# "GrabCut" — Interactive Foreground Extraction using Iterated Graph Cuts

Carsten Rother*            Vladimir Kolmogorov†            Andrew Blake‡
Microsoft Research Cambridge, UK

Figure 1: **Three examples of GrabCut**. The user drags a rectangle loosely around an object. The object is then extracted automatically.

## Abstract

The problem of efficient, interactive foreground/background segmentation in still images is of great practical importance in image editing. Classical image segmentation tools use either texture (colour) information, e.g. Magic Wand, or edge (contrast) information, e.g. Intelligent Scissors. Recently, an approach based on optimization by graph-cut has been developed which successfully combines both types of information. In this paper we extend the

free of colour blee[...]
degrees of interacti[...]
the labour-intensive[...]
background in a fe[...]

### 1.1    Previous [...]

In the following we describe briefly and compare several state of the art interactive tools for segmentation: Magic Wand, Intelligent [...]

SIGGRAPH2004

COMPUTER VISION AT MSRC            *Microsoft*

Microsoft® Office 2010

*Microsoft*®

# unwrap mosaics

Rav-Acha | Kohli | Rother | Fitzgibbon
http://research.microsoft.com/unwrap

[Shotton, Winn, Rother, Criminisi 06 + 08]
[Winn & Shotton 06]

[Shotton, Johnson, Cipolla 08]

COMPUTER VISION AT MSRC

Microsoft®

Ground truth | Entangled | Conventional

A. Montillo, J. Shotton, J. Winn, J. E. Iglesias, D. Metaxas, and A. Criminisi,
*Entangled Decision Forests and their Application for Semantic Segmentation of CT Images*,
in *Information Processing in Medical Imaging (IPMI)*, July 2011

*Microsoft*

# 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

Microsoft Research Cambridge & Xbox Incubation

## Abstract

*We propose a new method to quickly and accurately predict 3D positions of body joints from a single depth image, using no temporal information. We take an object recognition approach, designing an intermediate body parts representation that maps the difficult pose estimation problem into a simpler per-pixel classification problem. Our large* ... *dataset allows the classifier to* ... *nt to pose, body shape, clothing,* ... *nfidence-scored 3D proposals of* ... *ojecting the classification result* ... *frames per second on consumer* ... *shows high accuracy on both synthetic and real test sets, and investigates the effect of several training parameters. We achieve state of the art accu-*

**Figure 1. Overview.** From an single input depth image, a per-pixel body part distribution is inferred. (Colors indicate the most likely part labels at each pixel, and correspond in the joint proposals).

depth image ➡ body parts ➡ 3D joint proposals

**From**: Mark Finocchio
**To**: Jamie Shotton
**Date**: 11 Sept 2008
**Subject**:  Your computer vision expertise

Hi Jamie,

I work on Xbox Incubation and I noticed some work you've done on visual recognition using contours (http://jamie.shotton.org/work/research.html). I was hoping to be able to discuss an important scenario we are trying to solve with you. Would you be able to chat?

Thanks,

- Mark

THE CALL

*Microsoft*

**THE SCENARIO**

Microsoft®

Okada & Stenger 2008

Navaratnam *et al.* 2007

*Microsoft*

XBox prototype, Sept 2008

- Real time
- Accurate
- General poses

But…

- **Needs initialization**
- Limited body types
- Limited agility

[Hogg 1982]

Generative/Model-based

Discriminative/Regression

[Agarwal & Triggs 2004]
[Navaratnam & al 2007]

[Bourdev & Malik 09]

Detection

Tracking

[Gavrila 2000]

Whole

Parts

[Fischler & Elschlager 1973]

STATE OF THE ART

Microsoft

(a) Original picture.  (b) Differentiated picture.  (d) Rotated view.

1965. L. G. Roberts, **Machine Perception of Three Dimensional Solids**, in *Optical and electro-optical information processing*, J. T. Tippett (ed.), MIT Press.

MODEL-BASED VISION

*Microsoft*®

1980. J. O'Rourke and N. Badler. **Model-based image analysis of human motion using constraint propagation.** IEEE Trans. on Pattern Analysis and Machine Intelligence.

MODEL-BASED VISION

*Microsoft*®

# Model-based vision: a program to see a walking person

David Hogg

*For a machine to be able to 'see', it must know something about the object it is 'looking' at. A common method in machine vision is to provide the machine with general rather than specific knowledge about the object. An alternative technique, and the one used in this paper, is a model-based approach in which particulars about the object are given and this drives the analysis. The computer program described here, the WALKER model, maps images into a description in which a person is represented by the series of hierarchical levels, i.e. a person has an arm which has a lower-arm which has a hand. The performance of the program is illustrated by superimposing the machine-generated picture over the original photographic images.*

Keywords: vision, machine perception, WALKER model

## INTRODUCTION

Vision systems, both natural and artificial, require knowledge about the perceived objects, although the role played by this knowledge in the analytical process is unclear. Many techniques of machine vision seek to generate 3D structural descriptions without invoking object specific knowledge. An alternative is to adopt the 'model-based' approach wherein particular knowledge about the objects being sought drives the analysis.

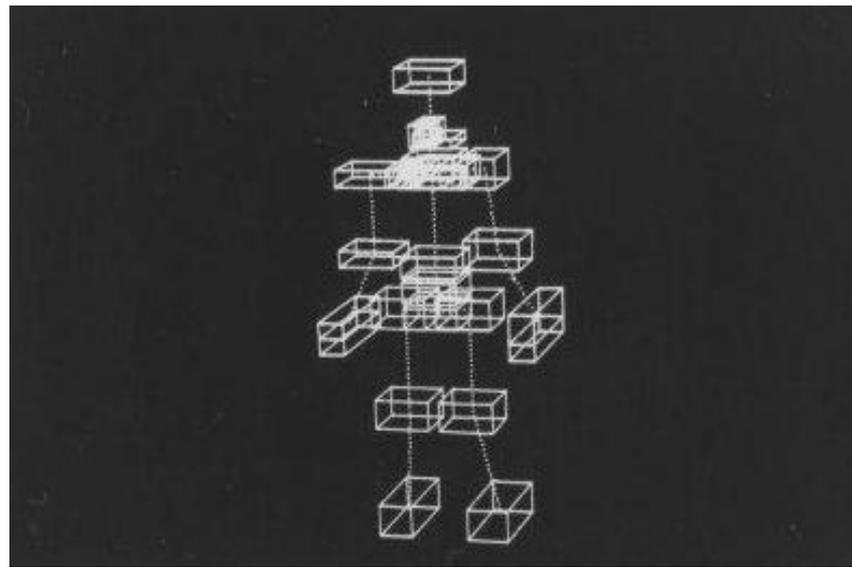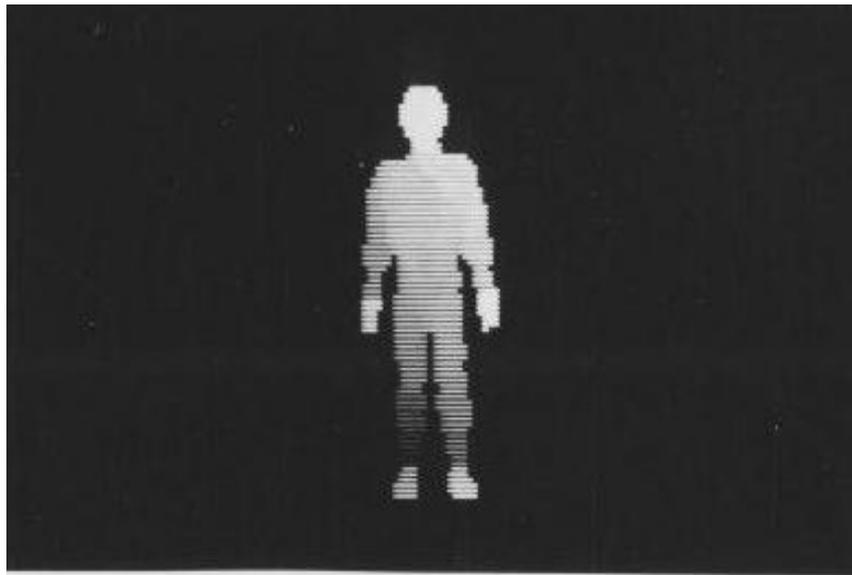This paper is concerned with a computer program that understands TV image sequences depicting a person walking through an arbitrary environment (Figure 1). The program maps given image sequences into a description in which the human body is represented by a collection of connected cylinders corresponding to its parts. It is supposed that such a 3D structural description would be both necessary and sufficient for many everyday tasks to be performed effectively. For example, touching someone's arm or deciding whether several people are marching in step all appear to require a grasp of 3D

School of Engineering and Applied Sciences, University of Sussex, Brighton, Sussex, UK

structure whether perceived visually or otherwise. Each output description is an instance of an abstract 3D model for a class of human walkers, henceforth called the WALKER model, itself an input to the program (Figure 2).

Descriptions generated by the program are sufficiently detailed to determine a pictorial reconstruction of the person from the perspective of the original imaging device. By superimposing these reconstructions over the original images a clear indication of the program's performance is visible to the human observer. When presented with the sequence depicted in Figure 1, the program generates as part of its output the sequence shown in Figure 3. The program copes with the enormous local ambiguity in an image by weighing evidence from across the image in support of a large number of possible interpretations. As a consequence, the program's performance should degrade gracefully for increasingly difficult image sequences in which the walker may be obscured or occulded to the camera.
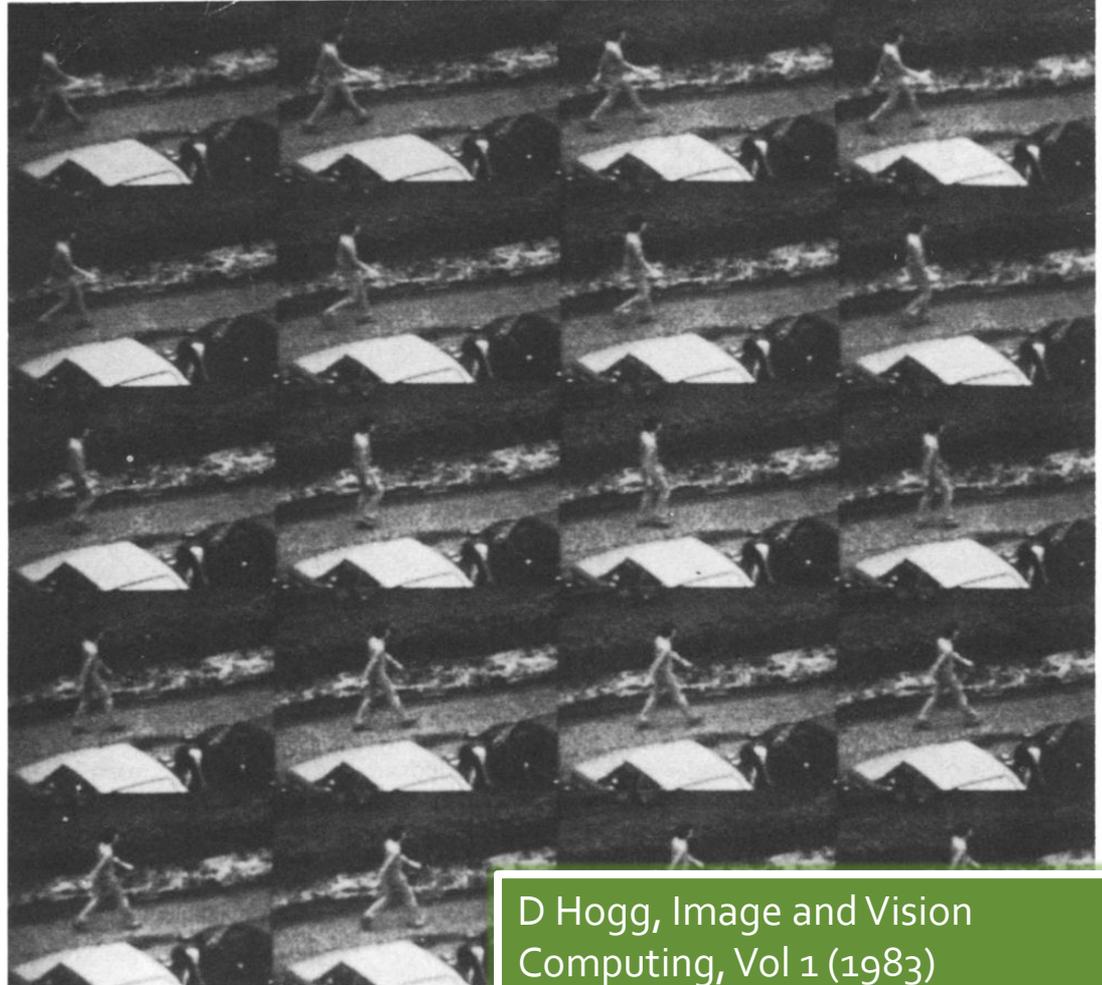
### Visual problem

The visual problem can be divided broadly into two parts; namely, what should be described and how can such descriptions be derived from a time-varying 2D image. It is impossible to divorce these two issues from one another since the difficulty of deriving a description from an image is bound to depend on the things being described. Moreover, certain representations may be required solely as intermediate descriptions for the interpretative process itself.

The question of what should be represented must depend on the visual system's function within a cognitive machine whose ultimate goal may be far removed from the visual world[1]. This paper takes a noncontroversial stand in accepting the usefulness of 3D structural descriptions as an interface to a larger system and instead concentrates on the second issue of how to generate such a description from an image.

### General-knowledge inference

Much of the current work in computer vision is concerned with the generation of 3D descriptions using only general-

# Model-based vision: a program to see a walking person

David Hogg

*For a machine to be able to 'see', it must know something about the object it is 'looking' at. A common method in machine vision is to provide the machine with general rather than specific knowledge about the object. An alternative technique, and the one used in this paper, is a model-based approach in which particulars about the object are given and this drives the analysis. The computer program described here, the WALKER model, maps images into a description in which a person is represented by the series of hierarchical levels, i.e. a person has an arm which has a lower-arm which has a hand. The performance of the program is illustrated by superimposing the machine-generated picture over the original photographic images.*

*Keywords: vision, machine perception, WALKER model*

## INTRODUCTION

Vision systems, both natural and artificial, require knowledge about the perceived objects, although the role played by this knowledge in the analytical process is unclear. Many techniques of machine vision seek to generate 3D structural descriptions without invoking object specific knowledge. An alternative is to adopt the 'model-based' approach wherein particular knowledge about the objects being sought drives the analysis.

This paper is concerned with a computer program that understands TV image sequences depicting a person walking through an arbitrary environment (Figure 1). The program maps given image sequences into a description in which the human body is represented by a collection of connected cylinders corresponding to its parts. It is supposed that such a 3D structural description would be both necessary and sufficient for many everyday tasks to be performed effectively. For example, touching someone's arm or deciding whether several people are marching in step all appear to require a grasp of 3D

School of Engineering and Applied Sciences, University of Sussex, Brighton, Sussex, UK

structure whether perceived visually or otherwise. Each output description is an instance of an abstract 3D model for a class of human walkers, henceforth called the WALKER model, itself an input to the program (Figure 2).

Descriptions generated by the program are sufficiently detailed to determine a pictorial reconstruction of the person from the perspective of the original imaging device. By superimposing these reconstructions over the original images a clear indication of the program's performance is visible to the human observer. When presented with the sequence depicted in Figure 1, the program generates as part of its output the sequence shown in Figure 3. The program copes with the enormous local ambiguity in an image by weighing evidence from across the image in support of a large number of possible interpretations. As a consequence, the program's performance should degrade gracefully for increasingly difficult image sequences in which the walker may be obscured or occulded to the camera.
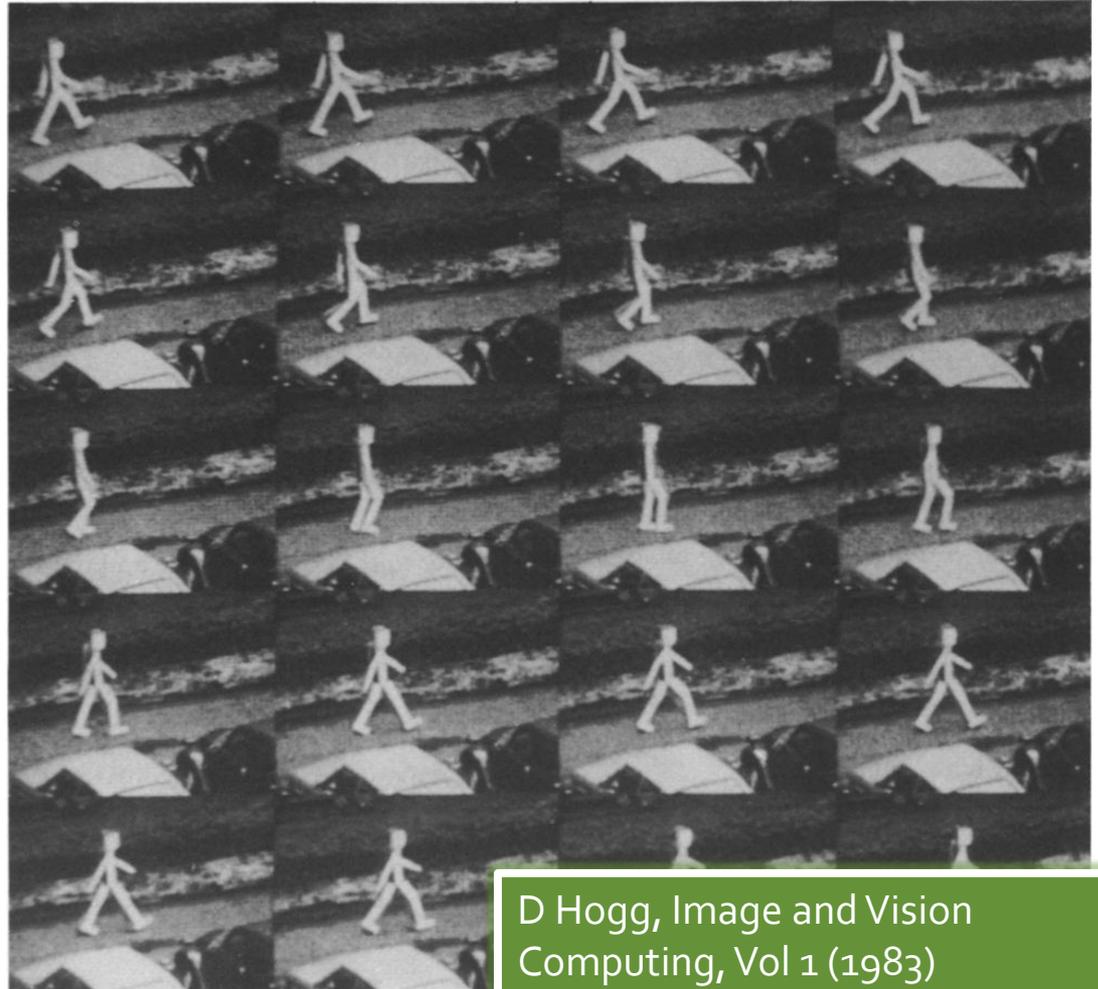
### Visual problem

The visual problem can be divided broadly into two parts; namely, what should be described and how can such descriptions be derived from a time-varying 2D image. It is impossible to divorce these two issues from one another since the difficulty of deriving a description from an image is bound to depend on the things being described. Moreover, certain representations may be required solely as intermediate descriptions for the interpretative process itself.

The question of what should be represented must depend on the visual system's function within a cognitive machine whose ultimate goal may be far removed from the visual world[1]. This paper takes a noncontroversial stand in accepting the usefulness of 3D structural descriptions as an interface to a larger system and instead concentrates on the second issue of how to generate such a description from an image.
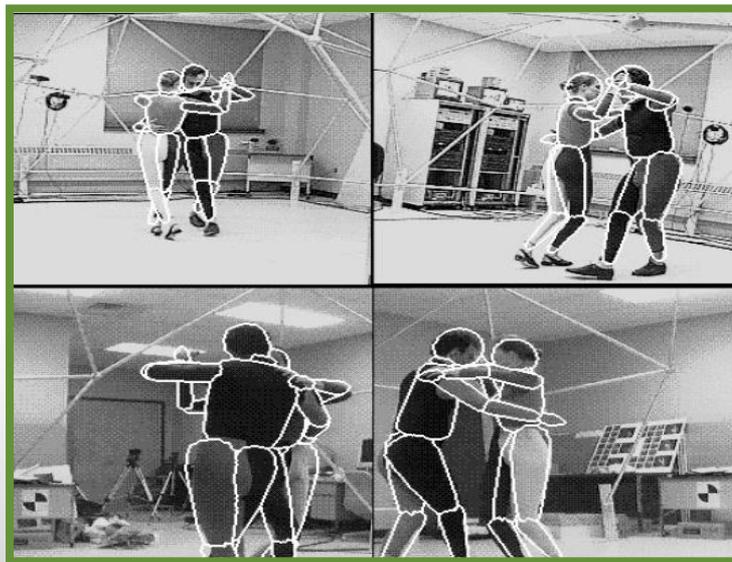
### General-knowledge inference

Much of the current work in computer vision is concerned with the generation of 3D descriptions using only general-

D Hogg, Image and Vision Computing, Vol 1 (1983)

# 3-D model-based tracking of humans in action: a multi-view approach

D.M. Gavrila and L.S. Davis
Computer Vision Laboratory, CfAR,
University of Maryland

MODEL-BASED VISION

- The "model" is *Forward/Generative/Graphical*
- Requiring *search* in many dimensions
  - say $10^{13}$ for the body

Resolved using
- (a) clever search: gradient descent and better
- (b) *temporal coherence*
  - Assume we were right in the *previous* frame
  - And search only "nearby" configurations in this

*Microsoft*®

# Exponential likelihood of failure

Assume 0.1% failure rate per frame

- After $n$ frames, chance of success = $0.999^n$

- At 30 frames per second, that's:
  - 3.0%        chance of failure after 1 second
  - 83.5%       chance of failure after 1 minute
  - 99.99%      chance of failure after 5 minutes

*Microsoft*®

# Exponential likelihood of failure

Assume 0.01% failure rate per frame

- After $n$ frames, chance of success = $0.9999^n$

- At 30 frames per second, that's:
  - 0.3%      chance of failure after 1 second
  - 16.5%    chance of failure after 1 minute
  - 59.3%    chance of failure after 5 minutes

- Need a method which works on a single frame
  - Single-frame methods all based on machine learning
  - So we'll need training data
  - Lots of training data
  - And will need to represent multiple hypotheses

*Microsoft*®

Paul A. Viola, Michael J. Jones
**Robust Real-Time Face Detection**
IEEE International Conference on Computer Vision, 2001

Microsoft

# Step Zero: Training data



$(\mathbf{z}^1, \boldsymbol{\theta}^1)$ | $(\mathbf{z}^2, \boldsymbol{\theta}^2)$ ... $(\mathbf{z}^i, \boldsymbol{\theta}^i)$ ... $(\mathbf{z}^N, \boldsymbol{\theta}^N)$

- Real home visits


- **Pose:** Motion capture
  - Standard "CMU" database
  - Custom database



- **Body size & shape:** Retargeting
  - Effects/Games industry tool: MotionBuilder

*Microsoft*®

**Actor wearing spherical markers**

**Observed by multiple cameras**

"MoCap"

**3D joint positions**

Computer Graphics

**Synthetic Depth Image**

*Microsoft*

- Standard motion capture datasets on the web

- Feed to *MotionBuilder* to generate 3D images

- Limited range of body types

*Microsoft*

Synthetic data:
realistic, but too clean

Artificially corrupted data

- depth resolution & noise
- rough edges
- missing pixels: hair/beards
- cropping & occlusions

*Microsoft*®

Image

Features

$$\boldsymbol{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix}$$

$f(z) \rightarrow \theta$

Joint angles

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

"Pose"

Andrew Blake, Kentaro Toyama,
**Probablistic tracking in a metric space**
IEEE International Conference on Computer Vision, 2001

# DETECTION VS. TRACKING

Microsoft

function

?

*Microsoft*®

function

Microsoft

function

Microsoft®

function

NEAREST NEIGHBOUR

*Microsoft*

Accuracy (1.0 is best)

Number of training images (log scale)

ALWAYS TRY NEAREST NEIGHBOUR FIRST

Microsoft®

Time taken: 500 milliseconds per frame

NEAREST NEIGHBOUR DIDN'T SCALE

*Microsoft*®

- Whole body $10^{12}$ poses (say)
- Four parts $4 \times 10^{3+\epsilon}$ poses
- But ambiguity increases

Microsoft®

# TRAINING DATA

Microsoft®

**Old (holistic) approach**

**New (parts) approach**

synthetic (held-out mocap poses)

real (from home visits)

TEST DATA

*Microsoft* 49

Input

Output

Input

Output

SLIDING WINDOW CLASSIFIER

Microsoft®

- Learn Prob(*body part* | *window*) from training data

EXAMPLE PIXEL 1: WHAT PART AM I?

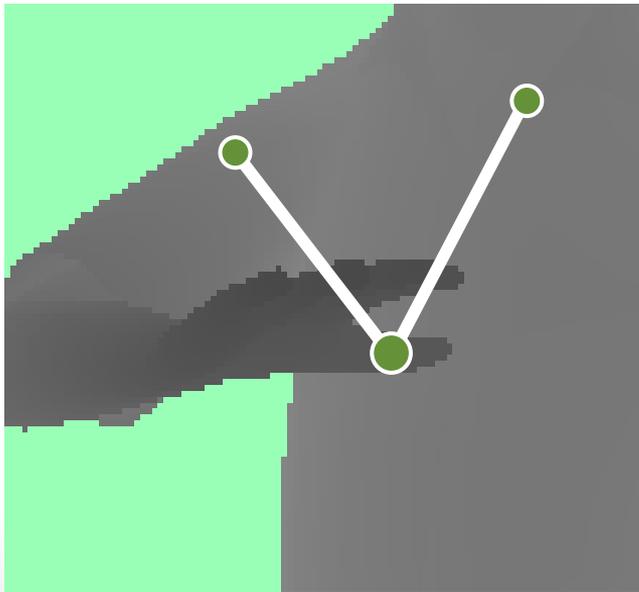EXAMPLE PIXEL 1: WHAT PART AM I?

$D_1 > 60$ mm

EXAMPLE PIXEL 2: WHAT PART AM I?

$D_1 > 60$ mm

$D_3 > 25$ mm

yes

no

yes

no

probability

probability

probability

head   l hand   r hand   l shoulder   r shoulder   chest   l elbow   r elbow
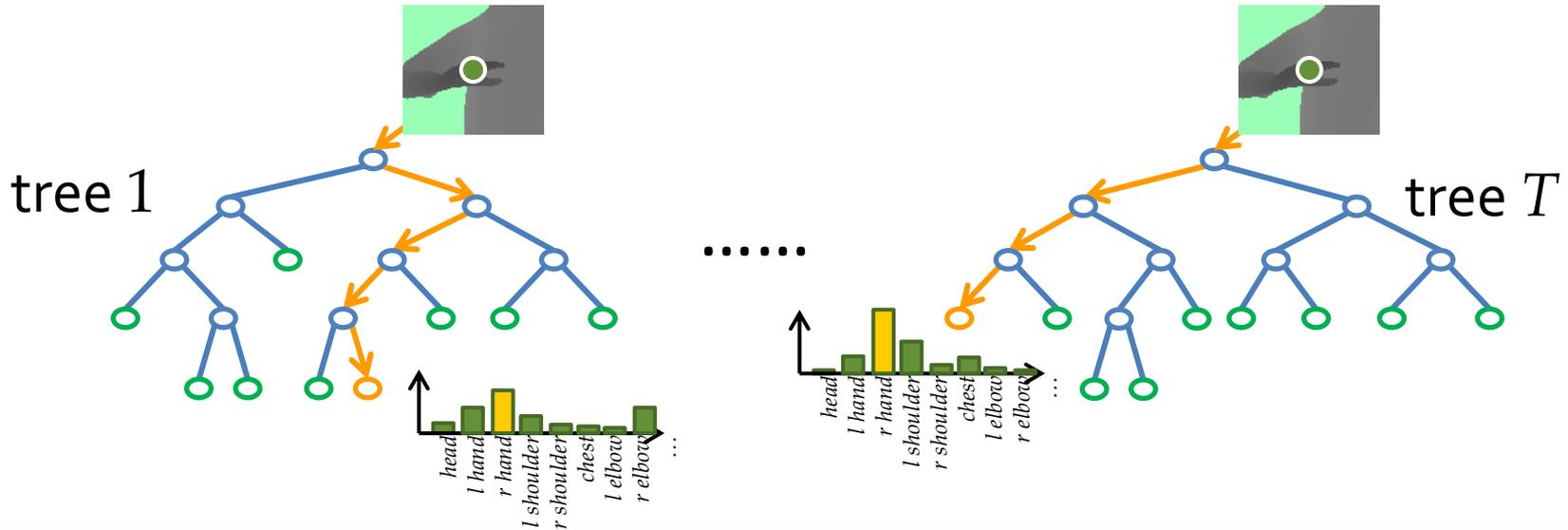
EXAMPLE PIXEL 2: WHAT PART AM I?

Microsoft

EXAMPLE PIXEL 2: WHAT PART AM I?

- **Same tree applied at every pixel**
- **Different pixels take different paths**
- **In practice, trees are much deeper**

$D_1 > 60$ mm

*yes*          *no*

$D_2 > 20$ mm          $D_3 > 25$ mm

*yes*          *no*          *yes*          *no*

*Microsoft*®

- A forest is an ensemble of trees:



tree 1 ...... tree $T$

- Helps avoid over-fitting during training
- Testing takes average of leaf nodes distributions

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

ground truth

inferred body parts (most likely)
1 tree          3 trees          6 trees

55%
50%
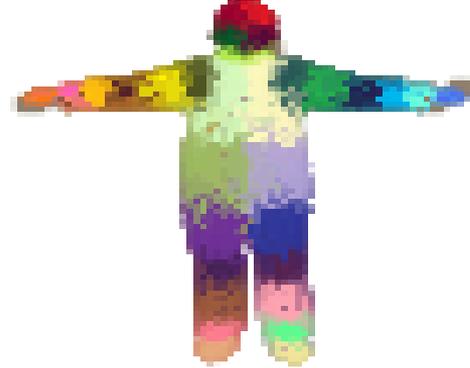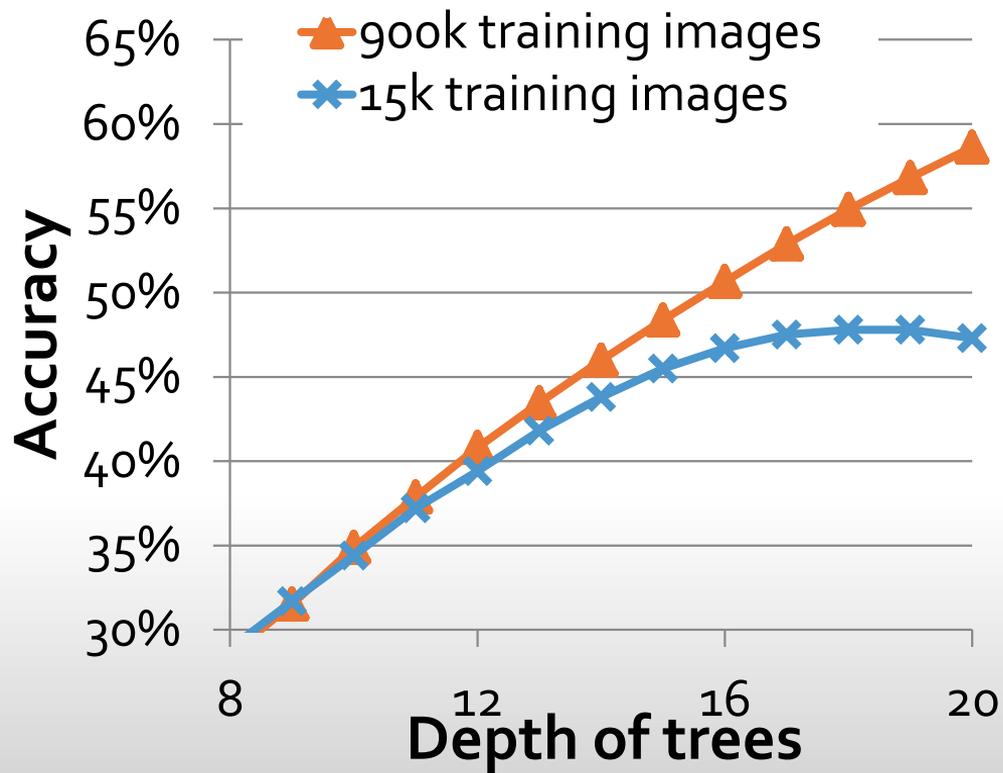45%
40%

Accuracy

1  2  3  4  5  6
Number of trees

*Microsoft*®

input depth     ground truth parts     inferred parts (soft)

depth 18
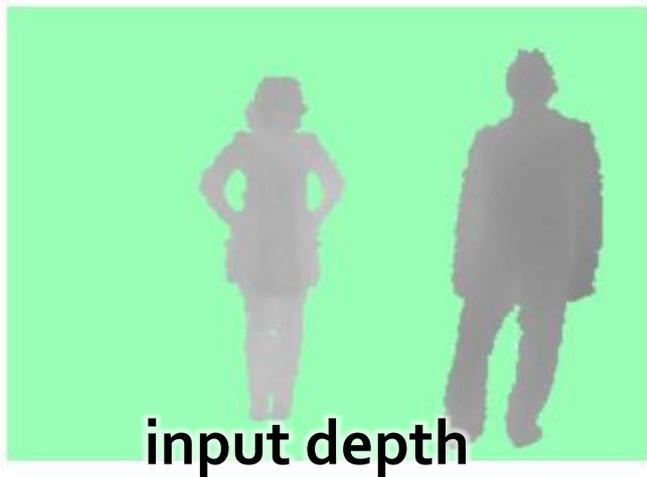
Microsoft

- Given
  - depth image
  - inferred body part probabilities

- Cluster high probability parts in 3D

hypothesized
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

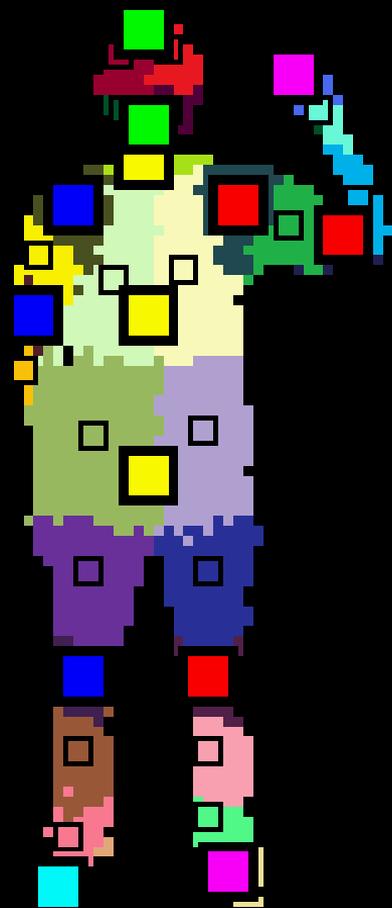MATCHING BODY *PARTS* IS BETTER

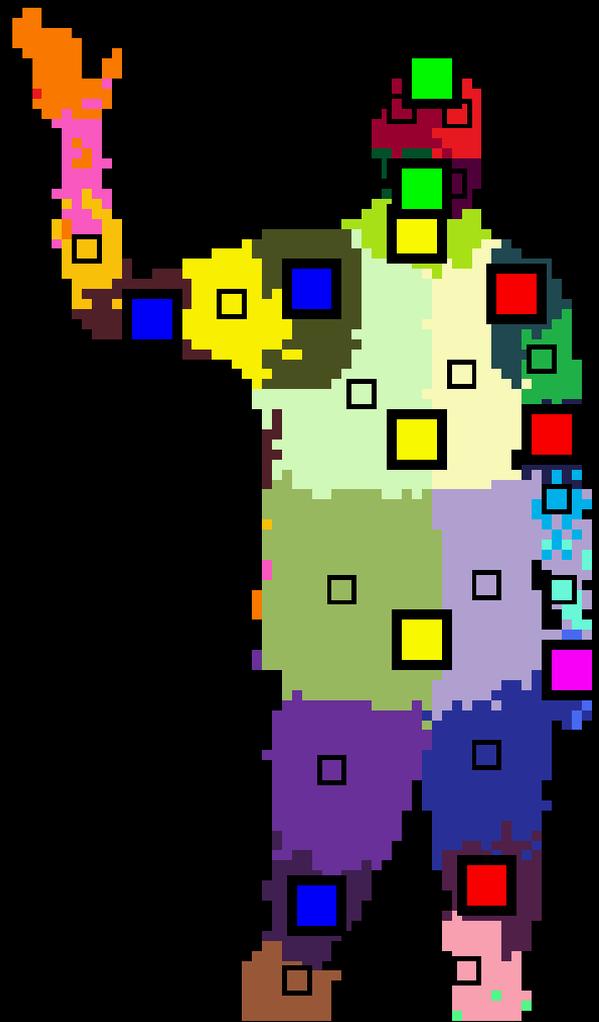Joint position hypotheses are not the whole story

Follow up with skeleton fitting incorporating
- Kinematic constraints (limb lengths etc)
- Temporal coherence (it's back!)

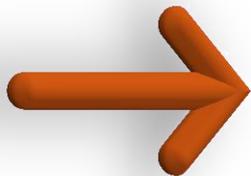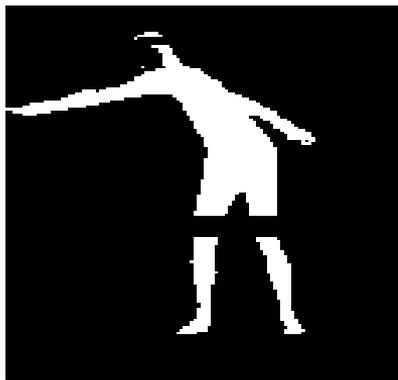And of course the incredible imagination of games designers…

# WRAPPING UP

Joint position hypotheses are not the whole story

Follow up with skeleton fitting incorporating
- Kinematic constraints (limb lengths etc)
- Temporal coherence (it's back!)

And of course the incredible imagination of games designers... and

YOU!

# WRAPPING UP
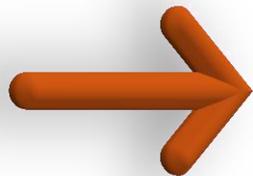
[Aside]

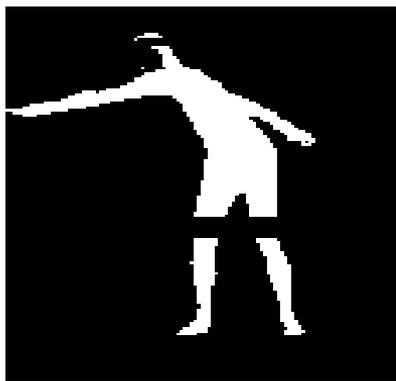# MULTIPLE HYPOTHESES AND REWRITING HISTORY

$$p(\boldsymbol{\theta}^{\mathrm{new}} | \mathbf{z}^{\mathrm{new}})$$
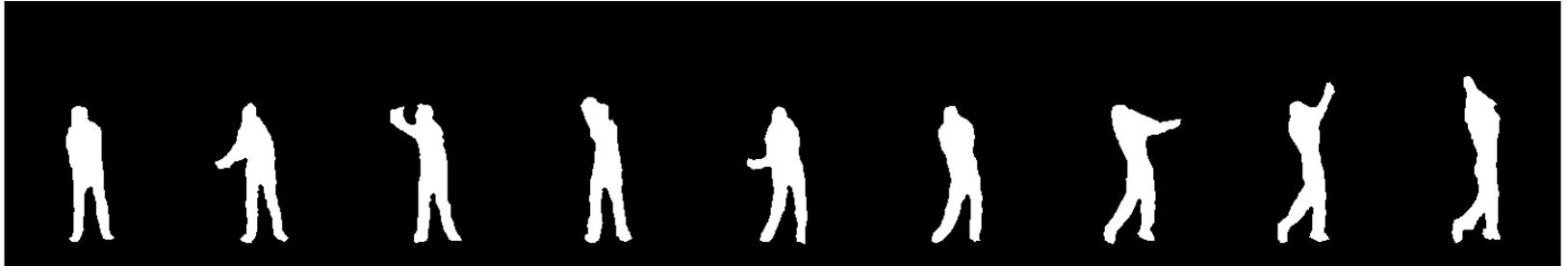
or ?

MULTIVALUED $F$:

Microsoft®

| 62 | 10 | 16 | 36 | 30 | 52 | 49 | 44 | 32 |
|----|----|----|----|----|----|----|----|----|
| 67 | 87 | 24 | 7  | 81 | 5  | 27 | 39 | 85 |
| 41 | 0  | 13 | 22 | 33 | 21 | 67 | 40 | 79 |
| 8  | 70 | 12 | 11 | 10 | 32 | 17 | 28 | 50 |
| 21 | 74 | 78 | 17 | 70 | 74 | 62 | 46 | 56 |
| 82 | 78 | 52 | 22 | 35 | 1  | 17 | 46 | 53 |

Microsoft®

Pose, θ

Image, *z*

Pose, θ

Image, **z**

Pose, $\theta$

Image, $z$

$p(\theta^{new}|z^{new})$

$z^{new}$

t=1

t=1

t=2

t=2

t=1

t=2

R Navaratnam, A Fitzgibbon, R Cipolla
**The Joint Manifold Model for Semi-supervised Multi-valued Regression**
IEEE International Conference on Computer Vision, 2007