# The Effects of Temperature and Top_p on AI Responses

**Temperature** is a parameter that controls the randomness and creativity of a language model's output. We can think of it as a "creativity thermostat." A low temperature (ex: 0.0 to 0.3) makes the model more confident and deterministic. It will consistently choose the most likely next word in a sequence. This is ideal for tasks that require factual accuracy and predictability, such as summarization or answering questions based on a provided text. On the other hand, a high temperature (ex: >0.7) increases randomness by giving less likely words a higher chance of being selected. This can lead to more diverse, creative, and sometimes surprising responses, making it useful for brainstorming, writing poetry, or creating varied content.

**Top_p**, also known as nucleus sampling, is another method for controlling the model's creativity, but it works by limiting the vocabulary pool. Instead of adjusting the probabilities of all words like temperature does, top_p dynamically selects a small set of the most probable words (the "nucleus") and has the model choose only from that set. For example, a top_p of 0.9 means the model will only consider the most likely words that make up the top 90% of the probability mass for the next word. This prevents the model from choosing highly improbable or nonsensical words, even when randomness is high. It is a common practice to adjust either temperature or top_p, but not both simultaneously, as they are different approaches to controlling the same outcome which is the predictability of the model's response.