# Multi-Aspect Mining of Complex Sensor Sequences

Takato Honda[1] , Yasuko Matsubara[1], Ryo Neyama[2], Mutsumi Abe[2], Yasushi Sakurai[1]

[1]AIRC-ISIR, Osaka University

[2]Toyota Motor Corporation

# Motivation

**Analysis of IoT sensor data, e.g., car**
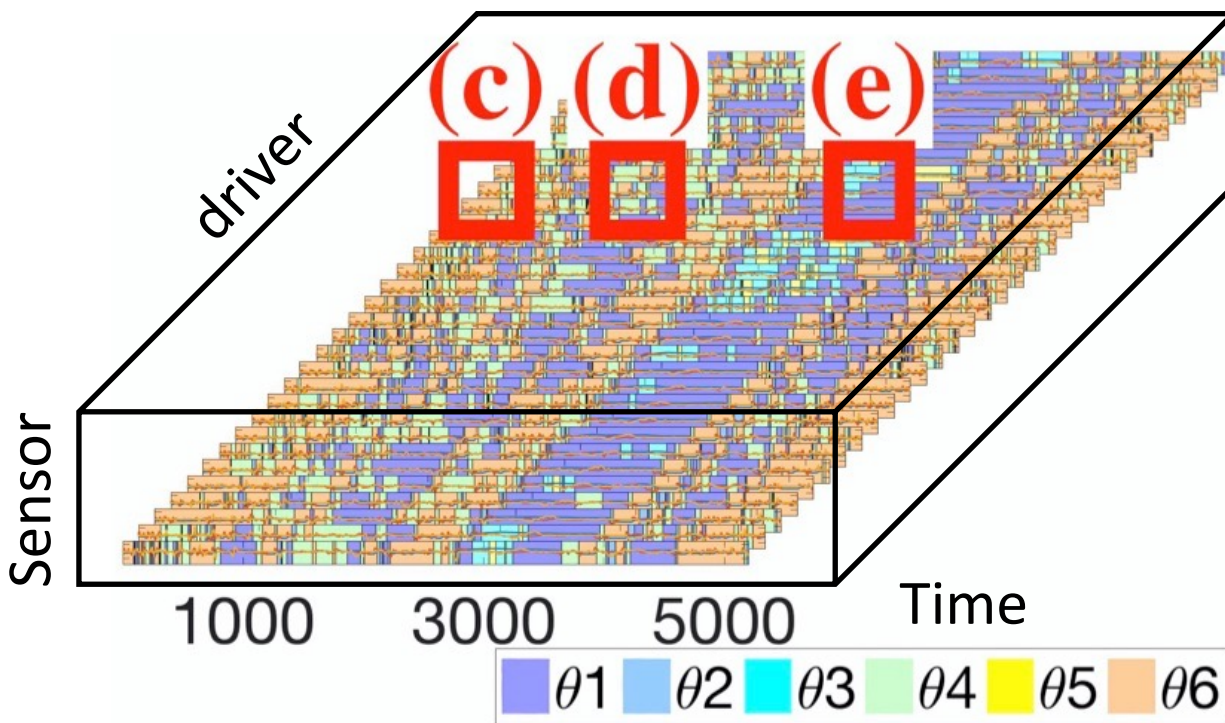 **- Advanced driving assistance service**

Risk

Fuel

GO!

Congestion

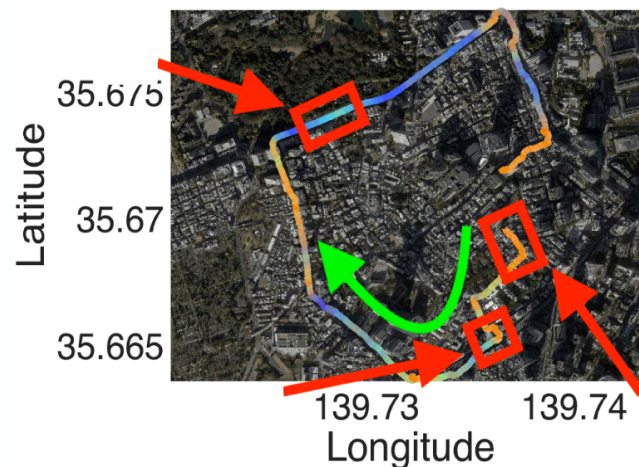# Motivation

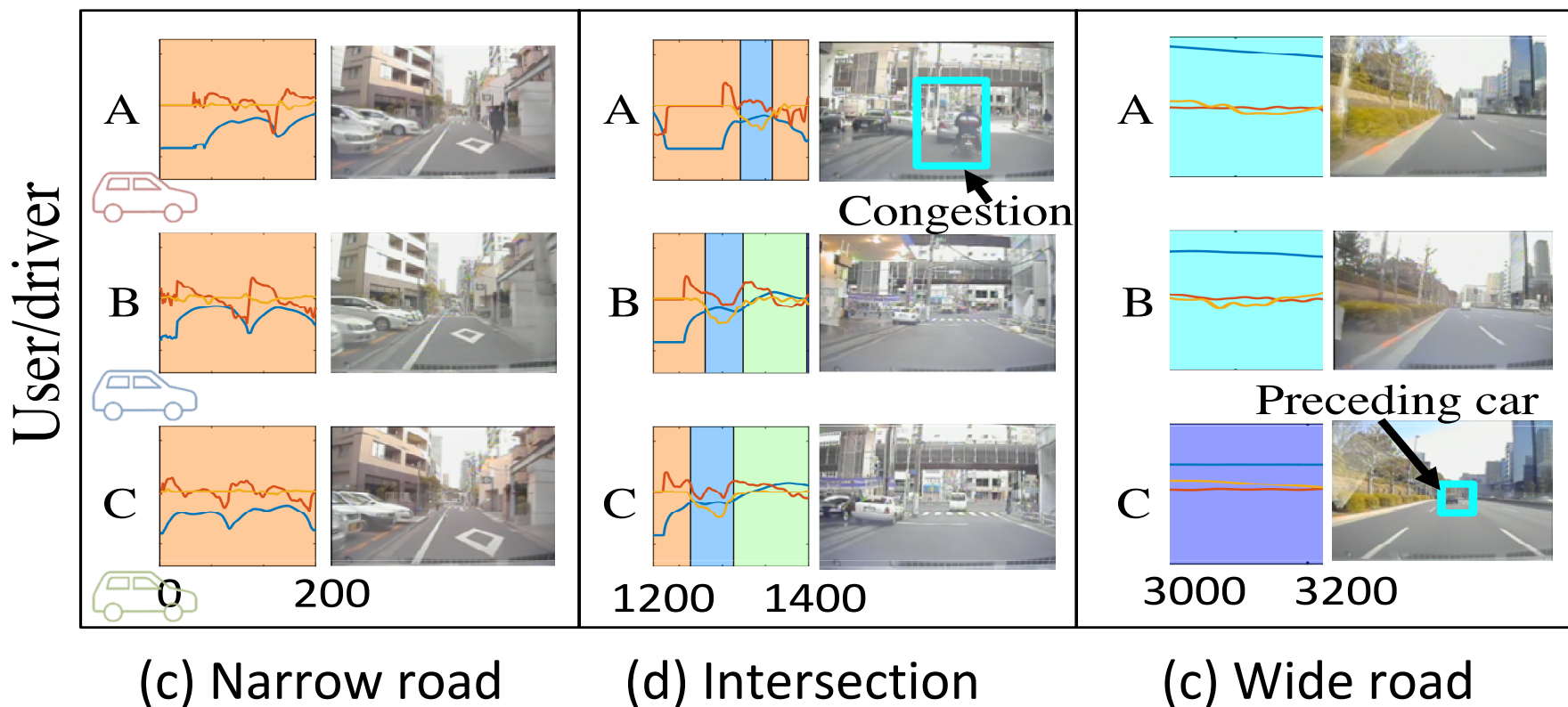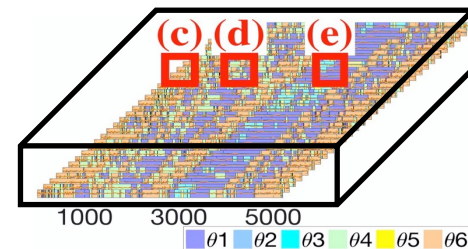## IoT sensor data is a tensor (sensor × driver × time)



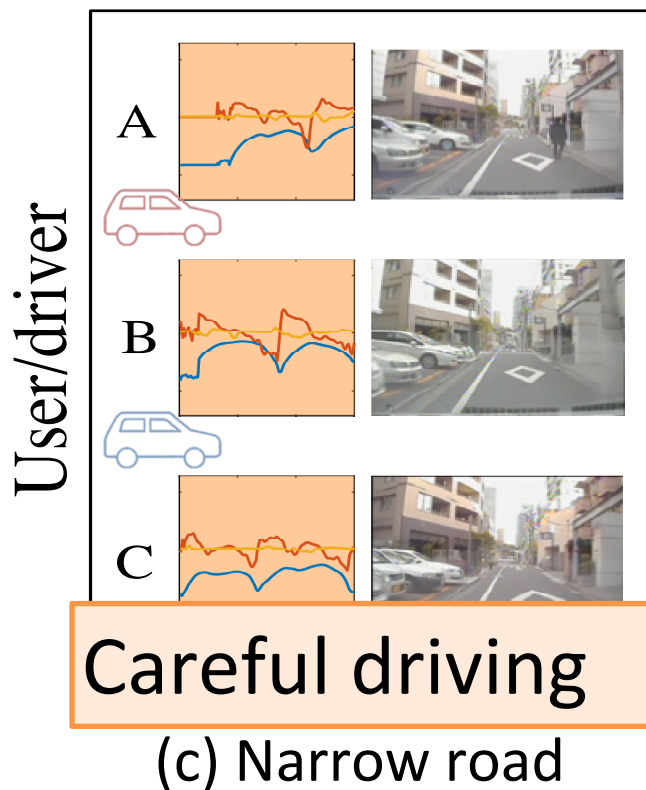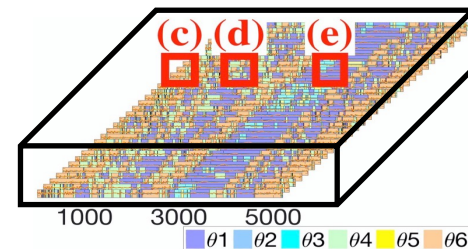(a) Time series tensor of automobile dataset

(b) On a map

# Motivation

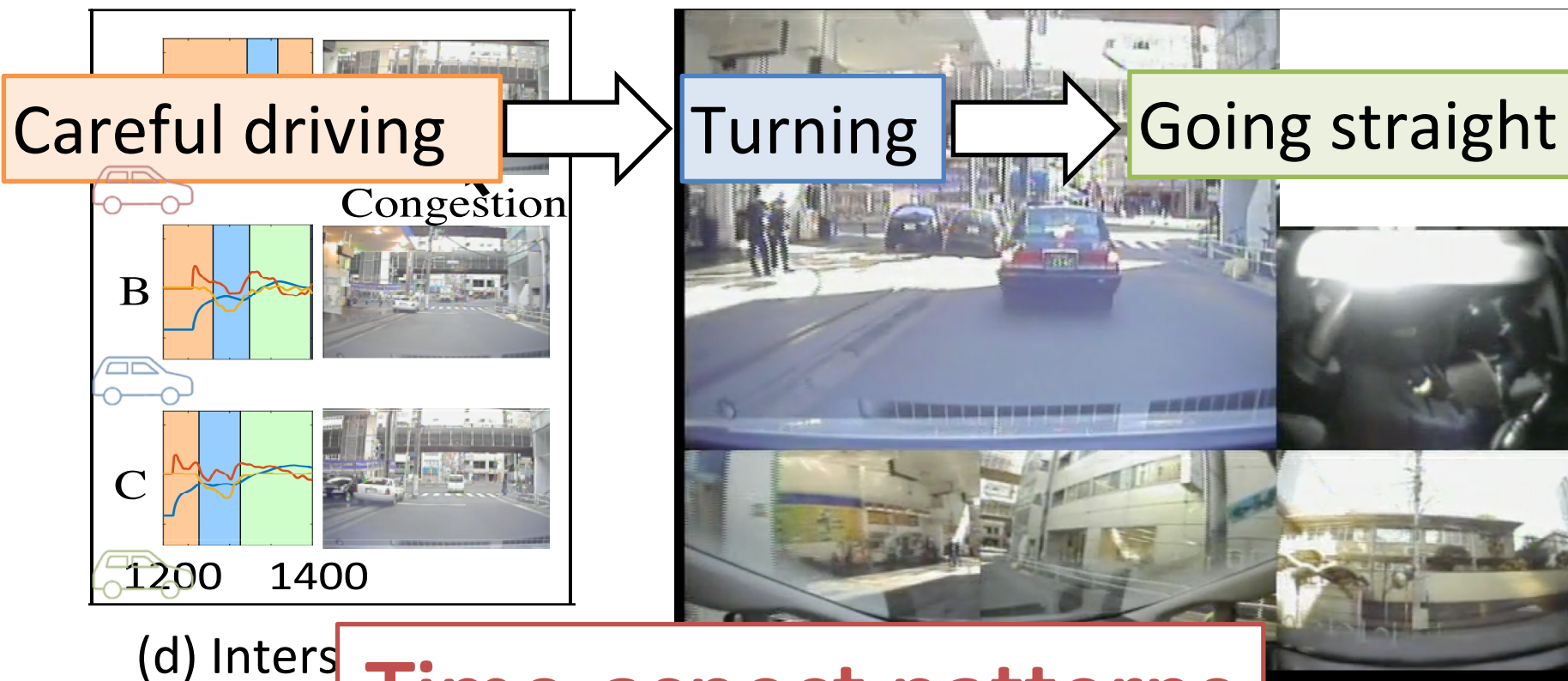**Tensor has multi-aspect patterns:**
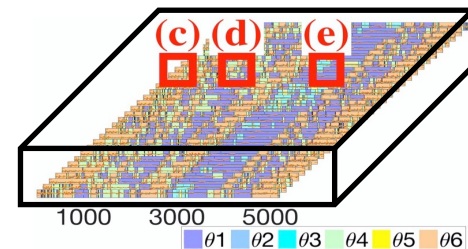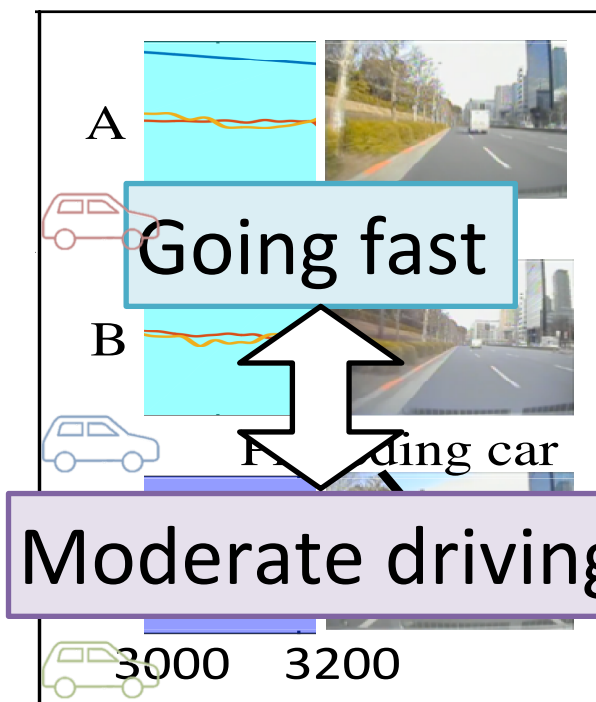**time-aspect** and **user-aspect**

© 2019 Takato Honda et al.

# Motivation

**Tensor has multi-aspect patterns:**
**time-aspect** and **user-aspect**
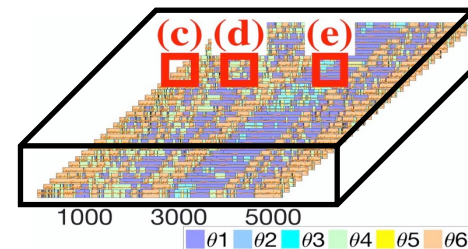


User/driver

A

B

C

Careful driving

(c) Narrow road

© 2019 Takato Honda et al.

**Tensor has multi-aspect patterns:**
**time-aspect** and **user-aspect**

1000  3000  5000

θ1 θ2 θ3 θ4 θ5 θ6

Careful driving ➡ Turning ➡ Going straight

Congestion

B

C

1200    1400

(d) Inters...

**Time-aspect patterns**

**Tensor has multi-aspect patterns:**
**time-aspect** and **user-aspect**



A

B

Going fast

Heading car

Moderate driving

3000    3200

(e) Wide

User-aspect patterns

# Motivation

**Given: Time-series tensor**
**(sensor ✕ user ✕ time)**
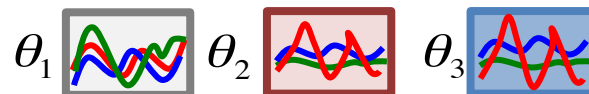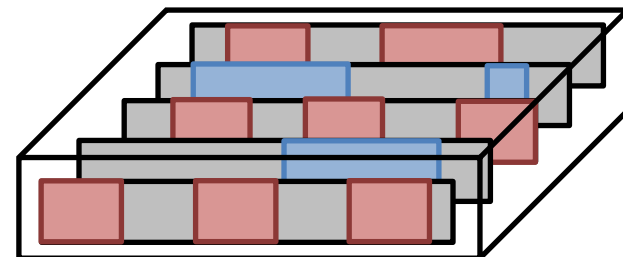


**Find: Multi-aspect patterns**
**(time and user-aspect)**

**Automatically & quickly**



$\theta_1$ $\theta_2$ $\theta_3$

# Outline

- **Motivation**
- **Problem definition**
- **Main ideas**
- **Algorithms**
- **Experiments**
- **Conclusions**

# Problem definition

## Key concepts

- **Tensor:** $\mathcal{X}$ — given
- **Segment:** $S$ — hidden
- **Regime:** $\Theta$ — hidden
- **Segment-membership:** $F$ — hidden

© 2019 Takato Honda et al.

# Problem definition

**Tensor** : $\mathcal{X} \in R^{d \times w \times n} = \{X_1, \dots, X_w\}$

given

# Problem definition

**Segment :** $S = \{s_1, \ldots, s_m\}$

hidden
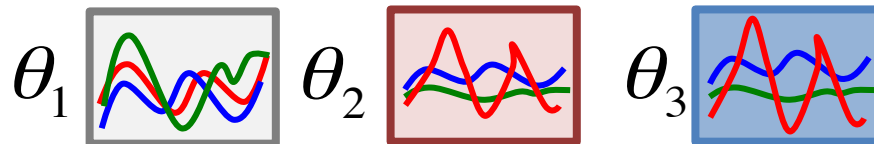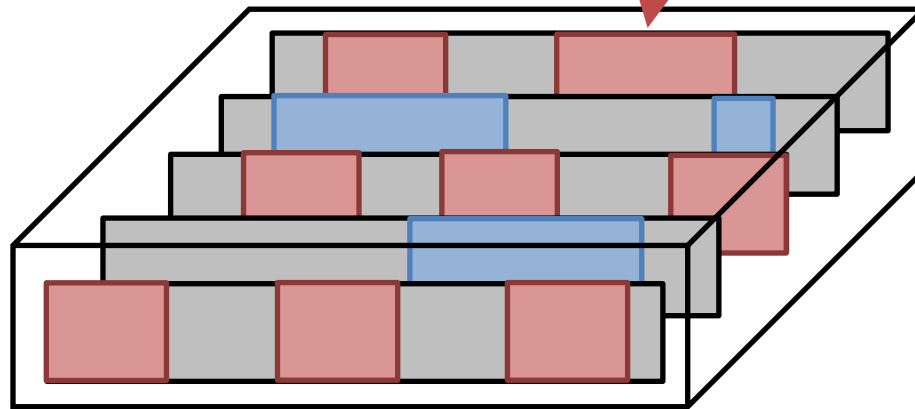
$s_i = \{t_s, t_e, userID\}$

**start position**

**end position**

**m = 25 segments**



$\theta_1$    $\theta_2$    $\theta_3$
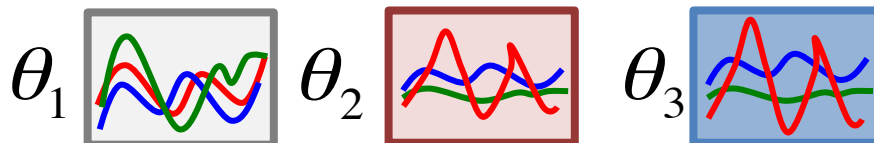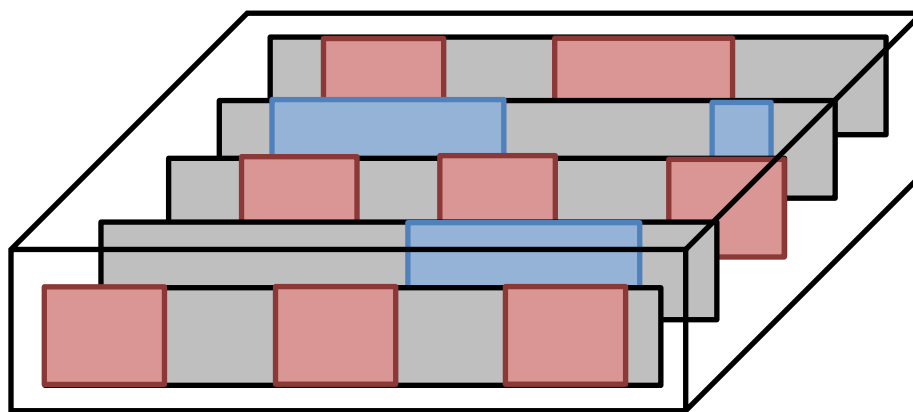
# Problem definition

**Regime:** $\Theta = \{\theta_1, \theta_2, \ldots, \theta_r, \Delta_{r \times r}\}$

hidden

$\theta_i = \{\pi, A, B\}$ (hidden Markov model)

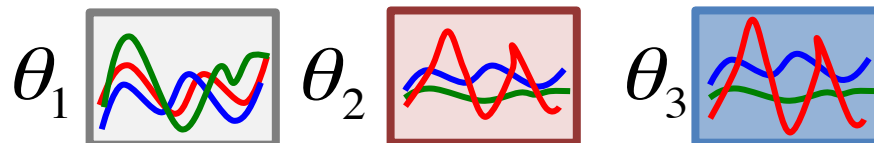**Initial prob.** **transition prob.** **output prob.**

**r = 3 regimes**

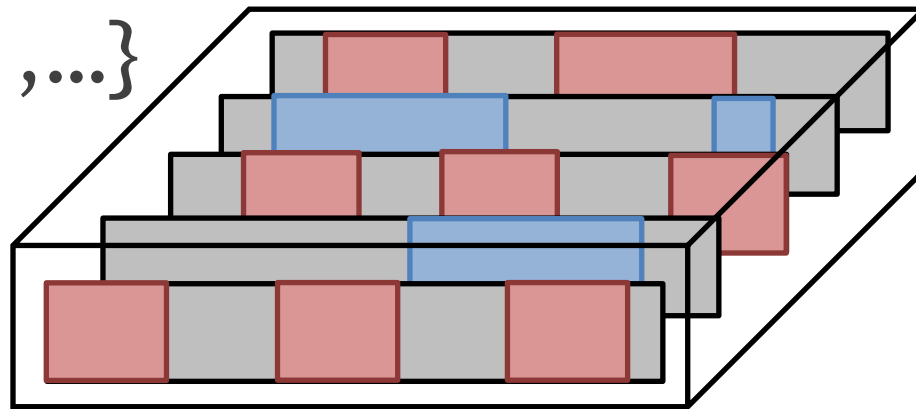$\theta_1$ $\theta_2$ $\theta_3$

**Membership:** $\quad F = \{f_1, f_2, \ldots, f_m\}$

hidden

$$1 \leq f_i \leq r$$

Example:
F = {1,2,1,2,1,...}



$\theta_1$  $\theta_2$  $\theta_3$

# Problem definition

**Given:** **tensor** $\mathcal{X}$

$$\mathcal{X} = \{X_1, \ldots, X_w\}$$

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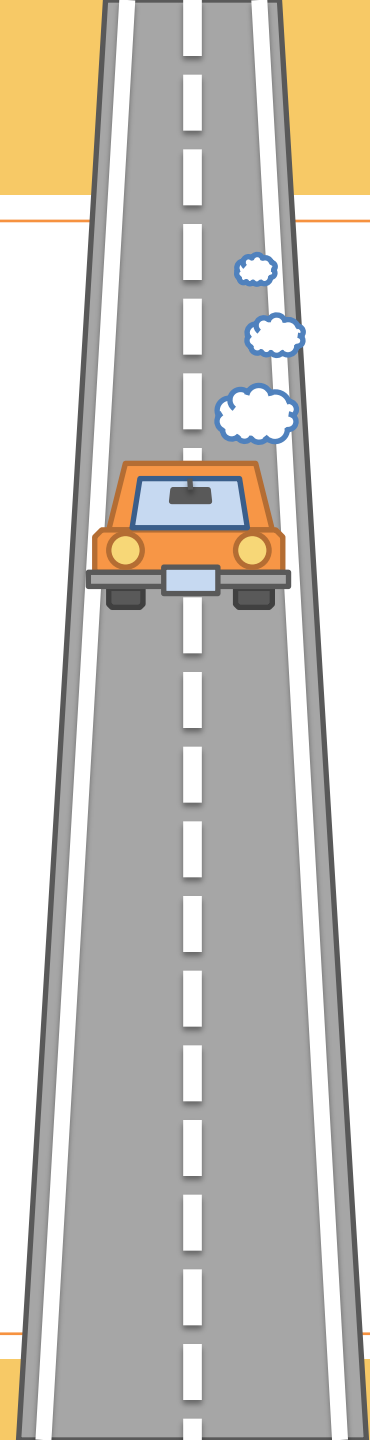
$$C = \{m, r, S, \Theta, F\}$$

**Automatically & quickly**

$\theta_1$ $\theta_2$ $\theta_3$

© 2019 Takato Honda et al.

# Outline

- **Motivation**
- **Problem definition**
- **Main ideas**
- **Algorithms**
- **Experiments**
- **Conclusions**

# Main ideas

**Goal:** compact description of

$$C = \{m, r, S, \Theta, F\}$$

**without user intervention**

**Challenges:**

**Q1. How to decide m and r automatically**

**Q2. How to find multi-aspect regimes**

# Main ideas

**Goal:** compact description of

$$C = \{m, r, S, \Theta, F\}$$

**without user intervention**

**Challenges:**

**Q1. How to decide m and r automatically**

**Idea 1: Model description cost**
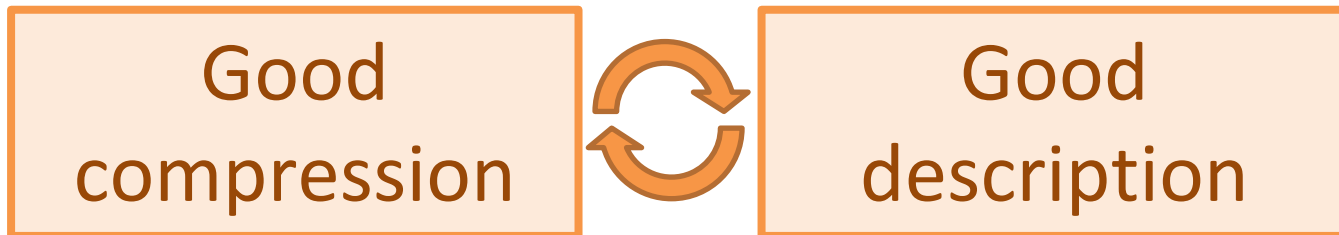
**Q2. How to find multi-aspect regimes**

**Idea 2: Multi-splitting algorithm**

# (1): model description cost

**Q1. How to decide # of regimes/segments?**

> **Idea 1: Model description cost**
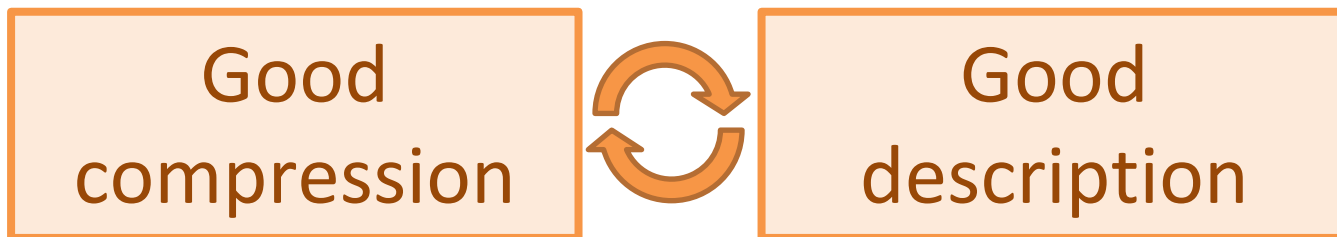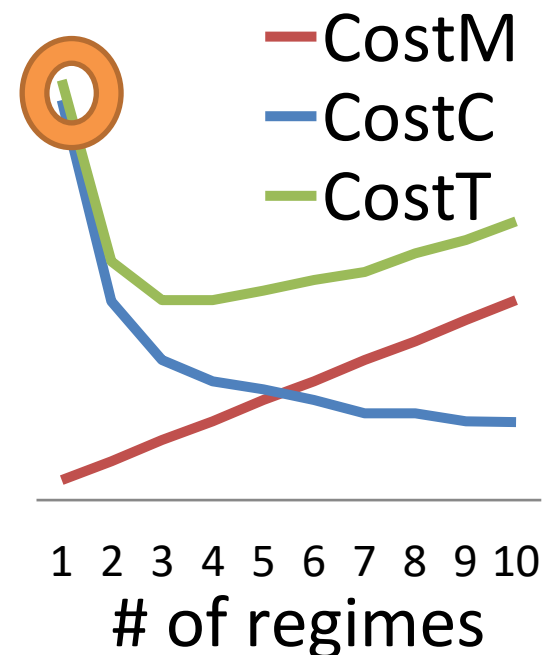> - **Minimize coding cost**
> - **Optimal # of segments/regimes**

Good compression ⟳ Good description

# (1): model description cost

**Idea: Minimize total cost**

$$\min \left( \boxed{\text{Cost}_M(\textcolor{red}{M})} + \boxed{\text{Cost}_c(\textcolor{blue}{X}|\textcolor{red}{M})} \right)$$

**Model cost**   **Coding cost**



— CostM
— CostC
— CostT

1 2 3 4 5 6 7 8 9 10
# of regimes

Good compression ⟳ Good description

© 2019 Takato Honda et al.
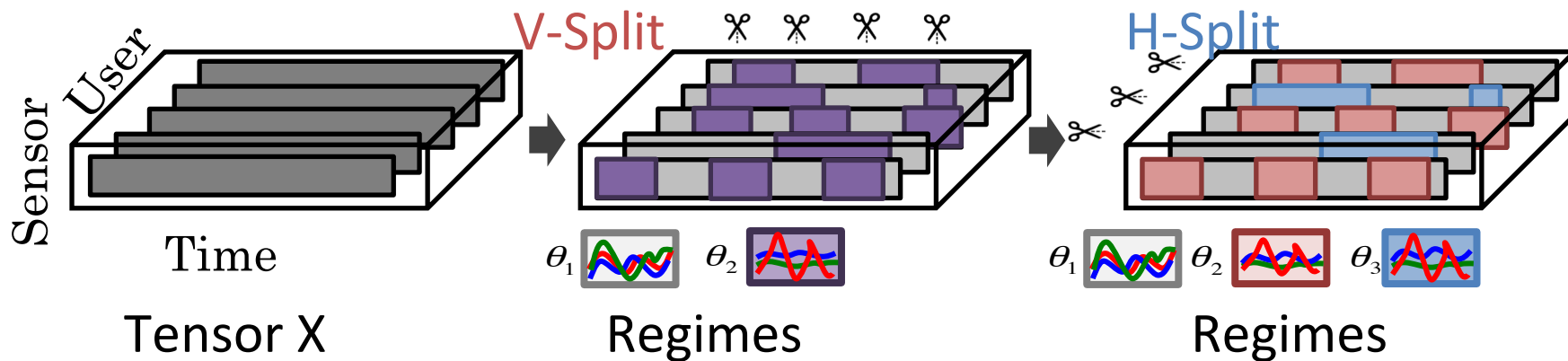
## Q2. How to find multi-aspect regimes?

**Idea 2: Multi-aspect splitting algorithm**

- **Find time-aspect transitions**
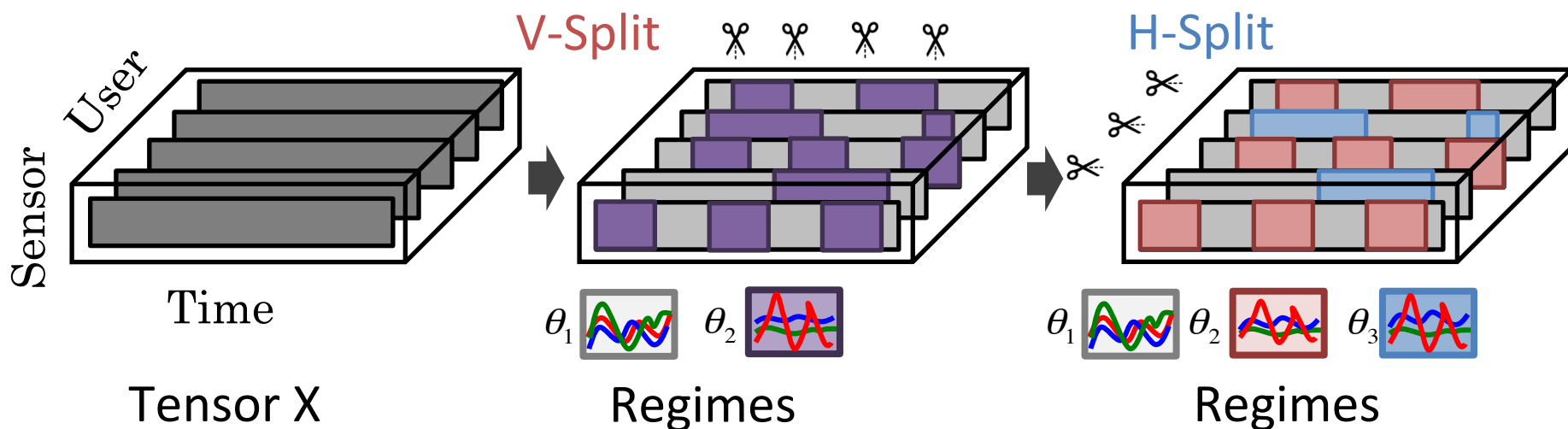- **And their differences between users**



Tensor X      Regimes      Regimes

# (2): Multi-aspect mining

**V-Split (vertical):**

split $\mathcal{X}$ into **time**-aspect

**H-Split (horizontal):**

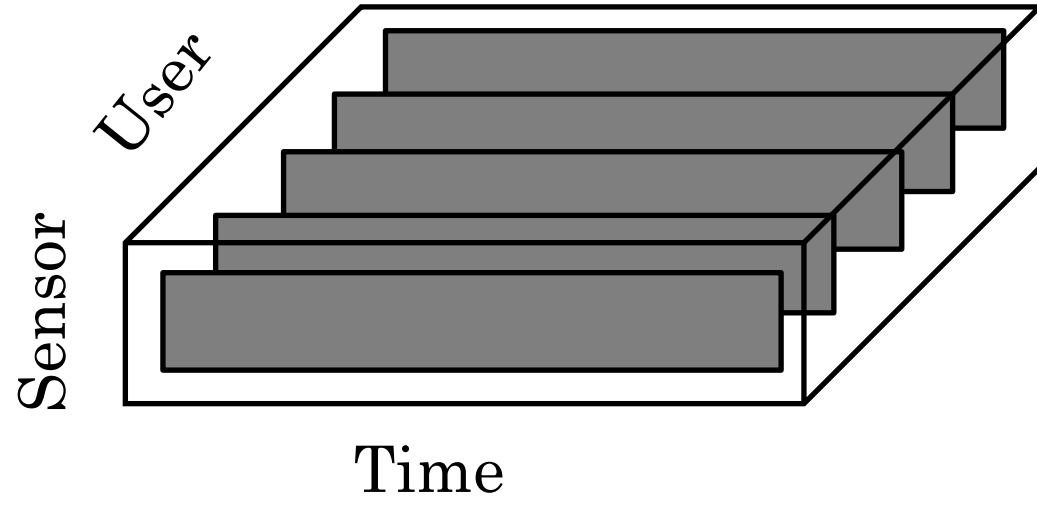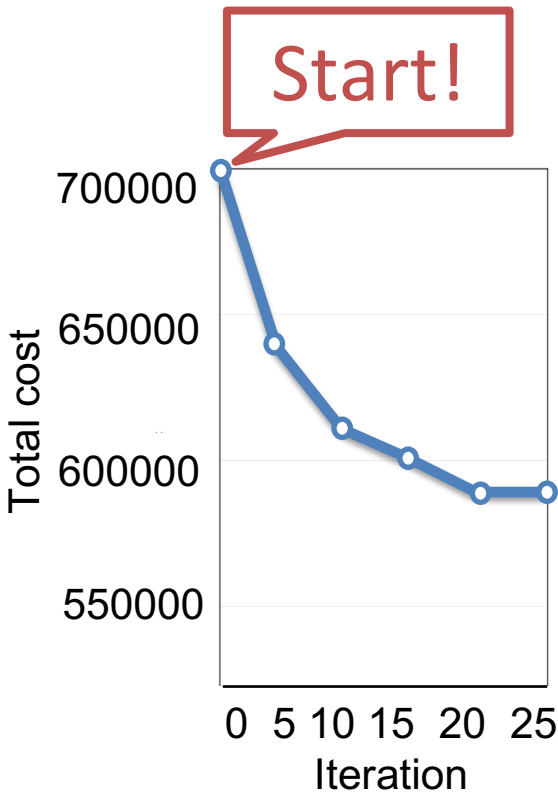split $\mathcal{X}$ into **user**-aspect

# Outline

- Motivation
- Problem definition
- Main ideas
- **Algorithms**
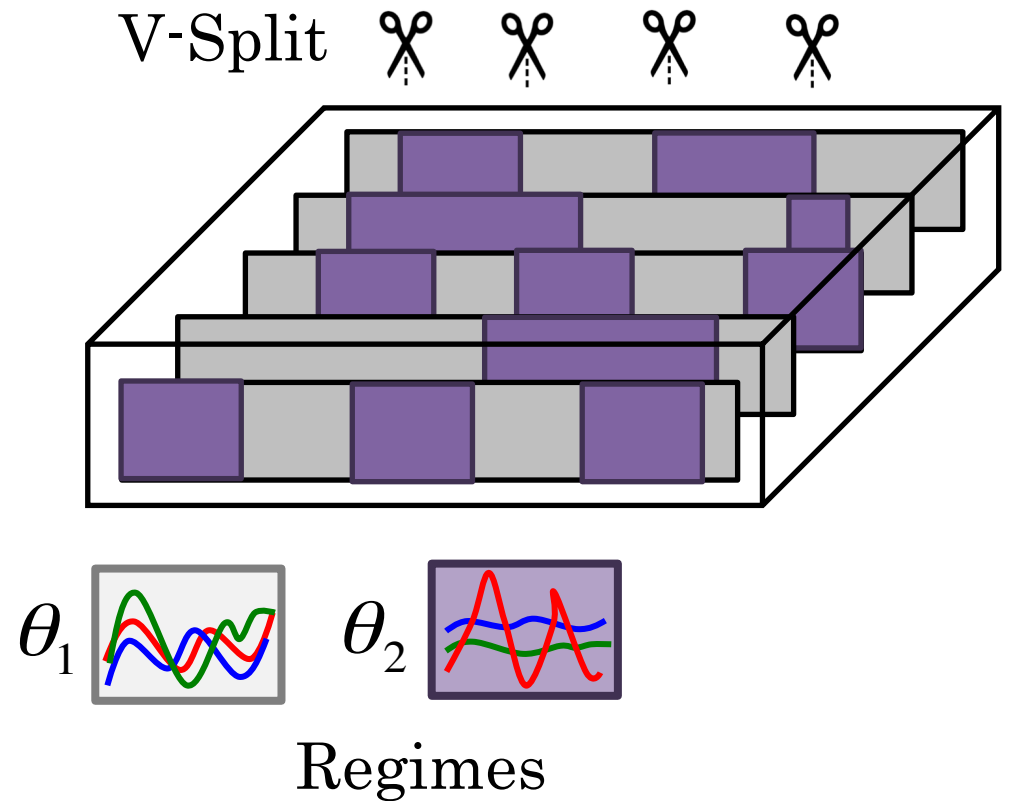- Experiments
- Conclusions

# Proposed algorithm

## Overview



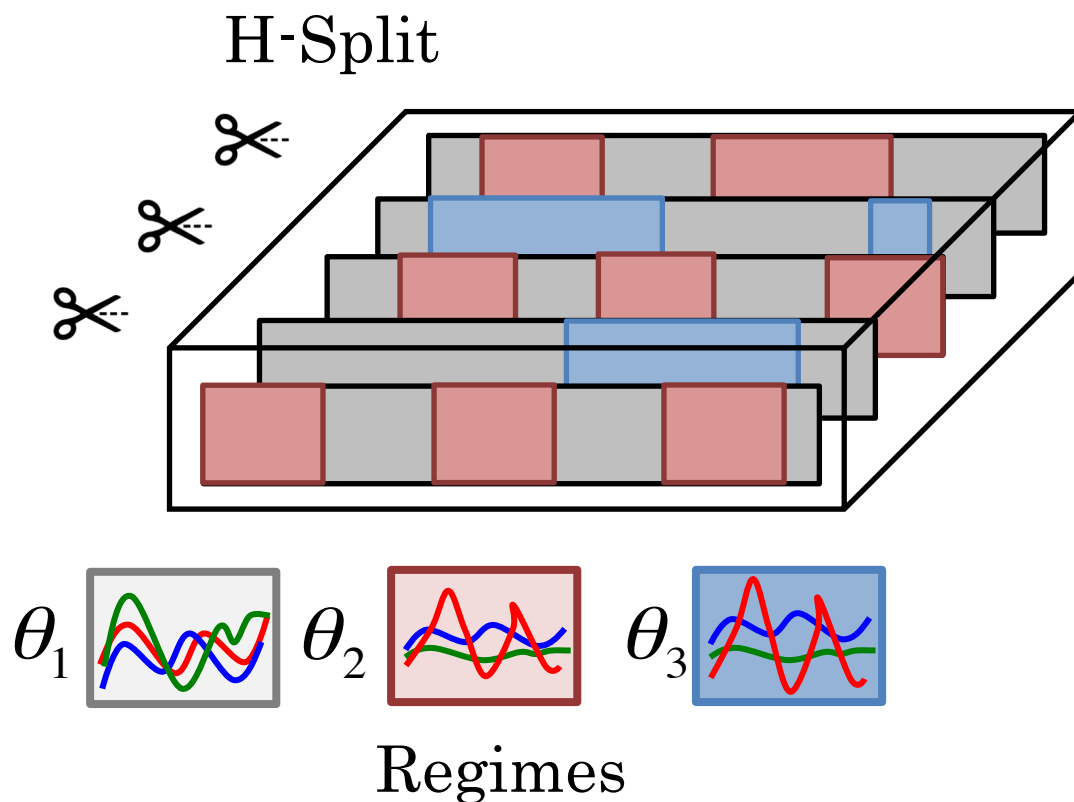Start!

Iteration 0 (r=1)

© 2019 Takato Honda et al.

# Proposed algorithm

## Overview



V-Split

$\theta_1$ $\theta_2$

Regimes

Iteration 1  (r=2)

# Proposed algorithm

## Overview

H-Split

$\theta_1$ $\theta_2$ $\theta_3$

Regimes

Iteration 2 (r=3)

© 2019 Takato Honda et al.

# Algorithms

## Algorithms of our method

CubeMarker

Outer loop

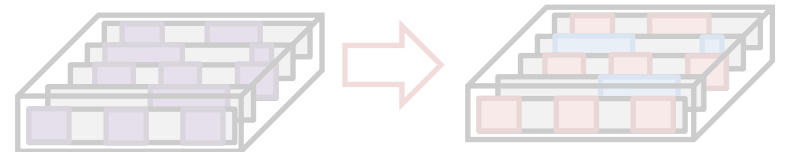**V-Split** Inner loop

- **V-Assignment**
- **ModelEstimation**



**Find time-aspect regime**

H-Split Inner loop

- H-Assignment
- ModelEstimation

Find user-aspect regime

Decide splitting algorithm
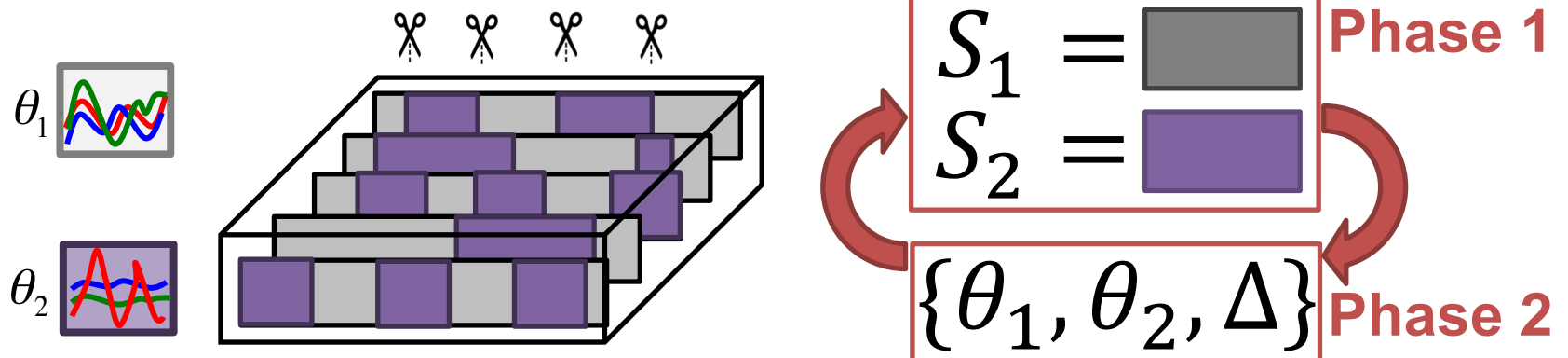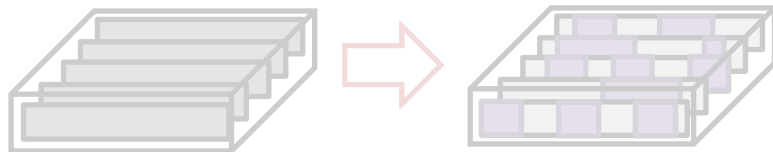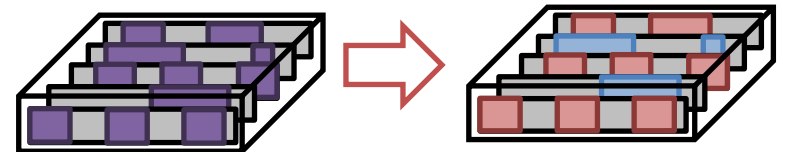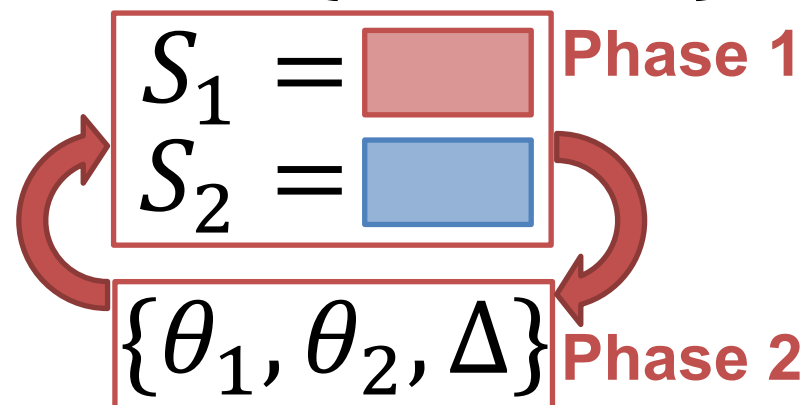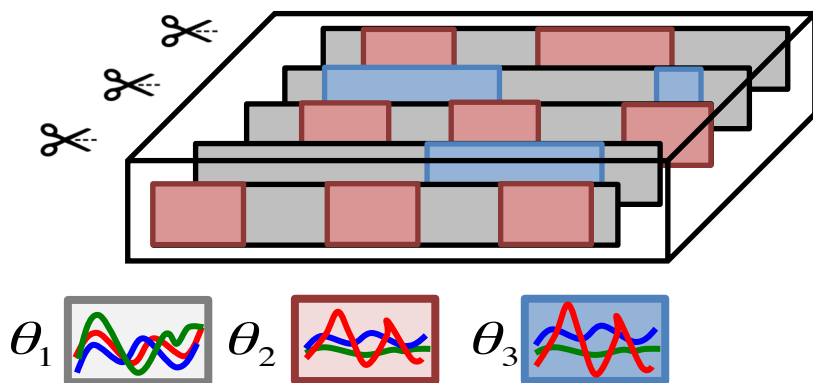
# V-Split

## Two phase iterative approach

- **Phase 1:** (V-Assignment)
  - Split segments into two groups: $S_1, S_2$

- **Phase 2:** (ModelEstimation)
  - Update model parameters: $\Theta = \{\theta_1, \theta_2, \Delta\}$



$$S_1 = \quad \text{Phase 1}$$
$$S_2 = \quad$$

$$\{\theta_1, \theta_2, \Delta\} \quad \text{Phase 2}$$

© 2019 Takato Honda et al.

# Algorithms

## Algorithms of our method

**CubeMarker**

**V-Split**  Inner loop

- V-Assignment

- ModelEstimation

Find time-aspect regime

**H-Split**  Inner loop

- **H-Assignment**

- **ModelEstimation**

**Find user-aspect regime**

*Decide splitting algorithm*

# H-Split

## Two phase iterative approach

- **Phase 1:** (H-Assignment)
  - **Split segments into two groups:** $S_1, S_2$
- **Phase 2:** (ModelEstimation)
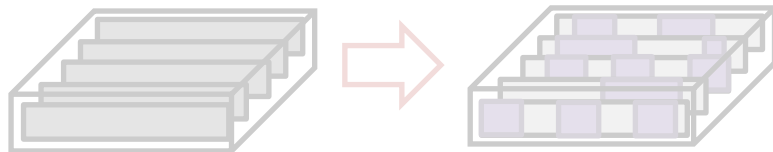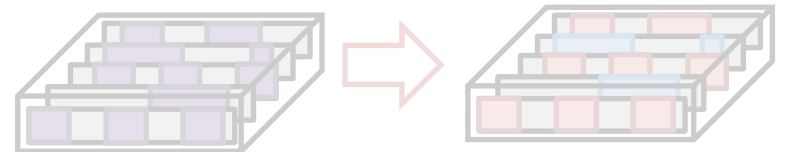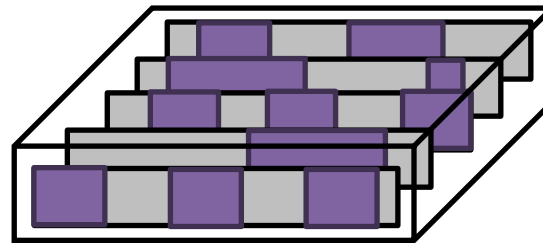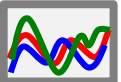  - **Update model parameters:** $\Theta = \{\theta_1, \theta_2, \Delta\}$



$\theta_1$ $\theta_2$ $\theta_3$

$S_1 =$ **Phase 1**
$S_2 =$

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

# H-Split

**Given:**

- tensor $\mathcal{X}$

- model parameter set $\quad \Theta = \{\theta_1, \theta_2, \Delta\}$

**Find:** two user-aspect regimes based on the similarity: $Cost_C(X_i|\theta_j)$

$$\mathcal{X}$$

$$\left.\{\theta_1, \theta_2, \Delta\}\right\}$$



$\theta_1$  $\theta_2$  $\theta_3$

# Algorithms

## Algorithms of our method

**CubeMarker**

**Outer loop**

**V-Split** | Inner loop

- **V-Assignment**
- **ModelEstimation**

Find time-aspect regime

**H-Split** | Inner loop

- **H-Assignment**
- **ModelEstimation**
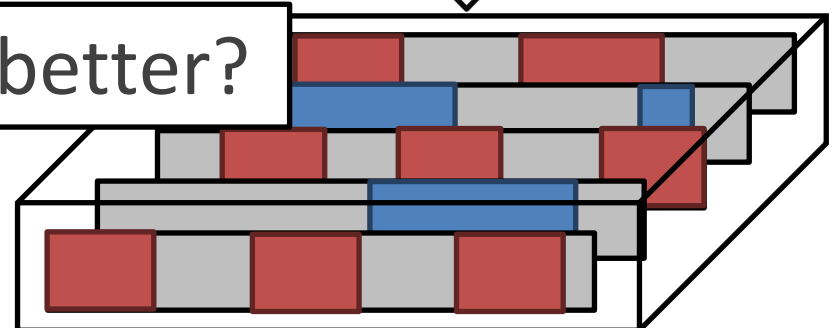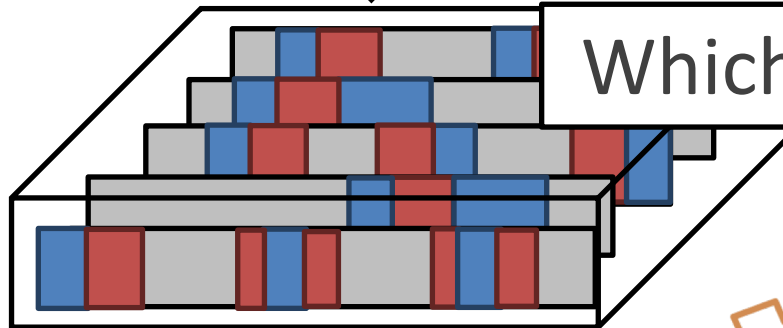
Find user-aspect regime
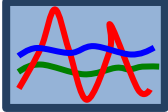
**Decide splitting algorithm**

# CubeMarker

$\theta_1$ [graph] $\theta_2$ [graph] Regimes

Tensor (cost: 687,395)

Which is better?

$\theta_1$ [graph] $\theta_2$ [graph] $\theta_3$ [graph]

V-Split result (cost: 673,255) vs.

$\theta_1$ [graph] $\theta_2$ [graph] [graph]

H-Split result (cost: 642,441)

© 2019 Takato Honda et al.

# CubeMarker

$\theta_1$ [graph] $\theta_2$ [graph] Regimes

Tensor (cost: 687,395)

$\theta_1$ [graph] $\theta_2$ [graph] $\theta_3$ [graph]

V-Split result (cost: 673,255) vs. H-Split result (cost: 642,441)

$\theta_1$ [graph] $\theta_2$ [graph] [graph]

© 2019 Takato Honda et al.

# Outline

- **Motivation**
- **Problem definition**
- **Main ideas**
- **Algorithms**
- **Experiments**
- **Conclusions**

# Experiments

## Q1. Effectiveness
**Can it help us understand the given tensor?**

## Q2. Scalability
**How does it scale in terms of computational cost?**

## Q3. Accuracy
**How well does it find segments and regimes?**

Competitors:

pHMM (SIGMOD'11)

AutoPlait (SIGMOD'14)

TICC (KDD'17)
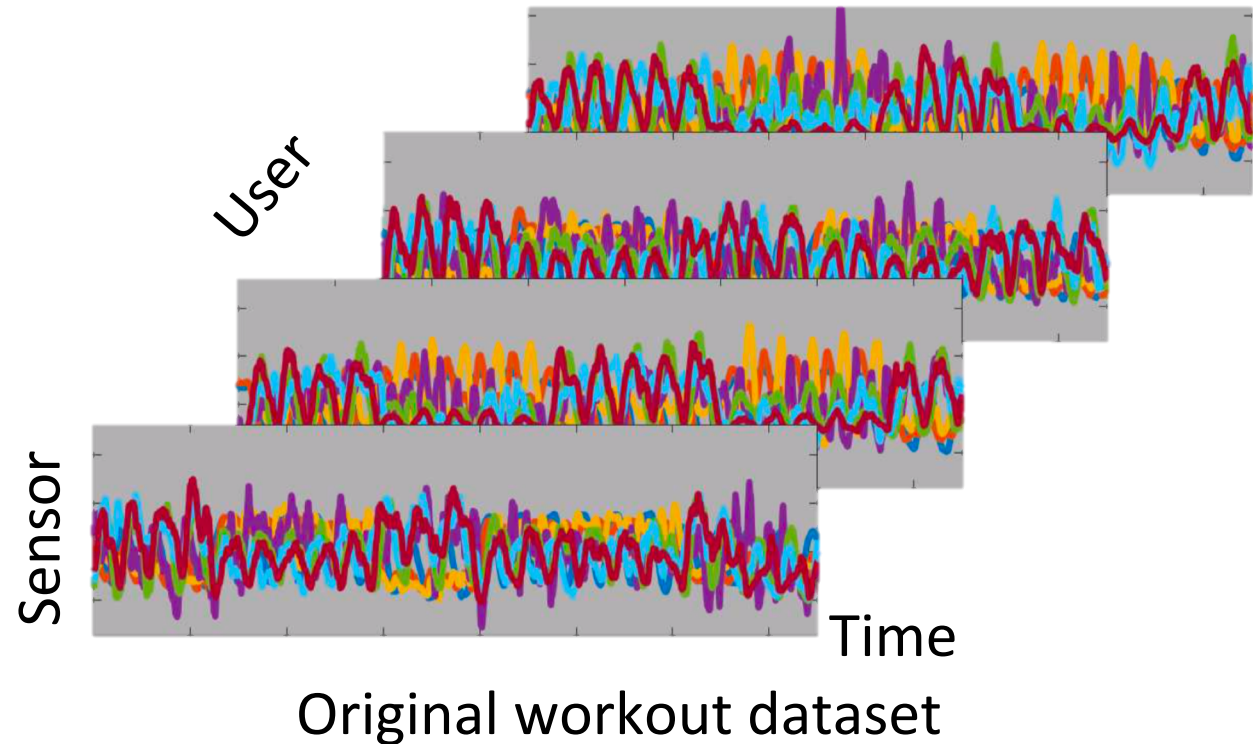
CubeMarker-V (naïve ver. of our method)

# Datasets

**Experiments on the 8 real-world datasets:**

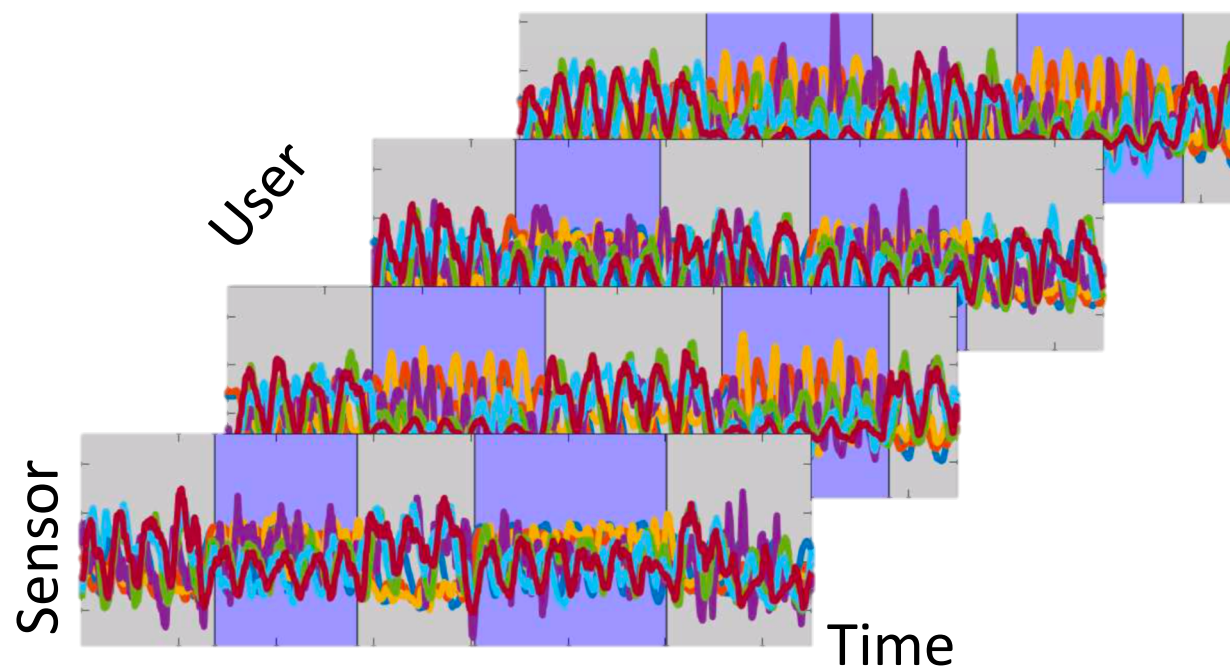| Dataset | Data size ($w \times n \times d$) |
|---|---|
| (#1) Workout | $182 \times 4000 \times 7$ |
| (#2) Tennis | $100 \times 4500 \times 7$ |
| (#3) Factory | $60 \times 3000 \times 7$ |
| (#4) Reading | $71 \times 10000 \times 5$ |
| (#5) Free throw | $170 \times 2000 \times 7$ |
| (#6) Automobile-Tokyo | $171 \times 2400 \times 3$ |
| (#7) Automobile-Expressway | $13 \times 9100 \times 3$ |
| (#8) Automobile-Togu | $32 \times 5200 \times 3$ |

Summary of the datasets
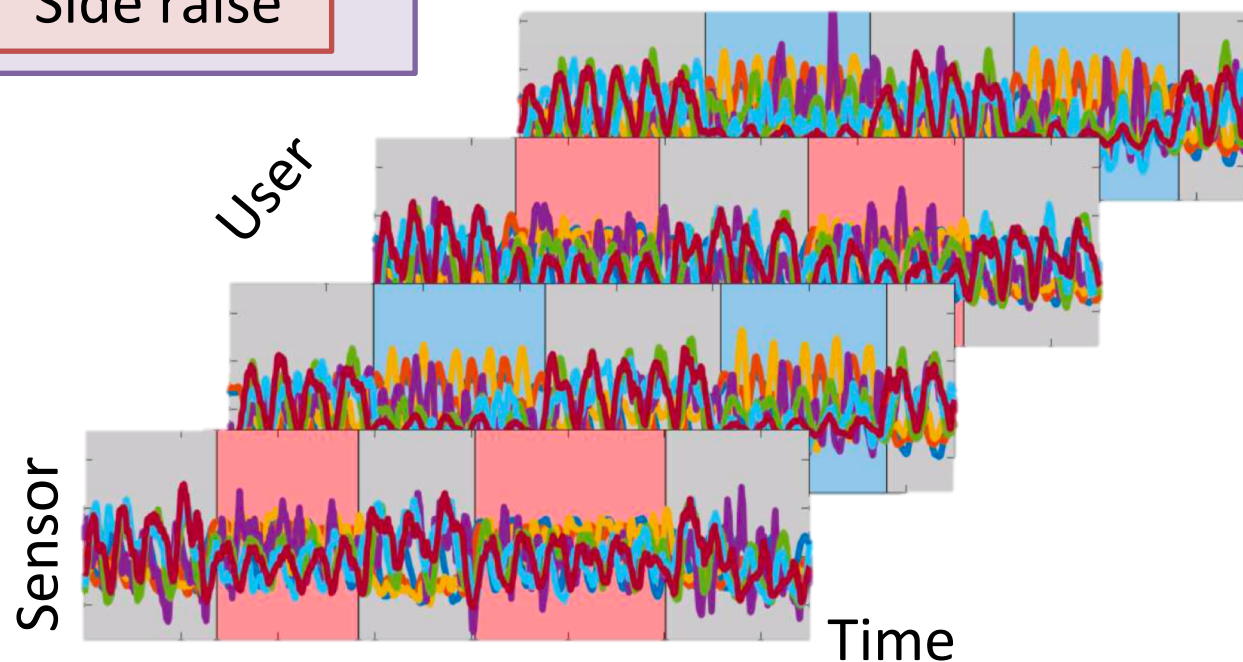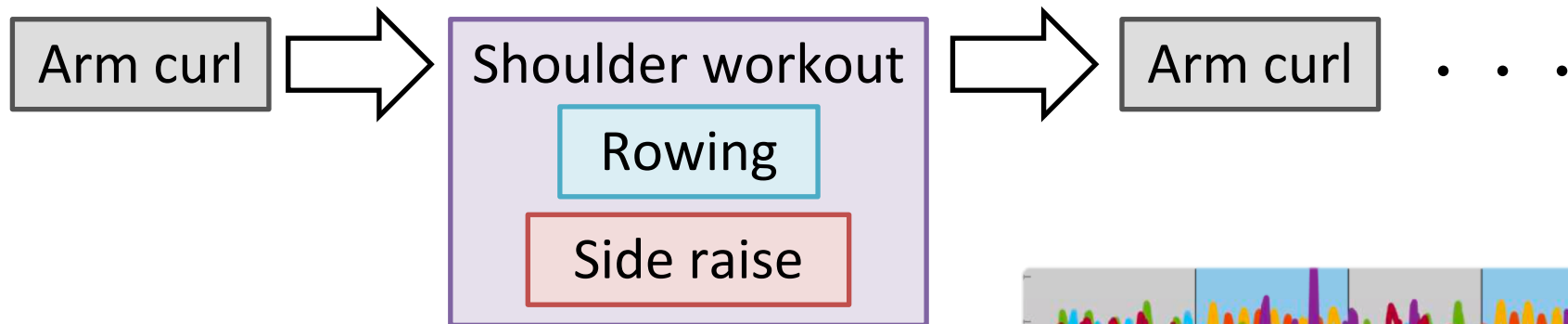
How many and what kind of patterns does it include?



Original workout dataset

Arm curl $\Rightarrow$ Shoulder workout $\Rightarrow$ Arm curl $\cdot\ \cdot\ \cdot$



Time-aspect patterns for a workout dataset

# Q1. Effectiveness - Workout

Arm curl $\Rightarrow$ Shoulder workout $\Rightarrow$ Arm curl $\cdot$ $\cdot$ $\cdot$

Rowing

Side raise

User

Sensor

Time

Multi-aspect patterns for a workout dataset

Basic pattern transitions: carrying ⟹ assembling ⟹ • • •

User-aspect pattern: Discarding defective products

One-shot outlier: stretch arms



500    1000    1500    Time

Regimes | θ1 | θ2 | θ3 | θ4

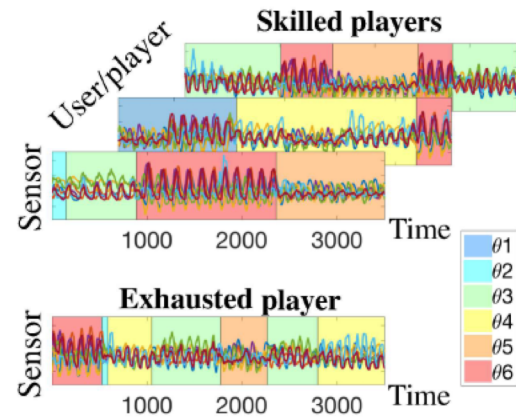θ1: discard  θ2: stretch  θ3: carry  θ4: assemble

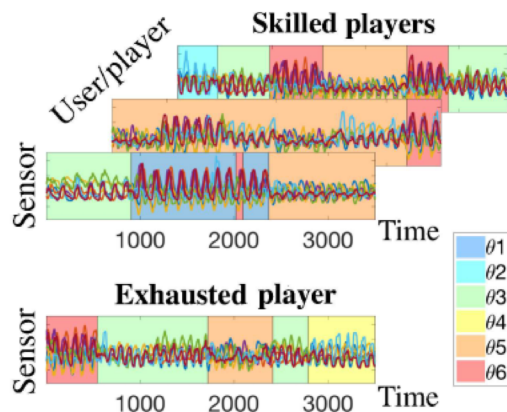Multi-aspect patterns for a factory workers
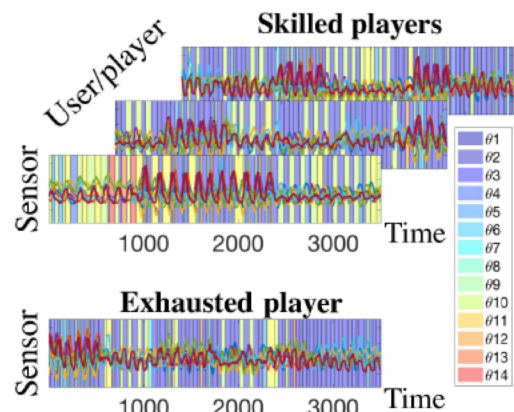
(a) **CubeMarker**
(**no** parameter setting)

(b-1) TICC ($\beta = 100, \lambda = 1000$)
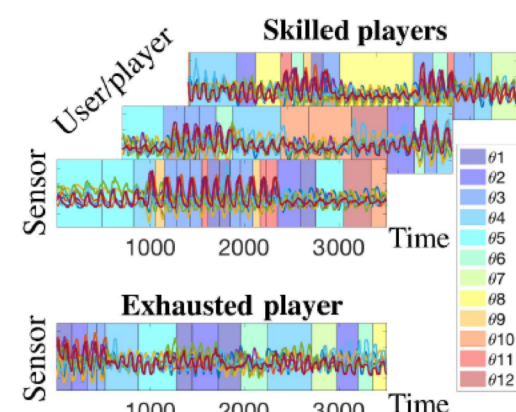(**need** parameter setting)

(b-2) TICC ($\beta = 600, \lambda = 1000$)
(**need** parameter setting)

(c) AutoPlait
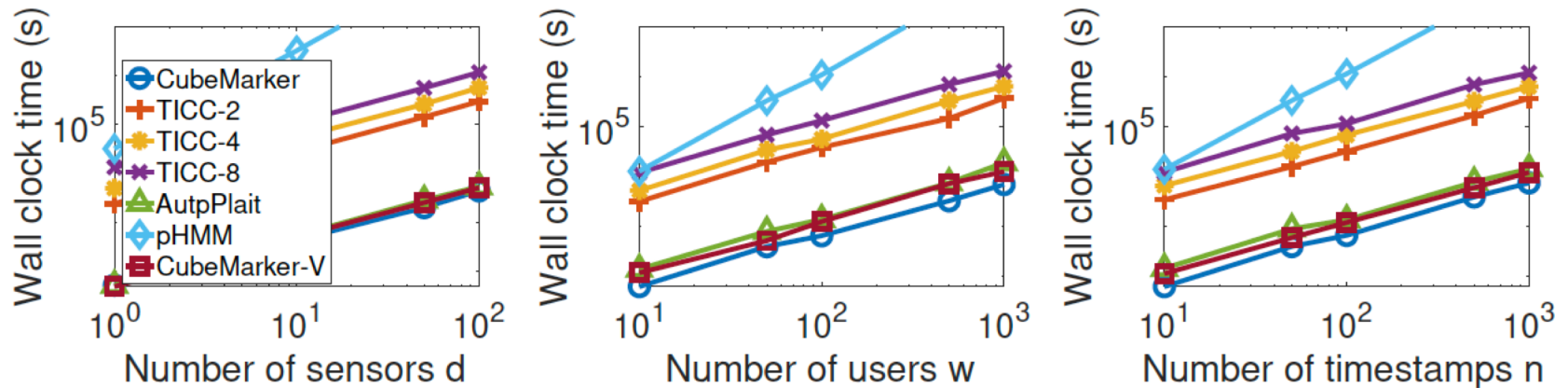(**no** parameter setting)

(d-1) pHMM ($\epsilon_r = 0.1, \epsilon_c = 0.8$)
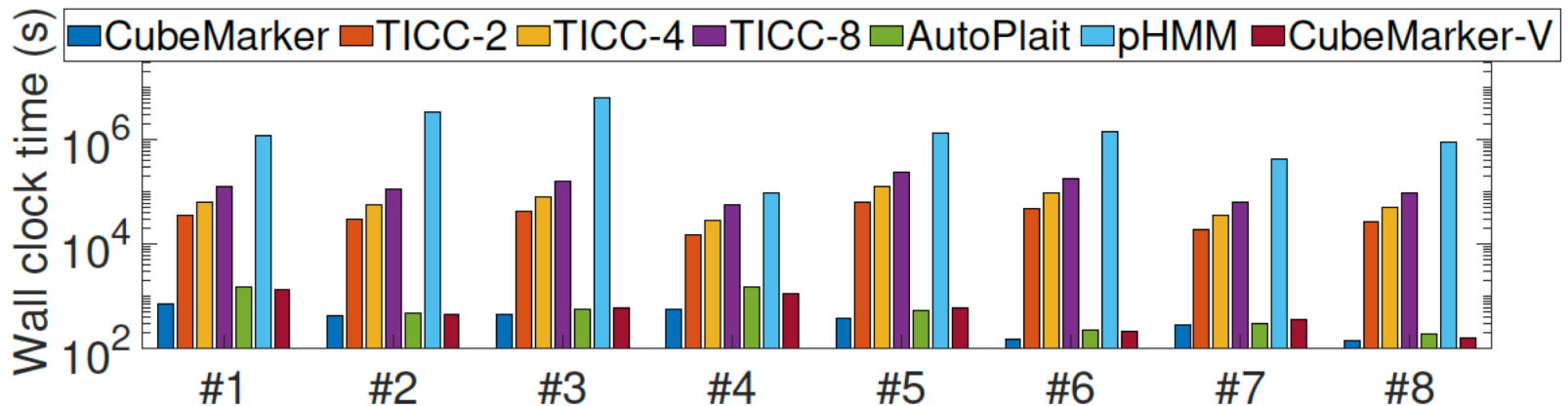(**need** parameter setting)

(d-2) pHMM ($\epsilon_r = 10, \epsilon_c = 0.8$)
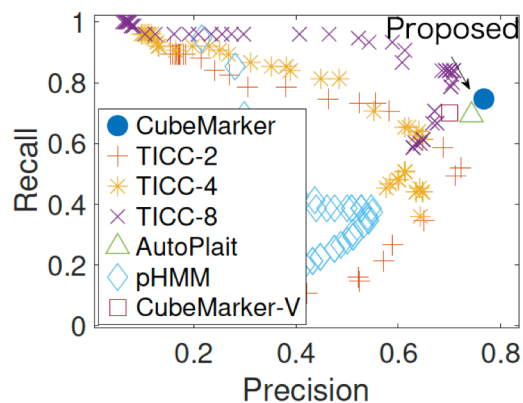(**need** parameter setting)

# Q2. Scalability

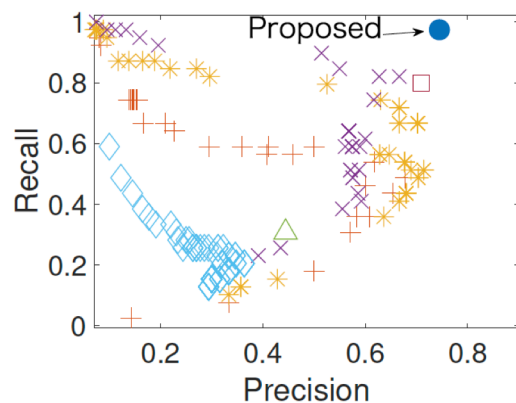Wall clock time v.s. dataset size for (#1) Workout ($O(dwn)$)
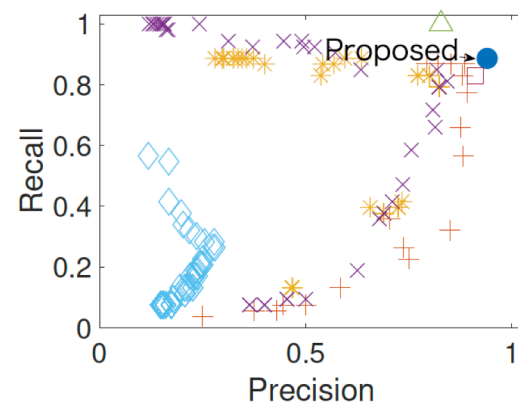


Wall clock time for each dataset (1700x faster than pHMM)

© 2019 Takato Honda et al.

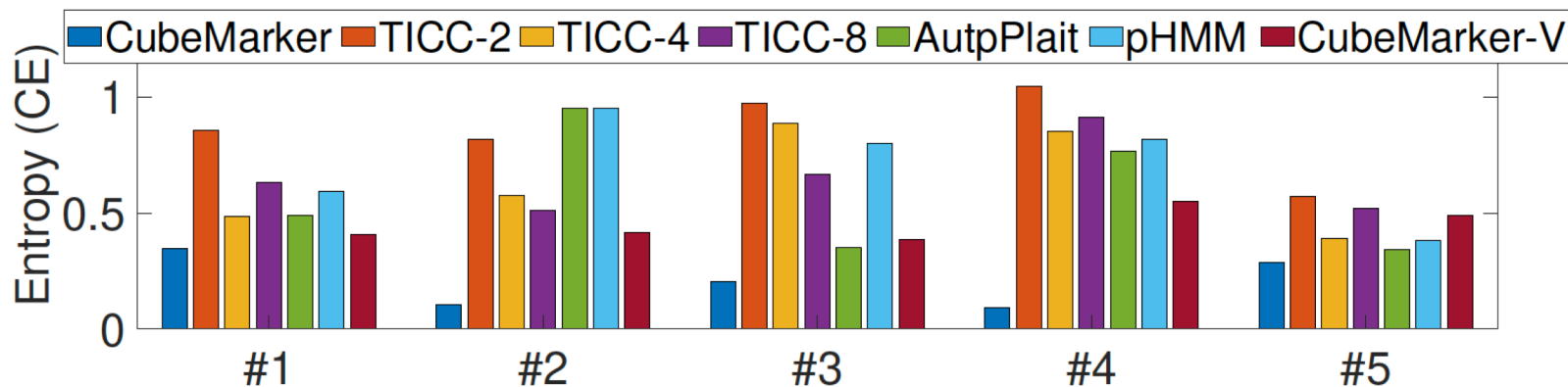# Q3. Accuracy (segment/regime)



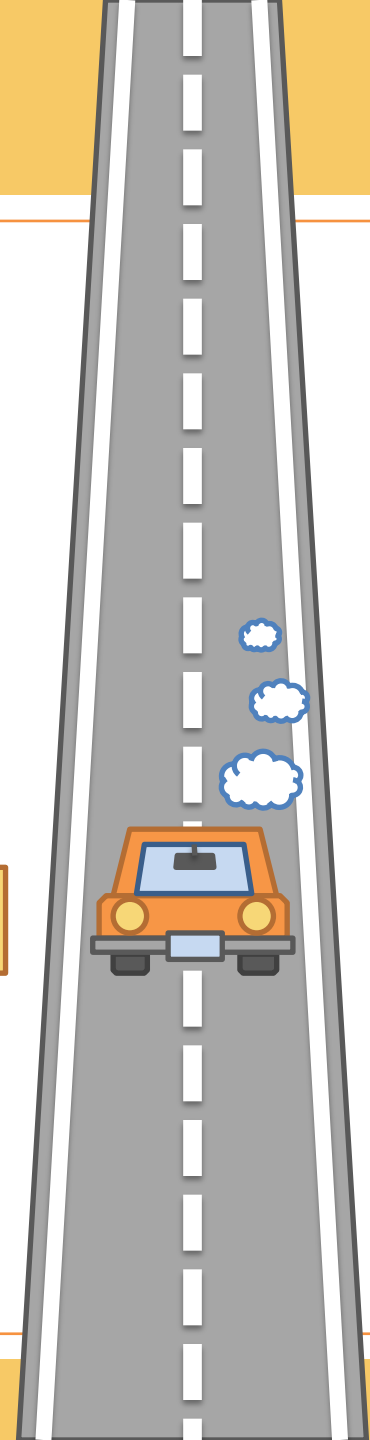(a) (#1) Workout     (b) (#2) Tennis     (c) (#3) Factory

Segmentation accuracy (top righter is better)
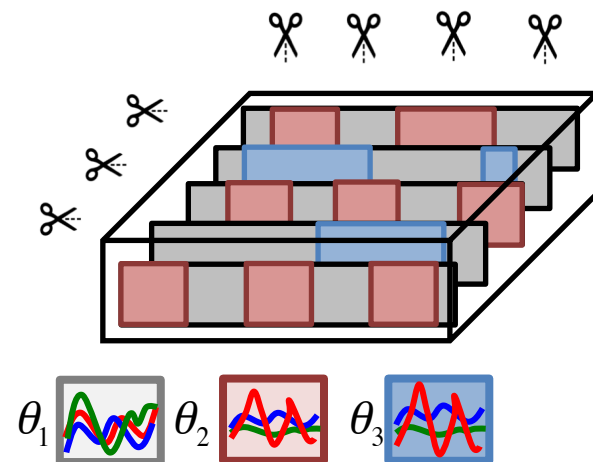


Regime clustering accuracy (lower is better)

# Outline

- **Motivation**
- **Problem definition**
- **Main ideas**
- **Algorithms**
- **Experiments**
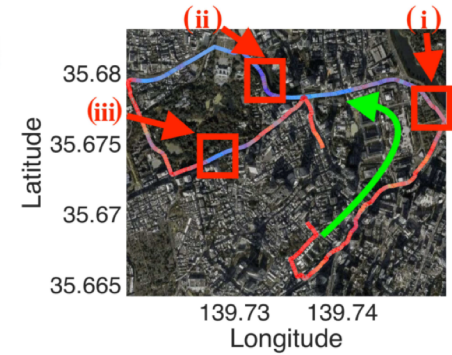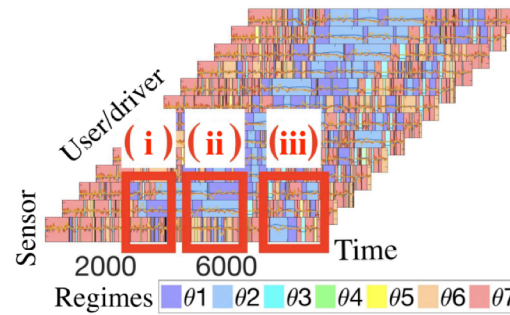- **Conclusions**

© 2019 Takato Honda et al.

# Conclusions

Our method has the following properties:

- **Effective**
  Find multi-aspect segments/regimes

- **Automatic**
  No magic numbers

- **Scalable**
  It scales linearly to the data size
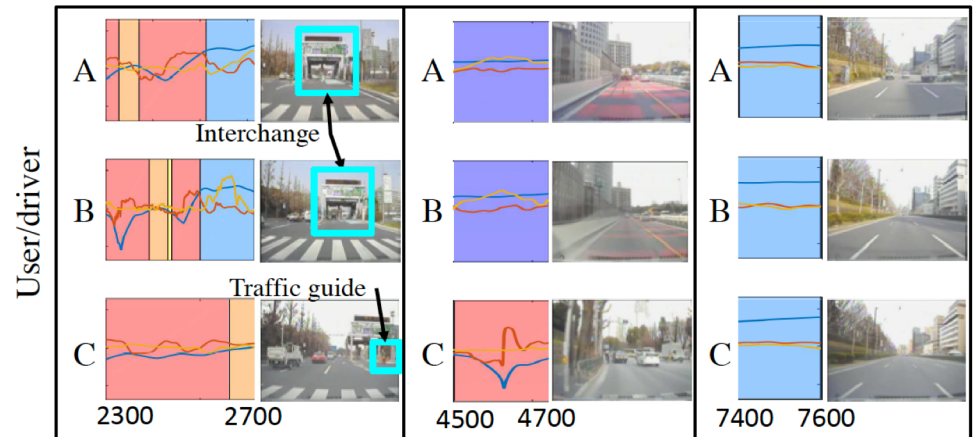
© 2019 Takato Honda et al.

# Thank you!



(a) Multi-aspect segmentation and summarization

(b) Representative driving behavior on a map

(c-i) Interchange (c-ii) Expressway (c-iii) Wide road
(c) User/driver-specific behavior at three different locations