

# MultiTask Learning

Hendrik Santoso Sugiarto

# Capaian Pembelajaran

- Multitask
  - Mengerti sistem pendekripsi objek
  - Mengerti sistem segmentasi objek
  - Mengerti sistem T5
- Siamese
  - Mengerti sistem pengenalan wajah
  - Mengerti sistem pengenalan tulisan
- Multiobjective
  - Mengerti sistem transfer style

# Multitask

# Deteksi Objek

- Lokalisasi vs deteksi vs segmentasi

~~Classification~~



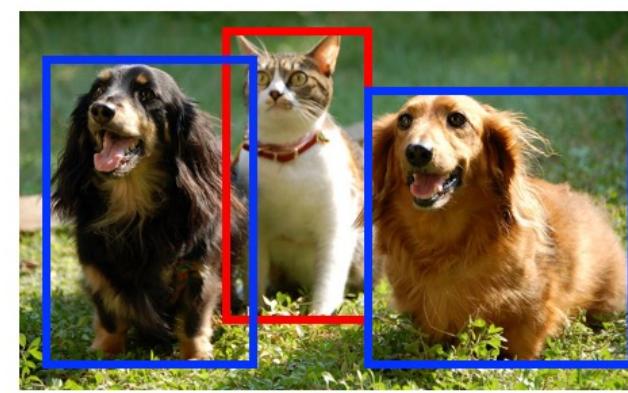
CAT

~~Classification + Localization~~



CAT

~~Object Detection~~



CAT, DOG

~~Instance Segmentation~~



CAT, DOG

Single object

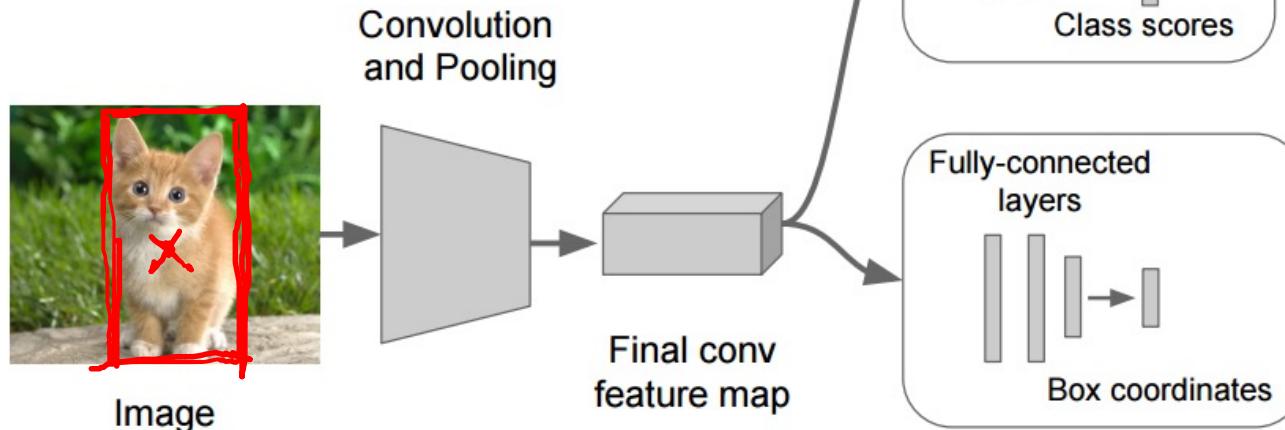
Multiple objects

# Lokalisasi

- Multitask Learning

- Prediksi:

- ① • Apakah ada objek?
- ② • Lokasi pusat kotak ( $x, y$ )
- ③ • Tinggi kotak
- ④ • Lebar kotak
- ⑤ • Kategori objek



“Classification head”

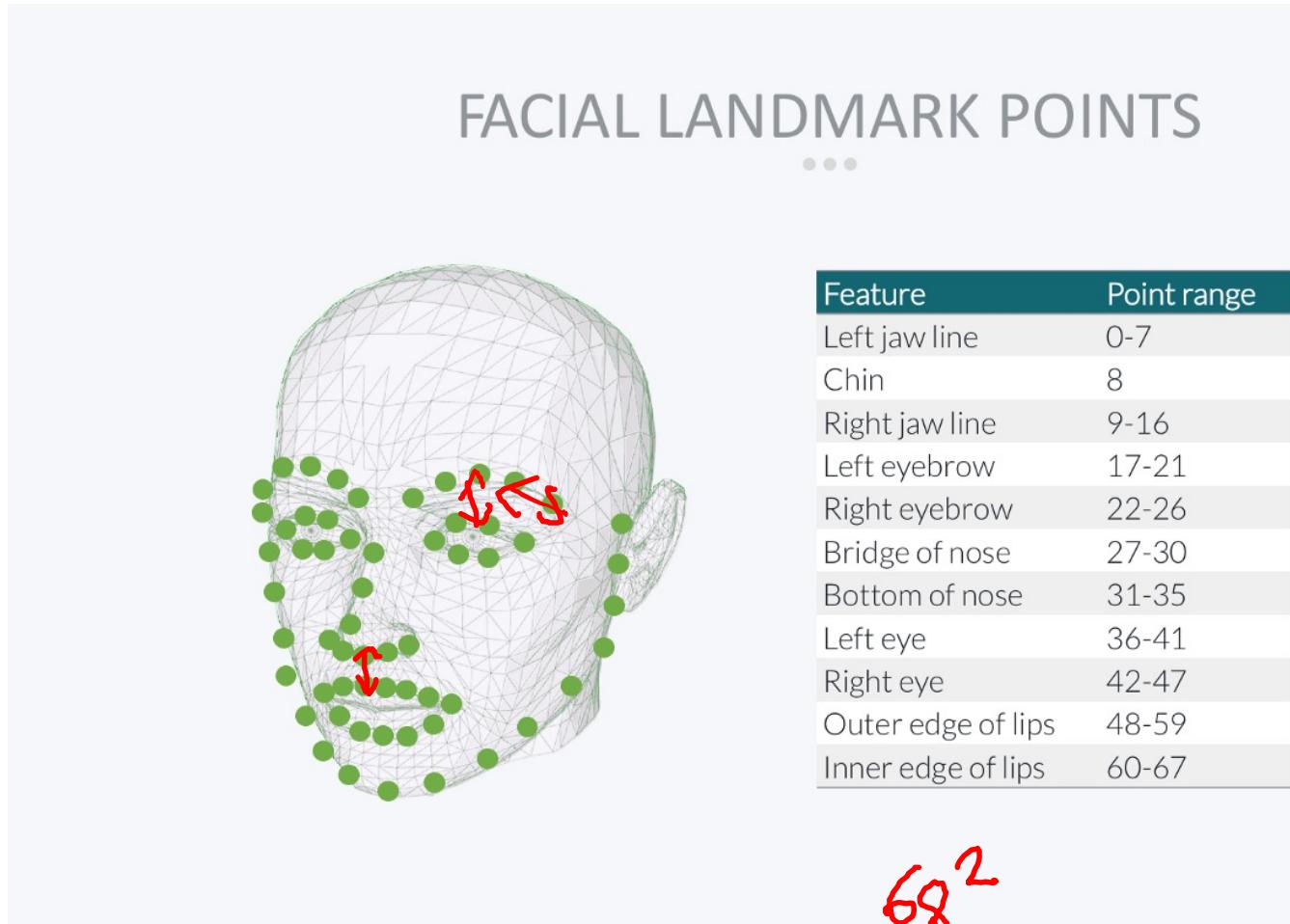
*cat dog bird .*

“Regression head”

$x, y$   
 $\Delta x, \Delta y$

# Deteksi Landmark

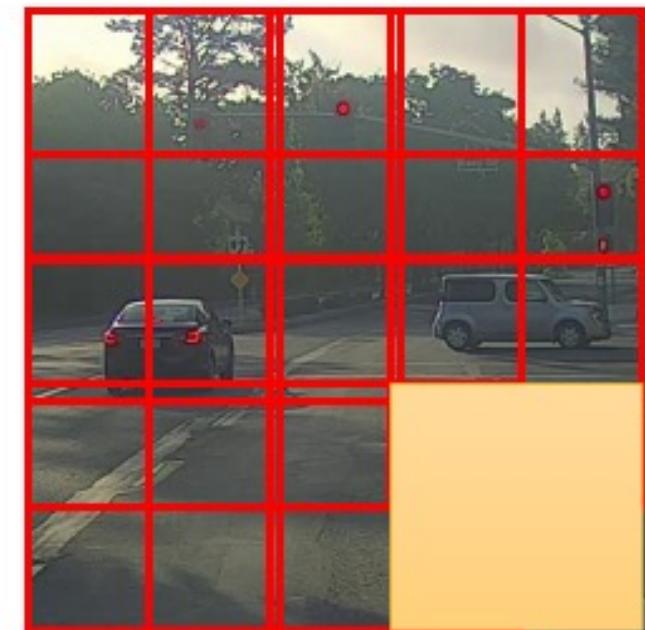
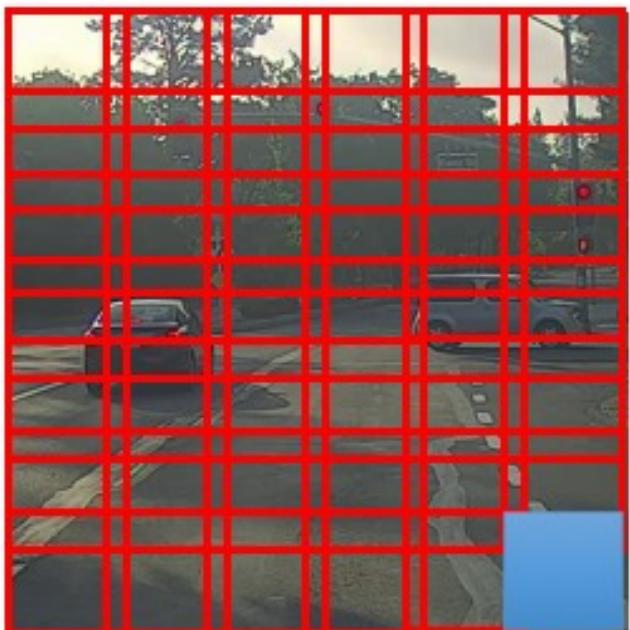
- Multitask learning dapat juga digunakan untuk memprediksi lokasi banyak poin sekaligus



Feature	Point range
Left jaw line	0-7
Chin	8
Right jaw line	9-16
Left eyebrow	17-21
Right eyebrow	22-26
Bridge of nose	27-30
Bottom of nose	31-35
Left eye	36-41
Right eye	42-47
Outer edge of lips	48-59
Inner edge of lips	60-67

68<sup>2</sup>

# Sliding Window

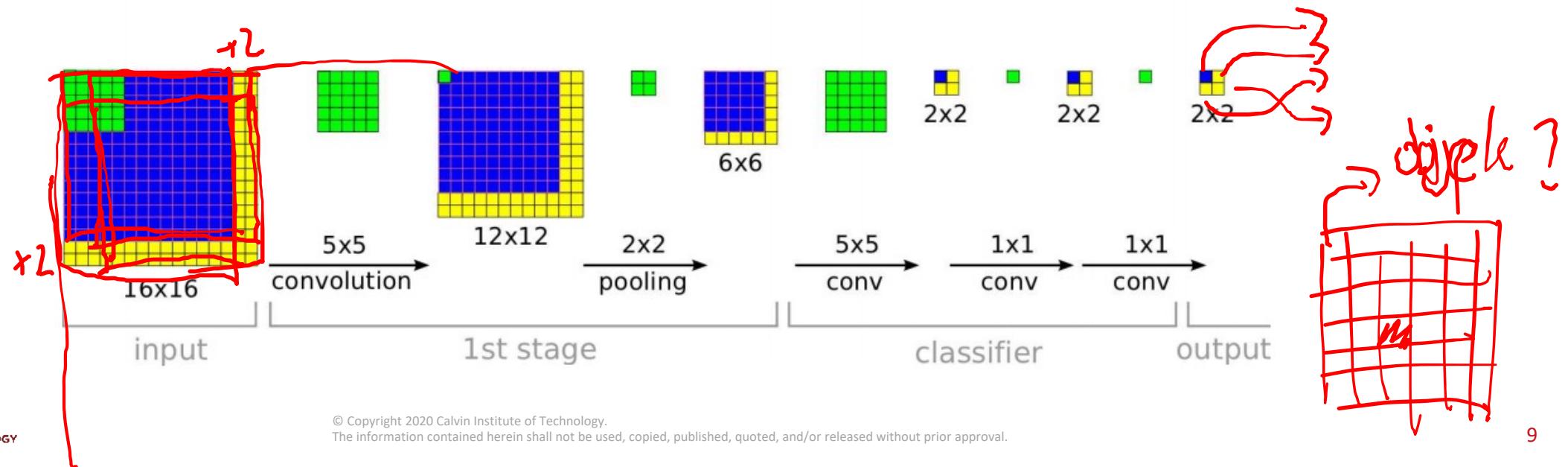
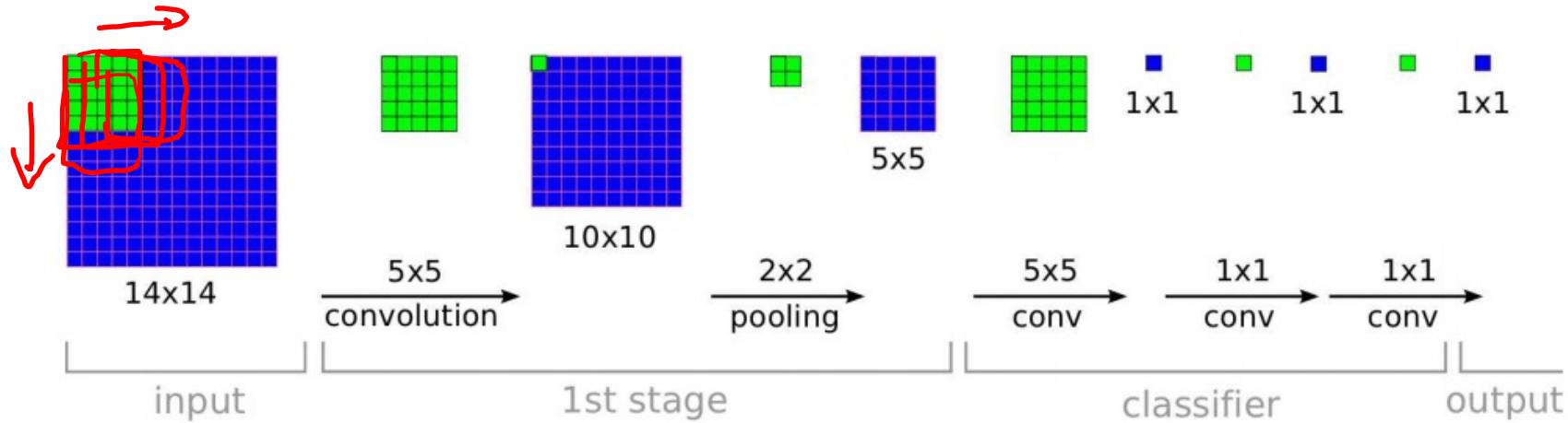


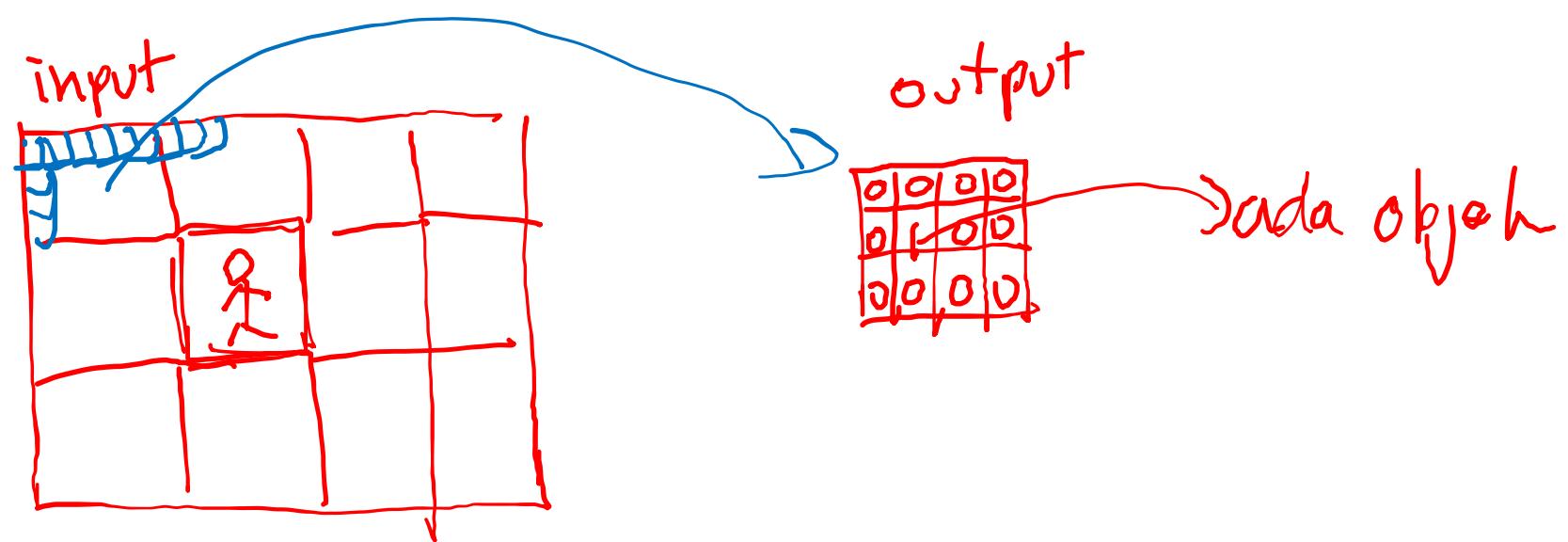
# Deteksi Objek dengan Sliding Window

- Latih model untuk mendeteksi objek tertentu
- Coba prediksi pada semua kemungkinan sliding window (dengan berbagai ukuran) untuk melihat apakah terdapat objek pada lokasi tersebut
- Problem: sangat lambat dan mahal



# Implementasi konvolusi terhadap sliding window

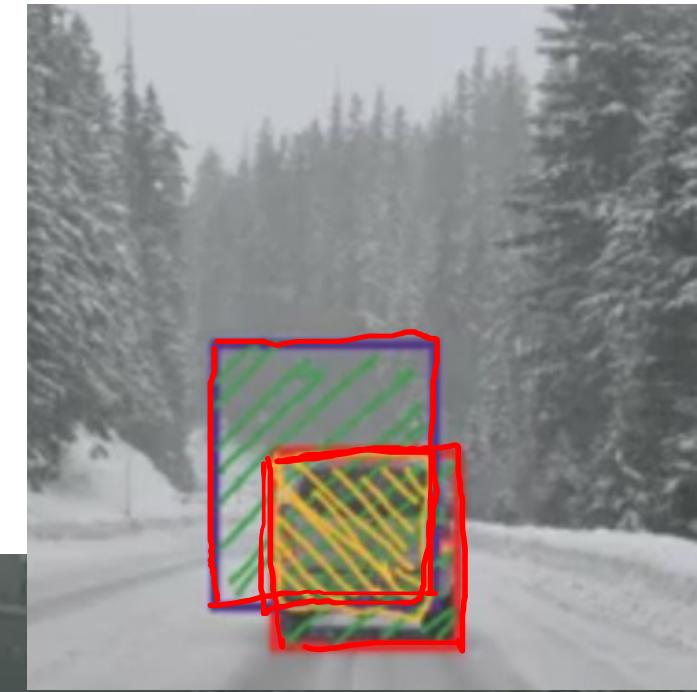
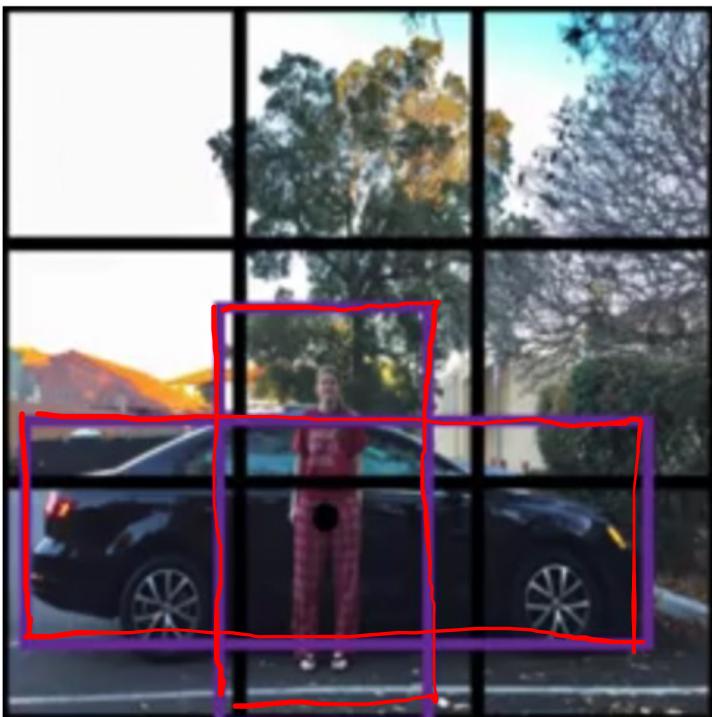




# YOLO (you only look once)

- Komponen:

- Konvolusi sliding windows
- Intersection over union
- Non-max suppression
- Anchor boxes

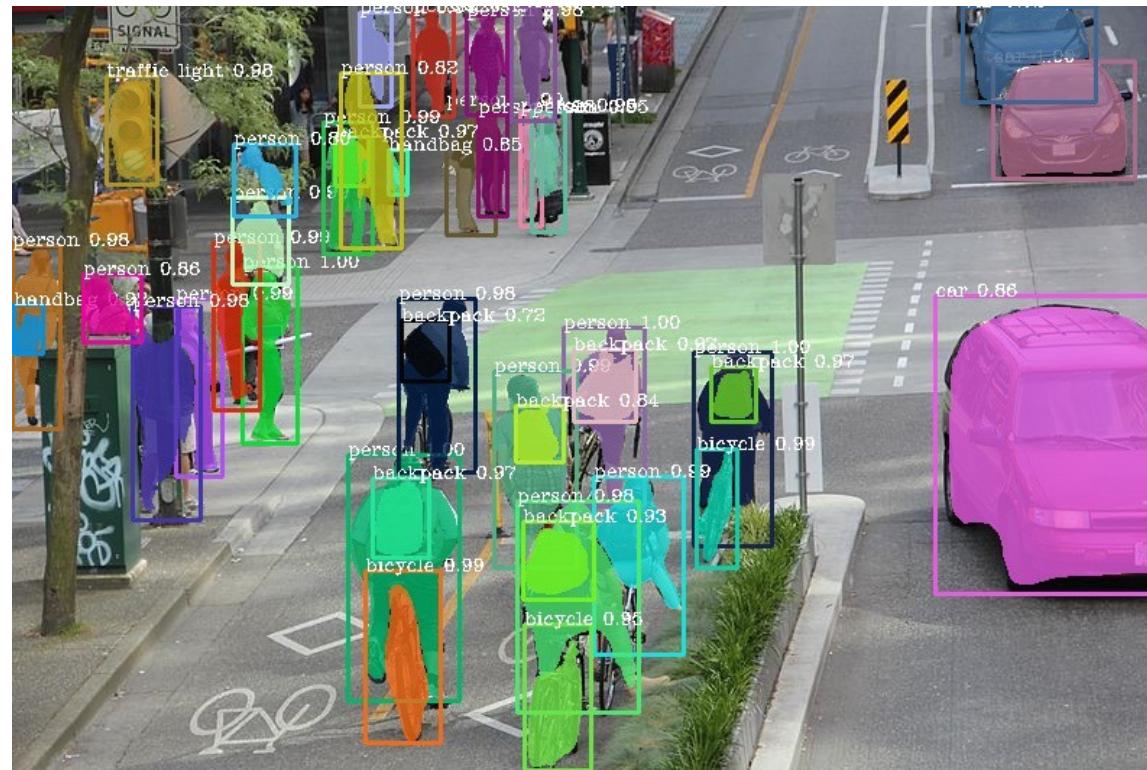


# YOLO loss function

$$\lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \quad \begin{array}{|l} \hline \text{1 when there is object, 0 when there is no object} \\ \hline \end{array}$$
$$+ \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[ (\sqrt{w_i} - \sqrt{\hat{w}_i})^2 + (\sqrt{h_i} - \sqrt{\hat{h}_i})^2 \right] \quad \begin{array}{|l} \hline \text{Bounding Box Location (x, y) when there is object} \\ \hline \text{Bounding Box size (w, h) when there is object} \\ \hline \end{array}$$
$$+ \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \quad \begin{array}{|l} \hline \text{Confidence when there is object} \\ \hline \end{array}$$
$$+ \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \quad \begin{array}{|l} \hline \text{1 when there is no object, 0 when there is object} \\ \hline \text{Confidence when there is no object} \\ \hline \end{array}$$
$$+ \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad \begin{array}{|l} \hline \text{Class probabilities when there is object} \\ \hline \end{array}$$

# Segmentasi

- Problem dari sliding window adalah banyak perhitungan terbuang untuk mendeteksi ruang kosong
- Kita bisa melakukan segmentasi terlebih dahulu sebelum mendeteksi



# Region Proposal

- [https://pytorch.org/tutorials/intermediate/torchvision\\_tutorial.html](https://pytorch.org/tutorials/intermediate/torchvision_tutorial.html)
- R-CNN (Region CNN):
  - Ajukan area.
  - Klasifikasi area yang diajukan satu-persatu
  - Output label + Bounding box
- Fast R-CNN:
  - Ajukan area.
  - Gunakan konvolusi dengan sliding windows untuk klasifikasi semua area yang diajukan
- Faster R-CNN:
  - Gunakan jejaring konvolusi untuk mengajukan area

# Seberapa efisien?

Faster!

FASTER!

Better!

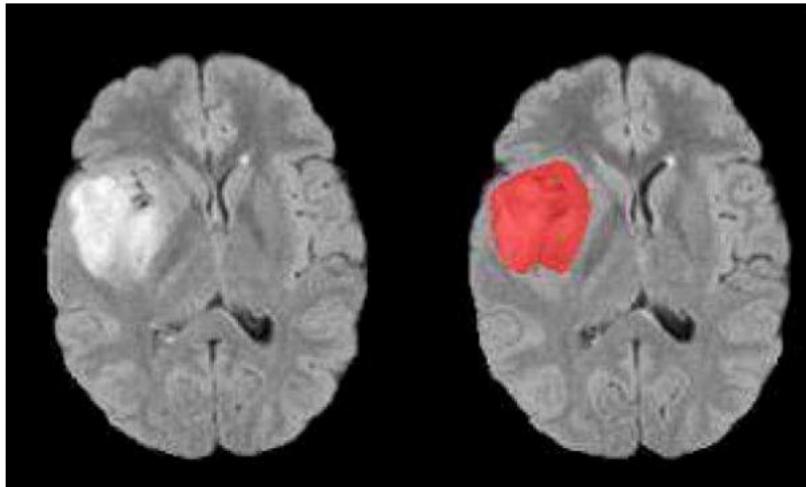
	R-CNN	Fast R-CNN
Training Time:	84 hours	<b>9.5 hours</b>
(Speedup)	1x	<b>8.8x</b>
Test time per image	47 seconds	<b>0.32 seconds</b>
(Speedup)	1x	<b>146x</b>
mAP (VOC 2007)	66.0	<b>66.9</b>

Using VGG-16 CNN on Pascal VOC 2007 dataset

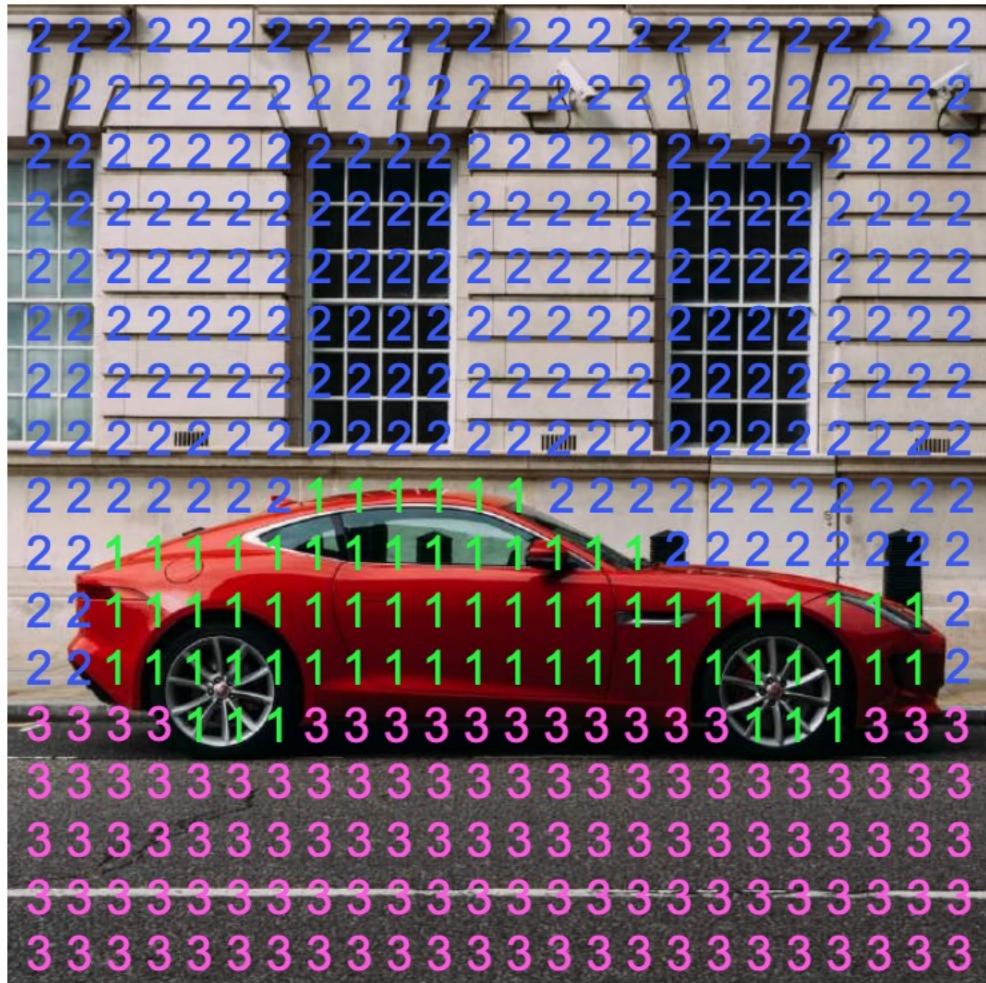
	R-CNN	Fast R-CNN	Faster R-CNN
Test time per image (with proposals)	50 seconds	2 seconds	<b>0.2 seconds</b>
(Speedup)	1x	25x	<b>250x</b>
mAP (VOC 2007)	66.0	<b>66.9</b>	<b>66.9</b>

# Semantic Segmentation

- Motivasi:
  - Deteksi area kanker

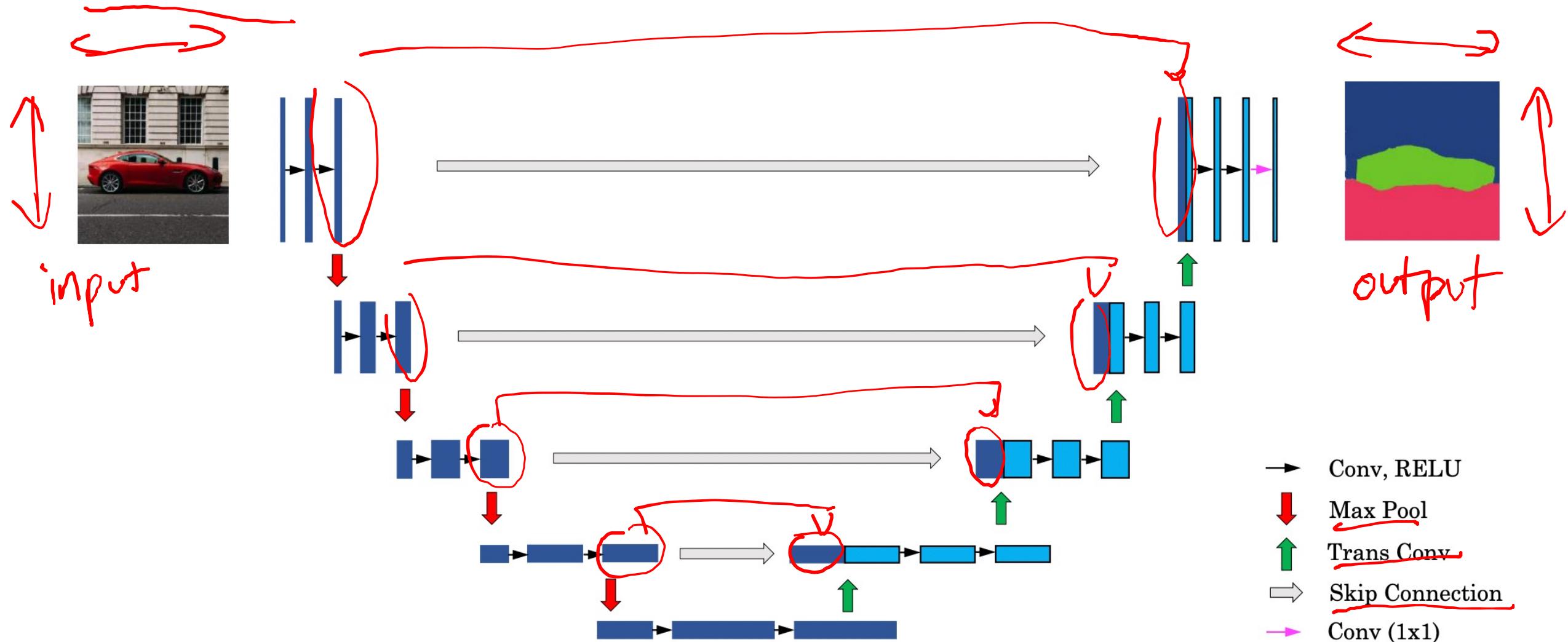


# Per pixel class label

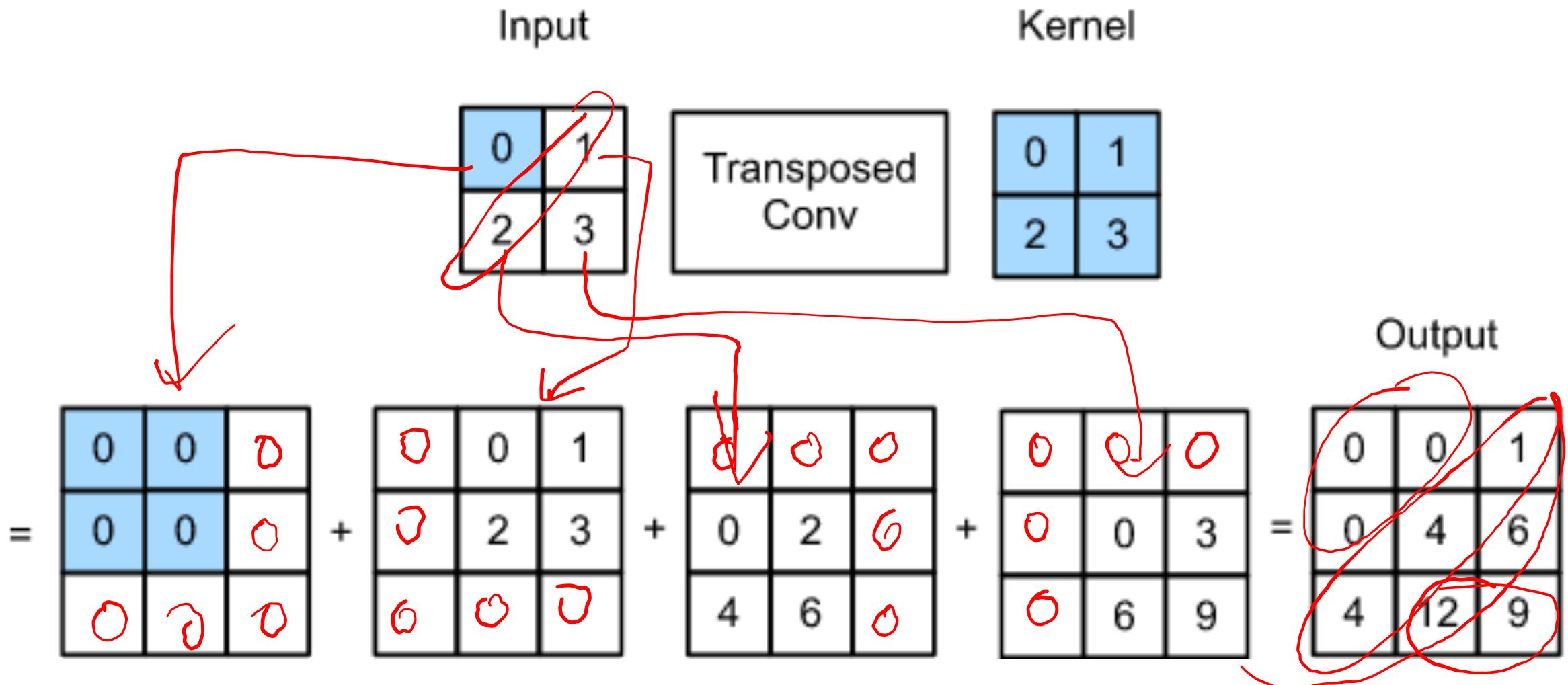


1. Car
  2. Building
  3. Road

# U-NET



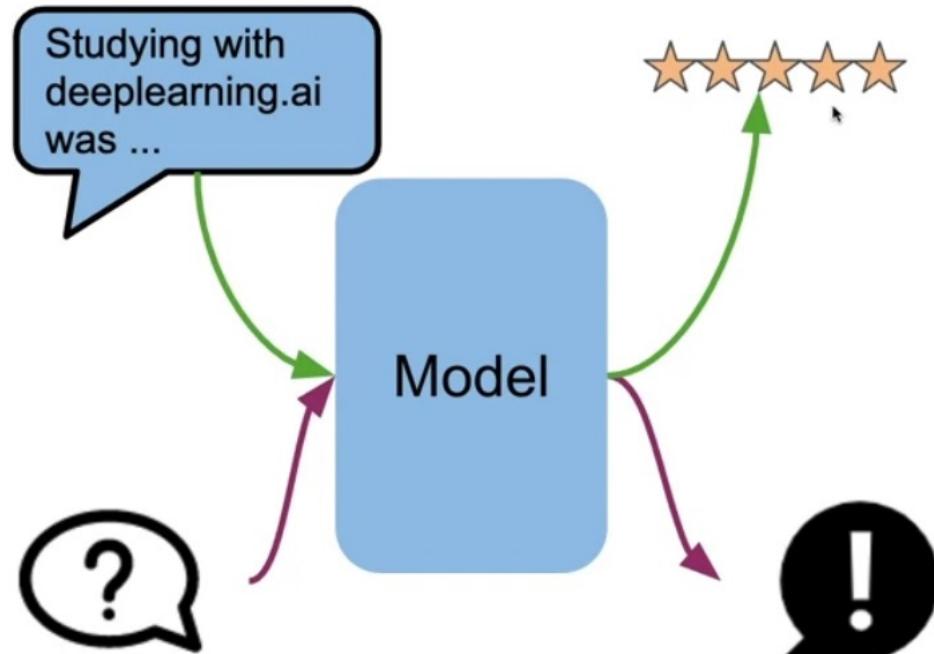
# Transpose Convolution



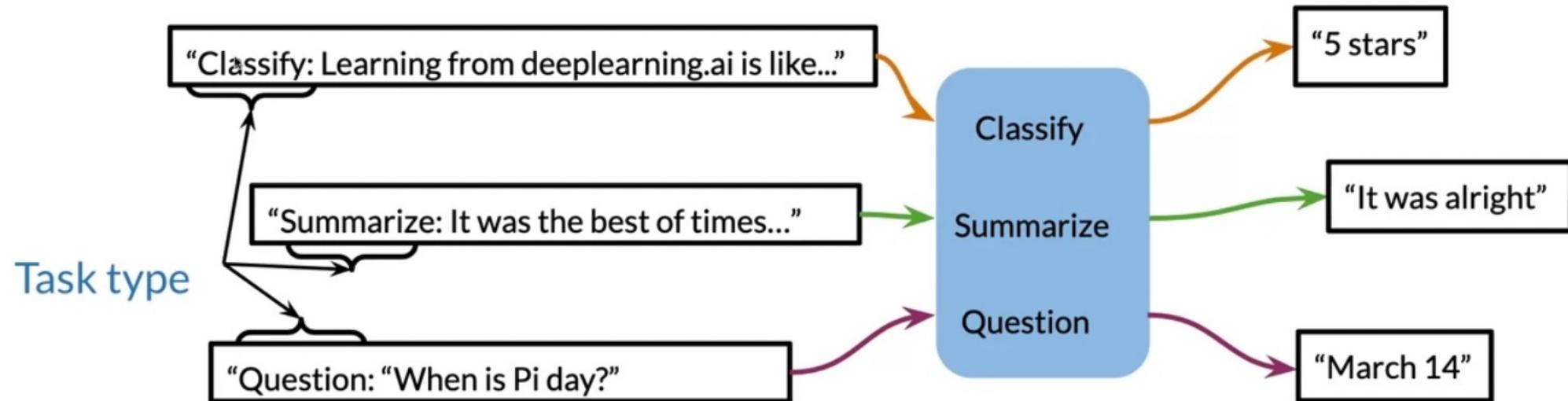
# T5: Encoder vs Encoder-Decoder



# T5: Multi-task



# T5: Text-to-text



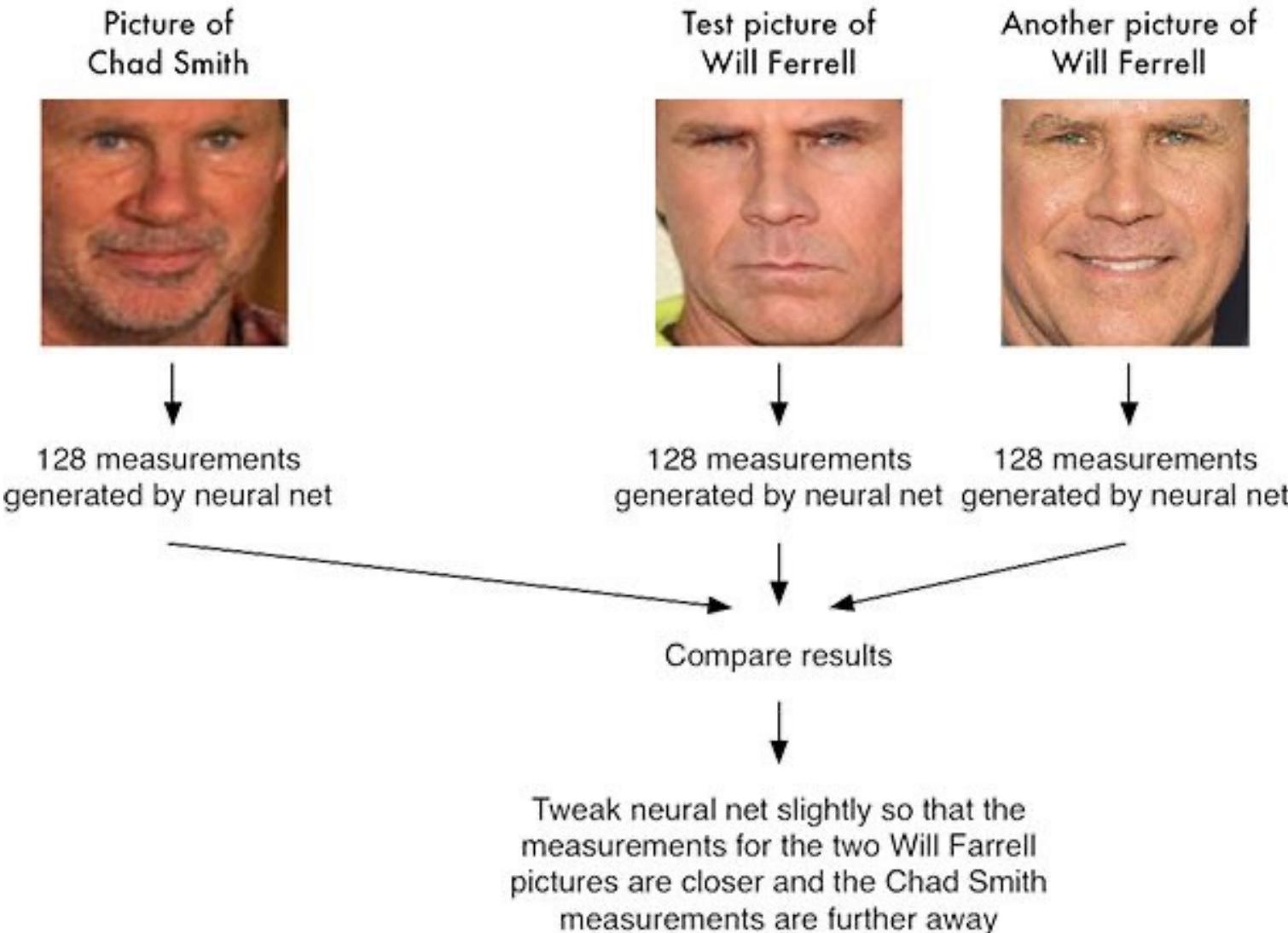
# Siamese



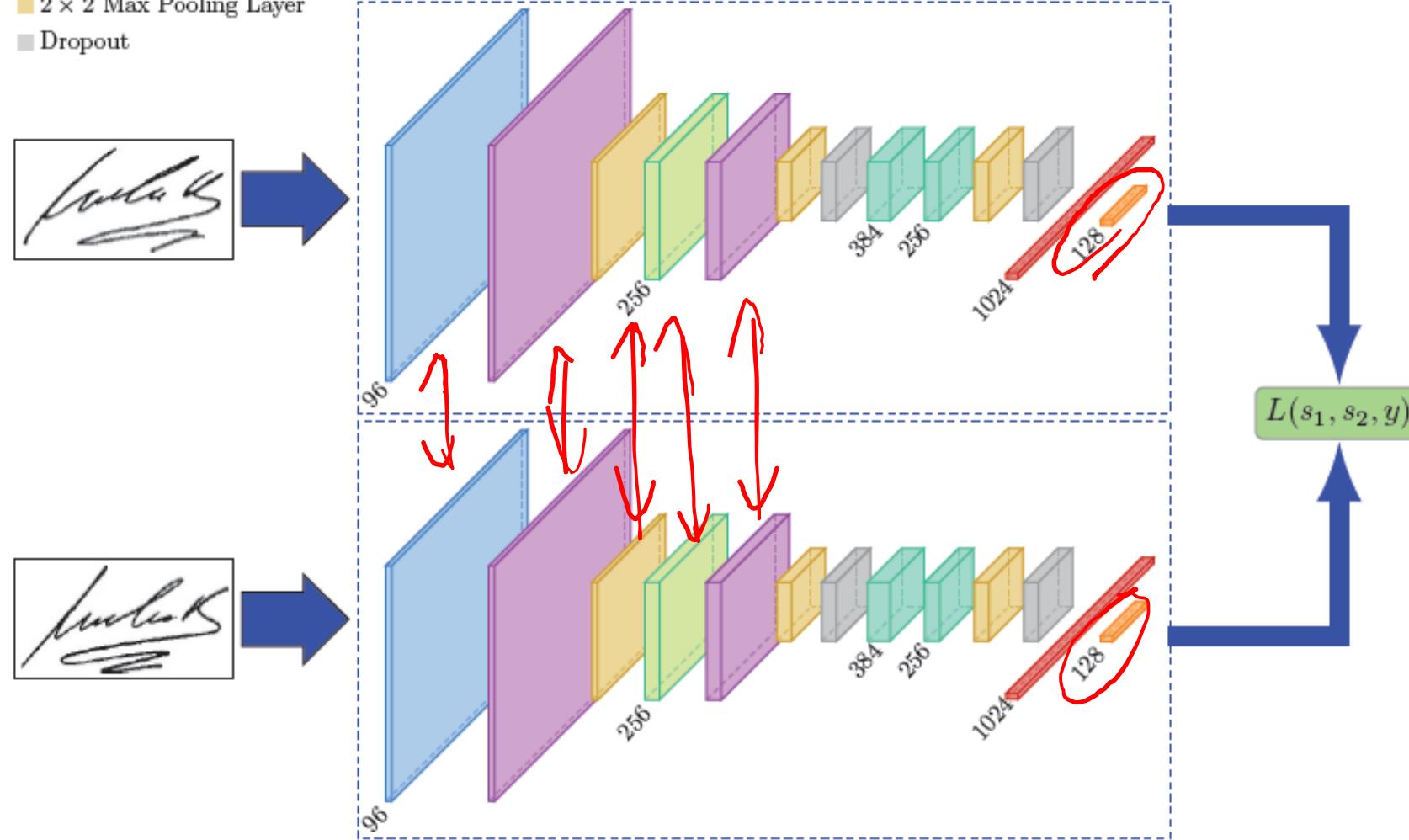
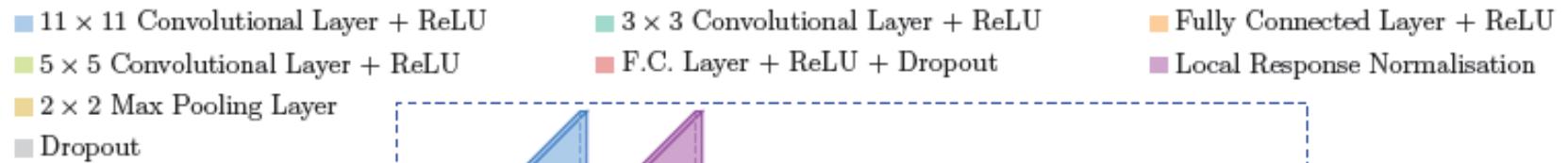
# Pengenalan wajah

- Problem: 1 shot learning
- Kita biasanya hanya mempunyai 1 foto per nama
- Menggunakan perbandingan pair (kemiripan antara 2 foto)

A single 'triplet' training step:

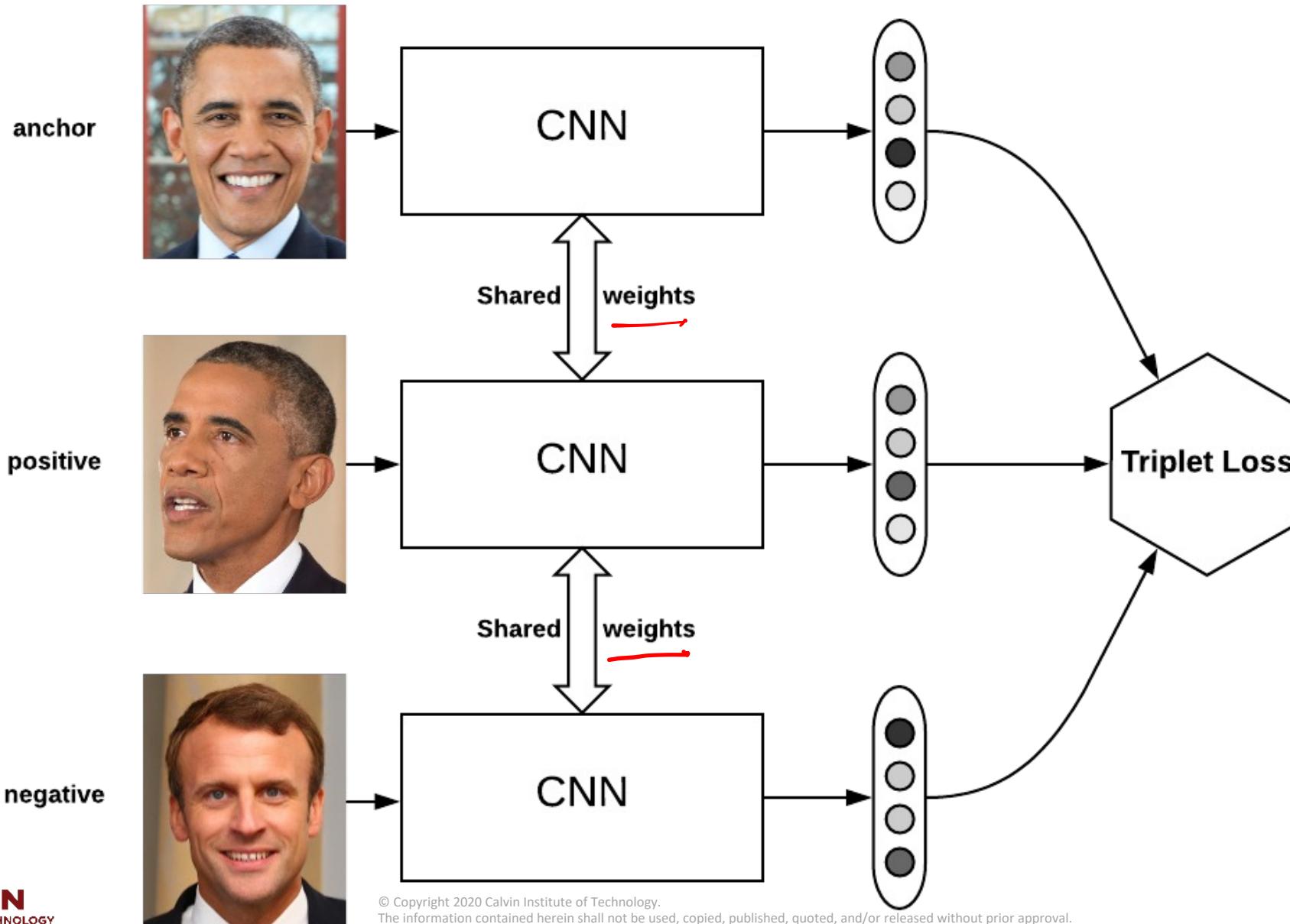


# Siamese Network



# Triplet Loss

Embeddings



# Duplikasi

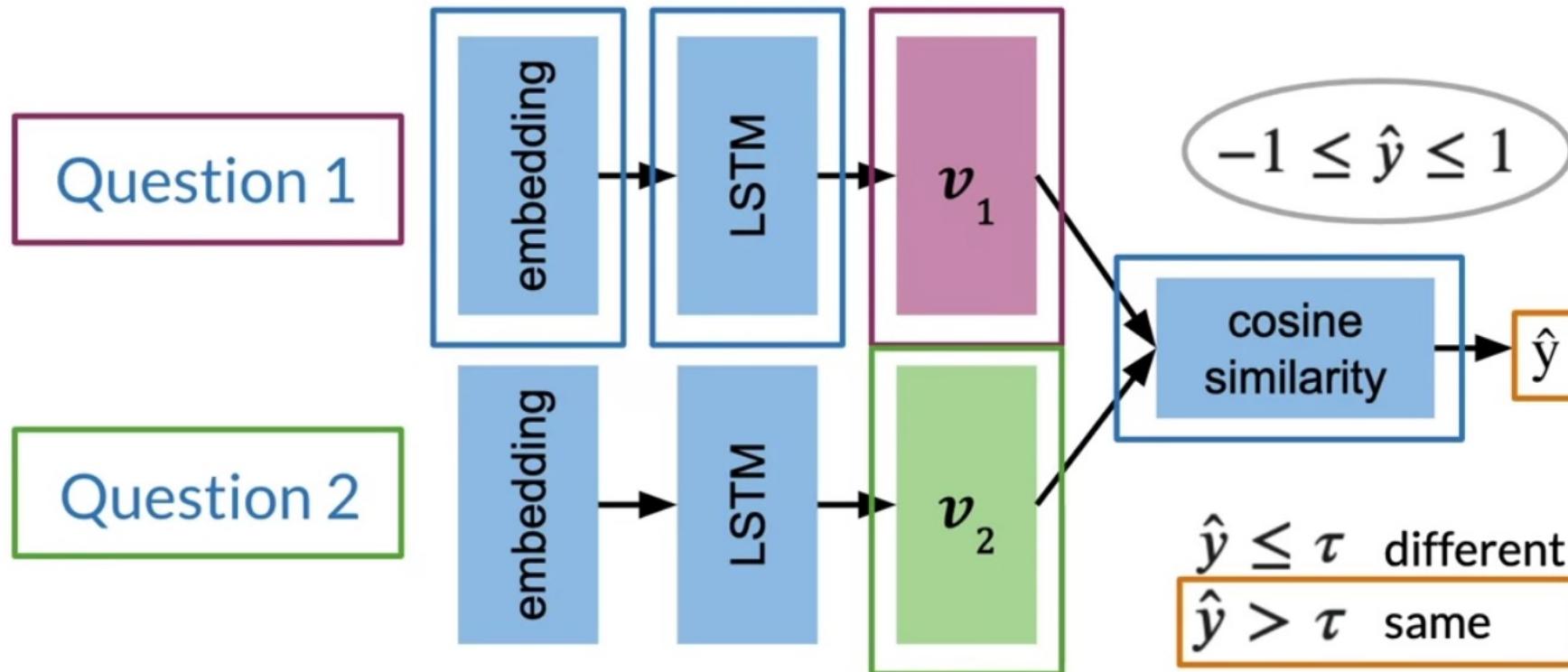
How old are you? =

What is your age?

Where are you from? ≠

Where are you going?

# Arsitektur



# Loss Function

How old are you?

Anchor

What is your age?

Positive

Where are you from?

Negative

$$\cos(v_1, v_2) = \frac{v_1 \cdot v_2}{\|v_1\| \|v_2\|}$$
$$s(v_1, v_2)$$

$$s(A, P) \approx 1$$

$$s(A, N) \approx -1$$

$$\text{diff} = s(A, N) - s(A, P)$$



# Uji Pemahaman

Apakah 2 subnetworks pada arsitektur Siamese share parameter yang sama?

- Benar ✓
- Salah



# Triplet Loss

## Triplet Loss

How old are you?

What is your age?

Where are you from?

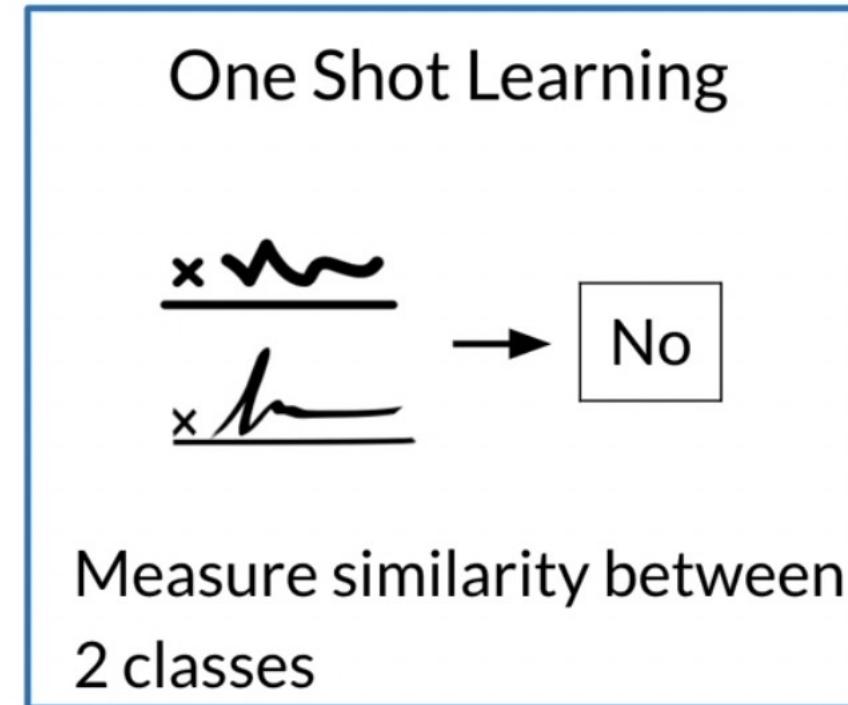
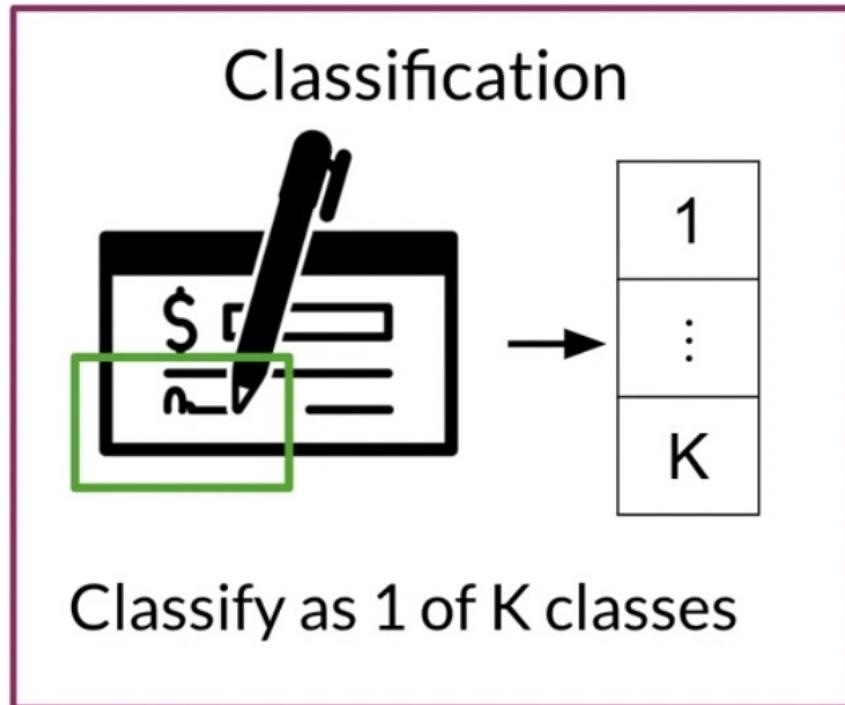
Simple loss:  
 $\text{diff} = s(A, N) - s(A, P)$

Anchor

Positive

Negative

# Classification vs One Shot Learning



# One Shot Learning

No need for retraining !



Learn a similarity score!

$s(sig1, sig2) > \tau$  0.9 ✓

$s(sig1, sig2) \leq \tau$  X

# Uji Pemahaman

- Manakah nilai similarity tertinggi?
  - Antara anchor dan positive ✓
  - Antara anchor dan negative
  - Antara positive dan negative

# Multiobjective

# Content + Style

Content image



+

Style image



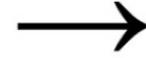
Output image



+



+



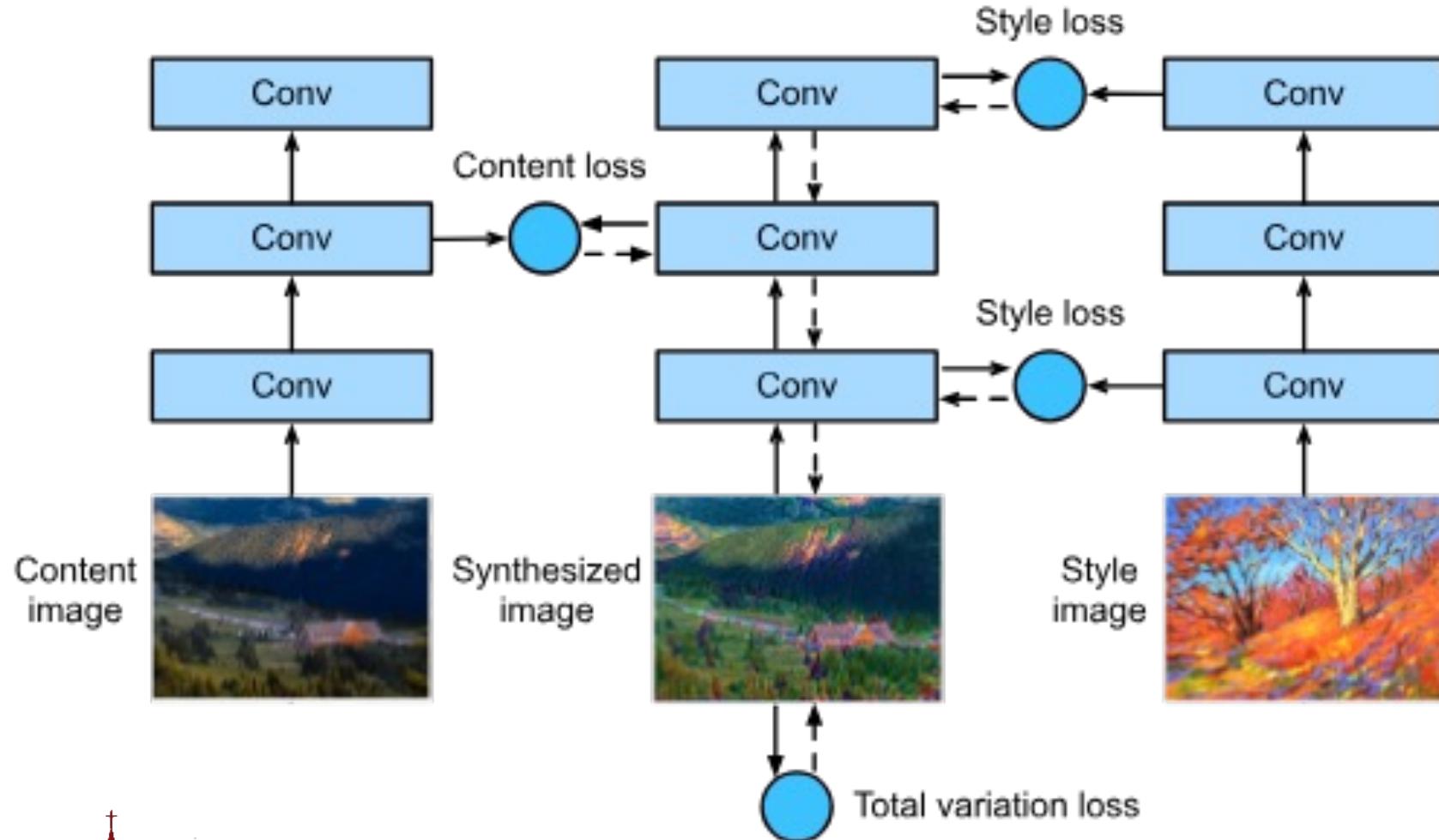
# Visualizing hidden layer

- <https://arxiv.org/pdf/1311.2901.pdf>



# Multi Objective Function

- $J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$



# Style Transfer Loss

$$J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$$

- Content loss function:

- Misalkan  $\underline{a^{(l)C}}$  dan  $\underline{a^{(l)G}}$  adalah layer aktivasi
- jika  $\underline{a^{(l)C}}$  dan  $\underline{a^{(l)G}}$  mirip, maka konten kedua gambar mirip

$$J_{content}(C, G) = \frac{1}{2} \left\| \underline{a^{(l)C}} - \underline{a^{(l)G}} \right\|^2$$

- Style loss function:

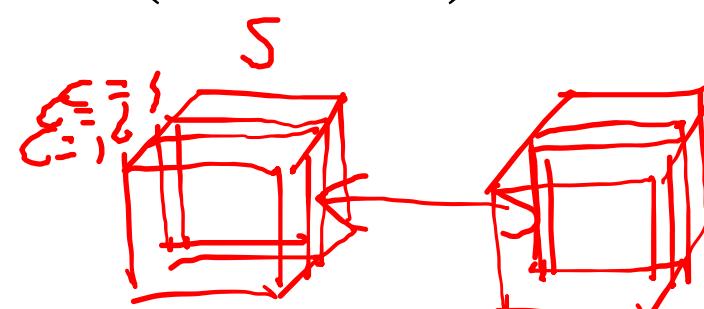
- Style didefinisikan sebagai korelasi antar aktivasi across channels

- Misalkan  $a_{i,j,k}^{(l)}$  merupakan aktivasi pada  $(i,j,k)$

- Maka  $G_{kk'}^{(l)} = \sum_{ij} \sum_{ij} a_{i,j,k}^{(l)} a_{i,j,k'}^{(l)}$

- $J_{style}(S, G) = \frac{1}{(2n_H^{(l)} n_W^{(l)} n_C^{(l)})^2} \sum_k \sum_{k'} \left( G_{kk'}^{(l)S} - G_{kk'}^{(l)G} \right)^2$

Gram Matrix



# GAN

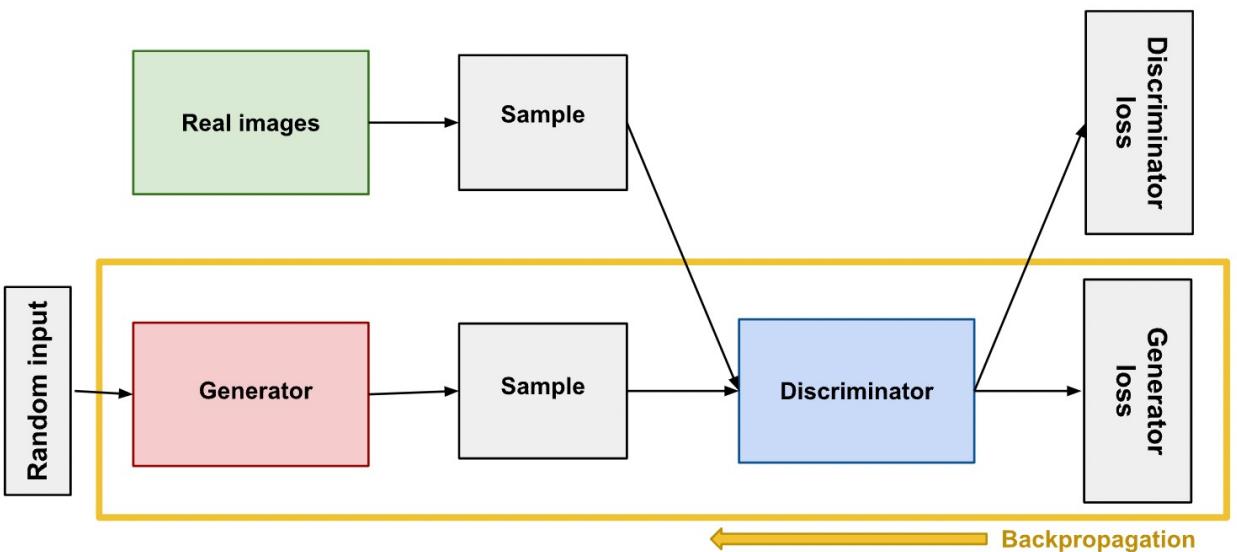


Figure 1: Backpropagation in generator training.

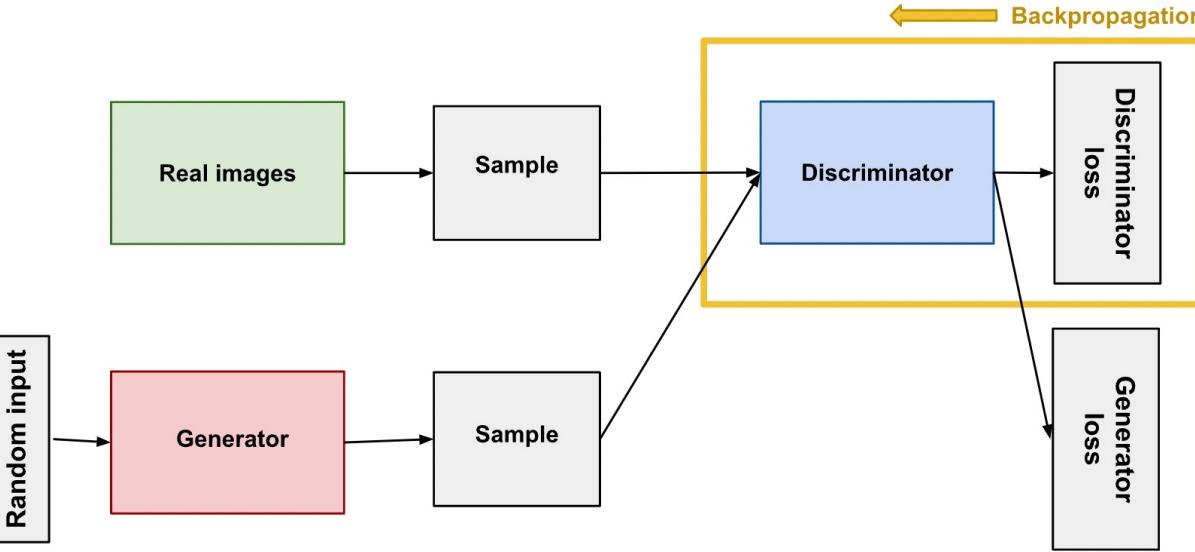


Figure 1: Backpropagation in discriminator training.

# Summary

- Masing-masing sebutkan konsep utama yang diingat

# Tuhan Memberkati

