



Search models, datasets, users...



bartowski/granite-3.0-1b-a400m-instruct-GGUF like 3

Text Generation GGUF language granite-3.0 Eval Results imatrix

conversational License: apache-2.0

Use this model

Model card Files Community

Downloads last month

0

GGUF

Model size 1.33B params Architecture granitemoe

3-bit Q3_K_S Q3_K_M Q3_K_L Q3_K_XL

4-bit IQ4_XS Q4_K_S Q4_0 Q4_0 Q4_0 Q4_K_M Q4_K_L

5-bit Q5_K_S Q5_K_M Q5_K_L

6-bit Q6_K Q6_K_L

8-bit Q8_0

16-bit F16

View +1 file

Inference Examples

Text Generation

Inference API (serverless) has been turned off for this model.

Model tree for bartowski/granite-3.0-1b-a400m-instruct-GGUF

Base model [ibm-granite/granite-3.0-1b-a400m-instruct](#)

Quantized this model

Evaluation results

pass@1 on IFEval 32.390
pass@1 on IFEval 6.170
pass@1 on AGI-Eval 20.350
pass@1 on AGI-Eval 32.000
pass@1 on AGI-Eval 12.210
pass@1 on OBQA 38.400
pass@1 on OBQA 47.550
pass@1 on OBQA 65.590
pass@1 on OBQA 61.170
pass@1 on OBQA 49.110

..... Expand 23 evaluations

View on Papers With Code

Edit model card

Llamacpp imatrix Quantizations of granite-3.0-1b-a400m-instruct

Using llama.cpp release b3930 for quantization.

Original model: <https://huggingface.co/ibm-granite/granite-3.0-1b-a400m-instruct>

All quants made using imatrix option with dataset from [here](#)

Run them in [LM Studio](#)

Prompt format

```
<|start_of_role|>system<|end_of_role|>{system_prompt}<|end_of_text|>
<|start_of_role|>user<|end_of_role|>{prompt}<|end_of_text|>
<|start_of_role|>assistant<|end_of_role|>
```

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

Filename	Quant type	File Size	Split	Description
granite-3.0-1b-a400m-instruct-f16.gguf	f16	2.67GB	false	Full F16 weights.
granite-3.0-1b-a400m-instruct-Q8_0.gguf	Q8_0	1.42GB	false	Extremely high quality, generally unneeded but max available quant.
granite-3.0-1b-a400m-instruct-Q6_K_L.gguf	Q6_K_L	1.11GB	false	Uses Q8_0 for embed and output weights. Very high quality, near perfect, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q6_K.gguf	Q6_K	1.10GB	false	Very high quality, near perfect, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q5_K_L.gguf	Q5_K_L	0.97GB	false	Uses Q8_0 for embed and output weights. High quality, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q5_K_M.gguf	Q5_K_M	0.96GB	false	High quality, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q5_K_S.gguf	Q5_K_S	0.93GB	false	High quality, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q4_K_L.gguf	Q4_K_L	0.83GB	false	Uses Q8_0 for embed and output weights. Good quality, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q4_K_M.gguf	Q4_K_M	0.82GB	false	Good quality, default size for must use cases, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q4_K_S.gguf	Q4_K_S	0.77GB	false	Slightly lower quality with more space savings, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q4_0_8_8.gguf	Q4_0_8_8	0.77GB	false	Optimized for ARM inference. Requires 'sve' support (see link below). <i>Don't use on Mac or Windows</i> .
granite-3.0-1b-a400m-instruct-Q4_0_4_8.gguf	Q4_0_4_8	0.77GB	false	Optimized for ARM inference. Requires 'i8mm' support (see link below). <i>Don't use on Mac or Windows</i> .
granite-3.0-1b-a400m-instruct-Q4_0_4_4.gguf	Q4_0_4_4	0.77GB	false	Optimized for ARM inference. Should work well on all ARM chips, pick this if you're unsure. <i>Don't use on Mac or Windows</i> .
granite-3.0-1b-a400m-instruct-Q4_0.gguf	Q4_0	0.77GB	false	Legacy format, generally not worth using over similarly sized formats
granite-3.0-1b-a400m-instruct-IQ4_XS.gguf	IQ4_XS	0.73GB	false	Decent quality, smaller than Q4_K_S with similar performance, <i>recommended</i> .
granite-3.0-1b-a400m-instruct-Q3_K_XL.gguf	Q3_K_XL	0.72GB	false	Uses Q8_0 for embed and output weights. Lower quality but usable, good for low RAM availability.
granite-3.0-1b-a400m-instruct-Q3_K_L.gguf	Q3_K_L	0.71GB	false	Lower quality but usable, good for low RAM availability.
granite-3.0-1b-a400m-instruct-Q3_K_M.gguf	Q3_K_M	0.66GB	false	Low quality.
granite-3.0-1b-a400m-instruct-IQ3_M.gguf	IQ3_M	0.61GB	false	Medium-low quality, new method with decent performance comparable to Q3_K_M.
granite-3.0-1b-a400m-instruct-Q3_K_S.gguf	Q3_K_S	0.60GB	false	Low quality, not recommended.

[🔗](#) 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 huggingface-cli installed:

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

Then, you can target the specific file you want:

```
huggingface-cli download bartowski/granite-3.0-1b-a400m-instruct-GGUF ·
```

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/granite-3.0-1b-a400m-instruct-GGUF ·
```

You can either specify a new local-dir (granite-3.0-1b-a400m-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 [on the original pull request](#)

To check which one would work best for your ARM chip, you can check [AArch64 SoC features](#) (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:

[llama.cpp feature matrix](#)

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 Vulkan, which is also AMD, so if you have an AMD card double check if you're using the rocBLAS build or the Vulkan 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