

Cheat Sheet: LLM Decoding Strategies

Decoding Strategy	What It Is	Deterministic or Stochastic	Advantages	Disadvantages
Greedy Decoding	Always picks the most probable next token.	<ul style="list-style-type: none"> Deterministic 	<ul style="list-style-type: none"> Efficient; Fast. 	<ul style="list-style-type: none"> Can be repetitive; Predictable.
Multinomial Sampling	Selects the next token based on probability distribution.	<ul style="list-style-type: none"> Stochastic 	<ul style="list-style-type: none"> Promotes diversity; Encourages creativity. 	<ul style="list-style-type: none"> Can produce irrelevant outputs; Unpredictable.
Top-k Sampling	Narrows choices to top-k tokens, then samples based on probability.	<ul style="list-style-type: none"> Stochastic 	<ul style="list-style-type: none"> Balances predictability and diversity. 	<ul style="list-style-type: none"> Risk of repetition; Finding optimal k can be challenging.
Top-p (Nucleus) Sampling	Chooses from a subset of tokens whose combined probability reaches or exceeds a threshold p	<ul style="list-style-type: none"> Stochastic 	<ul style="list-style-type: none"> Reduces randomness; Focuses on plausible outputs. 	<ul style="list-style-type: none"> Complexity in choosing p; May exclude relevant but less probable tokens.
Beam Search	Explores multiple high-probability paths simultaneously.	<ul style="list-style-type: none"> Deterministic 	<ul style="list-style-type: none"> Produces coherent text; Reduces nonsensical outputs. 	<ul style="list-style-type: none"> Computationally intensive; Slower generation times.
Beam Search with Multinomial Sampling	Combines beam and multinomial sampling for token selection within beams.	<ul style="list-style-type: none"> Stochastic 	<ul style="list-style-type: none"> Coherent yet creative outputs; Enhances diversity. 	<ul style="list-style-type: none"> Increased complexity; Slower and resource-intensive.
Contrastive Search	Maximizes differences between top candidate sequences.	<ul style="list-style-type: none"> Deterministic 	<ul style="list-style-type: none"> Boosts diversity; Reduces redundancy. 	<ul style="list-style-type: none"> Higher processing demands; May affect predictability/coherence.