**#1 Redundancy in Test-time Scaling**

(a) Best-of-N

Which is bigger, 3.11 or 3.8?

Okay, let me try to figure out which number is bigger between 3.11 and 3.8, Hmm, so both ... [Answer] 3.8 > 3.11.

Path 1

Which is bigger, 3.11 or 3.8?

Okay, let me think step by step. I need to figure out which number is bigger between 3.11 and 3.8, Hmm, let's start by ... [Answer] 3.8 > 3.11.

Path 2

Which is bigger, 3.11 or 3.8?
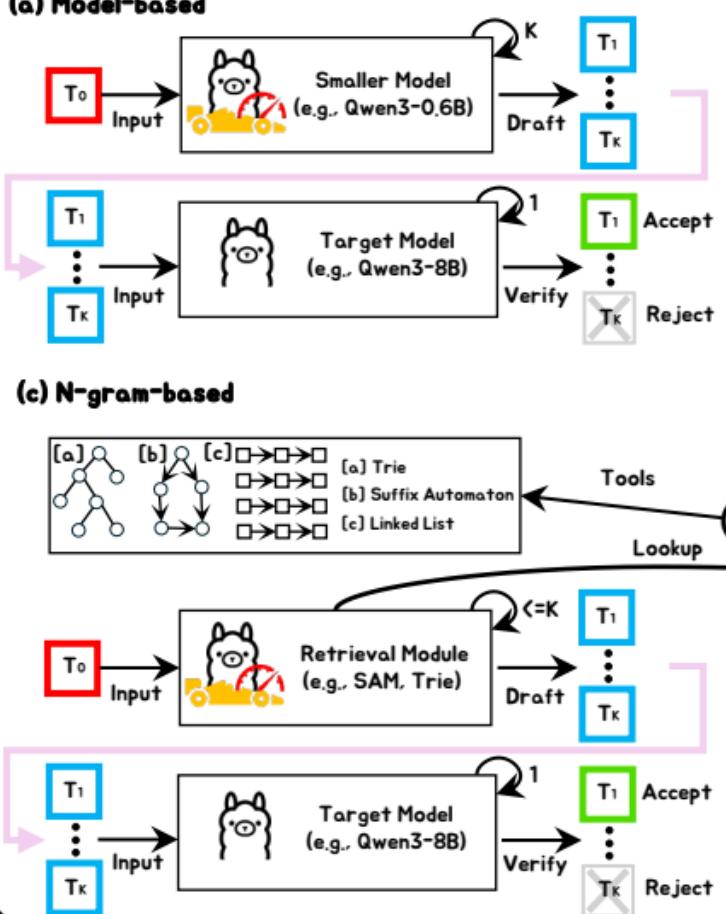
Okay, the first thing I should do is to figure out which number is bigger between 3.11 and 3.8, Hmm, the question... [Answer] 3.8 < 3.11.

Path N

(b) Multi-turn Thinking

Which is bigger, 3.11 or 3.8?

Okay, let me try to figure out which number is bigger between 3.11 and 3.8, Hmm, so both ... [Answer] 3.8 < 3.11.

Which is bigger, 3.11 or 3.8? Your previous answer is: [Answer] 3.8 < 3.11. Please re-answer.

Okay, let me first review my previous answer... Now I will try to figure out which number is bigger between 3.11 and 3.8, Hmm, I find that ... [Answer] 3.8 > 3.11.

Redundancy
Target Model
[Answer] Correct Answer
[Answer] Wrong Answer
Refine

**#2 Speculative Decoding**

(a) Model-based

$T_0$ Input → Smaller Model (e.g., Qwen3-0.6B) →K Draft $T_1$ ... $T_K$

$T_1$ ... $T_K$ Input → Target Model (e.g., Qwen3-8B) →1 Verify → Accept $T_1$ / Reject $T_K$

(b) Training-based

$F_0$ $T_0$ Input → Trainable Module (e.g., EAGLE) →K Draft $T_1$ ... $T_K$

$T_1$ ... $T_K$ Input → Target Model (e.g., Qwen3-8B) →1 Verify → Accept $T_1$ / Reject $T_K$

(c) N-gram-based

[a] Trie  [b] Suffix Automaton  [c] Linked List

Tools

Lookup

Extract

Okay, let me first review my previous answer... Now I will try to figure out which number is bigger between 3.11 and 3.8, Hmm, I find that ... [Answer] 3.8 > 3.11.

Redundancy Pattern

Okay,

try to figure out which number is bigger between 3.11 and 3.8, Hmm,

$T_0$ Input → Retrieval Module (e.g., SAM, Trie) →<=K Draft $T_1$ ... $T_K$

$T_1$ ... $T_K$ Input → Target Model (e.g., Qwen3-8B) →1 Verify → Accept $T_1$ / Reject $T_K$

Draft Model
$T_0$ Token
$F_0$ Feature