# Elastic Speculation Overview

Adaptive Draft Length + Confidence-Based Early Exit (KV Gating)

**Draft Model**

(EAGLE)

Generate K tokens

K draft tokens

**Target Model**

confidence scores

Verify in parallel

Accept prefix

**Early Exit Control**

Monitor draft confidence

If conf < threshold:

**Stop generating & gate KV**

K ∈ {5,10,15}

**Adaptive Control**

Adjust K based on acceptance rate

acceptance rate

gate signal

**KV Cache Writes**

Gate low-confidence

draft tokens

**~50% DRAM reduction**

**Examples:**

High accept → K=15

Low accept → K=5

**⬤ Adaptive Draft Length**

**~30–50% lower latency vs. fixed-length**

Dynamically adjusts speculation depth

**⬤ Confidence-Based Early Exit**

**~50% less KV-cache DRAM traffic**

Minimal latency impact (~1–3%)