

## LLM Terms (cont) . . .

- Open weights vs closed weights models: open weight models are those models whose learned weights are shared with everyone, meaning you can download and use it on your device and modify it. Closed weight is the opposite - you can send the LLM a request and it will use the weights locally but you cannot download or access the weights.

### Generation controls (Like temperature).

- Top-p: also called nucleus sampling, is a setting that guides how a LLM picks the next word by sorting all possible next words by probability, it then finds the smallest group of words whose combined chance adds to  $p$  (Ex 0.9) only those words stay in the running. Lower  $p$  keeps likely words (safer) higher  $p$  = creative & more risk of error.
- Frequency penalty: tells LLM "stop repeating yourself". as the model writes it keeps track of each word count. a pos freq penalty values lowers the chance of picking a word again if it has seen it many times while a value of 0 turns the rule off but too high = bad, good = 0.2
- presence penalty: forces model to choose words it has not chosen so far, each time a word already used it picked again the model gets a score cut, higher penalty = bigger cut pushing the model to pick new words while lower penalty lets it repeat words
- Stopping Criteria: tells LLM when to stop so it does not go forever. Common Rules: Max tokens, end of seq token Ex " $\text{\\n\\n}$ " if it hallucinates.
- Max Length: Max tokens a LLM model can use controls the length of output balancing cost, speed