



CTU

CZECH TECHNICAL
UNIVERSITY
IN PRAGUE

Deep Learning Essentials

9. Task-specific Architectures

Semantic segmentation, Object detection, generative models, ...

Lukáš Neumann

Adapted from [B3B33UROB](#) slides of Karel Zimmerman

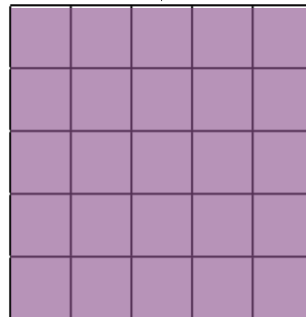
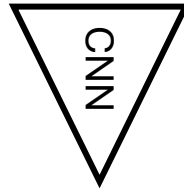
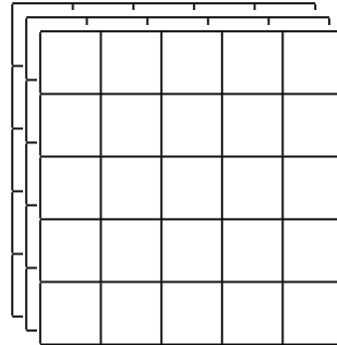
Semantic segmentation



- road
- sidewalk
- pedestrian
- traffic sign
- tree
- sky

Semantic segmentation

RGB image
(HxWx3)



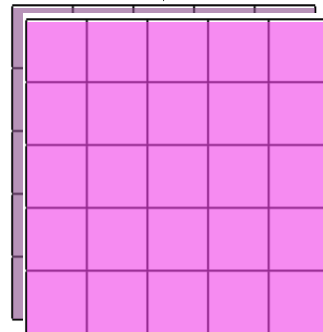
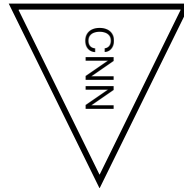
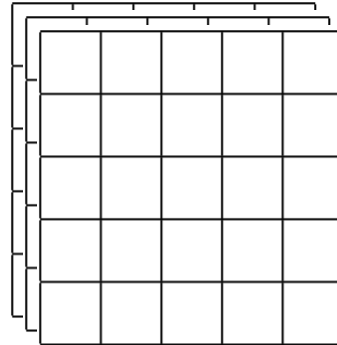
- road
- sidewalk
- pedestrian
- traffic sign
- tree
- sky

pixel-wise probability
of being **road**

channel 1

Semantic segmentation

RGB image
(HxWx3)



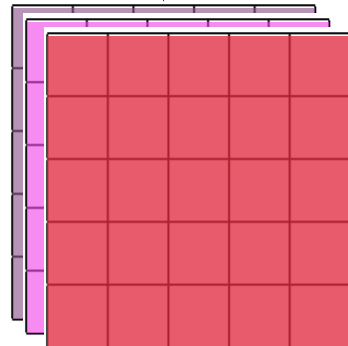
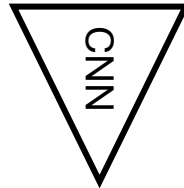
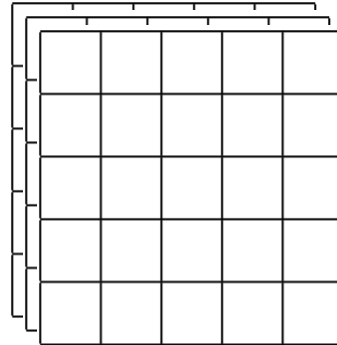
- road
- sidewalk
- pedestrian
- traffic sign
- tree
- sky

pixel-wise probability
of being **sidewalk**

channel 2

Semantic segmentation

RGB image
(HxWx3)

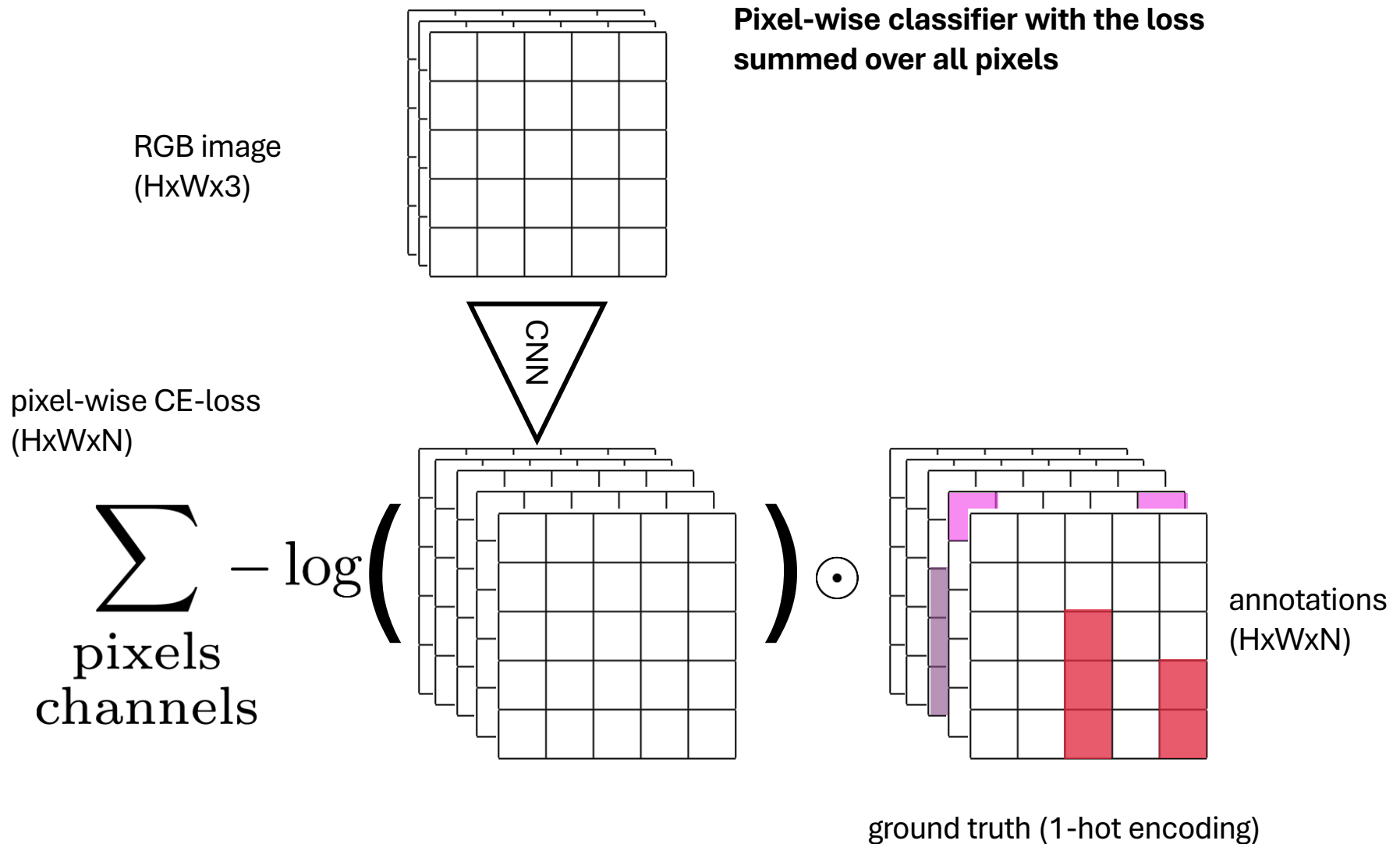


pixel-wise probability
of being **pedestrian**

channel 3

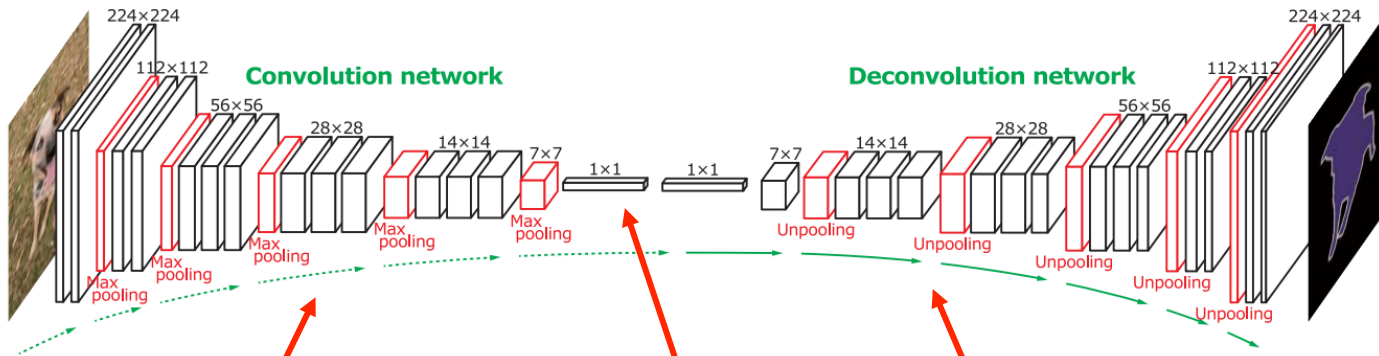
-  road
-  sidewalk
-  pedestrian
-  traffic sign
-  tree
-  sky

Semantic segmentation



Semantic segmentation

- UNet architecture



Convolution layers:

- decrease spatial resolution
- increase number of channels

Deconvolution layers

- Increase spatial resolution
- Decrease number of channels to K classes

Spatial context encoded in channels

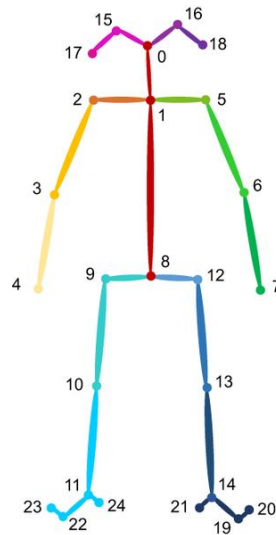
Pose estimation

- Estimate position of human body joints in 2D (3D)

input



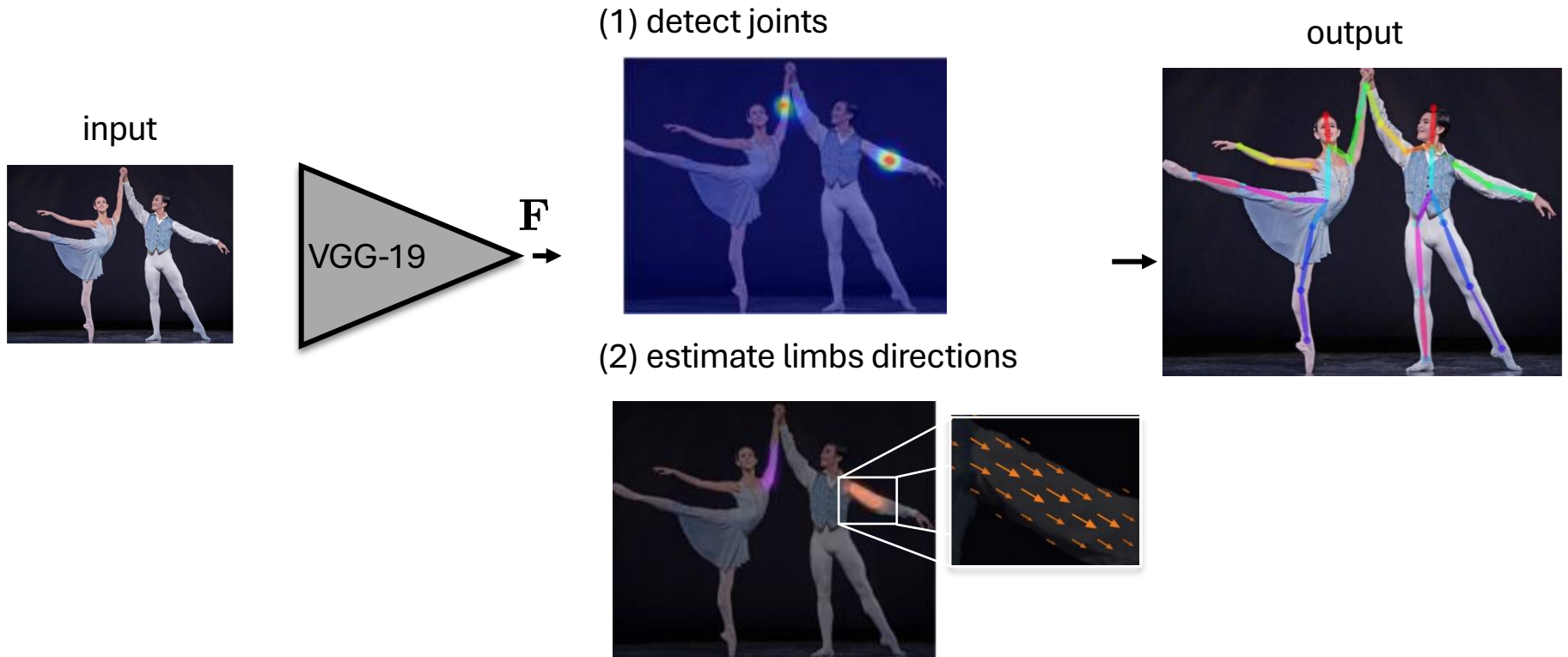
output



- | | |
|--------------------|----------------------|
| 0 - Nose | 13 - Left knee |
| 1 - Neck | 14 - Left ankle |
| 2 - Right shoulder | 15 - Right eye |
| 3 - Right elbow | 16 - Left eye |
| 4 - Right wrist | 17 - Right ear |
| 5 - Left shoulder | 18 - Left ear |
| 6 - Left elbow | 19 - Left big-toe |
| 7 - Left wrist | 20 - Left small-toe |
| 8 - Mid hip | 21 - Left heel |
| 9 - Right hip | 22 - Right big-toe |
| 10 - Right knee | 23 - Right small-toe |
| 11 - Right ankle | 24 - Right heel |
| 12 - Left hip | |

Pose estimation

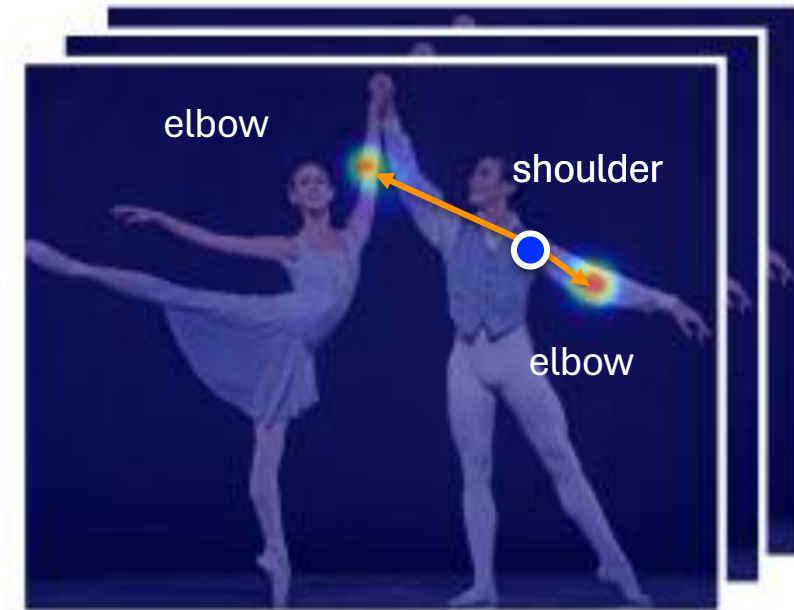
- OpenPose



Pose estimation

- There are inherent ambiguities how joint positions can be connected together

joints



PAFs

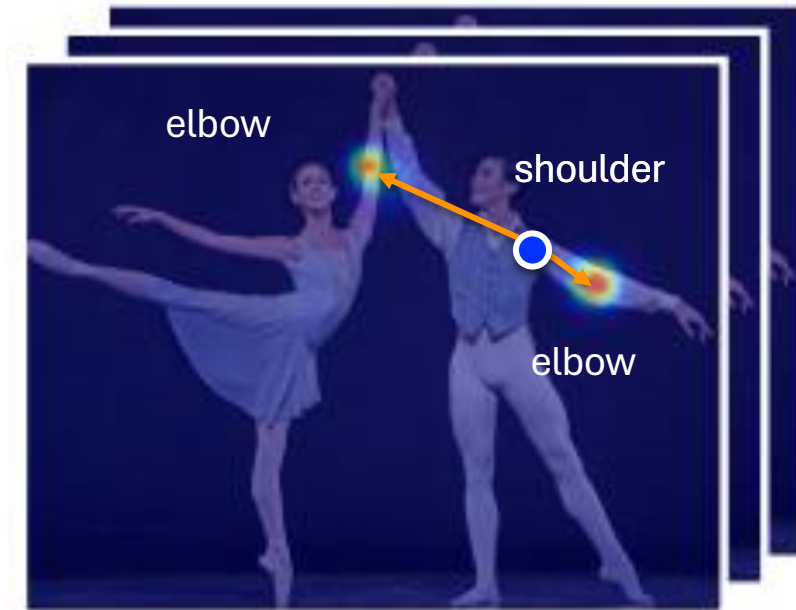


Pose estimation

- There are inherent ambiguities how joint positions can be connected together

joints

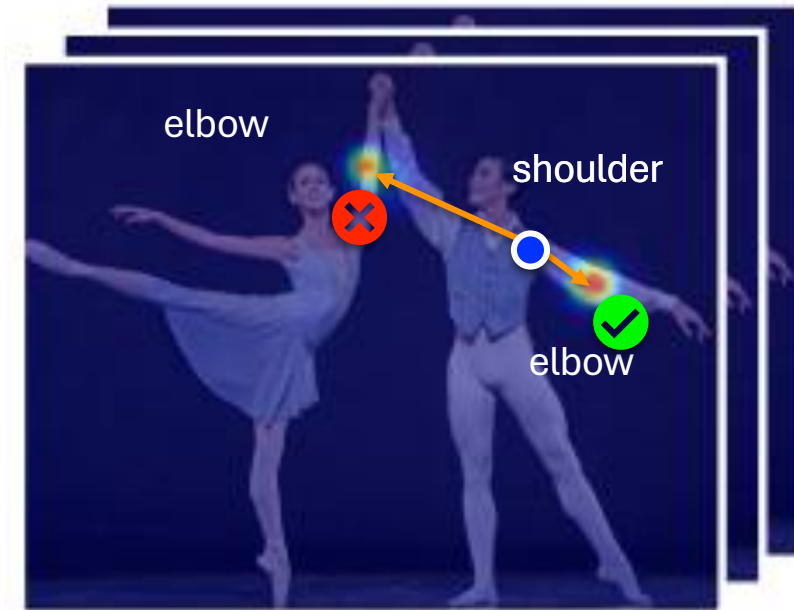
PAFs



Pose estimation

- Point Affinity Fields (PAFs) allow to correctly connect joints

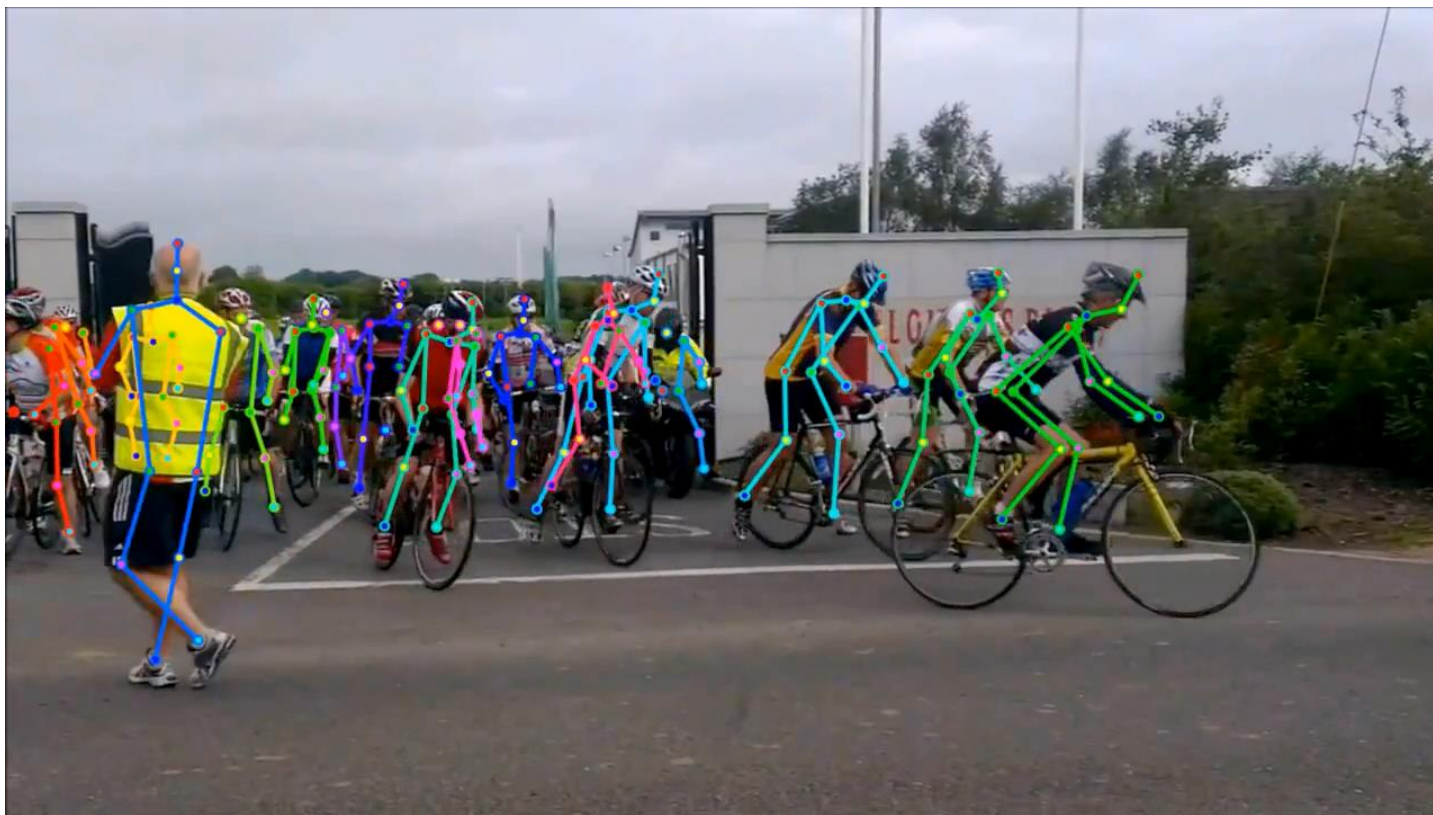
joints



PAFs



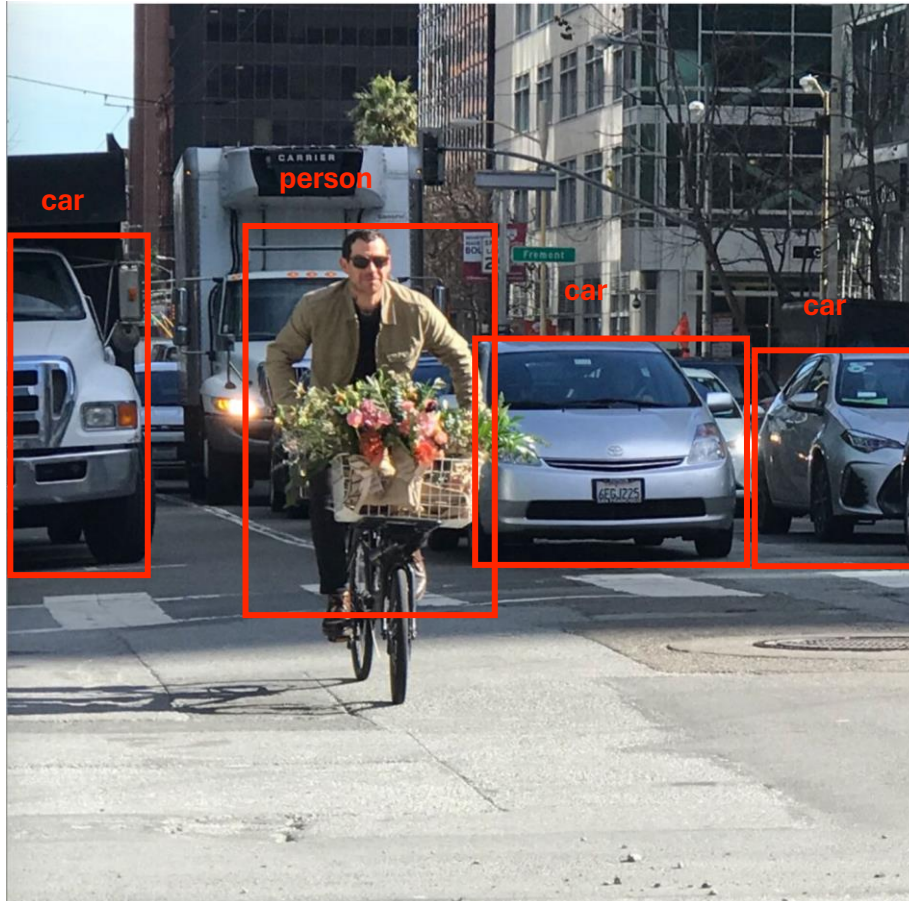
Pose estimation



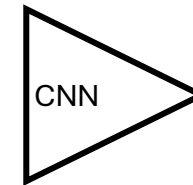
Object detection



Object detection

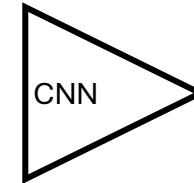


Object detection



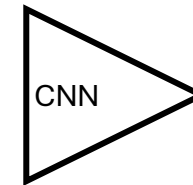
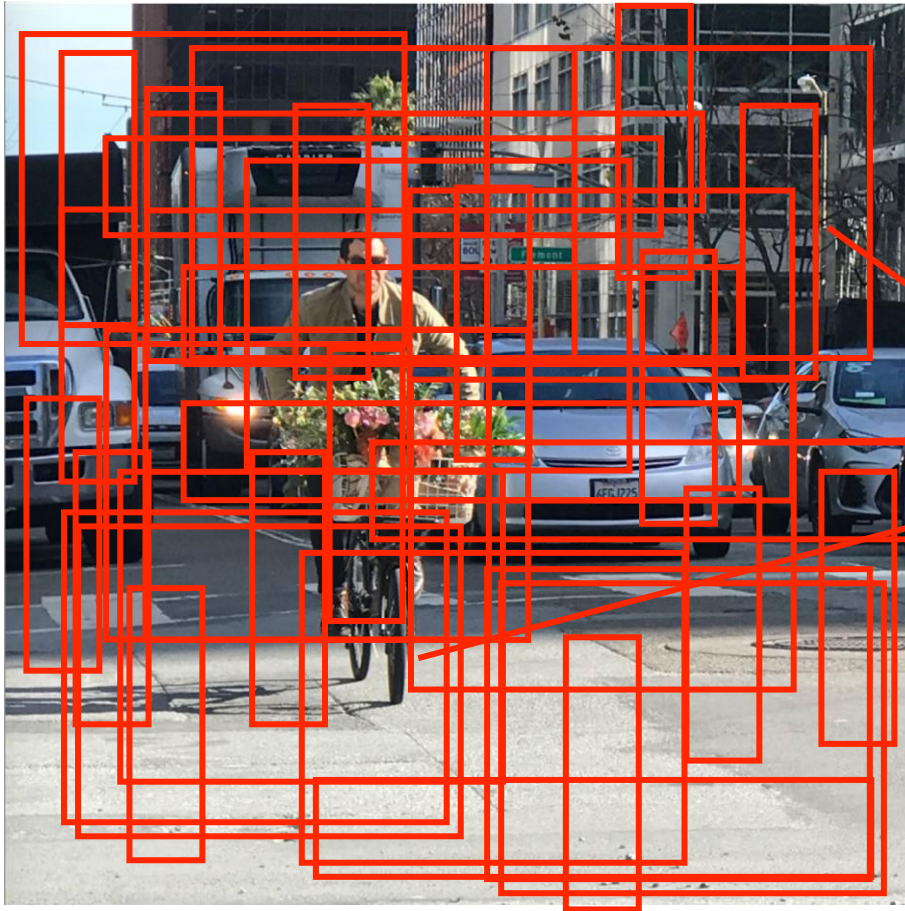
0.6	person
0.1	car
0.2	building
0.1	background

Object detection



0.1	person
0.1	car
0.2	building
0.6	background

Object detection



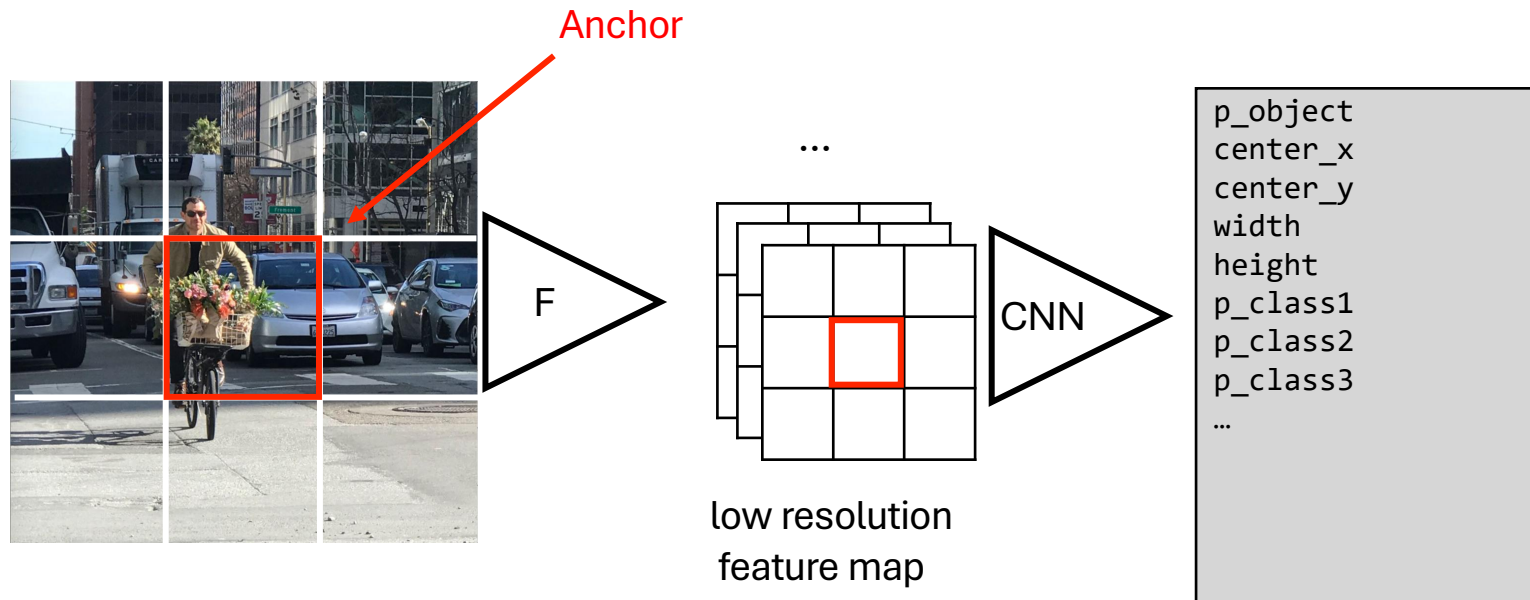
0.1	person
0.1	car
0.2	building
0.6	background

Classify all rectangles?

$H \times W \times \text{AspectRatio} \times \text{Scales} \times 0.001 \text{ sec} = \text{months}$

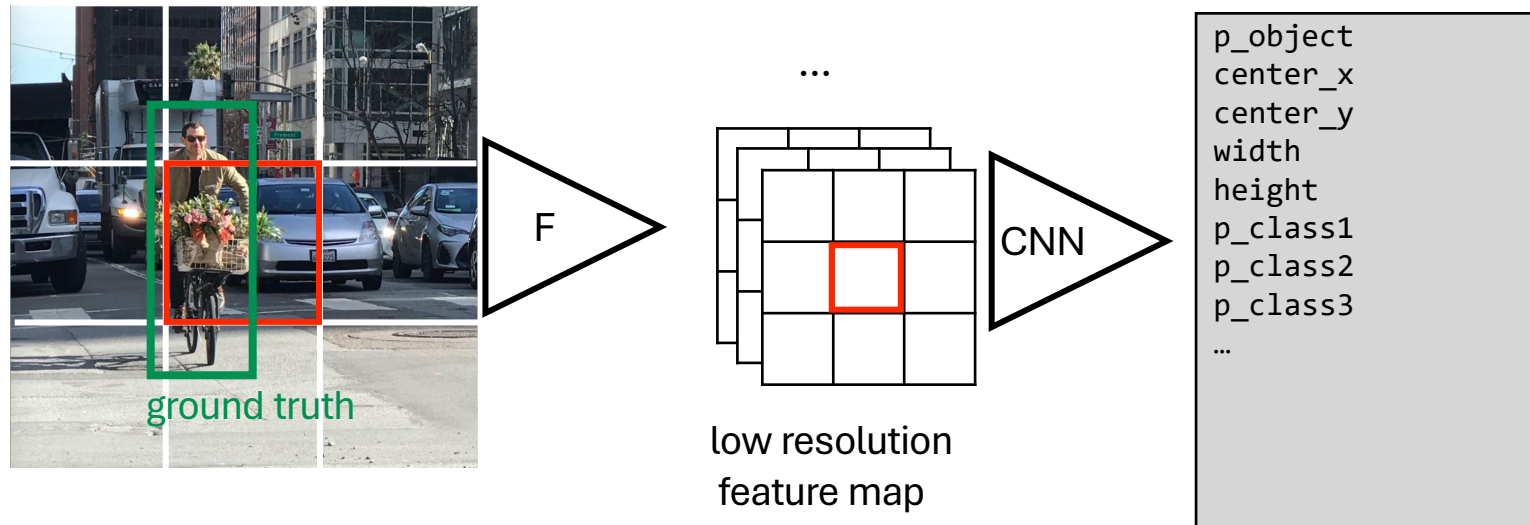
Object detection

- YOLO (You Only Look Once)
- Divide image into $M \times M$ sub images (corresponding to its receptive fields)
- Predict relative position, objectness, class for each patch



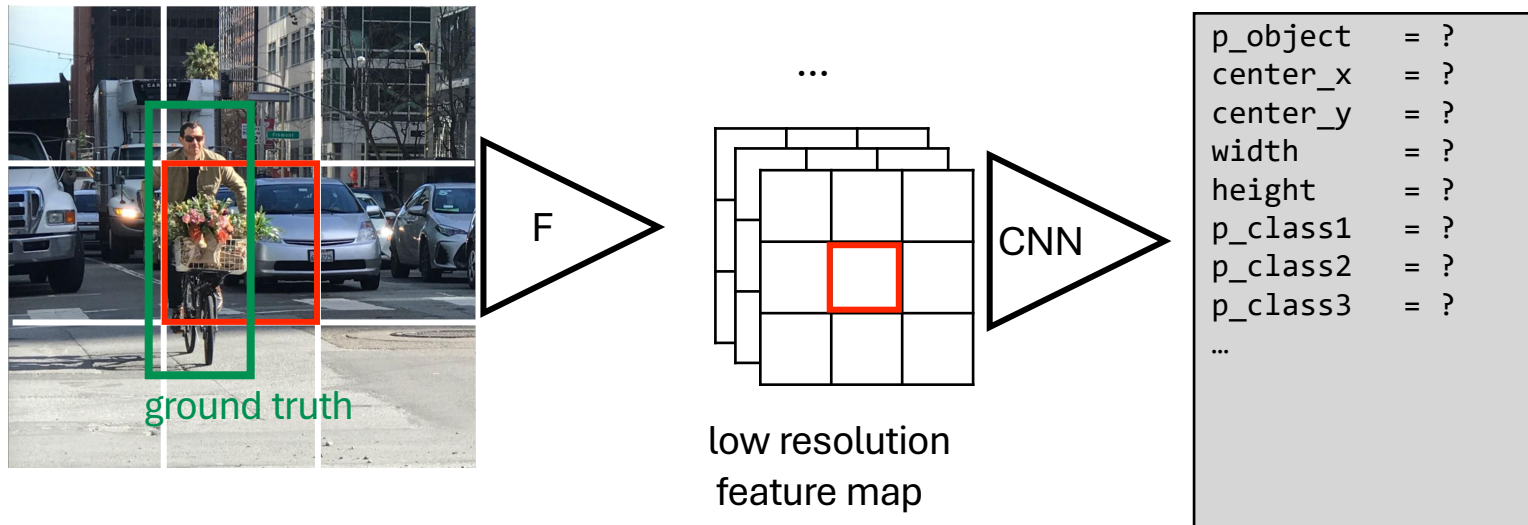
Object detection

- YOLO (You Only Look Once)
- Divide image into $M \times M$ sub images (corresponding to its receptive fields)
- Predict relative position, objectness, class for each patch



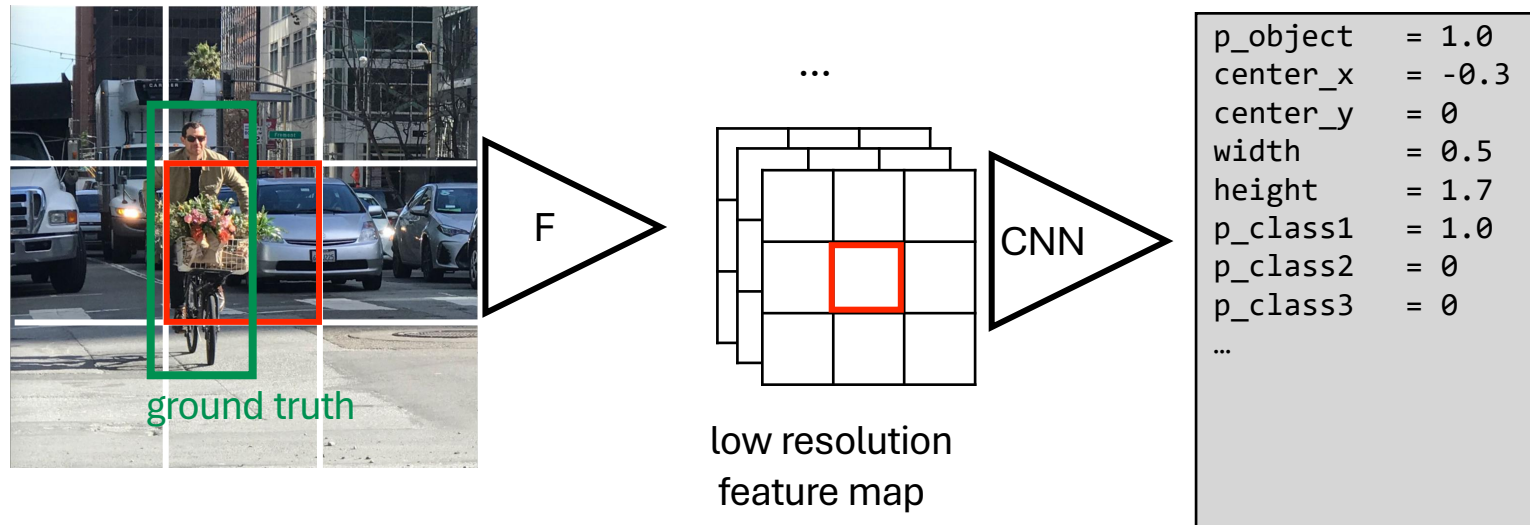
Object detection

- YOLO (You Only Look Once)
- Divide image into $M \times M$ sub images (corresponding to its receptive fields)
- Predict relative position, objectness, class for each patch



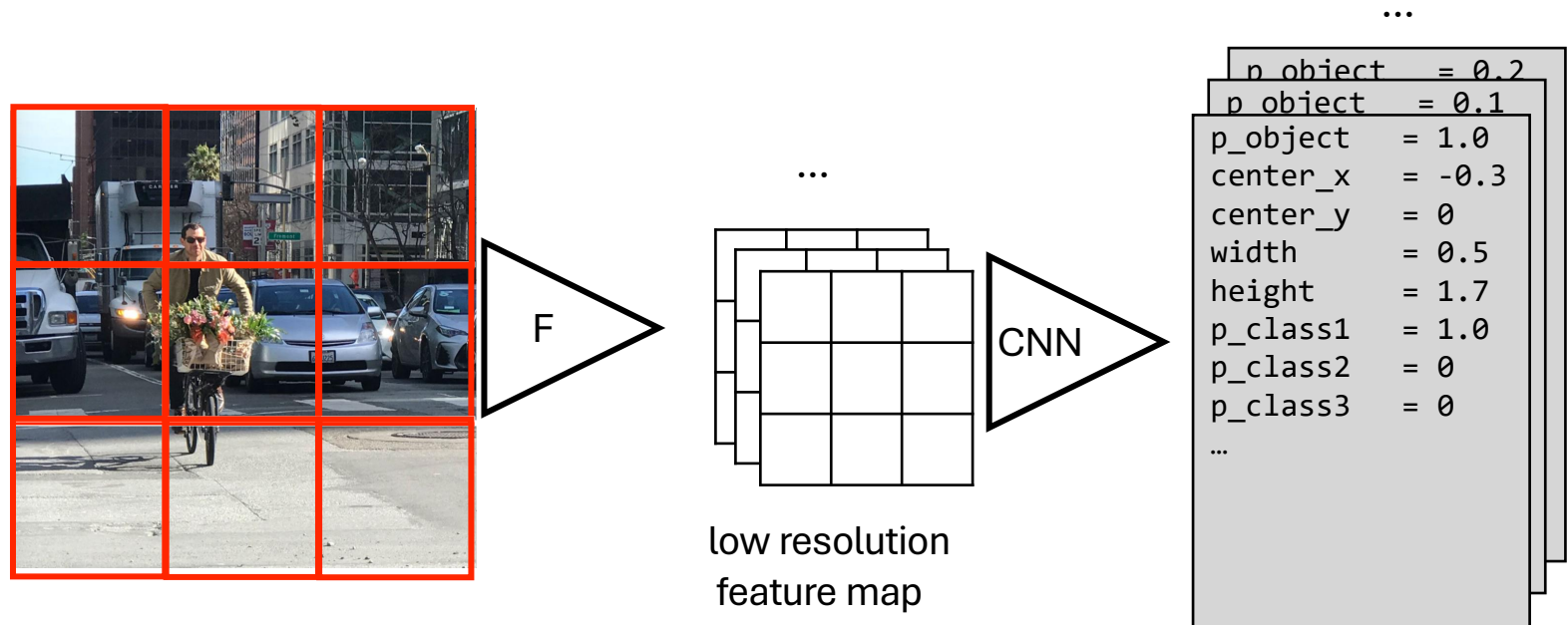
Object detection

- YOLO (You Only Look Once)
- Divide image into $M \times M$ sub images (corresponding to its receptive fields)
- Predict relative position, objectness, class for each patch



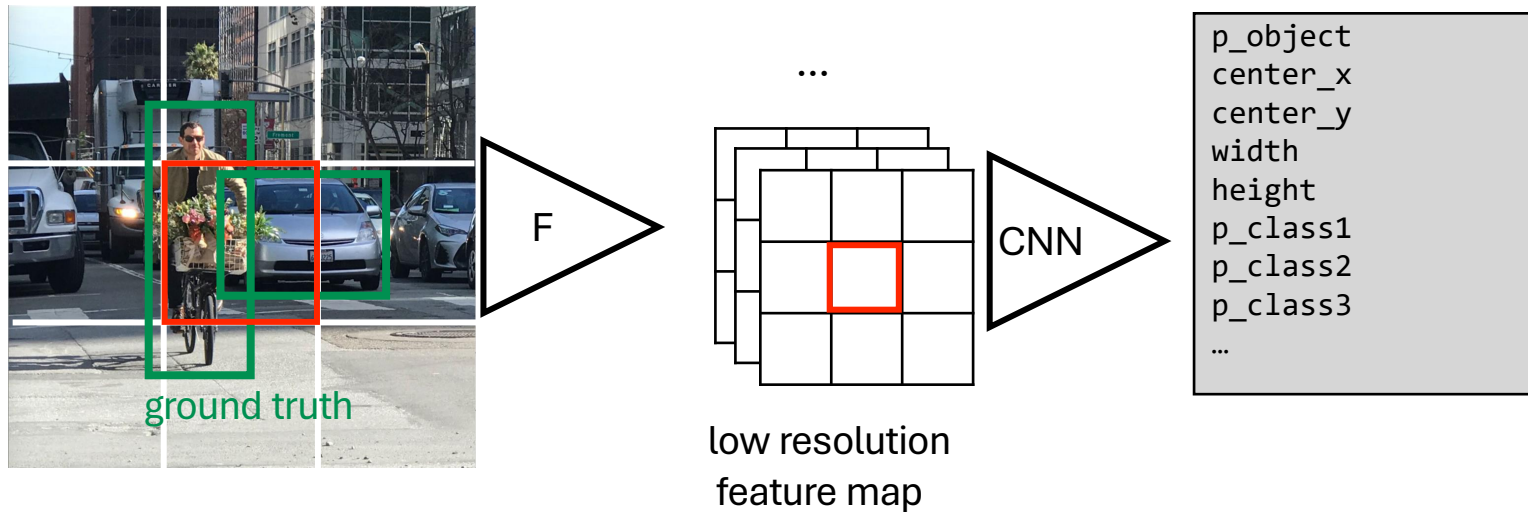
Object detection

- Each image patch has its own set of outputs



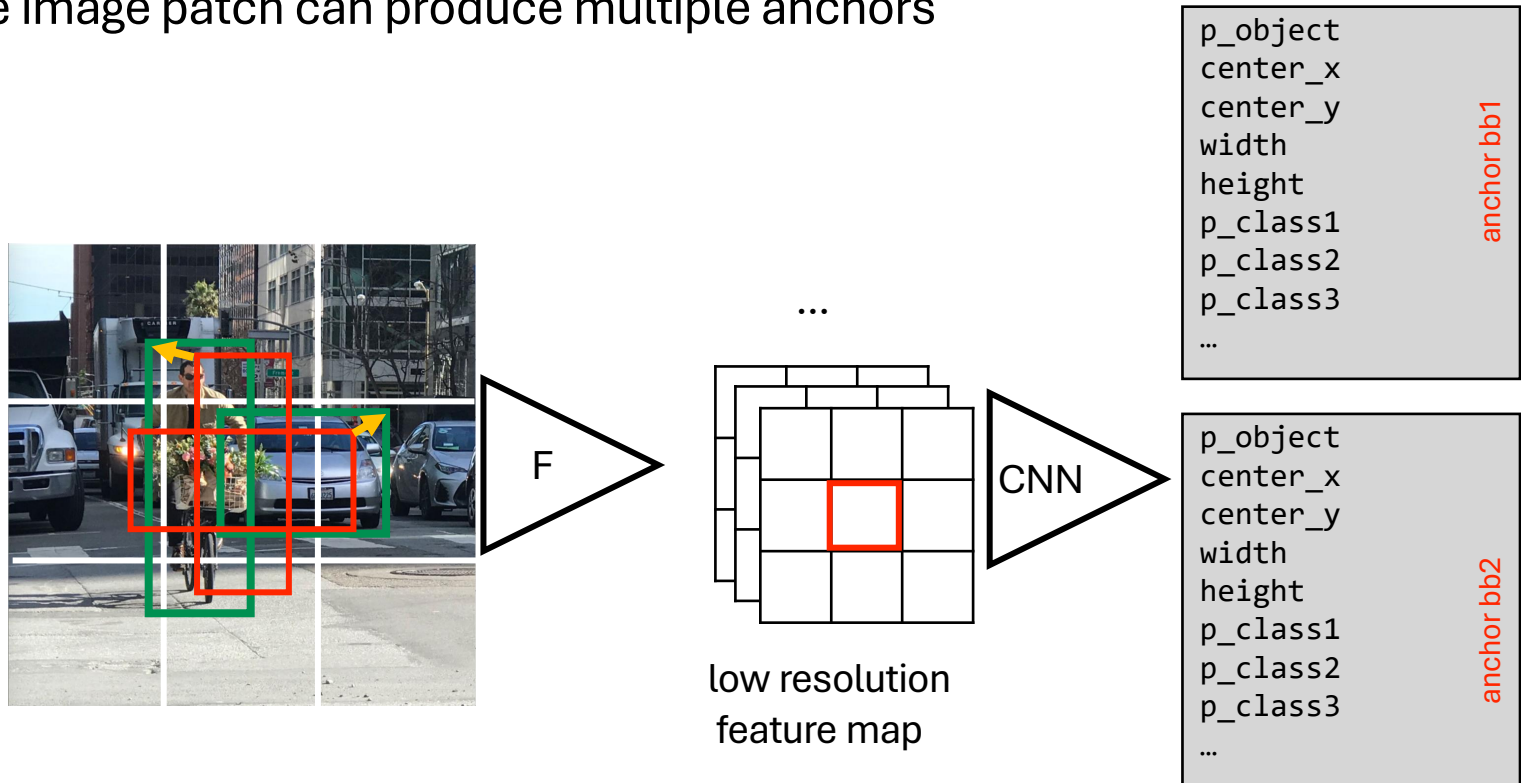
Object detection

- Do you see a problem?
- Multiple objects in a single patch!



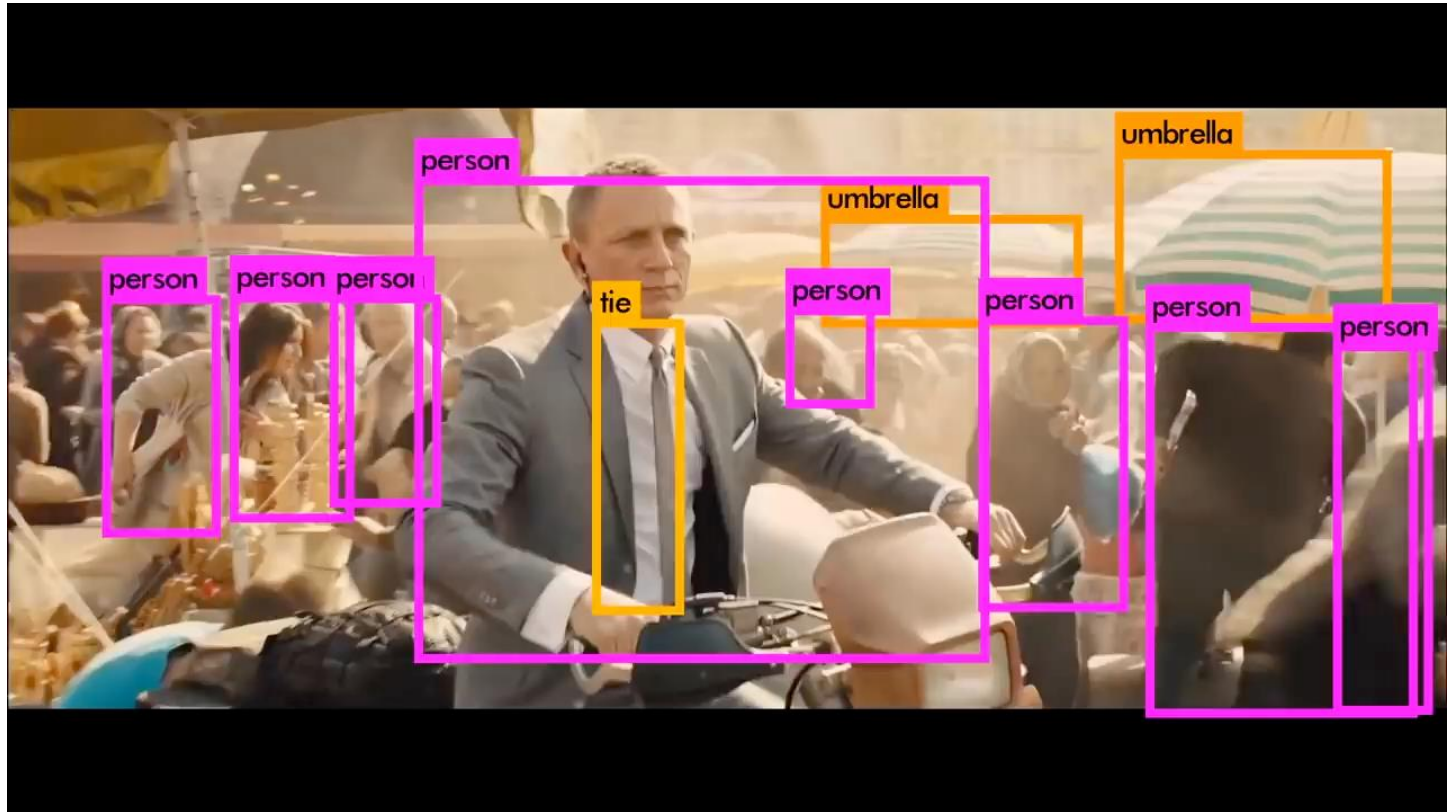
Object detection

- Multiple anchors with different aspect ratios
- One image patch can produce multiple anchors



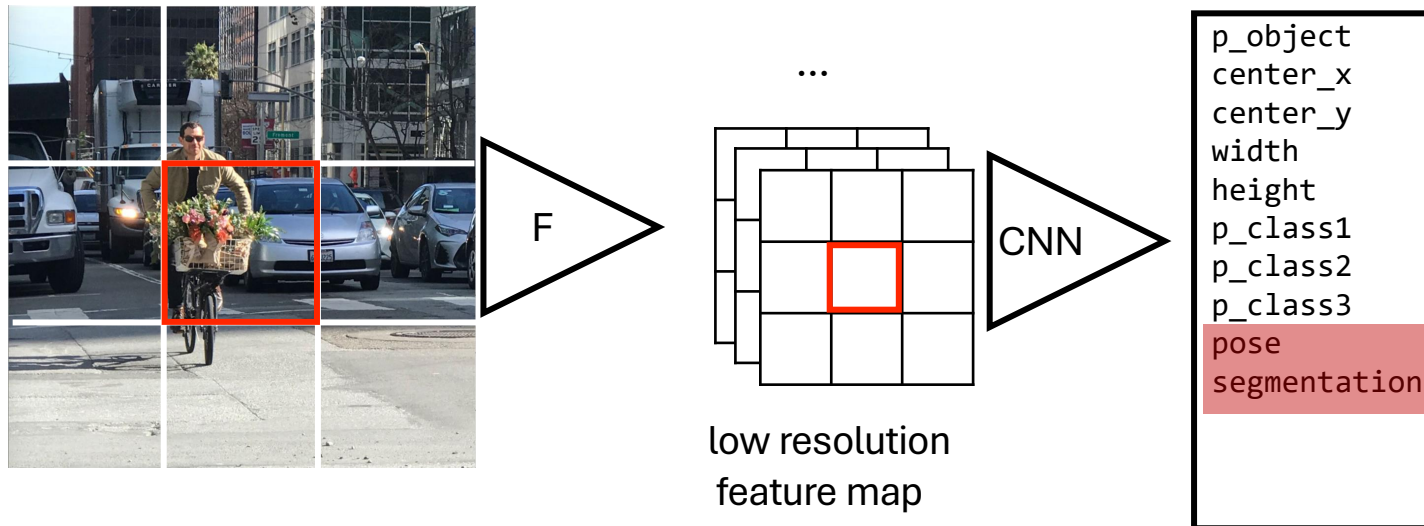
Object detection

- YOLOv2



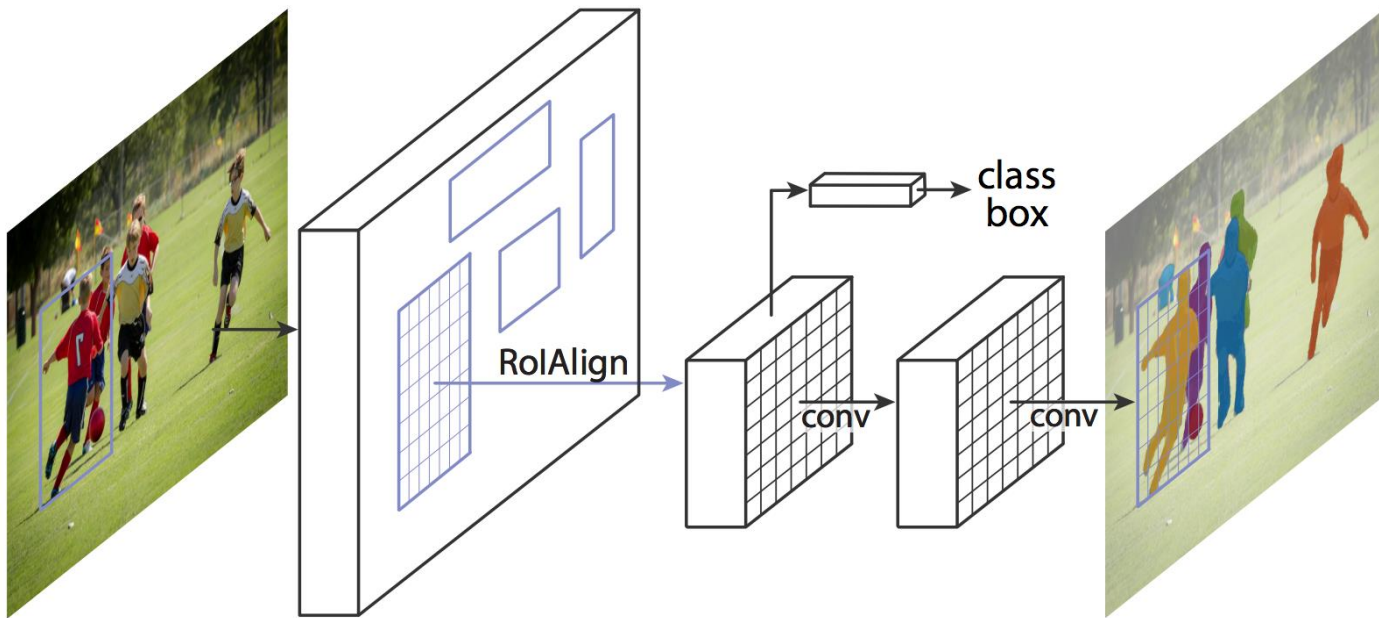
Instance segmentation

- Mask R-CNN
- Added pose and segmentation head for each detected object



Instance segmentation

- Mask R-CNN
- Added pose and segmentation head for each detected object



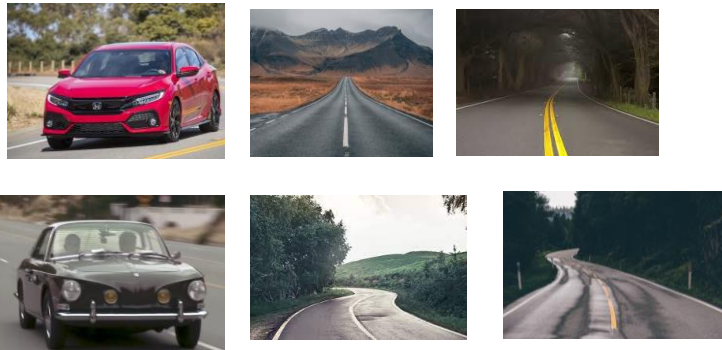
Instance segmentation

- Mask R-CNN
- Added pose and segmentation head for each detected object



Object detector evaluation

- Two classes – positive (object is present) & negative (object is not present)
- Example - Car detector
 - Car is present (positive class)
 - Car is not present = background only (negative class)



Object detector evaluation

Positive class



Negative class



Classifier
output:

Object detector evaluation

Positive class

Negative class

Classifier
output:



true positive (TP)

... classifier correctly found the object (e.g. car) = the **positive** class

Object detector evaluation

Classifier
output:

Positive class



Negative class



true positive (TP)

... classifier correctly found the object (e.g. car) = the **positive** class

true negative (TN)

... classifier correctly indicated the object is missing = the **negative** class

Object detector evaluation

Classifier
output:

Positive class



Negative class



true positive (TP)

... classifier correctly found the object (e.g. car) = the **positive** class

true negative (TN)

... classifier correctly indicated the object is missing = the **negative** class

false negative (FN)

... classifier **falsely** indicates **negative** class in presence of the object
→ **missed danger**

Object detector evaluation

Classifier
output:

Positive class



Negative class



- | | |
|---------------------|---|
| true positive (TP) | ... classifier correctly found the object (e.g. car) = the positive class |
| true negative (TN) | ... classifier correctly indicated the object is missing = the negative class |
| false negative (FN) | ... classifier falsely indicates negative class in presence of the object
→ missed danger |
| false positive (FP) | ... classifier falsely indicates positive class where the object is missing
→ false alarm |

Object detector evaluation

Classifier
output:

Positive class



Negative class



true positive (TP) = 1

true negative (TN) = 2

false negative (FN) = 1

false positive (FP) = 2

"1/3 of of samples classified as CARS are actually CARS"

$$\text{Precision (P)} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{1}{1 + 2} = 1/3$$

$$\text{Recall (R)} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{1}{1 + 1} = 1/2$$

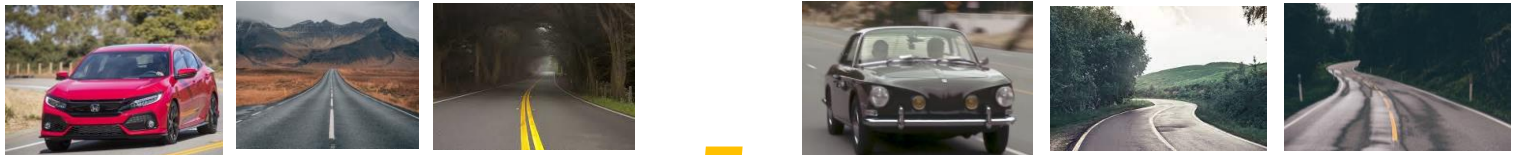
"1/2 of all CARS are discovered."

What is the best possible outcome? Oracle: Precision = Recall = 1

Object detector evaluation

Positive class

Negative class



Classifier
score:

0.9

0.5

0.1

-0.1

-0.4

-0.6

Threshold = 0

true positive (TP) = 1

true negative (TN) = 2

false negative (FN) = 1

false positive (FP) = 2

$$\text{Precision (P)} = \frac{TP}{TP + FP} = \frac{1}{1 + 2} = 1/3$$

$$\text{Recall (R)} = \frac{TP}{TP + FN} = \frac{1}{1 + 1} = 1/2$$

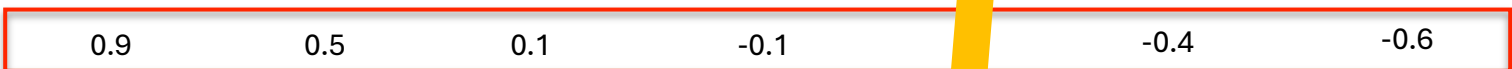
Object detector evaluation

Positive class

Negative class



Classifier
score:



Threshold = -0.2

true positive (TP) = 2

true negative (TN) = 2

false negative (FN) = 0

false positive (FP) = 2

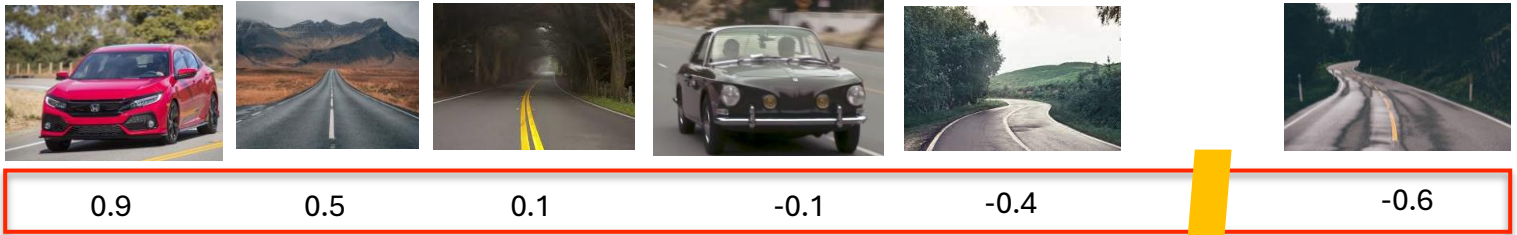
$$\text{Precision (P)} = \frac{\text{TP}}{\text{TP} + \text{FP}} = \frac{2}{2 + 2} = 1/2$$

$$\text{Recall (R)} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \frac{2}{2 + 0} = 1$$

Object detector evaluation

Positive class

Negative class



Classifier
score:

Threshold = -0.5

true positive (TP) = 2

true negative (TN) = 1

false negative (FN) = 0

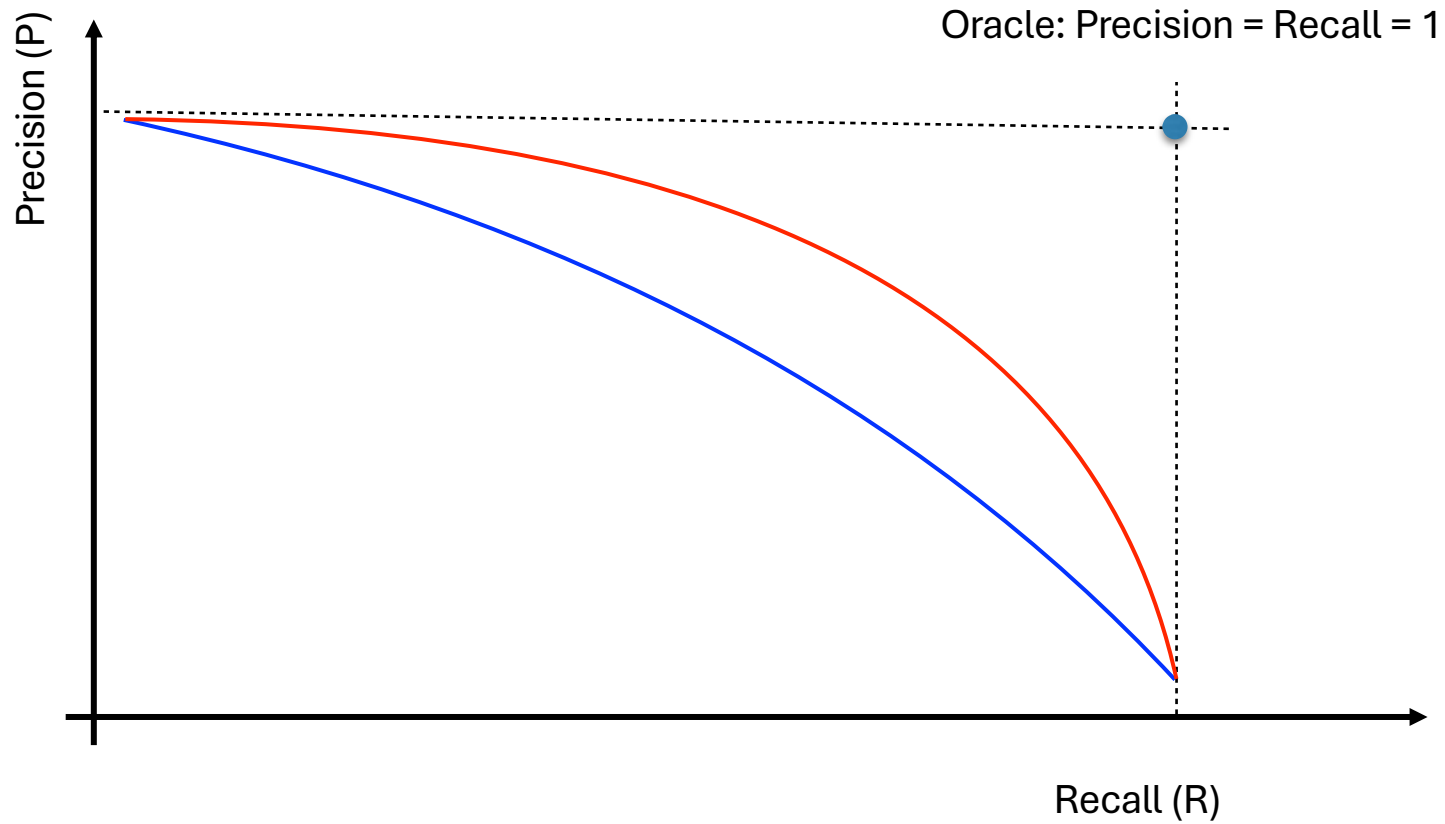
false positive (FP) = 3

$$\text{Precision (P)} = \frac{TP}{TP + FP} = \frac{2}{2 + 3} = 2/5$$

$$\text{Recall (R)} = \frac{TP}{TP + FN} = \frac{2}{2 + 0} = 1$$

Object detector evaluation

- PR-curve



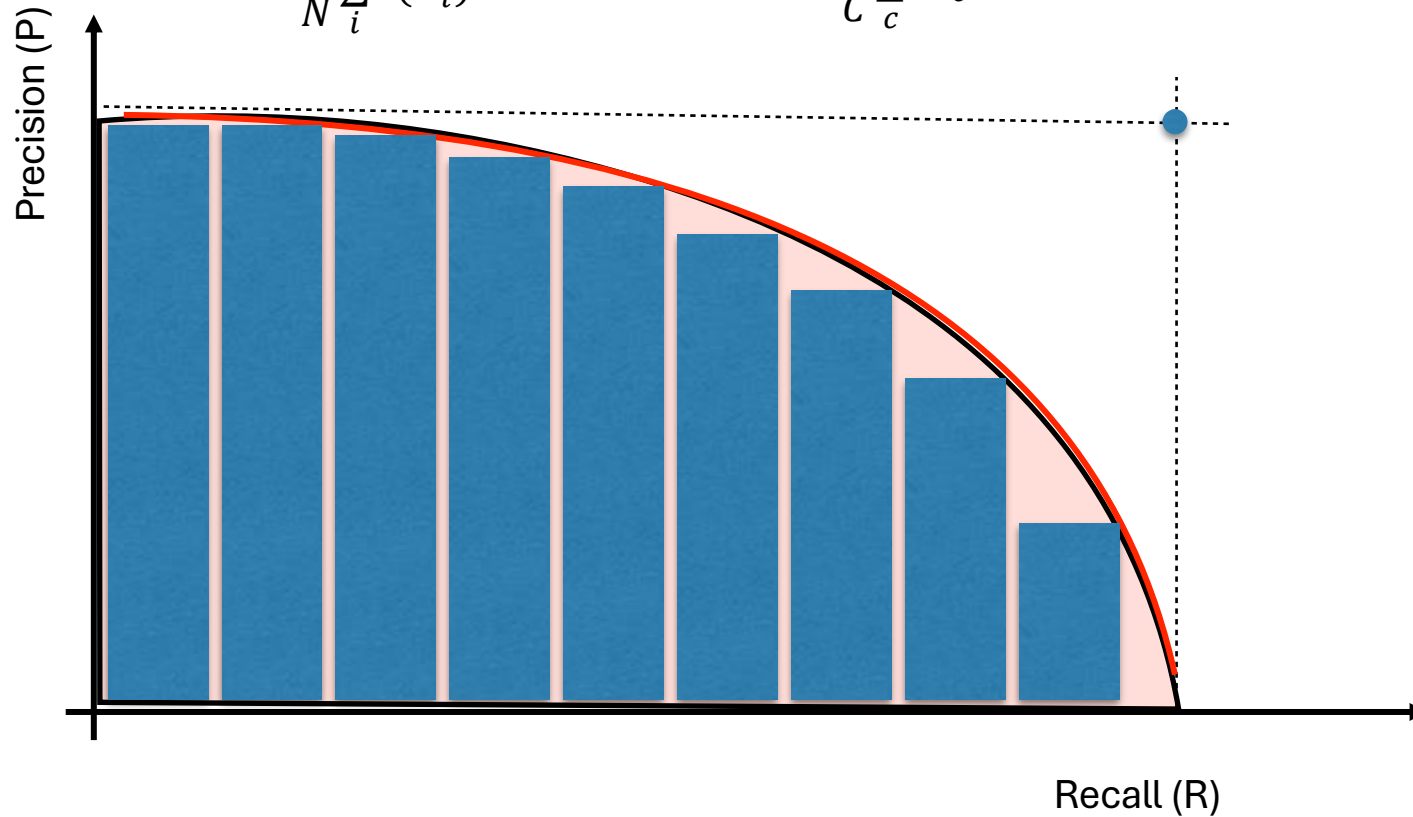
Object detector evaluation

Average Precision

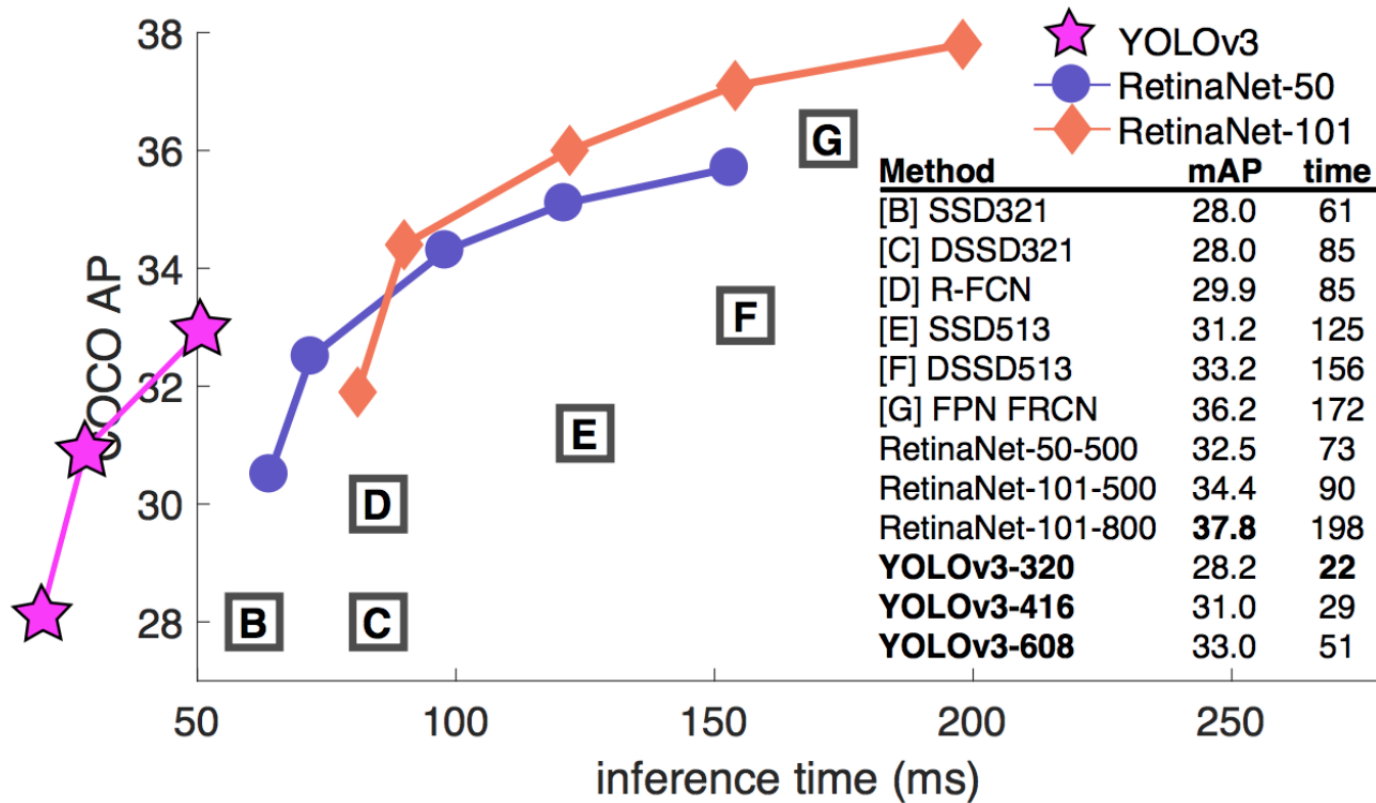
$$AP = \frac{1}{N} \sum_i P(R_i)$$

mean Average Precision

$$mAP = \frac{1}{C} \sum_c AP_c$$



Object detector evaluation



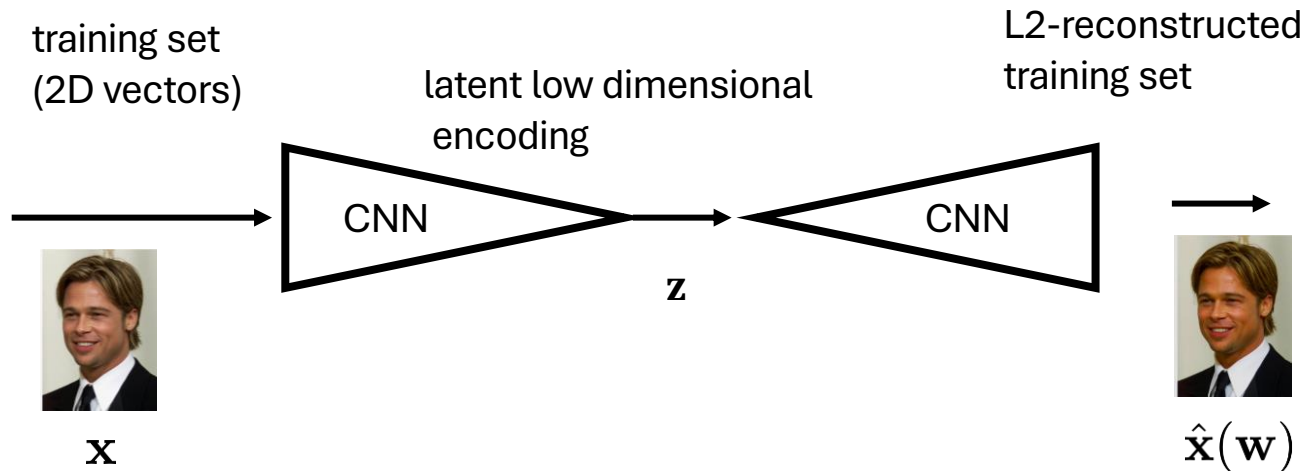
Generative models

- Variational Auto-Encoders (VAEs)
- Generative Adversarial Networks (GANs)
- Diffusion models

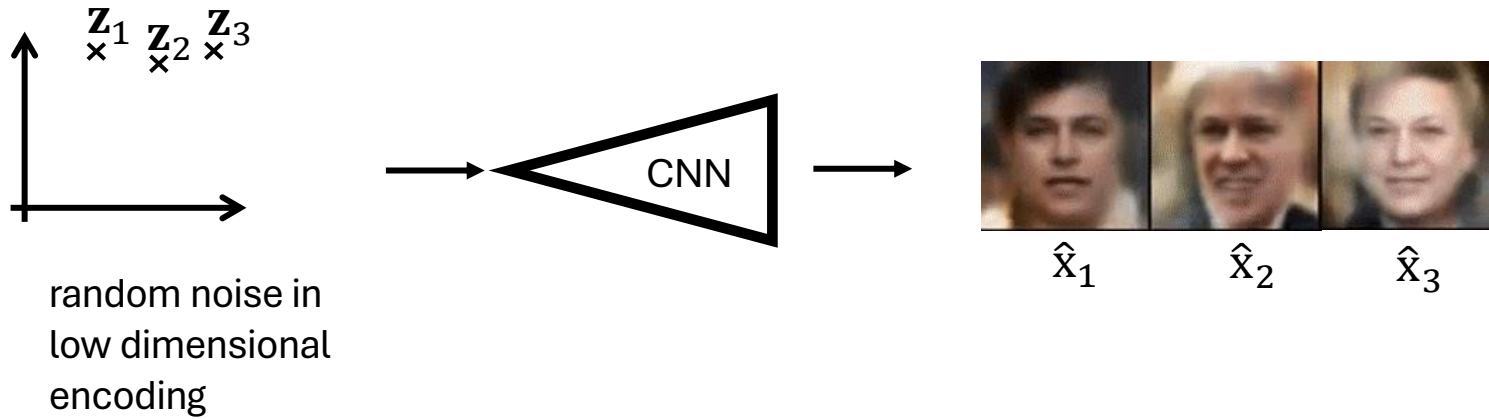
Variational Auto-Encoders (VAEs)

- Latent space \mathbf{z} is pushed towards a Gaussian distribution
- Learning the self-reconstruction with L2 reconstruction loss

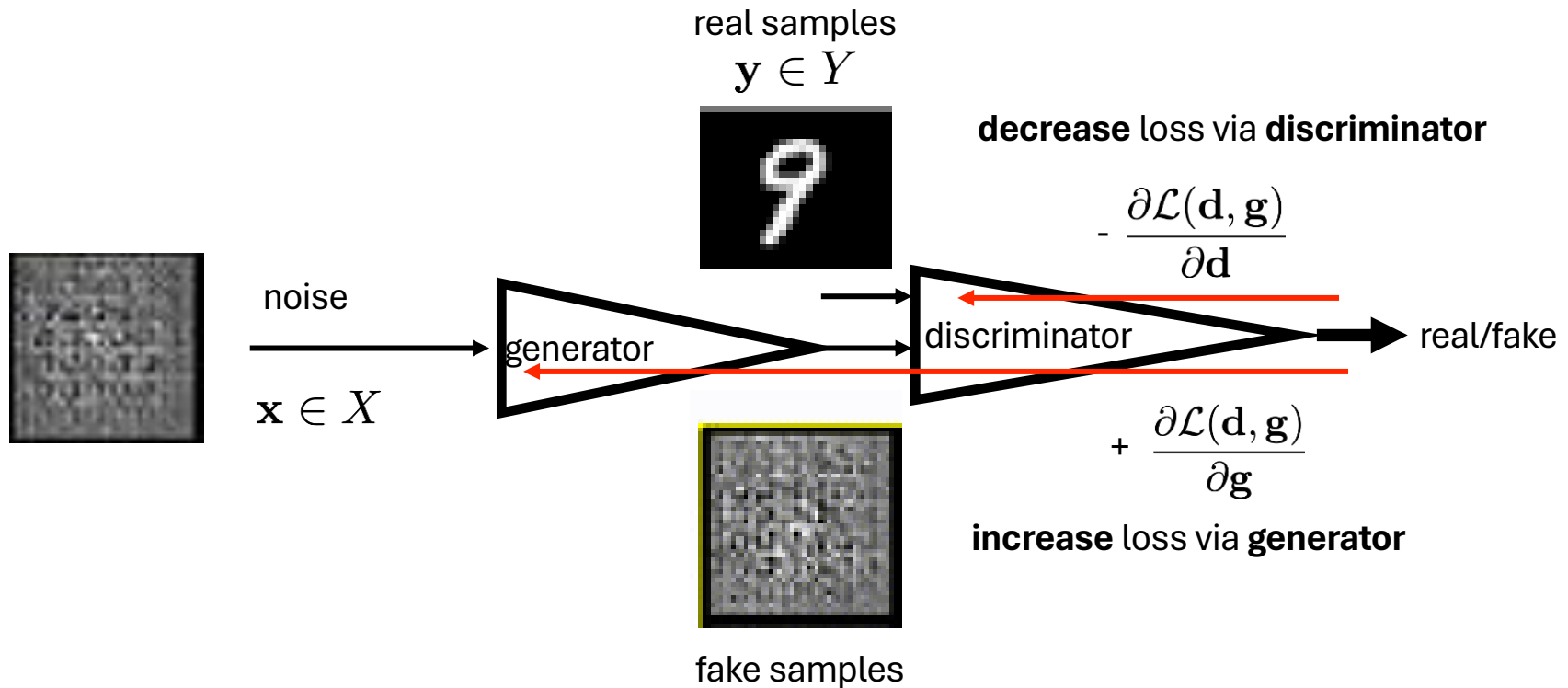
$$\arg \min_{\mathbf{w}} \|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{w})\|_2^2$$



Variational Auto-Encoders (VAEs)



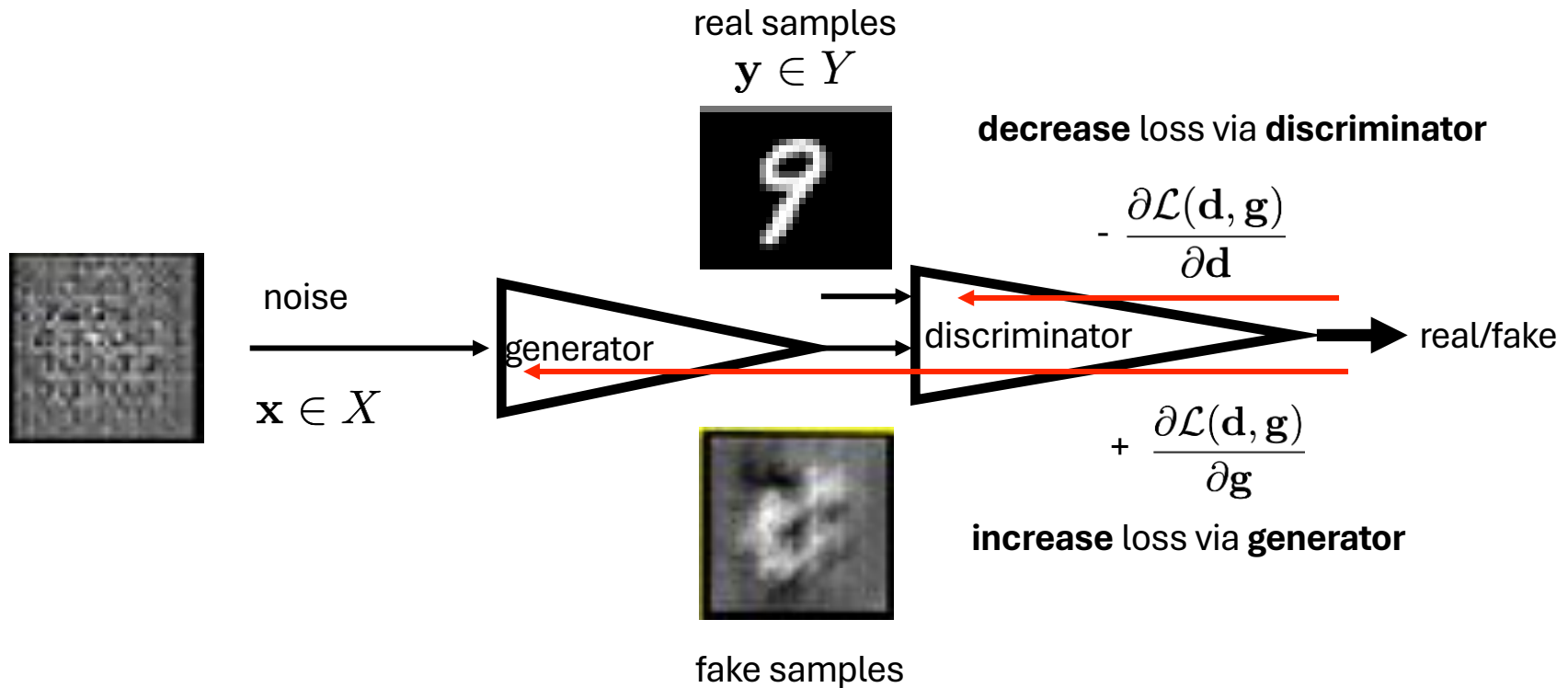
Generative Adversarial Networks (GANs)



Classification loss:

$$\mathcal{L}(d, g) = \sum_{x \in X} -\log(d(g(x))) + \sum_{y \in Y} -\log(1 - d(y))$$

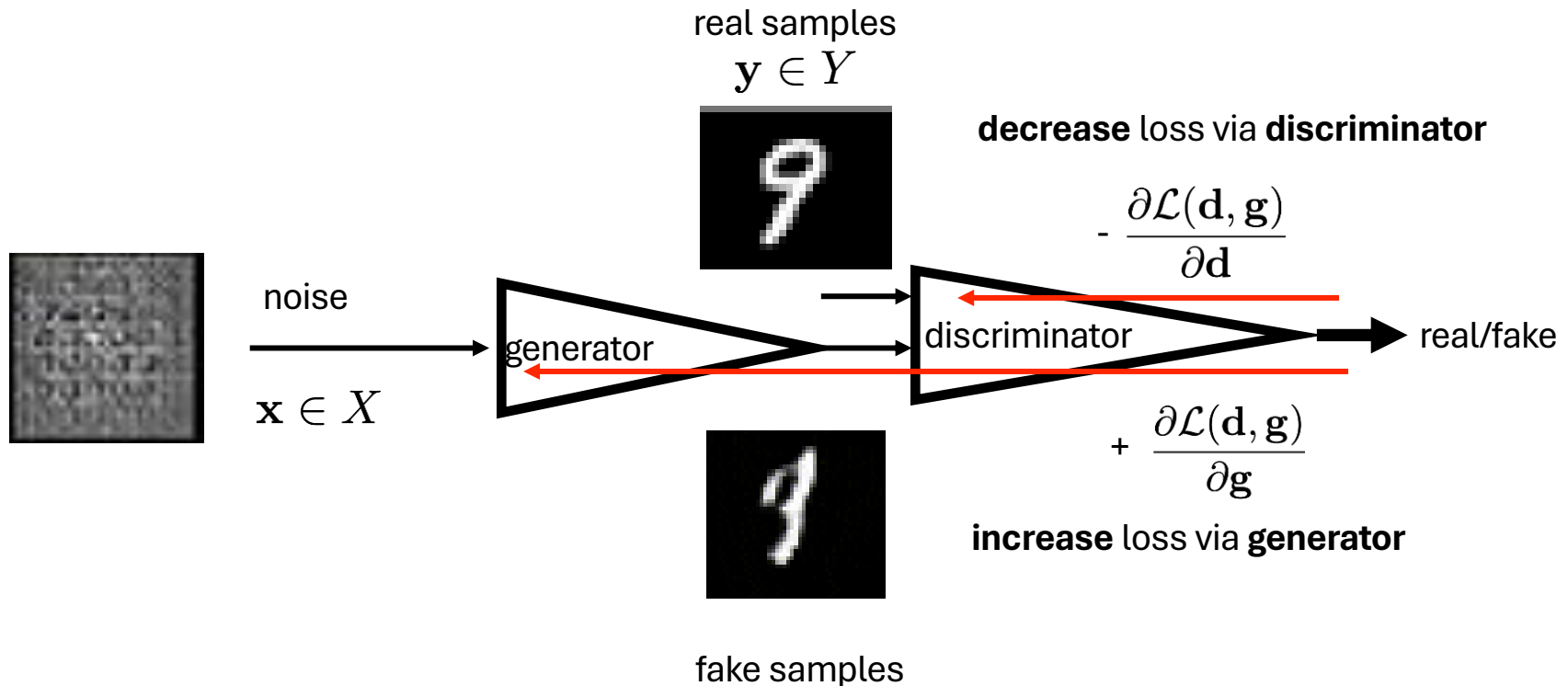
Generative Adversarial Networks (GANs)



Classification loss:

$$\mathcal{L}(d, g) = \sum_{x \in X} -\log(d(g(x))) + \sum_{y \in Y} -\log(1 - d(y))$$

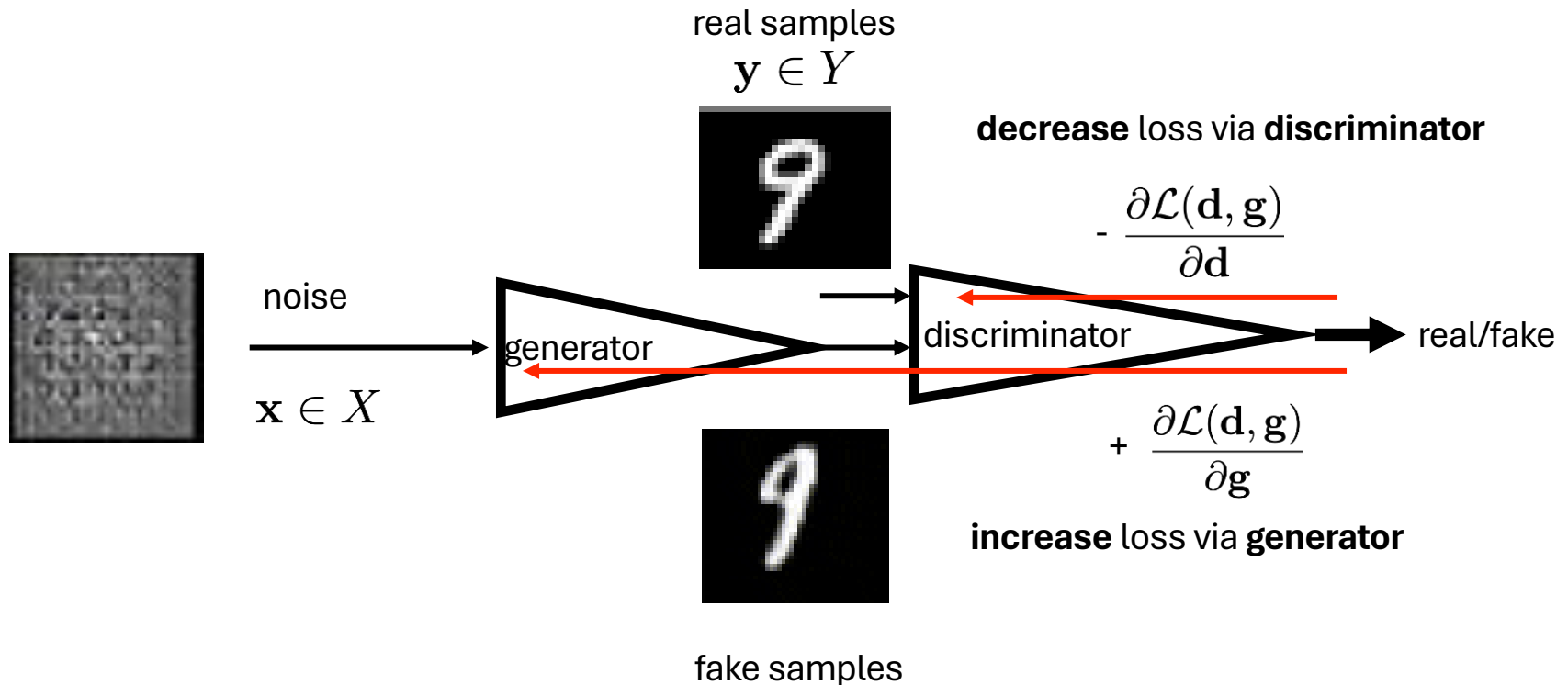
Generative Adversarial Networks (GANs)



Classification loss:

$$\mathcal{L}(\mathbf{d}, \mathbf{g}) = \sum_{\mathbf{x} \in X} -\log(\mathbf{d}(\mathbf{g}(\mathbf{x}))) + \sum_{\mathbf{y} \in Y} -\log(1 - \mathbf{d}(\mathbf{y}))$$

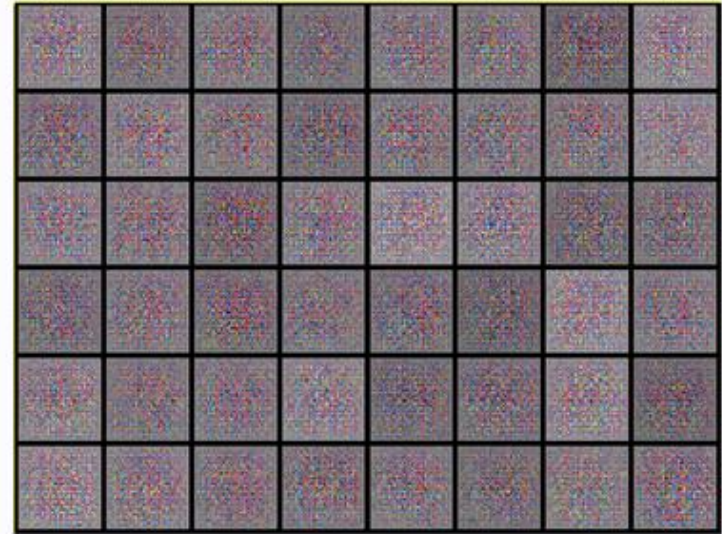
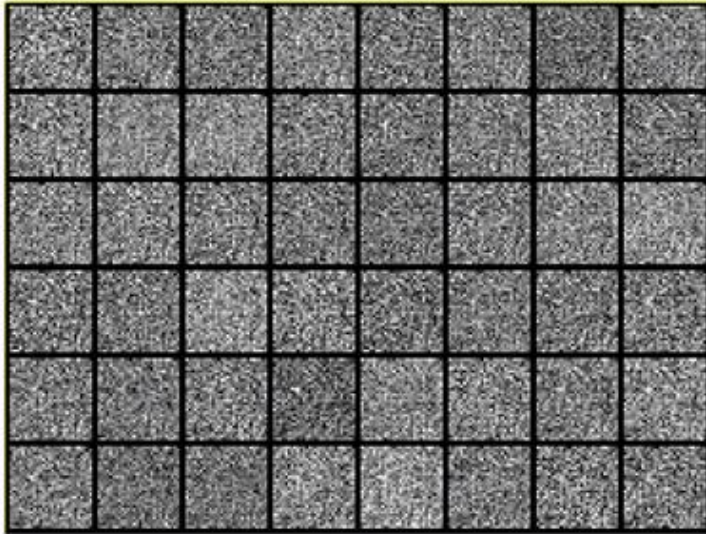
Generative Adversarial Nets



Classification loss:

$$\mathcal{L}(\mathbf{d}, \mathbf{g}) = \sum_{\mathbf{x} \in X} -\log(\mathbf{d}(\mathbf{g}(\mathbf{x}))) + \sum_{\mathbf{y} \in Y} -\log(1 - \mathbf{d}(\mathbf{y}))$$

Generative Adversarial Networks (GANs)

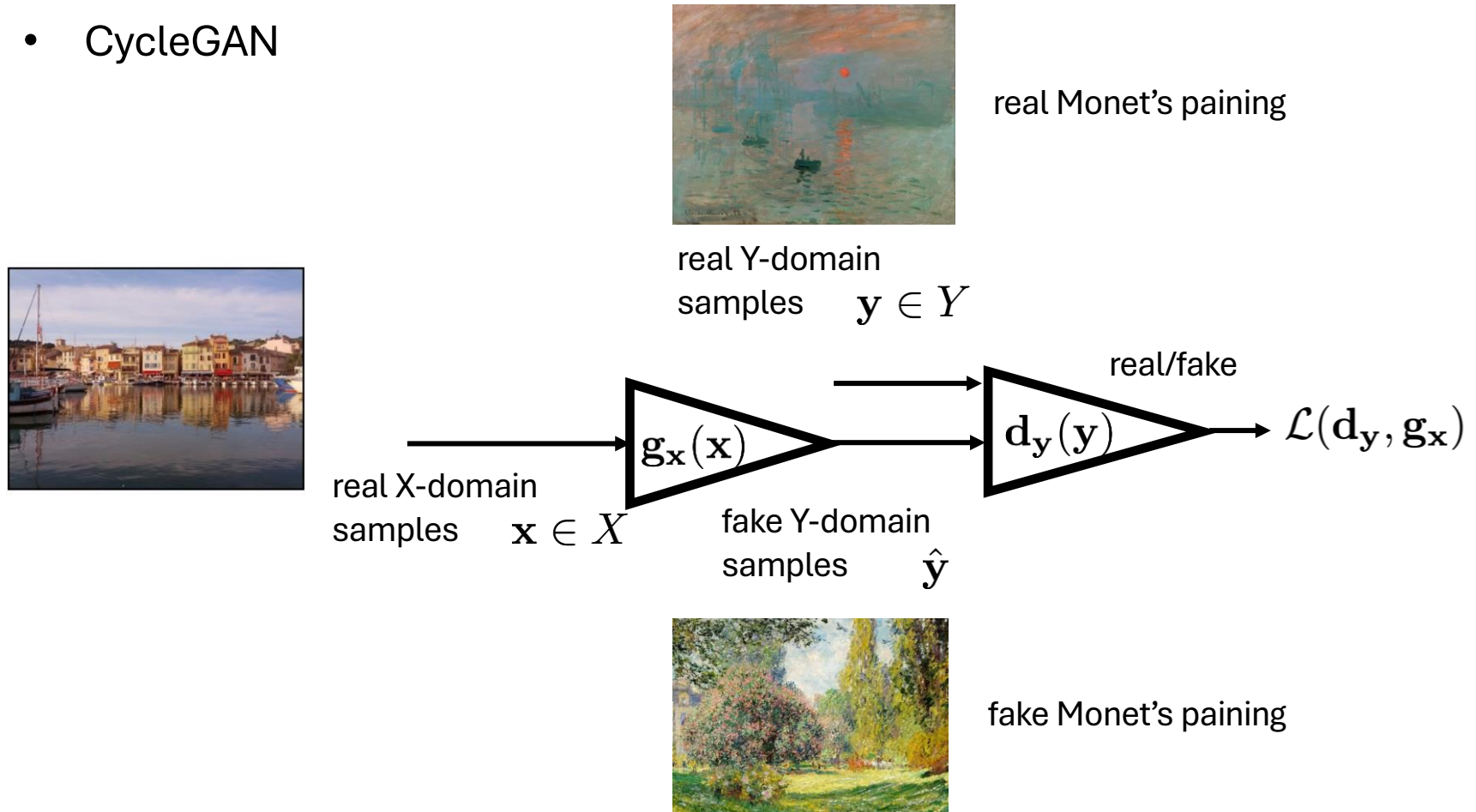


Generative Adversarial Networks (GANs)



Generative Adversarial Networks (GANs)

- CycleGAN



Generative Adversarial Networks (GANs)

- CycleGAN

$$|g_y(g_x(x)) - \hat{x}|$$

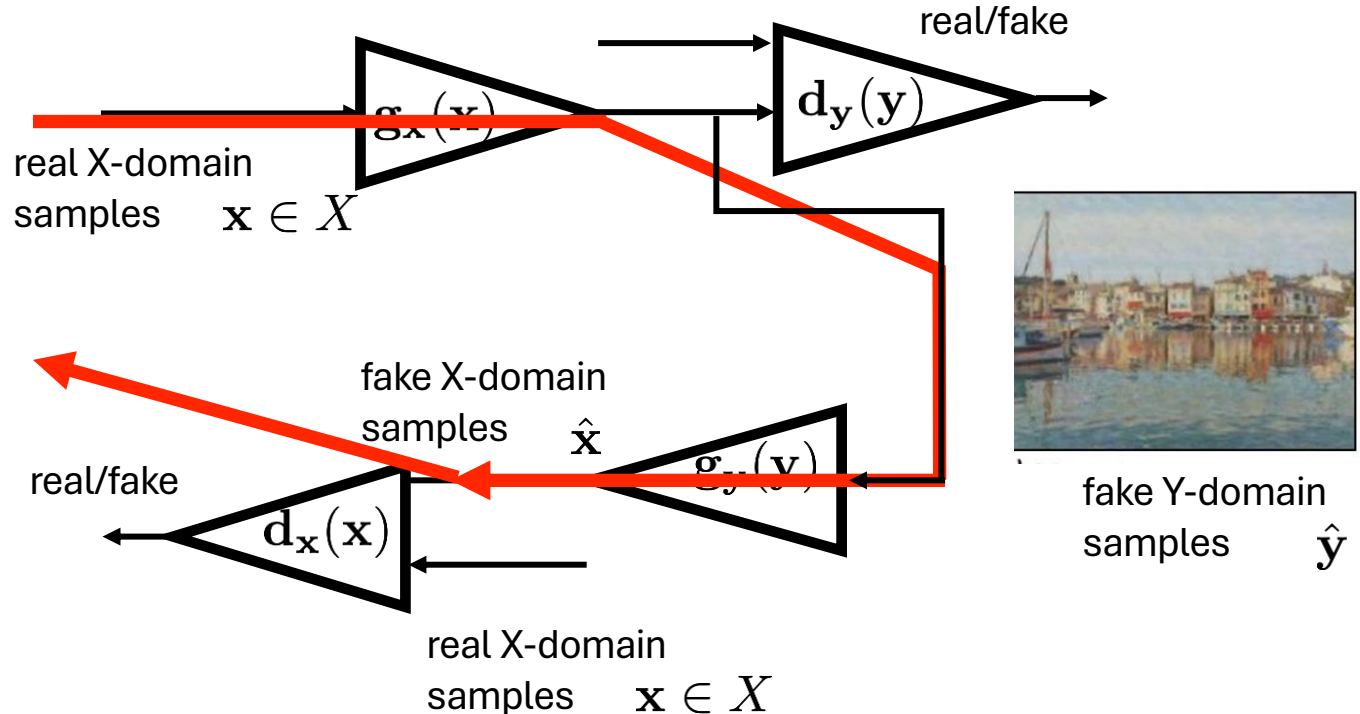


$$\mathcal{L}(d_x, g_y)$$



real Monet's painting

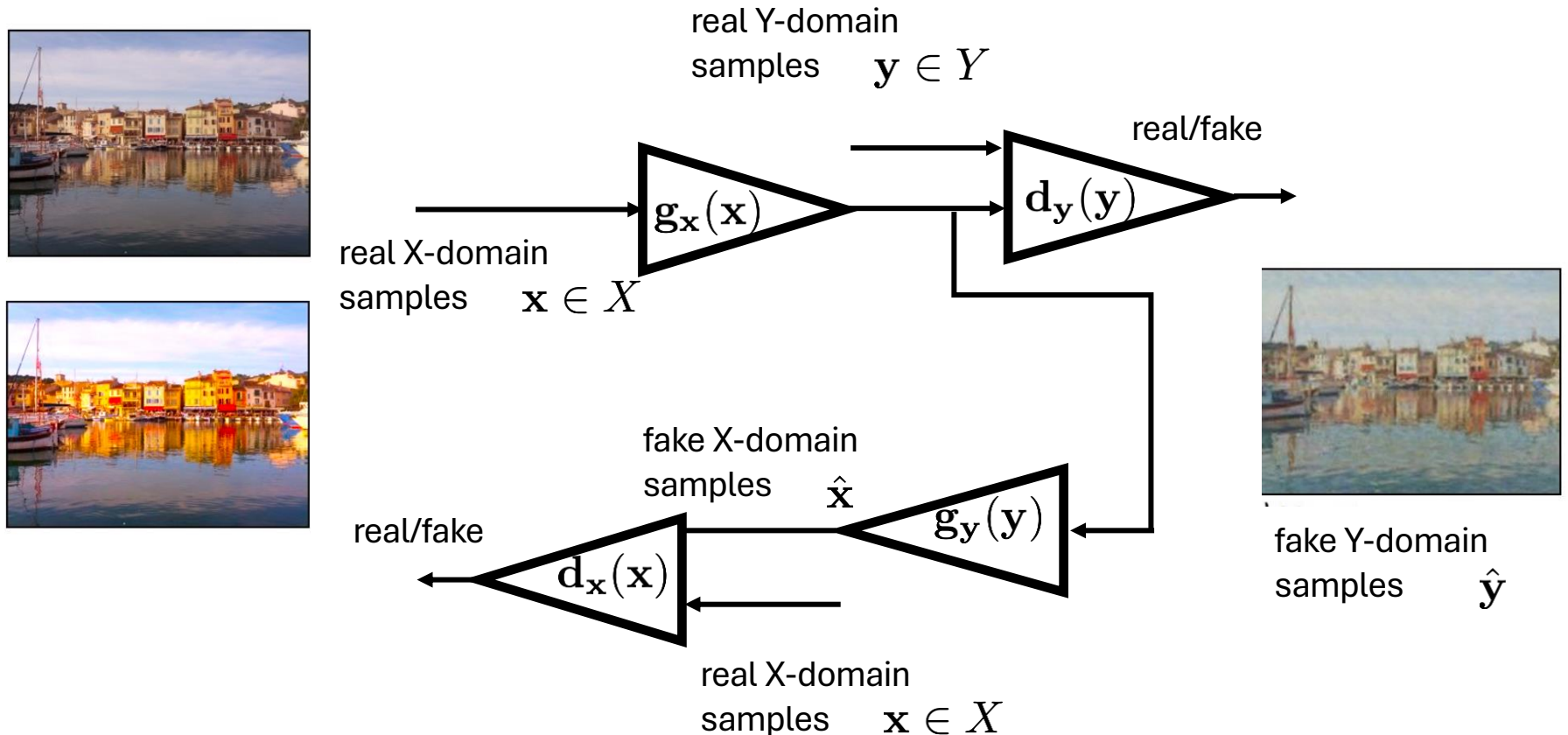
real Y-domain
samples $y \in Y$



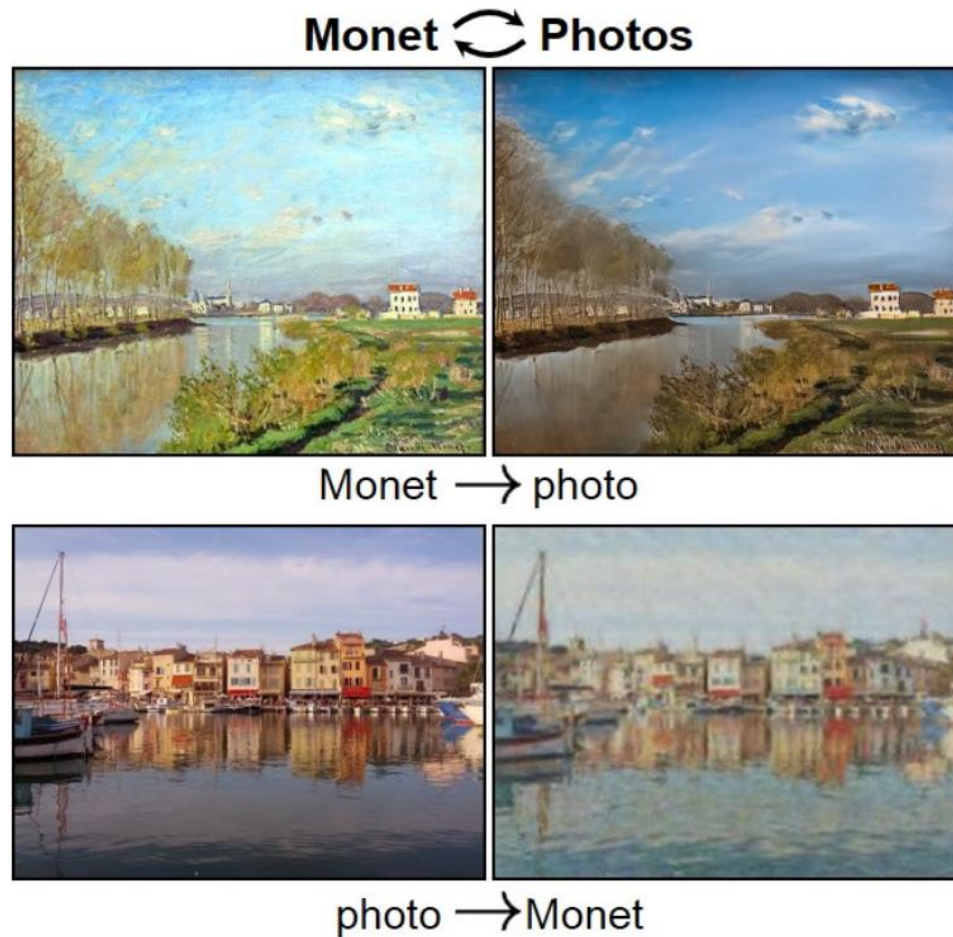
Generative Adversarial Networks (GANs)

- CycleGAN

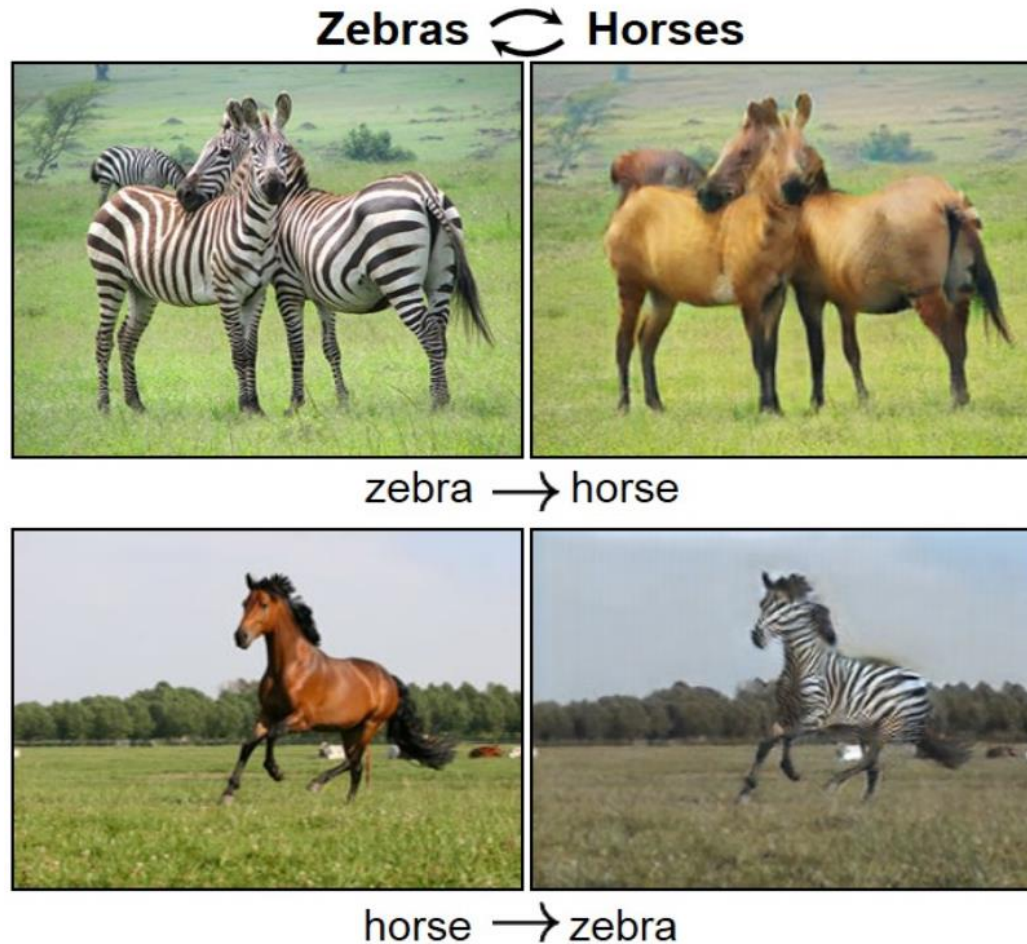
$$\mathcal{L}_{GAN}(\mathbf{d}_x, \mathbf{d}_y, \mathbf{g}_x, \mathbf{g}_y) = \mathcal{L}(\mathbf{d}_x, \mathbf{g}_y) + \mathcal{L}(\mathbf{d}_y, \mathbf{g}_x) + |\mathbf{g}_y(\mathbf{g}_x(\mathbf{x}) - \hat{\mathbf{x}}|$$



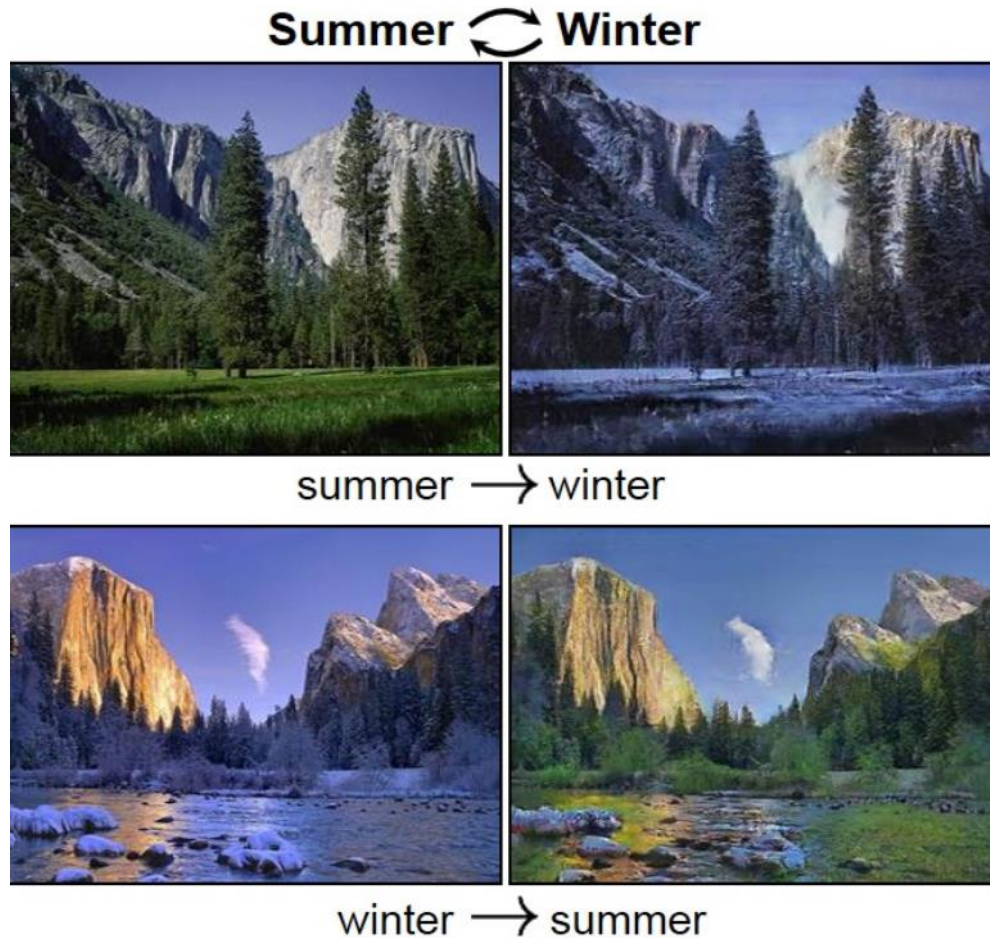
Generative Adversarial Networks (GANs)



Generative Adversarial Networks (GANs)



Generative Adversarial Networks (GANs)



Diffusion models

Particle diffusion

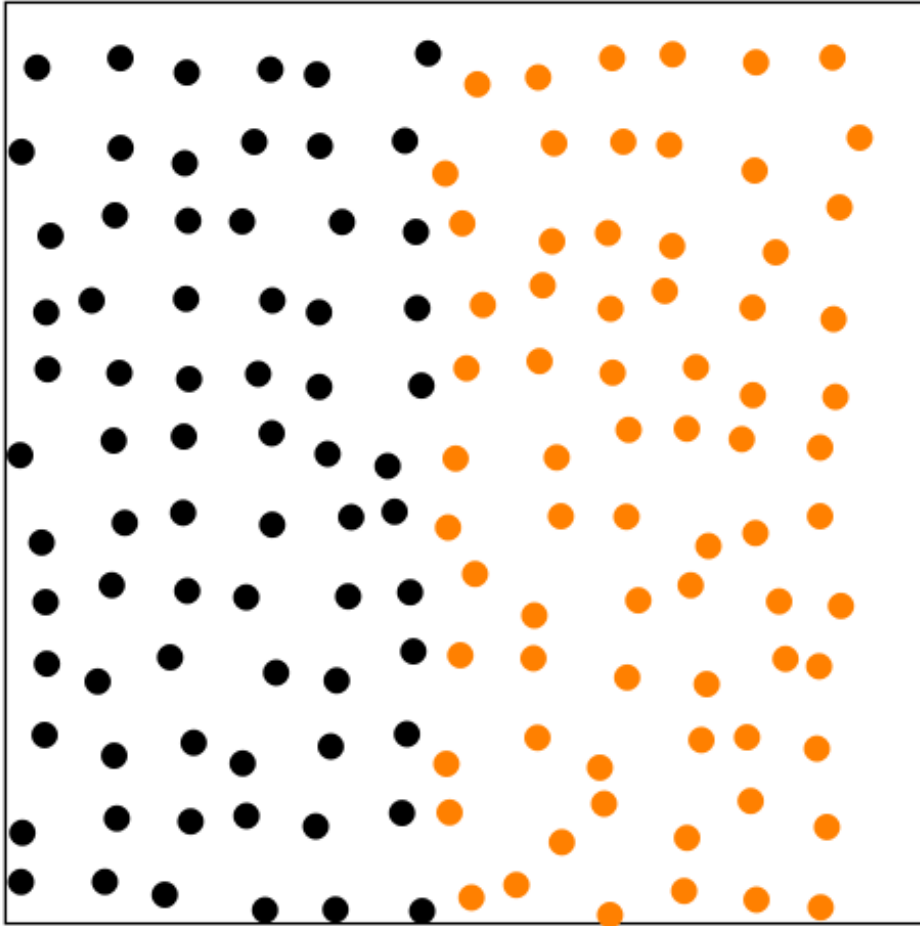


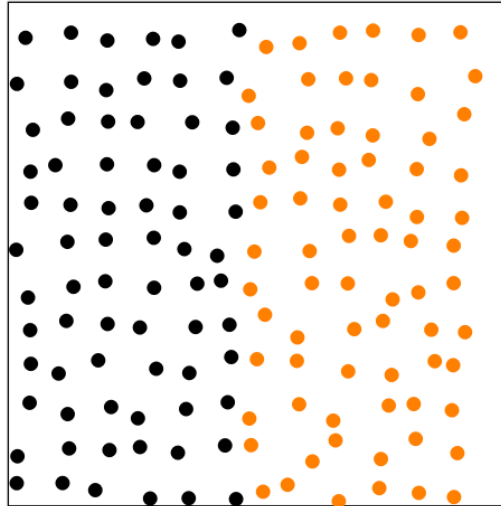
Image diffusion



Is the process reversible?

Diffusion models

Particle diffusion



Stochastic

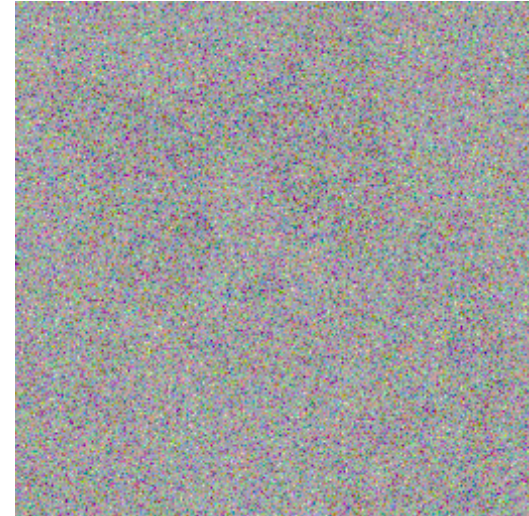
Entropy
increases

Partially
Reversible



$$\frac{\partial u}{\partial t} = D \nabla^2 u$$

Image diffusion

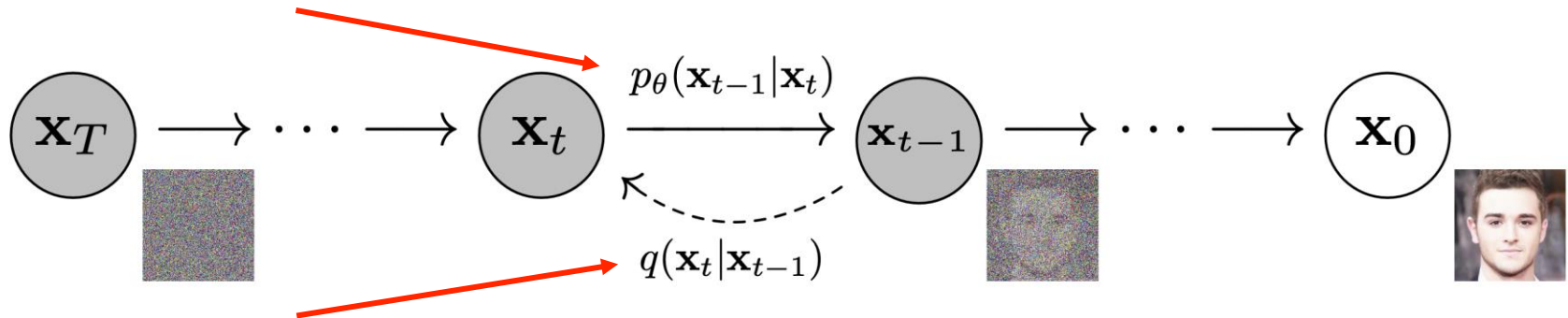


$$x_t = \sqrt{\alpha_t} x_{t-1} + \sqrt{1 - \alpha_t} \epsilon$$

Diffusion models

- Train a network to reverse the diffusion process

Reverse of the diffusion process to generate original data from the noise



Markov chain of diffusion steps in which we slowly and randomly add noise t

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\epsilon$$

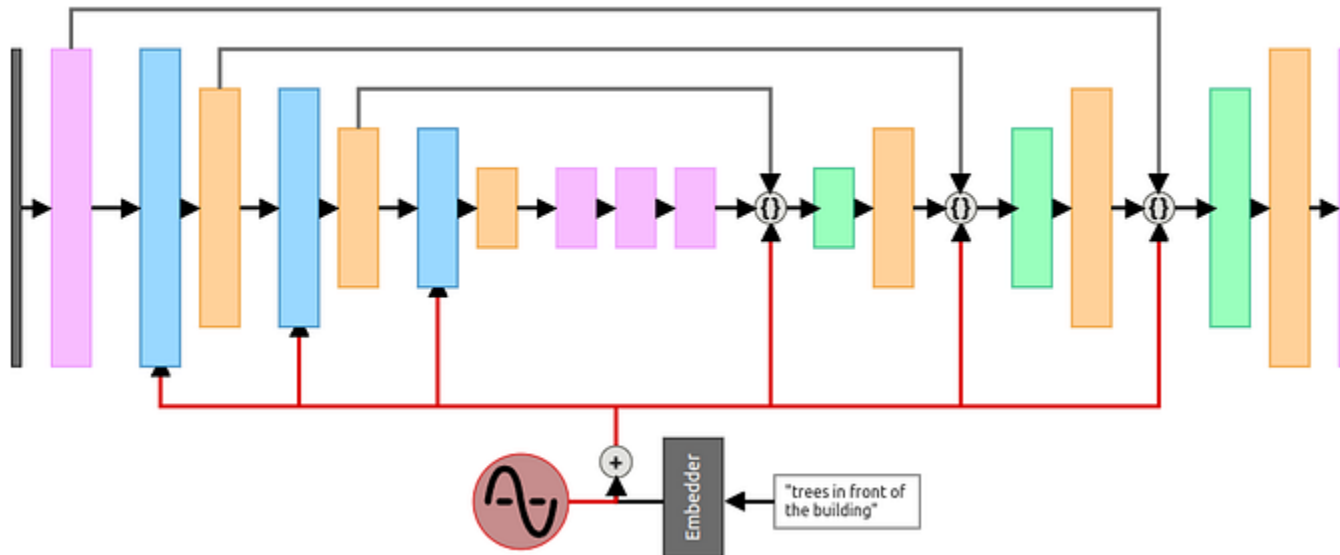
If noise is small the backward step has also “almost” gaussian distribution

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

→ train de-noising networks through L2-norm with the original image

Diffusion models

- Train a network to reverse the diffusion process
- An adapted U-Net architecture for denoising allows customizing the embedding
- Pros
 - High-quality imagery
- Cons
 - **Slow** (the generation process is iterative, thousands of steps needed)



Diffusion models

- Stable diffusion

Text-to-Image Synthesis on LAION. 1.45B Model.

'A street sign that reads
"Latent Diffusion" '



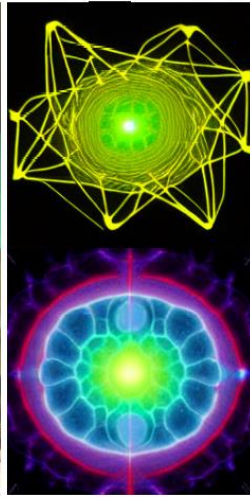
'A zombie in the
style of Picasso'



'An image of an animal
half mouse half octopus'



'An illustration of a slightly
conscious neural network'



'A painting of a
squirrel eating a burger'



'A watercolor painting of a
chair that looks like an octopus'



'A shirt with the inscription:
"I love generative models!" '



Competencies gained for the test

- Semantic segmentation
- Object detection, evaluation (Precision, Recall, AP)
- Generative models (GANs, Diffusion models)