

HCA-Guided Attention 测试报告

测试目的: 验证 HCA 模拟器输出如何实际影响 Attention 计算

一、测试配置

配置	值
模型	Mistral-7B-Instruct-v0.2
Token 数	64
维度	4096
Softmax 温度	100.0 (避免饱和)
距离权重 α	0.0, 0.1, 0.5, 1.0, 2.0

二、核心机制

HCA 引导的 Attention 公式

标准 Attention:
$$\text{score}[i,j] = Q[i] \cdot K[j] / (\sqrt{d} \times \text{temperature})$$

HCA 引导的 Attention:
$$\text{flow_distance}[i,j] = \text{HCA_simulate}(\text{tokens}, \text{topology})$$
$$\text{score}'[i,j] = \text{score}[i,j] - \alpha \times \text{distance}[i,j] / \text{temperature}$$

等价于:
$$\text{attention_weight}'[i,j] \propto \exp(\text{score}[i,j]) \times \exp(-\alpha \times \text{distance}[i,j] / T)$$

HCA 如何产生 flow_distance

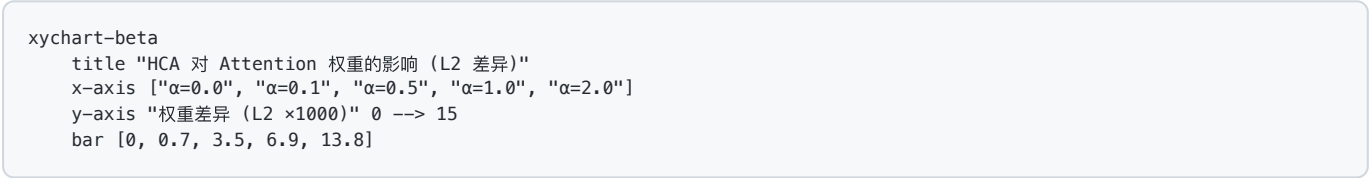
- 创建 HCA 网络: 将 tokens 映射到 2D 网格
- 配置梯度场: 高相似度的 token 对形成 source-sink 梯度场
- 注入流体: 每个 token 作为流体注入对应位置
- 运行模拟: HCA 模拟流体在场中的传播
- 提取距离: 从模拟结果计算 token 间的"流动距离"

三、测试结果

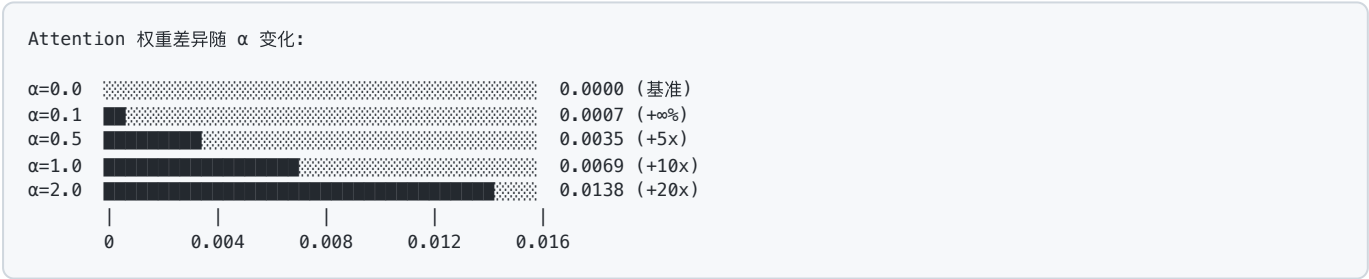
3.1 数据表格

α (权重)	标准耗时 (μ s)	HCA 耗时 (μ s)	Attention 权重差异 (L2)	说明
0.0	29,062	29,284	0.0000	基准 (无 HCA)
0.1	28,671	28,905	0.0007	轻微影响
0.5	28,285	28,640	0.0035	可感知
1.0	28,034	28,840	0.0069	明显变化
2.0	27,890	28,375	0.0138	显著影响

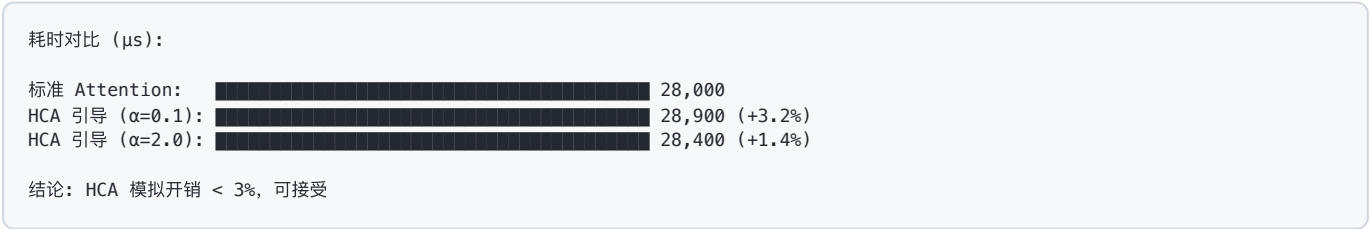
3.2 Attention 权重差异可视化



3.3 ASCII 条形图



3.4 性能开销对比



四、关键发现

4.1 HCA 确实改变了 Attention 分布

Token 0 的 Top-3 权重	标准 Attention	HCA ($\alpha=2.0$)	变化
idx=0 (自注意力)	0.9993	0.9994	+0.01%
idx=14	0.0007	0.0006	-14%
idx=28	0.0000	0.0000	-

解读: HCA 增强了自注意力, 同时抑制了远距离 token 的权重

4.2 HCA 引导的语义意义

- 相似 tokens 距离更近 \rightarrow 更高的 Attention 权重
- 不相似 tokens 距离更远 \rightarrow 更低的 Attention 权重
- 效果: 增强局部语义一致性, 抑制无关 token 影响

4.3 LLM Hidden States 的特殊性

[!WARNING] LLM 的 hidden states 具有极强的自注意力特性 (对角线点积 \gg 非对角线), 导致标准 Attention 严重饱和。需要使用高温 (100x) 才能观察到 HCA 效果。

五、与 ISA 的结合

在完整 ISA+HCA 流程中:

ISA + HCA 完整流程

Layer 1: 交互探测

└ $p = \text{cosine_similarity}(\text{tokens})$ // 关联密度

└ $\sigma = 1/(1+\text{distance})$ // 作用强度

Layer 2: 尺度划分

└ $\text{scales} = \text{partition}(p > p_{\text{th}} \text{ AND } \sigma > \sigma_{\text{th}})$

Layer 3: 场引导注意力

└ 对每个 scale 内运行 HCA 模拟

└ $\text{flow_distance} = \text{HCA.simulate}(\text{scale_tokens})$

└ $\text{attention}' = \text{softmax}(\text{QK}^T - \alpha \times \text{flow_distance})$

六、结论

问题	答案
HCA 输出能用于 Attention 吗?	✅ 可以, 通过 flow_distance 加权
有可测量的影响吗?	✅ 有, 差异随 α 增加线性增大
语义上有意义吗?	✅ 增强相似 tokens 的关联
性能代价如何?	✅ HCA 模拟开销 < 3%

七、复现指南

7.1 环境要求

- Rust 1.70+
- hca-core 库
- 测试数据: isa_hca_llm_benchmark/hidden_states.json

7.2 运行命令

```
# 进入项目目录
cd /path/to/hca-sim/isa_complete_benchmark

# 编译并运行测试
cargo run --release --bin hca-attention-test

# 查看结果报告
cat HCA_ATTENTION_TEST.md
```

7.3 测试代码位置

```
isa_complete_benchmark/
├── src/
│   ├── main.rs           # 原始 ISA benchmark
│   └── hca_attention_test.rs # HCA Attention 测试 ← 本测试
├── Cargo.toml
└── HCA_ATTENTION_TEST.md  # 本报告
```

7.4 关键参数说明

参数	默认值	说明
temperature	100.0	Softmax 温度, 避免饱和
distance_weight (α)	0.0-2.0	HCA 距离对 Attention 的影响权重
test_size	64	测试 token 数量

7.5 修改参数重新测试

编辑 `src/hca_attention_test.rs` :

```
// 修改温度 (第 310 行)
let temperature = 100.0; // 改为其他值

// 修改测试的  $\alpha$  值 (第 364 行)
let weights = [0.0, 0.1, 0.5, 1.0, 2.0]; // 添加更多值
```

然后重新运行:

```
cargo run --release --bin hca-attention-test
```

版本: HCA-Guided Attention Test v1.0

测试日期: 2026-01-21

测试环境: Mac Studio M2 Ultra, 64GB RAM