

Multi-LLM App - Prototipo Funcional

Descripción del Proyecto

La Multi-LLM App es una aplicación web que permite a los usuarios interactuar con múltiples modelos de lenguaje (LLMs) simultáneamente usando un único prompt. Las respuestas se muestran en formato de mosaicos para facilitar la comparación entre diferentes modelos de IA.

¿Qué he entendido de tu idea?

Tu concepto es crear una herramienta que permita:

- **Entrada única:** Un solo prompt que se envía a múltiples LLMs
- **Comparación visual:** Respuestas mostradas en mosaicos individuales
- **Eficiencia:** Obtener múltiples perspectivas de IA sin tener que escribir el mismo prompt varias veces
- **Flexibilidad:** Poder configurar qué LLMs usar según las necesidades

Arquitectura Implementada

Frontend (React + Vite)

- **Ubicación:** `/home/ubuntu/multi-llm-app/`
- **Tecnologías:** React 18, Vite, Tailwind CSS, shadcn/ui
- **Características:**
 - Interfaz responsiva y moderna
 - Área de texto para el prompt
 - Grid de mosaicos para mostrar respuestas

- Estados de carga y error
- Funcionalidad de copiar texto
- Indicadores visuales de estado

Backend (Flask)

- **Ubicación:** `/home/ubuntu/multi-llm-backend/`
- **Tecnologías:** Flask, Flask-CORS, aiohttp
- **Características:**
 - API RESTful para comunicación con LLMs
 - Soporte para llamadas asíncronas
 - Endpoints para simulación y APIs reales
 - Manejo de errores robusto
 - Configuración modular de LLMs

Funcionalidades Implementadas

✓ Funcionalidades Básicas

- ✓ ~~Interfaz de usuario intuitiva con mosaicos~~
- ✓ ~~Envío de prompts a múltiples LLMs~~
- ✓ ~~Visualización de respuestas en tiempo real~~
- ✓ ~~Estados de carga y completado~~
- ✓ ~~Funcionalidad de copiar respuestas~~
- ✓ ~~Comunicación frontend-backend~~

✓ Funcionalidades Técnicas

- ✓ ~~Arquitectura modular y escalable~~
- ✓ ~~Manejo de errores y timeouts~~
- ✓ ~~Soporte para CORS~~
- ✓ ~~Diseño responsivo~~
- ✓ ~~Simulación de respuestas para desarrollo~~

Funcionalidades Preparadas (Requieren API Keys)

- ☐ Integración real con OpenAI GPT-4
- ☐ Integración real con Google Gemini
- ☐ Integración real con Anthropic Claude
- ☐ Configuración de API keys por usuario

LLMs Soportados

Actualmente Simulados

1. **OpenAI GPT-4** - Simulación completa
2. **Google Gemini** - Simulación completa
3. **Anthropic Claude** - Simulación completa
4. **Meta Llama** - Simulación completa

Preparados para Integración Real

- OpenAI API (GPT-4)
- Google Generative AI API (Gemini)
- Anthropic API (Claude)

Estructura de Archivos

Plain Text

```
/home/ubuntu/
├── multi-llm-app/                                # Frontend React
│   ├── src/
│   │   ├── App.jsx                             # Componente principal
│   │   ├── App.css                             # Estilos
│   │   └── components/ui/                       # Componentes UI
│   ├── index.html                             # HTML principal
│   └── package.json                            # Dependencias
├── multi-llm-backend/                           # Backend Flask
│   ├── src/
│   │   ├── main.py                             # Punto de entrada
│   │   └── routes/
│   │       ├── llm.py                          # Rutas para LLMs
│   │       └── user.py                         # Rutas de usuario
│   ├── requirements.txt                        # Dependencias Python
│   └── venv/                                   # Entorno virtual
└── README_Multi_LLM_App.md                    # Esta documentación
```

Cómo Ejecutar el Proyecto

Prerequisitos

- Node.js 20+
- Python 3.11+
- pnpm (instalado)

Frontend

Bash

```
cd /home/ubuntu/multi-llm-app
pnpm run dev --host
# Acceder a: http://localhost:5173
```

Backend

Bash

```
cd /home/ubuntu/multi-llm-backend
source venv/bin/activate
python src/main.py
# API disponible en: http://localhost:5000
```

Endpoints de la API

`/api/llm/simulate` (POST)

Simula respuestas de múltiples LLMs para desarrollo y demo.

Request:

JSON

```
{
  "prompt": "Tu pregunta aquí",
  "llm_ids": ["openai", "gemini", "claude", "llama"]
}
```

Response:

JSON

```
{
  "success": true,
  "results": [
    {
      "id": "openai",
      "name": "OpenAI GPT-4",
      "status": "completed",
      "response": "Respuesta simulada..."
    }
  ]
}
```

`/api/llm/query` (POST)

Envía prompts a APIs reales de LLMs (requiere API keys).

`/api/llm/available` (GET)

Obtiene la lista de LLMs disponibles.

Próximos Pasos para Implementación Real

1. Configuración de API Keys

- Crear interfaz para gestionar API keys
- Implementar almacenamiento seguro de credenciales
- Validación de API keys

2. Integración con APIs Reales

- Activar endpoint `/api/llm/query`
- Configurar autenticación para cada LLM
- Implementar manejo de límites de tasa

3. Mejoras de UX

- Configuración de LLMs desde la interfaz
- Historial de conversaciones
- Exportación de respuestas
- Temas claro/oscuro

4. Optimizaciones

- Cache de respuestas

- Streaming de respuestas en tiempo real
- Compresión de datos
- Monitoreo de rendimiento

Tecnologías Utilizadas

Frontend

- **React 18:** Framework de UI
- **Vite:** Build tool y dev server
- **Tailwind CSS:** Framework de estilos
- **shadcn/ui:** Componentes UI
- **Lucide React:** Iconos

Backend

- **Flask:** Framework web de Python
- **Flask-CORS:** Soporte para CORS
- **aiohttp:** Cliente HTTP asíncrono
- **SQLAlchemy:** ORM (preparado para uso futuro)

Características del Prototipo

✨ Puntos Fuertes

- **Interfaz moderna y responsiva:** Diseño profesional con Tailwind CSS
- **Arquitectura escalable:** Fácil agregar nuevos LLMs
- **Comunicación asíncrona:** Respuestas simultáneas de múltiples LLMs

- **Manejo de errores:** Estados de error y carga bien definidos
- **Código limpio:** Estructura modular y bien documentada

Casos de Uso

- **Investigación:** Comparar respuestas de diferentes modelos
- **Desarrollo:** Evaluar calidad de respuestas para diferentes tareas
- **Educación:** Mostrar diferentes enfoques de IA
- **Productividad:** Obtener múltiples perspectivas rápidamente

Conclusión

El prototipo está completamente funcional y listo para uso en desarrollo. La arquitectura permite una fácil expansión para agregar más LLMs y funcionalidades. El código está bien estructurado y documentado, facilitando el mantenimiento y las mejoras futuras.

La aplicación demuestra exitosamente el concepto de "múltiples LLMs con un único prompt" y proporciona una base sólida para el desarrollo de una aplicación de producción.