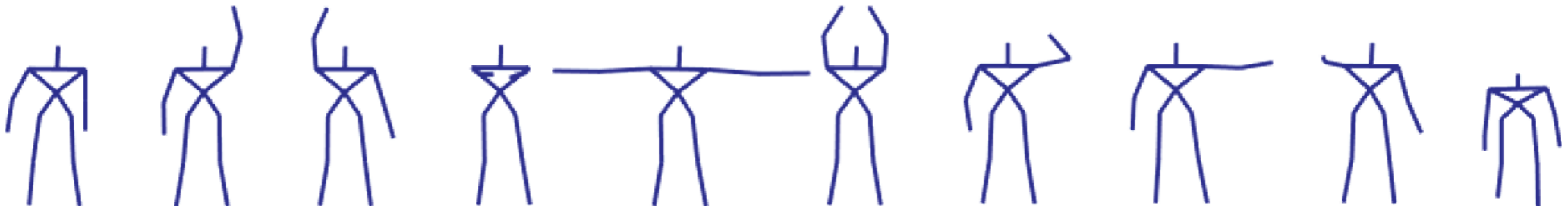# Simplified training for gesture recognition

Romain Faugeroux

*LIX, École Polytechnique*

Thales Vieira (presenter)

Dimas Martinez

*Mathematics, UFAL*

Thomas Lewiner

*Mathematics, PUC-Rio*

# Human Gesture Recognition

*Simplified training for gesture recognition*

# Human Gesture Recognition

*Simplified training for gesture recognition*

# Human Gesture Recognition

*Simplified training for gesture recognition*

# Human Gesture Recognition

*Simplified training for gesture recognition*

# Current Scenario


*Microsoft Kinect Sensor*

Popularization of real time depth sensors



Development of high quality Natural User Interfaces (NUI)

*Simplified training for gesture recognition*

# Current Scenario



Microsoft Kinect Sensor

Popularization of real time depth sensors

# Challenges



Gesture are culture specific

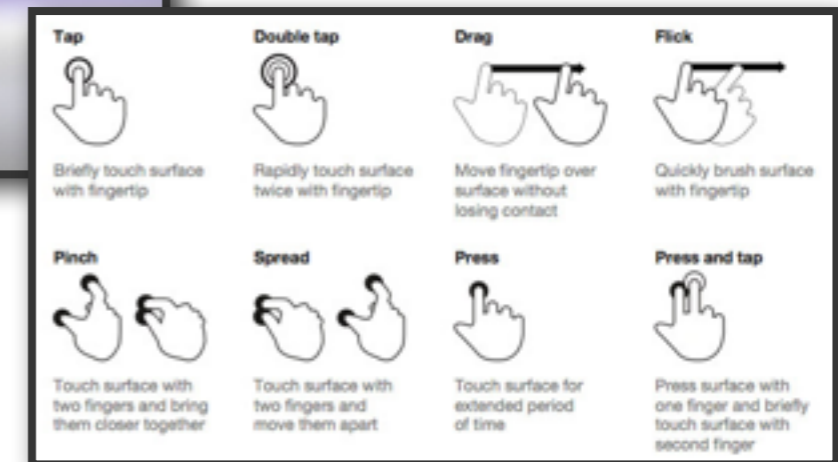Gestures can be performed at different speeds or sequences of poses

*Simplified training for gesture recognition*

# Challenges

Gesture are culture specific

Gestures can be performed at different speeds or sequences of poses

*Simplified training for gesture recognition*

# Challenges

Gesture are culture specific

Gestures can be performed at different speeds or sequences of poses

Solution: learning!

# Learning approaches
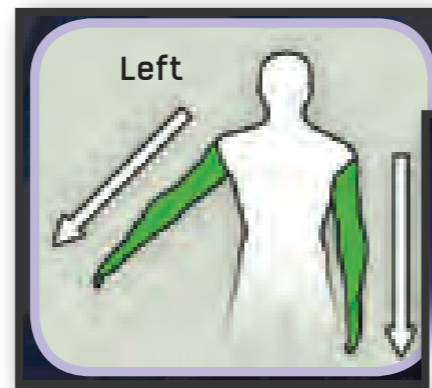
User learns from
  the machine

*Simplified training for gesture recognition*

# Learning approaches

User learns from
the machine

*Simplified training for gesture recognition*

# Learning approaches

User learns from
the machine



Left

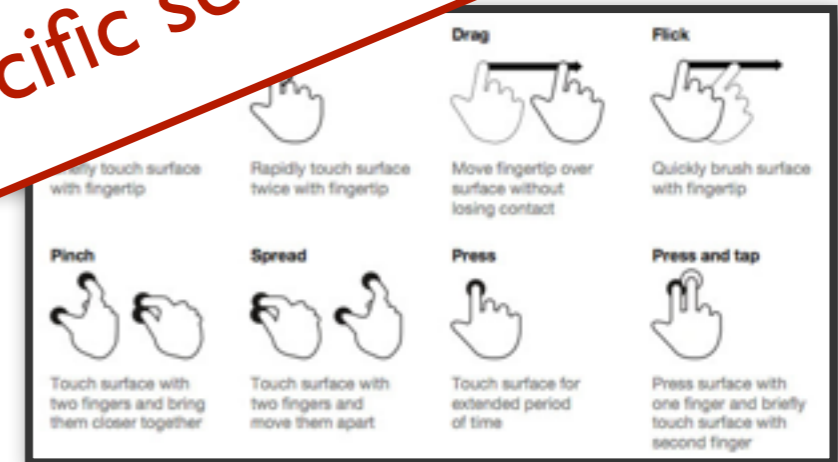Limited to a specific set of gestures

*Simplified training for gesture recognition*
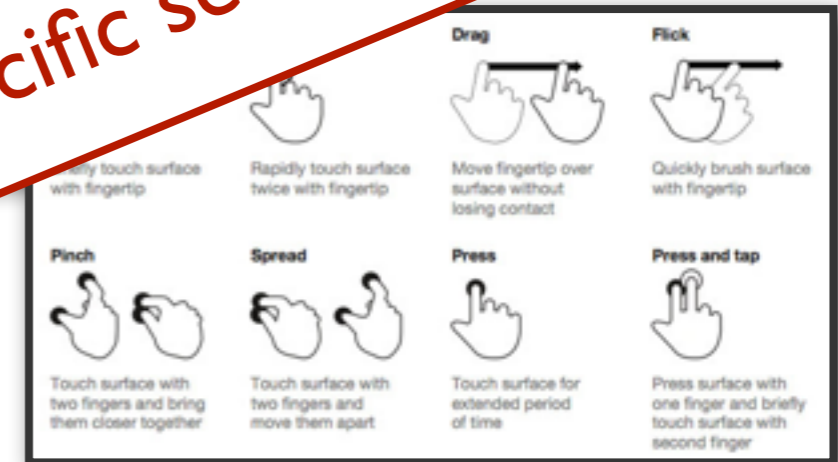
# Learning approaches

User learns from
the machine



Limited to a specific set of gestures

Machine learns
from the user

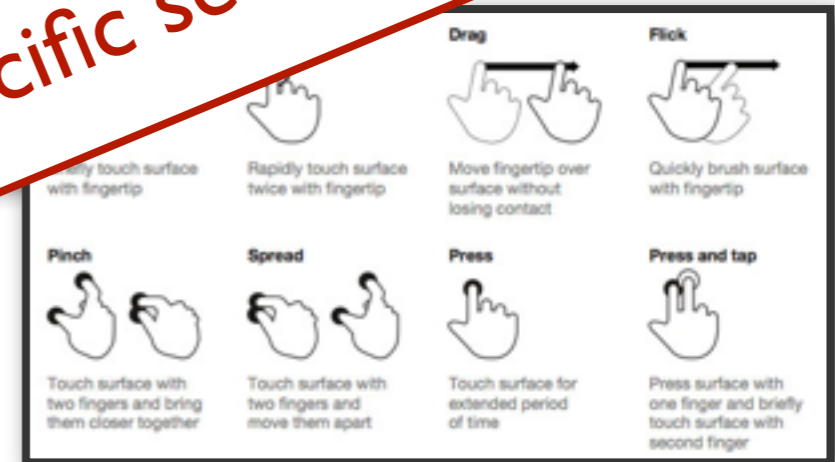*Simplified training for gesture recognition*

# Learning approaches

User learns from
the machine



Limited to a specific set of gestures

Machine learns
from the user

our focus

*Simplified training for gesture recognition*

# Learning approaches

User learns from
the machine



Limited to a specific set of gestures

Requires a training phase: often tedious
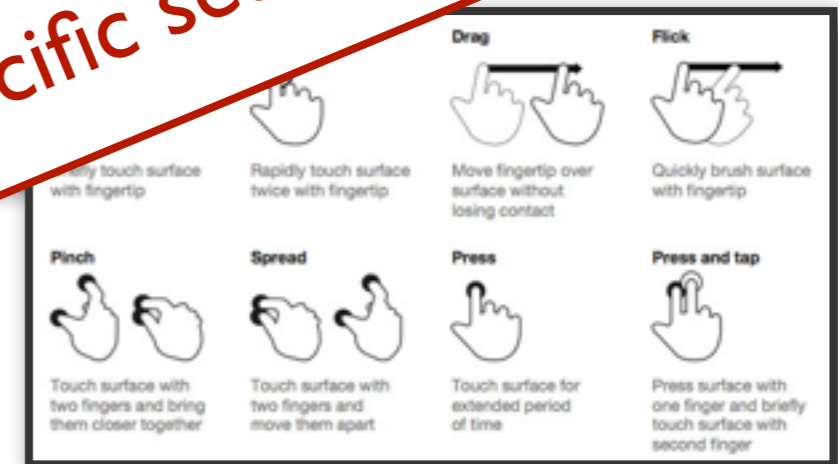
Machine learns
from the user

our focus

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Learning approaches

User learns from
the machine



Limited to a specific set of gestures

Requires a training phase: often tedious

Machine learns
from the user

our focus

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Related Work: global methods

space-time accumulation / spatio-temporal templates



Vieira *et al* (2014)



Bobick and Davis (2001)

Training phase required for gestures only

Requires more computing resources

*Simplified training for gesture recognition*

# Related Work: global methods

space-time accumulation / spatio-temporal templates



Vieira *et al* (2014)

**Gesture segmentation is manual!**

Bobick and Davis (2001)

Training phase required for gestures only

Requires more computing resources
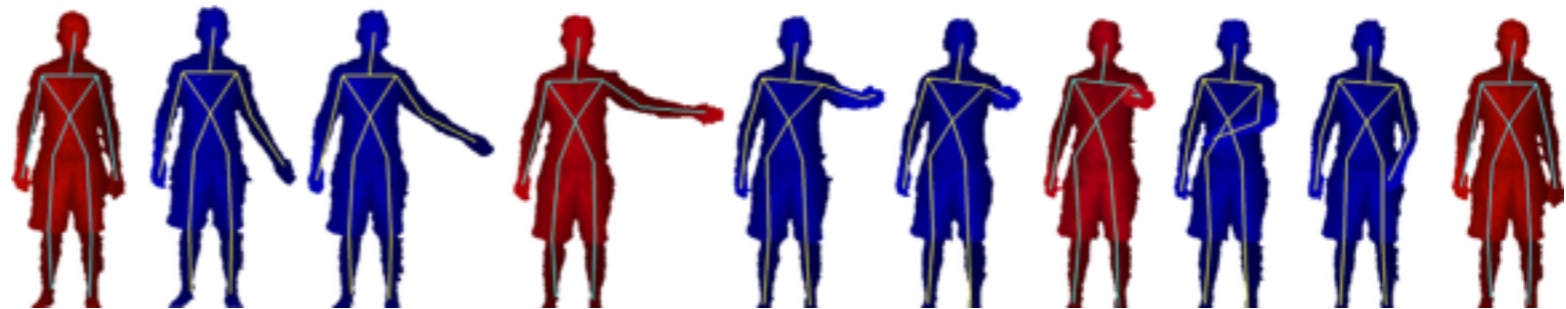
*Simplified training for gesture recognition*

# Related Work: local methods

Feature-based:  spatiotemporal interest points / **key poses**

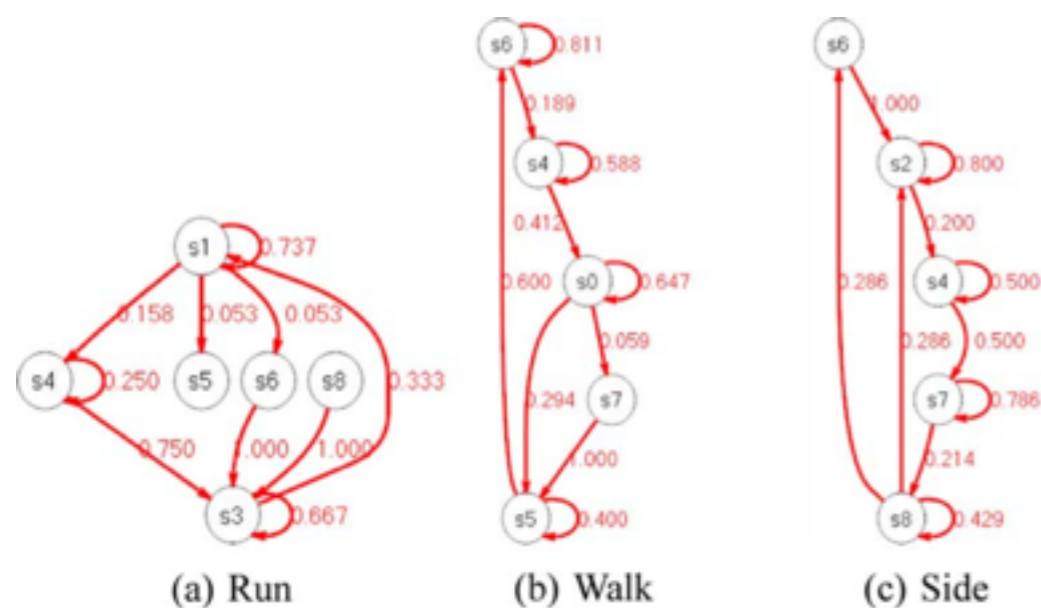salient postures
through clustering

manual key pose training



Miranda *et al* (2012)



(a) Run    (b) Walk    (c) Side

Li *et al* (2008)

Extremely fast recognition

Training phase usually required
for key pose learning

Discriminativity not considered for
trained gestures

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Related Work: local methods

Feature-based: spatiotemporal interest points / **key poses**

salient postures
through clustering

manual key pose training



Miranda *et al* (2012)

Li *et al* (2008)

(a) Run    (c) Side

**Gesture segmentation is manual!**

Extremely fast recognition

Training phase usually required
for key pose learning

Discriminativity not considered for
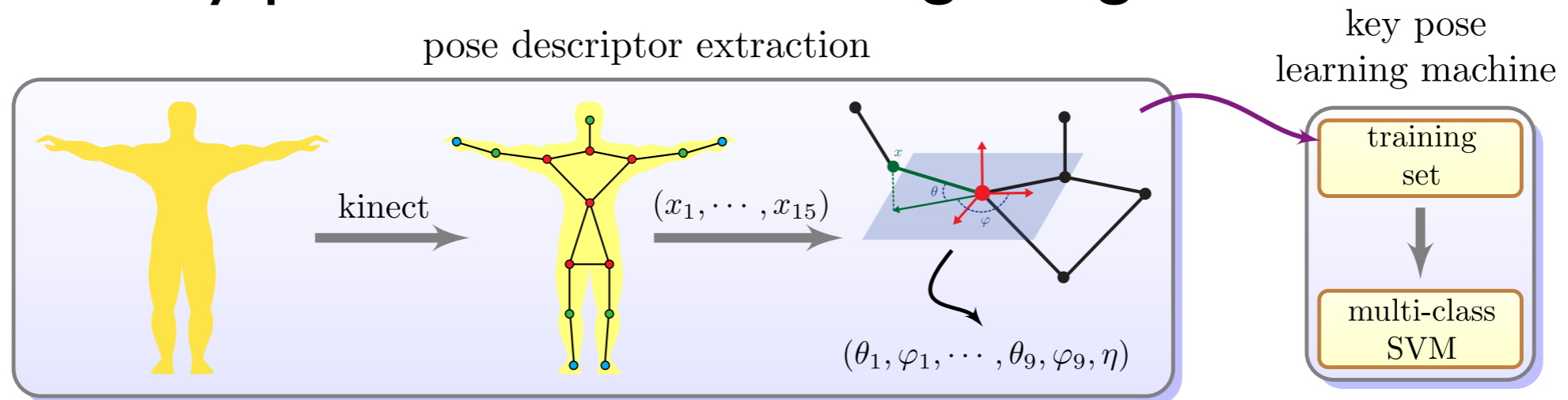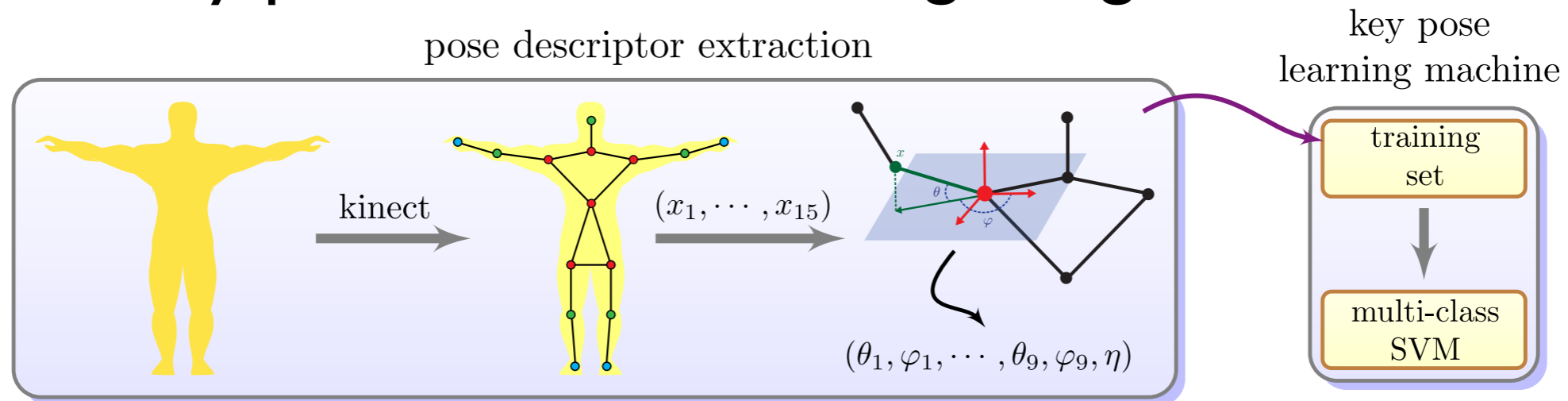trained gestures

*Simplified training for gesture recognition*

# Miranda *et al* (2012) training
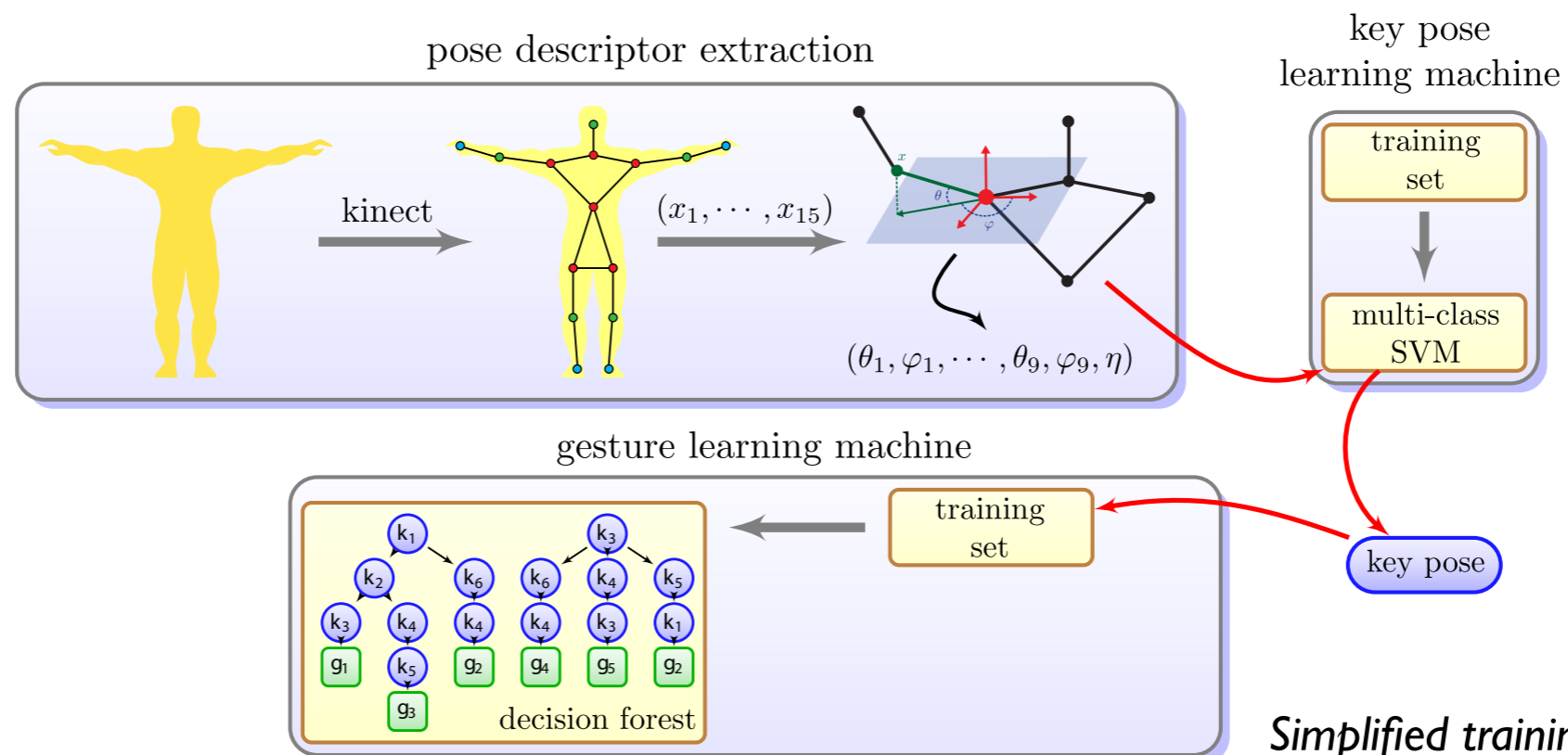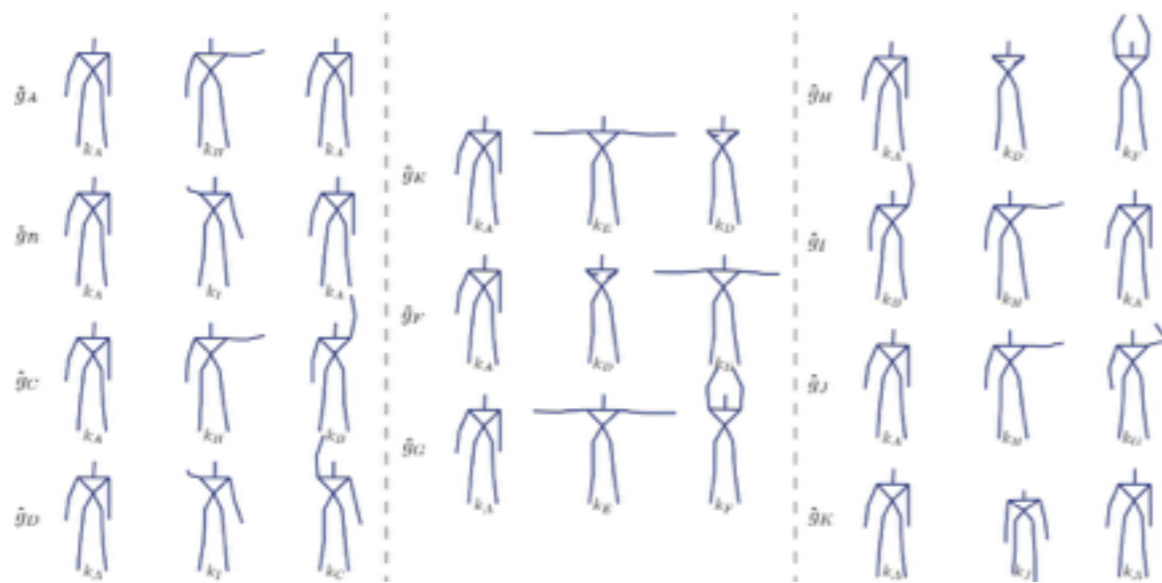
## 1 - manual key pose selection/training using SVMs



pose descriptor extraction

key pose
learning machine

kinect

$(x_1, \cdots, x_{15})$

$(\theta_1, \varphi_1, \cdots, \theta_9, \varphi_9, \eta)$

training set

multi-class SVM

*Simplified training for gesture recognition*

# Miranda *et al* (2012) training

## 1 - manual key pose selection/training using SVMs



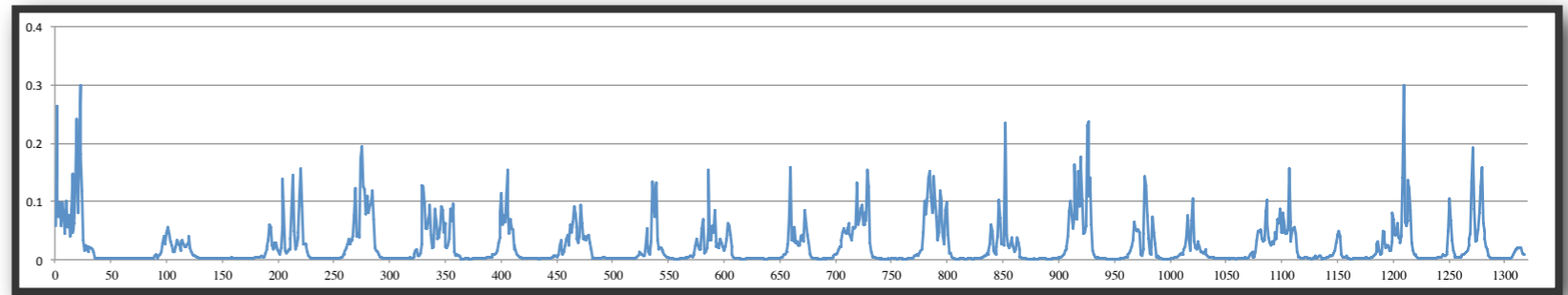pose descriptor extraction

key pose learning machine

$(x_1, \cdots, x_{15})$

$(\theta_1, \varphi_1, \cdots, \theta_9, \varphi_9, \eta)$

kinect

training set

multi-class SVM

## 2 - gesture training through key pose detection and decision forests



pose descriptor extraction

key pose learning machine

$(x_1, \cdots, x_{15})$

$(\theta_1, \varphi_1, \cdots, \theta_9, \varphi_9, \eta)$

kinect

training set

multi-class SVM

gesture learning machine

training set

key pose

decision forest

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Miranda *et al* (2012) training

1 - manual key pose selection/training using SVMs



2 - gesture training through key pose detection and decision forests



Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Contributions



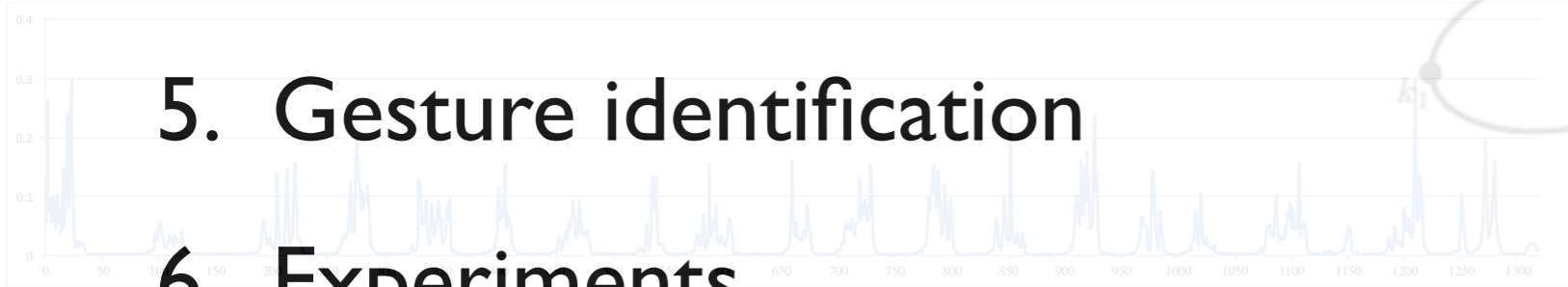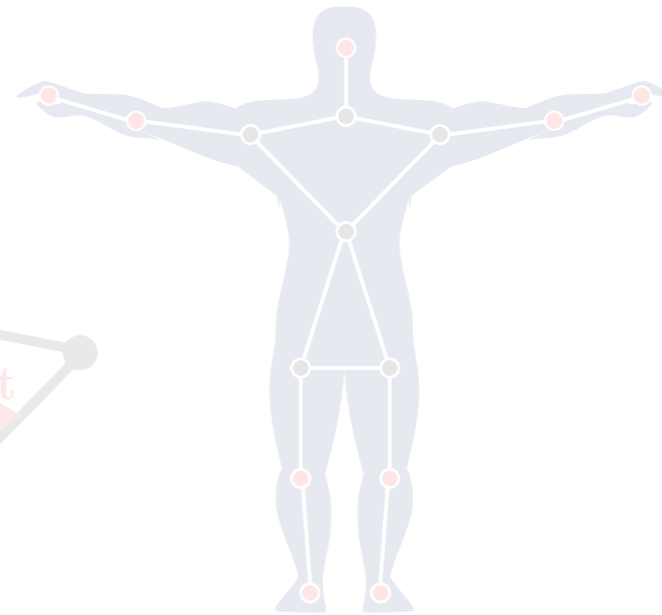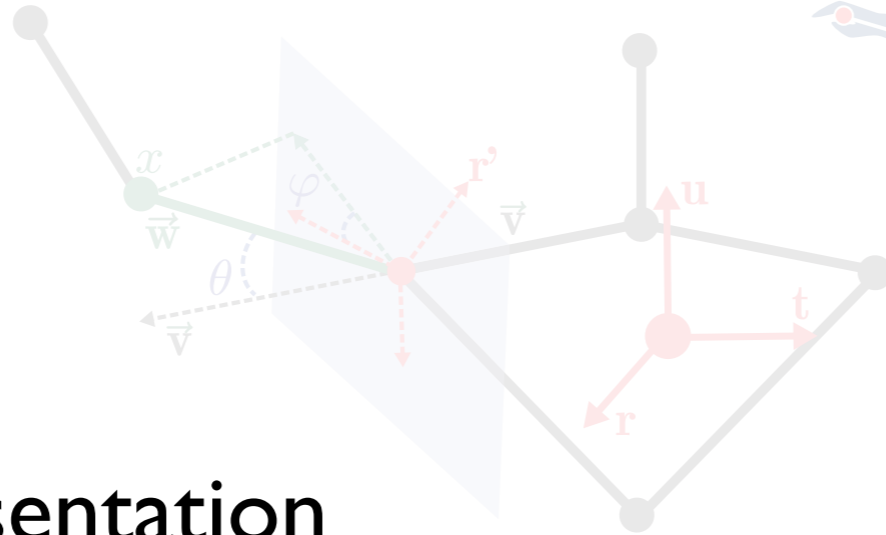Automatic gesture segmentation method

Automatic discriminative key pose selection method

Key pose based: Extremely fast recognition!

Minimal interaction training:  single record for all gestures!

*Simplified training for gesture recognition*

# Outline

1. Overview

2. Pose representation

3. Gesture segmentation

4. Discriminant key pose selection
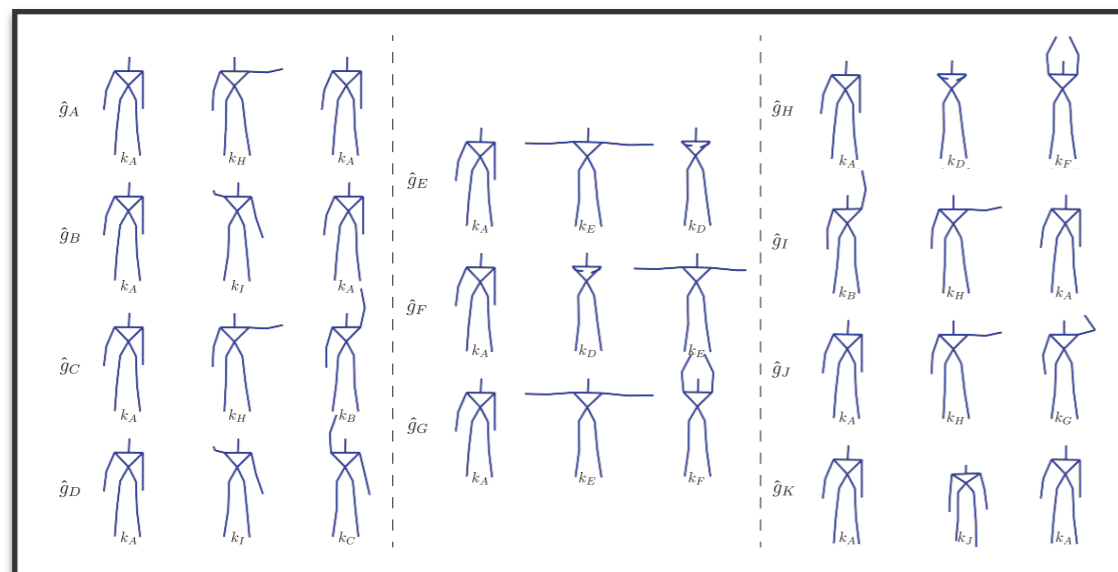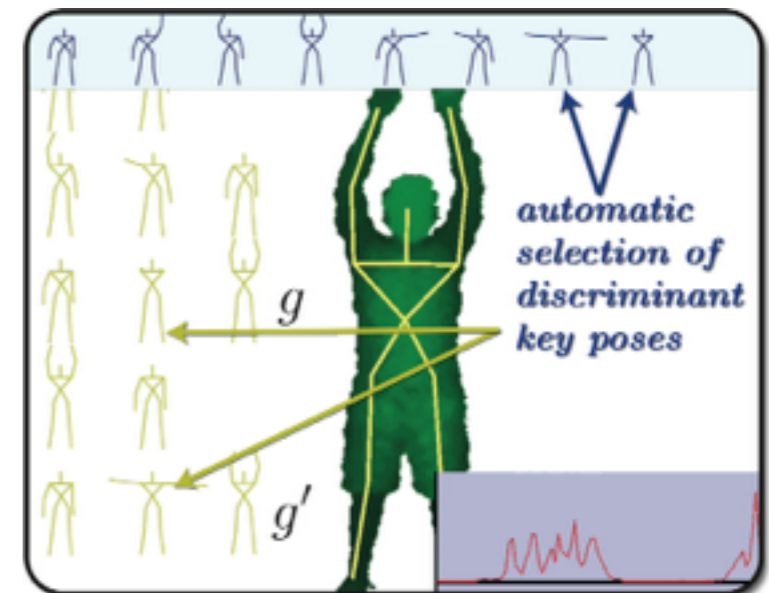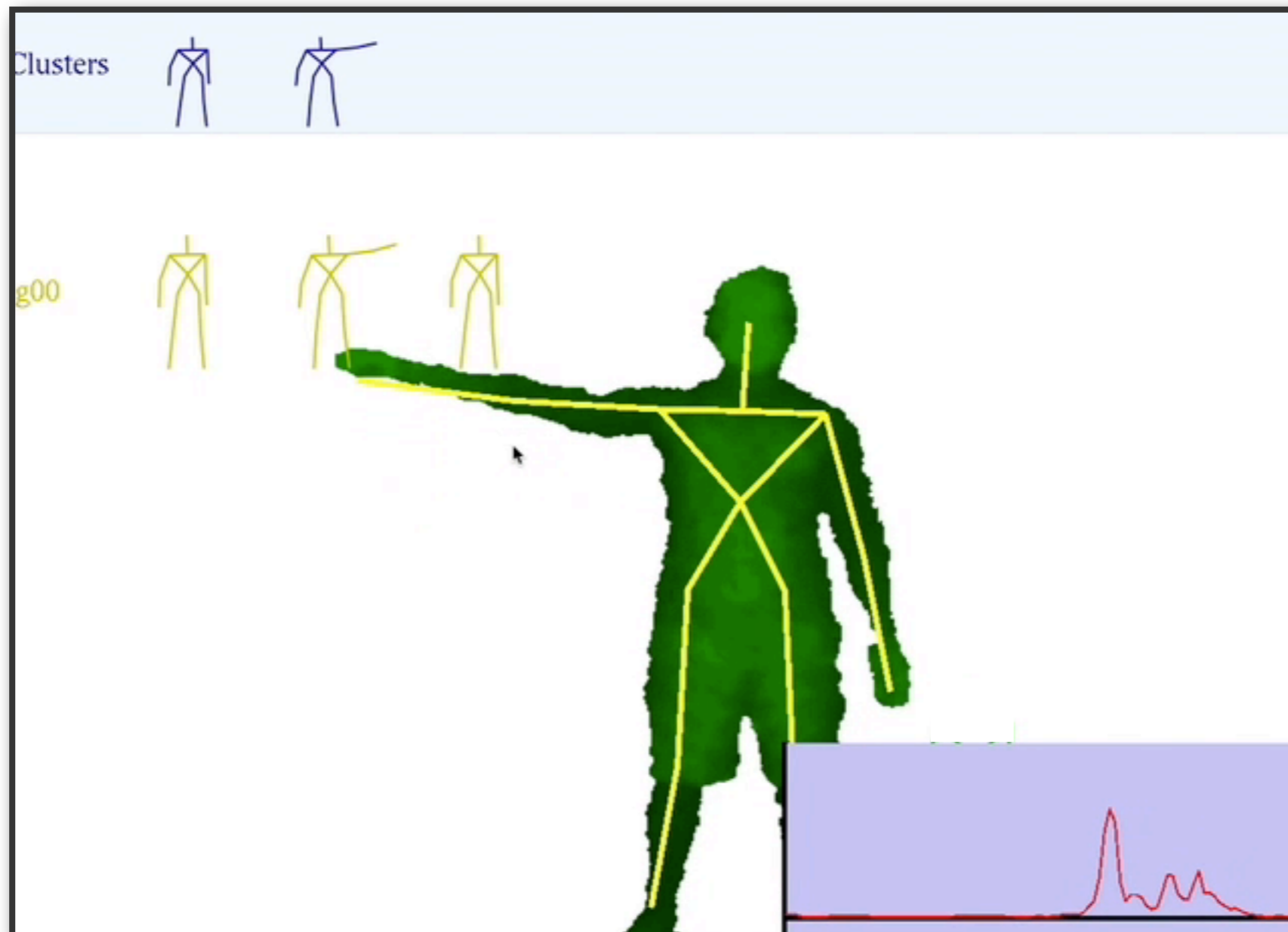
5. Gesture identification

6. Experiments

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Overview



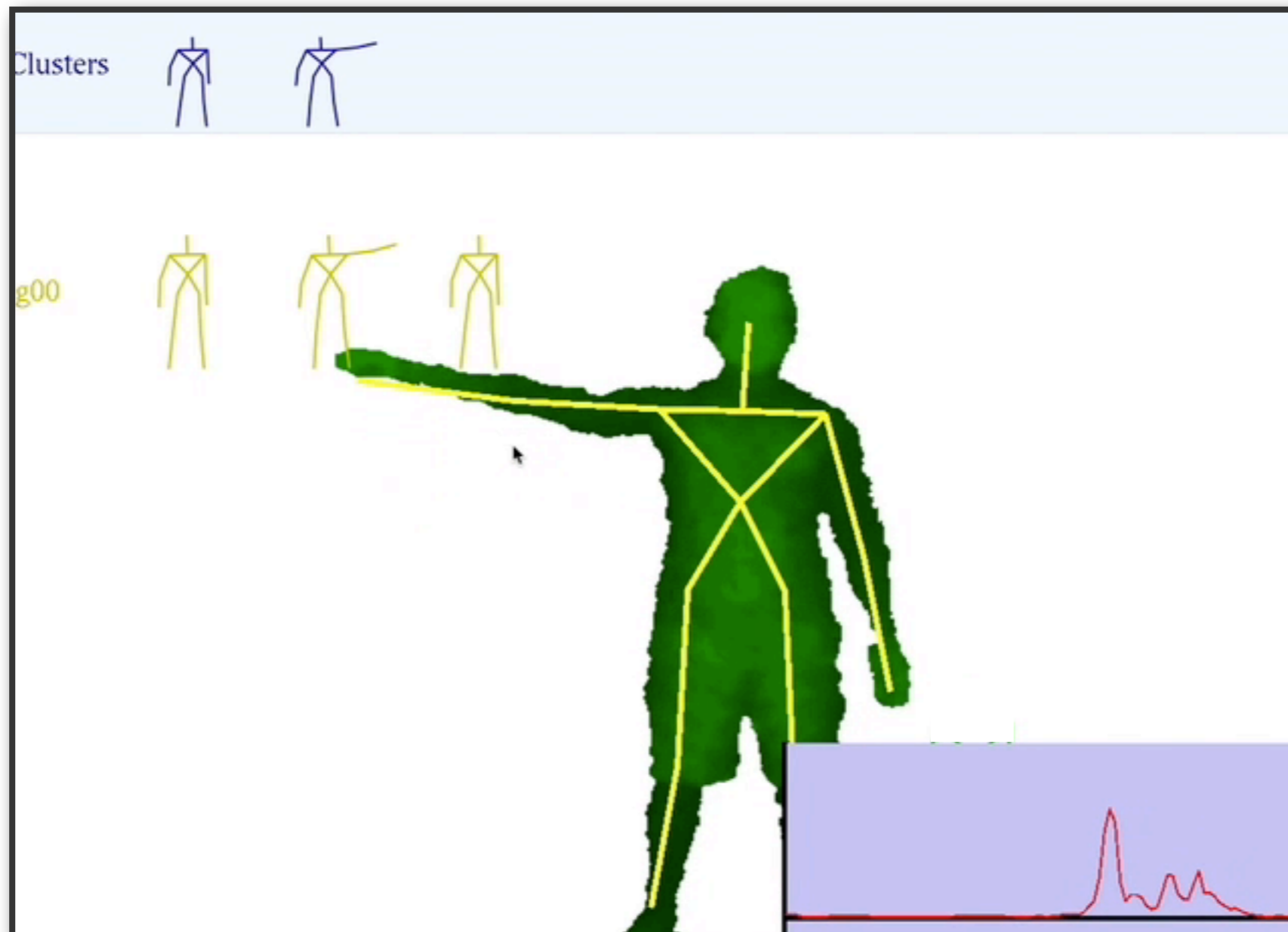automatic gesture segmentation

discriminant key pose selection

gesture identification

*Simplified training for gesture recognition*

# Overview

# Overview

*Simplified training for gesture recognition*
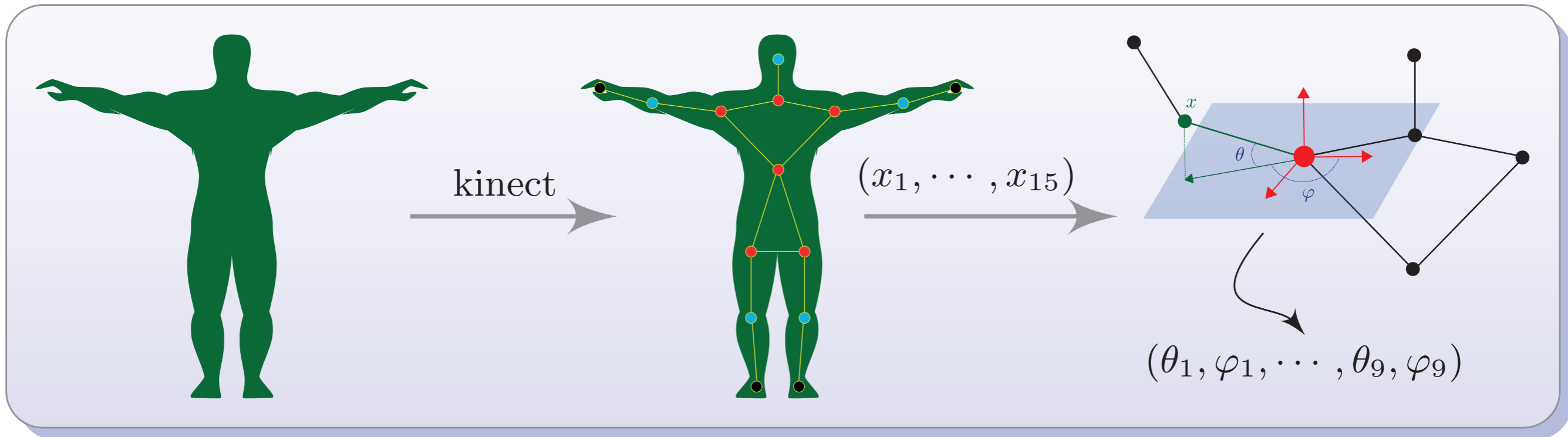
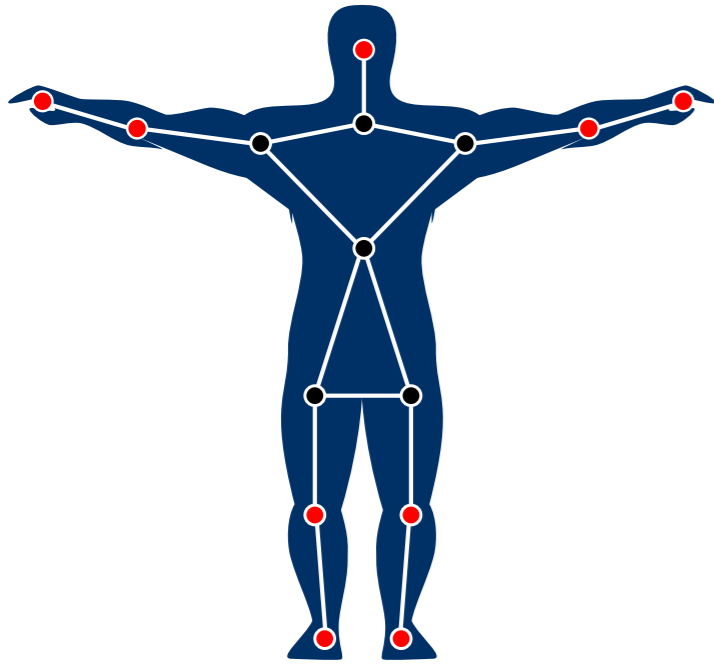# Pose representation



Real-time depth sensing system streaming depth data

OpenNI: public API to extract skeletons at 30fps
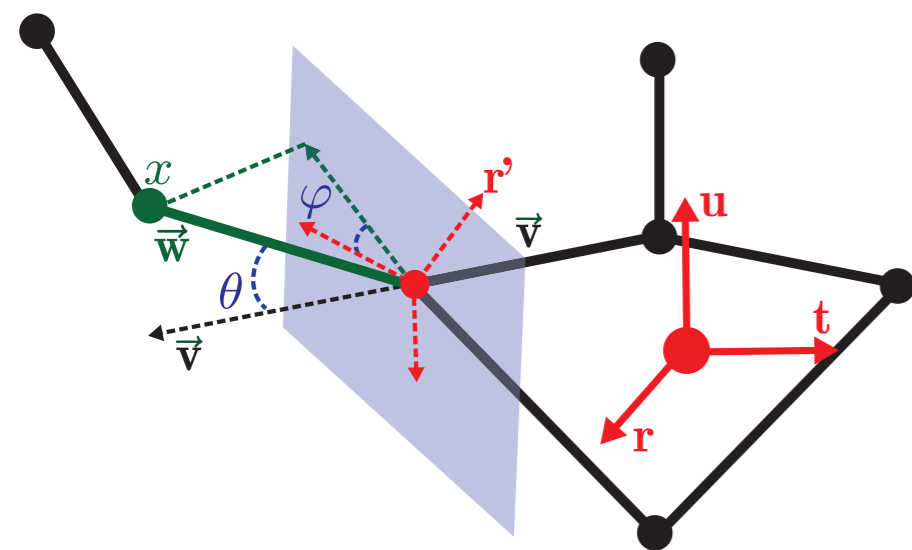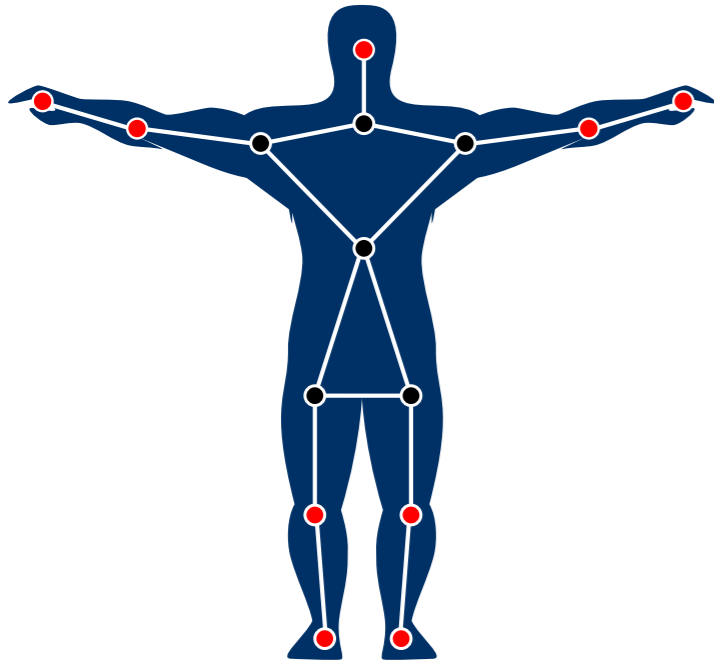
*Simplified training for gesture recognition*
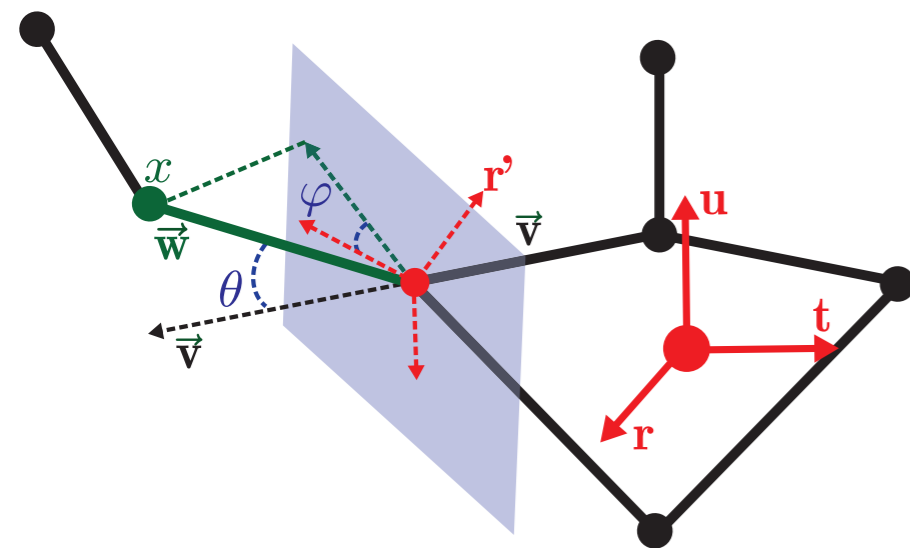
# Joint-angle pose descriptor

Miranda *et al* (2012): 9 relevant body joints converted to a list of spherical angles:

$$\text{pose:} \quad p \in \left(\mathbb{S}^2\right)^9$$

$$\text{gesture:} \quad \alpha : I \subset \mathbb{R} \mapsto \left(\mathbb{S}^2\right)^9$$

*Simplified training for gesture recognition*

# Joint-angle pose descriptor



Miranda *et al* (2012): 9 relevant body joints converted to a list of spherical angles:

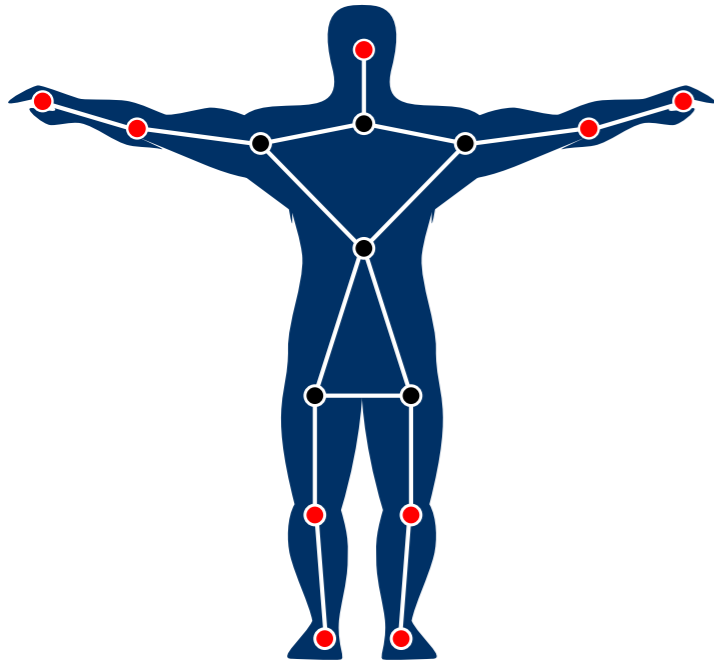$$\text{pose:} \quad p \in \left(\mathbb{S}^2\right)^9$$

$$\text{gesture:} \quad \alpha : I \subset \mathbb{R} \mapsto \left(\mathbb{S}^2\right)^9$$

Comparing joints l in distinct poses p and p':

$$\delta(p_l, p_l') = \arccos\left(\sin\theta_l \sin\theta_l' + \cos\theta_l \cos\theta_l' \cos|\varphi_l - \varphi_l'|\right)$$

*Simplified training for gesture recognition*

# Joint-angle pose descriptor



Miranda *et al* (2012): 9 relevant body joints converted to a list of spherical angles:
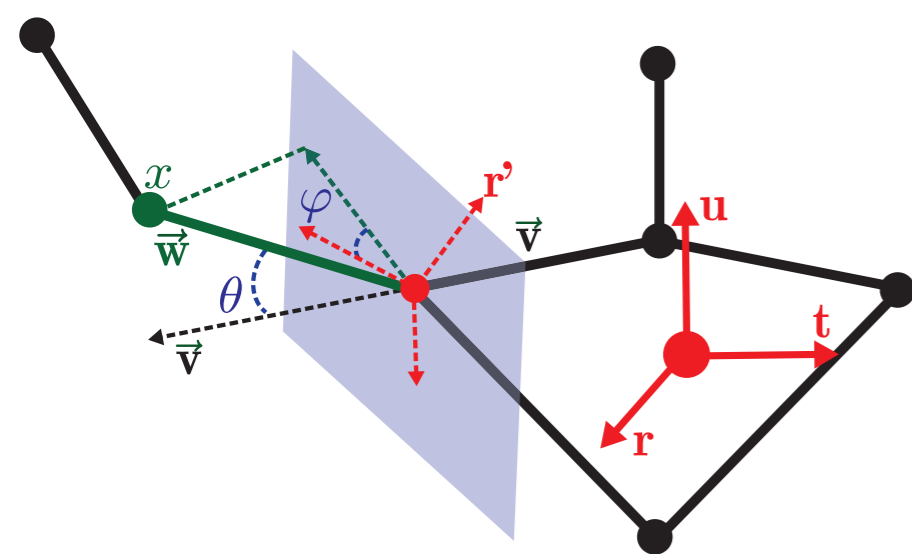
pose: $\quad p \in \left(\mathbb{S}^2\right)^9$

gesture: $\quad \alpha : I \subset \mathbb{R} \mapsto \left(\mathbb{S}^2\right)^9$



Comparing joints l in distinct poses p and p':

$$\delta(p_l, p_l') = \arccos\left(\sin\theta_l \sin\theta_l' + \cos\theta_l \cos\theta_l' \cos|\varphi_l - \varphi_l'|\right)$$

Distance between poses:

$$\Delta(p, p') = \sum_{l=1}^{9} \left[\delta\left(p_l, p_l'\right)\right]^2$$

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Gesture segmentation

Objective: avoid usual protocols requirements: neutral pose

*Simplified training for gesture recognition*

# Gesture segmentation

Objective: avoid usual protocols requirements: neutral pose

user inserts small pauses
in-between gestures

*Simplified training for gesture recognition*

# Gesture segmentation

Objective: avoid usual protocols requirements: neutral pose



user inserts small pauses in-between gestures
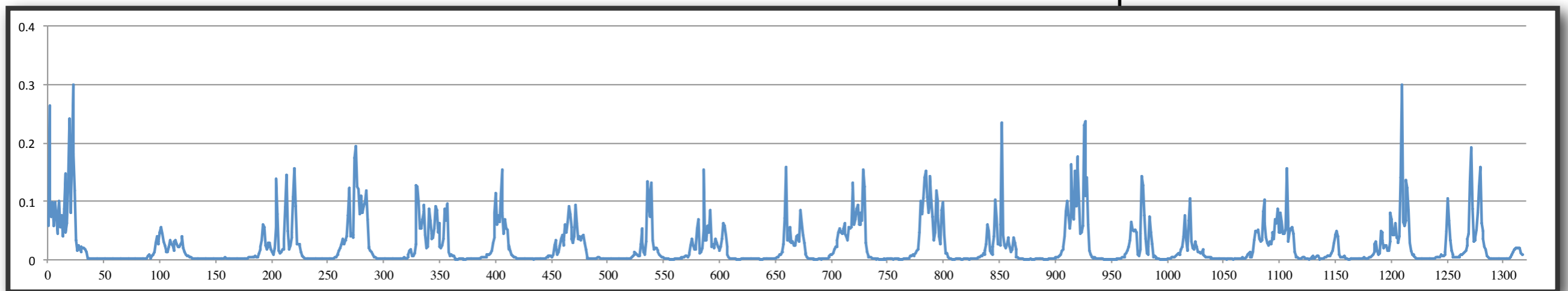
↓

depth sensor's random patterns generates rapid skeleton oscillations

*Simplified training for gesture recognition*

# Gesture segmentation

Objective: avoid usual protocols requirements: neutral pose



| user inserts small pauses in-between gestures |
| --- |

↓

| depth sensor's random patterns generates rapid skeleton oscillations |
| --- |

↓

| high curvature of the gesture curve |
| --- |

*Simplified training for gesture recognition*

# Gesture segmentation

Objective: avoid usual protocols requirements: neutral pose

user inserts small pauses in-between gestures



high curvature of the gesture curve

*Simplified training for gesture recognition*

# Curvature estimation

$$\alpha \colon I \subset \mathbb{R} \mapsto \left(\mathbb{S}^2\right)^9$$

*Simplified training for gesture recognition*

# Curvature estimation

$$\alpha: I \subset \mathbb{R} \mapsto \left(\mathbb{S}^2\right)^9$$

Pose encoded by cartesian coordinates
of 9 relevant joints:

$$\alpha: I \subset \mathbb{R} \mapsto \mathbb{R}^{27}$$

*Simplified training for gesture recognition*

# Curvature estimation

$$\alpha: I \subset \mathbb{R} \mapsto (\mathbb{S}^2)^9$$



Pose encoded by cartesian coordinates
of 9 relevant joints:

$$\alpha: I \subset \mathbb{R} \mapsto \mathbb{R}^{27}$$

First curvature:

$$\kappa(t) = \frac{\langle \alpha''(t), \mathbf{e}_2(t) \rangle}{\|\alpha'(t)\|^2}$$

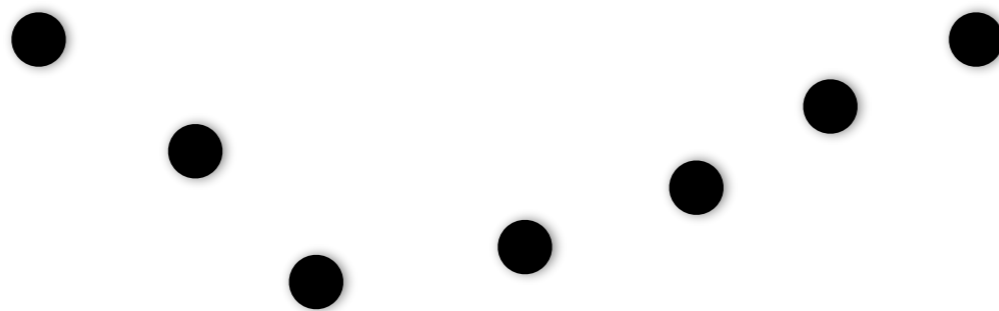where $\mathbf{e}_2(t)$ points in the direction of the first normal.

*Simplified training for gesture recognition*

# Curvature estimation

$$\alpha: I \subset \mathbb{R} \mapsto (\mathbb{S}^2)^9$$

Pose encoded by cartesian coordinates
of 9 relevant joints:

$$\alpha: I \subset \mathbb{R} \mapsto \mathbb{R}^{27}$$

First curvature:

$$\kappa(t) = \frac{\langle \alpha''(t), \mathbf{e}_2(t) \rangle}{\|\alpha'(t)\|^2}$$

where $\mathbf{e}_2(t)$ points in the direction of the first normal.

We need to estimate $\alpha'(t)$ and $\alpha''(t)$ in real time!
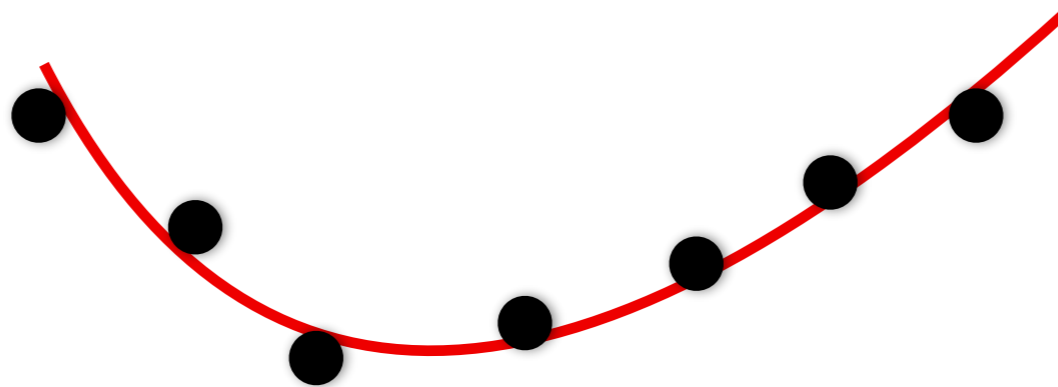
*Simplified training for gesture recognition*

# Parametric curve fitting

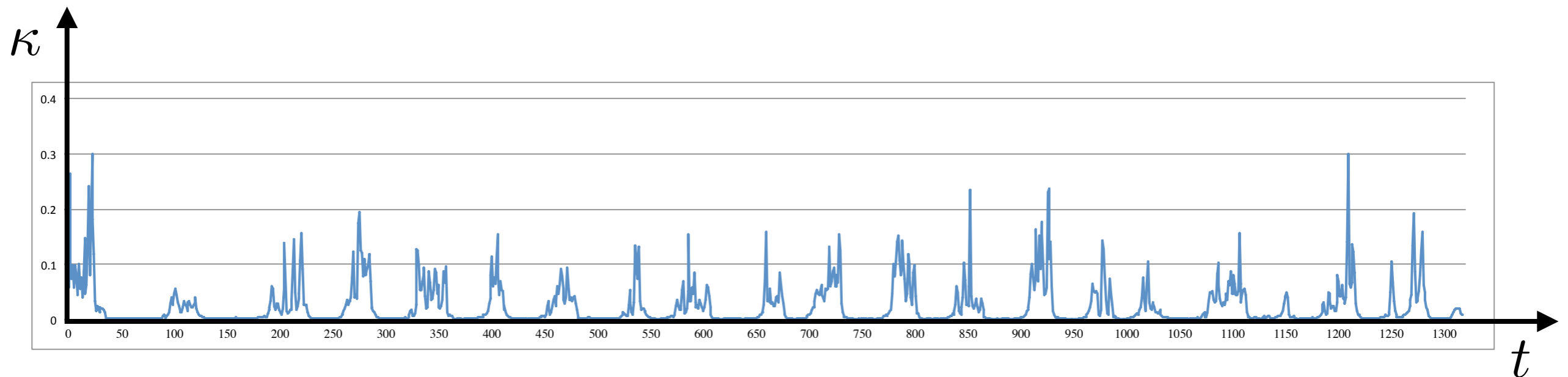Lewiner *et al* (2005): fit a portion of the gesture curve around $\alpha(t)$ to a parabola:

$$\tilde{\alpha}(s) = \alpha(t) + \tilde{\alpha}' \cdot s + \tfrac{1}{2} \cdot \tilde{\alpha}'' \cdot s^2$$

where $\tilde{\alpha}'$ and $\tilde{\alpha}''$ are estimates for the derivatives $\alpha'(t)$ and $\alpha''(t)$.
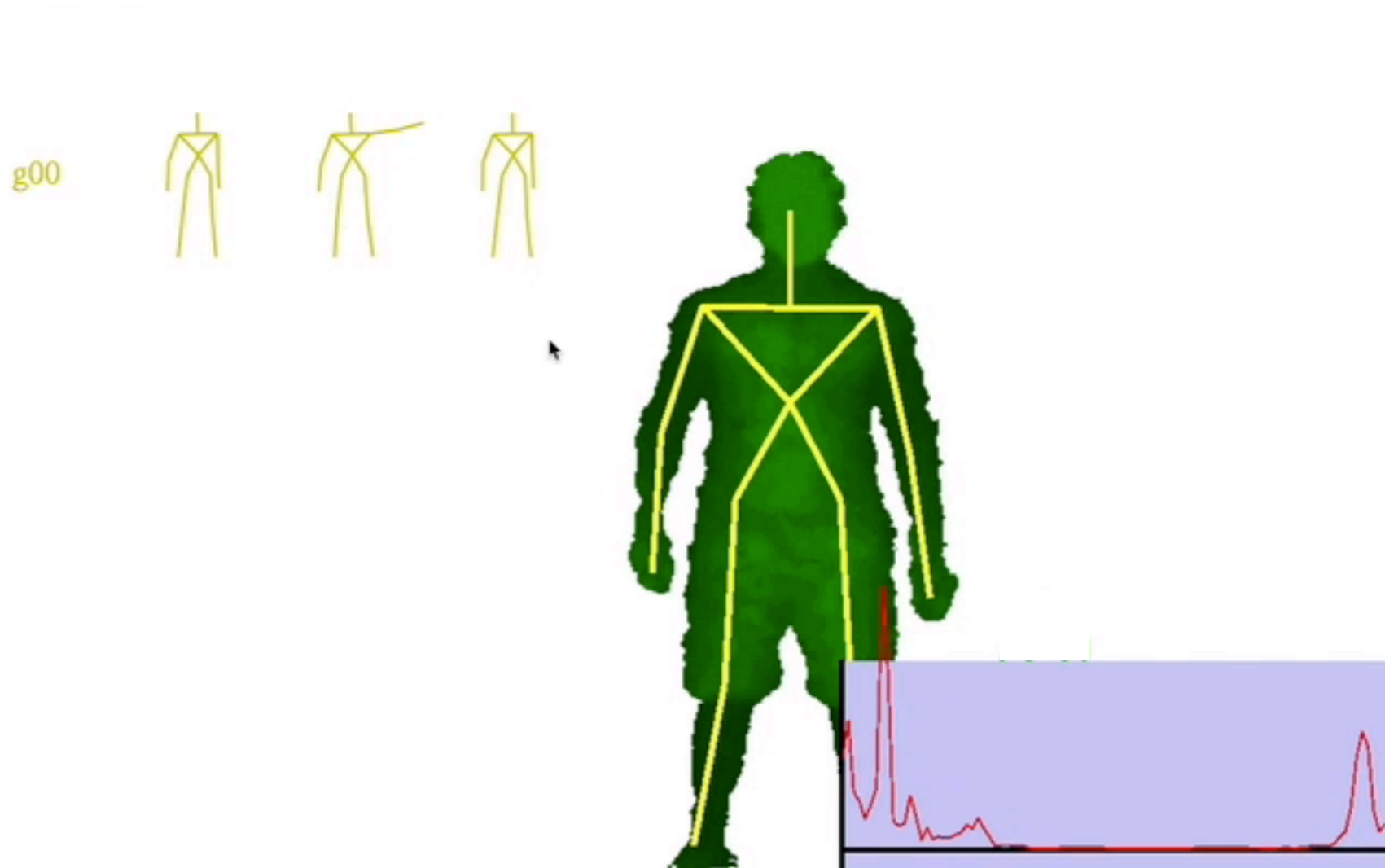
# Parametric curve fitting

Lewiner *et al* (2005): fit a portion of the gesture curve around $\alpha(t)$ to a parabola:

$$\tilde{\alpha}(s) = \alpha(t) + \tilde{\alpha}' \cdot s + \tfrac{1}{2} \cdot \tilde{\alpha}'' \cdot s^2$$

where $\tilde{\alpha}'$ and $\tilde{\alpha}''$ are estimates for the derivatives $\alpha'(t)$ and $\alpha''(t)$.

*Simplified training for gesture recognition*

# Parametric curve fitting

Lewiner *et al* (2005): fit a portion of the gesture curve around $\alpha(t)$ to a parabola:

$$\tilde{\alpha}(s) = \alpha(t) + \tilde{\alpha}' \cdot s + \tfrac{1}{2} \cdot \tilde{\alpha}'' \cdot s^2$$

where $\tilde{\alpha}'$ and $\tilde{\alpha}''$ are estimates for the derivatives $\alpha'(t)$ and $\alpha''(t)$.



Weighted least squares minimization: fast!

# Gesture segmentation



Simple thresholding

*Simplified training for gesture recognition*

# Gesture segmentation



Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Gesture segmentation

*Simplified training for gesture recognition*

# Outline

1. Overview

2. Pose representation

3. Gesture segmentation

4. <span style="color:red">Discriminant key pose selection</span>

5. Gesture identification

6. Experiments

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Gesture representation: key poses



Miranda *et al* (2012)



key pose set
$$\mathcal{K} = \{k_1, \ldots, k_n\}$$
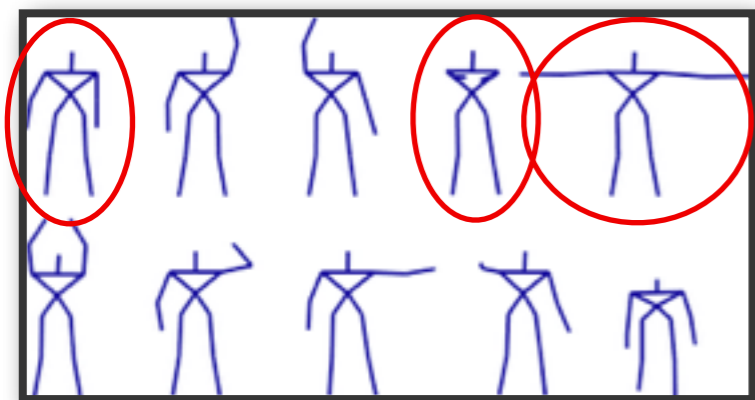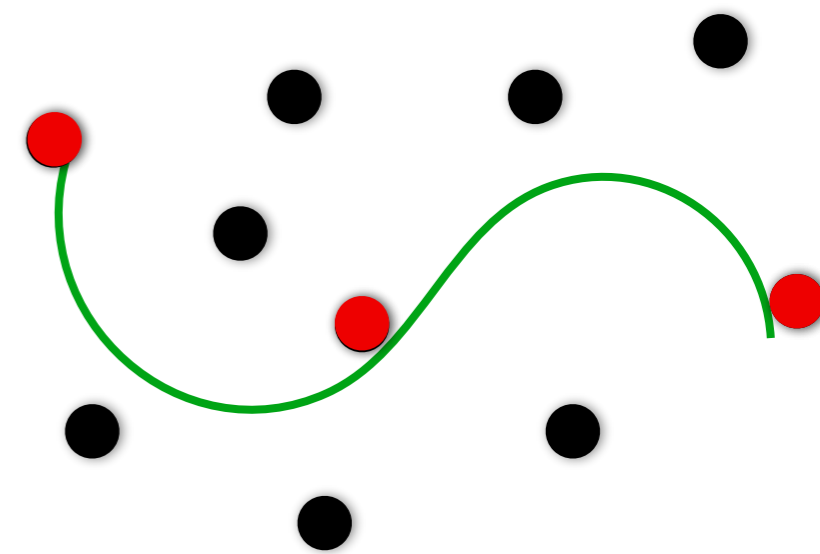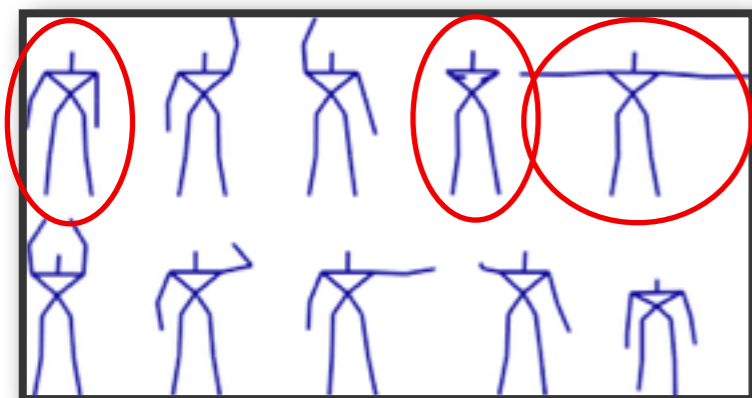
gesture
$$g = (p_1, p_2, \ldots)$$

gesture representation
$$\hat{g} = (k_i, k_j, \ldots)$$

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Gesture representation: key poses



Miranda *et al* (2012)
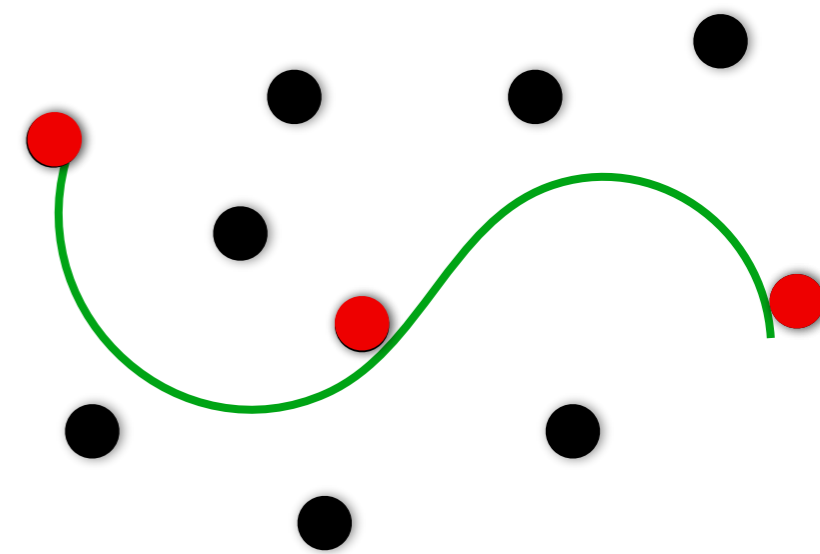


key pose set
$$\mathcal{K} = \{k_1, \ldots, k_n\}$$

gesture
$$g = (p_1, p_2, \ldots)$$

gesture representation
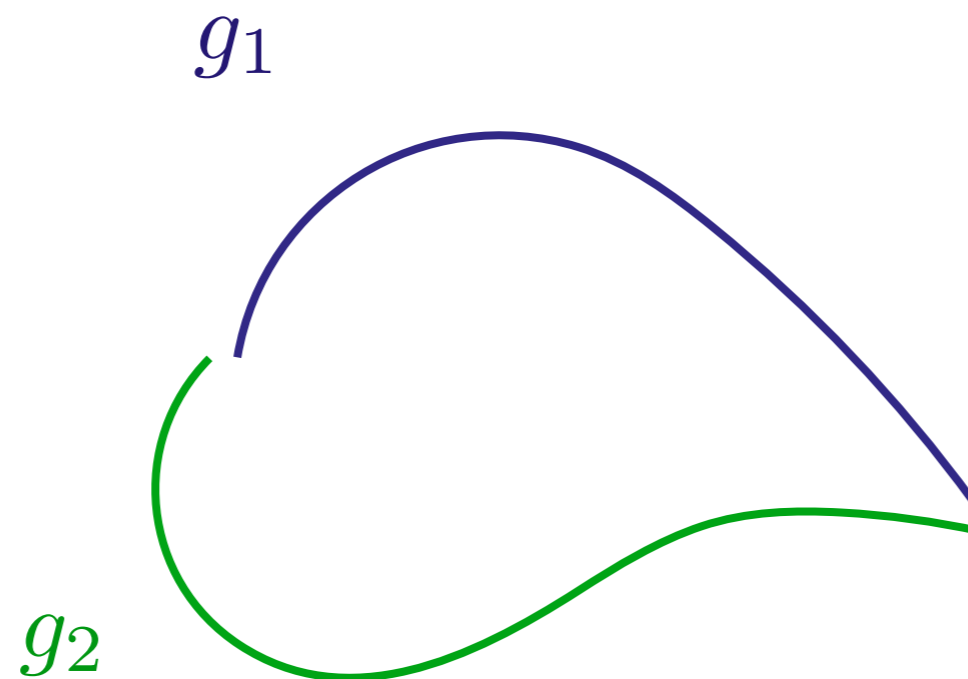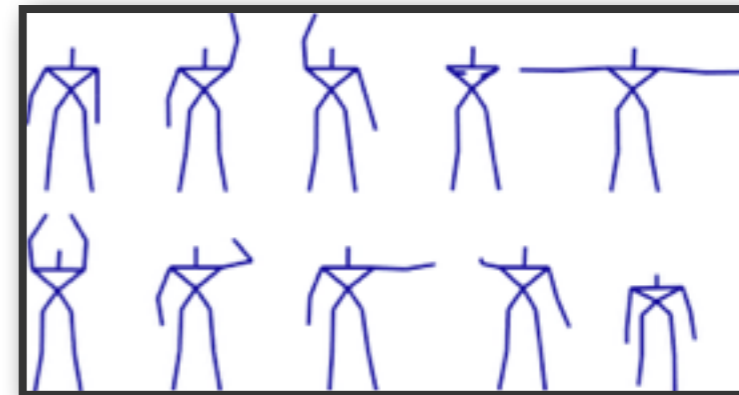$$\hat{g} = (k_i, k_j, \ldots)$$

Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Gesture representation: key poses



Miranda *et al* (2012)

$$\Delta(p, k_p) < \epsilon \quad ?$$

key pose set
$$\mathcal{K} = \{k_1, \ldots, k_n\}$$

gesture
$$g = (p_1, p_2, \ldots)$$

gesture representation
$$\hat{g} = (k_i, k_j, \ldots)$$

*Simplified training for gesture recognition*

# Gesture representation: key poses



Miranda *et al* (2012)

What is a good key pose set?

$\Delta(p, k_p) < \epsilon$ ?

key pose set
$\mathcal{K} = \{k_1, \dots, k_n\}$

gesture
$g = (p_1, p_2, \dots)$

gesture representation
$\hat{g} = (k_i, k_j, \dots)$
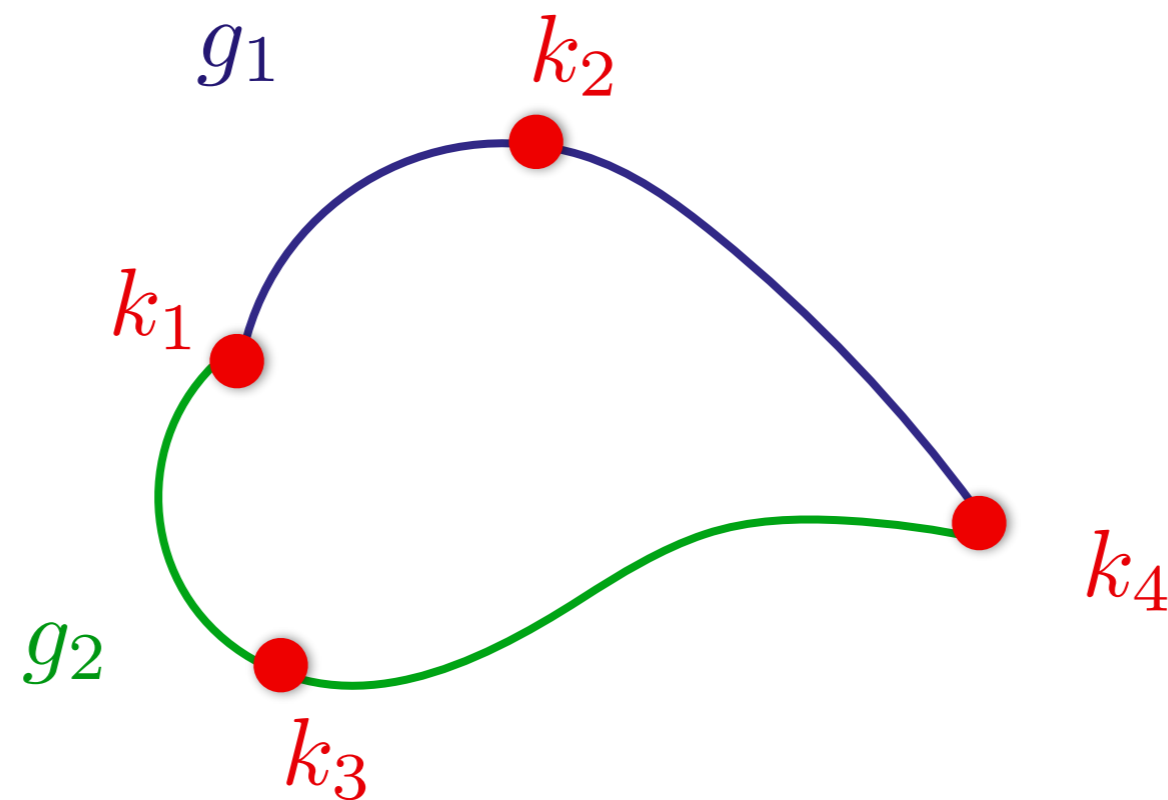
*Simplified training for gesture recognition*

# Ideal key pose set

✓ Concise (small)

✓ Discriminative (avoid ambiguity)



$g_1$

$g_2$

*Simplified training for gesture recognition*

# Ideal key pose set
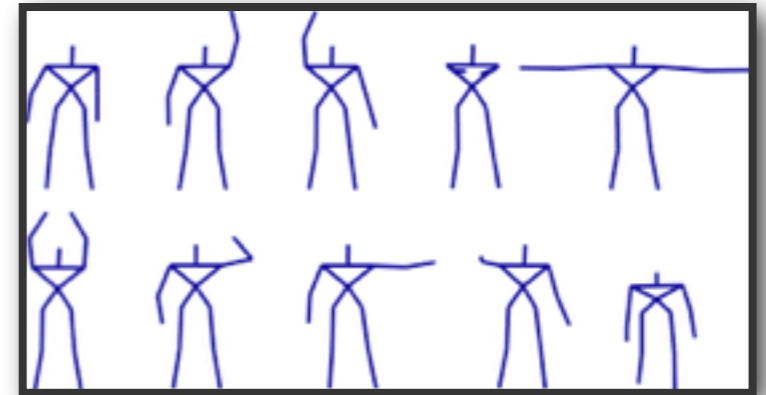
✓ Concise (small)

✓ Discriminative (avoid ambiguity)



$g_1$   $k_2$

$k_1$

$g_2$   $k_4$

$k_3$

*Simplified training for gesture recognition*

# Ideal key pose set

✓ Concise (small)

✓ Discriminative (avoid ambiguity)



$g_1$   $k_2$

$k_1$

$k_4$

$g_2$

$k_3$

Our solution: adaptive sampling
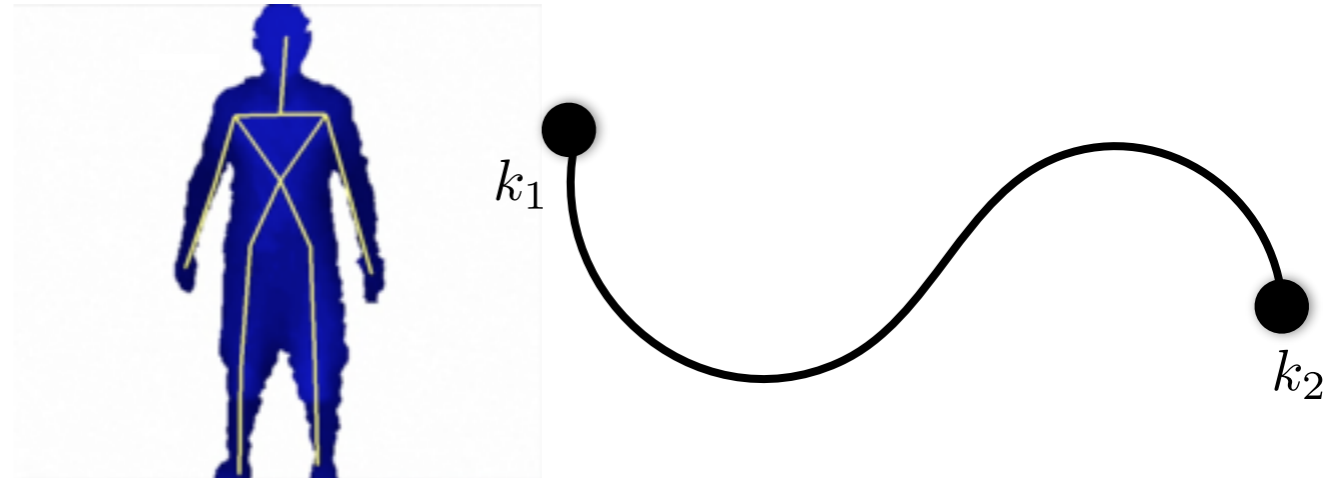
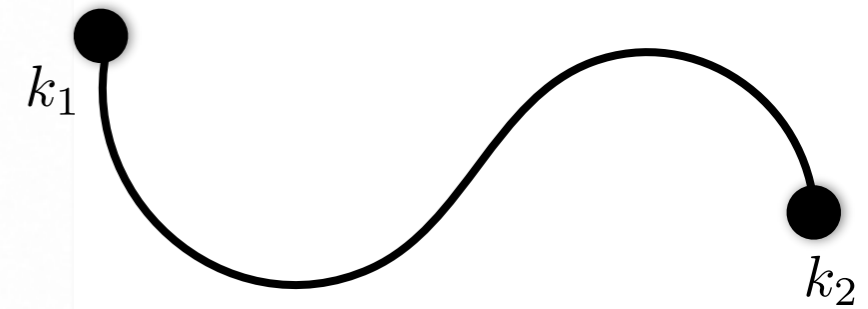*Simplified training for gesture recognition*

# Building a good key pose set



1 - initial / final gesture poses
$k_1$ and $k_2$ must be key poses

# Building a good key pose set

1 - initial / final gesture poses $k_1$ and $k_2$ must be key poses
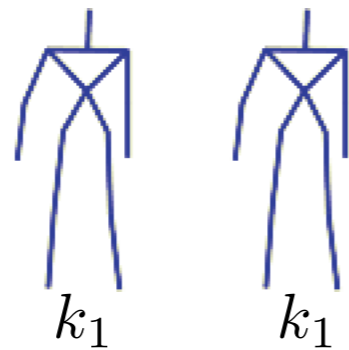


$k_1$

$k_2$

# Building a good key pose set

1 - initial / final gesture poses
$k_1$ and $k_2$ must be key poses

what if initial == final?



$$\hat{g} = (k_1, k_1)$$

*Simplified training for gesture recognition*

# Building a good key pose set

1 - initial / final gesture poses
$k_1$ and $k_2$ must be key poses

what if initial == final?
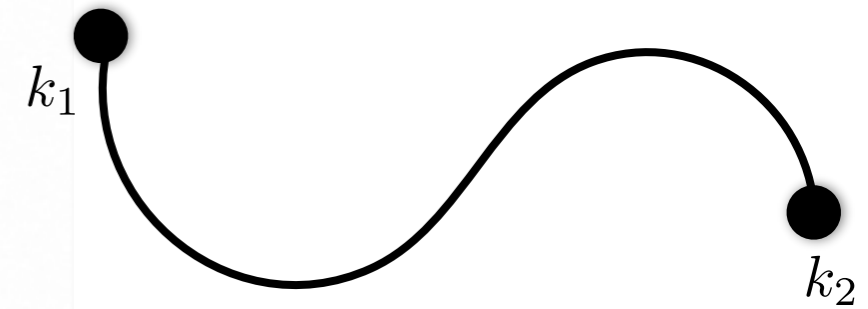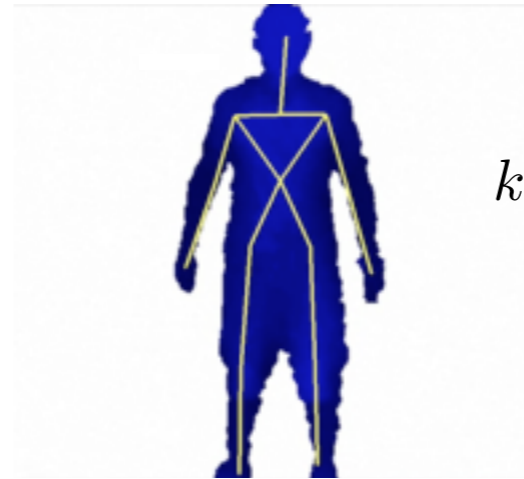


$k_1$    $k_1$
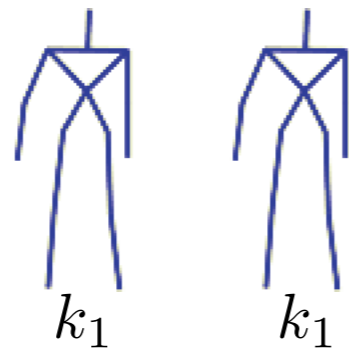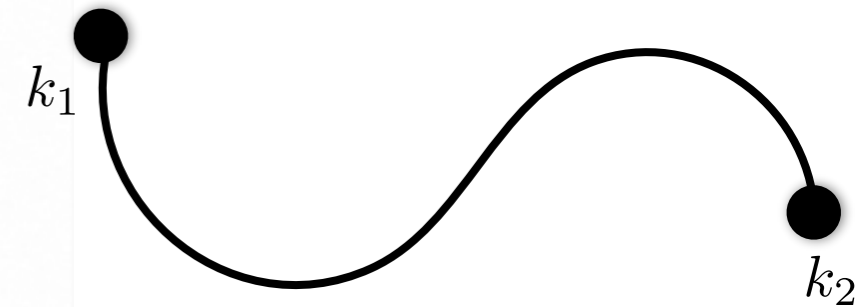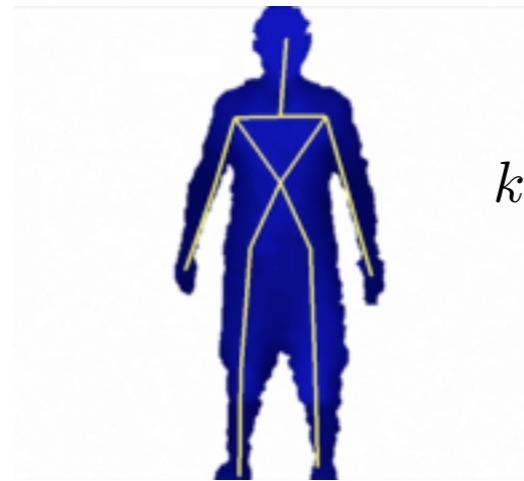
$$\hat{g} = (k_1, k_1)$$

# Building a good key pose set

1 - initial / final gesture poses
$k_1$ and $k_2$ must be key poses

what if initial == final?

$$\hat{g} = (k_1, k_2, k_1)$$

insert intermediate farthest pose:

$$k_2 = \operatorname*{argmax}_{p \in g} \Delta(p, k_1)$$

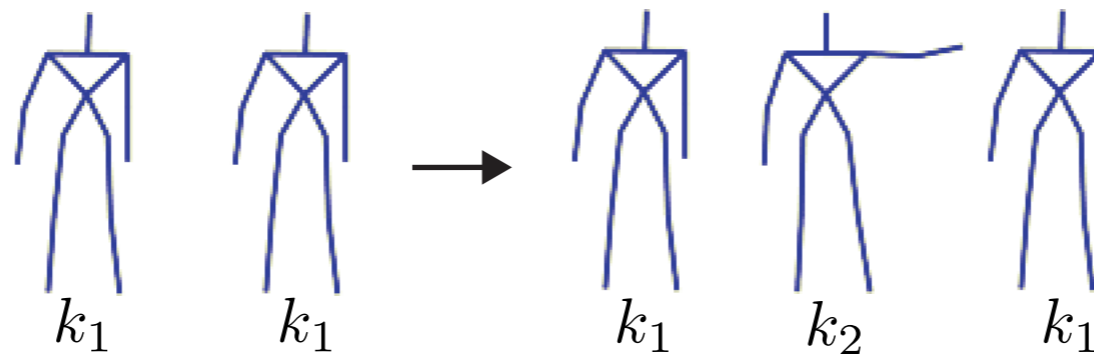*Simplified training for gesture recognition*
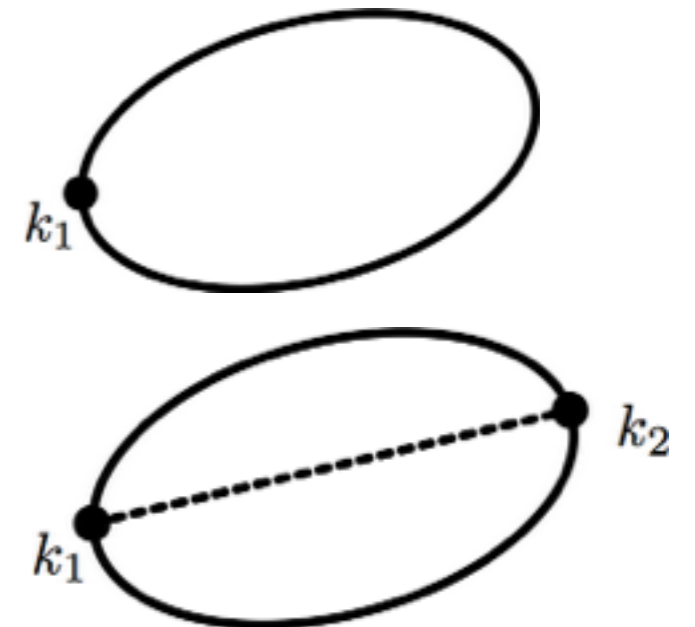
# Building a good key pose set

1 - initial / final gesture poses
$k_1$ and $k_2$ must be key poses

what if initial == final?
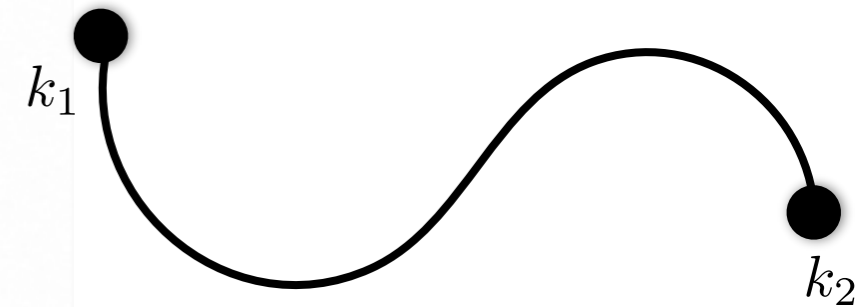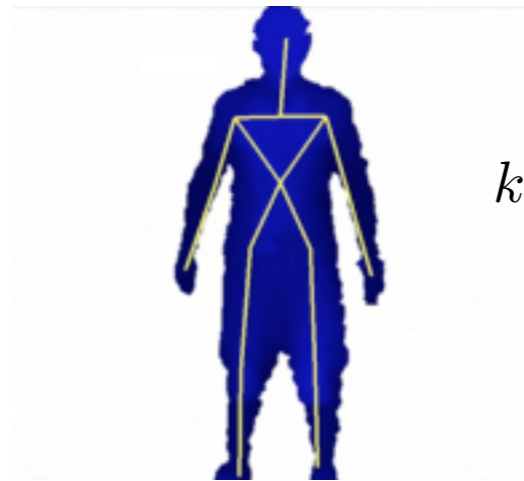
$$\hat{g} = (k_1, k_2, k_1)$$

insert intermediate farthest pose:

$$k_2 = \underset{p \in g}{\mathrm{argmax}}\; \Delta\left(p, k_1\right)$$

if $k_1$ == $k_2$ , discard. (static gesture)

Faugeroux et al., 2014

# Building a good key pose set: Discriminant poses

2 - similar representations for distinct gestures $g$ and $g'$ must be refined.

*Simplified training for gesture recognition*

# Building a good key pose set: Discriminant poses

2 - similar representations for distinct gestures $g$ and $g'$ must be refined.

# Building a good key pose set: Discriminant poses

2 - similar representations for distinct gestures $g$ and $g'$ must be refined.



Insert most discriminant pose in each gesture:

$$k_{1+1/2} = \operatorname*{argmax}_{p \in g} \min_{p' \in g'} \Delta\left(p, p'\right)$$

$$k'_{1+1/2} = \operatorname*{argmax}_{p' \in g'} \min_{p \in g} \Delta\left(p', p\right)$$



Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Discriminant poses: general case

Repeat the process above for every sub-sequence between successive key poses of $g$ and $g'$

Select the most distinctive pair

$$j = \operatorname*{argmax}_{i} \left\{ \min_{p' \in g'_i} \Delta\left(k_{i+\nicefrac{1}{2}}, p'\right) + \min_{p \in g_i} \Delta\left(k'_{i+\nicefrac{1}{2}}, p\right) \right\}$$



If $\Delta\left(k_{j+\nicefrac{1}{2}}, k'_{j+\nicefrac{1}{2}}\right) < \epsilon$ , give gestures identical label.

Iterate until all gestures have different representations

# Discriminant poses: general case

Repeat the process above for every sub-sequence between successive key poses of $g$ and $g'$

Select the most distinctive pair

$$j = \operatorname*{argmax}_{i} \left\{ \min_{p' \in g'_i} \Delta\left(k_{i+1/2}, p'\right) + \min_{p \in g_i} \Delta\left(k'_{i+1/2}, p\right) \right\}$$



If $\Delta\left(k_{j+1/2}, k'_{j+1/2}\right) < \epsilon$ , give gestures identical label.

Iterate until all gestures have different representations

# Outline

# Spurious gesture elimination

Transitions between gestures

Static gestures

*Simplified training for gesture recognition*

# Spurious gesture elimination

Transitions between gestures



Static gestures

*Simplified training for gesture recognition*

# Spurious gesture elimination

Transitions between gestures

Static gestures

*Simplified training for gesture recognition*

# Spurious gesture elimination

Transitions between gestures

Static gestures

*Simplified training for gesture recognition*

# Spurious gesture elimination

## Transitions between gestures



## Static gestures

*Simplified training for gesture recognition*

# Semiautomatic labeling

Similar gestures confirmation



$$\hat{g}_1 = (k_1, k_2, k_3)$$

$$\hat{g}_2 = (k_1, k_2, k_3)$$

"Is $g_1$ and $g_2$ performances of the same gesture?"

*Simplified training for gesture recognition*

# Semiautomatic labeling

Similar gestures confirmation



$$\hat{g}_1 = (k_1, k_2, k_3)$$
$$\hat{g}_2 = (k_1, k_2, k_3)$$

"Is $g_1$ and $g_2$ performances of the same gesture?"

Negative: force key-pose subdivision (ignore $\epsilon$)



$$\hat{g}_1 = (k_1, k_4, k_2, k_3)$$
$$\hat{g}_2 = (k_1, k_5, k_2, k_3)$$

*Simplified training for gesture recognition*

# Gesture recognition



key pose set
$$\mathcal{K} = \{k_1, \ldots, k_n\}$$

gesture set
$$\mathcal{G} = \{\hat{g}_1, \hat{g}_2, \ldots, \hat{g_m}\}$$

training set

*Simplified training for gesture recognition*

# Learning method?

Many alternatives: action graph, decision forests, bag of features, SVM, nearest neighbor classifier,…

# Learning method?

Many alternatives: action graph, decision forests, bag of features, SVM, nearest neighbor classifier,…

(Miranda *et al,* 2012): SVM + decision forest

<span style="color:red">nearest neighbor classifier</span>

$$\tilde{f}(p) = \begin{cases} k_p = \underset{k \in \mathcal{K}}{\operatorname{argmin}} \, \Delta\left(k, p\right) & \text{if } \Delta\left(k_p, p\right) < \epsilon, \\ -1 & \text{otherwise.} \end{cases}$$

*Simplified training for gesture recognition*

# Experiments

*Simplified training for gesture recognition*

# Experiment Setup

✓ Objective: comparison with Miranda *et al* (2012)

✓ Same set of 11 gestures

✓ Gestures briefly described to 10 inexperienced individuals

✓ Users should sequentially perform each gesture in a single record

✓ Unsuccessfully segmented gestures exceptionally retrained

*Simplified training for gesture recognition*

# Experiment Setup

✓ Objective: comparison with Miranda *et al* (2012)

✓ Same set of 11 gestures

✓ Gestures briefly described to 10 inexperienced individuals

✓ Users should sequentially perform each gesture in a single record

✓ Unsuccessfully segmented gestures exceptionally retrained



Faugeroux et al., 2014                                                                 *Simplified training for gesture recognition*

# Segmentation robustness

From 10 users recordings:

| gesture | id | segmentation accuracy |
|---|---|---|
| Turn Next Page | $\hat{g}_A$ | 10 |
| Turn Previous Page | $\hat{g}_B$ | 10 |
| Raise Right Arm | $\hat{g}_C$ | 10 |
| Raise Left Arm | $\hat{g}_D$ | 10 |
| Open Clap | $\hat{g}_E$ | 8 |
| Open Arms | $\hat{g}_F$ | 9 |
| Put Hands Up Lat. | $\hat{g}_G$ | 9 |
| Put Hands Up Front | $\hat{g}_H$ | 10 |
| Lower Right Arm | $\hat{g}_I$ | 8 |
| Bow | $\hat{g}_J$ | 6 |
| Goodbye | $\hat{g}_K$ | 7 |
| **average (%)** | | **88** |

over-segmentation:

*Simplified training for gesture recognition*

# Segmentation robustness

**Simplified training for gesture recognition**

*example: online segmentation and key pose selection*

# Segmentation robustness

**Simplified training for gesture recognition**

*example: online segmentation and key pose selection*

# Discriminant key pose selection

10 to 12 key poses per set



Key poses similar to manually designed key poses from Miranda *et al* (2012)

Miranda *et al* (2012) uses 11 key poses!

Bigger $\epsilon$ : less key poses, less accurate executions needed

Smaller $\epsilon$ : more key poses, more accurate executions needed

# Discriminant key pose selection

*Simplified training for gesture recognition*

# Gesture recognition

Each user executed each gesture 10 times

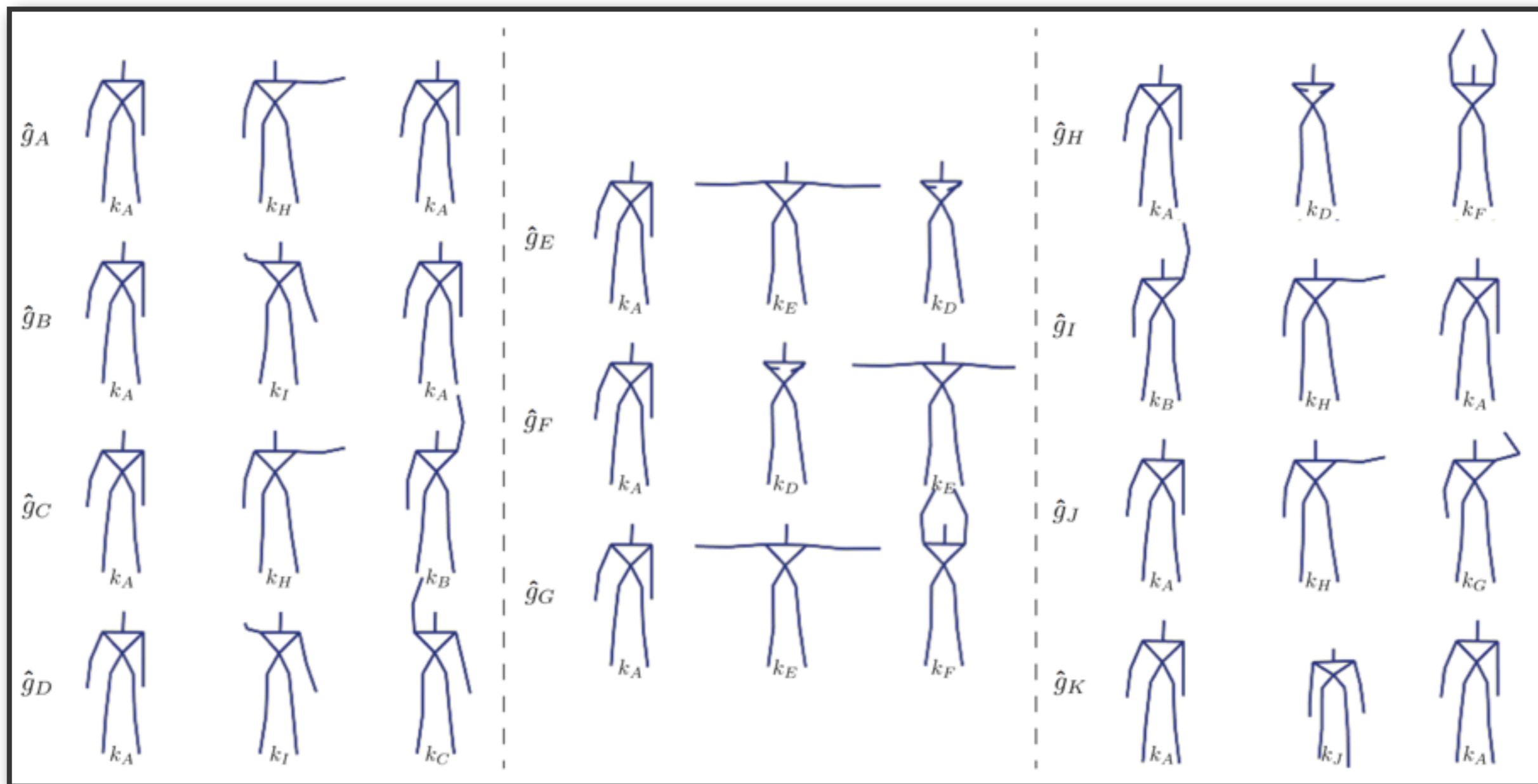| gesture | id | segmentation accuracy | recognized gestures per user | | | | | | | | | | ours (%) | [14] (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | $u_1$ | $u_2$ | $u_3$ | $u_4$ | $u_5$ | $u_6$ | $u_7$ | $u_8$ | $u_9$ | $u_{10}$ | | |
| Turn Next Page | $\hat{g}_A$ | 10 | 10 | 8 | 10 | 10 | 9 | 9 | 10 | 9 | 9 | 9 | 93 | 95 |
| Turn Previous Page | $\hat{g}_B$ | 10 | 10 | 9 | 10 | 10 | 9 | 6 | 9 | 9 | 9 | 10 | 91 | 95 |
| Raise Right Arm | $\hat{g}_C$ | 10 | 9 | 8 | 9 | 8 | 7 | 10 | 9 | 10 | 8 | 10 | 88 | 94 |
| Raise Left Arm | $\hat{g}_D$ | 10 | 10 | 10 | 9 | 10 | 9 | 9 | 10 | 9 | 10 | 9 | 95 | 94 |
| Open Clap | $\hat{g}_E$ | 8 | 10 | 10 | 10 | 9 | 9 | 8 | 10 | 9 | 8 | 10 | 93 | 99 |
| Open Arms | $\hat{g}_F$ | 9 | 9 | 9 | 10 | 8 | 9 | 10 | 10 | 9 | 8 | 9 | 91 | 97 |
| Put Hands Up Lat. | $\hat{g}_G$ | 9 | 10 | 10 | 10 | 10 | 10 | 9 | 10 | 10 | 10 | 10 | 99 | 100 |
| Put Hands Up Front | $\hat{g}_H$ | 10 | 10 | 9 | 7 | 9 | 10 | 9 | 10 | 10 | 10 | 9 | 93 | 96 |
| Lower Right Arm | $\hat{g}_I$ | 8 | 8 | 7 | 6 | 8 | 7 | 8 | 8 | 8 | 8 | 7 | 75 | 82 |
| Bow | $\hat{g}_J$ | 6 | 10 | 10 | 10 | 9 | 10 | 10 | 10 | 9 | 10 | 10 | 98 | 100 |
| Goodbye | $\hat{g}_K$ | 7 | 9 | 9 | 10 | 7 | 9 | 10 | 9 | 10 | 8 | 7 | 88 | 92 |
| average (%) | | 88 | 90 | 92 | 89 | 89 | 91 | 89 | 93 | 89 | 91 | 92 | | |



$k_A$  $k_B$  $k_C$  $k_D$  $k_E$  $k_F$  $k_G$  $k_H$  $k_I$  $k_J$
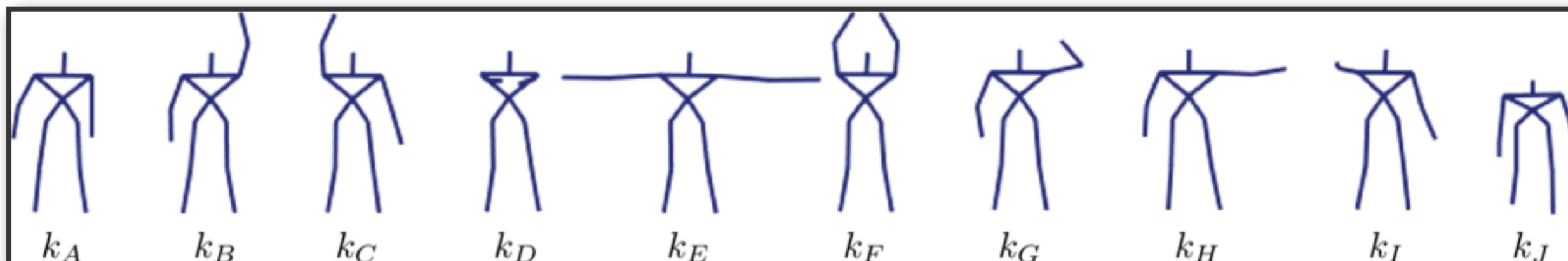
Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Gesture recognition

Each user executed each gesture 10 times

| gesture | id | segmentation accuracy | recognized gestures per user | | | | | | | | | | ours (%) | [14] (%) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | $u_1$ | $u_2$ | $u_3$ | $u_4$ | $u_5$ | $u_6$ | $u_7$ | $u_8$ | $u_9$ | $u_{10}$ | | |
| Turn Next Page | $\hat{g}_A$ | 10 | 10 | 8 | 10 | 10 | 9 | 9 | 10 | 9 | 9 | 9 | 93 | 95 |
| Turn Previous Page | $\hat{g}_B$ | 10 | 10 | 9 | 10 | 10 | 9 | 6 | 9 | 9 | 9 | 10 | 91 | 95 |
| Raise Right Arm | $\hat{g}_C$ | 10 | 9 | 8 | 9 | 8 | 7 | 10 | 9 | 10 | 8 | 10 | 88 | 94 |
| Raise Left Arm | $\hat{g}_D$ | 10 | 10 | 10 | 9 | 10 | 9 | 9 | 10 | 9 | 10 | 9 | 95 | 94 |
| Open Clap | $\hat{g}_E$ | 8 | 10 | 10 | 10 | 9 | 9 | 8 | 10 | 9 | 8 | 10 | 93 | 99 |
| Open Arms | $\hat{g}_F$ | 9 | 9 | 9 | 10 | 8 | 9 | 10 | 10 | 9 | 8 | 9 | 91 | 97 |
| Put Hands Up Lat. | $\hat{g}_G$ | 9 | 10 | 10 | 10 | 10 | 10 | 9 | 10 | 10 | 10 | 10 | 99 | 100 |
| Put Hands Up Front | $\hat{g}_H$ | 10 | 10 | 9 | 7 | 9 | 10 | 9 | 10 | 10 | 10 | 9 | 93 | 96 |
| Lower Right Arm | $\hat{g}_I$ | 8 | 8 | 7 | 6 | 8 | 7 | 8 | 8 | 8 | 8 | 7 | 75 | 82 |
| Bow | $\hat{g}_J$ | 6 | 10 | 10 | 10 | 9 | 10 | 10 | 10 | 9 | 10 | 10 | 98 | 100 |
| Goodbye | $\hat{g}_K$ | 7 | 9 | 9 | 10 | 7 | 9 | 10 | 9 | 10 | 8 | 7 | 88 | 92 |
| average (%) | | 88 | 90 | 92 | 89 | 89 | 91 | 89 | 93 | 89 | 91 | 92 | | |



Faugeroux et al., 2014

*Simplified training for gesture recognition*

# Performance

Real time training:

Segmentation + computing key pose representations

Real time gesture recognition.

Offline experiment:

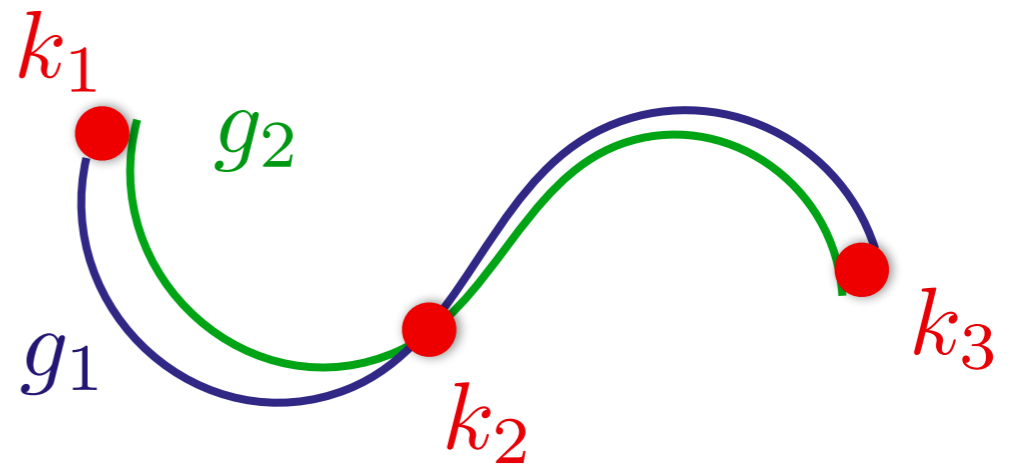Segmentation + computing key pose representations:

1239 frames (44.3 secs)

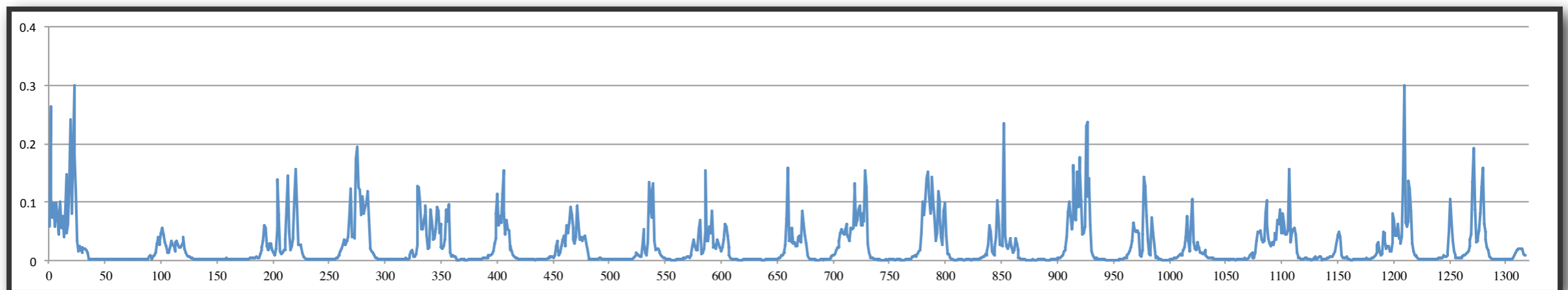20 gestures

Total time: 0.33 secs

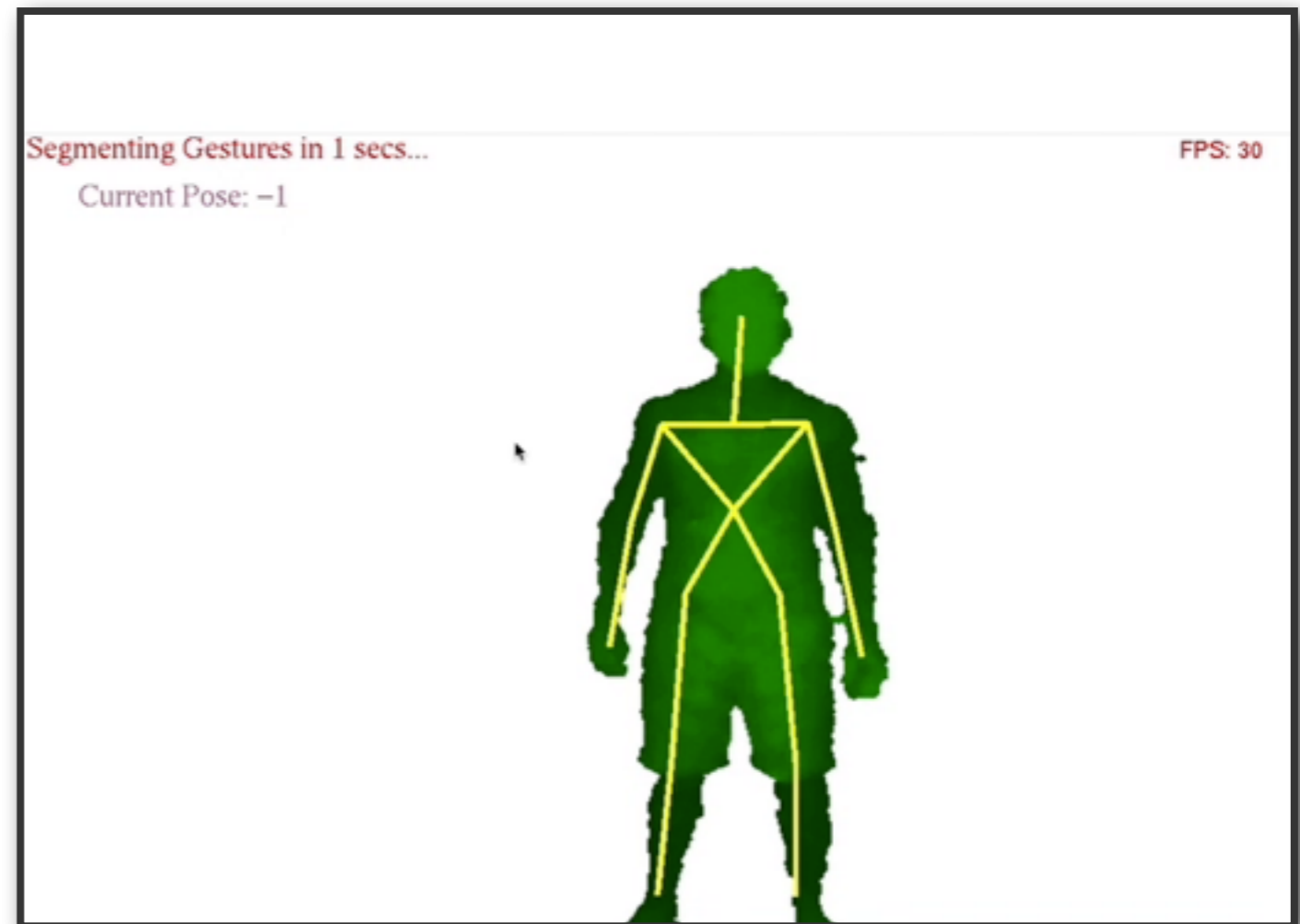# Limitations & Future work

✓ Execution speed not considered



✓ Fixed threshold $\epsilon$

✓ Only the first curvature is used for segmentation

*Simplified training for gesture recognition*

# Thank you for your attention!



## Questions?

*Simplified training for gesture recognition*
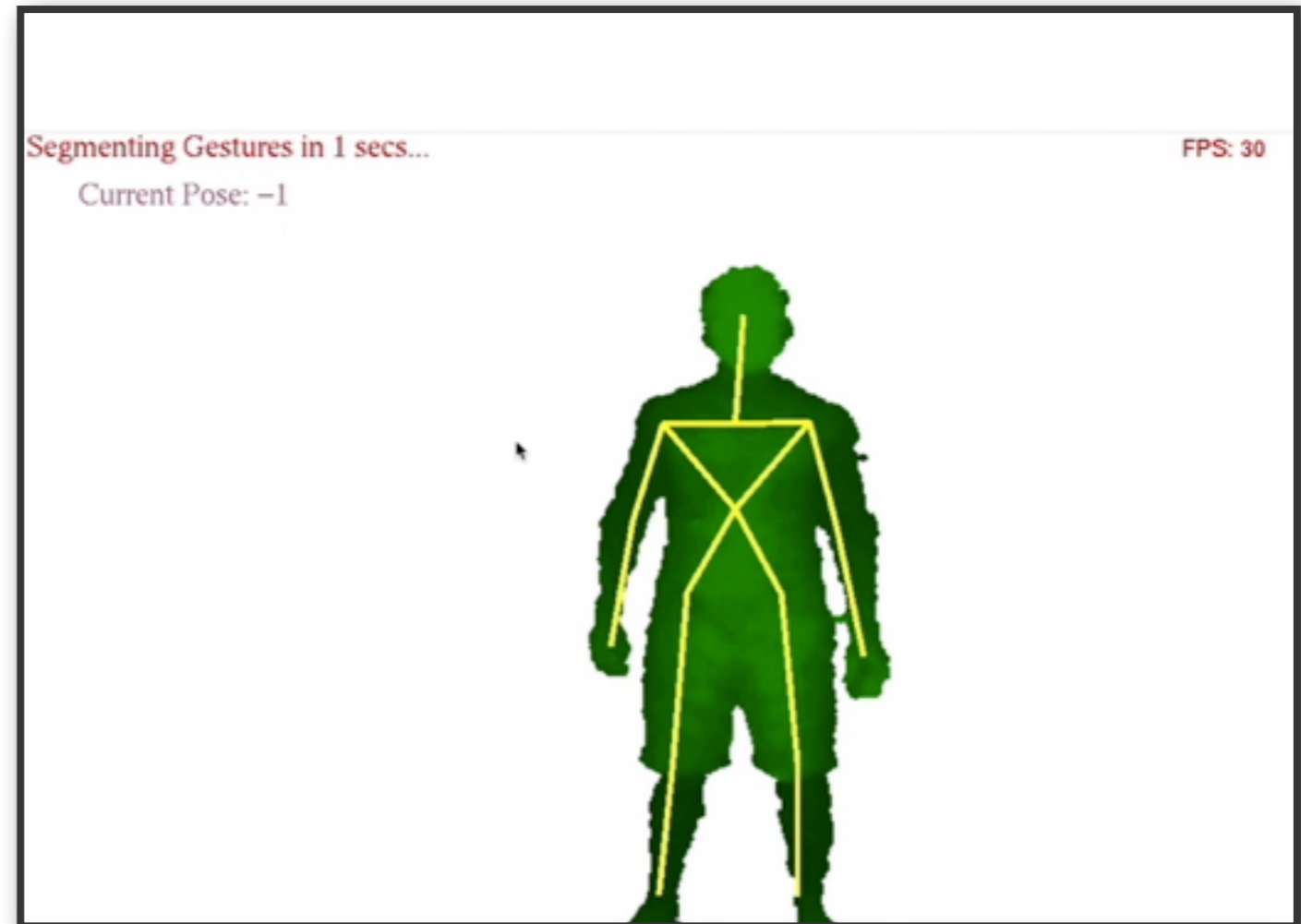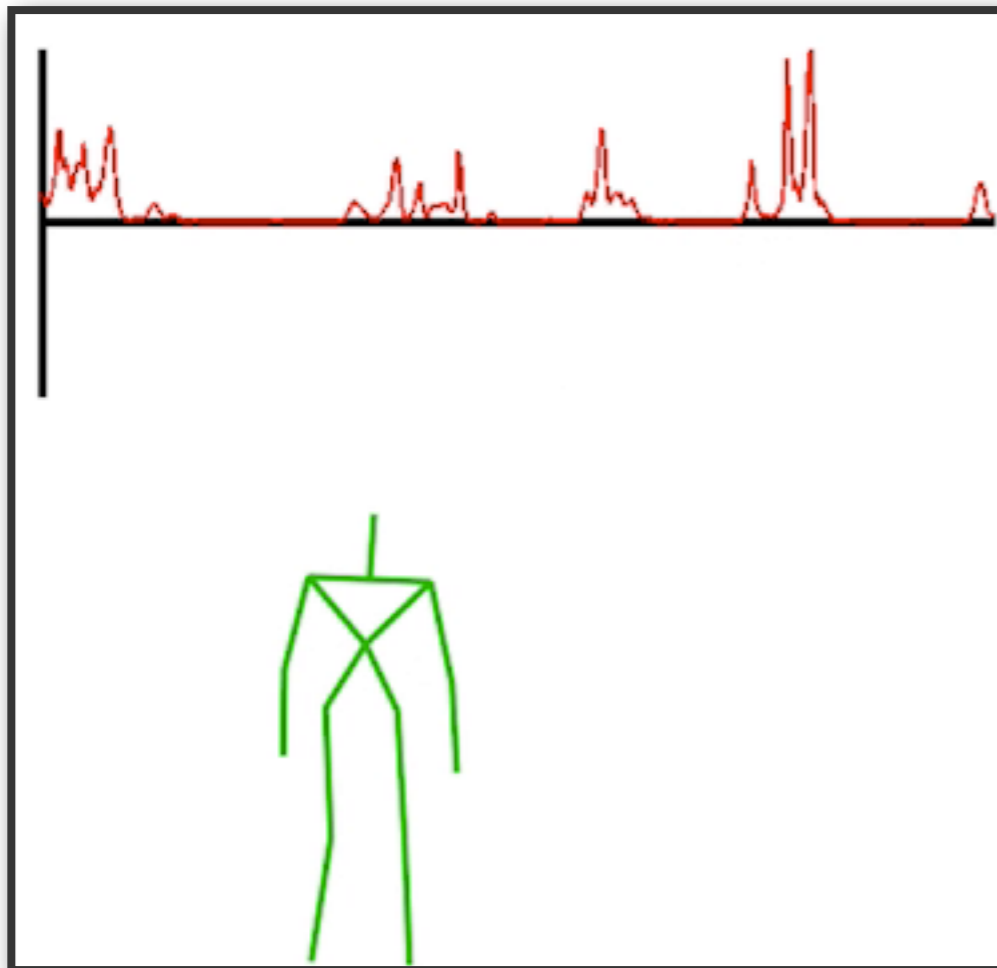
# Thank you for your attention!





# Questions?