

# 测试指标

## 推理性能测试指标

### 1. 吞吐量 (Throughput) :

- 测量每秒处理的 tokens 数量 (tokens per second) 。
- 在不同的 batch size 和序列长度下测试, 以评估模型的扩展能力。

### 2. 延迟 (Latency) :

- 测量生成单个 token 所需的时间。
- 对于实时应用, 低延迟是关键。

### 3. 内存使用 (Memory Utilization) :

- 监测显存使用情况, 确保推理过程中不会超出显存限制。
- 评估在不同序列长度和 batch size 下的内存需求。

### 4. 计算效率 (Compute Efficiency) :

- 评估 GPU 的利用率和计算资源使用效率。
- 确保内核执行的有效性和资源的充分利用。

### 5. 能耗 (Power Consumption) :

- 衡量每瓦特的推理性能 (Performance per Watt) 。
- 重要性在于降低运营成本和提高效率。

## 训练性能测试指标

### 1. 训练速度 (Training Speed) :

- 测量每秒处理的样本数 (samples per second) 。
- 在不同模型规模和 batch size 下测试。

### 2. 收敛速度 (Convergence Speed) :

- 评估模型达到特定精度或损失水平所需的时间。
- 观察不同优化器和学习率的影响。

### 3. 扩展性 (Scalability) :

- 测试在多 GPU 或分布式环境下的性能扩展能力。
- 测试数据并行和模型并行的效率。

#### 4. 稳定性 (Stability) :

- 长时间训练的稳定性测试, 观察是否存在崩溃或性能下降。

#### 5. 内存管理 (Memory Management) :

- 评估显存使用和管理效率, 包括检查内存泄漏和碎片化。

## 功能测试指标

#### 1. 兼容性测试:

- 确保 GPU 能够正确运行在目标软件栈 (如 VLLM、Deepspeed、CUDA、cuDNN、TensorRT) 上。
- 测试与不同深度学习框架 (如 TensorFlow、PyTorch) 的兼容性。

#### 2. 混合精度训练 (Mixed Precision Training) :

- 验证支持混合精度计算 (FP16/BF16) 的能力, 以提高训练速度和降低内存占用。

#### 3. 错误处理能力:

- 测试 GPU 在遇到错误时的响应能力, 如硬件故障、超时、数据传输错误等。

#### 4. 硬件加速功能:

- 验证特定硬件加速功能 (如张量核心) 的有效性和性能增益。

#### 5. 能效 (Energy Efficiency) :

- 测量在不同负载下的功耗, 并计算每瓦特性能。

#### 6. 模型精度 (Model Accuracy) :

- 确保在优化性能的过程中不牺牲模型的精度。
- 验证模型在训练和推理过程中的输出质量。