

ESCUELA POLITÉCNICA NACIONAL
FACULTAD DE INGENIERÍA EN SISTEMAS
METODOS NUMERICOS ICCD412

JUAN FRANCISCO PINTO ANDRANGO

Actividad Extracurricular 04 Costos relacionados a los modelos de Lenguaje

GR1CC

FECHA DE ENTREGA 05 DE NOVIEMBRE DEL 2025

Indicaciones

Investigue sobre las siguientes características para al menos 5 modelos de lenguaje comerciales (i.e. ChatGPT, Claude, Gemini, etc). Cree una tabla resumen

- ¿Qué es inferencia y entrenamiento, cuál es la diferencia?
- Modelo de GPU utilizado/s
- Costo del hardware (costo GPU x número de GPUs) en inferencia y entrenamiento
- Tiempo de entrenamiento
- Consumo energético (watts) en inferencia y entrenamiento.

Que es inferencia y entrenamiento, cual es la diferencia

la inferencia es un proceso con el cual podemos obtener conclusiones, realizar decisiones y interpretaciones con la información disponible como datos, hechos y evidencias. La inferencia es una habilidad importante para la toma de decisiones, resolución de conflictos y deducción de situaciones complejas. La inferencia es una etapa en la que el modelo aplica los conocimientos para entregar una respuesta o generar una predicción.

el entrenamiento son actividades de tipo sistemático y planificado para generar o mejorar capacidades específicas, tiene que ser una etapa en la que un modelo aprende en base a grandes cantidades de datos. Esta etapa es intensa y costosa.

Las inferencias se diferencian del entrenamiento por la aplicación y uso, la respuesta o predicción, un desarrollo en poco tiempo. En el entrenamiento se aprende y construye un modelo, trabaja con grandes volúmenes de datos, tiene que ser costoso y siempre precede a la implementación.

Modelos de lenguaje	GPU utilizados	Costo del Hardware	Tiempo de Entrenamiento	Consumo Energético Entrenamiento	Consumo Energético Inferencia
GPT-4 (OpenAI)	NVIDIA H100 o A1000	500 \$ millones USD	de 60 a 90 Días	25.000 - 50.000 MWh	Alto consumo por token de \$20/1M
DeepSeek-V3 (DeepSeek)	NVIDIA H800 (2.048 unidades)	5-20 \$ millones USD	60 días	10.7000 MWh	Bajo consumo por token de \$1/1 M
Claude 3 Opus (Anthropic)	NVIDIA H100 o A100	200 \$ millones USD	45 - 60 días	15.000 - 30.000 MWh	Alto consumo por token de \$ 75/1 M

| Gemini Ultra (Google)|Google TPUv4/v5 |N/A(Hardware propietario) |60 - 120 días | 30.000 - 60.000 MWh | Optimizado por TPU |Llama 3 70B (Meta/Open)| NVIDIA H100(16.000 unidades)| 10 - 20 millonesUSD|20dias|1.200MWh|Moderapuedeserauto – alojado||Mixtral8x22B(MinstralAI/Open)|NVIDIA H100|5 – 10 millones USD | 14 -21 días | 800 MWh | Bajo (optimizado para MoE) |

link del repositorio de Git-hub

