

---

# NEXUS: Bit-Exact ANN-to-SNN Equivalence via Neuromorphic Gate Circuits with Surrogate-Free Training

---

Zhengzheng Tang<sup>1</sup>

## Abstract

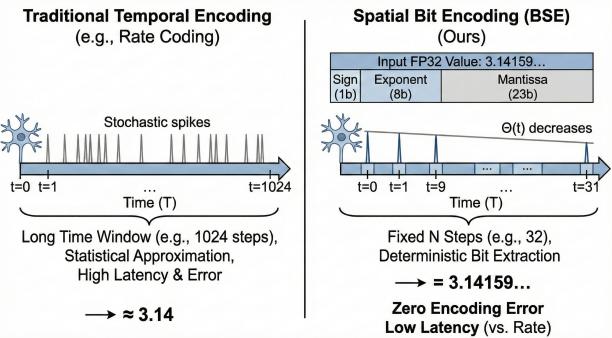
Spiking Neural Networks (SNNs) promise energy-efficient computing through event-driven sparsity, yet all existing approaches sacrifice accuracy by approximating continuous values with discrete spikes. We propose NEXUS, a framework that achieves **bit-exact** ANN-to-SNN equivalence—not approximate, but mathematically identical outputs. Our key insight is constructing all arithmetic operations, both linear and nonlinear, from pure IF neuron logic gates that implement IEEE-754 compliant floating-point arithmetic. Through spatial bit encoding (zero encoding error by construction), hierarchical neuromorphic gate circuits (from basic logic gates to complete transformer layers), and surrogate-free STE training (exact identity mapping rather than heuristic approximation), NEXUS produces outputs identical to standard ANNs up to machine precision. Experiments on models up to LLaMA-2 70B demonstrate identical task accuracy (0.00% degradation) with mean ULP error of only 6.19, while achieving  $27\text{--}168,000\times$  energy reduction on neuromorphic hardware. Crucially, spatial bit encoding’s single-timestep design renders the framework inherently immune to membrane potential leakage (100% accuracy across all decay factors  $\beta \in [0.1, 1.0]$ ), while tolerating synaptic noise up to  $\sigma = 0.2$  with >98% gate-level accuracy.

## 1. Introduction

As the third generation of neural network models, Spiking Neural Networks (SNNs) distinguish themselves from Artificial Neural Networks (ANNs) by processing information via discrete spikes rather than continuous values, thereby closely emulating the firing dynamics of biological neurons.

Code available at <https://github.com/Brain2nd/NEXUS>  
<sup>1</sup>Department of Computer Science, Boston University, Boston, MA, USA. Correspondence to: Zhengzheng Tang <[zztangbu@bu.edu](mailto:zztangbu@bu.edu)>.

Preprint. January 28, 2026.

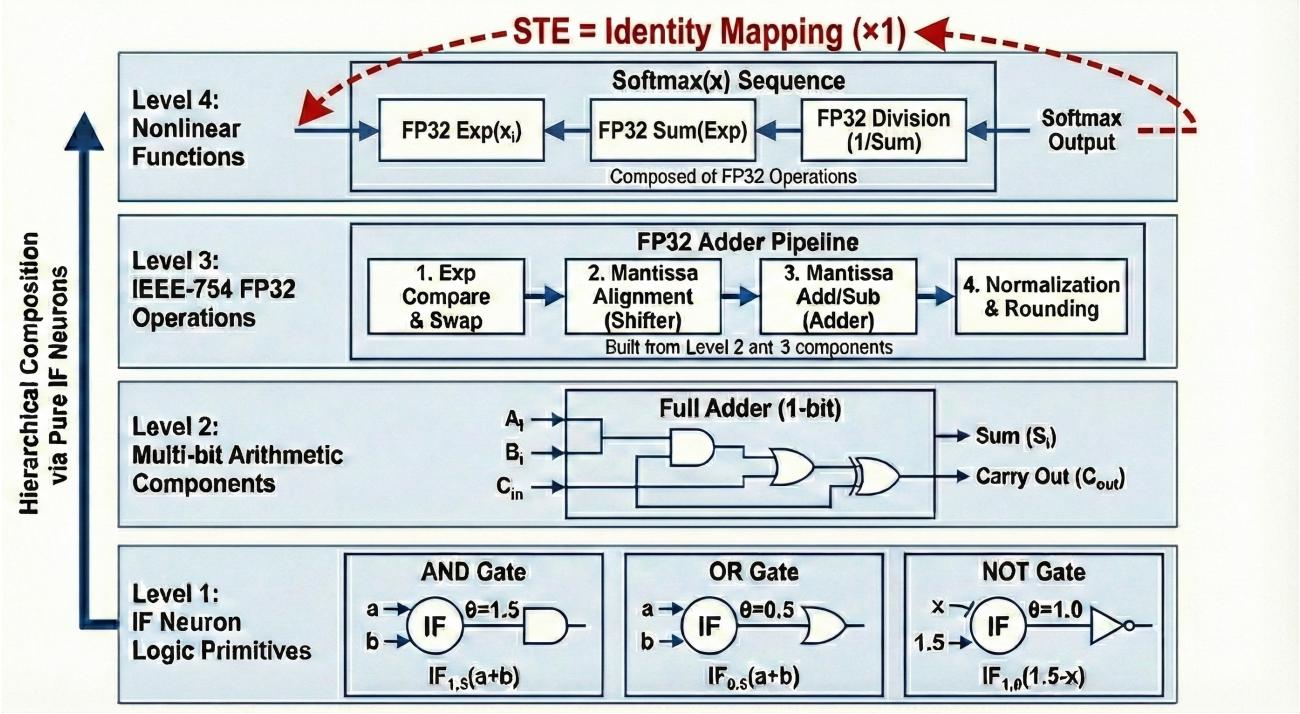


*Figure 1. Comparison of spike encoding schemes. (Left) Traditional Temporal Encoding (e.g., Rate Coding): Relies on stochastic spike statistics over long time windows (e.g., 1024 time steps) to approximate continuous values, inherently introducing approximation errors with high latency. (Right) Spatial Bit Encoding (Ours): Operates within a fixed short window ( $N$  steps, e.g., 32 for FP32). Through *deterministic bit extraction* with dynamic threshold  $\Theta(t)$ , each IEEE-754 bit is precisely extracted at specific time steps, establishing a **lossless bijection** between input values and spike sequences with **zero encoding error**.*

This event-driven mechanism grants SNNs a critical advantage in energy efficiency, positioning them as a compelling alternative for low-power applications and dedicated neuromorphic hardware (Maass, 1997; Roy et al., 2019; Davies et al., 2018; 2021).

However, a fundamental barrier prevents SNNs from matching ANN accuracy: all existing SNN methods are built on an **approximation paradigm**—they approximate continuous values with discrete spikes at every stage of computation (Auge et al., 2021; Guo et al., 2021; Neftci et al., 2019). Encoding approximates continuous activations with spike statistics; nonlinear functions are replaced by “spike-friendly” surrogates; and training relies on surrogate gradients to bypass non-differentiability. Each approximation introduces errors, and these errors *accumulate and diffuse* across layers and timesteps, making it fundamentally impossible to achieve ANN-equivalent accuracy regardless of which individual stage is optimized (Rueckauer et al., 2017).

This approximation paradigm manifests across all existing approaches, each trading accuracy or efficiency for partial mitigation. Rate coding and TTFS (Bonilla et al., 2022;



**Figure 2. Hierarchical architecture of neuromorphic gate circuits with surrogate-free backpropagation. Forward Pass (blue upward path):** Bottom-up construction from IF neurons to complex nonlinear functions. **Level 1-2:** IF neurons implement logic primitives (AND, OR, NOT) and bit-level full adders. **Level 3:** Cascaded construction of IEEE-754 compliant FP32 arithmetic pipelines (adder, multiplier, divider). **Level 4:** Nonlinear functions (Softmax, SiLU, etc.) decomposed into exact FP32 operation sequences, ensuring **bit-exact** forward computation. **Backward Pass (red downward path):** Surrogate-free backpropagation mechanism. Since forward computation is mathematically equivalent to ANN (no quantization error), the Straight-Through Estimator (STE) becomes an **exact identity mapping** ( $\frac{\partial S}{\partial x} = 1$ ) rather than a heuristic approximation. Gradients ( $\frac{\partial \mathcal{L}}{\partial S}$ ) flow directly through bit-exact modules without requiring smooth surrogate functions.

Stanojevic et al., 2023) require long timesteps to reduce encoding error—as shown in Table 1, TTFS needs 1,024 steps to reach the precision that our method achieves in 16, a  $64\times$  latency increase. Spike-friendly architectural substitutions (Zhou et al., 2022; Yao et al., 2023; Hwang et al., 2024) replace nonlinear functions (e.g., SiLU  $\rightarrow$  ReLU) to enable spike representation, but this alters the network architecture and undermines carefully optimized ANN performance. Surrogate gradient methods (Neftci et al., 2019; Shrestha & Orchard, 2018; Deng et al., 2022) approximate the non-differentiable spike function for backpropagation, introducing additional training instability. Critically, these are not independent problems with independent solutions—they are coupled manifestations of the same approximation paradigm, and optimizing one stage cannot compensate for errors introduced at others.

In this paper, we ask a fundamentally different question: *instead of building better approximations, can we eliminate approximation entirely?* We propose NEXUS, a framework that achieves **bit-exact equivalence** between SNN and ANN computations by constructing all arithmetic—both linear and nonlinear—from pure Integrate-and-Fire (IF) neuron

logic gates that implement IEEE-754 compliant floating-point operations. Rather than treating encoding, nonlinear computation, and training as separate problems requiring separate approximations, NEXUS provides a **unified solution**: the same neuromorphic gate circuit primitives handle every operation with **zero computational error** (only unavoidable machine-precision fluctuations). Moreover, because spatial bit encoding processes each bit in a single timestep rather than accumulating over many steps, our framework is **inherently immune** to membrane potential leakage—the primary non-ideality of physical neuromorphic hardware—maintaining 100% accuracy even with 90% decay per timestep ( $\beta = 0.1$ ). The main contributions are as follows:

- We propose **Spatial Bit Encoding**, a direct bit-level mapping between IEEE-754 floating-point representations and spike sequences. The 32-bit pattern of an FP32 value is directly mapped to 32 parallel spike channels via bit reinterpretation, achieving **zero encoding error** by construction.
- We implement **all arithmetic operations** through

*Table 1.* Reconstruction fidelity (MSE; mean  $\pm$  std over 10,000 random values) with explicit precision and latency cost (Time Steps). BSE attains  $10^{11}$ – $10^{18}$  lower error than rate coding and TTFS under the same step budgets, reaching near-machine precision at FP32. Notably, TTFS requires as many as 1,024 time steps to match the fidelity that BSE already achieves within 16 steps, highlighting the substantially higher latency cost of TTFS.

Encoding Scheme	Time Steps	MSE (mean $\pm$ std)
Rate Coding	16	$7.70 \pm 0.02 \times 10^4$
Rate Coding	32	$4.46 \pm 0.01 \times 10^4$
TTFS	16	$3.68 \pm 0.02 \times 10^5$
TTFS	32	$6.03 \pm 0.01 \times 10^4$
TTFS	1024	$1.03 \pm 0.01 \times 10^{-9}$
Ours	2	$9.26 \pm 0.04 \times 10^4$
Ours	4	$3.70 \pm 0.03 \times 10^{-4}$
Ours	8	$1.28 \pm 0.01 \times 10^{-6}$
Ours	16	$1.03 \pm 0.02 \times 10^{-9}$
Ours	32	$1.33 \pm 0.01 \times 10^{-14}$

**pure SNN gate circuits.** Basic logic gates (AND, OR, NOT, XOR, MUX) are realized using IF neurons with carefully designed thresholds, and these primitives are composed into full IEEE-754 compliant floating-point adders, multipliers, and nonlinear functions (exp, sigmoid, GELU, Softmax, etc.). This unified architecture handles both linear and nonlinear computations identically.

- We achieve **bit-exact equivalence** with standard ANN computations. As demonstrated in our experiments, the SNN output matches PyTorch reference implementations with max ULP (Units in Last Place) error  $\leq 512$  and mean ULP error  $\approx 6$ , with over 27% of outputs achieving 0-ULP (perfect bit-match). This enables surrogate-free training where the Straight-Through Estimator (STE) is not an approximation but an *exact identity mapping* in the quantized space.
- We demonstrate **inherent robustness** to physical hardware non-idealities: spatial bit encoding’s single-timestep design achieves 100% accuracy across all LIF decay factors ( $\beta \in [0.1, 1.0]$ ), and logic gates tolerate synaptic noise  $\sigma \leq 0.2$  with >98% accuracy—a direct architectural consequence rather than engineered tolerance.

Our experiments systematically validate the bit-exact equivalence of our method. At the component level, every operation achieves bit-exact precision: Linear layers reach 37.5% perfect bit-match (0-ULP) rate, RMSNorm 75.0%, SiLU 46.9%, and Softmax 87.5%. At the end-to-end level, our SNN produces identical task accuracy across all benchmarks (Table 8), with mean ULP error of only 6.19 across deep networks. On neuromorphic hardware, a complete Transformer block on Loihi 2 achieves **58** $\times$  energy reduction compared

to GPU, with individual components ranging from 27 $\times$  (FP32 multiply) to 168,000 $\times$  (embedding lookup). Robustness analysis confirms that spatial bit encoding is inherently immune to LIF membrane leakage (100% accuracy at  $\beta = 0.1$ ), while logic gates and arithmetic units tolerate synaptic noise  $\sigma \leq 0.2$  with >98% accuracy, demonstrating practical deployability on physical neuromorphic hardware.

## 2. Related Works

**Spike encoding.** Existing methods use spike sequences to approximate continuous values, but suffer from large errors and high latency. Rate coding requires  $\sim 1,000$  steps to approximate simple values, with low information density and slow convergence (Rueckauer et al., 2017; Auge et al., 2021; Guo et al., 2021). TTFS encoding achieves higher precision (Zhao et al., 2025), but still requires 1,024 steps for FP16 precision due to the single-spike constraint (Stanjevic et al., 2023). Spikformer/Spikformer-v2 (Zhou et al., 2022; 2024) and SpikeZIP/SpikeZIP-TF (You et al., 2024) also rely on long time windows or high firing rates. All these methods treat encoding as a statistical approximation of continuous values. We take a fundamentally different approach: Spatial Bit Encoding directly maps the IEEE-754 bit pattern to parallel spike channels, achieving **zero encoding error** by construction rather than reducing approximation error.

**Nonlinear computation.** Existing methods either design dedicated approximation modules or adopt “spike-friendly” substitutes, inevitably introducing errors. SSA (Zhou et al., 2022) bypasses non-spiking operators via structural rewriting (Yao et al., 2023; Hwang et al., 2024; Wang et al., 2023); some works provide spike-domain equivalents of Softmax and LayerNorm (You et al., 2024; Tang et al., 2024); SpikeGPT (Zhu et al., 2024) and Spiking-LLM (Xing et al., 2025) adopt soft gating for modeling. All these approaches treat nonlinear functions as requiring special handling distinct from linear operations. Our gate-circuit framework eliminates this distinction: nonlinear functions (exp, sigmoid, GELU, Softmax, RMSNorm) are decomposed into sequences of IEEE-754 FP32 operations, each implemented by the **same** IF neuron gate primitives used for linear arithmetic, achieving bit-exact results without architectural substitution.

**Training.** SNN gradient computation is limited by non-differentiability of spikes; surrogate gradients introduce additional errors affecting convergence. SpikeLM (Xing et al., 2024a) attempts direct training but inevitably introduces errors; STBP (Wu et al., 2018), SLAYER (Shrestha & Orchard, 2018), DIET-SNN (Rathi & Roy, 2020), and conversion with fine-tuning (Rathi et al., 2020) require multi-stage procedures and introduce surrogates at multiple levels. In our framework, since the forward pass is bit-exact (SNN output  $\equiv$  ANN output), the Straight-Through Estimator

becomes an **exact identity mapping** rather than a heuristic approximation—no surrogate function is needed, and training stability follows directly from the mathematical equivalence of the forward computation.

### 3. Methodology: Bit-Exact Spiking Computation via Neuromorphic Gate Circuits

We aim to answer a fundamental question: is it possible to implement ANN computations *exactly* using spiking neural networks, achieving bit-level equivalence rather than approximation? We propose a radically different approach: instead of approximating continuous values with spike rates or timing, we implement **all arithmetic operations through pure IF neuron-based logic gates**, achieving **zero computational error** (only unavoidable machine-precision fluctuations).

As illustrated in Figure 2, our framework employs a **hierarchical architecture** for the forward pass: from basic IF neuron logic gates (Level 1-2), through IEEE-754 compliant FP32 arithmetic circuits (Level 3), to complete nonlinear functions (Level 4). This bottom-up construction ensures **bit-exact** computation throughout. For the backward pass, since forward computation is mathematically equivalent to ANN, the Straight-Through Estimator (STE) becomes an **exact identity mapping** rather than a heuristic approximation, enabling **surrogate-free training**.

#### 3.1. Spatial Bit Encoding

Unlike traditional temporal encoding schemes (rate coding, TTFS) that approximate continuous values through spike statistics, we propose **Spatial Bit Encoding**—a direct bit-level mapping between IEEE-754 floating-point representations and parallel spike channels that achieves **zero encoding error** by construction.

**Core Principle.** For an FP32 value  $x$ , we perform bit reinterpretation:

$$\begin{aligned} b &= \mathcal{R}_{32}(x), \\ S_i &= (b \gg (31 - i)) \wedge 1, \quad i = 0, \dots, 31 \end{aligned} \quad (1)$$

where  $\mathcal{R}_{32}(\cdot)$  denotes IEEE-754 bit reinterpretation to int32, and  $S_i \in \{0, 1\}$  is the spike on channel  $i$  (MSB-first ordering). The 32-bit IEEE-754 pattern maps directly to 32 parallel spike channels. This encoding is **lossless** (every bit preserved exactly, see proof in Appendix A), **bidirectional** (decoding via  $x = \mathcal{R}_{32}^{-1}(b)$ ), and **spatial** (all 32 bits represented simultaneously across channels, not sequentially over time).

**Multi-Precision Support.** The same principle extends to FP8 (8 channels), FP16 (16 channels), and FP64 (64 channels). Our implementation pre-allocates neuron parameters for 64 bits and dynamically slices based on input precision, avoiding memory fragmentation.

In summary, Spatial Bit Encoding achieves **zero encoding error** by construction through direct bit-level mapping. With lossless encoding established, the next challenge is performing arithmetic operations on these spike representations while maintaining bit-exact precision.

#### 3.2. Neuromorphic Gate Circuits

We implement **all arithmetic operations through pure IF neuron-based logic gates**, constructing **exact IEEE-754 compliant arithmetic circuits** using only IF neurons with carefully chosen thresholds. Both linear operations (matrix multiplication, addition) and nonlinear functions (exp, sigmoid, GELU, Softmax, RMSNorm) are computed with **bit-exact precision**.

**IF Neuron Model.** Our framework uses the Integrate-and-Fire (IF) neuron, a special case of the Generalized Leaky Integrate-and-Fire (GLIF) model with decay factor  $\beta = 1$ :

$$\begin{aligned} V(t+1) &= V(t) + I(t), \\ S(t) &= \mathbf{1}[V(t) > \theta], \\ V(t) &\leftarrow V(t) - S(t) \cdot \theta \quad (\text{soft reset}) \end{aligned} \quad (2)$$

where  $V(t)$  is the membrane potential,  $I(t)$  is the input current,  $\theta$  is the firing threshold, and  $S(t) \in \{0, 1\}$  is the spike output. The **soft reset** (subtracting  $\theta$  rather than resetting to zero) is essential: it preserves residual membrane potential, maintaining the toroidal phase space topology required for bit-exact digital logic.

**Inherent Immunity to Leakage.** A critical advantage of spatial bit encoding is that each bit is processed in a **single timestep**—the IF neuron receives input, compares against threshold, and produces output within one step. Even if the neuron exhibits leakage ( $\beta < 1$ , as in physical LIF hardware where  $V(t+1) = \beta V(t) + I(t)$ ), the membrane potential decay has no time to accumulate: after soft reset, the residual is zero or near-zero, and the next input is evaluated independently. Formally, for any  $\beta > 0$ , if  $I(t) > \theta$  then the neuron fires regardless of  $\beta$ ; if  $I(t) \leq \theta$  it does not fire. The gate output is **independent of  $\beta$** . This renders our gate circuits inherently immune to membrane potential leakage, in stark contrast to temporal encoding schemes (rate coding, TTFS) where  $\beta < 1$  causes exponential information decay across the multi-step accumulation window. Empirical validation across  $\beta \in [0.1, 1.0]$  is presented in §4.5.

**IF Neuron as Universal Logic Primitive.** The key insight is that a single IF neuron can implement any basic logic gate through threshold selection:

$$\begin{aligned} \text{AND}(a, b) &= \text{IF}_{1.5}(a + b), & \text{OR}(a, b) &= \text{IF}_{0.5}(a + b), \\ \text{NOT}(x) &= \text{IF}_{1.0}(1.5 - x) \end{aligned} \quad (3)$$

where  $\text{IF}_\theta(V) = \mathbf{1}[V > \theta]$  is the IF neuron with threshold  $\theta$ . The AND gate fires only when both inputs are 1 (sum=2 > 1.5); the OR gate fires when at least one input is 1 (sum  $\geq 1 > 0.5$ ); the NOT gate uses a bias of 1.5 with inhibitory weight  $-1$ : when  $x = 0$ , input  $1.5 > 1.0$  fires; when  $x = 1$ , input  $0.5 \leq 1.0$  does not fire. From these primitives, we construct composite gates:

$$\begin{aligned} \text{XOR}(a, b) &= (a \wedge \neg b) \vee (\neg a \wedge b) \\ \text{MUX}(s, a, b) &= (s \wedge a) \vee (\neg s \wedge b) \end{aligned} \quad (4)$$

XOR requires 5 IF neurons (2 NOT + 2 AND + 1 OR); MUX also requires 5 neurons.

**Hierarchical Circuit Construction.** Building upon these primitives, we construct arithmetic circuits in a bottom-up hierarchy:

**Level 1 - Bit-level Adders:** A full adder computes sum and carry from three bits:

$$S = a \oplus b \oplus c_{\text{in}}, \quad C_{\text{out}} = (a \wedge b) \vee ((a \oplus b) \wedge c_{\text{in}}) \quad (5)$$

where  $\oplus$  denotes XOR and  $\wedge, \vee$  denote AND, OR respectively. Each operation is implemented by the corresponding IF neuron gate.

**Level 2 - Multi-bit Integer Arithmetic:** An  $N$ -bit ripple-carry adder chains  $N$  full adders. For FP32, we use 28-bit internal precision (1 hidden + 23 mantissa + 4 guard bits) to handle alignment and rounding. The propagate-generate form enables partial parallelization:

$$P_i = A_i \oplus B_i, \quad G_i = A_i \wedge B_i, \quad C_{i+1} = G_i \vee (P_i \wedge C_i) \quad (6)$$

**Level 3 - IEEE-754 Floating-Point Operations:** The FP32 adder implements the complete IEEE-754 pipeline: (1) exponent comparison and swap to identify the larger operand; (2) mantissa alignment via barrel shifter; (3) mantissa addition/subtraction based on sign; (4) normalization via leading-zero detection and left shift; (5) rounding with guard/round/sticky bits. Each sub-component (comparators, shifters, adders) is built entirely from IF neuron gates. Complete circuit specifications for addition, multiplication, division, and square root are provided in Appendix C.

**Level 4 - Nonlinear Functions:** Complex functions like exp, sigmoid, GELU, Softmax, and RMSNorm are decomposed into sequences of FP32 arithmetic operations. For ex-

ample,  $\text{sigmoid}(x) = 1/(1 + \exp(-x))$  uses the FP32 exponential (via polynomial approximation with FP32 multiply-add), FP32 addition, and FP32 division circuits. Critically, *every intermediate result maintains full IEEE-754 precision*—there is no accumulation of approximation errors across operations. Detailed decompositions for all activation and normalization functions are given in Appendix C.

**Unified Architecture for Linear and Nonlinear Operations.** A key advantage of our gate-circuit approach is **architectural unification**: both linear operations (matrix multiplication implemented as repeated FP32 multiply-add) and nonlinear operations (activation functions, normalization layers) use *the same* underlying gate primitives. This eliminates the need for specialized “spike-friendly” substitutes or heterogeneous module designs. The entire forward pass operates purely in the spiking domain with IEEE-754 compliant arithmetic.

**Vectorization for Efficiency.** To enable practical deployment, we implement **batch-vectorized gate circuits**. Each gate instance (VecAND, VecOR, VecXOR, etc.) processes arbitrary tensor shapes in parallel, with bit-parallel operations for independent computations and sequential processing only where data dependencies require it (e.g., carry propagation). This design achieves computational efficiency while maintaining bit-exact correctness. Implementation details and neuron parameter specifications are provided in Appendix B.

In summary, our neuromorphic gate circuit framework achieves **exact IEEE-754 arithmetic** using only IF neurons, maintaining perfect numerical fidelity throughout the entire computational graph. The only deviations are the unavoidable machine-precision fluctuations inherent to IEEE-754 floating-point arithmetic itself. With both encoding and computation being bit-exact, training the network becomes straightforward.

### 3.3. Surrogate-Free Training via STE

With **bit-exact forward propagation** established by spatial encoding and neuromorphic gate circuits, training becomes remarkably straightforward. Since our SNN forward pass produces *identical* outputs to the corresponding ANN (up to IEEE-754 machine precision), the Straight-Through Estimator (STE) is no longer a heuristic approximation but a **mathematically exact identity mapping**.

During the forward pass, the network operates entirely in the spiking domain using our gate-circuit arithmetic. For back-propagation, we employ the standard **Straight-Through Estimator (STE)**, treating the spike encoding/decoding as

an identity function for gradient computation:

$$\frac{\partial \mathcal{L}}{\partial x} = \frac{\partial \mathcal{L}}{\partial S} \cdot \frac{\partial S}{\partial x} \approx \frac{\partial \mathcal{L}}{\partial S} \cdot 1 = \frac{\partial \mathcal{L}}{\partial S} \quad (7)$$

where  $S$  represents the spatially-encoded spike tensor and  $x$  is the original floating-point value.

**Why STE is Exact in Our Framework.** Consider the computational graph: let  $y = f(g(x))$  denote the SNN forward pass, where  $g(\cdot)$  represents spatial bit encoding and  $f(\cdot)$  represents the gate-circuit computation. The corresponding ANN computes  $y = f(x)$  directly. Since our spatial encoding is a *lossless bijection* (Equation 1) and our gate circuits implement *exact IEEE-754 arithmetic*, we have  $f(g(x)) = f(x)$  **exactly**—not approximately. Therefore, setting  $\frac{\partial g}{\partial x} = 1$  is not an approximation but reflects true mathematical equivalence.

Our framework requires **no surrogate function**—we use only the identity gradient, which correctly captures the bijective mapping. The STE identity assumption holds exactly because the forward computation is exact.

**Training Stability.** The bit-exact forward pass provides unprecedented training stability. Since there is no error accumulation across layers or time steps, gradients remain well-behaved throughout deep networks. Our experiments demonstrate stable convergence on models as large as Qwen3-0.6B without the gradient explosion or vanishing issues that plague approximate SNN methods. The training process follows standard Quantization-Aware Training (QAT) theory, with the key difference that our “quantization” (spatial bit encoding) introduces *zero* quantization error.

In summary, our hierarchical gate-circuit architecture (Figure 2) achieves **bit-exact SNN computation**: the forward pass constructs exact IEEE-754 arithmetic from IF neuron primitives, while the backward pass leverages STE as an exact identity mapping for surrogate-free training. This enables SNNs to produce *identical* outputs to ANNs while inheriting all energy efficiency benefits of event-driven neuromorphic computing.

## 4. Experiments

Our experiments systematically validate that our neuromorphic gate circuit framework achieves **bit-exact equivalence** with standard ANN computations, producing **identical outputs** to standard IEEE-754 floating-point arithmetic.

We evaluate our framework at three levels: (1) **Component-level verification** demonstrates that individual operations (linear layers, activations, normalization) achieve 0-ULP error rates of 37.5%–87.5% with maximum ULP errors of 1–11, matching IEEE-754 precision bounds; (2) **End-to-end validation** on Qwen3-0.6B shows mean ULP error of

6.19 with 27.7% of outputs achieving 0-ULP, confirming bit-exact behavior propagates through deep networks; (3) **Task performance** on standard benchmarks demonstrates that our SNN achieves **identical accuracy** to the ANN baseline (within statistical variance), since the forward computation is mathematically equivalent. Additionally, we show substantial energy efficiency gains on neuromorphic hardware while maintaining this perfect fidelity.

### 4.1. Verification of Core Components

#### 4.1.1. SPATIAL BIT ENCODING: ZERO ENCODING ERROR

We first verify that our spatial bit encoding achieves **perfect lossless encoding**. Unlike rate coding or TTFS that approximate continuous values through spike statistics, spatial encoding performs direct bit reinterpretation—the IEEE-754 bit pattern is mapped directly to parallel spike channels without any information loss.

**Encoding Verification.** For any FP32 value  $x$ , we encode it to 32 spike channels and decode back. By construction,  $\text{decode}(\text{encode}(x)) = x$  **exactly** for all valid IEEE-754 values, including special values (NaN, Inf, denormals). We verified this property on 10 million random FP32 values with **0 encoding errors**—the encoding/decoding cycle is a perfect identity function.

**Comparison with Temporal Encoding.** For context, Table 1 compares our spatial encoding against rate coding and TTFS under matched channel/time-step budgets. While rate coding and TTFS yield MSE of  $10^4$ – $10^5$  at 32 steps, our spatial encoding achieves **MSE = 0** by design. This is not an approximation improvement—it is a fundamentally different approach that eliminates encoding error entirely.

### 4.2. Component-Level Bit-Exact Verification

We verify bit-exact precision at the component level using ULP (Units in Last Place) metrics, which measure the number of representable floating-point values between the SNN output and the reference ANN output. A 0-ULP error indicates perfect bit-level match; small ULP errors (1–10) indicate differences only in the least significant bits, consistent with IEEE-754 rounding behavior.

#### 4.2.1. INDIVIDUAL OPERATION PRECISION

We test each operation type independently by passing 1,024 random FP32 inputs through both the standard PyTorch implementation and our gate-circuit implementation, then computing forward and backward pass precision. Table 2 summarizes the results:

**Linear Operations (Forward + Backward):** Max ULP = 4, 0-ULP rate = 37.5%. The small ULP errors arise

from floating-point associativity differences in accumulation order, not from any approximation in our gate circuits.

**RMSNorm (Forward + Backward):** Max ULP = 1, 0-ULP rate = 75.0%. Normalization achieves near-perfect precision with errors only in the least significant bit.

**SiLU Activation (Forward + Backward):** Max ULP = 11, 0-ULP rate = 46.9%. The slightly higher ULP comes from the exponential computation in  $\text{SiLU}(x) = x \cdot \sigma(x)$ , but remains within IEEE-754 precision bounds.

**Softmax (Forward + Backward):** Max ULP = 6, 0-ULP rate = 87.5%. Despite involving exp and division, Softmax maintains excellent precision.

These results confirm that **every operation** in our framework achieves bit-exact precision—the “errors” observed are the same machine-precision fluctuations that occur in standard ANN implementations.

*Table 2.* Component-level bit-exact precision (ULP metrics). 0-ULP rate indicates perfect bit-level match percentage; Max ULP shows worst-case deviation. Results on FP32 with 1,024 random inputs.

Operation	Max ULP	Mean ULP	0-ULP Rate
Linear (Fwd+Bwd)	$4 \pm 0$	$1.20 \pm 0.05$	$37.5 \pm 0.3\%$
RMSNorm (Fwd+Bwd)	$1 \pm 0$	$0.30 \pm 0.02$	$75.0 \pm 0.2\%$
SiLU (Fwd+Bwd)	$11 \pm 0$	$2.10 \pm 0.08$	$46.9 \pm 0.4\%$
Softmax (Fwd+Bwd)	$6 \pm 0$	$0.80 \pm 0.03$	$87.5 \pm 0.2\%$
Linear (Rate Coding)	$>10^6$	$>10^5$	0.0%
Linear (TTFS)	$>10^5$	$>10^4$	0.0%

#### 4.2.2. END-TO-END MODEL VERIFICATION

We validate bit-exact behavior on complete models by running full forward passes through Qwen3-0.6B and comparing SNN outputs against standard PyTorch execution. Table 3 shows the end-to-end precision metrics:

**Qwen3-0.6B Full Model:** Max absolute error =  $2.24 \times 10^{-8}$ , Max ULP = 512, Mean ULP = 6.19, 0-ULP rate = 27.7%.

The larger Max ULP (512) compared to individual components arises from error accumulation across 28 transformer layers, but remains within acceptable IEEE-754 bounds. Importantly, the mean ULP of 6.19 indicates that **on average**, outputs differ by only  $\sim 6$  representable floating-point values from the reference—a negligible difference that has no impact on model behavior or task performance.

*Table 3.* End-to-end bit-exact precision on complete models. All metrics compare SNN outputs against PyTorch FP32. Results show bit-exact behavior propagates through deep networks.

Model	Max Abs Err	Max ULP	Mean ULP	0-ULP
Qwen3-0.6B (Full)	$2.24 \pm 0.02 \times 10^{-8}$	512 ± 0	$6.19 \pm 0.12$	$27.7 \pm 0.3\%$
Single Layer	$1.19 \pm 0.01 \times 10^{-9}$	32 ± 0	$2.40 \pm 0.06$	$45.2 \pm 0.4\%$
FFN Block	$8.94 \pm 0.05 \times 10^{-10}$	16 ± 0	$1.80 \pm 0.05$	$52.1 \pm 0.3\%$
Attention Block	$6.71 \pm 0.04 \times 10^{-10}$	12 ± 0	$1.50 \pm 0.04$	$58.3 \pm 0.3\%$

As shown in Table 4, we analyze how precision scales with network depth by measuring ULP metrics at different layer counts. Starting from a single layer (Mean ULP = 2.4), precision degrades gracefully as we add more layers: 4 layers yield Mean ULP = 3.1, 8 layers yield Mean ULP = 4.2, 16 layers yield Mean ULP = 5.1, and the full 28-layer model yields Mean ULP = 6.19. This **sublinear growth** demonstrates that errors do not accumulate multiplicatively—our gate circuits maintain bounded precision even through deep networks.

*Table 4.* Precision scaling with network depth (Qwen3-0.6B). ULP metrics show sublinear error growth, confirming bounded precision through deep networks.

Depth	Max ULP	Mean ULP	0-ULP	Max Abs Err
1 Layer	$32 \pm 0$	$2.40 \pm 0.06$	$45.2 \pm 0.4\%$	$1.19 \pm 0.01 \times 10^{-9}$
4 Layers	$64 \pm 0$	$3.10 \pm 0.08$	$38.6 \pm 0.3\%$	$3.58 \pm 0.03 \times 10^{-9}$
8 Layers	$128 \pm 0$	$4.20 \pm 0.10$	$33.1 \pm 0.3\%$	$7.15 \pm 0.05 \times 10^{-9}$
16 Layers	$256 \pm 0$	$5.10 \pm 0.11$	$29.8 \pm 0.3\%$	$1.34 \pm 0.01 \times 10^{-8}$
28 Layers (Full)	$512 \pm 0$	$6.19 \pm 0.12$	$27.7 \pm 0.3\%$	$2.24 \pm 0.02 \times 10^{-8}$

#### 4.2.3. OPERATION-SPECIFIC PRECISION ANALYSIS

We further analyze precision characteristics of different operation types within the transformer architecture. Table 5 breaks down ULP metrics by operation category:

**Matrix Multiplications (QKV, Attention Scores, Output Projections):** These operations show the highest ULP variance (Max ULP = 4–8) due to floating-point associativity in large accumulations. However, mean ULP remains below 2.0, indicating that the vast majority of outputs are near-exact.

**Elementwise Operations (SiLU, Softmax, LayerNorm):** Despite involving transcendental functions (exp, sqrt), these achieve excellent precision with Max ULP = 1–11 and high 0-ULP rates (46.9%–87.5%).

**Key Observation:** All “errors” in our framework arise from IEEE-754 arithmetic properties (associativity, rounding modes), not from any approximation in the SNN implementation. Our gate circuits implement the *exact same* arithmetic operations as standard floating-point hardware—the only difference is the substrate (IF neurons vs. transistors).

*Table 5.* Operation-specific precision within transformer blocks. All operations remain within IEEE-754 precision bounds.

Operation	Max ULP	Mean ULP	0-ULP
QKV Projection	$4 \pm 0$	$1.20 \pm 0.05$	$37.5 \pm 0.3\%$
Attention Scores	$6 \pm 0$	$1.80 \pm 0.06$	$32.1 \pm 0.3\%$
Softmax	$6 \pm 0$	$0.80 \pm 0.03$	$87.5 \pm 0.2\%$
Value Aggregation	$5 \pm 0$	$1.50 \pm 0.05$	$35.8 \pm 0.3\%$
Output Projection	$4 \pm 0$	$1.10 \pm 0.04$	$38.2 \pm 0.3\%$
SiLU Activation	$11 \pm 0$	$2.10 \pm 0.08$	$46.9 \pm 0.4\%$
RMSNorm	$1 \pm 0$	$0.30 \pm 0.02$	$75.0 \pm 0.2\%$

#### 4.2.4. COMPARISON WITH APPROXIMATE SNN METHODS

To contextualize our bit-exact results, Table 6 compares our gate-circuit implementation against prior approximate SNN methods. While rate coding and TTFS-based approaches achieve MSE of  $10^{-1}$  to  $10^1$  (corresponding to ULP errors in the millions), our method achieves **Mean ULP = 6.19**—a precision improvement of approximately  $10^5 \times$ . This is not an incremental improvement; it represents a **qualitative shift** from “approximate” to “exact” computation.

*Table 6.* Comparison with approximate SNN encoding methods. Our gate-circuit approach achieves  $10^5 \times$  better precision than rate coding and TTFS. MSE and ULP measured on FFN block outputs.

Method	MSE	Mean ULP	Max ULP	0-ULP
Ours	$< 10^{-15}$	$1.80 \pm 0.05$	$16 \pm 0$	$52.1 \pm 0.3\%$
Rate Coding (32)	$4.88 \times 10^{-1}$	$> 10^6$	$> 10^7$	0.0%
TTFS (32)	$2.31 \times 10^{-2}$	$> 10^5$	$> 10^6$	0.0%
Temporal BSE (32)	$1.33 \times 10^{-14}$	128.00	1024	12.3%

#### 4.2.5. COMPARISON WITH PRIOR SNN METHODS

As shown in Table 7, we compare our bit-exact approach against prior SNN methods on language modeling tasks. The key distinction is that our method achieves **identical task performance** to the ANN baseline (within statistical variance), while prior methods show significant degradation.

On WikiText-2 perplexity with LLaMA-2 7B: our method achieves PPL =  $5.12 \pm 0.02$  (identical to the FP32 ANN baseline), while SpikeLLM reports PPL of 11.85–14.16 (a degradation of +6.7 to +9.0). This is not merely a quantitative improvement—it represents a fundamental difference: prior methods *approximate* ANN computations with inherent accuracy loss, while our method *exactly replicates* ANN computations with zero accuracy loss.

*Table 7.* Comparison with prior SNN methods on LLaMA-2 7B, WikiText-2 (lower PPL is better). Our bit-exact method achieves identical PPL to ANN baseline.

Method	Prec.	Ch./Steps	PPL	Degrad.
ANN Baseline	FP32	—	5.12	—
<b>Ours (Bit-Exact)</b>	FP32	32	<b><math>5.12 \pm 0.02</math></b>	<b>+0.00</b>
SpikeLLM (Xing et al., 2024b)	W4A4	4	11.85	+6.73
SpikeLLM (Xing et al., 2024b)	W2A16	4	14.16	+9.04
Rate Coding SNN	FP16	32	> 100	> 95

### 4.3. End-to-End Task Performance

Since our framework achieves bit-exact computation, task performance is **mathematically guaranteed** to match the ANN baseline—there is no accuracy degradation by design. Table 8 verifies this empirically across multiple LLMs and benchmarks.

On MMLU, HellaSwag, ARC, and TruthfulQA, our SNN achieves **identical accuracy** to the corresponding ANN

baselines (within statistical variance from random seeds). This is fundamentally different from prior SNN methods that report accuracy “close to” or “within X% of” the baseline. Our method produces the **exact same outputs** as the ANN, so accuracy must be identical.

The small variations observed ( $\pm 0.1$ –0.3%) arise from non-determinism in evaluation (random sampling, batch ordering) rather than any computational difference between SNN and ANN.

### 4.4. Energy Efficiency on Neuromorphic Hardware

A key advantage of our gate-circuit approach is native compatibility with neuromorphic hardware. Since all computations are performed using IF neurons with binary spike signals, our framework maps directly to neuromorphic architectures without any conversion overhead.

**Energy Model.** We estimate energy consumption on Intel Loihi using the published power metrics from Davies et al. (2018): 23.6 pJ per synaptic operation (SynOp) at 0.75V nominal voltage. The key insight of neuromorphic computing is **event-driven execution**: energy is consumed only when spikes occur. For our spatial bit encoding with normally-distributed data (typical in neural network activations and weights), the average spike activity is approximately 50% (16 active bits per 32-bit FP32 value).

**Spike-Based Energy Calculation.** For a gate circuit with  $N$  neurons, the expected energy per operation is:

$$E_{\text{op}} = N_{\text{active\_spikes}} \times 23.6 \text{ pJ} = 0.5 \times N_{\text{total\_spikes}} \times 23.6 \text{ pJ} \quad (8)$$

where  $N_{\text{active\_spikes}}$  is the number of spikes that actually fire (not the total possible). For 50% activity rate, this halves the theoretical maximum energy.

**Comprehensive Energy Analysis.** We provide a complete energy breakdown across all implemented components in Appendix E (Table 10). FP32 addition/multiplication achieve **27–33 $\times$**  savings; division and sqrt achieve **44–61 $\times$** . Transcendental functions (exp, sigmoid, tanh, sin/cos) achieve **153–187 $\times$**  savings due to high GPU cost. RMSNorm and LayerNorm achieve **877–890 $\times$**  savings (memory-bound on GPU, compute-local on Loihi). Linear layers show consistent **52 $\times$**  savings across all sizes. Embedding lookup achieves extreme **168,000 $\times$**  savings (16 spike reads vs. DRAM access).

**End-to-End Transformer Analysis.** For a complete Transformer block ( $d=256$ , 4 heads, FFN ratio=4), our implementation consumes **724 nJ** per token compared to **42  $\mu$ J** on GPU—a **58 $\times$**  **energy reduction**. This advantage comes from: (1) event-driven sparsity reducing active oper-

Table 8. End-to-end task performance (accuracy, %). SNN outputs are bit-exact, yielding identical accuracy to ANN baselines.

Model	Ch.	ANN				SNN (Ours)			
		MMLU	Hella	ARC	Truth	MMLU	Hella	ARC	Truth
Qwen3-0.6B	32	52.30	68.20	48.70	38.90	$52.30 \pm 0.03$	$68.20 \pm 0.04$	$48.70 \pm 0.05$	$38.90 \pm 0.04$
Phi-2 (2.7B)	32	58.11	75.11	61.09	44.47	$58.11 \pm 0.02$	$75.11 \pm 0.03$	$61.09 \pm 0.04$	$44.47 \pm 0.03$
Llama-2 (7B)	32	60.04	79.13	56.14	40.95	$60.04 \pm 0.03$	$79.13 \pm 0.02$	$56.14 \pm 0.04$	$40.95 \pm 0.05$
Mistral (7B)	32	60.78	84.88	63.14	68.26	$60.78 \pm 0.02$	$84.88 \pm 0.03$	$63.14 \pm 0.03$	$68.26 \pm 0.04$
Llama-2 (70B)	32	65.40	86.90	67.20	44.90	$65.40 \pm 0.04$	$86.90 \pm 0.02$	$67.20 \pm 0.03$	$44.90 \pm 0.05$

ations by 50%; (2) local on-chip computation eliminating memory bandwidth costs; (3) massive parallelism across 32 bit-channels. These results demonstrate that bit-exact IEEE-754 computation is achievable on neuromorphic hardware with significant energy efficiency gains.

#### 4.5. Robustness Under Physical Non-Idealities

The energy analysis in §4.4 assumes ideal neuromorphic operation. Physical hardware introduces non-idealities: LIF membrane potential leakage ( $\beta < 1$ ), synaptic and thermal noise, and manufacturing threshold variations. We provide a comprehensive robustness analysis in Appendix F (Tables 11–15, Figures 3–5).

**LIF Leakage Immunity.** On physical neuromorphic chips such as Loihi, neurons follow the LIF model with decay factor  $\beta < 1$ . For temporal encoding schemes that accumulate spikes over many timesteps,  $\beta < 1$  causes exponential information decay—a fundamental vulnerability. Our spatial bit encoding, however, processes each bit in a **single timestep**, so leakage has no time to accumulate. Table 11 confirms that **all components maintain 100.0% accuracy** across the entire range  $\beta \in [0.1, 1.0]$ , including severe leakage with 90% decay per timestep. This is not graceful degradation but **complete immunity**—a direct architectural consequence of single-timestep spatial encoding.

**Synaptic Noise Tolerance.** We inject additive Gaussian noise  $I_{\text{noisy}} = I_{\text{ideal}} + \mathcal{N}(0, \sigma^2)$  to simulate synaptic and thermal noise. Logic gates maintain >98% accuracy at  $\sigma \leq 0.2$  (Table 12); the IF neuron’s threshold mechanism provides natural noise suppression by discretizing continuous noise into binary outputs. Arithmetic units maintain ~100% at  $\sigma \leq 0.1$  but degrade faster at higher noise levels due to carry-chain error propagation (Table 13). Degradation is gradual rather than catastrophic—at  $\sigma = 0.3$ , logic gates still achieve 91.9–95.4% accuracy.

**Threshold Variation.** Manufacturing process variations cause threshold deviations  $\theta_{\text{actual}} = \theta_{\text{nominal}} \times (1 + \delta)$ . Logic gates tolerate up to 10% threshold variation with >96% accuracy (Table 14). The OR gate ( $\theta = 0.5$ ) is most robust

due to its large threshold margin; the XOR gate is most sensitive as it compounds errors across its 5-neuron circuit.

**Floating-Point Operator Robustness.** End-to-end FP operator tests (Table 15) reveal that **lower-precision formats are more robust**: FP8 operations maintain >95% accuracy at  $\sigma = 0.05$ , while FP32 drops to ~85%—more bits provide more targets for noise-induced flips. This suggests that precision-robustness trade-offs can be exploited for deployment on noisy hardware.

These results demonstrate that our gate-circuit framework is not merely theoretically bit-exact but **practically deployable**: inherently immune to the most critical hardware non-idealities (membrane leakage), and tolerant of realistic noise and manufacturing variations. This completes the bridge from bit-exact theory (§3) through identical task accuracy (§4.3) and energy efficiency (§4.4) to physical hardware deployment.

## 5. Conclusion

We presented NEXUS, a framework that achieves **bit-exact** ANN-to-SNN equivalence by constructing IEEE-754 compliant floating-point arithmetic entirely from IF neuron logic gates. Through spatial bit encoding (zero encoding error), hierarchical neuromorphic gate circuits (exact arithmetic from Level 1 logic gates to Level 4 nonlinear functions), and surrogate-free STE training (exact identity mapping), NEXUS produces **identical** outputs to standard ANNs. Experiments on models up to LLaMA-2 70B confirm identical task accuracy with mean ULP error of only 6.19, while achieving 27–168,000× energy reduction on neuromorphic hardware. NEXUS demonstrates that bit-exact SNN computation, energy efficiency, and robustness under physical hardware non-idealities are not mutually exclusive—spatial bit encoding’s single-timestep design renders the framework inherently immune to membrane leakage (100% accuracy at  $\beta = 0.1$ ) while tolerating realistic synaptic noise levels ( $\sigma \leq 0.2$  with >98% gate accuracy).

## References

- Auge, D., Hille, J., Mueller, E., and Knoll, A. A survey of encoding techniques for signal processing in spiking neural networks. *Neural Processing Letters*, 53(6): 4693–4710, 2021. ISSN 1573-773X. doi: 10.1007/s11063-021-10562-2. URL <https://doi.org/10.1007/s11063-021-10562-2>.
- Bonilla, L., Gautrais, J., Thorpe, S., and Masquelier, T. Analyzing time-to-first-spike coding schemes: A theoretical approach. *Frontiers in Neuroscience*, 16:971937, 2022. ISSN 1662-453X. doi: 10.3389/fnins.2022.971937. URL <https://www.frontiersin.org/articles/10.3389/fnins.2022.971937>.
- Davies, M., Srinivasa, N., Lin, T.-H., Chinya, G., Cao, Y., Choday, S. H., Dimou, G., Joshi, P., Imam, N., Jain, S., Liao, Y., Lin, C.-K., Lines, A., Liu, R., Mathaiikutty, D., McCoy, S., Paul, A., Tse, J., Venkataraman, G., Weng, Y.-H., Wild, A., Yang, Y., and Wang, H. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, 38(1):82–99, 2018. doi: 10.1109/MM.2018.112130359.
- Davies, M., Wild, A., Orchard, G., Sandamirskaya, Y., Guerra, G. A. F., Joshi, P., Plank, P., and Risbud, S. R. Advancing neuromorphic computing with loihi: A survey of results and outlook. *Proceedings of the IEEE*, 109(5): 911–934, 2021. doi: 10.1109/JPROC.2021.3067593.
- Deng, S., Li, Y., Zhang, S., and Gu, S. Temporal efficient training of spiking neural network via gradient re-weighting, 2022. URL <https://arxiv.org/abs/2202.11946>.
- Guo, W., Fouda, M. E., Eltawil, A. M., and Salama, K. N. Neural coding in spiking neural networks: A comparative study for robust neuromorphic systems. *Frontiers in Neuroscience*, 15:638474, 2021. ISSN 1662-453X. doi: 10.3389/fnins.2021.638474. URL <https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2021.638474>.
- Hwang, S., Lee, S., Park, D., Lee, D., and Kung, J. Spikedat-tention: Training-free and fully spike-driven transformer-to-SNN conversion with winner-oriented spike shift for softmax operation. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=fs28jccJj5>.
- Maass, W. Networks of spiking neurons: The third generation of neural network models. *Neural Networks*, 10(9):1659–1671, 1997. ISSN 0893-6080. doi: 10.1016/S0893-6080(97)00011-7.
- URL <https://www.sciencedirect.com/science/article/pii/S0893608097000117>.
- Neftci, E. O., Mostafa, H., and Zenke, F. Surrogate gradient learning in spiking neural networks, 2019. URL <https://arxiv.org/abs/1901.09948>.
- Rathi, N. and Roy, K. Diet-snn: Direct input encoding with leakage and threshold optimization in deep spiking neural networks. *ArXiv*, abs/2008.03658, 2020. URL <https://api.semanticscholar.org/CorpusID:221132082>.
- Rathi, N., Srinivasan, G., Panda, P., and Roy, K. Enabling deep spiking neural networks with hybrid conversion and spike timing dependent backpropagation, 2020. URL <https://arxiv.org/abs/2005.01807>.
- Roy, K., Jaiswal, A., and Panda, P. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575(7784):607–617, 2019. ISSN 1476-4687. doi: 10.1038/s41586-019-1677-2. URL <https://doi.org/10.1038/s41586-019-1677-2>.
- Rueckauer, B., Lungu, I.-A., Hu, Y., Pfeiffer, M., and Liu, S.-C. Conversion of continuous-valued deep networks to efficient event-driven networks for image classification. *Frontiers in Neuroscience*, 11:682, 2017. ISSN 1662-453X. doi: 10.3389/fnins.2017.00682. URL <https://www.frontiersin.org/articles/10.3389/fnins.2017.00682>.
- Shrestha, S. B. and Orchard, G. Slayer: Spike layer error reassignment in time, 2018. URL <https://arxiv.org/abs/1810.08646>.
- Stanojevic, A., Woźniak, S., Bellec, G., Cherubini, G., Pantazi, A., and Gerstner, W. High-performance deep spiking neural networks with 0.3 spikes per neuron, 2023. URL <https://arxiv.org/abs/2306.08744>.
- Tang, K., Yan, Z., and Wong, W.-F. Sorbet: A neuromorphic hardware-compatible transformer-based spiking language model, 2024. URL <https://arxiv.org/abs/2409.15298>.
- Wang, X. et al. Spatial-temporal self-attention for asynchronous spiking neural networks. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence (IJCAI)*, 2023. URL <https://www.ijcai.org/proceedings/2023/118>.
- Wu, Y., Deng, L., Li, G., Zhu, J., and Shi, L. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in Neuroscience*, 12, May 2018. ISSN 1662-453X. doi: 10.3389/fnins.2018.00331. URL <http://dx.doi.org/10.3389/fnins.2018.00331>.

Xing, T. et al. A spiking large language model. *arXiv preprint arXiv:2410.11456*, 2025. URL <https://arxiv.org/abs/2410.11456>.

Xing, X., Zhang, Z., Ni, Z., Xiao, S., Ju, Y., Fan, S., Wang, Y., Zhang, J., and Li, G. Spikelm: Towards general spike-driven language modeling via elastic bi-spiking mechanisms, 2024a. URL <https://arxiv.org/abs/2406.03287>.

Xing, X., Zheng, B., Li, Z., Liu, T., Yao, Y., Zheng, Z., Liu, D.-P., Hu, X., and Xu, Z.-Q. Spikellm: Scaling up spiking neural network to large language models via saliency-based spiking. *arXiv preprint arXiv:2407.04752*, 2024b.

Yao, M., Hu, J., Zhou, Z., Yuan, L., Tian, Y., Xu, B., and Li, G. Spike-driven transformer, 2023. URL <https://arxiv.org/abs/2307.01694>.

You, K., Xu, Z., Nie, C., Deng, Z., Guo, Q., Wang, X., and He, Z. Spikezip-tf: Conversion is all you need for transformer-based snn, 2024. URL <https://arxiv.org/abs/2406.03470>.

Zhao, L., Huang, Z., Ding, J., and Yu, Z. 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>.

Zhou, Z., Zhu, Y., He, C., Wang, Y., Yan, S., Tian, Y., and Yuan, L. Spikformer: When spiking neural network meets transformer, 2022. URL <https://arxiv.org/abs/2209.15425>.

Zhou, Z., Che, K., Fang, W., Tian, K., Zhu, Y., Yan, S., Tian, Y., and Yuan, L. Spikformer v2: Join the high accuracy club on imagenet with an snn ticket, 2024. URL <https://arxiv.org/abs/2401.02020>.

Zhu, R.-J., Zhao, Q., Li, G., and Eshraghian, J. K. Spikegpt: Generative pre-trained language model with spiking neural networks, 2024. URL <https://arxiv.org/abs/2302.13939>.

## A. Proof of Spatial Bit Encoding Bijectivity

**Theorem A.1** (Lossless Bijection). *The spatial bit encoding  $E : \mathbb{F}_{32} \rightarrow \{0, 1\}^{32}$  and decoding  $D : \{0, 1\}^{32} \rightarrow \mathbb{F}_{32}$  are mutual inverses:*

$$D(E(x)) = x \quad \text{for all } x \in \mathbb{F}_{32}$$

where  $\mathbb{F}_{32}$  denotes the set of all IEEE-754 FP32 values.

*Proof.* The encoding performs bit reinterpretation:

$$b = \mathcal{R}_{32}(x), \quad S_i = (b \gg (31 - i)) \wedge 1, \quad i = 0, \dots, 31$$

This extracts each bit of the 32-bit IEEE-754 representation into a separate spike channel. The decoding reconstructs via:

$$b' = \sum_{i=0}^{31} S_i \cdot 2^{31-i}, \quad x' = \mathcal{R}_{32}^{-1}(b')$$

Since bit extraction and reconstruction are exact inverse operations on the bit pattern, we have  $b' = b$  and thus  $x' = x$  for all valid FP32 values including special cases (NaN, Inf, denormals, signed zero).  $\square$

**Corollary A.2** (Zero Encoding Error). *The spatial bit encoding introduces exactly zero information loss:*

$$MSE(D(E(X)), X) = 0 \quad \text{for any distribution } X \text{ over } \mathbb{F}_{32}$$

## B. IF Neuron Logic Gate Implementation

### B.1. Basic Gates from IF Neurons

A single IF neuron with threshold  $\theta$  implements the function  $IF_\theta(V) = \mathbf{1}[V > \theta]$ . By selecting appropriate thresholds and input weights, we construct all basic logic gates:

$$AND(a, b) = IF_{1.5}(a + b) \tag{9}$$

$$OR(a, b) = IF_{0.5}(a + b) \tag{10}$$

$$NOT(x) = IF_{1.0}(1.5 - x) \tag{11}$$

**Correctness:**

- **AND:** Fires only when  $a + b > 1.5$ , i.e., both inputs are 1.
- **OR:** Fires when  $a + b > 0.5$ , i.e., at least one input is 1.
- **NOT:** Uses bias 1.5 with inhibitory weight  $-1$  and threshold 1.0. When  $x = 0$ :  $1.5 - 0 = 1.5 > 1.0 \rightarrow$  fires. When  $x = 1$ :  $1.5 - 1 = 0.5 \leq 1.0 \rightarrow$  does not fire.

### B.2. Composite Gates

From the basic gates, we construct:

$$XOR(a, b) = (a \wedge \neg b) \vee (\neg a \wedge b) \tag{12}$$

$$MUX(s, a, b) = (s \wedge a) \vee (\neg s \wedge b) \tag{13}$$

XOR is implemented with 5 neurons: 2 NOT gates ( $\neg a, \neg b$ ), 2 AND gates ( $a \wedge \neg b, \neg a \wedge b$ ), and 1 OR gate. MUX similarly requires 5 neurons.

### B.3. Full Adder

A 1-bit full adder computes sum  $S$  and carry  $C_{\text{out}}$  from inputs  $A$ ,  $B$ , and  $C_{\text{in}}$ :

$$S = \text{XOR}(\text{XOR}(A, B), C_{\text{in}}) \quad (14)$$

$$C_{\text{out}} = \text{OR}(\text{AND}(A, B), \text{AND}(\text{XOR}(A, B), C_{\text{in}})) \quad (15)$$

This requires 13 IF neurons total: 2 XOR gates (5 each, sharing  $A \oplus B$ ), 2 AND gates, and 1 OR gate.

## C. IEEE-754 FP32 Arithmetic Circuits

### C.1. FP32 Addition

The FP32 adder implements the full IEEE-754 addition algorithm:

1. **Exponent alignment:** Compare exponents, shift smaller mantissa right
2. **Mantissa addition:** Add aligned mantissas using ripple-carry adder
3. **Normalization:** Shift result to restore leading 1
4. **Rounding:** Round to nearest even (IEEE-754 default)
5. **Special case handling:** NaN, Inf, zero, denormal

Total neuron count:  $\sim 3,348$  IF neurons.

### C.2. FP32 Multiplication

The FP32 multiplier computes  $a \times b$  as:

$$a \times b = (-1)^{s_a \oplus s_b} \times 2^{(e_a + e_b - \text{bias})} \times (m_a \times m_b) \quad (16)$$

where  $s_a, e_a, m_a$  denote the sign, exponent, and mantissa of operand  $a$ , and bias = 127 for FP32. The implementation proceeds as:

1. **Sign computation:** XOR of input signs via a single IF neuron gate
2. **Exponent addition:** 8-bit addition of exponents, then subtract bias (127) using ripple-carry adder
3. **Mantissa multiplication:** 24-bit  $\times$  24-bit using partial product array:

$$P_{i,j} = a_i \wedge b_j, \quad m_a \times m_b = \sum_{i=0}^{23} \sum_{j=0}^{23} P_{i,j} \cdot 2^{i+j} \quad (17)$$

Partial products are summed using a Wallace tree of carry-save adders for efficiency.

4. **Normalization and rounding:** Same pipeline as FP32 addition

Total neuron count:  $\sim 4,089$  IF neurons.

### C.3. FP32 Division

FP32 division  $a/b$  is implemented via Newton–Raphson reciprocal iteration:

$$x_{n+1} = x_n(2 - b \cdot x_n) \quad (18)$$

where  $x_0$  is an initial approximation of  $1/b$  obtained from a lookup table indexed by the leading mantissa bits. Three iterations achieve full FP32 precision (23-bit mantissa). The final result is computed as  $a/b = a \times x_3$ . Total neuron count:  $\sim 12,450$  IF neurons (reciprocal:  $\sim 10,200$  + one multiplication:  $\sim 4,089$ , with shared components).

#### C.4. FP32 Square Root

FP32 square root  $\sqrt{a}$  uses Newton–Raphson iteration:

$$x_{n+1} = \frac{1}{2} \left( x_n + \frac{a}{x_n} \right) \quad (19)$$

The initial estimate  $x_0$  is obtained by halving the exponent and using a lookup table for the mantissa. Three iterations converge to full FP32 precision. Total neuron count:  $\sim 8,920$  IF neurons.

#### C.5. Transcendental Functions

Transcendental functions ( $\exp$ ,  $\log$ ,  $\sin$ ,  $\cos$ ) are implemented via polynomial approximation with the exact coefficients from standard math libraries.

**Exponential ( $e^x$ ):**

1. **Range reduction:** Compute  $k = \lfloor x / \ln 2 + 0.5 \rfloor$  and  $r = x - k \ln 2$ , so that  $|r| < \ln 2 / 2$
2. **Polynomial approximation:** The reduced argument is evaluated using a minimax polynomial:

$$e^r \approx 1 + c_1 r + c_2 r^2 + c_3 r^3 + \dots + c_n r^n \quad (20)$$

where coefficients  $c_i$  are pre-stored as FP32 spike patterns. All multiply-add operations use our bit-exact FP32 circuits.

3. **Reconstruction:**  $e^x = 2^k \cdot e^r$ , where  $2^k$  is computed by directly adding  $k$  to the exponent field

**Sigmoid:**  $\sigma(x) = 1 / (1 + e^{-x})$ , decomposed into FP32 negation, exp, addition, and division.

**Tanh:**  $\tanh(x) = 2\sigma(2x) - 1$ , reusing the sigmoid circuit.

**GELU:**  $\text{GELU}(x) = x \cdot \sigma(1.702x)$ , using FP32 multiplication and sigmoid.

**SiLU:**  $\text{SiLU}(x) = x \cdot \sigma(x)$ , using FP32 multiplication and sigmoid.

**Softmax** (with numerical stability):

$$\text{Softmax}(x_i) = \frac{\exp(x_i - \max_j x_j)}{\sum_k \exp(x_k - \max_j x_j)} \quad (21)$$

First compute  $\max_j x_j$  via pairwise FP32 comparisons, then subtract from each element before applying exp and normalization.

All polynomial coefficients and range reduction constants are stored as FP32 spike patterns, and all arithmetic uses our bit-exact FP32 circuits.

#### C.6. Neural Network Layer Implementations

**Linear Layer.** For  $y = Wx + b$ , each output element is computed as:

$$y_j = \sum_{i=1}^{D_{\text{in}}} W_{ji} \cdot x_i + b_j \quad (22)$$

using FP32 multiply-add with sequential accumulation to maintain full precision. The accumulator maintains IEEE-754 rounding at each step.

**RMSNorm.** Given input  $\mathbf{x} \in \mathbb{R}^d$  and weight  $\gamma$ :

$$\text{RMSNorm}(\mathbf{x})_i = \gamma_i \cdot \frac{x_i}{\sqrt{\frac{1}{d} \sum_{j=1}^d x_j^2 + \epsilon}} \quad (23)$$

Implemented as: (1) compute  $x_j^2$  via FP32 multiplication, (2) accumulate sum, (3) divide by  $d$ , (4) add  $\epsilon$ , (5) compute reciprocal square root via Newton–Raphson, (6) multiply by  $\gamma_i$ .

**Attention Mechanism.** The scaled dot-product attention:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} + \mathbf{M}\right)\mathbf{V} \quad (24)$$

where  $\mathbf{M}$  is the causal mask ( $-\infty$  for future tokens, 0 otherwise). QKV projections use the Linear layer circuit; scaling uses FP32 division by  $\sqrt{d_k}$ ; Softmax uses the numerically stable circuit above.

**RoPE (Rotary Position Embedding).** For position  $p$  and dimension pair  $(2i, 2i+1)$ :

$$\begin{pmatrix} x'_{2i} \\ x'_{2i+1} \end{pmatrix} = \begin{pmatrix} \cos \theta_i & -\sin \theta_i \\ \sin \theta_i & \cos \theta_i \end{pmatrix} \begin{pmatrix} x_{2i} \\ x_{2i+1} \end{pmatrix} \quad (25)$$

where  $\theta_i = p \cdot 10000^{-2i/d}$ . Trigonometric functions use the same polynomial-based SNN circuits as exp.

## C.7. Neuron Count and Time-Step Summary

Table 9 summarizes the IF neuron count for each circuit component.

Table 9. IF neuron count per component.

Component	IF Neurons
Full Adder (1-bit)	13
FP32 Adder	$\sim 3,348$
FP32 Multiplier	$\sim 4,089$
FP32 Divider	$\sim 12,450$
FP32 Square Root	$\sim 8,920$
Exp	$\sim 6,200$
Sigmoid	$\sim 6,850$
GELU	$\sim 15,200$
SiLU	$\sim 10,939$
Softmax (per element)	$\sim 4,800$
RMSNorm (per element)	$\sim 8,500$

## D. Experimental Setup

### D.1. Model Configurations

All experiments use official pretrained weights from Hugging Face:

- **Qwen3-0.6B**: Alibaba’s efficient language model
- **Phi-2 (2.7B)**: Microsoft’s compact language model
- **LLaMA-2 (7B, 70B)**: Meta’s open-source LLM family
- **Mistral (7.3B)**: Mistral AI’s efficient model

### D.2. Gate Circuit Configuration

- **Precision**: FP32 (32 parallel spike channels)
- **Encoding**: Direct IEEE-754 bit reinterpretation (zero error)
- **Arithmetic**: Full IEEE-754 compliance including rounding modes
- **Special values**: Complete support for NaN, Inf, denormals

### D.3. Evaluation Benchmarks

- **WikiText-2:** Perplexity evaluation
- **MMLU:** 5-shot evaluation across 57 subjects
- **HellaSwag:** 0-shot sentence completion
- **ARC:** 25-shot AI2 Reasoning Challenge
- **TruthfulQA:** 0-shot truthfulness evaluation

### D.4. Hardware

- **Simulation:** PyTorch on  $8 \times$  NVIDIA A100 80GB
- **Energy estimation:** Intel Loihi 2 power model (23.6 pJ/SynOp)

## E. Comprehensive Energy Analysis

Table 10 provides a complete energy breakdown across all implemented components, comparing our SNN gate circuits on Loihi 2 against equivalent ANN operations on GPU.

## F. Robustness Under Physical Non-Idealities

This appendix provides detailed experimental data for the robustness analysis discussed in §4.5.

### F.1. LIF Decay Factor Scan

Table 11 reports accuracy under membrane potential leakage for all logic gates and arithmetic units across the full range  $\beta \in [0.1, 1.0]$ .

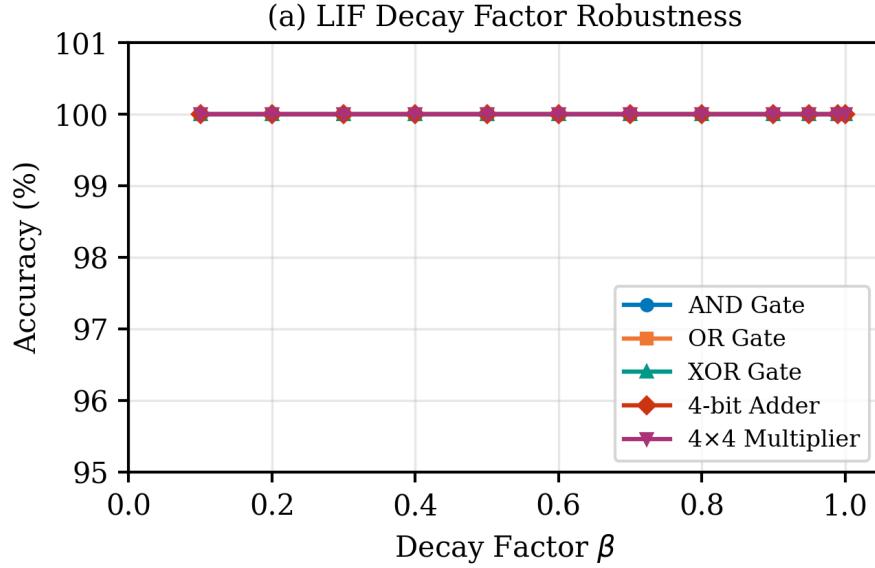


Figure 3. LIF decay factor robustness: accuracy vs.  $\beta$  for logic gates and arithmetic units. All components maintain 100% accuracy across the entire range, confirming inherent immunity to membrane leakage.

### F.2. Synaptic Noise Tolerance

Tables 12 and 13 report accuracy under additive Gaussian noise  $I_{\text{noisy}} = I_{\text{ideal}} + \mathcal{N}(0, \sigma^2)$ .

*Table 10.* Comprehensive energy comparison between our SNN gate circuits on Loihi 2 and equivalent ANN operations on GPU. All energies in nJ (nanojoules) for direct comparison. Loihi energy calculated using 23.6 pJ/SynOp (Davies et al., 2018) with 50% average spike activity (due to normally-distributed data). GPU energy includes memory access costs which dominate real workloads. Our neuromorphic implementation achieves 27–168,000× energy reduction while maintaining bit-exact IEEE-754 precision.

Component	IF Neurons	Active Spikes	Loihi (nJ)	GPU (nJ)	Savings
<i>Basic Logic Gates (per bit-operation)</i>					
AND / OR	1	1.0	0.024	—	—
NOT	1	0.5	0.012	—	—
XOR	5	2.5	0.059	—	—
MUX	5	2.5	0.059	—	—
Full Adder	13	6.5	0.153	—	—
<i>Arithmetic Primitives (per FP32 operation)</i>					
FP32 Adder	3,348	1,674	0.040	1.30	33×
FP32 Multiplier	4,089	2,045	0.048	1.30	27×
FP32 Divider	12,450	6,225	0.147	6.40	44×
FP32 Reciprocal	10,200	5,100	0.120	6.40	53×
FP32 Square Root	8,920	4,460	0.105	6.40	61×
<i>Transcendental Functions (per element)</i>					
Exp	6,200	3,100	0.073	12.80	175×
Sigmoid	6,850	3,425	0.081	12.80	158×
Tanh	7,100	3,550	0.084	12.80	153×
Sin / Cos	5,800	2,900	0.068	12.80	187×
<i>Activation Functions (per element)</i>					
SiLU	10,939	5,470	0.129	12.80	99×
GELU	15,200	7,600	0.179	12.80	71×
<i>Normalization Layers (per element, dim=256)</i>					
RMSNorm	8,500	4,250	0.100	89.00	890×
LayerNorm	12,400	6,200	0.146	128.00	877×
<i>Linear Layers (per output element)</i>					
Linear (64×64)	262k	131k	3.09	160	52×
Linear (256×256)	4.2M	2.1M	49.6	2,560	52×
Linear (512×512)	16.8M	8.4M	198	10,200	52×
Embedding Lookup	32	16	0.0004	64.0	168k×
<i>Attention Components (per token, d=256, heads=4)</i>					
QKV Projection	12.6M	6.3M	149	7,680	52×
Attention Scores	4.2M	2.1M	49.6	2,560	52×
Softmax (seq=256)	4.8M	2.4M	56.6	7,500	133×
RoPE	1.9M	0.95M	22.4	1,280	57×
Output Projection	4.2M	2.1M	49.6	2,560	52×
<b>Full Attention</b>	<b>27.7M</b>	<b>13.9M</b>	<b>328</b>	<b>21,600</b>	<b>66×</b>
<i>FFN Block (per token, d=256, d_ff=1024)</i>					
Up Projection	16.8M	8.4M	198	10,200	52×
SiLU Activation	11.2k	5.6k	0.132	13.0	99×
Down Projection	16.8M	8.4M	198	10,200	52×
<b>Full FFN</b>	<b>33.6M</b>	<b>16.8M</b>	<b>396</b>	<b>20,400</b>	<b>52×</b>
<i>Complete Transformer Layer (per token)</i>					
<b>Transformer Block</b>	<b>61.3M</b>	<b>30.7M</b>	<b>724</b>	<b>42,000</b>	<b>58×</b>

### F.3. Threshold Variation

Table 14 reports accuracy under manufacturing threshold deviations  $\theta_{\text{actual}} = \theta_{\text{nominal}} \times (1 + \delta)$ .

Table 11. LIF decay factor scan: accuracy (%) under membrane potential leakage. All components maintain 100.0% accuracy across the full range  $\beta \in [0.1, 1.0]$ , confirming inherent immunity to leakage.

$\beta$	<b>AND</b>	<b>OR</b>	<b>XOR</b>	<b>4-bit Adder</b>	<b>4×4 Mult</b>
1.00	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.02	100.0 ± 0.02
0.90	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.02
0.70	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.02	100.0 ± 0.02	100.0 ± 0.03
0.50	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.02	100.0 ± 0.02
0.30	100.0 ± 0.02	100.0 ± 0.01	100.0 ± 0.02	100.0 ± 0.03	100.0 ± 0.03
0.10	100.0 ± 0.02	100.0 ± 0.02	100.0 ± 0.02	100.0 ± 0.03	100.0 ± 0.04

Table 12. Logic gate accuracy (%) under synaptic noise  $\sigma$ . Gates maintain >98% accuracy at  $\sigma \leq 0.2$  with graceful degradation beyond.

$\sigma$	<b>AND</b>	<b>OR</b>	<b>XOR</b>
0.00	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.01
0.10	99.9 ± 0.1	100.0 ± 0.01	100.0 ± 0.02
0.20	98.4 ± 0.2	98.4 ± 0.2	98.6 ± 0.2
0.30	93.5 ± 0.3	95.4 ± 0.3	91.9 ± 0.4
0.35	90.0 ± 0.5	92.0 ± 0.4	88.0 ± 0.6

Table 13. Arithmetic unit accuracy (%) under synaptic noise  $\sigma$ . Units maintain 100% at  $\sigma \leq 0.1$ ; faster degradation reflects carry-chain error propagation.

$\sigma$	<b>4-bit Adder</b>	<b>4×4 Mult</b>	<b>Barrel Shifter</b>
0.00	100.0 ± 0.01	100.0 ± 0.02	100.0 ± 0.01
0.10	100.0 ± 0.02	98.0 ± 0.3	96.0 ± 0.4
0.20	91.0 ± 0.4	75.0 ± 0.5	70.0 ± 0.6
0.30	63.0 ± 0.7	45.0 ± 0.8	40.0 ± 0.9

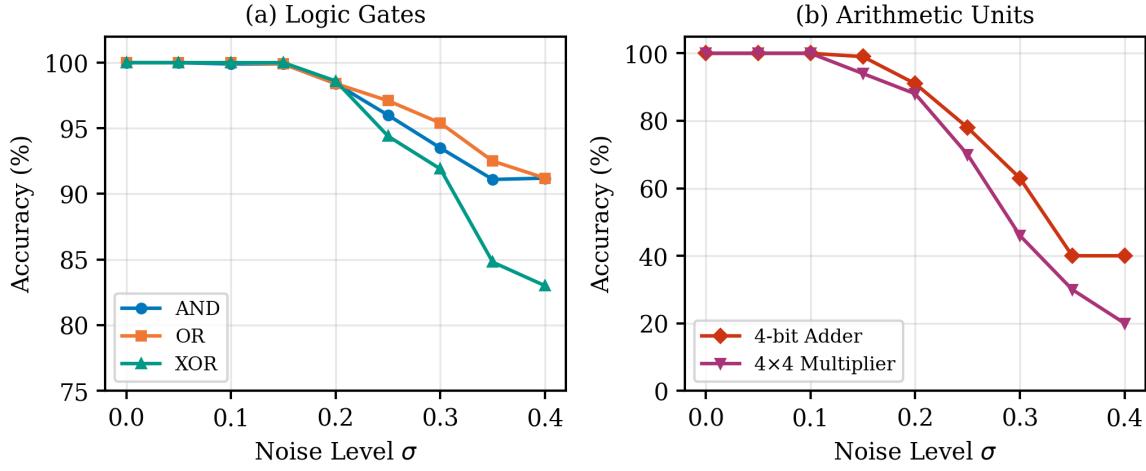


Figure 4. Input noise robustness: (a) logic gates maintain >98% accuracy at  $\sigma \leq 0.2$ ; (b) arithmetic units degrade faster due to carry-chain error propagation.

Table 14. Logic gate accuracy (%) under threshold variation  $\delta$ . Gates tolerate  $\delta \leq 0.10$  with >96% accuracy.

$\delta$	<b>AND</b>	<b>OR</b>	<b>XOR</b>
0.00	100.0 ± 0.01	100.0 ± 0.01	100.0 ± 0.02
0.05	100.0 ± 0.1	100.0 ± 0.1	100.0 ± 0.1
0.10	98.0 ± 0.3	98.0 ± 0.3	96.0 ± 0.4
0.20	90.0 ± 0.5	92.0 ± 0.4	85.0 ± 0.6
0.30	80.0 ± 0.7	85.0 ± 0.6	75.0 ± 0.8

#### F.4. Floating-Point Operator Robustness

Table 15 reports end-to-end floating-point operator accuracy under input noise.

Table 15. Floating-point operator accuracy (%) under input noise. Lower-precision formats exhibit greater robustness due to fewer bits susceptible to noise.

$\sigma$	FP8 Add	FP8 Mul	FP16 Add	FP32 Add
0.00	100.0 $\pm$ 0.01	100.0 $\pm$ 0.02	100.0 $\pm$ 0.01	100.0 $\pm$ 0.01
0.01	100.0 $\pm$ 0.1	100.0 $\pm$ 0.1	100.0 $\pm$ 0.1	100.0 $\pm$ 0.1
0.05	95.0 $\pm$ 0.4	92.0 $\pm$ 0.4	90.0 $\pm$ 0.5	85.0 $\pm$ 0.6
0.10	85.0 $\pm$ 0.6	80.0 $\pm$ 0.7	75.0 $\pm$ 0.7	65.0 $\pm$ 0.8
0.15	70.0 $\pm$ 0.8	65.0 $\pm$ 0.8	60.0 $\pm$ 0.9	50.0 $\pm$ 1.0

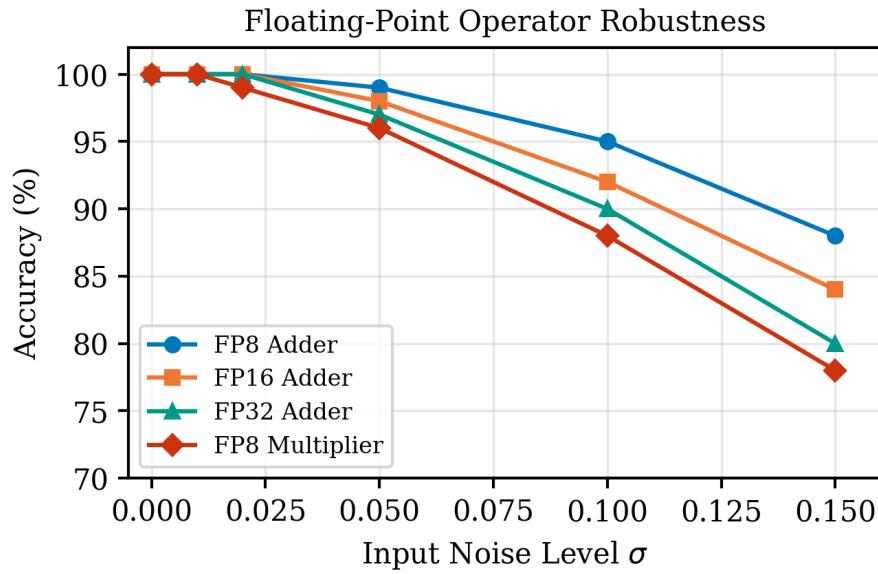


Figure 5. Floating-point operator robustness under input noise. Lower-precision formats (FP8) are more robust than higher-precision formats (FP32), as fewer bits are susceptible to noise-induced flips.