Media IC & System Lab
Graduate Institute of Electronics Engineering
National Taiwan University

# Accurate 6DoF Object Pose Estimation and Tracking

A Dissertation Defense

by
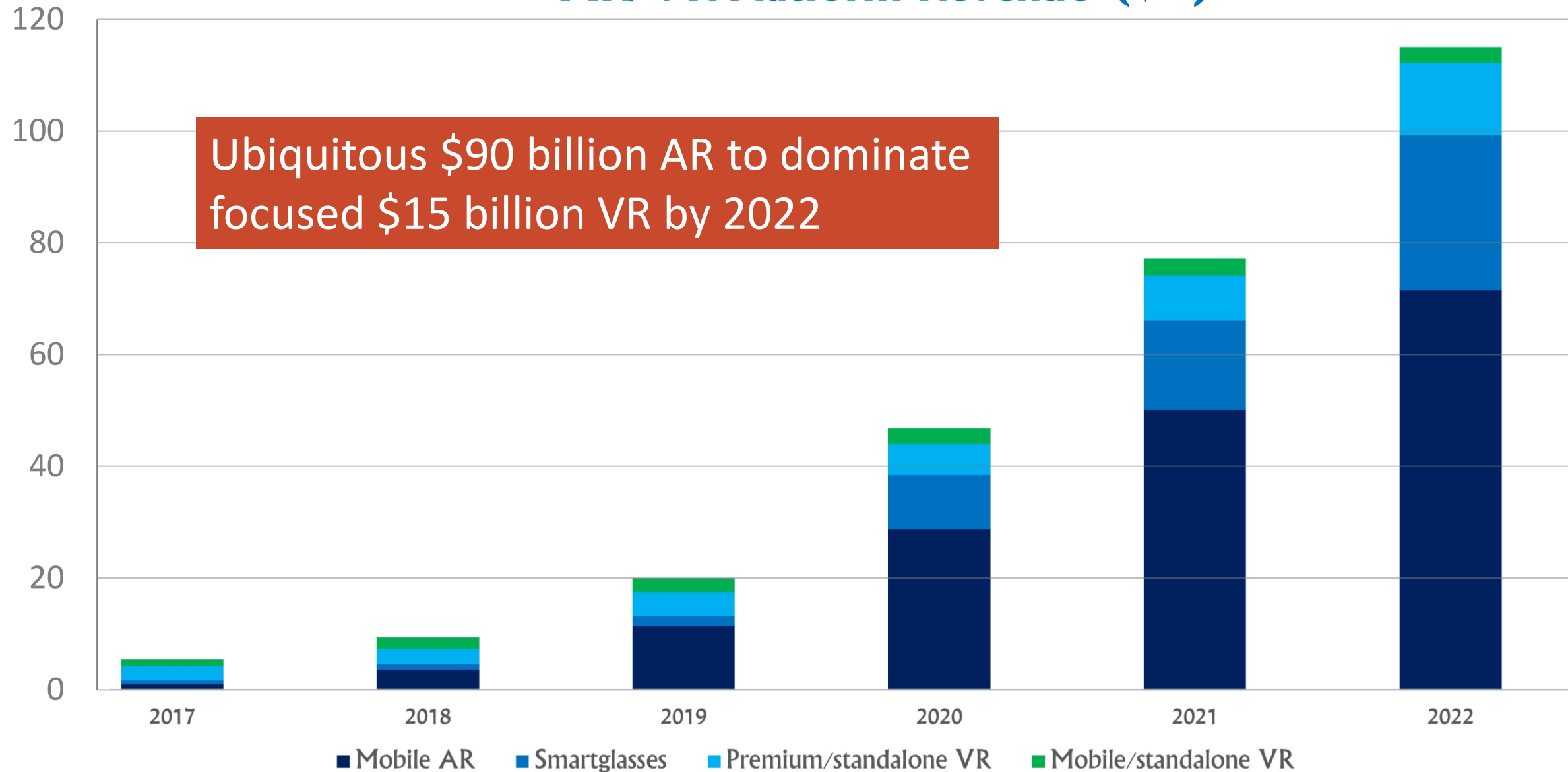
Po-Chen Wu

Advisor: Dr. Shao-Yi Chien

Co-Advisor: Dr. Ming-Hsuan Yang

# AR/VR Platform Revenue ($B)

Ubiquitous $90 billion AR to dominate focused $15 billion VR by 2022



Legend: ■ Mobile AR  ■ Smartglasses  ■ Premium/standalone VR  ■ Mobile/standalone VR

**Digi-Capital**™

2

# Augmented Reality



https://giphy.com/gifs/adweek-place-ar-4R63eQx8wyEda

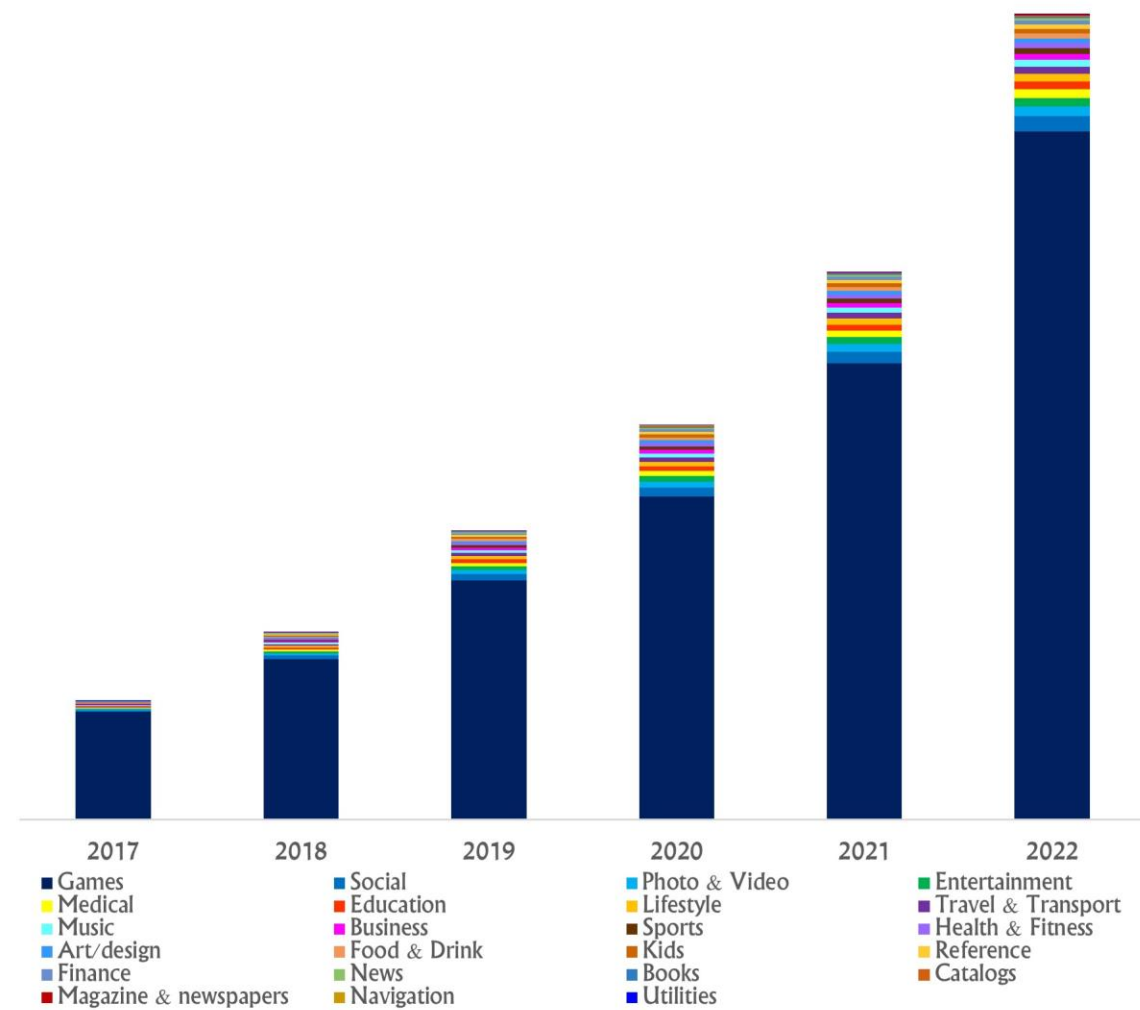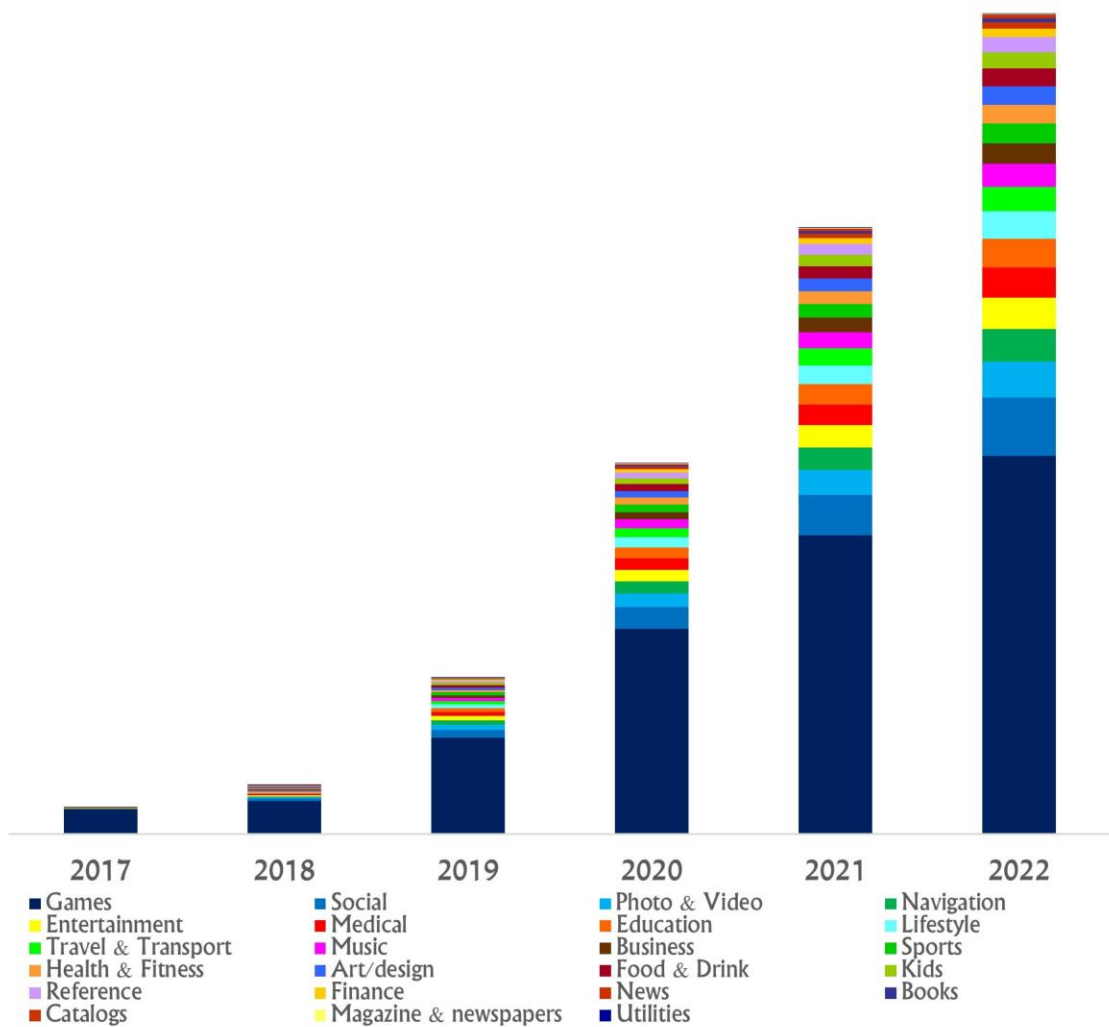# Virtual Reality



Tap me!

https://cdn.makeuseof.com/wp-content/uploads/2017/01/oculus-touch-gestures-gif.gif

# AR/VR App Store Category Revenue (IAP/Premium)

*(Note: scales on both charts are not the same)*

## AR Category (IAP/Premium) Revenue ($90B)        VR Category (IAP/Premium) Revenue ($15B)



**AR chart legend:**
- Games
- Entertainment
- Travel & Transport
- Health & Fitness
- Reference
- Catalogs
- Social
- Medical
- Music
- Art/design
- Finance
- Magazine & newspapers
- Photo & Video
- Education
- Business
- Food & Drink
- News
- Utilities
- Navigation
- Lifestyle
- Sports
- Kids
- Books

**VR chart legend:**
- Games
- Medical
- Music
- Art/design
- Finance
- Magazine & newspapers
- Social
- Education
- Business
- Food & Drink
- News
- Navigation
- Photo & Video
- Lifestyle
- Sports
- Kids
- Books
- Utilities
- Entertainment
- Travel & Transport
- Health & Fitness
- Reference
- Catalogs

Years: 2017, 2018, 2019, 2020, 2021, 2022

Digi-Capital™

4

# Google



https://inhabitat.com/ecouterre/your-surgeon-could-be-using-google-glass-in-the-operating-room/

https://www.theverge.com/2016/11/1/13480840/google-tango-lenovo-phab-2-pro-apps-games-release-date

# Facebook

https://www.independent.co.uk/life-style/gadgets-and-tech/news/oculus-rift-price-headset-and-computers-that-can-run-it-begin-at-1499-a6863676.html

http://immediatefuture.co.uk/blog/facebook-launches-ar-studio/

# Apple



https://mashable.com/2017/06/05/apple-arkit-hands-on/



 GLASS

Coming 2018

https://www.tomsguide.com/us/apple-ar-glasses-tim-cook,news-25964.html

# Microsoft



https://www.geeky-gadgets.com/microsoft-windows-10-vr-headset-22-11-2016/

http://fortune.com/2017/02/21/microsoft-hololens-update-delay/

8

# Magic Leap



https://www.engadget.com/2017/12/20/magic-leap-one-details-questions-dont-know/



https://www.engadget.com/2018/02/13/magic-leaps-ar-headsets-price-nba-deal/

GANZIN

https://www.youtube.com/watch?v=d9_kYHbEI5w

EYE TRACKING IC          SCENE CAMERA

EYE CAMERA                                    EYE CAMERA

https://www.youtube.com/watch?v=d9_kYHbEI5w

(02) 3366-3668

# Digi-Capital™ AR/VR Leaders*

*Includes funded/exited startups and selected corporates*

## Advertising/marketing

adtile, AdsReality, ADVR, SUBMISIVE, LIVE LIKE, Omnivirt, IMMERSAL, REALVISION, ZeniAd, retinad, AdvVr, blippAR, PLATTAR, AUGMENT, PRSONAS, advrty, immersv, vertebrae, VIRE, Xperiel, zappar

## Art/design

Mindesk, iris, artomatix

## Business

cluster., VIRTUALITICS

## Distribution

ConstructVR, Wevr, WEARVR, Sketchfab, UNIMERSIV, ROCK VR

## eCommerce

Imagine, CAPPASITY, VIVIDWORKS, InContext, spaceview, hexa, tylko

## Education

ARA, nearpod, cerevrum, PUBLIC3, serious:labs, mativision, IMMERSIVE EDUCATION, sdk, LVR, MEL Science, CURISCOPE, Jig SPACE, Learningbone, VISCOPIC, OpenSesame, LABSTER, STRIVR, VPS, DISCOVR LABS, BOULEVARD, zSpace, 开大

## Enterprise/B2B

UPSKILL, Mure, HYPER FAIR, SCOPE AR, taa tile, IMMERSION, soluis, FIELDBIT, nsiteVR, SPACESYS

## Entertainment

Madefire, TRIPP, SKYLIGHTS

## Games

Appliture, boom, bublar, Directive, ESCALATION, Fast Travel Games, Freeform Labs, MUNZEE, MULTIVERSE, MIDAS TOUCH INTERACTIVE, RESOLUTION, animosa BRANDS, CCP, fuzzycube, HARMONIX, metaverse, KENZAN, KINGDOMS CASTLES, FINAL LEVEL ZERO, VELAH, CYBERPONG VR, sólfar, SURVIOS, SVRVIVE SUPERSTAR, POLYARC, REALITEER, NIANTIC, nDreams, Owlchemy Labs, playful, PLAY FUSION, SHORTFUSE, TEMPLE GATES, umbra, STARBREEZE, VRANI, PROXY42, KITE & LIGHTNING, ROBLOX, VITO, inXile, WARDUCKS, LUDEN

## Kids

QuiverVision

## Lifestyle

Around Media, Habiteo, GeoCV, Vieweet, iStaging, Pixie, rooomy, urbanbase, MODSY

## Location based

Igloo, ENTER SPACE, HAVSON, VRSTUDIOS, IMAX, NOMADIC VR, DREAMSCAPE, VREC, hyperverse, 恐龙VR, LEKE VR, ZERO LATENCY, TWO BIT CIRCUS, FORCE FIELD, SPACES

## Medical

aira, appliedVR, AUGMEDIX, BIOLUCID, Luminate, echopixel, MEDICAL REALITIES, SURGICAL THEATER, KBVISION, psious, Fearless, VRPHOBIA, marion surgical, orosim, VIVID, mindmaze, oriense

## Music

TheWaveVR

## Navigation

Aero Glass, NAVISENS, mapbox, PARACOSM, WAYRAY

## News

87870.com, HAMMER & TUSK, UPLOAD

## Peripherals

REMORIA, 3EGway, Ximmerse, RAZER, VIRZOOM, VAQSO, UMBO, OSSIC, Myo, LIGHTFORM, SUBPAC, POLARIANT, TACTICAL HAPTICS, YOST, ultrahaptics, adabra, nod, here., CyberGlove, LEAP, BHYVE, Omnifinity, bhaptics, CYBERITH, KAT VR, NEURODIGITAL

## Photo/video

FOV, FUSION, scandy, 360fly, ALCHEMY VR, ARCHIACT, BLEND, briskeye, bubl, brabo, BITMOVIN, CONSTRUCT STUDIO, CONSTRUCT STUDIO, dRi, EEVO, FELIX & PAUL, LYTRO, entrypoint., GoPro, GIROPTIC, Infinite Mixed Reality, InstaVR, eko, JAUNT, KITE & LIGHTNING, kogeto, Kolor, livelke, LUCID, LUMIERE VR, matterport, LUNA, METTLE, nctech, NEXTVR, NGCODEC, OZO, otoy, Panocam, PARABLE, PENROSE STUDIOS, REALITIES.IO, RYOT, SPACEVR, SpinLe, Syreal, VRC, THE FUTURE, thinglink, TIpit, ORAH, VISBIT, VISUALPATHY, VOKE, Playhouse, zLense, EMERGENT, LITTLJSTAR

## Productivity

(logo)

## Smartglasses

ATHEER, PEPPER, DAQRI, EPSON, GDI, WaveOptics, KOPIN, magic leap, AVEGANT, CMOAR, Microsoft, REAL VIEW, RIDE ON, U, IMMY, VUZIX, ODG, SULON, Meta, SONY, Xikaku, mira, intel, recon

## Solutions/services

ATHASS, DIGITAL, LOOM.AI, pixelbug, spotscale, VNO, DoubleMe, DVERSE, FISHBOWL VR, HOLOGROUP, DIAKRIT, 8th WALL, BinaryVR, ARHT MEDIA, catchoom, BYOND, cognitive, Dacuda, DataMesh, Apmetrix, GAUDIO, iris, MELODY VR, OVA, snobal, smartvizx, VECTARY, mindshow, wrnch, DSCOPE, SIMPLYGON, Moback

## Social

FlirtAR, PLUTO VR, NEOS, SPACEOUT VR, CLEVR, Decentraland, SURREAL, VRCHAT, TIMEFIRE VR, HIGH FIDELITY

## Sports

HOLODIA, FIRSTVISION, REDD, eon sports, SoccerDream, BEYOND SPORTS

## Tech

VicoVR, 3DLOOK, actronika, AMD, amazon, ARVAD, crunchfish, CRYTEK, CUBICMOTION, DIGILENS, Dirac, DOLBY, DYSONICS, EONITE, EPIC GAMES, Gaitup, gestigon, G, DC INFINITY, MANOMOTION, Mint Mate, MOGEES, Movidius, nGRAIN, nitero, PERCEPTION NEURON, NVIDIA, occipital, ptc, REACH, SMI, SoftKinetic, VATJO, THEEYETRIBE, ThirdEye, unity, UNIVRSES, worldviz, VisiSonics, xPerception, SENNHEISER, IMPROBABLE, insightness, antilatency, THRIVE AUDIO, BrokenColors

## Travel/transport

ZeroLight, STURFEE, timelooper, GAMAR, CAR360

## Utilities

BAGEL LABS, bigscreen, Seene, visual camp, JANUSVR, nurulize, wakingapp, HASHPLAY, Vizor

## VR headset

NORAH, alcatel, altergaze, ARCHOS, ZEISS, 暴风魔镜, VUZIX, facebook, FOVE, G, XINGEAR, LG, airVR, MERGE, Microsoft, POWIS, RAZER, SONY, HTC, STARBREEZE STUDIOS, Dee Poon, wearality, evomade, VALVE, SAMSUNG, sensics, EYEDAK

Pose Estimation + Rendering

https://giphy.com/gifs/3d-LGP31CCUCylbO

Disney Research

# Six Degrees of Freedom (6DoF) (Rigid Body) Object Pose



Position (3DoF)

+

Orientation (3DoF)

# Problem Definition

# AR



https://www.microsoft.com/en-us/hololens

# VR



https://giphy.com/gifs/3d-simulator-bird-bl09MDvT4JIWc

# 6DoF Object Pose

- Model-Based Pose Estimation
- Model-Based Pose Tracking

# 6DoF Camera Pose

- Image-Based Camera Localization
- Simultaneous Localization and Mapping (SLAM)





http://www.societyofrobots.com/robotforum/index.php?topic=16714.0

# Camera Projection Model



$$\mathbf{p} \equiv (\mathbf{R}, \mathbf{t})$$

Camera Coordinate System

$\widehat{\mathbf{u}}_i : (\hat{u}_i, \hat{v}_i)$

$\mathbf{u}_i : (u_i, v_i)$

$\mathbf{x}_i : (x_i, y_i, z_i)$

Object Coordinate System

Image Plane

$$\begin{bmatrix} \hat{u}_i \\ \hat{v}_i \\ 1 \end{bmatrix} \sim \begin{bmatrix} hu_i \\ hv_i \\ h \end{bmatrix} = \mathbf{K}[\mathbf{R}|\mathbf{t}] \begin{bmatrix} x_i \\ y_i \\ z_i \\ 1 \end{bmatrix},$$

$$\mathbf{K} = \begin{bmatrix} f_x & 0 & x_0 \\ 0 & f_y & y_0 \\ 0 & 0 & 1 \end{bmatrix}, \mathbf{R} = \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix}, \mathbf{t} = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$$

Intrinsic Matrix          Rotation Matrix          Translation Vector

17

# Objective Functions



- Reprojection error

$$E_r(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\left((\hat{u}_i - u_i)^2 + (\hat{v}_i - v_i)^2\right)$$

$\hat{\mathbf{u}}_i: (\hat{u}_i, \hat{v}_i)$ is the observed point

- Appearance distance

$$E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}|I_c(\mathbf{u}_i) - O_t(\mathbf{x}_i)|$$

or

$$E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i) - O_t(\mathbf{x}_i)\right)^2$$



https://www.zuehlke.com/blog/apple-arkit-augmented-reality-erhaelt-schub/

18

# Feature-Based Approaches [1,2,3]

- Feature detection and matching
  - SIFT, SUFT, FAST, ORB, …

- Outlier removal
  - RANSAC, PROSAC, …

- Perspective-n-point (PnP) Algorithm
  - EPnP, OPnP, UPnP, …
  - $E_r(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\|\widehat{\mathbf{u}}_i - \mathbf{u}_i\|^2$

[1] Lepetit et al., "Point matching as a classification problem for fast and robust object pose estimation," CVPR, 2004.
[2] Collet et al., "Object Recognition and Full Pose Registration from a Single Image for Robotic Manipulation," ICRA, 2009.
[3] Tang et al., "A Textured Object Recognition Pipeline for Color and Depth Image Data," ICRA, 2012.

# Direct Approaches [1,2]

- Finding the best fit from numerous pre-determined candidates
  - $E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}|I_c(\mathbf{u}_i) - O_t(\mathbf{x}_i)|$
  - $E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i) - O_t(\mathbf{x}_i)\right)^2$

[1] Hinterstoisser et al., "Model Based Training, Detection and Pose Estimation of Texture-Less 3D Objects in Heavily Cluttered Scenes," ACCV, 2012.
[2] Tseng et al., "Direct 3D Pose Estimation of a Planar Target," WACV, 2016.

# Learning-Based Approaches

- ## Random forest [1]



- ## Deep neural network [2,3]

[1] Tejani et al., "Latent-Class Hough Forests for 3D Object Detection and Pose Estimation," ECCV, 2014.
[2] Kehl et al., "SSD-6D: Making RGB-Based 3D Detection and 6D Pose Estimation Great Again," CVPR, 2017.
[3] Rad et al., "BB8: A Scalable, Accurate, Robust to Partial Occlusion Method for Predicting the 3D Poses of Challenging Objects without Using Depth," ICCV, 2017.

# Objective

Accurate     Real-time     Low-cost

# Contributions

- ## OPT Dataset
  - We build a large-scale object pose tracking benchmark dataset consisting of RGB-D video sequences of 2D and 3D targets.

- ## DPE Algorithm
  - We propose a robust direct approach of 6DoF pose estimation for planar objects.

- ## DodecaPen
  - We develop a solution for real-time 6DoF tracking that achieves submillimeter accuracy.

Input Data

Approximate Pose Estimation

Pose Refinement

Object Pose Tracking

# OPT

Dataset

# 2D Objects

# 3D Objects

| Simple Geometry | Normal Geometry | Complex Geometry |
|---|---|---|
|  |  |  |

http://www.arzapstudio.com/wp-content/uploads/2016/12/Kinect-Banner-800x321.jpg

| Kinect V2 | RGB | Depth | Infrared |
|---|---|---|---|
| Resolution | 1920x1080 | 512x424 | 512x424 |
| Shutter Type | Rolling | Global | Global |

# KinectV2 Recorder

# Relative Rigid Transformation

Robotic Arm: KUKA KR 16-2 CR
- Payload: 16 kg
- Repeatability: 0.05 mm
- Max. reach: 1610 mm
- Number of axes: 6

http://img.directindustry.com/images_di/photo-mg/17587-2781073.jpg

# 7 Motion Patterns

## with 5 Speed Levels

# Take

# A Benchmark Dataset for 6DoF Object Pose Tracking

Po-Chen Wu[1]    Yueh-Ying Lee[1*]    Hung-Yu Tseng[1*]    Hsuan-I Ho[1*]    Ming-Hsuan Yang[2]    Shao-Yi Chien[1]

[1] Media IC & System Lab, National Taiwan University
[2] Vision and Learning Lab, University of California, Merced
(* indicates equal contribution)

## Abstract

Accurately tracking the six degree-of-freedom pose of an object in real scenes is an important task in computer vision and augmented reality with numerous applications. Although a variety of algorithms for this task have been proposed, it remains difficult to evaluate existing methods in the literature as oftentimes different sequences are used and no large benchmark datasets close to real scenarios are available. In this paper, we present a large object pose tracking benchmark dataset consisting of RGB-D video sequences of 6 targets with 6 objects are recorded under multiple lighting conditions, different motion patterns and speeds with the help of a programmable robotic arm. We present extensive quantitative evaluation results of the state-of-the-art methods on this benchmark dataset and discuss the potential research directions in this field.

http://media.ee.ntu.edu.tw/research/OPT/

# Marker Corner Localization

# Corner Localization Accuracy

# Ground-truth Pose



Pose Viewer

# Depth Calibration



Depth Image

Infrared Image

# Measured Depth Data



Infrared Image

Point Cloud (a)

Point Cloud (b)

# Masked Image

# Dataset Comparison

| Benchmark | Device | Mechanism | Pose Establishment | Video Clips | # 2D Targets | # 3D Targets | # Motion Patterns | # Frames |
|---|---|---|---|---|---|---|---|---|
| Lieberknecht [1] | Marlin F-080C | Handheld | Marker-based | Yes | 8 | -- | 5 | 48,000 |
| Gauglitz [2] | Fire-i | Manually Operated Contraption | Direct Alignment | Yes | 6 | -- | 16 | 6,889 |
| Hinterstoisser [3] | Kinect V1 | Handheld | Marker-based | No | -- | 15 | -- | 18,000 |
| Tejani [4] | Kinect V1 | Handheld | Marker-based | No | -- | 3 | -- | 5,229 |
| Brachmann [5] | Kinect V1 | Handheld | Marker-based | No | -- | 20 | 3 | 10,000 |
| Rennie [6] | Kinect V1 | Robotic Arm | Manual | No | -- | 24 | -- | 10,368 |
| Krull [7] | Kinect V1 | Handheld | ICP | Yes | -- | 3 | -- | 1,000 |
| Choi [8] | Synthetic | -- | Synthetic | Yes | -- | 4 | -- | 4,000 |
| Proposed | **Kinect V2** | **Programmable Robotic Arm** | **Checkerboard-based** | Yes | 6 | 3 | 23 | 100,956 |

1.  S. Lieberknecht, S. Benhimane, P. Meier, and N. Navab. A Dataset and Evaluation Methodology for Template-based Tracking. In ISMAR, 2009.
2.  S. Gauglitz, T. H¨ollerer, and M. Turk. Evaluation of Interest Point Detectors and Feature Descriptors for Visual Tracking. IJCV, 94(3):335– 360, 2011.
3.  S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab. Model Based Training, Detection and Pose Estimation of Texture-Less 3D Objects in Heavily Cluttered Scenes. In ACCV, 2012.
4.  A. Tejani, D. Tang, R. Kouskouridas, and T.-K. Kim. Latent-Class Hough Forests for Object Detection and Pose Estimation. In ECCV, 2014.
5.  E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother. Learning 6D Object Pose Estimation Using 3D Object Coordinates. In ECCV, 2014.
6.  C. Rennie, R. Shome, K. E. Bekris, and A. F. De Souza. A Dataset for Improved RGBD-based Object Detection and Pose Estimation for Warehouse Pick-and-Place. RAL, 1(2):1179–1185, 2016.
7.  A. Krull, F. Michel, E. Brachmann, S. Gumhold, S. Ihrke, and C. Rother. 6-DOF Model Based Tracking via Object Coordinate Regression. In ACCV, 2014.
8.  C. Choi and H. I. Christensen. RGB-D Object Tracking: A Particle Filter Approach on GPU. In IROS, 2013.

# Evaluation on 2D Datasets

# Evaluation on 3D Datasets

# Feature-Based Method

Feature Matching → Outlier Removal → PnP Algorithm



Rely on Natural Features

# Would Fail When...



Textureless



Blurry

Noisy
Depth
Data

# Direct Pose Estimation (DPE)

# Approximate Pose Estimation (APE)

- Branch-and-bound Algorithm in Pose Space
  - Find the pose with minimum appearance error $E_a$

$$E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}|I_c(\mathbf{u}_i) - O_t(\mathbf{x}_i)|$$



$E_a$

$\varepsilon \;\; \varepsilon'$

Pose Space (1D View)

# Approximate Pose Estimation (APE)

- Failure analysis



$E_a$

$\mathcal{E}$

Pose Space (1D View)

# Ɛ-covering Set

Pose Domain

# Pose Jumping



Original Video

Pose Estimation Result

# Pose Ambiguity

# Explanation of Pose Ambiguity

- Multiple local minima of a cost function (due to coplanar points)

# Pose Refinement (PR)

- Refine and disambiguate the approximately estimated pose

# Gauss-Newton Iteration

$$E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i(\mathbf{p})) - O_t(\mathbf{x}_i)\right)^2$$

- $\Delta\mathbf{p}^* = \underset{\Delta\mathbf{p}}{\mathrm{argmin}}\,\frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i(\mathbf{p}_c + \Delta\mathbf{p})) - O_t(\mathbf{x}_i)\right)^2$

$$\approx \underset{\Delta\mathbf{p}}{\mathrm{argmin}}\,\frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i(\mathbf{p}_c)) + \left.\frac{\partial I_c}{\partial\mathbf{p}}\right|_{\mathbf{p}=\mathbf{p}_c}\Delta\mathbf{p} - O_t(\mathbf{x}_i)\right)^2$$

Vectorization $\left.\dfrac{\partial \mathbf{I}_c}{\partial\mathbf{p}}\right|_{\mathbf{p}=\mathbf{p}_c} \equiv \mathbf{J}_c$ 

$E_a{}'(\mathbf{p}) = 0$ 

$$\mathbf{J}_c\Delta\mathbf{p} = \mathbf{O}_t - \mathbf{I}_c$$

$$\Delta\mathbf{p} = \left(\mathbf{J}_c^{\mathrm{T}}\mathbf{J}_c\right)^{-1}\mathbf{J}_c^{\mathrm{T}}(\mathbf{O}_t - \mathbf{I}_c)$$

# Jacobian Matrix $\mathbf{J}_c \equiv \left.\dfrac{\partial \mathbf{I}_c}{\partial \mathbf{p}}\right|_{\mathbf{p}=\mathbf{p}_c}$

- Chain rule

$$\mathbf{J}_c = \frac{\partial \mathbf{I}_c}{\partial \mathbf{p}} = \begin{bmatrix} \dfrac{\partial I_c(\mathbf{u}_1)}{\partial \mathbf{p}} \\ \vdots \\ \dfrac{\partial I_c(\mathbf{u}_n)}{\partial \mathbf{p}} \end{bmatrix}, \frac{\partial I_c}{\partial \mathbf{p}} = \frac{\partial I_c}{\partial \mathbf{u}} \begin{bmatrix} \dfrac{\partial \mathbf{u}}{\partial \mathbf{r}}, \dfrac{\partial \mathbf{u}}{\partial \mathbf{t}} \end{bmatrix} = \begin{bmatrix} \dfrac{\partial I_c}{\partial u}, \dfrac{\partial I_c}{\partial v} \end{bmatrix} \begin{bmatrix} \dfrac{\partial \mathbf{u}}{\partial \hat{\mathbf{x}}} \dfrac{\partial \hat{\mathbf{x}}}{\partial \hat{\mathbf{R}}} \dfrac{\partial \hat{\mathbf{R}}}{\partial \mathbf{r}}, \dfrac{\partial \mathbf{u}}{\partial \hat{\mathbf{x}}} \end{bmatrix},$$

$$\frac{\partial \mathbf{u}}{\partial \hat{\mathbf{x}}} = \begin{bmatrix} \dfrac{f_x}{\hat{z}} & 0 & -\dfrac{f_x \hat{x}}{\hat{z}^2} \\ 0 & \dfrac{f_y}{\hat{z}} & -\dfrac{f_y \hat{y}}{\hat{z}^2} \end{bmatrix}, \frac{\partial \hat{\mathbf{x}}}{\partial \hat{\mathbf{R}}} = \begin{bmatrix} x & y & 0 & 0 & 0 & 0 \\ 0 & 0 & x & y & 0 & 0 \\ 0 & 0 & 0 & 0 & x & y \end{bmatrix}, \partial \hat{\mathbf{R}} = \begin{bmatrix} R_{11} \\ R_{12} \\ R_{21} \\ R_{22} \\ R_{31} \\ R_{32} \end{bmatrix}, \partial \hat{\mathbf{x}} = \begin{bmatrix} \hat{x} \\ \hat{y} \\ \hat{z} \end{bmatrix} = \begin{bmatrix} R_{11} & R_{12} & t_x \\ R_{21} & R_{22} & t_y \\ R_{31} & R_{32} & t_z \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

# Jacobian Matrix $\mathbf{J}_c \equiv \left[\dfrac{\partial I_c}{\partial u}, \dfrac{\partial I_c}{\partial v}\right]\left[\dfrac{\partial \mathbf{u}}{\partial \hat{\mathbf{x}}}\dfrac{\partial \hat{\mathbf{x}}}{\partial \widehat{\mathbf{R}}}\dfrac{\partial \widehat{\mathbf{R}}}{\partial \mathbf{r}}, \dfrac{\partial \mathbf{u}}{\partial \hat{\mathbf{x}}}\right]$

- The rotation is parameterized as <span style="color:#29ABE2">rotation vector</span>

$$\mathbf{p} = \begin{bmatrix} \mathbf{r} \\ \mathbf{t} \end{bmatrix}, \mathbf{r} = \begin{bmatrix} r_x \\ r_y \\ r_z \end{bmatrix} \in \mathbb{R}^3, \mathbf{t} = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix} \in \mathbb{R}^3$$

- The derivative of $\mathbf{R}$ with respect to $\mathbf{r}$[1]:

$$\frac{\partial \mathbf{R}}{\partial r_a} = \frac{r_a[\mathbf{r}]_\times + [\mathbf{r} \times (\mathbf{I} - \mathbf{R})\mathbf{e}_i]_\times}{\|\mathbf{r}\|^2}\mathbf{R}, \qquad a = x, y, z$$

  - $\mathbf{I}$ : identity matrix

  - $\mathbf{e}_i$ : the $i$-th vector of the standard basis in $\mathbb{R}^3$

1.   G. Gallego and A. Yezzi, "A Compact Formula for the Derivative of a 3-D Rotation in Exponential Coordinates," JMIV, vol. 51, no. 3, pp. 378–384, 2015.

# Synthetic Dataset



Templates[1]

Background Images[2]

Test Images

1. Lieberknecht, Sebastian, et al. "A dataset and evaluation methodology for template-based tracking algorithms." ISMAR 2009
2. Jegou, Herve, Matthijs Douze, and Cordelia Schmid. "Hamming embedding and weak geometric consistency for large scale image search." ECCV, 2008.

# Evaluated Algorithms

1. SIFT + OPnP

2. SIFT + IPPE

3. ASIFT + OPnP

4. ASIFT + IPPE

5. APE (Approximate Pose Estimation)

6. DPE (Direct Pose Estimation)

1. Lowe, David G. "Distinctive image features from scale-invariant keypoints." IJCV, 2004.
2. Morel, Jean-Michel, and Guoshen Yu. "ASIFT: A new framework for fully affine invariant image comparison." SIIMS, 2009.
3. Zheng, Yinqiang, et al. "Revisiting the pnp problem: A fast, general and optimal solution." ICCV, 2013.
4. Collins, Toby, and Adrien Bartoli. "Infinitesimal plane-based pose estimation." IJCV, 2014

# Evaluated Metric

– Rotation error (degree)

$$E_r = \mathrm{acosd}\left(\frac{\mathrm{Tr}(\mathbf{R}^{\mathrm{T}} \cdot \widetilde{\mathbf{R}})}{2}\right)$$

– Translation error (%)

$$E_t = \frac{\|\widetilde{\mathbf{t}} - \mathbf{t}\|}{\|\widetilde{\mathbf{t}}\|} \times 100$$

– Success rate (%)

➢ The percentage of poses that $E_r < 20°$ and $E_t < 10\%$

# Evaluation Results

- Results with undistorted test images.



| Method | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR | $E_r$ | $E_t$ | SR |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SIFT+IPPE | 0.85 | 0.34 | 40.0 | 1.90 | 0.54 | 96.0 | 0.23 | 0.25 | 28.0 | 0.32 | 0.24 | 86.0 | 0.74 | 0.35 | 92.0 | 0.56 | 0.40 | 98.0 | 1.15 | 0.50 | 30.0 | 0.28 | 0.37 | 96.0 |
| SIFT+OPnP | 0.76 | 0.40 | 40.0 | 1.18 | 0.46 | 96.0 | 0.20 | 0.24 | 28.0 | 0.25 | 0.24 | 86.0 | 0.56 | 0.32 | 92.0 | 0.55 | 0.43 | 98.0 | 1.48 | 0.47 | 30.0 | 0.25 | 0.36 | 96.0 |
| ASIFT+IPPE | 9.70 | 2.92 | 20.0 | 2.96 | 0.81 | 94.0 | 1.48 | 0.43 | **100** | 1.65 | 0.51 | 94.0 | 1.59 | 0.57 | **100** | 1.29 | 0.34 | 98.0 | 2.17 | 0.52 | 52.0 | 1.96 | 0.36 | 90.0 |
| ASIFT+OPnP | 8.20 | 2.22 | 22.0 | 2.72 | 0.74 | **100** | 1.38 | 0.41 | **100** | 1.53 | 0.45 | 96.0 | 1.40 | 0.50 | 98.0 | 1.26 | 0.35 | **100** | 1.33 | 0.37 | 52.0 | 1.80 | 0.36 | 94.0 |
| APE | 1.10 | 0.33 | **100** | 1.44 | 0.42 | **100** | 0.90 | 0.47 | 98.0 | 2.56 | 1.23 | 94.0 | 1.03 | 0.35 | **100** | 1.63 | 0.49 | **100** | 1.96 | 0.91 | **100** | 1.57 | 0.68 | 98.0 |
| DPE | **0.39** | **0.17** | **100** | **0.42** | **0.24** | **100** | **0.16** | **0.14** | **100** | **0.16** | **0.12** | **98.0** | **0.21** | **0.16** | **100** | **0.21** | **0.11** | **100** | **0.15** | **0.14** | **100** | **0.17** | **0.13** | **100** |

$E_r(°)$

$E_t(\%)$

SR(%)

(a) Gaussian blur

(b) JPEG compression

(c) Intensity change

(d) Tilt angle

# Overall Evaluation

# Visual Tracking Dataset[1]



Templates ✕ Motion Patterns

- Unconstrained
- Panning
- Rotation
- Perspective Distortion
- Zoom
- Static Lighting
- Dynamic Lighting
- Motion Blur x9

1.    Gauglitz, Steffen, Tobias Höllerer, and Matthew Turk. "Evaluation of interest point detectors and feature descriptors for visual tracking." IJCV, 2011

Ground Truth

SIFT + OPnP

DPE

ASIFT + OPnP

# Overall Evaluation

# OPT Dataset (Proposed)



Templates

×

- Translation
- Zoom
- In-plane Rotation
- Out-of-plane Rotation
- Flashing Light
- Moving Light
- Free Motion

Motion Patterns

×

5

Speeds

Ground Truth

SIFT + OPnP

DPE

ASIFT + OPnP

# Overall Evaluation

# How About This?

# DPE Demo Video

# Average Runtime

- Average runtime (measured in seconds) using MATLAB.
  - Core i7-6700K 4.0 GHz processor
  - 32 GB RAM
  - NVIDIA GTX 970 GPU

- Numbers in parentheses denote the average runtime of the CUDA implementation.

| Dataset | SIFT-based Approach | ASIFT-based Approach | DPE | | |
| --- | --- | --- | --- | --- | --- |
| | | | APE | PR | Total |
| Synthetic | 7.446 | 10.912 | 10.549 (**1.505**) | 0.571 (**0.117**) | 11.120 (**1.622**) |
| VT | 3.618 | 15.814 | 17.920 (**1.217**) | 0.694 (**0.180**) | 18.615 (**1.397**) |
| OPT | 11.364 | 38.944 | 18.545 (**0.994**) | 0.214 (**0.088**) | 18.759 (**1.082**) |

# DodecaPen: Puppy



Input Frames

Pen-tip Trajectory

# How To Use?

Surface
Calibration

# Related Work

m-Sequence Projection

Sensor

**Lumitrack**
(*UIST 2013*)

Accuracy: 5mm

Sensor View

0   m-sequence in sensor memory   800

Sensor Position: **X=651**

Interaction space

10 cm

Pen

Sensor

Tablet surface

**IrPen**
*(CGA 2014)*

Accuracy: 10mm

# Light Chisel
(*PG 2015*)

**Accuracy: 2mm**

# DodecaPen
(*Proposed*)

Accuracy: 0.4mm

# How To Implement?

# Proposed 6DoF Pose Tracking System



Camera      DodecaPen

Input Frames

Inter-frame Corner Tracking (**ICT**)

Marker Intensity Normalization

*No*

Did **APE** Succeed?

*Yes*

Approximate Pose Estimation (**APE**)

Dense Pose Refinement (**DPR**)

Digital 2D Drawing

Output Pen-tip Trajectory (Based on DodecaPen Poses)

# Approximate Pose Estimation (APE)





- Marker Detection
- Minimize reprojection error $E_r(\mathbf{p})$ with P$n$P algorithm to get the initial pose $\mathbf{p}'$

$$E_r(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\|\widehat{\mathbf{u}}_i - \mathbf{u}_i(\mathbf{x}_i; \mathbf{p})\|^2$$

$\widehat{\mathbf{u}}$: detected point in the camera image
$\mathbf{x}$: point on the dodecahedron
$\mathbf{u}$: transformed $\mathbf{x}$ point in the camera image
$\mathbf{p}$: pose (including rotation matrix $\mathbf{R}$ and translation vector $\mathbf{t}$)

90

# Inter-frame Corner Tracking (ICT)

If APE does not succeed…

- Pyramidal Lucas-Kanade marker corner tracking
- P$n$P algorithm to get the initial pose $\mathbf{p}'$

91

# Marker Intensity Normalization



- We normalize the intensity to ensure intensity invariance before minimizing the residual between the model and image.

# Dense Pose Refinement (DPR)



- Minimize appearance distance $E_a(\mathbf{p})$ with Gauss Newton and backtracking line search (BLS) to get the final pose $\mathbf{p}^*$



$$E_a(\mathbf{p}) = \frac{1}{n} \sum_{i=1}^{n} \left( I_c(\mathbf{u}_i(\mathbf{p})) - O_t(\mathbf{x}_i) \right)^2$$

$I_c$: camera image

$O_t$: target object

$\mathbf{x}$: point on the dodecahedron

$\mathbf{u}$: transformed $\mathbf{x}$ point in the camera frame

$\mathbf{p}$: pose (including rotation matrix $R$ and translation vector $\mathbf{t}$)

# Gauss-Newton Iteration

$$E_a(\mathbf{p}) = \frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i(\mathbf{p})) - O_t(\mathbf{x}_i)\right)^2$$

- $\Delta\mathbf{p}^* = \underset{\Delta\mathbf{p}}{\operatorname{argmin}}\frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i(\mathbf{p}' + \Delta\mathbf{p})) - O_t(\mathbf{x}_i)\right)^2$

$$\approx \underset{\Delta\mathbf{p}}{\operatorname{argmin}}\frac{1}{n}\sum_{i=1}^{n}\left(I_c(\mathbf{u}_i(\mathbf{p}')) + \left.\frac{\partial I_c}{\partial \mathbf{p}}\right|_{\mathbf{p}=\mathbf{p}'}\Delta\mathbf{p} - O_t(\mathbf{x}_i)\right)^2$$

$$\Delta\mathbf{p} = \left(\mathbf{J}_c^{\mathrm{T}}\mathbf{J}_c\right)^{-1}\mathbf{J}_c^{\mathrm{T}}(\mathbf{O}_t - \mathbf{I}_c)$$

$$\mathbf{J}_c \equiv \left.\frac{\partial \mathbf{I}_c}{\partial \mathbf{p}}\right|_{\mathbf{p}=\mathbf{p}_c}$$

1. Chain rule
2. Rotation vector $\mathbf{r}$
3. $\frac{\partial \mathbf{R}}{\partial \mathbf{r}}$

# Backtracking Line Search (BLS)

- Gauss-Newton iteration does not always converge with a fixed step size since our least squares problem is nonlinear.

- We shrink $\Delta\mathbf{p}$ by $\Delta\mathbf{p} \leftarrow \alpha\Delta\mathbf{p}$ until it meets the *Armijo-Goldstein condition* below:

$$E_a(\mathbf{p}' + \Delta\mathbf{p}) \leq E_a(\mathbf{p}') + c\nabla E_a(\mathbf{p}')^{\mathrm{T}}\Delta\mathbf{p}$$

  – $\nabla E_a(\mathbf{p}')$ is the local function gradient
  – $\alpha = 0.5$ , $c = 10^{-4}$

# Masked Mipmaps

Mipmap Masks

Marker Mipmaps

∩

Masked Mipmaps

# Why Dodecahedron?

# Pose Jumping!

Multiple Candidates due to Coplanar Points

# Platonic Solid

## A Regular [...] olyhedron

| Tetrahedron | Cube | [...]ecahedron | Icosahedron |
|---|---|---|---|
| Four faces | Six fa[...] | [...] faces | Twenty faces |

# The Chosen One

Ideal ⟺ Real

😭😭 Fail 😭😭

# Dodecahedron Calibration

# Dodecahedron Calibration (DC)

- Determine the precise pose of each marker with respect to the dodecahedron

- One-time offline bundle adjustment

$$E_a(\{\mathbf{p}_j, \mathbf{p}_k\}) = \sum_i \sum_j \sum_k \left( I_c \left( \mathbf{u}_i(\mathbf{x}_i; \mathbf{p}_j; \mathbf{p}_k) \right) - O_t(\mathbf{x}_i) \right)^2$$

$I_c$: camera image

$O_t$: target object

$\mathbf{x}$: point on the dodecahedron

$\mathbf{u}$: transformed $\mathbf{x}$ point in the camera frame

$\mathbf{p}$: dodecahedron pose

$\mathbf{p}$: marker pose

# Pen-tip Calibration

# 6 DoF Pose Tracking

# Performance Analysis

# Evaluation with Synthetic Data

24 Motion Patterns
4 Varying Conditions

**Pen-tip Trajectory 01**

APE [5.074]
Proposed [0.244]
Ground truth

**Pen-tip Trajectory 02**

APE [5.423]
Proposed [0.380]
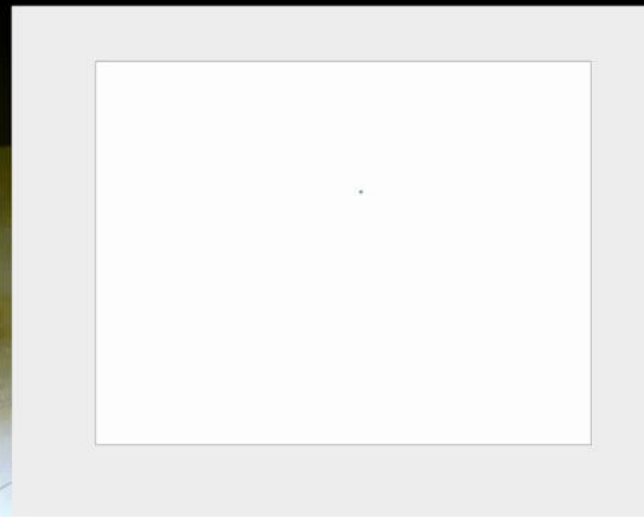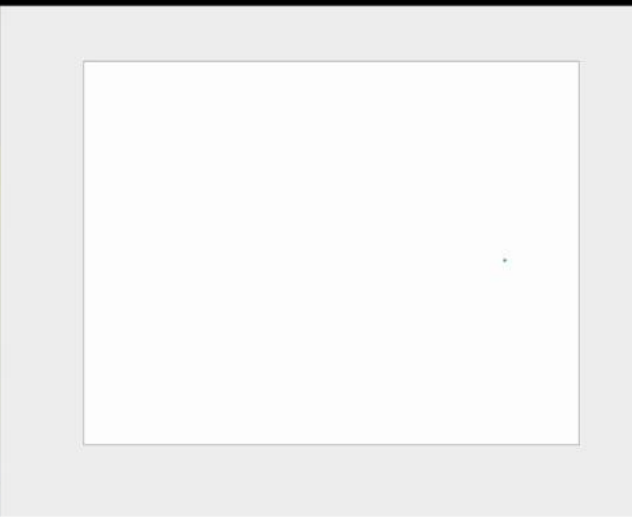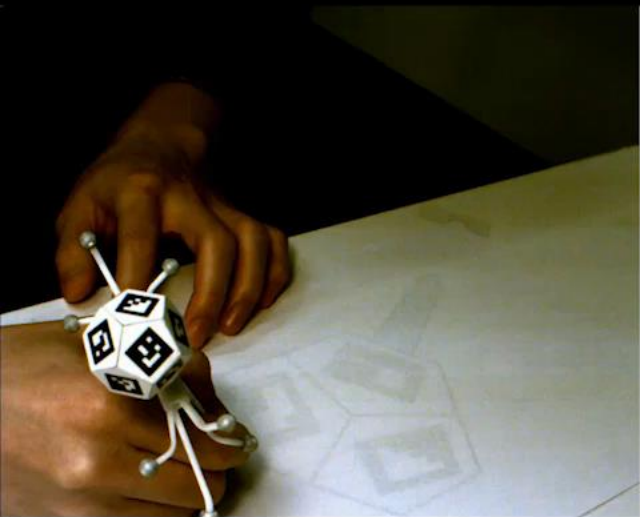Ground truth

# Evaluation with Real Data

4 Real Drawings
VS. Mocap System (16 Cameras)

Boba

Thumb

DodecaPen

UIST2017
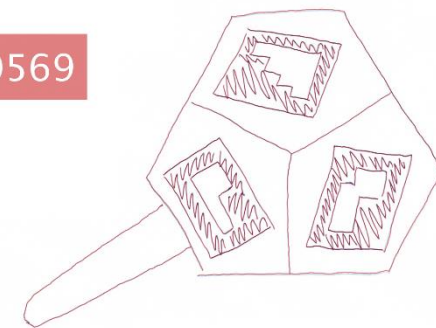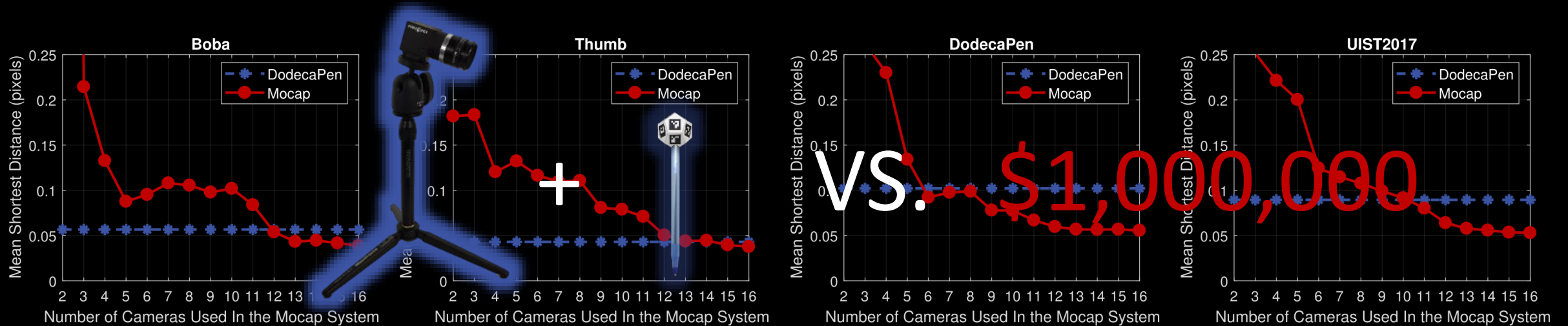
# DodecaPen VS. Mocap

Comparable to a mocap system
with 10 active cameras

# Main Achievements

1. ## OPT Dataset
   - A benchmark dataset for 6DoF object pose tracking

2. ## DPE Algorithm
   - A robust pose estimation method for planar objects

3. ## DodecaPen
   - A submillimeter-accurate 6DoF tracking solution

# Future Work

- Learning-based pose estimation followed by dense pose refinement for general objects.

- Marker-based accurate 6DoF pose estimation and tracking solution.

- Pose recovering for planar objects using depth information and filtering techniques.