# Learning Context for Collective Activity Recognition

**Wongun Choi**   **Khuram Shahid**   **Silvio Savarese**

**Department of Electirical and Computer Engineering, University of Michigan, Ann Arbor MI**

MICHIGAN
**VISION LAB**

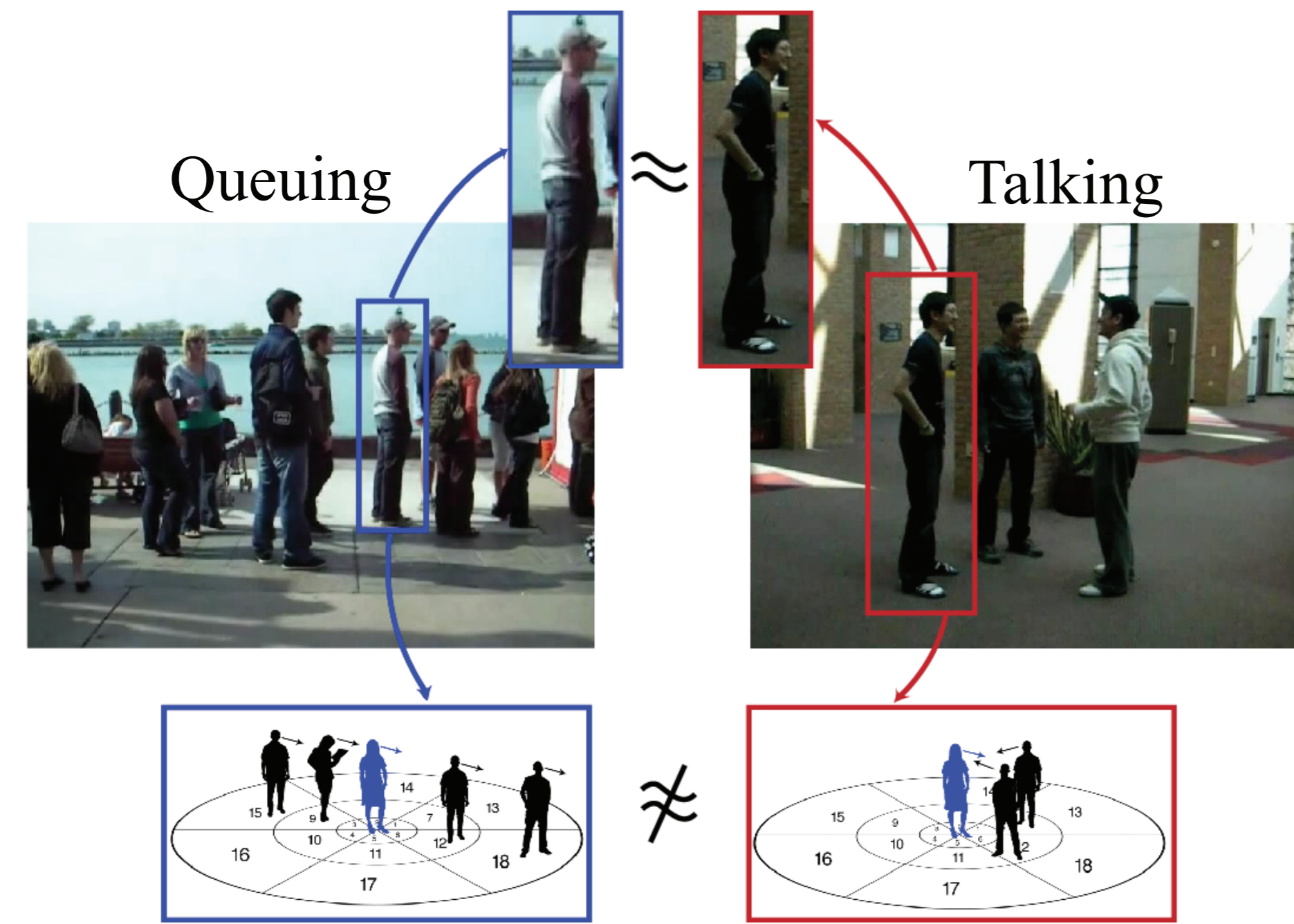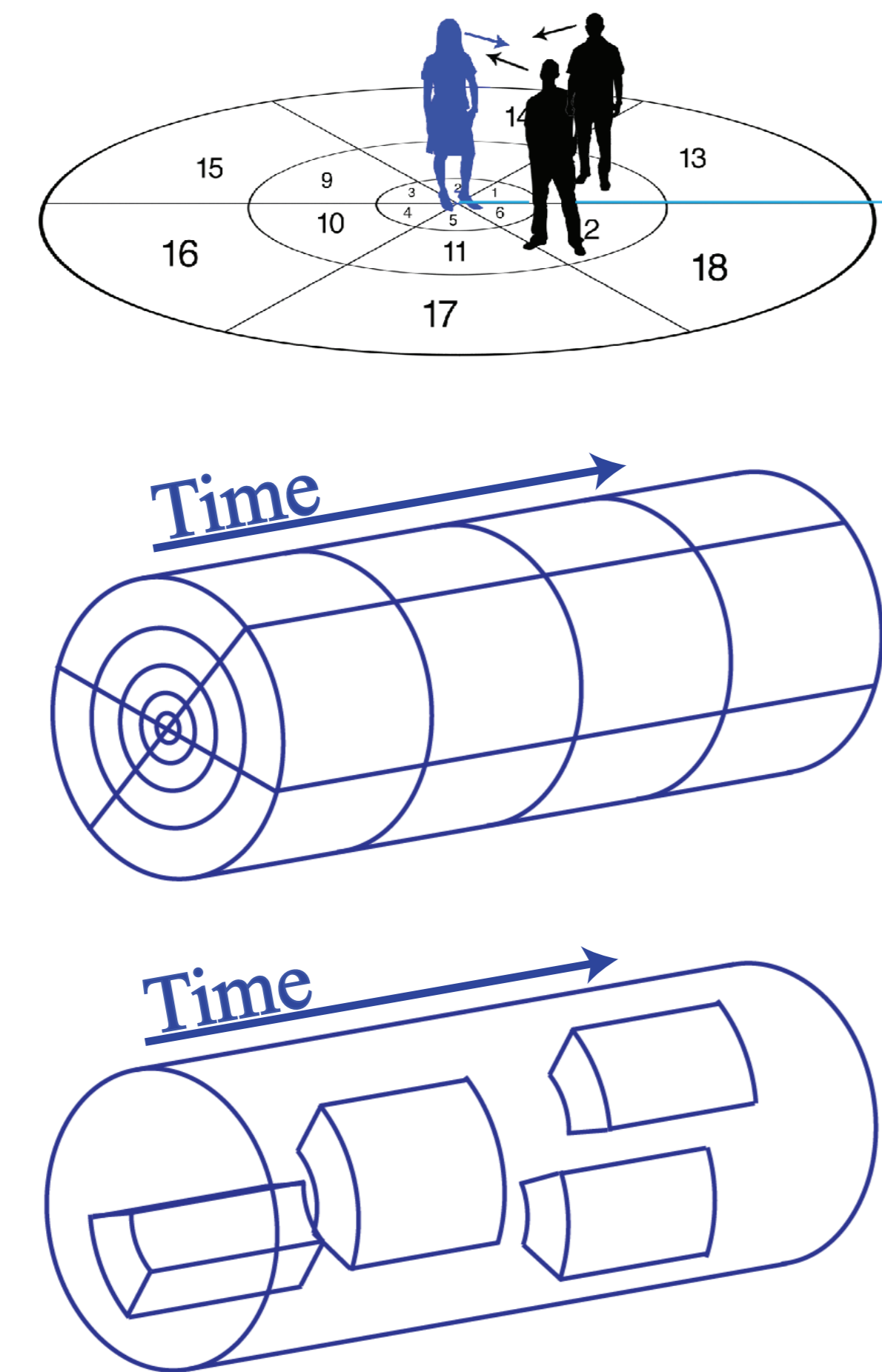THE UNIVERSITY OF MICHIGAN · 1817

## Introduction

- **Collective activity**
  - Activities that are defined by the interaction among people.
  - Cannot be characterized by single person's appearance.


Queuing ≈ Talking

- **Crowd Context**
  - Spatio-temporal context around one person.
  - Data-driven approach to learn the crowd context.

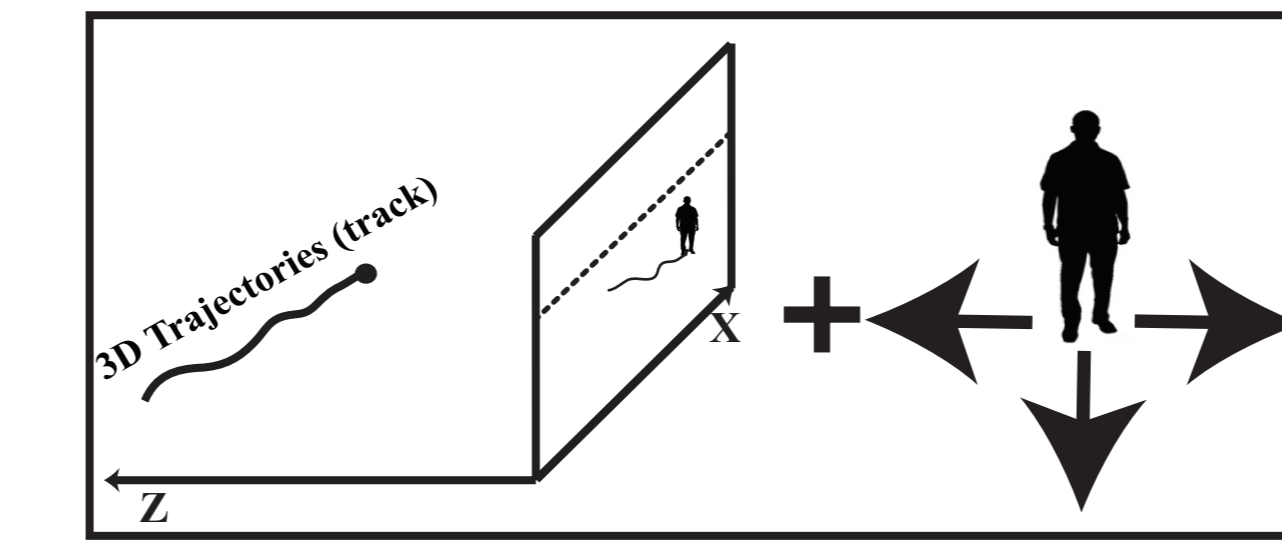- **Regularization by spatio-temporal consistency.**

## Crowd Context

- **STV [1]**
  - Appropriate for capturing spatio-temporal relationship.
  - Rigid structured descriptor
  - Susceptible to clutter
  - Several parameters to be tuned

- **RSTV**
  - Randomize the discretization in feature space.
  - Parameter free.
  - Structure learned from training data.
  - More flexible structure.
  - Higher scalablility.
  - Robustness under clutter.
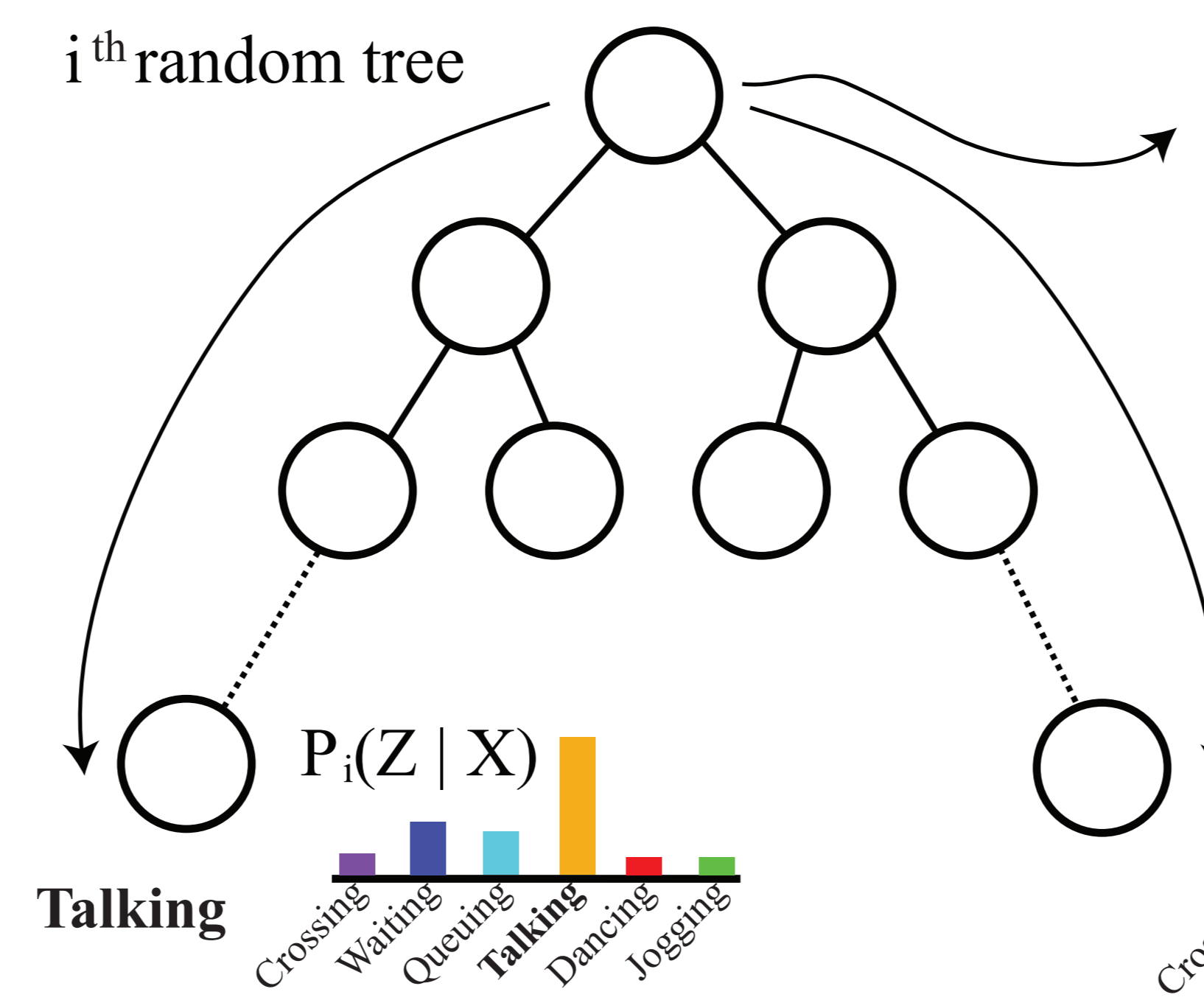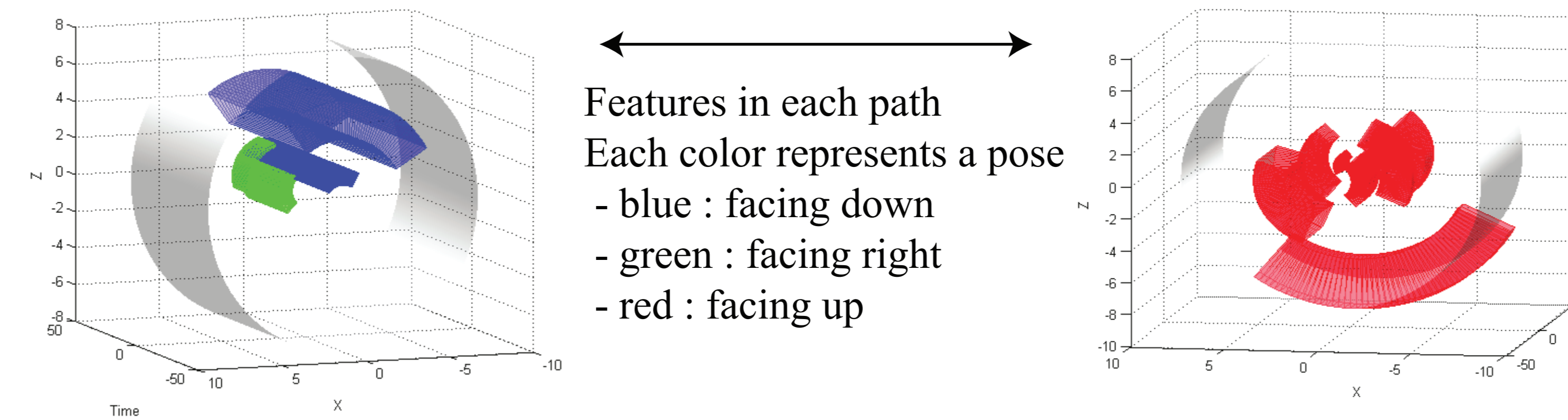
**CVPR 2011 Colorado Springs**

## Method

- **Inputs for algorithm**
  - 3D trajectories + velocity of people
  - N-directional pose classification



- **RSTV learning**
  - Divide training set into N random subset (bagging)
  - Train a tree with a subset of training data

$i$ th random tree

Randomly choose a set of features
- Spatial location (r, θ)
- Temporal support (t, Δt)
- Pose direction (p)
- Threshold (t)

$P_i(Z \mid X)$

Talking — Crossing Waiting Queuing **Talking** Dancing Jogging

$P_i(Z \mid X)$

Crossing **Waiting** Queuing Talking Dancing Jogging — Waiting

Features in each path
Each color represents a pose
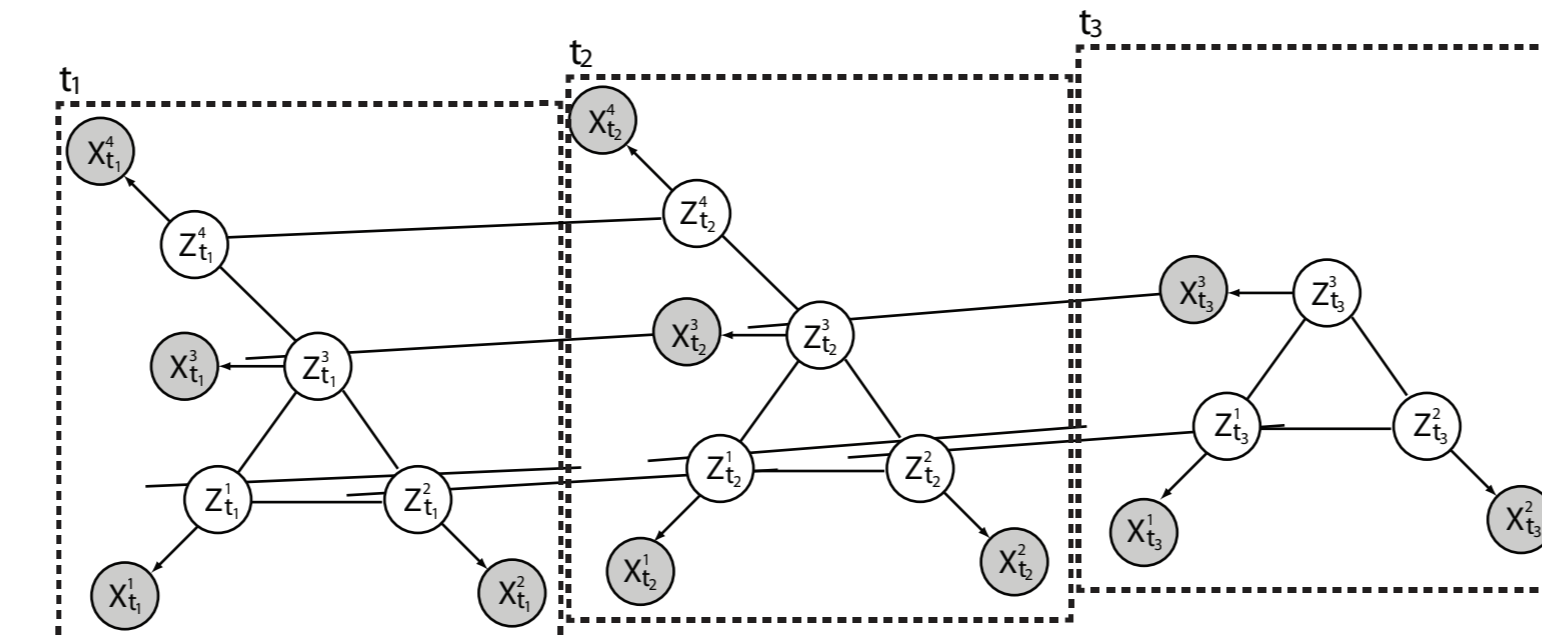- blue : facing down
- green : facing right
- red : facing up

- **Classification**
  - For each tree, find the leaf node for each example.
  - Compute sum of posterior probability

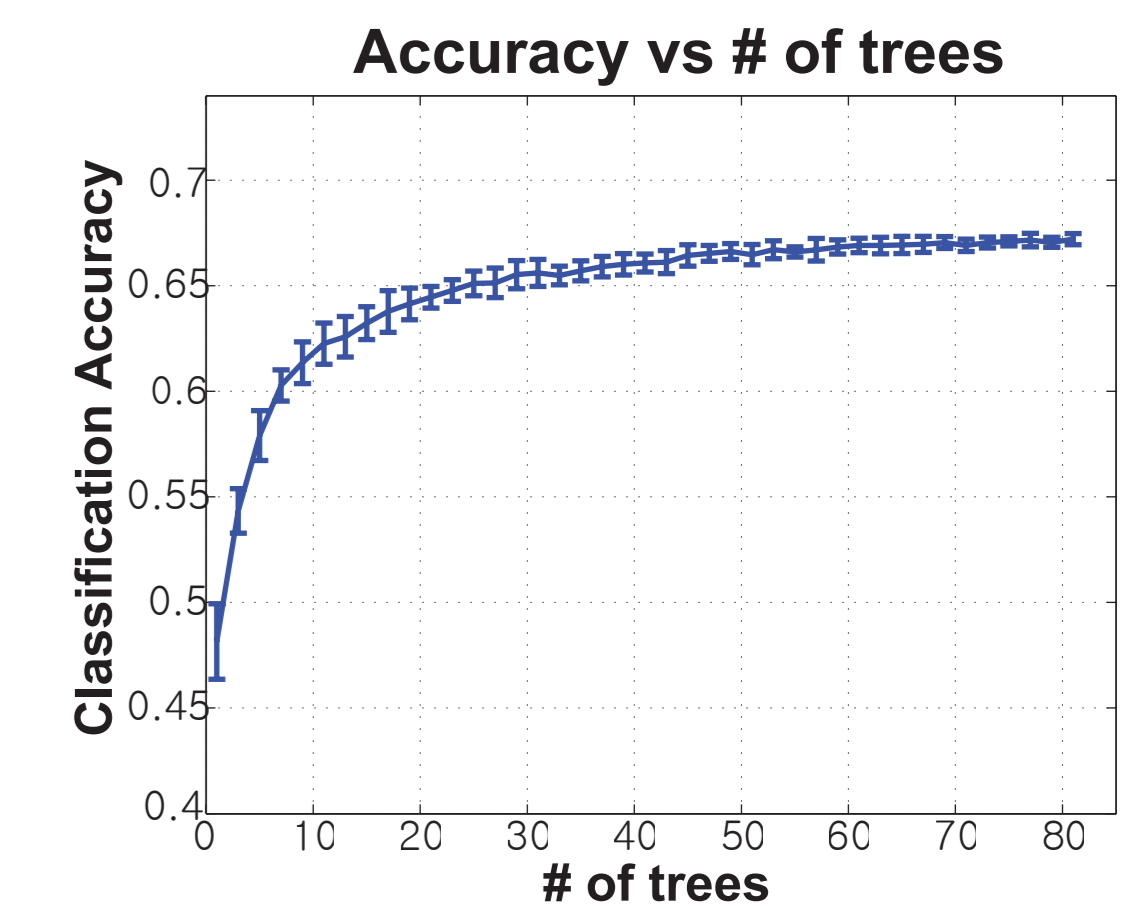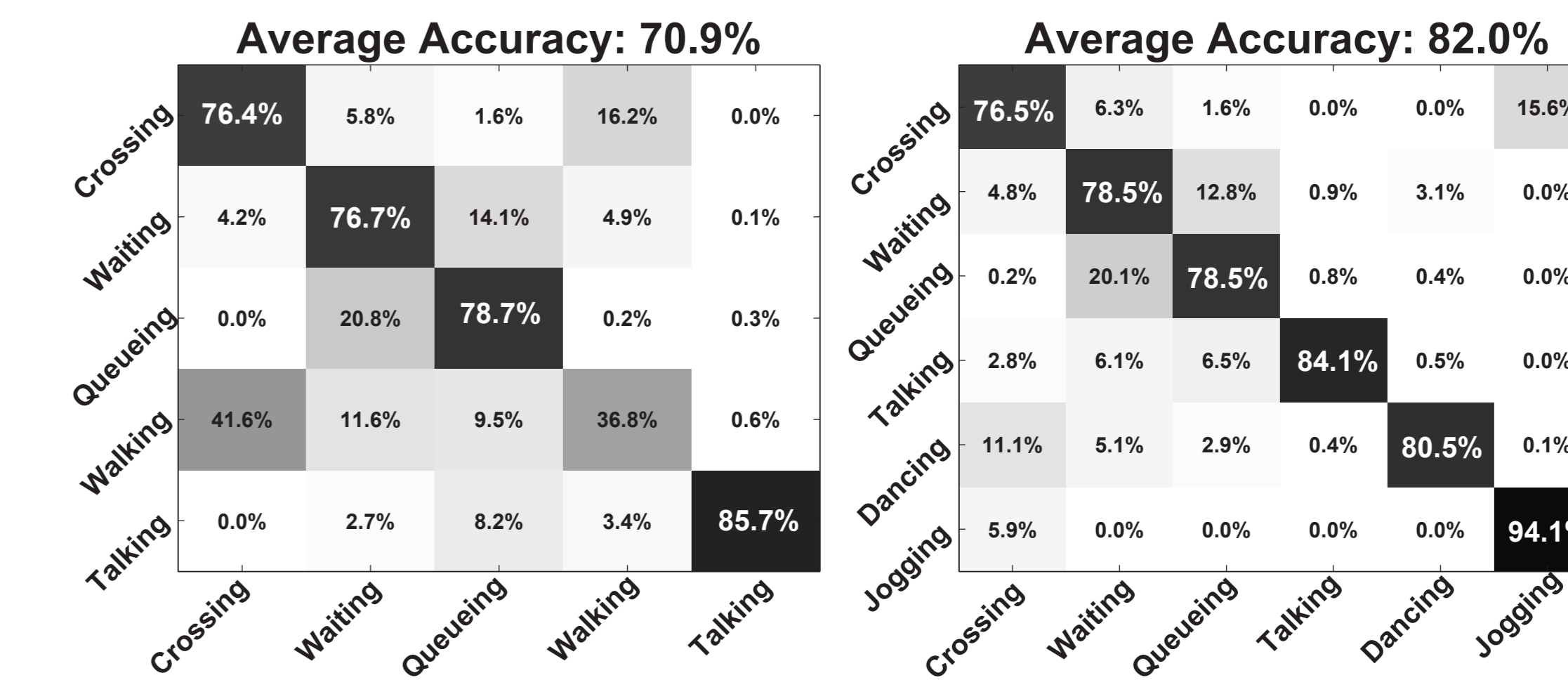$$\frac{1}{N}\sum_{i=1}^{N} P_i(Z \mid X)$$

- **MRF regularization**
  - Spatio-temporal context
  - Inference by Gibbs sampling

$$P(Z \mid X, p) \propto \prod_t \prod_i P(Z_t^i \mid X_t^i) \prod_t \prod_{(i,j)\in E_s} \Phi_S(Z_t^i, Z_t^j; p_t^i, p_t^j)$$
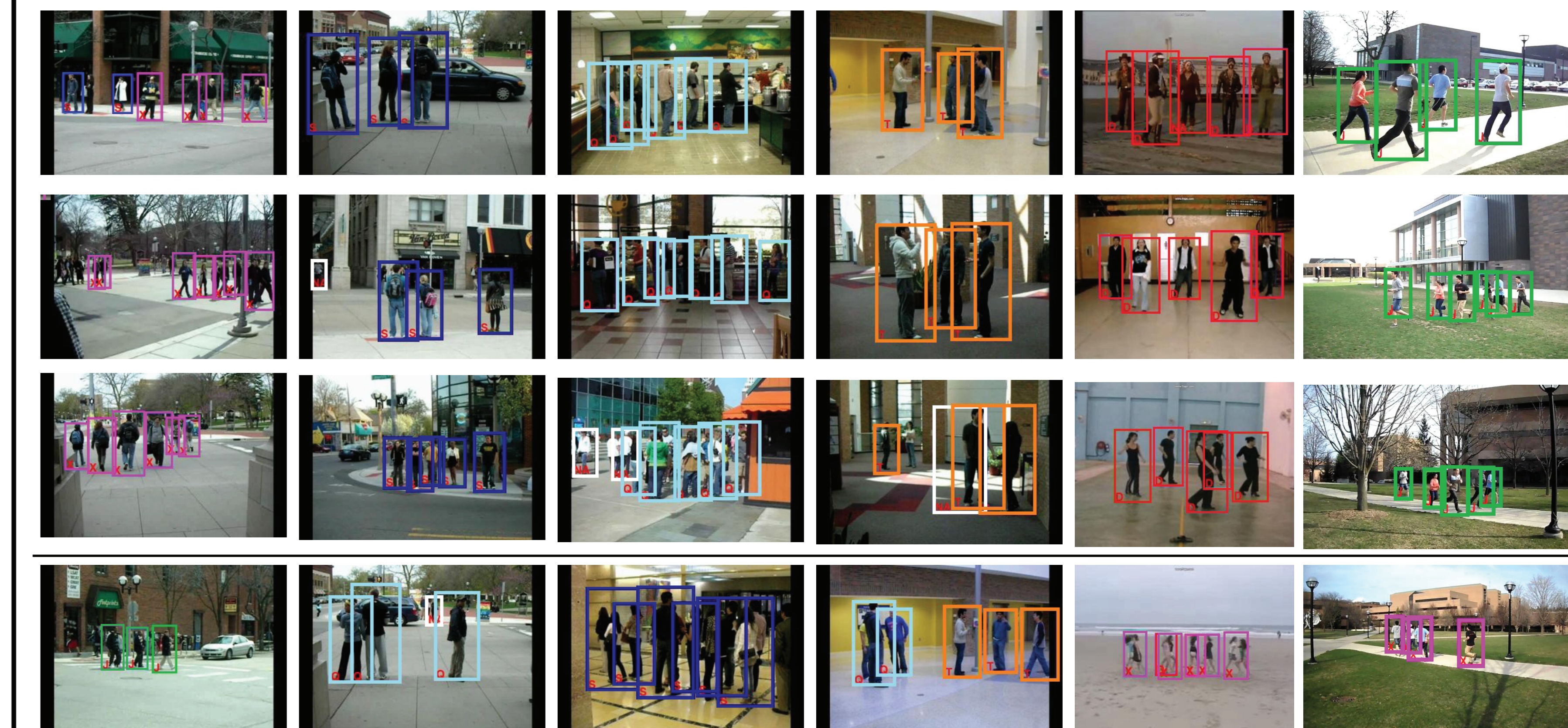
$$\prod_i \prod_t \Phi_T(Z_{t-1}^i, Z_t^i)$$

[1] W. Choi, K. Shahid, and S. Savarese. What are they doing? : Collective activity classification using spatio-temporal relationship among people. In Visual Surveillance Workshop, ICCV, 2009.
[2] T. Lan, Y. Wang, G. Mori, and S. Robinovitch. Retrieving actions in group contexts. In International Workshop on Sign Gesture Activity, 2010.

## Results

- **Classification Result**

**Average Accuracy: 70.9%**

| | Crossing | Waiting | Queuing | Walking | Talking |
|---|---|---|---|---|---|
| Crossing | 76.4% | 5.8% | 1.6% | 16.2% | 0.0% |
| Waiting | 4.2% | 76.7% | 14.1% | 4.9% | 0.1% |
| Queuing | 0.0% | 20.8% | 78.7% | 0.2% | 0.3% |
| Walking | 41.6% | 11.6% | 9.5% | 36.8% | 0.0% |
| Talking | 0.0% | 2.7% | 8.2% | 3.4% | 85.7% |

**Average Accuracy: 82.0%**

| | Crossing | Waiting | Queuing | Talking | Dancing | Jogging |
|---|---|---|---|---|---|---|
| Crossing | 76.5% | 6.3% | 1.6% | 0.0% | 0.0% | 15.6% |
| Waiting | 4.8% | 78.5% | 12.8% | 0.9% | 3.1% | 0.0% |
| Queuing | 0.2% | 20.1% | 78.5% | 0.8% | 0.4% | 0.0% |
| Talking | 2.8% | 6.1% | 6.5% | 84.1% | 0.5% | 0.0% |
| Dancing | 11.1% | 5.1% | 2.9% | 0.4% | 80.5% | 0.1% |
| Jogging | 5.9% | 0.0% | 0.0% | 0.0% | 0.0% | 94.1% |

Accuracy vs # of trees

| Dataset | 5 Activities | 6 Activities |
|---|---|---|
| AC [2] | 68.2% | - |
| STV [1] | 64.3% | - |
| STV+MC [1] | 65.9% | - |
| STV + RF | 64.4% | - |
| RSTV | 67.2% | 71.7% |
| RSTV + MRF | **70.9%** | **82.0%** |

- **Activity Segmentation**



- **Example Results**



## Conclusion

- RSTV enables more accurate classification results than state-of-the-art methods
- Capable of handling multiple activities in the scene.
- Enable segmentation of individuals into different collective activities.