# PyTorch live talks

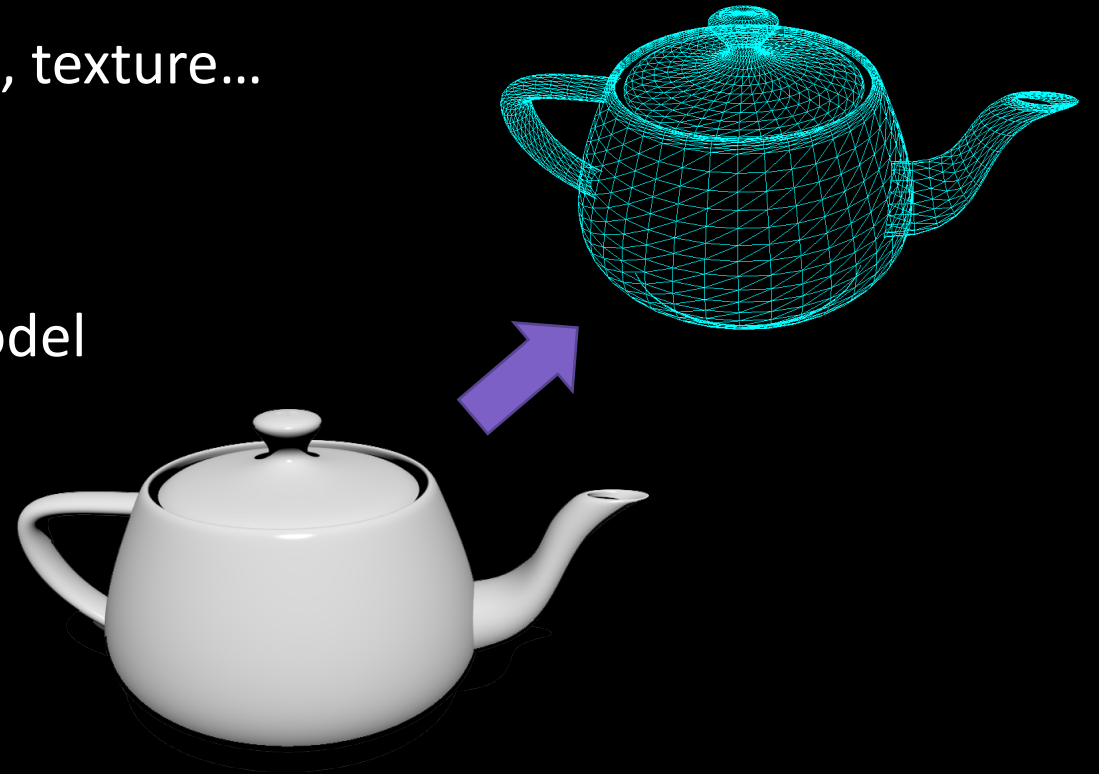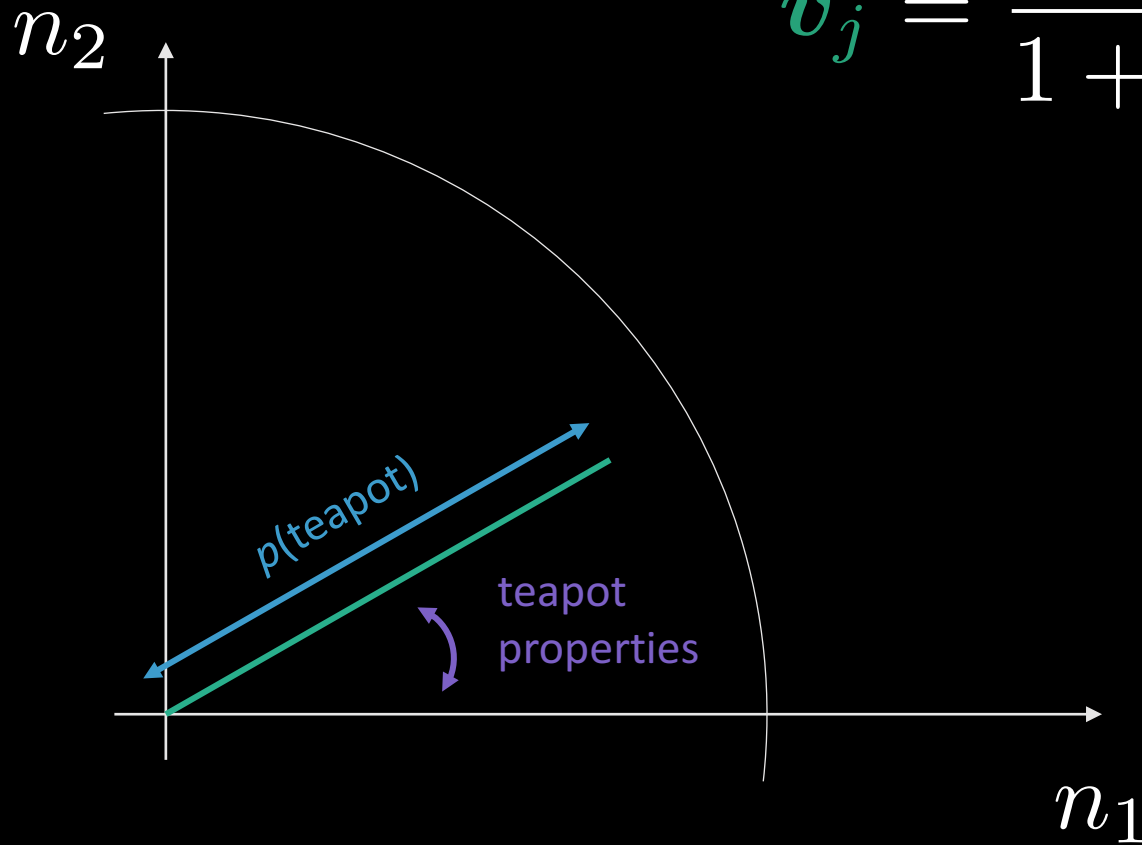Capsules and routing techniques

# Rationale

- CNN *can* deal with *translation*
- For other *affine transformations*
  - Exponential replicas of spatial feature detectors
  - Exponentially more labelled data


- Solution
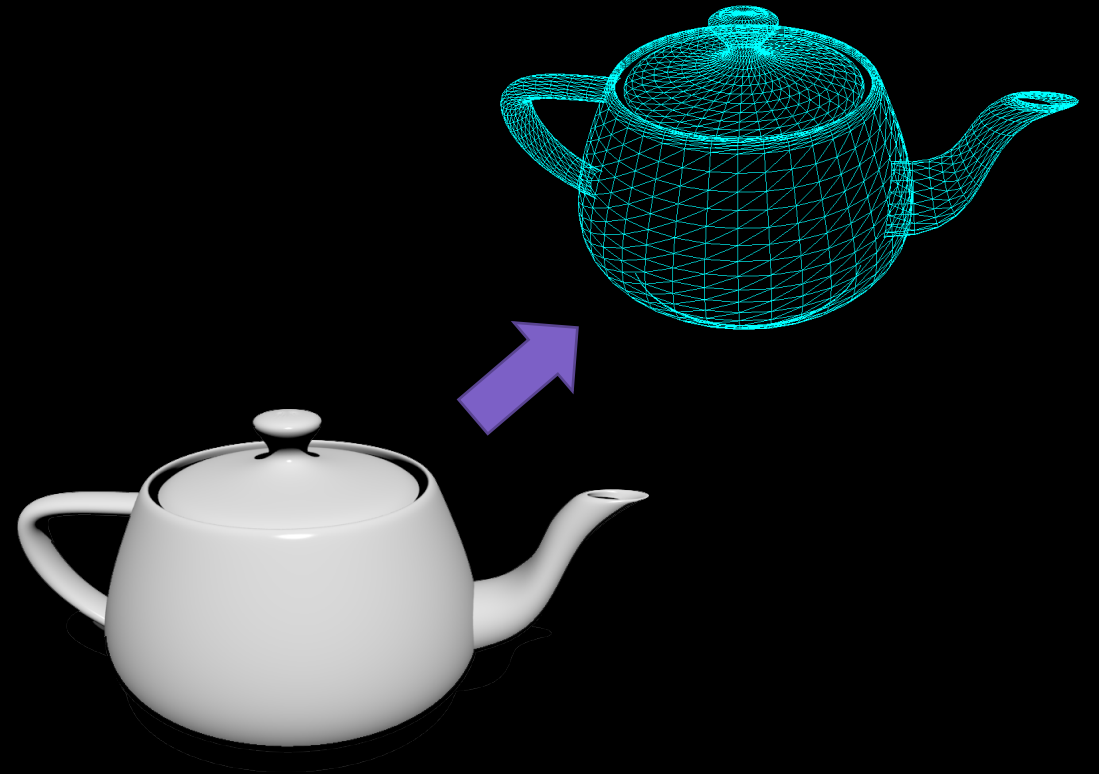  - Efficiently encode *viewpoint invariant knowledge*

# Capsules

- Groups of neurons characterising an entity in the image
- Properties include
  - Pose, deformation, velocity, albedo, hue, texture...
- Role
  - Invert the rendering process
  - 2D camera projection → 3D abstract model

# Capsules

$$\boldsymbol{v}_j = \frac{\|\boldsymbol{s}_j\|^2}{1 + \|\boldsymbol{s}_j\|^2} \frac{\boldsymbol{s}_j}{\|\boldsymbol{s}_j\|}$$



$n_2$

$n_1$

*p*(teapot)

teapot properties

# Capsules

$$v_j = \frac{\|s_j\|^2}{1 + \|s_j\|^2} \frac{s_j}{\|s_j\|}$$

Capsule output

$$s_j = \sum_i c_{ij} \hat{u}_{j|i}$$

Capsule input

$$\hat{u}_{j|i} = W_{ij} u_i$$

Prediction vectors

learnt by back-prop

$v_j$

Squash

$s_j$

Routing

$\hat{u}_{j|1}$  $\hat{u}_{j|2}$  $\hat{u}_{j|i}$

"Scale rotate"

$u_1$  $u_2$  $u_i$

Alfredo Canziani © 2017

# Dynamic routing

$$j = 1, \cdots, s_{\ell+1}$$

$$b_i \leftarrow 0, i = 1, \ldots, s_\ell$$

$$c_i \leftarrow \mathrm{softmax}(b_i)$$

$$s_j = \sum_i c_{ij} \hat{u}_{j|i}$$

$$v_j \leftarrow \mathrm{squash}(s_j)$$

$$b_{ij} \leftarrow b_{ij} + \hat{u}_{j|i}^\top v_j$$

$\ell + 1$

$v_j$

Squash

$s_j$

Routing

$\hat{u}_{j|1}$  $\hat{u}_{j|2}$  $\hat{u}_{j|i}$

$\ell$

$s_\ell$ : size of layer $\ell$

# Margin loss

$$\mathcal{L}_k{}^{(i)} = \begin{cases} \left[ \left( m_+ - \|\boldsymbol{v}_k^{(i)}\| \right)^+ \right]^2, & \text{digit } k \text{ preset} \\ \lambda \left[ \left( \|\boldsymbol{v}_k^{(i)}\| - m_- \right)^+ \right]^2, & \text{otherwise} \end{cases}$$
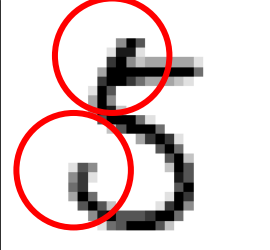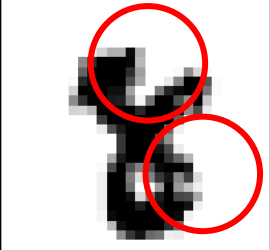
$$\mathcal{L}_\mu^{(i)} = \sum_k \mathcal{L}_k{}^{(i)}$$

$\uparrow$

margin

# CapsNet

$x^{(i)}$

ReLU Conv1 — 256

9X9

9X9 — 20

PrimaryCaps — 8

32

6

$W_{st} \in \mathbb{R}^{8 \times 16}$

16 — DigitCaps — 10

$\mathcal{L}_\mu^{(i)}$

10

16 — DigitCaps — 10

FC ReLU — 512

FC ReLU — 1024

FC Sigmoid — 784

$$\mathcal{L}^{(i)} = \mathcal{L}_\mu^{(i)} + \rho \cdot \mathcal{L}_\rho^{(i)}$$

margin        reconstruction

= 0 Masked        = Representation of the reconstruction target

# Results (I)



| $(l, p, r)$ | $(2, 2, 2)$ | $(5, 5, 5)$ | $(8, 8, 8)$ | $(9, 9, 9)$ | $(5, 3, 5)$ | $(5, 3, 3)$ |
|---|---|---|---|---|---|---|
| Input | | | | | | |
| Output | | | | | | |

# Results (II)

| Scale and thickness |  |
| --- | --- |
| Localized part |  |
| Stroke thickness |  |
| Localized skew |  |
| Width and translation |  |
| Localized part |  |