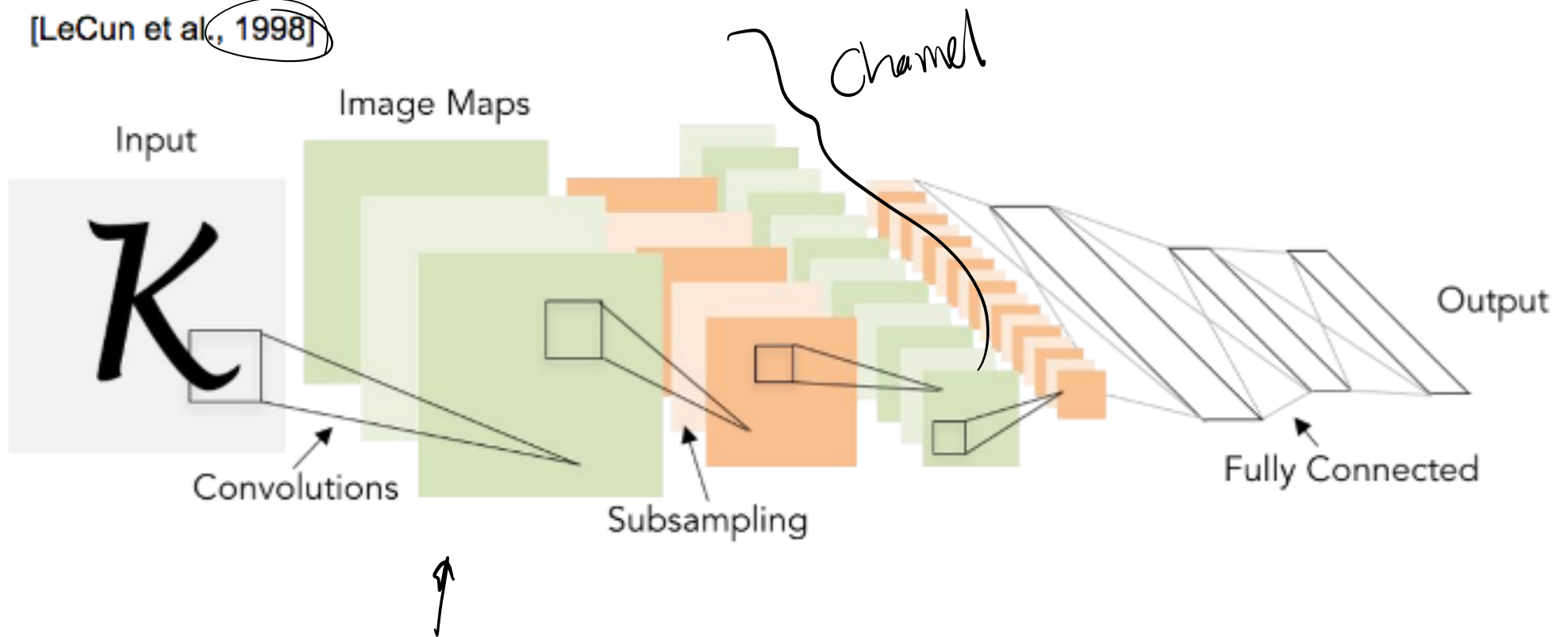
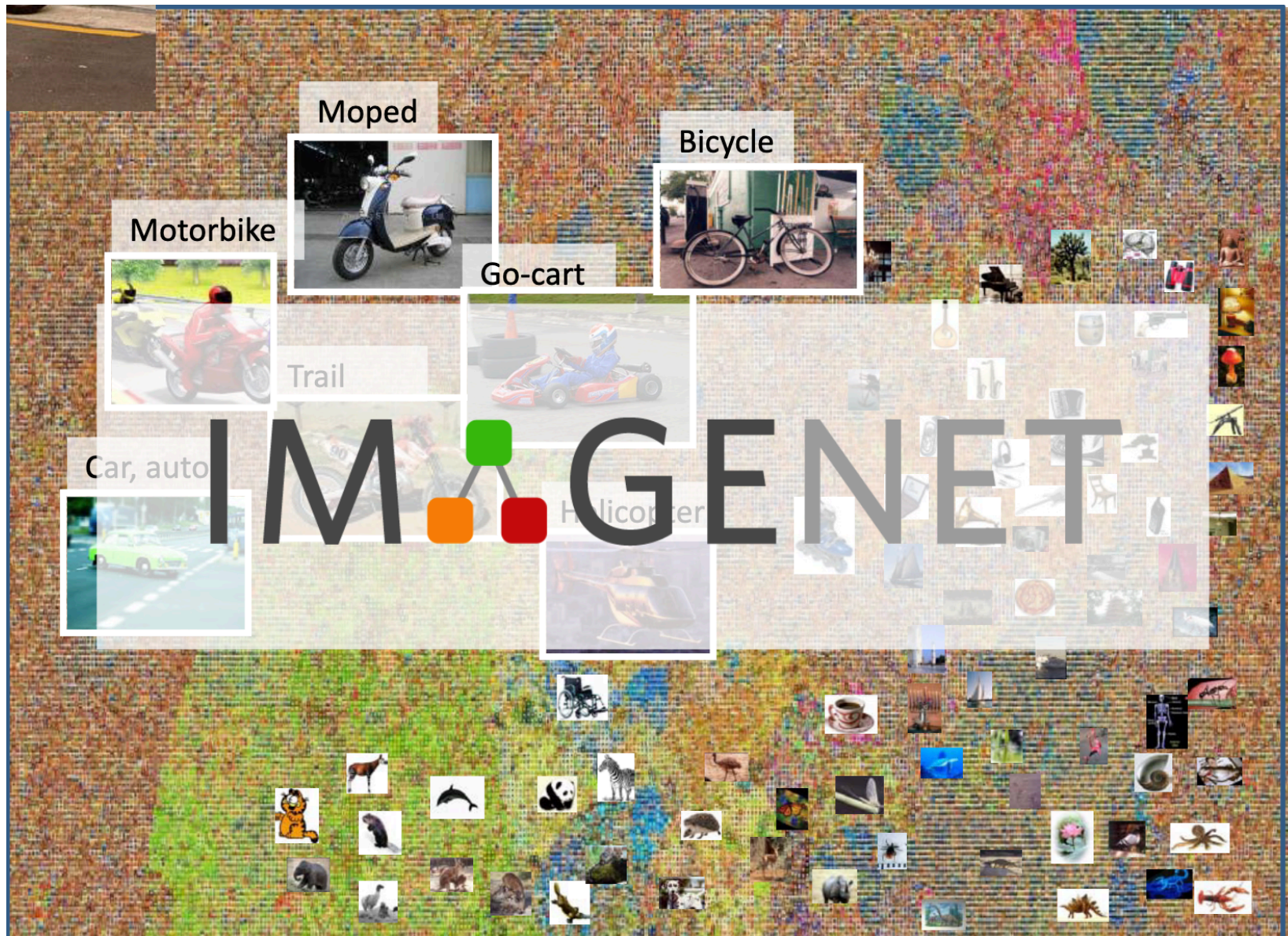




# Review: LeNet-5

[LeCun et al., 1998]





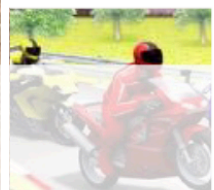
Moped



Bicycle



Motorbike



Go-cart



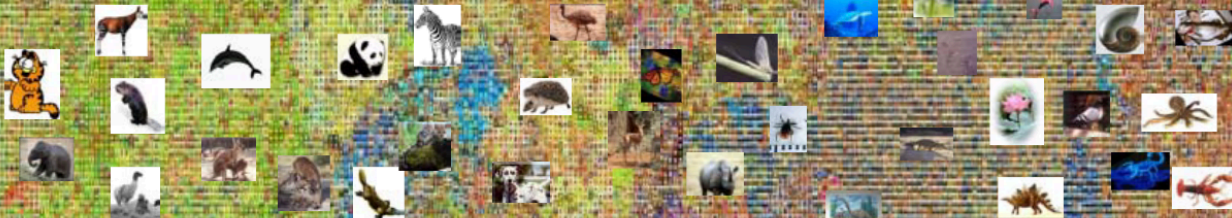
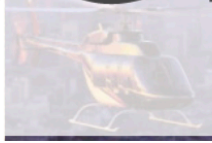
Trail

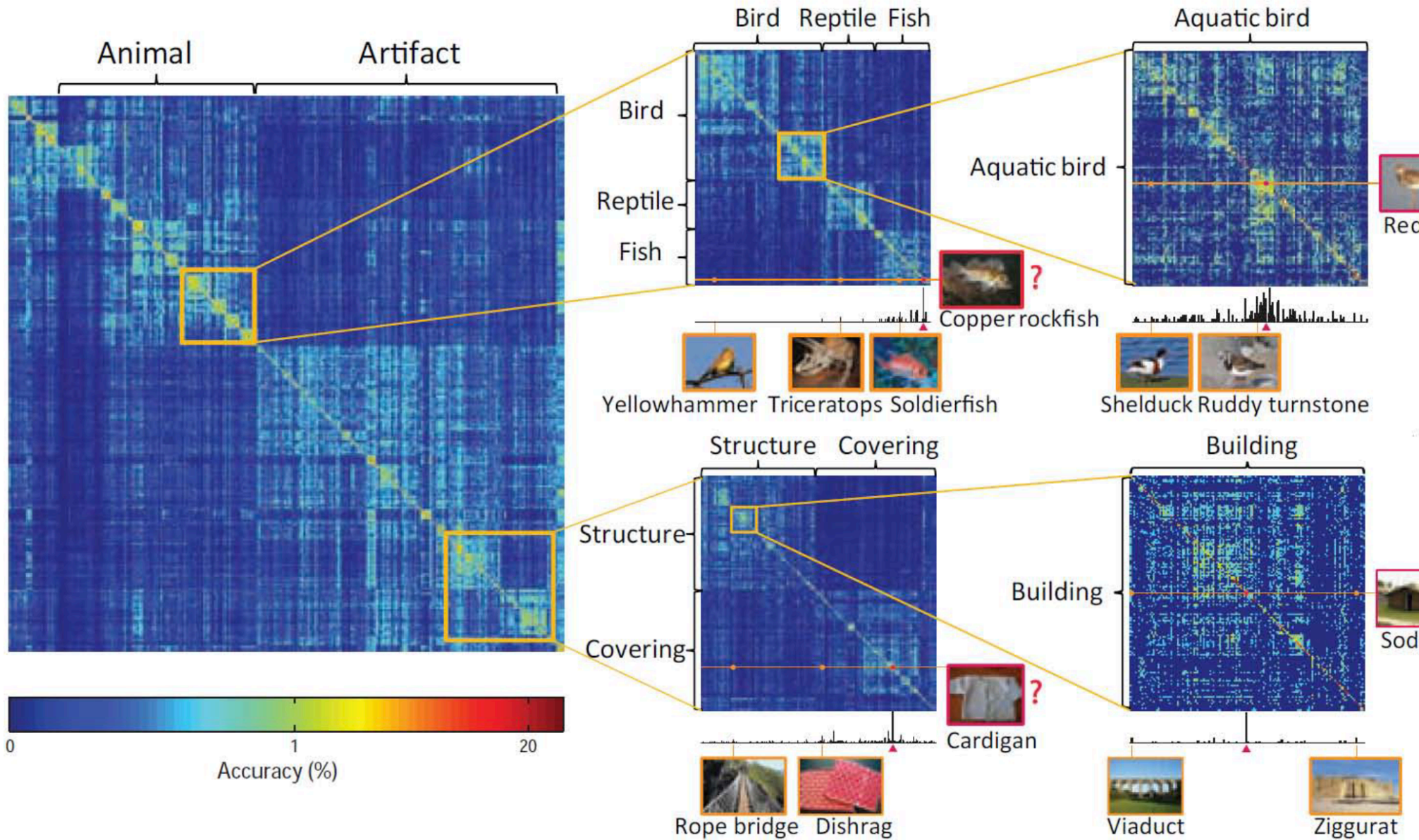
Car, auto



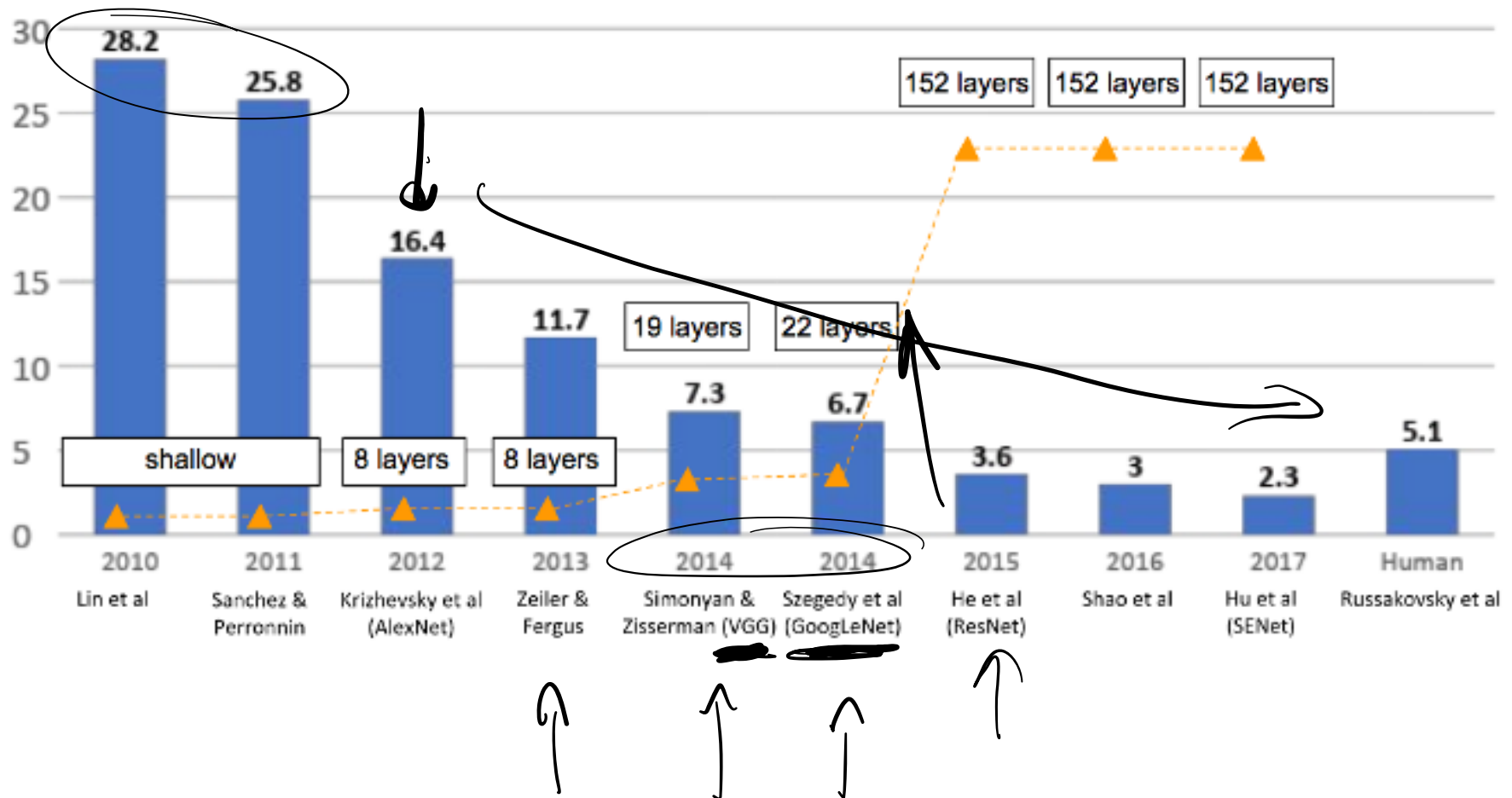
# IMAGENET

Helicopter





# ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

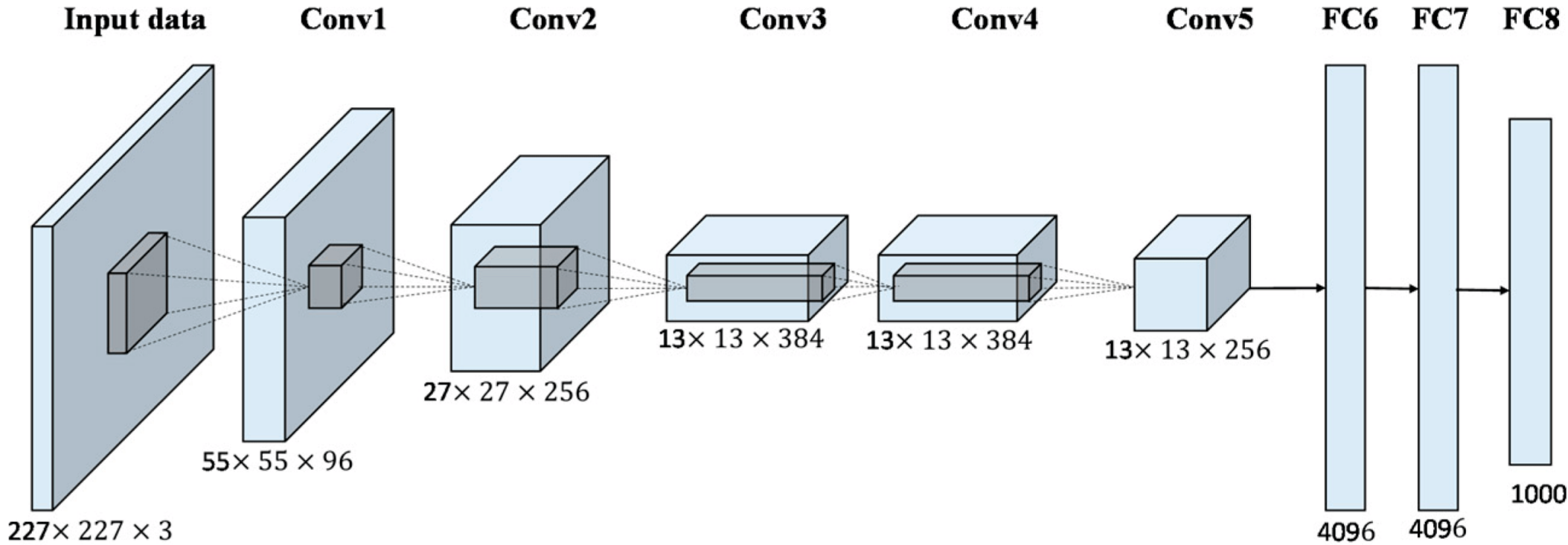




A still from the movie Inception showing Leonardo DiCaprio and Matt Damon in a car. DiCaprio is on the left, looking towards Damon on the right. The scene is dimly lit with light coming from a window behind them. The text "WE NEED TO GO DEEPER" is overlaid in white, bold, sans-serif font at the bottom of the image.

**WE NEED TO GO  
DEEPER**

# Alex Net

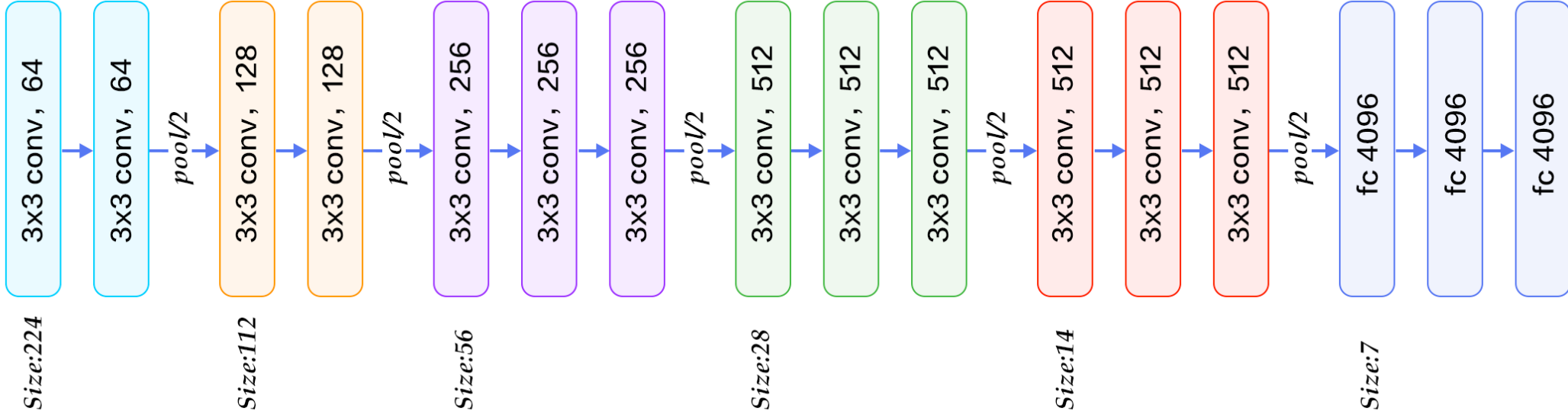




A still from the movie Inception showing Leonardo DiCaprio and Matt Damon in a car. DiCaprio is on the left, looking towards Damon on the right. The scene is dimly lit, with light coming from a window behind them. The text "WE NEED TO GO DEEPER" is overlaid in white, bold, sans-serif font at the bottom of the image.

**WE NEED TO GO  
DEEPER**

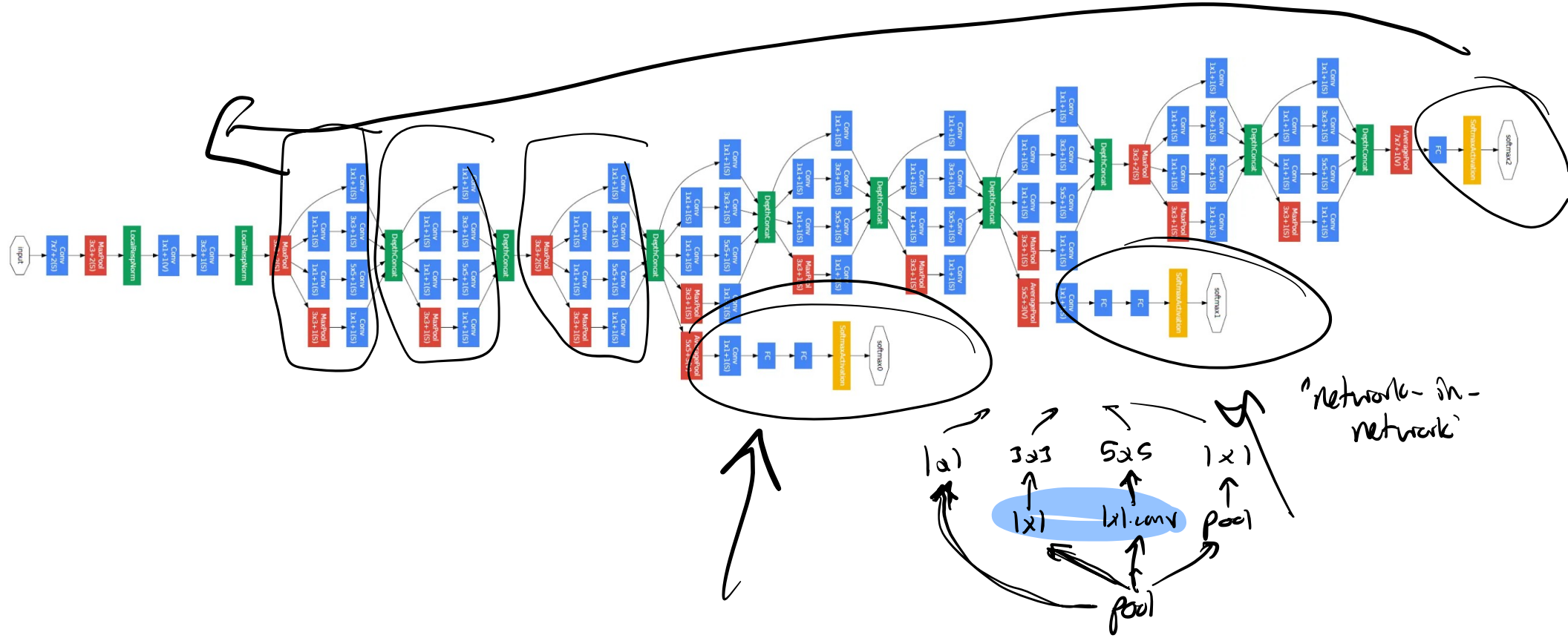
# VGG Net



A still from the movie Inception showing Leonardo DiCaprio and Matt Damon in a car. DiCaprio is on the left, looking towards Damon on the right. The scene is dimly lit with light coming from a window behind them.

**WE NEED TO GO  
DEEPER**

# GoogLeNet (Inception)

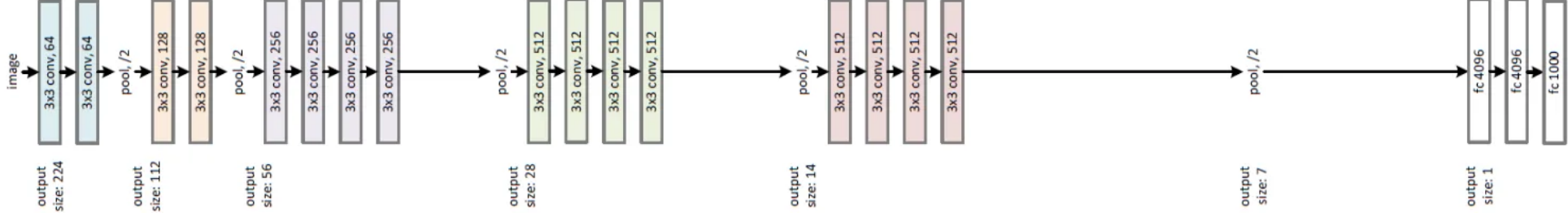




A still from the movie Inception showing Leonardo DiCaprio and Matt Damon in a car. DiCaprio is on the left, looking towards Damon on the right. The scene is dimly lit, with light coming from a window behind them. The text "WE NEED TO GO DEEPER" is overlaid in white, bold, sans-serif font at the bottom of the image.

**WE NEED TO GO  
DEEPER**

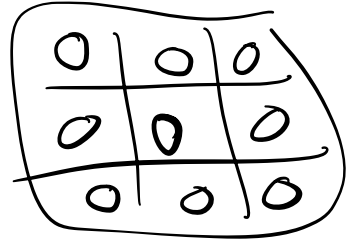
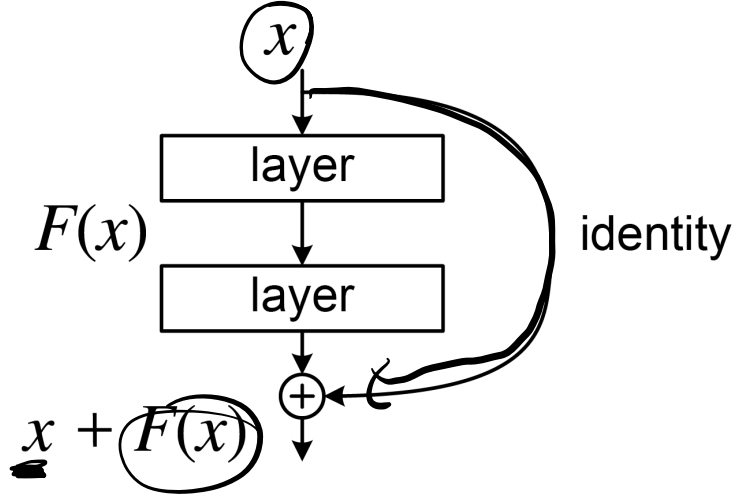
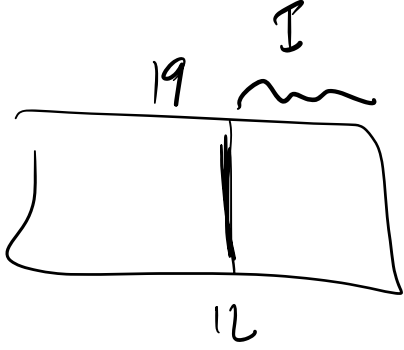
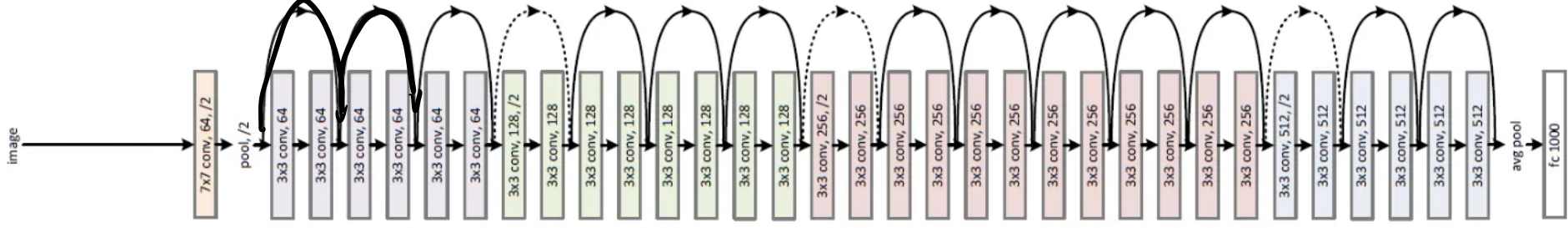
VGG-19



34-layer plain



34-layer residual



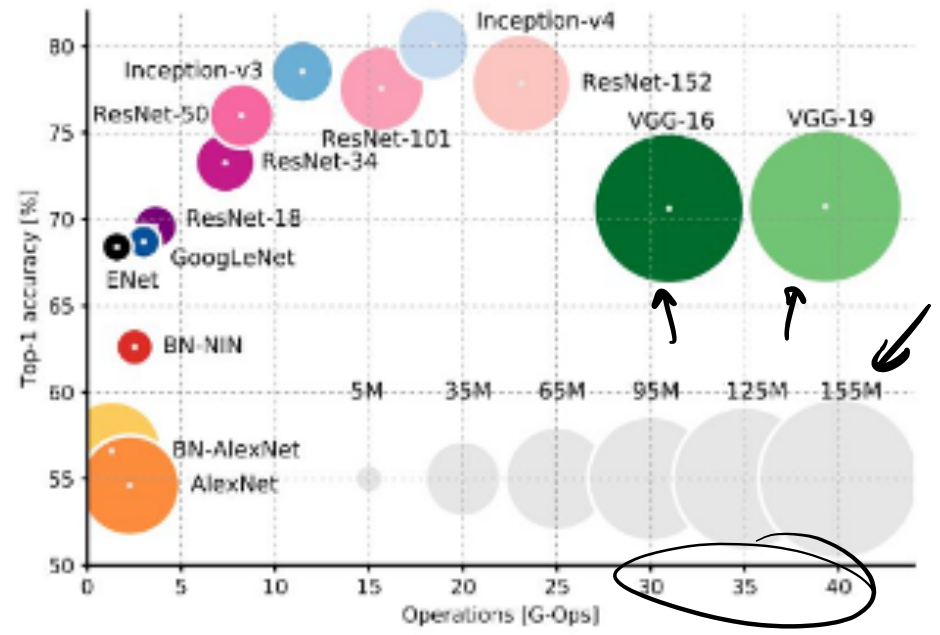
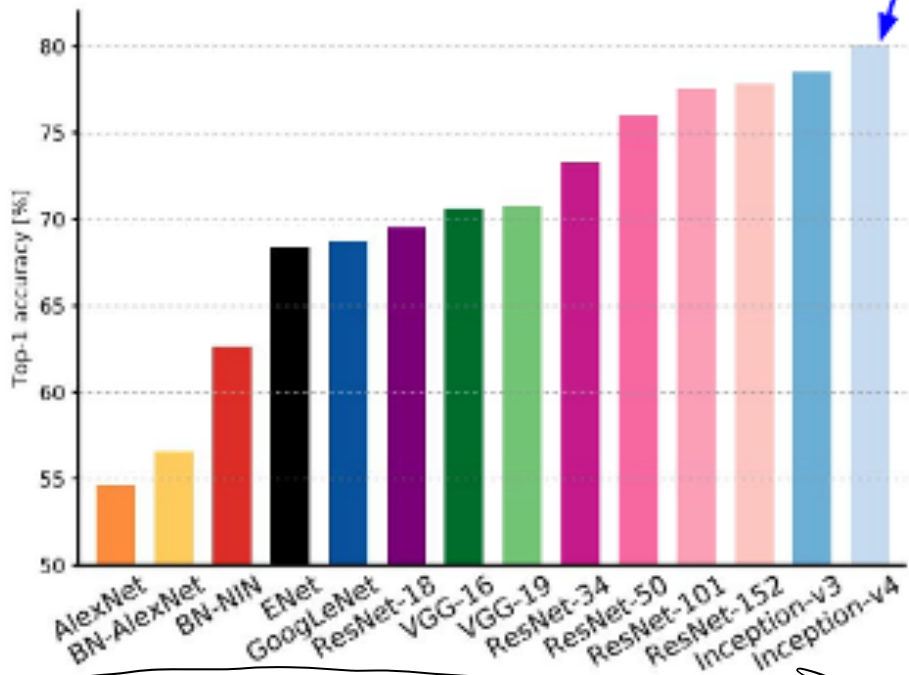


Do we?



# Comparing complexity...

Inception-v4: Resnet + Inception!



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

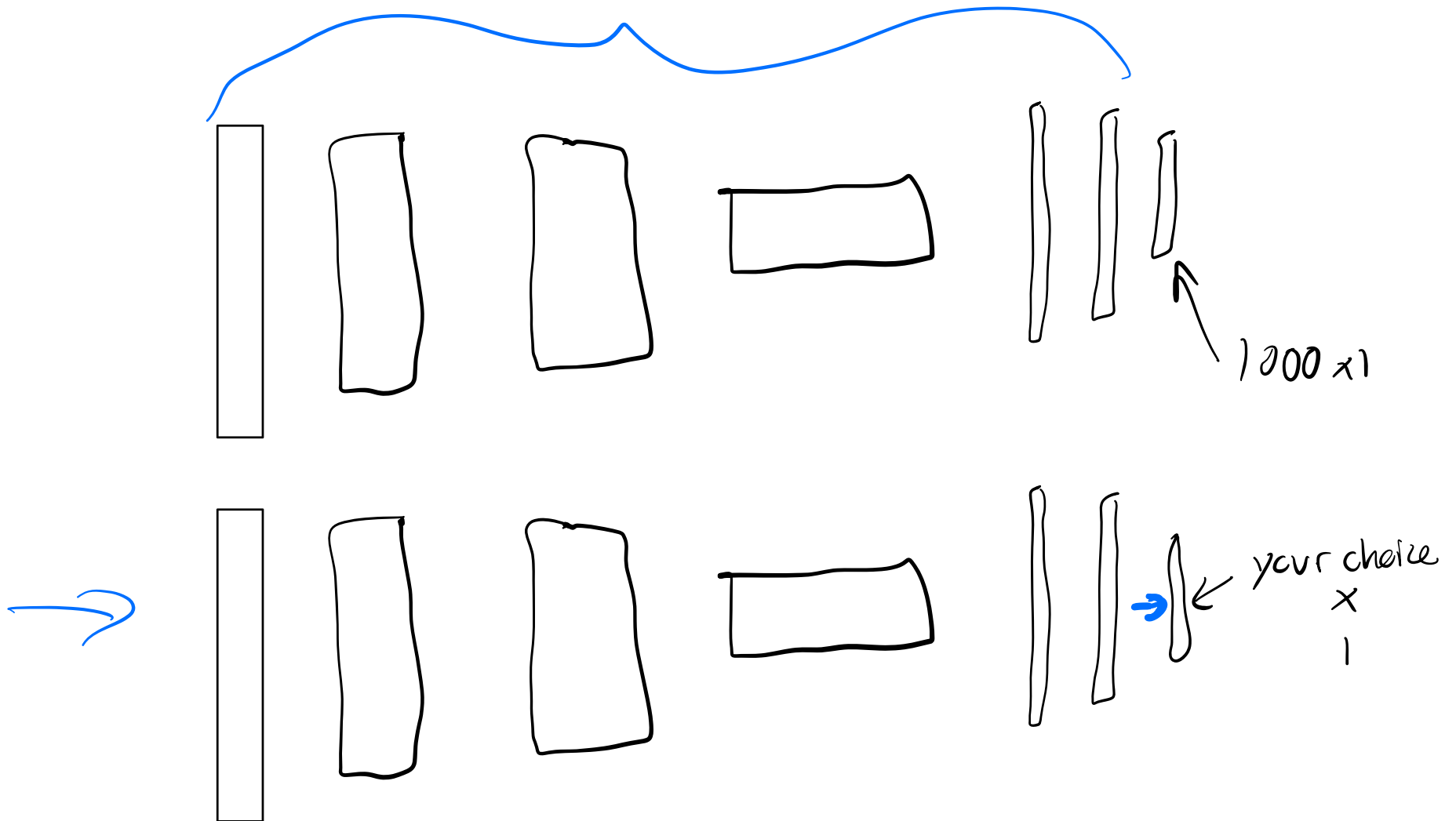
Figures copyright Alfredo Canziani, Adam Paszke, Eugenio Culurciello, 2017. Reproduced with permission.



Okay but the data...

*is expensive!*

# Transfer Learning / finetuning



# Unsupervised / self-supervised learning case study: SimCLR

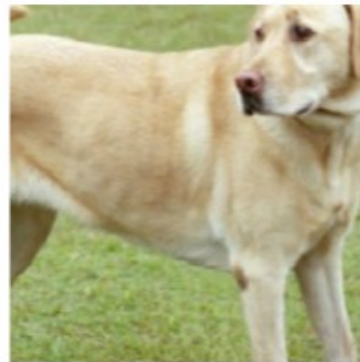
## A Simple Framework for Contrastive Learning of Visual Representations



(a) Original



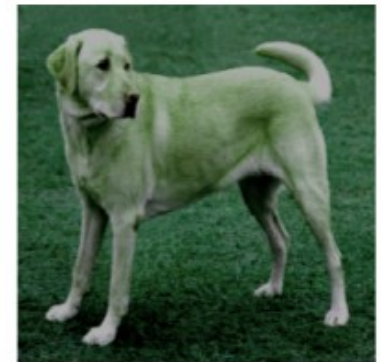
(b) Crop and resize



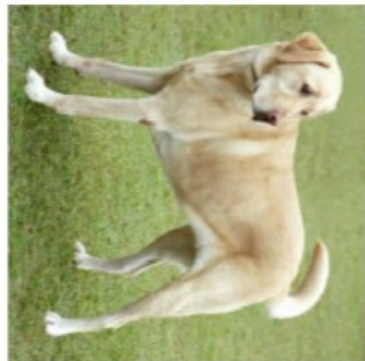
(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate  $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



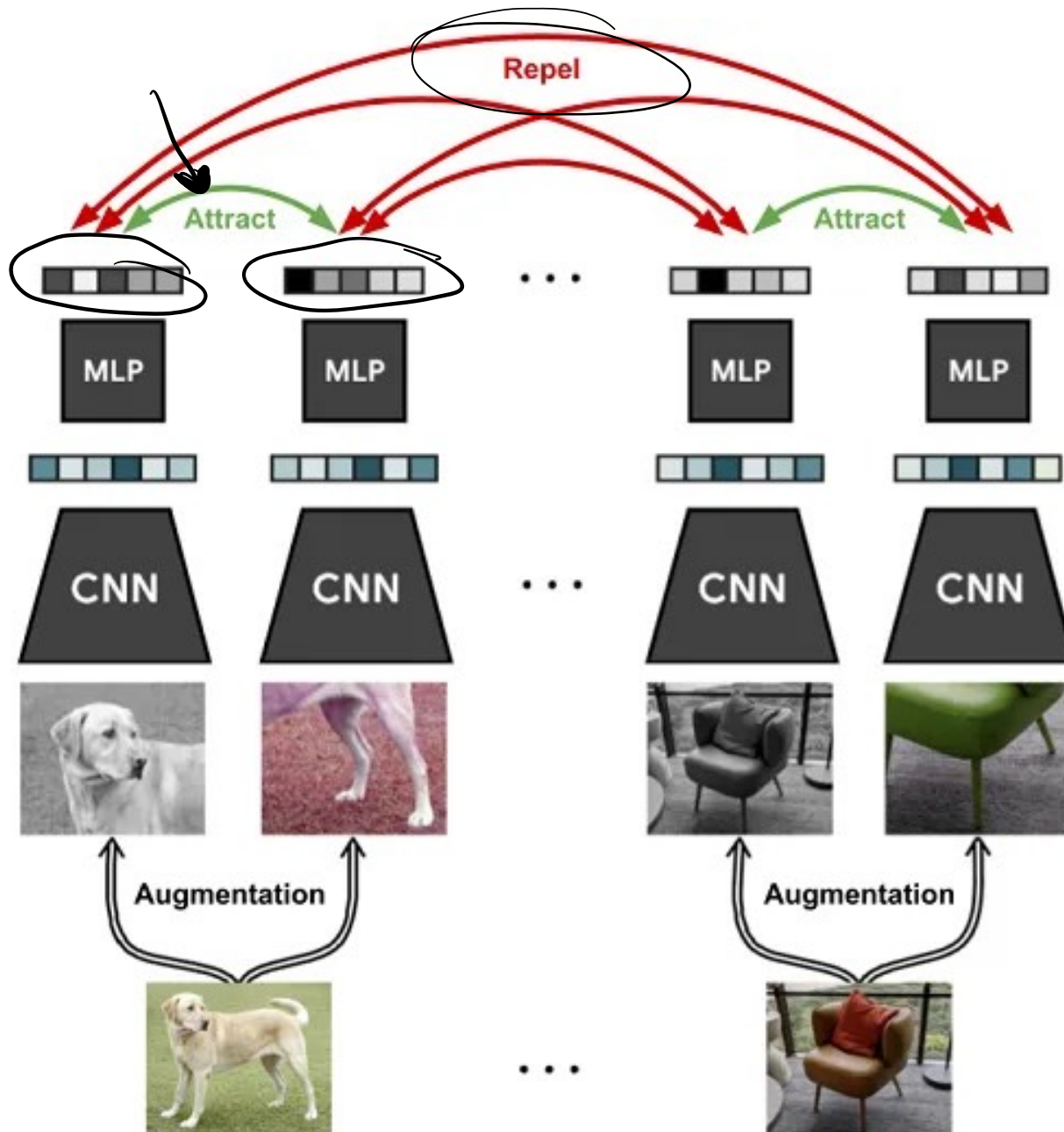
(i) Gaussian blur



(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



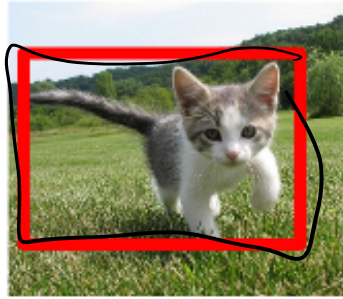
Semantic Segmentation



GRASS, CAT, TREE, SKY

No objects, just pixels

Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

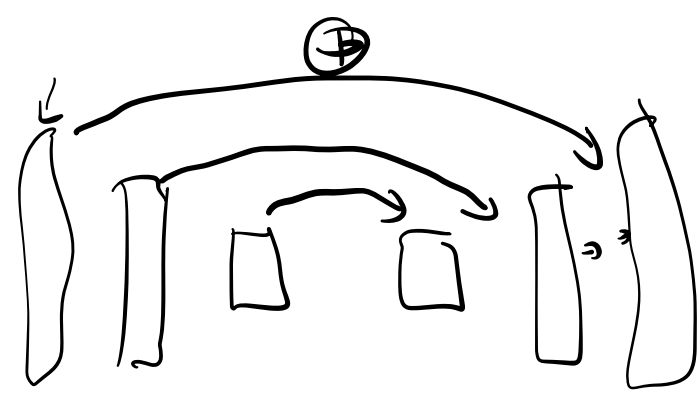
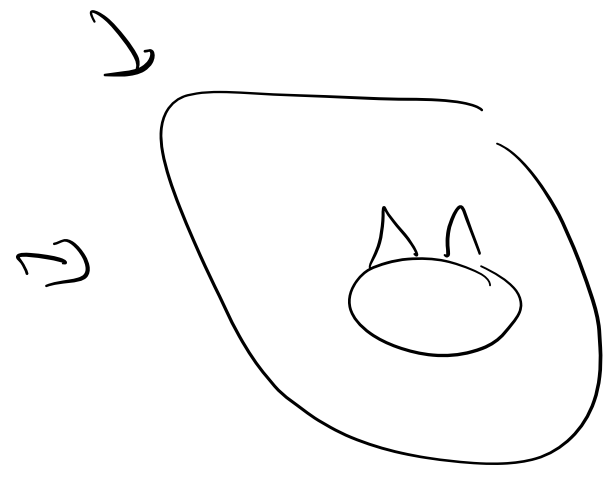
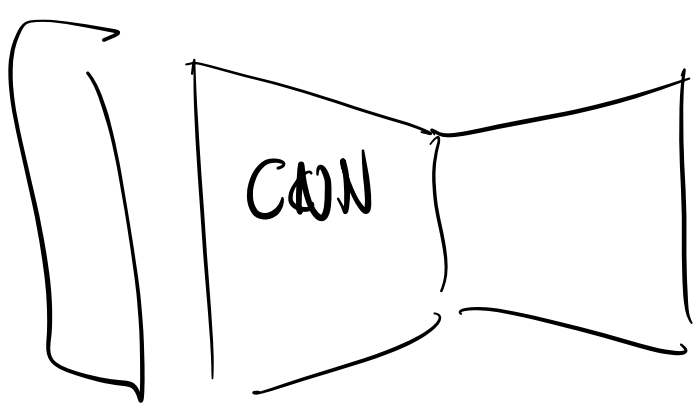
Instance Segmentation



DOG, DOG, CAT

*Panoptic Segmentation*

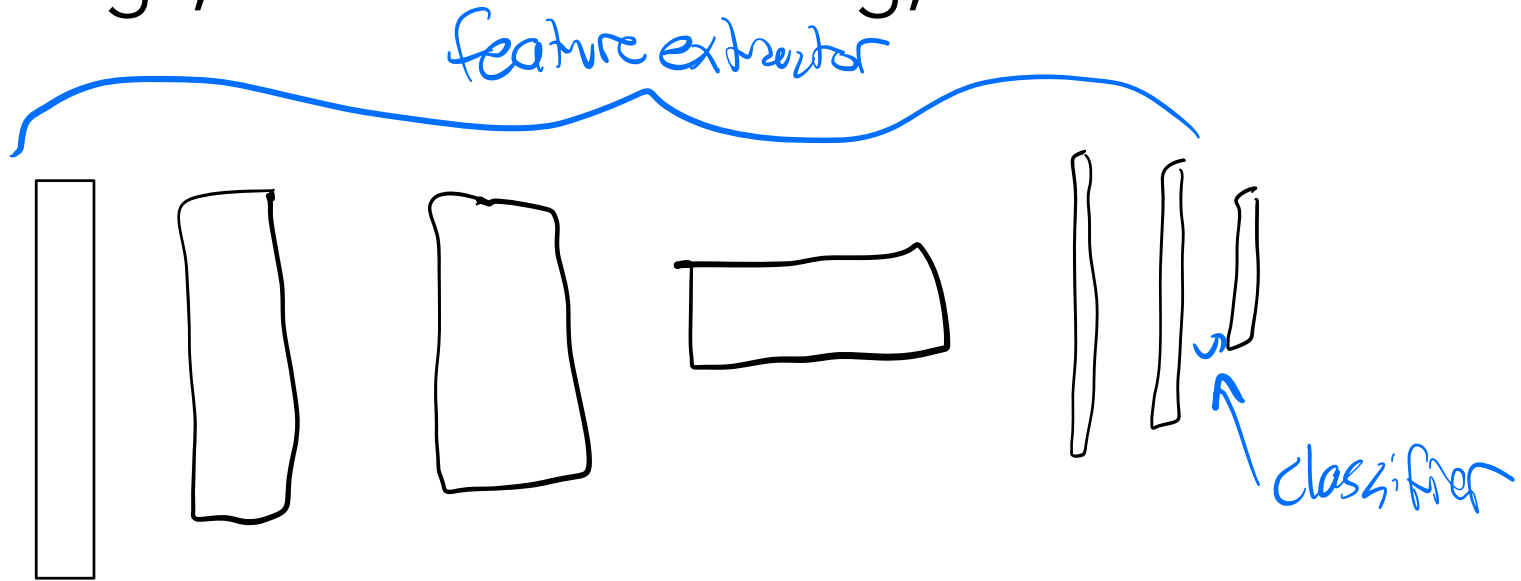






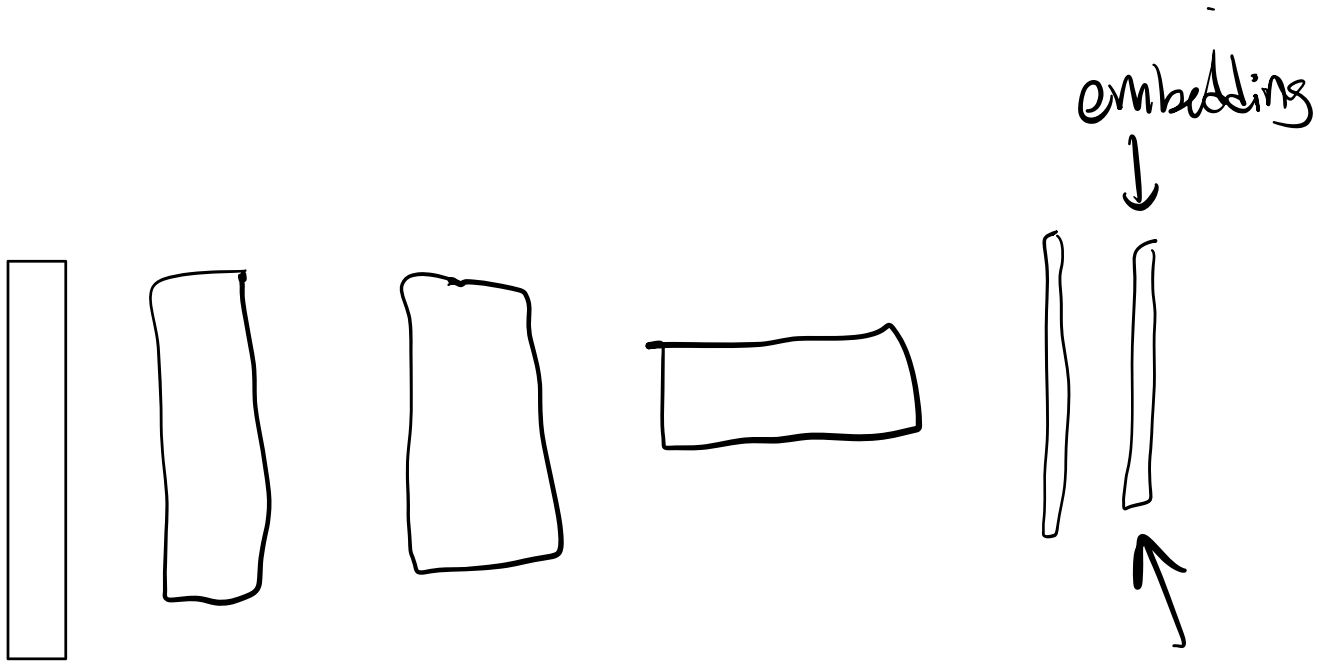
(Sharp?) left turn:

# Embeddings, Manifold Learning, and Autoencoders

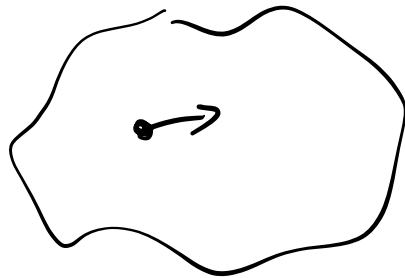


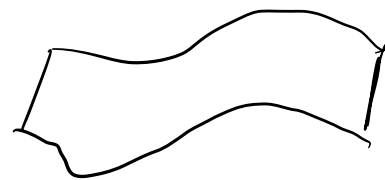
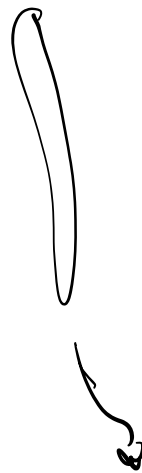
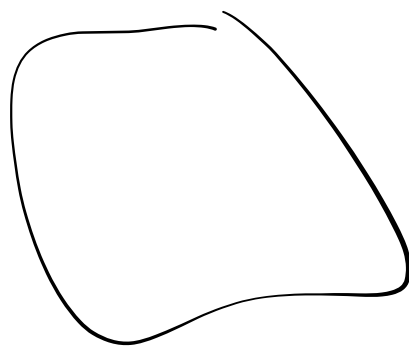
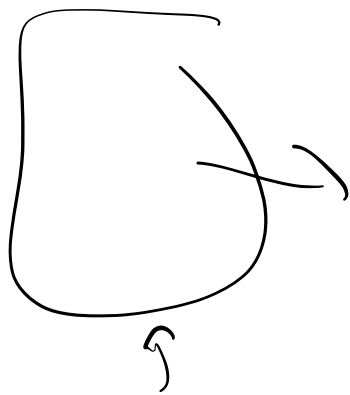
feature extraction  $\rightarrow$  classify

Img  $\rightarrow$  conv  $\rightarrow$  group edges into contours  
↓  
group into obj. parts  
↓  
SVM



"King" + "Woman" → "Queen"



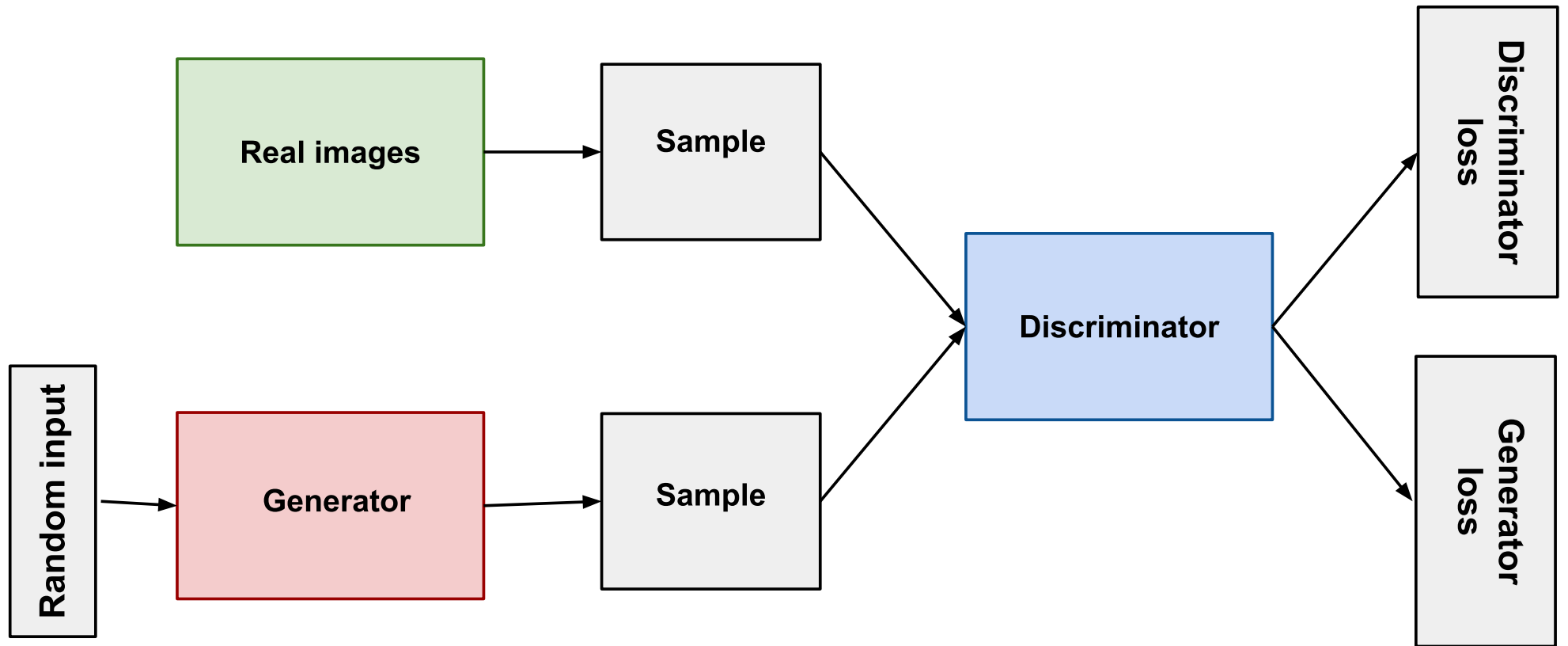




# Generative Modeling



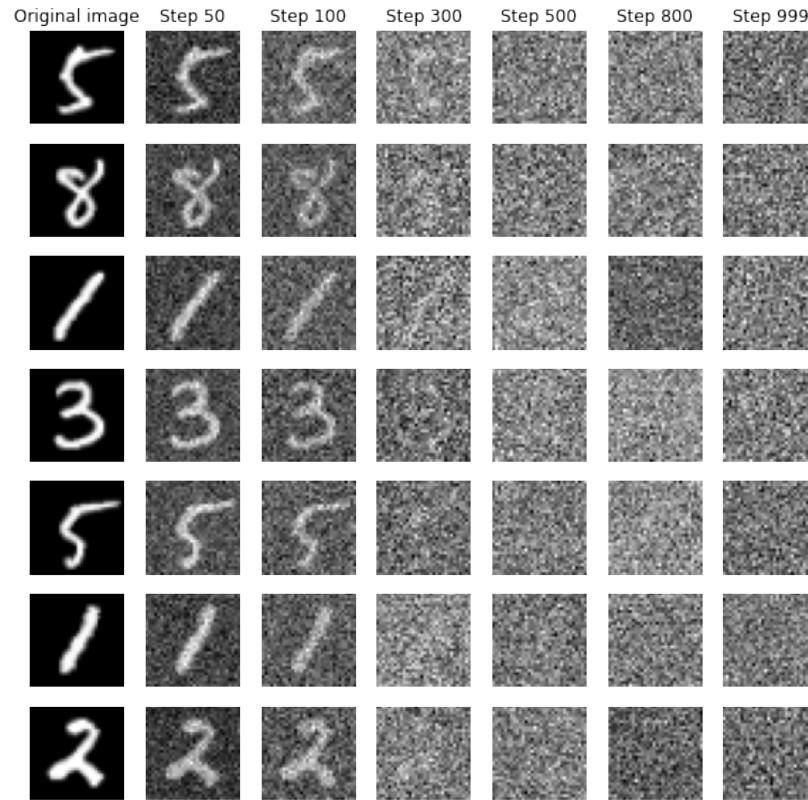
# Generative Adversarial Networks



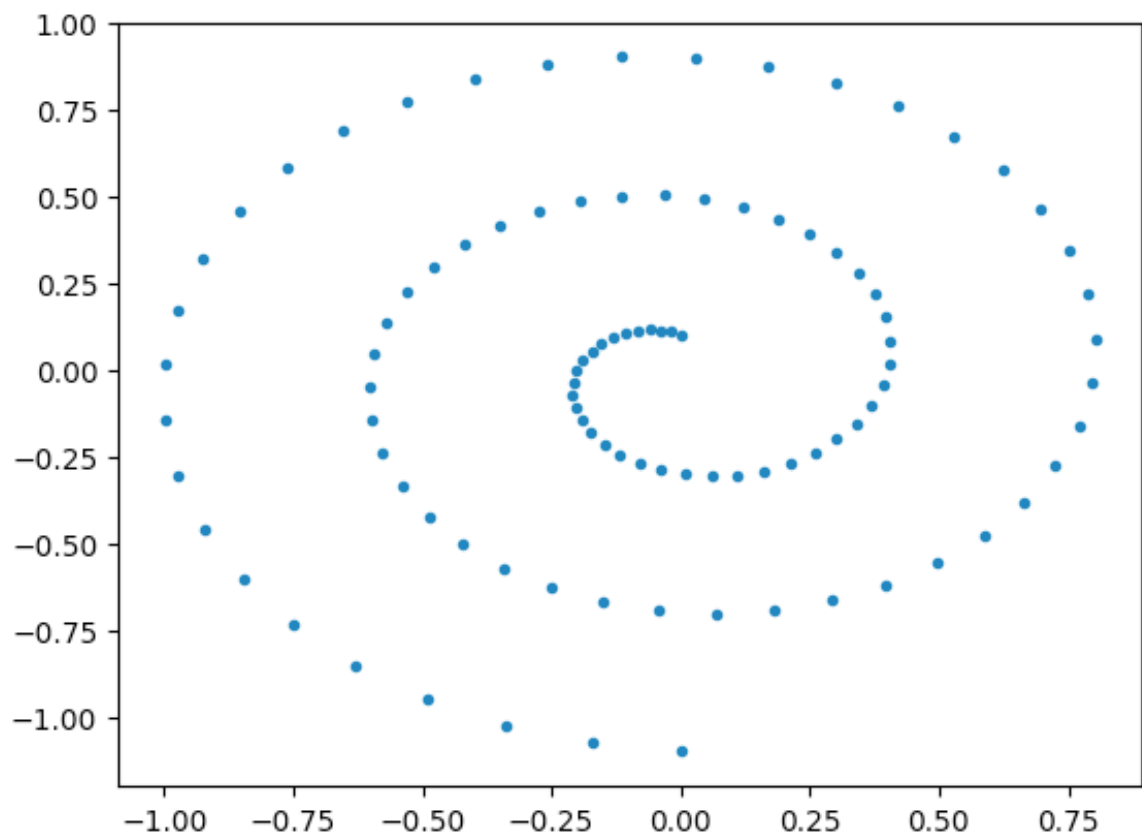




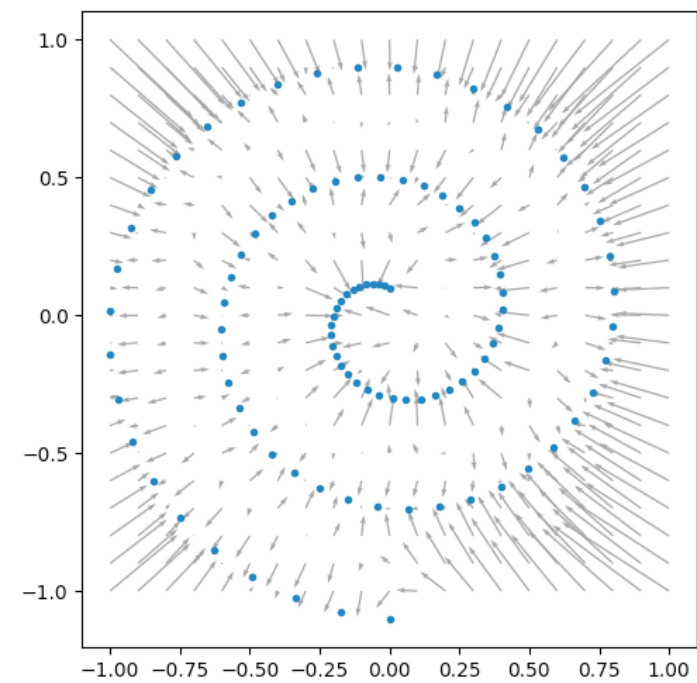
# Diffusion Models



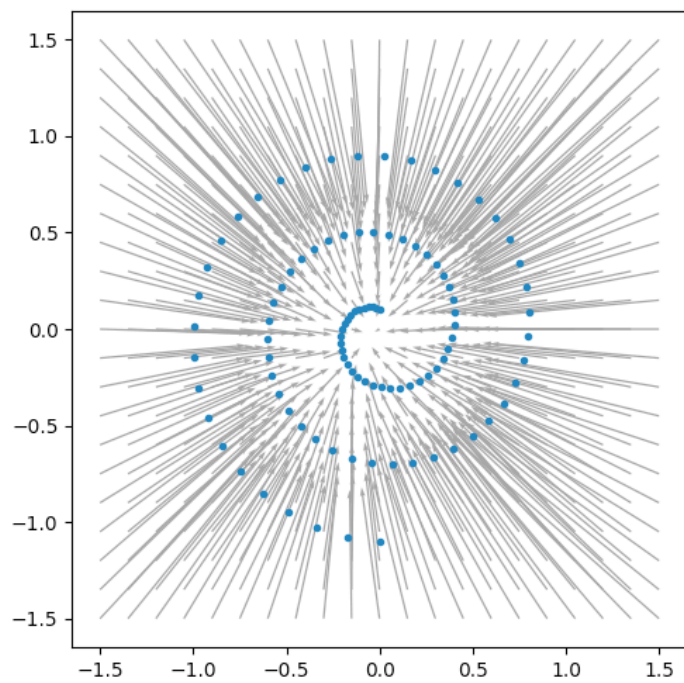
Some other good visuals: <https://www.chenyang.co/diffusion.html>



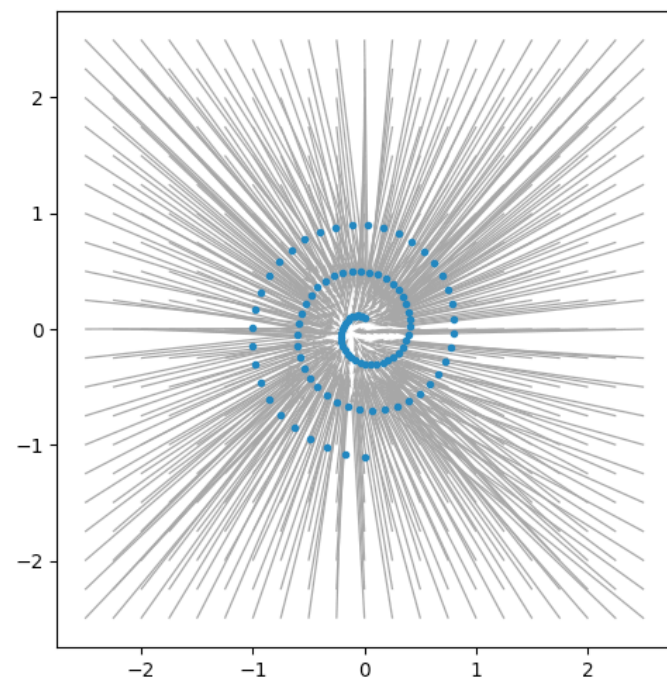
$\sigma = 0.1$



$\sigma = 0.5$



$\sigma = 1$









# Stable Diffusion

(without the text-conditioning)

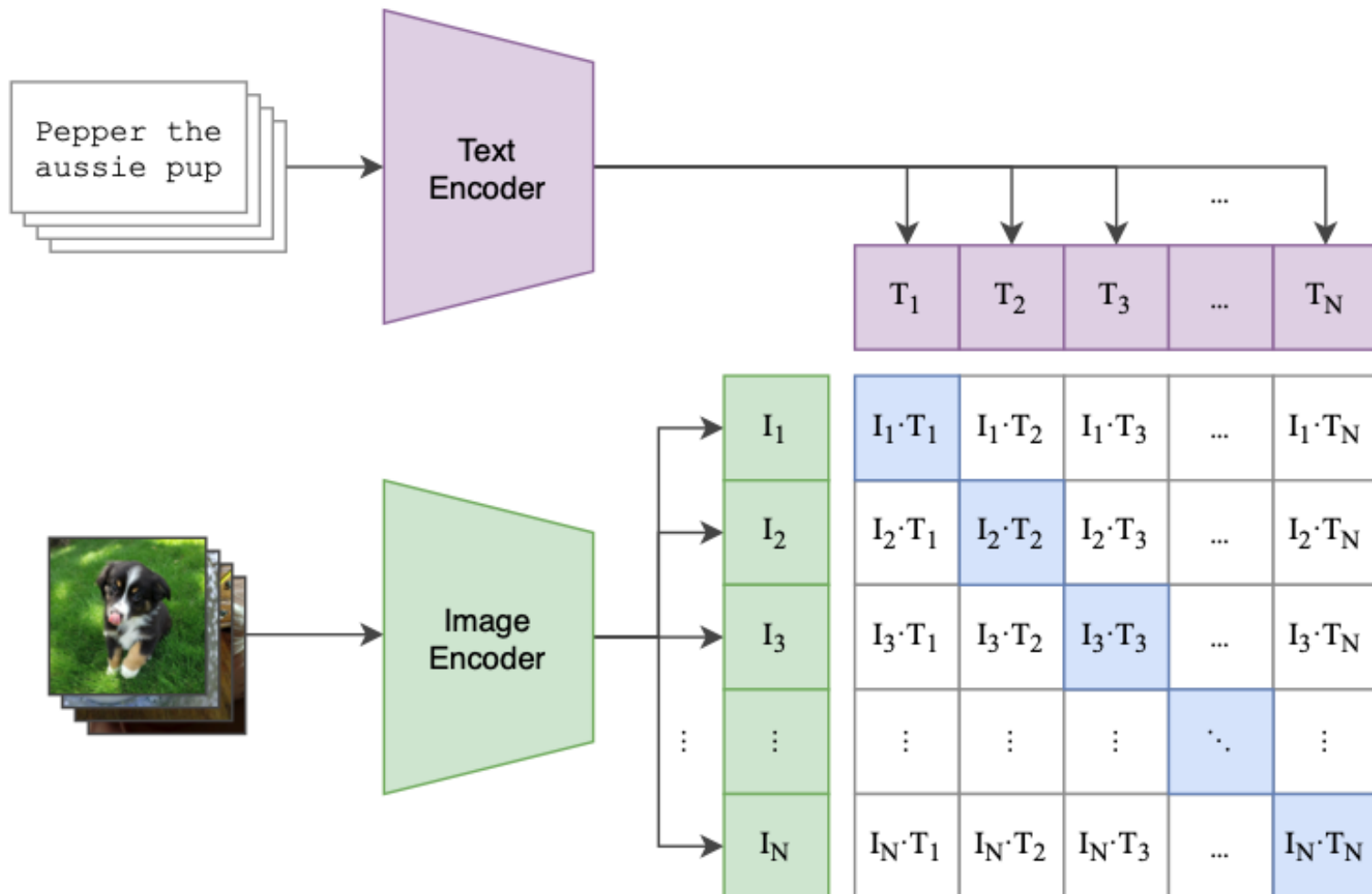




# Vision and Language

# Case study: CLIP

## (1) Contrastive pre-training



# unCLIP aka DALL-E 2

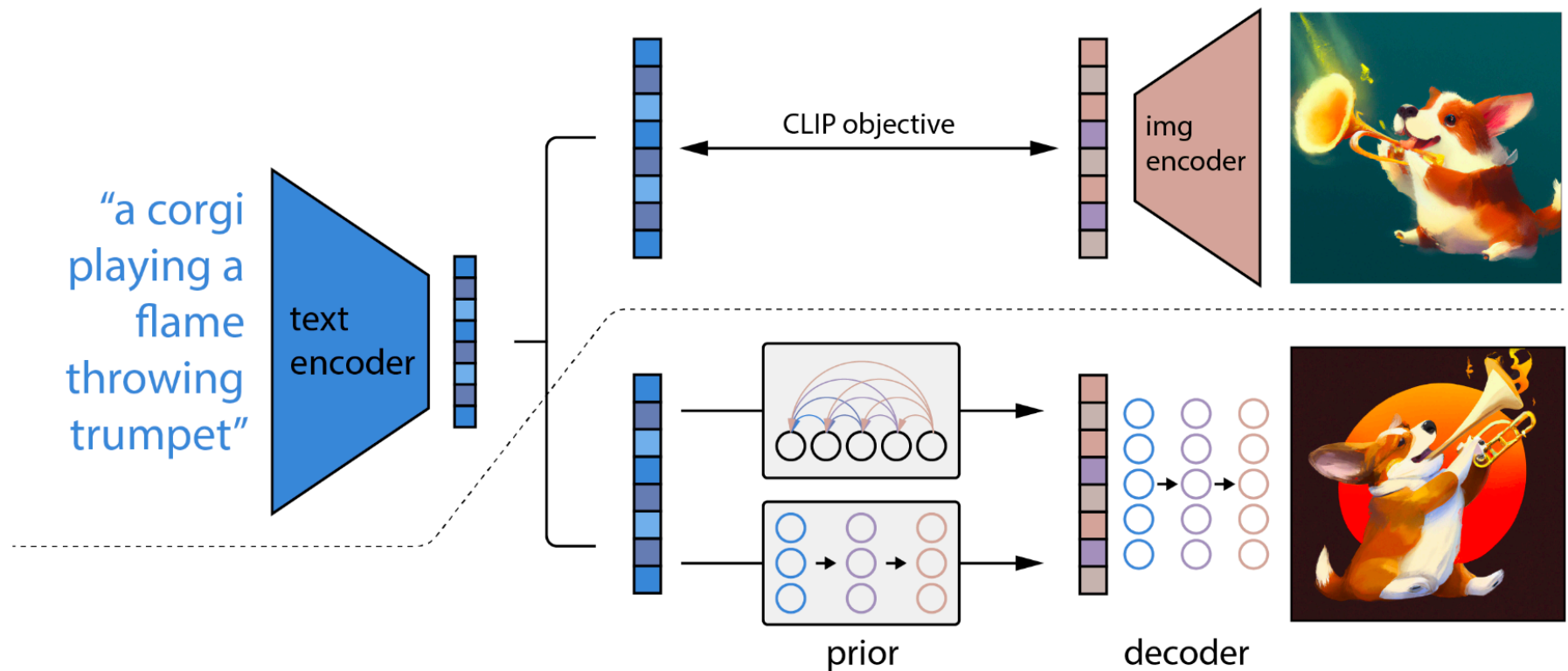


Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive 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)