

Tiny-LLM 项目报告

项目概述

Tiny-LLM 是一个基于 Rust 实现的高性能大语言模型推理引擎，支持 Llama 系列模型。该项目结合了 CUDA 加速、混合精度推理、Web 服务和多会话管理等功能，提供了一个完整的 LLM 推理解决方案。

功能特点

1. 高性能推理引擎

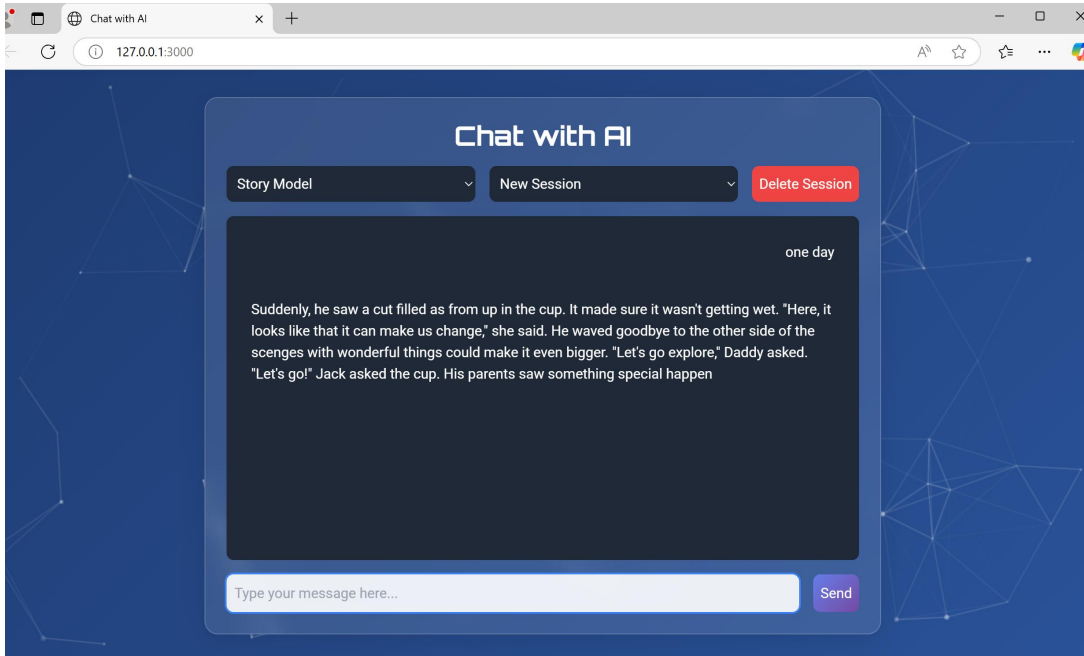
- 基于 Rust 实现的高效推理引擎
- 支持 Llama 系列模型
- 使用 KV Cache 优化推理性能
- 支持 top-k 和 top-p 采样策略

2. 混合精度推理

核心特性

- 模型参数使用 FP16 格式存储，减少内存占用
- 计算过程中根据需要转换为 FP32 格式，保证计算精度
- 支持参数精度转换
- 预留 FP16 计算接口
- 支持 FP16/FP32 之间的高效转换

此为基于CPU的混合精度推理的故事模型



实现细节

FP16 格式存储

```
// params.rs
pub struct LLamaParams<T> {
    pub embedding_table: Tensor<T>,
    pub rms_att_w: Vec<Tensor<T>>,
    // ... 其他参数
}

impl LLamaParams<f16> {
    pub fn from_safetensors(safetensor: &SafeTensors, config: &LlamaConfigJson) -> Self {
        // 从 safetensors 文件中加载 FP16 格式的参数
        let get_tensor = |name: &str| -> Tensor<f16> {
            // ...
        };
        // ...
    }
}
```

- 使用泛型 T 支持不同精度
- 实际使用 f16 类型存储模型参数
- 通过 safetensors 直接加载 FP16 格式数据

参数精度转换

```
// tensor.rs
impl Tensor<f16> {
    pub fn to_f32(&self) -> Tensor<f32> {
        let mut new_data = Vec::with_capacity(self.length);
        for x in self.data() {
            new_data.push(x.to_f32());
        }
        Tensor::new(new_data, &self.shape)
    }
}
```

- 提供了 FP16 到 FP32 的转换方法
- 在计算时按需转换为 FP32

CUDA 计算接口

```
// operators.rs
pub fn matmul_transb(c: &mut Tensor<f32>, beta: f32, a: &Tensor<f32>, b: &Tensor<f32>, alpha: f32) {
    // CUDA 计算目前只支持 FP32
    init_cuda();
    // ...
}

pub fn rms_norm(y: &mut Tensor<f32>, x: &Tensor<f32>, w: &Tensor<f32>, epsilon: f32) {
    // CUDA 计算目前只支持 FP32
    init_cuda();
    // ...
}
```

- 当前 CUDA 计算实现仅支持 FP32
- 预留了 FP16 计算的扩展空间

实际使用流程

```
// model.rs
impl Llama<f16> {
    pub fn from_safetensors(model_dir: impl AsRef<Path>) -> Self {
        // 加载 FP16 模型参数
        let params = LLamaParams::from_safetensors(&safetensor, &config);
        // ...
    }

    pub fn forward(&self, input: &Tensor<u32>, cache: &mut KVCache<f32>) -> Tensor<f32> {
        // 计算时转换为 FP32
        let params_f32 = self.params.to_f32();
        // 使用 FP32 进行计算
        // ...
    }
}
```

- 模型参数以 FP16 格式加载和存储
- 推理计算时转换为 FP32
- KV Cache 使用 FP32 格式

未来优化方向

虽然项目支持混合精度存储，但计算部分目前仍主要使用 FP32，之前尝试进行在cuda内核进行混合精度，并且将模型更多部分进行cuda的混合精度，但由于对 model.rs , parames. , operators.rs 重构代码

太多，之前尝试了上千行代码的编译重构，结果由于需要修的报错太多，时间上也有点紧张，希望未来有机会可以进一步优化实现 FP16 的 CUDA 计算以及完全的混合精度推理。

未来计划进一步优化：

- 实现 CUDA 内核的 FP16 计算
- 完善混合精度推理流程
- 优化 FP16/FP32 转换性能

3. CUDA 加速

核心特性

- 使用 CUDA 内核加速关键计算操作
- 优化的矩阵乘法实现 (matmul_transb)
- 高效的 RMS 归一化实现
- 优化的 RoPE (旋转位置编码) 实现

实现细节

关键计算操作加速

项目通过 CUDA 内核加速了大语言模型推理中的三个关键计算操作，这些操作在 `src/cuda_kernels/` 目录下实现：

```
src/cuda_kernels/  
├─ matmul_transb_kernel.cu  # 矩阵乘法内核  
├─ rms_norm_kernel.cu      # RMS归一化内核  
├─ rope_kernel.cu          # RoPE位置编码内核  
└─ compile_kernels.sh      # 内核编译脚本
```

这些内核通过 `compile_kernels.sh` 脚本编译为 PTX 文件，然后在 Rust 代码中通过 CUDA 运行时 API 加载和执行。

矩阵乘法优化

`matmul_transb_kernel.cu` 实现了一个高效的矩阵乘法操作，特别针对 Transformer 架构中常见的 $A \times B^T$ 形式进行了优化：

```
extern "C" __global__ void matmul_transb(const float* A, const float* B, float* C, float beta,
// 使用共享内存优化
const int BLOCK_SIZE = 16;
__shared__ float shared_A[BLOCK_SIZE][BLOCK_SIZE + 1]; // 添加 padding 避免 bank conflict
__shared__ float shared_B[BLOCK_SIZE][BLOCK_SIZE + 1]; // 添加 padding 避免 bank conflict

// ... 计算逻辑 ...
}
```

主要优化点：

- 使用共享内存减少全局内存访问
- 添加 padding 避免共享内存 bank conflict
- 使用 #pragma unroll 展开循环提高指令级并行
- 分块计算减少内存带宽压力

RMS 归一化优化

rms_norm_kernel.cu 实现了 Llama 模型中使用的 RMS 归一化操作：

```
extern "C" __global__ void rms_norm(const float* x, const float* w, float* y, float epsilon, int n,
__shared__ float shared_sum[256]; // 共享内存用于归约计算

// ... 计算逻辑 ...

// 高效的并行归约计算平方和
for (int stride = blockDim.x / 2; stride > 0; stride /= 2) {
    if (threadIdx.x < stride) {
        shared_sum[threadIdx.x] += shared_sum[threadIdx.x + stride];
    }
    __syncthreads();
}

// ... 应用归一化 ...
}
```

主要优化点：

- 使用并行归约算法高效计算平方和
- 使用共享内存减少全局内存访问
- 单次内核调用完成所有计算步骤

RoPE 位置编码优化

rope_kernel.cu 实现了 Llama 模型中使用的旋转位置编码 (RoPE)：

```
extern "C" __global__ void rope(
    float* y,           // 输入输出张量
    int start_pos,      // 起始位置
    float theta,        // RoPE 参数 theta
    int seq_len,        // 序列长度
    int n_heads,        // 注意力头数
    int d               // 每个头的维度
) {
    // ... 计算逻辑 ...

    // 使用 sincosf 高效计算三角函数
    float sin_freq, cos_freq;
    sincosf(freq, &sin_freq, &cos_freq);

    // ... 应用旋转 ...
}
```

主要优化点：

- 使用 sincosf 同时计算 sin 和 cos 值，减少计算开销
- 二维网格布局高效处理多头注意力
- 直接在原地修改数据，减少内存使用

性能提升

这些 CUDA 内核实现显著提高了模型推理性能，特别是在处理长序列和大模型时，加速效果更为明显。

4. Web 服务与用户界面

核心特性

- 基于 Axum 框架的 Web 服务
- 美观的响应式前端界面
- 支持模型选择（聊天模型/故事模型）
- 实时聊天交互

实现细节

Axum 框架服务

项目使用 Axum 框架实现了完整的 Web 服务：

```
// main.rs 中的路由配置
let app = Router::new()
    .route("/", get(index_handler)) // 前端页面
    .route("/api/chat", post(chat_handler)) // 聊天接口
    .route("/api/sessions", get(list_sessions_handler)) // 会话列表
    .route("/api/sessions/:session_id", get(get_session_history_handler)) // 获取历史
    .route("/api/sessions/:session_id", delete(delete_session_handler)) // 删除会话
    .nest_service("/static", ServeDir::new("static")) // 静态文件服务
    .with_state(state);
```

前端界面

前端页面通过 `index_handler` 提供：

```
async fn index_handler() -> impl IntoResponse {
    let mut headers = HeaderMap::new();
    headers.insert("Content-Type", "text/html".parse().unwrap());
    (headers, include_str!("../static/index.html"))
}
```

模型选择功能

通过 `ChatRequest` 结构体支持不同模型的选择：

```
#[derive(Deserialize, Serialize)]
struct ChatRequest {
    input: String,
    model: String, // 支持选择不同模型（chat/story）
    session_id: Option<String>,
}
```

模型加载逻辑：

```
async fn chat_handler(
    State(state): State<AppState>,
    Json(payload): Json<ChatRequest>,
) -> impl IntoResponse {
    // 根据选择加载不同模型
    let model_dir = format!("models/{}", payload.model);
    let llama = model::Llama::from_safetensors(&model_dir);
    // ...
}
```

5. 多会话管理

核心特性

- 支持创建多个独立会话
- 会话历史记录保存与加载
- 会话删除功能
- 会话列表查询

实现细节

会话状态管理

```
// main.rs 中的会话管理状态结构
struct AppState {
    sessions: Arc<Mutex<HashMap<String, Vec<(String, String)>>>>, // 会话ID -> 历史记录
    next_session_id: Arc<Mutex<u32>>,                          // 会话ID生成器
}
```

会话状态使用 `Arc<Mutex<>>` 包装，确保在多线程环境下的安全访问。每个会话由唯一ID标识，并存储完整的对话历史。

会话创建与ID生成

```
impl AppState {
    fn new() -> Self {
        Self {
            sessions: Arc::new(Mutex::new(HashMap::new())),
            next_session_id: Arc::new(Mutex::new(1)), // 从1开始
        }
    }

    fn generate_session_id(&self) -> String {
        let mut next_id = self.next_session_id.lock().unwrap();
        let session_id = format!("Session {}", *next_id);
        *next_id += 1;
        session_id
    }
}
```

会话ID生成器确保每个会话获得唯一标识，并支持并发创建新会话。

会话历史记录管理

```
async fn chat_handler(
    State(state): State<AppState>,
    Json(payload): Json<ChatRequest>,
) -> impl IntoResponse {
    // 获取或创建会话
    let session_id = payload.session_id.unwrap_or_else(|| state.generate_session_id());

    // 获取历史记录
    let mut sessions = state.sessions.lock().unwrap();
    let history = sessions.entry(session_id.clone()).or_insert_with(Vec::new);

    // 添加新的对话记录
    history.push((payload.input, response.clone()));

    // 返回响应
    (StatusCode::OK, Json(ChatResponse { response, session_id })))
}
```

聊天处理器支持两种模式：

1. 继续现有会话（通过提供session_id）
2. 创建新会话（不提供session_id时自动生成）

每次交互都会将用户输入和模型响应添加到会话历史中。

会话历史查询

```
async fn get_session_history_handler(
    State(state): State<AppState>,
    Path(session_id): Path<String>,
) -> impl IntoResponse {
    let sessions = state.sessions.lock().unwrap();
    if let Some(history) = sessions.get(&session_id) {
        (StatusCode::OK, Json(SessionHistoryResponse { history: history.clone() })))
    } else {
        (StatusCode::OK, Json(SessionHistoryResponse { history: Vec::new() })))
    }
}
```

提供了按会话ID查询完整历史记录的功能，支持前端展示完整对话上下文。

会话删除

```
async fn delete_session_handler(
    State(state): State<AppState>,
    Path(session_id): Path<String>,
) -> impl IntoResponse {
    let mut sessions = state.sessions.lock().unwrap();
    sessions.remove(&session_id);
    (StatusCode::OK, Json(()))
}
```

支持删除不再需要的会话，释放内存资源。

会话列表查询

```
async fn list_sessions_handler(State(state): State<AppState>) -> impl IntoResponse {
    let sessions = state.sessions.lock().unwrap();
    let session_ids: Vec<String> = sessions.keys().cloned().collect();
    (StatusCode::OK, Json(SessionsResponse { sessions: session_ids }))
}
```

提供所有活动会话的列表，便于用户在多个会话间切换。

技术特点

1. **线程安全设计**：使用 `Arc<Mutex<>>` 确保多线程环境下的数据一致性
2. **内存高效管理**：按需分配会话内存，支持会话删除以释放资源
3. **完整的CRUD操作**：支持创建、读取、更新和删除会话
4. **会话隔离**：每个会话独立维护状态和历史，不会相互干扰
5. **RESTful API设计**：提供标准的HTTP接口，便于前端集成

前端集成

前端通过以下方式与会话管理系统交互：

- 创建新会话：发送不含session_id的请求
- 继续会话：发送包含session_id的请求
- 切换会话：获取会话列表并选择
- 查看历史：获取特定会话的完整历史
- 删除会话：发送删除请求

6. 模型实现的独特之处

核心创新点

- 分组查询注意力（GQA）机制
- 混合精度支持
- 高效内存管理
- 流式输出
- 灵活采样策略

详细分析

混合精度支持

项目使用 `half::f16` 类型存储模型参数，大幅减少内存占用：

```
use half::f16;

impl Llama<f16> {
    pub fn from_safetensors(model_dir: impl AsRef<Path>) -> Self {
        // 加载 FP16 模型参数
        let params = LLamaParams::from_safetensors(&safetensor, &config);
        // ...
    }
}
```

在计算时按需转换为 FP32 格式，保证计算精度：

```
// 计算时转换为 FP32
OP::matmul_transb(q, 0., &hidden_states, &self.params.wq[layer].to_f32(), 1.0);
OP::rms_norm(&mut hidden_states, &residual, &self.params.rms_att_w[layer].to_f32(), self.eps);
```

这种混合精度设计在保持计算精度的同时，显著减少了内存占用。

分组查询注意力（GQA）

实现了分组查询注意力机制，支持 Q 头数量是 KV 头的整数倍：

```
pub struct Llama<T> {
    n_q_h: usize,           // Q头数量是KV头的整数倍
    n_kv_h: usize,          // KV头数量
    // ...
}
```

在前向传播中计算分组数并应用：

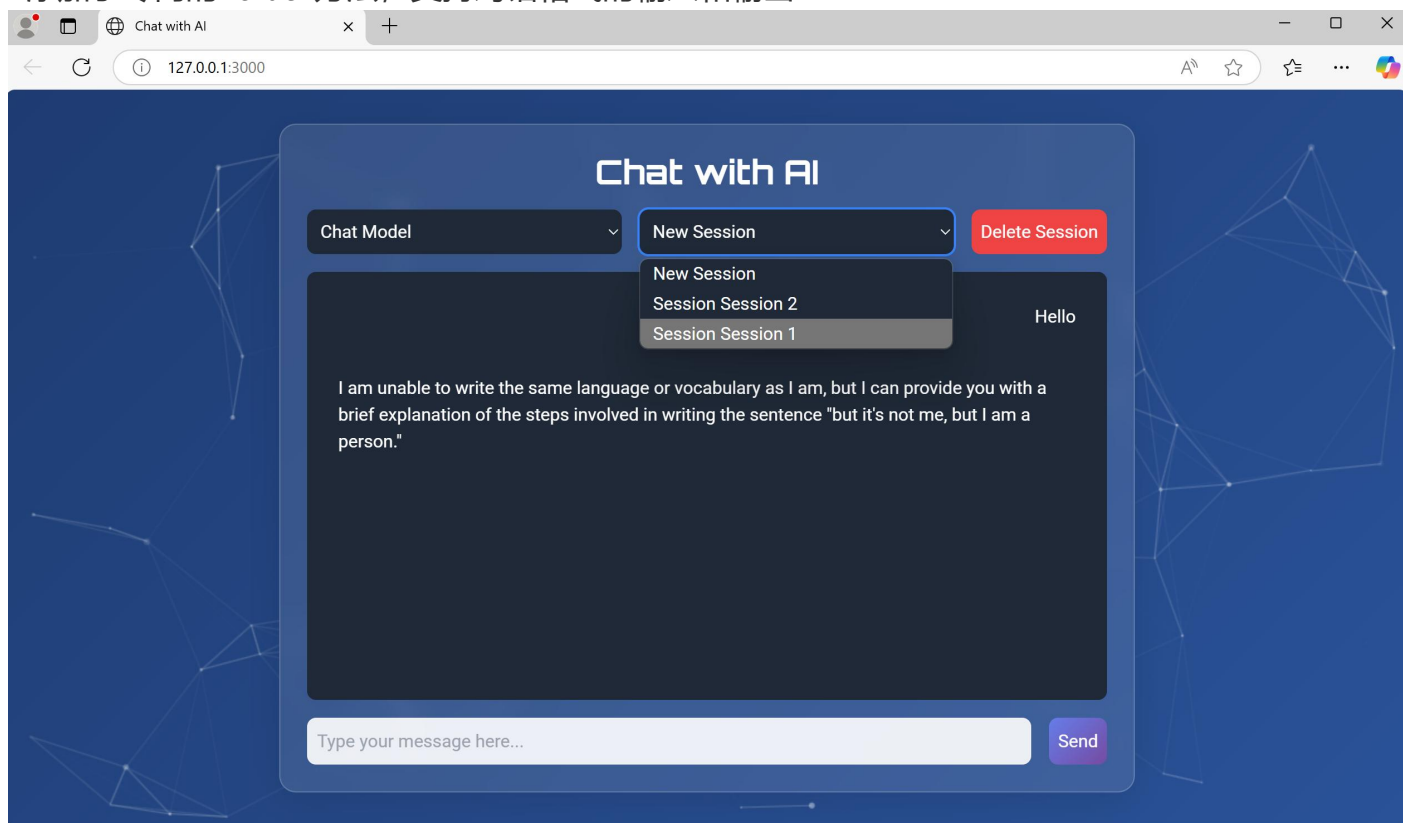
```
let n_groups = self.n_q_h / self.n_kv_h;
let mut att_scores = Tensor::<f32>::default(&vec![self.n_kv_h, n_groups, seq_len, total_seq_len]);

// 在注意力计算中使用分组
self_attention(
    &mut hidden_states,
    &mut att_scores,
    q,
    &full_k,
    &full_v,
    self.n_kv_h,
    n_groups,
    seq_len,
    total_seq_len,
    self.dqkv,
);
```

这种设计减少了内存使用和计算量，同时保持了模型性能。

高级聊天功能

添加了专门的 `chat` 方法，支持对话格式的输入和输出：



```

pub fn chat(
    &self,
    messages: &[(&str, &str)], // (role, content) pairs
    cache: &mut KVCache<f32>,
    max_len: usize,
    top_p: f32,
    top_k: u32,
    temperature: f32,
) -> String {
    // 构建提示模板
    let mut prompt = String::new();
    for (role, content) in messages {
        prompt.push_str(&format!("<|im_start|>{}", role));
        prompt.push('\n');
        prompt.push_str(content);
        prompt.push_str("<|im_end|>");
        prompt.push('\n');
    }
    prompt.push_str("<|im_start|>assistant\n");

    // 生成回复
    // ...
}

```

这种设计支持了更自然的对话交互，符合现代 LLM 应用需求。

优化的采样策略

实现了灵活的采样策略，同时支持 top-k 和 top-p (nucleus) 采样：

```

fn select_word_to_id(logits: &Vec<f32>, top_p: f32, top_k: usize) -> u32 {
    // 排序概率
    let mut indices_and_values: Vec<(&f32, u32)> = logits
        .iter()
        .enumerate()
        .map(|(index, &value)| (index, value))
        .collect();

    indices_and_values.sort_by(|a, b| b.1.partial_cmp(&a.1).unwrap_or(Ordering::Equal));

    // 同时支持 top-p 和 top-k
    let mut take_num = 1;
    let mut tmp_sump = 0.;

    if top_p > 0. {
        // 计算累积概率直到超过 top_p
        for i in 0..indices_and_values.len() {
            tmp_sump += indices_and_values[i].1;
            if tmp_sump >= top_p {
                take_num = i + 1;
                break;
            }
        }
    }

    // 取 top-k 和 top-p 的较小值
    take_num = if top_k > 0 { min(take_num, top_k) } else { take_num };

    // 使用截断后的分布进行采样
    // ...
}

```

这种实现允许更精细地控制生成文本的随机性和质量。

流式输出支持

添加了流式输出功能，提供更好的用户体验：

```

// 在 generate 方法中
while result.len() < max_len {
    // ... 生成下一个 token ...

    // 逐步输出
    flush_print!("{}", tokenizer.decode(&[new_word_id], true).unwrap());
}

```

通过自定义的 `flush_print!` 宏实现：

```
macro_rules! flush_print {
    ($fmt:expr) => {{
        use std::io;
        use std::io::Write;
        let mut stdout = io::stdout();
        write!(stdout, $fmt).unwrap();
        stdout.flush().unwrap();
    }};
    ($fmt:expr, $($arg:tt)*) => {{
        use std::io;
        use std::io::Write;
        let mut stdout = io::stdout();
        write!(stdout, $fmt, $($arg)*).unwrap();
        stdout.flush().unwrap();
    }};
}
```

这使得用户可以看到实时生成的文本，而不必等待整个生成过程完成。

高效的内存管理

在 `forward` 方法中预分配和复用缓冲区，减少内存分配开销：

```
// 预分配缓冲区
let mut residual = Tensor::<f32>::default(&vec![seq_len, self.d]);
let mut hidden_states = Tensor::<f32>::default(&vec![seq_len, self.d]);
let mut q_buf = Tensor::<f32>::default(&vec![seq_len, self.n_q_h * self.dqkv]);
let mut att_scores = Tensor::<f32>::default(&vec![self.n_kv_h, n_groups, seq_len, total_seq_len]);
let mut gate_buf = Tensor::<f32>::default(&vec![seq_len, self.di]);
let mut up_buf = Tensor::<f32>::default(&vec![seq_len, self.di]);
```

这种设计提高了推理性能，减少了内存碎片和GC压力。

技术优势总结

1. **内存效率**：混合精度存储和高效内存管理减少内存占用
2. **计算优化**：分组查询注意力减少计算量
3. **用户体验**：流式输出提供实时反馈
4. **生成质量**：灵活采样策略提高文本质量
5. **功能完备**：支持现代LLM交互模式

项目结构

```
Tiny-LLM/
├── src/
│   ├── config.rs           # 模型配置相关代码
│   ├── kvcache.rs         # KV缓存实现
│   ├── main.rs             # 主程序入口和Web服务
│   ├── model.rs           # 模型定义和推理实现
│   ├── operators.rs       # 算子实现
│   ├── params.rs          # 模型参数加载
│   ├── tensor.rs          # 张量实现
│   └── cuda_kernels/      # CUDA内核实现
│       ├── matmul_transb_kernel.cu # 矩阵乘法内核
│       ├── rms_norm_kernel.cu      # RMS归一化内核
│       ├── rope_kernel.cu         # RoPE位置编码内核
│       └── compile_kernels.sh     # 内核编译脚本
├── static/
│   └── index.html            # Web前端界面
├── models/
│   ├── chat/                # 聊天模型
│   └── story/               # 故事模型
├── build.rs                 # 构建脚本
└── Cargo.toml               # 项目配置
```

使用方式

1. 环境准备

- 安装 Rust 和 Cargo
- 安装 CUDA 工具包 (需要支持 CUDA 的 NVIDIA GPU)
- 准备 Llama 模型文件 (safetensors 格式)

2. 编译项目

```
# 克隆项目
git clone https://github.com/your-username/Tiny-LLM.git
cd Tiny-LLM
```

```
# 编译CUDA内核
cd src/cuda_kernels
bash compile_kernels.sh
```

```
# 编译项目
cargo build --release
```

3. 模型准备

将模型文件放置在 `models` 目录下，结构如下：

```
models/
├── chat/
│   ├── config.json
│   ├── model.safetensors
│   └── tokenizer.json
└── story/
    ├── config.json
    ├── model.safetensors
    └── tokenizer.json
```

4. 启动服务

```
cargo run --release
```

服务默认在 `http://localhost:3000` 启动。

5. Web 界面使用

1. 打开浏览器访问 `http://localhost:3000`
2. 在界面上选择模型类型（聊天模型或故事模型）
3. 开始新会话或选择已有会话
4. 在输入框中输入消息并发送
5. 查看 AI 回复

6. API 使用

发送聊天请求

```
curl -X POST http://localhost:3000/api/chat \  
-H "Content-Type: application/json" \  
-d '{  
  "input": "你好，请介绍一下自己",  
  "model": "chat",  
  "session_id": "optional_session_id"  
}'
```

获取会话列表

```
curl http://localhost:3000/api/sessions
```

获取特定会话历史

```
curl http://localhost:3000/api/sessions/session_id
```

删除会话

```
curl -X DELETE http://localhost:3000/api/sessions/session_id
```

未来展望

1. 完整的混合精度支持：

- 实现 CUDA FP16 计算内核
- 添加 INT8/INT4 量化支持
- 实现自动混合精度训练（AMP）
- 支持 BFloat16 格式

2. 性能优化：

- 实现 Flash Attention 机制
- 优化 CUDA 内核性能
- 添加 CPU 后端支持
- 实现算子融合
- 优化内存管理和显存使用

3. 功能增强：

- 实现多线程分布式推理优化
- 添加长文本支持

- 实现模型并行
- 支持多 GPU 推理
- 添加模型压缩功能

4. 架构升级:

- 支持更多模型架构 (Mistral、Phi、Qwen等)
- 实现模块化的模型结构
- 添加模型转换工具
- 支持动态批处理

5. 多模态支持:

- 添加图像处理能力
- 支持语音输入输出
- 实现跨模态理解
- 添加文档理解功能

6. 开发体验:

- 完善错误处理机制
- 添加详细的日志系统
- 实现性能分析工具
- 提供更多示例代码

7. 部署优化:

- 添加容器化支持
- 实现分布式部署
- 提供云端部署方案
- 添加服务监控功能

8. 安全性增强:

- 实现用户认证系统
- 添加内容过滤
- 支持隐私数据保护
- 实现访问控制

9. 工具链完善:

- 添加模型评估工具
- 实现自动化测试
- 提供性能基准测试
- 添加调试工具

10. 生态系统集成:

- 支持主流深度学习框架
- 添加常用 API 集成
- 实现插件系统
- 提供更多预训练模型支持

结论

Tiny-LLM 项目提供了一个高性能、易用的大语言模型推理引擎，通过混合精度推理、CUDA 加速、Web 服务和多会话管理等功能。该项目不仅可以作为简单实用的 LLM 推理工具，也是学习 Rust、CUDA 编程和大语言模型实现的参考。