

Neural Operators Can Play Dynamic Stackelberg Games

Guillermo A. Alvarez*

GUIALV@UMICH.EDU

*University of Michigan
Department of Mathematics
2074 East Hall, 530 Church Street Ann Arbor, Michigan, USA*

Ibrahim Ekren*

IEKREN@UMICH.EDU

*University of Michigan
Department of Mathematics
2074 East Hall, 530 Church Street Ann Arbor, Michigan, USA*

Anastasis Kratsios*†

KRATSIOA@MCMASTER.CA

*McMaster University and the Vector Institute
Department of Mathematics
1280 Main Street West, Hamilton, Ontario, L8S 4K1, Canada*

Xuwei Yang*

YANGX212@MCMASTER.CA

*McMaster University
Department of Mathematics
1280 Main Street West, Hamilton, Ontario, L8S 4K1, Canada*

Abstract

Dynamic Stackelberg games are a broad class of two-player games in which the leader acts first, and the follower chooses a response strategy to the leader’s strategy. Unfortunately, only stylized Stackelberg games are explicitly solvable since the follower’s best-response operator (as a function of the control of the leader) is typically analytically intractable. This paper addresses this issue by showing that the *follower’s best-response operator* can be approximately implemented by an *attention-based neural operator*, uniformly on compact subsets of adapted open-loop controls for the leader. We further show that the value of the Stackelberg game where the follower uses the approximate best-response operator approximates the value of the original Stackelberg game. Our main result is obtained using our universal approximation theorem for attention-based neural operators between spaces of square-integrable adapted stochastic processes, as well as stability results for a general class of Stackelberg games.

1. Introduction

In the classical formulation of Stackelberg games, there are generally two players: a leader (major) who moves first and a follower (minor) who then reacts. One is typically interested in studying the equilibrium of these games, in which both players cannot increase their utilities by (re)acting differently. The generic structure of these games has led their equilibria to become a powerful mathematical tool to describe the evolution of incentives in complex environments with economic applications ranging from control of disease transmission in epidemiology [Aurell et al. \(2022\)](#); [Hubert et al. \(2022\)](#), contract design [Conitzer and Sandholm \(2006\)](#); [Elie et al. \(2019\)](#); [Keppo et al. \(2024\)](#); [Hernández and Possamai \(2024\)](#); [Hernández et al. \(2024\)](#), advertising [He et al. \(2008\)](#), mobile network planning [Zheng et al. \(2018\)](#), economic behaviour in oligopolies [Carmona and Dayanikli \(2021\)](#), brokerage [Alvarez et al. \(2023\)](#), (re)insurance [Cao et al. \(2022\)](#); [Kroell et al. \(2023\)](#); [Ghossoub and Zhu \(2024\)](#), risk management [Bensalem et al. \(2020\)](#); [Li et al. \(2022\)](#), algorithmic auction/mechanism

*. Equal contribution, all authors are listed in alphabetic order.

†. Corresponding author.

design Conitzer and Sandholm (2006); Konrad and Leininger (2007); Dierks and Seuken (2022), security Kar et al. (2017); An et al. (2011), and green investments Zhang et al. (2023). Most of these applications consider *dynamic* Stackelberg games, where the game is played continuously in several rounds. Even though these games provide powerfully descriptive theoretical vehicles, the highly intertwined structure of Stackelberg equilibria can be challenging both numerically and analytically.

This paper shows that deep learning can provide a viable and generic computational vehicle by which dynamic Stackelberg games can be computationally solved. We exhibit a class of *neural operators* leveraging an attention mechanism which can approximately implement the follower’s best response map to arbitrary precision, uniformly on compact subsets of the leader’s strategies (defined below). Unlike most neural operator models, which focus on learning the solution map of PDEs Kovachki et al. (2021); Lanthaler et al. (2022a); Lee et al. (2023); Lanthaler and Stuart (2023); Kovachki et al. (2023); Benitez et al. (2023); Raonic et al. (2024); Bartolucci et al. (2024); Fanaskov and Oseledets (2024) or inverse-problems Calderon-Macias (1997); Molinaro et al. (2023); de Hoop et al. (2022, 2024), our neural operators are not defined between function spaces but between spaces of stochastic controls. Our attention mechanism mimics that of Vaswani et al. (2017) used in transformers Bahdanau et al. (2015) while reflecting the geometry of the input and output spaces of stochastic processes, and it extends the attention mechanisms of Acciaio et al. (2023). We note that, it is natural to consider compact sets of controls not only from approximation-theoretic vantage point but also from the control-theoretic perspective; this is because this guarantees the existence of a Stackelberg equilibrium under only minimal assumptions.

Here, we consider the following class of dynamic Stackelberg games, with stochastic effects, where both players (re)act in continuous time according to the following general dynamics

$$dX_t = f(X_t, u_t^0, u_t^1)dt + \sigma(X_t, u_t^0, u_t^1)dW_t$$

where $W \stackrel{\text{def.}}{=} (W_t)_{t \geq 0}$ is d -dimensional standard Brownian motions and $u^i \stackrel{\text{def.}}{=} (u_t^i)_{t \geq 0}$ are the (re)actions/strategies of each player, where $i = 0$ is the leader and $i = 1$ indexes the follower. Each player seeks to optimize their respective objective functions, one in which we impose minimal continuity requirements since we are not interested in analytic expression, which would require highly stylized assumptions on the dynamics and objective functions of all involved players. Rather, our goal is to show that a deep learning solution via neural operators is possible for a broad class of Stackelberg games lying outside the scope of these classical stylized settings. Our first main result (Theorem 7) shows under enough strong-convexity requirements on the utility of the follower the best response maps depend continuously on each leader’s actions, then there is a neural operator which can approximate the follower’s best response map, uniformly over any compact set of actions of the leader, to any given precision.

In general, it is well-known that the approximation of non-linear maps/operators between infinite-dimensional Hilbert spaces by deep learning models may be practically challenging due to necessarily slow convergence rates; see e.g. Lanthaler and Stuart (2023), which is effectively an exacerbated version of the curse of dimensionality known in the finite-dimensional setting, see e.g. Shen et al. (2022). Unlike the finite-dimensional setting, sufficient smoothness is insufficient to obtain reasonable convergence rates in general, Galimberti et al. (2022), and typically, one has to hope for favorable structures which can be exploited by the neural operator; see e.g. Marcati and Schwab (2023), to obtain fast convergence rates. Fortunately, we identify a set of structures that can be exploited by our neural operator for a class of Stackelberg games that encompass analytically-solvable linear-quadratic games. Our second main result (Theorem 11) shows that if the compact set of controls is compatible with the best-response map, then one may guarantee efficient convergence rates for neural operator approximations of the best response map for the follower. We conclude that neural operators can efficiently approximate the solutions to a wide class of Stackelberg games, encompassing those games which are solvable via classical analytic means.

Additionally, we identify an *unsupervised* objective function which provides a heuristic helping to detect the suitability of a neural operator approximating the best response map of the follower

(Theorem 8). Importantly, this heuristic objective function does not require observations of the true optimal response (which would be an unrealistic supervised problem).

We use control theoretical arguments to prove that the leader and the follower have optimal strategies with enough continuous dependence on one another to be *uniformly* approximable on compact sets of (re)actions. We also provide an illustrative counterexample showing that the best-response map might fail to be continuous without these strong-convexity requirements. Thus, without the strong-convexity assumption, discontinuous functions cannot be uniformly approximated by any continuous models due to the Uniform Limit Theorem; see e.g. (Munkres, 2000, Theorem 21.6).

On the technical front: our main control-theoretic contributions (Lemma B.5 together with 13) show that generically the optimal response of the follower is $1/2$ -Hölder continuous in the leader’s control if the problem of the follower is strongly convex. Our main approximation theoretic contribution is a quantitative universal approximation theorem (Theorem 12) showing that neural operators are capable of approximating Hölder continuous non-linear operators between space of square-integrable \mathbb{F} -adapted processes (open-loop controls). Together, these control-theoretic and approximation-theoretic are enough to show that the best-response map is approximable by neural operators. Moreover, for suitable sets of controls, our precise analysis both of the regularity of the best response map and the dependence of the neural operator complexity on the set of controls, allow us to conclude that polynomial approximation rates are possible (Theorem 11).

Our neural operators leverage an attention-like decoding layer, similar to transformers Vaswani et al. (2017), which allows for nonlinear decoding, unlike PCA-net Lanthaler (2023), the encoder-decoder models of Galimberti et al. (2022), and several others. The relationship between our infinite-dimensional analogue of the classical attention of Bahdanau et al. (2015) is also discussed. Attention mechanisms in operator learning, by now, have found common use in implementations; see e.g. the Galerkin transformers of Cao (2021) or the Continuum Attention Mechanism of Calvello et al. (2024).

Organization of Paper Following the literature review in Section 2, the remainder of our paper is organized as follows:

- (i) Section 2 and 3, respectively, review the literature and the necessary background in stochastic analysis and in deep approximation theory to formulate our main results.
- (ii) Section 4 contains our main results.
- (iii) Section 5 explains why our main results work by overviewing our proof strategy, during which we showcase our supporting technical results of independent technical interest.
- (iv) Section 6 showcases examples of Stackelberg games satisfying our convexity requirements.

All technical derivations are relegated to appendix B.

2. Related Literature

Stackelberg Games in Machine Learning In Stackelberg games, both players seek to maximize their gain while being fully rational. This characteristic “conditional” sequential structure is the hallmark challenge rendering Stackelberg games analytically intractable and the reason motivating significant attention from the machine learning community Reisinger and Zhang (2020); Ito et al. (2021); Gao et al. (2022); Haghtalab et al. (2023); Harris et al. (2023); Gerstgrasser and Parkes (2023); Dayanikli and Lauriere (2023), and their related FBSDEs Furuya and Kratsios (2024), looking for new algorithmic tools capable of solving this class of differential games. Nevertheless, there is currently no available deep learning model which is guaranteed to solve a Stackelberg game, much less in continuous time with stochastic effects.

Neural Operators The power of deep learning to solve previously intractable high-dimensional computational problems has motivated the deep learning community to extend these tools to the infinite-dimensional setting with models such as DeepONets Lu et al. (2019, 2021); Goswami et al. (2022) and a variety of *neural operator* architectures; e.g. Fourier Neural Operators Li et al. (2020); Kovachki et al. (2021); Li et al. (2023), graph neural operators Anandkumar et al. (2020), causal neural operators Galimberti et al. (2022), neural operator analogues of transformers Hao et al. (2023), convolutional neural operators Raonic et al. (2024), encoder-decoder models such as PCANet Lanthaler et al. (2022b), and a myriad of other models. The community has largely been motivated by demand in the scientific computing community, focusing on designing neural operators tailored which can learn to solve (i.e. learn the solution operator) to high-dimensional partial integro-differential equations (PIDEs) with applications ranging from physics and engineering De Ryck et al. (2024) to quantitative finance Acciaio et al. (2023). However, the full power of neural operators remains otherwise largely unexplored as the community has focused mainly on the approximation capacity of these models between function spaces with a view towards PIDEs. Here, we probe the limits of neural operators beyond PIDEs by showing that they can solve a broad class of problems at the intersection of game theory and stochastic analysis.

Notation

We use the following notations. For any $C \in \mathbb{N}_+$, we denote the C -simplex by $\Delta_C \stackrel{\text{def.}}{=} \{u \in [0, 1]^C : \sum_{c=1}^C u_c = 1\}$ and we define the associated softmax function by $\text{softmax}_C : \mathbb{R}^C \ni w \mapsto (e^{w_c} / \sum_{c=1}^C e^{w_c})_{c=1}^C \in \Delta_C$. Both in Δ_C and softmax_C , the subscript C will be suppressed when clear from the context. We will write $f \in \tilde{\mathcal{O}}(g)$ if $f \in \mathcal{O}(g \log^k(g))$ for some $k \in \mathbb{N}_+$.

3. Preliminaries

We now overview the background required to formulate the main results of our paper and formalize our neural operator model. We first overview the notion of a square-integrable predictable process from stochastic analysis and then the definition of a multilayer perceptron (MLP) from deep learning. Additional details are included in Appendix B.3.

3.1 Predictable Processes

Fix a time horizon $T > 0$ and let $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ be a filtered probability space whose filtration $\mathbb{F} \stackrel{\text{def.}}{=} (\mathcal{F}_t)_{0 \leq t \leq T}$ is generated by a d -dimensional Brownian motion $W. \stackrel{\text{def.}}{=} (W_t)_{0 \leq t \leq T}$; for some $d \in \mathbb{N}_+$. We consider the space \mathcal{H}_T^2 of all square-integrable \mathbb{F} -predictable processes on $(\Omega, \mathcal{F}, \mathbb{F}, \mathbb{P})$ whose elements $H. \stackrel{\text{def.}}{=} (H_t)_{0 \leq t \leq T} \in \mathcal{H}_T^2$ consist of d -dimensional \mathbb{F} -predictable processes for which the norm

$$\|H\|_{\mathcal{H}_T^2}^2 \stackrel{\text{def.}}{=} \mathbb{E} \left[\int_0^T |H_s|^2 ds \right],$$

is finite. Furthermore, \mathcal{H}_T^2 is a separable infinite-dimensional Hilbert space whose inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}_T^2}$, is given for each $H., \tilde{H}. \in \mathcal{H}_T^2$ by

$$\langle H., \tilde{H}. \rangle_{\mathcal{H}_T^2}^2 \stackrel{\text{def.}}{=} \mathbb{E} \left[\int_0^T H_s^\top \tilde{H}_s ds \right].$$

Similarly to the Fourier neural operator and the neural operators of, our approximation results will rely on an orthogonal basis¹ of \mathcal{H}_T^2 . We begin by constructing a relatively computationally

1. Some authors emphasize that “orthonormal basis” is rather a complete orthonormal system since it is not a linear algebraic basis. However, we follow the common abuse of terminology standard in functional analysis.

convenient basis of $L^2(\mathcal{F}_t)$ based on the Wiener Chaos decomposition. A key feature of the following orthonormal basis is that no iterated stochastic integrals need to be explicitly computed, as is the case with general Wiener Chaos decompositions; see Appendix B.3 for additional details.

For any $i \in \mathbb{N}$, the Hermite polynomials $(h_i)_{i \in \mathbb{N}}$ are the eigenfunctions of the generator $\frac{d^2}{dx^2} - x \frac{d}{dx}$ of the Ornstein-Uhlenbeck process $dX_t = -X_t + \sqrt{2}dW_t$. For each $i \in \mathbb{N}_+$, the i^{th} Hermite polynomial h_i is given recursively by Rodrigues' formula as

$$h_i(x) = \frac{(-1)^i}{i!} e^{x^2/2} \frac{d^i}{dx^i} e^{-x^2/2}, \quad h_0(x) = 1.$$

The Hermite polynomials allow us to define a family of random variables $\{u_{i,j,k}^t : j \in \mathbb{N}, i, k \in \mathbb{N}, \frac{k+1}{2^i} \leq 1\} \subset L^2(\mathcal{F}_t)$ where each $u_{i,j,k}^t$ is defined by

$$u_{i,j,k}^t \stackrel{\text{def.}}{=} \prod_{\tilde{j}=1}^j h_{\tilde{j}} \left(2^i W_{\frac{tk}{2^i}} - 2^{i+1} W_{\frac{t(1+2k)}{2^{i+1}}} + 2^i W_{\frac{t(k+1)}{2^i}} \right). \quad (1)$$

Using (1), we construct a predictable and dynamic version of the Wiener chaos decomposition. First recall the *Haar (wavelet) system* on the larger space $L^2([0, T])$ where $(\psi_{i,k})_{i,k \in \mathbb{N}; 0 \leq k, \frac{k+1}{2^i} \leq 1}$ and

$$\psi_{i,k}(t) \stackrel{\text{def.}}{=} 2^i \left(I_{[T \frac{k}{2^i}, T \frac{1+2k}{2^{i+1}}]}(t) - I_{[T \frac{1+2k}{2^{i+1}}, T \frac{k+1}{2^i}]}(t) \right)$$

which is a complete orthonormal basis of $L^2([0, T])$, see (Meyer, 1990, Chapter 3). The Haar wavelet system will allow us to activate/deactivate the random variables in Wiener Chaos, in (1), as a function of time. We focus on the collection of simple processes obtained as linear combinations of

$$\begin{aligned} \mathcal{S} &\stackrel{\text{def.}}{=} \left\{ u_{i,j,k}^{s_1, s_2} : i, j, k, s_1, s_2 \in \mathbb{N}, s_2 + 1 \leq 2^{s_1}, \frac{k+1}{2^i} \leq \frac{s_2}{2^{s_1}} \right\} \\ u_{i,j,k}^{s_1, s_2}(t, \omega) &\stackrel{\text{def.}}{=} \psi_{s_1, s_2}(t) \cdot u_{i,j,k}^T(\omega). \end{aligned} \quad (2)$$

It can be shown, see Lemma B.6 in the proofs section, that \mathcal{S} is an orthonormal basis of \mathcal{H}_T^2 .

Furthermore, \mathcal{S} has an elegant interpretation as a Haar wavelet expansion in time and the iterated Itô stochastic integral of a Haar wavelet expansion in space. The ability to explicitly compute the iterated Itô stochastic integrals of the Haar wavelet system in space makes this closed-form expansion particularly favourable, especially in higher dimensions where the computation of iterated stochastic integrals can be computationally intensive.

3.1.1 EXAMPLES OF COMPACT SETS OF SQUARE-INTEGRABLE CONTROLS

We will often be considering compact subsets of \mathcal{H}_T^2 whereon we frame the existence of Stackelberg equilibria exist and uniform approximation is possible. There are several other examples of a compact subset of \mathcal{H}_T^2 routinely encountered in the literature, for example, sets of processes which are Malliavin differentiable with uniformly bounded Malliavin Derivative, see (Baños et al., 2018, Corollary C.3.), or perturbations of closed-loop controls which are efficiently approximable (see Section 4.2). Two illustrative, but broad classes, of examples are now constructed; building on compactness results in classical function spaces.

Example 1 (Compactness Via Regularity of the Malliavin Derivative). *For given $C \geq 0$, and $\alpha \in (0, 1)$, define $\mathcal{K}_{\alpha, C} \subset \mathcal{H}_T^2$ as the set of $H \in \mathcal{H}_T^2$ so that*

$$\sup_{0 \leq s \leq t \leq T} \mathbb{E}[|D_s H_t|^2 + |H_t|^2] \leq C, \quad \sup_{0 \leq s < t \leq T} \frac{\mathbb{E}[|H_t - H_s|^2]}{|t - s|^\alpha} \leq C, \quad \text{and} \quad \sup_{r, 0 \leq s < t \leq T} \frac{\mathbb{E}[|D_t H_r - D_s H_r|^2]}{|t - s|^\alpha} \leq C.$$

For $H \in \mathcal{K}_{\alpha,C}$, thanks to (Nualart, 2006, Proposition 1.3.8) we can compute the Malliavin derivative $D_t \int_0^T H_r dW_r = H_t + \int_t^T D_t H_r dW_r$ for $0 \leq s < t \leq T$ so that we have the following estimate

$$\begin{aligned} & \frac{\mathbb{E}[|D_t \int_0^T H_r dW_r - D_s \int_0^T H_r dW_r|^2]}{|t-s|^\alpha} \\ & \leq 2 \frac{\mathbb{E}[|H_t - H_s|^2] + \mathbb{E}[|\int_s^t D_s H_r dW_r|^2] + \mathbb{E}[|\int_t^T D_t H_r - D_s H_r dW_r|^2]}{|t-s|^\alpha} \\ & \leq 2 \frac{\mathbb{E}[|H_t - H_s|^2]}{|t-s|^\alpha} + 2 \frac{\int_s^t \mathbb{E}[|D_s H_r|^2] dr}{|t-s|^\alpha} + 2 \int_t^T \frac{\mathbb{E}[|D_t H_r - D_s H_r|^2]}{|t-s|^\alpha} dr \leq 2C(1 + T^{1-\alpha} + T). \end{aligned}$$

Thanks to (Baños et al., 2018, Corollary C3), this estimate implies that the set of random variables $\{\int_0^T H_s dW_s : H \in \mathcal{K}_{\alpha,C}\}$ is relatively compact in the set of square integrable random variables. This implies by Ito's isometry that $\mathcal{K}_{\alpha,C}$ is relatively compact in \mathcal{H}_T^2 .

In a very special case of Example 1, one may require the predictable process H . to be a (non-random) smooth function. This next example shows precisely this, and it elucidates the link between classical families of smooth functions and compact sets of open-loop controls.

Example 2 (Compactness Via Continuously Differentiable Martingale Controls). Fix $T > 0$. Let $W^{1,2}([0,1])$ denote the Sobolev space on the unit interval $[0,1]$; and denote its norm by $\|\cdot\|_{W^{1,2}}$. By the Rellich-Kondrashov Theorem, see e.g. (Evans, 2010, Theorem 5.1), the set of function $\varsigma : [0,1] \rightarrow \mathbb{R}$ satisfying

$$\|\varsigma\|_{W^{1,2}} \leq 1 \tag{3}$$

is compact in $L^2([0,1])$. By the Ito isometry the set of \mathcal{F}_T -measurable random variables $\int_0^T \varsigma(t) dW_t$ is therefore compact in $L^2(\mathcal{F}_T)$. Since conditional expectations are non-expansive (1-Lipschitz) then the set $\mathcal{K} \subset \mathcal{H}_T^2$ of processes/open-loop controls $u \stackrel{\text{def.}}{=} (u_t)_{0 \leq t \leq T}$ of the form

$$u_t = \mathbb{E} \left[\int_0^T \varsigma(s) dW_s \middle| \mathcal{F}_t \right] = \int_0^t \varsigma(s) dW_s$$

where ς satisfies (3), is compact in \mathcal{H}_T^2 , and the last inequality held by the Martingale property of the Itô (stochastic) integral.

Example 3 (Deterministic Hölder Continuous Controls). Fix $T > 0$, $0 < \alpha \leq 1$, and consider the set \mathcal{K}_α of deterministic controls $u. = (u_t)_{t \geq 0}$ in \mathcal{H}_T^2 where $t \mapsto u_t$ is an α -Hölder function mapping $[0,T]$ to $[-1,1]^d$. In this case, for each $u., v. \in \mathcal{K}_\alpha$ we have

$$\mathbb{E} \left[\int_0^T |u_t - v_t|^2 \right]^{1/2} = \left(\int_0^T |u_t - v_t|^2 \right)^{1/2} \leq \max_{0 \leq t \leq T} |u_t - v_t|.$$

Therefore, the map $C([0,1], \mathbb{R}^d) \rightarrow \mathcal{H}_T^2$ is a 1-Lipschitz embedding when the domain is equipped with the uniform norm. By the Arzelà-Ascoli Theorem, we have that any set of uniformly bounded α -Hölder functions is relatively compact in $C([0,1], \mathbb{R}^d)$; thus, \mathcal{K}_α is relatively compact in \mathcal{H}_T^2 .

One can easily extend the construction in Example 2 to non-Martingale controls iterated integrals using the isometries between the space of symmetric functions in $L^2([0,1]^q)$, for any $q \in \mathbb{N}_+$, and the q^{th} Wiener Chaos (see e.g. (Nualart, 2006, Theorem 1.1.1)); however, do not do so for simplicity of presentation.

Example 4 (Conditioned Lipschitz Perturbations of Random Variables at Terminal Time). Fix $X \in L^2(\mathcal{F}_T)$, and fix $T > 0$. Consider the set \mathcal{X} of 1-Lipschitz function $f : \mathbb{R}^d \rightarrow [-1,1]^d$ which are

supported on the hypercube $[-1, 1]^d$; i.e. $f(x) = 0$ if $x \notin [-1, 1]^d$. By the Arzela-Ascoli Theorem, \mathcal{X} is relatively compact in $C(\mathbb{R}^d, \mathbb{R}^d)$. Since the map sending any $f \in C(\mathbb{R}^d, \mathbb{R}^d)$ to $f(X) \in L^2(\mathcal{F}_T)$ is 1-Lipschitz then the set of random variables $\{f(X) \in L^2(\mathcal{F}_T) : f \in \mathcal{X}\}$ is compact in $L^2(\mathcal{F}_T)$. As in Example (2), since conditional expectations are 1-Lipschitz then the set of controls $u \stackrel{\text{def.}}{=} (u_t)_{t \geq 0} \in \mathcal{K} \subset \mathcal{H}_T^2$ of the form

$$u_t = \mathbb{E}[f(X)|\mathcal{F}_t]$$

where $f \in \mathcal{X}$, is relatively compact in \mathcal{H}_T^2 . As a concrete example, one may take f to belong to the set of 1-Lipschitz ReLU Neural Networks with output restricted to belong to $[-1, 1]^d$; see e.g. (Hong and Kratsios, 2024, Theorem 1.1).

Remark 1 (Alternative Proofs of Compactness Directly Via Example 1). Example 4 can also be obtained from Example 1 upon adding Malliavin differentiability requirements on X and additional regularity. Examples 2 and 3 can have alternatively be obtained as straightforward consequences of Example 1. We opted for self-contained presentations for each example to illustrate various construction methods for compacta in \mathcal{H}_T^2 .

In practice, one often discretizes their space when implementing it on a digital machine. In these cases, the set of controls is finite and, therefore, compact.

Example 5 (Finite Sets of Controls). Let $I \in \mathbb{N}$ and $\mathcal{K} \stackrel{\text{def.}}{=} \{u_i\}_{i=1}^I \subset \mathcal{H}_T^2$. Then, \mathcal{K} is compact.

3.2 The Dynamic Stackelberg Game

In the previously introduced probability space $(\Omega, \mathcal{F}, \{\mathcal{F}\}_{0 \leq t \leq T}, \mathbb{P})$, we consider a Stackelberg game with a leader indexed with $i = 0$ and a follower indexed with $i = 1$. The state process of the game is described by the stochastic differential equation

$$dX_t = f(X_t, u_t^0, u_t^1)dt + \sigma(X_t, u_t^0, u_t^1)dW_t, \quad (4)$$

and $u_t^0 \in \mathbb{R}^{d_0}$ and $u_t^1 \in \mathbb{R}^{d_1}$ are the controls of the leader and the follower, respectively. The exact set of admissible controls will be provided below. We assume that a deterministic initial $X_0 \in \mathbb{R}^d$ is fixed and is known to both agents. Thus, we omit the dependence of various parameters on X_0 . The cost functionals of the two players are given by

$$J_0(u^0, u^1) = \mathbb{E} \left[\int_0^T L_0(X_t, u_t^0, u_t^1)dt + g_0(X_T) \right], \quad (5)$$

$$J_1(u^0, u^1) = \mathbb{E} \left[\int_0^T L_1(X_t, u_t^0, u_t^1)dt + g_1(X_T) \right], \quad (6)$$

where $L_i : \mathbb{R}^d \times \mathbb{R}^{d_0} \times \mathbb{R}^{d_1} \mapsto [0, \infty)$ and $g_i : \mathbb{R}^d \mapsto [0, \infty)$, $i = 0, 1$. We require the following regularity conditions of the involved functions.

Assumption 2 (Regularity Conditions). There exists a constant $K > 0$ such that for $h(x, u^0, u^1) = f(x, u^0, u^1)$, $\sigma(x, u^0, u^1)$, $L_i(x, u^0, u^1)$, and $g_i(x)$, $i = 0, 1$,

$$|h(x, u^0, u^1) - h(\tilde{x}, \tilde{u}^0, \tilde{u}^1)| \leq K(|x - \tilde{x}| + |u^0 - \tilde{u}^0| + |u^1 - \tilde{u}^1|). \quad (7)$$

We define

$$\mathcal{U}_i = \left\{ u : [0, T] \times \Omega \rightarrow \mathbb{R}^{d_i} \mid u(\cdot) \text{ is } \{\mathcal{F}\}_t\text{-adapted, } \mathbb{E} \int_0^T |u_t|^2 dt < \infty \right\}$$

and fix $\mathcal{K}_0 \subset \mathcal{U}_0$ so that \mathcal{K}_0 is the set of possible controls of the leader, \mathcal{U}_1 is the set of possible controls of the follower. The introduction of \mathcal{K}_0 is needed due to the fact that our operators in Subsection 3.3 will only perform optimization relative to a compact subset \mathcal{K}_0 of \mathcal{U}_0 .

For each $u^0 \in \mathcal{U}_0$, the set of best responses for the follower is

$$\mathcal{R}(u^0) = \{u \in \mathcal{U}_1 : J_1(u^0, u) \leq J_1(u^0, u^1), \forall u^1 \in \mathcal{U}_1\}. \quad (8)$$

Following the definitions in [Bensoussan et al. \(2015\)](#), we define adapted open-loop (AOL) responses of the follower to the controls of the leader by

$$\bar{\mathcal{U}}_1 = \left\{ u : [0, T] \times \Omega \times \mathcal{K}_0 \rightarrow \mathbb{R}^{d_1} \mid \forall u_0 \in \mathcal{K}_0, u(\cdot, u^0) \in \mathcal{U}_1 \right\}. \quad (9)$$

We recall that implicitly we make the assumption that the initial point X_0 of X is fixed throughout the paper. If this initial condition is not fixed, one has to allow the elements of $\bar{\mathcal{K}}_1$ to also depend on this initial condition as it is the case in [Bensoussan et al. \(2015\)](#).

We study the Stackelberg equilibria for the leader-follower problem; that is, a set of leader-follower strategies wherein the follower optimally responds to the preemptive optimal action of the leader in such a way that neither player can gain utility by perturbing their strategy. Formally, a Stackelberg equilibrium is defined as follows.

Definition 3 (Stackelberg Equilibrium). *A Stackelberg equilibrium (relative to $\mathcal{K}_0 \subset \mathcal{U}_0$) of the leader-follower game (4)-(5)-(6) is a pair $(u^{0,*}, U^*) \in \mathcal{K}_0 \times \bar{\mathcal{U}}_1$ such that $U^*(u^0) \in \mathcal{R}(u^0)$ for all $u^0 \in \mathcal{U}_0$ and*

$$J_0(u^{0,*}, U^*(u^{0,*})) \leq J_0(u^0, U^*(u^0)) \text{ for all } u^0 \in \mathcal{K}_0.$$

If it exists, a map U^ is called a best response map of the follower.*

In general, $\mathcal{R}(u^0)$ can be empty for some $u^0 \in \mathcal{K}_0$. If this happens, the equilibrium will not exist. However, if the problem of the follower is convex enough, $\mathcal{R}(u^0)$ is reduced to a point $U^*(u^0) \in \mathcal{U}_1$ that can be described by an FBSDE. Thus, the existence of a Stackelberg equilibrium is reduced to the optimization of $u^0 \in \mathcal{K}_0 \mapsto J_0(u^0, U^*(u^0))$ that we call the effective criterion of the leader. This can be done if \mathcal{K}_0 is compact or under additional structural assumption on the data as a form of control of solutions of FBSDEs. We will assume that the optimal response operator of the follower exists and possesses a minimal level of regularity.

Assumption 4 (Hölder Continuity of The Follower’s Best Response). *There exists a Hölder continuous mapping $u^0 \in \mathcal{K}_0 \mapsto U^*(u^0) \in \mathcal{U}_1$ so that $U^*(u^0) \in \mathcal{R}(u^0)$ for all $u^0 \in \mathcal{K}_0$.*

Remark 5. *Though several of our universal approximation results (Theorem 12) can be applied only while assuming continuity of the following best response, our quantitative estimates fundamentally rely on Hölder continuity since we use properties of doubling metrics, building on the method of [Kratsios et al. \(2023a\)](#), and metric snowflakes; see ([Weaver, 2018](#), page 66) for definitions.*

In Section 5 and 6, below, we show this Hölder dependence estimates of the optimal response if the optimization problem of the follower is strongly convex. In fact, we provide an example of a static game in Section A showing that even if the problem of the follower is only convex but not strongly convex, the optimal response of the follower will lack continuous dependence on the control of the leader. In such cases, our neural operator cannot approximate the optimal response. Thus, Assumption 4 is not merely an often satisfied and purely technical assumption.

3.3 Neural Operators

Fix an activation function $\sigma : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ and, for any $N \in \mathbb{N}_+$, define its componentwise composition with any vector $x \in \mathbb{R}^N$ with trainable parameter $\alpha \in \mathbb{R}^N$, by $\sigma_\alpha \bullet x \stackrel{\text{def.}}{=} (\sigma_{\alpha_i}(x_i))_{i=1}^d$. For the majority of our paper, we consider neural operators either with an unattainable activation function satisfying the condition of ([Kidger and Lyons, 2020](#)) or a trainable variant $\sigma \in C(\mathbb{R}^2)$ of the super-expressive activation function of [Zhang et al. \(2022\)](#); defined shortly. In the former case, we consider the following activation functions.

Example 6 (Kidger and Lyons (2020)-Type “Standard” Trainable Functions). *There is a non-affine $\sigma_0 \in C(\mathbb{R})$ such that: there exists some $t_0 \in \mathbb{R}$ at which σ_0 is differentiable and such that $\sigma_0(t_0)' \neq 0$. Define $\sigma \in C(\mathbb{R}^2)$ by $(\alpha, t) \mapsto \sigma_\alpha(t)$. Observe that the ReLU activation function falls into this class.*

Example 7 (Super-Expressive Activation with Neuron-Specific Skip-Connection). *We define the trainable variant of the activation function $\sigma : \mathbb{R}^2 \rightarrow \mathbb{R}$ mapping any $(\alpha, t) \in \mathbb{R}^2$ to*

$$\sigma_\alpha(t) \stackrel{\text{def.}}{=} \alpha t + (1 - \alpha) \begin{cases} |t \pmod{2}| & \text{if } t \geq 0 \\ \frac{t}{|t|+1} & \text{if } t < 0 \end{cases} \quad (10)$$

When $\alpha = 0$, then σ_0 coincides with the super-expressive activation function of Zhang et al. (2022). The parameter α allows us to apply a skip connection at any given specific neuron.

We now define the deep learning backbone of our neural operator model, namely the multilayer perceptrons (MLP). Fix input and output dimensions n and m . An MLP with (trainable) activation function σ , depth $J \in \mathbb{N}_+$, and width $W \in \mathbb{N}_+$ is a map $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with iterative representation,

$$\begin{aligned} \hat{f}_\theta(x) &\stackrel{\text{def.}}{=} x^{(J)} + c, \\ x^{(j+1)} &\stackrel{\text{def.}}{=} A^{(j)} \sigma_{\alpha^{(j)}} \bullet (x^{(j)} + b^{(j)}), \\ x^{(0)} &\stackrel{\text{def.}}{=} x. \end{aligned} \quad (11)$$

where for $j = 0, \dots, J-1$: $A^{(j)} \in \mathbb{R}^{d_{j+1} \times d_j}$, $\alpha^{(j)}, b^{(j)} \in \mathbb{R}^{d_{j+1}}$, and $n = d_0, \dots, m = d_J \leq W$. Denote the set of MLPs mapping \mathbb{R}^n to \mathbb{R}^m with depth at-most J and width at-most W by $\mathcal{NN}_{J,W;n,m}$.

The Attentional Neural Operator Model Neural operators (NOs) are natural, infinite dimensional extensions of classical (finite-dimensional) neural networks. We follow the general encoder-processor-decoder NO paradigm considered, e.g. in PCA-nets Lanthaler (2023), Castro (2023), or Kratsios et al. (2023b). Our neural operators model, illustrated in Figure 1, maps inputs and output in the space structure \mathcal{H}_T^2 , as opposed to standard neural operators which are maps functions on a Euclidean domain to functions in another Euclidean domain.

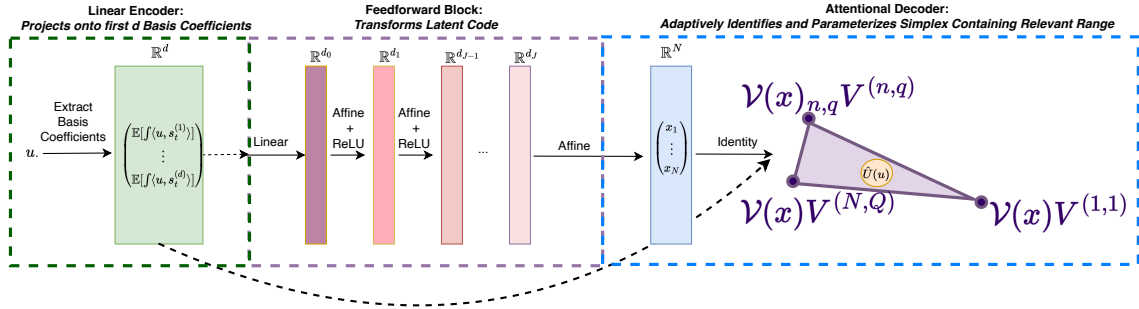


Figure 1: **Attentional Neural Operator Workflow:** Our *attentional* neural operator model maps controls u . to square-integrable \mathbb{F} -adapted processes $\hat{U}(u)$ in three phases. First, the (input) control is linearly projected onto the wavelet-like (in time) Wiener Chaos-like (in space) orthonormal basis of \mathcal{H}_T^2 . Next, the basis coefficients are transformed by a feedforward neural network (MLP). Lastly, the basis coefficients are used to identify extremal points in a simplex in \mathcal{H}_T^2 and the outputs of the MLP are used to parameterize a prediction in its relative interior.

The neural operators in Figure 1 can be formalized as follows.

Definition 6 (Attentional Neural Operator). *Fix a trainable activation function $\sigma \in C(\mathbb{R}^2)$, an encoding dimension $d \in \mathbb{N}_+$, a queries dimension $Q \in \mathbb{N}_+$, a values dimension $N \in \mathbb{N}_+$, depth and*

width parameters $J, W \in \mathbb{N}_+$, and an orthonormal basis \mathcal{S} of \mathcal{H}_T^2 .

The set $\mathcal{NO}_{N,Q,d,J,W:\mathcal{S}}^\sigma$ consists of all (non-linear) operators $U : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ admitting the following representation: for each $u. \in \mathcal{H}_T^2$

$$\begin{aligned}
 U(u.) &\stackrel{\text{def.}}{=} \mathcal{D}(f \circ \mathcal{E}(u.), \mathcal{E}(u.)) \\
 \mathcal{E}(u.) &\stackrel{\text{def.}}{=} \left(\mathbb{E} \left[\int_0^T \langle u_t, s_t^{(i)} \rangle dt \right] \right)_{i=1}^d \quad (\text{encoder}) \\
 \mathcal{D}(w, x) &\stackrel{\text{def.}}{=} \sum_{n=1}^N \text{softmax}(w)_n \sum_{q=1}^Q \mathcal{V}_{n,q}(x) V^{(n,q)}, \quad (\mathcal{H}_T^2\text{-attentional decoder})
 \end{aligned}$$

where $\{s^{(i)}\}_{i=1}^d, \{V^{(n,q)}\}_{n,q=1}^{N,Q} \subset \mathcal{S}$, $f \in \mathcal{NN}_{J,W:d,N}^\sigma$, $\mathcal{V} \in \mathcal{NN}_{J,W:d,N \times Q}$.

The number of (non-zero trainable) parameters defining the neural operators U is at most

$$\underbrace{JW^2}_{MLP(f)} + \underbrace{NQ}_{MLP(\mathcal{D})}$$

The map \mathcal{E} is called an encoder, \mathcal{D} is called an attention-based decoder, $(\mathcal{V}, \{V^{(n,q)}\}_{n,q=1}^{N,Q})$ are called values, and U is called a attentional neural operator.

We henceforth take \mathcal{S} to be (2) and denote its elements $\{s^{(i)}\}_{i=1}^\infty$. Thus, we write $\mathcal{NO}_{N,Q,d,J,W}^\sigma$ in place of $\mathcal{NO}_{N,Q,d,J,W:\mathcal{S}}^\sigma$.

Link Between Our Decoder, Attention, and Transformers Before moving on, we discuss the relationship between the decoder of our neural operator and the attention layer in Bahdanau et al. (2015) in the standard transformer networks of Vaswani et al. (2017). For simplicity, we focus on a single attention head, not the multi-head attention mechanism.

The standard attention mechanism attention maps N, d -dimensional vectors x_1, \dots, x_N , seen as a $N \times d$ matrix $x = (x_n)_{n=1}^d$ to another $N \times \tilde{d}$ matrix $\text{attention}(x)$ for some $\tilde{d} \in \mathbb{N}$. This attention mechanism operates in two phases, first, it extracts *contextual weights* w^x in the N -simplex via

$$w^x \stackrel{\text{def.}}{=} \text{softmax} \left(\left(\langle x_n Q, x_j K \rangle / \sqrt{\tilde{d}} \right)_{j=1}^{\tilde{d}} \right)_{n=1}^N \quad (12)$$

where the key and query matrices K and Q are $\tilde{d} \times d$ matrices. Using these weights, one defines a set of N values v_1, \dots, v_N , depending on the given input x , by

$$v_n^x \stackrel{\text{def.}}{=} V x_n \quad (13)$$

for $n = 1, \dots, N$. The attention mechanism then uses the weights w^x to tweak the contextual importance of each value v_1^x, \dots, v_N^x when generating a weighted prediction by

$$\text{attention}(x) \stackrel{\text{def.}}{=} \sum_{n=1}^N w_n^x v_n^x. \quad (14)$$

In this way, the standard attention mechanism parameterizes the interior of the convex hull of the values v_1^x, \dots, v_N^x using the softmax weights of any input.

If we streamline the contextual weight extraction step in (12), by simply an input weight $w \in \mathbb{R}^N$ which is mapped to contextual importance weights w by only using the softmax function itself. One can relax the dependence on the contextual values v_1^x, \dots, v_N^x in (13) to be vectors in the target space, i.e. $\mathbb{R}^{N \times \tilde{d}}$ for classical transformers and \mathcal{H}_T^2 for our neural operator, depending non-linearly on the given input. Here, we parameterize this non-linear dependence using a values neural network

\mathcal{V} depending on the encoding instead of a values matrix V ; doing so, one arrives at the following construction, which is precisely our \mathcal{H}_T^2 -attention

$$\mathcal{D}(w) \stackrel{\text{def.}}{=} \sum_{n=1}^N \underbrace{\text{softmax}(w)_n}_{\text{Contextual Weights (12)}} \underbrace{\sum_{q=1}^Q \mathcal{V}(u.) V_{\text{Contextual Values}}^{(n,q)}}_{n,q} \quad (13)$$

Thus, one can interpret our \mathcal{H}_T^2 -attention as an infinite-dimensional analogue of the standard attention mechanism of Bahdanau et al. (2015). Note that, since \mathcal{D} receives inputs from an MLP and since MLPs can approximately implement continuous functions then there is no need to rely on a vector of inner-products such as $(\langle Qx_n, Kx_j/\sqrt{d} \rangle_{j=1}^N)$, in (12), since that can be approximately implemented by the MLP in principle; by the universal approximation theorem.

Instead, we use the inner products to obtain low-dimensional approximate representations of any given input (control) in our encoding layer (encoder). Thus, one can view our neural operator as an infinite-dimensional take on the transformer network model.

4. Main Results

Our first result shows that that neural operators are rich enough to approximately minimize the response functional of the follower, to arbitrary precision on any given compact set of actions which the follower can take.

Theorem 7 (ε -Optimal Response Operators). *Under Assumptions 2 and 4, for each compact $\mathcal{K}_0 \subseteq \mathcal{U}_0$, and each $\varepsilon > 0$ there is an encoding dimension $d \in \mathbb{N}_+$ and a neural operator $\hat{U} \in \mathcal{NO} : \mathcal{U}_0 \rightarrow \mathcal{U}_1$ satisfying*

$$\sup_{u^0 \in \mathcal{K}_0} \|U^*(u^0) - \hat{U}(u^0)\|_{\mathcal{H}_T^2} \leq \varepsilon. \quad (15)$$

4.1 An Unsupervised Objective Function - Bypassing Computing U^*

Theorem 7 alone does not guarantee that approximations produce Stackelberg equilibria. Our second main result guarantees that approximately playing any such Stackelberg games is enough to yield approximate Stackelberg equilibria.

Additionally, theorem 7 guarantees that the solution operator to the Stackelberg game can be approximately implemented on compact sets, to arbitrary precision. However, it does not describe how one could train such a network in practice. Indeed it seems most natural to minimize the leader's loss J_0 for any given u^0 in the relevant compact set \mathcal{K}_0 . However, u^0 may not be *exactly implementable* in practice, instead one could consider minimizing J_0 where u^0 is replaced by a finite-dimensional (e.g. linear) approximation $\hat{u}^0 \stackrel{\text{def.}}{=} p_d(u^0)$ where $p_d : \mathcal{H}_T^2 \rightarrow \text{span}\{s_i\}_{i=1}^d$ is the orthogonal projection; i.e. $p_d(\sum_{i=1}^{\infty} \beta_i s_i) = \sum_{i=1}^d \beta_i s_i$ for all $u = \sum_{i=1}^{\infty} \beta_i s_i \in \mathcal{H}_T^2$. Thus, a natural and tractable objective would be to minimize

$$\min_U \min_{\hat{u}_d^0 \in p_d(\mathcal{K}_0)} J_0(\hat{u}_d^0, U(\hat{u}_d^0)) \quad (16)$$

when training the attentional neural operator where \hat{U} is minimized over a compact class of neural operators. Once a minimizer \hat{U} of (16) is identified, it can then be used to approximate the optimal action of the leader by minimizing the following objective

$$\min_{\hat{u}_d^0 \in p_d(\mathcal{K}_0)} J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0)) \quad (17)$$

Our following result suggests that (17) can be used, after training \hat{U} , as an *unsupervised* objective function, which is small only when \hat{U} is has correctly approximated to optimal response U^* . Unlike the supersized criterion in (15), which requires us to knowing pairs of u^0 and best responses $U^*(u^0)$, (15) can be minimized without having to first compute U^* . Moreover, \hat{U} is only optimized on a finite-dimensional subspace of our space of controls.

Theorem 8 (The Unsupervised Objective Function (17) Detects Optimality). *The following hold in the setting of Theorem 7.*

- (i) For every $\delta > 0$ there exist $\epsilon > 0$, and $d \in \mathbb{N}_+$ such that if \hat{U} satisfies (15), then the following hold

$$\sup_{u^0 \in \mathcal{K}_0} |J_0(u^0, U^*(u^0)) - J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))| < \delta. \quad (18)$$

- (ii) Moreover, there is a $\hat{u}_d^0 = \sum_{i=1}^d \beta_i s_i \in \mathcal{K}_0$ such that the pair $(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))$ is an ϵ -Stackelberg equilibrium in the sense that the pair $(\hat{u}_d^0, \hat{U}) \in \mathcal{K}_0 \times \bar{\mathcal{U}}_1$ satisfies for all $(u^0, u^1) \in \mathcal{K}_0 \times \mathcal{U}_1$ the inequalities

$$\begin{aligned} J_1(u^0, \hat{U}(u^0)) &\leq J_1(u^0, u^1) + \epsilon \\ J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0)) &\leq J_0(u_0, \hat{U}(u_0)) + \epsilon. \end{aligned}$$

As with classical uniform approximation results in deep learning, we worked on a compact in the space of square-integrable \mathbb{F} -predictable processes. Though this is standard in the approximation theory literature, it is not so in game theory. However, reducing the problem to compact \mathcal{K}_0 is asymptotically coherent in the sense that if we take compact sets $\mathcal{K}_n \subset \mathcal{U}_0$ so that $\cup_n \mathcal{K}_n = \mathcal{U}_0$. Then, for any $\delta_n \downarrow 0$ and (\hat{u}_d^n, \hat{U}^n) satisfying (18) with δ_n , we have that

$$\inf_{u^0 \in \mathcal{U}_0} J_0(u^0, U^*(u^0)) = \lim_{n \rightarrow \infty} J_0(\hat{u}_d^n, \hat{U}^n(\hat{u}_d^n)).$$

Thus, by taking the compact set \mathcal{K}_n and the NO larger, we approximate the optimum of the effective value of the leader $\inf_{u^0 \in \mathcal{U}_0} J_0(u^0, U^*(u^0))$. Note that if \mathcal{U}_0 is not compact or if the problem of the leader does not have additional properties such as convexity or coercivity of $u^0 \mapsto J_0(u^0, U^*(u^0))$, see e.g. (Dal Maso, 1993, Theorem 7.12), one cannot easily claim the existence of the optimizer of $\inf_{u^0 \in \mathcal{U}_0} J_0(u^0, U^*(u^0))$ or the convergence of the family (\hat{u}_d^n, \hat{U}^n) . However, these additional structural assumptions are not needed for the purposes of computationally approximating the optimal value of the leader.

More insight can be gained into the inner workings of these results by inspecting the steps in their derivations, at least at a high level. The next section shows the overarching structure of an argument which one can use to derive any such result.

4.2 Rates For Perturbations of Closed-Loop Solutions to Linearized Game

This section shows that, contrary to the general approximation rates for neural operators (see e.g. Lanthaler et al. (2022a) or Galimberti et al. (2022)), which may be highly inefficient due to the infinite-dimensional nature of the involved spaces, there are stylized conditions which allow for efficient approximation rates of the optimal response map by our attentional neural operator.

Building on the strategy of Shi et al. (2023), we focus on the case where the leader begins by considering a proxy/ansatz/pre-trained control, \bar{u} . Suppose that $\bar{u}_0 \in \mathcal{H}_T^2$ is a control for the *leader*. Suppose that the leader is comfortable deviating from their strategy to an open-loop control, but only marginally through “small residual perturbations”. If those deviations are not too large, efficient approximation rates are possible. More precisely, consider the following situation.

Assumption 9 (Perturbations of an Ansatz). Let $C \geq 0$, $r > 0$, $\bar{u} \in \mathcal{H}_T^2$, and suppose that Assumption 4 holds. Let \mathcal{K}_0 consist of all $u \in \mathcal{H}_T^2$ for which:

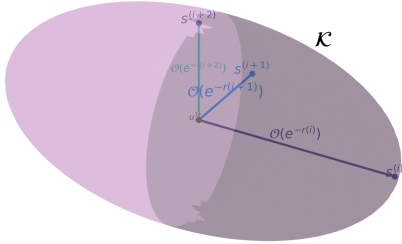
- (i) $|\langle u - \bar{u}, s_i \rangle| \leq C e^{-ri}$,
- (ii) $|\langle U^*(u - \bar{u}), s_i \rangle| \leq C e^{-ri}$.

Remark 10. Note that, for every $C, r \geq 0$, $\bar{u} \in \mathcal{H}_T^2$ the set \mathcal{K}_0 is non-empty since $\bar{u} \in \mathcal{K}_0$.

In both of the following examples, we consider the sets of “residual controls”. Let $C, r > 0$ and define $B_{C,r}(u)$ to be the “exponentially ellipsoidal” set of *open-loop* “residual controls” relative to a control $\bar{u} \in \mathcal{H}_T^2$ to consist of all $u \in \mathcal{H}_T^2$ such that $v \stackrel{\text{def.}}{=} u - \bar{u} = \sum_{i=1}^{\infty} \beta_i s_i \in \mathcal{H}_T^2$ satisfies

$$\max\{|\langle v, s_i \rangle|, |\langle U^*(v), s_i \rangle|\} \leq C e^{-ri} \quad (\forall i \in \mathbb{N}_+). \quad (19)$$

Illustrated by Figure 2, our primary example of a compact set satisfying the perturbation of the Ansatz assumption above is given by first solving a linear-quadratic proxy of the general (non-linear) Stackelberg game. Then, once an optimal *feedback* control for the leader is derived, for the linear-quadratic proxy of our (non-linear) Stackelberg game, we build \mathcal{K}_0 by adding residual open-loop controls in $B_{C,r}$. Intuitively, these residual loops provide added freedom to the closed-loop control optimizing the linear-quadratic proxy required when playing the Stackelberg game.



Explanation of Figure: The set \mathcal{K}_0 in Example 8 is a (compact) ellipsoidal region in \mathcal{H}_T^2 consisting of open-loop controls u which are a small perturbation of a base strategy/control $u^{1:*}$. The base strategy $u^{1:*}$ is the solution to a linearization, see (20)-(21), of the dynamic Stackelberg game. The perturbations u of $u^{1:*}$ in \mathcal{K}_0 are built by adding a small residual term in each of the basic directions $s^{(i)} \in \mathcal{S}$ where we allow for possibly large perturbations of the linearized strategy $u^{1:*}$ for the “low-frequency directions” (i.e. for small values of i) and much smaller perturbations of the linearized strategy for “high-frequency directions” (i.e. for small i). The key subtlety is that the size of the perturbations must decay exponentially in i .

Figure 2: The ellipsoidal compact set K of Example 8.

Example 8 (Perturbations of Feedback Control For Linearized Problem - Pt. I). Consider a finite subset $\{(x_n, v_n^0, v_n^1)\}_{n=1}^N \subset \mathbb{R}^d$ and let A, B_1, B_2, C, D_1, D_2 be matrices minimizing the following MSE problem over all matrices of compatible dimension

$$\sum_{n=1}^N \|f(x_n, v_n^0, v_n^1) - \underbrace{(Ax_n + B_1 v_n^0 + B_2 v_n^1)}_{\text{lin. approx. drift}}\|^2 + \|\sigma(x_n, v_n^0, v_n^1) - \underbrace{Cx_n + D_1 v_n^0 + D_2 v_n^1}_{\text{lin. approx. diff.}}\|^2.$$

Consider the (controlled) “linearized” state-space process $X^{\text{lin}} \stackrel{\text{def.}}{=} (X_t^{\text{lin}})_{t \geq 0}$ given by

$$dX_t^{\text{lin}} = (AX_t^{\text{lin}} + B_1 u_t^0 + B_2 u_t^1)dt + (CX_t^{\text{lin}} + D_1 u_t^0 + D_2 u_t^1)dW_t \quad (20)$$

Further, assume that $Q_1, Q_2, R_1, R_2, G_1, G_2$ are matrices minimizing the following MSE problem over all matrices of compatible dimension

$$\sum_{n=1}^N \sum_{i=1}^2 \|L_i(x_n, v_n^0, v_n^1) - \underbrace{((Q_i x_n)^\top x_n + (R_i v_n^i)^\top v_n^i)}_{\text{lin. approx. running cost}}\|^2 + \|\underbrace{g_i(x_n) - (G_i x_n)^\top x_n}_{\text{lin. approx. terminal}}\|^2.$$

Under (Yong, 2002, Assumptions (DI) and (H1)1035) the computations on (Yong, 2002, ages 1034 and 1035) together with (Yong, 2002, Proposition 2.2 and Theorem 2.), they imply that there is an optimal control $u^{0:\star}$ for the leader minimizing the approximate objective function

$$\mathbb{E} \left[\int_0^T L_0(X_t, u_t^0, u_t^1) dt + (G_i X_T)^\top X_T \right] \quad (21)$$

across all controls in \mathcal{H}_T^2 . Moreover, as shown in (Yong, 2002, Equation (5.12)), $u^{1:\star}$ is given by

$$u_t^{1:\star} = -R_1^{-1} B_1^\top p_t,$$

where $(p_t, q_t)_{t \geq 0}$ solve the FBSDE

$$\begin{aligned} dp_t &= p(A^\top p + Q_0 X_t) dt + q_t dW_t \\ p_T &= G_0 X_T. \end{aligned}$$

Let \mathcal{K}_0 be the set of controls for the leader $u \in \mathcal{H}_T^2$ for which there exists some open-loop “residual perturbations” $v \in B_{C,r}(u^{1:\star})$ such that

$$u = u^{1:\star} + v.$$

By construction, \mathcal{K}_0 satisfies Assumption 9.

The set of open-loop perturbations of the closed-loop optimal control for the linearized game, just described in Example 8, can be efficiently represented using few dimensions, and likewise for its image under the optimal response map U^\star . Nevertheless, the set need not be coverable by few controls; that is, it may still have high metric entropy (see e.g. Lorentz (1966) for an exposée). We, therefore, adopt a technique used in statistical learning/empirical process theory to construct classes of VC-dimension proportional to the metric entropy of high-dimensional balls (see e.g. (van der Vaart and Wellner, 1996, Theorem 2.7.11)). Namely, we postulate the existence of a latent unknown Lipschitz parameterization of a latent low-dimensional “manifold” of the set of open-loop controls \mathcal{K}_0 constructed in Example 8.

Example 9 (Perturbations of Feedback Control For Linearized Problem - Pt. II). *Fix a latent dimension $d \in \mathbb{N}_+$, a latent parameter space $B_d \stackrel{\text{def}}{=} \{x \in \mathbb{R}^d : \|x\| \leq 1\}$, and a latent parameterization given by a 1-Lipschitz map $\pi : \mathbb{R}^d \rightarrow \mathcal{K}_0$ with the property that $\pi(0) = u^{1:\star}$; i.e. parameterizing a latent low-dimensional structure which perturbs the optimal control for the linear-quadratic approximation of the general Stackelberg games. Let $K_{d,\pi} \stackrel{\text{def}}{=} \pi(B_d)$. By construction $K_{d,\pi} \subseteq \mathcal{K}_0$; thus, $K_{d,\pi}$ also satisfies Assumption 9. Finally, by (van der Vaart and Wellner, 1996, Theorem 2.7.11)² and (Lorentz et al., 1996, Proposition 15.1.3), we have that: for each $\varepsilon > 0$ there exists at-most $N \leq 3^d (\sqrt{d}2/\varepsilon)^d$ controls $\{u^{(n)}\}_{n=1}^N$ forming an ε -cover of $K_{d,\pi}$; that is,*

$$\sup_{u \in K_{d,\pi}} \min_{n=1, \dots, N} \mathbb{E} \left[\int_0^T \|v_t - u_t\|^2 dt \right] \leq \varepsilon. \quad (22)$$

In other words, if the latent dimension d is “small”, then $K_{d,\pi}$ is small in metric entropy. In turn, the parameter N in our transformer can be taken to also be proportionally small (see Theorem 11).

When the conditions of Assumption 9 are met, we are able to guarantee that the approximating attention neural operator \hat{U} to the best-response map U^\star , given by Theorem 7, is determined by relatively few parameters. This is the content of our last main result.

2. We have used the upper-bound of the ε -bracketing numbers on the $\varepsilon/2$ -covering numbers (see e.g. (van der Vaart and Wellner, 1996, page 84))

Theorem 11 (Efficient Approximation of the Best-Response Map). *Consider the setting of Theorem 7 and additionally suppose that \mathcal{K}_0 satisfies Assumption 9. For every $\varepsilon > 0$, the conclusion of Theorem 7 holds and there exists a neural operator $\hat{U} \in \mathcal{NO} : \mathcal{U}_0 \rightarrow \mathcal{U}_1$ satisfying the uniform estimate in (15) whose depth, width, decoding dimension, encoding dimension, and attentional complexity are bound above by the estimates in Table 1.*

Depth	Width	Decoding Dim. (Q)	Encoding Dim.	Att. Complexity (N)
$\mathcal{O}(\ln(\varepsilon^{-1/r}) \varepsilon^{-\ln(C)/r})$	$\mathcal{O}(\varepsilon^{-\ln(C)/r})$	$\mathcal{O}(\ln(\varepsilon^{-1/r}))$	$\mathcal{O}(\varepsilon^{1/(1-r)})$	$(\tilde{c} \varepsilon^{-1} \ln(\varepsilon^{-1/r})^{1/2})^{c(\ln(\varepsilon^{-1/r}))}$

Table 1: Parametric Complexity of Attentional Neural Operator Approximation of the Best Response Map U^* over the compact set \mathcal{K}_0 of Assumption 9, with trainable activation function σ of (7). Here, $C, r > 0$ are the constant defined in Assumption 9 and $c > 0$ is an absolute constant.

5. Overview of Proof for Theorems 7 and 8

The derivation of Theorem 7 is undertaken in two steps. First, one must show that, under mild conditions, the best response operator set \mathcal{R} is single-valued and there is a continuous selection therein, i.e., the best response operator U^* is a well-defined continuous non-linear operator. This is key since our neural operator architectures are continuous, and classes of continuous functions cannot uniformly approximate discontinuous functions by the Uniform Limit theorem, see e.g. (Munkres, 2000, Theorem 21.6).

Since we now know that the best response operator is well-defined and continuous, we need to show that our attentional neural operator class has the power to approximate it or, more generally, functions of the same regularity. This is guaranteed by the following Universal Approximation theorem, guaranteeing that our class of neural operators can approximate, uniformly on compacta, any continuous non-linear operator with square-integrable \mathbb{F} -adapted processes as inputs and outputs.

Theorem 12 (Universal Approximation for Operators Between \mathbb{F} -Adapted 2-Integrable Processes). *For every (non-empty) compact subset $\mathcal{K}_0 \subset \mathcal{H}_T^2$, each continuous “target” function $f : \mathcal{K}_0 \rightarrow \mathcal{H}_T^2$, and every “approximation error” $\varepsilon > 0$ there exists an attentional neural operator $\hat{F} : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ satisfying*

$$\max_{u \in \mathcal{K}_0} \|f(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} \leq \varepsilon.$$

Lemma B.5 shows that U^* is 1/2-Hölder continuous. Together, Lemma B.5 and Theorem 12 imply that the best response map can be approximated. The main step in using these results to deduce Theorem 8 lies in the continuous dependence of the leader’s response on the best response.

6. Examples of Stackelberg Games Satisfying Assumption 4

An easily verifiable sufficient condition guarantees that Assumption 4 holds. The condition requires that the Hamiltonian of the follower satisfies a basic level of strong convexity; where the Hamiltonian of the follower is given by

$$H_1(x, u^0, u^1, p_1, q_1) \stackrel{\text{def.}}{=} p_1^\top f(x, u^0, u^1) + \text{Tr}(q_1^\top \sigma(x, u^0, u^1)) + L_1(x, u^0, u^1).$$

Proposition 13 (Hamiltonian Strong Convexity Implies Hölder Continuity). *Assume Assumption 2 and that there exists $\kappa > 0$ such that the functions*

$$\begin{aligned} (x, u^1) \in \mathbb{R}^d &\mapsto H_1(x, u^0, u^1, p_1, q_1) - \frac{\kappa}{2}|u^1|^2 \\ x \in \mathbb{R}^d &\mapsto g_1(x) \end{aligned}$$

are convex for all values of other variables. Assume also that $\nabla_x L_1$ and $\nabla_x g_1$ are Lipschitz continuous in their variables. Then, $\mathcal{R}(u^0)$ is single valued and its unique element $\{U_t^*(u^0)\} = \mathcal{R}(u^0)$ have the following continuous dependence

$$\begin{aligned} & \frac{\kappa_v}{2} \mathbb{E} \left[\int_0^T |U_t^*(u^0) - U_t^*(\tilde{u}^0)|^2 dt \right] \\ & \leq J_1(\tilde{u}^0, U_t^*(\tilde{u}^0)) - J_1(u^0, U_t^*(u^0)) + J_1(u^0, U_t^*(\tilde{u}^0)) - J_1(\tilde{u}^0, U_t^*(\tilde{u}^0)) \end{aligned} \quad (23)$$

and the Assumption 4 holds.

Remark 14 (Deriving Alternative Sufficient Conditions). 1) Another case where one can check the Assumption 4 would be to directly rely on (Peng and Wu, 1999, Section 3) and check their Assumption H3.1 uniformly in u^0 and rely on the existence of solution to a fully coupled FBSDE. Note that the case in Proposition 13 deviates from (Peng and Wu, 1999, Section 3). Indeed, our proof is self-contained and only requires the existence of a solution for a BSDE (and not FBSDE) (41) for a given u^0, u^1, \tilde{u}^1 which is a trivial problem. Although it is out of scope of this work, one can then check that the optimally controlled state together with the solution to (41) leads to a solution to the same FBSDE.

The following is a broad class of functions satisfying the assumption of Proposition 13.

Example 10. Let A, A^σ be matrices of dimensions $d \times d$, B, B^σ be matrices of dimensions $d \times d^1$, $(C, C^\sigma, C^L) : \mathbb{R}^{d_0} \rightarrow \mathbb{R} \times \mathbb{R} \times \mathbb{R}_+$ and $(D, D^\sigma) : \mathbb{R}^{d_0} \rightarrow (\mathbb{R}^d)^2$ be Lipschitz functions. Let $L_1^1 : \mathbb{R}^{d \times d_1} \rightarrow \mathbb{R}$ convex in x and strongly convex in u^1 and $L_1^2 : \mathbb{R}^{d_0} \rightarrow \mathbb{R}_+$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}$ have Lipschitz gradients (in x). Define

$$\begin{aligned} f(x, u^0, u^1) & \stackrel{\text{def.}}{=} C(u^0) (Ax + Bu^1) + D(u^0) \\ \sigma(x, u^0, u^1) & \stackrel{\text{def.}}{=} C^\sigma(u^0) (A^\sigma x + B^\sigma u^1) + D^\sigma(u^0) \\ L_1(x, u^0, u^1) & \stackrel{\text{def.}}{=} C^L(u^0) L_1^1(x, u^1) + L_1^2(u^0). \end{aligned}$$

The Hamiltonian associated to these f, σ , and L_1 and the map g satisfy the assumption of Proposition 13.

In linear-quadratic Stackelberg games, e.g., Yong (2002), the coefficients $f(x, u^0, u^1)$ and $\sigma(x, u^0, u^1)$ are linear in (x, u^0, u^1) , and the costs $L_i(x, u^0, u^1)$ and $g_i(x)$ $i = 0, 1$ are quadratic forms of (x, u^0, u^1) . With an additional assumption that the weight matrix of the quadratic form for u^1 in L_1 is positive definite the assumption of Proposition 13 is satisfied.

7. Conclusion

We show that, the best response map for a broad class of dynamic Stackelberg games with random effects, can be approximated by neural operators (Theorem 7). This implies that Stackelberg games may be solved (approximate computation of equilibria in Theorem 8), and we can explicitly describe how they are played (approximate representation of best response map) without making highly stylized conditions needed in classical results needed to derive analytic expressions for the best response and the equilibrium itself (see e.g. Example 8). Thus, our approach allows one to specify realistic dynamics and action sets, and obtain an approximate solution.

We further showed that if the space of actions for the leader consists of perturbations of the optimal solution for a linearized version of the Stackelberg game, then one can approximate the best response map (for the general game) for the follower much more efficiently (Theorem 11).

Future Research

Our universal approximation theorem for non-linear operators between spaces of \mathbb{F} -adapted square-integrable processes, namely Theorem 12, is quantitative. Though, as with most approximation theorems between infinite-dimensional Hilbert spaces, the approximation rates are particularly insightful. Nevertheless, we find that if one uses a trainable variant of the “super-expressive” activation function of Zhang et al. (2022) (see (10)) when defining our neural operator and if the compact subset on which the non-linear operator is sufficiently close to being finite-dimensional (see Definition B.8) then exponential approximation rates can be achieved (see Table 2). The main challenge is then to verify that the non-linear operator being approximated and the relevant compact set of adapted open-loop controls satisfy the required compatibility conditions. Constructing examples of Stackelberg games with these properties, i.e. games whose best-response operator is efficiently approximable, is a highly non-trivial but interesting project and is the objective of our future research.

Nevertheless, we include the relevant quantitative universal approximation theorem and conditions for exponential approximation rates in the appendix of this paper so that these results may be used in other operator learning problem in stochastic analysis and its applications; e.g. to economics and finance.

A. Counterexamples

The effective optimization problem of the leader can be discontinuous

In Proposition 13, we show that a sufficient condition for the continuous dependence of the optimal response function U^* is the strong convexity of the problem of the follower. We now provide the example of a deterministic game that shows that both U^* and $u^0 \mapsto J_0(u^0, U^*(u^0))$ might not be continuous without the strong convexity assumption. Consider the deterministic single period game where the controls are $u^i \in [0, 1]$. Fix the loss functions $l_1(u^0, u^1) = u^0 u^1$ and $l_0(u^0, u^1) = -u^1$ so that the leader’s problem is

$$\inf_{u^0 \in [0, 1]} \inf_{u^1 \in \mathcal{R}(u^0)} l_0(u^0, u^1) \quad (24)$$

and $\mathcal{R}(u^0) \stackrel{\text{def.}}{=} \{u^1 \in [0, 1] : l_1(u^0, u^1) \leq \inf l_1(u^0, \cdot)\}$. Note that l_1 is convex in u_1 but it lacks the strong convexity in u_1 needed in Proposition 13. Clearly $\mathcal{R}(u^0) = \{0\}$ if $u^0 > 0$ and $\mathcal{R}(0) = [0, 1]$. Thus,

$$\inf_{u^1 \in \mathcal{R}(u^0)} l_0(u^0, u^1) = -1_{\{u^0=0\}}$$

which is discontinuous in u_0 .

B. Proofs

B.1 Continuous Dependence of J_i

For the next result, it is convenient to recall that the \mathcal{H}_T^∞ norm of a process X in \mathcal{H}_T^2 is given by

$$\|X\|_{\mathcal{H}_T^\infty} \stackrel{\text{def.}}{=} \mathbb{E} \left[\sup_{0 \leq t \leq T} |X_t| \right].$$

Lemma B.1. *Under Assumption 2, we have the following uniform continuity guarantees*

$$\left\| \|X^{u^0, u^1} - X^{u^0, \tilde{u}^1}\|^2 \right\|_{\mathcal{H}_T^\infty} \leq C \|u^1 - \tilde{u}^1\|_{\mathcal{H}_T^2}^2, \quad (25)$$

$$\left\| \|X^{u^0, u^1} - X^{\tilde{u}^0, u^1}\|^2 \right\|_{\mathcal{H}_T^\infty} \leq C \|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}^2, \quad (26)$$

where $C \geq 0$ is a constant that depends on T and the Lipschitz constant K in Assumption 2; moreover,

$$\|X^{u^0, u^1} - X^{u^0, \tilde{u}^1}\|_{\mathcal{H}_T^\infty} \leq C \|u^1 - \tilde{u}^1\|_{\mathcal{H}_T^2}, \quad (27)$$

$$\|X^{u^0, u^1} - X^{\tilde{u}^0, u^1}\|_{\mathcal{H}_T^\infty} \leq C \|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}. \quad (28)$$

Proof. We prove (25), noting that the derivation of (26) is carried out in a nearly identical similar manner. Let X and \tilde{X} be the strong solution of (4) under (u^0, u^1) and $(\tilde{u}^0, \tilde{u}^1)$, respectively, i.e., $X = X^{u^0, u^1}$, $\tilde{X} = X^{\tilde{u}^0, \tilde{u}^1}$. Denote $h(s) = h(X_s, u_s^0, u_s^1)$ and $\tilde{h}(s) = h(\tilde{X}_s, u_s^0, \tilde{u}_s^1)$, where h is as defined in Assumption 2.

By Jensen's inequality, we have

$$\begin{aligned} |X_r - \tilde{X}_r|^2 &= \left| \int_0^r f(s) - \tilde{f}(s) ds + \int_0^r \sigma(s) - \tilde{\sigma}(s) dW_s \right|^2 \\ &\leq 2 \left| \int_0^r f(s) - \tilde{f}(s) ds \right|^2 + 2 \left| \int_0^r \sigma(s) - \tilde{\sigma}(s) dW_s \right|^2 \\ &\leq 2 \int_0^r |f(s) - \tilde{f}(s)|^2 ds + 2 \left| \int_0^r \sigma(s) - \tilde{\sigma}(s) dW_s \right|^2 \end{aligned}$$

We deduce that

$$\sup_{0 \leq r \leq t} |X_r - \tilde{X}_r|^2 \leq 2 \int_0^t |f(s) - \tilde{f}(s)|^2 ds + 2 \sup_{0 \leq r \leq t} \left| \int_0^r \sigma(s) - \tilde{\sigma}(s) dW_s \right|^2. \quad (29)$$

By Burkholder-Davis-Gundy's inequality (Cohen and Elliott, 2015, Theorem 11.5.5), Tonelli's theorem (Cohen and Elliott, 2015, Theorem 1.4.6), and Assumption 2, we have

$$\begin{aligned} \mathbb{E} \sup_{0 \leq r \leq t} \left| \int_0^r \sigma(s) - \tilde{\sigma}(s) dW_s \right|^2 &\leq C \cdot \mathbb{E} \int_0^t |\sigma(s) - \tilde{\sigma}(s)|^2 ds \\ &\leq C \cdot \int_0^t \mathbb{E} \sup_{0 \leq r \leq s} |X_r - \tilde{X}_r|^2 + |u_s^1 - \tilde{u}_s^1|^2 ds. \end{aligned} \quad (30)$$

By Tonelli's Theorem and Assumption 2, we have

$$\begin{aligned} \mathbb{E} \int_0^t |f(s) - \tilde{f}(s)|^2 ds &\leq C \cdot \int_0^t |X_s - \tilde{X}_s|^2 + |u_s^1 - \tilde{u}_s^1|^2 ds \\ &\leq C \cdot \int_0^t \mathbb{E} \sup_{0 \leq r \leq s} |X_r - \tilde{X}_r|^2 + |u_s^1 - \tilde{u}_s^1|^2 ds \end{aligned} \quad (31)$$

Combining (29), (30), and (31), we obtain

$$\mathbb{E} \sup_{0 \leq r \leq t} |X_r - \tilde{X}_r|^2 \leq C \cdot \int_0^t \mathbb{E} \sup_{0 \leq r \leq s} |X_s - \tilde{X}_s|^2 ds + C \cdot \mathbb{E} \int_0^t |u_s^1 - \tilde{u}_s^1|^2 ds.$$

From the above inequality and Grönwall's inequality we obtain

$$\mathbb{E} \left[\sup_{0 \leq t \leq T} |X_t^{u^0, u^1} - X_t^{u^0, \tilde{u}^1}|^2 \right] \leq C \cdot \mathbb{E} \int_0^T |u_t^1 - \tilde{u}_t^1|^2 dt. \quad (32)$$

Arguing nearly identically, we may also obtain

$$\mathbb{E} \left[\sup_{0 \leq t \leq T} |X_t^{u^0, u^1} - X_t^{\tilde{u}^0, u^1}|^2 \right] \leq C \cdot \mathbb{E} \int_0^T |u_t^0 - \tilde{u}_t^0|^2 dt. \quad (33)$$

Note that the right-hand side of (32) (resp. (33)) is simply a scalar multiple (by a factor of $C \geq 0$) of the squared norm between u^1 and \tilde{u}^1 (resp. u^0 and \tilde{u}^0). By definition of the $\|\cdot\|_{\mathcal{H}_T^\infty}$ norm applied to the “squared difference processes” $|X^{u^0, u^1} - X^{\tilde{u}^0, \tilde{u}^1}|^2$ and $|X^{\tilde{u}^0, u^1} - X^{\tilde{u}^0, \tilde{u}^1}|^2$, (32) and (33) can be re-expressed as (25) and (26). (27) and (28) follows from (25) and (26) by Cauchy-Schwarz inequality; thus concluding our proof. \square

Lemma B.2. *The costs J_i , $i = 0, 1$ are Lipschitz in u^0 and u^1 such that for each $u^i, \tilde{u}^i \in \mathcal{U}_i \cap \mathcal{H}_T^2$, $i = 0, 1$,*

$$|J_i(u^0, u^1) - J_i(\tilde{u}^0, \tilde{u}^1)| \leq C (\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} + \|u^1 - \tilde{u}^1\|_{\mathcal{H}_T^2}), \quad (34)$$

where C is a constant depending on T and the Lipschitz constant K in Assumption 2.

Proof. For each $u^i, \tilde{u}^i \in \mathcal{U}_i \cap \mathcal{H}_T^2$, $i = 0, 1$, it follows from Assumption 2 that

$$\begin{aligned} & |J_i(u^0, u^1) - J_i(\tilde{u}^0, \tilde{u}^1)| \\ & \leq \mathbb{E} \left[\int_0^T |L_i(X_t^{u^0, u^1}, u_t^0, u_t^1) - L_i(X_t^{\tilde{u}^0, \tilde{u}^1}, \tilde{u}_t^0, \tilde{u}_t^1)| dt + |g_i(X_T^{u^0, u^1}) - g_i(X_T^{\tilde{u}^0, \tilde{u}^1})| \right] \\ & \leq \mathbb{E} \left[\int_0^T K(|X_t^{u^0, u^1} - X_t^{\tilde{u}^0, \tilde{u}^1}| + |u_t^0 - \tilde{u}_t^0| + |u_t^1 - \tilde{u}_t^1|) dt + K|X_T^{u^0, u^1} - X_T^{\tilde{u}^0, \tilde{u}^1}| \right] \\ & \leq K(1+T)\|X^{u^0, u^1} - X^{\tilde{u}^0, \tilde{u}^1}\|_{\mathcal{H}_T^\infty} + K(\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} + \|u^1 - \tilde{u}^1\|_{\mathcal{H}_T^2}), \end{aligned}$$

where the last inequality is due to the Cauchy-Schwarz inequality. By the triangle inequality and Lemma B.1, we have

$$\begin{aligned} \|X^{u^0, u^1} - X^{\tilde{u}^0, \tilde{u}^1}\|_{\mathcal{H}_T^\infty} & \leq 2 \left(\|X^{u^0, u^1} - X^{\tilde{u}^0, u^1}\|_{\mathcal{H}_T^\infty} + \|X^{\tilde{u}^0, u^1} - X^{\tilde{u}^0, \tilde{u}^1}\|_{\mathcal{H}_T^\infty} \right) \\ & \leq C (\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} + \|u^1 - \tilde{u}^1\|_{\mathcal{H}_T^2}). \end{aligned}$$

It then follows that

$$|J_1(u^0, u^1) - J_1(\tilde{u}^0, \tilde{u}^1)| \leq C (\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} + \|u^1 - \tilde{u}^1\|_{\mathcal{H}_T^2}). \quad \square$$

Lemma B.3. *Assume Assumption 2 and that for all $u^0 \in \mathcal{K}_0$, $\mathcal{R}(u^0)$ is not empty and choose $U^*(u^0) \in \mathcal{R}(u^0)$. Then, the map $u^0 \in \mathcal{K}_0 \mapsto J_1(u^0, U^*(u^0))$ is continuous; particularly, for each $u^0, \tilde{u}^0 \in \mathcal{U}_0 \cap \mathcal{H}_T^2$, it satisfies for the follower that*

$$|J_1(u^0, U^*(u^0)) - J_1(\tilde{u}^0, U^*(\tilde{u}^0))| \leq C \|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}. \quad (35)$$

The constant C is depend on T and the Lipschitz constant K in Assumption 2.

Remark B.4. *With the assumption of Lemma B.3, $u^0 \in \mathcal{K}_0 \mapsto J_1(u^0, U^*(u^0))$ is continuous but as shown in counterexample in Section A, $u^0 \in \mathcal{K}_0 \mapsto J_0(u^0, U^*(u^0))$ might fail to be continuous.*

Proof. By (34), we have for any $u \in \mathcal{U}_0 \cap \mathcal{H}_T^2$,

$$\begin{aligned} J_1(u^0, u) & \leq J_1(\tilde{u}^0, u) + |J_1(\tilde{u}^0, u) - J_1(u^0, u)| \\ & \leq J_1(\tilde{u}^0, u) + C\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}, \end{aligned}$$

and furthermore

$$\begin{aligned} J_1(u^0, U^*(u^0)) &= \inf_u J_1(u^0, u) \leq \inf_u J_1(\tilde{u}^0, u) + C\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} \\ &= J_1(\tilde{u}^0, U^{1,*}(\tilde{u}^0)) + C\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}. \end{aligned} \quad (36)$$

Exchanging the roles of $(u^0, U^{1,*}(u^0))$ and $(\tilde{u}^0, U^{1,*}(\tilde{u}^0))$ in (36), we have

$$J_1(\tilde{u}^0, U^*(\tilde{u}^0)) \leq J_1(u^0, U^*(u^0)) + C\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}. \quad (37)$$

Combining (36) and (37), we obtain (35). \square

Lemma B.5 (Hölder Regularity of the Leader's Utility on Optimal Response). *Under Assumptions 2 and 4, there exists $C, \alpha > 0$ depending only on K, T and the Holder norm in Assumption 4 so that for each $u^0, \tilde{u}^0 \in \mathcal{U}_0 \cap \mathcal{H}_T^2$, we have*

$$|J_0(u^0, U^*(u^0)) - J_0(\tilde{u}^0, U^*(\tilde{u}^0))| \leq \tilde{\omega}(\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2})$$

where $\tilde{\omega}(t) \stackrel{\text{def.}}{=} C \max\{|t|^\alpha, |t|\}$ for each $t \geq 0$.

Proof of Lemma B.5. By (36) and Assumption 4, we have

$$\begin{aligned} |J_0(u^0, U^*(u^0)) - J_0(\tilde{u}^0, U^*(\tilde{u}^0))| &\leq C\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} + C\|U^*(u^0) - U^*(\tilde{u}^0)\|_{\mathcal{H}_T^2} \\ &\leq C\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2} + (C\tilde{C})\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}^\alpha, \end{aligned} \quad (38)$$

for some constant $C \geq 0$ depending only on T, K and $\tilde{C}, \alpha \geq 0$ the Holder continuity parameters of the best response function given by Assumption 4. Define $C' \stackrel{\text{def.}}{=} C \max\{1, \tilde{C}\}/2$. Note that $a + b \leq 2 \max\{a, b\}$ for every $a, b \geq 0$ then (38) implies that

$$|J_0(u^0, U^*(u^0)) - J_0(\tilde{u}^0, U^*(\tilde{u}^0))| \leq C' \cdot \max\left\{\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}, \|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}^\alpha\right\}. \quad (39)$$

Define the modulus of continuity $\tilde{\omega}(t) \stackrel{\text{def.}}{=} C' \max\{|t|^\alpha, |t|\}$. Since $t \leq t^\alpha$ when $t \in [0, 1)$ and $t \geq t^\alpha$ when $t \in [1, \infty)$ then, (39) cleans up as

$$|J_0(u^0, U^*(u^0)) - J_0(\tilde{u}^0, U^*(\tilde{u}^0))| \leq \tilde{\omega}(\|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}). \quad (40)$$

Relabelling C as C' concludes our proof. \square

B.2 Proof of The Sufficient conditions for Assumption 4 In Proposition 13

Proof of Proposition 13. We fix $u^0 \in \mathcal{U}_0$ and choose $u^1, \tilde{u}^1 \in \mathcal{U}_1, \lambda \in [0, 1]$. Denote X_t^λ the state controlled by the pair $(u^0, u^\lambda) = (u^0, \lambda u^1 + (1 - \lambda)\tilde{u}^1)$ respectively. Thus,

$$\begin{aligned} &J_1(u^0, u^\lambda) - \lambda J_1(u^0, u^1) - (1 - \lambda)J_1(u^0, \tilde{u}^1) \\ &= \mathbb{E} \left[\int_0^T L_1(X_t^\lambda, u_t^0, u_t^\lambda) - \lambda L_1(X_t^1, u_t^0, u_t^1) - (1 - \lambda)L_1(X_t^0, u_t^0, \tilde{u}_t^1) dt \right] \\ &+ \mathbb{E} \left[g_1(X_T^\lambda) - \lambda g_1(X_T^1) - (1 - \lambda)g_1(X_T^0) \right] \\ &= \mathbb{E} \left[\int_0^T L_1(X_t^\lambda, u_t^0, u_t^\lambda) - \lambda L_1(X_t^1, u_t^0, u_t^1) - (1 - \lambda)L_1(X_t^0, u_t^0, \tilde{u}_t^1) dt \right] \\ &- \mathbb{E} \left[\lambda(g_1(X_T^1) - g_1(X_T^\lambda)) + (1 - \lambda)(g_1(X_T^0) - g_1(X_T^\lambda)) \right] \\ &\leq -\lambda \mathbb{E} \left[\int_0^T L_1(X_t^1, u_t^0, u_t^1) - L_1(X_t^\lambda, u_t^0, u_t^\lambda) dt + \nabla_x g_1(X_T^\lambda)^\top (X_T^1 - X_T^\lambda) \right] \\ &- (1 - \lambda) \mathbb{E} \left[\int_0^T L_1(X_t^0, u_t^0, \tilde{u}_t^1) - L_1(X_t^\lambda, u_t^0, u_t^\lambda) dt + \nabla g_1(X_T^\lambda)^\top (X_T^0 - X_T^\lambda) \right] \end{aligned}$$

where we used the the convexity of g_1 to obtain the last line. We now provide an upper bound for the last two terms. For fixed λ , by the definition of X^λ, u^λ and our assumptions on H_1 , the function $(y, z) \mapsto \nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, y, z)$ is uniformly Lipschitz continuous and $\nabla g_1(X_T^\lambda)$ is square integrable. Thus, there exists a unique solution $(Y_t^\lambda, Z_t^\lambda)$ to the BSDE

$$dY_t^\lambda = -\nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda)dt + Z_t^\lambda dW_t \quad (41)$$

$$Y_T^\lambda = \nabla g_1(X_T^\lambda) \quad (42)$$

so that by Ito's formula we have

$$\begin{aligned} & \mathbb{E} \left[\nabla g_1(X_T^\lambda)^\top (X_T^0 - X_T^\lambda) \right] = \mathbb{E} \left[Y_T^{\lambda \top} (X_T^0 - X_T^\lambda) \right] \\ & = \mathbb{E} \left[\int_0^T Y_t^{\lambda \top} (f(X_t^0, u_t^0, \tilde{u}_t^1) - f(X_t^\lambda, u_t^0, u_t^\lambda)) dt \right] \\ & \quad + \mathbb{E} \left[\text{Tr} \left(Z_t^{\lambda \top} (\sigma(X_t^0, u_t^0, \tilde{u}_t^1) - \sigma(X_t^\lambda, u_t^0, u_t^\lambda)) \right) dt \right] \\ & \quad - \mathbb{E} \left[\nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) (X_t^0 - X_t^\lambda) dt \right] \\ & = \mathbb{E} \left[\int_0^T H_1(X_t^0, u_t^0, \tilde{u}_t^1, Y_t^\lambda, Z_t^\lambda) - H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) dt \right] \\ & \quad - \mathbb{E} \left[\nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) (X_t^0 - X_t^\lambda) dt \right] \\ & \quad - \mathbb{E} \left[\int_0^T L_1(t, X_t^0, u_t^0, \tilde{u}_t^1) - L_1(t, X_t^\lambda, u_t^0, u_t^\lambda) dt \right]. \end{aligned}$$

Thus, we obtain

$$\begin{aligned} & \mathbb{E} \left[\nabla g_1(X_T^\lambda)^\top (X_T^0 - X_T^\lambda) \right] + \mathbb{E} \left[\int_0^T L_1(X_t^0, u_t^0, \tilde{u}_t^1) - L_1(X_t^\lambda, u_t^0, u_t^\lambda) dt \right] \\ & = \mathbb{E} \left[\int_0^T H_1(X_t^0, u_t^0, \tilde{u}_t^1, Y_t^\lambda, Z_t^\lambda) - H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) dt \right] \\ & \quad - \mathbb{E} \left[\nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) (X_t^0 - X_t^\lambda) dt \right] \end{aligned}$$

and

$$\begin{aligned} & \mathbb{E} \left[\nabla g_1(X_T^\lambda)^\top (X_T^1 - X_T^\lambda) \right] + \mathbb{E} \left[\int_0^T L_1(X_t^1, u_t^0, u_t^1) - L_1(X_t^\lambda, u_t^0, u_t^\lambda) dt \right] \\ & = \mathbb{E} \left[\int_0^T H_1(X_t^1, u_t^0, u_t^1, Y_t^\lambda, Z_t^\lambda) - H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) dt \right] \\ & \quad - \mathbb{E} \left[\nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) (X_t^1 - X_t^\lambda) dt \right]. \end{aligned}$$

Combining these equalities, we obtain

$$\begin{aligned} & -\lambda \mathbb{E} \left[\int_0^T L_1(X_t^1, u_t^0, u_t^1) - L_1(X_t^\lambda, u_t^0, u_t^\lambda) dt + \nabla_x g_1(X_T^\lambda)^\top (X_T^1 - X_T^\lambda) \right] \\ & - (1 - \lambda) \mathbb{E} \left[\int_0^T L_1(X_t^0, u_t^0, \tilde{u}_t^1) - L_1(X_t^\lambda, u_t^0, u_t^\lambda) dt + \nabla_x g_1(X_T^\lambda)^\top (X_T^0 - X_T^\lambda) \right] \\ & = \mathbb{E} \left[\int_0^T H_1(\lambda X_t^0 + (1 - \lambda) X_t^1, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) - \lambda H_1(X_t^1, u_t^0, u_t^1, Y_t^\lambda, Z_t^\lambda) \right. \\ & \quad \left. - (1 - \lambda) H_1(X_t^0, u_t^0, \tilde{u}_t^1, Y_t^\lambda, Z_t^\lambda) dt \right] \end{aligned}$$

$$\begin{aligned}
& - \mathbb{E} \left[\int_0^T (1 - \lambda) H_1(X_t^0, u_t^0, \tilde{u}_t^1, Y_t^\lambda, Z_t^\lambda) dt \right] \\
& + \mathbb{E} \left[\int_0^T H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) - H_1(\lambda X_t^0 + (1 - \lambda) X_t^1, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) dt \right] \\
& + \mathbb{E} \left[\int_0^T \nabla_x H_1(X_t^\lambda, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) (\lambda X_t^1 + (1 - \lambda) X_t^0 - X_t^\lambda) dt \right].
\end{aligned}$$

By the convexity of H_1 in x and strong convexity in u^1 we have that

$$\begin{aligned}
& \mathbb{E} \left[\int_0^T H_1(\lambda X_t^0 + (1 - \lambda) X_t^1, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) - \lambda H_1(X_t^1, u_t^0, u_t^\lambda, Y_t^\lambda, Z_t^\lambda) \right] \\
& - \mathbb{E} \left[\int_0^T (1 - \lambda) H_1(X_t^0, u_t^0, \tilde{u}_t^1, Y_t^\lambda, Z_t^\lambda) dt \right] \leq -\frac{\kappa\lambda(1 - \lambda)}{2} \mathbb{E} \left[\int_0^T |u_t^1 - \tilde{u}_t^1|^2 dt \right]
\end{aligned}$$

and

$$H_1(x, u^0, u^1, y, z) - H_1(\tilde{x}, u^0, u^1, y, z) + \nabla_x H_1(x, u^0, u^1, y, z)(\tilde{x} - x) \leq 0.$$

Thus,

$$J_1(u^0, u^\lambda) - \lambda J_1(u^0, u^1) - (1 - \lambda) J_1(u^0, \tilde{u}^1) \leq -\frac{\kappa\lambda(1 - \lambda)}{2} \mathbb{E} \left[\int_0^T |u_t^1 - \tilde{u}_t^1|^2 dt \right]$$

which is the strong convexity in u^1 .

Due to this strong convexity, for all u^0 , there exists a unique minimizer for the optimal response of the follower that we denote $U^*(u^0)$. To prove (23), we use the first order optimality condition for $U^{1,*}(u^0)$ which reads

$$J_1(u^0, u^1) - J_1(u^0, U^*(u^0)) \geq \frac{\kappa}{2} \mathbb{E} \left[\int_0^T |u_t^1 - U_t^*(u^0)|^2 dt \right]$$

which is (23) for $u^1 = U^*(\tilde{u}^0)$.

To conclude the proof of the Proposition it remains to prove the 1/2 Holder continuity of the best response function which allows us to verify Assumption 4. By (23), we have

$$\begin{aligned}
& \|U^*(u^0) - U^*(\tilde{u}^0)\|_{\mathcal{H}_T^2} \\
& \leq \frac{2}{\kappa^{1/2}} \left\{ |J_1(\tilde{u}^0, U^*(\tilde{u}^0)) - J_1(u^0, U^*(u^0))|^{1/2} + |J_1(u^0, U^*(\tilde{u}^0)) - J_1(\tilde{u}^0, U^*(\tilde{u}^0))|^{1/2} \right\}. \quad (43)
\end{aligned}$$

By (34) and (35), we have

$$|J_1(u^0, U^*(\tilde{u}^0)) - J_1(\tilde{u}^0, U^*(\tilde{u}^0))| \leq C \cdot \|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}, \quad (44)$$

$$|J_1(\tilde{u}^0, U^*(\tilde{u}^0)) - J_1(u^0, U^*(u^0))| < C \cdot \|u^0 - \tilde{u}^0\|_{\mathcal{H}_T^2}. \quad (45)$$

Then the desired result follows from (43), (44) and (45). \square

B.3 Additional Background on the Wiener Chaos

For any time $0 \leq t \leq T$, since the σ -algebra \mathcal{F}_t is generated by $\{W_s\}_{0 \leq s \leq t}$ then, any $u \in L^2(\mathcal{F}_t)$ admits the following *Wiener chaos expansion*

$$u = \mathbb{E}[u] + \sum_{i=0}^{\infty} \int_0^T \int_0^{s_n} \cdots \int_0^{s_1} f_i(s_n, \dots, s_1) dW_{s_1} \cdots dW_{s_n} \quad (46)$$

where, for each $i \in \mathbb{N}_+$, the deterministic functions f_i belongs to $L^2(\mathcal{S}_{i,T})$ where $\mathcal{S}_{i,T} \stackrel{\text{def.}}{=} \{(s_j)_{j=1}^i \in [0, T]^i : 0 < s_1 < \dots < s_i < T\}$.

We consider an alternative description of the Wiener chaos expansion of any random variable $u \in L^2(\mathcal{F}_t)$ for a given $0 \leq t \leq T$, which both extends more easily to multiple dimensions and does not require to lengthy computation of multiple iterated stochastic integrals. For any $i \in \mathbb{N}$, the Hermite polynomials $(h_i)_{i \in \mathbb{N}}$ are the eigenfunctions of the generator $\frac{d^2}{dx^2} - x \frac{d}{dx}$ of the Ornstein-Uhlenbeck process $dX_t = -X_t + \sqrt{2}dW_t$. For each $i \in \mathbb{N}_+$, the i^{th} Hermite polynomial h_i is given recursively by Rodrigues' formula as

$$h_i(x) = \frac{(-1)^i}{i!} e^{x^2/2} \frac{d^i}{dx^i} e^{-x^2/2} \quad h_0(x) = 1.$$

The multi-dimensional version of the i^{th} iterated integral in (46) is given by

$$\sum_{|\alpha|=i} \beta_{i,\alpha_j} \prod_{j=1}^{J_i} h_{\alpha_j} \left(\int_0^T \psi_{i,k}(s) dB_s \right) \quad (47)$$

where $\alpha \stackrel{\text{def.}}{=} (\alpha_1, \dots, \alpha_{J_i})$ is a multi-index consisting of positive integers with $i = |\alpha| \stackrel{\text{def.}}{=} \sum_{j=1}^{J_i} \alpha_j$, $\beta_{i,\alpha_j} \in \mathbb{R}$, and where $(\psi_{i,k})_{i,k \in \mathbb{Z}; 0 \leq k, \frac{k+1}{2^i} \leq t}$ is an orthonormal basis of $L^2([0, t])$. Here, we will consider the *Haar (wavelet) system* given by

$$\psi_{i,k}(s) \stackrel{\text{def.}}{=} 2^i \left(I_{[t \frac{k}{2^i}, t \frac{1+2k}{2^{i+1}})}(s) - I_{[t \frac{1+2k}{2^{i+1}}, t \frac{k+1}{2^i}]}(s) \right).$$

Since $L^2([0, t]) \subset L^2([0, T])$ we can replace the t in indicator functions with T and we can consider the *Haar (wavelet) system* on the larger space $L^2([0, T])$ where $(\psi_{i,k})_{i,k \in \mathbb{N}; 0 \leq k, \frac{k+1}{2^i} \leq 1}$ and

$$\psi_{i,k}(s) \stackrel{\text{def.}}{=} 2^i \left(I_{[T \frac{k}{2^i}, T \frac{1+2k}{2^{i+1}})}(s) - I_{[T \frac{1+2k}{2^{i+1}}, T \frac{k+1}{2^i}]}(s) \right).$$

By elementary functional analytic considerations, for each $0 \leq t \leq T$, $L^2(\mathcal{F}_t)$ is closure of the span on the of orthogonal $L^2(\mathcal{F}_t)$ random variables

$$\prod_{j=1}^j h_{\alpha_j} \left(\int_0^t \psi_{i,k}(s) dB_s \right) \quad (48)$$

where $j \in \mathbb{N}$, and $i, k \in \mathbb{N}; 0 \leq k, \frac{k+1}{2^i} \leq \frac{t}{T}$, for some $I \in \mathbb{N}$, $\beta_0, \beta_{1,\alpha_1}, \dots, \beta_{I,\alpha_{J_I}} \in \mathbb{R}$, and $(f_i)_{i \in \mathbb{N}}$ is an orthonormal basis of $L^2([0, t])$. Since $(\psi_{i,k})_{i,k \in \mathbb{N}; 0 \leq k, \frac{k+1}{2^i} \leq \frac{t}{T}}$ is piecewise constant, then the stochastic integrals in (48) simplify to

$$\int_0^t \psi_{i,k}(s) dB_s = 2^i W_{\frac{tk}{2^i}} - 2^{i+1} W_{\frac{t(1+2k)}{2^{i+1}}} + 2^i W_{\frac{t(k+1)}{2^i}}$$

where $T \frac{k+1}{2^i} \leq t$. Thus, the random variables $\{u_{i,j,k} : j \in \mathbb{N}, i, k \in \mathbb{N}, \frac{k+1}{2^i} \leq \frac{t}{T}\} \subset L^2(\mathcal{F}_t)$ where

$$u_{i,j,k}^{(t)} \stackrel{\text{def.}}{=} \prod_{j=1}^j h_{\alpha_j} \left(2^i W_{\frac{tk}{2^i}} - 2^{i+1} W_{\frac{t(1+2k)}{2^{i+1}}} + 2^i W_{\frac{t(k+1)}{2^i}} \right) \quad (49)$$

form an *orthonormal* basis of $L^2(\mathcal{F}_t)$. Conveniently, each of the $u_{i,j,k}$ can be computed without any explicit stochastic integration. We refer to [Nualart \(2006\)](#) for more details on Wiener Chaos.

Next, we will derive our universal approximation guarantees for our transformer model.

B.4 Proof of Universal Approximation Theorem 12

Our main universal approximation theorem relies on the following orthonormal basis of simple processes in \mathcal{H}_T^2 , defined by linear combinations of the elementary processes in (2).

Lemma B.6 (Orthonormal Basis of \mathcal{H}_T^2). *The collection of simple processes*
 $\mathcal{S} \stackrel{\text{def.}}{=} \{u_{i,j,k}^{s_1,s_2} : i, j, k, s_1, s_2 \in \mathbb{N}, s_2 + 1 \leq 2^{s_1}, \frac{k+1}{2^i} \leq \frac{s_2}{2^{s_1}}\}$ where

$$u_{i,j,k}^{s_1,s_2}(t, \omega) \stackrel{\text{def.}}{=} \psi_{s_1,s_2}(t) \cdot u_{i,j,k}^T(\omega)$$

is an orthonormal basis of \mathcal{H}_T^2 .

Proof of Lemma B.6. We denote the set of indices

$$\begin{aligned} \mathcal{I} &\stackrel{\text{def.}}{=} \{(i, j, k, s_1, s_2) \in \mathbb{N}^5 : \frac{k+1}{2^i} \leq \frac{s_2}{2^{s_1}} \leq 1, \frac{s_2+1}{2^{s_1}} \leq 1\}, \\ \mathcal{I}(s_1, s_2) &\stackrel{\text{def.}}{=} \{(i, j, k) \in \mathbb{N}^3 : (i, j, k, s_1, s_2) \in \mathcal{I}\}, \quad (s_1, s_2) \in \mathbb{N}^2. \end{aligned}$$

and observe the following equality

$$\mathcal{S} = \left\{ u_{i,j,k}^{s_1,s_2} \right\}_{(i,j,k,s_1,s_2) \in \mathcal{I}} \subset \mathcal{H}^2([0, T]).$$

It is easy to check that \mathcal{S} is orthonormal in $\mathcal{H}^2([0, T])$ as the set $\{\psi_{s_1,s_2}\}_{(s_1,s_2) \in \mathbb{N}^2}$ is orthonormal in $L^2([0, T])$, and for all $(s_1, s_2) \in \mathbb{N}^2$, $\left(u_{i,j,k}^T\right)_{(i,j,k) \in \mathcal{I}(s_1,s_2)}$ is orthonormal in $L^2(\mathcal{F}_T)$. It remains to show that $\mathcal{H}^2([0, T])$ is the closure of the span of \mathcal{S} .

We denote $\hat{\mathcal{H}}_0$ to the set of simple processes $Z \in \mathcal{H}^2([0, T])$ satisfying

$$Z \stackrel{\text{def.}}{=} \sum_{l=1}^m \xi_l \mathbb{1}_{[t_{2l-1}, t_{2l}]}, \quad (50)$$

where $m \in \mathbb{N}$, $(t_i)_{i=1}^{m+1}$ is a strictly increasing sequence of real numbers satisfying $t_1 = 0, t_{2m} < T$, and $\xi_l \in L^2(\mathcal{F}_{t_{2l-1}})$, for $1 \leq l \leq m$.

We notice that $\tilde{\mathcal{H}}_0 = \mathcal{H}^2([0, T])$. Indeed, we observe that for every simple process $Z := \sum_{l=1}^m \xi_l \mathbb{1}_{[t_l, t_{l+1}]}$, $t_1 = 0, t_{m+1} = T$ can be written as the $\mathcal{H}^2([0, T])$ -limit of the sequence $(Z^n)_{n \in \mathbb{N}} \subset \hat{\mathcal{H}}_0$:

$$Z_t^n = \sum_{l=1}^m \xi_l \mathbb{1}_{[t_l, t_{l+1} - \frac{1}{n}]}. \quad (51)$$

The previous shows that $\tilde{\mathcal{H}}_0 = \hat{\mathcal{H}}_0 = \mathcal{H}^2([0, T])$.

Next, assume that $Z \in \hat{\mathcal{H}}_0$ is such that: for all $U \in \mathcal{S}$

$$\int_0^T \mathbb{E} [U_t Z_t] dt = 0. \quad (52)$$

We notice that for any $\tilde{s}_1 \in \mathbb{N}$ sufficiently large there exists $\tilde{s}_2 \in \mathbb{N}$ satisfying

$$\begin{aligned} 1 + \tilde{s}_2 &\leq 2^{\tilde{s}_1}, \\ t_1 &< \frac{\tilde{s}_2 T}{2^{\tilde{s}_1}} \leq t_2 < \frac{(1 + 2\tilde{s}_2)T}{2^{\tilde{s}_1+1}} < \frac{(1 + \tilde{s}_2)T}{2^{\tilde{s}_1}} < t_3. \end{aligned} \quad (53)$$

Then, for any $(i, j, k) \in \mathcal{I}(\tilde{s}_1, \tilde{s}_2)$, we set

$$U \stackrel{\text{def.}}{=} u_{i,j,k}^{\tilde{s}_1, \tilde{s}_2} \in \mathcal{S}. \quad (54)$$

Plugging in the process (54) into (52), we obtain

$$\int_0^T \mathbb{E} [U_t Z_t] dt = 2^{\tilde{s}_1} \left(t_2 - \frac{\tilde{s}_2 T}{2^{\tilde{s}_1}} \right) \mathbb{E} [u_{i,j,k}^T \xi_1] = 0, \quad \forall (i, j, k) \in \mathcal{I}(\tilde{s}_1, \tilde{s}_2). \quad (55)$$

Finally, we will prove that the closure of the span of $\left(u_{i,j,k}^T \right)_{(i,j,k) \in \mathcal{I}(\tilde{s}_1, \tilde{s}_2)}$ contains $L^2(\mathcal{F}_{t_1})$. Using (53), we observe that for all $(i, k) \in \mathbb{N}^2$, $1 + k \leq 2^i$ the following inequality holds:

$$0 \leq \frac{t_1(k+1)}{T2^i} < \frac{\tilde{s}_2}{2^{\tilde{s}_1}} \leq 1.$$

Using that the set of the dyadic rationals is dense in $[0, 1]$, there exists a sequence $(i_n, k_n) \in \mathbb{N}^2$, $k_n + 1 \leq 2^{i_n}$, satisfying:

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{(k_n + 1)}{2^{i_n}} &= \frac{t_1(k+1)}{T2^i}, \\ \frac{t_1(k+1)}{T2^i} &\leq \frac{(k_n + 1)}{2^{i_n}} \leq \frac{\tilde{s}_2}{2^{\tilde{s}_1}}, \quad n \in \mathbb{N}. \end{aligned}$$

Hence, for all $j \in \mathbb{N}$, $\left(u_{i_n, j, k_n}^T \right)_{n \in \mathbb{N}} \subset \mathcal{S}$. By continuity of the paths of the Brownian motion and the Hermite polynomials, we obtain that $\left(u_{i_n, j, k_n}^T \right)_{n \in \mathbb{N}}$ converges \mathbb{P} -a.s. to $u_{i,j,k}^{(t_1)}$. Moreover, we observe that $\left(\xi_1 u_{i_n, j, k_n}^T \right)_{n \in \mathbb{N}}$ is uniformly integrable. Indeed, for $0 < \epsilon_0$ small enough, applying Hölder's inequality, there exist constants $p_1 := \frac{3}{2(1+\epsilon_0)}$, and $p_2 := \frac{1}{1-\frac{1}{p_1}}$ such that

$$\begin{aligned} &\sup_{n \in \mathbb{N}} \mathbb{E} \left[\left| \xi_1 u_{i_n, j, k_n}^T \right|^{1+\epsilon_0} \right] \\ &\leq \mathbb{E} \left[|\xi_1|^{(1+\epsilon_0)p_1} \right] \sup_{n \in \mathbb{N}} \mathbb{E} \left[\left| \prod_{\tilde{j}=1}^j h_{\tilde{j}} \left(2^{i_n} W_{\frac{T k_n}{2^{i_n}}} - 2^{i_n+1} W_{\frac{T(1+2k_n)}{2^{i_n+1}}} + 2^{i_n} W_{\frac{T(k_n+1)}{2^{i_n}}} \right) \right|^{p_2} \right] \\ &< \infty. \end{aligned}$$

We therefore deduce that

$$\mathbb{E} \left[\xi_1 u_{i,j,k}^{(t_1)} \right] = \lim_{n \rightarrow \infty} \mathbb{E} \left[\xi_1 u_{i_n, j, k_n}^T \right] = 0. \quad (56)$$

Using completeness of the basis $\{u_{i,j,k}^{(t_1)}\}_{(i,j,k) \in \mathbb{N}^3}$ in $L^2(\mathcal{F}_{t_1})$ we obtain that

$$\xi_1 = 0, \quad \mathbb{P}\text{-a.s.}$$

Finally, we repeat the same argument in (53) for every addend in (50), obtaining

$$Z = 0, \quad dt \otimes \mathbb{P} - a.e.$$

Hence,

$$\overline{\text{span}(\mathcal{S})}^\perp \cap \hat{\mathcal{H}}_0 = \{0\}, \quad (57)$$

where \mathcal{H}_0 denotes the set of simple processes in $\mathcal{H}^2([0, T])$. Finally, using the decomposition

$$\mathcal{H}^2([0, T]) = \overline{\text{span}(\mathcal{S})} \oplus \overline{\text{span}(\mathcal{S})}^\perp,$$

and (57), we have that $\hat{\mathcal{H}}_0 \subset \overline{\text{span}(\mathcal{S})}$. The later implies that $\overline{\text{span}(\mathcal{S})} = \mathcal{H}^2([0, T])$. \square

Generally, deep learning faces the curse of dimensionality in finite-dimensions; see e.g. [Lanthaler and Stuart \(2023\)](#). Nevertheless, the impact of infinite-dimensionality on the parametric complexity of deep learning models can be reduced by considering the following *trainable* version of the “super-expressive” activation function of [Zhang et al. \(2008\)](#) designed to exploit the bit-extraction mechanism of [Bartlett et al. \(2019\)](#). The next lemma provides quantitative rates for attentional neural operators with the activation function (10); since, in that case, they are not overwhelmingly large (as is the case approximation of general Lipschitz non-linear operators).

Lemma B.7 (Approximation Of Lipschitz Operators with By Attentional Neural Operator). *Let $\mathcal{K}_0 \subseteq \mathcal{H}_T^2$ be compact, $F : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ be an L -Lipschitz (non-linear) operator, and consider respective “dimension reduction” and “approximation” errors $\varepsilon_D, \varepsilon_A > 0$. There exists an attentional neural operator $\hat{F} : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ satisfying*

$$\sup_{u \in \mathcal{K}_0} \|F(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} \leq \varepsilon_D + \varepsilon_A.$$

Furthermore, the complexity of the neural operator \hat{F} is recorded in [Table 2](#).

We provide explicit quantitative parameter estimates in the special cases where the compact set \mathcal{K}_0 and the target neural operator are compatible. We consider the following notion of a small compact subset of a Banach space.

Definition B.8 ((r, f) -Exponentially Ellipsoidal). *Let $r > 0$, $f : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$, and fix an orthonormal basis $\{u_i\}_{i=0}^\infty$ of \mathcal{H}_T^2 . A subset $\mathcal{K} \subseteq \mathcal{H}_T^2$ is of (r, f) -exponential width if the following holds for each $u \in \mathcal{K}$:*

- (i) $u = \sum_{i=1}^\infty \beta_i u_i$ and $|\beta_i| \lesssim e^{-ri}$,
- (ii) $f(u) = \sum_{i=1}^\infty c_i u_i$ and $|c_i| \lesssim e^{-ri}$.

Exponentially ellipsoidal compact sets can be efficiently approximated by low-dimensional representations arising from projections onto the relevant basis. However, they may still be large in metric entropy (i.e. they may be difficult to cover by a few small metric balls). This is not the case if there is something akin to a latent “low-dimensional submanifold” on which the data/approximation is focused. The following definition makes this rigorous for our infinite-dimensional setting.

Definition B.9 ((r, f, d) -Exponential Manifold). *Let $r > 0$, $f : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$, fix an orthonormal basis $\{u_i\}_{i=0}^\infty$ of \mathcal{H}_T^2 , and let $\tilde{\mathcal{K}}$ be an (r, f) -exponentially ellipsoidal subset of \mathcal{H}_T^2 . A compact subset $\mathcal{K} \subseteq \tilde{\mathcal{K}}$ is said to be an (r, f, d) -Exponential Manifold if there exists a 1-Lipschitz “latent parameterization” map $\pi : \mathbb{R}^d \rightarrow \mathcal{H}_T^2$ and $\mathcal{K} = \pi(\{z \in \mathbb{R}^d : \|z\| \leq 1\})$.*

Notably, the map π need not be known nor be injective (as in [Kratsios et al. \(2024\)](#)), nor does it need to be inverted by the deep learning model (either explicitly or implicitly during the approximation theorem). Instead, it simply encodes (in a possibly non-linear way) a low-dimensional structure into the compact set of controls, allowing for an efficient approximation by controlling the entropy number, see e.g. [Carl \(1997\)](#); [Lorentz \(1966\)](#); [Petrova and Wojtaszczyk \(2023\)](#), of the compact set on which the approximation is performed number.

Param.	Example 7	Example 6
No. Param	$\mathcal{O}\left(\varepsilon_D^{-3\ln(C)/r} \ln(\varepsilon_D^{-1/r})^4\right)$	Finite
Depth	$\mathcal{O}\left(\ln(\varepsilon_D^{-1/r}) \varepsilon_D^{-\ln(C)/r}\right)$	Finite
Width	$\mathcal{O}(\varepsilon_D^{-\ln(C)/r})$	$\mathcal{O}(d^{d+1} N^{2d+2} \varepsilon_D^{-3d-3})$
Decoding Dim. (Q)	$\mathcal{O}(\ln(\varepsilon_D^{-1/r}))$	Finite
Encoding Dim. (d)	$\mathcal{O}(\varepsilon_D^{1/(1-r)})$	Finite
Att. Complexity (N)	$(\tilde{c} \varepsilon_A^{-1} \ln(\varepsilon_D^{-1/r})^{1/2})^{c(\ln(\varepsilon_D^{-1/r}))}$	Finite

Table 2: **Complexity of the neural operator.** *Case 1:* \mathcal{K} is an (r, f, d) -exponential manifold in controls in \mathcal{H}_T^2 and σ is the super-expressive activation function with neuron-specific skip-connections in (7); $c, \tilde{c} > 0$ are absolute constants.

Case 2: \mathcal{K} is an arbitrary compact subset of controls in \mathcal{H}_T^2 and σ is the standard non-trainable activation functions of Example (6).

The error $\varepsilon_D > 0$ in Table 2 expresses the ‘‘dimension reduction’’ error resulting from the encoding (\mathcal{E}) and decoding (\mathcal{D}) maps used in the definition of our attentional neural operator, in Definition 6. That is, ε_D expresses the error made in encoding infinite dimensional objects, namely \mathbb{F} -adapted processes, into finite-dimensional objects, namely vectors in some Euclidean spaces. Once the neural operator has implicitly transformed the approximation problem as an approximation problem between finite-dimensional spaces, it is approximated by the MLP (f) between the encoding and decoding layers of the attentional neural operator. The error $\varepsilon_A > 0$ in Table 2 expresses the error incurred in this finite-dimensional approximation step.

Proof of Lemma B.7. Fix respective ‘‘dimension reduction’’ and ‘‘approximation’’ errors $\varepsilon_D, \bar{\varepsilon}_A > 0$.

Step 1 - Finite Dimensional Encoding

Enumerate $\mathcal{S} = \{s_i\}_{i=1}^\infty$, where \mathcal{S} is as in Lemma B.6. For any $d \in \mathbb{N}$ (which we fix retroactively) define the 1-Lipschitz encoder $\mathcal{E}_d : \mathcal{H}_T^2 \rightarrow \mathbb{R}^d$ given, for each $u \in \mathcal{H}_T^2$, by

$$\mathcal{E}_d(u) \stackrel{\text{def.}}{=} (\langle u, s_j \rangle_{\mathcal{H}_T^2})_{j=1}^d.$$

Consider its right-inverse $\iota_d : \mathbb{R}^d \rightarrow \mathcal{H}_T^2$ given, for each $x \in \mathbb{R}^d$, by

$$\iota_d(x) \stackrel{\text{def.}}{=} \sum_{i=1}^d x_i s_i.$$

Observe that ι_d is an isometric embedding; we will come back to this point shortly.

- **\mathcal{K}_0 is the exponentially ellipsoidal \mathcal{K} in Definition B.8:** By the exponential decay condition in Definition B.8, we have that: for each $u \in \mathcal{K}$, with representation $u = \sum_{i=1}^\infty \beta_i s_i$ the following error estimate holds by orthonormality of the $(s_i)_{i=1}^\infty$

$$\begin{aligned}
\|u - \iota_d \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2}^2 &= \left\| \sum_{i=1}^\infty \beta_i s_i - \iota_d \circ \mathcal{E}_d(u) \right\|_{\mathcal{H}_T^2}^2 \\
&= \sum_{i=1}^\infty |\beta_i|^2 I_{i \leq d} \|s_i\|_{\mathcal{H}_T^2}^2 \\
&= \sum_{i=d+1}^\infty |\beta_i|^2 \\
&\leq \frac{C e^{-rd}}{1 - e^{-r}} \stackrel{\text{def.}}{=} \tilde{C}_K e^{-rd}
\end{aligned} \tag{58}$$

where $\tilde{C}_K \stackrel{\text{def.}}{=} C/(1 - e^{-r})$. Fix $\varepsilon_0 > 0$. We now can retroactively set $d \stackrel{\text{def.}}{=} \ln\left(\frac{2^r \tilde{C}_K^r}{(1 - e^{-r})^r} \varepsilon_0^{-r}\right) = \ln(C_K \varepsilon_0^{-r}) \in \mathcal{O}(\ln(\varepsilon_0^{-1/r}))$ where $C_K \stackrel{\text{def.}}{=} \left(\frac{C}{3L(1 - e^{-r})}\right)^{1/r} > 0$.

- **Exponential Sub-manifold:** If \mathcal{K}_0 satisfies Definition B.9, then this case is implied by the previous case (i.e. that of exponentially ellipsoidal compacta),
- **General \mathcal{K}_0 :** If \mathcal{K}_0 is general, then by the 1-bounded approximation property (e.g. see Szarek (1987)) of Hilbert spaces with orthonormal bases (which are simply Banach spaces with Schauder bases), for every $\varepsilon_0 > 0$ there is some $d \in \mathbb{N}$ for which $\sup_{x \in \mathcal{K}_0} \|x - \mathcal{E}_d(x)\| \leq \varepsilon_0$.

In each case, we have that

$$\sup_{u \in \mathcal{K}_0} \|u - \iota_d \circ \mathcal{E}_d(u)\| \leq \varepsilon_0. \quad (59)$$

Let $L \geq 0$ denote the optimal Lipschitz constant of F . We note that the map

$$f^{(1)} \stackrel{\text{def.}}{=} F \circ \iota_d : \mathbb{R}^d \rightarrow \mathcal{H}_T^2$$

is (L, α) -Hölder since F is (L, α) -Hölder and since is 1-Lipschitz. In particular, the Lipschitz constant of $f^{(1)}$ is independent of d .

Step 2 - Quantization of The Image and The Domain:

Since \mathcal{K}_0 is compact and since \mathcal{E}_d is continuous, then $\mathcal{E}_d(\mathcal{K}_0)$ is compact and thus $\mathcal{E}_d(\mathcal{K}_0)$ is totally bounded. Therefore, for every $\varepsilon_1 > 0$ there exists a finite subset $\{x_n\}_{n=1}^{N_{\varepsilon_1}} \subseteq \mathcal{E}_d(\mathcal{K}_0)$, of minimal cardinality $N \stackrel{\text{def.}}{=} N_{\varepsilon_1}$ (a so-called minimal ε_1 -net) such that:

$$\max_{x \in \mathcal{K}_0} \min_{n=1, \dots, N_{\varepsilon_1}} \|x - x_n\|_2 < \left(\frac{1}{L} \frac{\varepsilon_1}{3}\right)^{1/\alpha}. \quad (60)$$

Since, $f^{(1)}$ is an L -Lipschitz surjection of \mathcal{K}_0 onto $f(\mathcal{K}_0)$ then (60) implies that

$$\begin{aligned} \max_{x \in \mathcal{K}_0} \min_{n=1, \dots, N_{\varepsilon_1}} \|f^{(1)}(x) - f^{(1)}(x_n)\|_{\mathcal{H}_T^2} &\leq L \max_{x \in \mathcal{K}_0} \min_{n=1, \dots, N_{\varepsilon_1}} \|x - x_n\|_2^\alpha \\ &< L \left(\left(\frac{\varepsilon_1}{3L}\right)^{1/\alpha}\right)^\alpha = \frac{\varepsilon_1}{3}. \end{aligned} \quad (61)$$

By Lemma B.6, $(s_i)_{i=1}^\infty$ is an orthonormal basis of the Hilbert space \mathcal{H}_T^2 . Therefore, it realizes the 1-bounded approximation property. This means that since $F(\mathcal{K}_0)$ is compact then, for each $Q \in \mathbb{N}$

$$\max_{y \in F(\mathcal{K}_0)} \|y - P_Q(y)\|_{\mathcal{H}_T^2} \xrightarrow{Q \rightarrow \infty} 0 \text{ and } \|P_Q\|_{op} \leq 1 \quad (62)$$

where $P_Q : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ is the (rank Q) orthogonal projection operator of \mathcal{H}_T^2 onto $\text{span}(\{s_i\}_{i=1}^Q)$; that is, $P_Q(u) \mapsto \sum_{i=1}^Q \langle s_i, u \rangle_{\mathcal{H}_T^2} s_i$; where $\|\cdot\|_{op}$ denotes the operator norm. Thus, for each $\varepsilon_2 > 0$ (to be fixed retroactively) there exists some $Q \stackrel{\text{def.}}{=} Q_{\varepsilon_2} \in \mathbb{N}$ for which

$$\max_{y \in F(\mathcal{K}_0)} \|y - P_Q(y)\|_{\mathcal{H}_T^2} \leq \varepsilon_2.$$

In the special case where $F(\mathcal{K}_0)$ satisfies Definition (B.8) then, by a similar computation to Step 1 (58), we obtain the following bounds on $Q \stackrel{\text{def.}}{=} Q_{\varepsilon_2}$:

- **Exponentially Ellipsoidal \mathcal{K}_0 :** $Q \in \mathcal{O}(\ln(\varepsilon_D^{-1/r}))$,
- **Exponential Sub-manifold:** As before, if \mathcal{K}_0 satisfies Definition B.9, then this case is implied by the previous case,

- **General \mathcal{K}_0 :** $Q \rightarrow \infty$ as $\varepsilon_2 \rightarrow 0$.

To summarize this step, the set $\{x_n\}_{n=1}^N$ discretized $\mathcal{E}_d(\mathcal{K}_0)$ and the set $\{y_n\}_{n=1}^N$ discretized the “finitely parameterized” image of \mathcal{K}_0 under F .

Step 3 - Simplicialization of Target Function:

Our next objective is to replace the target function F , with a function which maps between $\mathcal{E}_d(\mathcal{K}_0)$ to an N -simplex and which, informally speaking, is an approximate continuous selection to the nearest neighbour problem

$$x \mapsto \operatorname{argmin}_{n=1, \dots, N} \|F(x) - y_n\|.$$

First, we construct an “projection-like/extremal” version of this solution, as in Step 4 of the proof of (Acciaio et al., 2023, Theorem 3.8). In the second step, we “mollify” that function to make it comparable with the softmax operation.

Fix $\varepsilon_3 > 0$. Let $\mathcal{P}_1(\{x_n\}_{n=1}^N, \mathcal{W}_1)$ denote the 1-Wasserstein space over $\{x_n\}_{n=1}^N$ with respect to the α -snowflaked of the Euclidean distance $\|\cdot\|_2^\alpha$ on the inherited finite set $\{x_n\}_{n=1}^N$; i.e. the metric $\mathbb{R}^d \times \mathbb{R}^d \ni (x, \tilde{x}) \mapsto \|x - \tilde{x}\|_2^\alpha$ restricted to the set $\{x_n\}_{n=1}^N$. By (Bruè et al., 2021, Theorem 3.2), there exists a Lipschitz map (a so-called weak random projection) $\Pi : (\mathcal{E}_d(\mathcal{K}_0), \|\cdot\|_2) \rightarrow \mathcal{P}_1(\{x_n\}_{n=1}^N, \mathcal{W}_1)$ with the property that: for each $x \in \mathcal{E}_d(\mathcal{K}_0)$ if $x \in \{x_n\}_{n=1}^N$ then $\Pi(x) = \delta_x$; where δ_x is the pointmass on x . Furthermore, the Lipschitz constant L_Π of Π is at-most $c \log_2(C_{(\{x_n\}_{n=1}^N, \|\cdot\|_2^\alpha)})$ where $c > 0$ is an absolute constant and $C_{(\{x_n\}_{n=1}^N, \|\cdot\|_2^\alpha)} > 0$ is the doubling constant of the set $\{x_n\}_{n=1}^N$ with respect to the α -snowflake of the Euclidean distance restricted to $\{x_n\}_{n=1}^N$. By (Robinson, 2011, Lemma 9.3), since the inclusion of $\{x_n\}_{n=1}^N$ into \mathbb{R}^d is an isometric embedding (with respect to the Euclidean distance on \mathbb{R}^d) then the doubling constant of $(\{x_n\}_{n=1}^N, \|\cdot\|_2)$ is no larger than that of $(\mathbb{R}^d, \|\cdot\|_2)$. By (Robinson, 2011, Lemma 9.3), the doubling constant of \mathbb{R}^d in the Euclidean metric is 2^{d+1} . Thus, the first statement in (Acciaio et al., 2023, Lemma 7.1) implies that doubling of $(\{x_n\}_{n=1}^N, \|\cdot\|_2^\alpha)$ is at-most equal to the doubling constant of \mathbb{R}^d to the power of $\lceil \frac{1}{\alpha} \rceil$. Hence, the Lipschitz constant L_Π of Π is

$$L_\Pi \leq c \log_2(C_{(\{x_n\}_{n=1}^N, \|\cdot\|_2^\alpha)}) \leq c \left\lceil \frac{1}{\alpha} \right\rceil \log_2(C_{(\{x_n\}_{n=1}^N, \|\cdot\|_2)}) \leq c \left\lceil \frac{1}{\alpha} \right\rceil (d+1) \leq \tilde{c} \left\lceil \frac{1}{\alpha} \right\rceil d \stackrel{\text{def.}}{=} C_\Pi \quad (63)$$

where $\tilde{c} \stackrel{\text{def.}}{=} 2 \max\{1, c\} > 0$.

As shown in (Acciaio et al., 2023, Equations (41)-(46)), the map $\iota_N : \mathcal{P}(\{x_n\}_{n=1}^N, \mathcal{W}_1) \ni \mu = \sum_{n=1}^N w_n \delta_{x_n} \rightarrow (w_n)_{n=1}^N \in (\Delta_N, \|\cdot\|_2)$ is $\frac{2}{\varepsilon_1}$ -Lipschitz since the minimal distance between any distinct pairs of points in $\{x_n\}_{n=1}^N$ is ε_1 . Together with the right-hand side of (63), this shows that the α -Hölder constant $L_{f^{(2)}}$ of composite map $f^{(2)} : \iota_N \circ \Pi_N : \mathcal{E}_d(\mathcal{K}_0) \rightarrow \Delta_N$ is bounded-above by

$$L_{f^{(2)}} \leq L_\Pi \frac{2}{\varepsilon_1} \leq \tilde{c} \left\lceil \frac{1}{\alpha} \right\rceil N \frac{2}{\varepsilon_1} \stackrel{\text{def.}}{=} \tilde{L}_{f^{(2)}}. \quad (64)$$

For $n = 1, \dots, N$, let $y_n \stackrel{\text{def.}}{=} P_{Q_{\varepsilon_3}}(f^{(1)}(x_n)) = \sum_{i=1}^{Q_{\varepsilon_3}} \langle f^{(1)}(x_n), s_i \rangle_{\mathcal{H}_T^2} s_i$. Define the Lipschitz map $\eta_N : \Delta_N \rightarrow \mathcal{H}_T^2$ by

$$\eta_N : w \mapsto \sum_{i=1}^N w_i y_n. \quad (65)$$

Note that $\eta_N : (\Delta_N, \|\cdot\|_2) \rightarrow (\mathcal{H}_T^2, \|\cdot\|_{\mathcal{H}_T^2})$ is $L_{\varepsilon_1, \varepsilon_2}^\eta$ -Lipschitz with optimal Lipschitz constant, which we denote by $\operatorname{Lip}(\eta_N) \geq 0$, bounded-above by the constant $L_{\varepsilon_1, \varepsilon_2}^\eta > 0$ defined by

$$|\eta_N(w) - \eta_N(v)| \leq \sum_{i=1}^N \|y_i\|_{\mathcal{H}_T^2} |w_i - v_i| \quad (66)$$

$$\begin{aligned}
&\leq (\text{diam}(f(\mathcal{E}_d(\mathcal{K}_0))) + 2\epsilon_2) \sum_{i=1}^N |w_i - v_i| \leq (L \text{diam}(\mathcal{K}_0)^\alpha + 2\epsilon_2) \sum_{i=1}^N |w_i - v_i| \\
&\leq (L \text{diam}(\mathcal{K}_0)^\alpha + 2\epsilon_2) \sqrt{N} \|w - v\|_2 \stackrel{\text{def.}}{=} L_{\epsilon_1, \epsilon_2}^\eta \|w - v\|_2
\end{aligned} \tag{67}$$

where $w, v \in \Delta_N$ are arbitrary and we have used the inequality $\|\cdot\|_1 \leq \sqrt{N} \|\cdot\|_2$ (on \mathbb{R}^N).

Next, we show that $\eta_N \circ f^{(2)}$ approximates F on $\mathcal{E}_d(\mathcal{K}_0)$.

Since \mathcal{H}_T^2 is a QAS space (see (Acciaio et al., 2023, Definition 3.4) with $p = 1$ and $C_\eta = 1$, as shown in (Acciaio et al., 2023, Example 5.1)) then Step 4 of the proof of (Acciaio et al., 2023, Theorem 3.8) holds unaltered in our setting (with $\mathcal{X} = \mathcal{E}_d(\mathcal{K}_0)$, $(\mathcal{Y}, d_Y = (\mathcal{H}_T^2, \|\cdot\|_{\mathcal{H}_T^2})$, the α -Hölder target function with respect to the α -Hölder seminorm $L > 0$ of $f^{(1)}$).

Set³ $\epsilon_1 \stackrel{\text{def.}}{=} \frac{\epsilon_D}{3 \cdot 3L C_\Pi} = \frac{\epsilon_D}{9L\bar{c}d} > 0$ and $\epsilon_2 \stackrel{\text{def.}}{=} \frac{\epsilon_D}{9} > 0$. Arguing identically to the (Acciaio et al., 2023, Proof of Theorem 3.8 - Step 3), following (Acciaio et al., 2023, Equation 51), we conclude that

$$\sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(1)}(x) - \eta_N \circ f^{(2)}(x)\|_{\mathcal{H}_T^2} \leq \frac{\epsilon_D}{3}. \tag{68}$$

where, for us, our improved estimate on C_Π was given in (63); as is summarized in (68).

We now modify the map $f^{(2)}$ so that it takes values in the range of the softmax function; that is, in the relative interior of the N -simplex, i.e., in the set $\text{int}(\Delta_N) \stackrel{\text{def.}}{=} \{w \in (0, 1)^N : \sum_{n=1}^N w_n = 1\}$. We subsequently associate $f^{(2)}$ to a map taking values \mathbb{R}^{N-1} . Finally, this latter map will be approximated by a neural network in step 3.

Fix $0 < \epsilon_3 \leq 1$, to be determined retroactively. Consider the 1-Lipschitz homotopy $H : [0, 1] \times \Delta_N \rightarrow \Delta_N$ given by $H(t, w) \stackrel{\text{def.}}{=} t(w - \bar{\Delta}_N) + \bar{\Delta}_N$, where⁴ $\bar{\Delta}_N = (1/N, \dots, 1/N) \in \Delta_N$ is the barycenter of the N -simplex. Observe that, for each $t \in [0, 1]$ we have $H(t, \Delta_N) \subset \text{int}(\Delta_N)$ and for each $0 \leq \epsilon_3 \leq \max_{w \in \Delta_N} \|w - \bar{\Delta}_N\|_2 1 - \frac{1}{N}$ there exists a $t_{\epsilon_3} \in [0, 1]$ satisfying⁵

$$\max_{w \in \Delta_N} \|H(t_{\epsilon_3}, w) - w\|_2 \leq \max_{w \in \Delta_N} \|H(t_{\epsilon_3}, w) - w\|_1 \leq \epsilon_3. \tag{69}$$

The right-hand inequality can be solved explicitly for the largest value of t_{ϵ_3} in $[0, 1]$; this is because $w^* \in \Delta_N$ given by $w_1^* = 1$ and $w_j^* = 0$ for $j = 2, \dots, N$ is a non-unique maximizer of $\max_{w \in \Delta_N} \|H(t_{\epsilon_3}, w) - w\|_1$. We therefore compute

$$\begin{aligned}
\|H(t_{\epsilon_3}, w) - w\|_1 &= \underbrace{|(t_{\epsilon_3}(1 - 1/N) + 1/N) - 1|}_{1^{st} \text{ component}} \\
&\quad + (N-1) \underbrace{|(t_{\epsilon_3}(0 - 1/N) + 1/N) - 0|}_{j > 1^{st} \text{ components}}
\end{aligned} \tag{70}$$

$$\begin{aligned}
&= |1 - t_{\epsilon_3}| |(1 - 1/N)| + (N-1) |1 - t_{\epsilon_3}| |1/N| \\
&= |1 - t_{\epsilon_3}| (N-1)/N + |1 - t_{\epsilon_3}| (N-1)/N \\
&= (1 - t_{\epsilon_3}) \frac{2(N-1)}{N}.
\end{aligned} \tag{71}$$

If $\epsilon_3 > 0$ is small enough,⁶ then setting the right-hand side of (71) equal to ϵ_3 and solving for t_{ϵ_3} yields $t_{\epsilon_3} = 1 - \frac{N\epsilon_3}{2(N-1)}$. For general values of ϵ_3 , we may set

$$t_{\epsilon_3} = 1 - \min \left\{ \frac{N\epsilon_3}{2(N-1)}, 1/2 \right\}. \tag{72}$$

3. In the notation of (Acciaio et al., 2023, Proof of Theorem 3.8 - Step 3), we have set $\bar{\epsilon}_A \stackrel{\text{def.}}{=} \epsilon_Q \stackrel{\text{def.}}{=} \epsilon_D/3$.

4. I.e. $\bar{\Delta}_N \stackrel{\text{def.}}{=} (1/N, \dots, 1/N)$ is the barycenter of the N -simplex Δ_N .

5. For the interested reader: we have just noted that the boundary of Δ_N is a \mathcal{Z} -set, in the sense of (van Mill, 2001, Section 5.1).

6. Namely, one needs that $0 < \epsilon_3 < \frac{2(N-1)}{N}$.

Consider the map

$$\rho \stackrel{\text{def.}}{=} \text{softmax}_N \circ W : \mathbb{R}^{N-1} \rightarrow \text{int}(\Delta_N) \quad (73)$$

where $W : \mathbb{R}^{N-1} \ni x \rightarrow (x_1, \dots, x_{N-1}, 1) \in \mathbb{R}^N$. A right-inverse of the smooth function $R : \text{int}(\Delta_N) \rightarrow \mathbb{R}^{N-1}$ given for each $y \in \text{int}(\Delta_N)$ by

$$R(y) \stackrel{\text{def.}}{=} ((\ln(y_i) - \ln(y_N) + 1))_{i=1}^N. \quad (74)$$

Finally, we define the ‘‘mollified simplicial target function’’ $f^{(3)} \stackrel{\text{def.}}{=} R \circ H(t_{\varepsilon_3}, \cdot) \circ f^{(2)} : \mathbb{R}^d \rightarrow \mathbb{R}^{N-1}$.

We therefore, have the following uniform estimate between $\rho \circ f^{(3)}$ and $f^{(2)}$

$$\max_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(2)}(x) - \rho \circ f^{(3)}(x)\|_{\mathcal{H}_T^2} = \max_{x \in \mathcal{E}_d(\mathcal{K})} \|f^{(2)}(x) - \rho \circ (R \circ H(t_{\varepsilon_3}, \cdot) \circ f^{(2)}(x))\|_{\mathcal{H}_T^2} \quad (75)$$

$$\begin{aligned} &= \max_{x \in \mathcal{E}_d(\mathcal{K})} \|f^{(2)}(x) - H(t_{\varepsilon_3}, \cdot) \circ f^{(2)}(x)\|_{\mathcal{H}_T^2} \\ &= \max_{w \in \Delta_N} \|w - H(t_{\varepsilon_3}, w)\|_{\mathcal{H}_T^2} \end{aligned} \quad (76)$$

$$\leq \varepsilon_3, \quad (77)$$

where we have use the fact that $f^{(2)}$ takes values in the N -simplex in deducing (76), and (77) held by (69). Combining our estimate in (68) with those in (75)-(77) yields

$$\begin{aligned} \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(1)}(x) - \eta_N \circ \rho \circ f^{(3)}(x)\|_{\mathcal{H}_T^2} &\leq \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(1)}(x) - \eta_N \circ f^{(2)}(x)\|_{\mathcal{H}_T^2} \\ &\quad + \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|\eta_N \circ f^{(2)}(x) - \eta_N \circ \rho \circ f^{(3)}(x)\|_{\mathcal{H}_T^2} \\ &\leq \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(1)}(x) - \eta_N \circ f^{(2)}(x)\|_{\mathcal{H}_T^2} \\ &\quad + \text{Lip}(\eta_N) \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(2)}(x) - \rho \circ f^{(3)}(x)\|_{\mathcal{H}_T^2} \\ &\leq \frac{\varepsilon_D}{3} + \text{Lip}(\eta_N) \varepsilon_3 \\ &\leq \frac{\varepsilon_D}{3} + ((L \text{diam}(\mathcal{K}_0)^\alpha + 2\varepsilon_2) \sqrt{N}) \varepsilon_3, \end{aligned} \quad (78)$$

where $\text{Lip}(\eta_N)$ denotes the optimal Lipschitz constant of the map η_N which we have bounded above by $L_{\varepsilon_1, \varepsilon_2}^\eta$, in (67). Combining the estimates in (61) and in (59) with those in (78)-(79) yields

$$\begin{aligned} \sup_{u \in \mathcal{K}_0} \|F(u) - \eta_N \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} &\leq \sup_{u \in \mathcal{K}_0} \|F(u) - F \circ \iota_d \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \\ &\quad + \sup_{u \in \mathcal{K}_0} \|F \circ \iota_d \circ \mathcal{E}_d(u) - \eta_N \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \\ &\leq L \sup_{u \in \mathcal{K}_0} \|u - \iota_d \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \\ &\quad + \sup_{u \in \mathcal{K}_0} \|F \circ \iota_d \circ \mathcal{E}_d(u) - \eta_N \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \\ &= L \sup_{u \in \mathcal{K}_0} \|u - \iota_d \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \\ &\quad + \sup_{u \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(1)}(x) - \eta_N \circ \rho \circ f^{(3)}(x)\|_{\mathcal{H}_T^2} \\ &\leq L\varepsilon_0 + \sup_{u \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(1)}(x) - \eta_N \circ \rho \circ f^{(3)}(x)\|_{\mathcal{H}_T^2} \\ &\leq L\varepsilon_0 + \frac{\varepsilon_D}{3} + ((L \text{diam}(\mathcal{K}_0)^\alpha + 2\varepsilon_2) \sqrt{N}) \varepsilon_3. \end{aligned} \quad (80)$$

Retroactively setting

$$\varepsilon_0 \stackrel{\text{def.}}{=} \frac{\varepsilon_D}{3L} \text{ and } \varepsilon_3 \stackrel{\text{def.}}{=} \frac{\varepsilon_D}{\min\{1, 3(L \text{diam}(\mathcal{K}_0)^\alpha)\sqrt{N}\}}, \quad (82)$$

implying that $\varepsilon_0 \in \mathcal{O}(\varepsilon_D)$ and $\varepsilon_3 \in \mathcal{O}(\varepsilon_D/\sqrt{N})$. Consequentially,

$$\sup_{u \in \mathcal{K}_0} \|F(u) - \eta_N \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \leq \varepsilon_D. \quad (83)$$

Next, we will obtain upper-bound the best α -Hölder constant of $f^{(3)}$ to obtain quantitative parameter estimates on our MLP, which will approximate $f^{(3)}$.

Step 4 - Computing the Regularity of the Surrogate Target Function To, apply a *quantitative* universal approximation theorem, we need a handle on the regularity of the target function being approximated. We first observe that, on the set $\Delta_N^{\varepsilon_2} \stackrel{\text{def.}}{=} H(t_{\varepsilon_2}, \Delta_N)$ the map R , defined in (74), is $L_{f^{(3)}\varepsilon_2, N}$ -Lipschitz with constant given by

$$L_{f^{(3)}\varepsilon_2, N} = \sup_{w \in \Delta_N^{\varepsilon_3}} \|\nabla R(w)\|_2 \leq \sup_{w \in \Delta_N^{\varepsilon_3}} \|\nabla R(w)\|_1 \quad (84)$$

$$= \sup_{w \in \Delta_N} \sum_{i=1}^N \frac{1}{|t_{\varepsilon_3}(w_i - 1/N) + 1/N|} \quad (85)$$

$$\begin{aligned} &\leq \sum_{i=1}^N \frac{1}{\min_{w \in \Delta_N} \lambda |t_{\varepsilon_3}(w_i - 1/N) + 1/N|} \\ &= \frac{N^2}{(1 - t_{\varepsilon_3})} \\ &= N \max\left\{2N, \frac{2(N-1)}{\varepsilon_3}\right\} \\ &\leq \frac{2N^2}{\min\{\varepsilon_3^{-1}, 1/2\}} \stackrel{\text{def.}}{=} \tilde{L}_{f^{(3)}\varepsilon_3, N} \end{aligned} \quad (86)$$

where ∇R denotes the Jacobian of R and where the inequality (84) holds by the Rademacher-Stephanov theorem, see e.g. (Federer, 1969, Theorems 3.1.6-3.1.9).

We thus conclude that, from the estimates in (64), (84)-(86), and the observation that H is 1-Lipschitz, that $f^{(3)}$ is $L_{f^{(3)}}$ -Lipschitz; where

$$\begin{aligned} L_{f^{(3)}} &\leq L_{f^{(2)}} L_{f^{(3)}\varepsilon_2, N} = \frac{\bar{c}N^3}{\varepsilon_1 \min\{\min\{1, \varepsilon_3\}, 1/2\}} \\ &= \frac{9\bar{c}dN^2}{\varepsilon_D \min\{\varepsilon_D/3, 1/2\}} \stackrel{\text{def.}}{=} \tilde{L}, \end{aligned} \quad (87)$$

where $\bar{c} \stackrel{\text{def.}}{=} 4\tilde{c} > 0$. We now construct our deep-learning approximation.

Step 5 - Neural Approximation of Surrogate Target Function $\hat{f}^{(3)}$:

We consider two cases; in the former, the activation function is smooth and in the latter, it is the trainable super-expressive activation (defined in (10)).

1. **Case 1 - σ as in Example 6:** By (Kratsios and Papon, 2022, Proposition 53), there is a MLP $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}^N$ with activation function $\sigma_0 \in C(\mathbb{R})$ (in the notation of Example 6) satisfying

$$\sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|\hat{f}(x) - f^{(3)}(x)\|_2 < \bar{\varepsilon}_A \quad (88)$$

with depth and width given by:

- **Width:** $d + N + 2$
- **Depth:** Finite, and if σ is non-affine and smooth then: $\mathcal{O}\left(N((1-d/4)N)^{2d/\alpha} (2C)^{2d} \varepsilon^{-2d/\alpha}\right)$

where we use the fact that the diameter of \mathcal{K}_0 is at-most $2C$.

2. If σ is the trainable super-expressive activation function in (10), then: for $i = 1, \dots, N$ (Gao et al., 2022, Theorem 1) there exists an MLP $\hat{f}_i : \mathbb{R}^d \rightarrow \mathbb{R}$ with activation function σ_0 satisfying

$$\max_{i=1, \dots, d} \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|\langle f^{(3)}(x), e_i \rangle - \hat{f}_i(x)\|_2 < \bar{\varepsilon}_A/N \quad (89)$$

where $\{e_i\}_{i=1}^N$ is the standard orthonormal basis of \mathbb{R}^N . Moreover, by (Gao et al., 2022, Theorem 1), and the remark directly after, the width, depth, and number of non-zero parameters determining each network is exactly

- **Width:** 11,
- **Depth:** $36(2d + 1)$,
- **No. Params:** $5437(d + 1)(2d + 1)$.

Since $\sigma_1(x) = x$, for all $x \in \mathbb{R}$, then the trainable activation function σ has the 1-identity requirement (see (Cheridito et al., 2021, Definition 4)) applies (Cheridito et al., 2021, Proposition 5) from which we conclude that there exists an MLP $\hat{f} : \mathbb{R}^d \rightarrow \mathbb{R}^N$ with activation function σ satisfying: for each $x \in \mathbb{R}^d$

$$\hat{f}(x) = \sum_{i=1}^N f_i(x) e_i,$$

furthermore, the with, depth, and number of non-zero determining \hat{f} are

- **Width:** $12N \in \mathcal{O}(N)$,
- **Depth:** $d(N - 1) + 36(2d + 1) \in \mathcal{O}(dN)$,
- **No. Params:** at-most $3738N^2(d^2 - 1)N(d + 1)(2d + 1) \in \mathcal{O}(N^3 d^4)$.

Consequentially, (89) implies that

$$\sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(3)}(x) - \hat{f}(x)\|_2 \leq \sum_{i=1}^N \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|\langle f^{(3)}(x), e_i \rangle - \hat{f}_i(x)\|_2 < N \frac{\bar{\varepsilon}_A}{N} = \bar{\varepsilon}_A \quad (90)$$

where e_i is the i^{th} standard basis vector in \mathbb{R}^N with 1 in the i^{th} coordinate and 0 otherwise. We are now ready to complete the proof by combining the estimates from the previous steps.

Step 6 - Putting it All Together:

Set $\hat{F} \stackrel{\text{def.}}{=} \eta_{N_{\varepsilon_1}} \circ \rho \circ \hat{f} \circ \mathcal{E}_d : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$. The estimates in (83) with those in (88) (resp. (90)) yield

$$\begin{aligned} \sup_{u \in \mathcal{K}_0} \|F(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} &\leq \sup_{u \in \mathcal{K}_0} \|F(u) - \eta_{N_{\varepsilon_1}} \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u)\|_{\mathcal{H}_T^2} \\ &\quad + \sup_{u \in \mathcal{K}_0} \|\eta_{N_{\varepsilon_1}} \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} \\ &\leq \varepsilon_D + \sup_{u \in \mathcal{K}_0} \|\eta_{N_{\varepsilon_1}} \circ \rho \circ f^{(3)} \circ \mathcal{E}_d(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} \\ &= \varepsilon_D + \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|\eta_{N_{\varepsilon_1}} \circ \rho_1 \circ \hat{f}(x) - \eta_{N_{\varepsilon_1}} \circ \rho_1 \circ \hat{f}(x)\|_{\mathcal{H}_T^2} \end{aligned}$$

$$\begin{aligned}
&= \varepsilon_D + \text{Lip}(\eta_{N_{\varepsilon_1}} \circ \rho) \sup_{x \in \mathcal{E}_d(\mathcal{K}_0)} \|f^{(3)}(x) - \hat{f}(x)\|_2 \\
&= \varepsilon_D + \text{Lip}(\eta_{N_{\varepsilon_1}} \circ \rho) \bar{\varepsilon}_A.
\end{aligned} \tag{91}$$

Since W is an isometric embedding and the softmax function is at-most 1-Lipschitz then ρ is at-most 1-Lipschitz and $\text{Lip}(\eta_{N_{\varepsilon_1}} \circ \rho) = \text{Lip}(\eta_{N_{\varepsilon_1}})$, which by (67) is at-most

$$\text{Lip}(\eta_{N_{\varepsilon_1}}) \leq (L \text{diam}(\mathcal{K}_0)^\alpha + 2\varepsilon_D/9)\sqrt{N}.$$

Since this upper-bound on $\text{Lip}(\eta_{N_{\varepsilon_1}} \circ \rho)$ depends only on ε_D and is independent of $\bar{\varepsilon}_A$. Consequently for the right-hand side of (91) can be made arbitrarily small by choosing ε_D and $\bar{\varepsilon}_A$ large enough.

It remains to bound N explicitly. There are three cases which we consider here, each of which corresponds to the respective assumptions made on \mathcal{K}_0 and its relationship to the target (non-linear) operator f :

1. **Exponentially Ellipsoidal:** Suppose that $\mathcal{K}_0 \subset \mathcal{H}_T^2$ is such that: for each $\mathcal{K}_0 \ni x = \sum_{i=1}^{\infty} \beta_i s_i$ we have that $|\beta_i| \leq C r^i$. Recall that, e.g. as noted on (Dumer et al., 2004, in Remark 1), that the ε_A -covering number of $p_d(\varepsilon_A^{-1} \cdot \mathcal{K}_0)$ (where $\delta\mathcal{K}_0$ denotes the δ -thickening of \mathcal{K}_0 in \mathbb{R}^d). Therefore, for each such $x \in \mathcal{K}_0$ we have that

$$\sum_{i=1}^d \frac{|\beta_i|^2}{\theta_i^2} \leq 1$$

where the scaling constants $(\theta_i)_{i=1}^{\infty}$ (independent of d) are given by

$$\theta_i = \left(r \left(1 + \left(\frac{C}{\varepsilon_A} \right)^2 \right)^{1/2} \right)^i. \tag{92}$$

Therefore, (Dumer et al., 2004, Theorem 2), with the description of $o(1)$ given in its proof on (Duchi et al., 2011, Equation (40)), by (92) we find that

$$\begin{aligned}
N &\leq \exp \left(\sum_{i=1}^d i \log (Cr (C^{-2} + \varepsilon_A^{-2})^{1/2}) \right) \\
&= \exp \left(\frac{d(d+1)}{2} \log (Cr (C^{-2} + \varepsilon_A^{-2})^{1/2}) \right) \\
&= \left(r (1 + (C/\varepsilon_A)^2)^{1/2} \right)^{\frac{d(d+1)}{2}}
\end{aligned} \tag{93}$$

Since, in this case, $d \in \mathcal{O}(\ln(\varepsilon_D^{-1/r}))$ then there exists some $C_1 > 0$ such that (93) reduces to

$$N \leq \left(r (1 + (C/\varepsilon_A)^2)^{1/2} \right)^{C_1 \ln(\varepsilon_D^{-1/r})^2} \tag{94}$$

2. **Exponential Manifold:** Suppose that \mathcal{K}_0 satisfies Definition B.9. For every $\tilde{\varepsilon}_A > 0$ (to be fixed momentarily), (Lorentz et al., 1996, Proposition 15.1.3), implies that $\tilde{\varepsilon}_A$ -covering number \tilde{N} of the Euclidean unit ball $B_d \stackrel{\text{def.}}{=} \{x \in \mathbb{R}^d : \|x\| \leq 1\}$ is bounded above and below by

$$2^{-d} (\sqrt{d}/\tilde{\varepsilon}_A)^d \leq \tilde{N} \leq 3^d (\sqrt{d}/\tilde{\varepsilon}_A)^d. \tag{95}$$

Since the ‘‘latent parameterization’’ map $\pi : \mathbb{R}^d \rightarrow \mathcal{H}_T^2$ was assumed to be 1-Lipschitz and maps onto \mathcal{K}_0 then, the image of every $\tilde{\varepsilon}_A$ of B_d under π must be $1 \cdot \tilde{\varepsilon}_A$ covering of \mathcal{K}_0 . Set $\tilde{\varepsilon}_A = \varepsilon_A$. Then, (95) implies that

$$N \leq (3\sqrt{d}/\tilde{\varepsilon}_A)^d = (\varepsilon_A^{-1} 3(c(\ln(\varepsilon_D^{-1/r}))^{1/2})^{c(\ln(\varepsilon_D^{-1/r}))})^d$$

where we have used the fact that \mathcal{K}_0 is contained in an (r, f) -exponentially ellipsoidal set to deduce that $d \leq c \ln(\varepsilon_D^{-1/r})$ for some absolute constant $c > 0$.

3. **General Case:** In the case of general \mathcal{K}_0 , (Acciaio et al., 2023, Lemma 7.1) and the upper-bound of 2^{d+1} on the doubling constant of $\mathcal{E}_d(\mathcal{K}_0)$, just prior to Equation (64), implies that

$$N \leq \left(2^{(d+1)}\right)^{\log_2(\text{diam}(\mathcal{K}_0)) - \frac{1}{\alpha} \log_2(\varepsilon_D/\tilde{L}) + \frac{1}{\alpha} \log_2(\tilde{c}d)}$$

for some absolute constant $\tilde{c} > 0$.

Setting $\varepsilon_A \stackrel{\text{def.}}{=} \bar{\varepsilon}_A / (c(L \text{diam}(\mathcal{K}_0)^\alpha + \varepsilon_D)\sqrt{N}) \in \mathcal{O}\left(\frac{\bar{\varepsilon}_A}{\varepsilon_D \sqrt{N}}\right)$ yields the conclusion.

Step 7 - Elucidating the Model

Let $V \in \mathbb{R}^{N \times Q}$ be such that, for $n = 1, \dots, N$ and $i = 1, \dots, Q$, $V_{n,q} \stackrel{\text{def.}}{=} \langle F \circ \iota_d(x_n), s_i \rangle_{\mathcal{H}_T^2}$. Then, (65) implies that: for each $w \in \Delta_N$

$$\eta(w) = \sum_{n=1}^N \left(\sum_{i=1}^Q \langle F \circ \iota_d(x_n), s_i \rangle_{\mathcal{H}_T^2} \right) = \sum_{n=1}^N \sum_{i=1}^Q w_n V_{n,q} \quad (96)$$

For either σ is smooth or σ as in (89), consider the MLP with σ activation function $\mathcal{V} : \mathbb{R}^d \rightarrow \mathbb{R}^{N \times Q}$ given for each $x \in \mathbb{R}^d$ by

$$\mathcal{V}(x) \stackrel{\text{def.}}{=} \mathbf{0}^{(1)} \sigma \bullet (\mathbf{0}^{(2)} x + \mathbf{0}^{(3)}) + V$$

where $\mathbf{0}^{(1)}$ is the $ND \times 1$ zero matrix, $\mathbf{0}^{(2)}$ is the $1 \times d$ zero matrix, and $\mathbf{0}^{(3)} = (0) \in \mathbb{R}$, and where we have identified $\mathbb{R}^{N \times D}$ with \mathbb{R}^{ND} . In either case, observe that the number of non-zero parameters defining \mathcal{V} are at-most ND .

For each $n = 1, \dots, N$ and $i = 1, \dots, Q$, we $V^{(n,i)} \stackrel{\text{def.}}{=} s_i$. By construction: for each $u \in \mathcal{H}_T^2$ and every $w \in \Delta_N$ we have that

$$\mathcal{D}(w, u) \stackrel{\text{def.}}{=} \sum_{n=1}^N w_n [\mathcal{V} \circ \mathcal{E}_d(u)]_n V^{(n,q)} = \eta_N(w). \quad (97)$$

Since the map W in the definition of ρ , see the line just below (73), was affine then our approximation $\hat{F} = \eta_{N_{\varepsilon_1}} \circ \rho \circ \hat{f} \circ \mathcal{E}_d : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ is of the form in Definition 6.

The number of non-zero parameters defining the model \hat{F} are, therefore, at-most

$$\underbrace{\text{Depth}(\hat{f}) \text{Width}(\hat{f})^2}_{\text{No. Param. } \hat{f}} + \underbrace{NQ}_{\text{No. Param. } \mathcal{V}} \quad (98)$$

where the depth and width of \hat{f} were computed in step 4. In particular, if σ is the trainable super-expressive activation function in (89) then the quantity in (98) is $\mathcal{O}(N(N^2 d^4 + Q))$. \square

Proof of Theorem 12. Fix a non-empty compact subset $\mathcal{K}_0 \subset \mathcal{H}_T^2$, a continuous function $f : \mathcal{K}_0 \rightarrow \mathcal{H}_T^2$, and a $\varepsilon > 0$.

Let $\mathcal{K}_0 \subset \mathcal{H}_T^2$ be non-empty and compact. We would like to use (Miculescu, 2002/03, Theorem 1) to reduce the problem of approximating f to ε precision, to the problem of approximating an $\varepsilon/2$ Lipschitz approximation of our target function f to $\varepsilon/2$ precision. Thus, we will be able to employ our technical approximation theorem for Lipschitz maps between \mathcal{H}_T^2 to itself, to deduce our conclusion. On a technical note, we do not argue on the domain \mathcal{H}_T^2 but rather on the compact subspace \mathcal{K}_0 , since all continuous functions are both bounded and uniformly continuous thereon; which then directly allows us to (Miculescu, 2002/03, Theorem 1) which only allows us to uniformly approximate bounded continuous functions on compacta.

Step 1 - Verification of Lipschitz Extension Property

In order to apply (Miculescu, 2002/03, Theorem 1) we will need to show that the pair $(\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2})$ and \mathcal{H}_T^2 have the so-called ‘‘Lipschitz extension property’’ (as named in (Miculescu, 2002/03, Theorem 1)). This means that the Lipschitz function from $(B, \|\cdot\|_{\mathcal{H}_T^2})$ to \mathcal{H}_T^2 , for any subset B of \mathcal{K}_0 , can be extended to a Lipschitz function of all of \mathcal{K}_0 with roughly the same Lipschitz constant.

Let $B \subseteq \mathcal{K}_0$. Note that, as $\mathcal{K}_0 \subset \mathcal{H}_T^2$ then, B is a subset of \mathcal{H}_T^2 . Since \mathcal{H}_T^2 is a separable Hilbert space then the extension theorem of (Benyamini and Lindenstrauss, 2000, Theorem 1.12): for every $L \geq 0$ and each L -Lipschitz (non-linear operator) $g : (B, \|\cdot\|_{\mathcal{H}_T^2}) \rightarrow \mathcal{H}_T^2$ there exists an L -Lipschitz extension $\tilde{G} : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$; i.e. \tilde{G} is L -Lipschitz

$$\tilde{G}|_B = g. \quad (99)$$

Since the restriction operator $\iota_{\mathcal{K}_0} : \mathcal{H}_T^2 \ni \tilde{g} \rightarrow \tilde{g}|_{\mathcal{K}_0} \in \mathcal{K}_0$ is 1-Lipschitz then the composite map $G \stackrel{\text{def.}}{=} \iota_{\mathcal{K}_0} \circ \tilde{G} = \tilde{G}|_{\mathcal{K}_0} : (\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2}) \rightarrow \mathcal{H}_T^2$ is L -Lipschitz. Furthermore, (99) and the inclusion of B in \mathcal{K}_0 imply that

$$G|_B = (\tilde{G}|_B)|_{\mathcal{K}_0} = g. \quad (100)$$

Thus, G is an L -Lipschitz extension of g to $(\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2})$. Thus, the pair $(\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2})$ and \mathcal{H}_T^2 has the Lipschitz extension property; thus (Miculescu, 2002/03, Theorem 1) implies that the space of Lipschitz functions from $(\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2})$ to \mathcal{H}_T^2 is *dense* in the space of uniformly continuous and bounded functions from $(\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2})$ to \mathcal{H}_T^2 with respect to the uniform norm.

Step 2 - $\varepsilon/2$ -Approximation of f by Lipschitz Maps

Since \mathcal{K}_0 is compact and f is continuous on \mathcal{K}_0 then f is uniformly continuous and bounded thereon. By (Miculescu, 2002/03, Theorem 1), we deduce that there exists a Lipschitz function $\tilde{f}_\varepsilon : (\mathcal{K}_0, \|\cdot\|_{\mathcal{H}_T^2}) \rightarrow \mathcal{H}_T^2$ satisfying

$$\max_{u \in \mathcal{K}_0} \|f(u) - \tilde{f}_\varepsilon(u)\|_{\mathcal{H}_T^2} \leq \varepsilon/2. \quad (101)$$

Again applying (Benyamini and Lindenstrauss, 2000, Theorem 1.12), we deduce that \tilde{f}_ε admits a Lipschitz extension $f_\varepsilon : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$, with the same Lipschitz constant. Since f_ε is a Lipschitz extension of \tilde{f}_ε , beyond \mathcal{K}_0 , then (101) implies that

$$\max_{u \in \mathcal{K}_0} \|f(u) - f_\varepsilon(u)\|_{\mathcal{H}_T^2} = \max_{u \in \mathcal{K}_0} \|f(u) - \tilde{f}_\varepsilon(u)\|_{\mathcal{H}_T^2} \leq \varepsilon/2. \quad (102)$$

Step 3 - $\varepsilon/2$ -Approximation of f_ε by Attentional Neural Operator

Since \mathcal{K}_0 is a compact subset of \mathcal{H}_T^2 and $f_\varepsilon : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ is Lipschitz then Lemma B.7 applies. Whence, there exists an attentional neural operator $\hat{F} : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ satisfying

$$\max_{u \in \mathcal{K}_0} \|f_\varepsilon(u) - \hat{F}(u)\|_{\mathcal{H}_T^2}. \quad (103)$$

Combining (102) and (103) yield

$$\max_{u \in \mathcal{K}_0} \|f(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} \leq \max_{u \in \mathcal{K}_0} \|f(u) - f_\varepsilon(u)\|_{\mathcal{H}_T^2} + \max_{u \in \mathcal{K}_0} \|f_\varepsilon(u) - \hat{F}(u)\|_{\mathcal{H}_T^2} \leq \varepsilon/2 + \varepsilon/2 = \varepsilon,$$

which concludes our proof. \square

B.5 Proof of Main Stackelberg Equilibria Results

Theorems 7 and 8 are two parts of a larger whole. As such, their derivation is most naturally merged into a single proof; we now do.

Joint Proof of Theorems 7, 8, and 11. Let $\hat{U} : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$ be a map, to be fixed retroactively. For any $d \in \mathbb{N}_+$, let $p_d : \mathcal{H}_T^2 \rightarrow \text{span}\{s_i\}_{i=1}^d$ be the orthogonal projection; i.e. $p_d(\sum_{i=1}^{\infty} \beta_i s_i) = \sum_{i=1}^d \beta_i s_i$ for all $u = \sum_{i=1}^{\infty} \beta_i s_i \in \mathcal{H}_T^2$. We will retroactively adjust d . Fix $u^0 \in \mathcal{K}_0$, denote $\hat{u}_d^0 \stackrel{\text{def.}}{=} p_d(u^0)$ and compute

$$\begin{aligned} |J_0(u^0, U^*(u^0)) - J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))| &\leq \underbrace{|J_0(u^0, U^*(u^0)) - J_0(\hat{u}_d^0, U^*(\hat{u}_d^0))|}_{\text{(I)}} \\ &\quad + \underbrace{|J_0(\hat{u}_d^0, U^*(\hat{u}_d^0)) - J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))|}_{\text{(II)}}. \end{aligned} \quad (104)$$

Step 1 - Bounding Term Term (I) By Lemma B.5, we can bound Term (I) from above by

$$\text{(I)} = |J_0(u^0, U^*(u^0)) - J_0(\hat{u}_d^0, U^*(\hat{u}_d^0))| \leq \tilde{\omega}(\|u^0 - \hat{u}_d^0\|_{\mathcal{H}_T^2}) \quad (105)$$

where $\tilde{\omega}(t) = C \max\{|t|, |t|^{1/2}\}$ for each $t \in \mathbb{R}$ and for some constant $C \geq 0$ depending only on T . Note that, $t \mapsto \tilde{\omega}(t)$ is continuous, monotonically increasing on $[0, \infty)$ and subjective, thus it is a homomorphism of $[0, \infty)$ to itself with continuous inverse given by

$$\tilde{\omega}(t) \stackrel{\text{def.}}{=} \begin{cases} (t/C)^2 & \text{if } 0 \leq t \leq 1 \\ t/C & \text{if } 1 \leq t. \end{cases}$$

We emphasize that, $\tilde{\omega}$ has range $[0, \infty)$.

Since \mathcal{H}_T^2 has the 1-bounded approximation property implemented by the (finite-rank) projection operators $(p_d)_{d \in \mathbb{N}_+}$ then choosing d large enough, we may ensure that $\sup_{u^0 \in \mathcal{K}_0} \|u^0 - \hat{u}_d^0\|_{\mathcal{H}_T^2} < \tilde{\omega}(\varepsilon/2)$. The right-hand side of (105) can be bounded-above as follows

$$\text{(I)} \leq \tilde{\omega}(\|u^0 - \hat{u}_d^0\|_{\mathcal{H}_T^2}) \leq \tilde{\omega}\left(\sup_{u^0 \in \mathcal{K}_0} \|u^0 - \hat{u}_d^0\|_{\mathcal{H}_T^2}\right) \leq \tilde{\omega}\left(\tilde{\omega}\left(\frac{\varepsilon}{2}\right)\right) = \frac{\varepsilon}{2}. \quad (106)$$

It remains to bound Term (II).

Step 2 - Bounding Term Term (II) By Lemma B.2, we find that

$$\text{(II)} = |J_0(\hat{u}_d^0, U^*(\hat{u}_d^0)) - J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))| \leq C \cdot \|U^*(\hat{u}_d^0) - \hat{U}(\hat{u}_d^0)\|_{\mathcal{H}_T^2}. \quad (107)$$

By our universal approximation theorem, in Lemma B.7, we have that there exists an attentional neural operator $\hat{U} : \mathcal{H}_T^2 \rightarrow \mathcal{H}_T^2$, as in Definition 6

$$\sup_{v \in p_d(\mathcal{K}_0)} \|U^*(v) - \hat{U}(v)\|_{\mathcal{H}_T^2} \leq \frac{\varepsilon}{2C} \quad (108)$$

where we have set $\varepsilon_D \stackrel{\text{def.}}{=} \varepsilon_A \stackrel{\text{def.}}{=} \varepsilon/4C$. Thus, the estimate in (108) implies that the right-hand side of (107) can be bounded above as follows

$$\text{(II)} \leq C \|\hat{U}(\hat{u}_d^0) - U^*(\hat{u}_d^0)\|_{\mathcal{H}_T^2} \leq C \sup_{v \in p_d(\mathcal{K}_0)} \|U^*(v) - \hat{U}(v)\|_{\mathcal{H}_T^2} \leq C \frac{\varepsilon}{2C} = \frac{\varepsilon}{2}. \quad (109)$$

Upon combining the estimates in (106) and in (109) we obtain the following upper-bound for the right-hand side of (104)

$$|J_0(u^0, U^*(u^0)) - J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))| \leq \text{(I)} + \text{(II)} \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon. \quad (110)$$

Step 3 Suppose additionally that $u^0 \in \mathcal{K}_0$ is such that $(u^0, U^*(u^0))$ is a (0-)Stackelberg equilibrium. If $(u^0, U^*(u^0))$ is a Stackelberg equilibrium, see Definition 3, then this pair is optimal for J_1 and the left-hand side of (110) becomes

$$0 \leq |J_0(u^0, U^*(u^0)) - J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0))| = J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0)) - J_0(u^0, U^*(u^0)). \quad (111)$$

Consequently, (111) can be rearranged yielding

$$J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0)) \leq J_0(u^0, U^*(u^0)) + \varepsilon. \quad (112)$$

Since the left-hand side of (112) is optimal, then taking infima overall v in \mathcal{H}_T^2 does not reduce it further. Thus, (112) becomes

$$\inf_{u^0 \in \mathcal{K}_0} J_0(u^0, U^*(u^0)) = J_0(u^0, U^*(u^0)). \quad (113)$$

Combining (112) and (113) yields

$$J_0(\hat{u}_d^0, \hat{U}(\hat{u}_d^0)) \leq J_0(u^0, U^*(u^0)) + \varepsilon = \inf_{u^0 \in \mathcal{K}_0} J_0(u^0, U^*(u^0)) + \varepsilon.$$

This concludes our proofs of Theorems 7 and 8.

To obtain the conclusion of Theorem 11, we first note that the non-linear operator U^* is 1/2-Hölder continuous. Since Example 9 showed that K is an exponential manifold (in the sense of Definition B.9) then the complexity estimates for the attentional neural operator \hat{U} selected in (108) must be as in Table 2. Consider the special case where $\varepsilon_D = \varepsilon_A$ completes the proof of Theorem 11. \square

C. Acknowledgments

A. Kratsios acknowledges financial support from an NSERC Discovery Grant No. RGPIN-2023-04482 and No. DGEER-2023-00230. A. Kratsios also acknowledges that resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute⁷. I. Ekren is partially funded by NSF grant DMS-2406240.

References

- Beatrice Acciaio, Anastasis Kratsios, and Gudmund Pammer. Designing universal causal deep learning models: The geometric (hyper) transformer. *Mathematical Finance*, 2023.
- Guillermo Alonso Alvarez, Sergey Nadtochiy, and Kevin Webster. Optimal brokerage contracts in almgren–chriss model with multiple clients. *SIAM Journal on Financial Mathematics*, 14(3): 855–878, 2023.
- Bo An, Milind Tambe, Fernando Ordonez, Eric Shieh, and Christopher Kiekintveld. Refinement of strong stackelberg equilibria in security games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25, pages 587–593, 2011.
- Anima Anandkumar, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Nikola Kovachki, Zongyi Li, Burigede Liu, and Andrew Stuart. Neural operator: Graph kernel network for partial differential equations. In *ICLR 2020 Workshop on Integration of Deep Neural Models and Differential Equations*, 2020.

⁷. <https://vectorinstitute.ai/partnerships/current-partners/>

- Alexander Aurell, Rene Carmona, Gokce Dayanikli, and Mathieu Lauriere. Optimal incentives to mitigate epidemics: a stackelberg mean field game approach. *SIAM Journal on Control and Optimization*, 60(2):S294–S322, 2022.
- David R. Baños, Sindre Duedahl, Thilo Meyer-Brandis, and Frank Proske. Construction of Malliavin differentiable strong solutions of SDEs under an integrability condition on the drift without the Yamada-Watanabe principle. *Ann. Inst. Henri Poincaré Probab. Stat.*, 54(3):1464–1491, 2018. ISSN 0246-0203,1778-7017. doi: 10.1214/17-AIHP845. URL <https://doi.org/10.1214/17-AIHP845>.
- Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *3rd International Conference on Learning Representations, ICLR 2015*, 2015.
- Peter L. Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight VC-dimension and pseudodimension bounds for piecewise linear neural networks. *J. Mach. Learn. Res.*, 20:Paper No. 63, 17, 2019. ISSN 1532-4435,1533-7928.
- Francesca Bartolucci, Emmanuel de Bezenac, Bogdan Raonic, Roberto Molinaro, Siddhartha Mishra, and Rima Alaifari. Representation equivalent neural operators: a framework for alias-free operator learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- J Benitez, Takashi Furuya, Florian Faucher, Anastasis Kratsios, Xavier Tricoche, and Maarten V de Hoop. Out-of-distributional risk bounds for neural operators with applications to the helmholtz equation. *arXiv preprint arXiv:2301.11509*, 2023.
- Sarah Bensalem, Nicolás Hernández Santibáñez, and Nabil Kazi-Tani. Prevention efforts, insurance demand and price incentives under coherent risk measures. *Insurance: Mathematics and Economics*, 93:369–386, 2020.
- Alain Bensoussan, Shaokuan Chen, and Suresh P Sethi. The maximum principle for global solutions of stochastic stackelberg differential games. *SIAM Journal on Control and Optimization*, 53(4):1956–1981, 2015.
- Yoav Benyamini and Joram Lindenstrauss. *Geometric nonlinear functional analysis. Vol. 1*, volume 48 of *American Mathematical Society Colloquium Publications*. American Mathematical Society, Providence, RI, 2000. ISBN 0-8218-0835-4. doi: 10.1090/coll/048. URL <https://doi.org/10.1090/coll/048>.
- Elia Bruè, Simone Di Marino, and Federico Stra. Linear Lipschitz and C^1 extension operators through random projection. *J. Funct. Anal.*, 280(4):Paper No. 108868, 21, 2021. ISSN 0022-1236,1096-0783. doi: 10.1016/j.jfa.2020.108868. URL <https://doi.org/10.1016/j.jfa.2020.108868>.
- Carlos Calderon-Macias. *Artificial neural systems for interpretation and inversion of seismic data*. The University of Texas at Austin, 1997.
- Edoardo Calvello, Nikola B Kovachki, Matthew E Levine, and Andrew M Stuart. Continuum attention for neural operators. *arXiv preprint arXiv:2406.06486*, 2024.
- Jingyi Cao, Dongchen Li, Virginia R Young, and Bin Zou. Stackelberg differential game for insurance under model ambiguity. *Insurance: Mathematics and Economics*, 106:128–145, 2022.
- Shuhao Cao. Choose a transformer: Fourier or galerkin. *Advances in neural information processing systems*, 34:24924–24940, 2021.

- Bernd Carl. Metric entropy of convex hulls in Hilbert spaces. *Bull. London Math. Soc.*, 29(4): 452–458, 1997. ISSN 0024-6093,1469-2120. doi: 10.1112/S0024609397003044. URL <https://doi.org/10.1112/S0024609397003044>.
- René Carmona and Gökçe Dayanıklı. Mean field game model for an advertising competition in a duopoly. *International Game Theory Review*, 23(04):2150024, 2021.
- Javier Castro. The kolmogorov infinite dimensional equation in a hilbert space via deep learning methods. *Journal of Mathematical Analysis and Applications*, 527(2):127413, 2023.
- Patrick Cheridito, Arnulf Jentzen, and Florian Rossmannek. Efficient approximation of high-dimensional functions with neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- Samuel N. Cohen and Robert J. Elliott. *Stochastic calculus and applications*. Probability and its Applications. Springer, Cham, second edition, 2015. ISBN 978-1-4939-2866-8; 978-1-4939-2867-5. doi: 10.1007/978-1-4939-2867-5. URL <https://doi.org/10.1007/978-1-4939-2867-5>.
- Vincent Conitzer and Tuomas Sandholm. Computing the optimal strategy to commit to. In *Proceedings of the 7th ACM conference on Electronic commerce*, pages 82–90, 2006.
- Gianni Dal Maso. *An introduction to Γ -convergence*, volume 8 of *Progress in Nonlinear Differential Equations and their Applications*. Birkhäuser Boston, Inc., Boston, MA, 1993. ISBN 0-8176-3679-X. doi: 10.1007/978-1-4612-0327-8. URL <https://doi.org/10.1007/978-1-4612-0327-8>.
- Gokce Dayanikli and Mathieu Lauriere. A machine learning method for stackelberg mean field games. *arXiv preprint arXiv:2302.10440*, 2023.
- Maarten de Hoop, J. Antonio Lara B., Anastasis Kratsios, Matti Lassas, and Takashi Furuya. Mixture of experts soften the curse of dimensionality in operator learning. *arXiv preprint arXiv:2404.09101*, 2024.
- Maarten V de Hoop, Matti Lassas, and Christopher A Wong. Deep learning architectures for non-linear operator functions and nonlinear inverse problems. *Mathematical Statistics and Learning*, 4(1):1–86, 2022.
- Tim De Ryck, Ameya D Jagtap, and Siddhartha Mishra. Error estimates for physics-informed neural networks approximating the navier–stokes equations. *IMA Journal of Numerical Analysis*, 44(1): 83–119, 2024.
- Ludwig Dierks and Sven Seuken. Cloud pricing: The spot market strikes back. *Management Science*, 68(1):105–122, 2022.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, 2011. ISSN 1532-4435.
- Ilya Dumer, Mark S. Pinsky, and Viacheslav V. Prelov. On coverings of ellipsoids in Euclidean spaces. *IEEE Trans. Inform. Theory*, 50(10):2348–2356, 2004. ISSN 0018-9448,1557-9654. doi: 10.1109/TIT.2004.834759. URL <https://doi.org/10.1109/TIT.2004.834759>.
- Romuald Elie, Thibaut Mastrolia, and Dylan Possamai. A tale of a principal and many, many agents. *Mathematics of Operations Research*, 44(2):440–467, 2019.
- Lawrence C. Evans. *Partial differential equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, second edition, 2010. ISBN 978-0-8218-4974-3. doi: 10.1090/gsm/019. URL <https://doi.org/10.1090/gsm/019>.

- VS Fanaskov and Ivan V Oseledets. Spectral neural operators. In *Doklady Mathematics*, pages 1–7. Springer, 2024.
- Herbert Federer. *Geometric measure theory*. Die Grundlehren der mathematischen Wissenschaften, Band 153. Springer-Verlag New York, Inc., New York, 1969.
- Takashi Furuya and Anastasis Kratsios. Simultaneously solving families of fbsdes with neural operators of logarithmic depth, constant width, and sub-linear rank. *arXiv preprint arXiv:2409.12335*, 2024.
- Luca Galimberti, Anastasis Kratsios, and Giulia Livieri. Designing universal causal deep learning models: The case of infinite-dimensional dynamical systems from stochastic analysis. *arXiv preprint arXiv:2210.13300*, 2022.
- Xiao-shan Gao, Shuang Liu, and Lijia Yu. Achieving optimal adversarial accuracy for adversarial deep learning using stackelberg games. *Acta Mathematica Scientia*, 42(6):2399–2418, 2022.
- Matthias Gerstgrasser and David C Parkes. Oracles & followers: Stackelberg equilibria in deep multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 11213–11236. PMLR, 2023.
- Mario Ghossoub and Michael B Zhu. Stackelberg equilibria with multiple policyholders. *Insurance: Mathematics and Economics*, 116:189–201, 2024.
- Somdatta Goswami, Minglang Yin, Yue Yu, and George Em Karniadakis. A physics-informed variational deeponet for predicting crack path in quasi-brittle materials. *Computer Methods in Applied Mechanics and Engineering*, 391:114587, 2022.
- Nika Haghtalab, Chara Podimata, and Kunhe Yang. Calibrated stackelberg games: Learning optimal commitments against calibrated agents. *Advances in Neural Information Processing Systems*, 36, 2023.
- Zhongkai Hao, Zhengyi Wang, Hang Su, Chengyang Ying, Yinpeng Dong, Songming Liu, Ze Cheng, Jian Song, and Jun Zhu. Gnot: A general neural operator transformer for operator learning. In *International Conference on Machine Learning*, pages 12556–12569. PMLR, 2023.
- Keegan Harris, Steven Wu, and Maria Florina Balcan. Stackelberg games with side information. In *Multi-Agent Security Workshop @ NeurIPS’23*, 2023. URL <https://openreview.net/forum?id=4RFv40DWkp>.
- Xiuli He, Ashutosh Prasad, and Suresh P Sethi. Cooperative advertising and pricing in a dynamic stochastic supply chain: Feedback stackelberg strategies. In *PICMET’08-2008 Portland International Conference on Management of Engineering & Technology*, pages 1634–1649. IEEE, 2008.
- Camilo Hernández and Dylan Possamai. Time-inconsistent contract theory. *Mathematical Finance*, 34(3):1022–1085, 2024.
- Camilo Hernández, Nicolás Hernández Santibáñez, Emma Hubert, and Dylan Possamai. Closed-loop equilibria for stackelberg games: it’s all about stochastic targets. *arXiv preprint arXiv:2406.19607*, 2024.
- Ruiyang Hong and Anastasis Kratsios. Bridging the gap between approximation and learning via optimal approximation by relu mlps of maximal regularity. *arXiv preprint arXiv:2409.12335*, 2024.

- Emma Hubert, Thibaut Mastrolia, Dylan Possamaï, and Xavier Warin. Incentives, lockdown, and testing: from thucydides’ analysis to the covid-19 pandemic. *Journal of mathematical biology*, 84(5):37, 2022.
- Kazufumi Ito, Christoph Reisinger, and Yufei Zhang. A neural network-based policy iteration algorithm with global H^2 -superlinear convergence for stochastic games on domains. *Found. Comput. Math.*, 21(2):331–374, 2021. ISSN 1615-3375,1615-3383. doi: 10.1007/s10208-020-09460-1. URL <https://doi.org/10.1007/s10208-020-09460-1>.
- Debarun Kar, Thanh H Nguyen, Fei Fang, Matthew Brown, Arunesh Sinha, Milind Tambe, and Albert Xin Jiang. Trends and applications in stackelberg security games. *Handbook of dynamic game theory*, pages 1–47, 2017.
- Jussi Keppo, Nizar Touzi, and Ruiting Zuo. Dynamic contracting in asset management under the investor-partner-manager relationship. *Operations Research*, 72(3):903–915, 2024.
- Patrick Kidger and Terry Lyons. Universal Approximation with Deep Narrow Networks. In Jacob Abernethy and Shivani Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pages 2306–2327. PMLR, 09–12 Jul 2020. URL <https://proceedings.mlr.press/v125/kidger20a.html>.
- Kai A Konrad and Wolfgang Leininger. The generalized stackelberg equilibrium of the all-pay auction with complete information. *Review of Economic Design*, 11:165–174, 2007.
- Nikola Kovachki, Samuel Lanthaler, and Siddhartha Mishra. On universal approximation and error bounds for Fourier neural operators. *J. Mach. Learn. Res.*, 22:Paper No. [290], 76, 2021. ISSN 1532-4435,1533-7928.
- Nikola Kovachki, Zongyi Li, Burigede Liu, Kamyar Azizzadenesheli, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Neural operator: Learning maps between function spaces with applications to pdes. *Journal of Machine Learning Research*, 24(89):1–97, 2023.
- Anastasis Kratsios and Léonie Papon. Universal approximation theorems for differentiable geometric deep learning. *J. Mach. Learn. Res.*, 23:Paper No. [196], 73, 2022. ISSN 1532-4435,1533-7928.
- Anastasis Kratsios, Chong Liu, Matti Lassas, Maarten V de Hoop, and Ivan Dokmanić. An approximation theory for metric space-valued functions with a view towards deep learning. *arXiv preprint arXiv:2304.12231*, 2023a.
- Anastasis Kratsios, Chong Liu, Matti Lassas, Maarten V de Hoop, and Ivan Dokmanić. Universal geometric deep learning via geometric attention. *arXiv preprint arXiv:2304.12231*, 2023b.
- Anastasis Kratsios, Takashi Furuya, Antonio Lara, Matti Lassas, and Maarten de Hoop. Mixture of experts soften the curse of dimensionality in operator learning. *arxiv*, 2024.
- Emma Kroell, Sebastian Jaimungal, and Silvana M Pesenti. Optimal robust reinsurance with multiple insurers. *arXiv preprint arXiv:2308.11828*, 2023.
- Samuel Lanthaler. Operator learning with PCA-Net: upper and lower complexity bounds. *Journal of Machine Learning Research*, 24(318):1–67, 2023.
- Samuel Lanthaler and Andrew M Stuart. The curse of dimensionality in operator learning. *arXiv preprint arXiv:2306.15924*, 2023.
- Samuel Lanthaler, Siddhartha Mishra, and George E Karniadakis. Error estimates for deepnets: A deep learning framework in infinite dimensions. *Transactions of Mathematics and Its Applications*, 6(1):tnac001, 2022a.

- Samuel Lanthaler, Roberto Molinaro, Patrik Hadorn, and Siddhartha Mishra. Nonlinear reconstruction for operator learning of pdes with discontinuities. *arXiv preprint arXiv:2210.01074*, 2022b.
- Jae Yong Lee, SungWoong CHO, and Hyung Ju Hwang. HyperdeepONet: learning operator with complex target function space using the limited resources via hypernetwork. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=0Aw6V3ZAhSd>.
- Guanguan Li, Qiqiang Li, Yi Liu, Huimin Liu, Wen Song, and Ran Ding. A cooperative stackelberg game based energy management considering price discrimination and risk assessment. *International Journal of Electrical Power & Energy Systems*, 135:107461, 2022.
- Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric partial differential equations. *arXiv preprint arXiv:2010.08895*, 2020.
- Zongyi Li, Daniel Zhengyu Huang, Burigede Liu, and Anima Anandkumar. Fourier neural operator with learned deformations for pdes on general geometries. *Journal of Machine Learning Research*, 24(388):1–26, 2023.
- G. G. Lorentz. Metric entropy and approximation. *Bull. Amer. Math. Soc.*, 72:903–937, 1966. ISSN 0002-9904. doi: 10.1090/S0002-9904-1966-11586-0. URL <https://doi.org/10.1090/S0002-9904-1966-11586-0>.
- George G. Lorentz, Manfred v. Golitschek, and Yuly Makovoz. *Constructive approximation*, volume 304 of *Grundlehren der mathematischen Wissenschaften [Fundamental Principles of Mathematical Sciences]*. Springer-Verlag, Berlin, 1996. ISBN 3-540-57028-4. doi: 10.1007/978-3-642-60932-9. URL <https://doi.org/10.1007/978-3-642-60932-9>. Advanced problems.
- Lu Lu, Pengzhan Jin, and George Em Karniadakis. Deeponet: Learning nonlinear operators for identifying differential equations based on the universal approximation theorem of operators. *arXiv preprint arXiv:1910.03193*, 2019.
- Lu Lu, Pengzhan Jin, Guofei Pang, Zhongqiang Zhang, and George Em Karniadakis. Learning nonlinear operators via deeponet based on the universal approximation theorem of operators. *Nature machine intelligence*, 3(3):218–229, 2021.
- Carlo Marcati and Christoph Schwab. Exponential convergence of deep operator networks for elliptic partial differential equations. *SIAM Journal on Numerical Analysis*, 61(3):1513–1545, 2023.
- Yves Meyer. *Ondelettes et opérateurs. I*. Actualités Mathématiques. [Current Mathematical Topics]. Hermann, Paris, 1990. ISBN 2-7056-6125-0. Ondelettes. [Wavelets].
- Radu Miculescu. Approximations by Lipschitz functions generated by extensions. *Real Anal. Exchange*, 28(1):33–40, 2002/03. ISSN 0147-1937,1930-1219. doi: 10.14321/realanalexch.28.1.0033. URL <https://doi.org/10.14321/realanalexch.28.1.0033>.
- Roberto Molinaro, Yunan Yang, Björn Engquist, and Siddhartha Mishra. Neural inverse operators for solving pde inverse problems. In *Proceedings of the 40th International Conference on Machine Learning*, pages 25105–25139, 2023.
- James R. Munkres. *Topology*. Prentice Hall, Inc., Upper Saddle River, NJ, second edition, 2000. ISBN 0-13-181629-2.

- David Nualart. *The Malliavin calculus and related topics*. Probability and its Applications (New York). Springer-Verlag, Berlin, second edition, 2006. ISBN 978-3-540-28328-7; 3-540-28328-5.
- Shige Peng and Zhen Wu. Fully coupled forward-backward stochastic differential equations and applications to optimal control. *SIAM Journal on Control and Optimization*, 37(3):825–843, 1999.
- Guergana Petrova and Przemysław Wojtaszczyk. Lipschitz widths. *Constructive Approximation*, 57(2):759–805, 2023.
- Bogdan Raonic, Roberto Molinaro, Tim De Ryck, Tobias Rohner, Francesca Bartolucci, Rima Alai-fari, Siddhartha Mishra, and Emmanuel de Bézenac. Convolutional neural operators for robust and accurate learning of pdes. *Advances in Neural Information Processing Systems*, 36, 2024.
- Christoph Reisinger and Yufei Zhang. Rectified deep neural networks overcome the curse of dimensionality for nonsmooth value functions in zero-sum games of nonlinear stiff systems. *Analysis and Applications*, 18(06):951–999, 2020.
- James C. Robinson. *Dimensions, embeddings, and attractors*, volume 186 of *Cambridge Tracts in Mathematics*. Cambridge University Press, Cambridge, 2011. ISBN 978-0-521-89805-8.
- Zuwei Shen, Haizhao Yang, and Shijun Zhang. Optimal approximation rate of ReLU networks in terms of width and depth. *J. Math. Pures Appl. (9)*, 157:101–135, 2022. ISSN 0021-7824. doi: 10.1016/j.matpur.2021.07.009. URL <https://doi.org/10.1016/j.matpur.2021.07.009>.
- Xiaofei Shi, Daran Xu, and Zhanhao Zhang. Deep learning algorithms for hedging with frictions. *Digital Finance*, 5(1):113–147, 2023.
- Stanisław J. Szarek. A Banach space without a basis which has the bounded approximation property. *Acta Math.*, 159(1-2):81–98, 1987. ISSN 0001-5962,1871-2509. doi: 10.1007/BF02392555. URL <https://doi.org/10.1007/BF02392555>.
- Aad W. van der Vaart and Jon A. Wellner. *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York, 1996. ISBN 0-387-94640-3. doi: 10.1007/978-1-4757-2545-2. URL <https://doi.org/10.1007/978-1-4757-2545-2>. With applications to statistics.
- Jan van Mill. *The infinite-dimensional topology of function spaces*, volume 64 of *North-Holland Mathematical Library*. North-Holland Publishing Co., Amsterdam, 2001. ISBN 0-444-50557-1.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Nik Weaver. *Lipschitz algebras*. World Scientific Publishing Co. Pte. Ltd., Hackensack, NJ, second edition, 2018. ISBN 978-981-4740-63-0.
- Jiongmin Yong. A leader-follower stochastic linear quadratic differential game. *SIAM Journal on Control and Optimization*, 41(4):1015–1041, 2002.
- Hui Zhang, Fengrui Zhang, Bing Gong, Xuan Zhang, and Yifan Zhu. The optimization of supply chain financing for bank green credit using stackelberg game theory in digital economy under internet of things. *Journal of Organizational and End User Computing (JOEUC)*, 35(3):1–16, 2023.

- Kai Zhang, Ivor W Tsang, and James T Kwok. Improved nyström low-rank approximation and error analysis. In *Proceedings of the 25th international conference on Machine learning*, pages 1232–1239, 2008.
- Shijun Zhang, Zuowei Shen, and Haizhao Yang. Deep network approximation: Achieving arbitrary accuracy with fixed number of neurons. *Journal of Machine Learning Research*, 23(276):1–60, 2022.
- Zijie Zheng, Lingyang Song, Zhu Han, Geoffrey Ye Li, and H Vincent Poor. A stackelberg game approach to proactive caching in large-scale mobile edge networks. *IEEE Transactions on Wireless Communications*, 17(8):5198–5211, 2018.