# Jaeger trace 文档

## 打点代码分支：

https://github.com/pyxyzc/sglang/tree/jeager

## 安装依赖服务

## 安装 启动 Jaeger （UI)

```
代码块
1   wget https://github.com/jaegertracing/jaeger/releases/download/v1.76.0/jaeger-
    1.76.0-linux-amd64.tar.gz
2
3   tar -xzvf jaeger-1.76.0-linux-amd64.tar.gz
4
5   cd jaeger-1.76.0-linux-amd64
6
7   # 使用绝对路径进行软连接保证全局可用，修改为jaeger-all-in-one的绝对路径
8   ln -s /sgl-workspace/jaeger/jaeger-1.76.0-linux-amd64/jaeger-all-in-one \
9   /usr/local/bin/jaeger-all-in-one
10
11  jaeger-all-in-one  --collector.otlp.enabled=true  --collector.otlp.grpc.host-
    port=0.0.0.0:4327  --collector.otlp.http.host-port=0.0.0.0:4328
```

报错 too many open files，具体报错信息：

```
{"level":"panic","ts":1765965126.7292247,"caller":"app/static_handler.go:46","msg":"Could not create static assets handler","error":"too many open files","stacktrace": ...}
```

```
panic: Could not create static assets handler
```

解决方法如下：

编辑 /etc/security/limits.conf

```
代码块
1   jaeger_user soft nofile 65536
2   jaeger_user hard nofile 65536
```

# 安装 启动 otel-collector (数据中转）

代码块

```
1   wget https://github.com/open-telemetry/opentelemetry-collector-
    releases/releases/download/v0.114.0/otelcol-contrib_0.114.0_linux_amd64.tar.gz
2
3   tar -xzvf otelcol-contrib_0.114.0_*.tar.gz
4
5   ln -s otelcol-contrib /usr/local/bin/otelcol-contrib
6
7
8   # 启动前先创建配置文件
9   otelcol-contrib  --config otel-collector.yaml
```

otel-collector.yaml

```
1   receivers:
2     otlp:
3       protocols:
4         grpc:
5           endpoint: 0.0.0.0:4317
6         http:
7           endpoint: 0.0.0.0:4318
8
9   exporters:
10    otlp:
11      endpoint: "localhost:4327"    # Jaeger OTLP gRPC 地址
12      tls:
13        insecure: true              # 本地环境常用
14
15  service:
16    pipelines:
17      traces:
18        receivers: [otlp]
19        exporters: [otlp]
20
```

# Sglang 启动参数

代码块

```
1   #!/bin/bash
```

```bash
 2
 3   # 模型路径（改成你本地的，比如 /data/models/Llama-2-7b-hf）
 4   MODEL_PATH=/home/models/QwQ-32B
 5
 6   # 端口号（可以改，比如 30000）
 7   PORT=30000
 8   # 卡数量
 9   TP=2
10
11   # 只使用
12   export CUDA_VISIBLE_DEVICES=0,1
13   # 启动命令
14   python python/sglang/launch_server.py \
15       --model-path $MODEL_PATH \
16       --tp $TP \
17       --port $PORT \
18       --enable-trace \
19       --otlp-traces-endpoint 127.0.0.1:4317 \
20       # --enable-metrics
```

## 同时启停脚本

需要将otel-collector启动配置文件放到sglang文件夹内，直接运行即可。

run_server_jaeger.sh

```bash
 1   #!/usr/bin/env bash
 2   set -e
 3
 4   #------------------------------
 5   # 端口检查函数
 6   #------------------------------
 7   check_and_kill_port() {
 8       local port=$1
 9       local pid
10
11       pid=$(lsof -t -i:"$port" || true)
12       if [[ -n "$pid" ]]; then
13           echo "Port $port is in use by PID $pid. Killing..."
14           kill -9 "$pid" || true
15           echo "Port $port cleared."
16       else
17           echo "Port $port is free."
18       fi
19   }
```
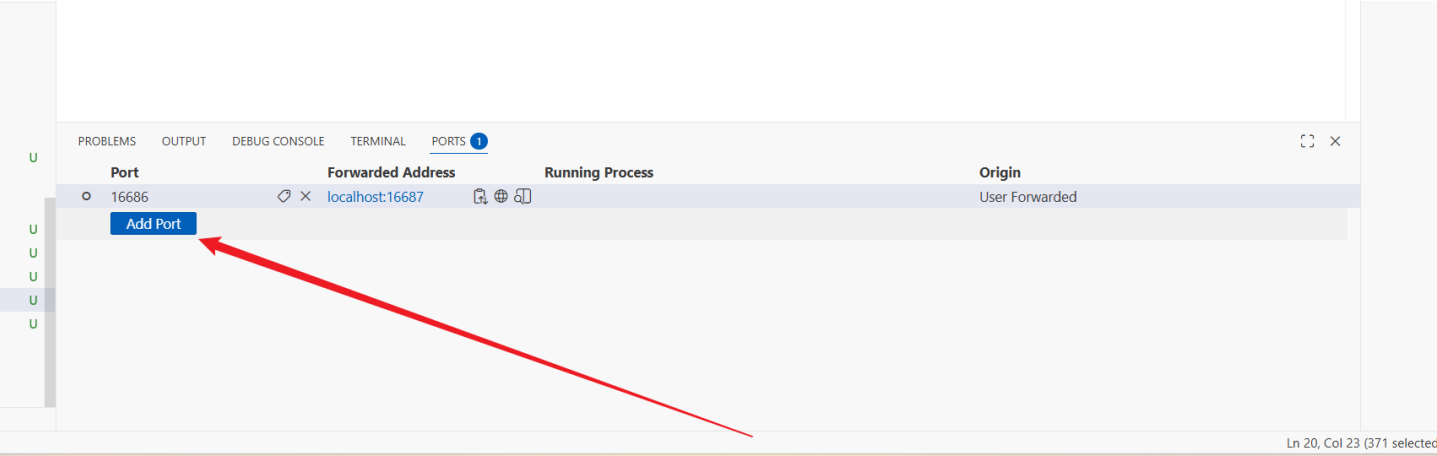
```bash
20
21   #------------------------------
22   # 启动前的环境变量
23   #------------------------------
24   export SGLANG_HICACHE_FILE_BACKEND_STORAGE_DIR=/sgl-workspace/hicachefile
25
26   MODEL_PATH=/home/models/QwQ-32B
27   PORT=30000
28   TP=2
29
30   export CUDA_VISIBLE_DEVICES=0,1
31
32   #------------------------------
33   # 启动前检查端口
34   #------------------------------
35   echo "Checking required ports..."
36   check_and_kill_port 4317   # otlp grpc
37   check_and_kill_port 4327   # jaeger grpc
38   check_and_kill_port 4328   # jaeger http
39   check_and_kill_port $PORT  # sglang server
40
41   #------------------------------
42   # 捕获退出信号，统一杀死所有子进程
43   #------------------------------
44   pids=()
45
46   cleanup() {
47       echo "Received stop signal. Killing all subprocesses..."
48       for pid in "${pids[@]}"; do
49           if kill -0 "$pid" 2>/dev/null; then
50               kill -9 "$pid" 2>/dev/null || true
51           fi
52       done
53       exit 1
54   }
55
56   trap cleanup SIGINT SIGTERM EXIT
57
58   #------------------------------
59   # 启动第 1 个: Jaeger
60   #------------------------------
61   jaeger-all-in-one \
62     --collector.otlp.enabled=true \
63     --collector.otlp.grpc.host-port=0.0.0.0:4327 \
64     --collector.otlp.http.host-port=0.0.0.0:4328 &
65   pids+=($!)
66   echo "Started Jaeger, PID=${pids[-1]}"
```

```
67
68    #-----------------------------
69    # 启动第 2 个: otelcol-contrib
70    #-----------------------------
71    otelcol-contrib --config otel-collector.yaml &
72    pids+=($!)
73    echo "Started otelcol-contrib, PID=${pids[-1]}"
74
75    #-----------------------------
76    # 启动第 3 个: sglang server
77    #-----------------------------
78    python python/sglang/launch_server.py \
79        --model-path $MODEL_PATH \
80        --page-size 128 \
81        --tp $TP \
82        --port $PORT \
83        --enable-hierarchical-cache \
84        --hicache-write-policy write_through \
85        --hicache-storage-backend file \
86        --hicache-storage-prefetch-policy wait_complete \
87        --enable-trace \
88        --otlp-traces-endpoint 127.0.0.1:4317 &
89    pids+=($!)
90    echo "Started sglang server, PID=${pids[-1]}"
91
92    #-----------------------------
93    # 等待任意子进程退出
94    #-----------------------------
95    echo "All processes started. Waiting for any to exit..."
96    wait -n
97
98    echo "One process exited. Killing all..."
99    cleanup
100
```
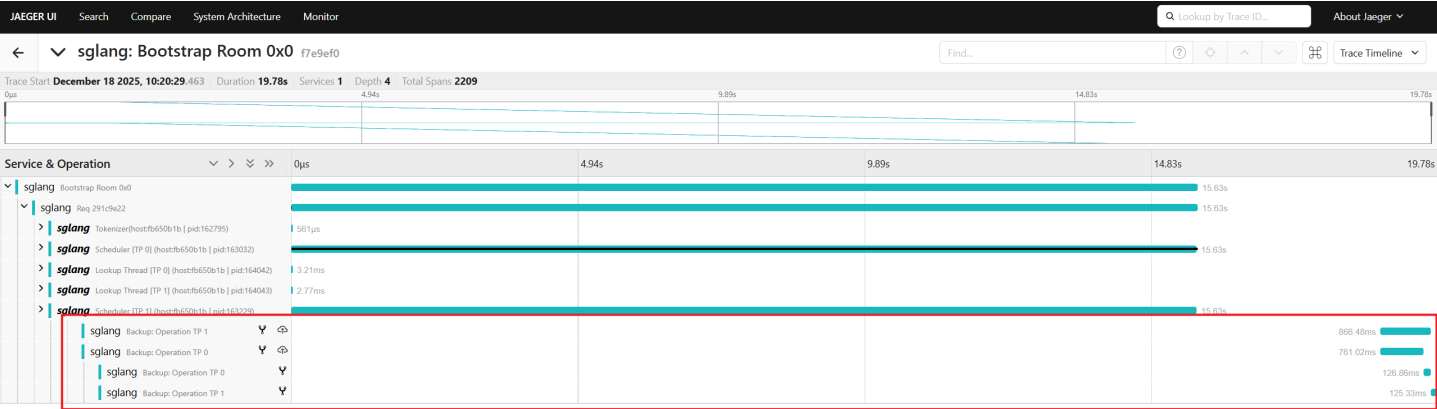
# 结果查看

## VScode 转发端口

`16686` 端口

# 查看打点数据！

## scheduler调度数据



## prefetch线程耗时

## lookup线程耗时



## dump线程耗时

对dump线程进行绑定后，可以记录dump耗时

# 添加打点教程

https://docs.sglang.io/references/production_request_trace.html

注意：

1. 和req相关的打点直接增加 `trace_slice_start` 、 `trace_slice_end` 即可