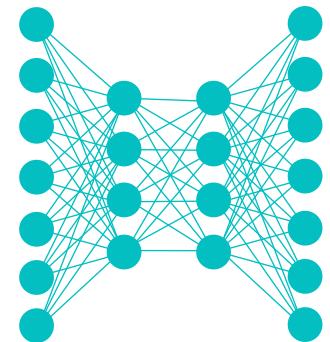Lecture Notes for

# Neural Networks and Machine Learning

Vision Transformers
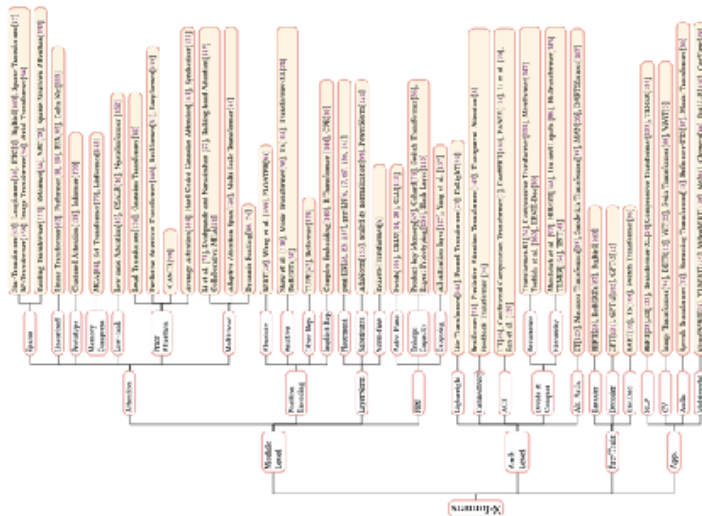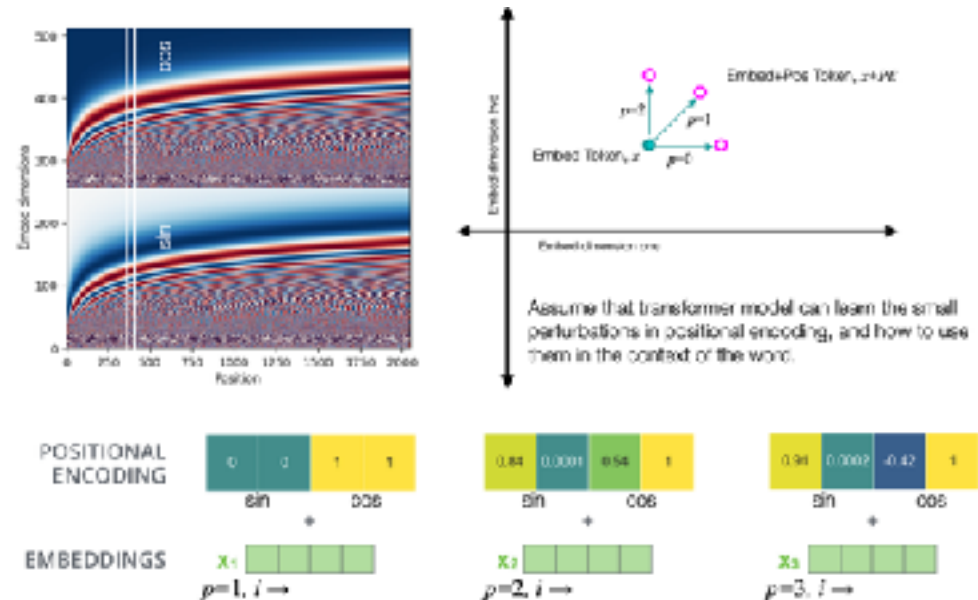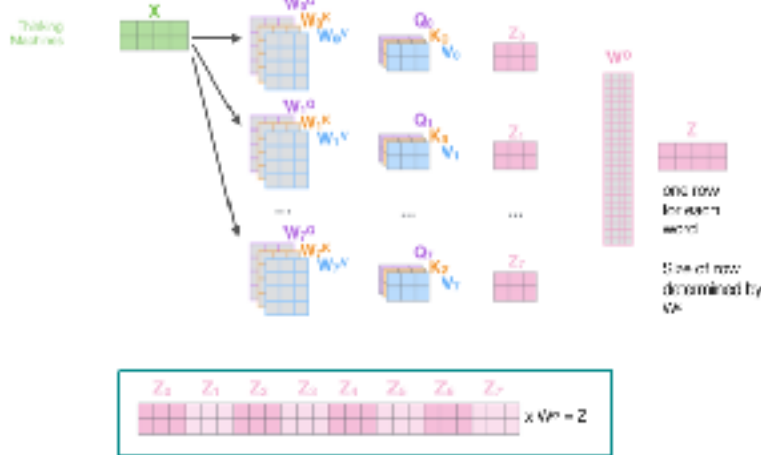Self-supervised Learning
Consistency Loss

# Logistics and Agenda

- Logistics
  - None!

- Agenda
  - Vision Transformers
  - Self Supervised Learning and Consistency Loss
- Next Time:
  - Paper Presentation: *Language Models are Few Shot Learners*
  - Multi-modal Learning
    - Techniques
    - Applications and domains
  - Multi-Task and Demo

# Last Time: Transformers

# Vision Transformers

# Vision Transformers

- Divide image into patches
  - Treat each patch as something to encode separately
  - Flatten each patch
  - Put through dense layer
- Add positional encoding based on position of patch
  - for 7x7 patch, there are 49 positions
- Put into transformer. Same as text transformers …
- **But you need a lot of data**
  - 14M or more images seems to be sweet spot

Position embedding similarity
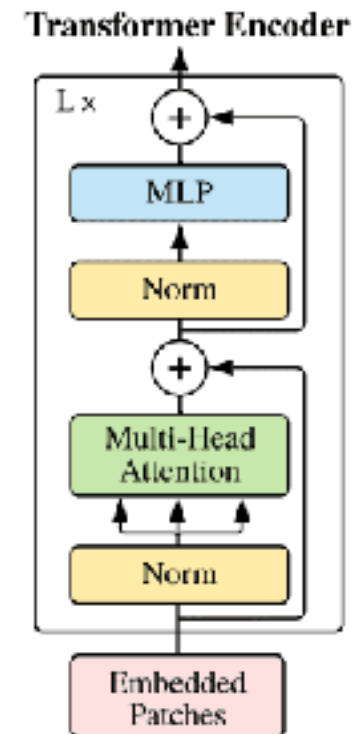
53

# Vision Transformers Video

# ViT Architectures

- *D* is size of patch embedding
- Uses skip connections (all size *D*)
- Multi-headed self attention (MSA) takes *D* input `patch_embed + pos_embed`
- Main difference in architectures
  - *L* blocks used (*i.e.*, "layers")
  - *H* heads in each layer (*i.e.*, "heads")
  - MLP head is final classifier

$$\mathbf{z}_0 = [\mathbf{x}_{\text{class}}; \mathbf{x}_p^1 \mathbf{E}; \mathbf{x}_p^2 \mathbf{E}; \cdots ; \mathbf{x}_p^N \mathbf{E}] + \mathbf{E}_{\text{pos}},$$
$$\mathbf{z}'_\ell = \text{MSA}(\text{LN}(\mathbf{z}_{\ell-1})) + \mathbf{z}_{\ell-1},$$
$$\mathbf{z}_\ell = \text{MLP}(\text{LN}(\mathbf{z}'_\ell)) + \mathbf{z}'_\ell,$$
$$\mathbf{y} = \text{LN}(\mathbf{z}_L^0)$$

**Transformer Encoder**



| Model | Layers | Hidden size $D$ | MLP size | Heads | Params |
|-------|--------|-----------------|----------|-------|--------|
| ViT-Base | 12 | 768 | 3072 | 12 | 86M |
| ViT-Large | 24 | 1024 | 4096 | 16 | 307M |
| ViT-Huge | 32 | 1280 | 5120 | 16 | 632M |

ResNet50: 23M

Dosovitskiy *et al.* (2021) An image is worth 16x16 words: Transformers for Image Recognition at Scale  https://arxiv.org/abs/2010.11929

# *What is the learnable class embedding?*



Transformer actually has a lot of outputs, same as num patch inputs

Class
Bird
Ball
Car
...

MLP Head

Each output is Conditioned on all other inputs!
*Does it matter the position of this class token?*
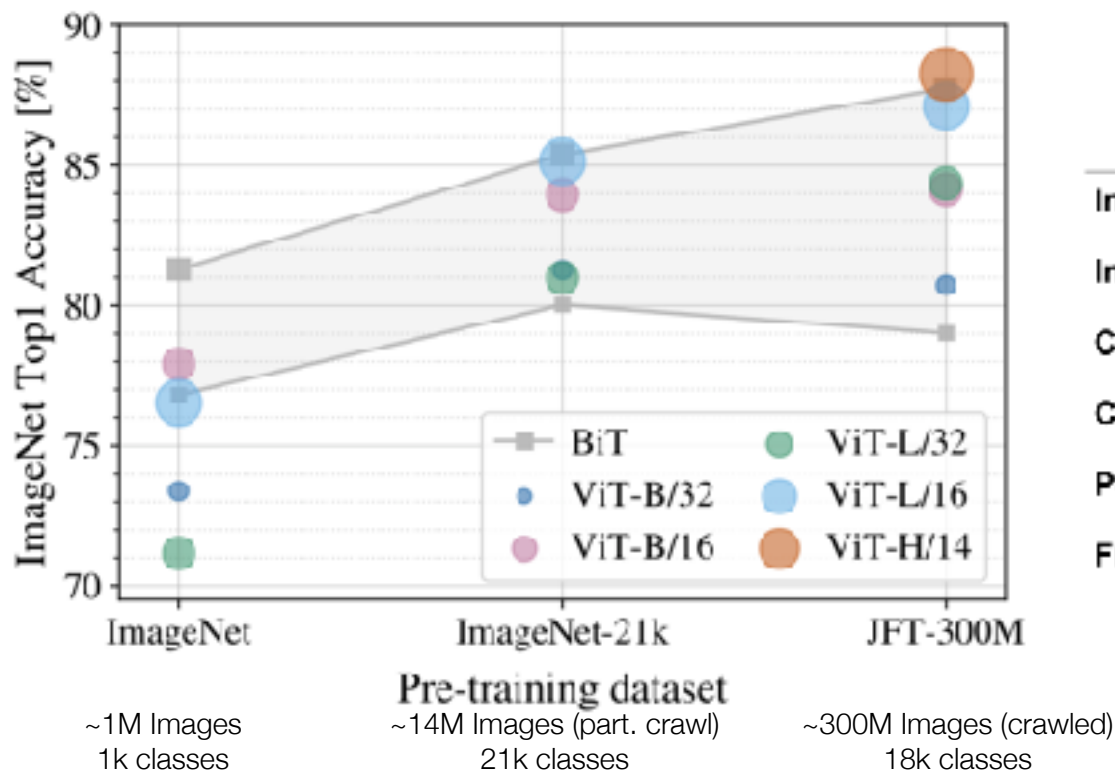
Transformer Encoder

Extra learnable [class] embedding

Linear Projection of Flattened Patches

# Do they work?

- Yes, but good luck getting weights for them or training them, even with the SuperPOD

- Less than 14M images for pre-training? Use ResNet.



**Transfer Learning From Huge ViT**

|  | ViT-H | Previous SOTA |
|---|---|---|
| ImageNet | 88.55 | 88.5 |
| ImageNet-ReaL | 90.72 | 90.55 |
| Cifar-10 | 99.50 | 99.37 |
| Cifar-100 | 94.55 | 93.51 |
| Pets | 97.56 | 96.62 |
| Flowers | 99.68 | 99.63 |

~1M Images
1k classes

~14M Images (part. crawl)
21k classes

~300M Images (crawled)
18k classes

# Self-Supervised Learning



From Yoshua Bengio

Three challenges for Deep Learning
- Deep Supervised Learning works well for perception
  - When labeled data is abundant.
- Deep Reinforcement Learning works well for action generation
  - When trials are cheap, e.g. in simulation.
- Three problems the community is working on:
- 1. Learning with fewer labeled samples and/or fewer trials
  - Self-supervised learning / unsup learning / learning to fill in the blanks
    - learning to represent the world before learning tasks
- 2. Learning to reason, beyond "system 1" feed forward computation
  - Making reasoning compatible with gradient-based learning.
- 3. Learning to plan complex action sequences
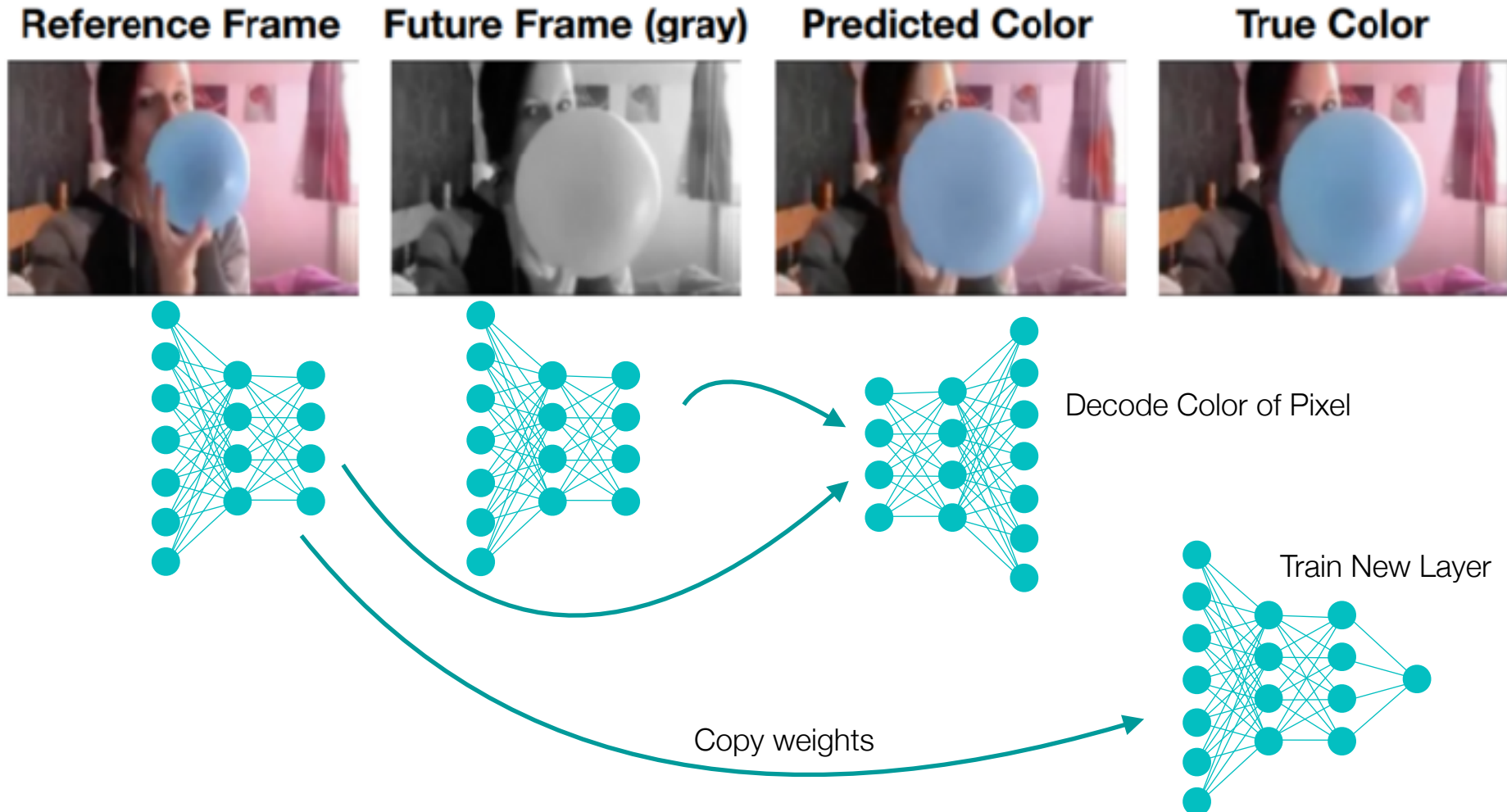  - Learning hierarchical representations of action plans

# Self-supervised Learning

- **Problem**: deep learning is not sample efficient
- **Idea**: learn about the world before learning the task
- **New Problem**: how do we learn about the world?
- **Solution**: transfer learning on toy problem
  - 1. train on auxiliary task that is easy to label
  - 2. throw away anything specific to auxiliary task
  - 3. train new network with task of interest, transferring knowledge (downstream task)
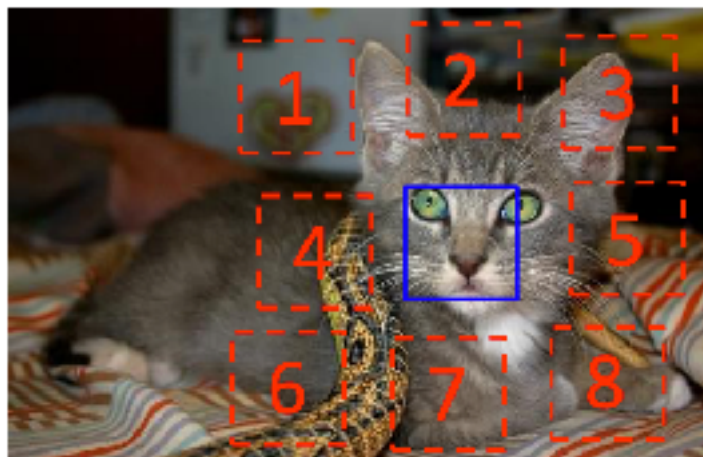  - 4. profit

# Examples of Self Supervised Learning



Reference Frame    Future Frame (gray)    Predicted Color    True Color

Decode Color of Pixel

Train New Layer

Copy weights

X = ( , ); Y = 3

**Unsupervised Visual Representation Learning by Context Prediction**

Carl Doersch[1,2]   Abhinav Gupta[1]   Alexei A. Efros[2]

[1] School of Computer Science
Carnegie Mellon University

[2] Dept. of Electrical Engineering and Computer Science
University of California, Berkeley

https://www.fast.ai/2020/01/13/self_supervised/

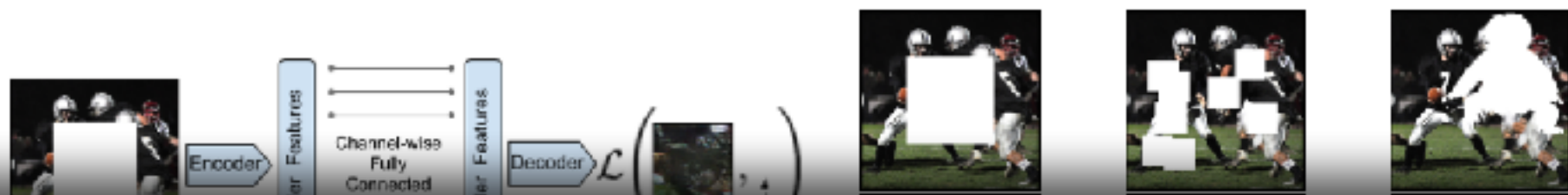| Pretraining Method | Supervision | Pretraining time | Classification | Detection | Segmentation |
|---|---|---|---|---|---|
| ImageNet [26] | 1000 class labels | 3 days | 78.2% | 56.8% | 48.0% |
| Random Gaussian | initialization | < 1 minute | 53.3% | 43.4% | 19.8% |
| Autoencoder | - | 14 hours | 53.8% | 41.9% | 25.2% |
| Agrawal et al. [1] | egomotion | 10 hours | 52.9% | 41.8% | - |
| Wang et al. [39] | motion | 1 week | 58.7% | 47.4% | - |
| Doersch et al. [7] | relative context | 4 weeks | 55.3% | 46.6% | - |
| Ours | context | 14 hours | 56.5% | 44.5% | 30.0% |



Doesn't always work to increase performance…

**Context Encoders: Feature Learning by Inpainting**

Deepak Pathak    Philipp Krähenbühl    Jeff Donahue    Trevor Darrell    Alexei A. Efros
University of California, Berkeley

https://www.fast.ai/2020/01/13/self_supervised/

# Consistency Loss



I'm from Canada, but live in the States now.

It took me a while to get used to writing boolean variables with an "Is" prefix, instead of the "Eh" suffix that Canadians use when programming.
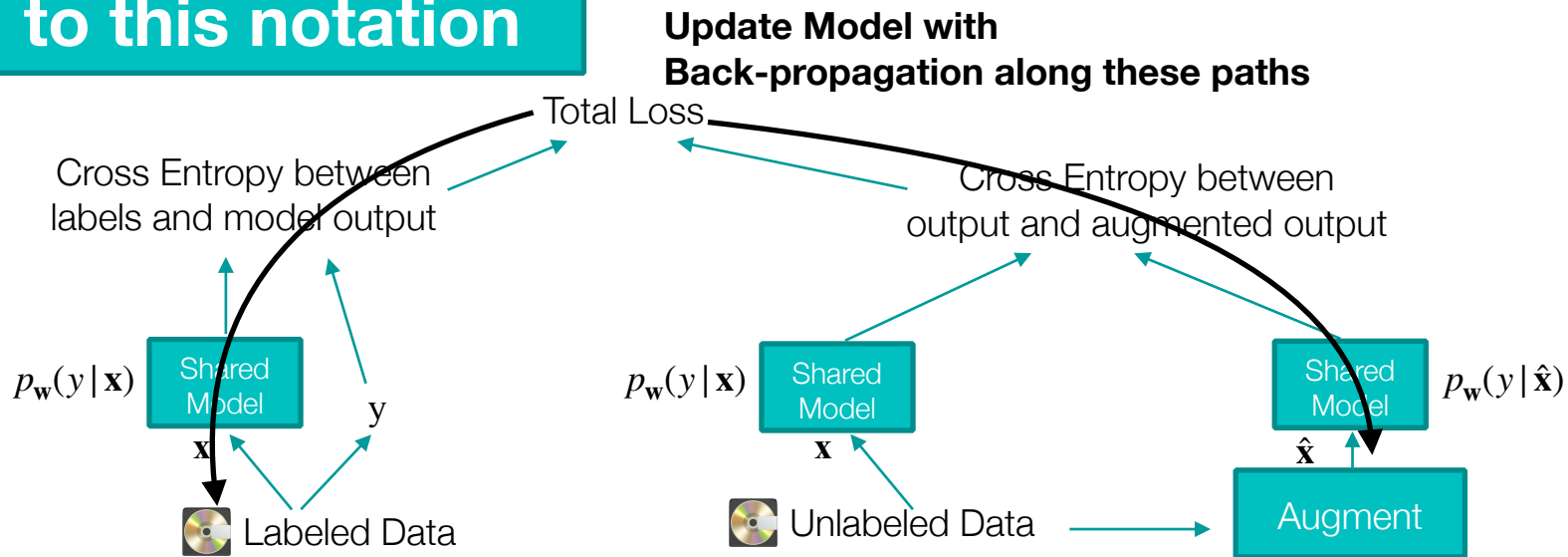
For example:

```
MyObj.IsVisible

MyObj.VisibleEh
```

# Unsupervised Consistency Loss

$$\min_{\mathbf{w}} \overbrace{\mathbf{E}_{\mathbf{x},y \in L}[-\log p_{\mathbf{w}}(y \,|\, \mathbf{x})]}^{\text{cross entropy}} + \lambda \overbrace{\mathscr{D}_{KL}\left(p_{\mathbf{w}}(y \,|\, \mathbf{x}) \,||\, p_{\mathbf{w}}(y \,|\, \hat{\mathbf{x}})\right)}^{\text{consistency in augmentation}}$$

**no** back prop      **yes** back prop

Neural Network approximates $p(y|\mathbf{x})$ by $\mathbf{w}$
Use labeled data to minimize network

Sample new $\mathbf{x}$ from unlabeled pool with function $q$
function $q$ is augmentation procedure
Minimize cross entropy of two models

**Get accustomed to this notation**

**Update Model with**
**Back-propagation along these paths**

Total Loss

Cross Entropy between
labels and model output

Cross Entropy between
output and augmented output

$p_{\mathbf{w}}(y \,|\, \mathbf{x})$  Shared Model

y

$p_{\mathbf{w}}(y \,|\, \mathbf{x})$  Shared Model

Shared Model  $p_{\mathbf{w}}(y \,|\, \hat{\mathbf{x}})$

$\mathbf{x}$

$\mathbf{x}$

$\hat{\mathbf{x}}$

Labeled Data

Unlabeled Data

Augment

# Unsupervised Consistency Loss

$$\min_{\mathbf{w}} \overbrace{\mathbf{E}_{\mathbf{x},y \in L}[-\log p_{\mathbf{w}}(y \,|\, \mathbf{x})]}^{\text{cross entropy}} + \lambda \overbrace{\mathscr{D}_{KL}\left(p_{\mathbf{w}}(y \,|\, \mathbf{x}) \,||\, p_{\mathbf{w}}(y \,|\, \hat{\mathbf{x}})\right)}^{\text{consistency in augmentation}}$$



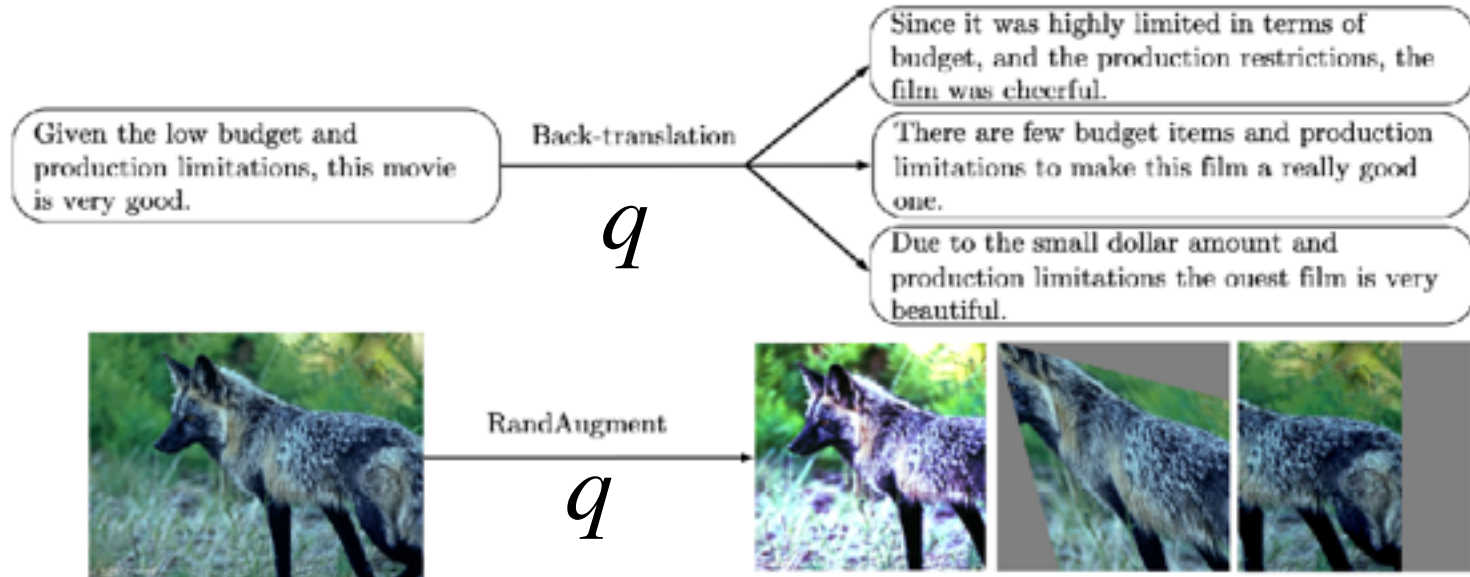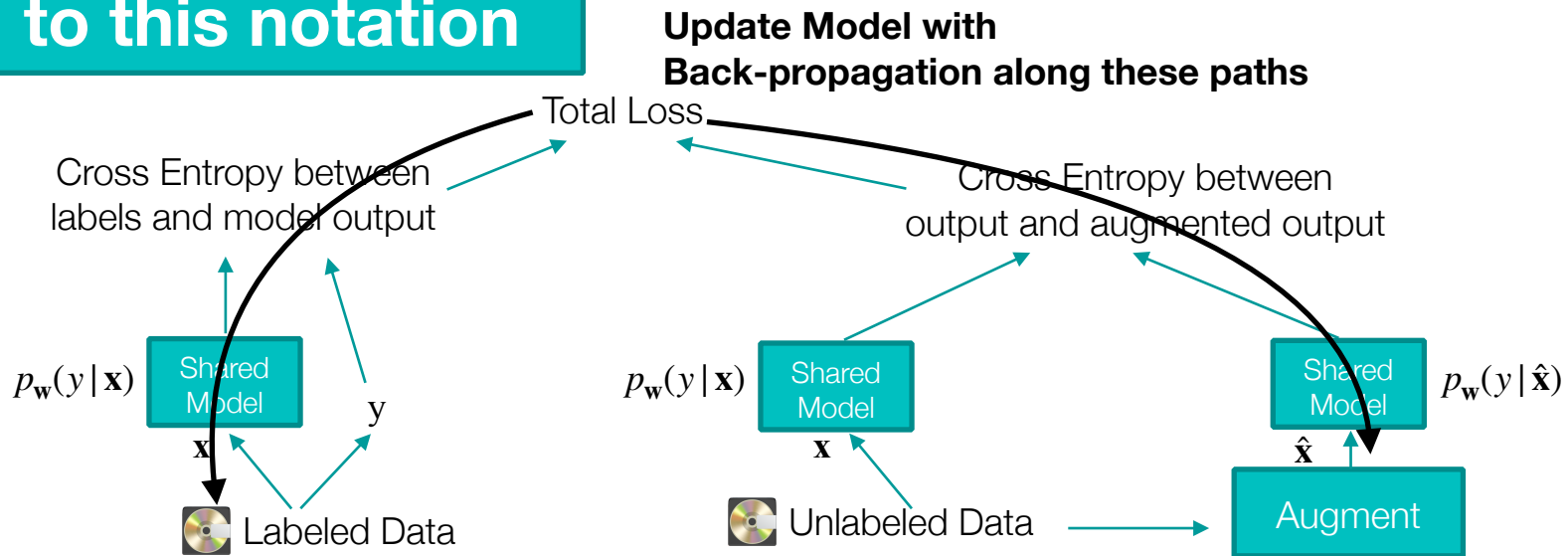Figure 2: Augmented examples using back-translation and RandAugment.

# Unsupervised Consistency Loss (review)

$$\underset{\mathbf{w}}{\min}\ \overbrace{\mathbf{E}_{\mathbf{x},y\in L}[-\log p_{\mathbf{w}}(y\,|\,\mathbf{x})]}^{\text{cross entropy}} + \lambda\ \overbrace{\mathscr{D}_{KL}\left(p_{\mathbf{w}}(y\,|\,\mathbf{x})\,||\,p_{\mathbf{w}}(y\,|\,\hat{\mathbf{x}})\right)}^{\text{consistency in augmentation}}$$

**no** back prop    **yes** back prop

Neural Network approximates $p(y|\mathbf{x})$ by $\mathbf{w}$
Use labeled data to minimize network

Sample new $\mathbf{x}$ from unlabeled pool with function $q$
function $q$ is augmentation procedure
Minimize cross entropy of two models

**Get accustomed to this notation**

**Update Model with
Back-propagation along these paths**

Total Loss

Cross Entropy between
labels and model output

Cross Entropy between
output and augmented output

$p_{\mathbf{w}}(y\,|\,\mathbf{x})$    Shared Model    $y$

$p_{\mathbf{w}}(y\,|\,\mathbf{x})$    Shared Model

Shared Model    $p_{\mathbf{w}}(y\,|\,\hat{\mathbf{x}})$

$\mathbf{x}$    Labeled Data

$\mathbf{x}$    Unlabeled Data

$\hat{\mathbf{x}}$    Augment

Unsupervised Data Augmentation (UDA) for Consistency Training, Xie et al., NeurIps 2019

$$\min_{\mathbf{w}} \overbrace{\mathbf{E}_{\mathbf{x},y \in L}[-\log p_{\mathbf{w}}(y \,|\, \mathbf{x})]}^{\text{cross entropy}} + \lambda \ \overbrace{\mathscr{D}_{KL}\left(p_{\mathbf{w}}(y \,|\, \mathbf{x}) \,||\, p_{\mathbf{w}}(y \,|\, \hat{\mathbf{x}})\right)}^{\text{consistency in augmentation}}$$

$$E[g] = \sum p(g) \cdot g \qquad \text{definition of expected value}$$

$$E[-\log p_{\mathbf{w}}(y \,|\, \mathbf{x})] = - \sum p(y) \cdot \log p_{\mathbf{w}}(y \,|\, \mathbf{x}) \qquad \text{insert -log probability, log likelihood}$$

$$NLL(y, p_{\mathbf{w}}(y \,|\, \mathbf{x})) = - \sum_{c} p(y = c) \cdot \log p_{\mathbf{w}}(y = c \,|\, \mathbf{x}) \qquad \text{negative log likelihood}$$

$$CE(f, g) = - \sum f(x) \cdot \log g(x) \qquad \text{cross entropy of two functions}$$

$$CE(y, p_{\mathbf{w}}(y \,|\, \mathbf{x})) = - \sum_{c} (y = c) \cdot \log p_{\mathbf{w}}(y = c \,|\, \mathbf{x}) \quad \text{if y=c is a probability, these are same equation}$$

```
cce = tf.keras.losses.CategoricalCrossentropy()
cce(y_true, y_pred)
```