Research

DEPTH, HUMAN POSE, AND CAMERA POSE

JAMIE SHOTTON











Kinect Adventures

- Depth sensing camera
- Tracks 20 body joints in real time
- Recognises your face and voice





Trial of "touchless" gaming technology in surgery

By Adam Brimelow Health Correspondent, BBC News

Doctors in London are trialling "touchless" technology, often used in TV games, to help them carry out delicate keyhole surgery.

The system allows them to manipulate images with their voice and hand-gestures rather than using a keyboard and mouse.

Surgeons say it gives them more control and avoids disruption.



Home US World Politics Business Sports Entertainment Health Digital Home ... msnbc.com

Child with autism connects with Kinect

When Kyle's father got Xbox's motion control system, he had no idea it would be a breakthrough for his boy

ect' surprisingly fun Bolow: 📴 Video 🤤 Discuss 🔤 Related By Wilson Rothman manbc.com Print | Font: A + updated 11/11/2010 7:18:30 PM ET

> John Yan reviews games for a site called Gaming Nexus, so despite his initial lack of enthusiasm in the Xbox 360 Kinect motion controller, he knew he'd have to buy one when they came out. After all, it wouldn't be fair to dump all the Kinect reviews on



Advances in Computer Vision and Pattern Recognition

Andrea Fossati Juergen Gall Helmut Grabner **Xiaofeng Ren** Kurt Konolige Editors



Consumer **Depth Cameras** for Computer Vision

Research Topics and Applications

D Springer

What the Kinect Sees



Structured light





Depth Makes Vision That Little Bit Easier



RGB

⊠ Only works well lit

⊠ Background clutter

Depth

✓ Works in low light

Background removal easier

Scale unknown

Calibrated depth readings

☑ Color and texture variation

Uniform texture





KINECTFusion



Joint work with Shahram Izadi, Richard Newcombe, David Kim, Otmar Hilliges, David Molyneaux, Pushmeet Kohli, Steve Hodges, Andrew Davison, Andrew Fitzgibbon. SIGGRAPH, UIST and ISMAR 2011.

KINECTFusion



ROADMAP



[CVPR 2012]



Scene Coordinate Regression [CVPR 2013]



THE VITRUVIAN MANIFOLD



Jonathan Taylor



Jamie Shotton



Toby Sharp



Andrew Fitzgibbon

CVPR 2012

Human Pose Estimation

Given some image input, recover the 3D human pose:



Joint positions and angles

In this work:

- Single frame at a time (no tracking)
- Kinect depth image as input (background removed)

Why is Pose Estimation Hard?





A Few Approaches

Regress directly to pose? e.g. [Gavrila '00] [Agarwal & Triggs '04]

Detect and assemble parts?

e.g. [Felzenszwalb & Huttenlocher '00] [Ramanan & Forsyth '03] [Sigal et al. '04]

Detect parts?

e.g. [Bourdev & Malik '09] [Plagemann et al. '10] [Kalogerakis et al. '10]



Per-Pixel Body Part Classification [Shotton *et al.* '11]



Background: Learning Body Parts for Kinect



____/

Synthetic Training Data



Train invariance to:



Depth Image Features

- Depth comparisons
 - very fast to compute

feature response $f(\mathbf{x}; \mathbf{v}) = d(\mathbf{x}) - d(\mathbf{x} + \Delta)$ image coordinate



scales inversely with depth

Background pixels d =large constant



Decision tree classification





Take (\mathbf{v}, θ) that maximises information gain:

$$\Delta E = -\frac{|S_{l}|}{|S_{n}|}E(S_{l}) - \frac{|S_{r}|}{|S_{n}|}E(S_{r})$$

Goal: drive entropy at leaf nodes to zero



A. Criminisi J. Shotton *Editors*

Decision Forests for Computer Vision and Medical Image Analysis

🖄 Springer

Decision Forests Book

- Theory Tutorial & Reference
- Practice Invited Chapters
- Software and Exercises
- Tricks of the Trade



no tracking or smoothing





inferred joint position hypotheses

no tracking or smoothing

Body Part Recognition in Kinect





Skeleton does not explain the depth data Limited ability to cope with hard poses



Human Skeleton Model

- Mesh is attached to a hierarchical skeleton
- Each limb l has a transformation matrix $T_l(\theta)$ relating its local coordinate system to the world:

 $T_{\text{root}}(\theta) = R_{\text{global}}(\theta)$ $T_{l}(\theta) = T_{\text{parent}(l)}(\theta)R_{l}(\theta)$

- $R_{l_{arm}}(\theta)$ $R_{\text{global}}(\theta)$
- $R_{\text{global}}(\theta)$ encodes a global scaling, translation and rotation
- $R_l(\theta)$ encodes a rotation and fixed translation relative to its parent
- 13 parameterized joints 🔘 using quaternions to represent unconstrained rotations
- This gives θ a total of 1 + 3 + 4 + 4 * 13 = 60 degrees of freedom

Linear Blend Skinning

Mesh in base pose θ_0

Each vertex *u*

- has position p in base pose θ_0
- is attached to *K* limbs $\{l_k\}_{k=1}^K$ with weights $\{\alpha_k\}_{k=1}^K$

In a new pose θ , the skinned position u of is:

$$M(u; \theta) = \sum_{k=1}^{K} \alpha_k T_{l_k}(\theta) T_{l_k}^{-1}(\theta_0) p$$
position in limb *l_k*'s coordinate system
position in world coordinate system

Test Time Model Fitting

• Assume each observation x_i is generated by a point on our model u_i



Note: simplified energy - more details to come

Optimizing $\min_{\theta} \min_{u_1...u_n} \sum_i d(x_i, M(u_i; \theta))$

- Alternating between pose θ and correspondences u₁, ... u_n
 ➢ Articulated Iterative Closest Point (ICP)
- Traditionally, start from initial θ
 - manual initialization
 - track from previous frame
- Could we instead infer initial correspondences u_1 , ... u_n discriminatively?
- And, do we even need to iterate?

One-Shot Pose Estimation: An Early Result

Can we achieve a good result without iterating between pose θ and correspondences $u_1, \dots u_n$?



test depth image ground truth correspondences

convergence visualization

From Body Parts to Dense Correspondences



Texture is mapped across body shapes and poses

The "Vitruvian Manifold" Embedding in 3D



Overview



Discriminative Model: Predicting Correspondences



Learning the Correspondences

• How to learn the mapping from depth pixels to correspondences?



• Train regression forest






Full Energy



 $\rho(e)$

 $e \stackrel{!}{=} 0$

 $\int c_t(\theta_0)$

 $c_s(\theta_0)$

- Term E_{vis} approximates hidden surface removal and uses robust error
- Gaussian prior term E_{prior}
- Self-intersection prior term *E*_{int} approximates interior volume

Energy is robust to noisy correspondences

- Correspondences far from their image points are "ignored"
- Correspondences facing away from the camera are "ignored"
 - avoids model getting stuck in front of the image pixels





"Easy" Metric: Average Joint Accuracy



Results on 5000 synthetic images

Hard Metric: "Perfect" Frame Accuracy













0.45m

D:

0.17m

Comparison



Sequence Result



Each frame fit independently: no temporal information used





Generalization to Multiple 3D/2D Views

- Easily extended to Q views where each view has
 - $-n_q$ correspondences per view
 - viewing matrix P_q to register the scene
- Can also extend to 2D silhouette views
 - let data points x_{ik} be 2D image coordinates
 - let P_q include a projection to 2D
 - minimize re-projection error

$$\min_{\theta} \sum_{q=1}^{Q} \sum_{i}^{n_{q}} d(x_{iq}, P_{q}M(u_{iq}; \theta))$$

Silhouette Experiment



Vitruvian Manifold: Summary

- Predict per-pixel image-to-model correspondences
 - train invariance to body shape, size, and pose

- "One-shot" pose estimation
 - fast, accurate
 - auto-initializes using correspondences



Scene Coordinate Regression Forests For Camera Relocalization In RGB-D Images

JAMIE SHOTTON BEN GLOCKER CHRISTOPHER ZACH SHAHRAM IZADI ANTONIO CRIMINISI ANDREW FITZGIBBON [CVPR 2013]



Know this



A world scene

Observe this



Where is this?



6D camera pose, *H* (camera to scene transformation)

Single RGB-D frame

APPLICATIONS

Lost or kidnapped robots



Improving KinectFusion

Augmented reality

TYPICAL APPROACHES TO CAMERA LOCALIZATION

- Tracking alignment relative to previous frame
- Key point detection → local descriptors → matching → geometric verification
 e.g. [Holzer et al. '12], [Winder & Brown '07], [Lepetit & Fua '06], [Irschara et al. '09]



- Whole key-frame matching e.g. [Klein & Murray 2008] [Gee & Mayol-Cuevas 2012]
- Epitomic location recognition

approximate

precise

e.g. [Besl & MacKay '92]

[Ni et al. 2009]

PROBLEMS IN REALWORLD CAMERA LOCALIZATION

- The real world is less exciting than vision researchers might like
 - > sparse interest points can fail





The real world is big



KEY IDEA: SCENE COORDINATE REGRESSION



Scene coordinate XYZ ⇔ RGB color space

KEY IDEA: SCENE COORDINATE REGRESSION

 Let each pixel predict direct correspondence to 3D point in scene coordinates:



Input RGB Input Depth







Desired Correspondences



3D model from KinectFusion (only used for visualization)

SCENE COORDINATE REGRESSION

- Offline approach to relocalization
 - observe a scene
 - train a regression forest
 - revisit the scene

p $\mathcal{M}_{l_1(p)}$

- Aim for really precise localization
 - e.g. suitable for AR overlays
 - from a single frame
 - without an explicit 3D model



[Bunny: Stanford]

SCENE COORDINATE REGRESSION (SCORE) FORESTS





Depth & RGB features



$$f_{\phi}^{\text{depth}}(\mathbf{p}) = D\left(\mathbf{p} + \frac{\boldsymbol{\delta}_1}{D(\mathbf{p})}\right) - D\left(\mathbf{p} + \frac{\boldsymbol{\delta}_2}{D(\mathbf{p})}\right)$$
$$f_{\phi}^{\text{da-rgb}}(\mathbf{p}) = I\left(\mathbf{p} + \frac{\boldsymbol{\delta}_1}{D(\mathbf{p})}, c_1\right) - I\left(\mathbf{p} + \frac{\boldsymbol{\delta}_2}{D(\mathbf{p})}, c_2\right)$$

Leaf Predictions

Forest Predictions

$$\mathcal{M}_l \subset \mathbb{R}^3$$
$$\mathcal{M}(\mathbf{p}) = \bigcup_t \mathcal{M}_{l_t(\mathbf{p})}$$

TRAINING A SCORE FOREST

Training Data

- RGB-D frames with known camera poses H
- Generate 3D pixel labels automatically:

 $\mathbf{m} = H\mathbf{x}$



Learning (standard)

- Greedily train tree
- Reduction in spatial variance objective:

$$Q(\mathcal{S}_n, \boldsymbol{\theta}) = V(\mathcal{S}_n) - \sum_{d \in \{L, R\}} \frac{|\mathcal{S}_n^d(\boldsymbol{\theta})|}{|\mathcal{S}_n|} V(\mathcal{S}_n^d(\boldsymbol{\theta}))$$

with $V(\mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{(\mathbf{p}, \mathbf{m}) \in \mathcal{S}} \|\mathbf{m} - \bar{\mathbf{m}}\|_2^2$

- Regression, not classification
- Mean shift to summarize distribution at leaf l into small set $\mathcal{M}_l \subset \mathbb{R}^3$

ROBUST CAMERA POSE OPTIMIZATION

Energy Function



Optimization

Preemptive RANSAC

```
[Nistér ICCV 2003]
```

- With pose refinement
 [Chum et al. DAGM 2003]
 - efficient updates to means & covariances used by Kabsch SVD
- Only a small subset of pixels used

INLYING FOREST PREDICTIONS



PREEMPTIVE RANSAC OPTIMIZATION



THE 7SCENES DATASET

Spatial	# Frames	
Extent	Train	Test
$3m^3$	4k	2k
$4m^3$	2k	2k
$2m^3$	1k	1k
$5.5m^{3}$	6k	4k
$6m^3$	4k	2k
$6m^3$	7k	5k
$5m^3$	2k	1k
	Spatial Extent 3m ³ 4m ³ 2m ³ 5.5m ³ 6m ³ 6m ³ 5m ³	$\begin{array}{c c c c c c c c c c c c c c c c c c c $

Dataset available from authors



BASELINES FOR COMPARISON

Sparse Key-Points (RGB only)

- ORB matching [Rublee et al. ICCV 2011]
 - FAST detector
 - Rotation aware BRIEF descriptor
 - Hashing for matching
- Geometric verification
 - RANSAC & perspective 3 point
 - Final refinement given inliers

Tiny-Image Key-Frames (RGB & Depth)

- Downsample to 40x30 pixels
- Blur
- Normalized Euclidean distance
- Brute-force search
- Interpolation of 100 closest poses

[Klein & Murray ECCV 2008] [Gee & Mayol-Cuevas BMVC 2012]

QUANTITATIVE COMPARISON

Metric:

Proportion of test frames with < 0.05m translational error and < 5° angular error

Results:	Baselines		Our Results		
Scene	Tiny-image RGB-D	Sparse RGB	Depth	DA-RGB	DA-RGB + D
Chess	0.0%	70.7%	82.7%	92.6%	91.5%
Fire	0.5%	49.9%	44.7%	82.9%	74.7%
Heads	0.0%	67.6%	27.0%	49.4%	46.8%
Office	0.0%	36.6%	65.5%	74.9%	79.1%
Pumpkin	0.0%	21.3%	58.6%	73.7%	72.7%
RedKitchen	0.0%	29.8%	61.3%	71.8%	72.9%
Stairs	0.0%	9.2%	12.2%	27.8%	24.4%

Choice of different image features

QUALITATIVE COMPARISON



ground truth

DA-RGB SCoRe forest

sparse baseline

closest training pose

QUALITATIVE COMPARISON



ground truth

DA-RGB SCoRe forest

sparse baseline

closest training pose

TRACK VISUALIZATION VIDEOS



ground truth

DA-RGB SCoRe forest

RGB sparse baseline

single frame at a time – no tracking

ARVISUALIZATION



[Bunny: Stanford]

single frame at a time – no tracking

SIMPLE ROBUST TRACKING

Add a single extra hypothesis to optimization: the result from previous frame

Single frame

	Our Results			Frame-to-Frame
Scene	Depth	DA-RGB	DA-RGB + D	Tracking
Chess	82.7%	92.6%	91.5%	95.5%
Fire	44.7%	82.9%	74.7%	86.2%
Heads	27.0%	49.4%	46.8%	50.7%
Office	65.5%	74.9%	79.1 %	86.8%
Pumpkin	58.6%	73.7%	72.7%	76.1%
RedKitchen	61.3%	71.8%	72.9%	82.4%
Stairs	12.2%	27.8%	24.4%	39.2%

AR VISUALIZATION WITH TRACKING



[Bunny: Stanford]

simple robust frame-to-frame tracking enabled

MODEL-BASED REFINEMENT

Model-based refinement

[Besl & McKay PAMI 1992]

- requires 3D model of scene
- run rigid ICP from our inferred pose between observed image and model


AR VISUALIZATION WITH TRACKING AND REFINEMENT



[Bunny: Stanford]

simple robust frame-to-frame tracking and ICP-based model refinement enabled

Fire Scene

SCoRe Forest (single frame at a time)

RGB inputdepth inputrend+ AR overlay+ AR overlayfrom

rendering of model from inferred pose



SCoRe Forest + simple robust frame-to-frame tracking

SCoRe Forest + simple robust frame-to-frame tracking + ICP refinement to 3D model





Pumpkin Scene

SCoRe Forest (single frame at a time)







SCoRe Forest simple robust frame-to-frame tracking ICP refinement to 3D model





SCENE RECOGNITION





SCENE COORDINATE REGRESSION - SUMMARY

Scene coordinate regression forests

- provide a single-step alternative to detection/description/matching pipeline
- can be applied at any valid pixel, not just at interest points
- allow accurate relocalization without explicit 3D model

Tracking-by-detection is approaching temporal tracking accuracy

WRAP UP

- Depth cameras are having huge impact
- Decision forests + big data





Unifying principal:

Per-pixel regression and per-image model fitting

Thank you!



With thanks to:

Andrew Fitzgibbon, Jon Taylor, Ross Girshick, Mat Cook, Andrew Blake, Toby Sharp, Pushmeet Kohli, Ollie Williams, Sebastian Nowozin, Antonio Criminisi, Mihai Budiu, Duncan Robertson, John Winn, Shahram Izadi

obertson, John Winn, Shahram Izadi

The whole Kinect team, especially: Alex Kipman, Mark Finocchio, Ryan Geiss, Richard Moore, Robert Craig, Momin Al-Ghosien, Matt Bronder, Craig Peeper



