



HAL
open science

Depth-Adaptive Neural Networks from the Optimal Control viewpoint

Joubine Aghili, Olga Mula

► **To cite this version:**

Joubine Aghili, Olga Mula. Depth-Adaptive Neural Networks from the Optimal Control viewpoint. 2020. hal-02897466v1

HAL Id: hal-02897466

<https://hal.science/hal-02897466v1>

Preprint submitted on 12 Jul 2020 (v1), last revised 19 Apr 2023 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Depth-Adaptive Neural Networks from the Optimal Control viewpoint

Joubine Aghili¹ and Olga Mula^{*,2}

¹CEREMADE, CNRS, UMR 7534, Université Paris-Dauphine, PSL University, 75016 Paris, France

²Inria Paris, Commedia Team, 2 rue Simone Iff, 75012 Paris, France

Abstract

In recent years, deep learning has been connected with optimal control as a way to define a notion of a continuous underlying learning problem. In this view, neural networks can be interpreted as a discretization of a parametric Ordinary Differential Equation which, in the limit, defines a continuous-depth neural network. The learning task then consists in finding the best ODE parameters for the problem under consideration, and their number increases with the accuracy of the time discretization. Although important steps have been taken to realize the advantages of such continuous formulations, most current learning techniques fix a discretization (i.e. the number of layers is fixed). In this work, we propose an iterative adaptive algorithm where we progressively refine the time discretization (i.e. we increase the number of layers). Provided that certain tolerances are met across the iterations, we prove that the strategy converges to the underlying continuous problem. One salient advantage of such a shallow-to-deep approach is that it helps to benefit in practice from the higher approximation properties of deep networks by mitigating over-parametrization issues. The performance of the approach is illustrated in several numerical examples.

Keywords: Deep Learning; Neural Networks; Continuous-Depth Neural Networks; Optimal Control

1. Introduction

1.1 Context

Neural networks produce structured parametric families of functions that have been studied for at least 70 years (see Hebb (1949); Rosenblatt (1958)). It is however only in the last decade that their popularity has surged. Thanks to the increase of computing power, and the development of easy-to-use computing tools for optimization and automatic differentiation, statistical learning of neural networks has produced state-of-the-art performance in a large variety of machine learning problems, from computer vision Krizhevsky et al. (2012) (e.g. self-driving cars, X-rays diagnosis,...) to natural language processing Wu et al. (2016) (e.g. Google translate, DeepL Translator,...) and reinforcement learning (e.g. superhuman performance at Go Silver et al. (2016)). Despite this great empirical success, neural networks are not entirely well-understood and there is a pressing need to provide:

*. This research was supported by the Emergence Project “Models and Measures” of the Paris City Council, {aghili,mula}@ceremade.dauphine.fr

- solid *mathematical foundations* to understand their approximation power and why they may outperform other classes of functions,
- *algorithms* to find systematically an appropriate architecture for each problem and to train the network in order to deliver in practice the approximation capabilities predicted by the theory, and to guarantee robustness over generalization errors.

On the first point, significant advances have been made recently on the approximation power of neural networks (see, e.g., Yarotsky (2017); Daubechies et al. (2019); Grohs et al. (2019); Gühring et al. (2019) for a selection of rigorous results). In particular, a few recent works have given theoretical evidence on the advantages of using deep versus shallow architectures for the approximation of certain relevant families of functions (see, e.g., Telgarsky (2015); Daubechies et al. (2019)). However, so far the obtained results do not seem to be informative on how to address the second point above, that is, how to build algorithms that discover automatically the right architecture for each problem, and that may allow to benefit from the theoretically high approximation properties of (deep) neural networks. As a result, their training remains a key open issue subject to very active research. The present work is a contribution in this direction, with special focus on understanding the underlying mechanisms of training efficiently deep architectures.

The most commonly applied training methods are based on the stochastic gradient descent (Robbins and Monro (1951); Bottou (2010)) and its variants (see, e.g., Duchi et al. (2011); Zeiler (2012)). It has the advantage of being easily implementable but it is difficult to tune properly in practice, especially in problems involving large data sets and using deep neural networks with many layers and coefficients which are very prone to over-parametrization issues. The difficulty in training deep networks raises the question regarding the benefits of applying a *shallow-to-deep adaptive strategy* in the network architecture (and a hence coarse-to-fine adaptive strategy in the number of coefficients) in order perform the training in a more robust and rapid manner. The underlying intuition is that shallow networks give poor approximation properties but converge faster than deeper networks which, on the contrary, have higher approximation power.

There is a large body of emerging works which explore numerically the potential of training deep neural networks with shallow-to-deep strategies (see, e.g., Chen et al. (2015); Wei et al. (2016); Li and Hoiem (2017); Cortes et al. (2017); Elsken et al. (2018)). In this work, we present an abstract setting and convergence results which are a first theoretical justification of this type of algorithmic approach. In order to understand the potential gain, it is necessary to define a notion of the continuous underlying objects that are approximated. This is the reason why we have chosen to work from a perspective which connects the task of learning with neural networks with optimal control and dynamical ODE systems. In this view, there is a notion of a continuous-depth neural network which is described by a continuous time-dependent ODE system. Its discretization fixes its architecture. For instance, ResNet He et al. (2016) can be regarded as an Explicit Euler scheme for the solution of certain dynamics. The weights involved in the ODE are seen as controls and correspond to the parameters of the network which have to be optimized. The process of learning these parameters can then be recast as an optimal control problem over the admissible controls where the cost function is the empirical risk (sometimes with an extra

regularization term). Necessary optimality conditions are then easily formulated through the Pontryagin Maximum Principle (PMP, Boltyanskii et al. (1960); Pontryagin (2018)).

1.2 Contribution of the present work

The optimal control point of view is appealing since it gives access to a fully continuous description of neural networks and of the underlying statistical learning process. The practical implementation involves two discretization errors:

1. **Sampling error:** The cost function of the fully continuous optimal control problem is the average of a loss function, sometimes also called *risk*. However, in practice, the average loss is replaced by an empirical mean which is built from N samples. So it is necessary to assess by how much the minimum the sampled mean deviates from the continuous underlying problem. We analyze this point in Theorem 6.
2. **Time-integration error:** To solve the optimal control problem on the empirical mean of the loss, one solves the resulting PMP with Newton-type iterations involving the solution of forward and backward dynamical systems with certain time integration schemes (which fix the neural network architecture). The usual approach is to first fix a discretization and then solve the discrete version of the PMP. However, the discretization errors accumulate at each iteration and make us deviate from the time-continuous version. It is thus necessary to examine how much the final output deviates from the time-continuous one.

In this context, the main contributions of the present work are:

1. Derivation of a full error analysis:

- We make a theoretical analysis on the interplay between the above two errors. The impact of the sampling error is analyzed in Theorem 6 and the discretization error in Theorem 4. From these results, it follows that to approximate the minimum of the underlying fully continuous problem at a certain accuracy, we not only need a minimal number of training samples, but we also need to solve the control problem on the empirical risk with some minimal accuracy. This requires, in turn, to use time-adaptive techniques. In our analysis, both errors are additive and can actually be optimally balanced. The final result is given in Corollary 7 and it is formulated in probability.
- To solve the control problem on the empirical risk at a given target accuracy, we show that we have to compute the forward and backward propagations at each step of the Newton-type algorithms at increasing accuracy. As our theoretical proof of Theorem 4 will illustrate, tightening the accuracy of the time integration is necessary to guarantee convergence to the exact continuous problem and justifies the coarse-to-fine, resp. shallow-to-deep, strategy. This yields neural networks with depths that are automatically adapted.

2. Practical coarse-to-fine adaptive algorithm:

- We next use the theory to derive actionable criteria to build time-integration schemes that are adaptively refined across the iterations. This has the advantage of reducing computational cost since early iterations use looser tolerances, thus avoiding unnecessary work due to oversolving, while later iterations use tighter tolerances to deliver accuracy. The coarse-to-fine strategy could also be understood as a dynamic regularization which helps to prevent over-parametrization.
- We illustrate the behavior of the algorithm in several benchmark examples.

1.3 Potential impact, limitations and extensions

We believe the results of this article can contribute to the following key topics in deep learning:

- Automate the selection of neural network architecture during the training,
- Enhance the interpretability of the generalization errors by connecting them with sampling and discretization errors,
- Provide a first theoretical justification of the works adopting the principle of shallow-to-deep training of deep neural networks.

Some limitations and natural questions that arise for future works are the following:

- Our theoretical results involve numerous bounds and some of them might be suboptimal. Certain constants involved are also difficult to estimate in practice. As a result, the theoretical convergence result cannot be used to derive fully explicit adaptive strategies for implementation. However it gives hints on how to proceed in practice, and we propose practical guidelines to perform adaptivity. In future works, it would be desirable to find a sharper analysis with constants that are easier to compute in order to obtain a fully closed pipeline. This may however require very different arguments as our current ones.
- As we will explain in the next section, the optimal control point of view only describes certain classes of neural networks since it introduces certain restrictions in the architecture. As a result, our adaptive setting cannot automate the training of all existing classes of neural networks. So a natural question is how to extend the current ideas to other training techniques that are not based on optimal control.
- One important question is whether the flow map of a given dynamical system can represent the data at hand. The same issue also arises in the classical approach of deep learning in terms of the optimality of the choice of the network. From the optimal control perspective, this representability problem can be viewed as a controllability problem. It would be interesting to transfer controllability results to the current deep learning framework to address this question.

1.4 Outline of the paper

In section 2, we define precisely what we understand by deep learning and we recall its connection with optimal control. Section 3 gives the optimal control setting which we use

in our subsequent developments. In particular, we formulate the fully continuous optimal control problem over the expectation of the loss and its sampled version involving an empirical mean of the loss. Section 4 recalls the Pontryagin’s Maximum Principle for the sampled problem. We also recall in this section the Extended Method of Successive Approximations (E-MSA), originally introduced in Li et al. (2018), to solve the sampled problem. The algorithm is iterative, and it requires to compute at each iteration a forward and a backward time propagation followed by a certain maximization over the control variables. In its original formulation, it is assumed that these steps are performed exactly. In Section 5, we formulate our main algorithm, which is an inexact version of the E-MSA that we call Adaptive MSA (A-MSA). It is based on approximately realizing each propagation and maximization within a certain accuracy. To ensure convergence to the exact solution, the accuracy needs to be tightened across the iterations. In section 6, we prove that A-MSA converges to a time-continuous solution of the sampled PMP. In section 7, we connect the solution of the sampled problem with the continuous one involving the exact average of the loss function. For this, we use results obtained in E et al. (2019). This last step allows to connect how much the solution of A-MSA deviates from the fully continuous one in terms of the number of samples and the accuracy of the discretization. In section 8, we give guidelines to implement A-MSA in practice, and section 9 illustrates the performance of the algorithm in numerical examples. We conclude the paper in Section 10.

2. From deep learning to optimal control

2.1 Statistical Learning Problems

We consider the following *statistical learning problem*: Assume that we are given a *domain set* $X \subseteq \mathbb{R}^n$ and a *label set* $Y \subseteq \mathbb{R}^k$, with $n, k \in \mathbb{N}$. Further assume that there exists an unknown probability distribution μ on $X \times Y$. Given a loss function $\mathcal{L} : Y \times Y \rightarrow \mathbb{R}^+$, the goal of the statistical learning problem is to find a function $v : X \rightarrow Y$, which we will call *prediction rule*, from a hypothesis class $\mathcal{V} \subset \{v : X \rightarrow Y\}$ such that the expected loss

$$\mathcal{J}(v) := \mathbb{E}_{(x,y) \sim \mu} \mathcal{L}(v(x), y)$$

is minimized over \mathcal{V} . In other words, the task is to find

$$v^* \in \arg \min_{v \in \mathcal{V}} \mathcal{J}(v).$$

In general, the probability distribution μ is unknown and we are only given a set \mathcal{S}_N of $N \in \mathbb{N}$ training samples

$$\mathcal{S}_N := \{(x_i, y_i)\}_{i=1}^N.$$

The most common choice is then to consider the uniform distribution $\mu_N : X \times Y \rightarrow \mathbb{R}^+$,

$$\mu_N(x, y) = \frac{1}{N} \sum_{i=1}^N \delta_{(x_i, y_i)}(x, y)$$

which yields the so-called *empirical loss*

$$\mathcal{J}_N(v) := \mathbb{E}_{(x,y) \sim \mu_N} \mathcal{L}(v(x), y) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(v(x_i), y_i)$$

and one finds v_N , an approximation of v , by minimizing over it,

$$v_N \in \arg \min_{v \in \mathcal{V}} \mathcal{J}_N(v). \quad (1)$$

In the following, the above optimization step will be called the *learning procedure*.

2.2 Neural Network Architectures and Deep Learning

While there exists numerous different architectures, we could simplify the discussion and say that neural networks are a class of functions with the basic general form

$$v : \mathbb{R}^n \rightarrow \mathbb{R}^k, \quad x \mapsto W_L \sigma(W_{L-1} \sigma(\dots (\sigma(W_1(x))))), \quad (2)$$

where:

- $L \in \mathbb{N}$ is the *depth*, i.e., the number of layers of the neural network,
- $W_\ell : \mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}$ are affine maps for $\ell = 1, \dots, L$. For consistency, we must set $N_0 = n$ and $N_L = k$ but the rest of the dimensions $N_\ell \in \mathbb{N}$ can be freely chosen. We have for all $x \in \mathbb{R}^{N_{\ell-1}}$, $W_\ell(x) = A_\ell(x) + b_\ell$ for a matrix $A_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ and a vector $b_\ell \in \mathbb{R}^{N_\ell}$.
- $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a *nonlinear activation function* which is applied coordinate-wise in (2). Some popular choices are the ReLU function, $\sigma(x) = \max(0, x)$, or the hyperbolic tangent, $\sigma(x) = \tanh(x)$.

Deep learning describes the range of learning procedures to solve statistical learning problems where the hypothesis class \mathcal{V} is taken to be a set of neural networks. One usually works with the class of neural networks with given depth L , activation function σ and fixed dimensions N_1, \dots, N_{L-1} for the affine mappings,

$$\mathcal{NN}(L, \sigma, N_1, \dots, N_{L-1}) := \{v : X \rightarrow Y : v(x) = W_L \sigma(W_{L-1} \sigma(\dots (\sigma(W_1(x)))))\}.$$

For this class, the task is to solve the empirical risk problem (1) associated to it.

2.3 Continuous Formulation as an Optimal Control Problem

Functions of the type (2) can be built by repeated composition of parametrizable functions

$$\phi_\ell : \mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}, \quad x \mapsto \phi_\ell(x) = \sigma(W_\ell(x)), \quad i = 1, \dots, L-1,$$

followed by an affine step, that is,

$$v = W_L \circ \phi_{L-1} \circ \dots \circ \phi_1, \quad \forall v \in \mathcal{NN}(L, \sigma, N_1, \dots, N_{L-1}).$$

This simple observation yields to numerous other possible architectures. The most relevant instance triggering the connection with optimal control are the so-called Residual Neural Networks (ResNet, He et al. (2016)). They correspond to the choice

$$\phi_\ell : \mathbb{R}^{N_{\ell-1}} \rightarrow \mathbb{R}^{N_\ell}, \quad x \mapsto \phi_\ell(x) = x + h\sigma(W_\ell(x)), \quad i = 1, \dots, L-1,$$

Network Architecture	Associated ODE discretization
ResNet He et al. (2016), RevNet Gomez et al. (2017)	Forward Euler
Polynet Zhang et al. (2017)	Approximation of Backward Euler
FractalNet Larsson et al. (2016)	Runge-Kutta order 2

Table 1: Some neural network classes and their associated ODE discretization.

for a certain parameter¹ $h > 0$.

If we now set $N_{L-1} = \dots = N_1 = N_0 = n$, the application of ϕ_ℓ can be interpreted as an Explicit Euler step from time $t_{\ell-1} = (\ell - 1)h$ to time $t_\ell = t_{\ell-1} + h$ of the dynamics

$$\dot{x}_t = f(x_t, \theta_t), \quad \text{with } x_0 = x,$$

where

$$f(x_t, \theta_t) = \sigma(W_t(x_t)), \quad \text{with } W_t(x) = A_t x + b_t,$$

and θ_t gathers the parameters upon which the dynamics depend. In our case,

$$\theta_t = (A_t, b_t).$$

Note that the matrices A and b are now time-dependent, hence the t -subscript.

In this view, ResNet functions can be interpreted as the output of performing $L - 1$ time steps of size h of the above dynamics, followed by an affine operation $W_L : \mathbb{R}^n \rightarrow \mathbb{R}^k$. This last operation is important since it maps the final output from the domain set in \mathbb{R}^n to the label set in \mathbb{R}^k . In fact, it can be replaced by any linear or nonlinear mapping $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$. Note also that h can be included as a parameter of the model class and that one can optimize over it as in Benning et al. (2019).

Following similar lines, one can view certain classes of neural networks as discretizations of an underlying continuous dynamical system. Table 1 lists some popular classes and the associated numerical scheme (see, e.g., Lu et al. (2018) for more details on these connections). We emphasize however that the connection with optimal control comes at the price of imposing that the input and intra-layer dimensions are equal ($N_{L-1} = \dots = N_1 = N_0$), thus it cannot be made for all classes of neural networks.

2.4 Related works

To the best of our knowledge, the connection between deep learning, dynamical systems and optimal control has been known since at least the 1980's, going back to the works of LeCun and Pineda in which the idea of back-propagation is connected to the adjoint variable arising in optimal control (see LeCun et al. (1988); Pineda (1988)). The approach has gained increasing attention in recent years. We refer to Weinan (2017); Haber and Ruthotto (2018); Li et al. (2018); Benning et al. (2019); Vialard et al. (2020) for some selected references. The E-MSA deep learning algorithm based on the optimal control framework has been introduced in Li et al. (2018) and is the main starting point for our developments on an adaptivity.

1. In fact, ResNets were originally defined for $h = 1$. The ‘‘augmented’’ definition with the parameter h was introduced in Haber and Ruthotto (2018).

A recent development on the topic of adaptivity has been made in Benning et al. (2019) where the idea is to fix a time-scheme with a *fixed* number of time steps, and include the time step as an additional parameter to be optimized. Note that this type of adaptivity lies in the time step sizes, and it is different than the one proposed here since we understand adaptivity as refining the time-discretization, thus increasing progressively the number of time steps.

3. Optimal control setting

We next define the notation for the optimal control setting which we will use in the rest of the paper. The notation is kept as consistent as possible with section 2 in order to further enhance the similarity between the usual deep learning approach and the present one.

Like in section 2, we assume that we are given a *domain set* $X \subset \mathbb{R}^n$ and a *label set* $Y \subset \mathbb{R}^k$, with $n, k \in \mathbb{N}$, and that there exists an unknown probability distribution μ on $X \times Y$ representing the distribution of the input-target pairs (x, y) . Consider now a set of admissible controls or training weights $\Theta \subseteq \mathbb{R}^m$. Usually, we set $\Theta = \mathbb{R}^m$ in deep learning but the present methodology allows to consider constraints we will denote $\Theta \subseteq [U_{\min}, U_{\max}]^m$. Fix $T > 0$ and let f (feed-forward dynamics), Φ (terminal loss function) and R (regularizer) be functions

$$f : \mathbb{R}^n \times \Theta \rightarrow \mathbb{R}^n, \quad \Phi : \mathbb{R}^n \times \mathbb{R}^k \rightarrow \mathbb{R}, \quad R : \Theta \rightarrow \mathbb{R}.$$

Let $L^\infty([0, T], \Theta)$ be the set of essentially bounded measurable controls. In the following, we use bold-faced letters for path-time quantities. For example, $\boldsymbol{\theta} := \{\theta_t : 0 \leq t \leq T\}$ for any $\boldsymbol{\theta} \in L^\infty([0, T], \Theta)$. For every control $\boldsymbol{\theta} \in L^\infty([0, T], \Theta)$ and every value of the random variable $x \in X$, we define the state dynamics $\mathbf{u}^{\boldsymbol{\theta}, x} := \{u_t^{\boldsymbol{\theta}, x} : 0 \leq t \leq T\}$ as the solution to the ordinary differential equation (ODE),

$$\begin{cases} \dot{u}_t^{\boldsymbol{\theta}, x} &= f(u_t^{\boldsymbol{\theta}, x}, \theta_t), \quad \forall t \in (0, T) \\ u_0^{\boldsymbol{\theta}, x} &= x \end{cases} \quad (3)$$

The ODE is stochastic and its only source of randomness is the initial condition x .

With this notation, the deep learning optimization problem can be posed as the optimal control problem of finding

$$J^* = \inf_{\boldsymbol{\theta} \in L^\infty([0, T], \Theta)} \mathcal{J}(\boldsymbol{\theta}), \quad \text{subject to (3)}, \quad (4)$$

where

$$\mathcal{J}(\boldsymbol{\theta}) := \mathbb{E}_{(x, y) \sim \mu} [\text{Loss}(x, y, \boldsymbol{\theta})],$$

and for any input-target pair $(x, y) \in \mathbb{R}^n \times \mathbb{R}^k$, the loss function is defined as

$$\text{Loss}(x, y, \boldsymbol{\theta}) := \Phi(u_T^{\boldsymbol{\theta}, x}, y) + \int_0^T R(\theta_t) dt$$

Note that the regularizer R could in general also depend on the state u and Φ is the actual loss function upon which we want to act. It plays the same role as the loss function \mathcal{L} of

section 2.1. We can relate them by introducing the mapping $g : \mathbb{R}^n \rightarrow \mathbb{R}^k$ from section 2.3, and setting

$$\Phi(u_T^{\theta,x}, y) = \mathcal{L}(g(u_T^{\theta,x}), y).$$

Like in the setting of section 2.1, we are only given a set of N samples $\mathcal{S}_N = \{(x_i, y_i)\}_{i=1}^N$ and the probability distribution μ is not exactly known. For a given control $\theta \in L^\infty([0, T], \Theta)$, each sample follows the dynamics

$$\begin{cases} \dot{u}_t^{\theta,i} &= f(u_t^{\theta,i}, \theta_t), \quad \forall t \in (0, T) \\ u_0^{\theta,i} &= x_i \end{cases} \quad (5)$$

for $i = 1, \dots, N$. We then perform empirical risk minimization taking, e.g., the uniform distribution $\mu_N = \frac{1}{N} \sum_{i=1}^N \delta_{(x_i, y_i)}$. The optimal control problem becomes

$$J_{\mathcal{S}_N}^* = \inf_{\theta \in L^\infty([0, T], \Theta)} \mathcal{J}_{\mathcal{S}_N}(\theta), \quad \text{subject to (5)}, \quad (6)$$

where

$$\mathcal{J}_{\mathcal{S}_N}(\theta) := \mathbb{E}_{(x,y) \sim \mu_N} \text{Loss}(x, y, \theta) = \frac{1}{N} \sum_{i=1}^N \Phi(u_T^{\theta,i}, y_i) + \int_0^T R(\theta_t) dt$$

Note that the solutions of the sampled optimal control problem (6) depend on the sample set \mathcal{S}_N and are therefore random variables. However, for \mathcal{S}_N fixed, problem (6) is deterministic and can thus be solved with deterministic optimal control techniques. In what follows, we adopt this viewpoint first to solve problem (6) at any target accuracy. By this we mean the following: for any given target accuracy $\varepsilon > 0$, we prove in Theorem 4 that it is possible to numerically compute a control θ_ε such that

$$\mathcal{J}_{\mathcal{S}_N}(\theta_\varepsilon) - J_{\mathcal{S}_N}^* \leq \varepsilon.$$

We then come back to the probabilistic point of view and prove in Corollary 7 that, provided that the number N of samples is sufficiently large, there exists a constant $C > 0$ such that

$$\mathcal{J}_{\mathcal{S}_N}(\theta_\varepsilon) - J^* \leq C\varepsilon.$$

with high probability. The derivation of the latter bound relies on the results of E et al. (2019).

4. Pontryagin's Maximum Principle and Method of Successive Approximations

In this section, we focus on the *sampled* optimal control problem (6). We first recall the necessary optimality conditions, usually known as the Pontryagin's Maximum Principle. We next introduce the main algorithm to solve the PMP which will be the starting point for our subsequent developments. In the following, the Euclidean norm of any vector $a \in \mathbb{R}^n$ is denoted by $\|a\|$ and the scalar product with any other vector $b \in \mathbb{R}^n$ is $a \cdot b$.

4.1 Pontryagin's Maximum Principle

We define the Hamiltonian $H : [0, T] \times \mathbb{R}^n \times \mathbb{R}^n \times \Theta \rightarrow \mathbb{R}$ as

$$H(t, u, p, \theta) := p \cdot f(u, \theta) - R(\theta). \quad (7)$$

The following classical Pontryagin's Maximum Principle (PMP) gives necessary optimality conditions to problem (6).

Theorem 1 *Let $\theta^* \in L^\infty([0, T], \Theta)$ be an optimal control to (6). For $i = 1, \dots, N$, let $u^{\theta^*, i}$ be the associated state process with initial condition x_i . There exists a co-state process $p^{\theta^*, i} \in L^\infty([0, T], \mathbb{R}^n)$ such that*

$$\begin{aligned} \dot{u}_t^{\theta^*, i} &= f(u_t^{\theta^*, i}, \theta_t^*), & u_0^{\theta^*, i} &= x_i, \\ \dot{p}_t^{\theta^*, i} &= -\nabla_u H(t, u_t^{\theta^*, i}, p_t^{\theta^*, i}, \theta_t^*), & p_T^{\theta^*, i} &= -\nabla_u \Phi(u_T^{\theta^*, i}, y_i), \end{aligned}$$

and, for each $t \in [0, T]$,

$$\frac{1}{N} \sum_{i=1}^N H(t, u_t^{\theta^*, i}, p_t^{\theta^*, i}, \theta_t^*) \geq \frac{1}{N} \sum_{i=1}^N H(t, u_t^{\theta^*, i}, p_t^{\theta^*, i}, \theta), \quad \forall \theta \in \Theta. \quad (8)$$

The proof of this theorem and its variants can be found in any optimal control theory reference (see, e.g., Bertsekas (1995); Clarke (2005); Athans and Falb (2013)). We omitted the case of an abnormal multiplier, which will not be considered in the following.

We emphasize that the PMP is only a necessary condition, so there can be cases where the solutions to the PMP are not global optima for (6). Nevertheless, in practice the PMP often gives good solution candidates, and when certain convexity assumptions are satisfied the PMP becomes sufficient (see Bressan and Piccoli (2007)). In the next section, we discuss the numerical methods that we take as a starting point to solve the PMP.

4.2 MSA and Extended-MSA

A classical algorithm to find θ^* and the corresponding forward and co-state dynamics $u^{\theta^*, i}$ and $p^{\theta^*, i}$ is the Method of Successive Approximations (see Chernousko and Lyubushin (1982)). It is a fixed-point method based on the following steps. Starting from an initial guess of the optimal control θ^0 , for each $k \geq 0$ we first solve for $i = 1, \dots, N$ the forward dynamics

$$\dot{u}_t^{\theta^k, i} = \nabla_p H(t, u_t^{\theta^k, i}, p_t^{\theta^k, i}, \theta_t^k) = f(u_t^{\theta^k, i}, \theta_t^k), \quad u_0^{\theta^k, i} = x_i. \quad (9)$$

The dynamics $u^{\theta^k, i}$ allows us to compute the backward dynamics

$$\dot{p}_t^{\theta^k, i} = -\nabla_u H(t, u_t^{\theta^k, i}, p_t^{\theta^k, i}, \theta_t^k), \quad p_T^{\theta^k, i} = -\nabla_u \Phi(u_T^{\theta^k, i}, y_i). \quad (10)$$

Finally, we update the control by using the maximization condition (8),

$$\theta_t^{k+1} \in \arg \max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N H(t, u_t^{\theta^k, i}, p_t^{\theta^k, i}, \theta), \quad \forall t \in [0, T]. \quad (11)$$

In this form, MSA converges only locally when it is initialized with a starting control guess that is sufficiently close to an optimal control θ^* . To overcome this limitation, an Extended MSA algorithm (E-MSA) based on an augmented Lagrangian strategy has been introduced in Li et al. (2018). This algorithm is the starting point for our subsequent developments. It works as follows: fix some $\rho > 0$ and define the augmented Hamiltonian

$$\tilde{H}(t, u, p, \theta, v, q) := H(t, u, p, \theta) - \frac{\rho}{2} \|v - f(u, \theta)\|^2 - \frac{\rho}{2} \|q + \nabla_u H(t, u, p, \theta)\|^2. \quad (12)$$

We can now formulate an extended PMP based on \tilde{H} .

Proposition 2 (Extended PMP, see Li et al. (2018)) *Let $\theta^* \in L^\infty([0, T], \Theta)$ be an optimal control to (6). For $i = 1, \dots, N$, let $\mathbf{u}^{\theta^*, i}$ be the associated state process with initial condition x_i . There exists a co-state process $\mathbf{p}^{\theta^*, i} \in L^\infty([0, T], \mathbb{R}^n)$ such that*

$$\begin{aligned} \dot{\mathbf{u}}_t^{\theta^*, i} &= \nabla_p \tilde{H}(t, \mathbf{u}_t^{\theta^*, i}, \mathbf{p}_t^{\theta^*, i}, \theta_t^*), & \mathbf{u}_0^{\theta^*, i} &= x_i, \\ \dot{\mathbf{p}}_t^{\theta^*, i} &= -\nabla_u \tilde{H}(t, \mathbf{u}_t^{\theta^*, i}, \mathbf{p}_t^{\theta^*, i}, \theta_t^*), & \mathbf{p}_T^{\theta^*, i} &= -\nabla_u \Phi(\mathbf{u}_T^{\theta^*, i}, \mathbf{y}_i), \end{aligned}$$

and, for each $t \in [0, T]$,

$$\frac{1}{N} \sum_{i=1}^N \tilde{H}(t, \mathbf{u}_t^{\theta^*, i}, \mathbf{p}_t^{\theta^*, i}, \theta_t^*) \geq \frac{1}{N} \sum_{i=1}^N \tilde{H}(t, \mathbf{u}_t^{\theta^*, i}, \mathbf{p}_t^{\theta^*, i}, \theta), \quad \forall \theta \in \Theta.$$

The E-MSA algorithm consists in applying the MSA algorithm with the augmented Hamiltonian \tilde{H} instead of with H . Since $\nabla_p \tilde{H} = \nabla_p H = f$ and $\nabla_u \tilde{H} = \nabla_u H$, steps (9) and (10) remain the same as in MSA and the maximization step (11) is replaced by

$$\theta_t^{k+1} \in \arg \max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \tilde{H}(t, \mathbf{u}_t^{\theta^k, i}, \mathbf{p}_t^{\theta^k, i}, \theta), \quad \forall t \in [0, T]. \quad (13)$$

It has been proven in Li et al. (2018) that if the parameter ρ is taken sufficiently large, the E-MSA algorithm converges to the set of solution of the extended PMP for any initial guess of the control θ^0 . However, this scheme cannot be realized in practice without discretizing. This introduces a perturbation of the time-continuous formulation of the algorithm since the resulting trajectories will be approximations of the continuous one. To guarantee convergence to the exact continuous solution, it is necessary to identify suitable approximation error tolerances that still guarantee convergence to the exact solution. This motivates to introduce an Adaptive MSA scheme that we describe in the next section. We prove that if the forward and backward propagations are solved with increasing accuracy at each step, then the algorithm converges towards a *continuous* solution of the Extended PMP.

5. A-MSA: Adaptive Method of Successive Approximations

The algorithm requires defining forward and backward time-integration schemes, which we introduce in section 5.1. We next present the algorithm and how predictions are made in sections 5.2 and 5.3. We end up by proving convergence in section 6.

5.1 Routines solve_fwd and solve_bwd

In the following, we work with time-integration routines `solve_fwd` and `solve_bwd` which can be cast in the abstract framework that we next describe. Relevant particular instances of it will be continuous Galerkin (cG) or Runge-Kutta collocation (RK-C) methods.

Setting: We describe the framework in the case of the forward dynamics

$$\begin{cases} \dot{u}_t &= f(u_t), \quad \forall t \in (0, T) \\ u_0 &= x \end{cases} \quad (14)$$

for a generic Lipschitz continuous function $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$. The backward dynamics can be deduced similarly.

Let $0 = t_0 < t_1 < \dots < t_L = T$ be a partition of $[0, T]$ and let $\mathcal{T}_\ell := (t_{\ell-1}, t_\ell]$, and $h_\ell := t_\ell - t_{\ell-1}$. For $q \in \mathbb{N}$, let $\mathcal{V}_q^{\mathcal{T}}$ be the space of continuous functions that are piecewise polynomials of degree q in the time mesh $\mathcal{T} = \cup_{\ell=1}^L \mathcal{T}_\ell$, i.e.,

$$\mathcal{V}_q^{\mathcal{T}} := \{v \in \mathcal{C}^0([0, T]; \mathbb{R}^n) : v|_{\mathcal{T}_\ell} \in \mathbb{P}^{(q)}(\mathcal{T}_\ell), \quad \ell = 1, \dots, L\},$$

where

$$\mathbb{P}^{(q)}(\mathcal{T}_\ell) := \{v \in \mathcal{C}^0(\mathcal{T}_\ell; \mathbb{R}^n) : \forall t \in \mathcal{T}_\ell, v(t) = \sum_{j=0}^q t^j v_j, v_j \in \mathbb{R}^n\}$$

is the space of polynomials of degree q in the interval \mathcal{T}_ℓ .

Introducing a *projection operator*

$$\Pi^{(q)} : \mathcal{C}^0([0, T]; \mathbb{R}^n) \mapsto \mathcal{V}_q^{\mathcal{T}}$$

the time discrete approximation U to the solution u is defined as follows: we seek $U \in \mathcal{V}_q^{\mathcal{T}}$ satisfying the initial condition $U(0) = x$ as well as

$$U'_t = \Pi^{(q-1)} f(U_t), \quad \forall t \in \mathcal{T}_\ell, \quad (15)$$

for $\ell = 1, \dots, L$. Note that since both terms belong to $\mathbb{P}^{(q-1)}(\mathcal{T}_\ell)$, (15) admits a Galerkin formulation

$$\int_{\mathcal{T}_\ell} v \cdot U'_t dt = \int_{\mathcal{T}_\ell} v \cdot \Pi^{(q-1)} f(U_t) dt, \quad \forall v \in \mathbb{P}^{(q)}(\mathcal{T}_\ell). \quad (16)$$

In the following, we use mainly (15) but (16) is of interest since it connects (RK-C) methods with (cG) methods. It is proven in Akrivis et al. (2009), that the continuous Galerkin method corresponds to the choice $\Pi^{(q-1)} := P^{(q-1)}$, with $P^{(q-1)}$ denoting the (local) L^2 orthogonal projection onto $\mathbb{P}^{(q-1)}(\mathcal{T}_\ell)$ for each ℓ . Runge-Kutta collocation methods with pairwise distinct nodes in \mathcal{T}_ℓ can be obtained by choosing $\Pi^{(q-1)} := I^{(q-1)}$ with $I^{(q-1)}$ denoting the interpolation operator by elements of $\mathbb{P}^{(q)}(\mathcal{T}_\ell)$ at the nodes $t_{\ell-1} + \alpha_i(t_\ell - t_{\ell-1})$, $i = 1, \dots, q$, $\ell = 1, \dots, L$, with appropriate weights $0 \leq \alpha_1 < \dots < \alpha_q \leq 1$. All RK-C methods with pairwise distinct nodes in $[0, 1]$ can be obtained by applying appropriate numerical quadrature to continuous Galerkin methods.

L^2 error estimation: One simple way of estimating the error between u and U is via the following Grönwall inequality. Denoting $e_t = u_t - U_t$, we have

$$\dot{e}_t = f(u_t) - \Pi^{(q-1)}f(U_t) = (f(u_t) - f(U_t)) + (f(U_t) - \Pi^{(q-1)}f(U_t))$$

Multiplying by e_t , we get

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} e_t^2 &= e_t (f(u_t) - f(U_t)) + e_t (f(U_t) - \Pi^{(q-1)}f(U_t)) \\ &\leq L e_t^2 + |e_t| |f(U_t) - \Pi^{(q-1)}f(U_t)| \\ &\leq (L + 1/2) e_t^2 + \frac{1}{2} |f(U_t) - \Pi^{(q-1)}f(U_t)|^2, \end{aligned}$$

where L is the Lipschitz constant on the second variable of g . By the Grönwall inequality, since $e(0) = 0$,

$$e_t^2 \leq \int_0^t |f(U_s) - \Pi^{(q-1)}f(U_s)|^2 e^{(2L+1)(t-s)} ds$$

which yields an L^2 estimate of the error

$$\begin{aligned} \|u - U\|_{L^2([0,T], \mathbb{R}^n)}^2 &\leq \int_0^T \int_0^t |f(U_s) - \Pi^{(q-1)}f(U_s)|^2 e^{(2L+1)(t-s)} ds dt \\ &\leq \frac{e^{(2L+1)T}}{2L+1} \|f - \Pi^{(q-1)}f\|_{L^2([0,T], \mathbb{R}^n)}^2 \end{aligned} \quad (17)$$

In the following, we denote Π_ζ a projector delivering an accuracy $\eta \geq 0$ for the involved dynamics

$$\|f - \Pi^{(q-1)}f\|_{L^2([0,T], \mathbb{R}^n)} \leq \eta.$$

From inequality (17), such a projector gives an accuracy in the solution U which is bounded by

$$\|u - U\|_{L^2([0,T], \mathbb{R}^n)} \leq \frac{e^{(2L+1)T/2}}{(2L+1)^{1/2}} \eta$$

The routines: In the following, the routine

$$\mathbf{u}^\eta = \text{solve_fwd}(\eta; f, x)$$

yields an evolution

$$\mathbf{u}^\eta = \{u_t^\eta : t \in [0, T]\}$$

that approximates the exact solution \mathbf{u} of the *forward* dynamics (14) with a projector Π_η . The routine

$$\mathbf{p}^\eta = \text{solve_bwd}(\eta; f, x)$$

works similarly for the *backward* dynamics $\dot{p}_t = f(p_t)$ with $p_T = x$.

5.2 The learning phase of the algorithm

The A-MSA algorithm consists in performing steps (9) and (10) of the original MSA but with a numerical time-integrators which give inexact trajectories. We carry the discussion assuming that we use the routines `solve_fwd` and `solve_bwd` introduced above. At every iteration k , both routines involves projectors Π_k that deliver an accuracy ε_k which is yet to be determined. The accuracy will be tightened as k increases in a way that still guarantees convergence to the exact, continuous solution (6) of the sampled PMP problem. We deduce its value later on from the convergence analysis of section 6.

Starting from an initial guess of the optimal control θ^0 , for each $k \geq 0$ we first solve at accuracy ε_k the forward dynamics

$$\mathbf{u}^{\theta^k, i, \varepsilon_k} = \text{solve_fwd} \left(\varepsilon_k; (t, u) \mapsto f(u, \theta_t^k), x_i \right), \quad \forall i = 1, \dots, N.$$

Satisfying this accuracy may require to adapt the time discretization mesh and/or the numerical scheme as discussed in the previous section.

We then use the final state $u_T^{\theta^k, i, \varepsilon_k}$ to compute an approximation of the *backward* dynamics at the same accuracy ε_k

$$\mathbf{p}^{\theta^k, i, \varepsilon_k} = \text{solve_bwd} \left(\varepsilon_k; (t, p) \mapsto -\nabla_u H(t, u_t^{\theta^k, i, \varepsilon_k}, p, \theta_t^k), -\nabla_u \Phi(u_T^{\theta^k, i, \varepsilon_k}, y_i) \right), \quad \forall i = 1, \dots, N.$$

This step may, again, require adaptivity to satisfy the target tolerance. To update the control, instead of finding the exact maximum like in (13), it is in fact sufficient to find controls that are away by a factor $0 < \gamma_k \leq 1$ of the maximum and for which $\gamma_k \rightarrow 1$ as $k \rightarrow \infty$. In other words, it suffices to find for all $t \in [0, T]$, a control $\theta_t^{k+1} \in \Theta$ such that

$$\frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1}) \geq \gamma_k \max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta), \quad \forall t \in [0, T]. \quad (18)$$

Note that when $\gamma_k = 1$, we fall back to the case where we look for the exact maximum. Algorithm 1 summarizes the whole procedure by describing a routine `A-MSA[S, τ]` which computes a control dynamics $\hat{\theta}$ that solves the extended PMP problem until the difference between two successive iterates is smaller than $\tau > 0$. The algorithm also yields a numerical scheme associated to the parameter $\hat{\varepsilon}$ which is the accuracy of the forward and backward solvers used in the last iteration. These elements are then used for prediction as explained in the next section.

The practical implementation of the algorithm requires to specify the tolerances ε_k . They can be theoretically derived from the convergence analysis of section 6. However, since the bounds of the analysis may not be sharp and that it involves quantities which are difficult to estimate in practice, devote section 8 to give indications on how to implement A-MSA in practice.

5.3 Predictions

The routine `A-MSA[S, τ]` gives a discrete control $\hat{\theta}$ and an accuracy $\hat{\varepsilon}$. We also have a certain numerical scheme for the forward propagator, which can be interpreted as the final

Algorithm 1 Learning Algorithm: A-MSA[\mathcal{S}_N, τ] \rightarrow $[\hat{\theta}, \hat{\varepsilon}]$

- 1: **Input:** $\mathcal{S}_N = \{x_i, y_i\}_{i=1}^N$, target tolerance τ
- 2: **Set of internal parameters:** θ^0 , ρ , maximum number of iterations k_{\max} .
- 3: $k \leftarrow 0$
- 4: $J, J_{\text{old}} \leftarrow J(\theta^0)$
- 5: $\Delta \leftarrow \tau + 1$
- 6: **while** $\Delta > \tau$ or $k \leq k_{\max}$ **do**
- 7: Update ε_k and γ_k
- 8: **for** i in $\{1, \dots, N\}$ **do** // parallel
- 9: $\mathbf{u}^{\theta^k, i, \varepsilon_k} = \text{solve_fwd}(\varepsilon_k; (t, u) \mapsto f(u, \theta_t^k), x_i)$
- 10: $\mathbf{p}^{\theta^k, i, \varepsilon_k} = \text{solve_bwd}(\varepsilon_k; (t, p) \mapsto -\nabla_x H(t, u_t^{\theta^k, i, \varepsilon_k}, p, \theta_t^k), -\nabla_u \Phi(u_T^{\theta^k, i, \varepsilon_k}, y_i))$
- 11: **end for**
- 12: **for all** $t \in [0, T]$ **do** // parallel
- 13: Find θ_t^{k+1} s.t.

$$\frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1}) \geq \gamma_k \max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta)$$

- 14: **end for**
 - 15: $J_{\text{old}} \leftarrow J$; $J \leftarrow J(\theta^{k+1})$; $\Delta \leftarrow J_{\text{old}} - J$
 - 16: **end while**
 - 17: $\hat{\theta} \leftarrow \theta^{k+1}$; $\hat{\varepsilon} \leftarrow \varepsilon_k$
 - 18: **Output:** $[\hat{\theta}, \hat{\varepsilon}]$
-

network architecture. As a result, for a given input data x , we can predict the output by first computing

$$\hat{u} = \text{solve_fwd} \left(\hat{\varepsilon}; (t, u) \mapsto f(u, \hat{\theta}), x \right),$$

and then taking

$$\hat{y} := g(\hat{u}_T)$$

as the approximation of the true y . In other words, our neural network is the mapping

$$\begin{aligned} \mathcal{NN} : \mathbb{R}^n &\rightarrow \mathbb{R}^k \\ x \rightarrow \mathcal{NN}(x) &:= g(\hat{u}_T) = g \left(\text{solve_fwd} \left(\hat{\varepsilon}; (t, u) \mapsto f(u, \hat{\theta}), x \right) (t = T) \right) \end{aligned}$$

We summarize the prediction pipeline in Algorithm 2.

Algorithm 2 Prediction algorithm: $\mathcal{NN}(x) \rightarrow \hat{y}$

- 1: **Input:** Observation $x \in \mathbb{R}^n$, Neural network architecture $[\hat{\theta}, \hat{\varepsilon}]$ (from A-MSA $[\mathcal{S}_N, \tau]$).
 - 2: $\hat{u} \leftarrow \text{solve_fwd} \left(\hat{\varepsilon}; (t, u) \mapsto f(u, \hat{\theta}), x \right)$
 - 3: **Output:** $g(\hat{u}_T)$
-

6. A priori convergence of A-MSA

In this section, we prove that, for suitably chosen tolerances $(\varepsilon_k)_k$, the sequence of controls $(\theta^k)_{k \geq 0}$ computed with A-MSA is a minimizing sequence for problem (6), that is,

$$\mathcal{J}_{\mathcal{S}_N}(\theta^k) \xrightarrow{k \rightarrow \infty} J_{\mathcal{S}_N}^*.$$

Our analysis is built upon the one presented in Li et al. (2018) for the exact E-MSA to which it is necessary to add perturbative arguments due to the inexact propagations.

We work with the same continuity assumptions as in Li et al. (2018) for the ideal case of exact propagations:

- (A1) Φ is twice continuously differentiable. Φ and $\nabla\Phi$ satisfy the following Lipschitz condition: there exists $K > 0$ such that

$$|\Phi(u) - \Phi(u')| + \|\nabla\Phi(u) - \nabla\Phi(u')\| \leq K\|u - u'\|, \quad \forall (u, u') \in \mathbb{R}^n.$$

- (A2) For all $t \in [0, T]$ and $\theta \in \Theta$, $u \rightarrow f(u, \theta)$ is twice continuously differentiable and satisfying the following Lipschitz condition: there exists $K' > 0$ such that

$$\|f(u, \theta) - f(u', \theta)\| + \|\nabla_u f(u, \theta) - \nabla_u f(u', \theta)\|_2 \leq K\|u - u'\|, \quad \forall (u, u') \in \mathbb{R}^n,$$

where $\|\cdot\|_2$ denotes the induced 2-norm.

Intermediate problem: We first study the following intermediate problem. Consider a fixed input-output pair (x, y) . For a given control $\theta \in L^\infty([0, T], \Theta)$, let $\mathbf{u}^{\theta, x}$ and $\mathbf{p}^{\theta, x}$ be the exact solutions of the forward and backward propagations using x and y in the initial and final conditions,

$$\begin{cases} \dot{u}_t^{\theta, x} &= f(u_t^{\theta, x}, \theta_t), & u_0^{\theta, x} &= x, \\ \dot{p}_t^{\theta, x} &= -\nabla_u H(t, u_t^{\theta, x}, p_t^{\theta, x}, \theta_t), & p_T^{\theta, x} &= -\nabla_u \Phi(u_T^{\theta, x}, y), \end{cases}$$

Let $\mathbf{u}^{\theta, x, \zeta}$ and $\mathbf{p}^{\theta, x, \zeta}$ be the outputs of the same propagation but with accuracy $\zeta > 0$,

$$\begin{cases} \mathbf{u}^{\theta, x, \zeta} &= \text{solve_fwd}(\zeta; (t, u) \mapsto f(u, \theta), x) \\ \mathbf{p}^{\theta, x, \zeta} &= \text{solve_bwd}(\zeta; (t, p) \mapsto -\nabla_u H(t, u_t^{\theta, x, \zeta}, p, \theta_t), -\nabla_u \Phi(u_T^{\theta, x, \zeta}, y)). \end{cases} \quad (19)$$

Similarly, $\mathbf{u}^{\varphi, x, \eta}$ and $\mathbf{p}^{\varphi, x, \eta}$ denote the η -accurate solutions but for another control $\varphi \in L^\infty([0, T], \Theta)$.

The loss function associated to $\mathbf{u}^{\theta, x, \zeta}$ is denoted

$$\text{Loss}(x, y, \theta, \zeta) := \Phi(u_T^{\theta, x, \zeta}, y) + \int_0^T R(\theta_t) dt,$$

and similarly for $\mathbf{u}^{\varphi, x, \eta}$. In addition, we define the quantity

$$\Delta H_{\varphi, \theta}^{x, \zeta}(t) := H(t, u_t^{\theta, \zeta}, p_t^{\theta, \zeta}, \varphi_t) - H(t, u_t^{\theta, \zeta}, p_t^{\theta, \zeta}, \theta_t), \quad (20)$$

which will play an essential role in what follows.

Lemma 3 is an important building-block in the proof of convergence of A-MSA. To not interrupt the flow of reading, we have deferred its proof to Appendix A.

Lemma 3 *Let $\theta, \varphi \in L^\infty([0, T], \Theta)$ be two controls and let $(x, y) \in \mathbb{R}^n \times \mathbb{R}^k$ be an input-output pair. Let $(\mathbf{u}^{\theta, x, \zeta}, \mathbf{p}^{\theta, x, \zeta})$ and $(\mathbf{u}^{\varphi, x, \eta}, \mathbf{p}^{\varphi, x, \eta})$ be solutions of the forward and backward associated problems with accuracy ζ and η like in (19). Then, there exist a constant $C > 0$ independent of θ, φ but dependent on $T, \|\Pi_\zeta\|$ and $\|\Pi_\eta\|$ such that, if $\eta \leq 1$,*

$$\begin{aligned} & \text{Loss}(x, y, \varphi, \eta) - \text{Loss}(x, y, \theta, \zeta) \\ & \leq - \int_0^T \Delta H_{\varphi, \theta}^{x, \zeta}(t) dt + C \left((\eta + \zeta)^2 + \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T])}^2 \right) \end{aligned} \quad (21)$$

Convergence of A-MSA: We can use Lemma 3 to prove our main convergence result which we give in the following theorem. It involves the quantity

$$\begin{aligned} \lambda_k^2 &:= \frac{1}{N} \sum_{i=1}^N \|\nabla_w(\Delta H_{\theta^{k+1}, \theta^k}^{i, \varepsilon_k})\|_{L^2([0, T])}^2 \\ &= \frac{1}{N} \sum_{i=1}^N \left(\int_0^T \|f(u_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1}) - f(u_t^{\theta^k, i, \varepsilon_k}, \theta_t^k)\|^2 dt \right. \\ & \quad \left. + \int_0^T \|\nabla_u H(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1}) - \nabla_u H(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^k)\|^2 dt \right) \end{aligned}$$

Theorem 4 Let $\delta > 0$ be a fixed constant and suppose that we run the A-MSA algorithm with the following accuracies:

- (ε_k) is a positive and decreasing sequence such that

$$0 \leq \varepsilon_k \leq \min \left(\varepsilon_{k-1}, \frac{\delta \lambda_k^2}{4C + \rho N^{-1/2} \lambda_k} \right), \quad \forall k \geq 1.$$

- The maximization step (18) is performed with the parameter

$$\gamma_k \geq \frac{|\sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1})/N|}{\delta \lambda_k^2 + |\sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1})/N|}, \quad \forall k \in \mathbb{N}.$$

Then, if $\rho > 2(C + 2\delta + \varepsilon_0)$, where C is the constant of Lemma 3, the sequence of controls (θ^k) of A-MSA is a minimizing sequence for problem (6), that is,

$$\mathcal{J}_{S_N}(\theta^k) \rightarrow_{k \rightarrow \infty} J_{S_N}^*.$$

Furthermore, $\lambda_k^2 \rightarrow 0$ as $k \rightarrow \infty$ and we have

$$\mathcal{J}_{S_N}(\theta^{k+1}) - \mathcal{J}_{S_N}(\theta^k) \leq -\kappa \lambda_k^2, \quad \kappa := |C + 2\delta - \frac{\rho}{2}| > 0, \quad \forall k \in \mathbb{N}.$$

Proof We start by applying inequality (21) from Lemma 3 to the controls $\varphi = \theta^{k+1}$ and $\theta = \theta^k$ of the A-MSA algorithm, with accuracies $\eta = \varepsilon_{k+1}$ and $\zeta = \varepsilon_k$ and for the samples $(x, y) = (x_i, y_i)$. We then take the average over i and obtain

$$\begin{aligned} \mathcal{J}_{S_N}(\theta^{k+1}) - \mathcal{J}_{S_N}(\theta^k) &= \frac{1}{N} \sum_{i=1}^N \text{Loss}(x_i, y_i, \theta^{k+1}, \varepsilon_{k+1}) - \text{Loss}(x_i, y_i, \theta^k, \varepsilon_k) \\ &\leq -\frac{1}{N} \sum_{i=1}^N \int_0^T \Delta H_{\theta^{k+1}, \theta^k}^{i, \varepsilon_k}(t) dt + C((\varepsilon_k + \varepsilon_{k+1})^2 + \lambda_k^2) \end{aligned} \quad (22)$$

We first do the proof when $\gamma_k = 1, \forall k \geq 1$. In this case, from the maximization step (18) in A-MSA,

$$\frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1}, \dot{u}_t^{\theta^k, i, \varepsilon_k}, \dot{p}_t^{\theta^k, i, \varepsilon_k}) \geq \frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^k, \dot{u}_t^{\theta^k, i, \varepsilon_k}, \dot{p}_t^{\theta^k, i, \varepsilon_k}).$$

Recalling definition (12) for \tilde{H} , we reassemble the terms in the above inequality to derive

$$\begin{aligned} & -\frac{1}{N} \sum_{i=1}^N \int_0^T \Delta H_{\theta^{k+1}, \theta^k}^{i, \varepsilon_k}(t) dt \\ & \leq -\frac{\rho}{2N} \sum_{i=1}^N \int_0^T \left(\|\dot{u}_t^{\theta^k, i, \varepsilon_k} - f(u_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1})\|^2 + \|\dot{p}_t^{\theta^k, i, \varepsilon_k} - \nabla_u H(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^{k+1})\|^2 \right) dt \\ & + \frac{\rho}{2N} \sum_{i=1}^N \int_0^T \left(\|\dot{u}_t^{\theta^k, i, \varepsilon_k} - f(u_t^{\theta^k, i, \varepsilon_k}, \theta_t^k)\|^2 + \|\dot{p}_t^{\theta^k, i, \varepsilon_k} - \nabla_u H(t, u_t^{\theta^k, i, \varepsilon_k}, p_t^{\theta^k, i, \varepsilon_k}, \theta_t^k)\|^2 \right) dt \end{aligned} \quad (23)$$

Since by construction

$$\dot{u}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} = \Pi_{\varepsilon_k} f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k)$$

and

$$\|f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - \Pi_{\varepsilon_k} f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k)\|_{L^2([0, T], \mathbb{R}^n)} \leq \varepsilon_k,$$

we have by Cauchy-Schwartz inequality

$$\begin{aligned} & \int_0^T \|\dot{u}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1})\|^2 dt \\ &= \int_0^T \|f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - \Pi f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k)\|^2 + \|f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1})\|^2 \\ & \quad - 2 \left\langle f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - \Pi f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k), f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \right\rangle dt \\ & \geq \int_0^T \|f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - \Pi f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k)\|^2 + \|f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1})\|^2 dt \\ & \quad - 2\varepsilon_k \left(\int_0^T \|f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1})\|^2 dt \right)^{1/2} \end{aligned}$$

Since a similar bound holds for $\int_0^T \|\dot{p}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} - \nabla_u H(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1})\|^2 dt$, we derive that

$$- \frac{1}{N} \sum_{i=1}^N \int_0^T \Delta H_{\boldsymbol{\theta}^{k+1}, \boldsymbol{\theta}^k}^{i, \varepsilon_k}(t) dt \leq -\frac{\rho}{2} \lambda_k^2 + \rho N^{-1/2} \varepsilon_k \lambda_k \quad (24)$$

Using (24), we can further bound (22) to obtain

$$\begin{aligned} \mathcal{J}_{S_N}(\boldsymbol{\theta}^{k+1}) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^k) & \leq -\frac{\rho}{2} \lambda_k^2 + \rho N^{-1/2} \varepsilon_k \lambda_k + C((\varepsilon_k + \varepsilon_{k+1})^2 + \lambda_k^2) \\ & \leq \left(C - \frac{\rho}{2}\right) \lambda_k^2 + \varepsilon_k (4C + \rho N^{-1/2} \lambda_k) \end{aligned} \quad (25)$$

where we have used that (ε_k) is a decreasing sequence.

Since, for a fixed $\delta > 0$, we have chosen ε_k such that $\varepsilon_k \leq \delta \lambda_k^2 / (4C + \rho N^{-1/2} \lambda_k)$, we infer that

$$\mathcal{J}_{S_N}(\boldsymbol{\theta}^{k+1}) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^k) \leq \left(C + \delta - \frac{\rho}{2}\right) \lambda_k^2.$$

As a result, fixing $\rho > 2(C + \delta)$,

$$\mathcal{J}_{S_N}(\boldsymbol{\theta}^{k+1}) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^k) \leq -\tilde{\kappa} \lambda_k^2, \quad \tilde{\kappa} := \left|C + \delta - \frac{\rho}{2}\right| > 0.$$

Moreover, we can rearrange and sum over $k = 0$ to K the above expression to get

$$\sum_{k=0}^K \lambda_k^2 \leq \tilde{\kappa}^{-1} (\mathcal{J}_{S_N}(\boldsymbol{\theta}^0) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^{K+1})) \leq \tilde{\kappa}^{-1} \left(\mathcal{J}_{S_N}(\boldsymbol{\theta}^0) - \inf_{\boldsymbol{\theta} \in L^\infty([0, T], \Theta)} \mathcal{J}_{S_N}(\boldsymbol{\theta}) \right).$$

Therefore $\sum_{k=0}^K \lambda_k^2 < +\infty$ which implies that $\lambda_k^2 \rightarrow_{k \rightarrow \infty} 0$. Therefore, the A-MSA converges to a solution of the extended PMP.

Let us consider now the general case $0 < \gamma_k \leq 1$. From the maximization step (18),

$$\frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \geq \gamma_k \max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \theta), \quad \forall t \in [0, T].$$

As a result, bound (23) receives an additional term and becomes,

$$\begin{aligned} & - \frac{1}{N} \sum_{i=1}^N \int_0^T \Delta H_{\boldsymbol{\theta}^{k+1}, \boldsymbol{\theta}^k}^{i, \varepsilon_k}(t) dt \\ & \leq - \frac{\rho}{2N} \sum_{i=1}^N \int_0^T \left(\| \dot{u}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \|^2 + \| \dot{p}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} - \nabla_u H(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \|^2 \right) dt \\ & + \frac{\rho}{2N} \sum_{i=1}^N \int_0^T \left(\| \dot{u}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} - f(u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) \|^2 + \| \dot{p}_t^{\boldsymbol{\theta}^k, i, \varepsilon_k} - \nabla_u H(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^k) \|^2 \right) dt \\ & + \frac{\gamma_k^{-1} - 1}{N} \left| \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \right|. \end{aligned}$$

It follows that bound (24) receives the same additional term

$$- \frac{1}{N} \sum_{i=1}^N \int_0^T \Delta H_{\boldsymbol{\theta}^{k+1}, \boldsymbol{\theta}^k}^{i, \varepsilon_k}(t) dt \leq - \frac{\rho}{2} \lambda_k^2 + \rho N^{-1/2} \varepsilon_k \lambda_k + \frac{\gamma_k^{-1} - 1}{N} \left| \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \right|.$$

Therefore (25) becomes

$$\mathcal{J}_{S_N}(\boldsymbol{\theta}^{k+1}) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^k) \leq \left(C - \frac{\rho}{2} \right) \lambda_k^2 + \varepsilon_k (4C + \rho N^{-1/2} \lambda_k) + \frac{\gamma_k^{-1} - 1}{N} \left| \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) \right|.$$

Choosing ε_k like before and setting

$$\gamma_k \geq \frac{\left| \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) / N \right|}{\delta \lambda_k^2 + \left| \sum_{i=1}^N \tilde{H}(t, u_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, p_t^{\boldsymbol{\theta}^k, i, \varepsilon_k}, \boldsymbol{\theta}_t^{k+1}) / N \right|}$$

yields

$$\mathcal{J}_{S_N}(\boldsymbol{\theta}^{k+1}) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^k) \leq \left(C + 2\delta - \frac{\rho}{2} \right) \lambda_k^2.$$

We conclude along the same lines as before by fixing now $\rho > 2(C + 2\delta)$, and inferring that

$$\mathcal{J}_{S_N}(\boldsymbol{\theta}^{k+1}) - \mathcal{J}_{S_N}(\boldsymbol{\theta}^k) \leq -\kappa \lambda_k^2, \quad \kappa := \left| C + 2\delta - \frac{\rho}{2} \right| > 0.$$

■

Before going to the next section, a few remarks on Theorem 4 are in order:

- Note that since the sequence of controls (θ^k) built by A-MSA is a minimizing sequence, it need not convergence in $L^\infty([0, T], \Theta)$ norm, nor does it necessarily contain a convergent subsequence. To guarantee that the minimizing sequence converges in norm, it is necessary to additional convexity and/or compactness properties.
- The convergence result is *a priori* in the sense that it cannot be directly implemented for the following reasons:
 - The bounds for (ε_k) and (γ_k) involve a constant C which is difficult to estimate in practice.
 - The bounds may be suboptimal due to the construction of the proof.
 - In practice, it is difficult to guarantee that the maximization step is performed within a certain fraction γ_k of the actual maximum.
- The main interest of the result lies in the fact that it reveals the necessity to tighten the accuracy in the time integration schemes and the maximization algorithms across the iterations in order to approach the continuous sampled optimal control problem (6).

7. Generalization error and convergence towards the fully continuous problem

The goal of this section is to connect the controls (θ^k) of the A-MSA Algorithm 1 with the PMP solutions of the fully continuous problem (4). For this, we rely on recent results from E et al. (2019) connecting the sampled PMP with the fully continuous PMP which we next briefly recall.

We assume in the following that θ^* and $\vartheta_{\mathcal{S}_N}^*$ are solutions of (4) and (6) such that the Hamiltonian step attains a maximum in the *interior* of Θ . Consequently, the continuous PMP solution θ^* satisfies

$$\mathbf{F}(\theta^*)_t := \mathbb{E}_\mu \nabla_\theta H(t, u_t^{\theta^*}, p_t^{\theta^*}, \theta_t^*) = 0$$

for a.e. $t \in [0, T]$ and where $\mathbf{F} : L^\infty([0, T], \Theta) \rightarrow L^\infty([0, T], \mathbb{R}^m)$. Similarly, an interior solution $\vartheta_{\mathcal{S}_N}^*$ of the sampled PMP is a random variable which satisfies

$$\mathbf{F}_{\mathcal{S}_N}(\vartheta_{\mathcal{S}_N}^*)_t := \frac{1}{N} \sum_{i=1}^N \nabla_\theta H(t, u_t^{\vartheta_{\mathcal{S}_N}^*, i}, p_t^{\vartheta_{\mathcal{S}_N}^*, i}, \vartheta_{N,t}) = 0.$$

Under these assumptions, the theorem below, proven in E et al. (2019), describes the convergence of an interior solution $\vartheta_{\mathcal{S}_N}^*$ of the first order condition of the sampled PMP to an interior solution θ^* of the continuous PMP. The additional local strong concavity assumption on the Hessian allows to guarantee that $\vartheta_{\mathcal{S}_N}^*$ is a global/local maximum of the sampled PMP. The result also guarantees convergence of loss function values.

Definition 5 For $\rho > 0$ and $x \in L^\infty([0, T], \Theta)$, define $S_\rho(x) := \{y \in L^\infty([0, T], \Theta) : \|x - y\| \leq \rho\}$. The mapping \mathbf{F} is said to be stable on $S_\rho(x)$ if there exists a constant $K_\rho > 0$ such that for all $y, z \in S_\rho(x)$,

$$\|y - z\|_{L^\infty([0, T], \Theta)} \leq K_\rho \|\mathbf{F}(x) - \mathbf{F}(y)\|_{L^\infty([0, T], \mathbb{R}^m)}.$$

Theorem 6 (Theorem 6, Corollary 1 and 2 from E et al. (2019)) *Let \mathbf{F} be a mapping which is stable on $S_\rho(\boldsymbol{\theta}^*)$ for some $\rho > 0$ and $\boldsymbol{\theta}^*$ a solution of $\mathbf{F} = 0$. Then there exists positive constants s_0, C, K_1, K_2 and $\rho_1 < \rho$ and a random variable $\boldsymbol{\vartheta}_{\mathcal{S}_N}^* \in S_{\rho_1}(\boldsymbol{\theta}^*)$ such that*

$$\mathbb{P}[\|\boldsymbol{\theta}^* - \boldsymbol{\vartheta}_{\mathcal{S}_N}^*\| \geq Cs] \leq 4 \exp\left(-\frac{Ns^2}{K_1 + K_2s}\right), \quad s \in (0, s_0]$$

and $\boldsymbol{\vartheta}_{\mathcal{S}_N}^* \rightarrow \boldsymbol{\theta}^*$ as $N \rightarrow \infty$ in probability. Moreover, there exists constants K'_1, K'_2 such that

$$\mathbb{P}[|\mathcal{J}(\boldsymbol{\theta}^*) - \mathcal{J}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*)| \geq s] \leq 4 \exp\left(-\frac{Ns^2}{K'_1 + K'_2s}\right), \quad s \in (0, s_0] \quad (26)$$

In addition, if for all $t \in [0, T]$, the Hessian $\mathbb{E}_\mu \nabla_{\theta, \theta}^2 H(u_t^{\boldsymbol{\theta}^*}, p_t^{\boldsymbol{\theta}^*}, \boldsymbol{\theta}_t^*) + \lambda_0 I$ is negative definite, then $\boldsymbol{\vartheta}_{\mathcal{S}_N}^*$ is also a strict local maximum of the sampled Hamiltonian $\vartheta \mapsto \frac{1}{N} \sum_{i=1}^N H(u_t^{\boldsymbol{\vartheta}_{\mathcal{S}_N}^*, i}, p_t^{\boldsymbol{\vartheta}_{\mathcal{S}_N}^*, i}, \vartheta)$. In particular, if the sampled Hamiltonian has a unique maximizer, then $\boldsymbol{\vartheta}_{\mathcal{S}_N}^*$ is a solution of the sampled PMP (8) with the same high probability.

Corollary 7 *If the sampled Hamiltonian has a unique maximizer $\boldsymbol{\vartheta}_{\mathcal{S}_N}^*$, then there exists constants \tilde{K}_1 and \tilde{K}_2 such that for any $\varepsilon > 0$, there exists $k(\varepsilon) \in \mathbb{N}$ such that for $k \geq k(\varepsilon)$, the control $\boldsymbol{\theta}^k$ at iteration k of A-MSA satisfies*

$$\mathbb{P}[|\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\theta}^k) - \mathcal{J}(\boldsymbol{\theta}^*)| \geq 3\varepsilon] \leq 4 \exp\left(-\frac{N\varepsilon^2}{\tilde{K}_1 + \tilde{K}_2\varepsilon}\right)$$

Proof We use that

$$\begin{aligned} \mathbb{P}[|\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\theta}^k) - \mathcal{J}(\boldsymbol{\theta}^*)| \geq 3\varepsilon] &\leq \mathbb{P}[|\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\theta}^k) - \mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*)| \geq \varepsilon] + \mathbb{P}[|\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*) - \mathcal{J}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*)| \geq \varepsilon] \\ &\quad + \mathbb{P}[|\mathcal{J}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*) - \mathcal{J}(\boldsymbol{\theta}^*)| \geq \varepsilon], \end{aligned}$$

and then bound each term as explained next. Let $\varepsilon > 0$ be fixed. With high probability, $\boldsymbol{\vartheta}_{\mathcal{S}_N}^*$ is a solution to the sampled PMP (6), that is $\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*) = \mathcal{J}_{\mathcal{S}_N}^*$. By Theorem 4, there exists $k(\varepsilon) \in \mathbb{N}$ such that for $k \geq k(\varepsilon)$, the control $\boldsymbol{\theta}^k$ at iteration k of A-MSA satisfies $|\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\theta}^k) - \mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*)| \leq \varepsilon$, which bounds the first term. For the second term, by the infinite dimensional Hoeffding's inequality (see (Pinelis and Sakhnenko, 1986, Corollary 2)), $|\mathcal{J}_{\mathcal{S}_N}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*) - \mathcal{J}(\boldsymbol{\vartheta}_{\mathcal{S}_N}^*)| \leq \varepsilon$ with probability $1 - \exp(-N\varepsilon^2/(\tilde{K}'_1 + \tilde{K}'_2\varepsilon))$. The third term is also bounded by ε with probability $4 \exp\left(-\frac{N\varepsilon^2}{K'_1 + K'_2\varepsilon}\right)$ thanks to (26). \blacksquare

8. Guidelines for numerical implementation

Algorithm 1 should in principle be implemented with tolerances ε_k and γ_k given in Theorem 6 to ensure convergence. These tolerances are in practice hard to guarantee, and the convergence analysis that leads to these quantities may be suboptimal. A trade-off between the theory and numerical implementation consists simply in replacing the forward and backward

propagations at prescribed accuracies by propagations where the time steps are progressively refined across the iterations. It is therefore necessary to prescribe a refinement strategy. For the Hamiltonian maximization step (line 13 of Algorithm 1), the best that one can do is to use an efficient optimizer of nonconvex problems.

In addition to these considerations, note that Algorithm 1 can be parallelized at several parts: for each sample i , the forward and backward propagations can be performed in parallel (see the for loop of line 8). Each propagation can additionally be parallelized by means of parallel in time algorithms as in Günther et al. (2020) (see Lions et al. (2001); Maday and Mula (2020) for some selected references on the topic). In addition, the Hamiltonian maximization step can also be parallelized (see for loop of line 12 of Algorithm 1). For our concern, we have chosen to parallelize the Maximization step, since it is way more expensive in terms of computational time than the propagation step.

9. Numerical experiments

In this section, we present some numerical experiments, aiming primarily at illustrating the behavior of the A-MSA algorithm. We give special focus on studying the performance in the learning phase in terms of the value of the loss function. We also discuss computing times and examine expressivity in terms of generalization errors. Note that the final quality of approximation depends on the number of samples, on the network depth, and on the optimizers and their starting guesses. We proceed by increasing levels of difficulty in the numerical examples in order to disentangle the effect of each of the different aspects (note however that their effects are usually mixed in non-trivial applications). The code to reproduce the results can be found at

<https://github.com/jaghili/amsa>.

The implementation has been done with Python 3 and NumPy. Particular attention has been paid to vectorize most of the operations in order to delegate as much as possible the looping to internal, highly optimized C and Fortran functions. We have observed that the time to perform the ODE propagations is negligible with respect to the Hamiltonian maximization (line 12 of Algorithm 1). We have thus parallelized this step with MPI.

9.1 Approximation of the sine function

In this example, inspired from (Li et al., 2018, Section 6), the task is to approximate the graph of the sine function,

$$y : X = [-\pi, \pi] \rightarrow Y = [-1, 1], \quad x \mapsto y(x) = \sin(x).$$

Note that here the dimensions n and k of the domain and label sets are $n = k = 1$. The continuous learning problem (see (4)) is to minimize the $L^2(X)$ approximation error over the controls $\theta \in L^\infty([0, T], \Theta)$, namely

$$\inf_{\theta \in L^\infty([0, T], \Theta)} \frac{1}{2} \int_X |y(x) - g(u_T^{\theta, x})|^2 dx.$$

For every point $x \in X$, $u_T^{\theta,x}$ is the final time state of the ODE

$$\begin{cases} \dot{u}_t^{\theta,x} &= \tanh(A_t \cdot u_t^{\theta,x} + b_t), \quad \forall t \in (0, T) \\ u_0^{\theta,x} &= x(1, \dots, 1)^T \in \mathbb{R}^d. \end{cases}$$

Here, $\theta_t = (A_t, b_t) \in \mathbb{R}^{d \times d} \times \mathbb{R}^d \sim \mathbb{R}^{d^2+d}$ are the parameters to optimize. We constraint them to lie in $[\alpha_{\min}, \alpha_{\max}] = [-1, 1]$ so $\Theta = [-1, 1]^{d^2+d}$. The function g is defined as

$$g: \mathbb{R}^d \rightarrow \mathbb{R}, \quad z = (z_1, \dots, z_d) \mapsto g(z) := \frac{1}{d} \sum_{i=1}^d z_i,$$

and note that we have chosen $\mu = dx$ as the Lebesgue measure, $f = \tanh$ as the activation function, $\Phi(u_T^{\theta,x}, y) = \frac{1}{2}|y(x) - g(u_T^{\theta,x})|^2$ and we do not have any regularization ($R = 0$).

In the sampled version which we consider in our experiments, we work with a set $\mathcal{S}_N = \{(x_i, y_i)\}_{i=1}^N$ of pairs generated with equidistant data points

$$x_i = -\pi + 2\pi \frac{i-1}{N-1}, \quad y_i = \sin(x_i), \quad i \in \{1, \dots, N\},$$

and perform empirical risk minimization taking the uniform distribution $\mu_N = \frac{1}{N} \sum_{i=1}^N \delta_{(x_i, y_i)}$ (see (6)). We solve this problem with the above described A-MSA algorithm where the final time is set to $T = 5$, the penalization parameter is set to $\rho = 5$, and the number of neurons per layer to $d = 3$. We use a basic explicit Euler scheme as a time integrator, which induces a ResNet architecture. We progressively refine the time discretization from $L = 3$ to $L = 32$ time steps (in other words, we increase the network depth from 3 to 32 layers) using three different strategies:

- (A1) *Abrupt refinement*: We refine from shallow ($L = 3$) to deep ($L = 32$) at iteration $k = 250$.
- (A2) *Fast refinement*: We add 10 layers every 50 iterations.
- (A3) *Slow refinement*: We add 10 layers every 100 iterations.

Our goal is to compare the behavior of these refinement strategies with the non-adaptive training of a shallow and a deep network having $L = 3$ and $L = 32$ layers respectively.

At each iteration $k \geq 0$ of A-MSA, the maximization step is performed with NumPy's L-BFGS-B optimizer, which is an iterative algorithm suited for non-convex optimization problems. The optimizer is initialized by picking the best candidate among a list Θ^0 of vectors, including the last best control $\theta_{\text{best}}^k = \arg \min_{0 \leq i < k} \mathcal{J}_{\mathcal{S}_N}(\theta^i)$ among all previous iterations before k and random perturbations of it, precisely

$$\Theta^0 := \bigcup_{q=0}^5 \bigcup_{i=1}^{25} \left\{ \theta_{\text{best}}^k + 10^{-2q} \mathbf{rand}(\alpha_{\min}, \alpha_{\max}), 10^{-2q} \mathbf{rand}(\alpha_{\min}, \alpha_{\max}) \right\}$$

where $\mathbf{rand}(\alpha_{\min}, \alpha_{\max})$ is a uniform real distribution function in the interval $[\alpha_{\min}, \alpha_{\max}]$. As a result, the final output of the algorithm is random due to the choice of the initial

guess. For this reason, we make $\bar{r} = 20$ runs of A-MSA, each one consisting in $k_{\max} = 800$ iterations.

In Figure 1, we first fix $N = 20$ samples and we plot the convergence history of the train loss $\mathcal{J}_{\mathcal{S}_N}^{(\text{run}=r)}(\boldsymbol{\theta}^k)$ (see red colors). We also plot the generalization error via the *test loss* across the iterations k (see blue colors). The test loss is the quantity $\mathcal{J}_{\tilde{\mathcal{S}}_N}(\boldsymbol{\theta}^k)$, defined in (6), where we use a test set $\tilde{\mathcal{S}}_N$ different from the training set \mathcal{S}_N . $\tilde{\mathcal{S}}_N$ is taken uniformly random at each iteration $1 \leq k \leq k_{\max}$. Since the history depends on each run r , the figure shows some statistics related to the repetitions. The continuous red curve shows the average value $\mathcal{J}_{\mathcal{S}_N, k}^{(\text{av})} := (1/\bar{r}) \sum_{r=1}^{\bar{r}} \mathcal{J}_{\mathcal{S}_N}^{(\text{run}=r)}(\boldsymbol{\theta}^k)$ over the runs and the red diffuse color shows its distribution around the average. The same presentation is given in blue for the generalization errors. We can first observe that the convergence history is, in average, relatively similar for all architectures. The average value of the loss function does not go below 10^{-2} for the shallow network (figure 1a), and the value reached by the deep network is only slightly better (figure 1b): at the end of the iterations, the cost function is in average below $7.1 \cdot 10^{-2}$. Similar values are reached by any of our three adaptive strategies (see figures 1d, 1e and 1f). As a consequence, if we take this average value as the performance indicator, we may conclude that the approximation power of the deep network is only marginally better than the shallow one. We may also conclude that the adaptive training strategy reduces to some extent the computational time for training (as we illustrate further on in section 9.4) but it does not bring extra approximation power. However, note that this indicator is not entirely appropriate because it does not correspond to any realized convergence history and, more importantly, because it does not inform about the best performance that we have at hand. It is crucial to remark that the variance around the average convergence value is significantly larger in the adaptive strategies compared to the non adaptive ones. This is particularly true for (A2) and (A3). Figure 1c shows that the minimal value $\min_{1 \leq k \leq k_{\max}} \mathcal{J}_{\mathcal{S}_N}^{(\text{run}=r)}(\boldsymbol{\theta}^k)$ reached during the $k_{\max} = 800$ iterations of each run is, in average, one order of magnitude lower than the one of the non-adaptive coarse and the non-adaptive deep network. This result illustrates that the adaptive strategy significantly contributes to reach better minima, thus allowing to benefit better from the higher approximation power of the deeper neural network.

Figure 2 shows the reconstruction of the sine function with the controls performing best in each training run, namely, realizing $\min_{1 \leq k \leq k_{\max}} \mathcal{J}_{\mathcal{S}_N}^{(\text{run}=r)}(\boldsymbol{\theta}^k)$ for every r . We see that the best over all realizations yields a very satisfactory approximation. Statistics on the generalization errors of these best controls are given in Figure 2c

We can next examine in Figure 3 the impact of the number N of data samples in the generalization errors of the best controls. As one can expect, the generalization errors decrease when the number N of samples increases. For $N \geq 20$, the errors with the shallow network stagnate before reaching 10^{-3} in average regardless of the number of samples. This illustrates the limitations in the approximation power of the shallow network in the current example. The non-adaptive deep network ($L = 32$) presents a behavior which is only better for a reduced number of samples $N \leq 8$ but overall its performance is relatively similar to the coarse neural network. This comes from the difficulty in finding good quality optimizers in networks involving many coefficients. If we now adaptively train the deep network with strategy (A2), the generalization errors are in general lower than the ones reached by the

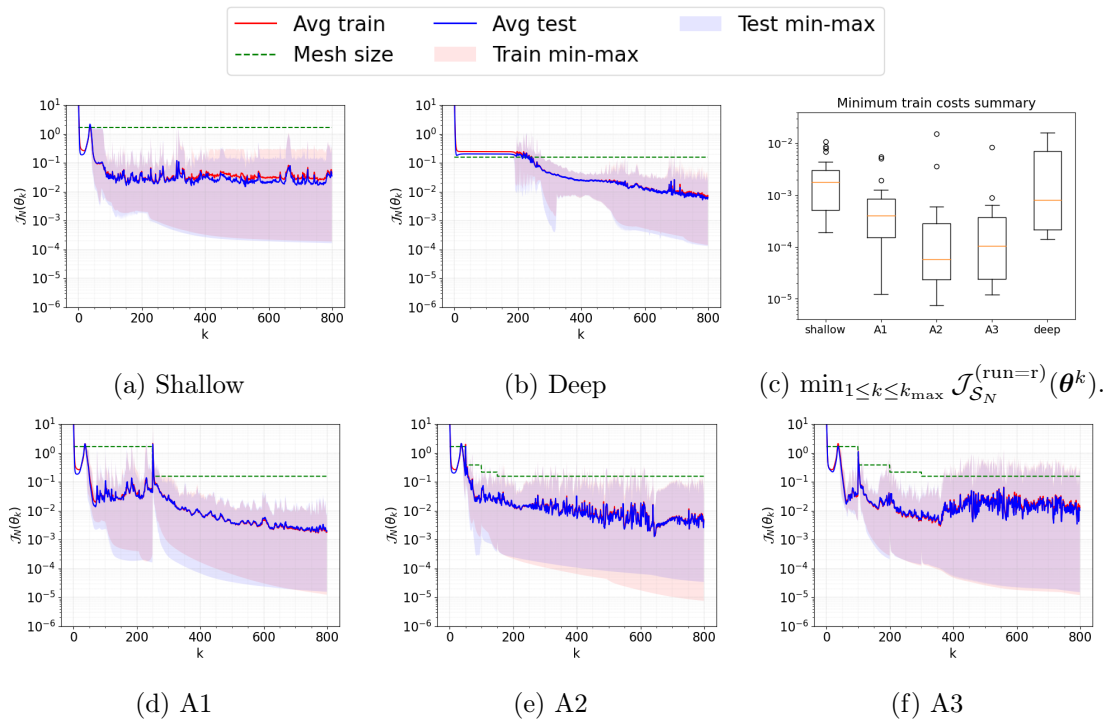


Figure 1: Sine function. Training loss $\mathcal{J}_{S_N}^{(\text{run}=\text{r})}(\theta_k)$ (red) and test loss $\mathcal{J}_{\tilde{S}_N}^{(\text{run}=\text{r})}(\theta_k)$ (blue). Here $N = 20$ and $\bar{r} = 20$ runs.

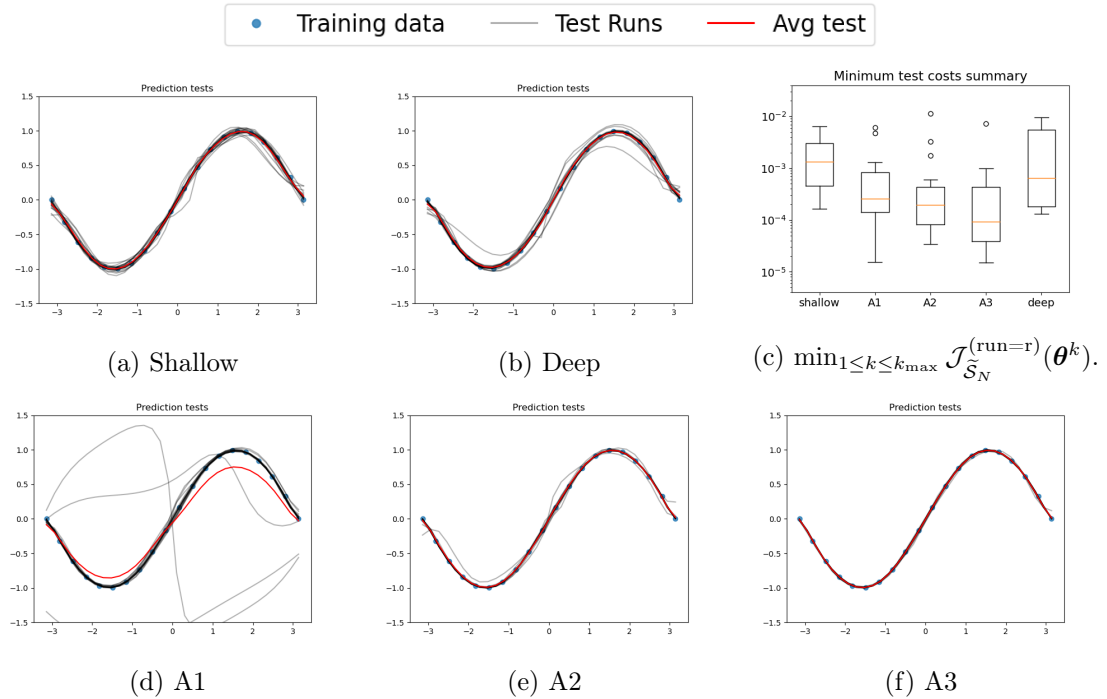
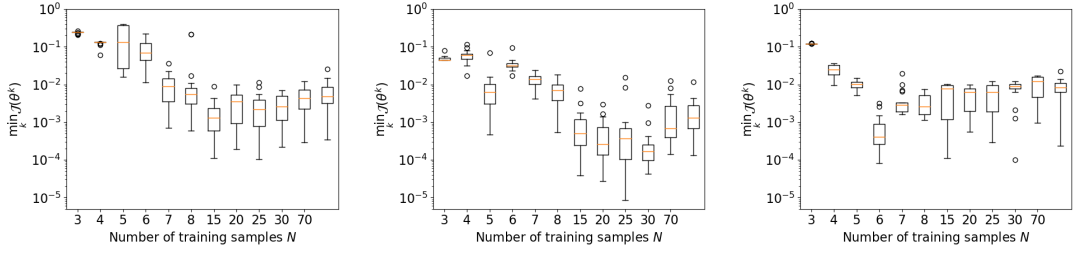


Figure 2: Sine function: Predictions of $\bar{r} = 20$ runs and average prediction.



(a) Shallow: 3 layers

(b) A2: 3 \rightarrow 32 layers.

(c) Deep: 32 layers

Figure 3: Sine case – Generalization errors vs N for coarse, adaptive (A2) and thin networks

non-adaptive strategies, especially when $N \geq 15$. The result illustrates that the adaptive strategy gives in average a better performance in terms of approximation. This is obtained at a reduced computing time as we show further on.

9.2 Noisy step function

We next examine the robustness of the method against noise. We consider the graph of the step function,

$$y : X = [-1, 1] \rightarrow Y = [-0.5, 0.5], \quad x \mapsto y(x) = \begin{cases} 0.5 & \text{if } x \leq 0 \\ -0.5 & \text{if } x > 0 \end{cases}$$

and we seek to approximate it when the given data are noisy. For each sample i , we get the pair (x_i, y_i) with

$$y_i = y(x_i) + \varepsilon_i, \quad i = 1, \dots, N,$$

and the ε_i follow a uniform distribution in $[-0.2, 0.2]$.

We consider the same setting as before ($T = 5$, $\rho = 5$, $d = 3$) and fix $N = 800$. Figure 4 shows the convergence history of the training with the shallow, deep and adaptively deep neural networks. We consider the same adaptive strategies as before. All approaches reach more or less the same minimal value in the training (see Figure 4c). However, we observe that (A2) and (A3) give better generalization errors than the non-adaptive deep architecture (see Figure 5c). Note however that this superiority is not very large and all methods produce excellent reconstruction results as Figure 5 illustrates. We think that this is due to the fact that in this case all methods find configurations that are very close the global optimum, which is here given by the approximation error of the mapping $x \mapsto y(x)$.

9.3 Classification

As a last example, we consider a simple 2D classification problem where the function to approximate is

$$y : X = [-1, 1]^2 \rightarrow Y = \{0, 1\}, \quad (x_1, x_2) \mapsto y(x_1, x_2) = \begin{cases} 1 & \text{if } x_1^2 + x_2^2 \leq (0.5)^2, \\ 0 & \text{otherwise.} \end{cases}$$

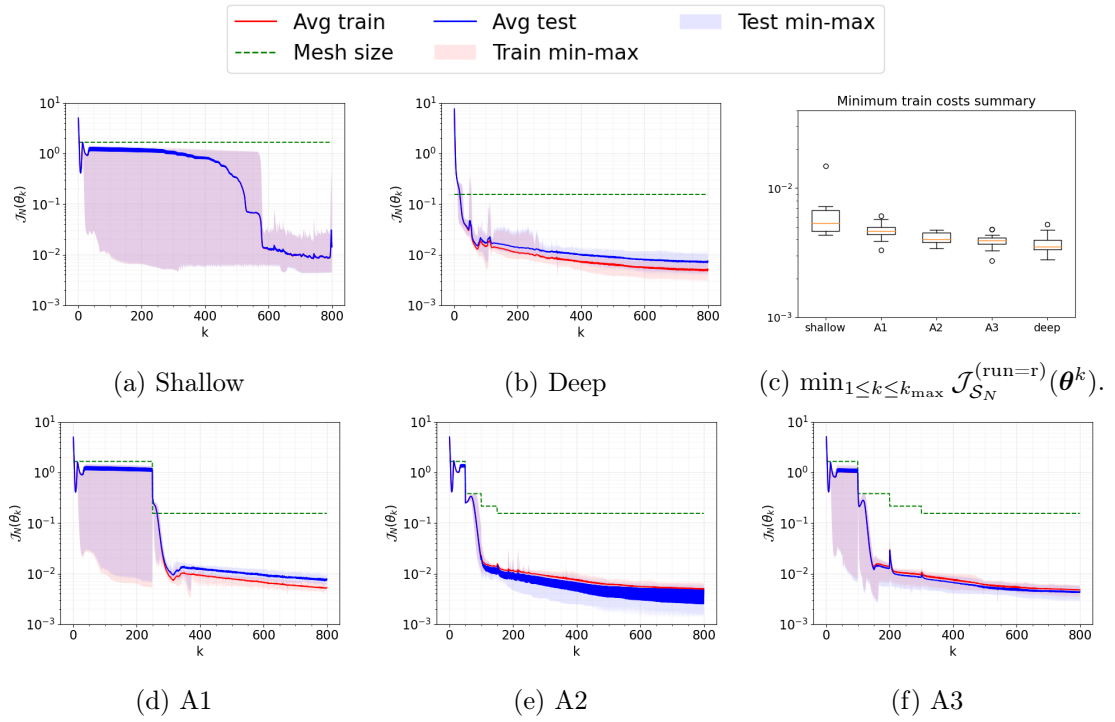


Figure 4: Noisy step: Loss $\mathcal{J}_N^{(\text{run}=\bar{r})}(\theta_k)$ for $N = 800$ and $\bar{r} = 20$ runs.

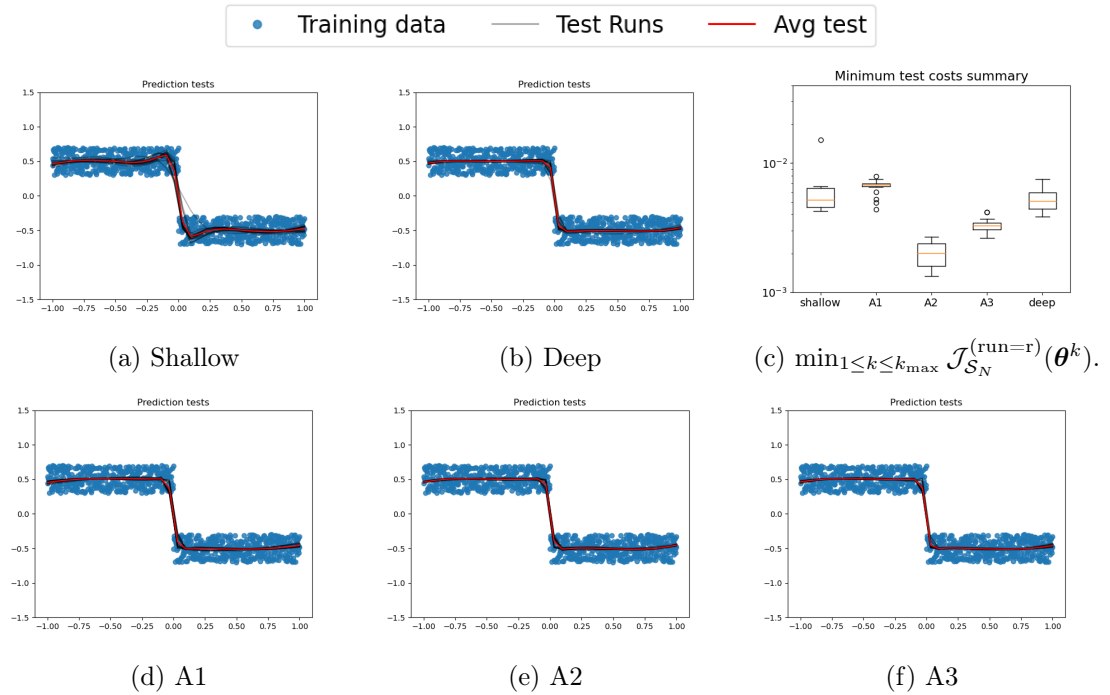


Figure 5: Noisy step: Predictions of some runs.

The setting is the same as in the above examples but with slight changes. Each layer has now $d = 2 \times 3 = 6$ neurons per layers, the constraint on the controls is set to $[\alpha_{\min}, \alpha_{\max}] = [-2, 2]$ and the output function g is defined as

$$g : \mathbb{R}^d \rightarrow \mathbb{R}, \quad z = (z_1, \dots, z_d) \mapsto g(z) := \mathcal{H} \left(\frac{1}{d} \sum_{i=1}^d z_i \right)$$

where \mathcal{H} is the usual Heaviside step function. The figure layout of the results is the same as in the other examples. Figures 6 and 7 show statistics on the loss function and examples of reconstruction. The blue-to-red scalar field in Fig 5 is the output of the trained network composed with a treshold filter $\mathcal{F} : x \mapsto \mathcal{H}(x - 0.5)$ from 1024 inputs points. The black crosses (\times) and white dots (\circ) are the $N = 800$ training points used for each case, the first refers to points with $z = 0$, the second to $z = 1$. Essentially, we observe that A-MSA yields the same quality of approximation as the non-adaptive deep neural network.

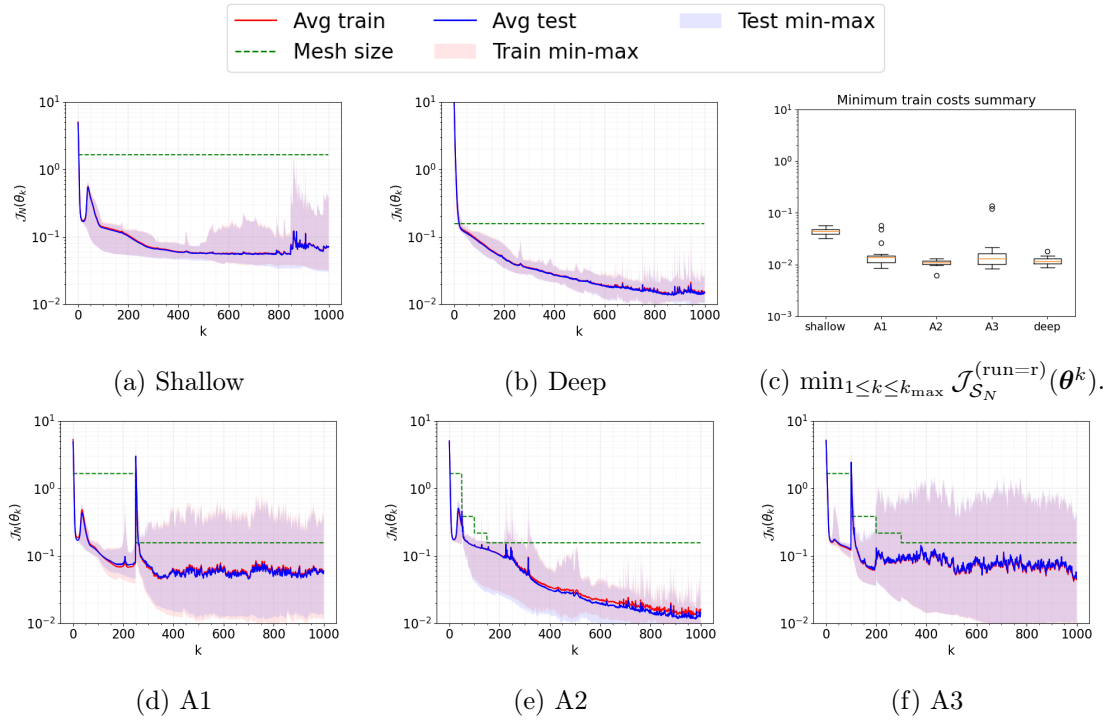


Figure 6: Classification: Loss $\mathcal{J}_N^{(\text{run}=r)}(\theta_k)$ for $N = 800$ and 20 runs.

9.4 Computational times

We finish the section on numerical tests by examining the computational time in the learning approach. We present runtimes only for the noisy step example for the sake of brevity (similar results are obtained for the other tests). We consider three different criteria:

- **Complexity index (see Figure 8a):** we estimate the number of operations by measuring the runtime of a sequential run.

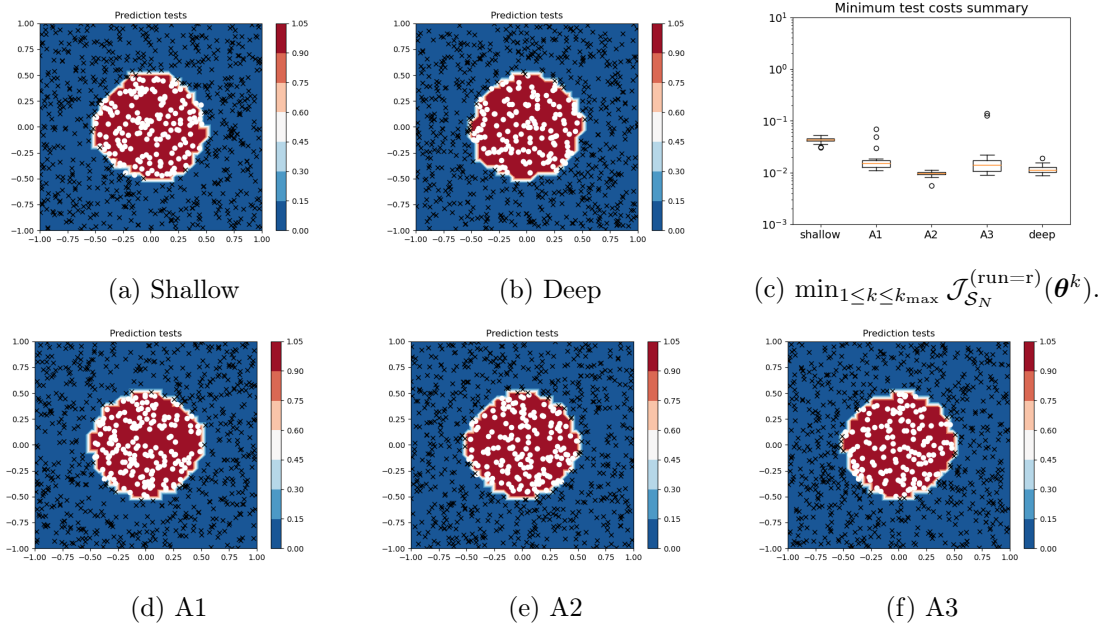


Figure 7: Classification: Predictions of $\bar{r} = 20$ runs and average prediction. Data labels: $\times : 0$, $\circ : 1$.

- **Parallel runtime index (see Figure 8b):** we measure the runtime of running in parallel the Hamiltonian maximization step of the learning algorithm. We place ourselves in a scenario with no constraints in the computing resources and allocate one processor per layer. Note that the number of layers is increased in the adaptive approach so we use more and more resources as the algorithm makes refinements.
- **Energy consumption index (see Figure 8a):** we estimate the total energy consumption of a parallel run as follows. For each iteration of the algorithm, we add the corresponding computing times of all the processors involved in that iteration. This gives an estimate in `cpu-seconds` of the energy consumption for each iteration. By summing over all iterations, we obtain the final estimate. This quantity is in fact equal to the cumulative sequential runtime which is given on Figure 8a.

Our runtimes have been measured on a cluster and we have observed that its occupation greatly impacts on the runtime results. To illustrate this issue, we show runtimes for the training of the shallow neural network in an empty node and a busy node, i.e. a node where others jobs are running in parallel. We have run the rest of the examples on empty nodes.

If we first consider Figure 8a, we see that the adaptive training strategy performs better than the non-adaptive one regarding the sequential runtimes and the energy consumption. This is comes as no surprise since the adaptive strategy performs less operations at the beginning of the iterations since they involve less layers. When computations are run in parallel, Figure 8b) shows that the run times are all very similar for all strategies. This is due to the fact that the tasks are well-balanced between processors since each of them

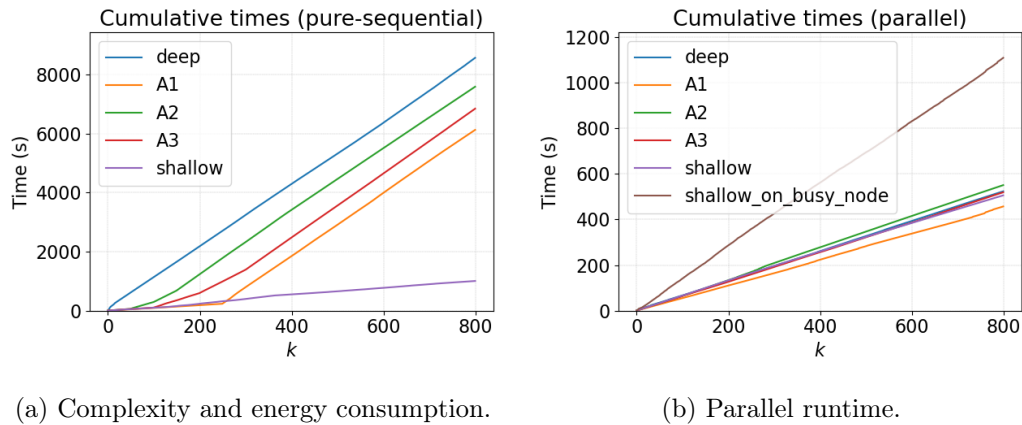


Figure 8: Computational times for the noisy step example.

does computations for one layer. We also observe that runtimes are greatly affected by the occupation of the cluster.

10. Conclusion

We have proven that the A-MSA algorithm converges to an underlying continuous learning problem in the limit of the time discretization and of the the number of samples. The convergence analysis requires that the time propagations and Hamiltonian maximizations are performed at each step with increasingly tight accuracy tolerances. Since a sharp estimation of these tolerances is difficult to obtain, we have implemented A-MSA with certain refinement strategies in order to examine its potential. The numerical experiments reveal that the adaptive strategy helps to benefit in practice from the higher approximation properties of deep networks by mitigating over-parametrization issues: our results show that adaptivity increases the chances to find better quality minimizers compared to the non-adaptive training of deep neural networks. In addition, it appears that the adaptive strategy is clearly more performant in terms of complexity and energy consumption compared to the non-adaptive training of a deep neural network.

Acknowledgments and Disclosure of Funding

This research has been funded by the Emergences Project of the Paris city council called “Models and Measures”.

References

G. Akrivis, C. Makridakis, , and R. Nochetto. Optimal order a posteriori error estimates for a class of runge–kutta and galerkin methods. *Numerische Mathematik*, 114(1):133, Aug 2009. ISSN 0945-3245. doi: 10.1007/s00211-009-0254-2. URL <https://doi.org/10.1007/s00211-009-0254-2>.

- M. Athans and P. L. Falb. *Optimal control: an introduction to the theory and its applications*. Courier Corporation, 2013.
- M. Benning, E. Celledoni, M. J. Ehrhardt, B. Owren, and C. Schönlieb. Deep learning as optimal control problems: Models and numerical methods. *Journal of Computational Dynamics*, 6(2):171–198, 2019. ISSN 2158-2491. doi: 10.3934/jcd.2019009. URL <http://aimsciences.org//article/id/37d95dff-e7a8-4aee-9ee6-cf35d0deb544>.
- D. P. Bertsekas. *Dynamic programming and optimal control*, volume 1. Athena scientific Belmont, MA, 1995.
- V. G. Boltyanskii, R. V. Gamkrelidze, and L. S. Pontryagin. The theory of optimal processes. i. the maximum principle. Technical report, TRW SPACE TECHNOLOGY LABS LOS ANGELES CALIF, 1960.
- L. Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010.
- A. Bressan and B. Piccoli. *Introduction to the mathematical theory of control*, volume 1. American institute of mathematical sciences Springfield, 2007.
- T. Chen, I. Goodfellow, and J. Shlens. Net2net: Accelerating learning via knowledge transfer. *arXiv preprint arXiv:1511.05641*, 2015.
- F. L. Chernousko and A. A. Lyubushin. Method of successive approximations for solution of optimal control problems. *Optimal Control Applications and Methods*, 3(2):101–114, 1982.
- F. Clarke. The maximum principle in optimal control, then and now. *Control and Cybernetics*, 34(3):709, 2005.
- C. Cortes, X. Gonzalvo, V. Kuznetsov, M. Mohri, and S. Yang. Adanet: Adaptive structural learning of artificial neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 874–883. JMLR. org, 2017.
- I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova. Nonlinear approximation and (deep) relu networks. *arXiv preprint arXiv:1905.02199*, 2019.
- J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(Jul):2121–2159, 2011.
- W. E, J. Han, and Q. Li. A mean-field optimal control formulation of deep learning. *Research in the Mathematical Sciences*, 6(1):10, 2019.
- T. Elsken, J. H. Metzen, and F. Hutter. Neural architecture search: A survey. *arXiv preprint arXiv:1808.05377*, 2018.
- A. Gomez, M. Ren, R. Urtasun, and R. Grosse. The reversible residual network: Backpropagation without storing activations. In *Advances in neural information processing systems*, pages 2214–2224, 2017.

- P. Grohs, D. Perekrestenko, D. Elbrächter, and H. Bölcskei. Deep neural network approximation theory. *arXiv preprint arXiv:1901.02220*, 2019.
- I. Gühring, G. Kutyniok, and P. Petersen. Error bounds for approximations with deep ReLU neural networks in $w^{s,p}$ norms. *arXiv preprint arXiv:1902.07896*, 2019.
- S. Günther, L. Ruthotto, J. B. Schroder, E. C. Cyr, and N. R. Gauger. Layer-parallel training of deep residual neural networks. *SIAM Journal on Mathematics of Data Science*, 2(1): 1–23, 2020.
- E. Haber and L. Ruthotto. Stable architectures for deep neural networks. *Inverse Problems*, 34(1):014004, 2018.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- D. Hebb. *The organization of behavior: a neuropsychological theory*. Wiley, 1949. ISBN 0-8058-4300-0.
- A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- G. Larsson, M. Maire, and G. Shakhnarovich. Fractalnet: Ultra-deep neural networks without residuals. *arXiv preprint arXiv:1605.07648*, 2016.
- Y. LeCun, D. Touresky, G. Hinton, and T. Sejnowski. A theoretical framework for back-propagation. In *Proceedings of the 1988 connectionist models summer school*, volume 1, pages 21–28. CMU, Pittsburgh, Pa: Morgan Kaufmann, 1988.
- Q. Li, L. Chen, C. Tai, and W. E. Maximum principle based algorithms for deep learning. *The Journal of Machine Learning Research*, 18(1):5998–6026, 2018.
- Z. Li and D. Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- J. Lions, Y. Maday, and G. Turinici. Résolution d’EDP par un schéma en temps pararéel. *C. R. Acad. Sci. Paris*, 2001. t. 332, Série I, p. 661-668.
- Y. Lu, A. Zhong, Q. Li, and B. Dong. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3282–3291, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL <http://proceedings.mlr.press/v80/lu18d.html>.
- Y. Maday and O. Mula. An adaptive parareal algorithm. *To appear in the Journal of Computational and Applied Mathematics*, 2020. URL <https://arxiv.org/abs/1909.08333>.

- F. J. Pineda. Generalization of back propagation to recurrent and higher order neural networks. In *Neural information processing systems*, pages 602–611, 1988.
- I. F. Pinelis and A. I. Sakhanenko. Remarks on inequalities for large deviation probabilities. *Theory of Probability & Its Applications*, 30(1):143–148, 1986.
- L. S. Pontryagin. *Mathematical theory of optimal processes*. Routledge, 2018.
- H. Robbins and S. Monro. A stochastic approximation method. *The annals of mathematical statistics*, pages 400–407, 1951.
- F. Rosenblatt. The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, pages 65–386, 1958.
- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- M. Telgarsky. Representation benefits of deep feedforward networks. *arXiv preprint arXiv:1509.08101*, 2015.
- F. Vialard, R. Kwitt, S. Wei, and M. Niethammer. A shooting formulation of deep learning, 2020.
- T. Wei, C. Wang, Y. Rui, and C. W. Chen. Network morphism. In *International Conference on Machine Learning*, pages 564–572, 2016.
- E. Weinan. A proposal on machine learning via dynamical systems. *Communications in Mathematics and Statistics*, 5(1):1–11, 2017. doi: 10.1007/s40304-017-0103-z.
- Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey, et al. Google’s neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144*, 2016.
- D. Yarotsky. Error bounds for approximations with deep ReLU networks. *Neural Networks*, 94:103–114, 2017.
- M. D. Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.
- X. Zhang, Z. Li, C. Change Loy, and D. Lin. Polynet: A pursuit of structural diversity in very deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 718–726, 2017.

Appendix A. Proof of Lemma 3

The proof of Lemma 3 makes use of certain bounds which we gather in the following Proposition 8. They are derived by an immediate extension of the proofs given in Li et al. (2018) so we just sketch the proof for self-completeness of the current work.

Proposition 8 We have the following bounds for all $t \in [0, T]$,

$$\|p_t^{\theta, \zeta}\| \leq K' := Ke^{\|\Pi_\zeta\|KT} \quad (27)$$

$$\|\delta u_t\| \leq K_{\delta u} := \left(\sqrt{T}(\eta + \zeta) + \int_0^T \|f(s, u_s^{\theta, \zeta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \right) e^{KT} \quad (28)$$

$$\begin{aligned} \|\delta p_t\| &\leq K_{\delta p} := \sqrt{T}e^{2KT}(1 + K(1 + K'T))(\eta + \zeta) \\ &\quad + e^{2KT}K(1 + K'T) \int_0^T \|f(s, u_s^{\theta, \zeta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \\ &\quad + e^{KT} \int_0^T \|\nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \varphi_s) - \nabla_u H(s, u_s^{\theta, \eta}, p_s^{\theta, \eta}, \theta_s)\| ds \end{aligned} \quad (29)$$

Proof [Proof of Proposition 8] Inequality (27) follows from the fact that

$$\dot{p}_t^{\theta, \zeta} = \Pi_\zeta \left(-p_t^{\theta, \zeta} \cdot \nabla_u f(u_t^{\theta, \zeta}, \theta_t) \right)$$

thus taking the scalar product with $p_t^{\theta, \zeta}$ and using the Cauchy-Schwarz inequality, we derive the evolution

$$\frac{1}{2} \frac{d}{dt} \|p_t^{\theta, \zeta}\|^2 \leq \|\Pi_\zeta\| \|\nabla_u f(u_t^{\theta, \zeta}, \theta_t)\|_2 \|p_t^{\theta, \zeta}\|^2, \quad \|p_T^{\theta, \zeta}\|^2 = \|\nabla_u \Phi(u_T^{\theta, \zeta}, y)\|^2$$

and the result follows by applying the Grönwall's inequality and the fact that $\|\nabla_u \Phi(u_T^{\theta, \zeta}, y)\| \leq K$ and $\|\nabla_u f(u_t^{\theta, \zeta}, \theta_t)\|_2 \leq K$ by hypothesis (A1) and (A2).

To derive inequality (28), we start from

$$\delta \dot{u}_t = \dot{u}_t^{\varphi, \eta} - \dot{u}_t^{\theta, \zeta} = \Pi_\eta f(u_t^{\varphi, \eta}, \varphi_t) - \Pi_\zeta f(u_t^{\theta, \zeta}, \theta_t).$$

Integrating between 0 and $t \leq T$,

$$\delta u_t = \int_0^t \Pi_\eta f(s, u_s^{\varphi, \eta}, \varphi_s) - \Pi_\zeta f(s, u_s^{\theta, \zeta}, \theta_s) ds$$

Taking norms,

$$\begin{aligned} \|\delta u_t\| &\leq \int_0^t \|f(s, u_s^{\varphi, \eta}, \varphi_s) - \Pi_\eta f(s, u_s^{\varphi, \eta}, \varphi_s)\| ds + \int_0^t \|f(s, u_s^{\theta, \zeta}, \theta_s) - \Pi_\zeta f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \\ &\quad + \int_0^t \|f(s, u_s^{\varphi, \eta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \\ &\leq \sqrt{T}(\eta + \zeta) + \int_0^t \|f(s, u_s^{\varphi, \eta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \varphi_s)\| ds + \int_0^T \|f(s, u_s^{\theta, \zeta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \\ &\leq \sqrt{T}(\eta + \zeta) + K \int_0^t \|\delta u_s\| ds + \int_0^T \|f(s, u_s^{\theta, \zeta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \end{aligned}$$

Thus by Grönwall's inequality and $t \leq T$,

$$\|\delta u_t\| \leq K_{\delta u} := \left(\sqrt{T}(\eta + \zeta) + \int_0^T \|f(s, u_s^{\theta, \zeta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \right) e^{KT}$$

Proceeding in a similar manner with $\delta \dot{p}_t$, by integrating between T and t , we derive

$$\int_T^t \delta \dot{p}_s ds = \delta p_t - \delta p_T = \int_T^t \Pi_\eta \nabla_u H(s, u_s^{\varphi, \eta}, p_s^{\varphi, \eta}, \varphi_s) - \Pi_\zeta \nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \theta_s) ds$$

and similarly as before, we infer that

$$\begin{aligned} \|\delta p_t\| &\leq \sqrt{T}(\eta + \zeta) + K\|\delta u_T\| + KK' \int_0^T \|\delta u_s\| + K \int_0^T \|\delta p_s\| ds \\ &\quad + \int_0^T \|\nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \varphi_s) - \nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \theta_s)\| ds \end{aligned}$$

Using bounds (28) and (29), and applying the Gronwall inequality and $t \leq T$, we derive the announced estimate

$$\|\delta p_t\| \leq K_{\delta p}$$

with

$$\begin{aligned} K_{\delta p} &= e^{KT} \left(\sqrt{T}(\eta + \zeta) + KK_{\delta u}(1 + K'T) + \int_0^T \|\nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \varphi_s) - \nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \theta_s)\| ds \right) \\ &= \sqrt{T}e^{2KT}(1 + K(1 + K'T))(\eta + \zeta) + e^{2KT}K(1 + K'T) \int_0^T \|f(s, u_s^{\theta, \zeta}, \varphi_s) - f(s, u_s^{\theta, \zeta}, \theta_s)\| ds \\ &\quad + e^{KT} \int_0^T \|\nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \varphi_s) - \nabla_u H(s, u_s^{\theta, \zeta}, p_s^{\theta, \zeta}, \theta_s)\| ds \end{aligned}$$

■

We can now prove Lemma 3.

Proof [Proof of Lemma 3] By definition (7) of the Hamiltonian, for any $\theta \in L^\infty([0, T], \Theta)$,

$$I(\mathbf{u}^\theta, \mathbf{p}^\theta, \theta) := \int_0^T \left(p_t^\theta \cdot f(u_t^\theta, \theta_t) - H(t, u_t^\theta, p_t^\theta, \theta_t) - L(\theta_t) \right) dt = 0.$$

Thus

$$\begin{aligned} 0 &= I(\mathbf{u}^{\varphi, \eta}, \mathbf{p}^{\varphi, \eta}, \varphi) - I(\mathbf{u}^{\theta, \zeta}, \mathbf{p}^{\theta, \zeta}, \theta) \tag{30} \\ &= \underbrace{\int_0^T p_t^{\varphi, \eta} \cdot f(u_t^{\varphi, \eta}, \varphi_t) - p_t^{\theta, \zeta} \cdot f(u_t^{\theta, \zeta}, \theta_t) dt}_{:= \mathcal{I}_1} - \underbrace{\int_0^T H(t, u_t^{\varphi, \eta}, p_t^{\varphi, \eta}, \varphi_t) - H(t, u_t^{\theta, \zeta}, p_t^{\theta, \zeta}, \theta_t) dt}_{:= \mathcal{I}_2} \\ &\quad - \underbrace{\int_0^T R(\varphi_t) - R(\theta_t) dt}_{:= \mathcal{I}_3} \end{aligned}$$

Denoting

$$e_t^{f,\varphi,\eta} := f(u_t^{\varphi,\eta}, \varphi_t) - \Pi_\eta f(u_t^{\varphi,\eta}, \varphi_t)$$

and similarly for $e_t^{f,\theta,\zeta}$, we have

$$\mathcal{I}_1 = \underbrace{\int_0^T p_t^{\varphi,\eta} \cdot \Pi_\eta f(u_t^{\varphi,\eta}, \varphi_t) - p_t^{\theta,\zeta} \cdot \Pi_\zeta f(u_t^{\theta,\zeta}, \theta_t) dt}_{:=\mathcal{I}_{1,1}} + \underbrace{\int_0^T p_t^{\varphi,\eta} \cdot e_t^{f,\varphi,\eta} - p_t^{\theta,\zeta} \cdot e_t^{f,\theta,\zeta} dt}_{:=\mathcal{I}_{1,2}}$$

In addition, since $\dot{u}_t^{\varphi,\eta} = \Pi_\eta f(u_t^{\varphi,\eta}, \varphi_t)$ and similarly for θ , it follows that

$$\mathcal{I}_{1,1} = \int_0^T p_t^{\varphi,\eta} \cdot \dot{u}_t^{\varphi,\eta} - p_t^{\theta,\zeta} \cdot \dot{u}_t^{\theta,\zeta} dt = \underbrace{\int_0^T p_t^{\theta,\zeta} \cdot \delta \dot{u}_t + \delta p_t \cdot \dot{u}_t^{\theta,\zeta} dt}_{:=\mathcal{I}_{1,1,1}} + \underbrace{\int_0^T \delta p_t \cdot \delta \dot{u}_t dt}_{:=\mathcal{I}_{1,1,2}}$$

where

$$\delta u_t := u_t^{\varphi,\eta} - u_t^{\theta,\zeta}, \quad \text{and} \quad \delta p_t := p_t^{\varphi,\eta} - p_t^{\theta,\zeta}.$$

By integration by parts, and the fact that $\dot{p}_t^{\theta,\zeta} = -\Pi_\zeta \nabla_u H(t, u_t^{\theta,\zeta}, p_t^{\theta,\zeta}, \theta_t)$, we have

$$\begin{aligned} \mathcal{I}_{1,1,1} &= p_t^{\theta,\zeta} \cdot \delta u_t \Big|_0^T + \int_0^T \delta p_t \cdot \dot{u}_t^{\theta,\zeta} - \dot{p}_t^{\theta,\zeta} \cdot \delta u_t dt \\ &= p_t^{\theta,\zeta} \cdot \delta u_t \Big|_0^T + \int_0^T \delta u_t \cdot \Pi_\zeta \nabla_u H(t, u_t^{\theta,\zeta}, p_t^{\theta,\zeta}, \theta_t) + \delta p_t \cdot \Pi_\zeta \nabla_p H(t, u_t^{\theta,\zeta}, p_t^{\theta,\zeta}, \theta_t) dt \\ &= p_t^{\theta,\zeta} \cdot \delta u_t \Big|_0^T + \int_0^T \Pi_\zeta \nabla_w H(t, w_t^{\theta,\zeta}, \theta_t) \cdot \delta w_t dt, \end{aligned}$$

where we have used the shorthand notation $w = (u, p)$ in the last line. This notation will also be used in what follows.

We next address the integral $\mathcal{I}_{1,1,2}$. By integration by parts,

$$\begin{aligned} \mathcal{I}_{1,1,2} &= \frac{1}{2} \int_0^T \delta p_t \cdot \delta \dot{u}_t dt + \frac{1}{2} \int_0^T \delta p_t \cdot \delta \dot{u}_t dt \\ &= \frac{1}{2} \delta p_t \cdot \delta u_t \Big|_0^T - \frac{1}{2} \int_0^T \delta \dot{p}_t \cdot \delta u_t dt + \frac{1}{2} \int_0^T \delta p_t \cdot \delta \dot{u}_t dt \\ &= \frac{1}{2} \delta p_t \cdot \delta u_t \Big|_0^T + \frac{1}{2} \int_0^T \left(\Pi_\eta \nabla_w H(t, w_t^{\varphi,\eta}, \varphi_t) - \Pi_\zeta \nabla_w H(t, w_t^{\theta,\zeta}, \theta_t) \right) \cdot \delta w_t dt \end{aligned}$$

Adding and subtracting $\nabla_w H(t, w_t^{\varphi,\eta}, \varphi_t)$ and $\nabla_w H(t, w_t^{\theta,\zeta}, \theta_t)$ to the last formula,

$$\begin{aligned} \mathcal{I}_{1,1,2} &= \frac{1}{2} \delta p_t \cdot \delta u_t \Big|_0^T + \frac{1}{2} \int_0^T \left(\Pi_\eta \nabla_w H(t, w_t^{\varphi,\eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi,\eta}, \varphi_t) \right) \cdot \delta w_t dt \\ &\quad + \frac{1}{2} \int_0^T \left(\nabla_w H(t, w_t^{\theta,\zeta}, \theta_t) - \Pi_\zeta \nabla_w H(t, w_t^{\theta,\zeta}, \theta_t) \right) \cdot \delta w_t dt \\ &\quad + \frac{1}{2} \int_0^T \left(\nabla_w H(t, w_t^{\varphi,\eta}, \varphi_t) - \nabla_w H(t, w_t^{\theta,\zeta}, \theta_t) \right) \cdot \delta w_t dt \end{aligned}$$

By applying Taylor's theorem around $w = w^{\theta, \zeta} = (u_t^{\theta, \zeta}, p_t^{\theta, \zeta})$ to the function $w \rightarrow \nabla_w H(t, w, \varphi_t)$, there exists $r_1(t) \in [0, 1]$ such that

$$\nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) = \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) + \delta w_t \cdot \nabla_w^2 H(t, w_t^{\theta, \zeta} + r_1(t)\delta w_t, \varphi_t)$$

Thus

$$\begin{aligned} \mathcal{I}_{1,1,2} &= \frac{1}{2} \delta p_t \cdot \delta u_t \Big|_0^T + \frac{1}{2} \int_0^T (\Pi_\eta \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t)) \cdot \delta w_t dt \\ &\quad + \frac{1}{2} \int_0^T (\nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) - \Pi_\zeta \nabla_w H(t, w_t^{\theta, \zeta}, \theta_t)) \cdot \delta w_t dt \\ &\quad + \frac{1}{2} \int_0^T (\nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) - \nabla_w H(t, w_t^{\theta, \zeta}, \theta_t)) \cdot \delta w_t dt \\ &\quad + \frac{1}{2} \int_0^T \delta w_t \cdot \nabla_w^2 H(t, w_t^{\theta, \zeta} + r_1(t)\delta w_t, \varphi_t) \cdot \delta w_t dt \end{aligned}$$

Since $\delta u_0 = 0$, the boundary terms from $\mathcal{I}_{1,1,1}$ and $\mathcal{I}_{1,1,2}$ are

$$\begin{aligned} \left(p_t^{\theta, \zeta} + \frac{1}{2} \delta p_t \right) \cdot \delta u_t \Big|_0^T &= \left(p_T^{\theta, \zeta} + \frac{1}{2} \delta p_T \right) \cdot \delta u_T \\ &= -\nabla \Phi(u_T^{\theta, \zeta}) \cdot \delta u_T - \frac{1}{2} \left(\nabla \Phi(u_T^{\varphi, \eta}) - \nabla \Phi(u_T^{\theta, \zeta}) \right) \cdot \delta u_T \end{aligned}$$

By applying Taylor's theorem around $z = u_t^{\theta, \zeta}$ to the function $z \rightarrow \nabla \Phi(z)$, there exists $r_2 \in [0, 1]$ such that

$$\nabla \Phi(u_T^{\varphi, \eta}) = \nabla \Phi(u_T^{\theta, \zeta}) + \delta u_T \cdot \nabla^2 \Phi(u_T^{\theta, \zeta} + r_2 \delta u_T).$$

Therefore

$$\left(p_t^{\theta, \zeta} + \frac{1}{2} \delta p_t \right) \cdot \delta u_t \Big|_0^T = -\nabla \Phi(u_T^{\theta, \zeta}) \cdot \delta u_T - \frac{1}{2} \delta u_T \cdot \nabla^2 \Phi(u_T^{\theta, \zeta} + r_2 \delta u_T) \cdot \delta u_T$$

Since, again by Taylor's theorem, there exists $r_3 \in [0, 1]$ such that

$$\Phi(p_T^{\varphi, \eta}) = \Phi(p_T^{\theta, \zeta}) + \delta u_T \cdot \nabla \Phi(u_T^{\theta, \zeta}) + \frac{1}{2} \delta u_T \cdot \nabla^2 \Phi(u_T^{\theta, \zeta} + r_3 \delta u_T) \cdot \delta u_T,$$

we finally get the expression for the the boundary terms from $\mathcal{I}_{1,1,1}$ and $\mathcal{I}_{1,1,2}$

$$\left(p_t^{\theta, \zeta} + \frac{1}{2} \delta p_t \right) \cdot \delta u_t \Big|_0^T = \Phi(u_T^{\theta, \zeta}) - \Phi(u_T^{\varphi, \eta}) - \frac{1}{2} \delta u_T \cdot \left(\nabla^2 \Phi(u_T^{\theta, \zeta} + r_2 \delta u_T) - \nabla^2 \Phi(u_T^{\theta, \zeta} + r_3 \delta u_T) \right) \cdot \delta u_T$$

We now turn to \mathcal{I}_2 . By Taylor's theorem, there exists $r_4(t) \in [0, 1]$ such that

$$\begin{aligned} \mathcal{I}_2 &= \int_0^T H(t, w_t^{\varphi, \eta}, \varphi_t) - H(t, w_t^{\theta, \zeta}, \theta_t) dt \\ &= \int_0^T \Delta H_{\varphi, \theta}^{x, \zeta}(t) dt + \int_0^T \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) \cdot \delta w_t dt \\ &\quad + \frac{1}{2} \int_0^T \delta w_t \cdot \nabla_w^2 H(t, w_t^{\theta, \zeta} + r_4(t)\delta w_t, \varphi_t) \cdot \delta w_t dt, \end{aligned}$$

where we remind that, as defined in equation (20), $\Delta H_{\varphi, \theta}^{x, \zeta}(t) := H(t, w_t^{\theta, \zeta}, \varphi_t) - H(t, w_t^{\theta, \zeta}, \theta_t)$. Injecting the above expressions for the terms of \mathcal{I}_1 and \mathcal{I}_2 into equation (30), passing the terms corresponding to the loss function to the left hand side, and adding and subtracting $\int_0^T \Pi_{\zeta} \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) \cdot \delta w_t dt$,

$$\text{Loss}(x, y, \varphi) - \text{Loss}(x, y, \theta) \quad (\text{Diff-Loss})$$

$$= \Phi(u_T^{\varphi, \eta}) - \Phi(u_T^{\theta, \zeta}) + \int_0^T (R(\varphi_t) - R(\theta_t)) dt$$

$$= \int_0^T \left(\Pi_{\zeta} \nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) - \Pi_{\zeta} \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) \right) \cdot \delta w_t dt \quad (\text{T1})$$

$$+ \int_0^T p_t^{\varphi, \eta} \cdot e_t^{f, \varphi, \eta} - p_t^{\theta, \zeta} \cdot e_t^{f, \theta, \zeta} dt \quad (\text{T2})$$

$$+ \frac{1}{2} \int_0^T \left(\Pi_{\eta} \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) \right) \cdot \delta w_t dt \quad (\text{T3})$$

$$+ \frac{1}{2} \int_0^T \left(\nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) - \Pi_{\zeta} \nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) \right) \cdot \delta w_t dt \quad (\text{T4})$$

$$+ \frac{1}{2} \int_0^T \left(\nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) - \nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) \right) \cdot \delta w_t dt \quad (\text{T5})$$

$$+ \frac{1}{2} \int_0^T \delta w_t \cdot \nabla_w^2 H(t, w_t^{\theta, \zeta} + r_1(t) \delta w_t, \varphi_t) \cdot \delta w_t dt \quad (\text{T6})$$

$$- \int_0^T \Delta H_{\varphi, \theta}^{x, \zeta}(t) dt \quad (\text{T7})$$

$$+ \int_0^T \left(\Pi_{\zeta} \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) - \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) \right) \cdot \delta w_t dt \quad (\text{T8})$$

$$- \frac{1}{2} \int_0^T \delta w_t \cdot \nabla_w^2 H(t, w_t^{\theta, \zeta} + r_4(t) \delta w_t, \varphi_t) \cdot \delta w_t dt \quad (\text{T9})$$

$$- \frac{1}{2} \delta u_T \cdot \left(\nabla^2 \Phi(u_T^{\theta, \zeta} + r_2 \delta u_T) - \nabla^2 \Phi(u_T^{\theta, \zeta} + r_3 \delta u_T) \right)$$

We next derive a bound for the difference $\text{Loss}(x, y, \varphi) - \text{Loss}(x, y, \theta)$. We proceed to bound all the terms by order of appearance in the above formula. We will sometimes use that by (28), (29) and Jensen's inequality,

$$\|\delta w_t\|^2 = \|\delta u_t\|^2 + \|\delta p_t\|^2 \lesssim (\eta + \zeta)^2 + \int_0^T \|\nabla_w H(s, w_s^{\theta, \zeta}, \varphi_s) - \nabla_w H(s, w_s^{\theta, \zeta}, \theta_s)\|^2 ds$$

which yields

$$\|\delta w_t\|_{L^2([0, T])} \lesssim (\eta + \zeta) + \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)}. \quad (32)$$

Term (T1)+(T5): By the Cauchy-Schwartz inequality and (32), we can bound term (T1) by

$$\begin{aligned} & \int_0^T \left(\Pi_\zeta \nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) - \Pi_\zeta \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t) \right) \cdot \delta w_t dt \\ & \leq \|\Pi_\zeta\| \|\nabla_w H(t, w_t^{\theta, \zeta}, \theta_t) - \nabla_w H(t, w_t^{\theta, \zeta}, \varphi_t)\|_{L^2([0, T])} \|\delta w_t\|_{L^2([0, T])} \\ & \lesssim (\eta + \zeta) \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T])} + \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)}^2 \end{aligned}$$

Thus

$$(T1) + (T5) \lesssim (\eta + \zeta) + (\eta + \zeta) \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T])} + \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T])}^2$$

Term (T2): Using (27) and the fact that $\|e_t^{f, \theta, \zeta}\|_{L^2([0, T], \mathbb{R}^n)} \leq \zeta$

$$\int_0^T p_t^{\varphi, \eta} \cdot e_t^{f, \varphi, \eta} - p_t^{\theta, \zeta} \cdot e_t^{f, \theta, \zeta} dt \lesssim \eta + \zeta$$

Terms (T3) and (T8): By the Cauchy-Schwarz inequality, we can bound term (T3)

$$\begin{aligned} & \left| \int_0^T (\Pi_\eta \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t)) \cdot \delta w_t dt \right| \\ & \leq \|\Pi_\eta \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t)\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)} \|\delta w_t\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)} \\ & \lesssim \eta \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)} \end{aligned}$$

where we have used that the propagator Π_η guarantees

$$\|\Pi_\eta \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t)\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)} \lesssim \eta$$

and finally

$$\left| \int_0^T (\Pi_\eta \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t) - \nabla_w H(t, w_t^{\varphi, \eta}, \varphi_t)) \cdot \delta w_t dt \right| \lesssim \eta(\eta + \zeta) + \eta \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T], \mathbb{R}^n \times \mathbb{R}^n)}$$

Term (T4): Using (27) and the fact that $\|e_t^{f, \theta, \zeta}\|_{L^2([0, T], \mathbb{R}^n)} \leq \zeta$, we have $\int_0^T p_t^{\varphi, \eta} \cdot e_t^{f, \varphi, \eta} - p_t^{\theta, \zeta} \cdot e_t^{f, \theta, \zeta} dt \lesssim \eta + \zeta$.

Term (T6), (T9) and (T10): By assumptions (A1), (A2), all second derivative terms are bounded element-wise by some constant K . Hence, we have $|\delta w_t \cdot A \cdot \delta w_t| \lesssim \|\delta w_t\|^2$ for A being a second derivative matrix.

Summary: By using the above inequalities, in formula (Diff-Loss), we deduce that there exists a constant $C > 0$ depending on T , $\|\Pi_\zeta\|$ and $\|\Pi_\eta\|$ such that, if $\eta \leq 1$,

$$\text{Loss}(x, y, \varphi) - \text{Loss}(x, y, \theta) \leq - \int_0^T \Delta H_{\varphi, \theta}^{x, \zeta}(t) dt + C \left((\eta + \zeta)^2 + \|\nabla_w(\Delta H_{\varphi, \theta}^{x, \zeta})\|_{L^2([0, T])}^2 \right)$$

Note that the term $(\eta + \zeta)^2$ accounts for the fact that we are not solving exactly the dynamics. ■