# People-Tracking-by-Detection and People-Detection-by-Tracking

Mykhaylo Andriluka     Stefan Roth     Bernt Schiele

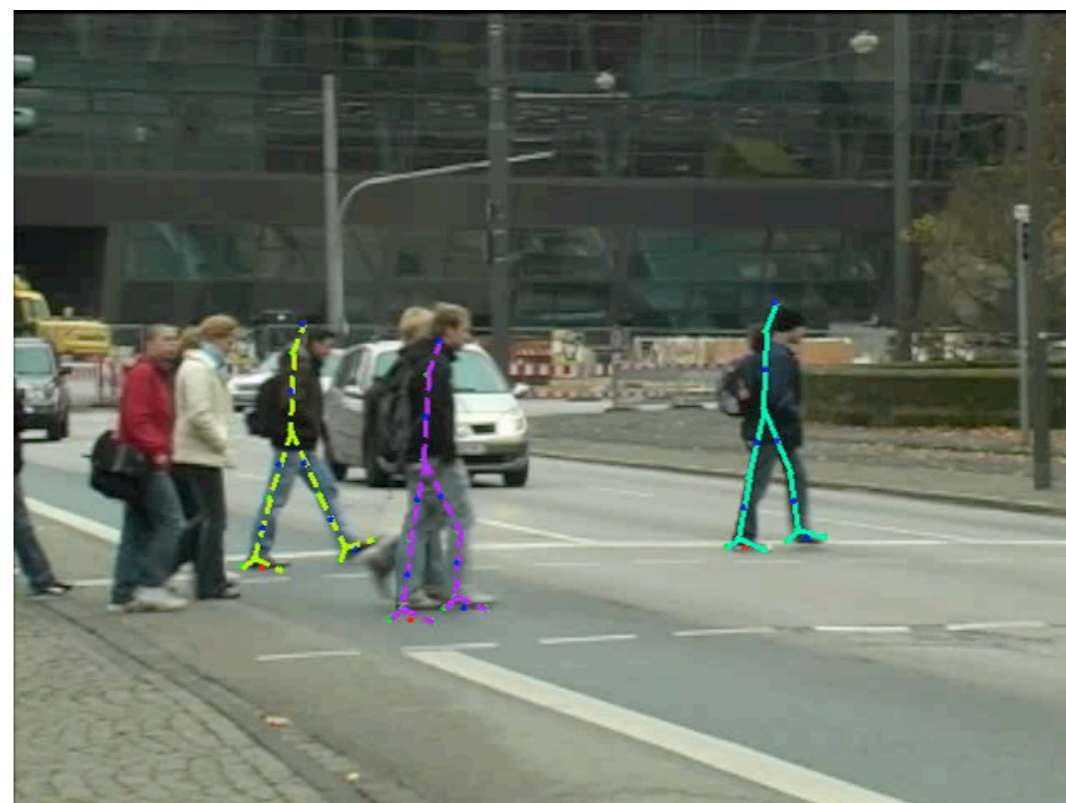Department of Computer Science
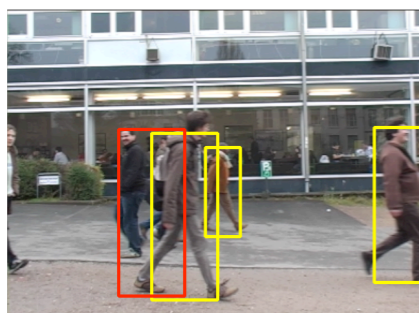TU Darmstadt

# Motivation

- Goal: Detection and tracking of people in complex scenes

- Challenges for detection:
  - ▶ Partial occlusions
  - ▶ Appearance variation
  - ▶ Data association difficult



- Challenges for tracking:
  - ▶ Dynamic backgrounds
  - ▶ Multiple people
  - ▶ Frequent long term occlusions

# **Motivation**
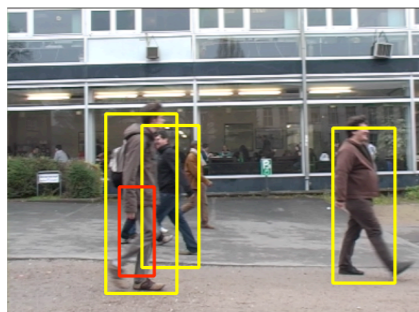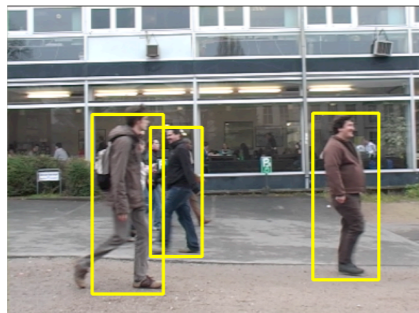
- Goal: Detection and tracking of people in complex scenes

- Challenges for detection:

  ▶ Partial occlusions

  ▶ Appearance variation

  ▶ Data association difficult

- Challenges for tracking:

  ▶ Dynamic backgrounds

  ▶ Multiple people

  ▶ Frequent long term occlusions

# Overview

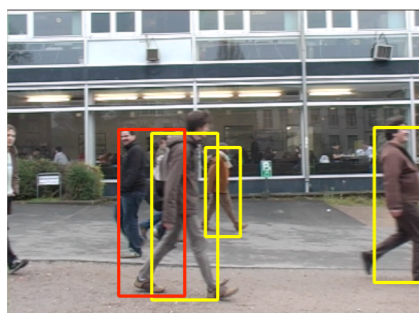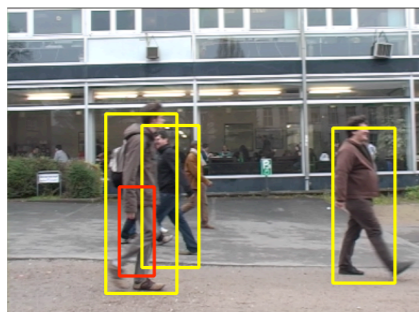Three stages of our multi-person detection and tracking system:
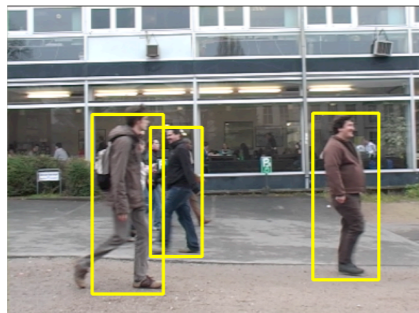
1. Single-frame detection

# Overview

Three stages of our multi-person detection and tracking system:

1. Single-frame detection

2. Tracklet detection

# Overview

Three stages of our multi-person detection and tracking system:
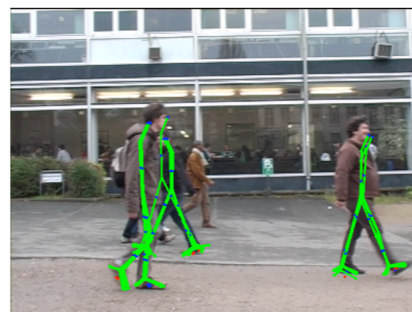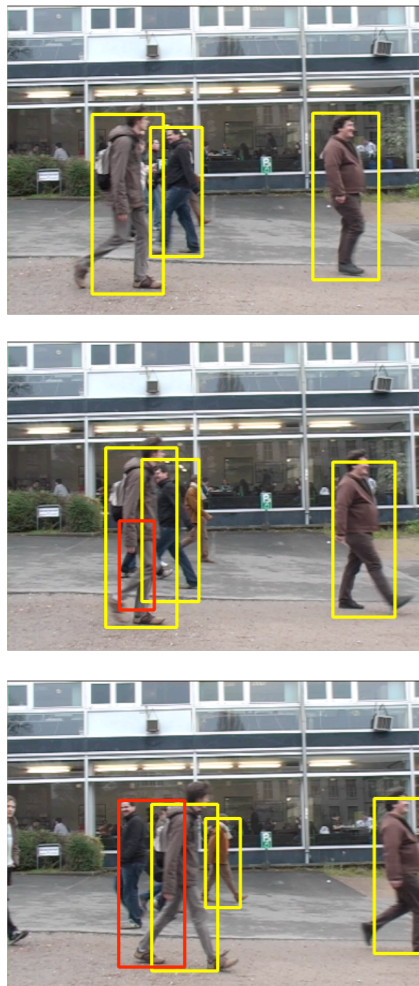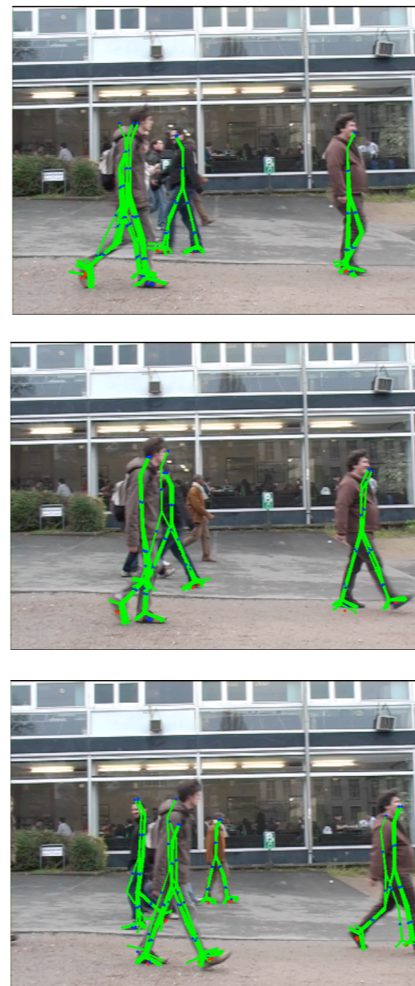


1. Single-frame detection

2. Tracklet detection
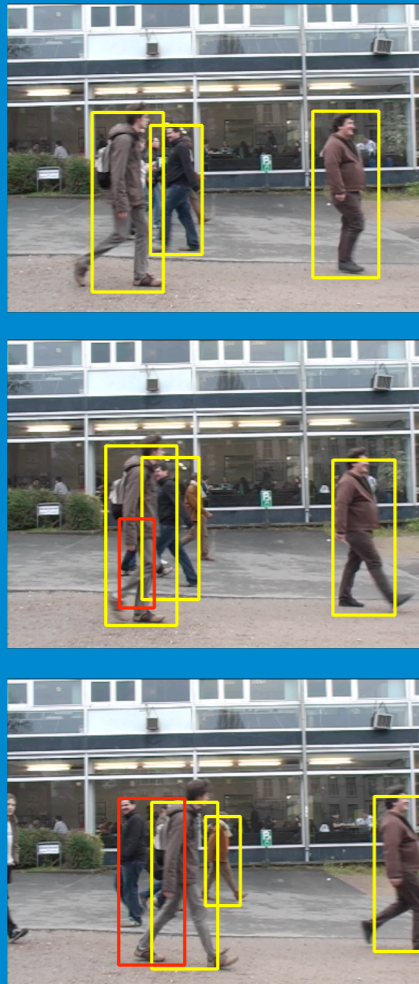
3. Tracking through occlusion

# Previous Work

- **People Detection & Tracking:**

  ▸ [Fossati et al., CVPR 2007]: 3D articulated tracking aided by detection, single person, ground plane needed.

  ▸ [Leibe et al., ICCV 2007]: Detection of tracking of multiple people, high viewpoint → no full-body occlusions.

  ▸ [Ramanan et al., PAMI 2007]: Appearance model learned from people detection, then used for tracking and data association.

  ▸ [Wu & Nevatia, IJCV 2007]: Use detection for tracking, works for multiple people → no articulations, detector not aided by tracking.

- **Here:**

  ▸ More people

  ▸ Significant, long-term full-body occlusions

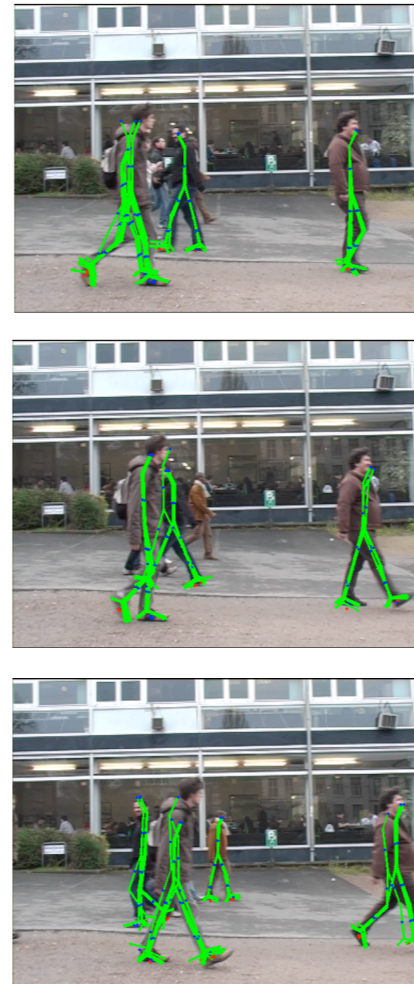  ▸ However: more restricted scenario (2-D, people in side views)

# Overview

Three stages of our multi-person detection and tracking system:



1. Single-frame detection

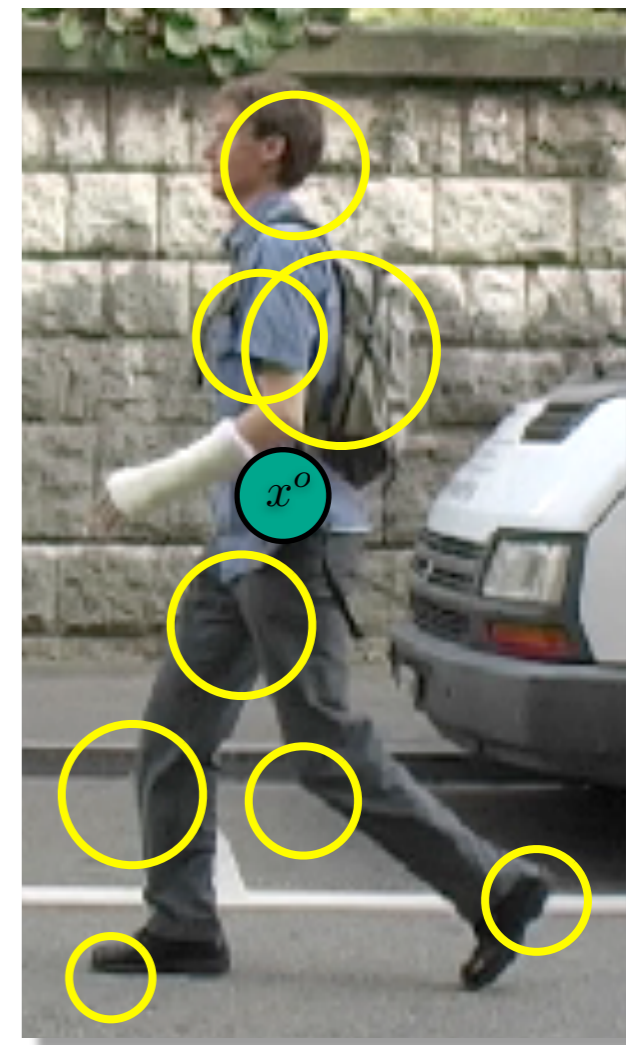2. Tracklet detection

3. Tracking through occlusion

# Single-frame Detector: partISM

- **Appearance of parts:**
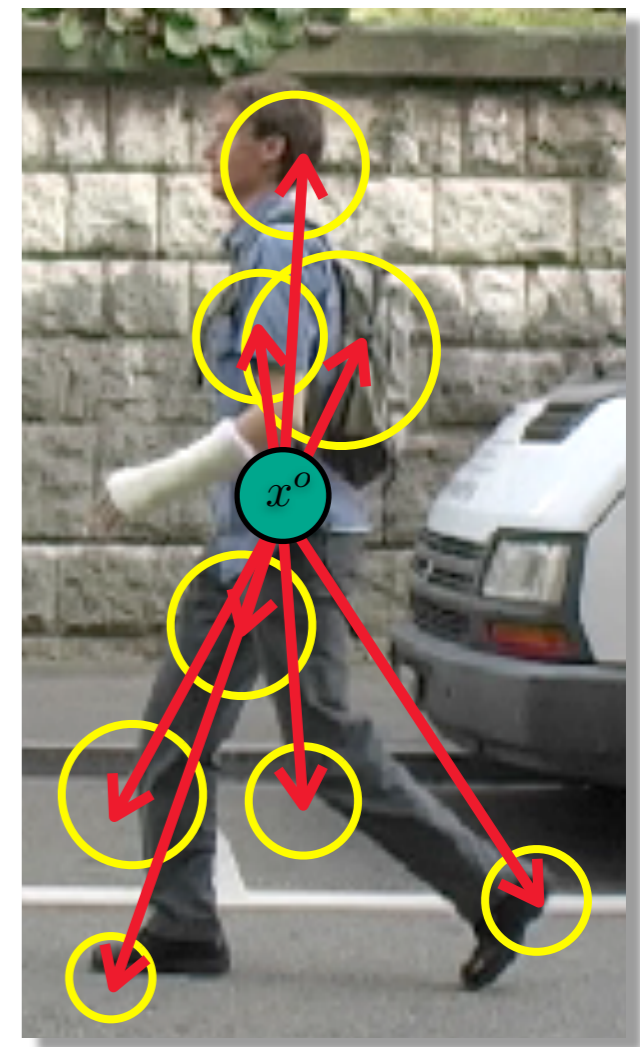  Implicit Shape Model (ISM)
  [Leibe, Seemann & Schiele, CVPR 2005]

# Single-frame Detector: partISM

- **Appearance of parts:**
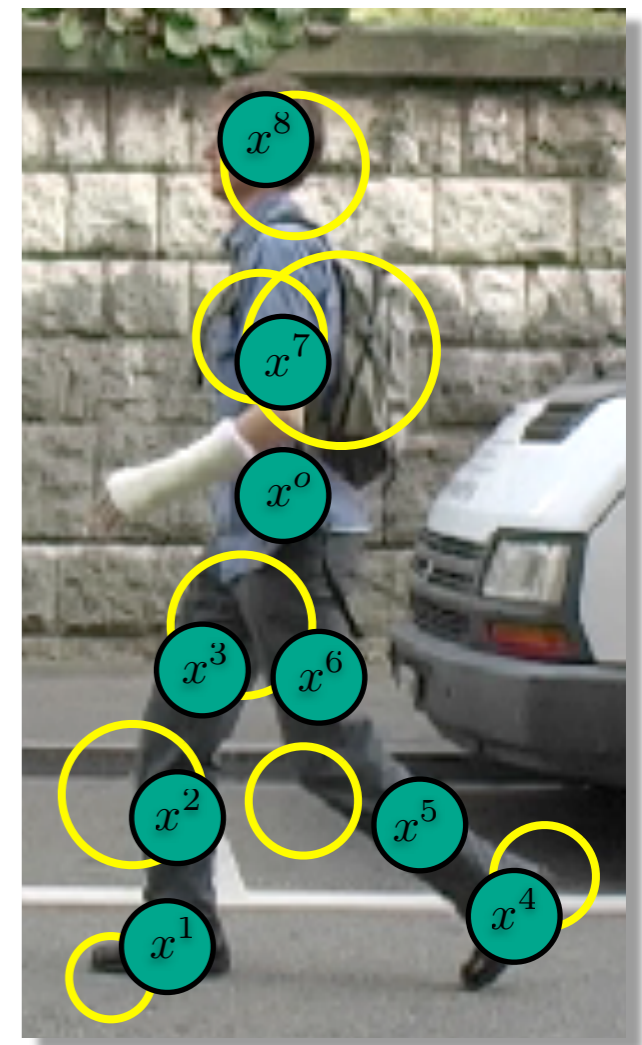  Implicit Shape Model (ISM)
  [Leibe, Seemann & Schiele, CVPR 2005]

# Single-frame Detector: partISM

- **Appearance of parts:**
  Implicit Shape Model (ISM)
  [Leibe, Seemann & Schiele, CVPR 2005]

# Single-frame Detector: partISM

- **Appearance of parts:**
  Implicit Shape Model (ISM)
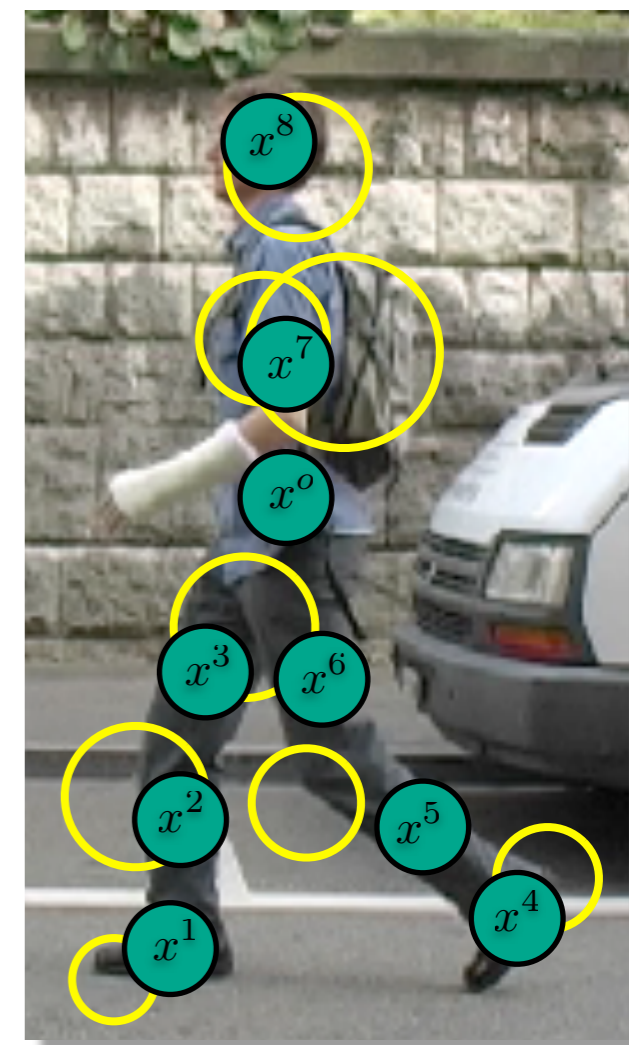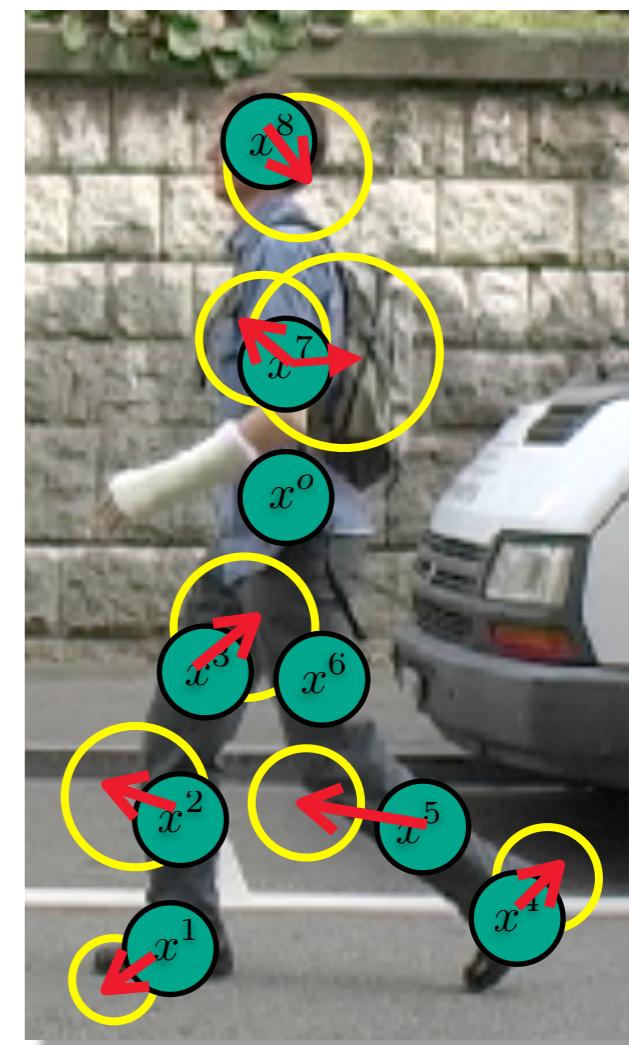  [Leibe, Seemann & Schiele, CVPR 2005]

# Single-frame Detector: partISM

- **Appearance of parts:**
  Implicit Shape Model (ISM)
  [Leibe, Seemann & Schiele, CVPR 2005]

- **Part decomposition and inference:**
  Pictorial structures model
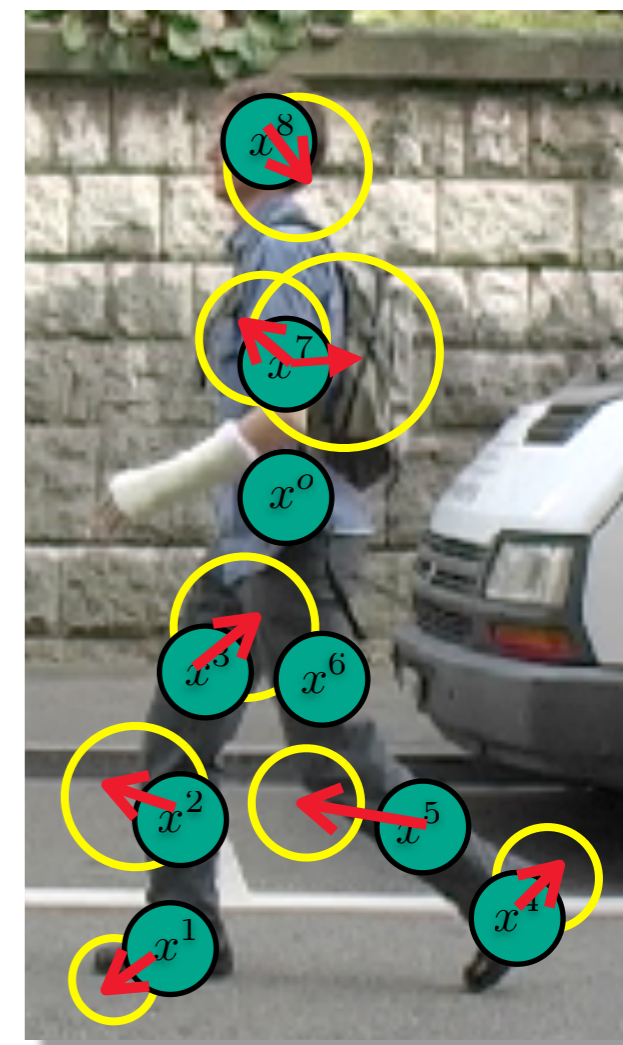  [Felzenszwalb & Huttenlocher, IJCV 2005]

# Single-frame Detector: partISM

- **Appearance of parts:**
  Implicit Shape Model (ISM)
  [Leibe, Seemann & Schiele, CVPR 2005]

- **Part decomposition and inference:**
  Pictorial structures model
  [Felzenszwalb & Huttenlocher, IJCV 2005]

# Single-frame Detector: partISM

- **Appearance of parts:**
  Implicit Shape Model (ISM)
  [Leibe, Seemann & Schiele, CVPR 2005]

- **Part decomposition and inference:**
  Pictorial structures model
  [Felzenszwalb & Huttenlocher, IJCV 2005]
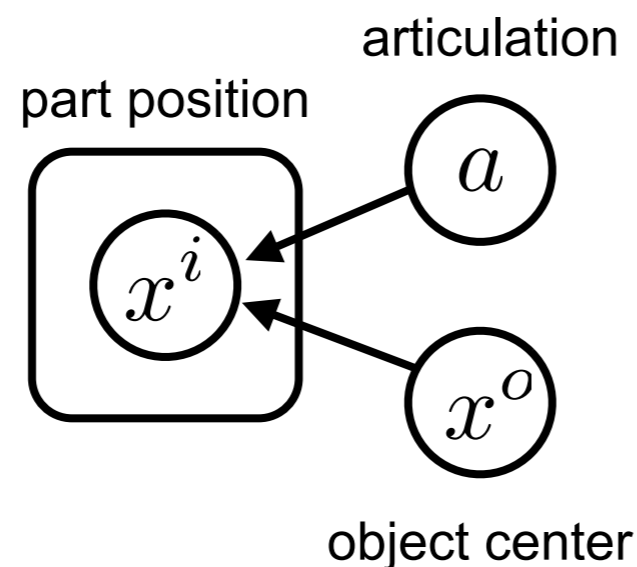
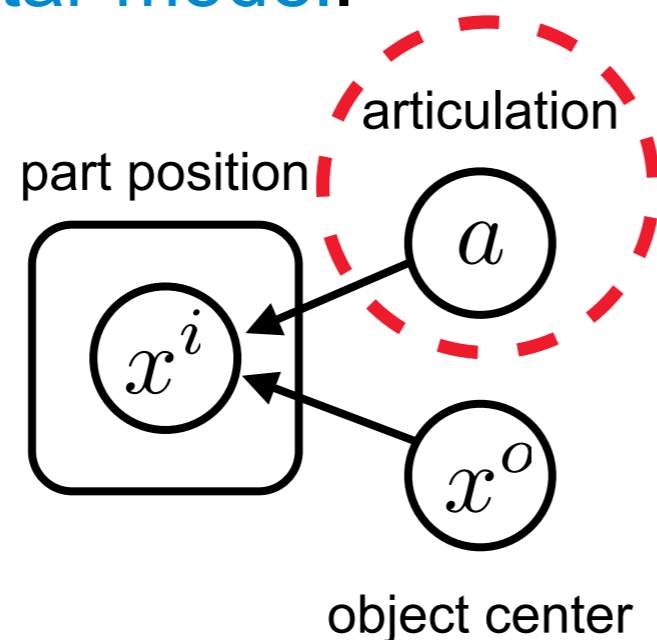$$p(L|E) \propto p(E|L)p(L)$$

Body-part positions    Image evidence

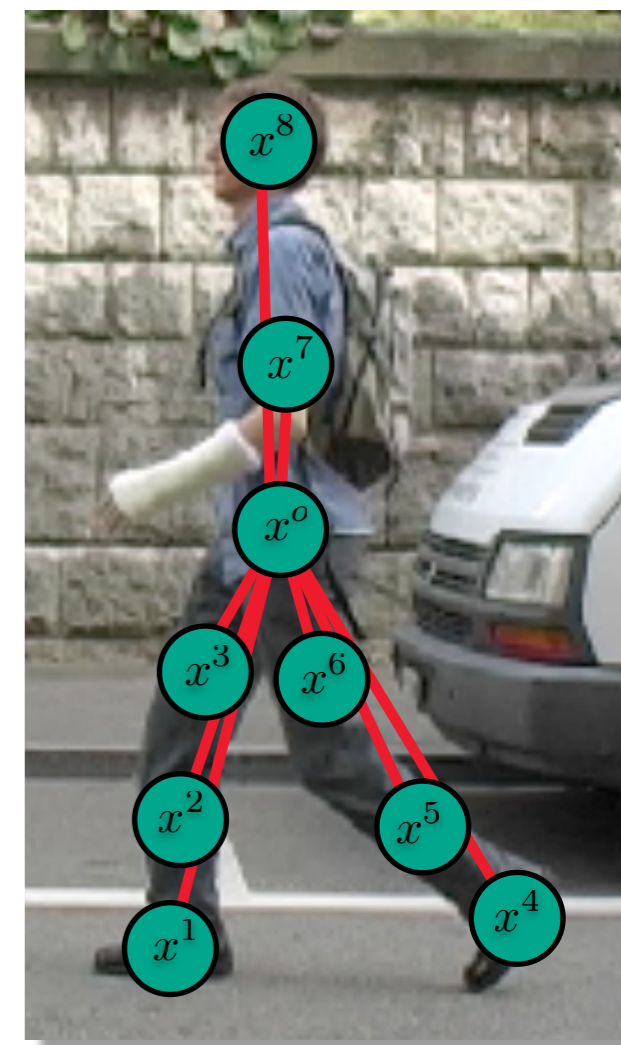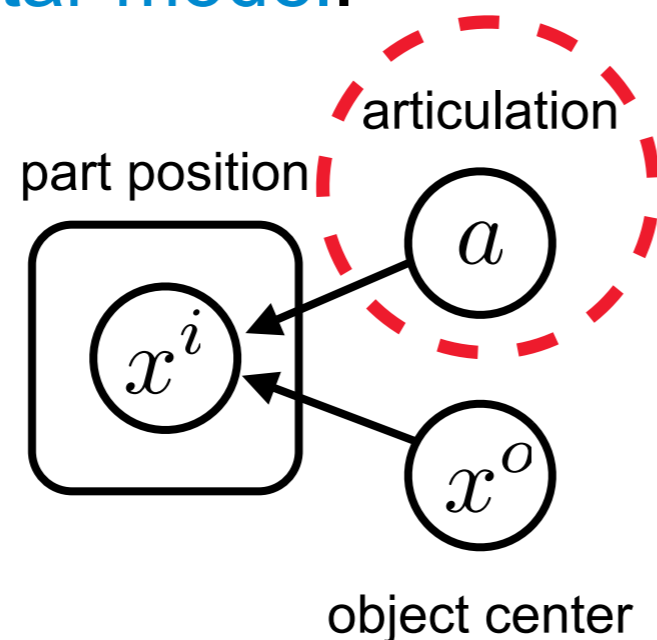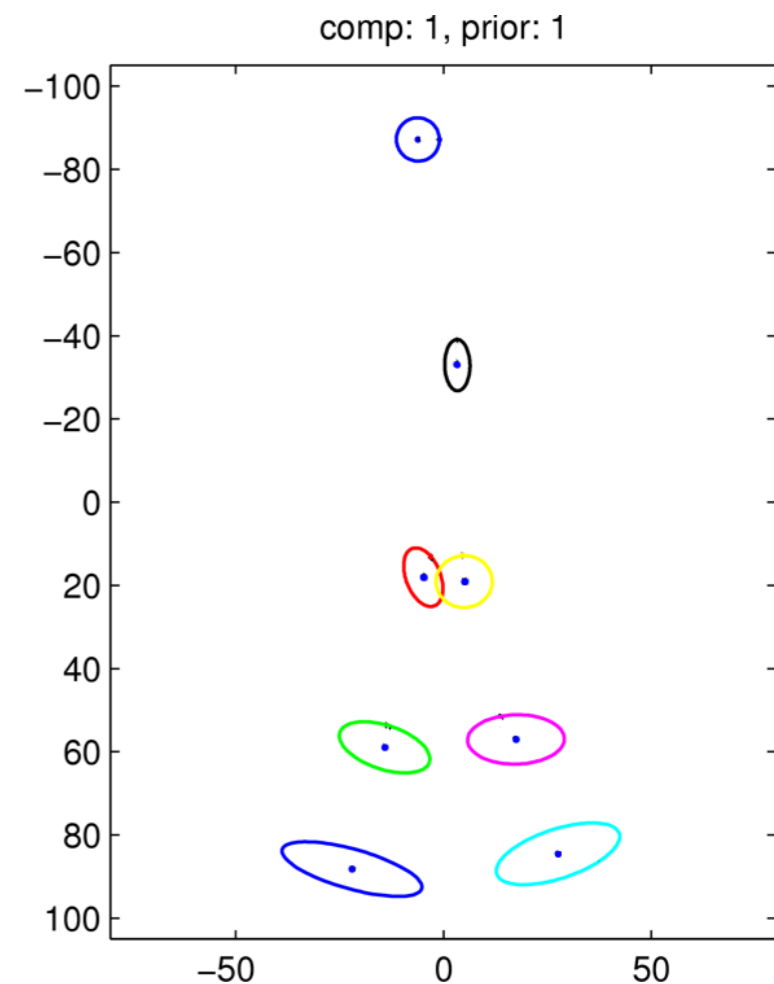# Part Decomposition

- $L = \{x^o, x^1, \ldots, x^8\}$ - configuration of body parts

- Structure of the prior distribution $p(L)$:

  ▸ Articulation variable $a$ models correlations between part positions.

  ▸ Given articulation, prior on configuration becomes a star model.



articulation

part position

$a$

$x^i$

$x^o$

object center

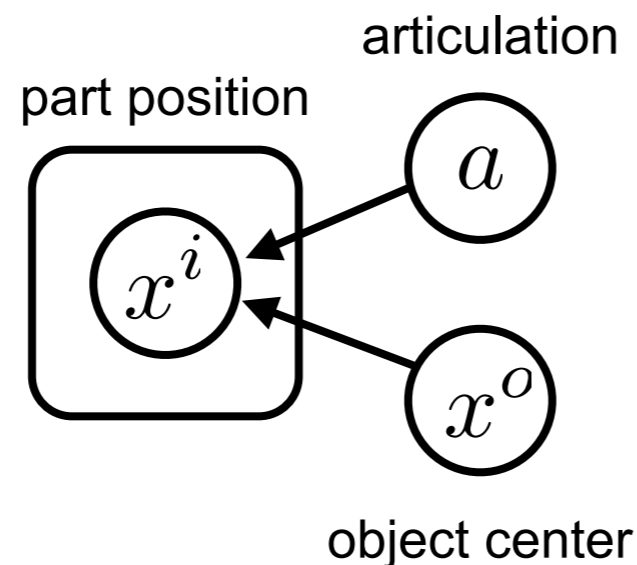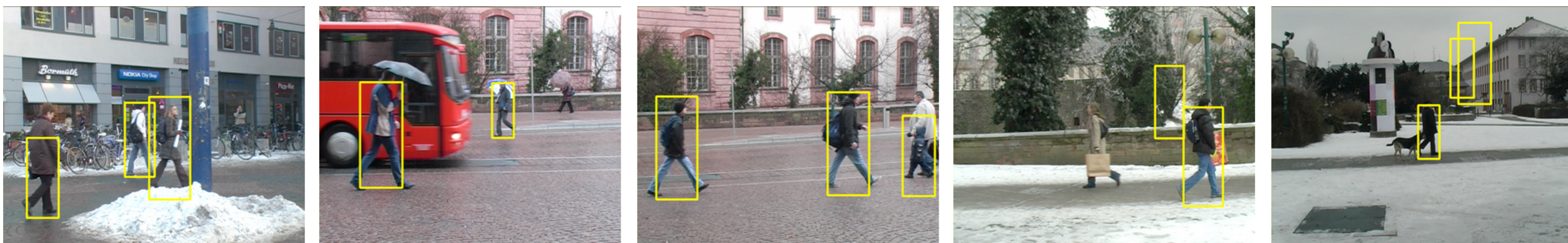# Part Decomposition

- $L = \{x^o, x^1, \ldots, x^8\}$ - configuration of body parts

- Structure of the prior distribution $p(L)$:

  - Articulation variable $a$ models correlations between part positions.

  - Given articulation, prior on configuration becomes a star model.

# Part Decomposition

- $L = \{x^o, x^1, \ldots, x^8\}$ - configuration of body parts

- Structure of the prior distribution $p(L)$:

  ▸ Articulation variable $a$ models correlations between part positions.

  ▸ Given articulation, prior on configuration becomes a star model.

# Part Decomposition

- $L = \{x^o, x^1, \ldots, x^8\}$ - configuration of body parts

- Structure of the prior distribution $p(L)$:

  ▸ Articulation variable $a$ models correlations between part positions.

  ▸ Given articulation, prior on configuration becomes a star model.

comp: 1, prior: 1

Covariance and mean part positions for $p(x^i | x^o)$.

articulation

part position

$a$

$x^i$

$x^o$

object center

# Single Frame Detection

- Detections at equal error rate:

# Single-frame Detection Results



TUD pedestrians data

No occlusions

- partISM clearly outperforms 4D-ISM [Seemann et al, DAGM'06].
- Outperforms HOG [Dalal&Triggs, CVPR'05] with much less training data (Note: we only use sideviews).
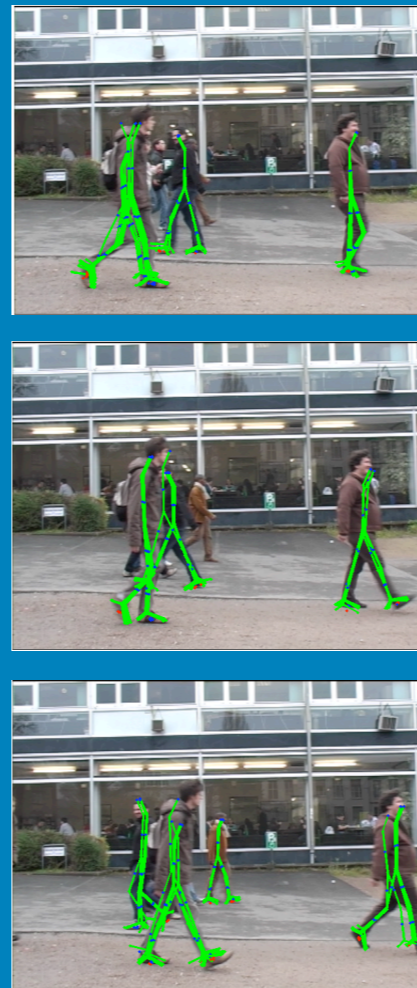
# Overview

TECHNISCHE
UNIVERSITÄT
DARMSTADT

Three stages of our multi-person detection and tracking system:

# Tracklet Detection in Short Subsequences

- Given: $E = [E_1, \dots, E_m]$



... 

overlapping subsequences

- Want:



- Posterior over positions and configurations:

# Tracklet Detection in Short Subsequences

- Given: $E = [E_1, \ldots, E_m]$



frame 1  frame 2  frame m
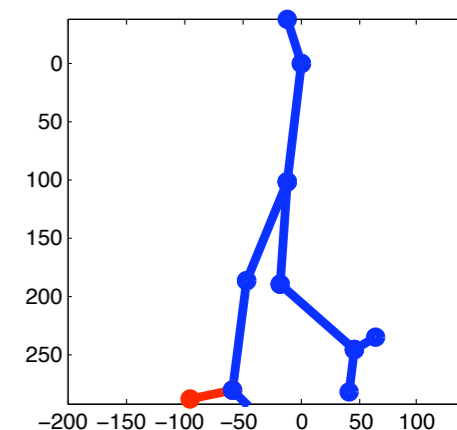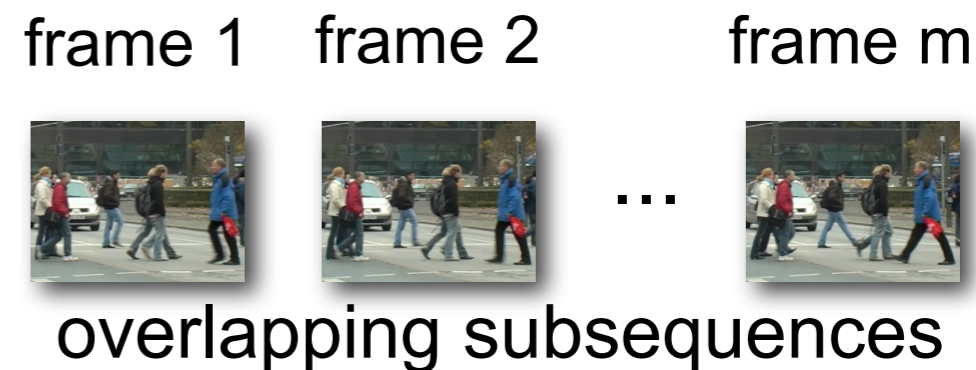
...

overlapping subsequences

- Want:



$$\mathbf{X}^{o*} = [\mathbf{x}_1^{o*}, \ldots, \mathbf{x}_m^{o*}]$$

body positions

- Posterior over positions and configurations:

# Tracklet Detection in Short Subsequences

frame 1    frame 2         frame m



... 

overlapping subsequences

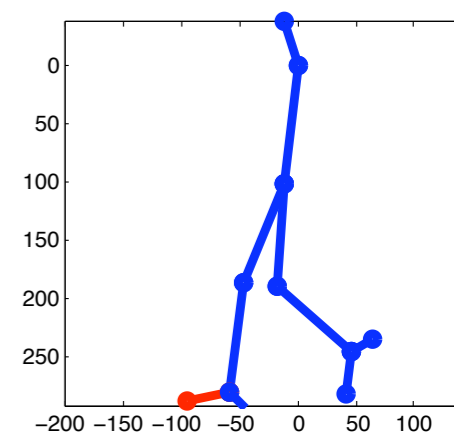- Given: $E = [E_1, \dots, E_m]$

- Want:



$$\mathbf{X}^{o*} = [\mathbf{x}_1^{o*}, \dots, \mathbf{x}_m^{o*}]$$
body positions

$$\mathbf{Y}^* = [\mathbf{y}_1^*, \dots, \mathbf{y}_m^*]$$
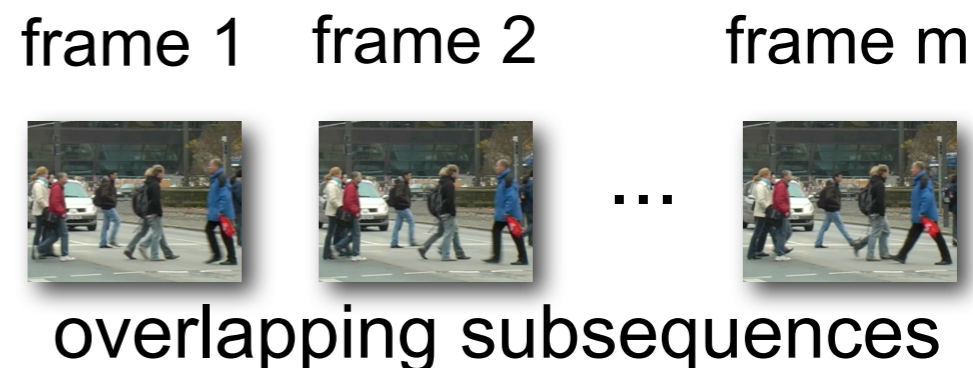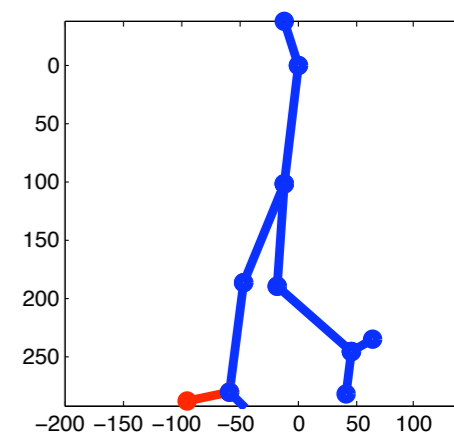body configurations
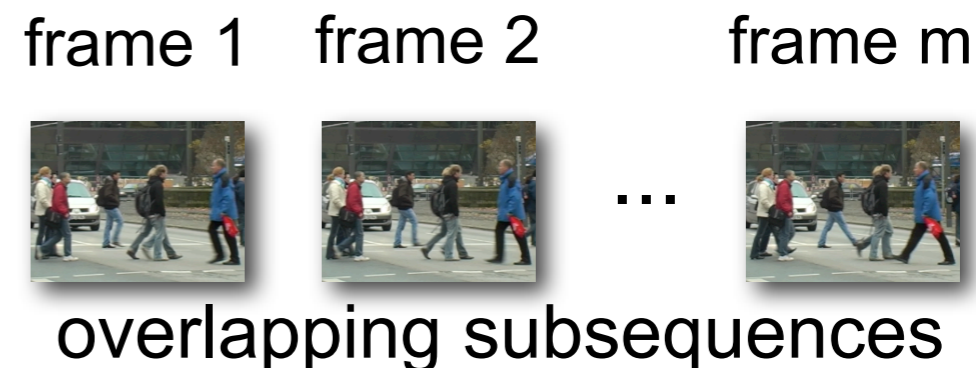


- Posterior over positions and configurations:

# Tracklet Detection in Short Subsequences

- Given: $E = [E_1, \ldots, E_m]$

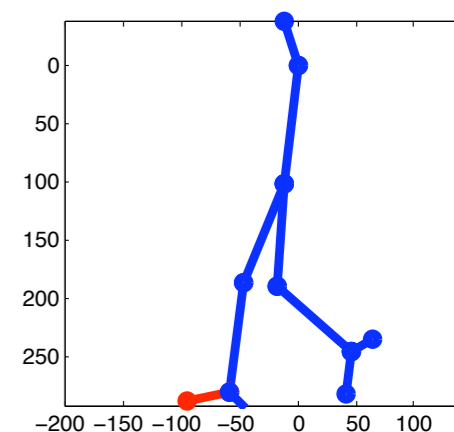overlapping subsequences

- Want:



$$\mathbf{X}^{o*} = [\mathbf{x}_1^{o*}, \ldots, \mathbf{x}_m^{o*}]$$
body positions

$$\mathbf{Y}^* = [\mathbf{y}_1^*, \ldots, \mathbf{y}_m^*]$$
body configurations



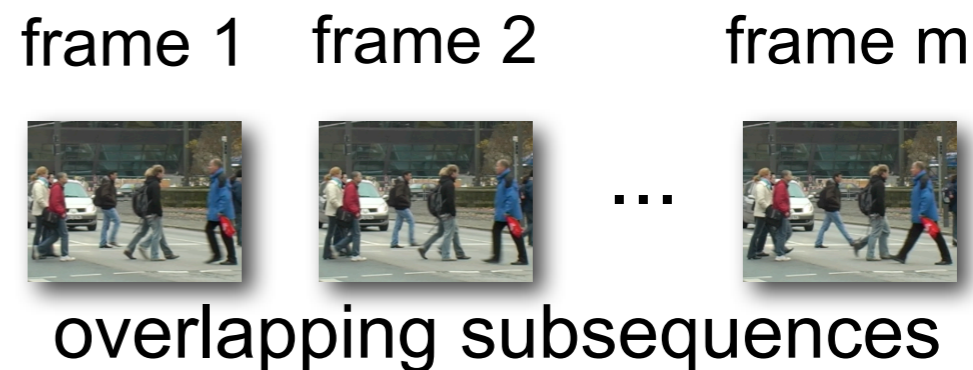- Posterior over positions and configurations:

$$p(\mathbf{X}^{o*}, \mathbf{Y}^*|E) \propto p(E|\mathbf{X}^{o*}, \mathbf{Y}^*)p(\mathbf{X}^{o*})p(\mathbf{Y}^*).$$

# Tracklet Detection in Short Subsequences

- Given: $E = [E_1, \ldots, E_m]$

frame 1   frame 2       frame m

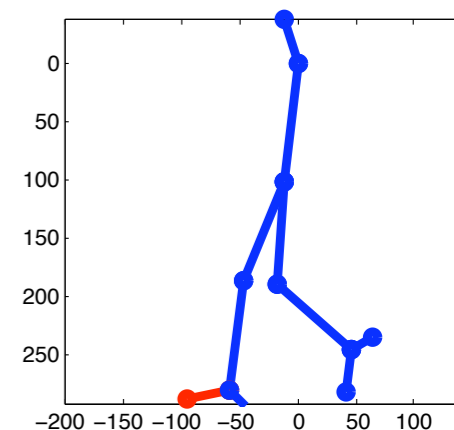

...

overlapping subsequences

- Want:

$$\mathbf{X}^{o*} = [\mathbf{x}_1^{o*}, \ldots, \mathbf{x}_m^{o*}]$$
body positions

$$\mathbf{Y}^* = [\mathbf{y}_1^*, \ldots, \mathbf{y}_m^*]$$
body configurations

- Posterior over positions and configurations:

$$p(\mathbf{X}^{o*}, \mathbf{Y}^*|E) \propto p(E|\mathbf{X}^{o*}, \mathbf{Y}^*)p(\mathbf{X}^{o*})p(\mathbf{Y}^*).$$

Likelihood model
(partISM)

# Tracklet Detection in Short Subsequences

- Given: $E = [E_1, \ldots, E_m]$



frame 1    frame 2     frame m

overlapping subsequences

- Want:



$$\mathbf{X}^{o*} = [\mathbf{x}_1^{o*}, \ldots, \mathbf{x}_m^{o*}]$$
body positions

$$\mathbf{Y}^* = [\mathbf{y}_1^*, \ldots, \mathbf{y}_m^*]$$
body configurations

- Posterior over positions and configurations:

$$p(\mathbf{X}^{o*}, \mathbf{Y}^* | E) \propto p(E | \mathbf{X}^{o*}, \mathbf{Y}^*) p(\mathbf{X}^{o*}) p(\mathbf{Y}^*).$$

Likelihood model
(partISM)

speed prior (Gaussian)

# Tracklet Detection in Short Subsequences

- Given: $E = [E_1, \ldots, E_m]$

frame 1    frame 2    frame m

  ... 

overlapping subsequences

- Want:

$$\mathbf{X}^{o*} = [\mathbf{x}_1^{o*}, \ldots, \mathbf{x}_m^{o*}]$$
body positions

$$\mathbf{Y}^* = [\mathbf{y}_1^*, \ldots, \mathbf{y}_m^*]$$
body configurations

- Posterior over positions and configurations:

$$p(\mathbf{X}^{o*}, \mathbf{Y}^* | E) \propto p(E | \mathbf{X}^{o*}, \mathbf{Y}^*) p(\mathbf{X}^{o*}) p(\mathbf{Y}^*).$$

Likelihood model
(partISM)

speed prior (Gaussian)

dynamical body model
(hGPLVM)

# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]
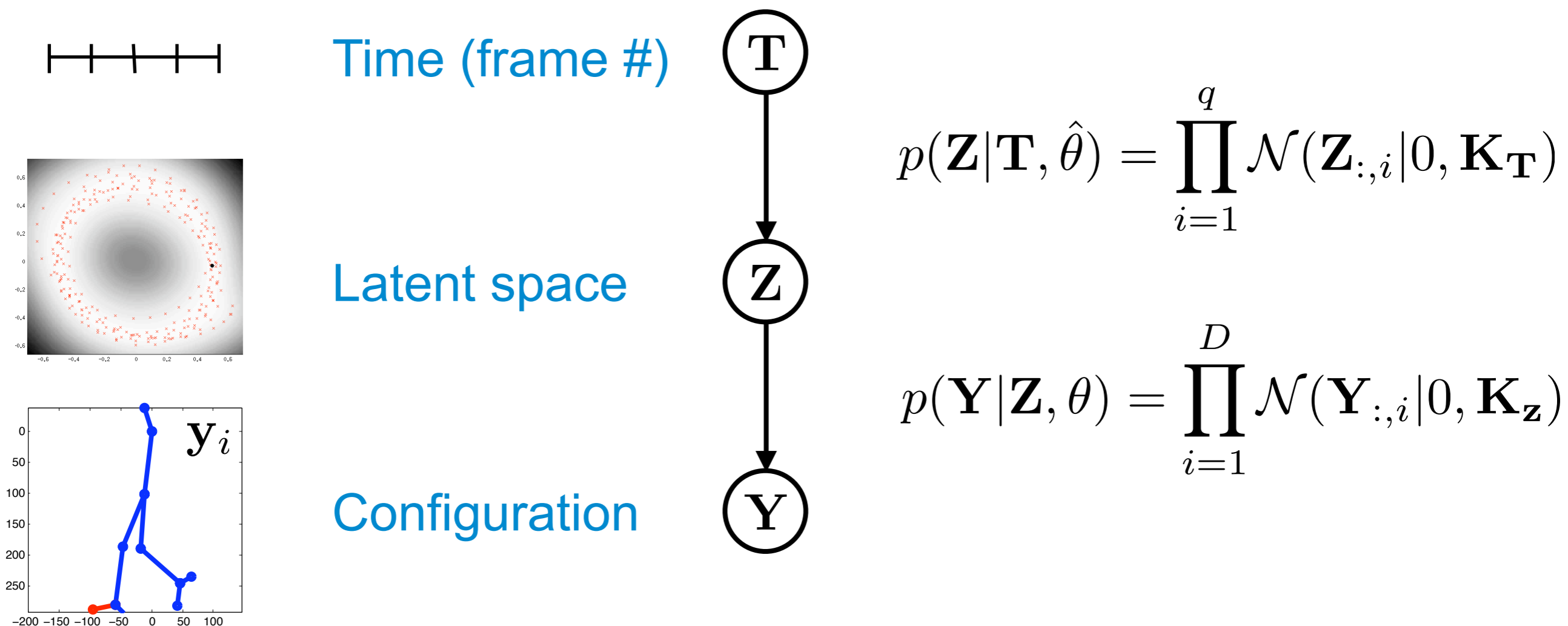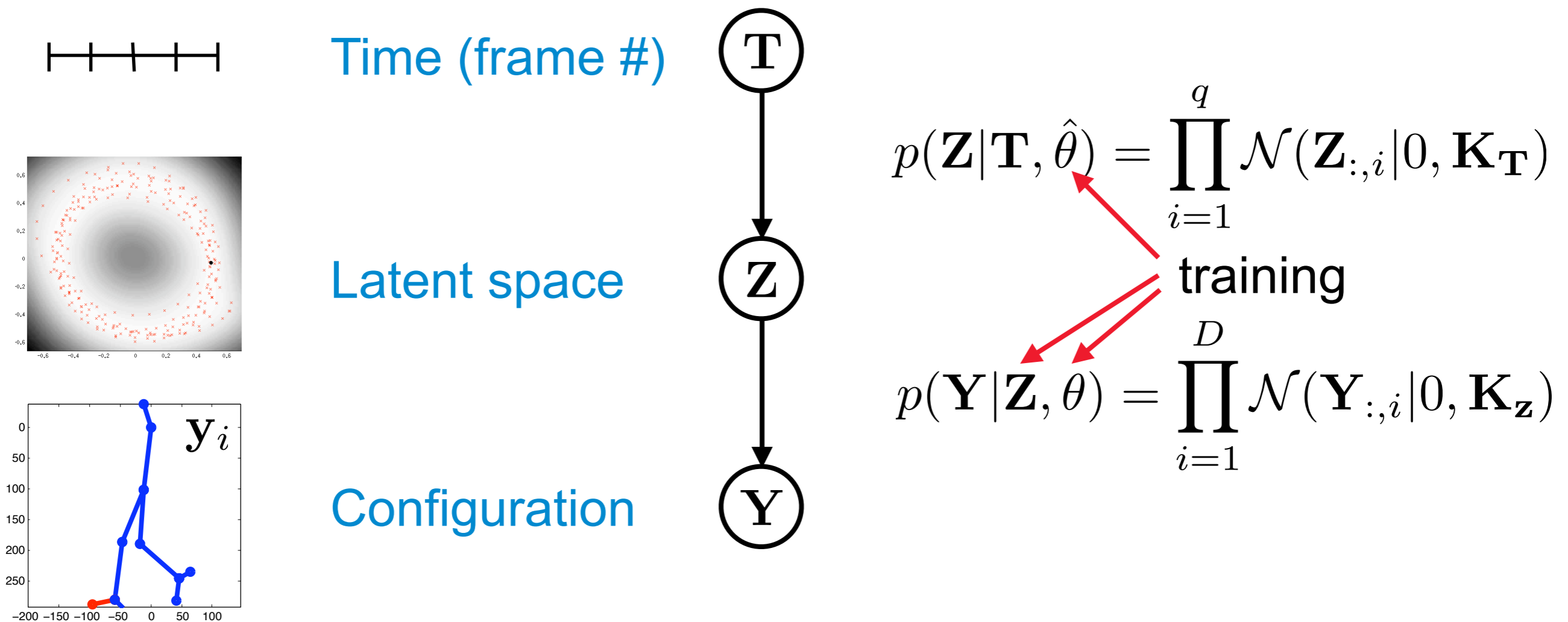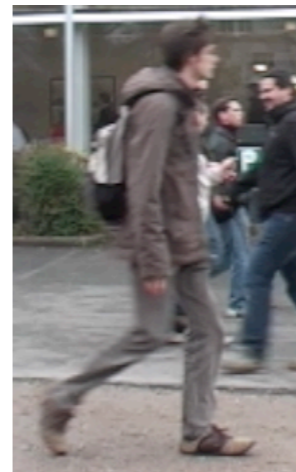
# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]



Configuration $\quad \textcircled{\mathbf{Y}} \quad \mathbf{Y} = [\mathbf{y}_i \in \mathbb{R}^D]$

# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]



Latent space

$\mathbf{Z}$

$\mathbf{Z} = [\mathbf{z}_i \in \mathbb{R}^q]$

$\mathbf{y}_i$

Configuration

$\mathbf{Y}$

$\mathbf{Y} = [\mathbf{y}_i \in \mathbb{R}^D]$

# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]

Time (frame #)     $\mathbf{T}$     $\mathbf{T} = [t_i \in \mathbb{R}]$

Latent space     $\mathbf{Z}$     $\mathbf{Z} = [\mathbf{z}_i \in \mathbb{R}^q]$

$\mathbf{y}_i$

Configuration     $\mathbf{Y}$     $\mathbf{Y} = [\mathbf{y}_i \in \mathbb{R}^D]$

# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]

Time (frame #)  $\mathbf{T}$

Latent space  $\mathbf{Z}$

$$p(\mathbf{Y}|\mathbf{Z}, \theta) = \prod_{i=1}^{D} \mathcal{N}(\mathbf{Y}_{:,i}|0, \mathbf{K_z})$$

$\mathbf{y}_i$

Configuration  $\mathbf{Y}$

# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]



Time (frame #)

$$p(\mathbf{Z}|\mathbf{T},\hat{\theta}) = \prod_{i=1}^{q} \mathcal{N}(\mathbf{Z}_{:,i}|0,\mathbf{K_T})$$

Latent space

$$p(\mathbf{Y}|\mathbf{Z},\theta) = \prod_{i=1}^{D} \mathcal{N}(\mathbf{Y}_{:,i}|0,\mathbf{K_z})$$

Configuration

$\mathbf{y}_i$

# Modeling Body Dynamics

- $\mathbf{Y}^*$ is very high-dimensional: Full body poses in $m$ frames.

- Model the body dynamics using hierarchical Gaussian process latent variable model (hGPLVM) [Lawrence&Moore, ICML 2007]

Time (frame #)  **T**

Latent space  **Z**

Configuration  **Y**

$$p(\mathbf{Z}|\mathbf{T},\hat{\theta}) = \prod_{i=1}^{q} \mathcal{N}(\mathbf{Z}_{:,i}|0,\mathbf{K_T})$$

training

$$p(\mathbf{Y}|\mathbf{Z},\theta) = \prod_{i=1}^{D} \mathcal{N}(\mathbf{Y}_{:,i}|0,\mathbf{K_z})$$

$\mathbf{y}_i$

# Tracklet Detection

- Tracklets are local maxima of:

$$p(\mathbf{X}^{o*}, \mathbf{Y}^* | E) \propto p(E | \mathbf{X}^{o*}, \mathbf{Y}^*) p(\mathbf{X}^{o*}) p(\mathbf{Y}^*).$$

- Local maxima can be found using standard non-linear optimization (e.g. conjugate gradients).

- **How can we provide good initial hypotheses for optimization?**

# Tracklet Detection

# Tracklet Detection

# Tracklet Detection

# Tracklet Detection



propagate detection

# Tracklet Detection



propagate detection

hGPLVM mean prediction

# Tracklet Detection



propagate detection

hGPLVM mean prediction

# Tracklet Detection



propagate detection

hGPLVM mean prediction

# Tracklet Detection



propagate
detection

hGPLVM
mean prediction

# Tracklet Detection



propagate detection

hGPLVM mean prediction

pose optimization

# Single-Frame Detector vs. Tracklet Detector

- ## At equal error rate:



partISM

Tracklet detector

▸ Fewer false positives.

▸ More robust detection of partially occluded people.

# Single-Frame Detector vs. Tracklet Detector

- **At equal error rate:**



- ▶ Fewer false positives.
- ▶ More robust detection of partially occluded people.

# Single-Frame Detector vs. Tracklet Detector

- ## At equal error rate:



partISM

Tracklet detector

- ▸ Fewer false positives.
- ▸ More robust detection of partially occluded people.

# Single-Frame Detector vs. Tracklet Detector

- ## At equal error rate:



- ▸ Fewer false positives.
- ▸ More robust detection of partially occluded people.

# Single-Frame Detector vs. Tracklet Detector

- **At equal error rate:**



▶ Fewer false positives.

▶ More robust detection of partially occluded people.

# Detection Performance



TUD campus data

With occlusions
(up to 50%)

- Significant improvement over single-frame detector.
  ▸ Also at high precision levels.

# Overview

Three stages of our multi-person detection and tracking system:



1. Single-frame detection

2. Tracklet detection

3. Tracking through occlusion

# Tracks from Overlapping Tracklets



$$t \qquad\qquad t+1 \qquad\qquad t+2 \qquad\qquad t+3$$

# Tracks from Overlapping Tracklets



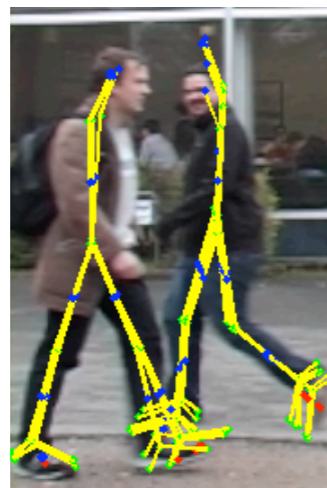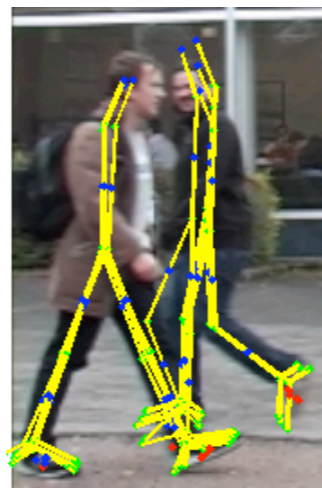$t$ $\qquad\qquad$ $t+1$ $\qquad\qquad$ $t+2$ $\qquad\qquad$ $t+3$ $\qquad$ ...

Candidate poses from all overlapping tracklets

# Tracks from Overlapping Tracklets



$t$        $t+1$        $t+2$        $t+3$

...

Candidate poses from all
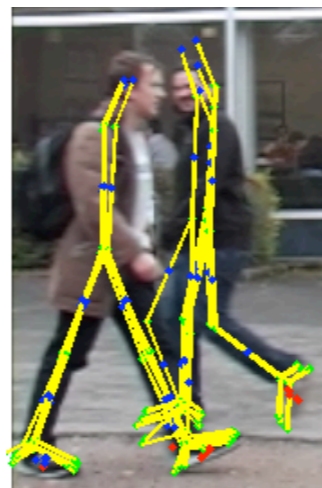overlapping tracklets

# Tracks from Overlapping Tracklets



$t$        $t+1$        $t+2$        $t+3$    ...
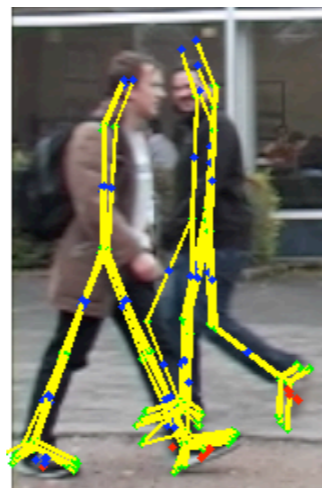
Candidate poses from all overlapping tracklets

# Tracks from Overlapping Tracklets



$t$         $t+1$         $t+2$         $t+3$

…

Candidate poses from all overlapping tracklets

# Tracks from Overlapping Tracklets



$t$          $t+1$          $t+2$          $t+3$      ...

Candidate poses from all overlapping tracklets

# Tracks from Overlapping Tracklets



$t$        $t+1$        $t+2$        $t+3$

Candidate poses from all
overlapping tracklets

# Tracks from Overlapping Tracklets



$t$         $t+1$         $t+2$         $t+3$
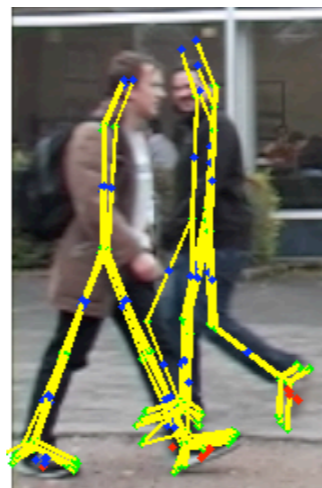
...

# Tracks from Overlapping Tracklets



$t$     $t+1$     $t+2$     $t+3$     ...
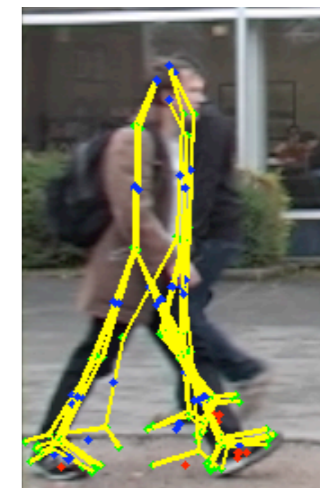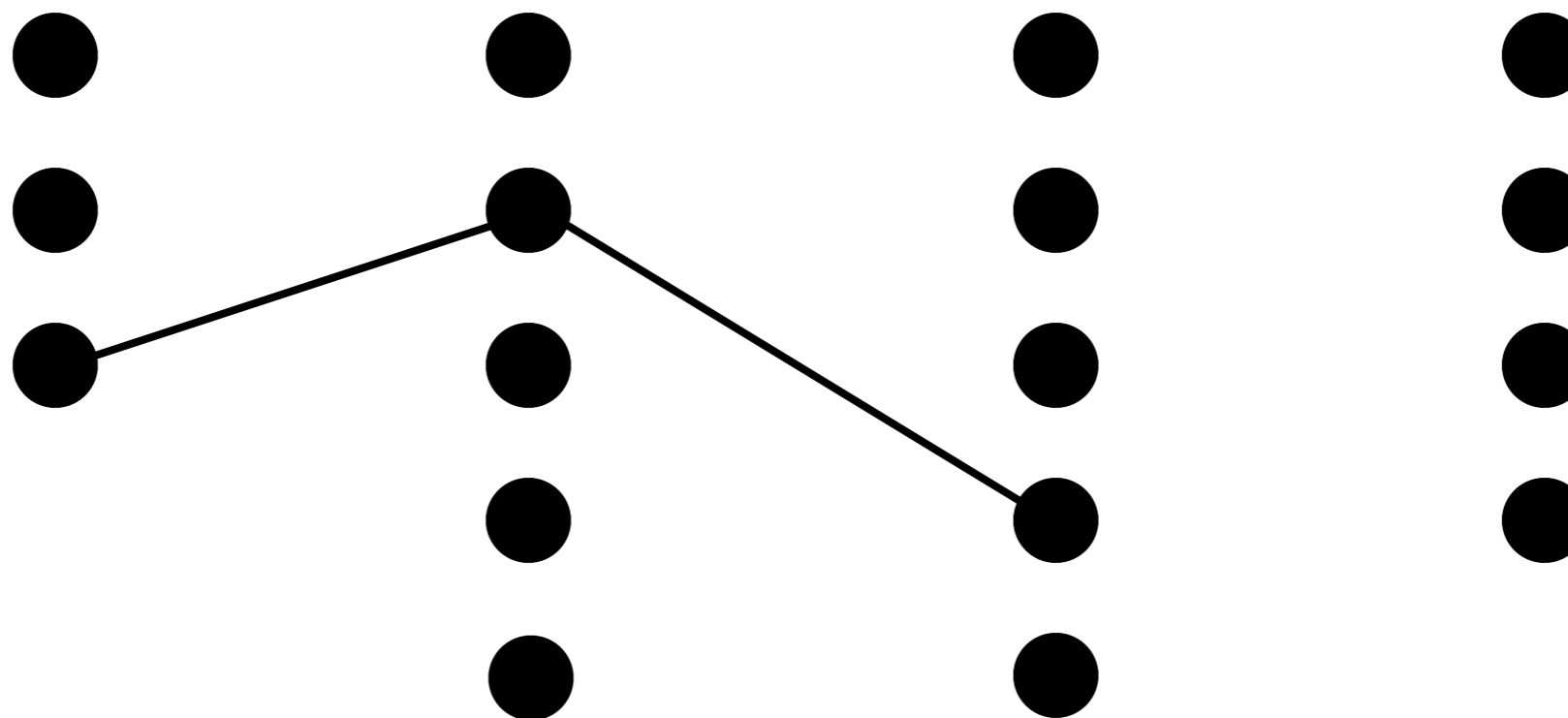
Viterbi
Decoding

# Tracks from Overlapping Tracklets



Viterbi Decoding

# Tracks from Overlapping Tracklets



Viterbi Decoding

$t$  $t+1$  $t+2$  $t+3$

...

# Tracks from Overlapping Tracklets



$t$  $t+1$  $t+2$  $t+3$

Viterbi Decoding

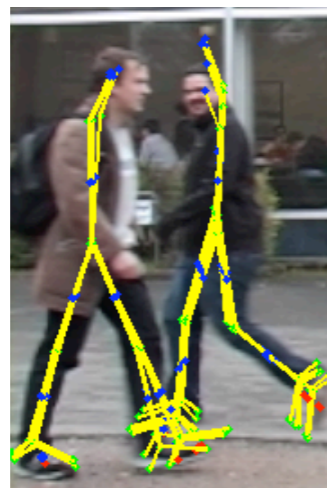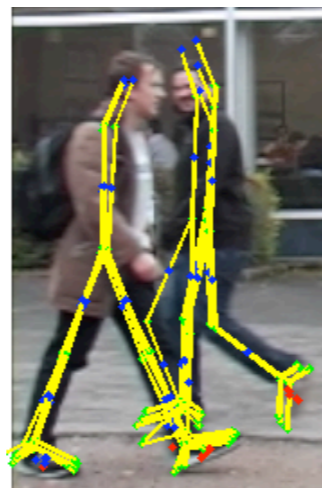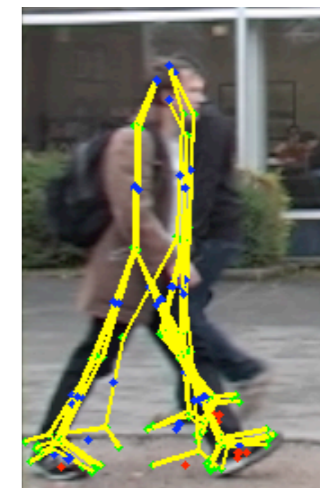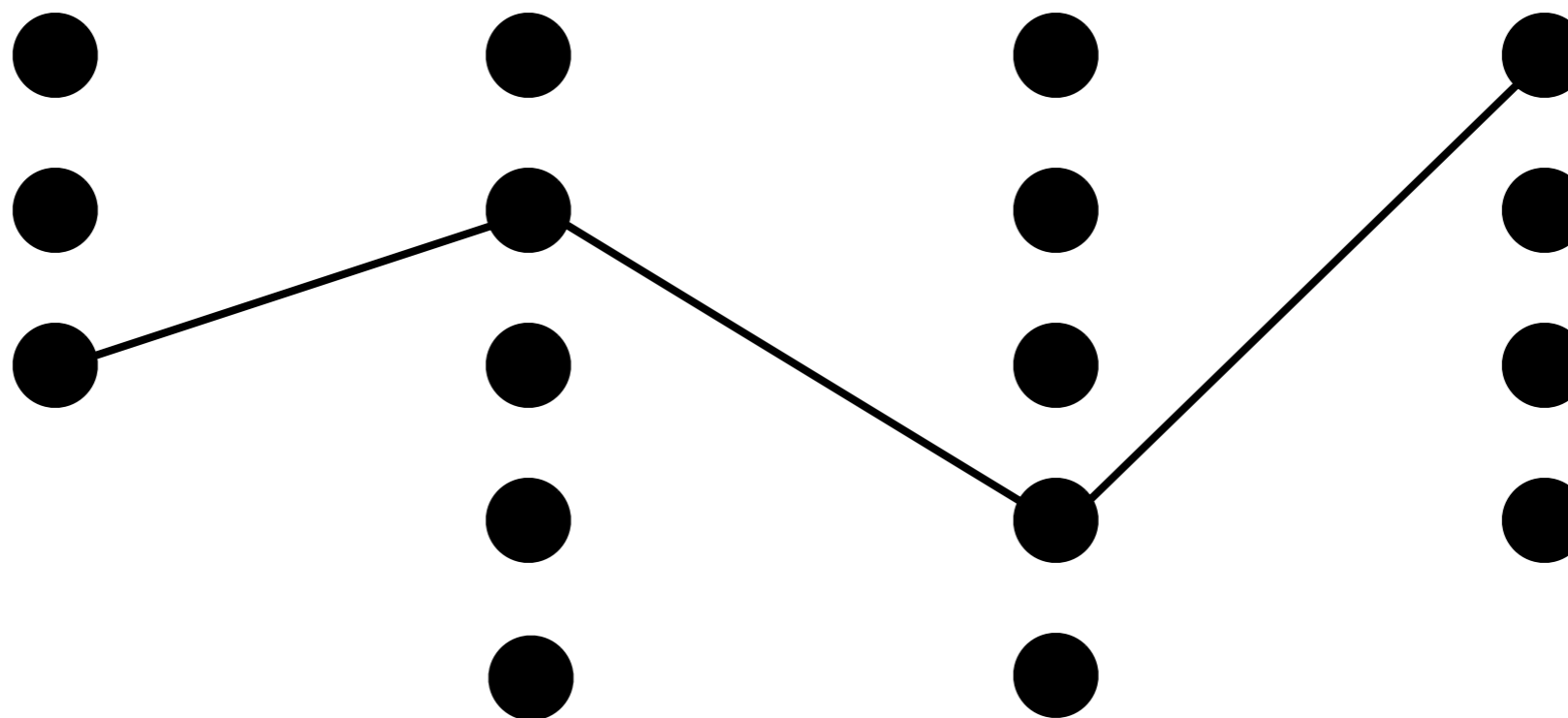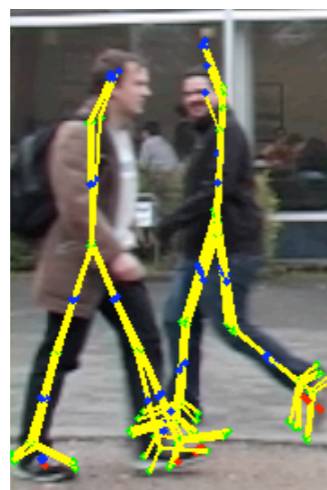# Tracks from Overlapping Tracklets



$t$       $t+1$       $t+2$       $t+3$

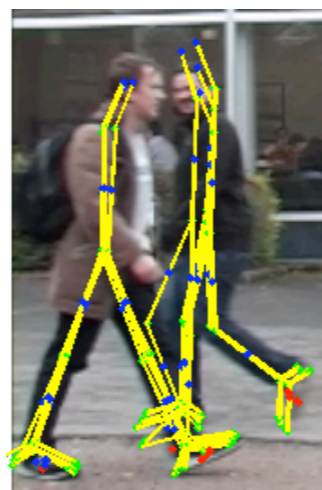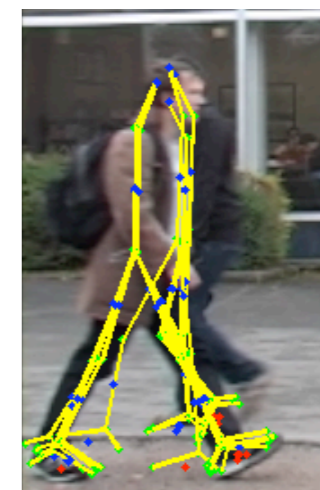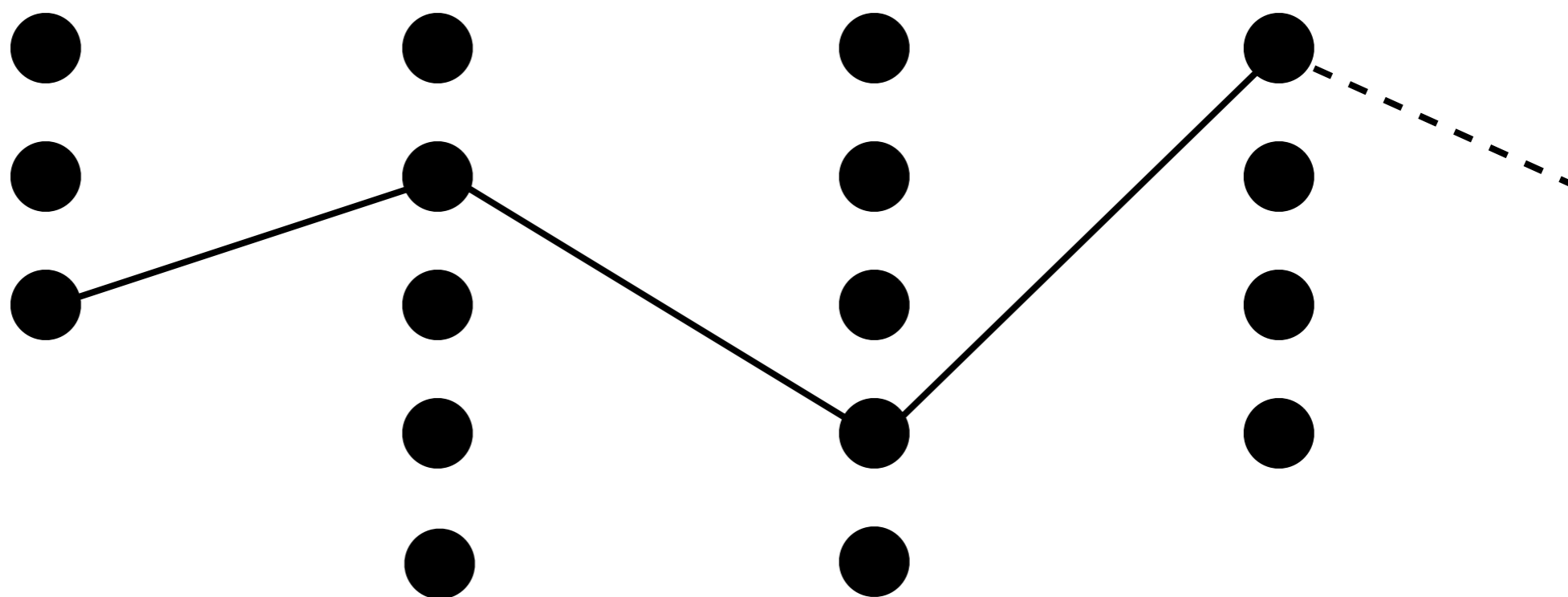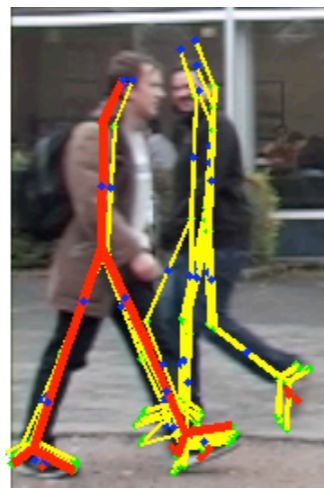Viterbi Decoding

# Tracks from Overlapping Tracklets



$t$       $t+1$       $t+2$       $t+3$
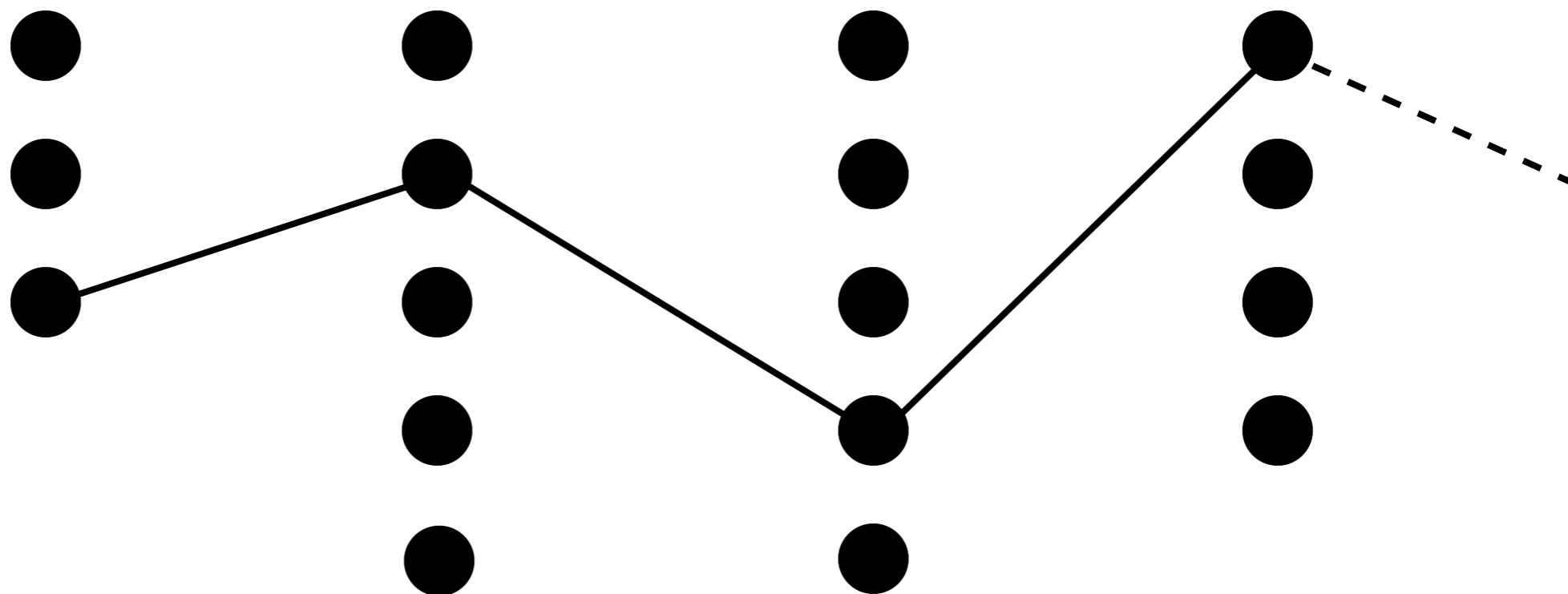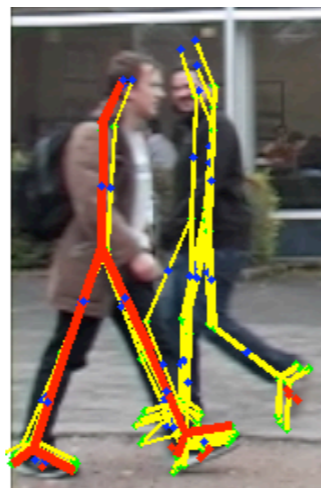
Viterbi Decoding

# Finding Multiple Tracks



$t$  $t+1$  $t+2$  $t+3$  ...

- Find the best track
- Remove its hypotheses
- Repeat
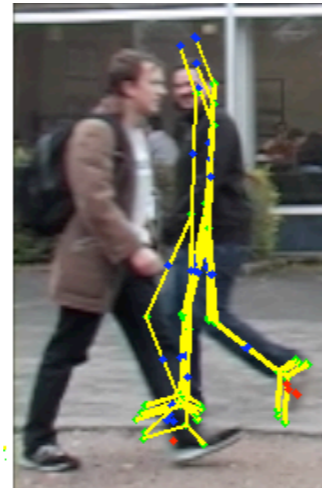
# Finding Multiple Tracks



- **Find the best track**
- **Remove its hypotheses**
- **Repeat**

# Finding Multiple Tracks



$t$       $t+1$       $t+2$       $t+3$

- Find the best track
- Remove its hypotheses
- Repeat

# Occlusion Event



$t$   $t+1$   $t+2$   $t+3$   ...

$t$       $t+1$       $t+2$       $t+3$

"bad" detections

# Occlusion Event



$t$  $\quad\quad$  $t+1$  $\quad\quad$  $t+2$  $\quad\quad$  $t+3$  $\quad$ ...

"bad" detections

# Occlusion Event



$t$     $t+1$     $t+2$     $t+3$     ...

"bad" detections

terminate if low-probability for any transition

# Appearance Model for Occlusion Recovery



- Extract person-specific appearance model for each limb:
  ▸ Color histogram.

- Require relatively accurate pose estimate:
  ▸ Pose from extracted tracks.

- Appearance comparison measure:
  ▸ Bhattacharyya distance.

# Appearance Model for Occlusion Recovery



- **Extract person-specific appearance model for each limb:**
  - ▸ Color histogram.

- **Require relatively accurate pose estimate:**
  - ▸ Pose from extracted tracks.

- **Appearance comparison measure:**
  - ▸ Bhattacharyya distance.

# Appearance Model for Occlusion Recovery



- Extract person-specific appearance model for each limb:
  ‣ Color histogram.

- Require relatively accurate pose estimate:
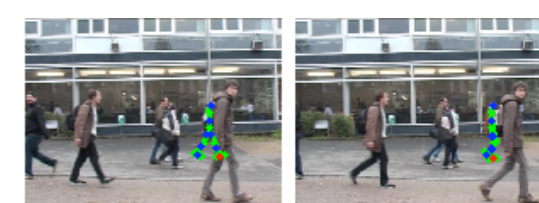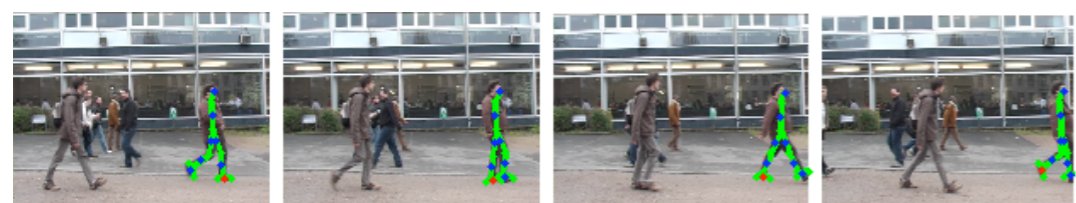  ‣ Pose from extracted tracks.

# Appearance Model for Occlusion Recovery



- **Extract person-specific appearance model for each limb:**
  - ▸ Color histogram.

- **Require relatively accurate pose estimate:**
  - ▸ Pose from extracted tracks.

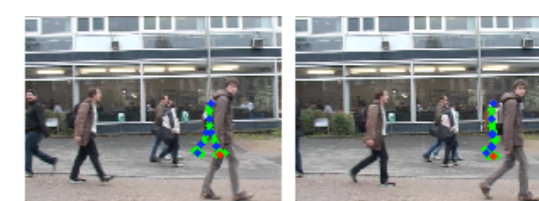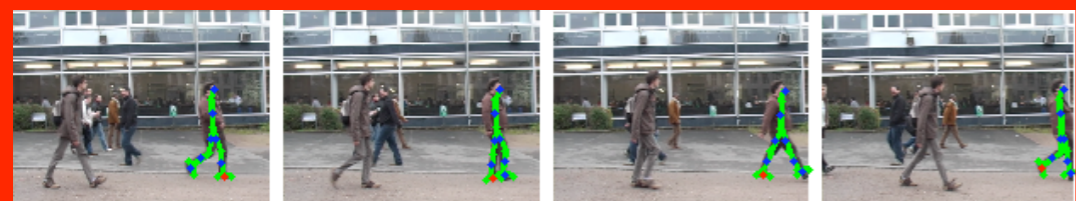- **Appearance comparison measure:**
  - ▸ Bhattacharyya distance.

# Occlusion Recovery



time ⟶

- Greedily link partial tracks based on:
  - ▶ Motion & articulation compatibility.
  - ▶ Plus appearance compatibility.

# Occlusion Recovery



time ⟶

- Greedily link partial tracks based on:
  - ▸ Motion & articulation compatibility.
  - ▸ Plus appearance compatibility.

# Occlusion Recovery



time ⟶

- Greedily link partial tracks based on:
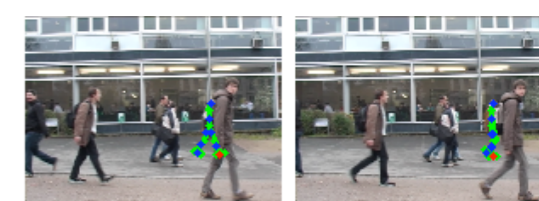  ▸ Motion & articulation compatibility.
  ▸ Plus appearance compatibility.

# Occlusion Recovery



time ⟶

- Greedily link partial tracks based on:
  - ▸ Motion & articulation compatibility.
  - ▸ Plus appearance compatibility.

# Occlusion Recovery



time ⟶

- Greedily link partial tracks based on:
    ▸ Motion & articulation compatibility.
    ▸ Plus appearance compatibility.

# Occlusion Recovery



time ⟶

- Greedily link partial tracks based on:
  - ▶ Motion & articulation compatibility.
  - ▶ Plus appearance compatibility.

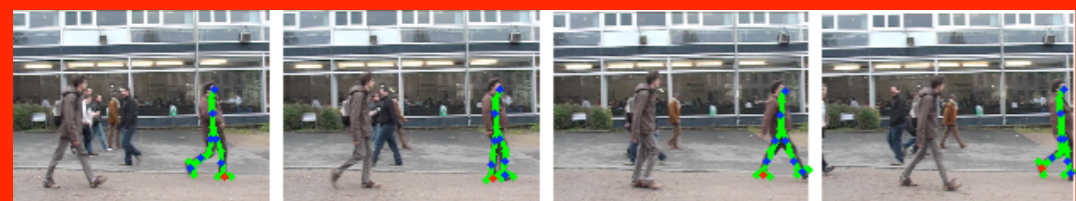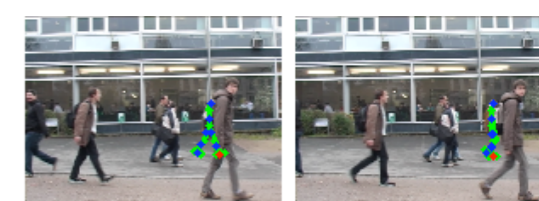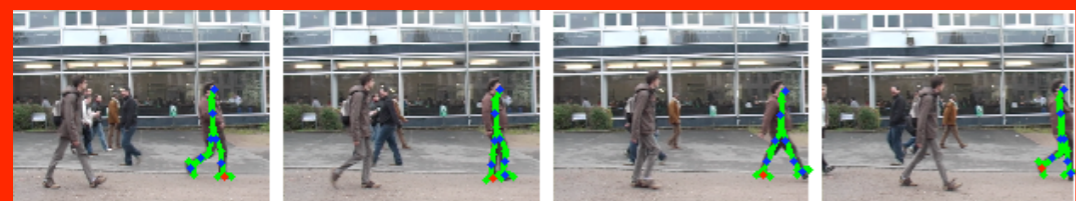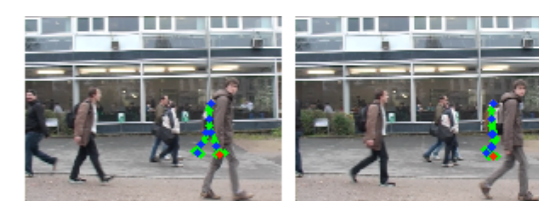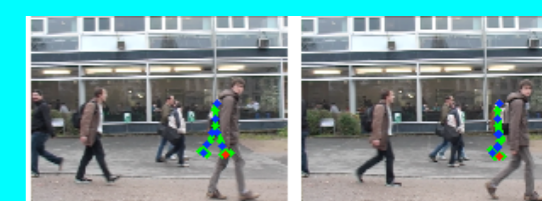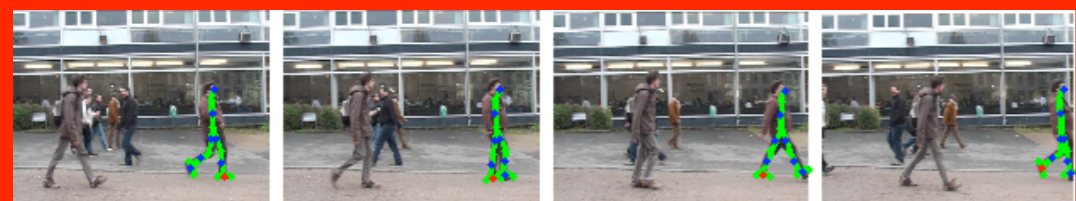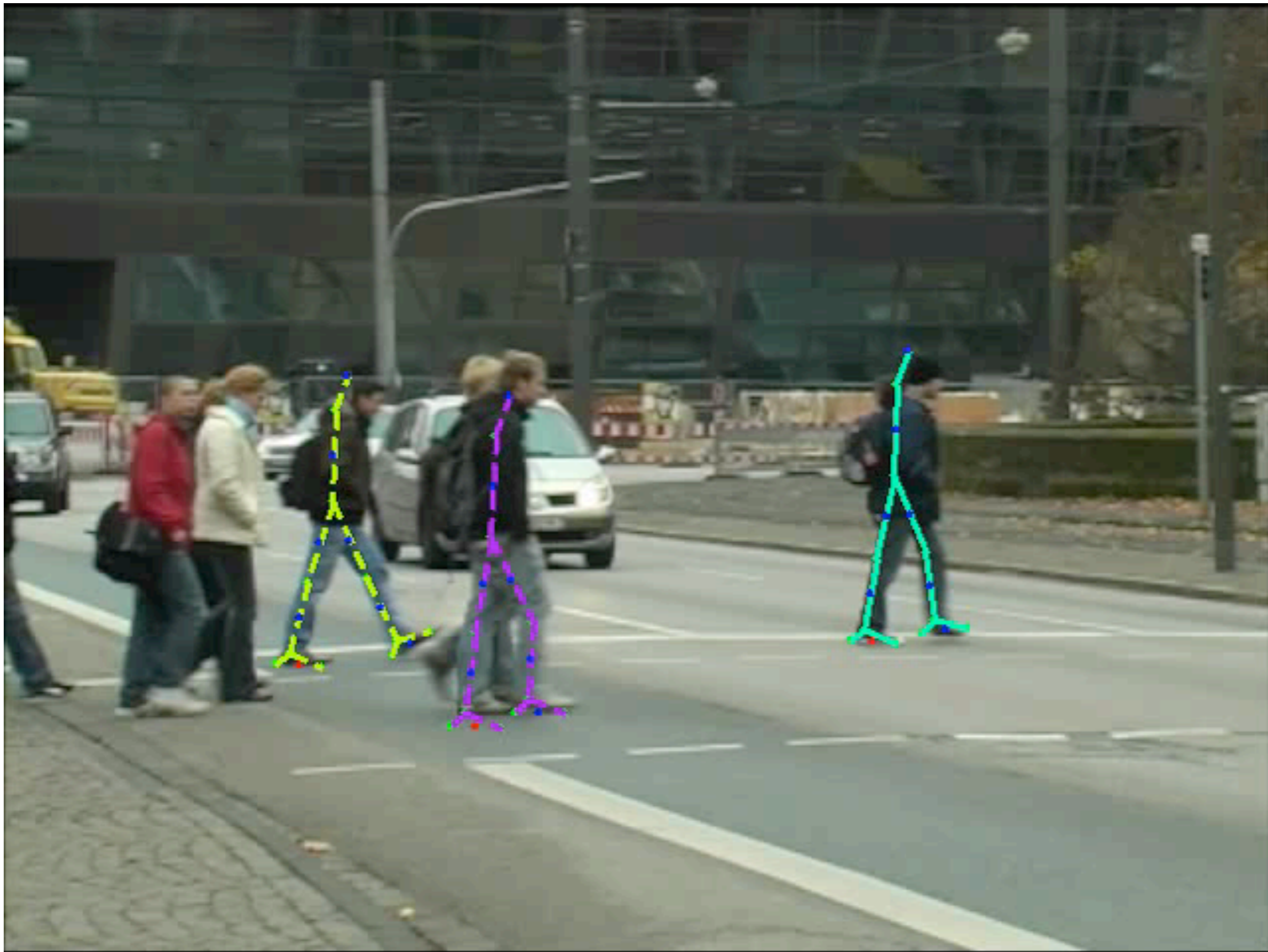# Occlusion Recovery



time →

- Greedily link partial tracks based on:
  - ▸ Motion & articulation compatibility.
  - ▸ Plus appearance compatibility.

# Summary

- partISM: Extended the ISM detection framework to part-based detection:

  ▸ Improved detection

  ▸ Basis for incorporating body dynamics.

- Incorporated temporal continuity in a "tracklet" detection framework:

  ▸ hGPLVM dynamics model.

  ▸ Improves occlusion robustness.

  ▸ Reduces false positives.

- Extracted and combined tracks across occlusion events:

  ▸ Person identification throughout entire sequences.

# Thanks!

- Acknowledgements:

  ‣ Neil Lawrence for his GPLVM code.

  ‣ Mario Fritz for helpful discussions.

  ‣ Partial funding through DFG GRK "Cooperative, Adaptive and Responsive Monitoring in Mixed Mode Environments"

  ‣ Travel funding from DFG.

- Data available at:

  `http://www.mis.informatik.tu-darmstadt.de/`