And now for something completely different!

# Learning-Based 3D

# Goal

- I'd like to answer: what is computer vision and where is it headed?

- In the process, I'd like to give you a sense of what computer vision is like.

# What is CV?

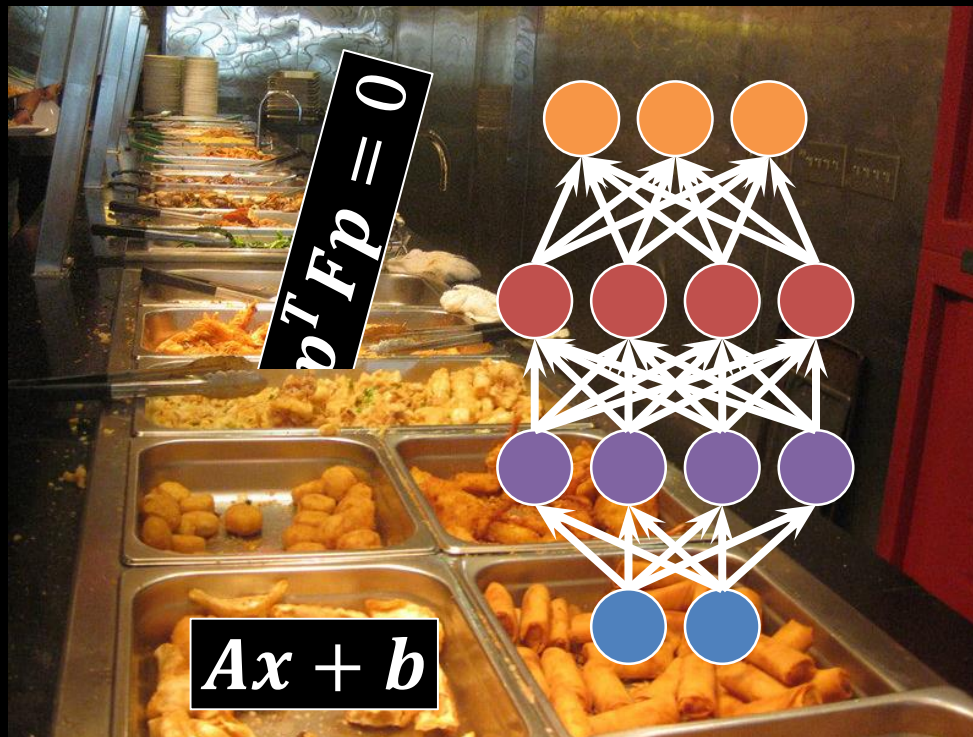## Get a computer to understand

# What is CV?

This could incorporate:

- Naming things (recognition)
- Reconstruction (geometry)
- Understanding opportunities for action (didn't cover – call me in 10 years)
- In the process, requires building up tools for processing images and fitting models

# Reality of "What is CV?"

Right now: most people would say a buffet of techniques & accumulated knowledge about geometry, pixels, data and learning.

# Don't Be Disappointed With The Buffet!

Understanding an image is incredibly difficult, involves much of your incredible brain, and is perhaps "AI-Complete"

Plus, email me in 15 years and see if we have non-buffet answers.

Get a computer to understand

# Where is it Headed?

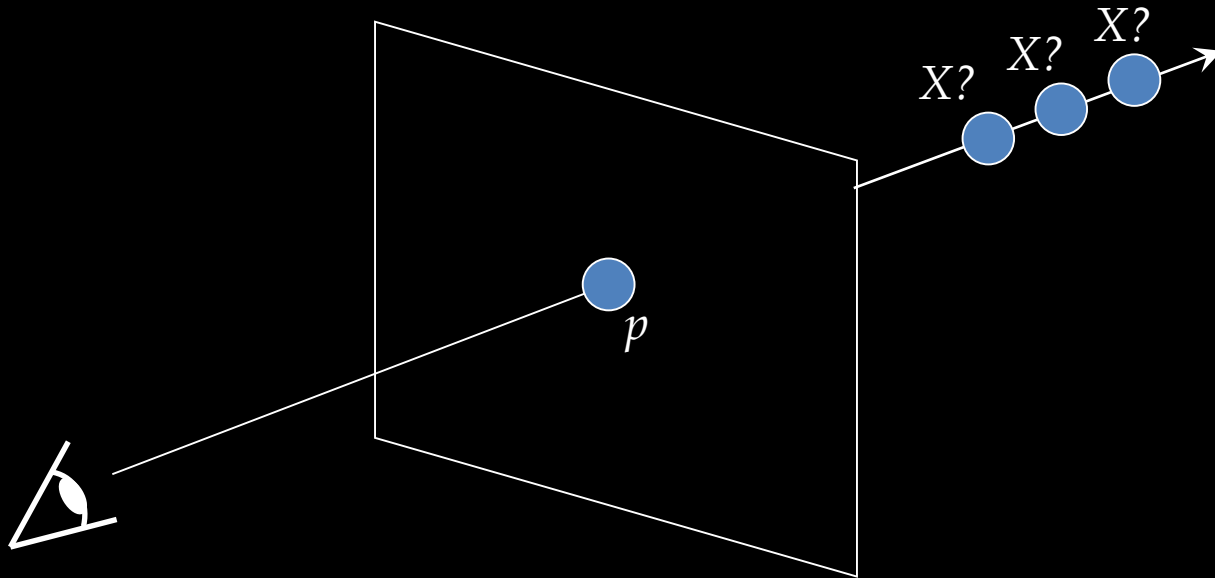3 topics that I am betting on or which other people are betting on:

- Learning and geometry (Today)

- Embodied agents (Thursday)

- Vision and language (Next Tuesday)

Some slides won't be posted since I'm borrowing heavily from others' current research slides that they've been generous to share.
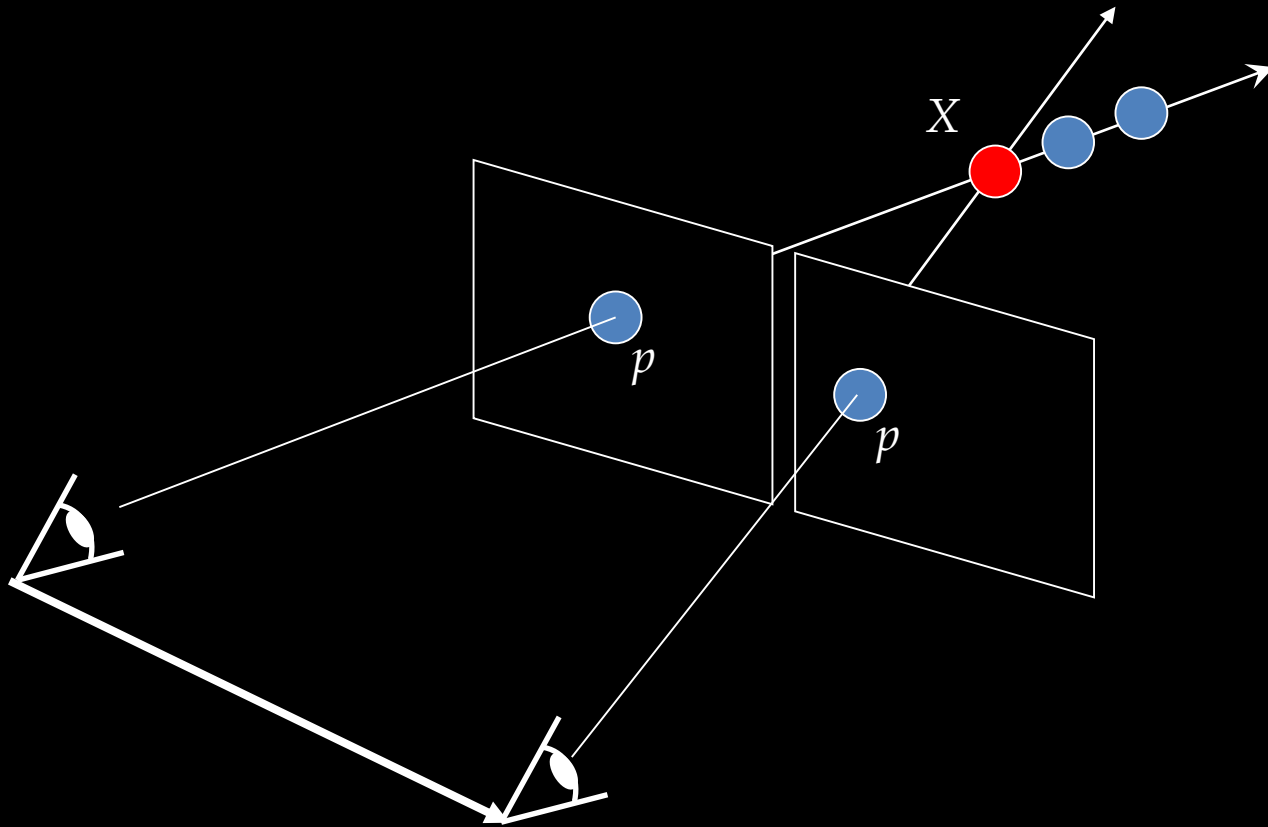
# In the Process

- I hope to give you a sense of how:
- Research in vision is conducted
- We think we know we've succeeded
- We think we know we're not fooling ourselves!

# Cues For 3D



- Given a *calibrated camera* and an image, we only know the ray corresponding to each pixel.
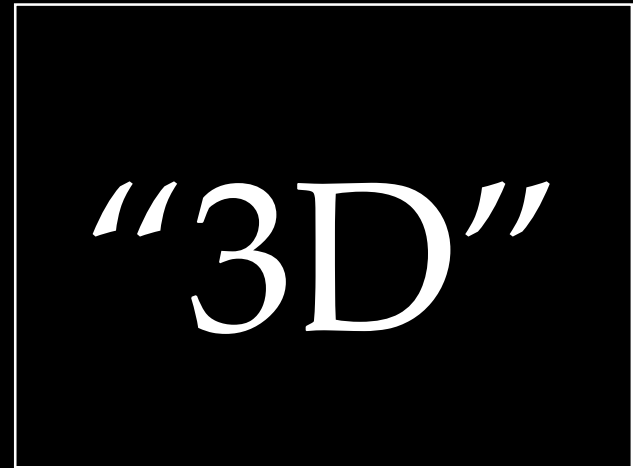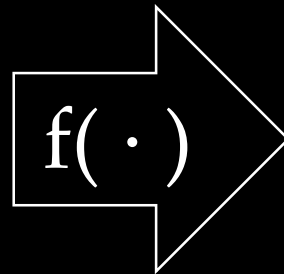- Nowhere near enough constraints for X

# Cues For 3D



- Stereo: given 2 calibrated cameras in different views and correspondences, can solve for X
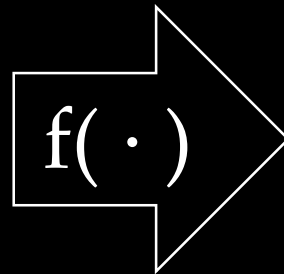
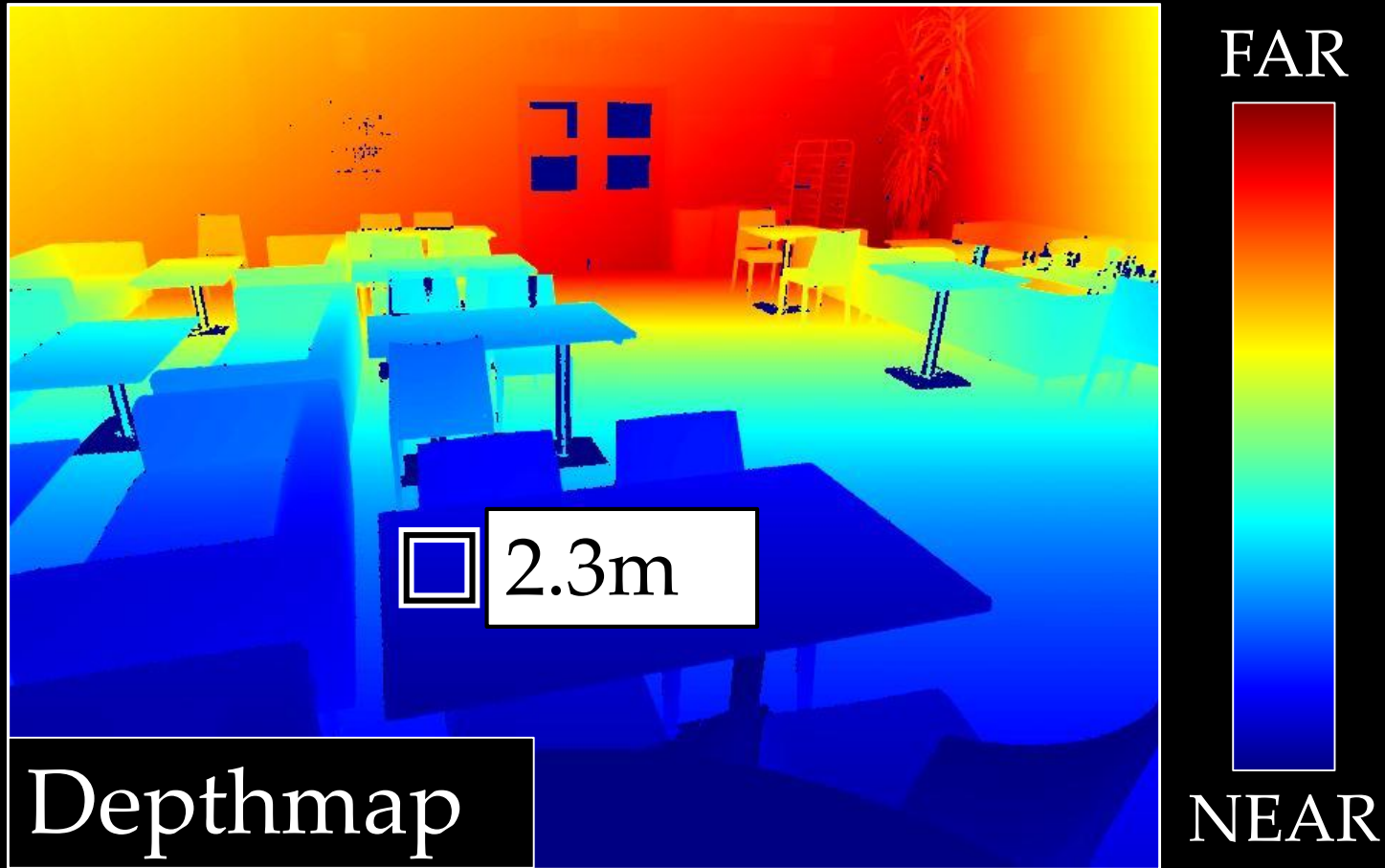# Cues For 3D

## Yet when you look at this…

# Pictorial Cues for 3D



f( · ) → "3D"

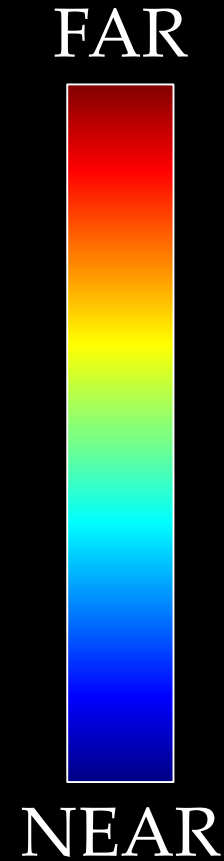Learned from Data

# Pictorial Cues for 3D



$f(\cdot) \Rightarrow$ "3D"

Learned from Data

# 3D Representations

# 3D Representations



$$z(u,v)$$

$$\partial$$

$$\left[\frac{\partial z}{\partial u}, \frac{\partial z}{\partial v}, -1\right]$$

FAR

NEAR

# 3D Representations



Surface Normals
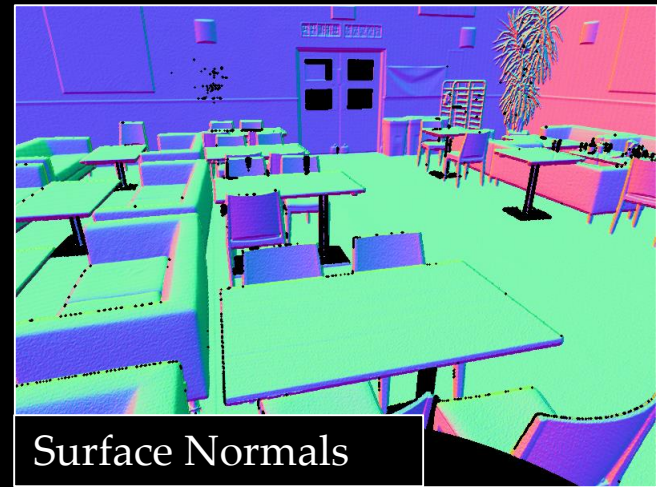
[0.06,0.99,0.12]

Room

Legend

# 3D Representations



f( · )

Surface Normals
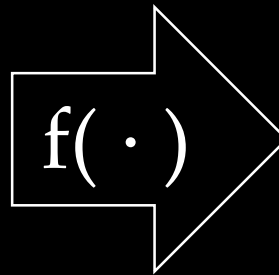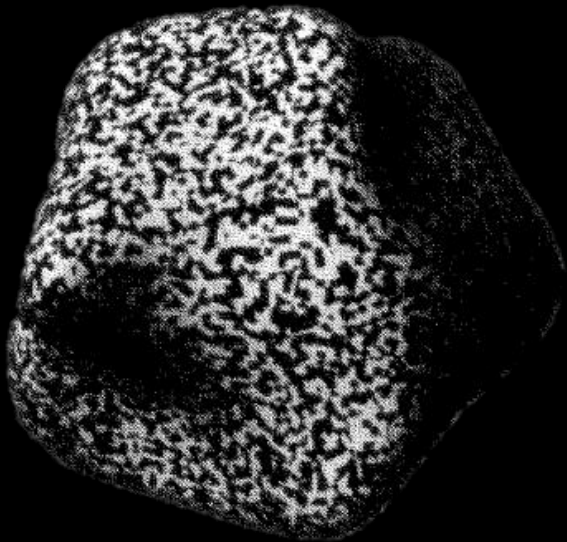
D.F. Fouhey, A. Gupta, M. Hebert. *Data-Driven 3D Primitives for Single Image Understanding.* ICCV 2013.
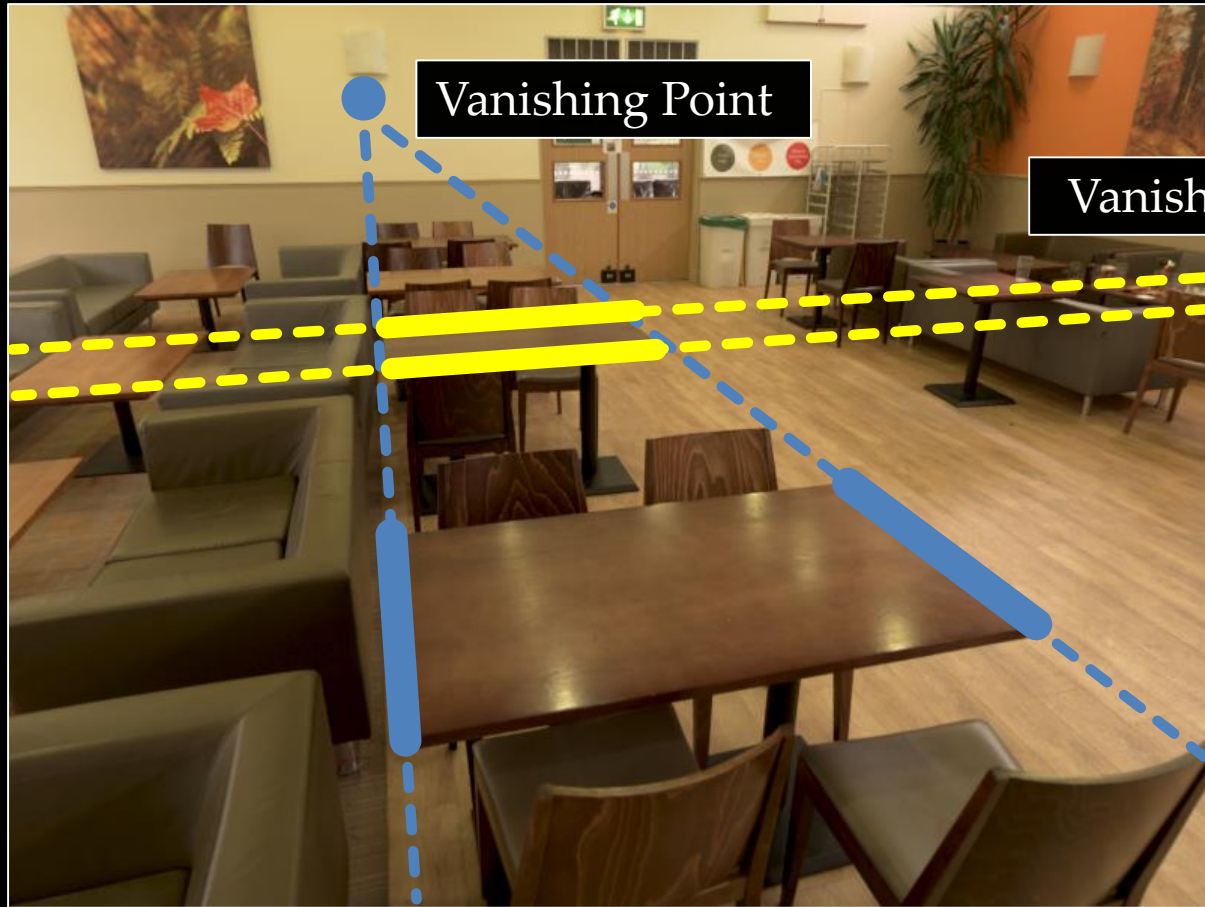
# Direct Cues For Normals



"Depth-difference judgments and attitude settings [surface normals] appear to be independent tasks."
-Koenderink, van Doorn, Kappers '96

Norman and Todd, *The Discriminability of Local Surface Structure*. Perception 1996
Koenderink, Van Doorn, Kappers. *Pictorial Surface attitude and Local Depth Comparisons*. Perception & Psychophysics, 1996
Johnston and Passmore, *Independent Encoding of Surface Orientation and Surface Curvature*. Vision Research, 1994
etc.

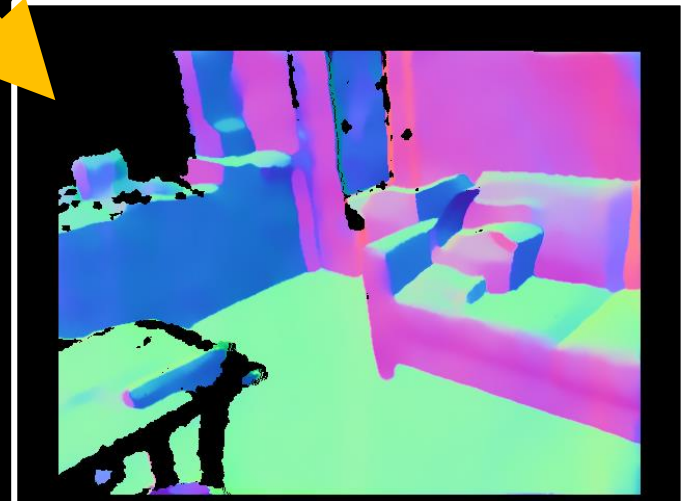# Direct Cues For Normals

# Direct Cues to Normals



Vanishing Point

Vanishing Point

# Comment – Representations

- For something as simple as whether to predict depth (z(u,v)) or the orientation of the plane ($\left[\frac{\partial z}{\partial u}, \frac{\partial z}{\partial v}, -1\right]$), there are different:

- Metrics (duh)

- Methods (hmm) – these typically aim to take advantage of special structure of the problem
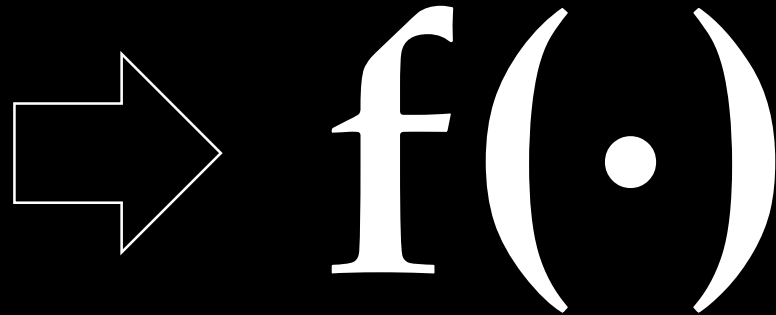
- Applications of techniques

# Surface Normals



## Color Image

## Normals

D. F. Fouhey, A. Gupta, M. Hebert. *Data-Driven 3D Primitives for Single Image Understanding.* ICCV 2013.
D. F. Fouhey, A. Gupta, M. Hebert. *Unfolding an Indoor Origami World.* ECCV 2014.
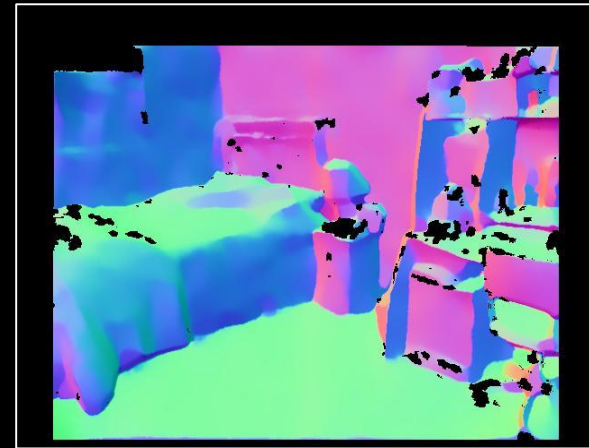X. Wang, D.F. Fouhey, A. Gupta. *Designing Deep Networks for Surface Normal Estimation.* CVPR 2015.

# Surface Normals



$$\mathbf{f}(\cdot)$$

<u>D. F. Fouhey</u>, A. Gupta, M. Hebert. *Data-Driven 3D Primitives for Single Image Understanding*. ICCV 2013.
<u>D. F. Fouhey</u>, A. Gupta, M. Hebert. *Unfolding an Indoor Origami World*. ECCV 2014.
X. Wang, <u>D.F. Fouhey</u>, A. Gupta. *Designing Deep Networks for Surface Normal Estimation*. CVPR 2015.
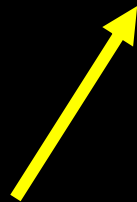
# Applying Deep Learning
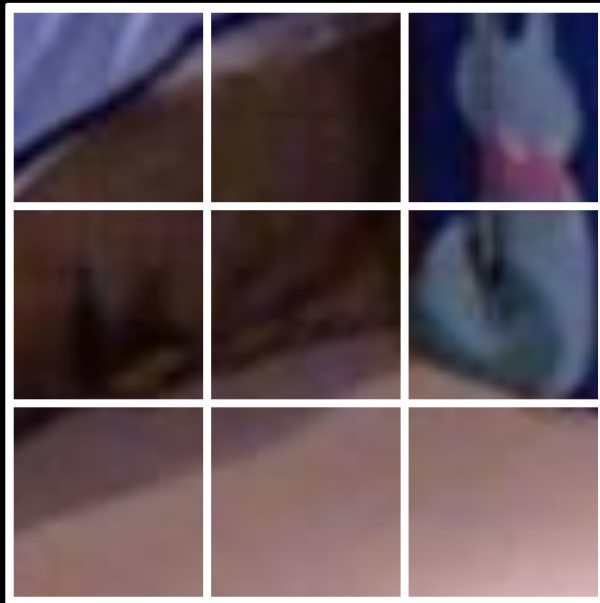
## Input

## Output



**CNN**

How do we incorporate constraints?

How do we represent the output?

X. Wang, <u>D.F. Fouhey</u>, A. Gupta. *Designing Deep Networks for Surface Normal Estimation.* CVPR 2015

# Representation and Objective

## Input

## Ground Truth



## Quantized Normals

Class 1: ↙     Class 2: ↑     …     Class K: ↘

Normal quantization scheme from Ladicky et al. 2014

# Results

# Results

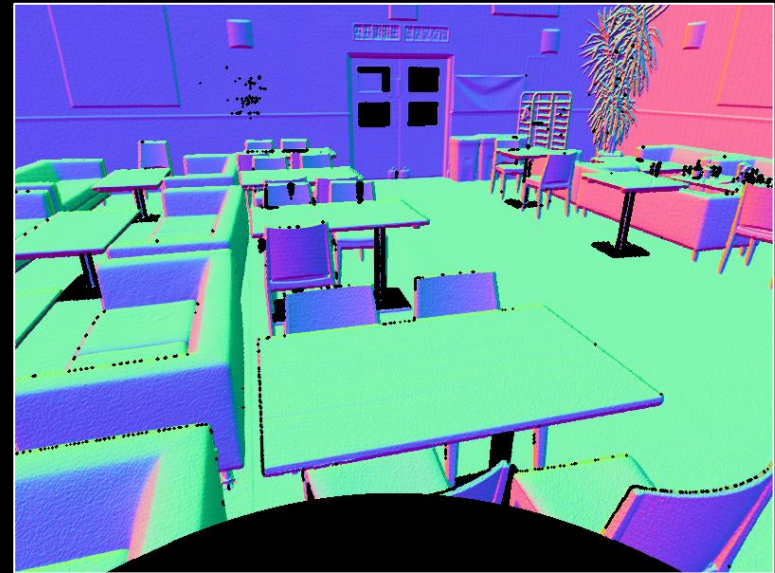| Input | Output | Input | Output |
|-------|--------|-------|--------|

# Results

Input　　　Output　　　Input　　　Output

# Comment – Picking Problems

- I'd show results, and the response from many people would be "sure, sure, neat but I'll just buy a Kinect"
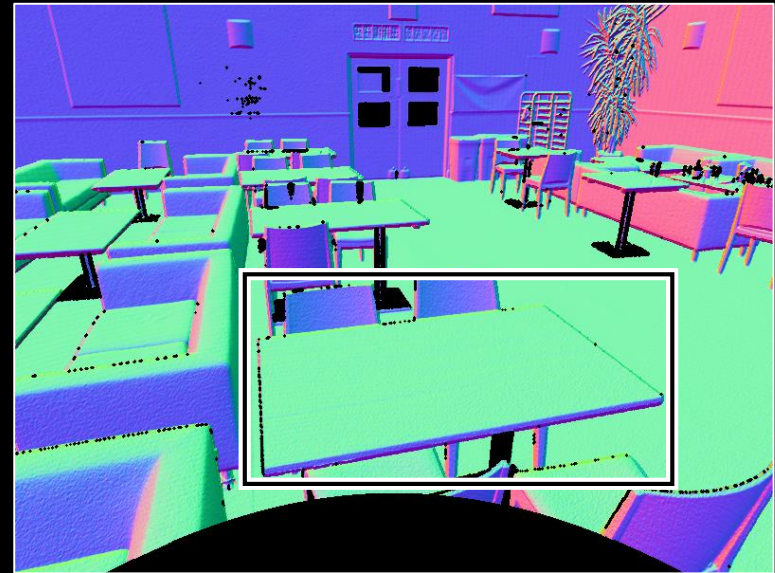
# 3D Representations

# 3D Representations



~$50K, 6.5 minutes an image

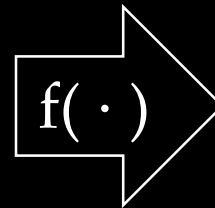Adams et al., Scientific Reports, 2016

# 3D Representations



How thick is the table?
What's behind it?
Is the chair attached to the table?
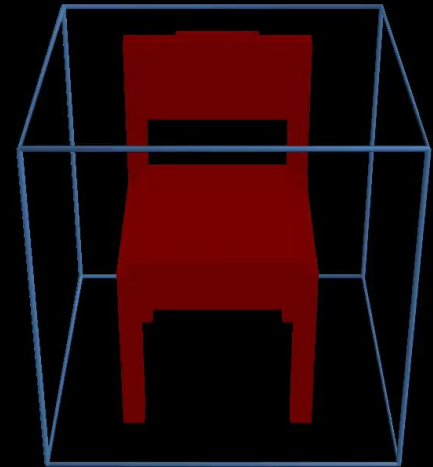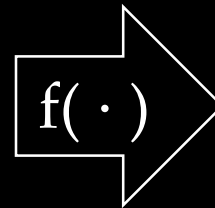
# 3D Representations



RGB Image

Voxels

R. Girdhar, <u>D. F. Fouhey</u>, M. Rodriguez, A. Gupta.
*Learning a predictable and generative vector representation for objects*. ECCV 2016
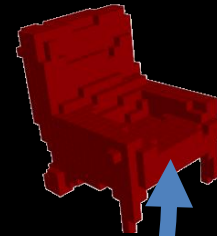Contemporary work also proposing to predict voxels: C. Choy et al. ECCV 2016.

# 3D Representations



RGB Image

Voxels

$f(\cdot)$

R. Girdhar, D. F. Fouhey, M. Rodriguez, A. Gupta.
*Learning a predictable and generative vector representation for objects*. ECCV 2016
Contemporary work also proposing to predict voxels: C. Choy et al. ECCV 2016.

# Approach



20x20x20
Voxel output
$S$

$S_{i,0}$  P(voxel i empty)   0.1

$S_{i,1}$  P(voxel i filled)   0.9

$S_i$

# Approach

224x224x3
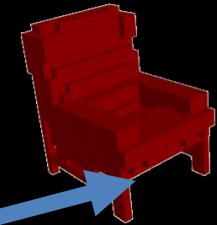Image Input
$I$

$$S = f(I; \theta)$$

20x20x20
Voxel output
$S$

20x20x20
Voxel Truth
$Y$
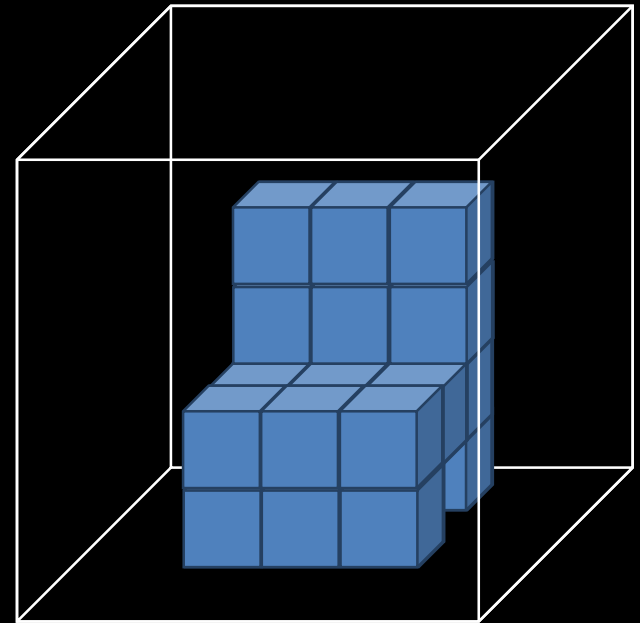
$$-\sum_{i=1}^{20^3} \log(S_{(i, Y_i)})$$

# Main Idea

Representation should satisfy

- **Generative in 3D:** should be able to construct objects

- **Predictable from 2D:** should be able to infer from an ordinary image

# Output Representation – Voxels

- Binary 3D Pixels
- Size fixed in advance
- Spatially organized

# Motivation

## How many couches are there really?

1737662031938094565999824459494356270619397861001172505471732865032623760224580084650943336301208543380031943621630075979872254724835986408433356854417101939662741313385571925863990067892927145547675001947961279645969066059766058736658595806001619985565113685309604009071992534506041686227703502285271246267285386268054188334701076510916419199007254159946899201122191709070235613544840470257137346516087775445798461110010594821321809566894441083157854016421880441787886298535922284673317305198107635595779448820162864939086315031011211661095716822957694703795145311052399652092453140826655185793355112915252303733164866977865323352062741492408134892018287738543530418555987093906754309603810722704323839135427021302024301866373218623310688617767802110828569845060500248953943201394358684846438433680024960899560464199640198775868455302077489943945015055881469790826298713660881217637905553645132439842440041476360402191364434103777980116087227171313236217001593357864456019476016940251078882930170581785626471754610263843434388748614065167671583732790323210962621265516202556666051857894632079443919057568868296675205530147243722453008787860917005634440791070990090033802303564619892603772739860232814440760827834068244717034998446429155877901463847580516635477753360218291710334110437969770421905196578617628042261474807555550852780628662686778424328514217905444070065811486319791485712994179639505792107199614224057680713352133248427093162050320783841687500910179645840602852401071615610199305056875023319605196226197093200883827976083431810104431171076945704867210395865501638889477089206526745122893895137023742284136605273617416043159302347321706676417294976882184360647907386625286437706439808510122321655834428195676716387657988975912495603567231757812214107093305855531027459888408998287964797402026449592170306443953289820794313437457625484027204707563385674951404429813592761132843332364065753355051237690077327370327532992465146575914511457917435677059343998713575588940361336452902960404986823380729513438228473074593730991070365767610344712409763107415328712004024783714365662404505561407611183224523961270833927279826288743741681844006492504983844337080564560942431478010803001668346156259756937153997400340269790302383010805303464513307820804391749208724895834408102637878891552851996724898933859202712442391408339177188452446496864505205821815101050847125828590768535580722988074767763478937 6
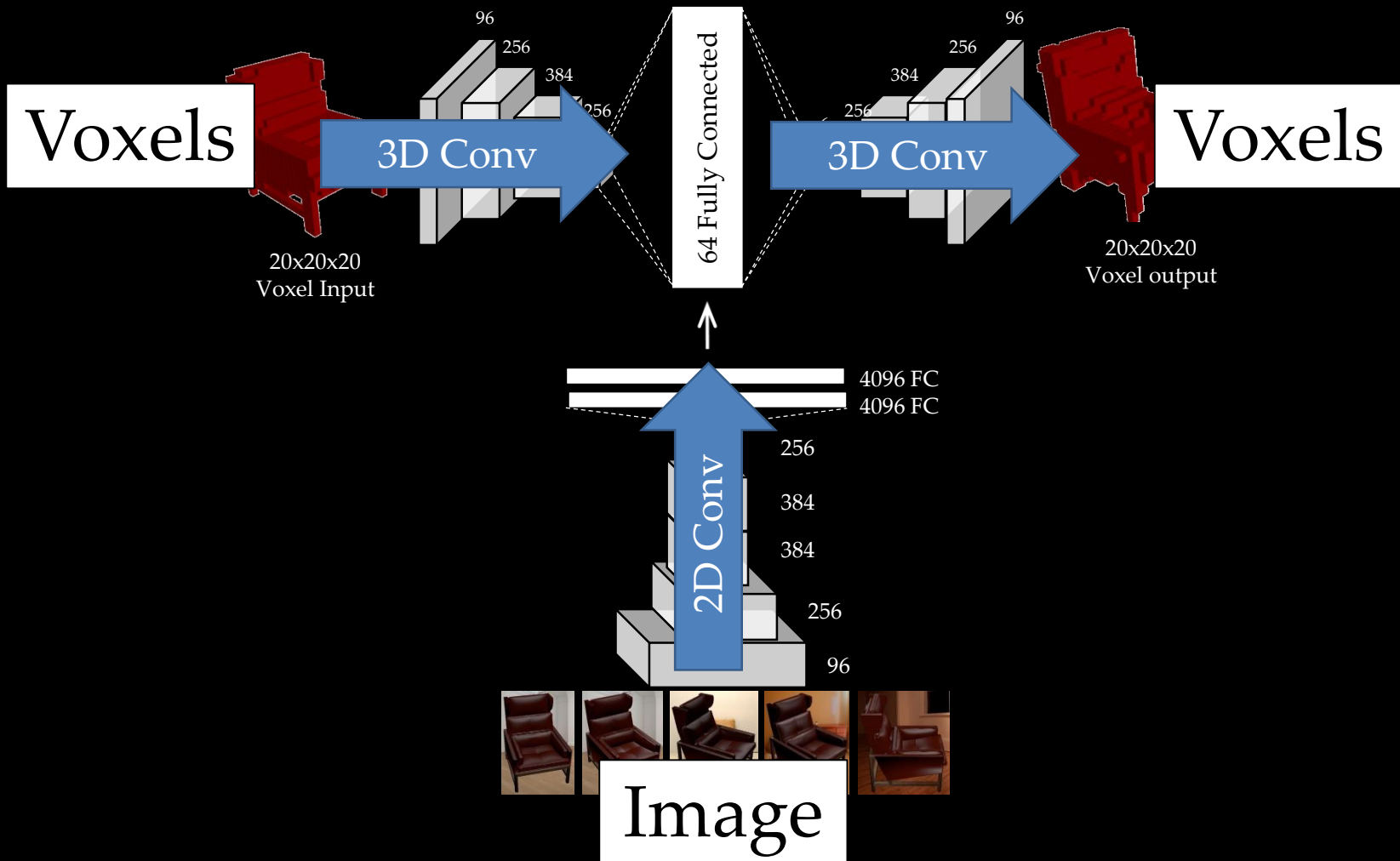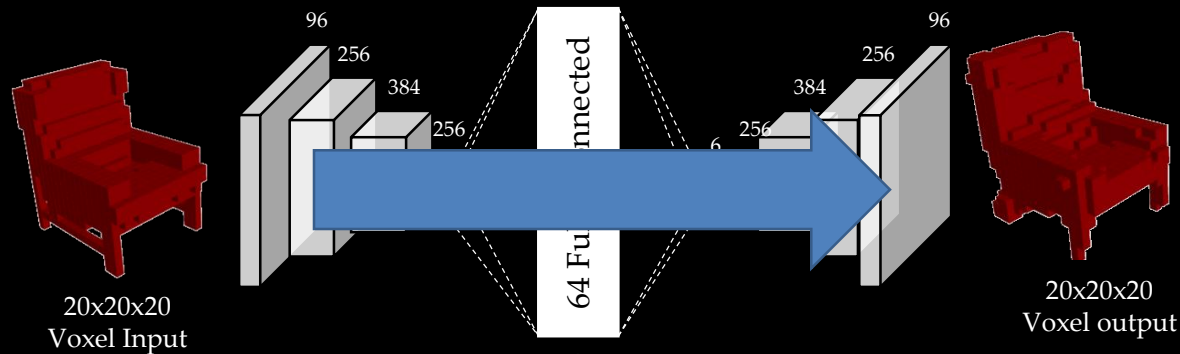
# Motivation

## How many couches are there really?

# Approach



Voxels

3D Conv

64 Fully Connected

3D Conv

Voxels

96
256
384
256

96
256
384
256

20x20x20
Voxel Input

20x20x20
Voxel output

4096 FC
4096 FC

256

384

384

256

96

2D Conv

Image

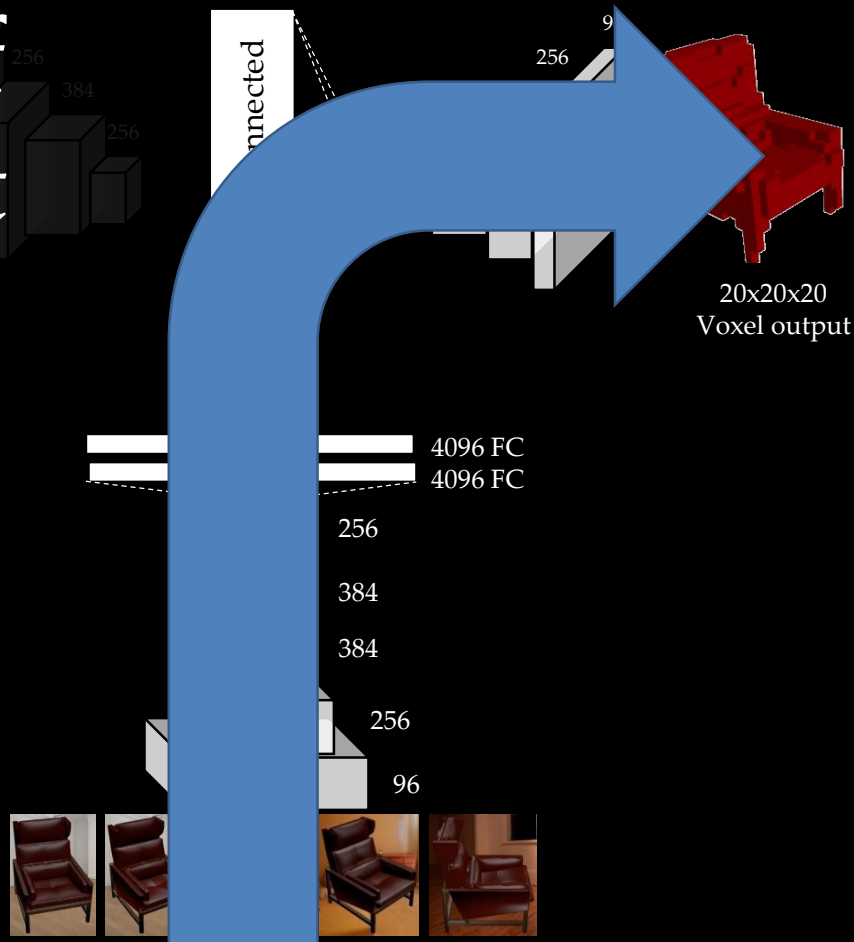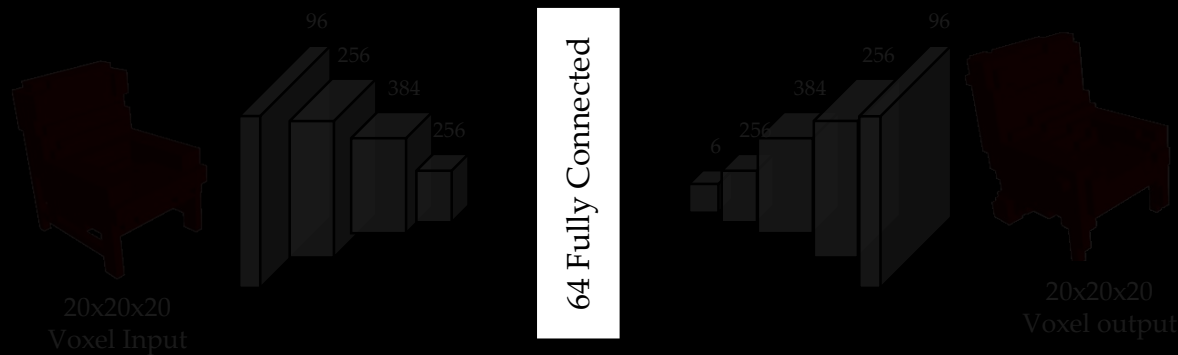# Approach



Turning off image branch
yields an autoencoder over
voxels

# Approach
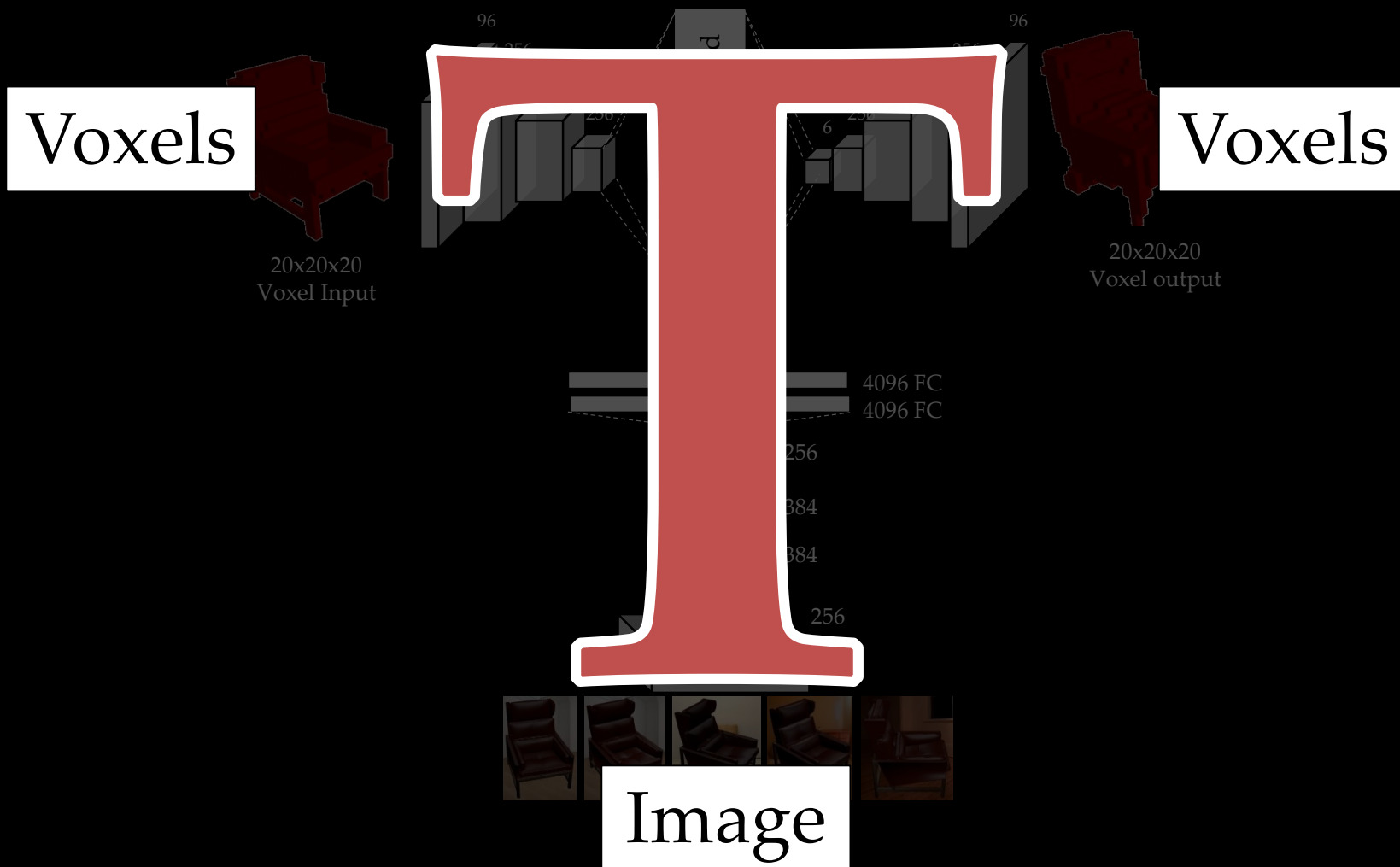
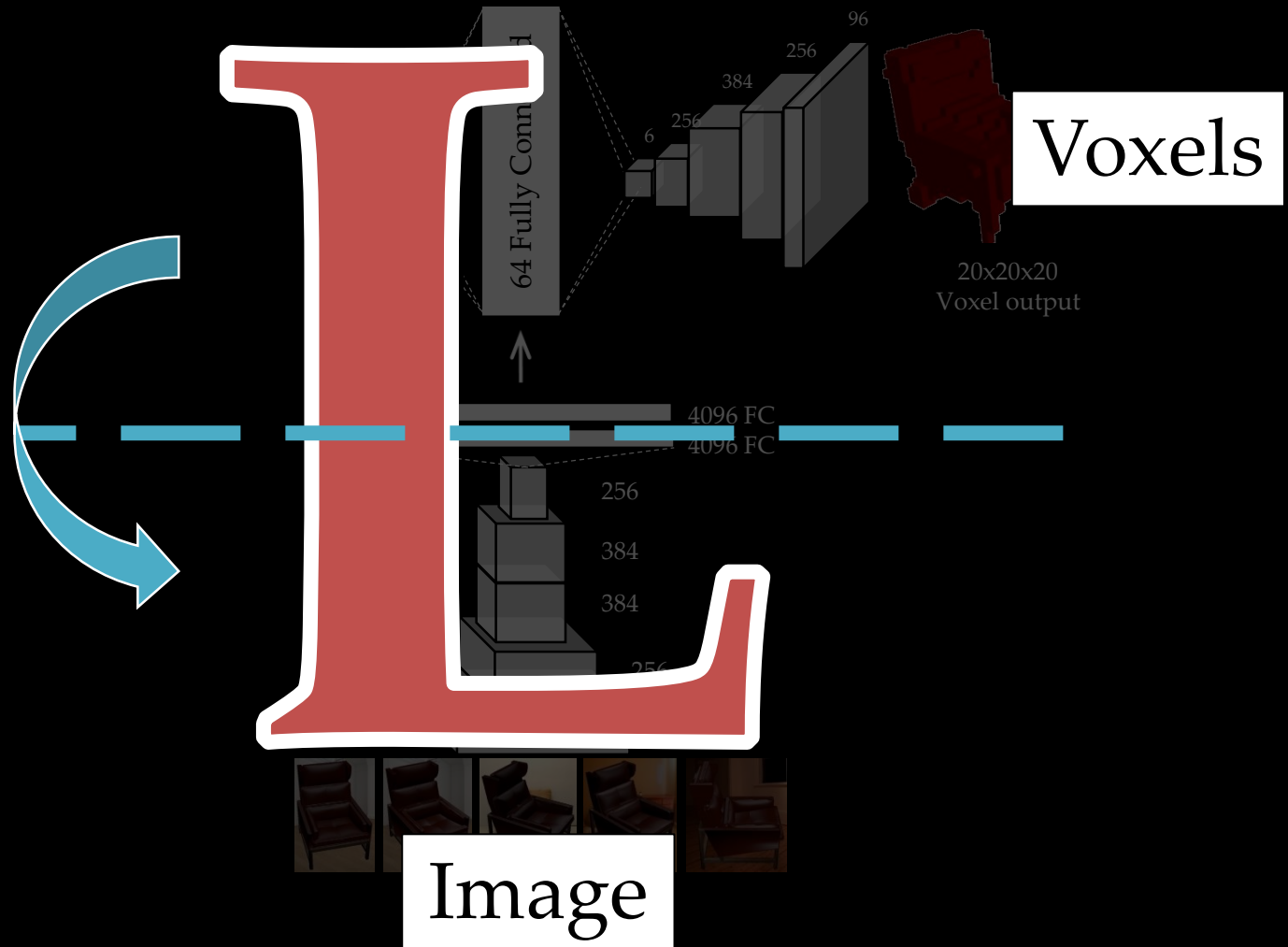Turning off voxel input yields an image to voxel predictor



256

384

256

connected

256

9

20x20x20
Voxel output

4096 FC
4096 FC

256

384

384

256

96

# Approach



20x20x20 Voxel Input

64 Fully Connected

20x20x20 Voxel output

96
256
384
256

6
256
384
256
96

4096 FC

Learned embedding parameterizes shape in a way that is:
(a) generative and predictable
(b) accessible from voxels and images

# TL-Network



Voxels

Voxels

20x20x20
Voxel Input

20x20x20
Voxel output

96

96

256

256

6

256

4096 FC
4096 FC

256

384

384

256

Image

# TL-Network

Voxels

20x20x20
Voxel output

64 Fully Connected

96
256
384
6
256

4096 FC
4096 FC

256
384
384
256

Image

# TL-Network



Voxels

20x20x20
Voxel output

64 Fully Connected

4096 FC
4096 FC

256

384

384

256

96

256

96

Image

# Training



20x20x20
Voxel Input

96

256

384

256

64 Fully Connected

6

256

384

256

96

20x20x20
Voxel output

4096 FC
4096 FC

256

384

384

256

96

# Training – Stage 1



20x20x20
Voxel Input

96
256
384
256
64 Fully Connected
6
256
384
256
96

20x20x20
Voxel output

Learned by cross-entropy loss

384
384
256
96

# Training – Stage 2



96
256
384
256

64 Fully Connected

20x20x20
Voxel Input

Frozen

96
256
384

20x20x20
Voxel output

4096 FC
4096 FC

256
384
384
256
96

Learned by
L2 Loss

Paired Image
+Voxels

# Training Data

- 5 Categories from ShapeNet (**no category labels used in learning**)

- Standard rendering techniques



Rendering techniques from Su et al. ICCV15

# Training – Stage 3



20x20x20
Voxel Input

96
256
384
256

64 Fully Connected

6
256
384
256
96

20x20x20
Voxel output

4096 FC
4096 FC

256
384
384
256
96

Learned by cross-entropy loss

# Experiments



Voxels

20x20x20
Voxel Input

96
256
384
256

64 Fully Connected

6
256
384
256
96

Voxels

20x20x20
Voxel output

4096 FC
4096 FC

256
384
384
256
96

Image

# Commentary – Experiments

- Main goals of any experiments: Empirically verify that we achieved what we said we would achieve.

- "In the computer field, the moment of truth is a running program; all else is prophecy" – Herbert Simon

# Experiments



Voxels

Voxels

20x20x20
Voxel Input

96
256
384
256

64 Fully Connected

6
256
384
256
96

20x20x20
Voxel output

Does it represent voxels well?

4096 FC
4096 FC
256
64
384
256
96

Image

# Visualization

# Quantifying Performance

Ground-Truth
Voxels

Predicted
P(Occupied)

# Voxel Representation



| Test Shape | PCA | TL-Network |

# Commentary – Baselines

- **Is 95% good or bad?**

- It depends! You might want to know: how well does something simple do? How well does a known method do?

- Typically comparison points: past methods, linear models, nearest neighbors.

- Considered embarrassing if someone later finds something simple that beats your complex method!

# Reconstruction Accuracy

Average precision

| PCA | TL |
| --- | --- |
| 96.8 | 97.6 |

Qualitatively a pretty big gap, but quantitatively not so. Because metric isn't quite right.

# Voxel Representation



**Voxels**

20x20x20
Voxel Input

96
256
384
256

f(voxels)

64 Fully Connected

f(voxels) ➡ $x \in R^{64}$

f(voxels) ➡ $x \in R^{64}$

Does this feature contain useful information for distinguishing these categories?

# Commentary – Alternate Tasks

- Convnets can have hundreds of million of degrees of freedom
- Biggest fear of many researchers: are we actually learning the thing we set out to learn or something entirely different?
- This can have profound issues, especially if you deploy this
- One solution: test the features on an entirely different task

# Voxel Representation

- Classification of 3D shape categories (e.g., toilet) on ModelNet40
- TL was <u>not</u> trained for this task; support vector fit on features

| No Class Info. | | Class Info. Used |
|---|---|---|
| PCA | TL | 3D ShapeNets |
| 68.4 | 74.4 | **<u>77.3</u>** |

Wu et al., 3D ShapeNets: A Deep Representation for Volumetric Shape Modeling, CVPR '15

# Experiments

Can you predict voxels from images?

64 Fully Connected

Voxels

20x20x20
Voxel output

4096 FC
4096 FC

256
384
384
256
96

6    256    384    256    96

Image

# Reconstructing Test Models

# Comments

- Train on synthetic images, test on synthetic images. **Any issues?**

- Is the network actually learning something or cheating by using cues left in by the renderer

- One solution: test on new, non-rendered data

# Reconstructing IKEA



Data from Lim et al., Parsing IKEA Objects: Fine Pose Estimation, ICCV 2013

# Baseline

## Current Setup

64 Fully Connected

6
256
384
256
96

**Voxels**

20x20x20
Voxel output

4096 FC
4096 FC

256
384
384
256
96

**Image**

# Baseline

Go directly
for voxels



Voxels

20x20x20
Voxel output

4096 FC
4096 FC

256

384

384

256

96

Image

# Comments – Baselines

- Not quite the right baseline: just tests whether the 3D convolutional structure is necessary.

- But you don't always get things right the first time around!

- Research is a process, and the real knowledge comes from multiple papers in a whole series, typically from different authors, not from just one paper

# Quantitatively

| | Direct to Voxels | | |
| --- | --- | --- | --- |
| | Conv4 | FC8 | TL Networks |
| CAD | 38.0 | 24.8 | **65.4** |
| IKEA | 31.1 | 19.8 | **38.3** |

CAD

IKEA

# Applying to Scenes

## Input: RGB Image

## Output: Voxels



## 1. Doesn't separate objects

S. Tulsiani, S. Gupta, D.F. **Fouhey**, A.A. Efros, J. Malik.
*Factoring Shape, Pose, and Layout from the 2D Image of a 3D Scene.* To appear at CVPR 2018.

# Applying to Scenes

## Input: RGB Image

## Output: Voxels



**2. Conflates shape and pose**

# Applying to Scenes

## Input: RGB Image

## Output: Factored

# Applying to Scenes

## Output: Factored

## Part 1: Layout

# Applying to Scenes



## Output: Factored



## Part 2: Per-Object



Voxels ($32^3$)

Scale    Rotation

Translation

# Approach

Image



Layout

Standard encoder/decoder
with skip connections

# Approach

Per-Object

Image



Bounding
Boxes

Voxels
$(32^3)$

Scale

Trans.

Rot.

# Approach

# Quantitative Results

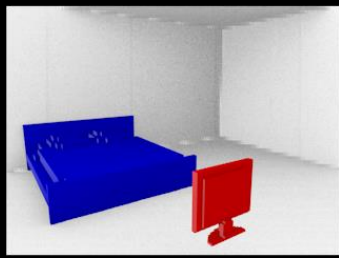|  | % Rot. < 30° | % Trans. < 1m |
|---|---|---|
| Base | **75.2** | **90.7** |
| No Context | 69.3 | 85.4 |
| Regress Rotation | 48.1 | -- |

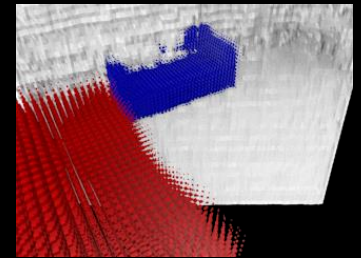# Results (SUNCG)

**Input**         **Ground-Truth**         **Prediction**

# Results (SUNCG)

Input      Ground-Truth      Prediction



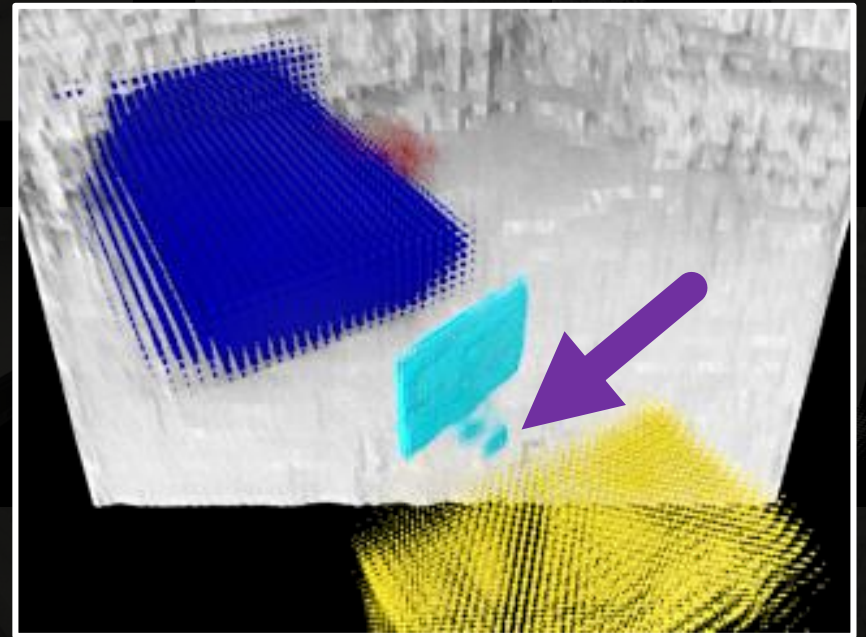SunCG: Song et al. CVPR 2017, rendered by Zhang et al. CVPR 2017

# Results (SUNCG)

Input            Ground-Truth                    Prediction
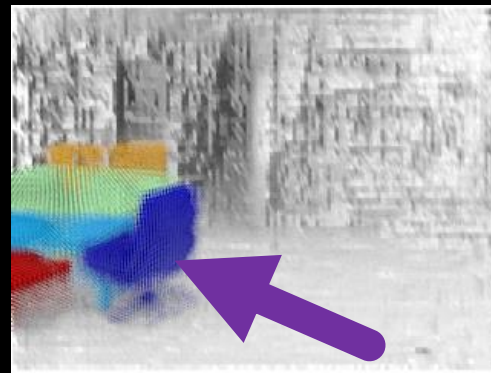


SunCG: Song et al. CVPR 2017, rendered by Zhang et al. CVPR 2017

# Results (NYUv2)



Input  Prediction  Other View

NYUv2: Silberman et al. ECCV 2012
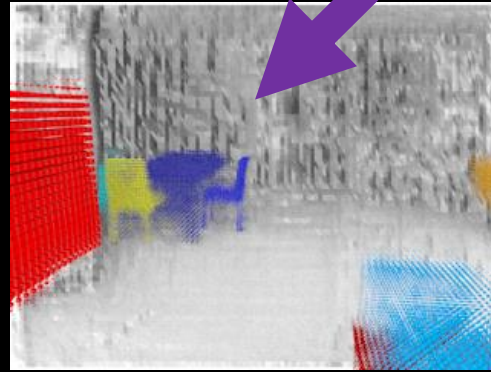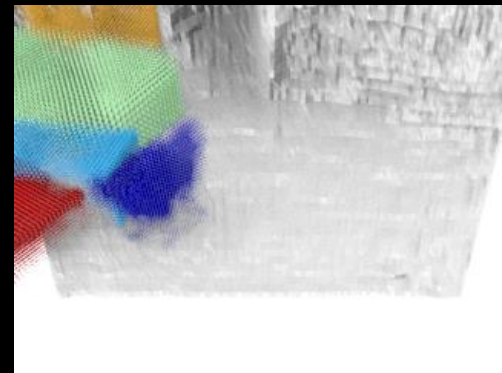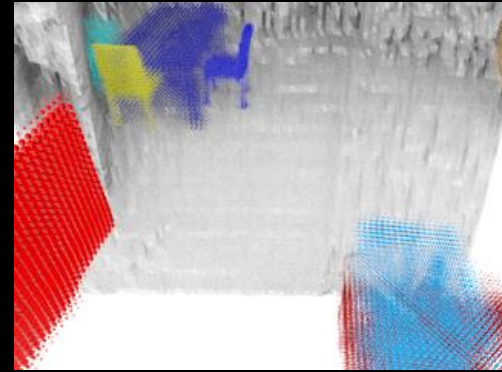
# Representational Benefits

Input RGB Image

Output