

A DISCUSSION ON THE DEFINITION OF (CAUSAL) FOUNDATION MODELS

In this paper, we focus on treatment effect estimation tasks (defined in Section 3.1). Our model is then tailored for generalizable zero-shot estimating average treatment effects. That is, given unseen datasets/contexts that contains observational records of covariates, treatments, and effects, we aim to estimate the underlying treatment effects using a forward pass of the underlying model.

This approach is inline with the definition of foundation models discussed in Bommasani et al. (2021): “any model that is trained on broad data (generally using self-supervision at scale) that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks”. Note that such *task-universality* of foundation models does not necessarily imply adaptability across different *machine learning formulations* (e.g., prediction, imputation, ATE, CATE, counterfactuals); instead, it can refer to adaptability across different *contexts* for a given task. This perspective is widely embraced by recent studies, such as those focusing on foundation models for tabular datasets (Zhang et al., 2023b), time series (Garza & Mergenthaler-Canseco, 2023; Das et al., 2023), and knowledge graphs (Galkin et al., 2023). These studies concentrate exclusively on a single type of task, but assess in-context generalization across datasets.

B EXTENDED RELATED WORKS

As our work also intersects with the literature on neural causal estimation methods, we provide a discussion in this section.

Neural Estimation Methods for Treatment Effects. Research in this direction employs deep learning methods to estimate treatment effects, typically relying on standard assumptions that ensure identifiability, similar to our setting. A prominent approach focuses on learning a representation of the covariates that is predictive of the outcome (Johansson et al., 2016; Shalit et al., 2017; Yao et al., 2018). Following this, several methods have been proposed to combine outcome models learned through neural networks with balanced propensity weights (Alaa et al., 2017; Schwab et al., 2018; Du et al., 2021). Semi-parametric estimation theory and doubly robust estimators have also been applied in neural estimation methods, e.g., using regularization (Shi et al., 2019) or shared representations (Chernozhukov et al., 2018). Another perspective of using neural network is to control for complex relationships and covariates. Kallus (2020a) extends adversarial covariate balancing (Kallus, 2020b) using flexible modeling with neural networks. Generative causal models have also been proposed to leverage the expressivity of neural networks to approximate structural causal models (Louizos et al., 2017; Kocaoglu et al., 2017; Alaa & Van Der Schaar, 2017; Yoon et al., 2018; Pawlowski et al., 2020; Xia et al., 2021, 2022), which then allows for the estimation of treatment effects. In addition, Xia et al. (2021) also proved that their proposed method can be used to test the identifiability of causal effect in terms of do-interventions (Pearl, 2009) in the general setting. Xia et al. (2022) extended such testing for counterfactual outcomes (Bareinboim et al., 2022). In (Melnichuk et al., 2022), the attention mechanism was employed to estimate treatment effect over time for a given unit. Concurrent to our work, Nilforoshan et al. (2023) proposed a meta-learning framework to learn causal effects of various structured treatments on the same population. Their method leverages information across different treatments, which allows for zero-shot learning on an unseen treatment. Our work can be viewed as orthogonal, as we focus on learning the causal effects of the same treatment across different populations.

C OMITTED PROOFS

C.1 DERIVATIONS OF EQ. (1) AND EQ. (2)

We first establish the conditional bias decomposition:

$$\begin{aligned} & \mathbb{E}(\hat{\tau} - \tau_{SATE} \mid \{\mathbf{X}_i, T_i\}_{i=1}^N) \\ &= \mathbb{E}\left(\sum_{i=1}^N \alpha_i W_i Y_i - \sum_{i=1}^N \frac{1}{N} (Y_i(1) - Y_i(0)) \mid \{\mathbf{X}_i, T_i\}_{i=1}^N\right) \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^N \alpha_i W_i \mathbb{E}(Y_i(T_i) \mid \mathbf{X}_i, T_i) + \sum_{i=1}^N \frac{1}{N} \mathbb{E}(Y_i(1) - Y_i(0) \mid \mathbf{X}_i, T_i) \\
&= \sum_{i=1}^N (\alpha_i W_i \mathbb{E}(Y_i(0) \mid \mathbf{X}_i) + \alpha_i T_i \mathbb{E}(Y_i(1) - Y_i(0) \mid \mathbf{X}_i)) + \sum_{i=1}^N \frac{1}{N} \mathbb{E}(Y_i(1) - Y_i(0) \mid \mathbf{X}_i) \\
&= \sum_{i=1}^N (\alpha_i T_i - \frac{1}{N}) \mathbb{E}(Y_i(1) - Y_i(0) \mid \mathbf{X}_i) + \sum_{i=1}^N \alpha_i W_i \mathbb{E}(Y_i(0) \mid \mathbf{X}_i),
\end{aligned}$$

where we use the assumption of consistency between observed and potential outcomes and non-interference between unit (SUTVA, Rubin (1990)) in the second equation and unconfoundedness in the third equation.

Formally, define a feature map $\phi : \mathbb{X} \rightarrow \mathcal{H}_\phi$, where \mathbb{X} is the support of covariates and \mathcal{H}_ϕ is some Hilbert space. The unit-ball RKHS is given by $\mathcal{F}_\phi = \{f : \mathbb{X} \rightarrow \mathbb{R} \mid \exists \theta \in \mathcal{H}_\phi, \text{ s.t. } f(x) = \langle \theta, \phi(x) \rangle, \forall x \in \mathbb{X} \text{ and } \|\theta\| \leq 1\}$. Recall that $\langle \cdot, \cdot \rangle$ denotes the inner product of Hilbert space \mathcal{H}_ϕ and $\|\cdot\|$ denotes the associated norm. The adversarial upper bound of the square of the second term in the conditional bias can be calculated via

$$\begin{aligned}
&\sup_{f \in \mathcal{F}_\phi} \left(\sum_{i=1}^N \alpha_i W_i f(\mathbf{X}_i) \right)^2 \\
&= \sup_{\theta \in \mathcal{H}_\phi, \|\theta\| \leq 1} \left(\sum_{i=1}^N \alpha_i W_i \langle \theta, \phi(\mathbf{X}_i) \rangle \right)^2 \\
&= \sup_{\theta \in \mathcal{H}_\phi, \|\theta\| \leq 1} \left(\left\langle \theta, \sum_{i=1}^N \alpha_i W_i \phi(\mathbf{X}_i) \right\rangle \right)^2 \\
&\leq \left\| \sum_{i=1}^N \alpha_i W_i \phi(\mathbf{X}_i) \right\|^2 = \boldsymbol{\alpha}^\top \mathbf{K}_\phi \boldsymbol{\alpha}.
\end{aligned}$$

Recall that $[\mathbf{K}_\phi]_{ij} = W_i W_j \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$. Therefore minimizing this adversarial loss subject to $\boldsymbol{\alpha} \in \mathbb{A}$ reduces to Eq. (1).

By evoking Theorem 1 in Tarr & Imai (2021), we have that Eq. (1) is equivalent to Eq. (2) for some $\lambda \geq 0$. However, the exact value of λ depends on \mathbf{K}_ϕ . For example, if \mathbf{K}_ϕ is such that the minimum value of Eq. (1) is 0, then $\lambda = 0$. This is because the minimizer of Eq. (1) would also be the minimizer under the unnormalized constraint (Eq. (2) with $\lambda = 0$), as $\boldsymbol{\alpha}^\top \mathbf{K}_\phi \boldsymbol{\alpha} \geq 0$ for any $\boldsymbol{\alpha} \in \mathbb{R}^N$.

Conversely, we can also show that $\lambda > 0$ if \mathbf{K}_ϕ is of full rank.

Lemma 1. *If \mathbf{K}_ϕ is of full rank, then $\lambda > 0$.*

Proof. From the proof of Theorem 1 in Tarr & Imai (2021), we know that $\lambda = 0$ only if $q_* = \min_{\mathbf{W}^\top \boldsymbol{\alpha} = 0, 0 \leq \alpha_i \leq 1, \alpha_i \neq 0} \frac{\sqrt{\boldsymbol{\alpha}^\top \mathbf{K}_\phi \boldsymbol{\alpha}}}{\mathbf{1}^\top \boldsymbol{\alpha} / 2}$ is zero. However, since \mathbf{K}_ϕ is of full rank, it is positive definite. Thus for any $\boldsymbol{\alpha} \neq \mathbf{0}$, there is $\boldsymbol{\alpha}^\top \mathbf{K}_\phi \boldsymbol{\alpha} > 0$. Therefore $q_* > 0$. Consequently, $\lambda > 0$. \square

C.2 DERIVATIONS OF EQ. (3) AND EQ. (4)

The dual form of Eq. (3) can be derived using its Lagrangian

$$L(\boldsymbol{\beta}, \beta_0, \boldsymbol{\xi}, \boldsymbol{\alpha}, \bar{\boldsymbol{\alpha}}) = \frac{\lambda}{2} \|\boldsymbol{\beta}\|^2 + \sum_{i=1}^N \xi_i + \sum_{i=1}^N \alpha_i \left(1 - \xi_i - W_i \left(\langle \boldsymbol{\beta}, \phi(\mathbf{X}_i) \rangle + \beta_0 \right) \right) - \sum_{i=1}^N \bar{\alpha}_i \xi_i,$$

where $\boldsymbol{\alpha} \succeq \mathbf{0}$ and $\bar{\boldsymbol{\alpha}} \succeq \mathbf{0}$. The primal form in Eq. (3) can be obtained by $\min_{\boldsymbol{\beta}, \beta_0, \boldsymbol{\xi}} \max_{\boldsymbol{\alpha} \succeq \mathbf{0}, \bar{\boldsymbol{\alpha}} \succeq \mathbf{0}} L(\boldsymbol{\beta}, \beta_0, \boldsymbol{\xi}, \boldsymbol{\alpha}, \bar{\boldsymbol{\alpha}})$. If we exchange $\min \max$ with $\max \min$, solving

$\min_{\beta, \beta_0, \xi_i}$ by setting the derivatives to zero leads to

$$\begin{aligned}\nabla_{\beta} L(\beta, \beta_0, \xi, \alpha, \bar{\alpha}) &= \lambda \beta - \sum_{i=1}^N \alpha_i W_i \phi(\mathbf{X}_i) = \mathbf{0}, \\ \nabla_{\beta_0} L(\beta, \beta_0, \xi, \alpha, \bar{\alpha}) &= - \sum_{i=1}^N \alpha_i W_i = 0, \\ \nabla_{\xi_i} L(\beta, \beta_0, \xi, \alpha, \bar{\alpha}) &= 1 - \alpha_i - \bar{\alpha}_i = 0, \quad \forall i \in [N].\end{aligned}$$

Plugging these in $L(\beta, \beta_0, \xi, \alpha, \bar{\alpha})$, we can reduce $\max_{\alpha \geq \mathbf{0}, \bar{\alpha} \geq \mathbf{0}} \min_{\beta, \beta_0, \xi_i} L(\beta, \beta_0, \xi, \alpha, \bar{\alpha})$ to Eq. (2). Thus it is the dual form of Eq. (3).

In addition, we can also derive Eq. (4). It is easy to check that Slater’s condition holds for the primal SVM problem in Eq. (3). Thus it satisfies strong duality. Therefore any optimal solutions to the primal-dual problems must satisfy the KKT condition $\lambda \beta^* = \sum_{j=1}^N \alpha_j^* W_j \phi(\mathbf{X}_j)$.

C.3 DERIVATIONS OF EQ. (6)

From the Taylor expansion

$$\begin{aligned}\exp(\mathbf{k}_i^\top \mathbf{k}_j / \sqrt{D}) &= \sum_{l=0}^{+\infty} \frac{1}{l!} (\mathbf{k}_i^\top \mathbf{k}_j / \sqrt{D})^l \\ &= \sum_{l=0}^{+\infty} \sum_{N_1 + \dots + N_D = l} \frac{([\mathbf{k}_i]_1^{N_1} \dots [\mathbf{k}_i]_D^{N_D}) ([\mathbf{k}_j]_1^{N_1} \dots [\mathbf{k}_j]_D^{N_D})}{D^{l/2} N_1! \dots N_D!},\end{aligned}$$

we have that $\exp(\mathbf{k}_i^\top \mathbf{k}_j / \sqrt{D}) = \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$ if

$$\phi(\mathbf{x}) = \left(\frac{[\mathbf{k}]_1^{N_1} \dots [\mathbf{k}]_D^{N_D}}{D^{l/2} (N_1! \dots N_D!)^{1/2}} \right)_{N_1 + \dots + N_D = l, l \in \mathbb{N}}. \quad (9)$$

Here \mathbf{k} denotes the key embedding of \mathbf{x} following the same transformation that \mathbf{k}_i is obtained from \mathbf{X}_i . Note that we allow the transformation to depend on \mathbf{X} , which corresponds to a data-dependent kernel.

Using this expression, the i -th output of the self-attention layer when $\mathbf{Q} = \mathbf{K}$ can be equivalently written as

$$\sum_{j=1}^N \frac{\exp(\mathbf{k}_i^\top \mathbf{k}_j / \sqrt{D})}{\sum_{j'=1}^N \exp(\mathbf{k}_i^\top \mathbf{k}_{j'} / \sqrt{D})} v_j = \sum_{j=1}^N \frac{\langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle}{h(\mathbf{X}_i)} v_i = \sum_{j=1}^N \frac{v_j}{h(\mathbf{X}_j)} \langle \phi(\mathbf{X}_j), \phi(\mathbf{X}_i) \rangle.$$

C.4 PROOF OF THEOREM 1

We first state its formal version:

Theorem 1. *If the covariates \mathbf{X} satisfy that $\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)$ are linearly independent, then Algorithm 1 recovers the optimal balancing weight at the global minimum of the penalized hinge loss in Eq. (7).*

In particular, the optimal solution α^ to Eq. (1), in which the feature function ϕ is defined using the optimal neural network parameters via Eq. (9), can be obtained using the optimal neural network parameters that minimize Eq. (7) via $\alpha_j^* = \lambda v_j / h(\mathbf{X}_j) W_j$.*

Proof. Denote $\beta = \sum_{j=1}^N \frac{v_j}{h(\mathbf{X}_j)} \phi(\mathbf{X}_j)$, then using Eq. (6), we can rewrite the loss function in Eq. (7) as

$$\mathcal{L}_{\theta}(\mathbb{D}) = \frac{\lambda}{2} \|\beta\|^2 + \sum_{i=1}^N [1 - W_i (\langle \beta, \phi(\mathbf{X}_i) \rangle + \beta_0)]_+.$$

Denote $\xi_i = [1 - W_i(\langle \boldsymbol{\beta}, \phi(\mathbf{X}_i) \rangle + \beta_0)]_+$, then minimizing $\mathcal{L}_\theta(\mathbb{D})$ can be equivalently written as

$$\begin{aligned} \min_{\boldsymbol{\theta}} \quad & \frac{\lambda}{2} \|\boldsymbol{\beta}\|^2 + \sum_{i=1}^N \xi_i, \\ \text{s.t.} \quad & W_i(\langle \boldsymbol{\beta}, \phi(\mathbf{X}_i) \rangle + \beta_0) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad \forall i \in [N]. \end{aligned}$$

Thus at the optimal $\boldsymbol{\theta}$, the corresponding $\boldsymbol{\beta}$ is also the optimal solution to

$$\begin{aligned} \min_{\boldsymbol{\beta}, \beta_0, \boldsymbol{\xi}} \quad & \frac{\lambda}{2} \|\boldsymbol{\beta}\|^2 + \sum_{i=1}^N \xi_i, \\ \text{s.t.} \quad & W_i(\langle \boldsymbol{\beta}, \phi(\mathbf{X}_i) \rangle + \beta_0) \geq 1 - \xi_i, \quad \xi_i \geq 0, \quad \forall i \in [N], \end{aligned}$$

where ϕ is defined using the optimal $\boldsymbol{\theta}$. This recovers the primal SVM problem. By the primal-dual connection proven in Appendix C.2, if we denote the optimal solution to the dual problem (which is Eq. (2)) as $\boldsymbol{\alpha}^*$, we have

$$\lambda \boldsymbol{\beta} = \sum_{j=1}^N \alpha_j^* W_j \phi(\mathbf{X}_j).$$

Consequently, by the definition of $\boldsymbol{\beta}$, we have

$$\sum_{j=1}^N \frac{\lambda v_j}{h(\mathbf{X}_j)} \phi(\mathbf{X}_j) = \sum_{j=1}^N \alpha_j^* W_j \phi(\mathbf{X}_j).$$

By the assumption that $\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)$ are linearly independent, we must have $\frac{\lambda v_j}{h(\mathbf{X}_j)} = \alpha_j^* W_j$ for all $j \in [N]$. Therefore $\alpha_j^* = \lambda v_j / h(\mathbf{X}_j) W_j$. \square

Remark 1. Note that when $\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)$ are linearly independent, the matrix $\mathbf{K}_\phi = [W_1 \phi(\mathbf{X}_1), \dots, W_N \phi(\mathbf{X}_N)]^\top [W_1 \phi(\mathbf{X}_1), \dots, W_N \phi(\mathbf{X}_N)]$ is of full rank. Thus by Lemma 1, there is $\lambda > 0$. Conversely, using a similar decomposition, we know that if $\hat{\mathbf{K}}_\phi = [\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)]^\top [\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)]$ is of full rank, then $\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)$ are linearly independent. Since $\hat{\mathbf{K}}_\phi = \exp(\mathbf{K} \mathbf{K}^\top / \sqrt{D})$, we have $\phi(\mathbf{X}_1), \dots, \phi(\mathbf{X}_N)$ linearly independent if \mathbf{K} is of row rank N . Thus the assumption on \mathbf{X} in Theorem 1 is satisfied when \mathbf{K} is of row rank N .

D ALTERNATIVE OBJECTIVES

Consider minimizing the square of both terms in the conditional bias, which we decompose into the following form

$$\begin{aligned} & (\mathbb{E}(\hat{\tau} - \tau_{SATE} \mid \{\mathbf{X}_i, T_i\}_{i=1}^N))^2 \\ & = \left(\sum_{i=1}^N \alpha_i W_i \mathbb{E}(Y_i(T_i) \mid \mathbf{X}_i, T_i) - \frac{1}{N} \sum_{i=1}^N \left(\mathbb{E}(Y_i(1) \mid \mathbf{X}_i) - \mathbb{E}(Y_i(0) \mid \mathbf{X}_i) \right) \right)^2. \end{aligned} \quad (10)$$

Denote the outcome models $\mathbb{E}(Y_i(1) \mid \mathbf{X}_i) = f_1(\mathbf{X}_i)$ and $\mathbb{E}(Y_i(0) \mid \mathbf{X}_i) = f_0(\mathbf{X}_i)$. We choose to minimize the above term in worst case over all possible potential outcome models $(f_0, f_1) \in \mathcal{F}_\phi^2$. Here the space \mathcal{F}_ϕ^2 is defined as $\mathcal{F}_\phi^2 = \{(f_0, f_1) \mid f_0 \in \mathcal{F}_\phi, f_1 \in \mathcal{F}_\phi\}$.

Suppose $f_0(x) = \langle \phi(x), \theta_0 \rangle$ and $f_1(x) = \langle \phi(x), \theta_1 \rangle$ for $\theta_0, \theta_1 \in \mathcal{H}_\phi$, $\|\theta_0\| \leq 1$, $\|\theta_1\| \leq 1$. We can bound Eq. (10) with respect to all outcome models in \mathcal{F}_ϕ^2 as

$$\begin{aligned} & \left(\sum_{i=1}^N \alpha_i W_i f_{T_i}(\mathbf{X}_i) - \frac{1}{N} \sum_{i=1}^N (f_1(\mathbf{X}_i) - f_0(\mathbf{X}_i)) \right)^2 \\ &= \left(\left\langle \sum_{i \in \mathbb{T}} \alpha_i W_i \phi(\mathbf{X}_i) - \frac{1}{N} \sum_{i \in [N]} \phi(\mathbf{X}_i), \theta_1 \right\rangle + \left\langle \sum_{i \in \mathbb{C}} \alpha_i W_i \phi(\mathbf{X}_i) + \frac{1}{N} \sum_{i \in [N]} \phi(\mathbf{X}_i), \theta_0 \right\rangle \right)^2 \\ &\leq 2 \left(\sum_{i \in \mathbb{T}} \alpha_i W_i \phi(\mathbf{X}_i) - \frac{1}{N} \sum_{i \in [N]} \phi(\mathbf{X}_i) \right)^2 + 2 \left(\sum_{i \in \mathbb{C}} \alpha_i W_i \phi(\mathbf{X}_i) + \frac{1}{N} \sum_{i \in [N]} \phi(\mathbf{X}_i) \right)^2 \end{aligned}$$

where the inequality uses Cauchy-Schwartz inequality. Minimizing this upper bound subject to $\alpha \in \mathbb{A}$ is equivalent to solving

$$\begin{aligned} \min_{\alpha} \quad & \alpha^\top \mathbf{G}_\phi \alpha + \alpha^\top \mathbf{g}_\phi, \\ \text{s.t.} \quad & \sum_{i \in \mathbb{T}} \alpha_i = \sum_{i \in \mathbb{C}} \alpha_i = 1, \quad \mathbf{0} \preceq \alpha \preceq \mathbf{1}. \end{aligned} \quad (11)$$

Here

$$\begin{aligned} [\mathbf{G}_\phi]_{i,j} &= \delta_{W_i=W_j} \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle, \\ [\mathbf{g}_\phi]_i &= -\frac{2}{N} \sum_{j=1}^N \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle. \end{aligned}$$

It is easy to show that $\mathbf{G}_\phi \succeq 0$ as it can be decomposed into two submatrixes which are positive semi-definite. In addition, as $\langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle = \exp(\mathbf{k}_i^\top \mathbf{k}_j / \sqrt{D}) > 0$, we know that $\mathbf{g}_\phi \prec \mathbf{0}$.

To come up with a consistent gradient-based solver, notice first that Eq. (11) is equivalent to the following unnormalized problem for some $\lambda, \mu \geq 0$

$$\begin{aligned} \min_{\alpha} \quad & \alpha^\top \mathbf{G}_\phi \alpha + 2\mu \cdot \mathbf{g}_\phi^\top \alpha - 2\lambda \cdot \mathbf{1}^\top \alpha, \\ \text{s.t.} \quad & \mathbf{W}^\top \alpha = 0, \quad \mathbf{0} \preceq \alpha \preceq \mathbf{1}. \end{aligned} \quad (12)$$

This can be shown similarly to the proof of Theorem 1 in Tarr & Imai (2021). We escape the details but provide the following main steps:

1. We first show that for some $\epsilon_\lambda, \epsilon_\mu \geq 0$, Eq. (12) is equivalent to

$$\begin{aligned} \min_{\alpha} \quad & \alpha^\top \mathbf{G}_\phi \alpha, \\ \text{s.t.} \quad & \mathbf{W}^\top \alpha = 0, \quad \mathbf{0} \preceq \alpha \preceq \mathbf{1}, \quad -\mathbf{g}_\phi^\top \alpha \geq \epsilon_\mu, \quad \mathbf{1}^\top \alpha \geq \epsilon_\lambda. \end{aligned}$$

2. Next, we show that the above problem is equivalent to

$$\begin{aligned} \min_{\alpha} \quad & \sqrt{\alpha^\top \mathbf{G}_\phi \alpha}, \\ \text{s.t.} \quad & \mathbf{W}^\top \alpha = 0, \quad \mathbf{0} \preceq \alpha \preceq \mathbf{1}, \quad -\mathbf{g}_\phi^\top \alpha \geq \epsilon_\mu, \quad \mathbf{1}^\top \alpha \geq \epsilon_\lambda, \end{aligned}$$

which is equivalent to

$$\begin{aligned} \min_{\alpha} \quad & \sqrt{\alpha^\top \mathbf{G}_\phi \alpha} + \nu_\mu \cdot \mathbf{g}_\phi^\top \alpha - \nu_\lambda \mathbf{1}^\top \alpha, \\ \text{s.t.} \quad & \mathbf{W}^\top \alpha = 0, \quad \mathbf{0} \preceq \alpha \preceq \mathbf{1}. \end{aligned}$$

for some $\nu_\lambda, \nu_\mu \geq 0$.

3. For some $\lambda \geq 0$, the above problem is equivalent to

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{\sqrt{\boldsymbol{\alpha}^\top \mathbf{G}_\phi \boldsymbol{\alpha}} + \nu_\mu \cdot \mathbf{g}_\phi^\top \boldsymbol{\alpha}}{\mathbf{1}^\top \boldsymbol{\alpha}}, \\ \text{s.t.} \quad & \mathbf{W}^\top \boldsymbol{\alpha} = 0, \quad \mathbf{0} \preceq \boldsymbol{\alpha} \preceq \mathbf{1}. \end{aligned}$$

Since this problem is scale-free, it is equivalent to

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \frac{\sqrt{\boldsymbol{\alpha}^\top \mathbf{G}_\phi \boldsymbol{\alpha}} + \nu_\mu \cdot \mathbf{g}_\phi^\top \boldsymbol{\alpha}}{\mathbf{1}^\top \boldsymbol{\alpha}}, \\ \text{s.t.} \quad & \sum_{i \in \mathbb{T}} \alpha_i = \sum_{i \in \mathbb{C}} \alpha_i = 1, \quad \mathbf{0} \preceq \boldsymbol{\alpha} \preceq \mathbf{1}, \end{aligned}$$

i.e.,

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \sqrt{\boldsymbol{\alpha}^\top \mathbf{G}_\phi \boldsymbol{\alpha}} + \nu_\mu \cdot \mathbf{g}_\phi^\top \boldsymbol{\alpha}, \\ \text{s.t.} \quad & \sum_{i \in \mathbb{T}} \alpha_i = \sum_{i \in \mathbb{C}} \alpha_i = 1, \quad \mathbf{0} \preceq \boldsymbol{\alpha} \preceq \mathbf{1}, \end{aligned}$$

4. Using similar arguments as above, one can show the above problem is equivalent to

$$\begin{aligned} \min_{\boldsymbol{\alpha}} \quad & \boldsymbol{\alpha}^\top \mathbf{G}_\phi \boldsymbol{\alpha} + \mathbf{g}_\phi^\top \boldsymbol{\alpha}, \\ \text{s.t.} \quad & \sum_{i \in \mathbb{T}} \alpha_i = \sum_{i \in \mathbb{C}} \alpha_i = 1, \quad \mathbf{0} \preceq \boldsymbol{\alpha} \preceq \mathbf{1}, \end{aligned}$$

for some $\mu \geq 0$.

The primal form of Eq. (12) can be written as

$$\begin{aligned} \min_{\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \beta_0, \boldsymbol{\xi}} \quad & \frac{1}{2} \|\boldsymbol{\beta}_1\|^2 + \frac{1}{2} \|\boldsymbol{\beta}_2\|^2 + \sum_{i=1}^N \xi_i, \\ \text{s.t.} \quad & \langle \boldsymbol{\beta}_1, \phi(\mathbf{X}_i) \rangle + \beta_0 \geq \lambda - \mu [\mathbf{g}_\phi]_i - \xi_i, \quad \forall i \in \mathbb{T} \\ & \langle \boldsymbol{\beta}_2, \phi(\mathbf{X}_i) \rangle - \beta_0 \geq \lambda - \mu [\mathbf{g}_\phi]_i - \xi_i, \quad \forall i \in \mathbb{C} \\ & \xi_i \geq 0, \quad \forall i \in [N]. \end{aligned}$$

Following similar derivations in Appendix C, we can write out an unconstrained loss function

$$\begin{aligned} \mathcal{L}_\theta(\mathbb{D}) = & \frac{1}{2} \left\| \sum_{j \in \mathbb{T}} \frac{v_j}{h(\mathbf{X}_j)} \phi(\mathbf{X}_j) \right\|^2 + \frac{1}{2} \left\| \sum_{j \in \mathbb{C}} \frac{v_j}{h(\mathbf{X}_j)} \phi(\mathbf{X}_j) \right\|^2 \\ & + \left[\lambda - \mu [\mathbf{g}_\phi]_{\mathbb{T}} - (\text{softmax}(\mathbf{K}_{\mathbb{T}} \mathbf{K}_{\mathbb{T}}^\top / \sqrt{D}) \mathbf{V}_{\mathbb{T}} + \beta_0) \right]_+ \\ & + \left[\lambda - \mu [\mathbf{g}_\phi]_{\mathbb{C}} - (\text{softmax}(\mathbf{K}_{\mathbb{C}} \mathbf{K}_{\mathbb{C}}^\top / \sqrt{D}) \mathbf{V}_{\mathbb{C}} - \beta_0) \right]_+, \end{aligned}$$

where the optimal $\boldsymbol{\alpha}^*$ solving Eq. (11) can be read off as $\alpha_i = \frac{v_i}{h(\mathbf{X}_i)}$.

For the conditional mean square error, under regularity constraints in Bennett & Kallus (2019), we can also use the same upper bound as above (up to an additive $\mathcal{O}(1/N)$ gap). Therefore the same derivation holds. However, as this loss function separates the treated group from the control group aside from sharing the constant intercept β_0 , it might not be preferable than the objective proposed in the main text.

E NON-BINARY TREATMENTS

Consider a generalization to the setting in Section 3.1, where the dataset $\mathbb{D} = \{(\mathbf{X}_i, \mathbf{T}_i, Y_i)\}_{i \in [N]}$ in which \mathbf{T}_i is a S -dimensional vector of multiple binary treatments. Let $Y_i^s(t)$ be the potential outcome of assigning treatment $[\mathbf{T}_i]_s = t$.

Assuming SUTVA ($Y_i = Y_i^s([\mathbf{T}]_s)$) and unconfoundedness. Denote $\mathbb{T}^s = \{i \in [N] : [\mathbf{T}]_s = 1\}$ and $\mathbb{C}^s = \{i \in [N] : [\mathbf{T}]_s = 0\}$. We consider weighted estimators in the form of

$$\hat{\tau}^s = \sum_{i \in \mathbb{T}^s} \alpha_i Y_i^s(1) - \sum_{i \in \mathbb{C}^s} \alpha_i Y_i^s(0)$$

for the sample average treatment of the s -th treatment

$$\tau_{SATE}^s = \frac{1}{N} \sum_{i=1}^N (Y_i^s(1) - Y_i^s(0)).$$

Following the same derivations in Section 3 and Appendix C, we can obtain a dual-SVM formulation to optimize α in the adversarial case. This dual-SVM formulation can then be transformed into its primal problem. As self-attention is implicitly implementing the predictor in the primal problem, we can then read off the optimal α^* by training this self-attention-based neural network with a penalized hinge loss.

However, as we would like to evaluate the sample average treatment for multiple treatments, we can actually aggregate S SVM problems together using the flexibility of self-attention layers. Namely, instead of consider a one-dimensional value vector \mathbf{V} in Section 3.2, we use $\mathbf{V} \in \mathbb{R}^{N \times S}$, where the s -th dimension corresponds to the s -th treatment. By minimizing the following loss function

$$\mathcal{L}_\theta(\mathbb{D}) = \frac{\lambda}{2} \sum_{s=1}^S \left\| \sum_{j=1}^N \frac{[\mathbf{V}]_{js}}{h(\mathbf{X}_j)} \phi(\mathbf{X}_j) \right\|^2 + \sum_{s=1}^S \left[\mathbf{1} - \mathbf{W}_{:,s} (\text{softmax}(\mathbf{K}\mathbf{K}^\top / \sqrt{D}) \mathbf{V}_{:,s} + \beta_0) \right]_+,$$

we can read off the optimal balancing weight α for the s -th treatment via $\lambda \cdot \mathbf{V}_{:,s} / h(\mathbf{X}) \mathbf{W}_{:,s}$.

F INDIVIDUAL TREATMENT EFFECT ESTIMATION

In this section, we further consider the problem of estimating **individual treatment effect (ITE)** in the binary treatment setup of Section 3. [Here we present one possible algorithmic approach to approximate ITEs with CInA.](#) Without loss of generality, suppose $T_1 = 1$ and we would like to estimate ITE on the first unit $\mathbb{E}(Y_1(1) - Y_1(0) | \mathbf{X}_1)$.

Denote the ‘‘counterfactual dataset’’ by replacing the first sample with $(\mathbf{X}_1, 0, \hat{Y}_1(0))$ as $\hat{\mathbb{D}}$, where $\hat{Y}_1(0)$ is a realization of $Y_1(0)$. Note that we do not have access to the value of $\hat{Y}_1(0)$. However, we do have access to the covariates and treatments of $\hat{\mathbb{D}}$. As these are all the required inputs to Algorithm 1, we can compute the optimal balancing weight for this counterfactual dataset \mathbb{D} , which we denote as $\hat{\alpha}$.

Notice that the sample average treatments of \mathbb{D} are $\hat{\mathbb{D}}$ should be the same, as they are defined for the same set of units. Therefore the two weighted estimators are approximating the same τ_{SATE} (or ATE when N increases) and thus

$$\begin{aligned} & \sum_{i \in \mathbb{T}} \alpha_i \mathbb{E}(Y_i(1) | \mathbf{X}_i) - \sum_{i \in \mathbb{C}} \alpha_i \mathbb{E}(Y_i(0) | \mathbf{X}_i) \\ & \approx \sum_{i \in \mathbb{T} \setminus \{1\}} \hat{\alpha}_i \mathbb{E}(Y_i(1) | \mathbf{X}_i) - \sum_{i \in \mathbb{C}} \hat{\alpha}_i \mathbb{E}(Y_i(0) | \mathbf{X}_i) - \hat{\alpha}_0 \mathbb{E}(\hat{Y}_1(0) | \mathbf{X}_1). \end{aligned}$$

Therefore we have the following approximation

$$\hat{\alpha}_1 \mathbb{E}(\hat{Y}_1(0) | \mathbf{X}_1) \approx -\alpha_1 Y_1(1) + \sum_{i \in \mathbb{T} \setminus \{1\}} (\hat{\alpha}_i - \alpha_i) Y_i(1) - \sum_{i \in \mathbb{C}} (\hat{\alpha}_i - \alpha_i) Y_i(0).$$

As we have access to all individual terms on the right, we can compute an approximation of $\mathbb{E}(Y_1(0) | \mathbf{X}_1)$, using this formula as long as $\hat{\alpha}_0 \neq 0$.²

²Once we have these estimands, policy evaluation can be done via plug-in estimations.

To enhance the robustness of this estimation, we can also compute this for units with covariates closed to \mathbf{X}_1 , e.g., using KNNs (Devroye et al., 1994; Li & Tran, 2009), which would give consistent estimations for conditional expectations. Algorithm 4 summarizes this procedure, where Algorithm 3 can be used instead of Algorithm 1 to estimate ITE in a zero-shot fashion.

Algorithm 4 CInA for ITE.

- 1: **Input:** Covariates \mathbf{X} and treatments \mathbf{W} .
 - 2: **Output:** Estimation of $\mathbb{E}(Y_1(1) - Y_1(0) \mid \mathbf{X}_1)$.
 - 3: Hyper-parameter: penalty weight $\lambda > 0$.
 - 4: Initialize $\tau = \emptyset$.
 - 5: **for** unit i with $\mathbf{X}_i \approx \mathbf{X}_1$ **do**
 - 6: Run Algorithm 1 on \mathbf{X}, \mathbf{W} to obtain α .
 - 7: Set $\hat{\mathbf{W}}$ to be \mathbf{W} except $\hat{W}_i = -W_i$.
 - 8: Run Algorithm 1 on $\mathbf{X}, \hat{\mathbf{W}}$ to obtain $\hat{\alpha}$.
 - 9: Let $\hat{\alpha}_i \mathbb{E}(\hat{Y}_i(1 - T_i) \mid \mathbf{X}_i) = -\alpha_i Y_i(T_i) + \sum_{j \neq i, T_j = T_i} (\hat{\alpha}_j - \alpha_j) Y_j(T_j) - \sum_{T_j \neq T_i} (\hat{\alpha}_j - \alpha_j) Y_j(T_j)$.
 - 10: Append $W_i \cdot (\mathbb{E}(\hat{Y}_i(1 - T_i) \mid \mathbf{X}_i) - Y_i(T_i))$ to τ if $\hat{\alpha}_i \neq 0$.
 - 11: **return** Average of τ .
-

G DATASET DETAILS

The details of the datasets for simulation A are provided in Section 5.1. We now provide the details of ER-5000 and the real-world datasets. Code for downloading and pre-processing these datasets will be provided upon publication.

ER-5000. Each of the ER-5000 datasets is generated following the structural causal model (SCM) framework. The detailed procedure is as follows. First, we sample a random directed acyclic graph (DAG) from the Erdős-Rényi random graph model (Erdős & Rényi, 1960) with edge probability sampled from 0.25 to 0.5. Then, Based on the sampled DAG, we sample the corresponding functional relationships using a linear weight sampler, with random weights sampled from a uniform distribution between 0 and 3. Next, a treatment node and effect node is randomly chosen. For each non-treatment node, we use additive gaussian random noise with standard deviation randomly sampled uniformly between 0.2 and 2. For treatment node, we specify a Bernoulli distribution with logit equal to the functional output of the corresponding node. Finally, we simulate each variable (in \mathbf{X}, T and Y) using the sampled DAG, functional relationships, and noises.

IHDP and IHDP-resampled. The Infant Health and Development Program (IHDP) dataset is a semi-dataset compiled by Hill (2011). We use the existing versions from Chernozhukov et al. (2022), which are sampled using the outcome model implemented as setting A in (Dorie, 2016). Each dataset comprises of 747 units and 25 covariates measuring the aspects of children and their mothers. For IHDP, the treatment group (139 out of 747 units) has been made imbalanced by removing a biased subset of the treated population. A total of 1000 datasets are used (following Shi et al. (2019)), where different datasets only differ in terms of outcome values. For IHDP-resampled, 100 datasets are used where the treatments are resampled by setting the propensity score to “True” in the (Dorie, 2016).

Twins. Introduced by Louizos et al. (2017), this is a semi-synthetic dataset based on the real data on twin births and twin mortality rates in the US from 1989 to 1991 (Almond et al., 2005). The treatment is “born the heavier twin”, which is simulated as a function of the GESTAT10 covariates. Therefore this dataset is confounded. After assigning the treatment for each pair of twins, the dataset is constructed by hiding the other twin. We downloaded the dataset and processed it following Neal et al. (2020).

LaLonde CPS and PSID. We also use the datasets from LaLonde (1986), in which the treatment is job training and the outcomes are income and employment status after training. The ground-truth average treatment effect is computed using a randomized study, where we use the observational data to estimate it. The observational data has multiple versions. We use both the PSID-1 and CPS-1 versions for our experiments (Dehejia & Wahba, 1999).

ACIC. The data for the 2018 Atlantic Causal Inference Conference competition (ACIC) (Shimoni et al., 2018) comprises of several semi-synthetic datasets derived from the linked birth and infant death (LBIDD) data (MacDorman & Atkinson, 1998). The data-generating process is described in (Shimoni et al., 2018). In our experiment, we use datasets containing $1k$ or $10k$ samples.³ In the experiments in Section 5, a total of 293 datasets (each of size $1k$) were used, where 93 were left out for testing. In Appendix I, we extend this to datasets of size $10k$, where a total of 288 datasets were used and 88 among these were left out for testing. We use datasets with polynomial link function for training and validation. For testing, we use datasets with exponential link functions thus creating a harder task for evaluating our methods.

H IMPLEMENTATION DETAILS

Code for our method will be released on GitHub upon publication. Below we describe the architecture, hyper-parameters, training procedures and other details of our method. We also provide the implementation details of the baselines. Finally, we discuss a new data augmentation technique that we observe to be helpful on certain datasets.

H.1 CINA

Pre-processing and Padding. For Algorithm 2, we might encounter multiple datasets with different number of samples. We wish them to share the same transformation from \mathbf{W}, \mathbf{K} to $\mathbf{V} \in \mathbb{R}^{N \times 1}$, where N is the number of units in the corresponding dataset. For this, we adopt similar pre-processing steps as in natural language. We pad all datasets to the same size (i.e., adding dummy units to smaller datasets) and save the masks that indicate these paddings. During back-propagation, we use this mask to make sure that the loss function is only computed using actual units.

Model Configurations. We describe the architecture used in Algorithm 2, as the single-dataset version uses the same components aside from parametrizing the values \mathbf{V} directly as learnable parameters. An illustration of the forward pass is provided in Figure 2.

For the transformation from covariates \mathbf{X} to keys \mathbf{K} , we implemented two versions: (1) an identical mapping followed by a batch-norm layer $\mathbf{K} = \text{bn}(\mathbf{X})$, (2) a projected mapping followed by a batch-norm layer $\mathbf{k}_i = \text{bn} \circ \text{relu} \circ \text{linear}(\mathbf{X}_i)$. In our first simulation study in Section 5.1, we observe that the projection to be marginally helpful and thus report all the results based on the identical mapping.

For the transformation from \mathbf{W}, \mathbf{K} to \mathbf{V} , we first embed $\mathbf{W}_i, \mathbf{k}_i$ into a 32-dimensional space using one layer of $\text{relu} \circ \text{linear}(\cdot)$. These two 32-dimensional vectors are then concatenated into a 64-dimensional vector following by a batch-norm layer. Denote these 64-dimensional embedding for each unit as $\mathbf{E} = [e_1, \dots, e_N]^\top$. We encode them into $N \times 1$ -dimensional outputs \mathbf{O} using a scaled product attention with value, key, query being linear transformations of \mathbf{E} . Notice that we read off the balancing weights via $\mathbf{V}/h(\mathbf{X})\mathbf{W}$ and $h(\mathbf{X}) \succ \mathbf{0}$. As the optimal weights $\alpha^* \succeq \mathbf{0}$, the values \mathbf{V} should have the same sign as \mathbf{W} in an element-wise fashion. Therefore to enforce this, we include another multiplier layer to obtain \mathbf{V} from the outputs \mathbf{O} , namely, $\mathbf{V} = \text{relu}(\mathbf{O}\mathbf{W})$.

Normalization. As the optimal balancing weights is in $\mathbb{A} = \{\mathbf{0} \preceq \alpha \preceq \mathbf{1}, \sum_{i \in \mathbb{T}} \alpha_i = \sum_{i \in \mathbb{C}} \alpha_i = 1\}$, we normalize the read-off balancing weights during inference. In particular, in Algorithm 1 and Algorithm 3, after setting $\alpha^* = \lambda \cdot \mathbf{V}/h(\mathbf{X})\mathbf{W}$, we project it into \mathbb{A} by taking $\max(\alpha^*, \mathbf{0})$ and normalizing the treated and control group to sum up to 1.

Hyper-parameters. For both Algorithm 1 and Algorithm 2, we search for the optimal penalty $\lambda > 0$ from range $[\lambda_{\min}, \lambda_{\max}]$ by exponentially increasing it from λ_{\min} to λ_{\max} . On the same dataset, this range remains the same for both algorithms (and all variations, if applicable). The following table summarizes the values of λ_{\min} to λ_{\max} for different datasets.

Training and Evaluations. For all the experiments, we use a cosine annealing schedule for the learning rate from l_{\max} to l_{\min} during the first half of the training epochs. Then the learning rate is fixed to l_{\min} for the second half of the training epochs. The exact values of l_{\max} and l_{\min} for

³In datasets with large sample sizes, techniques for efficient transformers (Child et al., 2019; Kitaev et al., 2020; Katharopoulos et al., 2020; Sun et al., 2023) can be applied to accelerate our method.

Dataset	λ_{\min}	λ_{\max}
Simulation A	1e-6	1e-2
Simulation B	1e-6	1e-2
IHDP	1	1000
IHDP-resampled	1e-5	1000
Twins	1e-8	1e-2
LaLonde CPS	1e-10	5e-6
LaLonde PSID	1e-10	5e-6
ACIC	1e-6	100

Table 1: Search range for λ in different datasets.

different datasets can be found in the codebase. For Algorithm 1, we train for 20,000 epochs on all datasets. For Algorithm 2, we train for 4,000 epochs on all datasets.

For evaluating the results of Algorithm 2, we choose the best hyper-parameters based on the mean absolute error on the validation sets of datasets and report the results on the testing sets of datasets. For evaluating the results of Algorithm 1, if the setting contains multiple datasets (Simulation A, Simulation B, IHDP-resampled, ACIC), we choose the best hyper-parameters based on the mean absolute error on the validation sets of datasets and report the results on the testing sets of datasets. Note that even though IHDP contains multiple datasets, they all share the same sets of covariates and treatments. Therefore we treat it the same as settings with one dataset for Algorithm 1. On these datasets (IHDP, Twins, LaLonde CPS, LaLonde PSID), we choose the best hyper-parameters based on the reported results.

H.2 BASELINES

IPW and Self-Normalized IPW. For both IPW and self-normalized IPW, we first standardized the covariates \mathbf{X} . Then we fit a random forest classifier on the data to predict propensity scores. The depth of the random forest classifier is chosen in the same way as the hyper-parameter λ is chosen in CInA, which we described above.

DML. For DML, we use the implementation of Battocchi et al. (2019). In particular, we consider three models: `LinearDML`, `CausalForestDML`, `KernelDML`. Similar as above, when a validation set of datasets is present, we report the results based on the best of these three models in terms of validation MAE. Otherwise we report based on the best performance on the reported dataset. However, in simulation A, we only use `LinearDML` as the outcome model is linear.

SVM. For this baseline, we first standardized the covariates \mathbf{X} . Then we solve the dual SVM problem in Eq. (2), where the kernel is defined using ϕ given in Eq. (9) on the standardized data. We use the support vector classifier (Pedregosa et al., 2011) with a precomputed kernel. The maximum number of iterations is capped with a hard limit of 50,000. The reported results are based on λ chosen in the same way as CInA described above.

H.3 DATASET AUGMENTATION

In our experiments in Section 5.1 and certain datasets in Section 5.3 using the multi-dataset version of CInA, we implemented a new type of data augmentation. As we observe that the network can learn how to balance on a set of datasets using very few training steps, we propose to reshuffle amongst different datasets in every epoch. This essentially creates a “new” set of datasets by combining units from different datasets. Intuitively, this augments the number of covariate balancing problems that the model has to learn to solve without actually needing to acquire more data. However, we note that this technique is only applied if different datasets from the same experiment share the same causal graph. If different datasets contain very different causal structures such as **ER-5000** in Section 5.2 and **ACIC** in Section 5.3, this shuffling is not used as it would create covariate balancing problem that does not aid learning. The main intuition is that if we reshuffle units among these datasets, units in a reshuffled dataset could follow different causal graphs, which means there is potentially no underlying causal structure that can explain the data.

I ADDITIONAL EMPIRICAL RESULTS

I.1 COMPARISON TO DRAGONNET AND RIESZNET

Method	Simulation-A	ER-5000	IHDP
Naive	0.172 ± 0.03	50.27 ± 5.97	0.259 ± 0.01
IPW	0.304 ± 0.03	27.42 ± 3.19	0.766 ± 0.02
Self-normalized IPW	0.158 ± 0.03	49.99 ± 5.88	0.141 ± 0.00
DML	0.094 ± 0.01	11.13 ± 3.17	0.585 ± 0.03
DragonNet	0.386 ± 0.01	11.21 ± 3.17	0.146 ± 0.01
RieszNet	0.045 ± 0.01	12.90 ± 4.54	0.110 ± 0.01
SVM	0.015 ± 0.00	11.09 ± 3.13	1.202 ± 0.05
Ours	0.126 ± 0.02	N/A	0.114 ± 0.01
Ours (ZS)	0.147 ± 0.01	11.50 ± 1.85	N/A
Ours (ZS-S)	N/A	2.66 ± 0.33	N/A
Mean	N/A	17.88 ± 1.83	N/A

Table 2: ATE MAE comparison of different methods on the "Simulation-A", "ER-5000", and "IHDP" datasets.

In this section, we further compare two additional baselines, DragonNet (Shi et al., 2019) and RieszNet (Chernozhukov et al., 2022), both of which were considered strong neural estimation methods for per-dataset causal inference. Results for **IHDP** dataset were directly cited from (Shi et al., 2019; Chernozhukov et al., 2022), following their best performing models. Furthermore, we also compare to **Simulation-A-Multi+OOD+diff.size**, and **ER-5000**, both are the most general synthetic settings in Section 5. On **Simulation-A-Multi+OOD+diff.size**, *CINA (ZS)* outperforms *DragonNet*, while *RieszNet* outperforms both *DragonNet* and *CINA (ZS)* method. On both **ER-5000** and **IHDP**, *CINA (ZS)* is on par with or outperforms *DragonNet* and *RieszNet*, while *CINA (ZS-S)* massively outperforms the other methods on **ER-5000**.

I.2 LARGER SCALE EXPERIMENTS ON 10k ACIC 2018, WITH CROSS-DATASET GENERALIZATION

Method	ATE MAE	Inference time on new data (s)	Pretraining time (s)
Naive	13.07 ± 8.25	0.005	N/A
IPW	10.29 ± 5.94	48.927	N/A
Self-normalized IPW	10.30 ± 5.90	49.322	N/A
DML	8.572 ± 8.96	7391.743	N/A
RieszNet	69.39 ± 31.9	8157.498	N/A
Ours (ZS)	1.460 ± 0.48	78.503	1800
Ours (ZS-S)	1.361 ± 0.42	77.546	1800
Ours (ZS-ER)	1.718 ± 0.74	78.085	1800
Ours (ZS-S-ER)	1.702 ± 0.74	77.947	1800

Table 3: Comparison of different methods on the 10k ACIC 2018 dataset.

To demonstrate the performance of our method on larger version of **ACIC 2018**, we produce additional experiment using the 10k-size datasets of ACIC (Shimoni et al., 2018), which is a commonly used scale considered in the literature (Shi et al., 2019; Mahajan et al., 2022). Note that instead of only selecting a subset of datasets in **ACIC 2018** as in (Shi et al., 2019; Mahajan et al., 2022), we make use of all datasets of size 10k generated by (Shimoni et al., 2018) that has polynomial link functions as training datasets, and all datasets of size 10k with exponential link functions as test datasets.

In this setting, we also compare two new variants of our method, *CINA (ZS-ER)* and *CINA (ZS-S-ER)*, that are fully trained on a larger-scale, 200-dimensional ER-5000 dataset Section 5.2 under both unsupervised and supervised settings, respectively. After pre-training, *CINA (ZS-ER)* and

CINA (ZS-S-ER) are applied directly to all ACIC 2018 test sets. This will help us to demonstrate whether the model can show generalization ability across datasets. All CINA-related methods are trained for a fixed time budget (1800 seconds), which is significantly shorter than the full training time of DML and RieszNet. As shown in Table 2, both *CINA (ZS)* and *CINA (ZS-S)* significantly outperforms all baselines. The *CINA (ZS-ER)* and *CINA (ZS-S-ER)* methods give marginally worse performance than *CINA (ZS)* and *CINA (ZS-S)*, but still out-performs the other baselines by a clear margin.