



AlexNet ~2012 62.3M params

LLM facts:

- BERT-L - 2018 340M params
- GPT-2 - 2019 1.5B params
- Gemini small - 1B?
- Llama3.1 8B params
- Llama3.1 70B params
- GPT-3 - 2020 175B params
- Llama3.1 408B params
- PaLM - 2022 540B params
- GPT-4o - 2023 1.8T params
- Claude 3.5 Sonnet estimated

Defining the size of model:

- Cost
- Capabilities
- RAM

Model weights 4 bytes per param

Adam optimizer 8 bytes per param

Gradients 4 bytes per param

Activation & temp 8 bytes per param

=> 4 bytes per param + 20 extra bytes per param

=> **20X more RAM** than space needed for the params

Huiku (~28B), Sonnet (~70B), and Opus (~2T)