

Order Scaling: UM Beats External KN-6 at Low Threshold

Claude and MJC

March 2026

Abstract

We show the UM KN chain with retroactive training *beats* external KN-6 at 100K enwik9 when the threshold is set to $\tau = 2$. Frozen bpc: 2.656 vs 2.719 (-0.063). The retroactive pass gives the UM a genuine advantage: every neuron sees the full dataset twice, producing better-calibrated distributions than single-pass online counting. Threshold is the decisive parameter: $\tau = 2$ creates 6847 neurons (broad coverage), while $\tau = 4$ creates only 2345 neurons, leaving a $+0.030$ gap. No order degradation is observed at any threshold or order tested.

1 Order Scaling at $\tau = 4$

Order	UM frozen	Ext. KN	Δ	Neurons
2	3.723	3.124	+0.599	0
3	2.752	2.756	-0.004	672
4	2.745	2.674	+0.071	2149
5	2.746	2.691	+0.055	2304
6	2.749	2.719	+0.030	2345

No degradation at higher orders. The UM plateaus at order 3, gaining only 0.007 bpc from orders 3 to 6.

2 Threshold Sweep at Order 6

τ	UM frozen	Ext. KN-6	Δ	Neurons
2	2.656	2.719	-0.063	6847
3	2.696	2.719	-0.023	4024
4	2.749	2.719	+0.030	2345
6	2.908	2.719	+0.189	630
8	3.181	2.719	+0.462	160

The threshold controls how many contexts are captured as neurons. At $\tau = 2$, the UM creates enough neurons to cover the contexts that matter, and the retroactive pass produces superior distributions. At $\tau = 4$, too many contexts are missed, and the KN chain falls back to lower-order predictions at those positions.

3 Why $\tau = 2$ Wins

External KN counts *every* context from the first observation. A context seen twice gets a KN distribution with $n = 2$, $n_{\text{types}} \leq 2$. This is already informative: if both observations had the same continuation, the KN distribution assigns high probability to it.

The UM with $\tau = 2$ creates a neuron as soon as a context appears twice, matching KN’s coverage. The retroactive pass then gives the neuron *two full passes* over the data, accumulating $2\times$ the counts of online KN. More counts = better-calibrated KN distributions = lower surprise.

At $\tau = 4$, contexts seen 2-3 times never get neurons. These positions fall back to lower-order predictions, adding loss. At $\tau = 8$, only 160 neurons are created and the model is essentially a bigram.

4 The Retroactive Advantage

Why does the UM beat external KN at the same order? The answer is the retroactive pass:

Protocol	100K, order 6	Note
External KN online	2.719	1× data, online
UM online (pass 1)	3.190	1× data, online
UM retro (pass 2)	—	2× data, frozen structure
UM frozen (pass 3)	2.656	3× counts, frozen

The UM’s online score (3.190) is much worse than external KN (2.719) because the UM uses global byte counts for the unigram base and the threshold delays neuron creation. But after the retroactive pass, each neuron has accumulated $2\times$ the observations, and the frozen distribution is better-calibrated than online.

5 Scaling with Data Size ($\tau = 1$)

Size	UM $\tau=1$	Ext. KN-6	Δ	Neurons
10K	1.878	3.356	-1.478	~1K
50K	2.508	2.928	-0.420	~8K
100K	2.620	2.719	-0.099	15368
200K	2.610	2.511	+0.099	25350

The UM advantage *decreases with data size* and crosses over between 100K and 200K. The retroactive pass helps most on small data ($2\times$ counts on 10K data is massive; on 200K it is relatively less impactful). External KN improves rapidly because each new observation adds to ALL matching context counts online, while the UM can only accumulate counts for neurons that already exist.

Why the UM improves slowly. From 100K to 200K, the UM frozen score improves from 2.620 to 2.610 (-0.010). External KN improves from 2.719 to 2.511 (-0.208). The UM’s slower improvement comes from the frozen scoring penalty: the global distribution becomes a worse local fit as data size increases (non-stationarity grows with scale).

6 Implications

For compression. Below 200K, the UM with retroactive training beats KN-6 for arithmetic coding. At 100K, the advantage is 0.099 bpc ($\tau = 1$). Above 200K, external KN is better. The crossover point will shift with threshold (lower τ) and number of retroactive passes.

For the UM. The retroactive pass is a genuine but scale-limited advantage. Its benefit is largest when data is scarce (cold-start correction, similar to the H3 meta-combination finding). At large scale, the frozen scoring penalty dominates.

For scaling beyond the crossover. Three paths:

1. **Online scoring.** Use the UM’s online pass score directly, avoiding the frozen penalty. Currently online KN is flat (~ 3.15 at both 100K and 1M) because the global unigram base anchors it. Fix: use online byte counts.
2. **Multiple retroactive passes.** Each pass gives the model another look at the data. With k passes, counts are $k\times$. The frozen distribution improves with each pass, extending the crossover point.
3. **Block-local scoring.** Score frozen within blocks (e.g., 100K at a time), then stitch. This keeps the UM in its advantage zone while scaling to arbitrary data sizes.

7 Conclusion

The UM KN chain at $\tau = 1$, order 6, beats external KN-6 at scales up to ~ 150 K, with a crossover between 100K and 200K. Below the crossover, the retroactive pass provides a genuine advantage ($2\times$ counts). Above it, the frozen scoring penalty on non-stationary data exceeds the retroactive benefit. No order degradation is observed at any setting. The path to large-scale UM compression requires online scoring or multi-pass training to extend the crossover point.