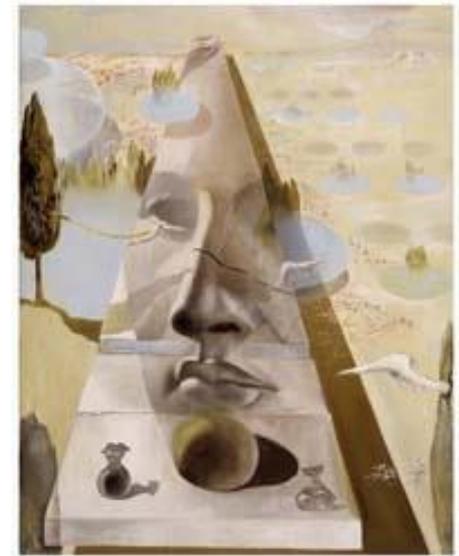


# CS231A

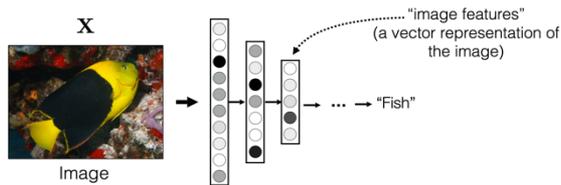
## Computer Vision: From 3D Reconstruction to Recognition

Representation Learning for Finding  
Correspondences and Depth Estimation

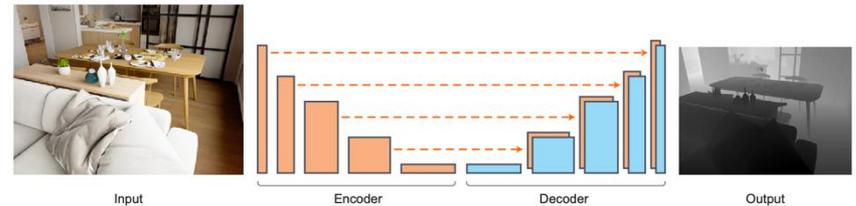


# Learning Goals for Upcoming Lectures

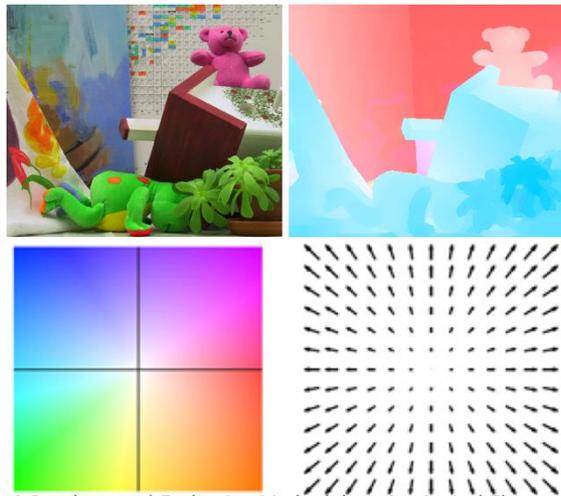
## Representations & Representation Learning



## Using Representation Learning for Depth Estimation and Finding Correspondences

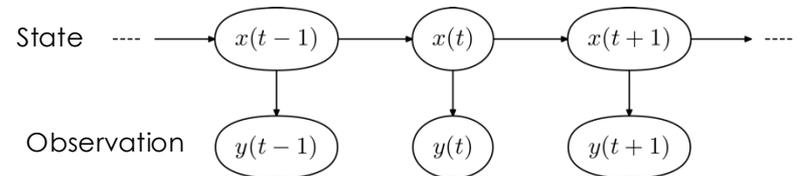


## Optical & Scene Flow

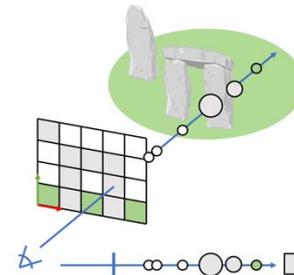


A Database and Evaluation Methodology for Optical Flow.  
Baker et al. IJCV. 2011

## Optimal Estimation



## Neural Radiance Fields



# Outline of the previous lecture

- What is a state? What is a representation?
- What are the different kinds of representations?
- How can we extract state from raw sensory data?
- How can we learn good representations from data?

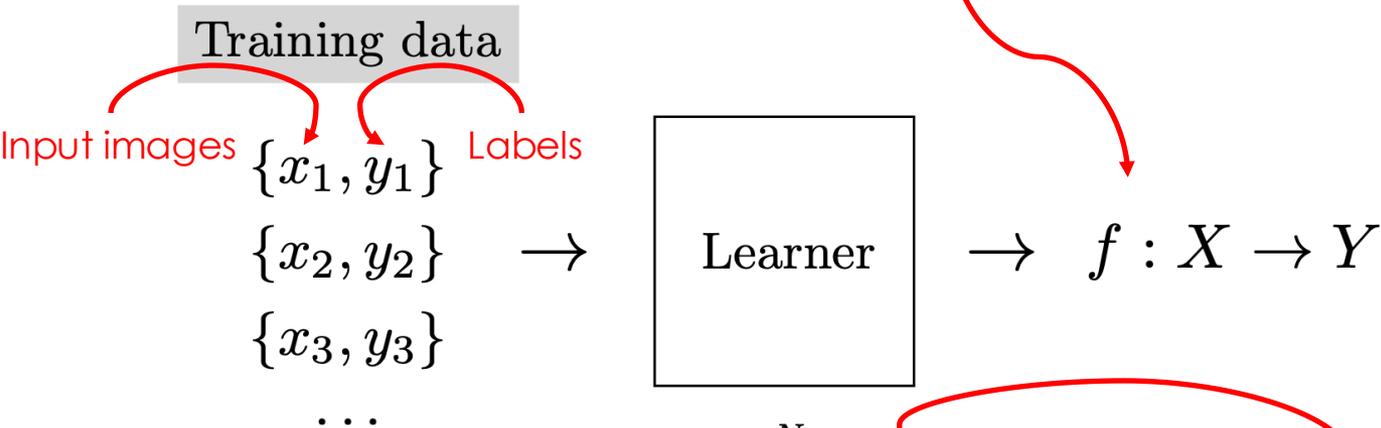
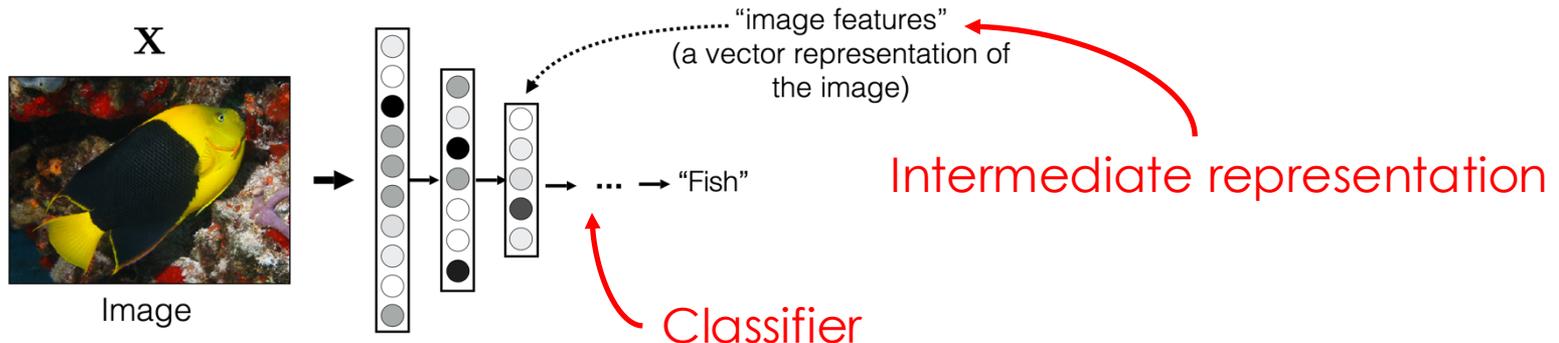
# Summary of what you learned

- **State:** Quantity that describes the most important aspect of a dynamical system at time  $t$
- **Representation:** data format of input or output including a low-dimensional representation of sensor data
  - Input/output/intermediate representation

# Summary of what you learned

- Learned versus interpretable representations
- Visualize learned representations
- How to learn representations?
  - Supervised
  - Unsupervised
  - Self-supervised

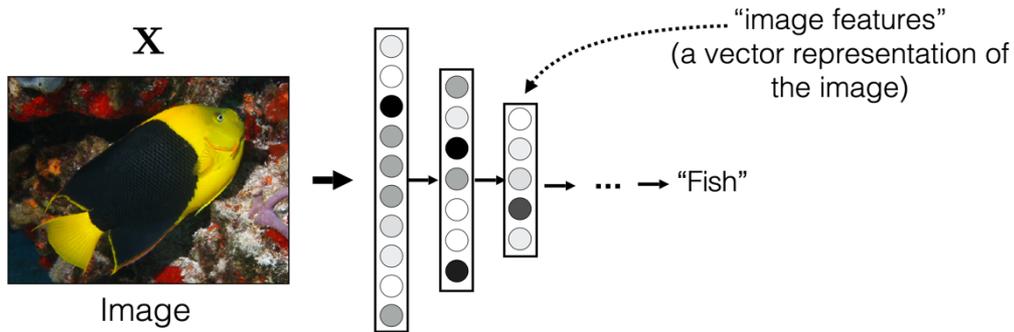
# Supervised learning of a representation



$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(x_i), y_i)$$

Loss function/Cost

# Learning without Labels



## Training data

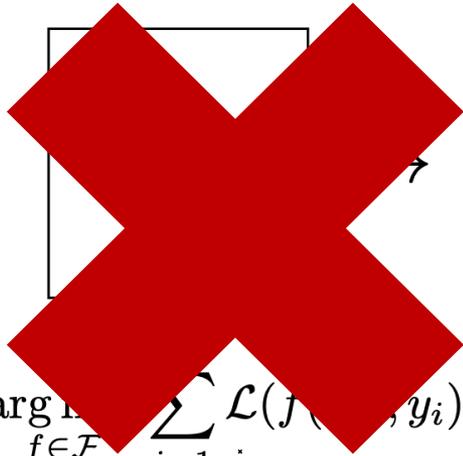
$\{x_1, y_1\}$

$\{x_2, y_2\}$

$\{x_3, y_3\}$

...

→

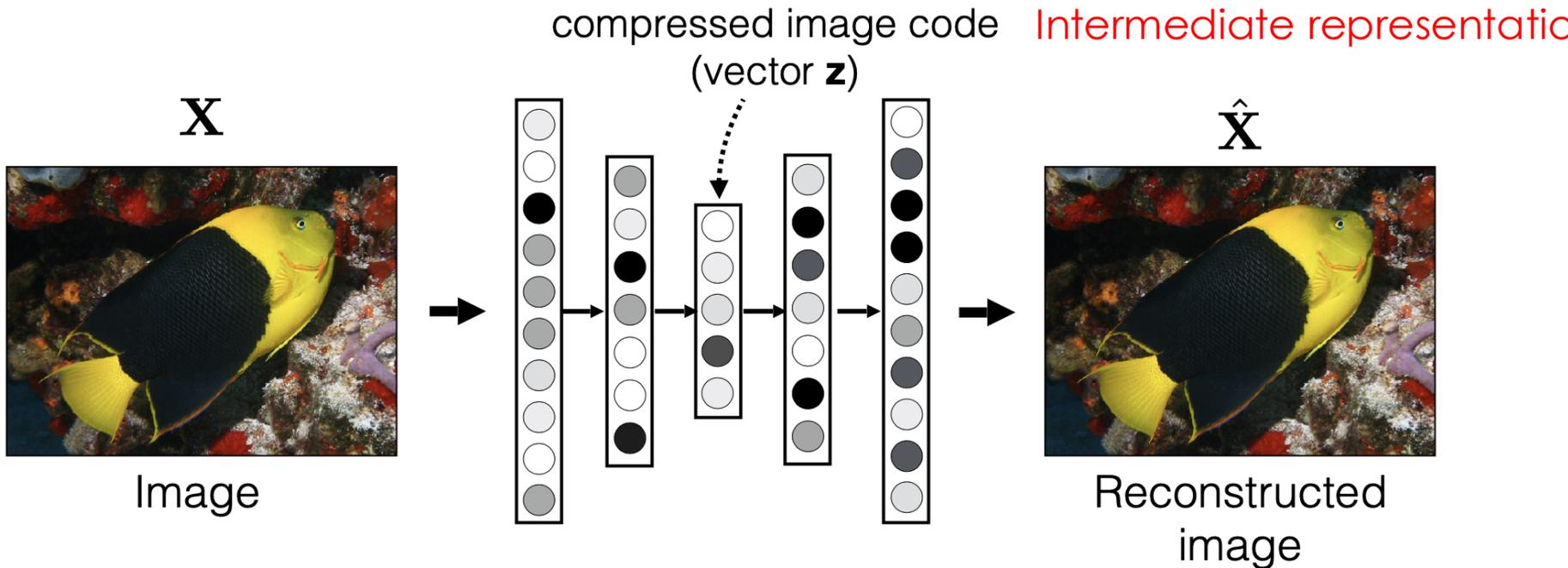


$f : X \rightarrow Y$

$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^n \mathcal{L}(f(x_i), y_i)$$

# Unsupervised Representation Learning

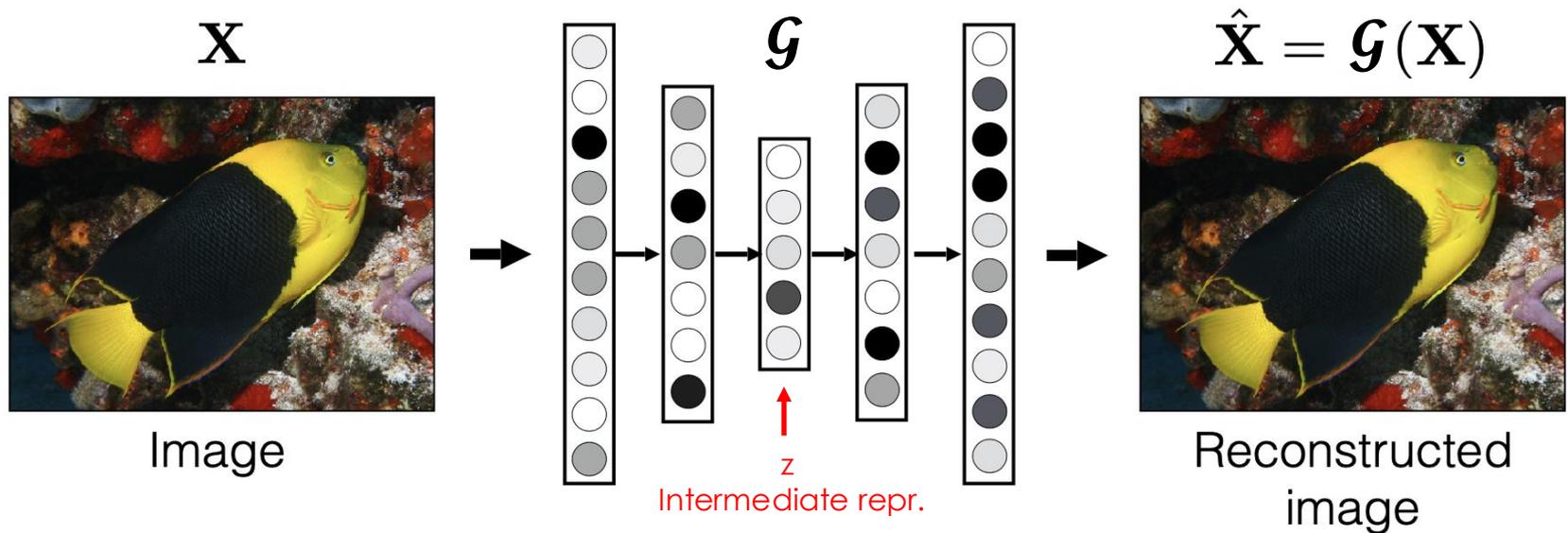
No category or symbolic label. Instead: learn to reconstruct.



One kind of unsupervised model: "Autoencoder"

[e.g., Hinton & Salakhutdinov, Science 2006]

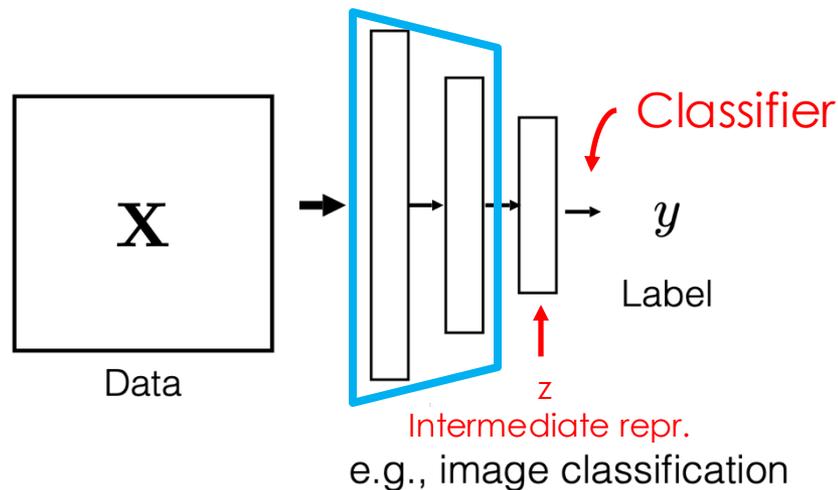
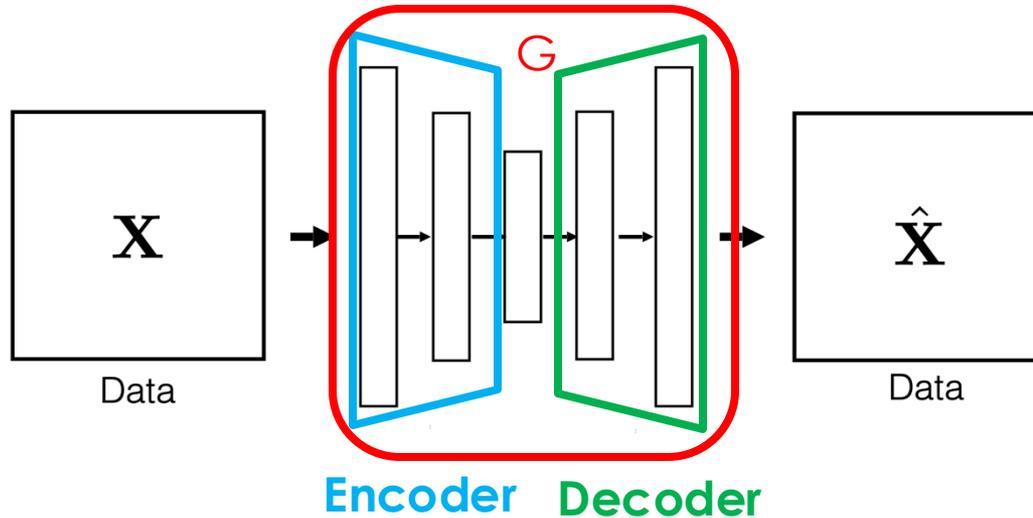
# Autoencoder



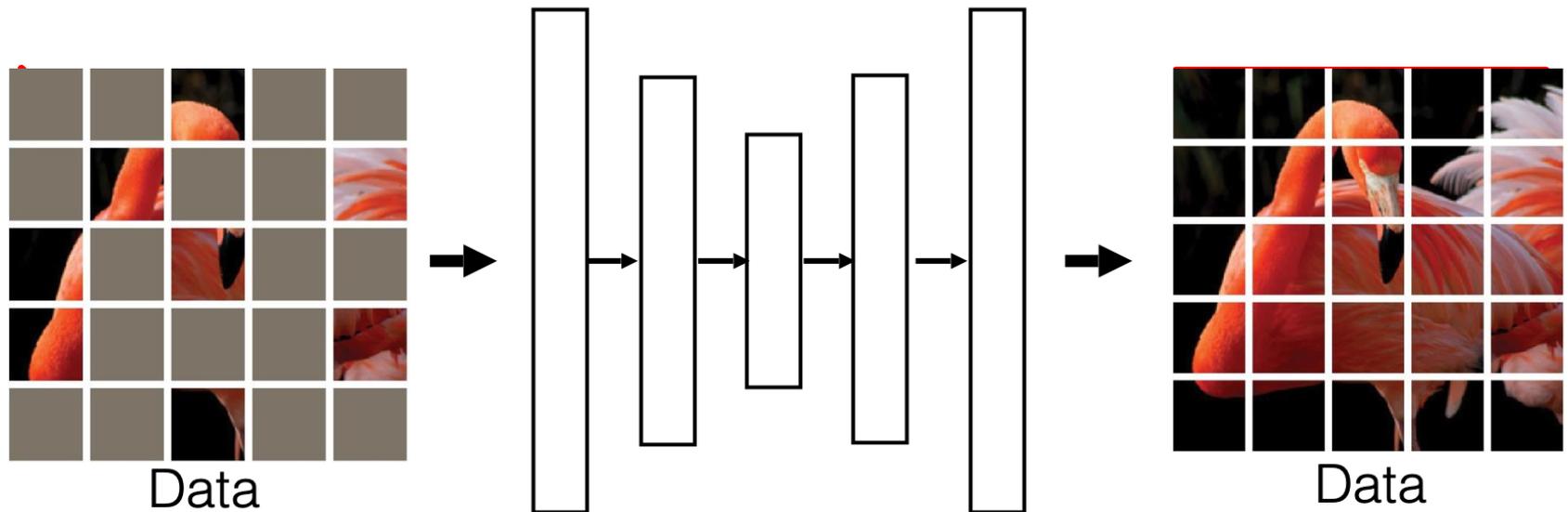
$$\arg \min_{\mathcal{G}} \mathbb{E}_{\mathbf{X}} [||\mathcal{G}(\mathbf{X}) - \mathbf{X}||]$$

Reconstruction loss to minimize by finding optimal  $\mathcal{G}$

# Data Compression & Task Transfer



# Self-Supervision



$$G(X) = \hat{X}$$
$$G(X_1) = \hat{X}_2$$

0 surveys completed



0 surveys underway

## Which of the options are supervised learning objectives in representation learning?

Classification loss

Image Reconstruction Loss

Object Detection loss

Depth Estimation Error from Stereo images compared to g...

Image Synthesis Error for right image of a stereo pair give...

Image synthesis Error of future images from current image

Semantic Segmentation Loss

## Which of the options are unsupervised learning objectives in representation learning?

Classification loss

Image Reconstruction Loss

Object Detection loss

Depth Estimation Error from Stereo images compared to g...

Image Synthesis Error for right image of a stereo pair give...

Image synthesis Error of future images from current image

Semantic Segmentation Loss

## Which of the options are self-supervised learning objectives in representation learning?

Classification loss

Image Reconstruction Loss

Object Detection loss

Depth Estimation Error from Stereo images compared to g...

Image Synthesis Error for right image of a stereo pair give...

Image synthesis Error of future images from current image

Semantic Segmentation Loss

# Representation Learning

Increasing level of difficulty

## Reinforcement Learning (Cherry)

Predicting a scalar reward given once in a while

A few bits for some samples

## Supervised Learning (Chocolate Coat)

Predicting category or vector of scalars per input as provided by human labels.

10-10k bits per sample

## Unsupervised / Self-Supervised Learning (Cake)

Predicting parts of observed input or predicting future observations or events

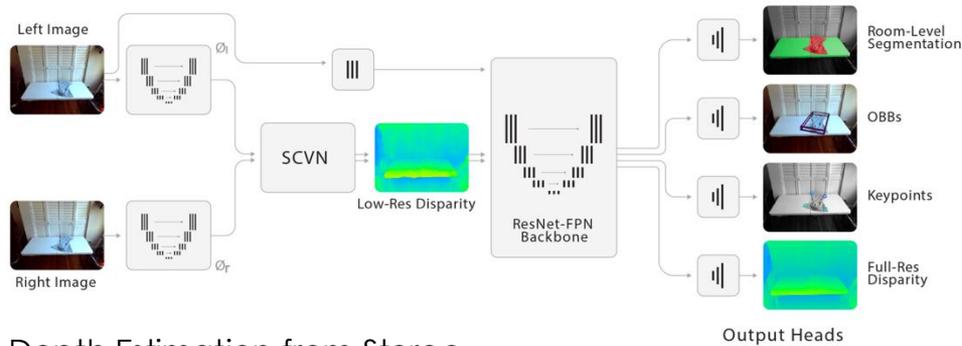
Millions of bits per sample



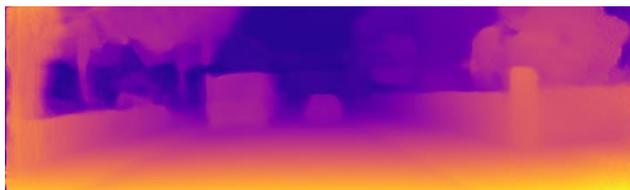
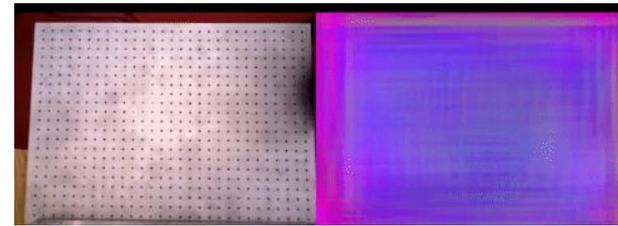
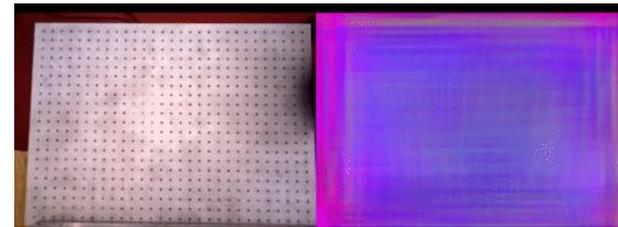
Visualisation Idea by Yann LeCun

Photo by [Kristina Paukshtite](#) from [Pexels](#)

# Let's use representation learning!



Depth Estimation from Stereo  
Supervised Learning



Monocular Depth Estimation  
Unsupervised Learning

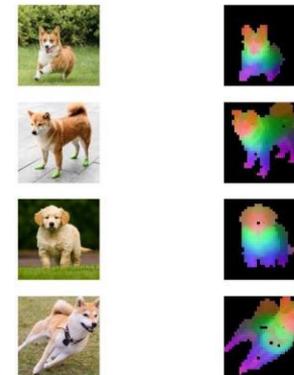
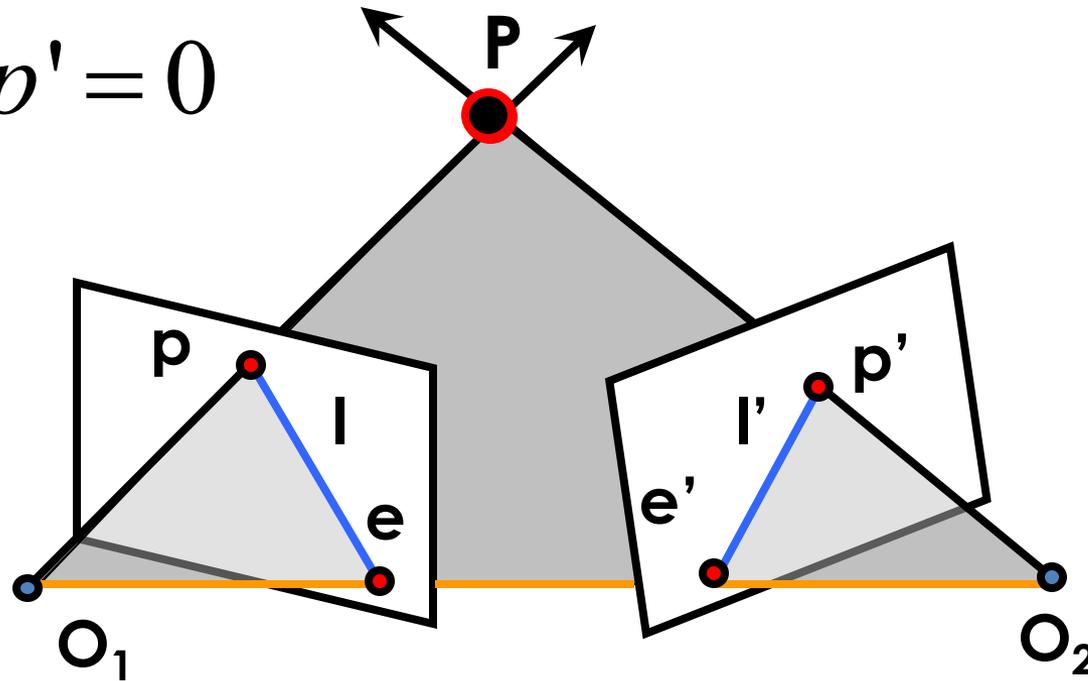


Image by Yunuk Cha.

Finding Correspondences across  
Frames  
Self-Supervised Learning

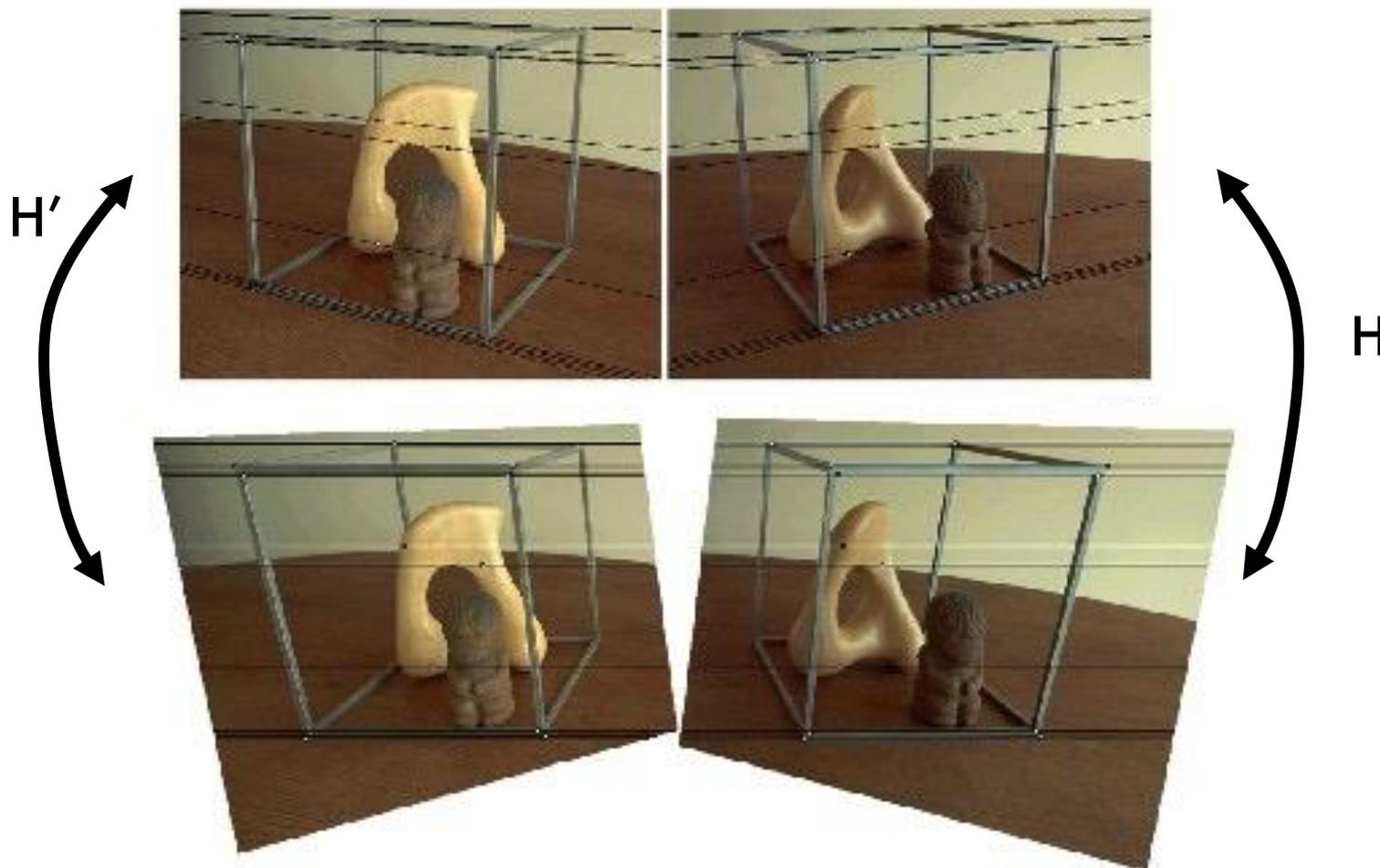
# Epipolar Constraint (Lecture 6)

$$p^T \square F p' = 0$$



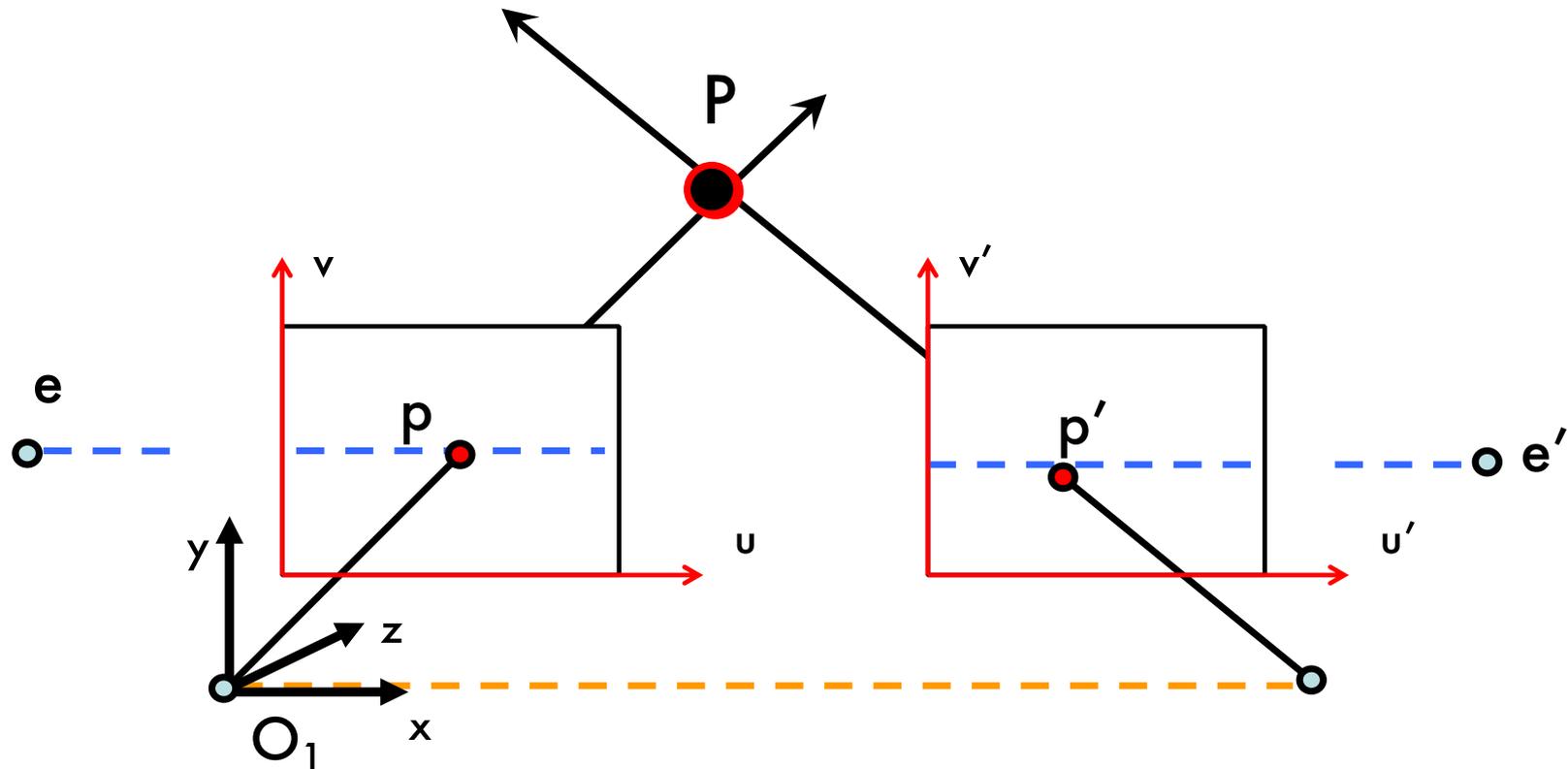
- $l = F p'$  is the epipolar line associated with  $p'$
- $l' = F^T p$  is the epipolar line associated with  $p$
- $F e' = 0$  and  $F^T e = 0$
- $F$  is 3x3 matrix; 7 DOF
- $F$  is singular (rank two)

# Rectification: making two images “parallel”



Courtesy figure S. Lazebnik

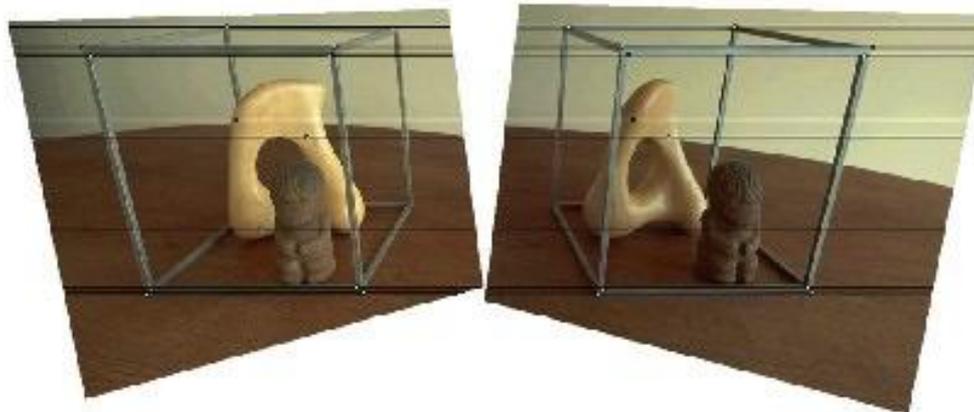
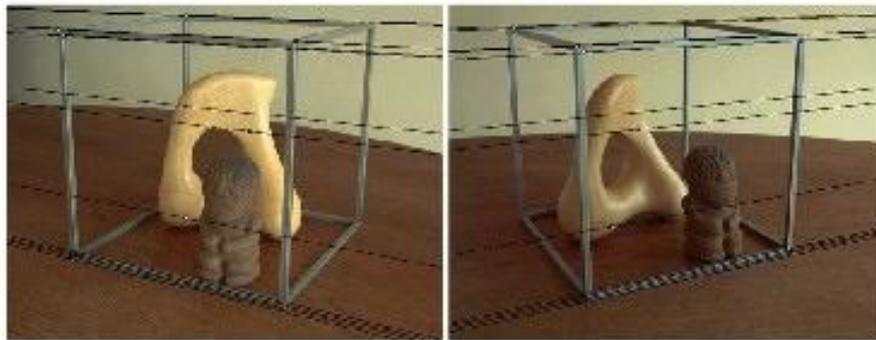
# Parallel image planes



Rectification: making two images “parallel”

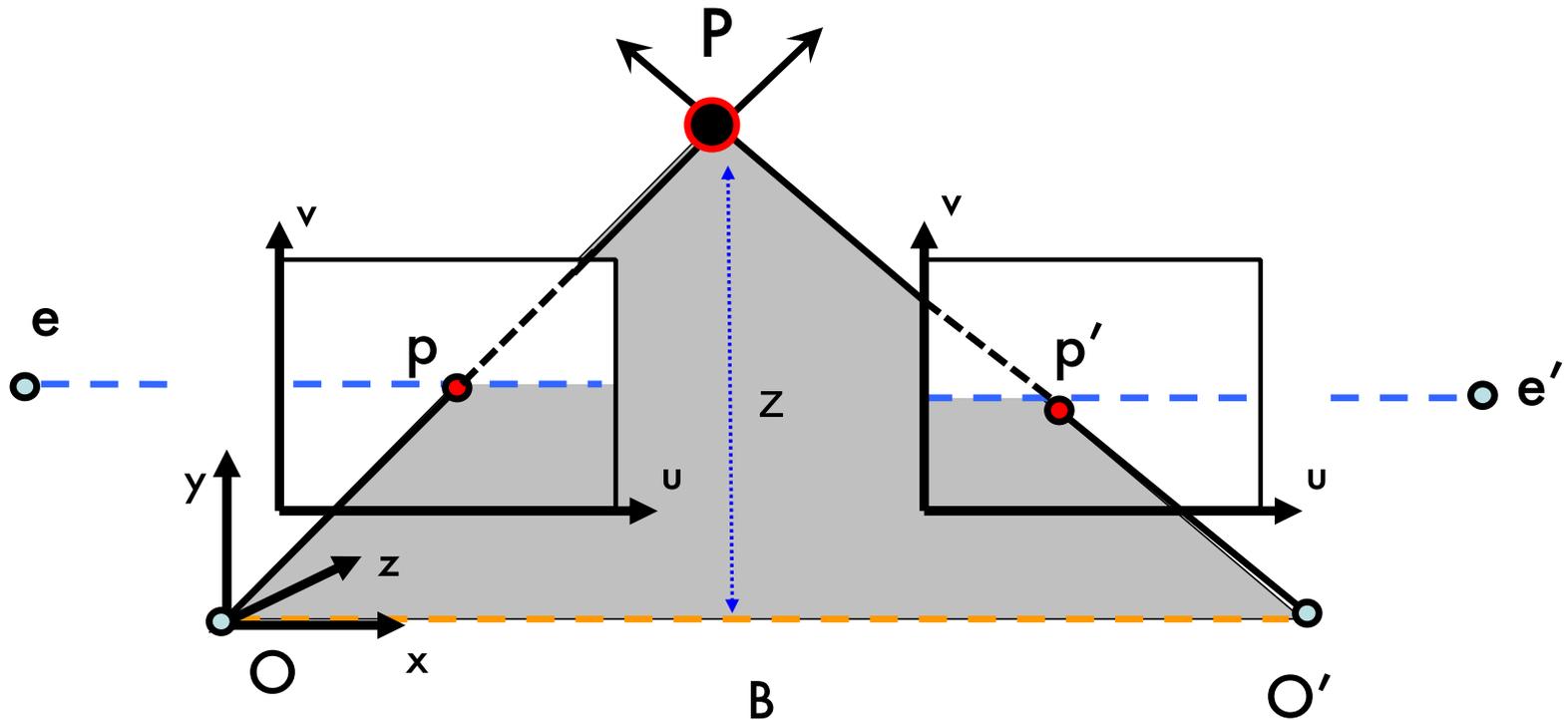
- Epipolar constraint  $\rightarrow v = v'$

# Why are parallel images useful?



- Makes triangulation easy
- Makes the correspondence problem easier

# Point triangulation



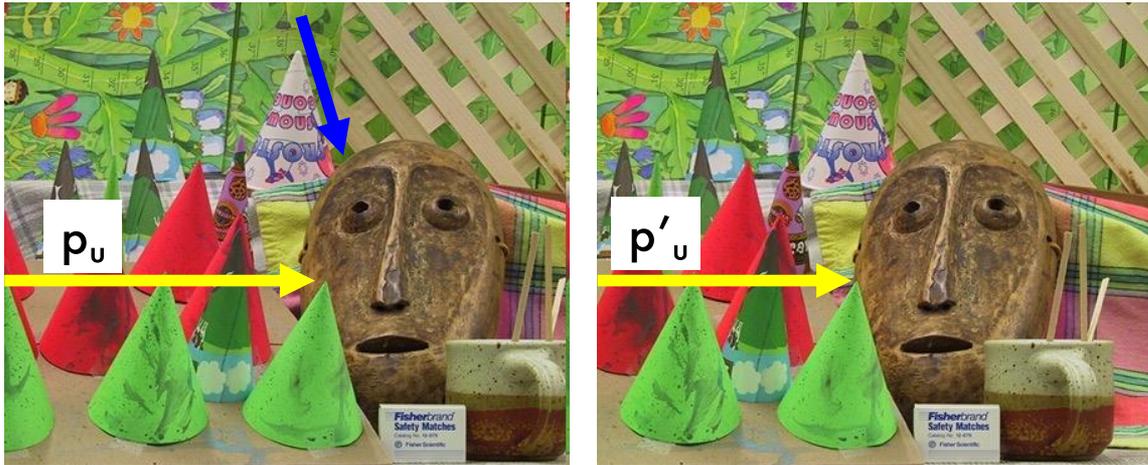
$$p = \begin{bmatrix} p_u \\ p_v \\ 1 \end{bmatrix} \quad p' = \begin{bmatrix} p'_u \\ p_v \\ 1 \end{bmatrix}$$

$$\text{disparity} = p_u - p'_u \propto \frac{B \cdot f}{z} \quad [\text{Eq. 1}]$$

Disparity is inversely proportional to depth  $z$ !

# Disparity maps

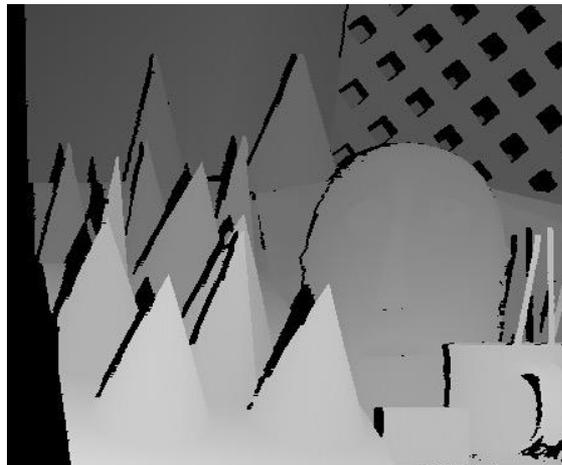
<http://vision.middlebury.edu/stereo/>



$$p_u - p'_u \propto \frac{B \cdot f}{z}$$

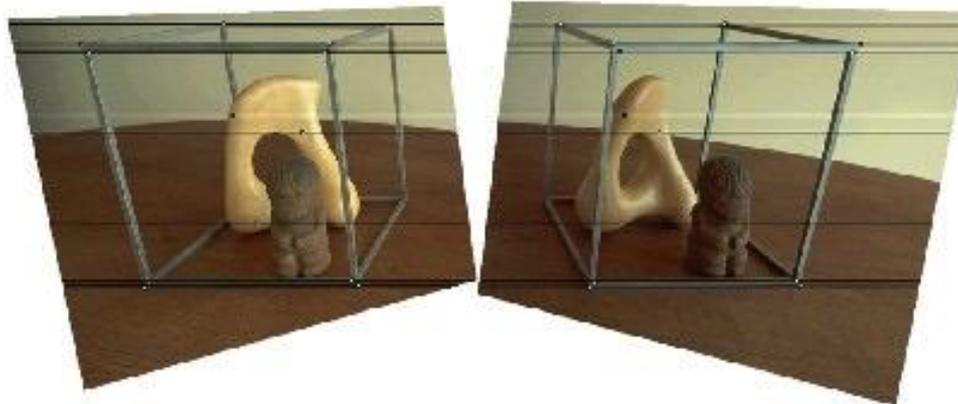
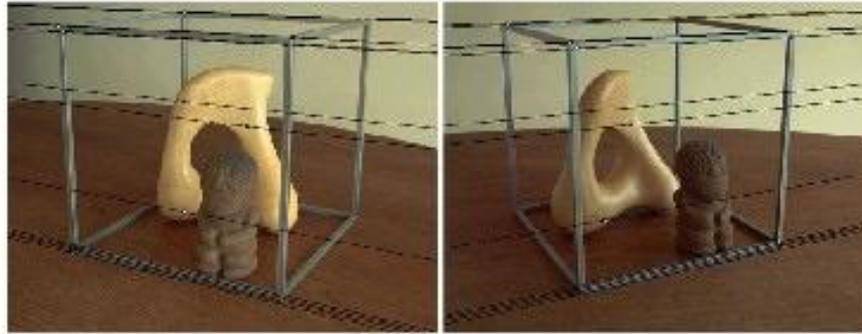
[Eq. 1]

Stereo pair



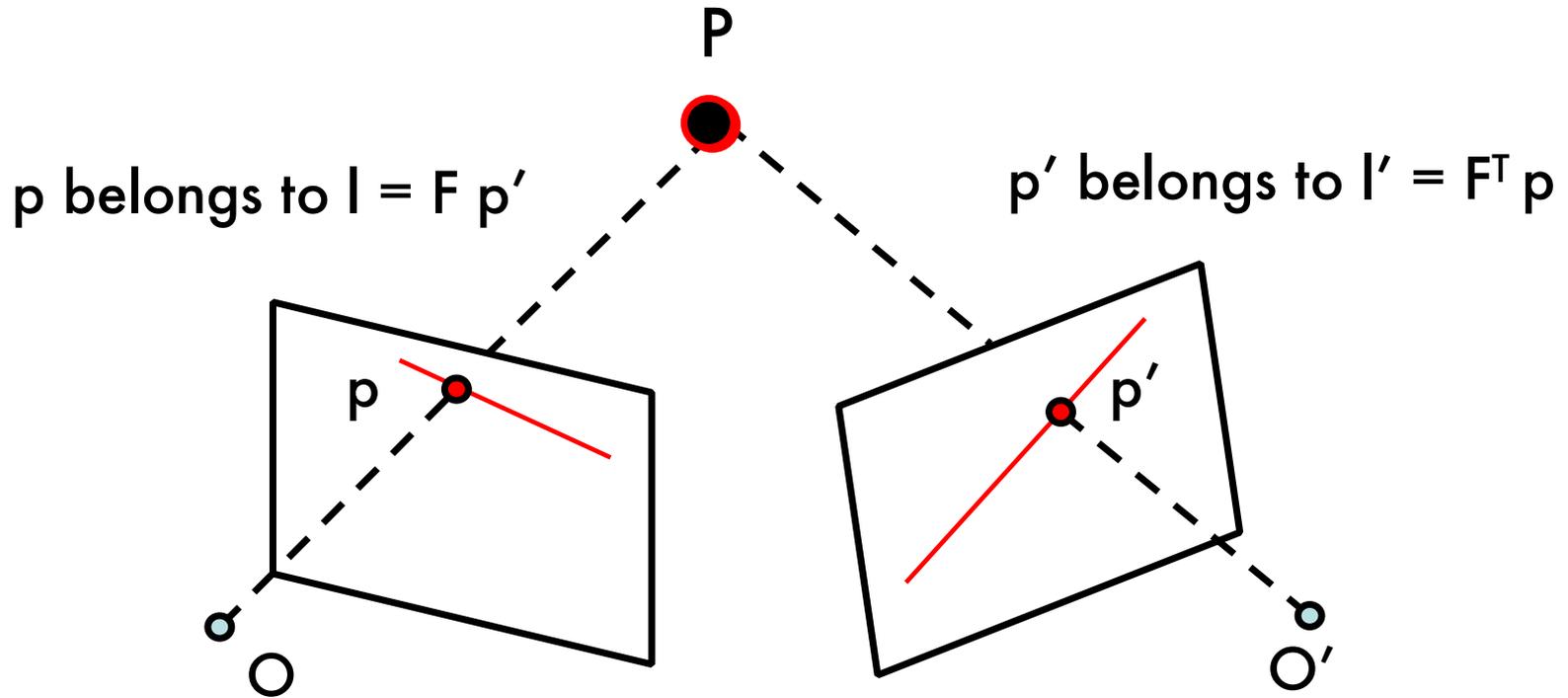
Disparity map / depth map

# Why are parallel images useful?



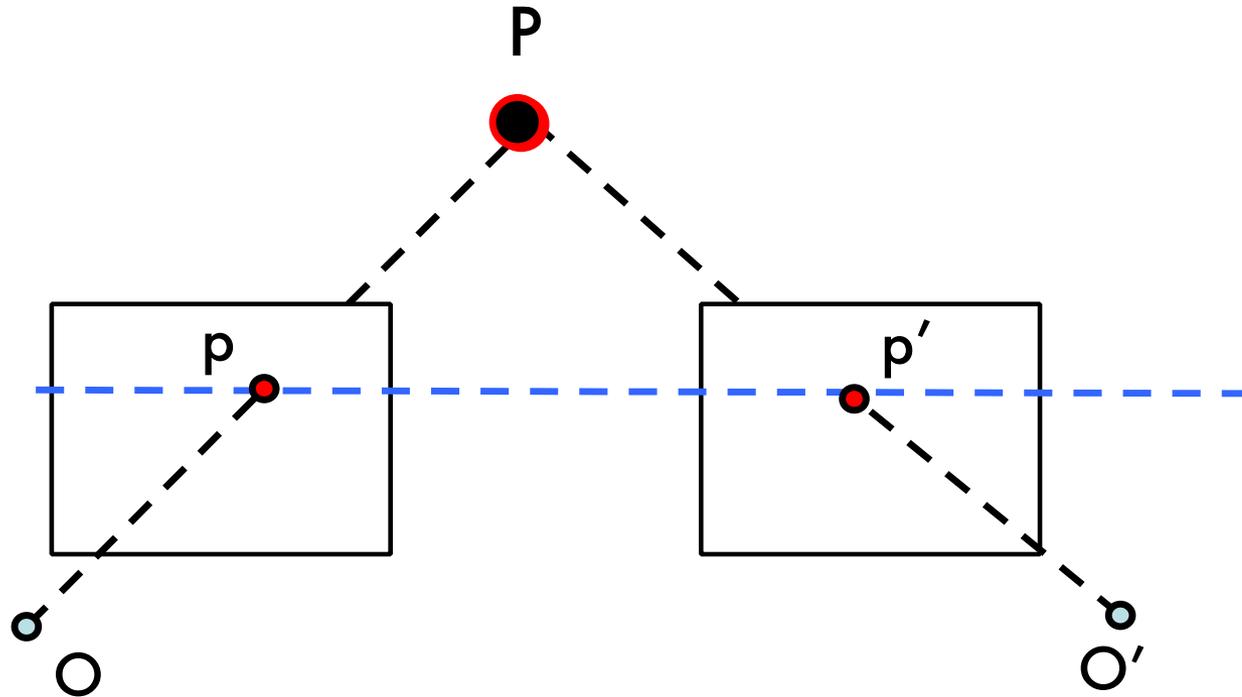
- Makes triangulation easy
- Makes the correspondence problem easier

# Correspondence problem



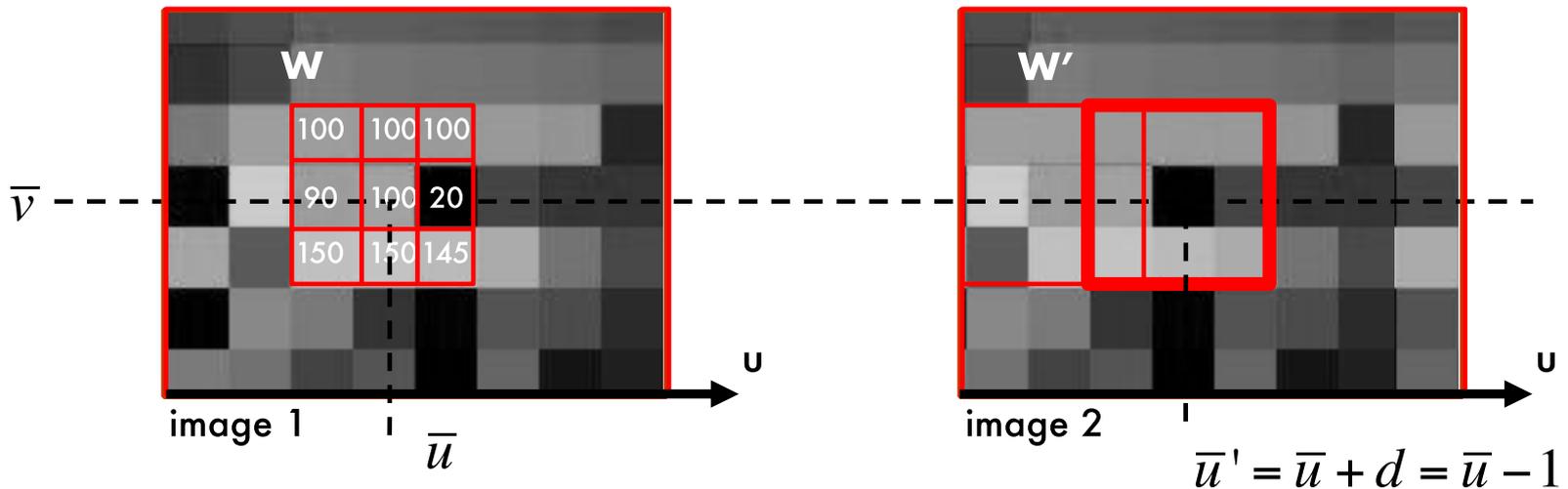
Given a point in 3D, discover corresponding observations in left and right images [also called binocular fusion problem]

# Correspondence problem



When images are rectified, this problem is much easier!

# Window-based correlation



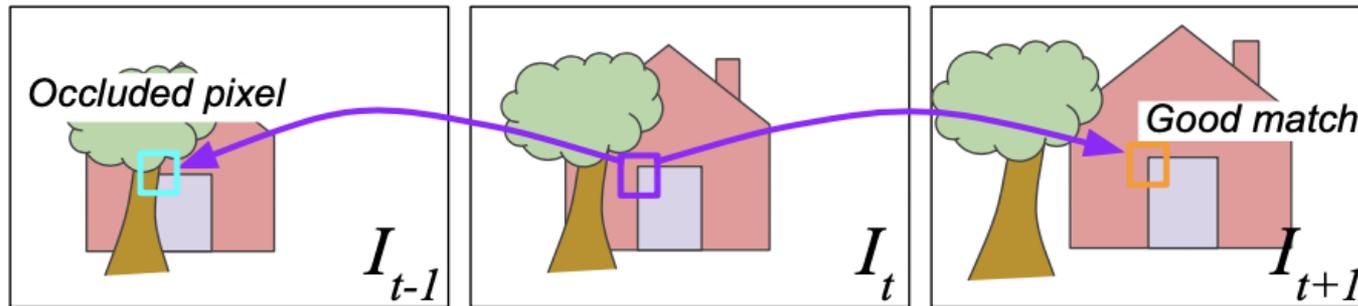
Example:  $\mathbf{W}$  is a 3x3 window in red

$\mathbf{w}$  is a 9x1 vector

$\mathbf{w} = [100, 100, 100, 90, 100, 20, 150, 150, 145]^T$

- Pick up a window  $\mathbf{W}$  around  $\bar{p} = (\bar{u}, \bar{v})$
- Build vector  $\mathbf{w}$
- Slide the window  $\mathbf{W}$  along  $v = \bar{v}$  in image 2 and compute  $\mathbf{w}'(u)$  for each  $u$
- Compute the dot product  $\mathbf{w}^T \mathbf{w}'(u)$  for each  $u$  and retain the max value

# The Correspondence Problem



Occlusions



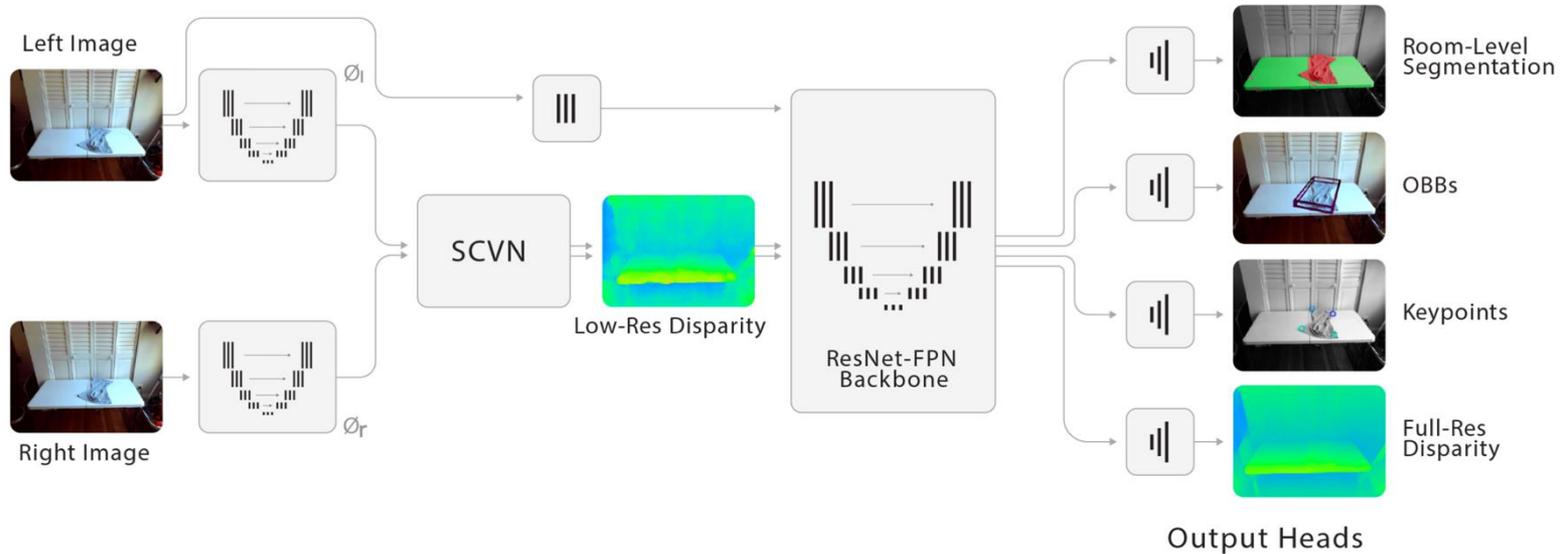
Repetitive Patterns



Hard to match pixels in these regions

Homogenous regions

# Can we learn a similarity function to find corresponding points?



## Supervised Learning of Disparity map.

SimNet: Enabling Robust Unknown Object Manipulation from Pure Synthetic Data via Stereo. CoRL '21. Kollar et al.



# Scaling Up Data Annotation

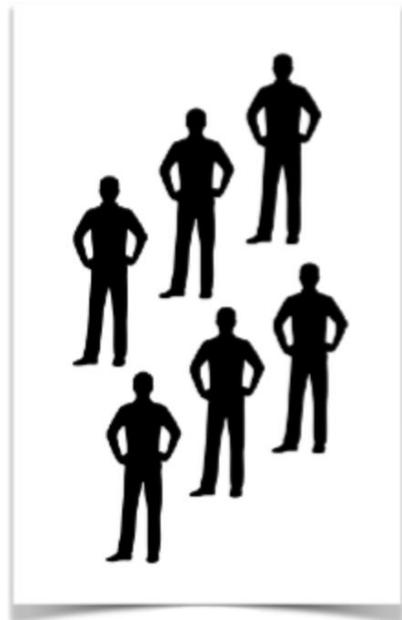


Buy Labeled Real Data

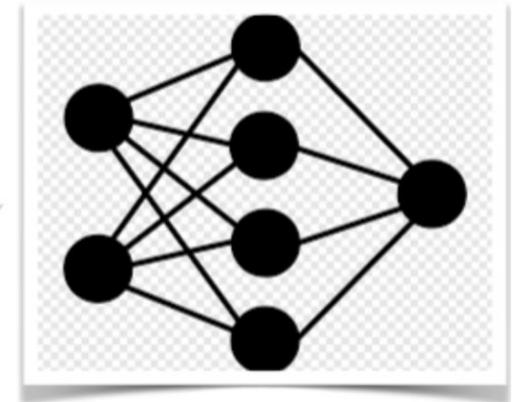


Buy Photorealistic Synthetic Scenes

# Challenges with Data Annotation



ML Engineer



Network

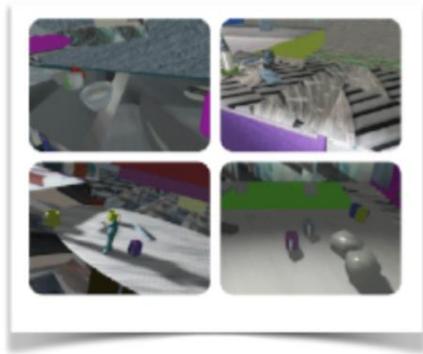
External Data Contractors

Iteration between contractors/engineers is quite painful

- >100K (USD) per cycle,
- >1 month lead time

***This hinders model prototyping and prevent progress.***

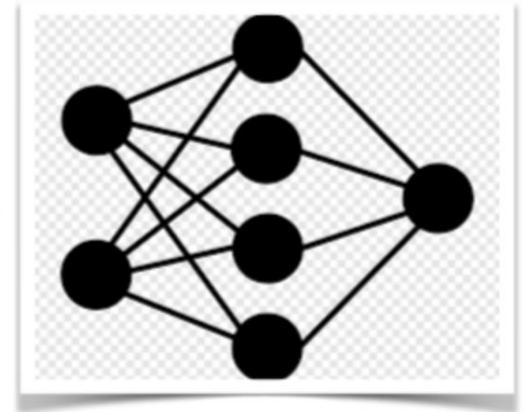
# What if data was programmable?



Procedural Simulator



ML Engineer



Network

# Image = Geometry + Appearance

An image is primarily composed of:

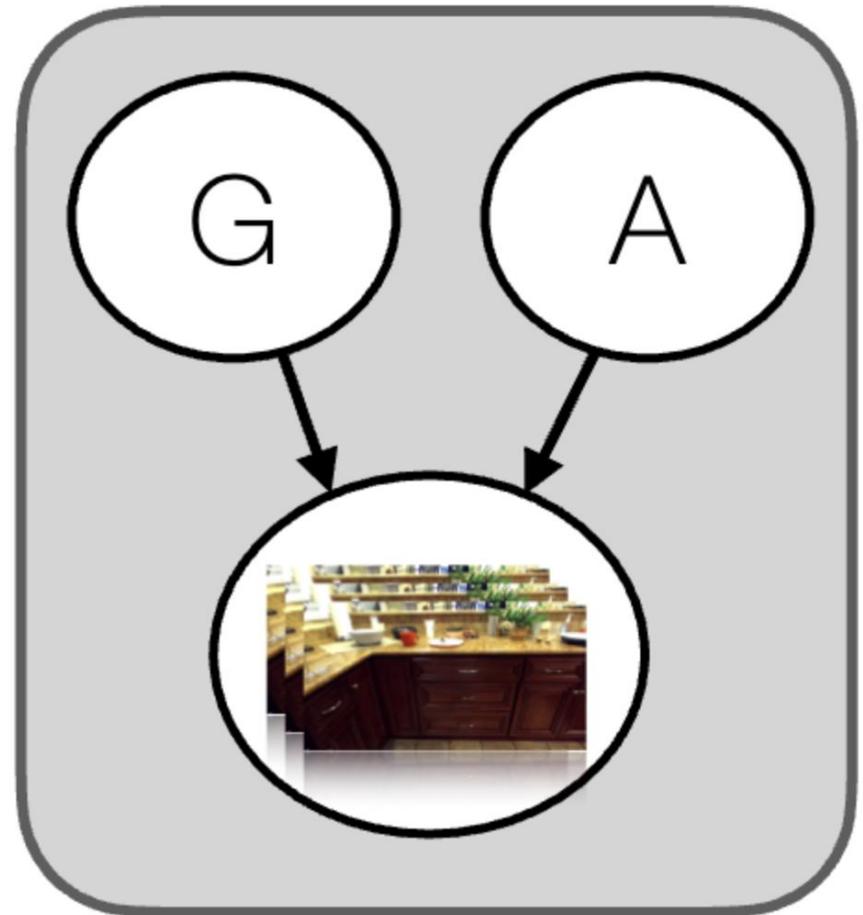
1. **G**eometry
2. **A**pppearance - materials, texture, lighting

Hypothesis:

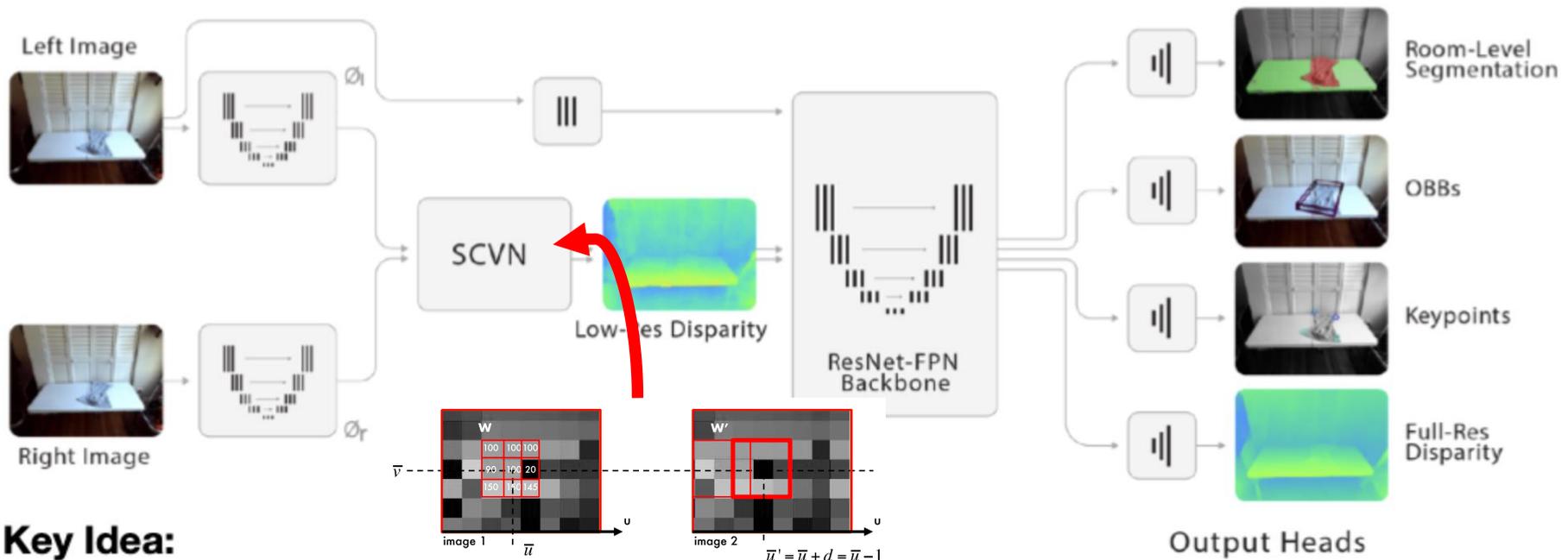
- Geometry **easy** to render/model
- Appearance **hard** to render/model

To achieve sim-to-real transfer, we need to:

1. Minimize the use of appearance features
2. Randomize over scene geometry



# Minimize Use of Appearance Features



## Key Idea:

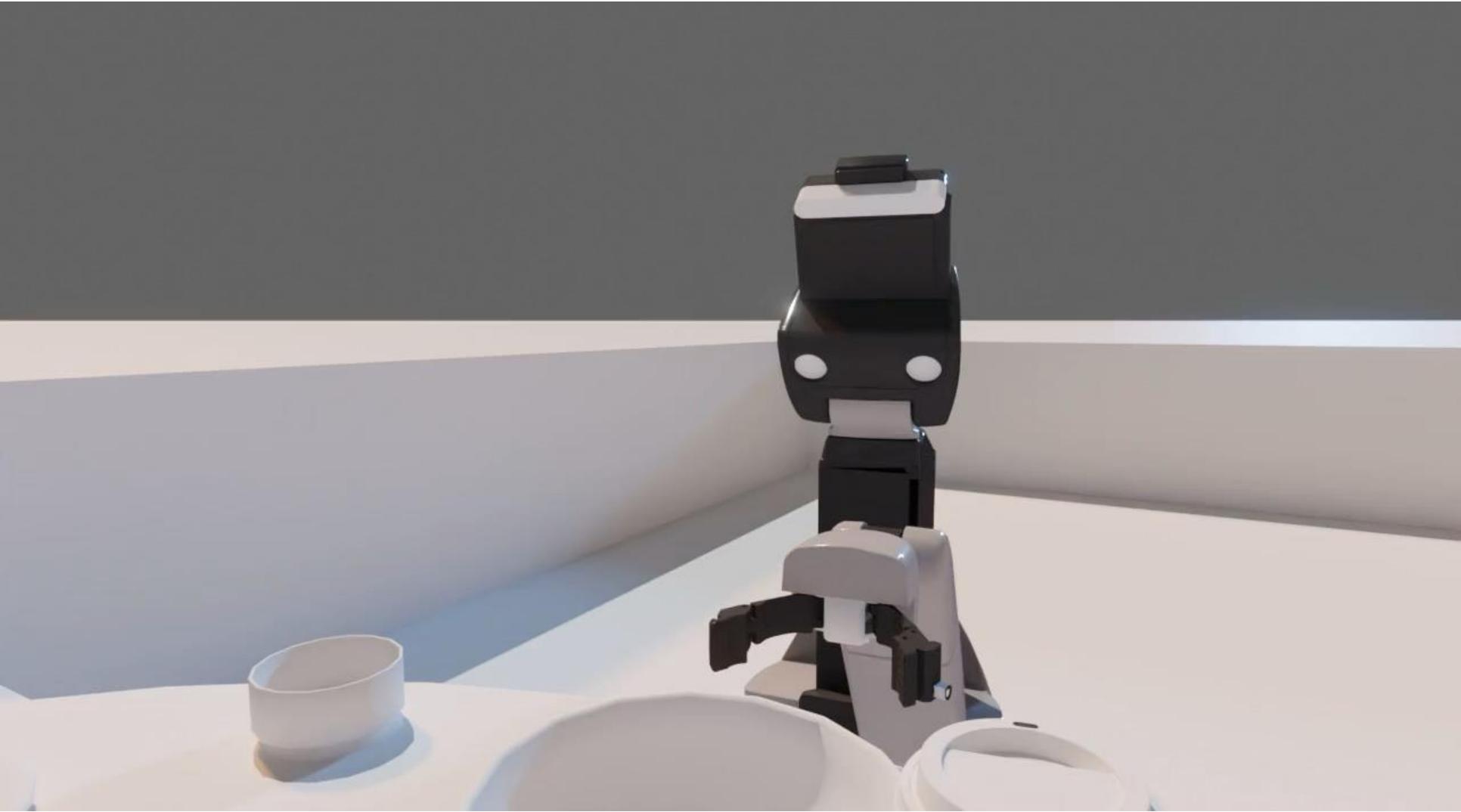
Add *explicit stereo reasoning* in the network to extract geometric features during inference.

$$c(\mathbf{x}_1, \mathbf{x}_2) = \sum_{\mathbf{o} \in [-k, k] \times [-k, k]} \langle \mathbf{f}_1(\mathbf{x}_1 + \mathbf{o}), \mathbf{f}_2(\mathbf{x}_2 + \mathbf{o}) \rangle$$

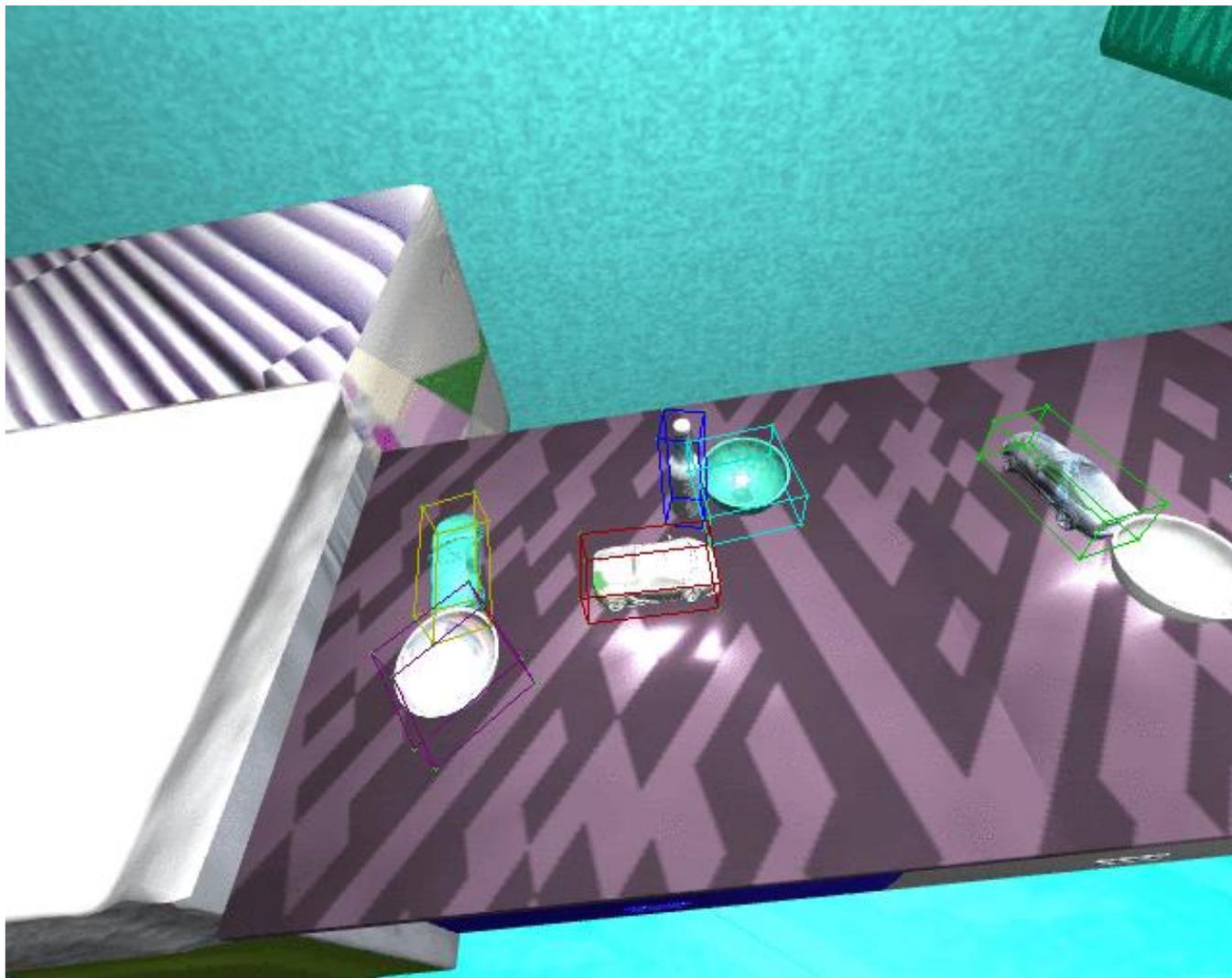
$\mathbf{f}_1$   
 $\mathbf{x}_1$   
 $k$

= Feature  
 = Image Patch location  
 =  $\frac{1}{2}$  Image Patch dimension

# Randomize over Geometry



# Densely Annotated Scenes



# 3D Scene Understanding

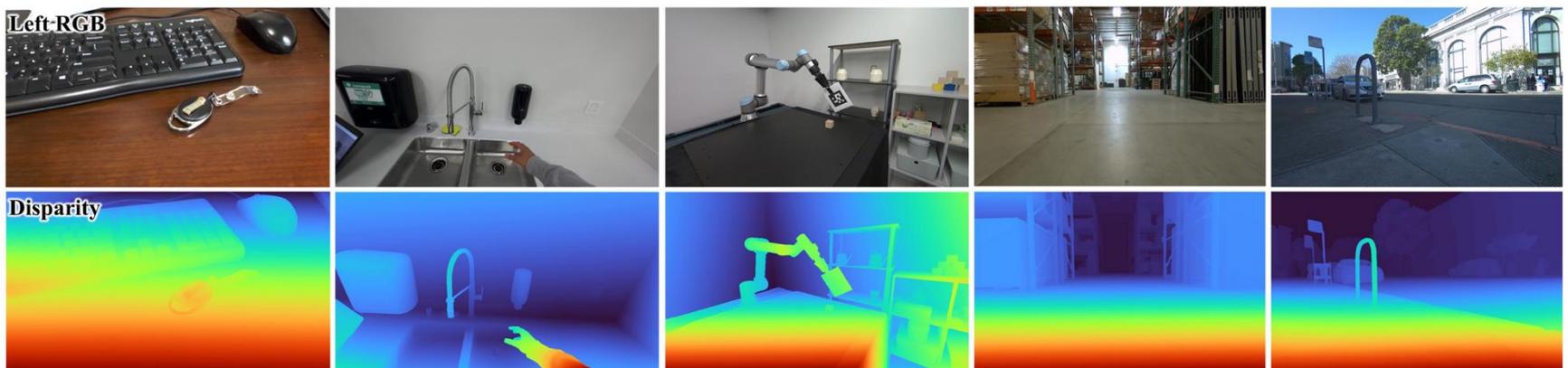


# TODO: Foundation Stereo

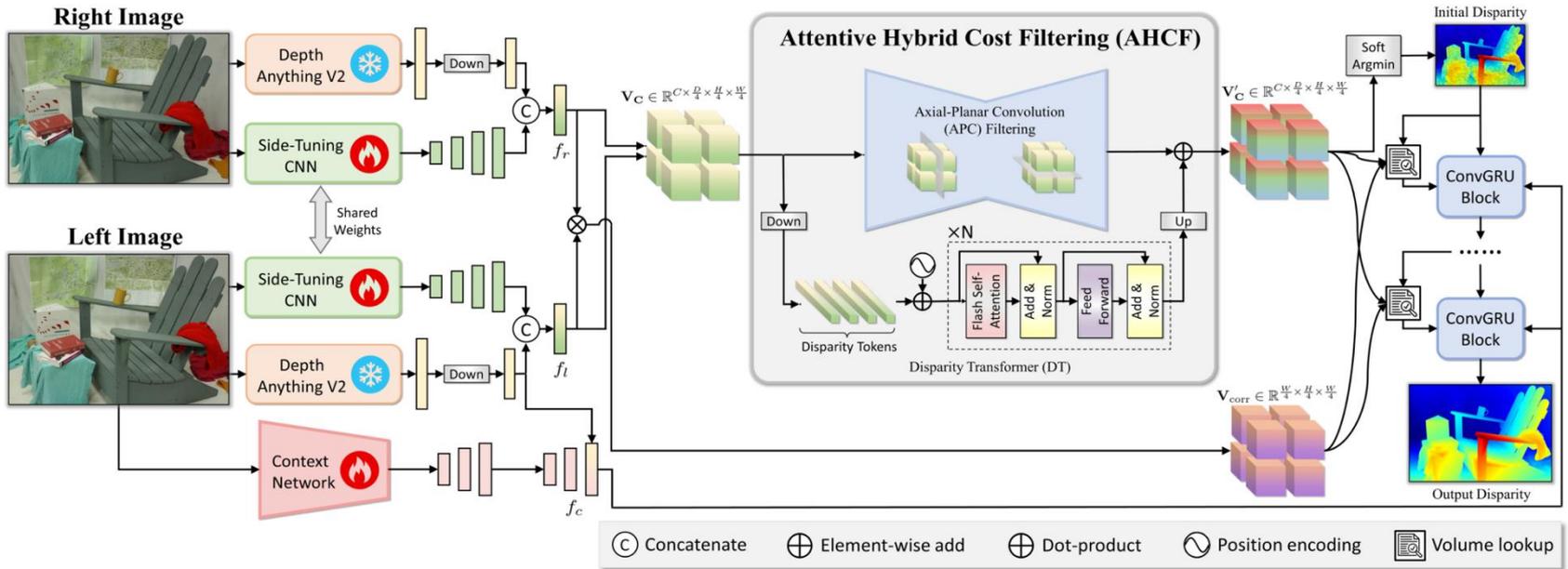
- Focus on correlation portion within Foundation Stereo
- Synthetic images for training
- What is different from earlier work?
  - Much more complex architecture and use of DepthAnything priors from monocular stereo

# FoundationStereo: Zero Shot Stereo Matching. Wen et al. CVPR 2025.

- Synthetic Training Images
- Much more complex architecture
- At the core: Feature Correlation
- Here: Group-wise Correlation



# Foundation Stereo



$$\mathbf{V}_{\text{gwc}}(g, d, h, w) = \left\langle \widehat{f}_{l,g}^{(4)}(h, w), \widehat{f}_{r,g}^{(4)}(h, w - d) \right\rangle,$$

$$\mathbf{V}_{\text{cat}}(d, h, w) = \left[ \text{Conv}(f_l^{(4)})(h, w), \text{Conv}(f_r^{(4)})(h, w - d) \right]$$

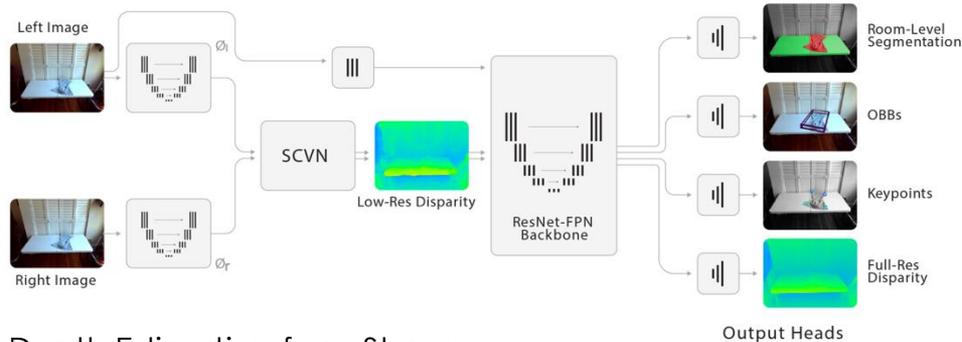
$$\mathbf{V}_{\mathbf{C}}(d, h, w) = [\mathbf{V}_{\text{gwc}}(d, h, w), \mathbf{V}_{\text{cat}}(d, h, w)] \quad (1)$$

# Foundation Stereo

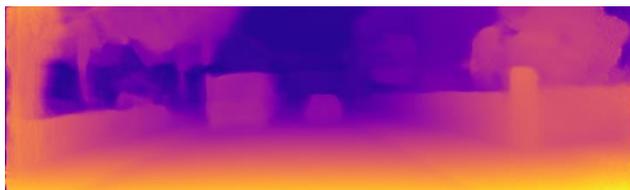
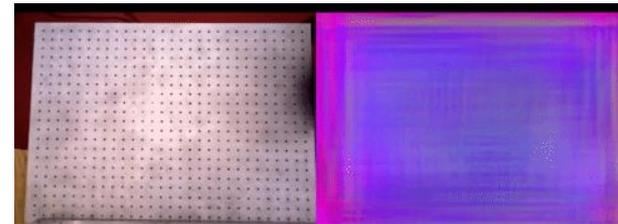
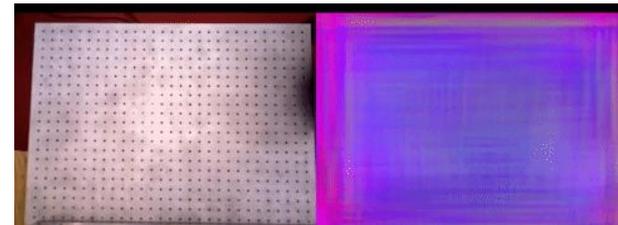


- Synthetic Training Images
- With many textures, geometries and poses
- Self-curation by discarding image pairs that are hard to predict

# Let's use representation learning!



Depth Estimation from Stereo  
Supervised Learning



Monocular Depth Estimation  
Unsupervised Learning



Image by Yunuk Cha.  
Finding Correspondences across  
Frames  
Self-Supervised Learning

# What if we don't need to form correspondences between images?

- Can we estimate depth from a single image?

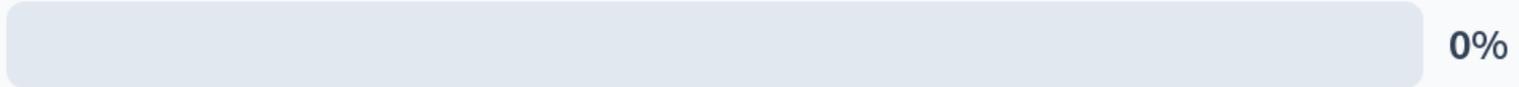
Can you estimate depth from a single RGB image?

Yes

No

## Can you estimate depth from a single RGB image?

Yes

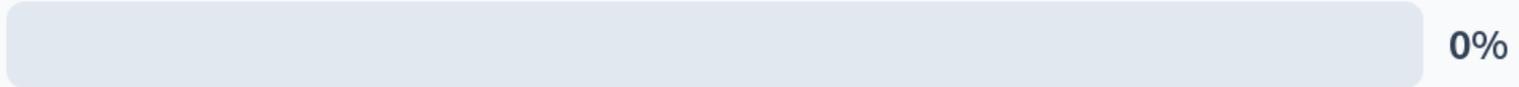


No



## Can you estimate depth from a single RGB image?

Yes



No



# Unsupervised Monocular Depth Estimation with Left-Right Consistency

Clément Godard<sup>1</sup>

Oisín Mac Aodha<sup>2</sup>

Gabriel J. Brostow<sup>1</sup>

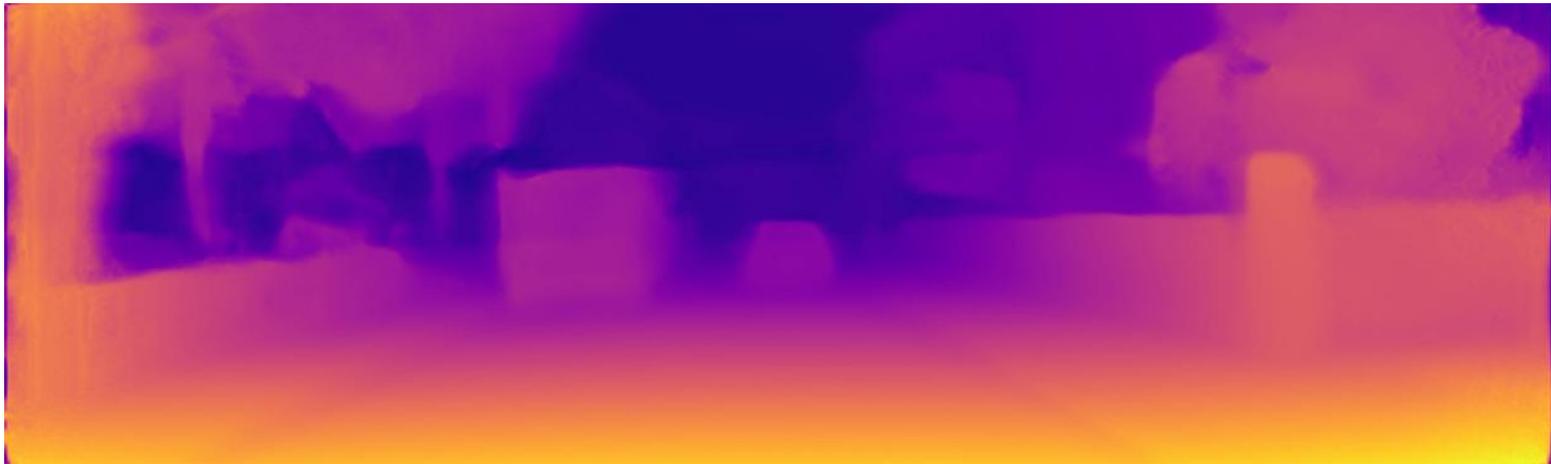
<sup>1</sup>University College London

<sup>2</sup>Caltech





Input image



**Result**

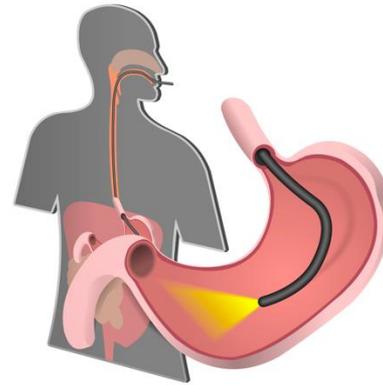
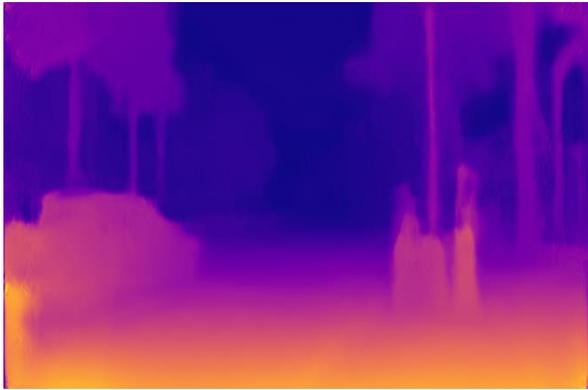
# Why depth?



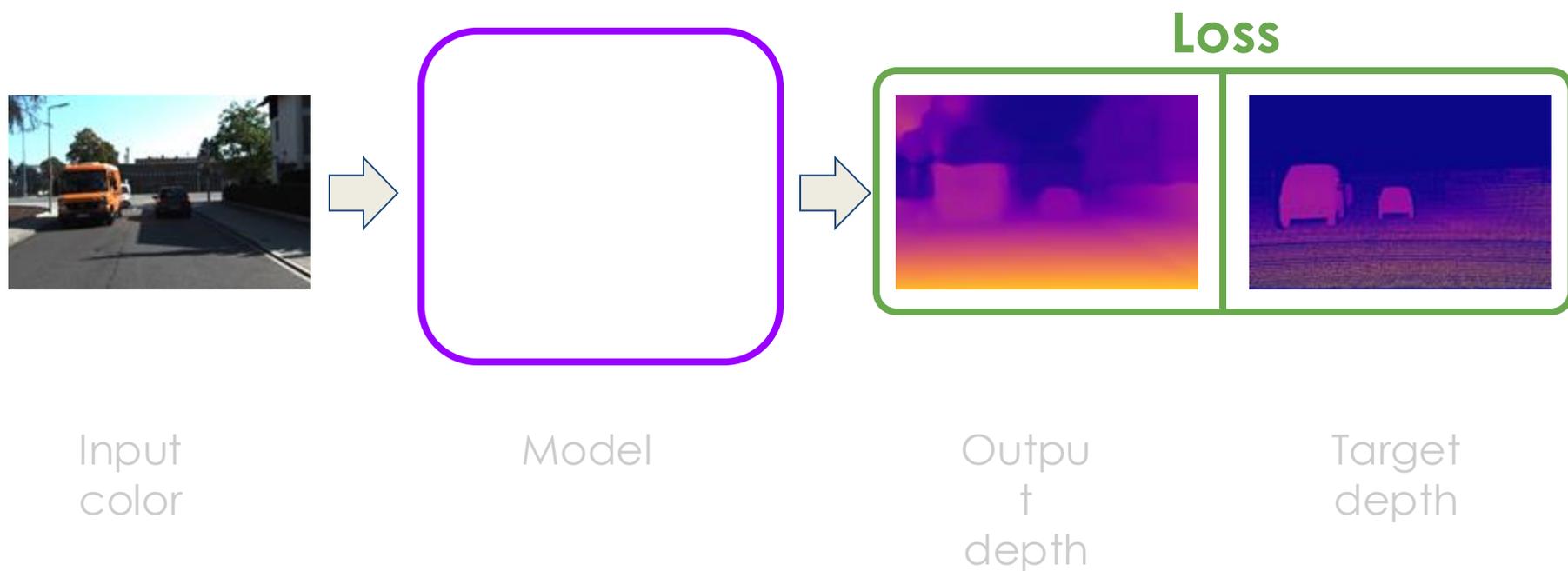
# How do we usually get depth?



# Why monocular?

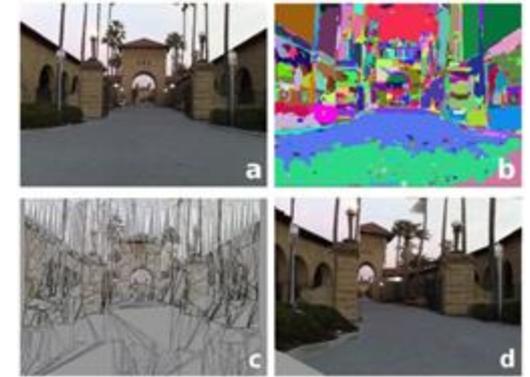


# Previous approaches - Supervised



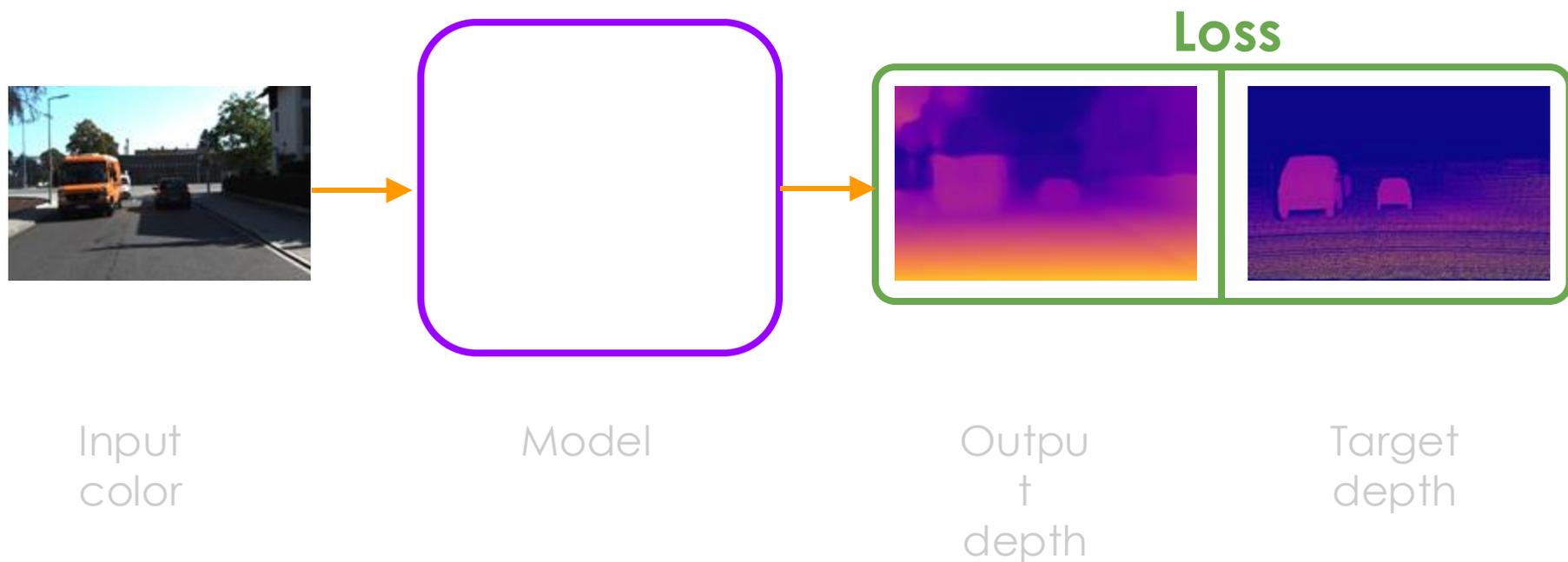
# Previous approaches - Supervised

Automatic Photo Pop-up [SIGGRAPH 05]

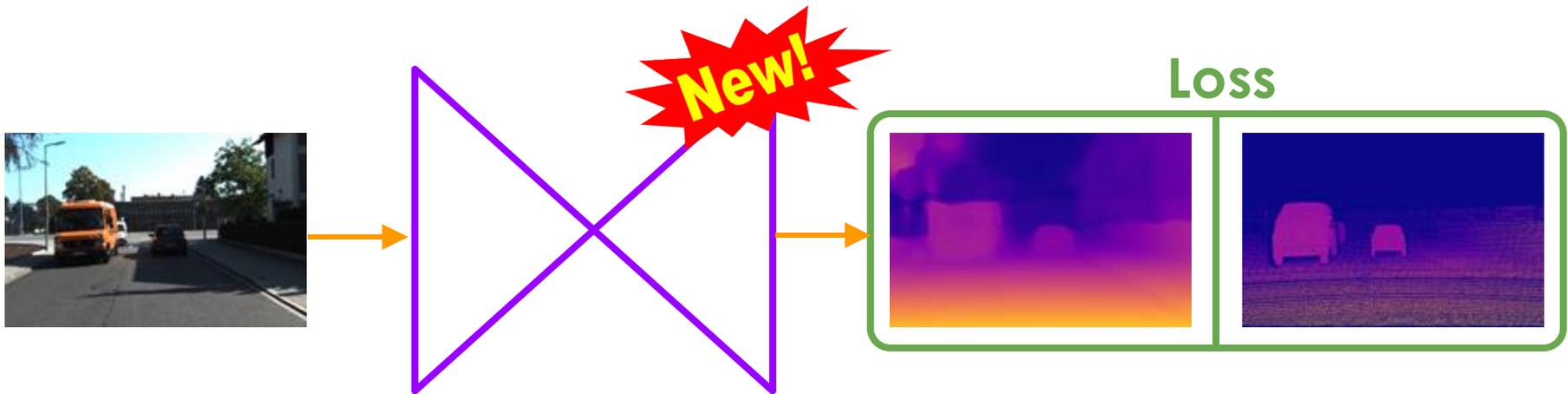


Make3D [PAMI 08]

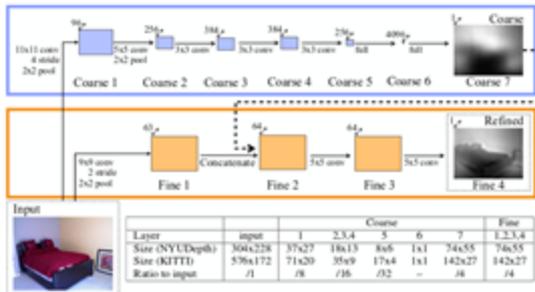
# Previous approaches - Supervised



# Previous approaches - Supervised



Input color



Output  
+  
depth

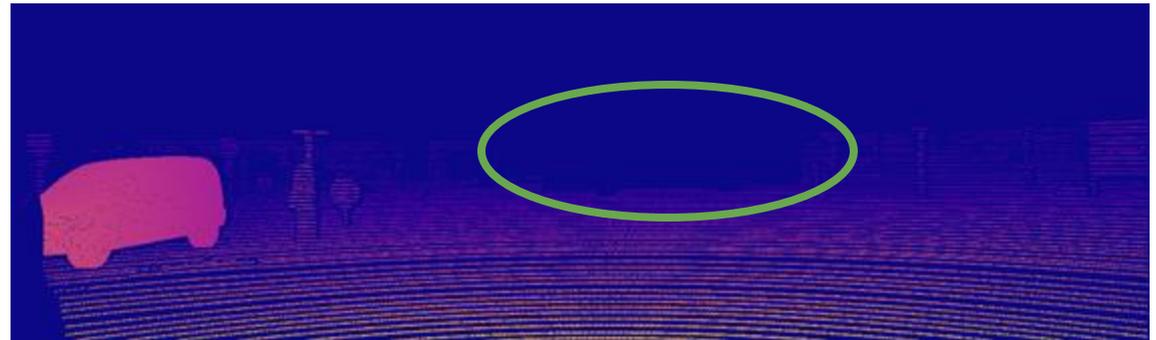
Target  
depth

Eigen et al. [NIPS 14]

Li et al., Laina et al., Cao et al., ...



# KITTI 2015



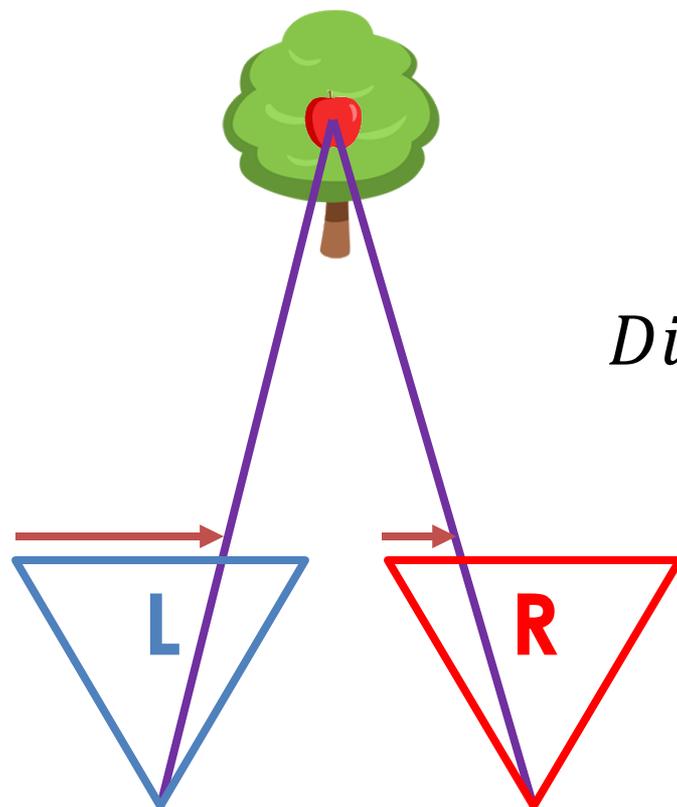
# IR Structured Light

- Does not work well outside





# Depth from stereo



$$\textit{Disparity} \propto \frac{1}{\textit{Depth}}$$

# Let's train with stereo data!



Point Grey



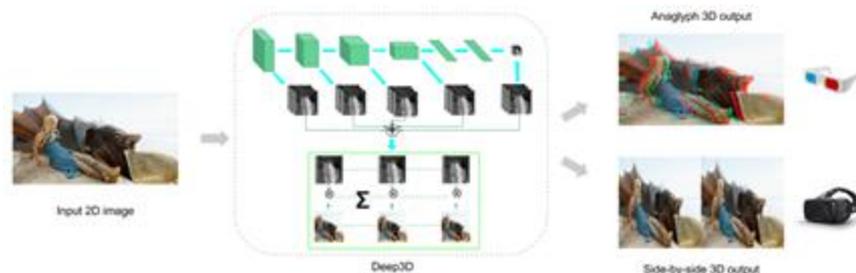
Apple



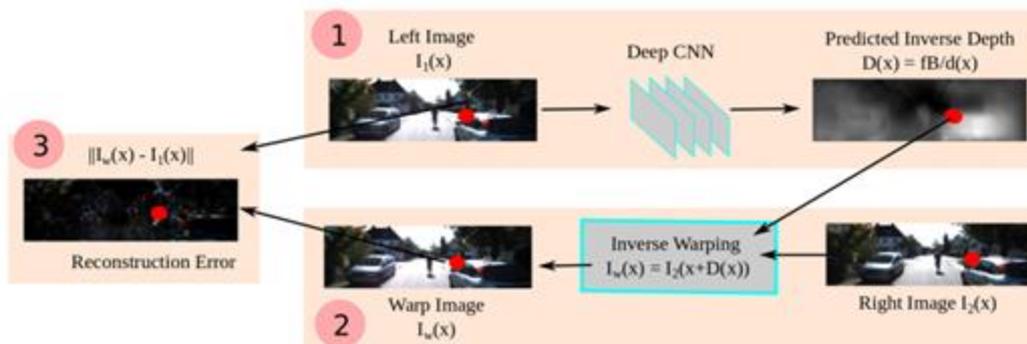
Stereolabs

# Previous approaches - Unsupervised

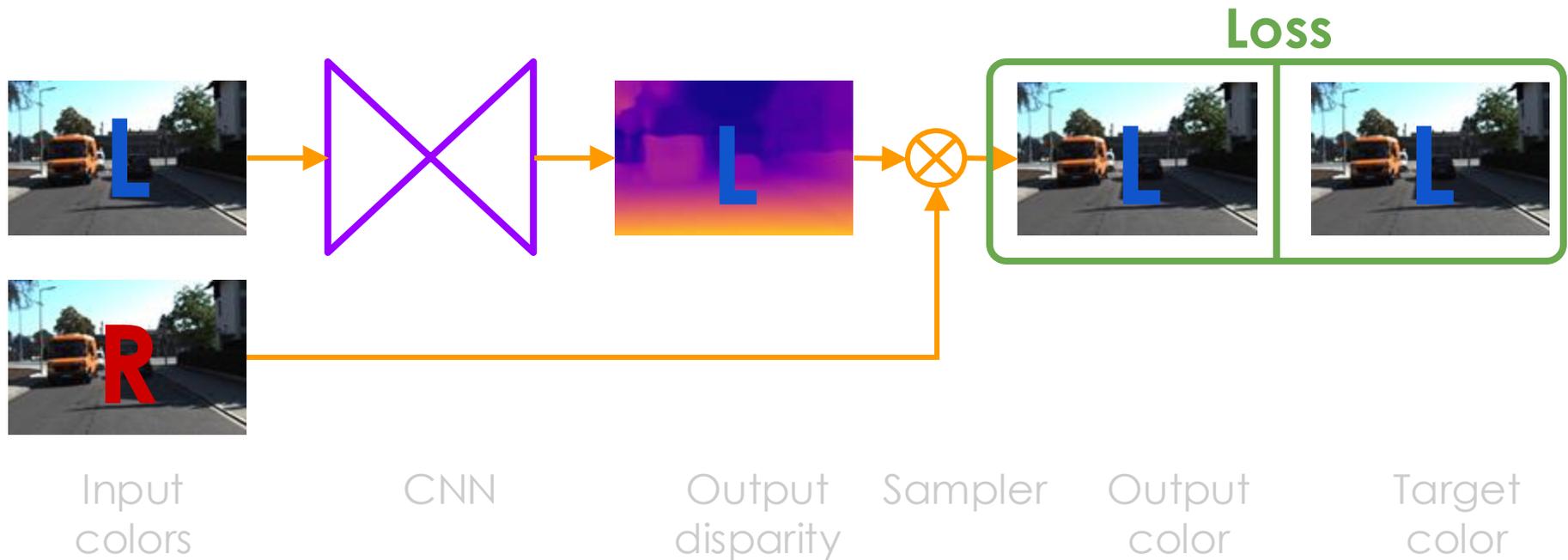
Deep3D  
Xie et al. [ECCV 16]



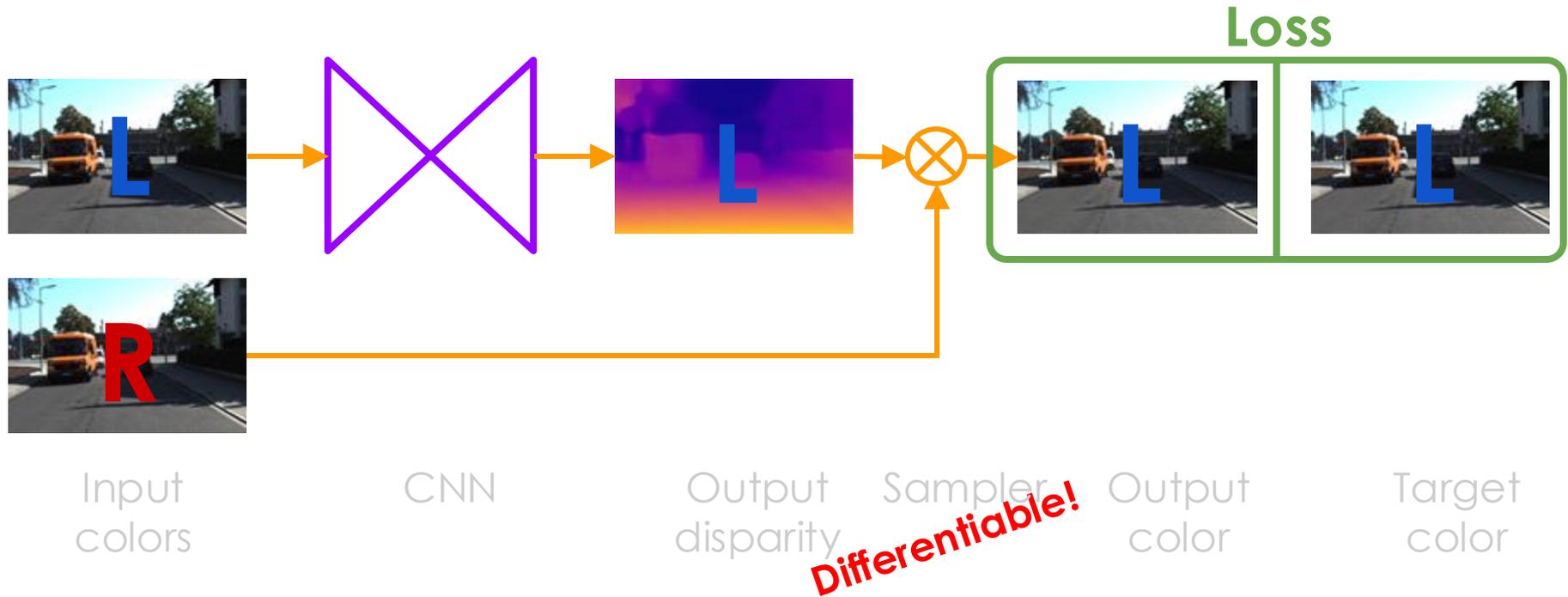
Garg et al. [ECCV 16]



# Unsupervised depth estimation - Concept

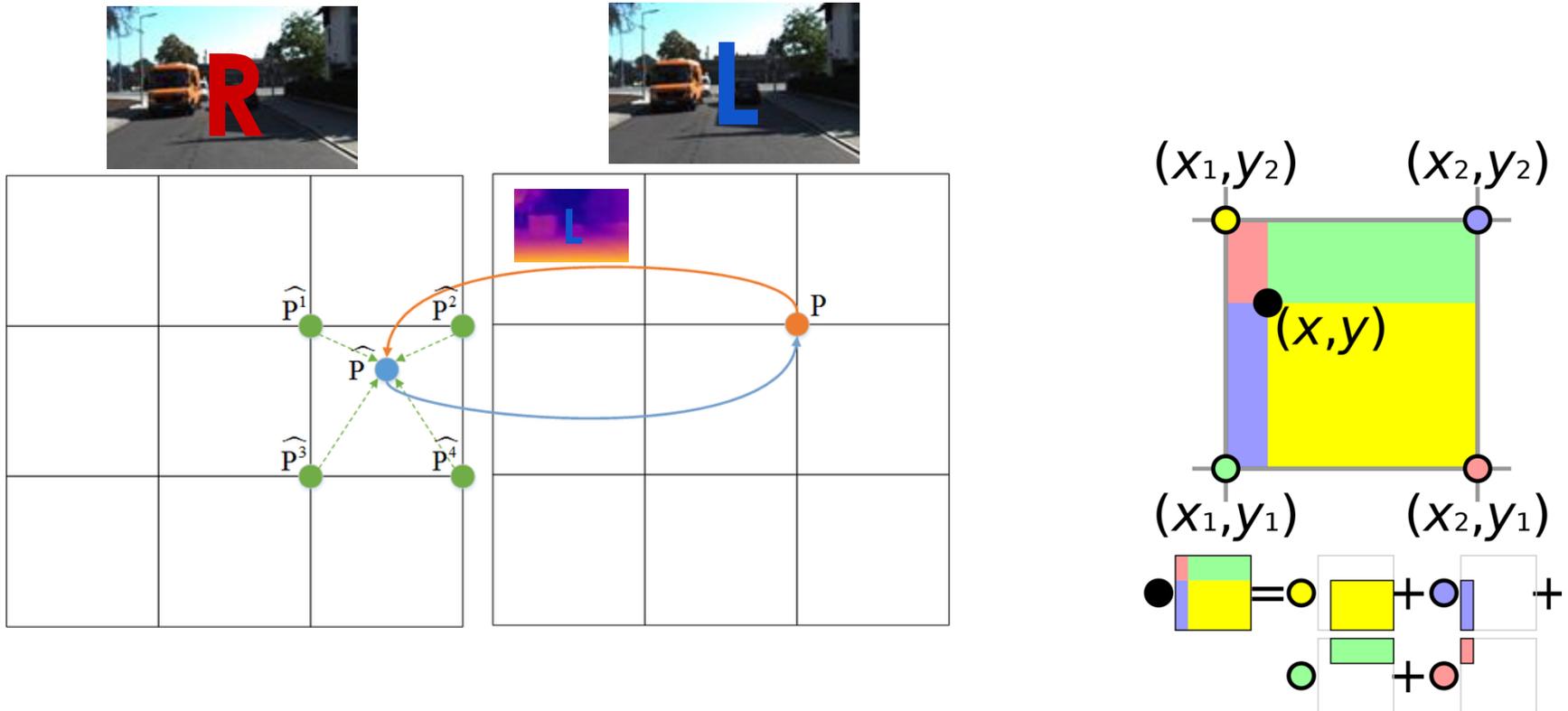


# Unsupervised depth estimation - Baseline



Spatial transformer networks, Jaderberg et al. [NIPS 15]

# Bilinear Sampling



Spatial transformer networks, Jaderberg et al. [NIPS 15]

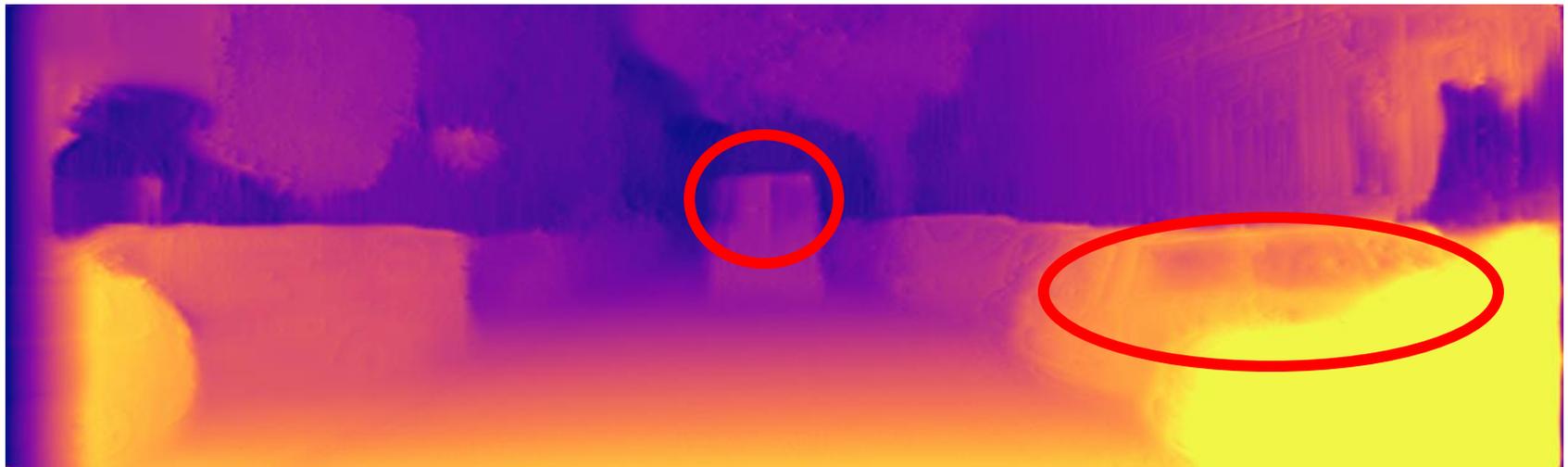
# Input



# Baseline

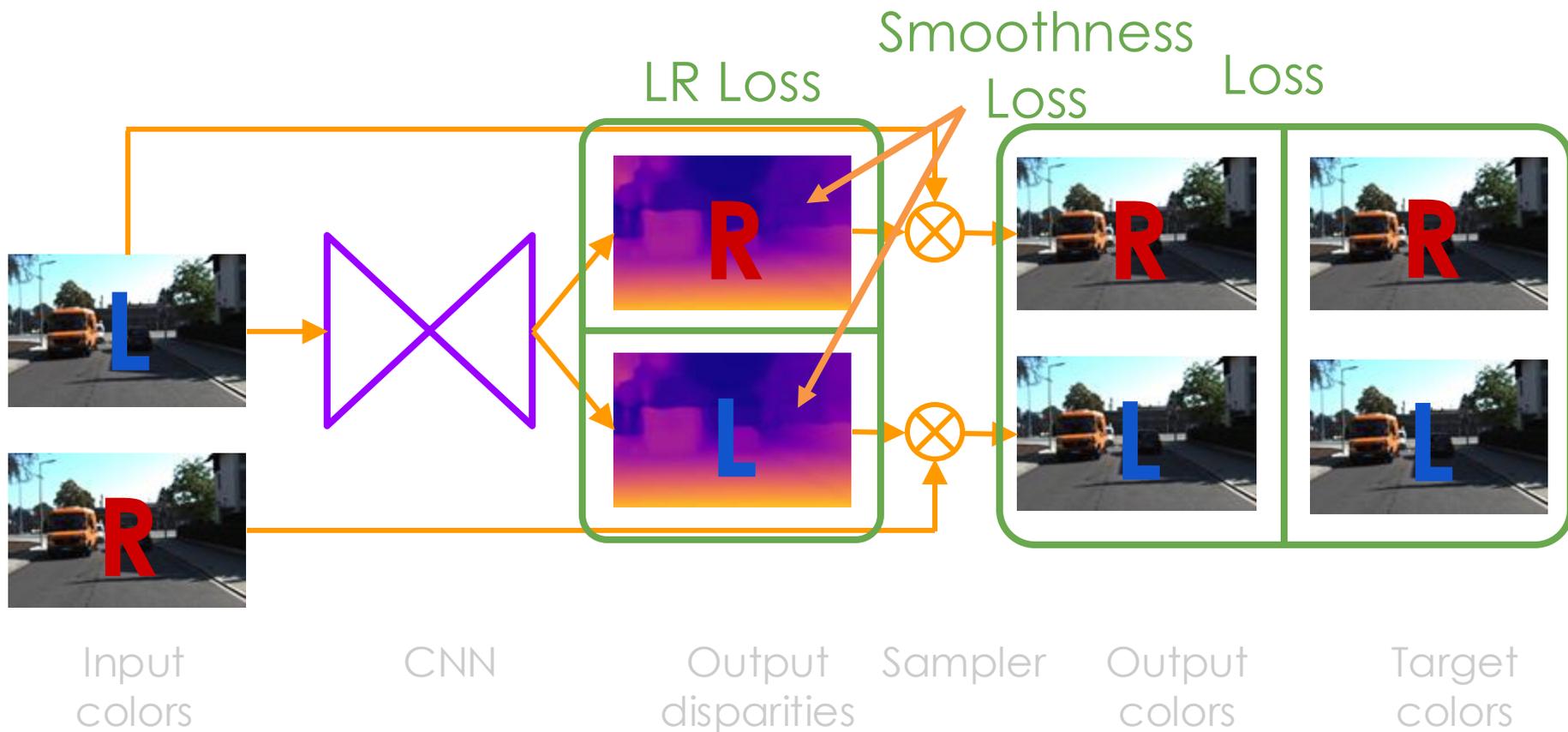


# This method



# Unsupervised depth estimation - **This method**

$$C_s = \alpha_{ap}(C_{ap}^l + C_{ap}^r) + \alpha_{ds}(C_{ds}^l + C_{ds}^r) + \alpha_{lr}(C_{lr}^l + C_{lr}^r).$$



# Reconstruction Loss

Complete Loss  $C_s = \alpha_{ap}(C_{ap}^l + C_{ap}^r) + \alpha_{ds}(C_{ds}^l + C_{ds}^r) + \alpha_{lr}(C_{lr}^l + C_{lr}^r)$ .

$N$  = # of pixels  
 $i, j$  = index of the pixels  
 $\alpha$  = weight

$$C_{ap}^l = \frac{1}{N} \sum_{i,j} \alpha \frac{1 - \text{SSIM}(I_{ij}^l, \tilde{I}_{ij}^l)}{2} + (1 - \alpha) \left\| I_{ij}^l - \tilde{I}_{ij}^l \right\|$$

Luminance   Contrast   Structure

SSIM: Structural Similarity Index =  $f(l(\mathbf{x}, \mathbf{y}), c(\mathbf{x}, \mathbf{y}), s(\mathbf{x}, \mathbf{y}))$

# Left-Right Consistency Loss

Complete Loss  $C_s = \alpha_{ap}(C_{ap}^l + C_{ap}^r) + \alpha_{ds}(C_{ds}^l + C_{ds}^r) + \alpha_{lr}(C_{lr}^l + C_{lr}^r)$

- $N$  = # of pixels  
 $i, j$  = index of the pixels  
 $d$  = disparity in left/right image

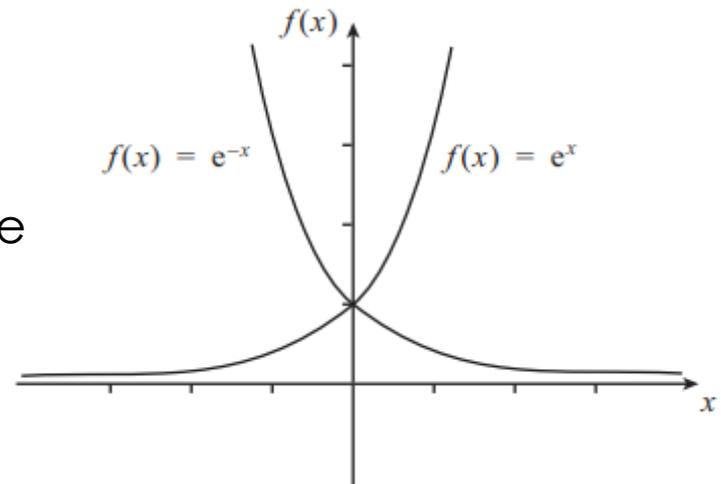
$$C_{lr}^l = \frac{1}{N} \sum_{i,j} \left| d_{ij}^l - d_{ij+d_{ij}^l}^r \right|$$

# Smoothness Loss

Complete Loss  $C_s = \alpha_{ap}(C_{ap}^l + C_{ap}^r) + \alpha_{ds}(C_{ds}^l + C_{ds}^r) + \alpha_{lr}(C_{lr}^l + C_{lr}^r)$ .

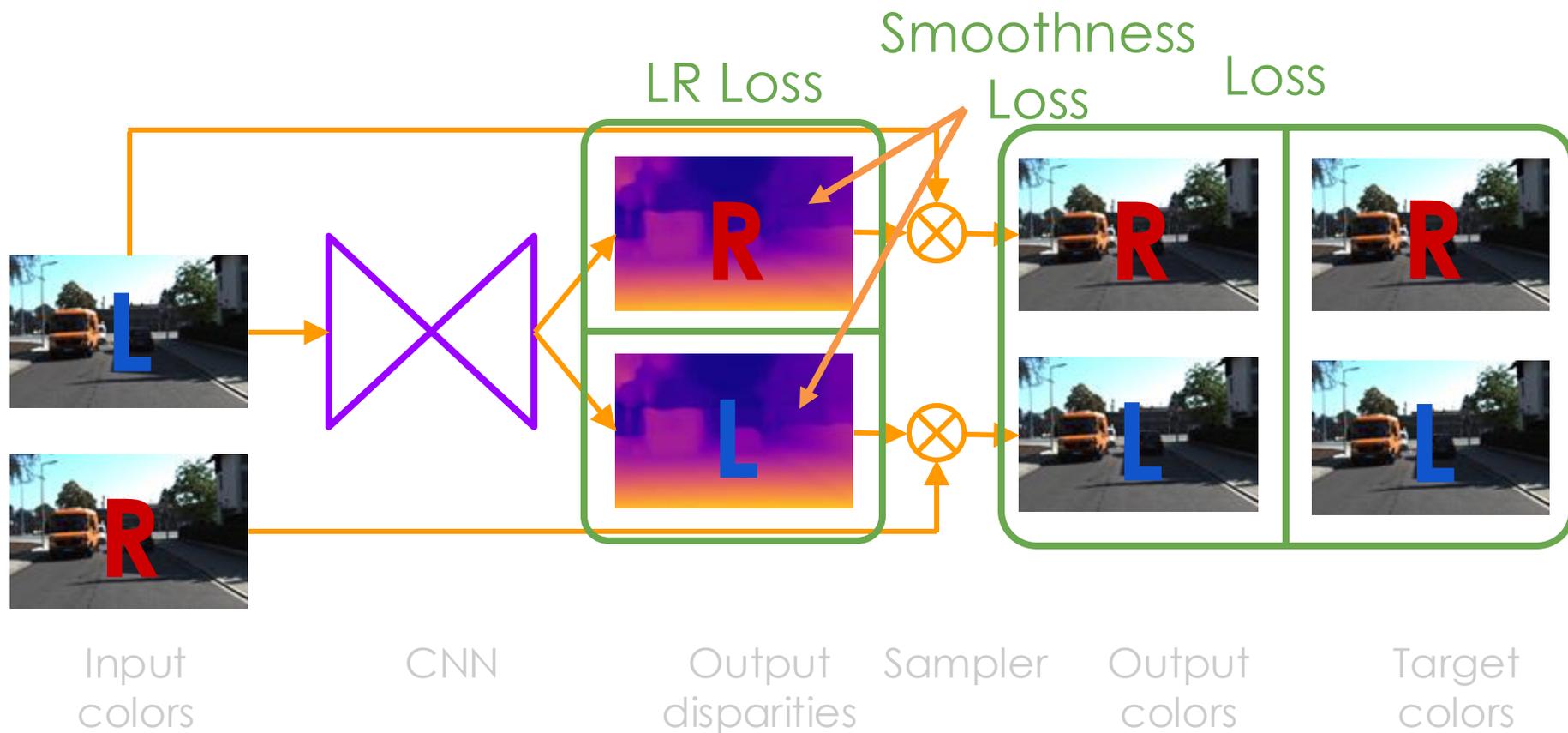
$$C_{ds}^l = \frac{1}{N} \sum_{i,j} |\partial_x d_{ij}^l| e^{-\|\partial_x I_{ij}^l\|} + |\partial_y d_{ij}^l| e^{-\|\partial_y I_{ij}^l\|}$$

- $N$  = # of pixels
- $i, j$  = index of the pixels
- $d$  = disparity in left/right image
- $\partial_x d_{i,j}^l$  = Gradient of disparity
- $\partial_x I_{i,j}^l$  = Image Gradient



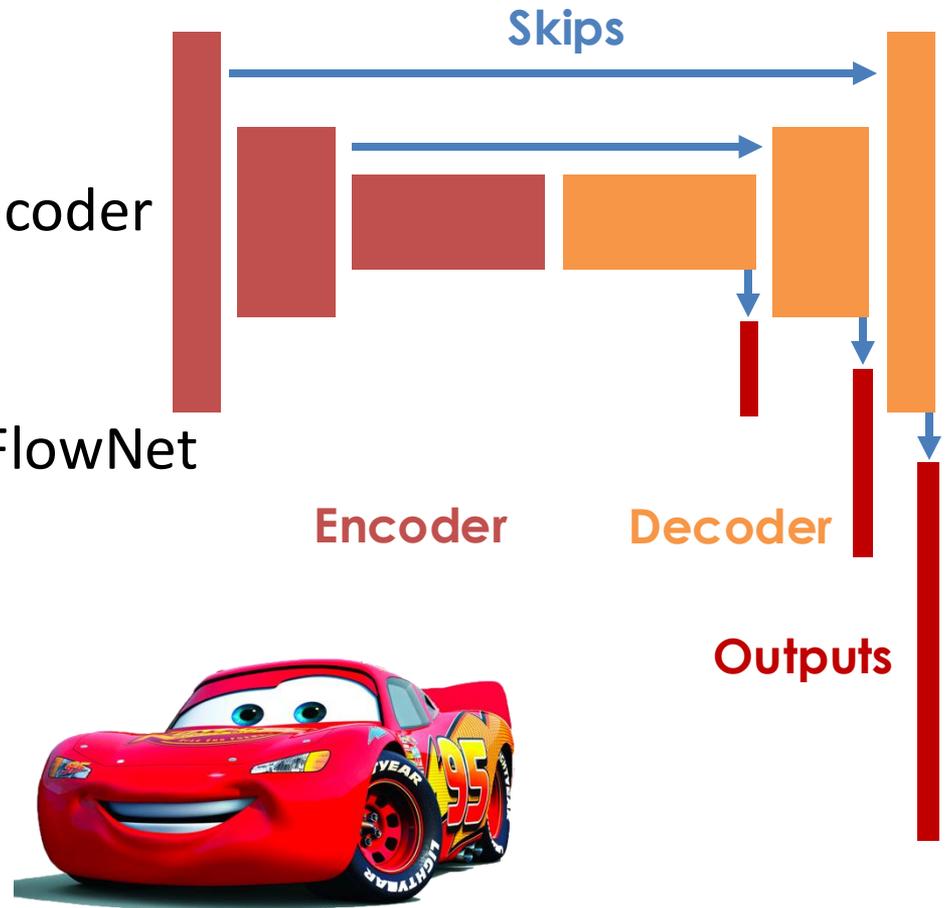
# Unsupervised depth estimation - **This method**

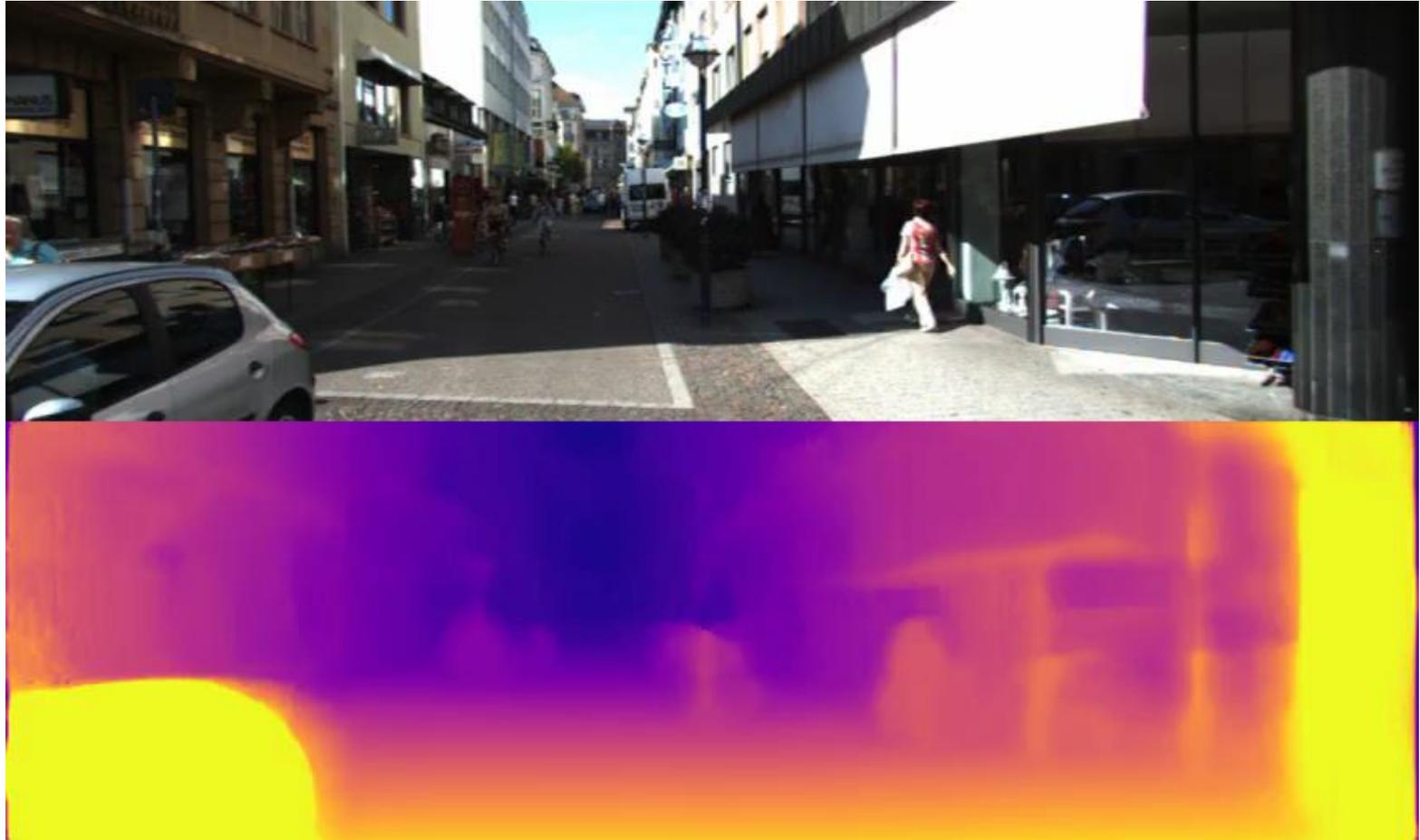
$$C_s = \alpha_{ap}(C_{ap}^l + C_{ap}^r) + \alpha_{ds}(C_{ds}^l + C_{ds}^r) + \alpha_{lr}(C_{lr}^l + C_{lr}^r).$$



# Architecture

- Fully convolutional
  - Choose your favorite encoder
- Skip connections
  - Similar to DispNet and FlowNet
- Multiscale generation
  - And Loss!
- Fast!
  - ~30fps on a Titan X

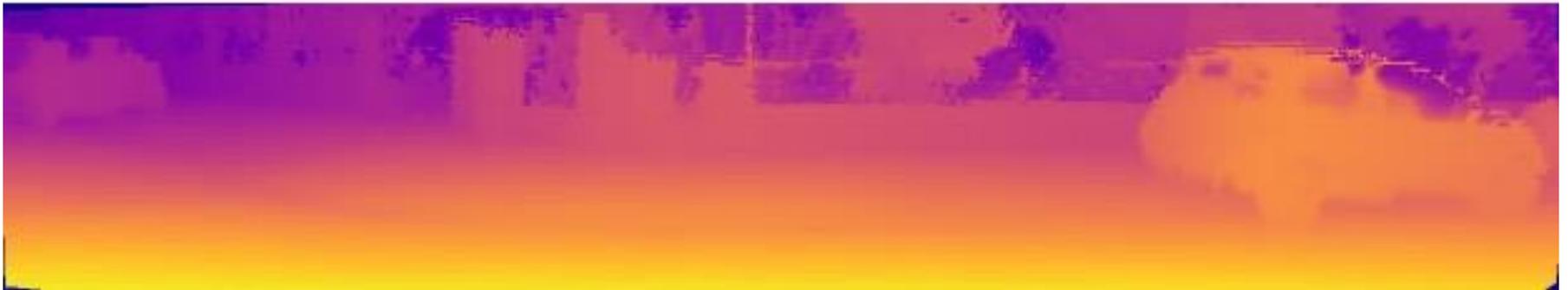




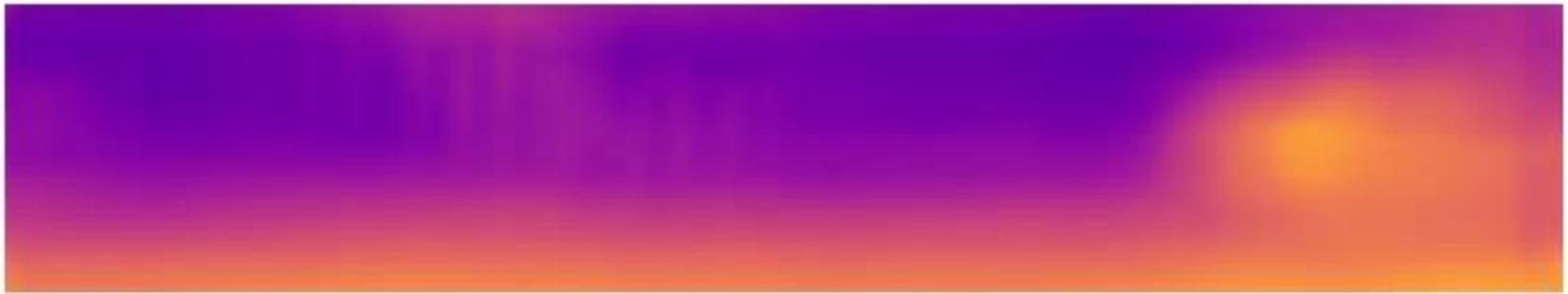
# KITTI – Input image



# KITTI – Ground truth depth



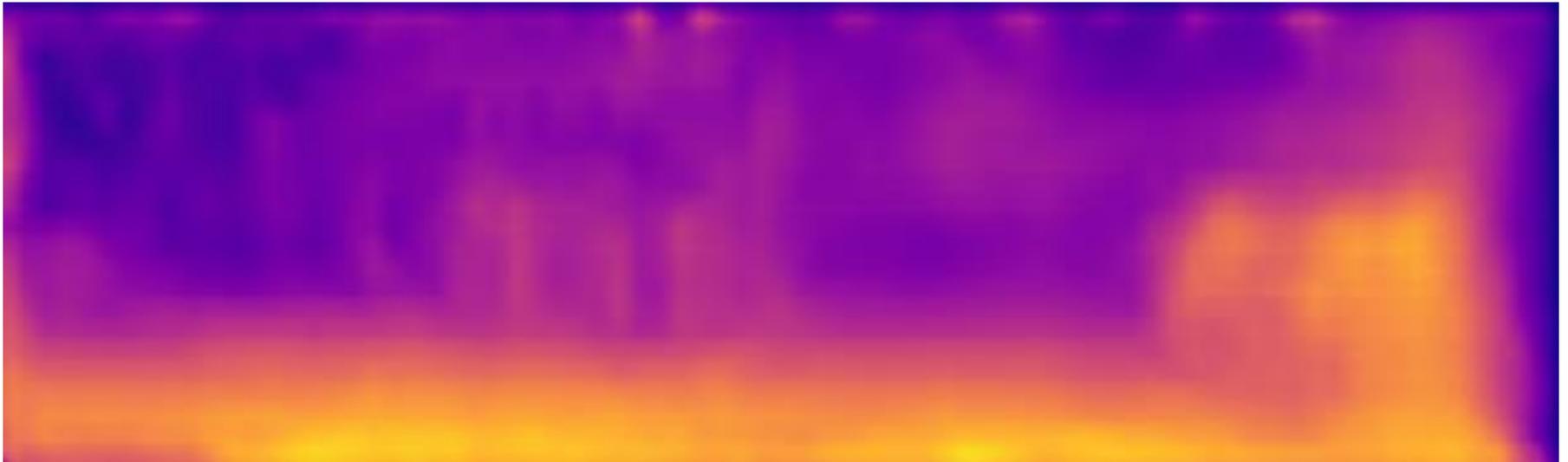
# KITTI – Eigen *et al.* [NIPS 14]



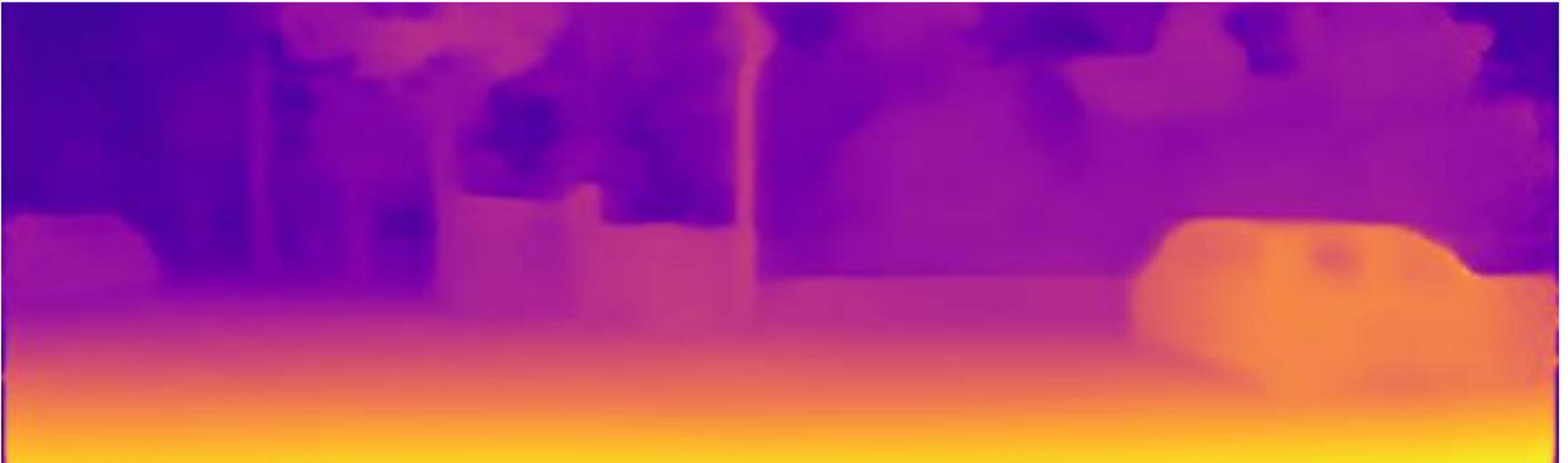
# KITTI – Liu *et al.* [CVPR 14]



# KITTI – Garg *et al.* [ECCV 16]



# KITTI – This work



# KITTI – Input image



# KITTI 2015

- All variants of our model beat previous supervised methods

Method	Supervised	Dataset	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Train set mean	No	K	0.361	4.826	8.102	0.377	0.638	0.804	0.894
Eigen et al. [10] Coarse °	Yes	K	0.214	1.605	6.563	0.292	0.673	0.884	0.957
Eigen et al. [10] Fine °	Yes	K	0.203	1.548	6.307	0.282	0.702	0.890	0.958
Liu et al. [36] DCNF-FCSP FT *	Yes	K	0.201	1.584	6.471	0.273	0.68	0.898	0.967
<b>Ours No LR</b>	No	K	0.152	1.528	6.098	0.252	0.801	0.922	0.963
<b>Ours</b>	No	K	0.148	1.344	5.927	0.247	0.803	0.922	0.964
<b>Ours</b>	No	CS + K	0.124	1.076	5.311	0.219	0.847	0.942	0.973
<b>Ours pp</b>	No	CS + K	0.118	0.923	5.015	0.210	0.854	0.947	<b>0.976</b>
<b>Ours resnet pp</b>	No	CS + K	<b>0.114</b>	<b>0.898</b>	<b>4.935</b>	<b>0.206</b>	<b>0.861</b>	<b>0.949</b>	<b>0.976</b>

# KITTI 2015

- All variants of our model beat the previous unsupervised method

Method	Supervised	Dataset	Abs Rel	Sq Rel	RMSE	RMSE log	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
Train set mean	No	K	0.361	4.826	8.102	0.377	0.638	0.804	0.894
Eigen et al. [10] Coarse °	Yes	K	0.214	1.605	6.563	0.292	0.673	0.884	0.957
Eigen et al. [10] Fine °	Yes	K	0.203	1.548	6.307	0.282	0.702	0.890	0.958
Liu et al. [36] DCNF-FCSP FT *	Yes	K	0.201	1.584	6.471	0.273	0.68	0.898	0.967
<b>Ours No LR</b>	No	K	0.152	1.528	6.098	0.252	0.801	0.922	0.963
<b>Ours</b>	No	K	0.148	1.344	5.927	0.247	0.803	0.922	0.964
<b>Ours</b>	No	CS + K	0.124	1.076	5.311	0.219	0.847	0.942	0.973
<b>Ours pp</b>	No	CS + K	0.118	0.923	5.015	0.210	0.854	0.947	<b>0.976</b>
<b>Ours resnet pp</b>	No	CS + K	<b>0.114</b>	<b>0.898</b>	<b>4.935</b>	<b>0.206</b>	<b>0.861</b>	<b>0.949</b>	<b>0.976</b>
Garg et al. [16] L12 Aug 8× cap 50m	No	K	0.169	1.080	5.104	0.273	0.740	0.904	0.962
<b>Ours cap 50m</b>	No	K	0.140	0.976	4.471	0.232	0.818	0.931	0.969
<b>Ours cap 50m</b>	No	CS + K	0.117	0.762	3.972	0.206	0.860	0.948	0.976
<b>Ours pp cap 50m</b>	No	CS + K	0.112	0.680	3.810	0.198	0.866	0.953	<b>0.979</b>
<b>Ours resnet pp cap 50m</b>	No	CS + K	<b>0.108</b>	<b>0.657</b>	<b>3.729</b>	<b>0.194</b>	<b>0.873</b>	<b>0.954</b>	<b>0.979</b>

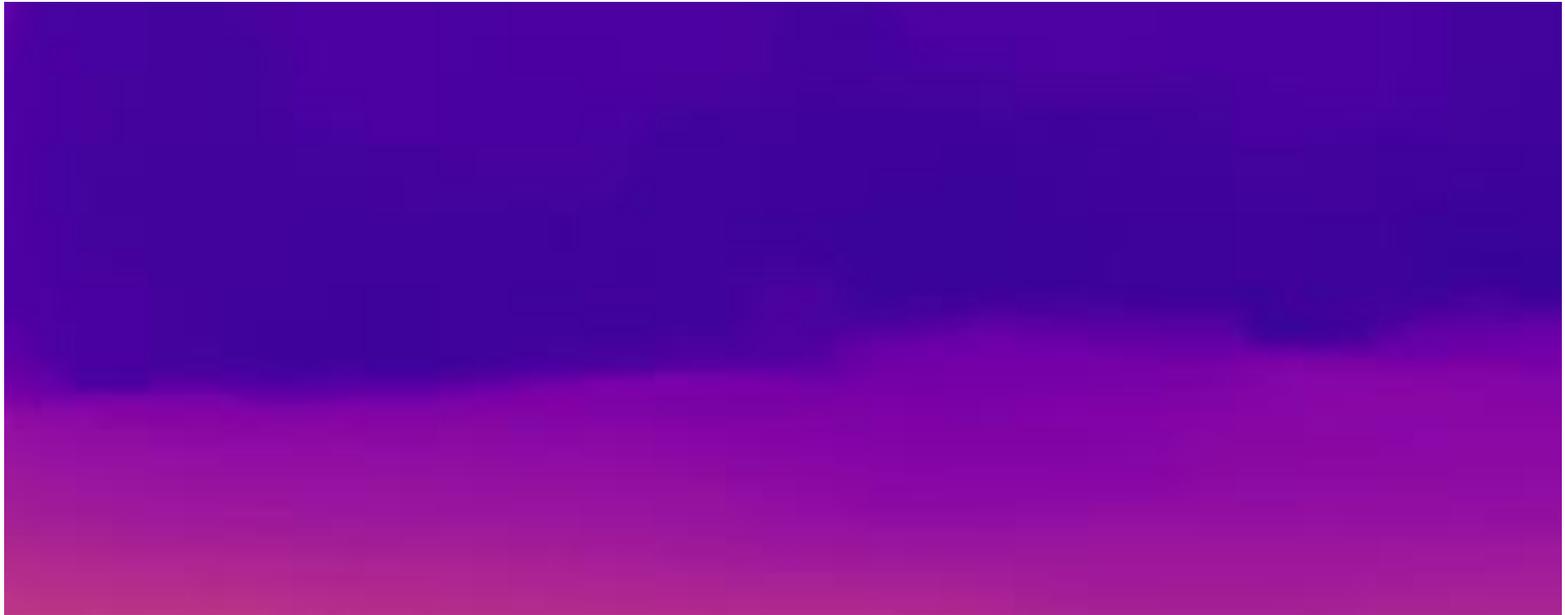
# Make3D – Input image



# Make3D – Ground truth depth



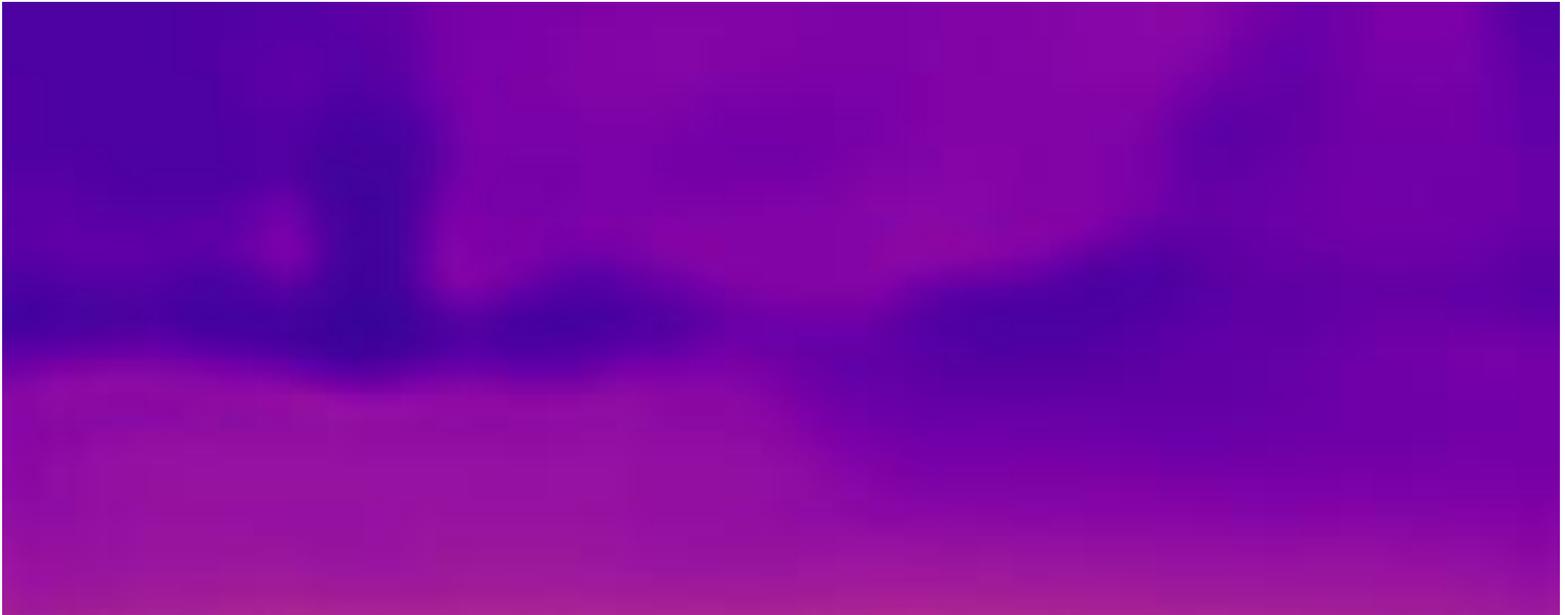
# Make3D – Karsch *et al.* [PAMI 14]



# Make3D – Liu *et al.* [CVPR 14]



# Make3D – Laina *et al.* [3DV 16]



# Make3D – This method



# Make3D – Input image





# Challenges

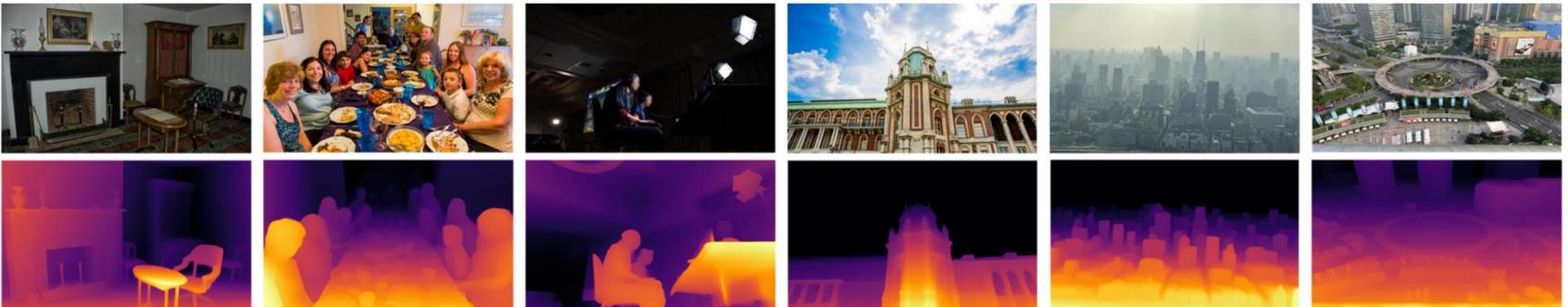
- Reprojection loss
  - Assumes lambertian world
  - More supervision?
    - Kuznietsov *et al.* [CVPR 17]
- Need calibrated data
  - Synced and rectified
  - Less supervision?
    - Zhou *et al.* [CVPR 2017]

# Conclusion

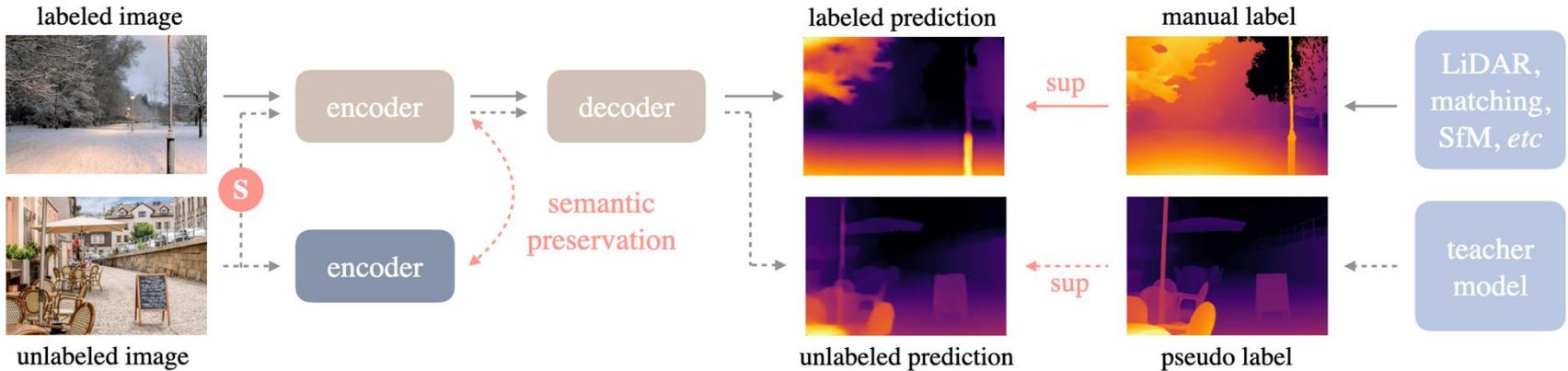
- We can get depth from a single photograph
- Self-supervision with stereo data
  - Cheap and scalable!
- Accurate
  - Beats fully-supervised methods on KITTI!

# Depth Anything: Unleashing the Power of Large-Scale Unlabeled Data. Yang et al. CVPR 2024.

- What is different?
  - Mix of labeled and unlabeled datasets
  - Teacher/student network
  - Encoder with strong semantic prior



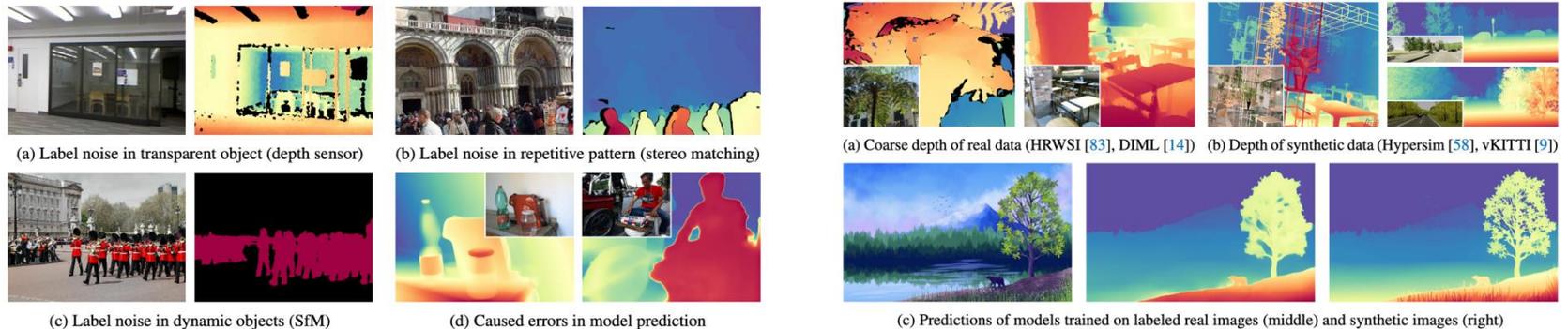
# Depth Anything



- 1.5M labeled images from multiple datasets
- 65M unlabeled images -> get pseudo ground truth from teacher network
- Train student network from scratch on pseudo ground truth
- Iterate: re-label images with new Teacher, train student
- Introduce strong perturbations (color distortions, color jittering, Gaussian blurring, CutMix)
- Semantic feature alignment based on Dino but with some slack

# Depth Anything V2. Yang et al.

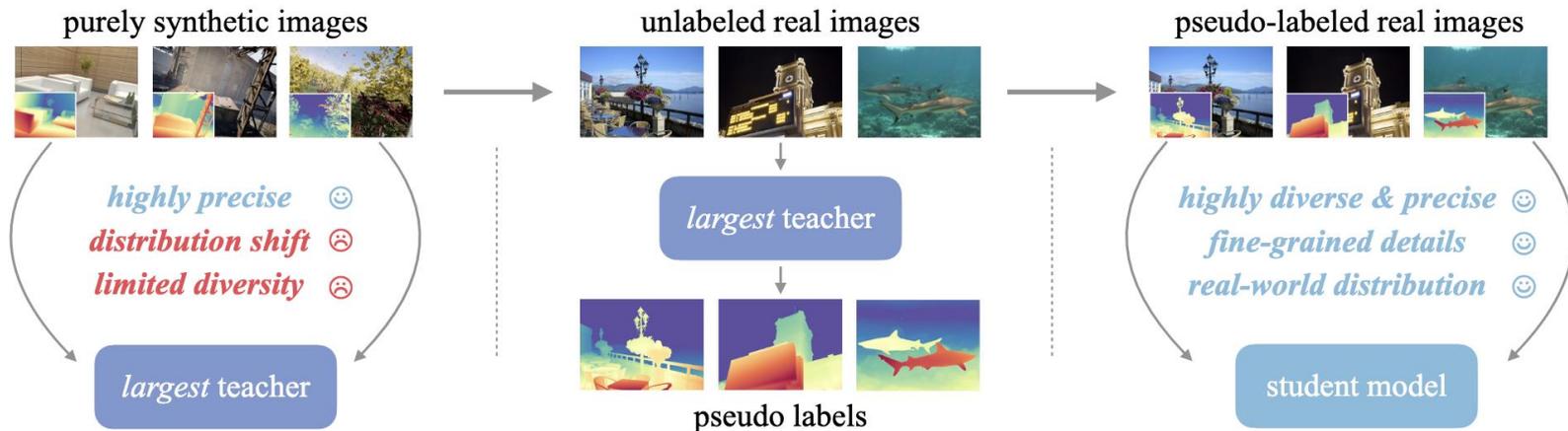
## NeurIPS 2024



- Teacher trained on synthetic images to mitigate label noise in real datasets
- Student trained on real images to bridge sim2real gap

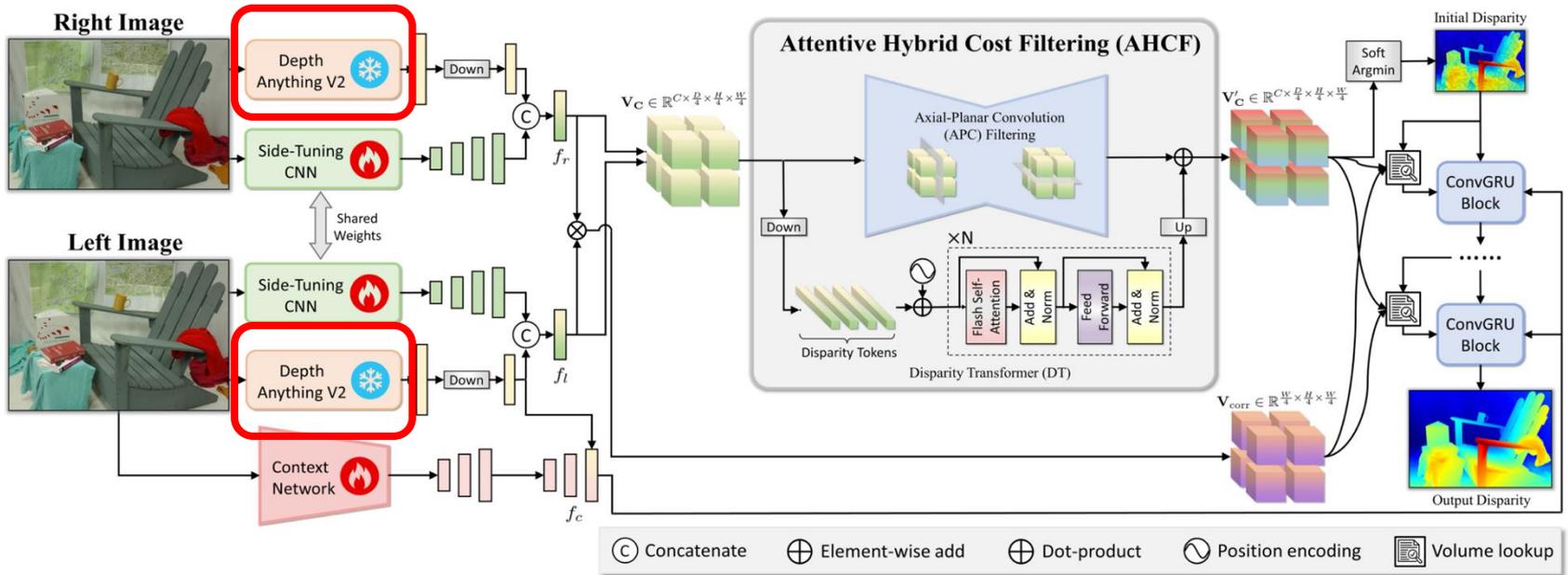
# Depth Anything V2. Yang et al.

## NeurIPS 2024



- Teacher trained on synthetic images to mitigate label noise in real datasets
- Student trained on real images to bridge sim2real gap

# Foundation Stereo

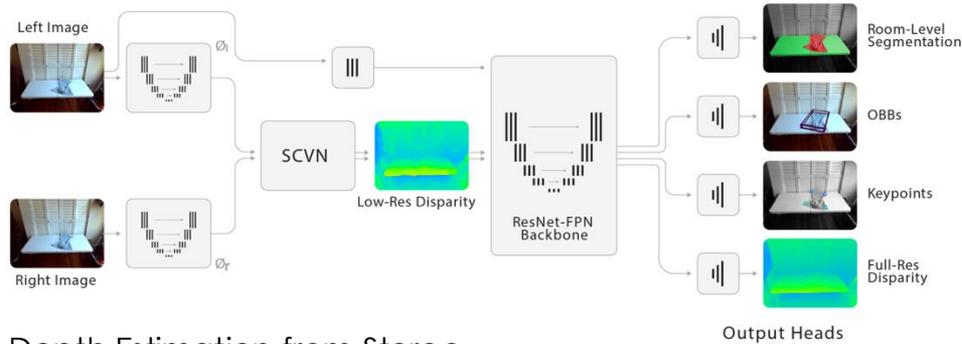


$$\mathbf{V}_{\text{gwc}}(g, d, h, w) = \left\langle \widehat{f}_{l,g}^{(4)}(h, w), \widehat{f}_{r,g}^{(4)}(h, w - d) \right\rangle,$$

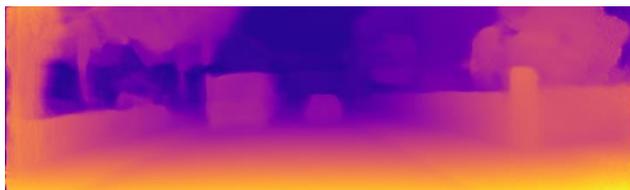
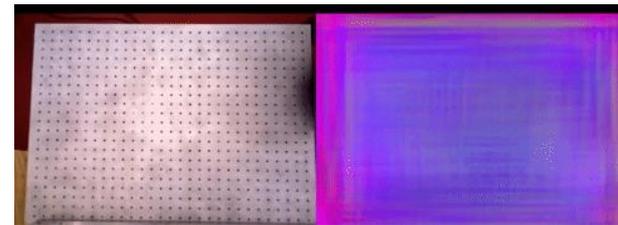
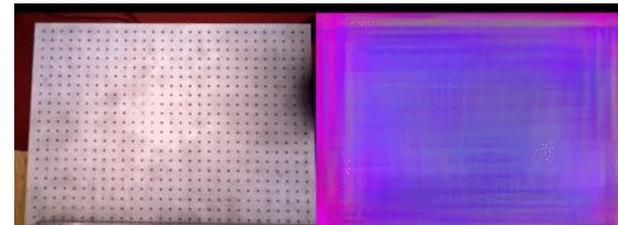
$$\mathbf{V}_{\text{cat}}(d, h, w) = \left[ \text{Conv}(f_l^{(4)})(h, w), \text{Conv}(f_r^{(4)})(h, w - d) \right]$$

$$\mathbf{V}_{\mathbf{C}}(d, h, w) = [\mathbf{V}_{\text{gwc}}(d, h, w), \mathbf{V}_{\text{cat}}(d, h, w)] \quad (1)$$

# Let's use representation learning!



Depth Estimation from Stereo  
Supervised Learning



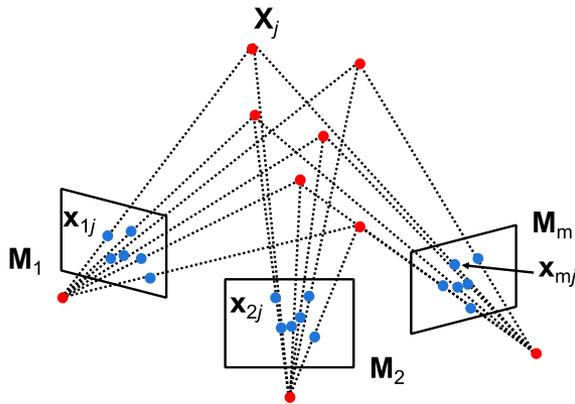
Monocular Depth Estimation  
Unsupervised Learning



Image by Yunuk Cha.  
Finding Correspondences across  
Frames  
Self-Supervised Learning

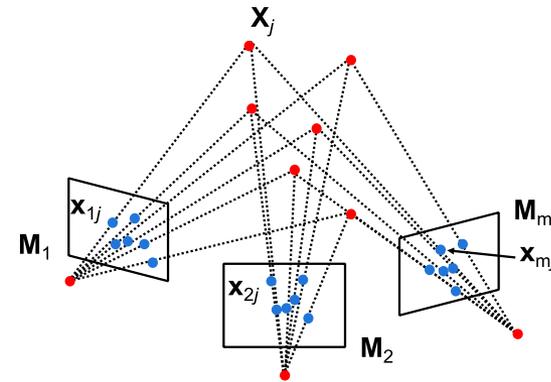
# Feature Tracking

## Structure From Motion Problem



Given  $m$  images of  ~~$n$  fixed~~ 3D points

$$\bullet x_{ij} = M_i X_j, \quad i = 1, \dots, m, \quad j = 1, \dots, n$$



From the  $m \times n$  observations  $x_{ij}$ , estimate:

- |   |           |
|---|-----------|
| <ul style="list-style-type: none"><li>• <math>m</math> projection matrices <math>M_i</math></li></ul> | motion    |
| <ul style="list-style-type: none"><li>• <math>n</math> 3D points <math>X_j</math></li></ul>           | structure |

# Problem statement

## Image sequence



Slide credit: Yonsei Univ.

Slides Adapted from CS131a.

# Problem statement

## Feature point detection

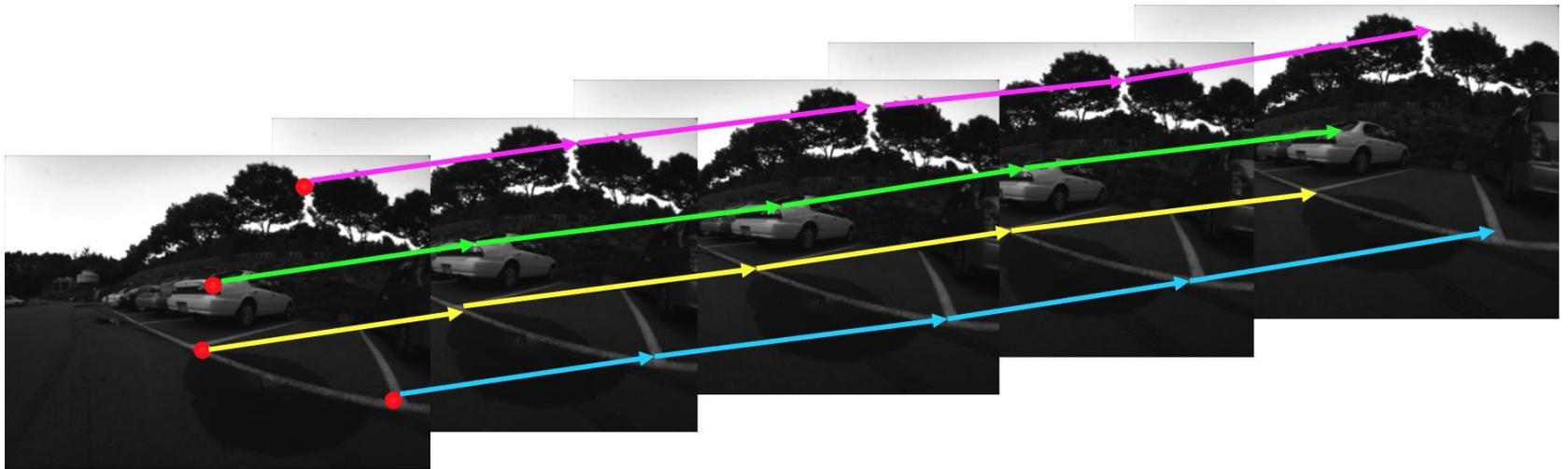


Slide credit: Yonsei Univ.

Slides Adapted from CS131a.

# Problem statement

## Feature point tracking



Slide credit: Yonsei Univ.

Slides Adapted from CS131a.

# Single object tracking



Slides Adapted from CS131a.

# Multiple object tracking



Slides Adapted from CS131a.

# Tracking with a fixed camera



Slides Adapted from CS131a.

# Tracking with a moving camera



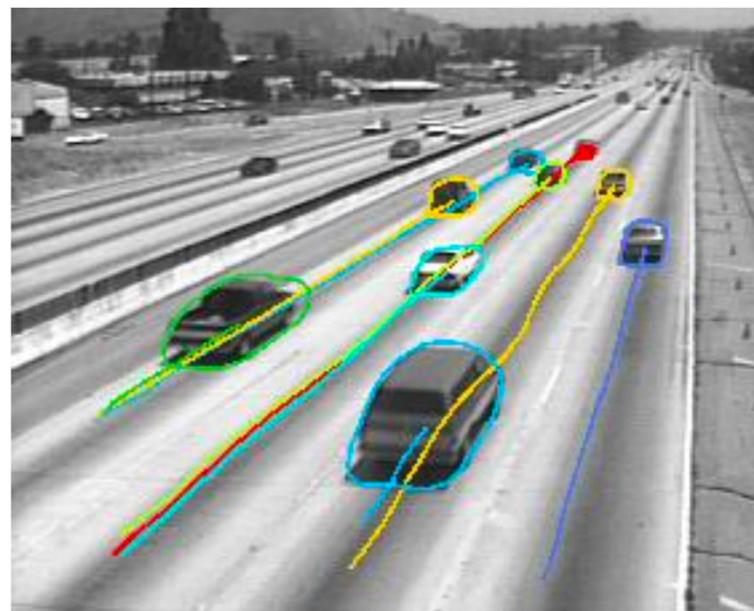
Slides Adapted from CS131a.

# Challenges in Feature Tracking

- Figure out which features can be tracked
  - Efficiently track across frames
- Some points may change appearance over time
  - e.g., due to rotation, moving into shadows, etc.
- Drift: small errors can accumulate as appearance model is updated
- Points may appear or disappear.
  - need to be able to add/delete tracked points.

# What are good features to track?

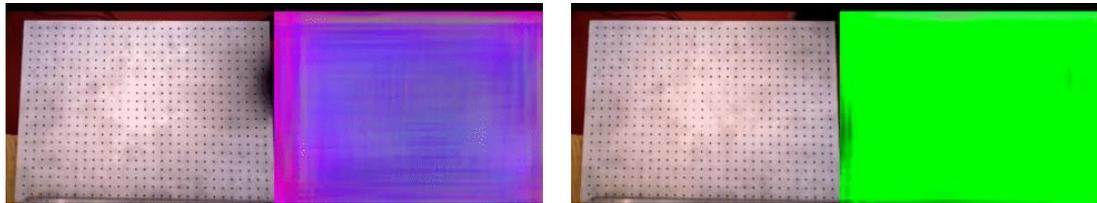
- Regions we can track easily and consistently
- Once we have the good features, we can use simple tracking methods
- Next Lecture: Optical Flow



# Dense Object Nets

Learning *Dense* Visual Object *Descriptors*  
*By and For* Robotic Manipulation. CORL 2018

Peter R. Florence, Lucas Manuelli, Russ Tedrake

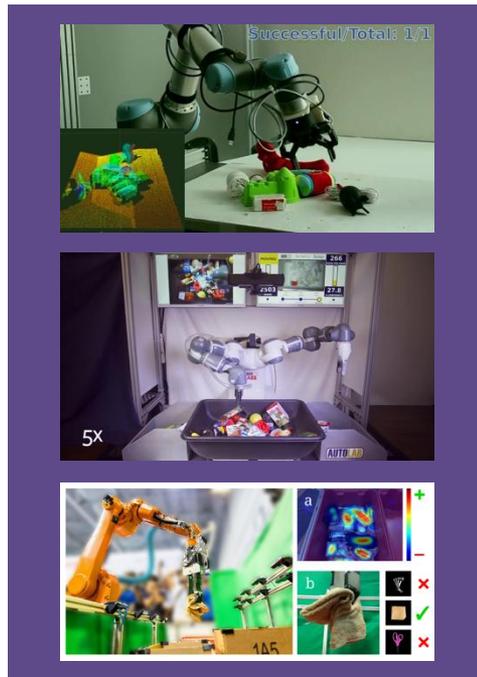


*Slides adapted from CS326 by Kevin Zakka and Sriram Somasundaram*

# Motivation



RL  
task-specific



Grasp feature-learning  
no task-specificity



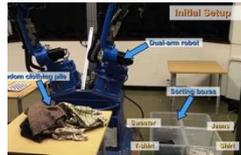
Grasp segmentation  
coarse  
no task-specificity

What is the right **object representation** for manipulation,  
and how can we **scalably** acquire it?

# Wish List



Deformable



Task agnostic

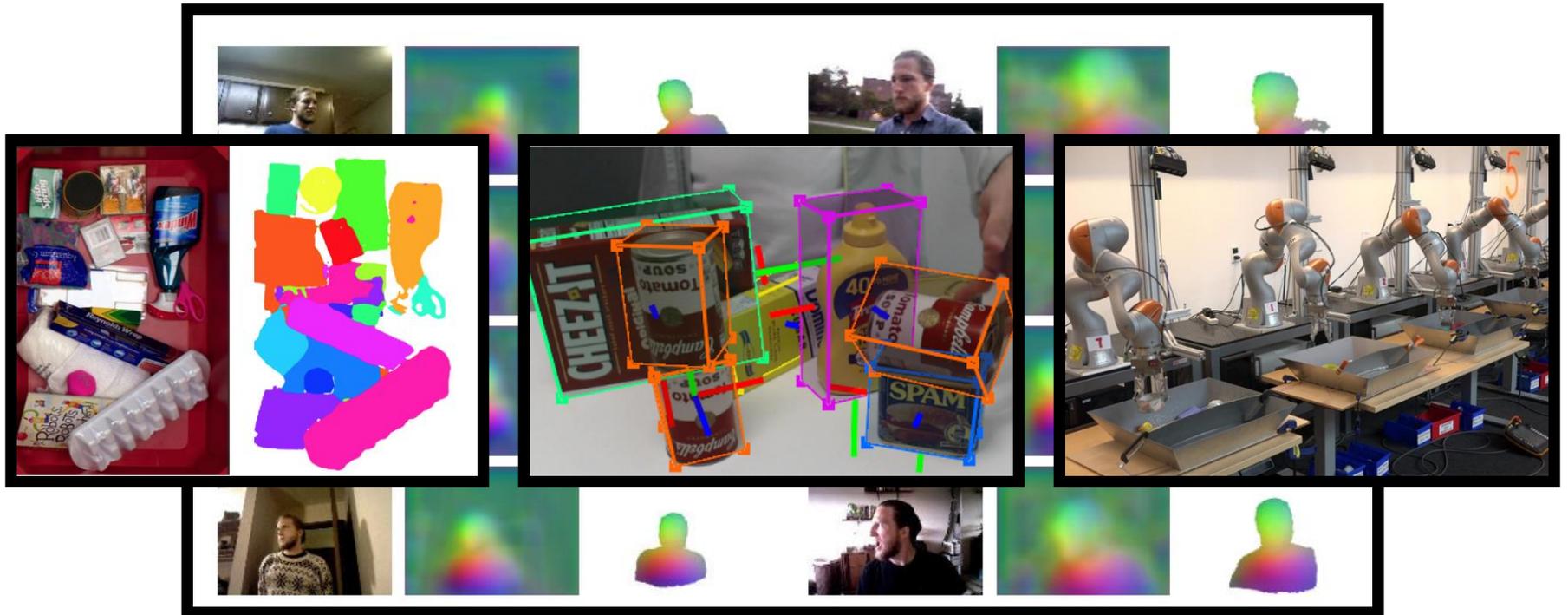


Self-supervised



3D perception

# Some Representations

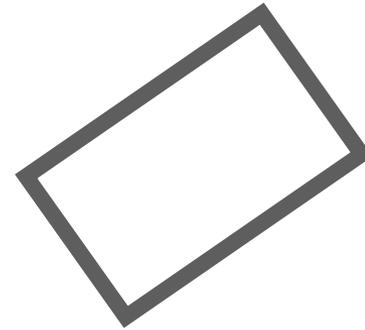
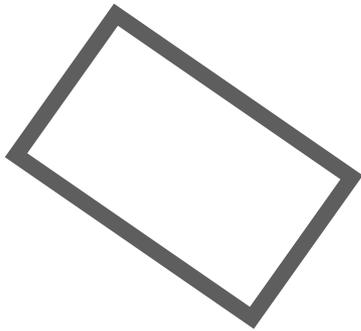


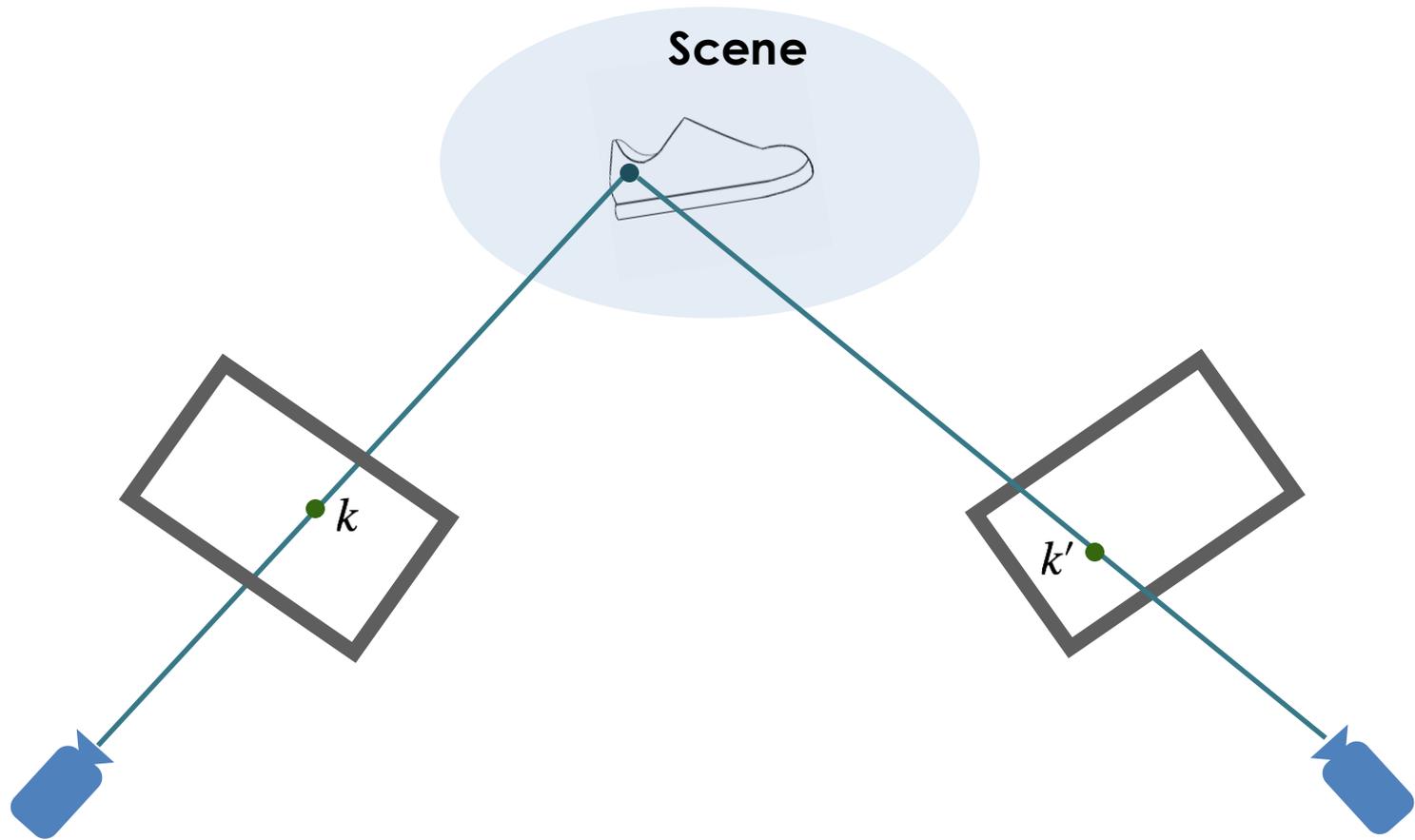
What exactly is a **descriptor**?

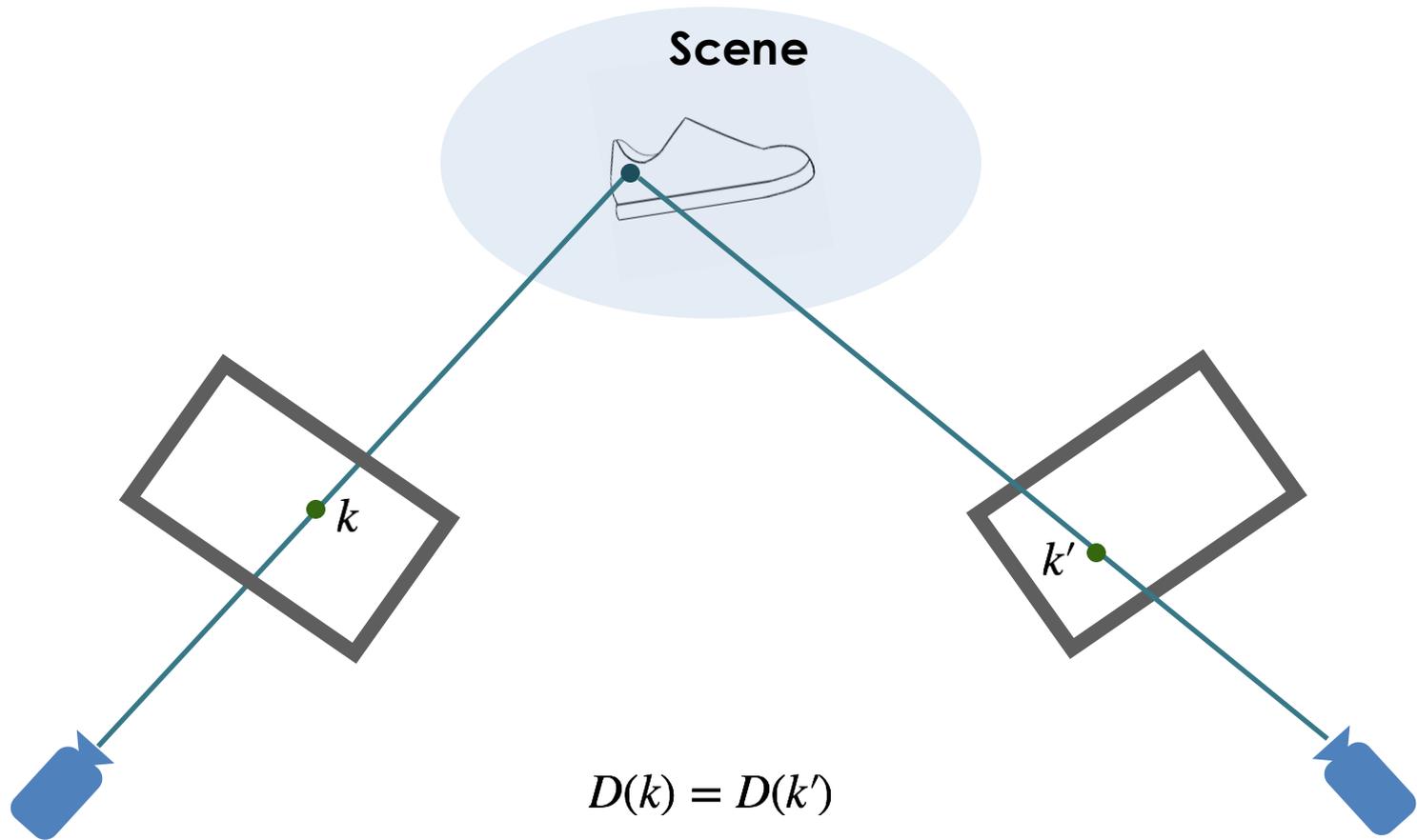


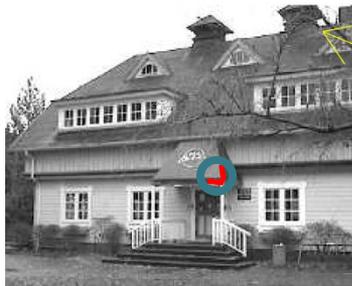
*The story goes that Takeo Kanade once told a young graduate student that the three most important problems in computer vision are: “correspondence, correspondence, correspondence!” - Wang et. al. 2019*

**Scene**

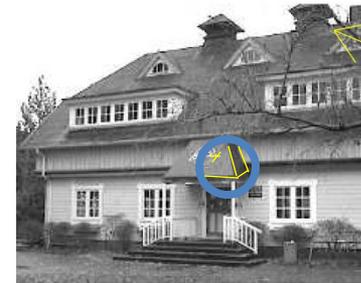








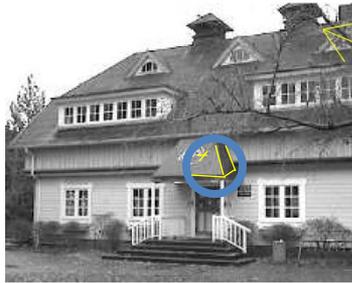
Feature detector



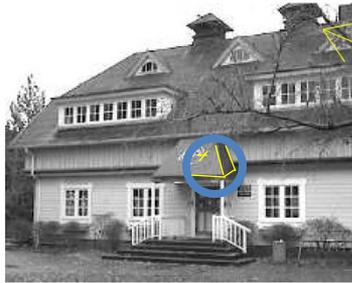
Feature descriptor

$$I \longrightarrow \{I[x_1, y_1], I[x_2, y_2], \dots\}$$

$$\text{Area around pixel } k \longrightarrow D(k)$$



Area around pixel  $k$   $\longrightarrow$   $D(k)$



Area around pixel  $k$   $\longrightarrow$   $D(k)$

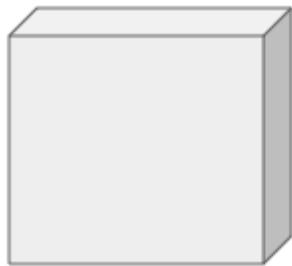
Area around pixel  $k$   $\longrightarrow$   $D(k)$

Features descriptors should be invariant under transformation

# Paper Overview

# Dense Descriptors

Input is an RGB image

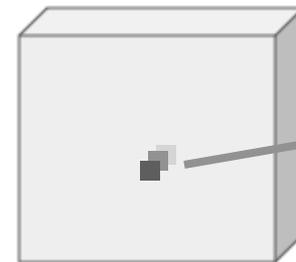


$$\mathbb{R}^{W \times H \times 3}$$

$f(\cdot)$



Output



D-dim descriptor  
for each pixel

$$\mathbb{R}^{W \times H \times D}$$

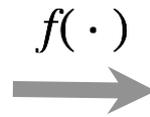
**Pay attention to the difference in Dimensionality**

# Dense Descriptors

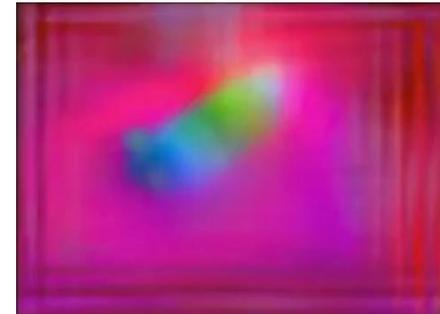
Input is an RGB image



$$\mathbb{R}^{W \times H \times 3}$$

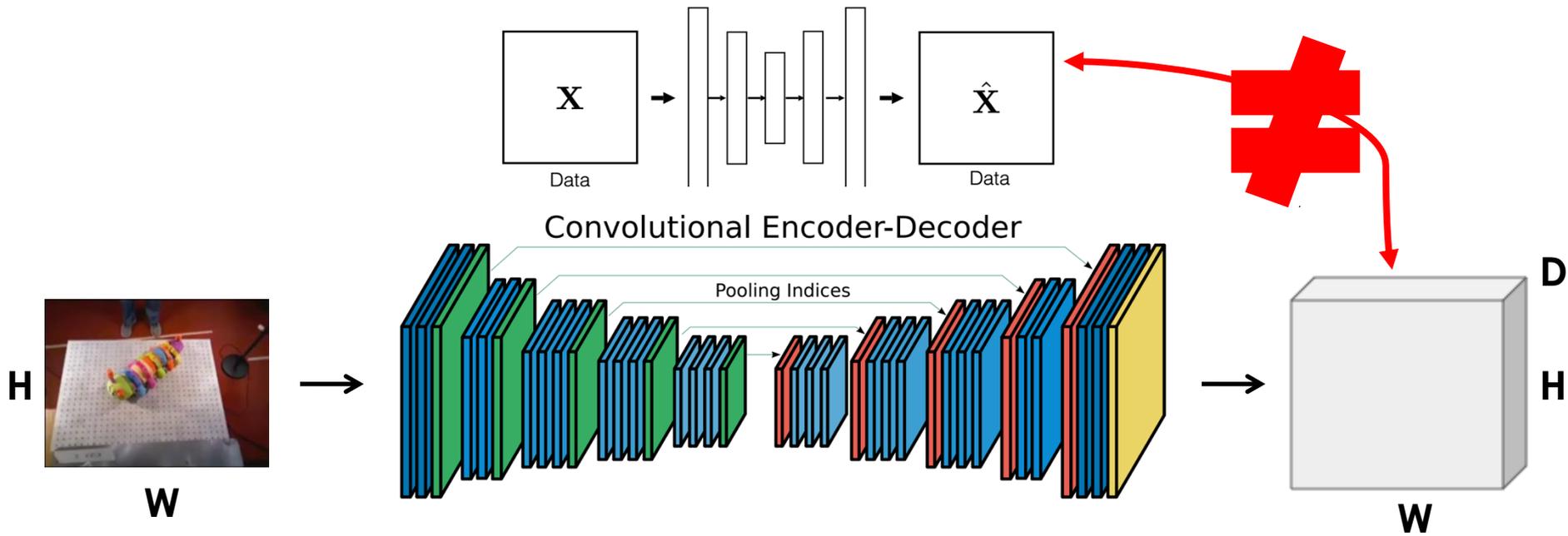


Output

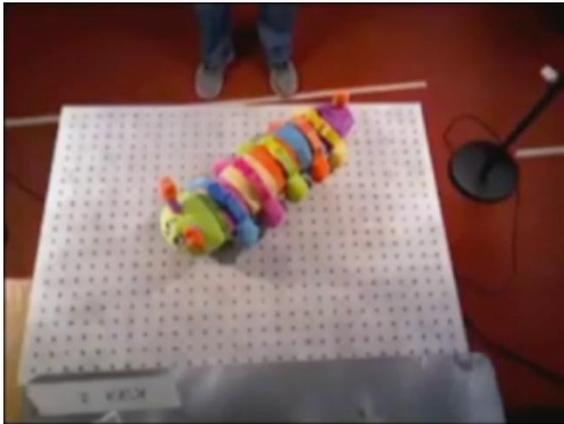


$$\mathbb{R}^{W \times H \times D}$$

# Network Architecture



# Pixelwise Contrastive Loss



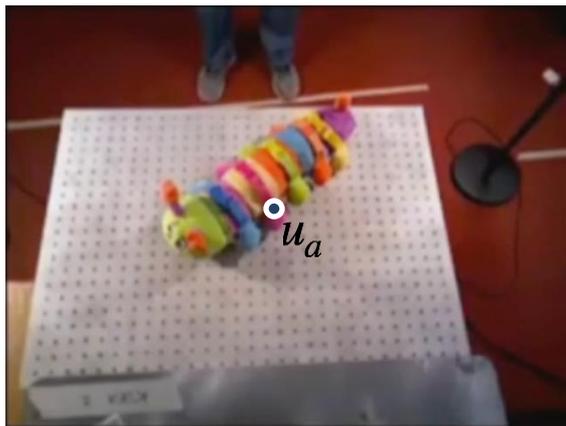
$I_a$



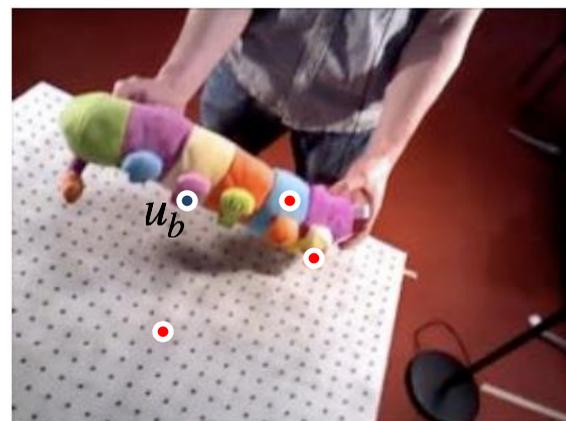
$I_b$

Hadsell et al., CVPR 2006

# Pixelwise Contrastive Loss



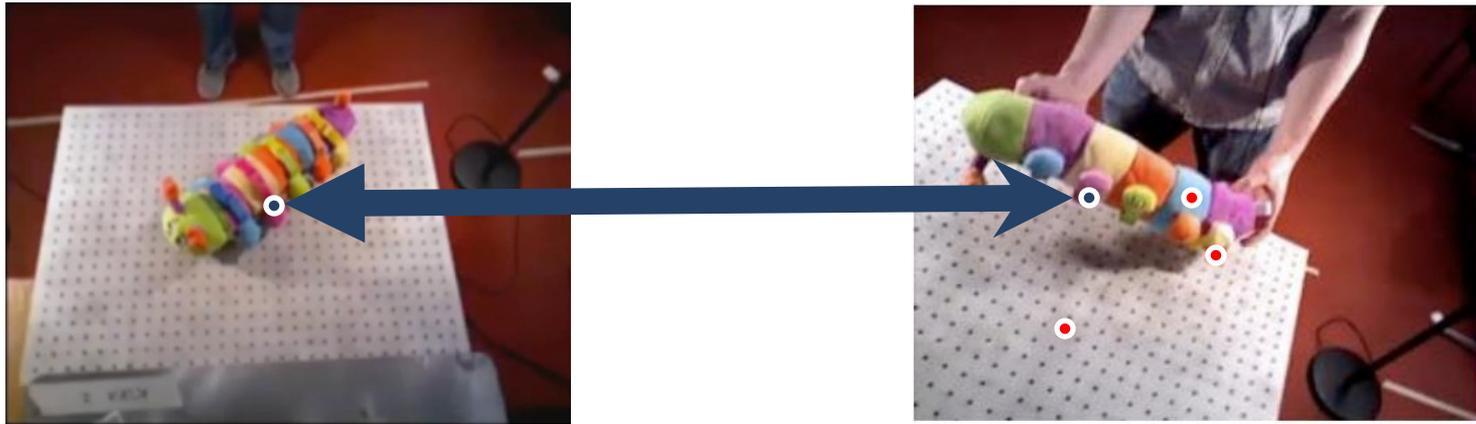
$I_a$



$I_b$

**Assumption: Ground truth Correspondences Given**

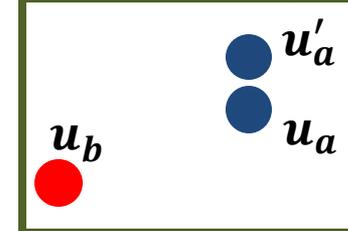
# Pixelwise Contrastive Loss - Matches



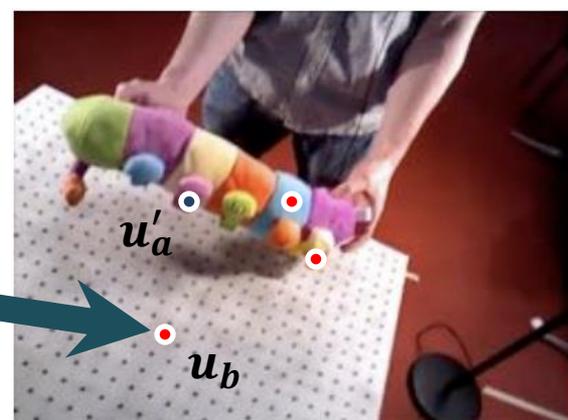
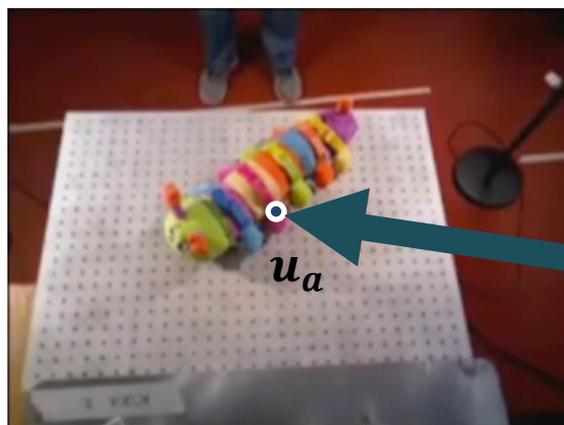
$$L_{\text{matches}}(I_a, I_b) = \frac{1}{N_{\text{matches}}} \sum_{N_{\text{matches}}} \underbrace{D(I_a, u_a, I_b, u_b)^2}_{\text{Distance in Descriptor Space}}$$

# Pixelwise Contrastive Loss – Non-Matches

Descriptor Space



- Training time

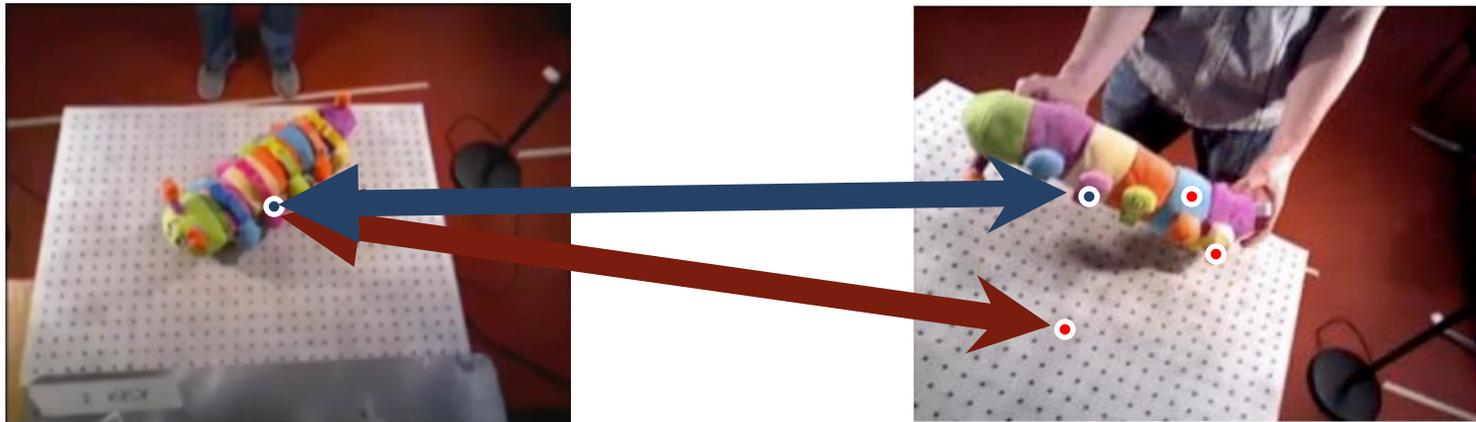


$$L_{\text{non-matches}}(I_a, I_b) = \frac{1}{N_{\text{non-matches}}} \sum_{N_{\text{non-matches}}} \max(0, M - \underbrace{D(I_a, u_a, I_b, u_b)}_{\text{Distance in Descriptor Space}})^2$$

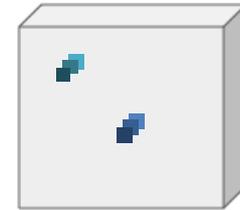
**Max distance**

You want this to be large for non-matches

# Pixelwise Contrastive Loss



$$L(I_a, I_b) = L_{\text{matches}}(I_a, I_b) + L_{\text{non-matches}}(I_a, I_b)$$

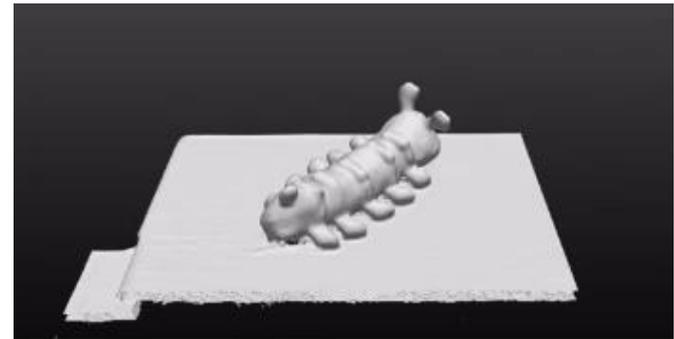
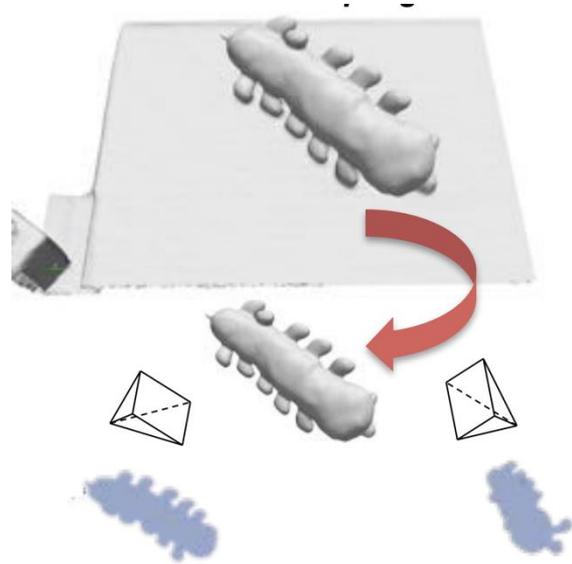


$\mathbb{R}^{W \times H \times D}$

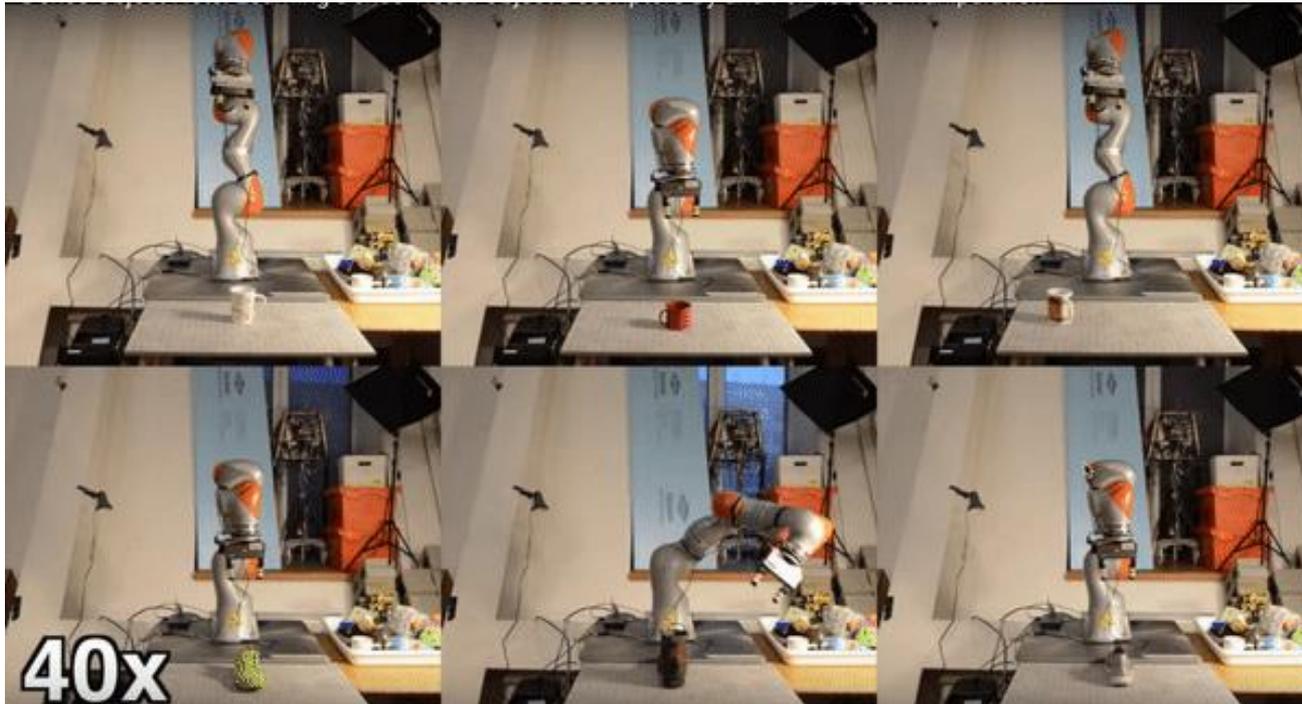
**Loss: Reconstruction + Contrastive Loss**

How can we generate these ground-truth correspondences?

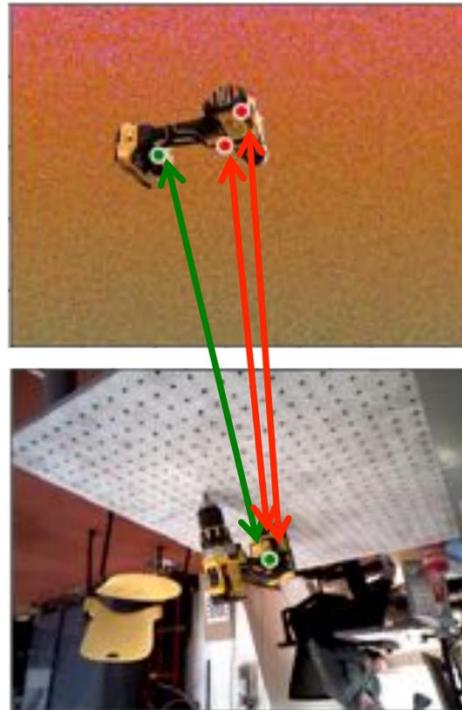
# 3D Reconstruction based Change Detection and Masked Sampling



# Autonomous and Self-supervised Data Collection



# Background Randomization



# Further Training Techniques

- Hard-Negative Scaling

$$L_{\text{non-matches}}(I_a, I_b) = \frac{1}{N_{\text{hard-negatives}}} \sum_{N_{\text{non-matches}}} \max(0, M - D(I_a, u_a, I_b, u_b)^2)$$

- Data Augmentation

# Results

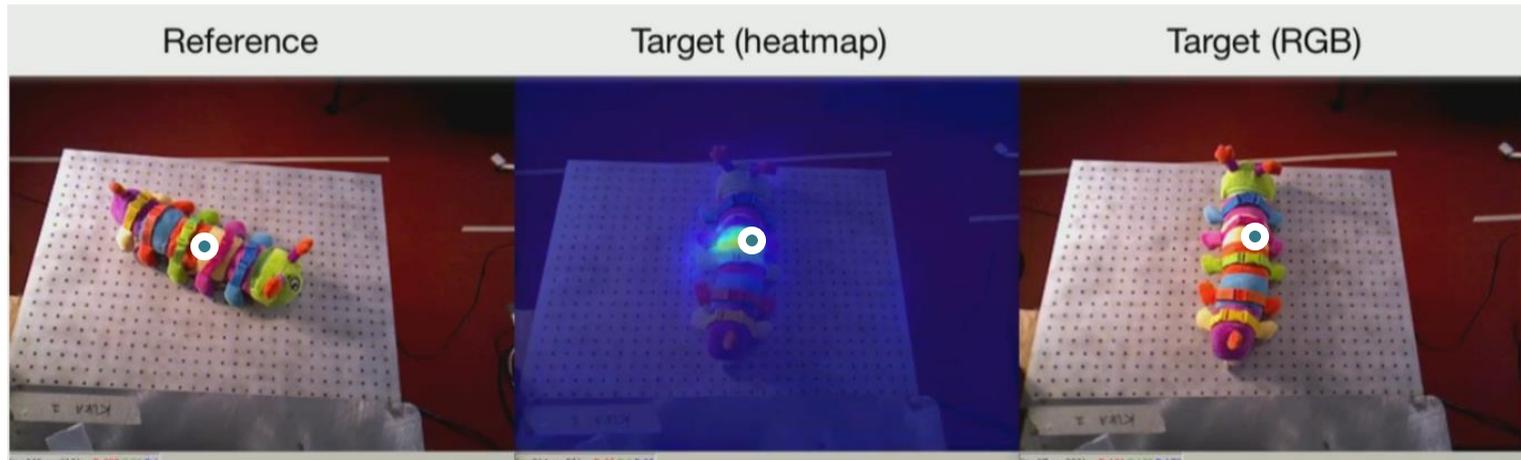
# Experiments

- 1) Can we acquire descriptors that can generalize across classes of objects and can be distinct for each object instance?
- 2) Can we apply dense descriptors to robotic manipulation?

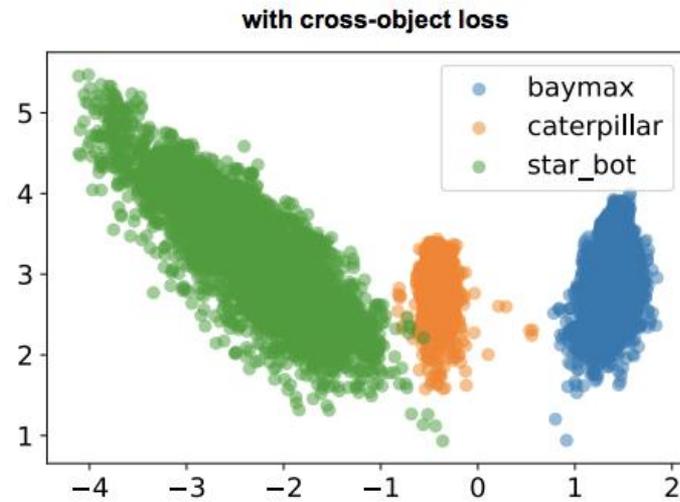
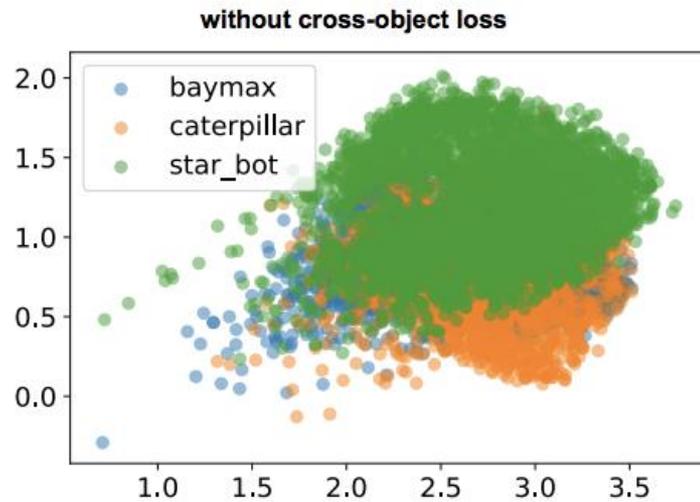
# Single Object



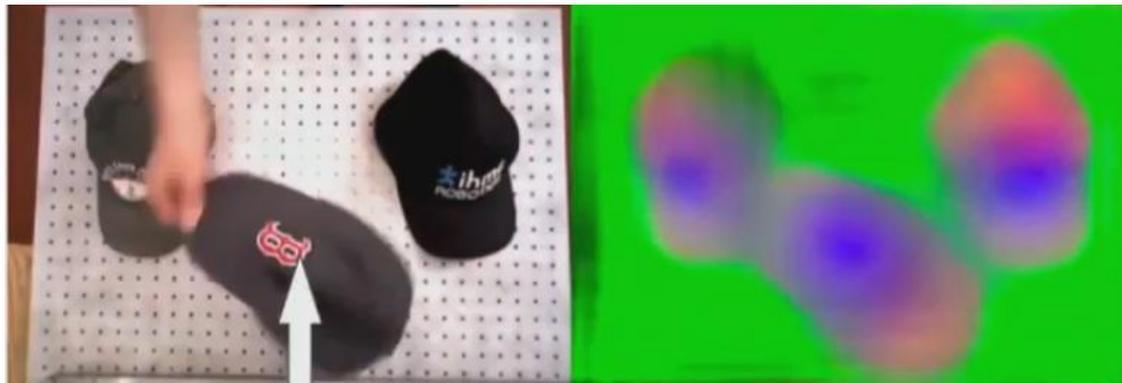
# Learned Dense Correspondences



# Multi-Object Unique Descriptors



# Class consistent descriptors



# Experiments

1) Can we acquire descriptors that can generalize across classes of objects and can be distinct for each object instance?

Yes

2) Can we apply dense descriptors to robotic manipulation?

# Results: Robotics Showcase

**Goal:** Perform grasps on target object with a real robot based on selected grasp points on a reference object

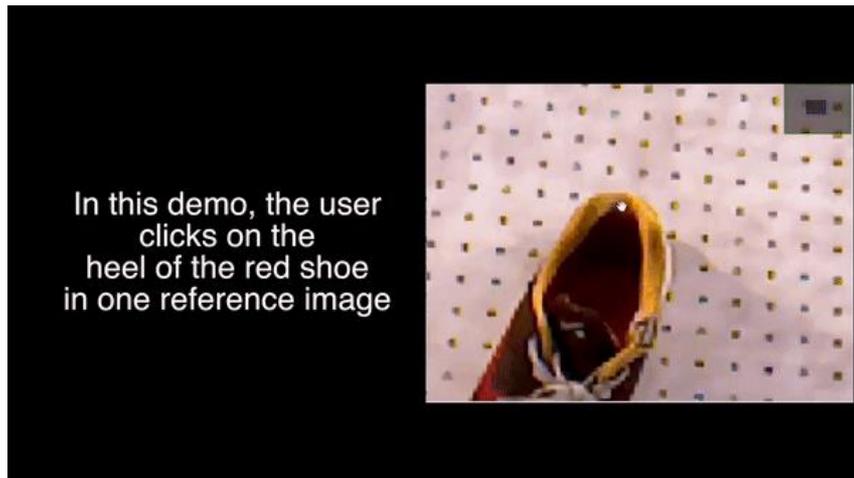
- Dense correspondence for manipulation
- Multi object dense descriptor manipulation
- Class consistent descriptor manipulation

# Dense Correspondence for Manipulation

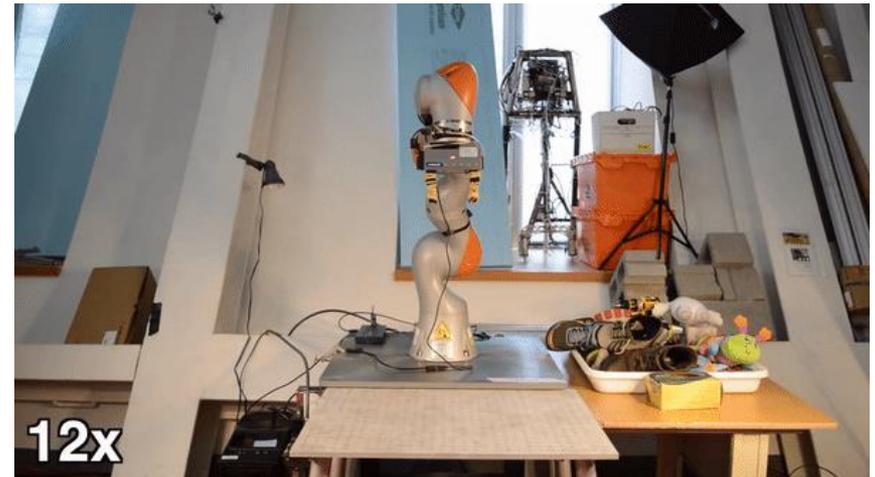
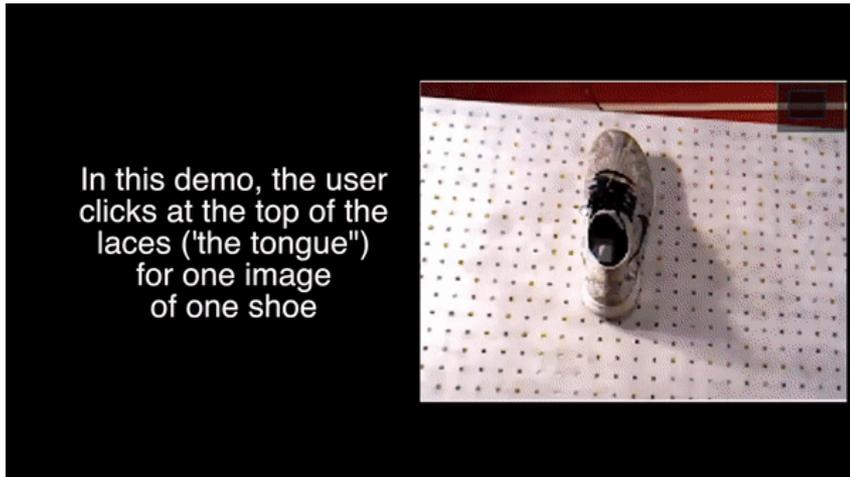
In this demonstration, the user clicks the *right* ear of the caterpillar in only one reference image



# Instance Specific Manipulation



# Class Consistent Manipulation



# Experiments

1) Can we acquire descriptors that can generalize across classes of objects and can be distinct for each object instance?

**Yes**

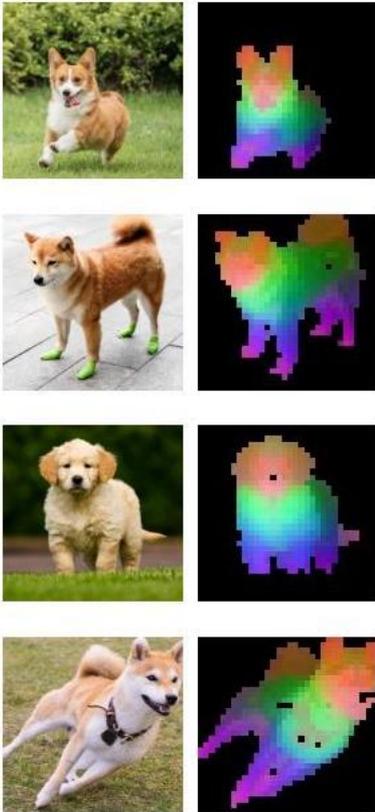
2) Can we apply dense descriptors to robotic manipulation?

**Yes, we can use descriptors to know where to grasp**

# Limitations

- Learned correspondences not always **unimodal**
- Training is finicky: **sensitive** to scale of match and non-matches
- Don't always have consistency **guarantee** (e.g. anthropomorphic toys)

# Emerging Properties in Self-Supervised Vision Transformers. Caron et al. ICCV. 2021



- Learns local features that are consistent across time and semantic classes
- Student/Teacher Architecture
- Uses Vision Transformer

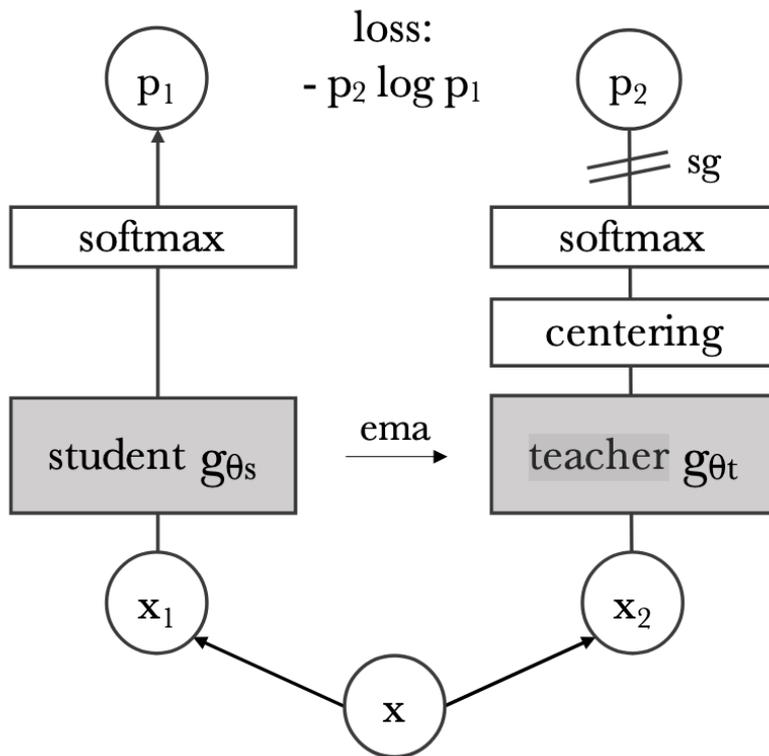
Image by Yunuk Cha.

# Background: Vision Transformer

- Image divided into Patches; patches are vectorized
- Fed into transformer with positional encoding

## Vision Transformers

# Emerging Properties in Self-Supervised Vision Transformers. Caron et al. ICCV. 2021

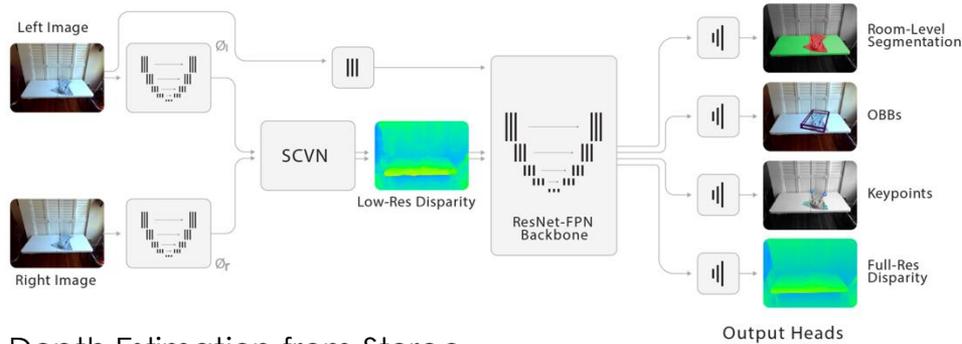


- Train student to match output of teacher
- Input: Images (here only 2)
- Output: Distribution over image categories
- Loss: Cross Entropy loss
- Teacher update with Exponential Moving Average (EMA) over past iterates of student

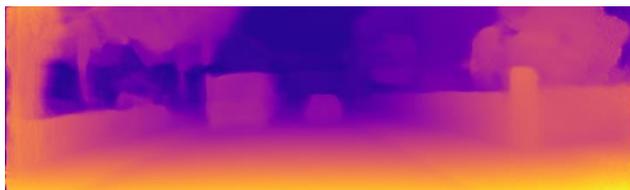
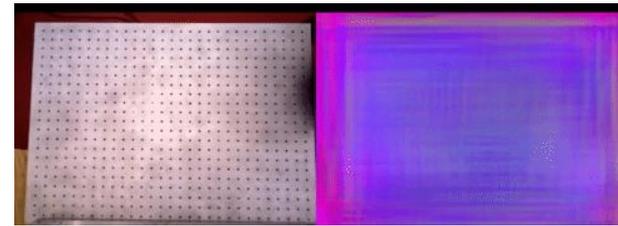
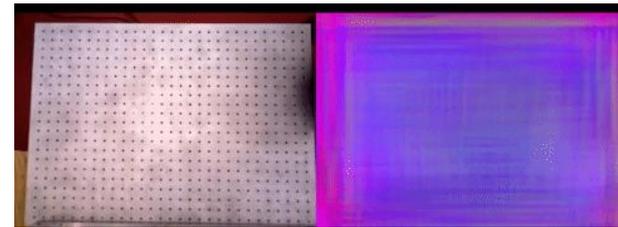
# Emerging Properties in Self-Supervised Vision Transformers. Caron et al. ICCV. 2021

- Local and global views extracted from images
  - crops have same class label as original image
- Teacher gets only global view
- Student gets both, global and local views

# Let's use representation learning!



Depth Estimation from Stereo  
Supervised Learning

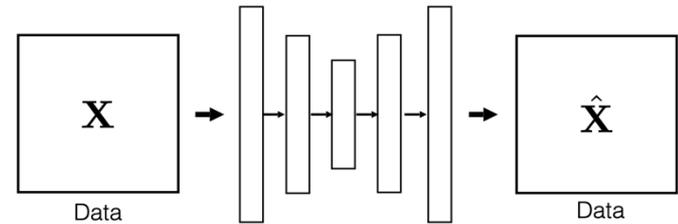


Monocular Depth Estimation  
Unsupervised Learning



Image by Yunuk Cha.  
Finding Correspondences across  
Frames  
Self-Supervised Learning

# Summary



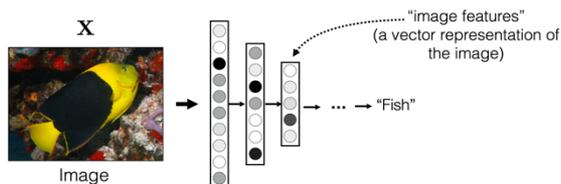
- Hourglass structure for representation learning
- Mapping either to depth or to descriptors
  - Descriptors are used for localization, robot grasping, tracking
  - Keypoint matching
- Training is supervised, unsupervised, self-supervised by exploiting structure that you know about the problem
  - Eases demand for ground truth labels
- Known structure can be used to generate training data (ground truth depth, correspondences, optical and scene flow, semantics, ...)

# Common themes in newer models

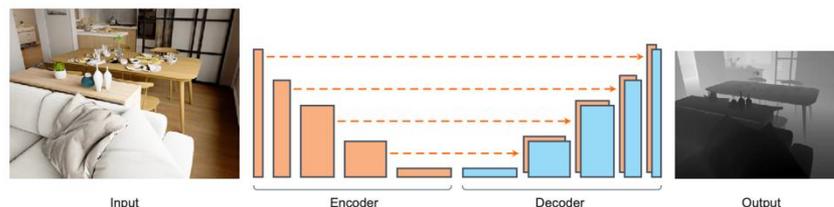
- Cost volume for depth estimation
- Implicit learning of correspondences
- Simulated training data
- Student / Teacher distillation
- Vision transformer

# Learning Goals for Upcoming Lectures

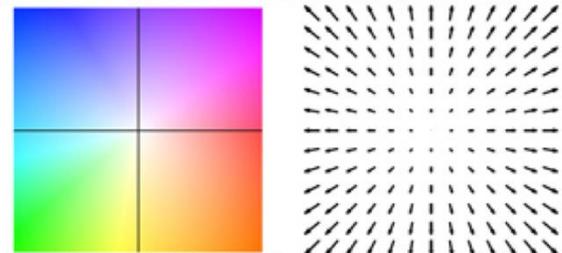
## Representations & Representation Learning



## Using Representation Learning for Depth Estimation and Finding Correspondences

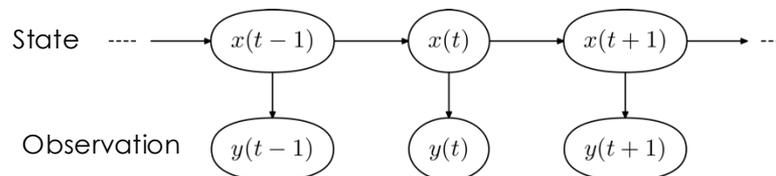


## Optical & Scene Flow

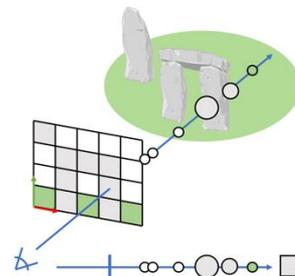


A Database and Evaluation Methodology for Optical Flow.  
Baker et al. IJCV. 2011

## Optimal Estimation



## Neural Radiance Fields



CS231

# Introduction to Computer Vision

Next Lecture:

Optical Flow and Scene Flow

