

Implementation Roadmap for Spiking Decision Transformer Novel Contributions

Prepared by Vishal Pandey

June 28, 2025

Abstract

This document provides a step-by-step LaTeX-formatted roadmap for integrating five novel modules into the Spiking Decision Transformer (SNN-DT) codebase. Each phase describes the high-level design, implementation steps, and validation strategy.

Contents

1	Phase 1: Adaptive, Data-Driven Temporal Windows	2
1.1	High-Level Design	2
1.2	Implementation Steps	2
1.3	Ablation & Validation	2
2	Phase 2: Hybrid Local Plasticity + Surrogate-Gradient	2
2.1	High-Level Design	2
2.2	Implementation Steps	2
2.3	Ablation & Validation	3
3	Phase 3: Spike-Domain Positional Encodings & Routing	3
3.1	Positional Spiking Codes	3
3.2	Dendritic-Style Routing	3
3.3	Implementation Steps	3
3.4	Ablation & Validation	3
4	Phase 4: Theoretical Convergence & Expressivity	3
4.1	Expressivity Theorem	3
4.2	Convergence Bound	3
4.3	Implementation Steps	4
4.4	Validation	4
5	Phase 5: Scalable, Sparse Spiking Attention	4
5.1	Locality-Sensitive Hashing (LSH)	4
5.2	Block-Sparse Attention	4
5.3	Implementation Steps	4
5.4	Ablation & Validation	4

1 Phase 1: Adaptive, Data-Driven Temporal Windows

1.1 High-Level Design

- Learn a per-token gate:

$$g_i = \sigma(w_g^\top x_i + b_g) \in [0, 1].$$

- Define token-specific window length:

$$T_i = \lceil T_{\max} \cdot g_i \rceil.$$

- Unroll each token’s LIF projection over its own T_i .

1.2 Implementation Steps

1. **Model:** In `SpikingSelfAttention`, add

```
self.window_gate = nn.Linear(hidden_dim, 1)
```

followed by a `sigmoid`.

2. Compute $g = \text{window_gate}(x)$ after embedding.

3. Refactor the time loop:

```
for i, xi in enumerate(tokens):  
    Ti = ceil(T_max * g[i])  
    for t in range(Ti):  
        # LIF projection ...
```

4. **Regularization:** Add penalty $\lambda \mathbb{E}[T_i]$ to the loss.

1.3 Ablation & Validation

- Plot average T_i vs. token index or return-to-go.
- Compare reward and spike counts to fixed- T .

2 Phase 2: Hybrid Local Plasticity + Surrogate-Gradient

2.1 High-Level Design

Combine global backprop with a local three-factor rule in the output LIF layer:

$$\Delta W_O \propto \underbrace{\sum_t \text{pre}(t) \text{post}(t)}_{\text{eligibility}} \times R_t.$$

2.2 Implementation Steps

1. Subclass Nornse’s `LIFCell` to accumulate eligibility traces.
2. After each trajectory, compute and apply:

$$W_O \leftarrow W_O + \eta_{\text{local}} e_{ij} G_t.$$

3. Normalize/clamp the local update to stabilize training.

2.3 Ablation & Validation

- Measure epochs to target reward with/without local plasticity.
- Report training curves for both variants.

3 Phase 3: Spike-Domain Positional Encodings & Routing

3.1 Positional Spiking Codes

Encode position via learned oscillators:

$$\text{pos_spike}_k(t) = \mathbf{1}(\sin(\omega_k t + \phi_k) > 0).$$

Learnable parameters: $\{\omega_k, \phi_k\}$.

3.2 Dendritic-Style Routing

After computing H heads' outputs $\{y_i^{(h)}(t)\}$, apply a routing MLP:

$$\alpha = \text{softmax}(W_{\text{route}} [y_i^{(1)}, \dots, y_i^{(H)}]).$$

Re-weight heads by α .

3.3 Implementation Steps

1. In embedding, generate `phase_spikes` alongside rate spikes.
2. In multi-head wrapper:

```
concat = torch.stack(head_outputs, dim=-1) # [..., H]
alpha = F.softmax(self.route_mlp(concat), dim=-1)
mixed = (concat * alpha).sum(-1)
```

3.4 Ablation & Validation

- Compare performance with/without phase coding.
- Visualize learned ω_k, ϕ_k .

4 Phase 4: Theoretical Convergence & Expressivity

4.1 Expressivity Theorem

Claim. For any dense attention matrix $A \in \mathbb{R}^{L \times L}$ and $\varepsilon > 0$, there exist spike trains of length $T = O(\log \frac{1}{\varepsilon})$ such that

$$\|\text{softmax}(\alpha S) - A\|_{\infty} < \varepsilon.$$

4.2 Convergence Bound

Under Lipschitz surrogate gradients, SNN-DT gradient descent converges to ANN-DT gradients as spike counts increase.

4.3 Implementation Steps

- Formalize assumptions (bounded weights, Lipschitz constant L_σ).
- Write proof sketch in a new `theory.tex` appendix.

4.4 Validation

Empirically plot $\|W_{\text{SNN}} - W_{\text{ANN}}\|$ vs. average spikes T .

5 Phase 5: Scalable, Sparse Spiking Attention

5.1 Locality-Sensitive Hashing (LSH)

Hash accumulated spike vectors $q_i \in \{0, 1\}^T$ into buckets; compute attention only within buckets.

5.2 Block-Sparse Attention

Divide the sequence into blocks of size B . Compute full attention within each block and use a global “summary” head across block means.

5.3 Implementation Steps

1. Implement LSH hashing: `bucket = torch.sign(random_proj @ qi)`
2. Refactor double loop:

```
for block in blocks:
    # intra-block attention
    # summary head attends across block means
```

5.4 Ablation & Validation

- Measure spikes & latency up to $N = 500$.
- Show returns remain within 95% of dense attention.

Version Control & Workflow

- Use separate feature branches: `adaptive-window`, `local-plasticity`, etc.
- After each merge, run benchmarks to catch regressions.
- Maintain consistent style (e.g. Google Python Style Guide).