

```
llama_init_from_model: KV self size = 392.00 MiB, K (q8_0): 136.00 MiB, V (f16): 256.00 MiB
llama_init_from_model: CPU output buffer size = 0.49 MiB
llama_init_from_model: CPU compute buffer size = 296.01 MiB
llama_init_from_model: graph nodes = 1030
llama_init_from_model: graph splits = 514 (with bs=512), 1 (with bs=1)
common_init_from_params: setting dry_penalty_last_n to ctx_size = 4096
common_init_from_params: warming up the model with an empty run - please wait ... (--no-warmup to disable)
main: llama threadpool init, n_threads = 32
```