



# ICML

International Conference  
On Machine Learning

## Geometry-Aware Decoding with Wasserstein- Regularized Truncation and Mass Penalties for LLMs

Top-W Decoding

Arash Gholami Davoodi  
Seyed Pouyan Mousavi

Navid Rezazadeh  
Pouya Pezeshkpour

# Current samplers: from rank to entropy

Probability-driven rules: each picks a different signal — rank, mass, or entropy.

## Greedy

$\operatorname{argmax}_i p_i$ : one-token confidence

No crop to tune:  
always take the highest-probability token.



## Min-p

$p_i \geq p_{base} * p_{max}$ : adds confidence

Adaptive threshold tied to the leading token; one scalar can miss uncertainty shape.



## What's missing?

All four samplers only see **probability**.

- **Duplicates count twice.**  
answer / solution / reply crowd the crop and starve other meanings.
- **Distinct ideas dropped.**  
a plausible proof or derivation falls just below the threshold and is lost.
- **No notion of distance.**  
the crop can't tell whether removed mass moved 0.1 or 10 in embedding space.

**The cost:** repetition, mode collapse, and lost diversity — even when temperature is raised.

### Next:

read **probability**, but also ask how far the **dropped mass had to travel** — Top-W.

## Top-k

$\operatorname{Top}_k(p)$ : rank cutoff

keep the k highest-probability tokens, Fixed width; simple but ignores confidence and distribution shape.



## Top-H

$\max_S \Gamma_S \quad \text{s.t.} \quad H(q) \leq \alpha H(p)$

**considers distribution uncertainty**

determine a subset  $S \subset V$  from which the next token will be sampled

$$\Gamma_S = \sum_i p_i 1_{v_i \in S}, \quad q_i = \begin{cases} \frac{p_i}{\Gamma_S} & v_i \in S \\ 0 & \text{otherwise} \end{cases}$$

# Why Top- $W$ changes the crop

Probability-only looks only at  $p_i$ . Top- $W$  still respects  $p_i$ , but also asks whether dropped tokens are nearby in embedding space.

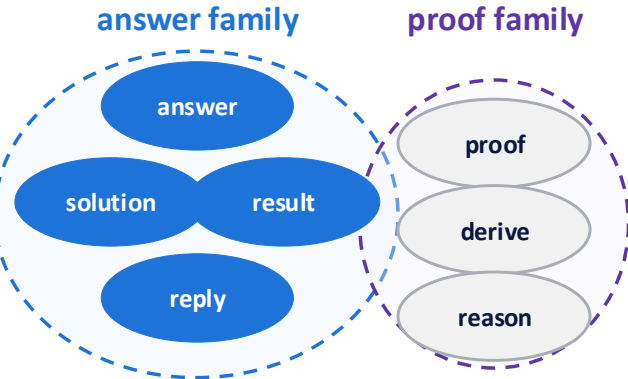
## Probability-only crop

top tokens can all come from one family

same probabilities

answer	0.18	kept
solution	0.14	kept
reply	0.12	kept
result	0.10	kept
proof	0.06	drop
reason	0.04	drop

embedding geometry



Only probability matters, so answer / solution / reply / result can crowd the crop.

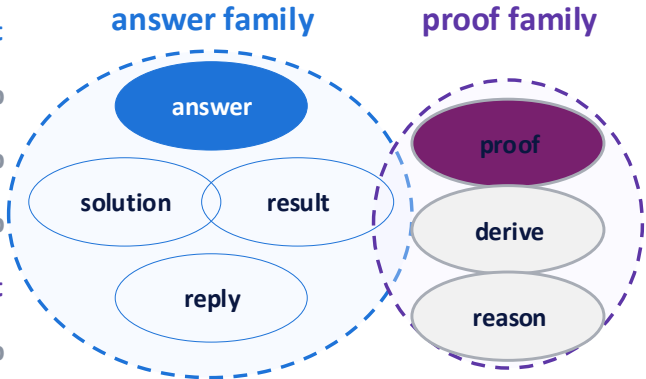
## Top- $W$ crop

geometry can keep coverage across families

same probabilities

answer	0.18	kept
solution	0.14	drop
reply	0.12	drop
result	0.10	drop
proof	0.06	kept
reason	0.04	drop

embedding geometry



Nearby neighbors are cheaper to drop; a distant proof-family token can stay.

# Top- $W$ : Optimization problem

## Top- $W$ objective:

$$\min_S F_{\lambda, \beta}(S) = \min_S \{ W_1(p, q_S) + \lambda H(q_S) - \beta \log \Gamma_S \}$$

$W_1(p, q_S)$ : Geometry cost

Cost of moving dropped mass in token-embedding space.

*Small when dropped tokens are near the kept ones — preserves meaning.*

$\lambda H(q_S)$ : Diversity

Entropy of the renormalized kept set.

*$\lambda$  controls sharpness vs. spread of the crop.*

$-\beta \log \Gamma_S$ : Mass reward

Reward for keeping more probability mass  $\Gamma_S$ .

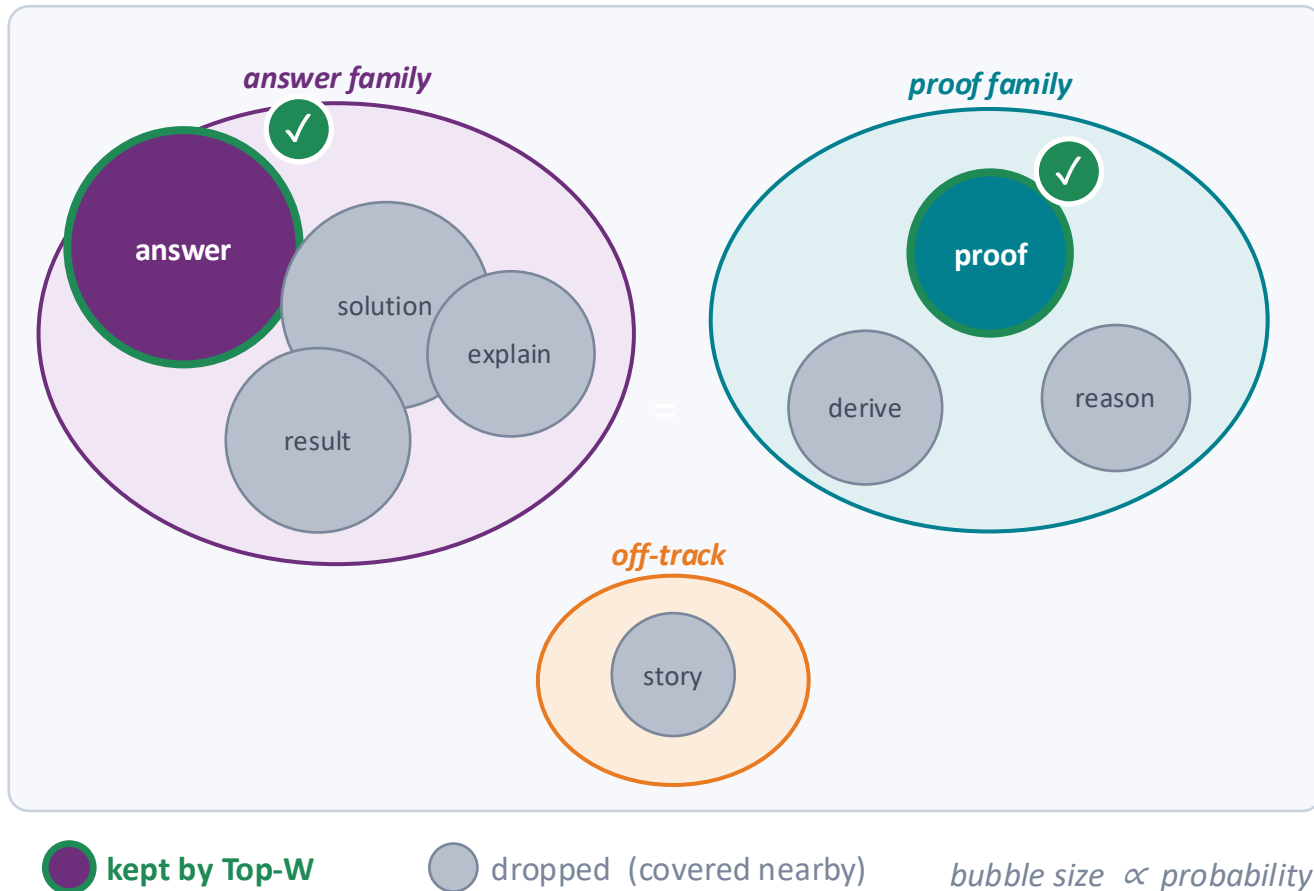
*$\beta$  controls how aggressive the truncation is.*

What we prove → what we get

Optimal crop = a prefix sorted by  $\phi_i = f_i + \lambda \log p_i$  (Theorem 3.4) **O(n) prefix scan**

# Top- $W$ : what the algorithm changes

Top- $W$  reads probability and embedding geometry — keeps tokens that **span distinct meanings**.



## 1 Score each token

$$\phi_i = f_i + \lambda \log p_i$$

f = embedding-distance term; p = probability.

## 2 Keep the top prefix

Sort tokens by  $\phi_i$ , the optimal subset, i.e., crop  $S^\star$  is a prefix (Thm 3.4).

## 3 Sample as usual

Outside  $S^\star \rightarrow -\infty$  logits; standard renormalize and sample inside.

**Not** a new distribution — a **geometry-aware crop mask**.

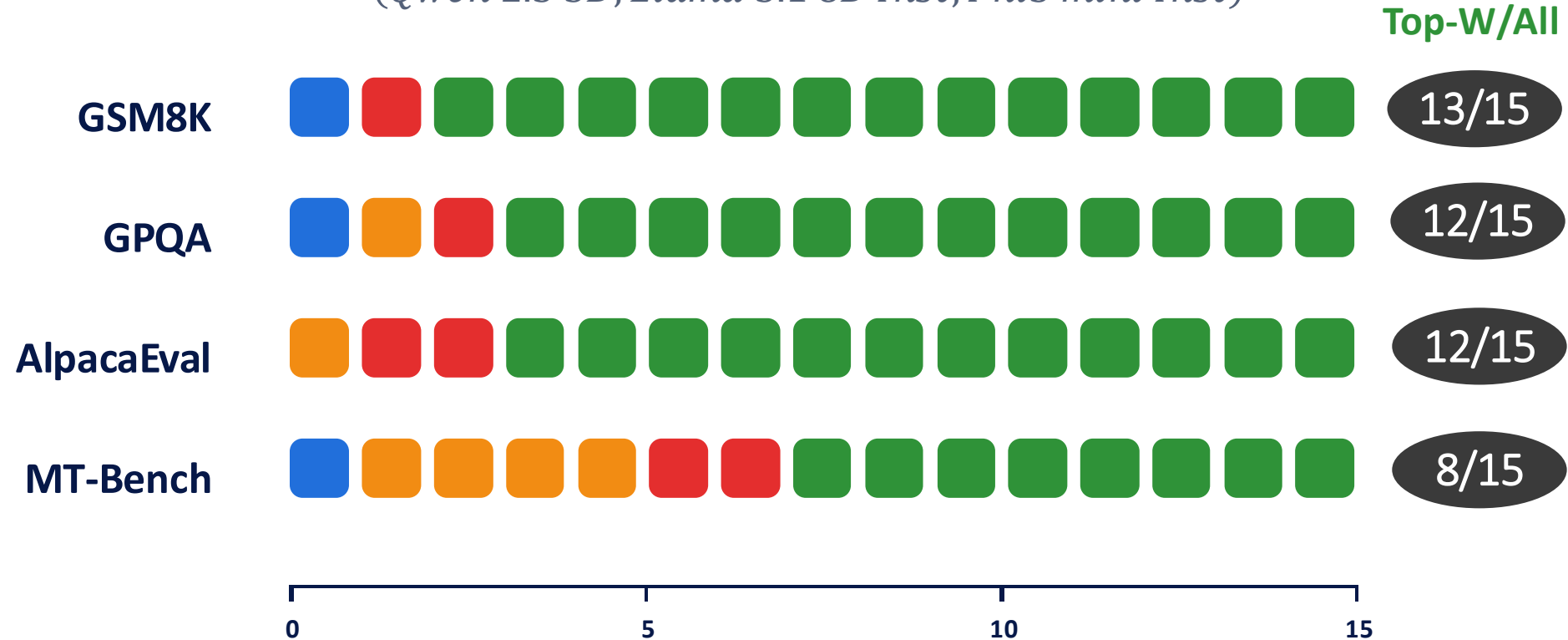
# Benchmark wins: each block is one setting

winning settings out of 15 tuples of

5 temperatures  $T \in (0.7, 1.0, 1.5, 2.0) \times$

3 models  $M \in$

*(Qwen 2.5 3B, Llama 3.1 8B Inst, Phi3 mini Inst)*



Higher temperature  $\Rightarrow$  larger Top-W gain (up to 33.7% gain).

Min-p and Top-p collapse as T rises (Top-p drops to single digits on GSM8K at T=2.0); Top-W stays stable across  $T \in \{0.5, 0.7, 1.0, 1.5, 2.0\}$ .

# Beyond accuracy: judged creative writing

LLM-as-a-judge rubric — diversity, originality, narrative flow, emotional impact, imagery.

CREATIVE WRITING — JUDGE WINS

12 / 27

(LLM, T, prompt) triplets won by Top-W with higher  $\beta$ .

*Min-p : 6   Top-p : 5   Top-H : 5*

*Larger  $\beta$  → broader support → more creative rubric gains.*

Creative-writing rubric wins (higher- $\beta$  Top-W)

