# ACHIEVE LATENCY-EFFICIENT TEMPORA-CODING SPIKING LLMs VIA DISCRETIZATION-AWARE CONVERSION

**Anonymous authors**Paper under double-blind review

### **ABSTRACT**

Large language models (LLMs) have achieved remarkable success while introducing critical energy bottlenecks that challenge sustainable deployment. Spiking neural networks (SNNs) provide a promising approach for energy-efficient spiking LLMs via ANN-to-SNN (A2S) conversion. Among various spike coding methods, time-to-first-spike (TTFS) coding is particularly appealing as it conveys information with a single spike, further reducing energy consumption. However, existing TTFS-based A2S conversion relies on continuous-time assumptions, requiring prohibitively large latencies (e.g., 4096 time steps) to approximate ANN's continuous values. This dependency leads to unacceptable inference delay in deep models, particularly LLMs, posing significant challenges for developing practical temporal-coding spiking LLMs. In this paper, we propose a discretization-aware theoretical framework that establishes a precise correspondence between discrete TTFS-based SNNs and ANNs. Our key insight reveals that conversion errors are bounded by latency-dependent terms. Motivated by these, we introduce the Quantization-Consistent ANN-to-SNN (QC-A2S) conversion, which integrates low-bit quantization with discretization-compatible TTFS neurons, achieving latency-efficient temporal-coding spiking LLMs. Comprehensive evaluation on LLaMA models demonstrates comparable performance with dramatically reduced latency.

## 1 Introduction

Large Language Models (LLMs) represent a paradigm shift in artificial intelligence, leveraging deep learning architectures trained on massive text corpora to capture intricate linguistic patterns, syntactic structures, and semantic relationships, thereby achieving remarkable capabilities in natural language understanding and generation (Zhang et al., 2022; Touvron et al., 2023; Achiam et al., 2023; Dubey et al., 2024). Most LLMs are built upon the Transformer architecture, which relies heavily on multi-head attention mechanisms and dense matrix multiplications, resulting in cubic computational complexity and substantial energy consumption during both training and inference (Vaswani et al., 2017; Zhao et al., 2023). Moreover, following the "scaling law", LLMs have grown from billions to trillions of parameters to achieve better performance, which further increases computational and storage demands (Chen et al., 2024a; Hoffmann et al., 2022). Consequently, the critical challenge facing the LLM community is developing approaches to reduce computational complexity and energy consumption while preserving model performance capabilities.

Spiking Neural Networks (SNNs) are biologically plausible computational models inspired by the mechanisms of neurons and synapses in the human brain (Maass, 1997; Roy et al., 2019). SNNs transmit and compute information asynchronously through discrete spike events rather than continuous-valued activation functions, demonstrating remarkable energy efficiency when implemented on specialized neuromorphic hardware (Yao et al., 2023; Zhou et al., 2022; Davies et al., 2018; Merolla et al., 2014). Consequently, developing **spiking LLMs** has emerged as a promising solution to address the substantial energy consumption challenges of LLMs. Currently, two primary approaches are used to develop spiking LLMs: direct training methods that incorporate surrogate gradients to address non-differentiability (Yao et al., 2023; Mukhoty et al., 2023; Zhou et al., 2024), and ANN-to-SNN (A2S) conversion methods that transfer pre-trained weights while preserving ap-

proximate equivalence through carefully designed techniques (Jiang et al., 2024; Chen et al., 2025a; Hao et al., 2023). Given the enormous computational and storage requirements of direct training for LLMs, practical spiking LLMs are predominantly achieved through A2S conversion for energy-efficient intelligent applications in resource-constrained environments (Xing et al., 2024a).

Beyond the rate coding commonly used in A2S conversion methods, recent neuroscience research has highlighted temporal-based spike coding that offer superiors energy efficiency advantages (Park et al., 2019; Zhang et al., 2019; Stanojevic et al., 2024). Temporal coding represents continuous values through precise spike timing rather than spike counts, suggesting that the representation of information depends on when the spikes occur (Gütig & Sompolinsky, 2006). Among various temporal codings, **time-to-first-spike (TTFS) coding** is particularly noteworthy, as it encodes information in the latency of a single spike, which substantially reduces energy consumption by minimizing spike counts (Park et al., 2020; Rueckauer & Liu, 2018).

Existing TTFS-based conversion methods underlying rely on continuous-time assumptions that directly approximate the continuous values of ANNs (Zhao et al., 2025; Stanojevic et al., 2024). However, practical hardware implementations impose discrete timing constraints through finite latency and clock granularity. Such discretization inevitably introduces conversion errors that severely compromise model accuracy. To mitigate the discretization-induced errors, existing methods require prohibitively large latency (e.g., 4096 time steps),

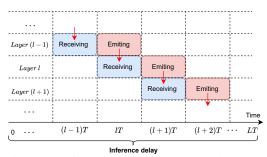


Figure 1: Inference delay across network layers.

causing extensive spike latency accumulation across network layers (Figure 1). This creates unacceptable inference delay in deep models, particularly for LLMs, posing significant challenges for developing practical **temporal-coding spiking LLMs**.

To address this fundamental challenge, we propose a discretization-aware theoretical framework that establishes a precise correspondence between discrete TTFS-based SNNs and ANNs. Our key theoretical insight reveals that conversion errors are formally bounded by latency-dependent terms, drawing a direct connection to quantization error bounds. Motivated by this equivalence, we introduce a paradigm shift from traditional continuous-approximation conversions to discrete-equivalent coversion. Specifically, we present the Quantization-Consistent ANN-to-SNN Conversion (QC-A2S), which integrates low-bit quantization with discretization-compatible TTFS neurons. QC-A2S leverages pre-quantized LLMs to inherently align with discrete spike dynamics, effectively mitigating conversion errors while achieving latency-efficient temporal-coding spiking LLMs. Comprehensive evaluation on LLaMA models demonstrates that our approach maintains comparable accuracy with dramatically reduced inference latency (Figure 2). The key contributions are summarized as follows:

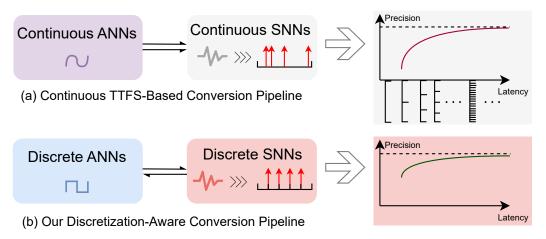


Figure 2: The framework of Quantization-Consistent ANN-to-SNN Conversion.

- We propose a discretization-aware theoretical framework for TTFS-based coding that identifies the fundamental discrepancy between continuous-time assumptions in prior TTFS methods and practical hardware constraints, revealing the formal equivalence between conversion errors and quantization error bounds.
- We present the QC-A2S framework, which represents a paradigm shift from traditional continuous-approximation conversions to discrete-equivalent transformation, enabling the first latency-efficient TTFS-based temporal-coding spiking LLMs.
- Extensive experiments on LLaMA models demonstrate that our framework successfully constructs temporal spiking LLMs with performance comparable to their original counterparts while achieving significant latency reduction.

### 2 RELATED WORKS

### 2.1 Spiking LLMs

The success of LLMs has motivated the development of SNN counterparts (spiking LLMs) that maintain energy efficiency while achieving comparable capabilities. Several approaches have emerged for creating spiking variants of transformer-based models (You et al., 2024; Zhou et al., 2022; 2023). SpikeGPT replaces traditional self-attention with Spiking RWKV mechanisms (Zhu et al., 2023). SpikingBERT employs a two-stage knowledge distillation method that utilizes pretrained BERT models as teachers to train spiking student architectures (Lv et al., 2023). Similarly, SpikingMiniLM builds upon BERT with parameter initialization and ANN-to-SNN distillation methods to achieve faster convergence during training. Recent work introduced SpikeLLM, scaling to 70 billion parameters through spike-driven quantization (Xing et al., 2024b;a). However, existing spiking LLMs rely exclusively on rate coding, where information is encoded through spike frequency. This leaves unexplored the potential of temporal-based spiking LLMs, which could achieve substantially lower energy consumption.

### 2.2 TEMPORAL-BASED A2S CONVERSIONS

While rate-based conversion methods have dominated ANN-to-SNN conversion research, temporal-based encoding approaches offer compelling advantages in terms of energy efficiency by leveraging precise spike timing rather than spike frequency. These methods include time-to-first spike (Thorpe et al., 2001), reverse coding (Zhang et al., 2019; Park et al., 2020), phase coding (Montemurro et al., 2008) and burst coding (Park et al., 2019). Among temporal coding schemes, time-to-first-spike (TTFS) coding has emerged as particularly promising, where each neuron emits at most one spike per time window with information encoded in the spike latency. Early TTFS-based conversion methods were developed by Rueckauer & Liu (2018) and further improved by Zhang et al. (2019) and Park et al. (2020), but these approaches introduced conversion errors across layers. A breakthrough came with Stanojevic et al. (2023; 2024), who demonstrated exact mapping from ReLU-based networks to SNNs using TTFS coding through a two-stage neuron activation process, achieving lossless conversion while maintaining energy benefits. Recently, Zhao et al. (2025) proposed TTFSFormer, the first TTFS-based conversion framework for Transformer architectures. However, existing TTFS-based conversion methods require extremely high latency to match continuous-time assumptions, preventing their implementation on large-scale models.

### 2.3 MODEL QUANTIZATION

Quantization has emerged as a critical technique for reducing model size and memory consumption, enabling efficient deployment of LLMs on resource-constrained devices (Shao et al., 2024), falling into two primary categories: quantization-aware training (QAT) (Liu et al., 2023) and post-training quantization (PTQ) (Xiao et al., 2023). QAT optimizes quantized weights during training using techniques like the straight-through estimator (Chen et al., 2024c; Du et al., 2024) but is computationally impractical for LLMs. PTQ has thus become the preferred approach, requiring only minimal calibration data while using dynamic activation quantization to address outlier-induced accuracy degradation (Frantar et al., 2023). Recent PTQ advances address outlier-induced errors using orthogonal transformations (QuaRot (Ashkboos et al., 2024), SpinQuant (Liu et al., 2024)) or dual

transformations (DuQuant (Lin et al., 2024)) to redistribute outliers across channels. However, these methods require computationally expensive per-token dynamic computation during inference. Pre-fixQuant (Chen et al., 2024b) offers an alternative by isolating token-wise outliers to enable efficient per-tensor static quantization, achieving comparable performance. While these quantization methods successfully achieve competitive performance with low-bit representations, energy consumption from dense matrix operations remains a fundamental barrier to edge deployment.

### 3 REVISITING TTFS-BASED ANN-TO-SNN CONVERSION

### 3.1 CONTINUOUS TTFS-BASED NEURONS

The activation process of continuous TTFS-based neurons is generally divided into two stages: the receiving phase and the firing phase (Zhao et al., 2025). At the i-th neuron in l-th layer, i=1,2,...,I and l=1,2,...,L. We denote the time range of the receiving phase as  $[t_{\rm recv}^{(l)},t_{\rm emit}^{(l)}]$ , and the emitting phase as  $[t_{\rm emit}^{(l)},t_{\rm end}^{(l)}]$ . With the initial membrane potential  $V(t_{\rm recv}^{(l)})=0$ , the continuous membrane potential dynamics are given by:

$$\frac{d}{dt}V(t) = \begin{cases}
\frac{1}{\tau_i^{(l)}} \left( \sum_j w_{ij}^{(l)} \eta_{ij}^{(l)} \left( t - t_j^{(l-1)} \right) + C_i^{(l)} \right), & t \in [t_{\text{recv}}^{(l)}, t_{\text{emit}}^{(l)}), \\
\psi_i^{(l)} \left( t - t_{\text{emit}}^{(l)} \right), & t \in [t_{\text{emit}}^{(l)}, t_{\text{end}}^{(l)}).
\end{cases} \tag{1}$$

The spike time  $t_j^{(l-1)}$  is received from the previous layer of the j-th input, while using the time range of the receiving phase from the previous layer as the time range for the firing phase of this layer, i.e.,  $t_{\rm recv}^{(l)} = t_{\rm emit}^{(l-1)}$  and  $t_{\rm emit}^{(l)} = t_{\rm end}^{(l-1)}$ ;  $w_{ij}^{(l)}$  are the weights; the input transform kernel function  $\eta_{ij}^{(l)}$  satisfies  $\eta_{ij}^{(l)}(u) = 0$ ,  $\forall u < 0$ ;  $\tau_i^{(l)} > 0$  is the time constant;  $C_i^{(l)}$  serves as a bias term; the output transform kernel function  $\psi_i^{(l)}$  is non-negative. Once the potential exceeds the threshold  $\theta_i^{(l)}$ , the neuron will emit a spike and record the spike firing time  $t_i^{(l)}$ . The relation between the spike time  $t_i^{(l)}$  and the corresponding activation value  $x_i^{(l)}$  of ANNs is:

$$x_i^{(l)} \tau_i^{(l)} = t_{\text{ref}}^{(l)} - t_i^{(l)}, \tag{2}$$

where  $t_{\mathrm{ref}}^{(l)}$  is the zero reference time. Therefore, the output range  $[a_i^{(l)}, b_i^{(l)}]$  can be expressed as:

$$a_i^{(l)} = \frac{1}{\tau_i^{(l)}} \left( t_{\text{ref}}^{(l)} - t_{\text{end}}^{(l)} \right), \qquad b_i^{(l)} = \frac{1}{\tau_i^{(l)}} \left( t_{\text{ref}}^{(l)} - t_{\text{emit}}^{(l)} \right). \tag{3}$$

We denote  $T^{(l)}=t_{\mathrm{end}}^{(l)}-t_{\mathrm{emit}}^{(l)}$  as the time window, and  $d_i^{(l)}=b_i^{(l)}-a_i^{(l)}$ .

### 3.2 PRACTICAL LIMITATIONS OF CONTINUOUS TTFS-BASED CONVERSION

The continuous TTFS-based conversion method (Zhao et al., 2025) establishes an equivalence between TTFS-based neurons and ANN neurons by modifying the input and output transform kernel functions, thereby enabling the mapping of TTFS-based SNNs to continuous ANNs:

**Theorems 4.1 and 4.3** in Zhao et al. (2025): Let  $f_{ij}:[a_i^{(l-1)},b_i^{(l-1)}]\to\mathbb{R}$  be differentiable functions and  $h:A\to\mathbb{R}$  be a differentiable monotone increasing function, and its inverse  $h^{-1}$  is well-defined on  $[a_i^{(l)},b_i^{(l)}]$ . If we let

$$\begin{split} &\eta_{ij}^{(l)}(s) = \begin{cases} f'_{ij} \bigg( \frac{s}{\tau_i^{(l-1)}} + a_i^{(l-1)} \bigg), & s \geq 0, \\ 0, & s < 0, \end{cases}, \ C_i^{(l)} = \sum_j w_{ij} \frac{f_{ij} \bigg( a_i^{(l-1)} \bigg)}{d_i^{(l-1)}}, \\ &\psi_i^{(l)}(s) = \frac{1}{\tau_i^{(l)} h' \bigg( h^{-1} \bigg( b_i^{(l)} - \frac{s}{\tau_i^{(l)}} \bigg) \bigg)}, \ \theta_i^{(l)} = h^{-1} (b_i^{(l)}) \end{split} \tag{4}$$

then the value  $x_i^{(l)}$  of ANNs represented by the output spike is

$$x_i^{(l)} = f^{(l)}(W^{(l)}; x_1^{(l-1)}, ..., x_I^{(l-1)}) = \operatorname{clip}\left(h\left(\sum_j w_{ij}^{(l)} f_{ij}\left(x_j^{(l-1)}\right)\right), a_i^{(l)}, b_i^{(l)}\right)$$
(5)

Although TTFS-based ANN-to-SNN conversion methods under continuous setting have been explored, their applications to LLMs remain limited in two aspects:

Infinite Clock Precision: For TTFS-based neurons under continuous setting, the spike time can be any real number (Stanojevic et al., 2023; Zhao et al., 2025; Stanojevic et al., 2024). At this point, the required clock precision is theoretically infinitely fine:  $\Delta t_{real} \rightarrow 0$ . However, electronic neuromorphic chips, which rely on discrete clock cycles, cannot provide infinitely fine clock precision (Deng et al., 2023). Consequently, TTFS coding based on continuous assumptions faces significant limitations in hardware implementations.

**Latency Overhead of Lossless Conversion:** In the continuous setting, TTFS-based lossless conversion methods establish **an equivalence between SNNs and continuous ANNs** and directly mapping the former to the latter. However, this process incurs extremely high latency (e.g., up to 4096 time steps), which propagates through the network and leads to prohibitively long inference delays.

In continuous settings, TTFS coding requires prohibitively high latency to achieve lossless conversion, resulting in excessively long inference delays for LLMs.

### 4 DISCRETIZATION-AWARE CONVERSION

In this section, we first construct discrete TTFS-based neurons to address the challenge of infinite clock precision. Next, rather than directly mapping TTFS-based SNNs to continuous ANNs in a continuous setting, we analyze the relationship between TTFS-based SNNs and discrete ANNs. We then examine the conversion error of discrete TTFS-based SNNs. Finally, we introduce the Quantization-Consistent ANN-to-SNN conversion method.

### 4.1 DISCRETE TTFS-BASED NEURONS

To overcome the challenge posed by infinite clock precision, we constructed a hardware-friendly discrete TTFS coding neuron model. Under the discrete time-step setting, the differential form of the original membrane potential equation can be approximated as follows:

$$\frac{d}{dt}V(t) = \frac{dV(t)}{dt_{real}(t)} \cdot \frac{dt_{real}(t)}{dt} \approx \frac{V(t+1) - V(t)}{t_{real}(t+1) - t_{real}(t)} \cdot \frac{d}{dt}t_{real}(t) = V(t+1) - V(t). \quad (6)$$

Building on the above discussion, we present a discretized version of TTFS-based neurons. At the i-th neuron in l-th layer, i=1,2,...,I and l=1,2,...,L. We denote the time range of the receiving phase as  $\{t_{\rm recv}^{(l)},\ldots,t_{\rm emit}^{(l)}\}$ , and the emitting phase as  $\{t_{\rm emit}^{(l)},\ldots,t_{\rm end}^{(l)}\}$ . With the initial membrane potential  $V(t_{\rm recv}^{(l)})=0$ , the discrete membrane potential dynamics are given by:

$$V(t+1) - V(t) = \begin{cases} \frac{1}{\tau_i^{(l)}} \left( \sum_j w_{ij}^{(l)} \eta_{ij}^{(l)} \left( t - t_j^{(l-1)} \right) + C_i^{(l)} \right) & t \in \{t_{\text{recv}}^{(l)}, \dots, t_{\text{emit}}^{(l)} - 1\}, \\ \psi_i^{(l)} \left( t - t_{\text{emit}}^{(l)} \right) & t \in \{t_{\text{emit}}^{(l)}, \dots, t_{\text{end}}^{(l)} - 1\}. \end{cases}$$
(7)

We denote  $T^{(l)}=t_{\mathrm{end}}^{(l)}-t_{\mathrm{emit}}^{(l)}$  as the time window, and  $d_i^{(l)}=b_i^{(l)}-a_i^{(l)}$ .

### 4.2 RELATIONSHIP BETWEEN DISCRETE TTFS-BASED SNNs AND ANNS

We theoretically establish the equivalence between TTFS-based SNNs and discrete ANNs. First, we determine the corresponding ANN function using the transform kernel functions and parameters of the TTFS-based neuron. For any TTFS-based neuron with fixed conversion functions and parameters, the corresponding ANN function can be identified:

**Theorem 1** For arbitrary fixed  $\eta_{ij}^{(l)}$ ,  $\psi_i^{(l)}$ ,  $C_i^{(l)}$  and  $\theta_i^{(l)}$  in SNNs with time window  $T^{(l)}$ , if we define  $S(t) = \sum_{v=0}^{t-t_{emit}^{(l)}-1} \psi_i^{(l)}(v)$  with  $t \in \left\{t_{emit}^{(l)}, \ldots, t_{end}^{(l)}\right\}$ , then the corresponding activation value of discrete ANNs is given by:  $x_i^{(l)} = f^{(l)}(W^{(l)}; x_1^{(l-1)}, \ldots, x_r^{(l-1)})$ (8)

$$=\frac{1}{\tau_i^{(l)}}\left(t_{\mathit{ref}}^{(l)}-S^{-1}\left(\theta_i^{(l)}+\Delta_i^{(l)}-\frac{1}{\tau_i^{(l)}}\sum_{j=1}^{I}\sum_{t=t_i^{(l-1)}}^{T_{\mathit{emit}}^{(l)}-1}w_{ij}^{(l)}\,\eta_{ij}^{(l)}\Big(x_j^{(l-1)}\tau_j^{(l-1)}+t-t_{\mathit{ref}}^{(l-1)}\Big)-T^{(l)}C_i^{(l)}\right)\right),$$

where  $W^{(l)} = (w_{ij}^{(l)})_{I \times I}$  is the weight matrix;  $\Delta_i^{(l)} \geq 0$  is a compensation constant, which is actually the difference between the  $\theta_i^{(l)}$  and the membrane potential at the spike time.

Next, we determine the corresponding transform kernel functions and parameters in the TTFS-based neuron using the ANN function. For any given fixed ANN function, the TTFS-based neuron with the corresponding transform kernel functions and parameters can be identified:

**Theorem 2** Let  $f_{ij}$  be a function with input set of discrete points between  $a_i^{(l-1)}$  and  $b_i^{(l-1)}$ , and h be a monotone increasing function with output set of discrete points between  $a_i^{(l)}$  and  $b_i^{(l)}$ . We denote  $u = t - t_j^{(l-1)}$  with  $t \in \left\{t_{recv}^{(l)}, \ldots, t_{emit}^{(l)}\right\}$ , and  $v = t - t_{emit}^{(l)}$  with  $t \in \left\{t_{emit}^{(l)}, \ldots, t_{end}^{(l)}\right\}$ . To represent the corresponding activation value of discrete ANNs:

$$x_i^{(l)} = f^{(l)}(W^{(l)}; x_1^{(l-1)}, ..., x_I^{(l-1)}) = \operatorname{clip}\left(h\left(\sum_j w_{ij}^{(l)} f_{ij}\left(x_j^{(l-1)}\right)\right), a_i^{(l)}, b_i^{(l)}\right). \tag{9}$$

we need to configure the SNN as follows:

$$\eta_{ij}^{(l)}(u) = \begin{cases} \tau_i^{(l-1)} \left( f_{ij} \left( \frac{u+1}{\tau_i^{(l-1)}} + a_i^{(l-1)} \right) - f_{ij} \left( \frac{u}{\tau_i^{(l-1)}} + a_i^{(l-1)} \right) \right) & u \ge 0, \\ 0 & u < 0. \end{cases}$$

$$\psi_i^{(l)}(v) = h^{-1}\left(b_i^{(l)} - \frac{v}{\tau_i^{(l)}}\right) - h^{-1}\left(b_i^{(l)} - \frac{v+1}{\tau_i^{(l)}}\right), \ C_i^{(l)} = \frac{\sum_j w_{ij}^{(l)} f_{ij}\left(a_i^{(l-1)}\right)}{d_i^{(l-1)}}, \ \theta_i^{(l)} = h^{-1}(b_i^{(l)}) + \Delta_i^{(l)}.$$

Furthermore, we demonstrate the equivalence between the discrete TTFS-based neuron and the quantization function:

**Corollary 1** *We define the processes of quantization and dequantization as follows:* 

$$\hat{\mathbf{X}}_i^{(l)} = \lambda_i^{(l)} \cdot \operatorname{clip}(\lfloor \frac{\mathbf{X}_i^{(l)}}{\lambda_i^{(l)}} \rfloor + z^{(l)}, 0, N) - \lambda_i^{(l)} \cdot z^{(l)}, \tag{10}$$

where  $\lambda_i^{(l)} = \frac{\max(\mathbf{X}_i^{(l)}) - \min(\mathbf{X}_i^{(l)})}{N}$  and  $z^{(l)} = -\lfloor \frac{\min(\mathbf{X}_i^{(l)})}{\lambda_i^{(l)}} \rfloor$  are scale and zero point values, respectively;  $\lfloor \cdot \rfloor$  denotes the floor operation;  $N = 2^n - 1$  denotes the quantization level and n denotes the quantization bits;  $\hat{\mathbf{X}}_i^{(l)}$  and  $\mathbf{X}_i^{(l)}$  are the dequantized and original tensor, respectively.

For a TTFS-based SNN defined in (7),  $\mathcal{H}$  is the Heaviside step function, if we set the  $\eta_{ij}^{(l)}$ ,  $\psi_i^{(l)}$ ,  $C_i^{(l)}$  and  $\theta_i^{(l)}$  as follow:

$$\eta_{ij}^{(l)}(u) = \mathcal{H}\left(\frac{u}{\tau_i^{(l-1)}} + a_i^{(l-1)}\right), \ \psi_i^{(l)}(v) = \frac{1}{\tau_i^{(l)}}, \ C_i^{(l)} = \sum_j \frac{a_i^{(l-1)}}{d_i^{(l-1)}} w_{ij}, \ \theta_i^{(l)} = b_i^{(l)}$$
(11)

and we let  $t_{emit}^{(l)} = 0$ ,  $t_{end}^{(l)} = N$ ,  $\tau_i^{(l)} = \frac{1}{\lambda_i^{(l)}}$ ,  $t_{end}^{(l)} - t_{ref}^{(l)} = z^{(l)}$ , and  $\mathbf{X}_i^{(l)} = \sum_{j=1}^{I} w_{ij}^{(l)} \, x_j^{(l-1)}$ . The output of spiking neural neuron and quantization function are equivalent, i.e,  $x_i^{(l)} = \hat{\mathbf{X}}_i^{(l)}$ .

325 326

327

328

330

331

332

337

338

339 340 341

342

343

344

345

346

347 348

349 350

351

352

353

354

355

356

357

358

359 360

361 362

364

365 366

367 368

369

370

371

372

373

374

375

376

377

### ERROR ANALYSIS FOR DISCRETE TTFS-BASED SNNS

In the continuous setting, although TTFS-based SNNs enable lossless conversion to ANNs, they require infinitely fine clock precision for hardware implementation and introduce significantly long inference delay in the network. We analyze the conversion error of discrete TTFS-based SNNs.

**Theorem 3** The error analysis of TTFS-based SNNs: Let  $T^{(l)}$  denotes the time window with the corresponding clock time constant  $\Omega$ , the derivatives of the function h and its inverse are bounded by  $G_1$  and  $G_2$ , I denotes the number of neurons in each layer of the network, and L denotes the

number of layers, 
$$T = \min \left\{ T^{(l)} \right\}_{l=1}^{L}$$
, and  $\tau = \max \left\{ \left\{ \tau_i^{(l)} \right\}_{i=1}^{L} \right\}_{l=1}^{L}$ ,  $\alpha_i^{(l)}$  is the corresponding

number of layers,  $T = \min \left\{ T^{(l)} \right\}_{l=1}^{L}$ , and  $\tau = \max \left\{ \left\{ \tau_i^{(l)} \right\}_{i=1}^{I} \right\}_{l=1}^{L}$ ,  $\alpha_i^{(l)}$  is the corresponding output of ANNs and  $\rho = \max_{\{i,l\}} \left\{ \left| \alpha_i^{(l)} - \frac{a_i^{(l)} + b_i^{(l)}}{2} \right| \right\}$ . The conversion error of the TTFS-based SNNs in car had be smalled. SNNs in can be bounded as:

$$\mathcal{E} \le LI \cdot \max\left(\rho - \frac{T}{2\tau}, 0\right) + \frac{LIG_1G_2\Omega}{T} \tag{12}$$

**Remark 1** In Theorem 3: The first term captures the clipping error in the TTFS-based SNNs, which can be eliminated by increasing the time window T. As T increases, the output range of TTFS-based SNNs expands. When this range encompasses the output of ANNs, the clipping error is eliminated; The second term reflects the quantization error, which can only be alleviated by increasing T. As T increases, the output range of TTFS-based SNNs becomes finer, facilitating better alignment between the output of ANNs and the discrete points of the SNNs' output, thereby reducing quantization error. Thus, achieving high accuracy TTFS-based SNNs necessitates sufficiently long time windows.

### QUANTIZATION-CONSISTENT ANN-TO-SNN CONVERSION

Our goal is to develop high-accuracy, low-latency temporal-coding spiking LLMs. Achieving high accuracy in temporal-coding spiking LLMs typically requires extending the time window, which in turn increases latency. This latency propagates through the network, leading to excessive inference delays. To address this challenge, we propose the Quantization-Consistent ANN-to-SNN (QC-A2S) conversion method, which leverages the equivalence between TTFS-based SNNs and discrete ANNs. Our approach combines low-bit quantization with discretization-compatible TTFS neurons, enabling low-latency temporal-coding spiking LLMs. Specifically, we first apply established techniques, such as post-training quantization, to minimize clipping and quantization errors, resulting in a low-bit, high-accuracy baseline model. We then map the quantized LLM to an equivalent spiking LLM, achieving a low-latency, high-accuracy temporal-coding spiking LLM.

### 5 EXPERIMENT

In this section, we conduct experiments to validate the effectiveness of our proposed method and compare its performance, computational count, and energy consumption with those of different approaches. Additionally, we conduct ablation studies on various latency.

### 5.1 IMPLEMENT DETAILS

**Datasets and Underlying Models** In the experiments, two types of benchmarks are used. For accuracy-oriented evaluation, five representative reasoning datasets are adopted, namely PIQA(Bisk et al., 2020), HellaSwag (Zellers et al., 2019), WinoGrande (Sakaguchi et al., 2021), ARC-Easy, and ARC-Challenge (Clark et al., 2018). **PIQA** targets physical commonsense reasoning in everyday scenarios, ARC-Easy and ARC-Challenge consist of science exam questions with varying difficulty levels, **HellaSwag** evaluates contextual understanding through plausible continuation tasks, and **WinoGrande** focuses on large-scale pronoun resolution for commonsense reasoning. For perplexity-oriented evaluation, we additionally use five widely adopted language modeling datasets, including C4 (Raffel et al., 2020), The Pile (Gao et al., 2021), Penn Treebank (PTB) (Marcus et al., 1993), WikiText-2 (Merity et al., 2017), and RedPajama (Together Computer, 2023). The datasets were preprocessed following standard practices, and data augmentation techniques were applied

where appropriate. In our study, all methods are applied to the LLaMA family of LLMs as the common backbone. We consider a range of representative models, including LLaMA-2-7B, LLaMA-2-13B, LLaMA-3-8B, and LLaMA-2-70B.

381 382

**Baselines** We compare our approach against several representative baselines that adapt large language models through either quantization or ANN-to-SNN conversion:

384 385

• **PrefixQuant** (Chen et al., 2025b) is a weight–activation quantization method that addresses token-wise outliers in the KV cache and employs lightweight blockwise training, achieving strong performance across different precision levels.

386 387 388

389

• SpikeLLM (Xing et al., 2024a) presents the first spiking LLMs by incorporating bioinspired spiking mechanisms with generalized integrate-and-fire neurons, yielding improvements in perplexity and reasoning accuracy compared to quantized LLMs.

390 391 392

• TTFSFormer (Zhao et al., 2025) applies time-to-first-spike coding to Transformers, extending TTFS neurons to handle nonlinear layers and achieving competitive accuracy with significantly reduced energy consumption.

393 394

396

397

398

399

**Experiment Configurations** All experiments were conducted on a server equipped with NVIDIA A100 GPUs (80 GB of memory), Intel Xeon CPUs, and 512 GB of RAM. The models were implemented in PyTorch 2.6 with CUDA 12.4 support. For fair comparison, all baseline methods were re-implemented or run using their officially released code under the same environment and hyperparameter settings whenever possible. In addition to the hardware information mentioned in the main text, we provide further details about the reproduction of baselines here. We adopt 8 bits for weight, 6 bits for activation quantization, i.e. W8A6, for SpikeLLM(Xing et al., 2024a) and PrefixQuant(Chen et al., 2025b), and use 8192 time precision for TTFSFormer (Zhao et al., 2025).

404

### 5.2 Main Results

410

411

412

Tables 1 and 2 report the accuracy and PPL metrics of all methods on the LLaMA-2-7B, LLaMA-2-13B, and LLaMA-3-8B models. The results indicate that: (i) temporally encoded spiking LLMs achieve performance comparable to quantized LLMs across all LLaMA models, providing further empirical evidence for the equivalence between TTFS-based SNNs and quantized ANNs; (ii) our method substantially outperforms TTFSFormer under low-latency settings, while TTFSFormer continues to exhibit unsatisfactory performance even at higher latencies, underscoring the excessive latency demands of continuous TTFS-based SNNs; and (iii) our model surpasses the state-of-the-art spiking LLM (SpikeLLM), further validating the effectiveness of the proposed approach.

413 414 415

416

417 418 Model

Method

Table 1: We report accuracy for WinoGrande and acc\_norm for HellaSwag, ArcC, ArcE, and PIQA on LLaMA-2-7B, LLaMA-2-13B, and LLaMA-3-8B models Precision WinoGrande HellaSwag ArcC

Δνσ

Model	Memou	Frecision	WilloGranue	HeliaSwag	AICC	AICE	FIQA	Avg.
2-7B	Baseline	FP16	69.22	76.00	46.33	74.62	79.11	69.06
	TTFSFormer	T=8192	50.04	25.49	<sup>-</sup> 2 <del>6</del> .88	26.81	50.82	36.01
	SpikeLLM	W8A6	65.51	73.61	42.49	70.16	75.41	65.44
	PrefixQ	W8A6	67.96	75.25	45.14	72.18	77.86	67.68
	Ours	T=64	69.06	75.30	44.97	71.97	78.13	67.89
3-8B	Baseline	FP16	72.69	79.16	53.41	77.69	80.69	72.73
	TTFSFormer	T=8192	52.41	26.86	$^{-}2\bar{5}.\bar{7}\bar{7}$	$\bar{2}\bar{4}.\bar{7}5^{-}$	51.09	36.18
	SpikeLLM	W8A6	58.25	59.28	32.34	53.37	68.66	54.38
	PrefixQ	W8A6	70.01	76.09	50.00	78.91	78.62	70.73
	Ours	T=64	70.56	75.89	51.28	78.16	77.64	70.71
2-13B	Baseline	FP16	72.38	79.38	49.06	77.53	80.52	71.77
	TTFSFormer	T=8192	48.70	26.29	<sup>-</sup> 2 <del>6</del> . <u>1</u> 1	$\bar{25}.\bar{72}$	51.25	35.61
	SpikeLLM	W8A6	68.03	76.76	44.88	73.32	77.48	68.09
	PrefixQ	W8A6	70.40	73.46	45.22	71.63	77.37	67.62
	Ours	T=64	70.32	73.47	45.31	71.68	77.48	67.65

Table 2: We report *Perplexity* for C4, Pile, PTB, WikiText2, and RedPajama on LLaMA-2-7B, LLaMA-2-13B, and LLaMA-3-8B models

Model	Method	Precision	C4	Pile	PTB	WikiText2	RedPajama	Avg.
	Baseline	FP16	6.98	4.63	37.92	5.47	5.61	12.12
	TTFSFormer	T=8192	>100	->100	- >100	>100	<del>-</del> - <del>-</del> 100	-5100
2-7B	SpikeLLM	W8A6	7.89	5.14	57.27	6.43	6.21	16.59
	PrefixQ	W8A6	7.24	4.76	56.83	5.68	5.84	16.07
	Ours	T=64	7.24	4.76	56.96	5.68	5.84	16.10
	Baseline	FP16	8.88	5.52	11.18	6.14	7.83	7.91
	TTFSFormer	T=8192	>100	->100	- > <del>1</del> 0 <del>0</del> -	>100	<del>-</del> - <del>-</del> 100	- > 100
3-8B	SpikeLLM	W8A6	> 100	> 100	> 100	>100	>100	> 100
	PrefixQ	W8A6	10.79	6.56	13.67	7.51	9.28	9.56
	Ours	T=64	10.81	6.57	13.70	7.53	9.30	9.58
	Baseline	FP16	6.47	4.34	50.94	4.88	5.19	14.36
2-13B	TTFSFormer	T=8192	>100	->100	- > <del>1</del> 0 <del>0</del> -	>100	>100	->100
	SpikeLLM	W8A6	7.16	4.74	62.07	5.81	5.53	17.06
	PrefixQ	W8A6	8.08	5.27	42.55	6.30	6.50	13.74
	Ours	T=64	8.08	5.27	42.68	6.30	6.51	13.77

### 5.3 Comparison of accuracy under different latency configurations

In Table 3, we conduct a detailed study of how latency influences the performance of temporally coded spiking LLMs using the LLaMA-2-7B model. The results reveal a clear trend: increasing latency consistently improves accuracy across all evaluated benchmarks. This indicates that longer time windows allow TTFS-based SNNs to better approximate the activations of ANNs, thereby reducing discretization-induced errors and enhancing representational fidelity. Such evidence provides strong empirical support for our theoretical analysis in Theorem 3 ( $\mathcal{E} \leq LI \cdot \max\left(\rho - \frac{T}{2\tau}, 0\right) + \frac{LIG_1G_2\Omega}{T}$ ), which establishes that achieving high accuracy in TTFS-based SNNs is inherently dependent on sufficiently long latency (T).

Table 3: The accuracy of our method for LLaMA-2-7B under different latency configurations (Latency corresponds to time window).

Latency	WinoGrande	HellaSwag	ARC-Challenge	ARC-Easy	PIQA	Average
$16(2^4)$	65.59	68.35	38.48	65.95	73.45	62.36
$64(2^6)$	69.06	75.30	44.97	71.97	78.13	67.89
$256(2^8)$	70.17	76.41	45.14	73.57	78.35	68.73
$1024(2^{10})$	70.72	76.49	45.56	73.99	78.13	68.98

## 6 Conclusion

LLMs have achieved remarkable success, but they also introduce severe energy bottlenecks that hinder their sustainable deployment. SNNs provide a promising pathway toward energy-efficient spiking LLMs through ANN-to-SNN conversion. Among various spike-coding schemes, TTFS coding is particularly appealing, as it conveys information with a single spike, thereby further reducing energy consumption. Existing TTFS-based A2S conversions depend on continuous-time assumptions and require prohibitively large latencies to approximate the continuous values of ANNs. This reliance results in unacceptable inference delays in deep models, particularly LLMs, creating significant obstacles to the development of practical temporal-coding spiking LLMs.

To overcome this challenge, we propose a discretization-aware theoretical framework that establishes a precise correspondence between discrete TTFS-based neurons and ANNs. Our key insight shows that conversion errors are constrained by latency-dependent terms. Building on this, we introduce the QC-A2S conversion method, which combines low-bit quantization with discretization-compatible TTFS neurons, enabling low-latency temporal-coding spiking LLMs.

### ETHICS STATEMENT

All participants in this work, as well as the paper submission, adhere to the ICLR Code of Ethics (https://iclr.cc/public/CodeOfEthics).

### REPRODUCIBILITY STATEMENT

We affirm that the results of this work are fully reproducible. Appendix D provides the theoretical proofs. Appendix B.1 details the experimental implementations, and the source code will be publicly released after publication of the paper.

### REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. arXiv preprint arXiv:2303.08774, 2023.
- Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian Croci, Bo Li, Pashmina Cameron, Martin Jaggi, Dan Alistarh, Torsten Hoefler, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. *Advances in Neural Information Processing Systems*, 37:100213–100240, 2024.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pp. 7432–7439, 2020.
- Lingjiao Chen, Jared Quincy Davis, Boris Hanin, Peter Bailis, Ion Stoica, Matei Zaharia, and James Zou. Are more llm calls all you need? towards scaling laws of compound inference systems. arXiv preprint arXiv:2403.02419, 2024a.
- Long Chen, Xiaotian Song, Andy Song, BaDong Chen, Jiancheng Lv, and Yanan Sun. Fas: Fast ann-snn conversion for spiking large language models. *arXiv preprint arXiv:2502.04405*, 2025a.
- Mengzhao Chen, Yi Liu, Jiahao Wang, Yi Bin, Wenqi Shao, and Ping Luo. Prefixquant: Static quantization beats dynamic through prefixed outliers in llms. *arXiv preprint arXiv:2410.05265*, 2024b.
- Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, and Ping Luo. Efficientqat: Efficient quantization-aware training for large language models. *arXiv preprint arXiv:2407.11062*, 2024c.
- Mengzhao Chen, Yi Liu, Jiahao Wang, Yi Bin, Wenqi Shao, and Ping Luo. Prefixquant: Eliminating outliers by prefixed tokens for large language models quantization, 2025b. URL https://arxiv.org/abs/2410.05265.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *ArXiv*, abs/1803.05457, 2018. URL https://api.semanticscholar.org/CorpusID:3922816.
- Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *Ieee Micro*, 38(1):82–99, 2018.
- Bin Deng, Yanrong Fan, Jiang Wang, and Shuangming Yang. Auditory perception architecture with spiking neural network and implementation on fpga. *Neural Networks*, 165:31–42, 2023.
- Dayou Du, Yijia Zhang, Shijie Cao, Jiaqi Guo, Ting Cao, Xiaowen Chu, and Ningyi Xu. Bitdistiller: Unleashing the potential of sub-4-bit llms via self-distillation. *arXiv preprint arXiv:2402.10631*, 2024.

- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pp. arXiv–2407, 2024.
  - Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan-Adrian Alistarh. Optq: Accurate post-training quantization for generative pre-trained transformers. In *11th International Conference on Learning Representations*, 2023.
  - Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2021.
  - Robert Gütig and Haim Sompolinsky. The tempotron: a neuron that learns spike timing—based decisions. *Nature neuroscience*, 9(3):420–428, 2006.
  - Zecheng Hao, Tong Bu, Jianhao Ding, Tiejun Huang, and Zhaofei Yu. Reducing ann-snn conversion error through residual membrane potential. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 11–21, 2023.
  - Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. arXiv preprint arXiv:2203.15556, 2022.
  - Yizhou Jiang, Kunlin Hu, Tianren Zhang, Haichuan Gao, Yuqian Liu, Ying Fang, and Feng Chen. Spatio-temporal approximation: A training-free snn conversion for transformers. In *The twelfth international conference on learning representations*, 2024.
  - Haokun Lin, Haobo Xu, Yichen Wu, Jingzhi Cui, Yingtao Zhang, Linzhan Mou, Linqi Song, Zhenan Sun, and Ying Wei. Duquant: Distributing outliers via dual transformation makes stronger quantized llms. *Advances in Neural Information Processing Systems*, 37:87766–87800, 2024.
  - Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023.
  - Zechun Liu, Changsheng Zhao, Igor Fedorov, Bilge Soran, Dhruv Choudhary, Raghuraman Krishnamoorthi, Vikas Chandra, Yuandong Tian, and Tijmen Blankevoort. Spinquant: Llm quantization with learned rotations. *arXiv preprint arXiv:2405.16406*, 2024.
  - Changze Lv, Tianlong Li, Jianhan Xu, Chenxi Gu, Zixuan Ling, Cenyuan Zhang, Xiaoqing Zheng, and Xuanjing Huang. Spikebert: A language spikformer trained with two-stage knowledge distillation from bert. 2023.
  - Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9):1659–1671, 1997.
  - Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of english: The penn treebank. *Computational Linguistics*, 19(2):313–330, 1993.
  - Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=Byj72udxe.
  - Paul A Merolla, John V Arthur, Rodrigo Alvarez-Icaza, Andrew S Cassidy, Jun Sawada, Filipp Akopyan, Bryan L Jackson, Nabil Imam, Chen Guo, Yutaka Nakamura, et al. A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science*, 345 (6197):668–673, 2014.
  - Marcelo A Montemurro, Malte J Rasch, Yusuke Murayama, Nikos K Logothetis, and Stefano Panzeri. Phase-of-firing coding of natural visual stimuli in primary visual cortex. *Current biology*, 18 (5):375–380, 2008.

- Bhaskar Mukhoty, Velibor Bojkovic, William de Vazelhes, Xiaohan Zhao, Giulia De Masi, Huan Xiong, and Bin Gu. Direct training of snn using local zeroth order method. *Advances in Neural Information Processing Systems*, 36:18994–19014, 2023.
  - Seongsik Park, Seijoon Kim, Hyeokjun Choe, and Sungroh Yoon. Fast and efficient information transmission with burst spikes in deep spiking neural networks. In *Proceedings of the 56th Annual Design Automation Conference 2019*, pp. 1–6, 2019.
  - Seongsik Park, Seijoon Kim, Byunggook Na, and Sungroh Yoon. T2fsnn: deep spiking neural networks with time-to-first-spike coding. In 2020 57th ACM/IEEE design automation conference (DAC), pp. 1–6. IEEE, 2020.
  - Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
  - Kaushik Roy, Akhilesh Jaiswal, and Priyadarshini Panda. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575(7784):607–617, 2019.
  - Bodo Rueckauer and Shih-Chii Liu. Conversion of analog to spiking neural networks using sparse temporal coding. In 2018 IEEE international symposium on circuits and systems (ISCAS), pp. 1–5. IEEE, 2018.
  - Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
  - Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omniquant: Omnidirectionally calibrated quantization for large language models. 2024. URL https://openreview.net/forum?id=8Wuvhh0LYW.
  - Ana Stanojevic, Stanisław Woźniak, Guillaume Bellec, Giovanni Cherubini, Angeliki Pantazi, and Wulfram Gerstner. An exact mapping from relu networks to spiking neural networks. *Neural Networks*, 168:74–88, 2023.
  - Ana Stanojevic, Stanisław Woźniak, Guillaume Bellec, Giovanni Cherubini, Angeliki Pantazi, and Wulfram Gerstner. High-performance deep spiking neural networks with 0.3 spikes per neuron. *Nature Communications*, 15(1):6793, 2024.
  - Simon Thorpe, Arnaud Delorme, and Rufin Van Rullen. Spike-based strategies for rapid processing. *Neural networks*, 14(6-7):715–725, 2001.
  - Together Computer. Redpajama: Reproducible pretraining data. https://www.together.xyz/blog/redpajama, 2023.
  - Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
  - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
  - Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pp. 38087–38099. PMLR, 2023.
- 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, 2024a.
- Xingrun Xing, Zheng Zhang, Ziyi Ni, Shitao Xiao, Yiming Ju, Siqi Fan, Yequan Wang, Jiajun Zhang, and Guoqi Li. Spikelm: Towards general spike-driven language modeling via elastic bi-spiking mechanisms. In *International Conference on Machine Learning*, pp. 54698–54714. PMLR, 2024b.

- 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.
  - Kang You, Zekai Xu, Chen Nie, Zhijie Deng, Qinghai Guo, Xiang Wang, and Zhezhi He. Spikeziptf: conversion is all you need for transformer-based snn. In *Proceedings of the 41st International Conference on Machine Learning*, pp. 57367–57383, 2024.
  - Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Annual Meeting of the Association for Computational Linguistics*, 2019. URL https://api.semanticscholar.org/CorpusID:159041722.
  - Lei Zhang, Shengyuan Zhou, Tian Zhi, Zidong Du, and Yunji Chen. Tdsnn: From deep neural networks to deep spike neural networks with temporal-coding. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pp. 1319–1326, 2019.
  - Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. Opt: Open pre-trained transformer language models. *arXiv preprint arXiv:2205.01068*, 2022.
  - Lusen Zhao, Zihan Huang, Jianhao Ding, and Zhaofei Yu. TTFSFormer: A TTFS-based lossless conversion of spiking transformer. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=mJAa823xKu.
  - Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv* preprint arXiv:2303.18223, 1(2), 2023.
  - Chenlin Zhou, Liutao Yu, Zhaokun Zhou, Zhengyu Ma, Han Zhang, Huihui Zhou, and Yonghong Tian. Spikingformer: Spike-driven residual learning for transformer-based spiking neural network. *arXiv preprint arXiv:2304.11954*, 2023.
  - Chenlin Zhou, Han Zhang, Liutao Yu, Yumin Ye, Zhaokun Zhou, Liwei Huang, Zhengyu Ma, Xiaopeng Fan, Huihui Zhou, and Yonghong Tian. Direct training high-performance deep spiking neural networks: a review of theories and methods. *Frontiers in Neuroscience*, 18:1383844, 2024.
  - Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang, Shuicheng Yan, Yonghong Tian, and Li Yuan. Spikformer: When spiking neural network meets transformer. *arXiv* preprint *arXiv*:2209.15425, 2022.
  - Rui-Jie Zhu, Qihang Zhao, Guoqi Li, and Jason K Eshraghian. Spikegpt: Generative pre-trained language model with spiking neural networks. *arXiv preprint arXiv:2302.13939*, 2023.

## A NOTIONS

Table 4: Symbol Definitions

Symbol	Definition	Symbol	Definition
l	Layer index	X	Inputs of QANN
i, j	Neuron index	$\hat{\mathbf{X}}$	Output of QANN
W	Weight matrix	$a_i^{(l)}$	Output lower bound of TTFS-based neuron
$t_{ m recv}$	Receiving time step	$b_i^{(l)}$	Output upper bound of TTFS-based neuron
$t_{ m emit}$	Emitting time step	I	The number of neurons in each layer
$t_{ m end}$	End time-step	$\theta$	Threshold
$\overline{\mathcal{H}}$	Heaviside function	C	Bias term in TTFS-based neuron
$\overline{t}$	Time step index	au	Time constant in TTFS-based neuron
$\overline{\eta}$	Input transform kernel	N	Quantization level
$\psi$	Output transform kernel	n	Quantization bits

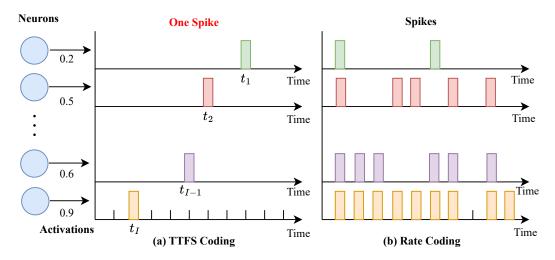


Figure 3: TTFS Coding vs. Rate Coding.

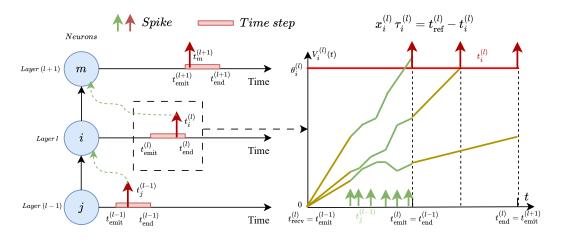


Figure 4: The process of TTFS-based spiking neural neurons.

## B EXPERIMENTAL SUPPLEMENTARY

### **B.1** EXPERIMENT CONFIGURATIONS

In addition to the hardware information mentioned in the main text, we provide further details about the reproduction of baselines here. We adopt 8 bits for weight, 6 bits for activation quantization, *i.e.* W8A6, for SpikeLLM(Xing et al., 2024a) and PrefixQuant(Chen et al., 2025b), and use 8192 time precision for TTFSFormer (Zhao et al., 2025).

### C CONVERSION ERROR

In this section, we provide a detailed analysis of the conversion error between the ANN and the converted TTFS-based SNN across layers. We assume that both the ANN and SNN receive the same input from layer l-1, i.e.,  $\alpha^{(l-1)}=x^{(l-1)}$ , and then analyze the error in layer l.

**ANN neurons.** For ANNs, the output  $\alpha^l$  of neurons in layer l is realized by a linear weighting  $W^{(l)}$  and nonlinear mappings  $f(\cdot)$ :

$$\alpha^{(l)} = f\left(W^{(l)}\alpha^{(l-1)}\right),\tag{13}$$

**SNN neurons.** For TTFS-based SNNs, we consider the relation between the spike time  $t^{(l)}$  of SNN and the corresponding activation value  $x^{(l)}$  of ANN:

$$x^{(l)} = \frac{1}{\tau^{(l)}} \left( t_{\text{ref}}^{(l)} - t^{(l)} \right). \tag{14}$$

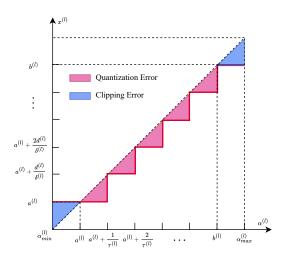


Figure 5: Clipping error and quantization error.

From Eqs.(13) and (14), along with the conditions  $V\left(t_{\mathrm{emit}}^{(l)}\right) = W^{(l)} \cdot f\left(x^{(l-1)}\right)$  and  $V(t^{(l)}-1) < \theta^{(l)} \leq V(t^{(l)}), \ t^{(l)} \in \{t_{\mathrm{emit}}^{(l)}, \dots, t_{\mathrm{end}}^{(l)}\}$ , it follows that a transformation between the temporal domain (relate to  $t^{(l)}$ ) and the numerical domain (relate to  $x^{(l)}$ ) enables the activation value  $a^{(l)}$  of analog neurons in the ANN to be mapped onto  $x^{(l)}$  in the TTFS-based SNN. Because the output ranges and types of SNNs and ANNs differ, conversion errors are generally unavoidable. During the ANN-to-SNN conversion, two primary sources of error, clipping error  $\mathcal{E}_{\mathrm{c}}^{(l)}$  and quantization error  $\mathcal{E}_{\mathrm{q}}^{(l)}$ , both of which contribute to the performance gap between ANNs and SNNs.

For layer l, the total error decomposes as:

$$\mathcal{E}^{(l)} = \mathcal{E}_{c}^{(l)} + \mathcal{E}_{q}^{(l)} \tag{15}$$

Clipping error. Clipping error denotes the error caused by different value ranges of ANNs and SNNs. For an temporal coding spiking neural neuron, when the time steps  $T^{(l)}$  are fixed, the output of SNN:  $x^{(l)}$  is in the range of  $\left[a^{(l)},b^{(l)}\right]$ , where  $a^{(l)}=\frac{t_{\mathrm{ref}}^{(l)}-t_{\mathrm{end}}^{(l)}}{\tau^{(l)}}$  and  $b^{(l)}=\frac{t_{\mathrm{ref}}^{(l)}-t_{\mathrm{emit}}^{(l)}}{\tau^{(l)}}$ . We define the  $\alpha_{max}$  as the maximum value in  $\alpha^{(l)}$ ,  $\alpha_{min}$  is the minimum value in in  $\alpha^{(l)}$ . Then the output  $\alpha \in \left[\alpha_{min},a^{(l)}\right]$  of ANNs will be mapped to the same value  $a^{(l)}$ , the output  $\alpha \in \left[b^{(l)},\alpha_{max}\right]$  of ANNs will be mapped to the same value  $a^{(l)}$ , which will cause conversion error named clipping error  $\mathcal{E}_{\mathrm{c}}$ .

Quantization error. The output spike time  $t^{(l)}$  is discrete, so the final output  $x^{(l)} = \frac{1}{\tau^{(l)}} \left( t_{\rm ref}^{(l)} - t^{(l)} \right)$  is also discrete, while the output activation value  $\alpha$  of the ANNs is continuous. Therefore, when mapping  $\alpha^{(l)}$  to  $x^{(l)}$ , there will be unavoidable error related to temporal resolution, named quantization error  $\mathcal{E}_{\rm q}$ . For example, when the output of ANNs satisfies  $\alpha \in \left[ \frac{t^{(l)}d^{(l)}}{T^{(l)}}, \frac{(t^{(l)}+1)d^{(l)}}{T^{(l)}} \right), t^{(l)} = t_{\rm ref}^{(l)} - t_{\rm end}^{(l)}, t_{\rm ref}^{(l)} - t_{\rm end}^{(l)} + 1, ..., t_{\rm ref}^{(l)} - t_{\rm emit}^{(l)} - 1$ , the corresponding mapped value of SNN will be  $\frac{t^{(l)}d^{(l)}}{T^{(l)}}$ .

**Lemma 1** Analysis for clipping error. For a target ANN's output  $\alpha_i^{(l)}$ , the clipping error between the output of ANN and SNN is:

$$\mathcal{E}_{c}^{(l)} = \begin{cases} \left\| \alpha_{i}^{(l)} - b_{i}^{(l)} \right\| & \text{if } \alpha_{i}^{(l)} > b_{i}^{(l)} \\ 0 & \text{if } \alpha_{i}^{(l)} \in [a_{i}^{(l)}, b_{i}^{(l)}] \\ \left\| a_{i}^{(l)} - \alpha_{i}^{(l)} \right\| & \text{if } \alpha_{i}^{(l)} < a_{i}^{(l)}, \end{cases}$$

$$(16)$$

where 
$$a_i^{(l)} = \frac{t_{\mathit{ref}}^{(l)} - t_{\mathit{end}}^{(l)}}{\tau_i^{(l)}}$$
 and  $b_i^{(l)} = \frac{t_{\mathit{ref}}^{(l)} - t_{\mathit{emit}}^{(l)}}{\tau_i^{(l)}}$ .

**Lemma 2** Upper bound for quantization error: In the theoretical analysis under the continuous setting of Theorem 4.1 and Theorem 4.3 in Zhao et al. (2025), we denote the output under continuous coding as  $y_i^{(l)} \in [a_i^{(l)}, b_i^{(l)}]$  corresponds to the ANN with continuous outputs, while in the practically deployable discrete coding scenario, the SNN output is denoted as  $x_i^{(l)} \in [a_i^{(l)}, b_i^{(l)}]$  corresponds to the ANN with discretized outputs. That is, there exists quantization error  $\mathcal{E}_q^{(l)}$  between TTFS coding in the continuous setting and its practical deployment. Let  $T^{(l)}$  be the time window and  $\Omega$  be the corresponding clock time, and the derivatives of the function h and its inverse are bounded by  $G_1$  and  $G_2$ . Then, the quantization error can be bounded as:

$$\mathcal{E}_q^{(l)} = \left\| x_i^{(l)} - y_i^{(l)} \right\| \le \frac{G_1 G_2 \Omega}{T^{(l)}}.$$
 (17)

### D PROOFS

### **proof 1** Proof of Theorem 1:

For arbitrary fixed  $\eta_{ij}^{(l)}$ ,  $\psi_i^{(l)}$ ,  $C_i^{(l)}$  and  $\theta_i^{(l)}$  in SNNs with time window  $T^{(l)}$ , in the receiving phase:

$$\begin{split} &V\left(t_{emit}^{(l)}\right)\\ &=\frac{1}{\tau_{i}^{(l)}}\sum_{t=T_{recv}^{(l)}}^{T_{emit}^{(l)}-1}\left(\sum_{j=1}^{I}w_{ij}^{(l)}\,\eta_{ij}^{(l)}\left(t-t_{j}^{(l-1)}\right)-C_{i}^{(l)}\right)\\ &=\frac{1}{\tau_{i}^{(l)}}\sum_{j=1}^{I}\sum_{t=t_{i}^{(l-1)}}^{T_{emit}^{(l)}-1}w_{ij}^{(l)}\,\eta_{ij}^{(l)}\left(x_{j}^{(l-1)}\tau_{j}^{(l-1)}+t-t_{ref}^{(l-1)}\right)+T^{(l)}C_{i}^{(l)}. \end{split} \tag{18}$$

In the emitting phase, let  $\Delta_i^{(l)} \geq 0$  is a compensation constant, which is actually the difference between the  $\theta_i^{(l)}$  and the membrane potential at the spike time. We can get:

$$V\left(t_{emit}^{(l)}\right) + \sum_{v=0}^{t_i^{(l)} - t_{emit}^{(l)} - 1} \psi_i^{(l)}(v) - \Delta_i^{(l)} = \theta_i^{(l)}. \tag{19}$$

We denote  $S(t) = \sum_{v=0}^{t-t_{emit}^{(l)}-1} \psi_i^{(l)}(v)$ , then:

$$S(t_i^{(l)}) = \theta_i^{(l)} + \Delta_i^{(l)} - V(t_{emit}^{(l)}).$$
 (20)

Then:

$$t_i^{(l)} = S^{-1} \left( \theta_i^{(l)} + \Delta_i^{(l)} - V \left( t_{\textit{emit}}^{(l)} \right) \right). \tag{21}$$

According to the relationship between  $x_i^{(l)}$  and  $t_i^{(l)}$ , we can get:

$$x_i^{(l)} = \frac{1}{\tau_i^{(l)}} \left( t_{ref}^{(l)} - S^{-1} \left( \theta_i^{(l)} + \Delta_i^{(l)} - V \left( t_{emit}^{(l)} \right) \right) \right). \tag{22}$$

let  $W^{(l)} = (w_{ij}^{(l)})_{I \times I}$  is the weight matrix:

$$\begin{split} x_i^{(l)} &= f^{(l)}(W^{(l)}; x_1^{(l-1)}, ..., x_I^{(l-1)}) \\ &= \frac{1}{\tau_i^{(l)}} \left( t_{\textit{ref}}^{(l)} - S^{-1} \left( \theta_i^{(l)} + \Delta_i^{(l)} - \frac{1}{\tau_i^{(l)}} \sum_{j=1}^{I} \sum_{t=t_j^{(l-1)}}^{T_{\textit{emit}}^{(l)} - 1} w_{ij}^{(l)} \, \eta_{ij}^{(l)} \left( x_j^{(l-1)} \tau_j^{(l-1)} + t - t_{\textit{ref}}^{(l-1)} \right) - T^{(l)} C_i^{(l)} \right) \right), \end{split}$$

**proof 2** Proof of Theorem 2:

Consider the potential change in the receiving stage.

$$V_{i}\left(t_{emit}^{(l)}\right)$$

$$= \frac{1}{\tau_{i}^{(l-1)}} \sum_{t=t_{emit}^{(l-1)}-1}^{t_{emit}^{(l-1)}-1} \sum_{j} w_{ij}^{(l)} \, \eta_{ij}^{(l)} \left(t - t_{j}^{(l-1)}\right) + C_{i}^{(l)}$$

$$= \frac{1}{\tau_{i}^{(l-1)}} \sum_{j} w_{ij}^{(l)} \sum_{u=0}^{t_{emid}^{(l-1)}-t_{j}^{(l-1)}-1} \eta_{ij}^{(l)}(u) + d_{i}^{(l-1)} \cdot C_{i}^{(l)}$$

$$= \sum_{j} w_{ij}^{(l)} \frac{1}{\tau_{i}^{(l-1)}} \sum_{u=0}^{t_{emid}^{(l-1)}-t_{ref}^{(l-1)}+\tau_{i}^{(l-1)}x_{j}^{(l-1)}-1} \eta_{ij}^{(l)}(u) + d_{i}^{(l-1)} \cdot C_{i}^{(l)}$$

$$= \sum_{j} w_{ij}^{(l)} \sum_{u=0}^{\tau_{i}^{(l-1)}\left(x_{j}^{(l-1)}-a_{i}^{(l-1)}\right)-1} \left(f_{ij}\left(\frac{u+1}{\tau_{i}^{(l-1)}}+a_{i}^{(l-1)}\right) - f_{ij}\left(\frac{u}{\tau_{i}^{(l-1)}}+a_{i}^{(l-1)}\right)\right) + d_{i}^{(l-1)} \cdot C_{i}^{(l)}$$

$$= \sum_{j} w_{ij}^{(l)} \left(f_{ij}(x_{j}^{(l-1)}) - f_{ij}(a_{i}^{(l-1)})\right) + d_{i}^{(l-1)} \cdot C_{i}^{(l)}$$

$$= \sum_{j} w_{ij}^{(l)} f_{ij}(x_{j}^{(l-1)}). \tag{24}$$

where the second equation uses  $u = t - t_j^{(l-1)}$ ; third equation uses  $x_j^{(l-1)} \, \tau_i^{(l-1)} = t_{\textit{ref}}^{(l-1)} - t_j^{(l-1)}$ ; fourth equation uses  $a_i^{(l-1)} = \frac{t_{\textit{ref}}^{(l-1)} - t_{\textit{end}}^{(l-1)}}{\tau_i^{(l-1)}}$ .

If the spike is emitted at time  $t_i^{(l)} \in \{t_{\textit{emit}}^{(l)}, t_{\textit{emit}}^{(l)} + 1, \dots, t_{\textit{end}}^{(l)}\}$ , i.e. the corresponding value  $x_i^{(l)} \in [a_i^{(l)}, b_i^{(l)}]$ . Then:

$$\theta^{(l)} = V(t_{emit}^{(l)}) + \sum_{v=0}^{t_i^{(l)} - t_{emit}^{(l)} - 1} \psi_i^{(l)}(v) - \Delta_i^{(l)}$$

$$= V(t_{emit}^{(l)}) + \sum_{v=0}^{t_{emit}^{(l)} - t_{emit}^{(l)} - \tau_i^{(l)} x_i^{(l)} - 1} \left( h^{-1}(b_i^{(l)} - \frac{v}{\tau_i^{(l)}}) - h^{-1}(b_i^{(l)} - \frac{v + 1}{\tau_i^{(l)}}) \right) - \Delta_i^{(l)}$$

$$= V(t_{emit}^{(l)}) + \sum_{v=0}^{\tau_i^{(l)} (b_i^{(l)} - x_i^{(l)}) - 1} \left( h^{-1}(b_i^{(l)} - \frac{v}{\tau_i^{(l)}}) - h^{-1}(b_i^{(l)} - \frac{v + 1}{\tau_i^{(l)}}) \right) - \Delta_i^{(l)}$$

$$= V(t_{emit}^{(l)}) + h^{-1}(b_i^{(l)}) - h^{-1}(x_i^{(l)}) - \Delta_i^{(l)}. \tag{25}$$

where the first equation uses  $v=t_i^{(l)}-t_{\mathit{emit}}^{(l)}$ , the second equation uses  $x_i^{(l)}\,\tau_i^{(l)}=t_{\mathit{ref}}^{(l)}-t_i^{(l)}$ , the third equation uses  $b_i^{(l)}=\frac{t_{\mathit{ref}}^{(l)}-t_{\mathit{emit}}^{(l)}}{\tau_i^{(l)}}$ .

Thus

$$h^{-1}(x_i^{(l)}) = V(t_{emit}^{(l)}).$$

which indicates that

$$x_i^{(l)} = h(V(t_{emit}^{(l)})).$$

If  $h(V(t_{\mathit{emit}}^{(l)})) > b_i^{(l)}$ , then  $V(t_{\mathit{emit}}^{(l)}) > h^{-1}(b_i^{(l)}) = \theta_i^{(l)}$ , which means that a spike is emitted once at  $t_{\mathit{emit}}^{(l)}$  representing the value  $\frac{t_{\mathit{ref}}^{(l)} - t_{\mathit{emit}}^{(l)}}{\tau_i^{(l)}} = b_i^{(l)}$ .

If  $h(V(t_{emit}^{(l)})) < a_i^{(l)}$ , then the potential at time  $t_{end}^{(l)}$  is:

$$\begin{split} &V(t_{\textit{emit}}^{(l)}) + \sum_{v=0}^{T^{(l)}} \psi_i^{(l)}(v) \\ &= V(t_{\textit{emit}}^{(l)}) + \sum_{v=0}^{T^{(l)}} \left( h^{-1}(b_i^{(l)} - \frac{v}{\tau_i^{(l)}}) - h^{-1}(b_i^{(l)} - \frac{v+1}{\tau_i^{(l)}}) \right) \\ &= V(t_{\textit{emit}}^{(l)}) + h^{-1}(b_i^{(l)}) - h^{-1}(a_i^{(l)}) \\ &< h^{-1}(b_i^{(l)}) = \theta_i^{(l)}. \end{split} \tag{26}$$

which means that there will be no spike, representing the value  $a_i^{(l)}$ .

### **proof 3** Proof of Lemma 2:

 According to Theorem 4.3 in Zhao et al. (2025): in the continuous setting, if the spike is emitted at time  $t_i^{(l)} \in [t_{emit}^{(l)}, t_{end}^{(l)}]$ , i.e. the corresponding value  $y_i^{(l)} \in [a_i^{(l)}, b_i^{(l)}]$ . Then

$$\theta_i^{(l)} = V(t_{emit}^{(l)}) + \int_0^{t_i^{(l)} - t_{emit}^{(l)}} \psi_i^{(l)}(v) \, dv \tag{27}$$

$$=V(t_{emit}^{(l)}) + \int_{0}^{t_{ref}^{(l)} - t_{emit}^{(l)} - \tau_{i}^{(l)} y_{i}^{(l)}} \frac{1}{\tau_{i}^{(l)}} (h^{-1})' \left(b_{i}^{(l)} - \frac{v}{\tau_{i}^{(l)}}\right) dv$$
 (28)

$$=V(t_{emit}^{(l)})-h^{-1}\left(b_{i}^{(l)}-\frac{v}{\tau_{i}^{(l)}}\right)\Big|_{0}^{\tau_{i}^{(l)}(b_{i}^{(l)}-y_{i}^{(l)})}$$
(29)

$$=V(t_{emit}^{(l)})-h^{-1}(y_i^{(l)})+h^{-1}(b_i^{(l)}).$$
(30)

Because  $\theta_i^{(l)} = h^{-1}(b_i^{(l)})$ :

$$h^{-1}(y_i^{(l)}) = V(t_{emit}^{(l)}), (31)$$

which indicates that:

$$y_i^{(l)} = h(V(t_{emit}^{(l)})).$$
 (32)

In the discrete setting, the spike is emitted at time  $t_i^{(l)} \in \{t_{emit}^{(l)}, t_{emit}^{(l)} + 1, \dots, t_{end}^{(l)}\}$ , the corresponding value  $x_i^{(l)} \in [a_i^{(l)}, b_i^{(l)}]$ . Let  $\Delta_i^{(l)} \geq 0$  is a compensation constant, which is actually the difference between the  $\theta_i^{(l)}$  and the membrane potential at the spike time. The following equation satisfies:

$$\theta_i^{(l)} = V(t_{emit}^{(l)}) + \int_0^{t_i^{(l)} - t_{emit}^{(l)}} \psi_i^{(l)}(v) \, dv - \Delta_i^{(l)}$$
(33)

$$=V(t_{\mathit{emit}}^{(l)})+\int_{0}^{t_{\mathit{ref}}^{(l)}-t_{\mathit{emit}}^{(l)}-\tau^{(l)}x_{i}^{(l)}}\frac{1}{\tau_{i}^{(l)}}(h^{-1})'\bigg(b_{i}^{(l)}-\frac{v}{\tau_{i}^{(l)}}\bigg)\,dv-\Delta_{i}^{(l)} \tag{34}$$

$$=V(t_{emit}^{(l)})-h^{-1}\left(b_{i}^{(l)}-\frac{v}{\tau_{i}^{(l)}}\right)\Big|_{0}^{\tau_{i}^{(l)}(b_{i}^{(l)}-x_{i}^{(l)})}-\Delta_{i}^{(l)} \tag{35}$$

$$=V(t_{\textit{emit}}^{(l)})-h^{-1}(x_i^{(l)})+h^{-1}(b_i^{(l)})-\Delta_i^{(l)}. \tag{36}$$

Because  $\theta_i^{(l)} = h^{-1}(b_i^{(l)})$ :

$$h^{-1}(x_i^{(l)}) = V(t_{qmit}^{(l)}) - \Delta_i^{(l)}, \tag{37}$$

which indicates that:

$$x_i^{(l)} = h\Big(V(t_{emit}^{(l)}) - \Delta_i^{(l)}\Big).$$
 (38)

The error of discrete coding in the continuous setting can be expressed as:

$$\|y_i^{(l)} - x_i^{(l)}\| = \|h(V(t_{emit}^{(l)})) - h(V(t_{emit}^{(l)}) - \Delta_i^{(l)})\|$$
(39)

By the mean value theorem, we obtain:

$$\|y_i^{(l)} - x_i^{(l)}\| = \||h'(\xi)| \cdot \Delta_i^{(l)}\|,$$
 (40)

where  $\xi \in \left[V(t_{\textit{emit}}^{(l)}) - \Delta_i^{(l)}, V(t_{\textit{emit}}^{(l)})\right]$ .

Furthermore, we examine  $\Delta_i^{(l)}$  to provide a more in-depth analysis of the error. We assume that the spike firing time corresponding precisely to the ANN output is denoted as  $[t]_i^{(l)}$ . Based on the characteristics of TTFS encoding, it follows that:

$$t_i^{(l)} - 1 \le [t]_i^{(l)} \le t_i^{(l)}. \tag{41}$$

Then  $\Delta_i^{(l)}$  can be represented as:

$$\Delta_{i}^{(l)} = \int_{[t]_{i}^{(l)} - t_{emit}^{(l)}}^{t_{i}^{(l)} - t_{emit}^{(l)}} \psi_{i}(s) \, ds = h^{-1}(t_{i}^{(l)} - t_{emit}^{(l)}) - h^{-1}([t]_{i}^{(l)} - t_{emit}^{(l)}). \tag{42}$$

By the mean value theorem, we obtain:

$$\left\| \Delta_i^{(l)} \right\| = \left\| \left| (h^{-1})'(\hat{t}_i) \right| \cdot (t_i^{(l)} - [t]_i^{(l)}) \right\|, \tag{43}$$

where  $\hat{t}_i \in \left[ [t]_i^{(l)} - t_{emit}^{(l)}, t_i^{(l)} - t_{emit}^{(l)} \right]$ .

Then the error  $\epsilon_i^{(l)}$  can be bounded by the following inequality:

$$\left\| y_i^{(l)} - x_i^{(l)} \right\| \le |h'(\xi)| \cdot \left| (h^{-1})'(\hat{t}_i) \right| \cdot |[t]_i^{(l)} - t_i^{(l)}| \tag{44}$$

By the definition of clock precision:  $\Delta t_{real} = t_{real}(t+1) - t_{real}(t)$ , where  $t_{real}(t) = t \cdot \Delta t_{real}$ , we obtain:

$$||y_{i}^{(l)} - x_{i}^{(l)}|| \leq |h'(\xi)| \cdot |(h^{-1})'(\hat{t}_{i})| \cdot \Delta t_{real}$$

$$= |h'(\xi)| \cdot |(h^{-1})'(\hat{t}_{i})| \cdot \frac{\Omega}{t_{end}^{(l)} - t_{emit}^{(l)}}$$

$$= |h'(\xi)| \cdot |(h^{-1})'(\hat{t}_{i})| \cdot \frac{\Omega}{T^{(l)}}$$

$$\leq \frac{G_{1}G_{2}\Omega}{T^{(l)}}$$
(45)

**proof 4** Proof of Theorem 3:

For clipping error, according to Lemma 1, we can get:

$$\mathcal{E}_{c}^{(l)} = \begin{cases} \left\| \alpha_{i}^{(l)} - b_{i}^{(l)} \right\| & \text{if } \alpha_{i}^{(l)} > b_{i}^{(l)} \\ 0 & \text{if } \alpha_{i}^{(l)} \in [a_{i}^{(l)}, b_{i}^{(l)}] \\ \left\| a_{i}^{(l)} - \alpha_{i}^{(l)} \right\| & \text{if } \alpha_{i}^{(l)} < a_{i}^{(l)}, \end{cases}$$

$$(46)$$

We define the center of the output interval of SNN as:

$$c_i^{(l)} = \frac{a_i^{(l)} + b_i^{(l)}}{2} \tag{47}$$

The clipping error can then be restated as follows:

$$\mathcal{E}_{c}^{(l)}(T^{(l)}) = \max\left(\left|\alpha_{i}^{(l)} - c_{i}^{(l)}\right| - \frac{T^{(l)}}{2\tau_{i}^{(l)}}, 0\right) \tag{48}$$

We take the derivative of  $T^{(l)}$  to get the sensitivity of the error  $\mathcal{E}_c^{(l)}(T^{(l)})$  to  $T^{(l)}$ :

$$\frac{d}{dt}\mathcal{E}_{c}^{(l)}(T^{(l)}) = \begin{cases}
-\frac{1}{2\tau_{i}^{(l)}}, & \left|\alpha_{i}^{(l)} - c_{i}^{(l)}\right| > \frac{T^{(l)}}{2\tau_{i}^{(l)}}, \\
0, & \left|\alpha_{i}^{(l)} - c_{i}^{(l)}\right| < \frac{T^{(l)}}{2\tau_{i}^{(l)}}.
\end{cases} (49)$$

Once clipping occurs, increasing  $T^{(l)}$  will reduce the error linearly with a constant slope of  $-\frac{1}{2\tau_i^{(l)}}$ ; within the valid interval, the error is unaffected by  $T^{(l)}$ .

For quantization error, according to Lemma 2, we can get:

$$\mathcal{E}_q^{(l)} \le \frac{G_1 G_2 \Omega}{T^{(l)}}.\tag{50}$$

For an L-layer network with I neurons in each layer, we can get:

$$\mathcal{E} = \sum_{i=1}^{I} \sum_{l=1}^{L} \left( \mathcal{E}_{c}^{(l)} + \mathcal{E}_{q}^{(l)} \right)$$

$$\leq \sum_{i=1}^{I} \sum_{l=1}^{L} \left( \max \left( \left| \alpha_{i}^{(l)} - c_{i}^{(l)} \right| - \frac{T^{(l)}}{2\tau_{i}^{(l)}} \right|, 0 \right) + \frac{G_{1}G_{2}\Omega}{T^{(l)}} \right)$$
(51)

Let 
$$T = \min \left\{ T^{(l)} \right\}_{l=1}^{L}$$
 and  $\tau = \max \left\{ \left\{ \tau_i^{(l)} \right\}_{i=1}^{I} \right\}_{l=1}^{L}$ :

$$\mathcal{E} \le LI \cdot \max\left( \left| \alpha_i^{(l)} - \frac{a_i^{(l)} + b_i^{(l)}}{2} \right| - \frac{T^{(l)}}{2\tau_i^{(l)}}, 0 \right) + \frac{LIG_1G_2\Omega}{T}$$
 (52)

**proof 5** *Proof of Corollary 1:* 

**Input transform:** The input of QANN at i-th neuron of l-th layer is  $\mathbf{X}_i^{(l)} = \sum_j w_{ij}^{(l)} x_j^{(l-1)} \in [a_i^{(l-1)}, b_i^{(l-1)}]$ . In order to approximate the input of QANN, based on Theorem 2, we set the kernel function  $\eta_{ij}^{(l)}$  and  $C_i^{(l)}$  as follows:

$$\eta_{ij}^{(l)}(u) = \mathcal{H}\left(\frac{u}{\tau_i^{(l-1)}} + a_i^{(l-1)}\right),$$

$$C_i^{(l)} = \sum_j \frac{a_i^{(l-1)}}{d_i^{(l-1)}} w_{ij},$$

Then, the membrane potential after reception is completed can be expressed as:

$$V(T_{emit}) = \sum_{i} w_{ij}^{(l)} x_j^{(l-1)} = \mathbf{X}_i^{(l)}.$$
 (53)

**Output transform:** In order to approximate the output of QANN at l-th layer of i-th neuron:  $\hat{\mathbf{X}}^{(l)} \in [a_i^{(l)}, b_i^{(l)}]$ , based on Theorem 2, we set the kernel function  $\psi_i^{(l)}$  and threshold as follows:

$$\psi_i^{(l)}(v) = \frac{1}{\tau_i^{(l)}}, \quad \theta_i^{(l)} = b_i^{(l)}$$
(54)

If the spike is emitted at time  $t \in \{t_{\textit{emit}}^{(l)}, t_{\textit{emit}}^{(l)} + 1, \dots, t_{\textit{end}}^{(l)}\}$ :

$$X_i^{(l-1)} + \frac{1}{\tau_i^{(l)}} \cdot t \ge \theta_i^{(l)}. \tag{55}$$

According to the definition of t, we can get:

$$t = \left[ (\theta_i^{(l)} - \mathbf{X}_i^{(l)}) \tau_i^{(l)} \right] \tag{56}$$

According to the rounding range of t, we add the clip function to get:

$$t = \operatorname{clip}\left(\left\lceil (\boldsymbol{\theta}_{i}^{(l)} - \mathbf{X}_{i}^{(l)}) \boldsymbol{\tau}_{i}^{(l)}\right\rceil, t_{\mathit{emit}}^{(l)}, t_{\mathit{end}}^{(l)}\right) \tag{57}$$

According to the relation between spike time and corresponding activation value and  $\theta_i^{(l)} = b_i^{(l)}$ :

$$x_i^{(l)} = \frac{1}{\tau_i^{(l)}} \left( t_{ref}^{(l)} - \text{clip}(\left[ (b_i^{(l)} - \mathbf{X}_i^{(l)}) \tau_i^{(l)} \right], t_{emit}^{(l)}, t_{end}^{(l)}) \right)$$
(58)

According to  $b_i^{(l)} = \frac{t_{\it ref}^{(l)} - t_{\it emit}^{(l)}}{ au_i^{(l)}}$ , we can get:

$$x_i^{(l)} = \frac{1}{\tau_i^{(l)}} \left( t_{ref}^{(l)} - \text{clip}(\left[ t_{ref}^{(l)} - t_{emit}^{(l)} - \mathbf{X}_i^{(l)} \tau_i^{(l)} \right], t_{emit}^{(l)}, t_{end}^{(l)}) \right)$$
(59)

Based on the relationship between the ceiling function and the floor function, we can derive the following:

$$x_i^{(l)} = \frac{1}{\tau_i^{(l)}} \operatorname{clip}(\left[\mathbf{X}_i^{(l)} \tau_i^{(l)}\right] - t_{ref}^{(l)} + t_{emit}^{(l)}, -t_{emd}^{(l)}, -t_{emit}^{(l)}) + \frac{1}{\tau_i^{(l)}} t_{ref}^{(l)}$$

$$(60)$$

Based on the properties of the floor function, we can conclude that:

$$x_{i}^{(l)} = \frac{1}{\tau_{i}^{(l)}} \text{clip}(\left[\mathbf{X}_{i}^{(l)} \tau_{i}^{(l)}\right] + t_{\textit{end}}^{(l)} - t_{\textit{ref}}^{(l)} + t_{\textit{emit}}^{(l)}, 0, t_{\textit{end}}^{(l)} - t_{\textit{emit}}^{(l)}) + \frac{1}{\tau_{i}^{(l)}} \left(t_{\textit{ref}}^{(l)} - t_{\textit{end}}^{(l)}\right) \tag{61}$$

Let 
$$t_{emit}^{(l)}=0$$
,  $t_{end}^{(l)}=N$ ,  $au_i^{(l)}=rac{1}{\lambda_i^{(l-1)}}$ ,  $t_{end}^{(l)}-t_{ref}^{(l)}=z^{(l)}$ , we can get:

$$x_i^{(l)} = \lambda_i^{(l)} \cdot \text{clip}(\left[\frac{\mathbf{X}_i^{(l)}}{\lambda_i^{(l)}}\right] + z^{(l)}, 0, N) - \lambda_i^{(l)} z^{(l)} = \hat{\mathbf{X}}^{(l)}$$
(62)

# E NONLINEAR OPERATIONS IN QC-A2S

**Corollary 2** (Construction of GELU) A TTFS-based neuron can be made equivalent to a discrete SiLU function with through the following configuration:

$$\eta_{ij}^{(l)}(u) = \mathbb{I}[u \geq 0] \cdot \tau_i^{(l-1)} \cdot \left( \left( \frac{u+1}{\tau_i^{(l-1)}} + a_i^{(l-1)} \right) \cdot \sigma \left( \frac{u+1}{\tau_i^{(l-1)}} + + a_i^{(l-1)} \right) - \left( \frac{u}{\tau_i^{(l-1)}} + + a_i^{(l-1)} \right) \cdot \sigma \left( \frac{u}{\tau_i^{(l-1)}} + a_i^{(l-1)} \right) \right)$$

$$C_i^{(l)} = \sum_j w_{ij}^{(l)} \frac{a_i^{(l-1)} \cdot \sigma(a_i^{(l-1)})}{d_i^{(l-1)}}, \ \psi_i^{(l)}(v) = \frac{1}{\tau_i^{(l)}}, \ \sigma(x) = \frac{1}{1 + e^{-x}}.$$
 (63)

**Corollary 3 (Construction of GELU)** A TTFS-based neuron can be made equivalent to a discrete GELU function with

$$\eta_{ij}^{(l)}(u) = \mathbb{I}[u \geq 0] \cdot \tau_i^{(l-1)} \cdot \left( \left( \frac{u+1}{\tau_i^{(l-1)}} + a_i^{(l-1)} \right) \cdot \Phi\left( \frac{u+1}{\tau_i^{(l-1)}} + + a_i^{(l-1)} \right) - \left( \frac{u}{\tau_i^{(l-1)}} + + a_i^{(l-1)} \right) \cdot \Phi\left( \frac{u}{\tau_i^{(l-1)}} + a_i^{(l-1)} \right) \right)$$

**Corollary 4 (Construction of Softmax)** The log-sum-exp of I inputs  $x_1, x_2, \dots, x_I$ , i.e.,

$$\log \sum_{j=1}^{I} e^{x_j},\tag{64}$$

can be calculated in a single neuron with

$$\eta_{ij}^{(l)}(u) = \tau_i^{(l-1)} \cdot \left( \exp\left(\frac{u+1}{\tau_i^{(l-1)}} + a_i^{(l-1)}\right) - \exp\left(\frac{u}{\tau_i^{(l-1)}} + a_i^{(l-1)}\right) \right). \tag{65}$$

$$C_i^{(l)} = \frac{I}{d_i^{(l-1)}} e^{a_i^{(l-1)}}, \ \psi_i^{(l)}(v) = \frac{1}{\tau_i^{(l)}} \exp\left(b_i^{(l)} - \frac{v}{\tau_i^{(l)}}\right). \tag{66}$$

With the log-sum-exp neuron, we can obtain the softmax operator. We can calculate the logarithm of softmax, i.e.

$$\log\left(\frac{e^{x_i}}{\sum_{j=1}^{I} e^{x_j}}\right) = x_i - \log\sum_{j=1}^{I} e^{x_j},\tag{67}$$

by subtracting the log-sum-exp from  $x_i$ . Finally, we can obtain the output after an exponential layer.

**Corollary 5 (Construction of RMSNorm)** *RMSNorm is a normalization method widely used in LLaMA architecture, which is a linear operation. RMSNorm is defined as:* 

$$RMSNorm(x_i) = \frac{x_i}{\sqrt{\frac{1}{I} \sum_{i=1}^{I} x_i^2}} \cdot \gamma + \beta.$$
 (68)

We first can obtain the  $\frac{1}{I}\sum_{i=1}^{I}x_{i}^{2}$  by a single neuron with

$$\eta_{ij}^{(1)}(u) = \tau_{ij}^{(0)} \left[ \left( \frac{u+1}{\tau_i^{(0)}} + a_i^{(0)} \right)^2 - \left( \frac{u}{\tau_i^{(0)}} + a_i^{(0)} \right)^2 \right], \ C_i^{(1)} = \frac{(a_i^{(0)})^2}{Id_i^{(0)}}, \ w^{(1)} = \frac{1}{I}$$
 (69)

$$\psi_i^{(1)}(v) = \frac{1}{\tau_i^{(1)}}, \ \theta_i^{(1)} = b_i^{(1)}. \tag{70}$$

Then, we can get  $\frac{1}{\sqrt{\frac{1}{t}\sum_{i=1}^{I}x_i^2}}$  with:

$$\eta_{ij}^{(2)}(u) = \tau_{ij}^{(1)} \left[ \left( \frac{u+1}{\tau_i^{(1)}} + a_i^{(1)} \right)^{-\frac{1}{2}} - \left( \frac{u}{\tau_i^{(1)}} + a_i^{(1)} \right)^{-\frac{1}{2}} \right], \ C_i^{(2)} = \frac{1}{Id_i^{(1)}(a_i^{(1)})^{\frac{1}{2}}}, \tag{71}$$

$$w^{(2)} = 1 \,\psi_i^{(l)}(v) = \frac{1}{\tau_i^{(l)}}, \,\theta_i^{(l)} = b_i^{(l)}. \tag{72}$$

Finally, multiply  $x_i$  with  $\frac{1}{\sqrt{\frac{1}{I}\sum_{i=1}^I x_i^2}}$ .

### F USE OF LLMS

In this work, LLMs are employed solely for polishing or grammar checking text that is originally written by us.