

# FOUNDATION MODELS FOR CAUSAL INFERENCE

FEDERICO BALDO

THE THIRD IPLESP WORKSHOP ON CAUSAL INFERENCE

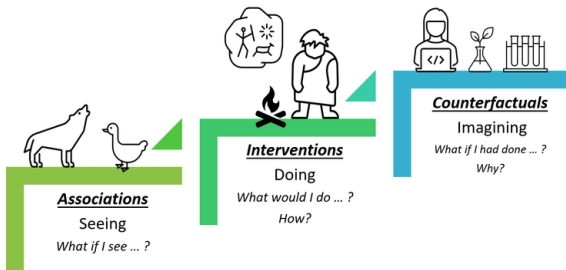
JUNE 2026



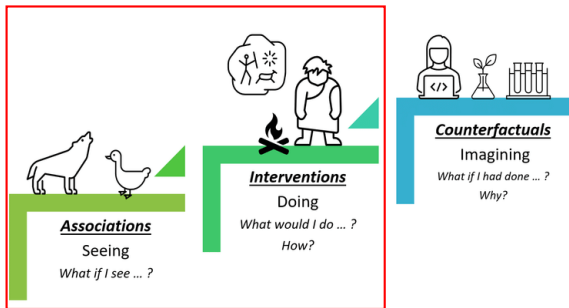
Do you know...

- What is a Neural Network? How it works?
- Transformers? Self-attention?
- What is a Foundation Model?
- TabPFN and how it works?

# THE LADDER OF CAUSATION



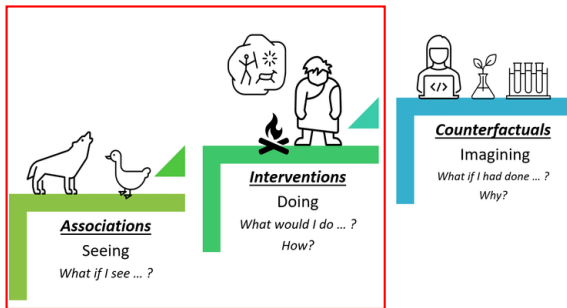
# THE LADDER OF CAUSATION



To climb the ladder:

- Assumptions: Causal DAG, Unconfoundedness.
- In practice, these assumptions are strong and often unverifiable.

# THE LADDER OF CAUSATION



To climb the ladder:

- Assumptions: Causal DAG, Unconfoundedness.
- In practice, these assumptions are strong and often unverifiable.

What if we cannot observe the DAG and we violate unconfoundedness?

# 1

## DEEP LEARNING

Given data  $\mathbb{D}$ , find a function  $f_w$  that captures some structure of  $p(\mathbb{D})$ .

Given data  $\mathbb{D}$ , find a function  $f_w$  that captures some structure of  $p(\mathbb{D})$ .

The structure we seek depends on the problem:

- **Supervised learning:**  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ , learn  $f_w : \mathcal{X} \rightarrow \mathcal{Y}$  to predict  $y$  from  $x$

Given data  $\mathbb{D}$ , find a function  $f_w$  that captures some structure of  $p(\mathbb{D})$ .

The structure we seek depends on the problem:

- **Supervised learning:**  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ , learn  $f_w : \mathcal{X} \rightarrow \mathcal{Y}$  to predict  $y$  from  $x$
- **Unsupervised learning:**  $\mathbb{D} = \{x_i\}_{i=1}^N$ , learn structure in  $p(x)$  (clustering, density estimation)

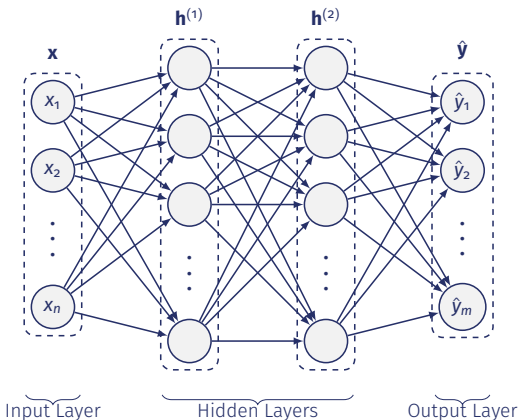
Given data  $\mathbb{D}$ , find a function  $f_w$  that captures some structure of  $p(\mathbb{D})$ .

The structure we seek depends on the problem:

- **Supervised learning:**  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ , learn  $f_w : \mathcal{X} \rightarrow \mathcal{Y}$  to predict  $y$  from  $x$
- **Unsupervised learning:**  $\mathbb{D} = \{x_i\}_{i=1}^N$ , learn structure in  $p(x)$  (clustering, density estimation)
- **Self-supervised learning:** labels derived from  $x$  itself, learn rich representations without human annotation



Artificial Neural Networks (ANNs) are organized in layers of interconnected nodes, inspired by the structure of the human brain.



Neural networks are **universal function approximators**:

## Theorem (Pinkus, 1999 (Simplified))

*For a neural network with a single hidden layer and a wide class of activation functions, for any continuous function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and any  $\varepsilon > 0$ , there exists  $M \in \mathbb{N}$  such that the network  $g$  with  $M$  hidden units satisfy:*

$$\max_{\mathbf{x} \in \Omega} |f(\mathbf{x}) - g(\mathbf{x})| < \varepsilon$$

Neural networks are **universal function approximators**:

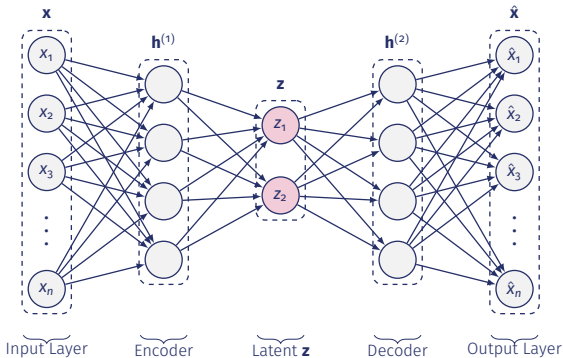
## Theorem (Pinkus, 1999 (Simplified))

*For a neural network with a single hidden layer and a wide class of activation functions, for any continuous function  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  and any  $\varepsilon > 0$ , there exists  $M \in \mathbb{N}$  such that the network  $g$  with  $M$  hidden units satisfy:*

$$\max_{\mathbf{x} \in \Omega} |f(\mathbf{x}) - g(\mathbf{x})| < \varepsilon$$

**Caveat:** the theorem guarantees existence, not learnability.

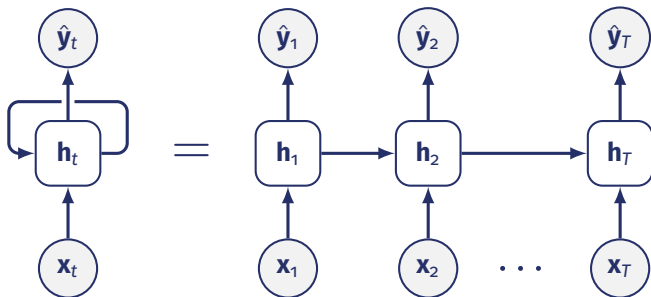
Autoencoders are a type of ANN architecture that learns to compress and reconstruct data.



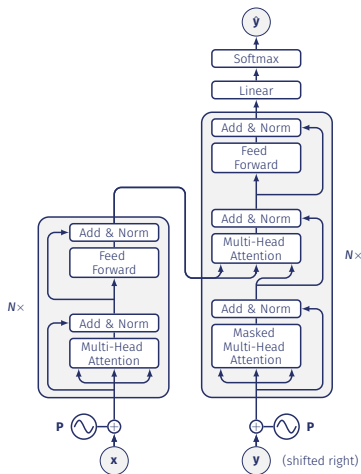
$$\mathcal{L}(w) = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|^2$$

# RECURRENT NEURAL NETWORKS (RNNs)

Recurrent Neural Networks (RNNs) are designed to handle sequential data by maintaining a hidden state that captures information from previous time steps.

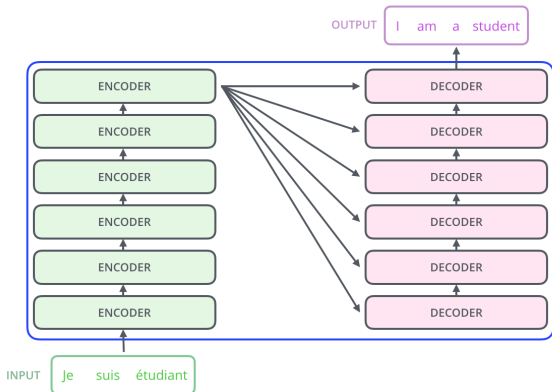


A Transformer<sup>1</sup> model is a neural network that learns the context of sequential data and generates new data out of it.



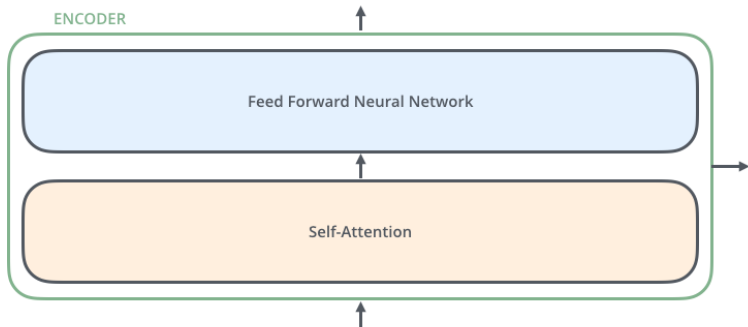
<sup>1</sup>Vaswani et al. "Attention is all you need." NeurIPS (2017).

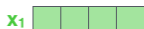
A Transformer<sup>1</sup> model is a neural network that learns the context of sequential data and generates new data out of it.



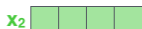
<sup>1</sup>Vaswani et al. "Attention is all you need." NeurIPS (2017).

We can simplify removing residual connections and layer normalization, but they are crucial for training deep transformers.

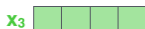




Je

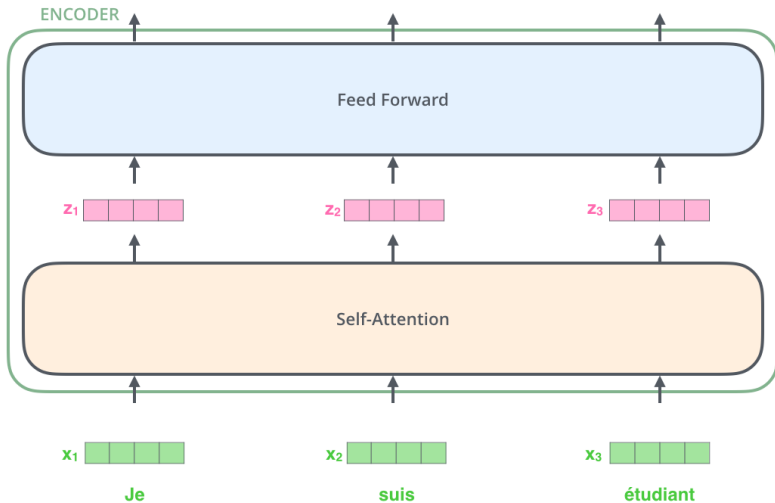


suis



étudiant

- The input sequence is tokenized and converted into embeddings, i.e., dense vector representations.
  - ▶ Frequency-based Embeddings (e.g., Bag of Words, TF-IDF)
  - ▶ Predictive Embeddings (e.g., Word2Vec, GloVe)
  - ▶ Contextual Embeddings (e.g., BERT, GPT)



**Intuition:** You are in a library looking for information.

- **Query:** the question you bring to the library
- **Key:** the label on each book's spine, what it is about
- **Value:** the content inside, what you actually read

You compare your question against every label, read the most relevant books more carefully, and aggregate what you learned.

**Intuition:** You are in a library looking for information.

- **Query:** the question you bring to the library
- **Key:** the label on each book's spine, what it is about
- **Value:** the content inside, what you actually read

You compare your question against every label, read the most relevant books more carefully, and aggregate what you learned.

In self-attention, each token plays all three roles simultaneously:

- **Q:** what this token is looking for
- **K:** what this token advertises about itself
- **V:** what this token contributes if selected

**Intuition:** You are in a library looking for information.

- **Query:** the question you bring to the library
- **Key:** the label on each book's spine, what it is about
- **Value:** the content inside, what you actually read

You compare your question against every label, read the most relevant books more carefully, and aggregate what you learned.

In self-attention, each token plays all three roles simultaneously:

- **Q:** what this token is looking for
- **K:** what this token advertises about itself
- **V:** what this token contributes if selected

“The animal didn't cross the street because **it** was too tired.”  
it queries all tokens → animal has the highest matching key → its value gets the most weight.

# MULTI-HEAD ATTENTION

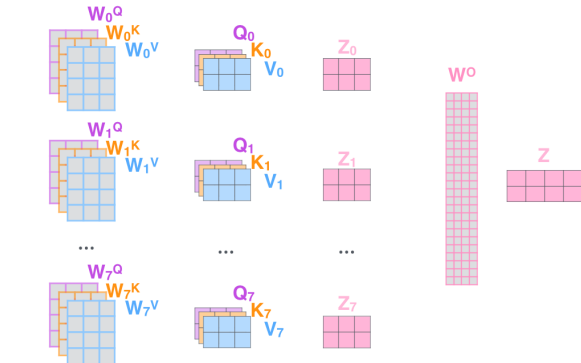
Multi-head attention captures different types of relationships within a sequence simultaneously.

- 1) This is our input sentence\*
- 2) We embed each word\*
- 3) Split into 8 heads. We multiply  $X$  or  $R$  with weight matrices
- 4) Calculate attention using the resulting  $Q/K/V$  matrices
- 5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer

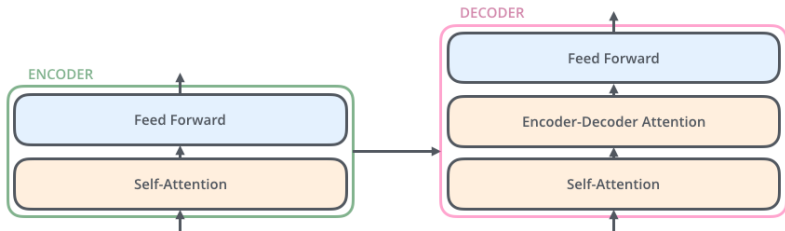
Thinking  
Machines



\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one



The decoder block is composed of layers similar to those in the encoder, plus a cross-attention layer.



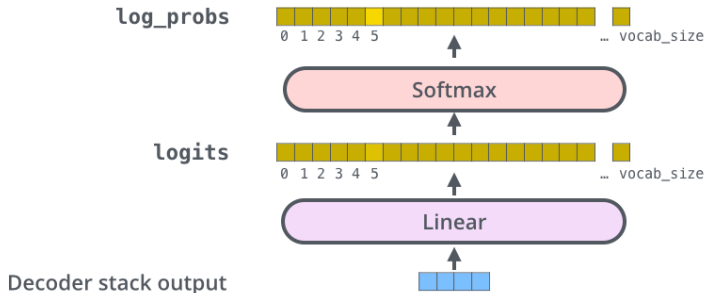
# TRANSFORMERS - OUTPUT

Which word in our vocabulary  
is associated with this index?

am

Get the index of the cell  
with the highest value  
(**argmax**)

5



A transformer only sees a sequence of vectors, any modality (data type) works as long as it can be tokenized and embedded into a common space.

- Each modality gets its own tokenization and encoder (patches, frames, rows...).
- A shared or aligned latent space connects modalities.
- Cross-attention lets the decoder attend to the encoded source modality.

# 2

## FOUNDATION MODEL FOR TABULAR DATA

A foundation model is a large neural network trained on broad, large-scale data that can be adapted to a wide range of downstream tasks.

- **Emergent capabilities:** abilities not explicitly trained for (reasoning, in-context learning) arise beyond certain scale thresholds.
- **Transfer:** a single pretrained model serves as the starting point for many tasks, drastically reducing the cost of specialisation.

Tabular data: Structured data organized in rows and columns, where each row is an observation and each column is a feature.

---

<sup>2</sup>Hollmann et al. "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second." ICLR 2023.

Tabular data: Structured data organized in rows and columns, where each row is an observation and each column is a feature.

Can we train Transformers with tabular data?

---

<sup>2</sup>Hollmann et al. "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second." ICLR 2023.

Tabular data: Structured data organized in rows and columns, where each row is an observation and each column is a feature.

Can we train Transformers with tabular data?

- Yes<sup>2</sup>, if we treat each row as a sequence of (feature, value) tokens, i.e.  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ .

---

<sup>2</sup>Hollmann et al. "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second." ICLR 2023.

Tabular data: Structured data organized in rows and columns, where each row is an observation and each column is a feature.

Can we train Transformers with tabular data?

- Yes<sup>2</sup>, if we treat each row as a sequence of (feature, value) tokens, i.e.  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ .

Where can we find a large sample of tabular datasets?

---

<sup>2</sup>Hollmann et al. "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second." ICLR 2023.

Tabular data: Structured data organized in rows and columns, where each row is an observation and each column is a feature.

Can we train Transformers with tabular data?

- Yes<sup>2</sup>, if we treat each row as a sequence of (feature, value) tokens, i.e.  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ .

Where can we find a large sample of tabular datasets?

- We can generate them !

---

<sup>2</sup>Hollmann et al. "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second." ICLR 2023.

Tabular data: Structured data organized in rows and columns, where each row is an observation and each column is a feature.

Can we train Transformers with tabular data?

- Yes<sup>2</sup>, if we treat each row as a sequence of (feature, value) tokens, i.e.  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ .

Where can we find a large sample of tabular datasets?

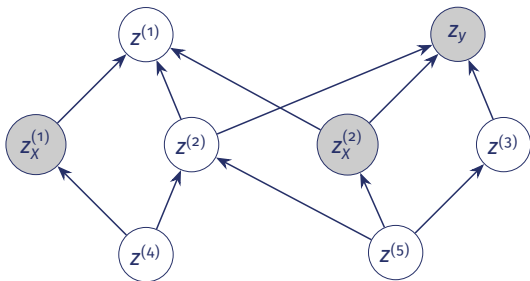
- We can generate them ! ...but we need a prior

---

<sup>2</sup>Hollmann et al. "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second." ICLR 2023.

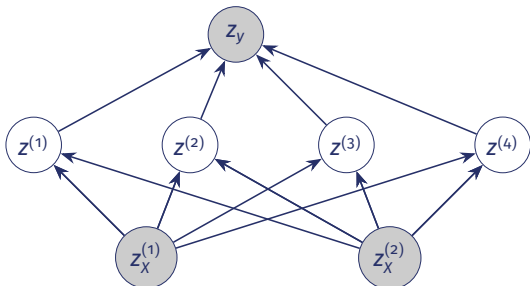
We generate synthetic datasets by sampling from a prior over data-generating processes:

1. Random Structural Causal Models (SCMs) as a generative process:
  - Select a subset of nodes as covariates,  $z_X^{(i)}$ , and a downstream node as target,  $z_y$ .



We generate synthetic datasets by sampling from a prior over data-generating processes:

2. Bayesian Neural Networks (BNNs) are neural models in which each weight has a distribution.



We generate synthetic datasets by sampling from a prior over data-generating processes:

- The generated SCMs are not training examples in the usual ML sense.

We generate synthetic datasets by sampling from a prior over data-generating processes:

- The generated SCMs are not training examples in the usual ML sense.
- Collectively, these generated SCMs define a prior over data-generating mechanisms.

We have a dataset,  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ , and we have an additional point  $x$ . We want to predict  $y$ :

We have a dataset,  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ , and we have an additional point  $x$ . We want to predict  $y$ :

$$p(y | x, \mathbb{D}) = \int p(y | x, \mathcal{M})p(\mathcal{M} | \mathbb{D}), d\mathcal{M}$$

- $\mathcal{M}$  is an SCM.
- The PPD weight all plausible models according to the data.

We have a dataset,  $\mathbb{D} = \{(x_i, y_i)\}_{i=1}^N$ , and we have an additional point  $x$ . We want to predict  $y$ :

$$p(y | x, \mathbb{D}) = \int p(y | x, \mathcal{M})p(\mathcal{M} | \mathbb{D}), d\mathcal{M}$$

- $\mathcal{M}$  is an SCM.
- The PPD weight all plausible models according to the data.

Given everything I have seen in  $\mathbb{D}$ , and given all the causal mechanisms that could have generated it, what is my best guess for  $y$ ?

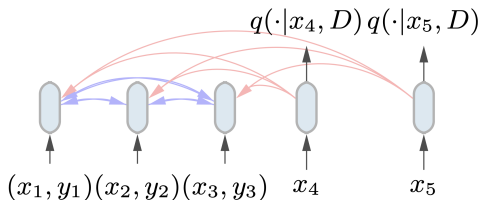
We can use a Transformer to approximate the PPD:

$$\mathcal{L}_{\text{PFN}} = \mathbb{E}_{\{(x,y)\} \cup \mathbb{D} \sim p_{\mathcal{M}}} [-\log q_w(\mathbf{y} \mid \mathbf{x}, \mathbb{D}_{\text{train}})]$$

- $q_w$  is the predictive distribution computed by the network with parameters  $w$ .
- Prediction requires only a **single forward pass**:  
 $q_w(\mathbf{y} \mid \mathbf{x}_{\text{test}}, \mathbb{D}_{\text{train}})$  — the entire training set is the context.

Because of self-attention!

- Each  $(x_i, y_i)$  pair is treated as a single **token**, so representations can attend to each other across the full dataset.
- The dataset represents the context to identify the right SCMs.
  - ▶ This does not require retraining!



*“Our work can be considered as rung 1.5 [...] we do not perform causal reasoning, but make association-based predictions on observational data assuming SCMs model common datasets well.”*

*“Our work can be considered as rung 1.5 [...] we do not perform causal reasoning, but make association-based predictions on observational data assuming SCMs model common datasets well.”*

*Since TabPFN is trained on data generated from SCMs, can we try estimate directly causal effects?*

Do-PFN<sup>3</sup> is a transformer that estimates causal effects in one forward pass from observational data alone, with no causal graph required.

- It approximates the **conditional interventional distribution (CID)** from observational data  $\mathbb{D}^{\text{obs}} = \{(t_i, \mathbf{x}_i, y_i)\}_{i=1}^N$

---

<sup>3</sup>Robertson, Reuter et al. “Do-PFN: In-Context Learning for Causal Effect Estimation.” NeurIPS 2025.

Do-PFN<sup>3</sup> is a transformer that estimates causal effects in one forward pass from observational data alone, with no causal graph required.

- It approximates the **conditional interventional distribution (CID)** from observational data  $\mathbb{D}^{\text{obs}} = \{(t_i, \mathbf{x}_i, y_i)\}_{i=1}^N$

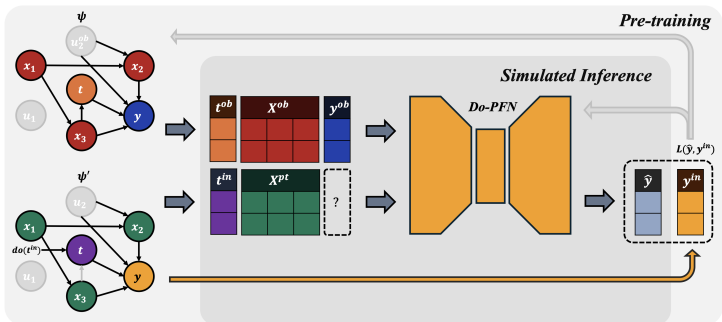
$$p(\tilde{y} \mid do(\tilde{t}), \mathbf{x}, \mathbb{D}^{\text{obs}}) = \int p(\tilde{y} \mid do(\tilde{t}), \mathbf{x}, \mathcal{M}) p(\mathcal{M} \mid \mathbb{D}^{\text{obs}}) d\mathcal{M}$$

---

<sup>3</sup>Robertson, Reuter et al. “Do-PFN: In-Context Learning for Causal Effect Estimation.” NeurIPS 2025.

Do-PFN<sup>3</sup> is a transformer that estimates causal effects in one forward pass from observational data alone, with no causal graph required.

- It approximates the **conditional interventional distribution (CID)** from observational data  $\mathbb{D}^{obs} = \{(t_i, x_i, y_i)\}_{i=1}^N$



<sup>3</sup>Robertson, Reuter et al. “Do-PFN: In-Context Learning for Causal Effect Estimation.” NeurIPS 2025.

We approximate the CID by training a Transformer on pairs of observational and interventional datasets:

$$\mathcal{L}(w) = \mathbb{E}_{(\tilde{t}, x, \tilde{y}) \sim p_{\mathcal{M}}^{do} \cup \mathbb{D}^{obs} \sim p_{\mathcal{M}}} [-\log q_w(\tilde{y} \mid do(\tilde{t}), x, \mathbb{D}^{obs})]$$

We approximate the CID by training a Transformer on pairs of observational and interventional datasets:

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{(\tilde{\mathbf{t}}, \mathbf{x}, \tilde{\mathbf{y}}) \sim p_{\mathcal{M}}^{do} \cup \mathbb{D}^{obs} \sim p_{\mathcal{M}}} [-\log q_{\mathbf{w}}(\tilde{\mathbf{y}} \mid do(\tilde{\mathbf{t}}), \mathbf{x}, \mathbb{D}^{obs})]$$

Performing Gradient Descent on  $\mathcal{L}(\mathbf{w})$  is equivalent to minimising the KL-divergence between  $q_{\mathbf{w}}$  and the true CID  $p(\tilde{\mathbf{y}} \mid \mathbf{x}, do(\tilde{\mathbf{t}}), \mathcal{M})$  (Proposition 1).

We approximate the CID by training a Transformer on pairs of observational and interventional datasets:

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{(\tilde{t}, x, \tilde{y}) \sim p_{\mathcal{M}}^{do} \cup \mathbb{D}^{obs} \sim p_{\mathcal{M}}} [-\log q_{\mathbf{w}}(\tilde{y} \mid do(\tilde{t}), x, \mathbb{D}^{obs})]$$

Performing Gradient Descent on  $\mathcal{L}(\mathbf{w})$  is equivalent to minimising the KL-divergence between  $q_{\mathbf{w}}$  and the true CID  $p(\tilde{y} \mid x, do(\tilde{t}), \mathcal{M})$  (Proposition 1).

## Important caveat

This does **not** mean every causal effect is identifiable. When hidden confounders are present, Do-PFN cannot resolve the ambiguity.

Does  $q_w$  recover the true interventional distribution as we observe more data?

Does  $q_w$  recover the true interventional distribution as we observe more data?

As  $|\mathbb{D}^{obs}| \rightarrow \infty$ , the posterior  $p(\mathcal{M} | \mathbb{D}^{obs})$  concentrates on the true SCM,  $\mathcal{M}_0$ :

Does  $q_w$  recover the true interventional distribution as we observe more data?

As  $|\mathbb{D}^{obs}| \rightarrow \infty$ , the posterior  $p(\mathcal{M} | \mathbb{D}^{obs})$  concentrates on the true SCM,  $\mathcal{M}_o$ :

$$p(\tilde{y} | do(\tilde{t}), x, \mathbb{D}^{obs}) \xrightarrow{|\mathbb{D}^{obs}| \rightarrow \infty} p(\tilde{y} | do(\tilde{t}), x, [\mathcal{M}_o])$$

where  $[\mathcal{M}_o]$  is the Markov-equivalence class of  $\mathcal{M}_o$ .

Does  $q_w$  recover the true interventional distribution as we observe more data?

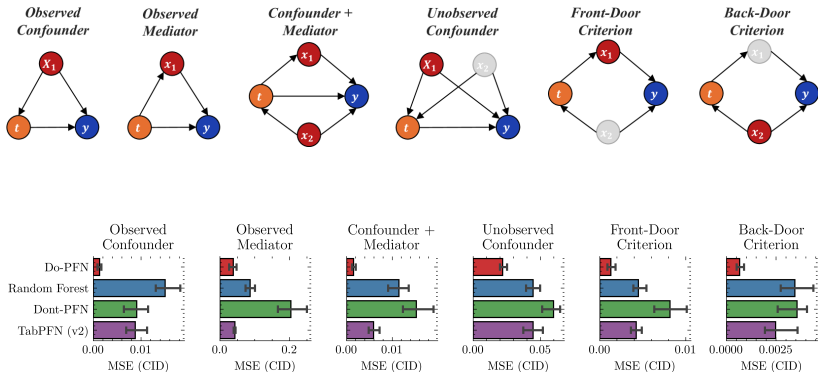
As  $|\mathbb{D}^{obs}| \rightarrow \infty$ , the posterior  $p(\mathcal{M} | \mathbb{D}^{obs})$  concentrates on the true SCM,  $\mathcal{M}_o$ :

$$p(\tilde{y} | do(\tilde{t}), x, \mathbb{D}^{obs}) \xrightarrow{|\mathbb{D}^{obs}| \rightarrow \infty} p(\tilde{y} | do(\tilde{t}), x, [\mathcal{M}_o])$$

where  $[\mathcal{M}_o]$  is the Markov-equivalence class of  $\mathcal{M}_o$ .

Why? Observational data alone cannot distinguish between SCMs in the same Markov-equivalence class

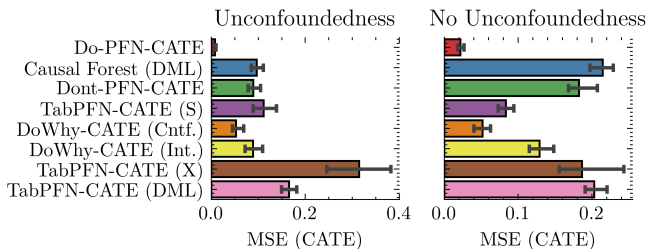
The graph structures of six causal case studies, requiring Do-PFN to automatically perform adjustment.



# DO-PFN UNDER CONFOUNDEDNESS

Do-PFN is trained on SCMs with and without hidden confounders.

- It cannot resolve unidentifiability.
- But when confounding is plausible, it outputs a wider predictive distribution rather than a falsely confident estimate.



- Transformers can process tabular data as **(feature, value)** token sequences.
- **TabPFN** generalizes across unseen datasets by training on synthetic SCMs.
- **Do-PFN** estimates causal effects from observational data with no causal graph required.
  - ▶ Non-identifiable instances produce higher uncertainty in the results

- Transformers can process tabular data as (**feature, value**) token sequences.
- **TabPFN** generalizes across unseen datasets by training on synthetic SCMs.
- **Do-PFN** estimates causal effects from observational data with no causal graph required.
  - ▶ Non-identifiable instances produce higher uncertainty in the results

## Limitations

- Sensitive to the choice of synthetic prior.
- Evaluation mostly on synthetic data, no real-world ground truth.
- Amortized inference does not provide formal guarantees.
- Many configurations were not considered, e.g., non-binary interventions.

## Books & Papers:

- "Deep Learning", Goodfellow and Bengio
- "Attention is all you need", Vaswani et al.
- "The Illustrated Transformer", Jay Alamar
- "TabPFN: A Transformer That Solves Small Tabular Classification Problems in a Second", Hollmann et al.
- "Do-PFN: In-Context Learning for Causal Effect Estimation", Robertson et al.

## Code:

- `tabpfn` python package
- `tabpfn` R code
- `dopf` python code

QUESTIONS?

