

PR #24149 完整报告

sgl-project/sglang

feat(bench): add SPEED-Bench dataset support to bench_serving

合并时间: 2026-05-29 05:37

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/24149>

执行摘要

- 一句话: 为 bench_serving 添加 SPEED-Bench 数据集支持, 用于推测解码基准测试
- 推荐动作: 值得关注。该 PR 展示了如何通过标准的 BaseDataset 接口扩展 bench_serving 的数据集类型, 是一种良好的插件式设计。对于计划集成其他标准化基准 (如 MMLU、HumanEval) 的开发者具有参考价值。review 中的优化建议和采纳过程也体现了代码质量意识。

功能与动机

PR body 明确指出: SPEED-Bench (SPEculative Evaluation Dataset) 是专为评估推测解码 (SD) 算法而设计的统一基准, 覆盖多样语义领域和真实服务场景。将其作为原生数据集选项, 可以直接通过 bench_serving 进行可重现的类别特定吞吐量基准测试, 无需自定义预处理脚本。

实现拆解

1. 创建 SpeedBenchDataset 数据集插件: 在 python/sglang/benchmark/datasets/speed_bench.py 中定义 SpeedBenchDataset 类 (继承 BaseDataset), 实现 from_args 方法解析 CLI 参数 (dataset_path, speed_bench_category, speed_bench_output_len, num_prompts), 实现 load 方法读取预下载的 JSONL 文件, 按类别过滤, 使用 tokenizer 应用 chat template 生成标准化的 DataRow 列表, 并处理采样和 shuffle。
2. 注册数据集: 在 python/sglang/benchmark/datasets/__init__.py 中导入 SpeedBenchDataset 并加入 DATASET_MAPPING, 映射字符串 'speed-bench'。
3. 添加 CLI 参数: 在 python/sglang/bench_serving.py 的 --dataset-name 选项的 choices 中添加 'speed-bench', 并新增两个专用参数: --speed-bench-category (可选值 low_entropy/mixed/high_entropy) 和 --speed-bench-output-len (默认 512)。
4. 编写单元测试: 在 test/registered/bench_fn/test_benchmark_datasets_api.py 中添加 5 个测试方法 (test_speed_bench_sampler, test_speed_bench_category_filter, test_speed_bench_output_len_override, test_speed_bench_empty_category_raises, test_speed_bench_no_path_raises), 覆盖正常采样、类别过滤、输出长度覆盖、空类别报错、缺少路径报错等场景。
5. 更新文档: 在 docs_new/docs/developer_guide/bench_serving.mdx 中添加 SPEED-Bench 的描述、CLI 参数说明和使用示例 (含推测解码服务启动和 bench_serving

命令)。

关键文件:

- `python/sglang/benchmark/datasets/speed_bench.py` (模块 数据集; 类别 `source`; 类型 `dependency-wiring`; 符号 `SpeedBenchDataset`, `from_args`, `load`): 新增的 SPEED-Bench 数据集插件, 核心变更文件, 包含 `SpeedBenchDataset` 类和 `from_args/load` 方法。
- `test/registered/bench_fn/test_benchmark_datasets_api.py` (模块 测试; 类别 `test`; 类型 `test-coverage`; 符号 `_write_speed_bench_jsonl`, `test_speed_bench_sampler`, `test_speed_bench_category_filter`, `test_speed_bench_output_len_override`): 单元测试文件, 覆盖 `SpeedBenchDataset` 的主要用例和错误路径, 确保插件正确性。
- `python/sglang/bench_serving.py` (模块 入口脚本; 类别 `source`; 类型 `core-logic`): 入口脚本, 添加了 `speed-bench` 到 `dataset choices` 并注册两个新 CLI 参数, 使用户可通过命令行直接指定。
- `python/sglang/benchmark/datasets/__init__.py` (模块 数据集; 类别 `source`; 类型 `dependency-wiring`): 数据集注册中心, 从映射中增加 `speed-bench` 条目, 使得 `get_dataset` 可按名称派发。
- `docs_new/docs/developer_guide/bench_serving.mdx` (模块 文档; 类别 `other`; 类型 `documentation`): 文档更新, 包含 SPEED-Bench 的说明、CLI 参数和使用示例, 帮助用户正确使用新功能。

关键符号: `SpeedBenchDataset`, `SpeedBenchDataset.from_args`,
`SpeedBenchDataset.load`, `test_speed_bench_sampler`,
`test_speed_bench_category_filter`, `test_speed_bench_output_len_override`,
`test_speed_bench_empty_category_raises`

关键源码片段

`python/sglang/benchmark/datasets/speed_bench.py`

新增的 SPEED-Bench 数据集插件, 核心变更文件, 包含 `SpeedBenchDataset` 类和 `from_args/load` 方法。

```
@classmethod
def from_args(cls, args: Namespace) -> 'SpeedBenchDataset':
    # 如果没有提供 dataset_path, 则抛出 ValueError
    if not args.dataset_path:
        raise ValueError(
            '--dataset-path must point to the SPEED-Bench JSONL file '
            '(run prepare_speed_bench.sh to generate it).'
        )
    # 使用 getattr 安全读取专用参数, 提供默认值 (category 默认 None 不过滤)
    return cls(
        dataset_path=args.dataset_path,
        category=getattr(args, 'speed_bench_category', None) or None,
        output_len=getattr(args, 'speed_bench_output_len', 512),
        num_requests=args.num_prompts,
```

```

)

def load(
    self, tokenizer: PreTrainedTokenizerBase, model_id=None
) -> List[DatasetRow]:
    # 先读取文件, 按类别过滤, 提取第一条用户消息
    unique_prompts = []
    with open(self.dataset_path, encoding='utf-8') as f:
        for line in f:
            row = json.loads(line)
            if self.category and row.get('category') != self.category:
                continue
            turns = row.get('turns', [])
            if not turns:
                continue
            unique_prompts.append(turns[0])
    if not unique_prompts:
        raise ValueError(
            f'No rows found in {self.dataset_path}'
            + (f' for category={self.category}' if self.category else '')
        )
    # 先 tokenize 唯一 prompt (而非对重复行多次 tokenize), 然后采样
    unique_dataset_rows: List[DatasetRow] = []
    for prompt_text in unique_prompts:
        try:
            prompt_ids = tokenizer.apply_chat_template(
                [{'role': 'user', 'content': prompt_text}],
                add_generation_prompt=True,
                tokenize=True,
            )
            prompt = tokenizer.decode(prompt_ids)
        except Exception:
            prompt_ids = tokenizer.encode(prompt_text)
            prompt = prompt_text
        unique_dataset_rows.append(
            DatasetRow(
                prompt=prompt,
                prompt_len=len(prompt_ids),
                output_len=self.output_len,
            )
        )
    # 采样 (可能 oversample), shuffle 以模拟真实分布
    if self.num_requests <= len(unique_dataset_rows):
        dataset_rows = random.sample(unique_dataset_rows, self.num_requests)
    else:
        dataset_rows = unique_dataset_rows * (
            self.num_requests // len(unique_dataset_rows) + 1
        )
    dataset_rows = dataset_rows[: self.num_requests]

```

```
        random.shuffle(dataset_rows)
    return dataset_rows
```

test/registered/bench_fn/test_benchmark_datasets_api.py

单元测试文件，覆盖 SpeedBenchDataset 的主要用例和错误路径，确保插件正确性。

```
def test_speed_bench_sampler(self):
    # 创建临时 JSONL 文件
    dataset_path = self._write_speed_bench_jsonl()
    args = make_args(
        dataset_name='speed-bench',
        dataset_path=dataset_path,
        num_prompts=3,
    )
    from sglang.benchmark.datasets.speed_bench import SpeedBenchDataset
    dataset = SpeedBenchDataset.from_args(args)
    rows = dataset.load(self.tokenizer)
    self.assertEqual(len(rows), 3)
    self.assertTrue(all(isinstance(row, DatasetRow) for row in rows))
    self.assertTrue(all(row.output_len == 512 for row in rows))
    self.assertTrue(all(row.prompt_len > 0 for row in rows))

def test_speed_bench_empty_category_raises(self):
    # 验证空类别时是否抛出 ValueError
    dataset_path = self._write_speed_bench_jsonl()
    args = make_args(
        dataset_name='speed-bench',
        dataset_path=dataset_path,
        num_prompts=100,
        speed_bench_category='nonexistent',
    )
    from sglang.benchmark.datasets.speed_bench import SpeedBenchDataset
    dataset = SpeedBenchDataset.from_args(args)
    with self.assertRaises(ValueError):
        dataset.load(self.tokenizer)
```

评论区精华

在 code review 中，gemini-code-assist[bot] 对 SpeedBenchDataset.load 方法提出了三点优化建议：（1）先 tokenize 唯一 prompt 再采样，避免冗余工作；（2）oversampling 后添加 shuffle，避免确定性重复序列；（3）打开文件时指定 encoding='utf-8' 以增加健壮性。作者 jmamou 在后续 commit (b5fc85c) 中全部采纳并修复。无其他争议。

- SpeedBenchDataset.load 优化建议 (performance): 作者 jmamou 在后续 commit 中全部采纳，并关闭了该 thread。

风险与影响

- 风险：风险较低。变更仅限 bench_serving 数据加载和 CLI 扩展，不涉及模型前向或核心调度。主要风险为：依赖外部 JSONL 文件格式与内容是否符合预期；若 dataset_path 未指定或文件损坏，load 方法抛出 ValueError 中断基准测试；测试覆盖了主流路径（类别过滤、输出长度覆盖、空类别等）。对运行时无影响。
- 影响：用户影响：为使用 bench_serving 进行推测解码评估的用户提供了便捷的标准化数据集选项，无需额外脚本即可复用 SPEED-Bench 的吞吐量子集（1K–32K 输入长度，按熵分类）。系统影响：无运行时影响，仅编译前数据准备步骤。团队影响：新增的数据集插件可被其他基准测试模块参考（约 102 行代码），维护负担低。
- 风险标记：新增外部数据依赖，文件格式验证

关联脉络

- 暂无明显关联 PR