From Structural Causal Models to Causal Abstraction Learning

Fabio Massimo Zennaro

University of Bergen

May 30, 2025



- 2 Causal Abstraction
- 3 Abstraction Learning

4 Current Developments

Assume we want to model a system.



Assume we want to model a system.



Assume we want to model a system.



Assume we want to model a system.



Assume we want to model a system.



Pri	or ,	、
		ODE
		Fitted Linear
		SCM
		Bayes Net
		Neural Net
		Data

Prior ,	↑	
	ODE Fitted Linear	 It discriminates <i>correlations</i> and <i>causes</i>.
	SCM	
	Bayes Net	
	Neural Net	
	· · · · · · · · · · · · · · · · · · ·	
	Data	

P۱	rior ,	`		
		ODE	٠	lt
		Fitted Linear		Cá
		SCM	۹	lt in
		Bayes Net		
		Neural Net		
		Data		

- It discriminates *correlations* and *causes*.
- It allows for reasoning about *interventions*.



- It discriminates *correlations* and *causes*.
- It allows for reasoning about *interventions*.
- It allows for reasoning about *counterfactuals*.



- It discriminates *correlations* and *causes*.
- It allows for reasoning about *interventions*.
- It allows for reasoning about *counterfactuals*.
- It implies a *causality ladder* of reasoning.

A Motivating Example

SCMs represent causal systems.



A Motivating Example

SCMs represent causal systems.



SCMs integrates a graphical model and probabilities distributions.

Structural Causal Models (SCMs) - Definition



Structural Causal Models (SCMs) - Definition

We express a **SCM** as $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$ [16, 17]:



X: set of *endogenous nodes* (S, T, C) representing variables of interest

Structural Causal Models (SCMs) - Definition



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U_S, U_T, U_C) representing stochastic factors

Structural Causal Models (SCMs) - Definition



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U_S, U_T, U_C) representing stochastic factors
- \mathcal{F} : Set of *structural functions* (f_S, f_T, f_C) describing the dynamics of each variable

Structural Causal Models (SCMs) - Definition



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U_S, U_T, U_C) representing stochastic factors
- *F*: Set of *structural functions* (*f_S*, *f_T*, *f_C*) describing the dynamics of each variable
- \mathcal{P} : Set of *distributions* (P_S, P_T, P_C) describing the random factors

Structural Causal Models (SCMs) - Definition

We express a **SCM** as $\mathcal{M} = \langle \mathcal{X}, \mathcal{U}, \mathcal{F}, \mathcal{P} \rangle$ [16, 17]:



- X: set of *endogenous nodes* (S, T, C) representing variables of interest
- U: Set of *exogenous nodes* (U_S, U_T, U_C) representing stochastic factors
- *F*: Set of *structural functions* (*f_S*, *f_T*, *f_C*) describing the dynamics of each variable
- \mathcal{P} : Set of *distributions* (P_S, P_T, P_C) describing the random factors

Every SCM \mathcal{M} implies a (joint) distribution $P_{\mathcal{M}}$: $P_{\mathcal{M}}(S, T, C)$

Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [16, 17]:



1

Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [16, 17]:



do(*T* = 1)

Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [16, 17]:



do(T=1)

2

Remove incoming edges in the intervened node

Structural Causal Models (SCMs) - Interventions

We can perform interventions on a causal model [16, 17]:



do(T=1)

- Remove incoming edges in the intervened node
- Set the value of the intervened node





An *intervention* ι defines a new **intervened model** \mathcal{M}_{ι} with new distributions.



 $P_{\mathcal{M}}$







2. Causal Abstraction

Levels of Abstraction

Systems may be represented at different levels of abstraction (LoA) [9].

Levels of Abstraction

Systems may be represented at different levels of abstraction (LoA) [9].

Thermodynamics example:

Low-level / Base model: Microscopic description **x**, **x**. High-level / Abstracted model: Macroscopic description P, T, V.

Levels of Abstraction

Systems may be represented at different levels of abstraction (LoA) [9].

Thermodynamics example:

Low-level / Base model: Microscopic description **x**, **x**. High-level / Abstracted model: Macroscopic description P, T, V.

LoA may be inaccessible, so we may want to *shift* among LoAs.

- We need a *mapping* between LoAs.
- We want the mapping to be consistent.
Abstraction (aka, multi-level modelling or multi-resolution modelling) aims at relating these levels.



Abstraction (aka, multi-level modelling or multi-resolution modelling) aims at relating these levels.



• It combines models from *different sources*.

Abstraction (aka, *multi-level modelling* or *multi-resolution modelling*) aims at relating these levels.



- It combines models from *different sources*.
- It aggregates information from *different resolutions*.

Abstraction (aka, *multi-level modelling* or *multi-resolution modelling*) aims at relating these levels.



- It combines models from *different sources*.
- It aggregates information from *different resolutions*.
- It allows for *computation with minimal effort*.

Lung cancer scenario example:



Lung cancer scenario example:



Lung cancer scenario example:



Lung cancer scenario example:



- The *transformation* approach [20, 2]
- The α -abstraction approach [19, 18]

• The Φ-abstraction approach [14, 15]

Causal Abstractions (CAs) - Definition

An α -abstraction $\langle R, a, \alpha_i \rangle$ [19, 18] is defined as:

Causal Abstractions (CAs) - Definition

An α -abstraction $\langle R, a, \alpha_i \rangle$ [19, 18] is defined as:



R: a set of *relevant* variables;



Causal Abstractions (CAs) - Definition

An α -abstraction $\langle R, a, \alpha_i \rangle$ [19, 18] is defined as:



- *R*: a set of *relevant variables*;
- *a*: a surjective function between *variables*;

Causal Abstractions (CAs) - Definition

An α -abstraction $\langle R, a, \alpha_i \rangle$ [19, 18] is defined as:



- *R*: a set of *relevant* variables;
- a: a surjective function between variables;
- *α_i*: a collection of surjective functions between *outcomes*.

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





We want an abstraction to guarantee *interventional consistency*.



• Ideally, mechanisms and abstractions *commute*.



- Ideally, mechanisms and abstractions *commute*.
- Otherwise, we compute an abstraction error as the *worst-case discrepancy* over all possible interventions:

$$E_{\alpha}(S',T') = \max_{\iota} D(\alpha_{T'} \cdot \mu, \nu \cdot \alpha_{S'})$$

Causal Abstractions (CAs) - Abstraction Error

Causal Abstractions (CAs) - Abstraction Error



Causal Abstractions (CAs) - Abstraction Error



Causal Abstractions (CAs) - Abstraction Error



Causal Abstractions (CAs) - Abstraction Error

An abstraction implies multiple causal mechanism diagrams:



A (global) abstraction error [19] $e(\alpha)$ is the maximum abstraction error over all diagrams.

$$\mathsf{e}(oldsymbol{lpha}) = \sup_{\mathbf{X}',\mathbf{Y}'\subseteq\mathcal{X}'} \mathsf{E}_{oldsymbol{lpha}}(\mathbf{X}',\mathbf{Y}')$$

3. Abstraction Learning

Joint work of FMZ, M. Drávucz, G. Apachitei, W.D. Widanage and T. Damoulas

Problem statement [25]

Given a partially define *abstraction* α in terms of $\langle R, a \rangle$ can I learn α_i as:

$$\min_{\alpha} e(\alpha)$$



Challenges [25]

(i) Multiple related problems



Challenges [25]

(i) Multiple related problems

(ii) Combinatorial optimization



Challenges [25]

- (i) *Multiple related* problems
- (ii) Combinatorial optimization
- (iii) Surjectivity constraints



 α

Challenges [25]

- (i) Multiple related problems
- (ii) Combinatorial optimization
- (iii) Surjectivity constraints
- Baselines: parallel or sequential approaches.

Abstraction Learning

Relaxation and parametrization [25]

We address (ii) combinatorial optimization by relaxing and parametrizing all α_i .

$$\min_{\alpha(\mathsf{W})} e(\alpha(\mathsf{W}))$$

 $\alpha_{S'}, \alpha_{T'}, \alpha_{C'} \in \mathbb{R}^{2 \times 2}$ $\begin{bmatrix} 0.7 & 1.2 \\ -0.2 & 3.3 \end{bmatrix}$

Abstraction Learning

Relaxation and parametrization [25]

We address (ii) combinatorial optimization by relaxing and parametrizing all α_i .

 $\alpha_{\mathcal{S}'}, \alpha_{\mathcal{T}'}, \alpha_{\mathcal{C}'} \in \mathbb{R}^{2 \times 2}$

 $\min_{\boldsymbol{\alpha}(\mathbf{W})} e(\boldsymbol{\alpha}(\mathbf{W})) \begin{bmatrix} 0.7 & 1.2 \\ -0.2 & 3.3 \end{bmatrix}$

We add *tempering* $t(W) = \frac{e^{\frac{W_i J}{T}}}{\sum_i e^{\frac{W_i J}{T}}}$ along the matrix columns to binarize them.

 $\mathcal{L}_1 : \min_{\alpha(\mathbf{W})} e(\alpha(t(\mathbf{W})))$

 $\alpha_{\mathcal{S}'}, \boldsymbol{\alpha_{\mathcal{T}'}}, \boldsymbol{\alpha_{\mathcal{C}'}} \in [0, 1]^{2 \times 2}$

$$t\left(\left[\begin{array}{rrr}0.7&1.2\\-0.2&3.3\end{array}\right]\right)=\left[\begin{array}{rrr}0.99&0.02\\0.01&0.98\end{array}\right]$$

Enforcing surjectivity [25]

We address (*iii*) surjective constraints through a *penalty function*:

$$\alpha_{S'}, \alpha_{T'}, \alpha_{C'} \in [0, 1]^{2 \times 2}$$

$$\mathcal{L}_2: \min_{\mathbf{W}} \sum_{\mathbf{W}} \sum_i \left(1 - \max_j t(W)_{ij} \right)$$

$$\begin{bmatrix} 0.99 & 0.02\\ 0.01 & 0.98 \end{bmatrix} \overset{\mathcal{L}_2}{\rightsquigarrow} (1-0.99) + (1-0.98)$$

Solution by gradient descent [25]

We address (i) multiple related problems by jointly solving all the problems via gradient descent:





Synthetic Experiments [25]

We evaluated our learning method:

- On multiple synthetic models;
- Against independent and sequential approach;
- Monitoring loss functions, L1-dist from ground truth, wall-clock time.



Real-World Experiments [25]

We want to model the stage of **coating** in lithium-ion battery manufacturing:

```
Mass Loading = f(input)
```

Experiments are costly, so we want to integrate data¹ collected by two groups running similar (but not identical) experiments:

LRCS (France)

WMG (UK)

Collection of few statistics in each a few stages of battery manufacturing [5].

Collection of detailed space- and time-dependent measurements during coating.

¹https://chemistry-europe.onlinelibrary.wiley.com/doi/full/10.1002/ batt.201900135 https://github.com/mattdravucz/jointly-learning-causal-abstraction/
Real-World Experiments [25]

We evaluated our learning method:

- \bullet Performing abstraction of data from base to abstracted (WMG \rightarrow LRCS);
- Evaluating change in performance using aggregated data when predicting *out-of-sample* (k).

	Training set	Test Set	MSE
(a)	$LRCS[CG \neq k]$	LRCS[CG = k]	1.86 ± 1.75
(b)	$LRCS[CG \neq k]$	LRCS[CG = k]	0.22 ± 0.26
	+ WMG		
(c)	$LRCS[CG \neq k]$	LRCS[CG = k]	1.22 ± 0.95
	$+ \operatorname{WMG}[CG \neq k]$	+ WMG[CG = k]	

Causality and abstraction may both play important role in modelling.

Causality and *abstraction* may both play important role in modelling.

Large space for conceptual and practical development of **causal abstraction frameworks**:

Causality and *abstraction* may both play important role in modelling.

Large space for conceptual and practical development of **causal abstraction frameworks**:

- Foundations of the framemorks
 - Category theory [15]
 - Measure theory [3]
 - Review [23]

Causality and *abstraction* may both play important role in modelling.

Large space for conceptual and practical development of **causal abstraction frameworks**:

- **Foundations** of the framemorks
 - Category theory [15]
 - Measure theory [3]
 - Review [23]
- Observation of these frameworks
 - Measures of abstraction [26]
 - Abstraction with soft interventions [12]
 - Cluster DAGs and do-calculus [1]
 - Causal bandits and abstraction [24]
 - Connection between frameworks [21]

Is Algorithmic and empirical development

- Learning with optimal transport [8]
- Learning linear abstraction [13]
- Target learning [11]
- Neural models [22]
- Riemannian optimization [6]

Algorithmic and empirical development

- Learning with optimal transport [8]
- Learning linear abstraction [13]
- Target learning [11]
- Neural models [22]
- Riemannian optimization [6]
- Applications of causal abstraction
 - Surrogate for agent-based models [7]
 - Explainable AI [10]
 - Visual coarsening [4]

Algorithmic and empirical development

- Learning with optimal transport [8]
- Learning linear abstraction [13]
- Target learning [11]
- Neural models [22]
- Riemannian optimization [6]
- Applications of causal abstraction
 - Surrogate for agent-based models [7]
 - Explainable AI [10]
 - Visual coarsening [4]

And connections to causal representation learning, reinforcement learning...

Algorithmic and empirical development

- Learning with optimal transport [8]
- Learning linear abstraction [13]
- Target learning [11]
- Neural models [22]
- Riemannian optimization [6]
- Applications of causal abstraction
 - Surrogate for agent-based models [7]
 - Explainable AI [10]
 - Visual coarsening [4]

And connections to causal representation learning, reinforcement learning...

More about causal abstraction:

https://github.com/FMZennaro/CausalAbstraction/

Thanks!

Thank you for listening!

References I

- G Anand, Swarnava Ghosh, Liwei Zhang, Angesh Anupam, Colin L Freeman, Christoph Ortner, Markus Eisenbach, and James R Kermode. Exploiting machine learning in multiscale modelling of materials. *Journal of The Institution of Engineers (India): Series D*, pages 1–11, 2022.
- [2] Sander Beckers and Joseph Y Halpern. Abstracting causal models. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, pages 2678–2685, 2019.
- [3] Simon Buchholz, Junhyung Park, and Bernhard Schölkopf. Products, abstractions and inclusions of causal spaces. arXiv preprint arXiv:2406.00388, 2024.
- [4] Krzysztof Chalupka, Pietro Perona, and Frederick Eberhardt. Visual causal feature learning. *arXiv preprint arXiv:1412.2309*, 2014.

References II

- [5] Ricardo Pinto Cunha, Teo Lombardo, Emiliano N Primo, and Alejandro A Franco. Artificial intelligence investigation of nmc cathode manufacturing parameters interdependencies. *Batteries & Supercaps*, 3(1):60–67, 2020.
- [6] Gabriele D'Acunto, Fabio Massimo Zennaro, Yorgos Felekis, and Paolo Di Lorenzo. Causal abstraction learning based on the semantic embedding principle. arXiv preprint arXiv:2502.00407, 2025.
- [7] Joel Dyer, Nicholas Bishop, Yorgos Felekis, Fabio Massimo Zennaro, Anisoara Calinescu, Theodoros Damoulas, and Michael Wooldridge. Interventionally consistent surrogates for agent-based simulators. arXiv preprint arXiv:2312.11158, 2023.
- [8] Yorgos Felekis, Fabio Massimo Zennaro, Nicola Branchini, and Theodoros Damoulas. Causal optimal transport of abstractions. arXiv preprint arXiv:2312.08107, 2023.

References III

- [9] Luciano Floridi. The method of levels of abstraction. *Minds and machines*, 18:303–329, 2008.
- [10] Atticus Geiger, Hanson Lu, Thomas Icard, and Christopher Potts. Causal abstractions of neural networks. Advances in Neural Information Processing Systems, 34:9574–9586, 2021.
- [11] Armin Kekić, Bernhard Schölkopf, and Michel Besserve. Targeted reduction of causal models. *arXiv preprint arXiv:2311.18639*, 2023.
- [12] Riccardo Massidda, Atticus Geiger, Thomas Icard, and Davide Bacciu. Causal abstraction with soft interventions. *arXiv preprint arXiv:2211.12270*, 2022.
- [13] Riccardo Massidda, Sara Magliacane, and Davide Bacciu. Learning causal abstractions of linear structural causal models. *arXiv preprint arXiv:2406.00394*, 2024.

References IV

- [14] Jun Otsuka and Hayato Saigo. On the equivalence of causal models: A category-theoretic approach. arXiv preprint arXiv:2201.06981, 2022.
- [15] Jun Otsuka and Hayato Saigo. The process theory of causality: an overview. 2022.
- [16] Judea Pearl. Causality. Cambridge University Press, 2009.
- [17] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Elements of causal inference: Foundations and learning algorithms. MIT Press, 2017.
- [18] Eigil F Rischel and Sebastian Weichwald. Compositional abstraction error and a category of causal models. arXiv preprint arXiv:2103.15758, 2021.
- [19] Eigil Fjeldgren Rischel. The category theory of causal models. 2020.

References V

- [20] Paul K Rubenstein, Sebastian Weichwald, Stephan Bongers, Joris M Mooij, Dominik Janzing, Moritz Grosse-Wentrup, and Bernhard Schölkopf. Causal consistency of structural equation models. In 33rd Conference on Uncertainty in Artificial Intelligence (UAI 2017), pages 808–817. Curran Associates, Inc., 2017.
- [21] Willem Schooltink and Fabio Massimo Zennaro. Aligning graphical and functional causal abstractions. *arXiv preprint arXiv:2412.17080*, 2024.
- [22] Kevin Xia and Elias Bareinboim. Neural causal abstractions. *arXiv* preprint arXiv:2401.02602, 2024.
- [23] Fabio Massimo Zennaro. Abstraction between structural causal models: A review of definitions and properties. In UAI 2022 Workshop on Causal Representation Learning, 2022.

References VI

- [24] Fabio Massimo Zennaro, Nicholas George Bishop, Joel Dyer, Yorgos Felekis, Ani Calinescu, Michael J Wooldridge, and Theodoros Damoulas. Causally abstracted multi-armed bandits. In *The 40th Conference on Uncertainty in Artificial Intelligence*, 2024.
- [25] Fabio Massimo Zennaro, Máté Drávucz, Geanina Apachitei, W. Dhammika Widanage, and Theodoros Damoulas. Jointly learning consistent causal abstractions over multiple interventional distributions. In 2nd Conference on Causal Learning and Reasoning, 2023.
- [26] Fabio Massimo Zennaro, Paolo Turrini, and Theo Damoulas. Quantifying consistency and information loss for causal abstraction learning. In Proceedings of the Thrity-Second International Conference on International Joint Conferences on Artificial Intelligence, 2023.