

---

# Deep Statistical Solvers

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

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

## 1 Introduction

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

31 Another important property of natural or artificial systems is that of invariance, a fundamental concept  
32 in science, allowing to generalize conclusions drawn from few observations, to whole invariance  
33 classes. This work focuses on permutation-invariant problems, which appear in simulations of  
34 complex systems of interacting and interchangeable entities [?] (*e.g.*, molecular dynamics, power  
35 grids, simulations of partial differential equations (PDEs) with finite elements). Invariance has  
36 made its way in machine learning, as illustrated by the success of Convolutional Neural Networks  
37 (CNN) [?, ?], and of Graph Neural Networks (GNN) [?, ?]. In particular, implementations of GNNs  
38 successfully handle materials dynamics simulations [?], power systems [?], interacting particles [?]

and classical [?] or quantum [?] 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  $\{1, \dots, n\}$ ;  $\Sigma_n$  is the set of permutations of  $[n]$ ; for any  $\sigma \in \Sigma_n$ , any set  $\Omega$  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  $\Omega$ ),  $\sigma \star \mathbf{m}$  is the matrix  $(m_{\sigma^{-1}(i)\sigma^{-1}(j)})_{i,j \in [n]}$ .

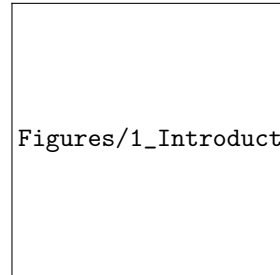
**Interaction Graphs** We call *Interaction Graph* a system of  $n \in \mathbb{N}$  interacting entities, or *nodes*, defined as  $\mathbf{G} = (n, \mathbf{A}, \mathbf{B})$ , where  $n$  is the size of  $\mathbf{G}$  (number of nodes),  $\mathbf{A} = (A_{ij})_{i,j \in [n]}$ ;  $A_{ij} \in \mathbb{R}^{d_A}$ ;  $d_A \geq 1$  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  $(n, \sigma \star \mathbf{A}, \sigma \star \mathbf{B})$ .

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-zero<sup>1</sup>. 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 is no confusion). A *loss function*  $\ell$  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  $\sigma \in \Sigma_n$ ,  $\mathcal{F}(\sigma \star \mathbf{G}) = \sigma \star \mathcal{F}(\mathbf{G})$ .

<sup>1</sup>A more rigorous definition of the actual underlying graph structure is deferred to Appendix A



Figures/1\_Introduction/graph.png

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

## 2.2 Problem Statement

**The Optimization Problem** In the remaining of the paper,  $\ell$  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  $\ell$  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  $\ell$  (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 [?] 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  $\ell$  (Figure 2.b), no such examples are needed.

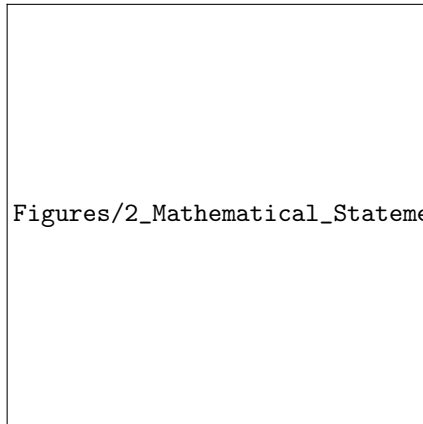


Figure 2: **Proxy approach (a) vs. DSS (b)**

## 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  $\ell$  is permutation-invariant and if for any  $\mathbf{G} \in \operatorname{supp}(\mathcal{D})$  there exists a unique minimizer  $\mathbf{U}^*(\mathbf{G})$  of problem (1), then  $\mathbf{U}^*$  is permutation-equivariant.*

Figures/3\_Graph\_Neural\_Network\_architecture/dense\_architecture.png

Figure 3: **Graph Neural Network implementation of a DSS**

129 Graph Neural Networks, introduced in [?], and further developed in [?, ?] (see also the recent surveys  
 130 [?, ?]), are a class of parameterized permutation-equivariant functions. They are hence natural  
 131 candidates to build SSP solutions, since Property 1 states that the ideal solver  $\mathbf{U}^*$  is permutation-  
 132 equivariant.

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

138 The overall architecture is described in Figure 3. All latent states in  $\mathbf{H}^0$  are initialized to a zero vector.  
 139 The *message passing* step performs  $\bar{k}$  updates on the latent state variable  $\mathbf{H}$  using  $\mathbf{M}_\theta^k$ , spreading  
 140 information using interaction coefficients  $\mathbf{A}$  and external inputs  $\mathbf{B}$  of  $\mathbf{G}$  (eq. 5–8). After each update,  
 141 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  
 142 output of the algorithm  $\hat{\mathbf{U}}$ . However, in order to robustify learning, all intermediate states  $\hat{\mathbf{U}}^k$  are  
 143 taken into account in the training loss through a discounted sum with hyperparameter  $\gamma \in [0, 1]$ :

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

144 **Message passing  $\mathbf{M}_\theta^k$**  For each node  $i$ , three different messages are computed,  $\phi_{\rightarrow, \theta}^k, \phi_{\leftarrow, \theta}^k, \phi_{\circ, \theta}^k$ ,  
 145 corresponding to outgoing, ingoing and self-loop links, respectively using trainable mappings

146  $\Phi_{\rightarrow,\theta}^k, \Phi_{\leftarrow,\theta}^k, \Phi_{\odot,\theta}^k$ , as follows:

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

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

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

147 Latent states  $H_i^k$  are then computed using trainable mapping  $\Psi_\theta^k$ , in a ResNet-like fashion:

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

148 **Decoding** The decoding step applies the same trainable mapping  $\Xi_\theta^k$  to every node:

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

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

152 **Inference Complexity** Assuming that each neural network block has a single hidden layer with  
 153 dimension  $d$ , that  $d \geq d_A, d_B, d_U$ , and denoting by  $m$  the average neighborhood size, one inference  
 154 has computational complexity of order  $\mathcal{O}(mn\bar{k}d^3)$ , scaling linearly with  $n$ . Furthermore, many  
 155 problems involve very local interactions, resulting in small  $m$ . However, one should keep in mind  
 156 that hyperparameters  $\bar{k}$  and  $d$  should be chosen according to the characteristics of distribution  $\mathcal{D}$ .

157 **Equivariance** The proposed architecture defines permutation-equivariant DSS, as proved in Appendix  
 158 B.1.

## 159 4 Deep Statistical Solvers are Universal Approximators for SSPs Solutions

160 This Section proves, heavily relying on [?], a Universal Approximation Theorem for the class of  
 161 DSSs with Lipschitz activation function (e.g. ReLU) in the space of the solutions of SSPs.

162 The space of Interaction Graphs is a metric space for the distance

$$163 \quad d(\mathbf{G}, \mathbf{G}') = \|\mathbf{A} - \mathbf{A}'\| + \|\mathbf{B} - \mathbf{B}'\| \text{ if } n = n' \text{ and } +\infty, \text{ otherwise}$$

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

167 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  
 168 since [?] that the set of neural networks with at least one hidden layer, an arbitrarily large amount of  
 169 hidden neurons, and an appropriate activation function, satisfies these conditions.

170 **Hypothesis space** Let  $\bar{k} \in \mathbb{N}$ . We denote by  $\bar{\mathcal{H}}^{\bar{k}}$  the set of graph neural networks defined in Section  
 171 3 such that  $\bar{k} \leq k, d \in \mathbb{N}$  and for any  $k = 1, \dots, \bar{k}$ , we consider all possible  $\Phi_{\rightarrow,\theta}^k, \Phi_{\leftarrow,\theta}^k \in \mathcal{H}_{d_A+2d}^d$ ,  
 172  $\Phi_{\odot,\theta}^k \in \mathcal{H}_{d_A+d}^d, \Psi_\theta^k \in \mathcal{H}_{d_B+4d}^d$  and  $\Xi_\theta^k \in \mathcal{H}_{d_U}^{d_U}$ .

173 **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  
 174 unweighted graph structure, as defined in Section 2.1. We will write  $\text{diam}(\mathbf{G})$  for  $\text{diam}(\tilde{\mathbf{G}})$ , the  
 175 diameter of  $\tilde{\mathbf{G}}$  [?].

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

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

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

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

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

188 **Sketch of the proof** (see Appendix B.3 for all details) Still following [?], we first prove a modified  
 189 version of the *Stone-Weierstrass theorem for equivariant functions*. This theorem guarantees that a  
 190 certain subalgebra of functions is dense in the set of continuous and permutation-equivariant functions  
 191 if it separates non-isomorphic Interaction Graphs. Following the idea of Hornik et al. [?], we extend  
 192 the hypothesis space to ensure closure under addition and multiplication. We then prove that the  
 193 initial hypothesis space is dense in this new subalgebra. Finally, we conclude the proof by showing  
 194 that the separability property mentioned above is satisfied by this newly-defined subalgebra.

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

$$198 \quad \forall \mathbf{G} \in \text{supp}(\mathcal{D}), \|\text{Solver}_\theta(\mathbf{G}) - \mathbf{U}^*(\mathbf{G})\| \leq \epsilon$$

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

## 203 5 Experiments

204 This section investigates the behavior and performances of DSSs on two SSPs. The first one amounts  
 205 to solving linear systems, though the distribution of problems is generated from a discretized Poisson  
 206 PDE. The second is the (non-quadratic) AC power flow computation. In all cases, the dataset is split  
 207 into training/validation/test sets, the hyperparameters that are not explicitly mentioned are found by  
 208 trial and errors using the validation set, and **all results presented are results on the test set**.

209 All trainings are performed with the Adam optimizer [?] with the standard hyperparameters of  
 210 TensorFlow 1.14 [?], running on an Nvidia GeForce RTX 2080 Ti. Gradient clipping was used to  
 211 avoid exploding gradient issues. **In the following, all experiments were repeated three times, with**  
 212 **the same datasets and different initialization seeds** (as reported in Tables 1 and 2). Our code is  
 213 made available in the supplementary materials, and links to the datasets can be found as references.

### 214 5.1 Solving Linear Systems from a Discretized PDE

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

221 For illustration purposes, the Poisson equation can be used to model a field of temperature. In Figure  
 222 4, the geometry (house profile) is shown in the Top Left. The result of the optimization is the field of  
 223 temperature everywhere in the house (shown in the Top Right).

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

$$226 \quad \ell(\mathbf{U}, \mathbf{G}) = \sum_{i \in [n]} \left( \sum_{j \in [n]} A_{ij} U_j - B_i \right)^2 \quad (10)$$

227 Our goal here is of course not to solve the Poisson equation, nor is it to propose a new competitive  
 228 method to invert linear systems. As a matter of fact, the proposed approach does not make use of the  
 229 linearity of the problem. Our goal is actually twofold: i) validate the DSS approach in high dimension  
 230 ( $n \approx 500$  nodes), and ii) analyze how DSS learns the distribution  $\mathcal{D}$ . Here, the distribution  $\mathcal{D}$  is  
 231 defined by the specific structure of linear systems that result from the discretization of the Poisson

Figures/6\_Experiments/poisson\_loss.png

Figure 4: **Intermediate losses and predictions** - Top left: the structure graph  $\tilde{\mathbf{G}}$  (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  $\hat{\mathbf{U}}^k$  are displayed for several values of  $k$ .

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 [?] 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 don't allow a precise control of  $n$ ).

The number of updates  $\bar{k}$  is set to 30 (average diameter size for meshes with 500 elements). Each NN block has one hidden layer of dimension  $d = 10$  and a leaky-ReLU non linearity. The complete DSS has 49,830 weights. Training is done for 280,000 iterations with batch size 100, and lasted 48h.

Two baseline methods are considered [?], 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 can decrease LU computational cost by a factor 6 [?]).

**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

Method	DSS (3 runs)			LU	BGS ( $10^{-3}$ )
Correlation w/ LU	99.99%			-	-
Time per instance (ms)	1.8			2.4	2.3
Loss $10^{th}$ percentile	$3.0 \cdot 10^{-4}$	$3.8 \cdot 10^{-4}$	$2.7 \cdot 10^{-4}$	$4.5 \cdot 10^{-27}$	$1.3 \cdot 10^{-3}$
Loss median	$6.0 \cdot 10^{-4}$	$1.3 \cdot 10^{-3}$	$6.9 \cdot 10^{-4}$	$6.1 \cdot 10^{-26}$	$1.7 \cdot 10^{-2}$
Loss $90^{th}$ percentile	$1.5 \cdot 10^{-3}$	$4.0 \cdot 10^{-3}$	$2.3 \cdot 10^{-3}$	$6.3 \cdot 10^{-25}$	$1.1 \cdot 10^{-1}$

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.

accuracy (a tunable parameter of BGS). Further work will explore how DSS scales up in much higher dimensions, in particular when LU becomes intractable.

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 model achieves a higher correlation with the LU solutions for data points with a lower number of nodes. 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.

Figures/6\_Experiments/shifting\_n.png

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 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 the violation of this law is directly used as loss function  $\ell$ . Moreover, some constraints over the states  $\mathbf{U}$  are here relaxed and included as an additional term of the loss (with factor  $\lambda$ ). 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  $|V_i|$ ,  $\theta_i$ ,  $P_{ij}$  (real part of power flow) and  $Q_{ij}$  (imaginary part). This problem is highly non-linear, and a substantial overview is provided in [?]. 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  $\ell$ :

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

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  $B_i^5$ .



Dataset	IEEE 14 nodes				IEEE 118 nodes			
Method	DSS (3 runs)			NR	DSS (3 runs)			NR
Correlation w/ NR	99.99%			-	99.99%			-
Time per instance (ms)	$1 \times 10^{-2}$			$2 \times 10^1$	$9 \times 10^{-2}$			$2 \times 10^1$
Loss $10^{th}$ percentile $\times 10^6$	25	34	41	$1.4 \times 10^{-6}$	0.53	13	1.7	$2.9 \times 10^{-8}$
Loss median $\times 10^6$	40	63	100	$2.1 \times 10^{-6}$	1.3	17	2.6	$4.2 \times 10^{-8}$
Loss $90^{th}$ percentile $\times 10^6$	100	150	440	$3.3 \times 10^{-6}$	3.4	25	4.8	$6.4 \times 10^{-8}$

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.

**Experimental conditions** Experiments are conducted on two standard benchmarks from the *Learning to Run a Power Network* competition [?]: IEEE case14 ( $n = 14$ ) [?], and IEEE case118 ( $n = 118$ ) [?]. 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),  $k$  was set to 10 (resp. 30). The number of weights is 1,722 for each of the  $\bar{k}$  ( $\mathbf{M}$ ,  $\mathbf{D}$ ) blocks, 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 (PandaPower implementation [?], on an Intel i5 dual-core (2.3GHz)). To the best of our knowledge, no GPU implementation was available, although recent work [?, ?] 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  $Q_{ij}$ ). In the 14 nodes case, correlations for  $V_i$  and  $\theta_i$  are also above 99.99%. In the 118 nodes case, those correlations are respectively around 99.9% and 95% (except for model #2, for which it is 99.7% and 83%). This might be caused by the presence of large but loosely interconnected communities in the graphs of case118, and will be the object of further investigations. 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 Deep Statistical Solvers 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, showing a good compromise between accuracy and speed in dimensions up to 500 on 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.

Future work will focus on incorporating discrete variables in the state space. 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.



## 332 A Underlying Graph Structure

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

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

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

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

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

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

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

## 348 B Mathematical proofs

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

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

### 353 B.1 Proof of equivariance of the proposed DSS architecture

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

355 *Proof.* Because the loss function  $\ell$  is permutation invariant, we only have to prove that eq. (8)-(9)  
 356 satisfy the permutation-equivariance property.

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

359 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  $\mathbf{G}$ .

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

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

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

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

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

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

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

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

361 All the above equalities are straightforward, except maybe eq. (18), which comes from the equivari-  
 362 ance property of the notion of neighborhood defined above by eq. (13). The same property follows  
 363 for  $\phi_{\leftarrow}^k$  by similar argument.

364 For  $\phi_{\odot}^k$ , eq. (7) gives

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

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

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

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

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

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

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

367 This concludes the proof.  $\square$

## 368 B.2 Proof of Property 1

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

372  
 373 Let  $\ell$  be a permutation-invariant loss function and  $\mathcal{D}$  a distribution such that for any  $\mathbf{G} \in \text{supp}(\mathcal{D})$   
 374 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  
 375 permutation.

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

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

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

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

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

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

### 378 B.3 Proof of Theorem 1

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

383 Let  $\mathcal{G}_{\text{eq.}} \subseteq \mathcal{G}$  be a set of compact, permutation-invariant Interaction Graphs. The compactness implies  
 384 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  
 385  $\mathcal{C}_{\text{eq.}}(\mathcal{G}_{\text{eq.}}, \mathcal{U})$  be the space of continuous functions from  $\mathcal{G}_{\text{eq.}}$  on  $\mathcal{U}$  that associate to any Interaction  
 386 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,  
 387 where  $(+, \cdot)$  are the usual addition and multiplication by a scalar, and  $\odot$  is the Hadamard product  
 388 defined by  $[(f \odot g)(x)]_i = [f(x)]_i \cdot [g(x)]_i$ . Its unit is the constant function  $\mathbf{1} = (1, \dots, 1)$ .

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

390 *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  
 391 following properties hold:*

- 392 • (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  
 393 isomorphic to  $\mathbf{G}'$ , and for all  $k \in [n], k' \in [n']$ , there exists  $f \in \mathcal{A}$  such that  $[f(\mathbf{G})]_k \neq$   
 394  $[f(\mathbf{G}')]_{k'}$ ;*
- 395 • (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  
 396 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  
 397  $f \in \mathcal{A}$  such that  $[f(\mathbf{G})]_k \neq [f(\mathbf{G})]_l$ .*

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

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

- 400 1. For the input space, we consider Interaction Graphs of the form  $(n, \mathbf{A}, \mathbf{B})$  with  $\mathbf{A} \in$   
 401  $(\mathbb{R}^{d_A})^{n^2}$  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  
 402 corresponding metrics are naturally different, although the difference is not critical for the  
 403 proof;
- 404 2. Similarly, we consider an output space with  $\mathbf{U} \in (\mathbb{R}^{d_U})^n$  instead of  $\mathbb{R}^n$ ;
- 405 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  
 406 an explicit form:  $\mathcal{G}_{\text{eq.}} := \{G \in \mathbb{R}^{n^d} | n \leq n_{\max}, \|G\| \leq R\}$  (which makes this modified  
 407 theorem more general).

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

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

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

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

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

422 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  
 423 3 such that  $\bar{k} \leq \bar{k}$ . Our goal is to prove that Theorem 2 can be applied to  $\mathcal{H}^{\bar{k}}$ .

424 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  
 425 Hadamard product. More formally:

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

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

428  $(\mathcal{H}^{\bar{k}\odot}, +, \cdot, \odot)$  is obviously a unital sub-algebra of  $(\mathcal{G}_{\text{eq}}, +, \cdot, \odot)$  (the constant function  $(1, \dots, 1)$   
 429 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  
 430 both separability hypotheses.

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

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

434 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,  
 435 there exists a sequence (*node, edge, node*) that only exists in  $\mathbf{G}$ . In Step 2, we build a continuous  
 436 function  $f^\dagger$  on  $\mathcal{G}$  that returns an indicator of the presence of this sequence in the input graph. In Step  
 437 3, we prove that there exists a function  $f_\theta \in \mathcal{H}^{\bar{k}\odot}$  that approximates well enough  $f^\dagger$ .

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

439 **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  
 440 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  
 441 inequality holds:

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

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

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

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

450 If  $n = n'$ , according to the separability hypothesis of  $\text{supp}(\mathcal{D})$ , there cannot exist  $i \neq l \in [n]$  that are  
 451 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  
 452  $m, m'$ ). Thus  $\alpha$  actually defines an injective mapping  $\chi : [n] \rightarrow [n]$  on the first component. Because  
 453  $n = n'$ , this mapping is also surjective and hence bijective. Due to the symmetry of  $i$  and  $j$ , we see  
 454 that the mapping on the second component  $\chi'$  defined by  $X$  is exactly  $\chi$ . Hence we have found a  
 455 permutation  $\chi \in \Sigma_n$  such that

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

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

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

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

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

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

461 All intermediate functions of DSSs  $\Phi_{\rightarrow}^k, \Phi_{\leftarrow}^k, \Phi_{\odot}^k, \Psi^k$  and  $\Xi^k$  (Sec-  
 462 tion 3) live in function spaces that satisfy the Universal Approxima-  
 463 tion Property (UAP). So let us consider now a space of continuous

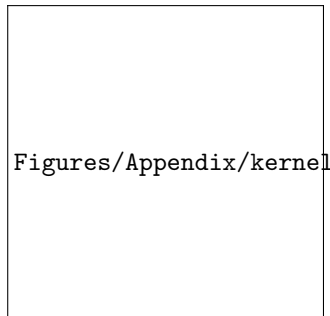


Figure 6: Kernel function

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  $\Phi_{\rightarrow}^{1\dagger}, \Phi_{\leftarrow}^{1\dagger}, \Phi_{\circ}^{1\dagger}$  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  $\sigma$  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

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

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

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

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  $\Delta$ , 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  $\mathcal{G} \circ \mathcal{F}$ .*

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

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

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

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

502 Then we have

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

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

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

503 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  $\mathcal{G} \circ \mathcal{F}$ .

504 □

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

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

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

511 *Proof.* of Theorem 3

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

513 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  $k \in [n], k' \in$   
514  $[n']$ , 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

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

515 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  
516 separability of  $\mathcal{H}^{\bar{k}}$  and furthermore,  $\mathcal{H}^{\bar{k}^\odot}$ . □

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

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

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

521 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  $S \in \mathbb{N}$ ,  $\{T_s\}_{s \in \{1, \dots, S\}} \in$   
522  $\mathbb{N}^S$ , 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 :

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

523 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  
524  $\{\Phi_{\rightarrow, \theta}^{k, s, t}, \Phi_{\leftarrow, \theta}^{k, s, t}, \Phi_{\odot, \theta}^{k, s, t}, \Psi_{\theta}^{k, s, t}, \Xi_{\theta}^{k, s, t}\}_{k \in [\bar{k}_{st}]}$ , as defined by eq. (5)-(9) and Figure 3 in Section 3.  $d_{st}$   
525 is the dimension of the latent states of *channel*  $f_{st}$ .



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

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

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

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

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

534 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

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

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

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

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

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

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

539 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 continu-  
 540 ous, there exists a mapping  $\Xi_{\theta}^{\bar{k}} \in \mathcal{H}_d^{d_U}$  such that :

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

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

544 □

545 We now have all necessary ingredients to prove Theorem 1.

546 *Proof.* According to the hypotheses of compactness and permutation-invariance on  $\text{supp}(\mathcal{D})$ , both  
 547 conditions of Theorem 2 are satisfied by  $\text{supp}(\mathcal{D})$ . Consider the subalgebra  $\mathcal{H}^{\bar{k} \odot}$  defined by eq. (33).  
 548 According to the hypothesis of separability of external inputs, the hypothesis of connectivity and  
 549 Theorem 3,  $\mathcal{H}^{\bar{k} \odot}$  satisfies the separability and self-separability conditions of Theorem 2. Applying  
 550 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  
 551 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. □

## 552 B.4 Proof of Corollary 1

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

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

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

556

□

## 557 C Linear Systems derived from the Poisson Equation

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

### 562 C.1 Data generation

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

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

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

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

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

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

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

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

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

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

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

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

### 581 C.2 Change of variables

582 Being able to properly normalize the input data of any neural network is a critical issue, and failing  
 583 to do so can often lead to gradient explosions and other training failures (more details on data  
 584 normalization in AppendixC.2). In the Poisson case study, the nodes at the boundary are constrained

Figures/Appendix/discretization.png

Figure 7: **Discretization of randomly generated domains**

(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  $\ell'$  (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  $\tau$ . 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  $\tau$ , they become totally uncorrelated for large values of  $\tau$  (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.

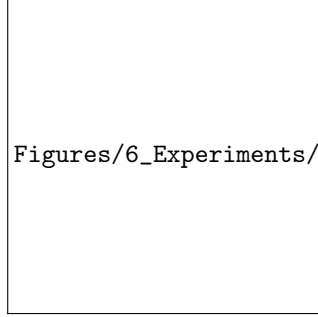


Figure 8: **Increasing noise variance  $\tau$** : Correlation (DSS, LU)

## 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 [?].

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  $\mathbf{j}$  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  $Q_{g,i}$ . 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  $\alpha$  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  $\alpha$ , 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  $\Phi_{\odot,\theta}^k$  take  $A_{ij}$  as input, and the functions  $\Psi_\theta^k$  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.