

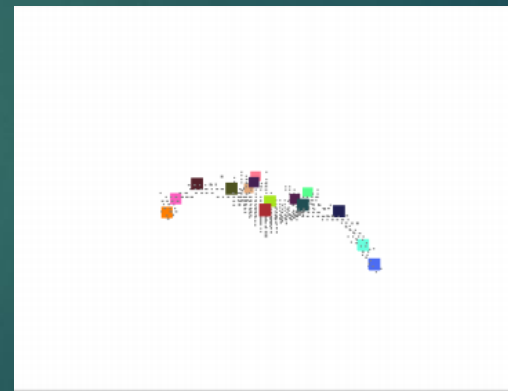
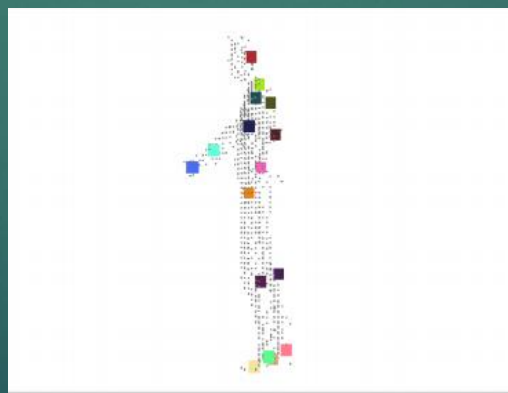
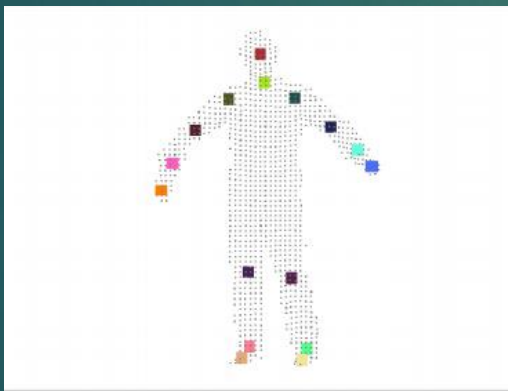
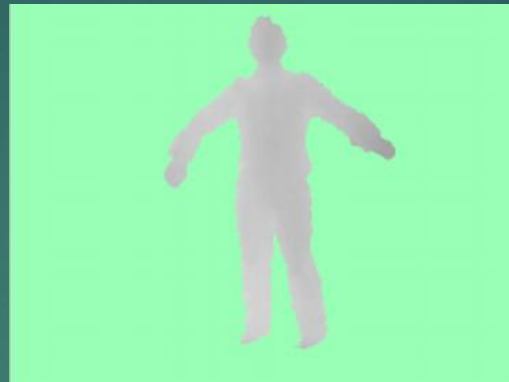


# Real-Time Human Pose Recognition in Parts from Single Depth Images

Jamie Shotton, Andrew Fitzgibbon, Mat Cook, Toby Sharp, Mark Finocchio, Richard  
Moore, Alex Kipman, Andrew Blake  
CVPR 2011

PRESENTER: AHSAN ABDULLAH

# PROBLEM

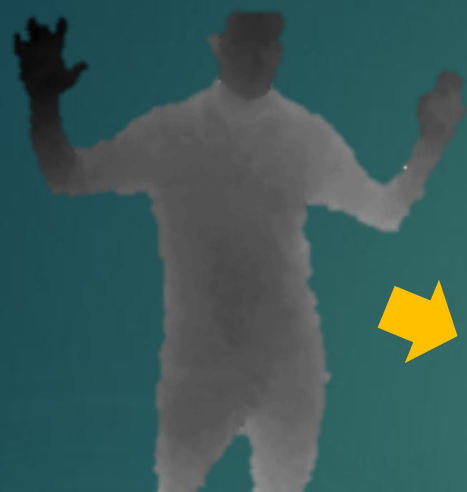


# APPROACH

- Partitioning into body parts helps localizing the joints



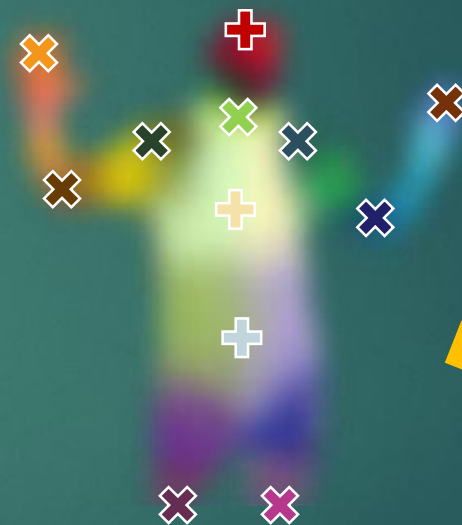
# PIPELINE



capture  
depth image &  
remove bg



infer  
body parts  
per pixel



cluster pixels to  
hypothesize  
body joint  
positions



fit model &  
track skeleton

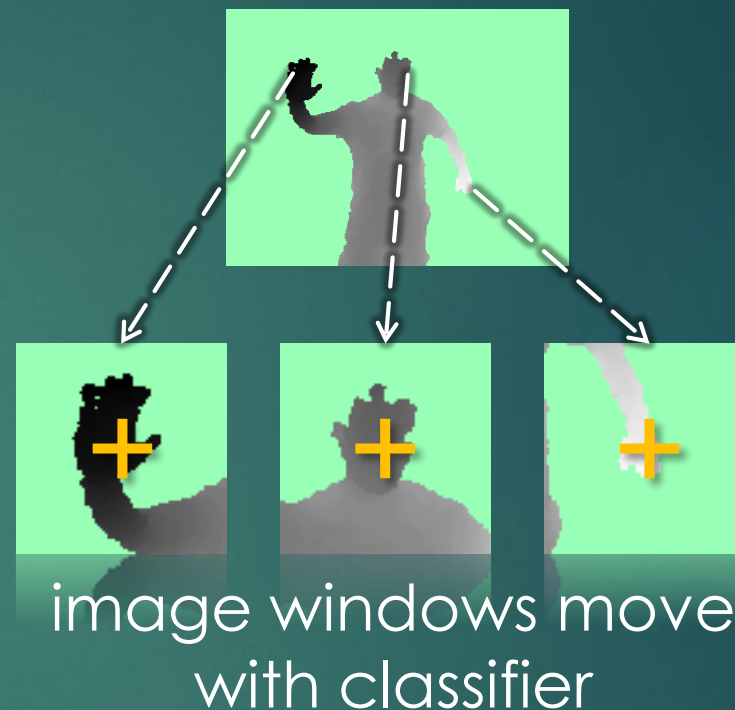
Design Goals

- Efficiency
- Robustness

# BODY PART CLASSIFICATION

- Compute  $P(c_i | w_i)$

- pixels  $i = (x, y)$
- body part  $c_i$
- image window  $w_i$



- Discriminative approach

- learn classifier  $P(c_i | w_i)$  from training data

# LEARNING DATA



**synthetic**  
*(train & test)*



**real**  
*(test)*



# LEARNING – DATA SYNTHESIS

Record MoCap  
500k frames  
distilled to 100k poses



Retarget to several models



Render (depth, body parts) pairs



# FEATURE SET

- Depth comparisons
  - very fast to compute

feature response  $f(I, \mathbf{x}) = d_I(\mathbf{x}) - d_I(\mathbf{x} + \Delta)$

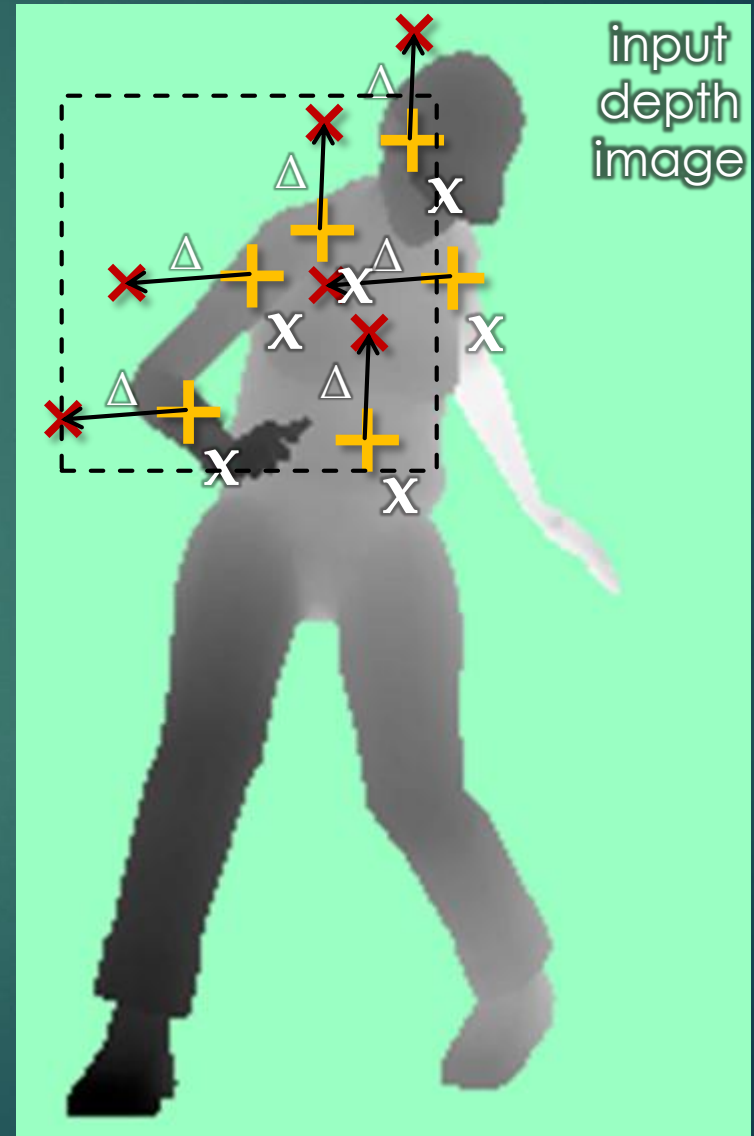
image coordinate  $\mathbf{x}$

image depth offset depth  $\Delta$

$$\Delta = \frac{\mathbf{v}}{d_I(\mathbf{x})}$$

scales inversely with depth

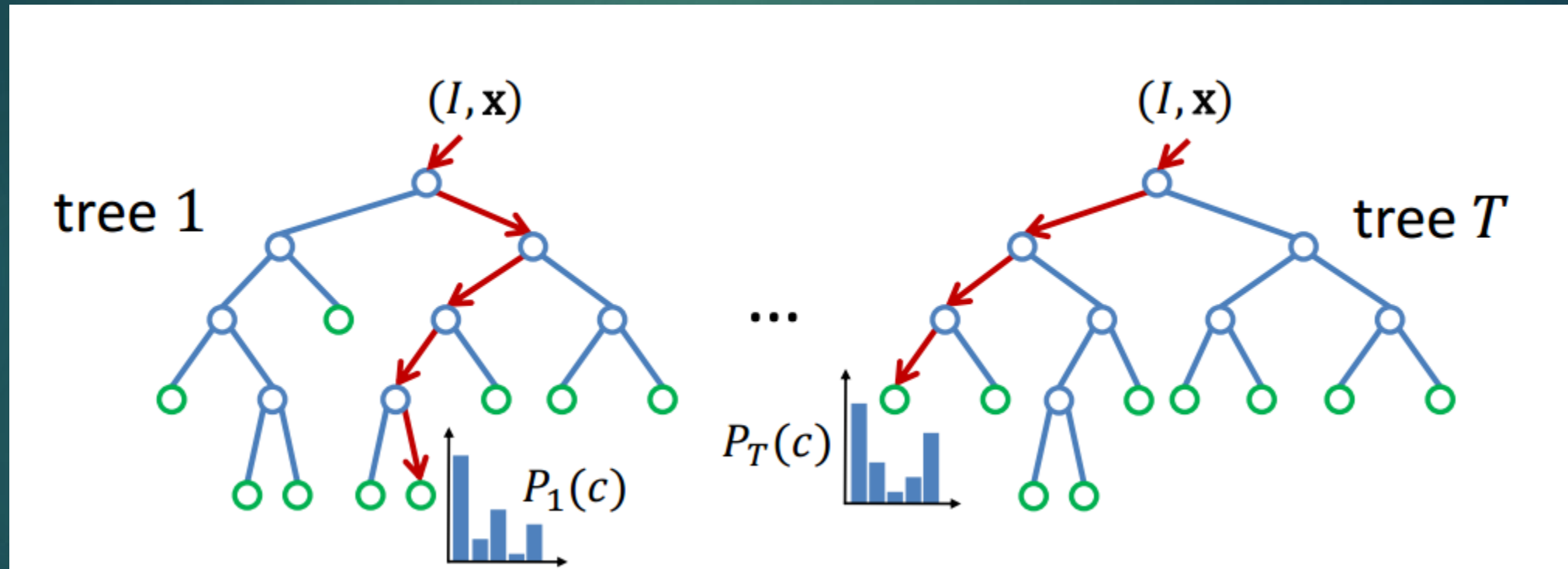
Background pixels  
 $d = \text{large constant}$



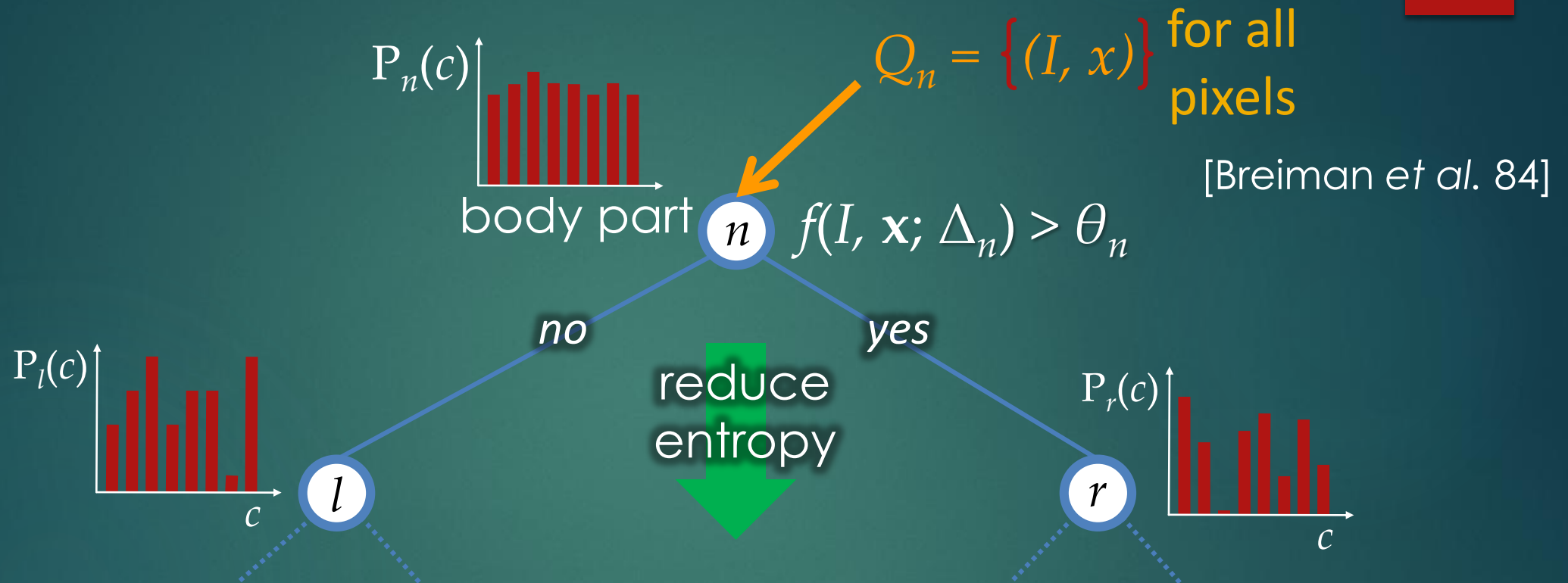


# DECISION FORESTS

- Aggregation of decision trees



# TRAINING DECISION TREES

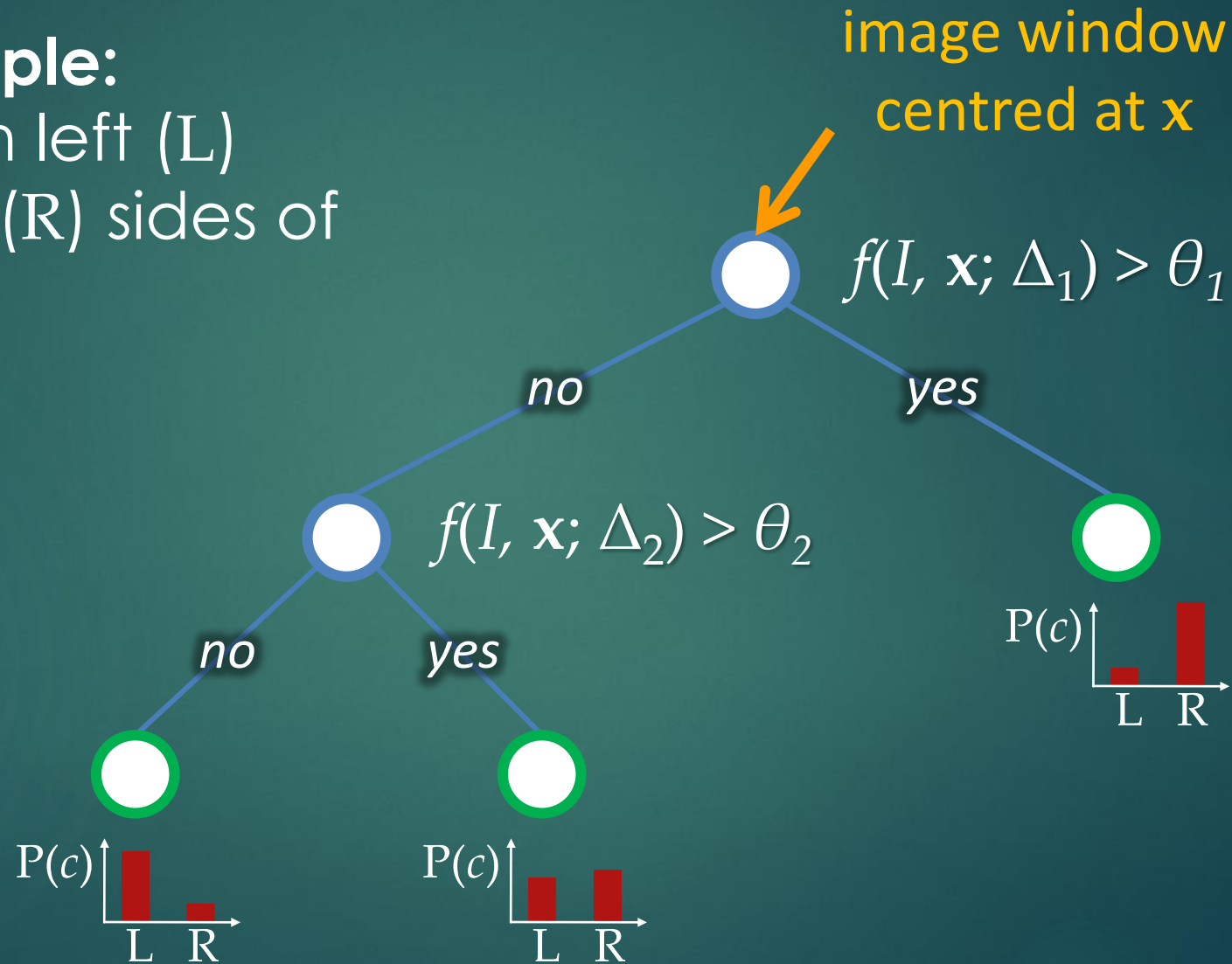


Take  $(\Delta, \theta)$  that maximises information gain

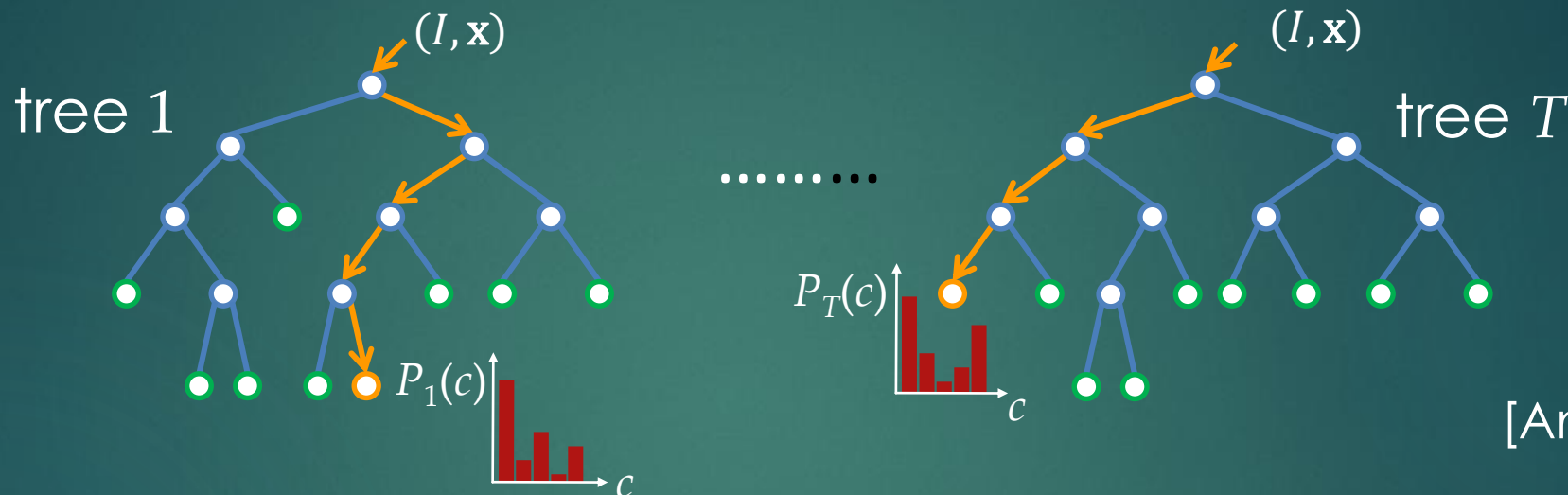
# DECISION TREE CLASSIFICATION

## Toy example:

Distinguish left (L) and right (R) sides of the body



# DECISION FOREST CLASSIFIER

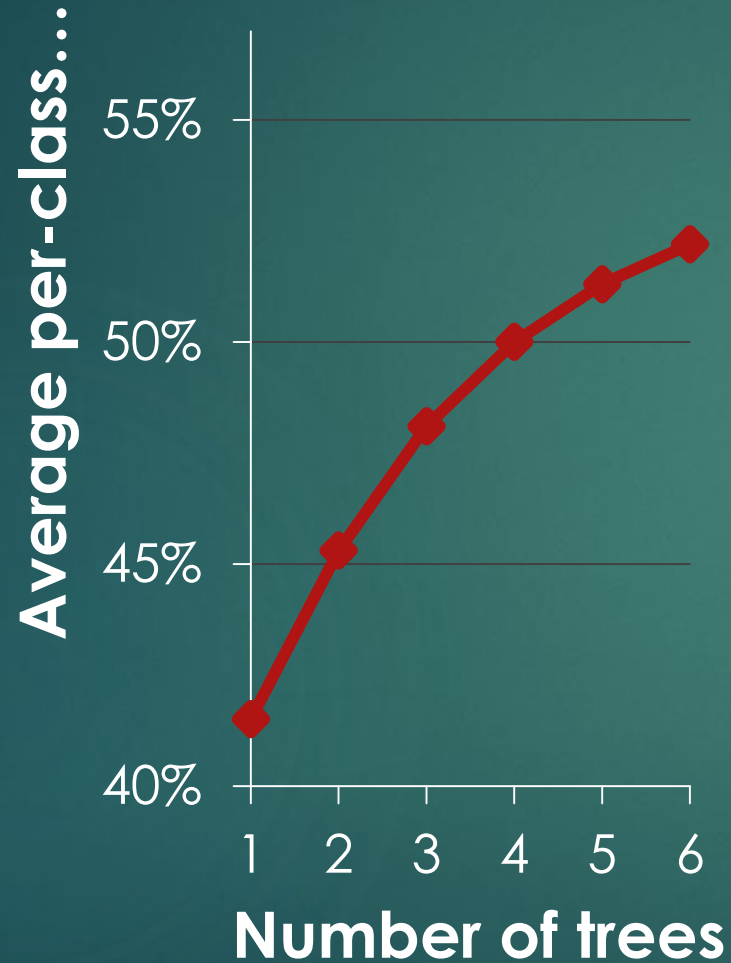


[Amit & Geman 97]  
[Breiman 01]  
[Geurts *et al.* 06]

- ▶ Trained on different random subset of images
  - ▶ “bagging” helps avoid over-fitting
- ▶ Average tree posteriors

$$P(c|I, \mathbf{x}) = \frac{1}{T} \sum_{t=1}^T P_t(c|I, \mathbf{x})$$

# NUMBER OF TREES



ground truth



inferred body parts (most likely)

1 tree



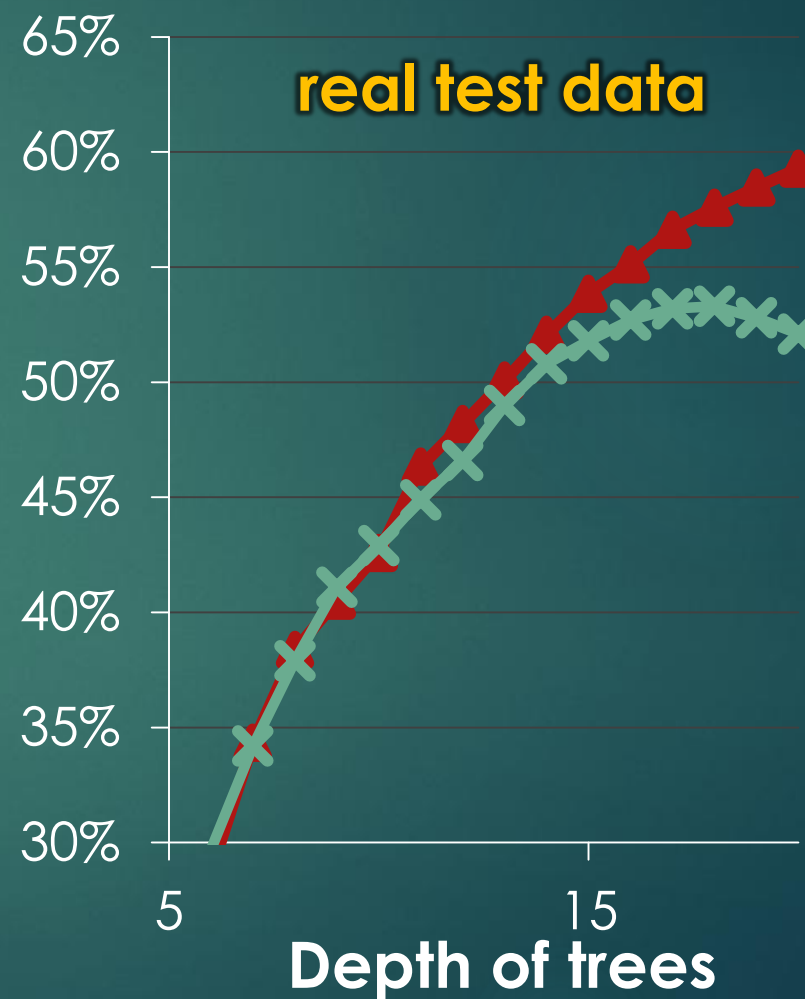
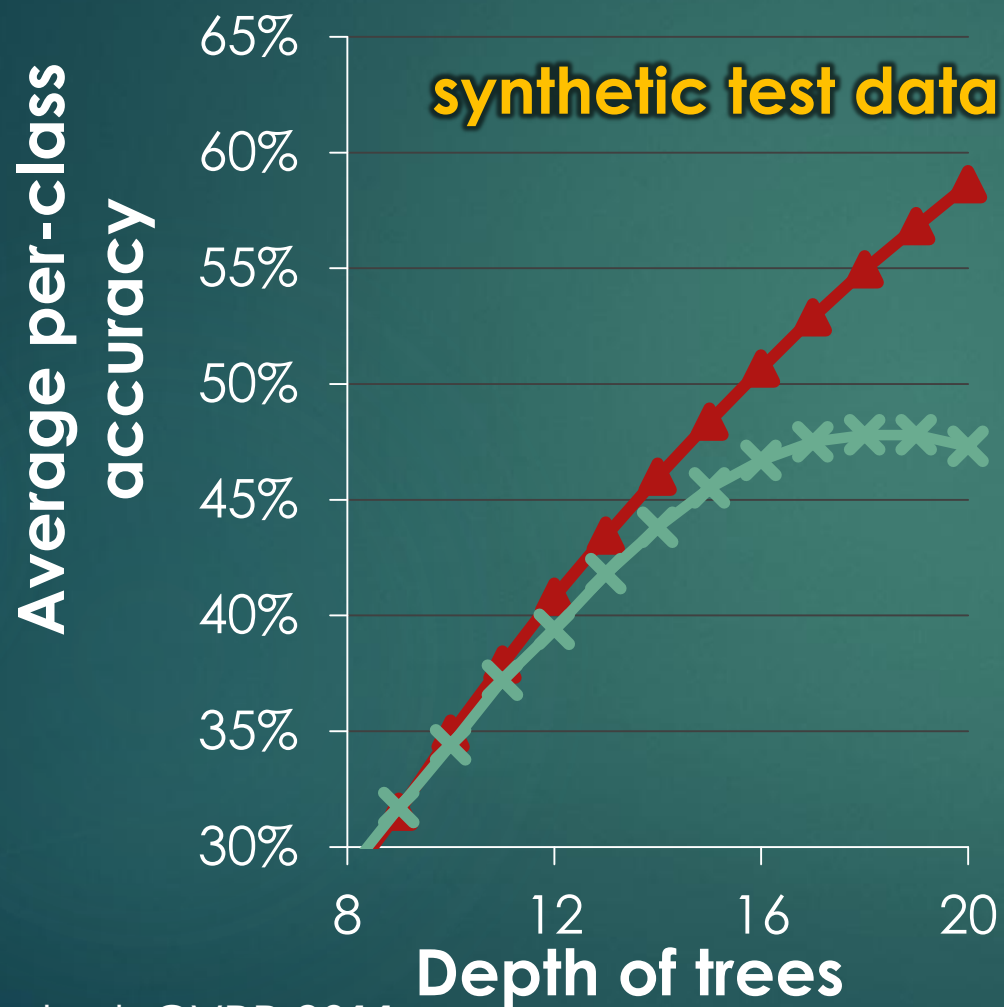
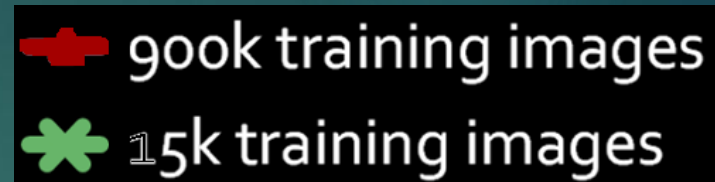
3 trees



6 trees



# TREE DEPTH





# Body parts to joint hypotheses

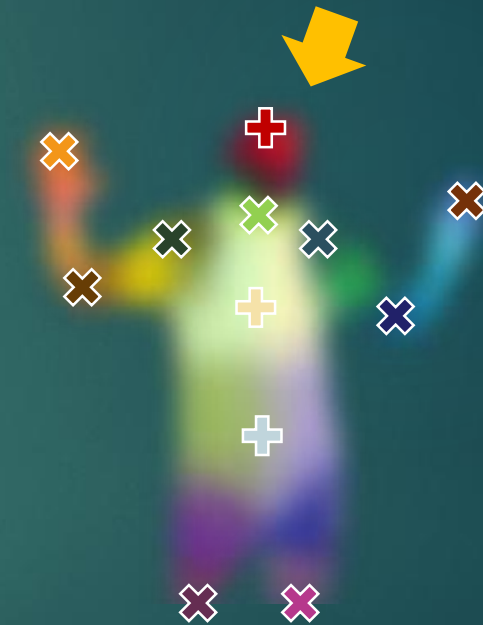
- Define 3D world space density

$$f_c(\hat{\mathbf{x}}) \propto \sum_{i=1}^N \underbrace{w_{ic}}_{\text{pixel weight}} \exp \left( - \left\| \frac{\underbrace{\hat{\mathbf{x}}}_{\text{3D coord}} - \underbrace{\hat{\mathbf{x}}_i}_{\text{3D coord of } i^{\text{th}} \text{ pixel}}}{\underbrace{b_c}_{\text{bandwidth}}} \right\|^2 \right)$$

$$w_{ic} = \underbrace{P(c|I, \mathbf{x}_i)}_{\text{inferred probability}} \cdot \underbrace{d_I(\mathbf{x}_i)^2}_{\text{depth at } i^{\text{th}} \text{ pixel}}$$

- Mean shift for mode detection

Shotton et. al. CVPR 2011



3. hypothesize body joints



...

**input depth**

**inferred body parts**



**front view**

**side view**

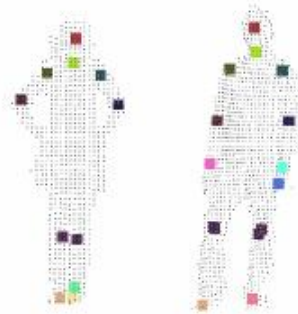
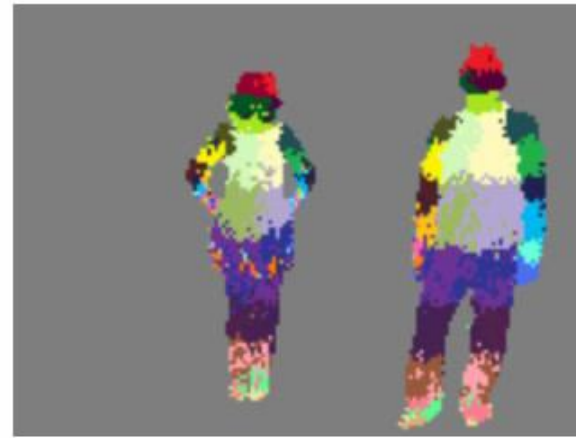
**top view**

**inferred joint positions**

**No tracking or smoothing**

**input depth**

**inferred body parts**



**front view**

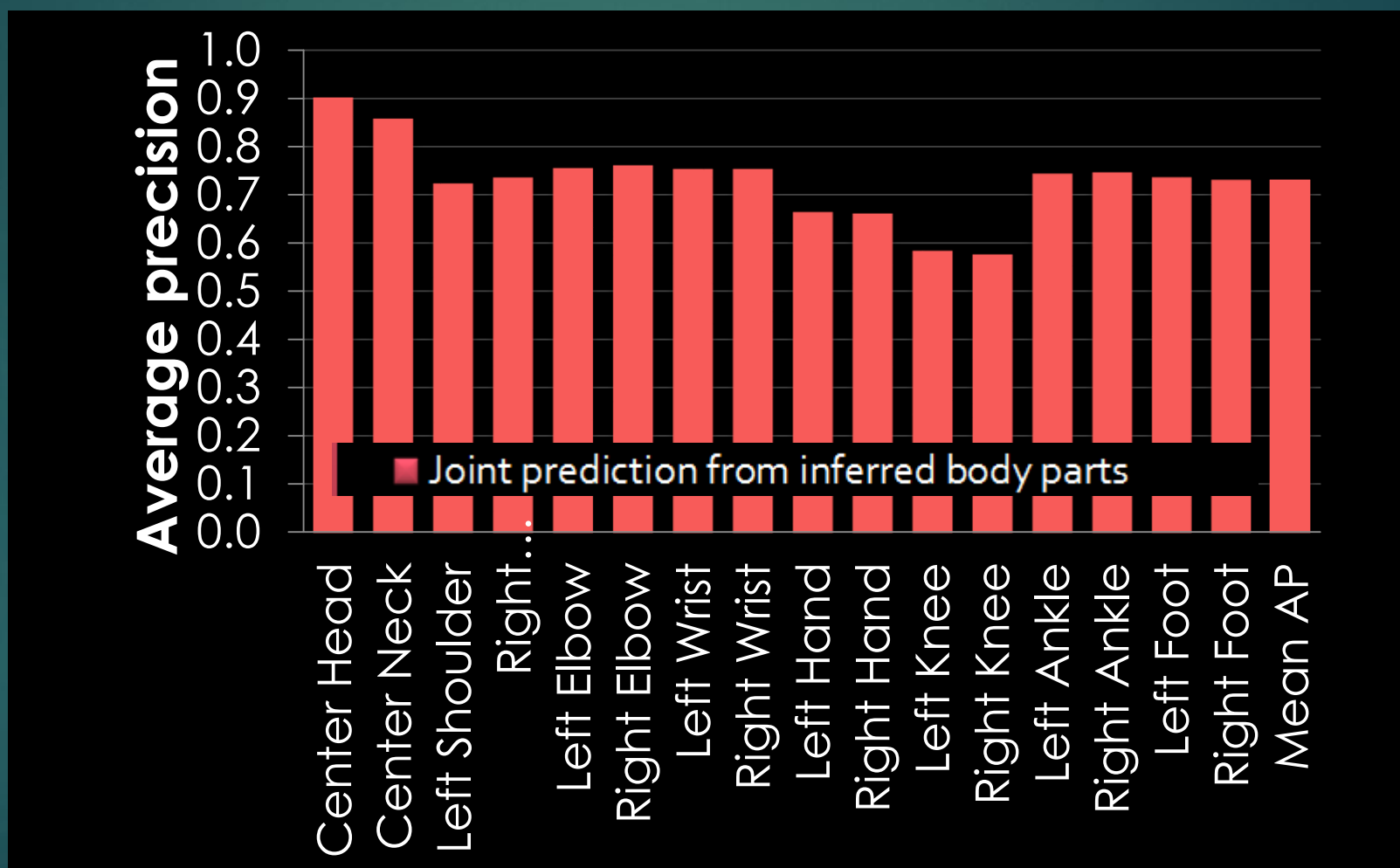
**side view**

**top view**

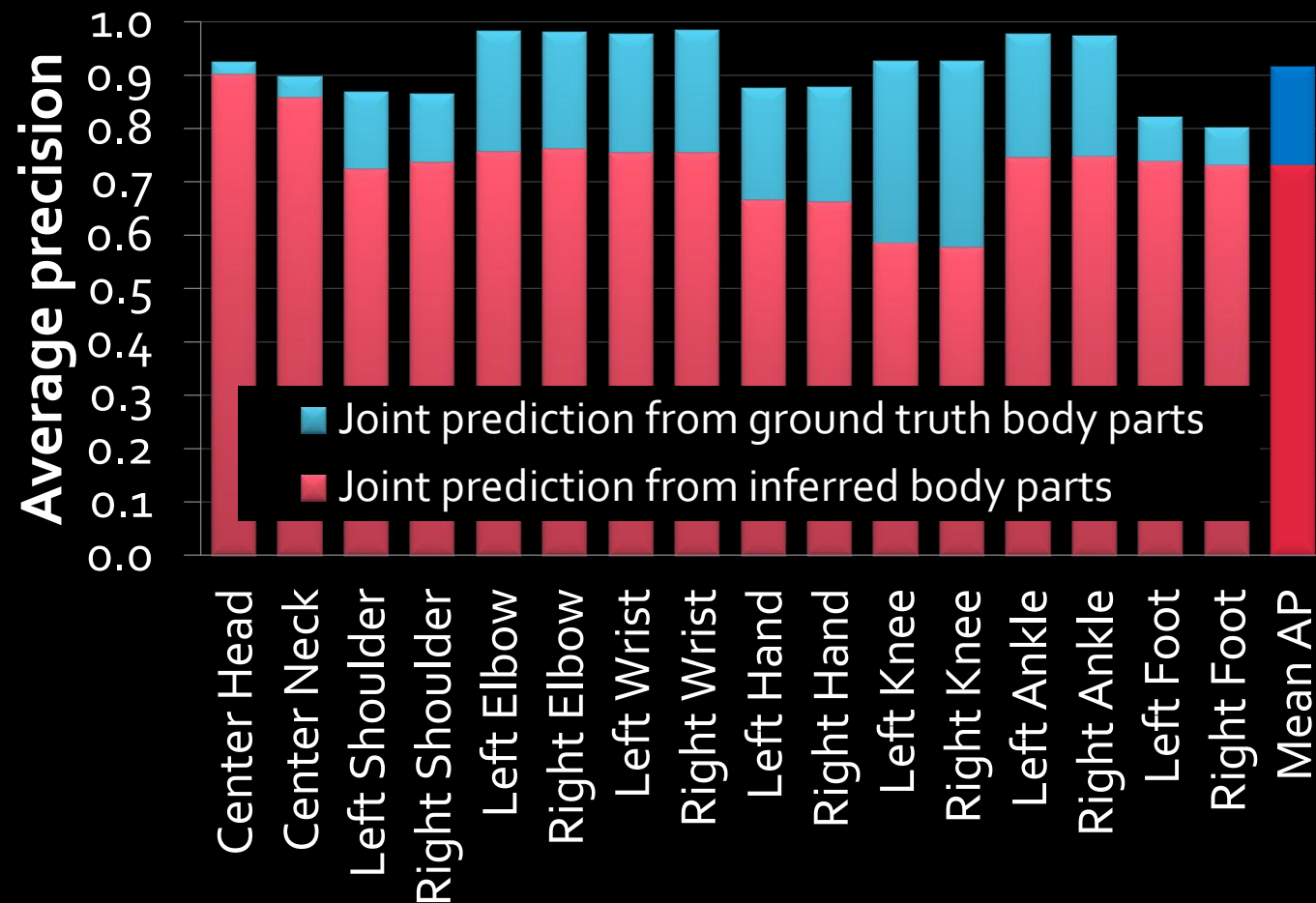
**inferred joint positions**

**No tracking or smoothing**

# JOINT PREDICTION ACCURACY



# JOINT PREDICTION ACCURACY



# ANALYSIS



- No temporal information
  - frame-by-frame
- Very fast
  - simple depth image feature
  - parallel decision forest classifier



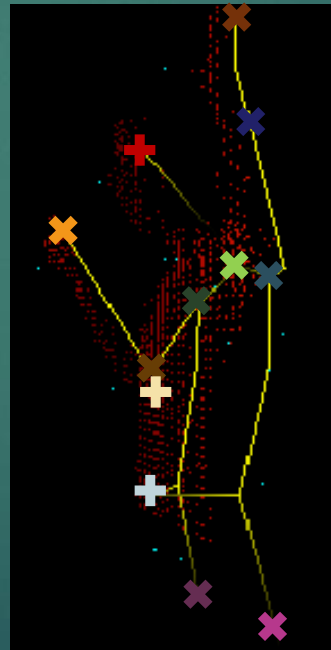
# KINECT SYSTEM

Uses...

- 3D joint hypotheses
- kinematic constraints
- temporal coherence

... to give

- full skeleton
- higher accuracy
- invisible joints
- multi-player



**4. track skeleton**



# SUMMARY



- Frame-by-frame gives robustness
- Body parts representation for efficiency
- Fast, simple machine learning
- Significant engineering to scale to a massive, varied training data set

# QUESTIONS

