# Computer vision systems that work: Movies, Games, Mixed Reality

# ANDREW FITZGIBBON

@awfidius







Overview Publications Events Career Opportunities Microsoft Research blog

"Again, [the Difference Engine] might act on other things beside *number*." — Ada Lovelace

The ADA group studies AI models and techniques that apply to complex real-world data: big data, small data, data with complex structures like trees and graphs. In short: "All Data AI".

Linking language and vision with data from structured domains such as programming languages, bioinformatics, chemistry, or the web, we explore all aspects of machine learning and related disciplines. We research in Bayesian inference, Deep Learning, Optimization and Programming Languages, with multiple prize-winning papers and research. And we have contributed to some of the most advanced technologies that Microsoft builds, with contributions to Xbox, Office, HoloLens, Bing, Visual Studio, and many others. Our superpower is our combined focus on fundamental research and real-world engineering, to push the boundaries of AI.



"Again, [the Difference Engine] might act on other things beside *number*." — Ada Lovelace

The ADA group studies AI models and techniques that apply to complex real-world data: big data, small data, data with complex structures like trees and graphs. In short: "All Data AI".



- School and university in Cork, Ireland
  First job: Water-taxi driver
- Master's and PhD in Edinburgh, Scotland
  Topic: 2D shapes, mainly the ellipse
- Researcher in Engineering, Oxford University
  3D vision—things that stay still
- Researcher and Scientist, Microsoft
  3D vision—things that move















## UNIVERSITY OF OXFORD ROBOTICS RESEARCH GROUP











G Johansson, 1973



G Johansson, 1973











#### 3D modelling from uncalibrated turntable sequences

### CASE STUDY (1998)



#### Automatic 2D point tracking







#### Automatically determined calibration and 3D structue







Automatically determined calibration and 3D structue









#### Automatically determined calibration and 3D structue









#### Automatically recovered model: polygon mesh Subpixel silhouettes for subvoxel marching cubes

















- "Let's make a computer vision startup!"
  - Somewhat unusual in 1999

- And nearly a very bad idea...
  - Industries that need 3D models are quite particular



#### Automatically determined calibration and 3D structure















Autosaved C:\Users\awf\AppData\Local\Temp\dance b5 solved augmented.bdp\_C\_autosave.bpj.bdp

Used : 172.9 Mb 🖳 Total : 1500.0 Mb







## How?

- 1. Gather 100 test videos
- 2. Build metrics
- 3. Implement all known algorithms
- 4. Fix the broken metrics
- 5. Add 100 more videos





*Troy (2004)* Warner Bros, MPC





- A computer program which makes inserting 3D objects easy
- Developed at Oxford University and at company "2d3"
- Now used in almost every movie
  - Lord of the Rings series
  - Harry Potter series
  - District 9

and, of course.... Bridget Jones's Diary



### BOUJOU





# 4 wide-angle tracking cameras






## 2004

# THINGS THAT MOVE





J. Shotton, J. Winn, C. Rother, A. Criminisi, *TextonBoost*: Joint Appearance, Shape and Context Modeling for Multi-Class Object Recognition and Segmentation. European Conference on Computer Vision, 2006



#### Real-Time Semantic Segmentation

#### File View Camera





=

24

## ☑ Wide range of motion

- **B**ut limited agility
- And not realtime



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









"We need a body tracker with

- ☑ All motions...
- ☑ All agilities...
- ✓ 10x Realtime...
- For multiple players...



- "We need a body tracker with
- All motions...
- ☑ All agilities...
- 🗹 10x Realtime...
- For multiple players...

... but you have got 3D 😳"





"... but you have got 3D <sup>(C)</sup>"







## inferred joint positions: no tracking or smoothing









## TRAINING DATA







#### 10 March 2011 Last updated at 11:09 GMT



Microsoft has sold more than 10 million Kinect sensor systems since launch on 4 November, and - according to Guinness World Records - is the fastest-selling consumer electronics device on record.

The sales figures outstrip those of both Apple's iPhone and iPad when launched, Guinness

![](_page_49_Picture_5.jpeg)

🗄 🖻 🧲 🖬 🔒

![](_page_50_Picture_0.jpeg)

Overview

<u>Products</u> Documentation Pricing Training Marketplace  $\vee$  Partners  $\vee$  Support  $\vee$  Blog More  $\vee$ 

Home / Products / Mixed Reality / Kinect DK

Solutions

## **Azure Kinect DK**

Developer kit with advanced AI sensors for building computer vision and speech models

![](_page_50_Picture_6.jpeg)

![](_page_50_Picture_7.jpeg)

#### Health and life sciences

Enhance physical therapy, improve and monitor

![](_page_50_Picture_10.jpeg)

#### Retail

Design a shopping experience that your

![](_page_50_Picture_13.jpeg)

## Logistics and manufacturing

Ship and receive items

![](_page_50_Picture_16.jpeg)

#### Robotics

Gain new environmental understanding with depth,

I was the one working on body tracking – yet my algorithms played no part

Jamie was working on farm animals – turned out to be the key

![](_page_51_Picture_2.jpeg)

![](_page_52_Picture_0.jpeg)

![](_page_53_Picture_0.jpeg)

![](_page_54_Picture_0.jpeg)

## Fitting subdivision surfaces to 2D data

![](_page_54_Picture_2.jpeg)

![](_page_55_Picture_0.jpeg)

## Fitting subdivision surfaces to 2D data

![](_page_55_Picture_2.jpeg)

![](_page_56_Figure_0.jpeg)

## Sidouette formation: 3D to 2D

![](_page_56_Picture_2.jpeg)

![](_page_57_Picture_0.jpeg)

![](_page_58_Picture_0.jpeg)

![](_page_59_Picture_0.jpeg)

![](_page_59_Picture_1.jpeg)

![](_page_60_Picture_0.jpeg)

## HoloLens 2

![](_page_61_Picture_1.jpeg)

## Making hand tracking accurate

Collecting ground truth Learning a better hand model Machine learning

![](_page_62_Picture_2.jpeg)

## 4D Ground Truth

148,000 hand poses

97 hand shapes

![](_page_63_Picture_3.jpeg)

### Synthetic training data

Synthetic training data are labelled images made using computer graphics.

Why use synthetic data?

- Clean labels without annotation noise or error
- Can make GT that is impossible to label by hand
- Easy to control variation in dataset

![](_page_64_Picture_6.jpeg)

## Synthetic Training Data for Hand Tracking

![](_page_65_Picture_1.jpeg)

Visible Light - RGB

Depth Camera

Ground Truth for Machine Learning

Making it fast on a (pair of) 500MHz machines with 128K RAM

- Optimize code
  Great, we went to Floptimal
- Whiteboard malloc
- Use simpler models?
  Great, but lose accuracy.

![](_page_66_Figure_4.jpeg)

Model-fitting iterations per second

### Making it fast on a (pair of) 500MHz machines with 128K RAM

- Optimize code
  Great, we went to Floptimal
- Whiteboard malloc
- Use simpler models?
  Great, but lose accuracy.

![](_page_67_Figure_4.jpeg)

![](_page_68_Picture_0.jpeg)

![](_page_69_Picture_2.jpeg)

**Overview** Publications Career Opportunities

"Again, [the Difference Engine] might act on other things beside number." — Ada Lovelace

The ADA group studies AI models and techniques that apply to complex real-world data: big data, small data, data with complex structures like trees and graphs. In short: "All Data AI".

Linking language and vision with data from structured domains such as programming languages, bioinformatics, chemistry, or the web, we explore all aspects of machine learning and related disciplines. We research in Bayesian inference, Deep Learning, Optimization and Programming Languages, with multiple prize-winning papers and research. And we have contributed to some of the most advanced technologies that Microsoft builds, with contributions to Xbox, Office, HoloLens, Bing, Visual Studio, and many others. Our superpower is our combined focus on fundamental research and real-world engineering, to push the boundaries of Al.

0.0	0 705	0.891	0.891	01	21		0 705	0.891	0.891	012
0 705	0.891	0.891	0 1 2 1	11	34	0.705	0.891	0 891	0 12	1 11
0.705	0.651	0.051	0.224	+1	58	0.703	0.051	0.051	0.12	4 41 5
0.564	0.673	0.107	0.334		18	0.564	0.673	0.107	0.334	+ ./ L
0.050	0.401	0.293	0.858	95	13	0.050	0.401	0.293	0.858	\$ 95 1
0.162	0.477	0.847	0.918	59	11	0.162	0.477	0.847	0.918	3 59 <u>t</u>
0.083	0.988	0.494	0.713	19	10	0.083	0.988	0.494	0.713	3 19 L
0.893	0.750	0.520	0.041	38	82	0.893	0.750	0.520	0.04	1 38 3
0.666	0.302	0.726	0.610	13	90	0.666	0.302	0.726	0.610	130
0.452	0.927	0.034	0.082	58	_	0.452	0.927	0.034	0.082	2 58
0.650	0.821	0.822	0.790			0.650	0.821	0.822	0.790	0
0.7	0.705	0.891 91 0.89 0.011	0.891 91 0.1 0.317	21	21 34	0.9	51 0.4 0.891	0.891	0.121	17 34 11 58
0.467	0.423	0.985	0.641	50	20	0.564	0.673	0.107	0.334	17 18
0.071	0.606	0.497	0.217	10	18	0.050	0.401	0.293	0.858	95 13
0.021	0.000	0.602	0.205	10	13	0.162	0.477	0.847	0.918	59 41
0.909	0.014	0.003	0.295	13	41	0.083	0.988	0.494	0.713	19 10
0.996	0.776	0.439	0.869	41	10	0.893	0.750	0.520	0.041	18 87
			111/10	10	82					14
0.717	0.726	0.897	0.149		10	0.666	0.302	0.726	0.610	13 bo
0.909	0.726	0.543	0.038	82	90	0.666	0.302	0.726	0.610	13 90
0.909	0.726 0.246 0.220	0.897 0.543 0.807	0.038	82 90	90	0.666	0.302	0.726	0.610	13 90 58

![](_page_70_Figure_1.jpeg)

Gen 2: Everything's a Tensor And frameworks constrain us to think this Gen3. Everything has its natural structure But we need frameworks to work with these structures

![](_page_70_Picture_4.jpeg)

way

![](_page_71_Picture_0.jpeg)

Three (well, four) computer vision systems that work

- 1. Pure geometry
- 2. Machine learning with huge training data
- 3. Geometry + deep learning, again with big data

#### The future:

- More hybrids
- More data
- More models
- More hardware