

# LOCALIZATION OF HUMANS IN IMAGES USING CONVOLUTIONAL NETWORKS

---

JONATHAN TOMPSON

ADVISOR: CHRISTOPH BREGLER



# OVERVIEW

Thesis goal:

State-of-the-art human pose recognition



from depth  
(hand tracking)



from RGB  
(body tracking)

Domain Specific Optimizations:

Kinematic structure in ConvNet (hybrid Graphical Model + ConvNet)

Motion feature inputs

Novel architecture to recover lost spatial precision due to pooling

# TALK OUTLINE

## 1. Hand Tracking

- a) J. Tompson, M. Stein, Y. LeCun, K. Perlin, "Real-Time Continuous Pose Recovery of Human Hands Using Convolutional Networks", ACM TOG/SIGGRAPH 2014

## 2. Body Tracking

- a) J. Tompson, A. Jain, Y. LeCun, C. Bregler, "Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation", NIPS 2014
- b) A. Jain, J. Tompson, Y. LeCun, C. Bregler, "MoDeep: A Deep Learning Framework Using Motion Features for Human Pose Estimation", ACCV 2014
- c) J. Tompson, R. Goroshin, A. Jain, Y. LeCun, C. Bregler, "Efficient Object Localization Using Convolution Networks", CVPR 2015

# TALK OUTLINE

## 1. Hand Tracking

- a) J. Tompson, M. Stein, Y. LeCun, K. Perlin, "Real-Time Continuous Pose Recovery of Human Hands Using Convolutional Networks", ACM TOG/SIGGRAPH 2014

## 2. Body Tracking

- a) J. Tompson, A. Jain, Y. LeCun, C. Bregler, "Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation", NIPS 2014
- b) A. Jain, J. Tompson, Y. LeCun, C. Bregler, "MoDeep: A Deep Learning Framework Using Motion Features for Human Pose Estimation", ACCV 2014
- c) J. Tompson, R. Goroshin, A. Jain, Y. LeCun, C. Bregler, "Efficient Object Localization Using Convolution Networks", CVPR 2015

# HAND POSE INFERENCE

Target: low-cost markerless mocap

Full articulated pose with high DoF

Real-time with low latency

## Challenges

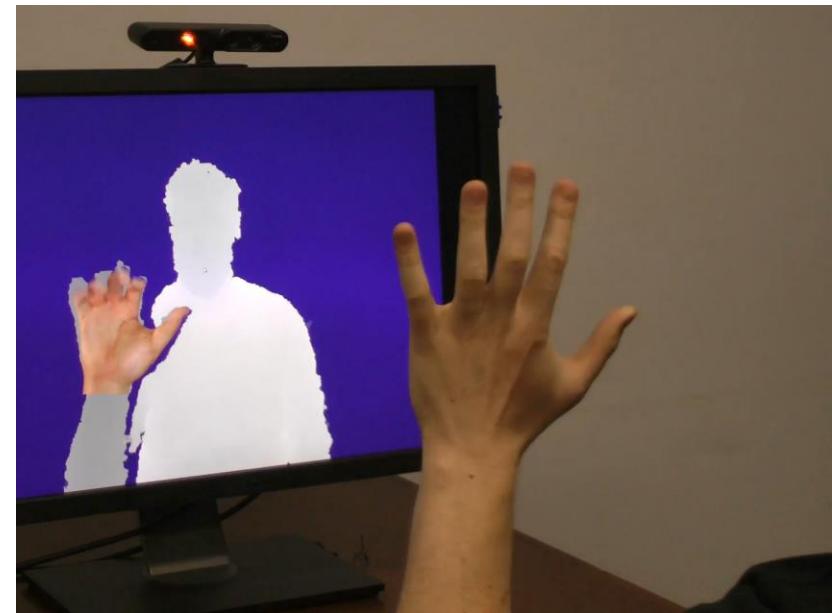
Many DoF contribute to model deformation

Constrained unknown parameter space

Self-similar parts

Self occlusion

Device noise



# PIPELINE

## Supervised learning based approach

Needs labeled dataset + machine learning

Existing datasets had limited pose information for hands

## Architecture

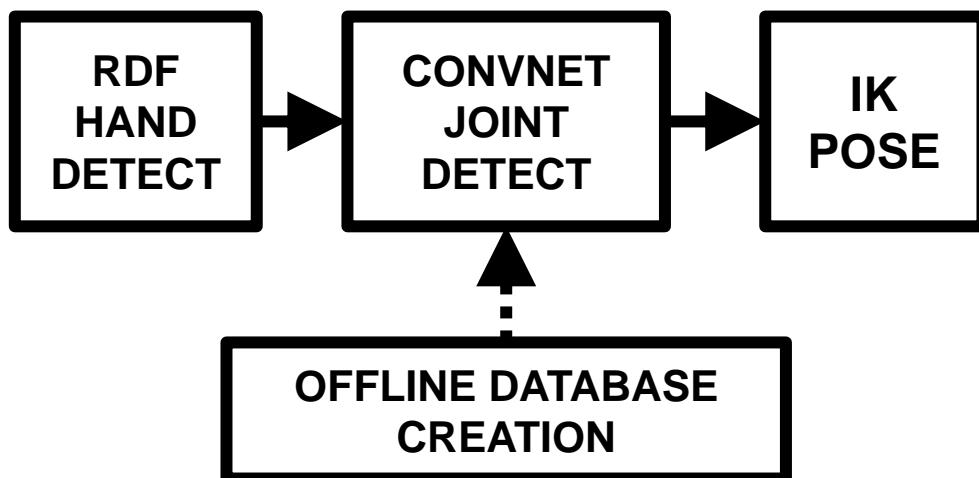
# PIPELINE

Supervised learning based approach

Needs labeled dataset + machine learning

Existing datasets had limited pose information for hands

Architecture



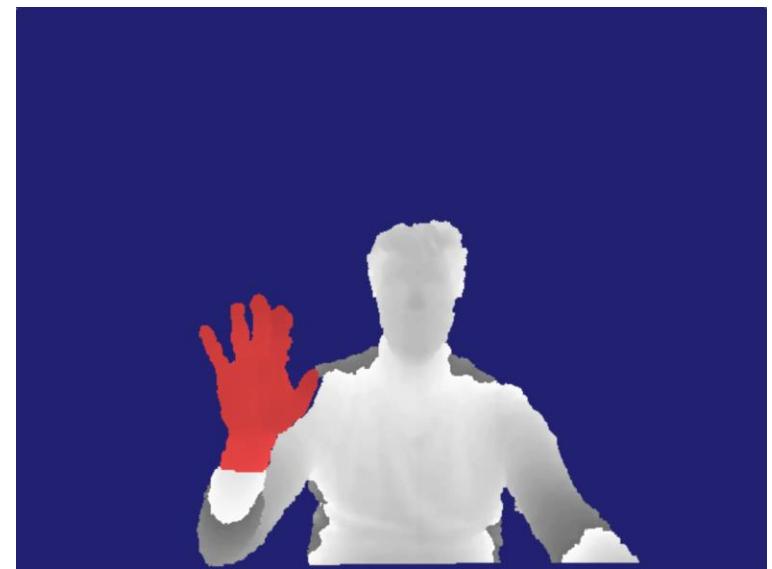
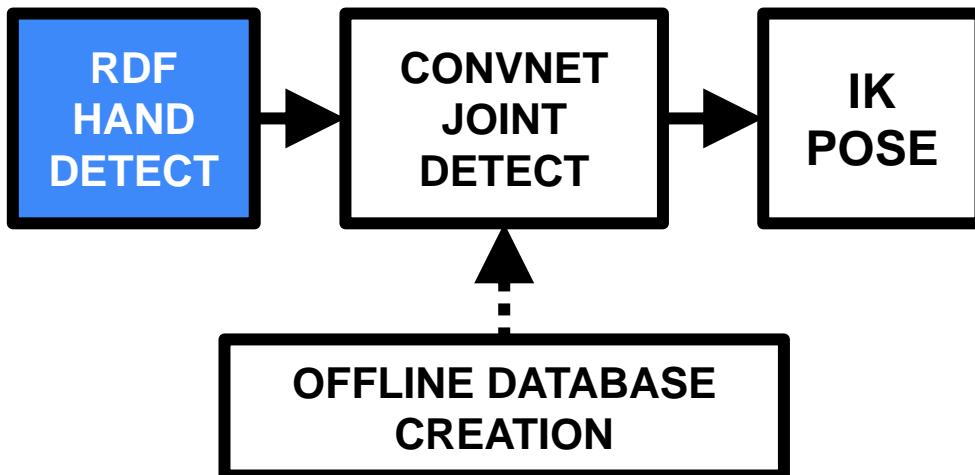
# PIPELINE

Supervised learning based approach

Needs labeled dataset + machine learning

Existing datasets had limited pose information for hands

Architecture



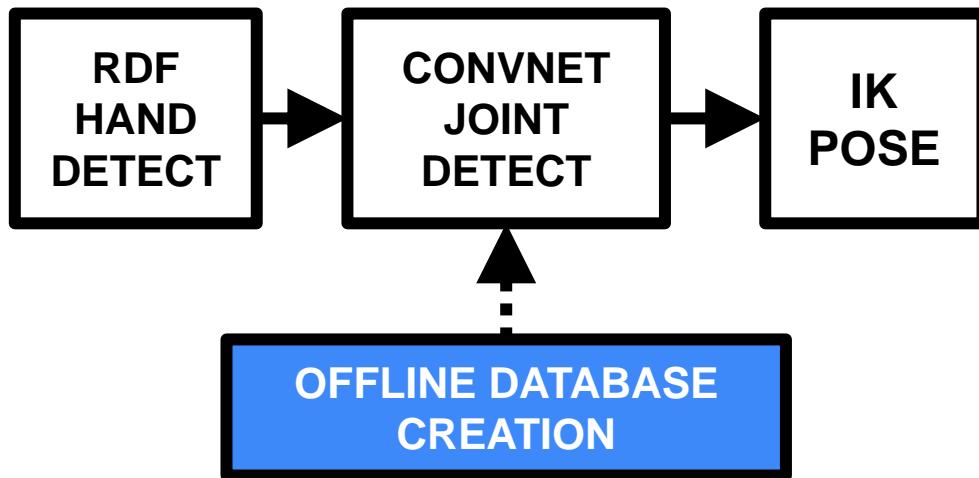
# PIPELINE

Supervised learning based approach

Needs labeled dataset + machine learning

Existing datasets had limited pose information for hands

Architecture



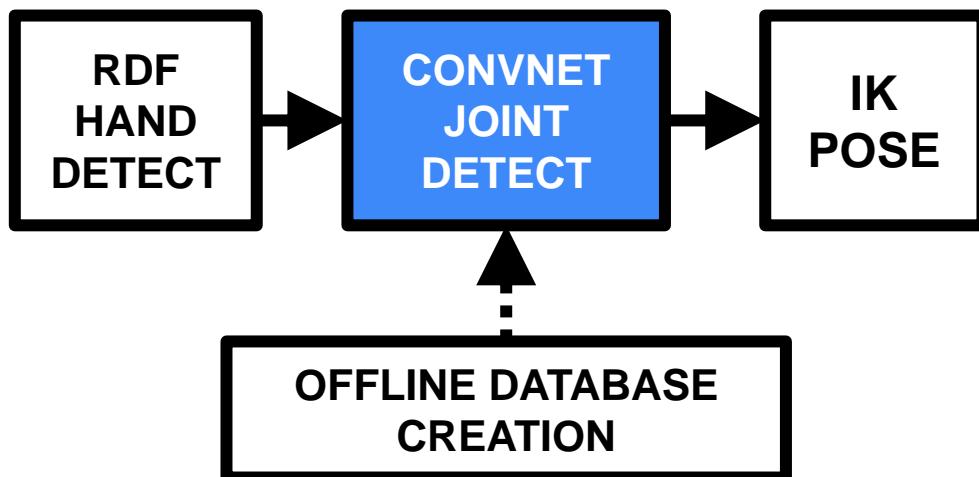
# PIPELINE

Supervised learning based approach

Needs labeled dataset + machine learning

Existing datasets had limited pose information for hands

Architecture



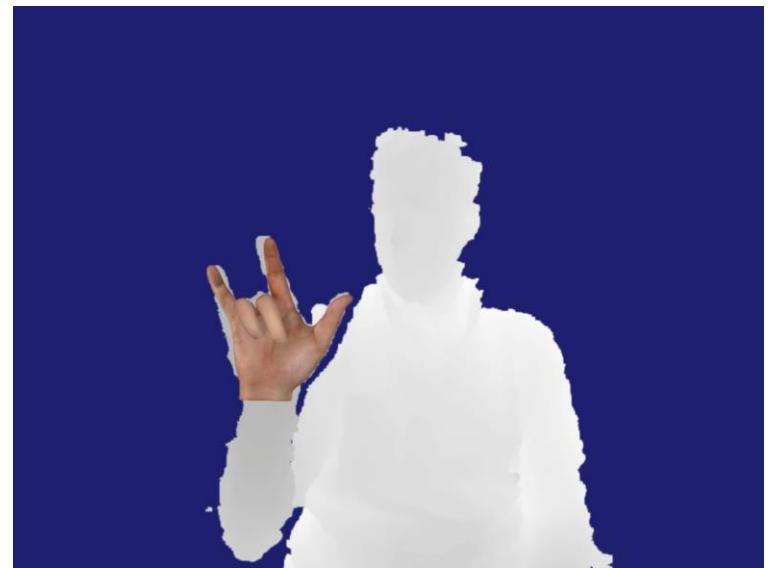
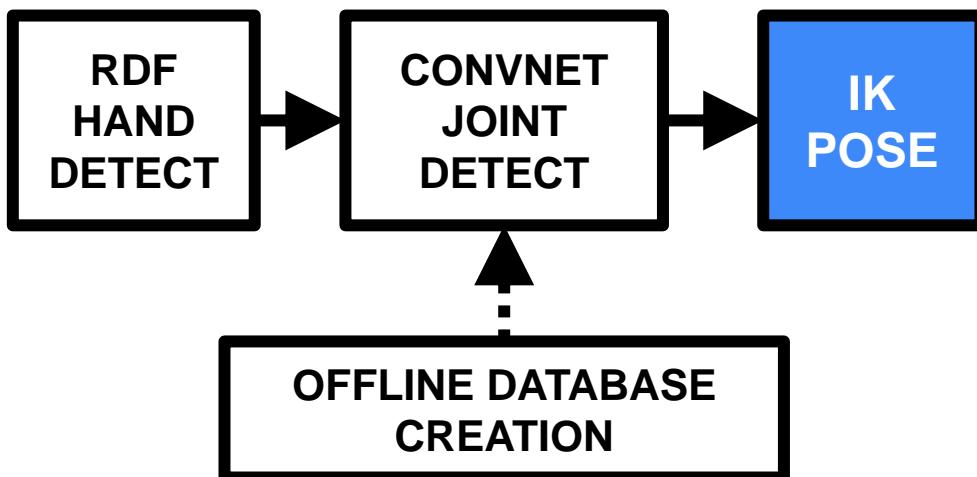
# PIPELINE

Supervised learning based approach

Needs labeled dataset + machine learning

Existing datasets had limited pose information for hands

Architecture



# RDF HAND DETECTION

Per-pixel binary classification → Hand centroid location



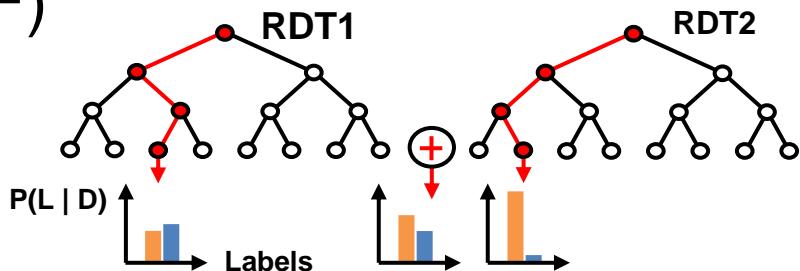
Randomized decision forest (RDF)

Shotton et al.[1]

Fast (parallel)

Generalize

12 hour training, depth 25



Dataset

7500 images (1000 for test)

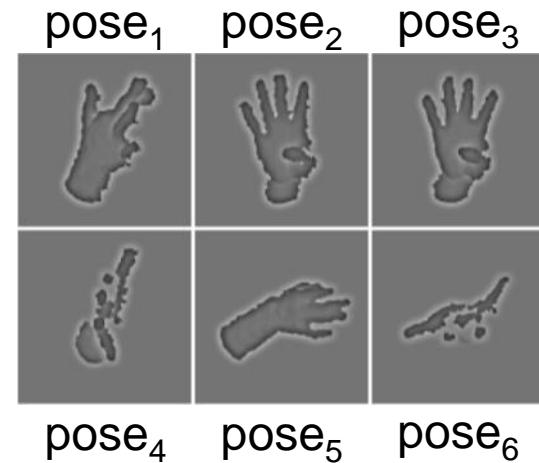
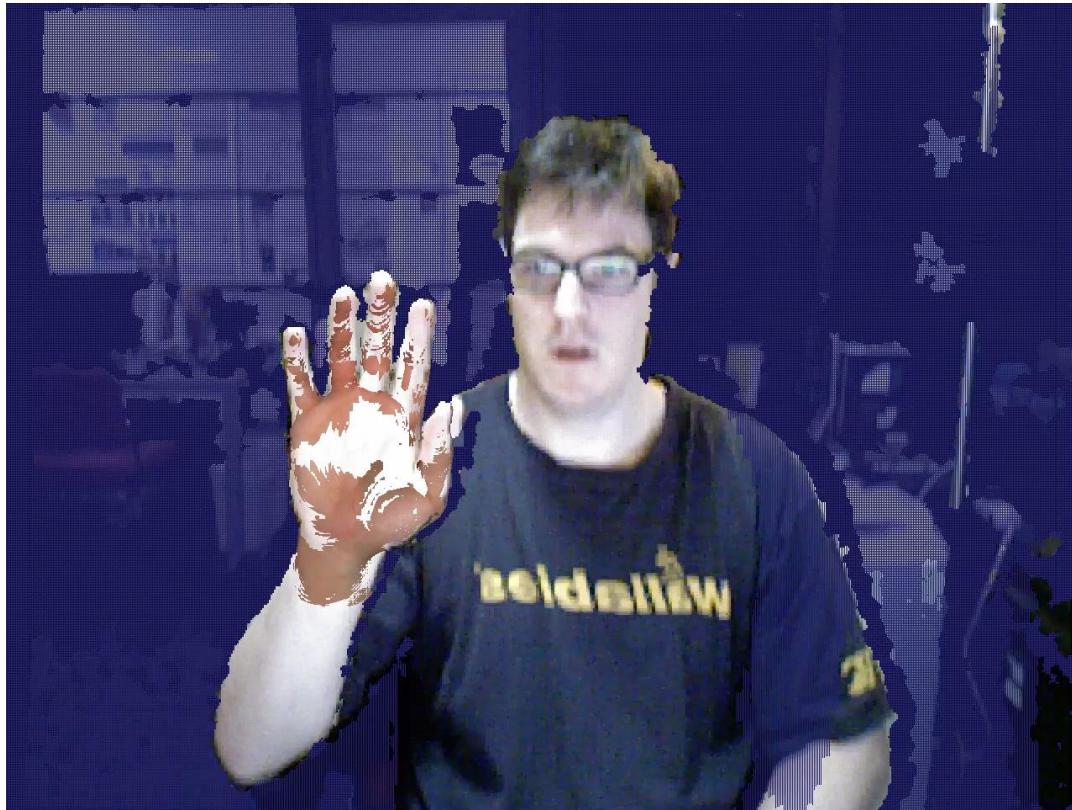


# DATASET CREATION

Goal: labeled RGBD images

$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



# DATASET CREATION

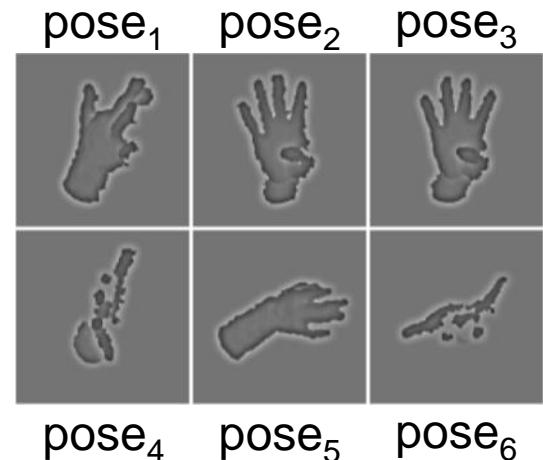
Goal: labeled RGBD images

$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise

Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

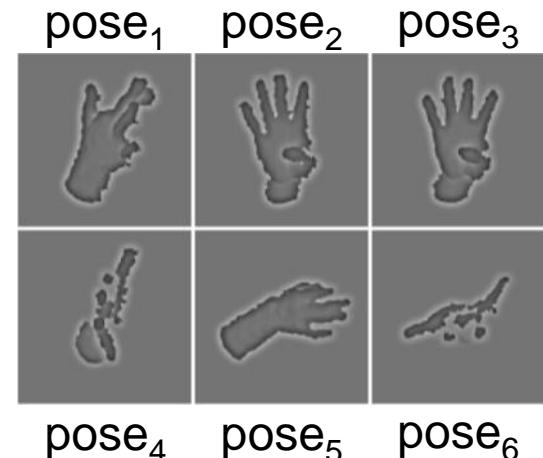


# DATASET CREATION

Goal: labeled RGBD images

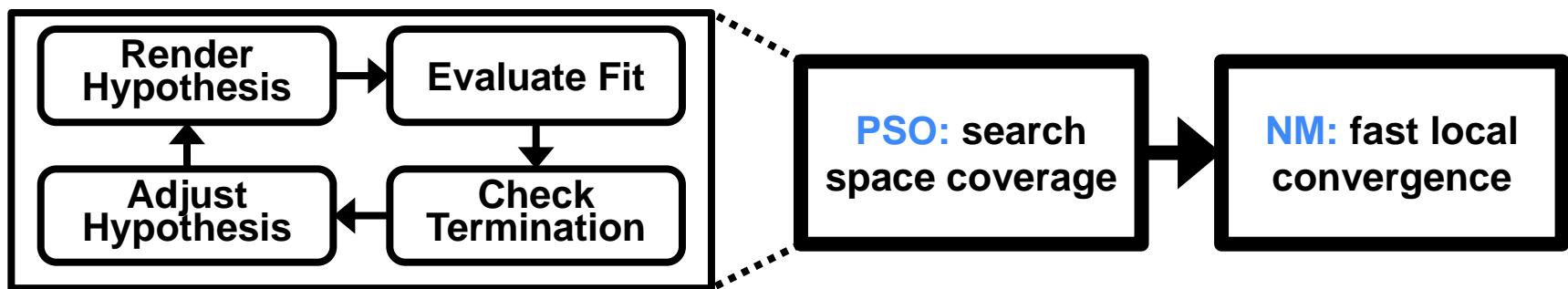
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

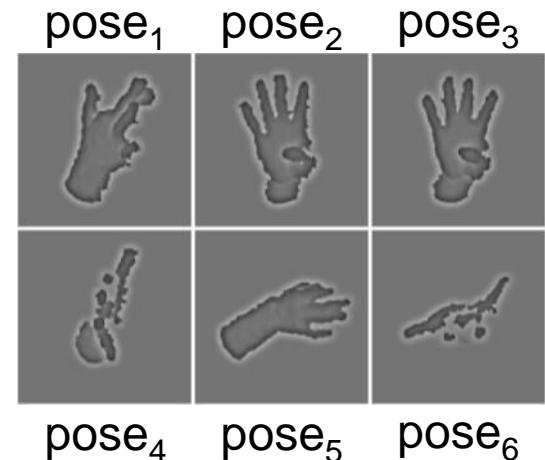


# DATASET CREATION

Goal: labeled RGBD images

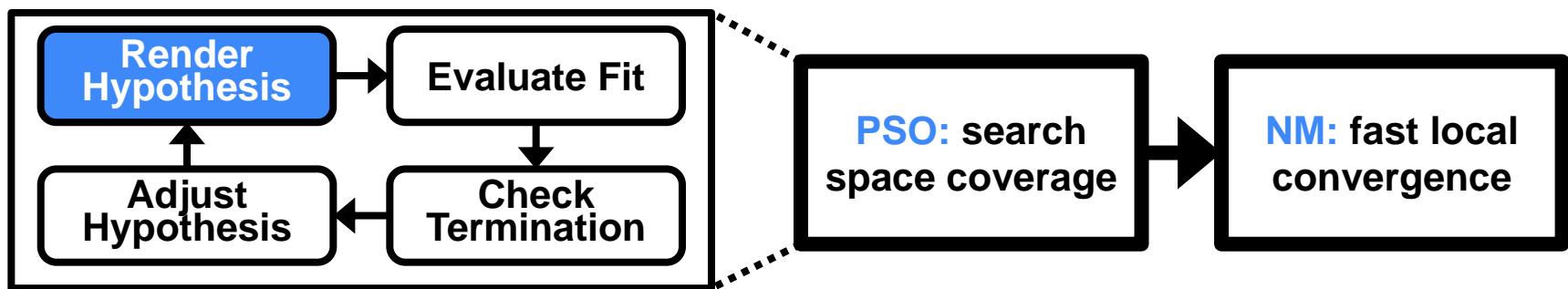
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

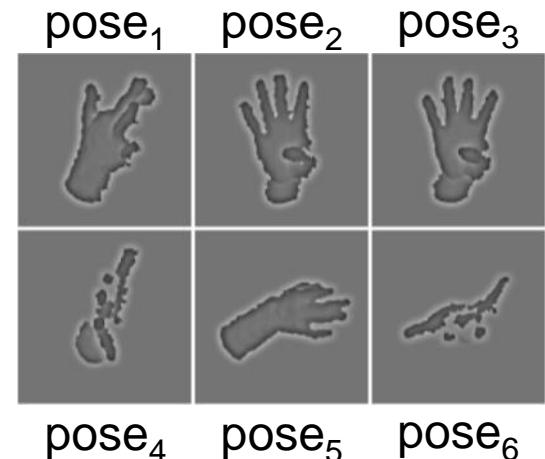


# DATASET CREATION

Goal: labeled RGBD images

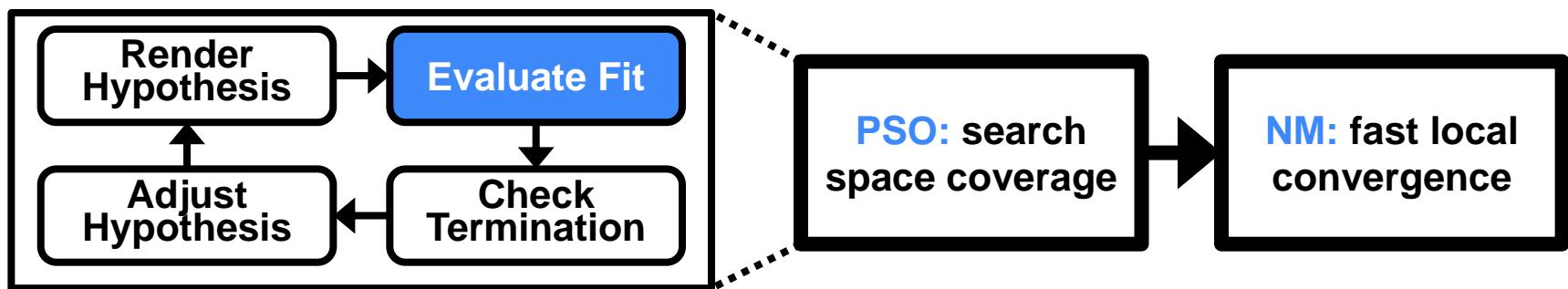
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

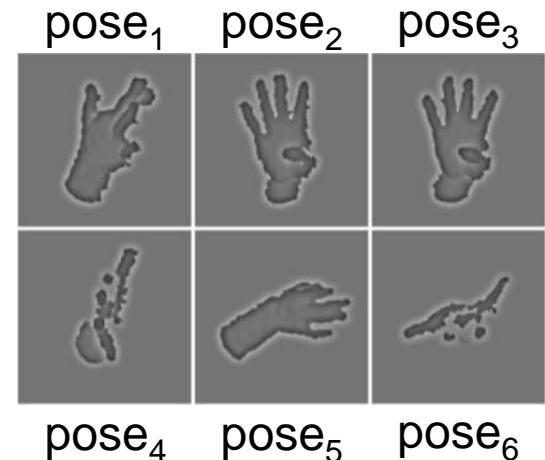


# DATASET CREATION

Goal: labeled RGBD images

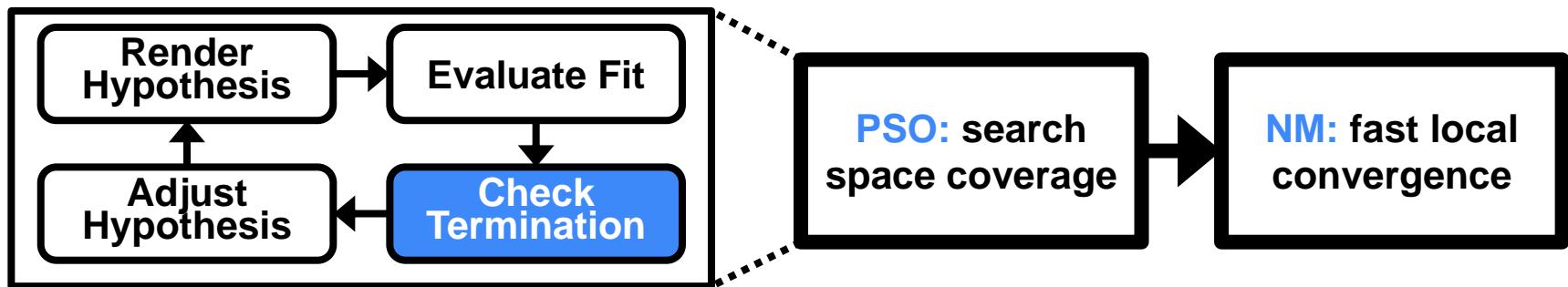
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

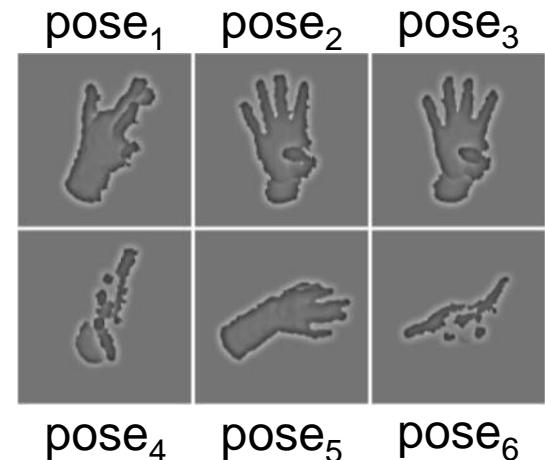


# DATASET CREATION

Goal: labeled RGBD images

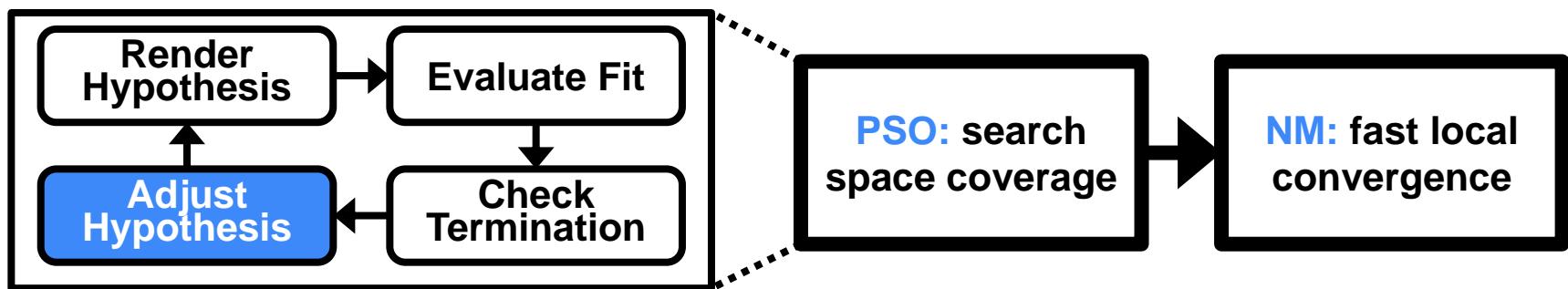
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

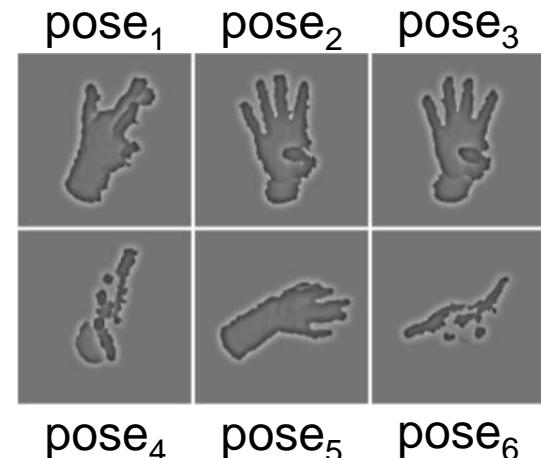


# DATASET CREATION

Goal: labeled RGBD images

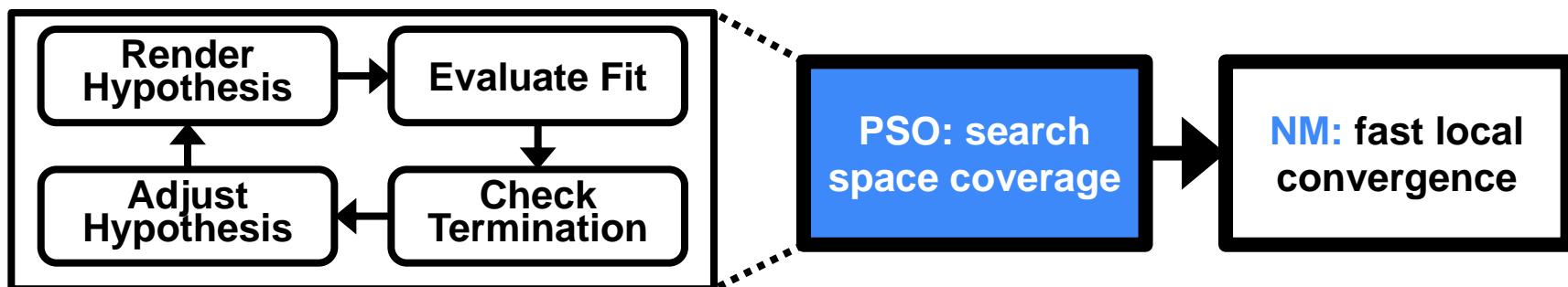
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements

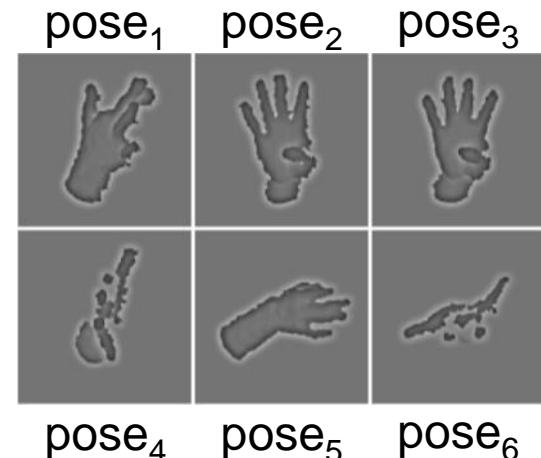


# DATASET CREATION

Goal: labeled RGBD images

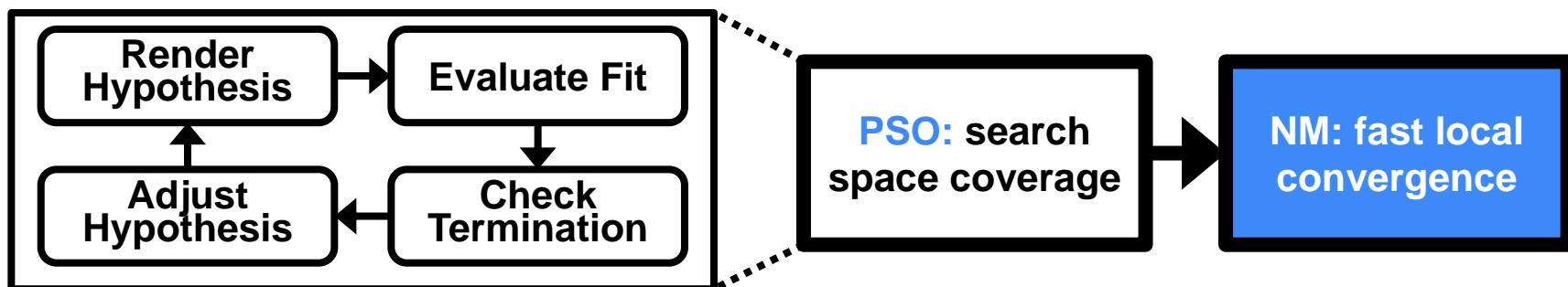
$\{\{D_1, \text{pose}_1\}, \{D_2, \text{pose}_2\}, \dots\}$ ,  $\text{pose}_i$  in  $R^{42}$

Synthetic data doesn't capture device noise



Analysis-by-synthesis Oikonomidis et al.<sup>[1]</sup>

With many improvements



# FEATURE DETECTION

CN has difficulty learning (U,V) positions directly

Require learned integration

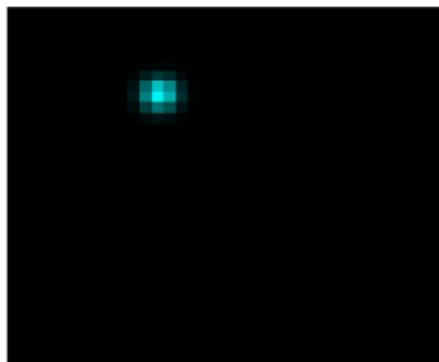
Possible in theory (**never works**)

Recast pose-recognition

Learn feature distributions

$$P_{\text{part1}}(x, y)$$

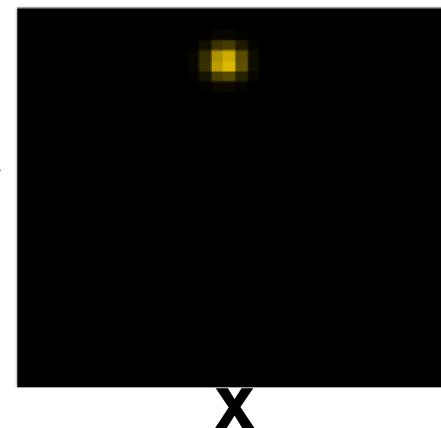
y



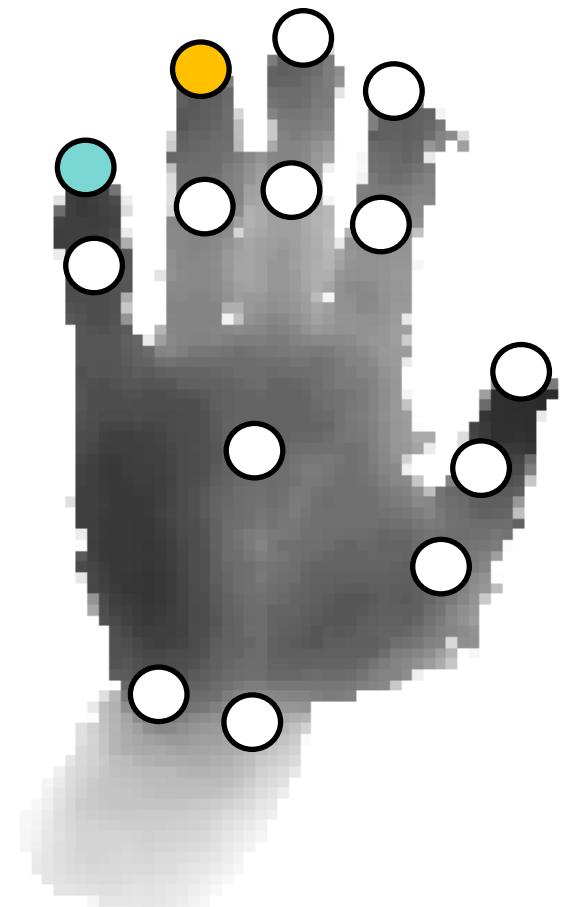
x

$$P_{\text{part2}}(x, y)$$

y



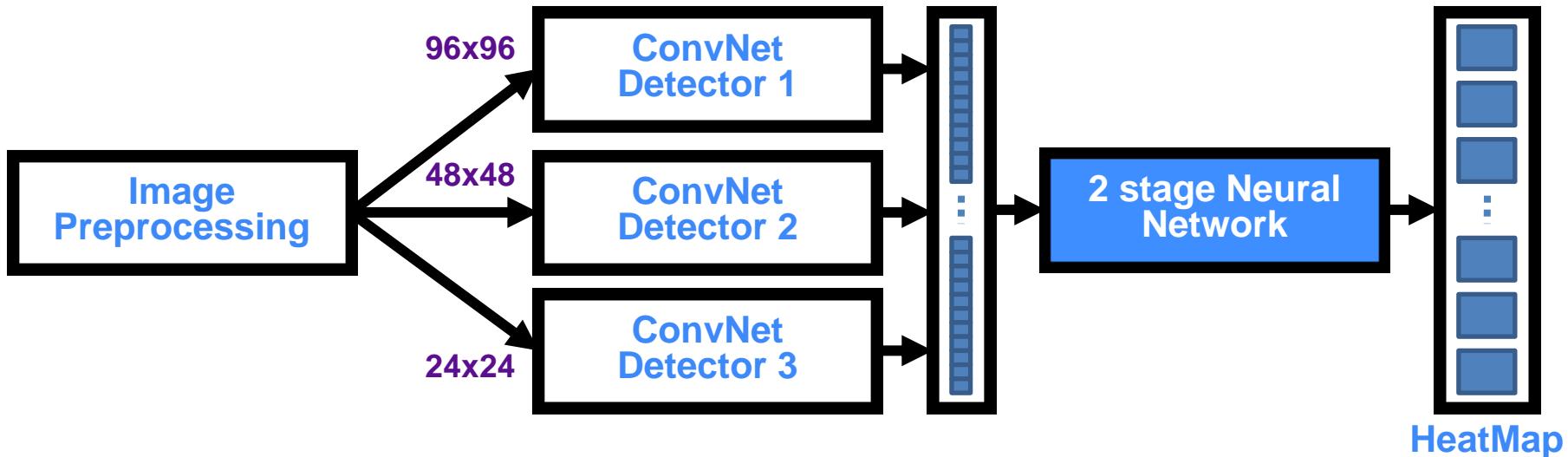
x



# DETECTION ARCHITECTURE

Inspired by Farabet et al. (2013)

Multi-resolution convolutional banks



# INFERRRED JOINT POSITIONS



# POSE RECOVERY (IK)

## Goal

2D heat-maps + 3D depth  $\rightarrow$  3D skeletal pose

## Solution

A variant of Inverse Kinematics (IK)

1. Fit a 2D Gaussian to heat-map (Levenberg-Marquardt)
2. Sample depth at Gaussian mean

Fit skeleton at UVD/UV locations using IK

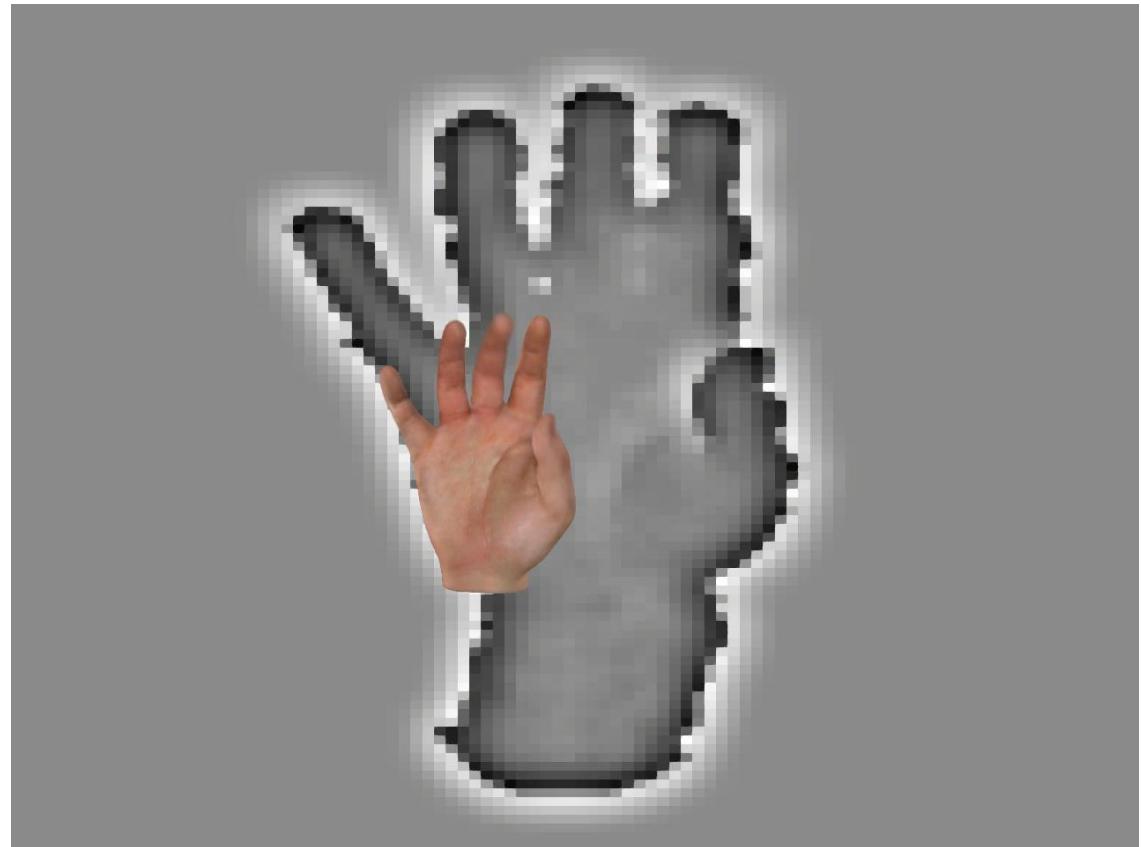
# POSE RECOVERY (IK)

Entire Pipeline: 24.9ms

DF: 3.4ms

CNN: 5.6ms

PSO pose: 11.2ms



# TALK OUTLINE

## 1. Hand Tracking

- a) J. Tompson, M. Stein, Y. LeCun, K. Perlin, "Real-Time Continuous Pose Recovery of Human Hands Using Convolutional Networks", ACM TOG/SIGGRAPH 2014

## 2. Body Tracking

- a) J. Tompson, A. Jain, Y. LeCun, C. Bregler, "Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation", NIPS 2014
- b) A. Jain, J. Tompson, Y. LeCun, C. Bregler, "MoDeep: A Deep Learning Framework Using Motion Features for Human Pose Estimation", ACCV 2014
- c) J. Tompson, R. Goroshin, A. Jain, Y. LeCun, C. Bregler, "Efficient Object Localization Using Convolution Networks", CVPR 2015

# HUMAN BODY POSE INFERENCE

## Motivation:

From my hand tracking work: do similar techniques work in RGB?

How about full-body tracking?

## Problem definition:

Track human body joints on a monocular RGB image

Arbitrary pose and background

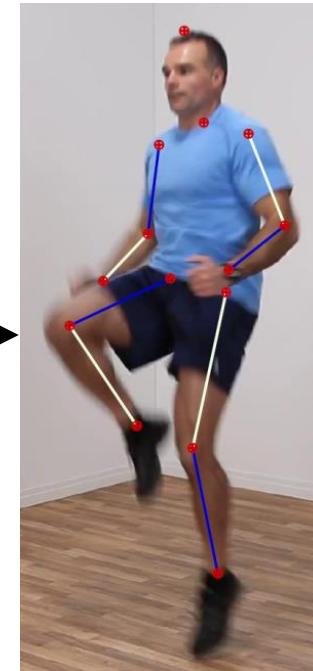
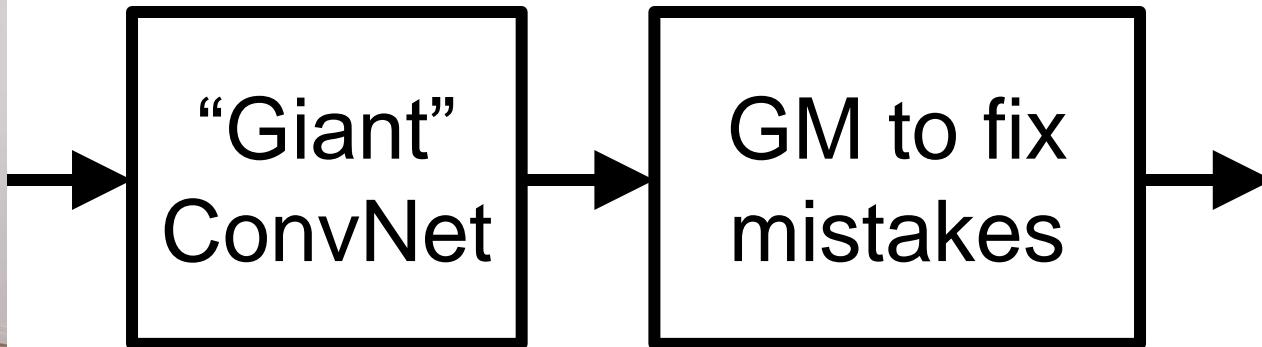
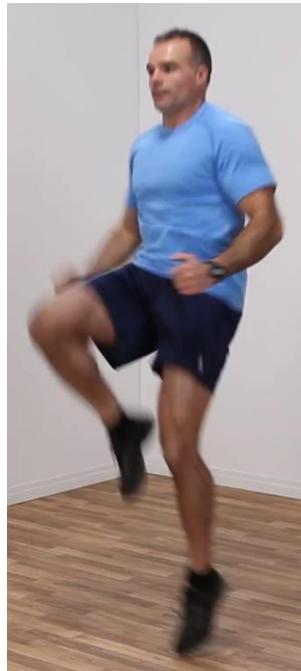
# MY BASIC IDEA

Two Parts:

ConvNet to track joints

Graphical Model to stitch them together

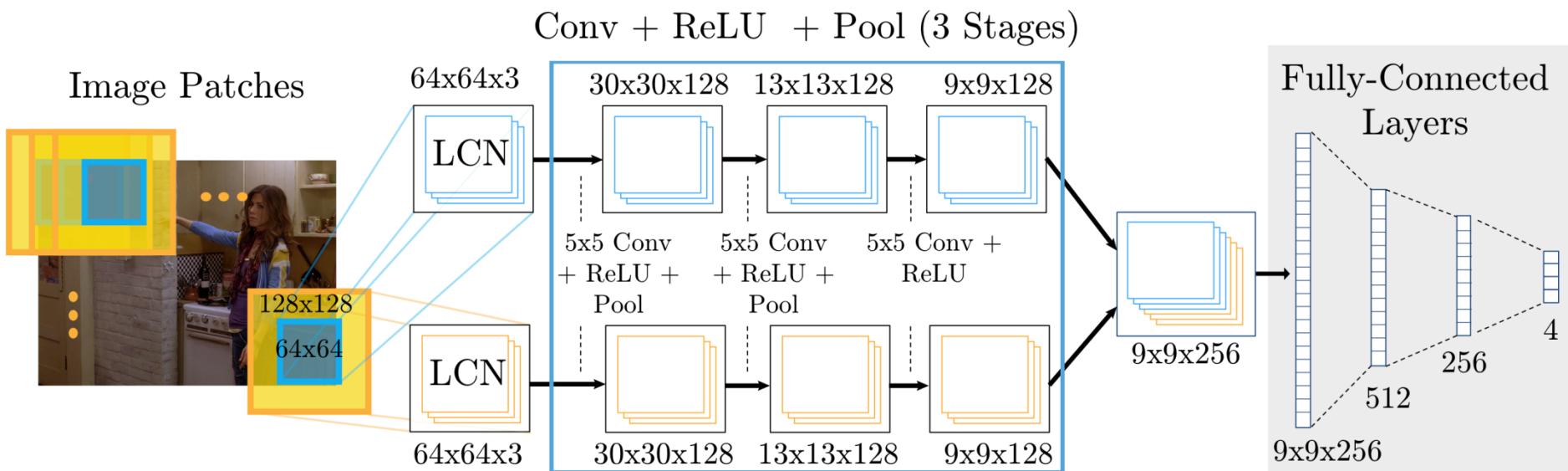
**JOINTLY TRAIN THEM!**



# PART DETECTOR

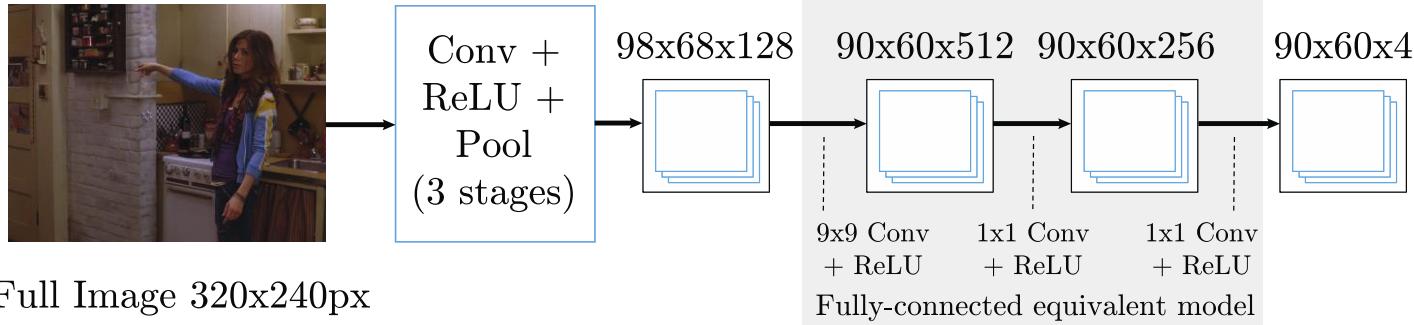
## ConvNet Architecture:

Multi-resolution sliding window with overlapping receptive fields

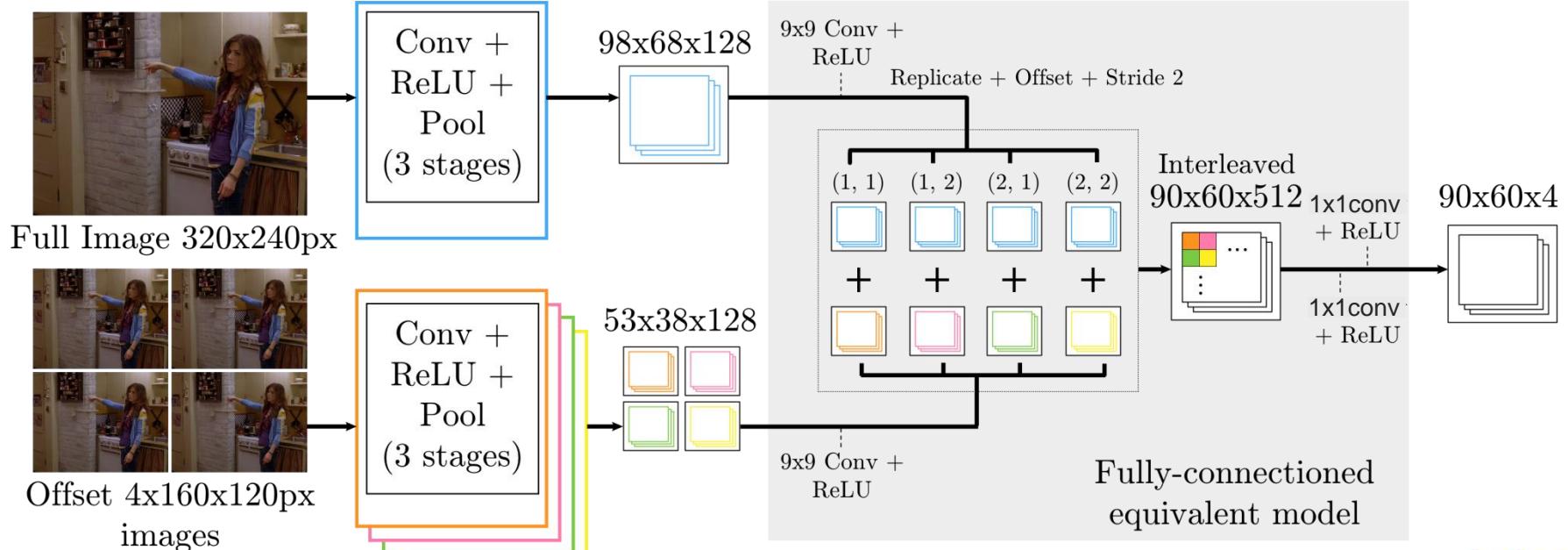


# PART DETECTOR

Efficient model (P. Sermanet and others):



Multi-resolution efficient model (my work):

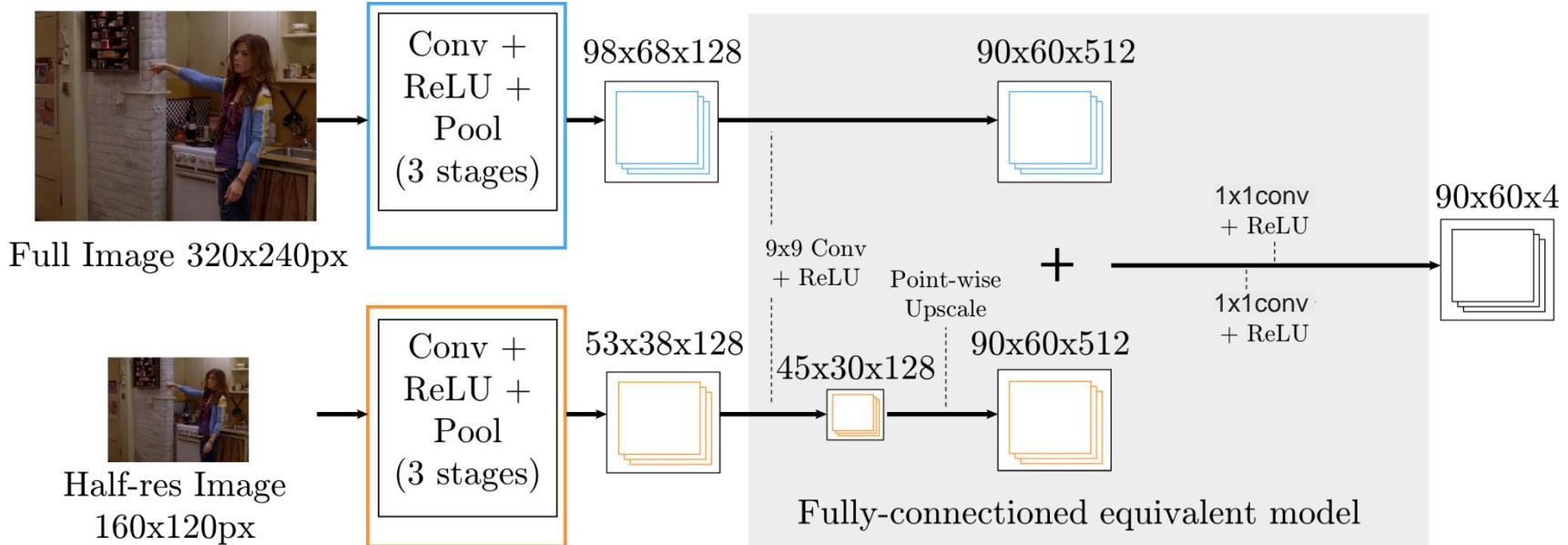


# PART DETECTOR

Simplified model:

Performance is close to the “full” model

Use this for real-time demo



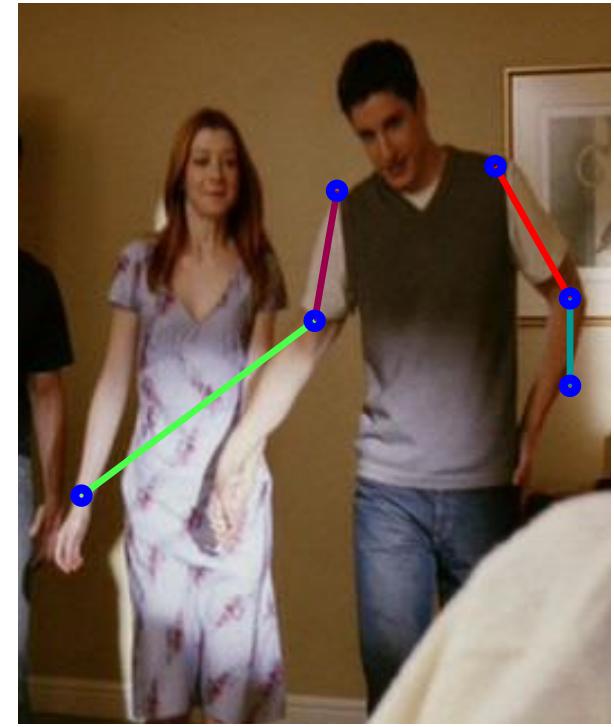
# PART DETECTOR RESULTS

What's wrong so far

Independent joint terms in objective function

We're hoping the network implicitly learns joint consistency

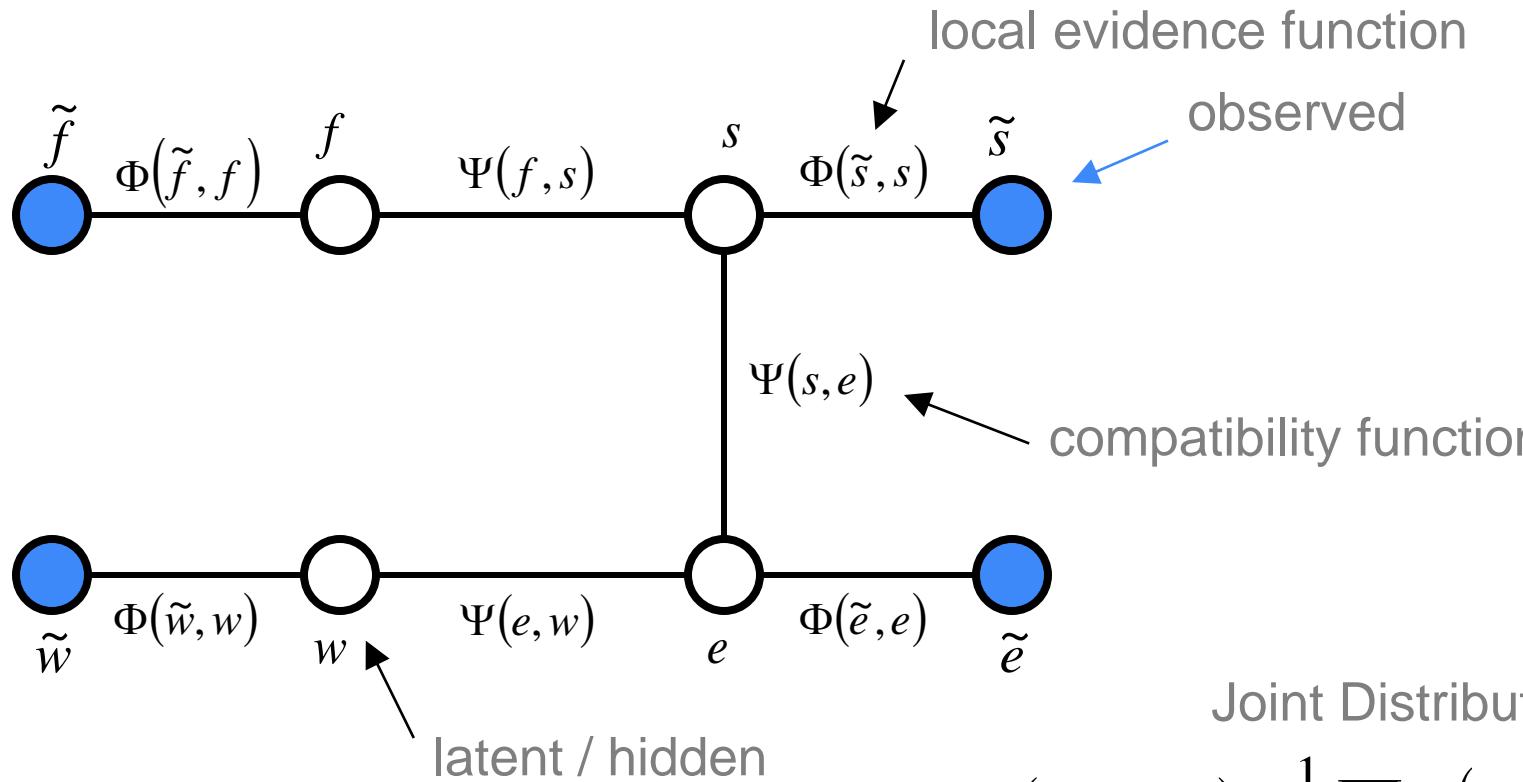
Failure cases are usually stupid:



# SPATIAL MODEL

Start with my GM from ICLR 2013 paper

MRF over spatial locations



Joint Distribution:

$$P(f, s, e, w) = \frac{1}{Z} \prod_{i,j} \Psi(x_i, x_j) \prod_i \Phi(x_i, \tilde{x}_i)$$

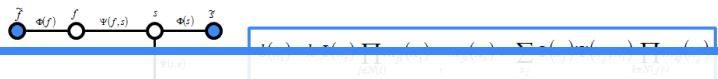
# SPATIAL MODEL

## SPATIAL MODEL

Start with my GM from ICLR 2013 paper

Use MRF over spatial locations

Belief Propagation to calculate marginals  $f, s, e, w$



$$\begin{aligned} b(f) &= k \Phi(f) \sum_s \Phi(s) \Psi(s, f) m_{es}(s) \\ b(f) &= k \Phi(f) \sum_s \Phi(s) \Psi(s, f) \sum_e \Phi(e) \Psi(e, s) m_{we}(e) \\ b(f) &= k \Phi(f) \sum_s \Phi(s) \Psi(s, f) \sum_e \Phi(e) \Psi(e, s) \sum_w \Phi(w) \Psi(w, e) \end{aligned}$$

NYU

## SPATIAL MODEL

Start with my GM from ICLR 2013 paper

Use MRF over spatial locations

Belief Propagation to calculate marginals  $f, s, e, w$



Ignore those terms  
The simplified star graph for each joint



$$b'(f) = k \Phi(f) \prod_i \sum_{x_i} \Phi(x_i) \Psi(x_i, f)$$

NYU

## SPATIAL MODEL

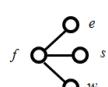
Start with my GM from ICLR 2013 paper

Use MRF over spatial locations

Belief Propagation to calculate marginals  $f, s, e, w$

Use simplified star graph for each joint

Replace tensor factor product & marginalization with convolution and conditional factor (NOT EQUIVALENT)



$$b'(f) = k \Phi(f) \prod_i \sum_{x_i} \Phi(x_i) \Psi(x_i, f)$$



$$b''(f) = k \Phi(f) \prod_i \Phi(x_i) * \Psi(f | x_i)$$

NYU

Wait, where did the partition function go?

$$b'''(\tilde{x}_i) \approx \frac{1}{k} \Phi(f) \prod_i (\Phi(\tilde{x}_i) * \Psi(\tilde{x}_i | \tilde{x}_1) + b(\tilde{x}_i | \tilde{x}_1))$$

One Reason why it's OK

We only care about Maximum Likelihood

→ Distribution Shape not important

NYU

# SPATIAL MODEL

Start with my GM from ICLR 2013 paper

... And approximate it

$$b(f) = \Phi(f) \prod_i (\Phi(x_i) * \Psi(f | x_i) + c(f | x_i))$$



$$\Phi(f)$$

$$* \quad \begin{matrix} \text{---} \\ \text{---} \end{matrix} = \quad \begin{matrix} \text{---} \\ \text{---} \end{matrix}$$

$$\Psi(f | f)$$



$$+ c(f | f)$$



$$b(f)$$



$$\Phi(f)$$

$$* \quad \begin{matrix} \text{---} \\ \text{---} \end{matrix} = \quad \begin{matrix} \text{---} \\ \text{---} \end{matrix}$$

$$\Psi(f | s)$$



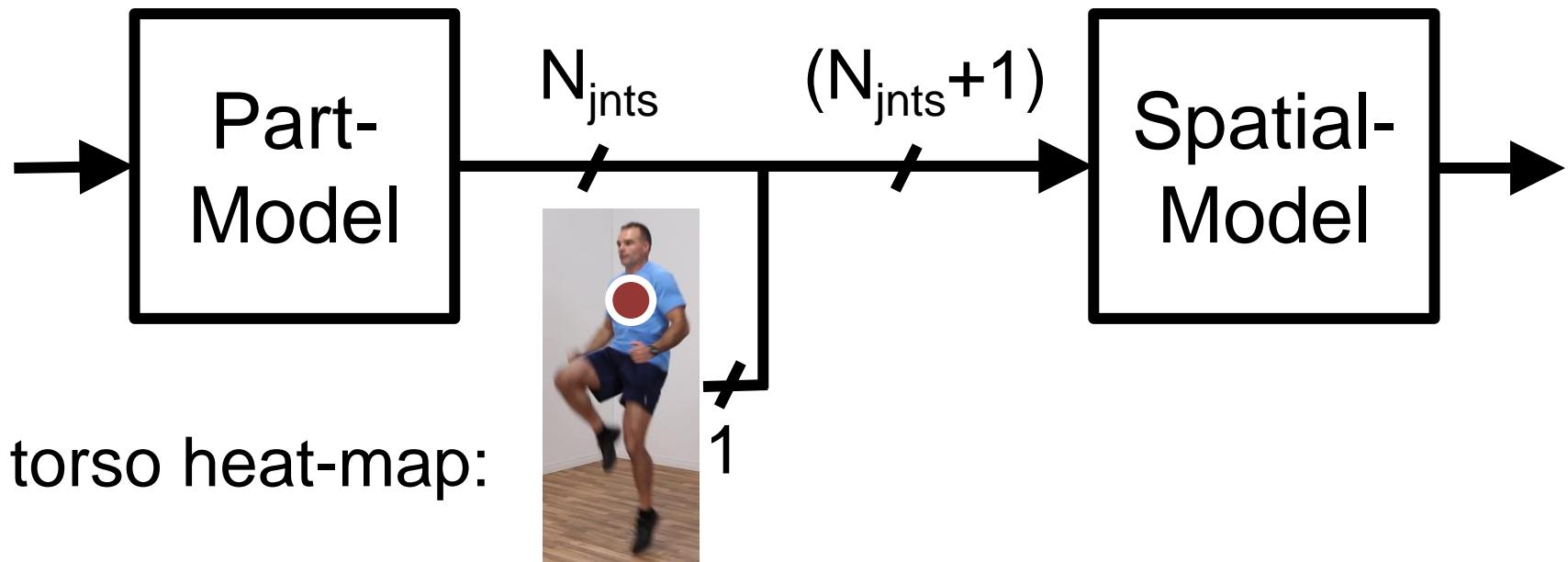
$$+ c(f | s)$$

$$x$$

# SPATIAL MODEL

Two additional details

1. Spatial model kernel size is 128x128! → Have to use FFT<sup>[1]</sup>
2. For standard datasets → add (noisy) torso location



torso heat-map:

[1] M. Mathieu, M. Henaff, and Y. LeCun. Fast training of convolutional networks through FFTs. 2013.

# JOINT TRAINING

## Joint training

Pre-train both models separately

Joint train (BPROP) through both models

# JOINT TRAINING

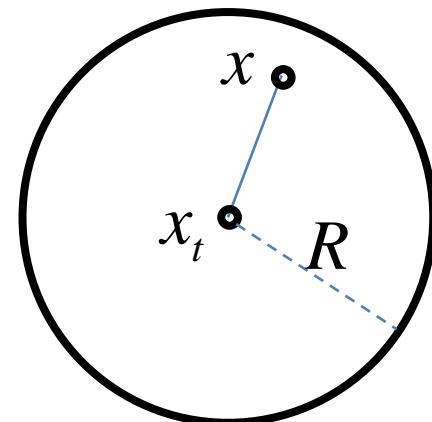
## Joint training

Pre-train both models separately

Joint train (BPROP) through both models

## Performance (Wrist):

$$\text{DetectionRate}(R) = \frac{100}{N} \sum_{t=1}^N \left( \frac{\|x - x^t\|_2}{(\text{torso height } t)/100} \leq R \right)$$



# JOINT TRAINING

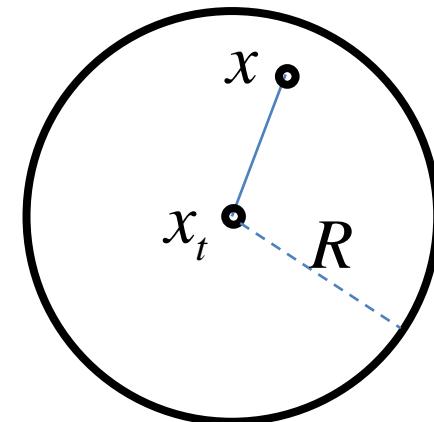
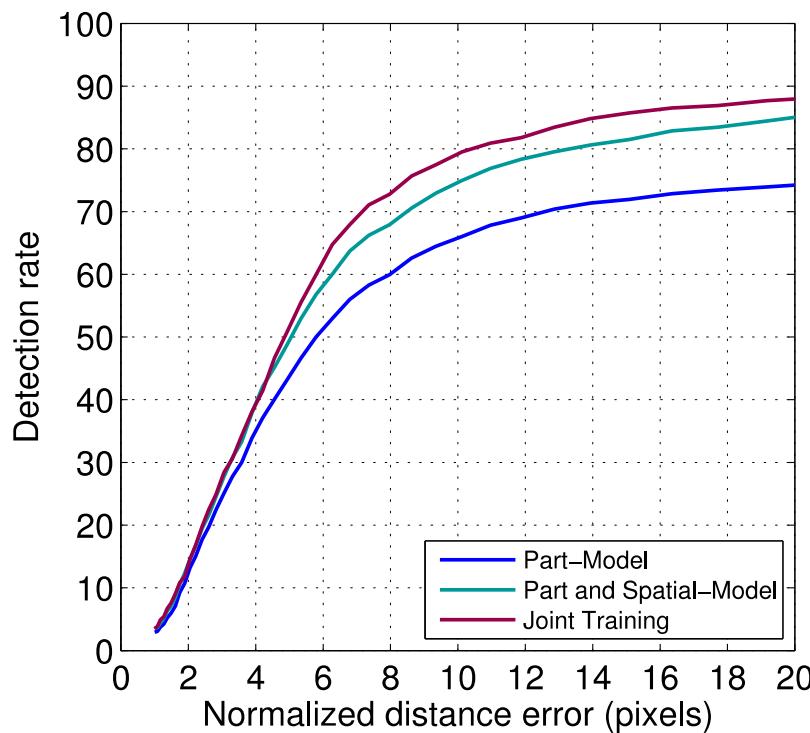
## Joint training

Pre-train both models separately

Joint train (BPROP) through both models

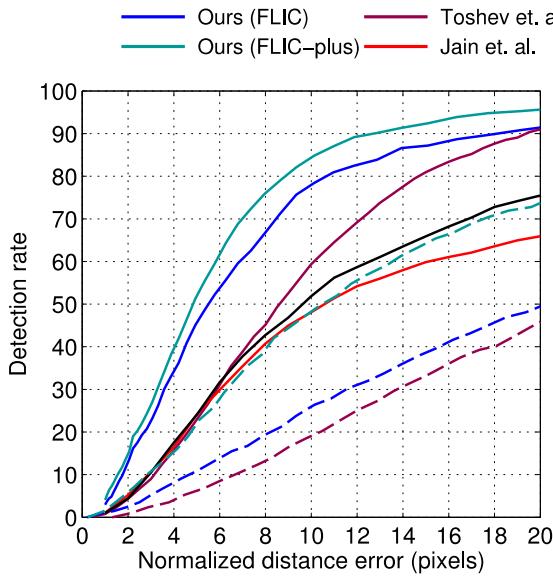
## Performance (Wrist):

$$\text{DetectionRate}(R) = \frac{100}{N} \sum_{t=1}^N \left( \frac{\|x - x^t\|_2}{(\text{torso height } t)/100} \leq R \right)$$

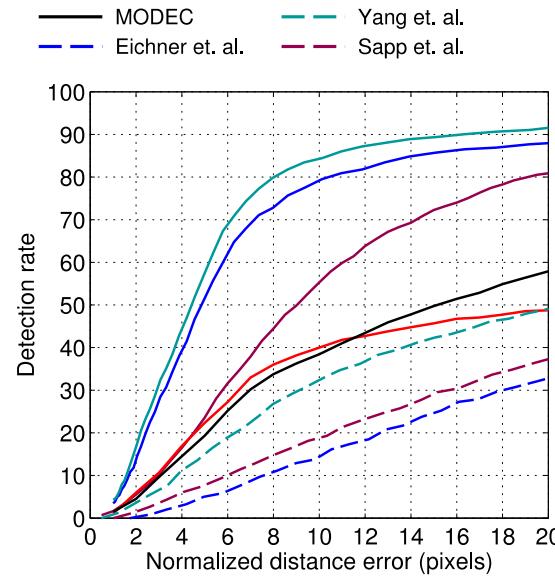


# RESULTS

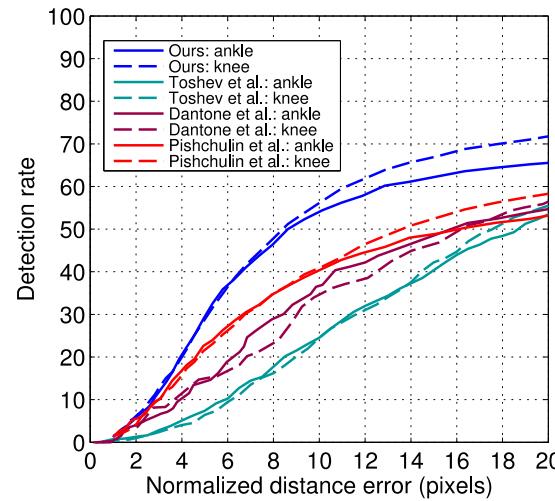
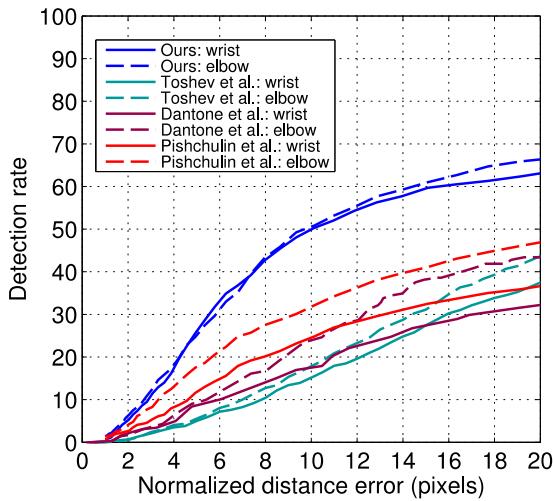
**FLIC<sup>(1)</sup>**  
Elbow



**FLIC<sup>(1)</sup>**  
Wrist



**LSP<sup>(2)</sup>**  
Arms



**LSP<sup>(1)</sup>**  
Legs

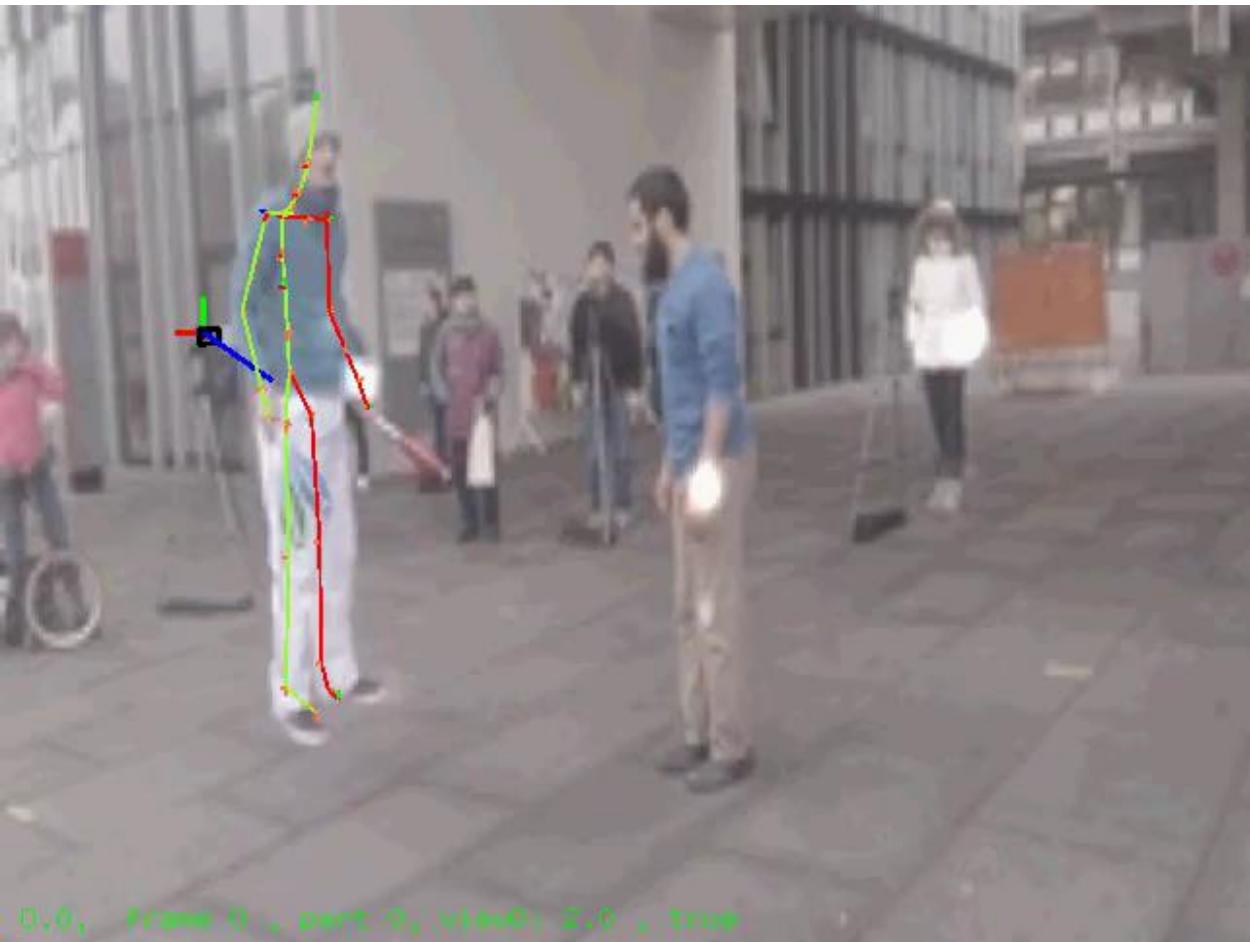
(1) B. Sapp and B. Taskar. MODEC: Multimodel decomposition models for human pose estimation. CVPR'13

(2) S. Johnson and M. Everingham. Learning Effective Human Pose Estimation for Inaccurate Annotation. CVPR'11

# RESULTS

## Joint work with MPII

Use our detector + analysis by synthesis technique



# RESULTS

## Joint work with MPII

Use our detector + analysis by synthesis technique



0.0, frame 0 , part 0, viewidx 2.0 , true

# TALK OUTLINE

## 1. Hand Tracking

- a) J. Tompson, M. Stein, Y. LeCun, K. Perlin, "Real-Time Continuous Pose Recovery of Human Hands Using Convolutional Networks", ACM TOG/SIGGRAPH 2014

## 2. Body Tracking

- a) J. Tompson, A. Jain, Y. LeCun, C. Bregler, "Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation", NIPS 2014
- b) A. Jain, J. Tompson, Y. LeCun, C. Bregler, "[MoDeep: A Deep Learning Framework Using Motion Features for Human Pose Estimation](#)", ACCV 2014
- c) J. Tompson, R. Goroshin, A. Jain, Y. LeCun, C. Bregler, "Efficient Object Localization Using Convolution Networks", CVPR 2015

# INCORPORATING TEMPORAL FEATURES

Single frame detection is overly constrained

Ambiguities can be resolved using multiple frames

Our ACCV Paper - A short investigation:

- How to use temporal features with a ConvNet detector
- Performance is improved using motion features

**motion  
features**



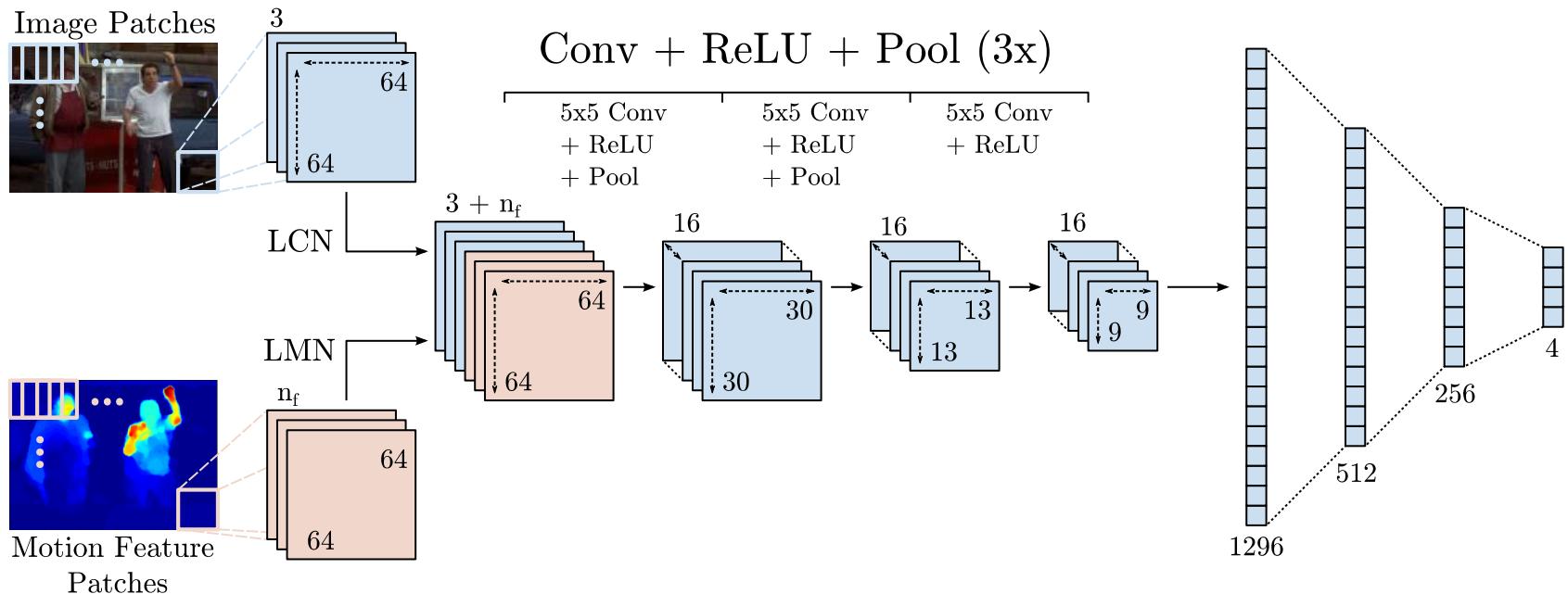
**baseline**



# MOTION FEATURES

Brainless:

Just concat some form of “motion feature” with the RGB



Question:

What features should we use?

# MOTION FEATURES

Try many and see what works:

RGB image pair

$$\{f_i, f_{i+\delta}\}$$

RGB difference

$$\{f_i, f_{i+\delta} - f_i\}$$

RGB + 2D Optical Flow

$$\{f_i, \text{FLOW}(f_{i+\delta}, f_i)\}$$

RGB + 2D Optical Flow Mag

$$\{f_i, \|\text{FLOW}(f_{i+\delta}, f_i)\|_2\}$$

... etc

# MOTION FEATURES

Try many and see what works:

RGB image pair

RGB difference

RGB + 2D Optical Flow

RGB + 2D Optical Flow Mag

$$\{f_i, f_{i+\delta}\}$$

$$\{f_i, f_{i+\delta} - f_i\}$$

$$\{f_i, \text{FLOW}(f_{i+\delta}, f_i)\}$$

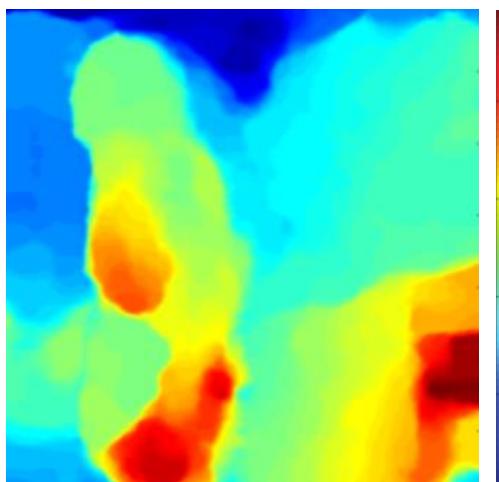
$$\{f_i, \|\text{FLOW}(f_{i+\delta}, f_i)\|_2\}$$

Annoying detail:

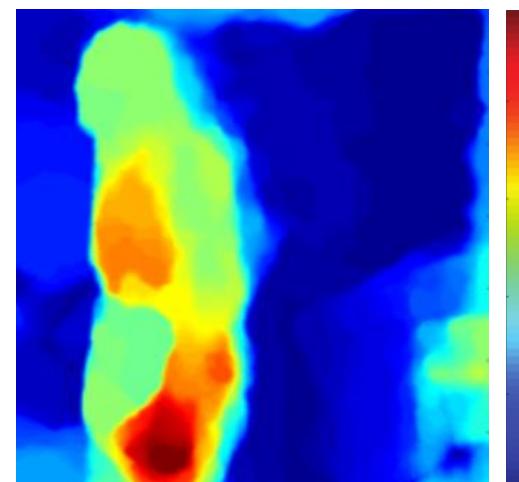
SIFT + RANSAC + Bundle adjustment → planar camera motion

Then subtract it out:

Optical flow mag



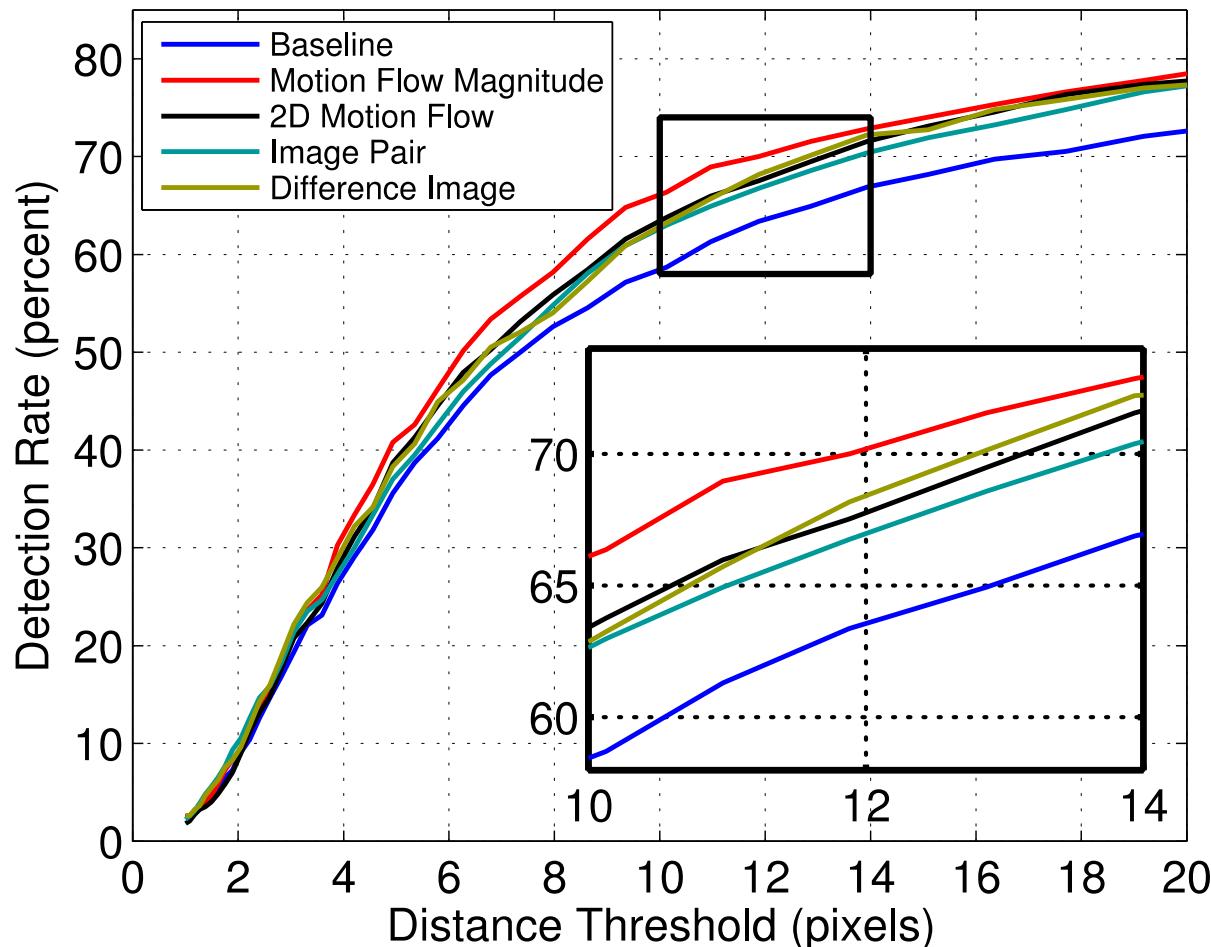
Optical flow mag  
w.o.  
camera



# MOTION FEATURES: RESULTS

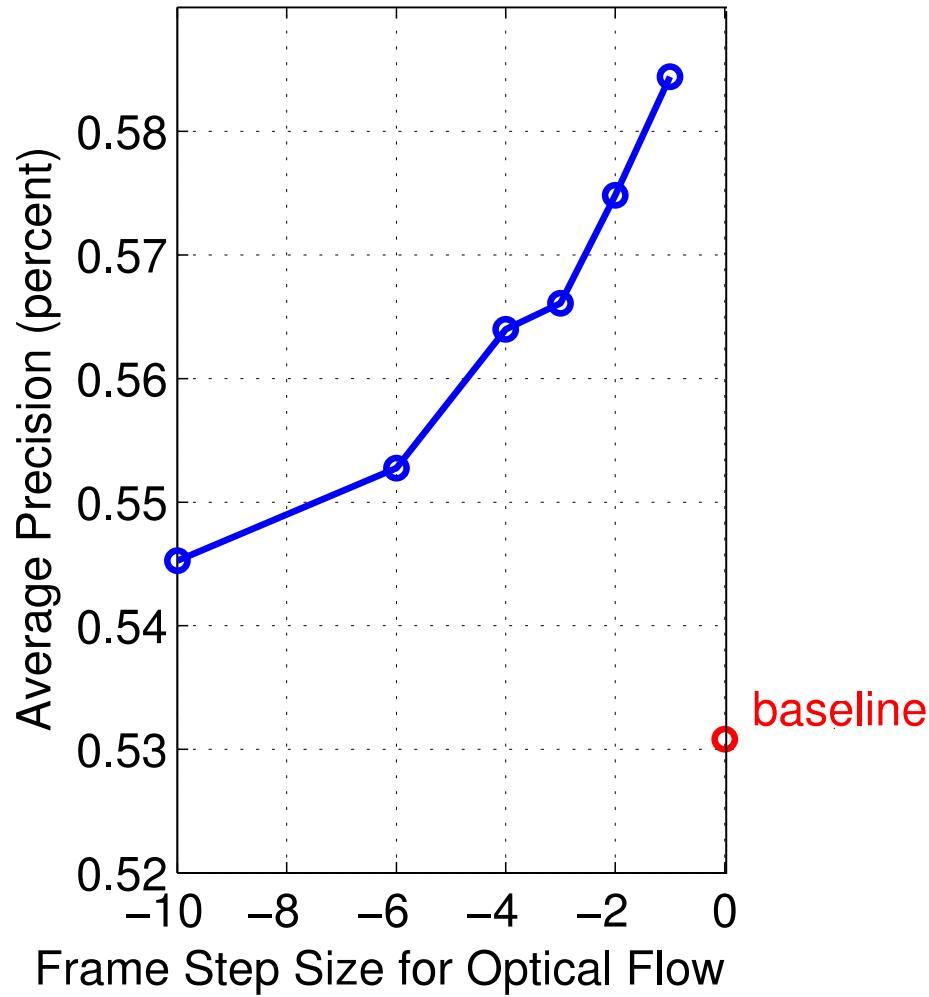
Unsurprising... Optical flow is the best

Magnitude prevents overtraining on motion direction



# MOTION FEATURES

How sensitive is it to time?



# TALK OUTLINE

## 1. Hand Tracking

- a) J. Tompson, M. Stein, Y. LeCun, K. Perlin, "Real-Time Continuous Pose Recovery of Human Hands Using Convolutional Networks", ACM TOG/SIGGRAPH 2014

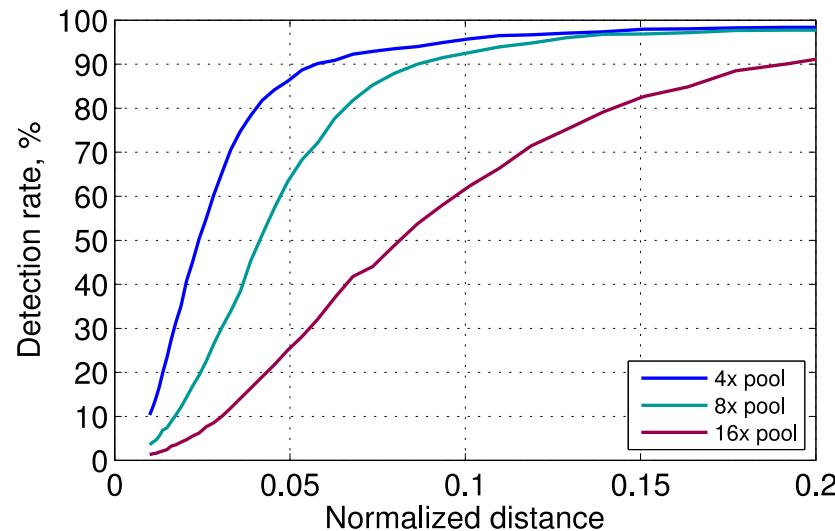
## 2. Body Tracking

- a) J. Tompson, A. Jain, Y. LeCun, C. Bregler, "Joint Training of a Convolutional Network and a Graphical Model for Human Pose Estimation", NIPS 2014
- b) A. Jain, J. Tompson, Y. LeCun, C. Bregler, "MoDeep: A Deep Learning Framework Using Motion Features for Human Pose Estimation", ACCV 2014
- c) J. Tompson, R. Goroshin, A. Jain, Y. LeCun, C. Bregler, "Efficient Object Localization Using Convolution Networks", CVPR 2015

# MOTIVATION

We use pooling

Too much: bad spatial accuracy



Too little: Over-trains, slow, not enough context

Idea:

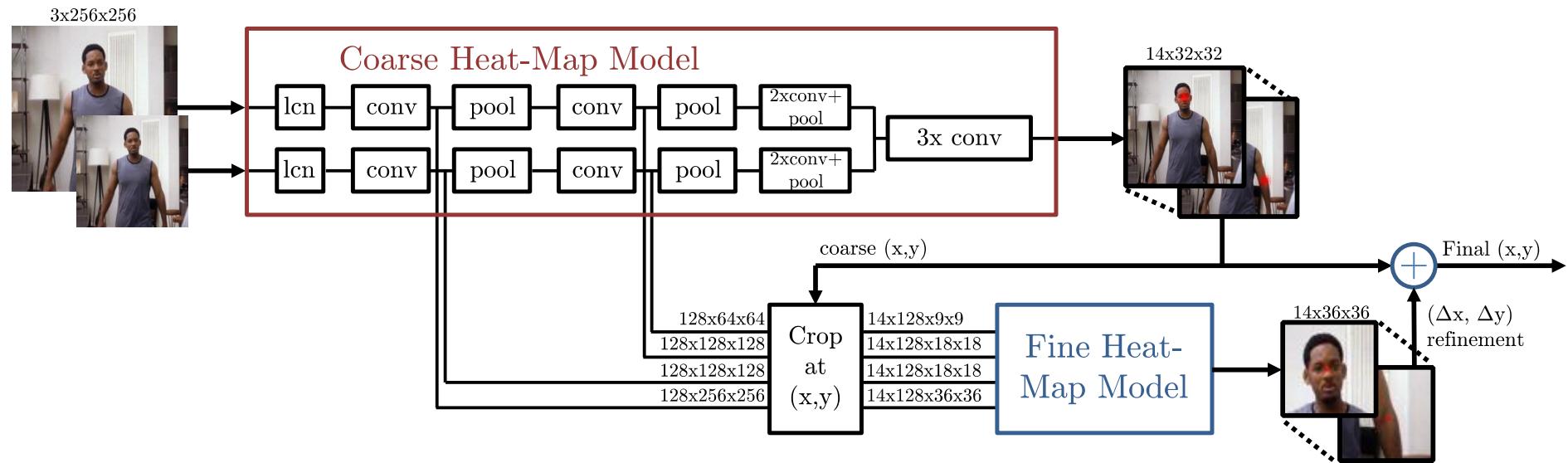
Use lots of pooling but recover spatial accuracy

A “smart” cascade (highly engineered)

# HIGH LEVEL ARCHITECTURE

Highly “engineered”

Coarse to Fine architecture with shared features



$$E_1 = \frac{1}{N} \sum_{j=1}^N \sum_{xy} \|H'_j(x, y) - H_j(x, y)\|^2$$

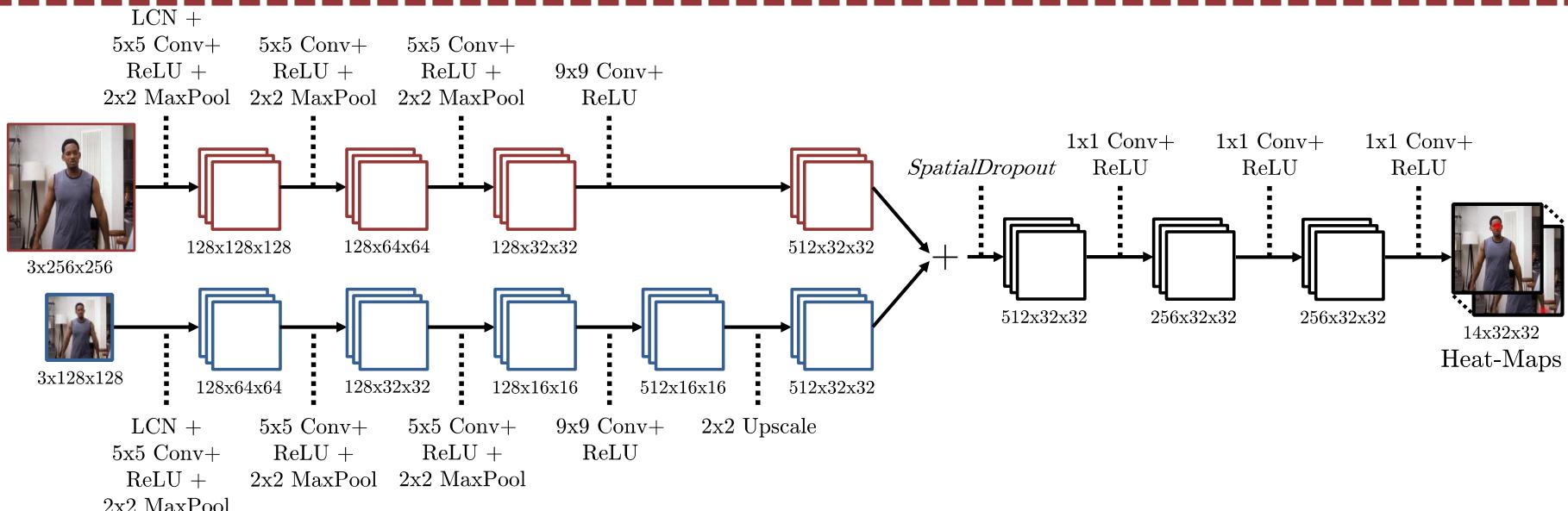
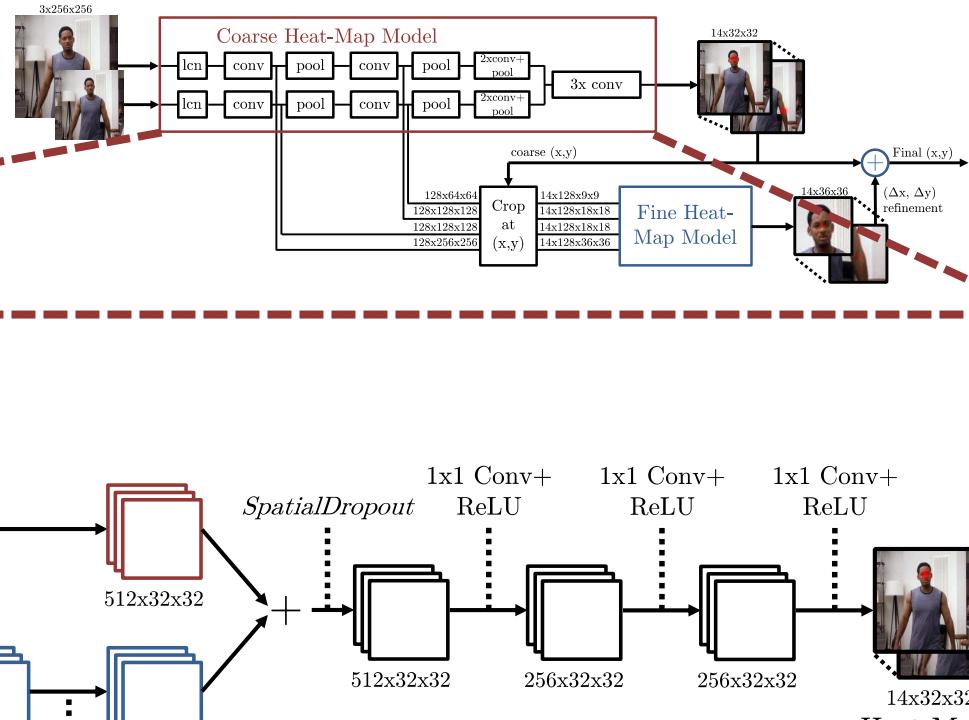
$$E_2 = E_1 + \lambda \frac{1}{N} \sum_{j=1}^N \sum_{x,y} \|G'_j(x, y) - G_j(x, y)\|^2$$

# COARSE MODEL

Straight from the NIPS paper

Many more features (more engineering)

Added “*Spatial Dropout*”



# CASCADE CROP

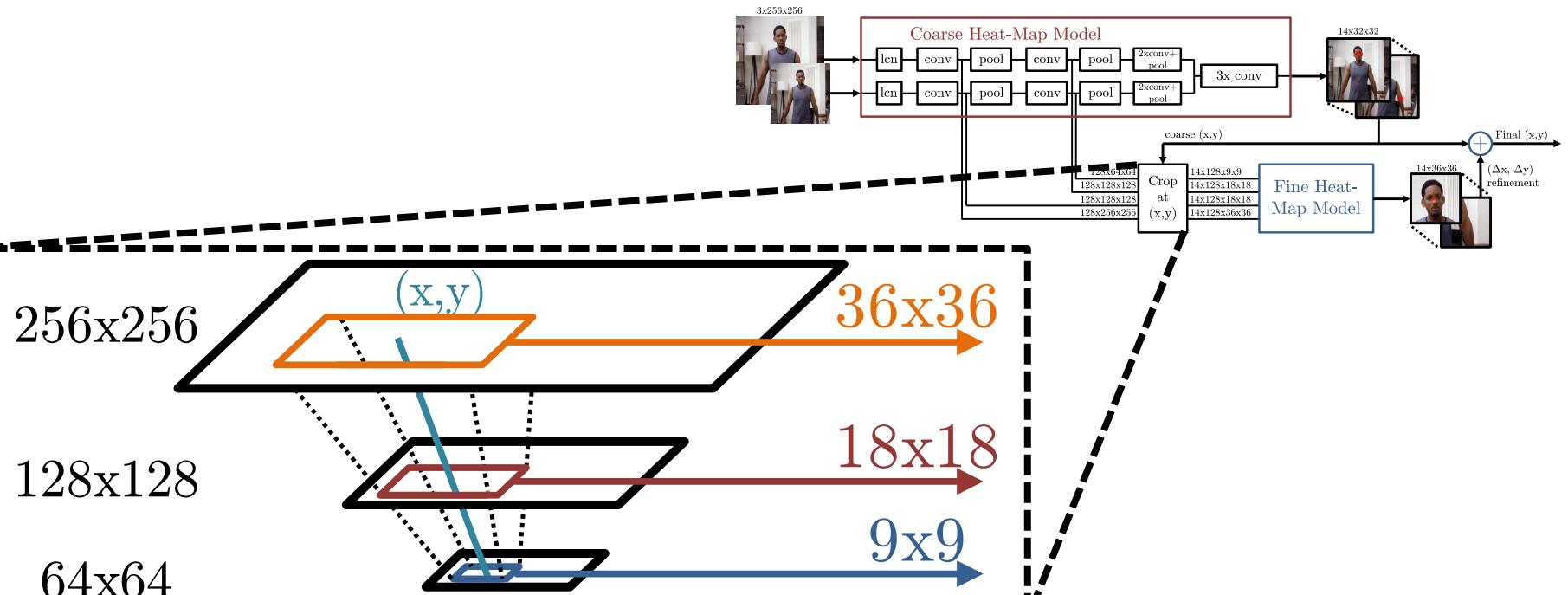
What do we want?

To refine the coarse heat-map by cropping around the approx. UV

What features do we need?

Early layers that haven't specialized

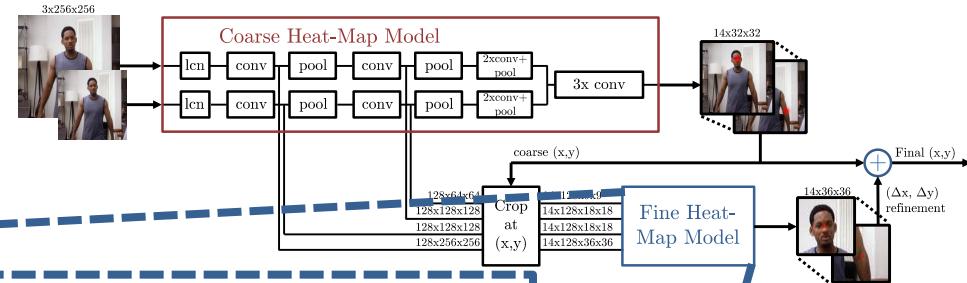
Consistent Spatial Context from all layers



# FINE MODEL

Inputs sampled from multiple locations → Need separate networks (Siamese)

Shared weights ( $L_1$  thru  $L_{n-1}$ )

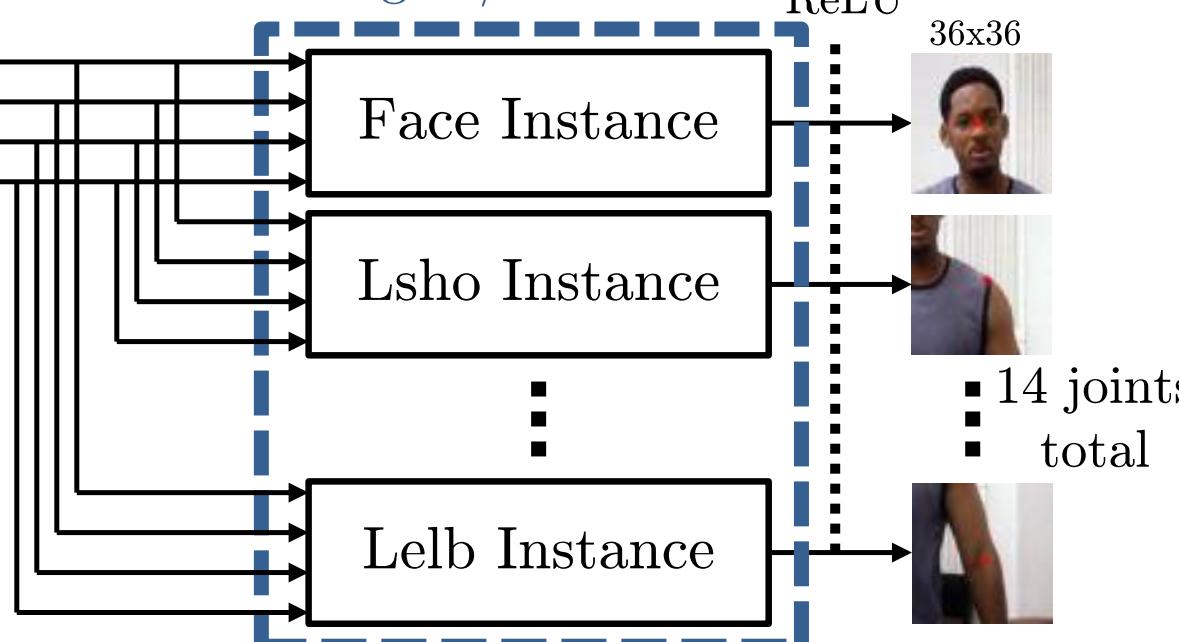


Shared weights/biases

1x1 Conv+  
ReLU



14 joints  
total

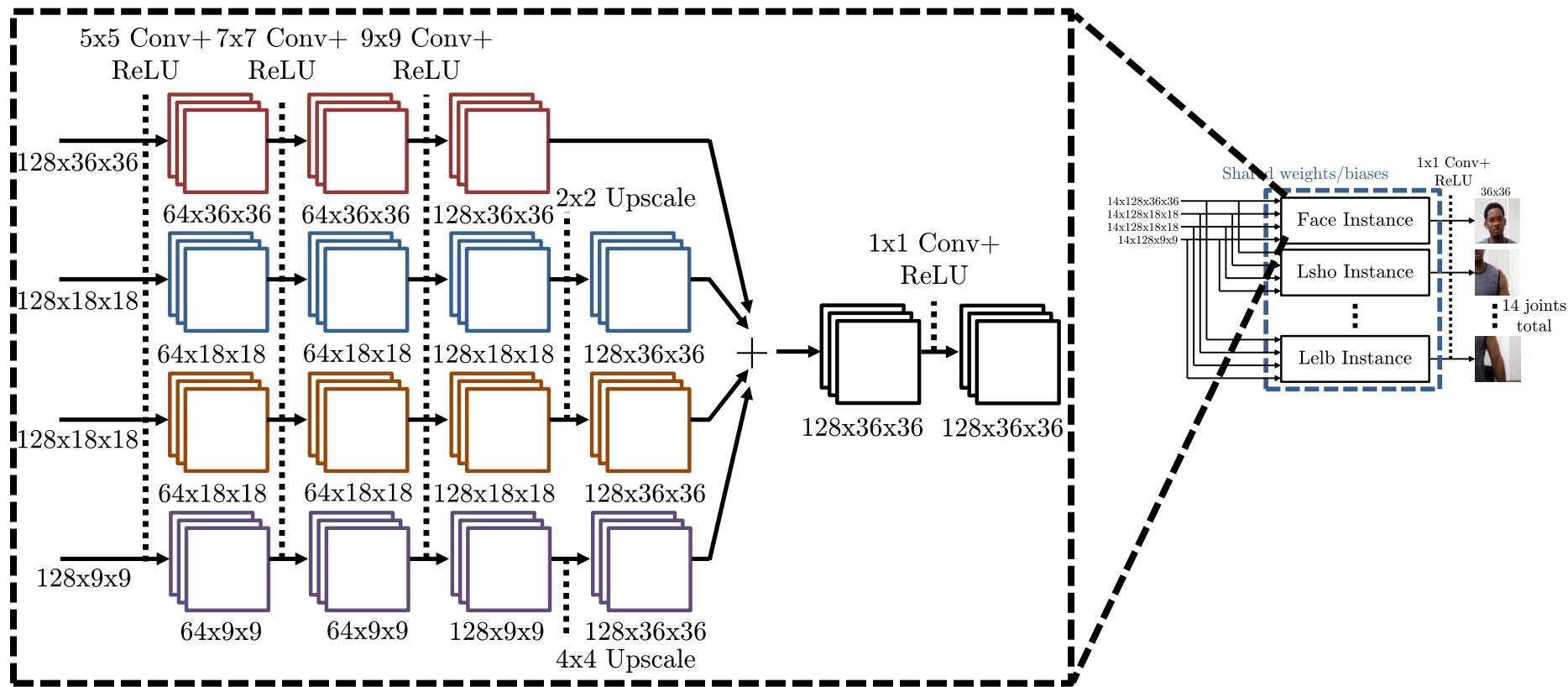


# FINE MODEL – REPLICATED INSTANCE

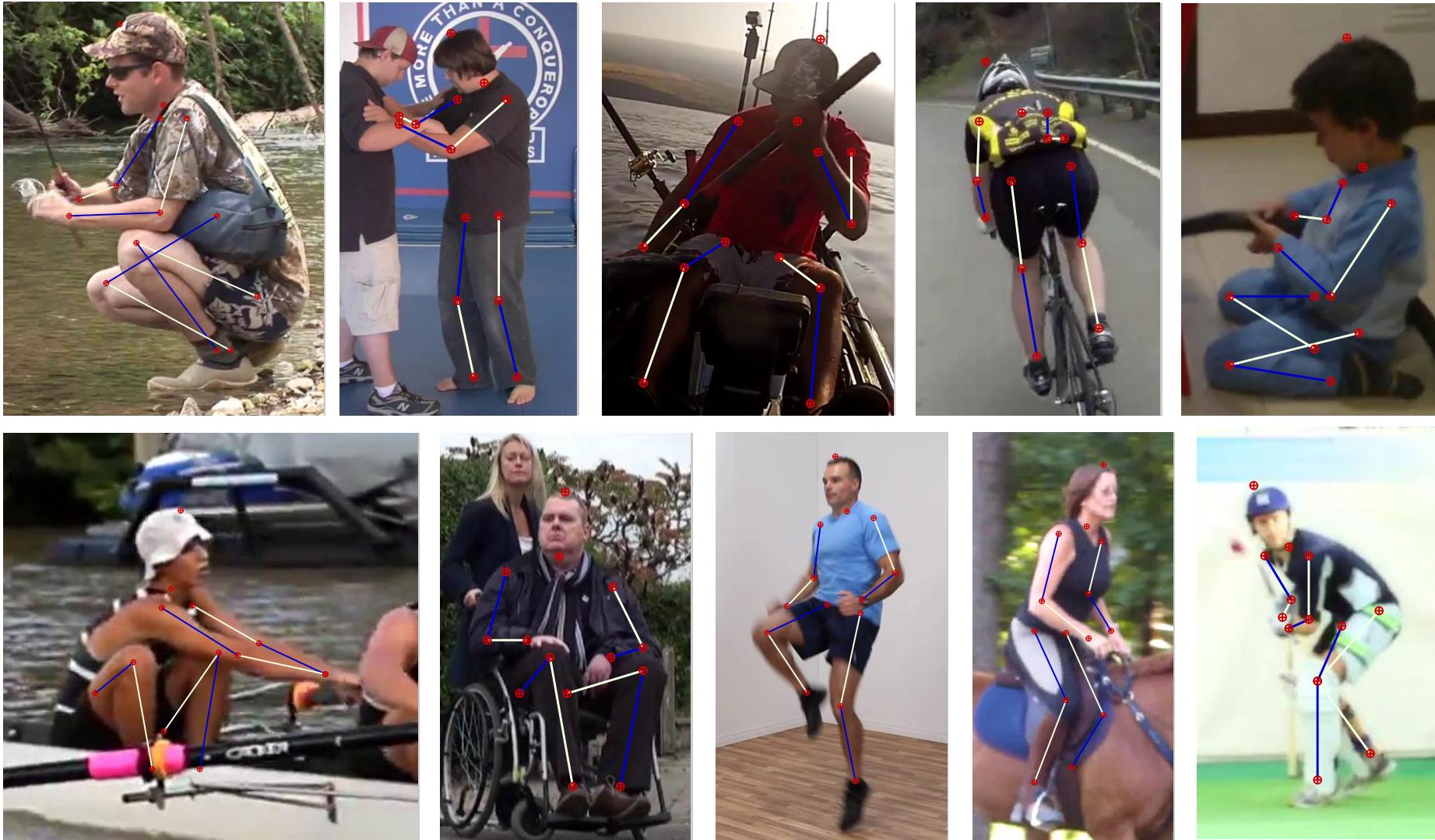
Use same strategy as the Coarse Model

Fully Convolutional (with 1x1 conv layers)

Up-sample to bring features into canonical resolution



# RESULTS



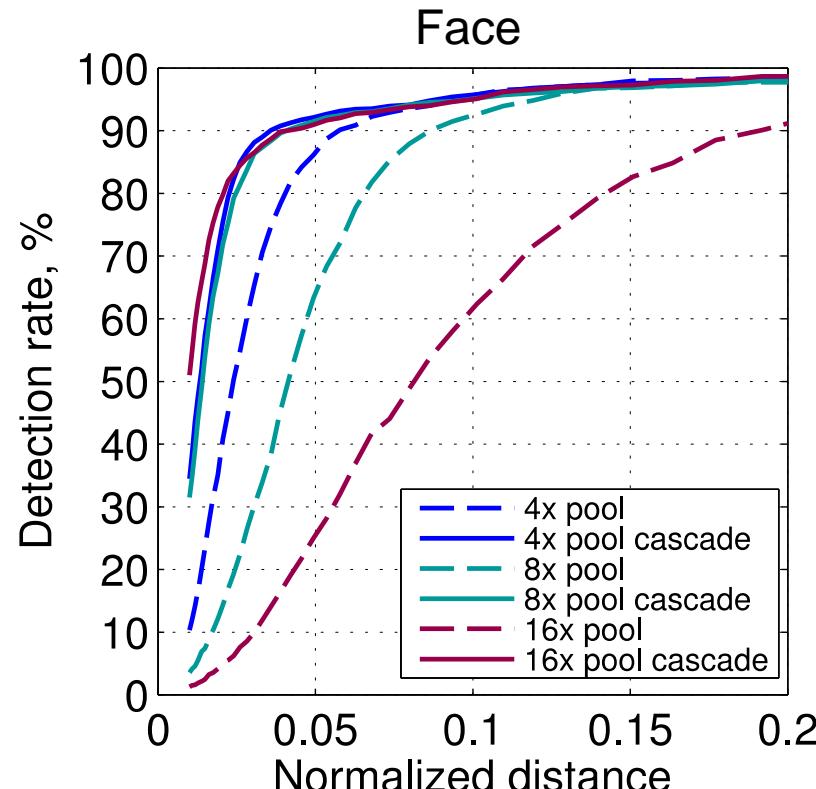
# RESULTS

## FLIC Dataset again

Small improvement for wrist:

Not enough context in fine model

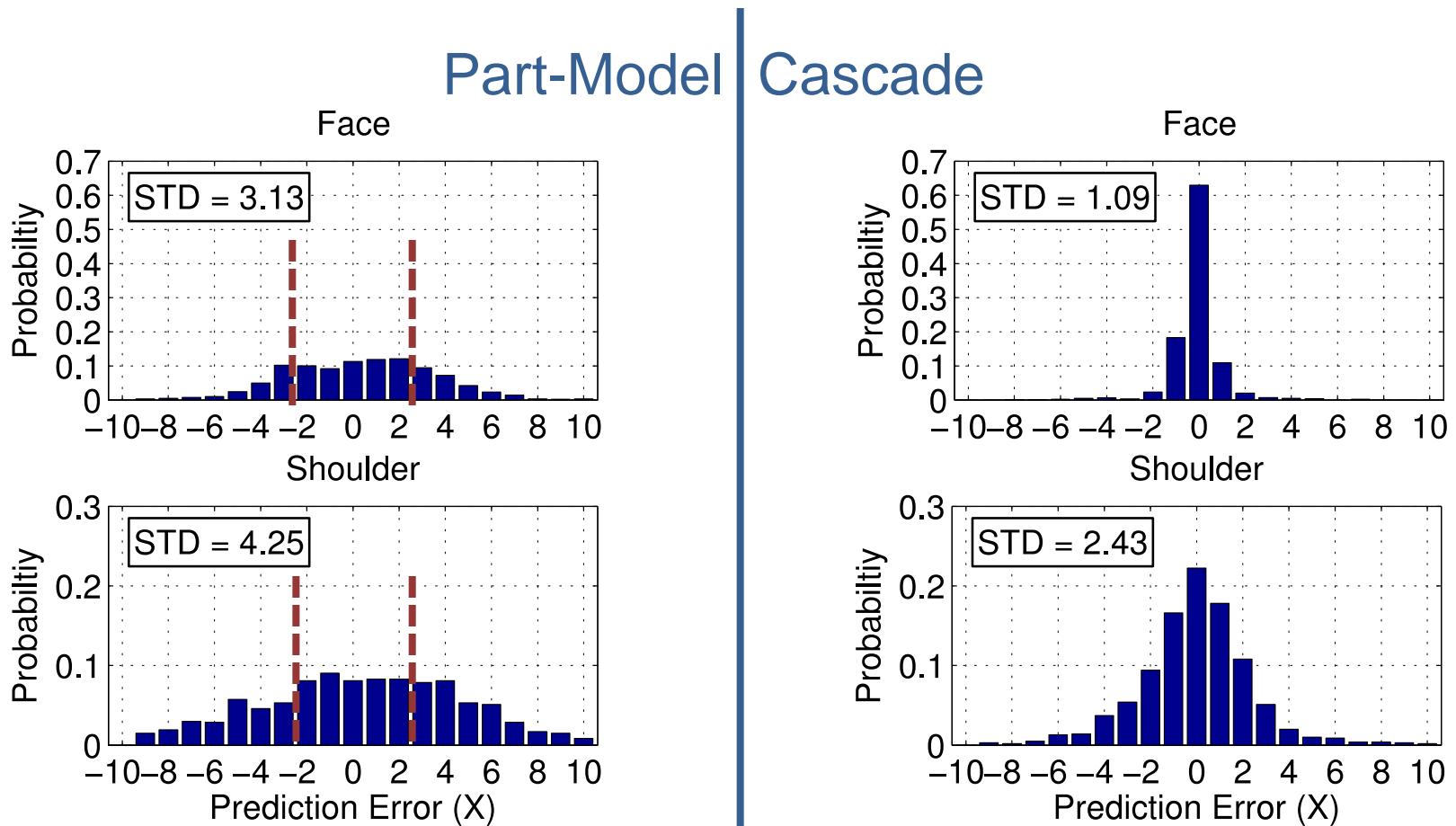
Too many mistakes in the coarse model



# How WELL ARE WE DOING?

## Evaluation of our model

Here we use a part model with 4x pooling ( $\pm 2$  pixel uncertainty)



# How WELL ARE WE DOING?

Informal study to evaluate human performance

1. Show users examples and explain desired joint location
2. Ask users to label 10 images from FLIC
3. Compare this to our performance



# How WELL ARE WE DOING?

## Comparing our ConvNet to Human Variance

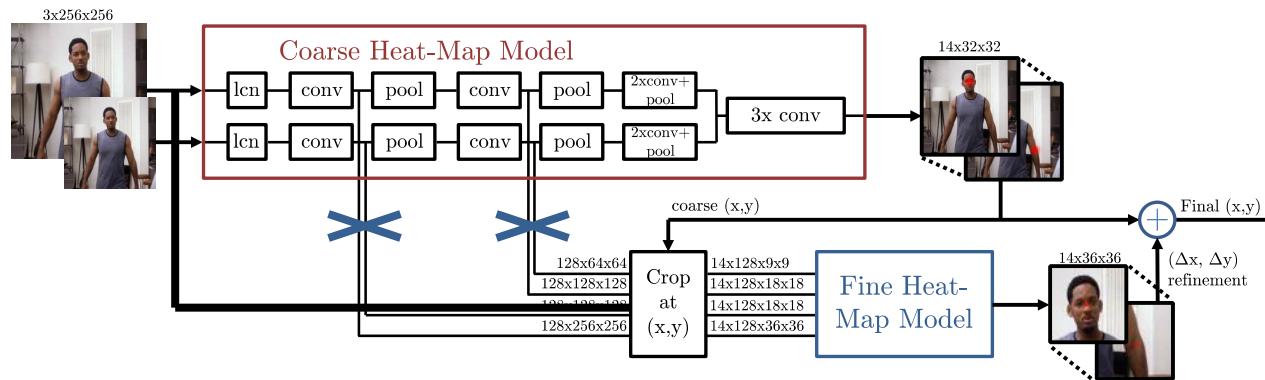
The 16x pooling module actually does quite well despite pooling.

	Face	Shoulder	Elbow	Wrist
Label Noise (10 images)	0.65	2.46	2.14	1.57
This work 4x (test-set)	1.09	2.43	2.59	2.82
This work 8x (test-set)	1.46	2.72	2.49	3.41
This work 16x (test-set)	1.45	2.78	3.78	4.16

# WELL, WHAT ABOUT A STANDARD CASCADE?

## Baseline experiment

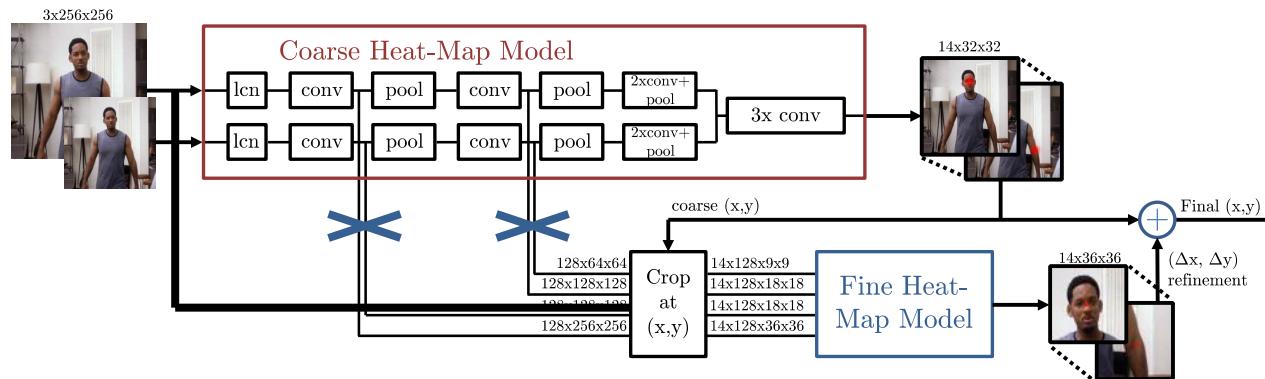
Construct a ConvNet that ONLY samples from the RGB



# WELL, WHAT ABOUT A STANDARD CASCADE?

## Baseline experiment

Construct a ConvNet that ONLY samples from the RGB

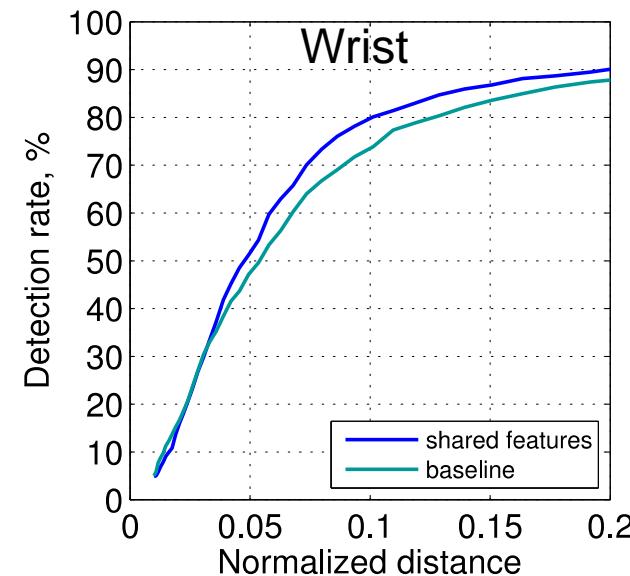


Shared features help

Regularization

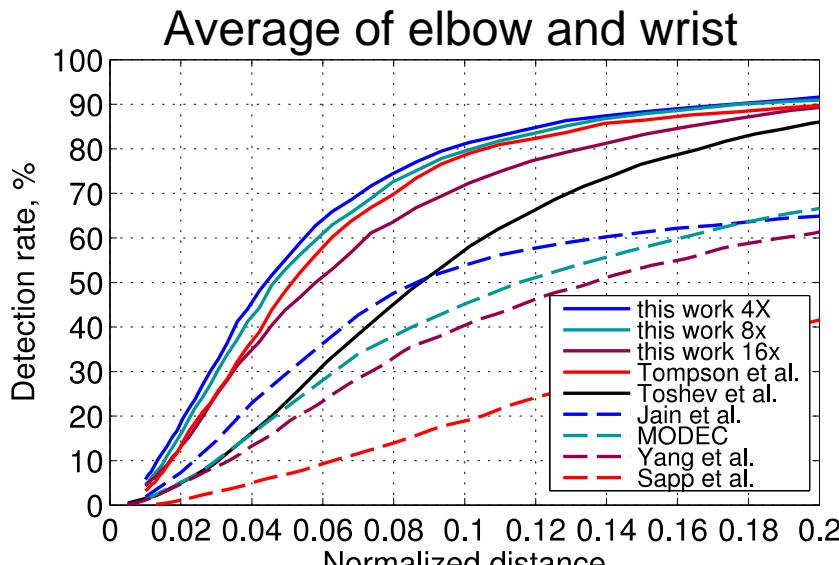
Shared capacity

More pooling → Higher delta

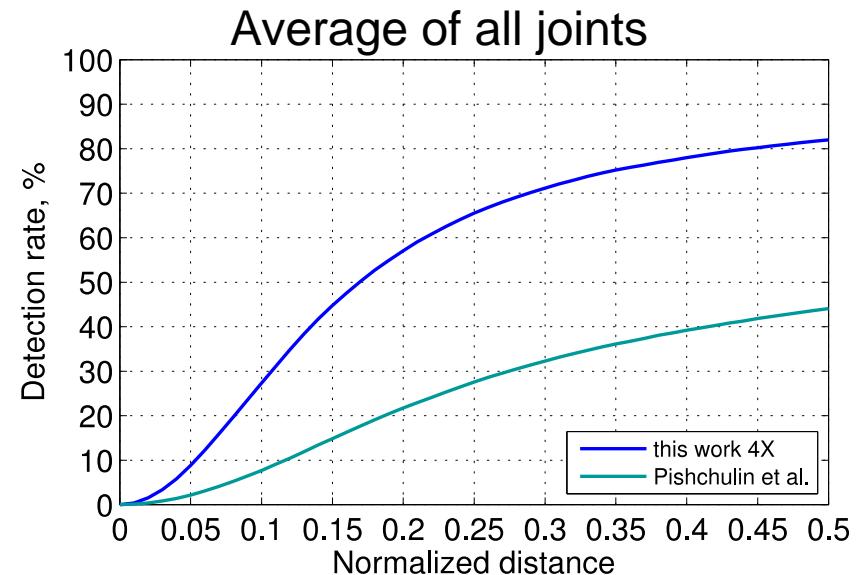


# COMPARISON WITH OTHER MODELS

## FLIC



## MPII



# CONCLUSION

# CONCLUSIONS

What have we shown?

ConvNets are \*really\* good at Object Recognition

Even with real-time constraints

# CONCLUSIONS

## What have we shown?

ConvNets are \*really\* good at Object Recognition

Even with real-time constraints

Simple networks in limited domains work well (hand tracking)

Can leverage “real” models to create amazing data

# CONCLUSIONS

## What have we shown?

ConvNets are \*really\* good at Object Recognition

Even with real-time constraints

Simple networks in limited domains work well (hand tracking)

Can leverage “real” models to create amazing data

For complex inputs we need domain optimizations

“belief propagation” within a neural network – enforce structure

Use multi-frame input – human motion is well defined and “learnable”

Multi-res input + Multi-res output – Offset spatial precision loss

# ACKNOWLEDGEMENTS

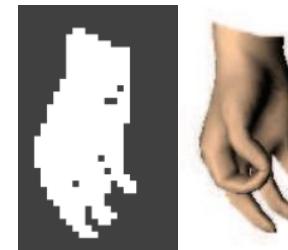
**QUESTIONS?**

**THANK YOU!**

# SOME RELATED WORK

Wang et al. & 3Gear (2009 to present)

Tiny images (nearest-neighbor)



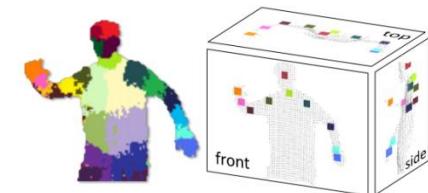
Oikonomidis et al. (2011,2012)

PSO search using synthetic depth images



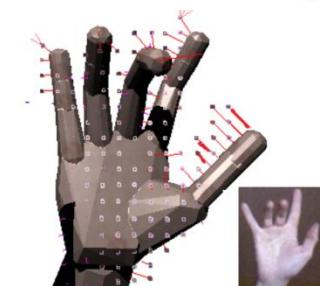
Shotton et al. (2011 and 2014-15)

Randomized Decision Forests



Melax et al. (2013)

Physics simulation (LCP)



Many more in the paper...

(Li, Weise, LeCun, Balan, Keskin, ...)

# FEATURE DETECTION

Infer 2D feature locations

Fingertips, palm, knuckles, etc.

Convolutional network (ConvNet) for feature inference

Efficient arbitrary function learner

Reasonably fast using modern GPUs

Self-similar features share learning capacity

# MULTI-RESOLUTION CONVNET

Downsampling (low pass) & local contrast normalization (high pass)

3 x banks with band-pass spectral density

CN convolution filter sizes constant

CN bandwidth context is high without the cost of large (expensive) filter kernels



# PSO/NM OBJECTIVE FUNCTION

$$F(C) = \sum_{s=1}^3 \left( \Delta(I_s, C) \right) + \Phi(C) + P(C)$$

L1 Depth comparison (multiple cameras)

$$\Delta(I_s, C) = \sum_{u,v} \min(|I_s(u, v) - R_s(C, u, v)|, d_{\max})$$

Coefficient prior (out-of-bound penalty)

$$\Phi(C) = \sum_{k=1}^n w_k [\max(C_k - C_{k\max}, 0) + \max(C_{k\min} - C_k, 0)]$$

Interpenetration constraint

Sum of bounding sphere inter penetrations

# IK OBJECTIVE FUNCTION

$$f(m) = \sum_{i=1}^n [\Delta_i(m)] + \Phi(C)$$

$$\Delta_i(m) = \begin{cases} \left\| (u, v, d)_i^t - (u, v, d)_i^m \right\|_2 & \text{If } d_i^t \neq 0 \\ \left\| (u, v)_i^t - (u, v)_i^m \right\|_2 & \text{otherwise} \end{cases}$$

$\Delta_i(m)$  is a L2 norm in 2D or 3D if there is depth image support for that pixel

Lots of problems... But it works

Use PrPSO to minimize  $f(m)$ : hard to parameterize and multi-modal (so gradient descent methods fail)

# SOME RELATED WORK

Andriluka et al. (CVPR 2009)

*Pictorial Structures Revisited: People Detection and Articulated Pose Estimation*

Shape context descriptors trained using ADABOOST + generative model and belief propagation

Felzenszwalb and Sirshick (PAMI 2010)

*Object Detection with Discriminatively Trained Part-Based Models*

Deformable Part Model (DPM) → “structure as latent variable”

Sapp and Taskar (CVPR 2013)

*Modec: Multimodal decomposable models for human pose estimation*

HOG + Graphical Model (and FLIC dataset)

Jain et al. (including me) (ICLR 2014)

*Learning Human Pose Estimation Features with Convolutional Networks*

Simple ConvNet with simple MRF spatial model

Toshev & Szegedy (CVPR 2014)

*DeepPose: Human pose estimation via deep neural networks*

Large ConvNet cascade to perform direct regression on UV locations

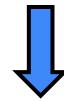
More in the paper (and more recent)...

(Pishchulin, Andriluka, LeCun, Johnson, Chen, ...)

# SPATIAL MODEL

Implement it as a network

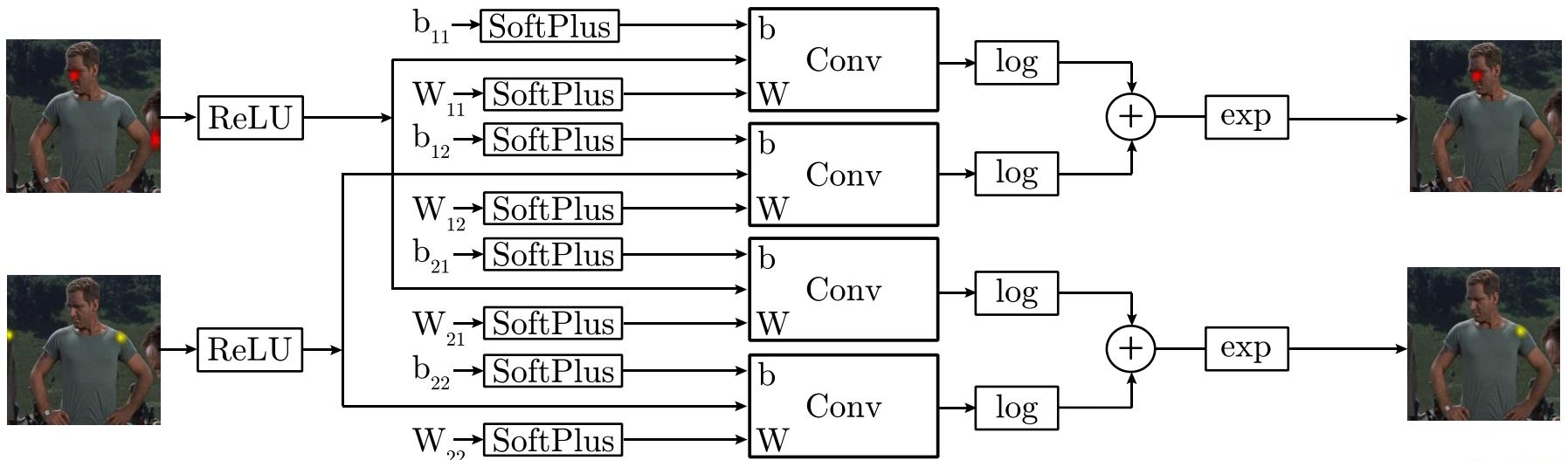
$$b(f) = \Phi(f) \prod_i (\Phi(x_i) * \Psi(f | x_i) + c(f | x_i))$$



$$\bar{e}_A = \exp \left( \sum_{v \in V} [\log (\text{SoftPlus}(e_{A|v}) * \text{ReLU}(e_v) + \text{SoftPlus}(b_{v \rightarrow A}))] \right)$$

where:  $\text{SoftPlus}(x) = 1/\beta \log(1 + \exp(\beta x))$ ,  $1/2 \leq \beta \leq 2$

$\text{ReLU}(x) = \max(x, \epsilon)$ ,  $0 < \epsilon \leq 0.01$



# JOINT TRAINING

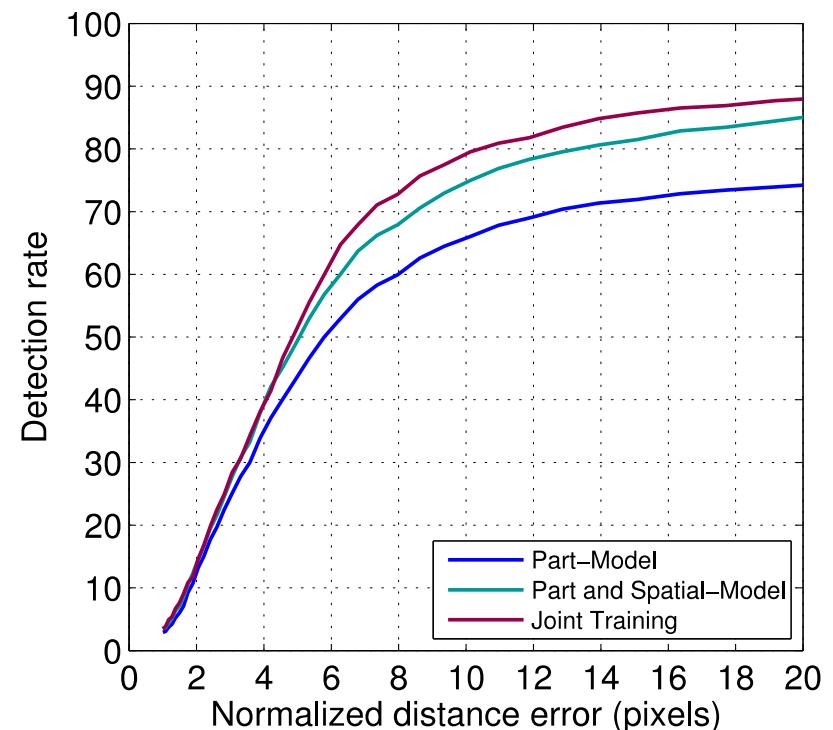
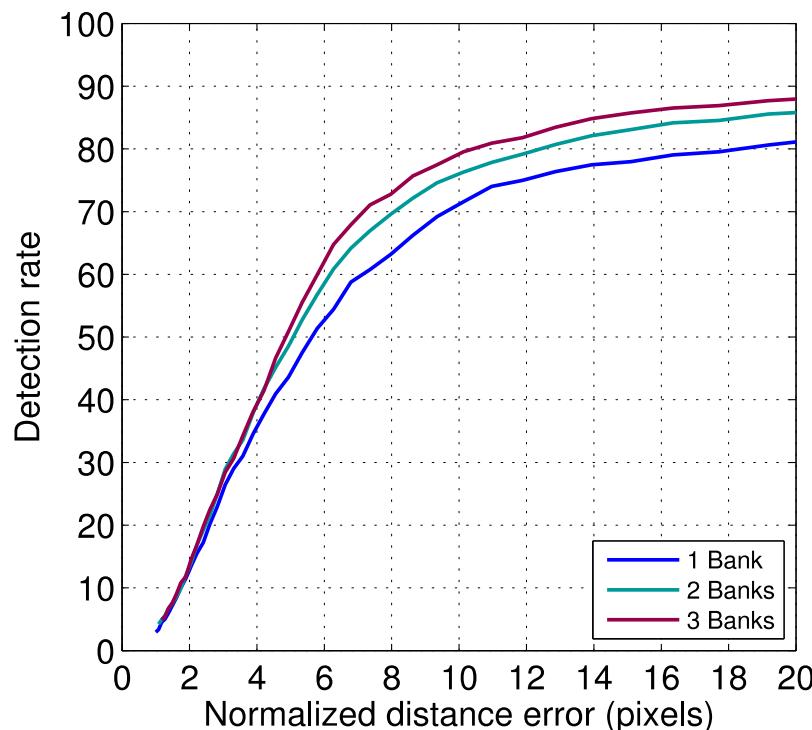
## Joint training

Pre-train both models separately

Joint train (BPROP) through both models

## Wrist Performance:

$$\text{DetectionRate}(R) = \frac{100}{N} \sum_{t=1}^N \left( \frac{\|x - x^t\|_2}{(\text{torso height } t)/100} \leq R \right)$$



# FLIC-PLUS

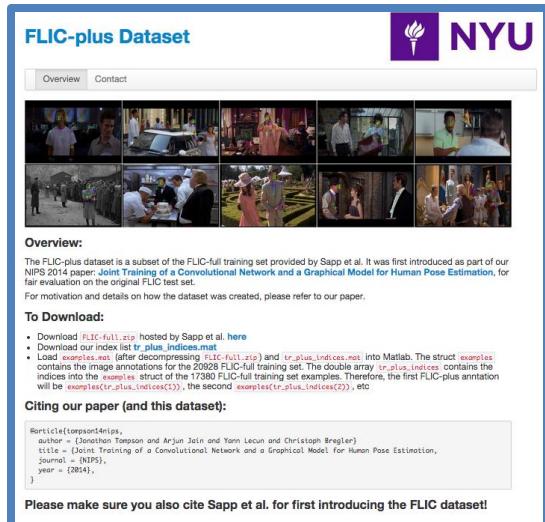
FLIC-Extended<sup>[1]</sup> is not fair!

20928 training / 1014 test samples

800 out of the 1014 test images (~80%) have an image in the training set that is at most 40 frames away

“FLIC-Plus”:

Use Mechanical Turk to find scenes in training and test sets  
Reject training set images from the same scene



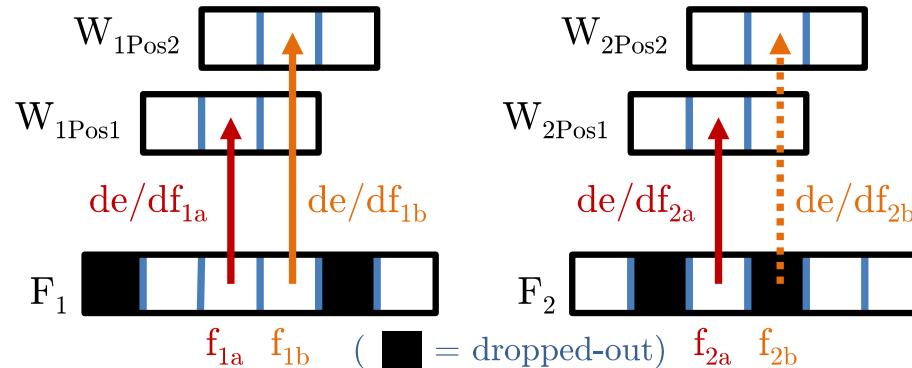
**cims.nyu.edu/~tompson/flic\_plus.htm**

# SPATIAL DROPOUT

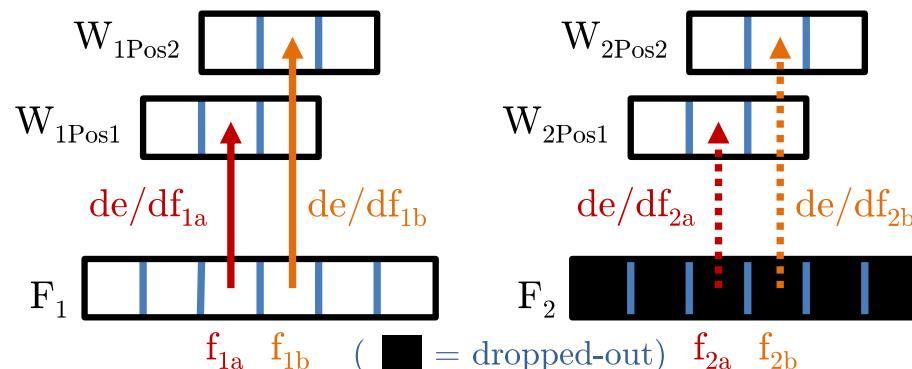
Dropout fails for fully-convolutional networks

Gradients are correlated because pixels are highly correlated

Dropout just scales the learning rate



Instead Dropout the entire feature



# How WELL ARE WE DOING?

## Time (ms) for FPROP

For 16x pooling: first 2 layers dominates runtime

	4x pool	8x pool	16x pool
Coarse-Model	140.0	74.9	54.7
Fine-Model	17.2	19.3	15.9
Cascade	157.2	94.2	70.6