

```
llm-broker git:(master) # bash tools/check_chat.sh

"text": "LLM Brokerは、複数の大規模言語モデル（LLM）を統合し、リクエストのルーティン  
簡単に言えば、LLMを「仲介」して最適なモデルへ接続させる仕組みです。",
"meta": {
  "profile": "text_extract",
  "backend": "llamacpp",
  "model": "gpt-oss-20b-Q4_K_M",
  "base_url": "http://192.168.0.193:8081",
  "temperature": 0.1,
  "elapsed_ms": 1908,
  "prompt_tokens": 88,
  "completion_tokens": 167,
  "total_tokens": 255,
  "tokens_per_sec": 87.52290416741549
}
```

LLM Broker

Description

LLM Runner は、LLM Broker を利用して、LLM を呼び出すためのライブラリです。

LLM Broker は、LLM Runner を利用して、LLM を呼び出すためのライブラリです。

Xeon W-1250P (6C/12T) RTX 3060 Core i5 13500 RTX A4000 llama.cpp Strix Point GPU 80 tok/sec 100 tok/sec

llama.cpp Strix Point GPU llama.cpp Strix Point GPU LLM Runner LLM Broker

LLM Broker

```
profiles:  
  text_extract:  
    backend: llamacpp  
    base_url: http://192.168.0.193:8081  
    model: gpt-oss-20b-Q4_K_M  
    temperature: 0.1  
    timeout_sec: 360  
  
  text_conversation:  
    backend: llamacpp  
    base_url: http://192.168.0.193:8081  
    model: gpt-oss-20b-Q4_K_M  
    temperature: 0.7  
    timeout_sec: 360  
  
  text_extract_ollama_fallback:  
    backend: ollama  
    base_url: http://192.168.0.192:11434  
    model: gpt-oss:20b  
    temperature: 0.1  
    timeout_sec: 360  
  
  text_conversation_ollama_fallback:  
    backend: ollama  
    base_url: http://192.168.0.192:11434  
    model: gpt-oss:20b  
    temperature: 0.7  
    timeout_sec: 360
```

LLM Broker
LLM Runner

llama.cpp
Ollama

```
→ llm-broker git:(master) X bash tools/check_chat.sh
{
  "text_extract_ollama_fallback": {
    "text": "LLM Brokerは、複数の大規模言語モデル（LLM）を統合し、リクエストのルー  
簡単に言えば、LLMを「仲介」して最適なモデルへ接続させる仕組みです。",
    "meta": {
      "profile": "text_extract",
      "backend": "llamacpp",
      "model": "gpt-oss-20b-Q4_K_M",
      "base_url": "http://192.168.0.193:8081",
      "temperature": 0.1,
      "elapsed_ms": 1908,
      "prompt_tokens": 88,
      "completion_tokens": 167,
      "total_tokens": 255,
      "tokens_per_sec": 87.52290416741549
    }
  }
}
→ llm-broker git:(master) X LLM_PROFILE=text_extract_ollama_fallback bash tools
{
  "text": "LLM Brokerは、複数の大規模言語モデル（LLMs）を統合・管理し、リクエスト  
ミドルウェアです。これにより、用途や性能要件に応じて最適なモデルを選択でき、シス  
"meta": {
  "profile": "text_extract_ollama_fallback",
  "backend": "ollama",
  "model": "gpt-oss:20b",
  "base_url": "http://192.168.0.192:11434",
  "temperature": 0.1,
  "elapsed_ms": 4124,
  "prompt_tokens": 89,
  "completion_tokens": 183,
  "total_tokens": 272,
  "tokens_per_sec": 49.918425654688335
}
}
```

Xeon W-1250P + RTX 3060 Core i5 13400 + RTX 2000 Ada PC LLM Runner @

Category

- 1.

Tags

1. AI
2. llama.cpp
3. LLM

4. LLM Broker
5. Ollama
6. ä,?é??å??å,?
7. å±±æç"ç??

Date Created

2026å'7æ??4æ?¥

Author

kazuo-tsubaki

default watermark