

课程大作业二基础题：GPU学习智能体（agent）

✓ 补充资料（二）

PagedAttention 的定义

✓ **定义：** PagedAttention 借鉴了操作系统中的分页（Paging）机制：将连续的逻辑内存（这里指 KV Cache）分割为固定大小的“块（Block）”，通过“块表（Block Table）”记录逻辑块与物理内存块的映射关系，从而允许逻辑上连续的 KV Cache 存储在物理上不连续的内存块中。

PagedAttention 的原理：KV Cache 的“块”划分

- ✓ PagedAttention将每个序列的KV Cache按固定大小（例如16或32个token）分割为多个块（Block）。每个块是内存管理的基本单元，物理上存储在GPU内存的连续空间中。
- ✓ 例如：一个长度为100 token的序列，若块大小为16，则需要7个块（前6块填满16 token，最后1块存4 token）。

PagedAttention 的原理：块表（Block Table）的映射作用

- ✓ 每个序列对应一个块表，记录该序列的 KV Cache 逻辑块与物理内存块的映射关系。例如，序列A的块表可能记录：“逻辑块 0 \rightarrow 物理块 5”、“逻辑块 1 \rightarrow 物理块 12”等。
- ✓ 块表的存在使得逻辑上连续的KV Cache（对模型而言，需要按序列顺序访问历史 token）可以被拆分为物理上不连续的块存储，无需依赖连续内存。

PagedAttention 的原理：注意力计算时的块访问

- ✓ 在LLM的注意力机制计算中（如自注意力公式），需要按序列顺序访问KV Cache的所有历史token：

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})V$$

- ✓ PagedAttention在计算时，通过块表找到当前序列对应的所有物理块，将这些块中的KV数据“拼接”为逻辑上连续的序列（无需实际复制内存，仅通过索引映射），再参与注意力计算。这一过程通过GPU高效的内存访问指令优化，几乎不增加额外开销。

PagedAttention 的原理：动态内存管理

- ✓ **块分配**：当序列生成新token时，若当前块已满，PagedAttention会从空闲块池中分配一个新的物理块，并更新块表。
- ✓ **块释放**：当序列结束（如生成 <EOS> 或被终止），其占用的物理块会被释放回空闲池，供其他序列复用。
- ✓ 这种动态分配机制避免了传统连续内存分配的“预分配冗余”和“扩容复制”问题，显著减少了内存浪费。

PagedAttention 的核心优势

- ✓ **内存利用率提升**：通过块级复用和碎片化减少，相同 GPU 内存可容纳更长的序列或更大的批次，内存利用率通常提升 2-4 倍。
- ✓ **支持动态序列长度**：无需为每个序列预分配最大长度内存，适应批次中序列长度差异大的场景（如实时对话中，不同用户的输入长度不同）。
- ✓ **高效批处理**：允许在有限内存中处理更多序列，提升推理吞吐量（Throughput）。
- ✓ **支持“换出”机制**：当内存不足时，可将暂时不用的块“换出”到 CPU 内存或磁盘，需要时再“换入”GPU 内存（类似虚拟内存），进一步突破 GPU 内存限制。

PagedAttention 与传统 KV Cache 管理的对比

特性	传统连续内存分配	PagedAttention
内存连续性要求	必须连续	允许不连续（通过块表映射）
内存利用率	低（易产生碎片和冗余）	高（块级复用，无连续内存依赖）
动态序列长度支持	差（扩容需复制数据）	好（动态分配 / 释放块）
最大批次 / 序列长度	受限于连续内存大小	受限于总空闲块数量

PagedAttention 的优化：核心机制优化

- ✓ **块大小与布局调整**：默认块大小为 16 token，支持自定义（e.g., 8—32）。优化时，根据序列长度和硬件（e.g., A100/H100 GPU）调整块大小，平衡并行度和碎片。特殊内存布局将键/值缓存分离（K-cache: [num_blocks, num_kv_heads, d_h / x, B, x]; V-cache: [num_blocks, num_kv_heads, d_h, B]），减少全局内存访问。
- ✓ **Copy-on-Write与共享**：允许多序列（如beam search或并行采样）共享物理块，通过引用计数跟踪。修改时复制新块，节省38%—55%内存。优化共享前缀（如聊天机器人系统提示），加速 1.67—3.58 倍。

PagedAttention 的优化：核心机制优化

- ✓ 驱逐与恢复策略：采用全或无驱逐（evict 整个序列组）。优化包括交换到 CPU（swapping）或重新计算（recomputation, overhead ~20%）。对于长序列，优先重新计算以减少 I/O。
- ✓ 连续批处理集成：结合连续批处理（dynamic batching），在生成阶段动态加入/移除序列，提升利用率 2—8 倍。

PagedAttention 的优化：硬件与内核优化

- ✓ 自定义 CUDA/PTX 内核：vLLM 使用模板化内核（paged_attention_kernel）融合 QK/Softmax/PV 阶段。优化包括 warp-level shuffle 归约、共享内存缓冲 Q/logits、异步拷贝（cp.async 在 Ampere+ GPU）。在 H100 上集成 Tensor Core，提升 20-26% 性能。
- ✓ Triton 与 TileLang 实现：简化编写，使用 tl.dot/tl.softmax 处理分页块，支持 AMD/NVIDIA 跨平台。Tiling 和 pipelining 隐藏延迟，提升 1.4-2x 速度。
- ✓ FlexAttention 集成：从零构建 vLLM，使用 PyTorch 的 FlexAttention 线性化分页逻辑，支持动态形状和混合精度。
- ✓ CXL-Attached 内存分配：针对长上下文微调，使用 CXL 附加内存扩展 KV 缓存，减少 OOM，提升微调效率。

PagedAttention 的优化：量化与压缩优

- ✓ **FP4/FP8量化**：使用NVFP4或MXFP4量化权重/激活/KV缓存（per-tensor 或 per-group 尺度），节省2—4x内存，支持更大批次。结合 llm-compressor 校准（512 样本），在 vLLM 中无缝集成，提升 1.6-3x 吞吐。
- ✓ **分组查询注意力 (GQA)**：替换多头注意力（MHA），减少 KV 头数（e.g., 8:1 比率），节省 ~75% KV 缓存内存，同时保持性能。
- ✓ **冗余层去除**：分析注意力层相似度，移除半数冗余层（Attention Drop），加速 48.4%—50% 而准确率仅降 2.4%。联合 MLP/Attention Drop 进一步优化。

PagedAttention 的优化：总结

优化类别	示例方法	收益示例
内存管理	分页块、Copy-on-Write	碎片 <4%，节省 38—55%
内核加速	CUDA/Triton/FlexAttention	1.4—2x 速度，跨平台
量化压缩	FP4/GQA/冗余层去除	2—4x 内存节省，1.6-3x TPS
混合架构	Star/Ring-Linear/Kimi Linear	6—11x 长序列加速，降 75% KV
整体系统	连续批处理/投机解码	2—8x 吞吐，成本降 50%