Dynamic Game Theoretic Neural Optimizer

Guan-Horng Liu¹² Tianrong Chen³ Evangelos A. Theodorou¹²

Abstract

The connection between training deep neural networks (DNNs) and optimal control theory (OCT) has attracted considerable attention as a principled tool of algorithmic design. Despite few attempts being made, they have been limited to architectures where the layer propagation resembles a Markovian dynamical system. This casts doubts on their flexibility to modern networks that heavily rely on non-Markovian dependencies between layers (e.g. skip connections in residual networks). In this work, we propose a novel dynamic game perspective by viewing each layer as a player in a dynamic game characterized by the DNN itself. Through this lens, different classes of optimizers can be seen as matching different types of Nash equilibria, depending on the implicit information structure of each (p)layer. The resulting method, called Dynamic Game Theoretic Neural Optimizer (DGNOpt), not only generalizes OCTinspired optimizers to richer network class; it also motivates a new training principle by solving a multi-player cooperative game. DGNOpt shows convergence improvements over existing methods on image classification datasets with residual and inception networks. Our work marries strengths from both OCT and game theory, paving ways to new algorithmic opportunities from robust optimal control and bandit-based optimization.

1. Introduction

Attempts from different disciplines to provide a fundamental understanding of deep learning have advanced rapidly in recent years. Among those, interpretation of DNNs as discrete-time nonlinear dynamical systems has received tremendous focus. By viewing each layer as a distinct time step, it motivates principled analysis from numerical equations (Weinan,

Proceedings of the 38th International Conference on Machine Learning, PMLR 139, 2021. Copyright 2021 by the author(s).

Terminology Mapping Deep Neural Net Training Optimizer Computation Order Layer Module (L, θ) Objective & Parameters Terminology Mapping Multi-Player Dynamic Game Solving Nash Equilibria Stage Sequence Player Payoff & Actions x_T $\min_{\theta} L(x_T, \theta)$

Figure 1. Dynamic game perspective of generic DNN training process, where we treat layer modules as players in a dynamic game and solve for the related Nash equilibria (Best viewed in color).

2017; Lu et al., 2017) to physics (Greydanus et al., 2019). For instance, casting residual networks (He et al., 2016) as a discretization of ordinary differential equations enables fundamental reasoning on the loss landscape (Lu et al., 2020) and inspires new architectures with numerical stability or continuous limit (Chang et al., 2018; Chen et al., 2018).

This dynamical system viewpoint also motivates control-theoretic analysis, which further recasts the network weight as control. With that, the training process can be viewed as an optimal control problem, as both methodologies aim to optimize some variables (weights v.s. controls) subjected to the chain structure (network v.s. dynamical system). This connection has lead to theoretical characterization of the learning process (Weinan et al., 2018; Hu et al., 2019; Liu & Theodorou, 2019) and practical methods for hyperparameter adaptation (Li et al., 2017b) or computational acceleration (Gunther et al., 2020; Zhang et al., 2019).

Development of algorithmic progress, however, remains relatively limited. This is because OCT-inspired training methods, by construction, are restricted to network class that resembles Markovian state-space models (Liu et al., 2021; Li & Hao, 2018; Li et al., 2017a). This raises questions of their flexibility and scalability to training modern architectures composed of complex dependencies between layers. It is unclear whether this interpretation of dynamical system and optimal control remains suitable, or how it should be adapted, under those cases.

In this work, we address the aforementioned issues using dynamic game theory, a discipline of interactive decision making (Yeung & Petrosjan, 2006) built upon optimal con-

¹Center for Machine Learning ²School of Aerospace Engineering ³School of Electrical and Computer Engineering, Georgia Institute of Technology, USA. Correspondence to: Guan-Horng Liu <ghliu@gatech.edu>.

trol and game theory. Specifically, we propose to treat each layer as a player in a dynamic game connected through the network propagation. The additional dimension gained from multi-player allows us to generalize OCT-inspired methods to accept a much richer network class. Further, introducing game-theoretic analysis, *e.g. information structure*, provides a novel algorithmic connection between different classes of training methods from a Nash equilibria standpoint (Fig. 1).

Unlike prior game-related works, which typically cast the whole network as a player competing over training iteration (Goodfellow et al., 2014; Balduzzi et al., 2018), the (p)layers in our dynamic game interact along the network propagation. This naturally leads to a coalition game since all players share the same objective. The resulting cooperative training scheme urges the network to yield group optimality, or Pareto efficiency (Pardalos et al., 2008). As we will show through experiments, this improves convergence of training modern architectures, as richer information flows between layers to compute the updates. We name our method Dynamic Game Theoretic Neural Optimizer (**DGNOpt**).

Notably, casting the network as a realization of the game has appeared in analyzing the convergence of Back-propagation (Balduzzi, 2016) or contribution of neurons (Stier et al., 2018; Ghorbani & Zou, 2020). Our work instead focuses on developing game-theoretic training methods and how they can be connected to, or generalize, existing optimizers. In summary, we present the following contributions.

- We draw a novel algorithmic characterization from the Nash equilibria perspective by framing the training process as solving a multi-player dynamic game.
- We propose DGNOpt, a game-theoretic optimizer that generalizes OCT-inspired methods to richer network class and encourages cooperative updates among layers with an enlarged information structure.
- Our method achieves competitive results on image classification with residual and inception nets, enabling rich applications from robust control and bandit analysis.

2. Preliminaries

Notation: Given a real-valued function \mathcal{F}_s indexed by $s \in \mathcal{S}$, we shorthand its derivatives evaluated on (x_s, θ_s) as $\nabla_{x_s} \mathcal{F}_s \equiv \mathcal{F}_x^s$, $\nabla_{x_s}^2 \mathcal{F}_s \equiv \mathcal{F}_{xx}^s$, and $\nabla_{x_s} \nabla_{\theta_s} \mathcal{F}_s \equiv \mathcal{F}_{x\theta}^s$, etc. Throughout this work, we will preserve $n \in \{1, \cdots, N\}$ as the player index and $t \in \{0, 1, \cdots, T-1\}$ as the propagation order along the network, or equivalently the stage sequence of the game (see Fig. 1). We will abbreviate them as $n \in [N]$ and $t \in [T]$ for brevity. Composition of functions is denoted by $f(g(\cdot)) \equiv (f \circ g)(\cdot)$. We use \dagger , \odot and \otimes to denote pseudo inversion, Hadamard and Kronecker product. A complete notation table can be found in Appendix A.

2.1. Training Feedforward Nets with Optimal Control

Let the layer propagation rule in feedforward networks (e.g. fully-connected and convolution networks) with depth T be

$$\boldsymbol{z}_{t+1} = f_t(\boldsymbol{z}_t, \boldsymbol{\theta}_t), \quad t \in [T], \tag{1}$$

where z_t and θ_t represent the vectorized hidden state and parameter at each layer t. For instance, $\theta_t := \text{vec}([\boldsymbol{W}_t, \boldsymbol{b}_t])$ for a fully-connected layer, $f_t(z_t, \theta_t) := \sigma(\boldsymbol{W}_t z_t + \boldsymbol{b}_t)$, with nonlinear activation $\sigma(\cdot)$. Equation (1) can be interpreted as a discrete-time Markovian model propagating the state z_t with the tunable variable θ_t . With that, the training process, *i.e.* finding optimal parameters $\{\theta_t : t \in [T]\}$ for all layers, can be described by Optimal Control Programming (OCP),

$$\min_{\theta_t: t \in [T]} L \coloneqq \left[\phi(\boldsymbol{z}_T) + \sum_{t=0}^{T-1} \ell_t(\theta_t) \right] \quad \textit{s.t. (1)}. \quad (2)$$

The objective L consists of a loss ϕ incurred by the network prediction z_T (e.g. cross-entropy in classification) and the layer-wise regularization ℓ_t (e.g. weight decay). Despite (2) considers only one data point z_0 , it can be easily modified to accept batch training (Weinan et al., 2018). Hence, minimizing L sufficiently describes the training process.

Equation (2) provides an OCP characterization of training feedforward networks. First, the optimality principles to OCP, according to standard optimal control theory, typically involve solving some time-dependent objectives recursively from the terminal stage T. Previous works have shown that these backward processes relate closely to the computation of Back-propagation (Li et al., 2017a; Liu et al., 2021). Further, the parameter update of each layer, $\theta_t \leftarrow \theta_t - \delta \theta_t$, can be seen as solving these layer-wise OCP objectives with certain approximations. To ease the notational burden, we leave a thorough discussion in Appendix B. We stress that this intriguing connection is, however, limited to the particular network class described by (1).

2.2. Multi-Player Dynamic Game (MPDG)

Following the terminology in Yeung & Petrosjan (2006), in a discrete-time N-player dynamic game, Player n commits to the action $\theta_{t,n}$ at each stage t and seeks to minimize

$$L_n(\bar{\theta}_n; \bar{\theta}_{\neg n}) := \left[\phi_n(\boldsymbol{x}_T) + \sum_{t=0}^{T-1} \ell_{t,n}(\theta_{t,1}, \dots, \theta_{t,N}) \right]$$
s.t. $\boldsymbol{x}_{t+1} = F_t(\boldsymbol{x}_t, \theta_{t,1}, \dots, \theta_{t,N}), \quad \theta_{t,n} \equiv \theta_{t,n}(\eta_{t,n}), \quad (3)$

where $\bar{\theta}_n \coloneqq \{\theta_{t,n} : t \in [T]\}$ denotes the action sequence for Player n throughout the game. The set $\neg n := \{i \in [N] : i \neq n\}$ includes all players except Player n. The key components that characterize an MPDG (3) are detailed as follows.

• Shared dynamics F_t . The stage-wise propagation rule for x_t , affected by actions across all players $\theta_{t,n}, \forall n$.

- Payoff/Cost L_n . The objective for each player that accumulates the costs $(\phi_n, \ell_{t,n})$ incurred at each stage.
- Information structure $\eta_{t,n}$. A set of information available to Player n at t for making the decision $\theta_{t,n}$.

The Nash equilibria $\{(\bar{\theta}_1^*, \cdots, \bar{\theta}_N^*)\}$ to (3) is a set of stationary points where no player has the incentive to deviate from the decision. Mathematically, this can be described by

$$L_n(\bar{\theta}_n^*; \bar{\theta}_{\neg n}^*) \le L_n(\bar{\theta}_n; \bar{\theta}_{\neg n}^*), \quad \forall n \in [N], \ \forall \bar{\theta}_n \in \bar{\Theta}_n,$$

where $\bar{\Theta}_n$ denotes the set of admissible actions for Player n. When the players agree to cooperate upon an agreement on a set of strategies and a mechanism to distribute the payoff/cost, a cooperative game (CG) of (3) will be formed. CG requires additional optimality principles to be satisfied. This includes (i) group rationality (GR), which requires all players to optimize their joint objective,

$$L^* := \min_{\bar{\theta}_1, \dots, \bar{\theta}_N} \sum_{n=1}^N L_n(\bar{\theta}_n; \bar{\theta}_{\neg n}), \tag{4}$$

and (ii) individual rationality (IR), which requires the cost distributed to each player from L^* be at most the cost he/she will suffer if plays against others non-cooperatively. Intuitively, IR justifies the participation of each player in CG.

3. Dynamic Game Theoretic Perspective

3.1. Formulating DNNs as Dynamic Games

In this section, we draw a novel perspective between the three components in MPDG (3) and the training process of generic (i.e. non-Markovian) DNNs. Given a network composed of the layer modules $\{f_i(\cdot, \theta_i)\}\$, where θ_i denotes the trainable parameters of layer f_i similar to (1), we treat each layer as a player in MPDG. The network can be converted into the form of F_t by indexing i := (t, n), where t represents the sequential order from network input to prediction, and n denotes the index of layers aligned at t. Fig. 2 demonstrates such an example for a residual block. When the network propagation collapses from multiple paths to a single one, we can consider either duplicated players sharing the same path or dummy players with null action space. Hence, w.l.o.g. we will treat N as fixed over t. Notice that the assignment i := (t, n) may not be unique. We will discuss its algorithmic implication later in §5.2.

Once the shared dynamics is constructed, the payoff L_n can be readily linked to the training objective. Since $\ell_{t,n}$ corresponds to the weight decay for layer $f_{t,n}$, it follows that $\ell_{t,n} := \ell_{t,n}(\theta_{t,n})$. Also, we will have $\phi_n := \phi$ whenever all (p)layers share the same task, 1 e.g. in classification. In short, the network architecture and training objective respectively characterize the structure of a dynamic game and its payoff.

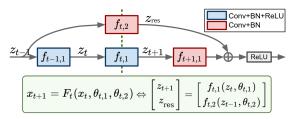


Figure 2. Example of representing a residual block as F_t . Note that x_t augments all hidden states across parallel paths.

3.2. Information Structure and Nash Optimality

We now turn into the role of information structure $\eta_{t,n}$. Standard game-theoretic analysis suggests that $\eta_{t,n}$ determines the type of Nash equilibria inherited in the MPDG (Petrosjan, 2005). Below we introduce several variants that are of our interests, starting from the one with the least structure.

Definition 1 (Open-loop Nash equilibrium (OLNE)). Let $\eta_{t,n}^{o} := \{x_0\}$ be the open-loop information structure. Then a set of action, $\{\theta_{t,n}^* : \forall t, n\}$, provides an OLNE to (3) if

$$\theta_{t,n}^* = \underset{\theta_{t,n}}{\operatorname{arg\,min}} \ H_{t,n}(\boldsymbol{x}_t, \boldsymbol{p}_{t+1,n}, \theta_{t,n}, \theta_{t,\neg n}^*), \tag{5}$$

where
$$\theta_{t,n}^* \equiv \theta_{t,n}^*(\eta_{t,n}^0)$$
 and $H_{t,n} := \ell_{t,n} + F_t^\mathsf{T} \boldsymbol{p}_{t+1,n}$

is the Hamiltonian for Player n at stage t. The co-state $p_{t,n}$ is a vector of the same size as x_t and can be simulated from the backward adjoint process, $(p_{t,n}, p_{T,n}) := (H_x^{t,n}, \phi_x^n)$.

The Hamiltonian objective $H_{t,n}$ varies for each (t,n) and depends on the proceeding stage via co-state $p_{t+1,n}$. When N=1, (5) degenerates to the celebrated Pontryagin principle (Pontryagin et al., 1962), which provides the necessary condition to OCP (2). This motivates the following result.

Proposition 2. Solving $\theta_{t,n}^* = \arg \min H_{t,n}$ with the iterative update, $\theta_{t,n} \leftarrow \theta_{t,n} - M^{\dagger} H_{\theta}^{t,n}$, recovers the descent direction of standard training methods. Specifically, setting

$$\boldsymbol{M} \coloneqq \left\{ \begin{array}{l} \boldsymbol{I} \\ \operatorname{diag} \left(\sqrt{H_{\theta}^{t,n} \odot H_{\theta}^{t,n}} \right) & yields \\ H_{\theta}^{t,n} H_{\theta}^{t,n}^{T} \end{array} \right\} \quad yields \left\{ \begin{array}{l} \operatorname{SGD} \\ \operatorname{RMSprop} \\ \operatorname{Gauss-Newton} \end{array} \right.$$

Proposition 2 provides a similar OCP characterization (c.f. §2.1) except for a more generic network class represented by F_t . It also gives our first game-theoretic interpretation of DNN training: standard training methods implicitly match an OLNE defined upon the network propagation. The proof (see Appendix C) relies on constructing a set of co-state $p_{t,n}$ such that $H_{\theta}^{t,n} \equiv \nabla_{\theta_{t,n}} H_{t,n}$ gives the exact gradient w.r.t. the parameter of layer $f_{t,n}$. The degenerate information structure $\eta_{t,n}^{O}$ implies that optimizers of this class utilize minimal knowledge available from the game (i.e. network) structure. This is in contrast to the following Nash equilibrium which relies on richer information.

¹One of the examples for multi-task objective is the *auxiliary loss* used in deep reinforcement learning (Jaderberg et al., 2016).

Table 1. Dynamic game theoretic perspective of DNN training.

Nash Equilibria	Information Structure	Optimality Principle	Class of Optimizer
OLNE	$\eta_{t,n}^{\mathrm{O}}$	$\min H_{t,n}$ in (5)	Baselines
FNE	$\eta_{t,n}^{ ext{C}}$	$\min Q_{t,n}$ in (6)	DGNOpt (ours)
GR	$\eta_{t,n}^{ ext{C-CG}}$	$\min P_t \text{ in (7)}$	DGNOpt (ours)

Definition 3 (Feedback Nash equilibrium (FNE)). Let $\eta_{t,n}^{\mathcal{C}} := \{ \boldsymbol{x}_s : s \leq t \}$ be the closed-loop information structure. Then a set of strategy, $\{\pi_{t,n}^* : \forall t, n\}$, is called a FNE to (3) *if it is the solution to the Isaacs-Bellman equation (6).*

$$V_{t,n}(\boldsymbol{x}_t) = \min_{\pi_{t,n}} \ Q_{t,n}(\boldsymbol{x}_t, \pi_{t,n}, \pi_{t,\neg n}^*),$$

$$V_{T,n} = \phi_n, \quad \textit{where} \quad Q_{t,n} := \ell_{t,n} + V_{t+1,n} \circ F_t$$
(6)

is the Isaacs-Bellman objective for Player n at stage t. Also, $\pi_{t,n} \equiv \theta_{t,n}(\boldsymbol{x}_t; \eta_{t,n}^C)$ denotes any arbitrary mapping from x_t to $\theta_{t,n}$, conditioned on the closed-loop structure $\eta_{t,n}^C$.

For the closed-loop information structure $\eta_{t,n}^{\mathbb{C}}$, each player has complete access to all preceding states until the current stage t. Consequently, it is preferable to solve for a state-dependent, i.e. feedback, strategy $\pi_{t,n}^*$ rather than a state-independent action $\theta_{t,n}^*$ as in OLNE. Similar to (5), the Isaacs-Bellman objective $Q_{t,n}$ is constructed for each (t,n), except now carrying a function $V_{t,n}(\cdot)$ backward from the terminal stage, rather than the co-state. This value function $V_{t,n}$ summarizes the optimal cost-to-go for Player nfrom each state x_t , provided all afterward stages are minimized accordingly. When N=1, (6) collapses to standard Dynamic Programming (DP; Bellman (1954)), which is an alternative optimality principle parallel to the Pontryagin. For nontrivial N>1, solving the FNE optimality (6) provides a game-theoretic extension for previous DP-inspired training methods, e.g. Liu et al. (2021), to generic (i.e. non-Markovian) architectures.

3.3. Cooperative Game Optimality

Now, let us consider the CG formulation. When a cooperative agreement is reached, each player will be aware of how others react to the game. This can be mathematically expressed by the following information structures,

$$\eta_{t,n}^{\text{O-CG}} \coloneqq \{\boldsymbol{x}_0, \boldsymbol{\theta}_{t,\neg n}^*\} \ \text{ and } \ \eta_{t,n}^{\text{C-CG}} \coloneqq \{\boldsymbol{x}_s, \pi_{t,\neg n}^* : s \leq t\},$$

which enlarge $\eta_{t,n}^{O}$ and $\eta_{t,n}^{C}$ with additional knowledge from other players, $\neg n$. We can characterize the inherited optimality principles similar to OLNE and FNE. Take $\eta_{t,n}^{\text{C-CG}}$ for instance, the joint optimization in GR (4) requires

Definition 4 (Cooperative feedback solution). A set of strategy, $\{\pi_{t,n}^* : \forall t, n\}$, provides an optimal feedback solution to the joint optimization (4) if it solves

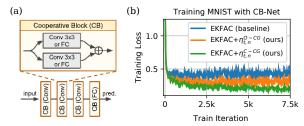


Figure 3. (a) The cooperative-block network and (b) its training performance on MNIST when the optimizer (EKFAC) is exposed to different information structure $\eta_{t,n}$. Note that by Proposition 2, the EKFAC baseline utilizes only the open-loop structure $\eta_{t,n}^{O}$.

$$W_{t}(\mathbf{x}_{t}) = \min_{\pi_{t,n}: n \in [N]} P_{t}(\mathbf{x}_{t}, \pi_{t,1}, \dots, \pi_{t,N}),$$
(7)
$$W_{T} = \sum_{n=1}^{N} \phi_{n}, \quad \text{where } P_{t} := \sum_{n=1}^{N} \ell_{t,n} + W_{t+1} \circ F_{t}$$

$$W_T = \sum_{n=1}^N \phi_n$$
, where $P_t := \sum_{n=1}^N \ell_{t,n} + W_{t+1} \circ F_t$

is the "group-rational" Bellman objective at stage t. $\pi_{t,n} \equiv$ $\theta_{t,n}(m{x}_t;\eta_{t,n}^{ extit{C-CG}})$ denotes arbitrary mapping from $m{x}_t$ to $m{ heta}_{t,n}$, conditioned on the cooperative closed-loop structure $\eta_{t,n}^{C-CG}$.

Notice that (7) is the GR extension of (6). Both optimality principles solve for a set of feedback strategies, except the former considers a joint objective P_t summing over all players. Hence, it is sufficient to carry a joint value function W_t backward. We leave the discussion on $\eta_{t,n}^{\text{O-CG}}$ in Appendix C.

To emphasize the importance of information structure, consider the architecture in Fig. 3a, where each pair of parallel layers shares the same input and output; hence the network resembles a two-player dynamic game with $F_t := f_{t,1} + f_{t,2}$. As shown in Fig. 3b, providing different information structures to the same optimizer, EKFAC (George et al., 2018), greatly affects the training. Having richer information tends to achieve better performance. Additionally, the fact that

$$\eta_{t,n}^{\text{O}} \subset \eta_{t,n}^{\text{C}} \subset \eta_{t,n}^{\text{C-CG}}$$
 and $\eta_{t,n}^{\text{O}} \subset \eta_{t,n}^{\text{O-CG}} \subset \eta_{t,n}^{\text{C-CG}}$ (8)

also implies an algorithmic connection between different classes of optimizers, which we will explore in §4.3.

Table 1 summarizes our game-theoretic analysis. Each information structure suggests its own Nash equilibria and optimality principle, which characterizes a distinct class of training methods. We already established the connection between baselines and $H_{t,n}$ in Proposition 2. In the next section, we will derive methods for solving $(Q_{t,n}, P_t)$.

4. Training DNN by Solving Dynamic Game

In this section, we derive a new second-order method, called Dynamic Game Theoretic Neural Optimizer (**DGNOpt**), that solves (6) and (7) as an alternative to training DNNs. While we will focus on the residual network for its popularity and algorithmic simplicity when deriving the analytic update, we stress that our methodology applies to other architectures. A full derivation is left in Appendix D.

4.1. Iterative Update via Linearization

Computing the game-theoretic objectives $(Q_{t,n}, P_t)$ requires knowing $(F_t, \ell_{t,n}, \phi_n)$. Despite they are well-defined from §3.1, carrying $Q_{t,n}$ or P_t as a stage-varying function is computationally impractical even on a relatively low-dimensional system (Tassa et al., 2012), let alone DNNs. Since the goal is to derive an incremental update given partial (e.g. mini-batch) data at each training iteration, we can consider solving them approximately via *linearization*.

Iterative methods via linearization have been widely used in real-time OCP (Pan et al., 2015; Tassa et al., 2014). We adopt a similar methodology for its computational efficiency and algorithmic connection to existing training methods (shown later). First, consider solving the FNE recursion (6) by $\pi_{t,n}^* \approx \arg\min Q_{t,n}$. We begin by performing second-order Taylor expansion on $Q_{t,n}$ w.r.t. to the variables that are *observable* to Player n at stage t according to $\eta_{t,n}^C$.

$$Q_{t,n} \approx \frac{1}{2} \begin{bmatrix} \mathbf{1} \\ \delta \mathbf{x}_t \\ \delta \boldsymbol{\theta}_{t,n} \end{bmatrix}^\mathsf{T} \begin{bmatrix} Q_{t,n} & Q_{\boldsymbol{x}}^{t,n}^\mathsf{T} & Q_{\boldsymbol{\theta}}^{t,n}^\mathsf{T} \\ Q_{\boldsymbol{x}}^{t,n} & Q_{\boldsymbol{x}\boldsymbol{x}}^{t,n} & Q_{\boldsymbol{\theta}\boldsymbol{x}}^{t,n}^\mathsf{T} \\ Q_{\boldsymbol{\theta}}^{t,n} & Q_{\boldsymbol{\theta}\boldsymbol{x}}^{t,n} & Q_{\boldsymbol{\theta}\boldsymbol{\theta}}^{t,n} \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \delta \mathbf{x}_t \\ \delta \boldsymbol{\theta}_{t,n} \end{bmatrix}$$

Note that $\delta\theta_{t,\neg n}$ does not appear in the above quadratic expansion since it is unobservable according to $\eta_{t,n}^{\mathsf{C}}$. The derivatives of $Q_{t,n}$ w.r.t. different arguments follow standard chain rule (recall $Q_{t,n} := \ell_{t,n} + V_{t+1,n} \circ F_t$), with the dynamics linearized at some fixed point $(\boldsymbol{x}_t, \theta_{t,n})$, e.g.

$$Q_{\theta}^{t,n} = \ell_{\theta}^{t,n} + {F_{\theta}^{t}}^\mathsf{T} V_{\boldsymbol{x}}^{t+1,n}, \; Q_{\theta \boldsymbol{x}}^{t,n} = {F_{\theta}^{t}}^\mathsf{T} V_{\boldsymbol{x} \boldsymbol{x}}^{t+1,n} F_{\boldsymbol{x}}^{t}.$$

The analytic solution to this quadratic expression is given by $\pi_{t,n}^* = \theta_{t,n} - \delta \pi_{t,n}^*$, with the incremental update $\delta \pi_{t,n}^*$ being

$$\delta \pi_{t,n}^* = \mathbf{k}_{t,n} + \mathbf{K}_{t,n} \delta \mathbf{x}_t.$$

$$\mathbf{k}_{t,n} \coloneqq (Q_{\theta\theta}^{t,n})^{\dagger} Q_{\theta}^{t,n} \quad \text{and} \quad \mathbf{K}_{t,n} \coloneqq (Q_{\theta\theta}^{t,n})^{\dagger} Q_{\theta\pi}^{t,n}$$
(9)

are called the open and feedback gains. The superscript † denotes the pseudo inversion. Note that $\delta\pi_{t,n}^*$ is only *locally* optimal around the region where the quadratic expansion remains valid. Since x_t augments preceding hidden states (e.g. $x_t := [z_t, z_{t-1}]^{\mathsf{T}}$ in Fig. 2), (9) implies that preceding hidden states contribute to the update via linear superposition.

Substituting the incremental update $\delta\pi_{t,n}^*$ back to the FNE recursion (6) yields the local expression of the value function $V_{t,n}$, which will be used to compute the preceding update $\delta\pi_{t-1,n}^*$. Since the computation depends on $V_{t,n}$ only through its local derivatives $V_{x}^{t,n}$ and $V_{xx}^{t,n}$, it is sufficient to propagate these quantities rather than the function itself. The propagation formula is summarized in (10). This procedure (line 4-7 in Alg. 1) repeats recursively backward from the terminal to initial stage, similar to Back-propagation.

$$V_{\boldsymbol{x}}^{t,n} = Q_{\boldsymbol{x}}^{t,n} - Q_{\boldsymbol{x}\theta}^{t,n} \mathbf{k}_{t,n}, \quad V_{\boldsymbol{x}}^{T,n} = \phi_{\boldsymbol{x}}^{n},$$

$$V_{\boldsymbol{x}\boldsymbol{x}}^{t,n} = Q_{\boldsymbol{x}\boldsymbol{x}}^{t,n} - Q_{\boldsymbol{x}\theta}^{t,n} \mathbf{K}_{t,n}, \quad V_{\boldsymbol{x}\boldsymbol{x}}^{T,n} = \phi_{\boldsymbol{x}\boldsymbol{x}}^{n}.$$
(10)

Algorithm 1 Dynamic Game Theoretic Neural Optimizer

```
1: Input: dataset \mathcal{D}, network F \equiv \{F_t : t \in [T]\}
 3:
         Compute x_t by propagating x_0 \sim \mathcal{D} through F

⊳ Solve FNE or GR

 4:
         for t = T - 1 to 0 do
            Solve the update \delta \pi_{t,n}^* with (9) or (11)
 5:
            Solve (V_x^{t,n}, V_{xx}^{t,n}) or (W_x^t, W_{xx}^t) with (10) or (25)
 6:
 7:
 8:
         Set x_0' = x_0
         for t = 0 to T-1 do

    □ Update parameter

 9:
            Apply \theta_{t,n} \leftarrow \theta_{t,n} - \delta \pi_{t,n}^*(\delta x_t) with \delta x_t = x_t' - x_t
10:
            Compute \mathbf{x}'_{t+1} = F_t(\mathbf{x}'_t, \theta_{t,1}, \cdots, \theta_{t,N})
11:
12:
         end for
13: until converges
```

Derivation for CG follows similar steps except we consider solving the GR recursion (7) by $\pi_{t,n}^* \approx \arg\min P_t$. Since all players' actions are now observable from $\eta_{t,n}^{\text{C-CG}}$, we need to expand P_t w.r.t. all arguments. For notational simplicity, let us denote $u \equiv \theta_{t,1}, v \equiv \theta_{t,2}$ in Fig. 2. In the case when each player minimizes P_t independently without knowing the other, we know the non-cooperative update for Player 2 admits the form² of $\delta v_t = \mathbf{I}_t + \mathbf{L}_t \delta x_t$. Now, the locally-optimal cooperative update for Player 1 can be written as

$$\delta \pi_{t,1}^* = \widetilde{\mathbf{k}}_t + \widetilde{\mathbf{K}}_t \delta \mathbf{x}_t, \text{ where}$$

$$\widetilde{\mathbf{k}}_t := \widetilde{P}_{uu}^{t\dagger} (P_u^t - P_{uu}^t \mathbf{I}_t), \ \widetilde{\mathbf{K}}_t := \widetilde{P}_{uu}^{t\dagger} (P_{ux}^t - P_{uu}^t \mathbf{L}_t).$$

Similar equations can be derived for Player 2. We will refer $\widetilde{P}_{uu}^t := P_{uu}^t - P_{uv}^t P_{vv}^t P_{vu}^t$ as the *cooperative precondition*. The update (11), despite seemly complex, exhibits intriguing properties. For one, notice that computing the cooperative open gain $\widetilde{\mathbf{k}}_t$ for Player 1 involves the non-cooperative open gain \mathbf{I}_t from Player 2. In other words, each player adjusts the strategy after knowing the companion's action. Similar interpretation can be drawn for the feedbacks $\widetilde{\mathbf{K}}_t$ and \mathbf{L}_t . Propagation of (W_x^t, W_{xx}^t) follows similarly as (10) once all players' updates are computed. We leave the complete formula in (25) (see Appendix D.1) since it is rather tedious.

Let us discuss the role of δx_t and how to compute them. Conceptually, δx_t can be any deviation away from the fixed point x_t where we expand the objectives, $Q_{t,n}$ or P_t . In MPDG application, it is typically set to the state difference when the parameter updates are applied until stage t,

$$\delta \boldsymbol{x}_t \coloneqq (F_{t-1} \circ \cdots \circ F_0)(\boldsymbol{x}_0, \{\theta + \delta \pi^*\}_{< t, \forall n}) - \boldsymbol{x}_t, (12)$$

where $\{\delta\pi^*\}_{< t, \forall n} := \{\delta\pi^*_{s,n} : s < t, \forall n\}$ collects all players' updates until stage t. In this view, the feedback compensates all changes, including those that may cause instability, cascading from the preceding layers; hence it tends to robustify the training process (Pantoja, 1988; Liu et al., 2021).

²Similar to (9), we have $\mathbf{I}_t := P_{vv}^{t} \stackrel{\dagger}{P}_v^t$ and $\mathbf{L}_t := P_{vv}^{t} \stackrel{\dagger}{P}_{vx}^t$.

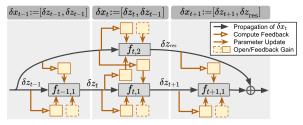


Figure 4. How $\delta \pi_{t,n}^*$ is applied (c.f. line 8-12 in Alg. 1) to a residual block. We compute δx_t with a forward propagation (12) and simultaneously update the parameter. The solid yellow box denotes the feedback dependent on the preceding hidden states $z_{s< t}$.

Alg. 1 presents the pseudo-code of DGNOpt, which consists of (i) the same forward propagation through the network (line 3), (ii) a distinct game-theoretic backward process that solves either FNE or GR optimality (line 4-7), and (iii) an additional forward pass that applies the feedback updates $\delta \pi_{t,n}^*$ (line 8-12; also Fig. 4). We stress that Alg. 1 accepts any generic DNN so long as it can be represented by F_t .

4.2. Curvature Approximation

Naively inverting the parameter curvature, i.e. $(Q_{\theta\theta}^{t,n})^{\dagger}$ and \widetilde{P}_{uu}^{t} , can be computationally inefficient and sometimes unstable for practical training. To mitigate the issue, we adopt curvature amortizations (Kingma & Ba, 2014; Hinton et al., 2012) used in DNN training. These methods naturally fit into our framework by recalling Proposition 2 that different baselines differ in how they estimate the curvature $H_{\theta\theta}^{t,n}$ for the preconditioned update $H_{\theta\theta}^{t,n}^{\dagger}H_{\theta}^{t,n}$. With this in mind, we can estimate the FNE parameter curvature $Q_{\theta\theta}^{t,n}$ with

$$Q_{\theta\theta}^{t,n} \approx Q_{\theta}^{t,n} Q_{\theta}^{t,n^{\mathsf{T}}} \text{ or } \operatorname{diag}(\sqrt{Q_{\theta}^{t,n} \odot Q_{\theta}^{t,n}}),$$
 (13)

which resembles the Gauss-Newton (GN) matrix or its adaptive diagonal matrix (as appeared in RMSProp and Adam).

As for \widetilde{P}_{uu}^t , which contains an *inner* inversion P_{vv}^t inside \widetilde{P}_{uu}^t , we propose a new approximation inspired by the Kronecker factorization (KFAC; Martens & Grosse (2015)). KFAC factorizes the GN matrix with two smaller-size matrices. We leave the complete discussion on KFAC, as well as the proof of the following result, in Appendix D.2.

Proposition 5 (KFAC for \widetilde{P}_{uu}^t). Suppose P_{uu}^t and P_{vv}^t are factorized with some vectors $\mathbf{z}_1, \mathbf{z}_2, \mathbf{g}_1, \mathbf{g}_2$ by

$$P_{uu}^{t} \approx \mathbb{E}[z_{1}z_{1}^{\mathsf{T}}] \otimes \mathbb{E}[g_{1}g_{1}^{\mathsf{T}}] =: A_{uu} \otimes B_{uu},$$

$$P_{uu}^{t} \approx \mathbb{E}[z_{2}z_{1}^{\mathsf{T}}] \otimes \mathbb{E}[g_{2}g_{1}^{\mathsf{T}}] =: A_{uu} \otimes B_{uu},$$

where the expectation is taken over the mini-batch data. Let $A_{uv} := \mathbb{E}[z_1 z_2^{\mathsf{T}}]$ and $B_{uv} := \mathbb{E}[g_1 g_2^{\mathsf{T}}]$, then the cooperative precondition matrix in (11) can be factorized by

$$\widetilde{P}_{uu}^{t} \approx \widetilde{A}_{uu} \otimes \widetilde{B}_{uu}$$

$$= (A_{uu} - A_{uv} A_{vv}^{\dagger} A_{uv}^{\mathsf{T}}) \otimes (B_{uu} - B_{uv} B_{vv}^{\dagger} B_{uv}^{\mathsf{T}}).$$
(14)

In practice, we set $(z_1, z_2, g_1, g_2) := (z_t, z_{t-1}, W_{z_t}^{t+1}, W_{z_{t-1}}^{t+1})$ for the residual block in Fig. 2 or 4. With Proposition 5, we can compute the update, take $\tilde{\mathbf{k}}_t$ for instance, by

$$\widetilde{\mathbf{k}}_{t} = \operatorname{vec}(\widetilde{B}_{uu}^{\dagger}(P_{u}^{t} - B_{uv}\operatorname{vec}^{\dagger}(\mathbf{I}_{t})A_{uv}^{\mathsf{T}})\widetilde{A}_{uu}^{-\mathsf{T}}), \quad (15)$$

where vec^{\dagger} is the inverse operation of vectorization (vec).

Another computation source comes from the curvature w.r.t. the MPDG state, *i.e.* $V_{xx}^{t,n}$ and W_{xx}^t . Here, we approximate them with low-rank matrices using either Gauss-Newton or their top eigenspace. These are rather reasonable approximations since it has been constantly observed that these Hessians are highly degenerate for DNNs (Wu et al., 2020; Papyan, 2019; Sagun et al., 2017). With all these, we are able to train modern DNNs by solving their corresponding dynamic games, (6) or (7), with a runtime comparable to other first and second-order methods (see Fig. 7).

4.3. Algorithmic Connection

Finally, let us discuss an intriguing algorithmic equivalence. Recall the subset relation among the information structures in (8). Manipulating these structures allows one to traverse between different game optimality principles. For instance, masking $\pi_{t,\neg n}^*$ in $\eta_{t,n}^{\text{C-CG}}$ makes it degenerate to $\eta_{t,n}^{\text{C}}$, which implies the FNE and GR optimality become equivalent. Through this lens, one may wonder if a similar algorithmic relation can be drawn for these iterative updates. This is indeed the case as shown below (proof left in Appendix D.3).

Theorem 6 (Algorithmic equivalence).

- (11) with $P_{uv}^t := \mathbf{0}$ gives (9)
- (9) with $(Q_{\theta x}^{t,n}, Q_{\theta \theta}^{t,n}) \coloneqq (\mathbf{0}, \mathbf{I})$ gives SGD
- (11) with $(P_{uv}^t, P_{ux}^t, P_{ux}^t) := (\mathbf{0}, \mathbf{0}, \mathbf{I})$ gives SGD

Setting $Q_{\theta\theta}^{t,n}$ and P_{uu}^t to other precondition matrices, similar to (13), recovers other baselines.

The intuition behind Theorem 6 is that when higher-order (>2) expansions are discarded, setting $P_{uv}^t \coloneqq \mathbf{0}$ completely blocks the communication between two players; therefore we effectively remove $\pi_{t,\neg n}^*$ from $\eta_{t,n}^{\text{C-CG}}$. Similarly, forcing $Q_{\theta x}^{t,n} \coloneqq \mathbf{0}$ prevents Player n from observing how changing x_t may affect the payoff, hence one can at best achieve the same OLNE optimality as baselines. Theorem 6 implies that (9) and (11) generalize standard updates to richer information structure; thereby creating more complex updates.

5. Experiment

5.1. Evaluation on Classification Datasets

Datasets and networks. We verify the performance of DGNOpt on image classification datasets as they are suitable

Table 2. Accuracy (%) of residual-based networks (averaged over 6 random seeds)

Dataset	Baselines (i.e. $\min H_{t,n}$ in OLNE)				DGNOpt	
	SGD	RMSProp	Adam	EKFAC	EMSA	(ours)
MNIST	98.65	98.61	98.49	98.77	98.25	98.76
SVHN	88.58	88.96	89.20	88.75	87.40	89.22
CIFAR10	82.94	83.75	85.66	85.65	75.60	85.85
CIFAR100	71.78	71.65	71.96	71.95	62.63	72.24

Table 3. Accuracy (%) of inception-based networks (averaged over 4 random seeds)

Dataset	SGD	Baselines (i.e RMSProp		*	E) EMSA	DGNOpt (ours)
MNIST	97.96	97.75	97.72	97.90	97.39	98.03
SVHN	87.61	86.14	86.84	88.89	82.68	88.94
CIFAR10	76.66	74.38	75.38	77.54	70.17	77.72

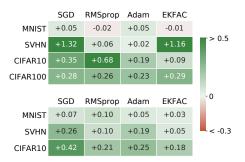


Figure 5. Accuracy (%) improvement (+) or degradation (-) when richer information structure, *i.e.* $\eta_{t,n}^{0} \rightarrow \{\eta_{t,n}^{C}, \eta_{t,n}^{C-CG}\}$, is used for each best-tuned baseline³ in Table 2 (upper) and Table 3 (bottom). Color bar is scaled for best view.

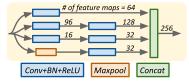


Figure 6. Architecture of the inception block.

testbeds for modern networks that contain non-Markovian dependencies. Specifically, we first consider residual-based networks given their popularity and our thorough discussions in §4. For larger datasets such as CIFAR10/100, we train ResNet18 with multi-stepsize learning rate decay. For MNIST and SVHN, we use residual networks composed of 3 residual blocks (see Fig. 2). Meanwhile, we also consider inception-based networks, which composed of an inception block (see Fig. 6) that resembles a 4-player dynamic game. All networks use ReLU activation and are trained with 128 batch size. Other setups are detailed in Appendix E.

Baselines. Motivated by our discussion in §3, we compare DGNOpt, which essentially solves FNE and GR, with methods involving OLNE either implicitly or explicitly. This includes standard training methods such as SGD, RMSprop, Adam, and EKFAC (George et al., 2018), which is an extension to the second-order method KFAC with eigenvalue-correction. To also compare against OCT-inspired methods, we include EMSA (Li et al., 2017a), which explicitly minimizes a modified Hamiltonian. Other OCT-based training methods mostly consider degenerate, *e.g.* discrete-weighted (Li & Hao, 2018) or Markovian (Liu et al., 2021), networks. In this view, DGNOpt generalizes those methods to both larger network class and richer information structure.

Performance and ablation study. Table 2 and 3 summarize the performance for the residual and inception networks. On most datasets, DGNOpt achieves competitive results against standard methods and outperforms EMSA by a large margin. Despite both originates from the OCT methodology, in practice EMSA often exhibits numerical instability for larger networks. On the contrary, DGNOpt leverages

iteration-based linearization and amortized curvature, which greatly stabilizes the training.

On the other hand, DGNOpt distinguishes itself from standard baselines by considering a larger information structure. To validate the benefit of having this additional knowledge during training, we conduct an ablation study using the algorithmic connection built in Theorem 6. Specifically, we measure the performance difference when the best-tuned baselines, *i.e.* the ones we report in Table 2 and 3, are further allowed to utilize higher-level information. Algorithmically, this can be done by running DGNOpt with the parameter curvature replaced by the precondition matrix of each baseline. For instance, replacing all $Q_{\theta\theta}^{t,n}$ with identity matrices Iwhile keeping other computation unchanged is equivalent to lifting SGD to accept the closed-loop structure $\eta_{t,n}^{\mathbb{C}}$. From Theorem 6, these two training procedures now differ only in the presence of $Q_{\theta x}^{t,n}$, which allows SGD to adjust its update based on the change of $x_t \in \eta_{t,n}^{\mathbb{C}}$. As shown in Fig. 5, enlarging the information structure tends to enhance the performance, or at least being innocuous. We highlight these improvements as the benefit gained from introducing dynamic game theory to the original OCT interpretation.

Overhead vs performance trade-offs. As shown in Fig. 7, DGNOpt enjoys a comparable runtime and memory complexity to standard methods on training ResNet18. Specifically, its per-iteration runtime is around ±40% compared to the second-order baseline, depending on the information structures (DGNOpt-FNE ν.s. DGNOpt-GR). In practice, these gaps tend to vanish for smaller networks. The overhead introduced by DGNOpt enables the computation of *feedback updates* using a richer information structure. From the OCT standpoint, the feedback is known to play a key role in compensating the unstable disturbance along the propagation. Particularly, when problems inherit chained

³The ablation analysis in Fig. 5 applies Theorem 6 to methods that solve the exact Hamiltonian; hence excludes EMSA since it instead considers a modified Hamiltonian (see Appendix E).

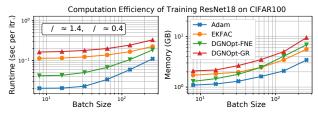


Figure 7. Our second-order method DGNOpt exhibits similar runtime ($\pm 40\%$) and memory ($\pm 30\%$) complexity compared to the second-order baseline EKFAC.

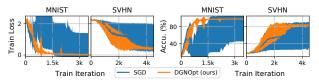


Figure 8. Training inception-based networks using larger step sizes, where MNIST (resp. SVHN) uses lr=1.0 (resp. lr=2.0).

constraints (e.g. DNNs), these feedback-enhanced methods often converge faster with superior numerical stability against standard methods (Murray & Yakowitz, 1984).

To validate the role of feedback in training modern DNNs, notice that one shall expect the effect of feedback becomes significant when a larger step size is taken. This is because (see (12)) larger $\delta\pi^*$ increases δx_t , which amplifies the feedback $\mathbf{K}\delta x_t$. Fig. 8 confirms our hypothesis, where we train the inception-based networks on MNIST and SVHN using relatively large learning rates. It is clear that utilizing feedback updates greatly stabilizes the training. While the SGD baseline struggles to make stable progress, DGNOpt converges almost flawlessly (with negligible overhead). As for well-tuned hyperparameter which often has a smaller step size, our ablation analysis in Fig. 5 suggests that having feedbacks throughout the stochastic training generally leads to better local minima.

5.2. Game-Theoretic Applications

Cooperative training with fictitious agents. Despite all the rigorous connection we have explored so far, it is perhaps unsatisfactory to see our multi-agent analysis degenerates when facing feedforward networks, since the number of player N becomes trivially 1. We can remedy this scenario by considering the following transformation.

$$F_t(\boldsymbol{z}_t, \theta_{t,1}, \dots, \theta_{t,N}) \coloneqq f_t(\boldsymbol{z}_t, \theta_t), \ \sum_{n=1}^N \theta_{t,n} = \theta_t \ (16)$$

In other words, we can divide the layer's weight (or player's action) into multiple parts, so that the MPDG framework remains applicable. Interestingly, the transformation of this kind resembles game-theoretic robust optimal control (Pan et al., 2015; Sun et al., 2018), where the controller (or player in our context) models external disturbances with *fictitious agents*, in order to enhance the robustness or convergence

Table 4. Convergence speed w.r.t. N. Numerical values report the training steps required to achieve certain accuracy on each dataset.

Achieved	Number of Player (N) in SGD				
Performance	1	2	4	6	
	5.14k	2.31k	1.25k	0.8k	
60% in CIFAR10	3.62k	2.97k	2.98k	5.83k	

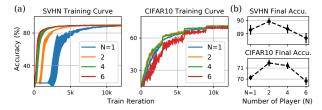


Figure 9. (a) Training curve and (b) final accuracy as we vary the number of player (N) as a hyperparameter of game-extended SGD.

of the optimization process.

Fig. 9 and Table 4 provide the training results when SGD presumes different numbers of players interacting in a feed-forward network consisting of 4 convolution and 2 fully-connected layers (see Appendix E). Notice that N=1 corresponds to the original method. For N>1, we apply the transformation (16) then solve for the cooperative update as in DGNOpt. We stress that these fictitious agents only appear during the training phase for computing the cooperative updates. At inference, actions from all players collapse back to θ_t by the summation in (16).

While it is clear that encouraging agents to cooperate during training can achieve better minima at a faster rate, having more agents, surprisingly, does not always imply better performance. In practice, the improvement can slow down or even degrade once N passes some critical values. This implies that N shall be treated as a hyper-parameter of these game-extended methods. Empirically, we find that $N{=}2$ provides a good trade-off between the final performance and convergence speed. We observe a consistent result for this setup on other optimizers (see Appendix E for EKFAC).

Adaptive alignment using multi-armed bandit. Finally, let us discuss an application of the bandit algorithm in our framework. In $\S 3.1$, we briefly mentioned that mapping from modern networks to the shared dynamics F_t most likely will not be unique. For instance (see Fig. 10a), placing the shortcut module of a residual block at different locations leads to different F_t ; hence results in different DGNOpt updates. This is a distinct feature arising exclusively from our MPDG framework, since these alignments are unrecognizable to standard baselines. It naturally raises the following questions: what is the optimal strategy to align the (p)layers of the network in our dynamic game, and how do different aligning strategies affect training?

Table 5. Training result (accuracy %) of Fig. 10 on two datasets.

Dataset	EKFAC	DGNO fixed	pt + Aligni random	ng Strategy adaptive
SVHN	87.49	88.20	88.12	88.33
CIFAR10	84.67	85.20	85.27	85.65

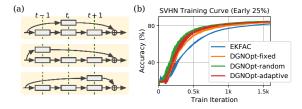


Figure 10. (a) Different alignments of the same residual block leads to distinct DGNOpt updates, yet they are unrecognizable to baselines. (b) Early phase of training with different aligning strategies.

To answer these questions, we compare the performance between three strategies, including (i) using a fixed alignment throughout training, (ii) random alignment at each iteration, and (iii) adaptive alignment using a multi-armed bandit. For the last case, we interpret pulling an arm as selecting one of the alignments and associate the round-wise reward with the validation accuracy at each iteration. Note that this is a non-stationary bandit problem since the reward distribution of each arm/alignment evolves as we train the network. We provide the pseudo-code of this procedure in Appendix E.

Fig. 10b and Table 5 report the results of DGNOpt using different aligning strategies. We also include the baseline when the information structure shrinks from $\eta_{t,n}^{\text{C-CG}}$ to $\eta_{t,n}^{\text{O}}$, similar to the ablation study in §5.1. In this case, all these DGNOpt variants degenerate to EKFAC. For the non-stationary bandit, we find EXP3++ (Seldin & Slivkins, 2014) to be sufficient in this application. While DGNOpt with fixed alignment already achieves faster convergence compared with the baseline, dynamic alignment using either random or adaptive strategy leads to further improvement (see Fig. 10b). Notably, having the adaptation throughout training also enhances the final accuracy. For CIFAR10 with ResNet18, the value is boost by 1% from baseline and 0.5% compared with the other two strategies. This sheds light on new algorithmic opportunities inspired by *architecture-aware* optimization.

6. Discussion

Comparison to Markovian-based OCT-inspired methods. As we briefly mentioned in §3.2, our DGNOpt (with FNE) can be seen as a game-theoretic extension of Liu et al. (2021), which is also an OCT-inspired method despite concerning only Markovian networks. It is natural to wonder whether these two methods are interchangeable since one can always force a non-Markovian system to be Markovian by lifting it into higher dimensions or aggregating the state.

Here, we stress that our DGNOpt differs from Liu et al. (2021) in many significant ways. For one, forming a Markovian chain by grouping the non-Markovian layers into higher dimensions leads to a *degenerate* information structure. The (p)layers inside each Markovian group, $\{f_{t,n}: t^- < t < t^+\}$, only have access to $x_{s \le t^-}$ rather than full latest information $x_{s \le t}$ as in DGNOpt, since their dependencies are discarded. From the Nash standpoint, this leads to degenerate backward optimality and update rules. Indeed, in the limit when we simply group the whole network as single-step dynamics, we will recover $\eta_{t,n}^{O} \triangleq \{x_0\}$ in baselines. In contrast, DGNOpt fully leverages the structural relation of the network, hence enables rich game-based applications, *e.g.* bandit or robust control, that are otherwise infeasible with Liu et al. (2021).

Degeneracy when partitioning parameters as players. In §5.2, we demonstrate a specific transformation, *i.e.* (16), that makes cooperative training possible for single-player feedforward networks while respecting our *layer-as-player* game formulation. This transformation may seem artificial at first glance compared to a naive alternative that directly partitions the parameters of each layer as distinct players. Unfortunately, the latter strategy yields *degenerate* cooperative optimality. To see it, notice that treating the $\pi_{t,n}$ appeared in the joint optimization (7) as the n^{th} -partitioned parameters of layer t is *equivalent* to solving the FNE optimality (6) with N=1 (so that the $\pi_{t,N=1}$ in (6) becomes the intact parameters of layer t). Hence, it collapses to the prior single-player non-cooperative method (Liu et al., 2021),

which converges much slower than our DGNOpt (Fig. 11). Our proposed transformation (16) enables collaborative control and may be naturally extended to other robust formulations, *e.g.* minimax adversarial training.

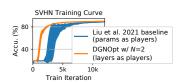


Figure 11. Convergence using different dynamic game formulations with the same setup as Fig. 9a.

7. Conclusion

In this work, we introduce a novel game-theoretic characterization by bridging the training process of DNN with a multiagent dynamic game. The inspired optimizer, DGNOpt, generalizes previous OCT-based methods to generic network class and encourages cooperative updates to improve the performance. Our work pushes forward principled algorithmic design from OCT and game theory.

Acknowledgements

G.H. Liu was supported by CPS NSF Award #1932068, and T. Chen was supported by ARO Award #W911NF2010151. The authors thank C.H. Lin, Y. Pan, C.W. Kuo, M. Gandhi, E. Evans, and Z. Wang for many helpful discussions.

References

- Balduzzi, D. Deep online convex optimization with gated games. *arXiv preprint arXiv:1604.01952*, 2016.
- Balduzzi, D., Racaniere, S., Martens, J., Foerster, J., Tuyls, K., and Graepel, T. The mechanics of n-player differentiable games. *arXiv preprint arXiv:1802.05642*, 2018.
- Bellman, R. The theory of dynamic programming. Technical report, Rand corp santa monica ca, 1954.
- Chang, B., Meng, L., Haber, E., Ruthotto, L., Begert, D., and Holtham, E. Reversible architectures for arbitrarily deep residual neural networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Chen, T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems*, pp. 6572–6583, 2018.
- George, T., Laurent, C., Bouthillier, X., Ballas, N., and Vincent, P. Fast approximate natural gradient descent in a kronecker factored eigenbasis. In *Advances in Neural Information Processing Systems*, pp. 9550–9560, 2018.
- Ghorbani, A. and Zou, J. Neuron shapley: Discovering the responsible neurons. *arXiv preprint arXiv:2002.09815*, 2020.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. In *Advances in neural* information processing systems, pp. 2672–2680, 2014.
- Greydanus, S., Dzamba, M., and Yosinski, J. Hamiltonian neural networks. In *Advances in Neural Information Processing Systems*, pp. 15353–15363, 2019.
- Gunther, S., Ruthotto, L., Schroder, J. B., Cyr, E. C., and Gauger, N. R. Layer-parallel training of deep residual neural networks. SIAM Journal on Mathematics of Data Science, 2(1):1–23, 2020.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Hinton, G., Srivastava, N., and Swersky, K. Neural networks for machine learning lecture 6a overview of mini-batch gradient descent. 2012.
- Hu, K., Kazeykina, A., and Ren, Z. Mean-field langevin system, optimal control and deep neural networks. *arXiv* preprint arXiv:1909.07278, 2019.

- Jaderberg, M., Mnih, V., Czarnecki, W. M., Schaul, T., Leibo, J. Z., Silver, D., and Kavukcuoglu, K. Reinforcement learning with unsupervised auxiliary tasks. arXiv preprint arXiv:1611.05397, 2016.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Li, Q. and Hao, S. An optimal control approach to deep learning and applications to discrete-weight neural networks. *arXiv* preprint arXiv:1803.01299, 2018.
- Li, Q., Chen, L., Tai, C., and Weinan, E. Maximum principle based algorithms for deep learning. *The Journal of Machine Learning Research*, 18(1):5998–6026, 2017a.
- Li, Q., Tai, C., and E, W. Stochastic modified equations and adaptive stochastic gradient algorithms. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2101–2110. JMLR. org, 2017b.
- Liu, G.-H. and Theodorou, E. A. Deep learning theory review: An optimal control and dynamical systems perspective. arXiv preprint arXiv:1908.10920, 2019.
- Liu, G.-H., Chen, T., and Theodorou, E. A. Ddpnopt: Differential dynamic programming neural optimizer. In *International Conference on Learning Representations*, 2021.
- Lu, Y., Zhong, A., Li, Q., and Dong, B. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations. *arXiv* preprint *arXiv*:1710.10121, 2017.
- Lu, Y., Ma, C., Lu, Y., Lu, J., and Ying, L. A mean-field analysis of deep resnet and beyond: Towards provable optimization via overparameterization from depth. *arXiv* preprint arXiv:2003.05508, 2020.
- Martens, J. and Grosse, R. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417, 2015.
- Murray, D. and Yakowitz, S. Differential dynamic programming and newton's method for discrete optimal control problems. *Journal of Optimization Theory and Applications*, 43(3):395–414, 1984.
- Pan, Y., Theodorou, E., and Bakshi, K. Robust trajectory optimization: A cooperative stochastic game theoretic approach. In *Robotics: Science and Systems*, 2015.
- Pantoja, J. F. A. Differential dynamic programming and newton's method. *International Journal of Control*, 47 (5):1539–1553, 1988.

- Papyan, V. Measurements of three-level hierarchical structure in the outliers in the spectrum of deepnet hessians. *arXiv* preprint arXiv:1901.08244, 2019.
- Pardalos, P. M., Migdalas, A., and Pitsoulis, L. Pareto optimality, game theory and equilibria, volume 17. Springer Science & Business Media, 2008.
- Petrosjan, L. A. Cooperative differential games. In *Advances in dynamic games*, pp. 183–200. Springer, 2005.
- Pontryagin, L. S., Mishchenko, E., Boltyanskii, V., and Gamkrelidze, R. The mathematical theory of optimal processes. 1962.
- Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian of over-parametrized neural networks. arXiv preprint arXiv:1706.04454, 2017.
- Seldin, Y. and Slivkins, A. One practical algorithm for both stochastic and adversarial bandits. In *International Con*ference on Machine Learning, pp. 1287–1295. PMLR, 2014.
- Stier, J., Gianini, G., Granitzer, M., and Ziegler, K. Analysing neural network topologies: a game theoretic approach. *Procedia Computer Science*, 126:234–243, 2018.
- Sun, W., Pan, Y., Lim, J., Theodorou, E. A., and Tsiotras, P. Min-max differential dynamic programming: Continuous and discrete time formulations. *Journal of Guidance*, *Control, and Dynamics*, 41(12):2568–2580, 2018.
- Tassa, Y., Erez, T., and Todorov, E. Synthesis and stabilization of complex behaviors through online trajectory optimization. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 4906–4913. IEEE, 2012.
- Tassa, Y., Mansard, N., and Todorov, E. Control-limited differential dynamic programming. In 2014 IEEE International Conference on Robotics and Automation (ICRA), pp. 1168–1175. IEEE, 2014.
- Weinan, E. A proposal on machine learning via dynamical systems. *Communications in Mathematics and Statistics*, 5(1):1–11, 2017.
- Weinan, E., Han, J., and Li, Q. A mean-field optimal control formulation of deep learning. *arXiv* preprint *arXiv*:1807.01083, 2018.
- Wu, Y., Zhu, X., Wu, C., Wang, A., and Ge, R. Dissecting hessian: Understanding common structure of hessian in neural networks. *arXiv preprint arXiv:2010.04261*, 2020.

- Yeung, D. W. and Petrosjan, L. A. Cooperative stochastic differential games. Springer Science & Business Media, 2006
- Zhang, D., Zhang, T., Lu, Y., Zhu, Z., and Dong, B. You only propagate once: Accelerating adversarial training via maximal principle. *arXiv preprint arXiv:1905.00877*, 2019.

Supplementary Material

A. Notation Summary

Table 6. Abbreviation.

OCT/OCP	Optimal Control Theory/Programming
MPDG	Multi-Player Dynamic Game
CG	Cooperative Game
OLNE	Open-loop Nash Equilibria
FNE	Feedback Nash Equilibria
GR	Group Rationality
IR	Individual Rationality

Table 7. Terminology mapping.

	MPDG	Training generic (non-Markovian) DNN	Ns
t	Stage order	Computation order from input to output	Layer index (t, m)
n	Player index	Index of parallel layers aligned at t	Layer index (t, n)
$f_{t,n}$	-	Layer module indexed by (t, n)	
F_t	Shared dynamics	Joint propagation rule of layers $\{f_{t,n}: n\}$	$\in [N]$
$\theta_{t,n}$	Action committed at stage t by Player n	Trainable parameter of layer $f_{t,n}$	
$\boldsymbol{z}_{t,n}$	-	Pre-activation vector of layer $f_{t,n}$	
$oldsymbol{x}_t$	State at stage t	Augmentation of pre-activation vectors of	flayers $\{f_{t,n}:n\in[N]\}$
$\ell_{t,n}$	Cost incurred at stage t for Player n	Weight decay for layer $f_{t,n}$	
ϕ_n	Cost incurred at final stage T for Player n	Lost w.r.t. network output (e.g. cross entre	opy in classification)

Table 8. Dynamic game theoretic terminology w.r.t. different optimality principles.

OLNE	$egin{array}{c} \eta_{t,n}^{\mathrm{O}} \ H_{t,n} \ oldsymbol{p}_{t,n} \end{array}$	Open-loop information structure Optimality objective (Hamiltonian) for OLNE Co-state at stage t for Player n
FNE	$ \begin{vmatrix} \eta_{t,n}^{C} \\ Q_{t,n} \\ V_{t,n} \\ \mathbf{k}_{t,n} \\ \mathbf{K}_{t,n} \end{vmatrix} $	Feedback information structure Optimality objective (Isaacs-Bellman objective) for FNE Value function for FNE Open gain of the locally optimal update for FNE Feedback gain of the locally optimal update for FNE
GR	$ \begin{vmatrix} \eta_{t,n}^{\text{O-CG}} \\ \eta_{t,n}^{\text{C-CG}} \end{vmatrix} $ $ P_t $ $ W_t $ $ \widetilde{\mathbf{k}}_t $ $ \widetilde{\mathbf{K}}_t $	Cooperative open-loop information structure Cooperative feedback information structure Optimality objective (group Bellman objective) for GR Value function for GR Open gain of the locally optimal update for GR Feedback gain of the locally optimal update for GR

B. OCP Characterization of Training Feedforward Networks

The optimality principle to OCP (2), or equivalently the training process of feedforward networks, can be characterized by Dynamic Programming (DP) or Pontryagin Principle (PP). We synthesize the related results below.

Theorem 7 (Bellman (1954); Pontryagin et al. (1962)).

(DP) Define a value function V_t computed recursively by the Bellman equation (17), starting from $V_T(z_T) = \phi(z_T)$,

$$V_t(\boldsymbol{z}_t) = \min_{\boldsymbol{\pi}_t} \ Q_t(\boldsymbol{z}_t, \boldsymbol{\theta}_t), \quad \text{where} \quad Q_t(\boldsymbol{z}_t, \boldsymbol{\theta}_t) \coloneqq \ell_t(\boldsymbol{\theta}_t) + V_{t+1}(f_t(\boldsymbol{z}_t, \boldsymbol{\theta}_t))$$

$$\tag{17}$$

is called the Bellman objective. $\pi_t \equiv \theta_t(z_t)$ is an arbitrary mapping from z_t to θ_t . Let π_t^* be the minimizer to (17), then $\{\pi_t^* : t \in [T]\}$ is the optimal feedback policy to (2).

(PP) The optimal trajectory $\pi_t^* \equiv \theta_t^*(z_t^*)$ along (17) obeys

$$\mathbf{z}_{t+1}^* = \nabla_{\mathbf{p}_{t+1}} H_t \left(\mathbf{z}_t^*, \mathbf{p}_{t+1}^*, \theta_t^* \right), \quad \mathbf{z}_0^* = \mathbf{z}_0,$$
 (18a)

$$\boldsymbol{p}_{t}^{*} = \nabla_{\boldsymbol{z}_{t}} H_{t} \left(\boldsymbol{z}_{t}^{*}, \boldsymbol{p}_{t+1}^{*}, \boldsymbol{\theta}_{t}^{*} \right), \qquad \boldsymbol{p}_{T}^{*} = \nabla_{\boldsymbol{z}_{T}} \phi \left(\boldsymbol{z}_{T}^{*} \right), \tag{18b}$$

$$\theta_t^* = \arg\min_{\theta_t} H_t \left(\mathbf{z}_t^*, \mathbf{p}_{t+1}^*, \theta_t \right), \tag{18c}$$

where (18b) is the adjoint equation for the co-state p_t^* and

$$H_t(\boldsymbol{z}_t, \boldsymbol{p}_{t+1}, \theta_t) \coloneqq \ell_t(\theta_t) + \boldsymbol{p}_{t+1}^\mathsf{T} f_t(\boldsymbol{z}_t, \theta_t)$$

is the discrete-time Hamiltonian.

Theorem 7 provides an OCP characterization of training feedforward networks. First, notice that the time-varying OCP objectives (Q_t, H_t) are constructed through some backward processes similar to the Back-propagation (BP). Indeed, one can verify that the adjoint equation (18b) gives the exact BP dynamics. Similarly, the dynamics of V_t in (17) also relate to BP under some conditions (Liu et al., 2021). The parameter update, $\theta_t \leftarrow \theta_t - \delta\theta_t$, for standard training methods can be seen as solving the discrete-time Hamiltonian H_t with different precondition matrices (Li et al., 2017a). On the other hand, DDPNOpt (Liu et al., 2021) minimizes the time-dependent Bellman objective Q_t with $\theta_t \leftarrow \theta_t - \delta\pi_t$. This elegant connection is, however, limited to the interpretation between feedforward networks and Markovian dynamical systems (1).

C. Missing Derivations in Section 3

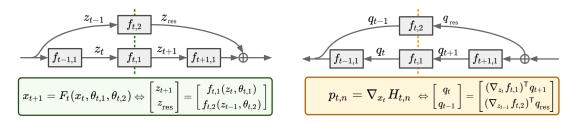


Figure 12. Forward propagation (left) and Back-propagation (right) of a residual block and how each quantity connects to OLNE optimality.

Proof of Proposition 2. Expand the expression of the Hamiltonian in OLNE:

$$H_{t,n}(\boldsymbol{x}_t, \boldsymbol{p}_{t+1,n}, \boldsymbol{\theta}_{t,1}, \cdots, \boldsymbol{\theta}_{t,N}) \coloneqq \ell_{t,n}(\boldsymbol{\theta}_{t,1}, \cdots, \boldsymbol{\theta}_{t,N}) + F_t(\boldsymbol{x}_t, \boldsymbol{\theta}_{t,1}, \cdots, \boldsymbol{\theta}_{t,N})^\mathsf{T} \boldsymbol{p}_{t+1,n},$$

where $p_{t,n}$ is the co-state whose dynamics obey

$$p_{t,n} = \nabla_{x_t} H_{t,n}, \quad p_{T,n} = \nabla_{x_T} \phi_n(x_T).$$

Recall §3.1 where we demonstrate that for training generic DNNs, one shall consider $\ell_{t,n} := \ell_{t,n}(\theta_{t,n})$ and $\phi_n := \phi$. Hence, the dynamics of $p_{t,n}$ become

$$\mathbf{p}_{t,n} = \nabla_{\mathbf{x}_t} H_{t,n}, \quad \mathbf{p}_{T,n} = \nabla_{\mathbf{x}_T} \phi(\mathbf{x}_T), \quad \text{where } H_{t,n} = \ell_{t,n} (\theta_{t,n}) + F_t(\mathbf{x}_t, \theta_{t,1}, \cdots, \theta_{t,N})^\mathsf{T} \mathbf{p}_{t+1,n}.$$
 (19)

Our goal is to show that (19) gives the exact Back-propagation dynamics. First, notice that the terminal condition of (19), i.e. $p_{T,n} = \nabla_{x_T} \phi$, is already the gradient w.r.t. the network output without any manipulation. Next, to show that $p_{t,n}$ corresponds to the Back-propagation at stage t, consider, for instance, the computation graphs of the residual block in Fig. 12, where we replot Fig. 2 together with its Back-propagation dynamic and denote q as the gradient w.r.t. the activation vector z. Then, it can be shown by induction that $p_{t,n}$ augments all "q"s aligned at stage t. Indeed, suppose $p_{t+1,n}$ is the augmentation of the Back-propagation gradients at stage t+1, i.e. $p_{t+1,n} := [q_{t+1}, q_{res}]^T$, then the co-state at the current stage t can be expanded as

$$\boldsymbol{p}_{t,n} = \nabla_{\boldsymbol{x}_t} H_{t,n} = \nabla_{\boldsymbol{x}_t} F_t^\mathsf{T} \boldsymbol{p}_{t+1,n} = \left[\begin{array}{cc} \nabla_{\boldsymbol{z}_t} f_{t,1} & \nabla_{\boldsymbol{z}_{t-1}} f_{t,1} \\ \nabla_{\boldsymbol{z}_t} f_{t,2} & \nabla_{\boldsymbol{z}_{t-1}} f_{t,2} \end{array} \right]^\mathsf{T} \left[\begin{array}{c} \boldsymbol{q}_{t+1} \\ \boldsymbol{q}_{\text{res}} \end{array} \right] = \left[\begin{array}{c} \nabla_{\boldsymbol{z}_t} f_{t,1}^\mathsf{T} \boldsymbol{q}_{t+1} \\ \nabla_{\boldsymbol{z}_{t-1}} f_{t,2}^\mathsf{T} \boldsymbol{q}_{\text{res}} \end{array} \right] = \left[\begin{array}{c} \boldsymbol{q}_t \\ \boldsymbol{q}_{t-1} \end{array} \right],$$

which augments all "q"s at stage t. Once we connect $p_{t,n}$ to the Back-propagation dynamics, it can be verified that

$$H_{\theta}^{t,n} \equiv \nabla_{\theta_{t,n}} H_{t,n} = \nabla_{\theta_{t,n}} \ell_{t,n} + \nabla_{\theta_{t,n}} F_t^{\mathsf{T}} \boldsymbol{p}_{t+1,n}.$$

is indeed the gradient w.r.t. the parameter $\theta_{t,n}$ of each layer $f_{t,n}$. Therefore, taking the iterative update $\theta_{t,n} \leftarrow \theta_{t,n} - H_{\theta}^{t,n}$ is equivalent to descending along the SGD direction, up to a learning rate scaling. Similarly, setting different precondition matrices M will recover other standard methods. Hence, we conclude the proof.

Optimality principle for $\eta_{t,n}^{\text{O-CG}}$. For the completeness, below we provide the optimality principle for the cooperative open-loop information structure $\eta_{t,n}^{\text{O-CG}}$.

Definition 8 (Cooperative optimality principle by $\eta_{t,n}^{\text{O-CG}}$). A set of strategy, $\{\theta_{t,n}^* : \forall t, n\}$, provides an open-loop optimal solution to the joint optimization (4) if

$$\begin{split} \boldsymbol{\theta}_{t,1}^*, \cdots, \boldsymbol{\theta}_{t,N}^* &= \mathop{\arg\min}_{\boldsymbol{\theta}_{t,n}: n \in [N]} \ \widetilde{H}_t(\boldsymbol{x}_t, \widetilde{\boldsymbol{p}}_{t+1}, \boldsymbol{\theta}_{t,1}, \cdots, \boldsymbol{\theta}_{t,N}), \\ \text{where} \quad \boldsymbol{\theta}_{t,n}^* &\equiv \boldsymbol{\theta}_{t,n}^*(\boldsymbol{\eta}_{t,n}^{\text{O-CG}}) \quad \text{and} \quad \widetilde{H}_t \coloneqq \sum_{n=1}^N \ell_{t,n} + F_t^\mathsf{T} \widetilde{\boldsymbol{p}}_{t+1} \end{split}$$

is the "group" Hamiltonian at stage t. Similar to OLNE, the joint co-state \widetilde{p}_t can be simulated by

$$\widetilde{\boldsymbol{p}}_t = \nabla_{\boldsymbol{x}_t} \widetilde{H}_t, \quad \widetilde{\boldsymbol{p}}_T = \sum_{n=1}^N \nabla_{\boldsymbol{x}_T} \phi_n.$$

In this work, we focus on solving the optimality principle inherited in $\eta_{t,n}^{\text{C-CG}}$ as a representative of the CG optimality. Since $\eta_{t,n}^{\text{C-CG}} \subset \eta_{t,n}^{\text{C-CG}}$, the latter captures richer information and tends to perform better in practice, as evidenced by Fig. 3.

D. Missing Derivations in Section 4

D.1. Complete Derivation of the Iterative Updates

Derivation of FNE update. Our goal is to approximately solve the Isaacs-Bellman recursion (6) only up to second-order. Recall that the second-order expansion of $Q_{t,n}$ at some fixed point $(\boldsymbol{x}_t, \theta_{t,n})$ takes the form

$$Q_{t,n} \approx \frac{1}{2} \begin{bmatrix} \mathbf{1} \\ \delta \mathbf{x}_{t} \\ \delta \theta_{t,n} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} Q_{t,n} & Q_{x}^{t,n^{\mathsf{T}}} & Q_{\theta}^{t,n^{\mathsf{T}}} \\ Q_{x}^{t,n} & Q_{tx}^{t,n} & Q_{\thetax}^{t,n^{\mathsf{T}}} \\ Q_{\theta}^{t,n} & Q_{\thetax}^{t,n} & Q_{\theta\theta}^{t,n^{\mathsf{T}}} \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \delta \mathbf{x}_{t} \\ \delta \theta_{t,n} \end{bmatrix}, \quad \text{where } Q_{\theta}^{t,n} = F_{\theta}^{t^{\mathsf{T}}} V_{xx}^{t+1,n} + \ell_{\theta}^{t,n} \\ Q_{\theta}^{t,n} = F_{\theta}^{t^{\mathsf{T}}} V_{xx}^{t+1,n} + \ell_{\theta\theta}^{t,n} \\ Q_{\theta}^{t,n} = F_{\theta}$$

follow standard chain rule (recall $Q_{t,n} := \ell_{t,n} + V_{t+1,n} \circ F_t$) with the linearized dynamics $F_{\theta}^t \equiv \nabla_{\theta_{t,n}} F_t$ and $F_{\boldsymbol{x}}^t \equiv \nabla_{\boldsymbol{x}_t} F_t$. The expansion (20) is a standard quadratic programming, and its analytic solution is given by

$$-\delta \pi_{t,n}^* \equiv -\delta \theta_{t,n}^*(\delta \boldsymbol{x}_t) = -(Q_{\theta\theta}^{t,n})^{\dagger} (Q_{\theta}^{t,n} + Q_{\theta\boldsymbol{x}}^{t,n} \delta \boldsymbol{x}_t) =: -(\mathbf{k}_{t,n} + \mathbf{K}_{t,n} \delta \boldsymbol{x}_t).$$

Substituting this solution back to the Isaacs-Bellman recursion gives us the local expression of $V_{t,n}$,

$$V_{t,n} \approx Q_{t,n} - \frac{1}{2} (Q_{\theta}^{t,n})^{\mathsf{T}} (Q_{\theta\theta}^{t,n})^{\dagger} Q_{\theta}^{t,n}. \tag{21}$$

Therefore, the local derivatives of $V_{t,n}$ can be computed by

$$\begin{aligned} V_{\boldsymbol{x}}^{t,n} &= Q_{\boldsymbol{x}}^{t,n} - Q_{\boldsymbol{x}\boldsymbol{\theta}}^{t,n}(Q_{\boldsymbol{\theta}\boldsymbol{\theta}}^{t,n})^{\dagger}Q_{\boldsymbol{\theta}}^{t,n} = Q_{\boldsymbol{x}}^{t,n} - Q_{\boldsymbol{x}\boldsymbol{\theta}}^{t,n}\mathbf{k}_{t,n} \\ V_{\boldsymbol{x}\boldsymbol{x}}^{t,n} &= Q_{\boldsymbol{x}\boldsymbol{x}}^{t,n} - Q_{\boldsymbol{x}\boldsymbol{\theta}}^{t,n}(Q_{\boldsymbol{\theta}\boldsymbol{\theta}}^{t,n})^{\dagger}Q_{\boldsymbol{\theta}\boldsymbol{x}}^{t,n} = Q_{\boldsymbol{x}\boldsymbol{x}}^{t,n} - Q_{\boldsymbol{x}\boldsymbol{\theta}}^{t,n}\mathbf{K}_{t,n}. \end{aligned}$$

Derivation of GR update. We will adopt the same terminology $u \equiv \theta_{t,1}, v \equiv \theta_{t,2}$. Following the procedure as in the FNE case, we can perform the second-order expansion of P_t at some fixed point (x_t, u, v) . The analytic solution to the corresponding quadratic programming is given by

$$-\begin{bmatrix} \delta \pi_{t,1}^* \\ \delta \pi_{t,2}^* \end{bmatrix} = -\begin{bmatrix} P_{uu}^t & P_{uv}^t \\ P_{vu}^t & P_{vv}^t \end{bmatrix}^{\dagger} \begin{pmatrix} P_{u}^t \\ P_{v}^t \end{bmatrix} + \begin{bmatrix} P_{ux}^t \\ P_{vx}^t \end{bmatrix} \delta x_t , \qquad (22)$$

where the block-matrices inversion can be expanded using the Schur complement.

$$\begin{bmatrix} P_{uu}^{t} & P_{uv}^{t} \\ P_{vu}^{t} & P_{vv}^{t} \end{bmatrix}^{\dagger} = \begin{bmatrix} \overbrace{(P_{uu}^{t} - P_{uv}^{t} P_{vv}^{t} P_{vu}^{t})}^{\widetilde{P}_{vv}^{t} + P_{vv}^{t}})^{\dagger} & -\widetilde{P}_{uu}^{t} P_{uv}^{t} P_{vv}^{t} \\ -\widetilde{P}_{vv}^{t} P_{vu}^{t} P_{uu}^{t} & (\underbrace{(P_{vv}^{t} - P_{vu}^{t} P_{uu}^{t} P_{uv}^{t})}^{-\widetilde{P}_{uu}^{t} + P_{uv}^{t}})^{\dagger} \end{bmatrix} . \tag{23}$$

Hence, (22) becomes

$$\begin{bmatrix} \delta \pi_{t,1}^* \\ \delta \pi_{t,2}^* \end{bmatrix} = \begin{bmatrix} \widetilde{P}_{uu}^t (P_u^t - P_{uv}^t P_{vv}^t P_v^t) \\ \widetilde{P}_{vv}^t (P_v^t - P_{vu}^t P_{uv}^t P_u^t) \end{bmatrix} + \begin{bmatrix} \widetilde{P}_{uu}^t (P_{ux}^t - P_{uv}^t P_{vv}^t P_{vx}^t) \\ \widetilde{P}_{vv}^t (P_{vx}^t - P_{vu}^t P_{uu}^t P_{ux}^t) \end{bmatrix} \delta x_t$$

$$= \begin{bmatrix} \widetilde{P}_{uu}^t (P_u^t - P_{uv}^t \mathbf{I}_t) \\ \widetilde{P}_{vv}^t (P_v^t - P_{vu}^t \mathbf{k}_t) \end{bmatrix} + \begin{bmatrix} \widetilde{P}_{uu}^t (P_{ux}^t - P_{uv}^t \mathbf{L}_t) \\ \widetilde{P}_{vv}^t (P_{vx}^t - P_{vu}^t \mathbf{K}_t) \end{bmatrix} \delta x_t$$

$$=: \begin{bmatrix} \widetilde{\mathbf{k}}_t \\ \widetilde{\mathbf{I}}_t \end{bmatrix} + \begin{bmatrix} \widetilde{\mathbf{K}}_t \\ \widetilde{\mathbf{L}}_t \end{bmatrix} \delta x_t,$$

where we denote the non-cooperative iterative update for Player 1 and 2 respectively by

$$\delta \boldsymbol{u}_t(\delta \boldsymbol{x}_t) = \mathbf{k}_t + \mathbf{K}_t \delta \boldsymbol{x}_t, \quad \text{where } \mathbf{k}_t \coloneqq P_{\boldsymbol{u}\boldsymbol{u}}^{t\ \dagger} P_{\boldsymbol{u}}^t \quad \text{and} \quad \mathbf{K}_t \coloneqq P_{\boldsymbol{u}\boldsymbol{u}}^{t\ \dagger} P_{\boldsymbol{u}\boldsymbol{x}}^t,$$
$$\delta \boldsymbol{v}_t(\delta \boldsymbol{x}_t) = \mathbf{I}_t + \mathbf{L}_t \delta \boldsymbol{x}_t, \quad \text{where } \mathbf{I}_t \coloneqq P_{\boldsymbol{v}\boldsymbol{v}}^{t\ \dagger} P_{\boldsymbol{v}}^t \quad \text{and} \quad \mathbf{L}_t \coloneqq P_{\boldsymbol{v}\boldsymbol{v}}^{t\ \dagger} P_{\boldsymbol{v}\boldsymbol{x}}^t.$$

Substituting this solution back to the GR Bellman equation gives the local expression of W_t ,

$$W_t \approx P_t - \frac{1}{2} \begin{bmatrix} P_u^t \\ P_v^t \end{bmatrix}^\mathsf{T} \begin{bmatrix} P_{uu}^t & P_{uv}^t \\ P_{vu}^t & P_{vv}^t \end{bmatrix}^\dagger \begin{bmatrix} P_u^t \\ P_v^t \end{bmatrix}. \tag{24}$$

Finally, taking the derivatives yields the formula for updating the derivatives of W_t ,

$$W_{\boldsymbol{x}}^{t} = P_{\boldsymbol{x}}^{t} - \frac{1}{2} \left(P_{\boldsymbol{x}\boldsymbol{u}}^{t} \widetilde{\mathbf{k}}_{t} + P_{\boldsymbol{x}\boldsymbol{v}}^{t} \widetilde{\mathbf{I}}_{t} + \widetilde{\mathbf{K}}_{t}^{\mathsf{T}} P_{\boldsymbol{u}}^{t} + \widetilde{\mathbf{L}}_{t}^{\mathsf{T}} P_{\boldsymbol{v}}^{t} \right) \quad \text{and} \quad W_{\boldsymbol{x}\boldsymbol{x}}^{t} = P_{\boldsymbol{x}\boldsymbol{x}}^{t,n} - P_{\boldsymbol{x}\boldsymbol{u}}^{t} \widetilde{\mathbf{K}}_{t} - P_{\boldsymbol{x}\boldsymbol{v}}^{t} \widetilde{\mathbf{L}}_{t}, \tag{25}$$

which is much complex than (10).

D.2. Kronecker Factorization and Proof of Proposition 5

We first provide some backgrounds for the Kronecker factorization (KFAC; Martens & Grosse (2015)). KFAC relies on the fact that for an affine mapping layer, i.e. $z_{t+1} = f_t(z_t, \theta_t) := W_t z_t + b_t$, $\theta_t := \text{vec}([W_t, b_t])$, the gradient of the training objective L w.r.t. the parameter θ_t admits a compact factorization,

$$\nabla_{\theta_t} L = \nabla_{\theta_t} f_t^\mathsf{T} \nabla_{\boldsymbol{z}_{t+1}} L = \boldsymbol{z}_t \otimes \nabla_{\boldsymbol{z}_{t+1}} L,$$

where \otimes denotes the Kronecker product. With this, the layer-wise Fisher information matrix, or equivalently the Gauss-Newton (GN) matrix, for classification can be approximated with

$$\mathbb{E}[\nabla_{\theta_t} L \nabla_{\theta_t} L^\mathsf{T}] = \mathbb{E}[(\boldsymbol{z}_t \otimes \nabla_{\boldsymbol{z}_{t+1}} L) (\boldsymbol{z}_t \otimes \nabla_{\boldsymbol{z}_{t+1}} L)^\mathsf{T}] \approx \mathbb{E}[\boldsymbol{z}_t \boldsymbol{z}_t^\mathsf{T}] \otimes \mathbb{E}[\nabla_{\boldsymbol{z}_{t+1}} L \nabla_{\boldsymbol{z}_{t+1}} L^\mathsf{T}].$$

We can adopt this factorization to our setup by first recalling from our proof of Proposition 2 (see Appendix C) that $(\nabla_{\theta_t} L, \nabla_{\mathbf{z}_{t+1}} L)$ are interchangeable with $(H_{\theta}^t, \mathbf{p}_{t+1})$, or equivalently $(H_{\theta}^t, H_{\mathbf{z}}^{t+1})$. Hence, the GN approximation of $\mathbb{E}[H_{\theta\theta}^t]$ can be factorized by

$$\mathbb{E}[H_{\theta}^{t}H_{\theta}^{t^{\mathsf{T}}}] \approx \mathbb{E}[z_{t}z_{t}^{\mathsf{T}}] \otimes \mathbb{E}[p_{t+1}p_{t+1}^{\mathsf{T}}] = \mathbb{E}[z_{t}z_{t}^{\mathsf{T}}] \otimes \mathbb{E}[H_{z}^{t+1}H_{z}^{t+1}]. \tag{26}$$

Equation (26) suggests that KFAC factorizes the parameter curvature with two smaller matrices using the activation state z_t and the derivative of *some optimality* (in this case the Hamiltonian H) w.r.t. z_{t+1} . The main advantage of this factorization is to exploit the following formula,

$$(\mathbf{A} \otimes \mathbf{B})^{\dagger} \operatorname{vec}(\mathbf{W}) = (\mathbf{A}^{\dagger} \otimes \mathbf{B}^{\dagger}) \operatorname{vec}(\mathbf{W}) = \operatorname{vec}(\mathbf{B}^{\dagger} \mathbf{W} \mathbf{A}^{-\mathsf{T}}), \tag{27}$$

which allows one to efficiently inverse the parameter curvature with two smaller matrices.

Now, let us proceed to the proof of Proposition 5. First notice that for the shared dynamics considered in Fig. 2, we have

$$F_t(oldsymbol{x}_t,oldsymbol{u},oldsymbol{v})\coloneqq \left[egin{array}{cc} f_{t,1}(oldsymbol{z}_1,oldsymbol{u}) \\ f_{t,2}(oldsymbol{z}_2,oldsymbol{v}) \end{array}
ight] = \left[egin{array}{cc} f_{t,1}(\cdot,oldsymbol{u}) & oldsymbol{0} \\ oldsymbol{0} & f_{t,2}(\cdot,oldsymbol{v}) \end{array}
ight] \left[egin{array}{cc} oldsymbol{z}_1 \\ oldsymbol{z}_2 \end{array}
ight],$$

which resembles the affine mapping concerned by KFAC. This motivates the following approximation,

$$\mathbb{E}[P_{\theta}^{t} P_{\theta}^{t^{\mathsf{T}}}] \approx \mathbb{E}[\boldsymbol{x}_{t} \boldsymbol{x}_{t}^{\mathsf{T}}] \otimes \mathbb{E}[W_{\boldsymbol{x}}^{t+1} W_{\boldsymbol{x}}^{t+1}]. \tag{28}$$

Similar to (26), this approximation (28) factorizes the GN matrix with the MPDG state x_t and the derivative of an optimality (in this case it becomes the GR value function W_{t+1}) w.r.t. x_{t+1} .

If we denote the derivatives w.r.t. the outputs of $f_{t,1}$ and $f_{t,2}$ by g_1 and g_2 , i.e. $W_x^{t+1} := [g_1, g_2]^\mathsf{T}$, and rewrite $x_t := [z_1, z_2]^\mathsf{T}$, then (28) can be expanded by

$$\mathbb{E}[\boldsymbol{x}_{t}\boldsymbol{x}_{t}^{\mathsf{T}}] = \begin{bmatrix} \mathbb{E}[\boldsymbol{z}_{1}\boldsymbol{z}_{1}^{\mathsf{T}}] & \mathbb{E}[\boldsymbol{z}_{1}\boldsymbol{z}_{2}^{\mathsf{T}}] \\ \mathbb{E}[\boldsymbol{z}_{2}\boldsymbol{z}_{1}^{\mathsf{T}}] & \mathbb{E}[\boldsymbol{z}_{2}\boldsymbol{z}_{2}^{\mathsf{T}}] \end{bmatrix} =: \begin{bmatrix} A_{\boldsymbol{u}\boldsymbol{u}} & A_{\boldsymbol{u}\boldsymbol{v}} \\ A_{\boldsymbol{v}\boldsymbol{u}} & A_{\boldsymbol{v}\boldsymbol{v}} \end{bmatrix}$$
$$\mathbb{E}[W_{\boldsymbol{x}}^{t+1}W_{\boldsymbol{x}}^{t+1}] = \begin{bmatrix} \mathbb{E}[\boldsymbol{g}_{1}\boldsymbol{g}_{1}^{\mathsf{T}}] & \mathbb{E}[\boldsymbol{g}_{1}\boldsymbol{g}_{2}^{\mathsf{T}}] \\ \mathbb{E}[\boldsymbol{g}_{2}\boldsymbol{g}_{1}^{\mathsf{T}}] & \mathbb{E}[\boldsymbol{g}_{2}\boldsymbol{g}_{2}^{\mathsf{T}}] \end{bmatrix} =: \begin{bmatrix} B_{\boldsymbol{u}\boldsymbol{u}} & B_{\boldsymbol{u}\boldsymbol{v}} \\ B_{\boldsymbol{v}\boldsymbol{u}} & B_{\boldsymbol{v}\boldsymbol{v}} \end{bmatrix}.$$

Their inverse matrices are given by the Schur component.

$$\begin{bmatrix} A_{uu} & A_{uv} \\ A_{vu} & A_{vv} \end{bmatrix}^{\dagger} = \begin{bmatrix} \widetilde{A}_{uu}^{\dagger} & -\widetilde{A}_{uu}^{\dagger} A_{uv} A_{vv}^{\dagger} \\ -\widetilde{A}_{vv}^{\dagger} A_{vu} A_{uu}^{\dagger} & \widetilde{A}_{vv}^{\dagger} \end{bmatrix}, \quad \text{where } \begin{cases} \widetilde{A}_{uu} \coloneqq A_{uu} - A_{uv} A_{vv}^{\dagger} A_{uv}^{\dagger} \\ \widetilde{A}_{vv} \coloneqq A_{vv} - A_{vu} A_{uu}^{\dagger} A_{vu}^{\dagger} \end{bmatrix}$$
$$\begin{bmatrix} B_{uu} & B_{uv} \\ B_{vu} & B_{vv} \end{bmatrix}^{\dagger} = \begin{bmatrix} \widetilde{B}_{uu}^{\dagger} & -\widetilde{B}_{uu}^{\dagger} B_{uv} B_{vv}^{\dagger} \\ -\widetilde{B}_{vv}^{\dagger} B_{vu} B_{uu}^{\dagger} & \widetilde{B}_{vv}^{\dagger} \end{bmatrix}, \quad \text{where } \begin{cases} \widetilde{B}_{uu} \coloneqq B_{uu} - B_{uv} B_{vv}^{\dagger} B_{uv}^{\dagger} \\ \widetilde{B}_{vv} \coloneqq B_{vv} - B_{vu} B_{uu}^{\dagger} B_{vu}^{\dagger} \end{cases}$$

With all these, the cooperative open gain can be computed with the formula (27),

$$\left(\mathbb{E}[\boldsymbol{x}_{t}\boldsymbol{x}_{t}^{\mathsf{T}}] \otimes \mathbb{E}[\boldsymbol{W}_{\boldsymbol{x}}^{t+1}\boldsymbol{W}_{\boldsymbol{x}}^{t+1}^{\mathsf{T}}]\right)^{\dagger} \operatorname{vec}\left(\begin{bmatrix} P_{\boldsymbol{u}}^{t} & \mathbf{0} \\ \mathbf{0} & P_{\boldsymbol{v}}^{t} \end{bmatrix}\right) = \operatorname{vec}\left(\begin{bmatrix} B_{\boldsymbol{u}\boldsymbol{u}} & B_{\boldsymbol{u}\boldsymbol{v}} \\ B_{\boldsymbol{v}\boldsymbol{u}} & B_{\boldsymbol{v}\boldsymbol{v}} \end{bmatrix}^{\dagger} \begin{bmatrix} P_{\boldsymbol{u}}^{t} & \mathbf{0} \\ \mathbf{0} & P_{\boldsymbol{v}}^{t} \end{bmatrix} \begin{bmatrix} A_{\boldsymbol{u}\boldsymbol{u}} & A_{\boldsymbol{u}\boldsymbol{v}} \\ A_{\boldsymbol{v}\boldsymbol{u}} & A_{\boldsymbol{v}\boldsymbol{v}} \end{bmatrix}^{-\mathsf{T}}\right). \tag{30}$$

Substituting (29) into (30), after some algebra we will arrive at the KFAC of the cooperative open gain suggested in (15).

$$\widetilde{\mathbf{k}}_{t} \approx \operatorname{vec}(\widetilde{B}_{uu}^{\dagger} P_{u}^{t} \widetilde{A}_{uu}^{-\mathsf{T}} + \widetilde{B}_{uu}^{\dagger} B_{uv} B_{vv}^{\dagger} P_{v}^{t} (\widetilde{A}_{uu}^{\dagger} A_{uv} A_{vv}^{\dagger})^{\mathsf{T}}) \\
= \operatorname{vec}(\widetilde{B}_{uu}^{\dagger} (P_{u}^{t} + B_{uv} B_{vv}^{\dagger} P_{v}^{t} A_{vv}^{-\mathsf{T}} A_{uv}^{\mathsf{T}}) \widetilde{A}_{uu}^{-\mathsf{T}}) \\
= \operatorname{vec}(\widetilde{B}_{uu}^{\dagger} (P_{u}^{t} + B_{uv} \operatorname{vec}^{\dagger} (\mathbf{I}_{t}) A_{uv}^{\mathsf{T}}) \widetilde{A}_{uu}^{-\mathsf{T}}), \tag{31}$$

where the last equality follows by another KFAC approximation $\mathbf{I}_t \approx (A_{\boldsymbol{v}\boldsymbol{v}} \otimes B_{\boldsymbol{v}\boldsymbol{v}})^\dagger \text{vec}(P_{\boldsymbol{v}}^t) = \text{vec}(B_{\boldsymbol{v}\boldsymbol{v}}^\dagger P_{\boldsymbol{v}}^t A_{\boldsymbol{v}\boldsymbol{v}}^{-\mathsf{T}})$. Finally, recalling the expression, $\widetilde{\mathbf{k}}_t \coloneqq \widetilde{P}_{\boldsymbol{u}\boldsymbol{u}}^t (P_{\boldsymbol{u}}^t - P_{\boldsymbol{u}\boldsymbol{v}}^t \mathbf{I}_t)$, from (11) and rewriting (31) by

$$\begin{split} \widetilde{\mathbf{k}}_t \approx & \mathrm{vec}(\widetilde{B}_{\boldsymbol{u}\boldsymbol{u}}^{\dagger}(P_{\boldsymbol{u}}^t + B_{\boldsymbol{u}\boldsymbol{v}}\mathrm{vec}^{\dagger}(\mathbf{I}_t)A_{\boldsymbol{u}\boldsymbol{v}}^{\mathsf{T}})\widetilde{A}_{\boldsymbol{u}\boldsymbol{u}}^{-\mathsf{T}}) \\ = & (\widetilde{A}_{\boldsymbol{u}\boldsymbol{u}} \otimes \widetilde{B}_{\boldsymbol{u}\boldsymbol{u}})^{\dagger}\mathrm{vec}(P_{\boldsymbol{u}}^t + B_{\boldsymbol{u}\boldsymbol{v}}\mathrm{vec}^{\dagger}(\mathbf{I}_t)A_{\boldsymbol{u}\boldsymbol{v}}^{\mathsf{T}}) \end{split}$$

imply the KFAC representation $\widetilde{P}_{uu}^t \approx \widetilde{A}_{uu} \otimes \widetilde{B}_{uu}$ in (14). Hence we conclude the proof.

D.3. Proof of Theorem 6

We first show that setting $P_{uv}^t := \mathbf{0}$ in the update (11) yields (9). To begin, observe that when P_{uv}^t vanishes, the cooperative gains $(\widetilde{\mathbf{k}}_t, \widetilde{\mathbf{K}}_t)$ appearing in (11) degenerate to $\widetilde{\mathbf{k}}_t = P_{uu}^t P_u^t$ and $\widetilde{\mathbf{K}}_t = P_{uu}^t P_{ux}^t$. Therefore, it is sufficient to prove the following result.⁴

Lemma 9. Suppose $Q_{t,n}$ in (6) and P_t in (7) are expanded up to second-order along the same local trajectory $\{(\boldsymbol{x}_t, \theta_{t,n}, \cdots, \theta_{t,N}) : \forall t \in [T]\}$, then we will have the following relations when $P_{\boldsymbol{u}\boldsymbol{v}}^t \coloneqq \boldsymbol{0}$ at all stages.

$$\forall t, \quad Q_{\theta}^{t,1} = P_{\mathbf{u}}^{t}, \quad Q_{\theta\theta}^{t,1} = P_{\mathbf{u}\mathbf{u}}^{t}, \quad Q_{\theta\mathbf{x}}^{t,1} = P_{\mathbf{u}\mathbf{x}}^{t}, \quad Q_{\theta}^{t,2} = P_{\mathbf{v}}^{t}, \quad Q_{\theta\theta}^{t,2} = P_{\mathbf{v}\mathbf{v}}^{t}, \quad Q_{\theta\mathbf{x}}^{t,2} = P_{\mathbf{v}\mathbf{x}}^{t}, \quad (32)$$

where $(u_t, v_t) \equiv (\theta_{t,1}, \theta_{t,2})$ denotes the actions for Player 1 and 2. Furthermore, we have

$$\forall t, \quad W_t = \sum_{n=1}^N V_{t,n}. \tag{33}$$

Proof. We will proceed the proof by induction. At the terminal stage T-1, we have

$$P_{T-1} = \sum_{n=1}^{2} \ell_{T-1,n} + W_T \circ F_{T-1} = \sum_{n=1}^{2} (\ell_{T-1,n} + \phi_n \circ F_{T-1}) = \sum_{n=1}^{2} Q_{T-1,n},$$

since $\phi_n = V_{T,n}$. This implies that when solving the second-order expansion for $\pi_{T-1,1}$ and $\pi_{T-1,2}$, we will have

$$\min_{\pi_{T-1,1},\pi_{T-1,2}} P_{T-1} = \min_{\pi_{T-1,1}} Q_{T-1,1} + \min_{\pi_{T-1,2}} Q_{T-1,2}$$

since the cross-correlation matrix P_{uv}^{T-1} is discarded. Therefore, all equalities in (32) hold at this stage. Furthermore, substituting $P_{uv}^{T-1} := \mathbf{0}$ into (24) yields the following GR value function

$$W_{T-1} = P_{T-1} - \frac{1}{2} \left((P_{\boldsymbol{u}}^{T-1})^{\mathsf{T}} (P_{\boldsymbol{u}\boldsymbol{u}}^{T-1})^{\dagger} P_{\boldsymbol{u}}^{T-1} + (P_{\boldsymbol{v}}^{T-1})^{\mathsf{T}} (P_{\boldsymbol{v}\boldsymbol{v}}^{T-1})^{\dagger} P_{\boldsymbol{v}}^{T-1} \right)$$
$$= \sum_{n=1}^{2} \left(Q_{T-1,n} - \frac{1}{2} (Q_{\theta}^{T-1,n})^{\mathsf{T}} (Q_{\theta\theta}^{T-1,n})^{\dagger} Q_{\theta}^{T-1,n} \right) = \sum_{n=1}^{2} V_{T-1,n}.$$

So (33) also holds. Now, suppose (32, 33) hold at t + 1, then

$$P_t = \sum_{n=1}^{2} \ell_{t,n} + W_{t+1} \circ F_t = \sum_{n=1}^{2} (\ell_{t,n} + V_{t+1,n} \circ F_t) = \sum_{n=1}^{2} Q_{t,n}.$$

Together with $P_{uv}^t := \mathbf{0}$, we can see that all equalities in (32) hold. Furthermore, it implies that

$$W_{t} = P_{t} - \frac{1}{2} \left(P_{\boldsymbol{u}}^{t \mathsf{T}} P_{\boldsymbol{u} \boldsymbol{u}}^{t \dagger} P_{\boldsymbol{u}}^{t} + P_{\boldsymbol{v}}^{t \mathsf{T}} P_{\boldsymbol{v} \boldsymbol{v}}^{t \dagger} P_{\boldsymbol{v}}^{t} \right) = \sum_{n=1}^{2} \left(Q_{t,n} - \frac{1}{2} (Q_{\theta}^{t,n})^{\mathsf{T}} (Q_{\theta\theta}^{t,n})^{\dagger} Q_{\theta}^{t,n} \right) = \sum_{n=1}^{2} V_{t,n}.$$

Hence we conclude the proof.

Next, we proceed to the second case, which suggests that running (9) with $(Q_{\theta x}^{t,n},Q_{\theta \theta}^{t,n})\coloneqq (\mathbf{0},\mathbf{I})$ yields SGD. Since the FNE update in this case degenerates to $\delta\pi_{t,n}^*=Q_{\theta}^{t,n}$, it is sufficient to prove the following lemma.

⁴We consider the two-player setup for simplicity, yet the methodology applies to the multi-player setup.

Lemma 10. Suppose $H_{t,n}$ in (5) and $Q_{t,n}$ in (6) are expanded up to second-order along the same local trajectory $\{(\boldsymbol{x}_t, \theta_{t,n}, \cdots, \theta_{t,N}) : \forall t \in [T]\}$, then we will have the following relations when $(Q_{\theta \boldsymbol{x}}^{t,n}, Q_{\theta \theta}^{t,n}) := (\boldsymbol{0}, \boldsymbol{I})$ for all stages.

$$\forall t, \quad Q_{\theta}^{t,n} = H_{\theta}^{t,n}, \quad V_{\boldsymbol{x}}^{t,n} = \boldsymbol{p}_{t,n}. \tag{34}$$

Proof. Again, we will proceed the proof by induction. First, notice that $V_x^{T,n} = \phi_x^n = p_{T,n}$ no matter whether or not $Q_{\theta x}^{t,n}$ and $Q_{\theta \theta}^{t,n}$ degenerate. At the terminal stage T-1, we have

$$Q_{\theta}^{T-1,n} = \ell_{\theta}^{T-1,n} + (F_{\theta}^{T-1})^{\mathsf{T}} V_{\boldsymbol{x}}^{T,n} = \ell_{\theta}^{T-1,n} + (F_{\theta}^{T-1})^{\mathsf{T}} \boldsymbol{p}_{T,n} = H_{\theta}^{T-1,n}.$$

Also, when $Q_{\theta x}^{T-1,n} := \mathbf{0}$, (10) becomes

$$V_{x}^{T-1,n} = Q_{x}^{T-1,n} = (F_{x}^{T-1})^{\mathsf{T}} V_{x}^{T,n} = (F_{x}^{T-1})^{\mathsf{T}} p_{T,n} = H_{x}^{T-1,n} = p_{T-1,n}.$$

Hence, (34) holds at T-1. Now, suppose these relations hold at t+1, then

$$Q_{\theta}^{t,n} = \ell_{\theta}^{t,n} + (F_{\theta}^{t})^{\mathsf{T}} V_{x}^{t+1,n} = \ell_{\theta}^{t,n} + (F_{\theta}^{t})^{\mathsf{T}} p_{t+1,n} = H_{\theta}^{t,n}$$

and similarly

$$V_{x}^{t,n} = Q_{x}^{t,n} = (F_{x}^{t})^{\mathsf{T}} V_{x}^{t+1,n} = (F_{x}^{t})^{\mathsf{T}} p_{t+1,n} = H_{x}^{t,n} = p_{t,n}.$$

Hence, we conclude the proof.

Finally, the last case follows readily by combining Lemma 9 and 10, so we conclude all proofs.

E. More on the Experiments

All experiments are run with Pytorch on the GPU machines, including GTX 1080 TI, GTX 2070, and TITAN RTX. We preprocessed all datasets with standardization. We also perform data augmentation when training CIFAR100. Below we detail the setup for each experiment.

Classification (Table 2 and 3). For CIFAR10 and CIFAR100, we use standard implementation of ResNet18 from https://pytorch.org/hub/pytorch_vision_resnet/. As for SVHN and MNIST, the residual network consists of 3 residual blocks. The residual block shares a similar architecture in Fig. 2 except with the identity shortcut mapping and without BN. We use 3×3 kernels for all convolution filters. The number of feature maps in the convolution filters is set to 12 and 16 respectively

Table 9. Hyper-parameter search in Table 2

SGD (ning Rate (LR)
Adam & RMSprop (EKFAC (7e-2, 5e-1) 7e-4, 1e-2) 1e-2, 3e-1)

for MNIST and SVHN. Meanwhile, the inception network consists of a convolution layer followed by an inception block (see Fig. 6), another convolution layer, and two fully-connected layers. Regarding the hyper-parameters used in baselines, we select them from an appropriate search space detailed in Table 9. We use the implementation in https://github.com/Thrandis/EKFAC-pytorch for EKFAC and implement our own EMSA in PyTorch since the official code released from Li et al. (2017a) does not support GPU parallelization.

Ablation study (**Fig. 5**) Each grid in Fig. 5 corresponds to a distinct combination of baseline and dataset. Its numerical value reports the performance difference between the following two training processes.

- Accuracy of the baseline run with the best-tuned configuration which we report in Table 2 and 3.
- Accuracy of DGNOpt with its parameter curvature set to the precondition matrix implied by the above best-tuned setup.

For instance, suppose the learning rate of EKFAC on MNIST is best-tuned to 0.01, then we simply set $Q_{\theta\theta}^{t,n} \approx 0.01 \times Q_{\theta\theta}^{t,n} Q_{\theta\theta}^{t,n}$ for all t. From Theorem 6, these two training procedures only differ in the presence of $Q_{\theta x}^{t,n}$, which allows EKFAC to adjust its update based on the change of $x_t \in \eta_{t,n}^{\mathbb{C}}$.

Runtime and memory complexity (Fig. 7). The numerical values are measured on the GTX 2070.

Feedback analysis (Fig. 8). We use the same inception-based network in Table 3.

Remark for EMSA (Footnote 3). Extended Method of Successive Approximations (EMSA) was originally proposed by Li et al. (2017a) as an OCP-inspired method for training *feedforward* networks. It considers the following minimization,

$$\theta_{t}^{*} = \arg\min H_{t}^{\rho} (\mathbf{z}_{t}, \mathbf{z}_{t+1}, \mathbf{p}_{t}, \mathbf{p}_{t+1}, \theta_{t}),$$
where $H_{t}^{\rho} (\mathbf{z}_{t}, \mathbf{z}_{t+1}, \mathbf{p}_{t}, \mathbf{p}_{t+1}, \theta_{t}) := H_{t} (\mathbf{z}_{t}, \mathbf{p}_{t+1}, \theta_{t}) + \frac{1}{2} \rho \|\mathbf{z}_{t+1} - f_{t}(\mathbf{z}_{t}, \theta_{t})\|_{2} + \frac{1}{2} \rho \|\mathbf{p}_{t} - \nabla_{\mathbf{z}_{t}} H_{t}\|_{2}$
(35)

essentially augments the original Hamiltonian H_t with the feasibility constraints on both forward states and backward co-states. EMSA solves the minimization (35) with L-BFGS per layer at each training iteration. In Table 2 and 3, we extend their formula to accept $H_{t,n}$. Due to the feasibility constraints, the resulting modified Hamiltonian $H_{t,n}^{\rho}$ depends additionally on x_{t+1} and $p_{t,n}$; hence being different from the original Hamiltonian $H_{t,n}$. As a result, the ablation analysis using Theorem 6 is not applicable for EMSA.

Cooperative training (Fig. 9, Fig. 11, and Table 4). The network consists of 4 convolutions followed by 2 fully-connected layers, and is activated by ReLU. We use 3×3 kernels with 32 feature maps for all convolutions and set the batch size to 128.

Adaptive alignment with bandit (Fig. 10 and Table 5). We use the same ResNet18 as in classification for CIFAR10, and a smaller residual network with 1 residual block for SVHN. The residual block shares the same architecture as in Fig. 2 except without BN. All convolution layers use 3×3 kernels with 12 feature maps. Again, the batch size is set to 128. Note that in this experiment we use a slightly larger learning rate compared with the one used in Table 2. While DGNOpt achieves better final accuracies for both setups, in practice, the former tends to amplify the stabilization when we enlarge the information structure during training. Hence, it differentiates DGNOpt from other baselines.

Alg. 2 presents the pseudo-code of how DGNOpt can be integrated with any generic bandit-based algorithm (marked as blue). For completeness, we also provide the pseudo-code of EXP3++ in Alg. 3. We refer readers to Seldin & Slivkins (2014) for the definition of $\xi_k(m)$ and η_k (do not confuse with $\eta_{t,n}$ in the main context).

```
Algorithm 2 DGNOpt with Multi-Armed Bandit (MAB)
```

```
Input: dataset \mathcal{D}, network \{f_i(\cdot, \theta_i)\}, number of alignments M
Initialize the multi-armed bandit MAB.init (M)
repeat
   Draw an alignment m \leftarrow \texttt{MAB.sample} ().
   Construct F \equiv \{F_t : t \in [T]\} according to m.
   Compute x_t by propagating x_0 \sim \mathcal{D} through F
   for t = T - 1 to 0 do
      Solve the update \delta \pi_{t,n}^* with (9) or (11)
      Solve (V_x^{t,n}, V_{xx}^{t,n}) or (W_x^t, W_{xx}^t) with (10) or (25)
   end for
   Set \boldsymbol{x}_0' = \boldsymbol{x}_0
  for t = 0 to T-1 do

    □ Update parameter

      Apply \theta_{t,n} \leftarrow \theta_{t,n} - \delta \pi_{t,n}^*(\delta x_t) with \delta x_t = x_t' - x_t
      Compute \boldsymbol{x}_{t+1}' = F_t(\boldsymbol{x}_t', \theta_{t,1}, \cdots, \theta_{t,N})
   end for
   Compute the accuracy r on validation set.
   Run MAB. update (r).
until converges
```

Algorithm 3 EXP3++ (Seldin & Slivkins, 2014)

```
\begin{split} &\forall m, L_k(m) = 0 \\ &\textbf{end function} \end{split} &\textbf{function sample ()} \\ &\forall m, \epsilon_k(m) = \min\{\frac{1}{2M}, \frac{1}{2}\sqrt{\frac{\ln M}{kM}}, \xi_k(m)\} \\ &\forall m, \rho_k(m) = e^{-\eta_k L_k(m)}/\sum_{m'} e^{-\eta_k L_k(m')} \\ &\forall m, \tilde{\rho}_k(m) = (1 - \sum_{m'} \epsilon_k(m'))\rho_k(m) + \epsilon_k(m) \\ &\textbf{Sample action according to } \tilde{\rho}_k(m) \end{split}
```

end function

function init (M)

 $(k,M) \leftarrow (1,M)$

```
 \begin{aligned} & \textbf{function} \text{ update } (r_k^m) \\ & \ell_k^m = (1-r_k^m)/\tilde{\rho}_k(m) \\ & L_{k+1}(m) = L_k(m) + \ell_k^m \\ & k \leftarrow k+1 \end{aligned}   & \textbf{end function}
```

Additional Experiments.

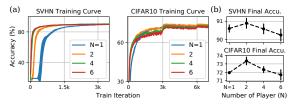


Figure 13. (a) Training curve and (b) final accuracy as we vary the number of player (N) as a hyper-parameter of game-extended EKFAC. Similar to Fig. 9, we also observe that N=2 gives the best final accuracy on both datasets.