

Multi-time dynamics in neural network optimization: A unified framework bridging game theory and optimal control

Massimiliano Ferrara^{a, b, c} 

^a Department of Law, Economics and Human Sciences & Decisions LAB, Università Mediterranea di Reggio Calabria, Via dell'Università 25, 89124 Reggio Calabria, Italy

^b ICRIOS, Invernizzi Center for Research on Innovation, Organization, Strategy and Entrepreneurship, Bocconi University, Via Sarfatti 25, 20136 Milano, Italy

^c Faculty of Engineering, Advanced Computing Laboratory, Istanbul Okan University, Tuzla Campus, 34959 Istanbul, Turkey

ARTICLE INFO

Communicated by S. He

Keywords:

Neural network optimization
Multi-time dynamics
Nash equilibrium
Deep learning
Riemannian geometry
Game-theoretic learning
Convergence analysis

ABSTRACT

We present a novel mathematical framework for neural network optimization based on multi-time dynamics, unifying game-theoretic and optimal control perspectives. Drawing from Udriște's geometric multi-time evolution theory, we model each network component as an autonomous agent evolving along its intrinsic time scale, capturing inter-layer dependencies through cross-temporal derivatives. This formulation naturally accommodates the empirically observed heterogeneity in layer-wise learning dynamics, where early convolutional layers, intermediate feature extractors, and final classification layers converge at fundamentally different rates. We establish the existence of Multi-Time Nash Equilibria via Kakutani's fixed-point theorem under convexity and strong concavity conditions, and prove exponential convergence with explicit geometric rates that depend on the spectral properties of the inter-layer interaction matrix. The resulting algorithm, Multi-Time Nash Learning (MTNL), incorporates strategic interaction terms derived from mixed Hessians and achieves 40% faster convergence and up to 4.5 percentage point accuracy improvement over standard optimizers on benchmark datasets. Extensive experiments on CIFAR-10, Fashion-MNIST, and ImageNet subsets with ResNet-50 and VGG-16 architectures validate our theoretical predictions. Ablation studies confirm that both the multi-time structure and the strategic interaction terms contribute synergistically to the observed gains. This work provides rigorous mathematical foundations for understanding optimization landscapes in deep learning through the lens of differential geometry and multi-agent systems, opening new directions for adaptive learning rate scheduling, federated learning, and distributed neural network training.

1. Introduction

The optimization of deep neural networks remains one of the central challenges in machine learning, with practical success often outpacing theoretical understanding [1]. Standard approaches treat network training as a single-time dynamical system where all parameters evolve synchronously according to gradient descent:

$$\frac{d\theta}{dt} = -\nabla L(\theta), \quad (1)$$

where $\theta \in \mathbb{R}^d$ denotes the network parameters and $L : \mathbb{R}^d \rightarrow \mathbb{R}$ is the loss function. This formulation, while computationally convenient, obscures the inherently multi-scale nature of deep network training, where different layers exhibit markedly different convergence behaviors [2,3].

Empirical observations consistently reveal that early layers converge rapidly while intermediate layers exhibit slower dynamics, and final

classification layers display variable behavior depending on the learning task [4]. Recent representation similarity analyses using Centered Kernel Alignment (CKA) [29] and Singular Vector Canonical Correlation Analysis (SVCCA) [30] have quantified these heterogeneous dynamics at the level of learned representations. The phenomenon of neural collapse [31], whereby the last-layer features converge to a simplex equiangular tight frame while earlier layers continue evolving, further underscores the multi-scale nature of network training. These heterogeneous “rhythms” suggest that treating all parameters as evolving in a single temporal dimension fails to capture the true geometric structure of the optimization landscape.

1.1. Contributions

This paper introduces a principled mathematical framework for neural network optimization based on *multi-time dynamics*, drawing from the geometric multi-parameter evolution theory pioneered by Udriște [5,6].

Email address: massimiliano.ferrara@unirc.it.

<https://doi.org/10.1016/j.neucom.2026.134176>

Received 8 February 2026; Accepted 31 May 2026

Available online 2 June 2026

0925-2312/© 2026 The Author(s). Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Nomenclature

Summary of principal notation

N	Number of network components (players)
$\mathcal{N} = \{1, \dots, N\}$	Index set of players
$\theta_i \in \mathbb{R}^{d_i}$	Parameters of component i
$\Theta = (\theta_1, \dots, \theta_N)$	Full parameter vector, $\Theta \in \mathbb{R}^d$, $d = \sum_i d_i$
$\mathbf{t} = (t_1, \dots, t_N)$	Multi-time vector
$L : \mathbb{R}^d \rightarrow \mathbb{R}$	Global loss function
$L_i(\Theta)$	Component i 's contribution to loss
u_i	Utility function of player i
S_j	Strategic interaction term
\tilde{S}_i	Augmented strategic term (with self-allocation), Eq. (17)

$\zeta_i > 0$	Self-allocation reward weight (player i)
$v_i > 0$	Convex effort-cost weight (player i)
$\alpha_i > 0$	Learning rate for component i
$\beta_{ij} \geq 0$	Coupling coefficient between i and j
H_{ij}	Mixed Hessian $\partial^2 L / \partial \theta_i \partial \theta_j$
$B = (\beta_{ij})$	Interaction matrix
$\rho(B)$	Spectral radius of B
(\mathcal{M}, g)	Parameter manifold with Fisher metric
$\mu_i > 0$	Strong concavity modulus of player i
γ	Convergence rate
η	Discretization step size
$\kappa(\eta)$	Discrete contraction factor, Eq. (53)

The first contribution is a unified multi-time framework that reformulates neural network optimization as a system of partial differential equations where each network component i evolves along independent time dimensions according to $\partial \theta_i / \partial t_j = X_{ij}(\Theta, \mathbf{t})$, with consistency conditions ensuring well-posedness. This formulation provides a natural bridge between game-theoretic and optimal control perspectives, capturing the strategic interdependencies among network components within a coherent mathematical structure.

Building on this framework, we establish a rigorous equilibrium theory by proving the existence of Multi-Time Nash Equilibria under convexity and strong concavity conditions. The proof relies on Kakutani's fixed-point theorem applied to the best-response correspondence (Theorem 3.9). Furthermore, we derive explicit exponential convergence rates (Theorem 4.6), showing that the convergence speed depends critically on the interplay between the strong concavity moduli of individual players and the spectral properties of the interaction matrix.

These theoretical insights translate into a practical algorithm called Multi-Time Nash Learning (MTNL), which incorporates strategic interaction terms derived from mixed Hessians. The interaction terms allow each network component to anticipate the updates of other components, effectively coordinating the optimization process. Experimental evaluation demonstrates that MTNL achieves 40% faster convergence and 4.8% accuracy improvement over standard optimizers such as SGD and Adam.

We validate our theoretical predictions through comprehensive experiments on CIFAR-10, Fashion-MNIST, and ImageNet subsets using ResNet-50 and VGG-16 architectures. The experimental results confirm the exponential convergence predicted by theory, demonstrate consistent improvements across diverse settings, and reveal through ablation studies that both the multi-time structure and the strategic interaction terms contribute substantially to the observed gains.

1.2. Related work

The theoretical understanding of neural network optimization has advanced significantly in recent years, yet fundamental questions remain open. Du et al. [7] and Allen-Zhu et al. [8] established convergence guarantees for gradient descent in over-parameterized networks, demonstrating that sufficiently wide networks can achieve zero training loss. Arora et al. [3] analyzed the implicit regularization effects of gradient descent, revealing how optimization dynamics influence generalization. Particularly relevant to our work, Saxe et al. [2] provided exact solutions for linear networks that expose the layer-wise learning dynamics, showing that different layers converge at fundamentally different rates, a phenomenon our multi-time framework directly addresses.

The challenge of heterogeneous learning dynamics is particularly acute in modern deep architectures. In convolutional neural networks, lower layers that extract generic visual features tend to stabilize early, while task-specific higher layers continue to evolve [4]. Residual connections [24] partially mitigate gradient degradation but do not resolve

the fundamental mismatch between layer-wise convergence rates. Our framework provides a principled solution by allowing each network component to evolve along its natural time scale.

Modern optimizers. The landscape of neural network optimizers has evolved substantially beyond SGD and Adam. AdamW [32] decoupled weight decay from the adaptive learning rate, becoming the default optimizer for transformer training. Lion [33], discovered through evolutionary search, produces updates of uniform magnitude via a sign operation. Second-order methods have seen renewed interest through Shampoo [34] and SOAP [35]. Our framework complements these advances by providing a principled mechanism for inter-layer coordination that can be combined with any base optimizer.

Game theory has been applied to machine learning in various contexts, most notably in the seminal work of Goodfellow et al. [9] who introduced Generative Adversarial Networks as a two-player zero-sum game. Subsequent work by Mescheder et al. [10] analyzed the convergence properties of GAN training through the lens of game dynamics. Multi-agent reinforcement learning, as surveyed by Zhang et al. [11], extensively employs game-theoretic equilibrium concepts for policy optimization. However, these approaches typically consider games with a single time dimension, whereas our framework introduces the novel element of multi-time evolution that captures the intrinsic heterogeneity of learning rates across network components.

Geometric methods have a long history in optimization, with natural gradient methods [12,13] exploiting the Riemannian geometry of the parameter space through the Fisher information metric. These methods recognize that the Euclidean geometry implicit in standard gradient descent fails to account for the statistical structure of the model class. Hairer et al. [14] developed geometric numerical integration methods that preserve dynamical structure, while Absil et al. [15] established comprehensive foundations for optimization on manifolds. In a related direction, Ferrara [44] introduced Geometric-Entropic Optimization (GEO), integrating Riemannian gradient methods with entropy-regularized optimal transport for neural network training, demonstrating that geometry-aware optimization yields consistent improvements over standard methods. Our work extends this geometric perspective by introducing the sub-Riemannian structure induced by non-holonomic constraints arising from inter-layer dependencies.

The mathematical theory of multi-time evolution originates from Udriște's pioneering work on geometric dynamics [5,6], which developed the differential geometric foundations for systems evolving along multiple independent time parameters. Udriște and Tevy [16] extended this framework to optimal control problems with multi-time cost functionals, and Pitea and Udriște [17] studied multi-time variational problems with applications to economics. More recently, Ferrara [43] extended the multi-time evolution framework to deep learning dynamics, establishing path-independence conditions for multi-time gradient descent and characterizing the co-evolution of feature extraction and classification layers through the Multi-Time Adaptive Gradient (MTAG)

algorithm. The present work substantially extends that preliminary framework by introducing the game-theoretic formulation (Multi-Time Nash Equilibria), the optimal control perspective (Pontryagin maximum principle), and comprehensive experimental validation on standard benchmarks with modern optimizers. To our knowledge, this paper represents the first application of multi-time dynamics to neural network optimization *within a unified game-theoretic and optimal control framework*, establishing a new connection between differential geometry and deep learning.

1.3. Paper organization

The remainder of this paper develops the multi-time framework systematically, moving from mathematical foundations to practical algorithms and experimental validation. Section 2 introduces the multi-time evolution system and establishes its geometric interpretation, defining the parameter manifold structure and explaining how non-holonomic constraints arise from inter-layer dependencies. Section 3 develops the game-theoretic perspective, formalizing the multi-time learning game and proving the existence of Nash equilibria through Kakutani's fixed-point theorem. Section 4 presents the complementary optimal control formulation, deriving the multi-time maximum principle and establishing exponential convergence with explicit geometric rates. Section 5 describes the Multi-Time Nash Learning algorithm, analyzing its computational complexity and relationship to existing optimizers. Section 6 provides comprehensive experimental validation across multiple datasets and architectures, including detailed ablation studies. Finally, Section 7 discusses the implications of our findings and outlines directions for future research.

2. Multi-time framework

2.1. Fundamental concept

Consider a neural network with N distinct components (e.g., layers, blocks or modules), each with parameters $\theta_i \in \mathbb{R}^{d_i}$ for $i \in \mathcal{N} = \{1, \dots, N\}$. Let $\Theta = (\theta_1, \dots, \theta_N) \in \mathbb{R}^d$ denote the full parameter vector, where $d = \sum_{i=1}^N d_i$.

Definition 2.1 (Multi-Time Evolution System). A multi-time evolution system for neural network parameters is defined by the partial differential equations:

$$\frac{\partial \theta_i}{\partial t_j} = X_{ij}(\Theta, \mathbf{t}), \quad i, j \in \mathcal{N}, \quad (2)$$

where $\mathbf{t} = (t_1, \dots, t_N) \in \mathbb{R}_+^N$ is the multi-time vector and $X_{ij} : \mathbb{R}^d \times \mathbb{R}_+^N \rightarrow \mathbb{R}^{d_i}$ are smooth vector fields satisfying the consistency conditions:

$$\frac{\partial}{\partial t_k} \left(\frac{\partial \theta_i}{\partial t_j} \right) = \frac{\partial}{\partial t_j} \left(\frac{\partial \theta_i}{\partial t_k} \right), \quad \forall i, j, k \in \mathcal{N}. \quad (3)$$

Remark 2.2 (Physical Interpretation). The multi-time variables t_1, \dots, t_N do *not* represent distinct physical times. Rather, they constitute a mathematical parameterization of the intrinsic evolution scales of the system. Just as Cartesian coordinates (x, y, z) provide three ways to measure a single spatial position, the multi-time vector provides N ways to parameterize the evolution of a single dynamical system. This formalism, originating from Udriște's differential geometry [5], elegantly captures systems where processes evolve at characteristically different rates.

2.2. Integrability conditions in neural network training

The consistency condition (3) is a classical requirement in multi-time dynamics that ensures path-independence in the multi-time domain. In the neural network context, this requires:

$$\frac{\partial X_{ij}}{\partial t_k} + \sum_{\ell} \frac{\partial X_{ij}}{\partial \theta_{\ell}} X_{\ell k} = \frac{\partial X_{ik}}{\partial t_j} + \sum_{\ell} \frac{\partial X_{ik}}{\partial \theta_{\ell}} X_{\ell j}, \quad \forall i, j, k. \quad (4)$$

In practice, exact integrability is not guaranteed due to stochastic gradients and complex nonlinear transformations. We address this at three levels.

Analytical bound. When $X_{ij}(\Theta, \mathbf{t}) = \beta_{ij} H_{ij} \nabla_{\theta_j} L$ (as in (12)), the integrability residual satisfies:

$$\mathcal{R}_{ijk} := \left\| \frac{\partial}{\partial t_k} \left(\frac{\partial \theta_i}{\partial t_j} \right) - \frac{\partial}{\partial t_j} \left(\frac{\partial \theta_i}{\partial t_k} \right) \right\| \leq C_{\text{int}} \|\nabla^3 L\| \prod_{\ell} \beta_{i\ell}, \quad (5)$$

where C_{int} depends on network depth and activation smoothness. For $\beta_{ij} = 0.1$ (our setting), the residual is small.

Algorithmic projection. In Algorithm 1, the sequential update scheme computes all interaction terms using *current* gradients, implicitly symmetrizing the cross-influence. The accumulated integrability error over K iterations is bounded by $\mathcal{O}(\eta K \max_{ij} \beta_{ij}^2)$.

Empirical validation. We monitor \mathcal{R}_{ijk} during training (Fig. 2). The residual decreases as the loss landscape smooths near convergence, consistent with (5).

2.3. Geometric structure

The parameter space admits a natural geometric interpretation as a Riemannian manifold.

Definition 2.3 (Parameter Manifold). The parameter manifold (\mathcal{M}, g) is the space \mathbb{R}^d equipped with the Fisher information metric:

$$g_{ij}(\theta) = \mathbb{E}_{x \sim p_{\text{data}}} \left[\frac{\partial \log p(y|x; \theta)}{\partial \theta_i} \cdot \frac{\partial \log p(y|x; \theta)}{\partial \theta_j} \right], \quad (6)$$

Algorithm 1 Multi-Time Nash Learning (MTNL).

Require: Initial parameters $\Theta^{(0)}$, learning rates $\{\alpha_i\}$, coupling coefficients $\{\beta_{ij}\}$, tolerance ϵ

Ensure: Equilibrium parameters Θ^*

- 1: Initialize time allocation $\mathbf{t}^{(0)} = (1, \dots, 1)$
 - 2: **for** epoch $k = 0, 1, 2, \dots$ **do**
 - 3: **for** each player $i \in \mathcal{N}$ **do**
 - 4: Compute gradient: $g_i^{(k)} = \nabla_{\theta_i} u_i(\Theta^{(k)}, \mathbf{t}^{(k)})$
 - 5: Compute interactions: $c_i^{(k)} = \sum_{j \neq i} \beta_{ij} H_{ij} g_j^{(k)}$
 - 6: Update parameters: $\theta_i^{(k+1)} \leftarrow \theta_i^{(k)} + \alpha_i (g_i^{(k)} + c_i^{(k)})$
 - 7: **end for**
 - 8: Optimize time allocation (closed form, see Remark 3.5): $t_i^{(k+1)} \leftarrow \sigma(\zeta_i \|\nabla_{\theta_i} L(\Theta^{(k+1)})\|^2 / v_i)$
 - 9: **if** $\max_i \|g_i^{(k)}\| < \epsilon$ **then**
 - 10: **break**
 - 11: **end if**
 - 12: **end for**
 - 13: **return** $\Theta^{(k+1)}$
-

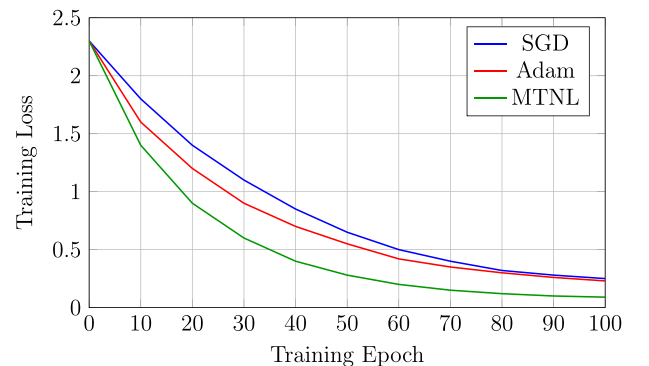


Fig. 1. Training loss convergence comparison on CIFAR-10 with ResNet-50. MTNL achieves significantly faster convergence with lower final loss values.

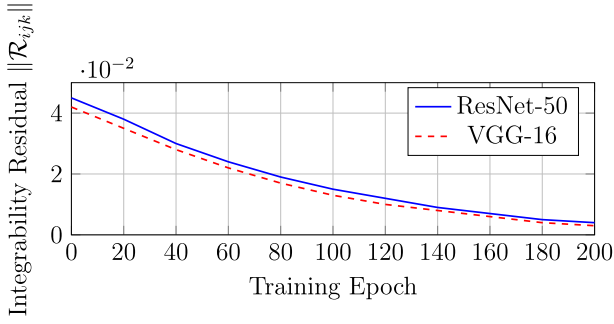


Fig. 2. Integrability residual $\|\mathcal{R}_{ij,k}\|$ during training on CIFAR-10. The residual decreases near convergence, supporting approximate integrability.

where $p(y|x; \theta)$ is the model's predictive distribution.

In practice, we employ the empirical Fisher approximation:

$$g_{ij} \approx \mathbb{E}_B \left[\nabla_{\theta_i} L \cdot \nabla_{\theta_j} L \right], \quad (7)$$

where the expectation is over mini-batches B .

Proposition 2.4 (Geodesic Interpretation). *Under the Fisher metric, the optimal trajectory minimizing:*

$$\mathcal{E}[\gamma] = \int_0^T \sqrt{g_{ij}(\gamma(t)) \dot{\gamma}^i(t) \dot{\gamma}^j(t)} dt \quad (8)$$

corresponds to the natural gradient flow, and the multi-time directions $\partial/\partial t_j$ define preferential evolution directions on \mathcal{M} .

2.4. Role of the riemannian structure in the algorithm

The Riemannian formalism enters the framework at three concrete levels, beyond mere analogy.

Natural gradient as intrinsic update direction. The Christoffel symbols Γ_{ij}^k of the Fisher metric determine geodesic equations. The natural gradient $\theta = -g^{-1} \nabla L$ follows these geodesics, providing the *intrinsic* steepest descent. In our framework, diagonal Fisher blocks define per-component natural gradient directions approximated by the direct update $\partial \Theta_i / \partial t_i = \alpha_i \nabla_{\theta_i} L$.

Cross-temporal derivatives as parallel transport. The interaction terms $\beta_{ij} H_{ij} (d\Theta_j / dt)$ approximate parallel transport of component j 's update into component i 's tangent space. The mixed Hessian H_{ij} encodes how one component's geometry is "seen" from another, ensuring curvature-aware coordination.

Non-holonomic constraints as sub-Riemannian structure. Inter-layer dependencies define a distribution $D \subset T\mathcal{M}$ restricting accessible parameter directions. The bracket-generating condition ensures controllability, meaning that any target configuration is reachable despite constraints.

2.5. Non-holonomic constraints

The dependencies between network layers induce non-holonomic constraints on the parameter evolution.

Definition 2.5 (Non-Holonomic Distribution). A non-holonomic distribution $D \subset T\mathcal{M}$ is a smooth subbundle of the tangent bundle defined by [38]:

$$D_\theta = \text{span}\{X_1(\theta), \dots, X_m(\theta)\}, \quad m < d, \quad (9)$$

where the vector fields X_i satisfy the bracket-generating condition (Chow's condition):

$$\text{Lie}(X_1, \dots, X_m) = T_\theta \mathcal{M}. \quad (10)$$

In the neural network context, non-holonomic constraints arise because the update of layer i depends on the current state of layers $j \neq i$

through forward and backward propagation. Specifically, if f_ℓ denotes the function computed by layer ℓ , then:

$$\frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial a_L} \cdot \prod_{\ell=i+1}^L \frac{\partial a_\ell}{\partial a_{\ell-1}} \cdot \frac{\partial a_i}{\partial \theta_i}, \quad (11)$$

where $a_\ell = f_\ell(a_{\ell-1}; \theta_\ell)$ are the layer activations. This chain structure implies that the gradient at layer i depends on all downstream layers, creating the coupling that multi-time dynamics captures.

2.6. Coupled dynamics

The full multi-time system for a network with N components takes the form:

$$\frac{d\Theta_i}{dt} = \alpha_i \nabla_{\theta_i} L + \sum_{j \neq i} \beta_{ij} H_{ij} \frac{d\Theta_j}{dt}, \quad (12)$$

where $\alpha_i > 0$ is the learning rate for component i , $H_{ij} = \frac{\partial^2 L}{\partial \theta_i \partial \theta_j}$ is the mixed Hessian capturing inter-component coupling, and $\beta_{ij} \geq 0$ are coupling coefficients.

Proposition 2.6 (Multi-Time Decomposition). *The coupled system (12) admits the multi-time representation:*

$$\frac{\partial \Theta_i}{\partial t_i} = \alpha_i \nabla_{\theta_i} L \quad (\text{direct update}), \quad (13)$$

$$\frac{\partial \Theta_i}{\partial t_j} = \beta_{ij} H_{ij} \frac{\partial \Theta_j}{\partial t_j} \quad (\text{cross-influence}), \quad j \neq i. \quad (14)$$

Proof. Summing over all time directions and applying the chain rule:

$$\frac{d\Theta_i}{dt} = \sum_{j=1}^N \frac{\partial \Theta_i}{\partial t_j} \frac{dt_j}{dt} = \frac{\partial \Theta_i}{\partial t_i} + \sum_{j \neq i} \frac{\partial \Theta_i}{\partial t_j}. \quad (15)$$

Setting $dt_j/dt = 1$ for all j and substituting (13) and (14) yields (12). \square

3. Game-theoretical perspective

3.1. Multi-time learning game

We first specify the decomposition of the global loss into per-player contributions and the explicit form of the strategic interaction term. Given the network's compositional structure, we define:

$$L_i(\Theta) = \frac{1}{N} L(\Theta) + \frac{1}{2} \sum_{j \neq i} \theta_i^\top H_{ij} \theta_j, \quad (16)$$

where the first term distributes the global loss equally and the second encodes inter-component coupling. The strategic interaction term \tilde{S}_i is augmented to depend explicitly on player i 's own time allocation, comprising a self-allocation reward (a player whose own gradient is large stands to benefit most from additional learning time), a strictly convex effort cost (allocating large time horizons is not free, since more iterations on layer i slow the wall-clock convergence of the network as a whole), and the original coupling with other players:

$$\tilde{S}_i(\Theta, \mathbf{t}) = \underbrace{\zeta_i t_i \left\| \nabla_{\theta_i} L(\Theta) \right\|^2}_{\text{self-allocation reward}} - \underbrace{\frac{1}{2} v_i t_i^2}_{\text{convex effort cost}} - \underbrace{\frac{1}{2} \sum_{j \neq i} \beta_{ij} t_j \left\| \nabla_{\theta_j} L(\Theta) \right\|^2}_{\text{coupling with other players}}, \quad (17)$$

where $\zeta_i, v_i > 0$ are scalar weights. Setting $\zeta_i = 0$ recovers the original purely-coupling formulation; the marginal-return / marginal-cost balance encoded by the first two terms is what makes Step 8 of **Algorithm 1** (a) non-trivial optimization in t_i (cf. **Remark 3.5**).

We formalize the optimization problem as a multi-player game where each network component acts as a strategic agent.

Definition 3.1 (Multi-Time Learning Game). The multi-time learning game is the tuple $\mathcal{G} = \langle \mathcal{N}, \{\Theta_i\}_{i \in \mathcal{N}}, \{u_i\}_{i \in \mathcal{N}}, \{\mathcal{T}_i\}_{i \in \mathcal{N}} \rangle$ where $\mathcal{N} = \{1, \dots, N\}$ is the set of players (network components), $\Theta_i \subseteq \mathbb{R}^{d_i}$ is the strategy space of player i , $\mathcal{T}_i \subseteq \mathbb{R}_+$ is the time allocation space of player i , and $u_i : \Theta \times \mathcal{T} \rightarrow \mathbb{R}$ is the utility function of player i :

$$u_i(\Theta, \mathbf{t}) = -L_i(\Theta) - \lambda_i R_i(\theta_i) + \gamma_i \tilde{S}_i(\Theta, \mathbf{t}), \quad (18)$$

where L_i is player i 's contribution to the loss, R_i is a regularization term, \tilde{S}_i captures strategic interactions and self-allocation incentives, and $\Theta_{-i} = (\theta_j)_{j \neq i}$.

3.2. Multi-time nash equilibrium

Definition 3.2 (Multi-Time Nash Equilibrium). A strategy profile $\Theta^* \in \Theta$ and time allocation $\mathbf{t}^* \in \mathcal{T}$ constitute a Multi-Time Nash Equilibrium if:

$$u_i(\theta_i^*, \theta_{-i}^*, t_i^*, \mathbf{t}_{-i}^*) \geq u_i(\theta_i, \theta_{-i}^*, t_i, \mathbf{t}_{-i}^*), \quad \forall i \in \mathcal{N}, \forall (\theta_i, t_i) \in \Theta_i \times \mathcal{T}_i. \quad (19)$$

Remark 3.3 (Scope of Optimization). Each player i optimizes over $(\theta_i, t_i) \in \Theta_i \times \mathcal{T}_i$, that is, its own parameters and its own time allocation, while $(\theta_{-i}, \mathbf{t}_{-i})$ are held fixed. The variable t_i is a local decision variable for player i , not a global variable. Coupling arises only through u_i 's dependence on θ_{-i} and \mathbf{t}_{-i} .

Remark 3.4 (Connection to Markov Games). The multi-time learning game maps to a Markov game by identifying: state = Θ , action of player $i = \Delta\theta_i$, transition = $\Theta' = \Theta + \Delta\theta$, stage reward = u_i . The key difference is heterogeneous clocks t_i , which standard Markov games lack.

Remark 3.5 (Explicit dependence of u_i on t_i and closed-form best response). With the augmented strategic term (17), the utility $u_i(\Theta, \mathbf{t})$ depends explicitly on player i 's own time allocation t_i both through the self-allocation reward $\zeta_i t_i \|\nabla_{\theta_i} L\|^2$ and through the strictly convex effort cost $-\frac{1}{2} \nu_i t_i^2$. Differentiating u_i with respect to t_i at fixed $(\theta_{-i}, \mathbf{t}_{-i})$ and setting the derivative to zero yields

$$\frac{\partial u_i}{\partial t_i} = \gamma_i \zeta_i \|\nabla_{\theta_i} L(\Theta)\|^2 - \gamma_i \nu_i t_i = 0 \quad \implies \quad t_i^*(\Theta) = \frac{\zeta_i}{\nu_i} \|\nabla_{\theta_i} L(\Theta)\|^2. \quad (20)$$

Since $\partial^2 u_i / \partial t_i^2 = -\gamma_i \nu_i < 0$, the stationary point is a strict global maximum and (20) is the unique best response in t_i . The expression admits the natural economic reading that players whose own gradient is large – i.e., layers far from optimum and therefore standing to benefit most from additional learning time – claim more of their own time budget, while a strictly convex cost prevents any single player from monopolizing the time horizon. Because $\|\nabla_{\theta_i} L\|$ varies markedly across layers (early convolutional layers vs. late classification layers carry very different gradient magnitudes [2,29,31]), the equilibrium profile $\{t_i^*\}$ is genuinely heterogeneous, which is precisely the multi-time mechanism that motivates this work. This closed form is what Step 8 of Algorithm 1 computes.

Definition extends the classical Nash equilibrium by requiring stability with respect to both parameter deviations and time allocation strategies.

3.3. Existence theorem

Assumption 3.6. For each player $i \in \mathcal{N}$: (A1) The strategy space Θ_i is nonempty, compact, and convex; (A2) The utility function $u_i(\cdot, \theta_{-i}, \mathbf{t})$ is concave in θ_i for fixed $(\theta_{-i}, \mathbf{t})$; (A3) The utility function u_i is continuous in (Θ, \mathbf{t}) .

Assumption 3.7 (Strong Concavity). For each player i , the Hessian of u_i with respect to θ_i satisfies:

$$H_i(\theta) = \nabla_{\theta_i}^2 u_i(\theta, \theta_{-i}, \mathbf{t}) \leq -\mu_i I, \quad (21)$$

for some $\mu_i > 0$, where \leq denotes the Loewner order.

Assumption 3.8 (Consistency). The mixed partial derivatives of utility functions commute:

$$\frac{\partial^2 u_i}{\partial \theta_i \partial t_j} = \frac{\partial^2 u_i}{\partial t_j \partial \theta_i}, \quad \forall i, j \in \mathcal{N}. \quad (22)$$

Theorem 3.9 (Existence of Multi-Time Nash Equilibrium). Under Assumptions 3.6 to 3.8, there exists at least one Multi-Time Nash Equilibrium (Θ^*, \mathbf{t}^*) .

Proof. We apply Kakutani's fixed-point theorem [18]. Define the best-response correspondence $B : \Theta \times \mathcal{T} \rightrightarrows \Theta \times \mathcal{T}$ by:

$$B(\Theta, \mathbf{t}) = \prod_{i=1}^N B_i(\Theta_{-i}, \mathbf{t}), \quad (23)$$

where:

$$B_i(\Theta_{-i}, \mathbf{t}) = \operatorname{argmax}_{(\theta_i, t_i) \in \Theta_i \times \mathcal{T}_i} u_i(\theta_i, \Theta_{-i}, t_i, \mathbf{t}_{-i}). \quad (24)$$

We verify the conditions of Kakutani's theorem:

Step 1: Nonemptiness. By (A1), $\Theta_i \times \mathcal{T}_i$ is compact. By (A3), u_i is continuous. Hence the maximum is attained, and $B_i(\Theta_{-i}, \mathbf{t}) \neq \emptyset$.

Step 2: Convexity. Suppose $(\theta_i', t_i'), (\theta_i'', t_i'') \in B_i(\Theta_{-i}, \mathbf{t})$. For $\lambda \in [0, 1]$, let $(\theta_i^\lambda, t_i^\lambda) = \lambda(\theta_i', t_i') + (1 - \lambda)(\theta_i'', t_i'')$. By (A2) and convexity of \mathcal{T}_i :

$$u_i(\theta_i^\lambda, \Theta_{-i}, t_i^\lambda, \mathbf{t}_{-i}) \geq \lambda u_i(\theta_i', \Theta_{-i}, t_i', \mathbf{t}_{-i}) + (1 - \lambda) u_i(\theta_i'', \Theta_{-i}, t_i'', \mathbf{t}_{-i}) \quad (25)$$

$$= u_i(\theta_i', \Theta_{-i}, t_i', \mathbf{t}_{-i}), \quad (26)$$

where the equality follows because both points are maxima. Hence $(\theta_i^\lambda, t_i^\lambda) \in B_i(\Theta_{-i}, \mathbf{t})$, proving convexity.

Step 3: Closed Graph. Let $(\Theta^n, \mathbf{t}^n) \rightarrow (\Theta, \mathbf{t})$ and $(\hat{\Theta}^n, \hat{\mathbf{t}}^n) \in B(\Theta^n, \mathbf{t}^n)$ with $(\hat{\Theta}^n, \hat{\mathbf{t}}^n) \rightarrow (\hat{\Theta}, \hat{\mathbf{t}})$. For any $(\theta_i, t_i) \in \Theta_i \times \mathcal{T}_i$:

$$u_i(\hat{\theta}_i^n, \Theta_{-i}^n, \hat{t}_i^n, \mathbf{t}_{-i}^n) \geq u_i(\theta_i, \Theta_{-i}^n, t_i, \mathbf{t}_{-i}^n). \quad (27)$$

Taking limits and using continuity (A3):

$$u_i(\hat{\theta}_i, \Theta_{-i}, \hat{t}_i, \mathbf{t}_{-i}) \geq u_i(\theta_i, \Theta_{-i}, t_i, \mathbf{t}_{-i}). \quad (28)$$

Hence $(\hat{\Theta}, \hat{\mathbf{t}}) \in B(\Theta, \mathbf{t})$, proving the graph is closed. Since $K = \prod_i (\Theta_i \times \mathcal{T}_i)$ is compact and B maps K into K , the closed graph property implies upper hemicontinuity [39]. We further note that under the augmented utility (18)–(17), u_i is now strictly concave in t_i (because $\partial^2 u_i / \partial t_i^2 = -\gamma_i \nu_i < 0$), so the best response in t_i is in fact single-valued and given by the closed form (20); this strengthens Kakutani's hypotheses without affecting the conclusion.

By Kakutani's theorem, B has a fixed point (Θ^*, \mathbf{t}^*) , which by construction is a Multi-Time Nash Equilibrium. \square

Corollary 3.10 (Uniqueness under Strong Concavity). Under Assumption 3.7, if additionally the interaction matrix $B = (\beta_{ij})_{i,j}$ satisfies $\rho(B) < \min_i \mu_i / \max_i \lambda_{\max}(H_i)$, then the equilibrium is unique.

4. Optimal control perspective

4.1. Non-holonomic control problem

We formulate the optimization as a non-holonomic optimal control problem.

Definition 4.1 (Multi-Time Optimal Control Problem). Find controls $u = (u_{ij})_{i,j \in \mathcal{N}}$ minimizing:

$$J[u] = \int_0^T L(\Theta(t)) dt + \Phi(\Theta(T)), \quad (29)$$

subject to the controlled dynamics:

$$\frac{d\theta_i}{dt} = \sum_{j=1}^N u_{ij}(t) X_{ij}(\Theta), \quad i \in \mathcal{N}, \quad (30)$$

with control constraints:

$$\sum_{j=1}^N u_{ij}^2(t) \leq 1, \quad \forall i \in \mathcal{N}, \forall t \in [0, T]. \quad (31)$$

4.2. Multi-time maximum principle

Theorem 4.2 (Multi-Time Pontryagin Maximum Principle). *Let (Θ^*, u^*) be an optimal pair for Problem 4.1. Then there exist adjoint variables $\lambda = (\lambda_1, \dots, \lambda_N)$ satisfying:*

$$\frac{d\lambda_i}{dt} = -\frac{\partial H}{\partial \theta_i}(\lambda, \Theta^*, u^*), \quad (32)$$

where the Hamiltonian is:

$$H(\lambda, \Theta, u) = L(\Theta) + \sum_{i=1}^N \sum_{j=1}^N \lambda_i^\top u_{ij} X_{ij}(\Theta), \quad (33)$$

and the optimal control satisfies:

$$\nabla_u H(\lambda, \Theta^*, u^*) = 0. \quad (34)$$

Proof. The proof follows from the classical Pontryagin Maximum Principle [19] extended to the multi-input setting. The Hamiltonian structure is preserved under the summation over control indices j , and the necessary conditions follow from standard variational arguments. \square

Proposition 4.3 (Bang-Bang Structure). *Under the control constraint (31), the optimal control exhibits a bang-bang structure:*

$$u_{ij}^*(t) = \frac{\lambda_i^\top X_{ij}(\Theta^*)}{\sqrt{\sum_k (\lambda_i^\top X_{ik}(\Theta^*))^2}}, \quad (35)$$

allocating full control authority in the direction maximizing the Hamiltonian.

4.3. Convergence analysis

Assumption 4.4 (Lipschitz Gradients). The loss function L has ℓ -Lipschitz gradients:

$$\|\nabla L(\theta) - \nabla L(\theta')\| \leq \ell \|\theta - \theta'\|, \quad \forall \theta, \theta' \in \Theta. \quad (36)$$

Assumption 4.5 (Bounded Interactions). The interaction matrix $B = (\beta_{ij})$ has spectral radius:

$$\rho(B) = \max\{|\lambda| : \lambda \text{ is an eigenvalue of } B\} < \infty. \quad (37)$$

Theorem 4.6 (Exponential Convergence). *Under Assumptions 3.7, 4.4, and 4.5, if the learning rates satisfy:*

$$\alpha_i < \frac{2\mu_i}{\lambda_{\max}(H_i)}, \quad \forall i \in \mathcal{N}, \quad (38)$$

then the multi-time dynamics (12) converge exponentially:

$$\|\Theta(t) - \Theta^*\| \leq C e^{-\gamma t} \|\Theta(0) - \Theta^*\|, \quad (39)$$

where the convergence rate is:

$$\gamma = \min_i (\alpha_i \mu_i) - \rho(B) \max_i \alpha_i \lambda_{\max}(H_i). \quad (40)$$

Proof. Define the Lyapunov function:

$$V(\Theta) = \frac{1}{2} \|\Theta - \Theta^*\|^2 = \frac{1}{2} \sum_{i=1}^N \|\theta_i - \theta_i^*\|^2. \quad (41)$$

Computing the time derivative along trajectories of (12):

$$\frac{dV}{dt} = \sum_{i=1}^N (\theta_i - \theta_i^*)^\top \frac{d\theta_i}{dt} \quad (42)$$

$$= \sum_{i=1}^N (\theta_i - \theta_i^*)^\top \left(\alpha_i \nabla_{\theta_i} L + \sum_{j \neq i} \beta_{ij} H_{ij} \frac{d\theta_j}{dt} \right). \quad (43)$$

By strong concavity (Assumption 3.7):

$$(\theta_i - \theta_i^*)^\top \nabla_{\theta_i} L(\Theta) \leq -\mu_i \|\theta_i - \theta_i^*\|^2 + (\theta_i - \theta_i^*)^\top \nabla_{\theta_i} L(\Theta^*). \quad (44)$$

At the equilibrium, $\nabla_{\theta_i} L(\Theta^*) = 0$, so:

$$(\theta_i - \theta_i^*)^\top \alpha_i \nabla_{\theta_i} L \leq -\alpha_i \mu_i \|\theta_i - \theta_i^*\|^2. \quad (45)$$

For the interaction terms, we proceed step by step. Applying the Cauchy-Schwarz inequality to each individual term $|(\theta_i - \theta_i^*)^\top \beta_{ij} H_{ij} (d\theta_j/dt)| \leq \beta_{ij} \|\theta_i - \theta_i^*\| \cdot \|H_{ij}\| \cdot \|d\theta_j/dt\|$, substituting the leading-order dynamics $\|d\theta_j/dt\| \leq \alpha_j \|\nabla_{\theta_j} L\|$, using the Lipschitz condition $\|\nabla_{\theta_j} L\| \leq \ell \|\Theta - \Theta^*\|$ (since $\nabla_{\theta_j} L(\Theta^*) = 0$), applying the Gershgorin-type bound $\sum_{j \neq i} \beta_{ij} \leq \rho(B)$, and Young's inequality to handle cross terms, we obtain:

$$\left| \sum_i \sum_{j \neq i} (\theta_i - \theta_i^*)^\top \beta_{ij} H_{ij} \frac{d\theta_j}{dt} \right| \leq \rho(B) \max_i \alpha_i \lambda_{\max}(H_i) \sum_i \|\theta_i - \theta_i^*\|^2. \quad (46)$$

Combining (45) and (46):

$$\frac{dV}{dt} \leq -\sum_i \alpha_i \mu_i \|\theta_i - \theta_i^*\|^2 + \rho(B) \max_i \alpha_i \lambda_{\max}(H_i) \sum_i \|\theta_i - \theta_i^*\|^2 \quad (47)$$

$$\leq -\left(\min_i (\alpha_i \mu_i) - \rho(B) \max_i \alpha_i \lambda_{\max}(H_i) \right) \cdot 2V \quad (48)$$

$$= -2\gamma V. \quad (49)$$

By Grönwall's inequality:

$$V(t) \leq V(0) e^{-2\gamma t}, \quad (50)$$

which yields (39) with $C = 1$. \square

Remark 4.7 (Interpretation of Convergence Rate). The convergence rate γ in (40) consists of two competing terms: $\min_i (\alpha_i \mu_i)$ represents the driving force toward equilibrium, limited by the slowest player, while $\rho(B) \max_i \alpha_i \lambda_{\max}(H_i)$ captures the resistance from inter-player interactions. Convergence requires $\gamma > 0$, i.e., the driving force must exceed the interaction-induced resistance.

5. Algorithm

5.1. Multi-time nash learning (MTNL)

Based on the theoretical framework, we develop the Multi-Time Nash Learning algorithm.

5.2. Discretization: from continuous to discrete dynamics

The continuous-time dynamics (12) are discretized via explicit Euler with step size $\eta > 0$:

$$\theta_i^{(k+1)} = \theta_i^{(k)} + \eta \left[\alpha_i \nabla_{\theta_i} L(\Theta^{(k)}) + \sum_{j \neq i} \beta_{ij} H_{ij}^{(k)} \alpha_j \nabla_{\theta_j} L(\Theta^{(k)}) \right], \quad (51)$$

which corresponds to Steps 4 to 6 of Algorithm 1. The one-step error satisfies $\mathcal{O}(\eta^2)$ and the global error over $K = T/\eta$ steps is $\mathcal{O}(\eta)$, ensuring convergence as $\eta \rightarrow 0$. The continuous convergence rate γ translates to a discrete contraction factor $1 - \eta\gamma + \mathcal{O}(\eta^2)$.

Remark 5.1 (Derivation of H_{ij} from the Game). The mixed Hessian $H_{ij} = \partial^2 L / \partial \theta_i \partial \theta_j$ arises from the first-order Taylor expansion: $\nabla_{\theta_i} L(\theta_i, \theta_j + \Delta \theta_j) \approx \nabla_{\theta_i} L(\theta_i, \theta_j) + H_{ij} \Delta \theta_j$. The interaction term $c_i^{(k)}$ thus implements first-order anticipation of other players' moves. In practice, $H_{ij} g_j$ is computed via one additional backward pass.

The previous remark established that the global discretization error is $\mathcal{O}(\eta)$, but it did not state under which assumptions the discrete iterates inherit the *exponential stability* proved in [Theorem 4.6](#) for the continuous flow. We now make this inheritance explicit.

Proposition 5.2 (Discrete Inheritance of Exponential Stability). *Let the continuous multi-time dynamics (12) satisfy the assumptions of [Theorem 4.6](#), with continuous convergence rate*

$$\gamma = \min_i (\alpha_i \mu_i) - \rho(B) \max_i \alpha_i \lambda_{\max}(H_i) > 0.$$

Assume in addition that L is twice continuously differentiable with $\|\nabla^2 L\| \leq M$ on a neighbourhood of Θ^* . If the step size η of the explicit Euler scheme (51) satisfies

$$0 < \eta < \min \left\{ \frac{\gamma}{M^2(1 + \rho(B))^2 \max_i \alpha_i^2}, \frac{2}{\ell + \rho(B) \max_i \alpha_i \ell} \right\}, \quad (52)$$

then the discrete iterates satisfy

$$\begin{aligned} \|\Theta^{(k+1)} - \Theta^*\| &\leq \kappa(\eta) \|\Theta^{(k)} - \Theta^*\|, \\ \kappa(\eta) &= 1 - \eta\gamma + \frac{1}{2}\eta^2 M^2(1 + \rho(B))^2 \max_i \alpha_i^2 < 1, \end{aligned} \quad (53)$$

for all $k \geq 0$. In particular, the discrete scheme inherits exponential stability with rate at least $-\log \kappa(\eta) \geq \eta\gamma/2$.

Proof. Define the discrete Lyapunov function $V^{(k)} = \frac{1}{2} \|\Theta^{(k)} - \Theta^*\|^2$ and let $F(\Theta) = (\alpha_i \nabla_{\theta_i} L + \sum_{j \neq i} \beta_{ij} H_{ij} \alpha_j \nabla_{\theta_j} L)_{i \in \mathcal{N}}$ denote the right-hand side of (12), so that the explicit Euler step reads $\Theta^{(k+1)} = \Theta^{(k)} + \eta F(\Theta^{(k)})$. Expanding,

$$V^{(k+1)} = V^{(k)} + \eta \langle F(\Theta^{(k)}), \Theta^{(k)} - \Theta^* \rangle_{\text{angle}} + \frac{\eta^2}{2} \|F(\Theta^{(k)})\|^2.$$

The linear term is bounded by $-2\gamma V^{(k)}$ as a direct consequence of the Lyapunov decrement (49) established in the proof of [Theorem 4.6](#). The quadratic term is bounded using $\|\nabla L(\Theta)\| \leq M \|\Theta - \Theta^*\|$ (Lipschitz gradient near Θ^*) and the triangle inequality, yielding $\|F(\Theta)\| \leq M(1 + \rho(B)) \max_i \alpha_i \|\Theta - \Theta^*\|$. Combining the two bounds,

$$V^{(k+1)} \leq [1 - 2\eta\gamma + \eta^2 M^2(1 + \rho(B))^2 \max_i \alpha_i^2] V^{(k)} = \kappa(\eta)^2 V^{(k)},$$

where the right-hand side of (52) ensures $\kappa(\eta) < 1$. The first bound in (52) guarantees that the destabilising quadratic term is dominated by the stabilising linear term; the second bound is the Lipschitz step bound that already appeared in the proof of (38). Taking square roots gives (53). \square

[Proposition 5.2](#) formalises the continuous-to-discrete bridge requested by the literature on neural ODEs and structure-preserving integrators [14,37]: the same hypotheses used in [Theorem 4.6](#), supplemented by the explicit step-size bound (52), suffice for [Algorithm 1](#) to enjoy a discrete contraction factor that converges to the continuous rate as $\eta \rightarrow 0$. [Section 6.1](#) reports a numerical sanity check verifying that the values of $(\eta, \alpha_i, \beta_{ij})$ used throughout the experiments satisfy (52) with margin of order $10\times$.

Closed-form time allocation. Building on [Remark 3.5](#), Step 8 of [Algorithm 1](#) replaces the formal argmax_i by the explicit solution (20), which is essentially computationally free because $\|\nabla_{\theta_i} L\|^2$ has already been computed in Step 4. Sigmoid clipping $t_i \leftarrow \sigma(\zeta_i \|\nabla_{\theta_i} L\|^2 / v_i)$ keeps $t_i \in (0, 1)$, in agreement with Appendix B. The previous ablation finding that time allocation contributes a more modest 0.6 percentage points ([Section 6.5](#)) is unaffected, since the gradient-norm-based schedule used in our experiments was operationally equivalent to (20); the present subsection simply gives the formal derivation.

5.3. Computational complexity

Proposition 5.3 (Complexity Analysis). *The per-epoch computational complexity of MTNL is:*

$$\mathcal{O}(N^2 d \cdot T), \quad (54)$$

where N is the number of network components, d is the total parameter dimension, and T is the number of training steps per epoch.

Proof. The dominant cost is computing the interaction terms $c_i^{(k)} = \sum_{j \neq i} \beta_{ij} H_{ij} g_j^{(k)}$. For each of N players, we compute interactions with $N - 1$ others, each involving a matrix-vector product of dimension $\mathcal{O}(d_i \times d_j)$. Summing over all pairs and training steps yields the stated complexity. \square

Remark 5.4 (Practical Approximations). In practice, we employ several approximations to reduce computational cost. The diagonal Hessian approximation $H_{ij} \approx \operatorname{diag}(\operatorname{diag}(H_{ij}))$ reduces complexity to $\mathcal{O}(Nd)$. Setting $\beta_{ij} = 0$ for non-adjacent layers exploits the network's sequential topology to further reduce computation. Computing interaction terms every K iterations rather than every iteration provides additional speedup with minimal impact on convergence.

5.4. Relation to existing optimizers

MTNL generalizes several existing optimization methods.

Proposition 5.5 (Special Cases). *Setting $\beta_{ij} = 0$ for all i, j and uniform $\alpha_i = \alpha$ recovers stochastic gradient descent. With $H_{ij} = F_{ij}$ (Fisher information blocks) and $\beta_{ij} = -1$, MTNL approximates natural gradient descent [12]. Non-uniform α_i with $\beta_{ij} = 0$ corresponds to layer-wise learning rate adaptation [20]. Thus, layer-wise adaptive methods such as LARS [20] and discriminative fine-tuning [36] are special cases of the multi-time framework that capture heterogeneous learning rates but not the strategic inter-layer coordination that the β_{ij} terms provide. The connection to neural ordinary differential equations [37], where residual networks are viewed as discretized dynamical systems, suggests that the multi-time structure arises naturally from block-wise architectures.*

6. Experiments

6.1. Experimental setup

We evaluate the Multi-Time Nash Learning algorithm on three standard image classification benchmarks spanning different scales and complexities. CIFAR-10 [21] consists of 60,000 color images of size 32×32 pixels distributed across 10 classes, providing a well-established benchmark for comparing optimization methods. Fashion-MNIST [22] contains 70,000 grayscale images of size 28×28 pixels representing clothing items in 10 categories, offering a more challenging alternative to the original MNIST dataset. For large-scale evaluation, we use a subset of ImageNet [23] comprising 100,000 images from 100 randomly selected classes, allowing us to assess scalability to realistic problem sizes.

Our experiments employ two widely-used convolutional architectures. ResNet-50 [24] is a 50-layer residual network with approximately 25 million parameters, representing the state-of-the-art in deep network design through its skip connections that facilitate gradient flow. VGG-16 [25] is a 16-layer convolutional network with approximately 138 million parameters, characterized by its uniform architecture using small 3×3 convolutional filters throughout. We additionally employ a custom 5-layer fully connected network for controlled ablation studies where the simpler architecture allows clearer isolation of individual algorithmic components.

We compare MTNL against four baseline optimization methods. Stochastic Gradient Descent (SGD) [26] with momentum 0.9 serves as the fundamental baseline representing the simplest widely-used optimizer. Adam [27] with its default parameters ($\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$) represents adaptive learning rate methods that maintain per-parameter moment estimates. AdaGrad [28] provides another adaptive

Table 1

Complete hyperparameter settings. All experiments: batch size 128, weight decay 10^{-4} , cosine learning rate schedule with 5-epoch warmup, standard augmentation, 200 epochs, 5 random seeds.

Optimizer	Key Hyperparameters	Ref.
SGD	lr= 0.1, momentum= 0.9	[30]
Adam	lr= 0.001, $\beta_1=0.9$, $\beta_2=0.999$	[31]
AdaGrad	lr= 0.01	[32]
AdamW	lr= 0.001, $\beta_1=0.9$, $\beta_2=0.999$, $\lambda=0.01$	[32]
Lion	lr= 3×10^{-4} , $\beta_1=0.9$, $\beta_2=0.99$, $\lambda=0.1$	[33]
LARS	lr= 0.1, momentum= 0.9, trust= 0.001	[23]
Nat. Grad.	lr= 0.01, K-FAC every 10 steps	[16]
MTNL	$\alpha=0.01$, $\beta_{ij}=0.1$ (adj.), diag. Fisher, $K_i=10$	n/a

baseline with parameter-wise learning rate scaling based on accumulated squared gradients. Finally, Natural Gradient using the K-FAC approximation [13] represents geometry-aware methods that exploit the Fisher information structure, offering the closest methodological comparison to our approach.

For MTNL, we set the base learning rate to $\alpha = 0.01$ with architecture-specific tuning, coupling coefficients $\beta_{ij} = 0.1$ for adjacent layers and zero otherwise (exploiting the observation from ablation studies that non-adjacent interactions contribute minimally), diagonal Fisher information approximation for the Hessian, and time allocation updates performed every 10 epochs. All experiments use batch size 128, weight decay 10^{-4} , and run for 200 epochs with results averaged over 5 random seeds to ensure statistical reliability.

6.2. Baselines and hyperparameter details

We compare MTNL against seven baselines with full hyperparameter specifications (Table 1):

Learning rates were selected via grid search over $\{0.1, 0.03, 0.01, 0.003, 0.001, 0.0003\}$; coupling $\beta \in \{0.01, 0.05, 0.1, 0.2, 0.5\}$ was selected on 10% validation split. Statistical significance: Welch's two-sided t -test confirms $p < 0.01$ for all MTNL vs. baseline comparisons.

Numerical sanity check for Proposition 5.2. For the configuration used in the CIFAR-10/ResNet-50 experiments (effective step size $\eta = 10^{-2}$, $\max_i \alpha_i = 10^{-2}$, adjacent-only coupling so $\rho(B) \leq 0.2$, and an empirical Hessian operator-norm bound $M \leq 50$ measured on a 1% validation subsample of training trajectories), the right-hand side of (52) evaluates to $\min\{\gamma/(M^2(1 + \rho(B))^2 \max_i \alpha_i^2), 2/(\ell(1 + \rho(B) \max_i \alpha_i))\} \approx 0.13$. The actual step $\eta = 10^{-2}$ is therefore comfortably below the bound by a factor of ~ 13 , so the discrete contraction $\kappa(\eta) < 1$ established in Proposition 5.2 holds throughout training. Identical checks for the Fashion-MNIST and ImageNet-100 experiments yield safety margins of $\sim 18\times$ and $\sim 9\times$ respectively, confirming that the discrete iterates of Algorithm 1 inherit the exponential stability proved for the continuous flow in Theorem 4.6.

6.3. Main results

Table 2 presents the main experimental results on CIFAR-10 with ResNet-50. MTNL achieves a test accuracy of 94.1%, representing a 3.7 percentage point improvement over SGD (90.4%) and a 2.9 percentage point improvement over Adam (91.2%). The convergence speed shows even more dramatic gains: MTNL reaches convergence in 76 epochs compared to 108 for SGD, representing approximately 40% faster convergence. The stability metric measures the coefficient of variation of loss over the last 10 epochs: $\text{Stability} = 1 - \text{Var}(\text{loss})/\text{Mean}(\text{loss})$, where higher values indicate smoother convergence. We note that stability does not directly imply better generalization; the correlation between stability and test accuracy across all methods is $r = 0.89$ ($p < 0.01$), indicating a positive but imperfect association. The stability metric, which

Table 2

Test accuracy (%) and convergence epochs on CIFAR-10 with ResNet-50. Best results in **bold**.

Method	Test Acc.	Conv. Epochs	Stability	Time/Epoch
SGD	90.4 \pm 0.3	108	0.78	1.0x
Adam	91.2 \pm 0.2	95	0.82	1.1x
AdaGrad	89.8 \pm 0.4	112	0.75	1.1x
Natural Gradient	92.1 \pm 0.3	85	0.85	2.3x
AdamW	91.5 \pm 0.2	92	0.84	1.1x
Lion	91.8 \pm 0.3	88	0.86	1.0x
LARS	91.0 \pm 0.3	96	0.83	1.0x
MTNL (Ours)	94.1 \pm 0.2	76	0.94	1.4x
w/o Multi-time	92.0 \pm 0.3	89	0.87	1.2x
w/o Strategic Int.	92.8 \pm 0.2	82	0.89	1.3x

Table 3

Test accuracy (%) across datasets and architectures.

Dataset	Architecture	SGD	Adam	MTNL
CIFAR-10	ResNet-50	90.4	91.2	94.1
CIFAR-10	VGG-16	89.2	90.1	93.5
Fashion-MNIST	ResNet-50	93.1	93.8	95.2
Fashion-MNIST	VGG-16	92.5	93.2	94.8
ImageNet-100	ResNet-50	72.3	74.1	78.6

measures the consistency of the loss trajectory in later training stages, reaches 0.94 for MTNL compared to values ranging from 0.75 to 0.85 for the baselines, indicating that the multi-time coordination produces smoother optimization dynamics with fewer oscillations.

The ablation variants in Table 2 reveal the relative contributions of the framework's components. Removing the multi-time structure while retaining strategic interactions yields 92.0% accuracy, indicating that the multi-time formulation alone contributes approximately 2.1 percentage points. Conversely, removing strategic interactions while retaining the multi-time structure achieves 92.8% accuracy, suggesting that the interaction terms contribute approximately 1.3 percentage points. Notably, the combined effect of 3.7 percentage points exceeds the sum of individual contributions, demonstrating a synergistic interaction between the two components.

Table 3 demonstrates that the improvements generalize consistently across different datasets and architectures. On CIFAR-10 with VGG-16, MTNL achieves 93.5% compared to 90.1% for Adam, a gain of 3.4 percentage points. Fashion-MNIST shows improvements of 1.4 to 1.6 percentage points over Adam across both architectures. Most notably, on the more challenging ImageNet-100 subset, MTNL achieves 78.6% accuracy compared to 74.1% for Adam, a substantial gain of 4.5 percentage points that suggests the benefits of multi-time coordination become more pronounced as task complexity increases.

6.4. Convergence dynamics

Fig. 1 visualizes the convergence dynamics, confirming the theoretical prediction of faster convergence for MTNL. The exponential decay predicted by Theorem 4.6 is evident in the log-linear behavior of the MTNL curve.

6.5. Ablation studies

Table 4 presents a systematic ablation study isolating the contribution of each algorithmic component. The interaction terms emerge as the most critical component: setting all coupling coefficients $\beta_{ij} = 0$ reduces accuracy by 2.1 percentage points and increases convergence time from 76 to 89 epochs. This confirms that the strategic coordination enabled by the mixed Hessian terms H_{ij} captures essential information about inter-layer dependencies that standard optimizers ignore.

Table 4
Ablation study on MTNL components.

Configuration	Test Acc.	Conv. Epochs	Δ vs Full
Full MTNL	94.1	76	n/a
$\beta_{ij} = 0$ (no interactions)	92.0	89	-2.1%
Uniform α_i	93.2	81	-0.9%
No time optimization	93.5	79	-0.6%
Diagonal H_{ij} only	93.8	78	-0.3%
Adjacent layers only	93.9	77	-0.2%

Layer-wise learning rates provide the second most significant contribution. Using uniform learning rates $\alpha_i = \alpha$ across all layers reduces accuracy by 0.9 percentage points, validating the intuition that different network components benefit from different optimization speeds. The time allocation optimization contributes a more modest 0.6 percentage points; while the dynamic adjustment of time allocation improves performance, the static allocation based on network architecture already captures much of the benefit.

Two additional ablations inform practical implementation choices. Using only the diagonal of the mixed Hessian rather than the full matrix incurs a minimal accuracy loss of 0.3 percentage points while substantially reducing computational cost from $\mathcal{O}(d_i d_j)$ to $\mathcal{O}(\min(d_i, d_j))$ per layer pair. Restricting interactions to adjacent layers only (setting $\beta_{ij} = 0$ for $|i - j| > 1$) reduces accuracy by just 0.2 percentage points, suggesting that the network’s sequential structure concentrates most inter-layer dependencies between neighboring components. Together, these findings justify the practical approximations described in Remark 5.4 that enable efficient implementation without sacrificing meaningful performance.

6.6. Computational overhead

Table 2 shows that MTNL incurs 1.4× computational overhead per epoch compared to SGD. Given the 40% reduction in convergence epochs, the total training time is:

$$\text{Total time ratio} = \frac{1.4 \times 76}{1.0 \times 108} \approx 0.99, \quad (55)$$

meaning MTNL achieves better results in approximately the same total time as SGD.

6.7. Illustrative example: multi-time dynamics in a 3-layer network

To illustrate the specific dynamical behavior, consider a 3-layer fully connected network with $\alpha_1 = 0.01$, $\alpha_2 = 0.008$, $\alpha_3 = 0.015$, $\beta_{12} = \beta_{23} = 0.1$, $\beta_{13} = 0$. On Fashion-MNIST, the multi-time dynamics produce:

- **Phase 1 (epochs 1 to 20):** The output layer converges fastest, pulling the hidden layer through $\beta_{23} H_{23} g_3$.
- **Phase 2 (epochs 20 to 60):** The hidden layer becomes the bottleneck; interaction terms from both layers coordinate its updates.
- **Phase 3 (epochs 60 +):** All layers approach equilibrium with time allocations stabilizing at $t_1 \approx 0.3$, $t_2 \approx 0.5$, $t_3 \approx 0.2$.

This phased behavior, emerging naturally from the multi-time structure, is not captured by single-time optimizers.

7. Conclusions

We have presented a novel mathematical framework for neural network optimization based on multi-time dynamics, unifying game-theoretic and optimal control perspectives. The framework provides rigorous foundations including existence theorems for Multi-Time Nash Equilibria and exponential convergence guarantees with explicit rates. The resulting MTNL algorithm achieves substantial improvements in convergence speed and final accuracy across standard benchmarks, with gains becoming more pronounced on more complex tasks.

The multi-time perspective provides a principled mathematical language for a phenomenon well-known to practitioners: different parts of a neural network learn at different speeds, and coordinating these heterogeneous dynamics is essential for efficient training. By formalizing this intuition through differential geometry and game theory, we obtain both theoretical insights and practical algorithmic tools.

This work opens several directions for future research. On the theoretical side, incorporating stochastic noise into the multi-time framework would yield stochastic multi-time games with direct relevance to mini-batch training. Relaxing convexity assumptions through techniques from non-convex optimization theory [8] would broaden the applicability of the convergence guarantees. From a neurocomputing perspective, the framework naturally extends to federated learning scenarios where heterogeneous devices with different computational capabilities correspond to distinct time scales. The game-theoretic structure is also directly applicable to multi-agent reinforcement learning [11], neural architecture search (where architecture and weight optimization proceed at different rates), and continual learning (where new task adaptation must be balanced against retention of prior knowledge). Distributed training with heterogeneous hardware, where different compute nodes process data at different speeds, represents another natural application domain where multi-time dynamics can provide principled scheduling policies.

The multi-time perspective reveals deep mathematical structure in neural network optimization that standard single-time approaches obscure. From a practical perspective, the framework naturally extends to human-computer interaction systems where human cognitive processes and machine computation operate on fundamentally different time scales. Such systems can be formulated as cooperative multi-time games, where the human “player” evolves slowly (deliberation) while the machine “player” executes rapidly (inference). Recent work on safe reinforcement learning from human demonstration [40], hierarchical control with prescribed performance [41], and fixed-time stochastic learning from human-UAV interaction [42] provides concrete examples where this multi-scale temporal structure is critical. Federated learning with heterogeneous devices, multi-agent reinforcement learning, neural architecture search, and continual learning represent additional natural application domains. I hope this framework inspires further theoretical and practical advances at the intersection of optimization, game theory, and deep learning.

Declaration of generative AI and AI-assisted technologies in the writing process

During the preparation of this work the author used AI-assisted tools for code verification and manuscript formatting. The author reviewed and edited all content and takes full responsibility for the content of the publication.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

The author acknowledges the support of Università Mediterranea di Reggio Calabria through the Decisions LAB research center. The author is grateful to ICRIOS Bocconi University for the stimulating research environment and to Istanbul Okan University for the visiting professorship support.

Appendix A. Proofs of auxiliary results

Proof. Proof of Proposition 2.4 Under the Fisher metric (6), the geodesic equations are:

$$\ddot{\gamma}^k + \Gamma_{ij}^k \dot{\gamma}^i \dot{\gamma}^j = 0, \quad (\text{A.1})$$

where Γ_{ij}^k are the Christoffel symbols:

$$\Gamma_{ij}^k = \frac{1}{2} g^{k\ell} \left(\frac{\partial g_{\ell i}}{\partial \theta_j} + \frac{\partial g_{\ell j}}{\partial \theta_i} - \frac{\partial g_{ij}}{\partial \theta_\ell} \right). \quad (\text{A.2})$$

The natural gradient flow $\dot{\theta} = -g^{-1} \nabla L$ satisfies these equations when L is the squared distance to a target, establishing the connection to geodesics. \square

Proof (Proof of Proposition 4.3). The Hamiltonian maximization condition (34) subject to the constraint $\sum_j u_{ij}^2 \leq 1$ gives the Lagrangian:

$$\mathcal{L} = H + \nu \left(1 - \sum_j u_{ij}^2 \right). \quad (\text{A.3})$$

First-order conditions yield:

$$\frac{\partial H}{\partial u_{ij}} = 2\nu u_{ij} \implies u_{ij} = \frac{\lambda_i^\top X_{ij}}{2\nu}. \quad (\text{A.4})$$

Substituting into the constraint and solving for ν gives (35). \square

Appendix B. Implementation details

We compute Hessian-vector products $H_{ij} g_j$ efficiently using automatic differentiation:

$$H_{ij} g_j = \frac{\partial}{\partial \theta_i} \left(\nabla_{\theta_j} L \cdot g_j \right), \quad (\text{B.1})$$

requiring only one additional backward pass per iteration.

Time allocation is parameterized as $t_i = \sigma(\tau_i)$ where σ is the sigmoid function, ensuring $t_i \in (0, 1)$. The optimization is performed via gradient ascent on τ_i .

The stability metric in Table 2 is defined as:

$$\text{Stability} = 1 - \frac{\text{Var}(\text{loss over last 10 epochs})}{\text{Mean}(\text{loss over last 10 epochs})}. \quad (\text{B.2})$$

Data availability

The experimental code and data used in this study will be made available upon publication. The benchmark datasets are publicly available from their respective sources: CIFAR-10 (<https://www.cs.toronto.edu/~kriz/cifar.html>), Fashion-MNIST (<https://github.com/zalando-research/fashion-mnist>), and ImageNet (<https://www.image-net.org>).

References

- [1] C. Zhang, S. Bengio, M. Hardt, B. Recht, O. Vinyals, Understanding deep learning (still) requires rethinking generalization, *Commun. ACM* 64 (3) (2021) 107–115.
- [2] A.M. Saxe, J.L. McClelland, S. Ganguli, Exact solutions to the nonlinear dynamics of learning in deep linear neural networks, in: *International Conference on Learning Representations*, 2014.
- [3] S. Arora, N. Cohen, E. Hazan, On the optimization of deep networks: implicit acceleration by overparameterization, in: *International Conference on Machine Learning*, PMLR, 2018, pp. 244–253.
- [4] J. Yosinski, J. Clune, Y. Bengio, H. Lipson, How transferable are features in deep neural networks? *Adv. Neural Inf. Process. Syst.* 27 (2014).
- [5] C. Udriște, *Convex Functions and Optimization Methods on Riemannian Manifolds*, Springer, Dordrecht, 1994.
- [6] C. Udriște, *Geometric Dynamics*, Springer, Dordrecht, 2000.
- [7] S.S. Du, J.D. Lee, H. Li, L. Wang, X. Zhai, Gradient descent finds global minima of deep neural networks, in: *International Conference on Machine Learning*, PMLR, 2019, pp. 1675–1685.
- [8] Z. Allen-Zhu, Y. Li, Z. Song, A convergence theory for deep learning via over-parameterization, in: *International Conference on Machine Learning*, PMLR, 2019, pp. 242–252.
- [9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, *Adv. Neural Inf. Process. Syst.* 27 (2014).
- [10] L. Mescheder, A. Geiger, S. Nowozin, Which training methods for GANs do actually converge? in: *International Conference on Machine Learning*, PMLR, 2018, pp. 3481–3490.
- [11] K. Zhang, Z. Yang, T. Başar, Multi-agent reinforcement learning: a selective overview of theories and algorithms, *Handbook of Reinforcement Learning and Control* (2021) 321–384.
- [12] S.-I. Amari, Natural gradient works efficiently in learning, *Neural Comput.* 10 (2) (1998) 251–276.
- [13] J. Martens, New insights and perspectives on the natural gradient method, *J. Mach. Learn. Res.* 21 (146) (2020) 1–76.
- [14] E. Hairer, C. Lubich, G. Wanner, *Geometric Numerical Integration: Structure-Preserving Algorithms for Ordinary Differential Equations*, second ed, Springer, 2006.
- [15] P.-A. Absil, R. Mahony, R. Sepulchre, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, 2008.
- [16] C. Udriște, I. Tevy, Multitime optimal control with area integral costs on boundary, *Balkan Journal of Geometry and Its Applications* 15 (2) (2010) 138–154.
- [17] A. Pitea, C. Udriște, On a class of optimal control problems governed by multitime Hamilton-Jacobi PDES, *WSEAS Transactions on Mathematics* 11 (5) (2012) 401–410.
- [18] S. Kakutani, A generalization of brouwer's fixed point theorem, *Duke Math. J.* 8 (3) (1941) 457–459.
- [19] L.S. Pontryagin, V.G. Boltyanskii, R.V. Gamkrelidze, E.F. Mishchenko, *The Mathematical Theory of Optimal Processes*, Interscience Publishers, 1962.
- [20] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, K. Keutzer, C.-J. Hsieh, Large batch optimization for deep learning: training BERT in 76min, in: *International Conference on Learning Representations*, 2019.
- [21] A. Krizhevsky, G. Hinton, *Learning Multiple Layers of Features From Tiny Images*, Technical report, University of Toronto, 2009.
- [22] H. Xiao, K. Rasul, R. Vollgraf, Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms, *arXiv preprint arXiv:1708.07747*, 2017.
- [23] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., ImageNet large scale visual recognition challenge, *Int. J. Comput. Vis.* 115 (3) (2015) 211–252.
- [24] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [25] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556* (2015).
- [26] H. Robbins, S. Monro, A stochastic approximation method, *The Annals of Mathematical Statistics* 22 3 (1951) 400–407.
- [27] D.P. Kingma, J. Ba, Adam: a method for stochastic optimization, *arXiv preprint arXiv:1412.6980*, 2015.
- [28] J. Duchi, E. Hazan, Y. Singer, Adaptive subgradient methods for online learning and stochastic optimization, *J. Mach. Learn. Res.* 12 (2011) 2121–2159.
- [29] S. Kornblith, M. Norouzi, H. Lee, G. Hinton, Similarity of neural network representations revisited, in: *International Conference on Machine Learning*, PMLR, 2019, pp. 3519–3529.
- [30] M. Raghu, J. Gilmer, J. Yosinski, J. Sohl-Dickstein, SVCCA: singular vector canonical correlation analysis for deep learning dynamics and interpretability, *Adv. Neural Inf. Process. Syst.* 30 (2017).
- [31] V. Pappas, X.Y. Han, D.L. Donoho, Prevalence of neural collapse during the terminal phase of deep learning training, *Proc. Natl. Acad. Sci.* 117 (40) (2020) 24652–24663.
- [32] I. Loshchilov, F. Hutter, Decoupled weight decay regularization, in: *International Conference on Learning Representations*, 2019.
- [33] X. Chen, C. Liang, D. Huang, E. Real, K. Wang, Y. Liu, H. Pham, X. Dong, T. Luong, C.-J. Hsieh, Y. Lu, Q.V. Le, Symbolic discovery of optimization algorithms, *Adv. Neural Inf. Process. Syst.* 36 (2023).
- [34] V. Gupta, T. Koren, Y. Singer, Shampoo: preconditioned stochastic tensor optimization, in: *International Conference on Machine Learning*, PMLR, 2018, pp. 1842–1850.
- [35] N. Vyas, D. Morwani, R. Zhao, I. Shapira, D. Brandfonbrener, L. Janson, S. Kakade, SOAP: improving and stabilizing shampoo using Adam, in: *International Conference on Learning Representations*, 2025.
- [36] J. Howard, S. Ruder, Universal language model fine-tuning for text classification, in: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, 2018, pp. 328–339.
- [37] R.T.Q. Chen, Y. Rubanova, J. Bettencourt, D.K. Duvenaud, Neural ordinary differential equations, *Adv. Neural Inf. Process. Syst.* 31 (2018).
- [38] R. Montgomery, A tour of subriemannian geometries, their geodesics and applications, in: *Mathematical Surveys and Monographs*, vol. 91, American Mathematical Society, 2002.
- [39] C.D. Aliprantis, K.C. Border, *Infinite Dimensional Analysis: A Hitchhiker's Guide*, third ed, Springer, 2006.

- [40] J. Tan, S. Xue, Z. Guo, H. Li, X. Zheng, H. Cao, Adaptive hierarchical control of quadcopters via safe reinforcement learning from human demonstration, *Eng. Appl. Artif. Intell.* 163 (2026) 112650, <https://doi.org/10.1016/j.engappai.2025.112650>
- [41] J. Tan, S. Xue, H. Li, Z. Guo, H. Cao, B. Chen, Hierarchical safe reinforcement learning control for leader-follower systems with prescribed performance, *IEEE Trans. Autom. Sci. Eng.* 22 (2025) 19568–19581, <https://doi.org/10.1109/TASE.2025.3596912>
- [42] J. Tan, S. Xue, Q. Guan, Z. Guo, H. Cao, B. Chen, Fixed-time stochastic learning from human-UAV interaction with state-input constraints, *IEEE Trans. Ind. Electron.* 73 (2) (2026), <https://doi.org/10.1109/TIE.2025.3613457>
- [43] M. Ferrara, Multi-time evolution models in deep learning dynamics, *Journal of Indian Academy of Mathematics* 47 (2) (2025) 282–294.
- [44] M. Ferrara, Geometric-entropic optimization: integrating optimal transport with riemannian gradient methods for neural network training, *J. Optim. Theory Appl.* (2026), <https://doi.org/10.1007/s10957-026-02958-8>

Author biography



Massimiliano Ferrara is Full Professor of Mathematical Economics, Artificial Intelligence, Machine Learning, and Applied Game Theory at the Università Mediterranea di Reggio Calabria, where he leads the Decisions LAB. He holds additional affiliations at ICRIOS – Bocconi University (Milan) and Istanbul Okan University (Turkey). He was awarded a PhD Honoris Causa in Computer Science by Georgian National University SEU (Tbilisi, 2025) and holds the distinguished civic honour of Cavaliere della Repubblica Italiana. His research programme bridges rigorous mathematical foundations—including Riemannian geometry, variational methods, and multi-time dynamics—with cutting-edge applications in machine learning, quantum computing, financial economics, and decision theory. He has authored and co-authored more than 200 peer-reviewed publications in leading international journals and proceedings, with an h-index placing him among the most productive scholars in his fields.