

DeepSeek-V3.2-Exp技术解读

从稀疏注意力到DSA

DeepSeek Sparse Attention

东川路第一可爱猫猫虫

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



DeepSeek-AI

research@deepseek.com

主要内容

- 我们为什么需要稀疏注意力
- 稀疏注意力 (Openai2019)
- 稀疏注意力的变体 (CVPR2024)
- DeepSeek Sparse Attention (DeepSeek2025)



我们为什么需要稀疏注意力

- Self attention的时间复杂度

对序列中任意两个token要计算attention得分
时间复杂度

$$O(N^2 D)$$

空间复杂度

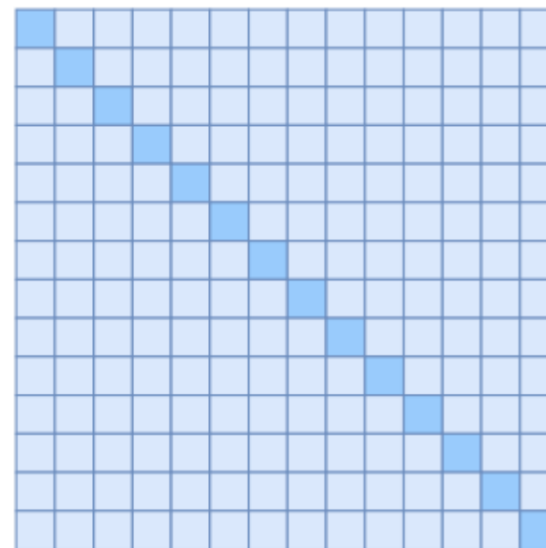
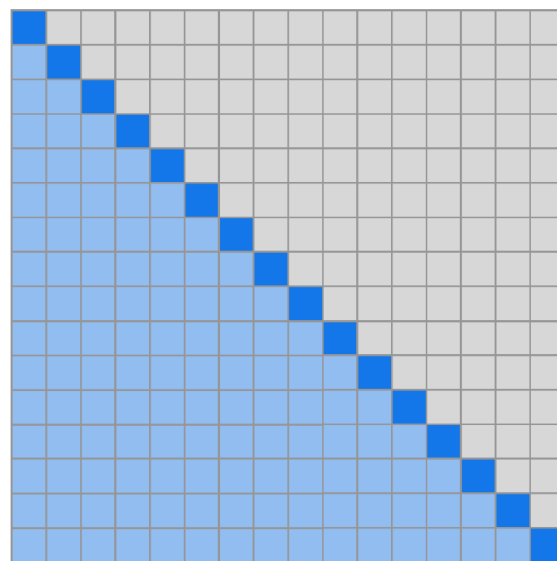
$$O(N^2)$$

- 显存和算力需求大
无法用于长序列



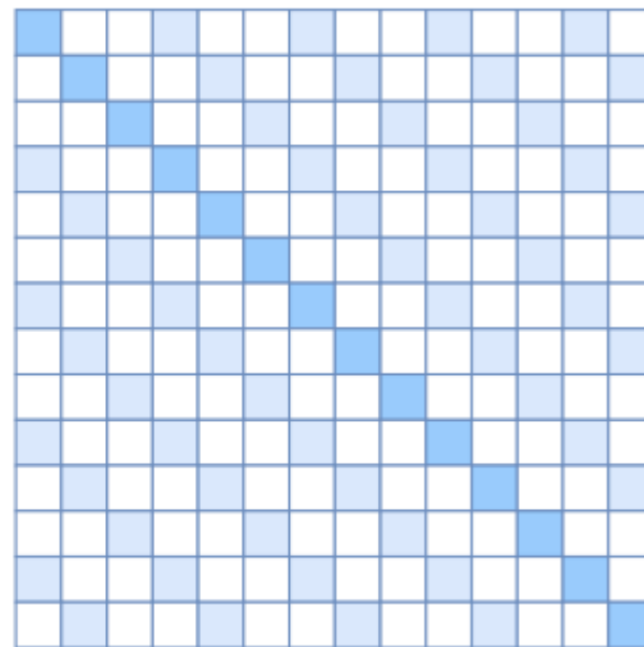
稀疏注意力 (Openai 2019)

- 《Generating Long Sequences with Sparse Transformers》
- 两种：strided跨步稀疏注意力和fixed固定稀疏注意力
- 跨步稀疏注意力=跳跃Atrous（局部头）+前后Local（稀疏头）
- Self Attention



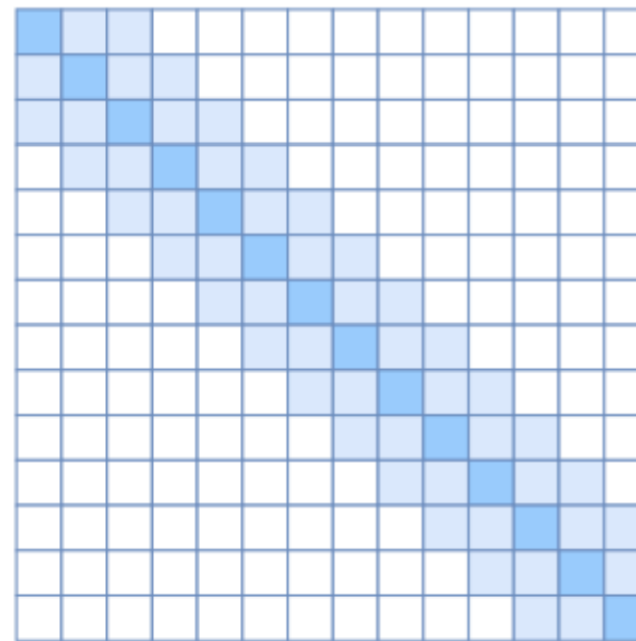
Atrous Self Attention (苏剑林 2019)

- 每个元素只跟它相对距离为 $k, 2k, 3k, \dots$ 的元素关联
- 至于其他的token, 我们不再计算注意力
- 每个元素只与 $\frac{n}{k}$ 个算相关性
- 这样, 时间和空间复杂度
都降低到原来的 $\frac{1}{k}$



Local Self Attention (苏剑林 2019)

- 每个元素只与前后 k 个元素以及自身有关联
- 超出这个范围的token不计算注意力
- 本质是加了权重的卷积



跨步稀疏注意力

= 跳跃Atrous (局部头) + 前后Local (稀疏头)

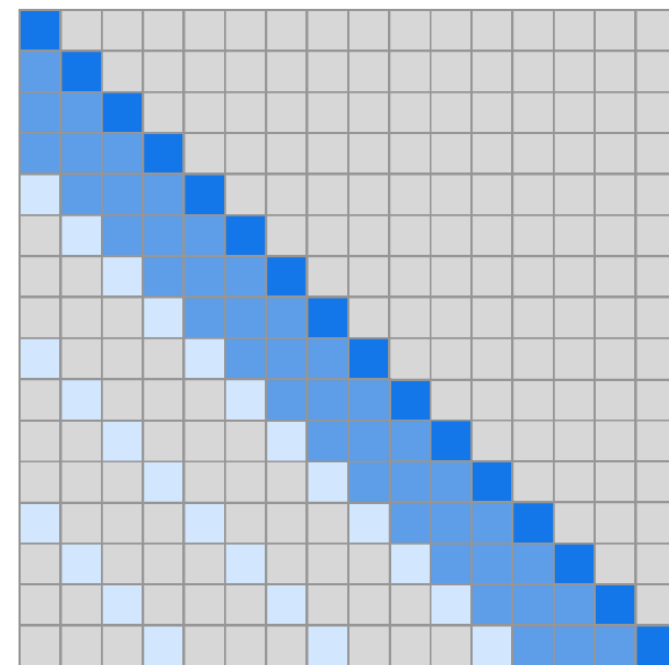
Formally, $A_i^{(1)} = \{t, t + 1, \dots, i\}$ for $t = \max(0, i - l)$ and $A_i^{(2)} = \{j : (i - j) \bmod l = 0\}$. This pattern can be visualized in Figure 3(b).

- 局部头

关注 $(i - l, i)$ 里的 token l 就是跨步的步长
若 $i - l < 0$, 则从序列位置 0 开始

- 稀疏头

关注 “与 i 的间隔为 l 的整数倍” 的历史 token
每隔 l 个位置选一个 token
跳跃 固定跳跃间隔



固定稀疏注意力

Formally, $A_i^{(1)} = \{j : (\lfloor j/l \rfloor = \lfloor i/l \rfloor)\}$, where the brackets denote the floor operation, and $A_i^{(2)} = \{j : j \bmod l \in \{t, t+1, \dots, l\}\}$, where $t = l - c$ and c is a hyperparameter.

- 局部头

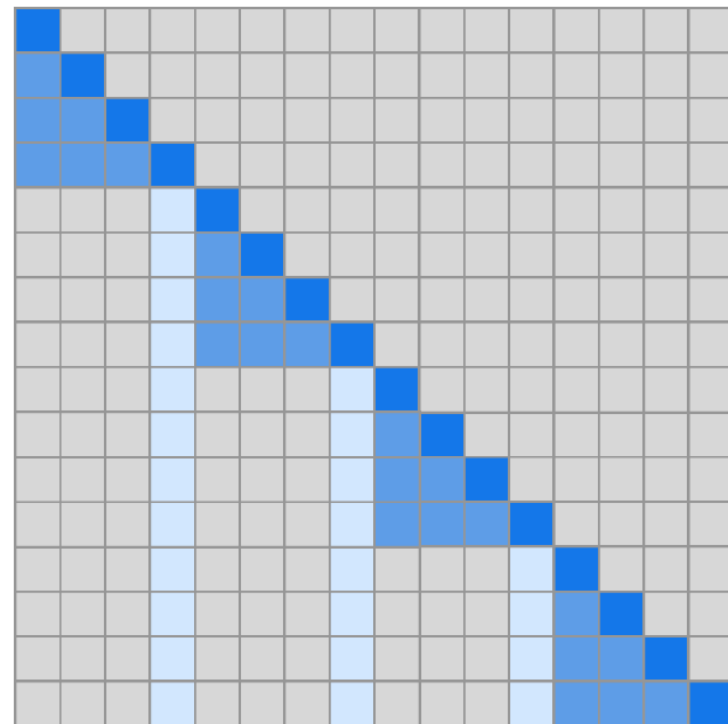
关注区块内的所有token

向下取整求出i所在的区块

- 稀疏头

不关注每个区块的所有 token

仅保留 最后c个关键 token



ACC-ViT (CVPR 2024)

- 成倍跳跃的稀疏注意力
- Atrous Convolution
- 空洞卷积

tioning of the featuremaps. For a given input featuremap $x \in \mathcal{R}^{c,h,w}$, where c, h, w refers to the number of channels, height, and width, respectively, we can compute windows at different levels of dilation as,

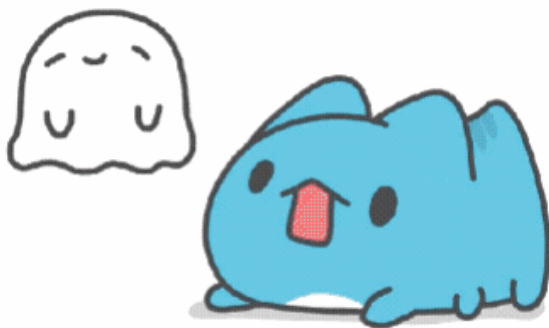
$$x_{atr-k}[c, i, j] = x[c, i + 2k - 1, j + 2k - 1] \quad (1)$$

For example,

$$x_{atr-1}[c, i, j] = x[c, i + 1, j + 1] \quad (2)$$

$$x_{atr-2}[c, i, j] = x[c, i + 3, j + 3] \quad (3)$$

$$x_{atr-3}[c, i, j] = x[c, i + 7, j + 7] \quad (4)$$



DeepSeek Sparse Attention

- Lightning Indexer 闪电索引器

轻量级计算token间的索引分数，筛选关键token

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

- 对每个查询token h_t 与他之前的所有token
- 用FP8精度

DeepSeek Sparse Attention



- fine-grained token selection 细粒度 token 选择

仅用选中的 k 个 token 来计算注意力

$$O(L^2) \rightarrow O(Lk)$$

- 计算注意力

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

把选中的 top- k 个分数对应的 c_s 与 h_t 求注意力分数

这里是 MQA 查询头 c_s 共享同一组 KV 条目

- 成功降低了计算复杂度