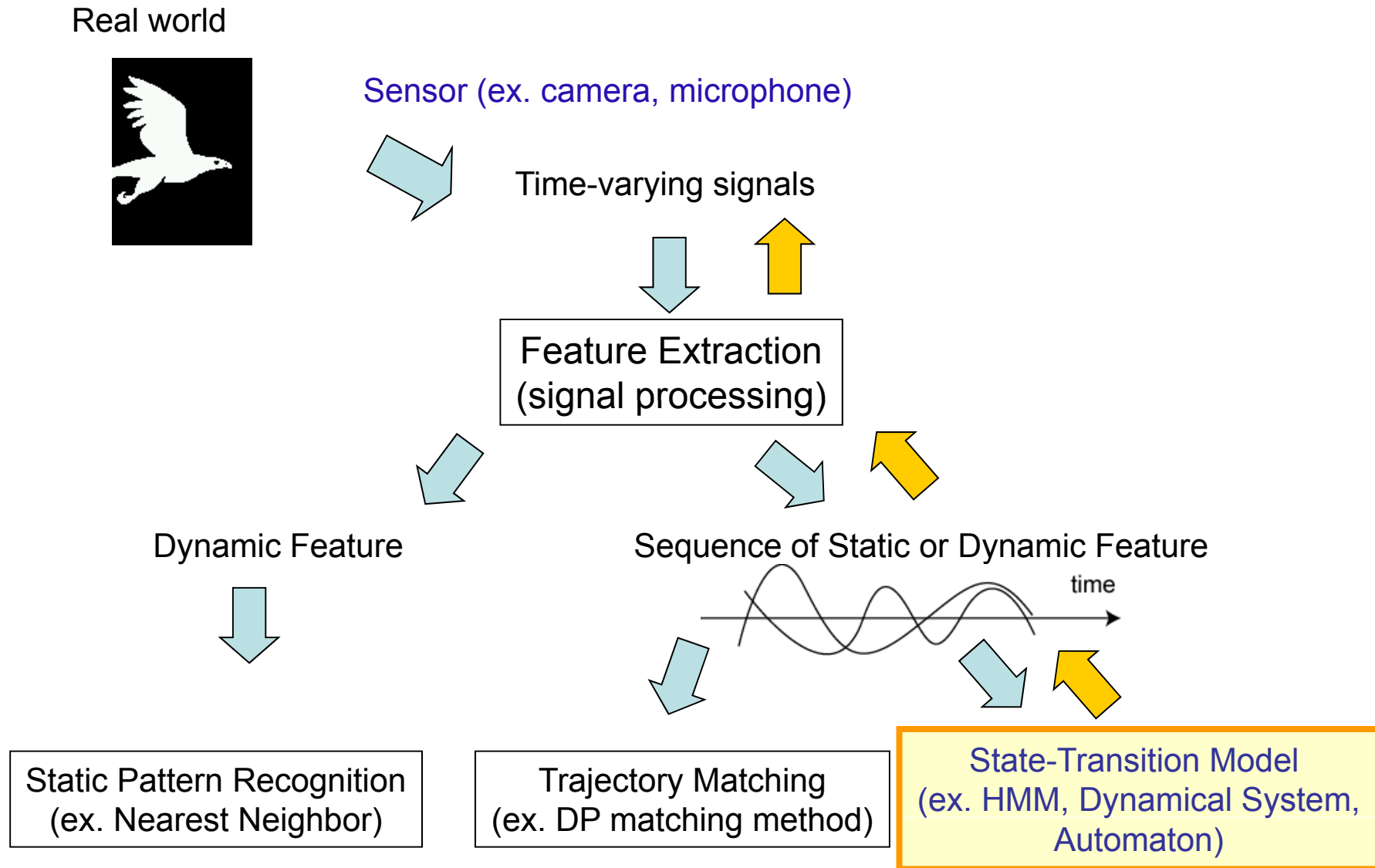




Interval-Based Hybrid Dynamical System for Modeling Dynamic Events and Structures

Hiroaki Kawashima

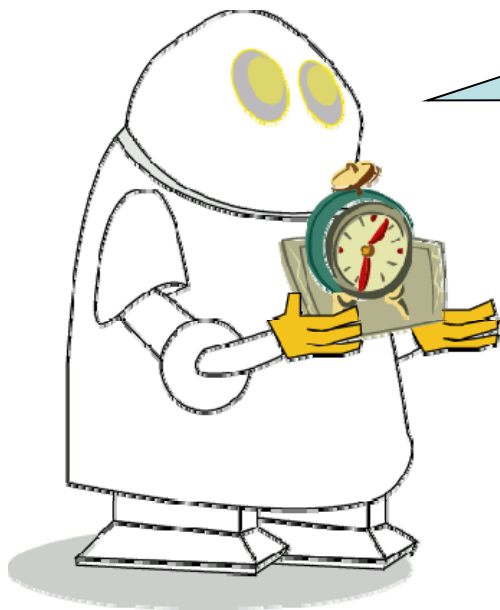
Event Recognition



Two Concepts of Time

Time flies
faster than it
used to...

?



Subjective time (*Kairos*)

Time has no metric: (\mathbb{N} , \leq)

Q: state set, E: input event set

Event sequence

S_i : Countable set $\mathbb{N} \rightarrow$ Finite set A

$A = \{\text{up}, \text{down}\}$



up \rightarrow down \rightarrow up \rightarrow down \dots
 0 1 2 3 \dots

Order of discrete events

Discrete-Event Systems

$M(\text{State}_{\text{now}}, \text{input event}) = \text{State}_{\text{next}}$

$M: Q \times A \rightarrow Q$

Automata, Petri nets
(Turing machine 1936)

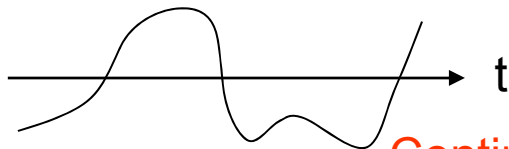
Objective (Physical) time (*Chronos*)

Time has metric : (\mathbb{R} , \leq , *dist*)

\mathbb{R}^n : state space

Signal data

$y(t)$: Real num. $\mathbb{R} \rightarrow$ Continuous space \mathbb{R}^k



Continuous change

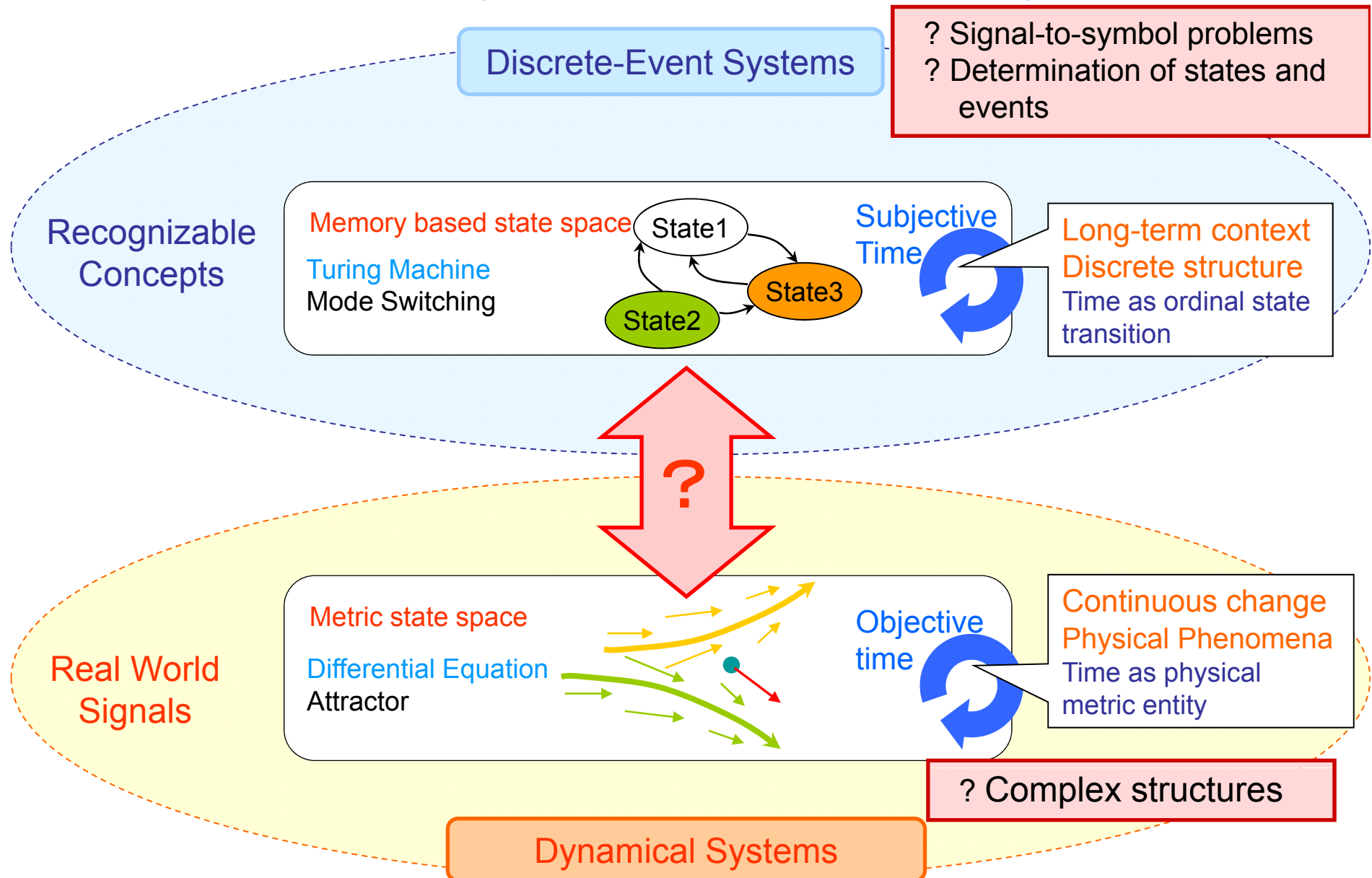
Dynamical Systems

$$\frac{d \text{State}(t)}{dt} = F(\text{state}(t))$$

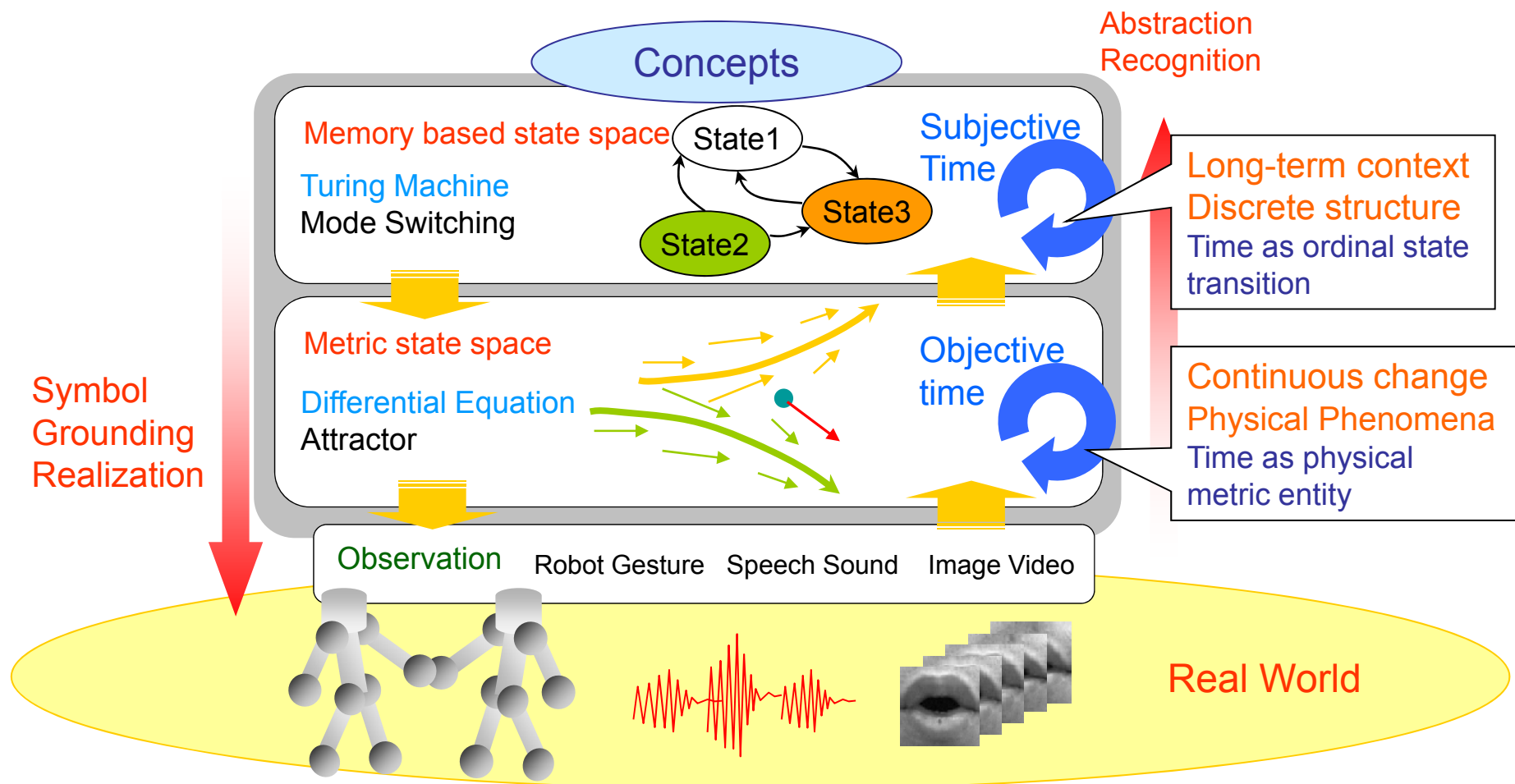
$F: \mathbb{R}^n \rightarrow \mathbb{R}^n$

Control systems, Neural networks
(Cybernetics 1947)

Advantages and Disadvantages



Hybrid Dynamical Systems



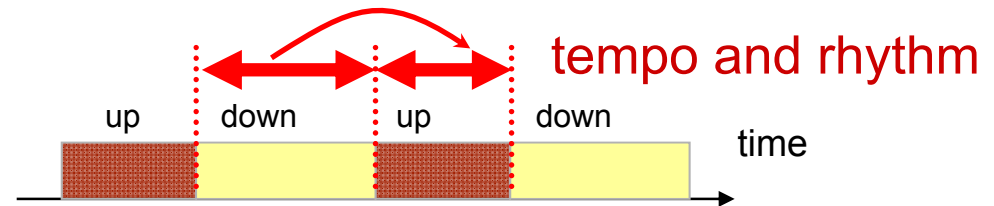
Existing Studies

- Computer vision
 - Hybrid dynamical models [C. Bregler 1997]
 - Multi-class condensation [B. North, A. Blake, M. Isard and J. Rittscher, 2000]
 - Switching linear dynamical systems [K.P. Murphy 1998, V. Pavlovic 1999]
- Speech recognition
 - Segment models [M. Ostendorf 1996]
- Computer graphics
 - Motion textures [Y. Li, T. Wang, H.Y. Shum 2002]
- Neural networks, Control theory, etc.
 - Piecewise linear models [R. Batruni 1991]
 - Switching space models [Z. Ghahramani 1996]
 - Piecewise affine maps [L. Breiman 1993]
 - Hybrid automata [R. Alur 1993]

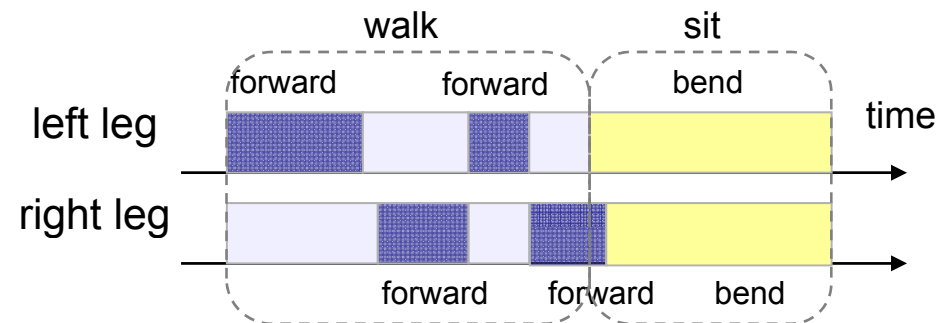
Integration of “subjective time” and “objective time”?

Dynamic Structures of Events

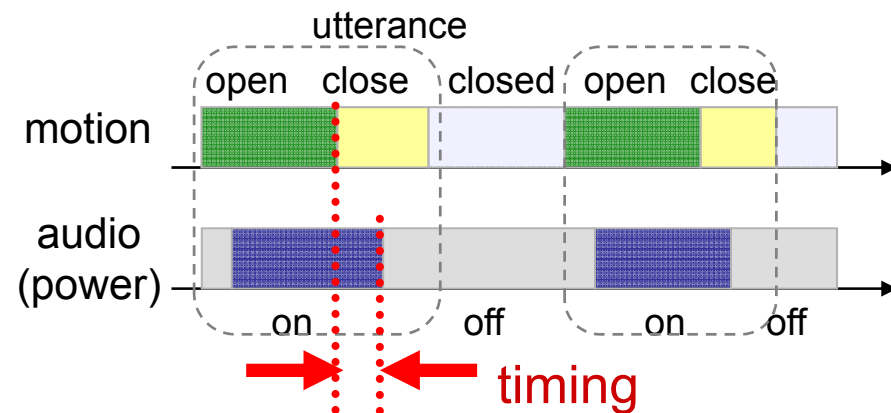
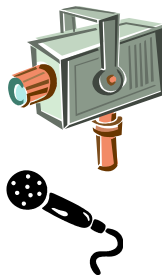
Sequential structures



Constituted by multiple objects

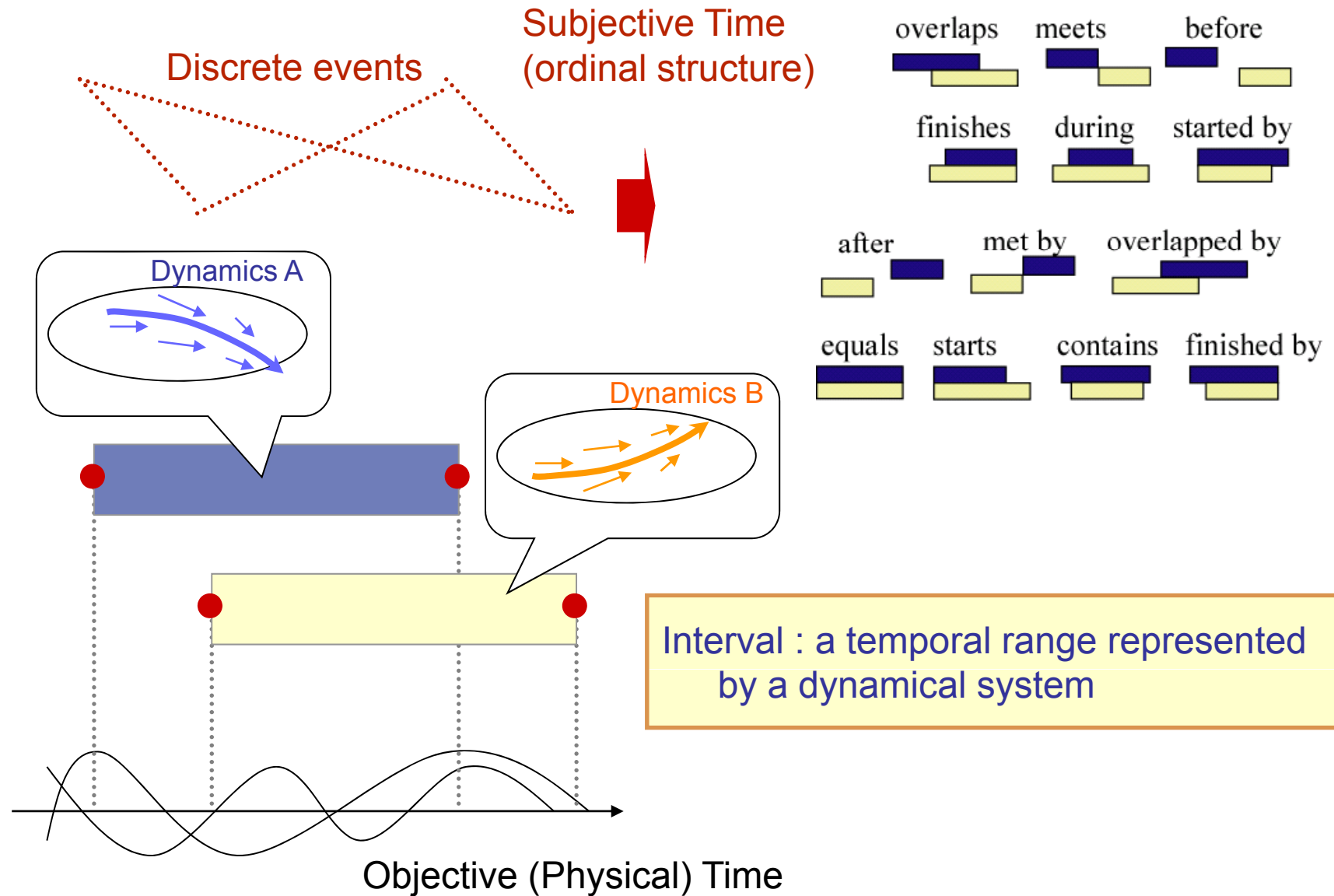


Observed as different media signals

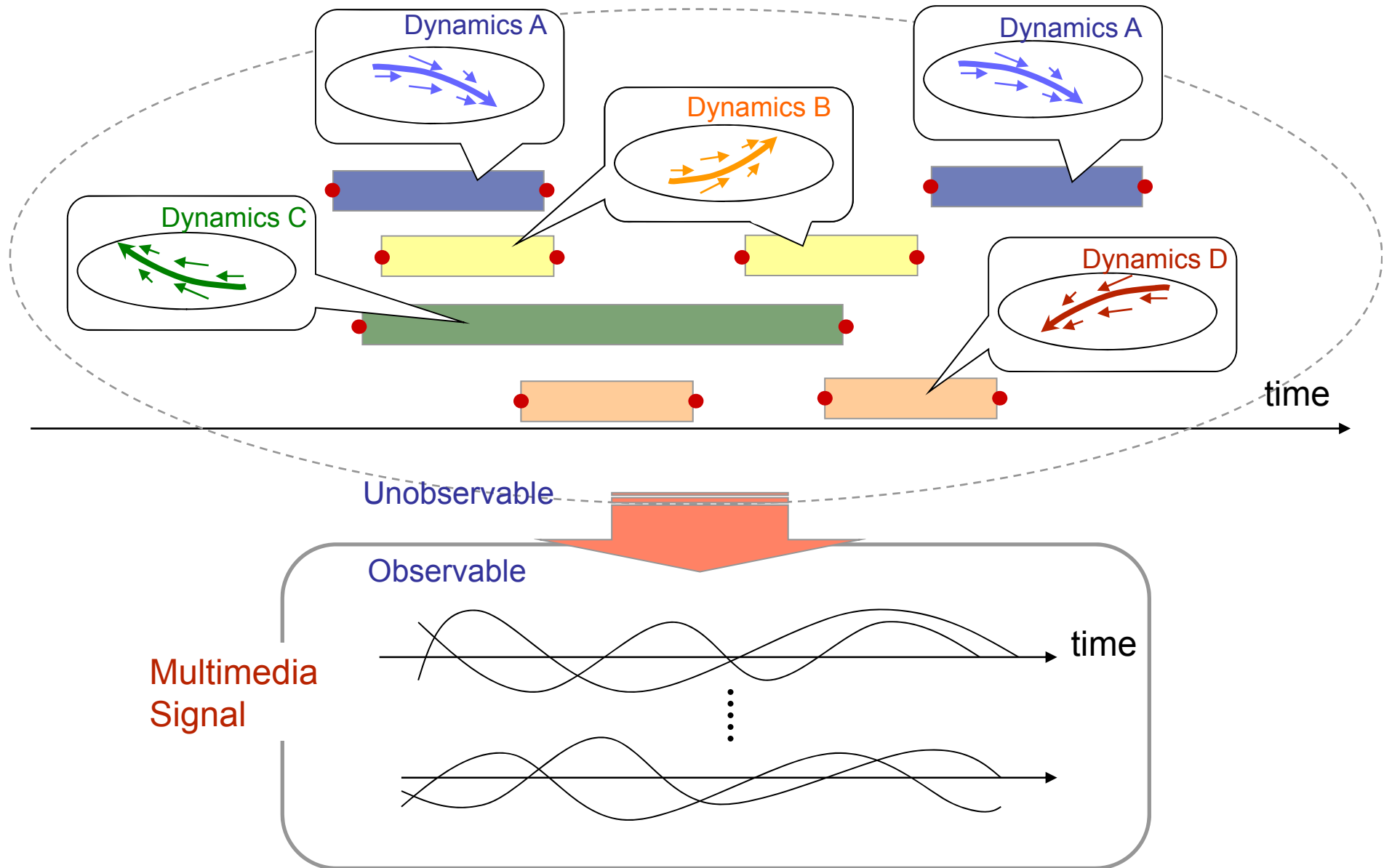


How to represent “a sense of time” of human?

Temporal Intervals



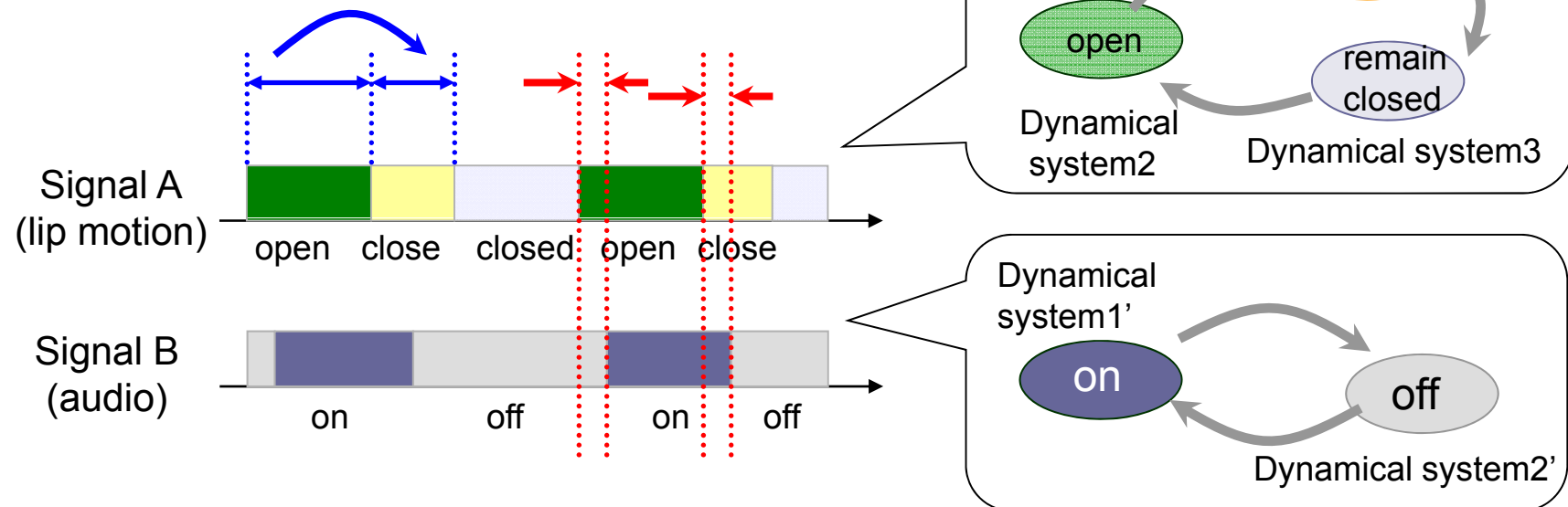
Orchestration of Dynamics



Interval-Based Hybrid Dynamical System

represents complex temporal structures with physical-time grounding

1. Interval-based state transition
→ to model rhythm and tempo of a single signal patterns
2. Timing structure model
→ to model timing structure in multiple signal patterns
3. Clustering of dynamical systems
→ to find a set of dynamics



Overview of the Thesis

Modeling Single-Channel Signals (Segmentation, Tempo, Rhythm)

Modeling Multi-Channel Signals (Timing Structure)

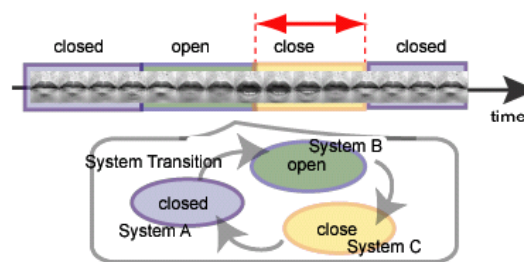
Single-Media Signal

Multimedia Signal

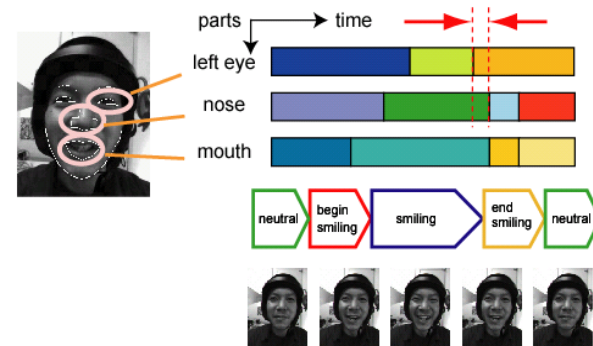
Modeling Single-Part Events in Signal Modality

Chapter 2
Architecture and the Inference Algorithm of
an Interval-Based Hybrid Dynamical System

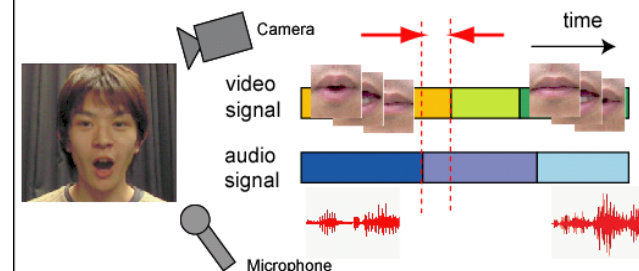
Chapter 3
Learning Algorithm of
an Interval-Based Hybrid Dynamical System



Chapter 4 Analysis of Multipart Events in Single Modality Based on Interval-Based Hybrid Dynamical Systems



Chapter 5 Modeling Multimodal Events Based on Interval-Based Hybrid Dynamical Systems

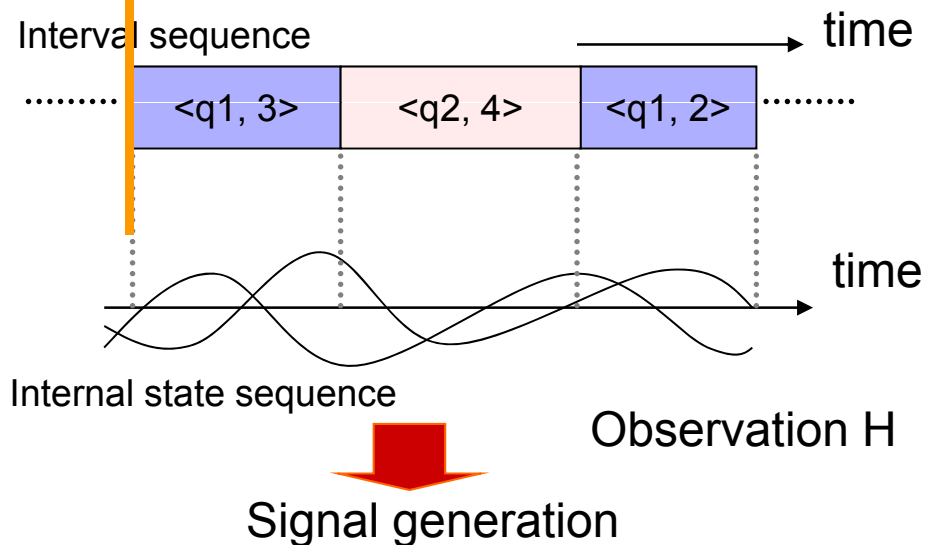
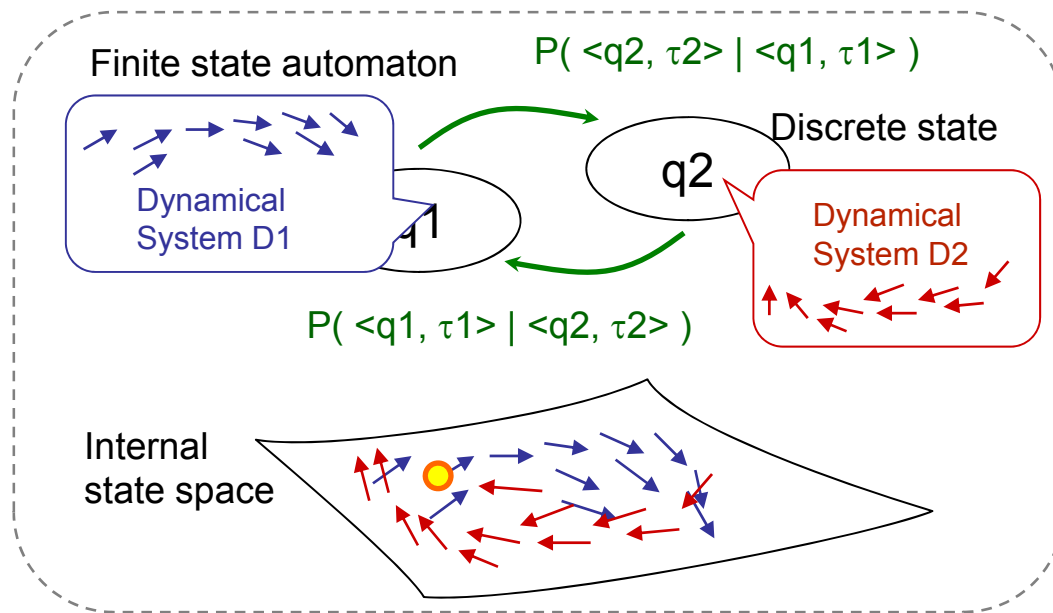




Chapter 2

Interval-Based Hybrid Dynamical System

Interval-Based Hybrid Dynamical System



Interval-based transition

$$P(I_k = <q_j, \tau> | I_{k-1} = <q_i, \tau_p>)$$

Intervals

$$<\underbrace{q_j}_{\text{state}}, \underbrace{\tau}_{\text{duration}}>$$

Linear dynamical systems

State transition

$$x_t = F^{(i)} x_{t-1} + g^{(i)} + w_t^{(i)}$$

Observation

$$y_t = H x_t + v_t$$

Linear Dynamical Systems

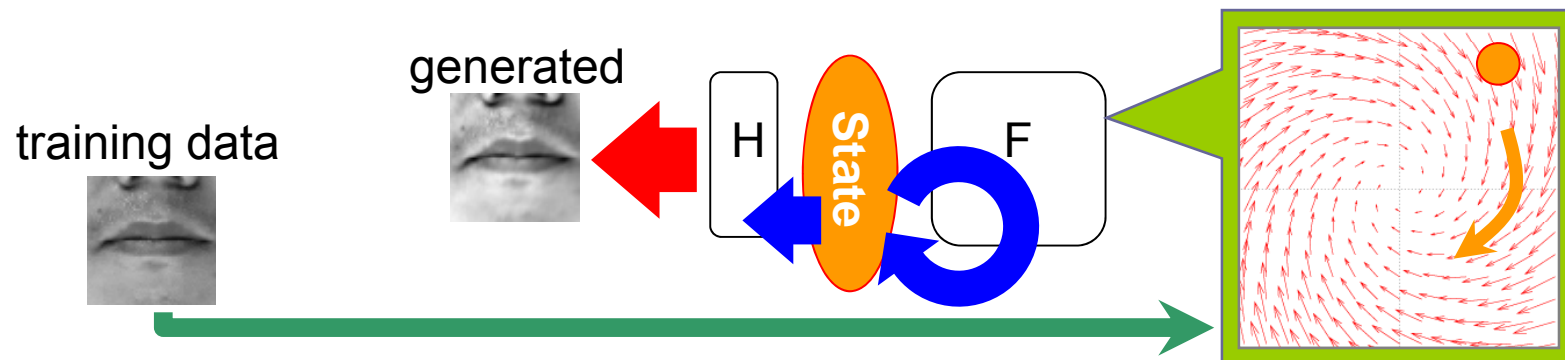
- Internal state $x \in \mathbb{R}^n$
- State transition
- Observation

$$x_t = \underbrace{F^{(i)}}_{\text{transition matrix}} x_{t-1} + \underbrace{g^{(i)}}_{\text{bias}} + \underbrace{w_t^{(i)}}_{\text{process noise}}$$

$$y_t = \underbrace{H}_{\text{observation matrix}} x_t + \underbrace{v_t}_{\text{observation noise}}$$

($w_t, v_t \sim \text{Gaussian}$)

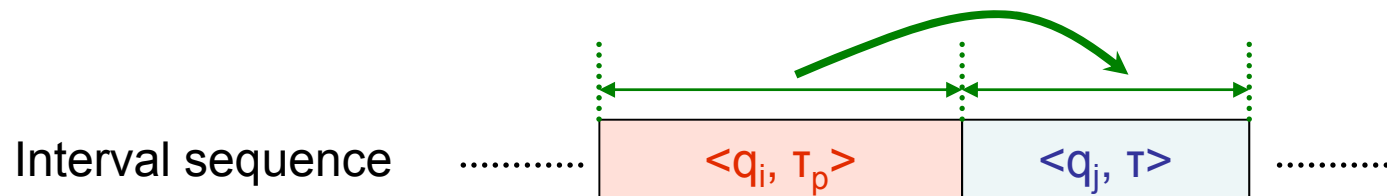
- Parameters: { F, g, H, noise covariance matrices, initial state }
- Generation of continuous change



Interval-Based State Transition

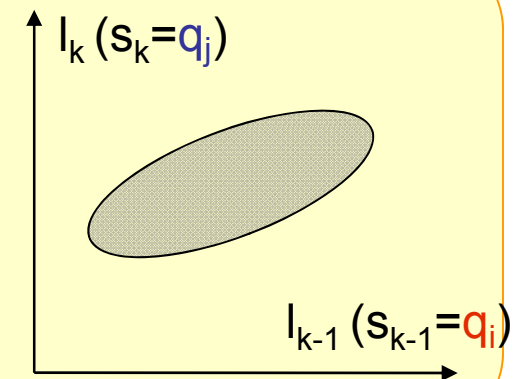
- Modeling relation of duration lengths between adjacent intervals
 - to represent tempo and rhythm
 - to enhance robustness of segmentation (top-down constraints)

$$P(I_k = \langle q_j, \tau \rangle \mid I_{k-1} = \langle q_i, \tau_p \rangle) = P(s_k = q_j \mid s_{k-1} = q_i) P(l_k = \tau \mid s_k = q_j, s_{k-1} = q_i, l_{k-1} = \tau_p)$$



$$P(l_k = \tau, l_{k-1} = \tau_p \mid s_k = q_j, s_{k-1} = q_i)$$

Assume Gaussian distribution

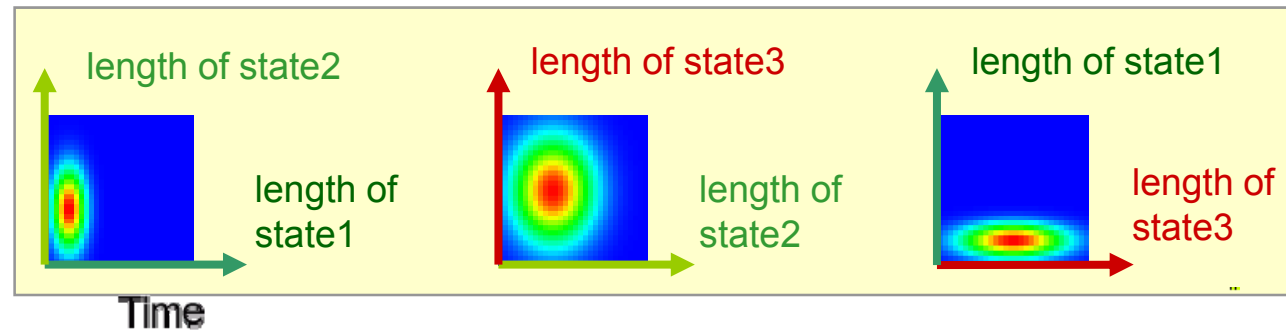


Example of Interval Sequence Generation

- Randomly generated interval sequences using manually given distributions

Three discrete states

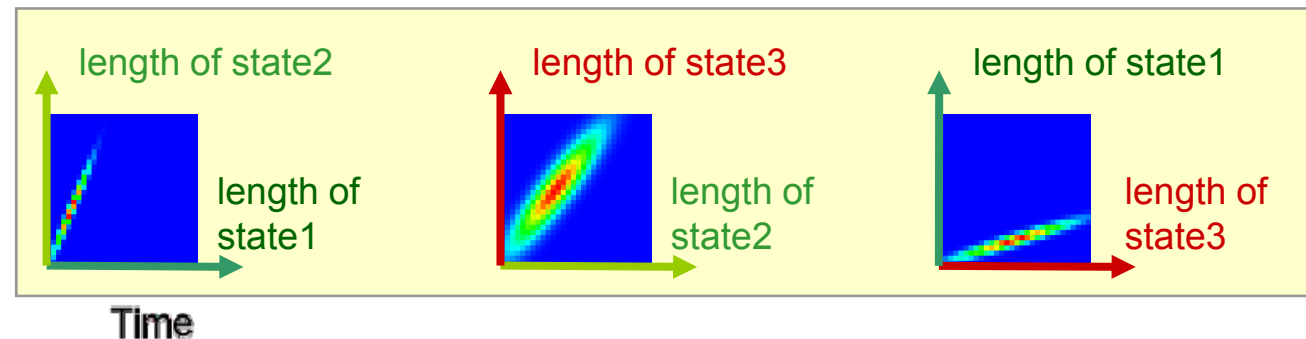
Correlation = 0



Generated Sequence



Correlation $\doteq 0.9$



Generated Sequence





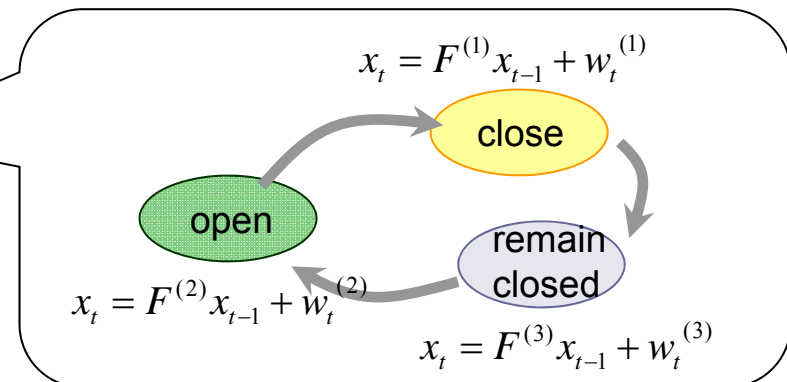
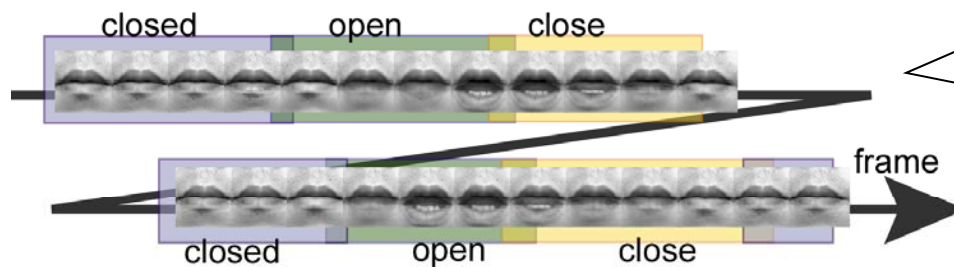
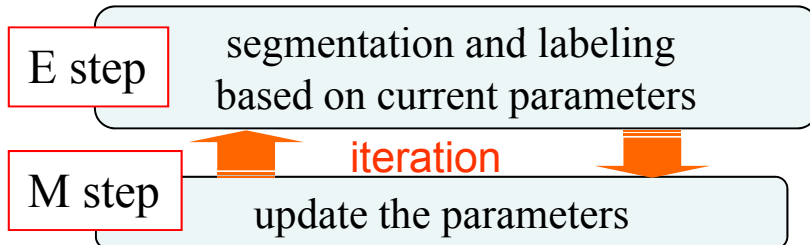
Chapter 3

Learning Method for the Interval-Based Hybrid Dynamical System

Difficulty of the Parameter Estimation

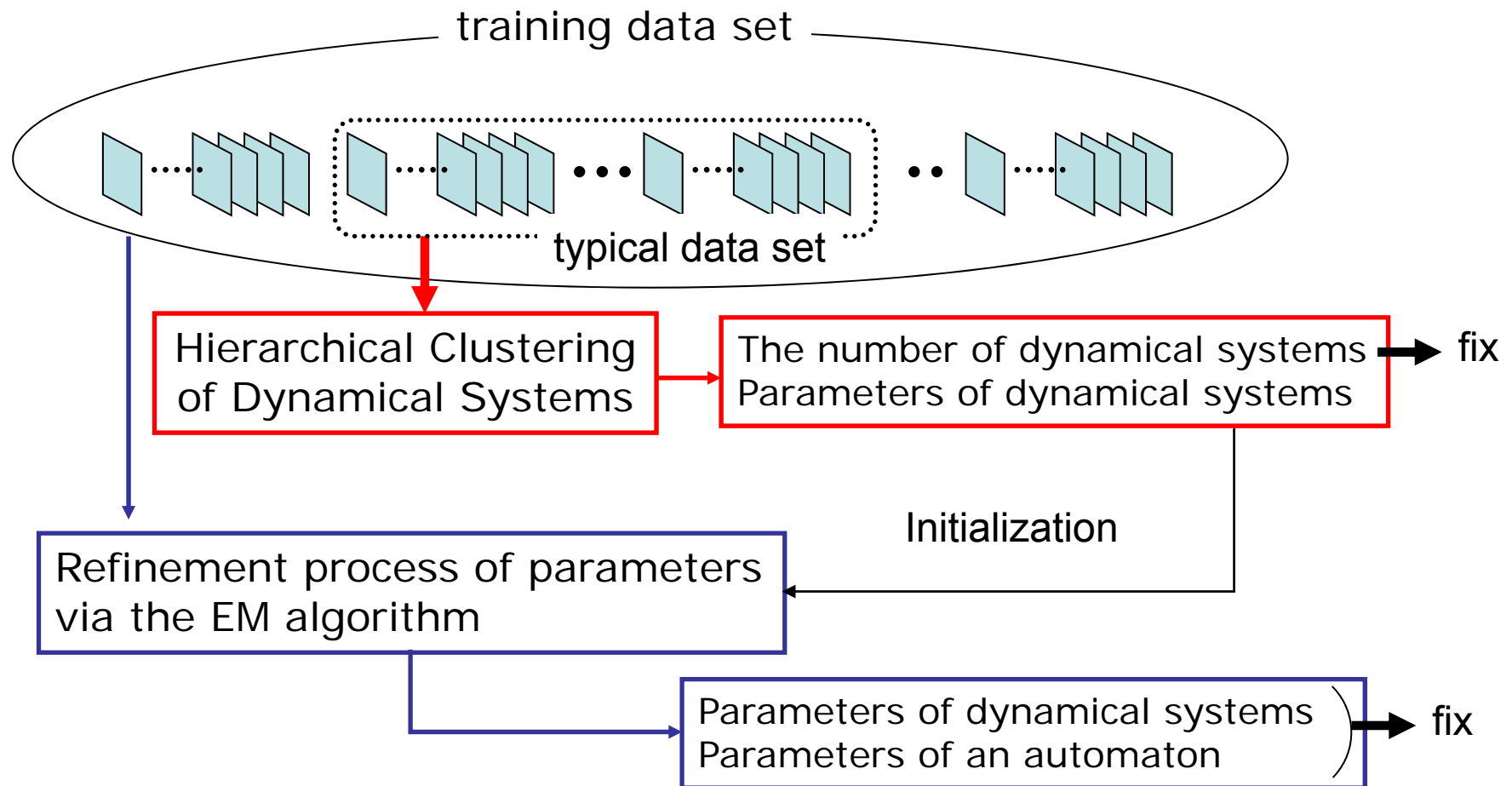
Assume that only a set of vector sequences is given

- Need to define a set of dynamics (primitives)
 - Defined manually in existing work
- Need to solve paradoxical nature of parameter estimation
 - Segmentation requires identified dynamical systems
 - Identification of dynamical systems requires segmentation
- Need to solve initialization problem of the EM algorithm
 - Strongly depends on its initial parameters



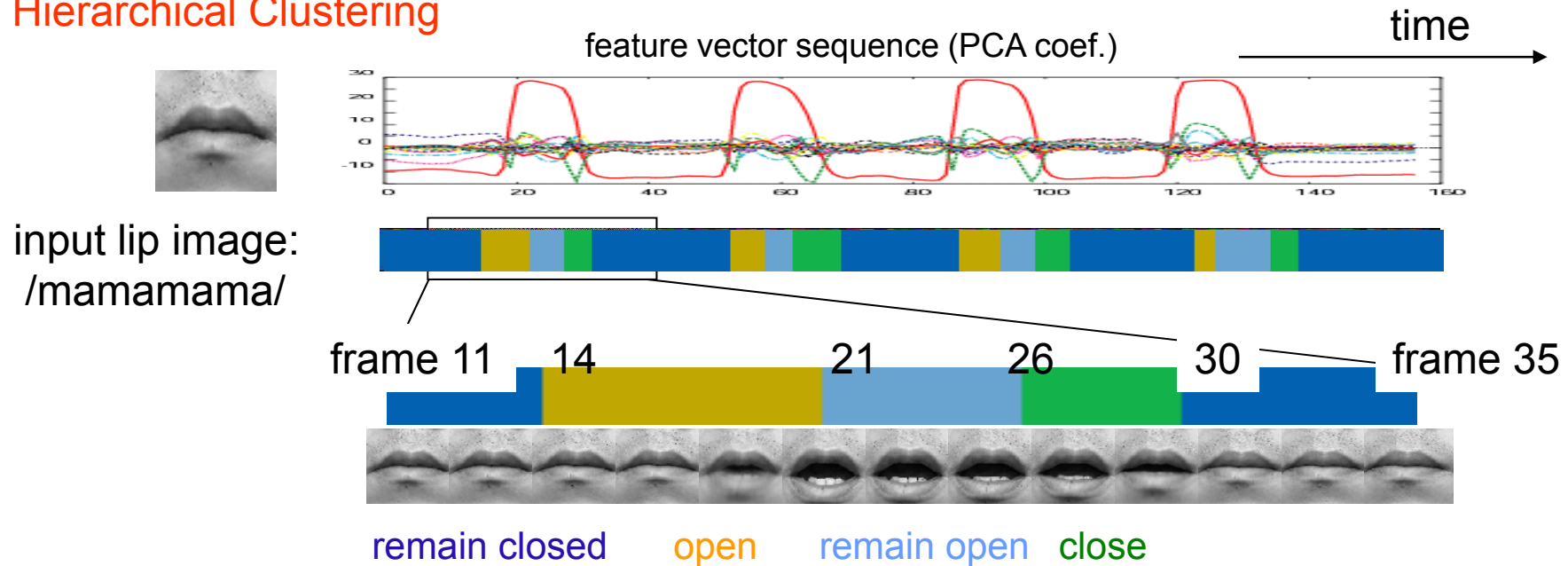
Overview of the Training

- Two-step learning method

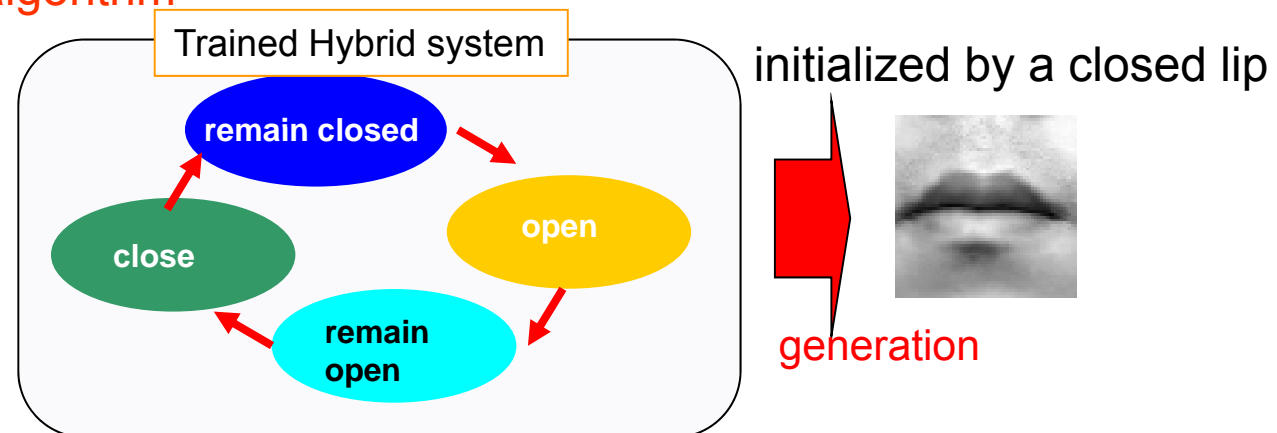


Example of Two-Step Learning Method

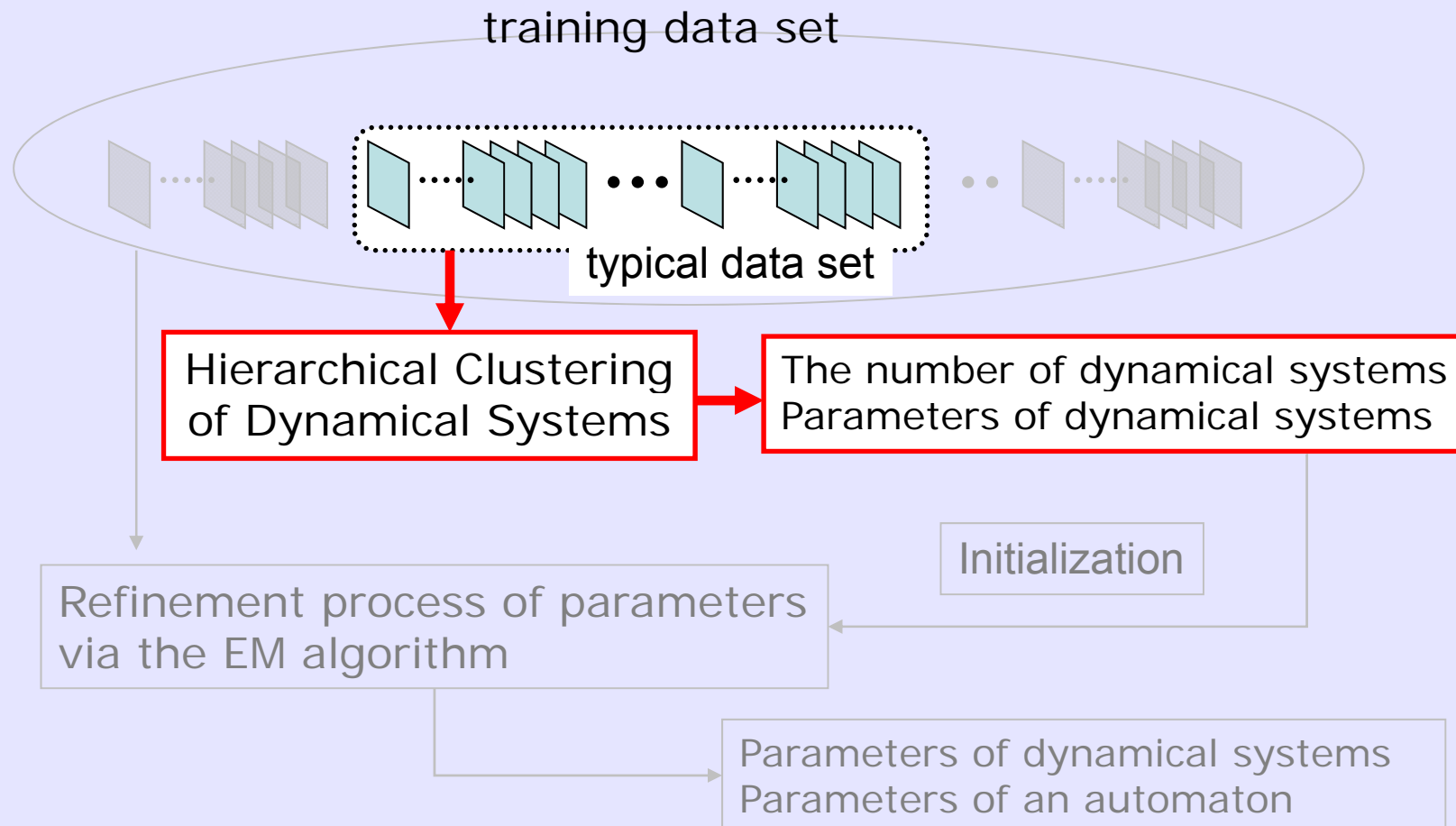
Hierarchical Clustering



Training via the EM algorithm



Hierarchical Clustering of Dynamical Systems



Constrained Linear System Identification

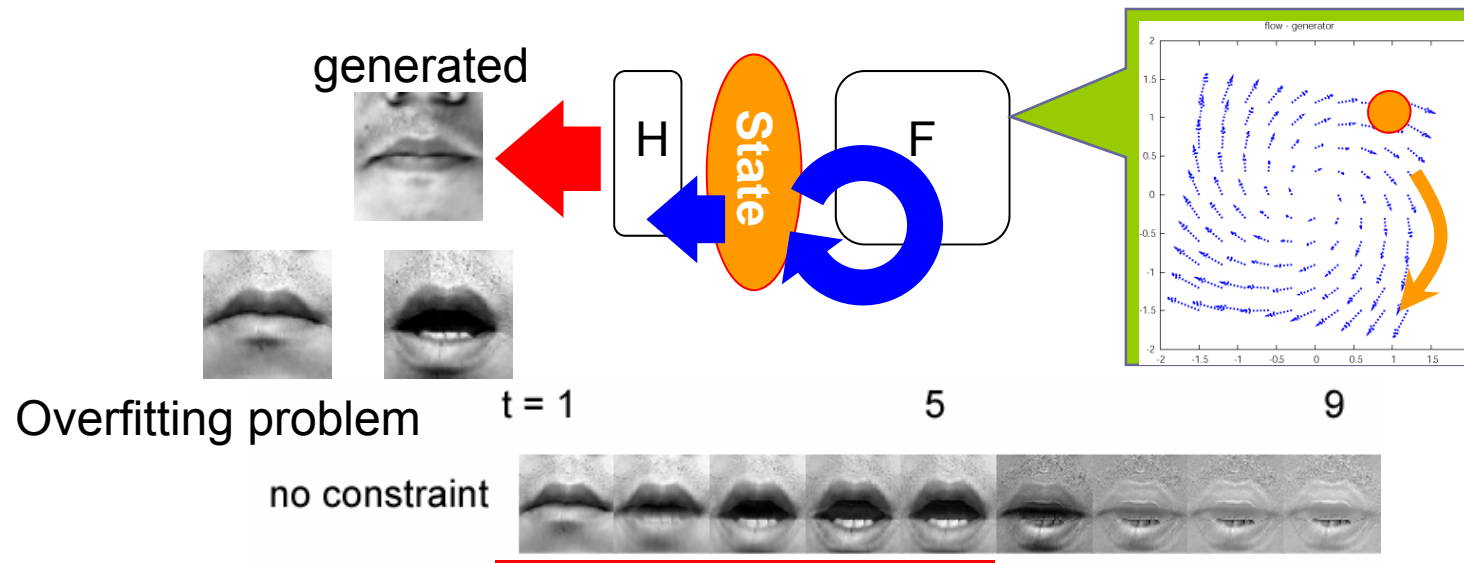
- State transition

$$x_t = \underbrace{F^{(i)}}_{\text{transition matrix}} x_{t-1} + \underbrace{g^{(i)}}_{\text{bias}} + \underbrace{w_t^{(i)}}_{\text{process noise}}$$

- Observation

$$y_t = \underbrace{H}_{\text{observation matrix}} x_t + v_t$$

($w_t, v_t \sim \text{Gaussian}$)
observation noise



System behavior is determined by transition matrix F

Need constraints on the transition matrix F

Class of Linear Dynamical Systems

- Temporal evolution of the state

$$x_t = F^t x_0 = c_1 \lambda_1^t e_1 + c_2 \lambda_2^t e_2 + \cdots + c_n \lambda_n^t e_n$$

$$F = E \Lambda E^{-1} = [e_1, \dots, e_n] \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} [e_1, \dots, e_n]^{-1}$$

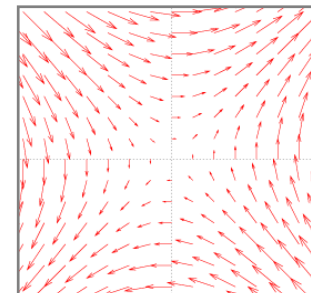
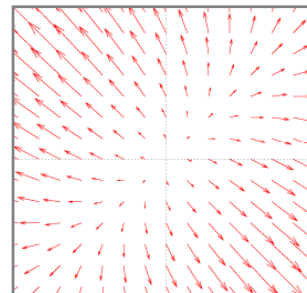
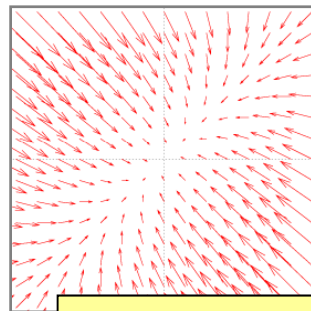
(ex.) $n = 2$

$|\lambda_1| < 1$ and $|\lambda_2| < 1$

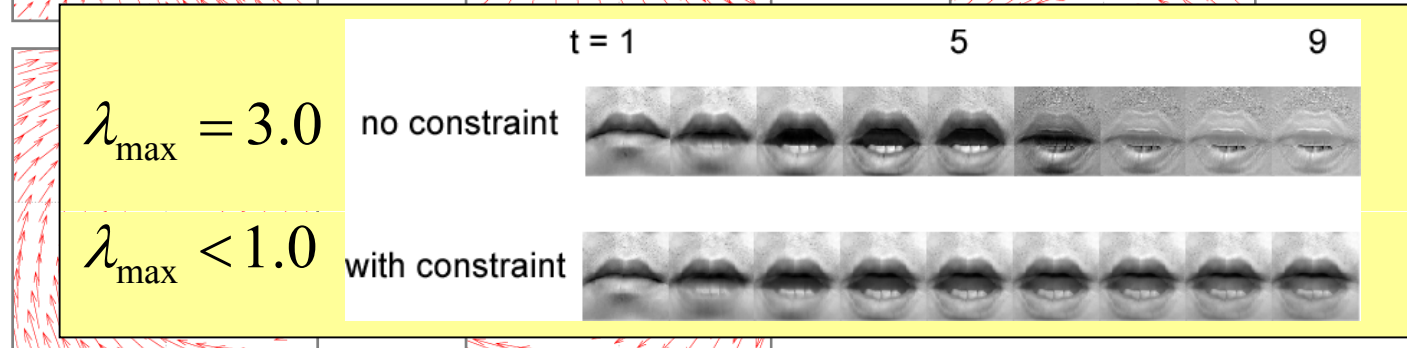
$|\lambda_1| > 1$ and $|\lambda_2| > 1$

$|\lambda_1| > 1$ and $|\lambda_2| < 1$

λ_1 and λ_2 is
positive real



λ_1 and λ_2 is
complex num.



Constrained Linear System Identification

- Eigenvalue constraints
 - Upper bound of eigenvalues is determined by

$$UpperBound = \max_r \sum_{c=1}^n |f_{rc}^{(i)}| = 1$$

$$F^{(i)} = \begin{bmatrix} f_{11}^{(i)} & \cdots & f_{1n}^{(i)} \\ \vdots & & \vdots \\ f_{n1}^{(i)} & \cdots & f_{nn}^{(i)} \end{bmatrix} u$$

The element $f_{rc}^{(i)}$ is circled in red in the original image.

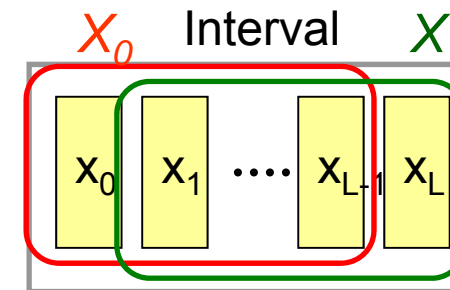
Deduced from
Gershgorin's theorem

- Estimation of transition matrix

Least squares problem

$$F^{(i)*} = \arg \min_{F^{(i)}} \| F^{(i)} X_0^{(i)} - X_1^{(i)} \|^2$$

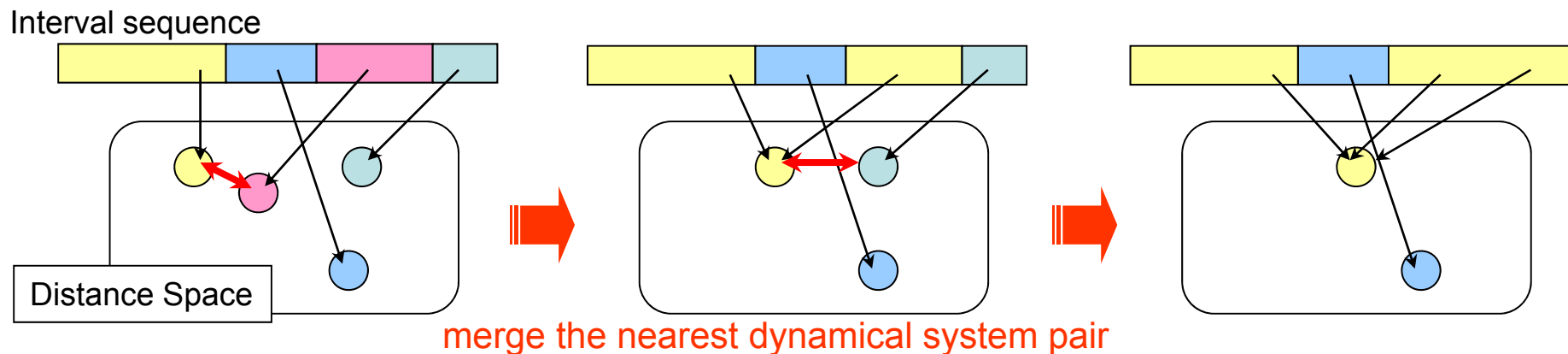
$$= \lim_{\delta^2 \rightarrow 0} X_1^{(i)} \underline{X_0^{(i)T} (X_0^{(i)} X_0^{(i)T} + \delta^2 I)^{-1}} = X_1^{(i)} \underline{X_0^{(i)+}} \quad (\text{pseudo-inverse})$$



Stop the limit before δ^2 converges to 0
 δ^2 controls the scale of matrix elements

Algorithm of Hierarchical Clustering

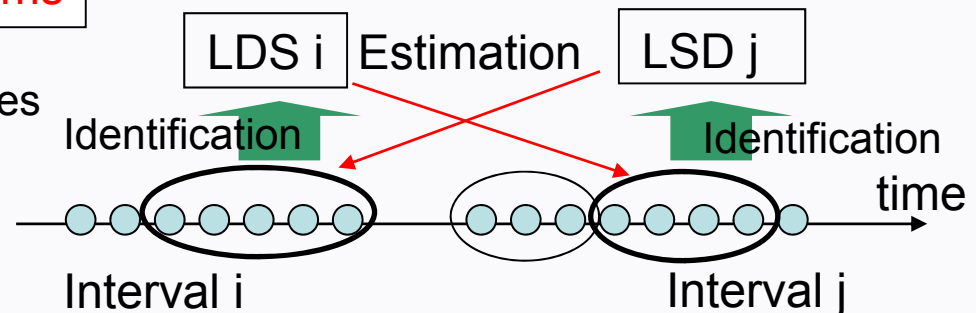
1. Divide the training data into short intervals (Initialization)
2. Identify the parameters of the dynamical systems from each interval
3. Calculate distances of all the dynamical system pairs
4. Merge the **nearest pair** of dynamical systems (intervals are also merged)
5. Identify the new dynamical system from the merged intervals
6. Repeat 3 to 5



Distance between dynamical systems

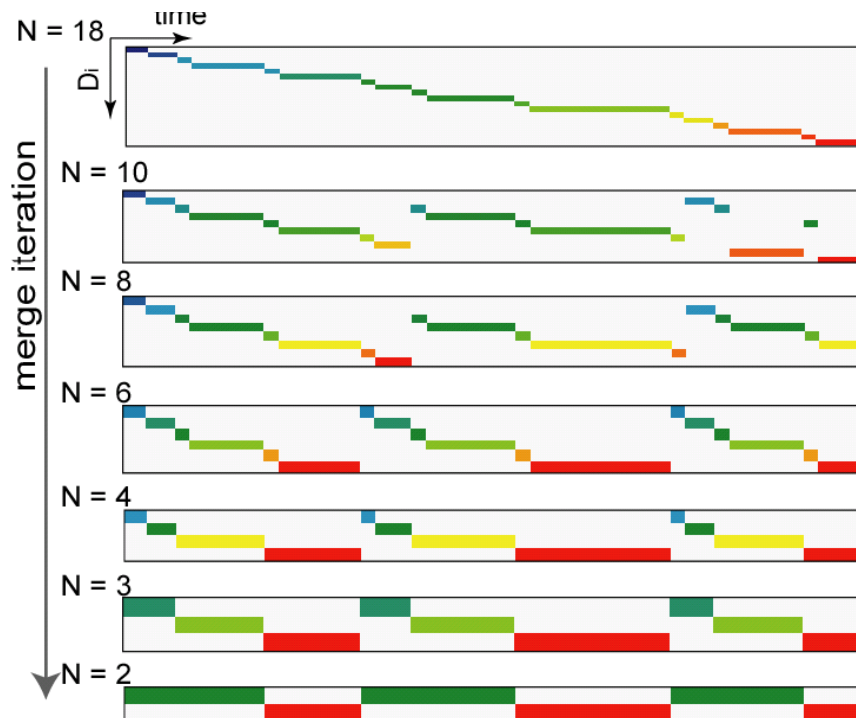
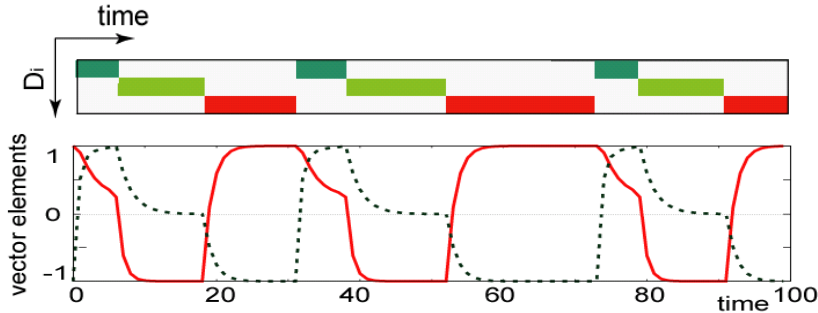
Average of Kullback-Leibler divergences

$$Dist(D_i, D_j) = \frac{KL(D_i \parallel D_j) + KL(D_j \parallel D_i)}{2}$$

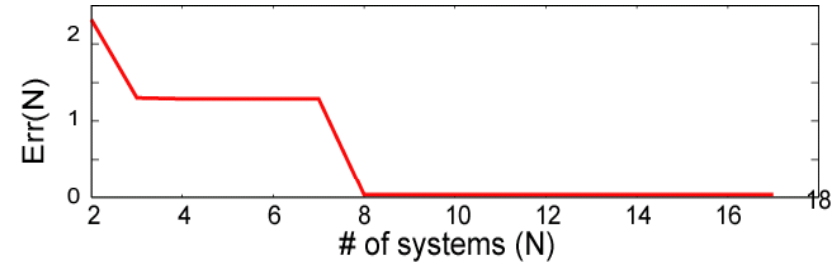


Simulation Result

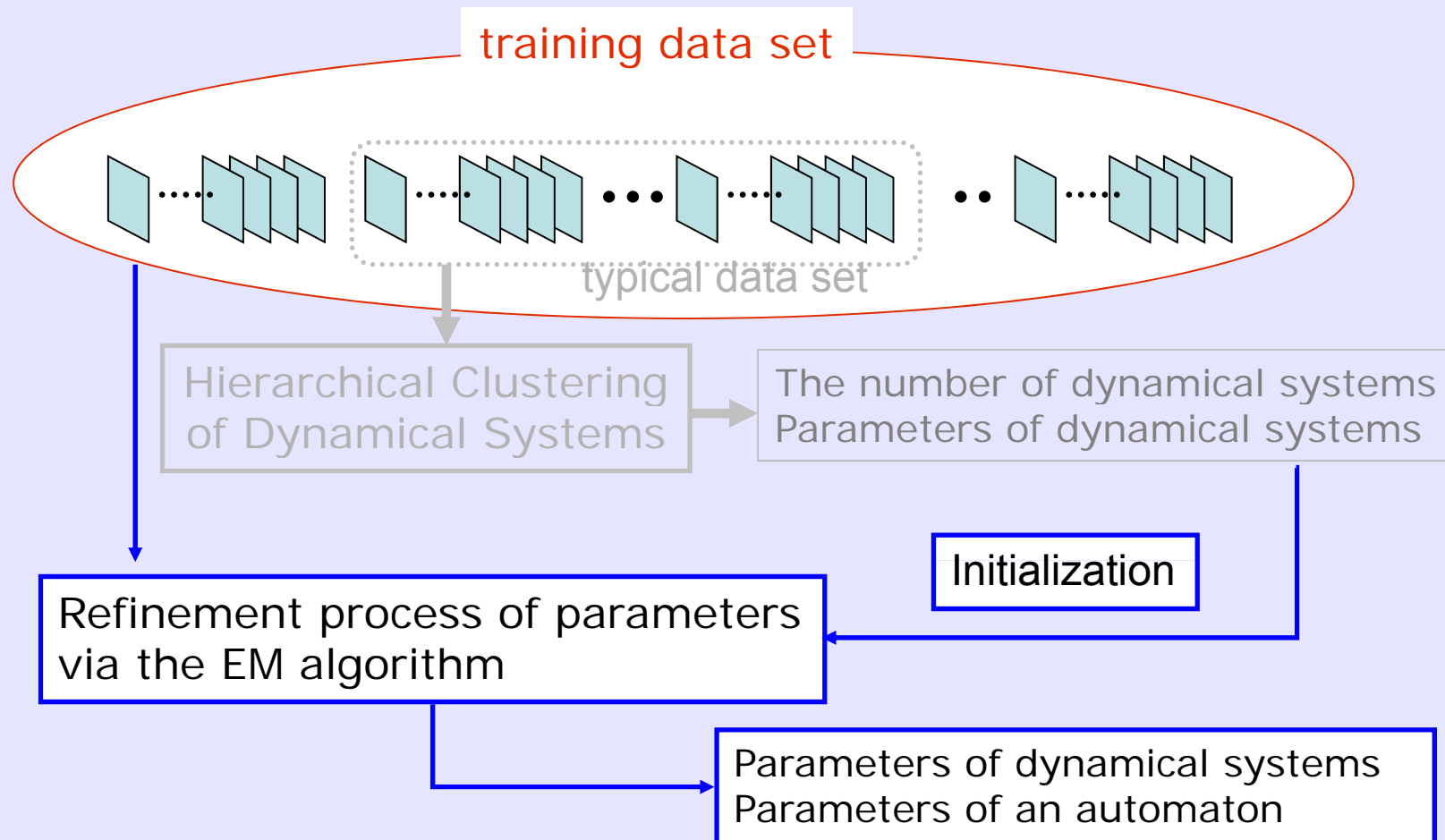
- Input data: generated from three dynamical systems



Prediction error of overall systems at each iteration step



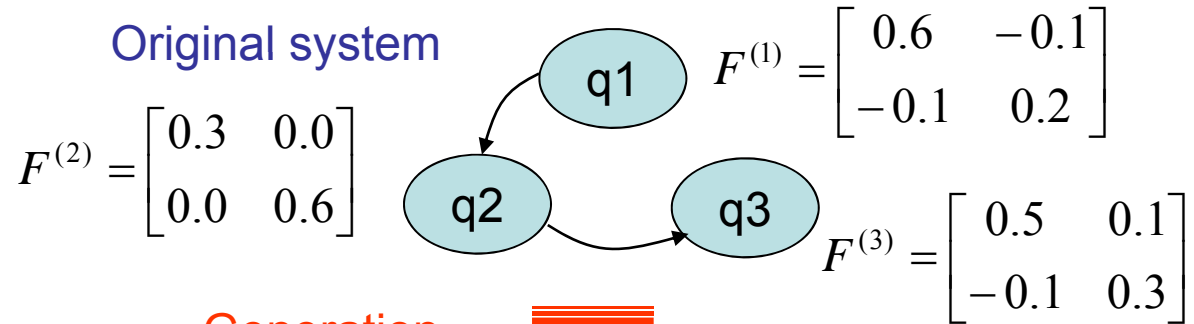
Refinement Process of All the Parameters



Evaluation based on Simulated Data

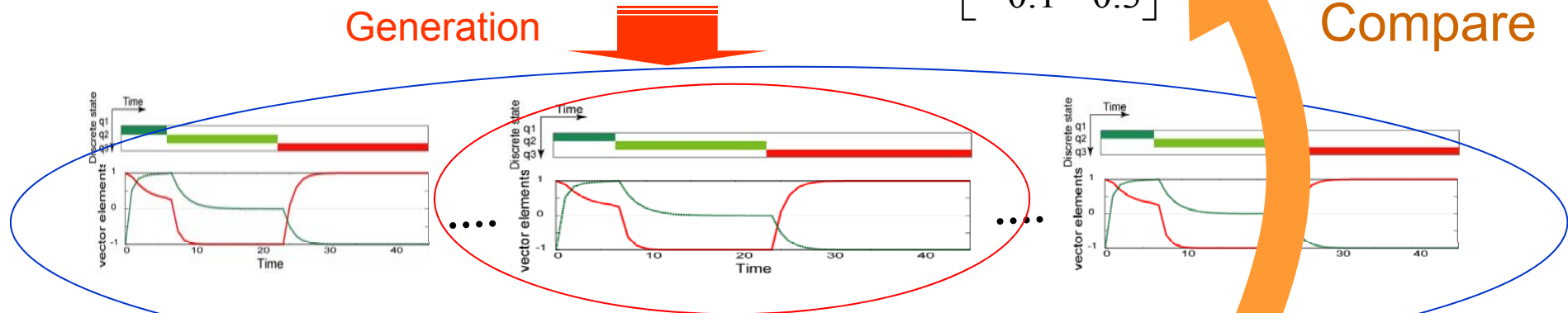
- Interval-Based HDS with known parameters

Original system



Generation

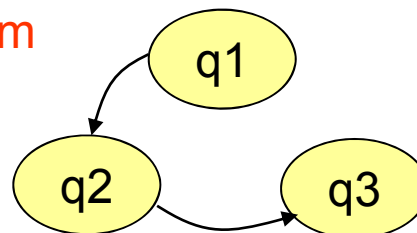
Compare



Used for clustering (one sequence)

Used for EM algorithm (10 sequences)

Trained system



Estimated parameters

Simulation Result

- Comparison between given and estimated parameters

Original (ground truth)

$$F^{(1)} = \begin{bmatrix} 0.60 & -0.10 \\ -0.10 & 0.20 \end{bmatrix} \quad F^{(2)} = \begin{bmatrix} 0.30 & 0.00 \\ 0.00 & 0.60 \end{bmatrix} \quad F^{(3)} = \begin{bmatrix} 0.50 & 0.10 \\ -0.10 & 0.30 \end{bmatrix}$$

Estimated parameters via clustering

$$F^{(1)} = \begin{bmatrix} 0.01 & -0.14 \\ -0.01 & 0.21 \end{bmatrix} \quad F^{(2)} = \begin{bmatrix} 0.86 & 0.30 \\ -0.21 & -0.06 \end{bmatrix} \quad F^{(3)} = \begin{bmatrix} 0.74 & -0.44 \\ -0.11 & 0.75 \end{bmatrix}$$

Estimated parameters via EM algorithm

$$F^{(1)} = \begin{bmatrix} 0.60 & -0.10 \\ -0.10 & 0.20 \end{bmatrix} \quad F^{(2)} = \begin{bmatrix} 0.32 & 0.02 \\ 0.06 & 0.52 \end{bmatrix} \quad F^{(3)} = \begin{bmatrix} 0.49 & 0.09 \\ -0.10 & 0.29 \end{bmatrix}$$

Discussion

- Interval-based hybrid dynamical system
 - Interval-based state transition to model tempo and rhythm
 - Linear dynamics to model continuously changing patterns
- Two-step learning method for the interval-based
 - Clustering of dynamical systems + EM algorithm
 - Constrained system identification based on eigenvalues



Chapter 4

Analysis of Timing Structures in Multipart Motion of Facial Expression

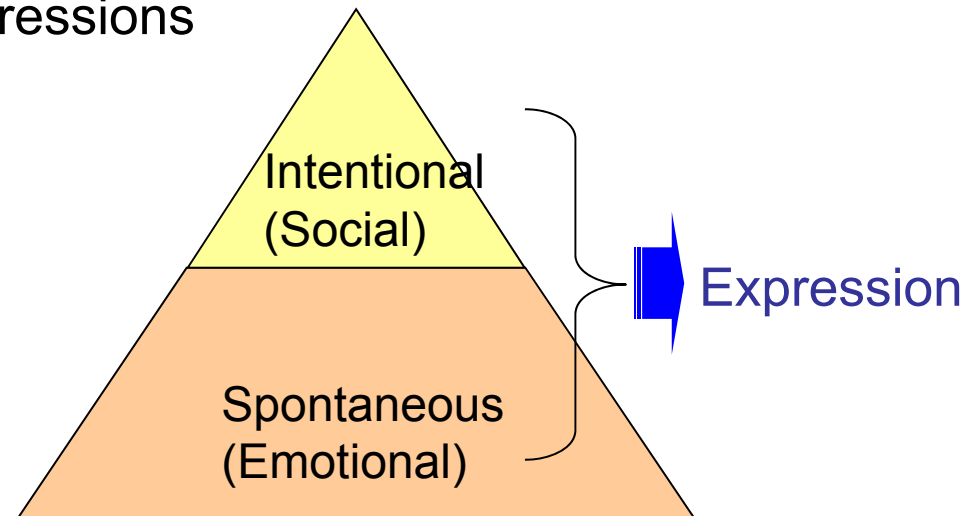
Facial Expression as Communication Protocol

- Communication via facial expressions

- Generation
 - Express internal state
- Recognition
 - Estimate internal state

- Acquisition of expressions

- Intrinsic
 - Smile, cry, surprise
- Learned from experience (parents)
 - Social contexts



Facial expression = Emotional category



Facial expression = Communication protocol

Related work

- FACS (Facial Action Coding System) (Ekman, et al.)
 - AU (Action Unit) : motion primitives in faces
 - Describe facial expressions based on combination of AU
(ex.) Surprise = AU1+2+5+26
 - Describes only emotional categories
 - {happiness, surprise, fear, anger, disgust, sadness}
- Problem: cannot describe dynamic structures
(synchronous/asynchronous motions, duration of motions, etc.)

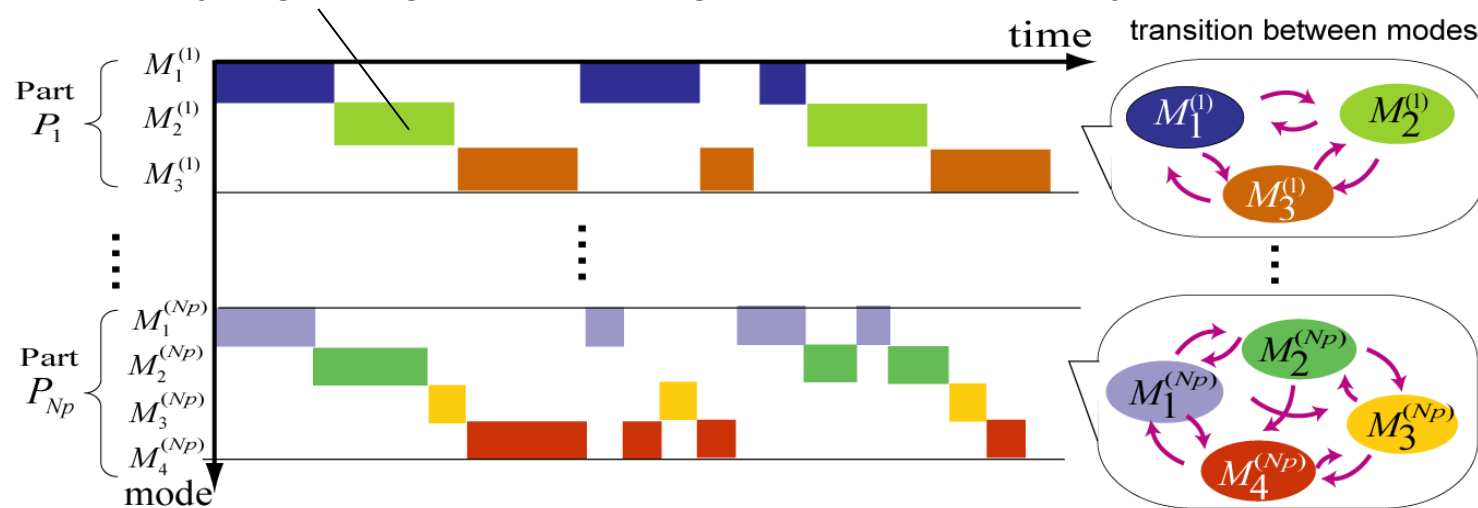
Psychological experiments

- Temporal difference of beginning time between eyes and mouth is important to discriminate social, pleasant, and unpleasant smiles (Nishio&Koyama1997)
- Human recognition of facial expressions depends on duration of motion (Ekman&Friesen1982, Kamachi2001, Krumhuber2005)

Facial Score: Interval-Based Facial Action Description

- Define facial parts – move independently
- Define **modes** - motion primitives (dynamics)

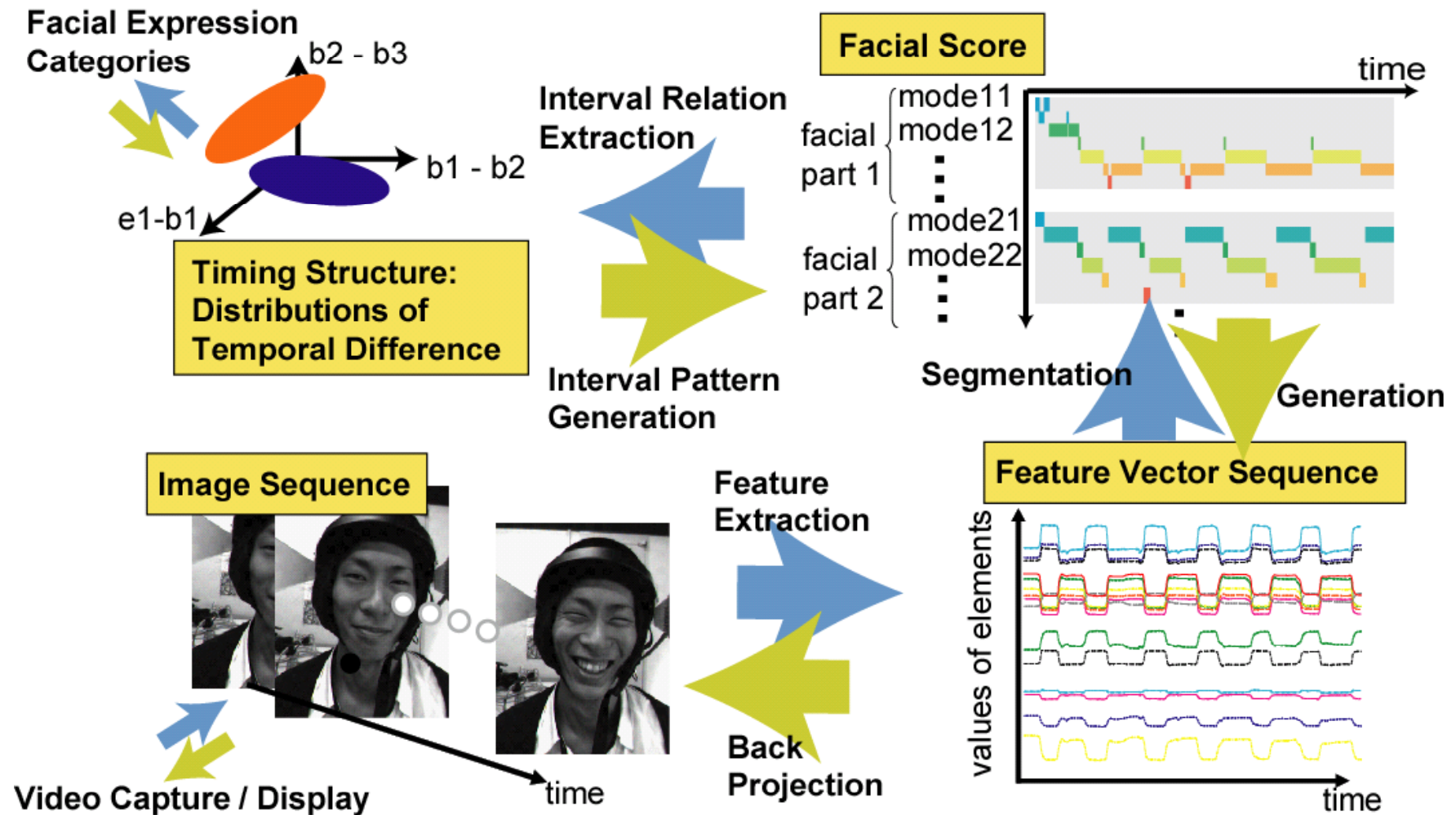
Interval: { beginning point, ending point , mode label}



Facial Score: {Interval set of part1, ... , Interval set of part N}

Represent timing structure among modes (dynamics)

Facial Expression Generation and Recognition



1. Definition of Facial Scores
2. Automatic Acquisition of Facial Scores
3. Evaluation

- Definition of parts
- Definition of modes

Facial Parts in Facial Scores

- Follows Ekman's definition
 - Parts = {left/right eyebrow, left/right eye, nose, mouth}
- Feature vector of each part
 - x,y coordinates of feature points

$$z^{(a)} = (x_1^{(a)}, y_1^{(a)}, \dots, x_{N_a}^{(a)}, y_{N_a}^{(a)})^T$$

$a \in \text{Parts}$

(dimensionality: each eyebrow : 10,
each eye: 16, nose: 22, mouth: 16)



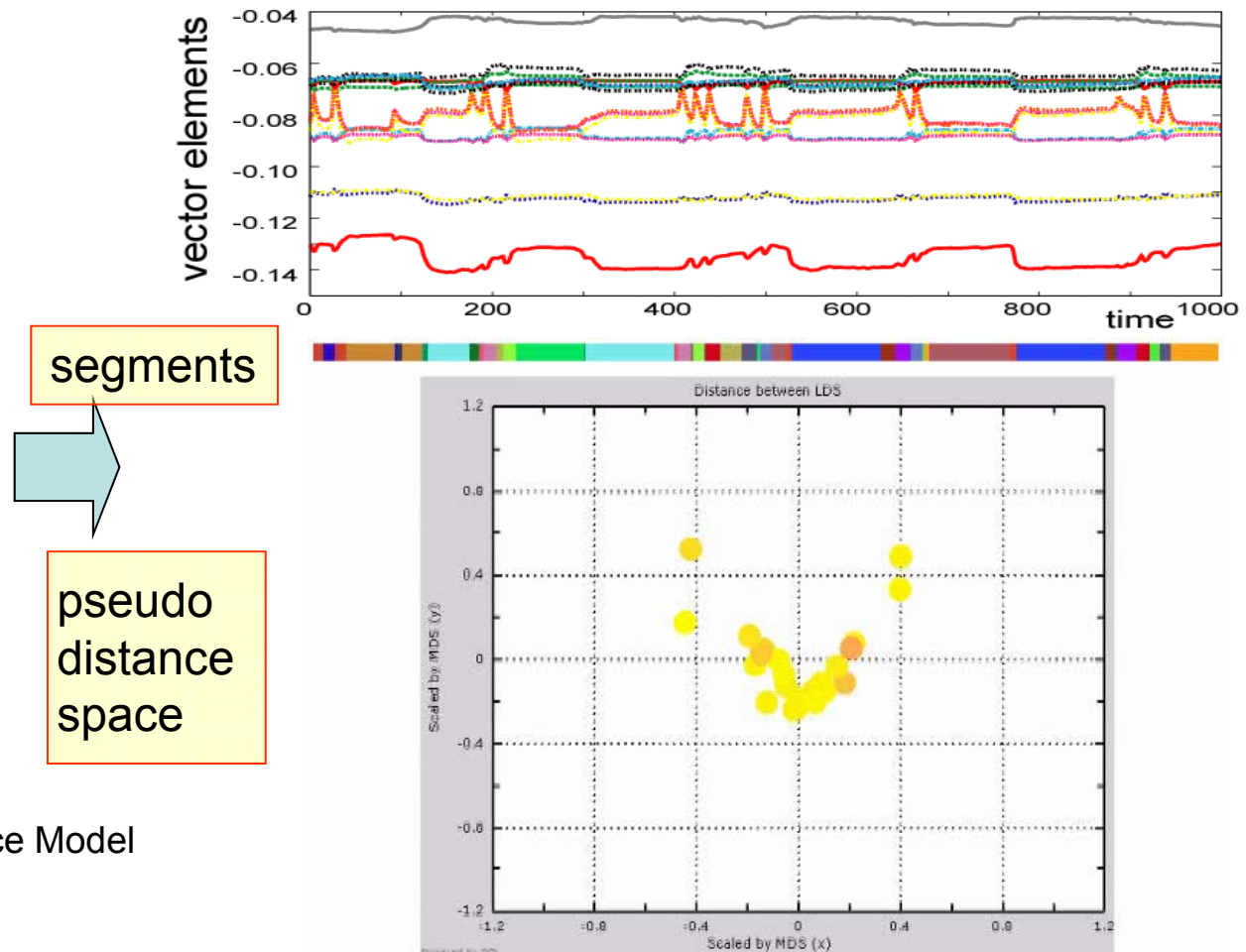
Active Appearance Model(AAM)
(Cootes 1998)

Facial Modes in Facial Scores

- Smiled four times
- Feature vector: x,y coordinates of feature points around right eye (eight points)
- Length: 1000 frames

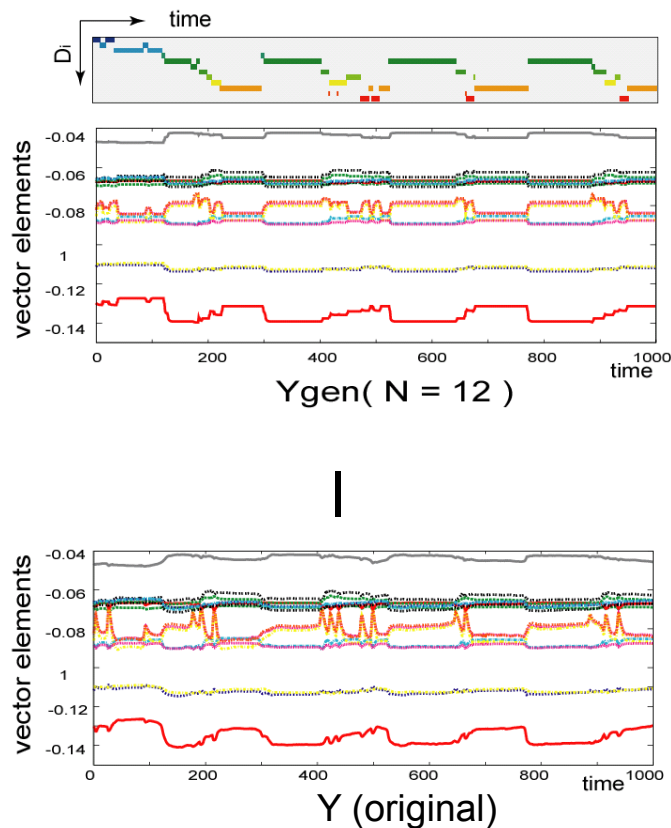


Tracked by Active Appearance Model
(Cootes 1998)
Thanks to Stegmann's AAM-API

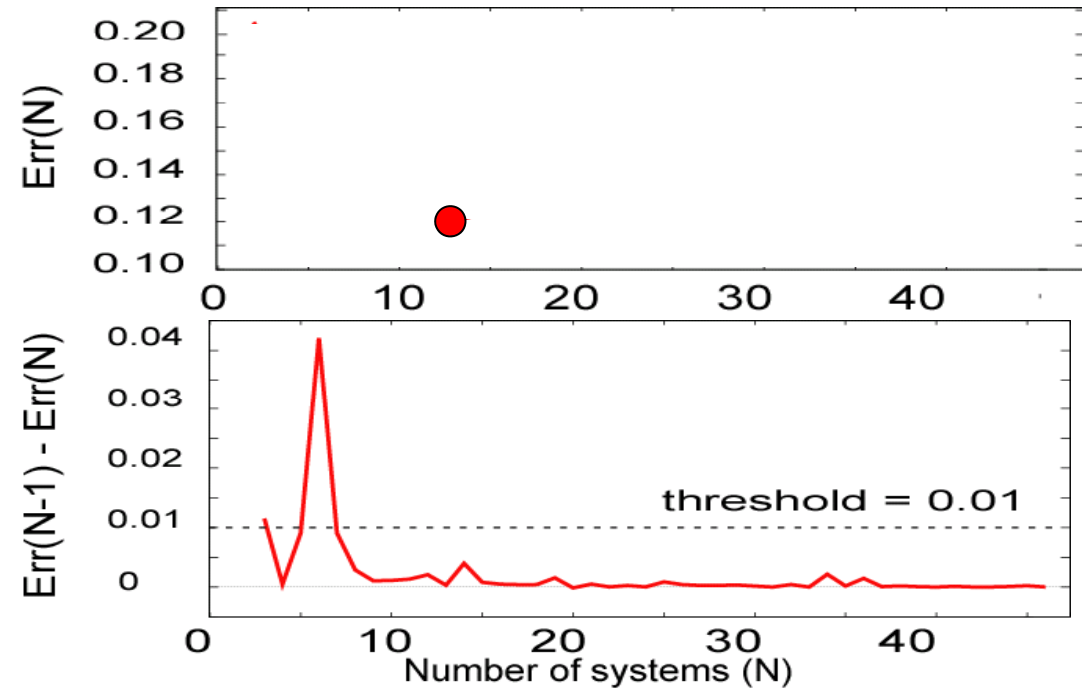


Determine the number of modes (dynamical systems)

- Find a rapid change of model fitting error curve



$$Err(N) = \sqrt{\sum_{t=1}^L \| \text{orig}(t) - \text{gen}^{(N)}(t) \|^2}$$

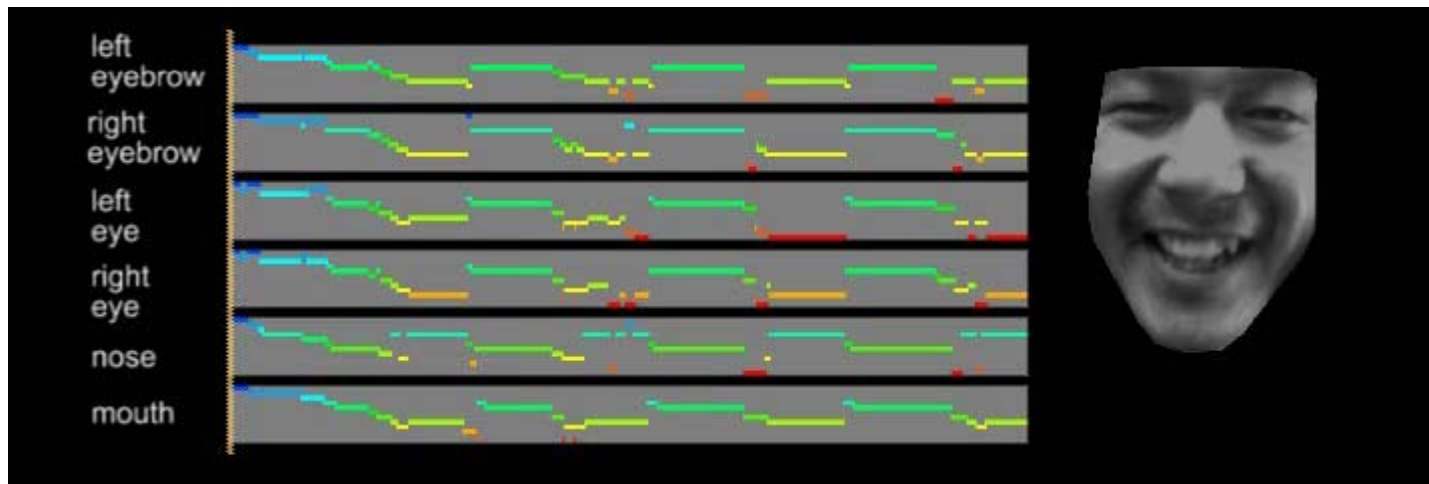


1. Definition of Facial Scores
2. Automatic Acquisition of Facial Scores
3. Evaluation

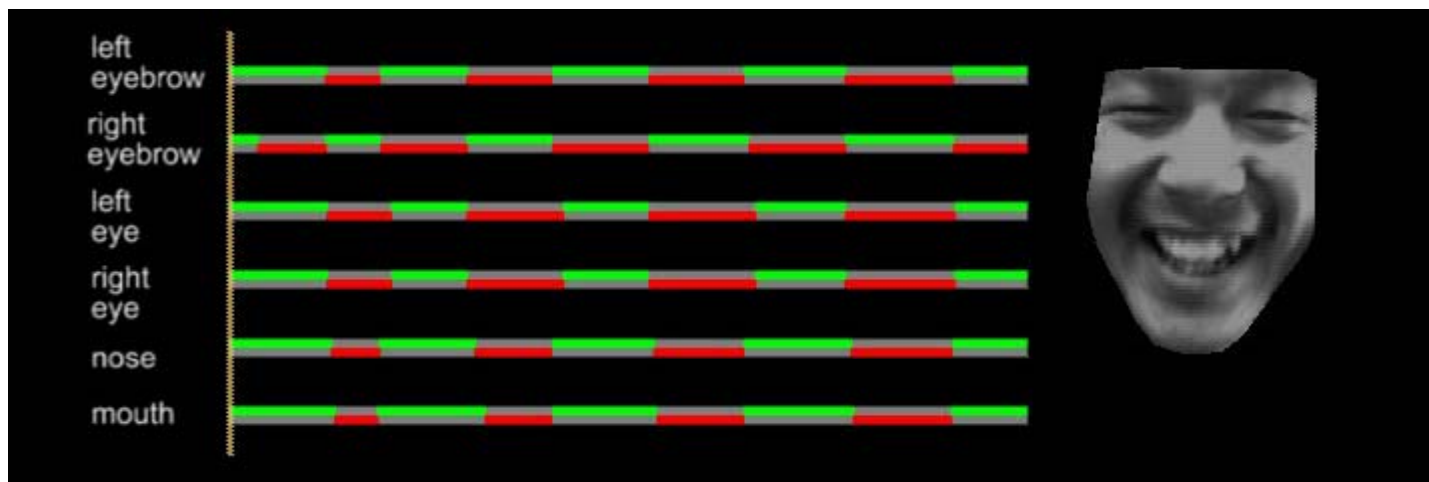
- Generation of expressions
- Discrimination of expressions

Generation of Facial Expression

Num. of modes: 12 in each parts



Num. of modes: 2 in each parts



Discrimination of facial expressions

- Subjects

- Intentional smile
- Spontaneous smile
- Six (male)
- about 30-50 times for each smile category

} separable?

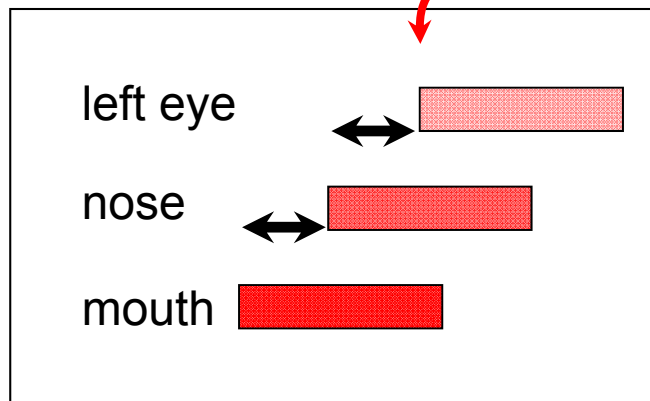
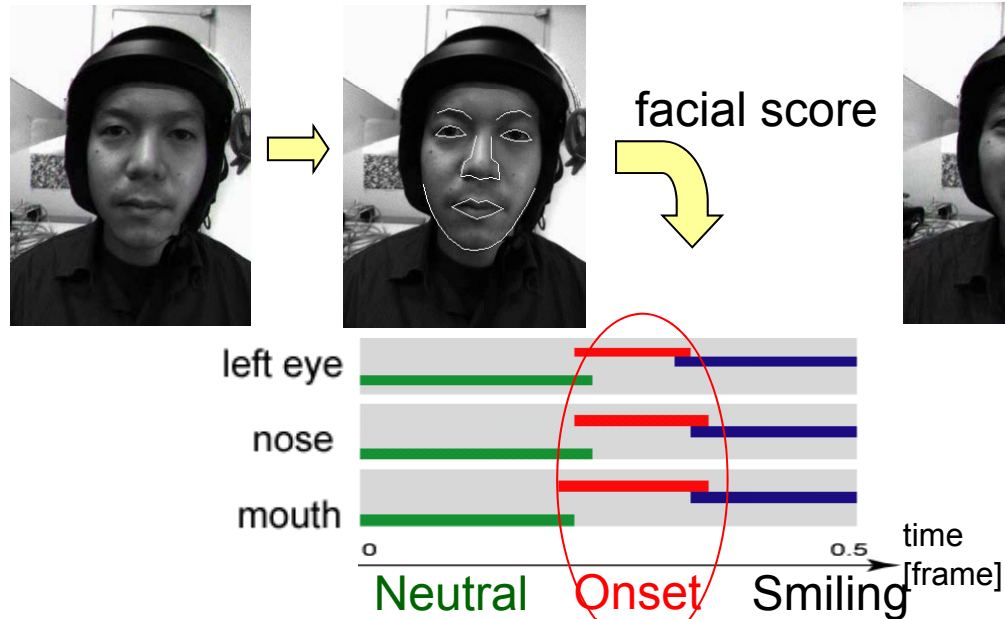


- Method

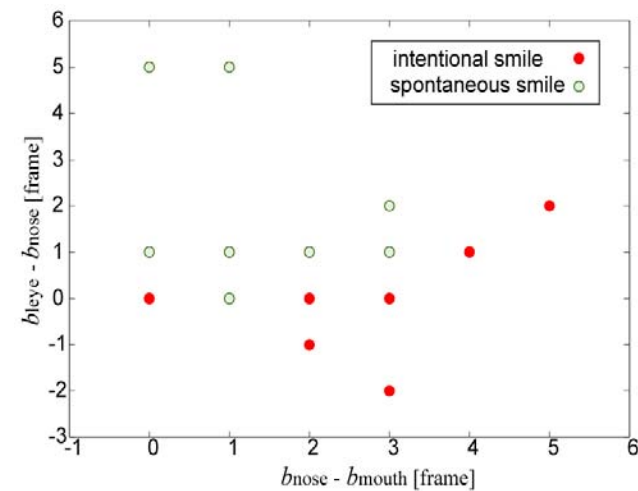
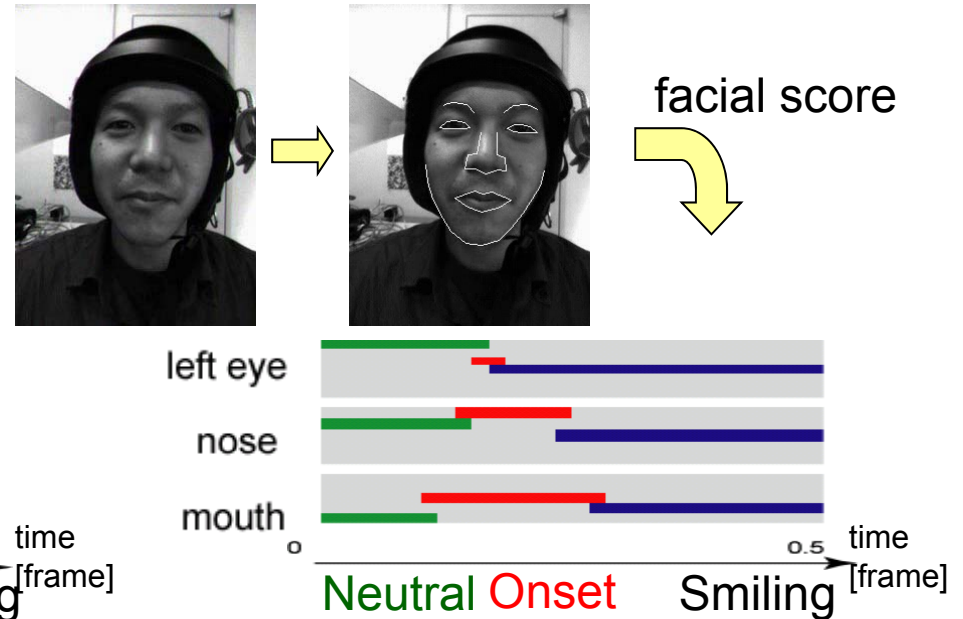
- Video data
 - VGA480x640(down sampling to 240x360) , 60fps
- Instruction of expressions
 - Start from neutral face
 - Intentional: make smile during watching a disgust movie
 - Spontaneous: watch Japanese stand-up comedy movies

Timing Structure in Facial Scores

Intentional

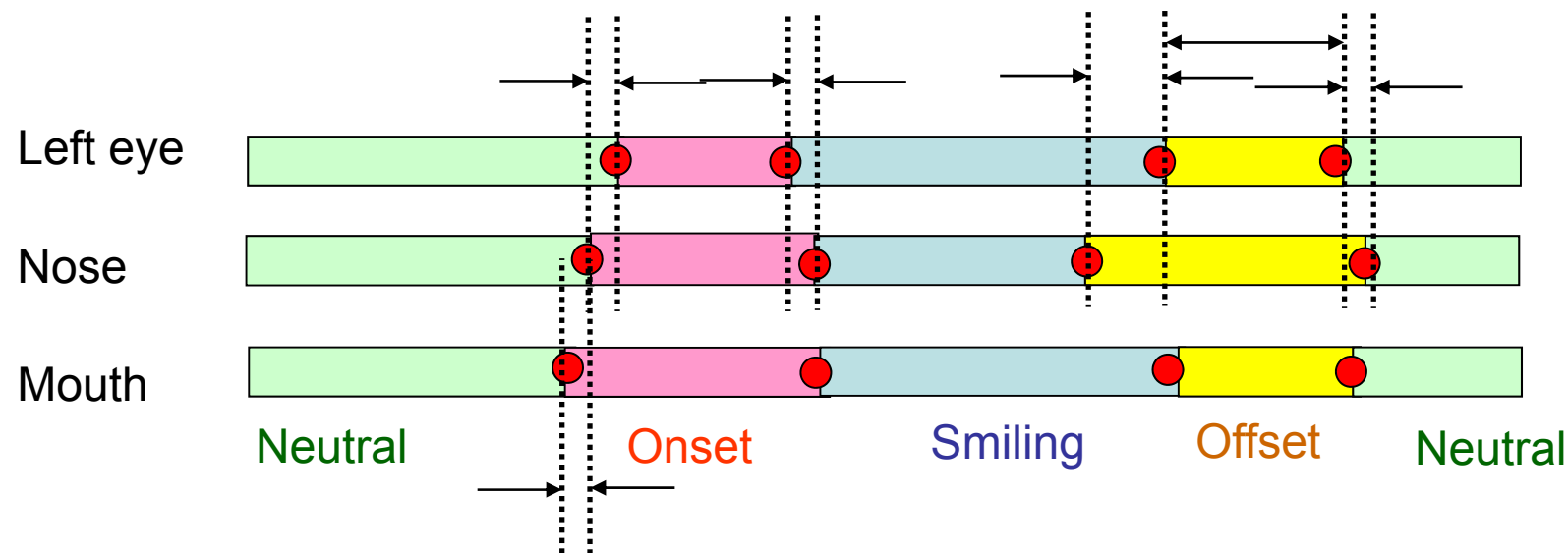


Spontaneous



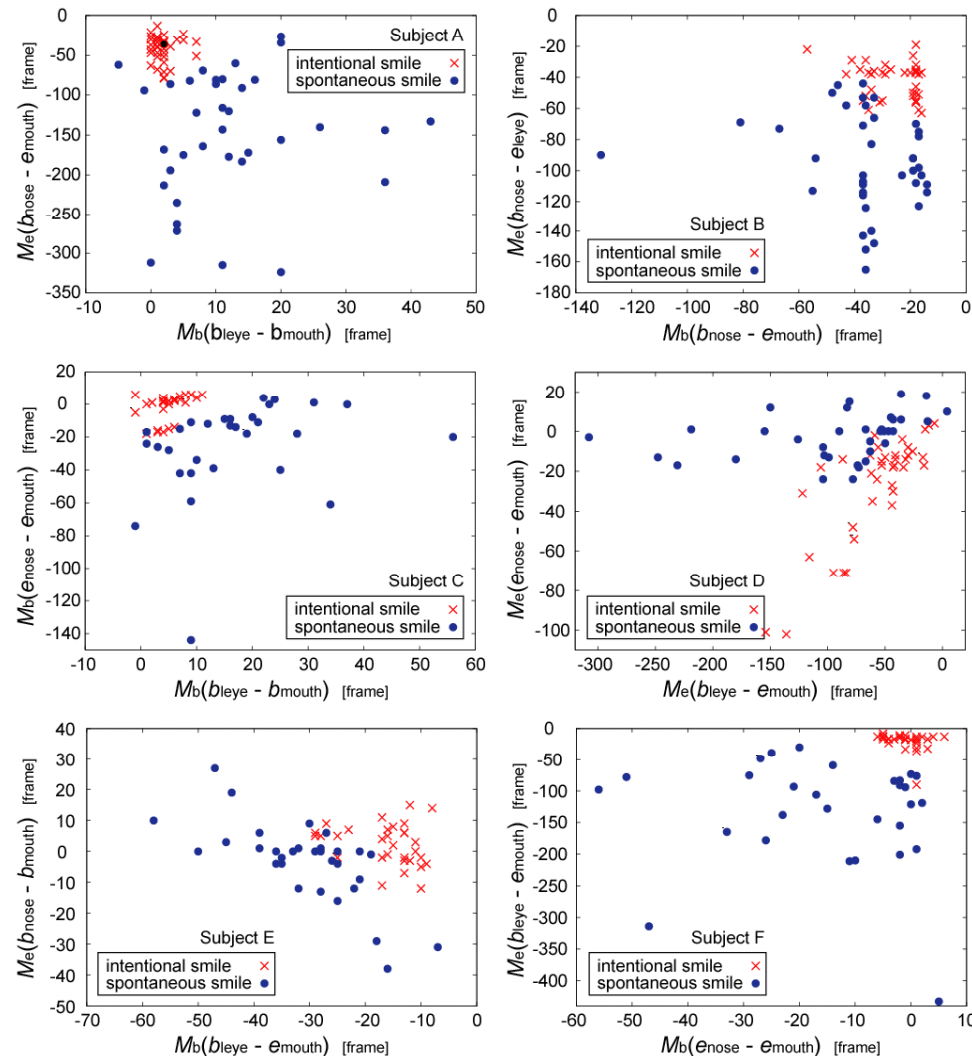
Extract Timing Structure from Facial Scores

1. Use temporal differences among the **beginning and ending points of “onset” and “offset” motion**
2. Calculate two-dimensional distributions using a combination of two temporal differences as the axes
3. Calculate distance between the distributions of two smiles for all the combinations in 2



Result of Discriminating Two Smiles

Chose two axes that provides the maximum distance between two distributions



Six subjects

intentional smile ×
spontaneous smile ●

Recognition rate of each subject
based on the support vector machine
(leave-one-out method)

subject	intentional (%)	spontaneous (%)
A	100	83.8
B	100	79.4
C	82.4	96.4
D	85.1	79.7
E	85.3	90.3
F	96.6	93.1

Discussion

- Analysis of timing structure in multipart motion of facial expression
 - Successfully discriminated and recognized intentional and spontaneous smiles

Future Work

1. Long term observation (video capturing)
 - Find expression categories in a bottom-up manner
2. Expression in a context
 - conversation, singing, watching movies
 - relation among multiple subjects
3. Personality
 - Common structure and modes
 - Specific structure and modes

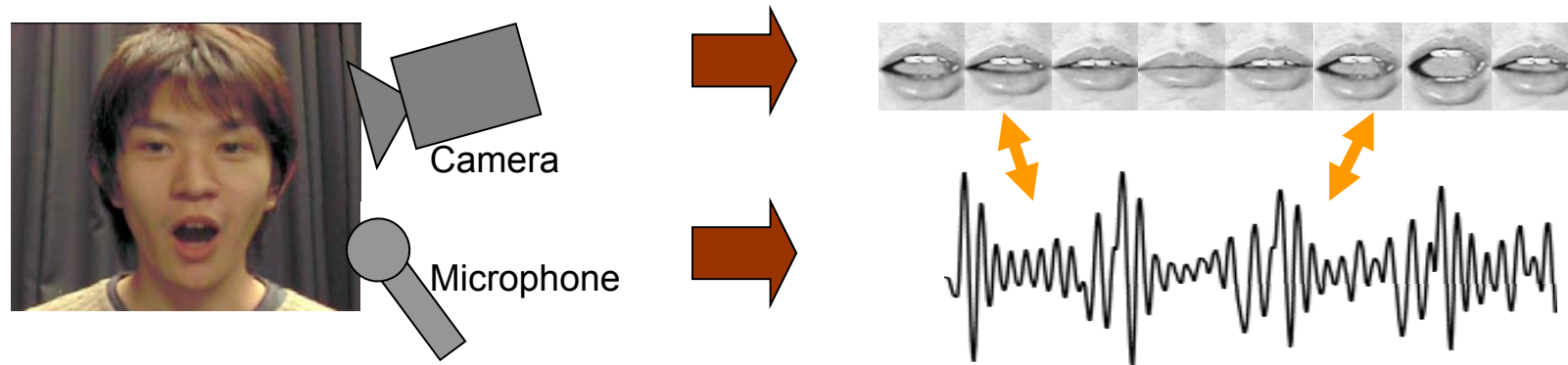


Chapter 5

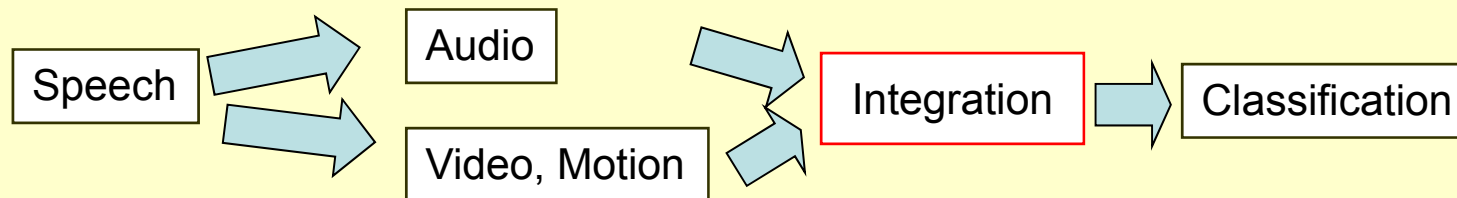
Modeling Timing Structures in Multimedia Signals

Temporal Relation in Multimedia Signals

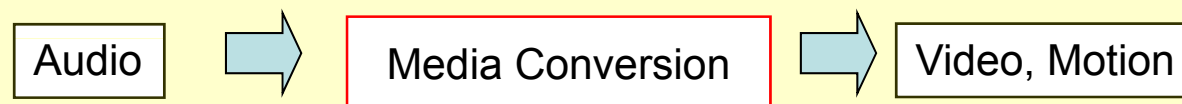
Multimedia signal



Recognition by multimedia integration (ex.) Audio-visual speech recognition

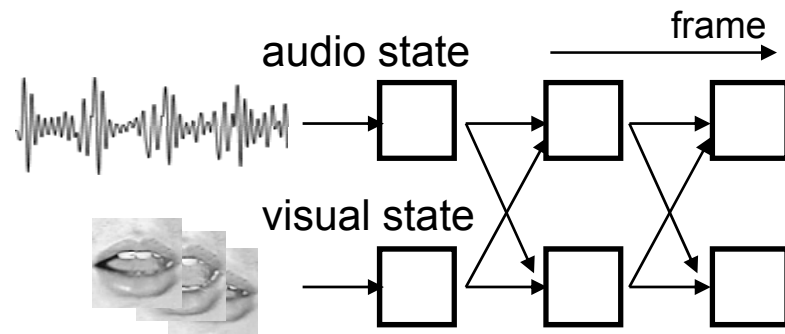


Media signal generation from another related signal (ex.) Lip sync.



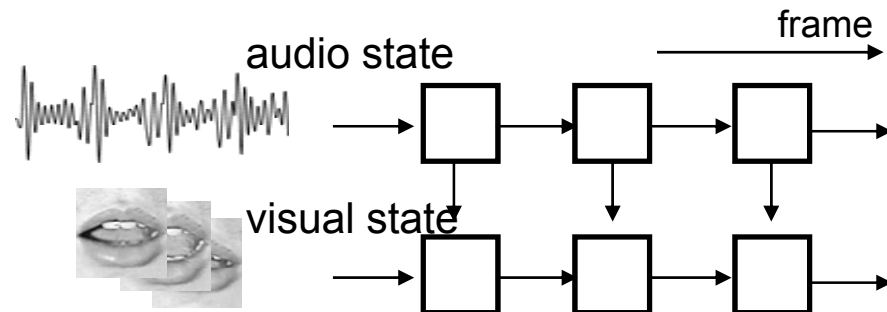
Related Work

State relation in adjacent frames



(ex.) Coupled HMM [Nefian, et al, ICASSP 2002]

Frame-wise state co-occurrence

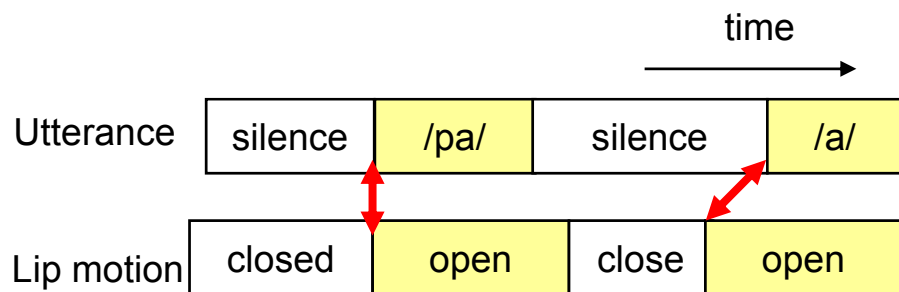


(ex.) Voice puppetry [Brand, SIGGRAPH 1999]

Frame-based

Open Issues

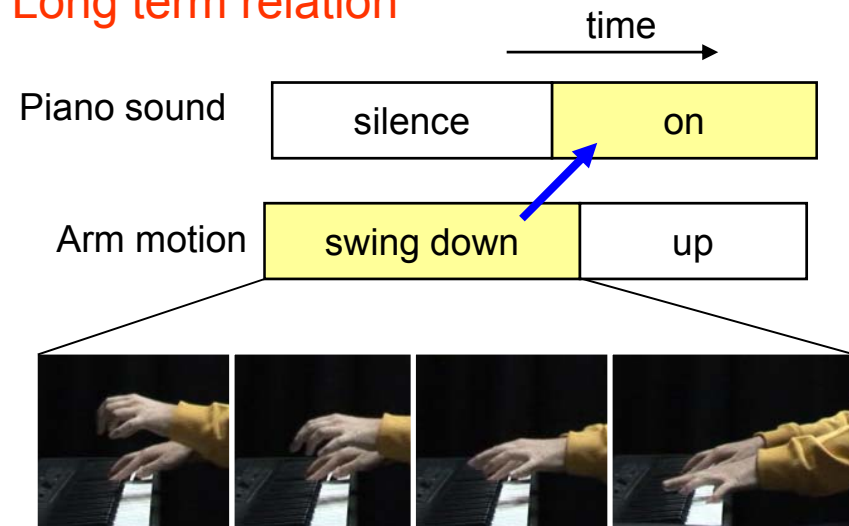
Synchronization mechanisms



strictly synchronized

loosely synchronized

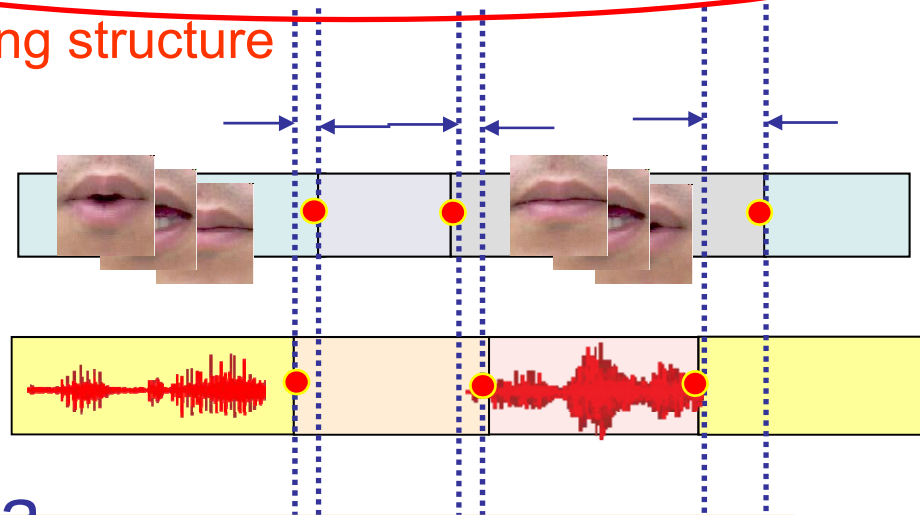
Long term relation



Timing Structure Model

There exists mutual dependency with organized temporal difference between signals

Timing structure



Key Idea

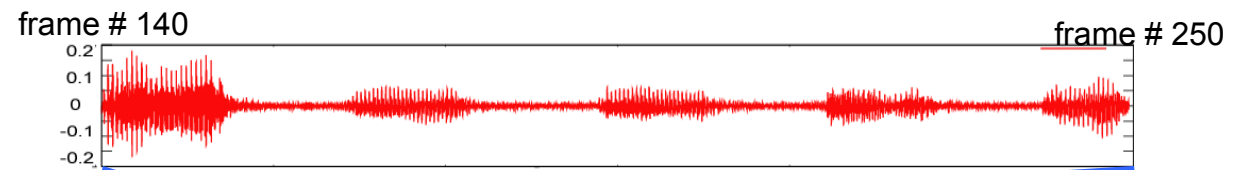
Directly model timing structure
using an interval based representation

1. How to divide signals into intervals?
2. How to model timing structure ?

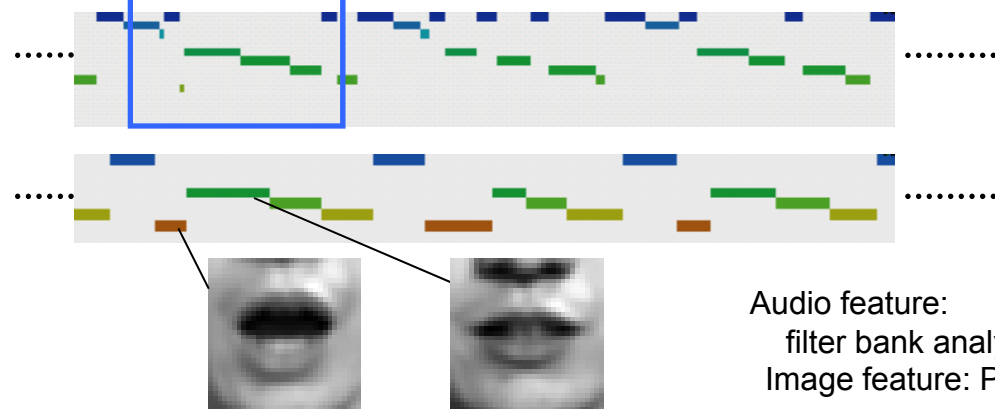
Interval-based HDS

Temporal Relation of Intervals

Input speech:
/aiueo/
nine times
continuously



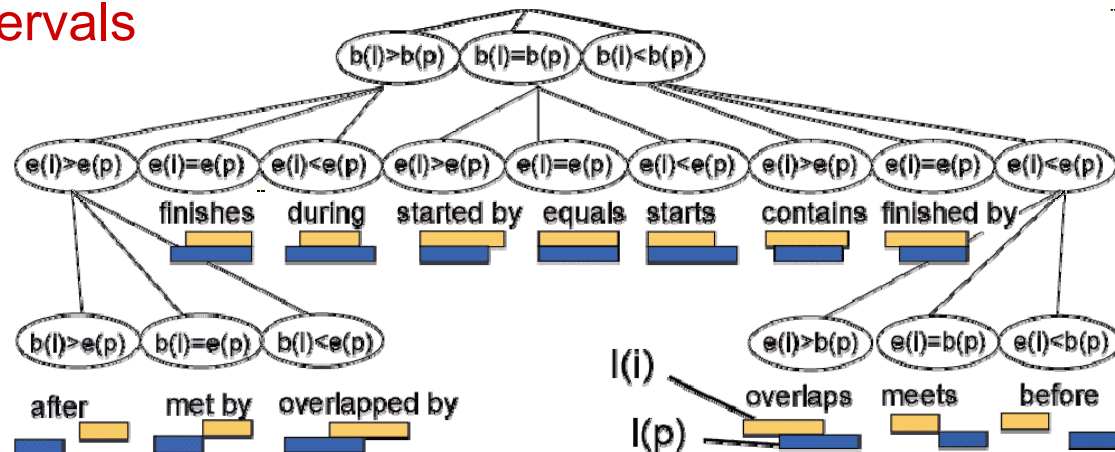
Utterance
Segmentation (Audio)
result
Lip motion
(Video)



Audio feature:
filter bank analysis
Image feature: PCA

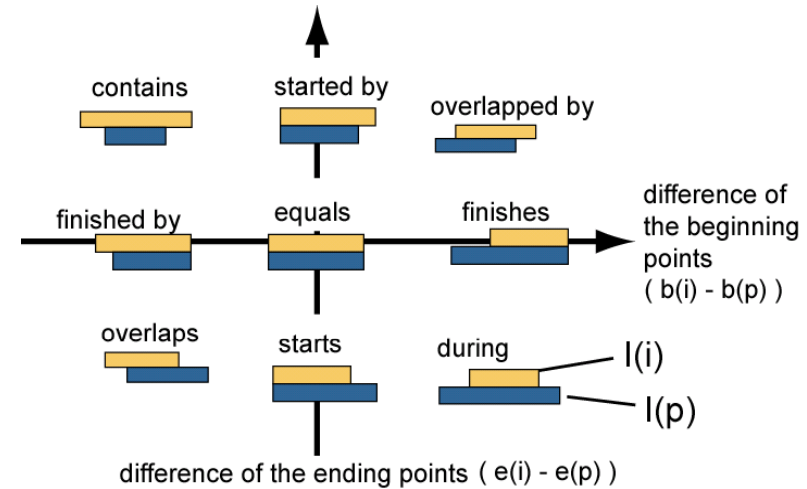
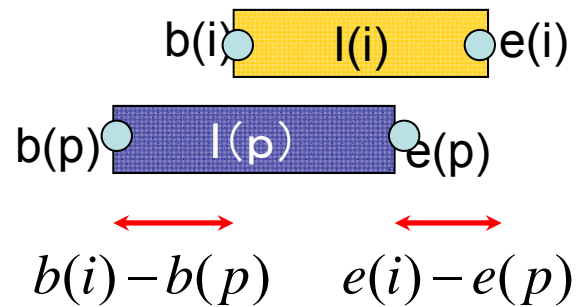
Temporal relation of two intervals
(without metric)

Focus on
overlapped
interval pairs

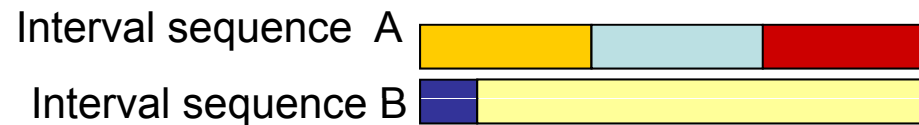


Metric Relation of Intervals

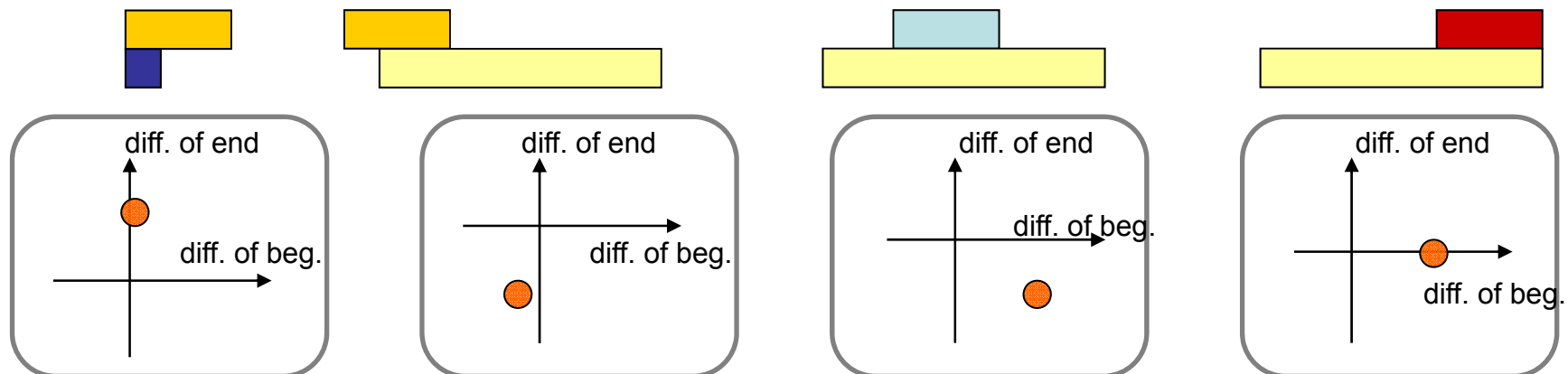
Overlapped interval pair



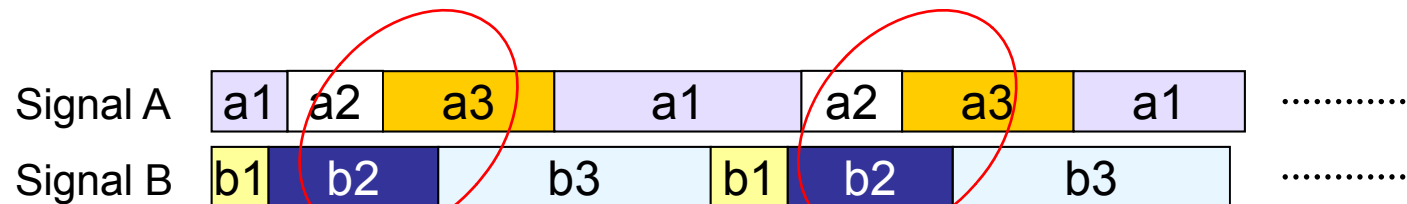
Example



Overlapped interval pairs

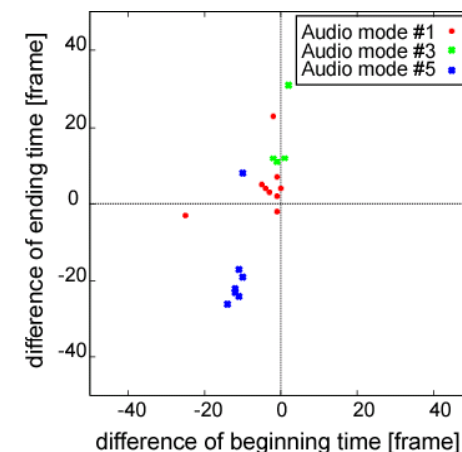
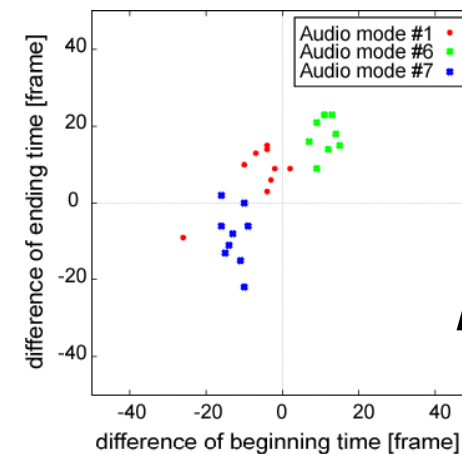
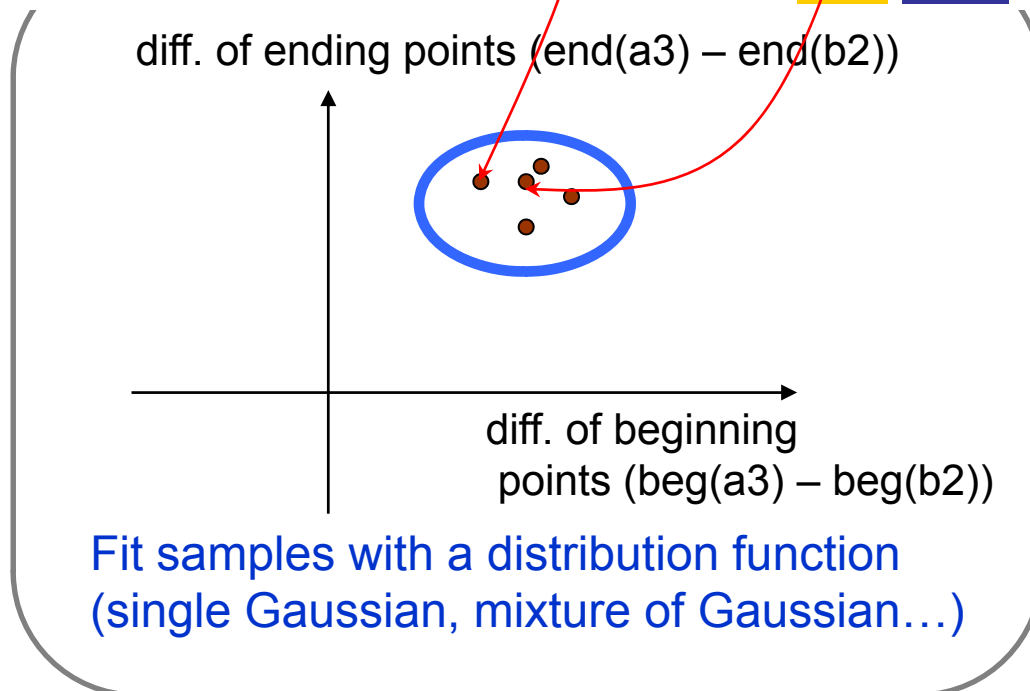


Timing Structure Model

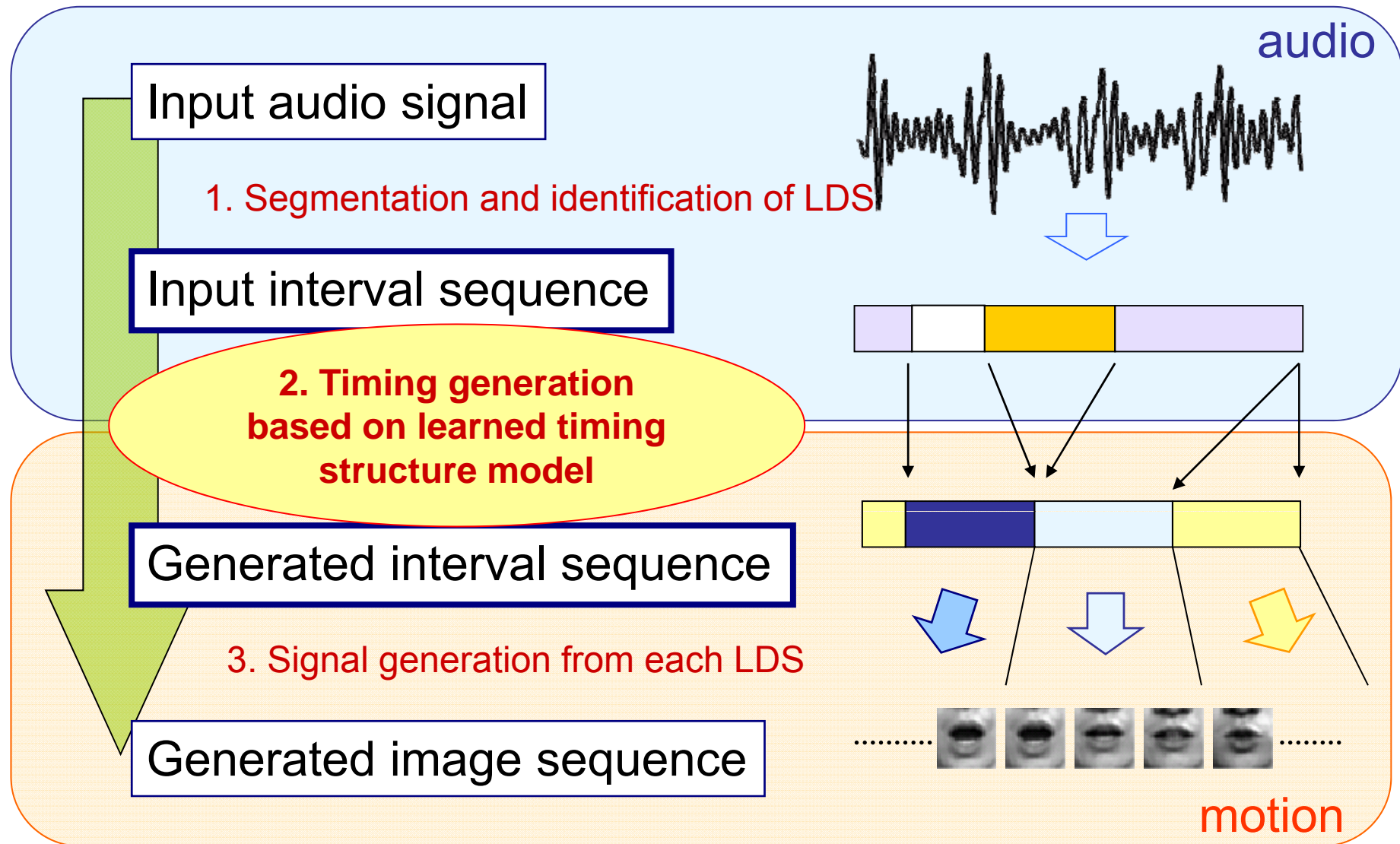


Learning: Prepare a distribution for each label pairs

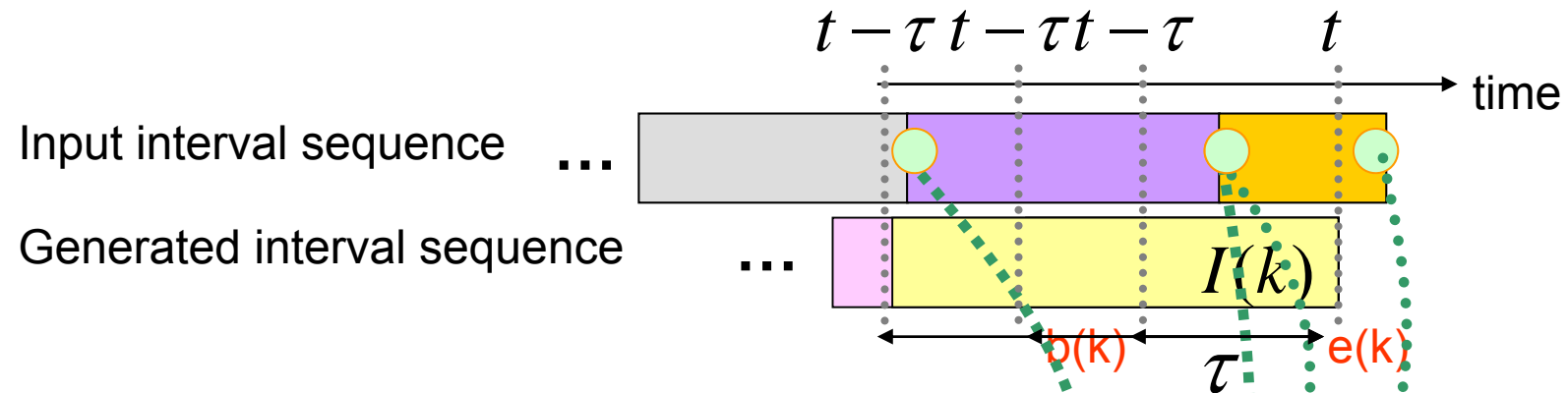
(ex.) **Temporal difference distribution** of pair (a3 b2)



Media Signal Conversion



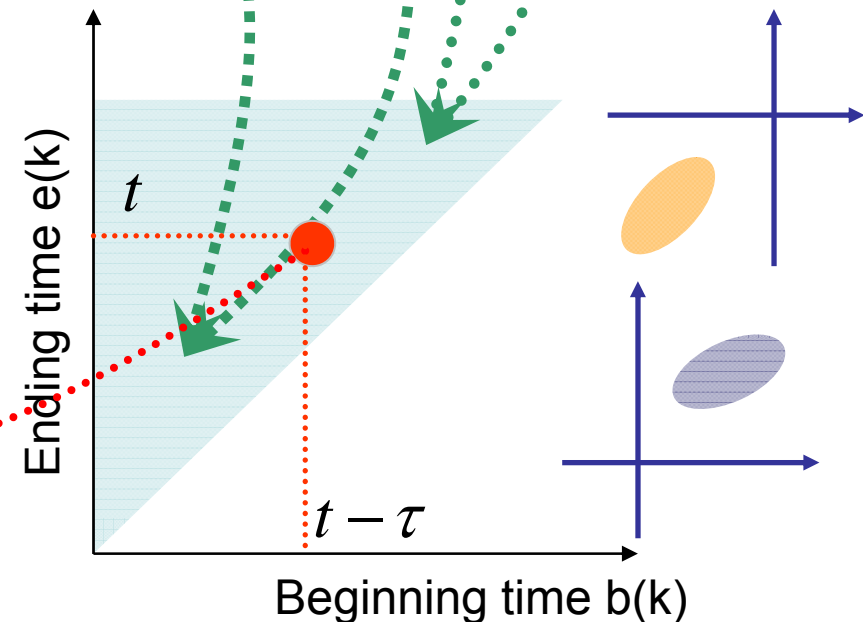
Timing Generation via DP Algorithm



Find maximum evaluate value
for every time t and labels

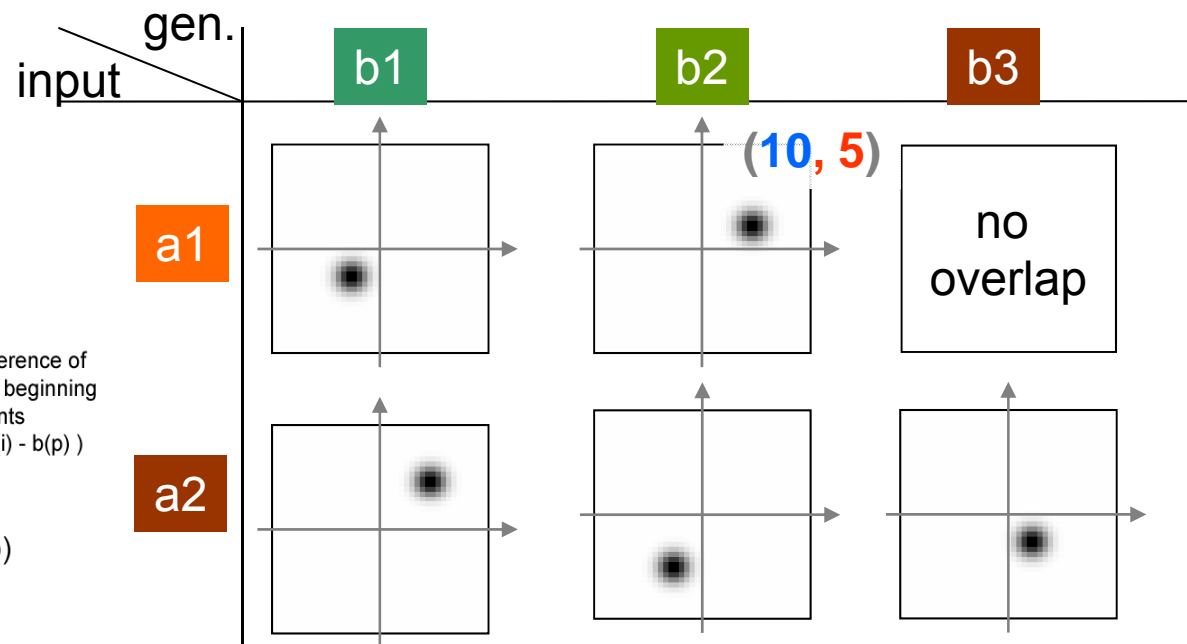
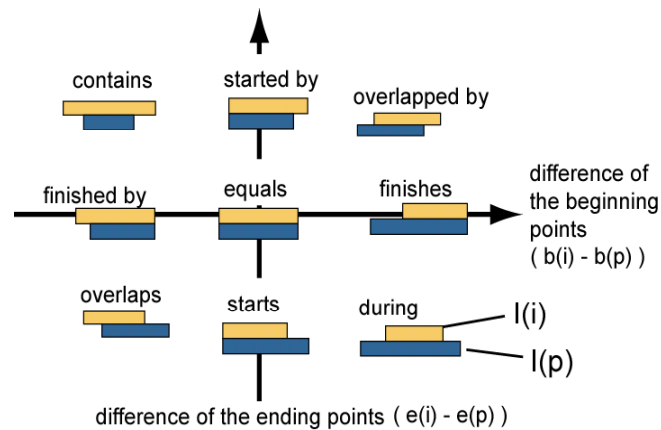
Maximized probability when the
interval ends at time t with label k

$$E(t, k) = \max_{\tau, p} \{ E(t, k | t - \tau, p) E(t - \tau, p) \}$$

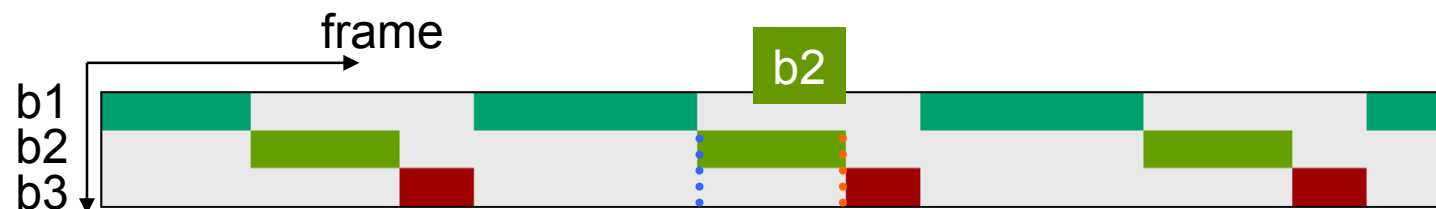


Verification of the Algorithm (Simulation)

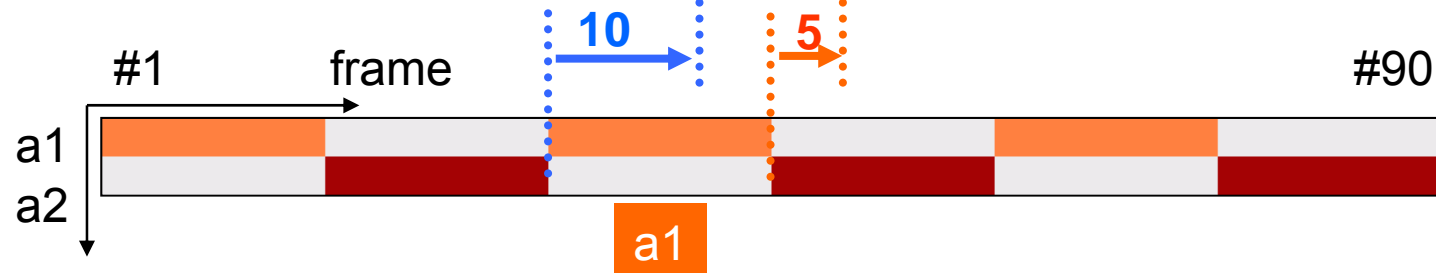
Temporal difference
distributions
(Gaussian functions)



Generated
interval seq.

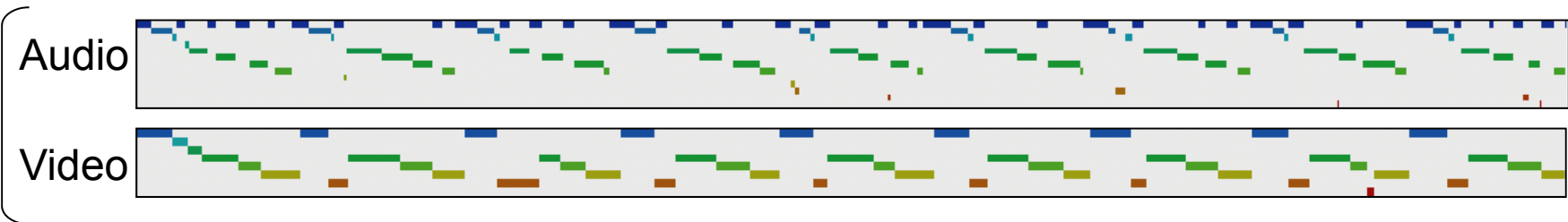


Input
interval seq.

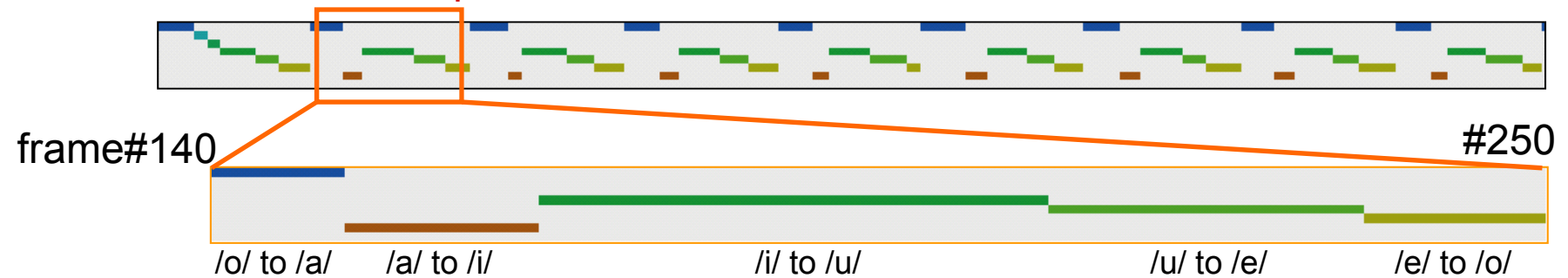


Lip Motion Generation from Audio Signal

Training data



Generated interval sequence



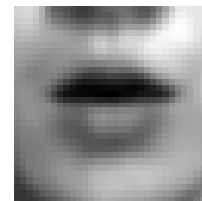
Generated image seq.



Original

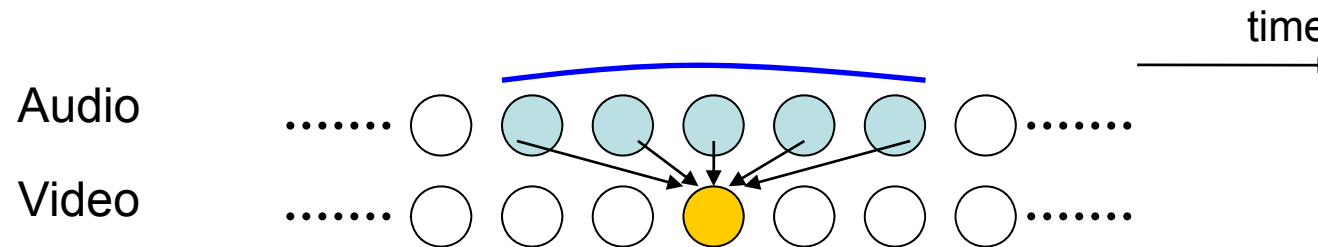


Generated

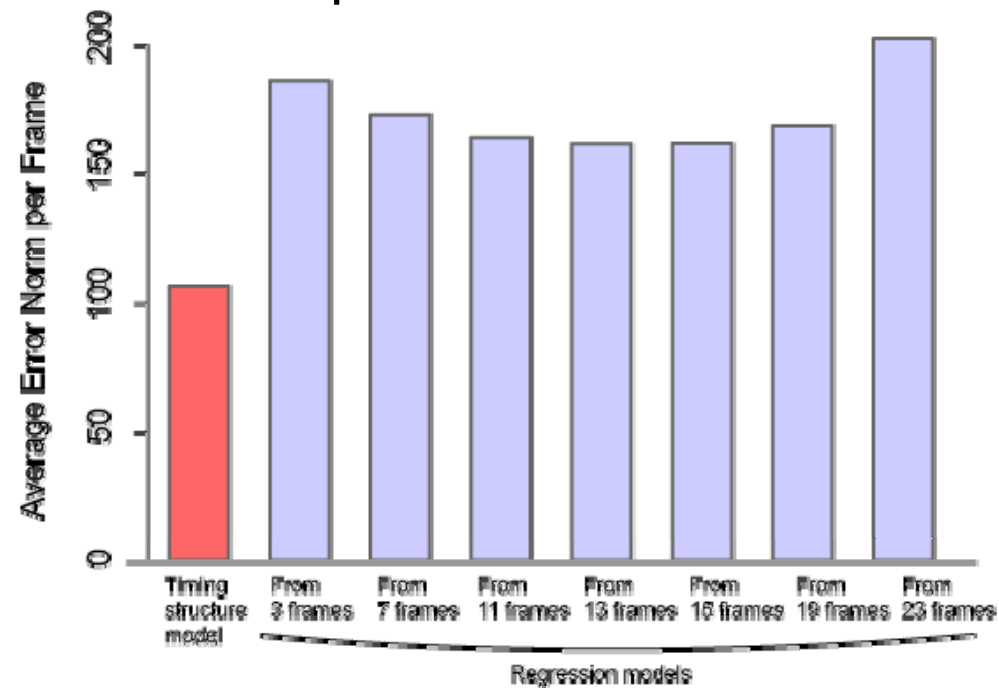


Comparison with Regression Models

- Linear regression models



- Average Error Norm per Frame



Pianist Motion Generation

Original



Generated



Discussion

- Timing structure model
 - Explicitly represents temporal metric relation between media signals
- Media conversion based on the timing structure model
 - Generates timing of one signal from other related media signals

Future work

- Apply to human-computer interaction (ex.)
 - audio-visual speech recognition
 - facial expression analysis
 - speaker detection in noisy environment
 - utterance timing generation for speech dialog system



Chapter 6

Conclusion

Summary

- Interval-based hybrid dynamical systems
 - integrate **discrete-event systems (subjective time)** and **dynamical systems (objective (physical) time)**
 - explicitly model temporal relations such as
 - tempo and rhythm in a signal
 - timing structure among different media signalsbased on *temporal intervals*
- Two-step learning method
 - Clustering of dynamical systems based on eigenvalue constraints
 - Refinement of parameters via the EM algorithm

Future Work

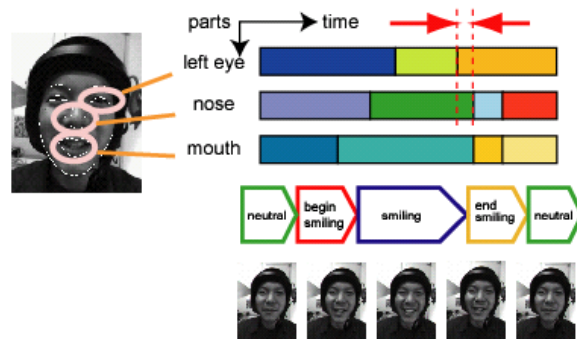
- Non-linear dynamical systems
 - Kernel method, neural networks
- Transition process between dynamics
 - Smooth signal generation
- Timing structure among more than three signals
 - Determination of causal relationship
 - Hidden interval sequence
- Hierarchical structures
 - Context-free grammar, hierarchical HMM
 - Variable length N-gram

Modeling Multiparty Interaction

Modeling Single Human Behavior from Multi-Channel Signals

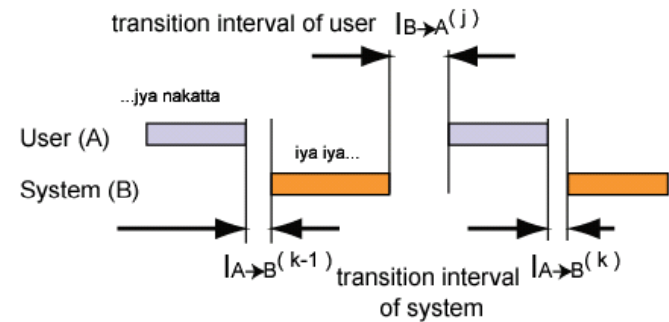
Single-Media Signal

Analysis of Multipart Events in Single Modality (Chapter 4)



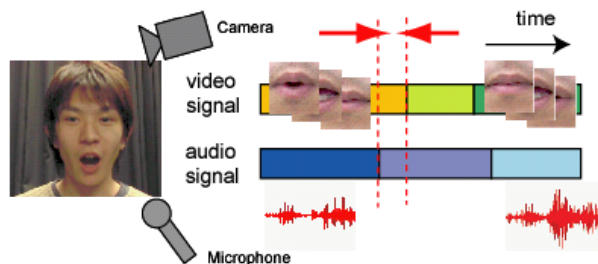
Modeling Multiparty Interaction

Modeling Utterance Timing in Speech Conversation



Multi-Media Signal

Modeling Multimodal Events (Chapter 5)



Modeling Multimodal Interaction

