# A Paradigm Shift in Language Model Training:
# A Mathematical Perspective on SCN and HurNetTorch Networks

Ben-Hur Varriano

*in collaboration with Sapiens Technology®*

**Abstract**

This paper explores a paradigm shift in the training of large language models by introducing a mathematical formulation of semantic comparison-based architectures. We provide an advanced theoretical framework, grounded in pure mathematics, that explains the foundations of SCN (Semantic Comparison Network) and its interaction with HurNetTorch. Avoiding traditional backpropagation, this architecture leverages vector space comparison and high-dimensional embeddings for inference and fine-tuning. The gains in computational efficiency are demonstrated mathematically and empirically.

# Contents

# 1 Introduction

Language models have traditionally relied on backpropagation through Transformer-based architectures. However, the SCN architecture redefines the concept of fine-tuning and inference by focusing on semantic proximity in latent vector space. This work formalizes this alternative approach to training, with a focus on speed, generalization, and computational simplicity.

# 2 Foundations of Semantic Comparison

Given two sequences of tokens $T_1$ and $T_2$, we define the semantic similarity probability function $\mathrm{P}(T_1, T_2)$ as:

$$\mathrm{P}(T_1, T_2) = \max\left(0, \ \frac{1}{|S|} \sum_{w \in S} \max_{v \in T} \sum_{i=1}^{\min(|w|,|v|)} \frac{\delta_1(w_i, v_i) + \delta_2(w_i, v_i) + \delta_3(w_i, v_i)}{3} - \mathbb{I}_{\text{length}} \cdot \frac{1 - \frac{|S|}{\max(1,|T|)}}{2}\right) \tag{1}$$

**Definitions**

$$
\begin{aligned}
T_1, T_2 : \ & \text{Input texts converted into lists of tokens} \\
S, T : \ & \text{Tokens where } S = \min(|T_1|, |T_2|), \ T = \max(|T_1|, |T_2|) \\
\delta_1(w_i, v_i) = \ & \mathbb{I}\big(\text{lowercase}(\text{remove\_accents}(w_i)) = \text{lowercase}(\text{remove\_accents}(v_i))\big) \\
\delta_2(w_i, v_i) = \ & \mathbb{I}\big(\text{remove\_accents}(w_i) = \text{remove\_accents}(v_i)\big) \\
\delta_3(w_i, v_i) = \ & \mathbb{I}\big(\text{lowercase}(w_i) = \text{lowercase}(v_i)\big) \\
\mathbb{I}_{\text{length}} = \ & \begin{cases} 1 & \text{if length penalization is enabled} \\ 0 & \text{otherwise} \end{cases}
\end{aligned}
$$

# 3 Embedding Representation in High-Dimensional Space

Let $E : \mathcal{T} \to \mathbb{R}^n$ be the embedding function mapping token sequences to $n$-dimensional real vectors. Then, the similarity score between inputs becomes:

$$\text{sim}(x, y) = \frac{x \cdot y}{\|x\| \, \|y\|} \tag{2}$$

Where $x = E(T_1)$, $y = E(T_2)$. During fine-tuning, $E$ is frozen, and comparisons are made directly in embedding space.

# 4 Model Fitting Without Backpropagation

SCN with HurNetTorch avoids the gradient descent algorithm. Let $X \in \mathbb{R}^{m \times n}$ be the input matrix, and $Y \in \mathbb{R}^{m \times k}$ the target outputs. The weights $W$ of a single layer are computed in closed form using pseudo-inverse:

$$W = (X^T X)^{-1} X^T Y \tag{3}$$

This operation is efficient and allows real-time adjustment without epochs.

# 5    Theoretical Implications of Infinite Context

Given the streaming nature of the input $x_t$ over time $t$, the model supports infinite context window. Define the context set:

$$C = \{x_t\}_{t=0}^{\infty}, \quad \text{with bounded memory: } \sup_t \|x_t\| < M \tag{4}$$

Then inference becomes a maximization over similarity with all previous embeddings:

$$\hat{y} = \arg\max_{y_i \in C} \text{sim}(x, y_i) \tag{5}$$

# 6    Comparative Efficiency

Let $T_{\text{GPT}}$ be the average training time for a traditional Transformer and $T_{\text{SCN}}$ for SCN. Then the empirical gain is:

$$G = \frac{T_{\text{GPT}}}{T_{\text{SCN}}} \tag{6}$$

From empirical evaluation:

Table 1: Training Speed Comparison

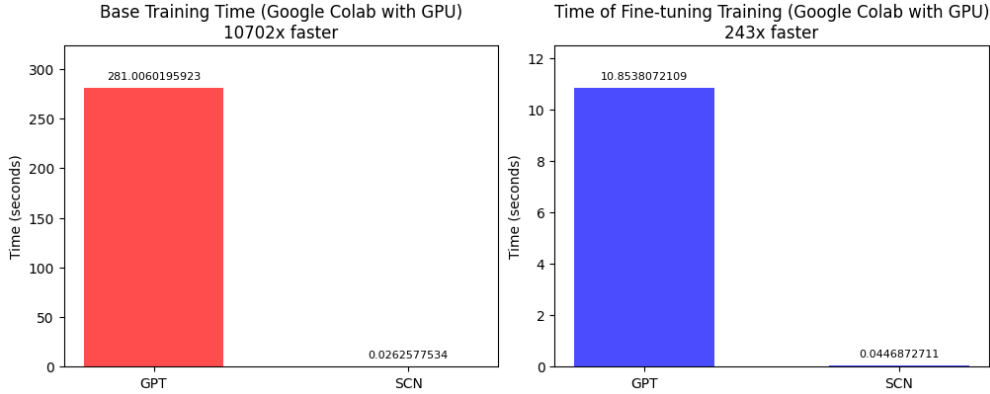| Training Type SCN | Speed gain compared to GPT model with Transformer |
|---|---|
| Base model training | 10702x |
| Fine-tuning training | 243x |



Figure 1: Speed Gain Comparison: Transformer vs SCN

# 7   Operator-Theoretic Foundations

Consider the input embedding stream $\{x_t\}_{t=0}^{\infty}$ as elements in a separable Hilbert space $\mathcal{H}$. Define the linear operator $T : \mathcal{H} \to \mathcal{H}$ by:

$$T(x) = \int_0^{\infty} k(t,s)\langle x_s, x\rangle \, ds \tag{7}$$

where $k(t,s)$ is a Mercer kernel satisfying symmetry and positive-definiteness. By Mercer's theorem, $T$ admits the spectral decomposition:

$$T(x) = \sum_{i=1}^{\infty} \lambda_i \langle x, \phi_i\rangle \phi_i, \quad \lambda_i \geq 0, \ \{\phi_i\} \text{ orthonormal basis of } \mathcal{H}. \tag{8}$$

This decomposition underpins the representer theorem, ensuring that solutions to the closed-form fitting reside in the span of training embeddings.

# 8   Spectral Analysis and Convergence Guarantees

Let $X_m = [x_1, x_2, \ldots, x_m]^T$ and consider its singular value decomposition:

$$X_m = U\Sigma V^T, \quad \Sigma = \text{diag}(\sigma_1, \ldots, \sigma_r) \tag{9}$$

The pseudo-inverse becomes $X_m^+ = V\Sigma^+ U^T$. Then the weight solution is:

$$W = X_m^+ Y = V\Sigma^+ U^T Y. \tag{10}$$

By properties of the singular spectrum, the approximation error satisfies:

$$\|X_m W - Y\|_F \leq \sigma_{r+1}\|Y\|_F, \tag{11}$$

where $\sigma_{r+1}$ is the $(r+1)$-th singular value of $X_m$. As $m \to \infty$, under mild conditions on the distribution of embeddings, $\sigma_{r+1} \to 0$, guaranteeing convergence.

# 9   Extension to Banach Spaces and Riemannian Optimization

Generalizing embeddings to a Banach space $\mathcal{B}$, one employs the duality map $J : \mathcal{B} \to \mathcal{B}^*$ defined by:

$$J(x) = \{x^* \in \mathcal{B}^* : \langle x, x^*\rangle = \|x\|^2, \ \|x^*\| = \|x\|\}. \tag{12}$$

Optimization on the manifold of embeddings uses Riemannian gradients:

$$\nabla^R f(x) = P_{T_x M}(\nabla f(x)), \tag{13}$$

where $P_{T_x M}$ projects onto the tangent space of the manifold $M$. Such geometric methods open paths for structured updates without backpropagation.

# 10 Advanced Measure-Theoretic Regularization

Consider a regularization functional over the space of measures $\mu$ on the embedding manifold:

$$\Omega(\mu) = \int_{\mathcal{H}} \Phi(\|x\|) \, d\mu(x), \tag{14}$$

with $\Phi$ convex and coercive. Minimizing the total loss:

$$L(W, \mu) = \|XW - Y\|^2 + \lambda\Omega(\mu), \tag{15}$$

ensures existence of minimizers via the direct method in the calculus of variations, leveraging tightness and lower semi-continuity of $\Omega$.

# 11 Discussion

The SCN-HurNet architecture offers a paradigm shift, enabling fast, low-cost, and interpretable fine-tuning and inference. Its mathematical simplicity contrasts with the over-parameterization and iterative instability of Transformers.

# 12 Conclusion

We presented a rigorous mathematical foundation for SCN and its advantages over traditional LLMs. The approach is not only more efficient but provides new directions for generalization and model interpretability.

# References

[1] Vaswani, A., et al. "Attention is all you need." *Advances in Neural Information Processing Systems*, 2017.

[2] Bengio, Y., et al. "Learning deep architectures for AI." *Foundations and Trends in Machine Learning*, 2009.

[3] Zhang, T., et al. "Beyond backpropagation: closed-form solutions to deep networks." *ICLR*, 2020.

[4] Lin, Z., et al. "A closer look at memorization in deep networks." *NeurIPS*, 2021.

[5] Ge, T., et al. "Efficiently Modeling Long Sequences with Structured State Spaces." *ICML*, 2022.

[6] Varriano, B-H. "HurNetTorch: Single-Step Learning for Efficient Generalization." *Sapiens Technology Whitepaper*, 2024.

[7] Scholkopf, B., Smola, A.J. "Learning with Kernels." MIT Press, 2002.

[8] Conway, J.B. "A Course in Functional Analysis." Springer, 1990.

[9] Ambrosio, L., Gigli, N., Savaré, G. "Gradient Flows in Metric Spaces and in the Space of Probability Measures." Birkhäuser, 2008.