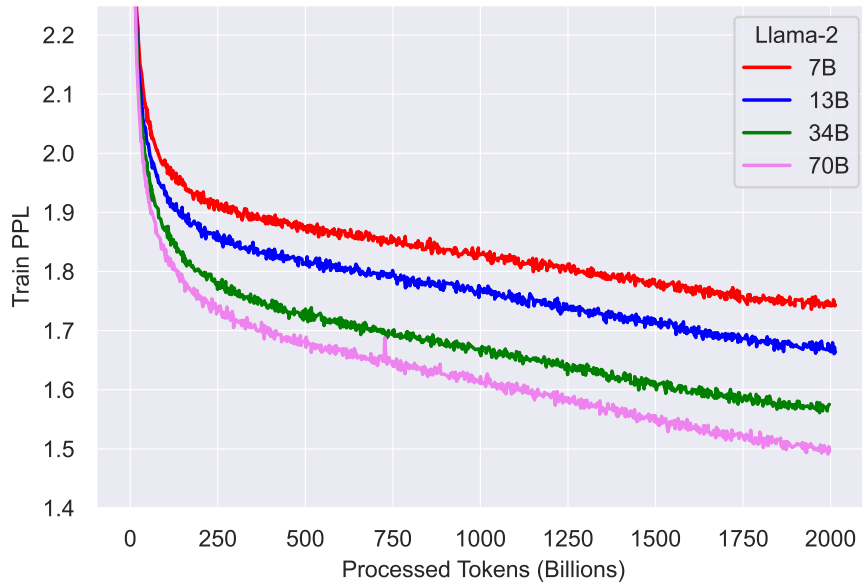


	Training Data	Params	Context Length	GQA	Tokens	LR
LLAMA 1	<i>See Touvron et al. (2023)</i>	7B	2k	$\times$	1.0T	$3.0 \times 10^{-4}$
		13B	2k	$\times$	1.0T	$3.0 \times 10^{-4}$
		33B	2k	$\times$	1.4T	$1.5 \times 10^{-4}$
		65B	2k	$\times$	1.4T	$1.5 \times 10^{-4}$
LLAMA 2	<i>A new mix of publicly available online data</i>	7B	4k	$\times$	2.0T	$3.0 \times 10^{-4}$
		13B	4k	$\times$	2.0T	$3.0 \times 10^{-4}$
		34B	4k	$\checkmark$	2.0T	$1.5 \times 10^{-4}$
		70B	4k	$\checkmark$	2.0T	$1.5 \times 10^{-4}$

**Table 1: LLAMA 2 family of models.** Token counts refer to pretraining data only. All models are trained with a global batch-size of 4M tokens. Bigger models — 34B and 70B — use Grouped-Query Attention (GQA) for improved inference scalability.



**Figure 5: Training Loss for LLAMA 2 models.** We compare the training loss of the LLAMA 2 family of models. We observe that after pretraining on 2T Tokens, the models still did not show any sign of saturation.

**Tokenizer.** We use the same tokenizer as LLAMA 1; it employs a bytepair encoding (BPE) algorithm (Sennrich et al., 2016) using the implementation from SentencePiece (Kudo and Richardson, 2018). As with LLAMA 1, we split all numbers into individual digits and use bytes to decompose unknown UTF-8 characters. The total vocabulary size is 32k tokens.

### 2.2.1 Training Hardware & Carbon Footprint

**Training Hardware.** We pretrained our models on Meta’s Research Super Cluster (RSC) (Lee and Sengupta, 2022) as well as internal production clusters. Both clusters use NVIDIA A100s. There are two key differences between the two clusters, with the first being the type of interconnect available: RSC uses NVIDIA Quantum InfiniBand while our production cluster is equipped with a RoCE (RDMA over converged Ethernet) solution based on commodity ethernet Switches. Both of these solutions interconnect 200 Gbps end-points. The second difference is the per-GPU power consumption cap — RSC uses 400W while our production cluster uses 350W. With this two-cluster setup, we were able to compare the suitability of these different types of interconnect for large scale training. RoCE (which is a more affordable, commercial interconnect network)