# Intelligent Perception

S. M. Ali Eslami

December 2016

**Underlying scene**

**Observation**

1. How should the scene be represented?

2. How should the representation be computed?
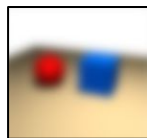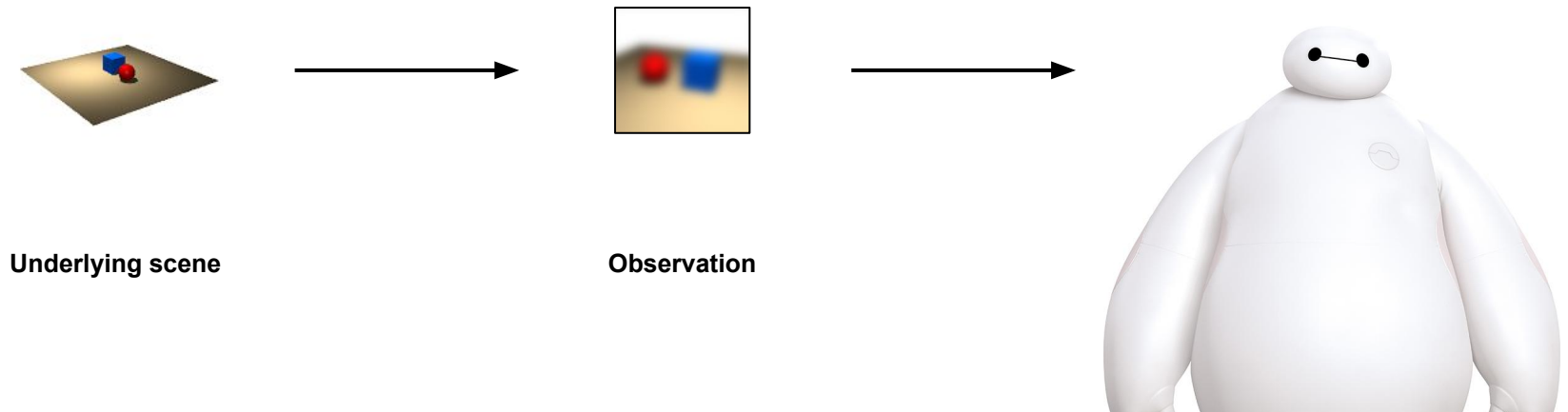
**Underlying scene**
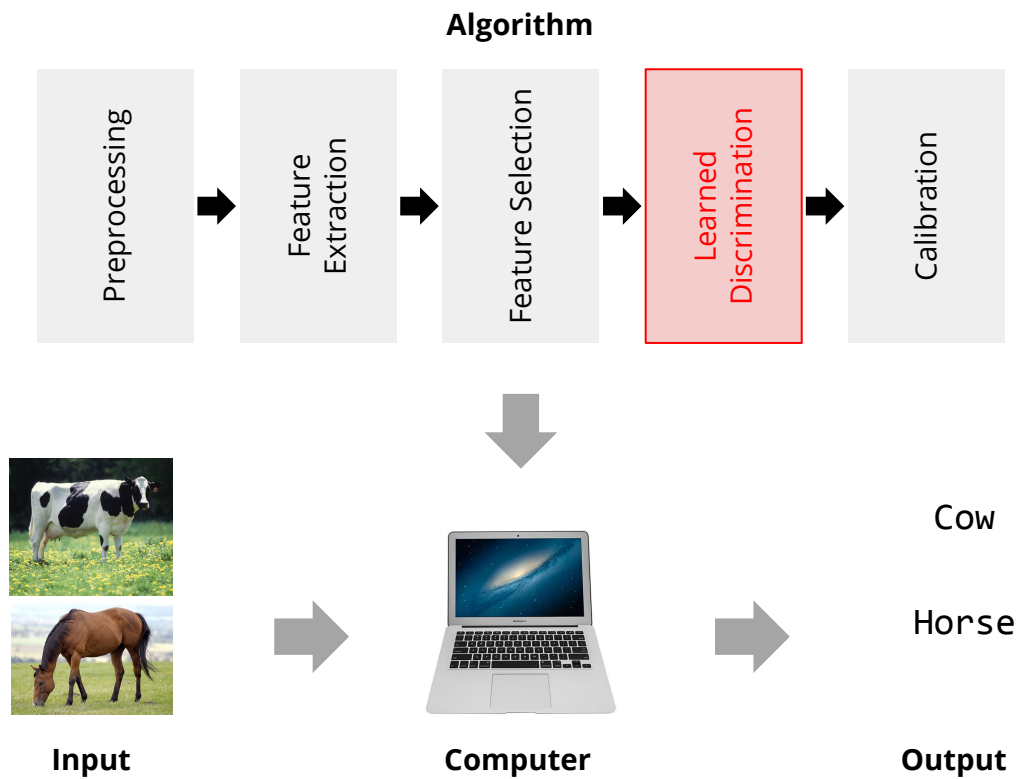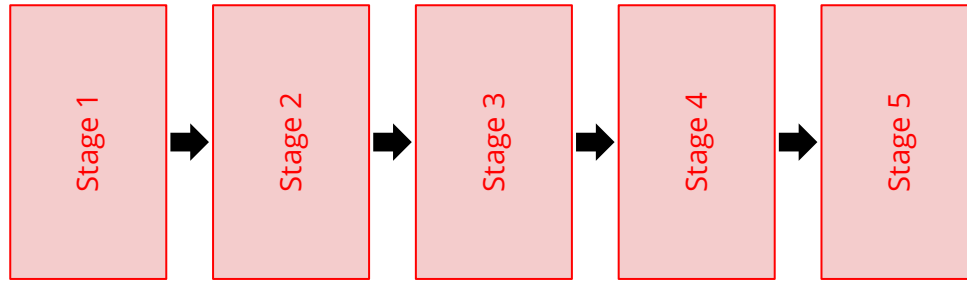
**Observation**

?

# Learning paradigms

horse



**Supervised
Learning**

**Algorithm**

Preprocessing → Feature Extraction → Feature Selection → Learned Discrimination → Calibration

**Input** → **Computer** → **Output**
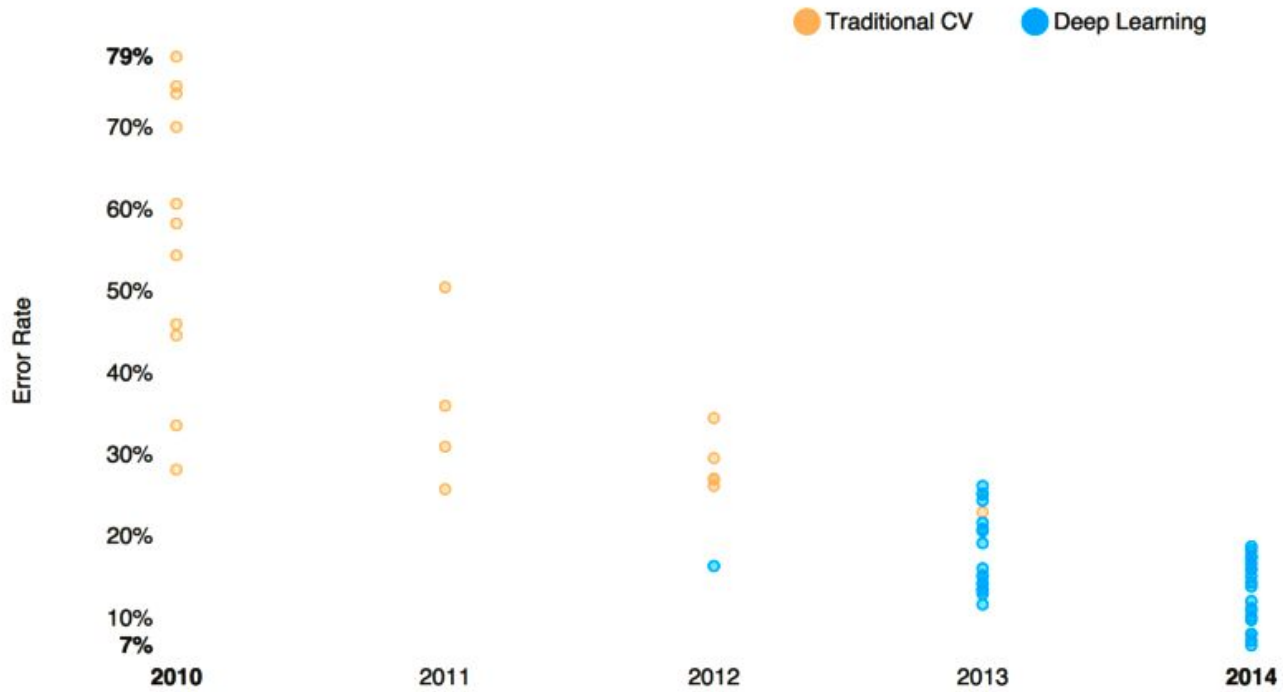
Cow

Horse

# Deep Supervised Learning

- Optimize directly for the end loss

- End-to-end training, no engineered inputs

- With enough data, learn a big non-linear function

- Supervised labeling is often enough for transferrable representations

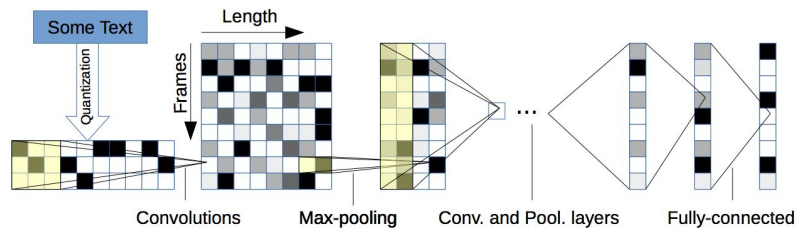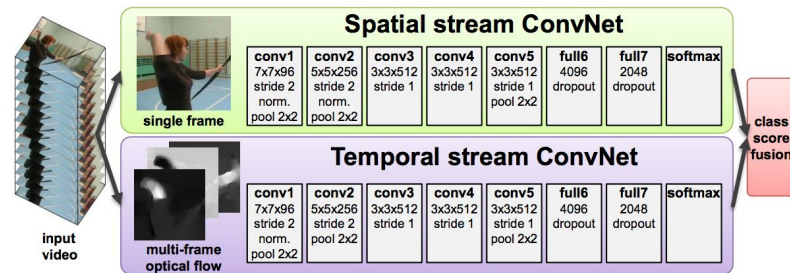- Large labeled dataset + big / deep neural network + GPUs

Clarifai (2014)

# Deep Supervised Learning

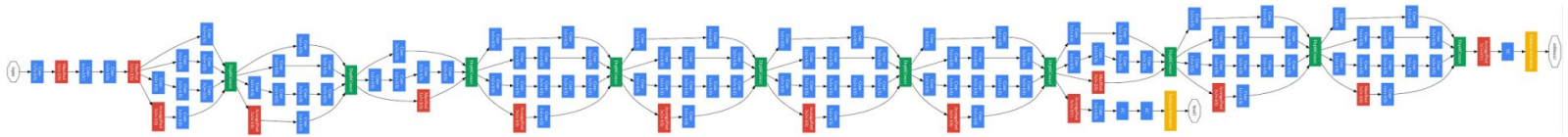## Text Classification



Zhang et al. (2015)

## Video Classification



Simonyan et al. (2014)

# Deep Supervised Learning

- Innovation continues
  - Inception (Szegedy et al., 2015)
  - Residual connections (He et al., 2015)
  - Batchnorm (Ioffe et al., 2015)
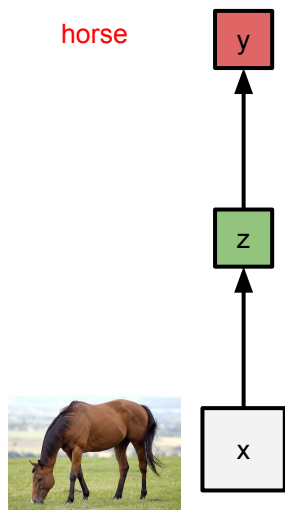
- Performance is continuously improving
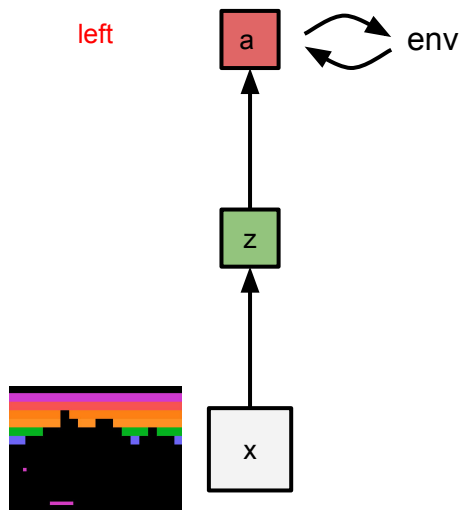


Szegedy et al., (2015)

**Where does the data come from?**

**What is the correct representation?**

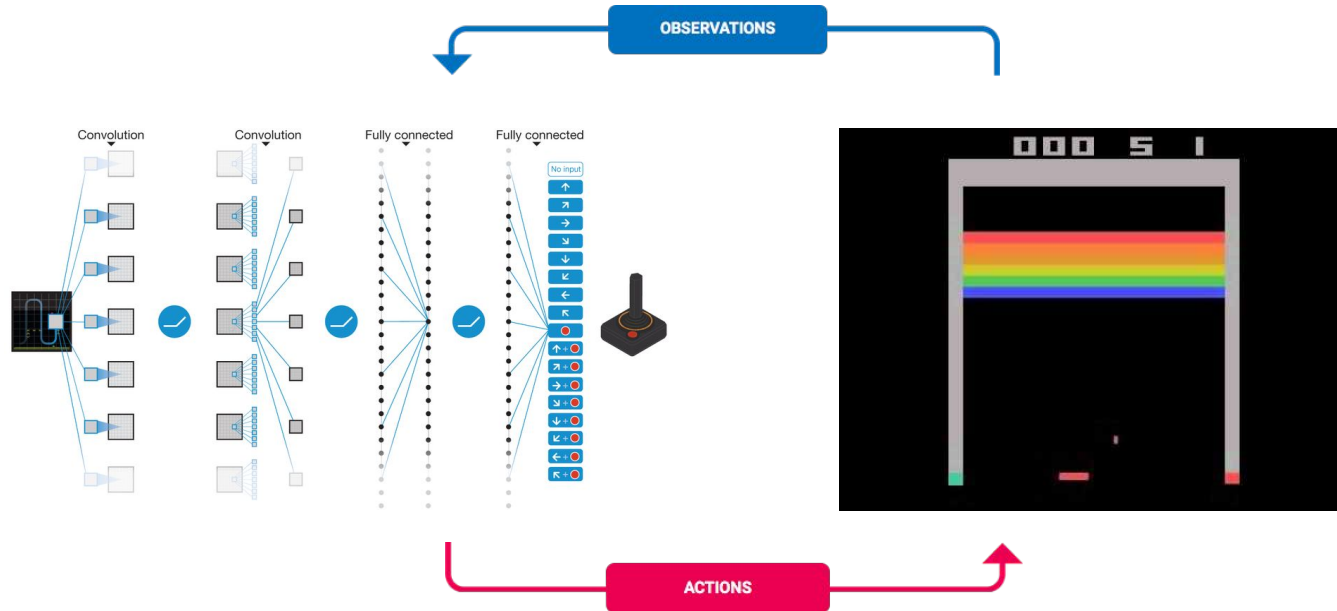# Learning paradigms



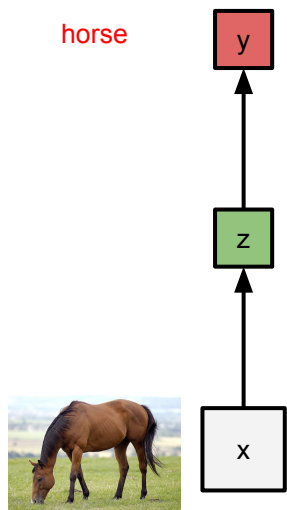**Supervised Learning**

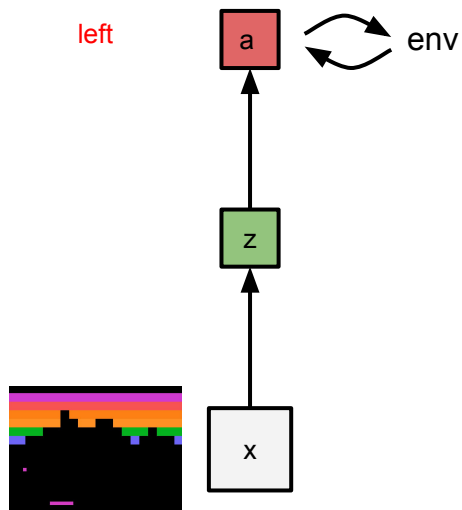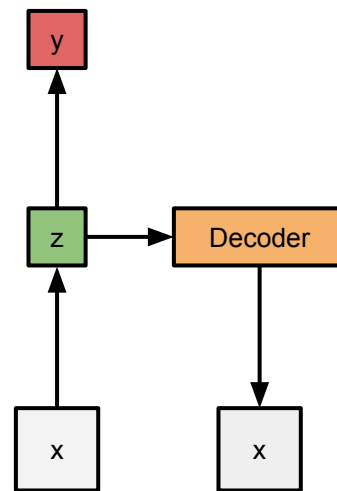**Reinforcement Learning**

# End-to-end reinforcement learning



Mnih et al. (2015)

# How much experience do we really need?
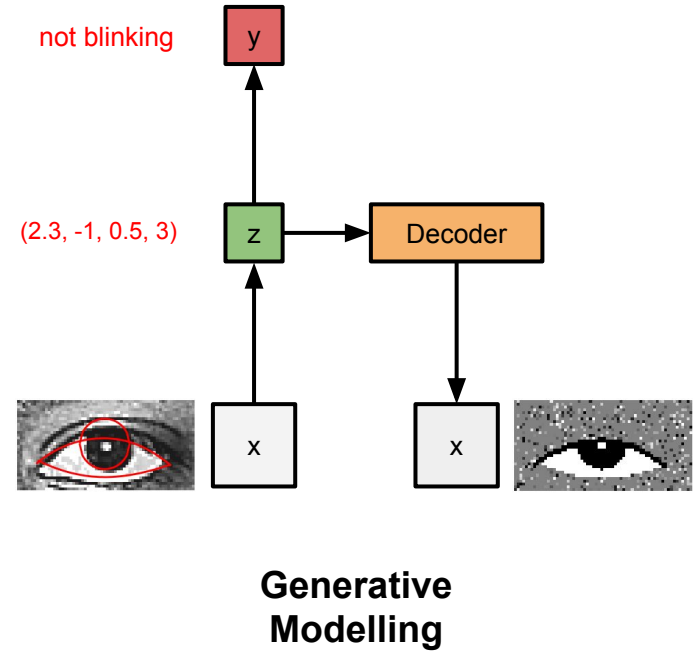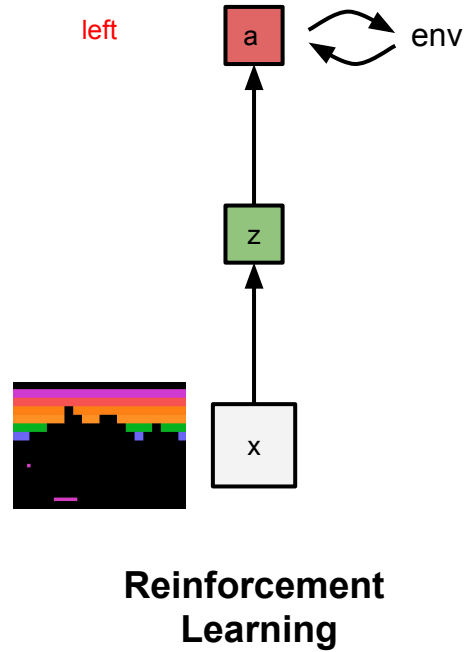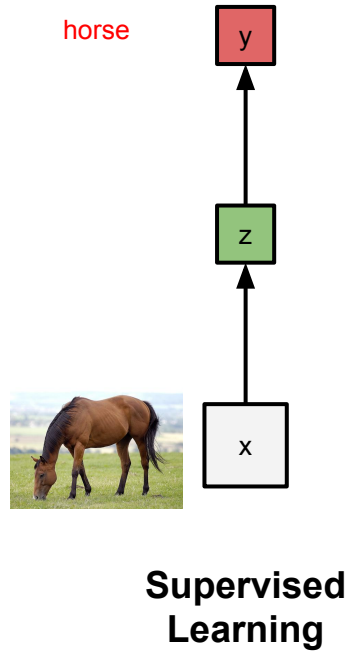
# Learning paradigms



**Supervised Learning**

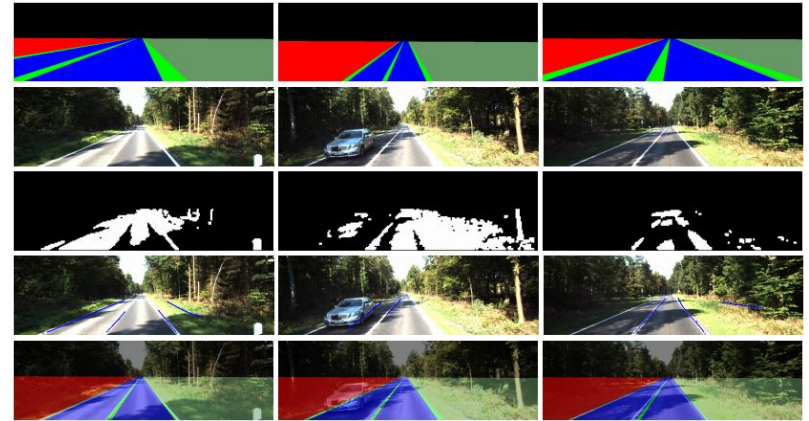**Reinforcement Learning**

**Generative Modelling**

# Learning paradigms



horse

y

z

x

**Supervised Learning**

left          env

a

z

x

**Reinforcement Learning**

not blinking          y

(2.3, -1, 0.5, 3)          z          Decoder

x          x
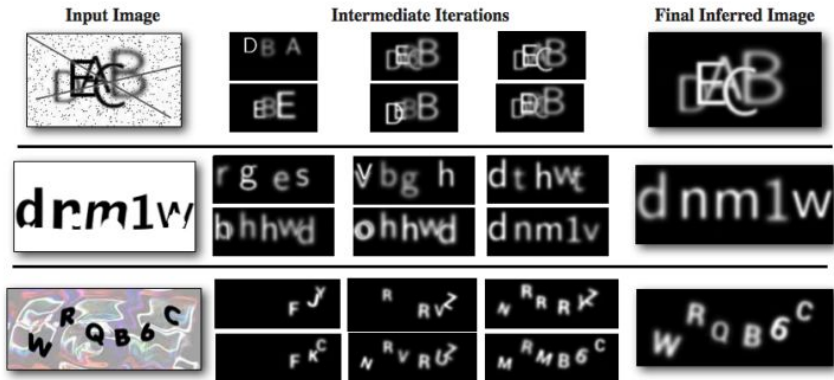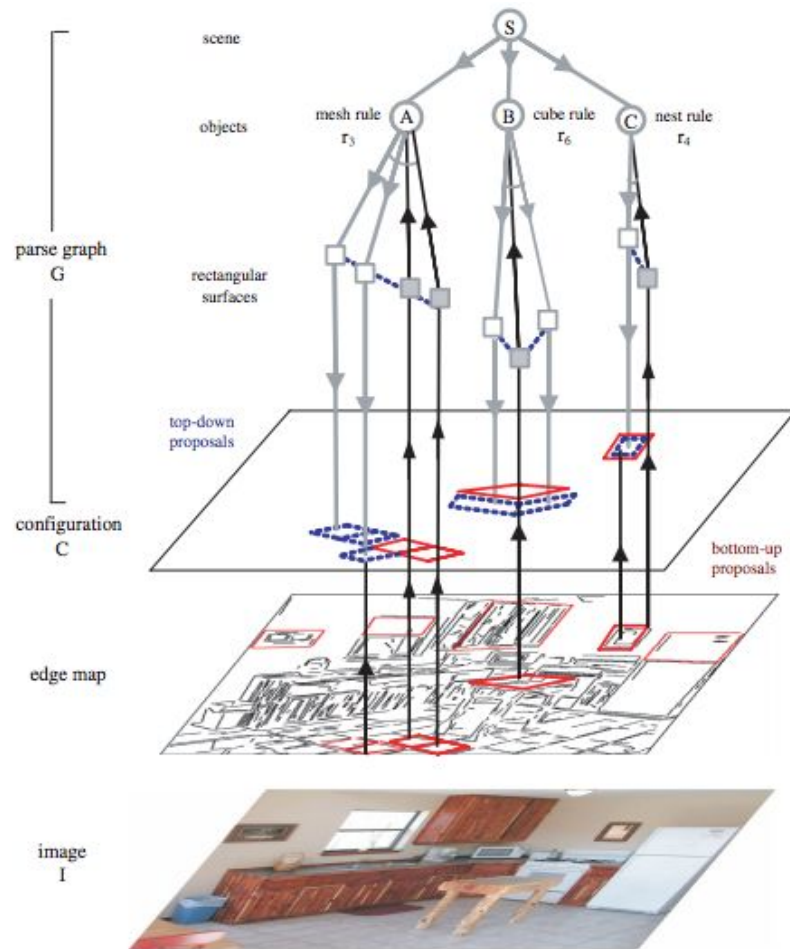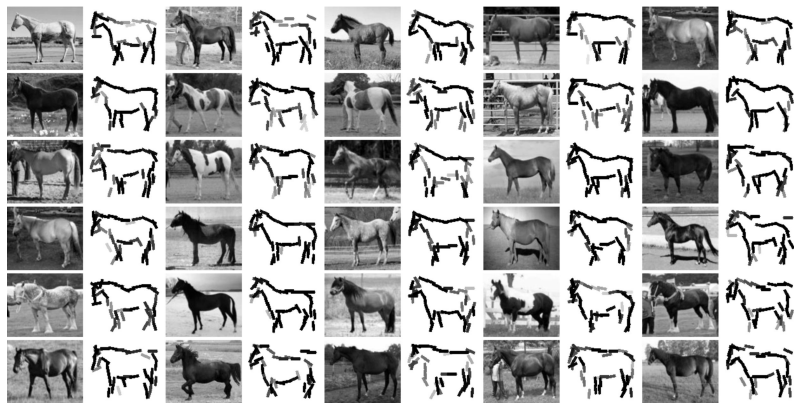
**Generative Modelling**

# Highly structured



General Purpose Graphics Programming

Vikash Mansinghka, Tejas D. Kulkarni, Yura N. Perov, and Joshua B. Tenenbaum (2013)
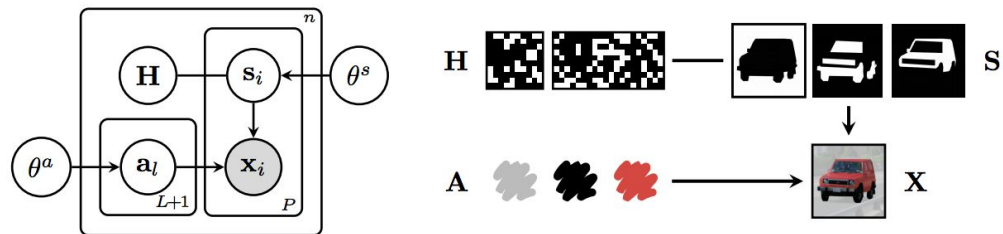
# Partially structured

A Stochastic Grammar of Images
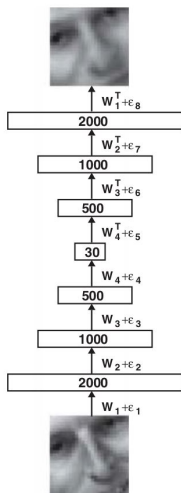
Song-Chun Zhu and David Mumford (2007)
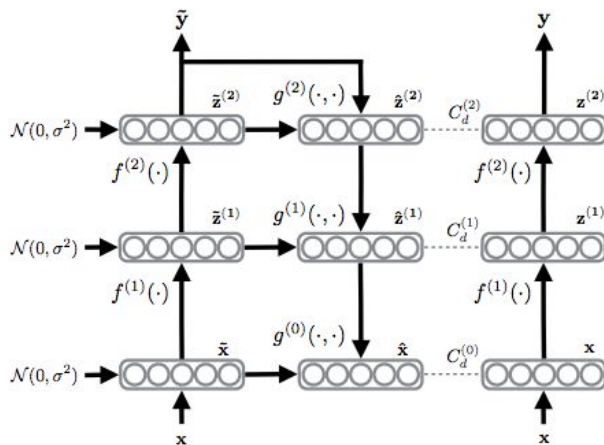
# Partially structured



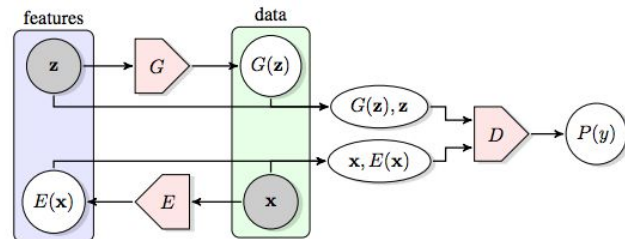S. M. Ali Eslami and Christopher K. I. Williams (2012)

# Fully unstructured



Geoffrey Hinton (2006)

Antti Rasmus et al. (2016)

Jeff Donahue et al. (2016)

# Attend, Infer, Repeat: Fast Scene Understanding with Generative Models

S. M. Ali Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, Koray Kavukcuoglu, Geoffrey Hinton
Neural Information Processing Systems (NIPS), 2016

# Motivation

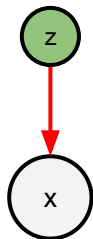To obtain **object-based** representations

To learn from orders-of-magnitude **less data**

**Cause**

blue brick

**Model**



**Image**

**Cause**

blue brick                    pile of bricks ⟵ not sufficient for
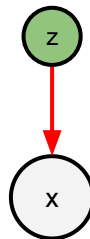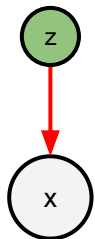                                              grasping
                                              counting
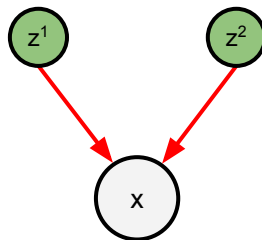                                              transfer
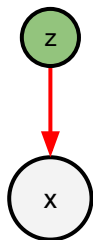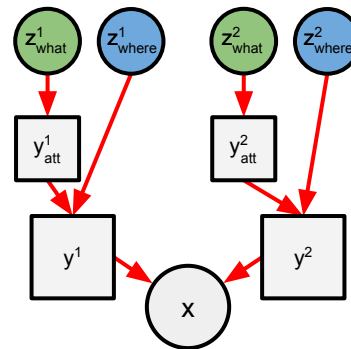                                              generalisation

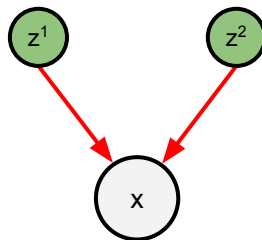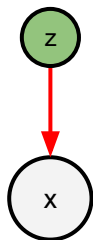**Model**



**Image**

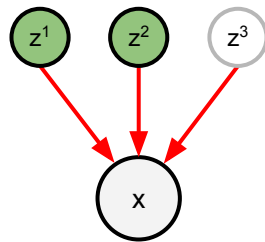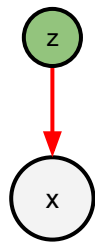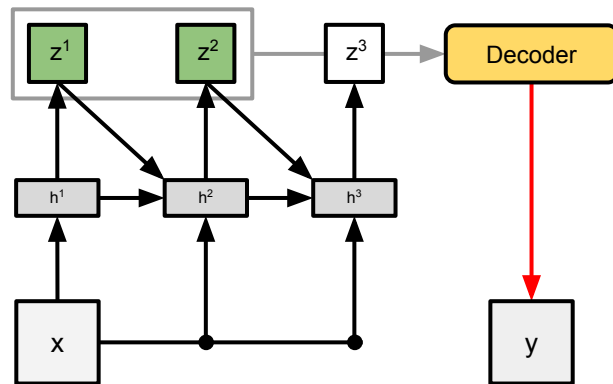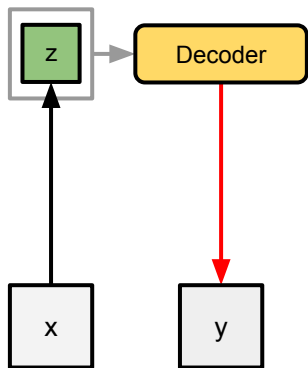| Cause | pile of bricks | blue brick | red brick | blue brick above | red brick below |

| Model | $z$ | $z^1$ $z^2$ | $z^1_{what}$ $z^1_{where}$ $z^2_{what}$ $z^2_{where}$ |

$y^1_{att}$  $y^2_{att}$

$y^1$  $y^2$

$x$  $x$  $x$

**Model**

**Inference Network**

**Model**

$z^1$ $z^2$ $z^3$

$x$

$z^1_{what}$ $z^1_{where}$ $z^2_{what}$ $z^2_{where}$

$y^1_{att}$ $y^2_{att}$

$y^1$ $y^2$

$x$

**Inference Network**

$z^1$ $z^2$ $z^3$ Decoder

$h^1$ $h^2$ $h^3$

$x$

$y$

focus on representation
not reconstruction

$z^1_{pres}$ $z^1_{what}$ $z^1_{where}$ $z^2_{pres}$ $z^2_{what}$ $z^2_{where}$ $z^3_{pres}$ $z^3_{what}$ $z^3_{where}$ Decoder

$h^1$ $h^2$ $h^3$

$x$

$y$
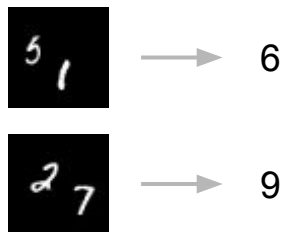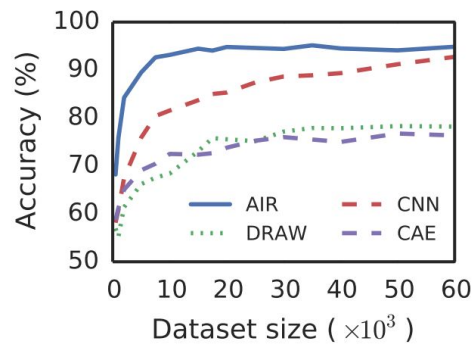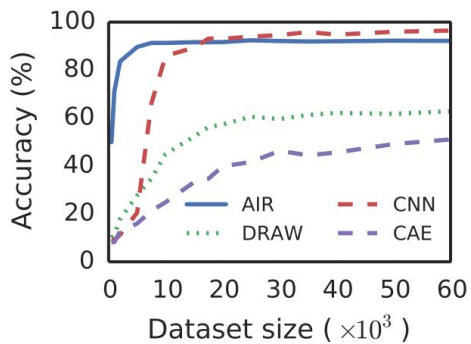
output is a set
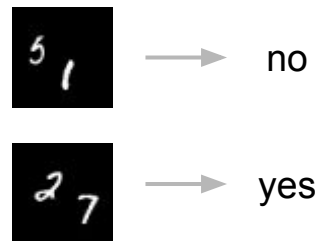order? count?

# Demo reel
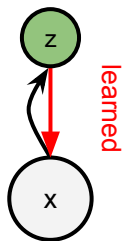
# Omniglot

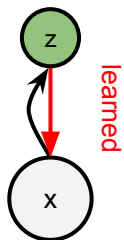# Representational power

Sum?



Increasing order?

# Additional structure

distributed **vector**
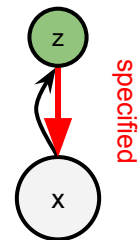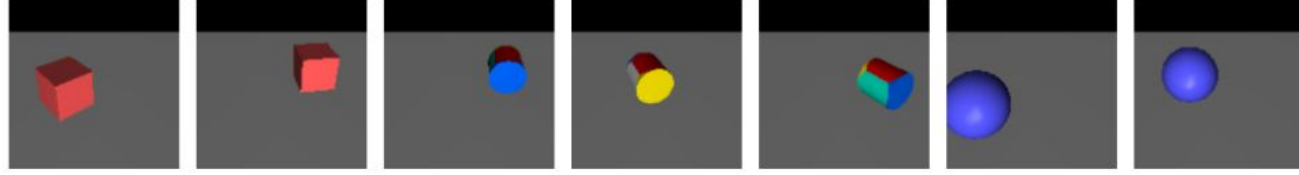that correlates
with blue brick

# Additional structure

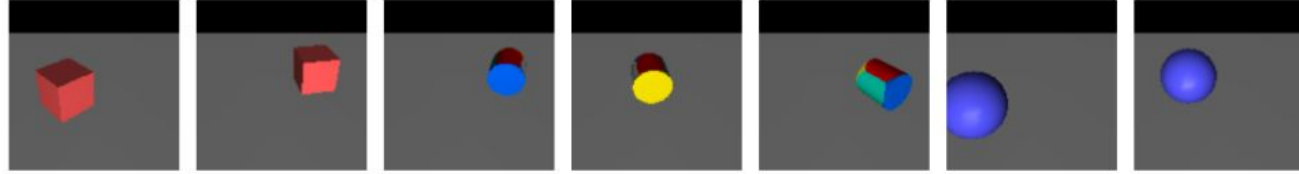distributed **vector**
that correlates
with blue brick

class=**brick**
colour=**blue**
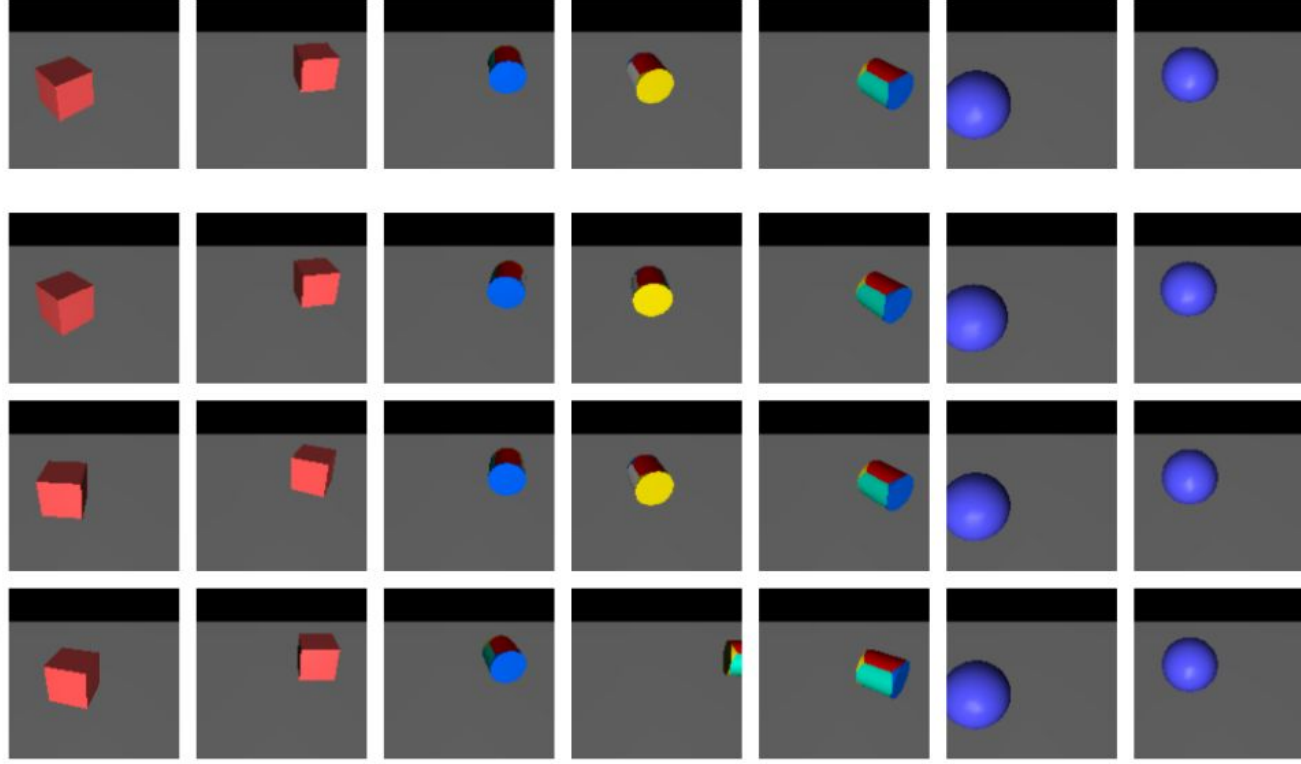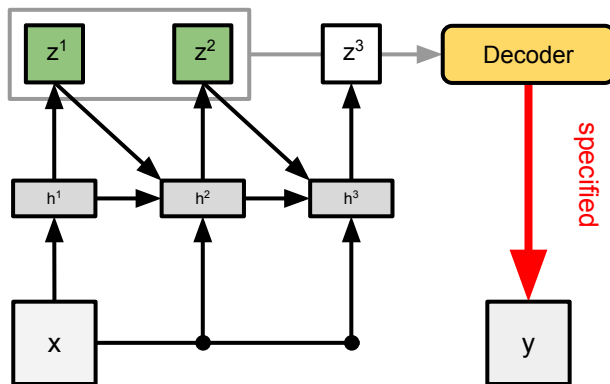position=**P**
rotation=**R**

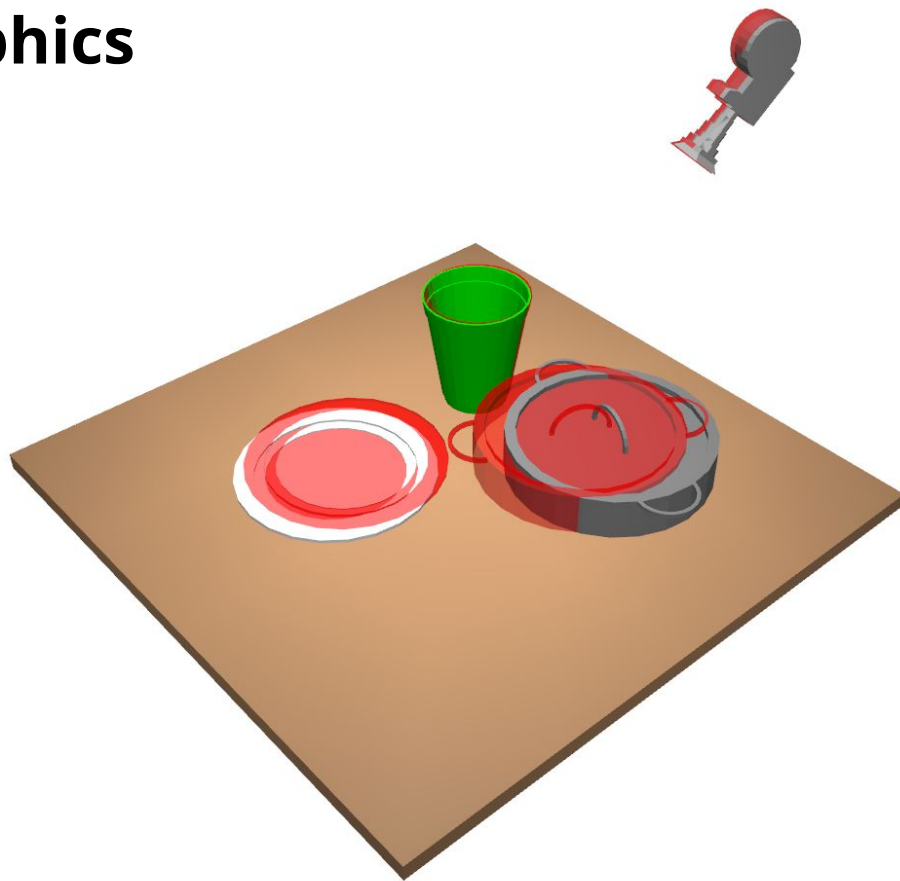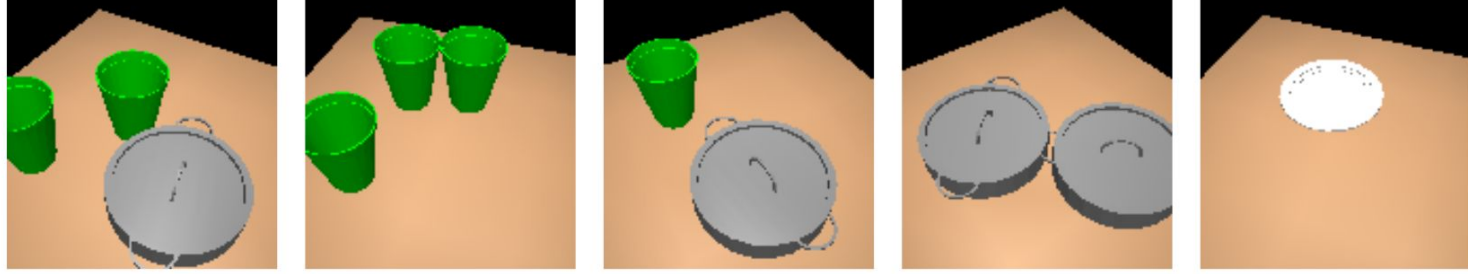*(a)* Data

*(b)* AIR

# Additional structure

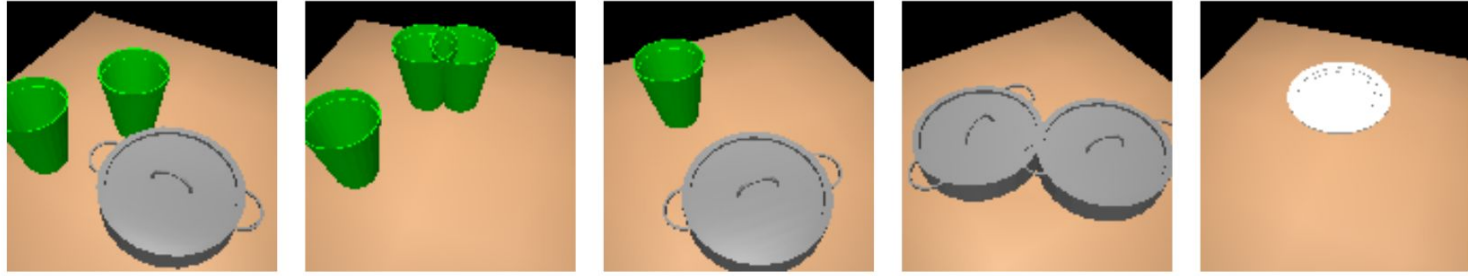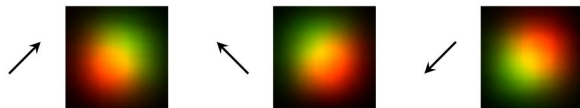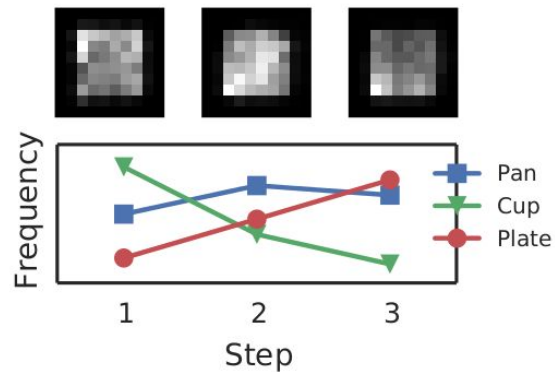# Inverse graphics

*(a)* Data

*(b)* Reconstruction

# Policy learning

# Unsupervised Learning of 3D Structure from Images

Danilo Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, Nicolas Heess
Neural Information Processing Systems (NIPS), 2016

# Motivation

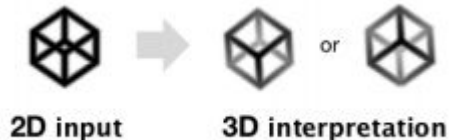To recover **3D structure** from **2D images**

To form **stable** representations, regardless of camera position

# Motivation

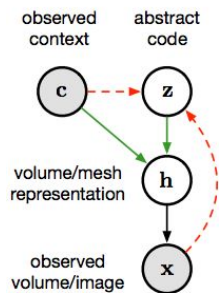To recover **3D structure** from **2D images**

To form **stable** representations, regardless of camera position

- Inherently ill-posed
  - All objects appear under self occlusion, infinite explanations
  - Therefore build statistical models to know what's likely and what's not

- Even with models, inference is intractable
  - Important to capture multi-modal explanations

- How are 3D scenes best represented?
  - Meshes or voxels?

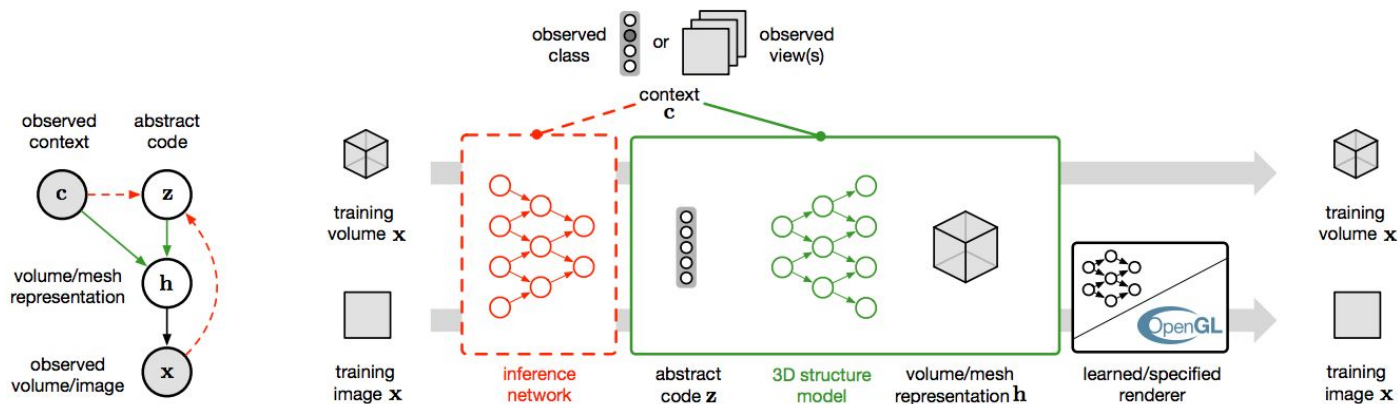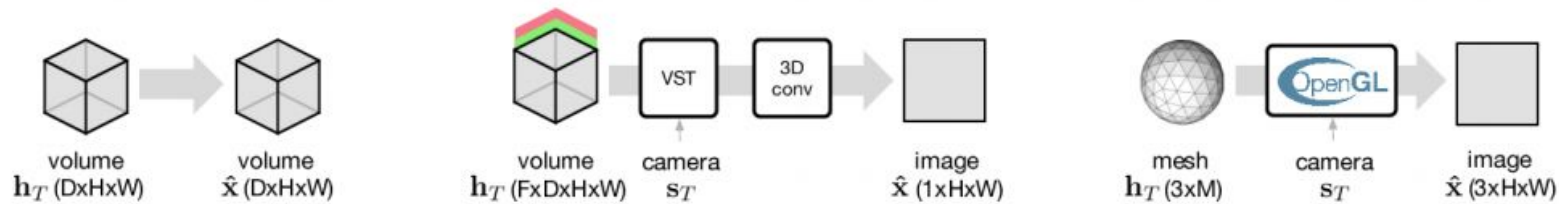- Where is training data collected from?



2D input → 3D interpretation

# Unsupervised Learning of 3D Structure from Images

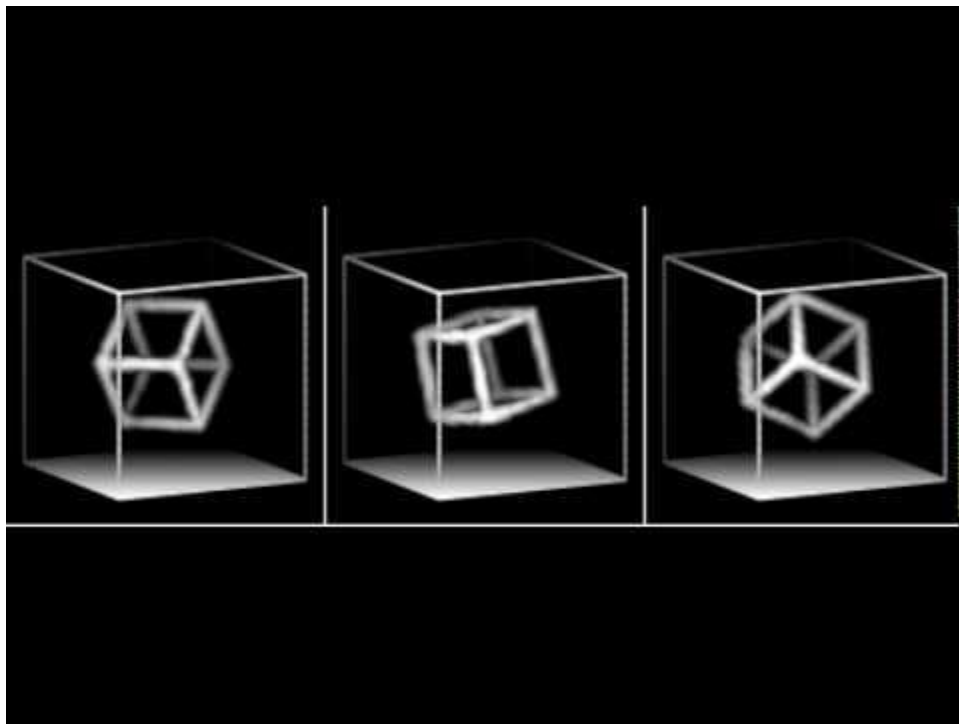# Unsupervised Learning of 3D Structure from Images

# Projection operators



volume
$\mathbf{h}_T$ (DxHxW)

volume
$\hat{\mathbf{x}}$ (DxHxW)

VST

3D conv

volume
$\mathbf{h}_T$ (FxDxHxW)

camera
$\mathbf{s}_T$

image
$\hat{\mathbf{x}}$ (1xHxW)

OpenGL

mesh
$\mathbf{h}_T$ (3xM)

camera
$\mathbf{s}_T$

image
$\hat{\mathbf{x}}$ (3xHxW)

# Unconditional samples

# Class-conditional samples

# Class-conditional samples

# Multi-modality of inference

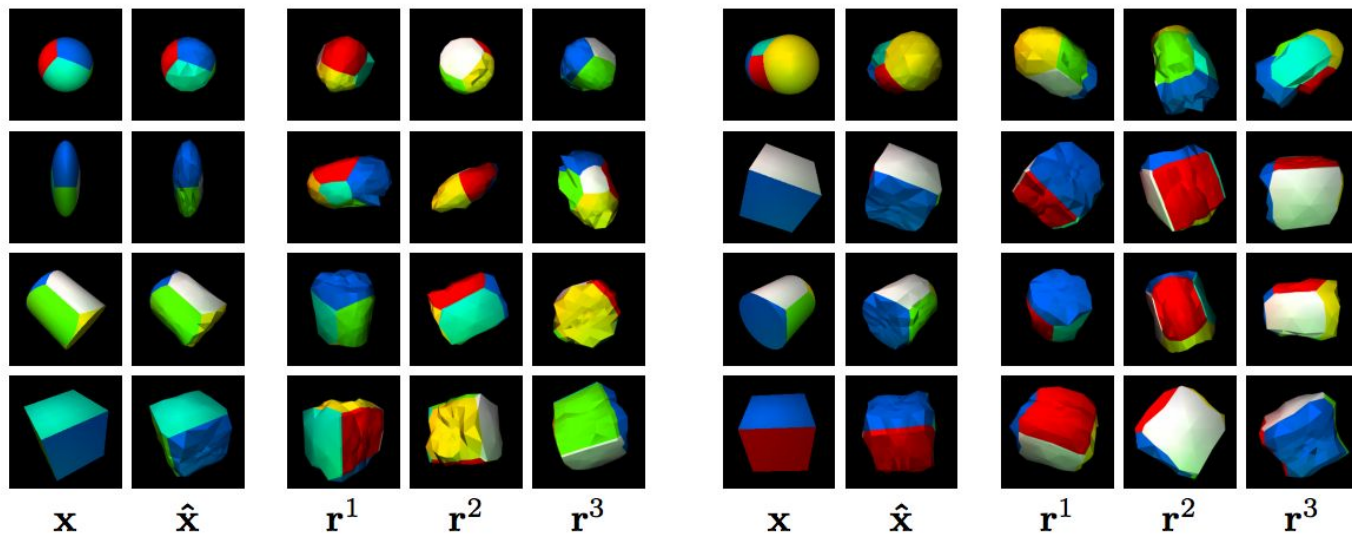# 3D structure from multiple 2D images

# Inferring object meshes



$\mathbf{x}$  $\hat{\mathbf{x}}$  $\mathbf{r}^1$  $\mathbf{r}^2$  $\mathbf{r}^3$   $\mathbf{x}$  $\hat{\mathbf{x}}$  $\mathbf{r}^1$  $\mathbf{r}^2$  $\mathbf{r}^3$

# Inferring object meshes

# Recap

- Deep Supervised Learning

- Deep Reinforcement Learning

- Model-based Methods

- Structured / Unstructured Generative Models

aeslami@google.com
arkitus.com