

# Real-time gesture recognition from depth data through key poses learning and decision forests

Leandro Miranda

Thales Vieira (presenter)

Dimas Martinez

*Mathematics, UFAL*

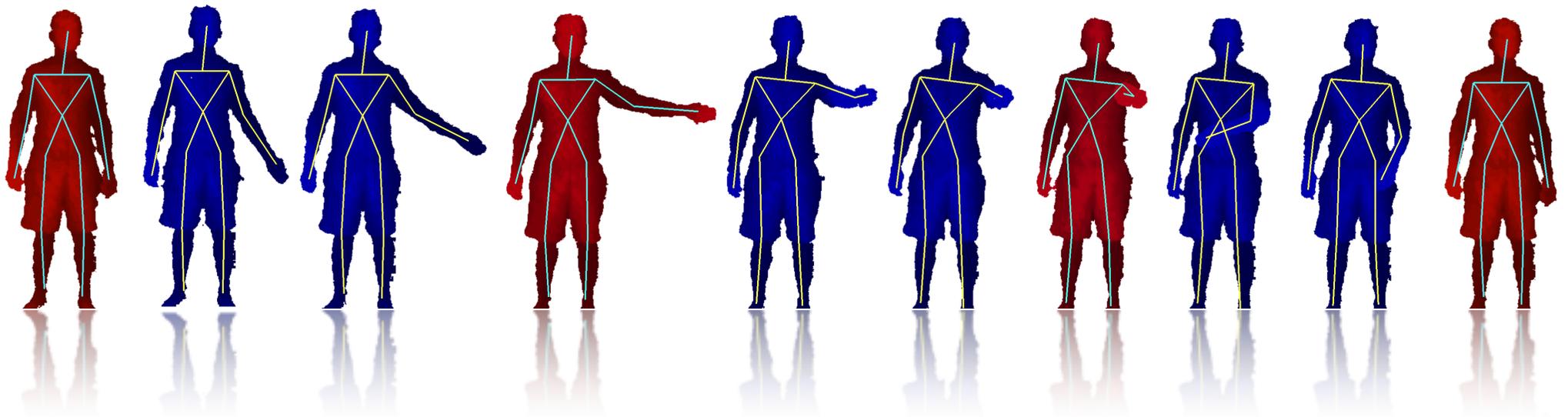
Thomas Lewiner

*Mathematics, PUC-Rio*

Antonio W. Vieira

Mario F. M. Campos

*Computer Science, UFMG*



# Human Gesture Recognition

# Human Gesture Recognition



# Human Gesture Recognition



# Human Gesture Recognition

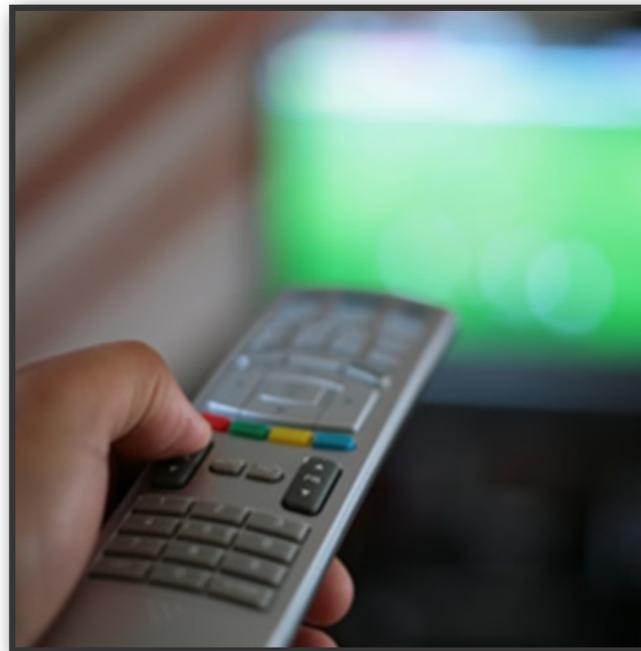


# Current Scenario

Popularization of real time depth sensors



*Microsoft Kinect Sensor*



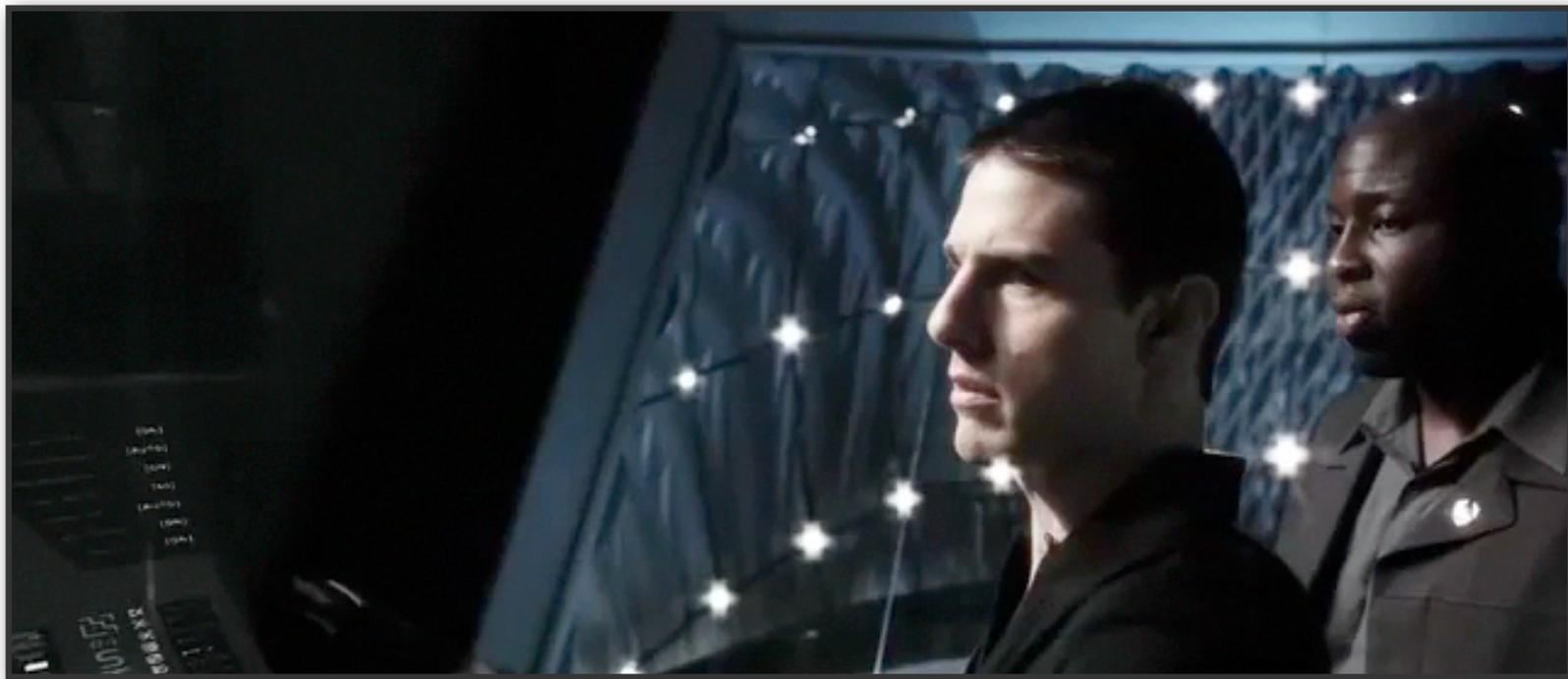
Development of high quality Natural User Interfaces (NUI)

# Current Scenario

Popularization of real time depth sensors

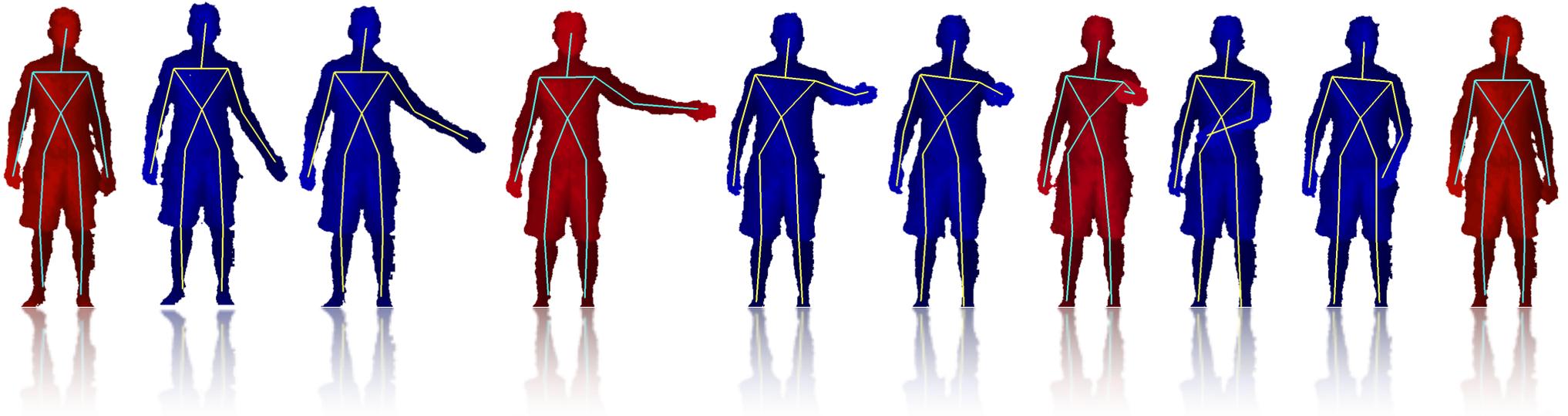


*Microsoft Kinect Sensor*



**Challenging task! Gestures performed at different speeds and/or sequence of poses**

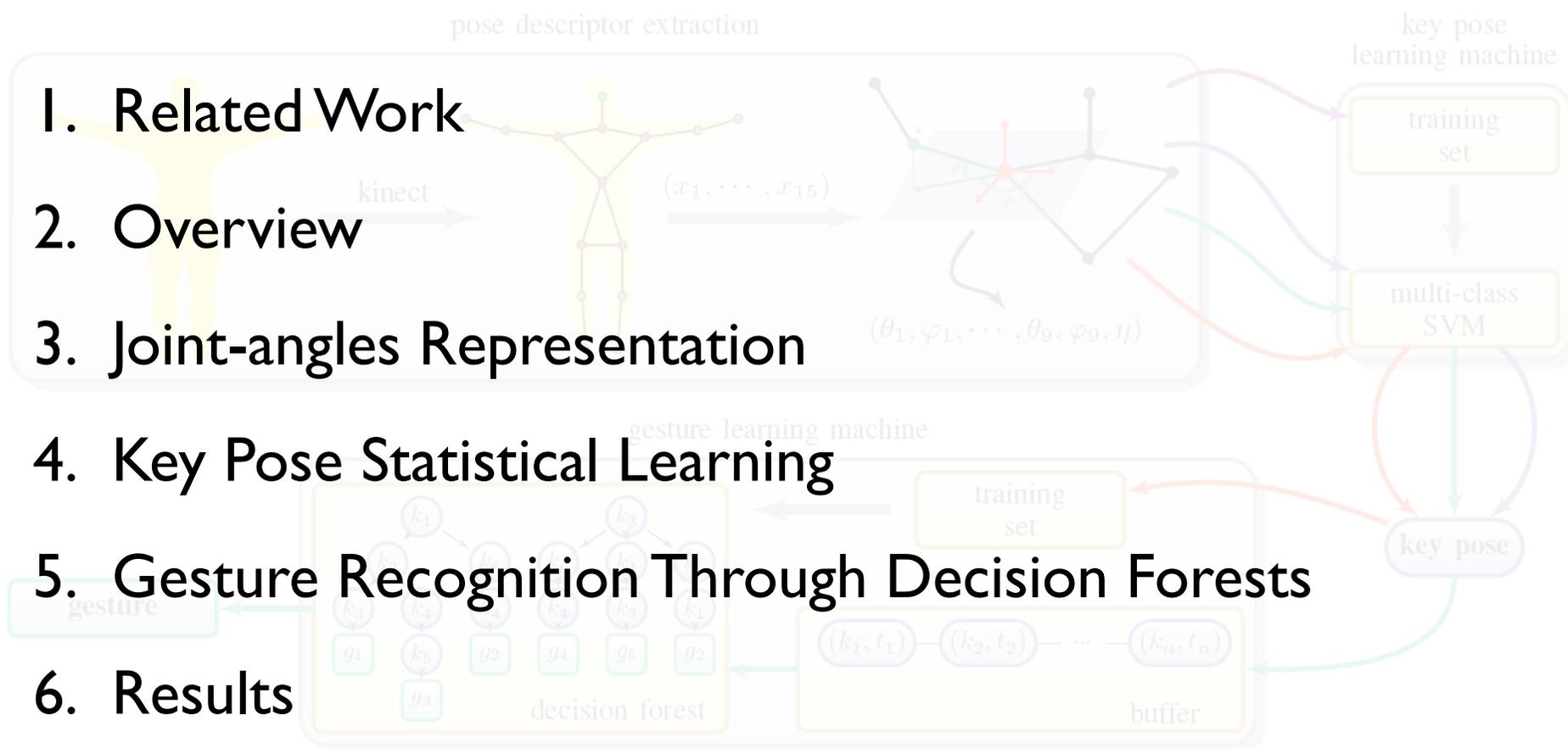
# Our approach: key poses learning



Gestures can be characterized by a few extreme poses!

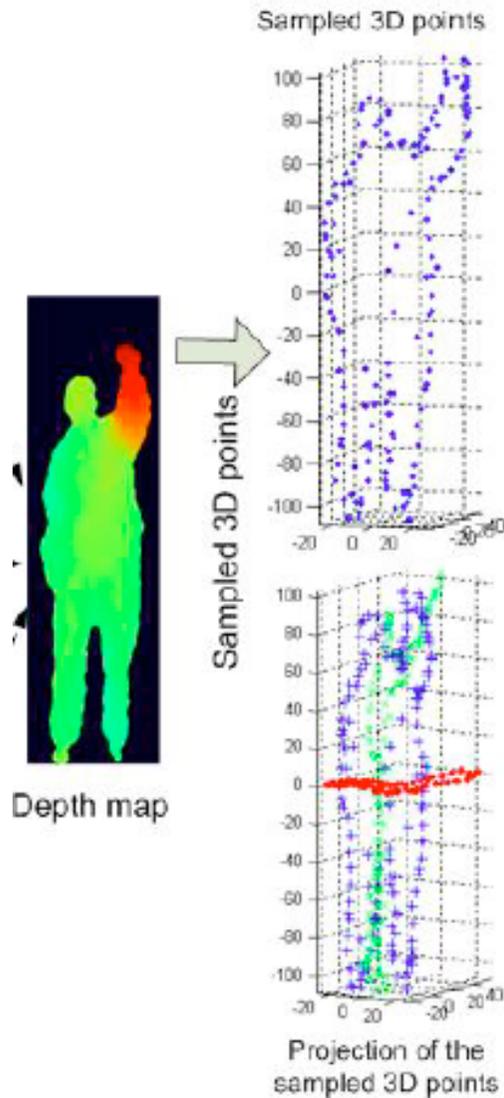
- ✓ Real-time gesture learning and recognition
- ✓ Ideal for the average inexperienced user

# Outline



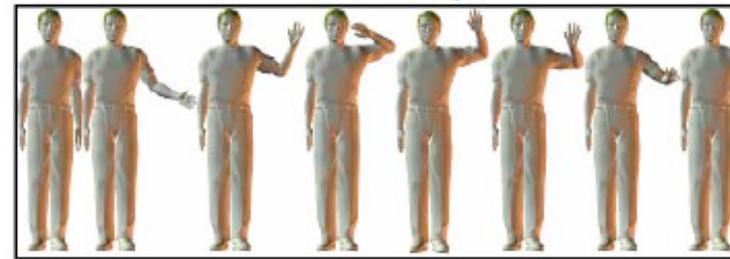
# Related Work

## Local methods



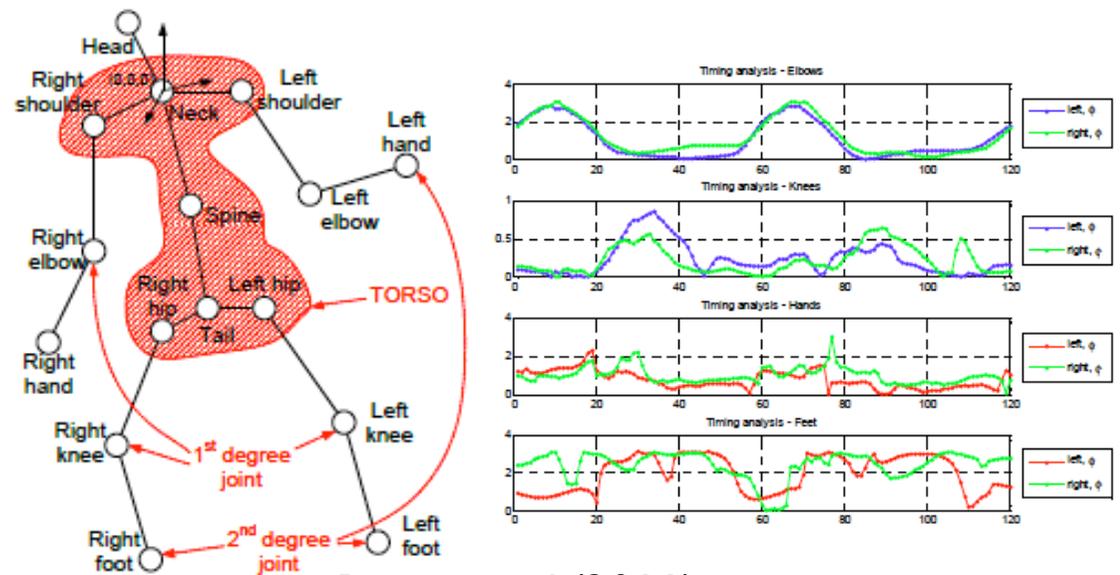
Li et al (2010)

## Global methods



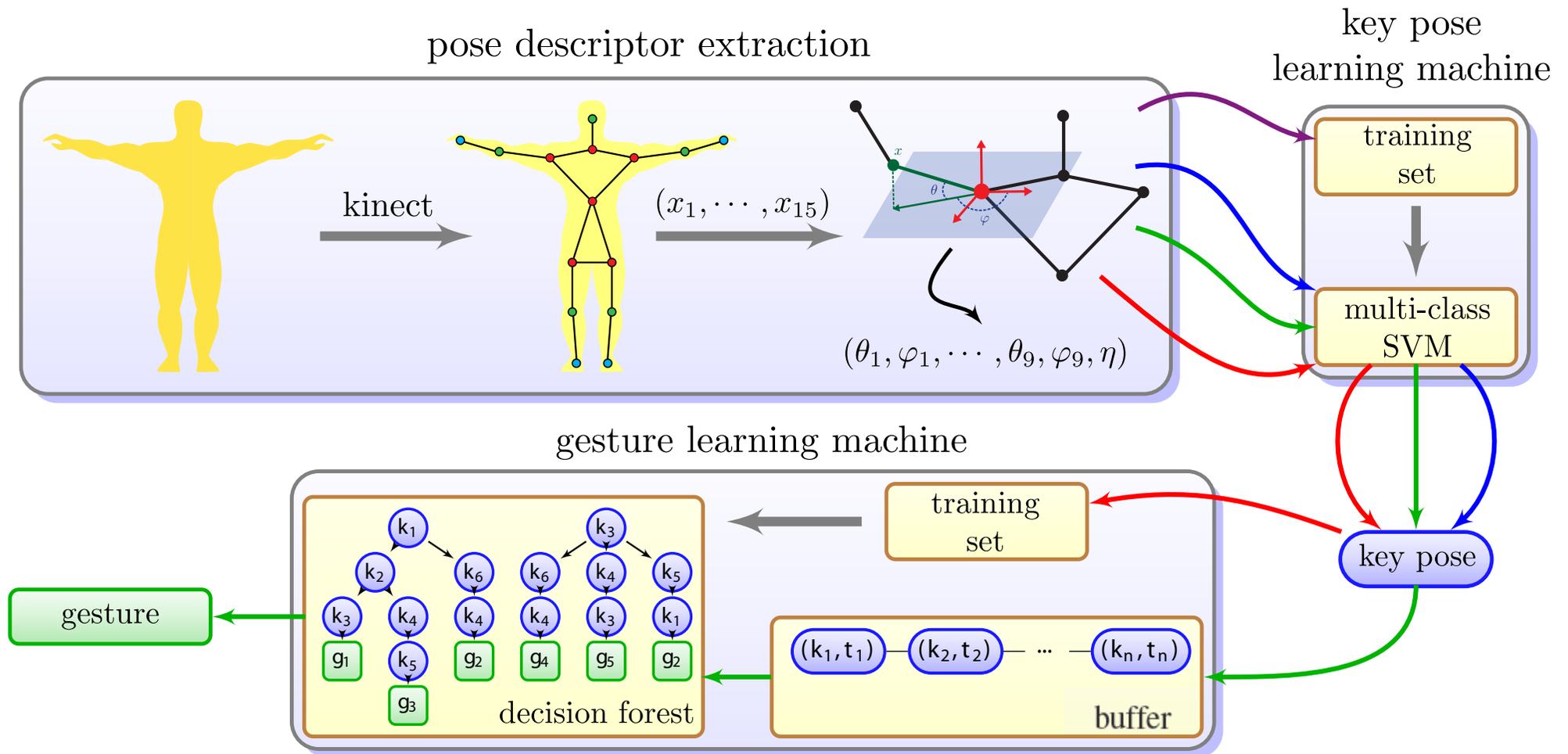
Lv and Nevatia (2007)

## Parametric methods

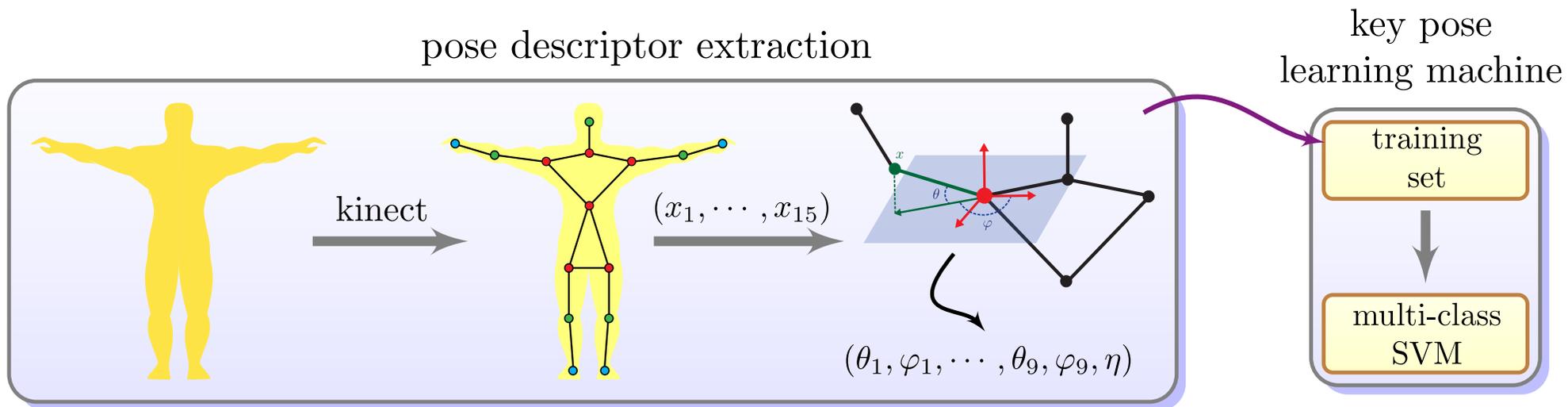


Raptis et al (2011)

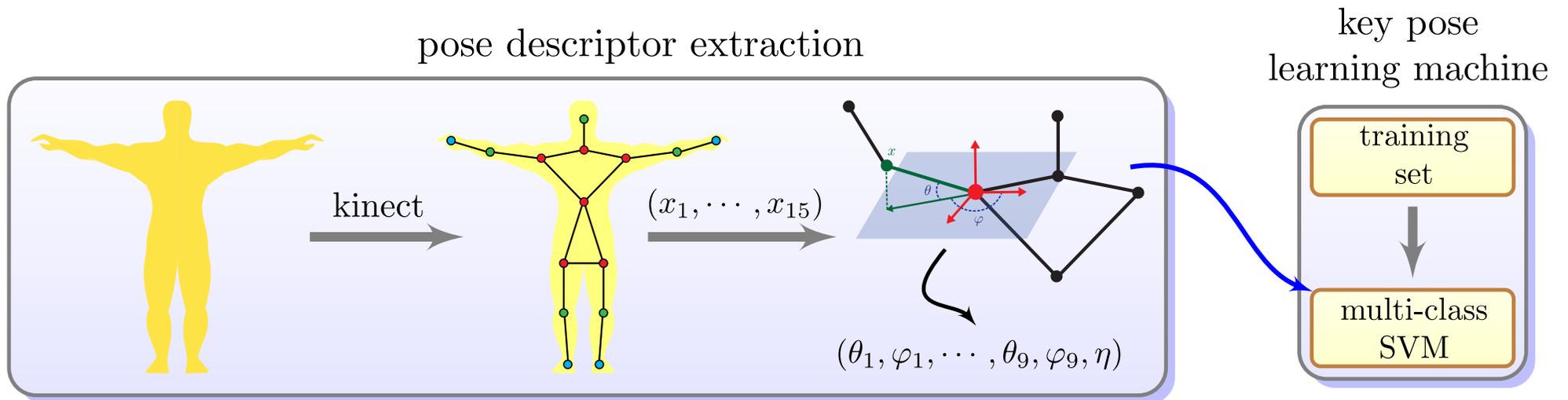
# Overview



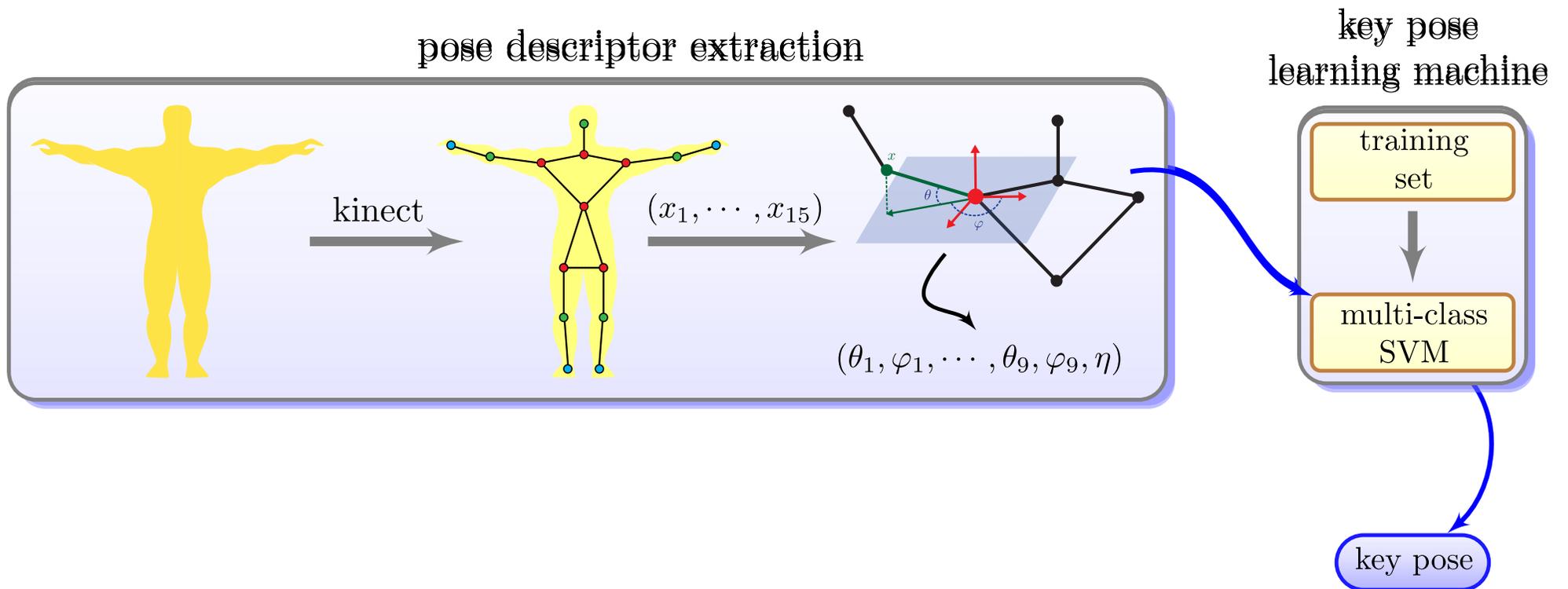
# Overview: training key poses



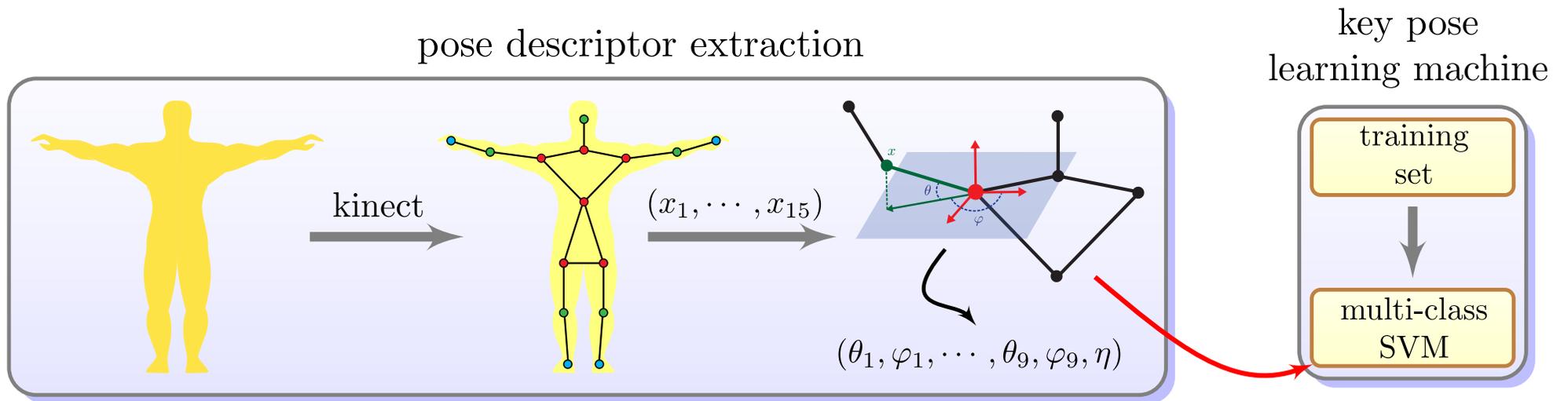
# Overview: recognizing key poses



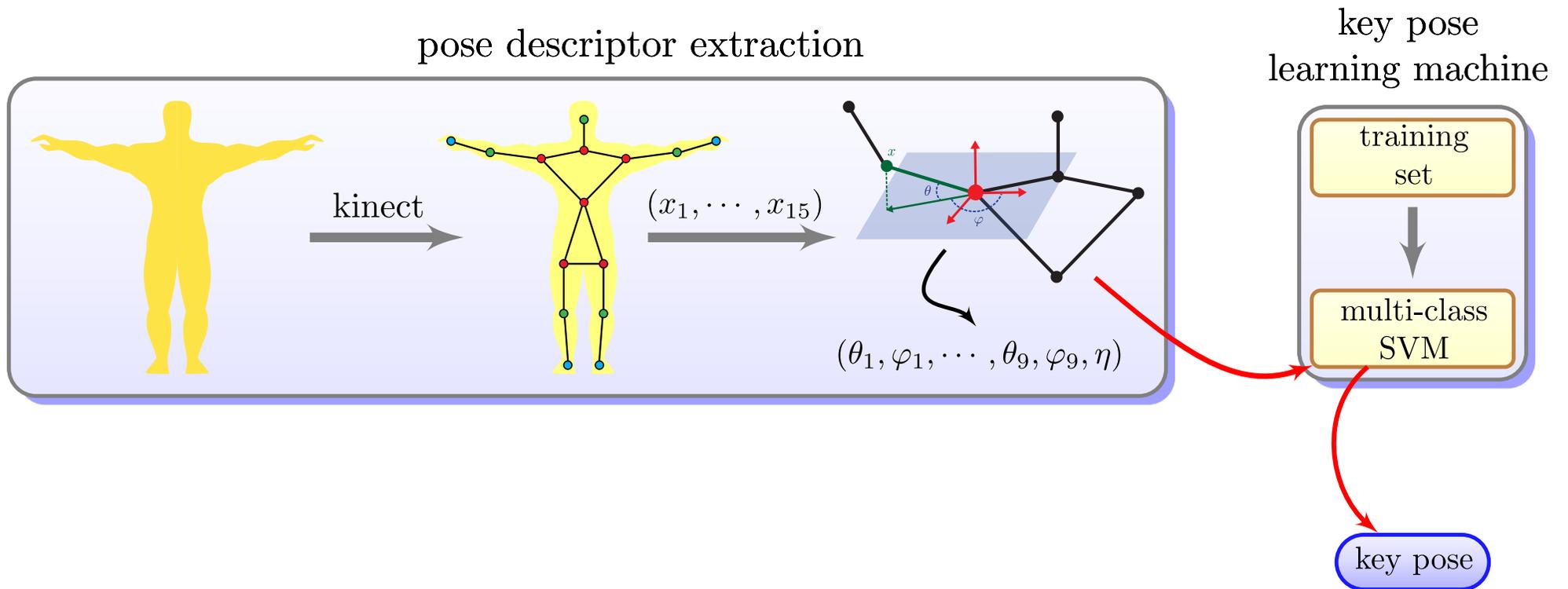
# Overview: recognizing key poses



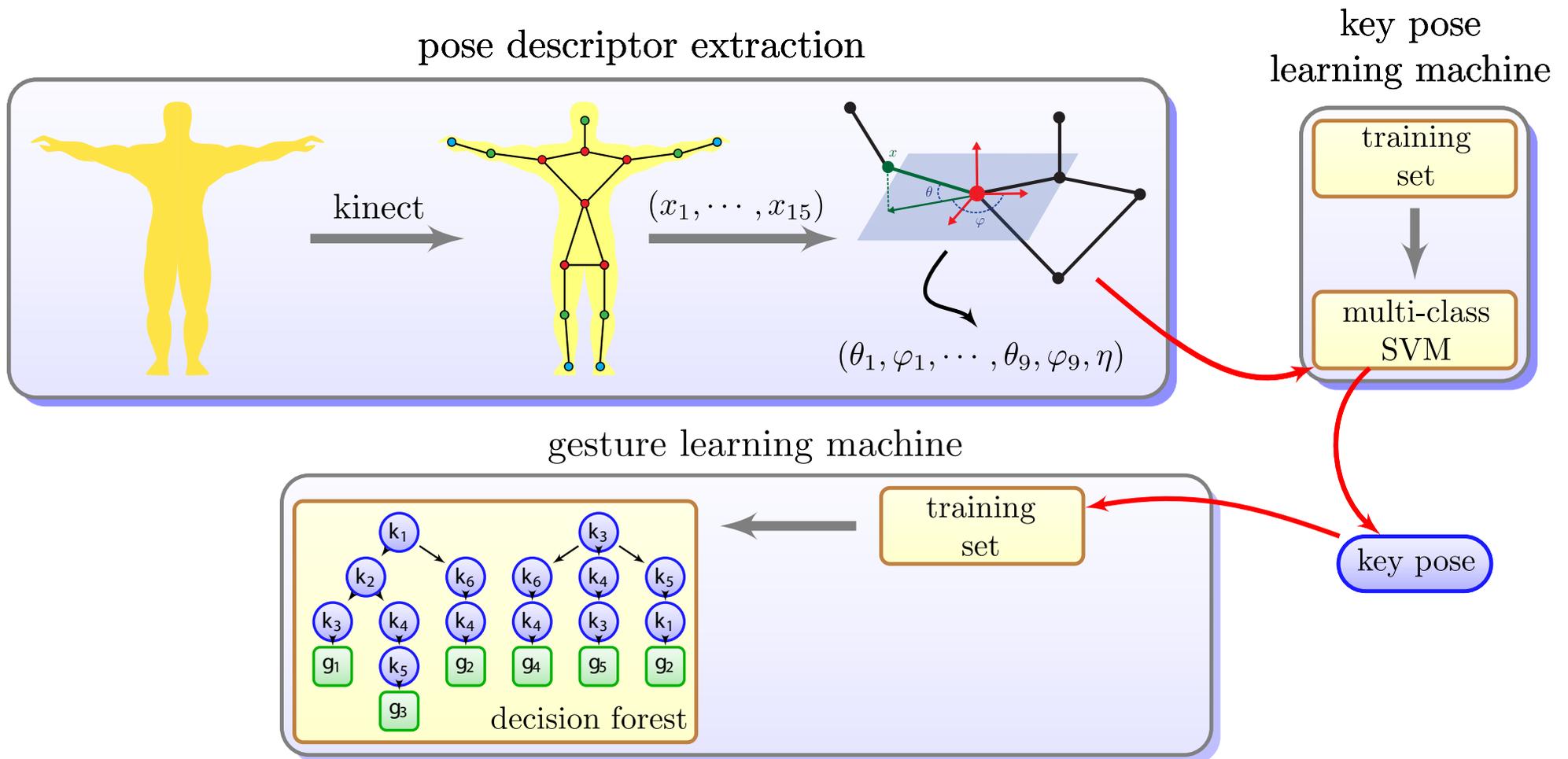
# Overview: training gestures



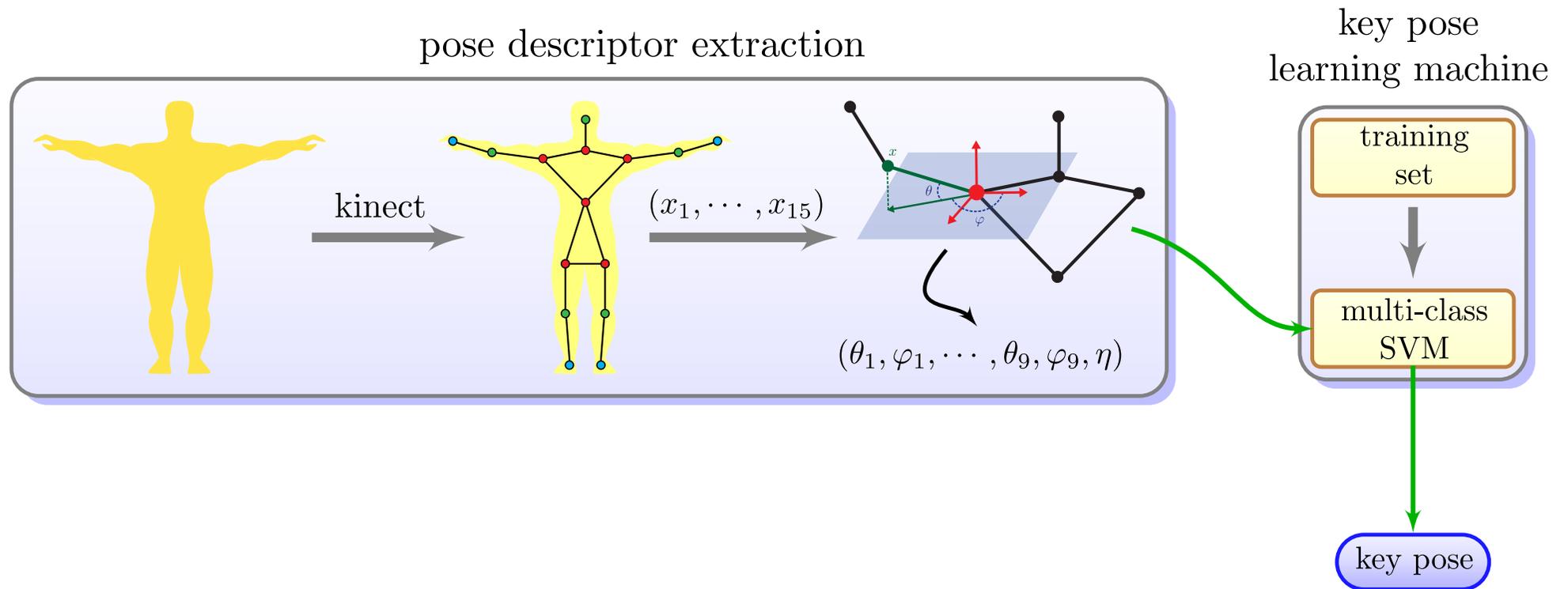
# Overview: training gestures



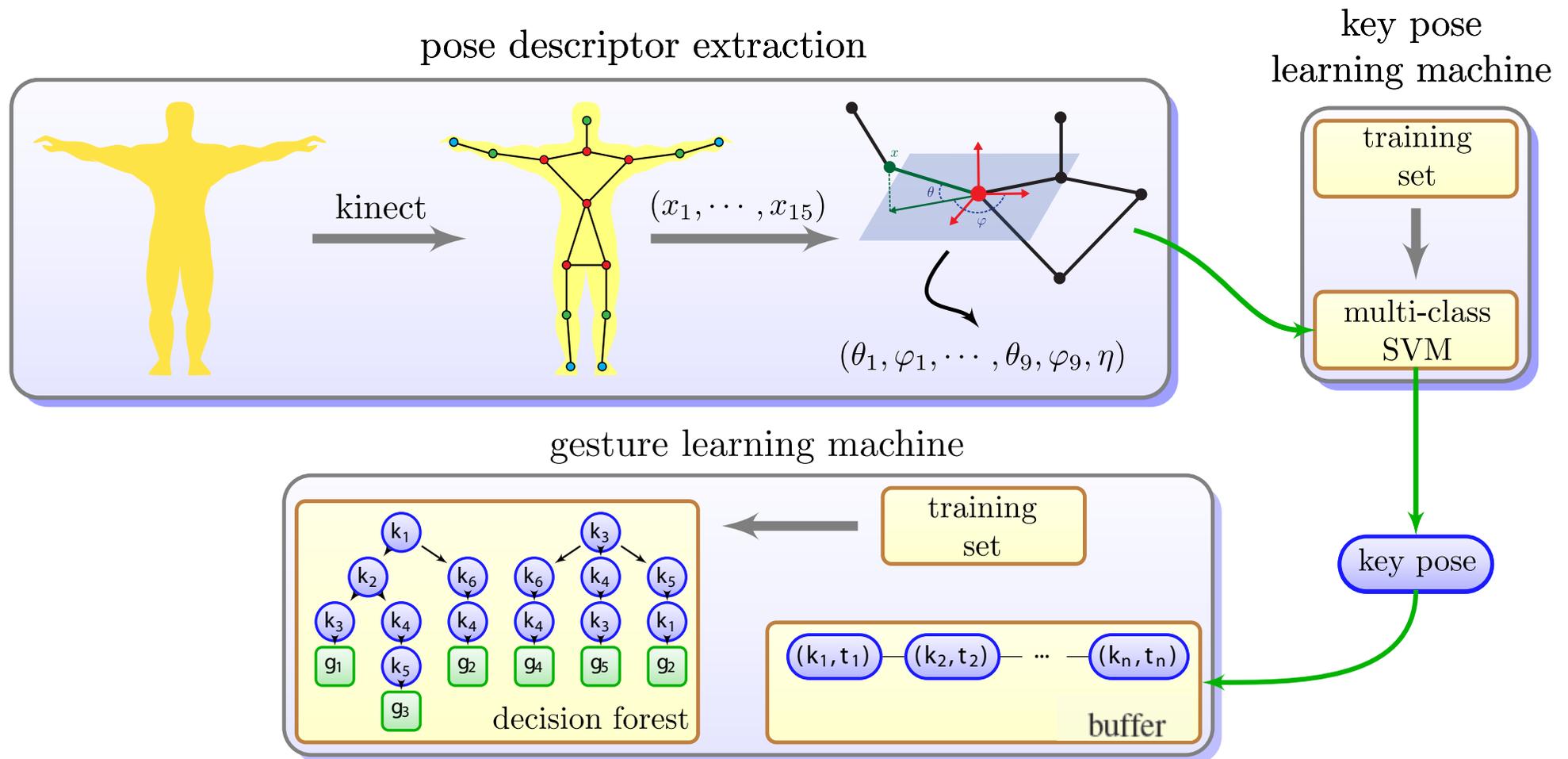
# Overview: training gestures



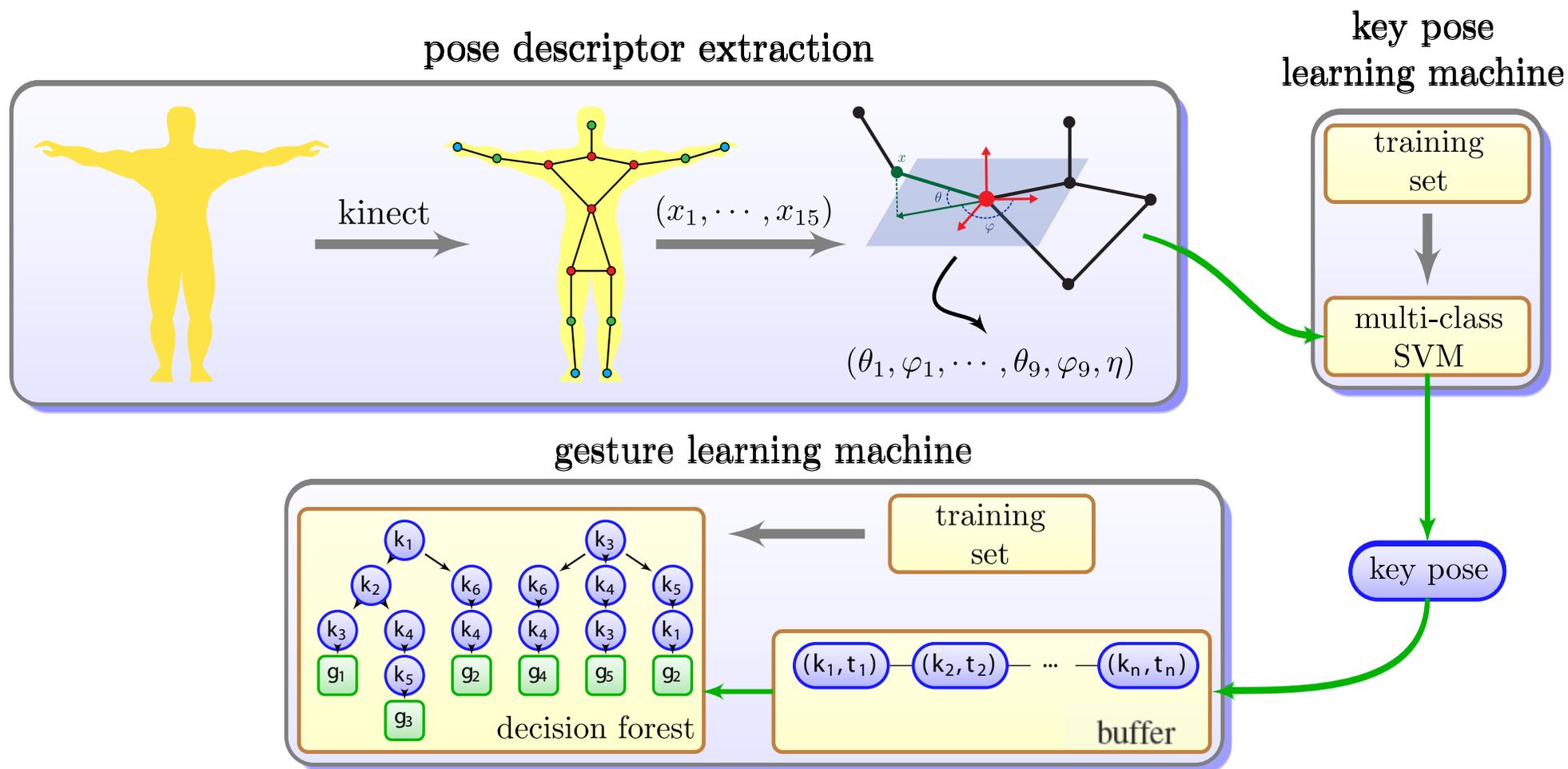
# Overview: recognizing gestures



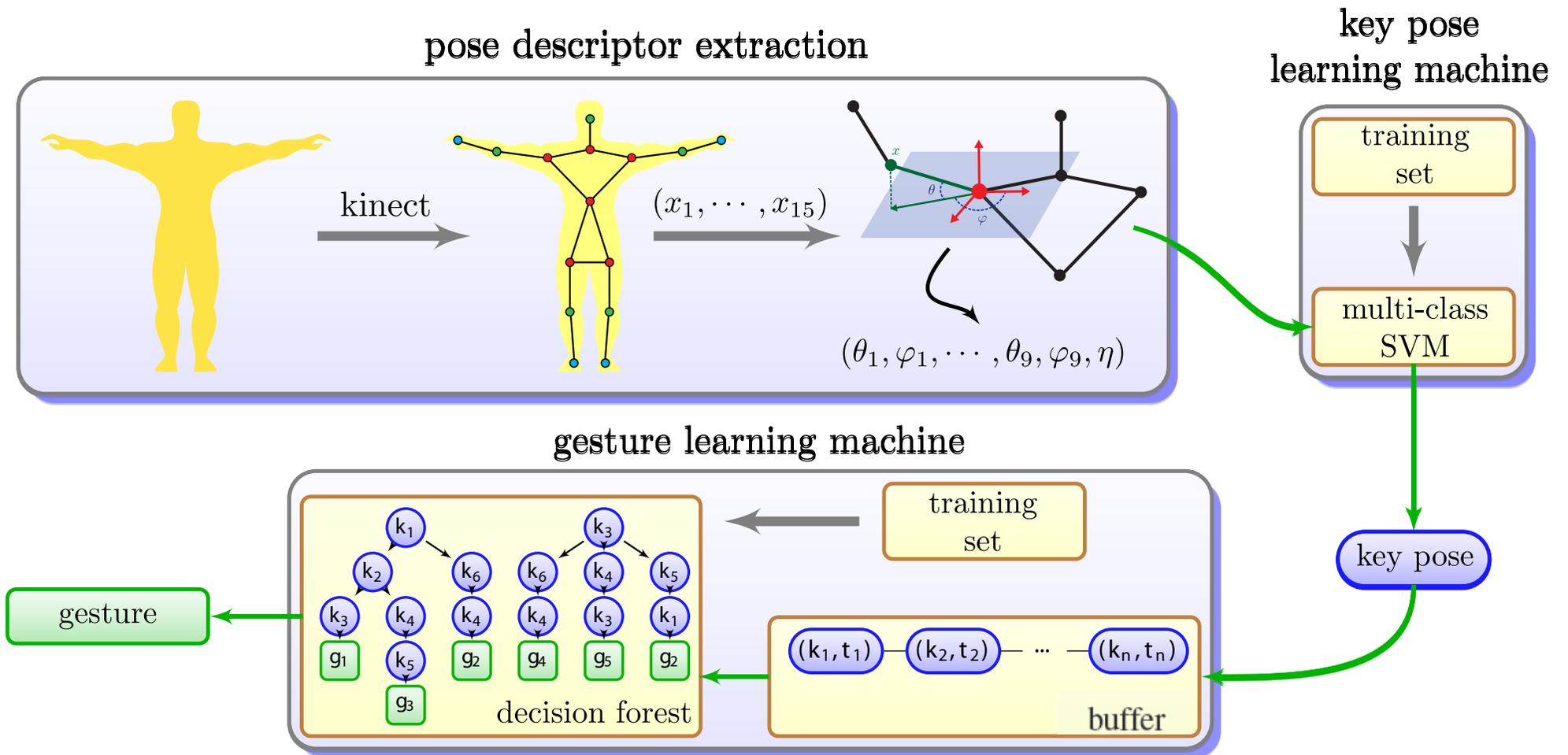
# Overview: recognizing gestures



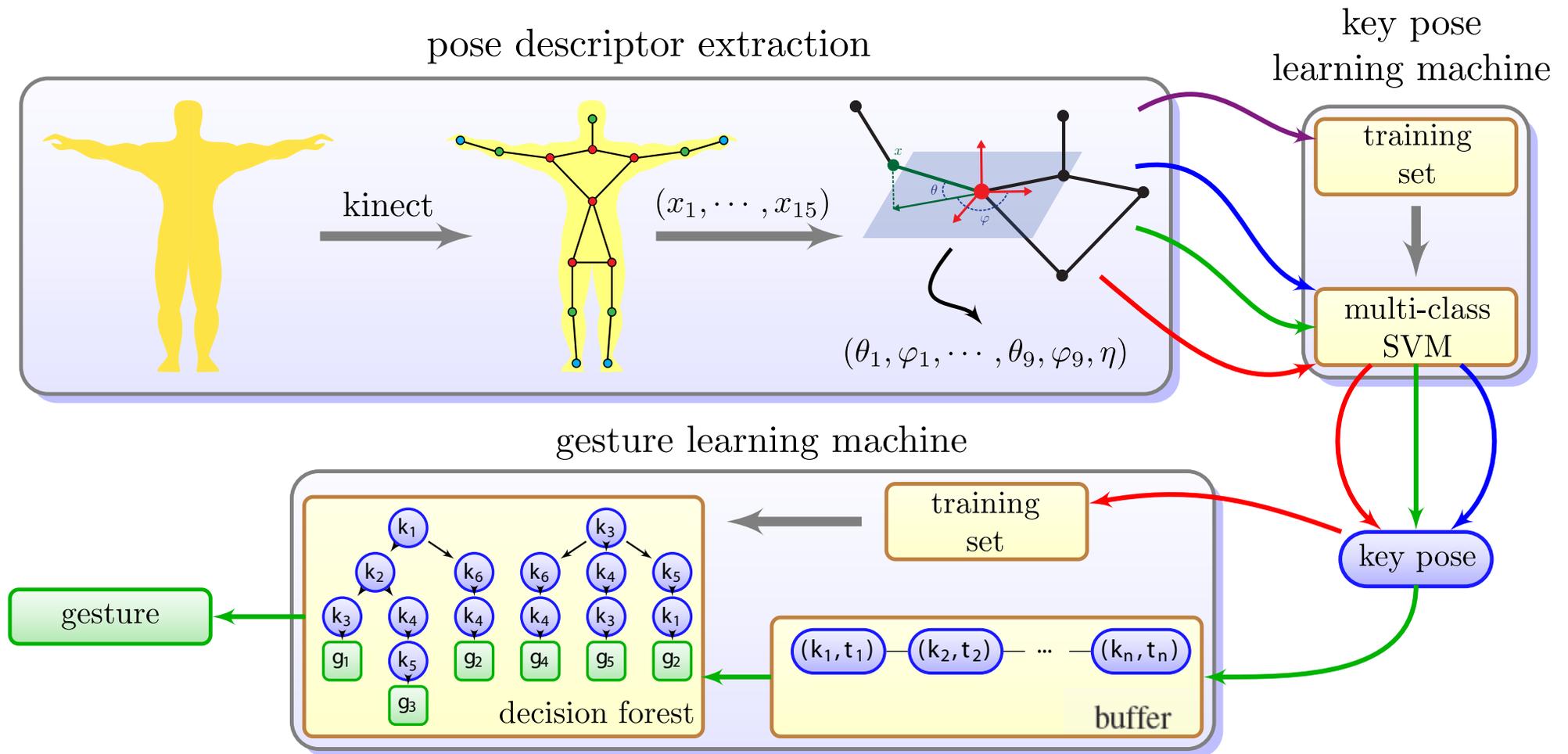
# Overview: recognizing gestures



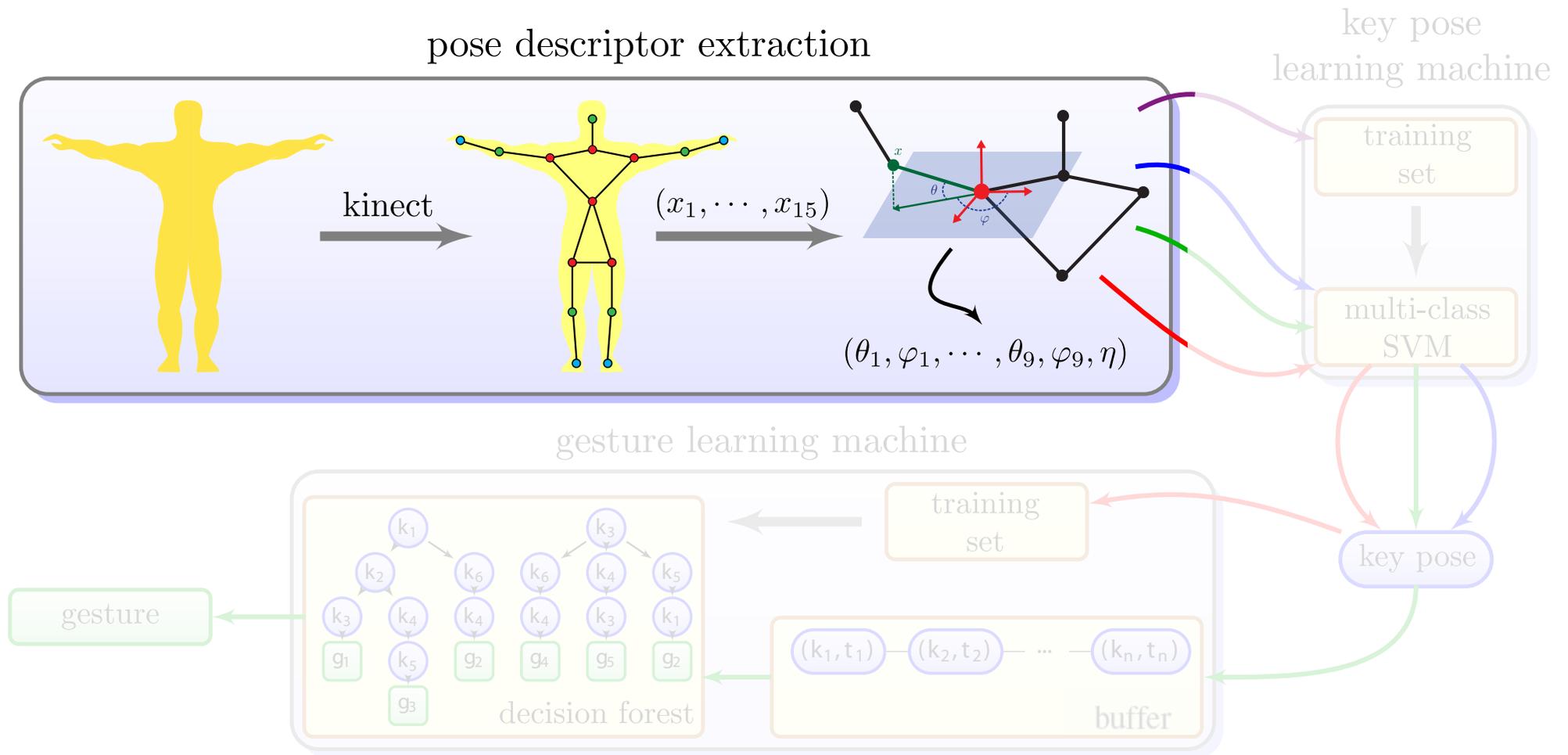
# Overview: recognizing gestures



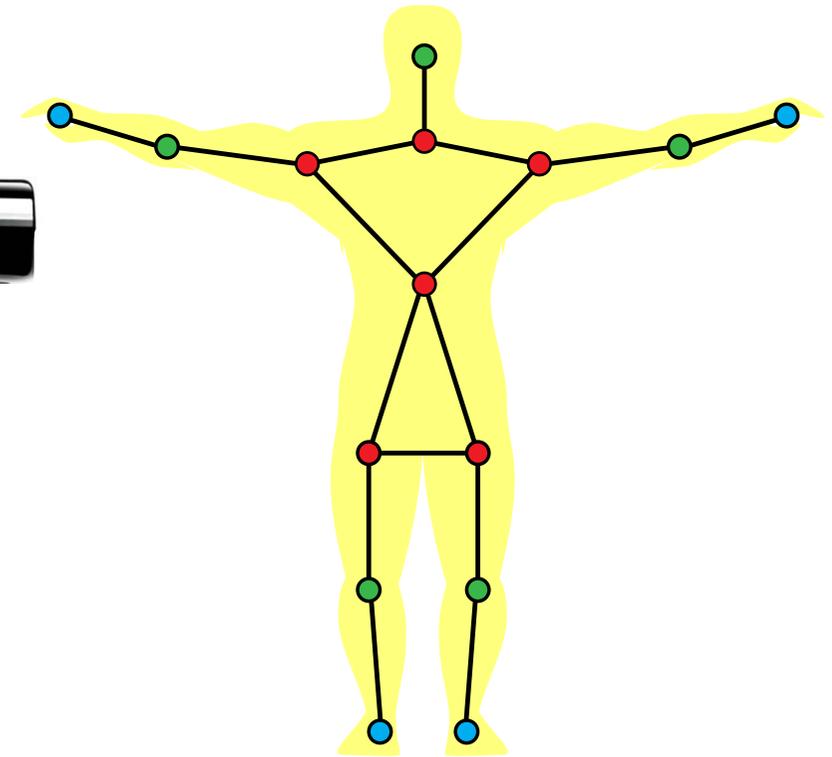
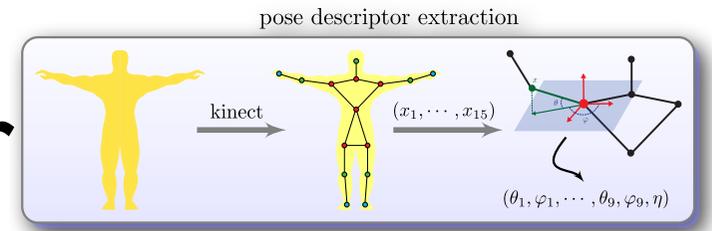
# Overview



# Overview

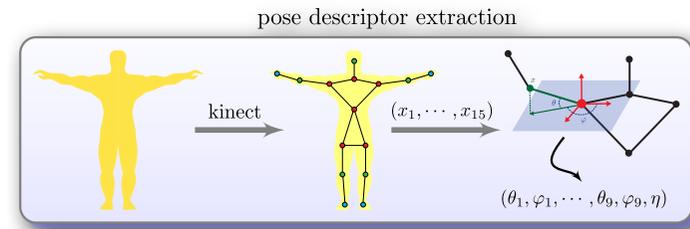


# Skeletons from Kinect Sensor



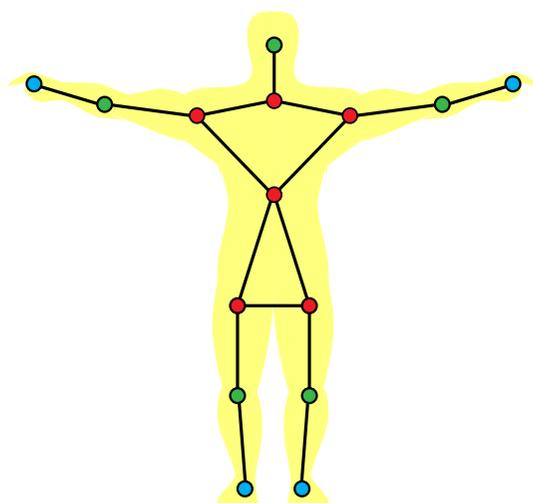
Real-time depth sensing system  
streaming depth data and skeletons at 30fps

# Joint-Angles Pose Descriptor

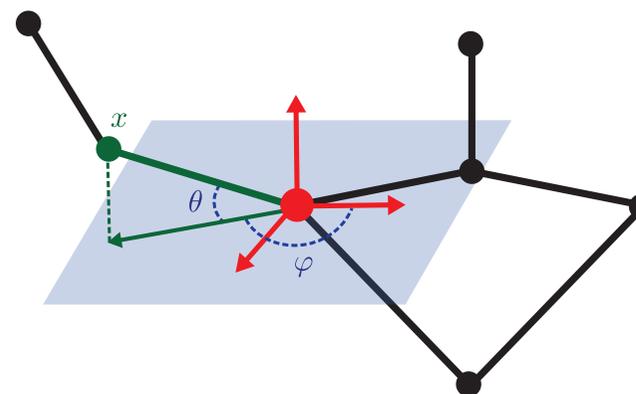


**Objective:** Concise and invariant representation of relevant pose information.

Improvement of Raptis *et al* (2011) local spherical coordinates.



$$(x_1, x_2, \dots, x_{15}) \in \mathbb{R}^{45}$$

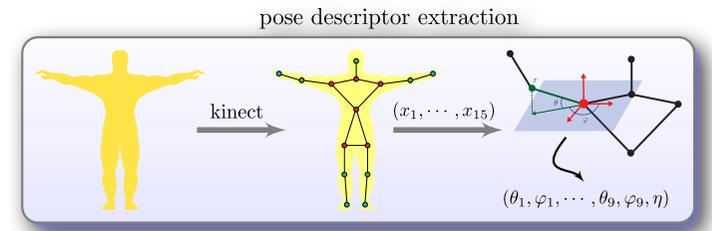


$$(\theta_1, \varphi_1, \dots, \theta_9, \varphi_9, \eta) \in \mathbb{R}^{19}$$

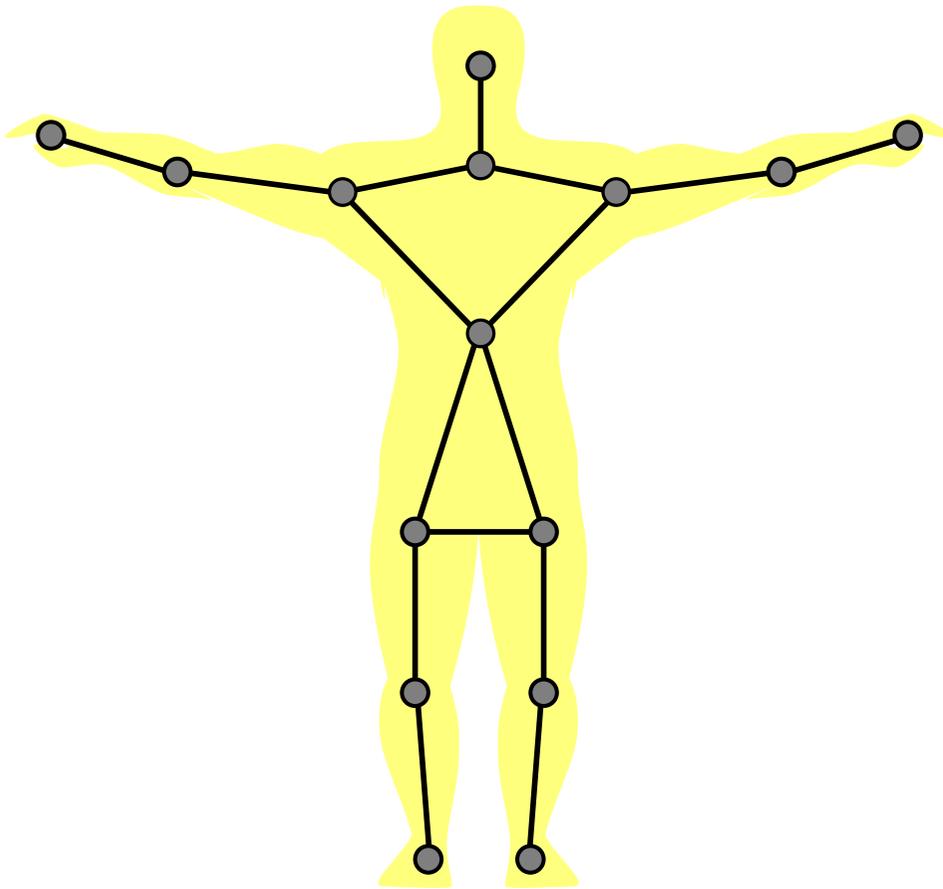
1st degree joints: elbows, knees and head

2nd degree joints: hands, feet.

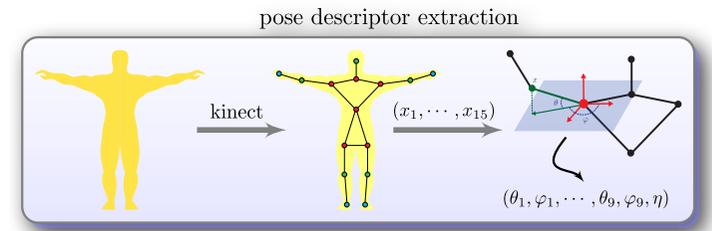
# How to compute the local bases?



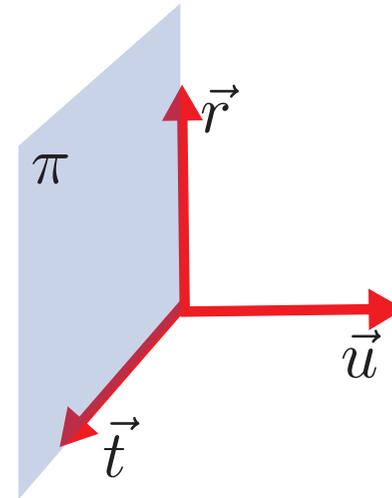
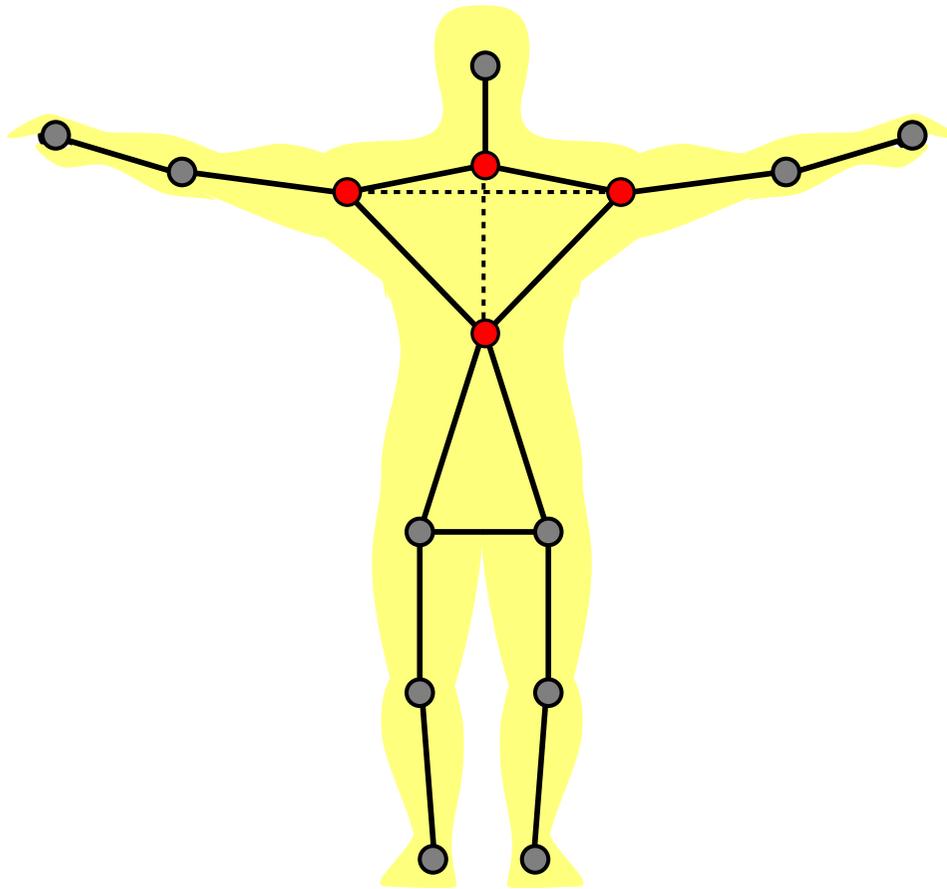
1st degree joints:



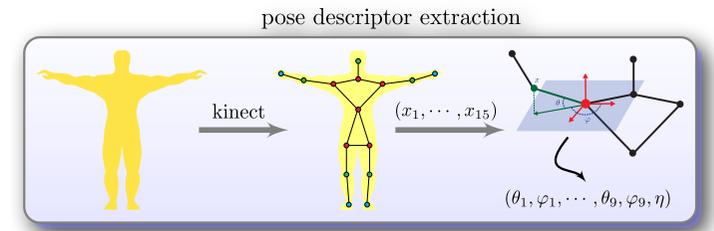
# How to compute the local bases?



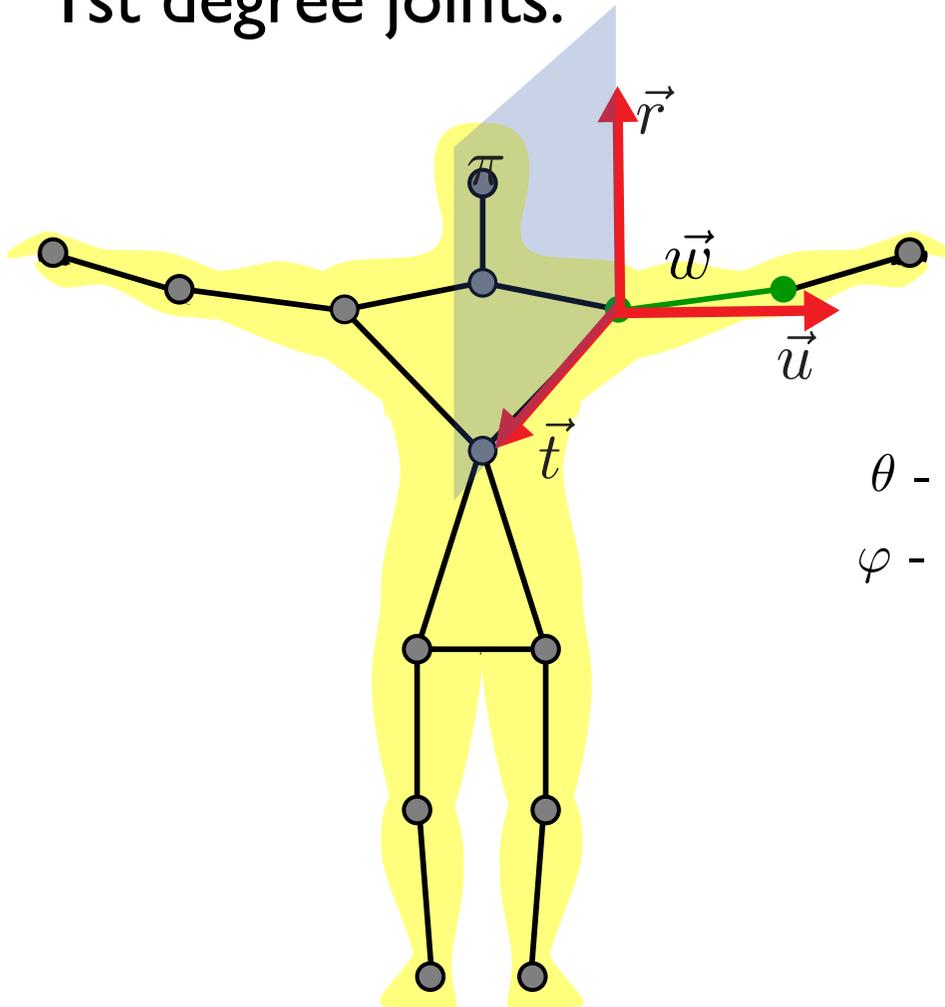
1st degree joints:



# How to compute the local bases?



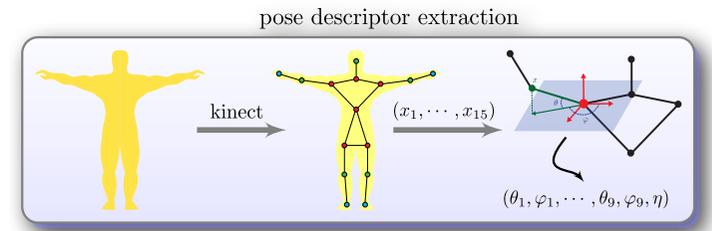
1st degree joints:



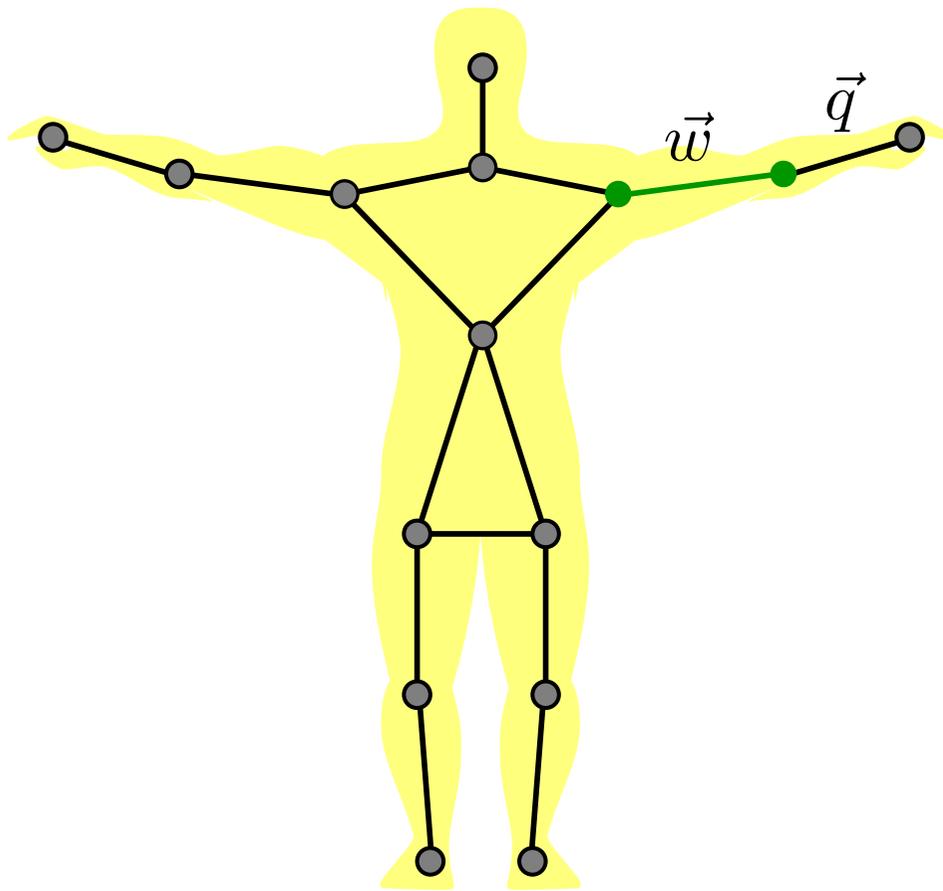
$\theta$  - angle between  $\vec{u}$  and  $\vec{w}$

$\varphi$  - angle between  $\vec{t}$  and the projection of  $\vec{w}$  in  $\pi$

# How to compute the local bases?



2nd degree joints:

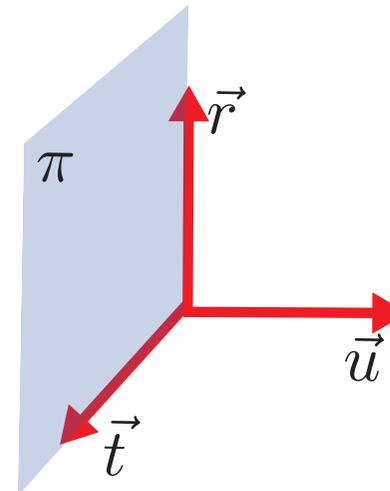


Rotate  $\{\vec{u}, \vec{r}, \vec{t}\}$  by

$$\beta = \arccos(\vec{w}, \vec{u})$$

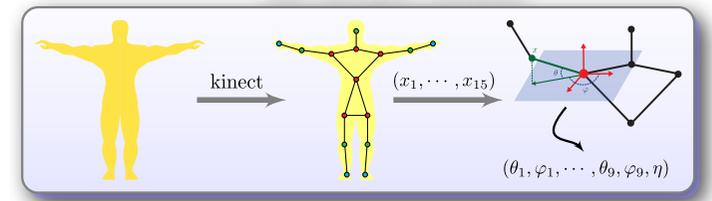
around

$$b = \vec{w} \times \vec{u}$$

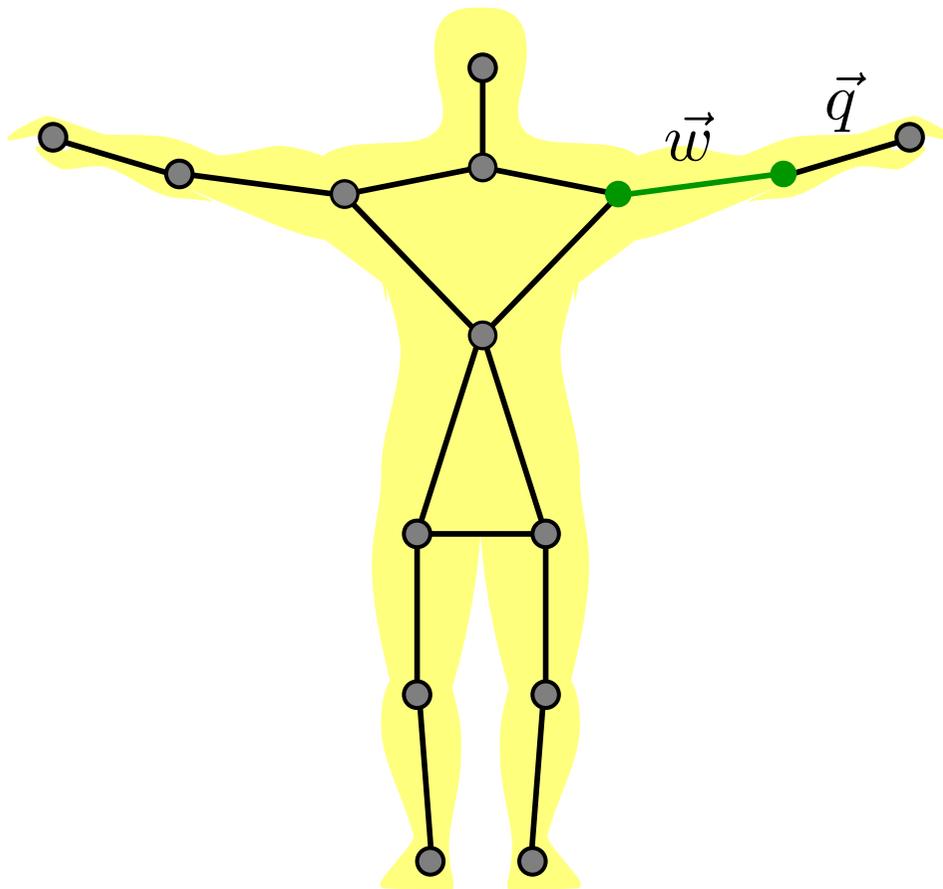


# How to compute the local bases?

pose descriptor extraction



2nd degree joints:

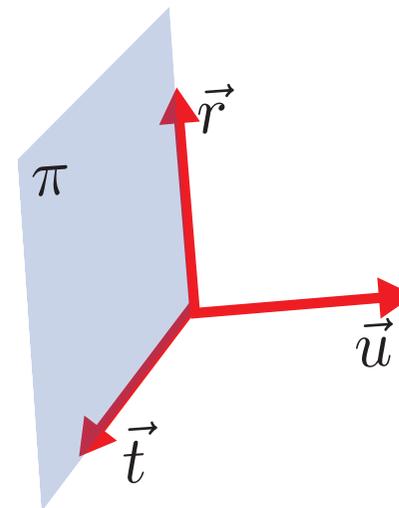


Rotate  $\{\vec{u}, \vec{r}, \vec{t}\}$  by

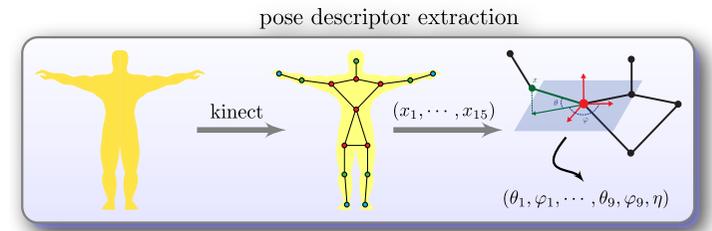
$$\beta = \arccos(\vec{w}, \vec{u})$$

around

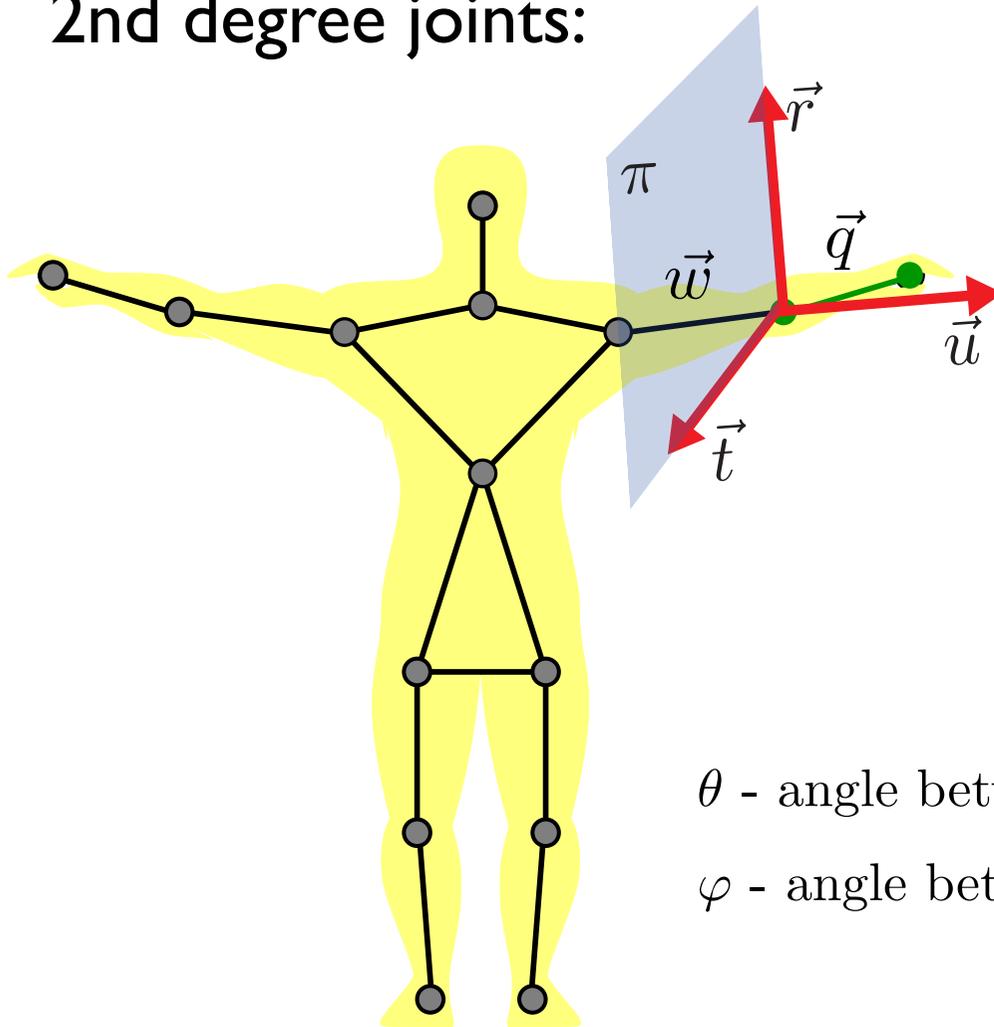
$$b = \vec{w} \times \vec{u}$$



# How to compute the local bases?



2nd degree joints:



Rotate  $\{\vec{u}, \vec{r}, \vec{t}\}$  by

$$\beta = \arccos(\vec{w}, \vec{u})$$

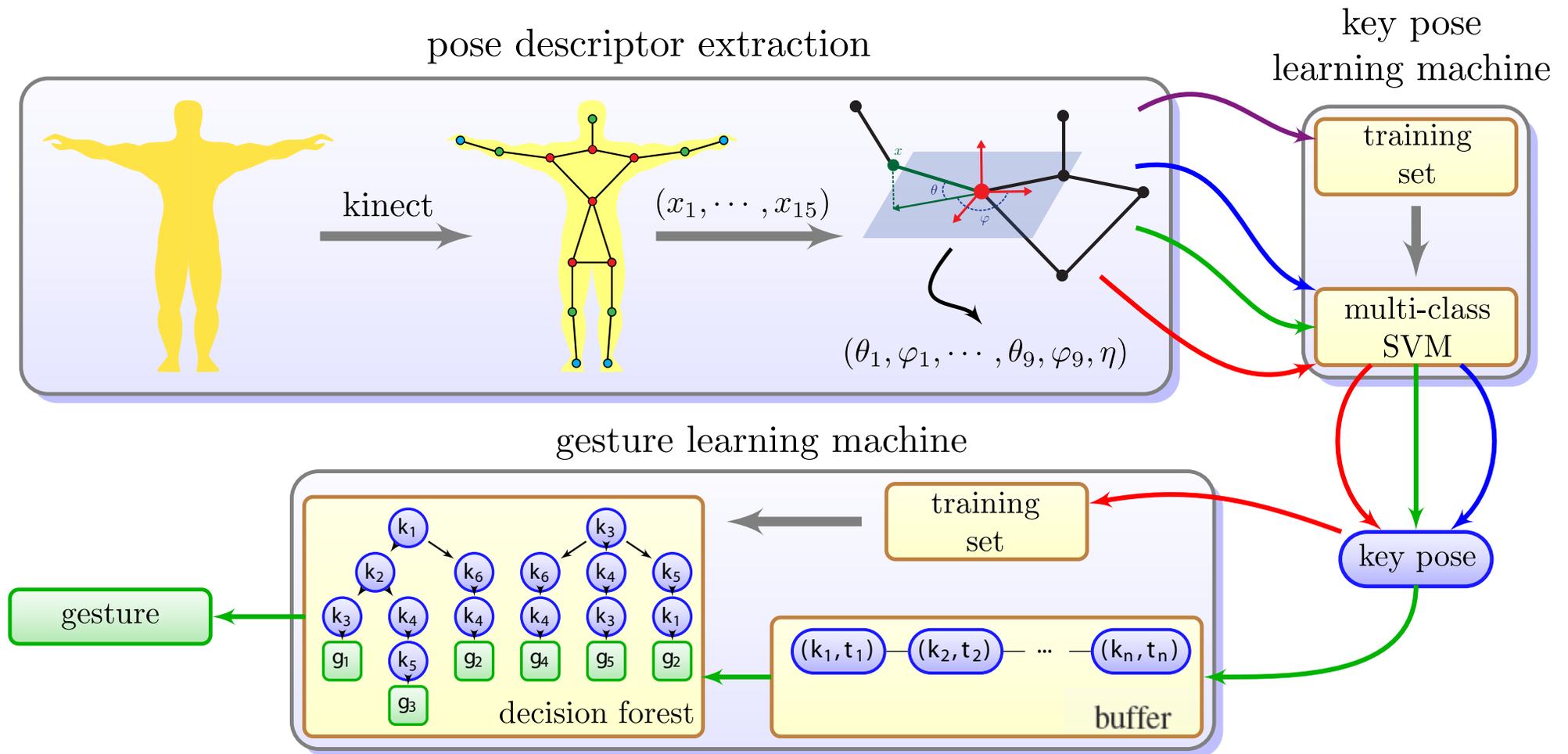
around

$$b = \vec{w} \times \vec{u}$$

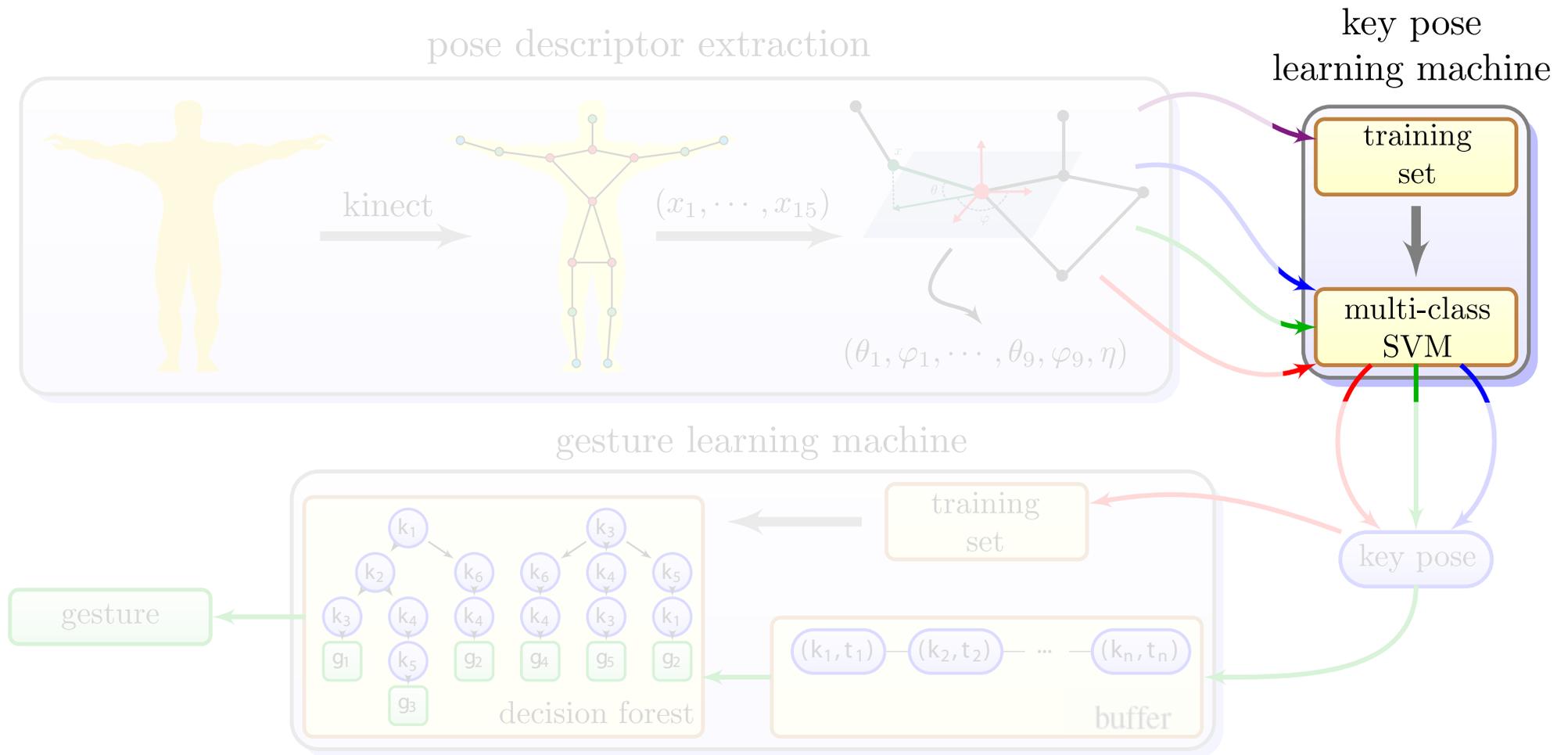
$\theta$  - angle between rotated  $\vec{u}$  and  $\vec{q}$

$\varphi$  - angle between rotated  $\vec{t}$  and the projection of  $\vec{q}$  in  $\pi$

# Overview

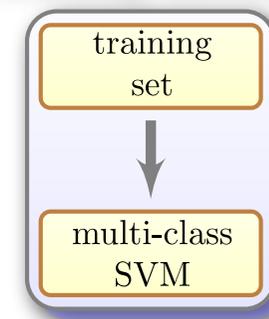


# Overview

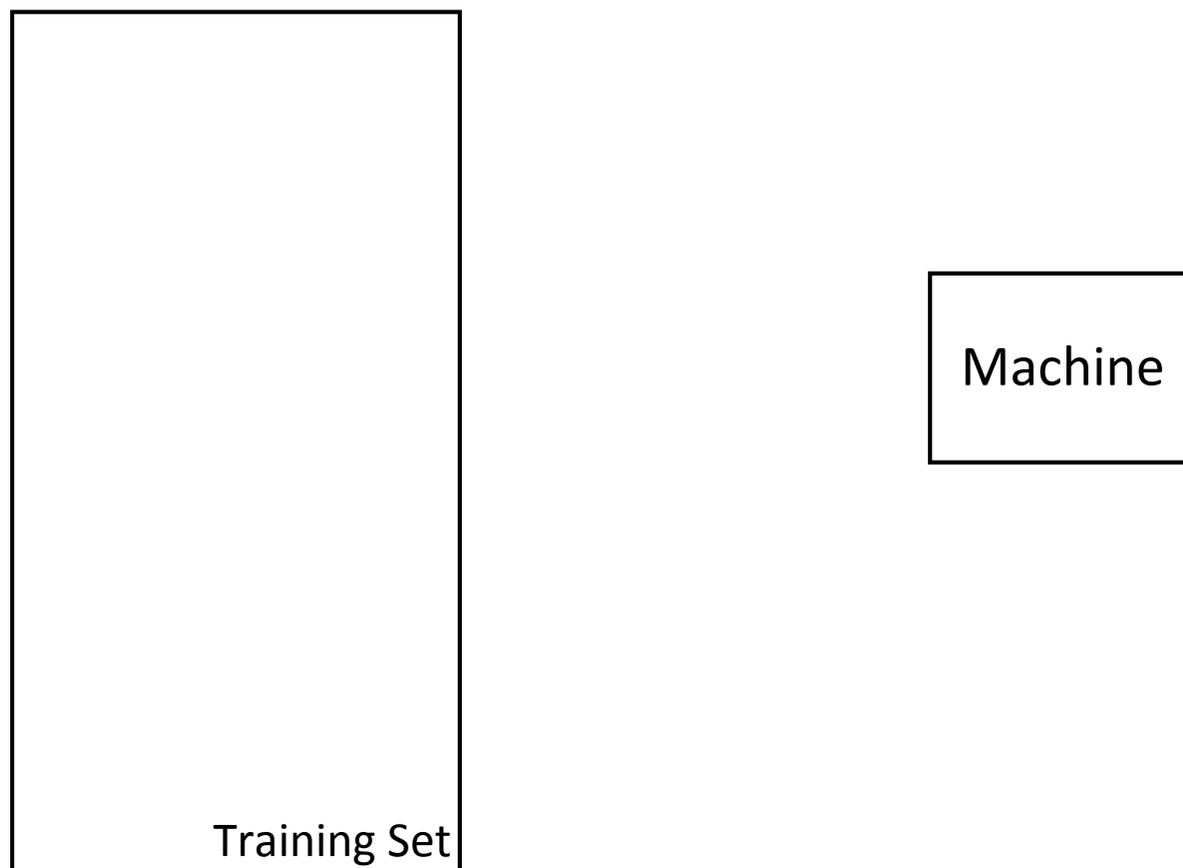


# Supervised Learning Machine

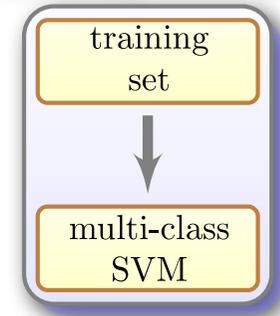
key pose  
learning machine



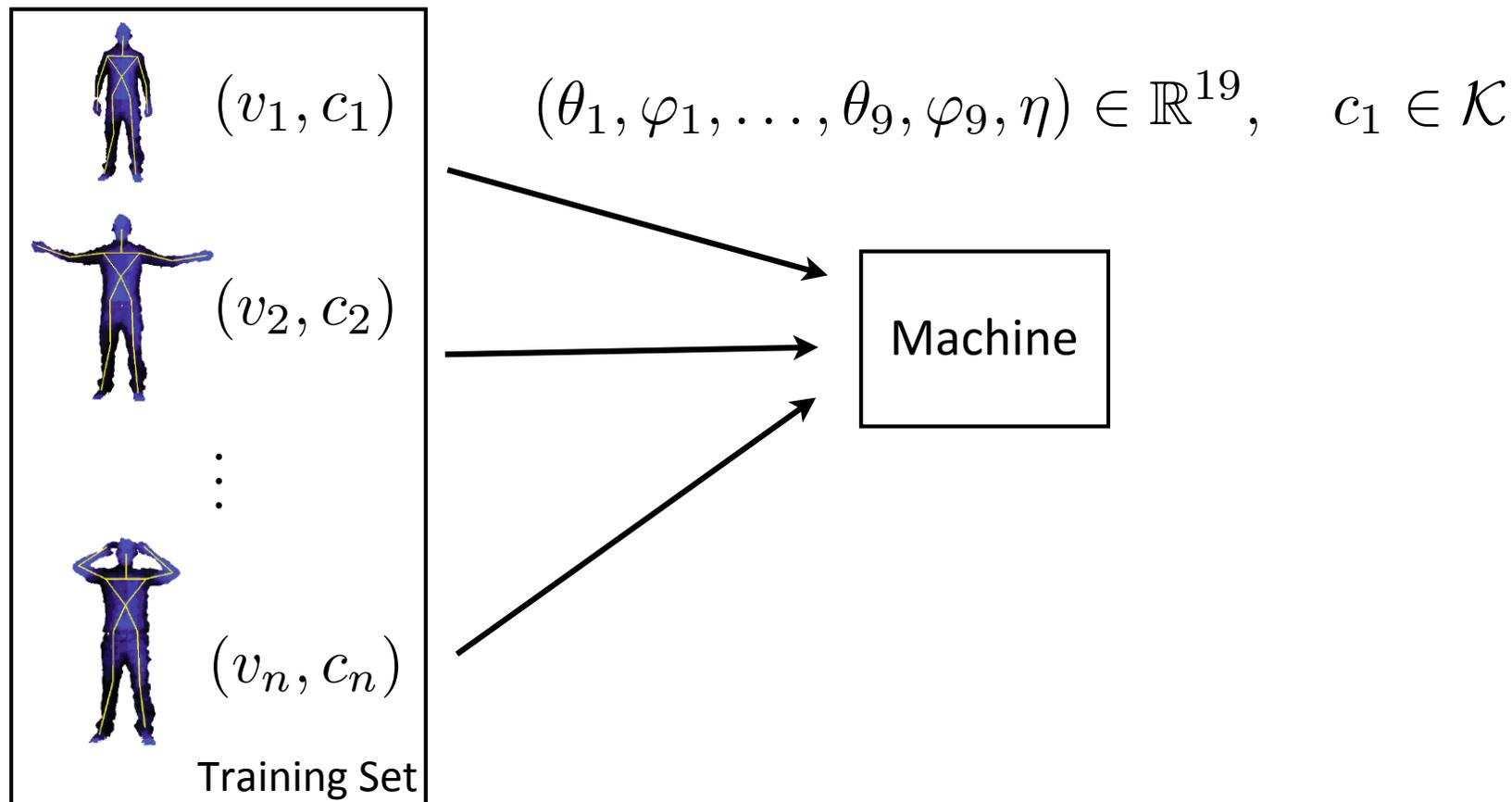
Predefined key pose classes:  $\mathcal{K} = \{k_1, k_2, \dots, k_{|\mathcal{K}|}\}$



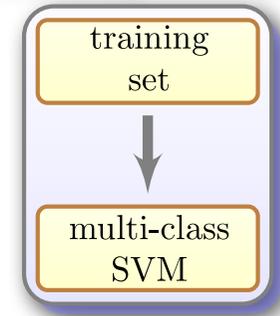
# Supervised Learning Machine



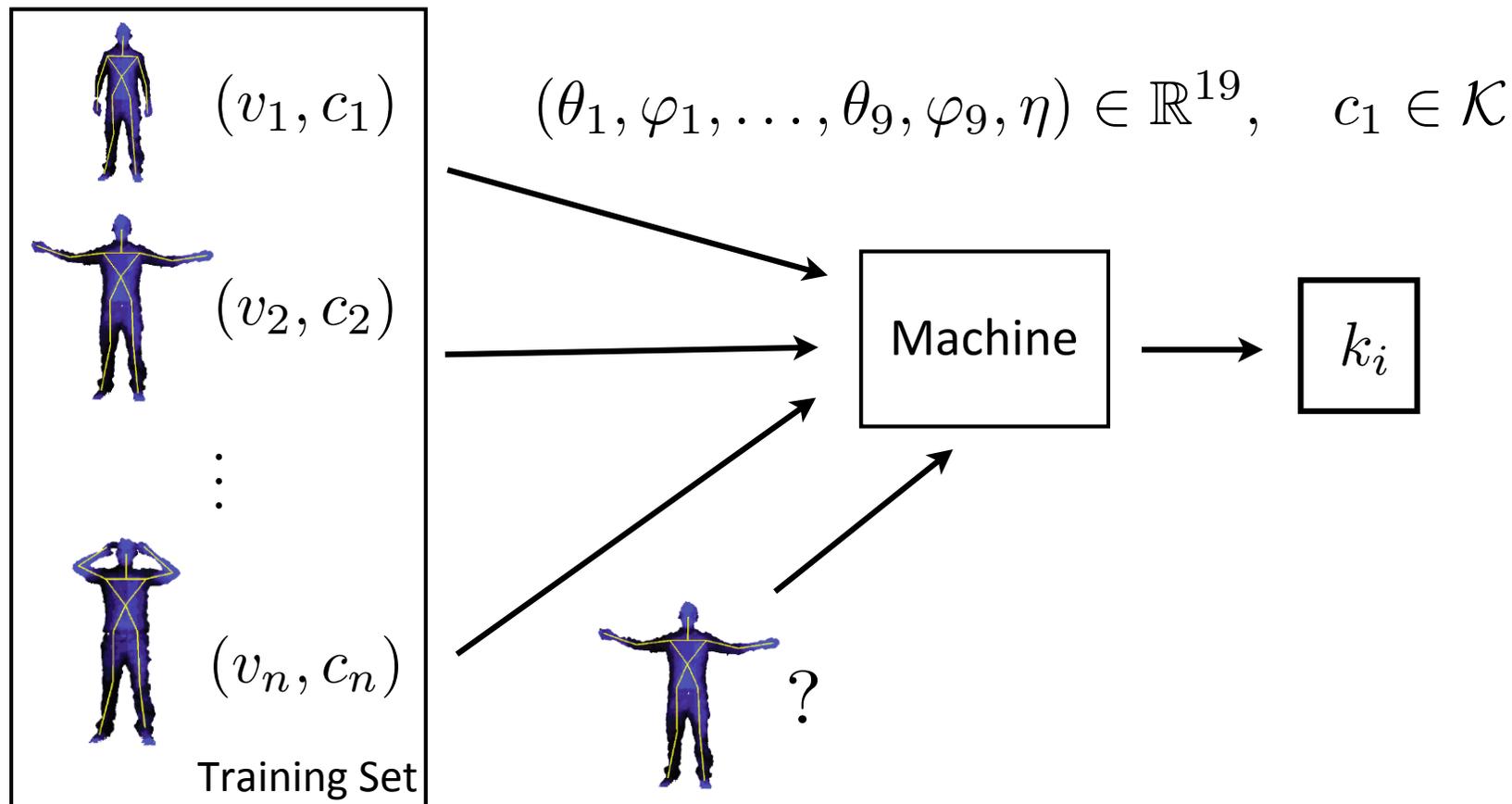
Predefined key pose classes:  $\mathcal{K} = \{k_1, k_2, \dots, k_{|\mathcal{K}|}\}$



# Supervised Learning Machine



Predefined key pose classes:  $\mathcal{K} = \{k_1, k_2, \dots, k_{|\mathcal{K}|}\}$



# Support Vector Machines (SVM)

Binary classifier

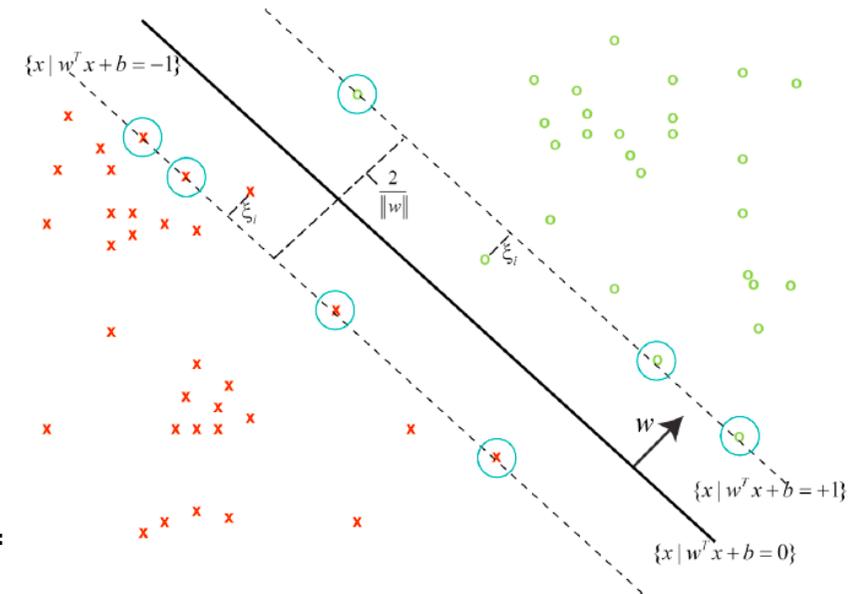
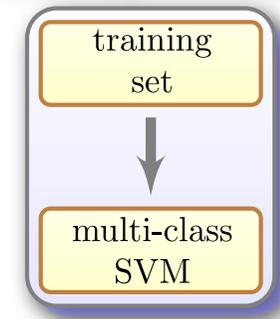
$$\hat{g} : \mathbb{R}^k \rightarrow \{-1, 1\}$$

$$v \rightarrow \text{sign}(\hat{f}(v)) = \{-1, 1\}$$

$$\hat{f}(v) = \sum_j \alpha_j s_j \langle \varphi(v_j), \varphi(v) \rangle + b$$

$$\text{MAX}_{w, \gamma} \quad \gamma - C \sum_{i=1}^l \varepsilon_i$$

$$\text{subject to} \quad y_i \langle w, \Phi(x_i) \rangle \geq \gamma - \varepsilon_i, \quad \varepsilon_i \geq 0, \quad \|w\|^2 =$$

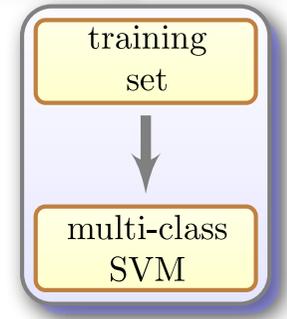


✓ Non-linear classification

✓ Efficiently computed for small training sets

# Multi-class SVM formulation

*One-versus-all* approach



One binary classifier for each key pose  $\mathbf{p} \in \mathcal{K}$ :

$$\hat{f}_{\mathbf{p}}(\mathbf{v}) = \sum_{j \in SV} \alpha_j \psi_{\mathbf{p}}(\mathbf{c}_j) \phi(\mathbf{v}_j, \mathbf{v}) + b,$$

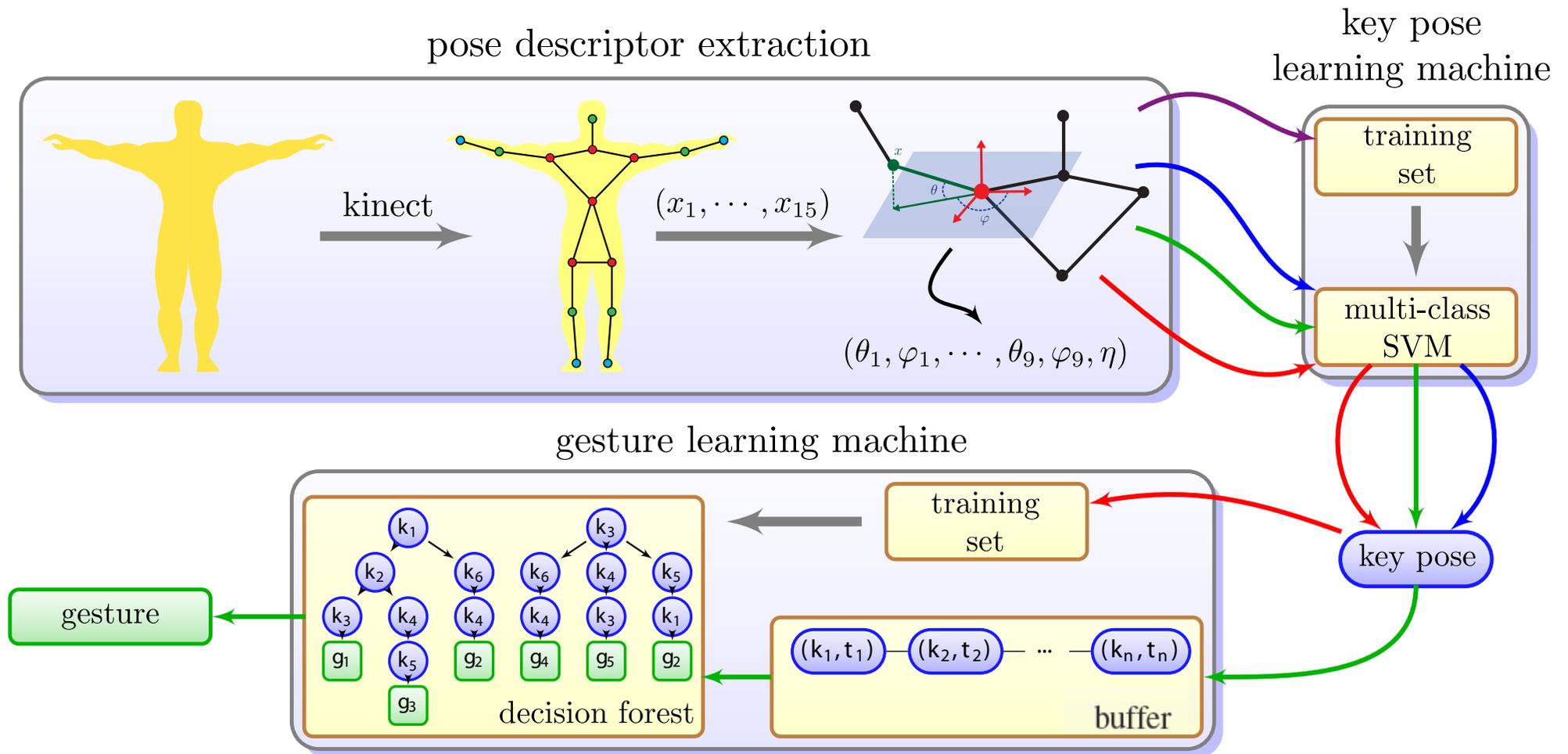
$$\text{where } \psi_{\mathbf{p}}(\mathbf{c}) = \begin{cases} 1 & \text{if } \mathbf{c} = \mathbf{p}, \\ -1 & \text{otherwise,} \end{cases}$$

$$\phi(\mathbf{v}_1, \mathbf{v}_2) = \exp\left(-\frac{\|\mathbf{v}_2 - \mathbf{v}_1\|^2}{2\sigma^2}\right)$$

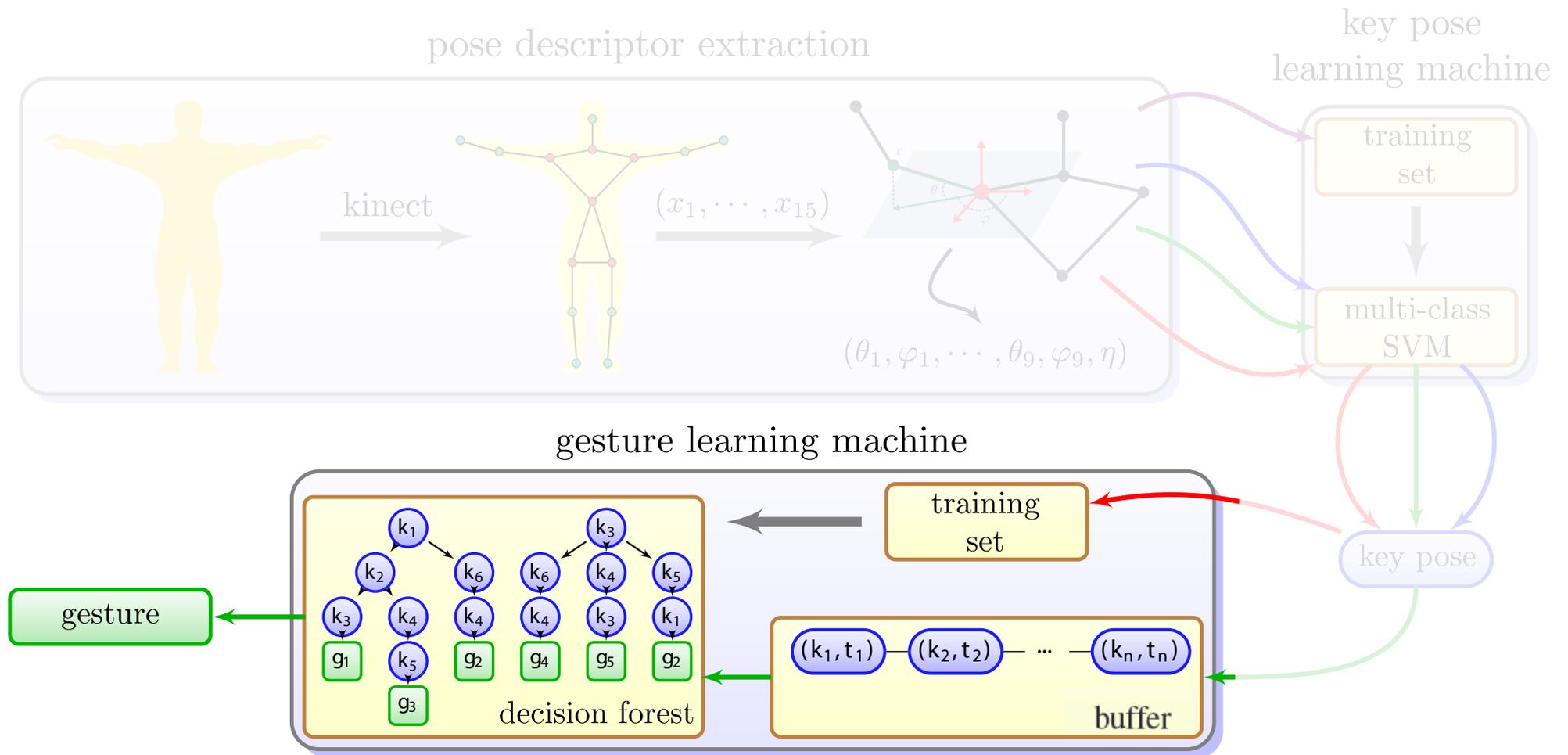
Voting process:

$$\hat{f}(\mathbf{v}) = \begin{cases} \mathbf{q} = \arg \max_{\mathbf{p}} \hat{f}_{\mathbf{p}}(\mathbf{v}) & \text{if } \hat{f}_{\mathbf{q}}(\mathbf{v}) > 0, \\ -1 & \text{otherwise.} \end{cases}$$

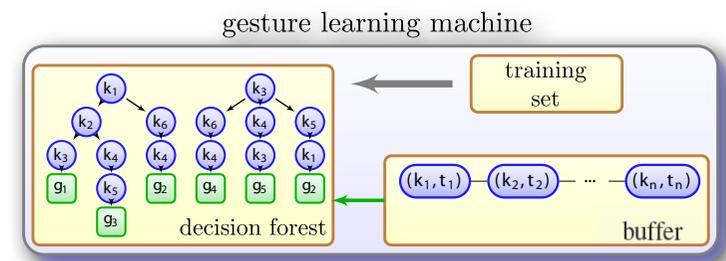
# Overview



# Overview

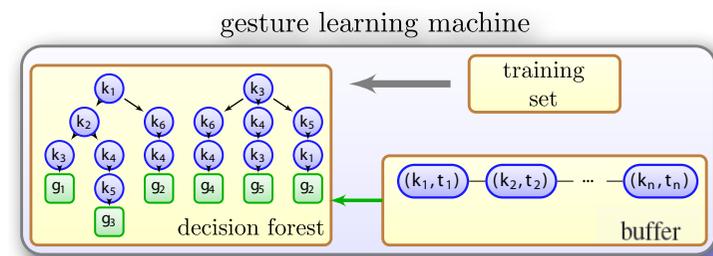


# Gestures as key pose sequences



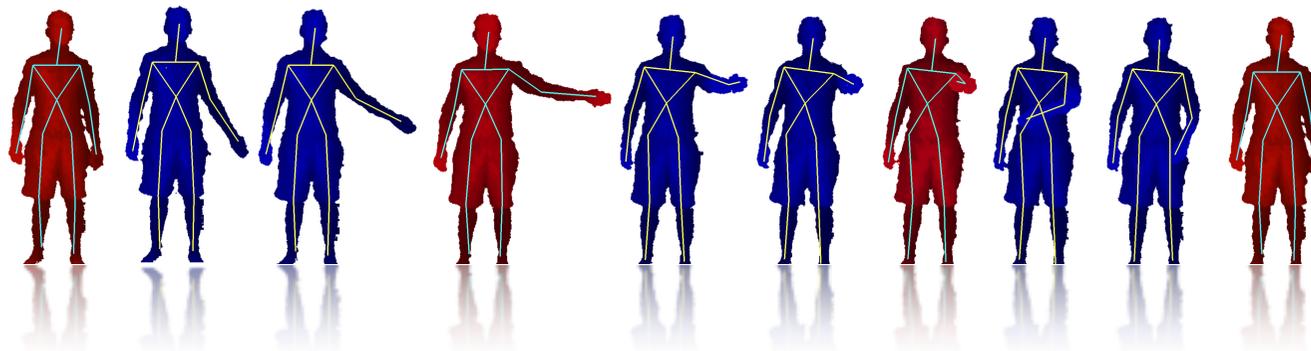
Gesture representation:  $g = \{k_1, k_2, \dots, k_{n_g}\}, \quad k_i \in \mathcal{K}.$

# Gestures as key pose sequences

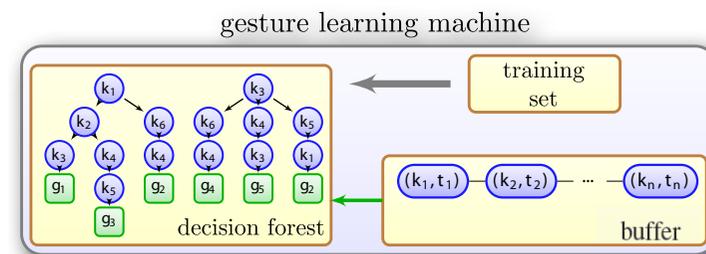


Gesture representation:  $g = \{k_1, k_2, \dots, k_{n_g}\}$ ,  $k_i \in \mathcal{K}$ .

Training session:

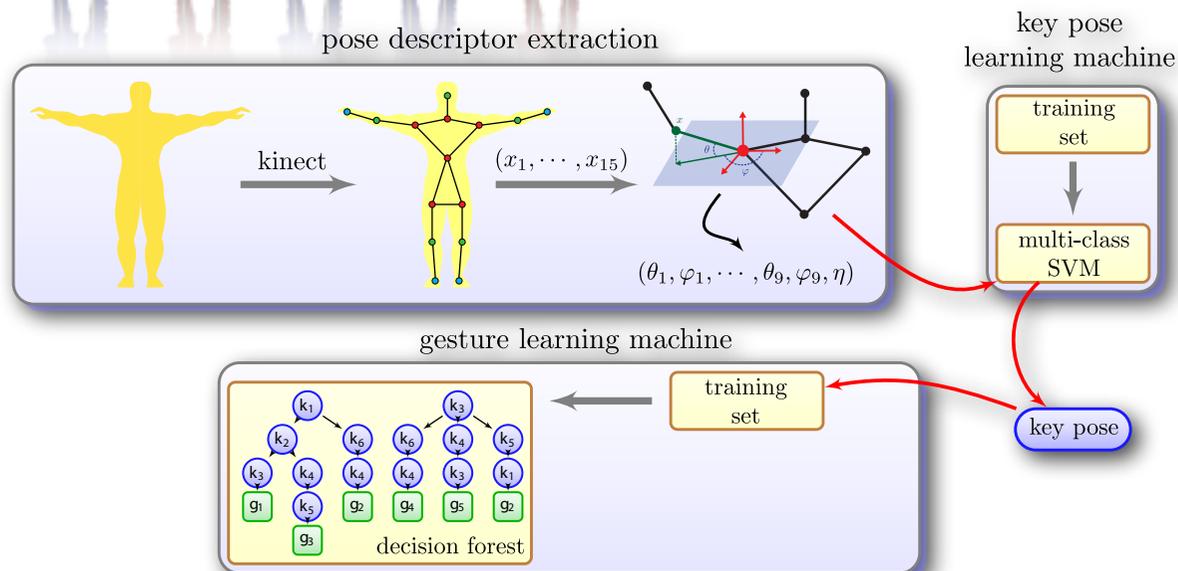
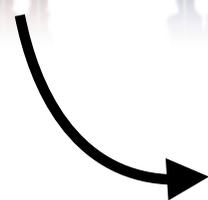
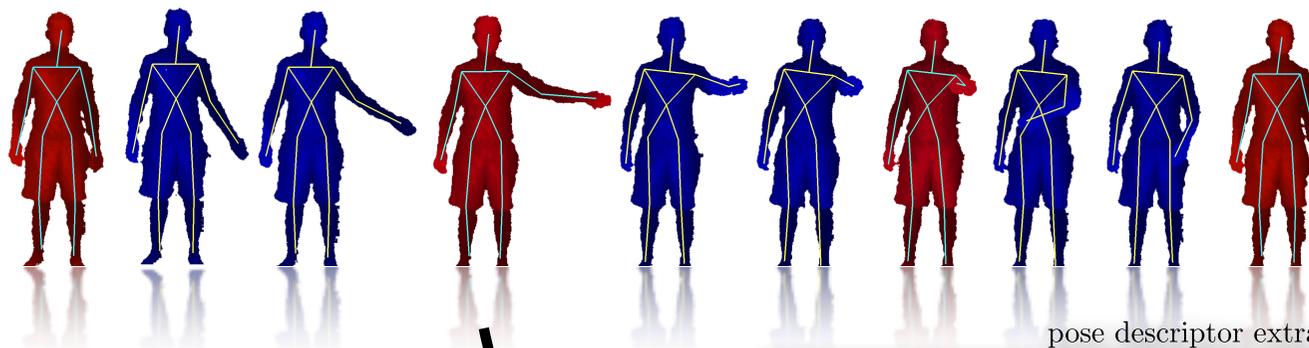


# Gestures as key pose sequences

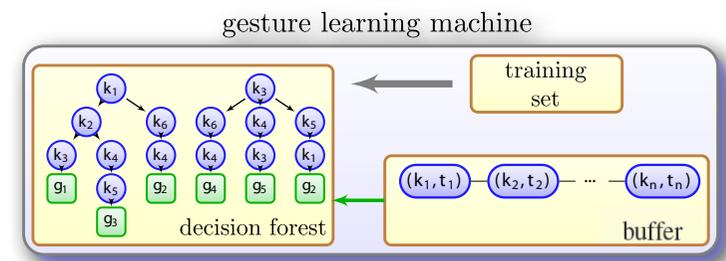


Gesture representation:  $g = \{k_1, k_2, \dots, k_{n_g}\}, k_i \in \mathcal{K}$ .

Training session:



# Decision Forests

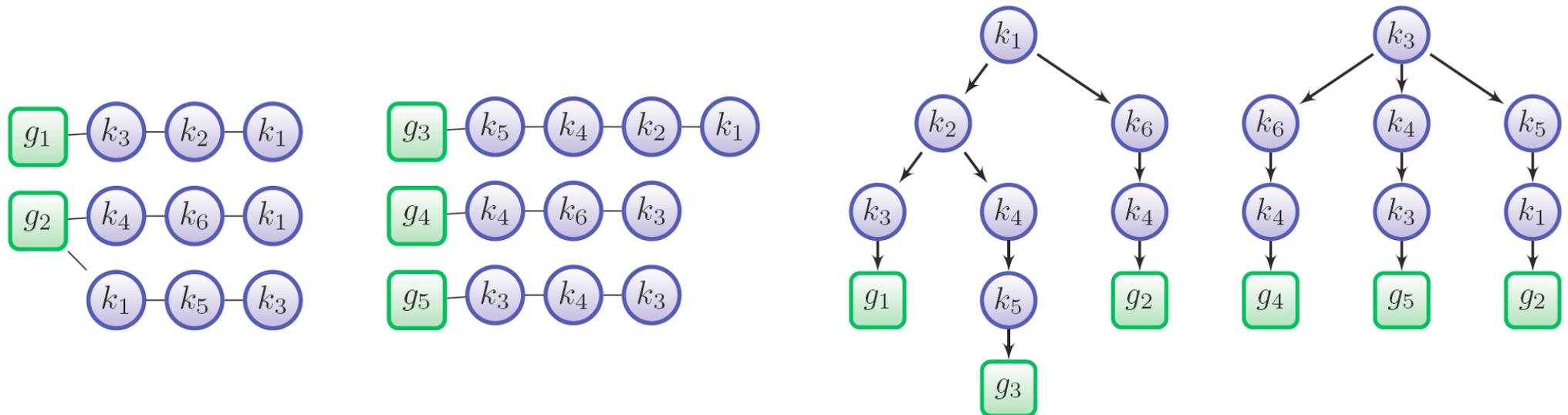


Each node represents a key pose

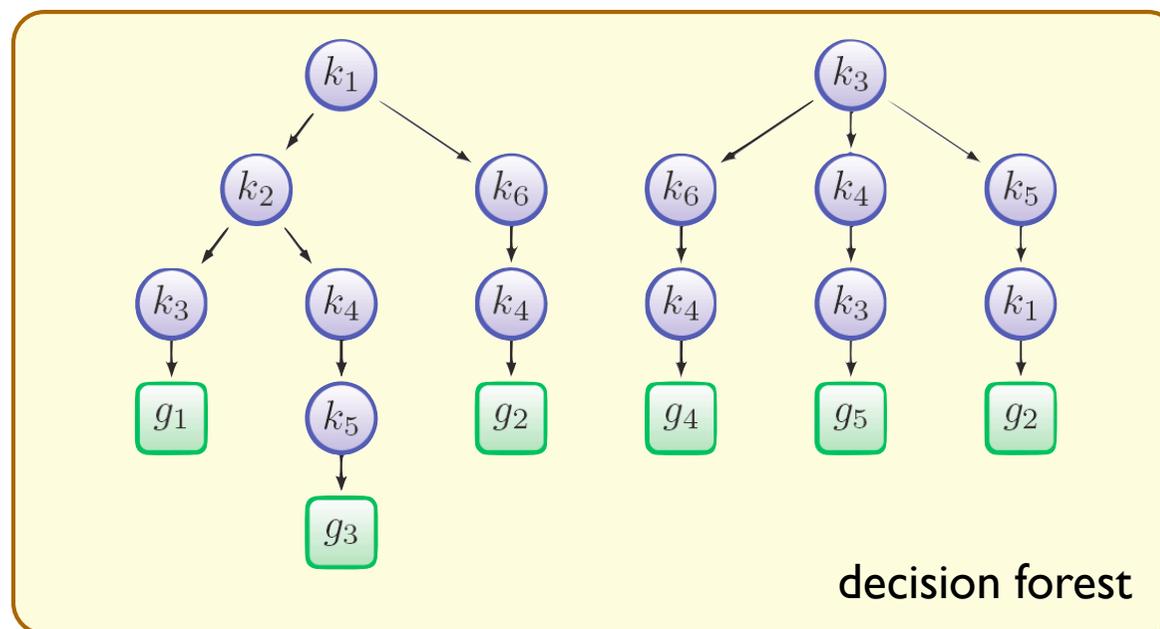
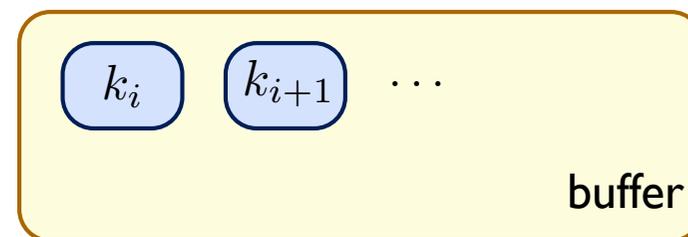
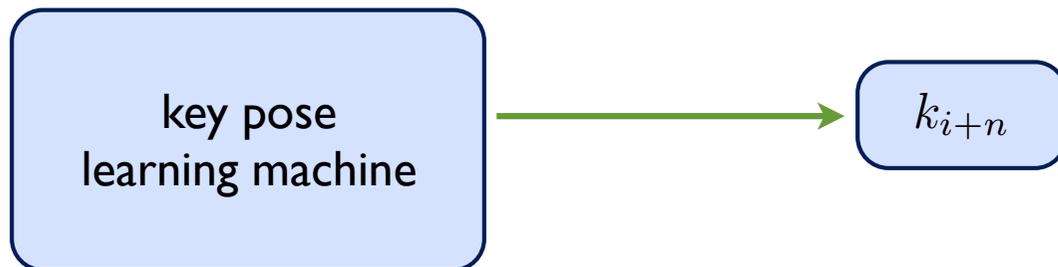
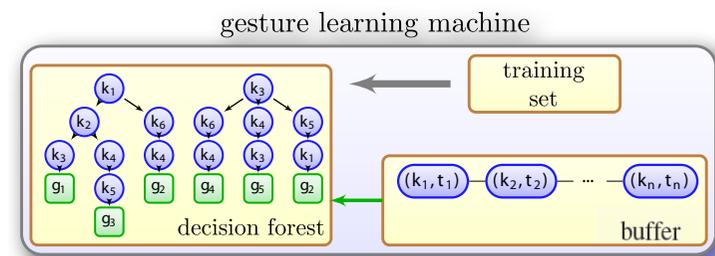
One tree per key pose

Each root-leaf path represents a gesture stored back-to-front

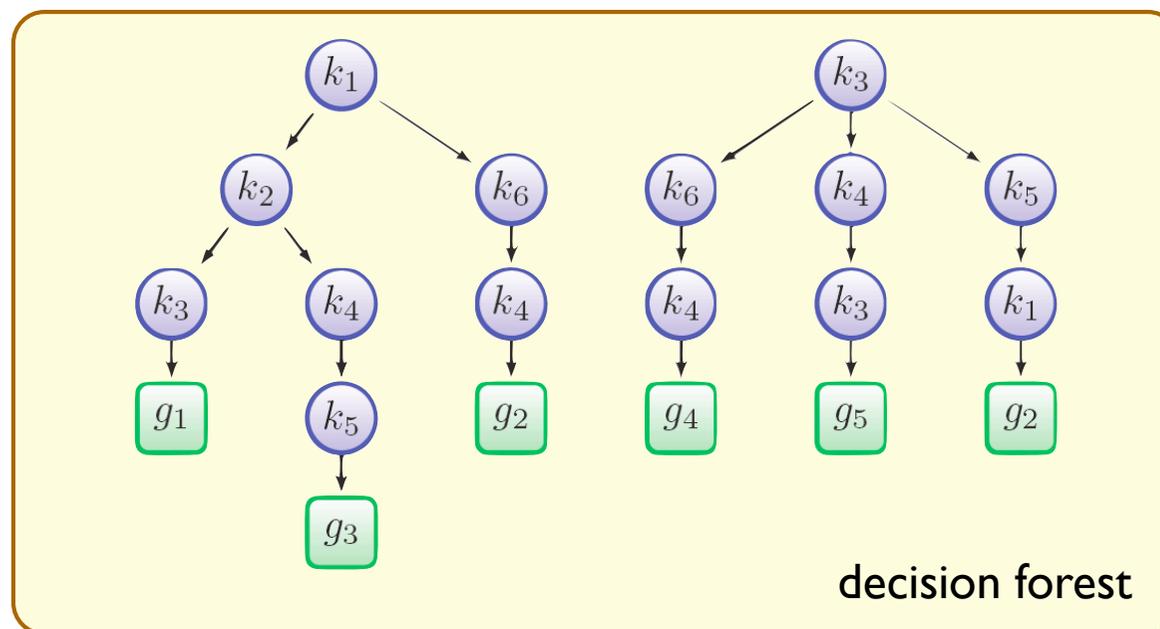
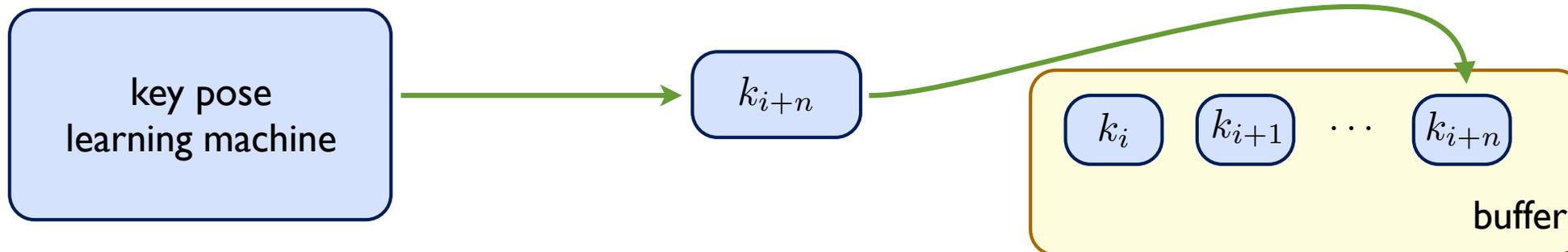
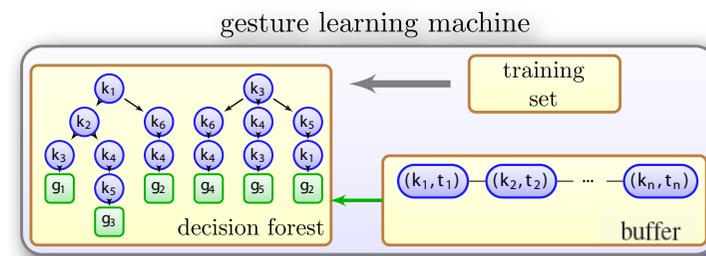
Two paths may represent the same gesture



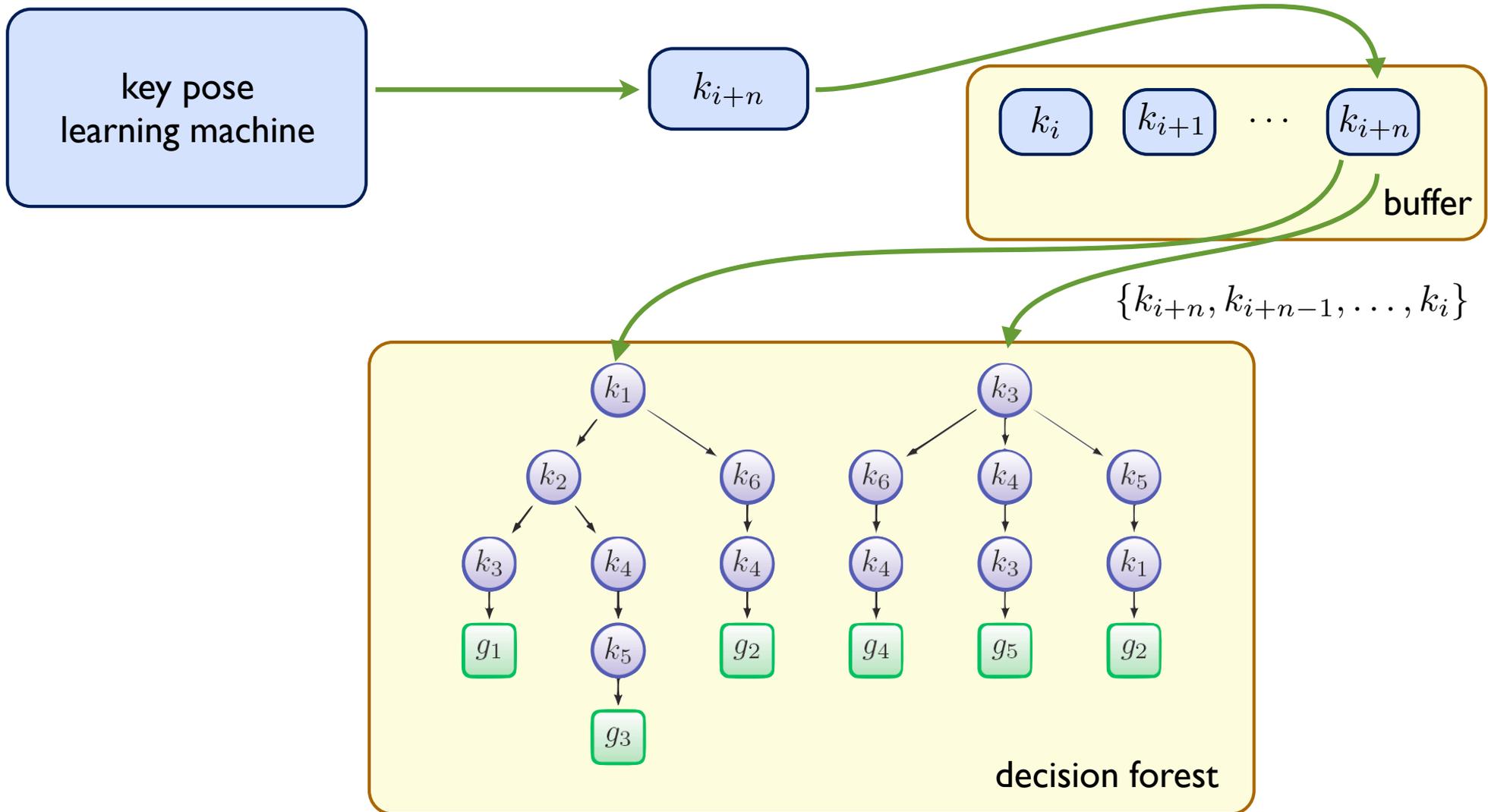
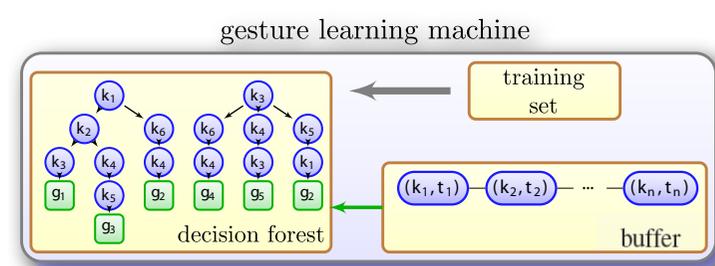
# Real-time gesture recognition



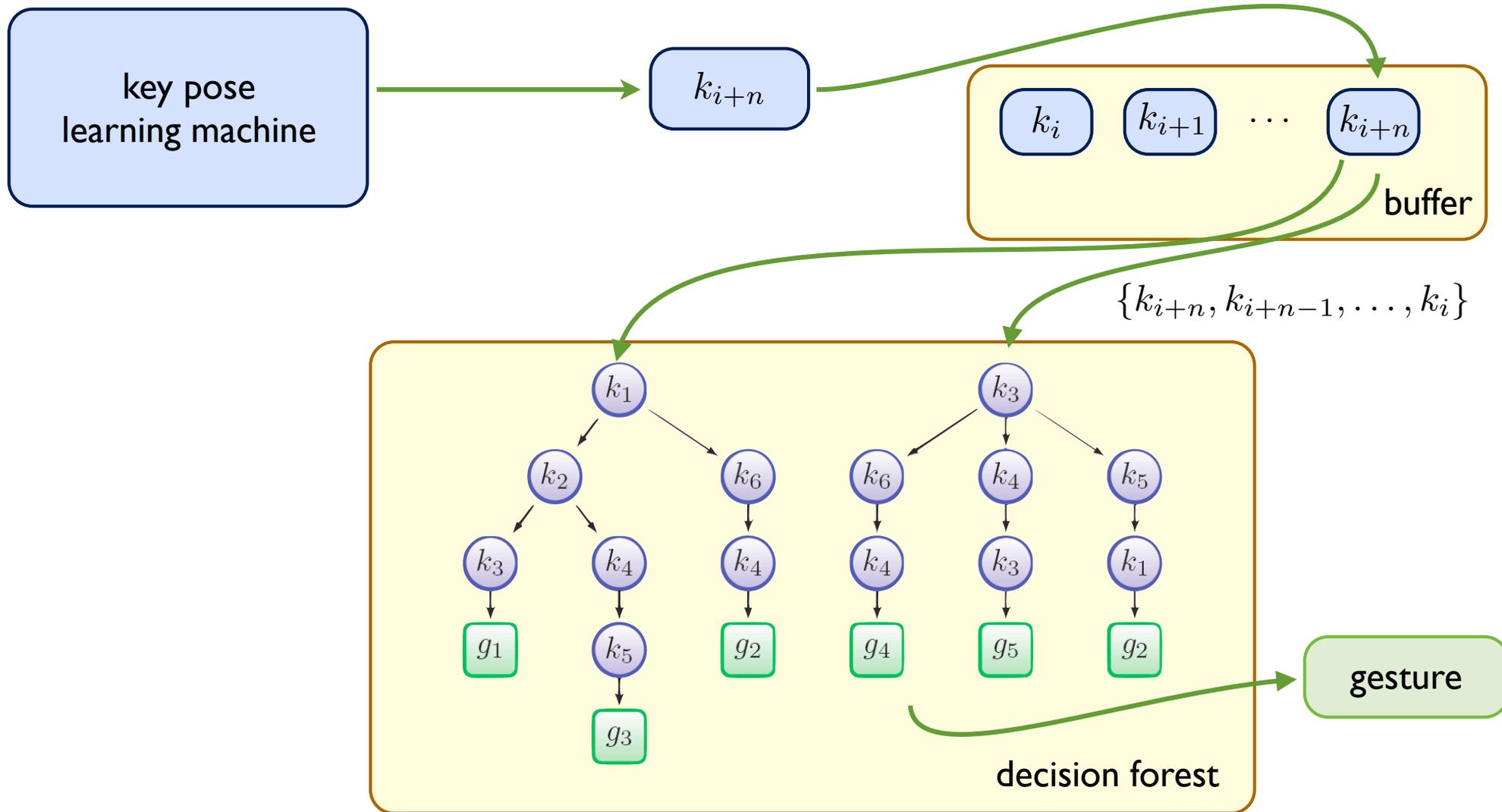
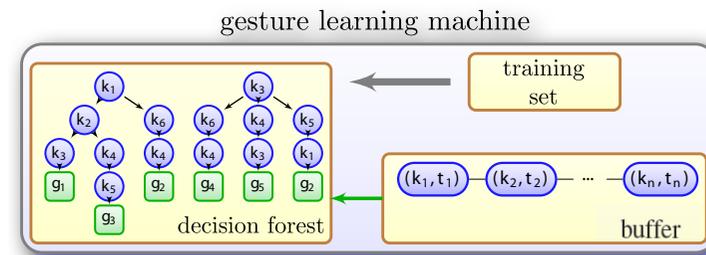
# Real-time gesture recognition



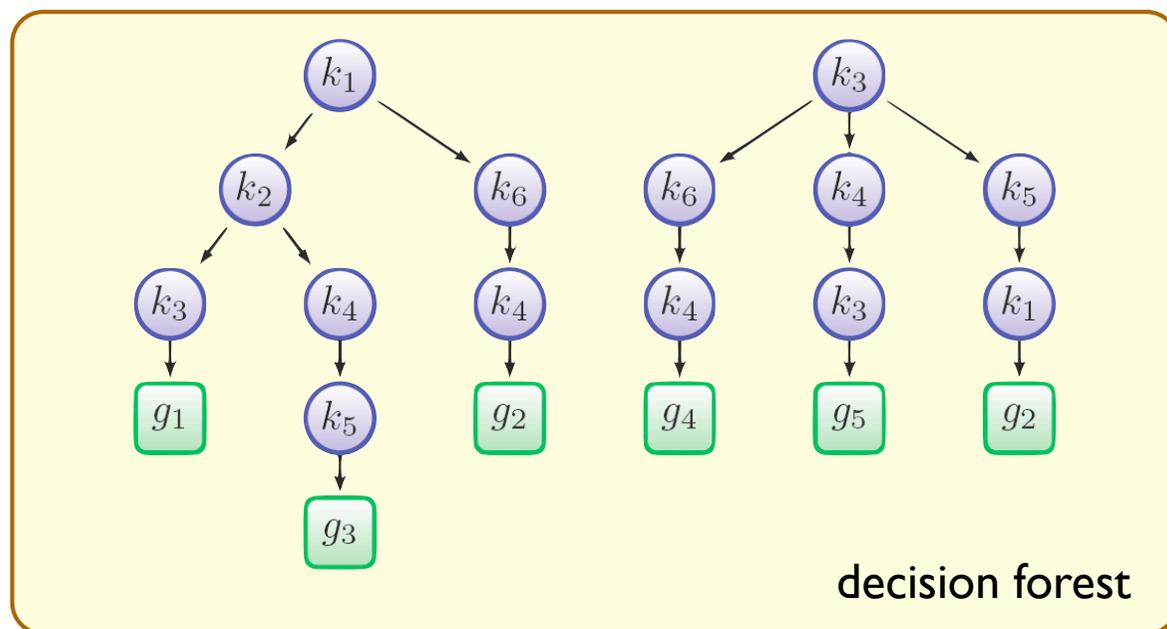
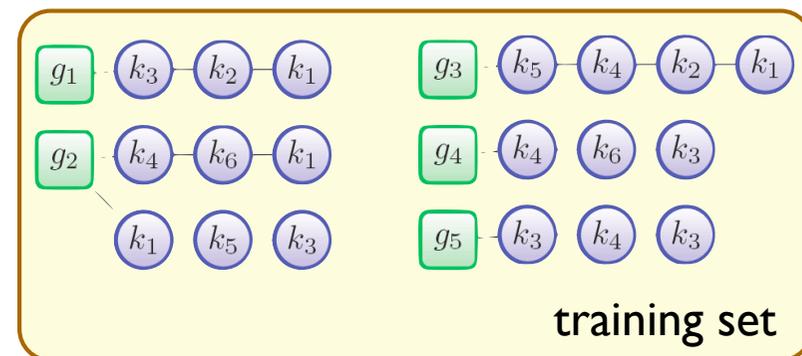
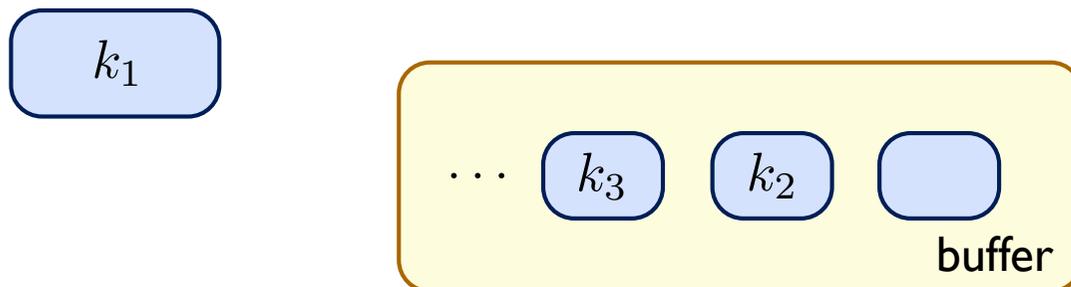
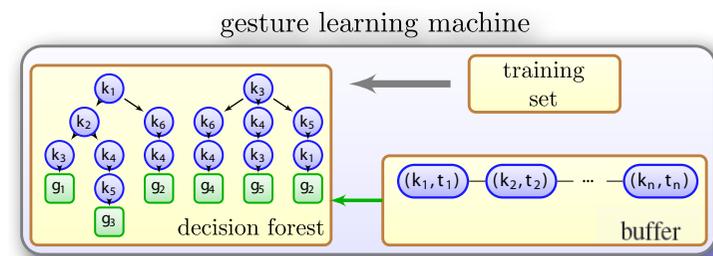
# Real-time gesture recognition



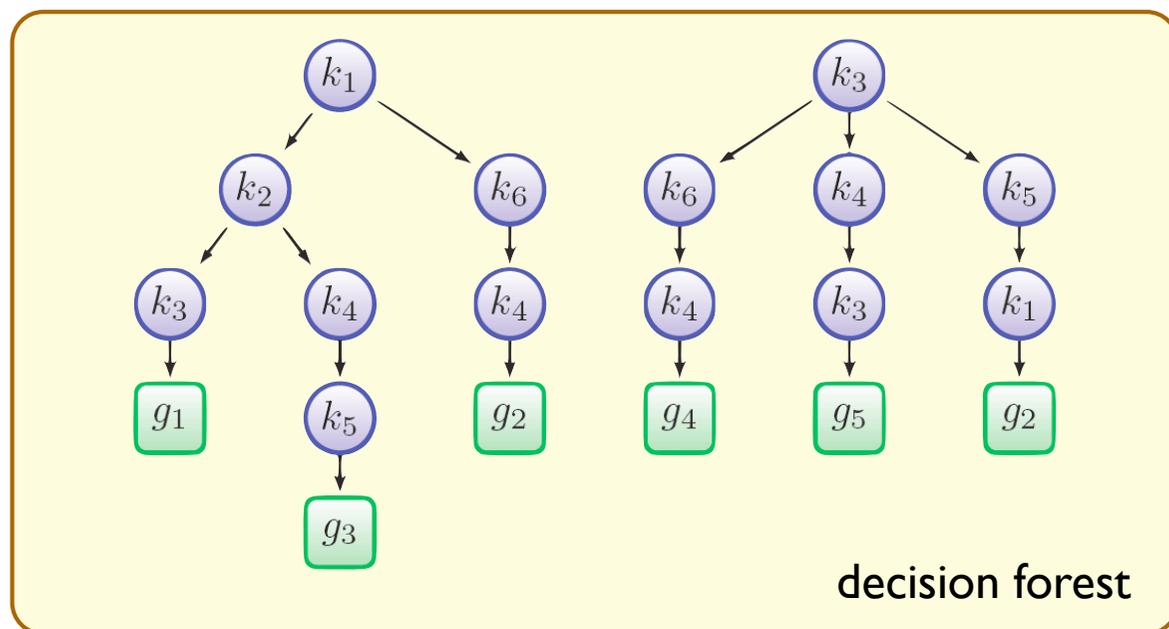
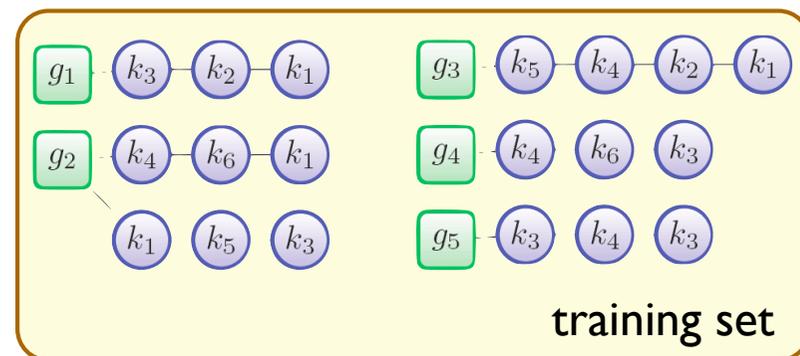
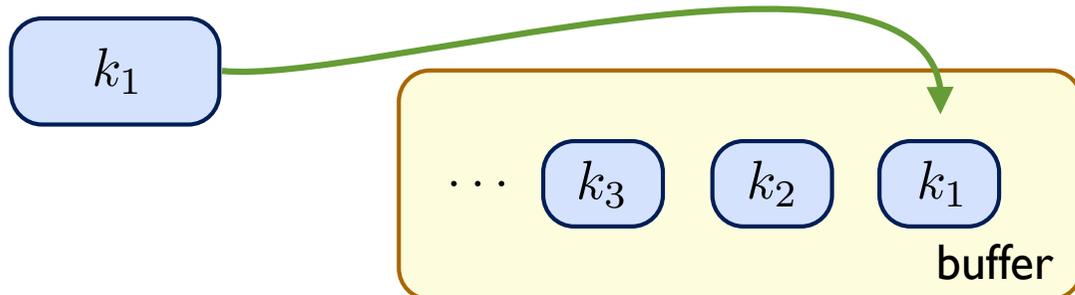
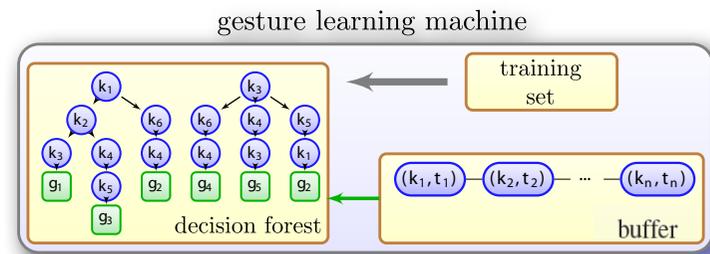
# Real-time gesture recognition



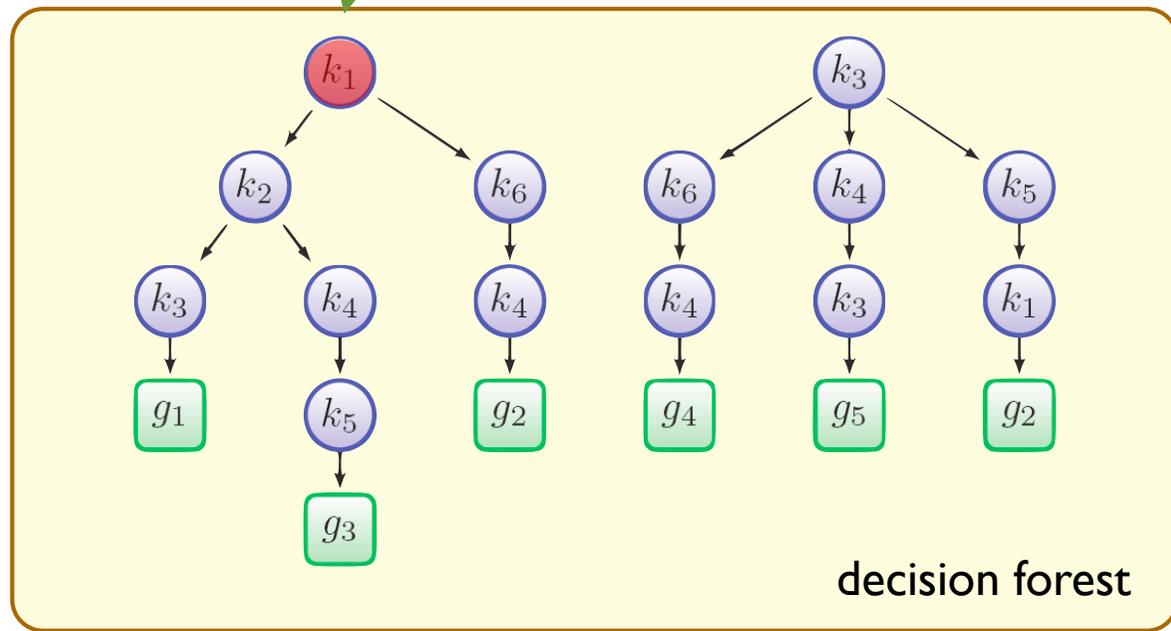
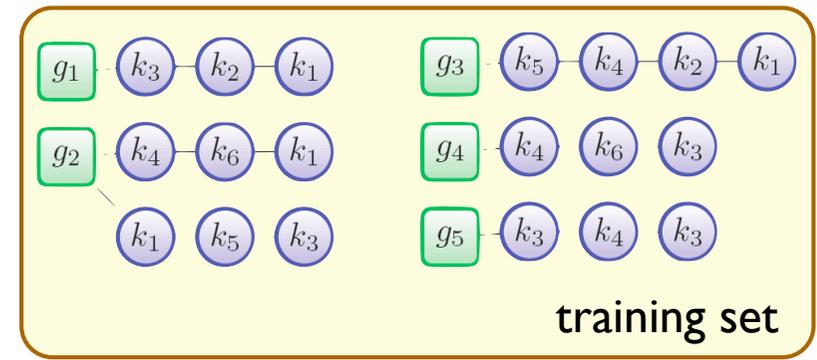
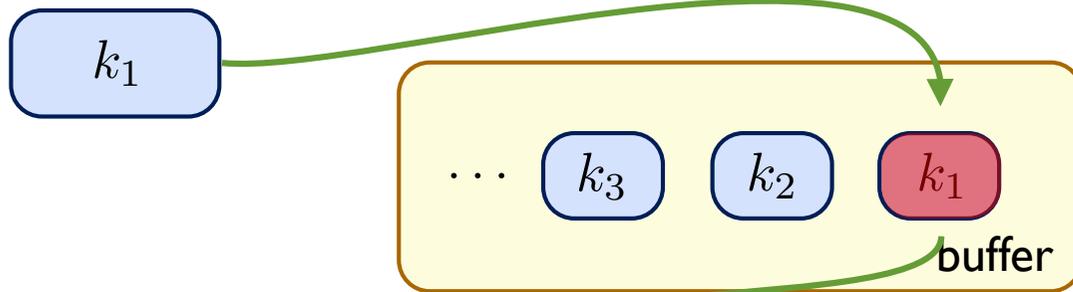
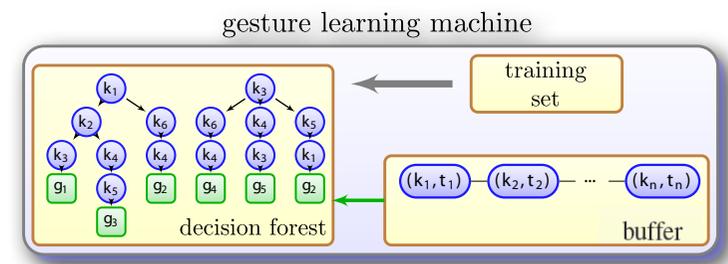
# Real-time gesture recognition: Example



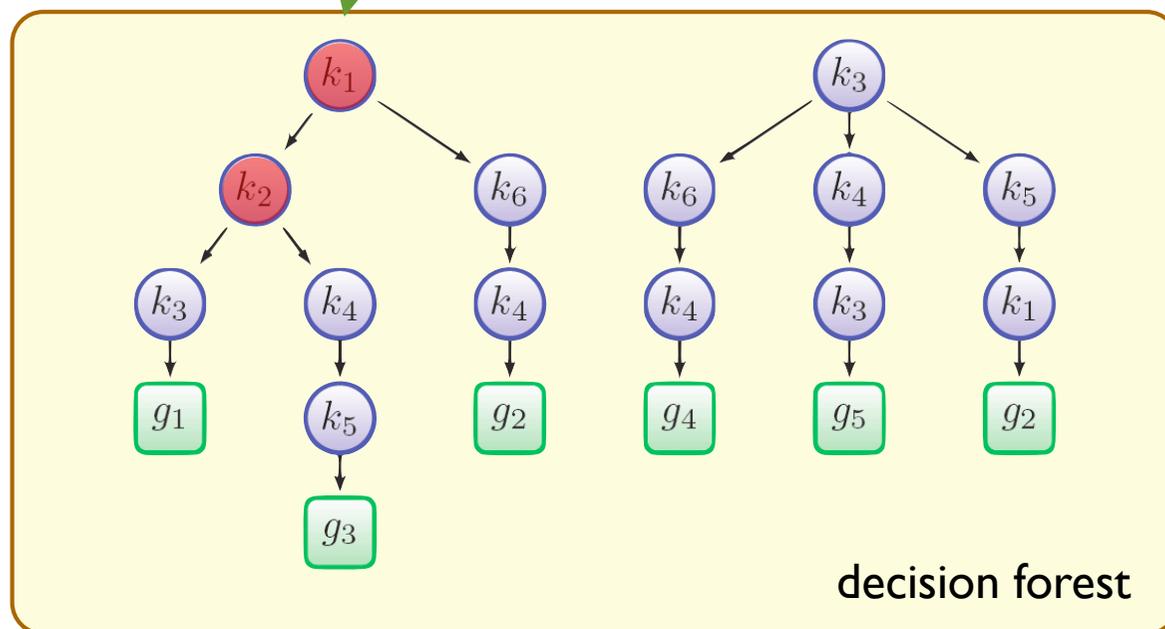
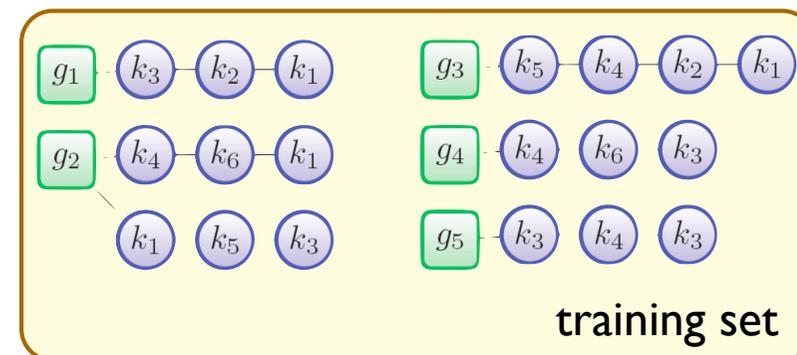
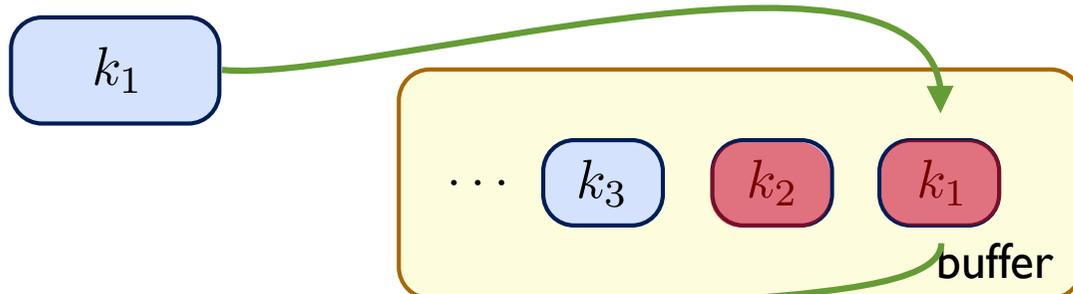
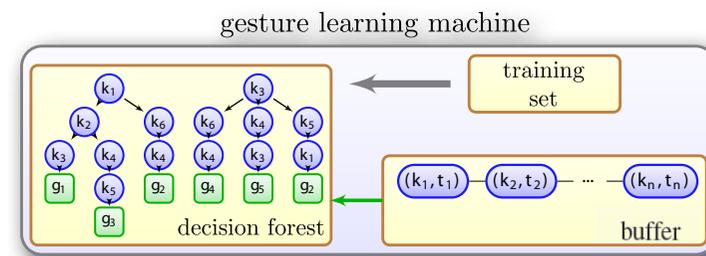
# Real-time gesture recognition: Example



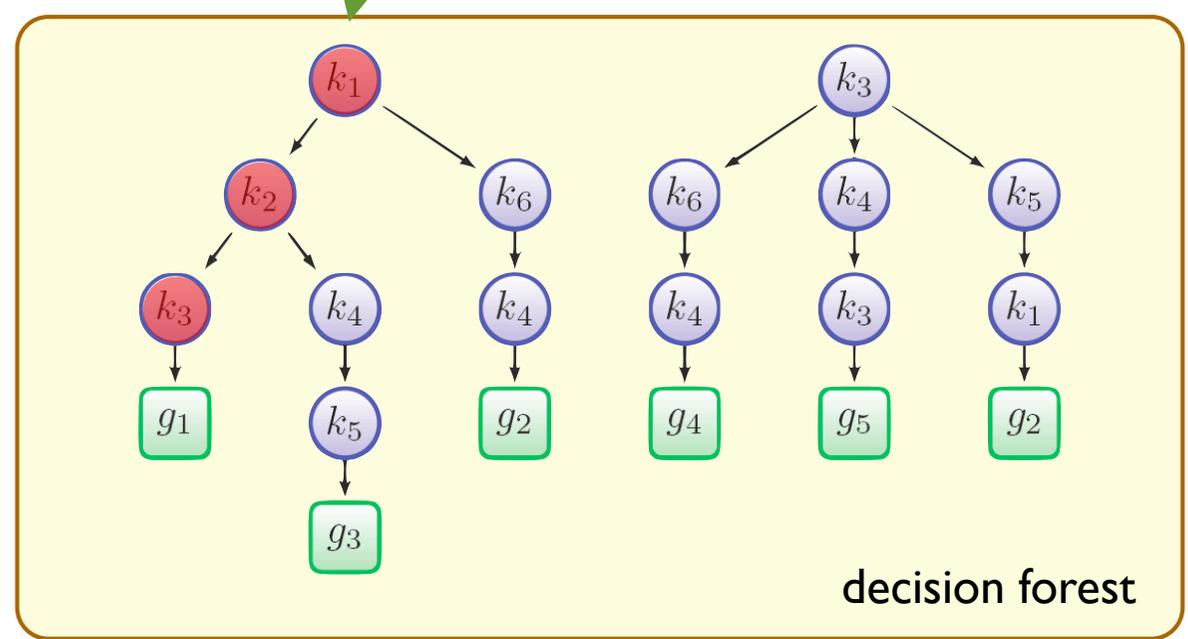
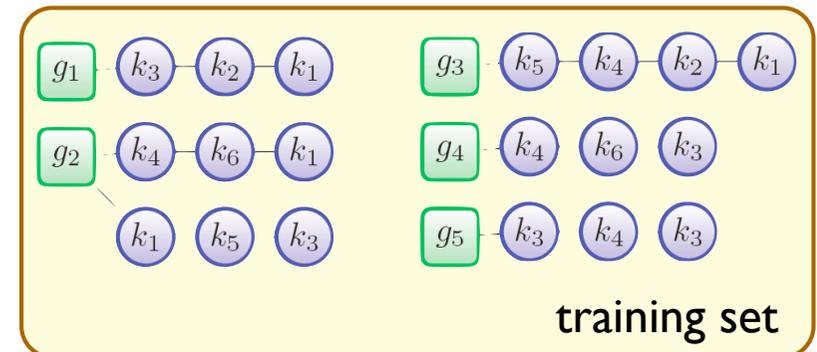
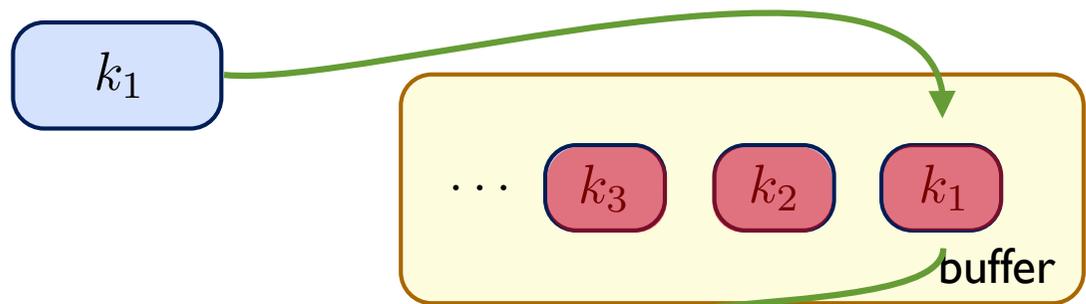
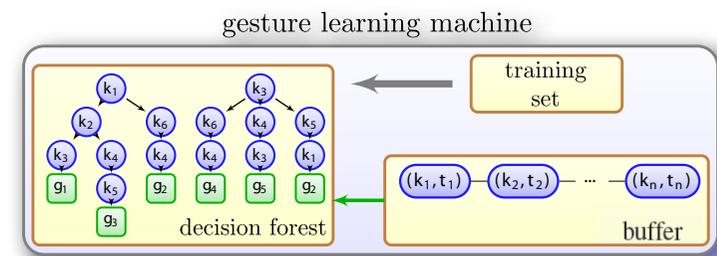
# Real-time gesture recognition: Example



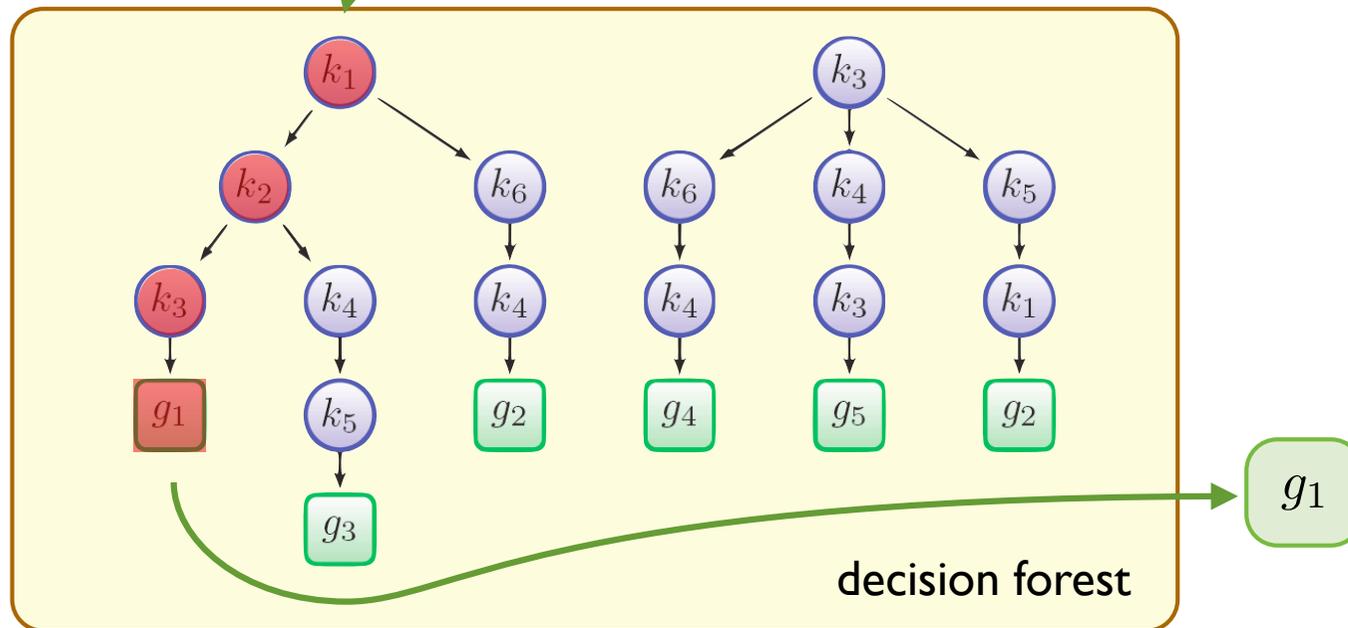
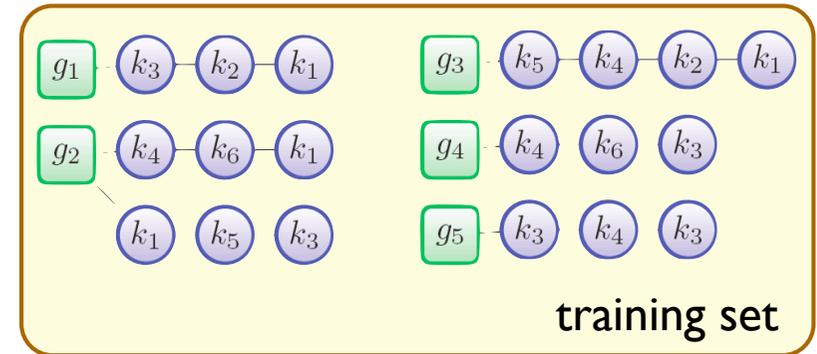
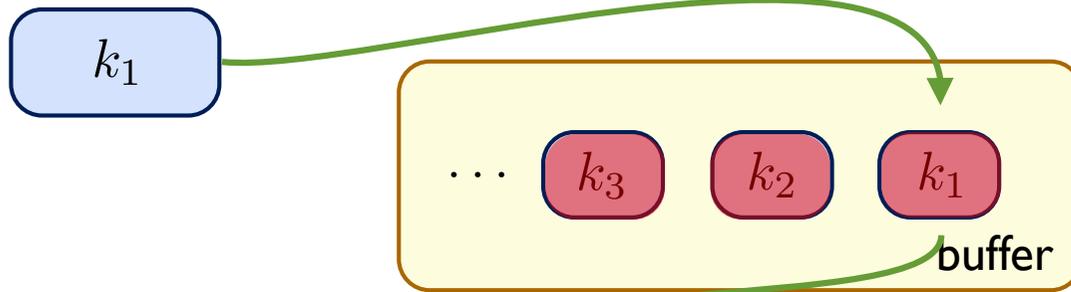
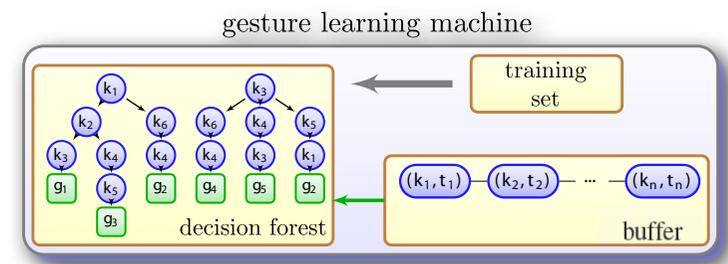
# Real-time gesture recognition: Example



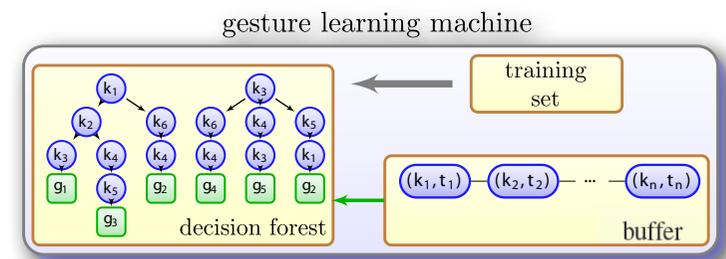
# Real-time gesture recognition: Example



# Real-time gesture recognition: Example



# Time constraints



Time vector: interval  $\mathbf{t} = [t_1, t_2, \dots, t_{n-1}]$   
between consecutive key poses

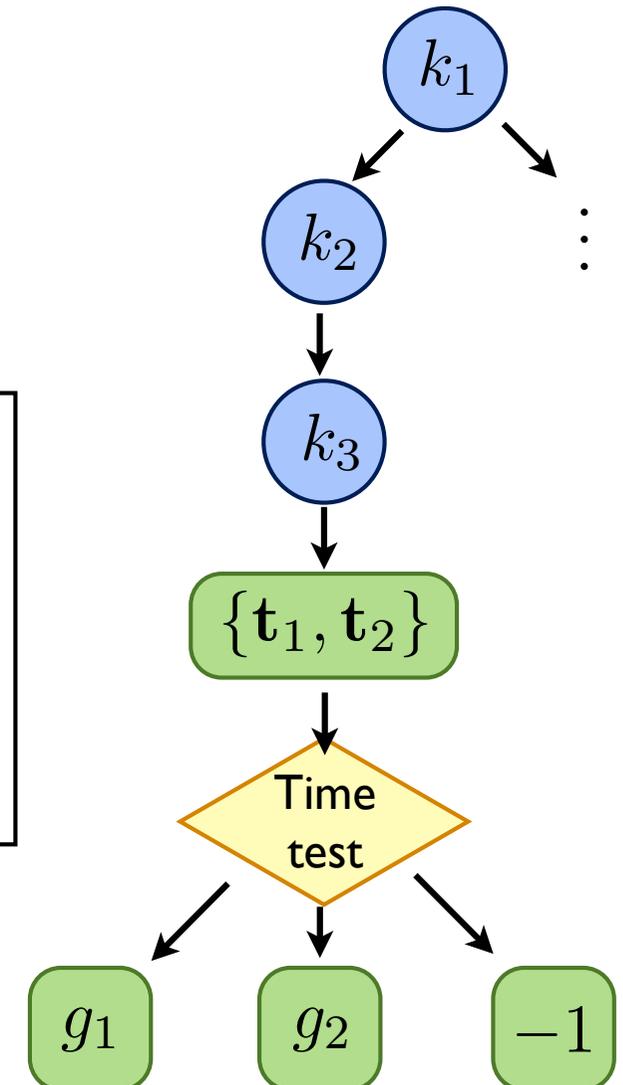
Time test

for each time vector  $t_i$  found on the leaf

if  $\|t_i - \mathbf{t}\|_\infty > T$

discard  $t_i$

return  $g_i$  that minimizes  $\|t_i - \mathbf{t}\|_1$



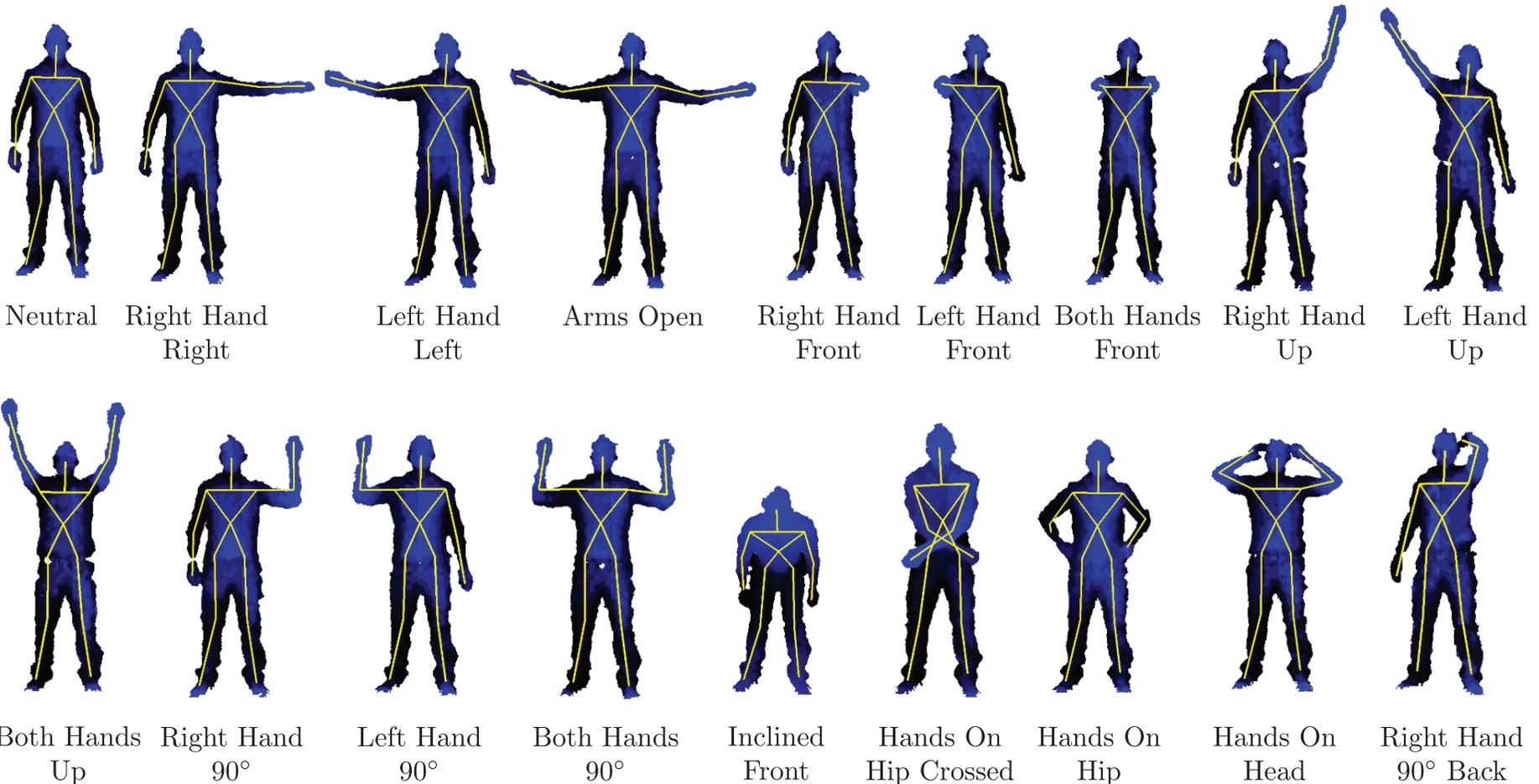
# Results

# Experiment Setup

One trainer

18 trained key poses (approx. 30 examples per key pose)

10 trained gestures (approx. 10 executions per gesture)



# Key pose recognition: robustness

10 inexperienced individuals performed trained key poses 10 times

key pose	id	recognized key poses per user											total (%)	
		$u_1$	$u_2$	$u_3$	$u_4$	$u_5$	$u_6$	$u_7$	$u_8$	$u_9$	$u_{10}^1$	$u_{10}^2$		
Neutral	$k_1$	10	10	10	10	10	10	10	10	10	10	10	10	<b>100.00</b>
Right Hand Right	$k_2$	10	10	10	10	10	10	10	10	10	10	10	8	<b>98.18</b>
Left Hand Left	$k_3$	10	10	10	9	10	10	9	10	10	10	10	10	<b>98.18</b>
Arms Open	$k_4$	10	10	10	7	10	10	10	9	10	7	10	10	<b>93.63</b>
Right Hand Front	$k_5$	10	10	10	10	10	10	10	10	10	8	7	10	<b>95.45</b>
Left Hand Front	$k_6$	10	10	9	10	10	10	10	10	10	10	10	10	<b>99.09</b>
Both Hands Front	$k_7$	10	10	10	10	10	10	10	10	10	10	10	10	<b>100.00</b>
Right Hand Up	$k_8$	10	10	10	10	10	10	10	10	10	10	10	10	<b>100.00</b>
Left Hand Up	$k_9$	10	10	10	10	10	9	10	10	10	9	10	10	<b>98.18</b>
Both Hands Up	$k_{10}$	10	10	10	10	10	10	10	10	10	10	10	10	<b>100.00</b>
Right Hand 90°	$k_{11}$	10	8	9	10	10	10	10	10	8	10	10	10	<b>95.45</b>
Left Hand 90°	$k_{12}$	10	10	10	10	10	6	10	10	10	5	10	10	<b>91.81</b>
Both Hands 90°	$k_{13}$	10	10	10	10	10	10	10	10	10	10	10	10	<b>100.00</b>
Inclined Front	$k_{14}$	8	10	10	10	10	8	10	10	10	5	7	10	<b>89.09</b>
Hands-on-Hip Crossed	$k_{15}$	7	8	6	8	8	10	10	10	8	10	8	10	<b>84.54</b>
Hand-On-Hip	$k_{16}$	10	10	10	10	10	10	10	9	10	10	10	10	<b>99.09</b>
Hands on Head	$k_{17}$	9	10	10	8	10	10	9	7	10	10	6	10	<b>90.00</b>
Right Hand 90° Back	$k_{18}$	8	10	9	6	7	7	7	10	10	3	8	10	<b>77.27</b>
<b>total (%)</b>		<b>95.5</b>	<b>97.7</b>	<b>96.1</b>	<b>93.3</b>	<b>97.2</b>	<b>94.4</b>	<b>97.2</b>	<b>97.2</b>	<b>97.7</b>	<b>87.2</b>	<b>91.11</b>		

Average recognition rate: 94.84%

# Key pose recognition: stability

*Out-of-sample tests:*

1. Remove 20% of training set data;
2. Compute SVM classifier;
3. Try to classify removed training data.

Results after 10 experiments:

False classifications: 4.16%

Unclassified key poses: 3.45%

# Key pose recognition



# Gesture recognition

10 inexperienced individuals performed trained gestures 10 times

<b>gesture</b>	<b>id</b>	<b>key pose seq.</b>	<b>rec. rate</b>
Open-Clap	$g_1$	$k_1, k_4, k_7$	<b>99%</b>
Open Arms	$g_2$	$k_1, k_7, k_4$	<b>96%</b>
Turn Next Page	$g_3$	$k_1, k_2, k_5, k_1$ $k_1, k_6, k_3, k_1$	<b>83%</b>
Turn Previous Page	$g_4$	$k_1, k_5, k_2, k_1$ $k_1, k_3, k_6, k_1$	<b>91%</b>
Raise Right Arm Laterally	$g_5$	$k_1, k_2, k_8$	<b>80%</b>
Lower Right Arm Laterally	$g_6$	$k_8, k_2, k_1$	<b>78%</b>
Good Bye ( $k_{11}$ time constraint: 1sec.)	$g_7$	$k_1, k_{11}$	<b>92%</b>
Japanese Greeting	$g_8$	$k_1, k_{14}, k_1$	<b>100%</b>
Put Hands Up Front	$g_9$	$k_1, k_5, k_{18}$ $k_1, k_5, k_8$ $k_1, k_5, k_{11}, k_8$ $k_1, k_8$	<b>96%</b>
Put Hands Up Laterally	$g_{10}$	$k_1, k_4, k_{10}$	<b>100%</b>

# Gesture recognition

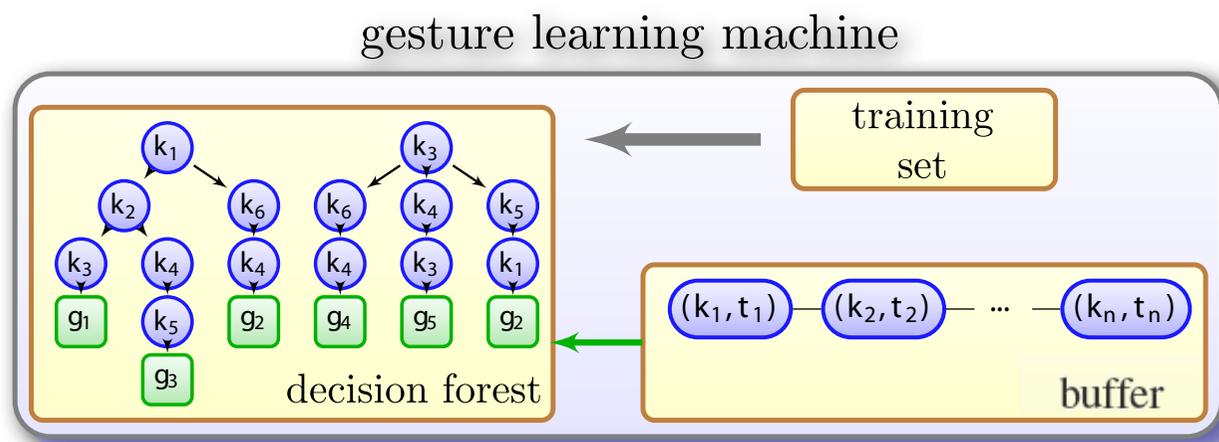


# Performance

Preprocessing bottleneck: computing SVM classifiers

For a training set of 2,000 key pose examples of 18 classes:  
18 functions were computed in 3.9 secs

Negligible performance during training/recognition phases



Usually very low tree depths

# Comparison

Dataset from Li *et al* (2010): 20 gestures, 10 individuals, 3 executions

AS1	AS2	AS3
Horizontal arm wave	High arm wave	High throw
Hammer	Hand catch	Forward kick
Forward punch	Draw x	Side kick
High throw	Draw tick	Jogging
Hand clap	Draw circle	Tennis swing
Bend	Two hand wave	Tennis serve
Pickup & throw	Side boxing	Pickup & throw

*Cross-subject test:*

Gesture subset	Li [10]	Vieira [15]	our method
AS1	72.9%	84.7%	<b>93.5%</b>
AS2	71.9%	81.3%	<b>52.0%</b>
AS3	79.2%	88.4%	<b>95.4%</b>
Average	74.7%	84.8%	<b>80.3%</b>

# Comparison

Dataset from Li *et al* (2010): 20 gestures, 10 individuals, 3 executions

AS1	AS2	AS3
Horizontal arm wave	High arm wave	High throw
Hammer	Hand catch	Forward kick
Forward punch	Draw x	Side kick
High throw	Draw tick	Jogging
Hand clap	Draw circle	Tennis swing
Bend	Two hand wave	Tennis serve
Pickup & throw	Side boxing	Pickup & throw

*Cross-subject test:*

Gesture subset	Li [10]	Vieira [15]	our method
AS1	72.9%	84.7%	93.5%
AS2	71.9%	81.3%	52.0%
AS3	79.2%	88.4%	95.4%
Average	74.7%	84.8%	80.3%

**Delicate gestures** →

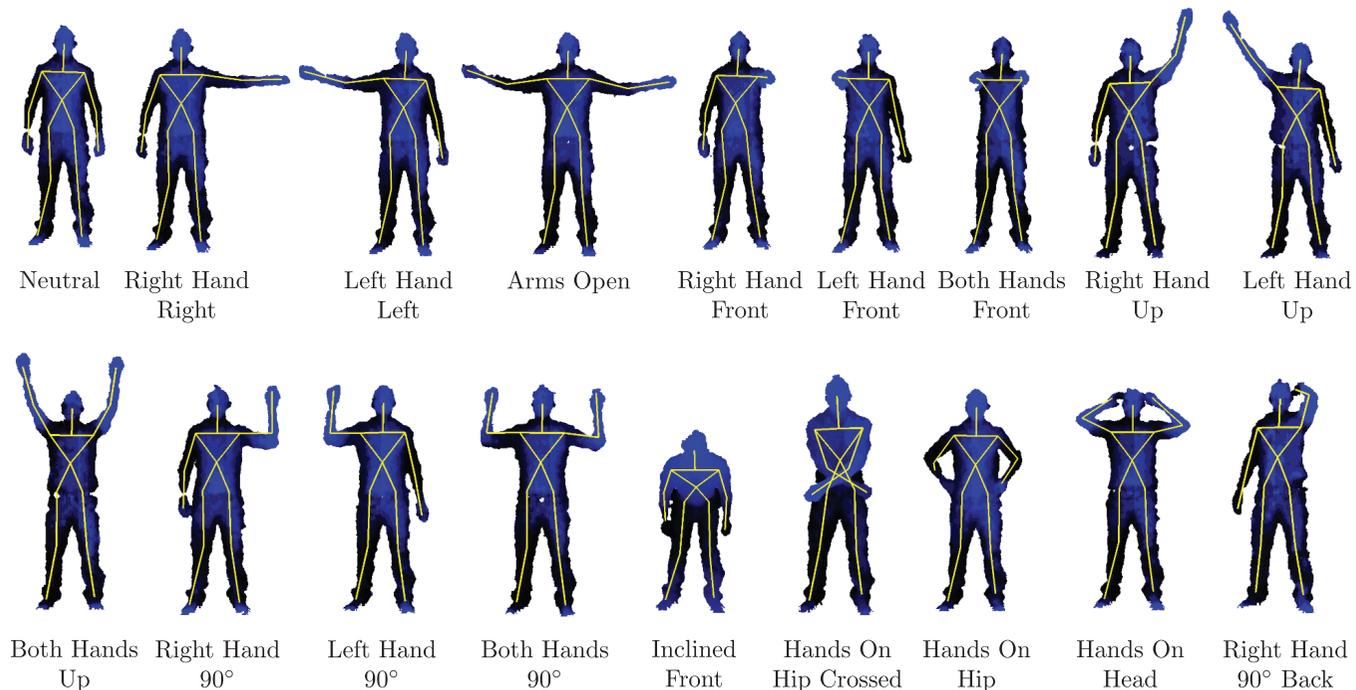
# Limitations

- Robustness issues

Skeleton tracking

Delicate gestures

- Key pose design not the friendliest solution



# Future Work

- ✓ Automatic key pose generation
- ✓ Work on skeleton tracking algorithms (More than 1 Kinect?)
- ✓ Improve time constrained gesture recognition
- ✓ Take into account key pose descriptor periodicity

# Thank you for your attention!

# Thank you for your attention!

# Thank you for your attention!

## Questions?