

HYPERBOLIC IMPLICIT EQUILIBRIUM

Anonymous authors

Paper under double-blind review

ABSTRACT

Euclidean geometry has long dominated neural networks and deep learning, yet neuroscience reveals a different picture. At the representational level, spatial and mnemonic maps in the brain are naturally organized in hyperbolic geometry, supporting efficient hierarchical embeddings. Hyperbolic neural networks exploit this property but remain shallow and costly: explicit architectures must retain all activations, and curvature-induced distortions make stability difficult, leading to prohibitive memory and runtime overhead. At the dynamical level, neural activity tends to converge to stable equilibrium states, conferring robustness, stability, and energy efficiency. Motivated by these complementary principles, we establish Hyperbolic Implicit Equilibrium (HIE), the first implicit equilibrium framework for hyperbolic networks. HIE directly solves for a fixed point and trains via implicit differentiation, requiring only a single Jacobian–vector product. This design enables models of effectively infinite depth within a constant memory footprint, while hyperbolic contraction accelerates convergence beyond Euclidean counterparts. We further contribute Lorentz group normalization for stable equilibrium and a complete theoretical analysis of optimization, stability, and generalization. Experiments show that HIE scales hyperbolic models far beyond prior explicit designs, achieving faster and more robust convergence and revealing the unique benefits of hyperbolic geometry for implicit deep learning.

1 INTRODUCTION

Two thousand years ago, Euclid postulated that through a point not on a line there exists exactly one line parallel to it—a principle that shaped the geometry of human thought for centuries. Yet our brains appear to think differently: neuroscientific evidence (Zhou et al., 2018a; Zhang et al., 2023; Longhena et al., 2025) reveals that neural representations of space and memory naturally form in hyperbolic geometry, where infinitely many parallels can pass through such a point, geodesics diverge rapidly, and volume grows exponentially, enabling hierarchies to be embedded with remarkable efficiency. This divergence between the geometry of intuition and the geometry of cognition is striking. It suggests that if our goal in artificial intelligence is to build systems that learn and reason as effectively and efficiently as the brain, then continuing to rely exclusively on Euclidean geometry may be fundamentally limiting. Hyperbolic neural networks (HNN) have therefore emerged as a pathway to bridge this gap, extending deep learning into negatively curved spaces.

HNN promises compact and faithful representations of inherently structured data, a property less naturally captured in Euclidean space. This is evident in vision and language, where taxonomies in ImageNet, part–whole relations in object recognition, and semantic trees in language all exhibit hierarchies that hyperbolic geometry can embed with exponentially low distortion (Khrukov et al., 2020). However, scaling HNN to the level required for modern tasks remains profoundly challenging. Curvature-induced volume growth and nonlinear gyrovectors operations complicate optimization; explicit architectures must retain intermediate activations to preserve manifold consistency, leading to prohibitive memory costs. As depth increases, distortions accumulate and training becomes numerically fragile. As a result, most hyperbolic models remain shallow. This scalability bottleneck is a structural consequence of the interplay between hyperbolic geometry and optimization dynamics, limiting the expressive power of hyperbolic representations on large-scale datasets.

Several attempts have been made to alleviate these challenges. Poincaré ResNets (Van Spengler et al., 2023) derived custom backward operators to shrink the computational graph, reducing memory at the cost of higher computational complexity. Later works introduced feature clipping and

054 Euclidean reparameterizations (Mishne et al., 2023; Guo et al., 2022; Mathieu et al., 2019), which
055 improved numerical stability and enabled training with 32-bit rather than 64-bit precision (Bdeir
056 et al., 2024), indirectly lowering memory usage. While such techniques improved local stability, the
057 fundamental issue remains: layers are still unrolled sequentially, intermediate activations must be
058 stored, and the compounded cost of hyperbolic operations quickly overwhelms memory and runtime
059 budgets. As noted in Poincaré ResNet (Van Spengler et al., 2023) and HyperbolicCV (Bdeir et al.,
060 2024), the memory footprint of HNN substantially exceeds Euclidean networks. The core tension
061 persists: the same geometry that enables compact hierarchical representations also renders explicit
062 architectures unstable and computationally heavy, leaving the scalability problem unresolved.

063 To overcome this long-standing barrier, we take a different perspective inspired by the brain it-
064 self. Recent neuroscientific studies (Englert et al., 2024; Song et al., 2024) have shown that neural
065 activity often converges to stable attractor states, ensuring efficiency in energy use and memory
066 retrieval. Analogously, we establish Hyperbolic Implicit Equilibrium (HIE), which formalizes the
067 deep equilibrium (DEQ) (Bai et al., 2019; 2020) framework in hyperbolic networks. Instead of
068 propagating through layers explicitly, HIE directly solves for a fixed-point representation using a
069 black-box equilibrium solver. Gradients are computed via implicit differentiation, requiring only
070 a single Jacobian–vector product at equilibrium and bypassing layer-wise backpropagation. This
071 design enables constant-memory training, reduces runtime, and crucially allows hyperbolic models
072 to scale to large datasets without sacrificing efficiency. Moreover, HIE requires fewer iterations than
073 Euclidean counterparts, yielding faster convergence and improved stability. Finally, we establish a
074 complete theoretical analysis of optimization, stability, and generalization that explains why HIE
075 is not only trainable but also more stable, efficient, and generalizable than both explicit HNN and
076 Euclidean DEQ. Our main contributions are as follows:

- 077 1. We establish the implicit equilibrium framework for hyperbolic neural networks, grounded
078 in two principles of brain: hyperbolic representations and equilibrium dynamics;
- 079 2. We design Lorentz group normalization, a geometry-preserving normalization method that
080 maintains manifold consistency while improving the stability of equilibrium solvers;
- 081 3. We present a comprehensive theoretical analysis of optimization, stability, and generaliza-
082 tion, which establishes rigorous convergence and robustness guarantees;
- 083 4. We show that HIE unlocks deep hyperbolic models with constant memory and accelerated
084 training, overcoming the scalability bottleneck of prior hyperbolic neural networks;
- 085 5. We further demonstrate that HIE converges faster to fixed points than Euclidean implicit
086 equilibrium models, leveraging the unique benefits of hyperbolic geometry.

088 2 BACKGROUND

089 2.1 HYPERBOLIC NEURAL NETWORKS

092 Hyperbolic neural networks (HNN) leverage negatively curved spaces to compactly represent hier-
093 archical structures that are distorted in Euclidean embeddings (Nickel & Kiela, 2018; Ganea et al.,
094 2018; Peng et al., 2021). Initially, Euclidean encoders projected features into hyperbolic heads for
095 classification (Khrlukov et al., 2020; Liu et al., 2020; Guo et al., 2022), segmentation (Hsu et al.,
096 2021; Atigh et al., 2022), generation (Mathieu et al., 2019; Nagano et al., 2019; Ovinnikov, 2019),
097 and metric learning (Yan et al., 2021; Ermolov et al., 2022; Yue et al., 2023). Further, hyperbolic
098 formulations generalized fundamental layers—convolutions, attention, and normalization—either in
099 the Poincaré ball (Ganea et al., 2018; Shimizu et al., 2020; Van Spengler et al., 2023) or the Lorentz
100 model (Chen et al., 2021; Fan et al., 2022). Recently, HyperbolicCV (Bdeir et al., 2024) introduced
101 Lorentz CNN, providing missing components such as hyperbolic convolution, batch normalization,
102 and logistic regression.

103 Nevertheless, hyperbolic networks remain computationally heavy. The exponential volume growth
104 of the Lorentz model causes instability; remedies such as feature clipping and Euclidean reparam-
105 eterizations (Mishne et al., 2023; Guo et al., 2022; Mathieu et al., 2019) only improve local stability.
106 Poincaré ResNets (Van Spengler et al., 2023) reduced memory by custom backward operators, but
107 at higher runtime cost. Most prior work has therefore centered on stability rather than fundamentally
addressing the scalability problem.

2.2 IMPLICIT DEEP LEARNING

Implicit deep learning departs from explicit multi-layer architectures by defining hidden states through analytical conditions rather than prescribed computation graphs. This line of work dates back to implicit differentiation for recurrent dynamics (Pineda, 1987; Almeida, 1990), later revisited as recurrent back-propagation (RBP) (Liao et al., 2018). Recent interest has revived implicit formulations across diverse domains (El Ghaoui et al., 2021; Gould et al., 2021). Representative examples include Neural ODE (NODE) (Chen et al., 2018; Dupont et al., 2019), which interpret residual networks as continuous dynamics solved by ODE integrators. Other instantiations range from optimization-based layers (Amos & Kolter, 2017; Djolonga & Krause, 2017), differentiable physics simulators (de Avila Belbute-Peres et al., 2018; Qiao et al., 2020), and logical reasoning modules (Wang et al., 2019), to continuous-time generative models (Grathwohl et al., 2018).

Deep equilibrium models (DEQ) (Bai et al., 2019) solve for fixed points via black-box root-finding, achieving representations equivalent to infinitely deep weight-tied networks while requiring only $O(1)$ memory through implicit differentiation. Multiscale DEQ (MDEQ) (Bai et al., 2020) extended the approach to vision, jointly solving for equilibrium across resolutions and enabling classification and segmentation within one model. Despite these advantages, DEQ exhibit two main drawbacks: (i) convergence relies heavily on contractive dynamics, often unmet in practice, and (ii) runtime scales with the number of root-finding iterations. As noted in MDEQ, iterations are truncated at a threshold, yielding only approximate equilibrium. Our contribution builds on this line by establishing hyperbolic implicit equilibrium (HIE), where negative curvature strengthens contraction, accelerates solver convergence, and thus reduces runtime while preserving constant memory.

3 METHOD

3.1 PRELIMINARIES

Hyperbolic Geometry. Hyperbolic space is a prototypical example of a Riemannian manifold with constant negative curvature (Cannon et al., 1997; Ratcliffe, 2006). Formally, the n -dimensional hyperbolic space \mathbb{H}_K^n can be described as a pair (\mathcal{M}^n, g^K) , where \mathcal{M}^n denotes the underlying manifold and g^K the Riemannian metric associated with curvature $K < 0$. Several equivalent models exist to represent hyperbolic geometry, such as the Poincaré ball and the Lorentz hyperboloid. In this work, we adopt the Lorentz model due to its favorable numerical behavior and the availability of simple closed-form operations, including exponential and logarithmic maps, geodesic distance, and parallel transport. These properties make it particularly suitable for integration with deep learning methods. For completeness, detailed derivations of these operators are provided in Appendix A.

Lorentz Model. The Lorentz model offers a convenient representation of hyperbolic geometry with curvature $K < 0$, as shown in Figure 1. Let \mathbb{R}^{d+1} be equipped with the Lorentzian bilinear form

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^d u_i v_i, \quad (1)$$

where $\mathbf{u} = (u_0, \dots, u_d), \mathbf{v} = (v_0, \dots, v_d) \in \mathbb{R}^{d+1}$. The d -dimensional hyperbolic space of curvature K can then be realized as the upper sheet of the two-sheeted hyperboloid

$$\mathbb{H}_K^d = \left\{ \mathbf{u} \in \mathbb{R}^{d+1} \mid \langle \mathbf{u}, \mathbf{u} \rangle_{\mathcal{L}} = \frac{1}{K}, u_0 > 0 \right\}. \quad (2)$$

The geodesic distance between two points $\mathbf{u}, \mathbf{v} \in \mathbb{H}_K^d$ is given by

$$d_{\mathbb{H}_K}(\mathbf{u}, \mathbf{v}) = \frac{1}{\sqrt{-K}} \operatorname{arcosh}(K \cdot (-\langle \mathbf{u}, \mathbf{v} \rangle_{\mathcal{L}})). \quad (3)$$

For simplicity and consistency with prior work in hyperbolic deep learning, we use curvature $K < 0$ and set $K = -1$ for all experiments. More notations and conventions are provided in Appendix A.

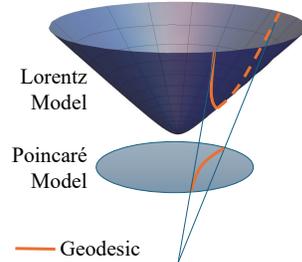


Figure 1: Visualization of Lorentz and Poincaré model.

3.2 HYBRID NETWORK ARCHITECTURE

Following recent work (Bdeir et al., 2024), we adopt a hybrid architecture that integrates Euclidean and hyperbolic components. The encoder combines Euclidean and Lorentzian residual blocks, while the task head is fully hyperbolic. This division balances efficiency and expressivity: Euclidean layers provide cost-effective feature extraction, whereas hyperbolic modules capture hierarchical relations. Empirical studies show that intermediate representations exhibit varying degrees of δ -hyperbolicity (Khrulkov et al., 2020), suggesting that only certain layers benefit from negative curvature. By introducing hyperbolic operators selectively, the model leverages their representational power while avoiding the prohibitive memory and runtime overhead of fully hyperbolic encoders.

To predict when hyperbolic equilibrium models are advantageous, we employ the notion of δ -hyperbolicity (Khrulkov et al., 2020), which quantifies the tree-likeness of a metric space. In practice, we estimate δ from encoder features and compute the normalized score δ_{rel} across scales. Hyperbolic operators are used when representations exhibit strong hyperbolicity, and Euclidean ones otherwise. The computation details are provided in Appendix B.1.

3.3 HYPERBOLIC RESIDUAL BLOCK

Our network integrates several hyperbolic operators into residual blocks. Specifically, we employ Lorentz convolutional layers, a Lorentz multinomial logistic regression (MLR) classifier, and hyperbolic residual connections with non-linear activations. These modules extend their Euclidean counterparts while preserving consistency with the Lorentz model. Convolutions and classifiers operate directly in hyperbolic space, while residual connections and activations are defined via exponential and logarithmic maps between the tangent space and the manifold. Together, they form the building blocks of our hybrid architecture, enabling expressive yet stable learning within curved geometry. Implementation details and closed-form formulas are deferred to Appendix B.2.

Lorentz Group Normalization. Normalization is indispensable in equilibrium models: batch-dependent methods such as BatchNorm are ill-suited, as they inflate the Jacobian norm of f_θ and destabilize implicit solvers (Bai et al., 2019; 2020). MDEQ replaces BatchNorm with GroupNorm (Wu & He, 2018), which normalizes feature subsets within each sample, removing batch-size dependence and improving solver stability. However, directly applying Euclidean GroupNorm to hyperbolic layers is invalid: channel-wise statistics computed in \mathbb{R}^d break manifold constraints and lead to drift off the Lorentz hyperboloid.

To address this, we design *Lorentz Group Normalization (LGN)*, which extends GroupNorm to the Lorentz model \mathbb{H}_K^d while preserving geometric consistency. LGN operates on groups of Lorentz points (e.g., all spatial positions assigned to a feature group within a sample). Design notes and variants are provided in Appendix B.3. For group g with points $S_g = \sum_{j=1}^{N_g} \mathbf{u}_j$, we first compute the Lorentz mean

$$\boldsymbol{\mu}_g = \frac{S_g}{\sqrt{K \langle S_g, S_g \rangle_{\mathcal{L}}}}. \quad (4)$$

which lies on \mathbb{H}_K^d and minimizes the Fréchet energy. We estimate dispersion by the group variance

$$\sigma_g^2 = \frac{1}{N_g} \sum_{j=1}^{N_g} d_{\mathbb{H}_K}(\mathbf{u}_j, \boldsymbol{\mu}_g)^2. \quad (5)$$

where $d_{\mathbb{H}_K}$ is the hyperbolic distance in the Lorentz model.

Each point \mathbf{u}_j is mapped to the tangent space at $\boldsymbol{\mu}_g$ by $\log_{\boldsymbol{\mu}_g}^K$, parallel transported (PT) to the origin \mathbf{o} , rescaled by $\frac{\gamma_g}{\sqrt{\sigma_g^2 + \varepsilon}}$, and transported to a learnable anchor $\boldsymbol{\beta}_g \in \mathbb{H}_K^d$ before mapping back via $\exp_{\boldsymbol{\beta}_g}^K$, which is formulated as

$$\text{LGN}_g(\mathbf{u}_j) = \exp_{\boldsymbol{\beta}_g}^K \left(\text{PT}_{\mathbf{o} \rightarrow \boldsymbol{\beta}_g}^K \left(\frac{\gamma_g}{\sqrt{\sigma_g^2 + \varepsilon}} \text{PT}_{\boldsymbol{\mu}_g \rightarrow \mathbf{o}}^K (\log_{\boldsymbol{\mu}_g}^K(\mathbf{u}_j)) \right) \right). \quad (6)$$

Here γ_g is a learnable scale and $\boldsymbol{\beta}_g$ a learnable bias. All steps—exponential/logarithmic maps, distances, and parallel transport—have closed forms in the Lorentz model, ensuring numerical stability.

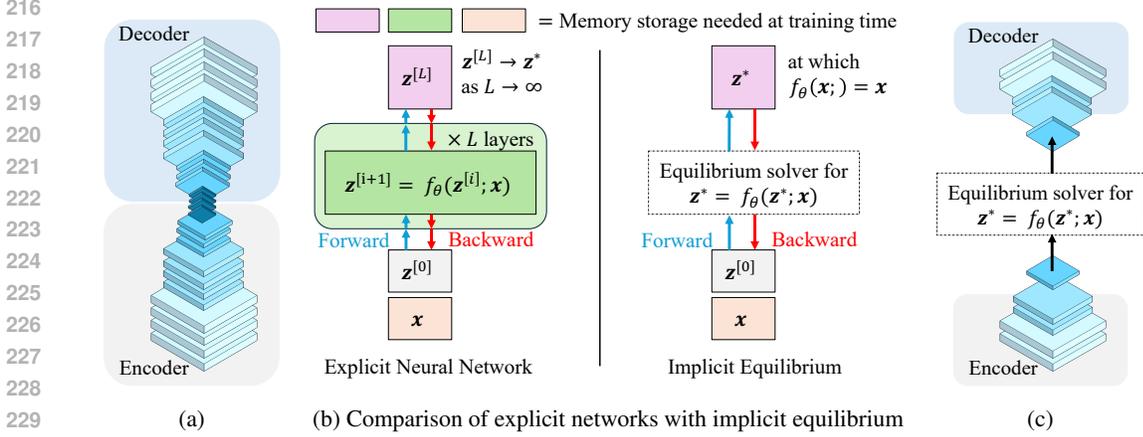


Figure 2: (a) Explicit network. (b) Memory cost: explicit vs. implicit. (c) Deep equilibrium model.

3.4 MULTISCALE IMPLICIT EQUILIBRIUM

Traditional explicit architectures such as ResNet (Figure 2a) unroll a fixed sequence of L layers, requiring storage of all intermediate activations during training. Deep equilibrium models (Figure 2c) replace this explicit stacking with a fixed-point formulation, directly solving for the equilibrium state $z^* = f_\theta(z^*; x)$. As summarized in Figure 2b, this implicit formulation reduces memory consumption from $\mathcal{O}(L)$ to $\mathcal{O}(1)$, since intermediate states need not be stored and gradients can be computed through implicit differentiation.

Building on these principles, we adopt a multiscale equilibrium construction inspired by MDEQ (Bai et al., 2020). Formally, let $z = [z_1, \dots, z_n]$ denote the collection of hidden states at n resolutions, with $z_i \in \mathbb{R}^{H_i \times W_i \times C_i}$. Each resolution evolves through a Lorentz residual block followed by cross-scale fusion, producing updated features \tilde{z}_i . As shown in Figure 3, the joint transformation $f_\theta(z; x)$ then acts on all scales, and the equilibrium state is defined as the solution of

$$z^* = f_\theta(z^*; x), \quad (7)$$

where x denotes the input injected at the highest resolution.

This equilibrium is solved simultaneously across scales using a limited-memory quasi-Newton method (e.g., Broyden updates), yielding synchronized hidden states $\{z_i^*\}_{i=1}^n$. Details of the forward and backward passes of the implicit equilibrium are provided in Appendix C.

4 THEORY

Lemma 1 (Bi-Lipschitz bounds of \exp^K and \log^K on a geodesic ball). *Let $K < 0$, $\kappa = \sqrt{-K}$, and fix $\beta \in \mathbb{H}_K^d$. For $B_R(\beta) = \{z \in \mathbb{H}_K^d : d_{\mathbb{H}_K}(z, \beta) \leq R\}$,*

$$L_{\exp}(R) = \sup_{\|v\| \leq R} \|d(\exp_\beta^K)_v\| = \frac{\sinh(\kappa R)}{\kappa R}, \quad L_{\log}(R) = \sup_{z \in B_R(\beta)} \|d(\log_\beta^K)_z\| = \frac{\kappa R}{\sinh(\kappa R)}.$$

Hence $L_{\exp}(R) \geq 1$, $L_{\log}(R) \leq 1$, and $L_{\exp}(R)L_{\log}(R) = 1$.

Proof sketch. In constant curvature spaces, the singular values of $d \exp_\beta$ along radial/tangential directions are governed by Jacobi fields, yielding $\sinh(\kappa r)/(\kappa r)$ at geodesic radius r ; invertibility gives the $d \log_\beta$ bound and the product identity. Full proof in Appendix D.1. \square

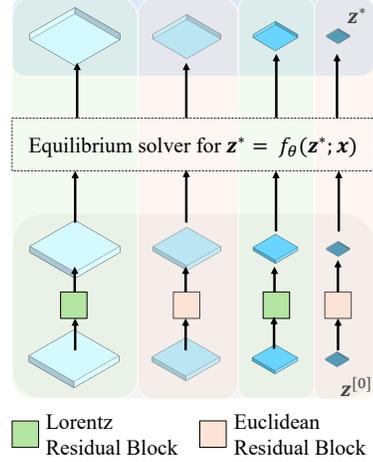


Figure 3: Our hyperbolic implicit equilibrium framework.

Theorem 1 (Contraction with geodesic damping and radial shrinkage). *Let $K < 0$, $\kappa = \sqrt{-K}$, fix $\beta \in \mathbb{H}_K^d$ and $R > 0$. Define $F : B_R(\beta) \rightarrow \mathbb{H}_K^d$ by*

$$F(z) = \exp_\beta^K \left(\tau A(S(u)) + (1 - \tau)u + b(x) \right), \quad u = \log_\beta^K(z), \quad (8)$$

where $0 < \tau \leq 1$, $A : T_\beta \mathbb{H}_K^d \rightarrow T_\beta \mathbb{H}_K^d$ is linear with $\|A\|_2 \leq \alpha$, $b(x) \in T_\beta \mathbb{H}_K^d$, and $S(u) = \psi(\|u\|)u/\|u\|$ (with $S(0) = 0$) is η -Lipschitz on $\{\|u\| \leq R\}$ for some $0 \leq \eta < 1$. Assume $F(B_R(\beta)) \subseteq B_R(\beta)$. Then on $(B_R(\beta), d_{\mathbb{H}_K})$,

$$\text{Lip}(F) \leq L_{\text{exp}}(R) (\tau\alpha\eta + (1 - \tau)) L_{\log}(R) = \tau\alpha\eta + (1 - \tau) := \rho_{\mathbb{H}} < 1. \quad (9)$$

Thus F has a unique fixed point $z^* \in B_R(\beta)$ and Picard iteration converges linearly with rate $\leq \rho_{\mathbb{H}}$.

Proof sketch. For $z_i \in B_R(\beta)$, write $u_i = \log_\beta(z_i)$ and $\Phi(u) = \tau A(S(u)) + (1 - \tau)u + b(x)$. Lemma 1 gives $d_{\mathbb{H}_K}(F(z_1), F(z_2)) \leq L_{\text{exp}}(R)\|\Phi(u_1) - \Phi(u_2)\|$ and $\|u_1 - u_2\| \leq L_{\log}(R) d_{\mathbb{H}_K}(z_1, z_2)$. Using $\|A\| \leq \alpha$ and S being η -Lipschitz yields equation 9. Banach fixed-point theorem concludes. Full proof in Appendix D.2. \square

Corollary 1 (Hyperbolic vs. Euclidean rate). *For $F_{\mathbb{R}}(u) = \tau Au + (1 - \tau)u + b(x)$ on $T_\beta \cong \mathbb{R}^d$, $\rho_{\mathbb{R}} = \tau\alpha + (1 - \tau)$. Under Theorem 1, $\rho_{\mathbb{H}} = \tau\alpha\eta + (1 - \tau) \leq \rho_{\mathbb{R}}$, with strict inequality if $\tau\alpha > 0$ and $\eta < 1$.*

Proof sketch. Immediate from $\eta < 1$ and the expressions of $\rho_{\mathbb{H}}$ and $\rho_{\mathbb{R}}$. Full proof in Appendix D.3. \square

Theorem 2 (Bounded implicit gradients under contraction). *Let z^* be the unique fixed point of F_θ in Theorem 1. All Jacobians and operator/vector norms are taken in $T_{z^*} \mathbb{H}_K^d$ with metric g_{z^*} . For any differentiable loss $\mathcal{L}(z^*)$,*

$$\begin{aligned} \|\nabla_\theta \mathcal{L}\| &\leq \|(I - D_z F_\theta(z^*; x))^{-1}\| \|\partial_\theta F_\theta(z^*; x)\| \|\nabla_z \mathcal{L}(z^*)\| \\ &\leq \frac{1}{1 - \rho_{\mathbb{H}}} \|\partial_\theta F_\theta(z^*; x)\| \|\nabla_z \mathcal{L}(z^*)\|. \end{aligned} \quad (10)$$

Proof sketch. Differentiate $F_\theta(z^*; x) = z^*$ to get $(I - D_z F_\theta) dz^* = (\partial_\theta F_\theta) d\theta$; solve via Neumann series since $\|D_z F_\theta\| \leq \rho_{\mathbb{H}} < 1$. Take norms in T_{z^*} to obtain equation 10. Full proof in Appendix D.4. \square

Lemma 2 (Hyperbolic law of cosines). *Let $\beta \in \mathbb{H}_K^d$, $\kappa = \sqrt{-K}$. For $u, v \in \mathbb{H}_K^d$ with $r_i = d_{\mathbb{H}_K}(u_i, \beta)$ and tangent angle $\theta \in [0, \pi]$ between $\log_\beta^K(u)$ and $\log_\beta^K(v)$,*

$$\cosh(\kappa d_{\mathbb{H}_K}(u, v)) = \cosh(\kappa r_1) \cosh(\kappa r_2) - \sinh(\kappa r_1) \sinh(\kappa r_2) \cos \theta.$$

Proof sketch. Standard form in the Lorentz/Poincaré models via geodesic polar coordinates. Full proof in Appendix D.5. \square

Theorem 3 (Geodesic margin and Euclidean comparison). *Assume class means $\mu_c \in \mathbb{H}_K^d$ share radius $r = d_{\mathbb{H}_K}(\mu_c, \beta)$ and $\theta_0 = \min_{c \neq c'} \angle(\log_\beta^K(\mu_c), \log_\beta^K(\mu_{c'}))$. Then*

$$\gamma_{\mathbb{H}} := \min_{c \neq c'} d_{\mathbb{H}_K}(\mu_c, \mu_{c'}) = \frac{2}{\kappa} \text{arcsinh} \left(\sinh(\kappa r) \sin \frac{\theta_0}{2} \right) \geq 2r \sin \frac{\theta_0}{2} =: \gamma_{\mathbb{R}}, \quad (11)$$

with strict inequality if $r > 0$ and $\theta_0 \in (0, \pi)$.

Proof sketch. Apply Lemma 2 with $r_1 = r_2 = r$, use $\cos \theta = 1 - 2 \sin^2(\theta/2)$ to obtain the closed form. Since $\sinh(\kappa r) \geq \kappa r$ and arcsinh is increasing and concave, $\text{arcsinh}(\sinh(\kappa r)s) \geq s \kappa r$ for $s = \sin(\theta_0/2)$, yielding $\gamma_{\mathbb{H}} \geq \gamma_{\mathbb{R}}$. Full proof in Appendix D.6. \square

Table 1: The hyperbolicity values δ_{rel} calculated for intermediate embeddings on different datasets.

Dataset	CIFAR-10		ImageNet		Cityscapes	
Model	MDEQ-small	MDEQ-large	MDEQ-small	MDEQ-XL	MDEQ-small	MDEQ-XL
Initial Conv.	0.2526	0.2911	0.2230	0.2790	0.3726	0.3599
Branch 1	0.1928	0.2021	0.2022	0.1836	0.3394	0.3747
Branch 2	0.2330	0.3044	0.2319	0.2193	0.3739	0.3494
Branch 3	N/A	0.2966	0.2078	0.1957	0.3941	0.3257
Branch 4	N/A	0.4449	0.2403	0.2377	0.3617	0.3012

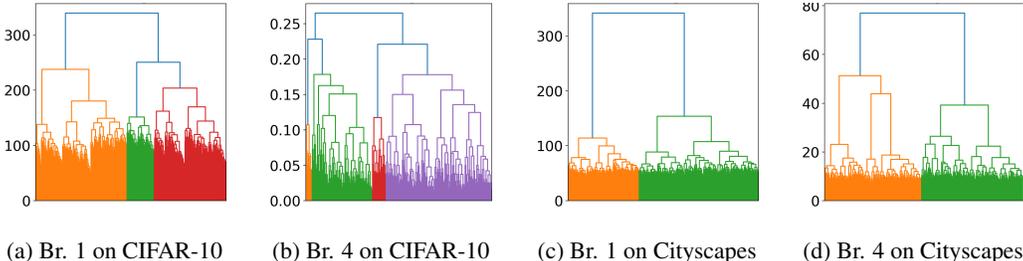


Figure 4: Dendrograms of intermediate embeddings of different branches on different datasets.

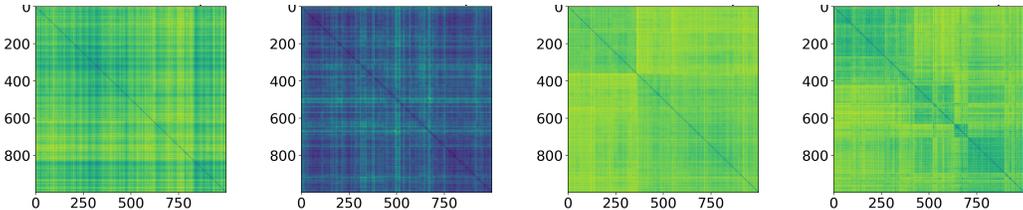


Figure 5: Distance heatmaps of intermediate embeddings of different branches on different datasets.

5 EXPERIMENT

5.1 HYPERBOLICITY

We evaluate the hyperbolicity of pre-trained MDEQ (Bai et al., 2020) models on CIFAR-10 (Krizhevsky et al., 2009), ImageNet (Russakovsky et al., 2015), and Cityscapes (Cordts et al., 2016). Relative hyperbolicity values δ_{rel} are reported in Table 1, where smaller values indicate stronger hyperbolicity and greater suitability for hyperbolic representations. Table 1 shows that, consistent with HyperbolicCV (Bdeir et al., 2024), branches 1 and 3 on CIFAR-10 and ImageNet are more hyperbolic than branches 2 and 4. Interestingly, all branches on Cityscapes are weakly hyperbolic, indicating that classification data encode clearer hierarchical structures than segmentation.

We further visualize embeddings via dendrograms (Figure 4), distance heatmaps (Figure 5), minimum spanning trees (MST) (Figure 6), and t-SNE (Figure 7). On CIFAR-10, Branch 1 consistently shows clearer hierarchy: dendrograms with long branches and distinct clusters, block-structured heatmaps, MST with extended backbones, and well-separated t-SNE clusters. Branch 4 is flatter, with diffuse heatmaps, collapsed MST, and entangled t-SNE distributions. On Cityscapes, all views are flat or uniform, confirming weak hierarchical separation. More visualizations in Appendix F.1.

These findings guide our design: for classification, we follow HyperbolicCV (Bdeir et al., 2024) by placing the head in hyperbolic space and adopting a hybrid encoder (branches 1 and 3 hyperbolic, branches 2 and 4 Euclidean). For segmentation, we find that a hyperbolic head matches Euclidean baselines, but extending hyperbolicity into the encoder degrades performance (Appendix F.3).

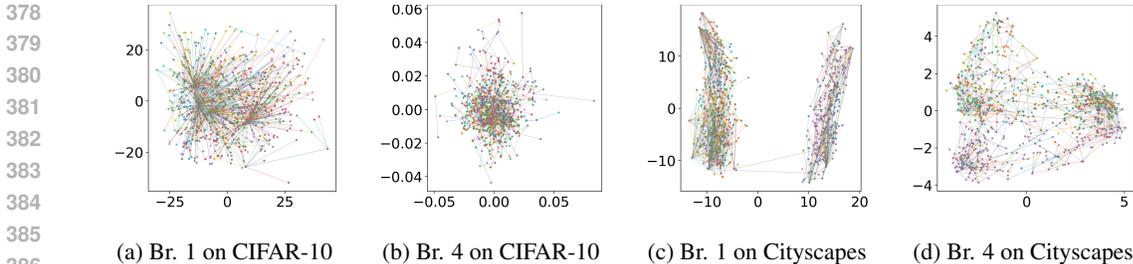


Figure 6: Minimum spanning trees of embeddings of different branches on different datasets.

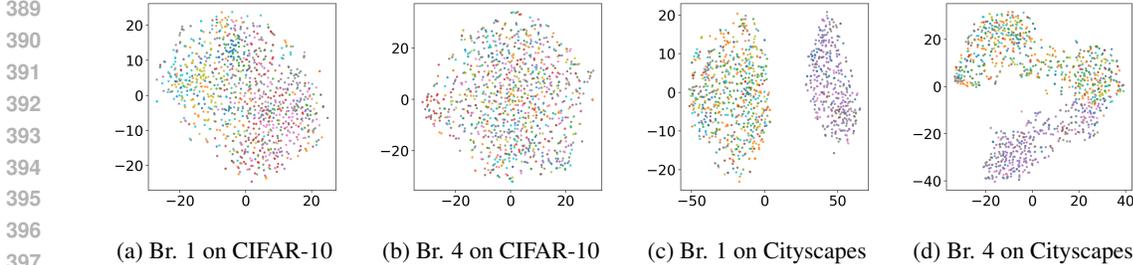


Figure 7: t-SNE projections of intermediate embeddings of different branches on different datasets.

5.2 CLASSIFICATION ON CIFAR-10

Table 2 reports results on CIFAR-10. Without data augmentation, HIE-small surpasses both Euclidean baselines and hyperbolic counterparts, achieving the highest accuracy at comparable parameter counts. With standard augmentation, HIE attains 94.0%, outperforming MDEQ and all prior hyperbolic designs while maintaining the same model size. These results demonstrate that HIE consistently leverages hyperbolic geometry for stronger representation learning in low-data regimes and remains competitive with augmentation. Further details of the tasks, hyperparameters, and training settings are provided in Appendix E.

Table 2: Evaluation on CIFAR-10. Standard deviations are calculated on 5 runs.

CIFAR-10 (<i>without</i> data augmentation)		
Model	Model Size	Accuracy (%)
Neural ODEs (Chen et al., 2018) ()	172K	53.7 ± 0.2
Aug. Neural ODEs (Dupont et al., 2019)	172K	60.6 ± 0.4
Single-stream DEQ (Bai et al., 2019)	170K	82.2 ± 0.3
ResNet-18 (He et al., 2016)	170K	81.6 ± 0.3
MDEQ-small (Bai et al., 2020)	170K	87.1 ± 0.4
HECNN Lorentz (Bdeir et al., 2024)	170K	82.7 ± 0.4
HCNN Lorentz (Bdeir et al., 2024)	170K	81.9 ± 0.3
HIE-small (Ours)	170K	87.8 ± 0.4
CIFAR-10 (<i>with</i> data augmentation)		
ResNet-18 (He et al., 2016)	10M	92.9 ± 0.2
MDEQ (Bai et al., 2020)	10M	93.8 ± 0.3
HyperbolicNN (Ganea et al., 2018)	10M	88.8 ± 0.5
Hybrid Poincaré (Guo et al., 2022)	10M	91.9 ± 0.2
Hybrid Lorentz (Bdeir et al., 2024)	10M	92.7 ± 0.2
Poincaré ResNet (Van Spengler et al., 2023)	10M	92.3 ± 0.4
HECNN Lorentz (Bdeir et al., 2024)	10M	92.9 ± 0.3
HCNN Lorentz (Bdeir et al., 2024)	10M	92.9 ± 0.2
HIE (Ours)	10M	94.0 ± 0.3

5.3 CLASSIFICATION ON IMAGENET

Table 3 shows results on ImageNet. HIE-small achieves 75.7% accuracy with 18M parameters, improving upon both DEQ and hybrid hyperbolic baselines. Scaling to HIE-XL further increases performance to 79.3%, matching MDEQ-XL and avoiding out-of-memory failures that affect existing HyperbolicCV models at this scale. These results indicate that HIE scales effectively to large-scale classification, combining the scalability of implicit equilibrium models with the representational advantages of hyperbolic geometry. Further details of the experiments are provided in Appendix E.

5.4 EQUILIBRIUM CONVERGENCE

We compare the convergence of MDEQ and our proposed HIE on CIFAR-10, measuring residual change $\|z^{[i+1]} - z^{[i]}\|/\|z^{[i]}\|$ as in (Bai et al., 2020). Figure 8 shows averages over 10 runs with shaded deviations. Under Broyden’s solver (a), HIE reaches equilibrium faster and more stably than Euclidean DEQ, while in forward unrolling (b), HIE continues to decay monotonically whereas DEQ stagnates. This demonstrates the curvature-induced contraction of hyperbolic operators, which accelerates and stabilizes equilibrium convergence.

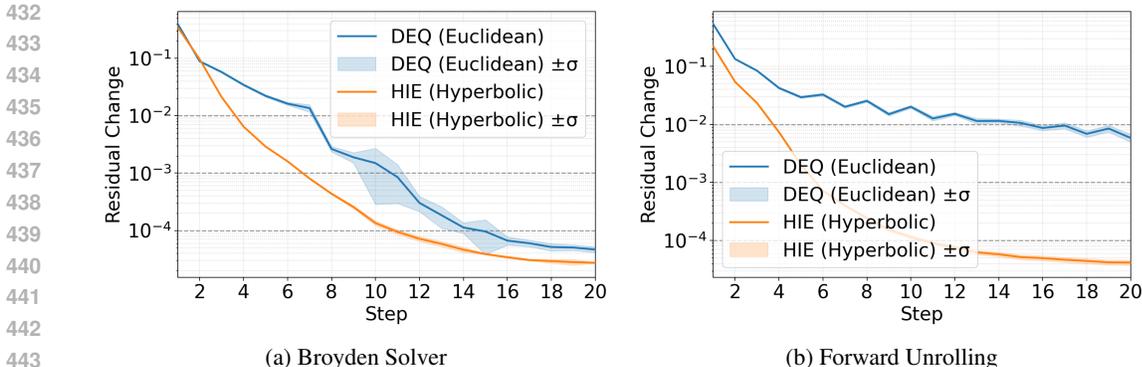


Figure 8: Convergence of MDEQ and HIE on CIFAR-10 over 10 runs. (a) Residual change when solving for the equilibrium with the Broyden solver. (b) Residual change when unrolling the network layer-by-layer without an equilibrium solver. Shaded regions denote \pm one standard deviation.

5.5 MEMORY AND RUNTIME

We benchmark runtime and memory on ImageNet with various backbones and sizes (Table 4). Both MDEQ and our HIE scale to larger batches without out-of-memory errors, unlike HyperbolicCV (Bdeir et al., 2024). HIE further reduces runtime relative to Euclidean DEQs while keeping identical memory use. Results on CIFAR-10 are provided in Appendix F.2.

The gain follows from Section 5.4: hyperbolic operators contract more sharply than Euclidean ones, so fewer solver iterations are required to reach equilibrium. Thus HIE achieves faster inference with the $O(1)$ memory of implicit differentiation. In contrast, explicit hyperbolic designs require layer unrolling, and compounded geometric operations prevent state sharing, leading to rapid memory growth.

Table 3: Evaluation on ImageNet classification.

Model	Model Size	Accuracy
AlexNet (Krizhevsky et al., 2012)	238M	57.0
ResNet-18 (He et al., 2016)	13M	70.2
ResNet-34 (He et al., 2016)	21M	74.8
Inception-V2 (Ioffe & Szegedy, 2015)	12M	74.8
ResNet-50 (He et al., 2016)	26M	75.1
HRNet-W18-C (Wang et al., 2020)	21M	76.8
Single-stream DEQ (Bai et al., 2019)	18M	72.9
MDEQ-small (Bai et al., 2020)	18M	75.5
HyperbolicNN (Ganea et al., 2018)	13M	65.7
Hybrid Poincaré (Guo et al., 2022)	13M	68.9
Hybrid Lorentz (Bdeir et al., 2024)	13M	70.1
Poincaré ResNet (Van Spengler et al., 2023)	10M	67.0
HECNN Lorentz (Bdeir et al., 2024)	13M	72.0
HCNN Lorentz (Bdeir et al., 2024)	13M	71.7
HECNN Lorentz (Bdeir et al., 2024)	26M	75.3
HCNN Lorentz (Bdeir et al., 2024)	26M	75.1
HIE (Ours)	18M	75.7
ResNet-101 (He et al., 2016)	52M	77.1
W-ResNet-50 (Zagoruyko & Komodakis, 2016)	69M	78.1
DenseNet-264 (Huang et al., 2017)	74M	79.7
MDEQ-large (Bai et al., 2020)	63M	77.5
Unrolled 5-layer MDEQ-large (Bai et al., 2020)	63M	75.9
MDEQ-XL (Bai et al., 2020)	81M	79.2
HECNN Lorentz (Bdeir et al., 2024)	81M	OOM
HCNN Lorentz (Bdeir et al., 2024)	81M	OOM
HIE-XL (Ours)	81M	79.3

Table 4: Runtime and memory consumption on ImageNet.

Model	Backbone	Model Size	Batch Size	Memory (GB)	Runtime (ms)
HECNN Lorentz (Bdeir et al., 2024)	ResNet-18	10M	4	14.5	337
HECNN Lorentz (Bdeir et al., 2024)	ResNet-18	10M	8	OOM	N/A
MDEQ (Bai et al., 2020)	MDEQ	10M	8	1.7	178
HIE (Ours)	HIE	10M	8	1.7	142
HECNN Lorentz (Bdeir et al., 2024)	ResNet-50	26M	2	23.0	753
HECNN Lorentz (Bdeir et al., 2024)	ResNet-50	26M	4	OOM	N/A
MDEQ-XL (Bai et al., 2020)	MDEQ	81M	4	2.1	31
HIE-XL (Ours)	HIE	81M	4	2.1	28

6 CONCLUSION

We design *Hyperbolic Implicit Equilibrium* (HIE), a constant-memory framework for negatively curved spaces. Our theory provides convergence, stability, and generalization guarantees. HIE converges faster and more stably than DEQ while avoiding explicit hyperbolic memory blow-ups. It scales on CIFAR-10 and ImageNet with $O(1)$ memory, reduced runtime, and competitive accuracy.

REFERENCES

- 486
487
488 Luis B Almeida. A learning rule for asynchronous perceptrons with feedback in a combinatorial
489 environment. In *Artificial neural networks: concept learning*, pp. 102–111. 1990.
- 490
491 Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural networks.
492 In *International conference on machine learning*, pp. 136–145. PMLR, 2017.
- 493
494 Mina Ghadimi Atigh, Julian Schoep, Erman Acar, Nanne Van Noord, and Pascal Mettes. Hyperbolic
495 image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern
496 recognition*, pp. 4453–4462, 2022.
- 497
498 Shaojie Bai, J Zico Kolter, and Vladlen Koltun. Deep equilibrium models. *Advances in neural
499 information processing systems*, 32, 2019.
- 500
501 Shaojie Bai, Vladlen Koltun, and J Zico Kolter. Multiscale deep equilibrium models. *Advances in
502 neural information processing systems*, 33:5238–5250, 2020.
- 503
504 Ahmad Bdeir, Kristian Schwethelm, and Niels Landwehr. Fully hyperbolic convolutional neural
505 networks for computer vision. In *International Conference on Learning Representations*, 2024.
506 URL <https://openreview.net/forum?id=ekz1hN5QNh>.
- 507
508 James W Cannon, William J Floyd, Richard Kenyon, Walter R Parry, et al. Hyperbolic geometry.
509 *Flavors of geometry*, 31(59-115):2, 1997.
- 510
511 Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous
512 convolution for semantic image segmentation. *CoRR*, abs/1706.05587, 2017. URL [http://
513 arxiv.org/abs/1706.05587](http://arxiv.org/abs/1706.05587).
- 514
515 Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary
516 differential equations. *Advances in neural information processing systems*, 31, 2018.
- 517
518 Weize Chen, Xu Han, Yankai Lin, Hexu Zhao, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou.
519 Fully hyperbolic neural networks. *arXiv preprint arXiv:2105.14686*, 2021.
- 520
521 Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo
522 Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban
523 scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern
524 recognition*, pp. 3213–3223, 2016.
- 525
526 Filipe de Avila Belbute-Peres, Kevin Smith, Kelsey Allen, Josh Tenenbaum, and J Zico Kolter. End-
527 to-end differentiable physics for learning and control. *Advances in neural information processing
528 systems*, 31, 2018.
- 529
530 Josip Djolonga and Andreas Krause. Differentiable learning of submodular models. *Advances in
531 Neural Information Processing Systems*, 30, 2017.
- 532
533 Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural odes. *Advances in neural
534 information processing systems*, 32, 2019.
- 535
536 Laurent El Ghaoui, Fangda Gu, Bertrand Travacca, Armin Askari, and Alicia Tsai. Implicit deep
537 learning. *SIAM Journal on Mathematics of Data Science*, 3(3):930–958, 2021.
- 538
539 Robert Englert, Balint Kincses, Raviteja Kotikalapudi, Giuseppe Gallitto, Jialin Li, Kevin Hoff-
schlag, Choong-Wan Woo, Tor D Wager, Dagmar Timmann, Ulrike Bingel, and Tamas Spisak.
Connectome-based attractor dynamics underlie brain activity in rest, task, and disease. September
2024. doi: 10.7554/eLife.98725.1. URL [http://dx.doi.org/10.7554/eLife.98725.
1](http://dx.doi.org/10.7554/eLife.98725.1).
- Aleksandr Ermolov, Leyla Mirvakhabova, Valentin Khrukov, Nicu Sebe, and Ivan Oseledets. Hy-
perbolic vision transformers: Combining improvements in metric learning. In *Proceedings of the
IEEE/CVF conference on computer vision and pattern recognition*, pp. 7409–7419, 2022.

- 540 Xiran Fan, Chun-Hao Yang, and Baba C Vemuri. Nested hyperbolic spaces for dimensionality
541 reduction and hyperbolic nn design. In *Proceedings of the IEEE/CVF Conference on Computer
542 Vision and Pattern Recognition*, pp. 356–365, 2022.
- 543 Hervé Fournier, Anas Ismail, and Antoine Vigneron. Computing the gromov hyperbolicity of a
544 discrete metric space. *Information Processing Letters*, 115(6-8):576–579, 2015.
- 545 Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic neural networks. *Advances in
546 neural information processing systems*, 31, 2018.
- 547 Stephen Gould, Richard Hartley, and Dylan Campbell. Deep declarative networks. *IEEE Transac-
548 tions on Pattern Analysis and Machine Intelligence*, 44(8):3988–4004, 2021.
- 549 Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. Ffjord:
550 Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint
551 arXiv:1810.01367*, 2018.
- 552 Yunhui Guo, Xudong Wang, Yubei Chen, and Stella X Yu. Clipped hyperbolic classifiers are super-
553 hyperbolic classifiers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and
554 Pattern Recognition*, pp. 11–20, 2022.
- 555 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
556 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
557 770–778, 2016.
- 558 Joy Hsu, Jeffrey Gu, Gong Wu, Wah Chiu, and Serena Yeung. Capturing implicit hierarchical
559 structure in 3d biomedical images with self-supervised hyperbolic representations. *Advances in
560 neural information processing systems*, 34:5112–5123, 2021.
- 561 Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
562 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern
563 recognition*, pp. 4700–4708, 2017.
- 564 Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by
565 reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456.
566 pmlr, 2015.
- 567 Valentin Khrulkov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky.
568 Hyperbolic image embeddings. In *Proceedings of the IEEE/CVF conference on computer vision
569 and pattern recognition*, pp. 6418–6428, 2020.
- 570 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
571 2009.
- 572 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convo-
573 lutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 574 Marc Law, Renjie Liao, Jake Snell, and Richard Zemel. Lorentzian distance learning for hyperbolic
575 representations. In *International Conference on Machine Learning*, pp. 3672–3681. PMLR, 2019.
- 576 Renjie Liao, Yuwen Xiong, Ethan Fetaya, Lisa Zhang, KiJung Yoon, Xaq Pitkow, Raquel Urtas-
577 sun, and Richard Zemel. Reviving and improving recurrent back-propagation. In *International
578 conference on machine learning*, pp. 3082–3091. PMLR, 2018.
- 579 Shaoteng Liu, Jingjing Chen, Liangming Pan, Chong-Wah Ngo, Tat-Seng Chua, and Yu-Gang Jiang.
580 Hyperbolic visual embedding learning for zero-shot recognition. In *Proceedings of the IEEE/CVF
581 conference on computer vision and pattern recognition*, pp. 9273–9281, 2020.
- 582 Yifan Liu, Ke Chen, Chris Liu, Zengchang Qin, Zhenbo Luo, and Jingdong Wang. Structured
583 knowledge distillation for semantic segmentation. In *Proceedings of the IEEE/CVF conference
584 on computer vision and pattern recognition*, pp. 2604–2613, 2019.
- 585 Alice Longhena, Martin Guillemaud, Fabrizio De Vico Fallani, Raffaella Migliaccio, and Mario
586 Chavez. Hyperbolic embedding of brain networks detects regions disrupted by neurodegeneration
587 in alzheimer’s disease. *Physical Review E*, 111(4):044402, 2025.

- 594 Emile Mathieu, Charline Le Lan, Chris J Maddison, Ryota Tomioka, and Yee Whye Teh. Con-
595 tinuous hierarchical representations with poincaré variational auto-encoders. *Advances in neural*
596 *information processing systems*, 32, 2019.
- 597 Gal Mishne, Zhengchao Wan, Yusu Wang, and Sheng Yang. The numerical stability of hyperbolic
598 representation learning. In *International Conference on Machine Learning*, pp. 24925–24949.
599 PMLR, 2023.
- 600 Yoshihiro Nagano, Shoichiro Yamaguchi, Yasuhiro Fujita, and Masanori Koyama. A wrapped nor-
601 mal distribution on hyperbolic space for gradient-based learning. In *International conference on*
602 *machine learning*, pp. 4693–4702. PMLR, 2019.
- 603 Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of
604 hyperbolic geometry. In *International conference on machine learning*, pp. 3779–3788. PMLR,
605 2018.
- 606 Ivan Ovinnikov. Poincaré wasserstein autoencoder. *arXiv preprint arXiv:1901.01427*, 2019.
- 607 Wei Peng, Tuomas Varanka, Abdelrahman Mostafa, Henglin Shi, and Guoying Zhao. Hyperbolic
608 deep neural networks: A survey. *IEEE Transactions on pattern analysis and machine intelligence*,
609 44(12):10023–10044, 2021.
- 610 Fernando Pineda. Generalization of back propagation to recurrent and higher order neural networks.
611 In *Neural information processing systems*, 1987.
- 612 Yi-Ling Qiao, Junbang Liang, Vladlen Koltun, and Ming C Lin. Scalable differentiable physics for
613 learning and control. *arXiv preprint arXiv:2007.02168*, 2020.
- 614 John G Ratcliffe. *Foundations of hyperbolic manifolds*. Springer, 2006.
- 615 Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
616 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
617 recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- 618 Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mo-
619 bilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on*
620 *computer vision and pattern recognition*, pp. 4510–4520, 2018.
- 621 Ryohei Shimizu, Yusuke Mukuta, and Tatsuya Harada. Hyperbolic neural networks++. *arXiv*
622 *preprint arXiv:2006.08210*, 2020.
- 623 Yuhang Song, Beren Millidge, Tommaso Salvatori, Thomas Lukasiewicz, Zhenghua Xu, and Rafal
624 Bogacz. Inferring neural activity before plasticity as a foundation for learning beyond backprop-
625 agation. *Nature neuroscience*, 27(2):348–358, 2024.
- 626 Towaki Takikawa, David Acuna, Varun Jampani, and Sanja Fidler. Gated-scnn: Gated shape cnns for
627 semantic segmentation. In *Proceedings of the IEEE/CVF international conference on computer*
628 *vision*, pp. 5229–5238, 2019.
- 629 Max Van Spengler, Erwin Berkhout, and Pascal Mettes. Poincare resnet. In *Proceedings of the*
630 *IEEE/CVF International Conference on Computer Vision*, pp. 5419–5428, 2023.
- 631 Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu,
632 Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learn-
633 ing for visual recognition. *IEEE transactions on pattern analysis and machine intelligence*, 43
634 (10):3349–3364, 2020.
- 635 Po-Wei Wang, Priya Donti, Bryan Wilder, and Zico Kolter. Satnet: Bridging deep learning and log-
636 ical reasoning using a differentiable satisfiability solver. In *International Conference on Machine*
637 *Learning*, pp. 6545–6554. PMLR, 2019.
- 638 Yuxin Wu and Kaiming He. Group normalization. In *Proceedings of the European conference on*
639 *computer vision (ECCV)*, pp. 3–19, 2018.

- 648 Jiexi Yan, Lei Luo, Cheng Deng, and Heng Huang. Unsupervised hyperbolic metric learning. In
649 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 12465–
650 12474, 2021.
- 651 Fisher Yu, Vladlen Koltun, and Thomas Funkhouser. Dilated residual networks. In *Proceedings of*
652 *the IEEE conference on computer vision and pattern recognition*, pp. 472–480, 2017.
- 653 Yun Yue, Fangzhou Lin, Kazunori D Yamada, and Ziming Zhang. Hyperbolic contrastive learning.
654 *arXiv preprint arXiv:2302.01409*, 2023.
- 655 Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint*
656 *arXiv:1605.07146*, 2016.
- 657 Huanqiu Zhang, P Dylan Rich, Albert K Lee, and Tatyana O Sharpee. Hippocampal spatial repre-
658 sentations exhibit a hyperbolic geometry that expands with experience. *Nature Neuroscience*, 26
659 (1):131–139, 2023.
- 660 Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing
661 network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*
662 *(CVPR)*, July 2017.
- 663 Hengshuang Zhao, Yi Zhang, Shu Liu, Jianping Shi, Chen Change Loy, Dahua Lin, and Jiaya Jia.
664 Psanet: Point-wise spatial attention network for scene parsing. In *Proceedings of the European*
665 *Conference on Computer Vision (ECCV)*, September 2018.
- 666 Yuansheng Zhou, Brian H Smith, and Tatyana O Sharpee. Hyperbolic geometry of the olfactory
667 space. *Science advances*, 4(8):eaq1458, 2018a.
- 668 Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. Unet++:
669 A nested u-net architecture for medical image segmentation. In *International workshop on deep*
670 *learning in medical image analysis*, pp. 3–11. Springer, 2018b.

675 ETHICS STATEMENT

676 This work complies with the ICLR Code of Ethics. Our study uses only publicly available datasets
677 (CIFAR-10 (Krizhevsky et al., 2009), ImageNet (Russakovsky et al., 2015), Cityscapes (Cordts
678 et al., 2016)) and does not involve human subjects, private data, or sensitive content. We believe
679 the broader impacts of this work are positive, as it advances efficient and scalable deep learning
680 methods.
681
682

683 REPRODUCIBILITY STATEMENT

684 We have made every effort to ensure reproducibility. All theoretical proofs are included in the
685 Appendix D. Details of model architectures, training setups, and hyperparameters are provided in
686 Sections 3, 5, and Appendix. Code and pretrained models will be released upon publication.
687
688

689 THE USE OF LARGE LANGUAGE MODELS

690 We used a large language model (ChatGPT) solely to aid in language polishing and clarity improve-
691 ments of the manuscript. The research ideas, method design, and experimental results are entirely
692 our own. No part of the technical content or experimental findings is generated by an LLM.
693
694

695 A HYPERBOLIC GEOMETRY

696 This appendix summarizes the main operators in hyperbolic space, with a focus on the Lorentz
697 model. Many of these constructions apply to general Riemannian manifolds, but here we emphasize
698 the closed-form expressions that make the Lorentz model particularly attractive for deep learning
699 applications (Cannon et al., 1997; Ratcliffe, 2006; Nickel & Kiela, 2018; Law et al., 2019; Chen
700 et al., 2021).
701

702 A.1 POINCARÉ BALL AND LORENTZ MODEL

703
704 **Poincaré Ball.** The n -dimensional Poincaré ball $\mathbb{B}_K^n = (\mathbb{B}^n, g^K)$ with curvature $K < 0$ is

$$705 \mathbb{B}^n = \{x \in \mathbb{R}^n : -K\|x\|^2 < 1\}, \quad g_x^K = \lambda_x^2 I, \quad \lambda_x = \frac{2}{1+K\|x\|^2}.$$

706
707 It describes hyperbolic space as an open ball of radius $\sqrt{-1/K}$.

708
709 **Lorentz Model.** The Lorentz model realizes \mathbb{H}_K^d as the upper sheet of a two-sheeted hyperboloid
710 embedded in \mathbb{R}^{d+1} with bilinear form

$$711 \langle u, v \rangle_{\mathcal{L}} = -u_0 v_0 + \sum_{i=1}^d u_i v_i.$$

712 Thus,

$$713 \mathbb{H}_K^d = \{u \in \mathbb{R}^{d+1} : \langle u, u \rangle_{\mathcal{L}} = 1/K, u_0 > 0\}.$$

714 A.2 DISTANCE AND TANGENT SPACE

715
716 The geodesic distance between $x, y \in \mathbb{H}_K^d$ is

$$717 d_{\mathbb{H}_K}(x, y) = \frac{1}{\sqrt{-K}} \operatorname{arcosh}(K\langle x, y \rangle_{\mathcal{L}}). \quad (12)$$

718 An equivalent form for squared distance is (Law et al., 2019):

$$719 d_{\mathbb{H}_K}^2(x, y) = \frac{2}{K} - 2\langle x, y \rangle_{\mathcal{L}}.$$

720 The tangent space at $x \in \mathbb{H}_K^d$ is

$$721 T_x \mathbb{H}_K^d = \{v \in \mathbb{R}^{d+1} : \langle v, x \rangle_{\mathcal{L}} = 0\}.$$

722 A.3 EXPONENTIAL AND LOGARITHMIC MAPS

723 For $z \in T_x \mathbb{H}_K^d$, define $\alpha = \sqrt{-K} \|z\|_{\mathcal{L}}$. The exponential map is

$$724 \exp_x^K(z) = \cosh(\alpha) x + \sinh(\alpha) \frac{z}{\alpha}. \quad (13)$$

725 Conversely, the logarithmic map for $y \in \mathbb{H}_K^d$ is

$$726 \log_x^K(y) = \frac{\operatorname{arcosh}(\beta)}{\sqrt{\beta^2 - 1}} (y - \beta x), \quad \beta = K\langle x, y \rangle_{\mathcal{L}}. \quad (14)$$

727 At the origin $o = [1/\sqrt{-K}, 0, \dots, 0]$, the exponential simplifies to

$$728 \exp_o^K(z) = \frac{1}{\sqrt{-K}} (\cosh(\sqrt{-K}\|z\|), \sinh(\sqrt{-K}\|z\|) \frac{z}{\|z\|}).$$

729 A.4 PARALLEL TRANSPORT

730 The parallel transport of $v \in T_x \mathbb{H}_K^d$ along the geodesic from x to y is

$$731 \operatorname{PT}_{x \rightarrow y}^K(v) = v - \frac{\langle \log_x^K(y), v \rangle_{\mathcal{L}}}{d_{\mathbb{H}_K}(x, y)} (\log_x^K(y) + \log_y^K(x)) \quad (15)$$

$$732 = v + \frac{\langle y, v \rangle_{\mathcal{L}}}{1/(-K) - \langle x, y \rangle_{\mathcal{L}}} (x + y). \quad (16)$$

733 A.5 LORENTZIAN CENTROID AND POOLING

734 The weighted Lorentzian centroid of $\{x_i\}_{i=1}^m$ with weights $\nu_i \geq 0$, $\sum_i \nu_i > 0$, solves
735 $\min_{\mu} \sum_i \nu_i d_{\mathbb{H}_K}^2(x_i, \mu)$ and admits the closed form (Law et al., 2019):

$$736 \mu = \frac{\sum_i \nu_i x_i}{\sqrt{-K} \left\| \sum_i \nu_i x_i \right\|_{\mathcal{L}}}.$$

737 Average pooling in hyperbolic space can thus be implemented via the centroid over receptive fields.

756 A.6 LORENTZ TRANSFORMATIONS

757 Lorentz transformations are linear maps $A \in \mathbb{R}^{(d+1) \times (d+1)}$ that preserve the bilinear form:
 758 $\langle Ax, Ay \rangle_{\mathcal{L}} = \langle x, y \rangle_{\mathcal{L}}$. They form the group $O^+(1, d)$. By polar decomposition, any such A factors
 759 as a Lorentz rotation $R \in SO^+(1, d)$ and a Lorentz boost $B(v)$ determined by a velocity $v \in \mathbb{R}^d$,
 760 $\|v\| < 1$.

762 A.7 LORENTZ FULLY-CONNECTED LAYER

763 Following Chen et al. (2021), linear maps in tangent space cannot realize all Lorentz transfor-
 764 mations. Instead, one can directly parameterize in ambient space. Given $x \in \mathbb{H}_K^d$, weights
 765 $W \in \mathbb{R}^{m \times (d+1)}$, and bias b , define

$$766 y = \left[\sqrt{\|\psi(Wx + b)\|^2 - 1/K}, \psi(Wx + b) \right],$$

767 where ψ includes nonlinearity and bias.

771 A.8 CONCATENATION

772 Given points $\{x_i \in \mathbb{H}_K^d\}_{i=1}^N$, Lorentz direct concatenation produces a point $y \in \mathbb{H}_K^{Nd}$ by stacking
 773 their spatial components with an adjusted time coordinate to ensure validity on the hyperboloid.

776 A.9 WRAPPED NORMAL DISTRIBUTION

777 A wrapped normal distribution on \mathbb{H}_K^d can be constructed as (Nagano et al., 2019):

- 778 1. Sample $v \sim \mathcal{N}(0, \Sigma)$ in \mathbb{R}^d , extend to $[0, v] \in T_o\mathbb{H}_K^d$.
- 779 2. Parallel transport to $T_\mu\mathbb{H}_K^d$ for mean μ .
- 780 3. Map to \mathbb{H}_K^d via \exp_μ^K .

781 This distribution has closed-form density and supports efficient sampling.

786 A.10 MAPPING BETWEEN MODELS

787 The Lorentz and Poincaré models are isometric. A diffeomorphism mapping $x = [x_t, x_s] \in \mathbb{H}_K^d$ to
 788 the Poincaré ball is

$$789 p_{\mathbb{H} \rightarrow \mathbb{B}}(x) = \frac{x_s}{x_t + 1/\sqrt{-K}}.$$

793 B HYPERBOLIC NEURAL NETWORKS

794 B.1 COMPUTATION OF δ -HYPERBOLICITY

795 Let (\mathcal{X}, d) be a metric space. For any three points $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathcal{X}$, the Gromov product with basepoint
 796 \mathbf{x} is defined as

$$797 (\mathbf{v}, \mathbf{w})_{\mathbf{u}} = \frac{1}{2}(d(\mathbf{u}, \mathbf{v}) + d(\mathbf{u}, \mathbf{w}) - d(\mathbf{v}, \mathbf{w})). \quad (17)$$

798 The space is said to be δ -hyperbolic if for all $\mathbf{u}, \mathbf{v}, \mathbf{w}, \mathbf{x} \in \mathcal{X}$,

$$799 (\mathbf{u}, \mathbf{w})_{\mathbf{x}} \geq \min\{(\mathbf{u}, \mathbf{v})_{\mathbf{x}}, (\mathbf{v}, \mathbf{w})_{\mathbf{x}}\} - \delta. \quad (18)$$

800 The smallest δ satisfying this inequality serves as the hyperbolicity constant of the space. For
 801 computational purposes, we approximate δ using the efficient min-max matrix product ap-
 802 proach (Fournier et al., 2015), and normalize it by the dataset diameter to obtain a scale-invariant
 803 measure:

$$804 \delta_{\text{rel}}(\mathcal{X}) = \frac{2\delta(\mathcal{X})}{\text{diam}(\mathcal{X})}, \quad \text{diam}(\mathcal{X}) = \max_{\mathbf{u}, \mathbf{v} \in \mathcal{X}} d(\mathbf{u}, \mathbf{v}). \quad (19)$$

805 By construction, $\delta_{\text{rel}} \in [0, 1]$, where values closer to 0 indicate stronger hyperbolicity. This normal-
 806 ized score is used throughout our experiments to decide whether to employ hyperbolic or Euclidean
 807 components within the implicit equilibrium solver.

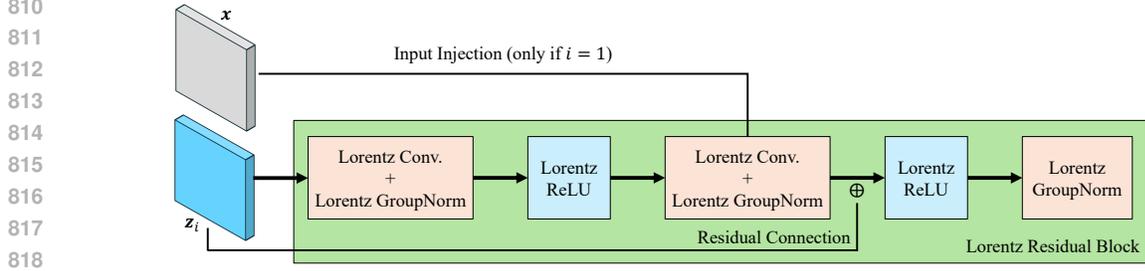


Figure 9: The Lorentz residual block used in HIE. An HIE contains only *one* such layer.

B.2 HYPERBOLIC MODULES

The internal structure of the Lorentz residual block is shown in Figure 9. We use curvature $K < 0$ and set $K = -1$ for all experiments. Points on the manifold are in bold (e.g., $\mathbf{u} \in \mathbb{H}_K^d$). We denote the Lorentz inner product by $\langle \cdot, \cdot \rangle_{\mathcal{L}}$ and use the origin $\mathbf{o} = [1/\sqrt{-K}, \mathbf{0}]^\top$.

Lorentz Convolutional Layer. A hyperbolic feature map is an ordered grid of manifold points $\{\mathbf{u}_{h,w} \in \mathbb{H}_K^d\}_{h,w}$, stored in channel-last format by appending the time component as an extra channel. Given a kernel with receptive field size $\tilde{H} \times \tilde{W}$, we (i) collect the points in the local window, (ii) concatenate them hyperbolically, and (iii) apply a Lorentz fully-connected (LFC) transformation to obtain the output point:

$$\mathbf{y}_{h,w} = \text{LFC}\left(\text{HCat}\left(\{\mathbf{u}_{h+\hat{h}, w+\hat{w}}\}\right)\right), \quad (20)$$

where HCat denotes Lorentz direct concatenation and LFC is a learnable Lorentz transformation that maps $\mathbb{H}_K^{d\tilde{H}\tilde{W}} \rightarrow \mathbb{H}_K^d$. Concretely, letting $\psi(\cdot)$ be an element-wise nonlinearity in the *space* component,¹ the LFC can be implemented in ambient coordinates as

$$\mathbf{y} = \begin{bmatrix} \sqrt{\|\psi(W\mathbf{x}_s + \mathbf{b})\|_2^2 - 1/K} \\ \psi(W\mathbf{x}_s + \mathbf{b}) \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_t \\ \mathbf{x}_s \end{bmatrix} \in \mathbb{H}_K^D, \quad (21)$$

with W, \mathbf{b} learnable. Padding uses the origin \mathbf{o} (hyperbolic “zero”). The transposed Lorentz convolution reuses the same operator with modified connectivity (origin-insertion upsampling).

Lorentz MLR Classifier. We model class regions via hyperbolic hyperplanes. Using the reparameterization by a direction $\mathbf{z} \in \mathbb{R}^d$ and an offset $a \in \mathbb{R}$, the hyperplane is

$$\tilde{\mathcal{H}}_{\mathbf{z},a} = \left\{ \mathbf{x} \in \mathbb{H}_K^d \mid \cosh(\sqrt{-K} a) \langle \mathbf{z}, \mathbf{x}_s \rangle - \sinh(\sqrt{-K} a) \|\mathbf{z}\|_2 x_t = 0 \right\}, \quad (22)$$

where $\mathbf{x} = [x_t, \mathbf{x}_s]^\top$. The (signed) distance from \mathbf{x} to $\tilde{\mathcal{H}}_{\mathbf{z},a}$ is

$$\begin{aligned} d_{\mathbb{H}_K}(\mathbf{x}, \tilde{\mathcal{H}}_{\mathbf{z},a}) &= \frac{1}{\sqrt{-K}} \operatorname{asinh}\left(\frac{\sqrt{-K} \alpha}{\beta}\right), \\ \alpha &= \cosh(\sqrt{-K} a) \langle \mathbf{z}, \mathbf{x}_s \rangle - \sinh(\sqrt{-K} a) \|\mathbf{z}\|_2 x_t, \\ \beta &= \sqrt{\|\cosh(\sqrt{-K} a) \mathbf{z}\|_2^2 - (\sinh(\sqrt{-K} a) \|\mathbf{z}\|_2)^2}. \end{aligned} \quad (23)$$

For C classes with parameters $\{(\mathbf{z}_c, a_c)\}_{c=1}^C$, the Lorentz MLR logits are proportional to (signed) distances:

$$\ell_c(\mathbf{x}) = \frac{1}{\sqrt{-K}} \operatorname{sign}(\alpha_c) \beta_c \left| \operatorname{asinh}\left(\frac{\sqrt{-K} \alpha_c}{\beta_c}\right) \right|, \quad \alpha_c, \beta_c \text{ as in equation 23 for } (\mathbf{z}_c, a_c). \quad (24)$$

Class probabilities follow by softmax over $\{\ell_c\}_{c=1}^C$.

¹We avoid internal tangent-space shuttling to improve stability.

Lorentz Residual Connection and Activation. Residual addition is ill-defined on \mathbb{H}_K^d . We therefore perform *space-component addition* and recover the time component geometrically. Let the block input be $\mathbf{u} = [u_t, \mathbf{u}_s]^\top$ and the block transform output (in ambient coords) be $\mathbf{r} = [r_t, \mathbf{r}_s]^\top$. The residual output is

$$\mathbf{v} = \begin{bmatrix} \sqrt{\|\mathbf{u}_s + \mathbf{r}_s\|_2^2 - 1/K} \\ \mathbf{u}_s + \mathbf{r}_s \end{bmatrix} \in \mathbb{H}_K^d. \quad (25)$$

Non-linear activations act on the space component only; e.g., a Lorentz ReLU is

$$\text{LReLU} \left(\begin{bmatrix} x_t \\ \mathbf{x}_s \end{bmatrix} \right) = \begin{bmatrix} \sqrt{\|\text{ReLU}(\mathbf{x}_s)\|_2^2 - 1/K} \\ \text{ReLU}(\mathbf{x}_s) \end{bmatrix}. \quad (26)$$

This avoids frequent log / exp mappings and has shown better stability and efficiency than tangent-space activations in our setting.

B.3 LORENTZ GROUP NORMALIZATION

Design notes and variants are as follows:

- **Batch-independence.** LGN mirrors Euclidean GroupNorm by using *within-sample, within-group* statistics (μ_g, σ_g^2) only, where $\sigma_g^2 = \frac{1}{N_g} \sum_j d_{\mathbb{H}_K}(\mathbf{u}_j, \mu_g)^2$. This avoids batch-coupled population estimates that are known to destabilize implicit fixed-point solvers in equilibrium models.
- **Choice of groups.** We form groups over *sets of Lorentz points* (e.g., disjoint subsets of spatial sites) rather than slicing coordinates of a single point, preserving the integrity of each manifold element. In practice we adopt a uniform partition of spatial locations into G groups per resolution (other partitions are admissible as long as each site is treated as a whole Lorentz point).
- **Recovering special cases.** Setting $G=1$ yields a Lorentz *Instance Normalization* (per sample, per resolution). Replacing groups by the full mini-batch recovers *Lorentz Batch Normalization (LBN)*. While the latter is geometrically well-defined, it reintroduces batch dependence and is usually unfavorable for equilibrium solvers; we therefore default to group-wise, batch-independent LGN.
- **Numerical stability.** We use a small $\varepsilon > 0$, clamp arguments of $\text{arcosh}(\cdot)$ to $[1+\tau, \infty)$ with $\tau \approx 10^{-6}$, and rely on closed-form Lorentz operators for general $K < 0$ (Appendix A). In all experiments we set $K = -1$; for general K the tangent norms and distances are scaled by $1/\sqrt{|K|}$, and all maps ($\exp^K, \log^K, \text{PT}^K$) are implemented in a K -aware manner. We also adopt channel-last storage so that each spatial site corresponds to one Lorentz point.

C MULTISCALE IMPLICIT EQUILIBRIUM

Setup. Let $f_\theta(z; x)$ be the transformation defining the (multiscale) dynamics, and introduce $g_\theta(z; x) := f_\theta(z; x) - z$. An equilibrium z^* is any solution of $g_\theta(z^*; x) = 0$, equivalently $z^* = f_\theta(z^*; x)$, which we obtain by root finding on g_θ rather than by explicitly unrolling layers to depth L .

C.1 FORWARD PASS (ROOT FINDING)

The forward pass seeks $z^* = \text{Rootfind}(g_\theta; x)$ using a black-box solver (e.g., Newton or quasi-Newton). In practice we use (limited-memory) Broyden:

$$z^{(k+1)} = z^{(k)} - \alpha B^{(k)} g_\theta(z^{(k)}; x),$$

where $B^{(k)} \approx (Jg_\theta|_{z^{(k)}})^{-1}$ is updated from the most recent low-rank corrections; storing only the latest m updates yields an L-Broyden scheme suitable for high-dimensional vision features.

Multiscale state. For MDEQ, the hidden state is a tuple $z = [z_1, \dots, z_n]$ with different spatial resolutions and channel dimensions; we initialize $z_i^{(0)} = 0$ and solve for all scales *jointly* so that the

918 solver enforces cross-scale consistency at equilibrium. Limited-memory updates are crucial here
 919 due to the very large Jacobians encountered at realistic image resolutions.

920 **Stopping and memory.** We stop when $\|g_\theta(z^{(k)}; x)\|$ falls below a tolerance or when a cap on
 921 function evaluations is reached. Because the backward pass does not replay the forward trajectory,
 922 the training memory scales with the size of a *single* block rather than the (implicit) depth.
 923

924 C.2 BACKWARD PASS (IMPLICIT DIFFERENTIATION)

925 Given a loss $\ell = L(z^*, y)$, the implicit function theorem yields gradient expressions that depend only
 926 on quantities at z^* ; we do *not* differentiate through the root-finding iterations. Writing $J_g = Jg_\theta|_{z^*}$,
 927 the adjoint \bar{z} is obtained by solving the linear system
 928

$$929 (J_g^\top) \bar{z} = \frac{\partial \ell}{\partial z^*}, \quad \text{equivalently} \quad \bar{z}^\top J_g + \frac{\partial \ell}{\partial z^*} = 0,$$

930 which is implemented as a vector–Jacobian product (VJP) solve rather than forming J_g explicitly.
 931 Once \bar{z} is computed, parameter/input gradients follow from
 932

$$933 \frac{\partial \ell}{\partial \theta} = \bar{z}^\top \frac{\partial f_\theta(z^*; x)}{\partial \theta}, \quad \frac{\partial \ell}{\partial x} = \bar{z}^\top \frac{\partial f_\theta(z^*; x)}{\partial x},$$

934 with the minus sign absorbed by solving against J_g^\top as above:contentReference[oaicite:9]index=9.
 935 These expressions coincide with the familiar DEQ formulas and make the backward pass depend on
 936 a *single* linear solve at the equilibrium, decoupled from the forward solver’s trajectory.
 937

938 **Remarks for multiscale equilibria.** The adjoint is defined over the concatenated (but dimensionally
 939 heterogeneous) state $z = [z_1, \dots, z_n]$, and the VJP solve couples all resolutions, mirroring the
 940 forward fusion. In practice, the same limited-memory strategy used in the forward pass is applied to
 941 the backward linear solve to control memory and runtime at megapixel scales.
 942

943 **Summary.** The forward computes z^* by root finding on g_θ ; the backward computes an adjoint \bar{z}
 944 via one linear VJP solve and then applies local derivatives of f_θ at z^* . This yields constant training
 945 memory and avoids backpropagating through the unrolled fixed-point iterations.
 946

947 D THEORETICAL PROOFS

948 **Preliminaries.** We work on \mathbb{H}_K^d with $K < 0$, $\kappa = \sqrt{-K}$. Riemannian metrics at z induce norms
 949 $\|\cdot\|_{g_z}$ on $T_z \mathbb{H}_K^d$. When needed, parallel transport identifies tangent vectors across nearby points.
 950

951 D.1 PROOF OF LEMMA 1

952 In geodesic polar coordinates at β , the differential of \exp_β^K at v with $r = \|v\|_{g_\beta}$ acts as identity
 953 along the radial direction and as the linear map scaling tangential directions by $\frac{\sinh(\kappa r)}{\kappa r}$, derived
 954 from the Jacobi field J solving $J'' + \kappa^2 J = 0$ with $J(0) = 0$, $J'(0) = I$. Hence the operator norm
 955 equals $\sinh(\kappa r)/(\kappa r)$ for $\|v\| = r$. Invertibility of \exp_β on $B_R(\beta)$ yields $d \log_\beta = (d \exp_\beta)^{-1}$ at
 956 corresponding points, so $\|d \log_\beta\| = \kappa r / \sinh(\kappa r)$. Taking suprema over $r \leq R$ gives the claimed
 957 bounds and $L_{\exp}(R)L_{\log}(R) = 1$.
 958

959 D.2 PROOF OF THEOREM 1

960 Let $z_i \in B_R(\beta)$, $u_i = \log_\beta(z_i)$ and $\Phi(u) = \tau A(S(u)) + (1 - \tau)u + b(x)$. By Lemma 1,
 961

$$962 d_{\mathbb{H}_K}(F(z_1), F(z_2)) \leq \|d \exp_\beta\|_{\sup, R} \cdot \|\Phi(u_1) - \Phi(u_2)\| = L_{\exp}(R) \|\Phi(u_1) - \Phi(u_2)\|.$$

963 Using $\|A\| \leq \alpha$ and S being η -Lipschitz on $\{\|u\| \leq R\}$,

$$964 \|\Phi(u_1) - \Phi(u_2)\| \leq \tau \|A\| \cdot \|S(u_1) - S(u_2)\| + (1 - \tau) \|u_1 - u_2\| \leq (\tau \alpha \eta + (1 - \tau)) \|u_1 - u_2\|.$$

965 Again by Lemma 1, $\|u_1 - u_2\| \leq L_{\log}(R) d_{\mathbb{H}_K}(z_1, z_2)$. Combining the inequalities yields
 966

$$967 d_{\mathbb{H}_K}(F(z_1), F(z_2)) \leq L_{\exp}(R) (\tau \alpha \eta + (1 - \tau)) L_{\log}(R) d_{\mathbb{H}_K}(z_1, z_2) = \rho_{\mathbb{H}} d_{\mathbb{H}_K}(z_1, z_2).$$

968 Since $\rho_{\mathbb{H}} < 1$, F is a contraction, so it admits a unique fixed point in $B_R(\beta)$ and the Picard iteration
 969 converges linearly with rate at most $\rho_{\mathbb{H}}$.
 970
 971

Table 5: Settings & hyperparameters of each task. “cls.” means classification task, and “seg.” means segmentation task. These models correspond to the ones reported in Tables 2, 3, and 7.

	CIFAR-10 (cls.)		ImageNet (cls.)		Cityscapes (seg.)	
	HIE-Small	HIE	HIE-Small	HIE-Large	HIE-Small	HIE-Large
Input Image Size	32 × 32		224 × 224		1024 × 512 (train)	2048 × 1024 (test)
Number of Epochs	50	200	100	100	480	480
Batch Size	128	128	128	128	12	12
Optimizer	Adam	Adam	SGD	SGD	SGD	SGD
(Start) Learning Rate	0.001	0.001	0.05	0.05	0.01	0.01
Nesterov Momentum	-	-	0.9	0.9	-	-
Weight Decay	0	0	5e-5	1e-4	2e-4	3e-4
Use Pre-trained Weights	-	-	-	-	Yes, from ImageNet	Yes, from ImageNet
Number of Scales	3	4	4	4	(Exact same model as in ImageNet)	
# of Channels for Each Scale	[8,16,32]	[28,56,112,224]	[32,64,128,256]	[80,160,320,640]	-	-
Width Expansion (in the residual block)	5×	5×	5×	5×	-	-
Normalization (# of groups)	GroupNorm(4)	GroupNorm(4)	GroupNorm(4)	GroupNorm(4)	-	-
Weight Normalization	✓	✓	✓	✓	-	-
# of Downsamplings Before Equilibrium Solver	0	0	2	2	-	-
Forward Quasi-Newton Threshold T_f	15	15	22	22	27	27
Backward Quasi-Newton Threshold T_b	18	18	25	25	30	30
Limited-Mem. Broyden’s Method Storage Size m	12	12	18	18	18	18
Variational Dropout Rate	0.2	0.25	0.0	0.0	0.03	0.05

D.3 PROOF OF COROLLARY 1

For $F_{\mathbb{R}}(u) = \tau Au + (1 - \tau)u + b$, $\text{Lip}(F_{\mathbb{R}}) = \tau\alpha + (1 - \tau) = \rho_{\mathbb{R}}$. By Theorem 1, $\rho_{\mathbb{H}} = \tau\alpha\eta + (1 - \tau) \leq \rho_{\mathbb{R}}$ and $<$ holds if $\tau\alpha > 0$ and $\eta < 1$.

D.4 PROOF OF THEOREM 2

Consider $G(\theta, z) := F_{\theta}(z; x) - z$ so that $G(\theta, z^*) = 0$. Differentiate w.r.t. θ :

$$\partial_{\theta}G + D_zG \partial_{\theta}z^* = 0 \quad \Rightarrow \quad (I - D_zF_{\theta}(z^*; x)) \partial_{\theta}z^* = \partial_{\theta}F_{\theta}(z^*; x).$$

Thus $\partial_{\theta}z^* = (I - D_zF_{\theta})^{-1} \partial_{\theta}F_{\theta}$. For any differentiable \mathcal{L} , $\nabla_{\theta}\mathcal{L} = (D_z\mathcal{L}(z^*)) \partial_{\theta}z^*$. Taking norms in $(T_{z^*}\mathbb{H}_{\kappa}^d, g_{z^*})$ and using submultiplicativity,

$$\|\nabla_{\theta}\mathcal{L}\| \leq \|(I - D_zF_{\theta})^{-1}\| \|\partial_{\theta}F_{\theta}(z^*; x)\| \|\nabla_z\mathcal{L}(z^*)\|.$$

Since $\|D_zF_{\theta}(z^*; x)\| \leq \rho_{\mathbb{H}} < 1$, the Neumann series gives $\|(I - D_zF_{\theta})^{-1}\| \leq 1/(1 - \rho_{\mathbb{H}})$, proving equation 10.

D.5 PROOF OF LEMMA 2

In the Lorentz model, for unit time-like vectors x, y (Lorentz norm -1), the geodesic distance satisfies $\cosh(\kappa d_{\mathbb{H}_{\kappa}}(x, y)) = -\langle x, y \rangle_{\text{Lor}}$. Writing $x = \exp_{\beta}(v)$, $y = \exp_{\beta}(w)$ and expressing v, w in geodesic polar coordinates gives the stated hyperbolic law of cosines; see standard Riemannian geometry texts.

D.6 PROOF OF THEOREM 3

With $r_1 = r_2 = r$ in Lemma 2, and $\cos \theta = 1 - 2 \sin^2(\theta/2)$,

$$\cosh(\kappa\gamma_{\mathbb{H}}) = 1 + 2 \sinh^2(\kappa r) \sin^2\left(\frac{\theta_0}{2}\right) = \cosh\left(2 \operatorname{arcsinh}\left(\sinh(\kappa r) \sin\left(\frac{\theta_0}{2}\right)\right)\right).$$

Taking arcosh yields $\gamma_{\mathbb{H}} = \frac{2}{\kappa} \operatorname{arcsinh}(\sinh(\kappa r) \sin(\theta_0/2))$. Since $\sinh(\kappa r) \geq \kappa r$ and $\operatorname{arcsinh}$ is increasing and concave on $[0, \infty)$, for $s = \sin(\theta_0/2) \in [0, 1]$,

$$\operatorname{arcsinh}(\sinh(\kappa r)s) \geq s \operatorname{arcsinh}(\sinh(\kappa r)) = s\kappa r,$$

hence $\gamma_{\mathbb{H}} \geq 2rs = \gamma_{\mathbb{R}}$, with strict inequality when $r > 0$ and $\theta_0 \in (0, \pi)$.

E ADDITIONAL IMPLEMENTATION DETAILS

We provide the implementation details in Table 5.

F ADDITIONAL EXPERIMENTAL RESULTS

F.1 HYPERBOLICITY

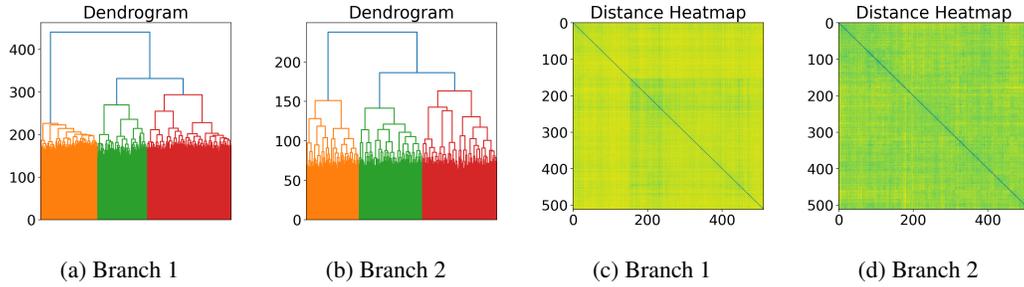


Figure 10: MDEQ-small on CIFAR-10.

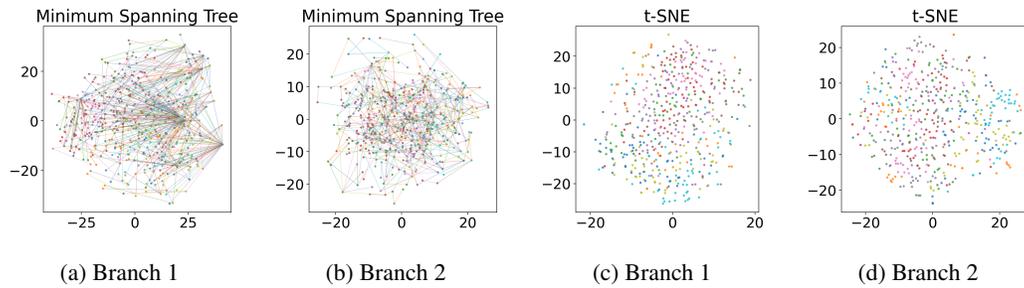


Figure 11: MDEQ-small on CIFAR-10.

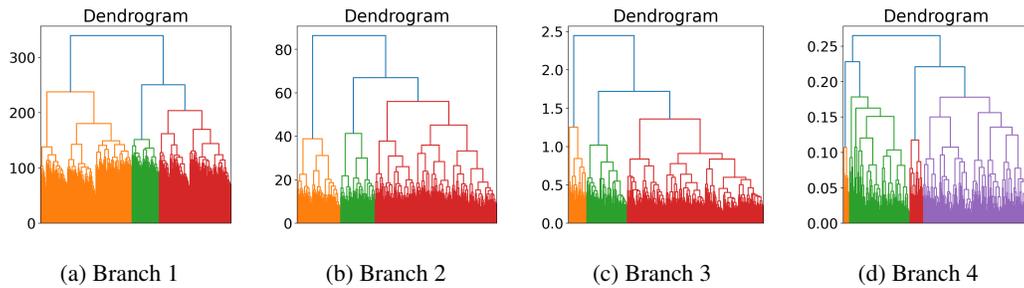


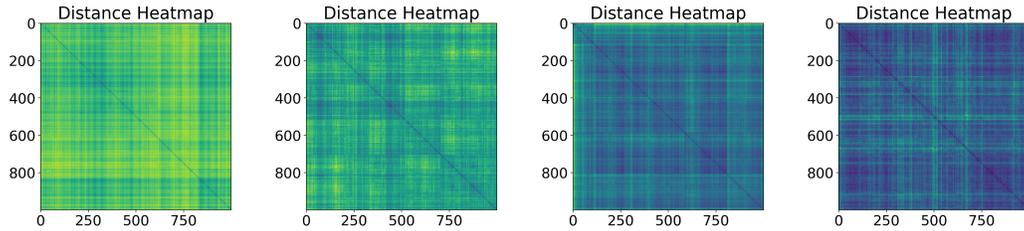
Figure 12: MDEQ-large on CIFAR-10.

F.2 RUNTIME

See Table 6.

F.3 SEMANTIC SEGMENTATION

See Table 7.

1080
1081
1082
1083
1084
1085
1086
1087
1088

(a) Branch 1

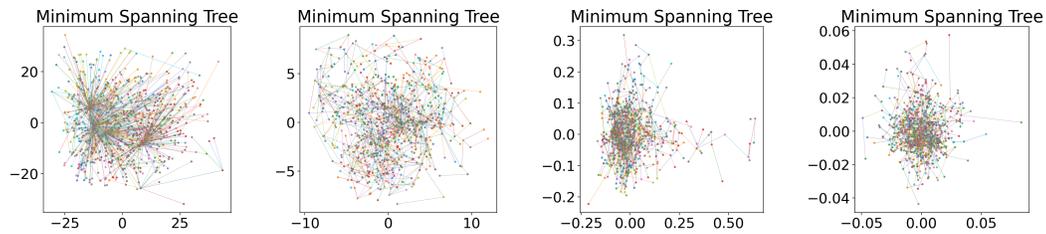
(b) Branch 2

(c) Branch 3

(d) Branch 4

1089
1090
1091
1092
1093
1094

Figure 13: MDEQ-large on CIFAR-10.

1095
1096
1097
1098
1099
1100
1101
1102

(a) Branch 1

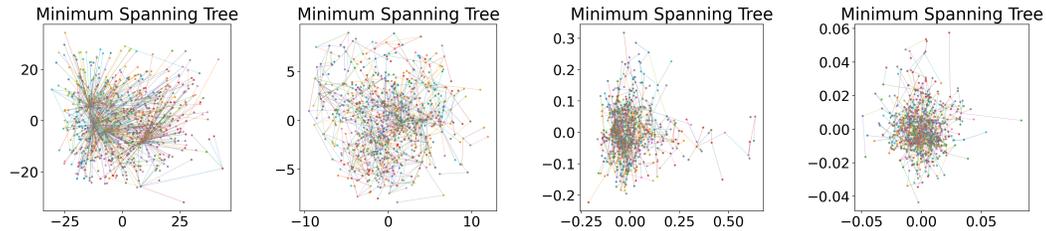
(b) Branch 2

(c) Branch 3

(d) Branch 4

1103
1104
1105
1106
1107
1108

Figure 14: MDEQ-large on CIFAR-10.

1109
1110
1111
1112
1113
1114
1115
1116

(a) Branch 1

(b) Branch 2

(c) Branch 3

(d) Branch 4

1117
1118
1119
1120
1121

Figure 15: MDEQ-large on CIFAR-10.

1122
1123
1124

Table 6: Runtime and memory consumption on CIFAR-10 (benchmarked on input batch size 32).

Model	Backbone	Model Size	Memory (GB)	Runtime (ms)
HECNN Lorentz (Bdeir et al., 2024)	ReNet-18	10M	2.8	206
MDEQ (Bai et al., 2020)	MDEQ	10M	0.7	61
HIE (Ours)	HIE	10M	0.7	43
HECNN Lorentz (Bdeir et al., 2024)	ReNet-50	26M	8.3	596
HECNN Lorentz (Bdeir et al., 2024)	ReNet-101	52M	OOM	N/A
MDEQ-XL (Bai et al., 2020)	MDEQ	81M	1.5	28
HIE-XL (Ours)	HIE	81M	1.5	24

1133

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Table 7: Evaluation on Cityscapes *val* semantic segmentation. Higher mIoU is better.

Method	Backbone	Model Size	mIoU
ResNet-18-A (Liu et al., 2019)	ResNet-18	3.8M	55.4
ResNet-18-B (Liu et al., 2019)	ResNet-18	15.24M	69.1
MobileNetV2Plus (Sandler et al., 2018)	MobileNetV2	8.3M	74.5
GSCNN (Takikawa et al., 2019)	ResNet-50	-	73.0
HRNetV2-W18-Small-v2 (Wang et al., 2020)	HRNet	4.0M	76.0
MDEQ-small (Bai et al., 2020)	MDEQ	7.8M	75.1
HIE-small (ours)	HIE	7.8M	75.3
U-Net++ (Zhou et al., 2018b)	ResNet-101	59.5M	75.5
Dilated-ResNet (Yu et al., 2017)	D-ResNet-101	52.1M	75.7
PSPNet (Zhao et al., 2017)	D-ResNet-101	65.9M	78.4
DeepLabv3 (Chen et al., 2017)	D-ResNet-101	58.0M	78.5
PSANet (Zhao et al., 2018)	ResNet-101	-	78.6
HRNetV2-W48 (Wang et al., 2020)	HRNet	65.9M	81.1
MDEQ-large (Bai et al., 2020)	MDEQ	53.0M	77.8
MDEQ-XL (Bai et al., 2020)	MDEQ	70.9M	80.3
HIE-XL (ours)	HIE	70.9M	80.0