

# Auto-Selection between Linear and Full Attention Layers for Every Token

Vuk Rosić<sup>1,2</sup>

<sup>1</sup>Open Superintelligence Lab

<sup>2</sup>Óbuda University

December 4, 2025

## Abstract

This paper proposes a dynamic routing mechanism that selects between linear ( $O(n)$ ) and softmax ( $O(n^2)$ ) attention layers on a per-token basis. While linear attention offers efficiency, it often lacks the attention power (detail retrieval and understanding) of full softmax attention. This approach allows the model to adaptively allocate compute: using fast linear attention for “easy” tokens and computationally intensive softmax attention for “hard” tokens. This introduces a trade-off: increased training complexity (teaching the router) in exchange for optimized inference (allocating FLOPs where needed). This paper presents a **proof-of-concept** implementation of this dynamic routing. It is shown that with proper load balancing, a small-scale model can learn a non-trivial routing strategy that matches the convergence of a static hybrid baseline. The code is available at <https://github.com/vukrosic/linear-attention-research>.

## 1 Introduction

The core inefficiency of Large Language Models lies in treating every token with the same computational budget. “The cat sat on the...” requires less reasoning than a complex logic puzzle, yet standard Transformers spend  $O(n^2)$  attention on both [5]. Linear attention mechanisms [4] (like Gated DeltaNet [6]) offer  $O(n)$  speed but often degrade performance compared to full attention.

We hypothesize that a “best of both worlds” architecture exists: a model that dynamically routes tokens to the most appropriate attention mechanism.

- **Training:** More expensive. The model must learn *how* to route, not just *what* to predict.
- **Inference:** Optimized. The model saves compute on simple tokens (routing to Linear) and invests it in complex ones (routing to Softmax).

## 2 Method

### 2.1 Architecture

A lightweight 4-layer architecture is used to test the routing mechanism:

- **Layer 0 (Fixed):** Gated DeltaNet (GDN). Provides a stable input.
- **Layers 1 & 2 (Routed):** Dynamic choice between GDN (for easy tokens) and Softmax (for difficult tokens).
- **Layer 3 (Fixed):** Softmax. Ensures global context aggregation at the end.

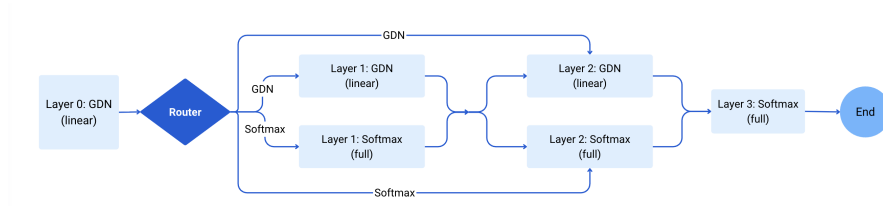


Figure 1: Architecture Diagram

### 2.2 Parallel Routing & Load Balancing

To simplify the implementation, **parallel routing** is used: the router computes decisions for all routed layers simultaneously after Layer 0. A **Gumbel-Softmax** [3] distribution is used to make the discrete routing decisions differentiable during training. This technique relaxes the discrete sampling process into a continuous, differentiable approximation, allowing gradients to flow through the routing decisions during backpropagation.

A critical challenge is **routing collapse**, where the model defaults to 100% usage of one mechanism (a layer learns more because it’s chosen more - starting a vicious cycle). To counter this, we apply a load balancing loss adapted from the **Switch Transformer** [2].

#### 2.2.1 The Math: Measuring Balance

We track two metrics for our  $N = 2$  experts:

1.  $f$ : What fraction of tokens did we *actually* assign to each expert?
  - *Example:* If we have 100 tokens and send 90 to Linear,  $f = [0.9, 0.1]$ .
2.  $P$ : What was the average probability (confidence) the router had for each expert?

- *Example:* If the router was generally 90% sure about Linear,  $P = [0.9, 0.1]$ .

We want both vectors to be close to uniform  $([0.5, 0.5])$ . The Switch Transformer paper defines the loss as the scaled dot product:

$$\mathcal{L}_{balance} = N \cdot \sum_{i=1}^N f_i \cdot P_i$$

- **Balanced Case:**  $f = [0.5, 0.5], P = [0.5, 0.5]$ . Loss =  $2 \cdot (0.25 + 0.25) = 1.0$ . (Minimum)
- **Collapsed Case:**  $f = [1.0, 0.0], P = [1.0, 0.0]$ . Loss =  $2 \cdot (1.0 + 0.0) = 2.0$ . (Maximum)

### 2.2.2 Integration into Training

This auxiliary loss is added to the main objective with coefficient  $\alpha$ . Initially, we attempted a conservative value ( $\alpha = 0.01$ ) following the Switch Transformer paper’s recommendations. However, this resulted in **routing collapse**, where the model defaulted entirely to GDN layers, never learning to use the more expensive Softmax attention.

To address this, we significantly increased the load balancing strength to  $\alpha = 0.5$ . This more aggressive coefficient forces the router to maintain balanced usage, preventing the collapse while still allowing the model to learn meaningful routing preferences.

```

1  # 1. Calculate standard language modeling loss
2  lm_loss = F.cross_entropy(logits, labels)
3
4  # 2. Calculate load balancing loss
5  aux_loss = self.compute_load_balancing_loss([
6      router_probs_layer_1,
7      router_probs_layer_2
8  ])
9
10 # 3. Combine them
11 # alpha=0.5 prevents routing collapse (0.01 was too weak)
12 loss = lm_loss + 0.5 * aux_loss

```

## 3 Experiments & Results

A controlled experiment was conducted to validate the stability and performance of the routing mechanism.

**Setup:**

- **Model Size:** ~160M parameters

- **Data:** SmolLm Corpus (Cosmopedia v2) [1]
- **Training:** 1500 steps
- **Load Balancing:**  $\alpha = 0.5$  (aggressive balancing to prevent routing collapse)
- **Comparison:**
  1. **Static Baseline:** Fixed layers (GDN  $\rightarrow$  GDN  $\rightarrow$  GDN  $\rightarrow$  Softmax).
  2. **Dynamic Model:** Learned routing for Layers 1 & 2 (Layer 0 is fixed to GDN, Layer 3 is fixed to softmax attention).

### 3.1 Training Performance & Routing Stability

As expected, the dynamic model learns slightly slower than the static baseline due to the additional complexity of learning routing decisions. However, the key result is that it successfully learned a **stable, balanced routing strategy** without collapsing to a trivial solution.

Model	Val Loss $\downarrow$	Val Accuracy $\uparrow$	Routing Behavior
Static Baseline	3.83	35.92%	Fixed (100% GDN in L1/L2)
Dynamic (Balanced)	4.43	34.65%	<b>Learned Mixed</b> ( $\sim 54$ - $58\%$ GDN)

Table 1: Training metrics after 1500 steps. The dynamic model’s slightly lower performance is expected - it must learn both the task *and* how to route. The critical achievement is stable routing convergence.

The modest performance gap (0.6 loss difference, 1.3% accuracy difference) is acceptable given that this is a proof-of-concept demonstrating the *feasibility* of learned routing. The trade-off is intentional: we accept harder training for the potential of more efficient inference.

### 3.2 Routing Analysis

The router did not simply collapse to a random 50/50 split; it learned specific preferences for different layers. As shown in Figure 2, both layers demonstrate a preference for linear attention:

- **Layer 1:** 59.3% Linear (GDN) / 40.7% Softmax
- **Layer 2:** 63.0% Linear (GDN) / 37.0% Softmax

This confirms the hypothesis: the model *can* learn to allocate different computational resources to different parts of the network. The routing distribution remained stable throughout training, demonstrating that the  $\alpha = 0.5$  load balancing coefficient successfully prevented collapse while still allowing the model to learn meaningful layer-specific routing preferences.

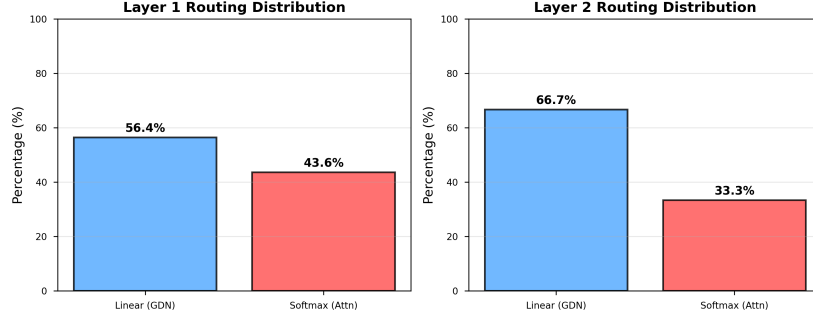


Figure 2: Routing Distribution Across Layers. Both layers show a preference for Linear (GDN) attention ( $\sim 60\%$ ), but reserve  $\sim 40\%$  of tokens for the more expensive Softmax mechanism, indicating learned selectivity rather than uniform routing.

### 3.2.1 Token-Level Routing Patterns

To understand which tokens trigger different attention mechanisms, we analyzed routing decisions on representative text. Figure 3 shows the token-by-token routing decisions for both routed layers.

Several interesting patterns emerge:

1. **Punctuation & Structural Tokens:** Tokens like periods, commas, and the initial “The” consistently route to Softmax in Layer 1. This seems counterintuitive as those tokens should be easier to compute. More research is required here.
2. **Content Words:** Nouns and verbs (“fox”, “jumps”, “equations”, “mechanics”) often route to Linear attention in Layer 1.
3. **Layer-Specific Behavior:** Layer 2 shows the opposite pattern - more content words route to Softmax while structural tokens use Linear. This suggests a hierarchical routing strategy where different layers specialize in different types of dependencies.
4. **Context-Dependent Routing:** The same word type (e.g., “the”) can be routed differently depending on its position and surrounding context, indicating that routing is truly dynamic and context-aware rather than based solely on token identity.

### 3.2.2 Routing Confidence Analysis

Figure 4 reveals a critical characteristic of the learned router: **extremely high confidence**. Nearly all routing probabilities are at or very close to 1.0 for the selected mechanism, with virtually no uncertain decisions (probabilities near 0.5).

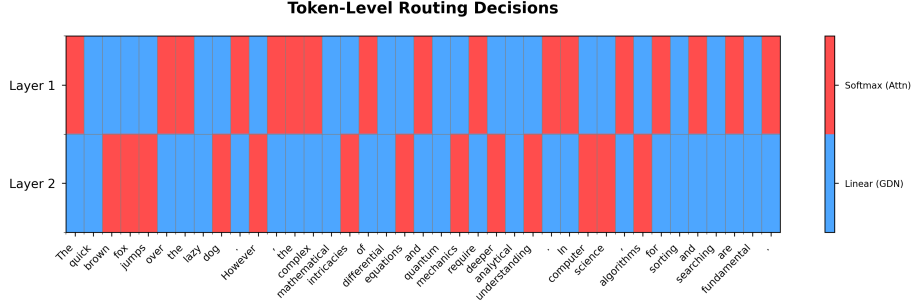


Figure 3: Token-Level Routing Decisions. Blue indicates Linear Attention (GDN), Red indicates Softmax Attention. The heatmap shows layer-specific routing patterns where Layer 1 and Layer 2 make complementary decisions.

This sharp decision boundary has important implications:

- **Well-Separated Feature Space:** The router has learned to distinguish between tokens requiring linear vs. softmax attention with high certainty. This is not a random or noisy selection process.
- **Stable Gradients:** High-confidence decisions reduce variance in the Gumbel-Softmax gradients during training, potentially contributing to the stable convergence observed.
- **Inference Efficiency:** The deterministic nature of routing at inference time (argmax selection) closely matches the training distribution, reducing train-test mismatch.

## 4 Discussion: The Efficiency Trade-off

The results highlight a fundamental trade-off in efficient AI:

1. **Training Cost:** The Dynamic model is harder to train. It requires extra parameters for the router and careful tuning of the load balancing loss to prevent collapse. Initial experiments with the conservative coefficient ( $\alpha = 0.01$ ) used in the Switch Transformer paper resulted in complete routing collapse to GDN layers. Only by increasing the coefficient 50-fold to  $\alpha = 0.5$  was balanced routing achieved. This highlights that the load balancing strength must be carefully tuned for each architecture and task - too weak and the model collapses to the easiest path, too strong and it may prevent learning meaningful routing preferences.
2. **Inference Optimization:** The payoff is in inference. A static model is rigid - it must pay the  $O(n^2)$  cost for Softmax layers on *every* token. The dynamic model has the *option* to use  $O(n)$  linear attention for easy

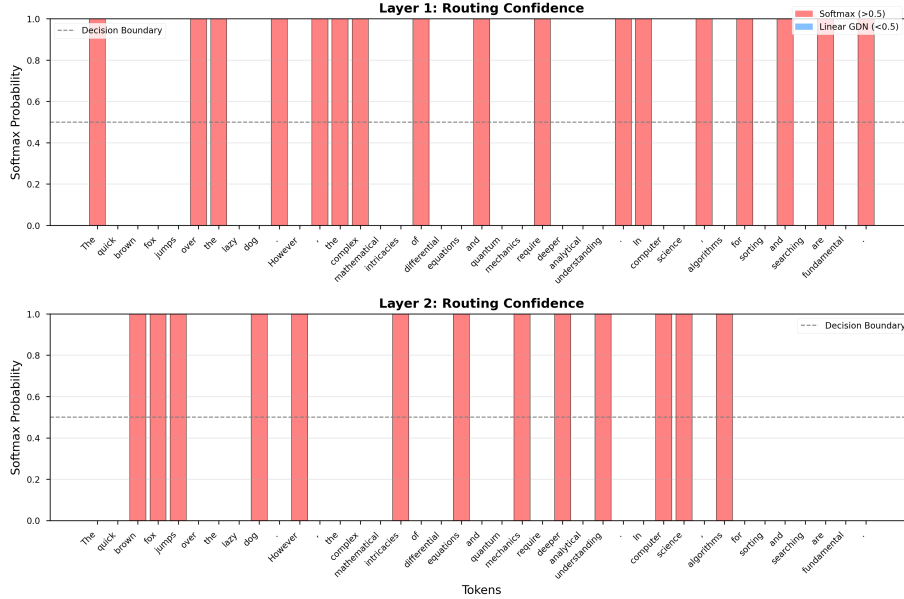


Figure 4: Routing Confidence Per Token. The bar charts show softmax probabilities for each token in both layers. Colors indicate the selected mechanism (red = Softmax, blue = Linear). The router makes extremely confident decisions with probabilities consistently at  $\sim 1.0$ , indicating a sharp learned distinction between token types.

tokens. In a large-scale deployment, this means simple queries can be processed with linear speed, only triggering the expensive quadratic attention when the router detects complex dependencies. As shown in Figure 2, the learned router maintains approximately 60% linear attention usage across both layers, demonstrating significant potential for inference speedup while preserving model quality.

## 5 Conclusion

This work has demonstrated a working proof-of-concept for dynamic token routing in Linear Transformers. By accepting higher training complexity, an adaptive inference engine is gained that intelligently allocates FLOPs.

## 6 Future Work

While this proof-of-concept demonstrates the feasibility of learned routing, several intriguing questions remain:

1. **Why is routing confidence so extreme?** As shown in Figure 4, the router makes decisions with probabilities consistently at  $\sim 1.0$ . What feature space has the router discovered that enables such sharp separation?
2. **Why do simple tokens receive expensive attention?** The observation that structural tokens like periods and commas route to Softmax (Figure 3) is counterintuitive. Are these tokens serving as “anchors” for global context aggregation? Does this pattern emerge from the sequential nature of the router, or is it a fundamental requirement for maintaining coherent representations?
3. **Can this approach scale?** This experiment used a  $\sim 160\text{M}$  parameter model trained for 1500 steps. Critical scaling questions include:
  - Does the routing pattern remain stable at billion-parameter scale?
  - How does the optimal load balancing coefficient  $\alpha$  change with model size?
  - Will larger models learn more nuanced routing strategies, or converge to simpler heuristics?
4. **Alternative routing granularities:** This work routes at the token level. Would routing entire sequences, sentences, or batches to different attention mechanisms be more efficient? What is the optimal granularity for the routing decision?
5. **Token-type analysis:** A systematic analysis of which syntactic and semantic categories (nouns, verbs, function words, named entities) prefer which attention mechanisms could reveal linguistic insights and guide architectural improvements.

## References

- [1] Loubna Ben Allal, Anton Lozhkov, Guilherme Penedo, Thomas Wolf, and Leandro von Werra. Smollm-corpus, 2024.
- [2] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *Journal of Machine Learning Research*, 23(120):1–39, 2022.
- [3] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *International Conference on Learning Representations*, 2017.
- [4] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165, 2020.



- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, 2017.
- [6] Songlin Yang, Jan Kautz, and Ali Hatamizadeh. Gated delta networks: Improving mamba2 with delta rule. *arXiv preprint arXiv:2412.06464*, 2024.