



Data-Driven and Knowledge-Based Causal Network Discovery for Identifying Differential Equations



Mitsuhiro Odaka^{1,2,3,4,*} Morgan Magnin^{3,2} Katsumi Inoue^{2,1,3}

¹ The Graduate University for Advanced Studies, SOKENDAI, Japan

² National Institute of Informatics, Japan

³ École Centrale de Nantes, France

⁴ Japan Society for the Promotion of Science (JSPS) Research Fellowships for Young Scientists



*odaka@nii.ac.jp

Table of contents

- ❑ Our motivation in scientific knowledge discovery: Uncovering dynamics
- ❑ Data-Driven & Knowledge-Based (DD-KB) integrated approach
- ❑ Gene network inference for constructing COVID-19 pathways
- ❑ Adversarial learning of causal networks + Equation discovery (ongoing)
- ❑ Related work
- ❑ Concluding remarks

Table of contents

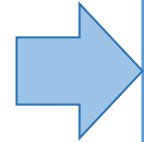
- ❑ **Our motivation in scientific knowledge discovery: Uncovering dynamics**
- ❑ **Data-Driven & Knowledge-Based (DD-KB) integrated approach**
- ❑ Gene network inference for constructing COVID-19 pathways
- ❑ Adversarial learning of causal networks + Equation discovery (ongoing)
- ❑ Related work
- ❑ Concluding remarks

Motivation: Uncovering dynamics

= Causal network discovery + Identifying governing equations

State variables

$$X = (x_1, x_2, \dots, x_n)$$



Uncovering dynamics

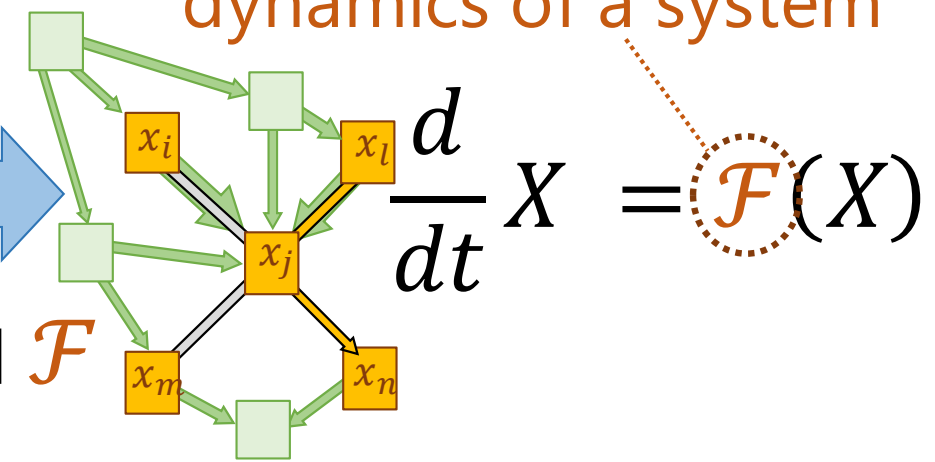


Equations governing the dynamics of a system

Identifying networks and \mathcal{F}



More profound understanding of the systems or the world



- Advocacy of the system control strategy
- Extrapolation by numerical simulation
- Comparison of different interventions

Data-Driven and Knowledge-Based (DD-KB) Integrated Approach

Learning and inference of causal networks and differential equations governing the dynamics of a system from **observed data & background knowledge**

Problem Settings

Inputs

- ☐ Observations \mathcal{O}

Continuous multivariate time series data $X(t) \in \mathbb{R}^d$

X : Observed variables t : Time index (arbitrary unit) $\in [0, \mathbb{R}^+]$

- ☐ Background knowledge \mathcal{B}

Outputs

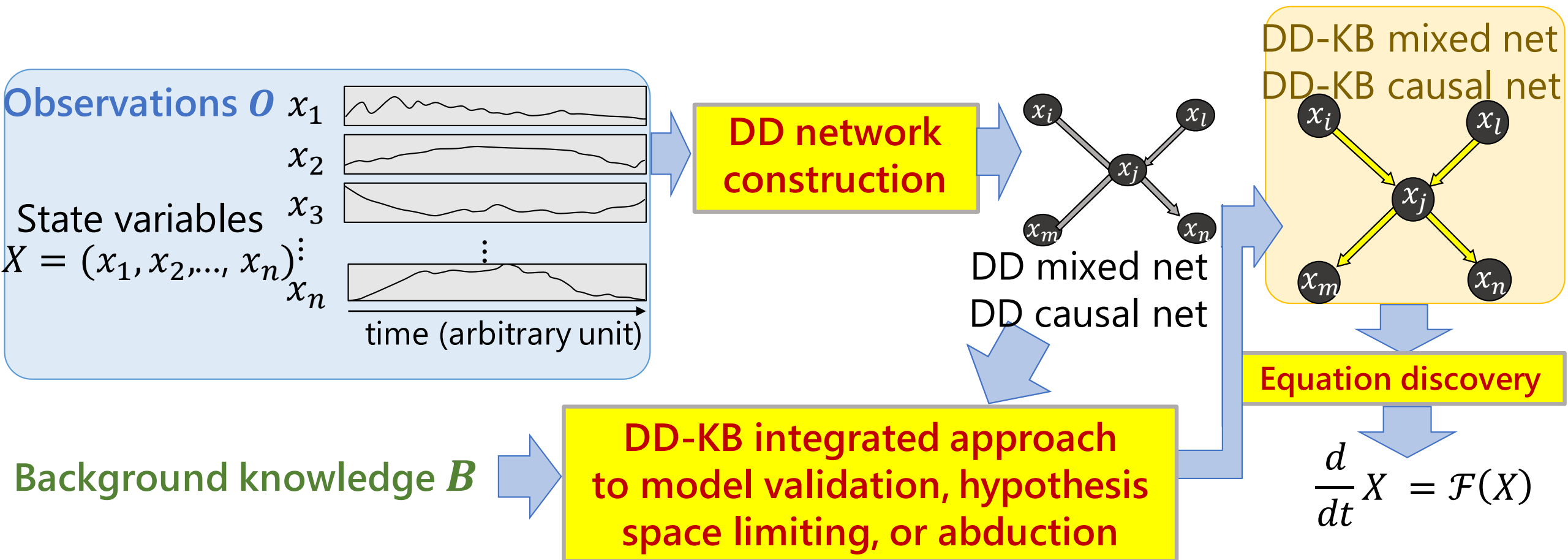
- ☐ Continuous deterministic dynamical system $\langle X, \dot{X}, \mathcal{F}, t \rangle$

\mathcal{F} : Dynamic constraints that define equations governing a dynamics of the system

- ☐ DD-KB network \mathcal{M}

(Directed/undirected mixed network, Causal network, etc.) 5

Data-Driven and Knowledge-Based (DD-KB) Integrated Approach



Aim 1: To apply DD-KB integrated approach to gene network inference for building new COVID-19 pathways

Aim 2: To develop method for learning causal network in continuous domain

Table of contents

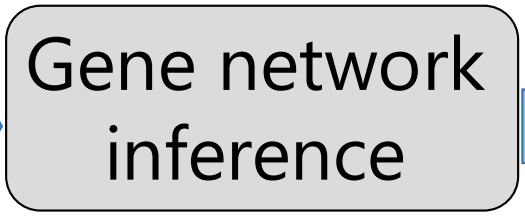
- ❑ Our motivation in scientific knowledge discovery: Uncovering dynamics
- ❑ Data-Driven & Knowledge-Based (DD-KB) integrated approach
- ❑ **Gene network inference for constructing COVID-19 pathways**
- ❑ Adversarial learning of causal networks + Equation discovery (ongoing)
- ❑ Related work
- ❑ Concluding remarks

Data-driven & knowledge-based (DD-KB) integrated approach to gene network inference for constructing COVID-19 pathways

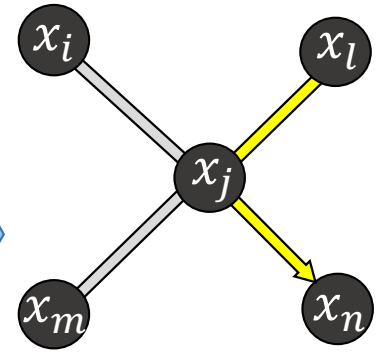
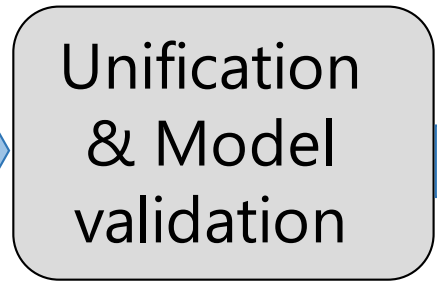
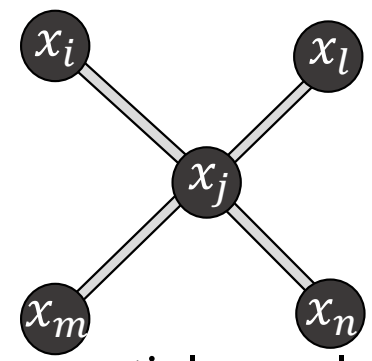
[Odaka et al. Preprint at *Research Square*. 2022]

[Odaka et al. *CAMDA*. 2022]

Observations O

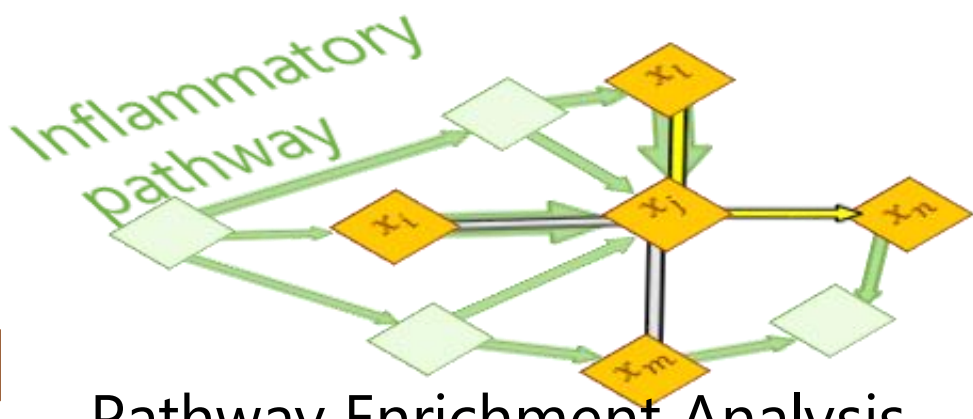
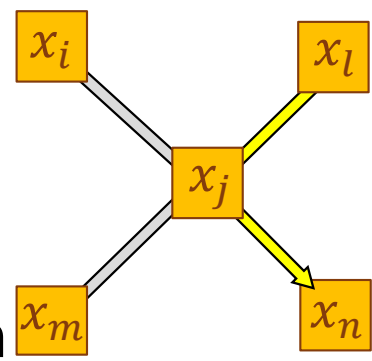
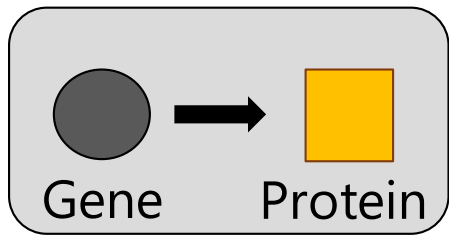
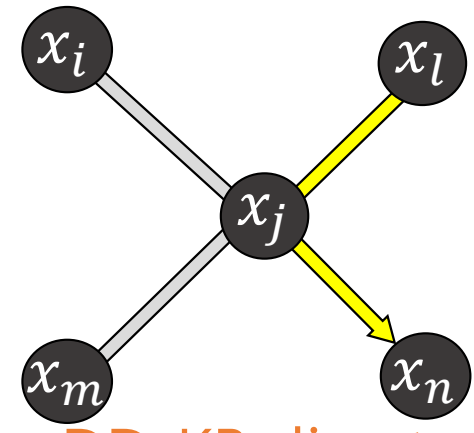


Gene expression profiles of bronchial lavage fluid samples from COVID-19 severe patients (n=10) & negative control (n=2) [Grant et al. *Nature*. 2021] (cells 15,481 × genes 21,819)



Background knowledge B

Pathway Commons ver 12 (Reactome, Panther, HumanCyc, BIND, MSigDB), BioGRID ver 4.4, STRING ver 11.5



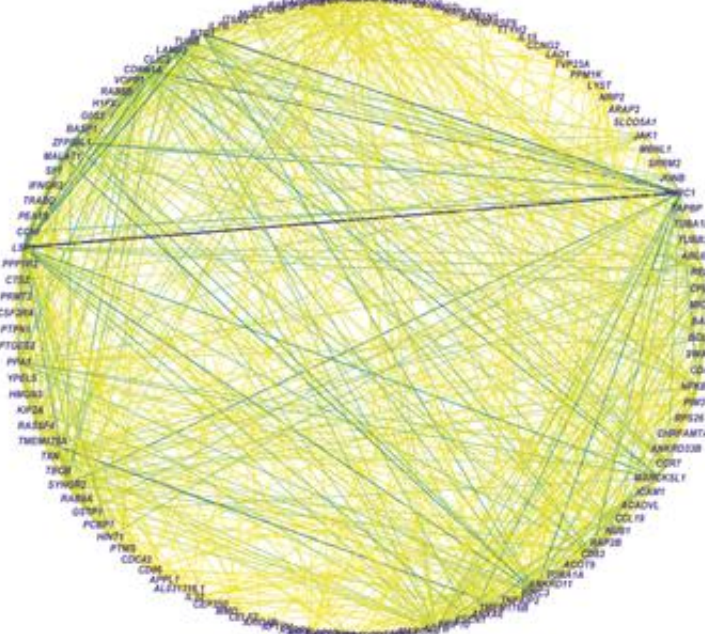
DD-KB directed / undirected mixed net

Mapping onto function-known KEGG pathways e.g., Activation of genes of inflammatory pathway w/ significant difference

DD-KB integrated approach to gene network inference for constructing COVID-19 pathways

Spurious correlation removal

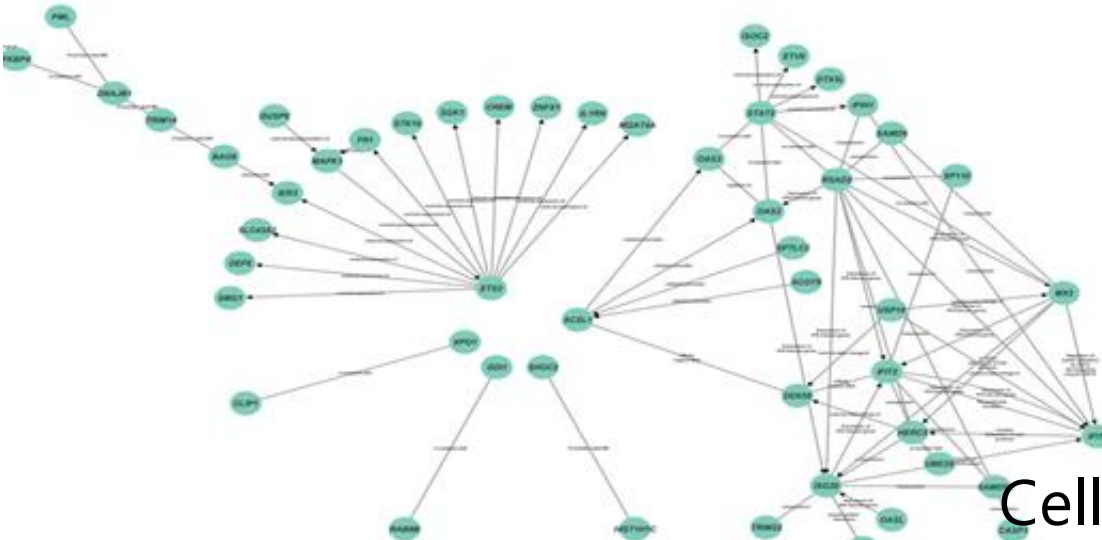
	Day	Full	Removed	Final
<i>ACTB</i>	1	5,671 (100%)	4,757 (84%)	914 (16%)
	5	6,328 (100%)	5,301 (84%)	1,027 (16%)
	10	5,050 (100%)	4,257 (84%)	793 (16%)
<i>ICAM1</i>	1	7,503 (100%)	6,309 (84%)	1,194 (16%)
	5	20,706 (100%)	18,914 (91%)	1,792 (9%)
	10	8,001 (100%)	6,748 (84%)	1,253 (16%)
<i>C15orf48</i>	1	9,453 (100%)	7,995 (85%)	1,458 (15%)
	5	13,530 (100%)	12,049 (89%)	1,481 (11%)
	10	8,001 (100%)	6,748 (84%)	1,253 (16%)



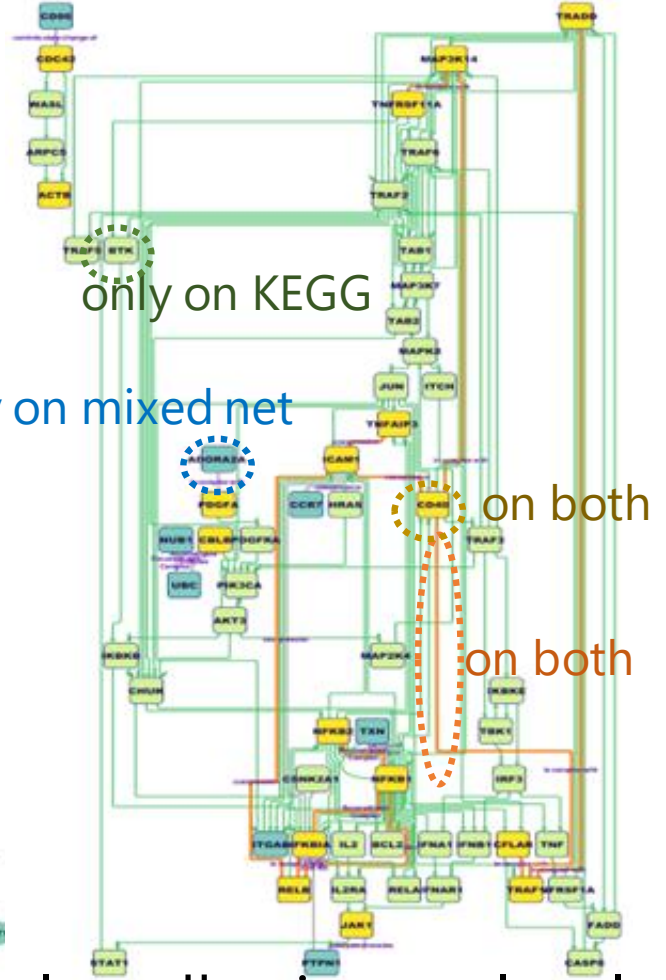
2nd-order partial correlation network

Model validation by knowledge

Source	Target	Weight	Relation	Knowledge bases
REPIN1	SNN	0.366753074	controls-expression-of	MSigDB
CCL22	CCR6	0.189063535	in-complex-with	PathwayCommons_Get_API
DUSP5	PTPN1	0.428685	neighbor-of	Panther
JAK1	PTPN1	0.564563581	catalysis-precedes	Panther
NFKBIA	PTPN1	0.580935037	Proximity Label-MS	BioGRID
LAMP3	RAB9A	0.487968997	Proximity Label-MS	BioGRID



DD-KB directed / undirected mixed net



Cellular adhesion molecule *ICAM1* pathways inferred from COVID-19 omics data

- (1) Pathways missing from the current C19DMap
- (2) Molecules supporting SARS-CoV-2 cell-to-cell transimission hypothesis [Odeka & Inoue *Heliyon* 2021]

Table of contents

- ❑ Our motivation in scientific knowledge discovery: Uncovering dynamics
- ❑ Data-Driven & Knowledge-Based (DD-KB) integrated approach
- ❑ Gene network inference for constructing COVID-19 pathways
- ❑ **Adversarial learning of causal networks + Equation discovery (ongoing)**
- ❑ Related work
- ❑ Concluding remarks

Learning differential equations via causal network discovery from multivariate time series

Observations O

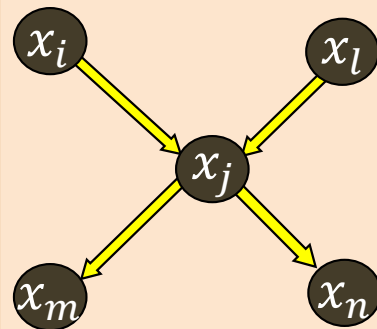
Causal network discovery

from multivariate time series by adversarial learning [Odaka et al. *JSAI*. 2023]

Input as candidate skeleton

Background knowledge B

- Networks purely from knowledge bases
- DD-KB directed / undirected mixed net



DD-KB causal net (DAG)

Imputing as inductive bias to limit hypothesis space of equation discovery

Equation discovery

e.g., Sparse Identification of Nonlinear Dynamics (SINDy) [Brunton et al. *PNAS*. 2016]

$$\frac{d}{dt} X = \mathcal{F}(X)$$

ODEs, Causal network

Model verification

- Sensitivity analysis, Stability analysis
- Reachability analysis, model checking
- Assessment w/ model selection criteria (AIC, BIC, MDL, etc.)

Verified ODEs, Causal network

General pros/cons in deep learning and analytical technique

	Deep learning	Analytical technique
Robustness to noise	●	▲
Scalability	●	▲
Interpretability	▲	●

To cover the pros and cons of each other, we combine **deep learning and analytical technique** in equation discovery.

Problem settings

Given: Observed data

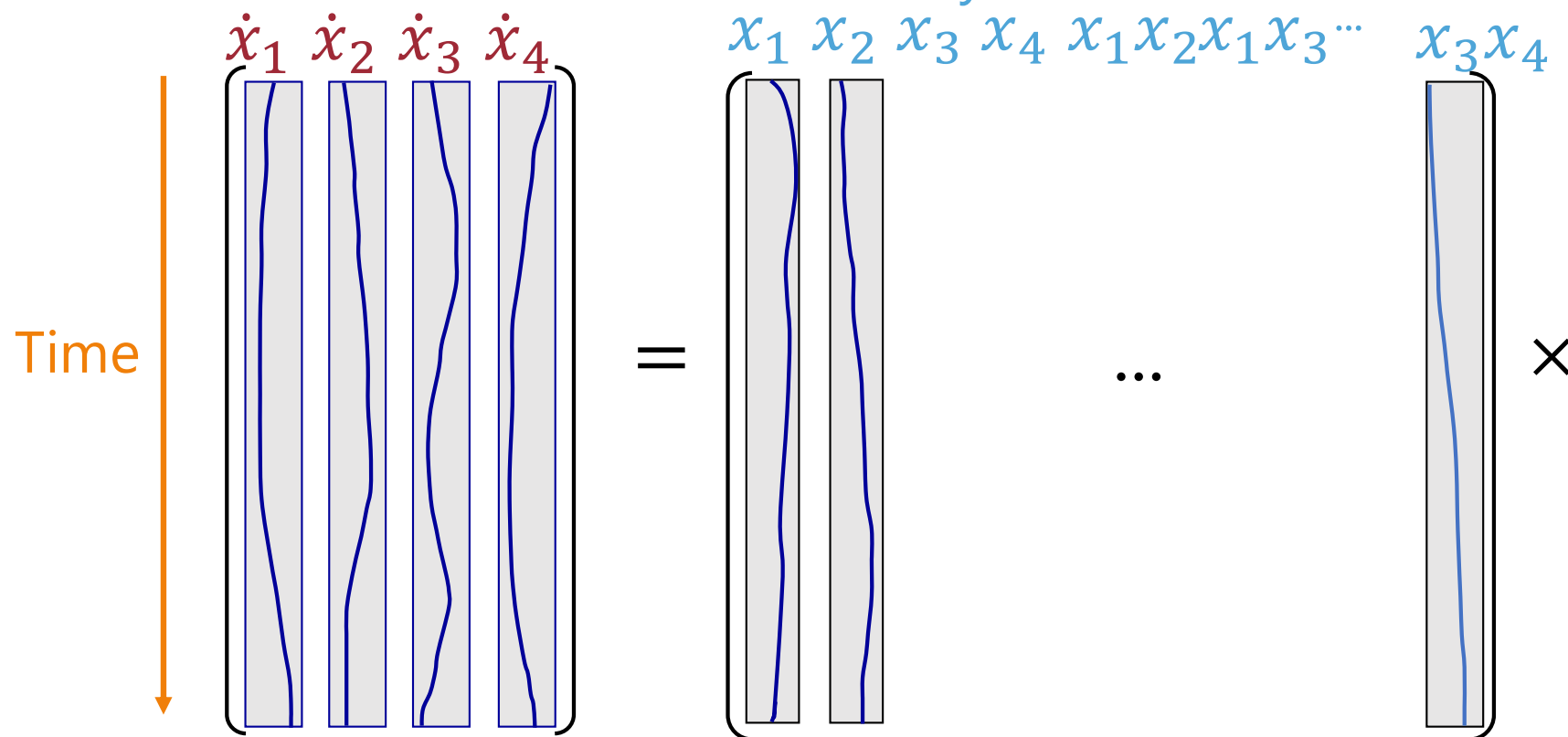
$$\frac{d}{dt} X = \Theta(X) E$$

Selecting variables amenable for the dynamics representation

Time series of endogenous vars' derivatives

Time series of endogenous vars & synthetic vars

Find:
Coefficient vector



	x_1	x_2	x_3	x_4
x_1	0	0.03	0.2	0.02
x_2	0.05	0	0	0.1
x_3	0	0	0	0.03
x_4	0	0	0.1	0
x_1x_2	0	0	0	0
x_1x_3	0	0	0	0
x_1x_4	0.01	0	0	0
x_2x_3	0	0	0	0
x_2x_4	0	0.02	0	0
x_3x_4	0	0	0	0

Existing equation discovery techniques have not often explicitly discovered causality. → Overfitting by dull terms

Problem settings

Given: Observed data

$$\frac{d}{dt} X = \Theta(X) E$$

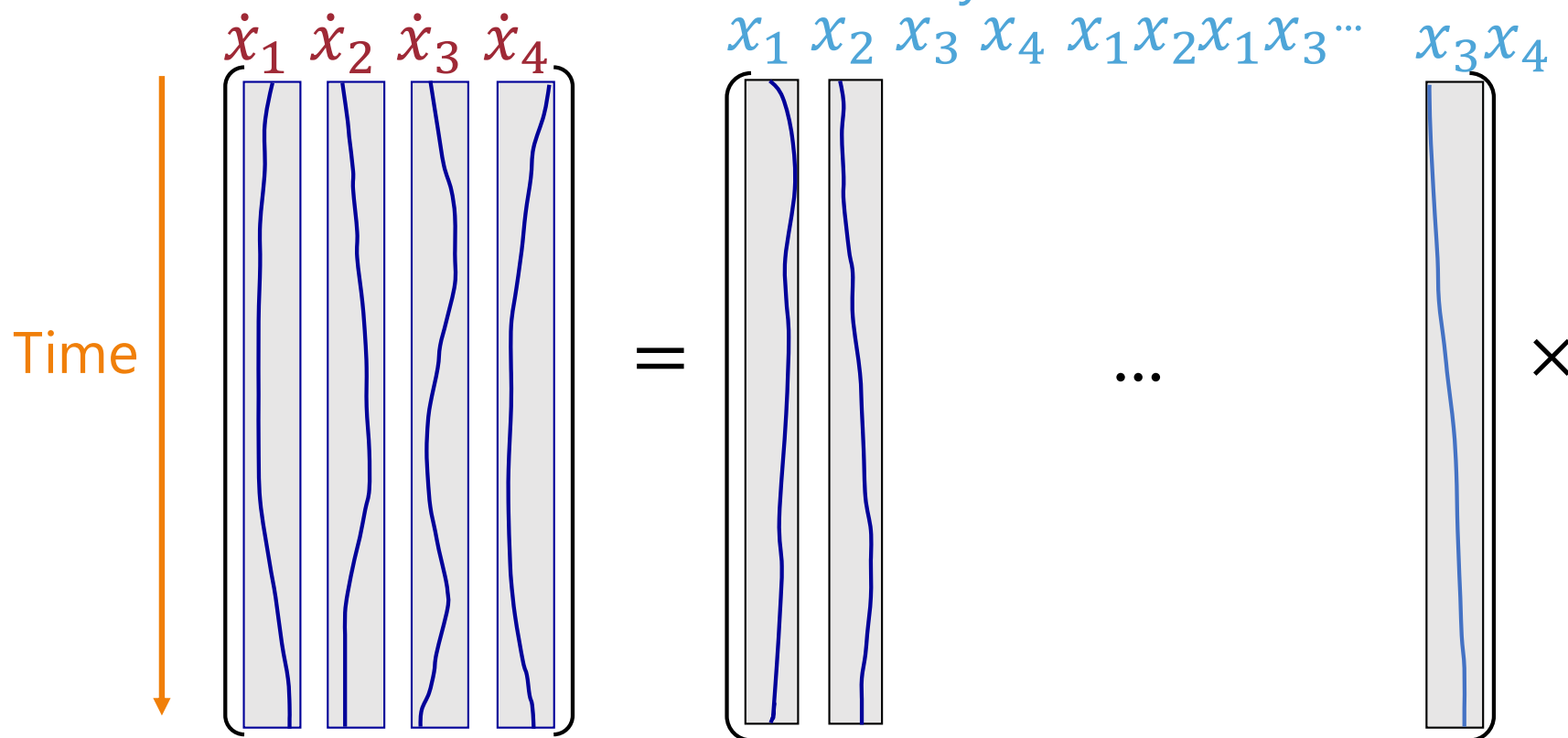
Selecting variables amenable for the dynamics representation

Time series of endogenous vars' derivatives

Time series of endogenous vars & synthetic vars

Find:

Coefficient vector



Find: Coefficient vector

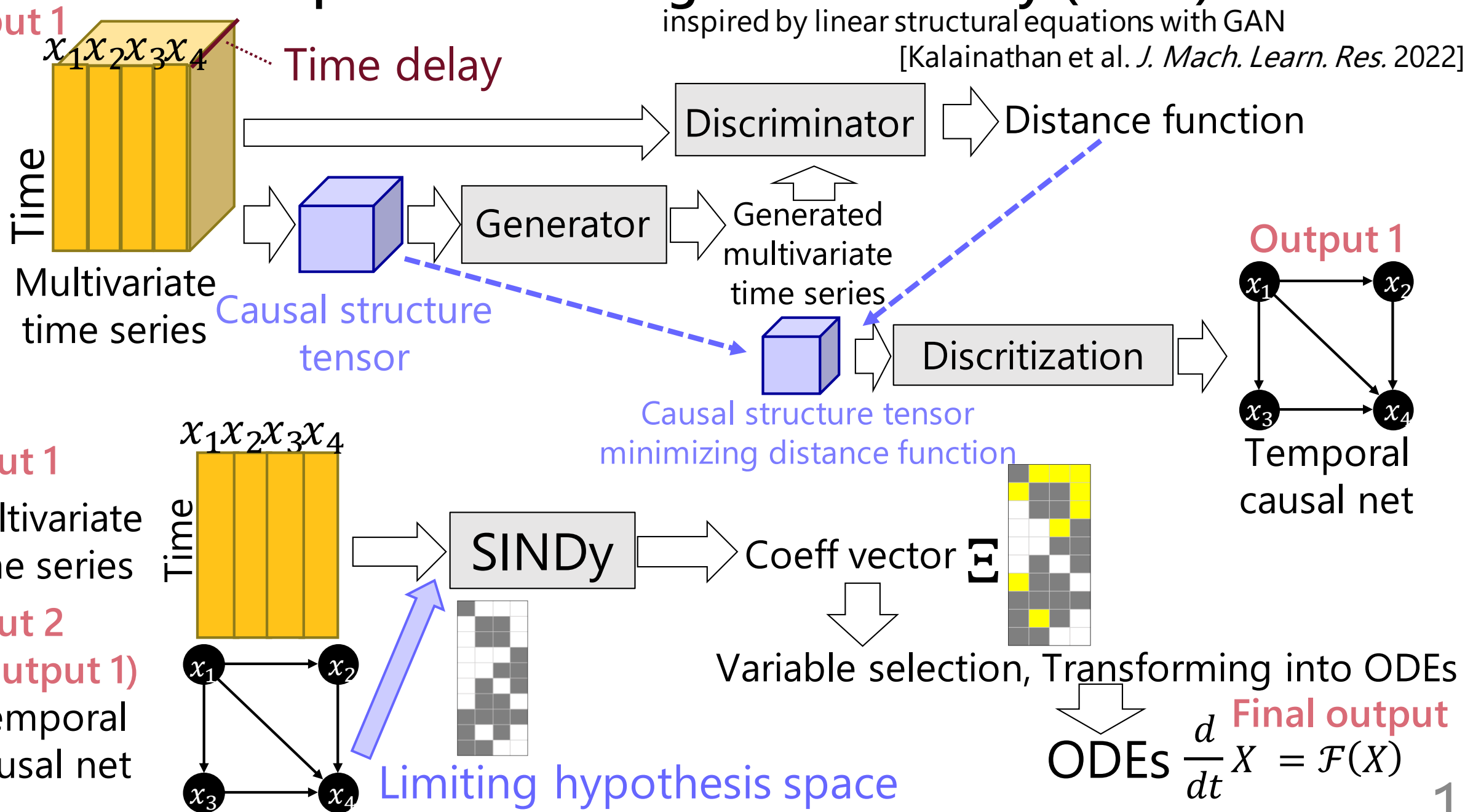
	x_1	x_2	x_3	x_4
x_1	0	0.03	0.2	0.02
x_2	0.05	0	0	0.1
x_3	0	0	0	0.03
x_4	0	0	0.1	0
x_1x_2	0	0	0	0
x_1x_3	0	0	0	0
x_1x_4	0.01	0	0	0
x_2x_3	0	0	0	0
x_2x_4	0	0.02	0	0
x_3x_4	0	0	0	0

Our method detects data-driven causality to filter out non-causal elements from parameter estimation.

Parsimonious Equation Learning with Causality (PELC)

inspired by linear structural equations with GAN

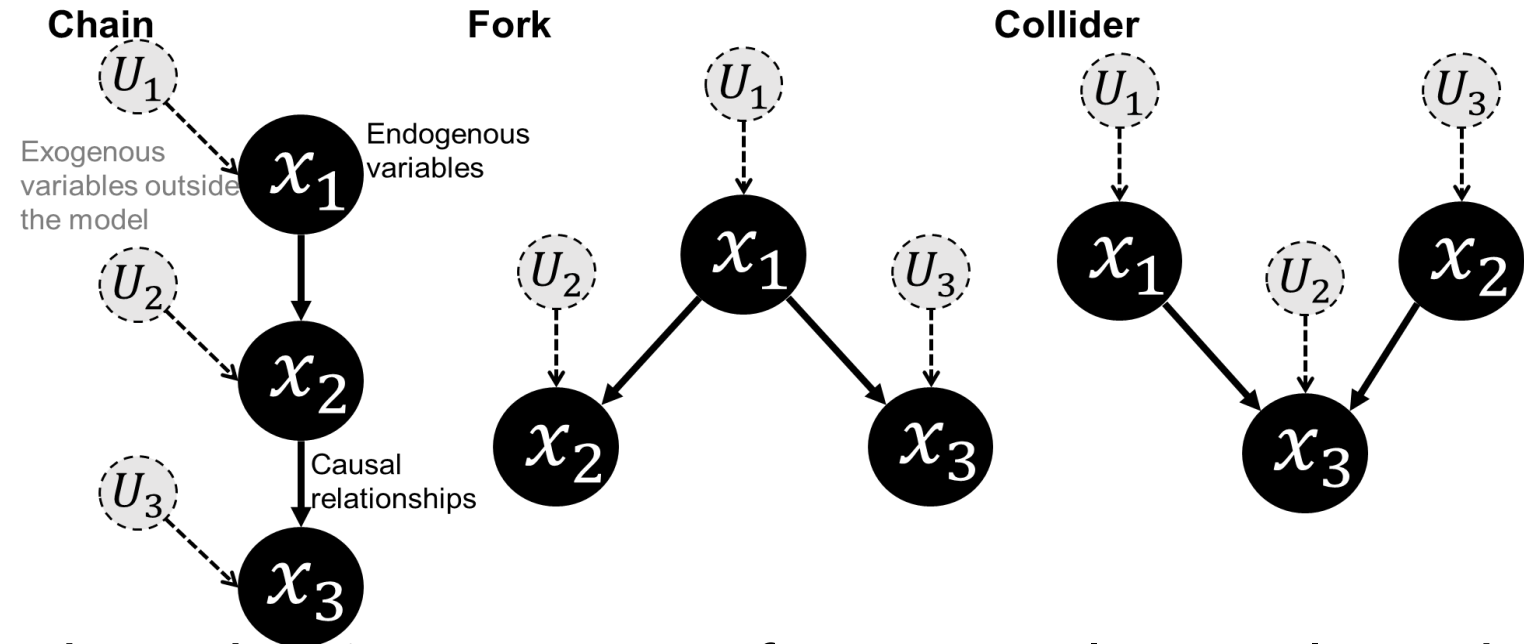
[Kalainathan et al. *J. Mach. Learn. Res.* 2022]



Experimental settings

Comparative techniques \

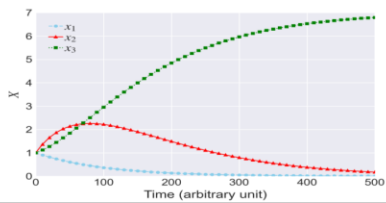
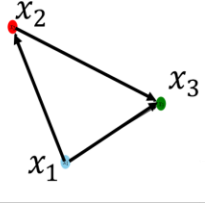
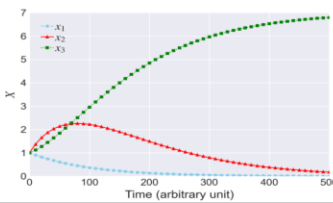
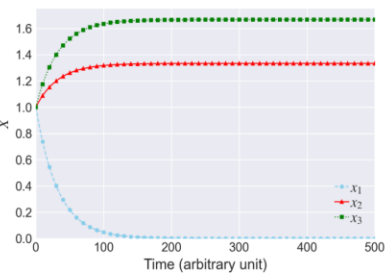
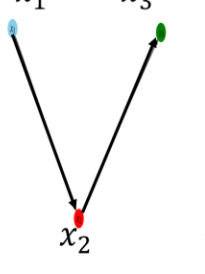
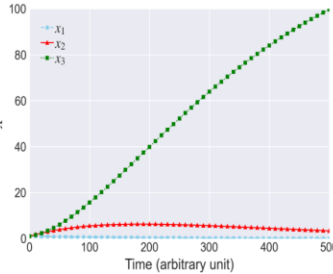
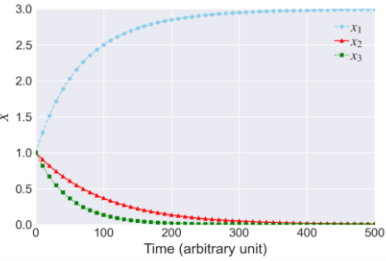
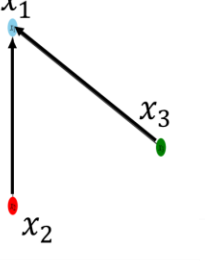
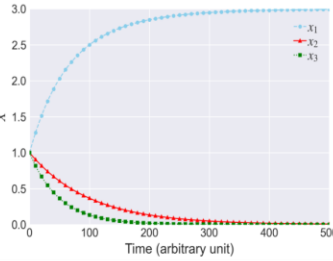
- SINDy (Sparse regression)
- VAR-LiNGAM (Linear structural equation with time delay)
(Vector Auto-Regression Linear Non-Gaussian Acyclic Model)
[Hyvärinen+*J. Mach. Learn. Res.*'10] Bootstrap sampling size: 100
- PELC (Proposed method)



\ are applied to time-series data generated by three ODEs.

Three basic patterns of structural causal models

Preliminary results (ongoing) and discussion

	Ground Truth		SINDy	VAR-LiNGAM	PELC		
	ODEs	Original dynamics	ODEs	ODEs	Causal network	ODEs	Reproduced dynamics
Chain	$\begin{aligned}\dot{x}_1 &= -0.01 x_1 \\ \dot{x}_2 &= 0.05 x_1 - 0.01 x_2 \\ \dot{x}_3 &= 0.01 x_2\end{aligned}$		$\begin{aligned}\dot{x}_1 &= -0.003x_2 \\ \dot{x}_2 &= 0.014x_1 - 0.003x_2 + 0.004x_3 \\ \dot{x}_3 &= 0.003x_2\end{aligned}$ RT: 0.00899	$\begin{aligned}\dot{x}_1 &= -0.01 x_1 \\ \dot{x}_2 &= 0.05 x_1 - 0.01 x_2 \\ \dot{x}_3 &= 0.01 x_2\end{aligned}$ RT: 4.279	 $\begin{aligned}\dot{x}_1 &= -0.01 x_1 \\ \dot{x}_2 &= 0.05 x_1 - 0.01 x_2 \\ \dot{x}_3 &= 0.01 x_2\end{aligned}$ RT: 160.392		
Fork	$\begin{aligned}\dot{x}_1 &= -0.03 x_1 \\ \dot{x}_2 &= 0.01 x_1 \\ \dot{x}_3 &= 0.02 x_1\end{aligned}$		$\begin{aligned}\dot{x}_1 &= -0.008x_1 \\ \dot{x}_2 &= 0.000 \\ \dot{x}_3 &= 0.005x_1\end{aligned}$ RT: 0.0107	$\begin{aligned}\dot{x}_1 &= -0.05 x_1 + 0.12 x_2 - 0.09 x_3 \\ \dot{x}_2 &= 0.13 x_1 - 0.59x_2 + 0.48x_3 \\ \dot{x}_3 &= -0.08 x_1 + 0.48x_2 - 0.38x_3\end{aligned}$ RT: 4.365	 $\begin{aligned}\dot{x}_1 &= -0.005x_1 \\ \dot{x}_2 &= 0.08x_1 - 0.005x_2 \\ \dot{x}_3 &= 0.04 x_2\end{aligned}$ RT: 159.674		
Collider	$\begin{aligned}\dot{x}_1 &= 0.01 x_2 + 0.02 x_3 \\ \dot{x}_2 &= -0.01 x_2 \\ \dot{x}_3 &= -0.02 x_3\end{aligned}$		$\begin{aligned}\dot{x}_1 &= -0.024x_2x_3 \\ \dot{x}_2 &= 0.000 \\ \dot{x}_3 &= 0.000\end{aligned}$ RT: 0.0137	$\begin{aligned}\dot{x}_1 &= 0.01 x_2 + 0.02 x_3 \\ \dot{x}_2 &= -0.01 x_2 \\ \dot{x}_3 &= -0.02 x_3\end{aligned}$ RT: 4.537	 $\begin{aligned}\dot{x}_1 &= 0.01 x_2 + 0.02 x_3 \\ \dot{x}_2 &= -0.01 x_2 \\ \dot{x}_3 &= -0.02 x_3\end{aligned}$ RT: 160.071		

RT: Run Time (s)

- ❑ SINDy did not reproduce the original dynamics in a pilot study (perhaps due to our failure of parameter setting ...)
- ❑ VAR-LiNGAM and PELC reproduced ODEs of chain and collider SCMs.
- ❑ PELC reproduced dynamics with a fewer variables than VAR-LiNGAM.

Summary of Parsimonious Equation Learning with Causality (PELC)

- ❑ PELC learns *causal structure tensor*, which represents linear structural equations with time delay.
- ❑ *Causal network topology* is incorporated into hypothesis space of equation discovery *as inductive bias*.
- ❑ Limiting hypothesis space by replacing non-causal elements of coefficient vectors with zero

Open questions

- ❖ Is it plausible to associate **causality** with **differential equations**?
- ❖ If so, what is the best definition of causality?
(Granger causality, transfer entropy, conditional probability, etc.)

Table of contents

- ❑ Our motivation in scientific knowledge discovery: Uncovering dynamics
- ❑ Data-Driven & Knowledge-Based (DD-KB) integrated approach
- ❑ Gene network inference for constructing COVID-19 pathways
- ❑ Adversarial learning of causal networks + Equation discovery (ongoing)
- ❑ **Related work**
- ❑ **Concluding remarks**

Related work on causal discovery

- ❑ Statistical approach (based on Structural Equation Modeling):
LiNGAM (Linear Non-Gaussian Acyclic Model) [Shimizu+J. Mach. Learn. Res.'06]
- ❑ Bayesian approach: Structure learning of Bayesian networks
- ❑ ML approach: DAG-GNN [Yu et al. *ICML* 2019]
Causal discovery w/ reinforcement learning [Zhu et al. *ICLR* 2020]

Related work on equation discovery

- ❑ Constraining hypothesis space with context-free grammar [Todorovski & Džeroski. *ICML*. 1997]
- ❑ Constraint system identification with human-computer communication [Stolle & Bradley. 2007]
- ❑ Process model induction from data + knowledge on observed behavior [Bridewell et al. *Mach.Learn.*2008]
- ❑ Genetic learning of free-form natural laws [Schmidt & Lipson. *Science*. 2009]
- ❑ Physics-informed learning from small data [Raissi et al. *J. Comput. Phys.* 2019]
- ❑ Learning symmetry and separability in data (AI-Feynman) [Udrescu & Tegmark. *Sci. Adv.* 2020]

Their relationships

- ❑ Transforming ODEs into deterministic SCMs via equilibrium equations under intervention
[Mooij et al. *UAI* 2013]
- ❑ Constructing SCMs by reproducing asymptotic behavior of ODEs under intervention
[Rubenstein et al. *UAI* 2018]
- ❑ 1-to-1 correspondence between ODEs and SCMs has not been adequately developed.

Concluding remarks

□ Summary of contributions

- Data-driven & knowledge-based (DD-KB) integrated approach to uncovering dynamics
- Case study: Gene network inference for constructing COVID-19 pathways
- PELC: Causal network discovery w/ GAN + Mapping discovered causal structure onto hypothesis space in equation discovery

□ Limitation

- Parameter sensitivity
- Generalization

□ Future work

- Sensitivity analysis
- Learning from partial or small data
- Neuro-symbolic AI:
Realizing causal & equation discovery in the same hypothesis space
- High-dimensional synthetic data experiments, Real-world dataset applications