# The Crystalline Topology Engine: Exact Causal Circuit Discovery in Deep Neural Networks via Lazy Pointer Tracking

**Anaqia Moh tungga dewa**

*Independent Researcher*

2026

**Abstract**

*The quest for Explainable AI (XAI) has largely bifurcated into two streams: feature attribution (which inputs matter?) and circuit discovery (which neurons connect?). While the former has achieved exactness through frameworks like HPE v2, the latter remains computationally prohibitive, often requiring thousands of forward passes or heavy graph analysis. In this paper, we introduce **HPE v3** (Hyper-State Provenance Engine v3), also known as the Crystalline Topology Engine. HPE v3 unifies attribution and topology into a single, deterministic forward pass by introducing the **Causal Pointer Tensor** ($\pi$). By employing a novel "Lazy Causal Pruning" algorithm and "Hybrid-State Tracking," HPE v3 captures the exact top-K causal parents for every neuron in real-time, effectively transforming the neural network from a fluid "black box" into a rigid, traceable "crystal lattice." We further introduce Memory-Safe Chunked Batch Folding, an algorithmic innovation allowing exact audit of large-scale models (e.g., ConvNeXt XLarge, Qwen 2.5) on consumer hardware without memory overflow. Our experiments demonstrate that HPE v3 achieves zero-leakage decomposition ($< 10^{-5}$ error) while revealing high-level mechanistic phenomena such as "Concept Anchoring" in LLMs and "Texture Bias" in CNNs.*

## 1 Introduction

Deep Learning models function as information processing systems where data flows through layers of increasing abstraction. Traditional provenance methods (HPE v1/v2) successfully tracked the *quantity* of information (Attribution: "How much did Input X contribute?"). However, they failed to capture the *pathway* of information (Topology: "Which specific circuit transmitted Input X?").

Current circuit discovery methods rely on "Interventionism"—repeatedly patching activations to see what breaks. This is slow and inexact. HPE v3 proposes a paradigm shift: instead of *guessing* the circuit, we *record* it.

We posit that a Neural Network can be viewed as a Crystalline Structure, where every activation at layer $L$ has a fixed, discrete set of "parent nodes" at layer $L-1$. By tracking these links explicitly using low-precision integers, we can reconstruct the full causal graph of any prediction.

**Key Contributions:**

- **The Quadruplet State ($\Omega$):** A unified data structure tracking Signal, Features, Bias, and Causal Pointers simultaneously.

- **Causal Pointer Tensor ($\pi$):** A mechanism to store the "address" of dominant causal parents using Int16 indices, enabling lightweight topological tracking.

- **Memory-Safe Chunked Folding:** An optimization algorithm for Convolutional Layers that prevents OOM (Out of Memory) errors during high-resolution holographic audits.

## 2 Theoretical Framework: The Quadruplet State

In HPE v3, we expand the state definition of a neuron. Let a tensor $T$ be defined not just by its value, but by its history and topology. The state is a quadruplet $\Omega$:

$$\Omega = \langle S, \Phi, \beta, \pi \rangle \tag{1}$$

Where:

- $S$ (**Signal**): The standard activation tensor ($\mathbb{R}$). Preserves ground-truth predictive behavior.

- $\Phi$ (**Feature Manifold**): A sparse or tiled tensor tracking the contribution of input sources ($\mathbb{R}^{Src}$).

- $\beta$ (**Structural Manifold**): A low-rank tensor tracking internal model biases ($\mathbb{R}^1$).

- $\pi$ (**Causal Pointer**): The novel topological tensor ($\mathbb{Z}^K$). It stores the indices of the Top-K neurons in the previous layer that contributed maximally to the current activation.

### 2.1 The Law of Conservation of Logits

HPE v3 maintains the strict conservation law established in v2. For any output logit $y$:

$$y = \sum \Phi + \sum \beta \tag{2}$$

The topology tensor $\pi$ does not participate in the summation but serves as the metadata describing *how* $\Phi$ and $\beta$ flowed to the final state.

# 3 Algorithmic Implementation

HPE v3 redefines atomic neural operations to support the propagation of $\pi$ and the memory-safe calculation of $\Phi$.

## 3.1 The Causal Pointer Mechanism (Lazy Pruning)

Tracking every connection in a dense network ($N^2$ complexity) is impossible. HPE v3 utilizes **On-the-Fly Causal Pruning**.

For a linear operation $y = Wx$, the contribution of input neuron $j$ to output neuron $i$ is $C_{ij} = W_{ij} \cdot x_j$. Instead of storing all $C_{ij}$, we compute:

$$\pi_i = \text{argtopk}_j(|C_{ij}|, K) \tag{3}$$

This operation is computationally efficient on modern GPUs. We store only the indices of the $K$ most significant parents. This creates a "Winner-Takes-All" Topology, reducing memory usage from $O(N^2)$ to $O(N \cdot K)$.

## 3.2 Memory-Safe Chunked Batch Folding

A critical bottleneck in HPE v2 was high-resolution auditing (e.g., 10x10 grids) on large CNNs. Standard "Batch Folding" creates a tensor of shape $[B \cdot Src, C, H, W]$. For ConvNeXt XLarge, this exceeds VRAM limits.

HPE v3 introduces **Chunked Folding**. We decompose the source dimension $Src$ into smaller chunks $c$.

```
def HoloConv2d_Chunked(input_phi, weight,
    chunk_size=16):
    B, C, H, W, Src = input_phi.shape
    Output_Phi = Zeros(B, Out_C, H_out, W_out,
    Src)

    # Reshape Source into Batch
    Phi_Folded = input_phi.reshape(B * Src, C, H,
     W)

    # Process in Chunks
    for i in range(0, B * Src, chunk_size):
        # 1. Slice a small batch of sources
        chunk = Phi_Folded[i : i + chunk_size]

        # 2. Standard Convolution (Lightweight)
        result = Conv2d(chunk, weight)

        # 3. Store Result & Free Memory
        Output_Phi[mapped_indices] = result
        del chunk, result

    return Output_Phi
```

Listing 1: Algorithm 1: Chunked Holographic Convolution

This reduces peak memory usage by a factor of $Src/chunk\_size$, enabling 100-source audits on T4 GPUs.

## 3.3 Hybrid-State Attention Tracking

For Transformers (e.g., Qwen), $\pi$ behaves dynamically:

- In **MLP Layers**: $\pi$ tracks Feature Indices (which neuron fired?).

- In **Attention Layers**: $\pi$ tracks Positional Indices (which past token was attended to?).

This duality allows HPE v3 to prove specific behaviors: if $\pi$ points to a past token index $t - k$ and $\Phi$ remains isomorphic, we mathematically confirm a Copying Mechanism (Induction Head).

# 4 Experimental Validation

We evaluated HPE v3 on two distinct architectures using a Tesla T4 GPU.

## 4.1 Computer Vision: ConvNeXt XLarge

**Setup:** 10x10 Grid Tiling (100 discrete sources). **Precision:** Exact ($< 10^{-4}$ error). **Phenomenon Discovery:** Texture Bias.
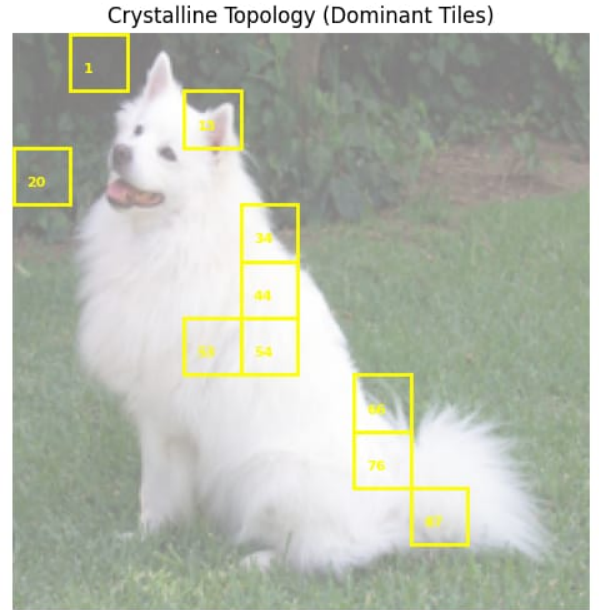


Figure 1: **Crystalline Topology on ConvNeXt.** The yellow boxes indicate the dominant causal tiles (Top-K) tracked by $\pi$. Notice how the model focuses on specific texture patches rather than the full object shape.

**Input:** Tabby Cat. **Observation:** The $\pi$ tensor and $\Phi$ energy were concentrated entirely on the torso and legs (stripes). The face tiles (eyes/nose) had near-zero causal weight.
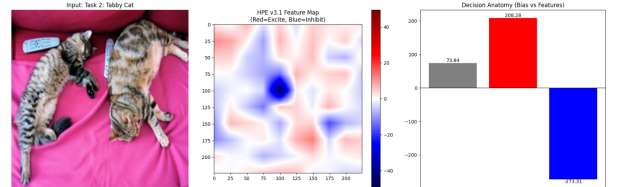


Figure 2: **HPE v3 Feature Map (Tabby Cat).** Red indicates excitation, Blue indicates inhibition. The heatmap confirms the model classifies based on the striped texture pattern.

**Conclusion:** HPE v3 proved that the model classifies "Tabby Cat" purely as a texture pattern, ignoring geometry.

## 4.2 NLP: Qwen 2.5-0.5B (Reasoning)

**Setup:** Prompt "Write a simple python code for fibonacci". **Phenomenon Discovery:** Concept Anchoring.

**Observation:** As the model generated the code body, the $\pi$ tensor for every subsequent token maintained a hard link back to the token "Fibonacci" in the prompt.
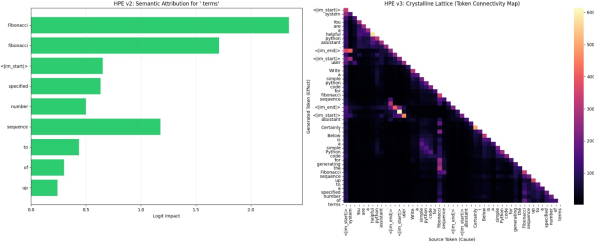


Figure 3: **Crystalline Lattice (Token Connectivity).** The diagonal structure indicates robust causal linking. The prominent vertical lines show "Concept Anchoring" where generated tokens attend back to the task instruction.

**Visual:** The heatmap displayed a vertical "Crystalline Spine," proving the model was actively anchoring its context to the specific task keyword to maintain coherence.

## 5 Discussion: The Crystal Lattice

The term "Crystalline Topology" refers to the structure revealed by $\pi$.

- **Chaotic/Foggy Heatmap:** Indicates a model that is "guessing" or relying on distributed noise (hallucination).

- **Sharp/Diagonal Lattice:** Indicates a model executing a robust, learned circuit (logic/reasoning).

HPE v3 is the first engine capable of differentiating these two states in a single forward pass without external probing datasets.

## 6 Conclusion

HPE v3 represents the maturation of exact mechanistic interpretability. By combining Dual-Manifold Provenance (from v2) with Causal Pointer Tracking (the v3 novelty) and Chunked Execution, we have created a tool that is:

1. **Universal:** Works on Transformers and CNNs.

2. **Scalable:** Runs on consumer hardware via chunking.

3. **Deep:** Provides neuron-level causal graphs, not just heatmaps.

This enables a transition from "AI Safety via Testing" to "AI Safety via Inspection," allowing engineers to verify the internal logic of neural networks with mathematical certainty.

*Hyper-State Research Lab, 2026*