

## A Appendix

### A.1 A Counter Example of a Basis of the Null Space

We consider a special case of Eq. (9), where  $a_i(\mathbf{x}) \equiv 0$ ,  $b_i(\mathbf{x}) \equiv 1$ ,  $g_i(\mathbf{x}) \equiv 0$ , and the dimension is  $d = 3$ .

$$n_1(\mathbf{x})p_{j1}(\mathbf{x}) + n_2(\mathbf{x})p_{j2}(\mathbf{x}) + n_3(\mathbf{x})p_{j3}(\mathbf{x}) = 0, \quad \mathbf{x} \in \gamma_i, \quad \forall i = 1, \dots, m_j. \quad (\text{A1})$$

And a counter example of  $\mathbf{B}(\mathbf{x})$  is given by

$$\begin{aligned} \mathbf{B}(\mathbf{x}) &= [\beta_1(\mathbf{x}), \beta_2(\mathbf{x})], \\ \beta_1(\mathbf{x}) &= [n_2(\mathbf{x}), -n_1(\mathbf{x}), 0]^\top, \\ \beta_2(\mathbf{x}) &= [n_3(\mathbf{x}), 0, -n_1(\mathbf{x})]^\top. \end{aligned} \quad (\text{A2})$$

One could verify that the above formula of  $\mathbf{B}(\mathbf{x})$  is a basis of the null space, if  $n_1(\mathbf{x}) \neq 0, \forall \mathbf{x} \in \gamma_i$ . For a special case where  $\gamma_i$  is a plane parallel to the  $x$ -axis, however, we have  $n_1(\mathbf{x}) \equiv 0, \forall \mathbf{x} \in \gamma_i$ . In this case,  $\beta_1(\mathbf{x}), \beta_2(\mathbf{x})$  are no longer linearly independent and cannot represent all possible solutions to  $(p_{j1}(\mathbf{x}), p_{j2}(\mathbf{x}), p_{j3}(\mathbf{x}))$ . Therefore, Eq. (A2) is not an admissible choice for  $\mathbf{B}(\mathbf{x})$ .

### A.2 A Basis of the Null Space in Low Dimensions

Let  $\tilde{\mathbf{n}} = (a_i, b_i \mathbf{n}) / \sqrt{a_i^2 + b_i^2}$ ,  $\tilde{g}_i = g_i / \sqrt{a_i^2 + b_i^2}$ , and  $\tilde{\mathbf{p}}_j = (u_j, \mathbf{p}_j)$ . Eq. (9) is equivalent to

$$\tilde{\mathbf{n}}(\mathbf{x}) \cdot \tilde{\mathbf{p}}_j(\mathbf{x}) = \tilde{g}_i(\mathbf{x}), \quad \mathbf{x} \in \gamma_i, \quad \forall i = 1, \dots, m_j. \quad (\text{A3})$$

For  $d = 1$ , we can rewrite Eq. (A3) as (the dimension of  $\tilde{\mathbf{p}}_j$  is  $d + 1$ )

$$\tilde{n}_1(\mathbf{x})\tilde{p}_{j1}(\mathbf{x}) + \tilde{n}_2(\mathbf{x})\tilde{p}_{j2}(\mathbf{x}) = \tilde{g}_i(\mathbf{x}), \quad \mathbf{x} \in \gamma_i, \quad \forall i = 1, \dots, m_j. \quad (\text{A4})$$

And we can find that the following basis is an acceptable one

$$\mathbf{B}(\mathbf{x}) = [\tilde{n}_2(\mathbf{x}), -\tilde{n}_1(\mathbf{x})]^\top, \quad (\text{A5})$$

since  $\mathbf{B}(\mathbf{x}) = 0 \Leftrightarrow \tilde{\mathbf{n}}(\mathbf{x}) = 0$ , and the latter contradicts the fact that  $\tilde{\mathbf{n}} \cdot \tilde{\mathbf{n}} = 1$ . Then, we can use  $\mathbf{B}$  to construct the general solution  $\tilde{\mathbf{p}}_j^{\gamma_i}$  under  $d = 1$ .

And for  $d = 2$ , Eq. (A3) becomes

$$\tilde{n}_1(\mathbf{x})\tilde{p}_{j1}(\mathbf{x}) + \tilde{n}_2(\mathbf{x})\tilde{p}_{j2}(\mathbf{x}) + \tilde{n}_3(\mathbf{x})\tilde{p}_{j3}(\mathbf{x}) = \tilde{g}_i(\mathbf{x}), \quad \mathbf{x} \in \gamma_i, \quad \forall i = 1, \dots, m_j. \quad (\text{A6})$$

An acceptable  $\mathbf{B}(\mathbf{x})$  is given by

$$\begin{aligned} \mathbf{B}(\mathbf{x}) &= [\beta_1(\mathbf{x}), \beta_2(\mathbf{x}), \beta_3(\mathbf{x})], \\ \beta_1(\mathbf{x}) &= [0, \tilde{n}_3(\mathbf{x}), -\tilde{n}_2(\mathbf{x})]^\top, \\ \beta_2(\mathbf{x}) &= [-\tilde{n}_3(\mathbf{x}), 0, \tilde{n}_1(\mathbf{x})]^\top, \\ \beta_3(\mathbf{x}) &= [\tilde{n}_2(\mathbf{x}), -\tilde{n}_1(\mathbf{x}), 0]^\top. \end{aligned} \quad (\text{A7})$$

We note that  $\beta_1(\mathbf{x}), \beta_2(\mathbf{x}), \beta_3(\mathbf{x})$  live in the null space and  $\text{rank}(\mathbf{B}(\mathbf{x})) = 2$ . So  $\mathbf{B}(\mathbf{x})$  contains a basis in the null space, which can be used to construct the general solution  $\tilde{\mathbf{p}}_j^{\gamma_i}$  under  $d = 2$ .

### A.3 Explanation for the General Solution

We first show how to find an admissible expression of  $\mathbf{B}(\mathbf{x})$  in arbitrary dimensions with respect to Eq. (A3) which is equivalent to original formulation of the BC (see Eq. (9)). We perform a Gram-Schmidt orthogonalization of  $\tilde{\mathbf{n}}$  (whose dimension is  $d + 1$ ) on each vector in the standard basis to get

$$\beta_k(\mathbf{x}) = \mathbf{e}_k - \frac{\mathbf{e}_k \cdot \tilde{\mathbf{n}}(\mathbf{x})}{\tilde{\mathbf{n}}(\mathbf{x}) \cdot \tilde{\mathbf{n}}(\mathbf{x})} \tilde{\mathbf{n}}(\mathbf{x}) = \mathbf{e}_k - (\mathbf{e}_k \cdot \tilde{\mathbf{n}}(\mathbf{x})) \tilde{\mathbf{n}}(\mathbf{x}), \quad k = 1, \dots, d + 1, \quad (\text{A8})$$

where  $[\mathbf{e}_1, \dots, \mathbf{e}_{d+1}] = \mathbf{I}_{d+1}$ , and obviously all  $\beta_k(\mathbf{x}), k = 1, \dots, d + 1$ , are in the  $\text{Null}(\tilde{\mathbf{n}}^\top)$ . We set  $\mathbf{B}(\mathbf{x}) = [\beta_1(\mathbf{x}), \dots, \beta_{d+1}(\mathbf{x})] = \mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x})\tilde{\mathbf{n}}(\mathbf{x})^\top$ . Furthermore, we can prove that  $\text{rank}(\mathbf{B}(\mathbf{x})) = d, \forall \mathbf{x} \in \gamma_i$  (see Lemma A.1). Therefore, for  $\forall \mathbf{x} \in \gamma_i$ ,  $\mathbf{B}(\mathbf{x})$  always contains a basis of  $\text{Null}(\tilde{\mathbf{n}}^\top)$ , and we consider such a  $\mathbf{B}(\mathbf{x})$  to be an ideal choice for the general solution  $\tilde{\mathbf{p}}_j^{\gamma_i}$ .

**Lemma A.1.**  $\text{rank}(\mathbf{B}(\mathbf{x})) = d$  holds for all  $\mathbf{x} \in \gamma_i$ , where  $\mathbf{B}(\mathbf{x}) = \mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x})\tilde{\mathbf{n}}(\mathbf{x})^\top$ .

*Proof.* For all  $\mathbf{x} \in \gamma_i$ , we have known that  $\mathbf{B} = \mathbf{I}_{d+1} - \tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top$ , where  $\tilde{\mathbf{n}} \cdot \tilde{\mathbf{n}} = \tilde{\mathbf{n}}^\top \tilde{\mathbf{n}} = 1$ , and

$$\mathbf{B}\tilde{\mathbf{n}} = (\mathbf{I}_{d+1} - \tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top)\tilde{\mathbf{n}} = \tilde{\mathbf{n}} - \tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top \tilde{\mathbf{n}} = \tilde{\mathbf{n}} - \tilde{\mathbf{n}} = \mathbf{0}. \quad (\text{A9})$$

Hence,  $\text{rank}(\mathbf{B}) \leq d$ . Besides, we notice that  $\mathbf{H} = \mathbf{I}_{d+1} - 2\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top$  is a Householder matrix, which is an invertible matrix, since

$$\mathbf{H}^\top \mathbf{H} = (\mathbf{I}_{d+1} - 2\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top)^2 = \mathbf{I}_{d+1} - 4\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top + 4\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top \tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top = \mathbf{I}_{d+1} - 4\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top + 4\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top = \mathbf{I}_{d+1}. \quad (\text{A10})$$

So  $\text{rank}(\mathbf{H}) = d + 1$ , and we have

$$d + 1 = \text{rank}(\mathbf{I}_{d+1} - 2\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top) \leq \text{rank}(\mathbf{I}_{d+1} - \tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top) + \text{rank}(\tilde{\mathbf{n}}\tilde{\mathbf{n}}^\top) = \text{rank}(\mathbf{B}) + 1, \quad (\text{A11})$$

which can deduce  $d \leq \text{rank}(\mathbf{B})$ . Therefore,  $\text{rank}(\mathbf{B}) = d$ .  $\square$

Finally, we show that the general solution in Eq. (10) satisfies the BC in Eq. (A3).

$$\begin{aligned} \tilde{\mathbf{n}}(\mathbf{x}) \cdot \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}) &= \tilde{\mathbf{n}}(\mathbf{x}) \cdot \mathbf{B}(\mathbf{x}) \text{NN}_j^{\gamma_i}(\mathbf{x}) + \tilde{\mathbf{n}}(\mathbf{x}) \cdot \tilde{\mathbf{n}}(\mathbf{x}) \tilde{g}_i(\mathbf{x}) \\ &= \tilde{\mathbf{n}}(\mathbf{x}) \cdot (\mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x})\tilde{\mathbf{n}}(\mathbf{x})^\top) \text{NN}_j^{\gamma_i}(\mathbf{x}) + \tilde{g}_i(\mathbf{x}) \\ &= \tilde{g}_i(\mathbf{x}), \end{aligned} \quad (\text{A12})$$

where we omit the trainable parameters for simplicity. Besides, the discussion of  $\mathbf{B}(\mathbf{x})$  in low-dimensional cases (i.e.,  $d = 1$  and  $d = 2$ , see Appendix A.3) is similar, and we will leave it to the reader.

#### A.4 Theoretical Guarantee of the Constructed Ansatz

In Appendix A.3, we have demonstrated that  $\mathbf{B}(\mathbf{x})$  contains a basis of the null space of the BC for  $\forall \mathbf{x} \in \gamma_i$  and the general solution in Eq. (10) satisfies the corresponding BC. In this subsection, we will show that our constructed ansatz in Eq. (11) is theoretically correct. We first prove that the ansatz in Eq. (11) satisfies all the BCs under the following assumptions.

**Assumption A.2.** The problem domain  $\Omega$  is bounded.

**Assumption A.3.** The shortest distance between  $\gamma_1, \dots, \gamma_{m_j}$  is greater than zero for  $j = 1, \dots, n$ .

**Assumption A.4.** All the extended distance functions  $l^{\partial\Omega}, l^{\gamma_i}, i = 1, \dots, m_j$  are continuous and satisfy that  $\min_{\mathbf{x} \in \partial\Omega \setminus \gamma_i} l^{\gamma_i}(\mathbf{x}) \geq C_i, \forall \mathbf{x} \in \gamma_i, i = 1, \dots, m_j$  for  $j = 1, \dots, n$ , where  $C_i$  is a positive constant.

**Theorem A.5.**  $\forall \epsilon > 0$ , there exists  $\beta_s^0 \in \mathbb{R}$ , such that

$$|\tilde{\mathbf{n}}(\mathbf{x}) \cdot (\hat{u}_j, \hat{\mathbf{p}}_j) - \tilde{g}_i(\mathbf{x})| < \epsilon, \quad (\text{A13})$$

holds for all  $\beta_s > \beta_s^0, \mathbf{x} \in \gamma_i, i = 1, \dots, m_j, j = 1, \dots, n$ , where  $\tilde{\mathbf{n}} = (a_i, b_i \mathbf{n}) / \sqrt{a_i^2 + b_i^2}$ ,  $\tilde{g}_i = g_i / \sqrt{a_i^2 + b_i^2}$ , and  $\mathbf{B}(\mathbf{x}) = \mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x})\tilde{\mathbf{n}}(\mathbf{x})^\top$ .

*Proof.* For any  $\mathbf{x} \in \gamma_i$ , we have  $l^{\partial\Omega}(\mathbf{x}) = 0$  according to the definition of the extended distance functions. Thus,  $(\hat{u}_j, \hat{\mathbf{p}}_j)$  is now equal to

$$(\hat{u}_j, \hat{\mathbf{p}}_j) = \sum_{k=1}^{m_j} \exp[-\alpha_k l^{\gamma_k}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_k}(\mathbf{x}), \quad (\text{A14})$$

where we omit the trainable parameters. Then, according to Assumptions A.2 ~ A.4, we can choose a sufficiently large  $\beta_s^i$  (see Eq. (12) for the relationship between  $\alpha_i$  and  $\beta_s$ ), such that

$$\begin{aligned}
|\tilde{\mathbf{n}}(\mathbf{x}) \cdot (\hat{u}_j, \hat{\mathbf{p}}_j) - \tilde{g}_i(\mathbf{x})| &= \left| \tilde{\mathbf{n}}(\mathbf{x}) \cdot \left( \sum_{k=1}^{m_j} \exp[-\alpha_k l^{\gamma_k}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_k}(\mathbf{x}) \right) - \tilde{g}_i(\mathbf{x}) \right| \\
&< \left| \tilde{\mathbf{n}}(\mathbf{x}) \cdot \exp[-\alpha_i l^{\gamma_i}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}) - \tilde{g}_i(\mathbf{x}) \right| \\
&\quad + \left| \tilde{\mathbf{n}}(\mathbf{x}) \cdot \sum_{k \neq i} \exp[-\alpha_k l^{\gamma_k}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_k}(\mathbf{x}) \right| \\
&\leq \left| \tilde{\mathbf{n}}(\mathbf{x}) \cdot \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}) - \tilde{g}_i(\mathbf{x}) \right| + \left| \sum_{k \neq i} \exp[-\alpha_k l^{\gamma_k}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_k}(\mathbf{x}) \right| \\
&= \left| \sum_{k \neq i} \exp[-\alpha_k l^{\gamma_k}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_k}(\mathbf{x}) \right| \\
&< \epsilon,
\end{aligned} \tag{A15}$$

where we note that  $l^{\gamma_i}(\mathbf{x}) = 0$  and  $l^{\gamma_k}(\mathbf{x}) > 0, \forall k \neq i$  for all  $\mathbf{x} \in \gamma_i$ . Let  $\beta_s^0 = \max\{\beta_s^i \mid i = 1, \dots, m_j\}$ , then according to the arbitrariness of  $j$ , we conclude that the theorem holds.  $\square$

Next, we will prove that our ansatz can approximate the solution to the PDEs with arbitrarily low errors under following assumptions in addition to Assumptions A.2 ~ A.4.

**Assumption A.6.** The solution to the PDEs  $u(\mathbf{x})$  is unique, bounded, and at least first order continuous by element.

**Assumption A.7.**  $a_i(\mathbf{x})$ ,  $b_i(\mathbf{x})$ , and  $g_i(\mathbf{x})$  are continuous (hence  $\tilde{g}_i(\mathbf{x})$  is continuous, too) in  $\gamma_i$  for  $i = 1, \dots, m_j, j = 1, \dots, n$ .

**Assumption A.8.** Since  $\mathbf{B}(\mathbf{x}) = \mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x})\tilde{\mathbf{n}}(\mathbf{x})^\top$  is a real symmetric matrix, we can perform an orthogonal diagonalization  $\mathbf{B}(\mathbf{x}) = \mathbf{P}(\mathbf{x})^\top \Lambda(\mathbf{x}) \mathbf{P}(\mathbf{x})$ , where  $\Lambda(\mathbf{x}) = \text{diag}(\lambda_1(\mathbf{x}), \dots, \lambda_d(\mathbf{x}), 0)$ ,  $\lambda_1(\mathbf{x}) > \dots > \lambda_d(\mathbf{x}) > 0$ . We assume that  $\tilde{\mathbf{n}}(\mathbf{x})$ ,  $\mathbf{P}(\mathbf{x})$ , and  $\Lambda(\mathbf{x})$  are piece-wise continuous by element in  $\gamma_i$  for  $i = 1, \dots, m_j, j = 1, \dots, n$ .

To begin with, we prove this lemma.

**Lemma A.9.**  $\forall \epsilon > 0$ , there exists  $\theta_j^{\gamma_i} \in \Theta_j^{\gamma_i}$ , such that

$$\|\tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}; \theta_j^{\gamma_i}) - \mathbf{q}(\mathbf{x})\|_1 < \epsilon, \tag{A16}$$

holds for all  $\mathbf{x} \in \gamma_i$  if  $\mathbf{q}(\mathbf{x})$  is continuous in  $\gamma_i$  and satisfies the BC (i.e.,  $\tilde{\mathbf{n}}(\mathbf{x}) \cdot \mathbf{q}(\mathbf{x}) = \tilde{g}_i(\mathbf{x}), \forall \mathbf{x} \in \gamma_i$ ), where  $\Theta_j^{\gamma_i}$  is the parameter space of the neural network  $\text{NN}_j^{\gamma_i}$ ,  $\|\cdot\|_1$  is the 1-norm of matrices (operator norm), and  $\tilde{\mathbf{p}}_j^{\gamma_i}$  as well as  $\mathbf{q}$  are both of dimension  $d+1$ . The above conclusion holds for all  $i = 1, \dots, m_j, j = 1, \dots, n$ .

*Proof.* From Eq. (A8) and Lemma A.1, we know that  $\mathbf{B}(\mathbf{x})$  contain a basis of  $\text{Null}(\tilde{\mathbf{n}}^\top)$ . Since  $\mathbf{q}(\mathbf{x})$  satisfies the BC, we can represent it as

$$\mathbf{q}(\mathbf{x}) = \mathbf{B}(\mathbf{x})\mathbf{r}(\mathbf{x}) + \tilde{\mathbf{n}}(\mathbf{x})\tilde{g}_i(\mathbf{x}). \tag{A17}$$

Then we will show that there exists a piece-wise continuous choice of  $\mathbf{r}(\mathbf{x})$ . We rewrite Eq. (A17) as

$$\mathbf{q}(\mathbf{x}) = \mathbf{P}(\mathbf{x})^\top \Lambda(\mathbf{x}) \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x}) + \tilde{\mathbf{n}}(\mathbf{x})\tilde{g}_i(\mathbf{x}). \tag{A18}$$

Since  $\mathbf{B}(\mathbf{x})$  has  $d+1$  column vectors, which is greater than the dimension of  $\text{Null}(\tilde{\mathbf{n}}^\top)$  (i.e.,  $d$ ), the choice of  $\mathbf{q}(\mathbf{x})$  is not unique. We can choose a particular  $\mathbf{q}(\mathbf{x})$  which satisfies that the last element of  $\mathbf{P}(\mathbf{x})\mathbf{r}(\mathbf{x})$  is zero (i.e.,  $\mathbf{P}(\mathbf{x})\mathbf{r}(\mathbf{x}) = [\dots, 0]^\top$ ). Next, we continue with the equivalent

transformation of Eq. (A18).

$$\begin{aligned}
& \mathbf{P}(\mathbf{x})^\top \Lambda(\mathbf{x}) \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x}) + \tilde{\mathbf{n}}(\mathbf{x}) \tilde{g}_i(\mathbf{x}) = \mathbf{q}(\mathbf{x}), \\
\iff & \mathbf{P}(\mathbf{x})^\top \Lambda(\mathbf{x}) \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x}) = \mathbf{q}(\mathbf{x}) - \tilde{\mathbf{n}}(\mathbf{x}) \tilde{g}_i(\mathbf{x}), \\
\iff & \Lambda(\mathbf{x}) \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x}) = \mathbf{P}(\mathbf{x}) (\mathbf{q}(\mathbf{x}) - \tilde{\mathbf{n}}(\mathbf{x}) \tilde{g}_i(\mathbf{x})), \\
\iff & \text{diag}(1, \dots, 1, 0) \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x}) = \Lambda'(\mathbf{x}) \mathbf{P}(\mathbf{x}) (\mathbf{q}(\mathbf{x}) - \tilde{\mathbf{n}}(\mathbf{x}) \tilde{g}_i(\mathbf{x})), \\
\iff & \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x}) = \Lambda'(\mathbf{x}) \mathbf{P}(\mathbf{x}) (\mathbf{q}(\mathbf{x}) - \tilde{\mathbf{n}}(\mathbf{x}) \tilde{g}_i(\mathbf{x})),
\end{aligned} \tag{A19}$$

where  $\Lambda'(\mathbf{x}) = \text{diag}(1/\lambda_1(\mathbf{x}), \dots, 1/\lambda_d(\mathbf{x}), 0)$ . The last equivalence holds because the last element of  $\mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x})$  is always zero. From Assumption A.7 and A.8, combining the above formula, we have that  $\mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x})$  is piece-wise continuous by element. Noticing that  $\mathbf{r}(\mathbf{x}) = \mathbf{P}(\mathbf{x})^\top \mathbf{P}(\mathbf{x}) \mathbf{r}(\mathbf{x})$ , we know that the  $\mathbf{r}(\mathbf{x})$  we chosen is also piece-wise continuous by element.

We notice that

$$\begin{aligned}
\|\mathbf{B}(\mathbf{x})\|_1 &= \|\mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x}) \tilde{\mathbf{n}}(\mathbf{x})^\top\|_1 \\
&\leq \|\mathbf{I}_{d+1}\|_1 + \|\tilde{\mathbf{n}}(\mathbf{x}) \tilde{\mathbf{n}}(\mathbf{x})^\top\|_1 \\
&\leq 1 + d + 1 \\
&= d + 2.
\end{aligned} \tag{A20}$$

According to the Universal Approximation of neural networks [2],  $\forall \epsilon > 0$ , there exists  $\boldsymbol{\theta}_j^{\gamma_i} \in \Theta_j^{\gamma_i}$ , such that

$$\|\text{NN}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i}) - \mathbf{r}(\mathbf{x})\|_1 < \frac{\epsilon}{d+2}, \tag{A21}$$

holds for all  $\mathbf{x} \in \gamma_i$ . Therefore,

$$\begin{aligned}
\|\tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i}) - \mathbf{q}(\mathbf{x})\|_1 &= \|\mathbf{B}(\mathbf{x}) (\text{NN}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i}) - \mathbf{r}(\mathbf{x}))\|_1 \\
&\leq \|\mathbf{B}(\mathbf{x})\|_1 \|\text{NN}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i}) - \mathbf{r}(\mathbf{x})\|_1 \\
&< \epsilon.
\end{aligned} \tag{A22}$$

According to the arbitrariness of  $i$  and  $j$ , we conclude that the lemma holds.  $\square$

Finally, we state the following theorem.

**Theorem A.10.**  $\forall \epsilon > 0$ , there exists  $\beta_s \in \mathbb{R}$ ,  $\boldsymbol{\theta}_{\text{main}} \in \Theta_{\text{main}}$ ,  $\boldsymbol{\theta}_j^{\gamma_i} \in \Theta_j^{\gamma_i}$ ,  $i = 1, \dots, m_j$ , such that

$$\|(\hat{u}_j, \hat{\mathbf{p}}_j) - (u_j, \nabla u_j)\|_1 < \epsilon, \tag{A23}$$

holds for all  $\mathbf{x} \in \Omega \cup \partial\Omega$ ,  $j = 1, \dots, n$ , where  $\Theta_*$  is the parameter space of the corresponding neural network,  $\|\cdot\|_1$  is the 1-norm. The ground truth solution to the PDEs is  $\mathbf{u}(\mathbf{x}) = (u_1(\mathbf{x}), \dots, u_n(\mathbf{x}))$ .

*Proof.* For  $\mathbf{x} \in \gamma_i$ ,  $i = 1, \dots, m_j$ ,  $(u_j, \nabla u_j)$  is continuous (according to Assumption A.6) and satisfies the BC (which the solution needs to meet). From Lemma A.9, the definition of  $(\hat{u}_j, \hat{\mathbf{p}}_j)$  in Eq. (11), and Assumptions A.2 ~ A.4, we can find  $\boldsymbol{\theta}_j^{\gamma_i} \in \Theta_j^{\gamma_i}$  and a large enough  $\beta_s^i$  such that Eq. (A23) holds for all  $\mathbf{x} \in \gamma_i$ .

Then we fix  $\beta_s = \max\{\beta_s^i \mid i = 1, \dots, m_j\}$  and  $\boldsymbol{\theta}_j^{\gamma_i} \in \Theta_j^{\gamma_i}$ ,  $i = 1, \dots, m_j$  (which are what we determined for  $\mathbf{x} \in \gamma_i$ ,  $i = 1, \dots, m_j$ ). From Assumption A.2, we have  $|l^{\partial\Omega}(\mathbf{x})| < C$ ,  $\forall \mathbf{x} \in \Omega$ , where  $C$  is a positive constant. For  $\mathbf{x} \in \Omega$ , according to the Universal Approximation Theorem of neural networks [1], there exists  $\boldsymbol{\theta}_{\text{main}} \in \Theta_{\text{main}}$  satisfying that

$$\left\| \frac{(u_j, \nabla u_j) - \sum_{i=1}^{m_j} \exp[-\alpha_i l^{\gamma_i}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i})}{l^{\partial\Omega}(\mathbf{x})} - \text{NN}_{\text{main}}(\mathbf{x}; \boldsymbol{\theta}_{\text{main}}) \right\|_1 < \frac{\epsilon}{C}. \tag{A24}$$

Therefore,

$$\begin{aligned}
& \|(\hat{u}_j, \hat{\mathbf{p}}_j) - (u_j, \nabla u_j)\|_1 \\
&= \left\| (u_j, \nabla u_j) - \sum_{i=1}^{m_j} \exp[-\alpha_i l^{\gamma_i}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i}) - l^{\partial\Omega}(\mathbf{x}) \text{NN}_{\text{main}}(\mathbf{x}; \boldsymbol{\theta}_{\text{main}}) \right\|_1 \\
&= \left\| \frac{(u_j, \nabla u_j) - \sum_{i=1}^{m_j} \exp[-\alpha_i l^{\gamma_i}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}; \boldsymbol{\theta}_j^{\gamma_i})}{l^{\partial\Omega}(\mathbf{x})} - \text{NN}_{\text{main}}(\mathbf{x}; \boldsymbol{\theta}_{\text{main}}) \right\|_1 |l^{\partial\Omega}(\mathbf{x})| \\
&< \frac{\epsilon}{C} \cdot C = \epsilon.
\end{aligned} \tag{A25}$$

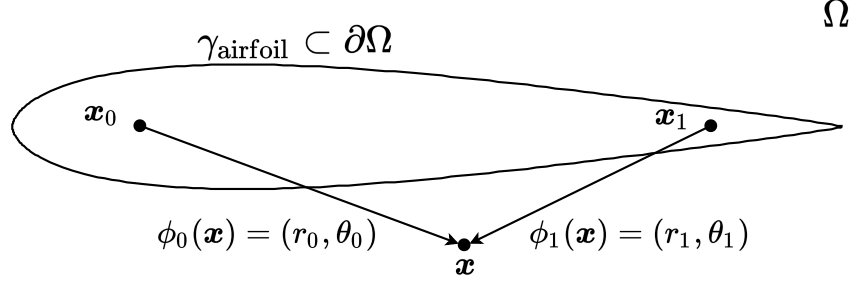


Figure A1: Illustration of the extension from  $\gamma_{\text{airfoil}}$  to  $\Omega \cup \partial\Omega$ .

According to the arbitrariness of  $j$ , we have proven this theorem.  $\square$

Besides, we note that it is easy to extend the above theorems to time-dependent cases (the ansatz is given in Appendix A.6), which will not be discussed separately here.

### A.5 Extension of the Parameter Functions in the BCs

In Eq. (9), it is noted that  $a_i$ ,  $b_i$ ,  $\mathbf{n}$  or  $g_i$  may be only defined at  $\gamma_i$ . But they are included in our ansatz (see Eq. (10) and Eq. (11)), which is defined in  $\Omega \cup \partial\Omega$ . So we need to extend their definition smoothly to  $\Omega \cup \partial\Omega$ , using interpolation or approximation via neural networks. We consider the airfoil boundary (i.e.,  $\gamma_{\text{airfoil}}$ ) in Section 5.3 as a motivating example.

Supposing  $f(\mathbf{x})$  is only defined in  $\gamma_{\text{airfoil}}$ , our task is to extend its definition to  $\Omega \cup \partial\Omega$ . As shown in Figure A1, we first place two reference points (i.e.,  $\mathbf{x}_0$  and  $\mathbf{x}_1$ ) on the front and rear half of the airfoil. For any  $\mathbf{x} \in \Omega \cup \partial\Omega$ , it can be expressed as polar coordinates with respect to  $\mathbf{x}_0$  and  $\mathbf{x}_1$ , respectively. We concatenate the two polar coordinates to form a new space. We next perform interpolation and approximation under the new space. This is because in the new space we can better characterize the shape of the airfoil. It is true that there are many ways for coordinate transformations, not limited to the example here.

As for the interpolation, we can sample several points at the  $\gamma_{\text{airfoil}}$  to obtain the dataset  $\{((\theta_0^{(i)}, \theta_1^{(i)}), f^{(i)})\}_{i=1}^N$ . For any  $\mathbf{x} \in \Omega \cup \partial\Omega$ , we generate the corresponding extended  $f(\mathbf{x})$  by interpolating in the dataset. The interpolation method used here depends on the smoothness requirements of the ansatz. In addition, the number of reference points can also be changed, and in experiments we found that only one reference point is enough.

Approximation via neural networks is a general method that does not require manual design. In this case, we can sample several points at the  $\gamma_{\text{airfoil}}$  to construct our dataset  $\{((\phi_0(\mathbf{x}^{(i)}), \phi_1(\mathbf{x}^{(i)})), f^{(i)})\}_{i=1}^N$ , followed by training a neural network on the dataset, i.e.  $\text{NN}(\phi_0(\mathbf{x}^{(i)}), \phi_1(\mathbf{x}^{(i)})) \approx f^{(i)}$ . For any  $\mathbf{x} \in \Omega \cup \partial\Omega$ , we take  $\text{NN}(\phi_0(\mathbf{x}), \phi_1(\mathbf{x}))$  as the corresponding extended  $f(\mathbf{x})$ . We can also train the neural network in the original space. However, experimental results show that training on the new space can achieve better results. The reason may be that the complex geometry become smoother and easier to learn in the new space.

It is worth noting that, in addition to the cases mentioned above, the extended distance functions  $l(\mathbf{x})$  (here we omit the superscript and see Eq. (4) for its definition) may also need to be handled similarly. Because for the complex geometry, the distance function can be very complex and we may want to replace it with a cheap surrogate model. The methods are similar, including approximating the distance function with a neural network, or constructing splines function [5].

### A.6 The Hard-Constraint Framework for Time-dependent PDEs

In this section, we consider the following time-dependent PDEs

$$\mathcal{F}[\mathbf{u}(\mathbf{x}, t)] = \mathbf{0}, \quad \mathbf{x} = (x_1, \dots, x_d) \in \Omega, t \in (0, T], \quad (\text{A26})$$

where  $t$  is the temporal coordinate, and the other notations are the same as those in Section 3.1. For each  $u_j, j = 1, \dots, n$ , we pose suitable boundary conditions (BCs)

$$a_i(\mathbf{x}, t)u_j + b_i(\mathbf{x}, t)(\mathbf{n}(\mathbf{x}) \cdot \nabla u_j) = g_i(\mathbf{x}, t), \quad \mathbf{x} \in \gamma_i, t \in (0, T], \quad \forall i = 1, \dots, m_j, \quad (\text{A27})$$

and an initial condition (IC)

$$u_j(\mathbf{x}, 0) = f_j(\mathbf{x}), \quad \mathbf{x} \in \Omega. \quad (\text{A28})$$

Following the pipeline described in Section 3.2, we can find the general solution  $\tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}, t)$  as

$$\tilde{\mathbf{p}}_j^{\gamma_i} = \mathbf{B}(\mathbf{x}, t) \text{NN}_j^{\gamma_i}(\mathbf{x}, t) + \tilde{\mathbf{n}}(\mathbf{x}, t) \tilde{g}_i(\mathbf{x}, t), \quad (\text{A29})$$

where  $\tilde{\mathbf{n}} = (a_i, b_i \mathbf{n}) / \sqrt{a_i^2 + b_i^2}$ ,  $\tilde{g}_i = g_i / \sqrt{a_i^2 + b_i^2}$ ,  $\text{NN}_j^{\gamma_i} : \mathbb{R}^{d+1} \rightarrow \mathbb{R}^{d+1}$  is a neural network, and  $\mathbf{B}(\mathbf{x}, t) = \mathbf{I}_{d+1} - \tilde{\mathbf{n}}(\mathbf{x}, t) \tilde{\mathbf{n}}(\mathbf{x}, t)^\top$ . And we omit the trainable parameters of neural networks for neatness.

Finally, we can construct our ansatz  $(\hat{u}_j, \hat{\mathbf{p}}_j)$  as

$$(u_j^\dagger, \hat{\mathbf{p}}_j) = l^{\partial\Omega}(\mathbf{x}) \text{NN}_{\text{main}}(\mathbf{x}, t) + \sum_{i=1}^{m_j} \exp[-\alpha_i l^{\gamma_i}(\mathbf{x})] \tilde{\mathbf{p}}_j^{\gamma_i}(\mathbf{x}, t), \quad \forall j = 1, \dots, n, \quad (\text{A30a})$$

$$\hat{u}_j = u_j^\dagger(\mathbf{x}, t) (1 - \exp[-\beta_t t]) + f_j(\mathbf{x}) \exp[-\beta_t t], \quad \forall j = 1, \dots, n, \quad (\text{A30b})$$

where  $u_j^\dagger$  is an intermediate variable that incorporates hard constraints in spatial dimensions,  $\text{NN}_{\text{main}} : \mathbb{R}^{d+1} \rightarrow \mathbb{R}^{d+1}$  is the main neural network,  $l^{\partial\Omega}, l^{\gamma_i}, i = 1, \dots, m_j$  are extended distance functions (see Eq. (4)),  $\alpha_i$  ( $i = 1, \dots, m_j$ ) is determined in Eq. (12), and  $\beta_t \in \mathbb{R}$  is a hyper-parameter of the ‘‘hardness’’ in the temporal domain.

## A.7 Supplements to the Theoretical Analysis

In this section, we first give some supplements to the problem setting in Section 4. Then we present the proof of Theorem 4.1. Finally, we will characterize the mechanism described in Section 4 with another tool, the condition number.

### A.7.1 Supplements to the Problem Setting

As mentioned in Section 4, we consider the following 1D Poisson’s equation

$$\Delta u(x) = -a^2 \sin ax, \quad x \in (0, 2\pi), \quad (\text{A31a})$$

$$u(x) = 0, \quad x = 0 \vee x = 2\pi, \quad (\text{A31b})$$

where  $a \in \mathbb{R}$  and  $u$  is the physical quantity of interest. Here we use a single-layer neural network of width  $K$  as our ansatz, i.e.,  $\hat{u} = \mathbf{c}^\top \sigma(\mathbf{w}x + \mathbf{b})$ , where  $\mathbf{c}, \mathbf{w}, \mathbf{b} \in \mathbb{R}^K$ ,  $\sigma$  is an element-wise activation function (for simplicity, we take  $\sigma$  as  $\tanh$ ). To study the impact of the *extra fields* alone, we train  $\hat{u}$  in a soft-constrained manner. For ease of discussion, we consider the loss function in continuous form

$$\mathcal{L}(\boldsymbol{\theta}) = \mathcal{L}_{\mathcal{F}}(\boldsymbol{\theta}) + \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta}) \approx \frac{1}{2\pi} \int_0^{2\pi} (\Delta \hat{u}(x) + a^2 \sin(ax))^2 dx + (\hat{u}(0))^2 + (\hat{u}(2\pi))^2, \quad (\text{A32})$$

where  $\boldsymbol{\theta} = (\mathbf{c}, \mathbf{w}, \mathbf{b})$  is a set of trainable parameters.

Let  $p = \nabla u = du/dx$ . We reformulate Eq. (A31) via the *extra fields* to obtain

$$\nabla p(x) = -a^2 \sin ax, \quad x \in (0, 2\pi), \quad (\text{A33a})$$

$$p(x) = \nabla u(x), \quad x \in (0, 2\pi), \quad (\text{A33b})$$

$$u(x) = 0, \quad x = 0 \vee x = 2\pi. \quad (\text{A33c})$$

Our ansatz becomes  $\hat{u} = \mathbf{c}^\top \sigma(\mathbf{w}x + \mathbf{b})$  and  $\hat{p} = \mathbf{c}_p^\top \sigma(\mathbf{w}x + \mathbf{b})$ , where  $\mathbf{c}_p \in \mathbb{R}^K$  is a weight vector with respect to the output  $\hat{p}$ . We can see that the loss term of the BC does not change while that of the PDE becomes

$$\tilde{\mathcal{L}}_{\mathcal{F}}(\tilde{\boldsymbol{\theta}}) \approx \frac{1}{2\pi} \int_0^{2\pi} [(\nabla \hat{p}(x) + a^2 \sin(ax))^2 + (\hat{p}(x) - \nabla \hat{u}(x))^2] dx, \quad (\text{A34})$$

where  $\tilde{\boldsymbol{\theta}} = (\mathbf{c}, \mathbf{w}, \mathbf{b}, \mathbf{c}_p)$  is a set of trainable parameters.

### A.7.2 Proof of Theorem 4.1

In this part, we provide detailed proof of Theorem 4.1.

We first derive the derivatives of the ansatz for the original PDEs (we recall that  $\sigma$  is tanh and we have  $\sigma' = 1 - \sigma^2$ )

$$\frac{d\hat{u}}{dx} = \mathbf{c}^\top \left[ (\mathbf{1} - \sigma^2(\mathbf{w}x + \mathbf{b})) \circ \mathbf{w} \right], \quad (\text{A35a})$$

$$\frac{d^2\hat{u}}{dx^2} = -2\mathbf{c}^\top \left[ \sigma(\mathbf{w}x + \mathbf{b}) \circ (\mathbf{1} - \sigma^2(\mathbf{w}x + \mathbf{b})) \circ \mathbf{w}^2 \right]. \quad (\text{A35b})$$

We will abbreviate  $\sigma(\mathbf{w}x + \mathbf{b})$  as  $\boldsymbol{\sigma}$ , and then we have

$$\frac{d\hat{u}}{dx} = \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}], \quad (\text{A36a})$$

$$\frac{d^2\hat{u}}{dx^2} = -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2]. \quad (\text{A36b})$$

Now, we can provide a bound for  $(\partial\mathcal{L}_{\mathcal{F}}/\partial\mathbf{c})^\top$  as

$$\begin{aligned} \left| \left( \frac{\partial\mathcal{L}_{\mathcal{F}}}{\partial\mathbf{c}} \right)^\top \right| &= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \frac{d^2\hat{u}}{dx^2} + a^2 \sin(ax) \right) \cdot \left( \partial \left( \frac{d^2\hat{u}}{dx^2} \right) / \partial\mathbf{c} \right) dx \right| \\ &= \frac{2}{\pi} \left| \int_0^{2\pi} \left( -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] + a^2 \sin(ax) \right) \right. \\ &\quad \left. \cdot [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] dx \right| \\ &\leq \frac{2}{\pi} \left( \int_0^{2\pi} \left( -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] + a^2 \sin(ax) \right)^2 dx \right)^{\frac{1}{2}} \\ &\quad \cdot \left( \int_0^{2\pi} [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2]^2 dx \right)^{\frac{1}{2}} \quad (\text{Cauchy - Schwarz}) \\ &\leq \frac{2}{\pi} \left( \int_0^{2\pi} (2|\mathbf{c}|^\top \mathbf{w}^2 + a^2)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{w}^4 dx \right)^{\frac{1}{2}} \\ &= 4(2|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \mathbf{w}^2 \\ &\leq 8(|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \mathbf{w}^2, \end{aligned} \quad (\text{A37})$$

where  $\leq$  between two vectors is an element-wise comparison. Thus,  $(\partial\mathcal{L}_{\mathcal{F}}/\partial\mathbf{c})^\top$  can be bounded by

$$\left| \left( \frac{\partial\mathcal{L}_{\mathcal{F}}}{\partial\mathbf{c}} \right)^\top \right| = \mathcal{O}(|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot \mathbf{w}^2. \quad (\text{A38})$$

Similarly, for  $(\partial\mathcal{L}_{\mathcal{F}}/\partial\mathbf{w})^\top$ , we have

$$\begin{aligned}
\left| \left( \frac{\partial\mathcal{L}_{\mathcal{F}}}{\partial\mathbf{w}} \right)^\top \right| &= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \frac{d^2\hat{u}}{dx^2} + a^2 \sin(ax) \right) \cdot \left( \partial \left( \frac{d^2\hat{u}}{dx^2} \right) / \partial\mathbf{w} \right) dx \right| \\
&= \frac{2}{\pi} \left| \int_0^{2\pi} \left( -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] + a^2 \sin(ax) \right) \right. \\
&\quad \left. \cdot \mathbf{c} \circ [2\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w} + x(\mathbf{1} - 3\boldsymbol{\sigma}^2) \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] dx \right| \\
&\leq \frac{2}{\pi} \left( \int_0^{2\pi} \left( -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] + a^2 \sin(ax) \right)^2 dx \right)^{\frac{1}{2}} \\
&\quad \cdot \left( \int_0^{2\pi} \mathbf{c}^2 \circ [2\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w} \right. \\
&\quad \left. + x(\mathbf{1} - 3\boldsymbol{\sigma}^2) \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2]^2 dx \right)^{\frac{1}{2}} \\
&\leq \frac{2}{\pi} \left( \int_0^{2\pi} (2|\mathbf{c}|^\top \mathbf{w}^2 + a^2)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{c}^2 \circ (\mathbf{w} + 2\pi\mathbf{w}^2)^2 dx \right)^{\frac{1}{2}} \\
&= 4(2|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot |\mathbf{c}| \circ (|\mathbf{w}| + 2\pi\mathbf{w}^2) \\
&\leq 16\pi(|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot |\mathbf{c}| \circ (|\mathbf{w}| + \mathbf{w}^2).
\end{aligned} \tag{A39}$$

Thus, the bound for  $(\partial\mathcal{L}_{\mathcal{F}}/\partial\mathbf{w})^\top$  is given by

$$\left| \left( \frac{\partial\mathcal{L}_{\mathcal{F}}}{\partial\mathbf{w}} \right)^\top \right| = \mathcal{O}(|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot (|\mathbf{c}| \circ |\mathbf{w}| \circ (|\mathbf{w}| + \mathbf{1})). \tag{A40}$$

And for  $(\partial\mathcal{L}_{\mathcal{F}}/\partial\mathbf{b})^\top$ , we have

$$\begin{aligned}
\left| \left( \frac{\partial\mathcal{L}_{\mathcal{F}}}{\partial\mathbf{b}} \right)^\top \right| &= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \frac{d^2\hat{u}}{dx^2} + a^2 \sin(ax) \right) \cdot \left( \partial \left( \frac{d^2\hat{u}}{dx^2} \right) / \partial\mathbf{b} \right) dx \right| \\
&= \frac{2}{\pi} \left| \int_0^{2\pi} \left( -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] + a^2 \sin(ax) \right) \right. \\
&\quad \left. \cdot [\mathbf{c} \circ (\mathbf{1} - 3\boldsymbol{\sigma}^2) \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] dx \right| \\
&\leq \frac{2}{\pi} \left( \int_0^{2\pi} \left( -2\mathbf{c}^\top [\boldsymbol{\sigma} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2] + a^2 \sin(ax) \right)^2 dx \right)^{\frac{1}{2}} \\
&\quad \cdot \left( \int_0^{2\pi} [\mathbf{c} \circ (\mathbf{1} - 3\boldsymbol{\sigma}^2) \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}^2]^2 dx \right)^{\frac{1}{2}} \\
&\leq \frac{2}{\pi} \left( \int_0^{2\pi} (2|\mathbf{c}|^\top \mathbf{w}^2 + a^2)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{c}^2 \circ \mathbf{w}^4 dx \right)^{\frac{1}{2}} \\
&= 4(2|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot |\mathbf{c}| \circ \mathbf{w}^2 \\
&\leq 8(|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot |\mathbf{c}| \circ \mathbf{w}^2.
\end{aligned} \tag{A41}$$

Thus, the bound for  $(\partial\mathcal{L}_{\mathcal{F}}/\partial\mathbf{b})^\top$  is given by

$$\left| \left( \frac{\partial\mathcal{L}_{\mathcal{F}}}{\partial\mathbf{b}} \right)^\top \right| = \mathcal{O}(|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot (|\mathbf{c}| \circ \mathbf{w}^2). \tag{A42}$$



Recalling that  $\boldsymbol{\theta} = (\mathbf{c}, \mathbf{w}, \mathbf{b})$ , from Eq. (A38), Eq. (A40), and Eq. (A42), we have

$$\left| (\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{F}})^{\top} \right| = \left| \left( \frac{\partial \mathcal{L}_{\mathcal{F}}}{\partial \boldsymbol{\theta}} \right)^{\top} \right| = \mathcal{O}(|\mathbf{c}|^{\top} \mathbf{w}^2 + a^2) \cdot (\mathbf{w}^2, |\mathbf{c}| \circ |\mathbf{w}| \circ (|\mathbf{w}| + 1), |\mathbf{c}| \circ \mathbf{w}^2). \quad (\text{A43})$$

In contrast, for the transformed PDE, we first derive the derivatives of the ansatz (i.e.,  $\hat{u} = \mathbf{c}^{\top} \sigma(\mathbf{w}x + \mathbf{b})$ ,  $\hat{p} = \mathbf{c}_p^{\top} \sigma(\mathbf{w}x + \mathbf{b})$ )

$$\frac{d\hat{u}}{dx} = \mathbf{c}^{\top} [(1 - \sigma^2(\mathbf{w}x + \mathbf{b})) \circ \mathbf{w}], \quad (\text{A44a})$$

$$\frac{d\hat{p}}{dx} = \mathbf{c}_p^{\top} [(1 - \sigma^2(\mathbf{w}x + \mathbf{b})) \circ \mathbf{w}]. \quad (\text{A44b})$$

We again abbreviate  $\sigma(\mathbf{w}x + \mathbf{b})$  as  $\boldsymbol{\sigma}$  to obtain

$$\frac{d\hat{u}}{dx} = \mathbf{c}^{\top} [(1 - \boldsymbol{\sigma}^2) \circ \mathbf{w}], \quad (\text{A45a})$$

$$\frac{d\hat{p}}{dx} = \mathbf{c}_p^{\top} [(1 - \boldsymbol{\sigma}^2) \circ \mathbf{w}]. \quad (\text{A45b})$$

We now compute a bound for  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{c})^{\top}$

$$\begin{aligned} \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{c}} \right)^{\top} \right| &= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \hat{p} - \frac{d\hat{u}}{dx} \right) \cdot \left( \partial \left( \frac{d\hat{u}}{dx} \right) / \partial \mathbf{c} \right) dx \right| \\ &= \frac{1}{\pi} \left| \int_0^{2\pi} \left( \mathbf{c}_p^{\top} \boldsymbol{\sigma} - \mathbf{c}^{\top} [(1 - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right) \cdot [(1 - \boldsymbol{\sigma}^2) \circ \mathbf{w}] dx \right| \\ &\leq \frac{1}{\pi} \left( \int_0^{2\pi} \left( \mathbf{c}_p^{\top} \boldsymbol{\sigma} - \mathbf{c}^{\top} [(1 - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} [(1 - \boldsymbol{\sigma}^2) \circ \mathbf{w}]^2 dx \right)^{\frac{1}{2}} \\ &\leq \frac{1}{\pi} \left( \int_0^{2\pi} (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^{\top} |\mathbf{w}|)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{w}^2 dx \right)^{\frac{1}{2}} \\ &= 2(\|\mathbf{c}_p\|_1 + |\mathbf{c}|^{\top} |\mathbf{w}|) |\mathbf{w}|, \end{aligned} \quad (\text{A46})$$

Thus, the bound for  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{c})^{\top}$  is given by

$$\left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{c}} \right)^{\top} \right| = \mathcal{O}(\|\mathbf{c}_p\|_1 + |\mathbf{c}|^{\top} |\mathbf{w}|) \cdot |\mathbf{w}|. \quad (\text{A47})$$

As for  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{w})^\top$ , we have

$$\begin{aligned}
& \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{w}} \right)^\top \right| \\
&= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \frac{d\hat{p}}{dx} + a^2 \sin(ax) \right) \cdot \left( \partial \left( \frac{d\hat{p}}{dx} \right) / \partial \mathbf{w} \right) dx + \int_0^{2\pi} 2 \left( \hat{p} - \frac{d\hat{u}}{dx} \right) \cdot \left( \frac{\partial \hat{p}}{\partial \mathbf{w}} - \partial \left( \frac{d\hat{u}}{dx} \right) / \partial \mathbf{w} \right) dx \right| \\
&= \frac{1}{\pi} \left| \int_0^{2\pi} \left( \mathbf{c}_p^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] + a^2 \sin(ax) \right) \cdot [\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ (\mathbf{1} - 2x\boldsymbol{\sigma} \circ \mathbf{w})] dx \right. \\
&\quad \left. + \int_0^{2\pi} \left( \mathbf{c}_p^\top \boldsymbol{\sigma} - \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right) \cdot [x\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) - \mathbf{c} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ (\mathbf{1} - 2x\boldsymbol{\sigma} \circ \mathbf{w})] dx \right| \\
&\leq \frac{1}{\pi} \left( \left( \int_0^{2\pi} \left( \mathbf{c}_p^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] + a^2 \sin(ax) \right)^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \cdot \left( \int_0^{2\pi} [\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ (\mathbf{1} - 2x\boldsymbol{\sigma} \circ \mathbf{w})]^2 dx \right)^{\frac{1}{2}} \\
&\quad \left. + \left( \int_0^{2\pi} \left( \mathbf{c}_p^\top \boldsymbol{\sigma} - \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right)^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \left. \cdot \left( \int_0^{2\pi} [x\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) - \mathbf{c} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ (\mathbf{1} - 2x\boldsymbol{\sigma} \circ \mathbf{w})]^2 dx \right)^{\frac{1}{2}} \right) \\
&\leq 4 \left( \left( \int_0^{2\pi} (|\mathbf{c}_p|^\top |\mathbf{w}| + a^2)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{c}_p^2 \circ (\mathbf{1} + |\mathbf{w}|)^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \left. + \left( \int_0^{2\pi} (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}|)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} [|\mathbf{c}_p| + |\mathbf{c}| \circ (\mathbf{1} + |\mathbf{w}|)]^2 dx \right)^{\frac{1}{2}} \right) \\
&= 8\pi \left( (|\mathbf{c}_p|^\top |\mathbf{w}| + a^2) \cdot [|\mathbf{c}_p| \circ (\mathbf{1} + |\mathbf{w}|)] + (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}|) \cdot [|\mathbf{c}_p| + |\mathbf{c}| \circ (\mathbf{1} + |\mathbf{w}|)] \right) \\
&\leq 40\pi (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot [\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})].
\end{aligned} \tag{A48}$$

Thus, we can bound  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{w})^\top$  by

$$\begin{aligned}
\left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{w}} \right)^\top \right| &\leq 40\pi (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot [\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})] \\
&= \mathcal{O}(\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot [\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})].
\end{aligned} \tag{A49}$$

And for  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{b})^\top$ , we have

$$\begin{aligned}
& \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{b}} \right)^\top \right| \\
&= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \frac{d\hat{p}}{dx} + a^2 \sin(ax) \right) \cdot \left( \partial \left( \frac{d\hat{p}}{dx} \right) / \partial \mathbf{b} \right) dx \right. \\
&\quad \left. + \int_0^{2\pi} 2 \left( \hat{p} - \frac{d\hat{u}}{dx} \right) \cdot \left( \frac{\partial \hat{p}}{\partial \mathbf{b}} - \partial \left( \frac{d\hat{u}}{dx} \right) / \partial \mathbf{b} \right) dx \right| \\
&= \frac{1}{\pi} \left| \int_0^{2\pi} \left( \mathbf{c}_p^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] + a^2 \sin(ax) \right) \cdot [-2\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \boldsymbol{\sigma} \circ \mathbf{w}] dx \right. \\
&\quad \left. + \int_0^{2\pi} \left( \mathbf{c}_p^\top \boldsymbol{\sigma} - \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right) \cdot [\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) + 2\mathbf{c} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \boldsymbol{\sigma} \circ \mathbf{w}] dx \right| \\
&\leq \frac{1}{\pi} \left( \left( \int_0^{2\pi} \left( \mathbf{c}_p^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] + a^2 \sin(ax) \right)^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \cdot \left( \int_0^{2\pi} [-2\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \boldsymbol{\sigma} \circ \mathbf{w}]^2 dx \right)^{\frac{1}{2}} \\
&\quad + \left( \int_0^{2\pi} \left( \mathbf{c}_p^\top \boldsymbol{\sigma} - \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right)^2 dx \right)^{\frac{1}{2}} \\
&\quad \cdot \left( \int_0^{2\pi} [\mathbf{c}_p \circ (\mathbf{1} - \boldsymbol{\sigma}^2) + 2\mathbf{c} \circ (\mathbf{1} - \boldsymbol{\sigma}^2) \circ \boldsymbol{\sigma} \circ \mathbf{w}]^2 dx \right)^{\frac{1}{2}} \Bigg) \\
&\leq \frac{2}{\pi} \left( \left( \int_0^{2\pi} (|\mathbf{c}_p|^\top |\mathbf{w}| + a^2)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{c}_p^2 \circ \mathbf{w}^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \left. + \left( \int_0^{2\pi} (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}|)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} [|\mathbf{c}_p| + |\mathbf{c}| \circ |\mathbf{w}|]^2 dx \right)^{\frac{1}{2}} \right) \\
&= 4 \left( (|\mathbf{c}_p|^\top |\mathbf{w}| + a^2) \cdot [|\mathbf{c}_p| \circ |\mathbf{w}|] + (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}|) \cdot [|\mathbf{c}_p| + |\mathbf{c}| \circ |\mathbf{w}|] \right) \\
&\leq 12 (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot [\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})].
\end{aligned} \tag{A50}$$

Thus, we can bound  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{b})^\top$  by

$$\begin{aligned}
& \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{b}} \right)^\top \right| \leq 12 (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot [\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})] \\
&= \mathcal{O}(\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot [\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})].
\end{aligned} \tag{A51}$$

For  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{c}_p)^\top$ , we have

$$\begin{aligned}
& \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{c}_p} \right)^\top \right| \\
&= \frac{1}{2\pi} \left| \int_0^{2\pi} 2 \left( \frac{d\hat{p}}{dx} + a^2 \sin(ax) \right) \cdot \left( \partial \left( \frac{d\hat{p}}{dx} \right) / \partial \mathbf{c}_p \right) dx + \int_0^{2\pi} 2 \left( \hat{p} - \frac{d\hat{u}}{dx} \right) \cdot \left( \frac{\partial \hat{p}}{\partial \mathbf{c}_p} \right) dx \right| \\
&= \frac{1}{\pi} \left| \int_0^{2\pi} \left( \mathbf{c}_p^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] + a^2 \sin(ax) \right) \cdot [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] dx \right. \\
&\quad \left. + \int_0^{2\pi} \left( \mathbf{c}_p^\top \boldsymbol{\sigma} - \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right) \cdot \boldsymbol{\sigma} dx \right| \\
&\leq \frac{1}{\pi} \left( \left( \int_0^{2\pi} \left( \mathbf{c}_p^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] + a^2 \sin(ax) \right)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}]^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \left. + \left( \int_0^{2\pi} \left( \mathbf{c}_p^\top \boldsymbol{\sigma} - \mathbf{c}^\top [(\mathbf{1} - \boldsymbol{\sigma}^2) \circ \mathbf{w}] \right)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \boldsymbol{\sigma}^2 dx \right)^{\frac{1}{2}} \right) \\
&\leq \frac{1}{\pi} \left( \left( \int_0^{2\pi} (|\mathbf{c}_p|^\top |\mathbf{w}| + a^2)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{w}^2 dx \right)^{\frac{1}{2}} \right. \\
&\quad \left. + \left( \int_0^{2\pi} (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}|)^2 dx \right)^{\frac{1}{2}} \cdot \left( \int_0^{2\pi} \mathbf{1} dx \right)^{\frac{1}{2}} \right) \\
&= 2 \left( (|\mathbf{c}_p|^\top |\mathbf{w}| + a^2) \cdot |\mathbf{w}| + (\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}|) \cdot \mathbf{1} \right) \\
&\leq 4 (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot \max(|\mathbf{w}|, \mathbf{1}).
\end{aligned} \tag{A52}$$

Thus, we can bound  $(\partial \tilde{\mathcal{L}}_{\mathcal{F}} / \partial \mathbf{c}_p)^\top$  by

$$\begin{aligned}
& \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \mathbf{c}_p} \right)^\top \right| \leq 4 (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot \max(|\mathbf{w}|, \mathbf{1}) \\
&= \mathcal{O}(\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot \max(|\mathbf{w}|, \mathbf{1}).
\end{aligned} \tag{A53}$$

Recalling that  $\tilde{\boldsymbol{\theta}} = (\mathbf{c}, \mathbf{w}, \mathbf{b}, \mathbf{c}_p)$ , from Eq. (A47), Eq. (A49), Eq. (A51), and Eq. (A53), noting that  $\|\mathbf{c}_p\|_1 + |\mathbf{c}|^\top |\mathbf{w}| \leq \|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2$ , we have

$$\begin{aligned}
& \left| (\nabla_{\tilde{\boldsymbol{\theta}}} \tilde{\mathcal{L}}_{\mathcal{F}})^\top \right| = \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \tilde{\boldsymbol{\theta}}} \right)^\top \right| = \mathcal{O}(\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot (|\mathbf{w}|, \\
& \quad \max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1}), \max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1}), \max(|\mathbf{w}|, \mathbf{1})).
\end{aligned} \tag{A54}$$

### A.7.3 Analysis via the Condition Number

In addition to providing bounds for the gradients of  $\mathcal{L}_{\mathcal{F}}$  and  $\tilde{\mathcal{L}}_{\mathcal{F}}$ , we can also use the condition number to characterize the sensitivity of their gradients with respect to the parameters of the neural network. Let  $\boldsymbol{\theta}^{(t)} = (\mathbf{c}^{(t)}, \mathbf{w}^{(t)}, \mathbf{b}^{(t)})$  and  $\tilde{\boldsymbol{\theta}}^{(t)} = (\mathbf{c}^{(t)}, \mathbf{w}^{(t)}, \mathbf{b}^{(t)}, \mathbf{c}_p^{(t)})$  are the parameters of the neural network in the  $t$ th step (before and after the reformulation). For simplicity, we introduce the following notations

$$\Delta \boldsymbol{\theta} = \boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}, \tag{A55a}$$

$$\Delta \tilde{\boldsymbol{\theta}} = \tilde{\boldsymbol{\theta}}^{(t+1)} - \tilde{\boldsymbol{\theta}}^{(t)}, \tag{A55b}$$

$$\Delta \mathcal{L}_{\mathcal{F}} = \mathcal{L}_{\mathcal{F}}(\boldsymbol{\theta}^{(t+1)}) - \mathcal{L}_{\mathcal{F}}(\boldsymbol{\theta}^{(t)}), \tag{A55c}$$

$$\Delta \tilde{\mathcal{L}}_{\mathcal{F}} = \tilde{\mathcal{L}}_{\mathcal{F}}(\tilde{\boldsymbol{\theta}}^{(t+1)}) - \tilde{\mathcal{L}}_{\mathcal{F}}(\tilde{\boldsymbol{\theta}}^{(t)}). \tag{A55d}$$

The condition numbers of  $\mathcal{L}_{\mathcal{F}}$  and  $\tilde{\mathcal{L}}_{\mathcal{F}}$  are defined as

$$\text{cond} = \frac{|\Delta \mathcal{L}_{\mathcal{F}}|}{\|\Delta \boldsymbol{\theta}\|_2}, \quad \tilde{\text{cond}} = \frac{|\Delta \tilde{\mathcal{L}}_{\mathcal{F}}|}{\|\Delta \tilde{\boldsymbol{\theta}}\|_2}. \quad (\text{A56})$$

Next we derive the bounds for  $\text{cond}$  and  $\tilde{\text{cond}}$ , respectively. We first consider  $\text{cond}$

$$\begin{aligned} \text{cond} &= \frac{|\Delta \mathcal{L}_{\mathcal{F}}|}{\|\Delta \boldsymbol{\theta}\|_2} \approx \left| \left( \frac{\partial \mathcal{L}_{\mathcal{F}}}{\partial \boldsymbol{\theta}} \right)^\top \cdot \Delta \boldsymbol{\theta} \right| / \|\Delta \boldsymbol{\theta}\|_2 \leq \left\| \left( \frac{\partial \mathcal{L}_{\mathcal{F}}}{\partial \boldsymbol{\theta}} \right)^\top \right\|_2 \\ &= \mathcal{O} \left( (|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot \|\mathbf{w}^2, |\mathbf{c}| \circ |\mathbf{w}| \circ (|\mathbf{w}| + \mathbf{1}), |\mathbf{c}| \circ \mathbf{w}^2\|_2 \right) \quad (\text{Eq. (A43)}) \\ &= \mathcal{O} \left( (|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot \|\mathbf{w}^2, |\mathbf{c}| \circ |\mathbf{w}| \circ (|\mathbf{w}| + \mathbf{1}), |\mathbf{c}| \circ \mathbf{w}^2\|_1 \right) \quad (\text{Equivalence of norms}) \\ &= \mathcal{O} \left( (|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot (\|\mathbf{w}^2\|_1 + \| |\mathbf{c}| \circ |\mathbf{w}| \circ (|\mathbf{w}| + \mathbf{1}) \|_1 + \| |\mathbf{c}| \circ \mathbf{w}^2 \|_1) \right) \\ &= \mathcal{O} \left( (|\mathbf{c}|^\top \mathbf{w}^2 + a^2) \cdot \|\max(|\mathbf{c}|, \mathbf{1}) \circ |\mathbf{w}| \circ \max(|\mathbf{w}|, \mathbf{1})\|_1 \right). \end{aligned} \quad (\text{A57})$$

Similarly, for  $\tilde{\text{cond}}$ , we have

$$\begin{aligned} \tilde{\text{cond}} &= \frac{|\Delta \tilde{\mathcal{L}}_{\mathcal{F}}|}{\|\Delta \tilde{\boldsymbol{\theta}}\|_2} \approx \left| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \tilde{\boldsymbol{\theta}}} \right)^\top \cdot \Delta \tilde{\boldsymbol{\theta}} \right| / \|\Delta \tilde{\boldsymbol{\theta}}\|_2 \leq \left\| \left( \frac{\partial \tilde{\mathcal{L}}_{\mathcal{F}}}{\partial \tilde{\boldsymbol{\theta}}} \right)^\top \right\|_2 \\ &= \mathcal{O} \left( (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot \|\mathbf{w}, \max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1}), \right. \\ &\quad \left. \max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1}), \max(|\mathbf{w}|, \mathbf{1})\|_2 \right) \quad (\text{Eq. (A54)}) \\ &= \mathcal{O} \left( (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot \|\mathbf{w}, \max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1}), \right. \\ &\quad \left. \max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1}), \max(|\mathbf{w}|, \mathbf{1})\|_1 \right) \quad (\text{Equivalence of norms}) \\ &= \mathcal{O} \left( (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot (\|\mathbf{w}\|_1 \right. \\ &\quad \left. + \|\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})\|_1 + \|\max(|\mathbf{c}|, |\mathbf{c}_p|) \circ \max(|\mathbf{w}|, \mathbf{1})\|_1 \right. \\ &\quad \left. + \|\max(|\mathbf{w}|, \mathbf{1})\|_1) \right) \\ &= \mathcal{O} \left( (\|\mathbf{c}_p\|_1 + \max(|\mathbf{c}|, |\mathbf{c}_p|)^\top |\mathbf{w}| + a^2) \cdot \|\max(|\mathbf{c}|, |\mathbf{c}_p|, \mathbf{1}) \circ \max(|\mathbf{w}|, \mathbf{1})\|_1 \right). \end{aligned} \quad (\text{A58})$$

From Eq. (A57) and Eq. (A58), we find that for the original PDEs, the condition number has a higher order relationship with respect to  $\boldsymbol{\theta}$ . If  $\boldsymbol{\theta}$  is large, the condition number can be very large, leading to oscillations in training. However, if  $\boldsymbol{\theta}$  is small, the condition number will also be very small, resulting in smaller changes of  $\boldsymbol{\theta}$  between adjacent iterations and therefore slower convergence. In contrast, after the reformulation, the condition number has a lower order relationship with respect to  $\tilde{\boldsymbol{\theta}}$ , which keeps the condition number more stable during training and alleviates this problem.

## A.8 Experimental details

In the following, we will briefly introduce some essential details of our experiments, including the experimental environment, hyper-parameters, construction of the ansatz, and details of the governing PDEs in each experiment. We first introduce the experimental environment, while other details are put into the subsections corresponding to the experiments.

**Experimental environment** We use PyTorch [4] as our deep learning library. And our codes for the physics-informed learning are based on DeepXDE [3]. We train all the models except domain decomposition based baselines (i.e., xPINN, FBPINN, and PFNN-2) on one NVIDIA TITAN Xp 12GB GPU, while the other three are trained on eight NVIDIA GeForce RTX 3090 24GB GPUs (since domain decomposition based models consist of several sub-networks and require more memory to be stored). The operating system is Ubuntu 18.04.5 LTS. If the analytical solution is unavailable, the ground truth solutions to the PDEs (i.e., the testing data) will be generated by COMSOL Multiphysics, a FEM commercial software. And we have put the generated testing data into the zip file.

### A.8.1 Simulation of a 2D battery pack (Heat Equation)

**Governing PDEs** The governing PDEs (along with boundary/initial conditions) are given by

$$\frac{\partial T}{\partial t} = k \Delta T(\mathbf{x}, t), \quad \mathbf{x} \in \Omega, t \in (0, 1], \quad (\text{A59a})$$

$$k(\mathbf{n}(\mathbf{x}) \cdot \nabla T(\mathbf{x}, t)) = h(T_a - T(\mathbf{x}, t)), \quad \mathbf{x} \in \gamma_{\text{outer}}, t \in (0, 1], \quad (\text{A59b})$$

$$k(\mathbf{n}(\mathbf{x}) \cdot \nabla T(\mathbf{x}, t)) = h(T_c - T(\mathbf{x}, t)), \quad \mathbf{x} \in \gamma_{\text{cell}, i}, t \in (0, 1], \quad i = 1, \dots, n_c, \quad (\text{A59c})$$

$$k(\mathbf{n}(\mathbf{x}) \cdot \nabla T(\mathbf{x}, t)) = h(T_w - T(\mathbf{x}, t)), \quad \mathbf{x} \in \gamma_{\text{pipe}, i}, t \in (0, 1], \quad i = 1, \dots, n_w, \quad (\text{A59d})$$

$$T(\mathbf{x}, 0) = T_0, \quad \mathbf{x} \in \Omega, \quad (\text{A59e})$$

where  $\mathbf{x} = (x_1, x_2)$ ,  $t$  are the spatial and temporal coordinates, respectively,  $T(\mathbf{x}, t)$  is the temperature over time,  $k = 1$  is the thermal conductivity,  $\Delta T = \partial^2 T / \partial x_1^2 + \partial^2 T / \partial x_2^2$ ,  $h = 1$  is the heat transfer coefficient,  $\nabla T = (\partial T / \partial x_1, \partial T / \partial x_2)$ ,  $T_a = 0.1$ ,  $T_c = 5$ ,  $T_w = 1$  are, respectively, the temperature of the air, the cells ( $n_c = 11$  cells of radius  $r_c = 1$ ), the cooling pipes ( $n_w = 6$  pipes of radius  $r_w = 0.4$ ),  $T_0 = 0.1$  is the initial temperature, and the geometry (i.e.,  $\Omega$ ,  $\gamma_{\text{outer}}$ , etc) is shown in Figure 3(a).

And the reformulated PDEs are (which is used by the proposed model, HC)

$$\frac{\partial T}{\partial t} = k(\nabla \cdot \mathbf{p}(\mathbf{x}, t)), \quad \mathbf{x} \in \Omega, t \in (0, 1], \quad (\text{A60a})$$

$$\mathbf{p}(\mathbf{x}, t) = \nabla T, \quad \mathbf{x} \in \Omega \cup \partial\Omega, t \in (0, 1], \quad (\text{A60b})$$

$$k(\mathbf{n}(\mathbf{x}) \cdot \mathbf{p}(\mathbf{x}, t)) = h(T_a - T(\mathbf{x}, t)), \quad \mathbf{x} \in \gamma_{\text{outer}}, t \in (0, 1], \quad (\text{A60c})$$

$$k(\mathbf{n}(\mathbf{x}) \cdot \mathbf{p}(\mathbf{x}, t)) = h(T_c - T(\mathbf{x}, t)), \quad \mathbf{x} \in \gamma_{\text{cell}, i}, t \in (0, 1], \quad i = 1, \dots, n_c, \quad (\text{A60d})$$

$$k(\mathbf{n}(\mathbf{x}) \cdot \mathbf{p}(\mathbf{x}, t)) = h(T_w - T(\mathbf{x}, t)), \quad \mathbf{x} \in \gamma_{\text{pipe}, i}, t \in (0, 1], \quad i = 1, \dots, n_w, \quad (\text{A60e})$$

$$T(\mathbf{x}, 0) = T_0, \quad \mathbf{x} \in \Omega, \quad (\text{A60f})$$

where  $\mathbf{p}(\mathbf{x}, t)$  is the introduced extra field.

**Construction of the Ansatz** Since the solution is a scalar function (i.e.,  $T(\mathbf{x}, t)$ ), we directly denote the solution by  $u(\mathbf{x}, t) = T(\mathbf{x}, t)$ . Let  $\tilde{\mathbf{p}}$  denote  $(u, \mathbf{p})$ . We first derive the general solutions at  $\gamma_{\text{outer}}, \gamma_{\text{cell}, i}, \gamma_{\text{pipe}, i}$ , respectively.

For  $\mathbf{x} \in \gamma_{\text{outer}}$ , we have  $a(\mathbf{x}) = h$ ,  $b(\mathbf{x}) = k$ ,  $g(\mathbf{x}) = hT_a$ . According to Eq. (10), the general solution  $\tilde{\mathbf{p}}^{\gamma_{\text{outer}}}$  is given by

$$\tilde{\mathbf{p}}^{\gamma_{\text{outer}}} = \mathbf{B}(\mathbf{x}) \mathbf{N} \mathbf{N}^{\gamma_{\text{outer}}}(\mathbf{x}, t) + \frac{(h, k\mathbf{n})}{\sqrt{h^2 + k^2}} \frac{hT_a}{\sqrt{h^2 + k^2}}, \quad (\text{A61})$$

where  $\mathbf{B}(\mathbf{x})$  is computed in Eq. (A7) (with  $\tilde{\mathbf{n}} = (\tilde{n}_1, \tilde{n}_2, \tilde{n}_3) = (h, k\mathbf{n}) / \sqrt{h^2 + k^2}$ ). And for  $\mathbf{x} \in \gamma_{\text{cell}, i}$  and  $\gamma_{\text{pipe}, i}$ , the derivation is similar, where we only need to change  $T_a$  to  $T_c$  and  $T_w$ , respectively.

Then, we gather all the general solutions computed to form our ansatz  $(\hat{u}, \hat{\mathbf{p}})$  according to Eq. (A30), where  $\{\gamma_{\text{outer}}, \gamma_{\text{cell}, 1}, \dots, \gamma_{\text{cell}, n_c}, \gamma_{\text{pipe}, 1}, \dots, \gamma_{\text{pipe}, n_w}\}$  are reordered as  $\{\gamma_i\}_{i=1}^{1+n_c+n_w}$  and  $f(\mathbf{x}) = T_0$ .

**Choices of Extended Distance Functions** For  $\gamma_{\text{cell}, i}$  and  $\gamma_{\text{pipe}, i}$ , since they are 2D circles, we can directly choose the extended distance functions  $l^{\gamma_{\text{cell}, i}}(\mathbf{x})$  and  $l^{\gamma_{\text{pipe}, i}}(\mathbf{x})$  as the distance between  $\mathbf{x}$  and the center minus the radius. For the rectangular  $\gamma_{\text{outer}}$ , supposing that it is given by  $[a_1, a_2] \times [b_1, b_2]$ , we construct the extended distance function  $l^{\gamma_{\text{outer}}}$  as follows

$$l^{\gamma_{\text{outer}}}(\mathbf{x}) = \text{SoftMin}(x_1 - a_1, a_2 - x_1, x_2 - b_1, b_2 - x_2), \quad (\text{A62})$$

where  $\mathbf{x} = (x_1, x_2)$ ,  $\text{SoftMin}$  is a differentiable version of  $\min$  function which is implemented by  $\text{LogSumExp}$  in PyTorch, i.e.,  $\text{SoftMin}(\mathbf{y}) = \text{LogSumExp}(-\beta \mathbf{y}) / (-\beta)$ ,  $\beta = 4$ . And the extended function  $l^{\partial\Omega}(\mathbf{x})$  is computed by taking the  $\text{SoftMin}$  of the distances to all the boundaries

$$l^{\partial\Omega}(\mathbf{x}) = \text{SoftMin}(l^{\gamma_{\text{outer}}}(\mathbf{x}), l^{\gamma_{\text{cell}, 1}}(\mathbf{x}), \dots, l^{\gamma_{\text{cell}, n_c}}(\mathbf{x}), l^{\gamma_{\text{pipe}, 1}}(\mathbf{x}), \dots, l^{\gamma_{\text{pipe}, n_w}}(\mathbf{x})). \quad (\text{A63})$$

**Implementation** All the models are trained for 5000 Adam iterations (with a learning rate scheduler of ReduceLRonPlateau from PyTorch and an initial learning rate of 0.01), followed by a L-BFGS optimization until convergence. Unless otherwise specified, the mean squared error (MSE) is used for the loss function and tanh is used for the activation function. And the hyper-parameters of each model are listed as follow

- **HC**: The main neural network is a multilayer perceptron (MLP) of size  $[3] + 4 \times [50] + [3]$  (which means 3 inputs, 4 hidden layers of width 50, and 3 outputs). The sub-networks (corresponding to Eq. (A60c), Eq. (A60d), and Eq. (A60e)) are all MLPs of size  $[3] + 3 \times [20] + [3]$ . And the hyper-parameters of “hardness” are  $\beta_s = 5$  and  $\beta_t = 10$ .
- **PINN**: The ansatz is an MLP of size  $[3] + 4 \times [50] + [1]$ .
- **PINN-LA**: The weights of the loss terms corresponding to the BCs are approximated by  $\hat{\lambda}_i = \max_{\theta_n} \{|\nabla_{\theta} \mathcal{L}_r(\theta_n)|\} / |\nabla_{\theta} \lambda_i \mathcal{L}_i(\theta_n)|$ . And the parameter of the moving average is  $\alpha = 0.1$ , which is recommended by the paper [7]. Besides, the parameters of the PINN are the same as above.
- **PINN-LA-2**: In our modified version, we approximate the weights of the loss terms as  $\hat{\lambda}_i = |\nabla_{\theta} \mathcal{L}_r(\theta_n)| / |\nabla_{\theta} \lambda_i \mathcal{L}_i(\theta_n)|$ . And the parameter of the moving average is also  $\alpha = 0.1$ . Here we replace the maximum with the mean to make the weights of the loss terms more stable during the training process.
- **FBPINN**: The domain of the problem is divided into  $4 \times 6 = 24$  subdomains by a regular grid. The size of the sub-network, an MLP, corresponding to each subdomain is  $[3] + 3 \times [30] + [1]$ . And the scale factor is  $\sigma = 0.4$ , chosen so that the window function is close to zero outside the overlapping region of the subdomains.
- **xPINN**: The domain of the problem is divided into  $4 \times 6 = 24$  subdomains by a regular grid. The size of the sub-network corresponding to each subdomain is  $[3] + 3 \times [30] + [1]$ . And the loss terms of the interface condition include average solution as well as residual continuity conditions.
- **PFNN**: The PFNN considers the variational formulation of the Eq. (A59) (i.e., Ritz formulation), and embed the initial condition (IC) into its ansatz (similar to Eq. (5)). And the size of the neural network (an MLP) is  $[3] + 4 \times [50] + [1]$ .
- **PFNN-2**: The PFNN-2 replaces a single neural network with a domain decomposition based neural network on the basis of PFNN. In the original literature [6], the domain is decomposed in a hard way (like xPINN). However, in our experiments (see Table 1), we find that the performance of hard decomposition is relatively poor, which is because new loss terms are needed to maintain the continuity of the ansatz at the interfaces between the sub-domains, which further aggravates the unbalanced competition. To overcome this, we instead employ a soft domain decomposition, as in FBPINN. See the parts of PFNN and FBPINN for the values of hyper-parameters.

### A.8.2 Simulation of an Airfoil (Navier-Stokes Equations)

**Governing PDEs** The governing PDEs are given by

$$\mathbf{u}(\mathbf{x}) \cdot \nabla \mathbf{u}(\mathbf{x}) = -\nabla p(\mathbf{x}) + v \nabla^2 \mathbf{u}(\mathbf{x}), \quad \mathbf{x} \in \Omega, \quad (\text{A64a})$$

$$\nabla \cdot \mathbf{u}(\mathbf{x}) = 0, \quad \mathbf{x} \in \Omega, \quad (\text{A64b})$$

$$\mathbf{u}(\mathbf{x}) = \mathbf{u}_0(\mathbf{x}), \quad \mathbf{x} \in \gamma_{\text{inlet}} \cup \gamma_{\text{top}} \cup \gamma_{\text{bottom}}, \quad (\text{A64c})$$

$$p(\mathbf{x}) = 1, \quad \mathbf{x} \in \gamma_{\text{outlet}}, \quad (\text{A64d})$$

$$\mathbf{n}(\mathbf{x}) \cdot \mathbf{u}(\mathbf{x}) = 0, \quad \mathbf{x} \in \gamma_{\text{airfoil}}, \quad (\text{A64e})$$

where  $\mathbf{u}(\mathbf{x}) = (u_1(\mathbf{x}), u_2(\mathbf{x}))$ ,  $p(\mathbf{x})$ ,  $v = 1/50$  are the velocity, pressure, and viscosity of the fluid, respectively,  $\mathbf{u}_0(\mathbf{x}) = (1, 0)$ , and the geometry of the problem (i.e.,  $\Omega$ ,  $\gamma_{\text{inlet}}$ , etc) is shown in Figure 3(b).

And the reformulated PDEs are (which is used by the proposed model, HC)

$$(\mathbf{u}(\mathbf{x}) - v\nabla) \cdot (\mathbf{p}_1(\mathbf{x}), \mathbf{p}_2(\mathbf{x})) = -\nabla p(\mathbf{x}), \quad \mathbf{x} \in \Omega, \quad (\text{A65a})$$

$$\nabla \cdot \mathbf{u}(\mathbf{x}) = 0, \quad \mathbf{x} \in \Omega, \quad (\text{A65b})$$

$$\mathbf{p}_1(\mathbf{x}) = \nabla u_1, \quad \mathbf{x} \in \Omega \cup \partial\Omega, \quad (\text{A65c})$$

$$\mathbf{p}_2(\mathbf{x}) = \nabla u_2, \quad \mathbf{x} \in \Omega \cup \partial\Omega, \quad (\text{A65d})$$

$$\mathbf{u}(\mathbf{x}) = \mathbf{u}_0(\mathbf{x}), \quad \mathbf{x} \in \gamma_{\text{inlet}} \cup \gamma_{\text{top}} \cup \gamma_{\text{bottom}}, \quad (\text{A65e})$$

$$p(\mathbf{x}) = 1, \quad \mathbf{x} \in \gamma_{\text{outlet}}, \quad (\text{A65f})$$

$$\mathbf{n}(\mathbf{x}) \cdot \mathbf{u}(\mathbf{x}) = 0, \quad \mathbf{x} \in \gamma_{\text{airfoil}}, \quad (\text{A65g})$$

where  $\mathbf{p}_1(\mathbf{x})$  and  $\mathbf{p}_2(\mathbf{x})$  are the introduced extra fields.

**Construction of the Ansatz** Here, the solution is  $(\mathbf{u}(\mathbf{x}), p(\mathbf{x}))$ . For  $p(\mathbf{x})$ , the general solution in  $\gamma_{\text{outlet}}$  is exactly  $p^{\gamma_{\text{outlet}}}(\mathbf{x}) = 1$ . And the ansatz for  $p$  is given by

$$\hat{p} = p^{\gamma_{\text{outlet}}}(\mathbf{x}) + l^{\gamma_{\text{outlet}}}(\mathbf{x})\text{NN}_{\text{main}}(\mathbf{x})[3], \quad (\text{A66})$$

where [3] means taking the third elements of the output of  $\text{NN}_{\text{main}}(\mathbf{x})$ .

For  $\mathbf{u}(\mathbf{x})$ , the general solution in  $\gamma_{\text{inlet}} \cup \gamma_{\text{top}} \cup \gamma_{\text{bottom}}$  is exactly  $\mathbf{u}^{\gamma_*}(\mathbf{x}) = \mathbf{u}_0(\mathbf{x})$ , where we define an alias  $\gamma_*$  for  $\gamma_{\text{inlet}} \cup \gamma_{\text{top}} \cup \gamma_{\text{bottom}}$ . In  $\gamma_{\text{airfoil}}$ , the general solution is given by

$$\mathbf{u}^{\gamma_{\text{airfoil}}} = \mathbf{B}(\mathbf{x})\text{NN}^{\gamma_*}(\mathbf{x}), \quad (\text{A67})$$

where  $\mathbf{B}(\mathbf{x}) = [n_2(\mathbf{x}), -n_1(\mathbf{x})]^\top$  according to Eq. (A5) and the output of  $\text{NN}^{\gamma_*}(\mathbf{x})$  is a scalar. Gathering  $\mathbf{u}^{\gamma_*}$  and  $\mathbf{u}^{\gamma_{\text{airfoil}}}$ , we then follow Eq. (11) to obtain the ansatz for  $\hat{\mathbf{u}}$

$$\hat{\mathbf{u}} = l^{\partial\Omega}(\mathbf{x})\text{NN}_{\text{main}}(\mathbf{x})[1 : 2] + \exp[-\alpha_{\gamma_*} l^{\gamma_*}(\mathbf{x})]\mathbf{u}^{\gamma_*}(\mathbf{x}) + \exp[-\alpha_{\gamma_{\text{airfoil}}} l^{\gamma_{\text{airfoil}}}(\mathbf{x})]\mathbf{u}^{\gamma_{\text{airfoil}}}(\mathbf{x}), \quad (\text{A68})$$

where [1 : 2] means taking the first two elements of the output of  $\text{NN}_{\text{main}}(\mathbf{x})$  and  $\alpha_{\gamma_*}$  as well as  $\alpha_{\gamma_{\text{airfoil}}}$  are similarly defined as in Eq. (12).

**Choices of Extended Distance Functions** For  $l^{\gamma_{\text{airfoil}}}(\mathbf{x})$ , a direct way is to calculate the distance between  $\mathbf{x}$  and the airfoil  $\gamma_{\text{airfoil}}$ . However, it may be very time consuming since the  $\gamma_{\text{airfoil}}$  is highly complicated. So we prefer to approximate the true distance with an MLP with 3 hidden layers of width 30. We train the neural network before training our main model with  $1024 \times 6$  points sampled in  $\Omega$  (5/6 of them are sampled in the bounding box of the airfoil, and the rest are sampled in  $\Omega$ ) along with their truth distances (which are computed by using the formula of the distance to a polygon) for 10,000 Adam epochs (with a learning rate scheduler of ReduceLR0nPlateau from PyTorch and an initial learning rate of 0.001). The loss function is a  $\ell_1$  loss. A polar coordinate transformation trick is utilized as in Appendix A.5.

And for  $\gamma_{\text{outlet}}$ , since it is a vertical line, for example,  $x_1 = a$ , we can compute the extended distance function as  $l^{\gamma_{\text{outlet}}}(\mathbf{x}) = a - x_1$ , where  $\mathbf{x} = (x_1, x_2)$ . And  $\gamma_*$  is an open rectangle, so we can compute  $l^{\gamma_*}(\mathbf{x})$  similarly to the case of the rectangle (see Eq. (A62)) while ignoring the right side. Besides,  $l^{\partial\Omega}(\mathbf{x})$  is still computed by taking the SoftMin of the distances to all the boundaries.

**Implementation** All the models are trained for 5000 Adam iterations (with a learning rate scheduler of ReduceLR0nPlateau from PyTorch and an initial learning rate of 0.001), followed by a L-BFGS optimization until convergence. And the hyper-parameters of each model are listed as follow

- **HC**: The main neural network is an MLP of size  $[2] + 6 \times [50] + [7]$ . The sub-network (corresponding to Eq. (A65g)) is an MLP of size  $[2] + 4 \times [40] + [1]$ . And the hyper-parameters of “hardness” is  $\beta_s = 5$ .
- **PINN**: The ansatz is an MLP of size  $[2] + 6 \times [50] + [3]$ .
- **PINN-LA**: The parameter of the moving average is  $\alpha = 0.1$ .
- **PINN-LA-2**: The parameter of the moving average is  $\alpha = 0.1$ .
- **FBPINN**: The domain of the problem is divided into  $3 \times 6 = 18$  subdomains by a regular grid. The size of the sub-network, an MLP, corresponding to each subdomain is  $[2] + 4 \times [30] + [3]$ . And the scale factor is  $\sigma = 0.2$ , chosen so that the window function is close to zero outside the overlapping region of the subdomains.



- **xPINN**: The domain of the problem is divided into  $3 \times 6 = 18$  subdomains by a regular grid. The size of the sub-network corresponding to each subdomain is  $[2] + 4 \times [30] + [3]$ . And the loss terms of the interface condition include average solution as well as residual continuity conditions.

### A.8.3 High-dimensional Heat Equation

**Governing PDEs** The governing PDEs are given by

$$\frac{\partial u}{\partial t} = k\Delta u(\mathbf{x}, t) + f(\mathbf{x}, t), \quad \mathbf{x} \in \Omega \subset \mathbb{R}^d, t \in (0, 1], \quad (\text{A69a})$$

$$\mathbf{n}(\mathbf{x}) \cdot \nabla u(\mathbf{x}, t) = g(\mathbf{x}, t), \quad \mathbf{x} \in \partial\Omega, t \in (0, 1], \quad (\text{A69b})$$

$$u(\mathbf{x}, 0) = g(\mathbf{x}, 0), \quad \mathbf{x} \in \Omega, \quad (\text{A69c})$$

where  $u$  is the quantity of interest,  $k = 1/d$ ,  $f(\mathbf{x}, t) = -k|\mathbf{x}|^2 \exp(0.5|\mathbf{x}|^2 + t)$ ,  $d = 10$ ,  $\Omega$  is a unit ball (i.e.,  $\Omega = \{|\mathbf{x}| \leq 1\}$ ), and  $g(\mathbf{x}, t) = \exp(0.5|\mathbf{x}|^2 + t)$  which is also the analytical solution to above PDEs.

And the reformulated PDEs are (which is used by the proposed model, HC)

$$\frac{\partial u}{\partial t} = k(\nabla \cdot \mathbf{p}(\mathbf{x}, t)) + f(\mathbf{x}, t), \quad \mathbf{x} \in \Omega \subset \mathbb{R}^d, t \in (0, 1], \quad (\text{A70a})$$

$$\mathbf{p}(\mathbf{x}, t) = \nabla u, \quad \mathbf{x} \in \Omega, t \in (0, 1], \quad (\text{A70b})$$

$$\mathbf{n}(\mathbf{x}) \cdot \mathbf{p}(\mathbf{x}, t) = g(\mathbf{x}, t), \quad \mathbf{x} \in \partial\Omega, t \in (0, 1], \quad (\text{A70c})$$

$$u(\mathbf{x}, 0) = g(\mathbf{x}, 0), \quad \mathbf{x} \in \Omega, \quad (\text{A70d})$$

where  $\mathbf{p}(\mathbf{x}, t)$  is the introduced extra field.

**Construction of the Ansatz** The solution to the PDEs is a scalar function  $u$  and there is only one boundary  $\partial\Omega = \{|\mathbf{x}| = 1\}$ . We now derive the general solution  $\mathbf{p}^{\partial\Omega}$  with respect to Eq. (A70c)

$$\mathbf{p}^{\partial\Omega} = \mathbf{B}(\mathbf{x})\text{NN}^{\partial\Omega} + \mathbf{n}(\mathbf{x})g(\mathbf{x}, t), \quad (\text{A71})$$

where  $\mathbf{B}(\mathbf{x}) = \mathbf{I}_d - \mathbf{n}(\mathbf{x})\mathbf{n}(\mathbf{x})^\top$ . And the ansatz  $(\hat{u}, \hat{\mathbf{p}})$  is given by

$$\hat{\mathbf{p}} = l^{\partial\Omega}(\mathbf{x})\text{NN}_{\text{main}}(\mathbf{x}, t)[1:d] + \mathbf{p}^{\partial\Omega}(\mathbf{x}, t), \quad (\text{A72a})$$

$$\hat{u} = \text{NN}_{\text{main}}(\mathbf{x}, t)[d+1](1 - \exp[-\beta_t t]) + g(\mathbf{x}, 0) \exp[-\beta_t t], \quad (\text{A72b})$$

where  $[1:d]$  means taking the first  $d$  elements of the output of  $\text{NN}_{\text{main}}(\mathbf{x}, t)$  while  $[d+1]$  means the last element.

**Choices of Extended Distance Functions** Since  $\partial\Omega$  is a ND sphere, we can compute  $l^{\partial\Omega}(\mathbf{x})$  by subtracting the distance between  $\mathbf{x}$  and the center from the radius (the symbol is different from the previous 2D circles, since  $\partial\Omega$  is the outer boundary).

**Implementation** All the models are trained for 5000 Adam iterations (with a learning rate scheduler of ReduceLROnPlateau from PyTorch and an initial learning rate of 0.01), followed by a L-BFGS optimization until convergence. And the hyper-parameters of each model are listed as follow

- **HC**: The main neural network is an MLP of size  $[11] + 4 \times [50] + [11]$ . The sub-network (corresponding to Eq. (A70c)) is an MLP of size  $[11] + 3 \times [20] + [10]$ . And the hyper-parameters of “hardness” is  $\beta_t = 10$  (here we only have one boundary, so we can construct our ansatz (in the spatial domain) as in Eq. (3) instead of Eq. (A30a) and  $\beta_s$  is no longer needed).
- **PINN**: The ansatz is an MLP of size  $[11] + 4 \times [50] + [1]$ .
- **PINN-LA**: The parameter of the moving average is  $\alpha = 0.1$ .
- **PINN-LA-2**: The parameter of the moving average is  $\alpha = 0.1$ .
- **PFNN**: The size of the neural network (an MLP) is  $[11] + 4 \times [50] + [1]$ .

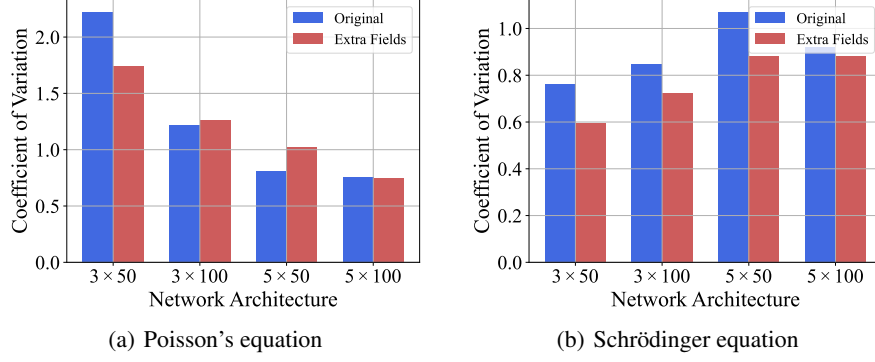


Figure A2: The CV of all the values of  $|\nabla_{\theta}\mathcal{L}_{\mathcal{F}}|$  and  $|\nabla_{\tilde{\theta}}\tilde{\mathcal{L}}_{\mathcal{F}}|$  during training.

#### A.8.4 Ablation Study: Extra fields

Here, our experiment is divided into two parts where we consider the Poisson's equation and the nonlinear Schrödinger equation, respectively. The relevant details are as follows

**Poisson's Equation** The governing PDEs are described as

$$\Delta u(x) = -a^2 \sin ax, \quad x \in (0, 2\pi), \quad (\text{A73a})$$

$$u(x) = 0, \quad x = 0 \vee x = 2\pi, \quad (\text{A73b})$$

where  $a = 2$  and  $u(x)$  is the physical quantity of interest.

And the reformulated PDEs are (corresponding to the *extra fields*)

$$\nabla p(x) = -a^2 \sin ax, \quad x \in (0, 2\pi), \quad (\text{A74a})$$

$$p(x) = \nabla u(x), \quad x \in (0, 2\pi), \quad (\text{A74b})$$

$$u(x) = 0, \quad x = 0 \vee x = 2\pi. \quad (\text{A74c})$$

where  $p(x)$  is the introduced extra field.

The two models (PINNs with and without the *extra fields*) are trained with  $N_f = 128$  collocation points and  $N_b = 2$  boundary points for 10,000 Adam iterations (with a learning rate of 0.001). And we have tested different network architectures, including  $[1] + 3 \times [50] + [\cdot]$ ,  $[1] + 3 \times [100] + [\cdot]$ ,  $[1] + 5 \times [50] + [\cdot]$ ,  $[1] + 5 \times [100] + [\cdot]$ , where for the PINN without the *extra fields*, the number of outputs is 1, and for the PINN with the *extra fields*, the number of outputs is 2.

**Schrödinger Equation** The governing PDEs are described as

$$i \frac{\partial h}{\partial t} + \frac{1}{2} \frac{\partial^2 h}{\partial x^2} + |h(x, t)|^2 h(x, t) = 0, \quad x \in (-5, 5), t \in (0, \pi/2], \quad (\text{A75a})$$

$$h(t, -5) = h(t, 5), \quad t \in (0, \pi/2], \quad (\text{A75b})$$

$$\frac{\partial h}{\partial x}(t, -5) = \frac{\partial h}{\partial x}(t, 5), \quad t \in (0, \pi/2], \quad (\text{A75c})$$

$$h(0, x) = 2 \operatorname{sech}(x), \quad x \in (-5, 5), \quad (\text{A75d})$$

where  $h(x, t)$  is the physical quantity of interest.

And the reformulated PDEs are (corresponding to the *extra fields*)

$$i \frac{\partial h}{\partial t} + \frac{1}{2} \frac{\partial p}{\partial x} + |h(x, t)|^2 h(x, t) = 0, \quad x \in (-5, 5), t \in (0, \pi/2], \quad (\text{A76a})$$

$$p(x, t) = \frac{\partial h}{\partial x}, \quad x \in [-5, 5], t \in (0, \pi/2], \quad (\text{A76b})$$

$$h(t, -5) = h(t, 5), \quad t \in (0, \pi/2], \quad (\text{A76c})$$

$$\frac{\partial h}{\partial x}(t, -5) = \frac{\partial h}{\partial x}(t, 5), \quad t \in (0, \pi/2], \quad (\text{A76d})$$

$$h(0, x) = 2 \operatorname{sech}(x), \quad x \in (-5, 5), \quad (\text{A76e})$$

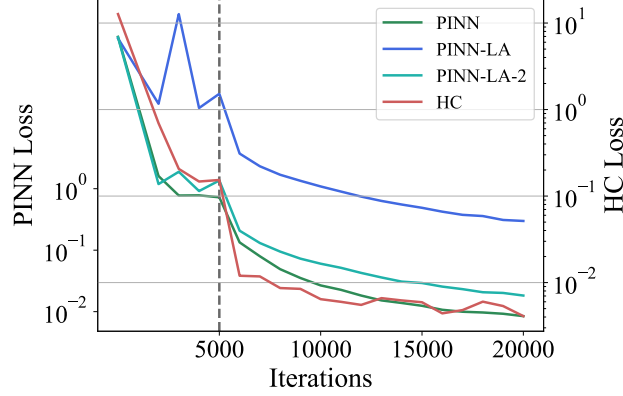


Figure A3: The convergence history of the simulation of a 2D battery pack.

where  $p(x)$  is the introduced extra field.

The two models (PINNs with and without the *extra fields*) are trained with  $N_f = 1000$  collocation points,  $N_b = 20$  boundary points, and  $N_i = 200$  initial points for 10,000 Adam iterations (with a learning rate of 0.001). And we have tested different network architectures, including  $[1] + 3 \times [50] + [\cdot]$ ,  $[1] + 3 \times [100] + [\cdot]$ ,  $[1] + 5 \times [50] + [\cdot]$ ,  $[1] + 5 \times [100] + [\cdot]$ , where for the PINN without the *extra fields*, the number of outputs is 2, and for the PINN with the *extra fields*, the number of outputs is 4.

**Experimental Results** We report the ratio of the the moving variance (MovVar) of  $|\nabla_{\theta} \mathcal{L}_{\mathcal{F}}|$  to that of  $|\nabla_{\tilde{\theta}} \tilde{\mathcal{L}}_{\mathcal{F}}|$  at each iteration during training, where the window size of the MovVar is 500 and after the MovVar, a moving average filter with a window size of 500 is applied. The results are shown in Figure 4. Besides, we also calculate the coefficient of variation (CV) of all the values of  $|\nabla_{\theta} \mathcal{L}_{\mathcal{F}}|$  and  $|\nabla_{\tilde{\theta}} \tilde{\mathcal{L}}_{\mathcal{F}}|$ , respectively. And we give the results in Figure A2. Using the CV as a criterion, we also find that the *extra fields* significantly reduces the gradient oscillations during training, especially for the complex nonlinear PDEs.

## A.9 Empirical Analysis of Convergence

In this subsection, we will empirically analyze the convergence of our method as well as some representative baselines in the context of the simulation of a 2D battery pack (see Section 5.2). We now report the training history with respect to iterations in Figure A3. The left axis shows the loss of PINNs (including PINN, PINN-LA, and PINN-LA-2) while the right axis shows the loss of our method, HC. The PINN loss is computed by adding up all the loss terms (including the losses of PDEs and BCs), where the loss weights are ignored for PINN-LA and PINN-LA-2. The first 5000 iterations are trained with Adam (separated by the gray dotted line), and the last 15000 are trained with L-BFGS.

From the results in Figure A3, we can see that the loss functions of all models drop significantly after switching to L-BFGS. This shows that L-BFGS can further promote convergence through utilizing

Table 1: Parallel experimental results of the simulation of a 2D battery pack (MAE of  $T$ )

	$t = 0$	$t = 0.5$	$t = 1$	average
PINN	$0.1232 \pm 0.0219$	$0.0417 \pm 0.0141$	$0.0263 \pm 0.0078$	$0.0499 \pm 0.0135$
PINN-LA	$0.1083 \pm 0.0266$	$0.0927 \pm 0.0372$	$0.1168 \pm 0.0739$	$0.0969 \pm 0.0385$
PINN-LA-2	$0.1065 \pm 0.0059$	$0.0322 \pm 0.0031$	<b><math>0.0200 \pm 0.0020</math></b>	$0.0400 \pm 0.0031$
FBPINN	$0.0763 \pm 0.0071$	$0.0258 \pm 0.0037$	$0.0205 \pm 0.0041$	$0.0318 \pm 0.0027$
xPINN	$0.2085 \pm 0.0252$	$0.1144 \pm 0.0194$	$0.1352 \pm 0.0241$	$0.1310 \pm 0.0194$
PFNN	<b><math>0.0000 \pm 0.0000</math></b>	$0.3769 \pm 0.0974$	$0.6012 \pm 0.2274$	$0.3522 \pm 0.1019$
PFNN-2	<b><math>0.0000 \pm 0.0000</math></b>	$0.3814 \pm 0.0381$	$0.5247 \pm 0.0394$	$0.3365 \pm 0.0236$
HC	<b><math>0.0000 \pm 0.0000</math></b>	<b><math>0.0244 \pm 0.0010</math></b>	$0.0226 \pm 0.0012$	<b><math>0.0219 \pm 0.0007</math></b>

Table 2: Parallel experimental results of the simulation of a 2D battery pack (MAPE of  $T$ )

	$t = 0$	$t = 0.5$	$t = 1$	average
PINN	$123.16 \pm 21.91\%$	$10.97 \pm 3.57\%$	$4.52 \pm 0.98\%$	$23.58 \pm 5.61\%$
PINN-LA	$108.14 \pm 26.57\%$	$24.15 \pm 8.07\%$	$17.98 \pm 8.97\%$	$33.64 \pm 8.89\%$
PINN-LA-2	$106.39 \pm 5.86\%$	$8.70 \pm 0.57\%$	<b><math>3.86 \pm 0.36\%</math></b>	$19.56 \pm 1.00\%$
FBPINN	$76.25 \pm 7.06\%$	$7.69 \pm 0.68\%$	$5.26 \pm 0.71\%$	$15.07 \pm 0.57\%$
xPINN	$208.36 \pm 25.21\%$	$26.25 \pm 4.99\%$	$18.15 \pm 3.05\%$	$49.60 \pm 7.37\%$
PFNN	<b><math>0.02 \pm 0.00\%</math></b>	$94.63 \pm 17.68\%$	$105.38 \pm 27.16\%$	$80.92 \pm 15.08\%$
PFNN-2	<b><math>0.02 \pm 0.00\%</math></b>	$71.39 \pm 6.98\%$	$82.04 \pm 8.90\%$	$61.65 \pm 4.64\%$
HC	<b><math>0.02 \pm 0.00\%</math></b>	<b><math>5.29 \pm 0.16\%</math></b>	$3.87 \pm 0.14\%$	<b><math>5.22 \pm 0.17\%</math></b>

Table 3: Parallel experimental results of the simulation of an airfoil (MAE)

	$u_1$	$u_2$	$p$
PINN	$0.4234 \pm 0.0809$	$0.0681 \pm 0.0162$	$0.3204 \pm 0.1404$
PINN-LA	$0.4467 \pm 0.0450$	$0.0630 \pm 0.0061$	$0.3028 \pm 0.0480$
PINN-LA-2	$0.4542 \pm 0.0875$	$0.0679 \pm 0.0111$	$0.3230 \pm 0.1115$
FBPINN	$0.3975 \pm 0.0221$	$0.0544 \pm 0.0030$	$0.2650 \pm 0.0059$
xPINN	$0.6942 \pm 0.0432$	$0.0581 \pm 0.0013$	$1.1587 \pm 0.1251$
HC	<b><math>0.2824 \pm 0.0215</math></b>	<b><math>0.0435 \pm 0.0024</math></b>	<b><math>0.2144 \pm 0.0114</math></b>

the information of the second derivatives of the loss function. However, we may not start with L-BFGS because it can easily lead to divergence. We consider Adam+L-BFGS to be a practical choice. Furthermore, we find that the convergence of PINNs is negatively affected by the tricks of learning rate annealing algorithm, especially the PINN-LA without our modification. The HC has the fastest convergence rate among all models. This means that the hard-constraint method or extra fields may be helpful in accelerating convergence.

#### A.10 Parallel Experiments

In this subsection, we revisit the three experiments in Section 5.2~5.4 and perform parallel tests in 5 runs to assess the significance of the results. We report the testing results (along with the 95% confidence intervals) in Table 1~6. From the results, we can see that our method, HC still outperforms all the other baselines. Besides, HC has the least variation, which shows that the hard-constraint methods can improve the stability of training.

#### A.11 Ethics Statement

PDEs have important applications in many fields, including applied physics, automobile manufacturing, economics, and the aerospace industry. Solving PDE via neural networks has attracted much attention in recent years, and it may be applied to the above fields in the future. Our method also belongs to this kind. However, the method for solving PDE based on neural networks has no theoretical explanation and safety guarantee for the time being. Applying such methods to security-sensitive domains may lead to unexpected incidents and the cause of the accident may be hard to diagnose. Possible solutions include developing alternatives with theoretical interpretability or using safeguards.

Table 4: Parallel experimental results of the simulation of an airfoil (WMAPE)

	$u_1$	$u_2$	$p$
PINN	$0.5358 \pm 0.1024$	$1.1709 \pm 0.2778$	$0.2921 \pm 0.1279$
PINN-LA	$0.5653 \pm 0.0570$	$1.0819 \pm 0.1048$	$0.2760 \pm 0.0437$
PINN-LA-2	$0.5747 \pm 0.1106$	$1.1670 \pm 0.1920$	$0.2944 \pm 0.1016$
FBPINN	$0.5030 \pm 0.0279$	$0.9347 \pm 0.0517$	$0.2416 \pm 0.0054$
xPINN	$0.8784 \pm 0.0546$	$0.9986 \pm 0.0225$	$1.0562 \pm 0.1140$
HC	<b><math>0.3573 \pm 0.0272</math></b>	<b><math>0.7472 \pm 0.0418</math></b>	<b><math>0.1954 \pm 0.0104</math></b>

Table 5: Parallel experimental results of the high-dimensional heat equation (MAE of  $u$ )

	$t = 0$	$t = 0.5$	$t = 1$	average
PINN	$0.0204 \pm 0.0148$	$0.0357 \pm 0.0104$	$0.1600 \pm 0.0600$	$0.0525 \pm 0.0173$
PINN-LA	$0.0430 \pm 0.0751$	$0.3039 \pm 0.6691$	$0.8011 \pm 1.7228$	$0.3464 \pm 0.7531$
PINN-LA-2	$0.0287 \pm 0.0670$	$0.2071 \pm 0.6524$	$0.5933 \pm 1.7225$	$0.2455 \pm 0.7433$
PFNN	<b>0.0000</b> $\pm 0.0000$	$0.0895 \pm 0.0727$	$0.2130 \pm 0.1790$	$0.0963 \pm 0.0788$
HC	<b>0.0000</b> $\pm 0.0000$	<b>0.0028</b> $\pm 0.0006$	<b>0.0046</b> $\pm 0.0008$	<b>0.0027</b> $\pm 0.0006$

Table 6: Parallel experimental results of the high-dimensional heat equation (MAPE of  $u$ )

	$t = 0$	$t = 0.5$	$t = 1$	average
PINN	$1.15 \pm 0.51\%$	$1.41 \pm 0.41\%$	$3.83 \pm 1.45\%$	$1.75 \pm 0.58\%$
PINN-LA	$2.86 \pm 5.07\%$	$12.06 \pm 26.54\%$	$19.35 \pm 41.65\%$	$11.46 \pm 24.75\%$
PINN-LA-2	$1.92 \pm 4.55\%$	$8.19 \pm 25.80\%$	$14.29 \pm 41.54\%$	$8.00 \pm 24.25\%$
PFNN	<b>0.00</b> $\pm 0.00\%$	$3.59 \pm 2.93\%$	$5.19 \pm 4.36\%$	$3.20 \pm 2.60\%$
HC	<b>0.00</b> $\pm 0.00\%$	<b>0.11</b> $\pm 0.03\%$	<b>0.11</b> $\pm 0.02\%$	<b>0.10</b> $\pm 0.02\%$

## References

- [1] Tianping Chen and Hong Chen. Universal approximation to nonlinear operators by neural networks with arbitrary activation functions and its application to dynamical systems. *IEEE Transactions on Neural Networks*, 6(4):911–917, 1995.
- [2] Bernardo Llanas, Sagrario Lantarón, and Francisco J Sáinz. Constructive approximation of discontinuous functions by neural networks. *Neural Processing Letters*, 27(3):209–226, 2008.
- [3] Lu Lu, Xuhui Meng, Zhiping Mao, and George Em Karniadakis. Deepxde: A deep learning library for solving differential equations. *SIAM Review*, 63(1):208–228, 2021.
- [4] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. 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 8024–8035. Curran Associates, Inc., 2019.
- [5] Hailong Sheng and Chao Yang. Pfn: A penalty-free neural network method for solving a class of second-order boundary-value problems on complex geometries. *Journal of Computational Physics*, 428:110085, 2021.
- [6] Hailong Sheng and Chao Yang. Pfn-2: A domain decomposed penalty-free neural network method for solving partial differential equations. *arXiv preprint arXiv:2205.00593*, 2022.
- [7] Sifan Wang, Yujun Teng, and Paris Perdikaris. Understanding and mitigating gradient flow pathologies in physics-informed neural networks. *SIAM Journal on Scientific Computing*, 43(5):A3055–A3081, 2021.