payload

mg

y

x

z

# Review: LeNet-5

Moped

Bicycle

Motorbike

Go-cart

Trail

Car, auto

Helicopter

IMAGENET

Animal  Artifact

Bird  Reptile  Fish

Aquatic bird

Bird

Reptile

Fish

Aquatic bird

Red

Copper rockfish

Yellowhammer  Triceratops  Soldierfish

Shelduck  Ruddy turnstone

Structure  Covering

Building

Structure

Building

Sod

Covering

Cardigan

Rope bridge  Dishrag

Viaduct  Ziggurat

Accuracy (%)
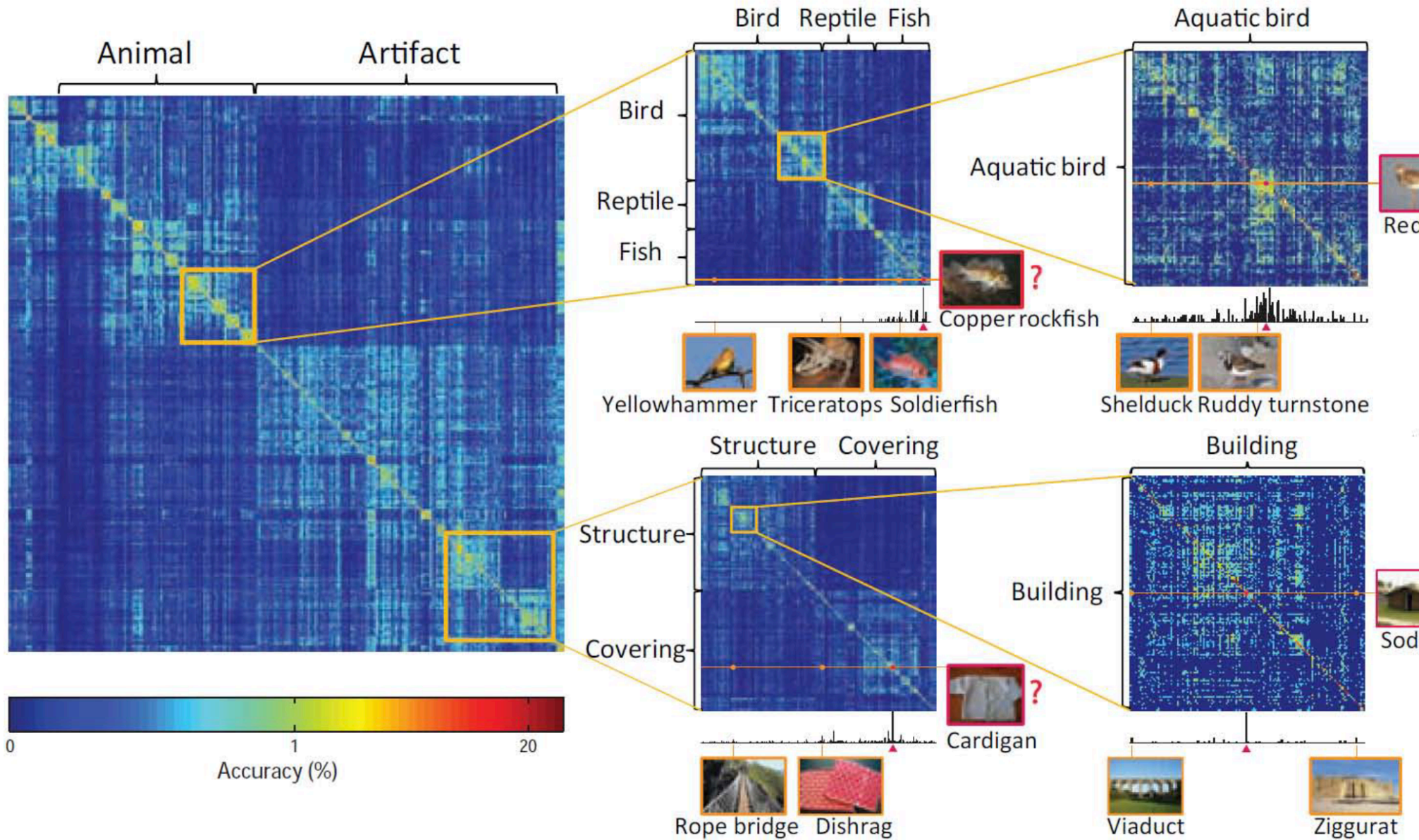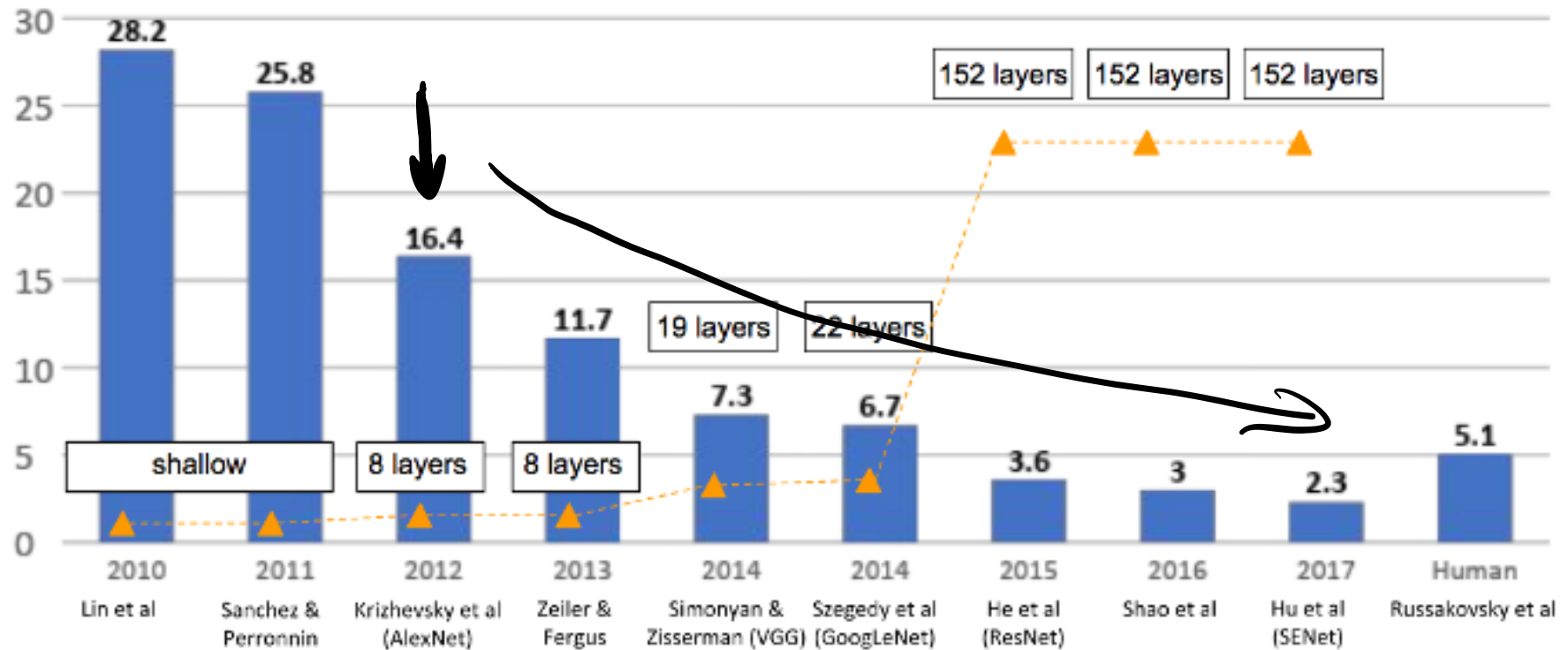
0  1  20

# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 28.2 | 25.8 | 16.4 | 11.7 | 7.3 | 6.7 | 3.6 | 3 | 2.3 | 5.1 |

shallow    8 layers    8 layers    19 layers    22 layers    152 layers    152 layers    152 layers

| 2010 | 2011 | 2012 | 2013 | 2014 | 2014 | 2015 | 2016 | 2017 | Human |
|---|---|---|---|---|---|---|---|---|---|
| Lin et al | Sanchez & Perronnin | Krizhevsky et al (AlexNet) | Zeiler & Fergus | Simonyan & Zisserman (VGG) | Szegedy et al (GoogLeNet) | He et al (ResNet) | Shao et al | Hu et al (SENet) | Russakovsky et al |

AlexNet

$||x|| \times 3$

$5 \times 5 \times 96$

Input data     Conv1     Conv2     Conv3     Conv4     Conv5     FC6   FC7   FC8



$13 \times 13 \times 384$    $13 \times 13 \times 384$    $13 \times 13 \times 256$

$27 \times 27 \times 256$

$55 \times 55 \times 96$

$227 \times 227 \times 3$

1000

4096    4096

# VGG (2014)  VGG19

WE NEED TO GO DEEPER

# GoogLeNet (2014)

## Inception "network-in-network"



concat

1x1    3x3    5x5    1x1

1x1    1x1    pool

in

WE NEED TO GO DEEPER

$$x + F(x)$$

$$F(x)$$

layer

layer

identity

$$x$$

Do we?

# Comparing complexity...

Inception-v4: Resnet + Inception!



An Analysis of Deep Neural Network Models for Practical Applications, 2017.

# Okay but the data...

# Transfer Learning / finetuning



fc3

feature extractor

1000 x 1

4096 x 1

fc3

4096

# Unsupervised / self-supervised learning case study: SimCLR



**A Simple Framework for Contrastive Learning of Visual Representations**

(a) Original    (b) Crop and resize    (c) Crop, resize (and flip)    (d) Color distort. (drop)    (e) Color distort. (jitter)

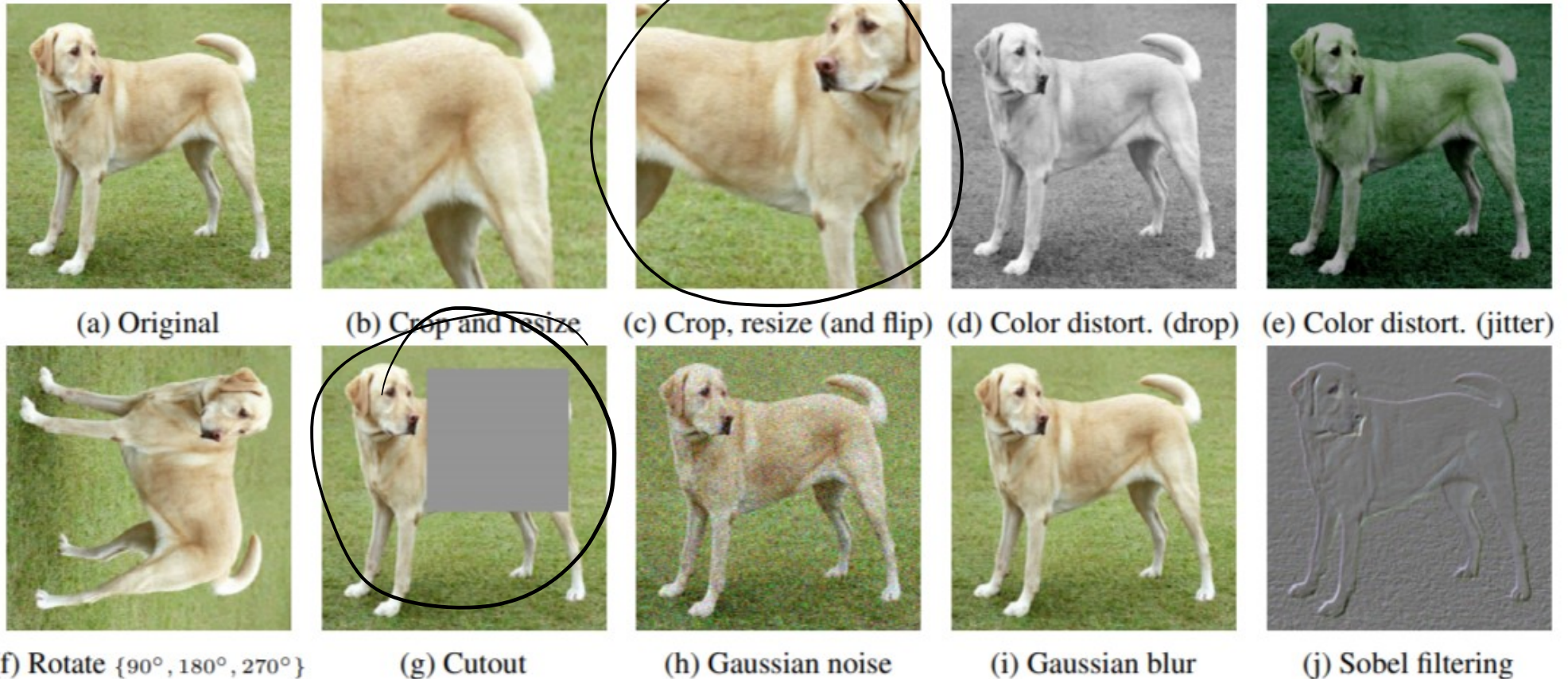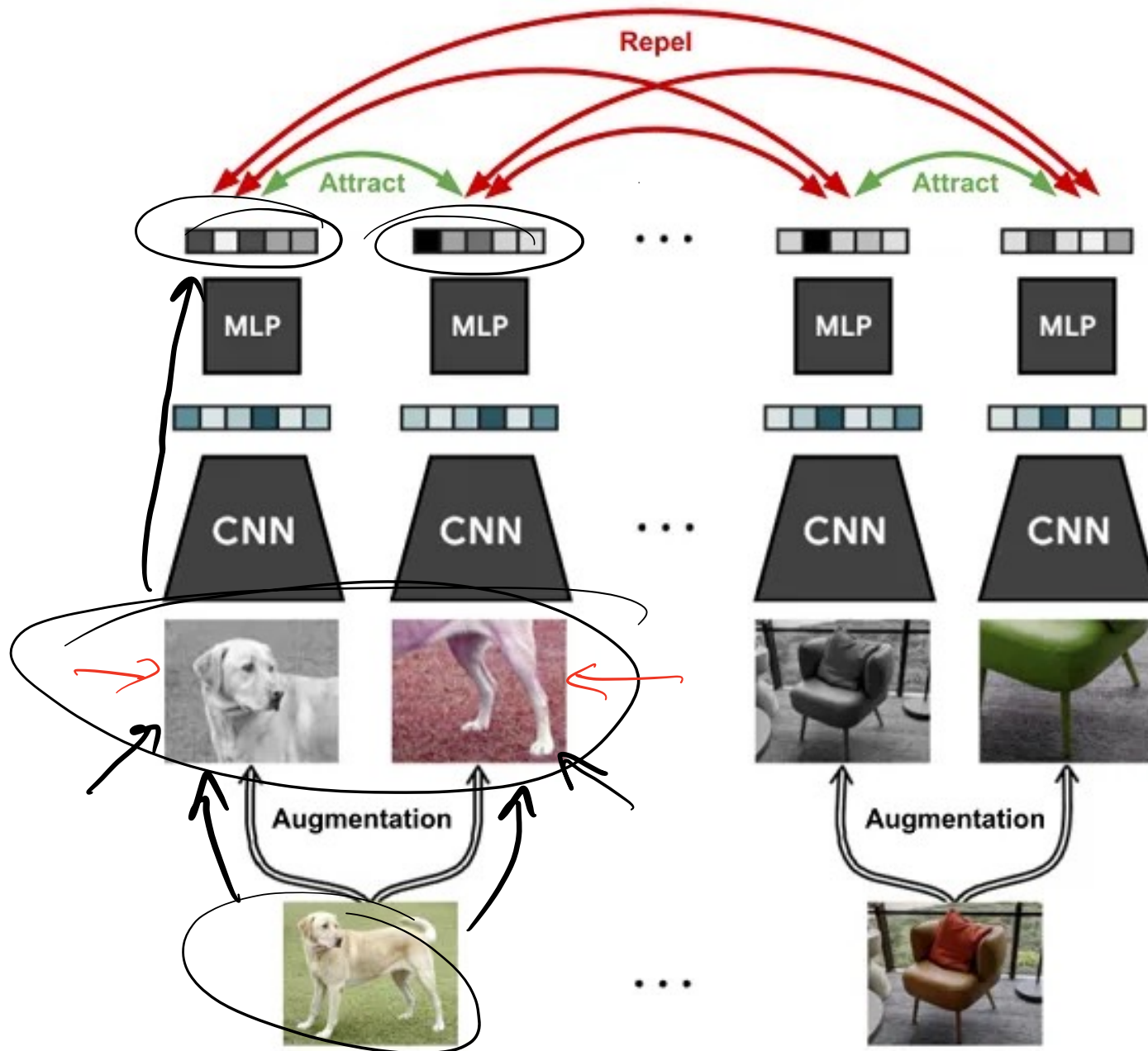(f) Rotate {90°, 180°, 270°}    (g) Cutout    (h) Gaussian noise    (i) Gaussian blur    (j) Sobel filtering

*Figure 4.* Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize)*, *color distortion*, and *Gaussian blur.* (Original image cc-by: Von.grzanka)
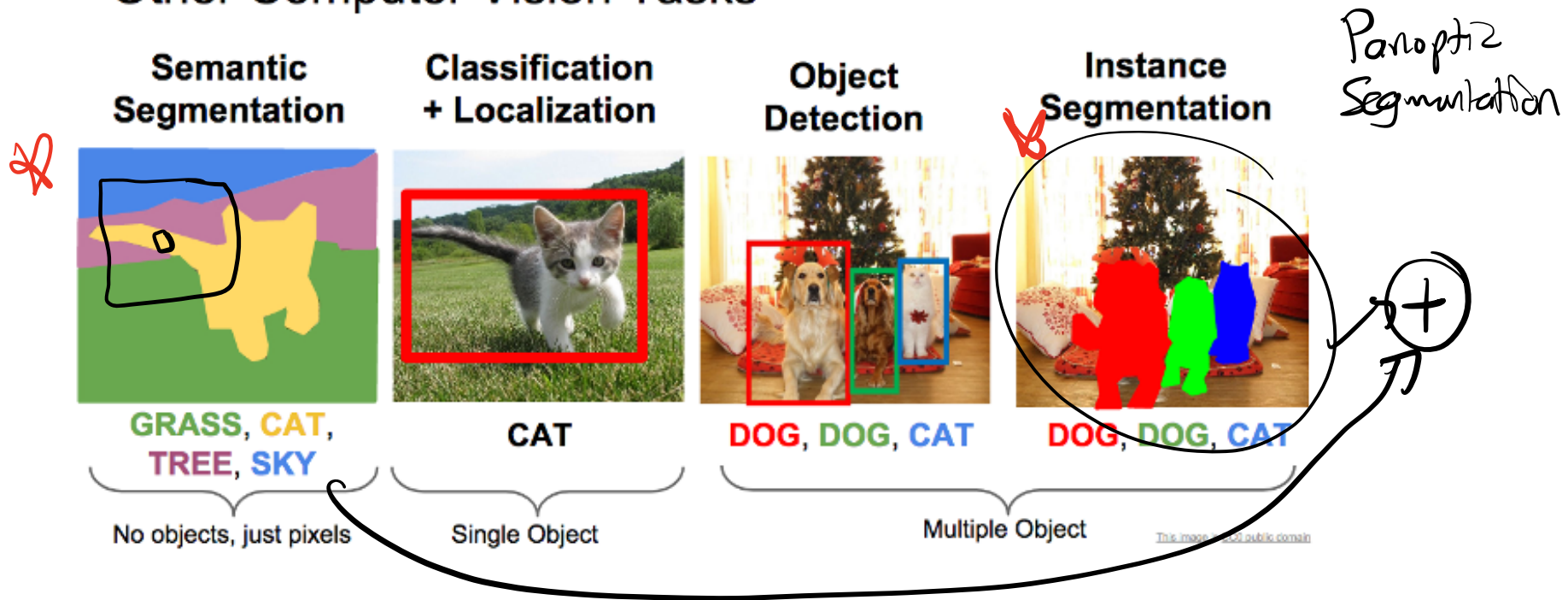
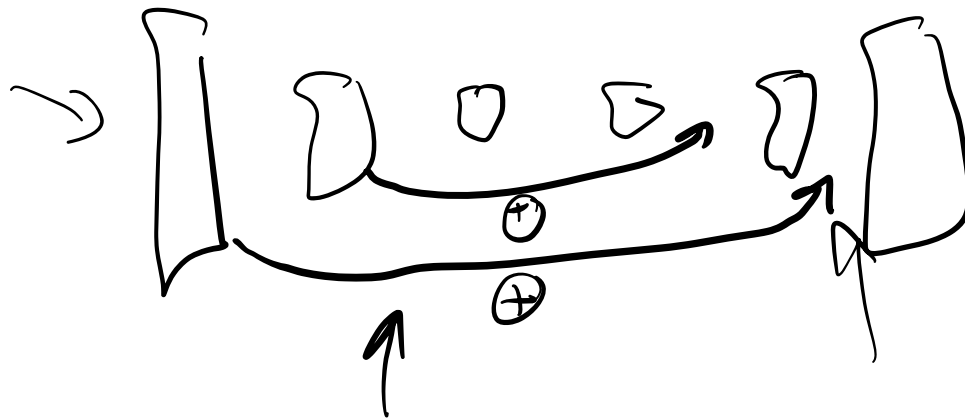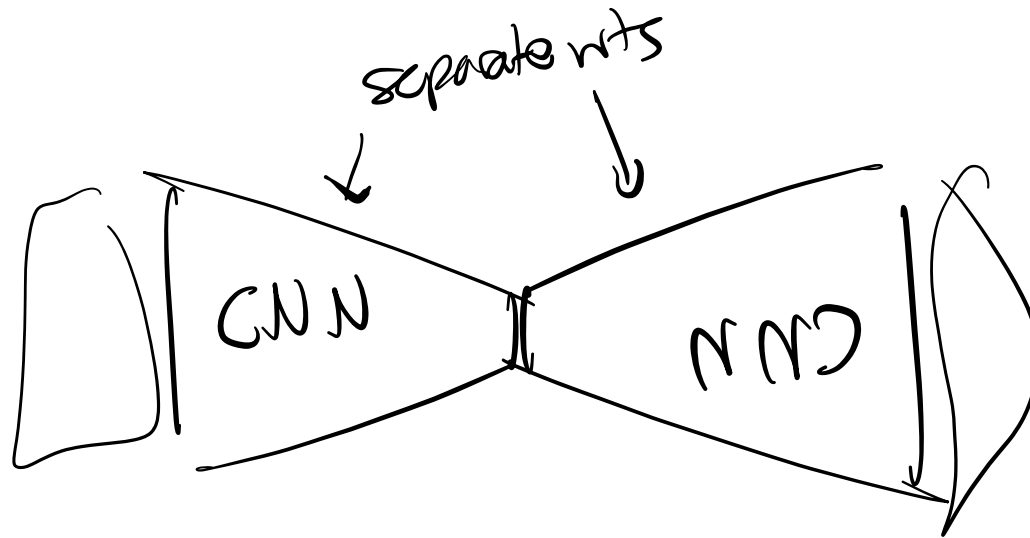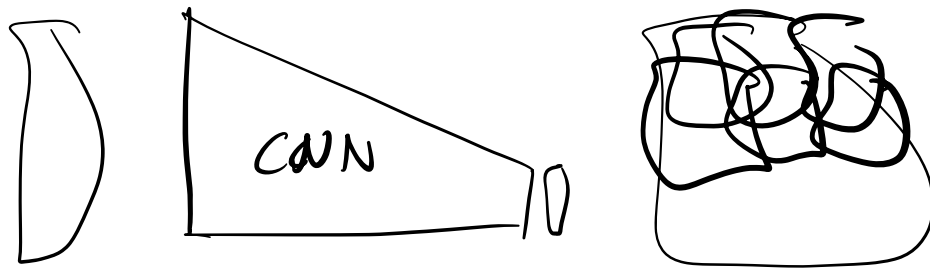# Unsupervised / self-supervised learning case study: SimCLR

# What about not image recognition?



**Other Computer Vision Tasks**

| Semantic Segmentation | Classification + Localization | Object Detection | Instance Segmentation |
|---|---|---|---|
| GRASS, CAT, TREE, SKY | CAT | DOG, DOG, CAT | DOG, DOG, CAT |
| No objects, just pixels | Single Object | Multiple Object | |

*Panoptic Segmentation*

This image is CC0 public domain

CNN
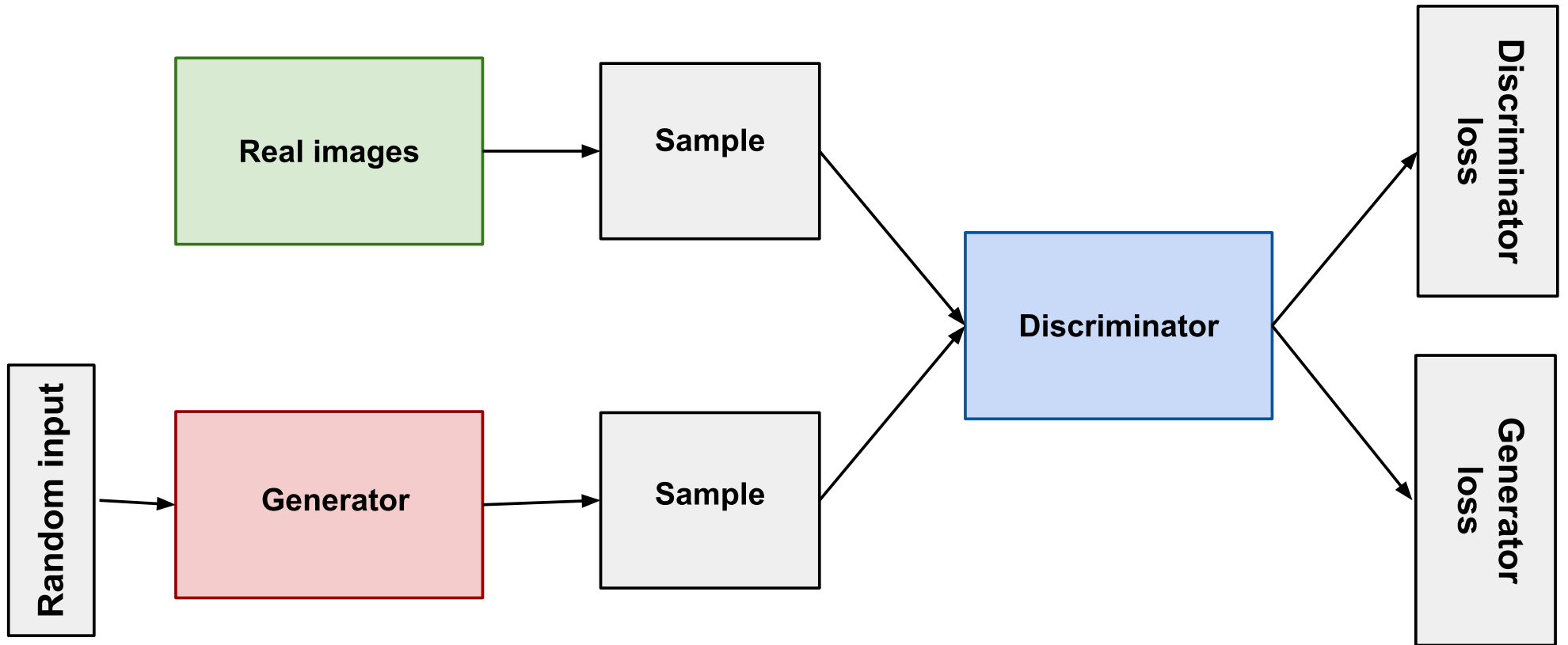
separate wts

CNN NNC

U Net

# (Sharp?) left turn:
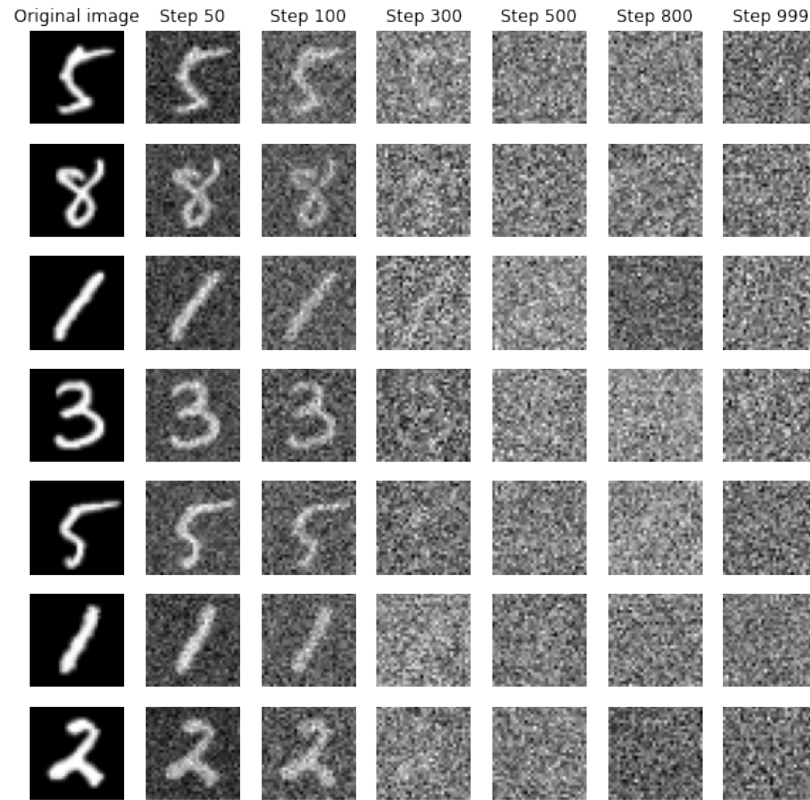# Embeddings, Manifold Learning, and Autoencoders
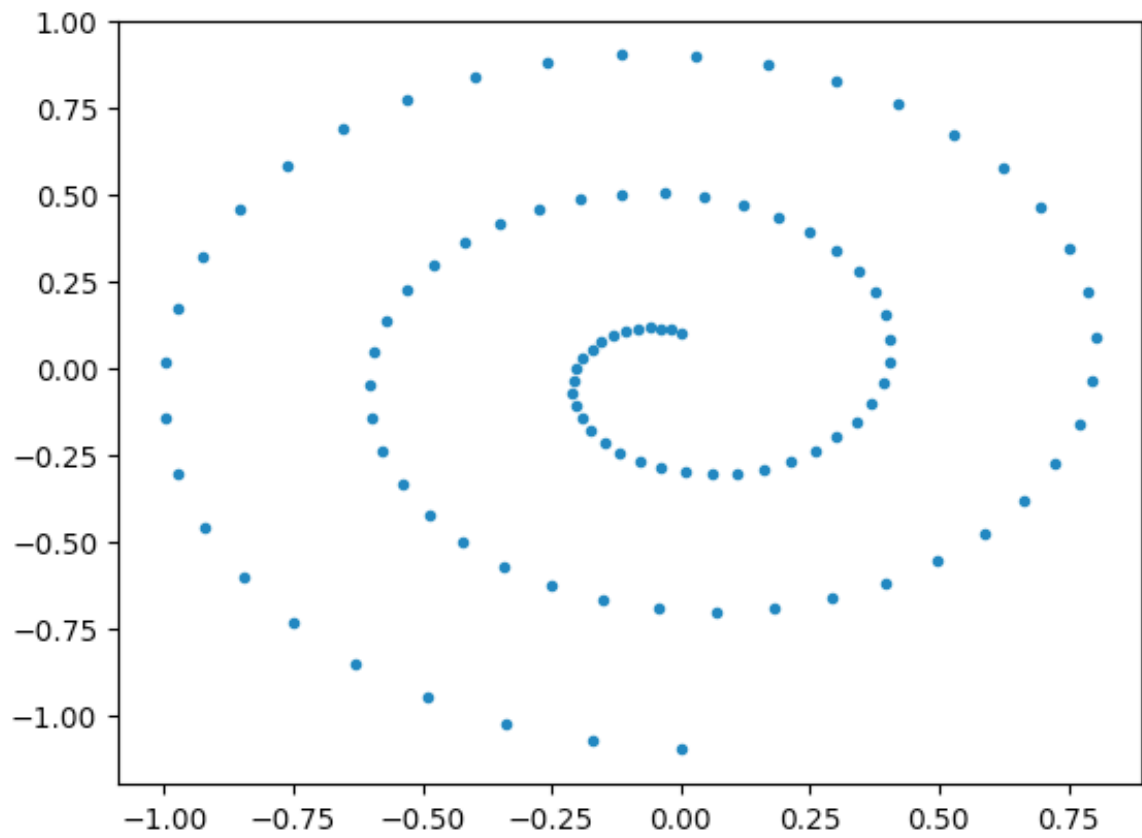
# Generative Modeling

# Generative Adversarial Networks

# Diffusion Models
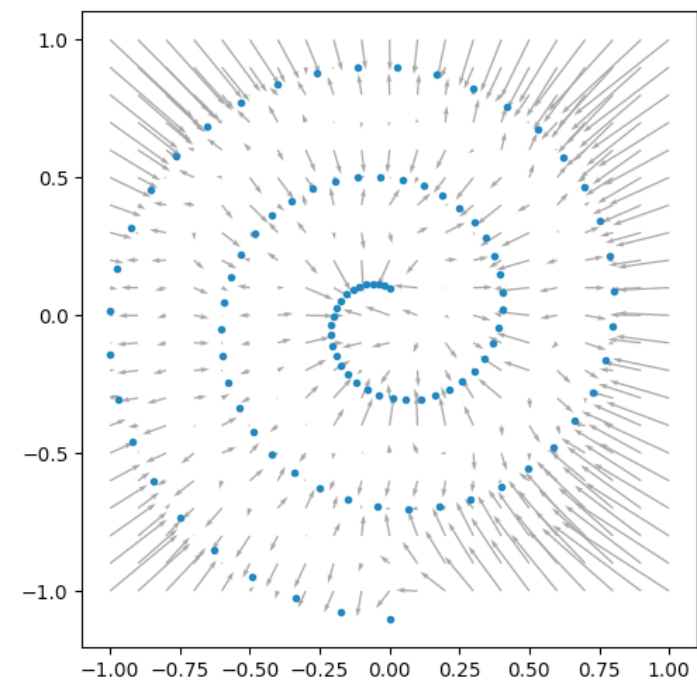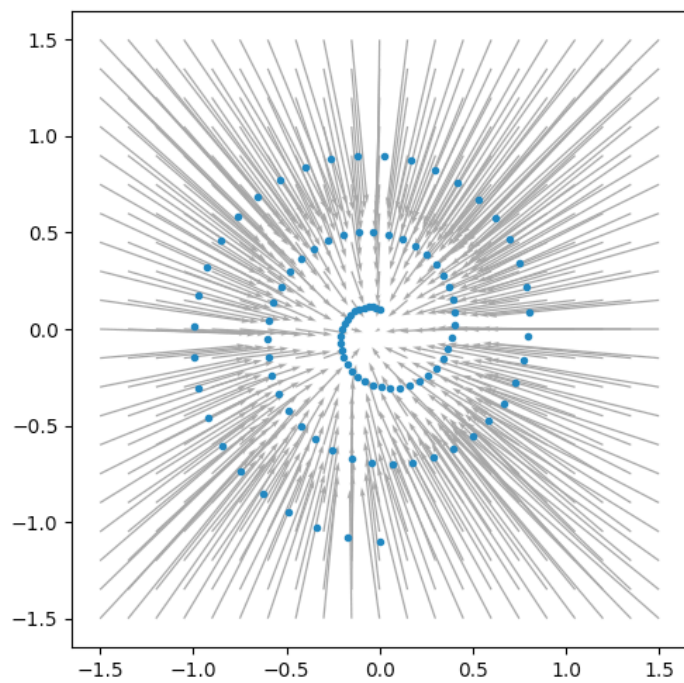


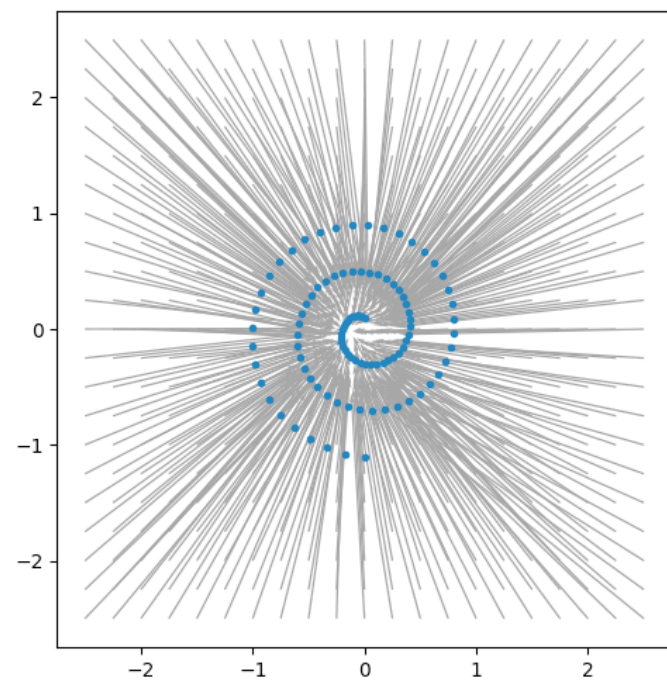Some other good visuals: https://www.chenyang.co/diffusion.html

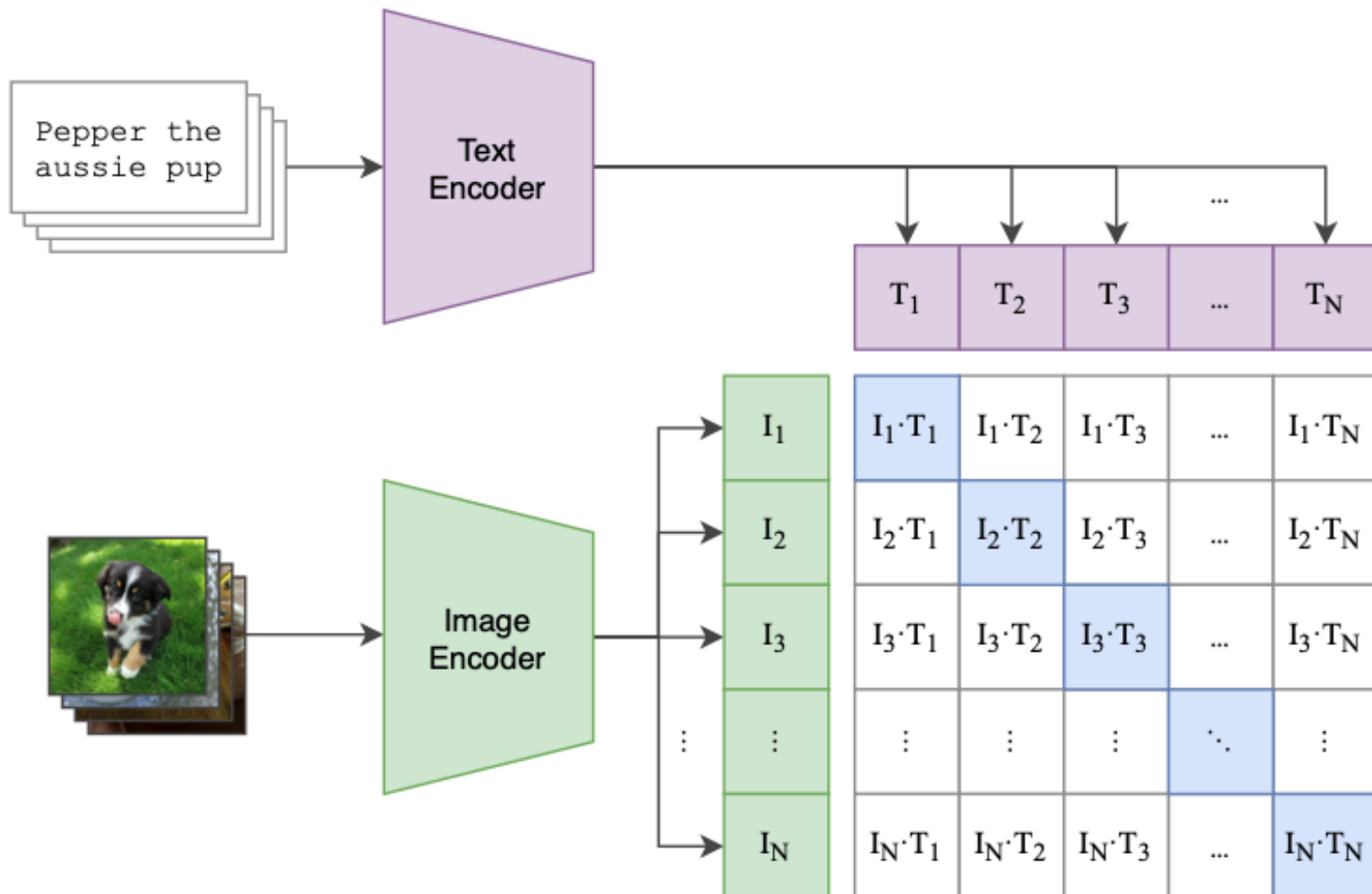$\sigma = 0.1$      $\sigma = 0.5$      $\sigma = 1$

# Stable Diffusion
# (without the text-conditioning)

# Vision and Language

# Case study: CLIP



(1) Contrastive pre-training
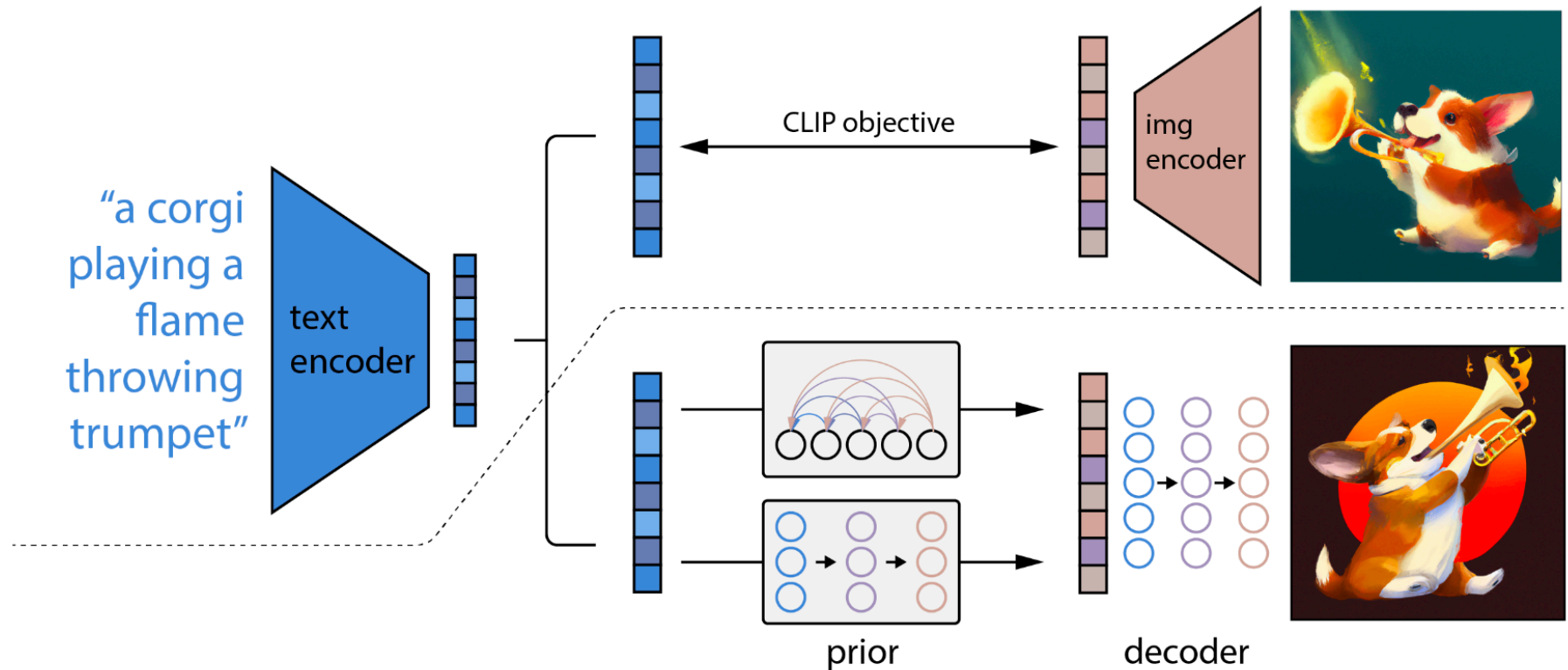
# unCLIP aka DALL-E 2



Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.

# Stable Diffusion
## (with the text-conditioning)