

PR #23753 完整报告

sgl-project/sglang

tokenizer: Add fastokens support

合并时间: 2026-04-29 02:43

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23753>

执行摘要

- 一句话: 新增 fastokens tokenizer 后端, 加速 TTFT
- 推荐动作: 该 PR 设计清晰, 值得精读。关键看点是: ①如何通过 monkey-patch 无缝替换 tokenizer 后端; ②如何在保证现有路径不变的前提下引入可选高速路径。对于希望了解 SGLang tokenizer 抽象层的开发者来说, 这是很好的学习材料。

功能与动机

PR body 明确指出 "Add fastokens tokenizer support in SGLang - enables significantly faster TTFT"。提供的基准测试显示, 在 Qwen3.5-35B-A3B 模型上, TTFT 中位数从 2765ms 降至 2178ms (降低约 21%), 输入吞吐量从 48k tok/s 提升至 58k tok/s。

实现拆解

1. 配置入口: 在 `server_args.py` 的 `ServerArgs` 中添加 `tokenizer_backend` 字段, 默认值为 "huggingface", 并注册 CLI 参数 `--tokenizer-backend`, 可选值 ["huggingface", "fastokens"]。
2. 核心 patching 函数: 在 `hf_transformers/tokenizer.py` 中添加 `_ensure_fastokens_patched()`, 使用模块级全局变量确保只执行一次 monkey-patch。通过 `import fastokens; fastokens.patch_transformers()` 替换 HuggingFace tokenizers 后端为 fastokens 的 `_TokenizerShim`。
3. Tokenzier 加载路径: 修改 `get_tokenizer()` 函数, 新增 `tokenizer_backend` 参数。当指定为 fastokens 时, 先调用 `_ensure_fastokens_patched()`, 然后在加载 tokenizer 后跳过 `_TOKENIZERS_BACKEND` 回退逻辑 (因为 fastokens 已经返回了正确的 shim)。加载失败时给出明确的错误信息并建议移除 `--tokenizer-backend=fastokens`。
4. Processor 加载路径: 在 `hf_transformers/processor.py` 的 `get_processor()` 中添加相似逻辑, 确保通过 processor 获取 tokenizer 时也能正确应用 fastokens 补丁, 并将 `tokenizer_backend` 参数透传给内部 `get_tokenizer()` 调用。
5. 调用点透传: 在 `scheduler.py`、`tokenizer_manager.py`、`tp_worker.py`、`detokenizer_manager.py`、`encode_receiver.py` 中, 将 `server_args.tokenizer_backend` 传递给 `get_tokenizer / get_processor` 调用。
6. 依赖与文档: 在 `pyproject.toml` 中将 fastokens 注册为可选依赖 (`[project.optional-dependencies] fastokens = ["fastokens"]`), 并加入 test 组以确保 CI

安装。更新 `server_arguments.mdx` 文档。

7. 测试配套：新增 `test/registered/unit/utils/test_hf_transformers_fastokens.py`，包含两个测试：验证 `tokenizer._tokenizer` 是 `_TokenizerShim` 实例，以及 `encode-decode` 往返正确性。测试注册为 CPU CI，依赖 `fastokens` 包时自动跳过。

关键文件：

- `python/sglang/srt/utils/hf_transformers/tokenizer.py`（模块 `Tokenizer`；类别 `source`；类型 `dependency-wiring`；符号 `_ensure_fastokens_patched`）：核心改动文件，添加 `_ensure_fastokens_patched` 函数并修改 `get_tokenizer` 以支持 `fastokens` 后端。
- `test/registered/unit/utils/test_hf_transformers_fastokens.py`（模块 `测试`；类别 `test`；类型 `test-coverage`；符号 `TestFastokensBackend`, `test_shim_is_applied`, `test_encode_decode_roundtrip`）：新测试文件，验证 `fastokens` 后端正确应用并确保 `encode-decode` 往返正确，注册为 CPU CI。
- `python/sglang/srt/utils/hf_transformers/processor.py`（模块 `处理器`；类别 `source`；类型 `dependency-wiring`）：需要在 `processor` 加载路径中也应用 `fastokens` 补丁，确保多模态模型使用 `processor` 时也能受益。
- `python/sglang/srt/server_args.py`（模块 `配置`；类别 `source`；类型 `core-logic`）：添加 `tokenizer_backend` 字段和 `--tokenizer-backend` CLI 参数，定义可选值。
- `python/sglang/srt/managers/scheduler.py`（模块 `调度器`；类别 `source`；类型 `core-logic`）：在调度器初始化 `tokenizer` 时传递 `tokenizer_backend` 参数，是集成点之一。
- `python/sglang/srt/managers/tokenizer_manager.py`（模块 `Tokenizer`；类别 `source`；类型 `core-logic`）：在多处初始化 `tokenizer/processor` 时传递 `tokenizer_backend` 参数。
- `python/sglang/srt/disaggregation/encode_receiver.py`（模块 `分离式`；类别 `source`；类型 `core-logic`）：分离式编码接收器需要初始化 `tokenizer`，需传递参数。
- `python/sglang/srt/managers/tp_worker.py`（模块 `工作器`；类别 `source`；类型 `core-logic`）：TP worker 也需要初始化 `tokenizer`。
- `python/sglang/srt/managers/detokenizer_manager.py`（模块 `去 Tokenzier`；类别 `source`；类型 `core-logic`）：`detokenizer` 管理器初始化 `tokenizer` 时需传递参数。
- `python/pyproject.toml`（模块 `依赖`；类别 `config`；类型 `configuration`）：添加 `fastokens` 可选依赖并加入 `test` 依赖组，确保 CI 安装。
- `docs_new/docs/advanced_features/server_arguments.mdx`（模块 `文档`；类别 `other`；类型 `core-logic`）：文档更新，记录 `--tokenizer-backend` 参数。

关键符号：`_ensure_fastokens_patched`, `get_tokenizer`, `get_processor`,
`TestFastokensBackend.test_shim_is_applied`,
`TestFastokensBackend.test_encode_decode_roundtrip`

关键源码片段

`python/sglang/srt/utils/hf_transformers/tokenizer.py`

核心改动文件，添加 `_ensure_fastokens_patched` 函数并修改 `get_tokenizer` 以支持 `fastokens` 后端。

```
# python/sglang/srt/utils/hf_transformers/tokenizer.py
```

```
_fasttokens_patched = False
```

```
def _ensure_fasttokens_patched():  
    """Monkey-patch transformers to use the fasttokens backend (once)."""  
    global _fasttokens_patched  
    if _fasttokens_patched:  
        return  
    try:  
        import fasttokens  
    except ImportError:  
        raise ImportError(  
            "The fasttokens package is required when --tokenizer-backend=fasttokens. "  
            "Install it with: pip install 'sglang[fasttokens]'"  
        ) from None  
  
    fasttokens.patch_transformers()  
    _fasttokens_patched = True  
    logger.info("fasttokens backend enabled - transformers patched successfully")
```

```
def get_tokenizer(  
    tokenizer_name: str,  
    *args,  
    tokenizer_mode: str = "auto",  
    trust_remote_code: bool = False,  
    tokenizer_revision: Optional[str] = None,  
    tokenizer_backend: str = "huggingface",  
    **kwargs,  
    ) -> Union[PreTrainedTokenizer, PreTrainedTokenizerFast]:  
    """Gets a tokenizer for the given model name via Huggingface."""  
    # Tiktoken format has its own backend — no fasttokens patching needed.  
    if tokenizer_name.endswith(".json"):  
        from sglang.srt.tokenizer.tiktoken_tokenizer import TiktokenTokenizer  
        return TiktokenTokenizer(tokenizer_name)  
  
    if tokenizer_backend == "fasttokens":  
        _ensure_fasttokens_patched()  
  
    if tokenizer_mode == "slow":  
        if kwargs.get("use_fast", False):  
            raise ValueError("Cannot use the fast tokenizer in slow tokenizer mode.")  
        kwargs["use_fast"] = False  
    elif tokenizer_mode == "auto":  
        if "use_fast" not in kwargs:  
            kwargs["use_fast"] = True
```

```

tokenizer_name = _resolve_tokenizer_name(tokenizer_name, kwargs)

common_kwargs = dict(
    trust_remote_code=trust_remote_code,
    tokenizer_revision=tokenizer_revision,
    clean_up_tokenization_spaces=False,
    **kwargs,
)

try:
    tokenizer = _auto_tokenizer_from_pretrained(
        tokenizer_name, *args, **common_kwargs
    )

    # 当使用 fasttokens 时，补丁后的 TokenizersBackend.from_pretrained 已经返回了
    # 包含 fasttokens shim 的 tokenizer，因此不需要重新解析回具体的 tokenizer 类。
    if (
        type(tokenizer).__name__ == _TOKENIZERS_BACKEND
        and tokenizer_backend != "fasttokens"
    ):
        tokenizer = _resolve_tokenizers_backend(
            tokenizer_name, *args, **common_kwargs
        )

    return _apply_post_load_fixes(tokenizer, tokenizer_name, tokenizer_revision)
except Exception as e:
    if tokenizer_backend == "fasttokens":
        raise RuntimeError(
            f"fasttokens failed to load tokenizer for {tokenizer_name!r}. "
            f"This model's tokenizer may not be supported by fasttokens — "
            f"see https://github.com/crusoecloud/fasttokens. "
            f"Re-run without --tokenizer-backend=fasttokens to use the default backend."
        ) from e
    raise

```

test/registered/unit/utils/test_hf_transformers_fasttokens.py

新测试文件，验证 fasttokens 后端正确应用并确保 encode-decode 往返正确，注册为 CPU CI。

```
# test/registered/unit/utils/test_hf_transformers_fasttokens.py
```

```

import unittest
from sglang.test.ci.ci_register import register_cpu_ci
from sglang.test.test_utils import (
    DEFAULT_SMALL_MODEL_NAME_FOR_TEST_QWEN,
    CustomTestCase,
)

```

```
TOKENIZER_MODEL = DEFAULT_SMALL_MODEL_NAME_FOR_TEST_QWEN
```

```
register_cpu_ci(est_time=30, suite="stage-a-test-cpu")
```

```

try:
    import fasttokens # noqa: F401
    HAS_FASTOKENS = True
except ImportError:
    HAS_FASTOKENS = False

@unittest.skipUnless(HAS_FASTOKENS, "fasttokens package not installed")
class TestFasttokensBackend(CustomTestCase):
    def test_shim_is_applied(self):
        # _TokenizerShim 是 fasttokens 的私有兼容 shim。
        # SGLang 集成依赖于 tokenizer._tokenizer 是该类的实例，
        # 以确认 fasttokens 已正确加载。
        from fasttokens._compat import _TokenizerShim
        from sglang.srt.utils.hf_transformers.tokenizer import get_tokenizer

        tokenizer = get_tokenizer(
            TOKENIZER_MODEL,
            tokenizer_backend="fasttokens",
        )
        backend = getattr(tokenizer, "_tokenizer", None)
        self.assertIsInstance(
            backend,
            _TokenizerShim,
            f"Expected tokenizer._tokenizer to be _TokenizerShim, "
            f"got {type(backend).__name__}",
        )

    def test_encode_decode_roundtrip(self):
        from sglang.srt.utils.hf_transformers.tokenizer import get_tokenizer

        tokenizer = get_tokenizer(
            TOKENIZER_MODEL,
            tokenizer_backend="fasttokens",
        )
        text = "Hello, world!"
        ids = tokenizer.encode(text, add_special_tokens=False)
        self.assertGreater(len(ids), 0)
        self.assertEqual(tokenizer.decode(ids, skip_special_tokens=True), text)

if __name__ == "__main__":
    unittest.main()

```

评论区精华

1. 版本依赖与错误处理：审查者 alexnails 要求添加版本控制并建议在模型不支持 fasttokens 时提供回退或清晰的错误提示。作者添加了 pyproject.toml 依赖，并在加载失败时抛出 RuntimeError，指导用户移除 --tokenizer-backend=fasttokens 回退到默认后端。

2. 测试模型选择：审查者指出测试中使用了硬编码的 Qwen/Qwen3-0.6B，建议使用工具函数 `DEFAULT_SMALL_MODEL_NAME_FOR_TEST_QWEN`。作者已采纳。
3. CI 依赖安装：审查者发现 CI 可能因为未安装 `fastokens` 而始终跳过测试，建议将 `sglang[fastokens]` 加入 `pyproject.toml` 的 `test` 依赖组。作者已添加。
4. 代码风格：审查者要求避免在终端日志中使用长破折号（em dash），作者已删除。
 - 依赖版本控制与错误回退 (design): 添加了依赖声明和友好的错误消息，未实现自动回退但提供了明确指引。
 - 测试模型选择与 CI 依赖 (testing): 作者采用 `DEFAULT_SMALL_MODEL_NAME_FOR_TEST_QWEN` 并在 `pyproject.toml` 的 `test` 组中添加 `fastokens` 依赖。
 - 代码风格：避免 em dash 等特殊字符 (style): 作者移除了 em dash 字符。

风险与影响

- 风险：
 1. 兼容性风险：fastokens 尚未覆盖所有模型 tokenizer，当模型不兼容时，用户会看到明确的 `RuntimeError` 并被告知回退到 HuggingFace 后端，不会静默失败。
 2. 回归风险：fastokens 仅在 `--tokenizer-backend=fastokens` 时启用，默认路径完全不变，不影响现有用户。但修改了 tokenizer 和 processor 的加载流程（新增参数、条件分支），需确保所有调用点都正确传递参数。已在多个调用点测试。
 3. 外部依赖风险：fastokens 是一个外部 Rust 库，版本更新可能引入 breaking changes。当前未固定版本，需关注上游变化。未来可考虑锁定最小版本。
 4. 性能风险：fastokens 声称更快，但若模型不兼容或边缘情况可能退化。不过由于是可选项，用户可自行评估。
 - 影响：对用户：提供显著的速度提升（TTFT 降低约 20%，吞吐量提升约 20%），且完全向后兼容。用户只需升级并安装可选依赖即可受益。对系统：新增一个可选依赖，不影响默认部署。对团队：需维护 fastokens 集成代码，关注上游变化。该模式为未来引入其他 tokenizer 后端（如 tokenizers Rust）提供了范例。
 - 风险标记：新外部依赖，可选功能未全覆盖，多个调用点需同步

关联脉络

- 暂无明显关联 PR