

PR #23506 完整报告

sgl-project/sglang

[gRPC] Native server: Rust crate (1/4)

合并时间: 2026-05-19 13:31

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/23506>

执行摘要

该 PR 引入了 `rust/sglang-grpc` crate，这是 SGLang 原生 gRPC 服务器实现的第一阶段。它建立了 Rust 侧与 Python 运行时的桥接，实现了 tonic gRPC 服务处理器、Rust 原生分词器封装以及请求序列化工具。当前 PR 仅为 crate 本身，不改变默认运行行为，但为后续集成 PR 奠定基础。

功能与动机

当前 SGLang 依赖于外部 gRPC 代理进行进程间通信，这增加了延迟和部署复杂度。原生进程内 gRPC 服务器旨在消除代理跳转，提供更低延迟和更高吞吐量。此 PR 是拆分 PR 栈的第一部分，专注于 Rust crate 结构 (PR #22907)。

实现拆解

1. 模块入口 (`lib.rs`) : 定义 `GrpcServerHandle` 和 `start_server` Python 函数。新增 `extract_tokenizer_info` 从 Python `RuntimeHandle` 提取 tokenizer 配置。启动后台 Tokio 运行时并绑定 `TcpListener`。
2. 请求桥接层 (`bridge.rs`) : `PyBridge` 管理每个请求的 `mpsc` 通道和回调状态。`ResponseChunk` 枚举描述响应片段。`TerminalError` 分类不可恢复错误。通过 `lock_or_recover` 安全获取互斥锁。提供 `abort()` 接口供服务器取消请求。
3. gRPC 服务处理器 (`server.rs`) : `SglangServiceImpl` 实现 tonic 服务 trait，涵盖文本生成、嵌入、分类、`tokenize/detokenize`、控制流等方法。使用 `pyerr_to_status` 映射 Python 错误，`RequestAbortGuard` 确保 Drop 时自动取消。`resolve_max_message_size` 通过环境变量配置大小限制。
4. Rust 原生分词器 (`tokenizers.rs`) : `TokenizerBackend` trait 定义 `encode/decode` 接口。`HuggingFaceTokenizerBackend` 包装 `tokenizers::Tokenizer` crate。`RustTokenizer::from_tokenizer_path` 在 `tokenizer.json` 存在且非 `slow` 模式时加载，否则静默回退 Python。
5. 工具函数 (`utils/request_utils.rs`) : `sampling_params_to_map` 等函数将 proto 结构转为 Python dict 兼容格式。
6. 测试覆盖: 覆盖错误映射、环境变量覆盖、JSON 编码回退等核心行为。

以下为 `server.rs` 和 `bridge.rs` 的核心实现 (完整代码见对应文件)。

`server.rs` 核心结构与辅助函数:

```

use std::pin::Pin;
use std::sync::Arc;
use tokio::sync::mpsc::Receiver;
use tokio::time::{Duration, timeout};
use tokio_stream::Stream;
use tonic::{Request, Response, Status};
use pyo3::PyErr;
use pyo3::exceptions::{PyTypeError, PyValueError};
use crate::bridge::{PyBridge, ResponseChunk, TerminalError};

pub type StreamResult<T> = Pin<Box<dyn Stream<Item = Result<T, Status>> + Send + 'static>>;
pub const DEFAULT_GRPC_MAX_MESSAGE_SIZE: usize = 64 * 1024 * 1024;

pub fn resolve_max_message_size() -> usize {
    match std::env::var("SGLANG_TONIC_PAYLOAD") {
        Ok(raw) => match raw.parse::<usize>() {
            Ok(n) if n > 0 => { tracing::info!("..."); n },
            _ => { tracing::warn!("..."); DEFAULT_GRPC_MAX_MESSAGE_SIZE }
        },
        _ => DEFAULT_GRPC_MAX_MESSAGE_SIZE,
    }
}

pub fn pyerr_to_status(err: PyErr, context: &str) -> Status {
    let is_client_error = Python::with_gil(|py| {
        err.is_instance_of::<PyValueError>(py) || err.is_instance_of::<PyTypeError>(py)
    });
    let msg = format!("{}", context, err);
    if is_client_error { Status::invalid_argument(msg) }
    else { Status::internal(msg) }
}

struct RequestAbortGuard {
    bridge: Arc<PyBridge>,
    rid: String,
    armed: bool,
}

impl RequestAbortGuard {
    fn new(bridge: Arc<PyBridge>, rid: impl Into<String>) -> Self { ... }
    fn disarm(&mut self) { self.armed = false; }
    fn abort_now(&mut self) { if self.armed { self.armed = false; spawn_abort(...); } }
}

impl Drop for RequestAbortGuard {
    fn drop(&mut self) { if self.armed { spawn_abort(...); } }
}

```

bridge.rs核心类型与桥接逻辑:

```

pub enum ResponseChunk {
    Data(ResponseData),

```

```

    Finished(ResponseData),
    Error(String),
}
impl ResponseChunk {
    pub fn is_terminal(&self) -> bool {
        matches!(self, Self::Finished(_) | Self::Error(_))
    }
}

pub struct ResponseData {
    pub text: Option<String>,
    pub output_ids: Option<Vec<i32>>,
    pub embedding: Option<Vec<f32>>,
    pub json_bytes: Option<Vec<u8>>,
    pub meta_info: HashMap<String, String>,
}

pub enum TerminalError {
    ChannelFull { rid: String },
    ClientDisconnected { rid: String },
    Aborted { rid: String },
}

fn lock_or_recover<'a, T>(mutex: &'a Mutex<T>, name: &'static str) -> MutexGuard<'a, T> {
    mutex.lock().unwrap_or_else(|poisoned| {
        tracing::warn!(mutex = name, "Recovering from poisoned gRPC bridge mutex");
        poisoned.into_inner()
    })
}

pub struct PyBridge {
    runtime_handle: PyObject,
    state: Arc<Mutex<BridgeState>>,
    rust_tokenizer: Option<RustTokenizer>,
    context_len: i32,
    response_channel_capacity: usize,
    tokio_handle: Handle,
}

impl PyBridge {
    pub fn new(runtime_handle: PyObject, rust_tokenizer: Option<RustTokenizer>,
              context_len: i32, response_channel_capacity: usize,
              tokio_handle: Handle) -> Self {
        debug_assert!(response_channel_capacity > 0);
        Self { ..., state: Arc::new(Mutex::new(BridgeState::default())), ... }
    }

    pub fn create_channel(&self, rid: &str) -> PyResult<Receiver<ResponseChunk>> {
        let (sender, receiver) = mpsc::channel(self.response_channel_capacity);
        let mut state = lock_or_recover(self.state.as_ref(), "state");
        if state.channels.contains_key(rid) {

```

```

        return Err(PyRuntimeError::new_err(format!("Duplicate rid: {}", rid)));
    }
    state.channels.insert(rid.to_string(), sender);
    state.terminal_errors.remove(rid);
    state.ready_callbacks.remove(rid);
    state.ready_signals.remove(rid);
    state.pending_sends.remove(rid);
    Ok(receiver)
}
}

```

评论区精华

- ishandhanani 质疑桥接层信号量的必要性，认为调度器已有流量控制；alexnails 最终移除信号量。
- JustinTong0323 要求将验证错误映射为 `INVALID_ARGUMENT` 而非 `INTERNAL`；alexnails 实现 `pyerr_to_status` 区分客户端 / 服务器错误。
- JustinTong0323 指出需要显式消息大小限制；alexnails 添加 64 MiB 默认值与 `SGLANG_TONIC_PAYLOAD` 环境变量覆盖。
- ishandhanani 发现 `meta_info` 序列化使用 `v.str()` 产生 Python repr 而非 JSON，影响客户端解析；暂未修复。
- JustinTong0323 强调 gRPC 服务缺少 API 认证；同意在后续 PR 中解决。
- ispobock 质疑使用独立 maturin wheel；alexnails 移除 wheel，回归嵌入式构建。

风险与影响

- 缺少认证：gRPC 端点当前无防护，若在启用了 `auth_key` 的部署中暴露，可能被未授权访问。当前 PR 仅为 `crate`，后续集成时需优先解决。
- 元数据序列化不一致：`meta_info` 字段输出 Python repr，客户端需要特殊处理。
- 桥接层锁争用：单 `Mutex` 保护所有通道，高并发下可能成为瓶颈。
- 默认消息大小 64MiB：虽然可覆盖，但过大可能导致内存压力。
- 序列化开销：`proto` → `JSON` → `Python dict` 路径增加 CPU 使用。
- 分词器静默回退：加载失败时静默回退 Python，可能隐藏配置错误。

关联脉络

- 此 PR 是 PR #22907 的拆解之一（第一阶段），后续将有 Python 入口集成、服务器启动参数绑定和集成测试 PR。
- 依赖 PR #22736 中的 `protobuf` 定义，并遵循 RFC #22558 的设计方向。
- 与近期重构 PR（如 #25724、#25727）无关，但 gRPC 服务器后续需适配调度器变更。