# Automatic Sequential Pattern Mining in Data Streams

Koki Kawabata*, Yasuko Matsubara & Yasushi Sakurai

ISIR-AIRC, Osaka University
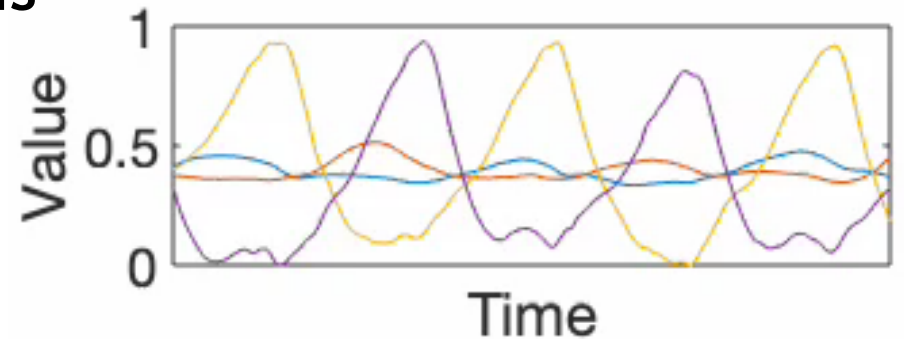
# Motivation

**Given:** time-evolving data streams
- e.g., IoT sensors/Web click logs
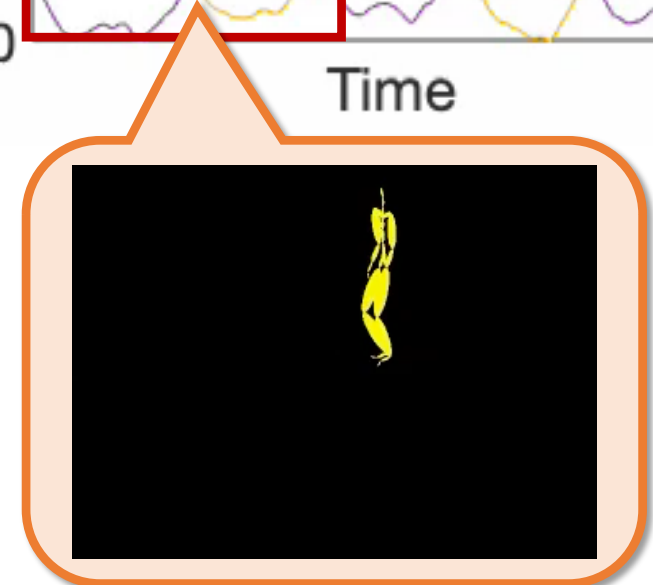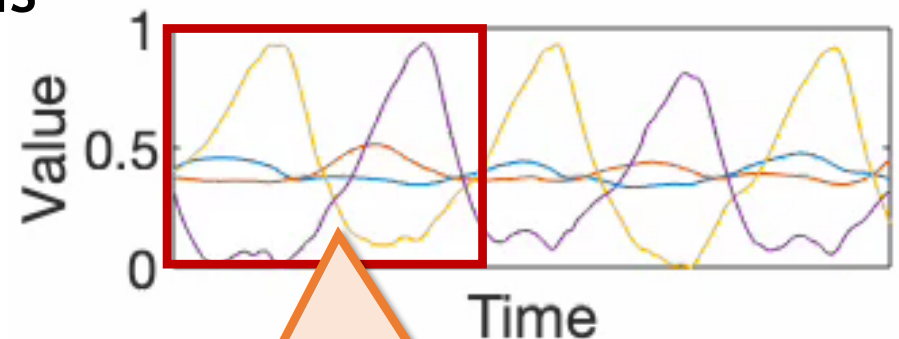- contain multiple patterns

# Motivation

**Given:** time-evolving data streams
- e.g., IoT sensors/Web click logs
- contain multiple patterns

**Answer:** the following questions:

1. What kind of patterns?
2. How many patterns?
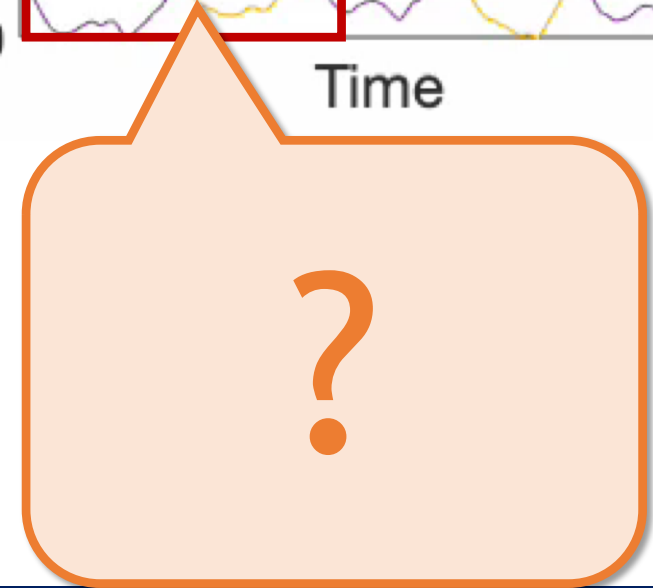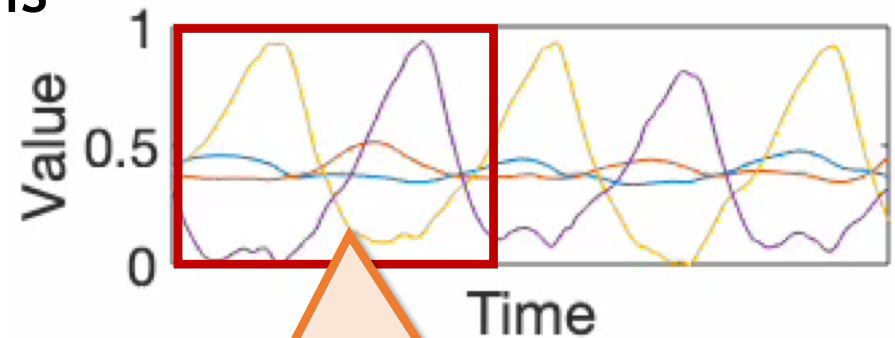3. When do patterns change?

# Motivation

**Given:** time-evolving data streams
- e.g., IoT sensors/Web click logs
- contain multiple patterns

**Requirements:**
- **Incremental**
  - We cannot access all historical data
- **Automatic**
  - # of patterns are unknown in advance
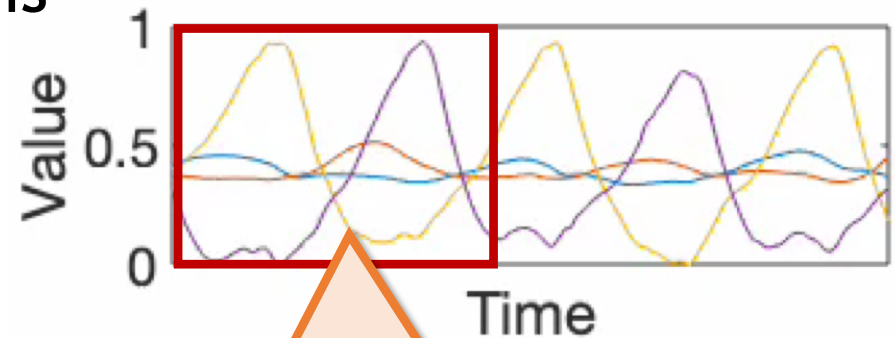  - without any parameter tunings

# Motivation

**Given:** time-evolving data streams
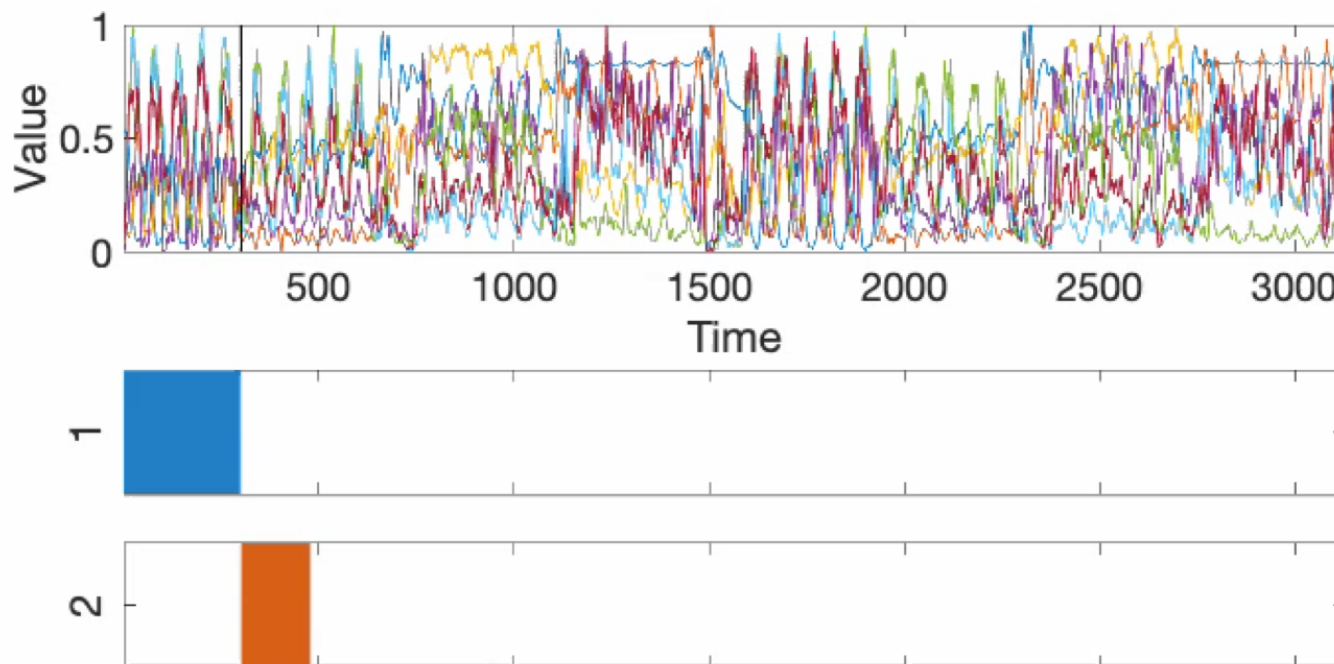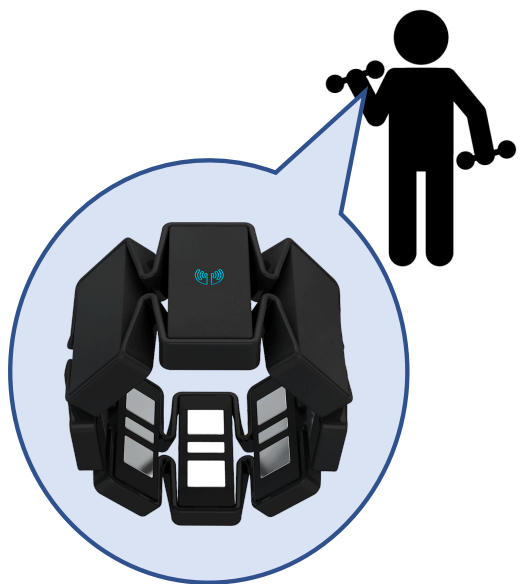- e.g., IoT sensors/Web click logs
- contain multiple patterns

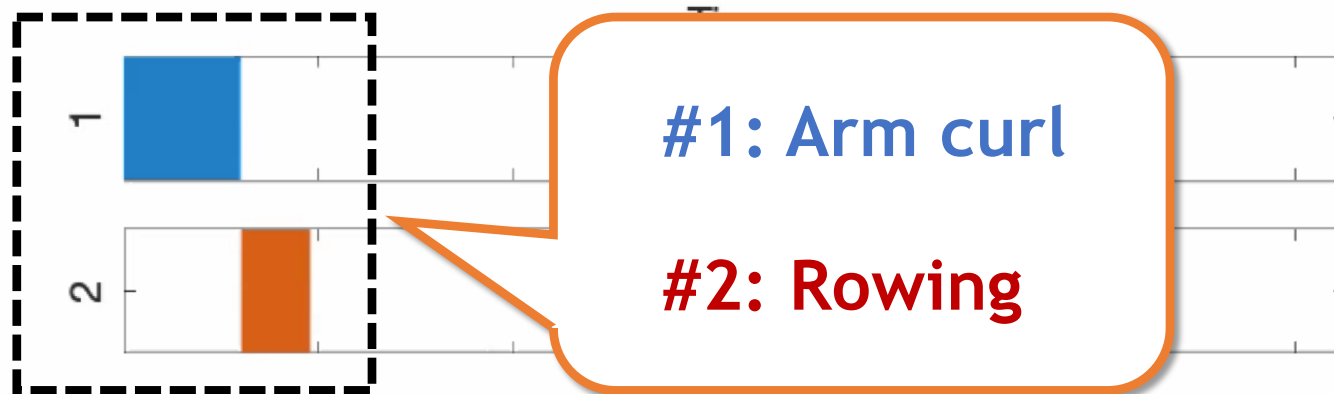**Requirements:**
- **Incremental**
  - We cannot access all historical data
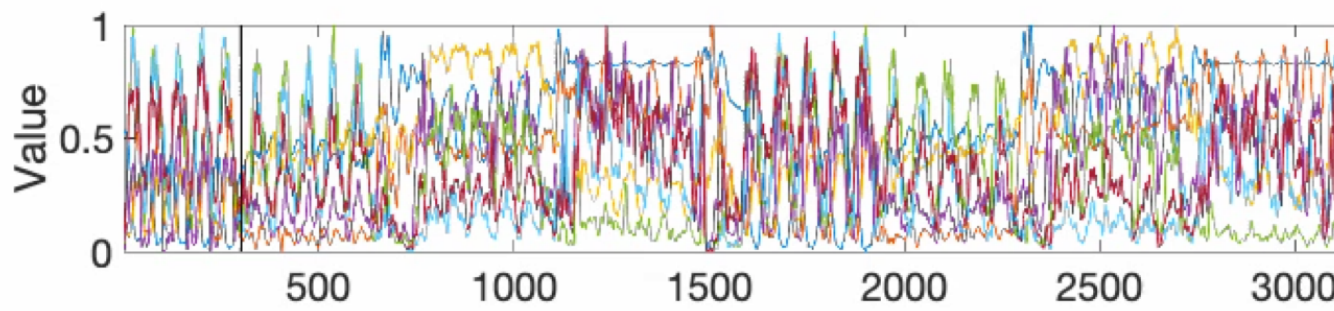
**StreamScope: automatic & incremental approach**

# Demo movie

# Demo movie



#1: Arm curl

#2: Rowing

# Demo movie

# Demo movie

#1: Arm curl

#2: Rowing

#3: Intervals

#4: side raise

#5: Push up

# Outline

1. ~~Motivation~~

2. **Problem definition**

3. Model

4. Streaming Algorithm

5. Experiments

6. Conclusions

# Problem definition

**Given:**
- Data stream $X$

**Find:**
1. Segment: $\mathcal{S}$
2. Regime: $\Theta$
3. Segment-membership: $\mathcal{F}$

# Problem definition

**Given** Data stream: set of d-dimensional vectors

$$X = \{x_1, \ldots, x_t\}$$

$d = 4$

# Problem definition

**Hidden**  Segment: start/end positions of each pattern

$$\mathcal{S} = \{s_1, \dots, s_m\}$$

$m = 8$

# Problem definition

**Hidden** Regime: segment groups

$$\Theta = \{\theta_1, \ldots, \theta_r, \Phi\}$$

$r = 4$

# Problem definition

**Hidden** Segment-membership: regime-assignment

$$\mathcal{F} = \{f(1), \dots, f(m)\}$$



e.g., $f(3) = 4$

$$\mathcal{F} = \{ \quad \textcolor{blue}{1}, \quad \textcolor{red}{2}, \quad \textcolor{purple}{4}, \quad \textcolor{red}{2}, \quad \textcolor{orange}{3}, \quad \textcolor{purple}{4}, \quad \textcolor{red}{2}, \quad \textcolor{blue}{1}, \quad \}$$

# Problem definition

**Given:** $d$-dimensional data stream $X = \{x_1, \ldots, x_t\}$

**Find:** compact description $\mathcal{C} = \{m, r, \mathcal{S}, \Theta, \mathcal{F}\}$ of $X$

$$\mathbf{=} \begin{cases} \bullet\, m \text{ segments } \mathcal{S} \\ \bullet\, r \text{ regimes } \Theta \\ \bullet \text{ segment-membership } \mathcal{F} \end{cases}$$

# Outline

1. ~~Motivation~~

2. ~~Problem definition~~

3. **Model**

4. Streaming Algorithm

5. Experiments

6. Conclusions

# Proposed model

**Goal:** find compact description C in **a streaming setting**

**Challenges:**
Q1. How can we represent regimes?
Idea (1): Hierarchical probabilistic model
Q2. How can we decide # of segments/regimes?
Idea (2): Model description cost

# Idea (1): hierarchical probabilistic model

Q. How to describe patterns?



Data stream

Model

stand ⟷ walk

run

Regimes

**Idea:** HMM-based probabilistic model
- **'within-regime'** transitions: A hidden Markov model $\theta = \{\pi, A, B\}$
- **'across-regime'** transitions: Regime transition matrix $\Phi = \{\phi_{ij}\}_{i,j=1}^{r}$

# Idea (1): hierarchical probabilistic model

**Full model** $\Theta = \{\theta_1, \ldots, \theta_r, \Phi\}$



Single HMM parameters:
$$\theta_i = \{\pi_i, A_i, B_i\}$$

Regimes

# Idea (1): hierarchical probabilistic model

**Full model** $\Theta = \{\theta_1, \ldots, \theta_r, \Phi\}$



Single HMM parameters:
$$\theta_i = \{\pi_i, A_i, B_i\}$$

Regime transition matrix:
$$\Phi = \{\phi_{ij}\}_{i,j=1}^{r}$$

# Idea (2): Incremental encoding scheme

Q. How to decide # of segments/regimes?

**Idea: Minimum description length (MDL)**
- Minimize the total description cost of a data stream
- Update 'optimal' # of segments/regimes

# Idea (2): Incremental encoding scheme

**Idea:** Minimize total encoding cost

$$\min \left( \boxed{\text{Cost}_M(C)} + \boxed{\text{Cost}_C(X|C)} \right)$$

Model cost      Coding cost

| Good compression | ⟳ | Good description |
|:---:|:---:|:---:|



— CostM
— CostC
— CostT

1 2 3 4 5 6 7 8 9 10

(# of r, m)

# Idea (2): Incremental encoding scheme

Q. How many new components does $\mathcal{C}$ need?



segment?
regime?
How many?

$X_c$

$\mathcal{C}$

Keep compact!

A segment    A state    A regime

# Idea (2): Incremental encoding scheme

Q. How many new components does $\mathcal{C}$ need?



$X_c$

segment?
regime?
How many?

$\mathcal{C}$

Keep compact!

$$\Delta Cost_T(X_c; \mathcal{C}) = \log^*(m_+) - \log^*(m) + \log^*(r_+) - \log^*(r)$$
$$+ \Delta Cost_M(\mathcal{S}) + \Delta Cost_M(\Theta) + \Delta Cost_M(\mathcal{F})$$
$$+ Cost_C(X_c | \Theta).$$

Details in paper

# Outline

1. ~~Motivation~~
2. ~~Problem definition~~
3. ~~Model~~
4. **Streaming Algorithm**
5. Experiments
6. Conclusions

# Streaming algorithms

- <u>Algorithms</u>

**Main** **StreamScope**
Optimize/update parameter set $\mathcal{C}$

1. **SegmentAssignment**
Identify regime transitions & segments

2. **RegimeGeneration**

Estimate new regimes $\theta$

# StreamScope

- <u>Overview</u>



Data stream $X$

$X_c$

$t \rightarrow$

**1. Keep current window:**
- The latest segment, $s_m$
- New observations, $x_t, \ldots$

$$X_c = s_m \cup x_t$$

$\mathcal{C}$

# StreamScope

- Overview



Data stream $X$

$X_c$

$t \rightarrow$

**1. Keep current window:**
- The latest segment, $s_m$
- New observations, $x_t, \ldots$

$$X_c = s_m \cup x_t$$

$\mathcal{C}$

**2. Update model set $\mathcal{C}$**
- Minimize $\Delta Cost_T(X_c | \mathcal{C})$
  Increase segments?
  (SegmentAssignment)
  vs.
  Increase states/regimes?
  (RegimeGeneration)

# StreamScope

- <u>Overview</u>



Data stream $X$

Segmented

$t \rightarrow$

$X_c$

**3. Update $X_c$**
- If pattern has changed

**1. Keep current window:**
- The latest segment, $s_m$
- New observations, $x_t, \dots$
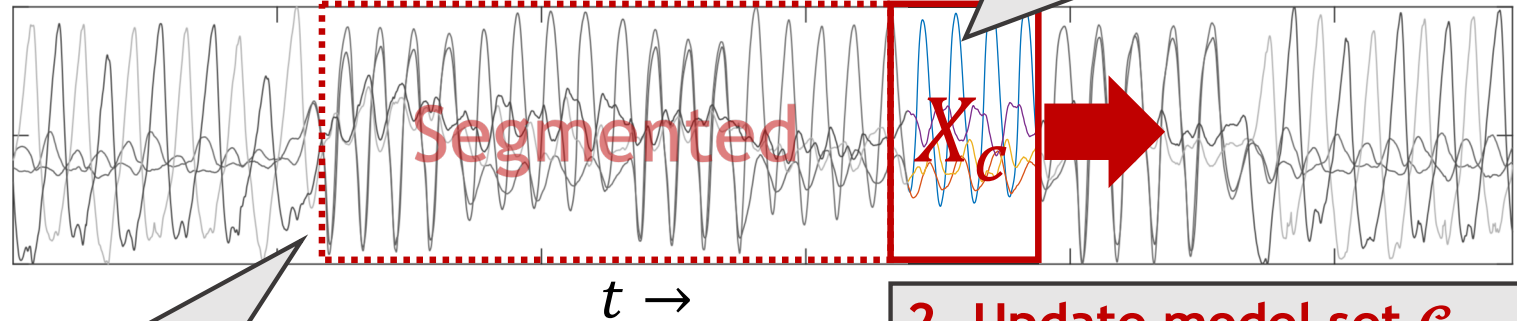
$$X_c = s_m \cup x_t$$

$\mathcal{C}$

**2. Update model set $\mathcal{C}$**
- Minimize $\Delta Cost_T(X_c | \mathcal{C})$
  Increase segments?
  (SegmentAssignment)
  vs.
  Increase states/regimes?
  (RegimeGeneration)

# 1. SegmentAssignment
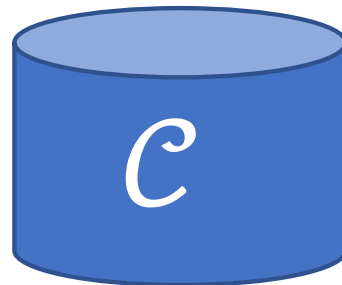
**Given:**

- Observation $x_t$
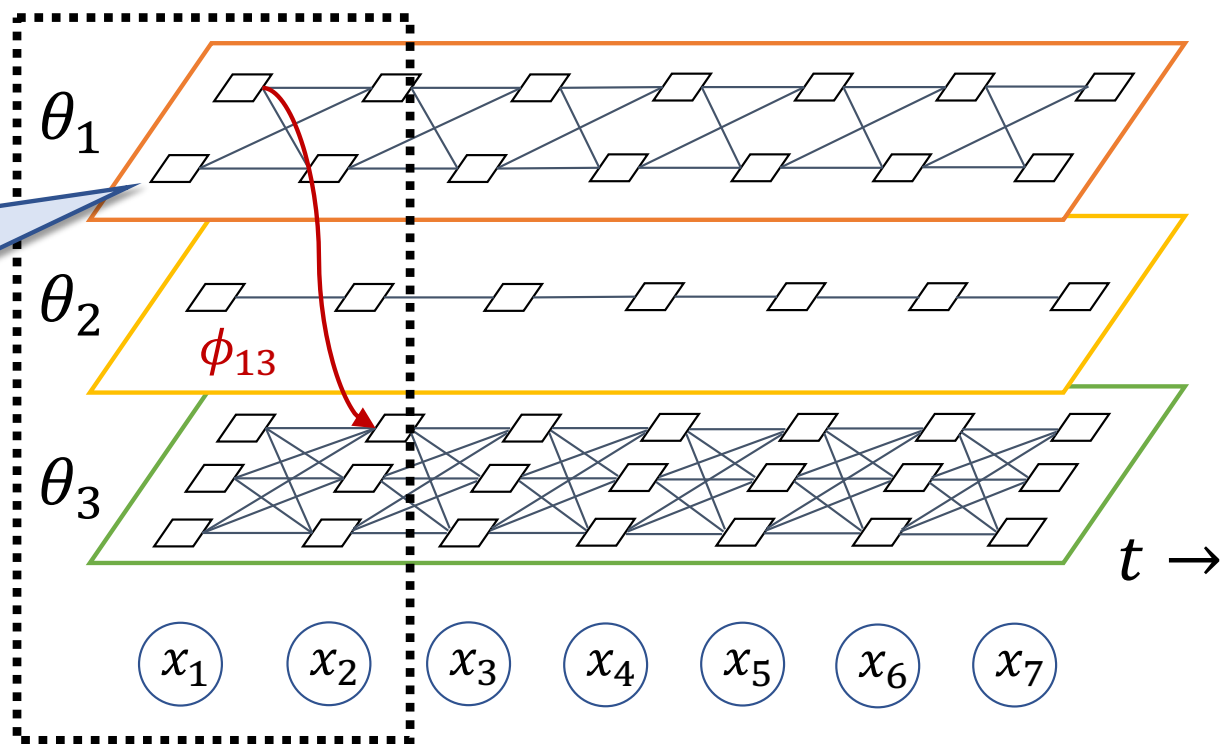- Model parameter set $\Theta = \{\theta_1, \ldots, \theta_r, \Phi\}$

**Find:**

- Optimal cut point between regimes: $\{m, \mathcal{S}, \mathcal{F}\}$

# 1. SegmentAssignment

Overview

Dynamic programing algorithm to compute $P(x_t|\Theta)$

$\theta_1$

$\theta_2$

$\phi_{13}$

$\theta_3$

$t \rightarrow$

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$  $x_6$  $x_7$

# 1. SegmentAssignment

## Overview

Dynamic programing algorithm to compute $P(x_t|\Theta)$

Keep all candidate cut points $\mathcal{L} = \{L_1, L_2, ...\}$

$\theta_1$

$\theta_2$

$\theta_3$

*switch??* $\quad \phi_{12} \quad L_2 = \{4, 2\}$

$\phi_{13} \quad L_2 = \{2, 3\} \quad$ *switch??*

$t \rightarrow$

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5 \quad x_6 \quad x_7$

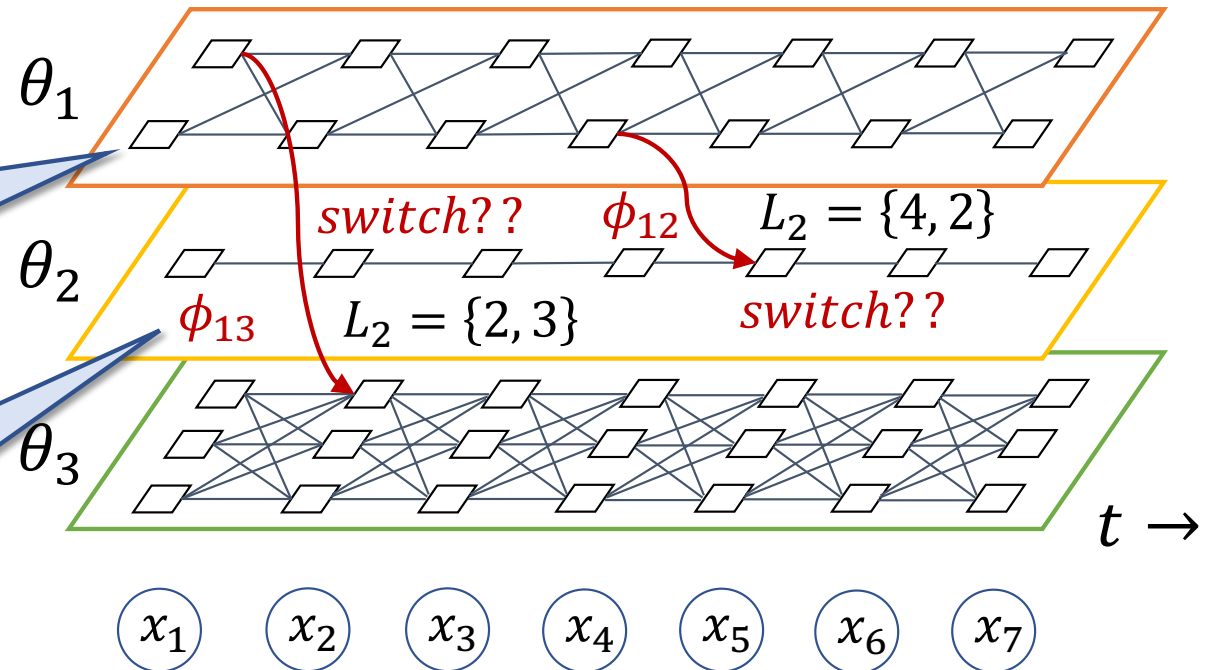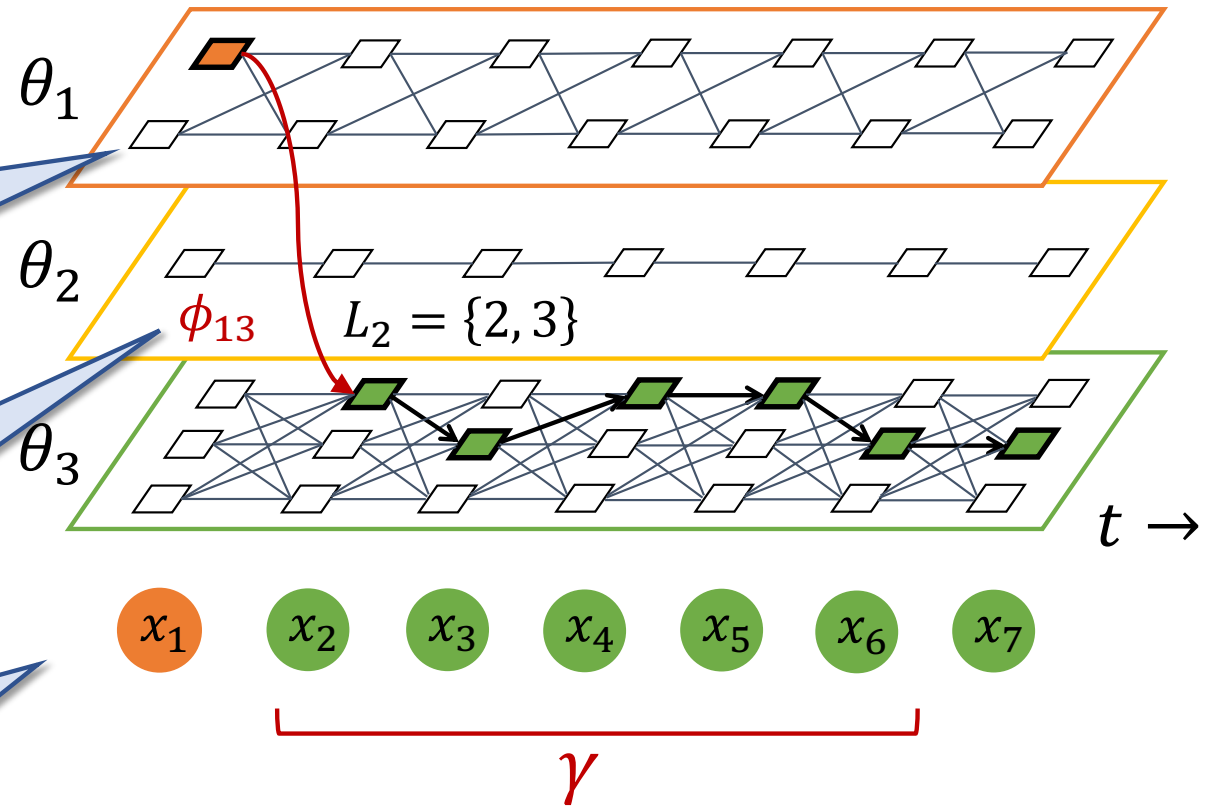# 1. SegmentAssignment

Overview



Dynamic programing algorithm to compute $P(x_t|\Theta)$

Keep all candidate cut points $\mathcal{L} = \{L_1, L_2, \dots\}$

$\gamma -$ guarantee: $\gamma \propto mean(|s|)$

$\theta_1$

$\theta_2$

$\phi_{13}$  $L_2 = \{2, 3\}$

$\theta_3$

$t \rightarrow$

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$  $x_6$  $x_7$

$\gamma$

# 2. RegimeGeneration
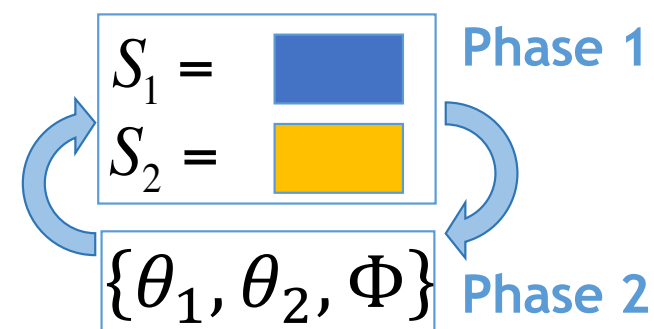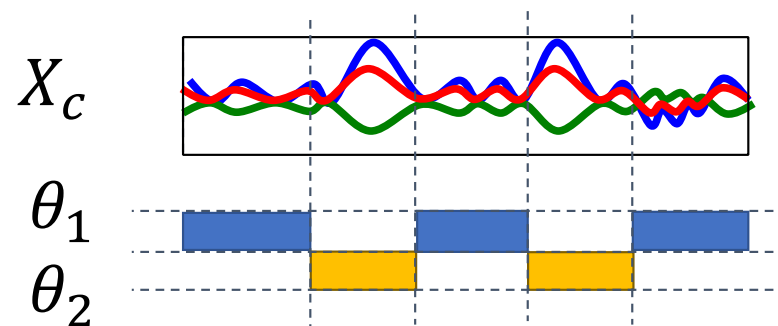
**Given:**

- Current window $X_c$

**Find:**

- New regimes: parameter set $\{m,\ r,\ \mathcal{S},\ \Theta,\ \mathcal{F}\}$ for $X_c$

# 2. RegimeGeneration

1. <u>Two phase iterative approach</u>

- **Phase1:** split segments into 2 groups

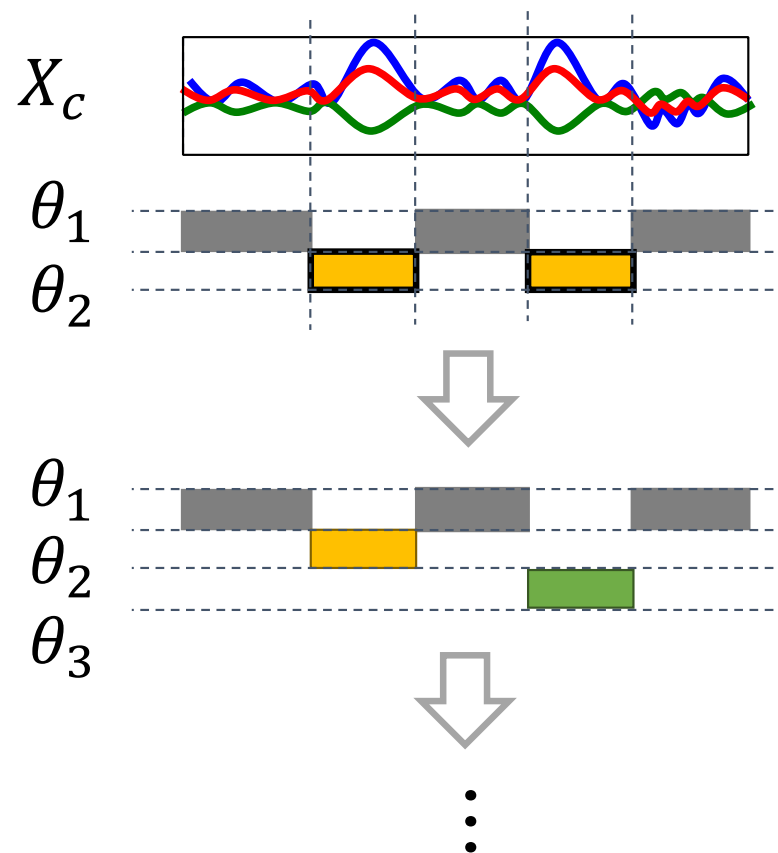- **Phase2:** update 2 model parameters

$X_c$

$\theta_1$
$\theta_2$

$S_1 =$      **Phase 1**

$S_2 =$

$\{\theta_1, \theta_2, \Phi\}$ **Phase 2**

# 2. RegimeGeneration

## 1. Two phase iterative approach

- **Phase1:** split segments into 2 groups
- **Phase2:** update 2 model parameters

## 2. Recursively split new regimes

- While total cost can be reduced



$X_c$

$\theta_1$

$\theta_2$

$\theta_1$

$\theta_2$

$\theta_3$

# Outline

1. ~~Motivation~~
2. ~~Problem definition~~
3. ~~Model~~
4. ~~Streaming Algorithm~~
5. **Experiments**
6. Conclusions

# Experiments

- We answer the following questions:

**Q1. Effectiveness:**
How successful is it in discovering patterns?
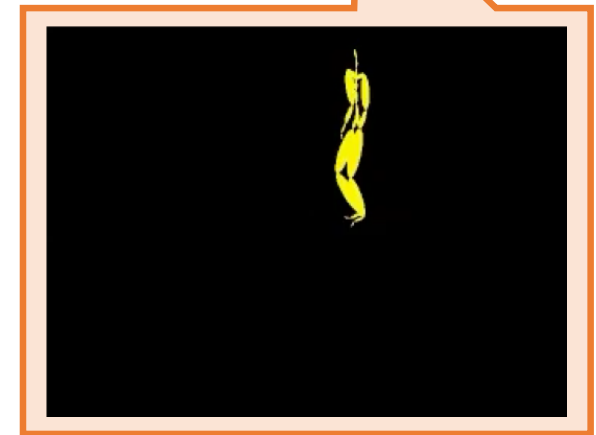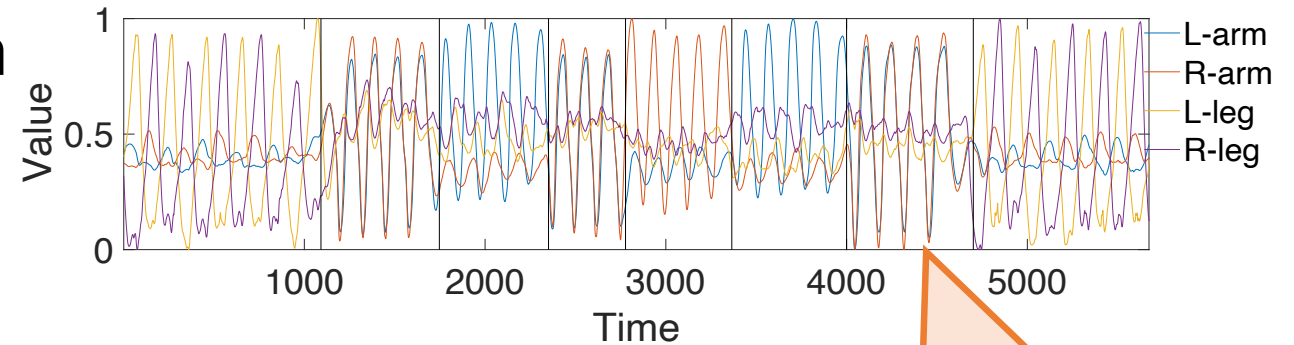
**Q2. Accuracy:**
How well does it find cut-points & regimes?

**Q3. Scalability:**
How does it scale in terms of time & memory consumption?

# Q1. Effectiveness - #Mocap

- MoCap sensor stream

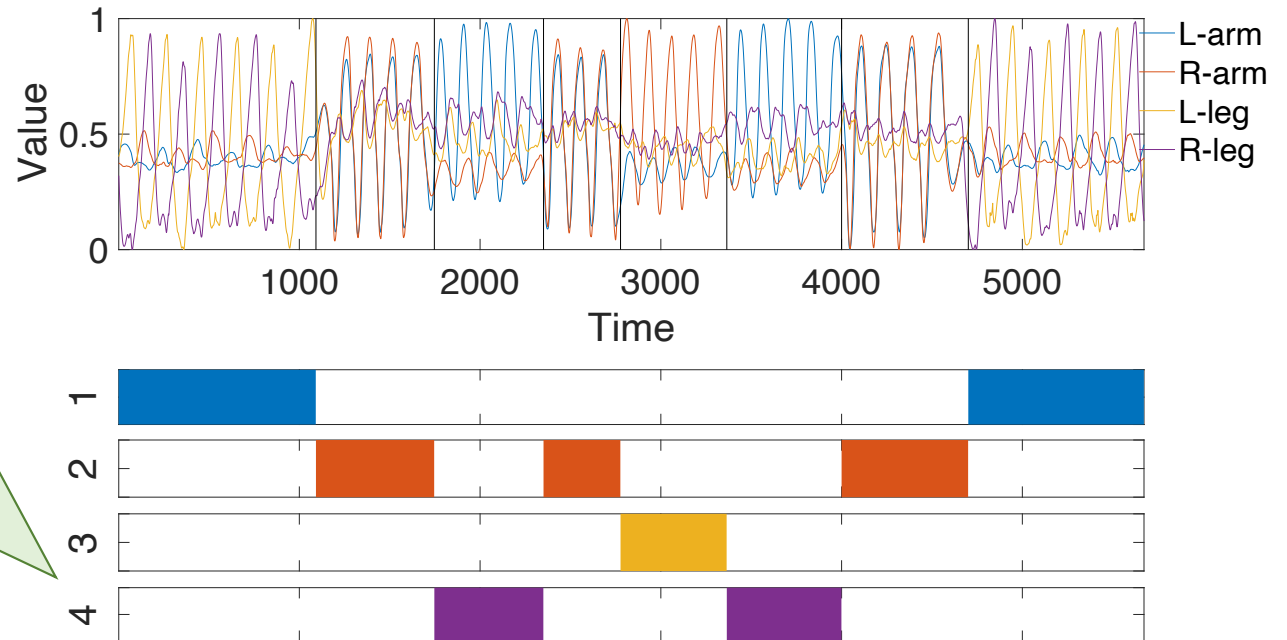# Q1. Effectiveness - #Mocap

- MoCap sensor stream



#1 Going straight

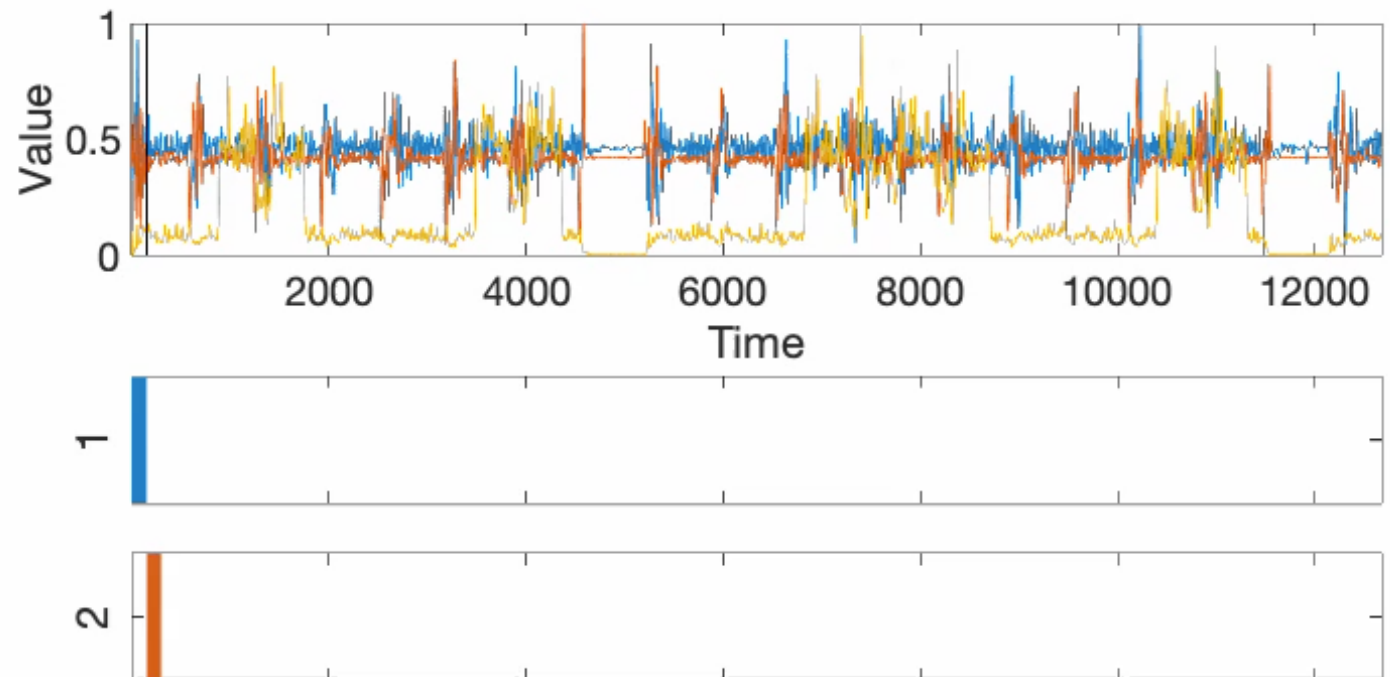#2 Stretching arms

#3 Stretching right arm

#4 Stretching left arm

StreamScope can find intuitive patterns **automatically**
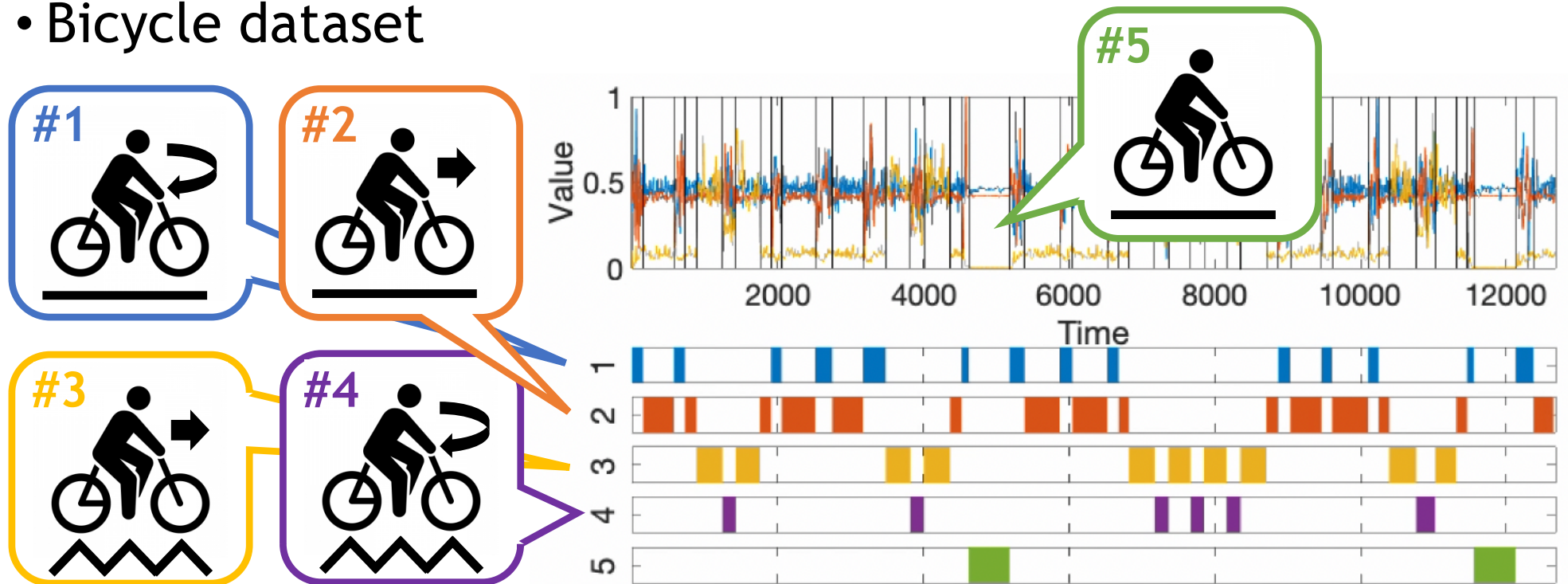
# Q1. Effectiveness - #Bicycle

- Bicycle dataset

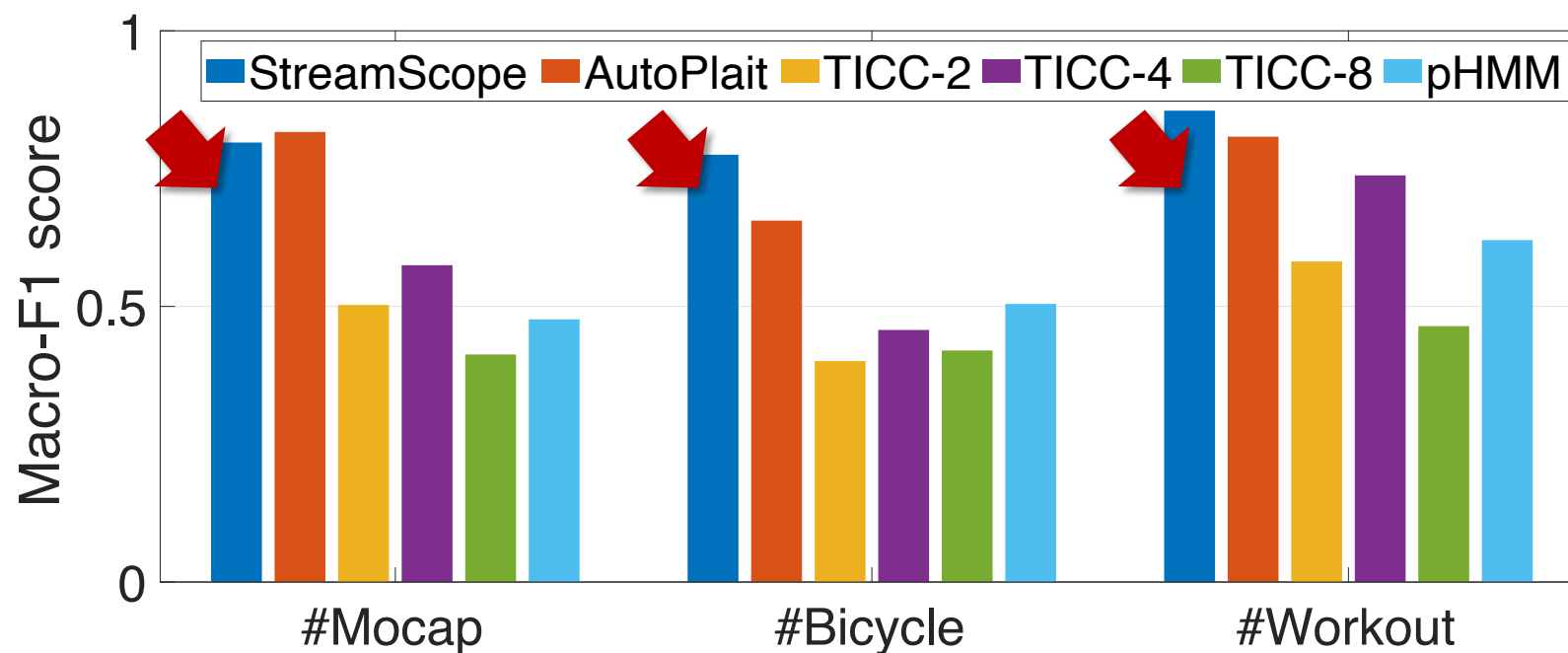

$\approx$

Acceleration
– (X,Y,Z)

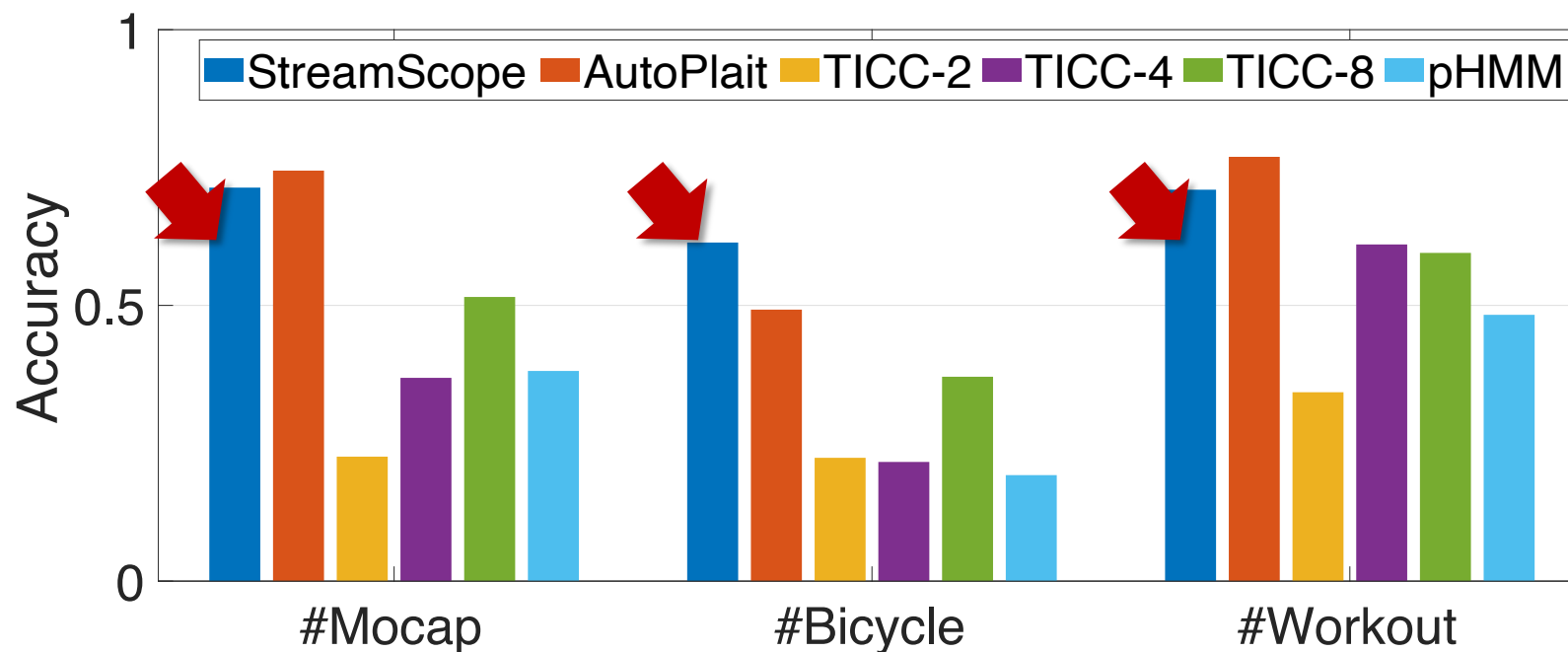# Q1. Effectiveness - #Bicycle

- Bicycle dataset

# Q2. Accuracy

- **Segmentation** accuracy (**higher is better**)



**Good accuracy** compared with other methods

Sakurai Lab. K.Kawabata et al.
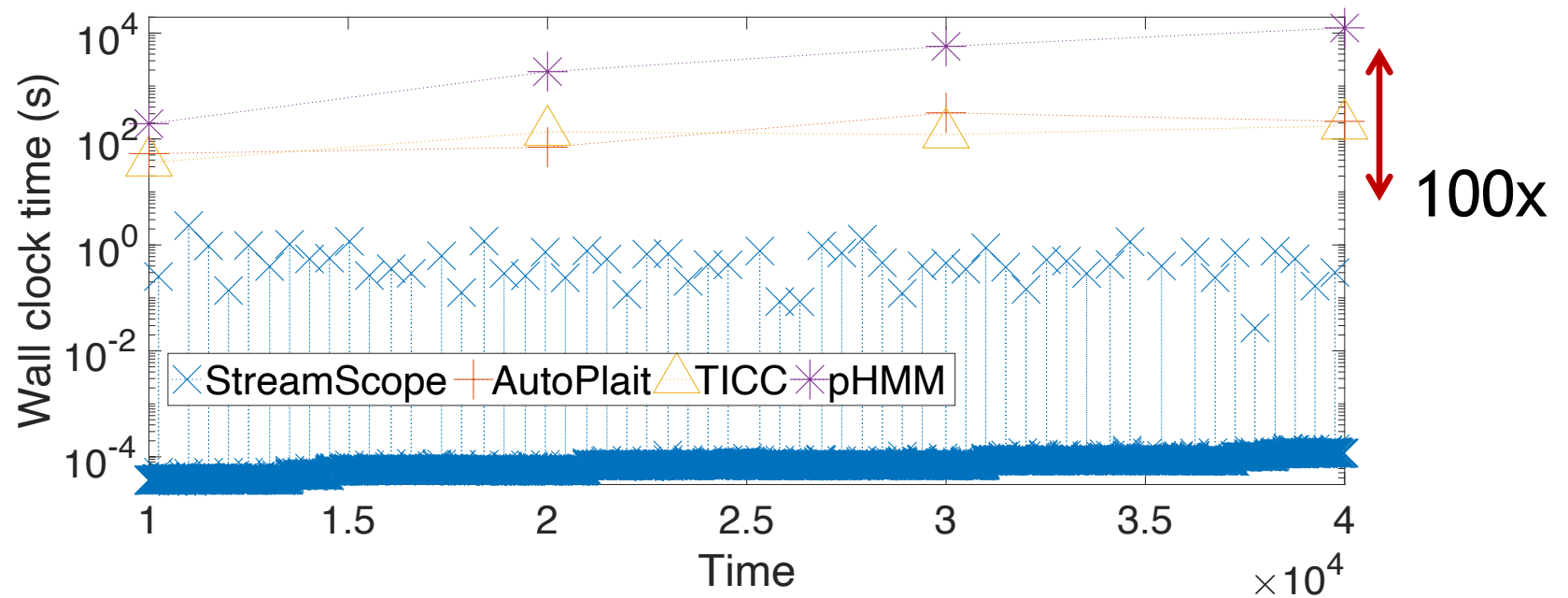
# Q2. Accuracy

- **Clustering** accuracy (**higher is better**)



**Good accuracy** compared with other methods
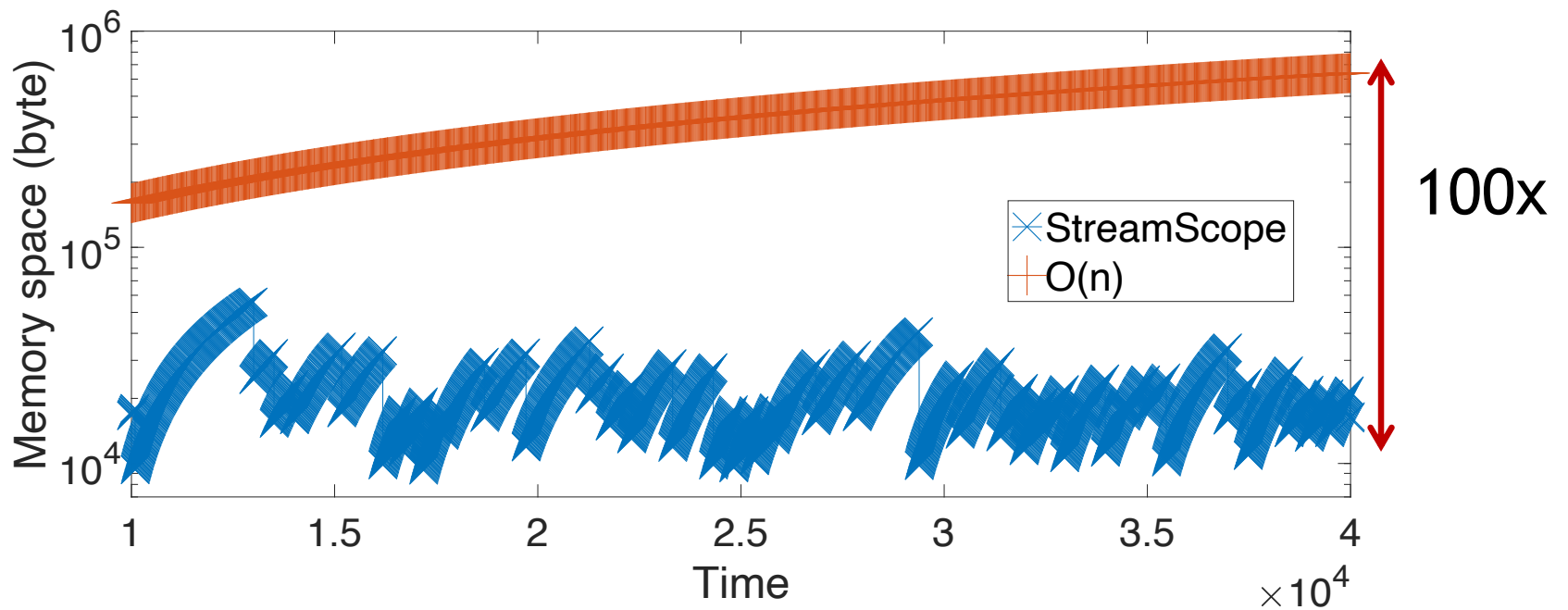
# Q3. Scalability

- Wall clock time vs. stream length



The complexity is **independent of the data length**

# Q3. Scalability

- Memory space vs. stream length



The complexity is **independent of the data length**

# Conclusions

StreamScope has the following advantages:

✓ **Effective:**
Find optimal segments/regimes

✓ **Adaptive:**
Automatic and incremental

✓ **Scalable:**
It does not depend on data length

# Thank you !