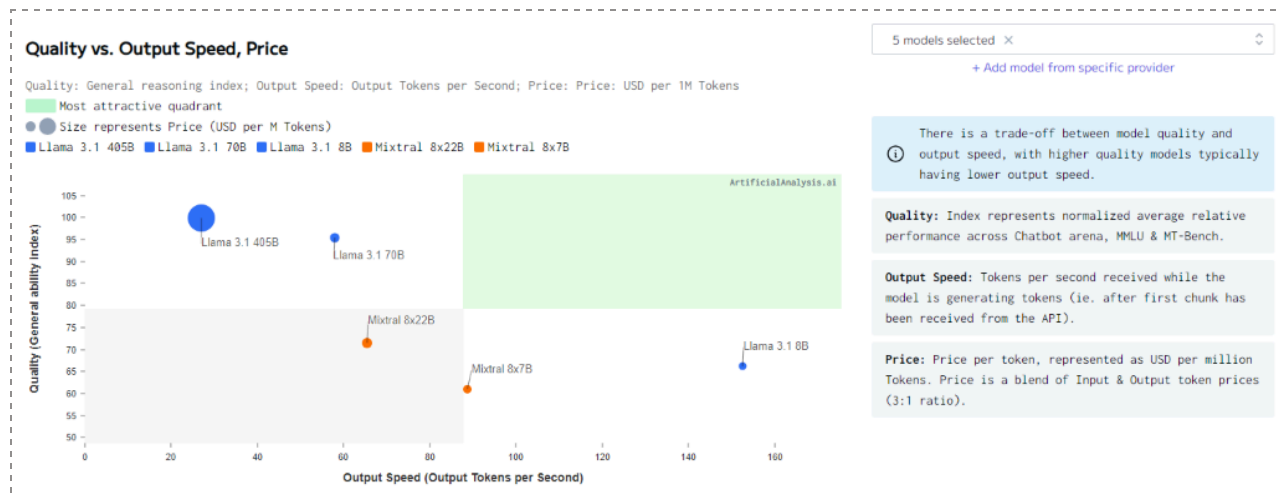# LLM, RAG – Llama Alternative

## Hacker GPT Model

Mixtral 8X22B is a Large Language Model(LLM) developed by Mistral AI that is designed to handle NLP tasks such as text generation, summarization, and conversational AI.

**Advantages of Mixtral:**

- Multilingual capabilities; fluent in English, Spanish, French, German and Italian.

- Offers a sparse mixture of experts(SoME) architecture that reduces computational cost during pre-training.



**Disadvantages of Mixtral 8X7B compared to Llama 3.1 8B**

- Output speed is slower than Llama 3.1 8B.

- Both Mixtral models costs more than Llama 3.1 8B per 1 million tokens.

- Token limit is smaller for Mixtral 8X22B model than any Llama 3 model.

- Requires more RAM vs the Llama 3.1 8B model. 16GB minimum vs 64GB minimum of RAM for optimal performance.

- When scaling up to try and improve accuracy, resource requirement of a 300 GB GPU running Mixtral 8X22B on cloud and consumer software make it very expensive and challenging.

# Context window

Context window: Tokens limit; Higher is better

| Model | Context Window |
|---|---|
| Llama 3.1 405B | 128k |
| Llama 3.1 70B | 128k |
| Llama 3.1 8B | 128k |
| Mixtral 8x22B | 65.4k |
| Mixtral 8x7B | 32.8k |