



Deep Statistical Solvers

Balthazar Donon, Zhengying Liu, Wenzhuo Liu, Isabelle Guyon, Antoine Marot, Marc Schoenauer

► To cite this version:

Balthazar Donon, Zhengying Liu, Wenzhuo Liu, Isabelle Guyon, Antoine Marot, et al.. Deep Statistical Solvers. NeurIPS, Dec 2020, Vancouver, Canada. hal-02974541v1

HAL Id: hal-02974541

<https://inria.hal.science/hal-02974541v1>

Submitted on 23 Oct 2020 (v1), last revised 2 Dec 2020 (v3)

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

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

Deep Statistical Solvers

Balthazar Donon

RTE R&D, INRIA, Université Paris-Saclay
balthazar.donon@rte-france.com

Zhengying Liu

Université Paris-Saclay, INRIA
zhengying.liu@inria.fr

Wenzhuo Liu

IRT SystemX
wenzhuo.liu@irt-systemx.fr

Isabelle Guyon

Université Paris-Saclay, INRIA, Chalearn
guyon@chalearn.org

Antoine Marot

RTE R&D
antoine.marot@rte-france.com

Marc Schoenauer

INRIA, Université Paris-Saclay
marc.schoenauer@inria.fr

Abstract

This paper introduces Deep Statistical Solvers (DSS), a new class of trainable solvers for optimization problems, arising from system simulations and others. The key idea is to learn a solver that generalizes to a given distribution of problem instances. This is achieved by directly using as loss the objective function of the problem, as opposed to most previous Machine Learning based approaches, which mimic the solutions attained by an existing solver. Though both types of approaches outperform classical solvers with respect to speed for a given accuracy, a distinctive advantage of DSS is that they can be trained without a training set of sample solutions. Focusing on use cases of systems of interacting and interchangeable entities (*e.g.* molecular dynamics, power systems, discretized PDEs), the proposed approach is instantiated within a class of Graph Neural Networks. Under sufficient conditions, we prove that the corresponding set of functions contains approximations to any arbitrary precision of the actual solution of the optimization problem. The proposed approach is experimentally validated on large linear problems, on demonstrating super-generalisation properties; And on AC power grid simulations, on which the predictions of the trained model have a correlation higher than 99.99% with the outputs of the classical Newton-Raphson method (known for its accuracy), while being 2 to 3 orders of magnitude faster.

1 Introduction

In many domains of physics and engineering, Deep Neural Networks (DNNs) have sped up simulations and optimizations by orders of magnitude, replacing some computational bricks based on first principles with data-driven numerical models – see *e.g.*, [1, 2, 3, 4]. However, in general, such data-driven approaches consist in training a *proxy* in a supervised way, to imitate solutions provided by some numerical solver. This is sometimes infeasible due to the high computational cost of existing simulators (*e.g.*, molecular dynamics, car crash simulations, computational fluid dynamics, and power grid simulation). Furthermore, such approaches ignore problem-specific considerations and may end up providing inconsistent solutions, failing to satisfy physical laws such as energy conservation (which can only be a posteriori checked, see *e.g.* [4]). In order to bypass this weakness, a growing body of work pushes towards an interplay between physics and Machine Learning [5, 6], *e.g.*, incorporating physical knowledge in the loss function during learning [7, 8].

Another important property of natural or artificial systems is that of invariance, a fundamental concept in science, allowing to generalize conclusions drawn from few observations, to whole invariance classes. This work focuses on permutation-invariant problems, which appear in simulations of complex systems of interacting and interchangeable entities [9] (e.g., molecular dynamics, power grids, simulations of partial differential equations (PDEs) with finite elements). Invariance has made its way in machine learning, as illustrated by the success of Convolutional Neural Networks (CNN) [10, 11], and of Graph Neural Networks (GNN) [12, 13]. In particular, implementations of GNNs successfully handle materials dynamics simulations [14], power systems [15], interacting particles [16] and classical [17] or quantum [18] chemistry. However, all of these works pertain to the *proxy approach* described above.

Our first contribution is to propose, at the interface of optimization and statistics, the Statistical Solver Problem (SSP), a novel formulation for learning to solve a whole class of optimization and system simulation problems. The resulting framework i) directly minimizes the global loss function of the problems during training, thus not requiring any existing solution of the problems at hand, and ii) directly incorporates permutation-invariance in the representation of the problems using a GNN-based architecture, called Deep Statistical Solver (DSS). Our second contribution is to prove that DSS satisfies some Universal Approximation property in the space of SSP solutions. The third contribution is an experimental validation of the approach.

The outline of the paper is the following. Section 2 sets the background, and defines SSPs. Section 3 introduces Deep Statistical Solvers. Section 4 proves the Universal Approximation property for permutation-invariant loss functions (and some additional hypotheses). Section 5 experimentally validates the DSS approach, demonstrating its efficiency w.r.t. state-of-the-art solvers, and unveiling some super-generalization capabilities. Section 6 concludes the paper.

2 Definitions and Problem Statement

This section introduces the context (notations and definitions) and the research goal of this work: The basic problem is, given a network of interacting entities (referred to later as Interaction Graph), to find a state of the network that minimizes a given loss function; From thereon, the main goal of this work is to learn a parameterized mapping that accurately and quickly computes such minimizing state for any Interaction Graph drawn from a given distribution.

2.1 Notations and Definitions

Notations Throughout this paper, for any $n \in \mathbb{N}$, $[n]$ denotes the set $\{1, \dots, n\}$; Σ_n is the set of permutations of $[n]$; for any $\sigma \in \Sigma_n$, any set Ω and any vector $\mathbf{x} = (x_i)_{i \in [n]} \in \Omega^n$, $\sigma \star \mathbf{x}$ is the vector $(x_{\sigma^{-1}(i)})_{i \in [n]}$; for any $\sigma \in \Sigma_n$ and any matrix $\mathbf{m} = (m_{ij})_{i,j \in [n]} \in \mathcal{M}_n(\Omega)$ (square matrices with elements in Ω), $\sigma \star \mathbf{m}$ is the matrix $(m_{\sigma^{-1}(i)\sigma^{-1}(j)})_{i,j \in [n]}$.

Interaction Graphs We call *Interaction Graph* a system of $n \in \mathbb{N}$ interacting entities, or *nodes*, defined as $\mathbf{G} = (n, \mathbf{A}, \mathbf{B})$, where n is the size of \mathbf{G} (number of nodes), $\mathbf{A} = (A_{ij})_{i,j \in [n]}$; $A_{ij} \in \mathbb{R}^{d_A}$; $d_A \geq 1$ represents the interactions between nodes, and $\mathbf{B} = (B_i)_{i \in [n]}$; $B_i \in \mathbb{R}^{d_B}$, $d_B \geq 1$ are some local external inputs at each node. Let \mathcal{G}_{d_A, d_B} be the set of all such Interaction Graphs and simply \mathcal{G} when there is no confusion. For any $\sigma \in \Sigma_n$ and any Interaction Graph $\mathbf{G} = (n, \mathbf{A}, \mathbf{B})$, $\sigma \star \mathbf{G}$ denotes the Interaction Graph $(n, \sigma \star \mathbf{A}, \sigma \star \mathbf{B})$.

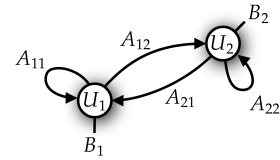


Figure 1: A sample Interaction Graph $(2, \mathbf{A}, \mathbf{B})$

Interaction Graphs can also be viewed as "doubly weighted" graphs, i.e., graphs with weights on both the edges (weights A_{ij}) and the nodes (weights B_i), considering that those weights are vectors. For a given \mathbf{G} , we will also consider the underlying undirected unweighted graph $\tilde{\mathbf{G}}$ for which links between nodes i and j exist iff either A_{ij} or A_{ji} is non-zero¹. We will use the notion of neighborhood induced by $\tilde{\mathbf{G}}$: $j \in \mathcal{N}(i; \mathbf{G})$ iff i and j are neighbors in $\tilde{\mathbf{G}}$ (and $\mathcal{N}^*(i; \mathbf{G})$ will denote $\mathcal{N}(i; \mathbf{G}) \setminus \{i\}$).

States and Loss Functions Vectors $\mathbf{U} = (U_i)_{i \in [n]}$; $U_i \in \mathbb{R}^{d_U}$, $d_U \geq 1$ represent *states* of Interaction Graphs of size n , where U_i is the state of node i . \mathcal{U}_{d_U} denotes the set of all such states \mathcal{U} when there

¹A more rigorous definition of the actual underlying graph structure is deferred to Appendix A

is no confusion). A *loss function* ℓ is a real-valued function defined on pairs (\mathbf{U}, \mathbf{G}) , where \mathbf{U} is a state of \mathbf{G} (i.e., of same size).

Permutation invariance and equivariance A loss function on Interaction Graph \mathbf{G} of size n is *permutation-invariant* if for any $\sigma \in \Sigma_n$, $\ell(\sigma \star \mathbf{U}, \sigma \star \mathbf{G}) = \ell(\mathbf{U}, \mathbf{G})$.

A function \mathcal{F} from \mathcal{G} to \mathcal{U} , mapping an Interaction Graph \mathbf{G} of size n on one of its possible states \mathbf{U} is *permutation-equivariant* if for any $\sigma \in \Sigma_n$, $\mathcal{F}(\sigma \star \mathbf{G}) = \sigma \star \mathcal{F}(\mathbf{G})$.

2.2 Problem Statement

The Optimization Problem In the remaining of the paper, ℓ is a loss function on Interaction Graphs $\mathbf{G} \in \mathcal{G}$ that is both continuous and permutation-invariant. The elementary question of this work is to solve the following optimization problem for a given Interaction Graph \mathbf{G} :

$$\mathbf{U}^*(\mathbf{G}) = \underset{\mathbf{U} \in \mathcal{U}}{\operatorname{argmin}} \ell(\mathbf{U}, \mathbf{G}) \quad (1)$$

The Statistical Learning Goal We are not interested in solving problem (1) for just ONE Interaction Graph, but in learning a parameterized *solver*, i.e., a mapping from \mathcal{G} to \mathcal{U} , which solves (1) for MANY Interaction Graphs, namely all Interaction Graphs \mathbf{G} sampled from a given distribution \mathcal{D} over \mathcal{G} . In particular, \mathcal{D} might cover Interaction Graphs of different sizes. Let us assume additionally that \mathcal{D} and ℓ are such that, for any $\mathbf{G} \in \operatorname{supp}(\mathcal{D})$ (the support of \mathcal{D}) there is a unique minimizer $\mathbf{U}^*(\mathbf{G}) \in \mathcal{U}$ of problem (1). The goal of the present work is to learn a single solver that best approximates the mapping $\mathbf{G} \mapsto \mathbf{U}^*(\mathbf{G})$ for all \mathbf{G} in $\operatorname{supp}(\mathcal{D})$. More precisely, assuming a family of solvers $Solver_\theta$ parameterized by $\theta \in \Theta$ (Section 3 will introduce such a parameterized family of solvers, based on Graph Neural Networks), the problem tackled in this paper can be formulated as a *Statistical Solver Problem* (SSP):

$$\text{SSP}(\mathcal{G}, \mathcal{D}, \mathcal{U}, \ell) \left\{ \begin{array}{l} \text{Given distribution } \mathcal{D} \text{ on space of Interaction Graphs } \mathcal{G}, \text{ space of states } \mathcal{U}, \\ \text{and loss function } \ell, \text{ solve } \theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} \mathbb{E}_{\mathbf{G} \sim \mathcal{D}} [\ell(Solver_\theta(\mathbf{G}), \mathbf{G})] \end{array} \right. \quad (2)$$

Learning phase In practice, the expectation in (2) will be empirically computed using a finite number of Interaction Graphs sampled from \mathcal{D} , by directly minimizing ℓ (i.e., without the need for any \mathbf{U}^* solution of (1)). The result of this empirical minimization is a parameter $\hat{\theta}$.

Inference The solver $Solver_{\hat{\theta}}$ can then be used, at inference time, to compute, for any $\mathbf{G} \in \operatorname{supp}(\mathcal{D})$, an approximation of the solution $\mathbf{U}^*(\mathbf{G})$

$$\hat{\mathbf{U}}(\mathbf{G}) = Solver_{\hat{\theta}}(\mathbf{G}) \quad (3)$$

Solving problem (1) has been replaced by a simple and fast inference of the learned model $Solver_{\hat{\theta}}$ (at the cost of a possibly expensive learning phase).

Discussion The SSP experimented with in Section 5.2 addresses the simulation of a Power Grid, a real-world problem for which the benefits of using the proposed approach becomes clear. Previous work [19] used a "proxy" approach, which consists in learning from known solutions of the problem, provided by a classical solver. The training phase is sketched on Figure 2.a. The drawback of such an approach is the need to gather a huge number of training examples (i.e., solutions of problem (1)), something that is practically infeasible for complex problems: either such solutions are too costly to obtain (e.g., in car crash simulations), or there is no provably optimal solution (e.g., in molecular dynamics simulations). In contrast, since the proposed approach directly trains $Solver_\theta$ by minimizing the loss ℓ (Figure 2.b), no such examples are needed.

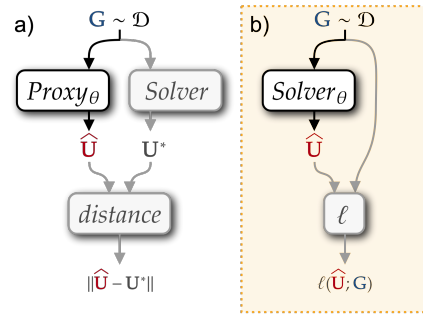


Figure 2: **Proxy approach (a) vs. DSS (b)**

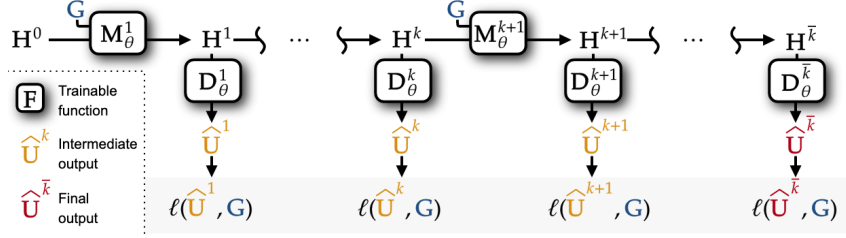


Figure 3: Graph Neural Network implementation of a DSS

3 Deep Statistical Solver Architecture

In this section, we introduce the class of Graph Neural Networks (GNNs) that will serve as DSSs. The intuition behind this choice comes from the following property (proof in Appendix B.2):

Property 1. *If the loss function ℓ is permutation-invariant and if for any $\mathbf{G} \in \text{supp}(\mathcal{D})$ there exists a unique minimizer $\mathbf{U}^*(\mathbf{G})$ of problem (1), then \mathbf{U}^* is permutation-equivariant.*

Graph Neural Networks, introduced in [20], and further developed in [21, 22] (see also the recent surveys [13, 23]), are a class of parameterized permutation-equivariant functions. Therefore, they seem to be good candidates to build SSP solutions, since Property 1 states that the ideal solver \mathbf{U}^* is permutation-equivariant (this will be confirmed by Corollary 1).

Overall architecture There are many possible implementations of GNNs, but whatever the chosen type, it is important to make room for information propagation throughout the whole network (see also Section 4). Hence the choice of an iterative process that acts on a latent state $\mathbf{H} \in \mathcal{U}_d$; $H_i \in \mathbb{R}^d$, $d \geq 1$ for \bar{k} iterations (d and \bar{k} are hyperparameters). For a node $i \in [n]$, the latent state H_i can be seen as an embedding of the actual state U_i .

The overall architecture is described in Figure 3. All latent states in \mathbf{H}^0 are initialized to a zero vector. The *message passing* step performs \bar{k} updates on the latent state variable \mathbf{H} using \mathbf{M}_{θ}^k , spreading information using interaction coefficients \mathbf{A} and external inputs \mathbf{B} of \mathbf{G} (eq. 5–8). After each update, latent state \mathbf{H}^k is decoded into a meaningful actual state $\hat{\mathbf{U}}^k$ (eq. 9). The last state $\hat{\mathbf{U}}^{\bar{k}}$ is the actual output of the algorithm $\hat{\mathbf{U}}$. However, in order to robustify learning, all intermediate states $\hat{\mathbf{U}}^k$ are taken into account in the training loss through a discounted sum with hyperparameter $\gamma \in [0, 1]$:

$$\text{Training Loss} = \sum_{k=1}^{\bar{k}} \gamma^{\bar{k}-k} \ell(\hat{\mathbf{U}}^k, \mathbf{G}) \quad (4)$$

Message passing \mathbf{M}_{θ}^k For each node i , three different messages are computed, $\phi_{\rightarrow, \theta}^k, \phi_{\leftarrow, \theta}^k, \phi_{\odot, \theta}^k$, corresponding to outgoing, ingoing and self-loop links, respectively using trainable mappings $\Phi_{\rightarrow, \theta}^k, \Phi_{\leftarrow, \theta}^k, \Phi_{\odot, \theta}^k$, as follows:

$$\phi_{\rightarrow, i}^k = \sum_{j \in \mathcal{N}^+(i; \mathbf{G})} \Phi_{\rightarrow, \theta}^k(H_i^{k-1}, A_{ij}, H_j^{k-1}) \quad \text{outgoing edges} \quad (5)$$

$$\phi_{\leftarrow, i}^k = \sum_{j \in \mathcal{N}^-(i; \mathbf{G})} \Phi_{\leftarrow, \theta}^k(H_i^{k-1}, A_{ji}, H_j^{k-1}) \quad \text{ingoing edges} \quad (6)$$

$$\phi_{\odot, i}^k = \Phi_{\odot, \theta}^k(H_i^{k-1}, A_{ii}) \quad \text{self loop} \quad (7)$$

Latent states H_i^k are then computed using trainable mapping Ψ_{θ}^k , in a ResNet-like fashion:

$$\mathbf{H}^k = \mathbf{M}_{\theta}^k(\mathbf{H}^{k-1}, \mathbf{G}) := (H_i^k)_{i \in [n]}, \text{ with } H_i^k = H_i^{k-1} + \Psi_{\theta}^k(H_i^{k-1}, B_i, \phi_{\rightarrow, i}^k, \phi_{\leftarrow, i}^k, \phi_{\odot, i}^k) \quad (8)$$

Decoding The decoding step applies the same trainable mapping Ξ_{θ}^k to every node:

$$\hat{\mathbf{U}}^k = \mathbf{D}_{\theta}^k(\mathbf{H}^k) = (\Xi_{\theta}^k(H_i^k))_{i \in [n]} \quad (9)$$

Training All trainable blocks $\Phi_{\rightarrow, \theta}^k, \Phi_{\leftarrow, \theta}^k, \Phi_{\odot, \theta}^k$ and Ψ_{θ}^k for the message passing phase, and Ξ_{θ}^k for the decoding phase, are implemented as Neural Networks. They are all trained simultaneously, backpropagating the gradient of the training loss of eq. (4) (see details in Section 5).

Number of propagation steps Our current implementation choice is to consider different neural network blocks at each propagation step. The underlying intuition is that the nature of information exchange does not have to be the same at the beginning and at the end of the process. This comes at the expense of a fixed amount of propagation steps \bar{k} . However, future work will include the investigation of a Recurrent Graph Neural Network architecture, drawing inspiration from [24]. This would allow for an adaptive number of steps \bar{k} , and a much lighter model.

Inference Complexity Assuming that each neural network block has a single hidden layer with dimension d , that $d \geq d_A, d_B, d_U$, and denoting by m the average neighborhood size, one inference has computational complexity of order $\mathcal{O}(mn\bar{k}d^3)$, scaling linearly with n . Furthermore, many problems involve very local interactions, resulting in small m . However, one should keep in mind that hyperparameters \bar{k} and d should be chosen according to the characteristics of distribution \mathcal{D} . If we can compute the maximal diameter of any $\mathbf{G} \in \text{supp}(\mathcal{D})$ (e.g., if \mathcal{D} is a database of the history of the Californian power grid), one should choose a larger value for \bar{k} (see Corollary 1). Similarly, if one is working with data that have very large d_A and d_B , one may want to choose a sufficiently large value for d to let information flow properly.

Equivariance The proposed architecture defines permutation-equivariant DSS (see Appendix B.1).

4 Deep Statistical Solvers are Universal Approximators for SSPs Solutions

This Section proves, heavily relying on work by [25], a Universal Approximation Theorem for the class of DSSs with Lipschitz activation function (e.g. ReLU) in the space of the solutions of SSPs. The space of Interaction Graphs is a metric space for the distance

$$d(\mathbf{G}, \mathbf{G}') = \|\mathbf{A} - \mathbf{A}'\| + \|\mathbf{B} - \mathbf{B}'\| \text{ if } n = n' \text{ and } +\infty, \text{ otherwise}$$

Universal Approximation Property Given metric spaces \mathcal{X} and \mathcal{Y} , a set of continuous functions $\mathcal{H} \subset \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$ is said to satisfy the *Universal Approximation Property* (UAP) if it is dense in the space of all continuous functions $\mathcal{C}(\mathcal{X}, \mathcal{Y})$ (with respect to the uniform metric).

Denote by $\mathcal{H}_{d_{in}}^{d_{out}}$ a set of neural networks from $\mathbb{R}^{d_{in}}$ to $\mathbb{R}^{d_{out}}$, for which the UAP holds. It is known since [26] that the set of neural networks with at least one hidden layer, an arbitrarily large amount of hidden neurons, and an appropriate activation function, satisfies these conditions.

Hypothesis space Let $\bar{k} \in \mathbb{N}$. We denote by $\mathcal{H}^{\bar{k}}$ the set of graph neural networks defined in Section 3 such that $\bar{k} \leq \bar{k}$, $d \in \mathbb{N}$ and for any $k = 1, \dots, \bar{k}$, we consider all possible $\Phi_{\rightarrow, \theta}^k, \Phi_{\leftarrow, \theta}^k \in \mathcal{H}_{d_A+2d}^d$, $\Phi_{\circ, \theta}^k \in \mathcal{H}_{d_A+d}^d$, $\Psi_{\theta}^k \in \mathcal{H}_{d_B+4d}^d$ and $\Xi_{\theta}^k \in \mathcal{H}_d^{d_U}$.

Diameter of an Interaction Graph Let $\mathbf{G} = (n, \mathbf{A}, \mathbf{B}) \in \mathcal{G}$, and let $\tilde{\mathbf{G}}$ be its undirected and unweighted graph structure, as defined in Section 2.1. We will write $\text{diam}(\mathbf{G})$ for $\text{diam}(\tilde{\mathbf{G}})$, the diameter of $\tilde{\mathbf{G}}$ [27].

Hypotheses over distribution \mathcal{D} We introduce the four following hypotheses over $\text{supp}(\mathcal{D})$:

- *Permutation-invariance.* For any $\mathbf{G} \in \text{supp}(\mathcal{D})$ and $\sigma \in \Sigma_n$, $\sigma \star \mathbf{G} \in \text{supp}(\mathcal{D})$;
- *Compactness.* $\text{supp}(\mathcal{D})$ is a compact subset of \mathcal{G} ;
- *Connectivity.* For any $\mathbf{G} \in \text{supp}(\mathcal{D})$, $\tilde{\mathbf{G}}$ has only one connected component;
- *Separability of external inputs.* There exist $\delta > 0$ such that for any $\mathbf{G} = (n, \mathbf{A}, \mathbf{B}) \in \text{supp}(\mathcal{D})$ and any $i \neq j \in [n]$, $\|B_i - B_j\| \geq \delta$.

The *compactness* implies that there is an upper bound \bar{n} over the size n of Interaction Graphs in $\text{supp}(\mathcal{D})$. Also, these hypotheses imply that there is a finite upper bound on the diameters of all \mathbf{G} s. In the following, Δ will denote such upper bound. We denote by $\mathcal{C}_{eq}(\text{supp}(\mathcal{D}))$ the set of continuous and permutation-equivariant functions over $\text{supp}(\mathcal{D})$.

Theorem 1. *Let \mathcal{D} be a distribution over \mathcal{G} for which the above hypotheses hold.*

Then if $\bar{k} \geq \Delta + 2$, $\mathcal{H}^{\bar{k}}$ is dense in $\mathcal{C}_{eq}(\text{supp}(\mathcal{D}))$.

Sketch of the proof (see Appendix B.3 for all details) Still following [25], we first prove a modified version of the *Stone-Weierstrass theorem for equivariant functions*. This theorem guarantees that a certain subalgebra of functions is dense in the set of continuous and permutation-equivariant

functions if it separates non-isomorphic Interaction Graphs. Following the idea of [26], we extend the hypothesis space to ensure closure under addition and multiplication. We then prove that the initial hypothesis space is dense in this new subalgebra. Finally, we conclude the proof by showing that the separability property mentioned above is satisfied by this newly-defined subalgebra.

Corollary 1. *Let \mathcal{D} be a distribution over \mathcal{G} for which the above hypotheses hold. Let ℓ be a continuous and permutation-invariant loss function such that for any $\mathbf{G} \in \text{supp}(\mathcal{D})$, problem (1) has a unique minimizer $\mathbf{U}^*(\mathbf{G})$, continuous w.r.t \mathbf{G} . Then $\forall \epsilon > 0$, $\exists \text{Solver}_\theta \in \mathcal{H}^{\Delta+2}$, such that*

$$\forall \mathbf{G} \in \text{supp}(\mathcal{D}), \|\text{Solver}_\theta(\mathbf{G}) - \mathbf{U}^*(\mathbf{G})\| \leq \epsilon$$

This corollary is an immediate consequence of Theorem 1 and ensures that there exists a DSS using at most $\Delta + 2$ propagation updates that approximates with an arbitrary precision for all $\mathbf{G} \in \text{supp}(\mathcal{D})$ the actual solution of problem (1). This is particularly relevant when considering large Interaction Graphs that have small diameters.

Discussion: This universal approximation theorem does not offer any guarantee of convergence toward the ideal solver \mathbf{U}^* – but there hardly exist such convergence guarantees in the field of Deep Learning. However, this non-trivial result provides a solid theoretical ground to the proposed approach by proving its consistency.

5 Experiments

This section investigates the behavior and performances of DSSs on two SSPs. The first one amounts to solving linear systems, though the distribution of problems is generated from a discretized Poisson PDE. The second is the (non-quadratic) AC power flow computation. With respect to the hypotheses of the theoretical results in Section 4, the continuity and the permutation invariance conditions are satisfied in both cases, while the uniqueness can only be proven for the linear system. However it is very likely to hold for many problems.

In all cases, the dataset is split into training/validation/test sets. All free hyperparameters² are tuned by trial and errors using the validation set, and **all results presented are results on the test set**. We also compare the DSS to the proxy approach: the architecture is strictly the same, but the loss function used during training is the distance to the “ground truth” (provided by the LU or Newton-Raphson methods). Training is performed using the Adam optimizer [28] with the standard hyperparameters of TensorFlow 1.14 [29], running on an Nvidia GeForce RTX 2080 Ti. Gradient clipping is used to avoid exploding gradient issues. **In the following, all experiments were repeated three times, with the same datasets and different random seeds** (as reported in Tables 1 and 2). In all experiment and for both the DSS and the proxy approaches, we only report the results of the worst of the three trained models. Our code is in the supplementary materials³, and links to the datasets are in references.

The main metrics for our experiments are the Pearson correlation and the normalised RMSE (NRMSE) with the output of the classical optimization method (*i.e.* LU in the linear case and Newton-Raphson in the AC Power Flow problems). The NRMSE is computed by dividing the RMSE by the difference of the highest and the lowest values (dividing by the mean for data centered around zero would make no sense). The value of the loss function ℓ is computed over the whole test set: the 10th and 90th percentiles as well as the median are reported.

5.1 Solving Linear Systems from a Discretized PDE

Problem, and goals of experiments The example SSP considered here comes from the Finite Element Method applied to solve the 2D Poisson equation, one of the simplest and most studied PDE in applied mathematics: the geometry of the domain of the equation is discretized into an unstructured mesh, and computing the vector \mathbf{U} of solution values at each node of the mesh amounts to solving a linear system $\mathbf{A}\mathbf{U} = \mathbf{B}$ obtained by assembling local equations [30]. \mathbf{A} and \mathbf{B} encode both the geometry of the problem and the boundary conditions.

For illustration purposes, the Poisson equation can be used to model a field of temperature. In Figure 4, the geometry (house profile) is shown in the Top Left. The result of the optimization is the field of

²Ranges of tested hyperparameters : $d \in [5, 20]$, $\bar{k} \in [10, 40]$, hidden layers $\in [1, 2]$, $\alpha \in [1\text{e-}1, 1\text{e-}4]$ (see Appendix E), non linearity $\in \{\tanh, \text{leaky_ReLU}\}$, $lr \in [1\text{e-}1, 1\text{e-}4]$, $\gamma \in \{0, 0.5, 0.9, 1\}$.

³code also available at <https://github.com/bdonon/DeepStatisticalSolvers>

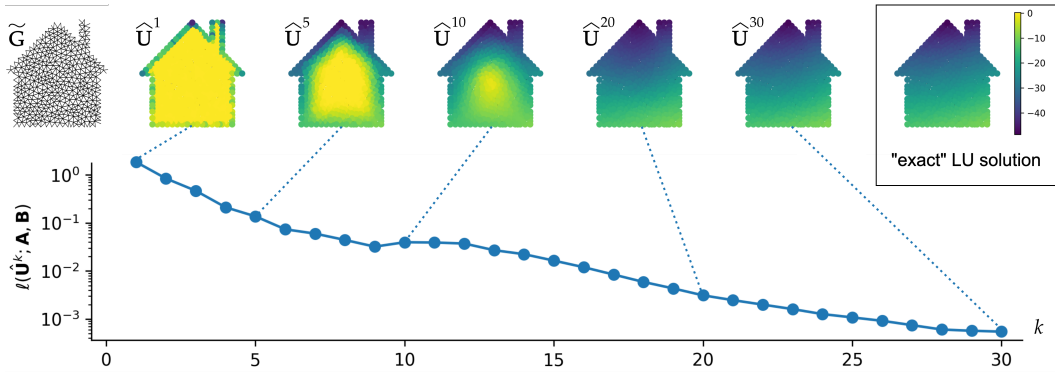


Figure 4: **Intermediate losses and predictions** - Top left: the structure graph \tilde{G} (the mesh); Top right: the LU solution; Bottom: evolution of the loss along the $\bar{k} = 30$ updates for a trained DSS, at inference time. The intermediate predictions \hat{U}^k are displayed for several values of k .

Method	DSS	Proxy	LU	BGS (tol=1e-3)
Correlation w/ LU	> 99.99%	> 99.99%	-	-
NRMSE w/ LU	1.6e-3	1.1e-3	-	-
Time per instance (ms)	1.8	1.8	2.4	2.3
Loss 10 th percentile	3.9e-4	7.0e-3	4.5e-27	1.3e-3
Loss 50 th percentile	1.2e-3	1.6e-2	6.1e-26	1.7e-2
Loss 90 th percentile	4.1e-3	4.0e-2	6.3e-25	1.1e-1

Table 1: **Solving specific linear systems** – for similar accuracy, DSS is faster than the iterative BGS, while highly correlated with the "exact" solution as given by LU.

temperature everywhere in the house (shown in the Top Right).

This problem is easily set as an SSP in which each node i corresponds to a node of the mesh, all parameters are scalars ($d_A = d_B = d_U = 1$), and the loss function is the following:

$$\ell(\mathbf{U}, \mathbf{G}) = \sum_{i \in [n]} \left(\sum_{j \in [n]} A_{ij} U_j - B_i \right)^2 \quad (10)$$

It is clearly permutation-invariant and satisfies both the unicity of the solution and the continuity conditions evoked in Corollary 1. Our goal here is of course not to solve the Poisson equation, nor is it to propose a new competitive method to invert linear systems. As a matter of fact, the proposed approach does not make use of the linearity of the problem. Our goal is actually twofold: i) validate the DSS approach in high dimension ($n \approx 500$ nodes), and ii) analyze how DSS learns the distribution \mathcal{D} . Here, the distribution \mathcal{D} is defined by the specific structure of linear systems that result from the discretization of the Poisson equation. In particular, we will carefully study the generalization capability of the learned model in terms of problem size, for similar problem structures.

Experimental conditions The dataset [31] consists of 96180/32060/32060 training/validation/test examples from the distribution generated from the discretization of the Poisson equation: randomly generated 2D geometries and random values for the second-hand function f and boundary condition g are used to compute the \mathbf{A} s and \mathbf{B} s. Their number of nodes n are around 500 (max 599) (automatic mesh generators do not allow a precise control of n).

The number of updates \bar{k} is set to 30 (average diameter size for the considered meshes). Each NN block has one hidden layer of dimension $d = 10$ and a leaky-ReLU non linearity; we have $\alpha = 1e-3$, $\text{lr} = 1e-2$ and $\gamma = 0.9$. The complete DSS has 49,830 weights. Training is done for 280,000 iterations (48h) with batch size 100.

Two baseline methods are considered [32], the direct LU decomposition, that could be considered giving the "exact" solution for these sizes of matrices, and the iterative Biconjugate Gradient Stabilized methods (BGS), with stopping tolerances of 10^{-3} . These algorithms are run on an Intel Xeon Silver 4108 CPU (1.80GHz) (GPU implementations were not available, they could decrease LU computational cost by a factor 6 [33]).

Results Table 1 displays comparisons between a trained DSS and the baselines. First, these results validate the approach, demonstrating that DSS can learn to solve 500 dimensional problems rather accurately, and in line with the "exact" solutions as provided by the direct method LU (99.99% correlation). Second, DSS is slightly but consistently faster than the iterative method BGS for similar accuracy (a tunable parameter of BGS). Further work will explore how DSS scales up in much higher

dimensions, in particular when LU becomes intractable. We observe similar results for the proxy approach. Figure 4 illustrates, on a hand-made test example (the mesh is on the upper left corner), how the trained DSS updates its predictions, at inference time, along the \bar{k} updates. The flow of information from the boundary to the center of the geometry is clearly visible.

But what did exactly the DSS learn? Next experiments are concerned with the super-generalization capability of DSSs, looking at their results on test examples sampled from distributions departing from the one used for learning.

Super-Generalization We now experimentally analyze how well a trained model is able to generalize to a distribution \mathcal{D} that is different from the training distribution. The same data generation process that was used to generate the training dataset (see above) is now used with meshes of very different sizes, everything else being equal. Whereas the training distribution only contains Interaction Graphs of sizes around 500, out-of-distribution test examples have sizes from 100 and 250 (left of Figure 5) up to 750 and 1000 (right of Figure 5). In all cases, the trained model is able to achieve a correlation with the "true" LU solution as high as 99.99%. Interestingly, the trained DSS achieves a higher correlation with the LU solutions for data points with a lower number of nodes, while the correlation of the proxy model decreases when n both increases and decreases. Further experiments with even larger sizes are needed to reach the upper limit of such a super generalization. Nevertheless, thanks to the specific structure dictated to the linear system by the Poisson equation, DSS was able to perform some kind of zero-shot learning for problems of very different sizes.

Other experiments (see Appendix C) were performed by adding noise to \mathbf{A} and \mathbf{B} . The performance of the trained model remains good for small noise, then smoothly degrades as the noise increases.

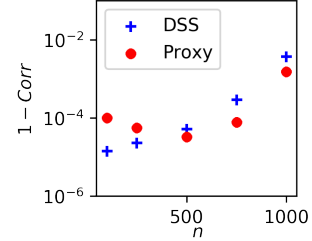


Figure 5: **Varying problem size n : Correlation (DSS, LU)**

5.2 AC power flow experiments

Problem and goals of experiments The second SSP example is the AC power flow prediction. The goal is to compute the steady-state electrical flows in a Power Grid, an essential part of real-time operations. Knowing the amount of power that is being produced and consumed throughout the grid (encoded into \mathbf{B} , and assumed to be consistent, i.e., production equates consumption), and the way power lines are interconnected, as well as their physical properties (encoded into \mathbf{A}), the goal is to compute the voltage defined at each electrical node $V_i = |V_i|e^{j\theta_i}$ (j denotes the imaginary unit), which we encode in the states \mathbf{U} . Kirchhoff's law (energy conservation at every node) governs this system, and its violation is directly used as loss function ℓ . Moreover, some constraints over the states \mathbf{U} are here relaxed and included as an additional term of the loss (with factor λ). One should also keep in mind that the main goal is to predict power flows, and not the voltages per se: Both aspects will be taken into account by measuring the correlation w.r.t $|V_i|$, θ_i , P_{ij} (real part of power flow) and Q_{ij} (imaginary part). This problem is highly non-linear, and a substantial overview is provided in [34]. This set of complex equations can be converted into a SSP using \mathbf{A} , \mathbf{B} and \mathbf{U} as defined above ($d_A = 2$, $d_B = 5$, $d_U = 2$), and loss function ℓ :

$$\begin{aligned} \ell(\mathbf{U}, \mathbf{G}) = & \sum_{i \in [n]} (1 - B_i^5) (-B_i^1 + U_i^1 \sum_{j \in [n]} A_{ij}^1 U_j^1 \cos(U_i^2 - U_j^2 - A_{ij}^2))^2 \\ & + \sum_{i \in [n]} B_i^3 (-B_i^2 + U_i^1 \sum_{j \in [n]} A_{ij}^1 U_j^1 \sin(U_i^2 - U_j^2 - A_{ij}^2))^2 + \lambda \sum_{i \in [n]} (1 - B_i^3) (U_i^1 - B_i^4)^2 \end{aligned} \quad (11)$$

More details about the conversion from classical power systems notations to this set of variables is provided in Appendix D. This loss is not quadratic, as demonstrated by the presence of sinusoidal terms. One can notice the use of binary variables B_i^3 and B_i^5 . Since both \mathbf{A} and \mathbf{B} vary across the dataset, the problem is largely non-linear with regards to the inputs.

Experimental conditions Experiments are conducted on two standard benchmarks from the *Learning to Run a Power Network* competition [35]: *IEEE case14* [36] ($n = 14$), and *IEEE case118* [37] ($n = 118$). In order to increase the diversity in terms of grid topology (encoded in \mathbf{A}), for each example, one (resp. two) randomly chosen power lines are disconnected with probability 25% (resp. 25%). For case14 (resp. case118), the dataset is split into 16064/2008/2008 (resp. 18432/2304/2304).

Each NN block has a single hidden layer of dimension $d = 10$ and a leaky-ReLU non linearity. For case14 (resp. case118), \bar{k} was set to 10 (resp. 30); we have $\alpha = 1e-2$, $lr = 3e-3$ and $\gamma = 0.9$ (resp. $\alpha = 3e-4$, $lr = 3e-3$ and $\gamma = 0.9$). The number of weights is 1,722 for each of the \bar{k} (\mathbf{M} , \mathbf{D}) blocks,

Dataset		IEEE 14 nodes			IEEE 118 nodes		
Method		DSS	Proxy	NR	DSS	Proxy	NR
Corr. w/ NR	$ V_i $	99.93%	> 99.99%	-	99.79%	> 99.99%	-
	θ_i	99.86%	> 99.99%	-	81.31%	> 99.99%	-
	P_{ij}	> 99.99%	> 99.99%	-	> 99.99%	> 99.99%	-
	Q_{ij}	> 99.99%	> 99.99%	-	> 99.99%	> 99.99%	-
NRMSE w/ NR	$ V_i $	2.0e-3	4.9e-4	-	1.4e-3	1.2e-3	-
	θ_i	7.1e-3	1.7e-3	-	5.7e-2	4.5e-3	-
	P_{ij}	6.2e-4	2.6e-4	-	1.0e-3	3.9e-4	-
	Q_{ij}	4.2e-4	2.0e-4	-	1.1e-4	1.7e-4	-
Inference time (ms)		1e-2	1e-2	2e1	9e-2	9e-2	2e1
Loss 10 th percentile		4.2e-6	2.3e-5	1.4e-12	1.3e-6	6.2e-6	2.9e-14
Loss 50 th percentile		1.0e-5	4.0e-5	2.1e-12	1.7e-6	8.3e-6	4.2e-14
Loss 90 th percentile		4.4e-5	1.2e-4	3.3e-12	2.5e-6	1.3e-5	6.4e-14

Table 2: **Solving specific AC power flow**— our trained DSS models are highly correlated with the Newton-Raphson solutions, while being 2 to 3 orders of magnitude faster.

hence 17, 220 (resp. 51, 660) in total. Training is done for 883, 000 (resp. 253000) iterations with batch size 1, 000 (resp. 500), and lasted 48h.

State-of-the-art AC power flow computation uses the Newton-Raphson method, used as baseline here ([38] implementation, on an Intel i5 dual-core (2.3GHz)). To the best of our knowledge, no GPU implementation was available, although recent work [39, 40] investigate such an avenue.

Results In both cases, correlations between power flows output by the trained DSSs and the Newton-Raphson method are above 99.99% (both real P_{ij} and imaginary Q_{ij}). The same can be said of the proxy models. However, one can observe a less satisfactory correlation in terms of $|V_i|$ and θ_i for the DSSs while the proxies maintain a correlation higher than 99.99%. This can be explained by the fact that the DSSs minimizes power mismatches while the proxies minimize the distance to the Newton-Raphson output in terms of $|V_i|$ and θ_i . However, this does not impact the quality of the power flow prediction, and one should keep in mind that the DSSs learn without any labels, contrarily to the proxies. Table 2 shows the huge acceleration of DSS (by two orders of magnitude) over Newton-Raphson, at the cost of an important decrease in accuracy, although both methods output very similar power flows (correlation higher than 99.99%).

6 Conclusions and Future Work

This paper proposes a novel paradigm that blends statistics and optimization, Statistical Solver Problems. In the SSP framework, a single solver is trained to solve a family of problem instances sampled from a given distribution of optimization problems, possibly arising from system simulations. Such training is performed by directly minimizing the loss of the optimization problems at hand. In particular, no existing solutions (obtained from costly simulations) are needed for training. The DSSs proposed in this paper, as a particular embodiment of the new proposed framework, is a class of Graph Neural Network, well suited to solving SSPs for which the loss function is permutation-invariant, and for which we theoretically prove some universal approximation properties.

The effectiveness of DSSs are experimentally demonstrated on two problems. Even though experiments on more complex problems are required, the proposed approach shows a good compromise between accuracy and speed in dimensions up to 500 on these two sample problems: solving linear systems, and the non-linear AC power flow. The accuracy on power flow computations matches that of state-of-the-art approaches while speeding up calculations by 2 to 3 orders of magnitude. Our DSS method could also be used as an initialization heuristic for classical optimization algorithms.

Future work will focus on incorporating discrete variables in the state space, and integrating constraints by casting them into a bilevel optimization problem (using two successive DSS that are trained jointly). Other avenues for research concern further theoretical improvements to investigate convergence properties of the DSS approach, in comparison to other solvers, as well as investigations on the limitations of the approach.

Broader Impact

This work introduces an original approach to solving permutation-invariant problems defined on a graph. The proposed approach is agnostic w.r.t. the practical problem it is applied to. As such, no direct poor societal consequences of this work are to be feared. However, and this is an issue that goes beyond this particular work, it can be applied to critical industrial problems, as demonstrated with the power grid experiments we use to illustrate and validate the approach in Section 5.2: in such context, it is important to ensure by design (*i.e.*, in the definition of the search space and the objective function) that the proper constraints are applied to avoid detrimental solutions. Being able to validate the obtained solution is a problem-specific issue, but validating the whole approach is the holy grail of such work, and is by now out of reach.

Funding Disclosure

The first author is partly funded by the ANR CIFRE contract 2018/0386.

Acknowledgement

The authors would like to thank Rémy Clément, Laure Crochepierre, Cédric Josz and the reviewers for their careful reviews and insightful suggestions; and Victor Berger for the fruitful discussions.

References

- [1] J. Tompson, K. Schlachter, P. Sprechmann, and K. Perlin. Accelerating eulerian fluid simulation with convolutional networks. *ArXiv: 1607.03597*, 2016.
- [2] M. F. Kasim, D. Watson-Parris, L. Deaconu, S. Oliver, P. Hatfield, D. H. Froula, G. Gregori, M. Jarvis, S. Khatiwala, J. Korenaga, J. Topp-Muggleston, E. Viezzer, and S. M. Vinko. Up to two billion times acceleration of scientific simulations with deep neural architecture search, 2020.
- [3] T. T. Nguyen. Neural network load-flow. *IEEE Proceedings - Generation, Transmission and Distribution*, 142:51–58(7), 1995.
- [4] S. Rasp, M. S. Pritchard, and P. Gentine. Deep learning to represent sub-grid processes in climate models. *PNAS*, 2018.
- [5] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová. Machine learning and the physical sciences. *Reviews of Modern Physics*, 91(4), Dec 2019.
- [6] J. Sirignano and K. Spiliopoulos. Dgm: A deep learning algorithm for solving partial differential equations. *Journal of Computational Physics*, 375:1339–1364, Dec 2018.
- [7] M. Raissi, P. Perdikaris, and G. E. Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- [8] L. Lu, X. Meng, Z. Mao, and G. E. Karniadakis. Deepxde: A deep learning library for solving differential equations. *ArXiv: 1907.04502*, 2019.
- [9] V. Vemuri. Modeling of complex systems: An introduction. *New York: Academic Press.*, 1978.
- [10] Y. LeCun, P. Haffner, L. Bottou, and Y. Bengio. Object recognition with gradient-based learning. In *Shape, Contour and Grouping in Computer Vision*, pages 319–. Springer-Verlag, 1999.
- [11] M. T. McCann, K. H. Jin, and M. Unser. Convolutional neural networks for inverse problems in imaging: A review. *IEEE Signal Processing Magazine*, 34(6):85–95, 2017.
- [12] P. W. Battaglia and al. Relational inductive biases, deep learning, and graph networks. *ArXiv: 1806.01261*, 2018.
- [13] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, page 1–21, 2020.
- [14] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. W. Battaglia. Learning to simulate complex physics with graph networks, 2020.
- [15] V. Bolz, J. Rueß, and A. Zell. Power flow approximation based on graph convolutional networks. In *18th IEEE International Conference On Machine Learning And Applications (ICMLA)*, pages 1679–1686, 2019.
- [16] T. N. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. S. Zemel. Neural relational inference for interacting systems. In *Proceedings of the 35th International Conference on Machine Learning, ICML, 2018*.
- [17] C. Chi, Y. Weike, Z. Yunxing, Z. Chen, and O. S. Ping. Graph networks as a universal machine learning framework for molecules and crystals. *Chemistry of Materials*, 31(9):3564–3572, 2019.
- [18] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. *ArXiv: 1704.01212*, 2017.
- [19] B. Donon, B. Donnot, I. Guyon, and A. Marot. Graph neural solver for power systems. In *International Joint Conference on Neural Networks (IJCNN)*, 2019.

- [20] M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In *IEEE International Joint Conference on Neural Networks Proceedings*, volume 2, pages 729–734, 2005.
- [21] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- [22] C. Gallicchio and A. Micheli. Graph echo state networks. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2010.
- [23] J. Zhou, G. Cui, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun. Graph neural networks: A review of methods and applications. *ArXiv: 1812.08434*, 2018.
- [24] A. Sanchez-Gonzalez, N. Heess, J. Tobias Springenberg, J. Merel, M. A. Riedmiller, R. Hadsell, and P. W. Battaglia. Graph networks as learnable physics engines for inference and control. *CoRR*, abs/1806.01242, 2018.
- [25] N. Keriven and G. Peyré. Universal invariant and equivariant graph neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 7092–7101. Curran Associates, Inc., 2019.
- [26] K. Hornik, M. B. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- [27] M. Newman. *Networks: An Introduction*. Oxford University Press, Inc., USA, 2010.
- [28] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization, 2014.
- [29] M. Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [30] P. G. Ciarlet. *The Finite Element Method for Elliptic Problems*. Mathematics and its Applications. 1978.
- [31] Authors. Dataset of linear systems built from the poisson equation. <http://doi.org/10.5281/zenodo.4024811>, 2020.
- [32] G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, third edition, 1996.
- [33] N. Galoppo, N. K. Govindaraju, M. Henson, and D. Manocha. LU-GPU: Efficient Algorithms for Solving Dense Linear Systems on Graphics Hardware. In *ACM/IEEE SC 2005 Conference*, pages 3–3. IEEE, 2005.
- [34] P. S. Kundur. Power system stability. In *Power System Stability and Control, Third Edition*, pages 1–12. CRC Press, 2012.
- [35] A. Marot, B. Donnot, C. Romero, L. Veyrin-Forrer, M. Lerousseau, B. Donon, and I. Guyon. Learning to run a power network challenge for training topology controllers. *ArXiv: 1912.04211*, 2019.
- [36] Authors. Dataset of ac power flow systems on the ieee case14 power grid. <http://doi.org/10.5281/zenodo.4024866>, 2020.
- [37] Authors. Dataset of ac power flow systems on the ieee case118 power grid. <http://doi.org/10.5281/zenodo.4024875>, 2020.
- [38] L. Thurner, A. Scheidler, F. Schafer, J. H. Menke, J. Dollichon, F. Meier, S. Meinecke, and M. Braun. pandapower - an open source python tool for convenient modeling, analysis and optimization of electric power systems. *IEEE Transactions on Power Systems*, 2018.
- [39] X. Su, C. He, T. Liu, and L. Wu. Full parallel power flow solution: A gpu-cpu-based vectorization parallelization and sparse techniques for newton–raphson implementation. *IEEE Transactions on Smart Grid*, 11(3):1833–1844, 2020.
- [40] D.-H. Yoon and Y. Han. Parallel power flow computation trends and applications: A review focusing on gpu. *Energies*, 13:2147, 05 2020.

A Underlying Graph Structure

We generalize the standard notion of neighborhood to the setting of Interaction Graphs and SSP $(\mathcal{G}, \mathcal{D}, \mathcal{U}, \ell)$. The intuitive way of defining neighbors of a node i is to look for nodes j such that $A_{ij} \neq 0$ or $A_{ji} \neq 0$. However this intuitive definition does not perfectly suit the case of SSPs as the properties of the loss function ℓ do impact the interactions between nodes.

For instance, let the loss function ℓ be defined by $\ell(\mathbf{U}, \mathbf{G}) = \sum_{i \in [n]} f(U_i)$, for some real-valued function f . In this case, there is no interaction between nodes when computing the state of the Interaction Graph, even though some coefficient A_{ij} may be non zero. We note in this case that:

$$\forall i \neq j, \forall \mathbf{U} \in \mathcal{U}, \frac{\partial^2 \ell}{\partial U_i \partial U_j}(\mathbf{U}, \mathbf{G}) = 0 \quad (12)$$

We thus propose the following definition of the neighborhood of a node i with respect to Interaction Graph \mathbf{G} and loss function ℓ :

$$\mathcal{N}(i; \mathbf{G}, \ell) = \left\{ j \in [n] \mid \exists \mathbf{U}, \frac{\partial^2 \ell}{\partial U_i \partial U_j}(\mathbf{U}, \mathbf{G}) \neq 0 \right\} \quad (13)$$

For a given class of SSP, the loss function ℓ does not change, so it will be omitted in the following, and we will write $\mathcal{N}(i; \mathbf{G})$.

One can observe that in the case of a quadratic optimization problem where $d_A = d_B = d_U = 1$ and $\ell(\mathbf{U}, \mathbf{G}) = \mathbf{U}^T \mathbf{A} \mathbf{U} + \mathbf{B}^T \mathbf{U}$, this notion of neighborhood is exactly that given in Section 2.1, and $\tilde{\mathbf{G}}$ is indeed the undirected graph defined by the non-zero entries of \mathbf{A} (or more precisely those of $(\mathbf{A} + \mathbf{A}^T)/2$ when \mathbf{A} is not symmetric).

B Mathematical proofs

In this section, we'll follow Keriven and Peyré [25] and use the notation $[\mathbf{G}]_i$ to denote the i^{th} component of any Interaction Graph or hyper-graph G . In the following, 'dense' means 'dense with respect to the uniform metric' by default. As a reminder, the uniform metric \bar{d} on function spaces given two metric spaces (X, d_X) and (Y, d_Y) is defined by

$$\bar{d}(f, f') := \sup_{x \in X} d_Y(f(x), f'(x)). \quad (14)$$

B.1 Proof of equivariance of the proposed DSS architecture

The following is a proof of the equivariance of the architecture proposed in Section 3.

Proof. Because the loss function ℓ is permutation invariant, we only have to prove that eq. (8)-(9) satisfy the permutation-equivariance property.

Let us prove by induction on k that \mathbf{H}^k is permutation-equivariant (by a slight abuse of notation in eq. (8), we consider the latent states \mathbf{H}^k as functions of \mathbf{G}), *i.e.* that $\mathbf{H}^k(\sigma \star \mathbf{G}) = \sigma \star \mathbf{H}^k(\mathbf{G})$.

For $k = 0$, it is clear that $\sigma \star \mathbf{H}^0 = (0, \dots, 0)_{i \in [n]} = \mathbf{H}^0$, which is independent of \mathbf{G} .

Now suppose the equivariance property holds for \mathbf{H}^{k-1} , then from eq. (5) comes

$$[\phi_{\rightarrow}^k(\sigma \star \mathbf{G})]_i = \sum_{j \in \mathcal{N}^*(i; \sigma \star \mathbf{G})} \Phi_{\rightarrow, \theta}^k(H_i^{k-1}(\sigma \star \mathbf{G}), (\sigma \star \mathbf{A})_{ij}, H_j^{k-1}(\sigma \star \mathbf{G})) \quad (15)$$

$$= \sum_{j \in \mathcal{N}^*(i; \sigma \star \mathbf{G})} \Phi_{\rightarrow, \theta}^k([\sigma \star \mathbf{H}^{k-1}(\mathbf{G})]_i, (\sigma \star \mathbf{A})_{ij}, [\sigma \star \mathbf{H}^{k-1}(\mathbf{G})]_j) \quad (16)$$

$$= \sum_{j \in \mathcal{N}^*(i; \sigma \star \mathbf{G})} \Phi_{\rightarrow, \theta}^k(H_{\sigma^{-1}(i)}^{k-1}(\mathbf{G}), A_{\sigma^{-1}(i)\sigma^{-1}(j)}, H_{\sigma^{-1}(j)}^{k-1}(\mathbf{G})) \quad (17)$$

$$= \sum_{\sigma^{-1}(j) \in \mathcal{N}^*(\sigma^{-1}(i); \mathbf{G})} \Phi_{\rightarrow, \theta}^k(H_{\sigma^{-1}(i)}^{k-1}(\mathbf{G}), A_{\sigma^{-1}(i)\sigma^{-1}(j)}, H_{\sigma^{-1}(j)}^{k-1}(\mathbf{G})) \quad (18)$$

$$= \sum_{j \in \mathcal{N}^*(\sigma^{-1}(i); \mathbf{G})} \Phi_{\rightarrow, \theta}^k(H_{\sigma^{-1}(i)}^{k-1}(\mathbf{G}), A_{\sigma^{-1}(i)j}, H_j^{k-1}(\mathbf{G})) \quad (19)$$

$$= [\phi_{\rightarrow}^k(\mathbf{G})]_{\sigma^{-1}(i)} \quad (20)$$

$$= [\sigma \star \phi_{\rightarrow}^k(\mathbf{G})]_i \quad (21)$$

All the above equalities are straightforward, except maybe eq. (18), which comes from the equivariance property of the notion of neighborhood defined above by eq. (13). The same property follows for ϕ_{\leftarrow}^k by similar argument.

For ϕ_{\odot}^k , eq. (7) gives

$$[\phi_{\odot}^k(\sigma \star \mathbf{G})]_i = \Phi_{\odot, \theta}^k([\mathbf{H}^{k-1}(\sigma \star \mathbf{G})]_i, (\sigma \star \mathbf{A})_{ii}) \quad (22)$$

$$= \Phi_{\odot, \theta}^k([\sigma \star \mathbf{H}^{k-1}(\mathbf{G})]_i, (\sigma \star \mathbf{A})_{ii}) \quad (23)$$

$$= \Phi_{\odot, \theta}^k(H_{\sigma^{-1}(i)}^{k-1}(\mathbf{G}), A_{\sigma^{-1}(i)\sigma^{-1}(i)}) \quad (24)$$

$$= [\phi_{\odot}^k(\mathbf{G})]_{\sigma^{-1}(i)} \quad (25)$$

$$= [\sigma \star \phi_{\odot}^k(\mathbf{G})]_i \quad (26)$$

This concludes the proof that \mathbf{H}_i^k is permutation equivariant for all k , and from eq. (8) we conclude that \mathbf{M}_{θ}^k is permutation-equivariant. Similar proof holds for \mathbf{D}_{θ}^k and $\hat{\mathbf{U}}^k$, which in turn prove that

$$\hat{\mathbf{U}}(\sigma \star \mathbf{G}) = \sigma \star \hat{\mathbf{U}}(\mathbf{G}). \quad (27)$$

This concludes the proof. \square

B.2 Proof of Property 1

In Section 3, Property 1 states that if the loss function ℓ is permutation-invariant and if for any $\mathbf{G} \in \text{supp}(\mathcal{D})$ there exists a unique minimizer $\mathbf{U}^*(\mathbf{G})$ of problem (1), then \mathbf{U}^* is permutation-equivariant.

Let ℓ be a permutation-invariant loss function and \mathcal{D} a distribution such that for any $\mathbf{G} \in \text{supp}(\mathcal{D})$ there is a unique solution \mathbf{U}^* of problem 1. Let $\mathbf{G} = (n, \mathbf{A}, \mathbf{B}) \in \text{supp}(\mathcal{D})$ and $\sigma \in \Sigma_n$ a permutation.

$$\ell(\sigma \star \mathbf{U}^*(\mathbf{G}), \sigma \star \mathbf{G}) = \ell(\mathbf{U}^*(\mathbf{G}), \mathbf{G}) \quad \text{by invariance of } \ell \quad (28)$$

$$= \min_{\mathbf{U} \in \mathcal{U}} \ell(\mathbf{U}, \mathbf{G}) \quad \text{by definition of } \mathbf{U}^* \quad (29)$$

$$= \min_{\mathbf{U} \in \mathcal{U}} \ell(\sigma \star \mathbf{U}, \sigma \star \mathbf{G}) \quad \text{by invariance of } \ell \quad (30)$$

$$= \min_{\mathbf{U} \in \mathcal{U}} \ell(\mathbf{U}, \sigma \star \mathbf{G}) \quad \text{by invariance of } \mathcal{U} \quad (31)$$

$$= \ell(\mathbf{U}^*(\sigma \star \mathbf{G}), \sigma \star \mathbf{G}) \quad \text{by definition of } \mathbf{U}^* \quad (32)$$

Moreover the uniqueness of the solution ensures that $\mathbf{U}^*(\sigma \star \mathbf{G}) = \sigma \star \mathbf{U}^*(\mathbf{G})$, which concludes the proof.

B.3 Proof of Theorem 1

We will now prove Theorem 1, the main result on DSSs, by closely following the approach of [25]. We will first prove a modified version of the Stone-Weierstrass theorem, and then verify that the defining spaces for Interaction Graphs indeed verify the conditions of this theorem by proving several lemmas (most importantly Theorem 3 on separability).

Let $\mathcal{G}_{\text{eq.}} \subseteq \mathcal{G}$ be a set of compact, permutation-invariant Interaction Graphs. The compactness implies that there exist $\bar{n} \in \mathbb{N}$ such that all graphs in $\mathcal{G}_{\text{eq.}}$ have an amount of nodes lower than $\bar{n} \in \mathbb{N}$. Let $\mathcal{C}_{\text{eq.}}(\mathcal{G}_{\text{eq.}}, \mathcal{U})$ be the space of continuous functions from $\mathcal{G}_{\text{eq.}}$ on \mathcal{U} that associate to any Interaction Graph $\mathbf{G} = (n, \mathbf{A}, \mathbf{B})$ one of its possible states $\mathbf{U} \in \mathbb{R}^n$. $(\mathcal{C}_{\text{eq.}}(\mathcal{G}_{\text{eq.}}, \mathcal{U}), +, \cdot, \odot)$ is a unital \mathbb{R} -algebra, where $(+, \cdot)$ are the usual addition and multiplication by a scalar, and \odot is the Hadamard product defined by $[(f \odot g)(x)]_i = [f(x)]_i \cdot [g(x)]_i$. Its unit is the constant function $\mathbf{1} = (1, \dots, 1)$.

Theorem 2 (Modified Stone-Weierstrass theorem for equivariant functions).

Let \mathcal{A} be a unital subalgebra of $\mathcal{C}_{\text{eq.}}(\mathcal{G}_{\text{eq.}}, \mathcal{U})$, (i.e., it contains the unit function $\mathbf{1}$) and assume both following properties hold:

- (Separability) *For all $\mathbf{G}, \mathbf{G}' \in \mathcal{G}_{\text{eq.}}$, with number of nodes n and n' such that \mathbf{G} is not isomorphic to \mathbf{G}' , and for all $k \in [n], k' \in [n']$, there exists $f \in \mathcal{A}$ such that $[f(\mathbf{G})]_k \neq [f(\mathbf{G}')]_{k'}$;*
- (Self-separability) *For all $n \leq \bar{n}$, $I \subseteq [n]$, $\mathbf{G} \in \mathcal{G}_{\text{eq.}}$ with n nodes, such that no isomorphism of \mathbf{G} exchanges at least one index between I and I^c , and for all $k \in I, l \in I^c$, there exists $f \in \mathcal{A}$ such that $[f(\mathbf{G})]_k \neq [f(\mathbf{G})]_l$.*

Then \mathcal{A} is dense in $\mathcal{C}_{\text{eq.}}(\mathcal{G}_{\text{eq.}}, \mathcal{U})$ with respect to the uniform metric.

This proof of Theorem 2 is almost identical to that of Theorem 4 in [25], with the following differences.

1. For the input space, we consider Interaction Graphs of the form $(n, \mathbf{A}, \mathbf{B})$ with $\mathbf{A} \in (\mathbb{R}^{d_A})^{n^2}$ and $\mathbf{B} \in (\mathbb{R}^{d_B})^n$, instead of hyper-graphs of the form \mathbb{R}^{n^d} for $d \in \mathbb{N}$. The corresponding metrics are naturally different, although the difference is not critical for the proof;
2. Similarly, we consider an output space with $\mathbf{U} \in (\mathbb{R}^{d_U})^n$ instead of \mathbb{R}^n ;
3. We only assume $\mathcal{G}_{\text{eq.}} \subseteq \mathcal{G}$ to be compact and permutation-invariant instead of a $\mathcal{G}_{\text{eq.}}$ with an explicit form: $\mathcal{G}_{\text{eq.}} := \{G \in \mathbb{R}^{n^d} | n \leq n_{\max}, \|G\| \leq R\}$ (which makes this modified theorem more general).

We shall then indicate how to bypass these differences one by one and then reuse the proofs in [25].

For 1, the only properties of the input space involved in [25] are the number of nodes, action of permutation and the metric (with the corresponding topology). For the first two points, everything is still applicable in our setting. For the topology, the difference is not critical either since we are actually considering the product space of two of metric spaces defined in [25] and all corresponding properties follow.

For 2, we can always reduce to the case with $d_U = 1$ then stack the resulting function d_U times to have the expected shape. This works seamlessly with Hadamard product and all properties related to density.

For 3, there is actually no dependency on the explicit form of $\mathcal{G}_{\text{eq.}}$ or \mathbf{G} in [25] (as for the case in 1). And the proof only relies on the upper bound on the number of nodes. So this generalization can be naturally obtained.

The detailed proof of Theorem 2 then follows the exact same procedure than that of Theorem 4 in [25], and we shall omit it here, referring the reader to [25] for all details.

Let $\bar{k} \in \mathbb{N}$, and, as defined in Section 4, let $\mathcal{H}^{\bar{k}}$ be the set of graph neural networks defined in Section 3 such that $\bar{k} \leq k$. Our goal is to prove that Theorem 2 can be applied to $\mathcal{H}^{\bar{k}}$.

Because $\mathcal{H}^{\bar{k}}$ is not an algebra, let us consider $\mathcal{H}^{\bar{k}\odot}$, the algebra generated by $\mathcal{H}^{\bar{k}}$ with respect to the Hadamard product. More formally:

$$\mathcal{H}^{\bar{k}\odot} = \left\{ \sum_{s=1}^S \bigodot_{t=1}^{T_s} c_{st} f_{st} \mid S \in \mathbb{N}, T_s \in \mathbb{N}, c_{st} \in \mathbb{R}, f_{st} \in \mathcal{H}^{\bar{k}} \right\}. \quad (33)$$

Note that the Hadamard product among f_{st} 's is well-defined since for a fixed input \mathbf{G} , all output values $f_{st}(\mathbf{G})$ take the same dimension - the size of \mathbf{G} .

$(\mathcal{H}^{\bar{k}\odot}, +, \cdot, \odot)$ is obviously a unital sub-algebra of $(\mathcal{G}_{\text{eq}}, +, \cdot, \odot)$ (the constant function $(1, \dots, 1)$ trivially belongs to $\mathcal{H}^{\bar{k}\odot}$). In order to apply Theorem 2 to $\mathcal{H}^{\bar{k}\odot}$, one needs to prove that it satisfies both separability hypotheses.

Let us first notice that the self-separability property is a straightforward consequence of the hypothesis of separability of external inputs on $\text{supp}(\mathcal{D})$. Hence we only need to prove the separability property:

Theorem 3. $\mathcal{H}^{\bar{k}\odot}$ satisfies the separability property of Theorem 2.

The proof consists of 3 steps. In step 1, we prove that for all $\mathbf{G}, \mathbf{G}' \in \text{supp}(\mathcal{D})$ that are not isomorphic, there exists a sequence (*node, edge, node*) that only exists in \mathbf{G} . In Step 2, we build a continuous function f^\dagger on \mathcal{G} that returns an indicator of the presence of this sequence in the input graph. In Step 3, we prove that there exists a function $f_\theta \in \mathcal{H}^{\bar{k}\odot}$ that approximates well enough f^\dagger .

For Step 1, we formally state it in the following lemma.

Lemma 1. Let $\mathbf{G} = (n, \mathbf{A}, \mathbf{B})$ and $\mathbf{G}' = (n', \mathbf{A}', \mathbf{B}')$ be in $\text{supp}(\mathcal{D})$ such that \mathbf{G} and \mathbf{G}' are not isomorphic and $n \geq n'$. Then there exist $i, j \in [n]$, $i \neq j$, such that, for all $i', j' \in [n']$, the following inequality holds:

$$(B_i, A_{ij}, B_j) \neq (B'_{i'}, A'_{i'j'}, B'_{j'}) \quad (34)$$

Proof. This lemma relies on the separability hypothesis of $\text{supp}(\mathcal{D})$ which states that there exists $\delta > 0$ such that for all $\mathbf{G} = (n, \mathbf{A}, \mathbf{B}) \in \text{supp}(\mathcal{D})$ and for all $i \neq j \in [n]$, $\|B_i - B_j\| \geq \delta$.

We shall use proof by contradiction: assume that for any $(i, j) \in [n]^2$ with $i \neq j$, there exists $\alpha(i, j) = (i', j') \in [n']^2$ such that $(B_i, A_{ij}, B_j) = (B'_{i'}, A'_{i'j'}, B'_{j'})$. Two cases must be distinguished, depending on whether $n < n'$ or $n = n'$.

If $n > n'$, then according to the pigeonhole principle, there exist two pairs $(i, j) \in [n]^2$ and $(l, m) \in [n]^2$ that have the same image by α , $(i', j') \in [n']^2$. Hence, $(B_i, A_{ij}, B_j) = (B'_{i'}, A'_{i'j'}, B'_{j'}) = (B_l, A_{lm}, B_m)$, which contradicts the separability hypothesis for \mathbf{G} .

If $n = n'$, according to the separability hypothesis of $\text{supp}(\mathcal{D})$, there cannot exist $i \neq l \in [n]$ that are mapped to the same $i' \in [n']$ (i.e. $\alpha(i, j) = (i', j')$ for some j, j' and $\alpha(l, m) = (i', m')$ for some m, m'). Thus α actually defines an injective mapping $\chi : [n] \rightarrow [n]$ on the first component. Because $n = n'$, this mapping is also surjective and hence bijective. Due to the symmetry of i and j , we see that the mapping on the second component χ' defined by X is exactly χ . Hence we have found a permutation $\chi \in \Sigma_n$ such that

$$B_i = B'_{\chi(i)} \quad (35)$$

$$A_{ij} = A'_{\chi(i)\chi(j)} \quad (36)$$

for any $(i, j) \in [n]^2$, which means that \mathbf{G} and \mathbf{G}' are isomorphic, contradicting the hypothesis, and thus completing the proof. \square

Let us now proceed with Step 2. For convenience, we shall use a continuous kernel function defined by

$$K_\epsilon(x) = \max(0, 1 - |x|/\epsilon) \quad (37)$$

for $\epsilon > 0$. Then we have $K_\epsilon(0) = 1$ and $K_\epsilon(x) = 0$ for $|x| > \epsilon$.

All intermediate functions of DSSs $\Phi_{\rightarrow}^k, \Phi_{\leftarrow}^k, \Phi_{\odot}^k, \Psi^k$ and Ξ^k (Section 3) live in function spaces that satisfy the Universal Approximation Property (UAP). So let us consider now a space of continuous

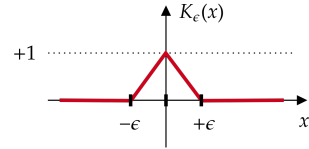


Figure 6: **Kernel function**

functions that share the same architecture than DSS, but in which all spaces of parameterized neural networks have been replaced by corresponding continuous function space. We denote this space by $\mathcal{H}^{\bar{k}\dagger}$ (by convention, a dagger(\dagger) added to a Neural Network block from Section 3 will refer to the corresponding continuous function space (e.g. $\Phi_{\rightarrow}^{k\dagger}$). We are now in position to prove the following lemma.

Lemma 2. *For any $\mathbf{G}, \mathbf{G}' \in \text{supp}(\mathcal{D})$ that are not isomorphic, there exists a function $f^\dagger \in \mathcal{H}^{\bar{k}\dagger}$ such that for any $k \in [n], k' \in [n']$, we have $[f^\dagger(\mathbf{G})]_k \neq [f^\dagger(\mathbf{G}')]_{k'}$.*

Proof. Without loss of generality, we suppose $n \geq n'$. According to Lemma 1, there exist $(i^\dagger, j^\dagger) \in [n]^2, i^\dagger \neq j^\dagger$, such that \mathbf{G} contains a sequence $(B_{i^\dagger}, A_{i^\dagger j^\dagger}, B_{j^\dagger})$ that does not appear in \mathbf{G}' .

We are going to construct a continuous function $f^\dagger : \text{supp}(\mathcal{D}) \rightarrow \mathcal{U}$ that will be an indicator of the presence of the above sequence in the graph, and such that $f^\dagger(\mathbf{G}) = (1, \dots, 1) \in \mathbb{R}^n$ and $f^\dagger(\mathbf{G}') = (0, \dots, 0) \in \mathbb{R}^{n'}$ (thus proving Lemma 2).

Let us first recall the architecture of DSS, as defined by eq. (5)-(9), and let us choose *continuous* functions $\Phi_{\rightarrow}^{1\dagger}, \Phi_{\leftarrow}^{1\dagger}, \Phi_{\circlearrowleft}^{1\dagger}$ and $\Psi^{1\dagger}$ such that $\mathbf{H}^{1\dagger}$ is defined by, for any $\mathbf{G}'' \in \text{supp}(\mathcal{D})$,

$$[\mathbf{H}^{1\dagger}(\mathbf{G}'')]_i = 2K_\epsilon(\|B''_i - B_{i^\dagger}\|) - K_\epsilon(\|B''_i - B_{j^\dagger}\|) \quad (38)$$

where $\epsilon = \|A_{i^\dagger j^\dagger} - A'_{\sigma(i^\dagger)\sigma(j^\dagger)}\|$ if \mathbf{B} and \mathbf{B}' are isomorphic through permutation σ and $\epsilon = \min_\sigma \max_i \|B_i - B'_{\sigma(i)}\|$ otherwise. This function allows us to identify whether the external input B''_i is close to one of B_{i^\dagger} or B_{j^\dagger} .

For $k = 2$, we define

$$\Phi_{\rightarrow}^{2\dagger}(h, a, h') = K_\epsilon(\|h - 2\| + \|a - A_{i^\dagger j^\dagger}\| + \|h' + 1\|) \quad (39)$$

and

$$\Phi_{\circlearrowleft}^{2\dagger}(h, a) = K_\epsilon(\|h - 1\| + \|a - A_{i^\dagger j^\dagger}\|). \quad (40)$$

Then we choose $\Psi^{2\dagger}$ such that

$$[\mathbf{H}^{2\dagger}]_i = \phi_{\circlearrowleft, i}^{2\dagger} + \phi_{\rightarrow, i}^{2\dagger} = \Phi_{\circlearrowleft}^{2\dagger}(H_i^{1\dagger}, A_{ii}) + \sum_{j \in \mathcal{N}^*(i; \mathbf{G})} \Phi_{\rightarrow}^{2\dagger}(H_i^{1\dagger}, A_{ij}, H_j^{1\dagger}) \quad (41)$$

According to the construction of $\mathbf{H}^{1\dagger}$ and $\mathbf{H}^{2\dagger}$, we have $[\mathbf{H}^{2\dagger}(\mathbf{G})]_i = 1$ if $i = i^\dagger$ and 0 otherwise. And $\mathbf{H}^{2\dagger}(\mathbf{G}') = (0, \dots, 0)$.

For $k \geq 3$, we let

$$[\mathbf{H}^{k+1\dagger}]_i = [\mathbf{H}^{k\dagger}]_i + \sum_{j \in \mathcal{N}^*(i; \mathbf{G})} [\mathbf{H}^{k\dagger}]_j \quad (42)$$

Thus if $\bar{k} \geq \Delta + 2$, we have $[\mathbf{H}^{\bar{k}\dagger}(\mathbf{G})]_i \geq 1$ for any $i \in [n]$, due to the connectivity and the fact that the diameter of \mathbf{G} is bounded by Δ , i.e. the propagation process described in eq. (42) reaches every node of \mathbf{G} . We have $[\mathbf{H}^{\bar{k}\dagger}(\mathbf{G})]_i \geq 1$ for any $i \in [n]$, and $[\mathbf{H}^{\bar{k}\dagger}(\mathbf{G}')]_i = 0$ for any $i \in [n']$

Finally for the decoder, we let

$$\Xi^{\bar{k}\dagger}(h) = \min(1, h) \quad (43)$$

and

$$[\hat{\mathbf{U}}^{\bar{k}\dagger}]_i = \Xi^{\bar{k}\dagger}(H_i^{\bar{k}\dagger}). \quad (44)$$

We have thus constructed a function f^\dagger such that $f^\dagger(\mathbf{G}) = (1, \dots, 1) \in \mathbb{R}^n$ and $f^\dagger(\mathbf{G}') = (0, \dots, 0) \in \mathbb{R}^{n'}$. Thus for any $k \in [n], k' \in [n']$, we have $[f^\dagger(\mathbf{G})]_k = 1 \neq 0 = [f^\dagger(\mathbf{G}')]_{k'}$, which concludes the proof. \square

Lemma 3. *Let X, Y, Z be three metric spaces. Let $\mathcal{F} \subseteq \mathcal{C}(X, Y)$ and $\mathcal{G} \subseteq \mathcal{C}(Y, Z)$ be two sets of continuous functions. And let $\mathcal{F}^\ell \subseteq \mathcal{F}, \mathcal{G}^\ell \subseteq \mathcal{G}$ be two subsets of Lipschitz functions that are dense in \mathcal{F} and \mathcal{G} respectively. Then $\mathcal{G}^\ell \circ \mathcal{F}^\ell := \{g \circ f | g \in \mathcal{G}^\ell, f \in \mathcal{F}^\ell\}$ is dense in $\mathcal{G} \circ \mathcal{F}$.*

Proof. Let $g \circ f$ be a continuous function in $\mathcal{G} \circ \mathcal{F}$, $\epsilon > 0$. Due to the density of \mathcal{G}^ℓ in \mathcal{G} , there exists $g^\ell \in \mathcal{G}^\ell$ such that

$$\bar{d}(g, g^\ell) < \frac{\epsilon}{2}. \quad (45)$$

Let L_{g^ℓ} be the Lipschitz constant of g^ℓ , the density of \mathcal{F}^ℓ in \mathcal{F} implies that there exists f^ℓ such that

$$\bar{d}(f, f^\ell) < \frac{\epsilon}{2L_{g^\ell}}. \quad (46)$$

Then we have

$$d_Z(g \circ f(x), g^\ell \circ f^\ell(x)) \leq d_Z(g \circ f(x), g^\ell \circ f(x)) + d_Z(g^\ell \circ f(x), g^\ell \circ f^\ell(x)) \quad (47)$$

$$< \frac{\epsilon}{2} + L_{g^\ell} d_Y(f(x), f^\ell(x)) \quad (48)$$

$$< \frac{\epsilon}{2} + L_{g^\ell} \frac{\epsilon}{2L_{g^\ell}} = \epsilon \quad (49)$$

for any $x \in X$. Thus $\bar{d}(g \circ f, g^\ell \circ f^\ell) < \epsilon$. Hence $\mathcal{G}^\ell \circ \mathcal{F}^\ell$ is dense in $\mathcal{G} \circ \mathcal{F}$. \square

Lemma 4. $\mathcal{H}^{\bar{k}}$ is dense in $\mathcal{H}^{\bar{k}^\dagger}$.

Proof. As functions in $\mathcal{H}^{\bar{k}}$ are composition of Lipschitz functions (neural network with linear transformation and Lipschitz activation as assumed), and all intermediate function spaces verify the Universal Approximation Property. We conclude immediately from using the definition of $\mathcal{H}^{\bar{k}^\dagger}$ and applying Lemma 3 consecutively. \square

We are ready to prove Theorem 3, i.e., that $\mathcal{H}^{\bar{k}^\odot}$ satisfies the separability hypothesis of Theorem 2.

Proof. of Theorem 3

It suffices to show the separability for $\mathcal{H}^{\bar{k}}$ since it is a subset of $\mathcal{H}^{\bar{k}^\odot}$.

Let $\mathbf{G}, \mathbf{G}' \in \text{supp}(\mathcal{D})$. According to Lemma 2, there exists $f^\dagger \in \mathcal{H}^{\bar{k}^\dagger}$ such that for any $k \in [n], k' \in [n']$, we have $[f^\dagger(\mathbf{G})]_k \neq [f^\dagger(\mathbf{G}')]_{k'}$. According to Lemma 4, there exists $f \in \mathcal{H}^{\bar{k}}$ such that

$$\bar{d}(f^\dagger, f) < \frac{1}{3}. \quad (50)$$

Then for any $k \in [n], k' \in [n']$, we have $[f(\mathbf{G})]_k > \frac{2}{3}$ and $[f(\mathbf{G}')]_{k'} < \frac{1}{3}$. This proves the separability of $\mathcal{H}^{\bar{k}}$ and furthermore, $\mathcal{H}^{\bar{k}^\odot}$. \square

Before being able to prove Theorem 1, we need the last following lemma.

Lemma 5. $\mathcal{H}^{\bar{k}}$ is dense in $\mathcal{H}^{\bar{k}^\odot}$.

Proof. We shall prove this result by explicitly constructing an approximation function in $\mathcal{H}^{\bar{k}}$ for a given function in $\mathcal{H}^{\bar{k}^\odot}$.

Let $f^\odot \in \mathcal{H}^{\bar{k}^\odot}$, and $\epsilon > 0$. By definition of $\mathcal{H}^{\bar{k}^\odot}$ in eq. (33), there exists $S \in \mathbb{N}$, $\{T_s\}_{s \in \{1, \dots, S\}} \in \mathbb{N}^S$, as well as $\{c_{st}\} \in \mathbb{R}$ and $\{f_{st}\} \in \mathcal{H}^{\bar{k}}$ for all (s, t) with $s \in [S], t \in [T_s]$, such that :

$$f^\odot = \sum_{s=1}^S \bigodot_{t=1}^{T_s} c_{st} f_{st} \quad (51)$$

Thus, for any (s, t) , there exists $\bar{k}_{st} \leq \bar{k}$, and $d_{st} \in \mathbb{N}$, such that f_{st} is composed of functions $\{\Phi_{\rightarrow, \theta}^{k, s, t}, \Phi_{\leftarrow, \theta}^{k, s, t}, \Phi_{\odot, \theta}^{k, s, t}, \Psi_{\theta}^{k, s, t}, \Xi_{\theta}^{k, s, t}\}_{k \in [\bar{k}_{st}]}$, as defined by eq. (5)-(9) and Figure 3 in Section 3. d_{st} is the dimension of the latent states of *channel* f_{st} .

The different channels can have different number of propagation updates \bar{k}_{st} , but they are all bounded by \bar{k} . Without loss of generality, we can assume that all \bar{k}_{st} are equal to \bar{k} by padding, when needed, exactly $\bar{k} - \bar{k}_{st}$ null operations Φ_{\rightarrow}^k , Φ_{\leftarrow}^k and Ψ^k before the actual ones.

Let $d = \sum_{s=1}^S \sum_{t=1}^{T_s} d_{st}$ be the cumulated dimensions of the different channels.

For each (s, t) , we introduce the matrix $W_{st} \in \{0, 1\}^{d_{st} \times d}$ which is defined by:

$$[W_{st}]_{ij} = \begin{cases} 1, & \text{if } \sum_{s'=1}^s \sum_{t'=1}^{T_{s'}} d_{s't'} + \sum_{t'=1}^{t-1} d_{st'} + i = j \\ 0, & \text{otherwise.} \end{cases} \quad (52)$$

Thus $W_{st} = [0, \dots, 0, I_{d_{st}}, 0, \dots, 0]$. Basically, when given a vector of dimension d , W_{st} will be able to select exactly the component that corresponds to the channel (s, t) , and will thus return a vector of dimension d_{st} .

Let us now define the functions $\{\Phi_{\rightarrow, \theta}^k, \Phi_{\leftarrow, \theta}^k, \Phi_{\odot, \theta}^k, \Psi_{\theta}^k, \Xi_{\theta}^k\}_{k \in [\bar{k}]}$ such that

$$\Phi_{\rightarrow, \theta}^k(H_i^{k-1}, A_{ij}, H_j^{k-1}) = \sum_{s=1}^S \sum_{t=1}^{T_s} W_{st}^\top \cdot \Phi_{\rightarrow, \theta}^{k, s, t}(W_{st} \cdot H_i^{k-1}, A_{ij}, W_{st} \cdot H_j^{k-1}) \quad (53)$$

$$\Phi_{\leftarrow, \theta}^k(H_i^{k-1}, A_{ij}, H_j^{k-1}) = \sum_{s=1}^S \sum_{t=1}^{T_s} W_{st}^\top \cdot \Phi_{\leftarrow, \theta}^{k, s, t}(W_{st} \cdot H_i^{k-1}, A_{ij}, W_{st} \cdot H_j^{k-1}) \quad (54)$$

$$\Phi_{\odot, \theta}^k(H_i^{k-1}, A_{ij}) = \sum_{s=1}^S \sum_{t=1}^{T_s} W_{st}^\top \cdot \Phi_{\odot, \theta}^{k, s, t}(W_{st} \cdot H_i^{k-1}, A_{ij}) \quad (55)$$

$$\Psi_{\theta}^k(H_i^{k-1}, B_i, \phi_{\rightarrow, i}^k, \phi_{\leftarrow, i}^k, \phi_{\odot, i}^k) = \sum_{s=1}^S \sum_{t=1}^{T_s} W_{st}^\top \cdot \Psi_{\theta}^{k, s, t}(W_{st} \cdot H_i^{k-1}, B_i, W_{st} \cdot \phi_{\rightarrow, i}^k, W_{st} \cdot \phi_{\leftarrow, i}^k, W_{st} \cdot \phi_{\odot, i}^k) \quad (56)$$

These functions, using eq. (5)-(8), define a function acting on a latent space of dimension d . Moreover, for any channel (s, t) and any node $i \in [n]$, we have $W_{st} \cdot H_i^{\bar{k}} = H_i^{\bar{k}, s, t}$.

We have thus built a function of $\mathcal{H}^{\bar{k}}$ that exactly replicates the steps performed on the different channels. Now, let us take a closer look at the decoding step.

Observing that the mapping from \mathbb{R}^d to \mathbb{R}^{d_U} , $h \mapsto \sum_{s=1}^S \bigodot_{t=1}^{T_s} c_{st} \Xi_{\theta}^{\bar{k}, s, t}(W_{st} h)$ is indeed continuous, there exists a mapping $\Xi_{\theta}^{\bar{k}} \in \mathcal{H}^{d_U}$ such that :

$$\|\Xi_{\theta}^{\bar{k}}(h) - \sum_{s=1}^S \bigodot_{t=1}^{T_s} c_{st} \Xi_{\theta}^{\bar{k}, s, t}(W_{st} h)\| \leq \epsilon \quad (57)$$

for any h in a compact of \mathbb{R}^d . The resulting function $f \in \mathcal{H}^{\bar{k}}$, composed of $\{\Phi_{\rightarrow, \theta}^k, \Phi_{\leftarrow, \theta}^k, \Phi_{\odot, \theta}^k, \Psi_{\theta}^k, \Xi_{\theta}^k\}_{k \in [\bar{k}]}$ using eq. (5)-(9), approximates f^{\odot} with precision less than ϵ , which concludes the proof. \square

We now have all necessary ingredients to prove Theorem 1.

Proof. According to the hypotheses of compactness and permutation-invariance on $\text{supp}(\mathcal{D})$, both conditions of Theorem 2 are satisfied by $\text{supp}(\mathcal{D})$. Consider the subalgebra $\mathcal{H}^{\bar{k} \odot}$ defined by eq. (33). According to the hypothesis of separability of external inputs, the hypothesis of connectivity and Theorem 3, $\mathcal{H}^{\bar{k} \odot}$ satisfies the separability and self-separability conditions of Theorem 2. Applying Theorem 2, it comes that $\mathcal{H}^{\bar{k} \odot}$ is dense in $\mathcal{C}_{\text{eq}}(\text{supp}(\mathcal{D}))$. Then according to Lemma 5, $\mathcal{H}^{\bar{k}}$ is dense in $\mathcal{H}^{\bar{k} \odot}$. We conclude that $\mathcal{H}^{\bar{k}}$ is dense in $\mathcal{C}_{\text{eq}}(\text{supp}(\mathcal{D}))$ by the transitivity property of density. \square

B.4 Proof of Corollary 1

Proof. Let $\epsilon > 0$. From Property 1, \mathbf{U}^* is permutation-equivariant. Moreover, by hypothesis, \mathbf{U}^* is continuous. Thus $\mathbf{U}^* \in \mathcal{C}_{\text{eq}}(\text{supp}(\mathcal{D}))$.

And from Theorem 1, we know that there exists a function $\text{Solver}_\theta \in \mathcal{H}^{\Delta+2}$ such that

$$\forall \mathbf{G} \in \text{supp}(\mathcal{D}), \|\text{Solver}_\theta(\mathbf{G}) - \mathbf{U}^*(\mathbf{G})\| \leq \epsilon \quad (58)$$

□

C Linear Systems derived from the Poisson Equation

This appendix details the experiments of Section 5.1: it presents the data generation process, and also explains the change of variables that was made to help normalizing the data (not mentioned in the main paper for space reason, as it does not change the overall conclusions of the experiments). Finally, we also discuss an additional super generalization experiment briefly cited in the paper.

C.1 Data generation

Initial problem Consider a Poisson's equation with Dirichlet condition on its boundary $\partial\Omega$:

$$-\Delta u = f \text{ in } \Omega$$

$$u|_{\partial\Omega} = g$$

where Ω a spatial domain in \mathbb{R}^2 , and $\partial\Omega$ its boundaries. The right hand side f is defined on Ω , and the Dirichlet boundary condition g is defined on $\partial\Omega$. x and y will denote the classical 2D coordinates.

Random geometries Random 2D domains Ω are generated from 10 points, randomly sampled in the unit square. The Bézier curve that passes through these points is created, and is further subsampled to obtain approximately 100 points in the unit square. These points defines a polygon, that is used as the boundary $\partial\Omega$. See the left part of Figure 7 to see four instances.

Random f and g Functions f and g are defined by the following equations:

$$f(x, y) = r_1(x - 1)^2 + r_2y^2 + r_3, \quad (x, y) \in \Omega \quad (59)$$

$$g(x, y) = r_4x^2 + r_5y^2 + r_6xy + r_7x + r_8y + r_9, \quad (x, y) \in \partial\Omega \quad (60)$$

in which parameters r_i are uniformly sampled between -10 and 10.

Discretization The random 2D geometries are discretized using Fenics' standard mesh generation method (see Figure 7-right).

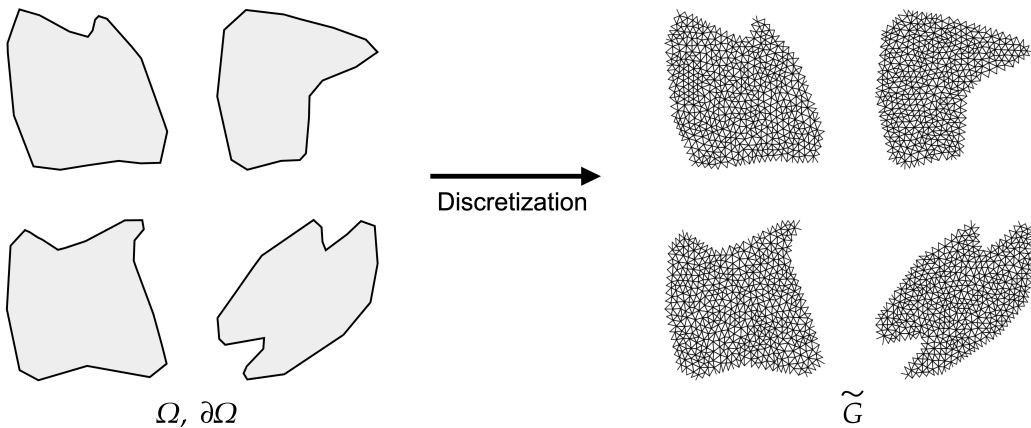


Figure 7: Discretization of randomly generated domains

Assembling The assembling step [30] consists in building a linear system from the partial differentiate equation and the discretized domain. The unknown are the values of the solution at the nodes of the mesh, and the equations are obtained by using the variational formulation of the PDE on basis functions with support in the neighbors of each node. This is also automatically performed using Fenics. The result of the assembling step is a square matrix \mathbf{A} and a vector \mathbf{B} , and the solution is the vector \mathbf{U} such that $\mathbf{A}\mathbf{U} = \mathbf{B}$. Thus, as stated in Section 5, in the framework of SSPs, an Interaction Graph is defined from the number of nodes of the mesh, the matrix \mathbf{A} and the vector \mathbf{B} , and the loss function is:

$$\ell(\mathbf{U}, \mathbf{G}) = \sum_{i \in [n]} (-B_i + \sum_{j \in [n]} A_{ij} U_j)^2 \quad (61)$$

C.2 Change of variables

Being able to properly normalize the input data of any neural network is a critical issue, and failing to do so can often lead to gradient explosions and other training failures (more details on data normalization in Appendix C.2). In the Poisson case study, the nodes at the boundary are constrained (*i.e.* $A_{ii} = 1$ and $A_{ij} = 0$ if $i \neq j$), and the interior nodes are not. Moreover, the coefficients of matrix \mathbf{A} at these interior nodes satisfy a conservation equality (*i.e.* $A_{ii} = -\sum_{j \in [n] \setminus \{i\}} A_{ij}$). As a consequence, the distributions of their respective B_i are very different, sometimes even with different orders of magnitude. It is then almost impossible to properly normalize those multimodal distribution.

In order to tackle this issue, we consider the following change of variable, changing \mathbf{A}, \mathbf{B} to \mathbf{A}', \mathbf{B}' , and modifying the loss function accordingly. For \mathbf{B} , we set the dimension $d_{B'}$ of \mathbf{B}' to 3 as follows:

$$B'_i = \begin{cases} [B_i, 0, 0] & \text{if node } i \text{ is not constrained} \\ [0, 1, B_i] & \text{otherwise} \end{cases} \quad (62)$$

The B_i 's for constrained and unconstrained nodes will hence be normalized independently.

Moreover, the information stored in the matrix \mathbf{A} is rather redundant. As mentioned, for constrained nodes $A_{ii} = 1$ and $A_{ij} = 0$ if $i \neq j$, whereas for unconstrained nodes $A_{ii} = -\sum_{j \in [n] \setminus \{i\}} A_{ij}$. Hence the diagonal information can always be retrieved from \mathbf{B} and the non diagonal elements of \mathbf{A} . We thus choose the following change of variable:

$$A'_{ij} = \begin{cases} A_{ij} & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases} \quad (63)$$

Finally, the loss function is transformed into the following function ℓ' (where $B_i'^p$ denotes the p^{th} component of vector B_i):

$$\ell'(\mathbf{U}, \mathbf{G}') = \sum_{i \in [n]} \left((1 - B_i'^2)(-B_i'^1) + B_i'^2(U_i - B_i'^3) + \sum_{j \in [n]} A'_{ij}(U_j - U_i) \right)^2 \quad (64)$$

One can easily check that this change of variables and of loss function defines the exact same optimization problem as in eq. (10), while allowing for an easier normalization, as well as a lighter sparse storage of \mathbf{A} .

C.3 Additional super generalization experiment

This appendix describes a second experiment regarding super-generalization. Figure 8 displays the results of the DSS model, learned without any noise, when increasing noise is added to the test examples, more and more diverging from the distribution of the training set (the graph size remains unchanged). Log-normal noise is applied to \mathbf{A} ($A_{ij} \exp(\mathcal{N}(0, \tau))$), and normal noise to \mathbf{B} ($B_i \mathcal{N}(1, \tau)$), for different values of noise variance τ . The correlation between the results of DSS and the 'ground truth', here given by the results of LU (solving the same noisy system). But although DSS results remain highly correlated with the ground truth for small values of

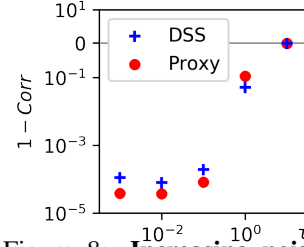


Figure 8: **Increasing noise variance τ :** Correlation (DSS, LU)

τ , they become totally uncorrelated for large values of τ (correlation close to 0): DSS has learned something specific to the distribution \mathcal{D} of linear systems coming from the discretized Poisson EDP. Further work will extend these results, analyzing in depth the specifics of the learned models.

D Power systems

This appendix gives more details about the AC power flow problem, and how it is converted into the DSS framework.

The AC power flow equations model the steady-state behavior of transportation power grids. They are an essential part of both real-time operation and long-term planning. A thorough overview of the domain is provided in [34].

Let's consider a power grid with n nodes. The voltage at every electrical node is a sinusoid that oscillates at the same frequency. However, each node has a distinct module and phase angle. Thus, we define the complex voltage at node i , $V_i = |V_i|e^{j\theta} \in \mathbb{C}$ (where j is the imaginary unit).

The admittance matrix $\mathbf{Y} = (Y_{ij})_{i,j \in [n]}$; $Y_{ij} \in \mathbb{C}$ defines the admittance of each power line of the network. The smaller $|Y_{ij}|$, the less nodes i and j are coupled. For $i, j \in [n]$, the coefficient Y_{ij} models the physical characteristics of the power line between nodes i and j (*i.e.* materials, length, etc.).

At each node i , there can be power consumption (houses, factories, etc.). The real part of the power consumed is denoted by $P_{d,i}$ and the imaginary part by $Q_{d,i}$. The subscript d stands for "demand". Additionally, there can also be power production (coal or nuclear power plants, etc.). They are very different from consumers, because they constrain the local voltage module. They are defined by $P_{g,i}$ and $V_{g,i}$. The subscript g stands for "generation". Nodes that have a producer attached to it are called "PV buses" and are denoted by $I_{PV} \subset [n]$. The nodes that are not connected to a production are called "PQ buses" and are denoted by $I_{PQ} \subset [n]$.

Moreover, one has to make sure that the global energy is conserved. There are losses at every power line that are caused by Joule's effect. The amount of power lost to Joule's effect being a function of the voltage at each node, it cannot be known before the voltage computation itself. Thus, to make sure that the production of energy equals the consumption plus the losses caused by Joule's effect, we need to be able to increase the power production accordingly. In this work we use the common "slack bus" approach which consists in increasing the production of a single producer so that global energy conservation holds. This node is chosen beforehand and we denote it by $i_s \in [n]$.

Thus the system of equations that govern the power grid is the following:

$$\forall i \in [n] \setminus \{i_s\}, \quad P_{g,i} - P_{d,i} = \sum_{j \in [n]} |V_i||V_j|(\operatorname{Re}(Y_{ij}) \cos(\theta_i - \theta_j) + \operatorname{Im}(Y_{ij}) \sin(\theta_i - \theta_j)) \quad (65)$$

$$\forall i \in I_{PQ}, \quad -Q_{d,i} = \sum_{j \in [n]} |V_i||V_j|(\operatorname{Re}(Y_{ij}) \sin(\theta_i - \theta_j) - \operatorname{Im}(Y_{ij}) \cos(\theta_i - \theta_j)) \quad (66)$$

$$\forall i \in I_{PV}, \quad |V_i| = V_{g,i} \quad (67)$$

The encoding into our framework requires a bit of work. For the coupling matrix we use $d_A = 2$ and $A_{ij} = [\operatorname{Re}(Y_{ij}), \operatorname{Im}(Y_{ij})]$. For the local input we take $d_B = 5$ and $B_i = [P_{g,i} - P_{d,i}, Q_{d,i}, 1(i \in I_{PQ}), V_{g,i}, 1(i = i_s)]$. Finally, for the state variable we use $d_U = 2$ and take $U_i = [|V_i|, \theta_i]$.

Taking the squared residual of eq. (65)-(67) and taking the sum over every node, we obtain the loss of eq. (11).

E Further implementation details

In this section we detail the implementation details that were made to robustify the training of the DSS. None of those changes alter the properties of the architecture.

Correction coefficient We introduce a parameter α that modifies eq. (42) in the following way:

$$H_i^k = H_i^{k-1} + \alpha \times \Psi_\theta^k(H_i^{k-1}, B_i, \phi_{\rightarrow,i}^k, \phi_{\leftarrow,i}^k, \phi_{\odot,i}^k) \quad (68)$$

Choosing a sufficiently low value of α , helps to keep the successive \bar{k} updates at reasonably low orders of magnitude.

Injecting existing solutions Depending on the problem at hand, it may be useful to initialize the predictions to some known value. This acts as an offset, that can help the training process to start not too far from the actual solutions. This offset is applied identically at every node, thus not breaking the permutation-equivariance of the architecture:

$$\hat{U}_i^k = U_{offset} + \Xi_\theta^k(H_i^k) \quad (69)$$

For instance, in the power systems application, it is known that the voltage module is commonly around 1.0, while the voltage angle is around 0. Thus we used $U_{offset} = [1, 0]$ (keeping in mind that $d_U = 2$). On the other hand, in the linear systems application, there is no reason to use such an offset, so we used $U_{offset} = [0]$ (keeping in mind that here $d_U = 1$). But in several contexts, there exists some fast inaccurate method that can give an approximate solution closer to the final one than $(0, \dots, 0)$.

Data normalization In addition to a potential change of variables (which helps disentangle multimodal distributions of the input data, see Appendix C.2), it is also critical to normalize the input Interaction Graph to help with the training of neural networks. Each function $\Phi_{\rightarrow, \theta}^k$, $\Phi_{\leftarrow, \theta}^k$ and $\Phi_{\circlearrowleft, \theta}^k$ take A_{ij} as input, and the functions Ψ_θ^k take b_i as input. We thus introduce hyperparameters $\mu_A, \sigma_A \in \mathbb{R}^{d_A}$ and $\mu_B, \sigma_B \in \mathbb{R}^{d_B}$ are used to create a normalized version of the data:

$$a_{ij} = \frac{A_{ij} - \mu_A}{\sigma_A} \quad (70)$$

$$b_i = \frac{B_i - \mu_B}{\sigma_B} \quad (71)$$

$\mathbf{g} = (\mathbf{a}, \mathbf{b})$ (with $\mathbf{a} = (a_{ij})_{i,j \in [n]}$ and $\mathbf{b} = (b_i)_{i \in [n]}$) is thus the normalized version of \mathbf{G} . We apply the DSS to this normalized \mathbf{g} and consider the loss $\ell(\text{Solver}_\theta(\mathbf{g}), \mathbf{G})$ instead of $\ell(\text{Solver}_\theta(\mathbf{G}), \mathbf{G})$.

Gradient clipping We sometimes observed (e.g., in the power systems experiments) some gradient explosions. The solution we are currently using is to perform some gradient clipping. Further work should focus on facilitating this training process automatically.