

# DeepSeek-V3.2-Exp: Boosting Long-Context Efficiency with DeepSeek Sparse Attention

DeepSeek-AI

research@deepseek.com

## Abstract

We introduce DeepSeek-V3.2-Exp, an experimental sparse-attention model, which equips DeepSeek-V3.1-Terminus with DeepSeek Sparse Attention (DSA) through continued training. With DSA, a fine-grained sparse attention mechanism powered by a lightning indexer, DeepSeek-V3.2-Exp achieves significant efficiency improvements in both training and inference, especially in long-context scenarios. The model checkpoints are available at <https://huggingface.co/deepseek-ai/DeepSeek-V3.2-Exp>.

## 1. Architecture

Compared with DeepSeek-V3.1-Terminus, the last version of DeepSeek-V3.1, the only architectural modification of DeepSeek-V3.2-Exp is the introduction of DeepSeek Sparse Attention (DSA) through continued training.

**Prototype of DSA.** The prototype of DSA primarily consists of two components: a lightning indexer and a fine-grained token selection mechanism.

The **lightning indexer** computes the index score  $I_{t,s}$  between the query token  $\mathbf{h}_t \in \mathbb{R}^d$  and a preceding token  $\mathbf{h}_s \in \mathbb{R}^d$ , determining which tokens to be selected by the query token:

$$I_{t,s} = \sum_{j=1}^{H^I} w_{t,j}^I \cdot \text{ReLU}(\mathbf{q}_{t,j}^I \cdot \mathbf{k}_s^I), \quad (1)$$

where  $H^I$  denotes the number of indexer heads;  $\mathbf{q}_{t,j}^I \in \mathbb{R}^d$  and  $w_{t,j}^I \in \mathbb{R}$  are derived from the query token  $\mathbf{h}_t$ ; and  $\mathbf{k}_s^I \in \mathbb{R}^d$  is derived from the preceding token  $\mathbf{h}_s$ . We choose ReLU as the activation function for throughput consideration. Given that the lightning indexer has a small number of heads and can be implemented in FP8, its computational efficiency is remarkable.

Given the index scores  $\{I_{t,s}\}$  for each query token  $\mathbf{h}_t$ , our **fine-grained token selection mechanism** retrieves only the key-value entries  $\{\mathbf{c}_s\}$  corresponding to the top-k index scores. Then, the attention output  $\mathbf{u}_t$  is computed by applying the attention mechanism between the query token  $\mathbf{h}_t$  and the sparsely selected key-value entries  $\{\mathbf{c}_s\}$ :

$$\mathbf{u}_t = \text{Attn}(\mathbf{h}_t, \{\mathbf{c}_s \mid I_{t,s} \in \text{Top-k}(I_{t,:})\}). \quad (2)$$

---

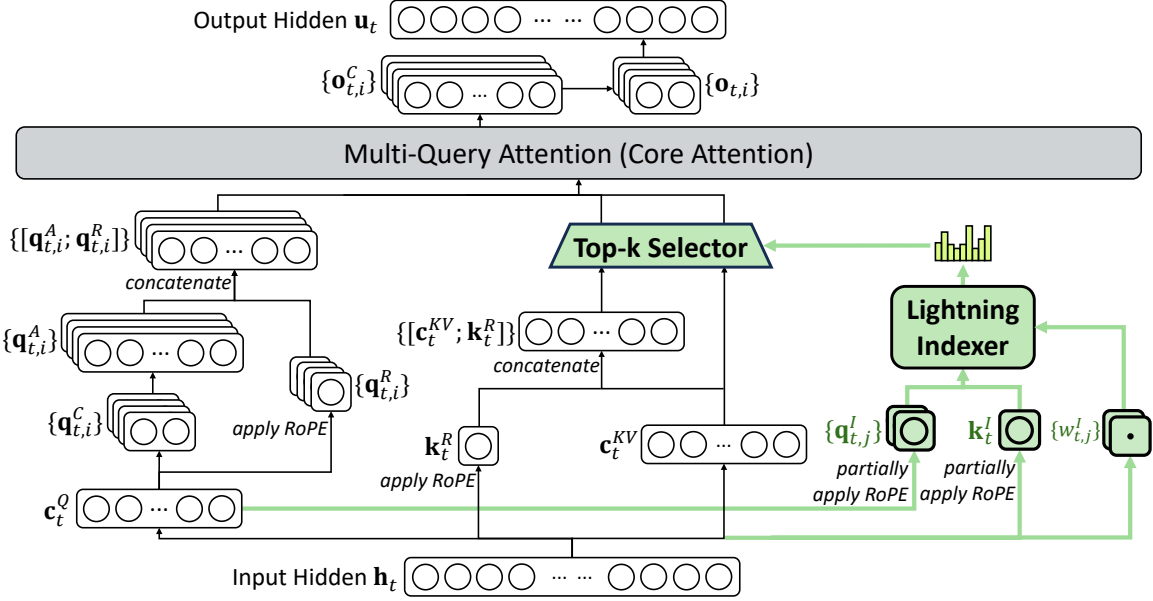


Figure 1 | Attention architecture of DeepSeek-V3.2-Exp, where DSA is instantiated under MLA. The green part illustrates how DSA selects the top-k key-value entries according to the indexer.

**Instantiate DSA Under MLA.** For the consideration of continued training from DeepSeek-V3.1-Terminus, we instantiate DSA based on MLA (DeepSeek-AI, 2024) for DeepSeek-V3.2-Exp. At the kernel level, each key-value entry must be shared across multiple queries for computational efficiency (Yuan et al., 2025). Therefore, we implement DSA based on the MQA (Shazeer, 2019) mode of MLA<sup>1</sup>, where each latent vector (the key-value entry of MLA) will be shared across all query heads of the query token. The DSA architecture based on MLA is illustrated in Figure 1. We also provide an open-source implementation of DeepSeek-V3.2-Exp<sup>2</sup> to specify the details unambiguously.

## 2. Training

Starting from a base checkpoint of DeepSeek-V3.1-Terminus, whose context length has been extended to 128K, we perform continued pre-training followed by post-training to create DeepSeek-V3.2-Exp.

### 2.1. Continued Pre-Training

The continued pre-training of DeepSeek-V3.2-Exp consists of two training stages. For both stages, the distribution of training data is totally aligned with the 128K long context extension data used for DeepSeek-V3.1-Terminus.

**Dense Warm-up Stage.** We first use a short warm-up stage to initialize the lightning indexer. In this stage, we keep dense attention and freeze all model parameters except for the lightning indexer. To align the indexer outputs with the main attention distribution, for the  $t$ -th query token, we first aggregate the main attention scores by summing across all attention heads.

<sup>1</sup>We illustrate the difference between the MQA and MHA modes of MLA in Appendix A.

<sup>2</sup><https://huggingface.co/deepseek-ai/DeepSeek-V3.2-Exp/tree/main/inference>