# Quantifying Consistency and Information Loss for Causal Abstraction Learning

F.M. Zennaro, P. Turrini, T. Damoulas
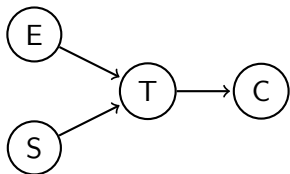
*University of Warwick*

32nd International Joint Conference on Artificial Intelligence
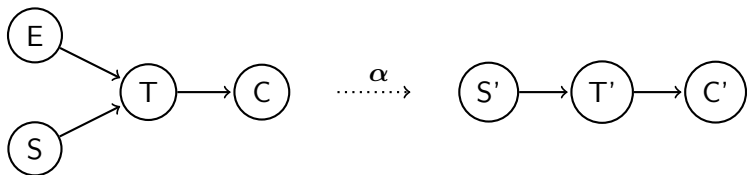
# Structural Causal Models

A *structural causal model* (SCM) $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$ is a mathematical object representing a causal system [2, 3].

A SCM is associated with a *directed acyclic graph* (DAG)

## Abstractions

The same causal system may be represented at different *levels of abstraction* [1].



Given two SCMs we want a formal **abstraction map** $\alpha$ between them.

- ✓ rely on multi-scale representations
- ✓ transfer data between different resolutions
- ✓ scale computational expense

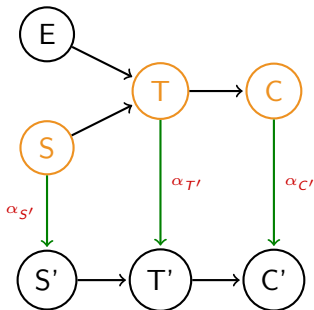An **abstraction** $\alpha$ is a tuple
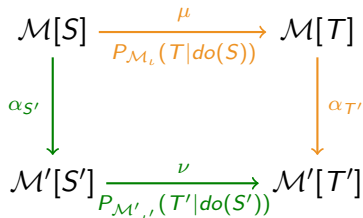
$$\langle R, a, \alpha_i \rangle$$

where:

- $R$ is a set of relevant nodes/variables;
- $a$ is a surjective function between *variables*;
- $\alpha_i$ is a collection of surjective functions between *outcomes*.

We evaluate the *quality* of an abstraction in terms of *interventional consistency*.

The **abstraction error** wrt $P(\mathbf{Y'}|do(\mathbf{X'}))$ is the maximum *distance between interventional distributions* in the base and abstracted model.

$$\mathcal{M}[S] \xrightarrow[\;P_{\mathcal{M}_\iota}(T|do(S))\;]{\mu} \mathcal{M}[T]$$

$$\alpha_{S'} \downarrow \qquad\qquad \downarrow \alpha_{T'}$$

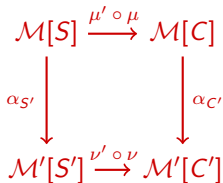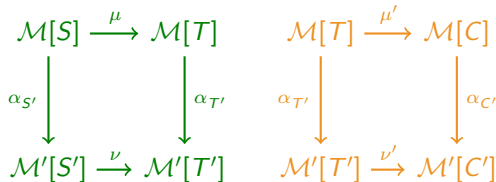$$\mathcal{M}'[S'] \xrightarrow[\;P_{\mathcal{M}'_{\iota'}}(T'|do(S'))\;]{\nu} \mathcal{M}'[T']$$

$$E(\alpha, S', T') = \max_{s \in \mathcal{M}[S]} D_{JSD}(\alpha_{T'} \cdot \mu, \nu \cdot \alpha_{S'})$$

An abstraction implies multiple *abstraction errors*.

**(Global) abstraction error** $e(\alpha)$ is the maximum abstraction error over all disjoint sets of variables.

$$\mathcal{M}[S] \xrightarrow{\mu} \mathcal{M}[T] \qquad \mathcal{M}[T] \xrightarrow{\mu'} \mathcal{M}[C]$$

$$\alpha_{S'} \downarrow \qquad \downarrow \alpha_{T'} \qquad \alpha_{T'} \downarrow \qquad \downarrow \alpha_{C'}$$

$$\mathcal{M}'[S'] \xrightarrow{\nu} \mathcal{M}'[T'] \qquad \mathcal{M}'[T'] \xrightarrow{\nu'} \mathcal{M}'[C']$$

$$\mathcal{M}[S] \xrightarrow{\mu' \circ \mu} \mathcal{M}[C]$$

$$\alpha_{S'} \downarrow \qquad \downarrow \alpha_{C'}$$

$$\mathcal{M}'[S'] \xrightarrow{\nu' \circ \nu} \mathcal{M}'[C']$$

$$e(\alpha) = \sup_{\mathbf{X}',\mathbf{Y}' \subseteq \mathcal{X}'} E(\alpha, \mathbf{X}', \mathbf{Y}')$$

# Generalizing Abstraction Error

The abstraction error can be
expressed more generally as:

$$E_{\boldsymbol{\alpha}}(\mathbf{X}', \mathbf{Y}') = \underset{x' \in \mathbf{X}'}{\mathrm{agg}} \; D(p, q)$$

$$e(\boldsymbol{\alpha}) = \underset{(\mathbf{X}', \mathbf{Y}') \in \mathcal{J}}{\mathrm{agg}} E_{\boldsymbol{\alpha}}(\mathbf{X}', \mathbf{Y}')$$

parametrized by aggregation
functions, distances, paths,
intervention sets, and pseudo-inverse.

$$
\begin{array}{ccc}
\mathcal{M}[S] & \xrightarrow{\mu} & \mathcal{M}[T] \\
\alpha_{S'} \Big\Updownarrow \alpha_{S'}^{+} & & \alpha_{T'} \Big\Updownarrow \alpha_{T'}^{+} \\
\mathcal{M}'[S'] & \xrightarrow{\nu} & \mathcal{M}'[T']
\end{array}
$$

# A new family of errors

**Interventional consistency (IC)**

$$\mathcal{M}[S] \xrightarrow{\mu} \mathcal{M}[T]$$

$\alpha_{S'}$ $\left(\quad\right)$ $\alpha_{S'}^{+}$ $\quad\alpha_{T'}$ $\left(\quad\right)$ $\alpha_{T'}^{+}$

$$\mathcal{M}'[S'] \xrightarrow{\nu} \mathcal{M}'[T']$$

*Consistency projected on the abstracted model.*

**Interventional information loss (IIL)**

$$\mathcal{M}[S] \xrightarrow{\mu} \mathcal{M}[T]$$

$\alpha_{S'}$ $\left(\quad\right)$ $\alpha_{S'}^{+}$ $\quad\alpha_{T'}$ $\left(\quad\right)$ $\alpha_{T'}^{+}$

$$\mathcal{M}'[S'] \xrightarrow{\nu} \mathcal{M}'[T']$$

*Loss in abstracting and reconstructing.*

# A new family of errors

**Interventional superresolution information loss (ISIL)**

$$\mathcal{M}[S] \xrightarrow{\ \mu\ } \mathcal{M}[T]$$

$$\alpha_{S'} \left(\ \right) \alpha_{S'}^{+} \qquad \alpha_{T'} \left(\ \right) \alpha_{T'}^{+}$$

$$\mathcal{M}'[S'] \xrightarrow{\ \nu\ } \mathcal{M}'[T']$$

*Loss in reconstructing and abstracting.*

**Interventional superresolution consistency (ISC)**

$$\mathcal{M}[S] \xrightarrow{\ \mu\ } \mathcal{M}[T]$$

$$\alpha_{S'} \left(\ \right) \alpha_{S'}^{+} \qquad \alpha_{T'} \left(\ \right) \alpha_{T'}^{+}$$

$$\mathcal{M}'[S'] \xrightarrow{\ \nu\ } \mathcal{M}'[T']$$

*Consistency projected on the base model.*

# In the paper...

https://arxiv.org/abs/2305.04357

- Properties of the errors (IC, IIL, ISIL, ISC)
- Discussion of other error measure parameters
- Algorithms for evaluating and learning abstractions
- Empirical evaluation

https://github.com/FMZennaro/CausalAbstraction/tree/main/
papers/2023-quantifying-consistency-and-infoloss

10/12

Thank you for your attention!

[1] Luciano Floridi. The method of levels of abstraction. *Minds and machines*, 18(3):303–329, 2008.

[2] Judea Pearl. *Causality*. Cambridge University Press, 2009.

[3] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: Foundations and learning algorithms*. MIT Press, 2017.

[4] Eigil Fjeldgren Rischel. The category theory of causal models. 2020.