

PR #20475 完整报告

sgl-project/sglang

feat: add Crusoe managed inference backend

合并时间: 2026-05-13 07:23

原文链接: <http://prhub.com.cn/sgl-project/sglang/pull/20475>

执行摘要

- 一句话: 新增 Crusoe 托管推理后端
- 推荐动作: 该 PR 设计清晰, 适合作为第三方后端集成的参考模式 (继承 + LazyImport)。建议合并后补充 CI 可运行的 mock 测试以覆盖回归。

功能与动机

PR 描述指出: 需要将 Crusoe managed inference 作为一等后端集成到 sglang 中, 利用其 OpenAI 兼容 API 提供无缝体验。

实现拆解

1. 在 `python/sglang/lang/backend/crusoe.py` 中定义 Crusoe 类, 继承自 OpenAI, 重写 `__init__` 以从环境变量或参数获取 API key, 并设置默认 `base_url`。
2. 在 `python/sglang/init.py` 中通过 LazyImport 注册 Crusoe 符号, 并加入 `__all__` 导出列表, 使用户可从 `sgl.Crusoe` 访问。
3. 新增 `test/manual/test_crusoe_backend.py`, 包含两类测试: `TestCrusoeBackend` (端到端调用, 需真实 API key) 和 `TestCrusoeBackendInit` (无网络依赖的初始化验证, 如缺失 key、显式 key、自定义 `base_url`)。

关键文件:

- `python/sglang/lang/backend/crusoe.py` (模块 Crusoe 后端; 类别 source; 类型 dependency-wiring; 符号 Crusoe, init) : 新增后端核心实现, 定义 Crusoe 类继承 OpenAI 并处理认证与端点配置。
- `test/manual/test_crusoe_backend.py` (模块 测试套件; 类别 test; 类型 test-coverage; 符号 TestCrusoeBackend, setUpClass, setUp, test_mt_bench) : 提供端到端测试和初始化验证, 保证后端可用性和正确性。
- `python/sglang/__init__.py` (模块 顶层接口; 类别 source; 类型 core-logic) : 将 Crusoe 注册为顶层可导入符号, 使用户可通过 `sgl.Crusoe` 访问。

关键符号: `Crusoe.init`, `TestCrusoeBackend.setUpClass`, `TestCrusoeBackend.test_mt_bench`, `TestCrusoeBackend.test_stream`, `TestCrusoeBackend.test_parallel_decoding`, `TestCrusoeBackend.test_parallel_encoding`, `TestCrusoeBackendInit.test_raises_without_api_key`,

TestCrusoeBackendInit.test_accepts_explicit_api_key,
TestCrusoeBackendInit.test_custom_base_url

关键源码片段

python/sglang/lang/backend/crusoe.py

新增后端核心实现，定义 Crusoe 类继承 OpenAI 并处理认证与端点配置。

```
import os
from typing import Optional

from sglang.lang.backend.openai import OpenAI
from sglang.lang.chat_template import ChatTemplate

CRUSOE_BASE_URL = 'https://managed-inference-api-proxy.crusoecloud.com/v1/'

class Crusoe(OpenAI):
    # SGLang backend for Crusoe managed inference.

    def __init__(
        self,
        model_name: str,
        api_key: Optional[str] = None,
        base_url: Optional[str] = None,
        chat_template: Optional[ChatTemplate] = None,
        **kwargs,
    ):
        # 优先使用传入的 api_key，否则从环境变量 CRUSOE_API_KEY 读取
        resolved_api_key = api_key or os.environ.get('CRUSOE_API_KEY')
        if not resolved_api_key:
            raise ValueError('Crusoe API key required. Pass api_key= or set CRUSOE_API_KEY.')

        # 调用 OpenAI 父类初始化
        super().__init__(
            model_name=model_name,
            chat_template=chat_template,
            api_key=resolved_api_key,
            base_url=base_url or CRUSOE_BASE_URL,
            **kwargs,
        )
```

test/manual/test_crusoe_backend.py

提供端到端测试和初始化验证，保证后端可用性和正确性。

```
# 手动测试文件，需要环境变量 CRUSOE_API_KEY
import unittest
from sglang import Crusoe, set_default_backend
from sglang.test.test_programs import (
```

```

    test_mt_bench, test_stream,
    test_parallel_decoding, test_parallel_encoding,
)
from sglang.test.test_utils import CustomTestCase

DEFAULT_CRUSOE_MODEL = 'nvidia/NVIDIA-Nemotron-3-Nano-30B-A3B'

class TestCrusoeBackend(CustomTestCase):
    # 端到端测试, 需要有效 API key
    backend = None

    @classmethod
    def setUpClass(cls):
        cls.backend = Crusoe(DEFAULT_CRUSOE_MODEL)

    def setUp(self):
        set_default_backend(self.backend)

    def test_mt_bench(self):
        test_mt_bench()

    def test_stream(self):
        test_stream()

    def test_parallel_decoding(self):
        test_parallel_decoding()

    def test_parallel_encoding(self):
        test_parallel_encoding()

class TestCrusoeBackendInit(CustomTestCase):
    # 无网络依赖的初始化验证

    def test_raises_without_api_key(self):
        import os
        key = os.environ.pop('CRUSOE_API_KEY', None)
        try:
            with self.assertRaises(ValueError):
                Crusoe(DEFAULT_CRUSOE_MODEL, api_key=None)
        finally:
            if key is not None:
                os.environ['CRUSOE_API_KEY'] = key

    def test_accepts_explicit_api_key(self):
        backend = Crusoe(DEFAULT_CRUSOE_MODEL, api_key='test-key')
        self.assertIsNotNone(backend)

```

```

def test_custom_base_url(self):
    backend = Crusoe(
        DEFAULT_CRUSOE_MODEL, api_key='test-key',
        base_url='https://managed-inference-api-proxy.crusoecloud.com/v1/',
    )
    self.assertIsNotNone(backend)

```

python/sglang/__init__.py

将 Crusoe 注册为顶层可导入符号，使用户可通过 `sgl.Crusoe` 访问。

```

# 在 LazyImport 块中添加 Crusoe 行
Anthropic = LazyImport('sglang.lang.backend.anthropic', 'Anthropic')
Crusoe = LazyImport('sglang.lang.backend.crusoe', 'Crusoe') # 新增
LiteLLM = LazyImport('sglang.lang.backend.litellm', 'LiteLLM')
OpenAI = LazyImport('sglang.lang.backend.openai', 'OpenAI')

```

在 `__all__` 列表中追加 'Crusoe'

```

__all__ = [
    'Engine',
    'Runtime',
    # ... 其他符号 ...
    'Anthropic',
    'Crusoe', # 新增条目
    'LiteLLM',
    'OpenAI',
    'VertexAI',
    'global_config',
    '__version__',
]

```

评论区精华

审核者 hnyls2002 首次要求添加测试文件 (CHANGES_REQUESTED)，作者随后补充测试后获得所有审核者批准。

- 请求添加测试文件 (testing): 作者添加测试后，hnyls2002 改为 APPROVED

风险与影响

- 风险:
 1. 依赖外部 API: Crusoe 端点的可用性和响应时间不在控制范围内。
 2. API 密钥安全: CRUSOE_API_KEY 可能被日志记录或泄露。
 3. 测试不可自动化: 手动测试需要实际网络连接和有效密钥，回归风险需人工验证。 - 影响: 用户: 获得直接调用 Crusoe 推理的能力。系统: 无核心组件改动，风险低。团队: 维护成本低，但需跟踪 Crusoe API 变更。 - 风险标记: 外部 API 依赖，API 密钥管理，测试非 CI 集成

关联脉络

- 暂无明显关联 PR