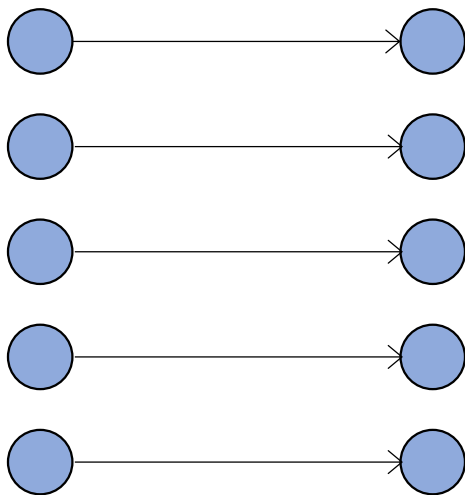# Application of Generative Models: X Learning

Hao Dong

Peking University

# From Data Point of View

Data in both input $x$ and output $y$
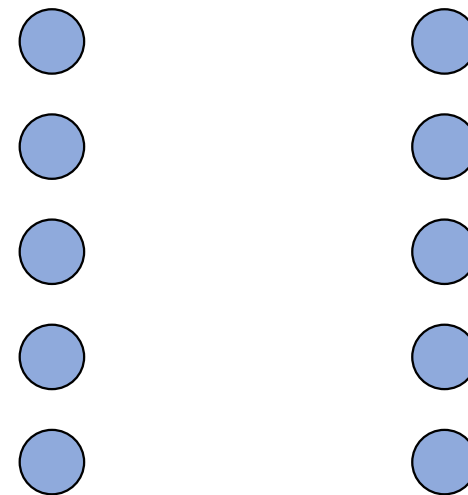with known mapping
(Learn the mapping $f$)



$$y = f(x)$$

**Supervised Learning**

- Image classification
- Object detection
- …

Data in both input $x$ and output $y$
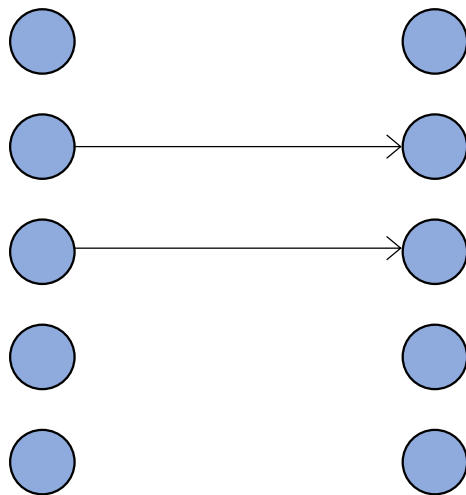(Learn the mapping $f$)



$$y = f(x)$$

**Unsupervised Learning**

- Autoencoder
  (when output is features)
- GANs
- …

# From Data Point of View

Data in both input $x$ and output $y$
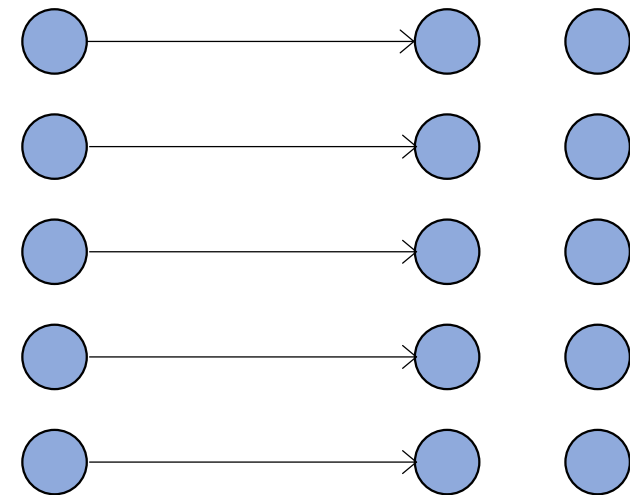with known partial mapping
(Learn the mapping $f$)



$$y = f(x)$$
**Semi-supervised Learning**
- ...

Data in both input $x$ and output $y$
with known mapping for $y$
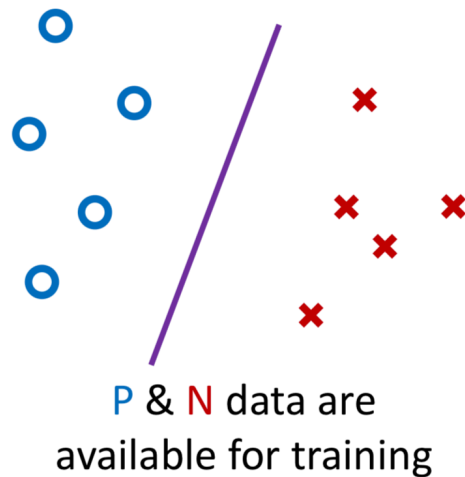(Learn the mapping $f$ for another output $y'$)



$$y' = f(x)$$
**Weakly-supervised Learning**
- Learn segmentation via classification
- ...

# From Data Point of View

**PN** learning
(i.e., supervised learning)

**PNU** learning
(i.e., semi-supervised learning)

**PU** learning
weakly-supervised learning

P & N data are
available for training

P, N & U data are
available for training

P & U data are
available for training

O : positive data          ✕ : negative data          □ : unlabeled data

From https://niug1984.github.io/paper/niu_tdlw2018.pdf

# From Mapping Point of View

Data in both input and output
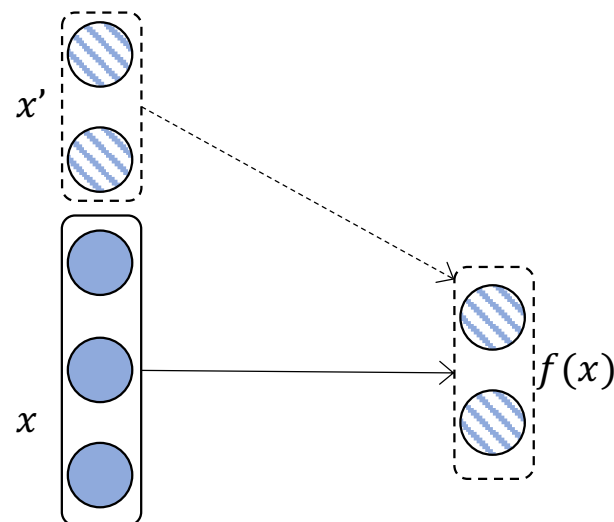(Learn the mapping $f, f'$)

$x$　　　　　　$y$

$y = f(x), x = f'(y)$
**(Unsupervised) Dual Learning**
- VAE
- CycleGAN
- …

Data in input $x, x'$ only
with known mapping $f'$
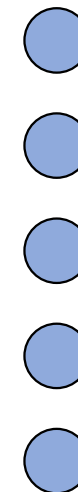(Learn the mapping $f$)

$x'$

$x$

$f(x)$

$x' = f(x)$
**Self-supervised Learning**
- Word2Vec
- Denoising Autoencoder
- …

Data in input only
with known inverse mapping $f'$
(Learn the mapping $f$ and output $y$)

$y = f(x), x = f'(y)$
**Self-augmented Learning**
- ?

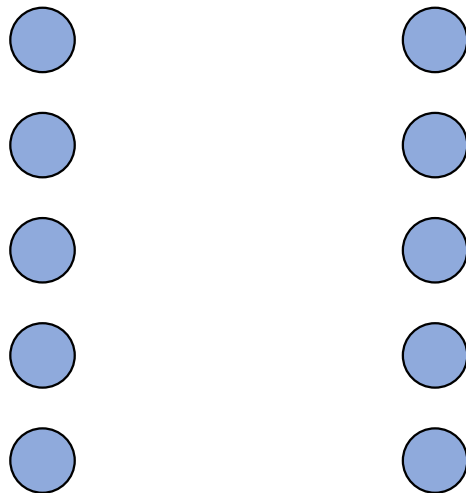# Application of Generative Models: Learning Methods

- Unsupervised Learning
- Semi-supervised Learning
- Weakly-supervised Learning
- Dual Learning
- Self-supervised Learning
- Self-augmented Learning

- **Unsupervised Learning**
- Semi-supervised Learning
- Weakly-supervised Learning
- Dual Learning
- Self-supervised Learning
- Self-augmented Learning

# Unsupervised Learning

Data in both input $x$ and output $y$
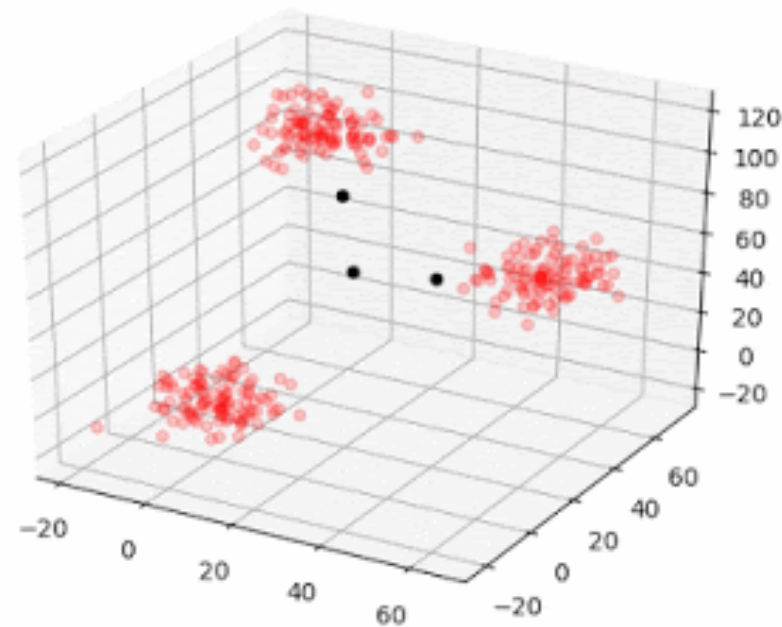(Learn the mapping $f$)

$y = f(x)$

**Unsupervised Learning**

- In practice, it is difficult to obtain a large amount of labeled data, but it is easy to get a large amount of unlabeled data.

- Learn a good feature extractor using unlabeled data and then learn the classifier using labeled data can improve the performance.
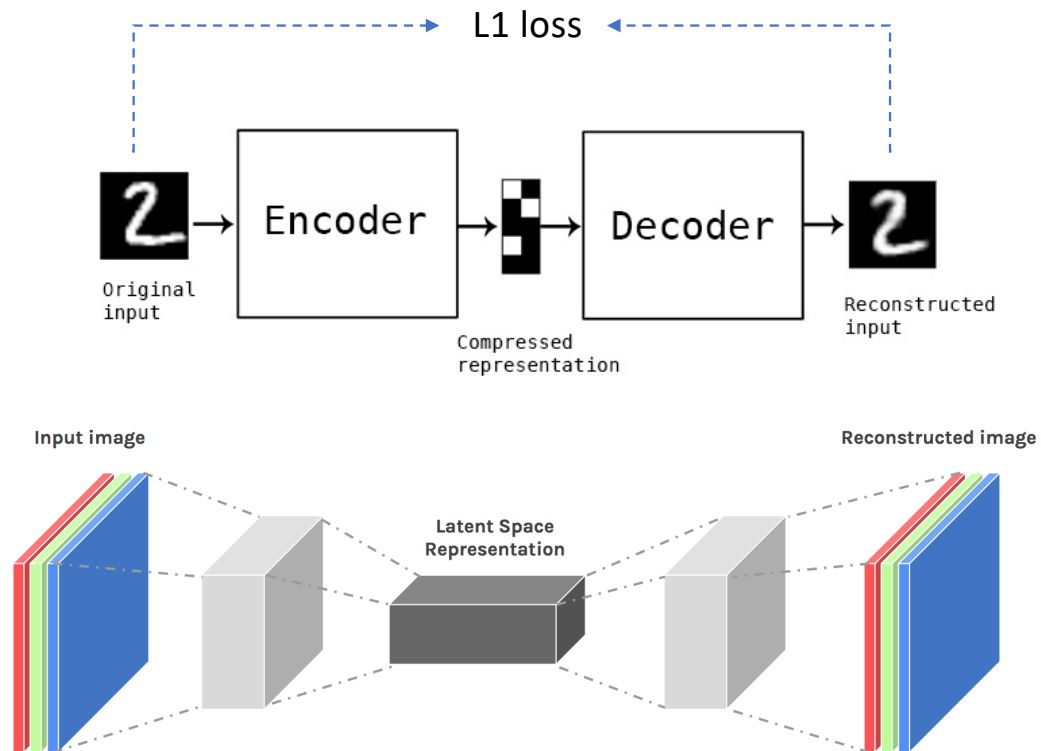
# Unsupervised Learning

- Unsupervised learning is about problems where we don't have labeled answers, such as clustering, dimensionality reduction, and anomaly detection.

- Clustering: EM

- Dimension Reduction: PCA

- ...

# Unsupervised Learning
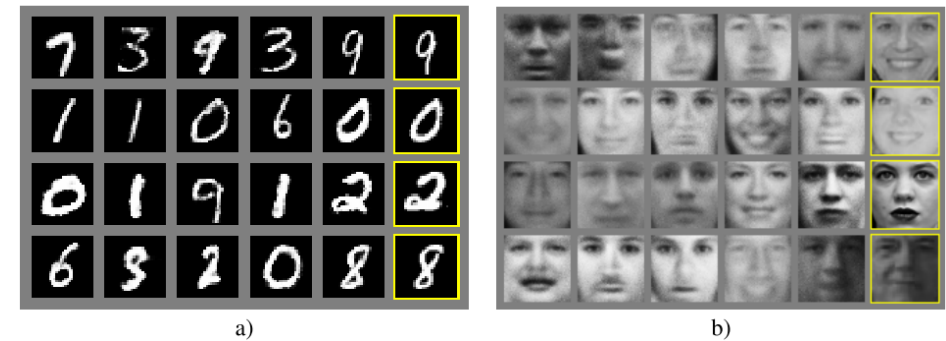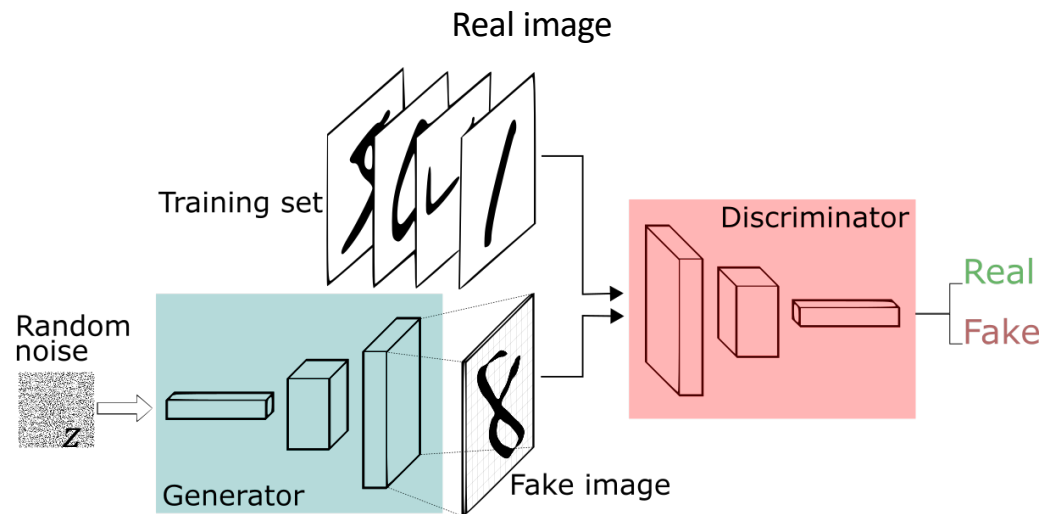
- **Autoencoder**      (when output is features)



Autoencoder: Encode the input image $x$ into a hidden state, then decode the latent space representation into a image $\bar{x}$. Then minimize the reconstruction loss between $x$ and $\bar{x}$.

# Unsupervised Learning

- **GANs**



Update the discriminator – ascending gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D\left(\boldsymbol{x}^{(i)}\right) + \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right)\right].$$

Update the generator – descending gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right).$$

https://pathmind.com/wiki/generative-adversarial-network-gan

# Unsupervised Learning

- **HoloGAN: learn the rotation concept**



HoloGAN: Unsupervised learning of 3D representations from natural images. NIPS 2019

# Unsupervised Learning

- **HoloGAN: How it works**



HoloGAN: Unsupervised learning of 3D representations from natural images. NIPS 2019

- Unsupervised Learning
- **Semi-supervised Learning**
- Weakly-supervised Learning
- Dual Learning
- Self-supervised Learning
- Self-augmented Learning

# Semi-supervised Learning

Data in both input $x$ and output $y$
with known partial mapping
(Learn the mapping $f$)



$y = f(x)$
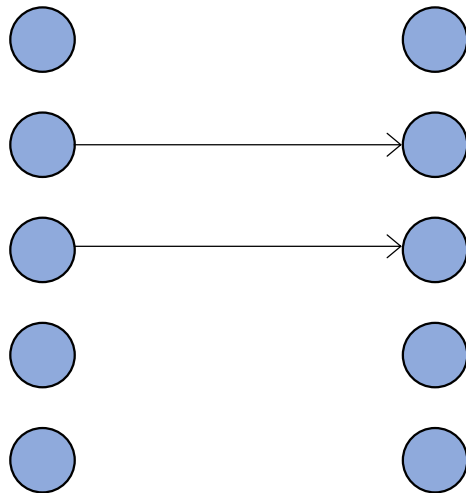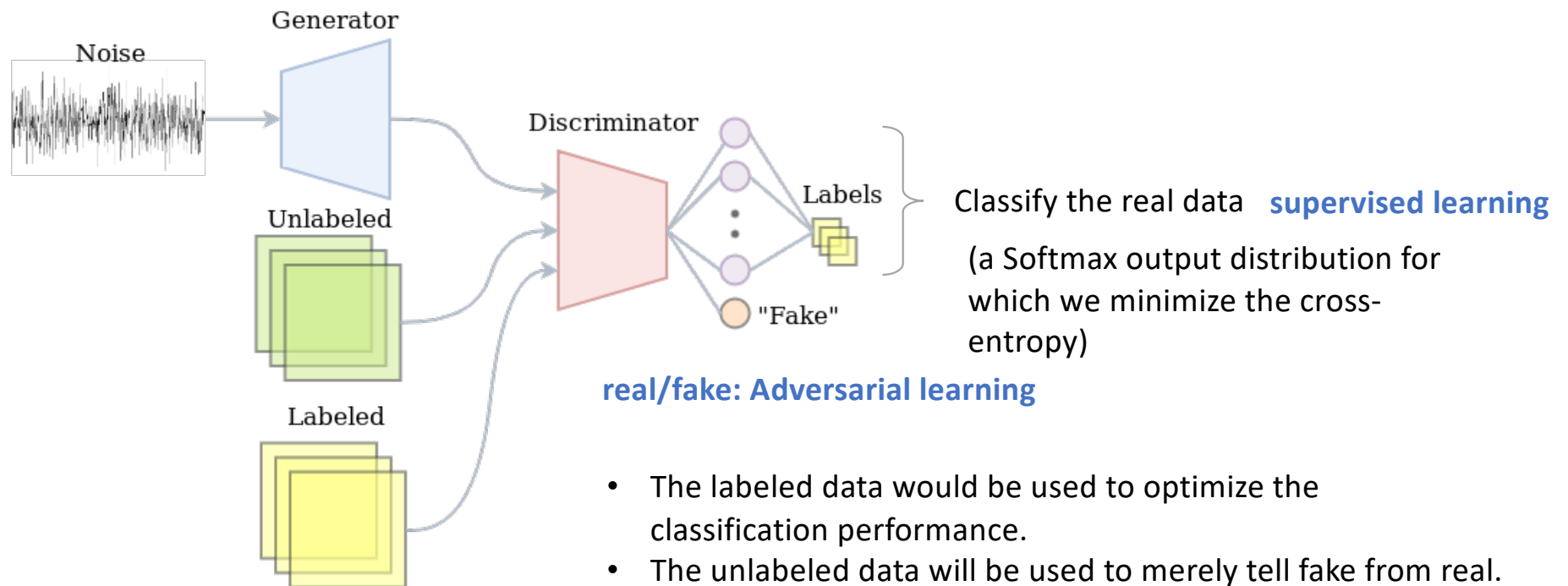**Semi-supervised Learning**

- **Motivation:**
  - Unlabeled data is easy to be obtained
  - Labeled data can be hard to get

- **Goal:**
  - Semi-supervised learning mixes labeled and labeled data to produce better models.

- **vs. Transductive Learning:**
  - Semi-supervised learning is eventually applied to the testing data
  - Transductive learning is only related to the unlabelled data

# Semi-supervised Learning

- **Semi-supervised GAN**



Classify the real data   **supervised learning**

(a Softmax output distribution for which we minimize the cross-entropy)

**real/fake: Adversarial learning**

- The labeled data would be used to optimize the classification performance.
- The unlabeled data will be used to merely tell fake from real.

https://jostosh.github.io/ssl-gan/

# Semi-supervised Learning

- **Semi-supervised GAN**

- Discriminator loss

Discriminator Output

| | cls 1 | cls 2 | ... | cls k |

fake

real

the probability of it being real:  $p(x) = \dfrac{Z(x)}{Z(x) + \exp(l_{fake})} = \dfrac{Z(x)}{1 + Z(x)}$

where $Z(x)$ is the sum of the unnormalized probabilities in softmax operation.
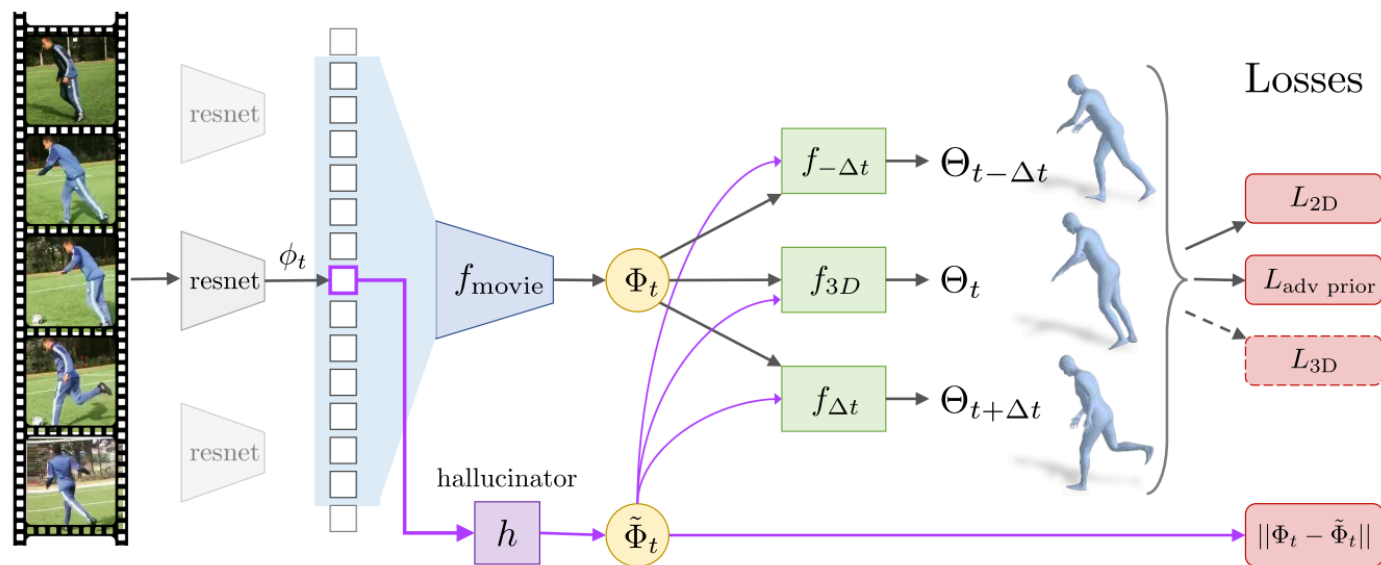$\log(Z(x)) = \text{logsumexp}(l_1, \ldots, l_k)$

Gradient descent:  $-\log(D(x)) - \log(1 - D(G(\mathbf{z})))$

$$= -\log\left(\frac{Z(x)}{1 + Z(x)}\right) - \log\left(1 - \frac{Z(G(\mathbf{z}))}{1 + Z(G(\mathbf{z}))}\right)$$

# Semi-supervised Learning

- **Example: 2D Video to 3D shape**

  The model can learn from videos with only 2D pose annotations in a semi-supervised manner.

$L_{2D}$, $L_{3D}$ : supervision from ground-truth



train a temporal encoder $f_{\text{movie}}$ that learns a representation of 3D human dynamics $\Phi_t$ over the **temporal window centered at frame t**

$L_{adv\,prior}$: each prior discriminator judge a corresponding joint rotation of the body model

$$\sum_k \left( D_k(\boldsymbol{\Theta}) - 1 \right)^2$$

make sure that the **hallucinator** can recover the current 3D mesh as well as its 3D past and future motion.

Learning 3D Human Dynamics from Video. *A. Kanazawa, J. Zhang et al*. CVPR, 2019

18

# Semi-supervised Learning

- **Example: 2D Video to 3D shape**

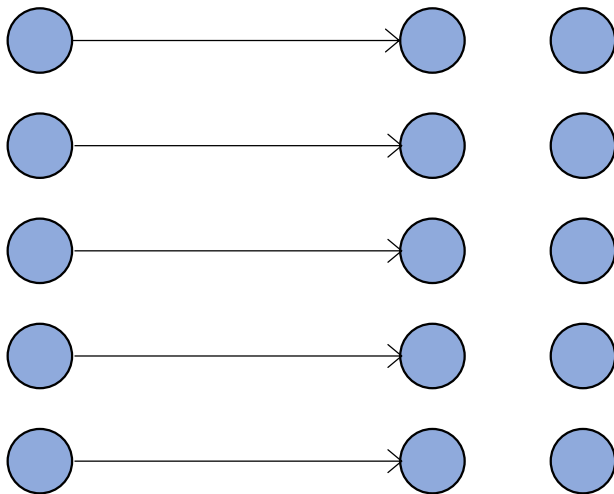From a single image, the model can recover the current 3D mesh as well as its 3D past and future motion.



$$L_t = L_{2D} + L_{3D} + L_{\text{adv prior}} + L_{\beta \text{ prior}}$$

$$L_{\text{const shape}} = \sum_{t=1}^{T-1} ||\beta_t - \beta_{t+1}||. \qquad L_{\text{temporal}} = \sum_t L_t + \sum_{\Delta t} L_{t+\Delta t} + L_{\text{const shape}}.$$

Learning 3D Human Dynamics from Video. *A. Kanazawa, J. Zhang et al*. CVPR, 2019

- Unsupervised Learning
- Semi-supervised Learning
- **Weakly-supervised Learning**
- Dual Learning
- Self-supervised Learning
- Self-augmented Learning

# Weakly-supervised Learning

Data in both input $x$ and output $y$
with known mapping for $y$
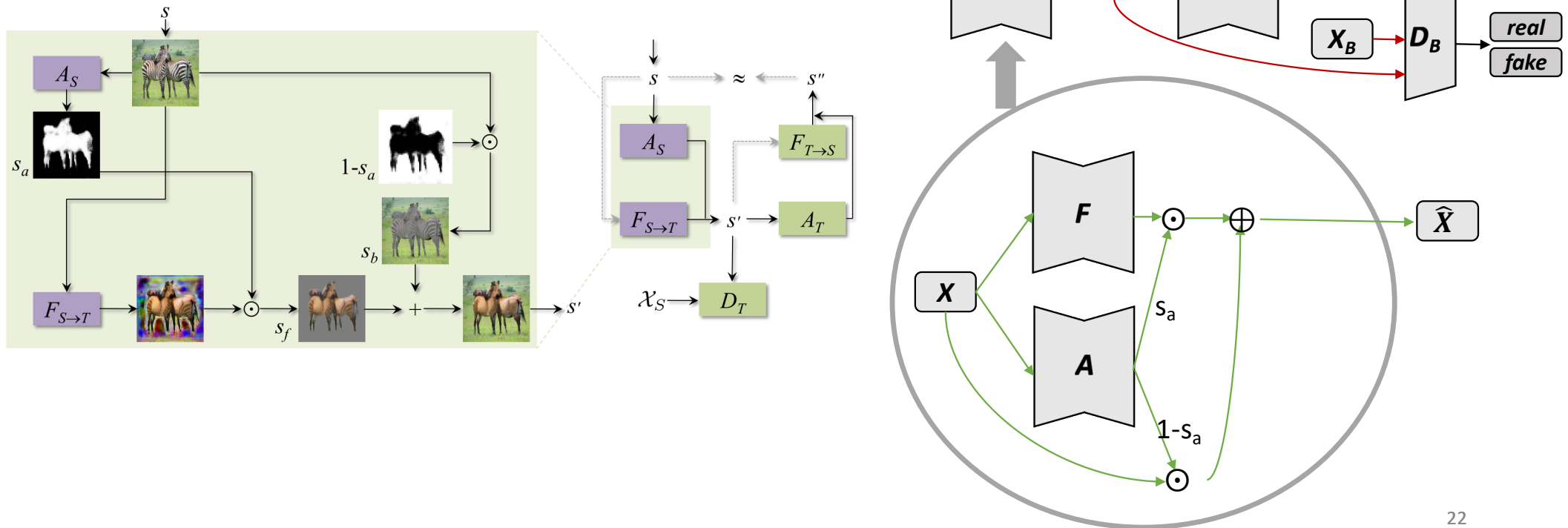(Learn the mapping $f$ for another output $y'$)



$$y' = f(x)$$
**Weakly-supervised Learning**

- Weakly supervised learning is a machine learning framework where the model is trained using examples that are only partially annotated or labeled.
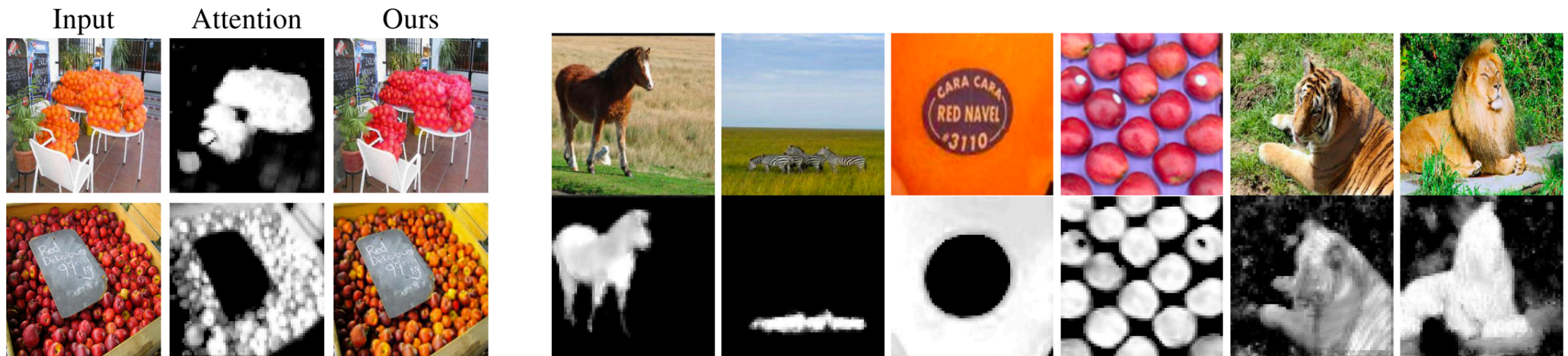
# Weakly-supervised Learning

- **Attention CycleGAN**
- Learn the segmentation via synthesis
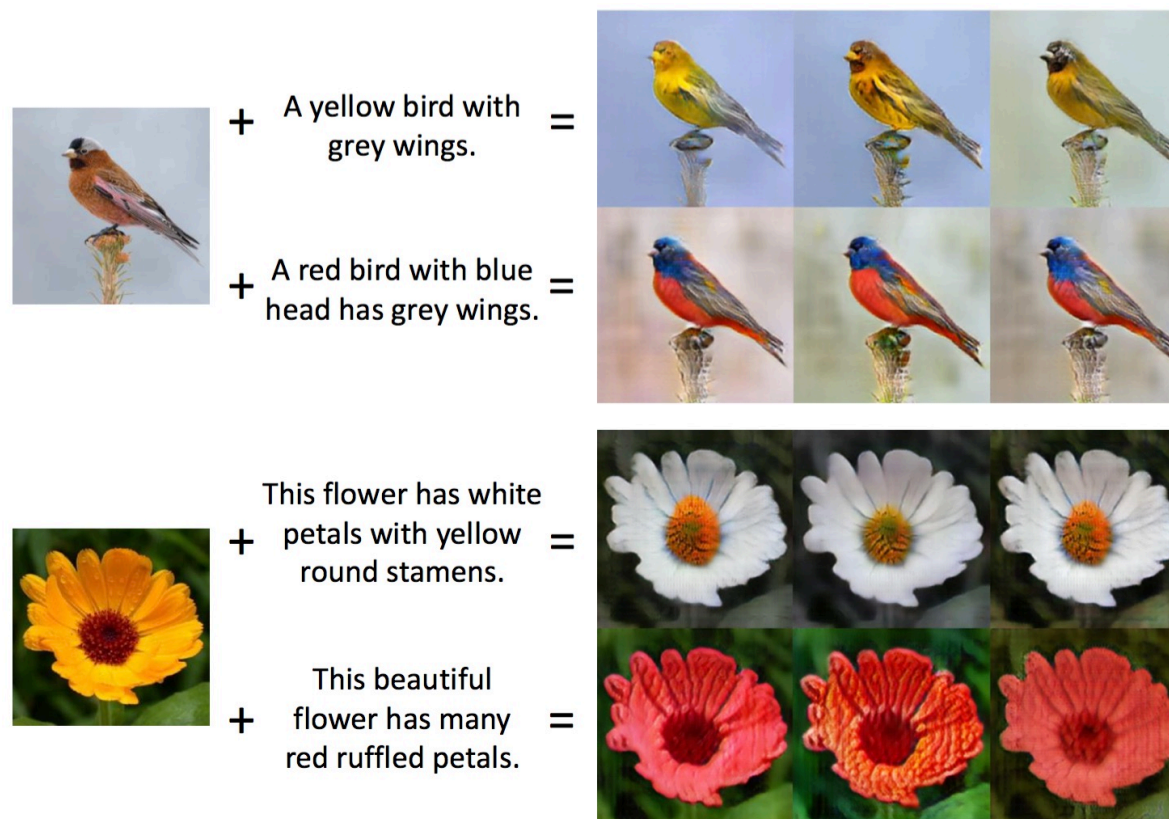
# Weakly-supervised Learning

- **Attention CycleGAN**
  - Learn the segmentation without segmentation masks
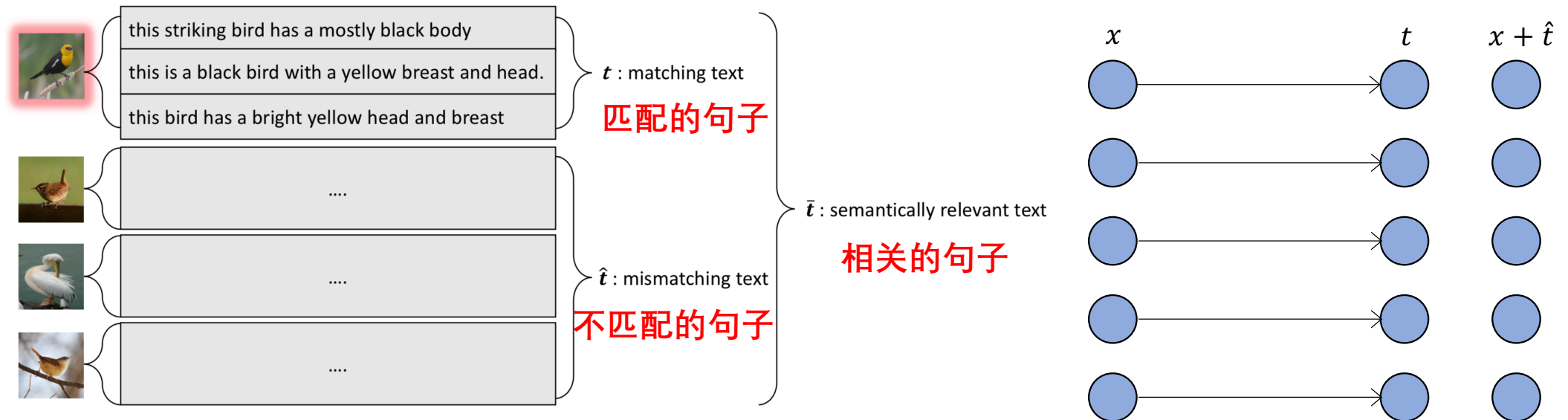
# Weakly-supervised Learning

- **Semantic Image Synthesis: Language Image Manipulation**



Semantic Image Synthesis via Adversarial Learning. *H. Dong, S. Yu et al. ICCV 2017.*
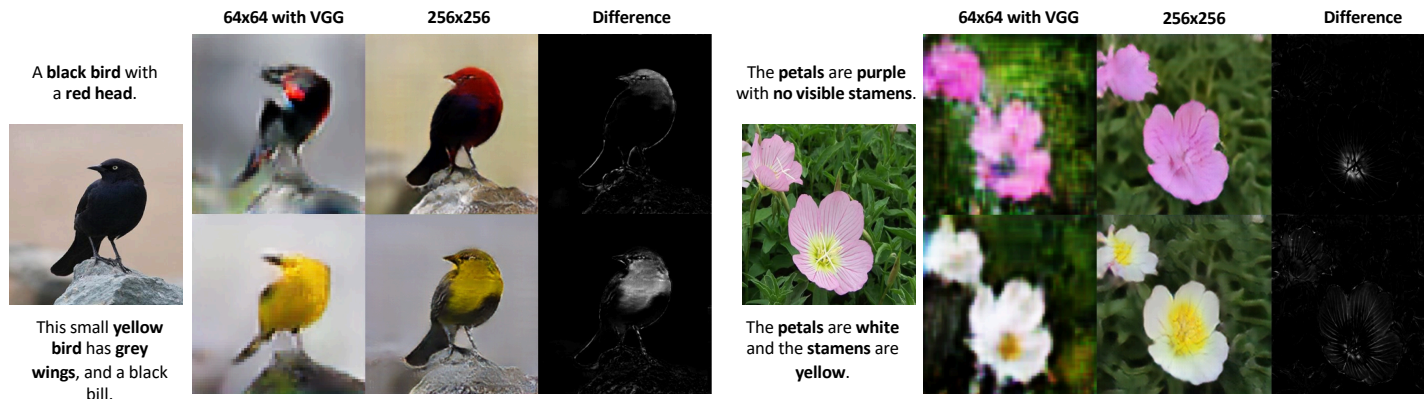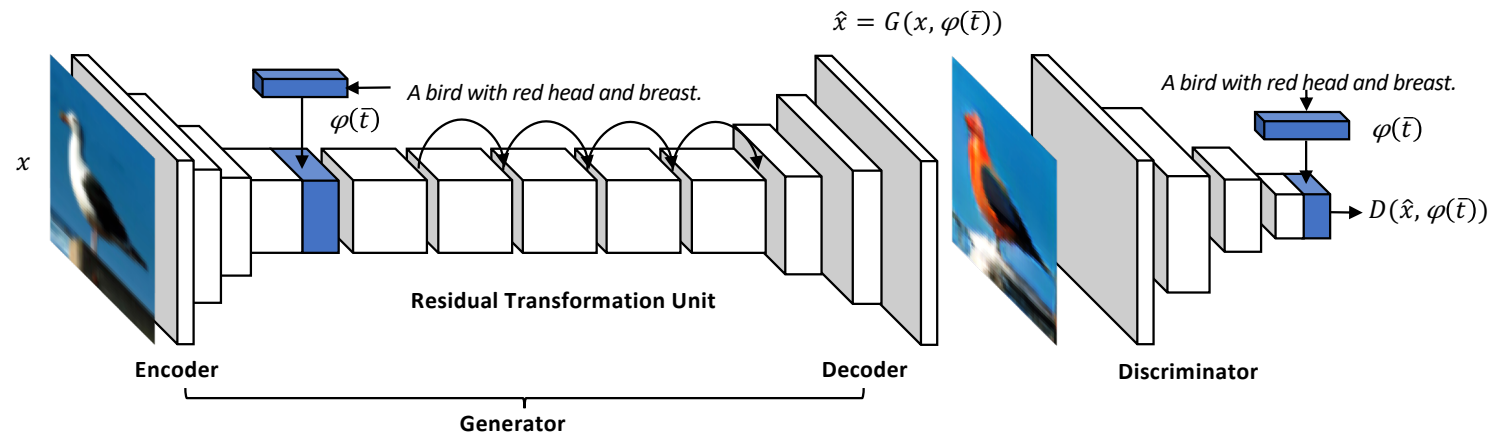
# Weakly-supervised Learning

- **Semantic Image Synthesis: Language Image Manipulation**



this striking bird has a mostly black body

this is a black bird with a yellow breast and head.

this bird has a bright yellow head and breast

$t$ : matching text

匹配的句子

....

....

....

$\hat{t}$ : mismatching text

不匹配的句子

$\bar{t}$ : semantically relevant text

相关的句子

$x$            $t$            $x + \hat{t}$

Semantic Image Synthesis via Adversarial Learning. *H. Dong, S. Yu et al. ICCV 2017.*

# Weakly-supervised Learning

- **Semantic Image Synthesis: Learn the segmentation via synthesis**



$$\hat{x} = G(x, \varphi(\bar{t}))$$

Semantic Image Synthesis via Adversarial Learning. *Dong, H., Yu, S., Wu, C., Guo, Y.* 2017. ICCV

- Unsupervised Learning
- Semi-supervised Learning
- Weakly-supervised Learning
- **Dual Learning**
- Self-supervised Learning
- Self-augmented Learning

# Dual Learning

Data in both input and output
(Learn the mapping $f, f'$)

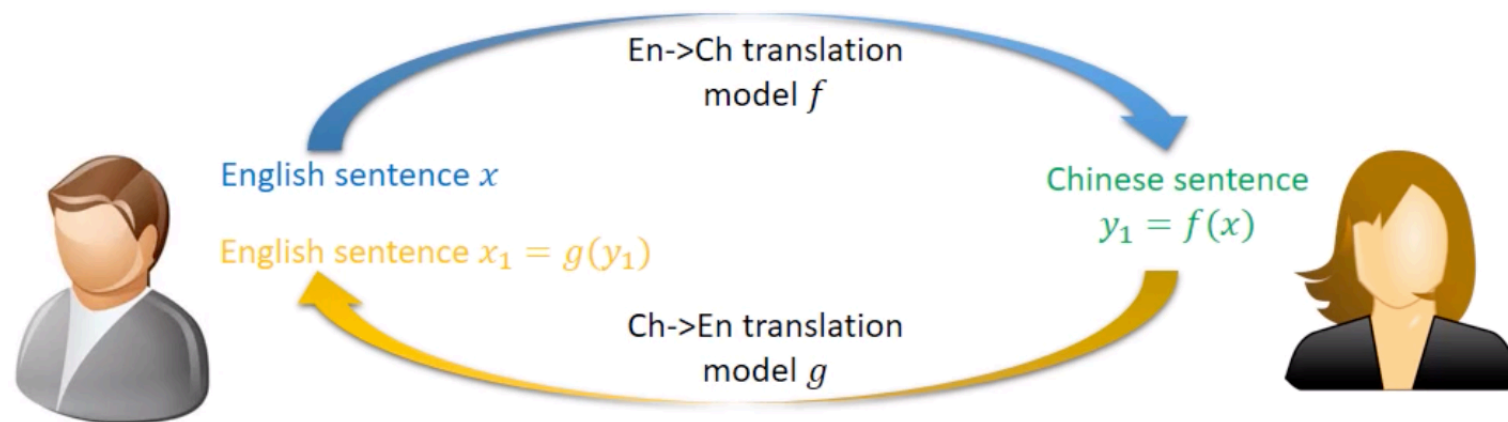$$x \qquad\qquad y$$

$$y = f(x), x = f'(y)$$
**(Unsupervised) Dual Learning**

- Motivation
  - Human label is expensive
  - No feedback if using unlabeled data

| Application | Primal Task | Dual (Inverse) Task |
|---|---|---|
| Machine translation | Translate language from A to B | Translate language from B to A |
| Speed processing | Speech to text (STT) | Text to speech (TTS) |
| Image understanding | Image captioning | Image generation |
| Conversation engine | Question | Answer |
| Search engine | Search | Query |

# Dual Learning

- **Language Translation**



En->Ch translation model $f$

English sentence $x$

Chinese sentence $y_1 = f(x)$

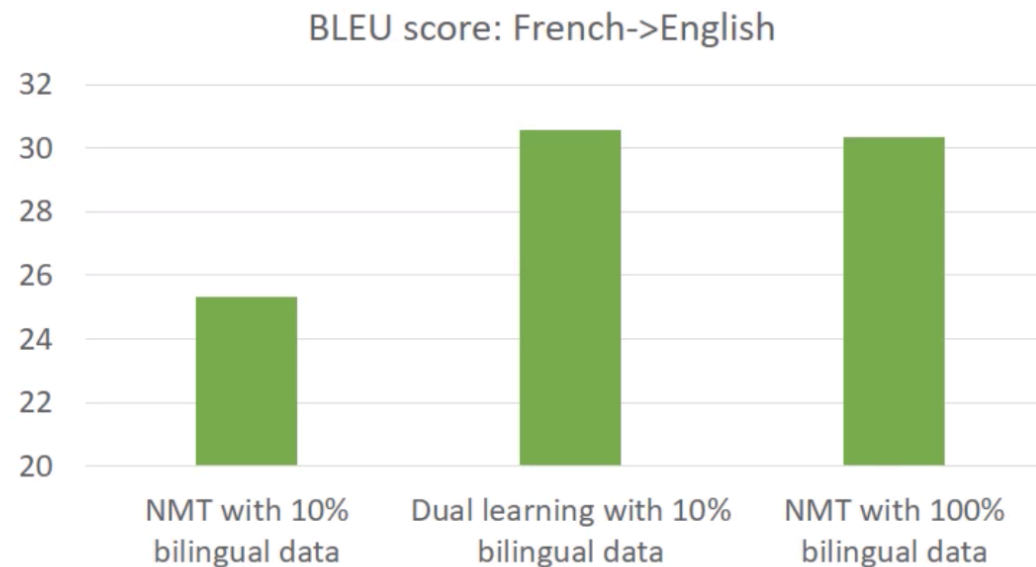English sentence $x_1 = g(y_1)$

Ch->En translation model $g$

Feedback signals during the loop:
- $s(x, x_1)$: BLEU score of $x_1$ given $x$
- $L(y)$ and $L(x_1)$: Likelihood and language model of $y_1$ and $x_1$

Reinforcement learning is used to improve the translation models from these feedback signals
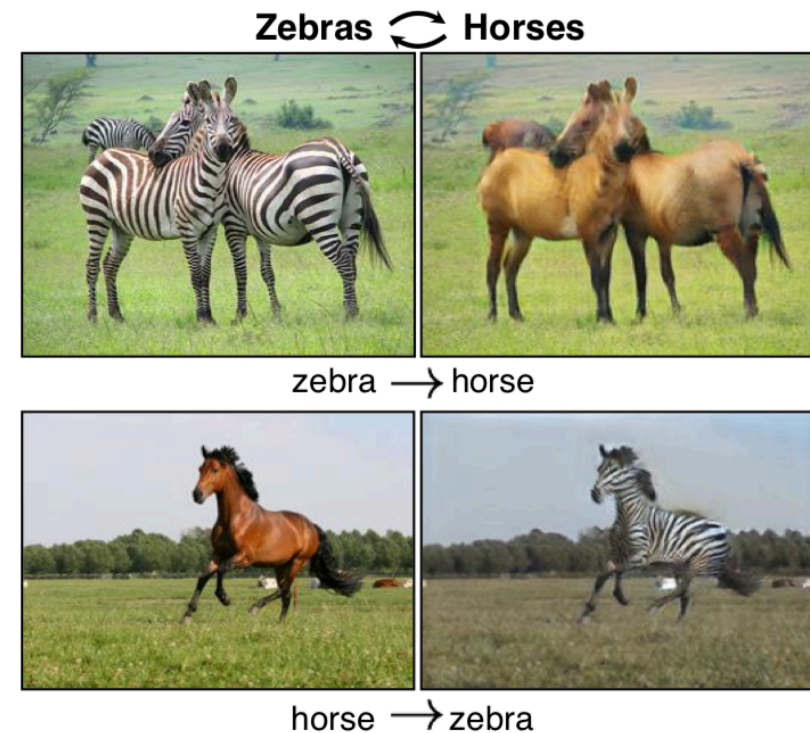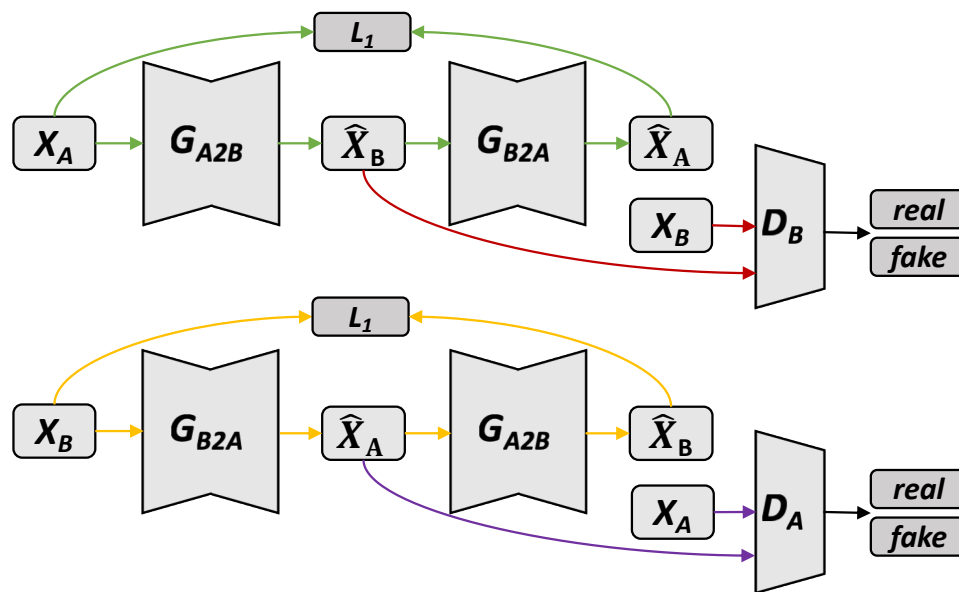
Dual Learning for Machine Translation.
*Yingce Xia, Di He, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, Wei-Ying Ma. NIPS, 2016*

# Dual Learning

- **Language Translation**

BLEU score: French->English



Starting from initial models obtained from only 10% bilingual data, dual learning can achieve similar accuracy as the NMT model learned from 100% bilingual data!
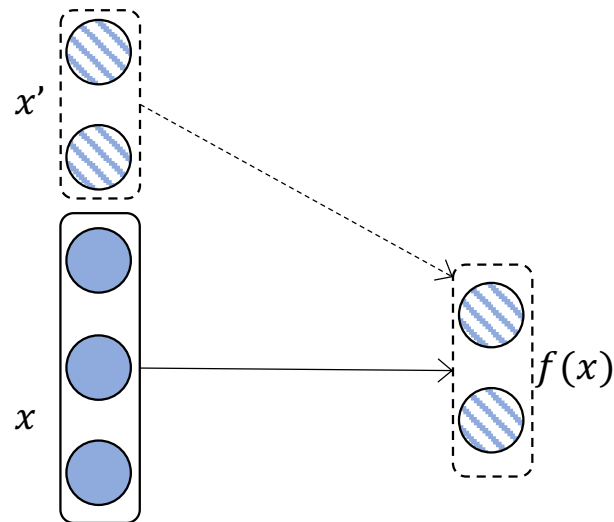
Dual Learning for Machine Translation.
Yingce Xia, Di He, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, Wei-Ying Ma. NIPS, 2016

# Dual Learning

- **Unpaired Image-to-Image Translation**



Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. *J. Zhu, T. Park et al. ICCV 2017.*

- Unsupervised Learning
- Semi-supervised Learning
- Weakly-supervised Learning
- Dual Learning
- **Self-supervised Learning**
- Self-augmented Learning

# Self-supervised Learning

Data in input $x, x'$ only
with known mapping $f'$
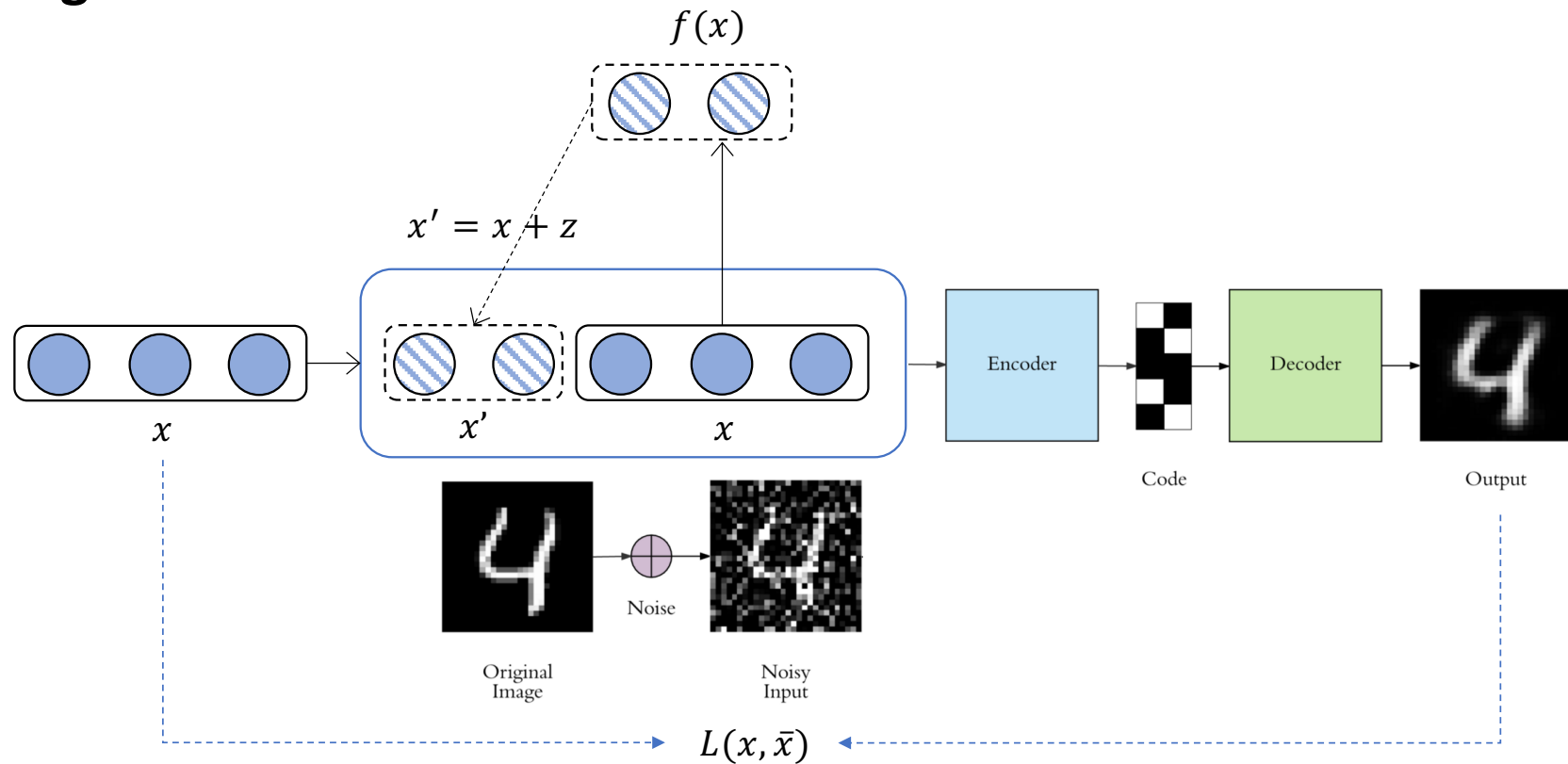(Learn the mapping $f$)

$x'$

$x$

$f(x)$

$$x' = f(x)$$

**Self-supervised Learning**

- Self-supervised learning is autonomous supervised learning, it learns to predict part of its input from other parts of its input.

- Examples: Word2Vec, Denoising Autoencoder

- Self-supervised vs. unsupervised learning: Self-supervised learning is like unsupervised Learning because the system learns without using explicitly-provided labels. It is different from unsupervised learning because we are not learning the inherent structure of data. Self-supervised learning, unlike unsupervised learning, is not centered around clustering and grouping, dimensionality reduction, recommendation engines, density estimation, or anomaly detection.
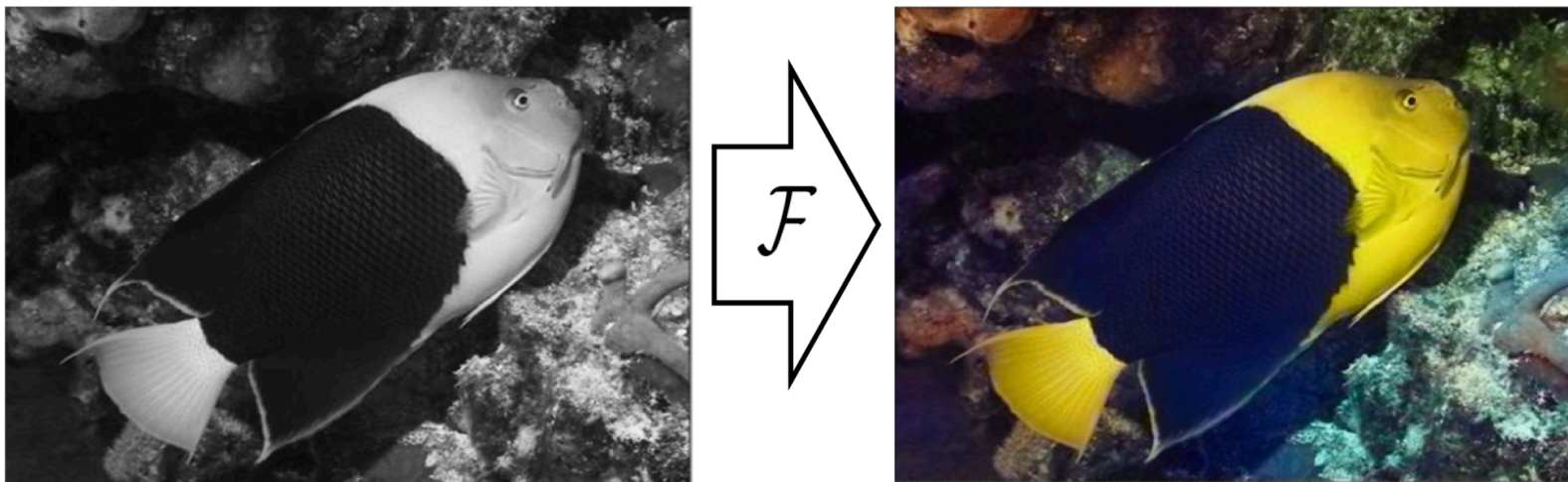
# Self-supervised Learning

- **Denoising Autoencoder**



$$f(x)$$

$$x' = x + z$$

$$x$$

$$x'$$ $$x$$

Encoder → Code → Decoder → Output

Original Image → Noise → Noisy Input

$$L(x, \bar{x})$$

Extracting and composing robust features with denoising autoencoders, Pascal Vincent etc, 2008

# Self-supervised Learning

- **Image Example: Colorization**
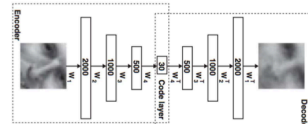


Colorful Image Colorization. *Zhang et al.,* ECCV 2016
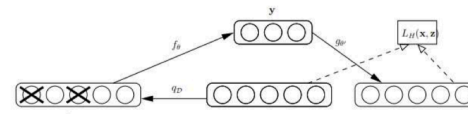
# Self-supervised Learning

- Image Examples



**Autoencoders**

Hinton & Salakhutdinov.
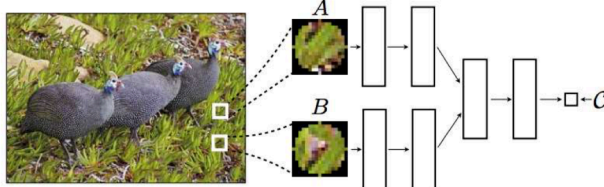Science 2006.

**Denoising Autoencoders**

Vincent *et al.* ICML 2008.

**Exemplar networks**

Dosovitskiy *et al.*, NIPS 2014
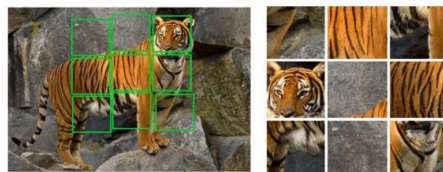
**Co-Occurrence**

Isola *et al.* ICLR Workshop 2016.

**Egomotion**

Agrawal *et al.* ICCV 2015   Jayaraman *et al.* ICCV 2015
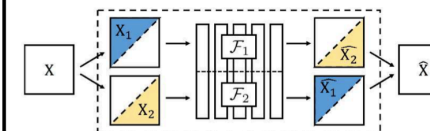
**Context**

Noroozi et al 2016          Pathak *et al.* CVPR 2016

**Split-brain auto-encoders**

Zhang *et al.* CVPR 2017

# Self-supervised Learning

- **Video Example**



- Videos contain
  - Color, Temporal info
- Possible proxy tasks
  - Temporal order of the frames
  - Optical flow: Motion of objects
  - …

# Self-supervised Learning
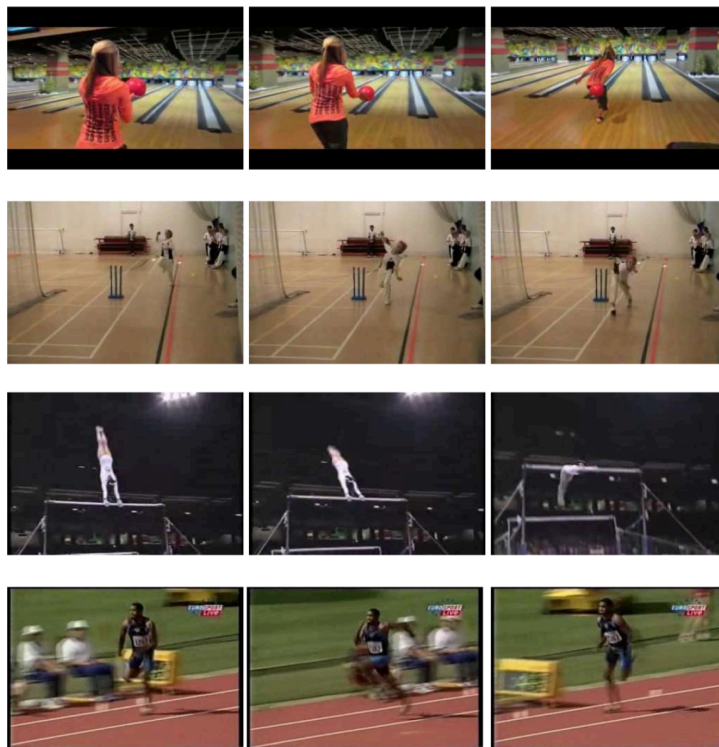
- Video Example: Shuffle and Learn

Given a start and an end, can this point lie in between?

Unsupervised Learning using Temporal Order Verification. *Ishan Misra, C. Lawrence Zitnick and Martial Hebert.* ECCV, 2016.

# Self-supervised Learning

- Video Example: Shuffle and Learn



Unsupervised Learning using Temporal Order Verification. *Ishan Misra, C. Lawrence Zitnick and Martial Hebert.* ECCV, 2016.

# Self-supervised Learning

- Video Example: Shuffle and Learn



Slide credit: Ishan Misra

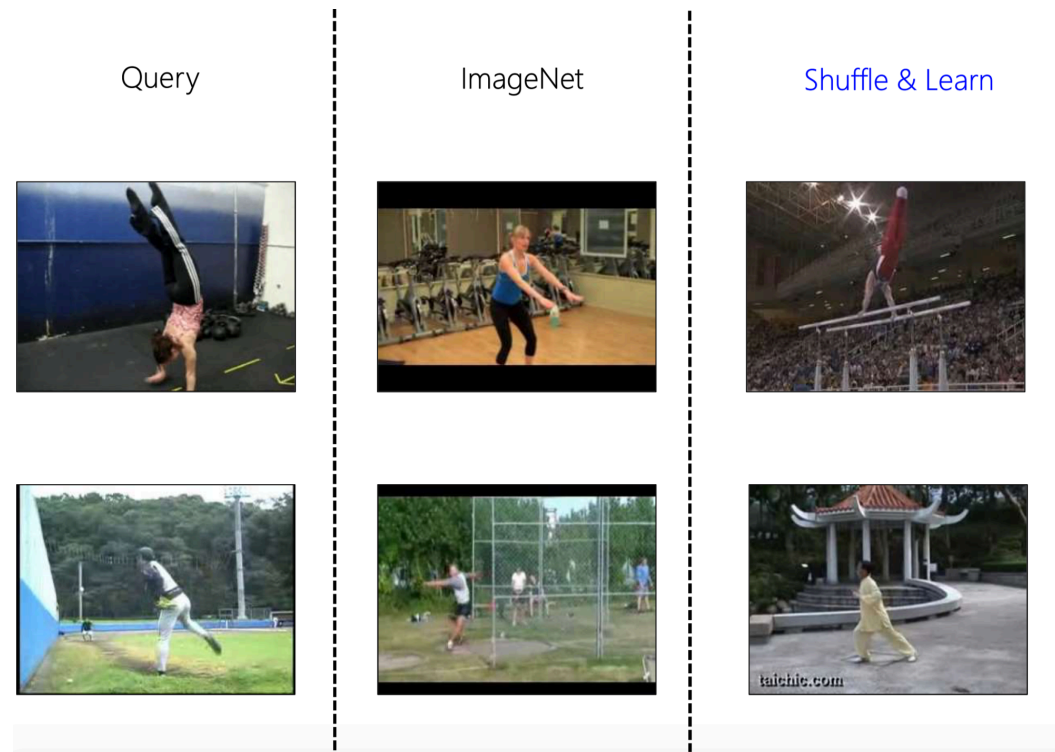Unsupervised Learning using Temporal Order Verification. *Ishan Misra, C. Lawrence Zitnick and Martial Hebert.* ECCV, 2016.
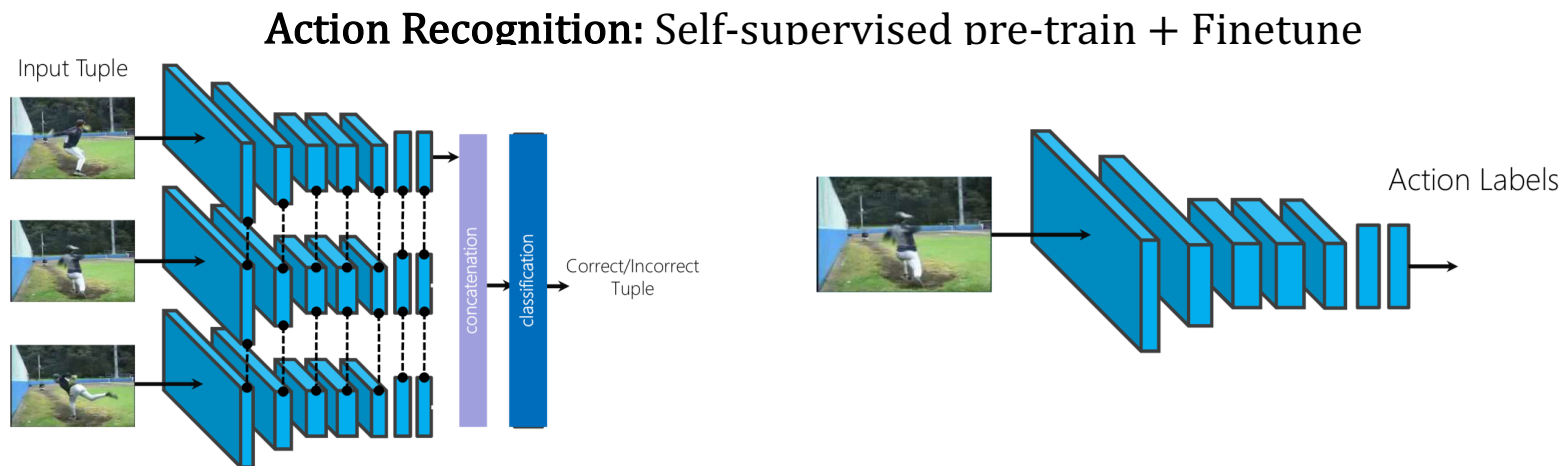
# Self-supervised Learning

- Video Example: Shuffle and Learn

**Image Retrieval:** Nearest Neighbors of Query Frame (FC5 outputs)
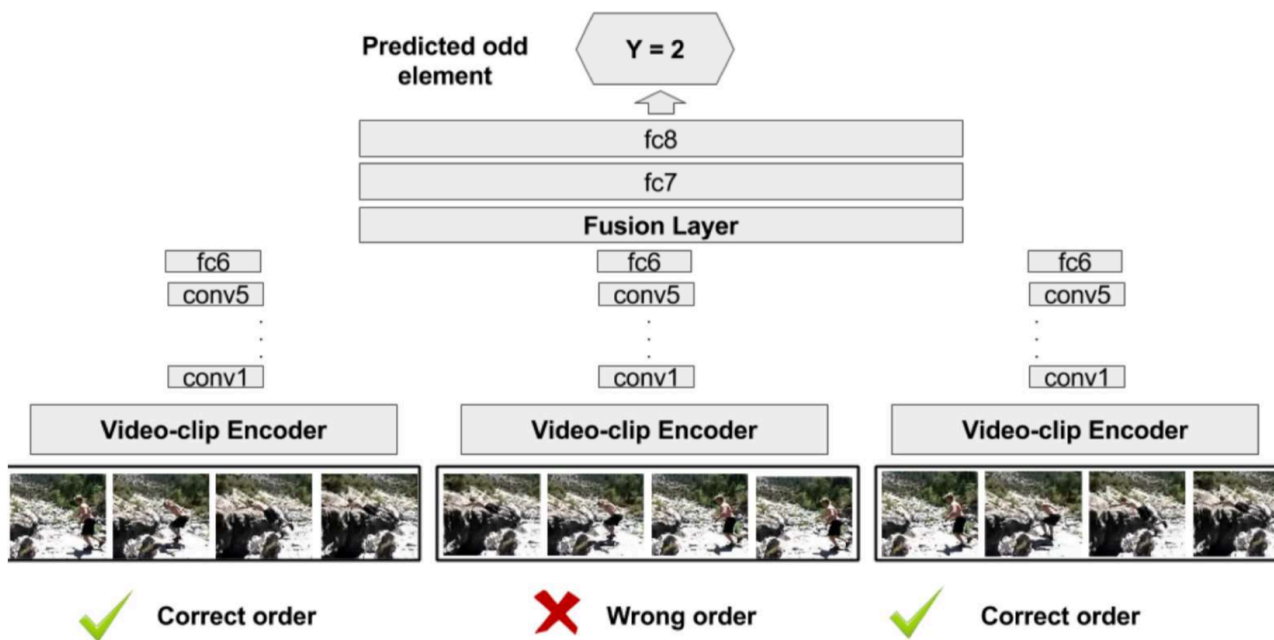
Unsupervised Learning using Temporal Order Verification. *Ishan Misra, C. Lawrence Zitnick and Martial Hebert.* ECCV, 2016.

# Self-supervised Learning

- Video Example: Shuffle and Learn

**Action Recognition:** Self-supervised pre-train + Finetune



| Dataset | Initialization | Mean Classification Accuracy |
|---------|----------------|------------------------------|
| UCF101 | Random | 38.6 |
| | Shuffle & Learn | 50.2 |
| | ImageNet pre-trained | **67.1** |

Unsupervised Learning using Temporal Order Verification. *Ishan Misra, C. Lawrence Zitnick and Martial Hebert.* ECCV, 2016.

# Self-supervised Learning

- Video Example: Odd-One-Out



| Initialization | Mean Classification Accuracy |
|---|---|
| Random | 38.6 |
| Shuffle and Learn | 50.2 |
| Odd-One-Out | 60.3 |
| ImageNet pre-trained | **67.1** |

Self-Supervised Video Representation Learning With Odd-One-Out Networks. *Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould*, ICCV 2017

# Self-supervised Learning

- Video Example: Learning the Arrow of Time

**Forward or backward plays?**



input motion

- Depending on the video, solving the task may require
(a) low-level understanding (e.g. physics)
(b) high-level reasoning (e.g. semantics)
(c) familiarity with very subtle effects
(d) camera conventions

- Input: optical flow in two chunks
- Final layer: global average pooling to allow class activation map (CAM)

Learning and Using the Arrow of Time. *Donglai Wei, Joseph Lim, Bill Freeman, Andrew Zisserman.* CVPR 2018
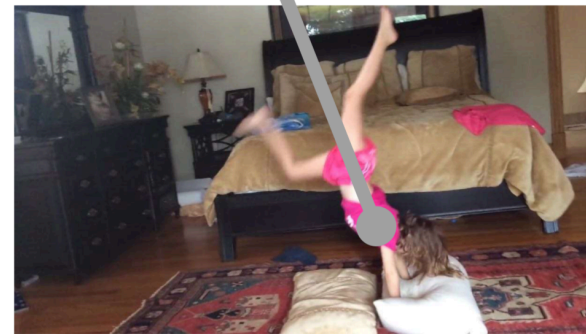
# Self-supervised Learning

- Video Example: Temporal Coherence of Color

Colorize all frames of a grey scale version using a reference frame
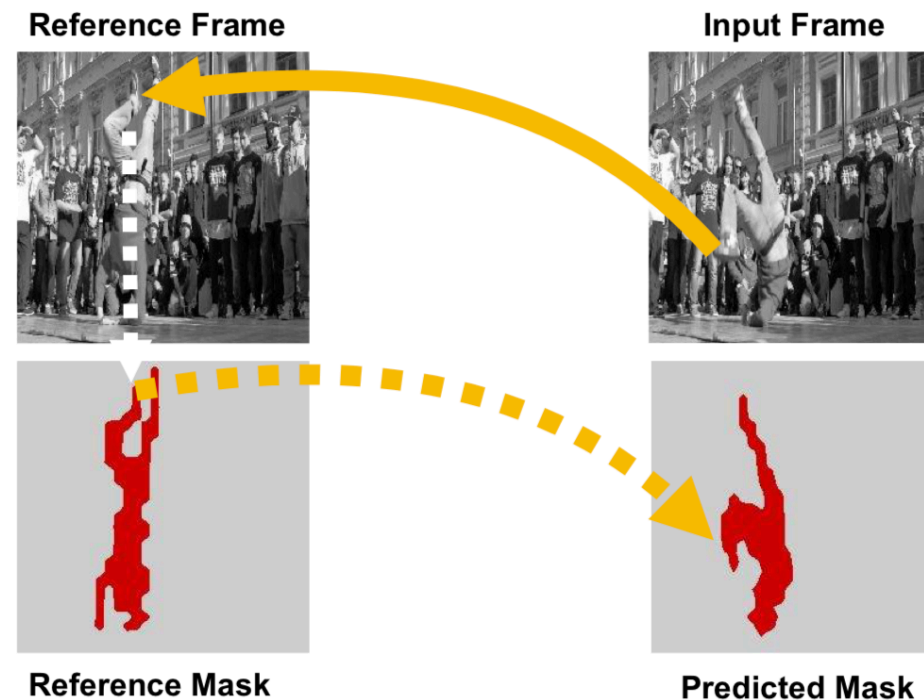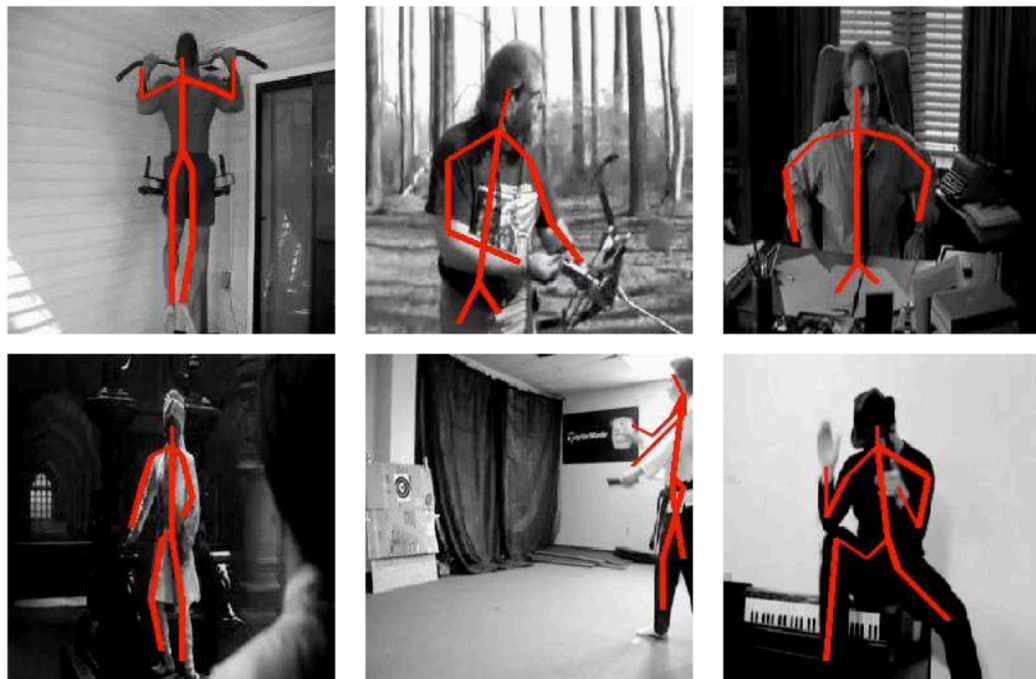


**Reference Frame**

**What color is that?**

Tracking Emerges by Colorizing Videos
*Vondrick, Shrivastava, Fathi, Guadarrama, Murphy,* ECCV 2018

# Self-supervised Learning

- Video Example: Temporal Coherence of Color

**Tracking Emerges:** Only the first frame is given, colors indicate different instances



Tracking Emerges by Colorizing Videos
*Vondrick, Shrivastava, Fathi, Guadarrama, Murphy*, ECCV 2018

# Self-supervised Learning

- Video Example: Temporal Coherence of Color

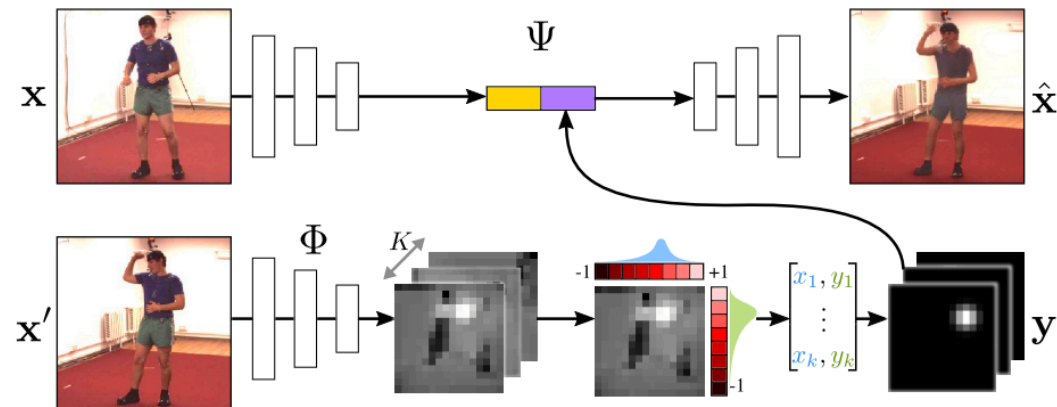**Segment Tracking:** Only the first frame is given, colors indicate different instances



Tracking Emerges by Colorizing Videos
*Vondrick, Shrivastava, Fathi, Guadarrama, Murphy*, ECCV 2018

# Self-supervised Learning

- Video Example: Temporal Coherence of Color

**Pose Tracking:** Only the skeleton in the first frame is given



Tracking Emerges by Colorizing Videos
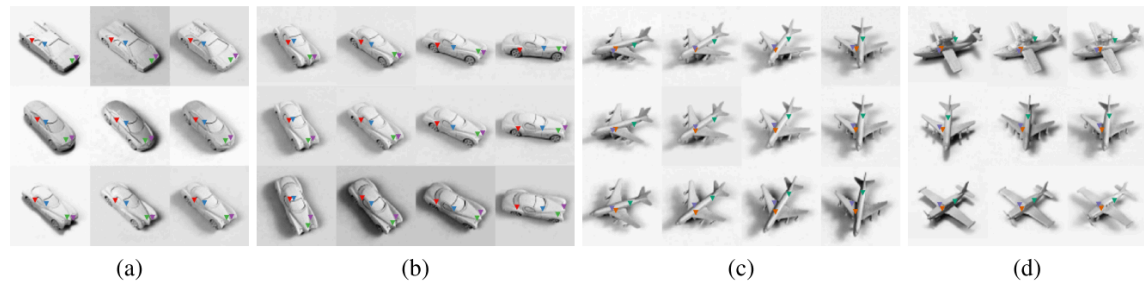*Vondrick, Shrivastava, Fathi, Guadarrama, Murphy*, ECCV 2018

# Self-supervised Learning

- Video Example: Temporal Coherence of Color

**Unsupervised Key-point Detection:** Only paired images of the same object is given



- Achieve retargeting
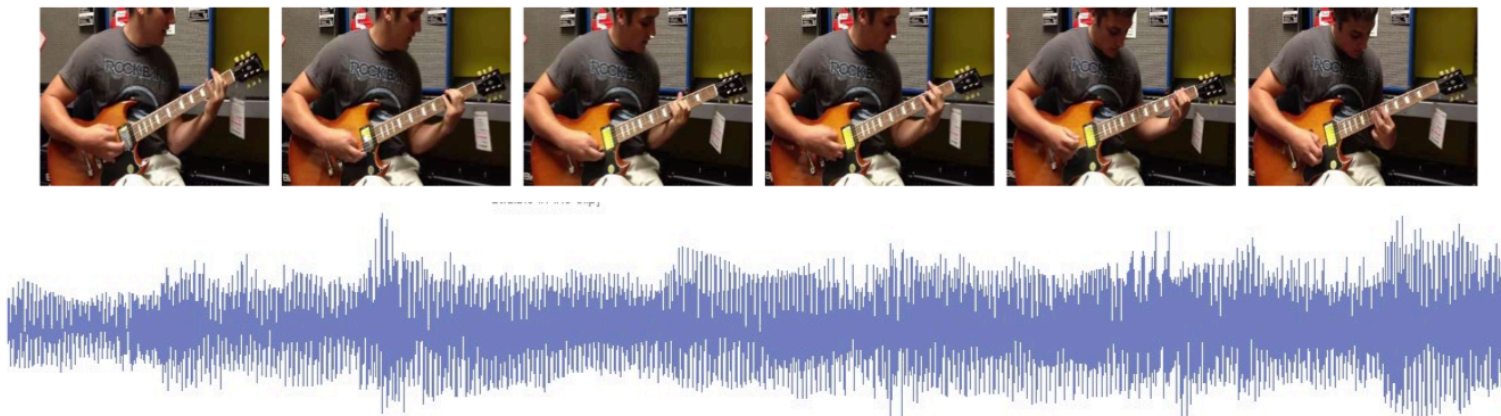- Disentangling Style and Geometry
- Invariant Localization

Unsupervised Learning of Object Landmarks through Conditional Image Generation
*Tomas Jakab, Ankush Gupta et al. NIPS, 2018.*

# Self-supervised Learning

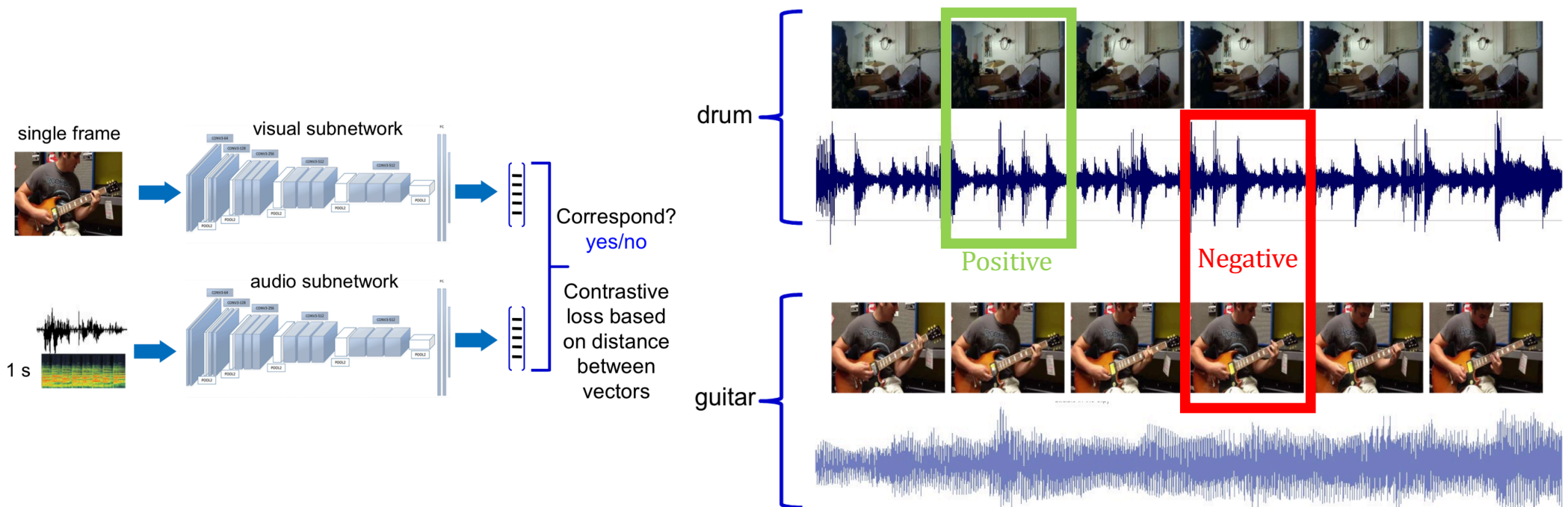- **Video + Sound Example**



- Sound and frames are:
    - Semantically consistent
    - Synchronized
- Two types of proxy task:
    - Predict audio-visual correspondence
    - Predict audio-visual synchronization
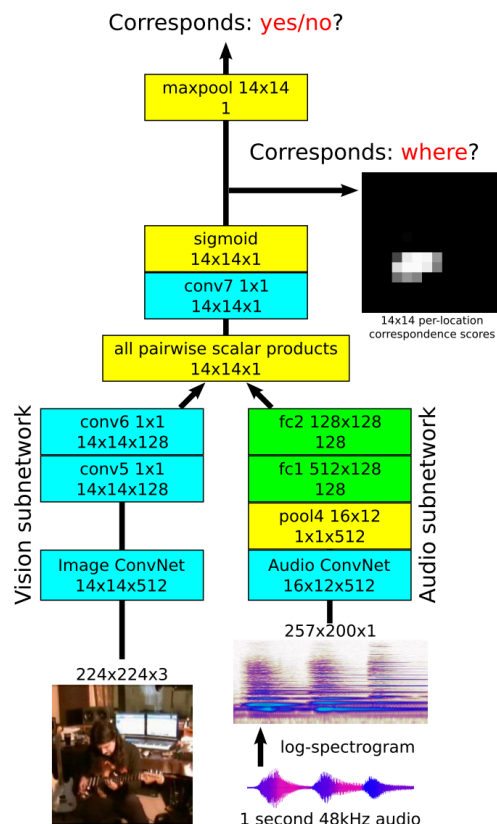
# Self-supervised Learning

- Video + Sound Example: Audio-Visual Co-supervision

Train a network to predict if image and audio clip correspond



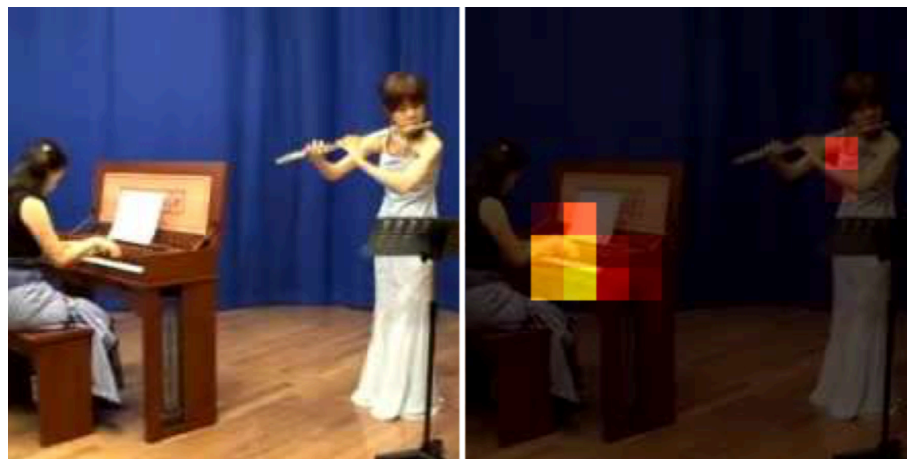Objects that Sound. *Arandjelović and Zisserman,* ICCV 2017 & ECCV 2018

# Self-supervised Learning

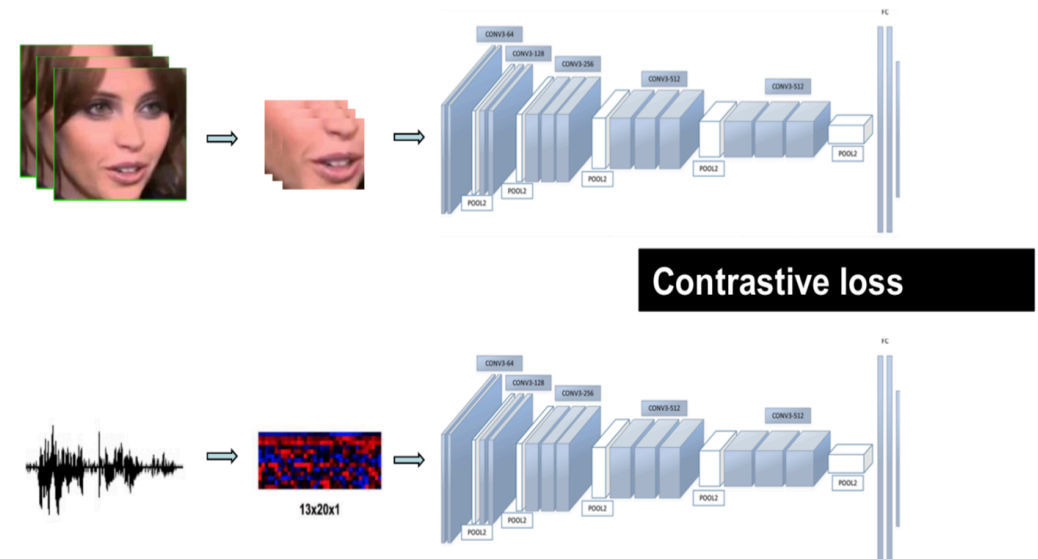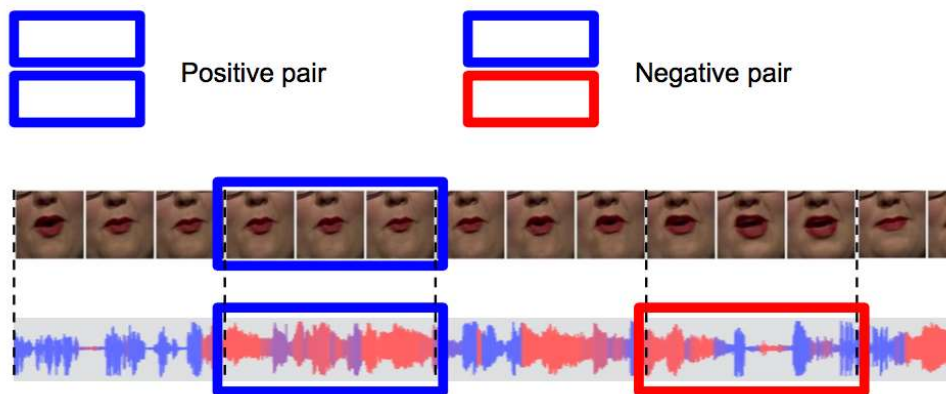- Video + Sound Example: Audio-Visual Co-supervision



- Learn good visual features
- Learn good audio features
- Learn aligned audio-visual embeddings
- Learn to localize objects that sound

- Using learned features
  - Sound classification
  - Query on image to retrieve audio
  - Localizing objects with sound

Objects that Sound. *Arandjelović and Zisserman* *(DeepMind, Ox)*, ICCV 2017 & ECCV 2018

# Self-supervised Learning

- Video + Sound Example: Audio-Visual Co-supervision



**Contrastive loss**

- Applications
  - Active speaker detection
  - Audio-to-video synchronization
  - Voice-over rejection
  - Visual features for lip reading

Out of time: Automatic lip sync in the wild. *Chung, Zisserman,* 2016

- Unsupervised Learning
- Semi-supervised Learning
- Weakly-supervised Learning
- Dual Learning
- Self-supervised Learning
- **Self-augmented Learning**

# Self-augmented Learning

Data in input only
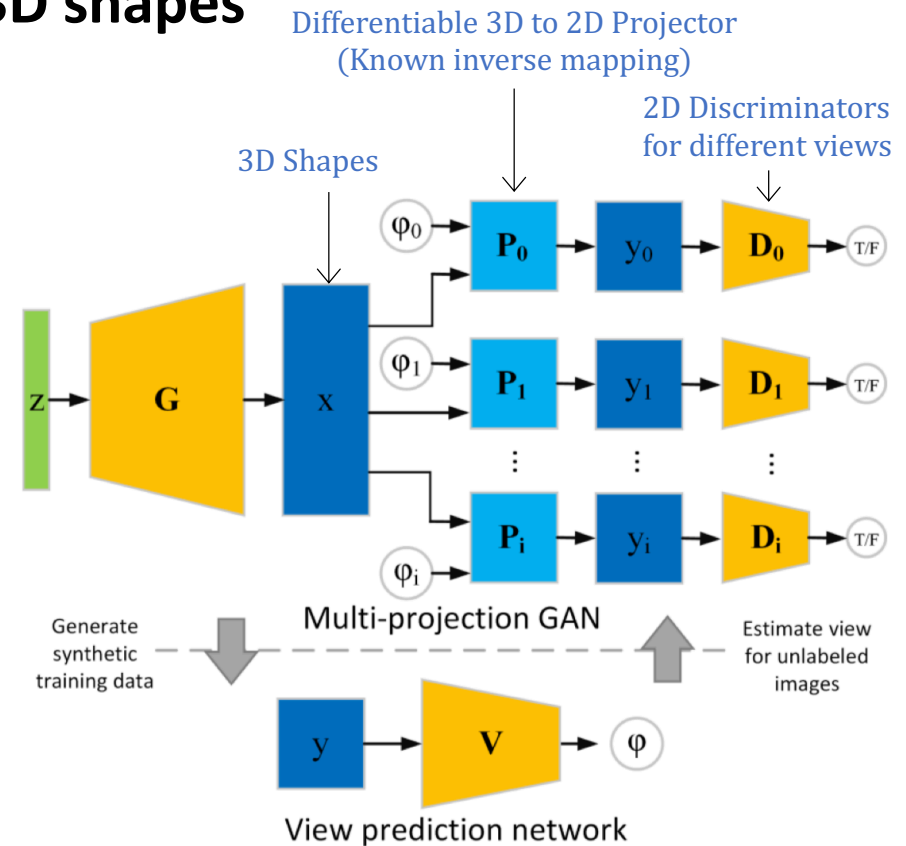with known inverse mapping $f'$
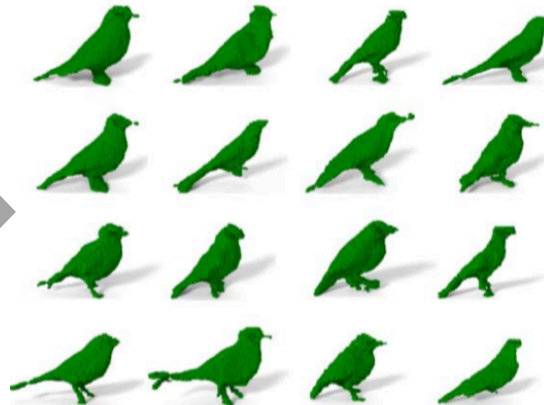(Learn the mapping $f$ and output $y$)

$y = f(x), x = f'(y)$
**Self-augmented Learning**

# Self-augmented Learning

- **Example: Unsupervised 2D images to 3D shapes**



Differentiable 3D to 2D Projector
(Known inverse mapping)

2D Discriminators for different views

3D Shapes

Synthesizing 3D Shapes from Unannotated Image Collections using Multi-projection Generative Adversarial Networks.
*Xiao Li, Yue Dong, Pieter Peers, Xin Tong*. CVPR, 2019

# Summary

- Unsupervised Learning
- Semi-supervised Learning
- Weakly-supervised Learning
- Dual Learning
- Self-supervised Learning
- Self-augmented Learning

# Thanks