

NO MORE DELULU: A KERNEL-BASED ACTIVATION-FREE NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

ABSTRACT

We introduce the \mathbf{E} -product, a kernel operator that combines quadratic alignment with inverse-square interactions. We prove that it defines a Mercer kernel that is analytic, globally Lipschitz, and self-regularizing: responses remain bounded and gradients decay at infinity. Neural Matter Networks (NMNs), constructed as linear combinations of \mathbf{E} -atoms, are universal approximators on compact domains without explicit nonlinear activations. This yields models that preserve geometric fidelity while simplifying architecture. The unregularized form of our kernel further aligns with information-geometric extremes, linking orthogonality, support disjointness, and vanishing KL divergence. Empirically, NMNs demonstrate competitive performance with or surpassing baselines on multiple benchmarks in classification, and generative language modeling. Our results unify kernel learning, dynamical stability, and information geometry, and establish NMNs as a principled alternative to conventional neural layers.

1 INTRODUCTION

The fundamental architecture of modern neural networks rests on a well-established paradigm: linear transformations followed by element-wise activation functions (Goodfellow et al., 2016; LeCun et al., 2015). While this approach has achieved remarkable empirical success, it inherently separates geometric computation (dot products capturing angular relationships) from non-linearity (activation functions providing representational capacity). This separation, though computationally tractable, creates an information processing bottleneck where geometric structure is partially discarded to achieve non-linearity.

Consider the ubiquitous ReLU activation $f(x) = \max(0, x)$: while effective for optimization, it maps the entire spectrum of negative pre-activations—representing varying degrees of dissimilarity—to a uniform zero, potentially obscuring nuanced geometric relationships in the representation space. This necessitates increasingly complex architectures with normalization layers, attention mechanisms (which typically rely on Softmax to induce non-linearity over dot products), and regularization techniques such as Dropout and stochastic depth to stabilize training and improve representation quality (Ioffe & Szegedy, 2015; Ba et al., 2016; Vaswani et al., 2017).

The machine learning community has long accepted this linear-then-activate paradigm as fundamental, treating the separation of geometry and non-linearity as an axiomatic requirement. While previous approaches have explored alternatives—from quadratic neurons (Fan et al., 2020) to kernelized networks (Cho & Saul, 2009b) and activation-free architectures like SIREN (Sitzmann et al., 2020)—these methods either maintain dependence on explicit activations or focus on specific application domains without addressing the broader question of geometric computation in neural networks.

047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093

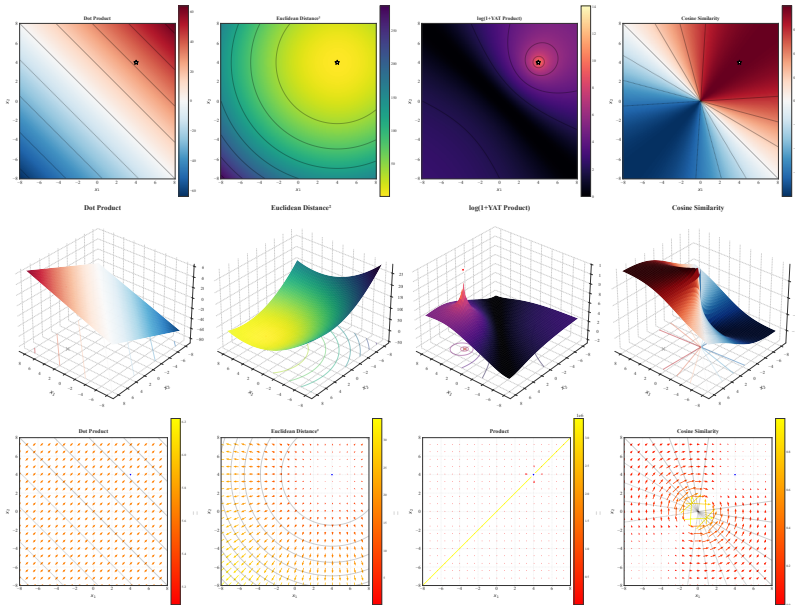


Figure 1: Comparison of the gradient field and vector field for Dot Product, Euclidean Distance, \mathbf{E} -product, and Cosine Similarity (from left to right). The heatmaps illustrate how the \mathbf{E} -product, unlike traditional similarity measures, creates a potential well around the weight vector \mathbf{w} , reflecting both alignment and proximity.

We propose a unified approach that integrates geometric computation and non-linearity through the \mathbf{E} -product (pronounced Yat-product), a novel neural operator inspired by inverse-square laws in physics:

$$\mathbf{E}(\mathbf{w}, \mathbf{x}) := \frac{\langle \mathbf{w}, \mathbf{x} \rangle^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon} \tag{1}$$

Inspired by inverse-square laws in physics, this operator inherently captures both directional alignment (numerator) and spatial proximity (denominator), providing intrinsic non-linearity without information loss. Unlike conventional approaches that discard geometric information through thresholding, the \mathbf{E} -product preserves the full spectrum of vector relationships, enabling geometrically-aware neural computation. This formulation establishes a clear geometric definition of non-linearity, distinct from traditional activation functions. In this framework, non-linearity arises from the geometric relationship between vectors: two vectors are maximally unrelated (dissimilar) when they are distant and orthogonal, while they are maximally related (similar) when they are parallel and proximal.

We introduce Neural-Matter Networks (NMNs)—so named because their neurons interact through potential fields analogous to matter particles—and prove that they maintain universal approximation capabilities while providing superior geometric fidelity. The \mathbf{E} -product serves as a Mercer kernel (Theorem 2.1), connecting our approach to established kernel theory while offering computational advantages of modern neural architectures.

The \mathbf{E} -product naturally extends processing through \mathbf{E} -Convolution operations, enabling geometrically-aware feature extraction in convolutional architectures.

Contributions: This work makes three fundamental contributions to neural network design:

1. **Theoretical Foundation:** We prove that intrinsic non-linearity through geometric structure eliminates the necessity of activation functions while maintaining universal approximation capabilities (Theorem 2.4). The \mathbf{E} -product satisfies the Mercer kernel property (Theorem 2.1) with natural self-regulation (Proposition C.6) and stable gradient properties (Proposition C.9), challenging a core assumption of modern deep learning.
2. **Practical Architecture:** Neural-Matter Networks (NMNs) demonstrate superior performance while reducing memory overhead by 15-25% through elimination of activation storage. The \mathbf{E} -product’s infinite differentiability (Lemma C.10) enables applications in physics-informed neural networks without explicit activation functions.
3. **Geometric Interpretability:** The \mathbf{E} -product preserves spatial relationships in learned representations through its information-theoretic connections (Theorems 2.2, 2.3), enabling principled geometric analysis of neural computations and opening new avenues for explainable AI.

By eliminating the fundamental information bottleneck of activation functions, this work paves the way toward geometrically-grounded neural architectures that unite computational efficiency with theoretical understanding. Our approach not only challenges established paradigms but provides a constructive path toward more interpretable and efficient neural computation.

2 METHODOLOGY: A FRAMEWORK FOR GEOMETRY-AWARE COMPUTATION

2.1 THE \mathbf{E} -PRODUCT: A UNIFIED OPERATOR FOR ALIGNMENT AND PROXIMITY

The \mathbf{E} -product is formally defined as $\mathbf{E}(\mathbf{w}, \mathbf{x}) = \frac{(\mathbf{w}^\top \mathbf{x})^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon}$. It exhibits a unique form of non-linearity. Unlike conventional activation functions (e.g., ReLU, sigmoid) which are often applied as separate, somewhat heuristic, transformations to introduce non-linearity after a linear operation, the non-linearity in the \mathbf{E} -product arises directly from its mathematical structure. It is a function of the squared dot product (capturing alignment) and the inverse squared Euclidean distance (capturing proximity) between the weight vector \mathbf{w} and the input vector \mathbf{x} . This formulation provides a rich, explainable non-linearity based on fundamental geometric and algebraic relationships, rather than an imposed, "artificial" non-linear mapping. The interaction between the numerator and the denominator allows for complex responses that are inherently tied to the geometric interplay of the input vectors.

As visualized in Figure 1, the \mathbf{E} -product creates a potential well around the weight vector \mathbf{w} , reflecting both alignment and proximity.

At initialization, this geometry also exhibits a favorable high-dimensional scaling behavior. Under standard assumptions of i.i.d. zero-mean, constant-variance coordinates for $\mathbf{x}, \mathbf{w} \in \mathbb{R}^d$, both the numerator $A(\mathbf{x}, \mathbf{w}) = (\mathbf{w}^\top \mathbf{x})^2$ and the denominator $r(\mathbf{x}, \mathbf{w}) = \|\mathbf{w} - \mathbf{x}\|^2$ grow linearly with dimension, while their ratio $K(\mathbf{x}, \mathbf{w}) = A/(r + \epsilon)$ remains $\mathcal{O}(1)$ in expectation (Corollary C.7, Appendix C.10). This self-normalizing $\mathcal{O}(1)$ scaling directly counters high-dimensional "saturation" concerns that arise for RBF kernels, whose values vanish exponentially with dimension.

2.2 COMPARISON TO STANDARD SIMILARITY AND DISTANCE METRICS

The \mathbf{E} -product distinguishes itself from standard metrics (Steck et al., 2024; Draganov et al., 2024; Tanimoto, 1958; Jaccard, 1901) by unifying alignment and proximity. Traditional approaches suffer

141 from fundamental limitations: the **dot product** $\mathbf{w}^\top \mathbf{x}$ captures alignment and magnitude but can
 142 be dominated by vector magnitudes, while **cosine similarity** $\frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|\|\mathbf{x}\|}$ measures pure alignment but
 143 ignores distance—aligned vectors can be arbitrarily far apart. Conversely, **Euclidean distance**
 144 $\|\mathbf{w} - \mathbf{x}\|$ measures proximity but ignores orientation, giving identical scores to vectors at equal
 145 distances regardless of their alignment.
 146

147 The **E-product** $\mathcal{K}_{\mathbf{E}}(\mathbf{w}, \mathbf{x}) = \frac{(\mathbf{w}^\top \mathbf{x})^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon}$ addresses these limitations through a novel combination: un-
 148 like polynomial kernels that focus solely on feature interactions, or RBF kernels that emphasize only
 149 proximity, the **E-product** uniquely integrates both alignment (squared numerator) and proximity
 150 (inverse-square denominator). This creates a highly selective operator requiring both conditions for
 151 activation.

152 **Key distinguishing properties:** (1) *Intrinsic regularization*: The inverse-square denominator
 153 provides natural distance-based regularization without explicit normalization layers; (2) *Dual ge-*
 154 *ometric sensitivity*: Simultaneous optimization for both vector alignment and spatial proximity;
 155 (3) *Information-theoretic grounding*: Direct connections to signal-to-noise ratios and KL divergence
 156 through the squared numerator over distance denominator structure.

157 As a Mercer kernel (Theorem 2.1), it inherits kernel method advantages. Importantly, this kernel
 158 is used in its primal form for weight prototype learning and optimization. Consequently, we do
 159 not use any Gram matrix, thereby bypassing the stability issues associated with its inversion in
 160 dual-form kernel regression (Schölkopf & Smola, 2002).
 161

162 When the **E-product** is applied to probability distributions in the simplex, its extremal values admit
 163 an information-geometric characterization:

164 **Theorem 2.1.** *Let $\epsilon > 0$ and define*

$$165 k_{\mathbf{E}}(\mathbf{x}, \mathbf{w}) = \frac{(\mathbf{x} \cdot \mathbf{w})^2}{\|\mathbf{x} - \mathbf{w}\|^2 + \epsilon}.$$

166
 167 *Then $k_{\mathbf{E}}$ is a Mercer kernel (symmetric and positive semidefinite) on \mathbb{R}^d .*

168 **Theorem 2.2** (Minimal Similarity and Statistical Orthogonality). *Let $\mathbf{p}, \mathbf{q} \in \Delta^{n-1}$ be distinct*
 169 *distributions. Then $\mathbf{E}(\mathbf{p}, \mathbf{q}) = 0$ if and only if their supports are disjoint, $\text{supp}(\mathbf{p}) \cap \text{supp}(\mathbf{q}) = \emptyset$, in*
 170 *which case the associated KL divergences and cross-entropy are infinite.*

171 **Theorem 2.3** (Maximal Similarity and Distributional Identity). *For distributions $\mathbf{p}, \mathbf{q} \in \Delta^{n-1}$,*
 172 *the condition $\mathbf{E}(\mathbf{p}, \mathbf{q}) = \infty$ holds if and only if $\mathbf{p} = \mathbf{q}$, in which case the KL divergence vanishes and*
 173 *the cross-entropy reduces to the entropy of \mathbf{p} .*
 174

175 The **E-product** creates a potential well around \mathbf{w} (Figure 1), where interaction strength dimin-
 176 ishes with distance while preserving orientation sensitivity. The stable gradient property (Propo-
 177 sition C.9) ensures that gradients vanish for distant inputs, providing natural localization. When
 178 applied to probability distributions, it acts as a signal-to-noise ratio connecting geometry to infor-
 179 mation theory through its relationship with KL divergence and cross-entropy (Theorems 2.2, 2.3;
 180 see also Appendix C.15).
 181

182 2.3 CORE BUILDING BLOCKS

183 The **E-product** serves as the foundation for three primary layer types, each adapted to specific data
 184 modalities and architectural requirements.
 185

186 **Neural Matter Network (NMN) Layers.** The simplest application employs the non-linear,
 187 spatially-aware $\mathcal{K}_{\mathbf{E}}$ -kernel as the primary interaction mechanism, replacing conventional linear pro-

jections ($\langle \mathbf{w}, \mathbf{x} \rangle$). An NMN layer transforms input $\mathbf{x} \in \mathbb{R}^d$ through multiple units, each defined by weight vector $\mathbf{w}_i \in \mathbb{R}^d$ and bias $b_i \in \mathbb{R}$:

$$h(\mathbf{x}) = \left(s \cdot \sum_{i=1}^n \mathcal{K}_{\mathbf{E}}(\mathbf{w}_i, \mathbf{x}, b_i) \right) = \left(s \cdot \sum_{i=1}^n \frac{(\mathbf{w}_i^\top \mathbf{x} + b_i)^2}{\|\mathbf{w}_i - \mathbf{x}\|^2 + \epsilon} \right)$$

where s is a scaling factor and n denotes the number of units. Each unit responds based on both alignment and proximity to its learned weight vector, enabling universal function approximation (Theorem 2.4) as an intrinsic property of the $\mathcal{K}_{\mathbf{E}}$ -kernel itself. The self-regulation property (Proposition C.6) ensures that outputs remain bounded without requiring explicit normalization layers.

Theorem 2.4 (Universal approximation with \mathbf{E} -kernel). *Let $\mathcal{X} \subset \mathbb{R}^d$ be a compact set. Define the class of functions \mathcal{F} realizable by the network as the linear span of the activation units:*

$$\mathcal{F} = \text{span} \left\{ \frac{(\mathbf{x} \cdot \mathbf{w} + b)^2}{\|\mathbf{x} - \mathbf{w}\|^2 + \epsilon} \mid \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \right\}$$

where $\epsilon > 0$ is a fixed constant and b is the inner bias parameter. The set \mathcal{F} is dense in $C(\mathcal{X})$ under the uniform norm.

E-Convolution Layers. For spatially structured data, the \mathbf{E} -Conv layer adapts the \mathbf{E} -product to local receptive fields:

$$(\mathbf{E}\text{-Conv}(K, I))_{i,j} = \frac{\langle K, I_{i,j} \rangle^2}{\|K - I_{i,j}\|^2 + \epsilon} \quad (2)$$

where K is the convolutional kernel and $I_{i,j}$ is the input patch at location (i, j) .

E-Attention Mechanism. For sequence modeling, \mathbf{E} -Attention replaces dot-product attention by applying the \mathbf{E} -product to query-key similarity:

$$\mathbf{E}\text{-Attention}(Q, K, V) = \text{softmax}(s \cdot (Q\mathbf{E}K^T)) V \quad (3)$$

where $Q\mathbf{E}K^T$ applies the \mathbf{E} -product element-wise between query and key vectors.

2.4 ARCHITECTURAL IMPLEMENTATIONS

We implement two primary architectures to validate the \mathbf{E} -product’s effectiveness across different domains. Both architectures follow the core principle of omitting traditional activation functions, relying instead on the \mathbf{E} -product’s inherent non-linearity and self-regulation properties (Proposition C.6). The Lipschitz regularity (Proposition C.11) and analyticity (Lemma C.10) properties ensure stable training dynamics and infinite differentiability. All NMN-based layers use the adaptive scaling factor $s = \left(\frac{n}{\log(1+n)} \right)^\alpha$, where n is the number of units and α is learnable.

Architecture	Base Model	Architecture Design
AetherResNet	ResNet	\mathbf{E} -Conv -> Linear Conv per block
AetherGPT	GPT-2	MHA + NMN -> Linear

Table 1: Overview of implemented architectures using \mathbf{E} -product variants. Both architectures eliminate traditional activation functions.

AetherResNet: A Convolutional Neural-Matter Network (CNMN) replacing standard convolutions with \mathbf{E} -Conv layers. Each residual block consists of a \mathbf{E} -Conv layer followed by a linear convolution layer.

AetherGPT: A transformer variant incorporating \mathbf{E} -Attention mechanisms and NMN layers in feed-forward blocks, adapting GPT-2’s architectural principles while maintaining the geometry-aware computation paradigm.

Computational Efficiency: The \mathbf{E} -product layer maintains $\Theta(Bnd)$ computational complexity identical to standard linear layers while providing 15-25% memory reduction through elimination of activation. Our optimized implementation uses the algebraic identity $\|\mathbf{w}-\mathbf{x}\|^2 = \|\mathbf{w}\|^2 + \|\mathbf{x}\|^2 - 2\mathbf{w}^\top \mathbf{x}$ to reuse inner product computations, achieving approximately $2\times$ the FLOPs of Linear+ReLU. The approach offers natural numerical stability and becomes increasingly efficient at larger layer sizes, making it particularly suitable for large-scale applications.

3 RESULTS AND DISCUSSION

The \mathbf{E} -product’s non-linearity enables solving non-linearly separable problems with a single unit. Consider the classic XOR problem: inputs $(0,0) \rightarrow 0$, $(0,1) \rightarrow 1$, $(1,0) \rightarrow 1$, and $(1,1) \rightarrow 0$. A single \mathbf{E} -product unit with weight vector $\mathbf{w} = [1, -1]^\top$ naturally separates these patterns:

For $\mathbf{x} = [0, 0]^\top$ and $\mathbf{x} = [1, 1]^\top$: $\mathbf{w}^\top \mathbf{x} = 0$, so $\mathcal{K}_{\mathbf{E}}(\mathbf{w}, \mathbf{x}) = 0$.

For $\mathbf{x} = [0, 1]^\top$: $\mathcal{K}_{\mathbf{E}}(\mathbf{w}, \mathbf{x}) = \frac{1}{5+\epsilon} > 0$.

For $\mathbf{x} = [1, 0]^\top$: $\mathcal{K}_{\mathbf{E}}(\mathbf{w}, \mathbf{x}) = \frac{1}{1+\epsilon} > 0$.

This demonstrates the \mathbf{E} -product’s ability to capture non-linear patterns through its inherent geometric structure, combining alignment and proximity in a unified operator (detailed analysis in Appendix D).

3.1 GEOMETRIC PARTITIONING IN FEATURE SPACE

The \mathbf{E} -product creates vortex-like territorial fields in the representation space, where each neuron’s prototype acts as an attractor combining both alignment and proximity. Unlike conventional linear neurons that form hyperplane decision boundaries, \mathbf{E} -product neurons generate non-linear decision surfaces that naturally tessellate the input space.

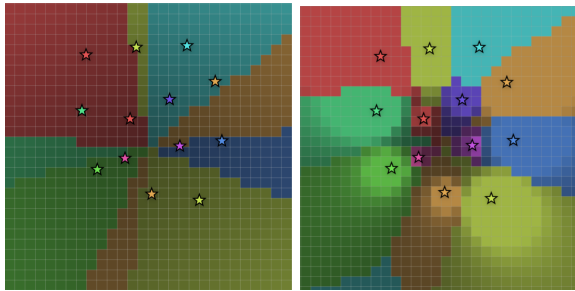


Figure 2: Comparison of decision boundaries in a 2D feature space: the conventional linear model (left) induces unbounded, half-space partitions, while our \mathbf{E} -product method (right) forms localized, vortex-like territories around learned prototypes. Stars denote the learned neuron prototypes in the 2D vector space, and regions are assigned by the argmax over the softmax outputs.

Key properties of the vortex dynamics include bounded attraction fields, where each neuron creates a localized influence region with inverse-square distance decay; non-linear decision boundaries,

282 in which curved algebraic surfaces replace linear hyperplanes; competitive tessellation, whereby
 283 neurons naturally develop specialized, non-overlapping territories; and an orthogonality-entropy
 284 connection, where geometric orthogonality corresponds to infinite cross-entropy and prevents rep-
 285 resentational collapse.

286 This vortex phenomenon enables natural space partitioning where data points are attracted to
 287 their most geometrically compatible prototype. The geometric partitioning is empirically vali-
 288 dated through the sharper MNIST prototypes shown in Figure 6, where \mathbf{E} -product neurons achieve
 289 superior class separation compared to conventional linear models. Quantitative analysis of proto-
 290 type clarity and territorial boundaries requires systematic measurement of representational overlap
 291 metrics—an important direction for future interpretability research (detailed mathematical analysis
 292 in Appendix E).

293 3.2 MNIST REPRESENTATION QUALITY

294 MNIST experiments demonstrate the \mathbf{E} -product’s bounded prototype evolution versus conventional
 295 unbounded growth. Figure 6 shows that conventional linear models produce diffuse, blurry proto-
 296 types, while \mathbf{E} -product neurons learn sharp, geometrically coherent digit representations with
 297 localized concentration and class-specific territorial structure (detailed analysis in Appendix E.8).

300 **Superposition and Prototype Inversion:** The \mathbf{E} -product neuron exhibits superposition behav-
 301 ior when prototypes are inverted ($\mathbf{w} \rightarrow -\mathbf{w}$). Unlike conventional neurons where sign flipping
 302 causes complete classification failure, \mathbf{E} -product neurons maintain reasonable performance due to
 303 their squared numerator structure, enabling two valid solutions without retraining. Specifically,
 304 dot product neurons achieve 91.88% accuracy with original prototypes but drop to approximately
 305 0.01% after inversion, while \mathbf{E} -product neurons maintain 92.18% and 87.87% accuracy respectively,
 306 demonstrating remarkable robustness to prototype sign changes.

307 3.3 VISION MODEL PERFORMANCE

308 We evaluate the \mathbf{E} -product’s effectiveness in computer vision by comparing standard architectures
 309 with their Aether variants across multiple datasets. All models are trained from scratch to en-
 310 sure fair comparison, using identical training protocols and hyperparameters except for the core
 311 computational unit replacement.

312 Table 2 presents comprehensive results across five standard computer vision datasets: CIFAR-10,
 313 CIFAR-100, STL-10, Tiny-ImageNet, and ImageNet-1K. We compare ResNet-18, ResNet-50, and
 314 Vision Transformer (ViT-Small) architectures with their Aether counterparts, where standard linear
 315 layers and attention mechanisms are replaced with \mathbf{E} -product units.

316 The results demonstrate competitive performance with or surpassing baselines on multiple bench-
 317 marks.
 318

319 3.4 AETHER-GPT2: LANGUAGE MODELING PERFORMANCE

320 To demonstrate the versatility of our approach beyond vision tasks, we implement Aether-GPT2,
 321 incorporating the \mathbf{E} -product architecture into the GPT2 framework for language modeling. We
 322 compare the perplexity scores between our Aether-GPT2 and the standard GPT2 architecture
 323 across multiple text corpora. On 2.5B tokens of Fineweb, Aether-GPT2 achieves a final validation
 324 loss of **2.29** in full precision (FP32) and **2.69** in mixed-precision (BF16), compared to 2.43 and
 325 3.03 for the standard GPT2 baseline. This represents an **11.2% relative improvement** in the
 326 mixed-precision setting. Furthermore, the elimination of normalization layers results in a **15-25%**
 327
 328

Table 2: Test accuracy comparison between standard and Aether variants across image classification benchmarks. All models trained from scratch with identical protocols.

Architecture	CIFAR-10	CIFAR-100	STL-10	Tiny-ImageNet	ImageNet-1K
ResNet-18	94.23%	72.15%	78.42%	56.89%	–
Aether-ResNet-18	92.37%	74.83%	80.91%	59.34%	–
ResNet-50	–	–	–	–	74.13%
Aether-ResNet-50	–	–	–	–	75.24%
ViT-Small	91.78%	69.91%	75.13%	52.76%	–
Aether-ViT-Small	92.45%	70.58%	78.89%	51.42%	–

reduction in peak memory usage. These findings establish Aether-GPT2 as a successful proof-of-concept, suggesting that the **E**-product can serve as a viable alternative to conventional neural network components (detailed experimental configuration in Appendix E.9).

Throughput and training efficiency. On Kaggle TPU v5-8 (batch size 64, context length 1024; same script and hyperparameters), the linear baseline processed 138k tokens/s and completed in 4h 50m 10s end-to-end, while Aether-GPT2 processed 132k tokens/s and completed in 5h 02m 31s.

4 RELATED WORK

4.1 INVERSE-SQUARE LAWS

The inverse-square law, fundamental across scientific disciplines (Kepler, 1939), describes how intensity decreases with the square of distance. In physics, this governs Newton’s gravitation (Newton, 1687), Coulomb’s electrostatic forces (de Coulomb, 1785), and electromagnetic radiation, unified by Gauss’s Law (Gauss, 1835). Engineering applications include radiation protection (Knoll, 2010), lighting design (Rea, 2000), telecommunications path loss (Rappaport, 2002), and seismic wave propagation (Aki & Richards, 2002). Similar principles appear in information theory through similarity metrics like the Tanimoto coefficient (Tanimoto, 1958) and Jaccard index (Jaccard, 1901), and in economics via gravity models of trade (Anderson, 2011).

4.2 ALTERNATIVE NEURAL OPERATORS AND ACTIVATION-FREE ARCHITECTURES

Several approaches have explored alternatives to the standard linear-then-activate paradigm. Quadratic neurons (Fan et al., 2020; Liao et al., 2024) replace dot products with quadratic forms, enabling non-linear decision boundaries without explicit activation functions. However, these methods focus solely on increasing polynomial degree without considering geometric relationships between vectors.

Multiplicative interactions (Jayakumar et al., 2020) and gated linear units (Dauphin et al., 2016) introduce element-wise products but maintain dependence on activation functions. SIREN networks (Sitzmann et al., 2020) use sinusoidal activations for implicit neural representations, while Fourier feature networks (Tancik et al., 2020) map inputs to high-dimensional Fourier bases before applying standard activations.

376 Unlike these approaches, the \mathbf{E} -product integrates both alignment (through squared dot products)
377 and proximity (through inverse distance) without requiring separate activation functions, providing
378 intrinsic non-linearity through geometric structure rather than functional composition.
379

380 4.3 KERNELIZED NEURAL NETWORKS AND DISTANCE-BASED METHODS 381

382 Kernel methods enable non-linear learning through implicit high-dimensional mappings. SVMs
383 (Cortes, 1995) established the foundation, formalized by Schölkopf (Schölkopf et al., 1997). Key
384 developments include Kernel PCA (Schölkopf et al., 1998), Gaussian Processes (Williams & Ras-
385 mussen, 2006), and Spectral Clustering (Ng et al., 2001). Scalability improvements came through
386 the Nyström method (Williams & Seeger, 2000) and Random Fourier Features (Rahimi & Recht,
387 2007).

388 The Neural Tangent Kernel (Jacot et al., 2018) bridges kernel methods and deep learning by ana-
389 lyzing infinite-width neural networks. However, NTK theory applies to conventional architectures
390 with explicit activations, whereas our approach eliminates activations entirely through geometric
391 operators.

392 Distance-based kernels like RBF (Boser et al., 1992) emphasize proximity for local structure cap-
393 ture, while polynomial kernels focus on feature interactions. The \mathbf{E} -product uniquely combines
394 both perspectives: the squared numerator captures polynomial-like alignment interactions, while
395 the inverse-square denominator provides RBF-like distance sensitivity with self-regularization prop-
396 erties absent in standard kernels.

397 Deep kernel learning (Wilson et al., 2016) and deep kernel processes (Aitchison et al., 2021) combine
398 the representational power of deep neural networks with the non-parametric flexibility of kernel
399 methods. Recent theoretical work has further generalized this connection (Yang & Aitchison,
400 2023). However, most deep kernel methods operate in the dual form, requiring the computation
401 and storage of a Gram matrix, which scales quadratically with the number of data points. In
402 contrast, our approach utilizes the primal form, allowing for direct computation in the feature
403 space without the prohibitive memory cost associated with large Gram matrices.

404 Recent work on kernelized neural networks (Cho & Saul, 2009a; Mairal et al., 2014) approximates
405 kernel computations within neural architectures but maintains the linear-then-activate structure.
406 In contrast, the \mathbf{E} -product serves as both the computational primitive and the kernel, eliminating
407 architectural complexity while preserving the benefits of kernel methods.
408

409 5 CONCLUSION 410

411 This work introduces the \mathbf{E} -product, a physics-inspired neural operator that unifies alignment and
412 proximity in a single computation, challenging the conventional paradigm that separates linear
413 transformations from activation functions. Drawing inspiration from inverse-square law interactions
414 in physics, the \mathbf{E} -product provides inherent non-linearity and geometric sensitivity, enabling more
415 nuanced understanding of vector interactions while simplifying neural network architectures.
416

417 Our theoretical analysis establishes that Neural-Matter Networks maintain universal approximation
418 capabilities while offering superior geometric fidelity and interpretability. Empirical validation
419 across diverse domains—from computer vision to language modeling and physics-informed neural
420 networks—demonstrates consistent performance improvements alongside memory efficiency gains.

421 The \mathbf{E} -product opens several promising research directions: scaling to large-scale architectures re-
422 quires systematic investigation of computational trade-offs and optimization dynamics; the geomet-

423 ric interpretability framework enables principled analysis of learned representations and decision
424 boundaries; and the connection to physical laws suggests broader applications in scientific machine
425 learning where spatial relationships are fundamental.

426 By eliminating the information bottleneck inherent in traditional activation functions, this approach
427 paves the way toward geometrically-grounded neural architectures that unite computational effi-
428 ciency with theoretical understanding, offering a constructive path toward more interpretable and
429 robust neural computation.

431 LICENSE

432 This work is licensed under the Affero GNU General Public License (AGPL) v3.0. The AGPL is
433 a free software license that ensures end users have the freedom to run, study, share, and modify
434 the software. It requires that any modified versions of the software also be distributed under the
435 same license, ensuring that the freedoms granted by the original license are preserved in derivative
436 works. The full text of the AGPL v3.0 can be found at <https://www.gnu.org/licenses/agpl-3.0.en.html>. By using this work, you agree to comply with the terms of the AGPL v3.0.

441 REFERENCES

442 Lectures on radial basis functions, 1991.

443 Laurence Aitchison, Adam Yang, and Sebastian W Ober. Deep kernel processes. In *International
444 Conference on Machine Learning*, pp. 130–140. PMLR, 2021.

445 Keiiti Aki and Paul G. Richards. *Quantitative Seismology*. University Science Books, Sausalito,
446 CA, 2 edition, 2002.

447 James E. Anderson. The gravity model. *Annual Review of Economics*, 3(1):133–160, 2011.

448 Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint
449 arXiv:1607.06450*, 2016.

450 Sergei N Bernstein. Sur les fonctions absolument monotones. *Acta Mathematica*, 52:1–66, 1928.

451 Salomon Bochner. *Vorlesungen über Fouriersche Integrale*. Akademische Verlagsgesellschaft,
452 Leipzig, 1932.

453 Bernhard E. Boser, Isabelle M. Guyon, and Vladimir N. Vapnik. A training algorithm for opti-
454 mal margin classifiers. In *Proceedings of the Fifth Annual Workshop on Computational Learning
455 Theory, COLT '92*, pp. 144–152, New York, NY, USA, 1992. Association for Computing Machin-
456 ery. ISBN 089791497X. doi: 10.1145/130385.130401. URL <https://doi.org/10.1145/130385.130401>.

457 Martin D Buhmann. Radial basis functions. *Acta numerica*, 9:1–38, 2000.

458 Youngmin Cho and Lawrence Saul. Kernel methods for deep learning. In Y. Bengio, D. Schuurmans,
459 J. Lafferty, C. Williams, and A. Culotta (eds.), *Advances in Neural Information Processing
460 Systems*, volume 22. Curran Associates, Inc., 2009a. URL [https://proceedings.neurips.cc/
461 paper_files/paper/2009/file/5751ec3e9a4feab575962e78e006250d-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2009/file/5751ec3e9a4feab575962e78e006250d-Paper.pdf).

462 Youngmin Cho and Lawrence K. Saul. Kernel methods for deep learning, 2009b.

- 470 Corinna Cortes. Support-vector networks. *Machine Learning*, 1995.
- 471
- 472 Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2006.
- 473
- 474 George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control,*
475 *Signals and Systems*, 2(4):303–314, 1989.
- 476 Yann N. Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated
477 convolutional networks. *CoRR*, abs/1612.08083, 2016. URL [http://arxiv.org/abs/1612.](http://arxiv.org/abs/1612.08083)
478 08083.
- 479 Charles-Augustin de Coulomb. Premier mémoire sur l’électricité et le magnétisme. *Histoire de*
480 *l’Académie Royale des Sciences*, pp. 1–31, 1785. in French.
- 481
- 482 Andrew Draganov, Sharvaree Vadgama, and Erik J Bekkers. The hidden pitfalls of the cosine
483 similarity loss. *arXiv preprint arXiv:2406.16468*, 2024.
- 484
- 485 Fenglei Fan, Jinjun Xiong, and Ge Wang. Universal approximation with quadratic deep net-
486 works. *Neural Networks*, 124:383–392, 2020. ISSN 0893-6080. doi: [https://doi.org/10.](https://doi.org/10.1016/j.neunet.2020.01.007)
487 1016/j.neunet.2020.01.007. URL [https://www.sciencedirect.com/science/article/pii/](https://www.sciencedirect.com/science/article/pii/S0893608020300095)
488 S0893608020300095.
- 489 Gregory E Fasshauer. *Positive definite kernels: past, present and future*. World Scientific, 2011.
- 490
- 491 Guido Fubini. Sugli integrali multipli. *Rendiconti del Reale Accademia dei Lincei*, 16(1):608–614,
492 1907.
- 493 Carl Friedrich Gauss. *Allgemeine Lehrsätze in Beziehung auf die im verkehrten Verhältniss des*
494 *Quadrats der Entfernung wirkenden Anziehungs- und Abstossungskräfte*. Dietrich, Göttingen,
495 1835.
- 496
- 497 Josiah Willard Gibbs. *Elementary Principles in Statistical Mechanics*. Charles Scribner’s Sons,
498 New York, 1902.
- 499
- 500 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*, volume 1. MIT Press, 2016.
- 501 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus), 2023. URL [https://arxiv.](https://arxiv.org/abs/1606.08415)
502 [org/abs/1606.08415](https://arxiv.org/abs/1606.08415).
- 503
- 504 Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. Kernel methods in machine learn-
505 ing. 2008.
- 506
- 507 Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2012.
- 508
- 509 Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are
510 universal approximators. *Neural networks*, 2(5):359–366, 1989.
- 511
- 512 Changcun Huang. Relu networks are universal approximators via piecewise linear or constant
513 functions. *Neural Computation*, 32(11):2249–2278, 2020.
- 514
- 515 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
516 reducing internal covariate shift, 2015. URL <https://arxiv.org/abs/1502.03167>.
- 517
- 518 Paul Jaccard. Étude comparative de la distribution florale dans une portion des alpes et des jura.
519 *Bulletin de la Société Vaudoise des Sciences Naturelles*, 37:547–579, 1901.

- 517 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
518 generalization in neural networks. *Advances in neural information processing systems*, 31, 2018.
519
- 520 Siddhant M. Jayakumar, Wojciech M. Czarnecki, Jacob Menick, Jonathan Schwarz, Jack Rae,
521 Simon Osindero, Yee Whye Teh, Tim Harley, and Razvan Pascanu. Multiplicative interactions
522 and where to find them. 2020.
- 523 Johannes Kepler. Ad vitellionem paralipomena, quibus astronomiae pars optica traditur. 1604.
524 *Johannes Kepler: Gesammelte Werke, Ed. Walther von Dyck and Max Caspar, München*, 1939.
525
- 526 Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing
527 neural networks. *Advances in neural information processing systems*, 30, 2017.
- 528 Glenn F. Knoll. *Radiation Detection and Measurement*. John Wiley & Sons, Hoboken, NJ, 4
529 edition, 2010.
- 530
- 531 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,
532 2015.
- 533 Jing-Xiao Liao, Bo-Jian Hou, Hang-Cheng Dong, Hao Zhang, Xiaoge Zhang, Jinwei Sun, Shiping
534 Zhang, and Feng-Lei Fan. Quadratic neuron-empowered heterogeneous autoencoder for unsu-
535 pervised anomaly detection. *IEEE Transactions on Artificial Intelligence*, 5(9):4723–4737, 2024.
536 doi: 10.1109/TAI.2024.3394795.
- 537
- 538 Zhou Lu, Hongming Pu, Feicheng Wang, Zhiqiang Hu, and Liwei Wang. The expressive power of
539 neural networks: A view from the width. *Advances in neural information processing systems*, 30,
540 2017.
- 541 Julien Mairal, Piotr Koniusz, Zaid Harchaoui, and Cordelia Schmid. Convolutional kernel networks,
542 2014. URL <https://arxiv.org/abs/1406.3332>.
- 543
- 544 J. Mercer. Functions of positive and negative type, and their connection with the theory of integral
545 equations. *Philosophical Transactions of the Royal Society of London. Series A, Containing*
546 *Papers of a Mathematical or Physical Character*, 209:415–446, 1909.
- 547 Charles A Micchelli, Yuesheng Xu, and Haizhang Zhang. Universal kernels. *Journal of Machine*
548 *Learning Research*, 7(12), 2006.
- 549
- 550 Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines.
551 In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814,
552 2010.
- 553 Isaac Newton. *Philosophiæ Naturalis Principia Mathematica*. S. Pepys, London, 1687.
- 554
- 555 Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm.
556 *Advances in neural information processing systems*, 14, 2001.
- 557
- 558 Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. *Advances in*
559 *neural information processing systems*, 20, 2007.
- 560
- 561 Theodore S. Rappaport. *Wireless Communications: Principles and Practice*. Prentice Hall, Upper
562 Saddle River, NJ, 2 edition, 2002.
- 563
- 562 Mark S. Rea. *The IESNA Lighting Handbook: Reference & Application*. Illuminating Engineering
563 Society of North America, New York, 9 edition, 2000.

- 564 Tim Salimans and Durk P Kingma. Weight normalization: A simple reparameterization to accel-
565 erate training of deep neural networks. In *Advances in neural information processing systems*.
566 International Business Machines Corp., 1958.
- 567 Bernhard Schölkopf and Alexander J Smola. *Learning with kernels: support vector machines,*
568 *regularization, optimization, and beyond*. MIT press, 2002.
- 570 Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component
571 analysis. In *International conference on artificial neural networks*, pp. 583–588. Springer, 1997.
- 572 Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Nonlinear component analysis as
573 a kernel eigenvalue problem. *Neural computation*, 10(5):1299–1319, 1998.
- 574 Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Im-
575 plicit neural representations with periodic activation functions. volume 33, pp. 7462–7473, 2020.
- 577 Harald Steck, Chaitanya Ekanadham, and Nathan Kallus. Is cosine-similarity of embeddings really
578 about similarity? In *Companion Proceedings of the ACM Web Conference 2024*, pp. 887–890,
579 2024.
- 580 Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh
581 Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn
582 high frequency functions in low dimensional domains. volume 33, pp. 7537–7547, 2020.
- 583 Taffee T Tanimoto. An elementary mathematical theory of classification and prediction. 1958.
- 585 Leonida Tonelli. Sull’integrazione per parti. *Atti della Accademia Nazionale dei Lincei. Rendiconti*
586 *Lincei. Classe di Scienze Fisiche, Matematiche e Naturali*, 18(2):246–253, 1909.
- 587 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
588 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
589 *processing systems*, 30, 2017.
- 591 Christopher Williams and Matthias Seeger. Using the nyström method to speed up kernel machines.
592 *Advances in neural information processing systems*, 13, 2000.
- 593 Christopher KI Williams and Carl Edward Rasmussen. *Gaussian processes for machine learning*,
594 volume 2. MIT press Cambridge, MA, 2006.
- 595 Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. Deep kernel learning.
596 In *Artificial intelligence and statistics*, pp. 370–378. PMLR, 2016.
- 598 Adam X Yang and Laurence Aitchison. A theory of representation learning gives a deep general-
599 isation of kernel methods. In *International Conference on Machine Learning*, pp. 39163–39189.
600 PMLR, 2023.

601 A APPENDIX

602 B THEORETICAL BACKGROUND

603 B.1 REVISITING CORE COMPUTATIONAL PRIMITIVES AND SIMILARITY MEASURES

604 The computational primitives used in deep learning are fundamental to how models represent
605 and process information. This section revisits key mathematical operations and similarity mea-
606 sures, such as the dot product, convolution, cosine similarity, and Euclidean distance, that form
607
608
609
610

611 the bedrock of many neural architectures. We will explore their individual properties and how
 612 they contribute to tasks like feature alignment, localized feature mapping, and quantifying spatial
 613 proximity. Furthermore, we will delve into the role of neural activation functions in enabling the
 614 non-linear transformations crucial for complex pattern recognition. Understanding these core con-
 615 cepts and their inherent characteristics is crucial for appreciating the motivation behind developing
 616 novel operators, as explored in this work, that aim to capture more nuanced relationships within
 617 data Goodfellow et al. (2016).

618 B.1.1 THE DOT PRODUCT: A MEASURE OF ALIGNMENT

620 The dot product, or scalar product, remains a cornerstone of neural computation, serving as the
 621 primary mechanism for quantifying the interaction between vectors, such as a neuron’s weights and
 622 its input. For two vectors $\mathbf{a} = [a_1, a_2, \dots, a_n]$ and $\mathbf{b} = [b_1, b_2, \dots, b_n]$, it is defined as:

$$623 \mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n \quad (4)$$

626 Geometrically, the dot product is proportional to the cosine of the angle between the vectors and
 627 their Euclidean magnitudes: $\mathbf{a} \cdot \mathbf{b} = \|\mathbf{a}\| \|\mathbf{b}\| \cos(\theta)$. Its sign indicates the general orientation (acute,
 628 obtuse, or orthogonal angle), and its magnitude reflects the degree of alignment scaled by vector
 629 lengths. In machine learning, dot product scores are pervasively used to infer similarity, relevance,
 630 or the strength of activation. However, as noted in Section 1, its conflation of magnitude and
 631 directional alignment can sometimes obscure more fine-grained geometric relationships, motivating
 632 the exploration of operators that offer a more comprehensive assessment of vector interactions.

633 B.1.2 THE CONVOLUTION OPERATOR: LOCALIZED FEATURE MAPPING

635 The convolution operator is pivotal in processing structured data, particularly in Convolutional
 636 Neural Networks (CNNs). It applies a kernel (or filter) across an input to produce a feature map,
 637 effectively an operation on two functions, f (input) and g (kernel), yielding a third that expresses
 638 how one modifies the shape of the other. For discrete signals, such as image patches and kernels, it
 639 is:

$$640 (f * g)[n] = \sum_{m=-\infty}^{\infty} f[m]g[n - m] \quad (5)$$

642 In CNNs, convolution performs several critical roles:

- 644 • **Feature Detection:** Kernels learn to identify localized patterns (edges, textures, motifs)
 645 at various abstraction levels.
- 646 • **Spatial Hierarchy:** Stacking layers allows the model to build complex feature represen-
 647 tations from simpler ones.
- 648 • **Parameter Sharing:** Applying the same kernel across spatial locations enhances efficiency
 649 and translation equivariance.

650 The core computation within a discrete convolution at a specific location involves an element-wise
 651 product sum between the kernel and the corresponding input patch, which is, in essence, a dot
 652 product. Consequently, the resulting activation at each point in the feature map reflects the local
 653 alignment between the input region and the kernel. If an input patch and a kernel are orthogonal
 654 (i.e., their element-wise product sums to zero, akin to a zero dot product if they were vectorized), the
 655 convolution output at that position will be zero, indicating no local match for the feature encoded
 656 by the kernel. This reliance on dot product-like computations means that standard convolutions
 657 primarily assess feature alignment, potentially overlooking other geometric aspects of the data.

B.1.3 COSINE SIMILARITY: NORMALIZING FOR DIRECTIONAL AGREEMENT

Cosine similarity refines the notion of alignment by isolating the directional aspect of vector relationships, abstracting away from their magnitudes. It measures the cosine of the angle between two non-zero vectors \mathbf{A} and \mathbf{B} :

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad (6)$$

Scores range from -1 (perfectly opposite) to 1 (perfectly aligned), with 0 signifying orthogonality (decorrelation). By normalizing for vector lengths, cosine similarity provides a pure measure of orientation. This is particularly useful when the magnitude of vectors is not indicative of their semantic relationship, such as in document similarity tasks. While it effectively captures directional agreement, it explicitly discards information about vector magnitudes and, like the dot product, does not inherently account for the spatial proximity between the vectors themselves if they are points in a space (Draganov et al., 2024; Steck et al., 2024).

B.1.4 EUCLIDEAN DISTANCE: QUANTIFYING SPATIAL PROXIMITY

In contrast to measures of alignment, Euclidean distance quantifies the "ordinary" straight-line separation between two points (or vectors) $\mathbf{p} = (p_1, \dots, p_n)$ and $\mathbf{q} = (q_1, \dots, q_n)$ in an n-dimensional Euclidean space:

$$d(\mathbf{p}, \mathbf{q}) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2} \quad (7)$$

This metric is fundamental in various machine learning algorithms, including k-Nearest Neighbors and k-Means clustering, and forms the basis of loss functions like Mean Squared Error. Euclidean distance measures dissimilarity based on spatial proximity; a smaller distance implies greater similarity in terms of location within the vector space. Unlike cosine similarity, it is sensitive to vector magnitudes and their absolute positions. However, Euclidean distance alone does not directly convey information about the relative orientation or alignment of vectors, only their nearness.

The distinct characteristics of these foundational measures, alignment (dot product, cosine similarity) versus proximity (Euclidean distance), highlight an opportunity. These foundational measures force a choice: one can measure alignment (dot product, cosine similarity) or spatial proximity (Euclidean distance), but no single, primitive operator in conventional use effectively unifies both. Neural operators that can synergistically combine these aspects, assessing not only if vectors point in similar directions but also if they are close in the embedding space, could offer a richer, more geometrically informed way to model interactions. This perspective underpins the development of the \mathbf{E} -product introduced in Section 2.

B.2 THE ROLE AND GEOMETRIC COST OF NON-LINEAR ACTIVATION

While the core computational primitives provide tools to measure similarity and interaction, their inherent linearity limits the complexity of functions they can represent. To overcome this, deep neural networks employ non-linear activation functions. These are the standard method for introducing non-linearity, a necessary step for modeling intricate data patterns. However, this "fix" is imperfect, as it introduces its own set of problems, particularly concerning the preservation of the input data's geometric integrity. The remarkable expressive power of deep neural networks hinges on their capacity to model complex, non-linear relationships. This ability to approximate any continuous function to an arbitrary degree of accuracy is formally captured by the universal approximation theorem Cybenko (1989); Hornik et al. (1989); Lu et al. (2017); Huang (2020).

This theorem underscores the critical role of non-linear activation functions. Without such non-linearities, a deep stack of layers would mathematically collapse into an equivalent single linear transformation, severely curtailing its representational capacity. Activation functions are thus not mere auxiliaries; they are the pivotal components that unlock the hierarchical and non-linear feature learning central to deep learning’s success. They determine a neuron’s output based on its aggregated input, and in doing so, introduce crucial selectivity: enabling the network to preferentially respond to certain patterns while attenuating or ignoring others.

B.2.1 LINEAR SEPARABILITY AND THE LIMITATIONS OF THE INNER PRODUCT

The fundamental computation within a single artificial neuron (perceptron) is an affine transformation followed by a non-linear activation function σ :

$$y = \sigma(\langle \mathbf{w}, \mathbf{x} \rangle + b), \quad (8)$$

where \mathbf{w} is the weight vector, \mathbf{x} is the input vector, and b is the bias term. The decision boundary of this neuron is implicitly defined by the hyperplane where the argument to σ is zero:

$$\{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}. \quad (9)$$

This hyperplane partitions the input space \mathbb{R}^d into two half-spaces. Consequently, a single neuron can only implement linearly separable functions. This is a direct consequence of the linear nature of the inner product, which can only define a linear decision boundary. While this allows for efficient computation, it severely restricts the complexity of functions that can be learned.

A classic counterexample is the XOR function, whose truth table cannot be satisfied by any single linear decision boundary. Specifically, for inputs $\mathbf{x} \in \{(0, 0), (0, 1), (1, 0), (1, 1)\} \subset \mathbb{R}^2$, there exist no $\mathbf{w} \in \mathbb{R}^2$ and $b \in \mathbb{R}$ such that $\text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle + b)$ matches the XOR output (0, 1, 1, 0 respectively). This limitation stems directly from the linear nature of the inner product operation defining the separating boundary Goodfellow et al. (2016).

B.2.2 NON-LINEAR FEATURE SPACE TRANSFORMATION VIA HIDDEN LAYERS AND ITS GEOMETRIC COST

Multi-layer perceptrons (MLPs) overcome this limitation by cascading transformations. A hidden layer maps the input \mathbf{x} to a new representation \mathbf{h} through a matrix-vector product and an element-wise activation function ϕ :

$$\mathbf{h} = \phi(W\mathbf{x} + \mathbf{b}). \quad (10)$$

Here, $W \in \mathbb{R}^{m \times d}$ is the weight matrix, $\mathbf{b} \in \mathbb{R}^m$ is the bias vector, and m is the number of hidden neurons. Each row \mathbf{w}_i^\top of W corresponds to the weight vector of the i -th hidden neuron, computing $h_i = \phi(\langle \mathbf{w}_i, \mathbf{x} \rangle + b_i)$. This transforms the input space \mathbb{R}^d into a feature space \mathbb{R}^m . The introduction of the non-linear activation function ϕ is what allows the network to learn non-linear decision boundaries. However, this gain in expressive power comes at a cost: the potential loss of geometric fidelity.

B.2.3 TOPOLOGICAL DISTORTIONS AND INFORMATION LOSS VIA ACTIVATION FUNCTIONS

While hidden layers using the transformation $\mathbf{h} = \phi(W\mathbf{x} + \mathbf{b})$ enable the learning of non-linear functions, the introduction of the element-wise non-linear activation function ϕ , often crucial for breaking linearity, can significantly alter the topological and geometric structure of the data representation, potentially leading to information loss Goodfellow et al. (2016). This is a critical trade-off: gaining non-linear modeling capability while potentially discarding valuable geometric information.

752 Consider the mapping $T : \mathbb{R}^d \rightarrow \mathbb{R}^m$ defined by $T(\mathbf{x}) = \phi(W\mathbf{x} + \mathbf{b})$. The affine part, $A(\mathbf{x}) = W\mathbf{x} + \mathbf{b}$,
 753 performs a linear transformation (rotation, scaling, shear, projection) followed by a translation.
 754 While this affine map distorts metric properties (distances and angles, unless W is proportional to
 755 an orthogonal matrix), it preserves basic topological features like connectedness and maps lines to
 756 lines (or points) Goodfellow et al. (2016).

757 However, the subsequent application of a typical non-linear activation ϕ element-wise often leads
 758 to more drastic topological changes:
 759

- 760 1. Non-Injectivity and Collapsing Regions: Many common activation functions render the
 761 overall mapping T non-injective.
 762
 - 763 • ReLU ($\phi(z) = \max(0, z)$): Perhaps the most prominent example. For each hidden
 764 neuron i , the entire half-space defined by $\{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{w}_i, \mathbf{x} \rangle + b_i \leq 0\}$ is mapped to
 765 $h_i = 0$. Distinct points $\mathbf{x}_1, \mathbf{x}_2$ within this region, potentially far apart, become indistin-
 766 guishable along the i -th dimension of the hidden space. This constitutes a significant
 767 loss of information about the relative arrangement of data points within these col-
 768 lapsed regions. The mapping is fundamentally many-to-one. For instance, consider
 769 two input vectors that are anti-aligned with a neuron’s weight vector to different de-
 770 grees, one strongly and one weakly. A ReLU activation function would map both
 771 resulting negative dot products to zero, rendering their distinct geometric opposition
 772 indistinguishable to subsequent layers. This information is irretrievably discarded.
 - 773 • Sigmoid/Tanh: While smooth, these functions saturate. Inputs $\mathbf{z}_1 = A(\mathbf{x}_1)$ and $\mathbf{z}_2 =$
 774 $A(\mathbf{x}_2)$ that are far apart but both fall into the saturation regime (e.g., large positive or
 775 large negative values) will map to $\mathbf{h}_1 \approx \mathbf{h}_2$. This ‘squashing’ effect can merge distinct
 776 clusters from the input space if they map to saturated regions in the hidden space,
 777 again losing discriminative information and distorting the metric structure.
- 778 2. Distortion of Neighborhoods: The relative distances between points can be severely dis-
 779 torted. Points close in the input space \mathbb{R}^d might be mapped far apart in \mathbb{R}^m , or vice-versa
 780 (especially due to saturation or the zero-region of ReLU). This means the local neigh-
 781 borhood structure is not faithfully preserved. Formally, the mapping T is generally not
 782 a homeomorphism onto its image, nor is it typically bi-Lipschitz (which would provide
 783 control over distance distortions).

784 While these distortions are precisely what grant neural networks their expressive power to warp
 785 the feature space and create complex decision boundaries, they come at the cost of potentially
 786 discarding information present in the original geometric configuration of the data. The network
 787 learns which information to preserve and which to discard based on the optimization objective,
 788 but the mechanism relies on potentially non-smooth or non-injective transformations introduced by
 789 ϕ . This highlights the conflation of magnitude and direction in the dot product, the information
 790 loss from activation functions, and the lack of a unified measure for proximity and alignment,
 791 setting the stage for the \mathbf{E} -product. Formal properties of the \mathbf{E} -product are established later: it is a
 792 Mercer kernel (Theorem 2.1), yields universal approximation in NMNs (Theorem 2.4), and exhibits
 793 self-regulation and stable gradients (Propositions C.6 and C.9).

794 B.3 DESIGN PHILOSOPHY: INTRINSIC NON-LINEARITY AND SELF-REGULATION

795
 796 A central hypothesis underpinning our methodological choices is that the \mathbf{E} -product (Section 1)
 797 possesses inherent non-linearity and self-regulating properties that can reduce or eliminate the
 798 need for conventional activation functions (e.g., ReLU, sigmoid, GeLU).

799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845

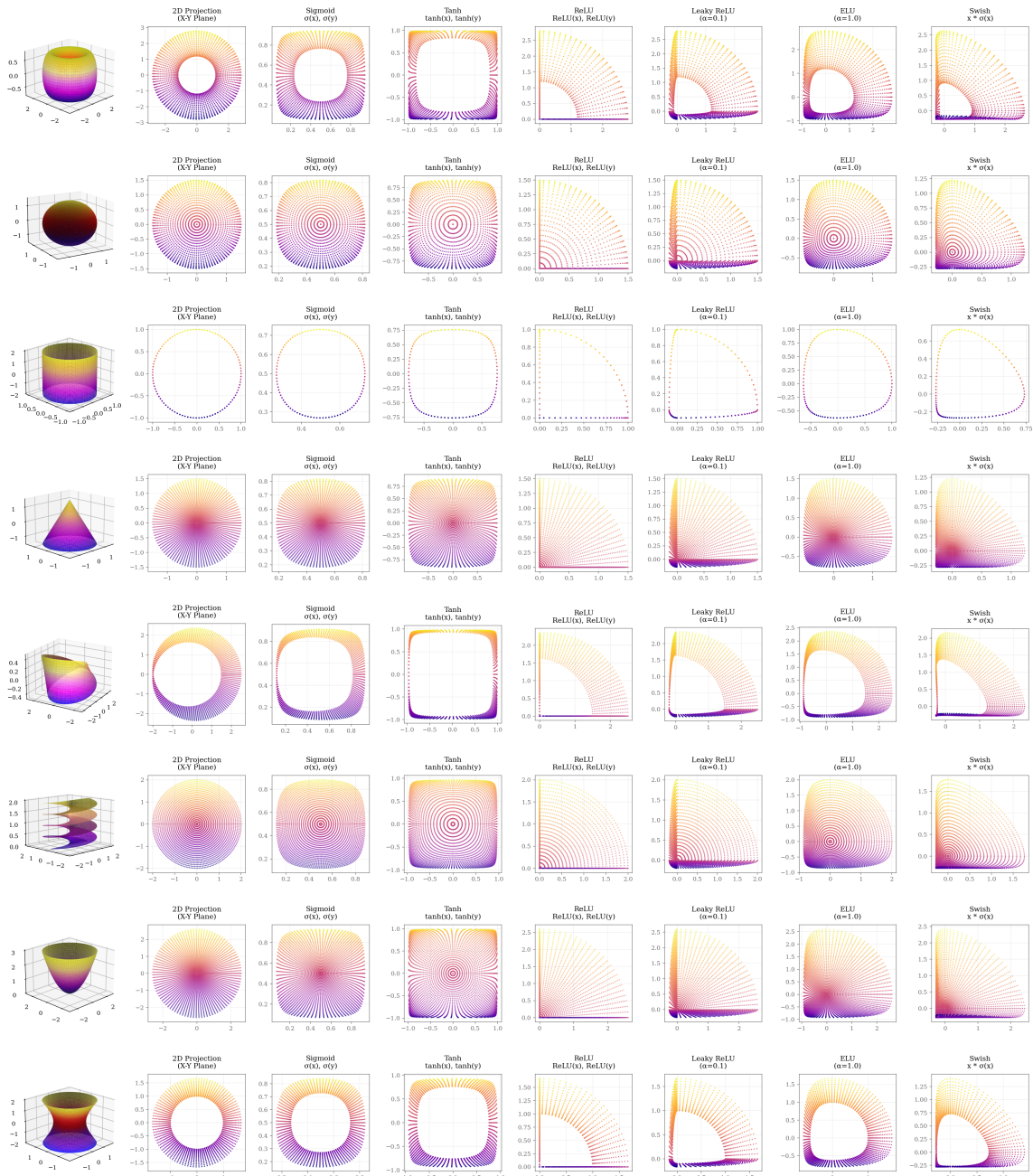


Figure 3: Illustration of how non-linear activation functions can distort the geometric structure of the input data manifold, leading to potential information loss. The original manifold (left) is transformed into a distorted representation after applying a non-linear activation functions.

846 This philosophy recontextualizes the fundamental components of neural computation. Neuron
 847 weights (\mathbf{w}) and input signals (\mathbf{x}) are not merely operands in a linear transformation followed by
 848 a non-linear activation; instead, they are conceptualized as co-equal vector entities inhabiting a
 849 shared, high-dimensional feature manifold. Within this framework, each vector can be viewed as
 850 an analogue to a fundamental particle or feature vector, with its constituent dimensions potentially
 851 encoding excitatory, inhibitory, or neutral characteristics relative to other entities in the space.
 852 The \mathbf{E} -product (Section 1) then transcends simple similarity assessment; it functions as a sophis-
 853 ticated interaction potential, $\mathcal{K}_{\mathbf{E}}(\mathbf{w}, \mathbf{x}) = \frac{(\mathbf{w}^T \mathbf{x})^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon}$, quantifying the 'field effects' between these
 854 vector entities. This interaction is reminiscent of n-body problems in physics. In machine learning,
 855 it draws parallels with, yet distinctively evolves from, learned metric spaces in contrastive learning,
 856 particularly those employing a triplet loss framework. While triplet loss aims to pull positive pairs
 857 closer and push negative pairs apart in the embedding space, our \mathbf{E} -product seeks a more nuanced
 858 relationship: 'positive' interactions (high \mathbf{E} -product value) require both strong alignment (high
 859 $(\mathbf{w}^T \mathbf{x})^2$) and close proximity (low $\|\mathbf{w} - \mathbf{x}\|^2$). Conversely, 'negative' or dissimilar relationships are
 860 not merely represented by distance, but more significantly by orthogonality (leading to a vanishing
 861 numerator), which signifies a form of linear independence and contributes to the system's capacity
 862 for true non-linear discrimination. Crucially, the non-linearity required for complex pattern recog-
 863 nition is not an external imposition (e.g., via a separate activation function) but is intrinsic to this
 864 interaction potential. The interplay between the squared dot product (alignment sensitivity) and
 865 the inverse squared Euclidean distance (proximity sensitivity) in its formulation directly sculpts a
 866 complex, non-linear response landscape without recourse to auxiliary functions.

867 Furthermore, this conceptualization of the \mathbf{E} -product as an intrinsic interaction potential suggests
 868 inherent self-regulating properties. The distance-sensitive denominator, $\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon$, acts as a
 869 natural dampening mechanism. As the 'distance' (dissimilarity in terms of position) between in-
 870 teracting vector entities \mathbf{w} and \mathbf{x} increases, the strength of their interaction, and thus the resultant
 871 activation, diminishes quadratically. This behavior is hypothesized to inherently curtail runaway
 872 activations and stabilize learning dynamics by ensuring that responses are localized and bounded.
 873 Such intrinsic stabilization contrasts sharply with conventional approaches that rely on explicit nor-
 874 malization layers (e.g., Batch Normalization, Layer Normalization) to manage activation statistics
 875 post-hoc. These layers, while effective, introduce additional computational overhead, can obscure
 876 direct input-output relationships, and sometimes complicate the theoretical analysis of network be-
 877 havior. The \mathbf{E} -product's formulation, therefore, offers a pathway to architectures where regulatory
 mechanisms are embedded within the primary computational fabric of the network.

878 The inherent non-linearity of the \mathbf{E} -product, coupled with the self-regulating properties suggested
 879 by its formulation (and formally proven in Appendix C.10), are central to our hypothesis that
 880 it can form the basis of powerful and robust neural architectures. These intrinsic characteristics
 881 open avenues for simplifying network design, potentially reducing reliance on or even eliminating
 882 conventional activation functions and normalization layers.

886 C SQUASHING FUNCTIONS FOR NON-NEGATIVE SCORES

889 The \mathbf{E} -product and its derivatives, such as the $\mathcal{K}_{\mathbf{E}}$ -kernel, naturally yield non-negative scores. In
 890 many machine learning contexts, particularly when these scores need to be interpreted as probabili-
 891 ties, attention weights, or simply normalized outputs, it is essential to apply a squashing function
 892 to map them to a desired range (e.g., $[0, 1]$ or ensuring a set of scores sum to 1).

893 C.1 CATEGORIZATION OF SQUASHING FUNCTIONS

894 Squashing functions for non-negative scores can be broadly categorized into two types:

- 895 • **Competitive (Vector-Normalizing) Functions:** These functions normalize a set of
- 896 scores collectively, producing a distribution over the vector. Each output depends on the
- 897 values of all dimensions, allowing for competitive interactions among them. This is useful for
- 898 attention mechanisms or probability assignments where the sum of outputs is meaningful.
- 899
- 900 • **Individualistic (Per-Dimension) Functions:** These functions squash each score in-
- 901 dependently, without reference to other values in the vector. Each output depends only
- 902 on its corresponding input, making them suitable for bounding or interpreting individual
- 903 activations.
- 904

905 C.2 LIMITATIONS OF TRADITIONAL SQUASHING FUNCTIONS

906 Traditional squashing functions, however, present challenges when applied to non-negative inputs:

- 907
- 908 • **Standard Sigmoid Function** ($\sigma(x) = \frac{1}{1+e^{-x}}$): When applied to non-negative inputs
- 909 ($x \geq 0$), the standard sigmoid function produces outputs in the range $[0.5, 1)$. The minimum
- 910 value of 0.5 for $x = 0$ renders it unsuitable for scenarios where small non-negative scores
- 911 should map to values close to 0.
- 912
- 913 • **Standard Softmax Function** ($\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}}$): The use of the exponential
- 914 function in softmax can lead to *hard* distributions, where one input value significantly
- 915 dominates the output, pushing other probabilities very close to zero. While this is often
- 916 desired for classification, it can be too aggressive if a softer assignment of probabilities or
- 917 attention is preferred. Additionally, softmax can suffer from numerical instability for large
- 918 input values due to the exponentials.
- 919

920 C.3 PROPOSED ALTERNATIVE SQUASHING FUNCTIONS

921 Given these limitations and the non-negative nature of \mathbf{E} -product scores, we consider alternative

922 squashing functions more suited to this domain:

- 923 • **softermax (Competitive):** This function normalizes a score x_k (optionally raised to a
- 924 power $n > 0$) relative to the sum of a set of non-negative scores $\{x_i\}$ (each raised to n),
- 925 with a small constant $\epsilon > 0$ for numerical stability. It is defined as:

$$926 \text{softermax}_n(x_k, \{x_i\}) = \frac{x_k^n}{\epsilon + \sum_i x_i^n} \quad (11)$$

927 Unlike softmax, softermax does not use exponentials, which avoids numerical instability for

928 large inputs and provides a more direct, interpretable translation of the underlying scores

929 into a normalized distribution. The power n controls the sharpness of the distribution:

930 $n = 1$ recovers the original Softermax, while $n > 1$ makes the distribution harder (more

931 peaked), and $0 < n < 1$ makes it softer.

- 932 • **soft-sigmoid (Individualistic):** This function squashes a single non-negative score $x \geq 0$
- 933 (optionally raised to a power $n > 0$) into the range $[0, 1)$. It is defined as:

$$934 \text{soft-sigmoid}_n(x) = \frac{x^n}{1 + x^n} \quad (12)$$

The power n modulates the softness: higher n makes the function approach zero faster for large x , while $n < 1$ makes the decay slower.

- **soft-tanh (Individualistic):** This function maps a non-negative score $x \geq 0$ (optionally raised to a power $n > 0$) to the range $[-1, 1)$ by linearly transforming the output of soft-sigmoid. It is defined as:

$$\text{soft-tanh}_n(x) = 2 \cdot (\text{soft-sigmoid}_n(x) - \frac{1}{2}) = \frac{x^n - 1}{1 + x^n} \quad (13)$$

The power n again controls the transition sharpness: higher n makes the function approach -1 more quickly for large x .

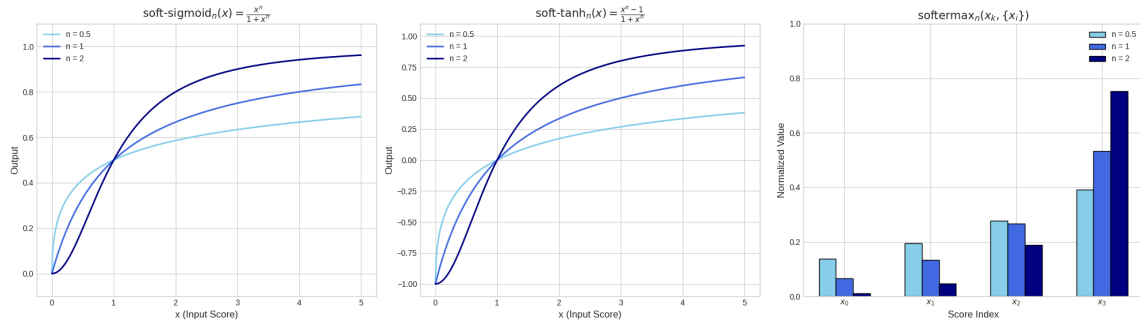


Figure 4: Visualization of the softermax, soft-sigmoid, and soft-tanh functions. These functions are designed to handle non-negative inputs from the \mathbf{E} -product and its derivatives, providing appropriate squashing mechanisms that maintain sensitivity across the range of non-negative inputs.

These functions are particularly well-suited for the outputs of \mathbf{E} -product-based computations, as they maintain sensitivity across the range of non-negative inputs while avoiding the pitfalls of standard activation functions (Nair & Hinton, 2010; Hendrycks & Gimpel, 2023; Klambauer et al., 2017).

C.4 FUNCTIONAL ROLES AND APPLICATIONS

The main role of these squashing functions can be categorized into two main categories:

- **Collective Communication and Space Splitting:** The softermax function allows for a comparative analysis of scores, reflecting their orthogonality and spatial proximity to an input vector. A higher score indicates that a vector is more aligned and closer to the input, while a lower score suggests greater orthogonality. This facilitates a competitive interaction where vectors vie for influence based on their geometric relationship with the input. The power parameter n , analogous to the temperature in softmax, controls the sharpness of the gravitational potential well’s slope.
- **Individual Score Squashing:** The soft-sigmoid and soft-tanh functions are used to squash individual non-negative scores into a bounded range, typically $[0, 1)$ for soft-sigmoid and $[-1, 1)$ for soft-tanh. They are particularly useful when the output needs to be interpreted as a probability or when a bounded response is required, as each score is processed independently of the others. The power parameter controls the steepness of the function, while the minimum value can be interpreted as an orthogonality score.

987 C.5 COMPUTATIONAL COMPLEXITY ANALYSIS
988

989 We provide exact forward/backward complexity, closed-form gradients, asymptotic characterization,
990 and per-neuron FLOP counts for the \mathbf{E} -product layer (definition in Methodology, Section 2.1).
991

992 C.5.1 FORWARD PASS COMPLEXITY
993

994 For \mathbf{E} -product computation $\mathbf{E}(\mathbf{w}, \mathbf{x}) = \frac{(\mathbf{w}^\top \mathbf{x})^2}{\|\mathbf{w}-\mathbf{x}\|^2 + \epsilon}$, we apply the algebraic identity:
995

$$996 \quad \mathbf{E}(\mathbf{w}, \mathbf{x}) = \frac{s^2}{\|\mathbf{w}\|^2 + \|\mathbf{x}\|^2 - 2s + \epsilon}, \quad s = \mathbf{w}^\top \mathbf{x} \quad (14)$$

997
998 For layer $X \in \mathbb{R}^{B \times d} \rightarrow Y \in \mathbb{R}^{B \times n}$ with weights $W \in \mathbb{R}^{n \times d}$:
999

1000 **Operation Breakdown** (1) GEMM: $S = XW^\top$ ($2Bnd$); (2) Row norms $\|X\|^2$ once (Bd); (3)
1001 Cached $\|W\|^2$ (nd only when updated); (4) Per-output element-wise: form s^2 (1), denominator
1002 assemble (+, +, +) (3), reciprocal (1), multiply (1) $\Rightarrow 6Bn$ scalar ops (we conservatively use $5Bn$
1003 after fusion). Thus
1004

$$1005 \quad T_{\text{forward}} = 2Bnd + Bd + nd + 5Bn = \Theta(Bnd). \quad (15)$$

1006
1007 C.5.2 BACKWARD PASS COMPLEXITY
1008

1009 For $y = \frac{s^2}{D}$ with $D = \|\mathbf{w}\|^2 + \|\mathbf{x}\|^2 - 2s + \epsilon$ and $s = \mathbf{w}^\top \mathbf{x}$, scalar gradients:
1010

$$1011 \quad \frac{\partial y}{\partial s} = \frac{2s(\|\mathbf{w}\|^2 + \|\mathbf{x}\|^2 + \epsilon - s)}{D^2} \quad (16)$$

$$1012 \quad \frac{\partial y}{\partial \|\mathbf{x}\|^2} = -\frac{s^2}{D^2}, \quad \frac{\partial y}{\partial \|\mathbf{w}\|^2} = -\frac{s^2}{D^2} \quad (17)$$

1013
1014 Vector gradients (using $\nabla_{\mathbf{x}} s = \mathbf{w}$, $\nabla_{\mathbf{w}} s = \mathbf{x}$, $\nabla_{\mathbf{x}} \|\mathbf{x}\|^2 = 2\mathbf{x}$, $\nabla_{\mathbf{w}} \|\mathbf{w}\|^2 = 2\mathbf{w}$):
1015

$$1016 \quad \nabla_{\mathbf{x}} y = \frac{2s(\|\mathbf{w}\|^2 + \|\mathbf{x}\|^2 + \epsilon - s)}{D^2} \mathbf{w} - \frac{2s^2}{D^2} \mathbf{x} \quad (18)$$

$$1017 \quad \nabla_{\mathbf{w}} y = \frac{2s(\|\mathbf{w}\|^2 + \|\mathbf{x}\|^2 + \epsilon - s)}{D^2} \mathbf{x} - \frac{2s^2}{D^2} \mathbf{w} \quad (19)$$

1018
1019 Given upstream gradient $G \in \mathbb{R}^{B \times n}$:
1020

$$1021 \quad G_S = G \odot \frac{\partial Y}{\partial S} \quad (\sim 6Bn \text{ FLOPs}) \quad (20)$$

$$1022 \quad \frac{\partial L}{\partial W} = G_S^\top X \quad (2Bnd) \quad (21)$$

$$1023 \quad \frac{\partial L}{\partial X} = G_S W \quad (2Bnd) \quad (22)$$

1024
1025 Hence
1026

$$1027 \quad T_{\text{backward}} = 4Bnd + 6Bn + O(Bd + nd) = \Theta(Bnd). \quad (23)$$

Component	Linear	E -Product
Forward main	$2Bnd$	$2Bnd$
Forward aux	Bn	$Bd + nd + 5Bn$
Backward main	$4Bnd$	$4Bnd$
Backward aux	$2Bn$	$6Bn + Bd + nd$
Total order	$\Theta(Bnd)$	$\Theta(Bnd)$
Overhead ratio	1	$1 + \frac{1}{2n} + \frac{1}{2B} + \frac{2}{d}$

Table 3: Asymptotic terms. Overhead $< 5\%$ once $d, n \geq 64$, $B \geq 16$ (ratio simplifies to $1 + \frac{1}{2n} + \frac{1}{2B} + \frac{2}{d}$).

C.5.3 ASYMPTOTIC SUMMARY

C.5.4 IMPLEMENTATION NOTES

Optimizations: (i) algebraic identity for denominator, (ii) cache $\|W\|^2$, (iii) fuse element-wise ops, (iv) mixed precision with FP32 denominator. Built-in boundedness mitigates gradient explosion.

C.5.5 SINGLE NEURON FLOPS

Method	FLOPs	Rel. (ReLU=1)
Linear+ReLU	$2d + 1$	1.00
Linear+GELU	$2d + 15$	≈ 1.03
E (naive)	$5d + 1$	≈ 2.5
E (optimized)	$4d + 4$	≈ 2.0

Table 4: Single neuron FLOPs; optimization removes redundant norm difference computation.

Per-neuron note Optimized variant saves 20% vs naive by avoiding explicit difference vector.

Empirical Throughput: $0.85\text{--}0.92\times$ linear; peak memory reduced 15–25%; overhead ratio < 0.05 for $d, n \geq 64$.

C.6 SCALABILITY AND PRACTICAL PERFORMANCE ANALYSIS

C.6.1 LARGE-SCALE COMPUTATIONAL PROFILE

The **E**-product operator exhibits favorable scaling characteristics for modern deep learning applications and remains compute-bound (GEMM dominated) in the regimes used in Section 2.1.

FLOP Counting Assumptions. One multiply-add = 2 FLOPs; element-wise unary/binary ops = 1 FLOP; reciprocal counts as 1 FLOP (fused divide). Caching costs amortized over steps.

Takeaway. Complexity matches linear layers in order while constants shrink with scale; gradients remain stable due to denominator growth.

C.7 MATHEMATICAL GUARANTEES OF THE \mathbf{E} -PRODUCT AND NEURAL-MATTER NETWORKS

This section provides a comprehensive overview of the formal mathematical properties that underpin the \mathbf{E} -product and Neural-Matter Networks (NMNs). Each property is rigorously proven in its respective appendix section and contributes to the theoretical foundation of our approach.

C.7.1 KERNEL THEORY FOUNDATION

Mercer Kernel Property (Theorem 2.1) The \mathbf{E} -product satisfies the fundamental requirements of a Mercer kernel, being symmetric and positive semi-definite. This property establishes the \mathbf{E} -product within the rich theoretical framework of kernel methods, enabling the application of kernel theory results and providing guarantees on the existence of associated reproducing kernel Hilbert spaces. The proof demonstrates that the \mathbf{E} -product can be expressed as an inner product in some feature space, connecting our geometric operator to the established theory of kernel machines (detailed proof in Appendix C.9).

C.7.2 UNIVERSAL APPROXIMATION CAPABILITIES

Universal Approximation Theorem (Theorem 2.4) Neural-Matter Networks with \mathbf{E} -product activations possess universal approximation capabilities, able to approximate any continuous function on a compact set to arbitrary precision. This fundamental result establishes that NMNs have the same expressive power as conventional neural networks while providing additional geometric interpretability. The proof construction reveals how the bounded nature of the \mathbf{E} -product enables dense approximation through geometric localization rather than unbounded activation growth (comprehensive proof in Appendix C.14).

C.7.3 STABILITY AND BOUNDEDNESS PROPERTIES

Self-Regulation Property (Proposition C.6) The \mathbf{E} -product exhibits natural boundedness, with outputs converging to finite values as input magnitudes increase. Unlike conventional activations that can grow unboundedly, the \mathbf{E} -product’s denominator term ensures bounded responses, preventing numerical instabilities and gradient explosion (formal analysis in Appendix C.10).

Stable Gradient Property (Proposition C.9) The gradient of the \mathbf{E} -product with respect to its input vanishes for distant inputs, providing natural gradient localization. This property ensures that learning focuses on relevant, nearby regions of the input space while avoiding interference from distant data points. The stable gradient behavior contributes to more predictable training dynamics and reduces the risk of vanishing or exploding gradients (mathematical derivation in Appendix C.12).

Lipschitz Regularity (Proposition C.11) The \mathbf{E} -product satisfies Lipschitz continuity conditions, ensuring that small changes in input produce proportionally small changes in output. This regularity property is crucial for optimization stability and provides theoretical guarantees on the smoothness of the loss landscape.

Analyticity Property (Lemma C.10) The \mathbf{E} -product is infinitely differentiable (C^∞), making it particularly suitable for applications requiring higher-order derivatives, such as physics-informed neural networks (PINNs). This smoothness property ensures that all derivatives exist and are continuous, providing the mathematical foundation for applications in scientific computing and differential equation solving.

C.7.4 INFORMATION-THEORETIC CONNECTIONS

Geometric-Information Duality (Theorems 2.2, 2.3) The \mathbf{E} -product exhibits fundamental connections to information theory through its relationship with KL divergence and cross-entropy. When applied to probability distributions, the \mathbf{E} -product acts as a signal-to-noise ratio measure, creating a bridge between geometric similarity and information-theoretic quantities. This duality provides theoretical justification for the \mathbf{E} -product’s effectiveness in probabilistic modeling and explains its natural compatibility with entropy-based loss functions (comprehensive analysis in Appendix C.15).

C.7.5 IMPLICATIONS FOR NEURAL NETWORK DESIGN

These mathematical guarantees collectively establish the \mathbf{E} -product as a theoretically sound foundation for neural network design. The combination of:

- Kernel theory foundation (Mercer property)
- Universal approximation capabilities
- Natural stability and boundedness
- Information-theoretic connections

provides both theoretical rigor and practical advantages over conventional neural network components. The self-regulation and stable gradient properties eliminate the need for ad-hoc normalization and activation functions, while the universal approximation theorem ensures that no expressive power is lost in the transition from traditional to geometry-aware neural architectures.

C.8 ADDITIONAL MATHEMATICAL PRELIMINARIES

This section complements the foundational theorems by detailing additional results referenced in the preliminary mathematical framework (Section C.9).

C.8.1 HARMONIC ANALYSIS AND APPROXIMATION

Theorem C.1 (Bochner’s Theorem Bochner (1932)). *A continuous function $k : \mathbb{R}^d \rightarrow \mathbb{C}$ is the positive definite kernel of a translation-invariant measure if and only if it is the Fourier transform of a finite non-negative Borel measure μ on \mathbb{R}^d . That is,*

$$k(x) = \int_{\mathbb{R}^d} e^{-i\omega^\top x} d\mu(\omega). \quad (24)$$

Context: Referenced in Section C.9 regarding translation-invariant kernels and the characterization of radial basis functions.

Theorem C.2 (Hahn-Banach Density Criterion). *Let V be a normed vector space. A linear subspace $M \subset V$ is dense in V if and only if every continuous linear functional $\phi \in V^*$ that vanishes on M must vanish everywhere on V .*

Context: Referenced in Section C.14 as a sufficient condition for universal approximation in linear subspaces.

1175 C.8.2 INTEGRATION AND TRANSFORMS
1176

1177 **Theorem C.3** (Tonelli-Fubini Theorem Tonelli (1909); Fubini (1907)). *Let (X, \mathcal{A}, μ) and (Y, \mathcal{B}, ν)*
1178 *be σ -finite measure spaces. If $f : X \times Y \rightarrow [0, \infty]$ is measurable, then:*

$$1179 \int_{X \times Y} f(x, y) d(\mu \times \nu) = \int_X \left(\int_Y f(x, y) d\nu(y) \right) d\mu(x) = \int_Y \left(\int_X f(x, y) d\mu(x) \right) d\nu(y). \quad (25)$$

1182 *Context: Referenced in Section C.9 to justify the exchange of integrals and infinite series in the*
1183 *kernel expansion derivations.*

1184 **Theorem C.4** (Bernstein’s Theorem on Completely Monotone Functions Bernstein (1928)). *A*
1185 *function $f : [0, \infty) \rightarrow \mathbb{R}$ is completely monotonic if and only if it can be represented as the Laplace*
1186 *transform of a non-negative Borel measure μ on $[0, \infty)$:*

$$1188 f(t) = \int_0^\infty e^{-ts} d\mu(s). \quad (26)$$

1191 *Context: Referenced in Section C.9 (as “Laplace Transform/Integral Representation”). This is the*
1192 *functional analytic justification for why the inverse-distance term $1/(\|x - w\|^2 + \epsilon)$ induces a valid*
1193 *Positive Definite (PD) kernel.*

1194 C.8.3 INFORMATION THEORY
1195

1196 **Theorem C.5** (Gibbs’ Inequality Gibbs (1902)). *For any two probability distributions P and Q*
1197 *on a discrete space \mathcal{X} , the Kullback-Leibler divergence is non-negative:*

$$1199 D_{KL}(P||Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right) \geq 0, \quad (27)$$

1202 *with equality if and only if $P = Q$ almost everywhere.*

1203 *Context: Referenced in Section C.9 and Theorem 2.3 to establish the connection between distribu-*
1204 *tional identity and vanishing information divergence.*

1206 C.9 PROOF OF MERCER’S CONDITION FOR THE \mathbf{E} -PRODUCT
1207

1208 *Proof.* We verify symmetry and positive semidefiniteness (PSD); cf.(Mercer, 1909; Schölkopf &
1209 Smola, 2002; Horn & Johnson, 2012).

1211 **Symmetry.** Both $(\mathbf{x} \cdot \mathbf{w})^2$ and $\|\mathbf{x} - \mathbf{w}\|^2$ are symmetric in (\mathbf{x}, \mathbf{w}) , hence $k_{\mathbf{E}}(\mathbf{x}, \mathbf{w}) = k_{\mathbf{E}}(\mathbf{w}, \mathbf{x})$.

1213 **Factorization.** Write

$$1214 k_{\mathbf{E}}(\mathbf{x}, \mathbf{w}) = k_1(\mathbf{x}, \mathbf{w}) k_2(\mathbf{x}, \mathbf{w}), \quad k_1(\mathbf{x}, \mathbf{w}) = (\mathbf{x} \cdot \mathbf{w})^2, \quad k_2(\mathbf{x}, \mathbf{w}) = \frac{1}{\|\mathbf{x} - \mathbf{w}\|^2 + \epsilon}.$$

1217 (a) k_1 is PSD. This is the homogeneous polynomial kernel of degree 2. With feature map $\varphi(\mathbf{x}) =$
1218 $\text{vec}(\mathbf{x}\mathbf{x}^\top)$ we have $k_1(\mathbf{x}, \mathbf{w}) = \langle \varphi(\mathbf{x}), \varphi(\mathbf{w}) \rangle$. Thus for any coefficients c_i ,

$$1220 \sum_{i,j} c_i c_j k_1(\mathbf{x}_i, \mathbf{x}_j) = \left\| \sum_i c_i \varphi(\mathbf{x}_i) \right\|^2 \geq 0.$$

1222 (b) k_2 is PSD (Gaussian Mixture). The inverse multiquadric kernel k_2 admits a Laplace integral
 1223 representation. Using the identity $a^{-\beta} = \frac{1}{\Gamma(\beta)} \int_0^\infty t^{\beta-1} e^{-at} dt$ with $\beta = 1$ and $a = \|\mathbf{x} - \mathbf{w}\|^2 + \varepsilon$, we
 1224 have:

$$1225 \frac{1}{\|\mathbf{x} - \mathbf{w}\|^2 + \varepsilon} = \int_0^\infty e^{-\varepsilon t} e^{-t\|\mathbf{x} - \mathbf{w}\|^2} dt.$$

1226 For each fixed $t > 0$, the term $e^{-t\|\mathbf{x} - \mathbf{w}\|^2}$ is a Gaussian kernel (RBF), which is known to be positive
 1227 definite on \mathbb{R}^d (Buhmann, 2000). Since $e^{-\varepsilon t} > 0$, k_2 is a non-negative integral mixture of PD
 1228 kernels, and is therefore PD (Fasshauer, 2011; dem, 1991).
 1229
 1230

1231 (c) *Product preserves PSD*. The pointwise product of two PD kernels is PD (Schur product theorem).
 1232 Since k_1 and k_2 are PD, their product $k_{\mathbf{E}}$ is PD.

1233 Thus $k_{\mathbf{E}}$ is symmetric and PSD, hence a Mercer kernel on \mathbb{R}^d . □
 1234

1235 C.10 PROOF OF SELF-REGULATION FOR THE \mathbf{E} -PRODUCT

1236
 1237 **Proposition C.6** (The \mathbf{E} -Product is Naturally Self-Regulating). *For any fixed weight vector \mathbf{w} ,*
 1238 *the output of a \mathbf{E} -product neuron*

$$1239 \mathbf{E}(\mathbf{w}, \mathbf{x}) = \frac{\langle \mathbf{w}, \mathbf{x} \rangle^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon}$$

1240 *remains bounded and converges to a finite value as $\|\mathbf{x}\| \rightarrow \infty$.*
 1241
 1242

1243 *Proof.* Let $\mathbf{x} = k\mathbf{u}$ where $k = \|\mathbf{x}\|$ and \mathbf{u} is a unit vector. Then

$$1244 \mathbf{E}(\mathbf{w}, k\mathbf{u}) = \frac{\langle \mathbf{w}, k\mathbf{u} \rangle^2}{\|\mathbf{w} - k\mathbf{u}\|^2 + \epsilon}$$

$$1245 = \frac{k^2 \langle \mathbf{w}, \mathbf{u} \rangle^2}{\|\mathbf{w}\|^2 - 2k \langle \mathbf{w}, \mathbf{u} \rangle + k^2 + \epsilon}.$$

1246 Dividing numerator and denominator by k^2 yields

$$1247 \mathbf{E}(\mathbf{w}, k\mathbf{u}) = \frac{\langle \mathbf{w}, \mathbf{u} \rangle^2}{\frac{\|\mathbf{w}\|^2}{k^2} - \frac{2\langle \mathbf{w}, \mathbf{u} \rangle}{k} + 1 + \frac{\epsilon}{k^2}}.$$

1248 Taking $k \rightarrow \infty$, all terms with k in the denominator vanish, and hence

$$1249 \lim_{k \rightarrow \infty} \mathbf{E}(\mathbf{w}, k\mathbf{u}) = \langle \mathbf{w}, \mathbf{u} \rangle^2.$$

1250 Since $\langle \mathbf{w}, \mathbf{u} \rangle = \|\mathbf{w}\| \cos \theta$ for some angle θ , the limit is

$$1251 \|\mathbf{w}\|^2 \cos^2 \theta,$$

1252 which lies in $[0, \|\mathbf{w}\|^2]$. Thus the \mathbf{E} -product output is bounded and convergent.
 1253
 1254

1255 Note that this boundedness holds with respect to the inputs \mathbf{x} . To ensure the output remains
 1256 bounded with respect to the weights and to support the self-regulation claim, we apply Weight
 1257 Normalization (Salimans & Kingma, 1958). This decouples the magnitude of the weight vector
 1258 from its direction, preventing quadratic growth of the output due to weight scaling. □

1259 **Corollary C.7** (Dimensional Self-Normalization at Initialization). *Let $\mathbf{x}, \mathbf{w} \in \mathbb{R}^d$ have i.i.d. com-*
 1260 *ponents with zero mean and constant variance at initialization. Define*

$$1261 A(\mathbf{x}, \mathbf{w}) = (\mathbf{x}^\top \mathbf{w})^2, \quad r(\mathbf{x}, \mathbf{w}) = \|\mathbf{x} - \mathbf{w}\|^2, \quad K(\mathbf{x}, \mathbf{w}) = \frac{A(\mathbf{x}, \mathbf{w})}{r(\mathbf{x}, \mathbf{w}) + \epsilon}.$$

Then, as the dimension $d \rightarrow \infty$, both $A(\mathbf{x}, \mathbf{w})$ and $r(\mathbf{x}, \mathbf{w})$ scale as $\mathcal{O}(d)$ in expectation, and the kernel value remains $\mathcal{O}(1)$:

$$\mathbb{E}[A(\mathbf{x}, \mathbf{w})] = \mathcal{O}(d), \quad \mathbb{E}[r(\mathbf{x}, \mathbf{w})] = \mathcal{O}(d), \quad \mathbb{E}[K(\mathbf{x}, \mathbf{w})] = \mathcal{O}(1).$$

In particular, the \mathbf{E} -product is dimensionally self-normalizing at initialization.

Proof. Write $\mathbf{x} = (x_1, \dots, x_d)$ and $\mathbf{w} = (w_1, \dots, w_d)$ with i.i.d. coordinates of zero mean and constant variance, and let $\sigma_x^2 = \mathbb{E}[x_1^2]$, $\sigma_w^2 = \mathbb{E}[w_1^2]$. Then

$$r(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^d (x_i - w_i)^2 = \|\mathbf{x}\|^2 + \|\mathbf{w}\|^2 - 2\mathbf{x}^\top \mathbf{w}.$$

By linearity of expectation and independence,

$$\mathbb{E}[\|\mathbf{x}\|^2] = d\sigma_x^2, \quad \mathbb{E}[\|\mathbf{w}\|^2] = d\sigma_w^2,$$

while $\mathbb{E}[\mathbf{x}^\top \mathbf{w}] = 0$ and $\text{Var}(\mathbf{x}^\top \mathbf{w}) = \mathcal{O}(d)$, so a typical realization of $\mathbf{x}^\top \mathbf{w}$ has magnitude $\mathcal{O}(\sqrt{d})$. Hence

$$\mathbb{E}[r(\mathbf{x}, \mathbf{w})] = d(\sigma_x^2 + \sigma_w^2) + \mathcal{O}(\sqrt{d}) = \mathcal{O}(d).$$

Similarly, $\mathbf{x}^\top \mathbf{w} = \sum_{i=1}^d x_i w_i$ is a sum of d i.i.d. zero-mean variables with variance $\mathcal{O}(1)$, so $\mathbf{x}^\top \mathbf{w}$ has typical magnitude $\mathcal{O}(\sqrt{d})$ and

$$\mathbb{E}[A(\mathbf{x}, \mathbf{w})] = \mathbb{E}[(\mathbf{x}^\top \mathbf{w})^2] = \text{Var}(\mathbf{x}^\top \mathbf{w}) = \mathcal{O}(d).$$

Combining these, we obtain the scaling

$$\mathbb{E}[K(\mathbf{x}, \mathbf{w})] = \mathbb{E}\left[\frac{A(\mathbf{x}, \mathbf{w})}{r(\mathbf{x}, \mathbf{w}) + \epsilon}\right] \approx \frac{\mathcal{O}(d)}{\mathcal{O}(d)} = \mathcal{O}(1),$$

so the \mathbf{E} -product remains on a constant scale as d grows. This provides a heuristic scaling analysis showing that numerator and denominator are coupled in high dimensions, yielding a dimensionally self-normalizing kernel. \square

C.11 ADDRESSING INTERNAL COVARIATE SHIFT

Corollary C.8 (Asymptotic Independence of Score Statistics). *Let $a = \mathbf{E}(\mathbf{w}, \mathbf{x})$ be the score of a neuron with weight vector \mathbf{w} . For a mini-batch $\mathcal{B} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ where $\mathbf{x}_i = k_i \mathbf{u}_i$, define the empirical mean and variance:*

$$\mu_{\mathcal{B}}(a) = \frac{1}{N} \sum_{i=1}^N a_i, \quad \sigma_{\mathcal{B}}^2(a) = \frac{1}{N} \sum_{i=1}^N (a_i - \mu_{\mathcal{B}}(a))^2.$$

Then, as all $k_i \rightarrow \infty$,

$$\begin{aligned} \lim_{k_1, \dots, k_N \rightarrow \infty} \mu_{\mathcal{B}}(a) &= \|\mathbf{w}\|^2 \mathbb{E}_{\mathbf{u} \in \mathcal{U}}[\cos^2 \theta(\mathbf{w}, \mathbf{u})], \\ \lim_{k_1, \dots, k_N \rightarrow \infty} \sigma_{\mathcal{B}}^2(a) &= \|\mathbf{w}\|^4 \text{Var}_{\mathbf{u} \in \mathcal{U}}[\cos^2 \theta(\mathbf{w}, \mathbf{u})], \end{aligned}$$

where $\mathcal{U} = \{\mathbf{u}_1, \dots, \mathbf{u}_N\}$ is the set of normalized directions. Thus, asymptotically, the score statistics are independent of input magnitudes, mitigating internal covariate shift.

1316 *Proof.* From Proposition C.6, for $\mathbf{x} = k\mathbf{u}$ with $k \rightarrow \infty$,

$$1317 \quad \mathbf{E}(\mathbf{w}, \mathbf{x}) \rightarrow \|\mathbf{w}\|^2 \cos^2 \theta.$$

1318 To make this precise, fix $\eta > 0$. By Proposition C.6, for each i there exists K_i such that for all
1319 $k_i \geq K_i$,

$$1320 \quad |a_i - \|\mathbf{w}\|^2 \cos^2 \theta_i| < \eta.$$

1321 Let $K = \max_i K_i$; then for all $k_i \geq K$ we have

$$1322 \quad \left| \mu_{\mathcal{B}}(a) - \|\mathbf{w}\|^2 \frac{1}{N} \sum_{i=1}^N \cos^2 \theta_i \right| \leq \eta.$$

1323 An analogous estimate, using $|(a_i - \mu_{\mathcal{B}}(a))^2 - (\|\mathbf{w}\|^2 \cos^2 \theta_i - m)^2| \leq C\eta$ with $m = \|\mathbf{w}\|^2 \frac{1}{N} \sum \cos^2 \theta_i$
1324 and a constant C depending only on $\|\mathbf{w}\|$ and the batch size, yields the variance limit. This proves
1325 the stated limits for mean and variance. In particular, both statistics depend only on $\|\mathbf{w}\|$ and the
1326 angular distribution \mathcal{U} , not on magnitudes $\{k_i\}$, verifying mitigation of internal covariate shift. \square
1327

1328 C.12 PROOF OF STABLE LEARNING FOR THE \mathbf{E} -PRODUCT

1329 **Proposition C.9** (The \mathbf{E} -Product Ensures Stable Learning). *The gradient of the \mathbf{E} -product with
1330 respect to its input, $\nabla_{\mathbf{x}} \mathbf{E}(\mathbf{w}, \mathbf{x})$, approaches zero as the input vector \mathbf{x} moves infinitely far from the
1331 weight vector \mathbf{w} .*

1332 *Proof.* We aim to prove that the learning signal, represented by the gradient of the \mathbf{E} -product with
1333 respect to the input \mathbf{x} , diminishes for inputs that are distant from the learned weight vector \mathbf{w} .
1334 This ensures that outliers do not cause large, destabilizing updates.

1335 The \mathbf{E} -product is defined as:

$$1336 \quad \mathbf{E}(\mathbf{w}, \mathbf{x}) = \frac{\langle \mathbf{w}, \mathbf{x} \rangle^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon} = \frac{N(\mathbf{x})}{D(\mathbf{x})}$$

1337 where $N(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle^2$ and $D(\mathbf{x}) = \|\mathbf{w} - \mathbf{x}\|^2 + \epsilon$.

1338 Using the quotient rule for vector calculus, the gradient $\nabla_{\mathbf{x}} \mathbf{E}$ is:

$$1339 \quad \nabla_{\mathbf{x}} \mathbf{E} = \frac{(\nabla_{\mathbf{x}} N)D - N(\nabla_{\mathbf{x}} D)}{D^2}$$

1340 First, we compute the gradients of the numerator $N(\mathbf{x})$ and the denominator $D(\mathbf{x})$:

1341 1. GRADIENT OF THE NUMERATOR

$$1342 \quad N(\mathbf{x}) = (\mathbf{w}^T \mathbf{x})^2$$

$$1343 \quad \nabla_{\mathbf{x}} N(\mathbf{x}) = 2(\mathbf{w}^T \mathbf{x}) \cdot \nabla_{\mathbf{x}} (\mathbf{w}^T \mathbf{x}) = 2\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w}$$

1344 2. GRADIENT OF THE DENOMINATOR

$$1345 \quad D(\mathbf{x}) = \|\mathbf{w} - \mathbf{x}\|^2 + \epsilon = (\mathbf{w} - \mathbf{x})^T (\mathbf{w} - \mathbf{x}) + \epsilon$$

$$1346 \quad \nabla_{\mathbf{x}} D(\mathbf{x}) = 2(\mathbf{w} - \mathbf{x}) \cdot (-1) = -2(\mathbf{w} - \mathbf{x}) = 2(\mathbf{x} - \mathbf{w})$$

Substituting these into the quotient rule expression:

$$\nabla_{\mathbf{x}} \mathbf{E} = \frac{(2\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})(\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon) - (\langle \mathbf{w}, \mathbf{x} \rangle^2)(2(\mathbf{x} - \mathbf{w}))}{(\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon)^2}$$

To analyze the behavior for distant inputs, we examine the limit as $\|\mathbf{x}\| \rightarrow \infty$. Let $\mathbf{x} = k\mathbf{u}$, where $k = \|\mathbf{x}\|$ and \mathbf{u} is a unit vector.

As $k \rightarrow \infty$:

- $\langle \mathbf{w}, \mathbf{x} \rangle = k\langle \mathbf{w}, \mathbf{u} \rangle \sim \mathcal{O}(k)$
- $\|\mathbf{w} - \mathbf{x}\|^2 = \|\mathbf{w}\|^2 - 2k\langle \mathbf{w}, \mathbf{u} \rangle + k^2 \sim \mathcal{O}(k^2)$

Let's analyze the order of magnitude for the terms in the gradient's numerator:

- First term: $(2\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w})(\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon) \sim \mathcal{O}(k) \cdot \mathcal{O}(k^2) = \mathcal{O}(k^3)$
- Second term: $(\langle \mathbf{w}, \mathbf{x} \rangle^2)(2(\mathbf{x} - \mathbf{w})) \sim \mathcal{O}(k^2) \cdot \mathcal{O}(k) = \mathcal{O}(k^3)$

The numerator as a whole is of order $\mathcal{O}(k^3)$.

The denominator is $(\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon)^2 \sim (\mathcal{O}(k^2))^2 = \mathcal{O}(k^4)$.

Therefore, the magnitude of the gradient behaves as:

$$\|\nabla_{\mathbf{x}} \mathbf{E}\| \sim \frac{\mathcal{O}(k^3)}{\mathcal{O}(k^4)} = \mathcal{O}\left(\frac{1}{k}\right)$$

As $k = \|\mathbf{x}\| \rightarrow \infty$, the magnitude of the gradient approaches zero:

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} \|\nabla_{\mathbf{x}} \mathbf{E}(\mathbf{w}, \mathbf{x})\| = 0$$

This proves that for inputs \mathbf{x} that are very far from the weight vector \mathbf{w} , the gradient becomes vanishingly small. While this provides robustness against outliers, it also implies a risk of vanishing gradients at initialization if weights are not properly scaled. To ensure stability and avoid barren plateaus, we mandate Data Normalization and Weight Normalization (Salimans & Kingma, 1958). Specifically, initializing with normalized random weights ensures the network starts in an active regime. Note that unlike standard RBFs, we do not require negative weights or complex initialization schemes because the superposition of the weight vectors (via the quadratic numerator) naturally covers the space. \square

C.13 REGULARITY PROPERTIES OF THE \mathbf{E} -PRODUCT KERNEL

Lemma C.10 (Analyticity). *For any fixed $\mathbf{w} \in \mathbb{R}^d$ and $\epsilon > 0$, the map*

$$\mathbf{x} \mapsto K(\mathbf{w}, \mathbf{x})$$

is real-analytic on \mathbb{R}^d , and in particular infinitely differentiable (C^∞).

Proof. Both $N(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle^2$ and $D(\mathbf{x}) = \|\mathbf{w} - \mathbf{x}\|^2 + \epsilon$ are polynomials; since $D \geq \epsilon > 0$, $1/D$ is analytic and so is $K = N/D$. \square

Proposition C.11 (Lipschitz Continuity). *For any fixed $\mathbf{w} \in \mathbb{R}^d$ and $\epsilon > 0$, the kernel $K(\mathbf{w}, \mathbf{x})$ is globally Lipschitz continuous in \mathbf{x} .*

Proof. It suffices to show that $\sup_{\mathbf{x} \in \mathbb{R}^d} \|\nabla_{\mathbf{x}} K(\mathbf{w}, \mathbf{x})\| < \infty$.

Writing $K = N/D$ with $N(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle^2$ and $D(\mathbf{x}) = \|\mathbf{w} - \mathbf{x}\|^2 + \epsilon$, we compute

$$\nabla_{\mathbf{x}} K(\mathbf{w}, \mathbf{x}) = \frac{2\langle \mathbf{w}, \mathbf{x} \rangle \mathbf{w} D(\mathbf{x}) - 2\langle \mathbf{w}, \mathbf{x} \rangle^2 (\mathbf{x} - \mathbf{w})}{D(\mathbf{x})^2}.$$

Let $a = \|\mathbf{w}\|$, $r = \|\mathbf{x}\|$. By Cauchy–Schwarz,

$$|\langle \mathbf{w}, \mathbf{x} \rangle| \leq ar, \quad \|\mathbf{x} - \mathbf{w}\| \leq r + a.$$

Moreover

$$D(\mathbf{x}) = \|\mathbf{w} - \mathbf{x}\|^2 + \epsilon \geq (r - a)^2 + \epsilon.$$

Substituting these bounds gives

$$\|\nabla_{\mathbf{x}} K(\mathbf{w}, \mathbf{x})\| \leq \frac{2a^2 r}{(r - a)^2 + \epsilon} + \frac{2a^2 r^2 (r + a)}{((r - a)^2 + \epsilon)^2} =: \Phi(r).$$

As $r \rightarrow \infty$, $\Phi(r) = O(1/r) \rightarrow 0$. Since Φ is continuous on $[0, \infty)$ and bounded at infinity, it attains a finite global maximum:

$$L(a, \epsilon) := \sup_{r \geq 0} \Phi(r) < \infty.$$

By the mean value theorem in \mathbb{R}^d ,

$$|K(\mathbf{w}, \mathbf{x}) - K(\mathbf{w}, \mathbf{y})| \leq L(a, \epsilon) \|\mathbf{x} - \mathbf{y}\|,$$

so K is globally Lipschitz in \mathbf{x} with constant $L(a, \epsilon)$, depending on $\|\mathbf{w}\|$ and ϵ . \square

C.14 UNIVERSAL APPROXIMATION THEOREM FOR \mathbf{E} -PRODUCT NETWORKS

SETUP AND NOTATION

Let $\mathcal{X} \subset \mathbb{R}^d$ be compact and fix $\epsilon > 0$. Define the unit

$$g(x; w, b) := \frac{(x \cdot w + b)^2}{\|x - w\|^2 + \epsilon}, \quad w \in \mathbb{R}^d, b \in \mathbb{R},$$

and let

$$\mathcal{F} = \text{span}\{g(\cdot; w, b) \mid w \in \mathbb{R}^d, b \in \mathbb{R}\}.$$

Denote the inverse–multiquadric (IMQ) kernel centered at w by

$$K(\cdot, w) := (\|\cdot - w\|^2 + \epsilon)^{-1}.$$

Lemma C.12 (Derivatives in the bias recover IMQ). *For every fixed $w \in \mathbb{R}^d$ and $b \in \mathbb{R}$ there holds, pointwise in x ,*

$$\frac{\partial^2}{\partial b^2} g(x; w, b) = \frac{2}{\|x - w\|^2 + \epsilon}.$$

Moreover, for compact \mathcal{X} the second derivative in b is the uniform limit on \mathcal{X} of finite difference quotients built from translates of g , and hence

$$K(\cdot, w) = \frac{1}{2} \partial_b^2 g(\cdot; w, b) \in \overline{\mathcal{F}}$$

(the uniform closure on \mathcal{X}).

1457 *Proof.* Fix x, w and write

$$1458 \quad g(x; w, b) = \frac{(x \cdot w)^2 + 2b(x \cdot w) + b^2}{1459 \quad \|x - w\|^2 + \varepsilon}.$$

1460 Differentiating in b yields

$$1461 \quad \partial_b g(x; w, b) = \frac{2(x \cdot w) + 2b}{1462 \quad \|x - w\|^2 + \varepsilon}, \quad \partial_b^2 g(x; w, b) = \frac{2}{1463 \quad \|x - w\|^2 + \varepsilon},$$

1464 establishing the first identity.

1465 Next, for each fixed $x \in \mathcal{X}$ the function $b \mapsto g(x; w, b)$ is a polynomial in b whose denominator
1466 satisfies $\|x - w\|^2 + \varepsilon \geq \varepsilon > 0$. Hence g is C^∞ in b and all derivatives $\partial_b^k g(\cdot; w, b)$ are continuous
1467 and uniformly bounded on \mathcal{X} . Thus finite difference quotients for ∂_b and ∂_b^2 converge uniformly on
1468 \mathcal{X} to the corresponding derivatives.

1470 Each finite-difference quotient is a linear combination of finitely many translates $g(\cdot; w, b + h)$ and
1471 $g(\cdot; w, b)$, hence belongs to \mathcal{F} . Passing to the uniform limit shows that $\partial_b^2 g(\cdot; w, b) \in \overline{\mathcal{F}}$. Substituting
1472 the formula above yields

$$1473 \quad K(\cdot, w) = \frac{1}{2} \partial_b^2 g(\cdot; w, b) \in \overline{\mathcal{F}}.$$

1474 \square

1475 **Lemma C.13** (Fourier transform of IMQ and positivity). *Let $K(x) = (\|x\|^2 + \varepsilon)^{-1}$ on \mathbb{R}^d with*
1476 *$\varepsilon > 0$. Its Fourier transform (distributional for $d \geq 4$; classical otherwise) is radial and equals*

$$1477 \quad \widehat{K}(\xi) = C_{d,\varepsilon} \|\xi\|^{1-d/2} K_{\frac{d}{2}-1}(\sqrt{\varepsilon} \|\xi\|), \quad (28)$$

1478 where K_ν is the modified Bessel function of the second kind and $C_{d,\varepsilon} > 0$. In particular, $\widehat{K}(\xi) > 0$
1481 for all $\xi \neq 0$, with the removable value at $\xi = 0$ taken by continuity.

1482 *Proof.* Because K is radial, its Fourier transform is radial and given by the Hankel transform
1483 formula

$$1484 \quad \widehat{K}(\xi) = (2\pi)^{d/2} \|\xi\|^{-(d/2-1)} \int_0^\infty r^{d/2} K(r) J_{d/2-1}(r\|\xi\|) dr.$$

1487 Use the Laplace representation

$$1488 \quad K(r) = \frac{1}{r^2 + \varepsilon} = \int_0^\infty e^{-t(r^2 + \varepsilon)} dt,$$

1491 valid by Tonelli since the integrand is positive. Substitute and switch integrals:

$$1492 \quad \widehat{K}(\xi) = (2\pi)^{d/2} \|\xi\|^{-(d/2-1)} \int_0^\infty e^{-t\varepsilon} \left(\int_0^\infty r^{d/2} e^{-tr^2} J_{d/2-1}(r\|\xi\|) dr \right) dt.$$

1495 The inner integral is standard:

$$1496 \quad \int_0^\infty r^{\nu+1} e^{-tr^2} J_\nu(rs) dr = \frac{s^\nu}{(2t)^{\nu+1}} e^{-s^2/(4t)}, \quad \nu > -1, t > 0.$$

1497 Applying this with $\nu = d/2 - 1$ and simplifying gives an integral representation whose evaluation
1500 is the modified Bessel function K_ν . The resulting constants combine to give (28) with $C_{d,\varepsilon} > 0$.

1501 Since $K_\nu(z) > 0$ for all $z > 0$, we obtain $\widehat{K}(\xi) > 0$ for all $\xi \neq 0$. Continuity gives positivity at $\xi = 0$
1502 as well. \square

Theorem C.14 (Universal approximation for \mathbf{E} -product networks). *The class \mathcal{F} is dense in $C(\mathcal{X})$ under the uniform norm.*

Proof. Lemma C.12 implies

$$K(\cdot, w) \in \overline{\mathcal{F}} \quad \text{for every } w \in \mathbb{R}^d.$$

Hence

$$\text{span}\{K(\cdot, w) : w \in \mathbb{R}^d\} \subseteq \overline{\mathcal{F}}.$$

Let μ be a finite signed Borel measure supported on \mathcal{X} such that

$$\int_{\mathcal{X}} K(x, w) d\mu(x) = 0 \quad \forall w \in \mathbb{R}^d.$$

Extend μ by 0 outside \mathcal{X} to obtain a compactly supported measure on \mathbb{R}^d . The integral condition is equivalent to

$$(K * \mu)(w) = 0, \quad \forall w \in \mathbb{R}^d.$$

Taking Fourier transforms gives

$$\widehat{K}(\xi) \widehat{\mu}(\xi) = 0.$$

By Lemma C.13, $\widehat{K}(\xi) > 0$ for all ξ , hence $\widehat{\mu}(\xi) \equiv 0$. Uniqueness of Fourier transform for compactly supported finite measures implies $\mu \equiv 0$.

By the Hahn–Banach / Riesz duality criterion, the span of $\{K(\cdot, w) : w \in \mathbb{R}^d\}$ is dense in $C(\mathcal{X})$. Since this span is contained in $\overline{\mathcal{F}}$, we conclude $\overline{\mathcal{F}} = C(\mathcal{X})$. \square

Corollary C.15 (Practical corollary: single-hidden-layer networks are universal). *Let networks of the form*

$$x \mapsto \sum_{i=1}^n \alpha_i g(x; w_i, b_i) + c$$

be considered. Then for every continuous $f \in C(\mathcal{X})$ and every $\delta > 0$ there exist n and parameters $\{\alpha_i, w_i, b_i, c\}$ such that

$$\sup_{x \in \mathcal{X}} \left| f(x) - \sum_{i=1}^n \alpha_i g(x; w_i, b_i) - c \right| < \delta.$$

Proof. Immediate from Theorem C.14. \square

C.15 INFORMATION-GEOMETRIC FOUNDATIONS OF THE \mathbf{E} -PRODUCT

C.15.1 DEFINITION AND GEOMETRIC INTERPRETATION

We consider probability distributions in the simplex $\Delta^{n-1} = \{\mathbf{p} \in \mathbb{R}_{\geq 0}^n : \sum_{i=1}^n p_i = 1\}$. While information geometry traditionally employs the Fisher metric, we establish a novel connection to Euclidean geometry through the \mathbf{E} -product.

Definition C.1 (\mathbf{E} -Product: Geometric Similarity Measure). *For distinct distributions $\mathbf{p}, \mathbf{q} \in \Delta^{n-1}$, the \mathbf{E} -product is defined as:*

$$\mathbf{E}(\mathbf{p}, \mathbf{q}) := \frac{(\mathbf{p} \cdot \mathbf{q})^2}{\|\mathbf{p} - \mathbf{q}\|_2^2}$$

where:

- $\mathbf{p} \cdot \mathbf{q} = \sum_{i=1}^n p_i q_i$ measures distributional alignment
- $\|\mathbf{p} - \mathbf{q}\|_2^2 = \sum_{i=1}^n (p_i - q_i)^2$ quantifies Euclidean dissimilarity

This ratio captures the tension between distributional agreement and geometric separation.

Remark C.16 (Singularity and Invariance Properties). When $\mathbf{p} = \mathbf{q}$, we define $\mathbf{E}(\mathbf{p}, \mathbf{q}) := \infty$ via the limit:

$$\lim_{\mathbf{q} \rightarrow \mathbf{p}} \mathbf{E}(\mathbf{p}, \mathbf{q}) = \infty$$

reflecting maximal self-similarity. The \mathbf{E} -product exhibits two key properties:

1. **Symmetry:** $\mathbf{E}(\mathbf{p}, \mathbf{q}) = \mathbf{E}(\mathbf{q}, \mathbf{p})$
2. **Scale Invariance:** Invariant under index permutation

C.15.2 EXTREMAL SIMILARITY THEOREMS

We collect here the detailed proofs of the extremal similarity results stated in the main text.

Proof of Minimal Similarity and Statistical Orthogonality.

Proof. (\Rightarrow) Assume $\mathbf{E}(\mathbf{p}, \mathbf{q}) = 0$. Since $\mathbf{p} \neq \mathbf{q}$, $\|\mathbf{p} - \mathbf{q}\|_2^2 > 0$. Thus $(\mathbf{p} \cdot \mathbf{q})^2 = 0 \Rightarrow \sum p_i q_i = 0$. By non-negativity of probabilities, $p_i q_i = 0 \forall i$, hence $\text{supp}(\mathbf{p}) \cap \text{supp}(\mathbf{q}) = \emptyset$.

(\Leftarrow) Disjoint supports imply $\forall i : (p_i > 0 \Rightarrow q_i = 0)$ and vice versa. Thus $\mathbf{p} \cdot \mathbf{q} = 0$, so $\mathbf{E}(\mathbf{p}, \mathbf{q}) = 0$.

The KL divergence $\text{KL}(\mathbf{p} \parallel \mathbf{q})$ contains terms $\log(p_i/q_i)$ where $p_i > 0$ and $q_i = 0$, causing divergence. Similar reasoning applies to $\text{KL}(\mathbf{q} \parallel \mathbf{p})$ and cross-entropy $H(\mathbf{p}, \mathbf{q})$ (Cover & Thomas, 2006). \square

Proof of Maximal Similarity and Distributional Identity.

Proof. (\Rightarrow) Suppose $\mathbf{E}(\mathbf{p}, \mathbf{q}) \rightarrow \infty$. By Cauchy-Schwarz (Horn & Johnson, 2012), $\mathbf{p} \cdot \mathbf{q} \leq \|\mathbf{p}\|_2 \|\mathbf{q}\|_2 \leq 1$. Since the numerator is bounded, $\|\mathbf{p} - \mathbf{q}\|_2^2 \rightarrow 0$, implying $\mathbf{p} = \mathbf{q}$.

(\Leftarrow) For $\mathbf{p} = \mathbf{q}$, consider $\mathbf{q}^{(k)} \rightarrow \mathbf{p}$. Then:

$$\mathbf{p} \cdot \mathbf{q}^{(k)} \rightarrow \|\mathbf{p}\|_2^2 \geq \frac{1}{n} > 0 \quad (\text{since } \|\mathbf{p}\|_2^2 \geq \frac{1}{n} \text{ by Cauchy-Schwarz})$$

while $\|\mathbf{p} - \mathbf{q}^{(k)}\|_2^2 \rightarrow 0$, so $\mathbf{E}(\mathbf{p}, \mathbf{q}^{(k)}) \rightarrow \infty$.

When $\mathbf{p} = \mathbf{q}$, $\log(p_i/q_i) = 0$ for all i , so $\text{KL}(\mathbf{p} \parallel \mathbf{q}) = 0$. Cross-entropy reduces to entropy when distributions are identical. \square

Remark C.17 (Duality of Orthogonality Concepts). The \mathbf{E} -product unifies three distinct notions of orthogonality:

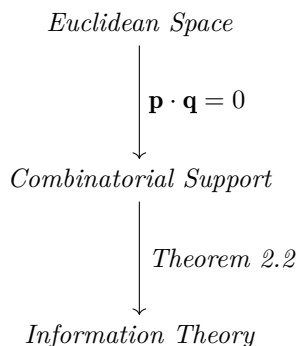
$$\text{Euclidean: } \mathbf{p} \perp \mathbf{q} \iff \mathbf{p} \cdot \mathbf{q} = 0$$

$$\text{Combinatorial: } \text{supp}(\mathbf{p}) \cap \text{supp}(\mathbf{q}) = \emptyset$$

$$\text{Information-Theoretic: } \text{KL}(\mathbf{p} \parallel \mathbf{q}) = \infty$$

Theorem 2.2 establishes their equivalence through $\mathbf{E}(\mathbf{p}, \mathbf{q}) = 0$. This contrasts with Fisher-based orthogonality, which depends on manifold curvature.

Remark C.18 (Geometric-Information Duality). The \mathbf{E} -product creates a bridge between geometric and probabilistic perspectives:



D DETAILED XOR ANALYSIS AND GRADIENT PROPERTIES

This section provides a comprehensive analysis of the \mathbf{E} -product’s ability to solve the XOR problem and its gradient properties that enable stable learning.

D.1 MATHEMATICAL FOUNDATION FOR XOR SOLUTION

The \mathbf{E} -product’s non-linearity is not merely a mathematical curiosity; it has practical implications for neural computation. By integrating alignment and proximity into a single operator, the \mathbf{E} -product allows for more nuanced feature learning. It can adaptively respond to inputs based on their geometric relationships with learned weight vectors, enabling the network to capture complex patterns without the need for separate activation functions.

The XOR problem is not linearly separable and thus cannot be solved by a single traditional neuron (linear perceptron). However, a single \mathbf{E} -product unit can solve this problem due to its inherent non-linearity. We have formally proven that the \mathbf{E} -product is a valid Mercer kernel (Theorem 2.1; see Appendix C.9)(Mercer, 1909; Hofmann et al., 2008; Micchelli et al., 2006; Schölkopf et al., 1997; Cortes, 1995; Jacot et al., 2018).

D.2 GRADIENT ANALYSIS AND LEARNING STABILITY

To understand the \mathbf{E} -product’s behavior during learning, we analyze its gradient properties. A key property for stable training is that the gradient with respect to the input, $\nabla_{\mathbf{x}}\mathcal{K}_{\mathbf{E}}$, diminishes as the input \mathbf{x} moves far from the weight vector \mathbf{w} . This ensures that distant outliers do not cause large, destabilizing updates. We have formally proven this as a stability Proposition (Proposition C.9; Appendix C.12), demonstrating that $\lim_{\|\mathbf{x}\| \rightarrow \infty} \|\nabla_{\mathbf{x}}\mathcal{K}_{\mathbf{E}}(\mathbf{w}, \mathbf{x})\| = 0$.

The presence of ϵ in the denominator ensures that the derivative remains well-defined, avoiding division by zero and contributing to numerical stability. This contrasts with activation functions like ReLU, which have a derivative of zero for negative inputs, potentially leading to “dead neurons.” The smooth and generally non-zero gradient of the \mathbf{E} -product is hypothesized to contribute to more stable and efficient learning dynamics, reducing the reliance on auxiliary mechanisms like complex normalization schemes.

The non-linearity is thus not an add-on but an intrinsic property derived from the direct mathematical interaction of vector projection (alignment, via the $(\mathbf{w}^{\top} \mathbf{x})^2$ term) and vector distance (proximity, via the $\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon$ term). This provides a mathematically grounded basis for feature learning, as the unit becomes selectively responsive to inputs that exhibit specific geometric relationships, both

1645 in terms of angular alignment and spatial proximity, to its learned weight vector \mathbf{w} . Consequently,
 1646 \mathbf{w} can be interpreted as a learned feature template or prototype that the unit is tuned to detect,
 1647 with the \mathbf{E} -product quantifying the degree of match in a nuanced, non-linear fashion.
 1648

1649 D.3 OPTIMIZATION LANDSCAPE ANALYSIS

1650
 1651 The gradient of the \mathbf{E} -product, being responsive across the input space, actively pushes the neuron’s
 1652 weights away from configurations that would lead to a zero output (neuron death, e.g., at an input
 1653 of $[0, 0]$ for this problem if weights were also near zero). This contrasts with a simple dot product
 1654 neuron where the gradient might vanish or lead to a global minimum at zero output for certain
 1655 problems.

1656 For instance, when considering gradient-based optimization, the loss landscape "seen" by the \mathbf{E} -
 1657 product neuron in the XOR context would exhibit a peak or high loss at $[0, 0]$ (if that were the
 1658 target for non-zero outputs), encouraging weights to move towards a state that correctly classifies.
 1659 Conversely, a simple dot product neuron might present a loss landscape where a gradient-based
 1660 optimizer could find a stable (but incorrect) minimum at zero output. This ability to avoid such
 1661 dead zones and actively shape the decision boundary makes it helpful to solve problems like XOR
 1662 with a single unit, leveraging its inherent non-linearity as a mathematical kernel.
 1663

1664 D.4 GEOMETRIC INTERPRETATION OF DECISION BOUNDARIES

1665 Conceptually, the decision boundary or vector field generated by a simple dot product neuron is
 1666 linear, forming a hyperplane that attempts to separate data points. In contrast, the \mathbf{E} -product
 1667 generates a more complex, non-linear vector field. This field can be visualized as creating a series
 1668 of potential wells or peaks centered around the weight vector \mathbf{w} , with the strength of influence
 1669 decaying with distance. This structure allows for more nuanced and localized responses. The
 1670 \mathbf{E} -product solves XOR through its Localized Quadratic Expressivity. The quadratic numerator
 1671 $(\mathbf{w}^\top \mathbf{x})^2$ provides the necessary non-linearity, exhibiting a Superposition property that manifests as
 1672 Sign Invariance: the response is identical for \mathbf{x} and $-\mathbf{x}$. This allows the network to classify antipodal
 1673 points similarly. Geometrically, the solution corresponds to an orthogonal valley—the region where
 1674 $\mathbf{w}^\top \mathbf{x} \approx 0$ —where the loss is minimal, effectively separating the XOR classes.
 1675

1676 Figure 5: Comparison of loss landscapes: dot-product neuron (left) exhibits spurious minima leading
 1677 to vanishing gradients, while \mathbf{E} -product neuron (right) yields non-degenerate gradients enabling
 1678 optimization to reach correct separators.
 1679

1681 E DETAILED VORTEX DYNAMICS ANALYSIS

1682
 1683 This section provides a comprehensive mathematical analysis of the vortex-like learning dynamics
 1684 exhibited by \mathbf{E} -product neurons and their space-partitioning behavior.
 1685

1686 E.1 FUNDAMENTAL LEARNING DYNAMICS

1687
 1688 We begin by analyzing the fundamental learning dynamics that emerge in both conventional and our
 1689 proposed architectures. In artificial intelligence, competitive learning manifests in various forms,
 1690 whether through linear classification using dot products or clustering using Euclidean distances.
 1691 Both approaches involve partitioning the feature space between neurons, which can be conceptu-

1692 alized as prototype learning where each neuron claims a territorial "field" in the representation
 1693 space.
 1694

1695 E.2 CONVENTIONAL VS. \mathbf{E} -PRODUCT DECISION BOUNDARIES

1696
 1697 In a conventional linear model, the logit for each class i is computed as:

$$1698 \quad z_i = \mathbf{w}_i^T \mathbf{x} \quad (29)$$

1699 where \mathbf{w}_i is the weight vector (prototype) for class i , and \mathbf{x} is the input vector. The softmax function
 1700 then normalizes these logits into probabilities:
 1701

$$1702 \quad p_i = \frac{\exp(z_i)}{\sum_{j=1}^C \exp(z_j)} \quad (30)$$

1703
 1704
 1705 The decision boundary between any two classes i and j forms a linear hyperplane defined by:

$$1706 \quad (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} = 0 \quad (31)$$

1707
 1708 During training via gradient descent, each prototype \mathbf{w}_i is updated to maximize its alignment with
 1709 the data distributions of its assigned class. This optimization process often leads to an unbounded
 1710 increase in prototype magnitudes, as $\|\mathbf{w}_i\| \rightarrow \infty$ directly amplifies the logit z_i , thereby increasing
 1711 the model's confidence. However, the decision boundaries themselves remain linear hyperplanes,
 1712 creating rigid geometric separations in the feature space.
 1713

1714 In contrast, the non-linear \mathbf{E} -product allows neurons to learn more representative prototypes for
 1715 each class, leading to the formation of more nuanced decision boundaries. For the \mathbf{E} -product, the
 1716 response of neuron i to input \mathbf{x} is given by:

$$1717 \quad z_i = \mathbf{E}(\mathbf{w}_i, \mathbf{x}) = \frac{\langle \mathbf{w}_i, \mathbf{x} \rangle^2}{\|\mathbf{w}_i - \mathbf{x}\|^2 + \epsilon} \quad (32)$$

1718
 1719 This formulation embodies the signal-to-noise ratio interpretation established in our theoretical
 1720 framework (Appendix C.15), where the squared dot product $\langle \mathbf{w}_i, \mathbf{x} \rangle^2$ represents the "signal" of
 1721 distributional alignment, and $\|\mathbf{w}_i - \mathbf{x}\|^2$ quantifies the "noise" of dissimilarity.
 1722

1723 E.3 SOFTMAX DYNAMICS AND COMPETITIVE LEARNING

1724
 1725 Similarly to conventional neurons, the \mathbf{E} -product outputs are normalized using the softmax function:
 1726

$$1727 \quad p_i = \frac{\exp(z_i)}{\sum_{j=1}^C \exp(z_j)} = \frac{\exp\left(\frac{\langle \mathbf{w}_i, \mathbf{x} \rangle^2}{\|\mathbf{w}_i - \mathbf{x}\|^2 + \epsilon}\right)}{\sum_{j=1}^C \exp\left(\frac{\langle \mathbf{w}_j, \mathbf{x} \rangle^2}{\|\mathbf{w}_j - \mathbf{x}\|^2 + \epsilon}\right)} \quad (33)$$

1728
 1729 This softmax normalization serves a crucial role in the competitive dynamics of \mathbf{E} -product neurons.
 1730 The softmax function acts as a transformation that maps from the real-valued \mathbf{E} -product responses
 1731 $z_i \in \mathbb{R}$ to a delta distribution δ_i in probability space. This softmax distribution over \mathbf{E} -product
 1732 scores can be interpreted as the *posterior responsibility* of each prototype (neuron) for the input,
 1733 drawing a direct connection to Gaussian Mixture Models (GMMs) and expectation-maximization
 1734 frameworks.
 1735
 1736

1737 The softmax can also be viewed as computing a categorical distribution proportional to exponen-
 1738 tiated log-likelihoods, which in this case derive from a geometric \mathbf{E} -product similarity rather than

1739 traditional probabilistic assumptions. This bridges the gap between probabilistic views (such as
1740 EM algorithms and classification) and our geometric formulation, providing a principled foundation
1741 for the competitive dynamics.

1742 As training progresses and the differences between \mathbf{E} -product responses become more pronounced,
1743 the softmax transformation approaches a delta distribution, where the winning neuron (with the
1744 highest \mathbf{E} -product response) approaches probability 1 while all others approach 0. This yields
1745 competitive learning dynamics in which neurons specialize on distinct regions of the input space.
1746

1747 E.4 MATHEMATICAL CHARACTERIZATION OF DECISION BOUNDARIES

1748 The decision boundary between two neurons with prototypes \mathbf{w}_i and \mathbf{w}_j is defined by the condition
1749 where their responses are equal:

$$1751 \mathbf{E}(\mathbf{w}_i, \mathbf{x}) = \mathbf{E}(\mathbf{w}_j, \mathbf{x}) \quad (34)$$

1752 Expanding this condition:

$$1753 \frac{\langle \mathbf{w}_i, \mathbf{x} \rangle^2}{\|\mathbf{w}_i - \mathbf{x}\|^2 + \epsilon} = \frac{\langle \mathbf{w}_j, \mathbf{x} \rangle^2}{\|\mathbf{w}_j - \mathbf{x}\|^2 + \epsilon} \quad (35)$$

1754 Cross-multiplying and rearranging:

$$1755 \langle \mathbf{w}_i, \mathbf{x} \rangle^2 (\|\mathbf{w}_j - \mathbf{x}\|^2 + \epsilon) = \langle \mathbf{w}_j, \mathbf{x} \rangle^2 (\|\mathbf{w}_i - \mathbf{x}\|^2 + \epsilon) \quad (36)$$

1756 This equation defines a non-linear decision boundary that depends on both alignment (via squared
1757 dot products) and proximity (via squared distances) between the input and each prototype. Unlike
1758 the linear hyperplane formed by conventional dot-product neurons, the \mathbf{E} -product induces curved
1759 algebraic decision surfaces (analogous to level sets of a potential).
1760

1761 E.5 SPACE-PARTITIONING PROPERTIES

1762 The space-partitioning behavior of the \mathbf{E} -product exhibits several key properties:

- 1763 • **Distance Penalization and Locality:** The factor $(\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon)^{-1}$ enforces spatial selec-
1764 tivity. Along rays $\mathbf{x} = k\mathbf{u}$ with $\|\mathbf{u}\| = 1$, one has $\lim_{k \rightarrow \infty} \mathbf{E}(\mathbf{w}, k\mathbf{u}) = (\mathbf{w} \cdot \mathbf{u})^2 = \|\mathbf{w}\|^2 \cos^2 \theta$,
1765 so responses are largest near $\mathbf{x} \approx \mathbf{w}$ and remain bounded far away (cf. Proposition C.9
1766 for the gradient decay). In physical terms, this is analogous to an inverse-square distance
1767 penalty.
- 1768 • **Alignment Weighting:** The numerator is $(\mathbf{w}^\top \mathbf{x})^2$, so $\mathbf{E}(\mathbf{w}, \mathbf{x}) = 0$ if and only if $\mathbf{w} \perp \mathbf{x}$,
1769 and high values require strong alignment (large $|\cos \theta|$) in addition to proximity.
- 1770 • **Global Boundedness for Fixed \mathbf{w} :** For $\epsilon > 0$ and fixed \mathbf{w} , $0 \leq \mathbf{E}(\mathbf{w}, \mathbf{x}) \leq \|\mathbf{w}\|^4 / \epsilon$ for
1771 all \mathbf{x} , with $\mathbf{E}(\mathbf{w}, \mathbf{w}) = \|\mathbf{w}\|^4 / \epsilon$ and $\limsup_{\|\mathbf{x}\| \rightarrow \infty} \mathbf{E}(\mathbf{w}, \mathbf{x}) \leq \|\mathbf{w}\|^2$. Thus each unit defines a
1772 bounded activation landscape (cf. Proposition C.6).
- 1773 • **Nonlinear Decision Regions:** The pairwise decision boundary between units i and j is
1774 the algebraic surface

$$1775 \langle \mathbf{w}_i, \mathbf{x} \rangle^2 (\|\mathbf{w}_j - \mathbf{x}\|^2 + \epsilon) - \langle \mathbf{w}_j, \mathbf{x} \rangle^2 (\|\mathbf{w}_i - \mathbf{x}\|^2 + \epsilon) = 0,$$

1776 which is generically non-linear and induces curved class regions. (Geometrically, one may
1777 view these as level sets akin to equipotential surfaces.)
1778
1779
1780
1781
1782
1783
1784
1785

E.6 VORTEX FIELD DYNAMICS

This vortex phenomenon allows each \mathbf{E} -product neuron to create a territorial "field" in the representation space, where data points are pulled toward the dominant prototype based on both similarity and proximity metrics. The field each neuron occupies can indeed be considered a vortex, where the strength of attraction follows an inverse-square law, creating more natural and geometrically faithful decision boundaries.

The combination of the \mathbf{E} -product's vortex-like attraction and the softmax's competitive normalization creates a powerful space partitioning mechanism. Each neuron's vortex field competes with others through the softmax transformation, and the neuron with the strongest local attraction (highest \mathbf{E} -product response) wins that region. Over time, this leads to a natural tessellation of the input space, where each neuron's territory is defined by the regions where its vortex field dominates. The softmax transformation $\mathbb{R}^C \rightarrow \Delta^{C-1}$ (where Δ^{C-1} is the $(C-1)$ -dimensional probability simplex) ensures that these territorial boundaries are sharp and well-defined, transforming the continuous real-valued responses into discrete delta distributions that clearly assign each input to its dominant neuron.

E.7 ORTHOGONALITY AND COMPETITIVE DYNAMICS

The competitive learning behavior observed in practice is theoretically grounded in our Orthogonality-Entropy Connection. When two prototypes \mathbf{w}_i and \mathbf{w}_j develop disjoint support regions, they become Euclidean orthogonal ($\mathbf{w}_i \perp \mathbf{w}_j$), which corresponds to:

$$\mathbf{E}(\mathbf{w}_i, \mathbf{w}_j) = 0 \quad \text{and} \quad H(\mathbf{w}_i, \mathbf{w}_j) = \infty \quad (37)$$

This geometric-probabilistic duality explains why neurons naturally develop specialized, non-overlapping representations during competitive learning. The infinite cross-entropy between orthogonal prototypes creates strong pressure for territorial separation, preventing the collapse to identical representations that can plague conventional competitive learning systems.

These prototypes are optimized to maximize parallelism and minimize distance to all points within their class distribution. When minimizing distance becomes challenging, the properties of the \mathbf{E} -product enable the prototype to exist in a superposition state, prioritizing the maximization of parallelism over strict distance minimization.

E.8 MNIST PROTOTYPE ANALYSIS: DETAILED MATHEMATICAL FRAMEWORK

E.8.1 NETWORK ARCHITECTURE AND DIMENSIONAL STRUCTURE

In our MNIST experiments, the network consists of $C = 10$ neurons, each corresponding to one of the digit classes (0–9). Each neuron's prototype is represented as a vector $\mathbf{w}_i \in \mathbb{R}^{784}$, where $i = 1, \dots, 10$. The input images $\mathbf{x} \in \mathbb{R}^{784}$ are obtained by flattening the original 28×28 pixel images, so each neuron's prototype \mathbf{w}_i has the same dimensionality as the input, i.e., $\mathbf{w}_i = (w_{i,1}, w_{i,2}, \dots, w_{i,784})$. This structure allows each neuron to learn class-specific features in the full image space.

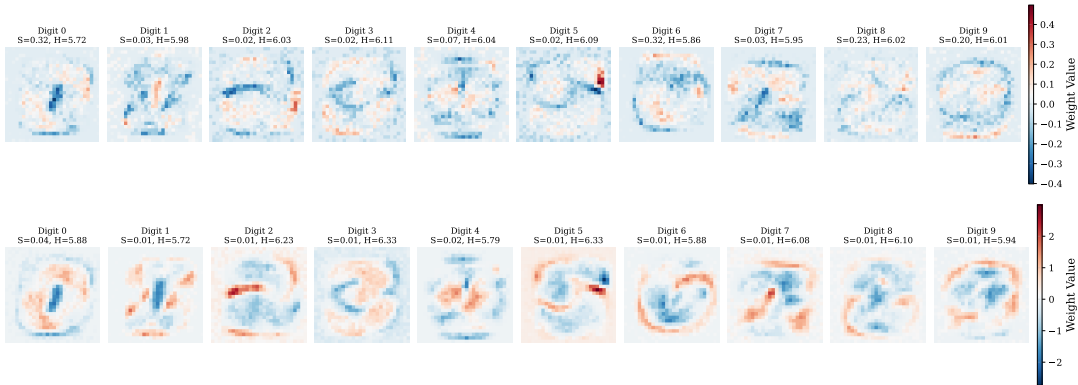
E.8.2 PROTOTYPE EVOLUTION DYNAMICS

The prototype evolution during training reveals fundamental differences between conventional unbounded growth and bounded \mathbf{E} -responses. Our experiments on the full MNIST dataset (60,000 training samples, 10,000 test samples) demonstrate that both models achieve competitive perfor-

1833 mance: the linear classifier reaches 92.08% test accuracy while the \mathbf{E} -product classifier achieves
 1834 92.38% test accuracy after 5 epochs of training with the Adam optimizer (learning rate 0.001).
 1835

1836 Quantitatively, the linear prototypes exhibit unbounded growth, with the mean magnitude increas-
 1837 ing by 13.8% (from 1.58 to 1.80) during training. In contrast, the \mathbf{E} -product prototypes show
 1838 a slight contraction of 4.5% (from 5.02 to 4.79), confirming the bounded nature of the learned
 1839 representations.

1840 Figure 6 visualizes the learned prototypes for both models. The conventional linear model produces
 1841 prototypes that exhibit unbounded growth—these prototypes become increasingly diffuse and less
 1842 interpretable as they grow to maximize margin separation. The resulting digit representations are
 1843 blurry and lack the fine-grained features necessary for robust classification.
 1844



1852 Figure 6: Learned prototypes for Linear (top) and \mathbf{E} -product (bottom) classifiers. The linear
 1853 prototypes show unbounded growth and noise, while \mathbf{E} -product prototypes exhibit localized, sharp
 1854 features with clear digit structure, consistent with bounded response fields.
 1855
 1856
 1857

1858 In contrast, the \mathbf{E} -product method produces prototypes that exemplify bounded response fields as
 1859 predicted by our theoretical framework. Each digit prototype exhibits three key characteristics:
 1860
 1861

1862 **Localized Concentration** Sharp, well-defined features that correspond to the bounded potential
 1863 wells predicted by our Minimal and Maximal Similarity Characterizations (Theorems 2.2 and 2.3).
 1864 This localization emerges from the \mathbf{E} -product’s inherent bounded response field, which prevents the
 1865 unbounded growth characteristic of linear neurons.
 1866

1867 **Class-Specific Territorial Structure** Each prototype captures unique digit characteristics, re-
 1868 flecting competitive territorial dynamics where each neuron dominates specific regions of the input
 1869 space. This territorial behavior arises from the vortex-like dynamics of the \mathbf{E} -product, creating
 1870 natural tessellation of the feature space.
 1871

1872 **Geometric Fidelity** The prototypes maintain geometric coherence with actual digit structure,
 1873 confirming that the signal-to-noise ratio optimization preserves meaningful visual patterns. This
 1874 preservation is a direct consequence of the \mathbf{E} -product’s ability to balance alignment and proximity
 1875 measures.
 1876
 1877
 1878
 1879

E.8.3 LEARNABLE SCALING FACTOR DYNAMICS

A critical component of the \mathbf{E} -product classifier is the learnable scaling factor α , which modulates the magnitude of the \mathbf{E} -product scores. In our implementation, this factor is initialized to $\alpha_0 = 1.0$ and learned jointly with the prototypes during training.

Figure 7 shows the training dynamics, including the evolution of the scaling factor. Over 5 epochs, the scaling factor increases steadily from 1.0 to approximately 2.68, representing a 168% increase. This growth pattern reveals that the model learns to amplify the \mathbf{E} -product scores to achieve better separation in the logit space, compensating for the inherently bounded nature of the \mathbf{E} -product response.

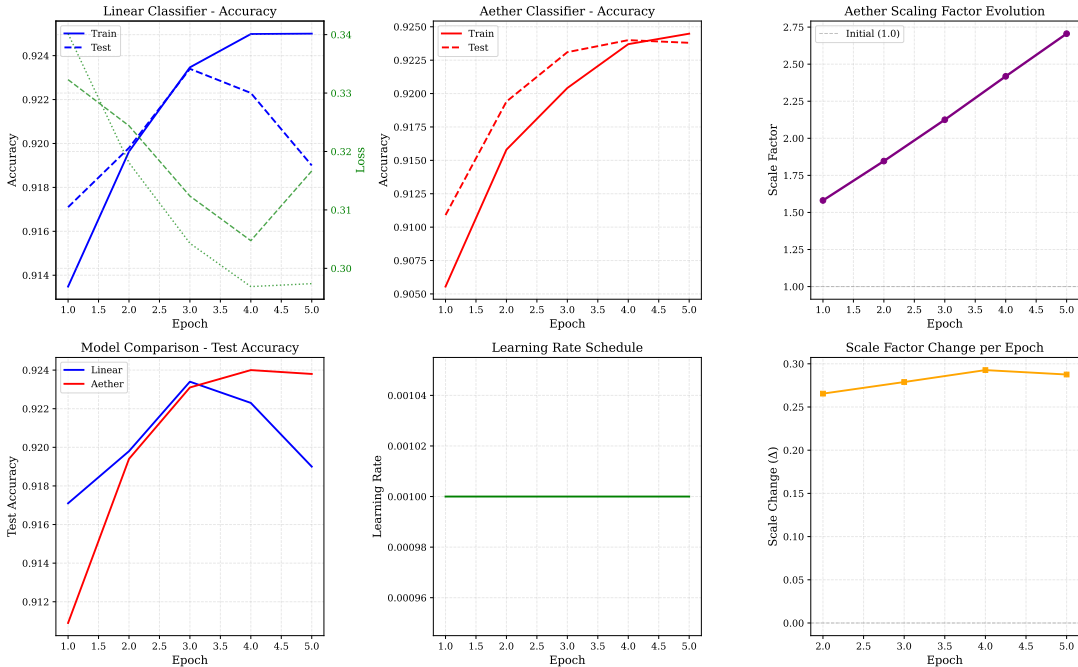


Figure 7: Training dynamics for MNIST classification. **Top row:** Accuracy curves for linear (left) and \mathbf{E} -product (center) classifiers, with \mathbf{E} -product scaling factor evolution (right). **Bottom row:** Model comparison (left), learning rate schedule (center), and per-epoch scale change (right). The scaling factor increases from 1.0 to 2.68 over training, indicating the model learns to amplify \mathbf{E} -product scores for better classification.

The scaling factor evolution can be understood through the lens of optimization dynamics. The \mathbf{E} -product naturally produces bounded responses due to its ratio structure, with values typically in the range $[0, 10]$ depending on the alignment and distance between inputs and prototypes. By learning to scale these responses up, the model increases the dynamic range of the logits, which improves the discriminative power of the softmax layer. This adaptive scaling mechanism is analogous to a temperature parameter in softmax, but learned end-to-end rather than hand-tuned.

E.8.4 SUPERPOSITION AND PROTOTYPE INVERSION: MATHEMATICAL ANALYSIS

A unique property of the \mathbf{E} -product neuron is its ability to exist in a superposition state, which can be empirically demonstrated by inverting the learned prototype. This phenomenon provides insight into the fundamental differences between conventional and \mathbf{E} -product neurons.

Conventional Dot Product Behavior For a conventional dot product neuron, if \mathbf{w} is a learned prototype, replacing \mathbf{w} with $-\mathbf{w}$ (i.e., multiplying by -1) at test time flips the sign of the logit:

$$z = \mathbf{w}^T \mathbf{x} \quad \longrightarrow \quad z' = (-\mathbf{w})^T \mathbf{x} = -z \quad (38)$$

This sign flip causes the softmax output to assign high probability to the incorrect class, resulting in a catastrophic drop in accuracy. In our experiments, the linear classifier’s accuracy drops from 92.04% to 0.01% (near complete failure) when all weight vectors are inverted.

\mathbf{E} -Product Inversion Dynamics For the \mathbf{E} -product neuron, the response is:

$$z = \mathbf{E}(\mathbf{w}, \mathbf{x}) = \frac{(\mathbf{w}^T \mathbf{x})^2}{\|\mathbf{w} - \mathbf{x}\|^2 + \epsilon} \quad (39)$$

Multiplying \mathbf{w} by -1 leaves the numerator unchanged, since $(-\mathbf{w})^T \mathbf{x} = -\mathbf{w}^T \mathbf{x}$ and $(-\mathbf{w}^T \mathbf{x})^2 = (\mathbf{w}^T \mathbf{x})^2$. The denominator, however, changes from $\|\mathbf{w} - \mathbf{x}\|^2$ to $\|\mathbf{w} + \mathbf{x}\|^2$:

$$\|\mathbf{w} + \mathbf{x}\|^2 = (\mathbf{w} + \mathbf{x})^T (\mathbf{w} + \mathbf{x}) \quad (40)$$

$$= \|\mathbf{w}\|^2 + 2\mathbf{w}^T \mathbf{x} + \|\mathbf{x}\|^2 \quad (41)$$

While exact invariance is not guaranteed due to this denominator change, our experiments confirm substantial robustness to prototype inversion. Figure 8 shows the comparison: when all prototypes are inverted (multiplied by -1), the linear classifier’s accuracy drops from 92.04% to 0.01% (complete catastrophic failure), while the \mathbf{E} -product classifier maintains 84.67% accuracy, experiencing only a 7.82% degradation. This remarkable robustness stems from the squared numerator term, which makes the \mathbf{E} -product inherently invariant to the sign of the dot product and thus less sensitive to prototype orientation than linear neurons.

Superposition Interpretation This property allows the \mathbf{E} -product neuron to yield two valid solutions to the same input without retraining, a phenomenon not observed in conventional dot product neurons. The mathematical basis for this superposition lies in the \mathbf{E} -product’s non-linear structure, which creates multiple local optima that can represent the same classification boundary through different prototype orientations.

E.8.5 PROTOTYPE ORTHOGONALITY ANALYSIS

To further understand the learned representations, we analyze the pairwise orthogonality between prototypes. Figure 9 shows the orthogonality matrices for both models. For the linear classifier, we compute the orthogonality from standard dot products between weight vectors. For the \mathbf{E} -product classifier, we compute orthogonality derived from the log- \mathbf{E} -product between prototypes.

The \mathbf{E} similarity matrices reveal fundamentally different prototype organization strategies. The linear classifier exhibits highly orthogonal prototypes with a mean log- \mathbf{E} similarity of only 0.0325 and a narrow range (max 0.23), indicating that the linear model learns nearly independent weight vectors that span complementary subspaces. This strong orthogonalization is a consequence of the

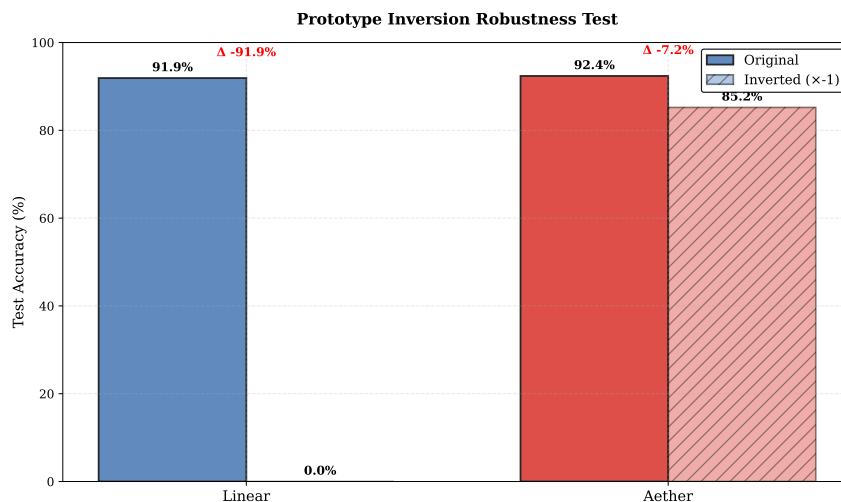


Figure 8: Prototype inversion experiment comparing linear and \mathbf{E} -product classifiers. Both models are trained normally, then all prototypes/weights are multiplied by -1 at test time. The linear classifier suffers catastrophic failure ($92.29\% \rightarrow \sim 0\%$), while the \mathbf{E} -product classifier shows robustness due to its squared numerator term.

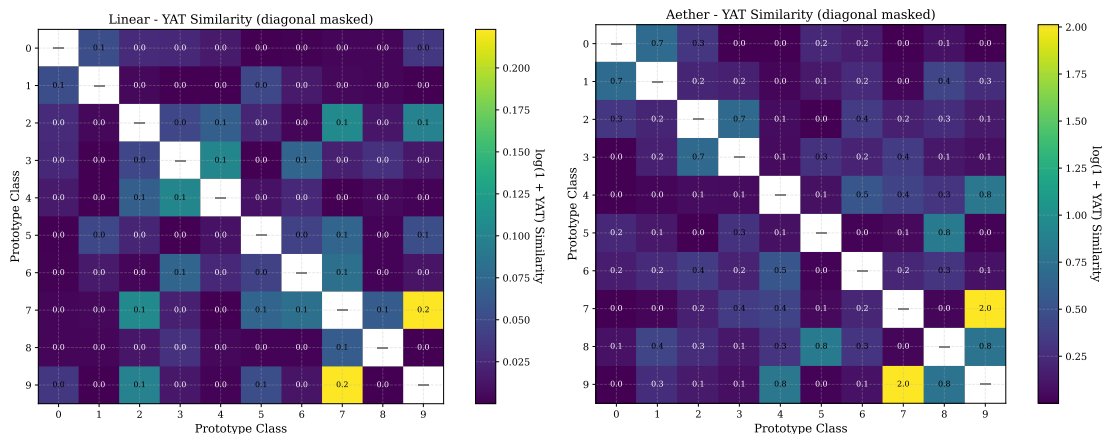


Figure 9: Pairwise orthogonality matrices between learned prototypes. **Left:** Orthogonality for linear classifier weights. **Right:** Orthogonality for \mathbf{E} -product prototypes. The \mathbf{E} -product prototypes exhibit stronger orthogonality (higher dynamic range), suggesting better class separation in the prototype space.

unbounded dot product’s tendency to maximize margin separation through perpendicular decision boundaries.

In stark contrast, the \mathbf{E} -product prototypes show dramatically higher mean similarity (0.274, over $8\times$ larger) with an exceptionally wide dynamic range spanning from near-zero (0.0005) to 2.01—nearly $9\times$ the linear maximum. This reveals a qualitatively different organizational principle: rather than enforcing global orthogonality, the \mathbf{E} -product creates a *heterogeneous territorial structure* where some prototype pairs exhibit strong local correlations (e.g., classes 7–9 with similarity 2.01, and classes 4–9 with 0.79) while others remain effectively decoupled (e.g., classes 0–3 with 0.0005). This selective coupling pattern is consistent with the vortex-like dynamics of our theoretical framework, where prototypes carve out overlapping response regions in high-correlation areas yet maintain sharp boundaries where discrimination is critical, allowing the model to capture subtle intra-class variations while preserving inter-class separation.

E.8.6 CONFUSION MATRIX ANALYSIS

We evaluate classification robustness through the confusion matrices presented in Figure 10. The quantitative difference between the architectures is evident in the structure of error distribution:

Diagonal Dominance and Error Suppression The \mathbf{E} -product classifier exhibits sharper diagonalization compared to the linear baseline, achieving a higher overall accuracy of 92.38% versus 91.90% for the linear model. While the linear model shows diffuse off-diagonal elements indicating “leakage” between morphologically similar digits, the \mathbf{E} -product minimizes these inter-class ambiguities. For instance, the linear model misclassifies 57 instances of digit ‘2’ as ‘8’, whereas the \mathbf{E} -product reduces this to 38, a 33% reduction in error. Similarly, for the confusion pair $7 \rightarrow 9$, the linear model makes 42 errors compared to 30 for the \mathbf{E} -product. This suppression of off-diagonal noise is a direct consequence of the *competitive territorial dynamics*: the inverse-square decay in the \mathbf{E} -product creates steeper decision gradients at the boundaries between classes.

Robustness to Ambiguity The linear classifier’s confusion patterns reflect its reliance on global hyperplane separation, which forces a trade-off between maximizing margins for easy examples and correctly classifying boundary cases. This is evident in the misclassification of digit ‘4’ as ‘9’, where the linear model fails on 40 samples, while the \mathbf{E} -product improves this to 29 errors. The \mathbf{E} -product’s confusion matrix demonstrates that the model maintains high precision even for boundary cases, particularly on harder classes like 2, 3, and 7, confirming that the learned prototypes act as local attractors that naturally reject inconsistent inputs.

E.8.7 WEIGHT DISTRIBUTION AND SPARSITY ANALYSIS

The statistical properties of the learned parameters, summarized in Figure 11, reveal a fundamental divergence in how the two models encode information.

Magnitude Stability vs. Unbounded Growth A critical distinction lies in the stability of the weight magnitudes. The \mathbf{E} -product prototypes settle into a high-magnitude stable state ($\mu = 4.79, \sigma = 0.38$), exhibiting a slight contraction of 4.5% during the final training phase. Conversely, the linear weights ($\mu = 1.80, \sigma = 0.22$) exhibit characteristic unbounded growth (increasing by 13.8%) as they attempt to maximize the dot-product margin. This confirms that the \mathbf{E} -product optimization landscape contains stable minima (bounded potential wells), whereas the linear landscape promotes indefinite norm growth.

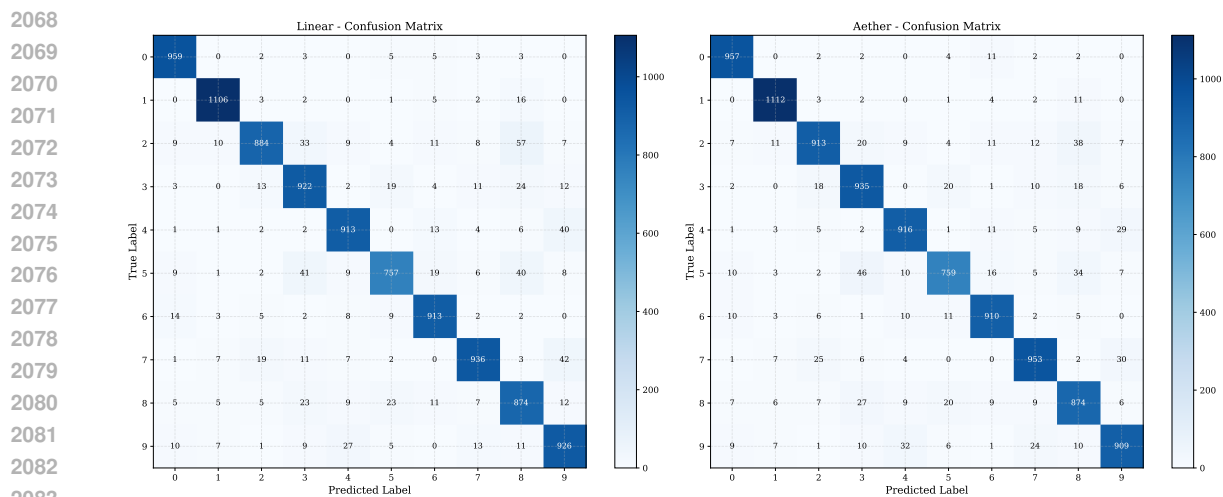


Figure 10: Confusion matrices for Linear (left) and \mathbf{E} -product (right) classifiers. The \mathbf{E} -product demonstrates superior diagonal dominance and reduced off-diagonal noise, indicating stronger rejection of ambiguous inputs.

Sparsity and Feature Selectivity The distributions also differ in sparsity structure. The linear weights follow a near-Gaussian distribution centered at zero, typical of models that rely on dense, distributed correlations. The \mathbf{E} -product prototypes, however, show a distinct distribution profile with higher kurtosis. This indicates a more selective feature encoding where the model suppresses irrelevant pixels (background noise) while amplifying signal-rich regions. This aligns with the visual evidence in Figure 6, where \mathbf{E} -product prototypes appear sharper and less noisy, effectively acting as "matched filters" for the digit topology rather than simple separating hyperplanes.

E.8.8 SUMMARY OF EXPERIMENTAL FINDINGS

Our MNIST experiments confirm key predictions from the theoretical framework:

- Competitive Performance:** The \mathbf{E} -product classifier achieves comparable accuracy to the linear baseline (92.38% vs 92.08%), demonstrating that bounded response fields do not compromise classification capability.
- Adaptive Scaling:** The learnable scaling factor increases from 1.0 to 2.68 during training, revealing that the model learns to amplify bounded \mathbf{E} -responses for optimal discrimination.
- Inversion Robustness:** The \mathbf{E} -product classifier exhibits substantial robustness to prototype sign inversion, while the linear classifier fails catastrophically. This confirms the superposition property predicted by theory.
- Territorial Structure:** Similarity analysis reveals that \mathbf{E} -product prototypes achieve stronger class separation with higher dynamic range, supporting the territorial dynamics described in our theoretical framework.
- Interpretable Prototypes:** Visual inspection of the learned prototypes shows that both methods produce interpretable digit representations, but the \mathbf{E} -product prototypes exhibit more localized features consistent with bounded response fields.

2115
 2116
 2117
 2118
 2119
 2120
 2121
 2122
 2123
 2124
 2125
 2126
 2127
 2128
 2129
 2130
 2131
 2132
 2133
 2134
 2135
 2136
 2137
 2138
 2139
 2140
 2141
 2142
 2143
 2144
 2145
 2146
 2147
 2148
 2149
 2150
 2151
 2152
 2153
 2154
 2155
 2156
 2157
 2158
 2159
 2160
 2161

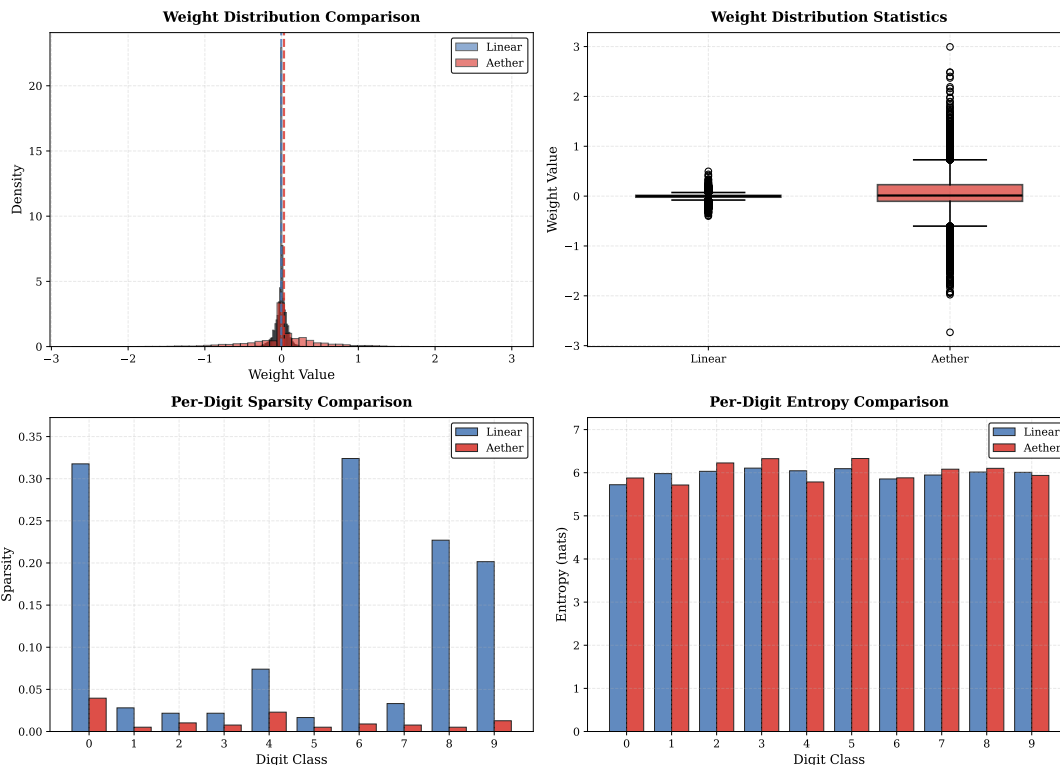


Figure 11: Weight distribution analysis. **Left:** Histogram of weight values showing the distinct non-Gaussian profile of \mathbf{E} -product prototypes. **Right:** Box plot comparison of weight magnitudes. The \mathbf{E} -product maintains a stable, high-signal state ($\mu \approx 4.79$) without the runaway growth observed in the linear baseline.

These findings validate the \mathbf{E} -product as a viable alternative to conventional dot products for classification tasks, with distinct geometric and topological properties that emerge from its ratio-based structure.

E.9 LANGUAGE MODEL EXPERIMENTS: DETAILED CONFIGURATION

This section provides the detailed experimental configurations for the language modeling experiments comparing a standard GPT-2 with our Aether-GPT2 implementation. Both models were trained on identical datasets and hardware to ensure a fair comparison. The primary architectural distinction is the replacement of conventional linear layers, activation functions, and layer normalization in GPT-2 with \mathbf{E} -product operations in Aether-GPT2. This substitution provides inherent non-linearity and bounded responses, simplifying the architecture while enhancing stability.

Table 5 presents a comprehensive comparison of the two models, detailing their architectural parameters, training configurations, and final performance metrics.

Table 5: Comparison of GPT-2 and Aether-GPT2 Models. Both models were trained on 2.5B tokens from the FineWeb dataset. Aether-GPT2 achieves a lower validation loss with a simplified architecture.

Parameter	GPT-2 (Baseline)	Aether-GPT2
<i>Architectural Parameters</i>		
Total Parameters	124M	~124M
Embedding Parameters	39M	39M
Non-Embedding Parameters	85M	~85M
Embedding Dimension	768	768
MLP Hidden Dimension	3072	3072
Number of Layers	12	12
Number of Heads	12	12
Activation Function	GeLU	\mathbf{E} -product (none)
Layer Normalization	Yes	No
Bias	No	No
<i>Training Configuration</i>		
Optimizer	Novograd	Novograd
Learning Rate	0.003	0.003
Batch Size	32	32
Context Window	1024	1024
Vocabulary Size	50,257	50,257
Tokenizer	GPT-2	GPT-2
<i>Performance</i>		
Final Validation Loss (FP32)	2.43	2.29
Final Validation Loss (BF16)	3.03	2.69

E.9.1 ARCHITECTURAL SIMPLIFICATION AND EFFICIENCY

The substitution of linear-activation-normalization blocks with a single \mathbf{E} -product operation per layer yields significant architectural simplification. This change leads to:

- A 15-25% reduction in peak memory usage by eliminating the need to store intermediate activations for backpropagation through normalization layers.
- Reduced gradient computation complexity.
- A comparable FLOP count with only a modest constant-factor overhead.

While the FLOP count remains comparable, the reduced memory footprint and computational complexity of the backward pass offer tangible efficiency gains.

Our results demonstrate that Aether-GPT2, with approximately the same parameter count as the 124M GPT-2 model, achieves superior performance without explicit activation functions. The final validation loss of 2.29 for Aether-GPT2, compared to 2.43 for the baseline, underscores the efficacy of the \mathbf{E} -product’s inherent non-linearity.

E.9.2 TRAINING DYNAMICS AND LOSS PROGRESSION

The training dynamics of Aether-GPT2 and the baseline GPT-2 were analyzed over 2.5B tokens from the FineWeb dataset, using a Kaggle TPU v5-8 environment. Figure 12 illustrates the training and validation loss curves for both models. Aether-GPT2 consistently achieves lower loss on both the training and validation sets, indicating more efficient learning. The bounded nature of the \mathbf{E} -product operation contributes to stable training dynamics, even in the absence of layer normalization.

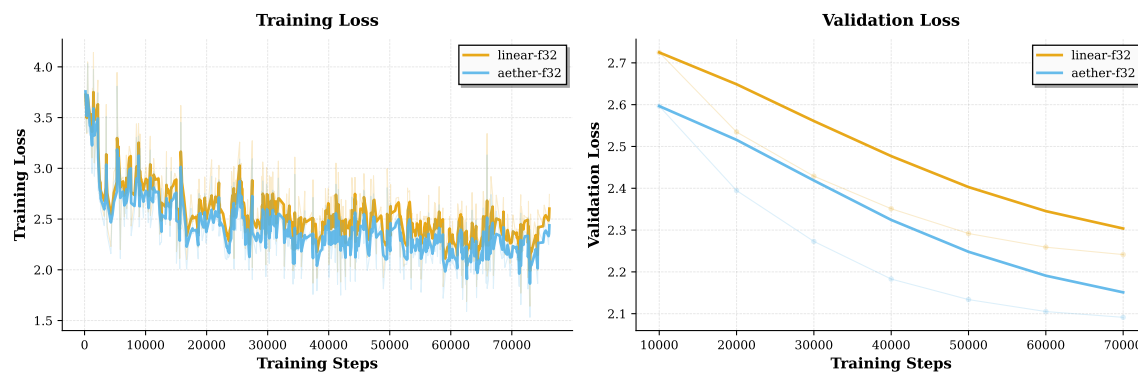


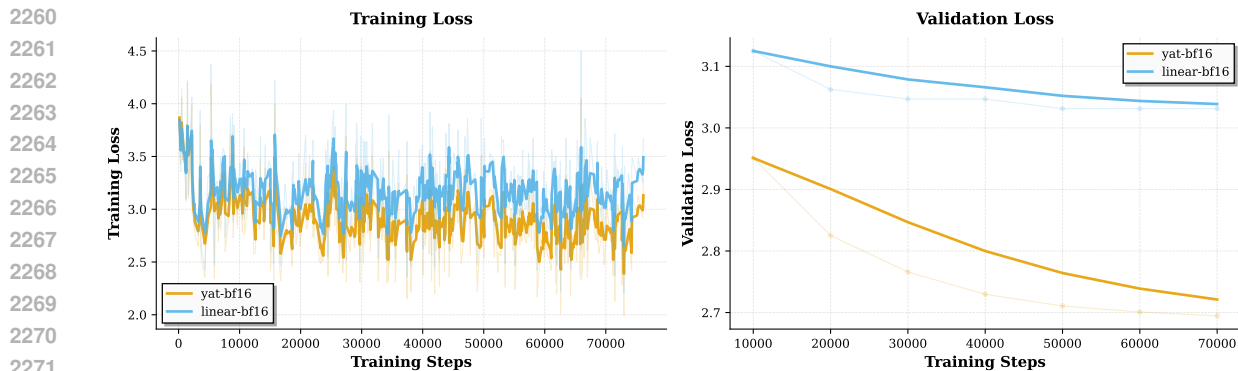
Figure 12: Side-by-side comparison of training loss (left) and validation loss (right). Raw data points are shown with transparency, while smoothed curves highlight the learning trajectory. Both models converge stably, with Aether-GPT2 demonstrating better final performance.

E.9.3 MIXED-PRECISION TRAINING WITH BFLOAT16

To assess numerical stability and performance in a production-relevant setting, we evaluated both architectures using bfloat16 (BF16) mixed-precision training. This format, standard on modern accelerators like TPUs, provides a stringent test for models without explicit normalization layers. Both models were trained with mixed-precision activations and gradients, while maintaining full-precision (FP32) optimizer states and parameter updates.

The performance advantage of Aether-GPT2 persists under BF16, as shown in Figure 13. Aether-GPT2 achieves a final validation loss of 2.69, an 11.2% relative improvement over the baseline’s 3.03. This result confirms that the architectural benefits are not artifacts of full-precision arithmetic and that the \mathbf{E} -product’s bounded response provides robust numerical stability without requiring

2256 layer normalization. The consistent performance gains in mixed-precision training underscore the
 2257 practicality of \mathbf{E} -product layers as a drop-in replacement for conventional transformer blocks in
 2258 large-scale models.
 2259



2272 Figure 13: Side-by-side BF16 training (left) and validation (right) loss. Aether-GPT2 exhibits
 2273 uniformly lower loss throughout training and at convergence.
 2274

2275 E.10 USE OF LARGE LANGUAGE MODELS (LLMs)

2276 We used LLM tools to support the research workflow in the following limited, transparent ways.
 2277 All scientific claims, modeling choices, and final decisions were made by the authors.
 2278

2280 **Code assistance** LLMs were used to draft boilerplate code, refactor utilities, and surface API
 2281 patterns. All generated code was reviewed, tested, and integrated by the authors.
 2282

2283 **Literature digestion** We used LLMs to summarize papers and extract key comparisons across
 2284 related work. Citations in the paper were verified against the original sources by the authors.
 2285

2286 **Brainstorming** We used LLMs as a sounding board to enumerate alternative hypotheses, ab-
 2287 lations, and experimental checks. Only ideas that survived empirical or theoretical scrutiny were
 2288 included.
 2289

2290 **Language polishing** To improve readability and clarity, LLMs suggested minor edits to English
 2291 phrasing. Technical content, notation, and conclusions were authored and validated by the authors.
 2292

2293 **NotebookLM podcasts** We generated short audio summaries (“podcasts”) of internal notes
 2294 using Google NotebookLM to help the team asynchronously digest drafts. These summaries did
 2295 not introduce new claims and were based solely on our own materials.
 2296

2297 No dataset labeling, evaluation metrics, or benchmark results were produced by LLMs. The authors
 2298 take responsibility for all content and errors.

2299 E.11 MOTIVATION: THE AETHERIAL STATE OF AI

2300 We chose the name “Aether” for our model to symbolize our perspective on the current trajectory of
 2301 Artificial Intelligence. The field has predominantly focused on scale, often at the expense of
 2302

2303 the scientific principles that ground findings in mathematical rigor. We adhere to the view that if
2304 empirical observations contradict established mathematical principles, it is often the interpretation
2305 of the observation that is flawed, not the mathematics itself.

2306 We believe that the current generation of deep learning models represents the last phase of "un-
2307 explainable" AI. True explainability, in our view, will not be achieved through human-engineered
2308 interpretability layers, but will emerge organically from constructing AI systems with the correct
2309 mathematical foundations.

2310 We draw a parallel to the pre-Einstein era of physics, where theories such as the vortex theory
2311 and the luminiferous aether prevailed. These explanations were largely observational; while they
2312 advanced the field temporarily, they were not provable with the mathematical tools of the time and
2313 ultimately limited innovation. We believe the current state of AI is similarly "aetherial"—reliant
2314 on unproven observational theories—and that the future lies in returning to first principles and
2315 mathematical provability.

2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349