

Pose Estimation and Segmentation of People in 3D Movies

Karteeek Alahari

Guillaume Seguin

Josef Sivic

Ivan Laptev

WILLOW Team - Inria / École normale supérieure / CNRS - Paris, France

Goal

Obtain
- pixel-wise **segmentation**
- **pose** estimates
- depth **ordering**
for **multiple people** in **stereo movies**



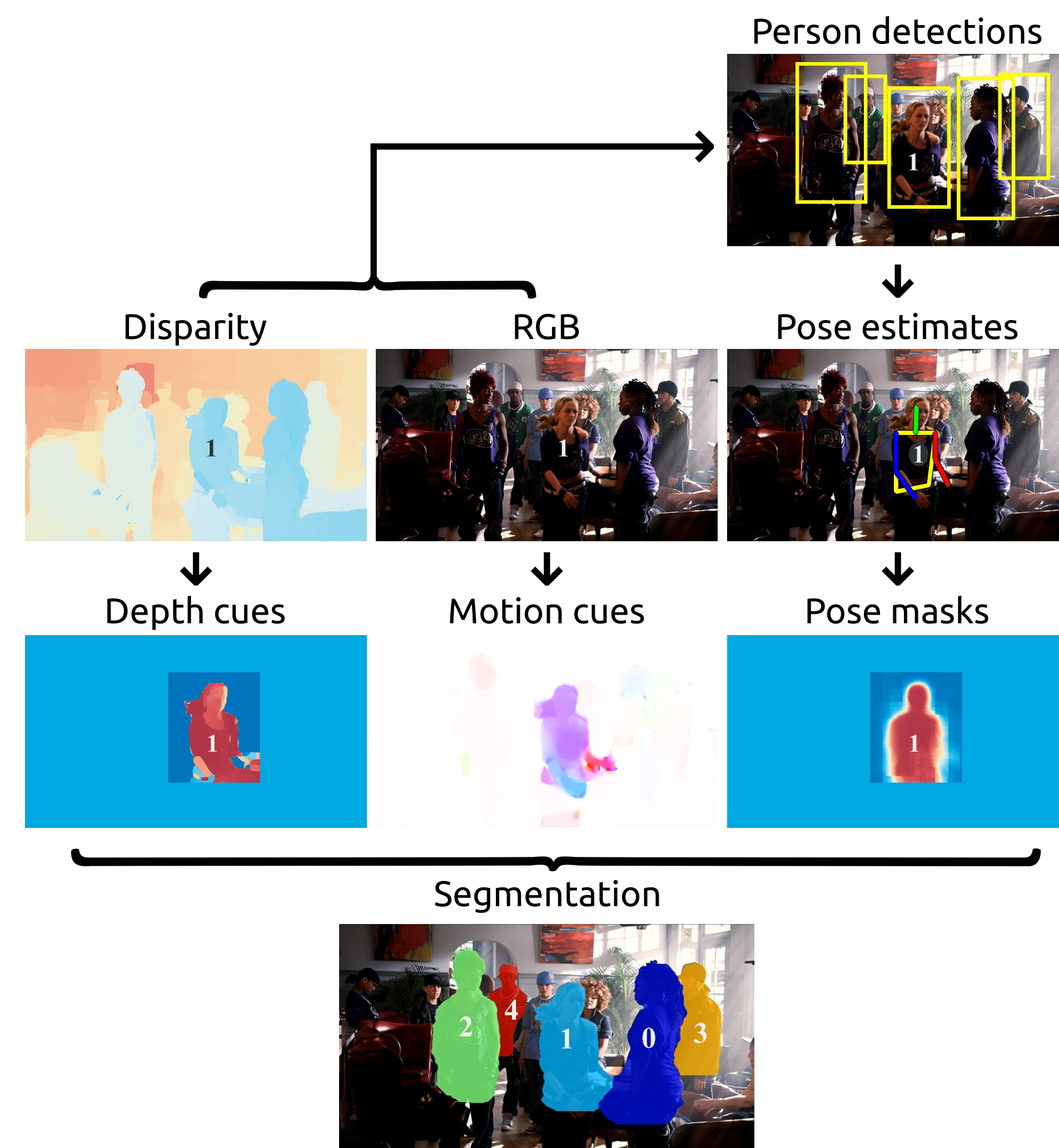
Motivation

- Many movies are now available in 3D (stereo)
- Develop a **mid-level representation** of stereoscopic video for recognition, editing, navigation
- Investigate **benefits of disparity** cues for segmentation and pose estimation
- Collect (noisy) **training data** of segmented people for monocular video

Contributions

1. A **multi-person segmentation model** for stereo videos
- combines multiple cues: disparity, colour, motion
- incorporates learnt **pose-specific** segmentation **masks**
- explicitly represents depth ordering and occlusions
2. A **new annotated Inria 3DMovie Dataset** for person detection, pose estimation and segmentation in indoor and outdoor scenes

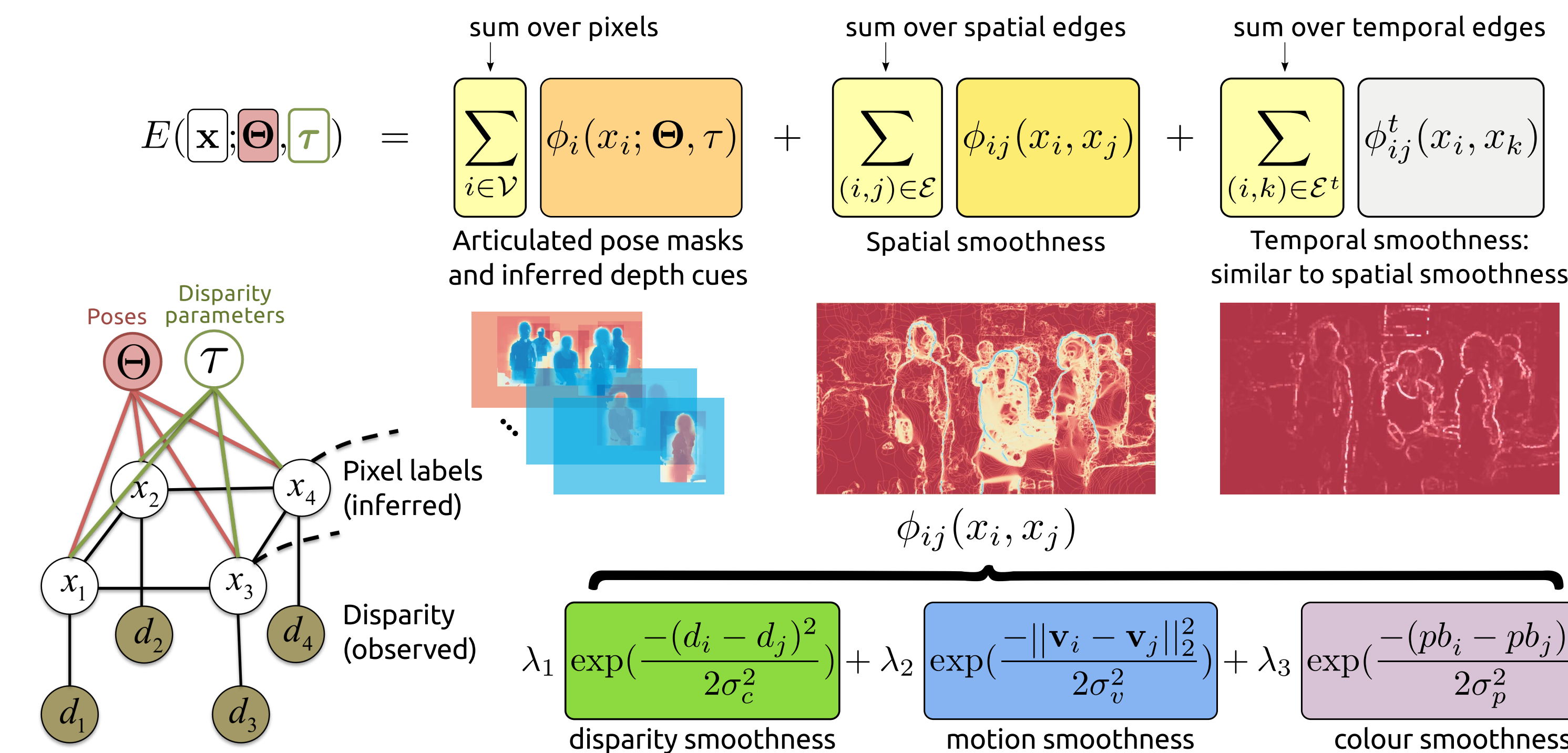
Overview



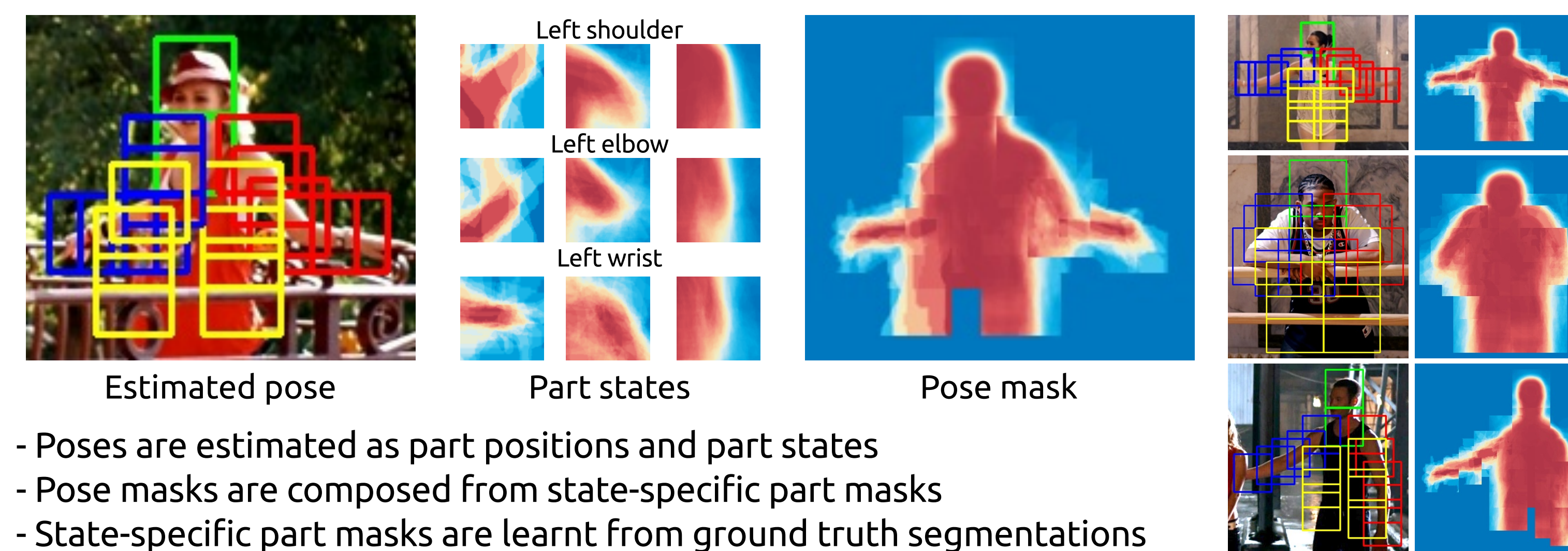
Related work

- [1] Shotton *et al.*, CVPR '11
- [2] Sheasby, Valentin, Crook and Torr, ACCV '12
- [3] Yang, Hallman, Ramanan and Fowlkes, CVPR '10
- [4] Yang and Ramanan, CVPR '11
- [5] Eichner, Marin-Jimenez, Zisserman and Ferrari, IJCV '12

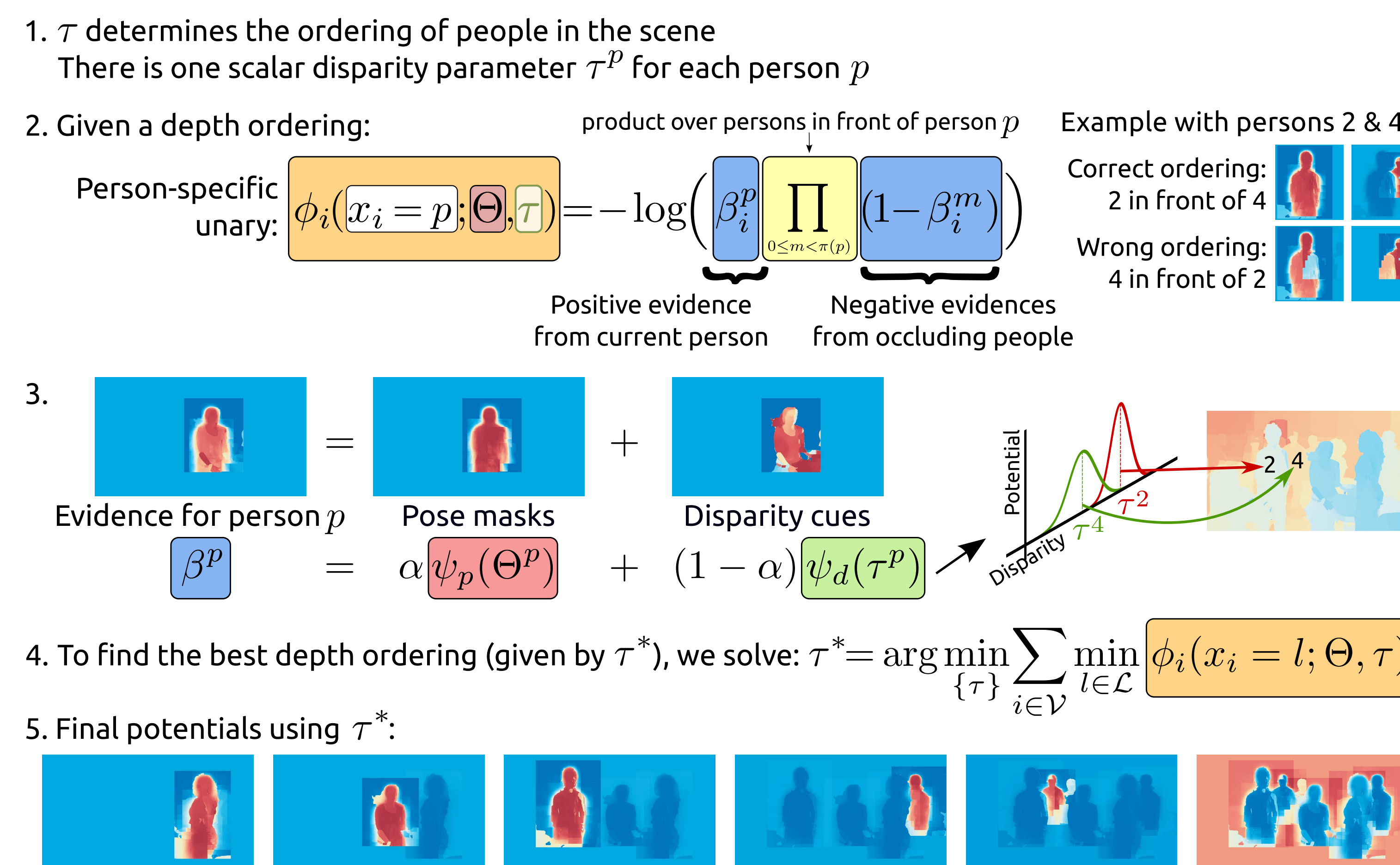
Model



Articulated pose masks



Inferring depth ordering

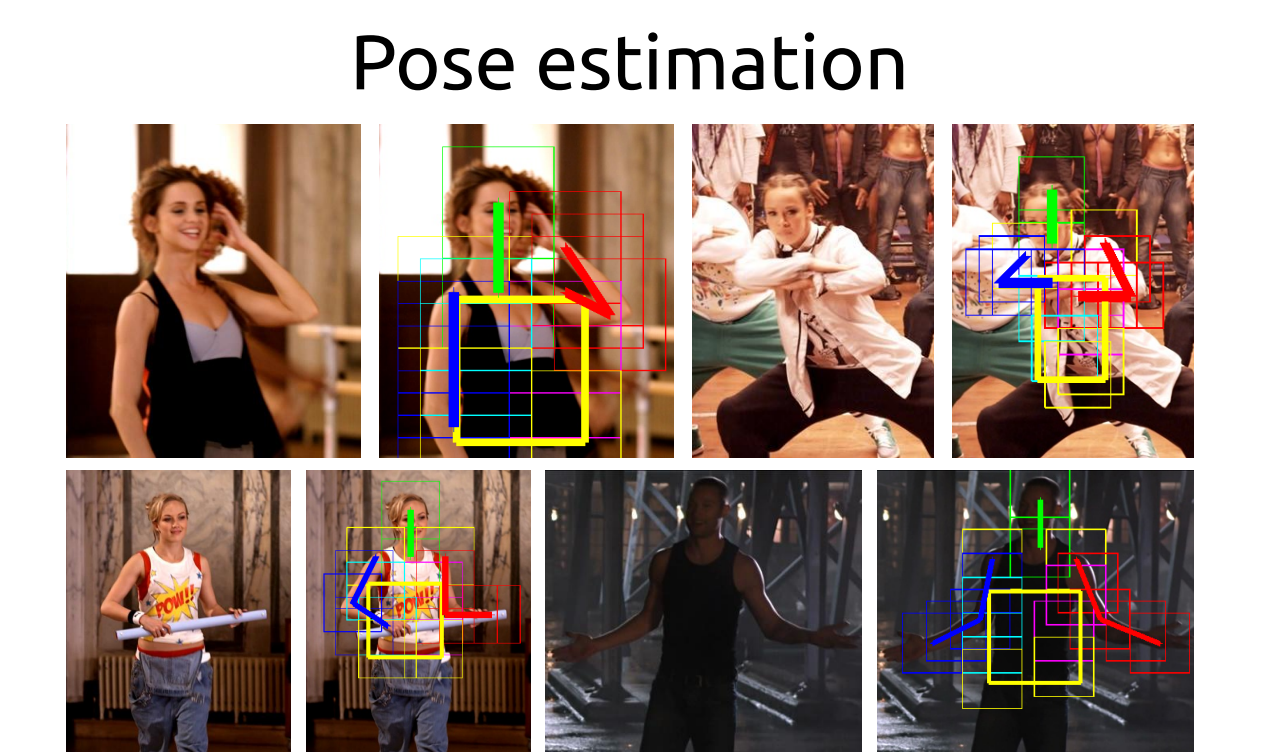
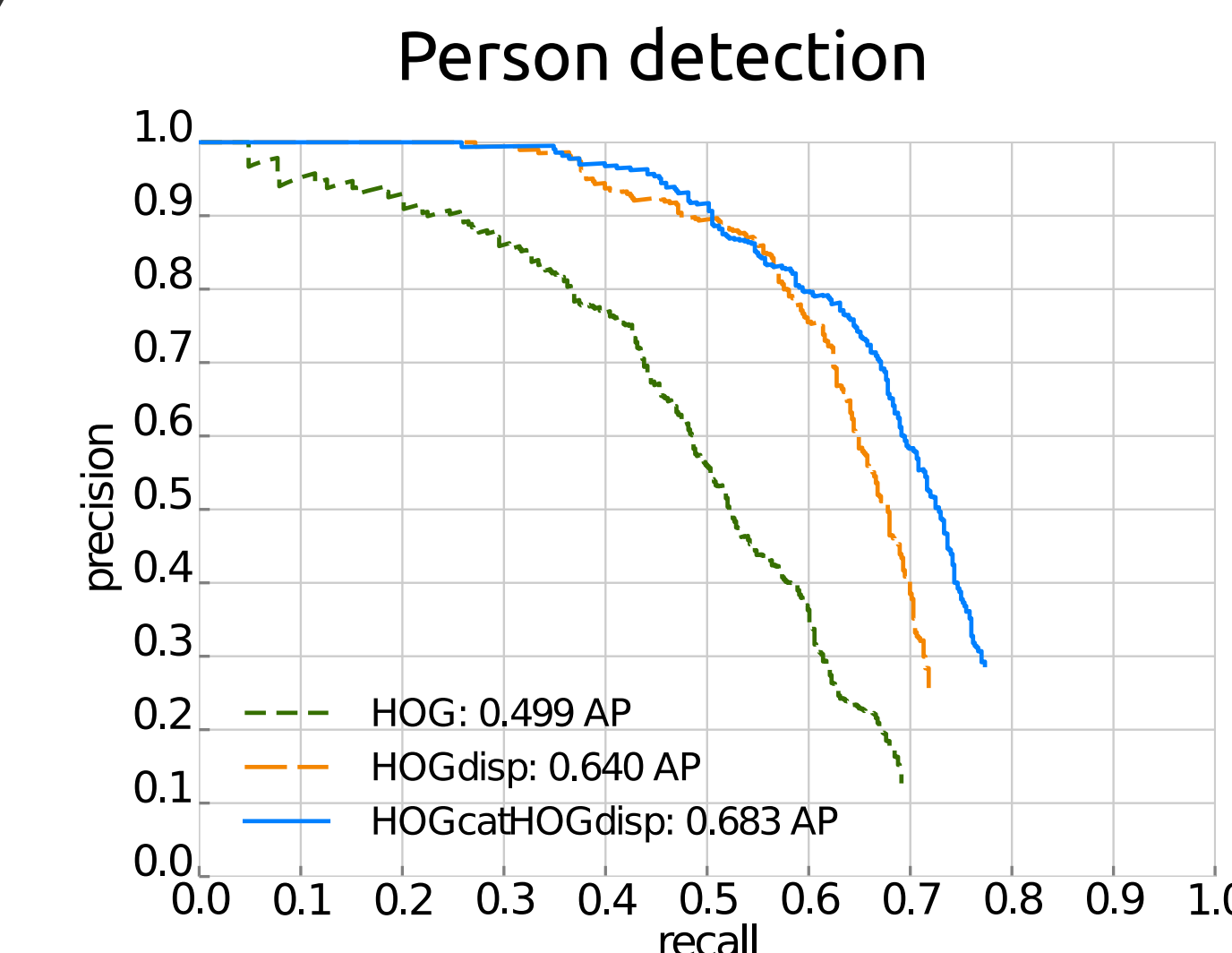


Inria 3DMovie dataset

- Annotated stereo pairs from movies "StreetDance 3D" and "Pina"
 - 440 training stereo pairs ; 36 test video sequences — 2727 pairs
 - **Labelling**: 686 person segmentations, 587 poses, 1158 person boxes
- <http://www.di.ens.fr/willow/research/stereoseg/>



Person detection & Pose estimation

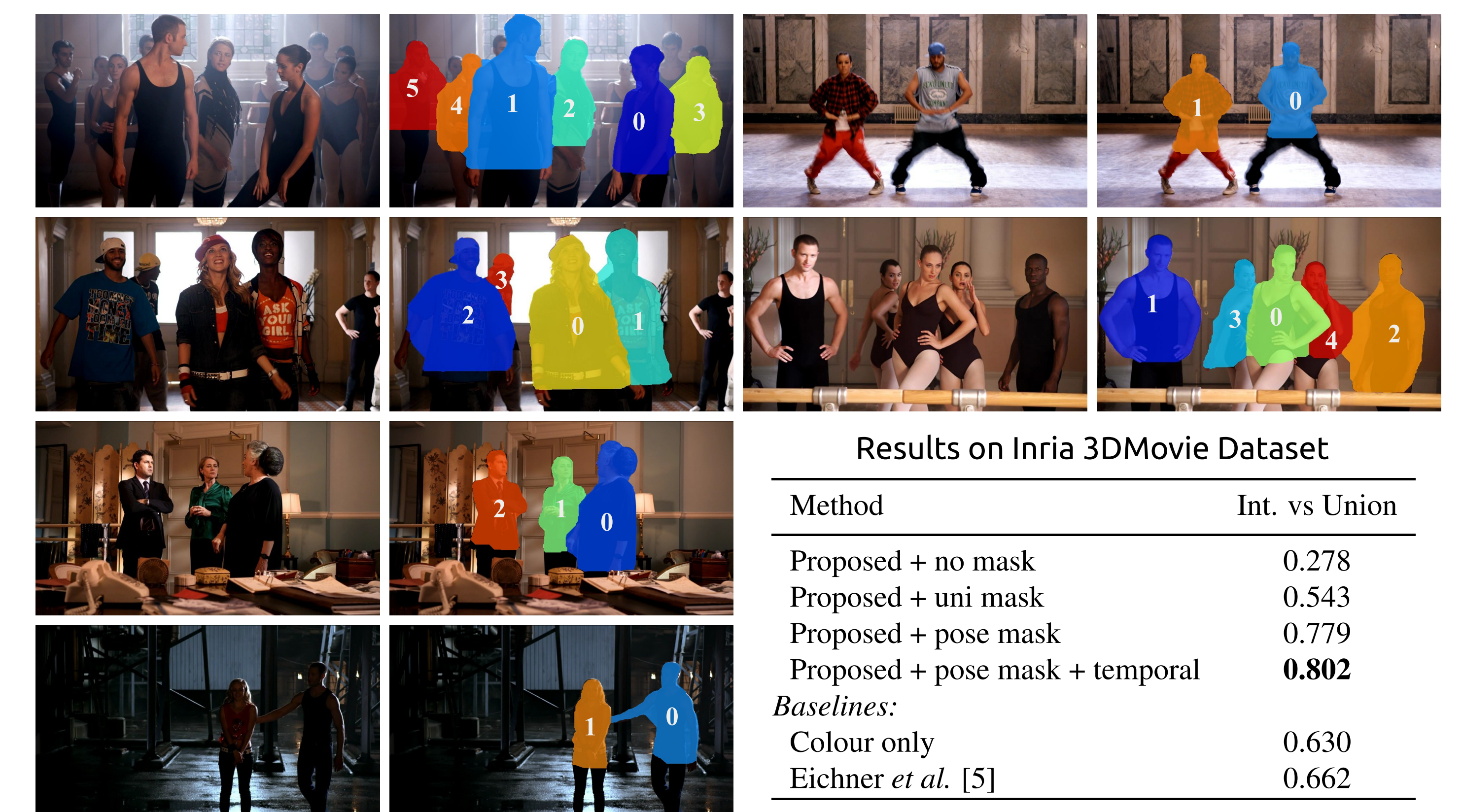


	[4]*	HOG	HOGdisp	HOGcomb
Head & Torso	0.989	0.989	0.991	0.998
Upper arms	0.839	0.856	0.869	0.889
Lower arms	0.518	0.559	0.535	0.594
Global	0.782	0.802	0.799	0.827

* This model was trained on the Buffy dataset
We use the PCP measure from Eichner *et al.* [5] which evaluates the percentage of correctly estimated body parts

- HOG: HOG on RGB only
- HOGdisp: HOG on disparity only
- HOGcatHOGdisp: concatenation of both

Person segmentation



Results on H2view [2]

Method	Int. vs Union
Upper body segmentation:	
Sheasby <i>et al.</i> [2]	0.735
Proposed + pose mask	0.814
Proposed + pose mask + temporal	0.825
Full body segmentation:	
Sheasby <i>et al.</i> [2]	0.692
Proposed + upper pose mask	0.706

Evaluation: the intersection vs. union measure between each detected person segmentation and the corresponding ground truth