

Challenge #26: BrainChip's IP for Targeting AI Applications at the Edge

Course: ECE 410/510 • **Spring 2025**

Topic: BrainChip's Temporal Event-based Neural Network (TENN) vs GPUs and Other Neuromorphic Approaches

Summary

After listening to the EETimes podcast featuring BrainChip engineers and executives, this write-up compares their **Temporal Event-based Neural Networks (TENN)** and **Akida chip architecture** to traditional **GPUs** and other **neuromorphic systems** like LSTMs and Transformers. Key differentiators, efficiency metrics, and their implications for edge AI are discussed.

1. Architectural Differences

- **BrainChip's Akida/TENN Architecture:**
 - Combines **event-based processing** and **state space modeling**.
 - Uses **internal states** (like memory) to process temporal data efficiently.
 - Based on **Legendre polynomial projections** to represent smooth, continuous temporal behavior.
- **GPUs:**
 - Suited for large-scale parallel computation.
 - Power-hungry and memory-intensive.
 - Optimized for traditional deep learning models (e.g., CNNs, Transformers).

- **Other Neuromorphic Chips (e.g., Loihi):**

- Spike-based, biologically inspired.
- Often emphasize event-driven communication but can lack training flexibility or standardization.

2. Efficiency and Memory Footprint

- **TENN:**

- Achieves **90%+ sparsity** in activation.
- Requires **less memory** due to internal state summarization.
- Efficient **spatio-temporal compression** avoids buffering multiple frames.
- Power usage is often **1/10th of conventional models**.

- **GPUs:**

- Require buffering, especially for temporal tasks.
- Less suited for sparse event-based workloads.
- Optimized for full-batch training, not low-power inference.

3. Training Methods and Ease

- **TENN:**

- Training is done in **convolutional (parallel) mode**.
- After training, models are folded into **recurrent form** for efficient inference.
- Compatible with transformer-like training pipelines but deploys like an RNN.
- **LSTM / RNNs:**
 - Training is often **sequential**, limiting parallelism.
 - Can become unstable, especially with long contexts.
- **Transformers:**
 - Fast, parallelizable training.
 - Require large amounts of memory and data.
- **Mamba:**
 - Also uses state-space modeling.
 - Larger internal states → better for GPUs, but poor edge compatibility.

4. Edge AI Workload Suitability

TENN excels in:

- Audio denoising
- Keyword spotting

- Eye and gesture tracking (event-based vision)
- Biomedical wearables (e.g., heart rate monitoring)
- Real-time ASR (Automatic Speech Recognition)
- Lightweight Large Language Models (Tiny LLaMA/LLaMA 3.2 class)

Advantages:

- **Causal inference:** Real-time predictions without future context.
- Tiny models with **competitive accuracy** across tasks.
- No need to be domain experts to get great results (as seen in CVPR competitions).

5. Causal vs. Non-Causal Inference

- **TENN supports both** causal (real-time) and non-causal (look-ahead) modes.
- For **edge deployments**, causal mode is essential (e.g., speech-to-text live transcription).
- TENN hits the **Pareto frontier** between **latency** and **accuracy**, making it ideal for real-time embedded applications.

6. Specific Advantages over LSTM, Transformer, Mamba

Model	Pros	Cons	TENN Advantage
LSTM	Compact, temporal	Unstable, hard to train	TENN is stable & trainable via conv layers

Transformer	Parallel training, high accuracy	High memory, slow inference at edge	TENN trains like Transformer, deploys like RNN
Mamba	State-space, good LLM performance	Large internal states, not edge-friendly	TENN uses small state banks , efficient for edge LLMs

Event-Based Hardware Synergy

- TENN is **sparse by nature**, making it highly compatible with BrainChip's **Akida hardware**.
- On-chip processing of **non-zero events only** reduces unnecessary computation.
- Supports analog signals, spikes, and words—unified by the same TENN framework.

Future Vision

"We're neuro-inspired, not neuro-identical." — Dr. John Tapson

- **Akida 2.0** will support 1-bit, 4-bit, and 8-bit modes.
- BrainChip will likely move to **16-bit neuromorphic architectures** to support richer algorithms.
- Future lies in **programmable state-aware hardware**—blending the strengths of both neuromorphic and traditional computation.

Final Thoughts

BrainChip's approach marks a **shift from biologically literal neuromorphics to engineering-optimized brain-inspired AI**. Their ability to run efficient, recurrent, and causal models on-chip—with performance rivaling or surpassing much larger models—is a clear advantage for edge AI.

Their strategy of training like a Transformer but deploying like an RNN is particularly novel and effective.