naoki88@sanken.osaka-u.ac.jp
yasuko@sanken.osaka-u.ac.jp

Ren Fujiwara
SANKEN, Osaka University
r-fujiwr88@sanken.osaka-u.ac.jp

Yasushi Sakurai
SANKEN, Osaka University
yasushi@sanken.osaka-u.ac.jp

## Motivation - Given: Multivariate data streams

e.g., Spread of infectious diseases, coronavirus (COVID-19)

### Challenges

How can we discover time-changing causal relationships?

(a) Snapshot takes on January 8, 2021

**Given:** Multivariate Data stream

i.e., $X = \{x(1), ..., x(t_c), ...\}$

(b) Snapshot taken on May 19, 2022

**Goal:** Achieve all of the followings
- **Find** distinct dynamical patterns / **regimes**
- **Discover** causal relationships, which changes over time / **time-evolving causality**
- **Forecast** an $l_s$-steps ahead future values

✨ **ModePlait: novel streamimg method**

## Proposed Model - ModePlait

**Key Concepts** - Our model is designed based on SEM

Exogenous variables evolve over time / **inherent signals**



Hankel matrix, Self-dynamics factor set $\mathcal{D}_{(i)}$

### Main idea (P1): Lat...

$e_{(i)}$ (inherent signals E), Mixing matrix

Augmentaion

$e_{(i)}(t) = g...$

$i$-th univariate inherent signal

Time-de... embedd...

$g(e_{(i)}(t)) := (e_{(i)}(t), e_{(i)}(t-1), ..., e_{(i)}(t-h...$

$\mathcal{D}_{(i)} = \{\Phi_{(i)}, \Lambda_{(i)}\}$ / **self-dyna...**

### Main idea (P2): Dynamic...

Describe distinct dynamical pa...

$$s_{(i)}(t+1) = \Lambda_{(i)} s_{(i)}(t) \quad (1 \leq i < d)$$
$$e_{(i)}(t) = g^{-1}(\Phi_{(i)} s_{(i)}(t)) \quad (1 \leq i \leq d)$$
$$v(t) = W^{-1} e(t)$$

Estimated value, Mixing matrix

(P1) Collection of $d$ self-dynamic factor sets

$\theta = \{W, \mathcal{D}_{(1)}, ..., \mathcal{D}_{(d)}\}$ / **regime**, $\Theta = \{\theta^1, ..., \theta^R\}$ / **regime set**

$\mathcal{B} = \{B^1, ..., B^R\}$ / **time-evolving causality**

## Optimization algorithm

**Given:**
- Multivariate data Stream $X$

### Estimate:
- Full parameter set
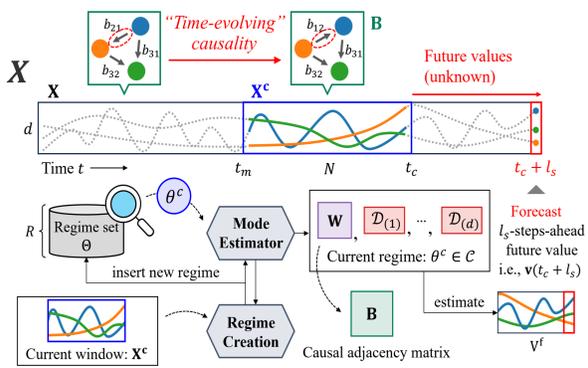  $\mathcal{F} = \{\Theta, \Omega\}$, $\Omega$: update param
- Model candidate
  $\mathcal{C} = \{\theta^c, \omega^c, S_{en}^c\}$
- Time evolving causality
  $\mathcal{B} = \{B^1, ..., B^R\}$, $R$: # of regimes
- $l_s$-steps ahead future value $v(t_c + l_s)$, $t_c$: current time point
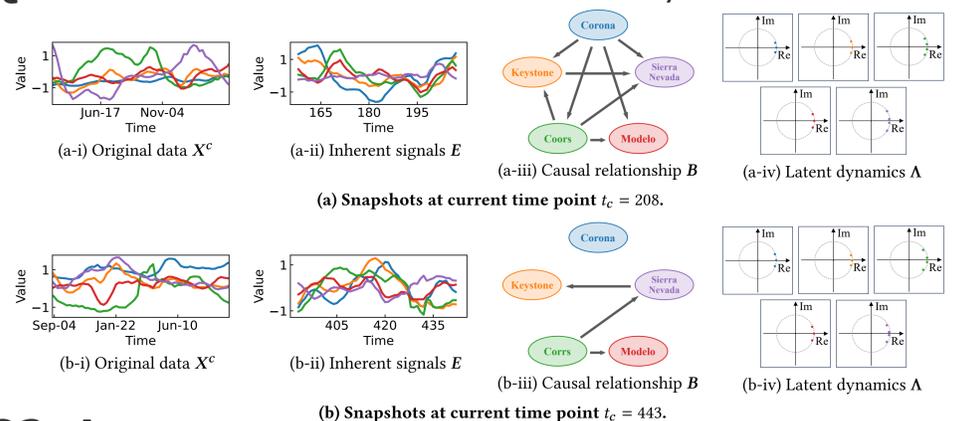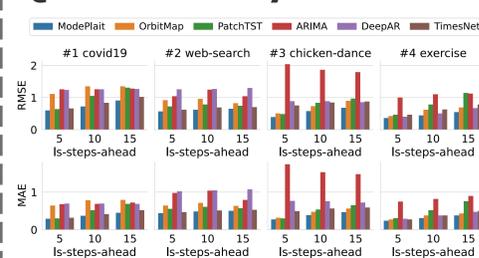
## Experiments - Answer the essential questions

**Datasets**: we used the following four real datasets

*(#1) covid19*: was obtained from Google COVID-19 Open Data [9].

*(#2) web-search*: consists of web-search counts on Google [10].

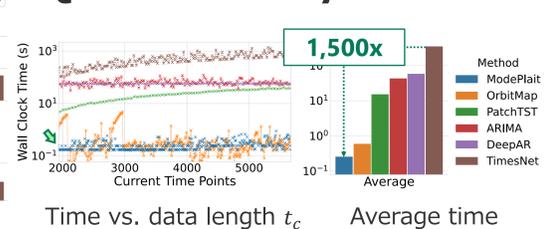*(#3) chicken-dance, (#4) exercise*: were obtained from the CMU motion capture database [4].

### Q1. Effectiveness - web-click activity stream

Interpretability of modes

### Q2. Accuracy

### Q3. Scalability

**1,500x**

Time vs. data length $t_c$    Average time

Method: ModePlait, OrbitMap, PatchTST, ARIMA, DeepAR, TimesNet

**performs its competitors**

Future values (unknown)

Forecast $l_s$-steps-ahead future value i.e., $v(t_c + l_s)$

Causal adjacency matrix

**Other downstream tasks:** we detect anomalous activities based on time-evolving causality in data streams.