

Aplicaciones LLM

Latencia y velocidad de las LLM Apps

Latencia: velocidad de la LLM App

- Clave: ¿puedes permitirte que el usuario no tenga una experiencia rápida?
 - Normalmente puedes permitirte si los usuarios son los empleados de una empresa.
 - Normalmente NO puedes permitirte si los usuarios son clientes.
- Apps que necesitan low latency (alta velocidad):
 - conversational agents, virtual agents y chatbots.
 - content personalization y recommendation systems.
- Apps que pueden funcionar bien con high latency (baja velocidad):
 - research (legal, market, etc).
 - creative writing y content generation.