



URAI 2017

14th International Conference on Ubiquitous Robots and Ambient Intelligence
June 28 - July 1, 2017 / Maison Glad Jeju, Jeju, Korea

Synthetic Learning Set for Object Pose Estimation: *Initial Experiments*

JOO-HAENG LEE, WOO-HAN YUN, JAEYEON LEE, JAEHONG KIM

HUMAN-MACHINE INTERACTION LAB

ETRI



Outline

GIVEN:

- **Target objects** for **pose estimation** in robotic manipulation task
- A small set of real images for training

PROBLEM:

- To obtain a **data set** with **quantitative diversity** in condition and **qualitative precision** in annotation.
- While considering **constraints** in human *labor* and *error*, production cost, physical limits in environments, and deploy *time*.

PROPOSED SOLUTION:

- To synthesize a learning set using **computer graphics**: SLS (Synthetic Learning Set)
- **Diversity**: pose (translation+rotation), illumination, ...
- **Precise** and **rich** annotation: 3D coord, normal, depth, camera extrinsic, ...

RESULTS

- Synthesis pipeline based on open source tools: **Blender**
- **Synthetic Learning Set** (SLS): 516,672 images = 23 objects * 5,616 viewpoints * 4 types (RGB, xyz, normal, depth)
- **Initial experiments** of applying SLS on pose estimation
- Soon open at **github**: <https://github.com/joohaeng/SLS4ML>

Background

- Object **pose estimation** for robotic manipulation.
 - Diverse **objects**: shape, physical property,
 - Complex and dynamic **environments**
 - Machine learning requires a nice **data set** with **quality annotation** and/or **applicable in simulation**



Amazon Robotics Challenge 2016—Winner Team & Contest Objects

joohaeng@etri.re.kr

Examples of Acquiring a Real Learning Set

A subset of
APC 2015 objects



Manual acquisition



Labor
Cost
Error
Precision
Diversity

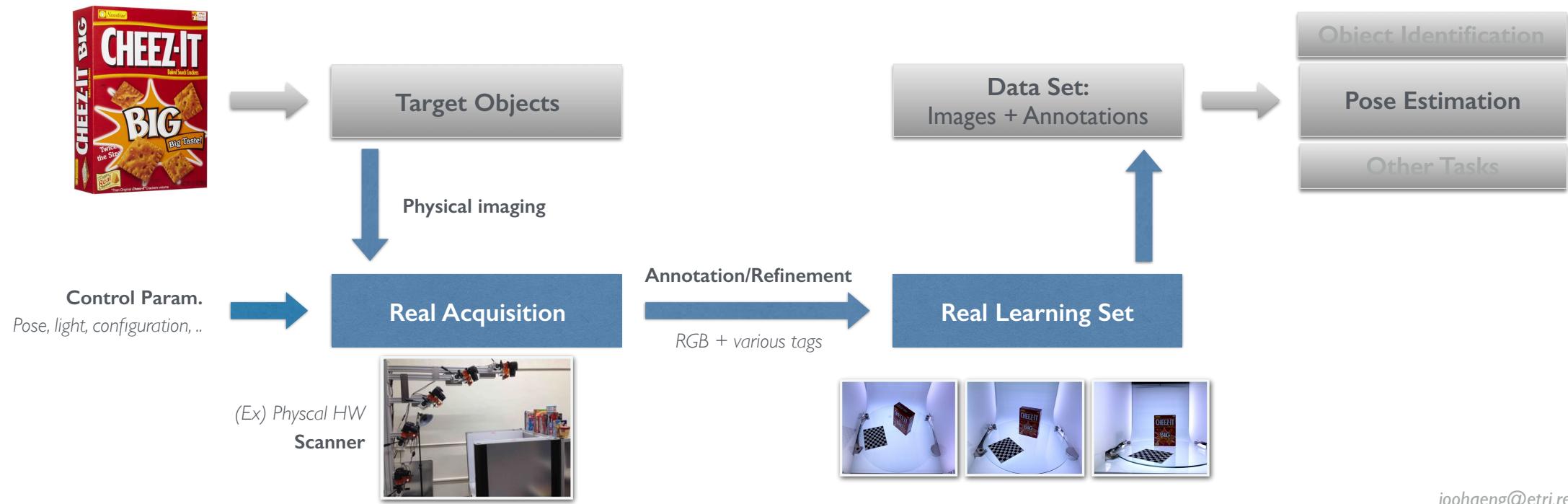
...

joohaeng@etri.re.kr

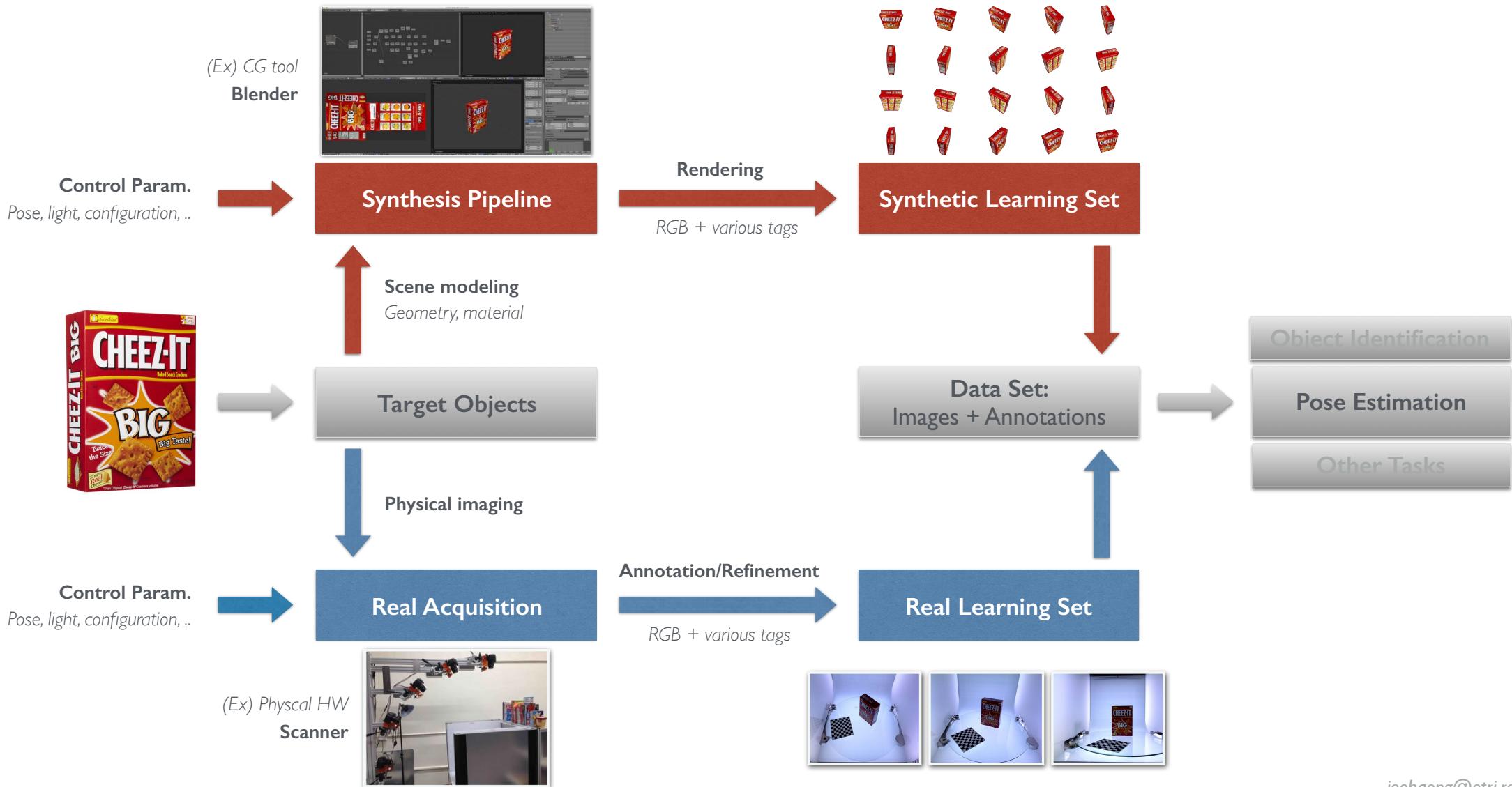
Proposed Solution: Synthetic Learning Set



Proposed Solution: Synthetic Learning Set



Proposed Solution: Synthetic Learning Set



Related Work

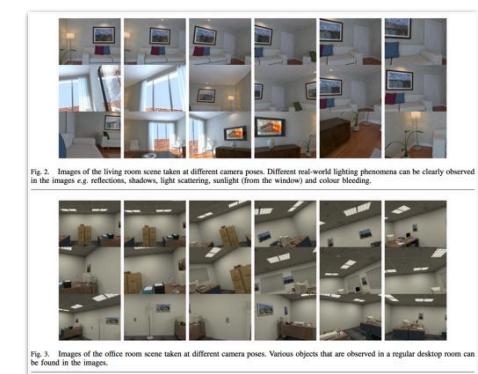
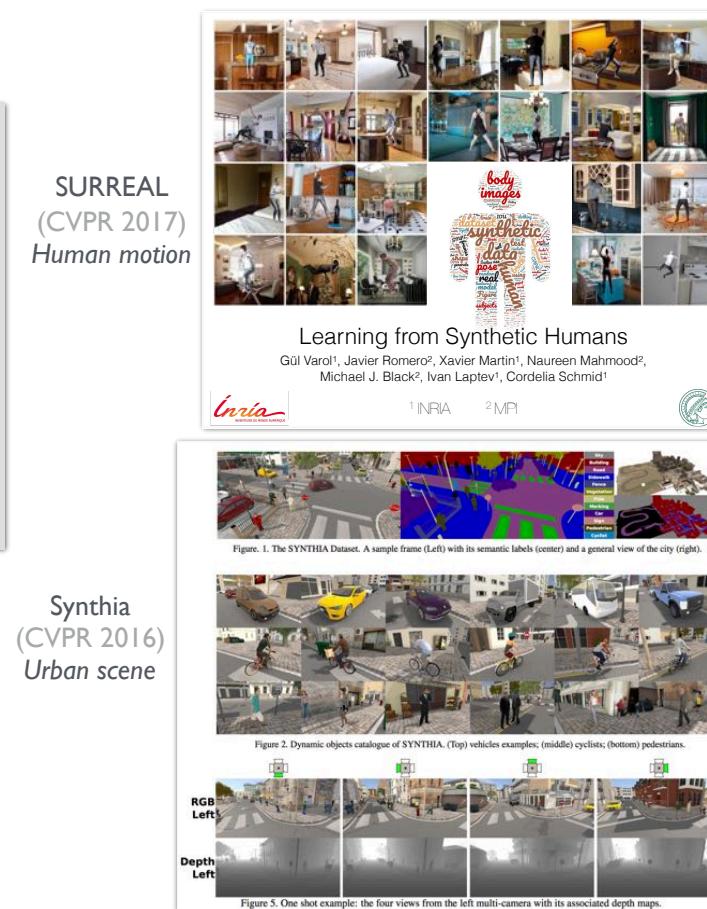
- A list of synthetic dataset and tools for computer vision
 - <https://github.com/unrealcv/synthetic-computer-vision>

Related Work

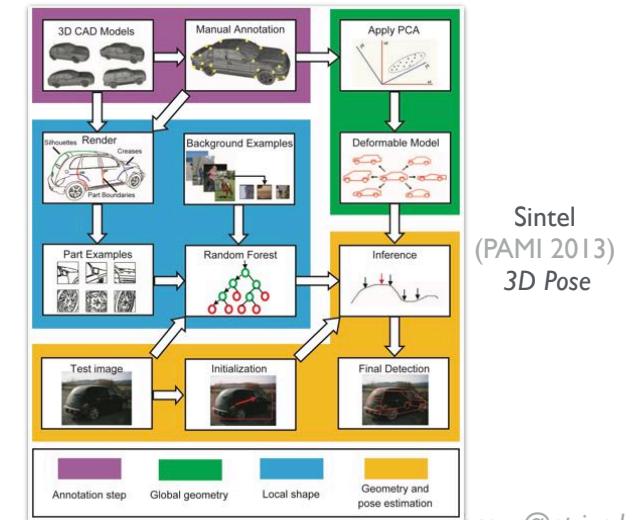
- A list of synthetic dataset and tools for computer vision
 - <https://github.com/unrealcv/synthetic-computer-vision>

Synthetic image dataset

- SURREAL
- Virtual KITTI
- Synthia
- Sintel, A synthetic dataset for optical flow
- SceneFlow
- 4D Light Fields
- ICL-NUIM dataset
- Driving in the Matrix



ICL-NUIM
(ICRA 2014)
SLAM



Sintel
(PAMI 2013)
3D Pose

joohaeng@etri.re.kr

Related Work

- A list of synthetic dataset and tools for computer vision
 - <https://github.com/unrealcv/synthetic-computer-vision>

Synthetic image dataset

- SURREAL
- Virtual KITTI
- Synthia
- Sintel, A synthetic dataset for optical flow
- SceneFlow
- 4D Light Fields
- ICL-NUIM dataset
- Driving in the Matrix

Tools

- Render SMPL human bodies on [Blender](#), see CVPR2017
- Render for CNN, based on [Blender](#), see ICCV2015
- UETorch, based on [UE4](#), see ICML2016
- UnrealCV, based on [UE4](#), see ArXiv
- VizDoom, based on Doom, see ArXiv
- OpenAI Universe, see project page
- [Blender](#) addon for 4D light field rendering, see project page

Resources

[Virtual/Augmented Reality for Visual Artificial Intelligence \(ECCV 2016 Workshop\)](#)

[Role of Simulation in Computer Vision \(ICCV 2017 Workshop\)](#)

[Virtual Reality Meets Physical Reality \(Siggraph Asia 2016 workshop\)](#)

[CVPR 2017 Workshop THOR Challenge](#)

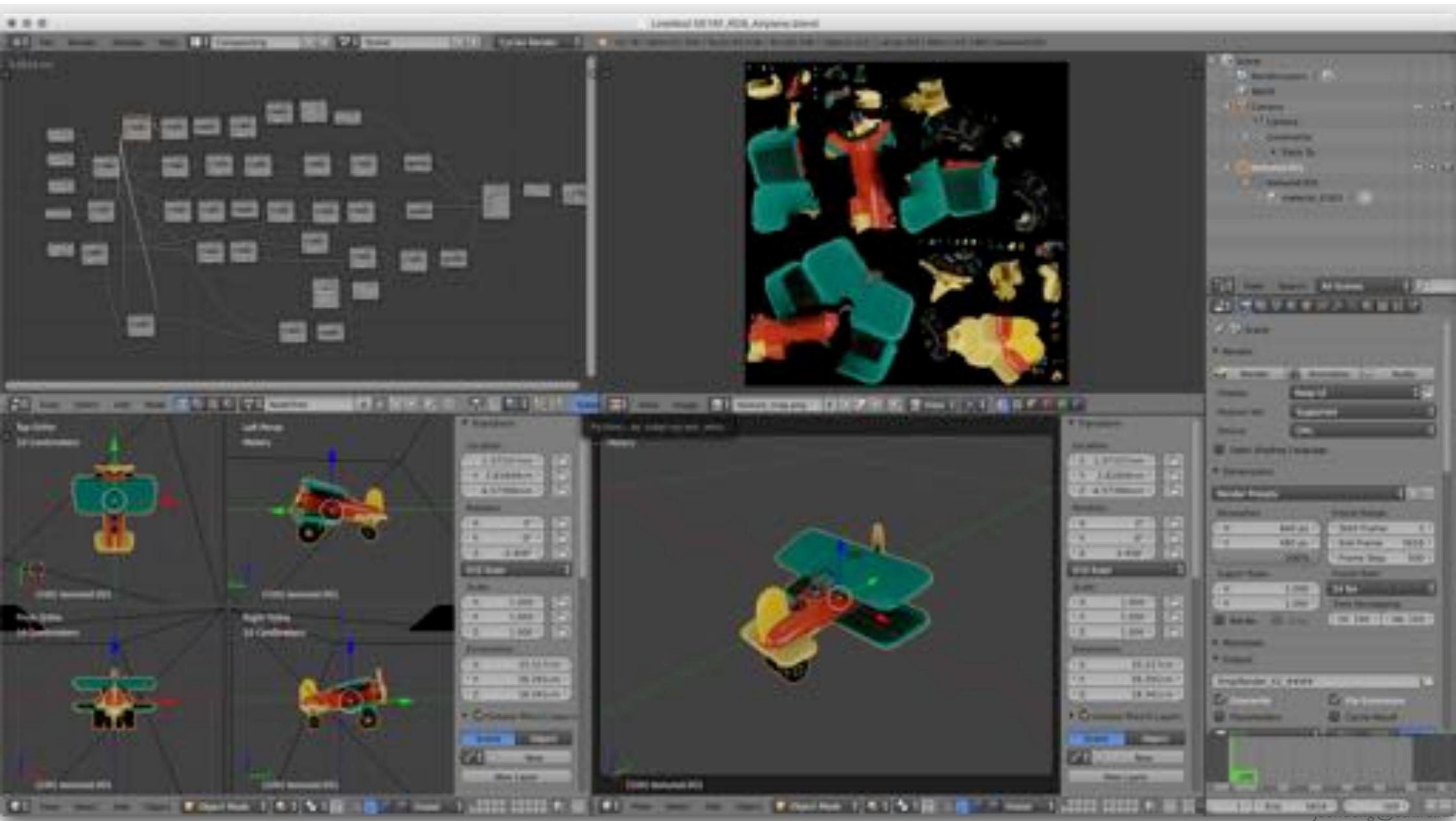
3D Model Repository

- ShapeNet
- 3dscan
- seeing3Dchairs

Synthesis Pipeline

- **Blender**

- ... **free** and **open** source 3D creation suite.
- It supports the **entirety of the 3D pipeline**—modeling, rigging, animation, simulation, rendering, compositing and motion tracking, even video editing and game creation.



Scene Modeling for SLS

- **Objects**
 - Geometry, material (RGB texture, BRDF), deformation, ...
- **Camera**
 - Pose, intrinsics, ...
- **Light**
 - Env map lighting

Object Modeling

- Manual modeling
 - **Cuboid**
 - Free-form
- Automatic
 - 3D scanner
- Open data
 - **YCB data set**



Cuboid-like Objects

9 objects from
APC 2015

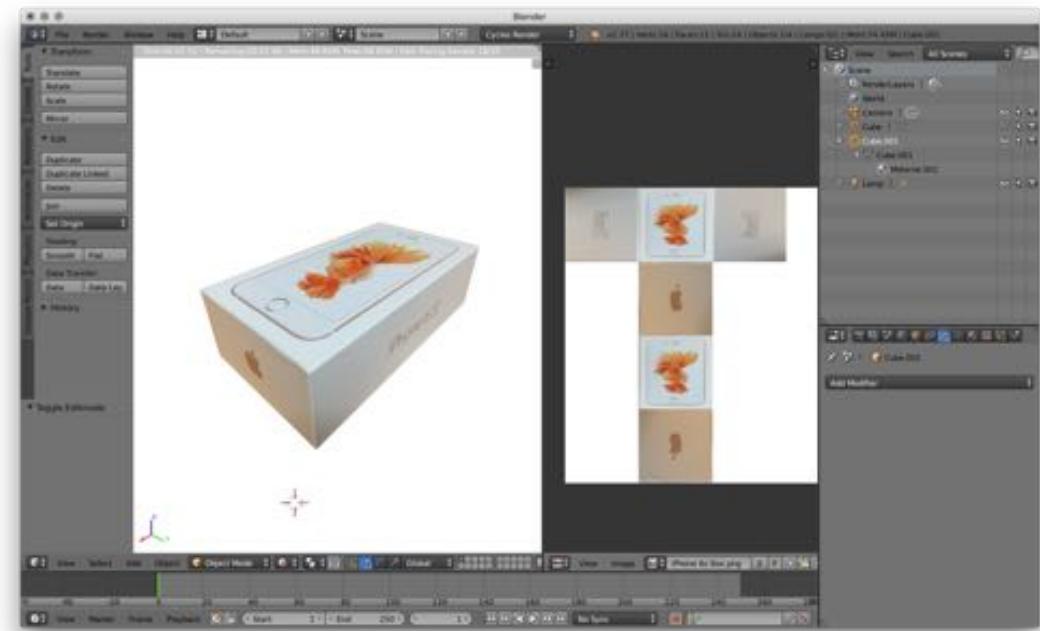
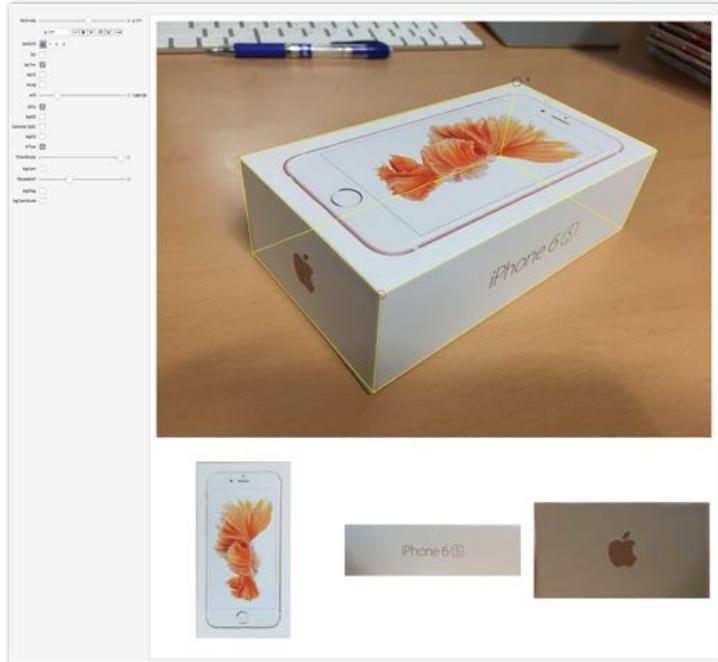


5 objects from
Korean grocery



joohaeng@etri.re.kr

Modeling — Cuboids



joohaeng@etri.re.kr

Free-form Objects

9 objects from
YCB Benchmark



YCB Benchmark Dataset

50.87.248.224

YCB Benchmarks – Object and Model Set

Benchmarking for robotic manipulation



Home News Object Set Object Models Protocols & Benchmarks Results & Records Publications Forums

Not Found

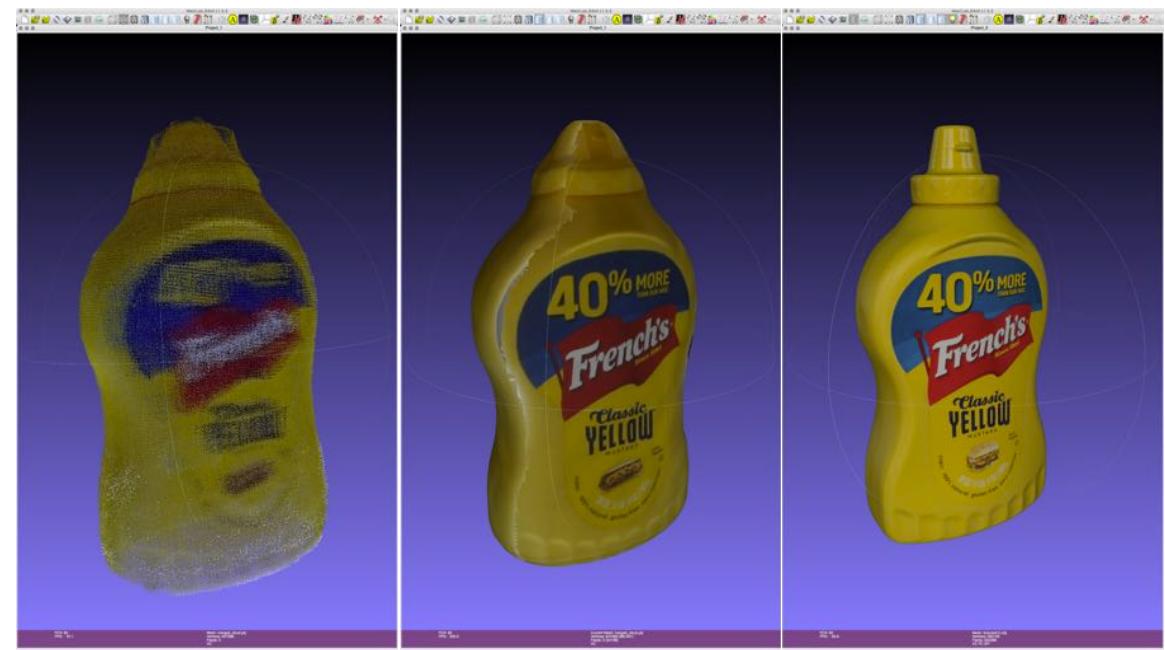
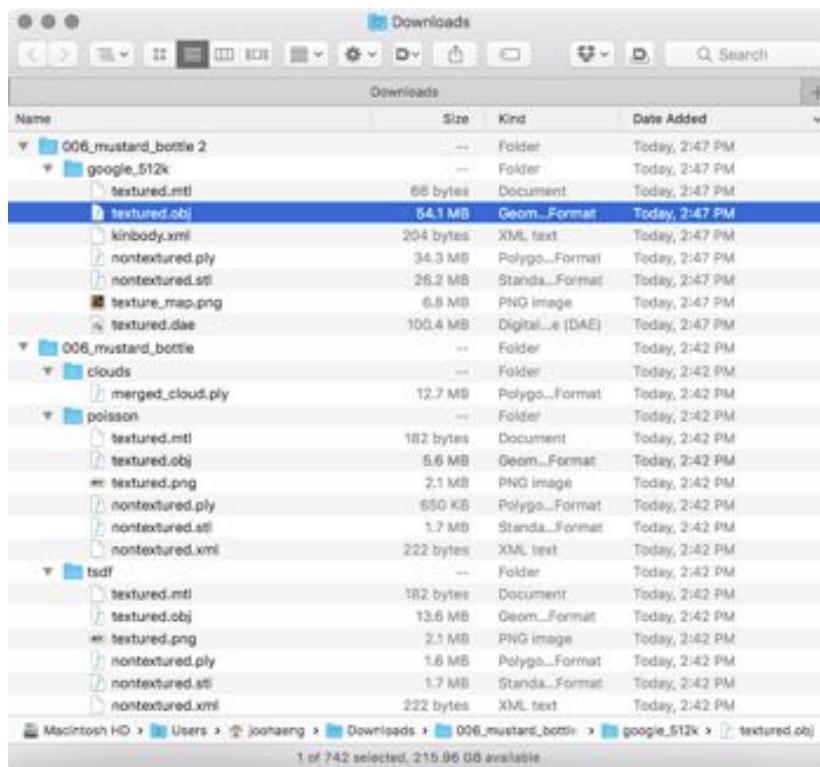
Apologies, but the page you requested could not be found. Perhaps searching will help.

ycb-benchmarks.s3-website-us-east-1.amazonaws.com

Objects

Object Name	Processed Data	Raw RGB	Raw RGBD	16k laser scan	64k laser scan	512k laser scan
001_chips_can	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	N/A [†]	N/A [†]	N/A [†]
002_master_chef_can	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
003_cracker_box	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
004_sugar_box	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
005_tomato_soup_can	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
006_mustard_bottle	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
007_tuna_fish_can	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
008_pudding_box	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
009_gelatin_box	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
010_potted_meat_can	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
011_banana	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
012_strawberry	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
013_apple	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
014_lemon	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
015_peach	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
016_pear	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
017_orange	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k
018_plum	Processed (.tgz)	RGB (.tgz)	RGBD (.tgz)	16k	64k	512k

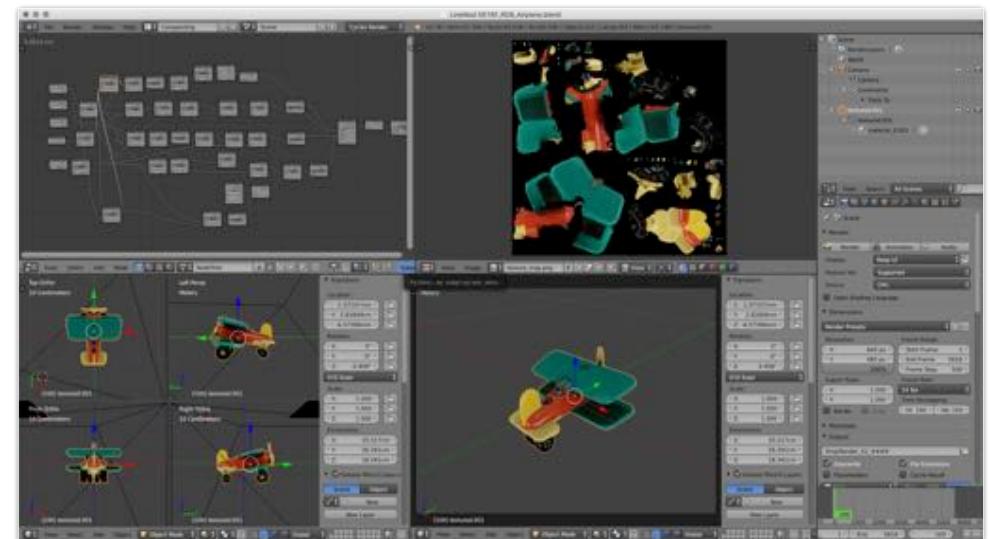
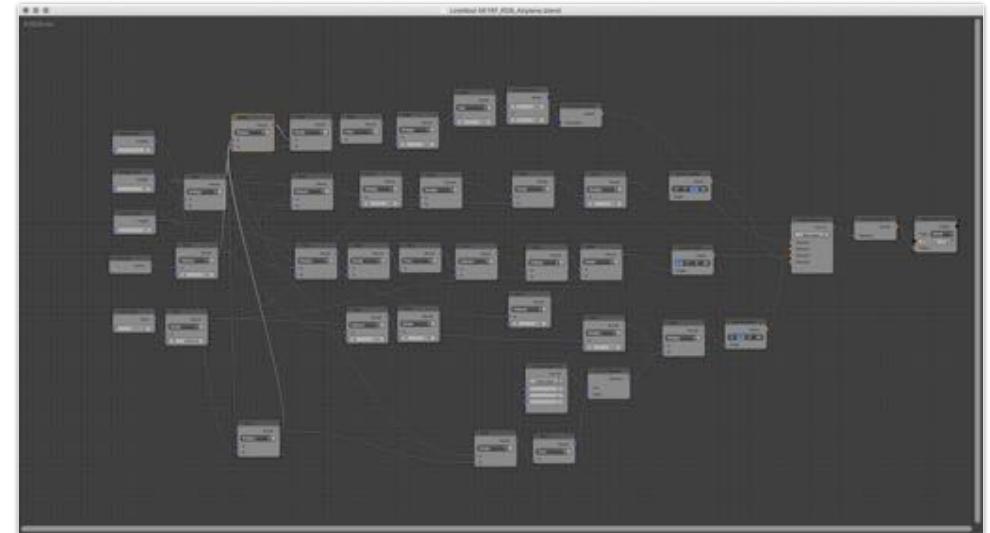
YCB Benchmark Dataset



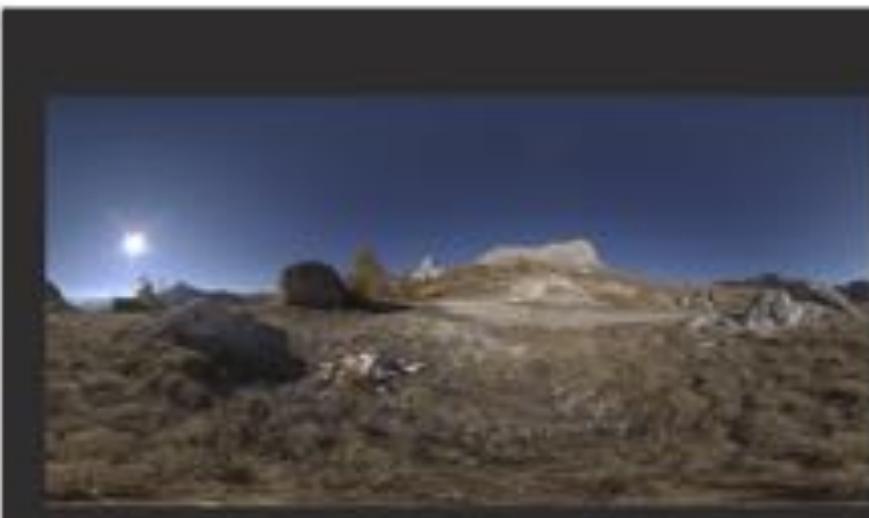
joohaeng@etri.re.kr

Camera

- **5,616 poses**
 - **Pan:** 24 samples in $[0, 2\pi]$
 - **Tilt:** 13 samples in $[-\pi, \pi]$
 - **Yaw:** 3 samples in $[-\pi/4, \pi/4]$
 - **Zoom:** 6 samples in $[0.5, 2]$ m



Lighting





Rendering

- **RGB + XYZ + Depth + Normal** for each camera pose
 - Mask is encoded in alpha channel (A).
 - Other types are rendered using **special shaders**
 - Pixel-wise correspondence between RGB and annotations



RGB (.PNG)



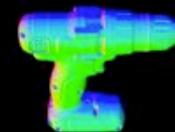
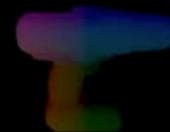
XYZ (.EXR)



Normal (.EXR)



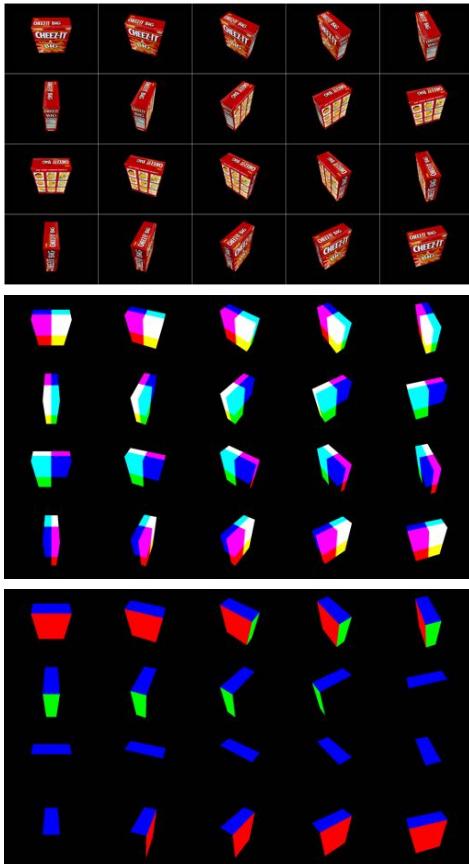
Depth (.EXR)



Application in Pose Estimation

Application in Pose Estimation

Synthetic Learning Set



Feature matching Specialized for PnP



Feature DB with object ID and annotations



Application in Pose Estimation

Synthetic Verification Set



Feature matching

Specialized for PnP



Feature DB with object ID and annotations

Pose estimation Solve PnP

Verification with ref.
 $X = [R|t]$

Application in Pose Estimation

Online Test



Feature DB with object ID and annotations

Feature matching Specialized for PnP



Pose estimation Solve PnP





Kinect v2:
프로그램 화면

Sony Cyber-shot:
외부 카메라

Application in Pose Estimation

ETRI Electronics and Telecommunications Research Institute

Object recognition and pose estimation for modular manipulation system: overview and initial results

Woo-han Yun, Jaeyeon Lee, Joo-Haeng Lee, and Jaehong Kim
HMI Research Group, ETRI

Introduction

- Task adaption of robots becomes an important factor. Modular manipulation system could be a solution.
- In modular manipulation system, recognition engine gives object instance ID, object location and pose information.
- Local feature based approach is selected because of its simplicity and small training dataset requirement.
- We sketch the overview of the process and show initial results using 10 cuboid-type object instances.

Object recognition engine

(Top) 10 cuboid-type objects
(Bottom) Depth-based background subtraction

- Data generation
 - Select 10 cuboid-type objects.
 - Scanning six surfaces and registering them.
 - 2D synthetic images are generated at intervals of 18 degree in yaw, roll, and pitch.
 - On synthetic images, extract local feature, SURF, and store two positional information (2D position in image – 3D position in object-centric coordinate)
- Depth-based background subtraction.
 - To reduce local features in noisy background and extract feature points on the foreground.
 - Depth data in each pixel is modeled by unimodal Gaussian distribution.
 - After subtraction, apply post-processing such as removing small blobs.

Object recognition engine

- Local feature matching
 - Extract SURF feature in query image.
 - Find the correspondence between features in query image and stored database.
 - Calculate Euclidean distance for each pair.
 - Select distinctive correspondence ($D1/D2 < 0.65$ or $D1 < 0.25$).
 - Features having same 3D position are considered as one feature having minimum distance.
- PnP problem solving
 - After matching, distinctive correspondences (2D position in image coor. - 3D position in object-centric coor) are obtained.
 - With these correspondences, find the final object instance and pose estimation by solving PnP problem with RANSAC.

Input query image (left) and its results (right)

Object detection/recognition rate and rotation/translation error

Object	Detection rate(%)	Recognition rate(%)	Rotation error(*)	Translation error(mm)
cheezit	98	100	16.07	2.63
genuine	82	96	0.89	1.20
crayola	98	100	9.15	8.46
mark	82	94	1.76	0.62
expo	92	100	11.55	5.55
ace	98	100	3.07	1.94
chiffon	78	94	5.07	4.81
chococo	96	98	5.17	3.95
moncher	98	100	2.15	3.85
waffle	64	98	4.03	3.96
Average	90.60	98	5.89	3.72

Experimental result

- For train, 10 objects x 231 images were used.
- Object recognition test
 - 400 real images (1 or 2 objects in the shelf) were tested.
- Pose estimation test
 - 1,000 synthetic images are tested (for exact groundtruth).
 - Repeated texture pattern degraded the accuracy.

URAI 2017, Jeju, Korea

Application in Object Detection

- Toward **end-to-end network** for pose estimation **with SLS**
 - No explicit feature extraction, at least.



Ongoing Work — **Complex Environments**

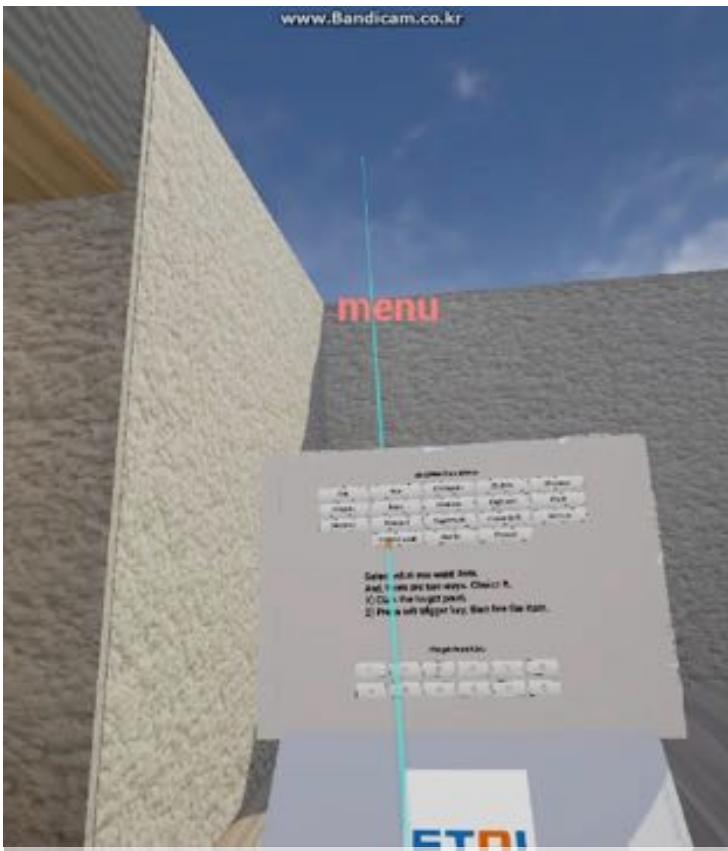


Ongoing Work — VR Simulation (UE4 + HTC Vive)

“I’m your arm.”



Interactive placement of objects



Camera-Projector simulation



Summary

GIVEN:

- **Target objects** for **pose estimation** in robotic manipulation task
- A small set of real images for training

PROBLEM:

- To obtain a **data set** with **quantitative diversity** in *condition* and **qualitative precision** in *annotation*.
- While considering **constraints** in human *labor* and *error*, production *cost*, physical limits in *environments*, and *deploy time*.

PROPOSED SOLUTION:

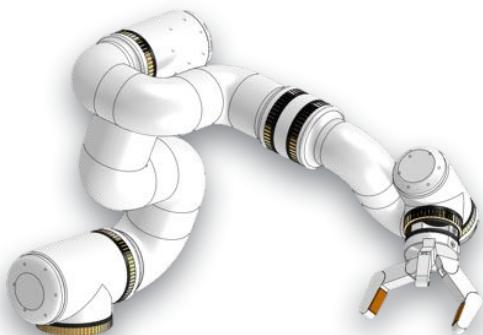
- To synthesize a learning set using **computer graphics**: SLS (Synthetic Learning Set)
- **Diversity**: pose (translation+rotation), illumination, ...
- **Precise** and **rich** annotation: 3D coord, normal, depth, camera extrinsic, ...

RESULTS

- Synthesis pipeline based on open source tools: **Blender**
- **Synthetic Learning Set** (SLS): 516,672 images = 23 objects * 5,616 viewpoints * 4 types (RGB, xyz, normal, depth)
- **Initial experiments** of applying SLS on pose estimation
- Soon open at **github**: <https://github.com/joohaeng/SLS4ML>

Acknowledgement

This research was supported by the MOTIE under the Robot industry fusion core technology development project (10048920) supervised by the KEIT.



Development of **Modular Manipulation System**
Capable of Self-Reconfiguration of Control and
Recognition System for Expansion of Robot
Applicability



Q&A

