

A Spike-Gated Residual Unit for Information Flow Control in Transformers

Anonymous CVPR submission

Paper ID 25

Abstract

001 *Large Language Models deliver state-of-the-art perfor-*
002 *mance but incur substantial computational and energy*
003 *costs due to their dense, fixed-depth Transformer archi-*
004 *tectures. We introduce a spike-gated residual unit that*
005 *enables dynamic depth adaptation in pretrained Trans-*
006 *formers without modifying backbone weights. By insert-*
007 *ing lightweight spike-gated unit into each residual block*
008 *and optimizing a sparsity-aware objective with a con-*
009 *trollable regularization coefficient $\lambda_{sparsity}$, the model*
010 *learns to selectively bypass redundant layers during infer-*
011 *ence. Using a gate-only fine-tuning strategy on TinyLlama-*
012 *1.1B, our unit reduces average active depth to 76.5%,*
013 *achieving a $1.31\times$ theoretical efficiency gain while main-*
014 *taining competitive performance ($0.93\times$ of the original*
015 *zero-shot HellaSwag accuracy compared to the dense*
016 *TinyLlama-1.1B baseline). Layer-wise analysis reveals*
017 *an emergent hierarchy in which early and late layers re-*
018 *main critical while intermediate blocks are selectively sup-*
019 *pressed under sparsity pressure. These results demon-*
020 *strate that spike-gated residual unit provides a practical*
021 *software-level pathway toward energy-aware Transformer*
022 *inference, enabling controllable trade-offs along an effi-*
023 *ciency–intelligence Pareto frontier without requiring full re-*
024 *training or specialized neuromorphic hardware. To support*
025 *transparency and reproducibility, the complete implemen-*
026 *tation is available at [https://anonymous.4open.](https://anonymous.4open.science/r/Gated-Residual-Unit-7BC5/)*
027 *science/r/Gated-Residual-Unit-7BC5/.*

028 1. Introduction

029 Large Language Models (LLMs) have revolutionized AI but
030 suffer from extreme computational and energy demands.
031 For example, training the 175-billion-parameter GPT-3 re-
032 quired on the order of 175 years on a single high-end GPU.
033 Modern LLMs drive unprecedented capabilities, but their
034 dense, feedforward design means every layer and neuron is
035 active for each token, leading to “wasted computation” on
036 easy inputs and a critical electricity demand for training and
037 inference [8, 13]. This rising energy burden is now recog-

nized as a fundamental challenge: as [13] notes, the “im- 038
ense electricity demand” of LLMs is an “urgent” obstacle 039
to sustainable deployment. By contrast, the human brain 040
operates on only ~ 20 W of power and uses sparse, event- 041
driven spike signals. Spiking Neural Networks (SNNs) seek 042
to mimic this efficiency: when implemented on neuromor- 043
phic hardware, SNNs can replace dense multiply-add oper- 044
ations with binary spike computations [28]. SNNs are “an 045
energy efficient alternative to traditional neural networks” 046
under event-driven execution [11]. Driven by these insights, 047
we explore a paradigm shift toward sparse, spiking commu- 048
nication within the Transformer block, aiming to achieve a 049
significant reduction in operational energy footprints while 050
maintaining a high ceiling for linguistic performance. 051

Recent works have begun exploring spiking Transform- 052
ers as a path to efficiency. In General Language Un- 053
derstanding Evaluation (GLUE) benchmarks, [26] notes 054
that “SNNs have emerged as a promising candidate for 055
energy-efficient LLM inference,” and introduces techniques 056
such as masked spike encodings and analog compute-in- 057
memory that improve accuracy while reducing energy con- 058
sumption by more than $2\times$ on standard tasks. Simi- 059
larly, specialized hardware accelerators propose execut- 060
ing Transformer blocks using spiking neurons: a recent design 061
for an “energy-efficient spike transformer accelerator” re- 062
ports matching the performance of a conventional Trans- 063
former while consuming only a fraction of the power [3]. 064
On the software side, a recent review work demonstrates 065
a complete neuromorphic LLM by transforming a 1.5B- 066
parameter model into spike-train dynamics. This archi- 067
tecture eliminates all matrix multiplications and achieves 068
approximately $20\times$ greater energy efficiency compared to 069
GPU baselines [25]. Even in vision tasks, transformer- 070
based SNNs have attained state-of-the-art accuracy. For ex- 071
ample, *Meta-SpikeFormer* achieved 80% top-1 accuracy on 072
ImageNet using purely spiking layers, outperforming prior 073
SNN models while relying solely on sparse binary opera- 074
tions [25, 28, 29]. 075

In parallel, a separate line of work studies conditional 076
computation in Transformers. Standard LLMs employ a 077
fixed stack of layers for all inputs; however, multiple studies 078

079 note that this design wastes computation on “trivial” tokens
080 and lacks flexibility for complex reasoning [8]. Techniques
081 such as early exiting, layer skipping, and router-based gat-
082 ing have therefore been developed so that simpler inputs ac-
083 tivate fewer layers. For example, mixture-of-experts (MoE)
084 architectures route each token through only a small subset
085 of expert subnetworks. The sparsely-gated MoE model en-
086 abled effective models with hundreds of billions of param-
087 eters by activating only a limited number of experts per token
088 [21]. Similarly, [8] trains per-layer routers to skip or repeat
089 Transformer blocks, demonstrating that early layers are al-
090 most always utilized, middle layers can be skipped for eas-
091 ier tasks, and late layers may be repeated for more difficult
092 inputs. These dynamic-depth approaches yield substantial
093 computational gains. Spike-driven Transformer models, for
094 instance, report approximately 51% faster training and 70%
095 faster inference by exploiting sparsity[28]. In effect, such
096 methods trace a Pareto frontier between accuracy and com-
097 pute, enabling selection of a “sweet spot” in which easy in-
098 puts consume minimal layers while challenging inputs still
099 receive full processing.

100 The Spike-Gated Residual Unit (S-GRU) bridges the gap
101 between the representational power of high-performance
102 Transformers and the inherent energy efficiency of SNNs.
103 By introducing per-layer gating units into a pretrained
104 Transformer backbone, we enable the model to apply a
105 learned activation threshold that determines whether a spe-
106 cific layer is executed or remains dormant for a given for-
107 ward pass. This framework implements a form of *learned*
108 *architectural sparsity* inspired by event-driven computation.
109 Unlike purely theoretical sparsification approaches, we rig-
110 orously evaluate our modified architecture against industry-
111 standard dense baselines.

112 Our contributions are threefold:

- 113 1. **S-GRU Architecture:** We introduce the S-GRU, a
114 novel neuromorphic-inspired module that enables dy-
115 namic, event-driven layer skipping within pretrained
116 Transformer backbones.
- 117 2. **Depth Optimization:** We propose a sparsity-aware
118 training framework using a regularization penalty
119 $\lambda_{sparsity}$, allowing the model to self-optimize its com-
120 putational depth based on input complexity.
- 121 3. **Cross-Domain Validation:** We demonstrate that S-
122 GRU achieves superior efficiency-performance trade-
123 offs against SOTA dense (TinyLlama, Qwen-2.5, Phi-
124 2) and sparse (Lightweight MoE) baselines across
125 WikiText-2, HellaSwag, C4, and LAMBADA.

126 The remainder of this paper is organized as follows. Sec-
127 tion 2 reviews prior work related to SNNs and efficient
128 Transformer architectures. Section 3 details the S-GRU,
129 including the spiking-gated mechanism and the sparsity-
130 aware objective function. Section 4 presents our empiri-
131 cal evaluation, highlighting the trade-offs between language

modeling performance and compute depth. Finally, the im- 132
plications of these findings, future works and limitations are 133
discussed in Section 6, and Section 7 summarizes the paper. 134

2. Related Works 135

SNNs have emerged as a compelling alternative for energy- 136
efficient deep learning, leveraging sparse binary signals to 137
minimize computational overhead. Recently, the integra- 138
tion of SNNs with Transformer architectures has yielded 139
significant performance gains in computer vision[22] and 140
robotics[5]. Early efforts, such as those by [30] and 141
[33], laid the foundation for combining the long-range 142
dependency modeling of Transformers with the event- 143
driven nature of SNNs. These architectures typically re- 144
place standard Softmax attention and Multi-Layer Percep- 145
trons (MLPs) with spike-driven components utilizing Leaky 146
Integrate-and-Fire (LIF) neurons [16, 32]. The primary 147
advantage of these Spike-based Transformers lies in their 148
hardware-friendly operations; for instance, spike-driven 149
self-attention mechanisms can replace energy-intensive 150
floating-point multiplications with bitwise AND operations 151
and sparse accumulations [23]. Subsequent research has 152
explored architectural refinements to further enhance effi- 153
ciency, including the use of Spiking Wavelet transforms 154
[4], MoE [34], and specialized spatial-temporal attention 155
modules [15]. Furthermore, advancements in ANN-to-SNN 156
conversion techniques [9] have bridged the performance 157
gap between dense and spiking models. Collectively, these 158
methods achieve substantial efficiency, with reports of faster 159
training and faster inference speeds relative to conventional 160
dense Transformers by maximizing event-driven dataflows. 161

The extension of the spiking paradigm to LLMs has re- 162
cently gained momentum as a viable path toward energy- 163
efficient inference. Notable advancements include [25], 164
which mathematically converts 1.5B-parameter Transform- 165
ers into neuromorphic networks to eliminate matrix multi- 166
plications, and [24], which scales SNNs to the LLM regime 167
using saliency-based spiking to maintain performance. Par- 168
allel to these structural spiking transitions, the field of 169
conditional computation has explored dynamic depth as a 170
means of reducing computational waste. Recent strategies 171
for adaptive layer skipping [17, 18] and sublayer skipping 172
[7] demonstrate that not all tokens require the full depth 173
of a model, with certain middle layers often being redun- 174
dant for simpler inputs. Various approaches have emerged 175
to govern these skipping decisions, ranging from residual 176
gates that learn “what layers to skip when” [14] to Markov 177
Decision Policies that enable input-aware dynamic skipping 178
[27]. While some methods investigate the theoretical con- 179
ditions for skipping in multimodal contexts [6] or focus on 180
architectural efficiency for general inference [19], many re- 181
quire significant retraining or specialized hardware. Recent 182
efforts have introduced continuous gating units, such as the 183

184 Evaluator Adjuster Unit and Gated Residual Connections to
 185 dynamically modulate feature flow and improve contextual
 186 adaptability[2]. However, these mechanisms focus primar-
 187 ily on enhancing representation quality through soft scaling
 188 and often require training models from scratch. In contrast,
 189 our *S-GRU* introduces a lightweight software-level gating
 190 solution that retrofits onto pretrained LLMs. By incorporat-
 191 ing an explicit penalty $\lambda_{sparsity}$ into the training loss, we
 192 leverage spiking-inspired sparsity to navigate a controllable
 193 Pareto frontier between inference speed and model accuracy
 194 without the need for neuromorphic hardware.

195 3. Spike-Gated Residual Unit Architecture

196 The S-GRU architecture introduces a dynamic, input-
 197 dependent gating mechanism designed to optimize the com-
 198 putational efficiency of LLMs. Unlike static pruning or
 199 fixed layer-dropping techniques, S-GRU adapts its depth
 200 per forward pass by evaluating the importance of each resid-
 201 ual block in real-time. By augmenting a pre-trained trans-
 202 former backbone with lightweight spike-gated units, the
 203 model learns to selectively bypass layers that contribute
 204 minimally to the final output, thereby reducing the total
 205 Flops required for inference without necessitating a full re-
 206 training of the underlying model weights.

207 3.1. The Spike-Gated Unit

208 The core component of our architecture is the spike-gated
 209 Unit. In a standard Transformer, each layer l performs
 210 a transformation f_l on the input hidden states \mathbf{x}_l ($\mathbf{x}_l \in$
 211 $\mathbb{R}^{B \times d_{model}}$) followed by a residual addition:

$$212 \quad \mathbf{x}_{l+1} = \mathbf{x}_l + f_l(\mathbf{x}_l). \quad (1)$$

213 As shown in Fig.1, in the S-GRU, the scalar gate $g_l \in [0, 1]$
 214 is produced by a spike-gated unit. The residual connection
 215 is modified as follows:

$$216 \quad \mathbf{x}_{l+1} = \mathbf{x}_l + g_l \cdot (f_l(\mathbf{x}_l) - \mathbf{x}_l). \quad (2)$$

217 The spike-gated unit, first compresses the input sequence
 218 through global average pooling across the sequence dimen-
 219 sion to form a context vector $\mathbf{s} \in \mathbb{R}^{d_{model}}$:

$$220 \quad \mathbf{s} = \frac{1}{B} \sum_{t=1}^B \mathbf{x}_{l,t}, \quad (3)$$

221 where $\mathbf{x}_{l,t} \in \mathbb{R}^{d_{model}}$ represents the hidden state vec-
 222 tor at layer l and sequence position t , B denotes the block
 223 size (sequence length), ($B = 256$). t is the index represent-
 224 ing a specific token position within the block, ranging from
 225 1 to B , d_{model} is the hidden dimension of the transformer
 226 backbone ($d_{model} = 2048$). $\mathbf{s} \in \mathbb{R}^{d_{model}}$ is the resulting
 227 context vector that provides a global summary of the input

sequence, serving as the input to the spike-gated unit which
 utilizes a MLP to determine the dynamic scalar gate.

$$g_l = \sigma((\text{SiLU}(\mathbf{s}\mathbf{W}_1 + \mathbf{b}_1))\mathbf{W}_2 + b_2), \quad (4)$$

228 where $\mathbf{W}_1 \in \mathbb{R}^{d_{model} \times (h \cdot d_{model})}$ is the weight matrix of the
 229 first linear layer, which projects the context vector into a
 higher-dimensional latent space. $\mathbf{b}_1 \in \mathbb{R}^{h \cdot d_{model}}$ is the bias
 vector for the first linear layer, h is the hidden multiplier
 that determines the width of the gate’s hidden layer. In our
 implementation, $h = 2$, resulting in a hidden dimension of
 $2 \times 2048 = 4096$, $\text{SiLU}(\cdot)$ is the Sigmoid Linear Unit ac-
 tivation function (also known as Swish), defined as $f(x) =$
 $x \cdot \sigma(x)$ which introduces non-linearity to allow the gate to
 learn complex skipping patterns, $\mathbf{W}_2 \in \mathbb{R}^{(h \cdot d_{model}) \times 1}$ is the
 weight matrix of the second linear layer, which reduces the
 hidden representation down to a single scalar logit, $b_2 \in \mathbb{R}$
 is the scalar bias for the final layer, initialized to a high value
 (e.g., 5.0) to ensure the gate starts in an “open” state and
 $\sigma(\cdot)$ is the Sigmoid function, which ensures that g_l remains
 within the unit interval, effectively acting as a “soft” skip
 connection that the model can optimize during the sparsity-
 aware training phase.

249 3.2. Sparsity-Aware Training

250 The training objective of S-GRU is to minimize the standard
 251 language modeling loss while simultaneously maximizing
 252 the sparsity of the network. To achieve this without degrad-
 253 ing the pre-trained capabilities of the model, we employ a
 254 Gate-Only fine-tuning strategy.

255 During this phase, the parameters of the pre-trained
 256 transformer backbone θ_{LLM} are frozen, and only the pa-
 257 rameters of the spike-gated unit’s θ_G are updated:

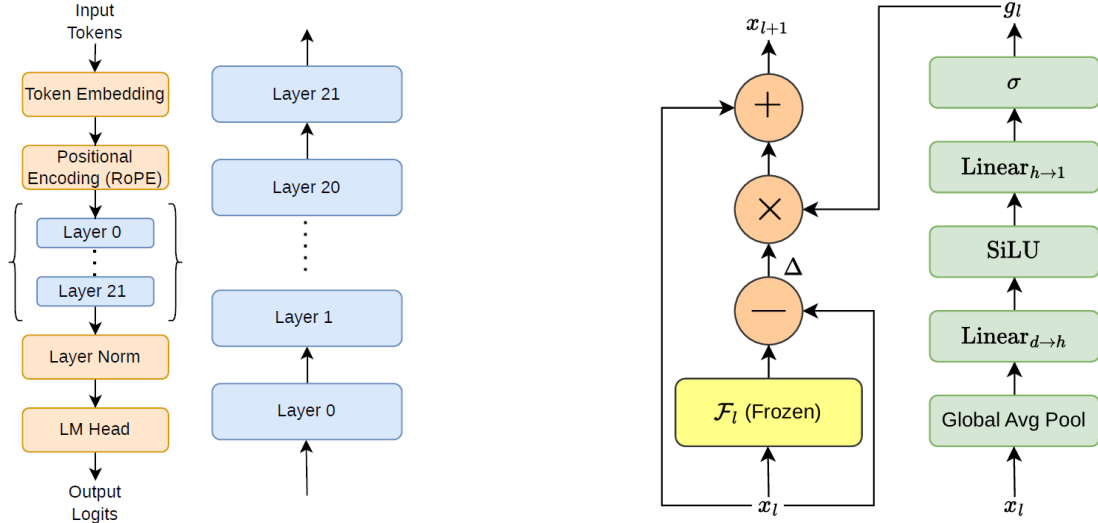
$$258 \quad \nabla_{\theta} \mathcal{L} = \nabla_{\theta_G} \mathcal{L}, \quad \text{subject to } \frac{\partial \mathcal{L}}{\partial \theta_{LLM}} = 0, \quad (5)$$

259 where $\nabla_{\theta} \mathcal{L}$ is the Total Gradient Vector with respect to all
 260 parameters in the network, $\nabla_{\theta_G} \mathcal{L}$ is the Gate Gradient term
 261 signifies that the optimizer only computes and applies up-
 262 dates to the parameters within the gated units, specifically
 263 the MLP weights and biases $\{\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, b_2\}$, $\frac{\partial \mathcal{L}}{\partial \theta_{LLM}}$
 264 which is the partial derivative of the loss with respect to
 265 the transformer backbone. By setting this term to zero, we
 266 freeze the original model weights and ∇ is the Gradient Op-
 267 erator, which collects all partial derivatives into a vector.

268 The model is optimized using a composite loss function
 269 \mathcal{L}_{total} , which introduces a penalty on the gate activations.
 270 We define the total objective as:

$$271 \quad \mathcal{L}_{total} = \mathcal{L}_{LM} + \lambda_{sparsity} \cdot \left(\frac{1}{L} \sum_{l=1}^L \bar{g}_l \right), \quad (6)$$

272 where L is the total number of decoder layers, \bar{g}_l is the mean
 273 gate activation for layer l across the batch and sequence,



(a) Overall language model architecture showing the 22 sequential decoder layers (highlighted in light blue), where the S-GRU is applied to the output of each of the decoder layer.

(b) Spiking Residual and dynamic gating module used for each decoder layer.

Figure 1. S-GRU integrated architecture. Residual modification used in each layer: the input x_l is passed through the frozen transformer block \mathcal{F}_l to compute a residual $\Delta = \mathcal{F}_l(x_l) - x_l$, which is modulated by a learned spike gate $g_l \in [0, 1]$ generated by an MLP, producing the output $x_{l+1} = x_l + g_l \cdot \Delta$.

274 and $\lambda_{sparsity}$ is a hyperparameter controlling the trade-off
 275 between model performance and computational efficiency.
 276 Let \mathcal{L}_{LM} denote the cross-entropy loss for causal language
 277 modeling.

$$278 \quad \mathcal{L}_{LM} = -\frac{1}{T} \sum_t t = 1^T \log P(x_t | x_{<t}; \theta), \quad (7)$$

279 which is the classic Average Negative Log-Likelihood
 280 (NLL) over a sequence of tokens. By including the gate
 281 values directly in the loss function, the optimizer is incen-
 282 tivated to drive g_l toward zero. Because the sigmoid activa-
 283 tion is continuous and differentiable, the gradients from the
 284 sparsity penalty flow directly into the spike-gate weights, al-
 285 lowing the model to learn which layers are redundant for the
 286 target domain. This formulation ensures that a layer is only
 287 bypassed if the resulting increase in \mathcal{L}_{LM} is smaller than
 288 the decrease in the sparsity penalty, effectively performing
 289 a cost-benefit analysis of each layer’s contribution during
 290 the forward pass.

291 3.3. Gate Initialization

292 The initialization of the spike-gate parameters is critical
 293 for stabilizing the transition from a dense pre-trained back-
 294 bone to a sparse gated architecture. Our MLP configura-
 295 tion consists of four primary learnable parameters: $\theta_G =$
 296 $\{\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, b_2\}$. Standard random initialization (e.g.,
 297 Xavier or Kaiming) [20] for these parameters typically re-
 298 sults in a pre-activation logit $z \approx 0$, yielding gate outputs

299 $g_l \approx 0.5$ at the onset of training. Such a state would imme-
 300 diately bypass 50% of the network’s processing capacity,
 301 leading to catastrophic forgetting and gradient instability.

302 To mitigate this, we implement a High-Bias Start strat-
 303 egy specifically targeting the final projection layer. While
 304 the input projection parameters $\{\mathbf{W}_1, \mathbf{b}_1\}$ are initialized us-
 305 ing standard Kaiming uniformity to facilitate feature dis-
 306 covery, the output weights are set to zero, $\mathbf{W}_2 = \mathbf{0}$, and
 307 the bias is initialized as $b_2 = \text{gate_init_bias}$. By setting
 308 the final weight matrix \mathbf{W}_2 to zero, we decouple the gate’s
 309 output from the input hidden states \mathbf{s} during the initial for-
 310 ward pass. Given the sigmoid activation function $\sigma(z) =$
 311 $(1 + e^{-z})^{-1}$, setting a high initial bias (e.g., $b_2 = 5.0$) en-
 312 sures that the gate activation g_l saturates near unity:

$$313 \quad g_l = \sigma((\text{SiLU}(\mathbf{s}\mathbf{W}_1 + \mathbf{b}_1))\mathbf{W}_2 + b_2) = \sigma(0 + 5.0) \approx 0.993. \quad (8)$$

314 This initialization ensures that at step 0, the Spikeformer
 315 behaves nearly identically to the original dense backbone,
 316 preserving its full predictive power. As training progresses,
 317 the gradients from the sparsity penalty $\lambda_{sparsity}$ compete
 318 with the language modeling loss to pull these biases down-
 319 ward. This allows the model to “prune” layers gracefully,
 320 only reducing g_l for a specific layer once the network has
 321 learned that the layer’s residual contribution is redundant
 322 for the target objective.

323 4. Experiments

324 4.1. Setup

325 To evaluate the efficiency and effectiveness of the S-GRU,
326 we perform experiments using a dense-to-sparse transfor-
327 mation on a representative small-scale Large Language
328 Model.

329 **Datasets and Tokenization.** Our primary training and
330 evaluation corpuses are *WikiText-2*, *HellaSwag*, *C4* and
331 *LAMBADA*, tokenized into sequences of 256 tokens using
332 the standard Llama-3 tokenizer. For the sparsity-aware
333 training phase, we utilize a subset of 8,000 training ex-
334 amples and 1,000 validation examples. To evaluate the
335 preservation of downstream reasoning and long-range de-
336 pendency retention under dynamic depth constraints, we
337 further benchmark the model on the *HellaSwag* common-
338 sense reasoning task (100 zero-shot samples), the *LAM-*
339 *BADA* narrative context dataset, and the *C4* (Colossal Clean
340 Crawled Corpus) validation set.

341 **Model Backbone.** We employ TinyLlama-1.1B as our
342 base transformer architecture. This model consists of 22
343 layers with a hidden dimension of $d_{model} = 2048$. Follow-
344 ing the methodology described in Section-3, we inject 22
345 spike-gate units (one per residual block) utilizing the MLP
346 configuration.

347 **Training Configuration.** We train only the gate param-
348 eters θ_G for 400 steps using the AdamW optimizer. Based on
349 the stabilization strategy in 3.2, we set the initial gate bias
350 to 3.0. The learning rate is set to 1×10^{-3} with zero weight
351 decay to prevent the gates from being prematurely forced to
352 zero by the optimizer’s regularization rather than the spar-
353 sity loss. All experiments are conducted in BFLoat16 pre-
354 cision to match the original model’s training regime.

355 **Hardware.** All experiments were performed on a single
356 NVIDIA RTX 5060 GPU. Given that the language model
357 weights are frozen, the memory overhead remains minimal,
358 requiring less than 7 GB of VRAM during the gating opti-
359 mization phase.
360

361 4.2. Comparative Analysis

362 We evaluate the S-GRU’s performance by comparing its
363 perplexity and zero-shot accuracy against the original dense
364 TinyLlama-1.1B backbone and a static layer-dropping
365 baseline. Our results indicate that at a sparsity constraint of
366 $\lambda_{sparsity} = 0.5$, the S-GRU achieves a perplexity of 12.46,
367 representing a controlled degradation of 2.28 units com-
368 pared to the dense baseline’s perplexity of 10.18. Notably,

on the *HellaSwag* benchmark, the S-GRU retains 90.9% of
the original reasoning performance (41.0% vs. 44.0% base-
line), significantly outperforming the static baseline which
suffers from a catastrophic drop in accuracy when an equiv-
alent number of layers are removed. This gap suggests that
the spike-gate units successfully identify and preserve “crit-
ical” layers required for complex reasoning while bypass-
ing redundant computations in simpler contexts, effectively
shifting the Pareto frontier of the accuracy-efficiency trade-
off.

5. Ablation Study

To characterize the relationship between computational cost
and model fidelity, we performed a sweep of the sparsity
penalty across $\lambda_{sparsity} \in \{0.0, 0.2, 0.5, 1.0\}$. The ob-
served learning trajectories across these coefficients reveal
a distinct, non-linear transition in the emergent architecture
of the S-GRU. For $\lambda_{sparsity} = 0.0$, the model retains near-
maximal density with an overall mean gate activation of
0.98, successfully preserving the original dense backbone
state. As $\lambda_{sparsity}$ increases to 0.2, the gates begin a se-
lective suppression of specific layer blocks, a trend that ac-
celerates significantly at $\lambda_{sparsity} = 0.5$. In this regime,
the overall mean gate activation drops to 0.76, and the per-
centage of layers exceeding the 0.95 activation threshold
collapses toward zero. Finally, at $\lambda_{sparsity} = 1.0$, the
model reaches an extreme minimalist state, identifying the
absolute lower bound of depth required to maintain func-
tional coherence. This shift from dense to sparse process-
ing follows a clear trade-off curve, which is visualized in
the sparsity-perplexity relationship shown in Fig. 3. De-
tailed analysis of the specific layer-wise behaviors and the
transition dynamics within this curve is provided in the Dis-
cussion section (Section 6).

6. Discussion

Architectural Efficiency The results presented in Ta-
ble 1 validate that S-GRU achieves a superior balance be-
tween linguistic modeling capacity and architectural effi-
ciency compared to established baselines. Most notably,
at a mean gate activation of 76.5%, our model achieves
a WikiText-2 perplexity of 12.45. This performance effec-
tively matches the density-heavy Qwen-2.5-1.5B (12.49
PPL) and surpasses the larger Phi-2 (13.14 PPL), despite
S-GRU utilizing 26% and 59% fewer active parameters per
forward pass, respectively. While we observe a marginal
degradation in HellaSwag reasoning accuracy (41.0%) re-
lative to the fully-dense TinyLlama baseline (44.0%), the S-
GRU compensates with a significant theoretical throughput
increase of $1.31\times$. These findings indicate that the spike-
gate units do not merely perform a crude truncation of the
network, but instead autonomously re-configure the 1.1B

Table 1. **Comprehensive Evaluation of S-GRU Against Baseline LLMs.** We compare models across general reasoning (HellaSwag), efficiency (Active Layers, Speedup), and linguistic generalization (WikiText-2, C4, and LAMBADA). Our method significantly improves efficiency while maintaining competitive perplexity across diverse corpora.

Model	Params	General Performance		Efficiency		Perplexity (PPL) ↓		
		HellaSwag ↑	Wiki-2 ↓	Active Layers ↓	Speedup ↑	C4	LAMBADA	Avg.
TinyLlama-1.1B [31]	1.1B	44.0%	10.18	100%	1.00x	11.00	22.30	14.49
Phi-2 [12]	2.7B	52.0%	13.14	100%	0.41x*	16.71	37.02	22.29
Qwen-2.5-1.5B [10]	1.5B	52.0%	12.49	100%	0.73x*	19.03	28.85	20.12
DeepSeek-MoE [1]	1.3B	46.0%	11.26	100%	1.00x*	14.84	23.34	16.48
S-GRU (Ours)	1.1B	41.0%	12.45	76.5%	1.31x	16.50	28.68	19.21

*Theoretical speedup relative to TinyLlama-1.1B base throughput.

Evolution of Learned Sparsity ($\lambda = 0.5$)

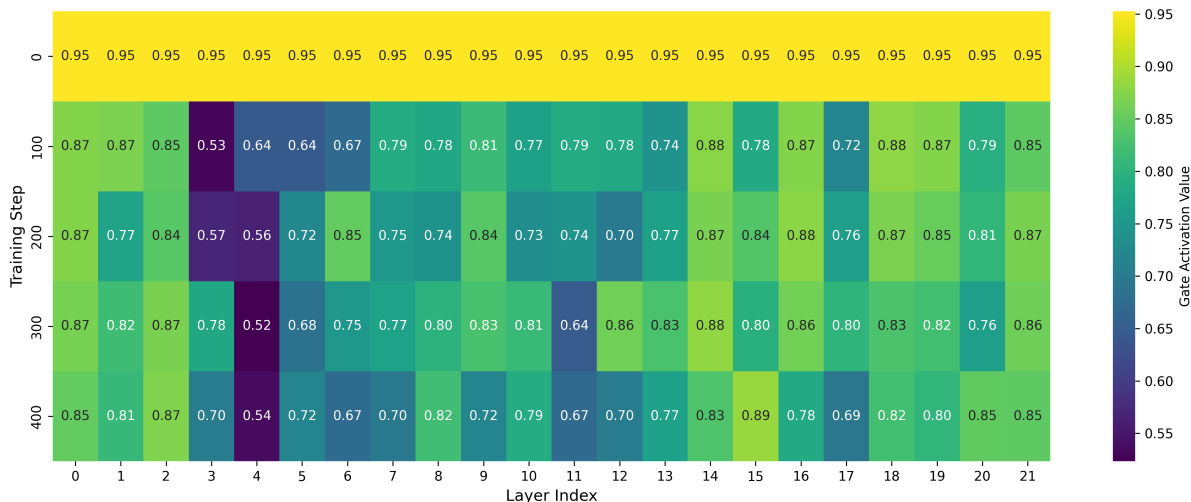


Figure 2. **Evolution of Architectural Sparsity** ($\lambda_{sparsity} = 0.5$). The heatmap displays the mean gate activation for each of the 22 layers over 400 training steps. Darker regions indicate layers the model has learned to bypass, achieving a 23.5% reduction in compute depth.

419 parameter backbone into a more potent, sparse representation that rivals the performance of significantly larger, fully-
420 dense models.
421

422 **Cross-Domain Robustness and Generalization.** To
423 evaluate the robustness of our sparse representation beyond standard benchmarks, we report linguistic generaliza-
424 tion across diverse corpora in Table 1. On the C4 dataset,
425 our framework achieves a perplexity of 16.5, notably out-
426 performing the significantly larger Qwen-2.5-1.5B (19.03
427 PPL) and maintaining parity with Phi-2 (16.71 PPL). These
428 results suggest that the S-GRU mechanism retains high-
429 fidelity world knowledge even under strict sparsity con-
430 straints. Furthermore, in long-range context modeling via
431 the LAMBADA benchmark, our model yields a perplexity
432 of 28.68, successfully bridging the performance gap
433 between established dense baselines and specialized effi-
434

cient architectures like the Lightweight MoE [1]. While
TinyLlama maintains a lead in raw next-token prediction
on these datasets, the competitive performance of S-GRU
across both general web text and long-context scenarios
confirms that the learned spiking dynamics generalize ef-
fectively to out-of-distribution data, rather than overfitting
to the specific structures of the training corpus.

443 **Layer wise Gate Dynamics.** The layer-wise gate dynam-
444 ics across the tested $\lambda_{sparsity}$ spectrum reveal a structured,
445 non-uniform pruning process that prioritizes specific archi-
446 tectural “anchors.” As illustrated in the gate activation heat
447 plot (Fig. 2), the model exhibits an emergent hierarchy of
448 layer importance. At $\lambda_{sparsity} = 0.0$, all 22 layers main-
449 tain near-maximal saturation ($\approx 0.95 - 0.98$), confirming
450 that every block remains functionally active. However, as

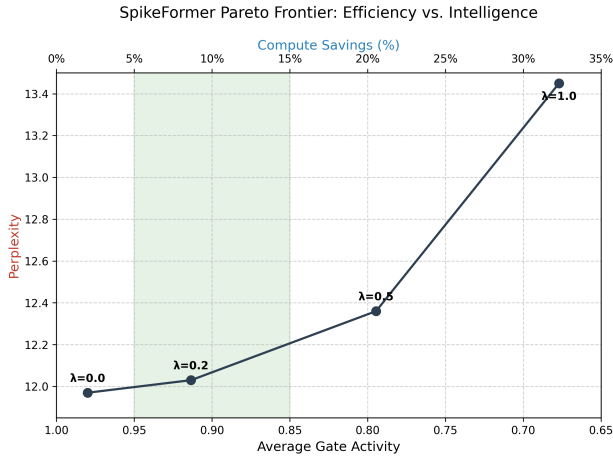


Figure 3. **Efficiency-Intelligence Pareto Frontier for S-GRU.** The plot illustrates the trade-off between model perplexity and average gate activity (compute savings) across varying sparsity constraints ($\lambda_{sparsity}$). The shaded green region highlights the optimal operating range where significant compute savings (5%–15%) are achieved with minimal impact on perplexity. As sparsity increases ($\lambda_{sparsity} = 1.0$), compute savings reach $\sim 33\%$ at the cost of higher perplexity.

451 the sparsity pressure increases to $\lambda_{sparsity} = 0.2$, we observe
 452 the first signs of selective thinning, particularly in the
 453 middle-to-late intermediate blocks, while the initial and final
 454 layers (such as Layer 0 and Layer 21) remain robustly
 455 active. By the critical $\lambda_{sparsity} = 0.5$ threshold, this behavior
 456 culminates in a “survivor” architecture: the overall mean gate
 457 activation collapses to 0.76, yet the gates consistently sustain
 458 higher coefficients for the deep residual units. This indicates
 459 that the S-GRU identifies the final stages of the transformer
 460 as the most critical for reconciling high-level semantic representa-
 461 tions before the language modeling head. In the extreme
 462 $\lambda_{sparsity} = 1.0$ regime, even these anchors begin to fade,
 463 suggesting that the model has reached its absolute structural
 464 limit, operating effectively as a shallow feature extractor rather
 465 than a deep reasoning engine.
 466

467 **Limitations and Future Work.** Despite the promising
 468 theoretical efficiency of the S-GRU, several limitations remain
 469 to be addressed in future work. While our model achieves a
 470 $1.31\times$ theoretical speedup, realizing these gains on standard
 471 deep learning hardware remains a significant challenge. Current
 472 GPU kernels are optimized for dense matrix multiplications;
 473 consequently, the dynamic skipping of individual residual units
 474 often results in negligible wall-clock time savings due to kernel
 475 launch overheads and the lack of hardware-level support for
 476 conditional branching in massive parallel arrays. To bridge this
 477 gap, future research

478 should focus on deploying the developed design on specialized
 479 CUDA kernels or hardware-aware compilers that can “short-circuit”
 480 the computation of skipped layers at the register level. Furthermore,
 481 while the current gating mechanism utilizes a sequence-wide summary,
 482 exploring token-level dynamic sparsity could further refine the
 483 model’s efficiency by skipping computations for “easy” tokens (e.g.,
 484 punctuation or common stop words). Beyond language modeling,
 485 we envision that the S-GRU framework could be generalized to
 486 other modalities, such as vision or multi-modal transformers,
 487 where dynamic depth remains a critical challenge. Finally, extend-
 488 ing this framework to larger models, such as Llama-3-70B, would
 489 determine if the emergent layer-redundancy observed in TinyLlama
 490 scales to architectures with deeper residual hierarchies.
 491
 492

7. Conclusion

493 In this work, we introduced a S-GRU framework that enables
 494 dynamic depth adaptation in pretrained Transformer-based large
 495 language models. By integrating lightweight spike-gated units
 496 into each residual block and optimizing a sparsity-aware objective
 497 with a controllable regularization coefficient $\lambda_{sparsity}$, S-GRU
 498 autonomously learns to bypass redundant layers while preserving
 499 essential reasoning capacity. Through gate-only fine-tuning on
 500 TinyLlama-1.1B, we demonstrated that the model can reduce
 501 average active depth to 76.5%, achieving a $1.31\times$ theoretical
 502 efficiency gain with minimal degradation in perplexity and zero-
 503 shot HellaSwag accuracy. Our empirical analysis reveals a
 504 structured layer-importance hierarchy, where early and late
 505 blocks remain critical while intermediate layers are selectively
 506 suppressed under sparsity pressure. These results validate that
 507 learned residual gating provides a practical and stable pathway
 508 toward energy-aware Transformer inference without requiring full
 509 model retraining or specialized neuromorphic hardware.
 510
 511
 512

References

- 513
- 514 [1] Damai Dai, Chengqi Deng, Chenggang Zhao, RX Xu,
 515 Huazuo Gao, Deli Chen, Jiashi Li, Wangding Zeng, Xingkai
 516 Yu, Yu Wu, et al. Deepseekmoe: Towards ultimate expert
 517 specialization in mixture-of-experts language models. In
 518 *Proceedings of the 62nd Annual Meeting of the Association
 519 for Computational Linguistics (Volume 1: Long Papers)*,
 520 pages 1280–1297, 2024. 6
 - 521 [2] Sahil Rajesh Dhayalkar. Dynamic context adaptation and
 522 information flow control in transformers: Introducing the
 523 evaluator adjuster unit and gated residual connections. *arXiv
 524 preprint arXiv:2405.13407*, 2024. 3
 - 525 [3] Congpeng Du, Qi Wen, Zhiqiang Wei, and Hao Zhang. En-
 526 ergy efficient spike transformer accelerator at the edge. *Intelli-
 527 gent Marine Technology and Systems*, 2(1):24, 2024. 1
 - 528 [4] Yuetong Fang, Ziqing Wang, Lingfeng Zhang, Jiahang Cao,
 529 Honglei Chen, and Renjing Xu. Spiking wavelet transformer.

- 530 In *European conference on computer vision*, pages 19–37. Springer, 2024. 2
- 531
- 532 [5] Maoxu Gao, Zhixing Hou, Di Zhu, and Chio-In Jeong. An accuracy-improved low-computational-cost spiking trans- 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586
- [5] Maoxu Gao, Zhixing Hou, Di Zhu, and Chio-In Jeong. An accuracy-improved low-computational-cost spiking transformer network for robotic arm manipulation. In *2025 IEEE International Conference on Integrated Circuits, Technologies and Applications (ICTA)*, pages 294–295. IEEE, 2025. 2
- [6] Max Hartman, Vidhata Jayaraman, Moulik Choraria, Akhil Bhimaraju, and Lav R Varshney. Skip-it? theoretical conditions for layer skipping in vision-language models. *arXiv preprint arXiv:2509.25584*, 2025. 2
- [7] Zhuomin He, Yizhen Yao, Pengfei Zuo, Bin Gao, Qinya Li, Zhenzhe Zheng, and Fan Wu. Adaskip: Adaptive sublayer skipping for accelerating long-context llm inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 24050–24058, 2025. 2
- [8] Ahmed Heakl, Martin Gubri, Salman Khan, Sangdoon Yun, and Seong Joon Oh. Dr. llm: Dynamic layer routing in llms. *arXiv preprint arXiv:2510.12773*, 2025. 1, 2
- [9] Zihan Huang, Xinyu Shi, Zecheng Hao, Tong Bu, Jianhao Ding, Zhaofei Yu, and Tiejun Huang. Towards high-performance spiking transformers from ann to snn conversion. In *Proceedings of the 32nd ACM international conference on multimedia*, pages 10688–10697, 2024. 2
- [10] Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024. 6
- [11] Phu Khanh Huynh, M Lakshmi Varshika, Ankita Paul, Murat Isik, Adarsha Balaji, and Anup Das. Implementing spiking neural networks on neuromorphic architectures: A review. *arXiv preprint arXiv:2202.08897*, 2022. 1
- [12] Mojan Javaheripi, Sébastien Bubeck, Marah Abdin, Jyoti Aneja, Sebastien Bubeck, Caio César Teodoro Mendes, Weizhu Chen, Allie Del Giorno, Ronen Eldan, Sivakanth Gopi, et al. Phi-2: The surprising power of small language models. *Microsoft Research Blog*, 1(3):3, 2023. 6
- [13] Zhenya Ji and Ming Jiang. A systematic review of electricity demand for large language models: evaluations, challenges, and solutions. *Renewable and Sustainable Energy Reviews*, 225:116159, 2026. 1
- [14] Filipe Laitenberger, Dawid Kopiczko, Cees GM Snoek, and Yuki M Asano. What layers when: Learning to skip compute in llms with residual gates. *arXiv preprint arXiv:2510.13876*, 2025. 2
- [15] Donghyun Lee, Yuhang Li, Youngeun Kim, Shiting Xiao, and Priyadarshini Panda. Spiking transformer with spatial-temporal attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13948–13958, 2025. 2
- [16] Yudong Li, Yunlin Lei, and Xu Yang. Spikeformer: Training high-performance spiking neural network with transformer. *Neurocomputing*, 574:127279, 2024. 2
- [17] Yijin Liu, Fandong Meng, and Jie Zhou. Accelerating inference in large language models with a unified layer skipping strategy. *arXiv preprint arXiv:2404.06954*, 2024. 2
- [18] Xuan Luo, Weizhi Wang, and Xifeng Yan. Adaptive layer-skipping in pre-trained llms. *arXiv preprint arXiv:2503.23798*, 2025. 2
- [19] Dheemanth Manur. Adaptive skip for efficient llm inference. Master’s thesis, University of California, Riverside, 2025. 2
- [20] Yu Pan, Zeyong Su, Ao Liu, Wang Jingquan, Nannan Li, and Zenglin Xu. A unified weight initialization paradigm for tensorial convolutional neural networks. In *International conference on machine learning*, pages 17238–17257. PMLR, 2022. 4
- [21] Noam Shazeer, Azalia Mirhoseini, Krzysztof Maziarz, Andy Davis, Quoc Le, Geoffrey Hinton, and Jeff Dean. Outrageously large neural networks: The sparsely-gated mixture-of-experts layer. *arXiv preprint arXiv:1701.06538*, 2017. 2
- [22] Xinyu Shi, Zecheng Hao, and Zhaofei Yu. Spikingresformer: Bridging resnet and vision transformer in spiking neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5610–5619, 2024. 2
- [23] Ziqing Wang, Yuetong Fang, Jiahang Cao, Qiang Zhang, Zhongrui Wang, and Renjing Xu. Masked spiking transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1761–1771, 2023. 2
- [24] Xingrun Xing, Boyan Gao, Zheng Zhang, David A Clifton, Shitao Xiao, Li Du, Guoqi Li, and Jiajun Zhang. Spikellm: Scaling up spiking neural network to large language models via saliency-based spiking. *arXiv preprint arXiv:2407.04752*, 2024. 2
- [25] Han Xu, Xuerui Qiu, Yunhui Xu, Mohammed E Elbity, Peng Zhou, Yang Tian, Rui-Jie Zhu, Jiahong Zhang, Shaowei Gu, Yuqi Pan, et al. Neuromorphic spike-based large language model. *National Science Review*, 13(4):nwaf551, 2026. 1, 2
- [26] Zhanguo Yan, Kaiwen Tang, Zixuan Zhu, Zhenyu Bai, Qianhui Liu, and Weng-Fai Wong. Matterhorn: Efficient analog sparse spiking transformer architecture with masked time-to-first-spike encoding. *arXiv preprint arXiv:2601.22876*, 2026. 1
- [27] Ning Yang, Fangxin Liu, Junjie Wang, Tao Yang, Kan Liu, Haibing Guan, and Li Jiang. Dash: Input-aware dynamic layer skipping for efficient llm inference with markov decision policies. *arXiv preprint arXiv:2505.17420*, 2025. 2
- [28] Man Yao, Jiakui Hu, Zhaokun Zhou, Li Yuan, Yonghong Tian, Bo Xu, and Guoqi Li. Spike-driven transformer. *Advances in neural information processing systems*, 36:64043–64058, 2023. 1, 2
- [29] Man Yao, Jiakui Hu, Tianxiang Hu, Yifan Xu, Zhaokun Zhou, Yonghong Tian, Bo Xu, and Guoqi Li. Spike-driven transformer v2: Meta spiking neural network architecture inspiring the design of next-generation neuromorphic chips. *arXiv preprint arXiv:2404.03663*, 2024. 1
- [30] Jiqing Zhang, Bo Dong, Haiwei Zhang, Jianchuan Ding, Felix Heide, Baocai Yin, and Xin Yang. Spiking transformers for event-based single object tracking. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*, pages 8801–8810, 2022. 2

- 643 [31] Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu.
644 Tinyllama: An open-source small language model. *arXiv*
645 *preprint arXiv:2401.02385*, 2024. 6
- 646 [32] Chenlin Zhou, Liutao Yu, Zhaokun Zhou, Zhengyu Ma, Han
647 Zhang, Huihui Zhou, and Yonghong Tian. Spikingformer:
648 Spike-driven residual learning for transformer-based spiking
649 neural network. *arXiv preprint arXiv:2304.11954*, 2023. 2
- 650 [33] Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang,
651 Shuicheng Yan, Yonghong Tian, and Li Yuan. Spikformer:
652 When spiking neural network meets transformer. *arXiv*
653 *preprint arXiv:2209.15425*, 2022. 2
- 654 [34] Zhaokun Zhou, Yijie Lu, Yanhao Jia, Kaiwei Che, Jun
655 Niu, Liwei Huang, Xinyu Shi, Yuesheng Zhu, Guoqi Li,
656 Zhaofei Yu, et al. Spiking transformer with experts mix-
657 ture. *Advances in Neural Information Processing Systems*,
658 37:10036–10059, 2024. 2