

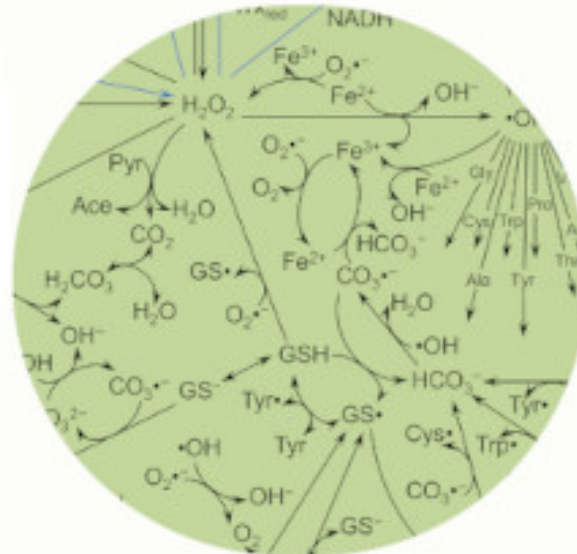
Discovering Causality in Event Time-Series

bionetgen/
bionetgen



msneddon/nfsim

A general-purpose, stochastic, biochemical
reaction simulator for large reaction networks



Pavel Loskot

pavelloskot@intl.zju.edu.cn



ZJU-UIUC INSTITUTE

Zhejiang University-University of Illinois at Urbana-Champaign Institute
浙江大学伊利诺伊大学厄巴纳香槟校区联合学院

The Seventh International Conference on Advances in Signal, Image and Video
Processing
SIGNAL 2022

May 22, 2022 to May 26, 2022 - Venice, Italy



ABOUT ME



Pavel Loskot joined the ZJU-UIUC Institute as Associate Professor in January 2021. He received PhD in Wireless Communications from University of Alberta, Canada, and MSc and BSc in Radioelectronics and Biomedical Electronics, respectively, from the Czech Technical University of Prague. He is Senior Member of the IEEE, Fellow of the HEA, and the Recognized Research Supervisor of the UKCGE in the UK.

In the past 25 years, he was involved in numerous industrial and academic collaborative projects in the Czech Republic, Finland, Canada, UK, Turkey, and China. These projects concerned mainly wireless and optical telecommunication networks, but also genetic circuits, air transport services, and renewable energy systems. This experience allowed him to truly understand the interdisciplinary workings, and crossing the disciplines boundaries.

His current research focuses on statistical signal processing, classical machine learning, and importing methods from Telecommunication Engineering and Computer Science to model and analyze systems more efficiently and with greater information power.

OBJECTIVES

- describe dynamic systems by their event history
- discover causality among events being categorical random variables
- implement automated causal event discovery for interpretable stochastic simulations of biochemical reaction networks

OUTLINE

- events in state-space model of dynamic systems
- causality as conditional probabilities
- modified matrix profile analysis
- simulation software and numerical examples

INTRODUCTION

Causality

- key in scientific reasoning: hypothesis discovery and testing
- key in engineering reasoning: prescriptive statistics
- statistical associations are neither necessary nor sufficient to infer causality

Literature

- discovering causality in time-series data
 - structural causal models (SCMs)
 - Granger and intervention causality, average causal effect (ACE)
- matrix profile analysis of time-series data
 - discovering motifs and discords, identifying changes in statistics

New problem

- causality in categorical time-series
 - event history of a dynamic system
- causality between events newly defined here as
 - nearly certain and uncertain conditional event sequences
- a versatile causal analysis framework, which can be fully automated
 - event sets and multi-sets and their distances and equivalences
 - modified matrix profile

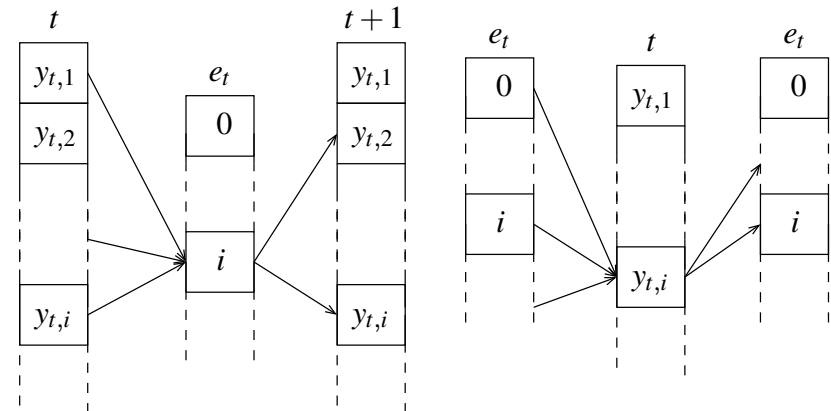
SYSTEM MODEL

Dynamic system

- state-space representation

$$\mathbf{z}_t = e_t(\mathbf{z}_{t-1}, \mathbf{z}_{t-2}, \dots)$$

$$\mathbf{y}_t = O(\mathbf{z}_t, \mathbf{z}_{t-1}, \dots)$$

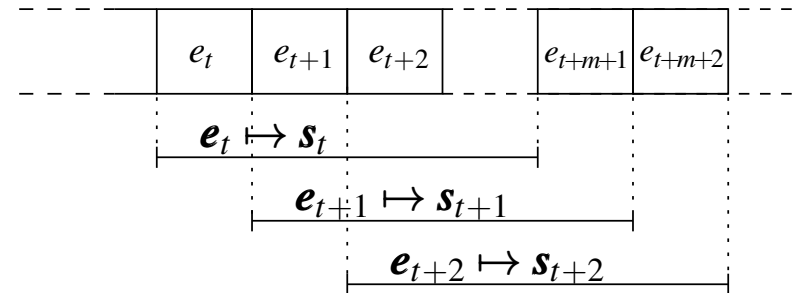


- assumptions

(1) system is memoryless:

$$\mathbf{z}_{t+1} = e_t(\mathbf{z}_t) \text{ and } \mathbf{y}_t = O(\mathbf{z}_t)$$

(2) observations are perfect: $\mathbf{y}_t = \mathbf{z}_t$



Theorem Given observations \mathbf{y}_t at time t , the ordering of $(m + 1) > 0$ events in sequence, $(e_t, e_{t+1}, \dots, e_{t+m})$, does not affect the observation \mathbf{y}_{t+m} at time $(t + m)$.

Consequently

- event ordering in shorter sequences is irrelevant
- event sequences can be converted to sets or multi-sets

ANALYSIS OF EVENT TIME-SERIES

Process

- sliding window partitioning to create sequences \mathbf{e}_t of N events each
- convert sequences \mathbf{e}_t to (multi-) sets \mathbf{s}_t
- compute empirical conditional probabilities $\Pr(\mathbf{e}_j|\mathbf{e}_i)$ or $\Pr(\mathbf{s}_j|\mathbf{s}_i)$
 → in general, $\Pr(\mathbf{e}_j|\mathbf{e}_i) \neq \Pr(\mathbf{e}_i|\mathbf{e}_j)$

Definition 1 The event sequences \mathbf{e}_i and \mathbf{e}_j , $j > i$, have a cause-effect relationship, provided that their conditional probability, $\Pr(\mathbf{e}_j|\mathbf{e}_i) \rightarrow 1$.

Definition 2 The event sequences \mathbf{e}_i and \mathbf{e}_j , $j > i$, have no cause-effect relationship, provided that $\Pr(\mathbf{e}_j|\mathbf{e}_i) \rightarrow 0$ and $\Pr(\mathbf{e}_i|\mathbf{e}_j) \rightarrow 0$.

Practical implementation

- split every $\mathbf{e}_t = \mathbf{e}_i \cup \mathbf{e}_j$, so that $\mathbf{e}_i \cap \mathbf{e}_j = \emptyset$, $|\mathbf{e}_i| = N_1$, $|\mathbf{e}_j| = N_2$, and $N_1 + N_2 = N$
 → referred to as left and right sub-sequences

Event equivalence The event (sub-) sequences \mathbf{e}_i and \mathbf{e}_j are said to be equivalent, provided that their distance, $d(\mathbf{e}_i, \mathbf{e}_j) = 0$.

→ enables to make $\Pr(\mathbf{e}_j|\mathbf{e}_i)$ closer to 1

ANALYSIS OF EVENT TIME-SERIES (2)

Distance between two event sequences

$$d(\mathbf{e}_i, \mathbf{e}_j) = d_0 - |\mathbf{s}_i \cup \mathbf{s}_j|$$

$$d(\mathbf{e}_i, \mathbf{e}_j) = d_0 - |\mathbf{s}_i \cap \mathbf{s}_j|$$

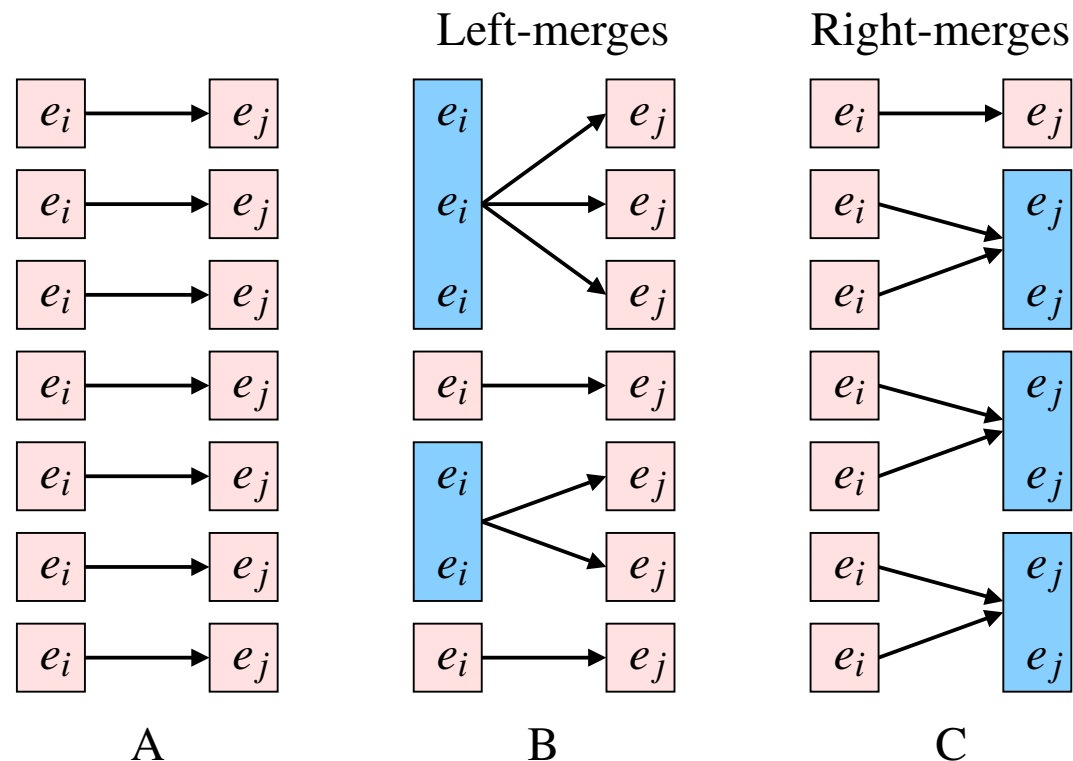
$$d(\mathbf{e}_i, \mathbf{e}_j) = d_0 - (|\mathbf{s}_i| + |\mathbf{s}_j|)$$

$$d(\mathbf{e}_i, \mathbf{e}_j) = d_0 - \max(|\mathbf{s}_i|, |\mathbf{s}_j|)$$

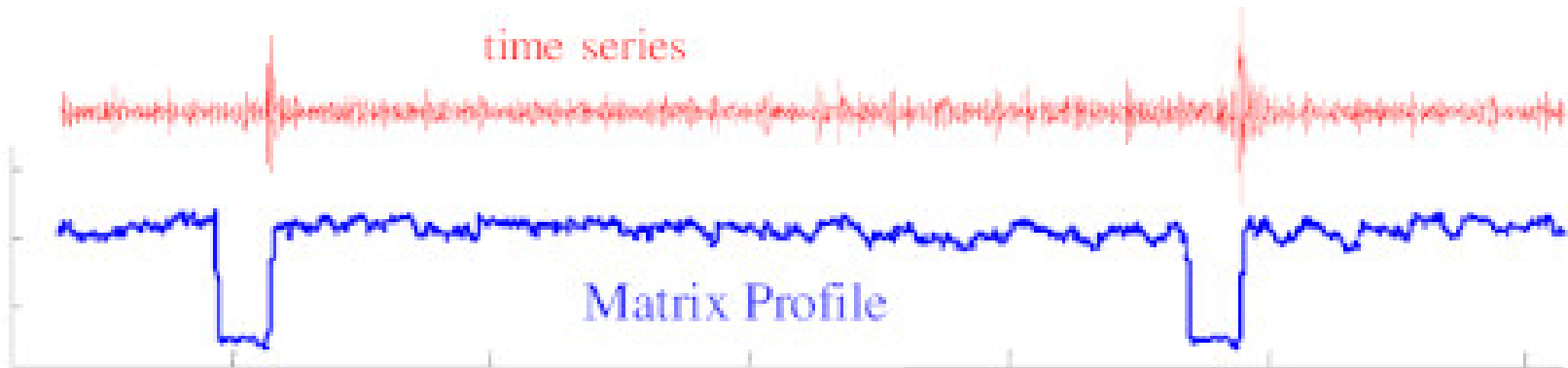
$$d(\mathbf{e}_i, \mathbf{e}_j) = \max(|\mathbf{s}_i|, |\mathbf{s}_j|) - \min(|\mathbf{s}_i|, |\mathbf{s}_j|)$$

$$d(\mathbf{e}_i, \mathbf{e}_j) = \min(|\mathbf{s}_i \setminus \mathbf{s}_j|, |\mathbf{s}_j \setminus \mathbf{s}_i|)$$

Merging left/right event sub-sequences



MATRIX PROFILE ANALYSIS



Canonical matrix profile

- given sliding window sub-sequences of length m , shows the minimum distance from any other distinct sub-sequence
→ the profile strongly affected by actual choice of m
- efficient implementations exist in different languages
→ e.g. stumpy in Python

Modified matrix profile

- for categorical event sequences, already defined various distance metrics
→ can be generalized to any other pairwise metrics
- can consider both the minimum/maximum distance as well as its multiplicity

NUMERICAL EXAMPLES

Biochemical reaction network

- a dynamic system described by Chemical Master Equation (CME)
- stochastic simulations are often used to solve CME
 - time evolution of chemical species counts
- chemical reactions (i.e., events) are occurring sequentially at random
- the goal is to analyze the event history (time-series)
 - it is readily available in simulations by recording all reaction events
 - this is an alternative approach to analyzing chemical species counts

Procedure

- select a specific biochemical system to simulate
- obtain the SBML model of this system and simulate it
 - BioNetGen and NFsim open source software
- record and process the reaction events
 - Python and Bash scripts
- the whole procedure can be largely automated
 - generated 160 diverse plots for 12 experiments across 9 biochemical models in less than 1 hour with minimum manual intervention

NUMERICAL EXAMPLES (2)

Simulated biochemical model

- antigen receptor signaling regulating the activity and fate of the B-cells
- 32 molecule types, 158 reaction rules, 129 model parameters
- full extracted model: 1,124 chemical species and 24,390 chemical reactions
- 100 s of simulation time resulted in 3,634,390 reaction events
→ can be divided in blocks, e.g. once per 1 s

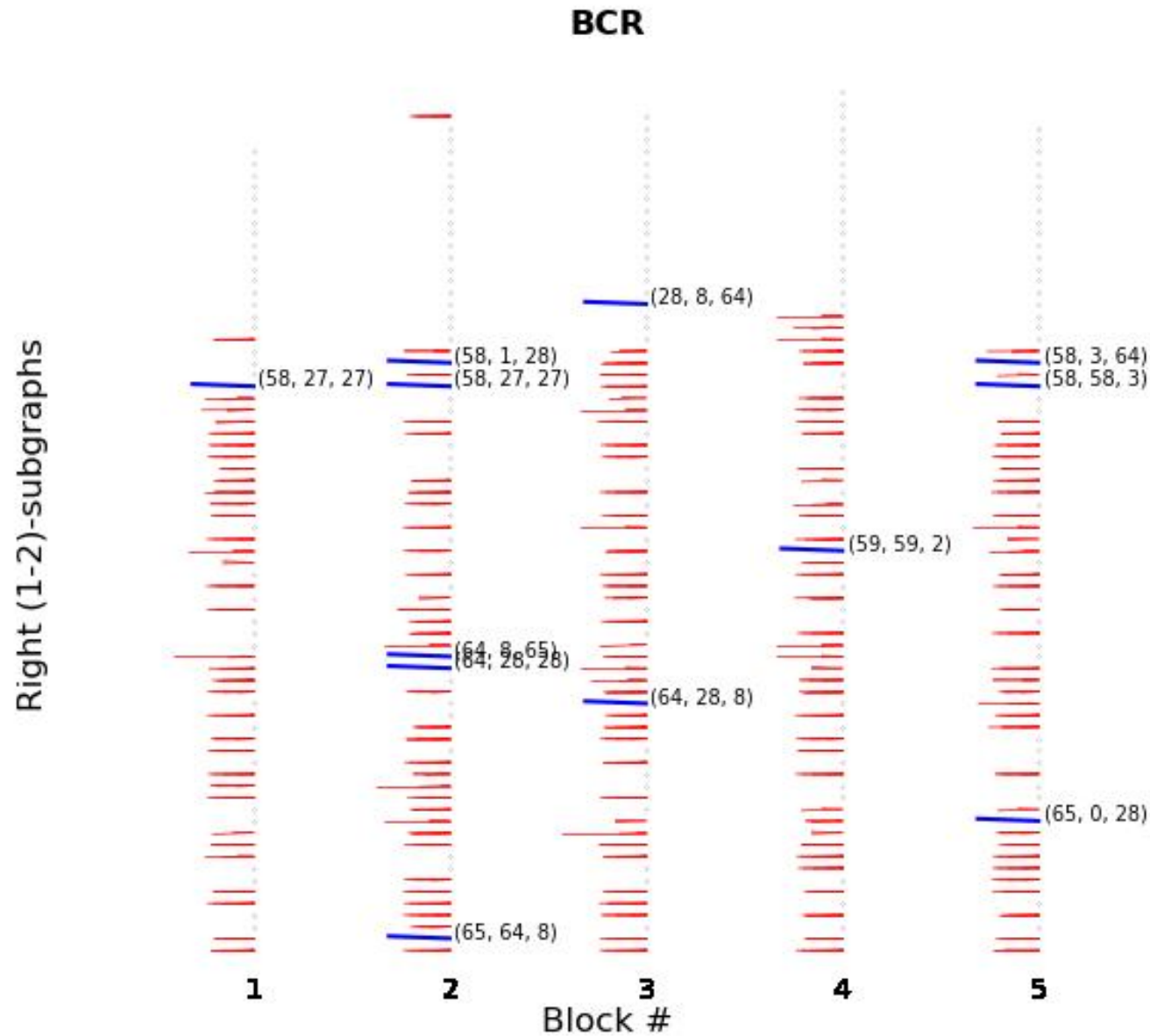
NFsim modifications

- upgraded source code last updated in 2011 to more recent C++ compiler
→ resolved compiling errors, removed most warnings
- added event recording functionality
→ written in text or binary file

BioNetGen

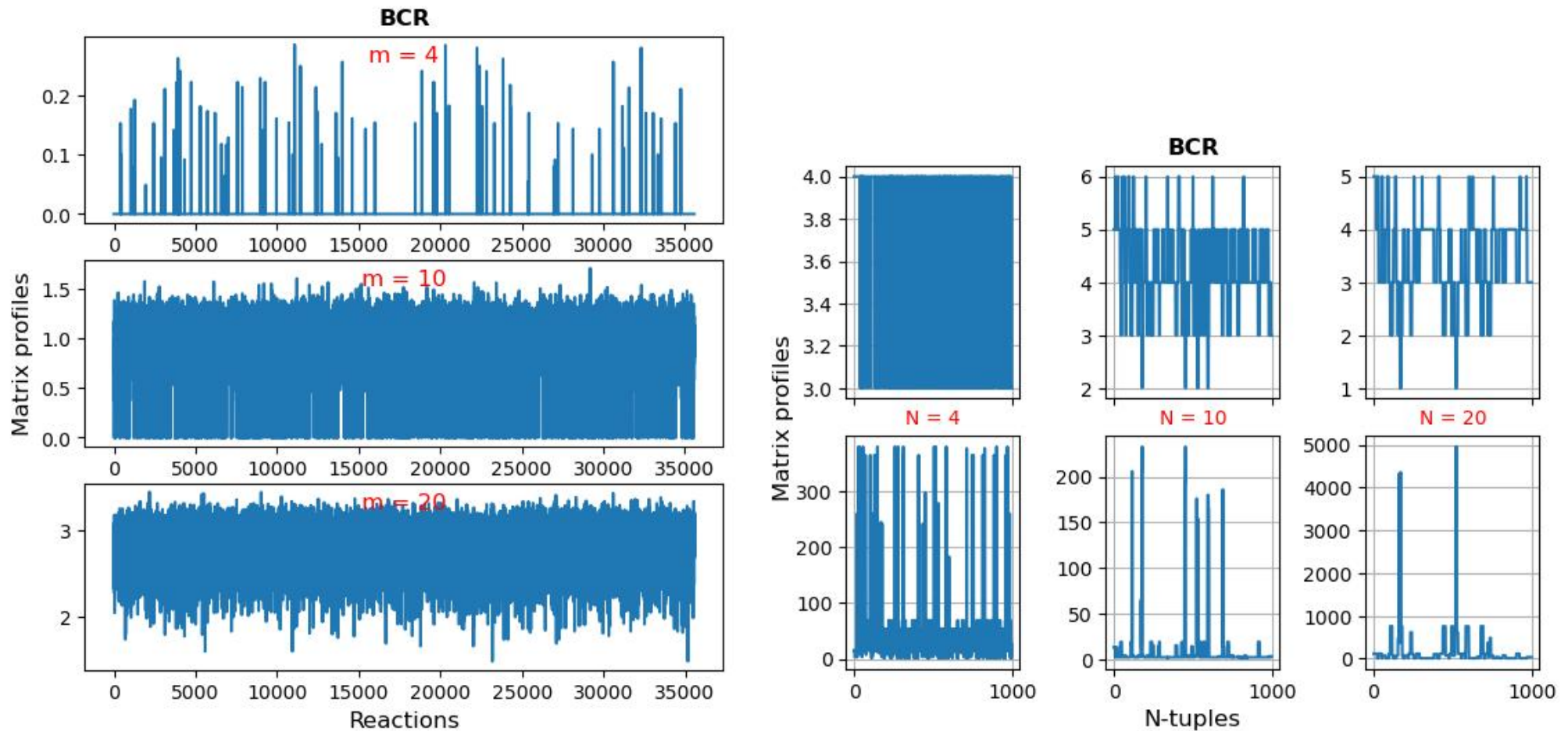
- biochemical model described in human-friendly bngl text file
- a Perl script in BioNetGen processes bngl file into SBML file
- SBML, the model description in XML, is simulated in NFsim
→ BioNetGen acts as a user-friendly interface for NFsim

NUMERICAL EXAMPLES (3)



- $N = 3$ with $|\mathbf{e}_j| = 2$ and $|\mathbf{e}_i| = 1$, and assuming $\Pr(\mathbf{e}_i | \mathbf{e}_j) \equiv P_{j,i}$
- Black: $P_{j,i} < 0.1$, Red: $P_{j,i} < 0.9$, Blue: $P_{j,i} \geq 0.9$

NUMERICAL EXAMPLES (4)



- Left: canonical matrix profile with minimum distances $d(\mathbf{e}_i, \mathbf{e}_j) = d_0 - |\mathbf{s}_i \cup \mathbf{s}_j|$
- Right: modified matrix profile with maximum distances $d(\mathbf{e}_i, \mathbf{e}_j) = \min(|\mathbf{s}_i \setminus \mathbf{s}_j|, |\mathbf{s}_j \setminus \mathbf{s}_i|)$

CONCLUSION

Key points

- behavior of dynamic systems accurately described by the events history
- ordering of events locally unimportant, can assume (multi-) sets instead
- causality defined here as conditional nearly certain or uncertain events
- distance metrics between event sequences used for matrix profile analysis

Future work

- causality of sub-sampled event sequences
- causality of longer event sequences
- automated interpretation of discovered causally related events

Thank you!

pavelloskot@intl.zju.edu.cn