



**Nombre:** Camila Caicedo

**Curso:** GR1CC

**Fecha de entrega:** 05/11/2025

## Actividad extracurricular 4: Costos relacionados a los modelos de lenguaje pendiente

|                                   | ChatGPT  | Grok  | Gemini  | Copilot   | Titan   |
|-----------------------------------|--|---|---|---|---|
| <b>Inferencia / entrenamiento</b> | <b>Entrenamiento</b><br>o: usa alta precisión y grandes lotes.<br><b>Inferencia:</b> usa baja precisión y es continua.           | <b>Entrenamiento</b><br>o: con datos sintéticos y paralelismo extremo.<br><b>Inferencia:</b> distribuida con baja latencia. | <b>Entrenamiento</b><br>o: multimodal intensivo.<br><b>Inferencia:</b> optimizada para apps móviles y búsqueda. | <b>Entrenamiento</b><br>o: masivo.<br><b>Inferencia:</b> optimizada para productividad (código, texto). | <b>Entrenamiento</b><br>o: con alto paralelismo.<br><b>Inferencia:</b> optimizada para bajo costo y alta velocidad. |
| <b>Modelo de GPU utilizado</b>    | NVIDIA H100, A100  | NVIDIA H100   | Google TPU v4/v5e   | Azure NVIDIA A100/H100 clusters   | AWS Trainium  |
| <b>Costo del hardware</b>         | Aproximadamente 300 – 600 millones USD   | Aproximadamente 200–500 millones USD  | Aproximadamente 10–20 millones USD  | Aproximadamente 900 millones – 3 mil millones USD   | Aproximadamente 30–100 millones USD   |
| <b>Tiempo de entrenamiento</b>    | Varios meses   | Varios meses  | Fine-tuning continuo  | Semanas a meses   | Meses   |
| <b>Consumo energético (watts)</b> | <b>Entrenamiento</b><br>o:<br>Aproximadamente 7–15 MW (30–90 días) → cientos MWh<br><b>Inferencia:</b> ~0.01–0.5 Wh por consulta | <b>Entrenamiento</b><br>o: 20–25 MW (~10–25 GWh).<br><b>Inferencia:</b> MW escala global                                    | <b>Entrenamiento</b><br>o: varios MW (decenas GWh)<br><b>Inferencia:</b> MW total a escala global               | <b>Entrenamiento</b><br>o: MWh por corrida<br><b>Inferencia:</b> Wh por consulta; uso total en MW       | <b>Entrenamiento</b><br>o: MW por pod (MWh totales)<br><b>Inferencia:</b> MW en despliegue global                   |



- DIFERENCIA ENTRE ENTRENAMIENTO E INFERNCE

- **Inferencia:** proceso mediante el cual un modelo de IA entrenado genera nuevos resultados razonando y haciendo predicciones sobre nuevos datos, clasificando las entradas y aplicando el conocimiento aprendido en tiempo real.
- **Entrenamiento:** es la primera fase de un modelo de IA, puede implicar un proceso de ensayo y error, o un proceso de mostrar al modelo ejemplos de las entradas y salidas deseadas, o ambos.
- La inferencia no puede darse sin entrenamiento.

Fuentes de consulta:

- <https://www.cloudflare.com/learning/ai/inference-vs-training/>