

CosyPose: Consistent multi-view multi-object 6D pose estimation

[arXiv:2008.08465](https://arxiv.org/abs/2008.08465)

Yann Labbé^{1,2}, Justin Carpentier^{1,2}, Mathieu Aubry⁴, Josef Sivic^{1,2,3}

¹Inria



²DI ENS, PSL



³CIIRC, CTU in Prague

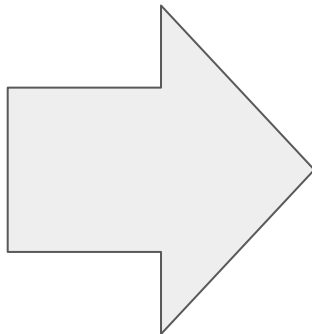
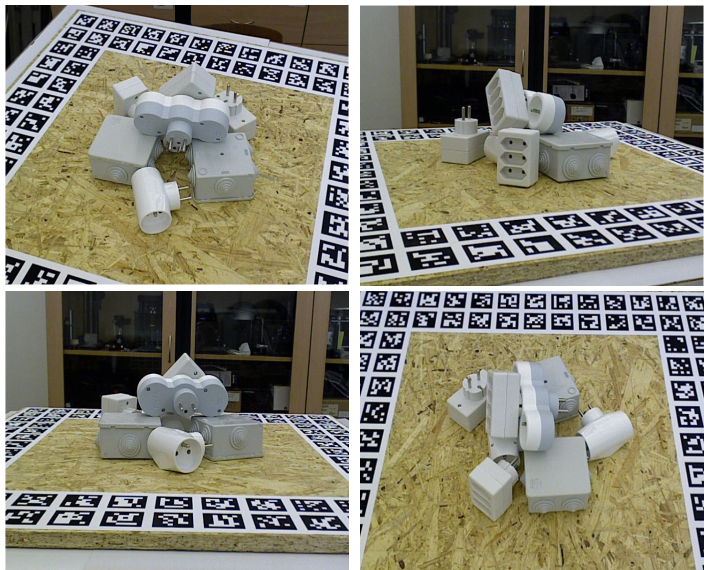


⁴ENPC

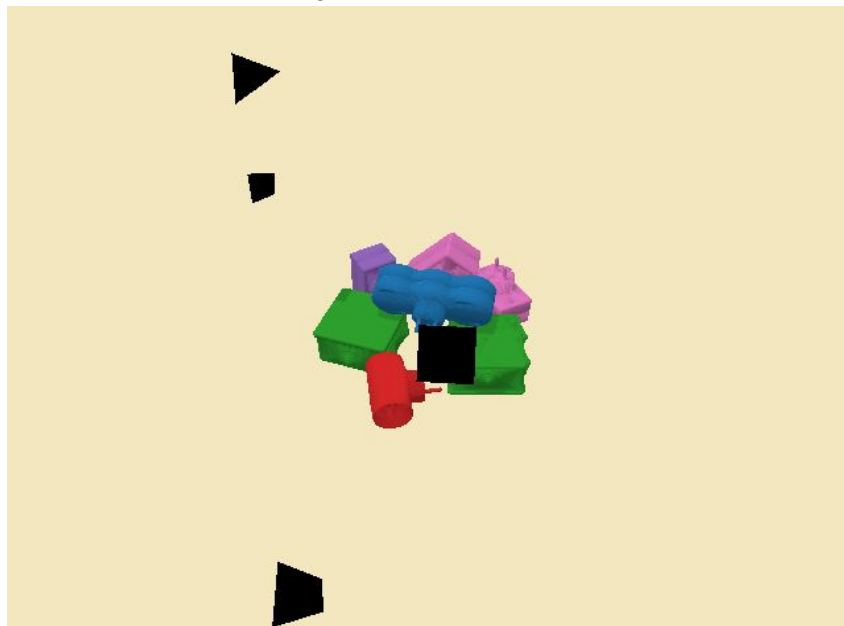


Multi-view 6D pose estimation

Input images



Output 3D scene

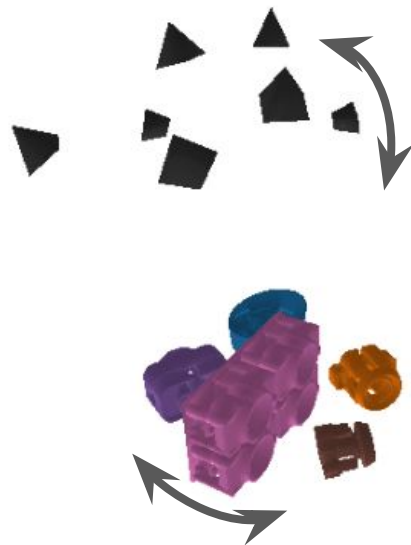
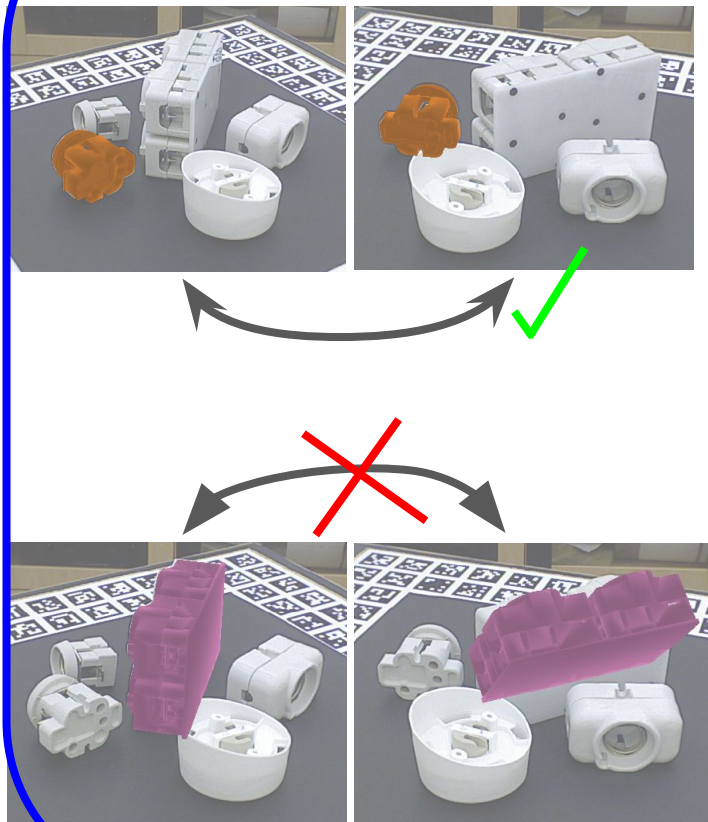


CosyPose: Approach overview

Single-view 6D pose estimation

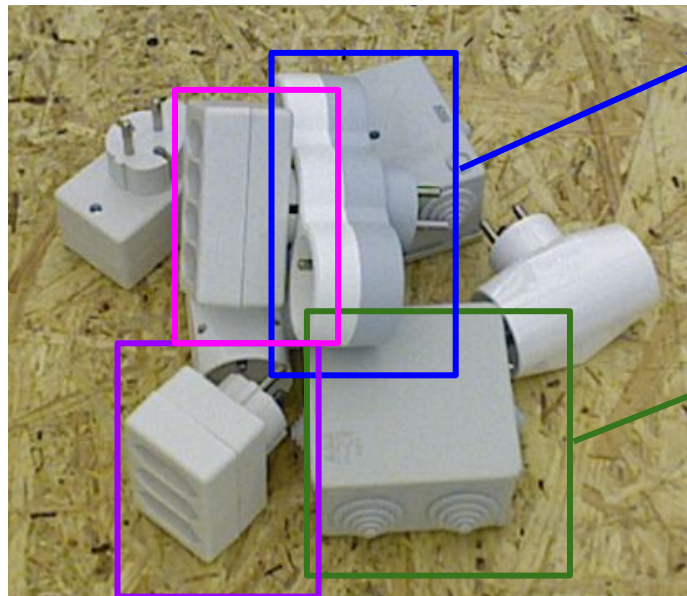


Robust multi-view multi-object reconstruction



...

Single-view CosyPose



Input RGB image

(only 3 networks trained per dataset)

2D detection



6D pose estimation

Coarse network

Refiner network

Coarse network

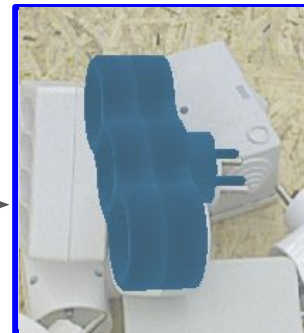
Refiner network

6D pose estimation

Coarse network

Refiner network

6D pose



Pose estimation networks

DeepIM, Li et al, ECCV 2018

- + Network
- + Rotation parametrization
- + Loss
- + Data augmentation

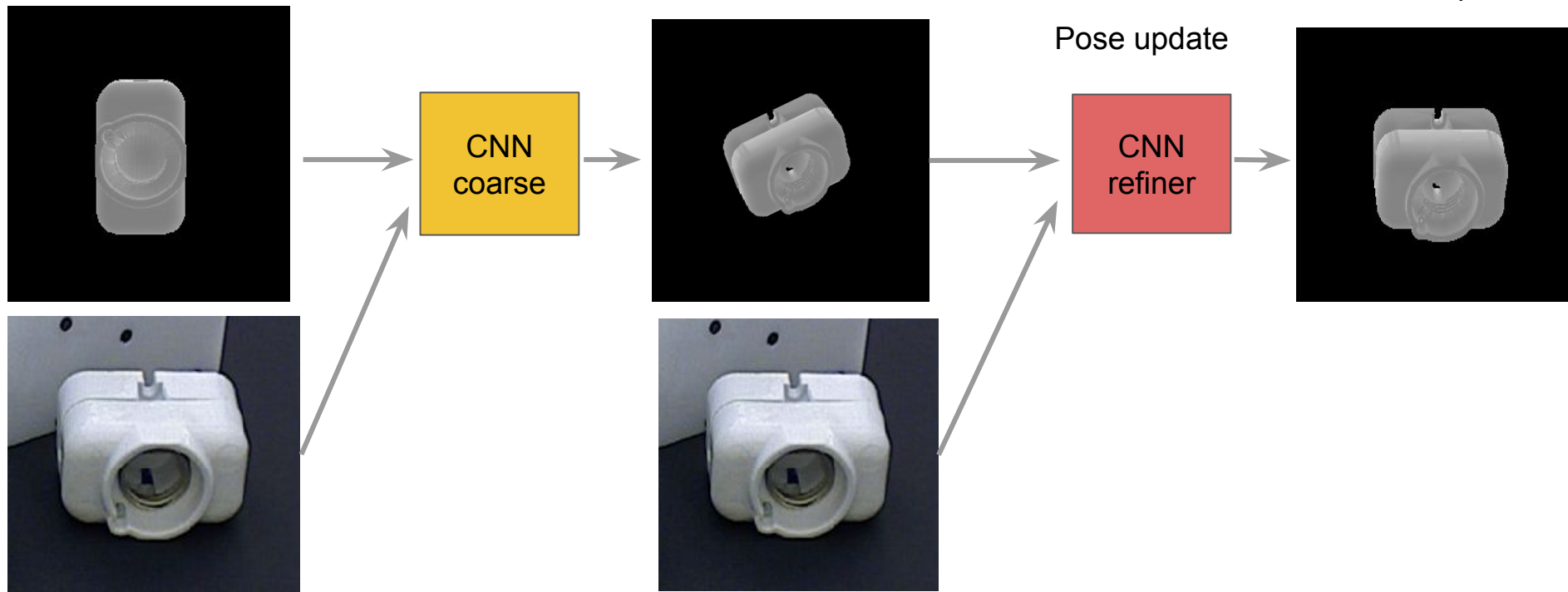


Input “canonical” pose

Input “coarse” pose

(details in the paper [arXiv:2008.08465](https://arxiv.org/abs/2008.08465))

“Refined” pose



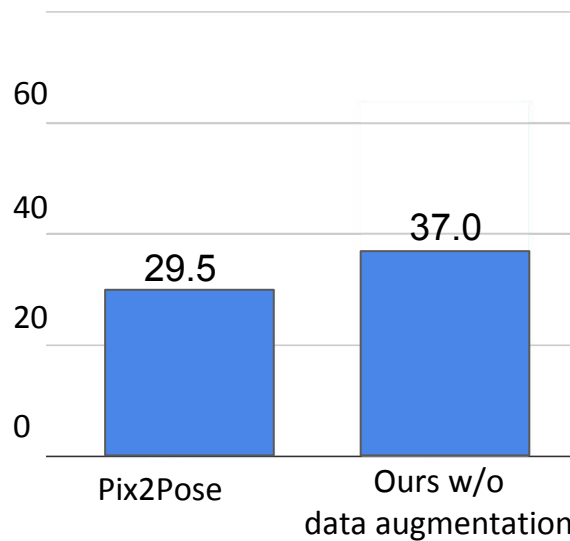
Key ingredients

Without data augmentation



$$e_{\text{vsd}} < 0.3$$

T-LESS

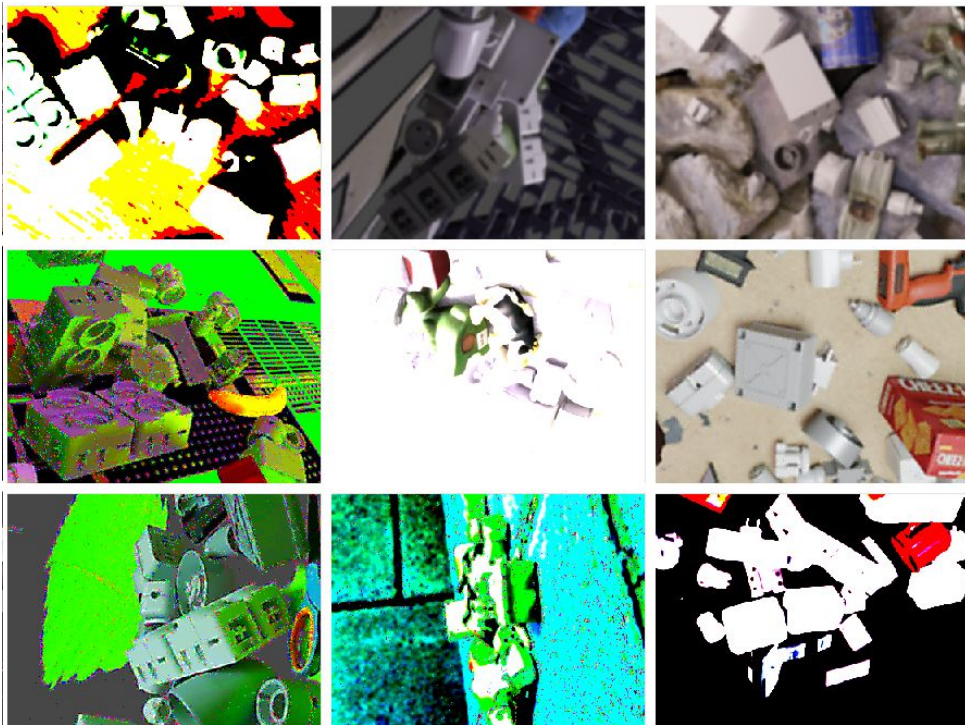


Pix2Pose, Park et al, ICCV 2019

(more ablations in the [paper](#),
Sec 3 Table 1b)

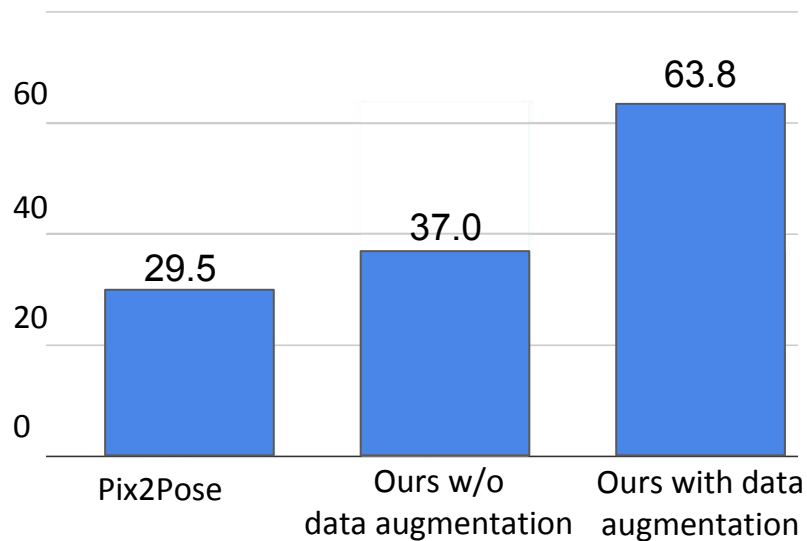
Key ingredients

With data augmentation



$$e_{\text{vsd}} < 0.3$$

T-LESS



Pix2Pose, Park et al, ICCV 2019

(more ablations in the [paper](#),
Sec 3 Table 1b)

+ Access to a GPU cluster*
training 1 pose network: **~10 hours on 32 GPUs**

*Jean-zay, French national cluster managed by GENCI-IDRIS

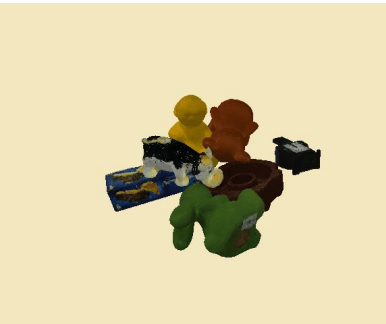
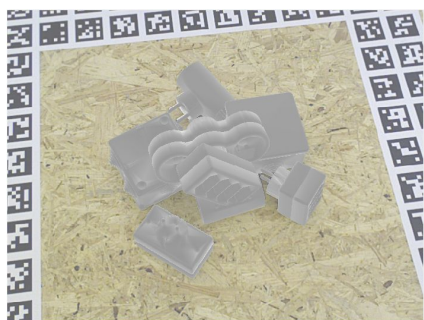
Input image



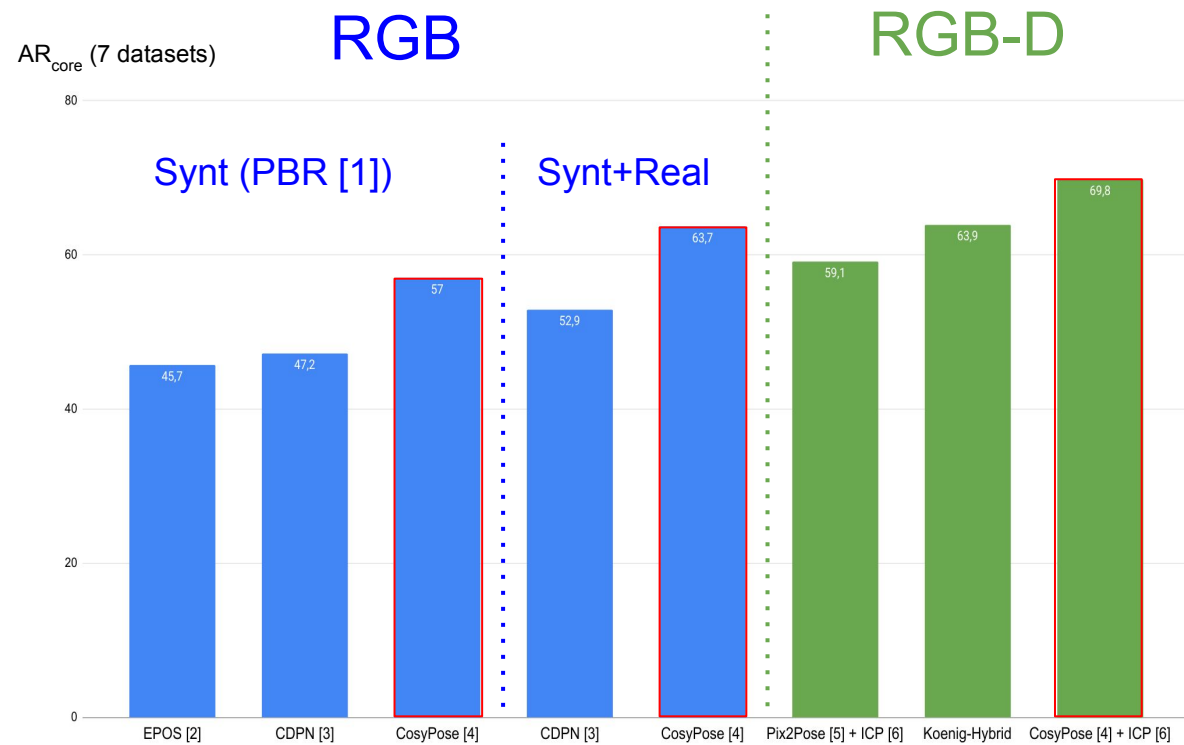
Predicted poses



3D visualization



BOP20 results



[1] BlenderProc: Denninger, Sundermeyer, Winkelbauer, Olefir, Hodan, Zidan, Elbadrawy, Knauer, Katam, Lodhi in RSS workshops.

[2] EPOS, Hodan et al, CVPR 2020

[3] CDPN, Li et al, ICCV 2019

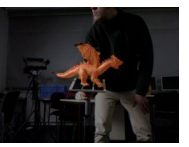
[4] **CosyPose, Labbé et al, ECCV 2020**

[5] Pix2Pose, Park et al, ICCV 2019

[6] <https://github.com/kirumang/Pix2Pose>

+ running time < 0.5s per image

Code



- State-of-the-art pre-trained models for multiple datasets
- RGB single-view and multi-view modular framework
- Full training code

```
# One or multiple RGB images
images = ...

# Define camera intrinsic parameters
intrinsics = ...

# 2D detections
detections = detector.get_detections(images)

# Single-view pose estimation
pose_predictions = pose_predictor.get_predictions(images, intrinsics, detections, n_refiner_iterations=4)

# Object-level multi-view reconstruction (if multiple images available)
scene_state = multiview_predictor.predict_scene_state(pose_predictions, intrinsics)
```



<https://github.com/ylabbe/cosypose>

CosyPose: Consistent multi-view multi-object 6D pose estimation

[arXiv:2008.08465](#)

Yann Labbé^{1,2}, Justin Carpentier^{1,2}, Mathieu Aubry⁴, Josef Sivic^{1,2,3}

¹Inria

²DI ENS, PSL

³CIIRC, CTU in Prague

⁴ENPC



<https://github.com/ylabbe/cosypose>