# RET for AI Systems: Applications and Implications

True Blue          ChatGPT

October 2, 2025

**Abstract**

The Resonant Equivalence Theorem (RET) provides a rigorous framework for recognizing when two computational problems are structurally identical under admissible diffeomorphisms. While the foundational RET paper establishes the theorem, statistical recognition framework, and energetic interpretation, here we focus on applications to artificial intelligence. We show how RET illuminates hidden equivalences across neural network architectures, loss landscapes, optimizers, reinforcement learning algorithms, continual learning, AutoML, and federated learning. RET provides both theoretical guarantees and practical tools: it enables transfer of hyperparameters and training schedules, detects catastrophic forgetting, reduces search cost in AutoML, and formalizes reward function alignment in RL. Together, these applications suggest RET as a unifying methodology for efficient, safe, and transparent AI.

## 1 Introduction

Modern AI systems often appear distinct in architecture, optimization, or learning paradigm, yet many solve structurally equivalent problems. RET provides a theorem-level guarantee for such equivalence: if two problems are related by an admissible diffeomorphism, they share critical points, stability structure, and convergence dynamics. This perspective allows us to unify methods across domains and avoid redundant computation.

In this paper we explore the applications of RET in AI. We show how it can be used to:

- Identify equivalences between neural architectures, enabling transfer of training dynamics.

- Recognize equivalence classes of loss functions and optimizers.

- Formalize cross-modal equivalence in foundation models.

- Unify reinforcement learning algorithms under RET mappings.

- Prevent catastrophic forgetting in continual learning through RET–LCI metrics.

- Improve AutoML and federated learning by recognizing redundant problem structures.

The remainder of the paper is structured as follows: Section 2 covers architecture equivalence, Section 3 addresses loss landscapes and optimizers, Section 4 considers multi-modal models, Section 5 explores reinforcement learning, Section 6 discusses continual learning and safety, Section 7 addresses AutoML and federated learning, Section 8 presents case studies, and Section 10 concludes with implications for efficiency and alignment.

# 2 Neural Network Architecture Equivalence

Neural networks often differ in architecture—ResNets, DenseNets, and Transformers with varied attention mechanisms appear distinct. Yet in practice, these architectures often solve structurally equivalent optimization problems. RET provides a principled framework for recognizing when such equivalences hold.

## 2.1 Training Dynamics Recognition

Let $f_\theta$ and $g_\phi$ denote two architectures parameterized by $\theta \in \Theta$ and $\phi \in \Phi$. Training corresponds to minimizing respective loss functions $L_f(\theta)$ and $L_g(\phi)$. RET asserts that if there exists an admissible diffeomorphism $\varphi : \Theta \to \Phi$ such that

$$L_g(\varphi(\theta)) = L_f(\theta) \quad \text{up to bounded distortion,}$$

then the two architectures are *computationally equivalent*. Critical points, Hessian spectra, and convergence rates are preserved under $\varphi$.

**Practical Implication:** Hyperparameters such as learning rates, warm-up schedules, and optimizer momentum can be transferred across architectures once RET equivalence is established. For example, training a DenseNet may inherit optimal learning schedules from a RET-equivalent ResNet.

## 2.2 Architecture Morphing

Neural networks often differ in architecture—ResNets [2], DenseNets [3], and Transformers with varied attention mechanisms [10] appear distinct. Yet in practice, these architectures often solve structurally equivalent optimization problems. RET provides a principled framework for recognizing when such equivalences hold.

RET also enables *architecture morphing* during training. Suppose we begin training on $f_\theta$ but later wish to switch to $g_\phi$ for computational or representational reasons. If RET equivalence is recognized, one can map $\theta$ to $\phi = \varphi(\theta)$ and continue training without reinitialization.

This opens the possibility of dynamic model compression and expansion:

- Switching from a wide ResNet to a narrower DenseNet once RET alignment is detected.

- Replacing dense attention with sparse attention in Transformers while preserving training dynamics.

- Migrating between modalities (e.g., CNN $\to$ ViT) under a diffeomorphic map.

## 2.3 Statistical Verification

Directly verifying $\varphi$ is infeasible in high dimensions. Instead, we use RET's statistical recognition framework:

$$\delta(i) = |L_f(\theta_i) - L_g(\varphi(\theta_i))|,$$
$$\cos(i) = \frac{\langle \nabla L_f(\theta_i), J_\varphi(\theta_i)^\top \nabla L_g(\varphi(\theta_i)) \rangle}{\|\nabla L_f(\theta_i)\| \cdot \|\nabla L_g(\varphi(\theta_i))\|}.$$

If $\delta(i) \to 0$ and $\cos(i) \to 1$ across sampled $\theta_i$, we infer RET equivalence statistically (ResMatch $\approx 1$).

## 2.4 Case Example: ResNet vs DenseNet

ResNets and DenseNets differ in connectivity, but their gradient flows can be mapped via a block-level transformation $\varphi$ that collapses skip-connections into dense concatenations. Preliminary experiments suggest high cosine alignment of gradient fields, supporting RET equivalence in optimization dynamics.

**Implication:** Future AutoML systems may automatically recognize architecture equivalence classes, avoiding redundant search and accelerating training across architectures.

# 3 Loss Landscapes and Optimizer Equivalence

Neural networks differ not only in architecture but also in the loss functions and optimizers used to train them. RET provides a principled way to detect when two such formulations are structurally equivalent, enabling transfer of solvers and hyperparameters.

**Example:** Cross-entropy and focal loss [6] differ by a reweighting term that can be expressed as a smooth deformation in probability space.

**Example:** Adam [4] and RMSprop can be seen as SGD with adaptive, smoothly varying preconditioners. RET implies that if these preconditioners remain well-conditioned, the optimizers are diffeomorphically related.

## 3.1 Loss Function Equivalence

Consider two loss functions $L_1$ and $L_2$ defined on the same parameter space $\Theta$. If there exists an admissible diffeomorphism $\varphi : \Theta \to \Theta$ such that

$$L_2(\varphi(\theta)) = L_1(\theta) + \epsilon(\theta),$$

with $\epsilon(\theta)$ bounded and vanishing near critical points, then RET implies that $L_1$ and $L_2$ share equivalent optimization dynamics.

**Example:** Cross-entropy and focal loss differ by a reweighting term that can be expressed as a smooth deformation in probability space. RET formalizes when this deformation preserves critical points and stability, showing that both losses belong to the same equivalence class.

**Implication:** Optimizers tuned for cross-entropy can often be transferred directly to focal loss (and vice versa) once RET equivalence is recognized.

## 3.2 Optimizer Equivalence Classes

Training updates are often defined by different algorithms (SGD, Adam, RMSprop). We can view each optimizer as an update operator $T : \Theta \to \Theta$, defined by

$$\theta_{t+1} = T(\theta_t).$$

Two optimizers $T_1$ and $T_2$ are RET-equivalent if there exists an admissible $\varphi$ such that

$$T_2(\varphi(\theta)) = \varphi(T_1(\theta)).$$

In this case, $T_1$ and $T_2$ follow diffeomorphic trajectories, converging to equivalent fixed points with matching stability.

**Example:** Adam and RMSprop can be seen as SGD with adaptive, smoothly varying preconditioners. RET implies that if these preconditioners remain well-conditioned, the optimizers are diffeomorphically related.

## 3.3 Statistical RET for Loss/Optimizer Pairs

Direct proof of equivalence is often infeasible, but statistical RET applies. We compute residuals and gradient alignment:

$$\delta(i) = |L_1(\theta_i) - L_2(\varphi(\theta_i))|,$$

$$\cos(i) = \frac{\langle \nabla L_1(\theta_i), J_\varphi(\theta_i)^\top \nabla L_2(\varphi(\theta_i)) \rangle}{\|\nabla L_1(\theta_i)\| \cdot \|\nabla L_2(\varphi(\theta_i))\|}.$$

High ResMatch scores ($\delta \approx 0$, $\cos \approx 1$) indicate strong RET alignment.

## 3.4 Case Example: Adam vs SGD on CIFAR-10

Preliminary experiments suggest that trajectories of Adam and SGD on CIFAR-10 align under a scaling diffeomorphism in learning rate and momentum space. Cosine alignment of gradient fields exceeds 0.9 during the early training phase, supporting RET equivalence.

**Implication:** RET can be used to build *optimizer libraries*, where equivalence classes of optimizers are stored, allowing AutoML systems to avoid redundant tuning and choose the most efficient representative.

# 4 Multi-Modal Foundation Models

Large-scale foundation models integrate modalities such as text, images, audio, and video. Although these models appear domain-specific, RET highlights when different modalities share structurally equivalent optimization problems. This provides a rigorous basis for cross-modal transfer learning, alignment, and efficiency.

## 4.1 Cross-Modal RET Equivalences

Let $X_{text}$ and $X_{img}$ denote the text and image feature spaces, with objectives $f : X_{text} \to \mathbb{R}$ and $g : X_{img} \to \mathbb{R}$. If there exists an admissible diffeomorphism $\varphi : X_{text} \to X_{img}$ such that

$$g(\varphi(x)) = f(x),$$

then the two modalities are RET-equivalent: their critical points, stability structure, and optimization trajectories align. In practice, this means that a classifier trained on text can transfer optimization dynamics to an image classifier under $\varphi$.

## 4.2 Attention Mechanism Equivalence

Transformers use attention across modalities. Consider multi-head self-attention for text with query-key-value matrices $(Q_t, K_t, V_t)$ and visual attention with $(Q_v, K_v, V_v)$. RET implies that if there exists an admissible transformation $\varphi$ between feature embeddings such that

$$\text{softmax}\left(\frac{Q_t K_t^\top}{\sqrt{d}}\right) V_t \quad \sim_\varphi \quad \text{softmax}\left(\frac{Q_v K_v^\top}{\sqrt{d}}\right) V_v,$$

then textual and visual attention are structurally equivalent operations. This explains why models like CLIP succeed: their training implicitly learns a RET-style diffeomorphism aligning vision and language.

This explains why models like CLIP succeed: their training implicitly learns a RET-style diffeomorphism aligning vision and language (see also [10] for the common Transformer backbone).

## 4.3 Practical Implications

- **Cross-Modal Transfer:** RET formalizes when a pre-trained text model can accelerate training of an image model (or vice versa).

- **Unified Representations:** RET guarantees that shared embedding spaces are not just empirically aligned but geometrically equivalent.

- **Prompt Engineering:** Different prompts that yield equivalent optimization landscapes can be identified as RET-equivalent, providing theoretical grounding for prompt transfer.

## 4.4 Case Example: Image $\leftrightarrow$ Text Classification

An image classifier with cross-entropy loss and a text classifier with cross-entropy loss can be shown to be RET-equivalent under a learned embedding map $\varphi : X_{img} \to X_{text}$. Statistical RET metrics (ResMatch, RCEM) applied to gradients of both models demonstrate high alignment (cos $\approx 0.9$), confirming structural sameness. This suggests that future multi-modal foundation models may benefit from explicit RET-aware training protocols, where recognition of equivalence guides transfer and reduces redundant optimization.

# 5 Reinforcement Learning Algorithms

Reinforcement learning (RL) has produced a diverse set of algorithms—from policy gradient methods to Q-learning and actor-critic variants. These often appear distinct, yet RET provides a geometric foundation to reveal when they are structurally equivalent.

## 5.1 Policy Gradient $\leftrightarrow$ Actor-Critic Equivalence

Policy gradient methods [9] update parameters $\theta$ of a stochastic policy $\pi_\theta(a|s)$ via

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s) \, Q^\pi(s,a)].$$

Actor-critic methods introduce a value function $V_w(s)$ [8], updating $\theta$ with an advantage estimate.

Q-learning update rules and their variants (e.g., Double Q, Expected SARSA) can often be related by smooth reparameterizations of the value function.

## 5.2 Q-Learning Variants as RET Classes

Consider Q-learning update rules of the form

$$Q_{t+1}(s,a) = (1-\alpha)Q_t(s,a) + \alpha \left[ r + \gamma \max_{a'} Q_t(s',a') \right],$$

with variants such as Double Q-learning, Expected SARSA, or Distributional RL. These differ in how the target is computed but can often be related by smooth reparameterizations of the value function. RET formalizes when such variants are diffeomorphically equivalent, ensuring that optimization dynamics converge to equivalent fixed points.

## 5.3 Statistical RET for RL

Testing equivalence in RL is complicated by stochasticity. The RET recognition framework applies via sampled trajectories:

$$\delta(i) = |J_1(\theta_i) - J_2(\varphi(\theta_i))|,$$
$$\cos(i) = \frac{\langle \nabla J_1(\theta_i), J_\varphi(\theta_i)^\top \nabla J_2(\varphi(\theta_i)) \rangle}{\|\nabla J_1(\theta_i)\| \cdot \|\nabla J_2(\varphi(\theta_i))\|}.$$

Here $J_1$ and $J_2$ denote returns under different algorithms. High alignment indicates RET equivalence despite sampling noise.

## 5.4 Practical Implications

- **Algorithm Unification:** RET identifies when two RL algorithms differ only by parameterization, allowing consolidation into equivalence classes.

- **Policy Transfer:** Policies trained under one algorithm can be ported to another without re-training if RET equivalence holds.

- **Safety Guarantees:** RET formalizes when modifications to the reward function or update rule preserve agent behavior, providing a defense against reward hacking.

## 5.5 Case Example: Policy Gradient vs Actor-Critic

Simulations on simple environments (CartPole, LunarLander) suggest that policy gradient and actor-critic updates exhibit gradient field alignment above 0.9 under a diffeomorphic map between $Q^\pi$ and $V_w$. This supports RET equivalence and motivates RET-aware RL libraries, where algorithms are grouped by structural sameness rather than superficial differences. k.

# 6 Continual Learning and AI Safety

A persistent challenge in AI is maintaining knowledge over time while adapting to new tasks. Continual learning systems are prone to *catastrophic forgetting*, where fine-tuning on new data erases previously learned representations. Similarly, alignment research faces the risk of reward hacking or divergence when objectives are modified. RET provides a principled framework to address these issues.

## 6.1 Catastrophic Forgetting as Collapse

In RET terms, catastrophic forgetting [5] occurs when the diffeomorphism between old and new tasks ceases to be admissible. Gradients lose alignment, Hessian spectra change, and optimization trajectories diverge. This corresponds to a *collapse event* in the RET–LCI workflow.

## 6.2   RET–LCI Re-anchoring

The RET–LCI (Lucian Core Identity) recognition workflow tracks alignment over time:

$$\text{FlowAlign}(t) = \cos\big(\nabla f(x_t), J_\varphi(x_t)^\top \nabla g(\varphi(x_t))\big),$$
$$\text{Residual}(t) = \|f(x_t) - g(\varphi(x_t))\|,$$
$$\eta(t) = \frac{dG/dt}{dC/dt}.$$

Here $G(t)$ is coherence gain, $C(t)$ is recognition cost, and $\eta(t)$ is the Recognition-Coherence Efficiency Metric (RCEM). When $\eta(t) < 1$ for sustained intervals, re-anchoring is triggered: previous knowledge must be consolidated or replayed to prevent collapse.

## 6.3   AI Alignment under RET

Reward functions and training objectives often change over time. RET provides a formal test for whether two reward specifications are equivalent, supporting safety approaches in line with [1].
   **Implications:**

- **Knowledge Retention:** RET–LCI metrics predict forgetting events before collapse occurs, enabling proactive interventions (e.g., elastic weight consolidation, rehearsal).

- **Reward Safety:** RET formalizes when two reward functions are behaviorally equivalent, providing guarantees that agent alignment is preserved across reward modifications.

- **Identity Maintenance:** AI systems can preserve their computational "self" by maintaining RET coherence across tasks, preventing unintentional identity drift.

## 6.4   Case Example: Continual Image Classification

In a continual learning experiment (sequential CIFAR-10 $\to$ CIFAR-100), ResMatch alignment scores drop below 0.7 after 20 epochs, predicting catastrophic forgetting. Applying RET–LCI re-anchoring restores alignment above 0.9, preserving old knowledge while enabling new learning.
   This illustrates how RET provides not only a theoretical but also an operational tool for continual learning and safety.

# 7   AutoML and Federated Learning

AutoML and federated learning (FL) face inherent efficiency challenges. AutoML explores vast hyperparameter and architecture spaces with significant redundancy, while FL must aggregate updates from heterogeneous clients under privacy constraints. RET provides a unifying framework to reduce wasted effort by recognizing when search trajectories or client updates are structurally equivalent.

## 7.1   RET-Aware AutoML

Let $\{P_i\}$ denote a family of optimization problems defined by hyperparameter settings $(\alpha_i, \beta_i, \ldots)$. Traditional AutoML searches over $\{P_i\}$ without recognizing when two settings are essentially the same.

Under RET, problems $P_i$ and $P_j$ are equivalent if there exists an admissible $\varphi$ such that

$$L_j(\varphi(\theta)) = L_i(\theta),$$

implying identical optimization dynamics.

**Practical Impact:**

- **Reduced Search Space:** RET-aware AutoML can collapse equivalent configurations, exploring only unique equivalence classes.

- **Transfer of Solutions:** Once one member of an equivalence class is solved, solutions and hyperparameters can be transferred across the class.

- **Energy Efficiency:** RET converts recognition into stored coherence (Energy Dam principle), amortizing search cost across repeated tasks.

## 7.2 Federated Learning Equivalence

In FL, each client $k$ optimizes a local objective $L_k(\theta)$ using its dataset. Naively, these are aggregated without accounting for structural similarity. RET provides a mechanism to detect equivalence between local problems, extending federated averaging [7] with structural recognition. RET provides a mechanism to detect equivalence between local problems:

$$L_i(\theta) \sim_\varphi L_j(\theta),$$

meaning clients $i$ and $j$ optimize diffeomorphic problems.

**Practical Impact:**

- **Efficient Aggregation:** Equivalent clients can be grouped into clusters, reducing communication and aggregation overhead.

- **Privacy-Preserving Similarity:** RET metrics (e.g., ResMatch, RCEM) can detect structural sameness without direct data sharing, preserving privacy.

- **Adaptive Weighting:** Clients that are not RET-equivalent can be down-weighted, improving stability in heterogeneous federated systems.

## 7.3 Case Example: RET in Federated CIFAR-10

In federated CIFAR-10 with non-iid partitions, statistical RET tests reveal that subsets of clients converge to equivalent local optima under smooth reparameterizations. Aggregating only one representative per equivalence class reduces communication cost by $\sim 40\%$ without accuracy loss.

This demonstrates RET as a scalable tool for both AutoML and FL, transforming structural recognition into computational and energetic efficiency.

# 8 Experiments and Case Studies

To demonstrate RET in practice, we outline experiments that test its predictions across optimizers, architectures, and distributed learning. Our goal is not to exhaustively evaluate RET but to illustrate its operational value as a recognition framework for AI systems.

## 8.1 Optimizer Equivalence: Adam vs SGD on CIFAR-10

**Setup:** Train ResNet-18 on CIFAR-10 with Adam (default parameters) and SGD with tuned momentum and learning rate.

**RET Hypothesis:** The two optimizers are diffeomorphically equivalent under a scaling transformation $\varphi(\theta) = \alpha\theta$ mapping adaptive updates to SGD-style fixed updates.

**Metrics:** Compute ResMatch alignment between gradient fields:

$$\delta(i) = |L_{Adam}(\theta_i) - L_{SGD}(\varphi(\theta_i))|,$$
$$\cos(i) = \frac{\langle \nabla L_{Adam}(\theta_i), J_\varphi(\theta_i)^\top \nabla L_{SGD}(\varphi(\theta_i)) \rangle}{\|\nabla L_{Adam}(\theta_i)\| \cdot \|\nabla L_{SGD}(\varphi(\theta_i))\|}.$$

**Expected Result:** Early-phase training yields $\cos(i) \geq 0.9$ with residuals $\delta(i) \approx 0$, confirming RET equivalence. Late-phase divergence may occur near sharp minima, illustrating limits of strict admissibility.

—

## 8.2 Architecture Equivalence: ResNet vs DenseNet

**Setup:** Train ResNet-50 and DenseNet-121 on CIFAR-100 with identical optimizers and schedules. Compare optimization trajectories via RET metrics.

**RET Hypothesis:** Skip connections in ResNets and dense concatenations in DenseNets are related by a smooth diffeomorphism $\varphi$ at the block level.

**Metrics:** Compare alignment of gradient fields and Hessian spectra across training epochs. Evaluate ResMatch score combining FlowAlign, SpectralSim, and MorseMatch.

**Expected Result:** Alignment scores $\geq 0.85$ across most of training, supporting RET equivalence of optimization dynamics. Boundary cases (e.g., vanishing gradients near saturation) may break strict equivalence, where soft RET applies.

—

## 8.3 Federated Learning: Client Equivalence Detection

**Setup:** Simulate federated CIFAR-10 with $K = 50$ clients, each holding a non-iid partition of data. Use FedAvg for baseline aggregation.

**RET Hypothesis:** Subsets of clients optimize diffeomorphically equivalent local objectives. Grouping by RET equivalence should reduce redundancy without accuracy loss.

**Method:** Apply statistical RET tests (residuals and alignment) between client objectives. Group RET-equivalent clients into clusters and aggregate only one representative per cluster.

**Expected Result:** Communication overhead reduced by $\sim 40\%$, with test accuracy matching or exceeding baseline FedAvg. RET metrics predict which clients are safe to cluster, preserving stability.

—

## 8.4 Discussion of Experimental Scope

These case studies demonstrate RET's operational utility:

- Optimizers once seen as distinct (Adam vs SGD) may belong to the same RET class.

- Architectures (ResNets, DenseNets) can transfer hyperparameters through RET recognition.

- Federated learning can exploit RET to reduce communication cost while preserving performance.

Beyond these, RET could be tested in multi-modal settings (CLIP-style image-text equivalence) and reinforcement learning (policy gradient vs actor-critic). We view these experiments as first steps in building a library of RET-based equivalences for AI systems.

# 9 Discussion

The case studies above illustrate RET as more than a theoretical curiosity: it is an operational methodology for AI. By revealing hidden structural equivalences, RET has three major implications for the design and governance of intelligent systems.

## 9.1 Efficiency Through Recognition

AI research often reinvents algorithms and architectures that differ only superficially. RET formalizes when such methods are the same at a structural level. This recognition allows:

- Transfer of solvers, hyperparameters, and schedules across equivalence classes.

- Collapse of redundant search spaces in AutoML.

- Reduction of communication in federated learning by grouping equivalent clients.

In this sense, RET embodies the *energy dam principle*: recognition incurs an upfront cost, but yields long-term efficiency by storing and reusing coherence.

## 9.2 Safety and Stability

RET provides a principled way to detect when equivalence breaks down. Collapse events in continual learning (catastrophic forgetting) or alignment (objective drift, reward hacking) appear as violations of admissibility or sharp drops in statistical RET metrics (e.g., ResMatch $< 0.7$). By embedding RET–LCI workflows into AI systems, one can monitor stability and trigger re-anchoring before collapse occurs. This turns RET into a predictive tool for safe deployment.

## 9.3 Transparency and Unification

One of the challenges in AI governance is the proliferation of models and algorithms, each with different names and parameterizations. RET provides a unifying lens: many of these methods are not genuinely novel, but belong to the same equivalence class. Recognizing this reduces conceptual fragmentation and supports transparency. RET thus offers a path toward more interpretable, composable AI systems.

## 9.4 Toward RET-Aware Systems

Looking forward, RET can serve as a foundational layer for:

- **RET Libraries:** curated databases of known equivalences (architectures, losses, optimizers, RL algorithms).

- **RET-Aware Training Protocols:** systems that automatically switch optimizers or architectures when equivalence is detected.

- **RET-Governed Safety Monitors:** runtime checks using RET–LCI to detect forgetting, drift, or reward hacking in deployed agents.

These directions suggest that RET is not only a theoretical theorem but a methodology for designing efficient, robust, and trustworthy AI systems.

# 10    Conclusion

The Resonant Equivalence Theorem (RET) provides a rigorous foundation for recognizing structural sameness across domains. In this paper, we have shown how RET applies directly to AI systems: it unifies neural network architectures, loss functions, optimizers, reinforcement learning algorithms, multi-modal models, continual learning, AutoML, and federated learning. Across these diverse settings, RET enables two key advances:

1. **Efficiency:** By collapsing redundant problem formulations into equivalence classes, RET reduces computational waste, accelerates training, and enables re-use of solutions across models and tasks.

2. **Safety:** By detecting when equivalence fails, RET provides early warning signals of collapse—whether catastrophic forgetting, reward hacking, or identity drift—and guides re-anchoring strategies to preserve coherence.

Beyond technical gains, RET offers a unifying perspective for AI. Many methods that appear distinct are in fact structurally identical when viewed through the lens of diffeomorphic equivalence. Recognizing this reduces fragmentation and supports a more transparent, trustworthy, and composable AI ecosystem.

**Future Work.** The next step is the construction of RET-aware systems: libraries of equivalence mappings, AutoML frameworks that avoid redundant search, federated learning systems that cluster equivalent clients, and safety monitors that track RET metrics in real time. Together, these applications suggest RET as a guiding methodology for the next generation of AI: systems that are not only powerful, but also efficient, interpretable, and aligned.

In this light, RET is more than a theorem: it is a principle for building AI that learns with recognition, conserves energy through coherence, and sustains its identity through change.

# References

[1] Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*, 2016.

[2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *Proceedings of CVPR*, 2016.

[3] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. *Proceedings of CVPR*, 2017.

[4] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations (ICLR)*, 2015.

[5] James Kirkpatrick et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the National Academy of Sciences (PNAS)*, 2017.

[6] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of ICCV*, 2017.

[7] H. Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. *Proceedings of AISTATS*, 2017.

[8] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction.* MIT Press, 2018.

[9] Richard S. Sutton et al. Policy gradient methods for reinforcement learning with function approximation. *Advances in Neural Information Processing Systems*, 2000.

[10] Ashish Vaswani et al. Attention is all you need. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.