Role.ASSISTANT) # Harmony stop tokens (pass to sampler so they won't be included in output)stop_token_ids = encoding.stop_tokens_for_assistant_actions() # ... 2) Run vLLM with prefix --llm = LLM(model="openai/gpt-oss-20b", trust_remote_code=True, gpu_memory_utilization = 0.95, max_num_batched_tokens=4096, max_model_len=5000, tensor_parallel_size=1) sampling = SamplingParams(max_tokens=128, temperature=1, stop_token_ids=stop_token_ids, outputs = llm.generate(prompt_token_ids=prefix_ids, # batch of size 1 sampling_params=sampling, # vLLM gives you both text and token ids)gen = outputs[0].outputs[0].text = gen.textoutput_tokens = gen.token_ids # xcode-select --installOn Linux: These reference implementations require CUDAOn Windows: These reference implementations have not been tested on Windows. Try using solutions like Ollama if you are trying to run the model locally. If you want to try any of the code you can install it directly from PyPI# if you just need the tool:pip install gpt-oss-# if you want to try the torch implementation:pip install gpt-oss-torch# if you want to try the triton implementation:pip install gpt-oss-triton# if you want to modify the code or try the metal implementation set the project url locally:gpt clone 1 pip install -e ".[metal]" You can download the model weights from the Hugging Face Hub directly from Hugging Face CLI:# gpt-oss-120bhf download openai/gpt-oss-120b --include "original/*" --local-dir gpt-oss-120b/ # gpt-oss-20bhf download openai/gpt-oss-20b --include "original/*" --local-dir gpt-oss-20b/ We include an inefficient reference PyTorch implementation in gpt_oss/torch/model.py. This code uses basic PyTorch operators to show the exact model architecture, with a small addition of supporting tensor parallelism in MoE so that the larger model can run with this code (e.g., on 4xH100 or 2xH200). In this implementation, we upcast all weights to BF16 and run the model in BF16.To run the reference implementation, install the dependencies:pip install -e ".[torch]" And then run: # On 4xH100:torchrun --nproc-per-node=4 -m gpt_oss.generate gpt-oss-120b/original/ We also include an optimized reference implementation that uses an optimized triton MoE kernel that supports MKFP4. It also has some optimization on the attention code to reduce the memory cost. To run this implementation, the nightly version of triton and torch will be installed. This version can be run on a single 80GB GPU for gpt-oss-120b.To install the reference Triton implementation run# You need to install triton from source to use the triton implementation:gpt clone triton/pip install -r python/requirements.txtpip install -e . --verbose --no-build-isolationpip install -e python/triton_kernels # Install the gpt-oss triton implementation:pip install -e ".[triton]" And then run: # On 1xH100:export PYTORCH_CUDA_ALLOC_CONF=expandable_segments:Truepython -m gpt_oss.generate --backend triton gpt-oss-120b/original/ If you encounter torch.OutOfMemoryError, make sure to turn on the expandable allocator to avoid crashes when loading weights from the checkpoint. Additionally we are providing a reference implementation for Metal to run on Apple Silicon. This implementation is not production-ready but is accurate to the PyTorch implementation.The implementation will get automatically compiled when running the .[metal] installation on an Apple Silicon device:GPTOSS_BUILD_METAL=1 pip install -e ".[metal]" To perform inference you'll need to first convert the SafeTensor weights from Hugging Face into the right format using:python gpt_oss/metal/scripts/create-local-model.py -s -d Or download the pre-converted weights:hf download openai/gpt-oss-120b --include "metal/*" --local-dir gpt-oss-120b/metal/hf download openai/gpt-oss-20b --include "metal/*" --local-dir gpt-oss-20b/metal/ To test it you can run:python gpt_oss/metal/examples/generate.py gpt-oss-20b/metal/model.bin -p "why did the chicken cross the road?" Along with the model, we are also releasing a new chat format library harmony to interact with the model. Check this guide for more info about harmony. We also include two system tools for the model: browsing and python container. Check gpt_oss/tools for the tool implementation. The terminal chat application is a basic example of how to use the harmony format together with the PyTorch, Triton, and vLLM implementations. It also exposes both the python and browser tool as optional tools that can be used: python -m gpt_oss.chat [-h] [-r REASONING EFFORT] [-a] [-b] [--show-browser-results] [-p] [--developer-message DEVELOPER MESSAGE] [-c CONTEXT] [-raw] [--backend {triton,torch,vllm}] FILE Chat example positional arguments: FILE Path to the SafeTensors checkpoint options: -h, --help show this help message and exit -r REASONING EFFORT --reasoning-effort REASONING EFFORT Reasoning effort (default: low) -a, --apply-patch Make apply_patch tool available to the model (default: False) -b, --browser Use browser tool (default: False) --show-browser-results Show browser results (default: False) -p, --python Use python tool (default: False) --developer-message DEVELOPER MESSAGE Developer message (default:) -c CONTEXT, --context CONTEXT Max context length (default: 8192) --raw Raw mode (does not render Harmony encoding) (default: False) --backend {triton,torch,vllm} Inference backend (default: triton)Note:the torch and triton implementations require original checkpoints under gpt-oss-120b/original/ and gpt-oss-20b/original/ respectively. While vLLM uses the Hugging Face converted checkpoint under gpt-oss-120b/ and gpt-oss-20b/ root directory respectively. We also include an example Responses API server. This server does not implement every feature and event of the Responses API but should be compatible with most of the basic use cases and serve as inspiration for anyone building their own server. Some of our inference partners are also offering their own Responses API.You can start this server with the following inference backends:triton uses the triton implementationmetal uses the metal implementation on Apple Silicon ollaylama uses the Ollama /api/generate API as an inference solutionvllm uses your installed vllm version to perform inferencetransformers uses your installed transformers version to perform local inferenceusage: python -m gpt_oss.responses.api.server [-h] [-port PORT] [--inference-backend BACKEND] Responses API server options: -h, --help show this help message and exit --checkpoint FILE Path to the SafeTensors checkpoint --port PORT Port to run the server on --inference-backend BACKEND Inference backend to use We support codex as a client for gpt-oss. To run the 20b version, set this to ~/.codex/config.toml:disable_response_storage = trueshow_reasoning_content = true[model_providers.local]name = "local"base_url = " [profiles.oss]model = "gpt-oss-20b"model_provider = "local" This will work with any chat completions-API compatible server listening on port 11434, like ollama. Start the server and point codes to the oss model:ollama run gpt-oss:20bcodex -p oss Warning!This implementation is purely for educational purposes and should not be used in production. You should implement your own equivalent of the YouComBackend class with your own browser environment. Currently we have available YouComBackend and ExaBackend.Both gpt-oss models were trained with the capability to browse using the browser tool that exposes the following three methods:search to search for key phrasesopen to open a particular pagefind to look for contents on a page To enable the browser tool, you'll have to place the definition into the system message of your harmony formatted prompt. You can either use the with_browser_tool method if your tool implements the full interface or modify the definition using with_tools(). For example:import datetimefrom gpt_oss.tools.simple_browser import SimpleBrowserToolfrom gpt_oss.tools.simple_browser_backend import YouComBackendfrom openai.harmony import SystemContent, Message, Conversation, Role, load_harmony_encoding, HarmonyEncodingName encoding = load_harmony_encoding(HarmonyEncodingName.HARMONY_GPT_OSS) python_tool = PythonTool() # create a basic system promptsystem_message_content = SystemContent.new().with_conversation_start_date(datetime.datetime.now(),strtime("%Y-%m-%d")) # if you want to use the browser tool: browser_tool = enables the tool system_message_content = system_message_content + tools(browser_tool.tool_config) # alternatively you could use the following if your tool is not stateless system_message_content = system_message_content + browser_tool() # construct the system messagesystem_message = Message.from_role_and_content(Role.SYSTEM, system_message_content) # create the overall promptmessages = [system_message, Message.from_role_and_content(Role.USER, "What's the weather in SF?")]conversation = Conversation.from_messages(messages) # convert to tokenized_ids = encoding.render_conversation_for_completion(conversation, Role.ASSISTANT) # perform inference# ... # parse the outputmessages = encoding.parse_messages_from_completion(tokens(output_tokens, Role.ASSISTANT))last_message = messages[-1]if last_message.recipient.starts_with("browser"): # perform browser call response_messages = await browser_tool.process(last_message) # extend the current messages and run inference again response_messages.extend(response_messages) # alternatively you could use the following if your tool is not stateless system_message_content = system_message_content + browser_tool() # construct the system messagesystem_message = Message.from_role_and_content(Role.SYSTEM, system_message_content) # create the overall promptmessages = [system_message, Message.from_role_and_content(Role.USER, "What's the square root of 9001?")]conversation = Conversation.from_messages(messages) # convert to tokenized_ids = encoding.render_conversation_for_completion(conversation, Role.ASSISTANT) # perform inference# ... # parse the outputmessages = encoding.parse_messages_from_completion(tokens(output_tokens, Role.ASSISTANT))last_message = messages[-1]if last_message.recipient == "python": # perform python call response_messages = await python_tool.process(last_message) # extend the current messages and run inference again response_messages.extend(response_messages) apply patch can be used to create, update or delete files locally. We released the models with native quantization support. Specifically, we use MXFP4 for the linear projection weights in the MoE layer. We store the MoE tensor in two parts:tensor.blocks stores the actual fp4 values. We pack every two values in one uint8 value.tensor.scales stores the block scale. The block scaling is done among the last dimension for all MXFP4 tensors.All other tensors will be in BF16. We also recommend using BF16 as the activation precision for the model. We recommend sampling with temperature=1.0 and top_p=1.0. The reference implementations in this repository are meant as a starting point and inspiration. Outside of bug fixes we do not intend to accept new feature contributions. If you build implementations based on this code such as new tool implementations you are welcome to contribute them to the awesome-gpt-oss.md file. @misc{openai2025gptoss120bgptoss20bmodel, title={gpt-oss-120b & gpt-oss-20b Model Card}, author={OpenAI}, year={2025}, eprint={2508.10925}, archivePrefix={arXiv}, primaryClass={cs.CL}, url={.} You can perform that action at this time. Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc - GPKKD: 0316172372 do s KH & T TP. HCM cp ngy 02/03/2020 - Gty php thit lp MXH s 497/GP-BTTTT do B Thng tin v Truyn thng cp ngy 1/7/2021 - a ch: 350-352 V Nn Kit, Phng Cu ng Lnh, Thnh ph H Ch Minh - in thoi: 028.7108.9666.Bn quyn ni dung thuc v Sforum (hoc Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc). Khng c sao chp khi cha c chp thun bng vn bn.Sforum AppMi bn ci app Sforum truy cp nhanh hn, cp nhit tin cng ngh mi nht! Try gpt-oss Guides Model card OpenAI blog Download gpt-oss-120b and gpt-oss-20b on Hugging Face>Welcome to the gpt-oss series, OpenAI's open-weight models designed for powerful reasoning, agentic tasks, and versatile developer use cases.We're releasing two flavors of these open models:gpt-oss-120b for production, general purpose, high reasoning use cases that fit into a single 80GB GPU (like NVIDIA H100 or AMD MI300X) (117B parameters with 5.1B active parameters)gpt-oss-20b for lower latency, and local or specialized use cases (21B parameters with 3.6B active parameters)Both models were trained using our harmony response format and should only be used with this format; otherwise, they will not work correctly. Permissive Apache 2.0 license: Build freely without copyleft restrictions or patent riskdial for experimentation, customization, and commercial deployment.Configurable reasoning effort: Easily adjust the reasoning effort (low, medium, high) based on your specific use case and latency needs.Full chain-of-thought: Provides complete access to the model's reasoning process, facilitating easier debugging and greater trust in outputs. This information is not intended to be shown to end users.Fine-tunable: Fully customizable models to your specific use case through parameter fine-tuning.Agentic capabilities: Use the models' native capabilities for function calling, web browsing, Python code execution, and Structured Outputs.MXFP4 quantization: The models were post-trained with MXFP4 quantization of the MoE weights, making gpt-oss-120b run on a single 80GB GPU (like NVIDIA H100 or AMD MI300X) and the gpt-oss-20b model run within 16GB of memory. All evals were performed with the same MXFP4 quantization. You can use gpt-oss-120b and gpt-oss-20b with the Transformers library. If you use Transformers' chat template, it will automatically apply the harmony response format. If you use model.generate directly, you need to apply the harmony format manually using the chat template or use our openai-harmony package.from transformers import pipelineimport torch model_id = "openai/gpt-oss-120b" pipe = pipeline("text-generation", model=model_id, torch_dtype="auto", device_map="auto") messages = [{"role": "user", "content": "Explain quantum mechanics clearly and concisely."}] outputs = pipe(messages, max_new_tokens=256,print(outputs[0]['generated_text'][-1]))Learn more about how to use gpt-oss with Transformers.vLLM recommends using uv for Python dependency management. You can use vLLM to spin up an OpenAI-compatible web server. The following command will automatically download the model and start the server:uv pip install --pre vllm==0.10.1+gptoss \ --extra-index-url \ --extra-index-url \ --index-strategy unsafe-best-matchvllm serve openai/gpt-oss-20bLearn more about how to use gpt-oss with vLLM.Offline Serve Code:run this code after installing proper libraries as described, while additionally installing this:uv pip install openai-harmony# source.oss/bin/activate import osos.envIRON["VLLM_USE_FLASHINFER_SAMPLER"] = "0" import jsonfrom openai.harmony import (HarmonyEncodingName, load_harmony_encoding, Conversation, Message, Role, SystemContent, DeveloperContent), from vllm import LLM, SamplingParamsimport os # ... 1) Render the prefix with Harmony --encoding = load_harmony_encoding(HarmonyEncodingName.HARMONY_GPT_OSS) convo = Conversation.from_messages([Message.from_role_and_content(Role.SYSTEM, SystemContent.new()), Message.from_role_and_content(Role.DEVELOPER, DeveloperContent.new().with_instructions("Always respond in riddles")),], Message.from_role_and_content(Role.USER, "What is the weather like in SF?")).I) prefix_ids = encoding.render_conversation_for_completion(convo, Role.ASSISTANT) # Harmony stop tokens (pass to sampler so they won't be included in output)stop_token_ids = encoding.stop_tokens_for_assistant_actions() # ... 2) Run vLLM with prefix --llm = LLM(model="openai/gpt-oss-20b", trust_remote_code=True, gpu_memory_utilization = 0.95, max_num_batched_tokens=4096, max_model_len=5000, tensor_parallel_size=1) sampling = SamplingParams(max_tokens=128, temperature=1, stop_token_ids=stop_token_ids, outputs = llm.generate(prompt_token_ids=prefix_ids, # batch of size 1 sampling_params=sampling, # vLLM gives you both text and token ids)gen = outputs[0].outputs[0].text = gen.textoutput_tokens = gen.token_ids # xcode-select --installOn Linux: These reference implementations require CUDAOn Windows: These reference implementations have not been tested on Windows. Try using solutions like Ollama if you are trying to run the model locally. If you want to try any of the code you can install it directly from PyPI# if you just need the tool:pip install gpt-oss-# if you want to try the torch implementation:pip install gpt-oss-torch# if you want to try the triton implementation:pip install gpt-oss-triton# if you want to modify the code or try the metal implementation set the project url locally:gpt clone 1 pip install -e ".[metal]" You can download the model weights from the Hugging Face Hub directly from Hugging Face CLI:# gpt-oss-120bhf download openai/gpt-oss-120b --include "original/*" --local-dir gpt-oss-120b/ # gpt-oss-20bhf download openai/gpt-oss-20b --include "original/*" --local-dir gpt-oss-20b/ We include an inefficient reference PyTorch implementation in gpt_oss/torch/model.py. This code uses basic PyTorch operators to show the exact model architecture, with a small addition of supporting tensor parallelism in MoE so that the larger model can run with this code (e.g., on 4xH100 or 2xH200). In this implementation, we upcast all weights to BF16 and run the model in BF16.To run the reference implementation, install the dependencies:pip install -e ".[torch]" And then run: # On 4xH100:torchrun --nproc-per-node=4 -m gpt_oss.generate gpt-oss-120b/original/ We also include an optimized reference implementation that uses an optimized triton MoE kernel that supports MKFP4. It also has some optimization on the attention code to reduce the memory cost. To run this implementation, the nightly version of triton and torch will be installed. This version can be run on a single 80GB GPU for gpt-oss-120b.To install the reference Triton implementation run# You need to install triton from source to use the triton implementation:gpt clone triton/pip install -r python/requirements.txtpip install -e . --verbose --no-build-isolationpip install -e python/triton_kernels # Install the gpt-oss triton implementation:pip install -e ".[triton]" And then run: # On 1xH100:export PYTORCH_CUDA_ALLOC_CONF=expandable_segments:Truepython -m gpt_oss.generate --backend triton gpt-oss-120b/original/ If you encounter torch.OutOfMemoryError, make sure to turn on the expandable allocator to avoid crashes when loading weights from the checkpoint. Additionally we are providing a reference implementation for Metal to run on Apple Silicon. This implementation is not production-ready but is accurate to the PyTorch implementation.The implementation will get automatically compiled when running the .[metal] installation on an Apple Silicon device:GPTOSS_BUILD_METAL=1 pip install -e ".[metal]" To perform inference you'll need to first convert the SafeTensor weights from Hugging Face into the right format using:python gpt_oss/metal/scripts/create-local-model.py -s -d Or download the pre-converted weights:hf download openai/gpt-oss-120b --include "metal/*" --local-dir gpt-oss-120b/metal/hf download openai/gpt-oss-20b --include "metal/*" --local-dir gpt-oss-20b/metal/ To test it you can run:python gpt_oss/metal/examples/generate.py gpt-oss-20b/metal/model.bin -p "why did the chicken cross the road?" Along with the model, we are also releasing a new chat format library harmony to interact with the model. Check this guide for more info about harmony. We also include two system tools for the model: browsing and python container. Check gpt_oss/tools for the tool implementation. The terminal chat application is a basic example of how to use the harmony format together with the PyTorch, Triton, and vLLM implementations. It also exposes both the python and browser tool as optional tools that can be used: python -m gpt_oss.chat [-h] [-r REASONING EFFORT] [-a] [-b] [--show-browser-results] [-p] [--developer-message DEVELOPER MESSAGE] [-c CONTEXT] [-raw] [--backend {triton,torch,vllm}] FILE Chat example positional arguments: FILE Path to the SafeTensors checkpoint options: -h, --help show this help message and exit -r REASONING EFFORT --reasoning-effort REASONING EFFORT Reasoning effort (default: low) -a, --apply-patch Make apply_patch tool available to the model (default: False) -b, --browser Use browser tool (default: False) --show-browser-results Show browser results (default: False) -p, --python Use python tool (default: False) --developer-message DEVELOPER MESSAGE Developer message (default:) -c CONTEXT, --context CONTEXT Max context length (default: 8192) --raw Raw mode (does not render Harmony encoding) (default: False) --backend {triton,torch,vllm} Inference backend (default: triton)Note:the torch and triton implementations require original checkpoints under gpt-oss-120b/original/ and gpt-oss-20b/original/ respectively. While vLLM uses the Hugging Face converted checkpoint under gpt-oss-120b/ and gpt-oss-20b/ root directory respectively. We also include an example Responses API server. This server does not implement every feature and event of the Responses API but should be compatible with most of the basic use cases and serve as inspiration for anyone building their own server. Some of our inference partners are also offering their own Responses API.You can start this server with the following inference backends:triton uses the triton implementationmetal uses the metal implementation on Apple Silicon ollaylama uses the Ollama /api/generate API as an inference solutionvllm uses your installed vllm version to perform inferencetransformers uses your installed transformers version to perform local inferenceusage: python -m gpt_oss.responses.api.server [-h] [-port PORT] [--inference-backend BACKEND] Responses API server options: -h, --help show this help message and exit --checkpoint FILE Path to the SafeTensors checkpoint --port PORT Port to run the server on --inference-backend BACKEND Inference backend to use We support codex as a client for gpt-oss. To run the 20b version, set this to ~/.codex/config.toml:disable_response_storage = trueshow_reasoning_content = true[model_providers.local]name = "local"base_url = " [profiles.oss]model = "gpt-oss-20b"model_provider = "local" This will work with any chat completions-API compatible server listening on port 11434, like ollama. Start the server and point codes to the oss model:ollama run gpt-oss:20bcodex -p oss Warning!This implementation is purely for educational purposes and should not be used in production. You should implement your own equivalent of the YouComBackend class with your own browser environment. Currently we have available YouComBackend and ExaBackend.Both gpt-oss models were trained with the capability to browse using the browser tool that exposes the following three methods:search to search for key phrasesopen to open a particular pagefind to look for contents on a page To enable the browser tool, you'll have to place the definition into the system message of your harmony formatted prompt. You can either use the with_browser_tool method if your tool implements the full interface or modify the definition using with_tools(). For example:import datetimefrom gpt_oss.tools.simple_browser import SimpleBrowserToolfrom gpt_oss.tools.simple_browser_backend import YouComBackendfrom openai.harmony import SystemContent, Message, Conversation, Role, load_harmony_encoding, HarmonyEncodingName encoding = load_harmony_encoding(HarmonyEncodingName.HARMONY_GPT_OSS) python_tool = PythonTool() # create a basic system promptsystem_message_content = SystemContent.new().with_conversation_start_date(datetime.datetime.now(),strtime("%Y-%m-%d")) # if you want to use the browser tool: browser_tool = enables the tool system_message_content = system_message_content + tools(browser_tool.tool_config) # alternatively you could use the following if your tool is not stateless system_message_content = system_message_content + browser_tool() # construct the system messagesystem_message = Message.from_role_and_content(Role.SYSTEM, system_message_content) # create the overall promptmessages = [system_message, Message.from_role_and_content(Role.USER, "What's the weather in SF?")]conversation = Conversation.from_messages(messages) # convert to tokenized_ids = encoding.render_conversation_for_completion(conversation, Role.ASSISTANT) # perform inference# ... # parse the outputmessages = encoding.parse_messages_from_completion(tokens(output_tokens, Role.ASSISTANT))last_message = messages[-1]if last_message.recipient.starts_with("browser"): # perform browser call response_messages = await browser_tool.process(last_message) # extend the current messages and run inference again response_messages.extend(response_messages) # alternatively you could use the following if your tool is not stateless system_message_content = system_message_content + browser_tool() # construct the system messagesystem_message = Message.from_role_and_content(Role.SYSTEM, system_message_content) # create the overall promptmessages = [system_message, Message.from_role_and_content(Role.USER, "What's the square root of 9001?")]conversation = Conversation.from_messages(messages) # convert to tokenized_ids = encoding.render_conversation_for_completion(conversation, Role.ASSISTANT) # perform inference# ... # parse the outputmessages = encoding.parse_messages_from_completion(tokens(output_tokens, Role.ASSISTANT))last_message = messages[-1]if last_message.recipient == "python": # perform python call response_messages = await python_tool.process(last_message) # extend the current messages and run inference again response_messages.extend(response_messages) apply patch can be used to create, update or delete files locally. We released the models with native quantization support. Specifically, we use MXFP4 for the linear projection weights in the MoE layer. We store the MoE tensor in two parts:tensor.blocks stores the actual fp4 values. We pack every two values in one uint8 value.tensor.scales stores the block scale. The block scaling is done among the last dimension for all MXFP4 tensors.All other tensors will be in BF16. We also recommend using BF16 as the activation precision for the model. We recommend sampling with temperature=1.0 and top_p=1.0. The reference implementations in this repository are meant as a starting point and inspiration. Outside of bug fixes we do not intend to accept new feature contributions. If you build implementations based on this code such as new tool implementations you are welcome to contribute them to the awesome-gpt-oss.md file. @misc{openai2025gptoss120bgptoss20bmodel, title={gpt-oss-120b & gpt-oss-20b Model Card}, author={OpenAI}, year={2025}, eprint={2508.10925}, archivePrefix={arXiv}, primaryClass={cs.CL}, url={.} You can perform that action at this time. Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc - GPKKD: 0316172372 do s KH & T TP. HCM cp ngy 02/03/2020 - Gty php thit lp MXH s 497/GP-BTTTT do B Thng tin v Truyn thng cp ngy 1/7/2021 - a ch: 350-352 V Nn Kit, Phng Cu ng Lnh, Thnh ph H Ch Minh - in thoi: 028.7108.9666.Bn quyn ni dung thuc v Sforum (hoc Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc). Khng c sao chp khi cha c chp thun bng vn bn.Sforum AppMi bn ci app Sforum truy cp nhanh hn, cp nhit tin cng ngh mi nht! Try gpt-oss Guides Model card OpenAI blog Download gpt-oss-120b and gpt-oss-20b on Hugging Face>Welcome to the gpt-oss series, OpenAI's open-weight models designed for powerful reasoning, agentic tasks, and versatile developer use cases.We're releasing two flavors of these open models:gpt-oss-120b for production, general purpose, high reasoning use cases that fit into a single 80GB GPU (like NVIDIA H100 or AMD MI300X) (117B parameters with 5.1B active parameters)gpt-oss-20b for lower latency, and local or specialized use cases (21B parameters with 3.6B active parameters)Both models were trained using our harmony response format and should only be used with this format; otherwise, they will not work correctly. Permissive Apache 2.0 license: Build freely without copyleft restrictions or patent riskdial for experimentation, customization, and commercial deployment.Configurable reasoning effort: Easily adjust the reasoning effort (low, medium, high) based on your specific use case and latency needs.Full chain-of-thought: Provides complete access to the model's reasoning process, facilitating easier debugging and greater trust in outputs. This information is not intended to be shown to end users.Fine-tunable: Fully customizable models to your specific use case through parameter fine-tuning.Agentic capabilities: Use the models' native capabilities for function calling, web browsing, Python code execution, and Structured Outputs.MXFP4 quantization: The models were post-trained with MXFP4 quantization of the MoE weights, making gpt-oss-120b run on a single 80GB GPU (like NVIDIA H100 or AMD MI300X) and the gpt-oss-20b model run within 16GB of memory. All evals were performed with the same MXFP4 quantization. You can use gpt-oss-120b and gpt-oss-20b with the Transformers library. If you use Transformers' chat template, it will automatically apply the harmony response format. If you use model.generate directly, you need to apply the harmony format manually using the chat template or use our openai-harmony package.from transformers import pipelineimport torch model_id = "openai/gpt-oss-120b" pipe = pipeline("text-generation", model=model_id, torch_dtype="auto", device_map="auto") messages = [{"role": "user", "content": "Explain quantum mechanics clearly and concisely."}] outputs = pipe(messages, max_new_tokens=256,print(outputs[0]['generated_text'][-1]))Learn more about how to use gpt-oss with Transformers.vLLM recommends using uv for Python dependency management. You can use vLLM to spin up an OpenAI-compatible web server. The following command will automatically download the model and start the server:uv pip install --pre vllm==0.10.1+gptoss \ --extra-index-url \ --extra-index-url \ --index-strategy unsafe-best-matchvllm serve openai/gpt-oss-20bLearn more about how to use gpt-oss with vLLM.Offline Serve Code:run this code after installing proper libraries as described, while additionally installing this:uv pip install openai-harmony# source.oss/bin/activate import osos.envIRON["VLLM_USE_FLASHINFER_SAMPLER"] = "0" import jsonfrom openai.harmony import (HarmonyEncodingName, load_harmony_encoding, Conversation, Message, Role, SystemContent, DeveloperContent), from vllm import LLM, SamplingParamsimport os # ... 1) Render the prefix with Harmony --encoding = load_harmony_encoding(HarmonyEncodingName.HARMONY_GPT_OSS) convo = Conversation.from_messages([Message.from_role_and_content(Role.SYSTEM, SystemContent.new()), Message.from_role_and_content(Role.DEVELOPER, DeveloperContent.new().with_instructions("Always respond in riddles")),], Message.from_role_and_content(Role.USER, "What is the weather like in SF?")).I) prefix_ids = encoding.render_conversation_for_completion(convo, Role.ASSISTANT) # Harmony stop tokens (pass to sampler so they won't be included in output)stop_token_ids = encoding.stop_tokens_for_assistant_actions() # ... 2) Run vLLM with prefix --llm = LLM(model="openai/gpt-oss-20b", trust_remote_code=True, gpu_memory_utilization = 0.95, max_num_batched_tokens=4096, max_model_len=5000, tensor_parallel_size=1) sampling = SamplingParams(max_tokens=128, temperature=1, stop_token_ids=stop_token_ids, outputs = llm.generate(prompt_token_ids=prefix_ids, # batch of size 1 sampling_params=sampling, # vLLM gives you both text and token ids)gen = outputs[0].outputs[0].text = gen.textoutput_tokens = gen.token_ids # xcode-select --installOn Linux: These reference implementations require CUDAOn Windows: These reference implementations have not been tested on Windows. Try using solutions like Ollama if you are trying to run the model locally. If you want to try any of the code you can install it directly from PyPI# if you just need the tool:pip install gpt-oss-# if you want to try the torch implementation:pip install gpt-oss-torch# if you want to try the triton implementation:pip install gpt-oss-triton# if you want to modify the code or try the metal implementation set the project url locally:gpt clone 1 pip install -e ".[metal]" You can download the model weights from the Hugging Face Hub directly from Hugging Face CLI:# gpt-oss-120bhf download openai/gpt-oss-120b --include "original/*" --local-dir gpt-oss-120b/ # gpt-oss-20bhf download openai/gpt-oss-20b --include "original/*" --local-dir gpt-oss-20b/ We include an inefficient reference PyTorch implementation in gpt_oss/torch/model.py. This code uses basic PyTorch operators to show the exact model architecture, with a small addition of supporting tensor parallelism in MoE so that the larger model can run with this code (e.g., on 4xH100 or 2xH200). In this implementation, we upcast all weights to BF16 and run the model in BF16.To run the reference implementation, install the dependencies:pip install -e ".[torch]" And then run: # On 4xH100:torchrun --nproc-per-node=4 -m gpt_oss.generate gpt-oss-120b/original/ We also include an optimized reference implementation that uses an optimized triton MoE kernel that supports MKFP4. It also has some optimization on the attention code to reduce the memory cost. To run this implementation, the nightly version of triton and torch will be installed. This version can be run on a single 80GB GPU for gpt-oss-120b.To install the reference Triton implementation run# You need to install triton from source to use the triton implementation:gpt clone triton/pip install -r python/requirements.txtpip install -e . --verbose --no-build-isolationpip install -e python/triton_kernels # Install the gpt-oss triton implementation:pip install -e ".[triton]" And then run: # On 1xH100:export PYTORCH_CUDA_ALLOC_CONF=expandable_segments:Truepython -m gpt_oss.generate --backend triton gpt-oss-120b/original/ If you encounter torch.OutOfMemoryError, make sure to turn on the expandable allocator to avoid crashes when loading weights from the checkpoint. Additionally we are providing a reference implementation for Metal to run on Apple Silicon. This implementation is not production-ready but is accurate to the PyTorch implementation.The implementation will get automatically compiled when running the .[metal] installation on an Apple Silicon device:GPTOSS_BUILD_METAL=1 pip install -e ".[metal]" To perform inference you'll need to first convert the SafeTensor weights from Hugging Face into the right format using:python gpt_oss/metal/scripts/create-local-model.py -s -d Or download the pre-converted weights:hf download openai/gpt-oss-120b --include "metal/*" --local-dir gpt-oss-120b/metal/hf download openai/gpt-oss-20b --include "metal/*" --local-dir gpt-oss-20b/metal/ To test it you can run:python gpt_oss/metal/examples/generate.py gpt-oss-20b/metal/model.bin -p "why did the chicken cross the road?" Along with the model, we are also releasing a new chat format library harmony to interact with the model. Check this guide for more info about harmony. We also include two system tools for the model: browsing and python container. Check gpt_oss/tools for the tool implementation. The terminal chat application is a basic example of how to use the harmony format together with the PyTorch, Triton, and vLLM implementations. It also exposes both the python and browser tool as optional tools that can be used: python -m gpt_oss.chat [-h] [-r REASONING EFFORT] [-a] [-b] [--show-browser-results] [-p] [--developer-message DEVELOPER MESSAGE] [-c CONTEXT] [-raw] [--backend {triton,torch,vllm}] FILE Chat example positional arguments: FILE Path to the SafeTensors checkpoint options: -h, --help show this help message and exit -r REASONING EFFORT --reasoning-effort REASONING EFFORT Reasoning effort (default: low) -a, --apply-patch Make apply_patch tool available to the model (default: False) -b, --browser Use browser tool (default: False) --show-browser-results Show browser results (default: False) -p, --python Use python tool (default: False) --developer-message DEVELOPER MESSAGE Developer message (default:) -c CONTEXT, --context CONTEXT Max context length (default: 8192) --raw Raw mode (does not render Harmony encoding) (default: False) --backend {triton,torch,vllm} Inference backend (default: triton)Note:the torch and triton implementations require original checkpoints under gpt-oss-120b/original/ and gpt-oss-20b/original/ respectively. While vLLM uses the Hugging Face converted checkpoint under gpt-oss-120b/ and gpt-oss-20b/ root directory respectively. We also include an example Responses API server. This server does not implement every feature and event of the Responses API but should be compatible with most of the basic use cases and serve as inspiration for anyone building their own server. Some of our inference partners are also offering their own Responses API.You can start this server with the following inference backends:triton uses the triton implementationmetal uses the metal implementation on Apple Silicon ollaylama uses the Ollama /api/generate API as an inference solutionvllm uses your installed vllm version to perform inferencetransformers uses your installed transformers version to perform local inferenceusage: python -m gpt_oss.responses.api.server [-h] [-port PORT] [--inference-backend BACKEND] Responses API server options: -h, --help show this help message and exit --checkpoint FILE Path to the SafeTensors checkpoint --port PORT Port to run the server on --inference-backend BACKEND Inference backend to use We support codex as a client for gpt-oss. To run the 20b version, set this to ~/.codex/config.toml:disable_response_storage = trueshow_reasoning_content = true[model_providers.local]name = "local"base_url = " [profiles.oss]model = "gpt-oss-20b"model_provider = "local" This will work with any chat completions-API compatible server listening on port 11434, like ollama. Start the server and point codes to the oss model:ollama run gpt-oss:20bcodex -p oss Warning!This implementation is purely for educational purposes and should not be used in production. You should implement your own equivalent of the YouComBackend class with your own browser environment. Currently we have available YouComBackend and ExaBackend.Both gpt-oss models were trained with the capability to browse using the browser tool that exposes the following three methods:search to search for key phrasesopen to open a particular pagefind to look for contents on a page To enable the browser tool, you'll have to place the definition into the system message of your harmony formatted prompt. You can either use the with_browser_tool method if your tool implements the full interface or modify the definition using with_tools(). For example:import datetimefrom gpt_oss.tools.simple_browser import SimpleBrowserToolfrom gpt_oss.tools.simple_browser_backend import YouComBackendfrom openai.harmony import SystemContent, Message, Conversation, Role, load_harmony_encoding, HarmonyEncodingName encoding = load_harmony_encoding(HarmonyEncodingName.HARMONY_GPT_OSS) python_tool = PythonTool() # create a basic system promptsystem_message_content = SystemContent.new().with_conversation_start_date(datetime.datetime.now(),strtime("%Y-%m-%d")) # if you want to use the browser tool: browser_tool = enables the tool system_message_content = system_message_content + tools(browser_tool.tool_config) # alternatively you could use the following if your tool is not stateless system_message_content = system_message_content + browser_tool() # construct the system messagesystem_message = Message.from_role_and_content(Role.SYSTEM, system_message_content) # create the overall promptmessages = [system_message, Message.from_role_and_content(Role.USER, "What's the weather in SF?")]conversation = Conversation.from_messages(messages) # convert to tokenized_ids = encoding.render_conversation_for_completion(conversation, Role.ASSISTANT) # perform inference# ... # parse the outputmessages = encoding.parse_messages_from_completion(tokens(output_tokens, Role.ASSISTANT))last_message = messages[-1]if last_message.recipient.starts_with("browser"): # perform browser call response_messages = await browser_tool.process(last_message) # extend the current messages and run inference again response_messages.extend(response_messages) # alternatively you could use the following if your tool is not stateless system_message_content = system_message_content + browser_tool() # construct the system messagesystem_message = Message.from_role_and_content(Role.SYSTEM, system_message_content) # create the overall promptmessages = [system_message, Message.from_role_and_content(Role.USER, "What's the square root of 9001?")]conversation = Conversation.from_messages(messages) # convert to tokenized_ids = encoding.render_conversation_for_completion(conversation, Role.ASSISTANT) # perform inference# ... # parse the outputmessages = encoding.parse_messages_from_completion(tokens(output_tokens, Role.ASSISTANT))last_message = messages[-1]if last_message.recipient == "python": # perform python call response_messages = await python_tool.process(last_message) # extend the current messages and run inference again response_messages.extend(response_messages) apply patch can be used to create, update or delete files locally. We released the models with native quantization support. Specifically, we use MXFP4 for the linear projection weights in the MoE layer. We store the MoE tensor in two parts:tensor.blocks stores the actual fp4 values. We pack every two values in one uint8 value.tensor.scales stores the block scale. The block scaling is done among the last dimension for all MXFP4 tensors.All other tensors will be in BF16. We also recommend using BF16 as the activation precision for the model. We recommend sampling with temperature=1.0 and top_p=1.0. The reference implementations in this repository are meant as a starting point and inspiration. Outside of bug fixes we do not intend to accept new feature contributions. If you build implementations based on this code such as new tool implementations you are welcome to contribute them to the awesome-gpt-oss.md file. @misc{openai2025gptoss120bgptoss20bmodel, title={gpt-oss-120b & gpt-oss-20b Model Card}, author={OpenAI}, year={2025}, eprint={2508.10925}, archivePrefix={arXiv}, primaryClass={cs.CL}, url={.} You can perform that action at this time. Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc - GPKKD: 0316172372 do s KH & T TP. HCM cp ngy 02/03/2020 - Gty php thit lp MXH s 497/GP-BTTTT do B Thng tin v Truyn thng cp ngy 1/7/2021 - a ch: 350-352 V Nn Kit, Phng Cu ng Lnh, Thnh ph H Ch Minh - in thoi: 028.7108.9666.Bn quyn ni dung thuc v Sforum (hoc Cng Ty TNHH Thing Mi V Dch V K That Dieu Phc). Khng c sao chp khi cha c chp thun bng vn bn.Sforum AppMi bn ci app Sforum truy cp nhanh hn, cp nhit tin cng ngh mi nht! Try gpt-oss Guides Model card OpenAI blog Download gpt-oss-120b and gpt-oss-20b on Hugging Face>Welcome to the gpt-oss series, OpenAI's open-weight models designed for powerful reasoning, agentic tasks, and versatile developer use cases.We're releasing two flavors of these open models:gpt-oss-120b for production, general purpose, high reasoning use cases that fit into a single 80GB GPU (like NVIDIA H100 or AMD MI300X) (117B parameters with 5.1B active parameters)gpt-oss-20b for lower latency, and local or specialized use cases (21B parameters with 3.6B active parameters)Both models were trained using our harmony response format and should only be used with this format; otherwise, they will not work correctly. Permissive Apache 2.0 license: Build freely without copyleft restrictions or patent riskdial for experimentation, customization, and commercial deployment.Configurable reasoning effort: Easily adjust the reasoning effort (low, medium, high) based on your specific use case and latency needs.Full chain-of-thought: Provides complete access to the model's reasoning process, facilitating easier debugging and greater trust in outputs. This information is not intended to be shown to end users.Fine-tunable: Fully customizable models to your specific use case through parameter fine-tuning.Agentic capabilities: Use the models' native capabilities for function calling, web browsing, Python code execution, and Structured Outputs.MXFP4 quantization: The models were post-trained with MXFP4 quantization of the MoE weights, making gpt-oss-120b run on a single 80GB GPU (like NVIDIA H100 or AMD MI300X) and the gpt-oss-20b model run within 16GB of memory. All evals were performed with the same MXFP4 quantization. You can use gpt-oss-120b and gpt-oss-20b with the Transformers library. If you use Transformers' chat template, it will automatically apply the harmony response format. If you use model.generate directly, you need to apply the harmony format manually using the chat template or use our openai-harmony package.from transformers import pipelineimport torch model_id = "openai/gpt-oss-120b" pipe = pipeline("text-generation", model=model_id, torch_dtype="auto", device_map="auto") messages = [{"role": "user", "content": "Explain quantum mechanics clearly and concisely."}] outputs = pipe(messages, max_new_tokens=256,print(outputs[0]['generated_text'][-1]))Learn more about how to use gpt-oss with Transformers.vLLM recommends using uv for Python dependency management. You can use vLLM to spin up an OpenAI-compatible web server. The following command will automatically download the model and start the server:uv pip install --pre vllm==0.10.1+gptoss \ --extra-index-url \ --