

PR #21667 完整报告

sgl-project/sglang

Unify GSM8K eval path to Chat API for regression CI readiness

合并时间: 2026-04-02 08:12

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/21667>

执行摘要

本次 PR 将三个独立的 GSM8K 评估实现统一至单一入口点 `run_eval.py`，通过 Chat API 的 Completion 接口路由，为所有 NVIDIA 测试文件启用 `dump_metric` 自动覆盖（包括得分、延迟、输出吞吐量）。这是 CI 回归测试基础设施的关键第一步，旨在消除代码冗余、提升测试可观测性，同时保持行为与旧实现完全一致。变更涉及 79 个文件，大量测试迁移，但 AMD/NPU 路径通过弃用警告保留以确保兼容性。

功能与动机

SGLang 当前存在 8 个独立的评估系统 (Issue #21046)，导致显著的代码重复和维护负担。以 GSM8K 为例，竟有 5 个不同实现，包括 `few_shot_gsm8k.py`、`simple_eval_gsm8k.py` 和 `few_shot_gsm8k_engine.py`，各有不同的接口和指标收集能力。PR body 明确指出：“After this PR: one entry point (`run_eval`), one args format, one set of metric keys, same Completion API behavior.” 这直接支持 Issue #21157 中“基于回归的 CI 检查”的 Phase 1，为后续自动基线比较奠定基础。

实现拆解

框架层变更

文件	关键改动	目的
<code>accuracy_test_runner.py</code>	删除 <code>_run_few_shot_eval</code> 函数，移除 GSM8K 特殊路由	将 GSM8K 评估统一至 <code>run_eval</code> 路径
<code>run_eval.py</code>	使 <code>model</code> 参数可选 (<code>getattr(args, "model", None)</code>)，为 Completion API 添加默认停止词	支持服务器自动检测并匹配旧行为
<code>eval_accuracy_kit.py</code>	更新 <code>GSM8KMixin.test_gsm8k</code> 使用新参数： <code>base_url</code> 、 <code>eval_name="gsm8k"</code> 、 <code>api="completion"</code>	统一 Mixin 接口，影响众多测试类
<code>few_shot_gsm8k.py</code>	添加 <code>DeprecationWarning</code> ，提示迁移至 <code>run_eval</code>	提供平滑过渡，保留 AMD/NPU 使用

测试迁移

约 72 个 NVIDIA 测试文件完成以下转换：

- 导入从 `from sglang.test.few_shot_gsm8k import run_eval` 改为 `from sglang.test.run_eval import run_eval`
- 参数从 `SimpleNamespace(host=..., port=...)` 转为 `SimpleNamespace(base_url=..., eval_name="gsm8k", api="completion", ...)`
- 指标键从 `metrics["accuracy"]` 转为 `metrics["score"]`
- 所有准确性阈值保持不变，确保评估结果一致

评论区精华

由于没有正式 review 讨论，但从 55 次 commit 历史中可提炼关键设计交锋：

兼容性权衡：在 commit '66d15c8f' 中，作者决定“Keep few_shot_gsm8k files for AMD/NPU tests, add deprecation warnings”，这反映了对团队自主权的尊重——AMD 和 NPU 测试由各自团队维护，不在此次 PR 中强制迁移，避免了破坏性变更。API 细节优化：commit 'dfbedab7' 添加了 `max_tokens=512` 和默认停止词 `["Question", "Assistant:", "<lseparatorl>"]` 到 Completion API，以确保与旧 `few_shot_gsm8k.py` 行为完全匹配，体现了对边缘用例的细致处理。

风险与影响

技术风险：

1. 回归风险：大规模文件修改可能引入参数错误，例如遗漏 `base_url` 设置或 `model` 属性缺失（commit '6db5913b' 曾修复此类问题）。
2. 性能偏差：新的 Chat API 路径可能引入额外网络开销，尽管 PR 声称行为一致，但需在 CI 中监控 `dump_metric` 收集的延迟指标。
3. 兼容性裂痕：AMD/NPU 测试保留旧路径，长期可能加剧代码分歧，增加跨平台维护成本。

影响评估：

- 正面：统一入口点减少代码重复，为 CI 回归测试提供标准化指标收集，提升团队效率。
- 负面：短期需 AMD/NPU 团队响应弃用警告进行迁移，否则可能错过 `dump_metric` 覆盖。
- 范围：影响所有 NVIDIA 测试的 GSM8K 评估，但对最终用户透明。

关联脉络

本 PR 是更大基础设施演进的一部分：

- 直接关联 Issue：#21046（整合碎片化评估系统）和 #21157（回归 CI 检查路线图），本 PR 实现了两者的初期目标。
- 历史 PR 趋势：近期 PR 如 #21873（添加评估数据集下载超时）和 #21830（修复 CI 测试稳定性）显示项目正加强 CI 可靠性，而本 PR 的 `dump_metric` 覆盖是这一趋势的延续。
- 未来方向：统一 GSM8K 路径后，预计类似变更将扩展至 MMLU、MMMU 等其他数据集，逐步实现评估系统的全面标准化。