# Real-time 6 DOF Pose Estimation with Limited Priors

Student: Dhruv Sheth

Collaborators: Aurelio Noca, Jonathan Becktor

Co Mentor: Dr. Ersin Das

Mentor: Prof. Joel Burdick

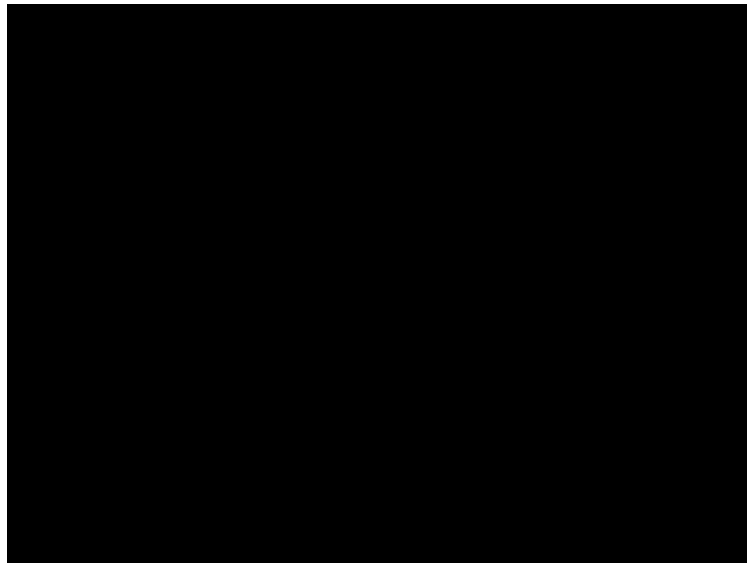Caltech

# Agenda

- Background and Motivation
- Objectives
- Methods
    - Simulation
    - Perception
- Results and Discussion
- Future Work
- Acknowledgements
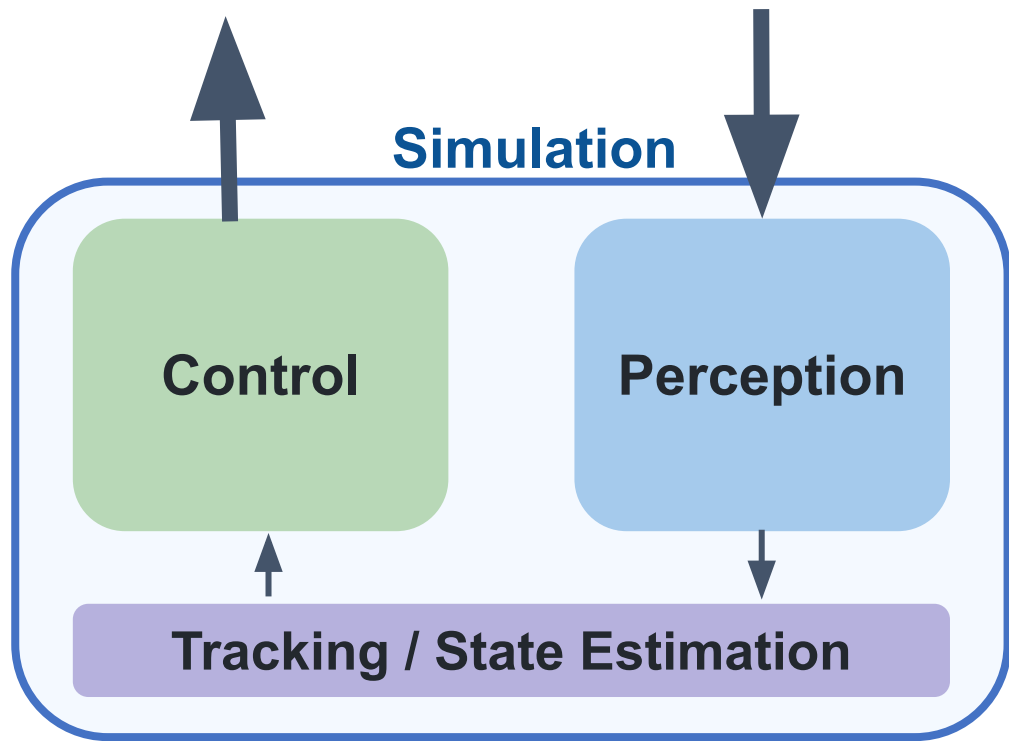- Q&A

Caltech

# Background

**DARPA LINC Phase 1**

- Precise payload placement on ships
- Respond to events not predicted at design time
- Robustness to:
  - uncertain state estimates
  - perception error in dynamic environments
- Design ML models with **safety guarantees**

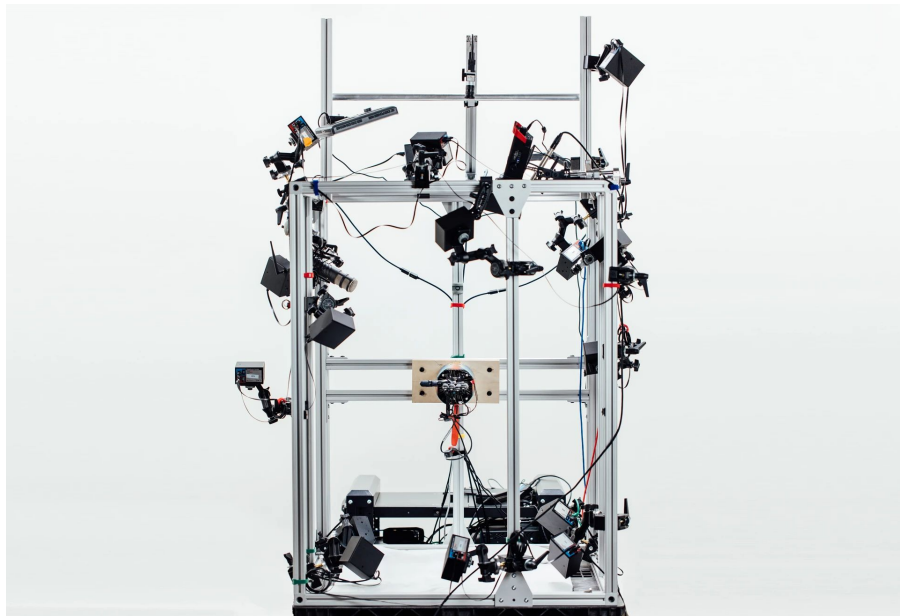**tldr; this is a complex problem**
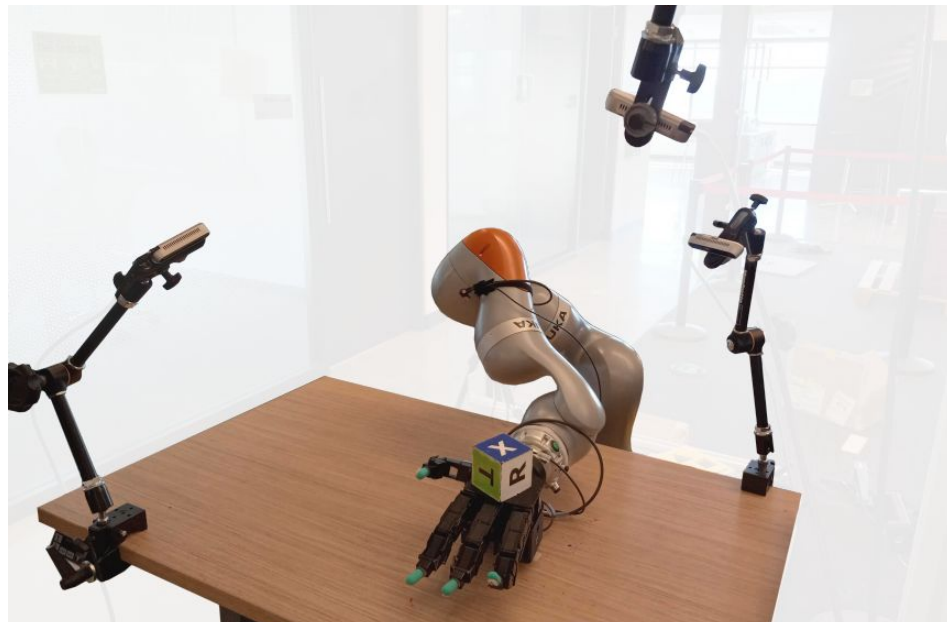


from Burdick Lab slack

# Project Objectives

- Implement a robust 6 DoF pose estimation method for objects with unknown geometry
    - Robust: Make it recover from errors in complete occlusion
    - Optimize: Make it run real time on AGX Orin with low compute
- Develop a high-fidelity simulation stack for testing control+perception
- Extend the algorithm to multiple cameras
- Create a ROS2 wrapper for plug and play support

Caltech

# Other Applications



*Dactyl lab setup with Shadow Dexterous Hand, PhaseSpace motion tracking cameras, and Basler RGB cameras*
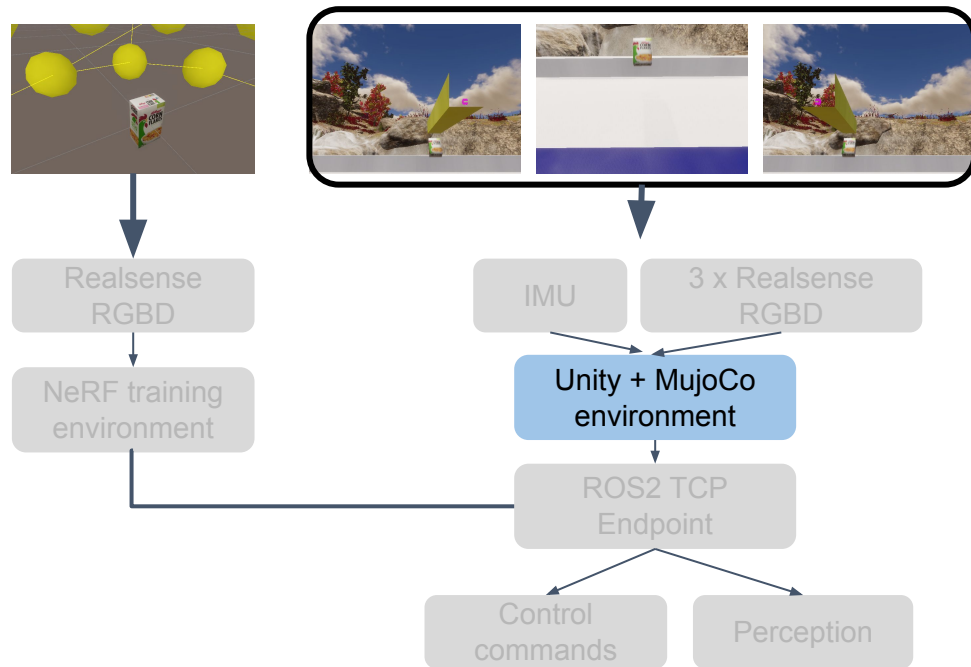OpenAI Dactyl Experiment, 2018



*DeXtreme: Transfer of Agile In-Hand Manipulation from Simulation to Reality, 2024*
OpenAI hand replication experiment

Caltech
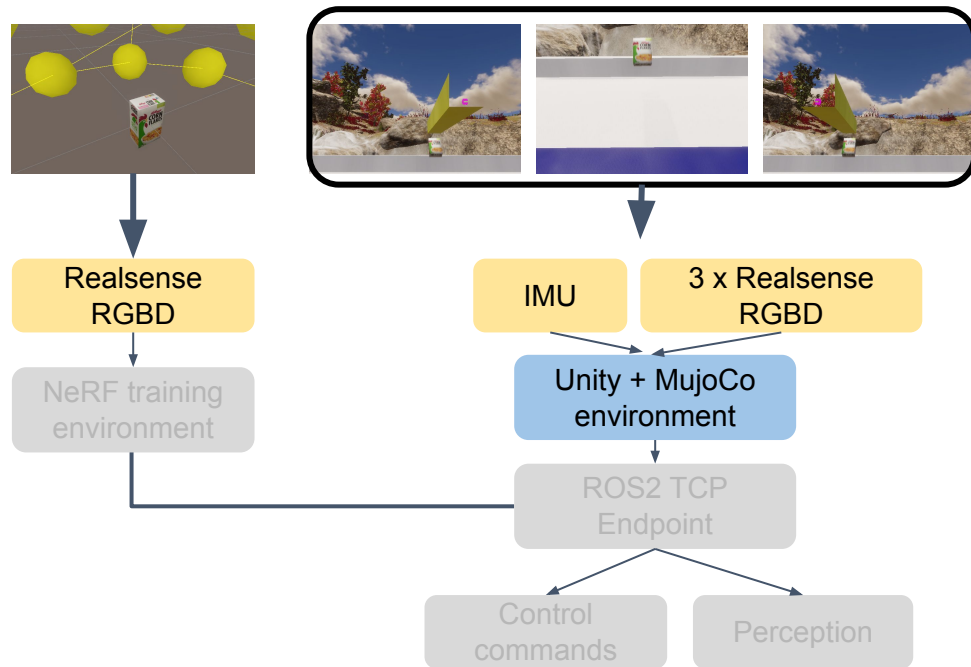
# Methods

**Caltech**

# Simulation Stack

- High Fidelity environment replicating testbed
- Unity + MuJoCo integration
  - MuJoCo physics: control
  - Unity photorealism: perception



Realsense RGBD

NeRF training environment

IMU

3 x Realsense RGBD

Unity + MujoCo environment

ROS2 TCP Endpoint

Control commands
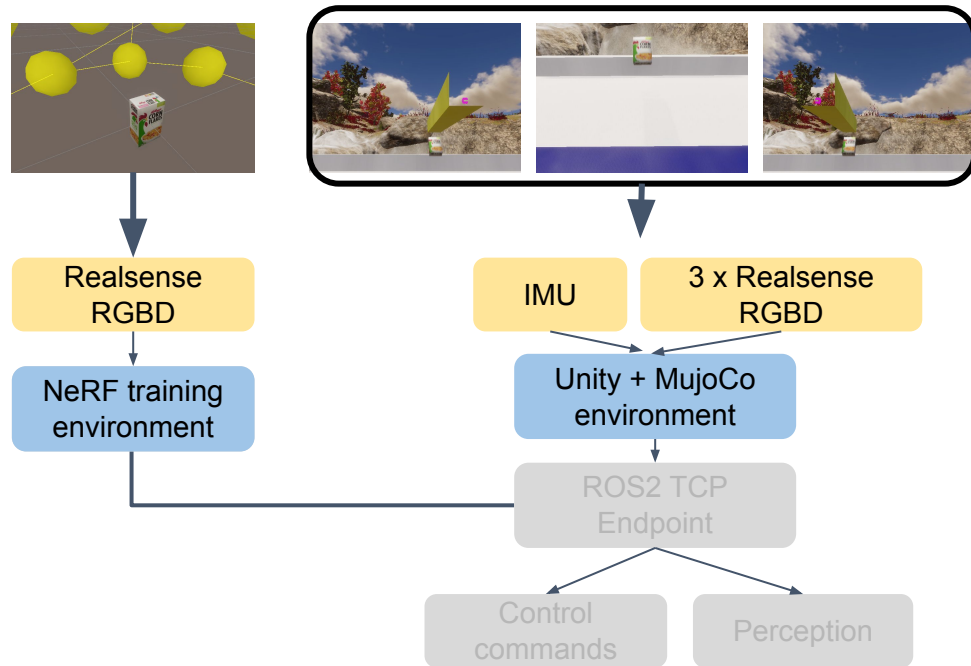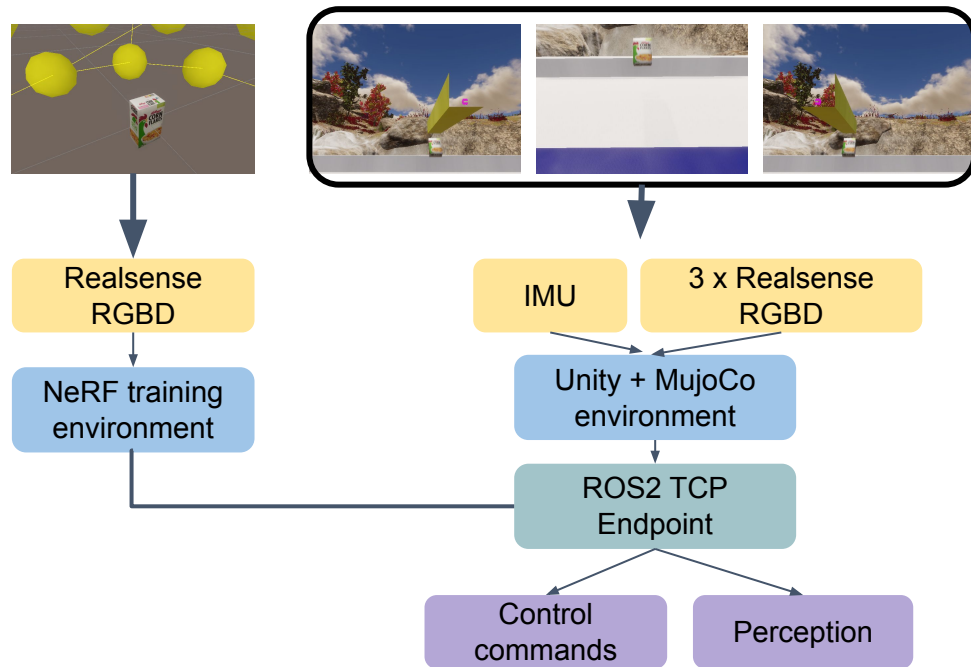
Perception

**Caltech**

# Simulation Stack

- High Fidelity environment replicating testbed
- Unity + MuJoCo integration
  - MuJoCo physics: control
  - Unity photorealism: perception
- RGBD + IMU sensor data integration for testing



Realsense RGBD

IMU

3 x Realsense RGBD

NeRF training environment

Unity + MujoCo environment

ROS2 TCP Endpoint

Control commands

Perception
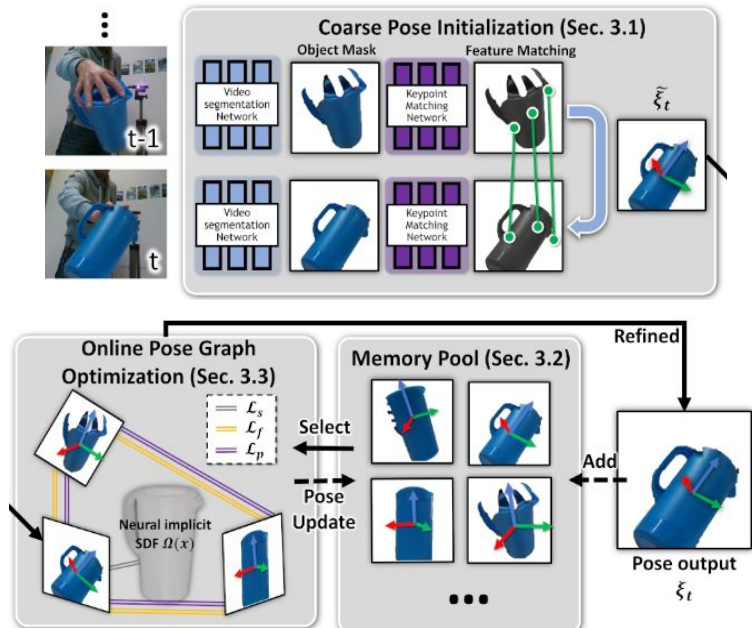
**Caltech**

# Simulation Stack

- High Fidelity environment replicating testbed
- Unity + MuJoCo integration
  - MuJoCo physics: control
  - Unity photorealism: perception
- RGBD + IMU sensor data integration for testing
- NeRF training environment



**Caltech**

# Simulation Stack

- High Fidelity environment replicating testbed
- Unity + MuJoCo integration
  - MuJoCo physics: control
  - Unity photorealism: perception
- RGBD + IMU sensor data integration for testing
- NeRF training environment
- ROS2 wrapper
  - communication with sim
  - swapping with hardware



**Caltech**

# Perception Stack

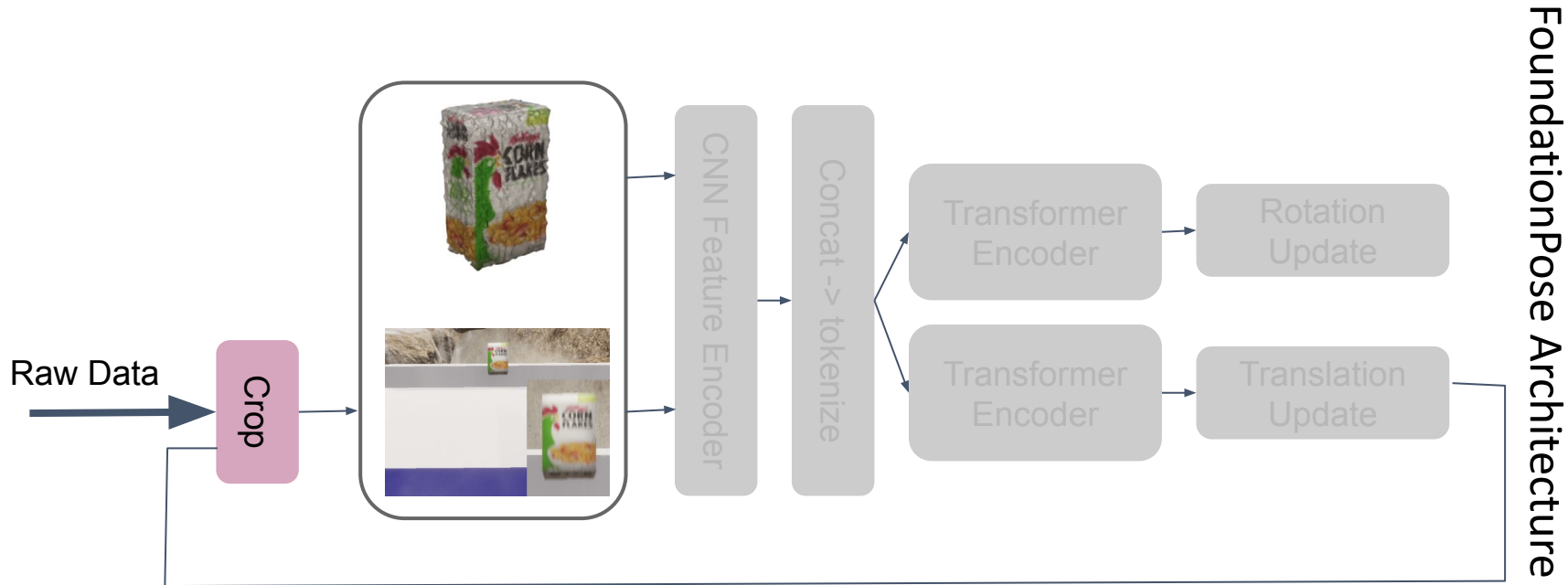**Previous Methods**

- Feature Tracking approaches
  - Classical Approaches:
    - Inverse SLAM
    - LoFTR feature tracking
  - Deep Learning:
    - BundleSDF
- Transformer approaches:
  - FoundationPose
    - Not suited for longer videos



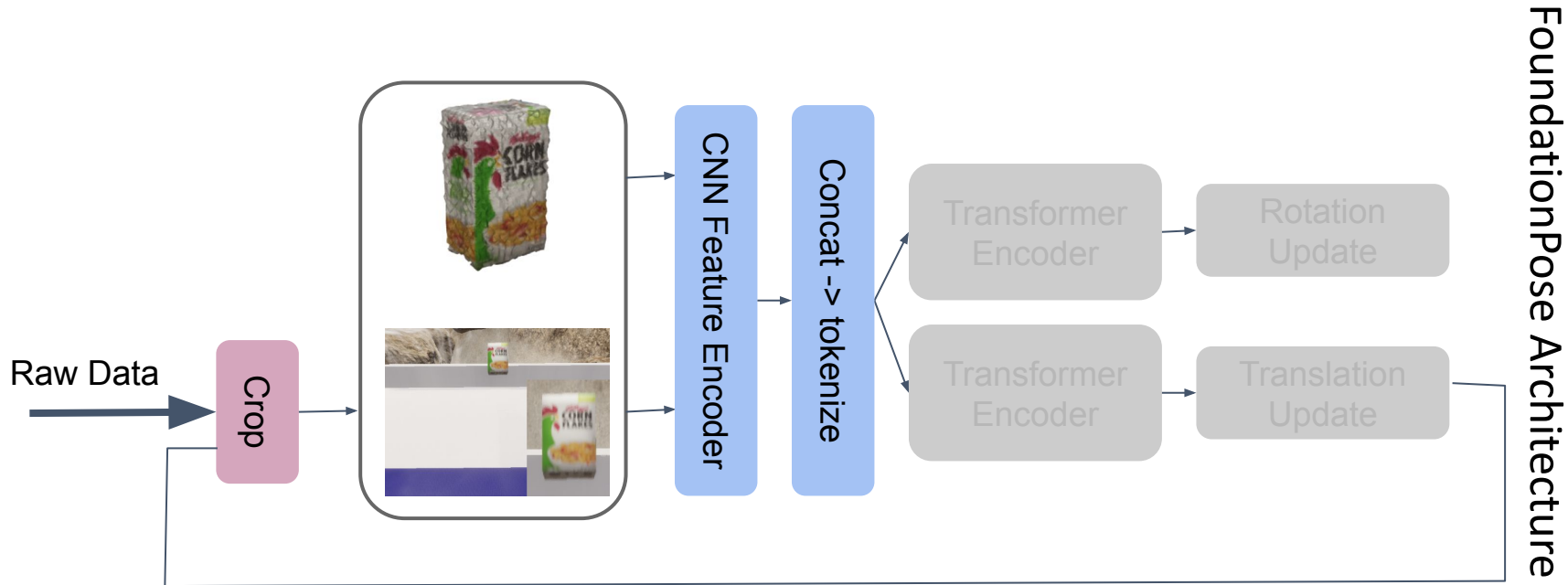*BundleSDF: Neural 6-DoF Tracking and 3D Reconstruction of Unknown Objects*
NVIDIA, 2023
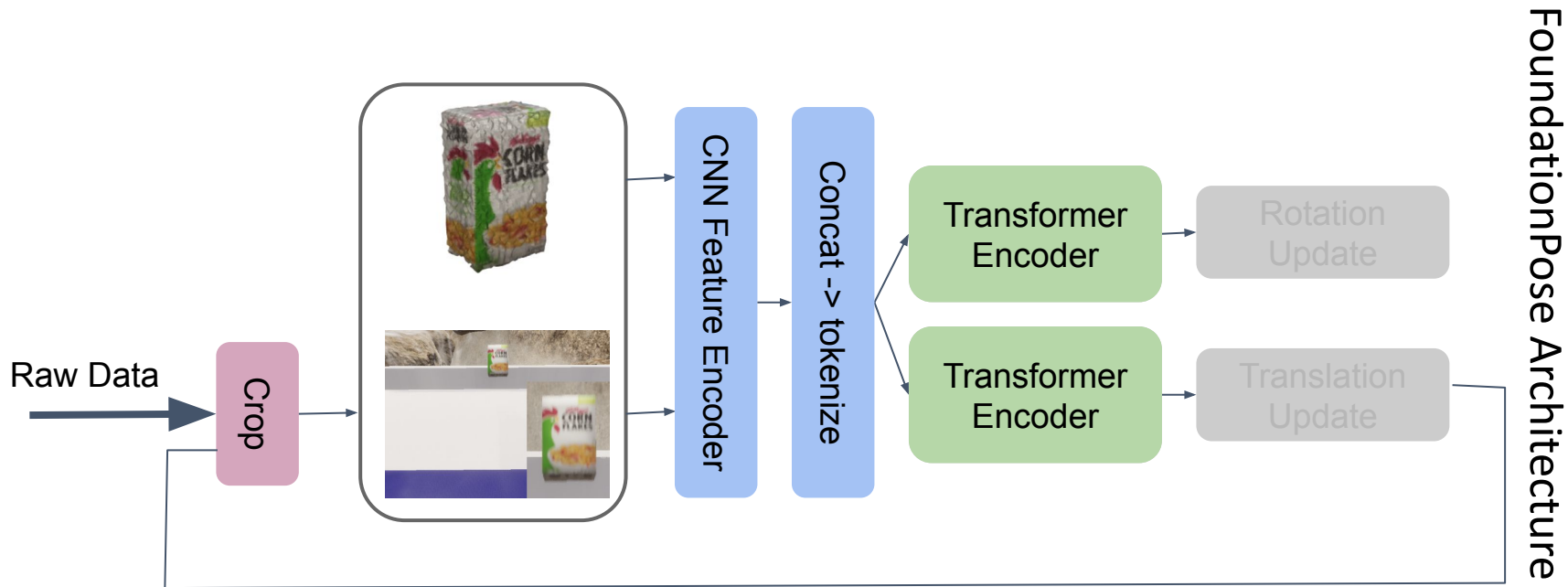
**Caltech**

# Perception Stack

**FoundationPose**

# Perception Stack

**FoundationPose**



Raw Data

Crop

CNN Feature Encoder

Concat -> tokenize

Transformer Encoder

Rotation Update

Transformer Encoder

Translation Update

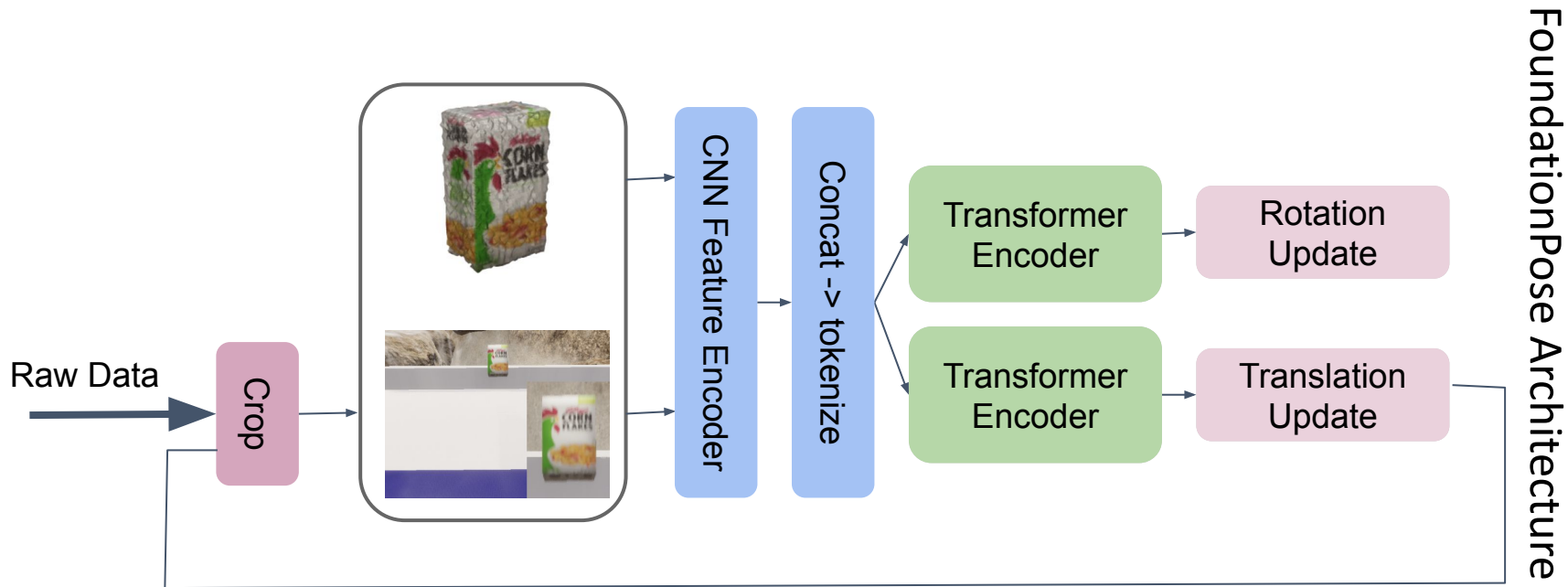FoundationPose Architecture

**Caltech**

# Perception Stack

**FoundationPose**

# Perception Stack

**FoundationPose**

# Perception Stack

**Video Object Segmentation**

- Cutie: Putting the Object Back Into Video Object Segmentation
- Why segmentation?
    - 2D data is easier than 3D data
    - Current 2D models outperform 3D
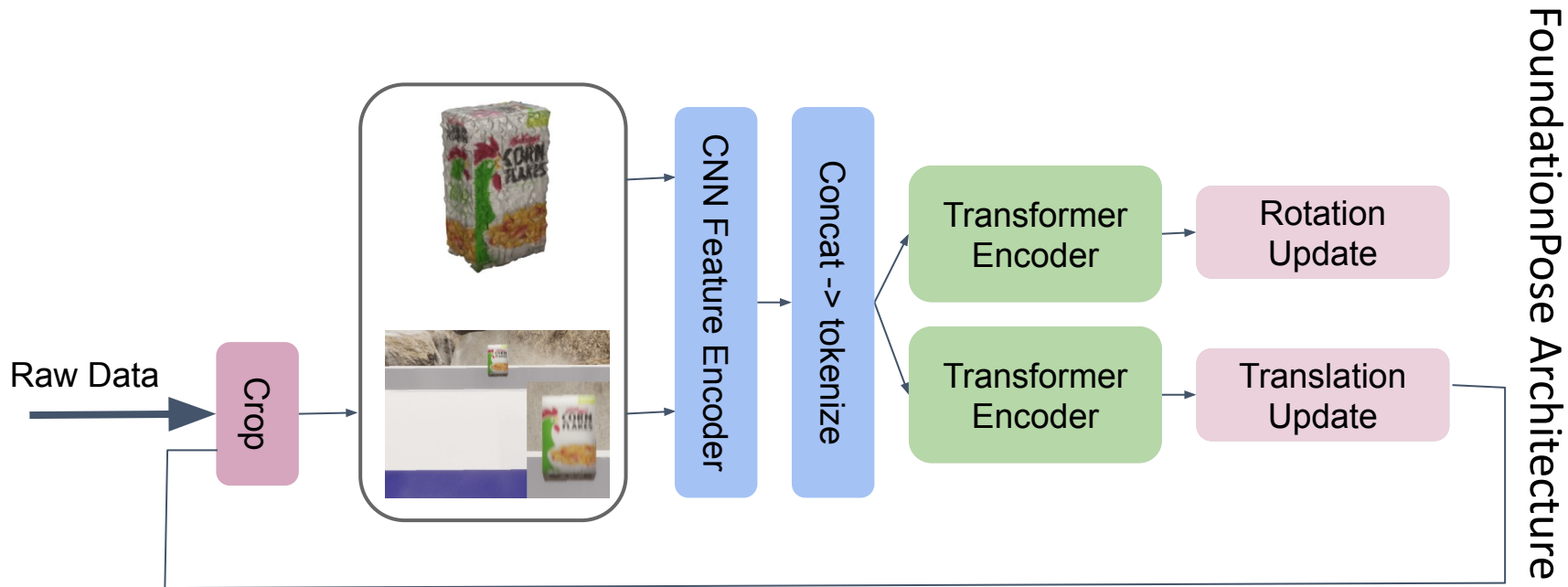    - Provides a good coarse input

# Perception Stack

**Video Object Segmentation**

- Cutie: Putting the Object Back Into Video Object Segmentation
- Why segmentation?
    - 2D data is easier than 3D data
    - Current 2D models outperform 3D
    - Provides a good coarse input
- Why Cutie?
    - Memory model
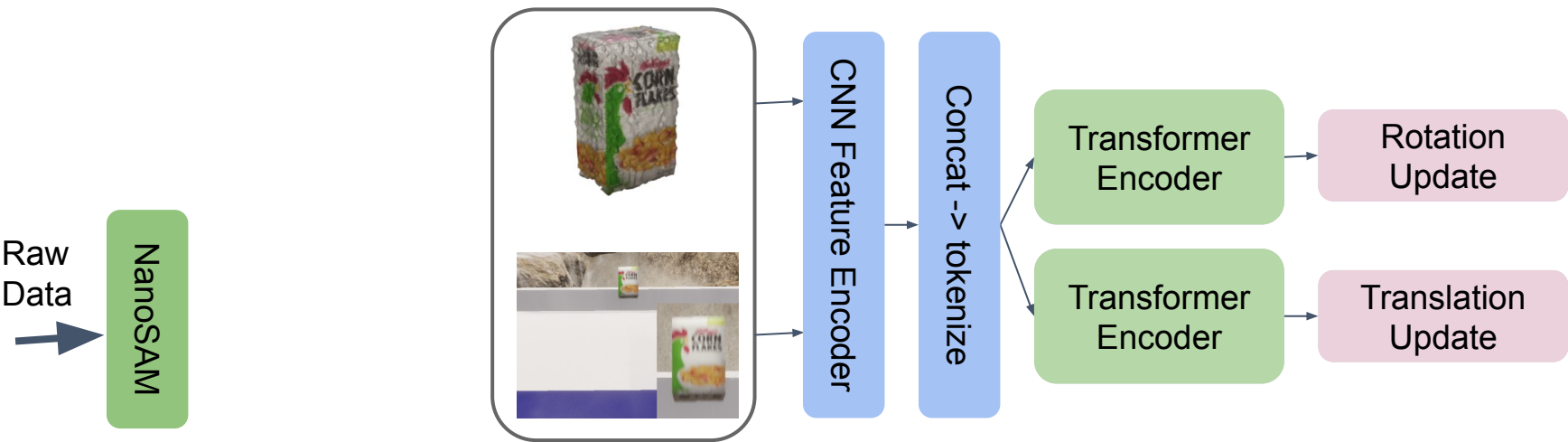    - Long horizon segmentation
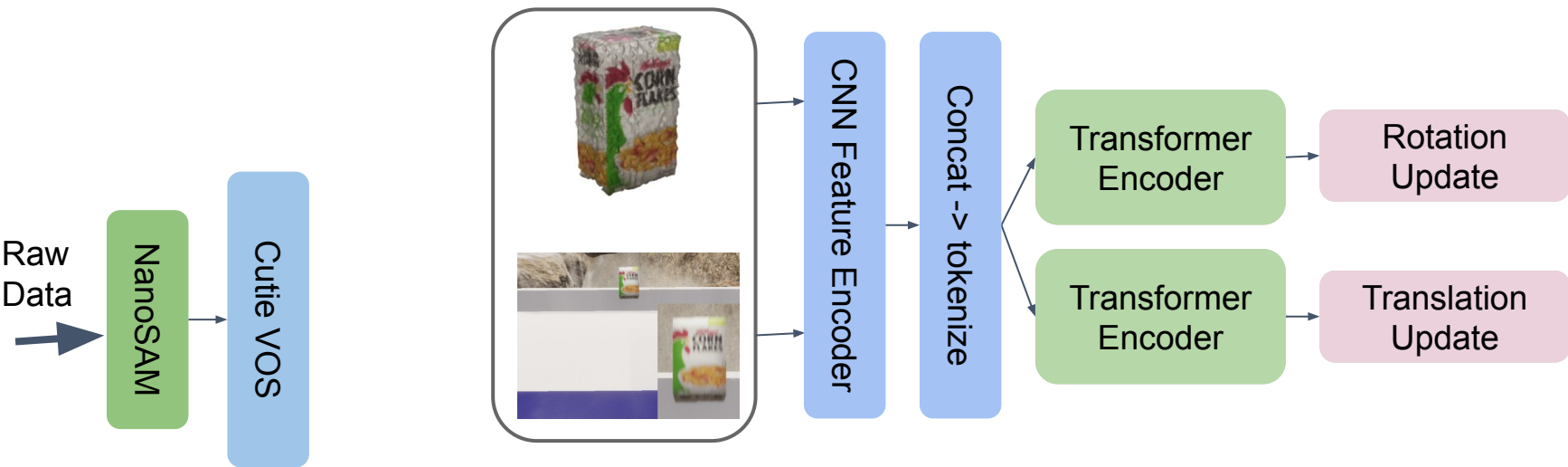    - Recovers from occlusions

**Caltech**

# Perception Stack

**FoundationPose**

# Perception Stack
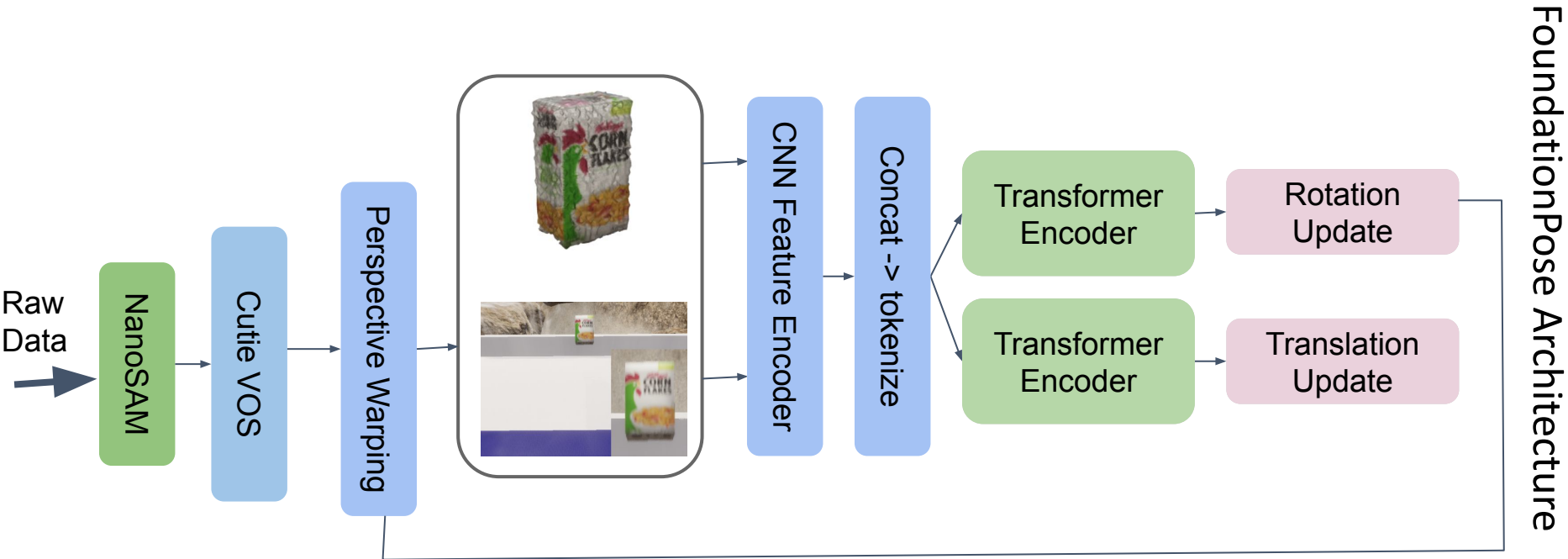
**FoundationPose Improvement**

# Perception Stack

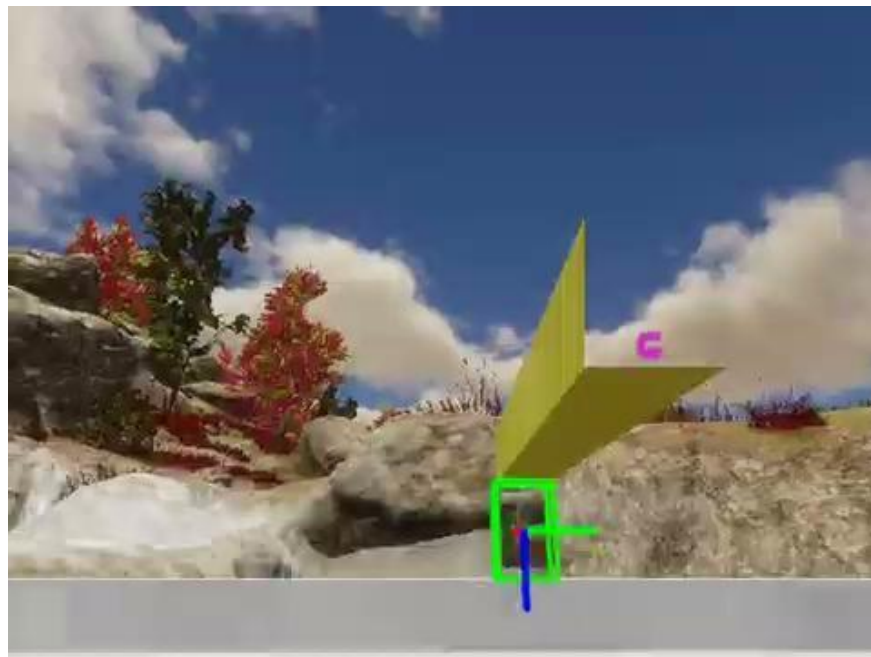**FoundationPose Improvement**
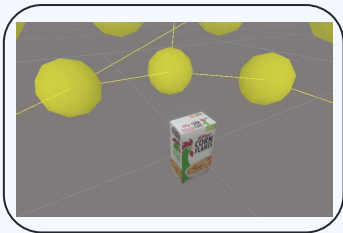
# Perception Stack

**FoundationPose Improvement**

# Perception Stack

FoundationPose + VOS results

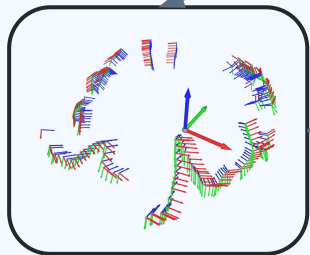- High Disturbance
- Low visibility of object



**Caltech**

# Perception Stack

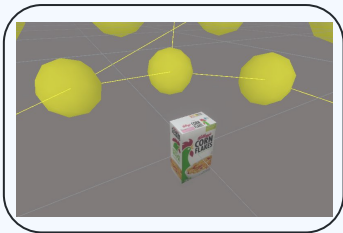**FoundationPose pipeline**
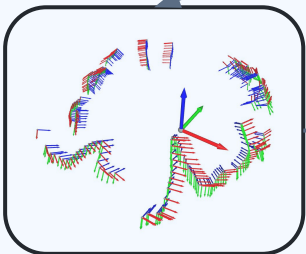


NeRF 3D data in sim



COLMAP camera pose



NeRF 3D reconstruction

**priors**

Caltech

# Perception Stack

**FoundationPose pipeline**



NeRF 3D data in sim

COLMAP camera pose

NeRF 3D reconstruction

**priors**

FoundationPose Improvement

pose estimate

# Perception Stack

**Multi-Cam integration & deployment**

- Under complete occlusion in one camera:
  - Switch between cameras with priority ranking
  - mTC to know when object recovered from occlusion
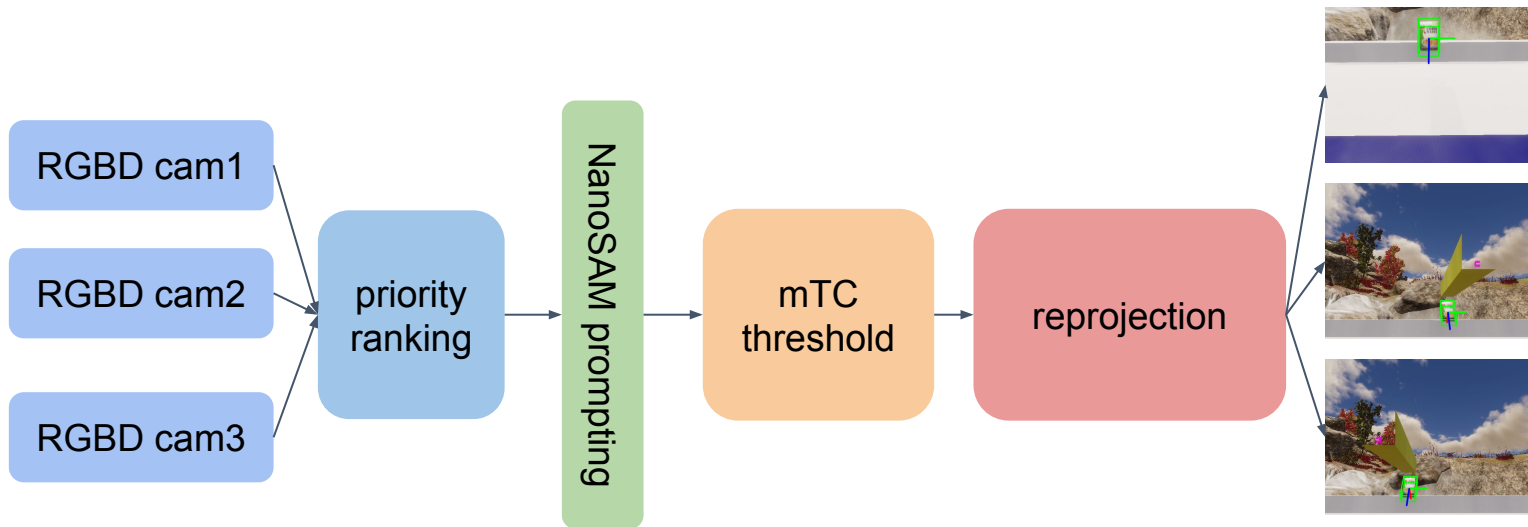
# Perception Stack

**Multi-Cam integration & deployment**

- Under complete occlusion in one camera:
    - Switch between cameras with priority ranking
    - mTC to know when object recovered from occlusion
- Implement ROS2 interface for:
    - Multicam data
    - Pose tracking
    - Simulation to perception

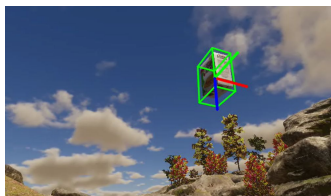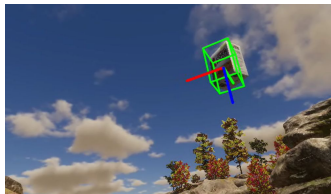**Caltech**

# Perception Stack

**Multi-Cam integration & deployment**

- Under complete occlusion in one camera:
    - Switch between cameras with priority ranking
    - mTC to know when object recovered from occlusion
- Implement ROS2 interface for:
    - Multicam data
    - Pose tracking
    - Simulation to perception
- TensorRT Optimization on VOS
    - 20% improvement
        - 3 modules

**Caltech**

# Perception Stack

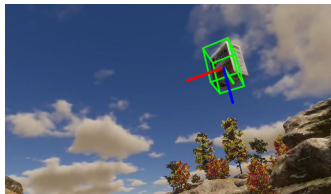**Safety-Critical use**

Which pose estimate do we trust more?
How we do know?





Caltech

# Perception Stack

## Safety-Critical use

Which pose estimate do we trust more?
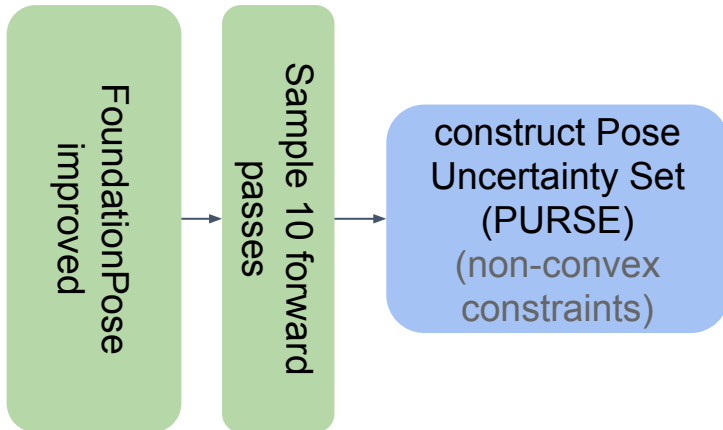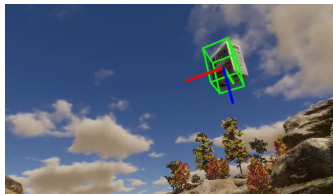How we do know?



FoundationPose improved → Sample 10 forward passes

Caltech

# Perception Stack

**Safety-Critical use**

Which pose estimate do we trust more?
How we do know?



FoundationPose improved

Sample 10 forward passes

construct Pose
Uncertainty Set
(PURSE)
(non-convex
constraints)

**Caltech**

# Perception Stack

**Safety-Critical use**

Which pose estimate do we trust more?
How we do know?



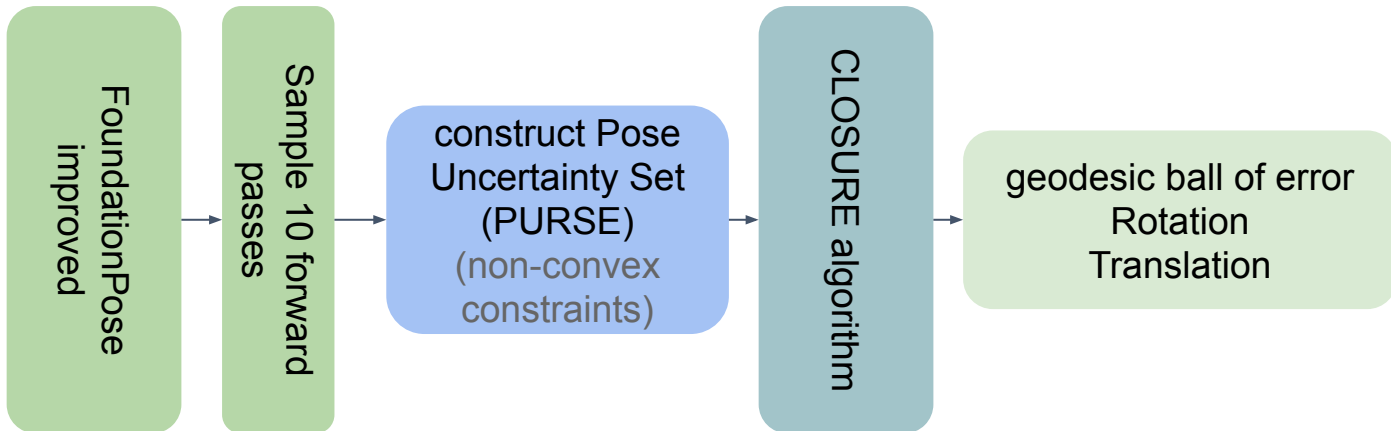FoundationPose improved → Sample 10 forward passes → construct Pose Uncertainty Set (PURSE) *(non-convex constraints)* → CLOSURE algorithm → geodesic ball of error Rotation Translation
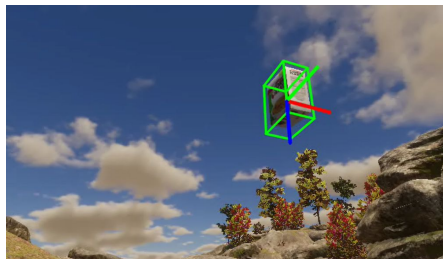
**Caltech**

# Perception Stack
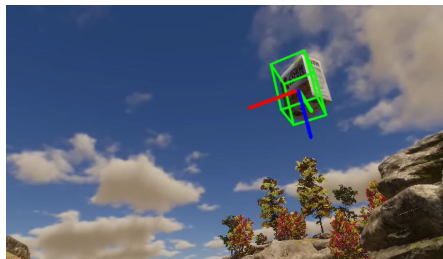
**Safety-Critical use**
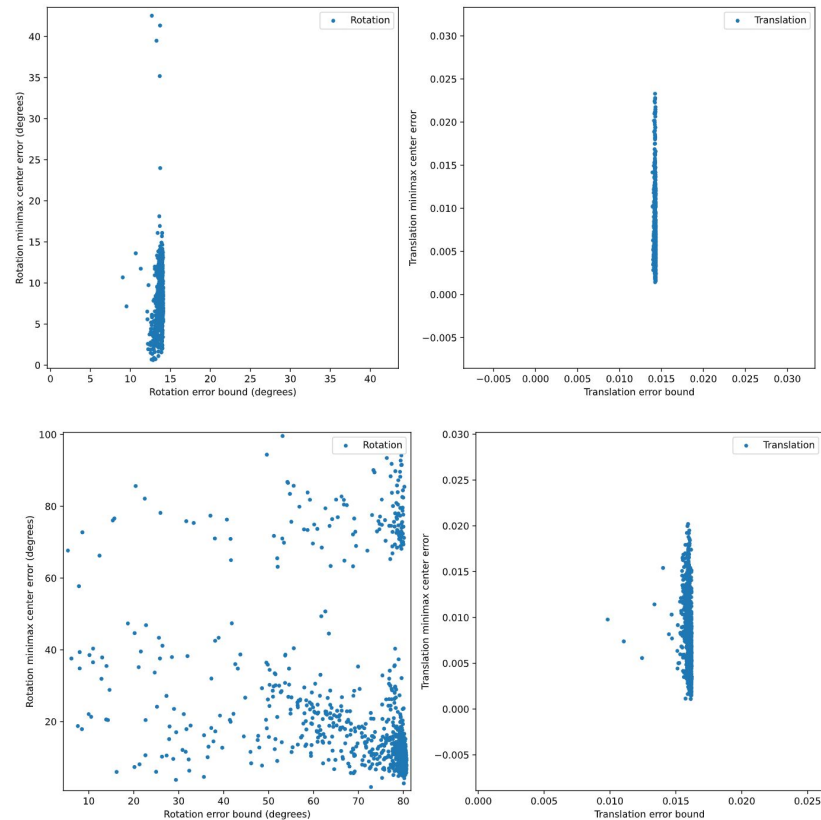
Results from **CLOSURE**



complete NeRF
reconstruction



partial NeRF
reconstruction

Actual Error

Predicted Error Bound

# Results and Conclusion

- FoundationPose + VOS pose estimation
- Multi-Cam integration to FoundationPose
- ROS2 interface for perception
- High-fidelity simulation environment for testing
- Optimization for low compute
- Uncertainty Quantification for safety

**Caltech**

# Future Work

- More Benchmarks
- Use less priors (No 3D reconstruction of object)
- Integrate a state estimator for filtering

**Caltech**

# Acknowledgements

- Prof. Joel Burdick for giving me the opportunity to SURF with the lab, resources and guidance along the way
- Dr. Ersin Das for his guidance and support and help whenever needed
- Aurelio, Jonathan and Thomas for help with code, hardware and debugging
- Dr. Jane Chen for SURF fellowship to support my research this summer

Caltech

# Q&A

Caltech

caltech.edu