# Deep neural nets for human pose estimation in videos

Tomas Pfister, James Charles, Andrew Zisserman
Department of Engineering Science
University of Oxford
http://www.robots.ox.ac.uk/~vgg

# Aim:

Estimate 2D upper body joint positions (wrist, elbow, shoulder, head) with high accuracy in real-time

# Outline

- Two types of loss functions for pose estimation

    - Coordinate net

    - Heatmap net

- Optical flow for pose estimation in videos

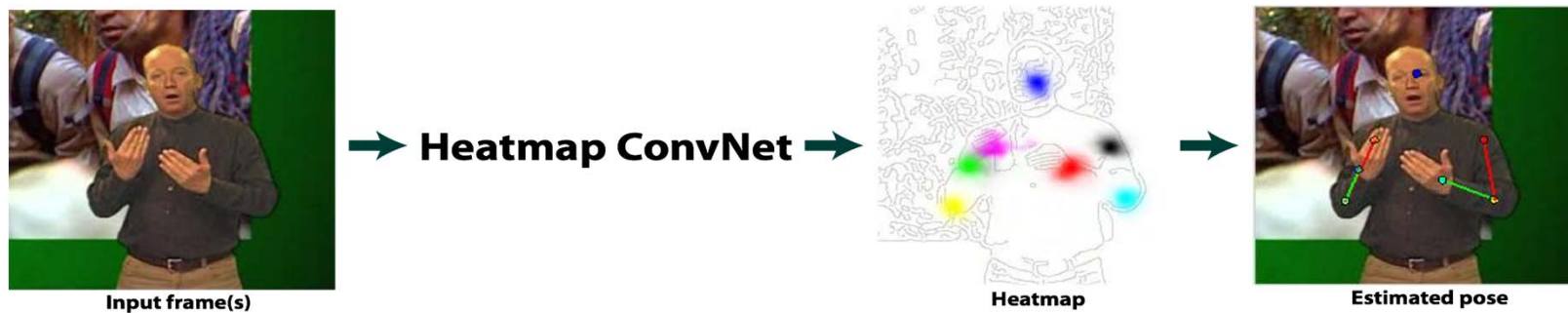- Results (cf state of the art)

# Method overview: single frame learning

## 1. Coordinate Net



**Coordinate ConvNet**

Input frame(s) → Estimated pose

e.g. DeepPose CVPR14, Pfister et al ACCV14

## 2. Heatmap Net



**Heatmap ConvNet**

Input frame(s) → Heatmap → Estimated pose

e.g. Jain et al ICLR14, Tompson et al CVPR15

# Coordinate Net: regress joint positions



**Pose ConvNet**

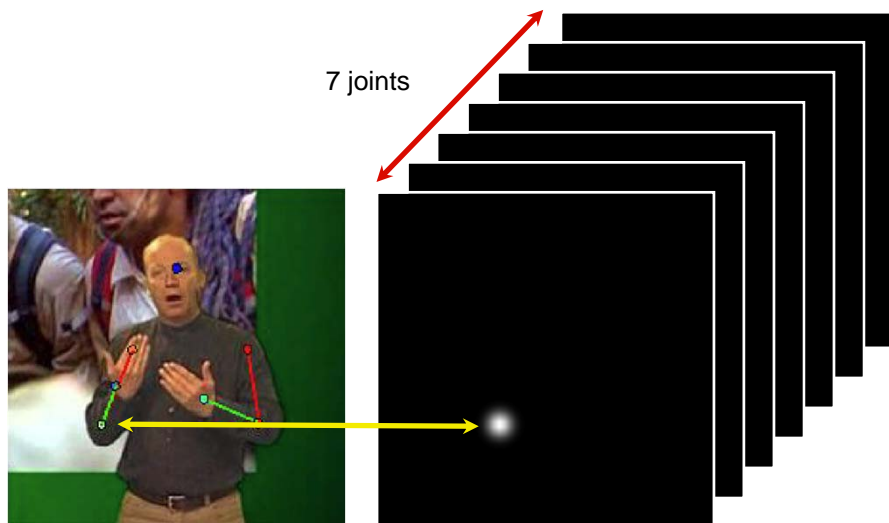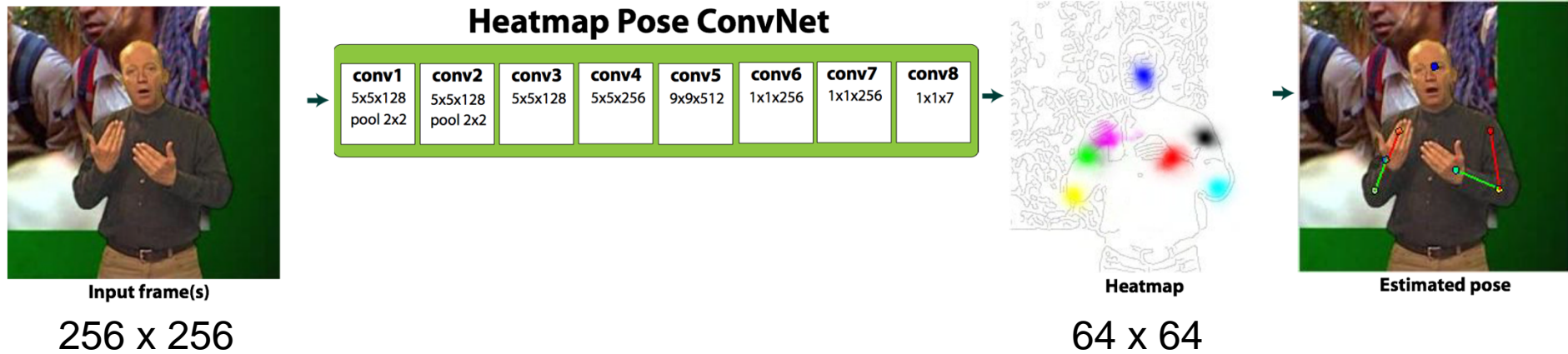| conv1 | conv2 | conv3 | conv4 | conv5 | full6 | full7 | regress |
|-------|-------|-------|-------|-------|-------|-------|---------|
| 7x7x96 norm pool 3x3 | 5x5x256 pool 2x2 | 3x3x512 | 3x3x512 | 3x3x512 pool 3x3 | 4096 dropout | 4096 dropout | 14 |

Input frame

Estimated pose

Training loss: L2 on joint positions

OverFeat like architecture
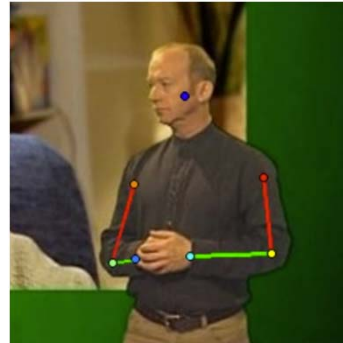
# Heatmap Net: regress heatmap for each joint



**Heatmap Pose ConvNet**

| conv1 5x5x128 pool 2x2 | conv2 5x5x128 pool 2x2 | conv3 5x5x128 | conv4 5x5x256 | conv5 9x9x512 | conv6 1x1x256 | conv7 1x1x256 | conv8 1x1x7 |

Input frame(s)

256 x 256

Heatmap

64 x 64

Estimated pose

7 joints

Represent joint position by Gaussian

Training loss: L2 on pixels

# Comparison

**Regression target**



**Coordinate Net**     Coordinates

**Heatmap Net**     Heatmap

# BBC sign language videos data set

**Training:**
15 videos each 0.5-1hr long, all frames annotated

**Testing:**
5 videos, 200 annotated frames per video

**Extended Training:**
72 videos with noisy automated annotations

Training set

Testing set
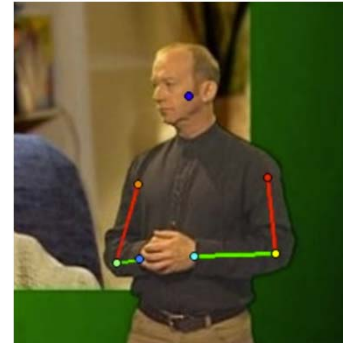
# Results on architecture comparison



- Heatmap net superior to coordinate net
- Performance of coordinate net saturates with more training data

**Evaluated on BBC Pose**

# Why is the heatmap network superior?

1. Can represent multimodal estimates, so can model uncertainty/confidence

2. In training there is an error signal from every pixel, so better smoothing for back propagation

Also, it is easier to visualize (and understand) what is being learnt
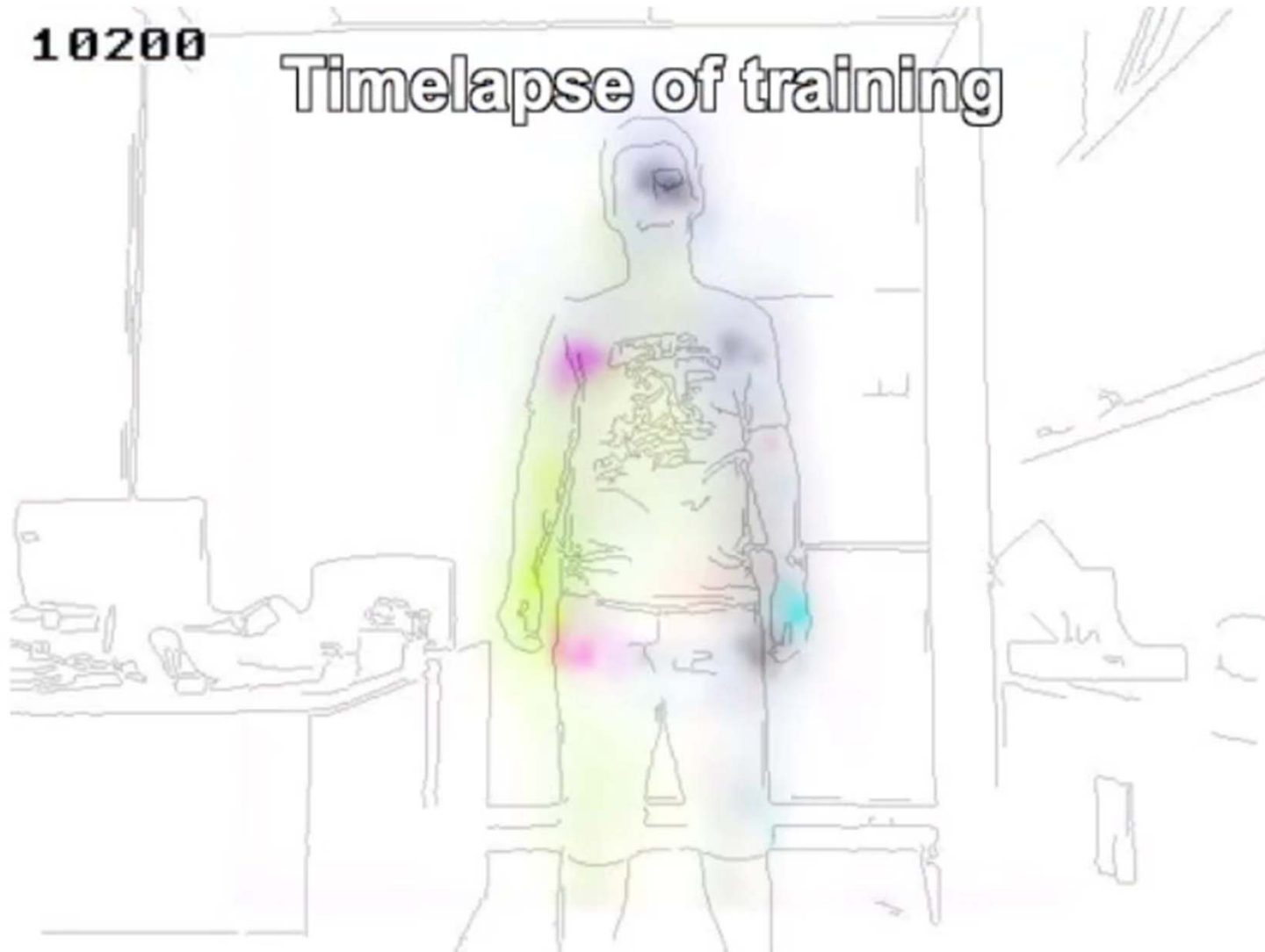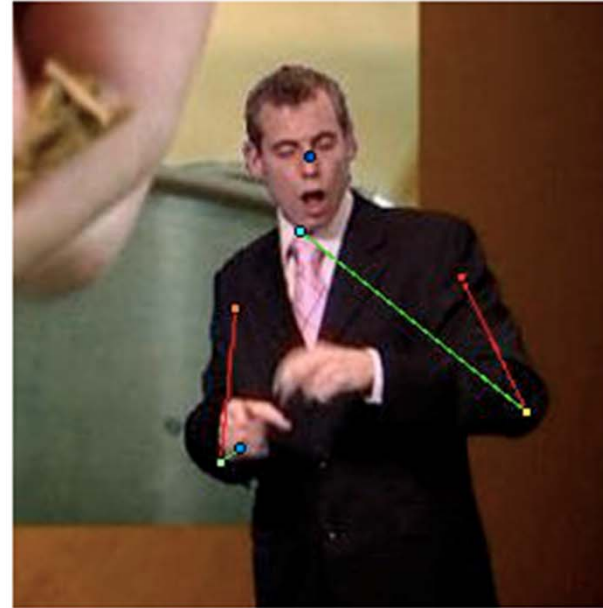
Regression target
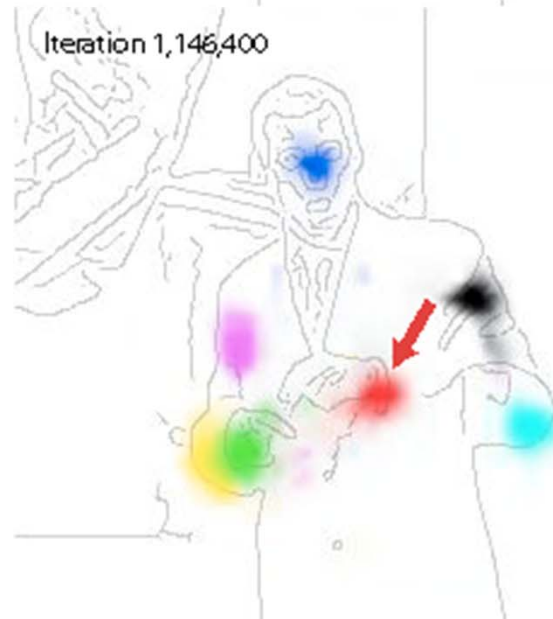


Coordinate Net

Coordinates



Heatmap Net

Heatmap

Timelapse of training

# Multiple modes example



Iteration 262,400

Iteration 1,146,400

early in training

late in training

# What do the layers learn?

**Three randomly selected activations from each layer**



Input frame

Edges

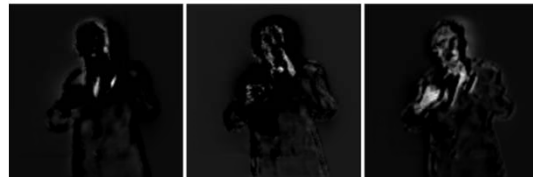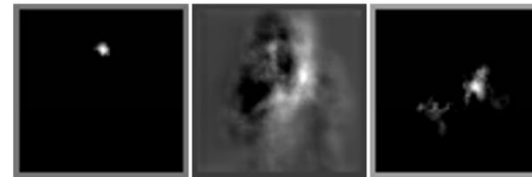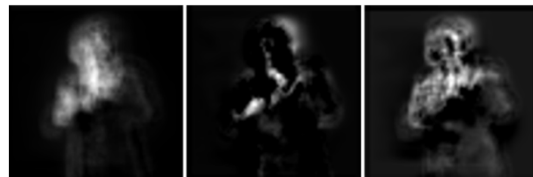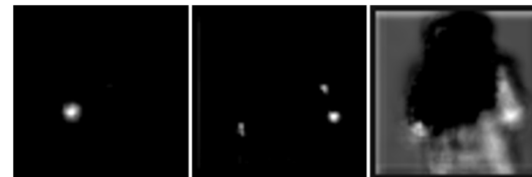Body parts (some)

conv1

conv2

conv3

conv4

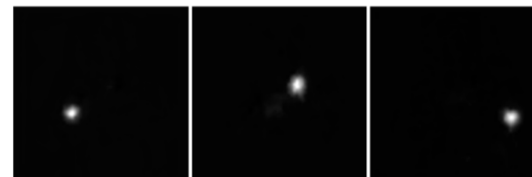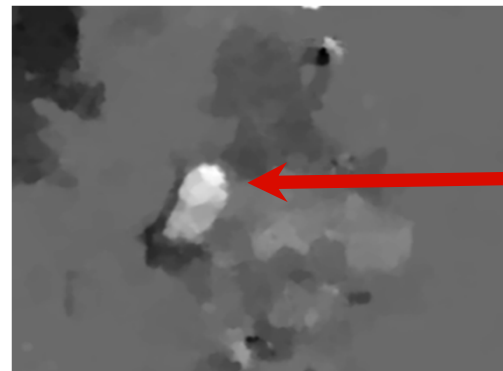conv5

conv6

conv7

conv8 (output)

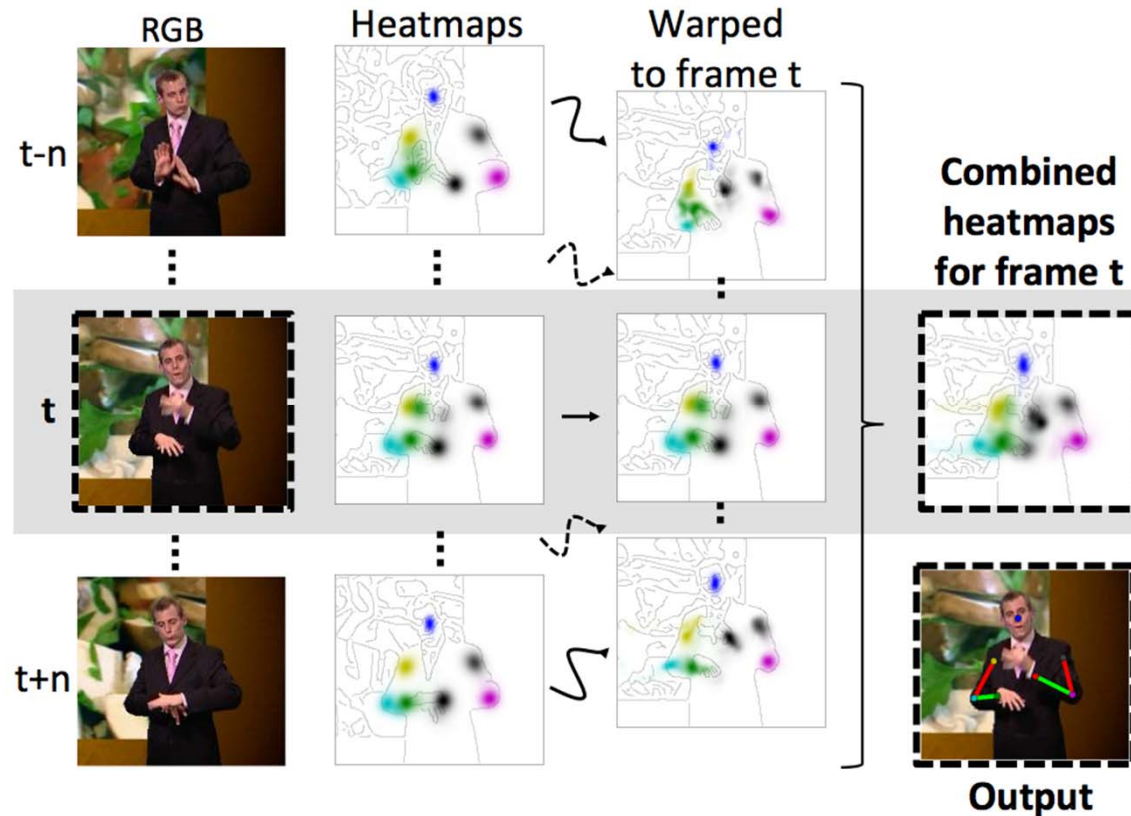# Learning from videos

• Temporal information



– How do we learn from temporal information with a ConvNet?

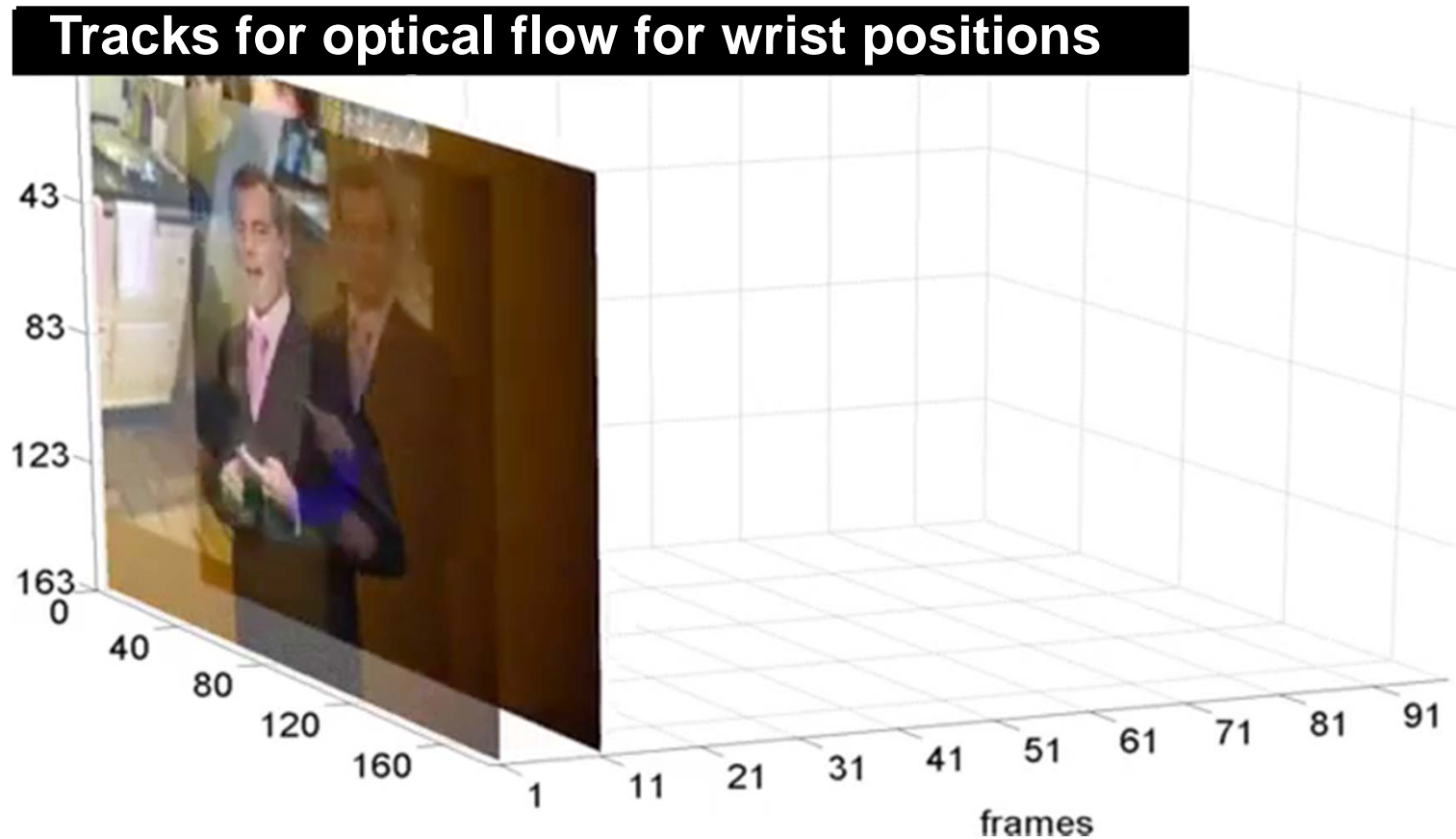

Hand moving
in x direction

# Late fusion using flow

**Warp the heatmaps** from previous/next frames & **combine**



Cf S. Zuffi et al., Estimating human pose with flowing puppets. Proc. ICCV, 2013
Charles et al., Upper Body Pose Estimation with Temporal Sequential Forests, BMVC 2014
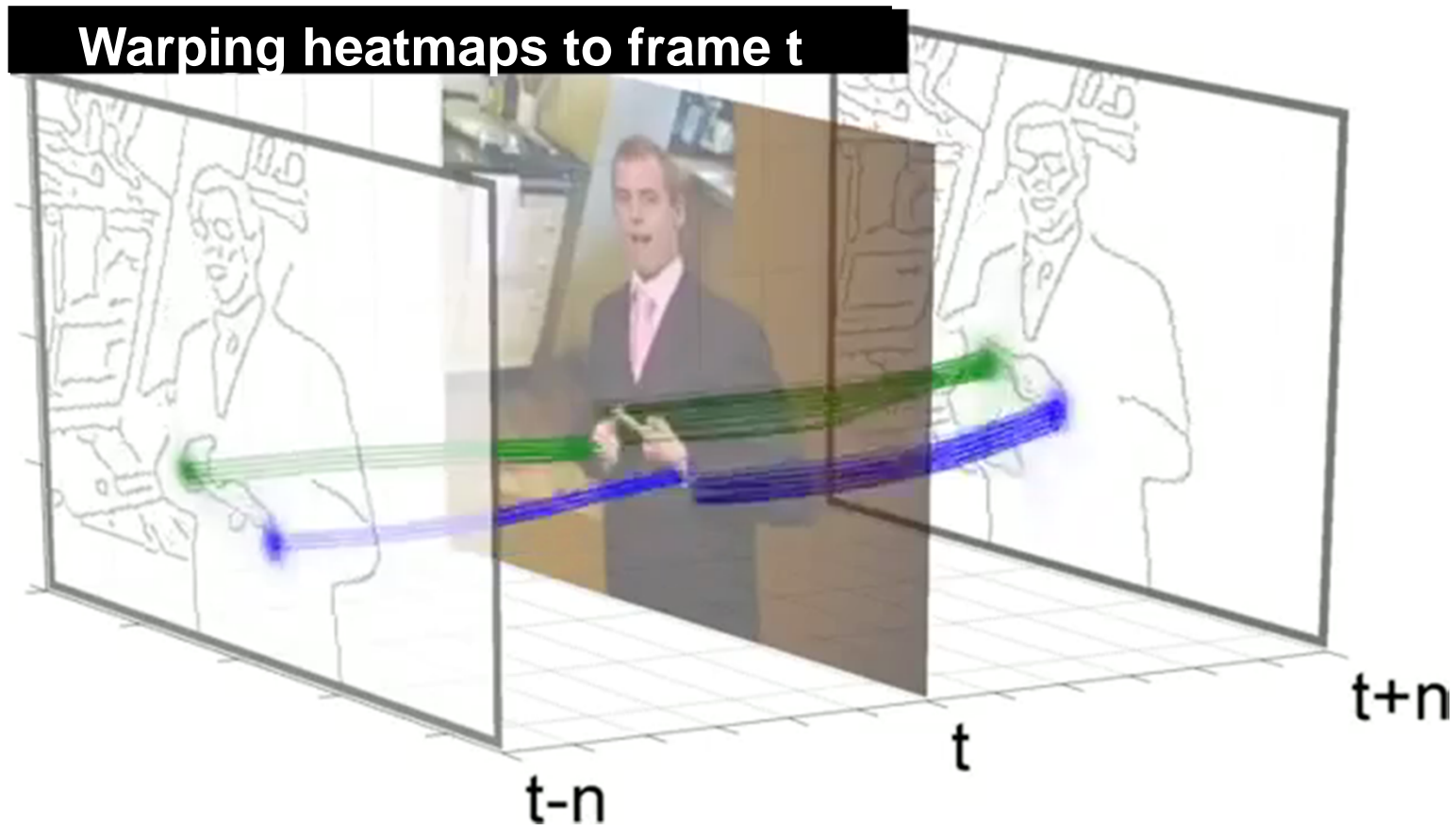
# Optical flow
# Example: Heatmap Net & Optical flow



**Tracks for optical flow for wrist positions**

Flow: Brox et al GPU flow from OpenCV, or FastDeepFlow

# Optical flow
# Example: Heatmap Net & Optical flow
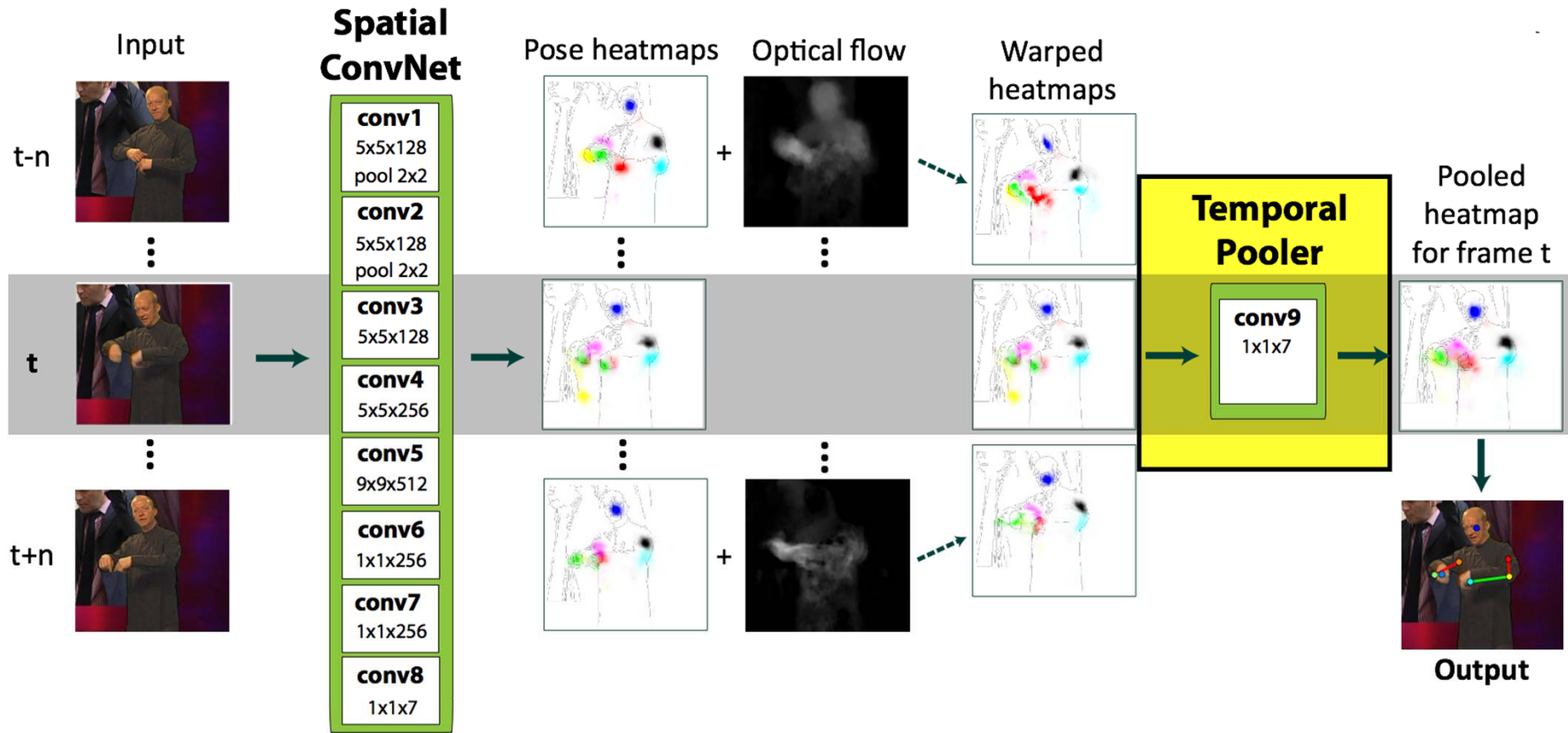


Warping heatmaps to frame t

# Flowing ConvNets

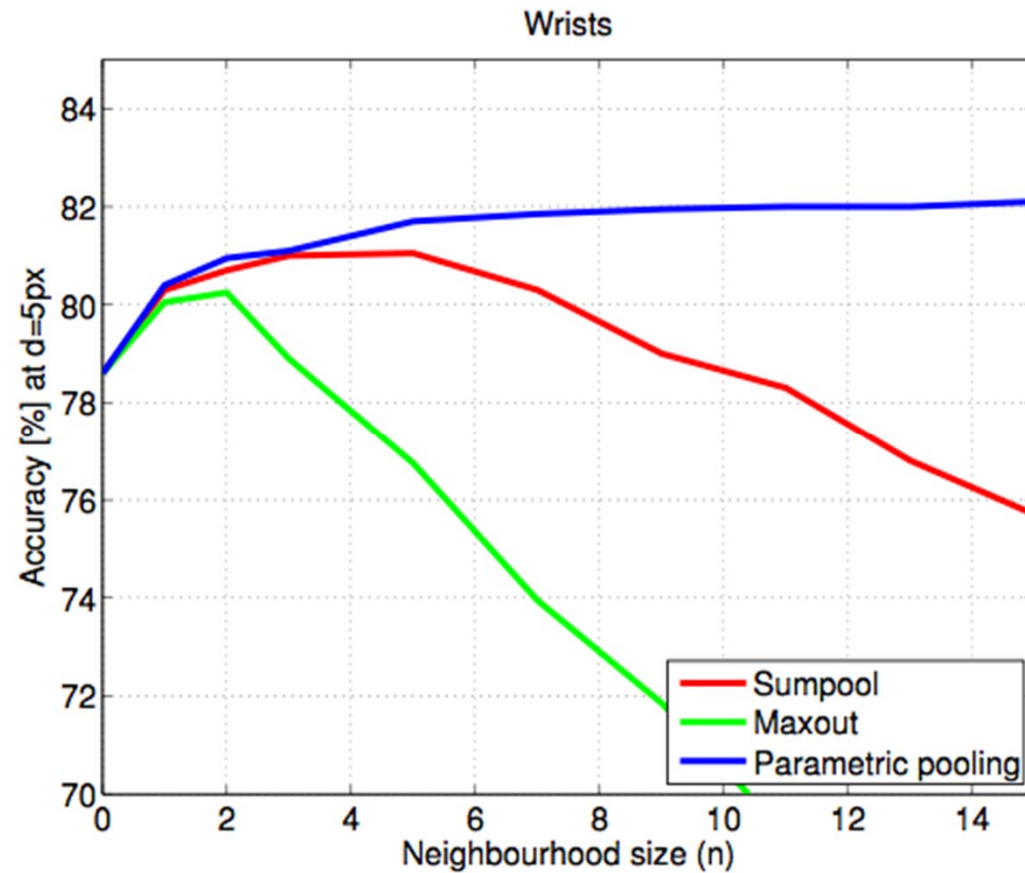- Learn the **pooling** of the warped heatmaps

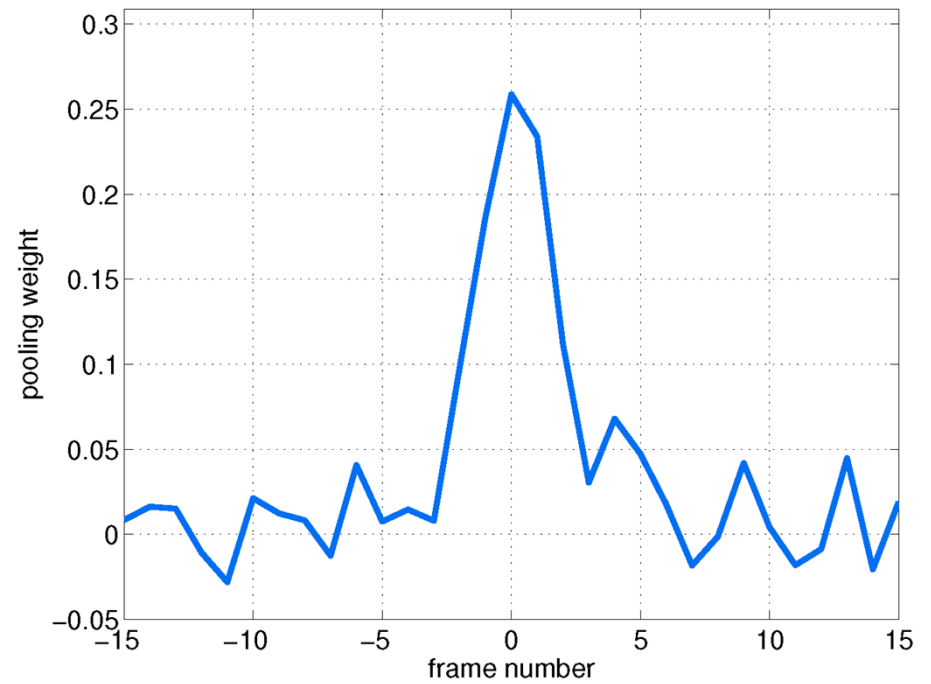# Results: with/without optical flow

# Results
# Comparison of pooling types

# Results
# Learnt optical flow pooling weights
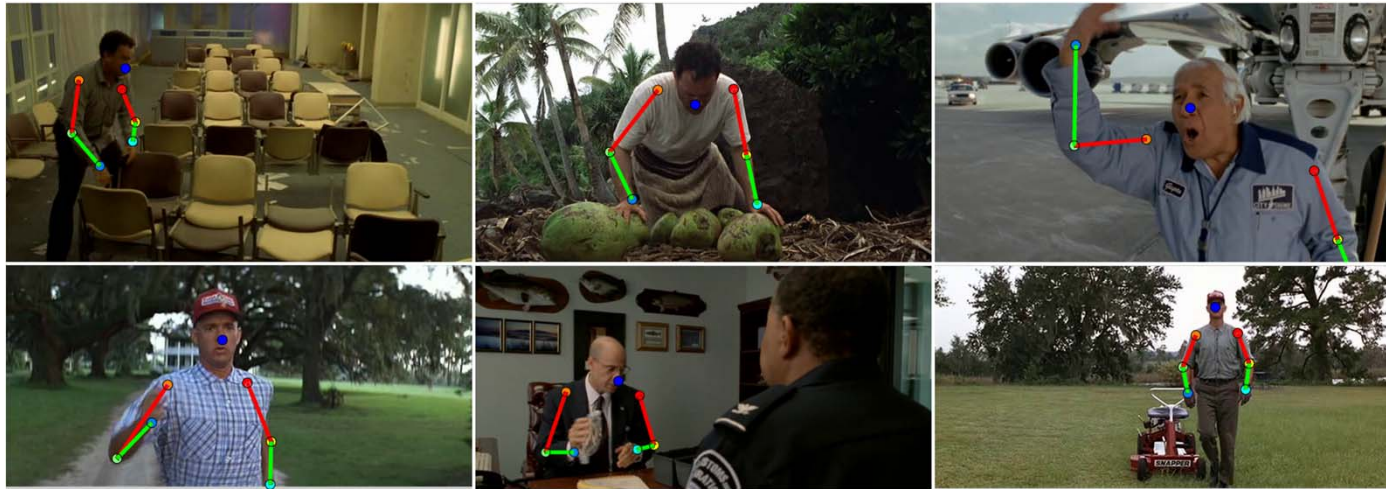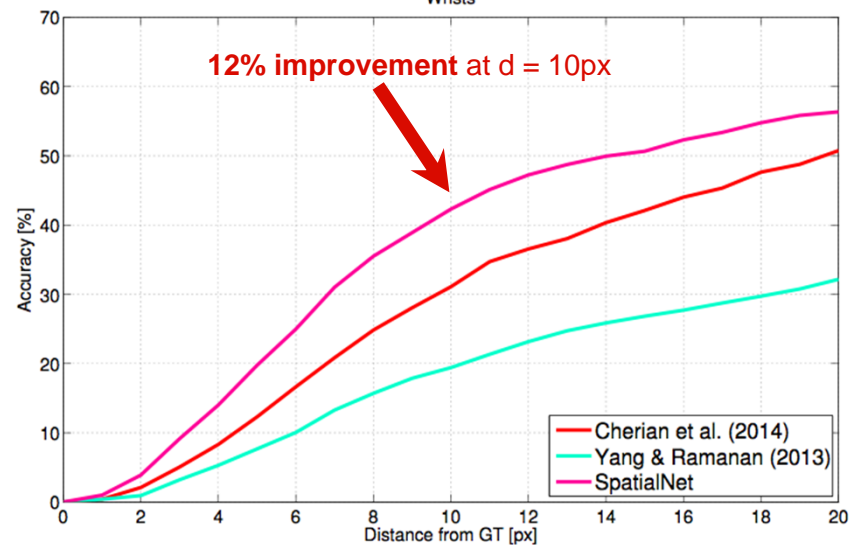
# Results
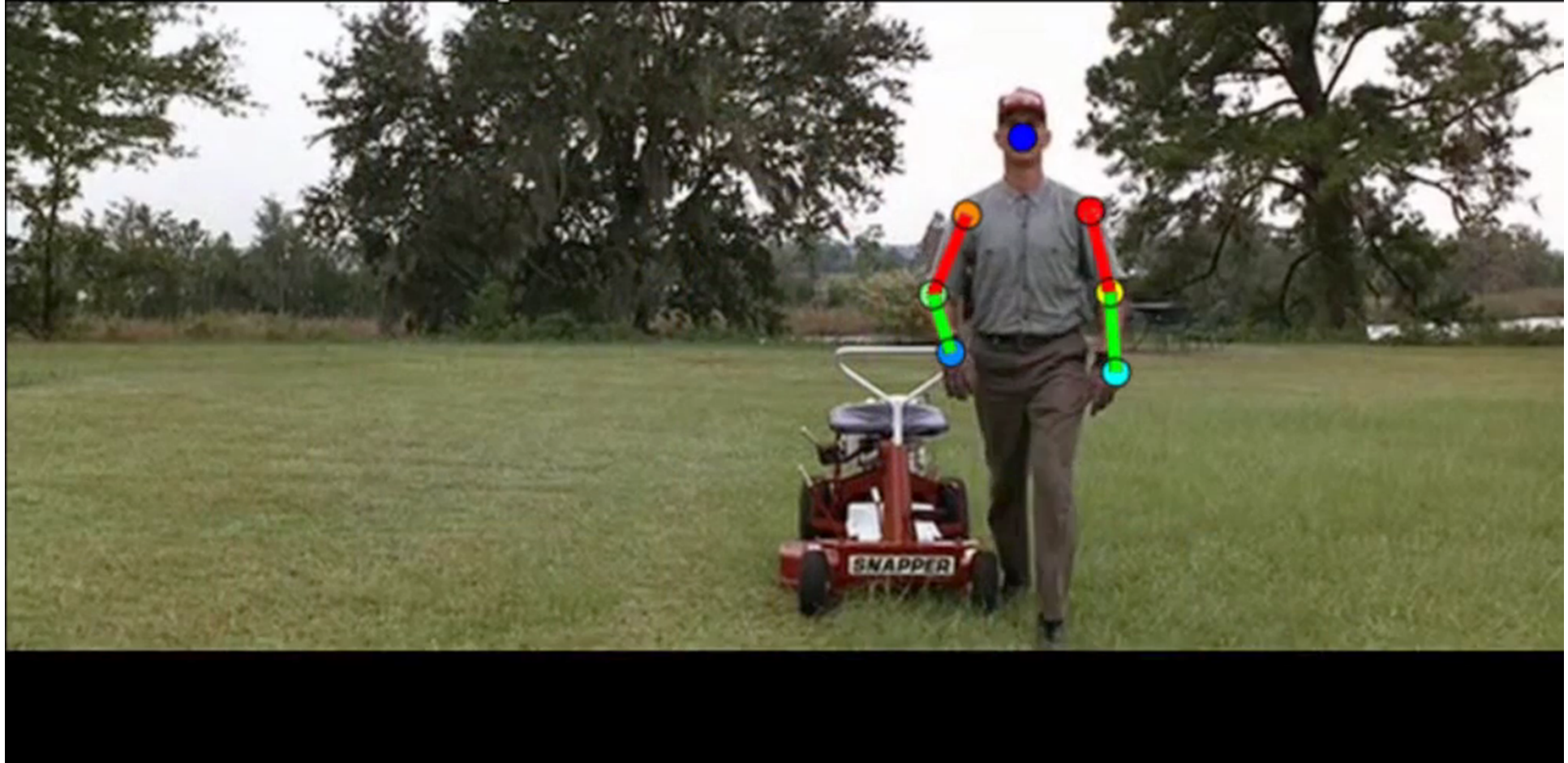# Comparison to the state of the art

**Poses in the Wild**

# Results: Example pose estimation



Output: Poses in the Wild

50fps on 1 GPU without optical flow, 5fps with optical flow

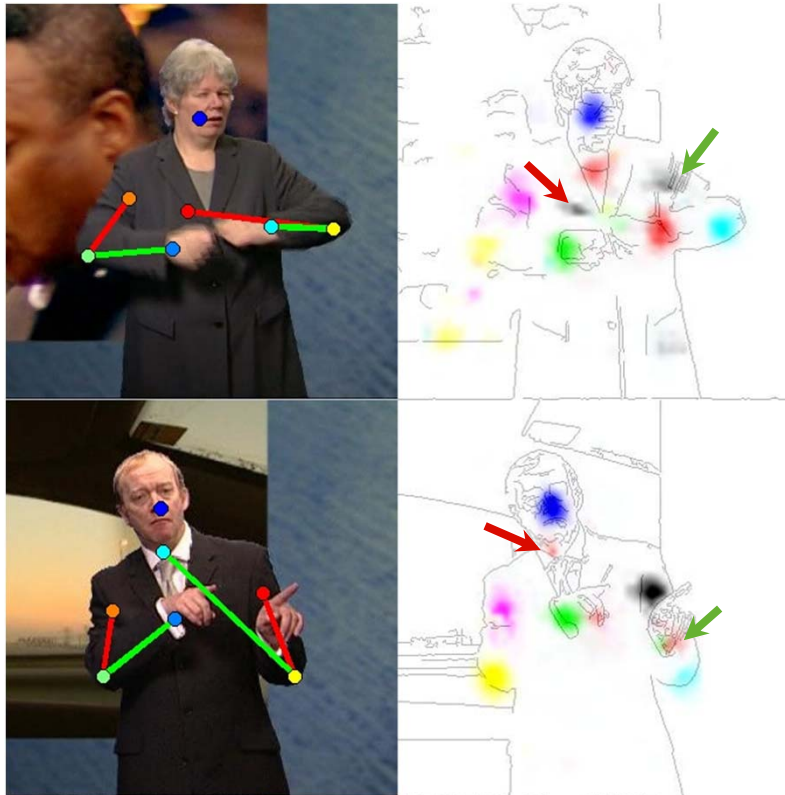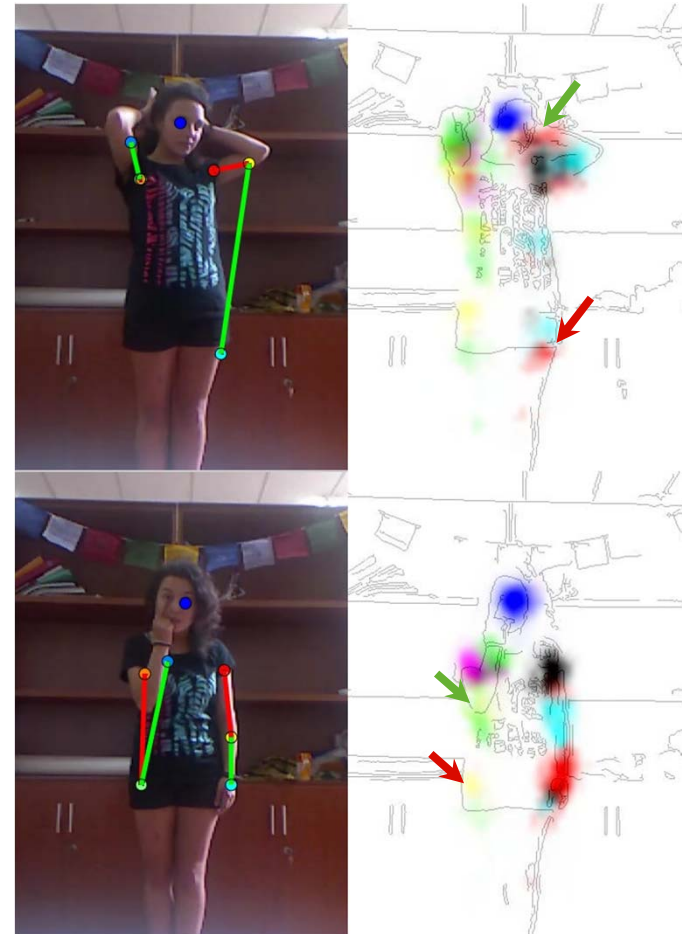# Results
# Failure cases
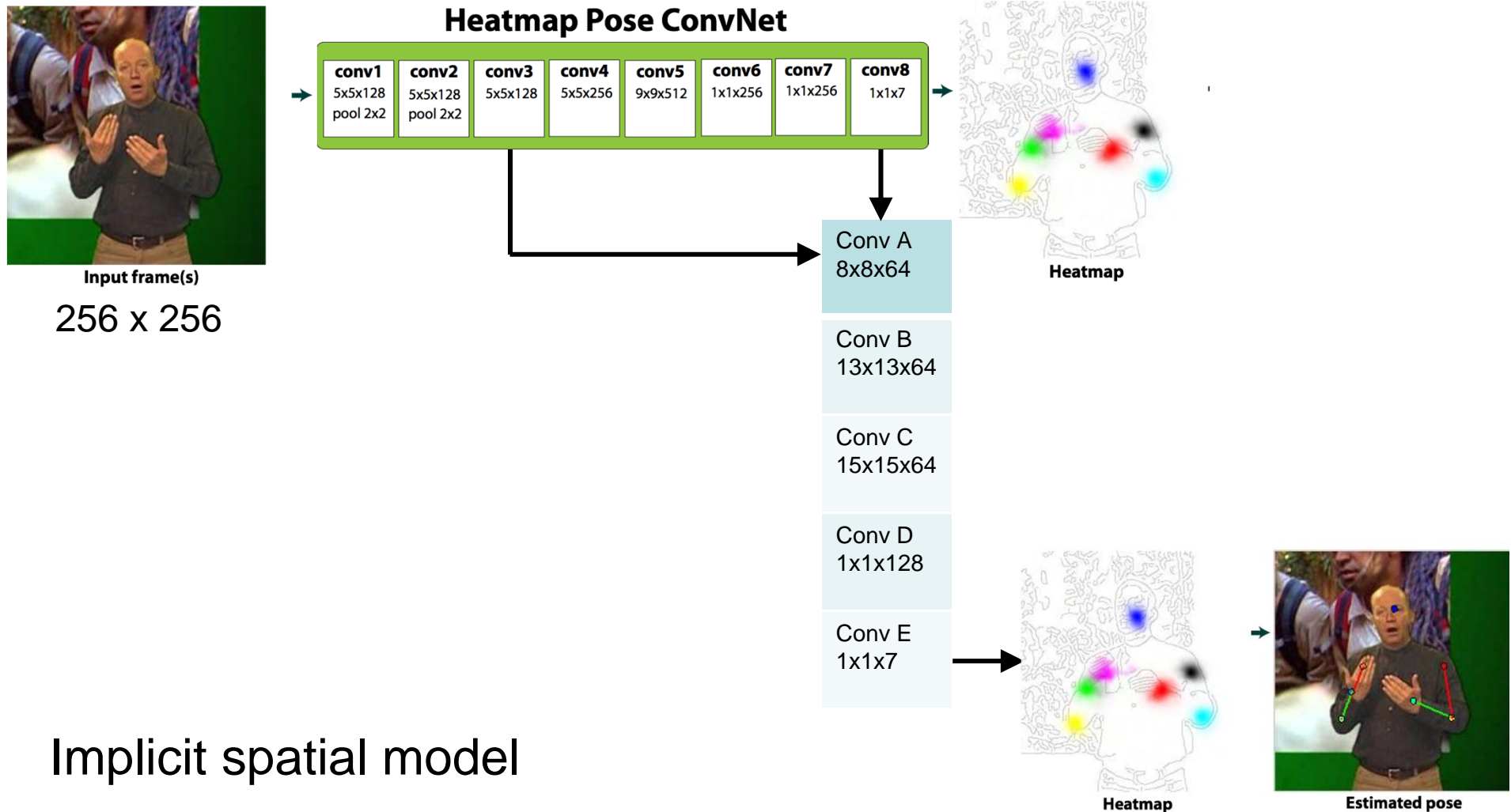
**Main failure case:** Picking the wrong mode

**BBC Pose**



Correctable with a **spatial model**

**ChaLearn**

# Additional Pooling Fusion Layers



256 x 256

Implicit spatial model

# Results: Additional Pooling Fusion Layers



Poses in the Wild

Heat map CNNs

Wrists

with fusion and flow

with fusion

original

Cherian et al. (2014)
Yang & Ramanan (2013)
SpatialNet Fusion Flow
SpatialNet Fusion
SpatialNet

Accuracy [%]

Distance from GT [px]

# Results: Additional Pooling Fusion Layers

**FLIC:** single image predictions



Average PCK for wrist & elbow

Legend:
- Toshev et al.
- Jain et al.
- MODEC
- Yang et al.
- Sapp et al.
- Tompson et al. *
- Chen et al. *
- Ours

# Summary

- Deep Heatmap ConvNet achieves state of the art with implicit spatial models

- Performance improved by optical flow pooling

- Futures:
  - Robust regression
  - Data dependent flow channel pooling
  - More training data