

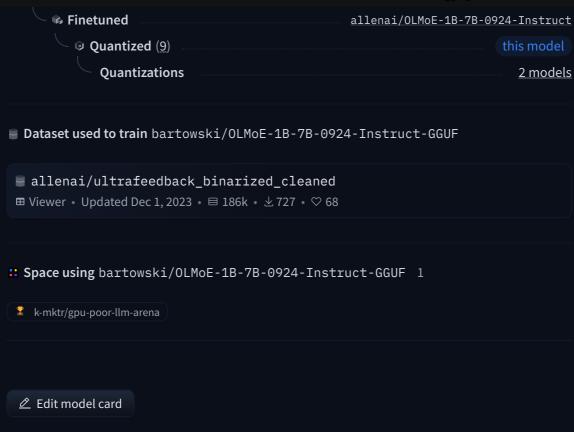
Inference Examples ①

Text Generation

This model does not have enough activity to be deployed to Inference API (serverless) yet. Increase its social visibility and check back later, or deploy to <u>Inference Endpoints (dedicated)</u> instead.

Model tree for bartowski/OLMoE-1B-7B-0924-Instruct-GGUF

Base model allenai/OLMoE-1B-7B-0924
Finetuned



■ Ø Llamacpp imatrix Quantizations of OLMoE-1B-7B-0924-Instruct

Using <u>llama.cpp</u> release <u>b3772</u> for quantization.

Original model: https://huggingface.co/allenai/OLMoE-1B-7B-0924-Instruct

All quants made using imatrix option with dataset from here

Run them in LM Studio

Prompt format

```
<|endoftext|><|system|>
{system_prompt}
<|user|>
{prompt}
<|assistant|>
```

Our Download a file (not the whole branch) from below:

Filename	Quant type	File Size	Split	Description
OLMoE-1B-7B-0924- Instruct-f16.gguf	f16	13.84GB	false	Full F16 weights.
OLMoE-1B-7B-0924- Instruct-Q8 0.gguf	Q8_0	7.36GB	false	Extremely high quality, generally unneeded but max available quant.
OLMoE-1B-7B-0924- Instruct-Q6 K L.gguf	Q6_K_L	5.73GB	false	Uses Q8_0 for embed and output weights. Very high quality, near perfect, <i>recommended</i> .
OLMoE-1B-7B-0924- Instruct-Q6 K.gguf	Q6_K	5.68GB	false	Very high quality, near perfect, recommended.
OLMoE-1B-7B-0924- Instruct-Q5 K L.gguf	Q5_K_L	4.99GB	false	Uses Q8_0 for embed and output weights. High quality, <i>recommended</i> .
OLMoE-1B-7B-0924- Instruct-Q5 K M.gguf	Q5_K_M	4.93GB	false	High quality, recommended.
OLMoE-1B-7B-0924- Instruct-Q5 K S.gguf	Q5_K_S	4.78GB	false	High quality, recommended.
OLMoE-1B-7B-0924- Instruct-Q4 K L.gguf	Q4_K_L	4.29GB	false	Uses Q8_0 for embed and output weights. Good quality, <i>recommended</i> .
OLMoE-1B-7B-0924- Instruct-Q4 K M.gguf	Q4_K_M	4.21GB	false	Good quality, default size for must use cases, <i>recommended</i> .
OLMoE-1B-7B-0924- Instruct-Q4 K S.gguf	Q4_K_S	3.96GB	false	Slightly lower quality with more space savings, <i>recommended</i> .
OLMoE-1B-7B-0924- Instruct-Q4 0.gguf	Q4_0	3.94GB	false	Legacy format, generally not worth using over similarly sized formats
OLMoE-1B-7B-0924- Instruct- Q4 0 8 8.gguf	Q4_0_8_8	3.93GB	false	Optimized for ARM inference. Requires 'sve' support (see link below).
OLMoE-1B-7B-0924- Instruct- Q4 0 4 8.gguf	Q4_0_4_8	3.93GB	false	Optimized for ARM inference. Requires 'i8mm' support (see link below).
OLMoE-1B-7B-0924- Instruct-	Q4_0_4_4	3.93GB	false	Optimized for ARM inference. Should work well on all ARM chips, pick this if

Filename	Quant type	File Size	Split	Description
<u>Q4 0 4 4.gguf</u>				you're unsure.
OLMoE-1B-7B-0924- Instruct-IQ4_XS.gguf	IQ4_XS	3.72GB	false	Decent quality, smaller than Q4_K_S with similar performance, recommended.
OLMoE-1B-7B-0924- Instruct-Q3 K XL.gguf	Q3_K_XL	3.70GB	false	Uses Q8_0 for embed and output weights. Lower quality but usable, good for low RAM availability.
OLMoE-1B-7B-0924- Instruct-Q3 K L.gguf	Q3_K_L	3.61GB	false	Lower quality but usable, good for low RAM availability.
OLMoE-1B-7B-0924- Instruct-Q3 K M.gguf	Q3_K_M	3.34GB	false	Low quality.
OLMoE-1B-7B-0924- Instruct-IQ3 M.gguf	IQ3_M	3.08GB	false	Medium-low quality, new method with decent performance comparable to Q3_K_M.
OLMoE-1B-7B-0924- Instruct-Q3 K S.gguf	Q3_K_S	3.02GB	false	Low quality, not recommended.
OLMoE-1B-7B-0924- Instruct-IQ3_XS.gguf	IQ3_XS	2.87GB	false	Lower quality, new method with decent performance, slightly better than Q3_K_S.
OLMoE-1B-7B-0924- Instruct-Q2 K L.gguf	Q2_K_L	2.66GB	false	Uses Q8_0 for embed and output weights. Very low quality but surprisingly usable.
OLMoE-1B-7B-0924- Instruct-Q2 K.gguf	Q2_K	2.56GB	false	Very low quality but surprisingly usable.
OLMoE-1B-7B-0924- Instruct-IQ2 M.gguf	IQ2_M	2.33GB	false	Relatively low quality, uses SOTA techniques to be surprisingly usable.

⊘ Embed/output weights

Some of these quants (Q3_K_XL, Q4_K_L etc) are the standard quantization method with the embeddings and output weights quantized to Q8_0 instead of what they would normally default to.

Some say that this improves the quality, others don't notice any difference. If you use these models PLEASE COMMENT with your findings. I would like feedback that these are actually used and useful so I don't keep uploading quants no one is using.

Thanks!

Downloading using huggingface-cli

First, make sure you have hugginface-cli installed:

```
pip install -U "huggingface_hub[cli]"
```

Then, you can target the specific file you want:

```
huggingface-cli download bartowski/OLMoE-1B-7B-0924-Instruct-GGUF --inc
```

If the model is bigger than 50GB, it will have been split into multiple files. In order to download them all to a local folder, run:

```
huggingface-cli download bartowski/OLMoE-1B-7B-0924-Instruct-GGUF --inc
```

You can either specify a new local-dir (OLMoE-1B-7B-0924-Instruct-Q8_0) or download them all in place (./)

⊘ Q4_0_X_X

These are *NOT* for Metal (Apple) offloading, only ARM chips.

If you're using an ARM chip, the Q4_0_X_X quants will have a substantial speedup. Check out Q4_0_4_4 speed comparisons <u>on the original pull request</u>

To check which one would work best for your ARM chip, you can check <u>AArch64 SoC</u> <u>features</u> (thanks EloyOn!).

Which file should I choose?

A great write up with charts showing various performances is provided by Artefact2 here

The first thing to figure out is how big a model you can run. To do this, you'll need to figure out how much RAM and/or VRAM you have.

If you want your model running as FAST as possible, you'll want to fit the whole thing on your GPU's VRAM. Aim for a quant with a file size 1-2GB smaller than your GPU's total VRAM.

If you want the absolute maximum quality, add both your system RAM and your GPU's VRAM together, then similarly grab a quant with a file size 1-2GB Smaller than that total.

Next, you'll need to decide if you want to use an 'I-quant' or a 'K-quant'.

If you don't want to think too much, grab one of the K-quants. These are in format 'QX_K_X', like Q5_K_M.

If you want to get more into the weeds, you can check out this extremely useful feature chart:

<u>llama.cpp feature matrix</u>

But basically, if you're aiming for below Q4, and you're running cuBLAS (Nvidia) or rocBLAS (AMD), you should look towards the I-quants. These are in format IQX_X, like IQ3_M. These are newer and offer better performance for their size.

These I-quants can also be used on CPU and Apple Metal, but will be slower than their K-quant equivalent, so speed vs performance is a tradeoff you'll have to decide.

The I-quants are *not* compatible with Vulcan, which is also AMD, so if you have an AMD card double check if you're using the rocBLAS build or the Vulcan build. At the time of writing this, LM Studio has a preview with ROCm support, and other inference engines have specific builds for ROCm.

Credits

Thank you kalomaze and Dampf for assistance in creating the imatrix calibration dataset

Thank you ZeroWw for the inspiration to experiment with embed/output

Want to support my work? Visit my ko-fi page here: https://ko-fi.com/bartowski



Company

TOS

Privacy

About

Jobs

Website

Models

Datasets

Spaces

Pricing

Docs

© Hugging Face