# **Causal Representation Meets Stochastic Modeling** under Generic Geometry

# Anonymous Author(s)

Affiliation Address email

# **Abstract**

We investigate the identifiability problem of latent stochastic processes characterized by high dynamics that occur in continuous time with varying intensities (e.g., a multivariate Hawkes process), and we provide the corresponding identifiability theory. Building on this theoretical foundation, we implement MUTATE, a variational autoencoder framework with a time-adaptive transition module to evaluate stochastic dynamics on both synthetic stochastic processes and real-world biological signal data. This work advances causal representation learning theory by extending it to continuous-time and stochastic settings via weak topology and algebraic signature, highlighting the importance of this approach in addressing scientific questions, such as the accumulation of mutations in genomics and the mechanisms driving neuron spike triggers in response to time-varying dynamics.

# Introduction

1

2

3

4

8

9

10

11

12

17

21

22

23

24

26

27

28

29

30

31

33

34

Inferring causal relationships among variables from observations capitalizes the potential of machine 13 learning to advance scientific discovery, as it reveals underlying mechanisms that are not identifiable 14 from observational distributions alone [1]. However, because of limited data sources and the challenge 15 of interpreting high-dimensional perceptual data, causal variables with their graphical structure are 16 often unknown and thus can be learned in a non-interpretable manner, which causes the difficulty of 18 identifiability [2, 3]. Recently, a growing number of studies on the disentanglement of latent causal representation advance the identifiability guarantee and propose methods for latent causal variable 19 estimation. Seminal works among them establish identifiability by leveraging sufficient variability in 20 latent distribution due to multiple-source data [4, 5], auxiliary variable [2, 6, 7, 8], or intervention to a latent causal graph [9, 10, 11, 12, 13].

Most concurrent work aforementioned aims at recovering the latent causal variables in the time-series; A common condition they heavily hinge on is step-wise conditional independence [4, 5], termed as the number of time lags in their claims, among variables. These works mainly address cases when the mixing is assumed invertible, such that latent variables can be recovered up to component-wise indeterminacy. However, the latent dynamics driven by a stochastic process or system of stochastic differential equations are less explored. For example, in biology, fatal diseases such as cancer are principally caused by cumulative multiple mutations found in driver genes as the colonial expansion proceeds. In neuroscience, the latent event dynamics trigger visible biological signals [14, 15]. Finding cancer-associated mutational genes and tracking their behavior through their representation has been given much more paramount importance in recent few decades [16, 17, 18]. Therefore, a formal theoretical guarantee for its identifiability is missing for stochastic causal dynamics and their intervention effects.

This paper aims to establish the identifiability of latent spaces governed by stochastic dynamics driven 35 by the intensity  $\lambda_t$ . Our main results demonstrate that such stochastic processes are compatible with

the current causal representation learning framework and converge to an equivalent infinite-order 37

- INAR process, provided that an integrated fluctuation term is added. Building on this foundation, the 38
- identifiability of dynamic processes is ensured by precisely controlling the geometry of the latent 39
- space via the subtle algebraic structure of cumulants. 40

#### Causal disentanglement in stochastic process 41

#### 2.1 Generative model for stochastic process 42

Let  $O_t \in \mathbb{R}^n$  be observable data, and  $Z_t \in \mathbb{R}^p$  be a latent process.  $O_t$  is being generated from 43

- latent point processes  $Z_t$  through an unknown mixing function f.  $A^{\otimes d}$  denotes the tensor/Kronecker 44
- product. In a time process,  $\Phi$  denotes the transition operator (e.g., an autoregressive coefficient matrix 45
- or continuous kernel matrix) and symbol \* denotes the convolution operator. We formally set up our 46
- problem to learn dynamics of  $Z_t$  from  $O_t$  in Def. 2.1. 47
- **Definition 2.1.** Consider a series of latent processes that are stochastic and self-exciting as time 48
- increases, that is,  $Z_t := N_t$ , the cumulative counting process at time t. Suppose we only have access
- to  $O_t$  without knowledge about  $Z_t$ . The conditional intensity  $\lambda_t$  of  $N_t$  [19] and generative model of
- $O_t$  is written as: 51

61

$$O_t = f(N_t(\Delta)), \ \lambda_t = h(u + \Phi_t \star dN_t) \tag{1}$$

- Then, the objective is to recover f,  $\lambda_t$  as well as its causal structure  $\Phi$ 52
- $h(\cdot)$  is an unknown function reflecting how  $\Phi_t$  and  $dN_t$  are mixed under a regular convolution. A 53
- simple choice  $h(\cdot) = x$  immediately reduces a convolution to linear point processes.  $\phi$  is one element 54
- of  $\Phi_t$ , an integrated kernel matrix and  $dN_t^i$  the measure of a counting process  $N_t^j$ , as well as the integral measure in Itô calculus. The counting process  $N_t^i$  and the conditional intensity  $\lambda_t^i$  satisfy:  $N_{t+\Delta t}^i N_t^i = N_{\Delta t}^i = dN_t^i$  and  $\lambda_t^i = \frac{\mathbb{E}[dN_t^i|\mathcal{F}_t]}{dt}$ . For a point process to be well-defined, non-trivial constraints are needed, one of which is the stationary condition, an assumption widely adopted in 55
- 56
- 57
- 58
- most stochastic process literature to ensure a unique process. We summarize the necessary conditions 59
- to define an inhomogeneous point process in A.4. 60

### 2.2 A graph isomorphism to causal kernels

- The seminal result in [20], establish the weak convergence of continuous stochastic processes under
- the corresponding weak topology. In particular, for any compact interval [a, b], a subsequence INAR 63
- process converges to the point process  $N_t$ . This convergence is crucial for constructing a causal graph 64
- structure compatible with the term causal representation in our paper. 65
- **Lemma 1** (Bounding Point Process in Variational Approximation). Let  $N_t \in \mathbb{R}^p$  be a multivariate 66
- point process whose conditional intensity function is governed by a convolution structure described 67
- in Eq. (1). Suppose the noise term  $\epsilon_t$  is mean-zero and mutually independent, then the intensity model 68
- admits the following weak convergence:

$$Z^{(\Delta)} := \text{INAR}(p) \xrightarrow{w} N, \quad p \to \infty$$
 (2)

- We then show that  $\Phi$  in INAR(p) admits an augmented DAG structure  $\mathscr{G}_K$  in Lem. 2. 70
- **Lemma 2.** Given a bipartite graph of the proposed point process with  $\Phi$ , it admits a kernel DAG, denoted by  $\mathscr{G}_K$ , corresponding to a matrix  $\mathcal{M}_{\mathscr{G}_K} \in \mathbb{R}^{2p \times 2p}$  such that  $\mathbb{I}_{2p} \mathcal{M}_{\mathscr{G}_K}$  is invertible. Consequently, its inverse can be expressed as a finite order k expansion of  $\mathcal{M}_{\mathscr{G}_K}$ , 71
- 72

$$(\mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_K})^{-1} = \sum_{i=0}^k \mathcal{M}_{\mathscr{G}_K}^i, \quad k \le p, \quad \mathcal{M}_{\mathscr{G}_K} = \begin{bmatrix} \mathcal{M}_{\mathscr{G}_K}[U] & \Phi \\ 0 & \mathcal{M}_{\mathscr{G}_K}[V] \end{bmatrix}$$

where k corresponds to the length of the longest path in the DAG and  $\mathcal{M}_{\mathscr{G}_K}[U] = \mathcal{M}_{\mathscr{G}_K}[V] = \mathbf{0}_{p \times p}$ 

#### Recovery from algebraic signature of mixed manifolds 75

The endowed topological structure in  $\mathcal{M}_{\mathscr{G}}$  leads to a spectrum of polynomials  $q_i$ , and an algebraic 76

- variety  $V(\mathcal{I})$  is associated with the generating ideals  $\mathcal{I}$  of those polynomials. We emphasize that the
- dimension of  $V(\mathcal{I})$  determines the identifiability through algebraic quantities.

### 79 3.1 Algebraic structure for stochastic causal representation

Cumulant is an important algebraic signature of its geometric property, as a full order cumulant 80 precisely encodes the entire distribution, including the component-wise and time-wise dependency 81 among variables. This enables the fine-grained mathematical nature of intervention effects beyond traditional mean and variance shifts. Under generic (non-Gaussian) conditions, the d-th order cumu-83 lant of a random variable X = As admits closed-form expressions as  $\kappa_d(X) = \sum_{j=1}^p \kappa_d(s) (A_j)^{\otimes d}$ , 84 which is a secant variety  $\sigma_k(X)$  that views each tensor factor  $(A_j)^{\otimes d}$  as indeterminate. If the matrix 85 A is generic in an open set, the Kruskal rank condition is satisfied and thus  $\sum_{j=1}^p \kappa_d(s)(A_j)^{\otimes d}$  has a 86 unique decomposition. This means the idea  $\mathcal{I}: \langle \kappa_d(X) - \sum_{j=1}^p \kappa_d(s) (A_j)^{\otimes d} \rangle$  has dimension zero. 87 We connect our reasoning to this and illustrate how high-order cumulants can capture sufficient 88 statistical variability in the system, even under temporally independent Gaussian noise. In particular, 89 we adapt the setting in [21], modifying it to allow additive noise that is independent over time. To this 90 end, we establish a result for full recovery of INAR processes in the asymptotic regime  $p \to \infty$ , by 91 formulating identifiability conditions through algebraic-geometric constraints imposed on the latent 92 space. 93

### 3.2 Generic identification from algebraic structure

Let  $J_f$  be the Jacobian matrix of f and  $K=J_f(\mathbb{I}_p-\Phi)^{-1}$ .  $J_{\mathscr{G}}$  and  $K_{\mathscr{G}}$  are augmented matrices by filling  $J_f,K$  in a larger matrix to match the dimension of  $\mathcal{M}$ . We present the following assumption.

# Assumption 1.

98

99

100

101

102

106

107

108

109

110

111

112

113

116

118

119

120

- 1. f is a generic  $C^d$  map with a full rank  $J_f$  almost surely.
- 2. There exist at least p nonzero tensors  $\bar{\kappa}_d(\Delta O_t)$  for  $d \in D := \{d \mid \bar{\kappa}_{d+1}(\Delta O_t) = 0\}$ , where  $\bar{\kappa}_d(\Delta O_t)$  is the difference cumulant truncated at the first-order Taylor expansion of  $O_t := f(Z_t)$ .
  - 3. The ideal  $\mathcal{I}^*$ :  $\langle J_{\mathscr{G}} K_{\mathscr{G}}^{(k)}(\mathbb{I}_{2p} \mathcal{M}^{(k)}), k = 1, 2, \cdots, p \rangle$  has a zero-dimensional associated variety  $\mathcal{V}(\mathcal{I}^*)$

Theorem 1. Under Assum. 1, the latent sources with their causal structure are identifiable up to permutation and component-wise scaling.

The  $C^d$  assumption is strictly weaker than requiring f to be a diffeomorphism, since even in the presence of directional collapse within the latent space, the cumulants may still faithfully transmit the non-redundant dependency structure to the observed domain. The core idea of our proof leverages the propagation behavior of algebraic structure under nonlinear transformations, allowing us to identify latent structure via the observed partial geometric-algebraic information on the mixed manifold  $O_t$ . That is, the cumulant hierarchy of each observed component has finite depth d, and the non-vanishing cumulants up to this order are sufficiently rich to ensure identifiability via tensor decomposition. The proposed rank condition (1) is classic and results in a generically unique decomposition of d order tensor  $\kappa_d(O_t)$ , which uniquely recovers the component  $v_i$  up to permutation and rescaling. Condition (2) ensures we can find such p different tensors so that  $J_{\mathscr{G}}$  and  $(\mathbb{I}_{2p} - \mathcal{M})^{-1}$  can be further disentangled up to the same indeterminacy. The proof strategy converts the identifiability problem into the precise geometry of latent manifolds associated with the time-dependent process. See the proof in B.1.

# 4 Recovery from MUTATAE

### 4.1 Architecture of MUTATE

Building upon our identifiability theory, we formally introduce MUTATE (**MU**lti -**T**ime **A**daptive
Transition Encoder), a novel estimation framework for latent multivariate self-exciting point processes.

MUTATE is designed as a causal representation learning architecture capable of modeling continuoustime stochastic dynamics. Importantly, the framework is modular and can be readily adapted to
other types of stochastic processes with suitable modifications. Unlike prior frameworks that rely
primarily on conditional independence to enforce latent structure, our approach accounts for the

nature of progressively adaptive stochastic processes. In such systems, the filtration  $\mathcal{F}_T$ , which captures the intrinsic history of the process, is defined as  $\sigma\left(\bigcup_{0 < t < T} \sigma(Z_t^{\Delta})\right)$  and grows strictly over 127 128 time (Figure. 1). This dynamically expanding information structure poses unique challenges for both 129 identifiability and representation learning, which MUTATE is explicitly designed to address. 130

**Time adaptive transition module.** We first 131 employ an encoder  $q_\phi(Z_t^\Delta|X_t^\Delta)$  to learn the estimated latents  $Z_t^{(\Delta)} \sim q_\phi(Z_t^\Delta \mid X_t^\Delta)$  as commonly applied in representation learning frame-132 133 134 works. Recall that the latent process is modeled 135 as  $Z_t = \Phi \star Z_t + R_t$ , where  $\Phi$  denotes a global 136 convolution kernel and  $R_t = U + \epsilon_t$  is a resid-137 ual process. Under our weak convergence con-138 dition,  $Z_t$  receives a well-structured representa-139 tion  $Z_t = (\mathbb{I} - \Phi)^{-1} \star (U + \epsilon_t)$ , which is also 140

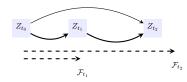


Figure 1: visualization of information loss in increasing filtration

a  $(U + \epsilon_t)$ -measurable process with tractable Power Spectrum Density (PSD):  $S_{Z_t^{(\Delta)}}(w)=(I-\Phi)^{-1}\Sigma_{R_t}(I-\Phi)^{-H}$  , where the baseline U142 is treated as a learnable parameter in the model and  $A^H$  is the Hamilton conjugate transpose of 143 A. w is a continuous frequency variable. The inverse mapping  $f_Z^{-1}$  is an encoder with training parameters of neural networks. Importantly, the learned functional f maps the observation  $Z_t$  to 144 145 a space of independent varying noise through the designated PSD decomposition that enforces  $\Sigma$ 146 to be diagonal and recursively infers  $H^{\dagger} = (\mathbb{I} - \Phi)^{-1}$ . Then, the evaluated prior from the PSD 147 module is sent to calculate the KL divergence. Decomposing S(w), each of the transitions satisfies  $\log p(Z_t^{(\Delta)}|\mathcal{F}_{t-}) = \log p[(I-\Phi)\star R_t^{\Delta}]$ , which is the main part of the latent prior estimation. Our model is trained based on the Variational Auto-encoding framework. Therefore, we aim to max-149 150 imize the  $\log$  likelihood of observation  $\log p_{data}(X)$  through the rule of the ELBO lower bound: 151  $ELBO = -\mathcal{L}_{recon} - \alpha \mathcal{L}_{KL}.$ 152

# 4.2 Simulation Study

141

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

To validate our identifiability results, we evaluate against several representative baselines, including TDRL [4], BetaVAE [22], SlowVAE [23], and PCL [7]. Among them, PCL and TDRL incorporate temporal dependencies by leveraging historical information and explicitly enforcing conditional independence among latent variables to recover underlying dynamics. In contrast, BetaVAE and SlowVAE assume independent latent components and disregard any time-delayed mechanisms. A detailed simulation procedure is included in D.1.

Performance of all baselines and our model is shown in Table 1 and extended results are reported in Table A2. During training, both BetaVAE and SlowVAE tend to converge prematurely, typically reaching a local optimum within the first epoch and triggering early stopping. This behavior highlights their limitations in modeling temporal structures essential for identifying latent event-driven processes. TDRL performs reasonably when the lag module is set to a longer one (we use L=9 in experiments) since it can harness shorter temporary contextual information. It is noticed that our identifiability can be readily applied to the prior framework by either adding the domain index in synthetic datasets or modulating the distribution shifts that change pairs of edges in the latent space. However, we also realize that the fully non-parametric setting is hard to interpret since our identifiability avoids such a case.

Table 1: MCC Scores with standard deviations for five kernels

Method	Ave.	Exponential	Powerlaw	Rectangular	nonlinear	nonparametric
TDRL [4]	0.599	$0.593 \pm 0.028$	$0.609\pm0.043$	$0.618 \pm 0.056$	$0.556 \pm 0.016$	<b>0.616</b> ±0.043
BetaVAE [22]	0.141	$0.153\pm0.863$	$0.128 \pm 0.077$	$0.128\pm0.078$	$0.146 \pm 0.108$	$0.149\pm0.096$
SlowVAE [23]	0.115	$0.108\pm0.075$	$0.104\pm0.073$	$0.104\pm0.073$	$0.126\pm0.074$	$0.131\pm0.076$
PCL [7]	0.375	$0.395 \pm 0.034$	$0.330 \pm 0.029$	$0.330\pm0.029$	$0.414 \pm 0.028$	$0.404 \pm 0.028$
MUTATE(ours)	0.837	$0.853 \pm 0.218$	<b>0.938</b> ±0036	<b>0.879</b> ±0.102	<b>0.921</b> ±0.029	<b>0.598</b> ±0.013

# References

- [1] Judea Pearl. Causality. Cambridge university press, 2009.
- [2] Aapo Hyvärinen and Petteri Pajunen. Nonlinear independent component analysis: Existence and uniqueness results. *Neural Networks*, 12(3):429–439, April 1999.
- 174 [3] Peter Sorrenson, Carsten Rother, and Ullrich Köthe. Disentanglement by Nonlinear ICA with General Incompressible-flow Networks (GIN), January 2020.
- [4] Weiran Yao, Guangyi Chen, and Kun Zhang. Temporally Disentangled Representation Learning.
   Advances in Neural Information Processing Systems, 35:26492–26503, December 2022.
- [5] Xiangchen Song, Weiran Yao, Yewen Fan, Xinshuai Dong, Guangyi Chen, Juan Carlos Niebles,
   Eric Xing, and Kun Zhang. Temporally Disentangled Representation Learning under Unknown
   Nonstationarity. Advances in Neural Information Processing Systems, 36:8092–8113, December
   2023.
- [6] Aapo Hyvarinen and Hiroshi Morioka. Unsupervised Feature Extraction by Time-Contrastive
   Learning and Nonlinear ICA. Advances in Neural Information Processing Systems, 29, 2016.
- [7] Aapo Hyvarinen and Hiroshi Morioka. Nonlinear ICA of Temporally Dependent Stationary
   Sources. In *Artificial Intelligence and Statistics*, pages 460–469. PMLR, April 2017. ISSN: 2640-3498.
- [8] Aapo Hyvarinen, Hiroaki Sasaki, and Richard Turner. Nonlinear ICA Using Auxiliary Variables
   and Generalized Contrastive Learning. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 859–868. PMLR, April 2019. ISSN: 2640-3498.
- [9] Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional Causal Representation Learning. In *International Conference on Machine Learning*, pages 372–407. PMLR, July 2023. ISSN: 2640-3498.
- [10] Chandler Squires, Anna Seigal, Salil S Bhate, and Caroline Uhler. Linear causal disentanglement
   via interventions. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt,
   Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference* on Machine Learning, volume 202 of Proceedings of Machine Learning Research, pages
   32540–32560. PMLR, 23–29 Jul 2023.
- [11] Yibo Jiang and Bryon Aragam. Learning Nonparametric Latent Causal Graphs with Unknown
   Interventions. Advances in Neural Information Processing Systems, 36:60468–60513, December
   200
- [12] Simon Bing, Urmi Ninad, Jonas Wahl, and Jakob Runge. Identifying Linearly-Mixed Causal
   Representations from Multi-Node Interventions. In *Causal Learning and Reasoning*, pages
   843–867. PMLR, March 2024. ISSN: 2640-3498.
- Simon Buchholz, Goutham Rajendran, Elan Rosenfeld, Bryon Aragam, Bernhard Schölkopf,
   and Pradeep Ravikumar. Learning Linear Causal Representations from Interventions under
   General Nonlinear Mixing. Advances in Neural Information Processing Systems, 36:45419–45462, December 2023.
- Patricia Reynaud-Bouret and Sophie Schbath. Adaptive estimation for Hawkes processes; application to genome analysis. *The Annals of Statistics*, 38(5), October 2010.
- [15] Lars Lorch, Andreas Krause, and Bernhard Schölkopf. Causal modeling with stationary diffusions. In Sanjoy Dasgupta, Stephan Mandt, and Yingzhen Li, editors, *Proceedings of The 27th International Conference on Artificial Intelligence and Statistics*, volume 238 of *Proceedings of Machine Learning Research*, pages 1927–1935. PMLR, 02–04 May 2024.
- 214 [16] Ali Torkamani and Nicholas J. Schork. Identification of rare cancer driver mutations by network reconstruction. *Genome Research*, 19(9):1570–1578, September 2009.

- [17] Matthew H. Bailey, Collin Tokheim, Eduard Porta-Pardo, Sohini Sengupta, and *et al.* Bertrand.
   Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell*, 173(2):371–385.e18, April 2018.
- 219 [18] Mona Nourbakhsh, Kristine Degn, Astrid Saksager, Matteo Tiberti, and Elena Papaleo. Predic-220 tion of cancer driver genes and mutations: the potential of integrative computational frameworks. 221 *Briefings in Bioinformatics*, 25(2):bbad519, January 2024.
- Emmanuel Bacry and Jean-Francois Muzy. Second order statistics characterization of Hawkes processes and non-parametric estimation, January 2014.
- [20] Matthias Kirchner. Hawkes and INAR(\$\infty\$) processes. *Stochastic Processes and their Applications*, 126(8):2494–2525, August 2016. arXiv:1509.02007 [math].
- [21] Paula Leyes Carreno, Chiara Meroni, and Anna Seigal. Linear causal disentanglement via
   higher-order cumulants. arXiv preprint arXiv:2407.04605, 2024.
- [22] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick,
   Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a
   constrained variational framework. In *International conference on learning representations*,
   2017.
- [23] David Klindt, Lukas Schott, Yash Sharma, Ivan Ustyuzhaninov, Wieland Brendel, Matthias
   Bethge, and Dylan Paiton. Towards Nonlinear Disentanglement in Natural Data with Temporal
   Sparse Coding, March 2021. arXiv:2007.10930 [stat].
- Patrick Billingsley. *Convergence of probability measures*. Wiley series in probability and statistics Probability and statistics section. Wiley, New York Weinheim, 2. ed edition, 1999.
- [25] Emmanuel Bacry, Iacopo Mastromatteo, and Jean-François Muzy. Hawkes Processes in Finance.

  Market Microstructure and Liquidity, 01(01):1550005, June 2015.
- 239 [26] Alan G Hawkes. Spectra of some self-exciting and mutually exciting point processes. 1971.
- [27] Daryl J. Daley and David Vere-Jones. An introduction to the theory of point processes. 1:
   Elementary theory and methods. Springer, New York NY, 2. ed., 2. corr. print edition, 2005.
- [28] Weiran Yao, Yuewen Sun, Alex Ho, Changyin Sun, and Kun Zhang. Learning Temporally
   Causal Latent Processes from General Temporal Data, February 2022. arXiv:2110.05428 [cs, stat].
- [29] Pierre Brémaud and Laurent Massoulié. Stability of nonlinear Hawkes processes. *The Annals of Probability*, 24(3), July 1996.
- 247 [30] Erhan Cinlar and RA Agnew. On the superposition of point processes. *Journal of the Royal Statistical Society: Series B (Methodological)*, 30(3):576–581, 1968.
- [31] Susan L Albin. On poisson approximations for superposition arrival processes in queues.
   Management Science, 28(2):126–137, 1982.
- Payam Dibaeinia and Saurabh Sinha. Sergio: a single-cell expression simulator guided by gene regulatory networks. *Cell systems*, 11(3):252–271, 2020.
- [33] Haixu Wu, Jiehui Xu, Jianmin Wang, and Mingsheng Long. Autoformer: Decomposition
   Transformers with Auto-Correlation for Long-Term Series Forecasting. Advances in Neural
   Information Processing Systems, 34:22419–22430, December 2021.
- Z56 [34] Yong Liu, Tengge Hu, Haoran Zhang, Haixu Wu, Shiyu Wang, Lintao Ma, and Mingsheng
   Long. iTransformer: Inverted Transformers Are Effective for Time Series Forecasting, March
   258 2024. arXiv:2310.06625 [cs].
- Ruichu Cai, Zhifang Jiang, Zijian Li, Weilin Chen, Xuexin Chen, Zhifeng Hao, Yifan Shen,
   Guangyi Chen, and Kun Zhang. From Orthogonality to Dependency: Learning Disentangled
   Representation for Multi-Modal Time-Series Sensing Signals, May 2024.

- [36] Xiangchen Song, Zijian Li, Guangyi Chen, Yujia Zheng, Yewen Fan, Xinshuai Dong, and Kun Zhang. Causal temporal representation learning with nonstationary sparse transition. In
   A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang, editors,
   Advances in Neural Information Processing Systems, volume 37, pages 77098–77131. Curran Associates, Inc., 2024.
- [37] Kun Zhang, Shaoan Xie, Ignavier Ng, and Yujia Zheng. Causal Representation Learning from
   Multiple Distributions: A General Setting, February 2024.
- [38] Massil Achab, Emmanuel Bacry, Stéphane Gaïffas, Iacopo Mastromatteo, and Jean-François
   Muzy. Uncovering Causality from Multivariate Hawkes Integrated Cumulants. 2018.
- [39] Vanessa Didelez. Graphical Models for Marked Point Processes Based on Local Independence.
   Journal of the Royal Statistical Society Series B: Statistical Methodology, 70(1):245–264,
   February 2008.
- [40] Phillip Lippe, Sara Magliacane, Sindy Löwe, Yuki M. Asano, Taco Cohen, and Efstratios
   Gavves. Causal Representation Learning for Instantaneous and Temporal Effects in Interactive
   Systems. 2023.
- [41] Dingling Yao, Caroline Muller, and Francesco Locatello. Marrying Causal Representation
   Learning with Dynamical Systems for Science. 2024.
- <sup>279</sup> [42] Amir Mohammad Karimi Mamaghan, Andrea Dittadi, Stefan Bauer, Karl Henrik Johansson, and Francesco Quinzan. Diffusion-Based Causal Representation Learning. 26(7):556, 2024.
- 281 [43] Alexander Sokol. Intervention in Ornstein-Uhlenbeck SDEs, August 2013.
- 282 [44] Stephan Bongers, Tineke Blom, and Joris M. Mooij. Causal Modeling of Dynamical Systems,283 March 2018.
- [45] Stephan Bongers, Tineke Blom, and Joris M. Mooij. Causal Modeling of Dynamical Systems,
   March 2022. arXiv:1803.08784 [cs].
- 286 [46] Philip Boeken and Joris M. Mooij. Dynamic Structural Causal Models, June 2024.
- <sup>287</sup> [47] Alexander Sokol and Niels Richard Hansen. Causal interpretation of stochastic differential equations. *Electronic Journal of Probability*, 19(none), January 2014. arXiv:1304.0217 [math].
- [48] Ivana Bozic, Tibor Antal, Hisashi Ohtsuki, Hannah Carter, Dewey Kim, Sining Chen, Rachel
   Karchin, Kenneth W. Kinzler, Bert Vogelstein, and Martin A. Nowak. Accumulation of driver
   and passenger mutations during tumor progression. *PNAS*, oct 2010.
- Johannes G. Reiter, Alvin P. Makohon-Moore, Jeffrey M. Gerold, Alexander Heyde, Marc A. Attiyeh, Zachary A. Kohutek, Collin J. Tokheim, Alexia Brown, Rayne M. DeBlasio, Juliana Niyazov, Amanda Zucker, Rachel Karchin, Kenneth W. Kinzler, Christine A. Iacobuzio-Donahue, Bert Vogelstein, and Martin A. Nowak. Minimal functional driver gene heterogeneity among untreated metastases. *Science*, 361(6406):1033–1037, September 2018.
- 297 [50] Ali Torkamani and Nicholas J. Schork. Identification of rare cancer driver mutations by network reconstruction. *Genome Research*, 19(9):1570–1578, September 2009.
- 299 [51] Benjamin J. Raphael, Jason R. Dobson, Layla Oesper, and Fabio Vandin. Identifying driver mu-300 tations in sequenced cancer genomes: Computational approaches to enable precision medicine. 301 *Genome Medicine*, 6(1):1–17, December 2014.
- Ruth Nussinov, Chung-Jung Tsai, and Hyunbum Jang. A New View of Activating Mutations in Cancer. *Cancer Research*, 82(22):4114–4123, November 2022.

304	Su	Supplement to						
305 306	"(	"Causal Representation Meets Stochastic Modeling"						
307	Appendix organization:							
308								
309	A	Useful Lemmata	9					
310		A.1 Preliminary lemmas	ç					
311		A.2 Background of Hawkes Process	ç					
312		A.3 Cumulants and Tensors	10					
313	В	B Proof of Identifiability Theory						
314		B.1 Proof of Thm. 1	10					
315	C	Proof of Supportive Results	17					
316		C.1 Proof of Lem. 1	17					
317	D	Detailed MUTATE Configuration	18					
318		D.1 Simulation Regime	18					
319		D.2 Prior decomposition of time-adaptive module	19					
320		D.3 Explicit control for convolution prior	21					
321		D.4 Extended Results	22					
322	E	Related Work	22					

323 F Conclusion and Limitation

# **Useful Lemmata**

325

#### Preliminary lemmas 326

**Lemma A.1** (Weak Convergence [24]). Let (S, S) be a Polish space equipped with its Borel  $\sigma$ -327 algebra, and let  $\{Z_n\}_{n\in\mathbb{N}}$  and Z be S-valued random elements defined on a common probability 328 space. Then the sequence  $\{Z_n\}$  converges in distribution (i.e., weakly) to Z, denoted  $Z_n \Rightarrow Z$ , if 329 330

$$\lim_{n \to \infty} \mathbb{E}[f(Z_n)] = \mathbb{E}[f(Z)]$$

 $\lim_{n\to\infty}\mathbb{E}[f(Z_n)]=\mathbb{E}[f(Z)]$  for all bounded continuous functions  $f:S\to\mathbb{R}.$ 331

The proof and demonstration of this lemma is classic in basic probability that we omit here. The 332 weak convergence, in most cases, corresponds to the convergence of finite dimension distribution of a 333 process or a variable. 334

**Lemma A.2** (Tightness of the Measure  $\mathbb{P}_{Z(\Delta)}$ ). Let  $\{Z_n^{(\Delta)}\}_{n\in\mathbb{N}}$  be a sequence of S-valued random elements (e.g., stochastic processes or path evaluations) indexed by  $\Delta$  and defined on a Polish space 335 336 S with Borel  $\sigma$ -algebra. Then the sequence of corresponding probability measures  $\{\mathbb{P}_{Z_{\bullet}^{(\Delta)}}\}$  is tight. 337 In particular, any subsequence admits a further weakly convergent subsequence. 338

Tightness of a sequence of probability measures ensures the existence of well-behaved subsequences: 339 every subsequence admits a further weakly convergent subsequence. This property is particularly 341 useful in Polish spaces, where tightness is equivalent to relative compactness (precompactness) under the weak topology. However, it is important to note that precompactness does not imply full 342 compactness; in general, a tight sequence need not converge without an additional uniqueness or 343 limit identification argument. Thus, tightness provides necessary control over subsequential behavior, 344 but does not guarantee full convergence of the entire sequence. 345

**Lemma A.3** (Higher-Order Moment Bound Implies Lower-Order Bounds). Let  $\{Z_n\}_{n\in\mathbb{N}}$  be a 346 sequence of real-valued random variables defined on a common probability space. Fix an integer 347 d > 0. Suppose there exists a constant C > 0 such that 348

$$\sup_{n\in\mathbb{N}}\mathbb{E}[|Z_n|^d]\leq C.$$

 $\sup_{n \in \mathbb{N}} \mathbb{E}[|Z_n|^d] \leq C.$  Then for any  $0 , there exists a constant <math>C_p > 0$  such that  $\sup_{n \in \mathbb{N}} \mathbb{E}[|Z_n|^p] \leq C_p.$ 

$$\sup_{n\in\mathbb{N}}\mathbb{E}[|Z_n|^p]\leq C_p.$$

#### 350 A.2 Background of Hawkes Process

**Assumption A.4** (Stability and stationary Increment, Proposition 1 in [25]). The process  $N_t$  has 351 asymptotically stationary increments, and intensity  $\lambda_t$  is asymptotically stationary if the kernel 352 satisfies the assumption: 353

$$\rho_{\Phi(t)} = \|\Phi(t)\| = \int_0^t |\Phi(t)| \ dt \text{ has spectral radium smaller than 1}$$
 (A1)

Asm. A.4 gives a necessary condition so that the point process has stable, stationary increments 354 in its intensity. In particular, it means the entire process tends to be stable with an unknown but 355 fixed expectation of the conditional intensity  $\mathbb{E}[\lambda_t^i] = \Lambda^i$ . Restricted by the stationary increment 357 assumption, the existence of the corresponding process is ensured by Lem. 3. To illustrate those conditions, we show a simpler version kernel in Example 1. 358

Lemma 3 (Proposition 6 in [20]). If all conditions and results in Asm. A.4 hold almost everywhere, 359 there exists only one determined process whose dynamics match observations with regard to  $\Lambda^i$ . 360

Example 1. Consider a point process whose kernel functions relay causal influence with an exponential 361 decay to other processes. The generating process thus be accordingly 362

$$\lambda_t^i = u^i + \sum_{i=1}^p \int_0^t \alpha^{ij} e^{-\beta(t-t')} dN_{t'}^j$$

shows the exponential kernel triggers influences that are sustaining but decaying as time proceeds. 363 Technically, the induced causal influences, although decaying from inside the system dynamics, will not disappear unless the causal strength  $\alpha = 0$  for all j.

#### A.2.1 Remarks on the filtration

366

In probability theory, the filtration  $\mathcal{F}_t$  is defined as the smallest  $\sigma$ -algebra that renders the intensity 367 process  $\lambda_t$  to be  $\mathcal{F}_t$ -adapted and measurable. This filtration is constructed by the minimal closure 368 under set operations (e.g., union, intersection) over past events, ensuring that  $\lambda_t$  evolves consistently 369 with the observable history [26, 27]. Therefore, for any filtration as its internal history, we have 370  $\mathcal{F}_s \subseteq \mathcal{F}_t$ , for  $s \leq t$ . Note that the filtration  $\mathcal{F}_t$  may theoretically differ from the intrinsic history 371  $\mathcal{H}_t$ , which introduces additional challenges in the evaluation and modeling of point processes. For a comprehensive discussion on scenarios where  $\mathcal{F}_t$  and  $\mathcal{H}_t$  are defined differently, we refer 372 373 the interested reader to [27]. We occasionally overload the notation  $dN_t^i$ , which represents an 374 integral element in stochastic calculus, to distinguish it from its deterministic counterpart. Despite 375 potential similarities in notation, they are fundamentally different: while standard calculus considers 376 infinitesimal increments over fixed mesh widths (e.g., dg(x) as  $\Delta t \to 0$ ), the increment  $dN_t^i$  is a 377 random variable governed by the stochastic process. Specifically, its realization at each infinitesimal 378 interval is drawn from a Bernoulli process with intensity  $\lambda_t^i$ , such that  $\mathbb{P}(dN_t^i > 0 \mid \mathcal{F}_t) = \lambda_t^i, dt$ . In 379 contrast to deterministic differentials,  $dN_i^i$  encapsulates the uncertainty of event occurrences within 380 each interval. The kernel matrix  $\Phi_t$  consists of time-decaying kernel functions that transmit the 381 influence of past events across processes. It captures both time-delayed and causal dependencies, and 382 plays a central role in modeling self-exciting or mutually-exciting dynamics. 383

# 384 A.3 Cumulants and Tensors

Cumulant tensor notation. The d-th order cumulant tensor of a random vector  $X \in \mathbb{R}^p$  is denoted  $\kappa_d(X) \in \mathbb{R}^{p \times \cdots \times p}$ , and is symmetric in all modes. In ICA and CRL settings, cumulants of independent components often admit a CP form:

$$\kappa_d(X) = \sum_{r=1}^R \lambda_r \cdot v_r^{\otimes d},$$

where  $v_r \in \mathbb{R}^p$  and  $\lambda_r \in \mathbb{R}$ . This structure enables identifiability of latent sources from cumulant information.

Tensor notation and operations. We denote an order d tensor as  $\mathcal{T} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_d}$ . The outer product  $u^{(1)} \otimes \cdots \otimes u^{(d)} \in \mathbb{R}^{I_1 \times \cdots \times I_d}$  produces a rank-1 tensor with entries:

$$\mathcal{T}_{i_1,\dots,i_d} = u_{i_1}^{(1)} \cdots u_{i_d}^{(d)}.$$

Given a tensor  $\mathcal{T} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$  and a matrix  $U \in \mathbb{R}^{J \times I_n}$ , the *mode-n* product  $\mathcal{T} \times_n U \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N}$  is defined as:

$$(\mathcal{T} \times_n U)_{i_1,\dots,i_{n-1},j,i_{n+1},\dots,i_N} = \sum_{i_n=1}^{I_n} \mathcal{T}_{i_1,\dots,i_N} \cdot U_{j,i_n}.$$

# 394 B Proof of Identifiability Theory

#### 395 B.1 Proof of Thm. 1

# 396 B.1.1 Useful Lemmas

To potentially identify any latent components of dynamics, we must introduce tensor algebra beyond our current setting as we present the following important results.

Corollary B.1 (CP decomposition). Let  $\mathcal{T} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$  be an order-N tensor. We say that  $\mathcal{T}$  admits an exact rank-R Canonical Polyadic (CP) decomposition if there exist component vectors  $a_r^{(n)} \in \mathbb{R}^{I_n}$  for each  $r = 1, \dots, R$ ,  $n = 1, \dots, N$ , such that:

$$\mathcal{T} = \sum_{r=1}^{R} a_r^{(1)} \otimes a_r^{(2)} \otimes \cdots \otimes a_r^{(N)} = [\![A^{(1)}, A^{(2)}, \dots, A^{(N)}]\!],$$

Where  $A^{(n)}=[a_1^{(n)}\ a_2^{(n)}\ \cdots\ a_R^{(n)}]\in\mathbb{R}^{I_n imes R}$  are the factor matrices.

Corollary B.2. Let  $X^{(1)}, X^{(2)}, \dots, X^{(n)} \in \mathbb{R}^p$  be independent random vectors with nonzero d-th order cumulants, such that each admits the form

$$\kappa_d(X^{(i)}) = \lambda_i \cdot v_i^{\otimes d}, \quad \text{for } i = 1, \dots, n,$$

- with  $v_i \in \mathbb{R}^p$  and  $\lambda_i \in \mathbb{R} \setminus \{0\}$ . Let  $\mathcal{T} := \kappa_d(X^{(1)} + \dots + X^{(n)}) \in \mathbb{R}^{p \times \dots \times p}$  be the d-th order cumulant tensor of their sum.
- 407 Assume that the matrix  $V = [v_1 \ v_2 \ \cdots \ v_n] \in \mathbb{R}^{p \times n}$  satisfies

$$\operatorname{krank}(V) \ge \left\lceil \frac{2n + (d-1)}{d} \right\rceil.$$

408 Then the CP decomposition

$$\mathcal{T} = \sum_{i=1}^{n} \lambda_i \cdot v_i^{\otimes d}$$

409 is unique up to scaling and permutation.

#### 410 **B.1.2 Proof of Thm. 1**

We prove Thm. 1 by showing that the tuple  $(f,\Phi,U)$  is identifiable up to a component-wise transformation and permutation. Our proof is based on the dimension of the associated variety defining special hypersurfaces in a polynomial ring  $K^n$ . We study the nonlinear propagation of the cumulant structure to find the identifiability conditions for INAR processes. Given a generic nonlinear f, its exact cumulant  $\kappa_d(O_t)$  follows an order d expansion of with Bell polynomial coefficients. For a smooth map  $f: Z_t \to f(Z_t)$ , we can construct  $O: = f(Z_{t+\Delta t})$  using Taylor expansion:

$$f(Z_{t+\Delta t}) = f(Z_t) + \frac{\partial}{\partial Z_t} f(Z_t) \Delta Z_t + \frac{1}{2} \frac{\partial^2}{\partial Z_t^2} f(Z_t) (\Delta Z_t)^2 + o(\Delta Z_t^3)$$

At this time, the expansion has rather abnormal behavior, as the order can be extremely large. However, the truncated expansion at order 1 has intriguing theoretical attractions. Let higher-order components be  $\mathcal{R}(1)$ , and recall  $Z_t = H \star \epsilon_t$ , we obtain the truncated differential process  $\Delta \tilde{f}(Z_t)$ , denoted as:

$$\Delta f(Z_t) - \mathcal{R}(1) = J_f \sum_{k=1}^t H_{t-s} \epsilon_s \tag{A1}$$

We treat all quantities appearing in Eq. (A1) as indeterminates in a polynomial ring

$$R = k[\{\Delta f(Z_t)\}_t, \mathcal{R}(1), J_f, \{H_{t-s}\}_s, \{\epsilon_s\}_s],$$

where k is a base field such as  $\mathbb R$  or  $\mathbb C$ . For each time index t, the defining polynomial is  $g_t:=$   $\Delta f(Z_t)-\mathcal R(1)-J_f\sum_{s=1}^t H_{t-s}\epsilon_s\in R$ . This polynomial generates the principal ideal  $\mathcal I_t=$   $\langle g_t\rangle\subset R$ , and considering all time indices  $t=1,2,\ldots,T$ , we obtain the global ideal  $\mathcal I=$   $\langle g_1,g_2,\ldots,g_T\rangle\subset R$ . The corresponding affine variety is then

$$V(\mathcal{I}) = \Big\{ (Z_t, \Delta f(Z_t), J_f, H_{t-s}, \epsilon_s, \mathcal{R}(1)) \in k^N \ \Big| \ g_t = 0 \text{ for all } t \Big\}.$$

It is evident that  $V(\mathcal{I})$  is positive-dimensional, since the defining relations do not specify finitely many points. To obtain more structure, we consider higher-order statistics. In particular, the d-th order cumulant tensor of the transformed increments takes the form

$$\kappa_d \left( \Delta \tilde{f}(Z_t) \right) = \sum_{s=1}^t \sum_{j=1}^p \kappa_d^{(j)}(\epsilon) \cdot \left( J_f H_{t-s}^{(:,j)} \right)^{\otimes d}.$$

This expression shows that the cumulant naturally defines a point in the projective tensor space

$$\mathbb{P}(V^n \otimes V^n \otimes \cdots \otimes V^n),$$

where the number of tensor factors equals d. Hence, while the affine variety  $V(\mathcal{I})$  is too large to give identifiability, the cumulant tensors lift the problem into a projective geometric setting, where connections to secant varieties of the Veronese embedding provide a natural framework for studying uniqueness and decomposition. So far, the generic mixing f is preserved by its Jacobian matrix  $J_f$ ; hence, identifying  $J_f$  is equivalent to the recovery of f up to a constant. Now, we are ready to prove our main theorem. Without loss of generality, we write  $J_f$  as F since they behave the same way in an algebraically closed field.

Step 1: Uniqueness of mixing kernel  $F(\mathbb{I}_p-\Phi)^{-1}$  We prove this supporting result via an extension of Proposition 3.1 in [21]. In the classical linear source decomposition (LSD) setting, the d-th order cumulant of X admits the following tensor decomposition:  $\kappa_d(X) = \sum_{i=1}^q \kappa_d(\epsilon_i) \cdot (B_i)^{\otimes d}$  under the assumption that the components of  $\epsilon$  are non-Gaussian with non-vanishing d-order cumulants, and that multiple interventions are available. The sufficient order d cumulant of each  $Z_t$  for a fixed  $t=t_i$  is

$$\kappa_d(Z_t) = \kappa_d[(I - \Phi)^{-1} \star \epsilon_t] = \kappa_d[\sum_{k=1}^t H_{t-s} \epsilon_s]$$

for each s, the linear transformation  $H_{t-s}$  results in a multi-linear transformation of their cumulants

$$\kappa_d(H_{t-s}\epsilon_k) = (H_{t-s})^{\otimes d} \mathcal{C}_{\epsilon_s}^d = \sum_{i=1}^p \kappa_d(\epsilon_s^i) (H_{t-s})_j^{\otimes d}$$

The full order d cumulant of the mixed manifold  $O_t$  is

$$\kappa_{d}(O_{t}) = \kappa_{d}(FZ_{t}, FZ_{t}, \dots, FZ_{t}) \\
d \text{ times} = F^{\otimes d} \cdot \kappa_{d}(Z_{t}, Z_{t}, \dots, Z_{t}) \\
= \underbrace{F \otimes F \otimes \dots \otimes F}_{d \text{ times}} \cdot \kappa_{d} \left( \sum_{s_{1}=1}^{t} H_{t-s_{1}} \epsilon_{s_{1}}, \sum_{s_{2}=1}^{t} H_{t-s_{2}} \epsilon_{s_{2}}, \dots, \sum_{s_{d}=1}^{t} H_{t-s_{d}} \epsilon_{s_{d}} \right) \\
= \underbrace{F \otimes F \otimes \dots \otimes F}_{d \text{ times}} \cdot \sum_{s_{1}=1}^{t} \dots \sum_{s_{d}=1}^{t} \kappa_{d} \left( H_{t-s_{1}} \epsilon_{s_{1}}, H_{t-s_{2}} \epsilon_{s_{2}}, \dots, H_{t-s_{d}} \epsilon_{s_{d}} \right) \\
= \underbrace{F \otimes F \otimes \dots \otimes F}_{d \text{ times}} \cdot \sum_{s_{1}=1}^{t} \dots \sum_{s=1}^{t} \sum_{j=1}^{p} \kappa_{d} \left( H_{t-s}^{(:,j)} \epsilon_{s}^{(j)}, H_{t-s}^{(:,j)} \epsilon_{s}^{(j)}, \dots, H_{t-s}^{(:,j)} \epsilon_{s}^{(j)} \right) \\
= \underbrace{F \otimes F \otimes \dots \otimes F}_{d \text{ times}} \cdot \left( \sum_{s=1}^{t} \sum_{j=1}^{p} \kappa_{d}^{(j)} (\epsilon) \cdot \underbrace{H_{t-s}^{(:,j)} \otimes H_{t-s}^{(:,j)} \otimes \dots \otimes H_{t-s}^{(:,j)}}_{d \text{ times}} \right) \\
= \left( \sum_{s=1}^{t} \sum_{j=1}^{p} \kappa_{d}^{(j)} (\epsilon) \cdot \underbrace{F \otimes F \otimes \dots \otimes F}_{d \text{ times}} \cdot \left( H_{t-s}^{(:,j)} \right)^{\otimes d} \right) \\
= \sum_{s=1}^{t} \sum_{s=1}^{p} \kappa_{d}^{(j)} (\epsilon) \cdot \left( F H_{t-s}^{(:,j)} \right)^{\otimes d} \right) \tag{A3}$$

Unlike in a time-free process, the joint cumulant of a time process is of order d that is coupled with the number of time lags:

$$\kappa_d(O_{t_1}, \dots, O_{t_d}) = \sum_{s=1}^{\min(t_1, \dots, t_d) - 1} \sum_{j=1}^p \kappa_d^{(j)}(\epsilon) \cdot \bigotimes_{\ell=1}^d \left( FH_{t_{\ell} - s}^{(:,j)} \right) \\
= \sum_{j=1}^p \sum_{s=1}^{\min(t_1, \dots, t_d) - 1} \kappa_d^{(j)}(\epsilon) \cdot \bigotimes_{\ell=1}^d \left( FH_{t_{\ell} - s}^{(:,j)} \right) \tag{A4}$$

We denote the Fourier transform of x(t) with respect to time t as  $\mathcal{F}[x](\omega)$ . Using the convolution theorem and linearity of the Fourier transform, we have:

$$\mathcal{F}[\kappa_d(O_t)](\omega) = \mathcal{F}\left[\sum_{s\geq 0}^t \sum_{j=1}^p \kappa_d^{(j)}(\epsilon) \cdot \left(FH_{t-s}^{(:,j)}\right)^{\otimes d}\right]$$

$$= \mathcal{F} \left[ \int_{0}^{t} \sum_{j=1}^{p} \kappa_{d}^{(j)}(\epsilon) \cdot \left( FH(t-s)^{(:,j)} \right)^{\otimes d} ds \right]$$

$$= \left( \sum_{j=1}^{p} \kappa_{d}^{(j)}(\epsilon) \cdot \mathcal{F} \left[ \left( FH^{(:,j)} \right)^{\otimes d} \star \mathcal{U}(t-s) \right] \right)$$

$$= \left( \sum_{j=1}^{p} \kappa_{d}^{(j)}(\epsilon) \cdot \mathcal{F} \left[ \left( FH^{(:,j)} \right)^{\otimes d} \right] (\omega) \cdot \left( \pi \delta(\omega) + \frac{1}{i\omega} \right) \right)$$
(A5)

Since  $\delta(\omega)$  vanishes everywhere except at  $\omega = 0$ , multiplying it by  $\omega$  gives zero, Eq. (A5) yields

$$i\omega \mathcal{F}[\kappa_d(O_t)](\omega) = \left(i\omega \sum_{j=1}^p \kappa_d^{(j)}(\epsilon) \cdot \mathcal{F}\left[\left(FH^{(:,j)}\right)^{\otimes d}\right](\omega) \cdot \left(\pi\delta(\omega) + \frac{1}{i\omega}\right)\right)$$
$$= \left(\sum_{j=1}^p \kappa_d^{(j)}(\epsilon) \cdot \mathcal{F}\left[\left(FH^{(:,j)}\right)^{\otimes d}\right](\omega) \cdot i\omega \cdot \left(\pi\delta(\omega) + \frac{1}{i\omega}\right)\right).$$

450 which reads

$$\left(\sum_{j=1}^{p} \kappa_{d}^{(j)}(\epsilon) \cdot \mathcal{F}\left[\left(FH^{(:,j)}\right)^{\otimes d}\right](\omega) \cdot i\omega \cdot \left(\pi\delta(\omega) + \frac{1}{i\omega}\right)\right) = \left(\sum_{j=1}^{p} \kappa_{d}^{(j)}(\epsilon) \cdot \mathcal{F}\left[\left(FH^{(:,j)}\right)^{\otimes d}\right](\omega) \cdot 1\right)$$
(A6)

By assuming non-Gaussianity in  $\epsilon_t$ , for each j, Eq. (A6), hence  $i\omega\mathcal{F}[\kappa_d(O_t)](\omega)$  has a unique decomposition of the summation of rank-1 tensor. Therefore, each column of the sub-linear mixing transferring matrix  $\mathcal{F}\left[\left(FH^{(:,j)}\right)\right]$  is theoretically recovered up to a scaling and permutation  $\pi$  if all assumptions made are satisfied in  $\Phi$ . This indicates that even if one needs to calculate the tensor decomposition unnecessarily, such uniqueness guarantees the possibility of further disentanglement. In the sequel,  $\mathcal{F}\left[\left(FH^{(:,j)}\right)\right]DP$  is available; we thus obtain the unique indeterminacy as an immediate result of Lemma. B.3.

Lemma B.3. Consider  $\mathbb{F}$  is an algebraically closed field, the unknown indeterminacy DP is preserved in  $\mathbb{F}$ , that is the following relation

$$\begin{split} FH_{\tau}^{(:,j)} &= \hat{F}\hat{H}_{\tau}^{(:,j)} \, D_{j,\tau} P_{j,\tau}, \\ \textit{with} \, (P_{j,\tau}, D_{j,\tau}) &\in \{(P_j, D_j) \mid \mathcal{F}[FH^{(:,j)}] \, D_j P_j = \mathcal{F}[\hat{F}\hat{H}^{(:,j)}]\}. \end{split}$$

Proof. Let  $\mathbb F$  be a field and  $n\geq 1$  an integer. The general linear group of degree n over  $\mathbb F$  is

$$\operatorname{GL}_n(\mathbb{F}) := \{ A \in M_n(\mathbb{F}) \mid \det(A) \neq 0 \}.$$

Equivalently, if V is a n-dimensional vector space over  $\mathbb{F}$ ,

$$GL(V) := Aut_{\mathbb{F}}(V) = \{ T : V \to V \text{ linear isomorphisms } \},$$

and any choice of basis identifies GL(V) with  $GL_n(\mathbb{F})$ . It is obvious that  $\mathcal{F}[\mathbb{F}]$  is exactly a subgroup 462 of  $GL(\mathbb{F})$  as the group  $GL_n(\mathbb{F})$  satisfies: i) It is precisely the set of all invertible linear transformations 463 (invertible matrices). ii) If  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{C}$ , then  $\mathrm{GL}_n(\mathbb{F})$  is an open subset of  $M_n(\mathbb{F})$  since  $\mathrm{GL}_n(\mathbb{F}) =$ 464  $\det^{-1}(\mathbb{F}\setminus\{0\})$ , and it is a Lie group. By the definition of kernel matrix, one notes i) trivially holds 465 due to the maximal spectrum being less than 1. For ii), det<sup>-1</sup> denotes the preimage of the open set 466  $\mathbb{F}\setminus\{0\}$  under det. Cutting the one-dimensional line at 0 produces two open intervals (for  $\mathbb{F}=\mathbb{R}$ ) or 467 a punctured plane (for  $\mathbb{F} = \mathbb{C}$ ), hence the preimage is open in  $M_n(\mathbb{F})$ . Therefore, the permutation 468 and scaling must be preserved in  $M \in \mathbb{R}^{p \times p}$ . 469

So far, the original kernel mixing matrix  $FH_{\tau}$  is recovered up to the same permutation and scaling for any  $\tau$ . In the sequel, what needs to be proved is the recovery of the causal structure as well as its full parameter space. Our proof focuses on the polynomial system and its associated ideal  $\mathcal{I}$  generated by the multi-linear constrained polynomial system.

# Step 2: Uniqueness of causal kernel $\Phi$

*Example* 2. Consider a kernel matrix  $\Phi \in \mathbb{R}^{2 \times 2}$  with internal arrows allowed:

$$\Psi_U = \begin{bmatrix} 0 & \psi_{12} \\ 0 & 0 \end{bmatrix}, \quad \Psi = \begin{bmatrix} \phi_{11} & \phi_{12} \\ \phi_{21} & \phi_{22} \end{bmatrix}, \Psi_V = \begin{bmatrix} 0 & \psi_{34}' \\ 0 & 0 \end{bmatrix}$$

where  $\Psi_U, \Psi_V$  encodes the internal arrows of the sub-graph  $\mathcal{G}_U$  and  $\mathcal{G}_V$ ;  $\Phi$  encodes time-delayed 476 kernel effects from s to t.

The corresponding expanded kernel matrix  $\mathcal{M} \in \mathbb{R}^{4 \times 4}$  is

$$\mathcal{M} = \begin{bmatrix} 0 & \psi_{12} & \phi_{11} & \phi_{12} \\ 0 & 0 & \phi_{21} & \phi_{22} \\ 0 & 0 & 0 & \psi'_{34} \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

*Proof.* Let  $\mathcal{M}_{\mathcal{G}_k} \in \mathbb{R}^{2p \times 2p}$  be the kernel matrix associated with the bipartite graph, where the 479 variables are partitioned into two subsets U and V. Consider a topological ordering where all nodes 480 in U precede those in V. 481

Since edges from V to U are forbidden by the bipartite structure, no element in the lower-left block 482 of  $\mathcal{M}$  can be nonzero. Moreover, edges within U are not allowed, so the diagonal and upper-left 483 block corresponding to U have zeros on the diagonal. The edges within V form a DAG, and under a 484 topological ordering of V, the corresponding block in  $\mathcal{M}$  is strictly upper-triangular. Therefore,  $\mathcal{M}$ 485 as a whole is strictly upper-triangular, which implies that it represents a DAG. 486

Because  $\mathcal{M}_{\mathcal{G}_K}$  is strictly upper-triangular, it is nilpotent. Let k denote the length of the longest directed path in the DAG. Then  $\mathcal{M}_{\mathcal{G}_K}^{k+1}=0$ , and the inverse of  $\mathbb{I}-\mathcal{M}_{\mathcal{G}_K}$  can be expressed as a finite 487 488 sum of powers of  $\mathcal{M}$ : 489

$$(\mathbb{I} - \mathcal{M}_{\mathcal{G}_K})^{-1} = \sum_{i=0}^k \mathcal{M}_{\mathcal{G}_K}^i.$$

Each term  $\mathcal{M}_{G_K}^i$  corresponds to contributions from paths of length i in the DAG. This shows that the 490 inverse is fully determined by path products up to the longest path length k, completing the proof.  $\Box$ 491

By Lem. 2, identifying the space of all parameters in F and  $\Phi$  is equivalent to solving the following 492 ideal  $\mathcal{I}: \langle F_{\mathscr{G}} - F_{\mathscr{G}}(\mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_K})^{-1}(\mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_K}) \rangle$ . 493

In latent causal models,  $F_{\mathscr{G}}$ ,  $\mathcal{M}_{\mathscr{G}_K}$  are filled as indeterminates that need to be recovered, where  $F_{\mathscr{G}}$  is the expanded linear mixing obtained by filling  $F \in \mathbb{R}^{n \times p}$  in a larger block-diagonal matrix of size 494 495  $2n \times 2p$ , denoted by  $F_{\mathscr{G}}$ . Remember, we have  $F(\mathbb{I} - \Phi(\tau))^{-1}$  to be unique due to decomposition up to a scaling and permutation. We write  $H_{\mathscr{G}} = (\mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_K})^{-1}$ , then 496 497

$$F_{\mathscr{G}}(\mathbb{I}_{2p}-\mathcal{M}_{\mathscr{G}_K})^{-1}=\begin{bmatrix}F & \mathbf{0}\\ \mathbf{0} & F\end{bmatrix}\begin{bmatrix}H_{\mathscr{G}}(U) & H_{\mathscr{G}}(\Phi)\\ \mathbf{0} & H_{\mathscr{G}}(V)\end{bmatrix}=\begin{bmatrix}FH_{\mathscr{G}}(U) & FH(\Phi)\\ \mathbf{0} & FH_{\mathscr{G}}(V)\end{bmatrix}=K_{\mathscr{G}}$$

Under INAR, the diagonal blocks of  $K_{\mathscr{G}}$  are 0 matrix. Therefore, we have the ideal:  $\mathcal{I}: \langle F_{\mathscr{G}} -$ 498  $K_{\mathscr{G}}(\mathbb{I}_{2p}-\mathcal{M}_{\mathscr{G}_K})$  where  $K_{\mathscr{G}}$  is known because FH is known, so not considered as indeterminate 499 and thus does not contribute any degrees in the dimension of  $\mathcal{I}$ . Clearly, under passively observational 500 settings, recovery of full models is never possible as the current  $\mathcal{I}$  must be positive dimensional, 501 leading to no fixed points defined in the associated variety V. This leads to a central goal to find 502 the number of contexts that indicate sufficient variability or interventional settings, to recover the 503 parameter space. To this end, we need to first discuss important properties of  $F_{\mathscr{G}}$ . 504

**Lemma B.4.**  $F \in \mathbb{R}^{n \times p}$  is a generic full-rank matrix. Then  $F_{\mathscr{G}}$  is full rank with  $\operatorname{rank}(F_{\mathscr{G}}) =$  $2 \cdot \operatorname{rank}(F) = 2 \min(n, p)$ , and it is not generic in an open dense subset of  $\mathbb{R}^{2n \times 2p}$  due to the 506 additional linear constraints imposed by the block-diagonal structure. Consequently, F.g. belongs to a proper linear subvariety of  $\mathbb{R}^{2n \times 2p}$  defined by

$$\mathcal{V}_{I_{\star}} := \left\{ B \in \mathbb{R}^{2n \times 2p} : B = \begin{pmatrix} F' & 0 \\ 0 & F' \end{pmatrix}, \ F' \in \mathbb{R}^{n \times p} \right\}.$$

Therefore, for any square matrix  $A^{2p}$  with rank $(A) \leq r$ , rank $(F_{\mathscr{G}}A) \leq r$ 

505

507

510 *Proof.* This proof is trivial by linear algebra.

Note we have p processes; we generally assume we obtain different contextual information from at least K environments (i.e., K = p), which is a mild condition in causal representation learning. Each

sub-ideal  $\mathcal{I}_k: \langle F_\mathscr{G} - K_\mathscr{G}(\mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_K}) \rangle$  constitutes a polynomial system, denoted by  $\mathcal{S}^k_{\mathbb{P}}$  such that:

$$F_{\mathscr{G}} + K_{\mathscr{G}}^{(k)} \mathcal{M}_{\mathscr{G}_K}^{(k)} = K_{\mathscr{G}}^{(k)}. \tag{A7}$$

There are  $2n \times 2p$  indeterminates for  $F_{\mathscr{G}}$  and  $2|e(\mathcal{G})|$  for  $\mathcal{M}$  since each environment k introduces a new  $\mathcal{M}_{k,j}$ . Considering all K ideals  $(\mathcal{I}_0, \mathcal{I}_2, \cdots, \mathcal{I}_K)$ , we obtain the union of all K varieties:

$$V(\mathcal{I}) = \left\{ \left( F_{\mathscr{G}}, \mathcal{M}_{\mathscr{G}_{K}}^{(0)}, \dots, \mathcal{M}_{\mathscr{G}_{K}}^{(K)} \right) \middle| F_{\mathscr{G}} - K_{\mathscr{G}}^{(k)} \left( \mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_{K}}^{(k)} \right) = 0, \ \forall k \in K \right\}.$$
 (A8)

Adding polynomial constraints by subtracting Eq.(A7) for 0 from that for k obtains:

$$V(\mathcal{I}^{\star}) = \left\{ V(\mathcal{I}) \mid K_{\mathscr{G}}^{(k)} \mathcal{M}_{\mathscr{G}_{K}}^{(k)} - K_{\mathscr{G}}^{(0)} \mathcal{M}_{\mathscr{G}_{K}}^{(0)} - \left( K_{\mathscr{G}}^{(k)} - K_{\mathscr{G}}^{(0)} \right) = 0, \ \forall k \in K \right\}.$$
 (A9)

that induces a coordinate ring  $R/\mathcal{I}$  in a polynomial ring  $R=k(F_{\mathscr{G}_{i,j}},\mathcal{M}_{i,j}^k)$ . The order of the coordinate ring is the dimension of the original variety. We use  $\left(\begin{array}{c} & \\ & \end{array}\right)$  to represent the blocked system so as not to confuse with the matrix bracket. For each  $m\in[n], j\in[p]$ , and  $i\in\operatorname{ch}_{\mathcal{G}}(j)$ , we have  $|\operatorname{ch}_{\mathcal{G}}(j)|$  columns for  $\mathcal{M}_{\mathcal{G}}$ .  $F_v,\mathcal{M}_v^{(k)}$  write entries  $f_{i,j},\mathcal{M}_{m,i}^k$  as vectors, which describe the polynomial constraints as a linear system:

In INAR  $(\infty)$ , each  $\mathcal{M}_{t-s}$  preserves all paths  $\{(\varphi(j) \to \varphi(i) | Z^{\Delta}_{\varphi(j),t-s} \to Z^{\Delta}_{\varphi(i),t}, \mathcal{M}_{i,j} \neq 0\}$  from  $\Phi_{t-s}$  through an isomorphism  $\varphi$ . Therefore, entry  $(\mathbb{I}_{2p} - \mathcal{M})^{-1}_{i,j}$  is the product of  $\mathcal{M}_{n,m}$  for path  $j \to m \to n \to i$ . We drop the time index and graph label whenever the context is clear. Such  $(\mathbb{I}_{2p} - \mathcal{M})^{-1}_{i,j}$  admits the representation

$$(\mathbb{I}_{2p} - \mathcal{M})^{-1} = \mathbb{I} + \mathcal{M}.$$

Results from [21] are applied to get  $\operatorname{rank}(\mathbb{I}_{2p}-\mathcal{M}^{(k)})^{-1}-(\mathbb{I}_{2p}-\mathcal{M}^0)^{-1}\leq 1$ . Using Lem. B.4, we obtain  $\operatorname{rank}(F_{\mathscr{G}}(\mathbb{I}_{2p}-\mathcal{M}^{(k)})^{-1}-F_{\mathscr{G}}(\mathbb{I}_{2p}-\mathcal{M}^0)^{-1})\leq 1$ . Therefore, the left part of Eq.(A10) has columns that are multiples of each other:

$$(K_{\mathscr{A}}^{(k)} - K_{\mathscr{A}}^{(0)})_{l,j} = (K_{\mathscr{A}}^{(0)})_{l,k} \Delta_{k,i}, \quad \Delta := (I - \mathcal{M}^k)^{-1} - (I - \mathcal{M}^0)^{-1}$$
(A11)

We examine the non-zero sub-blocks of the lower-right block  $\star$  in Eq.(A10), which has size |de(j)|  $\text{ch}(j)| \times |\text{ch}(j)|$ . Following the convention, we represent the sunblocks as M[j] and choose the smaller blocks  $[K_{\mathscr{G}}^{(k)} - K_{\mathscr{G}}^{(0)}]$  corresponding to the size of M[j] and write it as b[j]. The dimension of variety  $V(\mathcal{I}^{\star})$  is the dimension of the points  $(\mathcal{M}_{i,j}^0)$ ,  $i \in \text{ch}(j)$  that satisfy:

$$M[j](\mathcal{M}_{i,j}^0) = b[j] \tag{A12}$$

The variety is a null set when the above constraints lead to no solutions. Therefore, we require  $\operatorname{rank}(M[j]) = \operatorname{rank}(M[j]|b[j])$ .  $[M[j] \mid b[j]]$  is the common augmented matrix to check the stability of a polynomial equation system. We conclude our proof by making a formal statement about the dimension of  $\mathcal{V}(\mathcal{I}^*)$  in the next lemma.

Lemma B.5. For generic F and FH arising from the cumulant decomposition, the full generating model is identifiable if and only if the variety  $V(\mathcal{I}^*)$  has dimension zero, that is,

$$\dim(\mathcal{V}(\mathcal{I}^*)) = \sum_{j=1}^q \operatorname{ch}(j) - \operatorname{rank}(M[j]) = 0.$$

*Proof.* Recall the kernel-delayed DAG structure. For the left subset  $U_s$ , each node j has outgoing edges only to nodes i in the right subset  $V_t$ , all of which are direct children of j. By construction, no edges exist within  $U_s$  or within  $V_t$ . Consequently, for each j, we have  $M[j] = \emptyset$ , since de(j) = ch(j).

Step 3: Independent conditions for each  $\Phi$  under the generic F When  $\mathcal{M}_{t-s}$  is fully recovered, the full matrix  $\mathcal{F}_{\mathscr{G}}$  can be obtained as

$$\mathcal{F}_{\mathscr{G}} = K_{\mathscr{G}}(\mathbb{I}_{2p} - \mathcal{M}_{\mathscr{G}_K}).$$

If  $F_v$  is injective, identification of  $F_v$  (and hence  $\mathcal{F}_\mathscr{G}$  and F) is equivalent to identifying each  $\mathcal{M}$  individually, due to the direct multiplication of  $F^{-1}$  (or the pseudo-inverse  $F^{\dagger}$ ) with K. However, identification of the full generating model is hindered by the genericity of F. Even if  $K_\mathscr{G}$  is unique up to the usual indeterminacies, recovering other kernel matrices  $\mathcal{M}_{t-s'}$  requires analogous identifiability conditions for each individual kernel matrix.

Under  $k \in 1, 2, ..., K$  distinct contexts—each introducing sufficient variability in the distribution, or ensuring that each lag k receives at least one intervention that shifts the downstream mechanism—the full latent structure is identifiable up to the same indeterminacy.

Recovery of baseline U Once the full causal structure is recovered up to a scaling and permutation matrix, lower level moments of  $\mathbb{E}[F(\mathbb{I}_p - \Phi)^{-1} \star (U + \epsilon_t)]$  are can be directly computed to find U up to the same indeterminacy.

### **B.1.3** Identifiability under Gaussian Noise

556

567

Now we focus on the case where non-Gaussianity does not hold for the entire time process. When the noise is Gaussian,  $\kappa_d(O_t)=0$  for all  $d\geq 3$ , which leads to a  $\mathbf{0}\in\mathbb{R}^{p^{\times d}}$  (the d-th order zero tensor). The solution of decomposition is infinite, thus FH cannot be recovered up to a column scaling and permutation, nor can the latent transition graph  $\mathcal{G}$ . We argue that preserving only d- order cumulant of order  $d\leq 2$  is a minimal building block for identification.  $\kappa_2(O_t)$  is an order 2 variance-covariance matrix. Let  $M_1, M_2, \ldots, M_T \in \mathbb{R}^{p \times p}$  be a collection of order-2 tensors. We define the order-3 tensor  $\mathcal{X} \in \mathbb{R}^{T \times p \times p}$  via concatenation along the first mode (tensor slices):

$$\mathcal{K}_c = \operatorname{con}(M_1, M_2, \dots, M_T), \quad \text{where } \mathcal{X}_{t,::} = M_t.$$
(A13)

By standard tensor algebra, an order-3 tensor  $\mathcal{X} \in \mathbb{R}^{T \times p \times p}$  can be reshaped or flattened into a higher-order tensor, under a specific indexing scheme. More generally, given a desired tensor order d, and assuming  $T = p^{d-2}$ , we define a transformation:

$$\mathcal{T}: \mathbb{R}^{T \times p \times p} \rightarrow \mathbb{R}^{p^d}, \; \kappa_d(\varepsilon_t) \in \mathbb{R}^{\underbrace{p \times \cdots \times p}_{d \; \text{times}}}$$

$$\kappa_2(\varepsilon_t) = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \cdots & \sigma_{pp} \end{bmatrix} \in \mathbb{R}^{p \times p}, \kappa_3(X_t) = \begin{bmatrix} \mathbf{0}_{p \times p} \\ \mathbf{0}_{p \times p} \\ \vdots \\ \mathbf{0}_{p \times p} \end{bmatrix} \in \mathbb{R}^{p \times p \times p}$$

that re-indexes the tensor slices  $M_t$  to fill the missing indices of an order d cumulant tensor. The replacing and re-indexing rule is illustrated in an order 3 tensor as a real plane, assuming d-1 is the maximal order such that  $\kappa_d=0$ .

Under this transformation, each slice  $M_t$  is interpreted as contributing to a specific mode configuration of the higher-order tensor. That is, the tensor  $\mathcal X$  is "lifted" into a d-way tensor by embedding each  $p \times p$  matrix slice as filling in the cumulant entries with fixed positions in the first d-2 indices corresponding to  $t \in \{1,\ldots,T\}$ , and varying the remaining two indices over  $p \times p$ . This leads to the same form as Eq. (A4) where all  $\kappa_d(\epsilon_t) \neq 0$ . Under Corollary B.1 and B.2, the new tensor has a unique decomposition of rank-1 tensor summation. To be specific, we assume  $v_i$  has no pair of columns to be collinear. This ensures the identification of  $F(I-\Phi)^{-1}$  and restricts F to be injective to only  $\operatorname{span}(H_i)$ .

The sequential steps are the same for non-Gaussian noise since the construction of the ideal  $\mathcal{I}^*$  associated with its variety is not influenced by  $\epsilon$  once FH is fixed up to a permutation and scaling.

Consequently, when Gaussianity is assumed, the distribution of  $O_t$  is fully 581 characterized by the first two cumulants. This property implies that the entire cumulant expansion, 582 and hence any higher-order dependency, collapses at second order. In this sense, the Gaussian 583 distribution is the unique fixed point of the cumulant hierarchy at order two. In temporal parametric 584 transition [28], a widely known condition to ensure the component-wise identifiability of the latent 585 process  $Z_t$  is to require that the driving noise  $\epsilon$  is not a fully isotropic Gaussian. That is, the Gaussian 586 noise distribution must shift under either intervention [13] or exhibit heterogeneity in its variance. 587 This is because all cumulants of order p > 2, which encode the exact causal dependencies, vanish for 588 Gaussian noise. As a result, sufficient variability can only arise from changes in the second-order 589 cumulant. Our results reflect that non-isotropic Gaussian noise is a necessary but not a sufficient 590 condition for full identifiability of the time-delayed generative model and parameters. 591

# **B.1.4** Identifying causal structure

Our proof is constructive: with the minimal hierarchy  $\mathcal{H}^{\mathcal{C}}$  by soft interventions such that no fixed value of entry  $(i \to j)$  in  $\Phi(w)$  induces a dependence removal, the transitive closure  $\bar{\mathcal{G}}$  of the ground truth process with its causal structure can be recovered up to the trivial transformation aforementioned. If a time process has no instantaneous influence it must have  $TC(\mathcal{G}) = \mathcal{G}$ , we can recover the original process and its causal structure up to a scaling and permutation  $\pi$ .

# C Proof of Supportive Results

# 599 C.1 Proof of Lem. 1

592

598

603

604

605

606

607

This lemma significantly constitutes the reasoning chains that lead to our identifiability results. We restate the original statement to cover more details and background in point process and theory of weak topology and convergence.

**Lemma C.1** (Bounding Point Process in intensity, constructive). For a measurable mapping  $N^{\Delta}$ :  $(\Omega, \mathcal{F}) \to (M_p, \mathcal{M})$  such that  $\omega \mapsto N(\omega)$  is a point process at scale  $\Delta$ . Let  $\Delta$  be the control operator for any subsequence of its point process. Consider  $A \in \mathcal{B}$  generated by the topology  $\mathcal{M}_p := \mathcal{B}(M_p)$ .  $\lambda$  and  $\epsilon$  is defined on this metric space. If  $\lambda$  satisfies the stationary increment condition, then we can establish the weak convergence of the constructed equivalent class:

$$\sum_{k: k\Delta \in A} \lambda_k^{(\Delta)} + \epsilon_k^{(\Delta)} \overset{w}{\Rightarrow} N(A) \ \ \textit{for} \ \lambda_k^{\Delta} = \lim_{\Delta \to 0} \lim_{\delta \to 0} \frac{\mathbb{E}[dN^\Delta | \mathcal{F}]}{\delta}$$

Proof. We start the proof with a trivial case. If  $\delta = \Delta_1 = 1$ , the condition always trivially holds. In this case, we only need to show  $N_t^{(\Delta_1)} = \lambda_t^{(\Delta_1)} + R_t$  by simply using the tower rule. Therefore, our proof gives more attention to the non-trivial case for  $\delta \neq 1$ .

611 *Case* 2:  $\delta \in (0, \Delta_1)$ 

The reasoning of this case becomes more complicated if the time step operator used for generating subsequences proportionally shrinks to a sufficiently small unit in  $(0, \Delta_1)$ . We rewrite the approximating sequence N to leverage the metricizability of the space. Since we work in a Polish space, the Borel  $\delta$ -algebra is countably generated and the space is separable and metrizable. Given a measurable set  $A \in \mathcal{B}$ , and a metric  $\rho$ , define the open  $\delta$ -neighborhood as:

$$\mathcal{A} = A^{\delta} := \{ x \in \mathbb{R}^d : \rho(x, A) < \delta \}$$

By outer regularity of Borel probability measures on Polish spaces, for every  $\epsilon>0$ , there exists a countable collection of open sets  $\{A_i\}_{i\in\mathbb{N}}$  such that  $\bigcup_i A_i \supset \mathcal{A}$  and  $\sum_i \mu(A_i\setminus A) < \epsilon$ . This allows us to approximate any compact subset from outside using open sets with arbitrarily small excess mass and ensures the approximating sequence is defined on a non-decreasing base. We paraphrase the convergence as

$$\sum_{k:k\Delta \in A} \lim_{i \to \infty} \frac{\mathbb{E}[N^{\Delta}(A_i)|\mathcal{F}]}{|A_i|} + \epsilon_k^{(\Delta)} \stackrel{w}{\Rightarrow} N(A) \quad \text{for } \Delta \to 0$$
 (A1)

The equation above is adapted from the continuous-time intensity for point processes. However, it requires us to work with two limit conditions for  $A_i$  with the 1/k closed ball shrinking to zero

measure and for the subsequence operator  $\Delta$  approaching to 0. A common method is to ensure dominated and uniform convergence of the limit. To harness information regarding the intensity in our convergence to a more generalized process, we work with only  $\Delta$  to induce the same time scale of intensity function. Therefore, we have the equivalent condition

$$\sum_{k:k\Delta\in A} \frac{\mathbb{E}\left[\sum_{k=1}^{\infty} Z_k^{\Delta} - \sum_{k=1}^{\infty} Z_{k-1}^{\Delta}|\mathcal{F}\right]}{|A_i|} + \epsilon_k^{(\Delta)} = \sum_{k:k\Delta\in A} \frac{\mathbb{E}\left[Z^{\Delta}(\Delta)|\mathcal{F}\right]}{\Delta} + \epsilon_k^{(\Delta)} \stackrel{w}{\Rightarrow} N(A) \quad \text{for } \Delta \to 0$$
(A2)

We remove the limit condition as it is clear that  $|A_i|$  is of measure zero when  $\Delta=0$ , which ensures the alignment between our topological property and plausibility to analyze only subsequences in the sequel. According to Lemma.2 by [20], for any compact interval  $[a,b]^{(\delta)}$  with the number of bins  $[b-a]/\delta$ ,  $\mathbb{E}[N^{(\delta)}([a,b])] < (b-a+2)(I-G^{(\delta)}(a,b))^{-1}\Lambda$  where  $G(a,b)=\int_a^b \Phi(s)ds$  is a solution of the stochastic differential equation systems

$$\mathbb{E}[\lambda([a,b])] = \mathbb{E}[u + G(a,b)\Lambda], \text{ for } \mathbb{E}[\lambda(a,b)] = \Lambda$$

Note that, by reapplying tower rule, Eq. (A1) implies:

$$\lim_{\delta \to 0} \mathbb{E}[N_A^{(\delta \in (0,1))}] \to \lim_{\delta \to 0} \mathbb{E}[\frac{\mathbb{E}[N_A^{(\delta)}|\mathcal{F}_t]}{\delta}]$$

Next, we show the necessity of tightness of the corresponding probability measure  $\mathbb{P}^{\Delta}$  for the left-hand of Eq. (A2) to achieve the desired convergence. Without loss of generality, we consider a nonparametric intensity function  $\lambda_t = \psi(u + \int \phi(t-s)Z^{\Delta}(s) \, ds)$ . Consequently,  $\mathbb{E}[\lambda_t] = \Lambda$  and  $\mathbb{E}[\lambda_t] = \mathbb{E}[\psi(u + \int \phi(t-s)Z^{\Delta}(s) \, ds)]$ . Assume that  $\psi$  is  $\alpha$ -Lipschitz and  $\alpha \|\phi\|_1 < 1$  [29], so the mapping  $F(\Lambda) = \psi(u + \|\phi\|_1 \Lambda)$  is a contraction on  $\mathbb{R}_+$ . By Banach's fixed-point theorem, there exists a unique solution  $\Lambda^*$  to the equation:

$$\Lambda^{\star} = \psi(u + \|\phi\|_1 \Lambda^{\star})$$

Formally, this can be rearranged as:

$$\psi^{-1}(\Lambda^*) - \|\phi\|_1 \Lambda^* = u \implies \Lambda^* = (\mathrm{id} - \|\phi\|_1 \cdot \psi^{-1})^{-1} (-u)$$

provided that id  $-\|\phi\|_1 \cdot \psi^{-1}$  is invertible on the image of  $\psi$ .

To control the tail probability, we apply Markov's inequality:

$$\mathbb{P}\left(\sum_{k: k, \Delta \in A} \frac{\mathbb{E}[Z^{\Delta}(\Delta)|\mathcal{F}]}{\Delta} + \epsilon_k^{(\Delta)} > M_{\varepsilon}\right) \leq \frac{\mathbb{E}[\sum_k \Lambda^{\Delta}]}{M_{\varepsilon}} \leq \frac{(b - a + 2\delta) \cdot \Lambda^{\star}}{M_{\varepsilon}}$$

643 Here, we define:

$$M_\varepsilon := \frac{(b-a+2\delta)\cdot \Lambda^\star}{\varepsilon} \quad \text{where } \Lambda^\star = \psi(u+\|\phi\|_1\Lambda^\star)$$

This choice ensures the upper bound remains within the prescribed  $\varepsilon$ -level for all  $\Delta \in (0, \Delta_1)$ . Since the only thing we need is the precompactness, we will not establish any tighter bound. Tightness of measure indicates we can always find a subsequence  $\lambda_{k_n}^{\Delta} + \epsilon_{k_n}^{\Delta}$  in  $\lambda_k^{\Delta} + \epsilon_k^{\Delta}$  converges weakly to a sequence  $\lambda^* + \epsilon^*$ . This weak convergence of subsequences, however, cannot control the limit uniqueness for each sequence. Therefore, we also should further control the limiting behavior of each sequence by uniform convergence of the characteristic functional defined by the approximating process and the target process.

# **D Detailed MUTATE Configuration**

# 652 D.1 Simulation Regime

We simulate multivariate point processes and their converging equivalent class  $Z_t$  extensively studied in our identifiability theory. We sample all point processes using the Poisson Superposition method

(rejection sampling from the upper bound of conditional intensity [30, 31]) in order to mimic highly 655 dynamic changes in conditional intensity, to capture denser information contained in stochastic 656 processes. Then we create corresponding converging classes as a proof-of-concept validation: A total 657 of 20,000 latent trajectories are sampled for each of the five kernel functions—exponential, power-law, 658 rectangular, simple nonlinear, and flexible mixing—under two noise regimes: heterogeneous noise 659 and Gaussian mixture noise. To illustrate the latent events underlying the unstructured data, we also 660 simulate stochastic dynamics for biological data using SERGIO [32], a GRN-guided gene expression 661 simulator used in Lorch et al.'s [15] causal modeling as well. All observation  $O_t$  is obtained from 662 latents  $Z_t$  through MLP and LeakyReLU nonlinearity mixing. 663

We demonstrate the generative process for INAR equivalent classes. For a fair comparison to those baselines mainly addressing step-wise conditional independence, we generate for both timestep dynamics and denser dynamics by changing the setup to very short kernel effects with  $\tau \in (0.001, 0.01) = t - t'$ . We generate stochastic point processes from three basic kernel response functions:

$$\begin{split} \phi_{\text{exponential}}(t) &= \alpha e^{-\beta t'}, \alpha \sim \text{uniform}(0.1, 0.5) \text{ and } \beta \sim \text{uniform}[0.5, 2) \\ \phi_{\text{powerlaw}}(t) &= \frac{\alpha}{(t+c)^{\beta}} \cdot \mathbf{1}, \alpha \sim \text{uniform}(0.5, 1.2), \beta \sim \text{uniform}[0.1, 0.8) \text{ and } \gamma \in \text{uniform}(1, 3, 1.8) \\ \phi_{\text{rectangular}}(t) &= \frac{1}{T-T'} \cdot \mathbf{1}_{\{t' \leq T\}} \end{split}$$

The baseline intensity  $u_0$  is sampled from uniform(0,1,0.2). All parameters of the basic kernel are uniformly sampled by ensuring  $\alpha < \beta$  in exponential responses,  $\alpha < \gamma$  in power-law response, respectively, to satisfy the stationary increment condition such that  $|\phi| < 1$ . In the simulation, we also consider two extreme cases for simple nonlinear intensity and nonparametric intensity. We construct the conditional intensity function by mixing latent features through a linear transformation followed by a non-linear activation. Specifically, we first compute a log-linear intensity using the expression

$$\lambda_t = \log(1 + \exp(z_\ell - r_\ell[:, \Delta, :]))$$

that ensures positivity and controls the scale of the output through a smoothed ReLU (i.e., softplus). In an alternative setting (kernel == "np"), we learn the intensity function using a small neural network (MLP): a two-layer perceptron with ReLU activation, ending in a Softplus to maintain positive outputs. This setup enables flexible, data-driven modeling of intensity dynamics beyond purely additive or linear forms. We define the mixing intensity function using a two-layer feedforward neural network with ReLU and Softplus activations. Formally, the architecture is given by:

$$\lambda_t = \sigma_+ \left( W_2 \cdot \text{ReLU}(W_1 \lambda_t(l) + b_1) + b_2 \right), \tag{A1}$$

681 where

682

683

684

685

693

- $\lambda_t(l) \in \mathbb{R}^d$  is the input linear basic intensity at time t,
- $W_1 \in \mathbb{R}^{64 \times d}$ ,  $b_1 \in \mathbb{R}^{64}$  are the weights and bias of the first layer,
- $W_2 \in \mathbb{R}^{d \times 64}, \ b_2 \in \mathbb{R}^d$  are the weights and bias of the second layer,
  - $\sigma_+(x) := \log(1 + e^x)$  denotes the Soft-plus activation.

This design ensures the output  $\lambda_t$  remains strictly positive and can model complex dependencies in the latent dynamics while maintaining numerical stability.

We model the transformation from the latent variable  $Z_t \in \mathbb{R}^d$  to the observational space via a multi-layer mixing network. Specifically, for each layer  $l=1,\ldots,L-1$ , the transformation is given by  $Z_t^{(l)} = \mathbf{A}^{(l)} \cdot \sigma_{\text{leaky}}(Z_t^{(l-1)})$ , where  $\mathbf{A}^{(l)} \in \mathbb{R}^{d \times d}$  is an orthogonal mixing matrix and  $\sigma_{\text{leaky}}$  denotes the leaky ReLU activation with slope  $\alpha=0.2$ . The initial input is  $Z_t^{(0)}=Z_t$ , and the final output  $Z_t^{(L-1)}$  represents the observation-space signal.

# D.2 Prior decomposition of time-adaptive module

Without loss of generality, we consider non-finite steps for a latent stochastic generative process, as discussed in Lem. 1, where  $\Delta t \to 0$ . This induces an equivalence that the intrinsic history—the filtration  $\mathcal{F}_t := \sigma\left(\bigcup_{0 < t < T} \sigma(Z_t^\Delta)\right)$ —ensures that the process  $Z_t^{(\Delta)}$  is  $\mathcal{F}_t$ -adaptive and measurable.

697 We decompose the ELBO objective as follows:

$$\begin{split} \text{ELBO} &= \log p(O) - D_{\text{KL}}(q_{\phi}(Z|O) \| p(Z)) \\ &= \mathbb{E}_{z \sim q(Z_t|O_t)}[\log p(O_t|Z_t)] + \mathbb{E}_{z \sim q(Z_t|O_t)}\left[\log \frac{q(Z_t|O_t)}{p(Z_t)}\right] \\ &= \mathbb{E}_{z \sim q(Z_t|O_t)}[\log p(O_t|Z_t)] - \mathbb{E}_{z \sim q(Z_t|O_t)}\left[\log q(Z_t|O_t) - \log p(Z_t)\right] \\ &= \mathbb{E}_{z \sim q(Z_t|O_t)}\left[\log p(O_t|Z_t) - \log q(Z_t|O_t)\right] + \mathbb{E}_{z \sim q(Z_t|O_t)}\left[\log p(Z_t)\right] \\ &= \mathbb{E}_{z \sim q(Z_t|O_t)}\left[\log p(O_t|Z_t) - \log q(Z_t|O_t)\right] + \mathbb{E}_{z \sim q(Z_t|O_t)}\left[\sum_{\mathcal{F}_0^+} \log p(Z_t^{(\Delta)}|\mathcal{F}_t)\right] \end{split}$$

The reason we can segment the increasing filtration in the last term is due to the nice property of  $\mathcal{F}_t$ -measurable sequence. We can show that filtration of  $Z_t|Z_s,R_t$  and  $Z_t|R_s$  is equal because it is well known that any p-order INAR sequence with stationary increments admits a moving average (MA) representation. Further construction of their filtration  $\tilde{\mathcal{F}}_t(\text{resp}.R_{s< t})$  and  $\mathcal{F}_t(\text{resp}.Z_{s< t},R_t)$  can show

$$\tilde{\mathcal{F}}_t = \mathcal{F}_t$$

We prove the result in the sequel. For  $\tilde{\mathcal{F}}_t, Z_t$  is a measurable function for s < t. By causality of the convolution kernel  $\Psi = (I - \Phi)^{-1}$  satisfying  $\Psi_\tau = 0$  for  $\tau < 0$ , which indicates  $Z_t \in \sigma(R_s: s < t)$ . Then, we construct another filtration  $\tilde{\mathcal{F}}_s: \sigma(R_u: u < s)$ . By adaptivity  $\tilde{\mathcal{F}}_s: \sigma(R_u: u < s) \subseteq \tilde{\mathcal{F}}_t: \sigma(R_u: u < t)$ . Therefore,  $Z_s$  is also  $\sigma(R_u: u < t)$ -measurable. Since the minimal  $\sigma$ -algebra of the original  $\mathcal{F}_t$ -measurable function must be contained in its  $\sigma$ -algebra, we have  $\sigma(Z_s) \subseteq \sigma(R_u: u < t)$  and  $\sigma(\bigcup_{s \le t} \sigma(Z_s)) \subseteq \sigma(R_u: u < t)$ . For  $\mathcal{F}_t, R_t = Z_t - \Psi \star Z_t$  so  $R_t$  is  $\sigma(Z_s: s \le t)$ -measurable.

Therefore, by a similar construction, it is evident that  $\sigma(\bigcup_{s < t} \sigma(R_s)) \subseteq \sigma(Z_s : s \le t)$ . Therefore,

710 because  $\tilde{\mathcal{F}}_t \subseteq \mathcal{F}_t$  and  $\mathcal{F}_t \subseteq \tilde{\mathcal{F}}_t$ , there must be  $\tilde{\mathcal{F}}_t = \mathcal{F}_t$ .

Following this set-up, the prior becomes:

$$Z_t \mid \mathcal{F}_t \sim \mathcal{N}\left(\begin{bmatrix} u_1(t) \\ u_2(t) \\ \vdots \\ u_p(t) \end{bmatrix} \sum_{t' < t} (I - \Phi)^{-1}, \sum_{t' < t} (I - \Phi)^{-1} \Sigma_{t'} (I - \Phi)^{-T} \right)$$

The latents are generated by  $Z_t = (I - \Phi) \star R_t$ , where  $R_t$  is modeled as isotropic Gaussian noise with mean U and variance  $\Sigma$ . Note that the variance matrix  $\Sigma_{Z_t}$  is zero for any  $t - t' \neq 0$ . By the Wiener-Khinchin Theorem, we have the covariance matrix  $C_{Z_t}(0) = \frac{1}{N} \sum_{k=0}^{N-1} S_z(w_k)$ , we drop the sub-index  $Z_t$  whenever no confusion is caused. Now we can derive the decomposition of the convolution prior as

$$\mathbb{E}_{z \sim q(Z_t|O_t)} \left\{ \sum_{\mathcal{F}_0^+, Z_t}^{\mathcal{F}_T} \log p(Z_t^{(\Delta)}|\mathcal{F}_t) \right\} \\
= \mathbb{E}_{z \sim q(Z_t|O_t)} \left\{ \sum_{\mathcal{F}_0^+, Z_t}^{\mathcal{F}_T} \log p \left[ (I - \Phi_t) \star \hat{R}_t^{(\Delta)} \right] \right\} \\
= \mathbb{E}_{z \sim q(Z_t|O_t)} \left\{ \sum_{\mathcal{F}_0^+, Z_t}^{\mathcal{F}_T} \log p \left[ (I - \Phi_t) \star \hat{R}_t^{(\Delta)} \right] \right\} \\
= \mathbb{E}_{z \sim q(Z_t|O_t)} \left\{ \sum_{\mathcal{F}_0^+, Z_t}^{\mathcal{F}_T} \log p \left[ \mathcal{N}(\hat{U}_R, \sum H_t^{-1} \sum_{\hat{R}_t'} H_t^{-T}) \right] \right\}$$

$$= \mathbb{E}_{z \sim q(Z_{t}|O_{t})} \left\{ \sum_{\mathcal{F}_{0}^{+}, Z_{t}}^{\mathcal{F}_{T}} \log p \left[ \mathcal{N}(\hat{U}_{R}, \sum_{t} H_{t}^{-1}(PSD_{\hat{Z}_{t}}) \Sigma_{\hat{R}_{t}'} H_{t}^{-T}(PSD_{\hat{Z}_{t}})) \right] \right\}$$

$$= \mathbb{E}_{z \sim q(Z_{t}|O_{t})} \left\{ \sum_{\mathcal{F}_{0}^{+}, Z_{t}}^{\mathcal{F}_{T}} \log p \left[ \mathcal{N}(\hat{U} \sum_{t' < t} \underbrace{(I - \Phi(t - t'))^{-1}}_{(1 - \Phi(t)(w_{k})^{-1} \Sigma(1 - \Phi)(w_{k})^{-H} = S_{Z}(w_{k})}, \frac{1}{N} \sum_{k=0}^{N-1} S_{Z_{t}}(w_{k})) \right] \right\}$$

$$= \mathbb{E}_{z \sim q(Z_{t}|O_{t})} \left\{ \sum_{\mathcal{F}_{0}^{+}, Z_{t}}^{\mathcal{F}_{T}} \log p \left[ \mathcal{N}(\hat{U} \sum_{\tau > 0, w = 0} (I - \Phi(\tau))^{-1} e^{-jw\tau}, \frac{1}{N} \sum_{k=0}^{N-1} S_{Z_{t}}(w_{k})) \right] \right\}$$

$$= \mathbb{E}_{z \sim q(Z_{t}|O_{t})} \left\{ \sum_{\mathcal{F}_{0}^{+}, Z_{t}, N \in (N_{0}, T)}^{\mathcal{F}_{T}} \log p \left[ \mathcal{N}(\hat{U} PSD_{Z_{t}}(H(0)), \frac{1}{N} \sum_{k=0}^{N-1} S_{Z_{t}}(w_{k})) \right] \right\}$$

$$(A3)$$

# 717 D.3 Explicit control for convolution prior

718

719

720

721

722

723

724

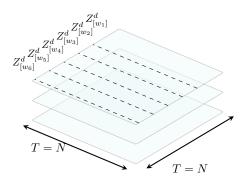


Figure A2: Visually Time-adaptive PSD Computation

**Encoder-PSD flow** As shown in Eq. (A3), a key component of our module is to efficiently compute the decomposition of PSD matrix. However, under milder regularity conditions, the PSD decomposition is not unique, thus only be recovered up to the minimal-phase. By the Theorem, the energy of the time domain and frequency domain is equivalent. Therefore, the encoded distribution is not sufficient to decompose the PSD matrix for which a reparameterization is needed. An encoder receives a T-length sequence  $O_t$  and returns the latent variable vector. Fast Fourier Transformation converts the latent sequence to a vector of equal length up to t:

$$[Z_{\mathcal{F}_0}, Z_{\mathcal{F}_t}, \cdots, Z_T] \Rightarrow \{[Z[f_0], Z[f_k], \cdots, Z[K]] | K = 0, 1, 2, \cdots, T\}$$

And the flow method is enforced by solving the following Wilson Factorization optimization problem for each  $[Z[f_0], Z[f_k], \cdots, Z[K]]$ , finding the transfer matrix

$$H^{\dagger} = \arg\min_{\Sigma_t = \sigma^2 I} PSD(Z_t) - H^{\dagger} \Sigma_t H^{\dagger H}$$

That is then sent to evaluate the true prior distribution, supporting the joint optimization of all loss components.

The summation of kernel products and integrated noise variables is guaranteed to converge to the true time-adaptive process under  $\mathcal{F}_t$ , provided that the time discretization is sufficiently dense. The

latent variable  $Z_t^{\Delta}$  is sampled from the encoder distribution  $q_{\phi}$  and passed to the PSD decomposition 731 module to compute the frequency-domain representation of the full kernel matrix  $F_w[1-\Phi_t]$  and 732 the power spectral density  $S_{R_*}$ . We further remark that the key step, spectrum decomposition, is 733 completed for the entire encoded trajectory  $\hat{Z}_{t_0:T}$  , and the prior structure is ensured by segmenting 734 filtration. This features the major difference in prior work that recursively constructs an equal-length 735 sliding window for each latent. Filtration segmentation can work with causal masks that a more 736 expressive encoder leverages. Note that transformer modules are not a required component for 737 shorter sequences, i.e. T < 100. However, when the sequence is extremely long, as simulated in the 738 conventional class of stochastic point processes, a transformer can be used in place of a common 739 MLP encoder to learn much more expressive latent embeddings by utilizing the filtration attention 740 from arbitrarily long past events.

Overall training loss. To encourage sparsity in transferring kernels, we follow the widely used 742 penalty to jointly optimize: 743

$$\mathcal{L}_{Total} = \mathcal{L}_{Recon} - \beta \mathcal{L}_{KLD} - \gamma |\Phi| - \omega \mathcal{L}_{PSD}$$
 (A4)

This training objective ensures the learned latent process is driven by a family of generalized white 744 processes, as, in the Encoder-PSD flow, the decomposition is enforced by the prescribed isotropic 745 noise, which omits any discriminator module as used in [28]. The coefficients in sparsity loss and 746 PSD accuracy are registered as tunable hyperparameters.

#### **D.4 Extended Results** 748

741

	Table A	2: Reporting ti	ic best peric	illiance for	cacii basciii	ic	
Method	Metric	Kernel Ave.	Exp	Power.	Rect.	Nonlin.	Nonpar.
TDRL	MCC	0.657	0.629	0.653	0.773	0.584	0.644
	$\mathcal{L}_{vae}$	0.449	0.308	0.302	0.302	0.871	0.461
BetaVAE	MCC	0.419	0.395	0.414	0.420	0.433	0.433
	$\mathcal{L}_{vae}$	9.480	8.538	7.533	8.424	11.683	11.220
SlowVAE	MCC	0.410	0.384	0.405	0.420	0.425	0.412
	$\mathcal{L}_{vae}$	362.890	395.107	448.105	452.472	238.520	280.247
PCL	MCC	0.440	0.469	0.379	0.430	0.474	0.449
	$\mathcal{L}_{vae}$ (train)	0.693	0.693	0.694	0.693	0.693	0.693
MUTATE	MCC	0.811	0.922	0.784	0.964	0.885	0.501
	$\mathcal{L}_{vae}$	0.670	0.448	0.508	0.253	0.942	1.201

Table A2: Reporting the best performance for each baseline

#### **Related Work** ${f E}$ 749

751

752

753

754

755

756

757

758

759

760

761

762

Causal disentanglement and learning time series. Although estimating and predicting time series is a classical problem in both traditional statistics and modern machine learning, representation learning has opened new avenues for leveraging latent information to better characterize time series data [33, 34]. Recently, learning causal representations in time series has become a foundational approach for enabling new scientific discoveries. This line of research primarily focuses on establishing identifiability of causal latent variables by exploiting nonstationary data [28, 4] and modular distribution shifts [5, 35] with sparsity constraints [36, 37] on the latent transition. Those works solve the identifiability problem of latent causal models by leveraging sufficient variability that can come from proper interventions or passive distribution shifts. Another line of research focuses on learning the underlying causal graph among latent variables

**Learning Causality in Stochastic Processes.** While learning causality remains a considerably more challenging task than causal discovery or representation learning, several efforts have been made to bridge these areas. Here we review existing approaches that link causal learning with stochastic modeling. Our scope is not limited to causal representation learning with stochastic processes, but extends to a broader set of problems that are closely related to either domain.

One representative direction in causal learning for dynamical systems is the study of Granger causality—a broader and looser notion compared to strictly structured causal models [38]. It is widely acknowledged that full causal recovery in such systems is impossible. Consequently, even the most recent work on stochastic processes can only determine whether a point process a is Granger-causal or non-causal with respect to another process b, typically formalized through *local independence* and the  $\delta$ -separation rule [39]. Another active line of work concerns identifiability in dynamical systems [40]. However, to the best of our knowledge, none of these models provides provable guarantees for highly dynamical systems such as self-exciting or more general stochastic processes.

Connections between causal representation and dynamical systems have also been explored through ordinary differential equations (ODEs) [41]. Technically, these approaches recover only a set of parameters that are difficult to interpret as causal in the latent space, or at best allow stochastic dynamics in the observed variables. More recently, causal diffusion models have been proposed [42, 15], yet they largely treat diffusion as a standard denoising process and thus do not permit a well-structured stochastic latent causal representation.

Another important line of research investigates interventions on stochastic processes and the corresponding post-intervention distributions, which serve as the basis for causal inference [43, 44, 45, 46, 15]. The first attempt to introduce a causal interpretation into stochastic differential equations (SDEs) was made by the authors of [47], where interventions are defined as the removal of single variables in SDEs. They showed that causal principles in SDEs can be formalized as interventions, with the resulting post-interventional distribution identifiable via the infinitesimal generator. However, such interventions are too restrictive to capture more complex dynamical scenarios. Following this initial line of work, [44] further develops methods for estimating stationary causal models by minimizing the deviation of stationarity of diffusion.

**Stochastic representation in biological science.** Prior work on the dynamics of cancer genomics has investigated modeling the underlying stochastic processes, typically under the assumption that all mutations can be identified through observable changes in protein binding and synthesis. A seminal line of studies focuses on methods and conditions under which mutation rates are treated as fixed for each driver mutation during tumor progression. Under these assumptions, the evolutionary dynamics can be effectively modeled using a linear Moran process with fixed population size [48, 49]. One branch of this literature aims to identify driver mutations that are directly manifested in observations. However, given the limited prior knowledge, mutation interactions are unlikely to be strictly linear or fixed. As a result, due to the inherent stochasticity and dynamical nature of cancer development, most driver mutations remain latent and their patterns are not readily discernible in protein sequences [50, 51, 52]. More recently, [52] reformulated this problem by introducing a framework to distinguish between weak and strong driver mutations to better characterize cancer progression. For example, in several cell cycles, a normal cell must accumulate multiple mutations in tumor-susceptibility genes to trigger oncogenesis. This process is inherently stochastic, and many mutational events may be dependent, self-exciting, or regulated by other processes. Identifying latent processes underlying disease-specific mutations and recovering their causal relationships is therefore crucial for computational biology and the planning of sequential cancer treatment regimens.

# **F** Conclusion and Limitation

Our paper makes extensions of causal representation learning framework to stochastic causal dynamics (i.e., multivariate Hawkes Processes), a topic not yet covered in current CRL literature. We propose a new perspective that a branch of stochastic processes can be viewed as the corresponding equivalent class through INAR representation and weak convergence. Under those conditions, we show that the latent stochastic process can be identified up to a component-wise transformation and a scaling permutation matrix. Our theoretical result bridges the gap between stochastic modeling and causal representation. We also propose a novel framework to learn the time-adaptive transition dynamics to accurately estimate the latent processes. However, our work avoids the most complicated case for a fully nonparametric kernel, which, most of the time, can be replaced with a simpler kernel. Future direction may include solving this condition and causal representation learning for stochastic differential processes that manifest in rich scientific questions.