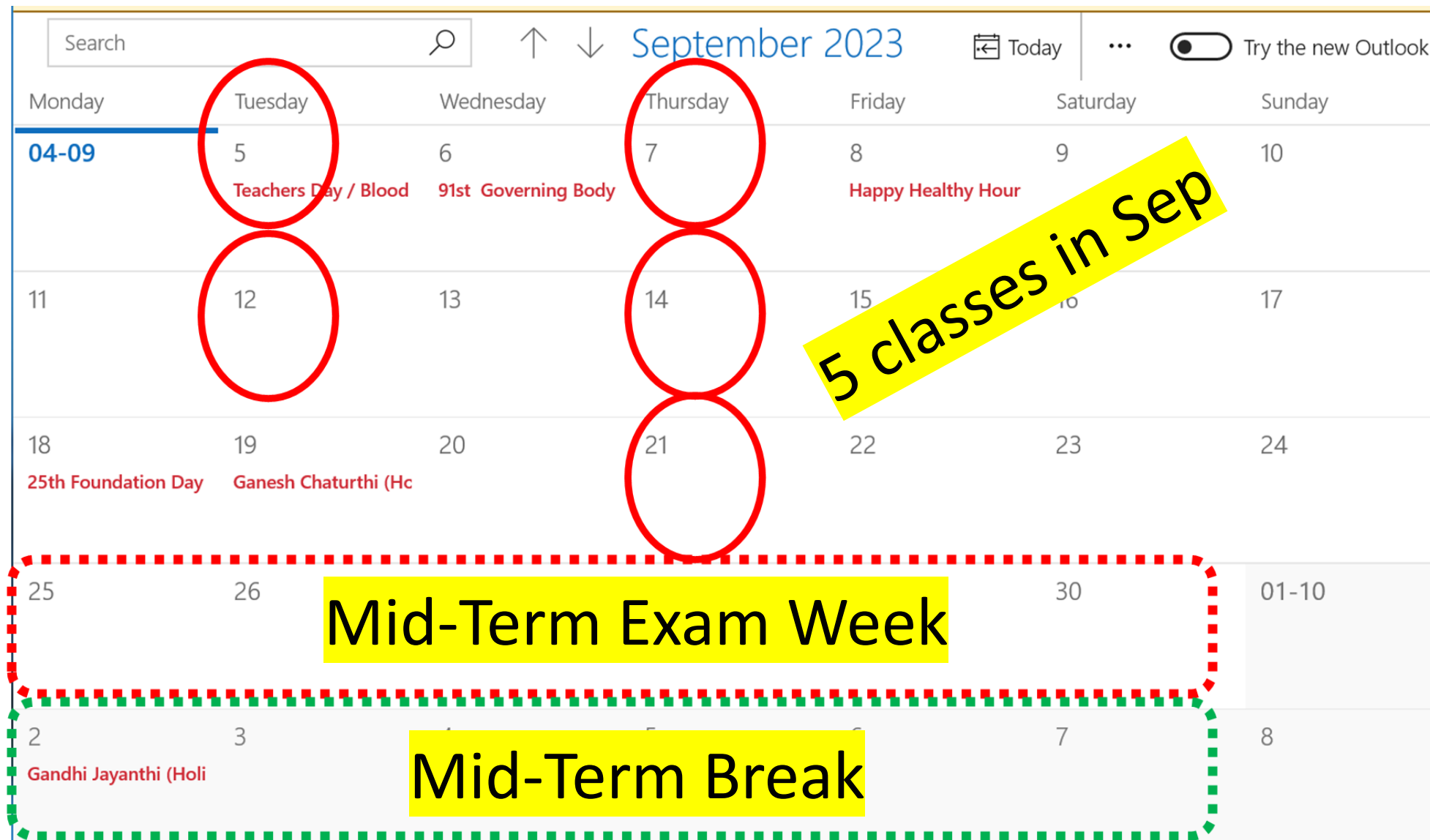


Pretext tasks

1. Predicting Rotations

September Schedule



1

2

SELF-PREDICTION

INNATE RELATIONSHIP
(Context-based)

1. ROTATION

2. RELATIVE POSITION

IMAGE

3

CONTRASTIVE LEARNING

INTER-SAMPLE
CLASSIFICATION

- 1. Instance Discrimination
- 2. SimCLR [Contrastive Loss]
- 3. Theory – Guarantees / Bounds

IMAGE

4

CONTRASTIVE LEARNING

INTER-SAMPLE
CLASSIFICATION

Contrastive Predictive Coding
(CPC), [NCE, InfoNCE Loss]

AUDIO/
SPEECH

5

SELF-PREDICTION

GENERATIVE
(VAE)

- 1. AE – Variational Bayes
- 2. VQ-VAE + AR

IMAGE

AUDIO/
SPEECH

6

SELF-PREDICTION

GENERATIVE
(AR)

- 1. AR-LM – GPT
- 2. Masked-LM – BERT

LANGUAGE

7

SELF-PREDICTION

MASKED-GEN
(Masked LM for ASR)

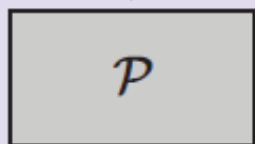
- 1. Wav2Vec / 2.0
- 2. HuBERT

AUDIO/
SPEECH

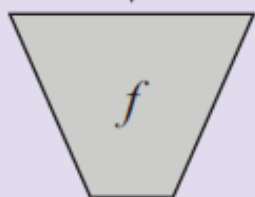
Self-Supervised

Unlabeled
Data Set

x



x, z



\hat{z}

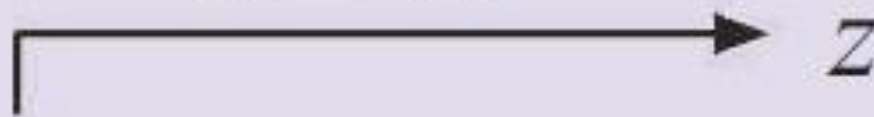


\mathcal{L}



Transformation Prediction

$z = 90^\circ$



z

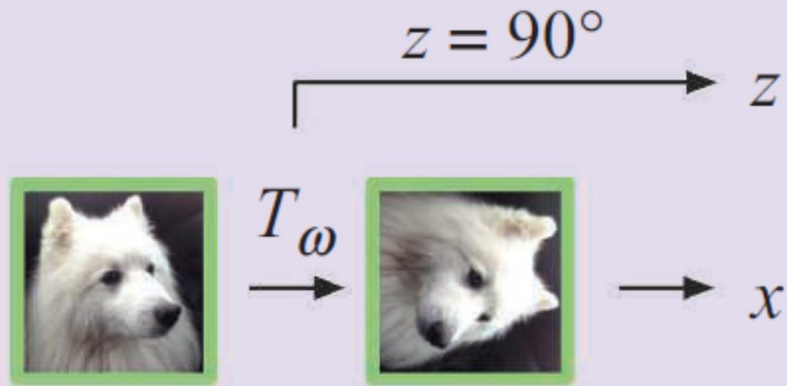


T_ω



x

Transformation Prediction



Algorithm 2. The pseudolabel generation process \mathcal{P} for TP.

Input: Unlabeled data set $D_s = \{x_i^{(s)}\}_{i=1}^M$.

for i from 1 to M **do**

 Sample $\omega \sim \Omega$

$x_i \leftarrow T_\omega(x_i^{(s)})$

 ▷ Apply transformation to raw input

$z_i \leftarrow \omega$

end for

Output: $\{x_i, z_i\}_{i=1}^M$.

$$\theta^* = \operatorname{argmin}_{\theta, \gamma} \sum_{(x_i, z_i) \in \mathcal{P}(D_s)} \mathcal{L}_{CE}(k_\gamma(h_\theta(x_i)), z_i).$$

UNSUPERVISED REPRESENTATION LEARNING BY PREDICTING IMAGE ROTATIONS

Spyros Gidaris, Praveer Singh, Nikos Komodakis

University Paris-Est, LIGM

Ecole des Ponts ParisTech

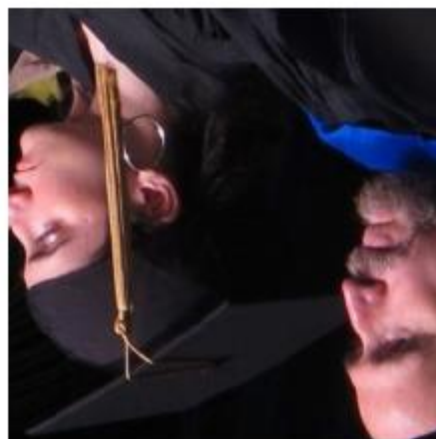
`{spyros.gidaris,praveer.singh,nikos.komodakis}@enpc.fr`



90° rotation



270° rotation



180° rotation

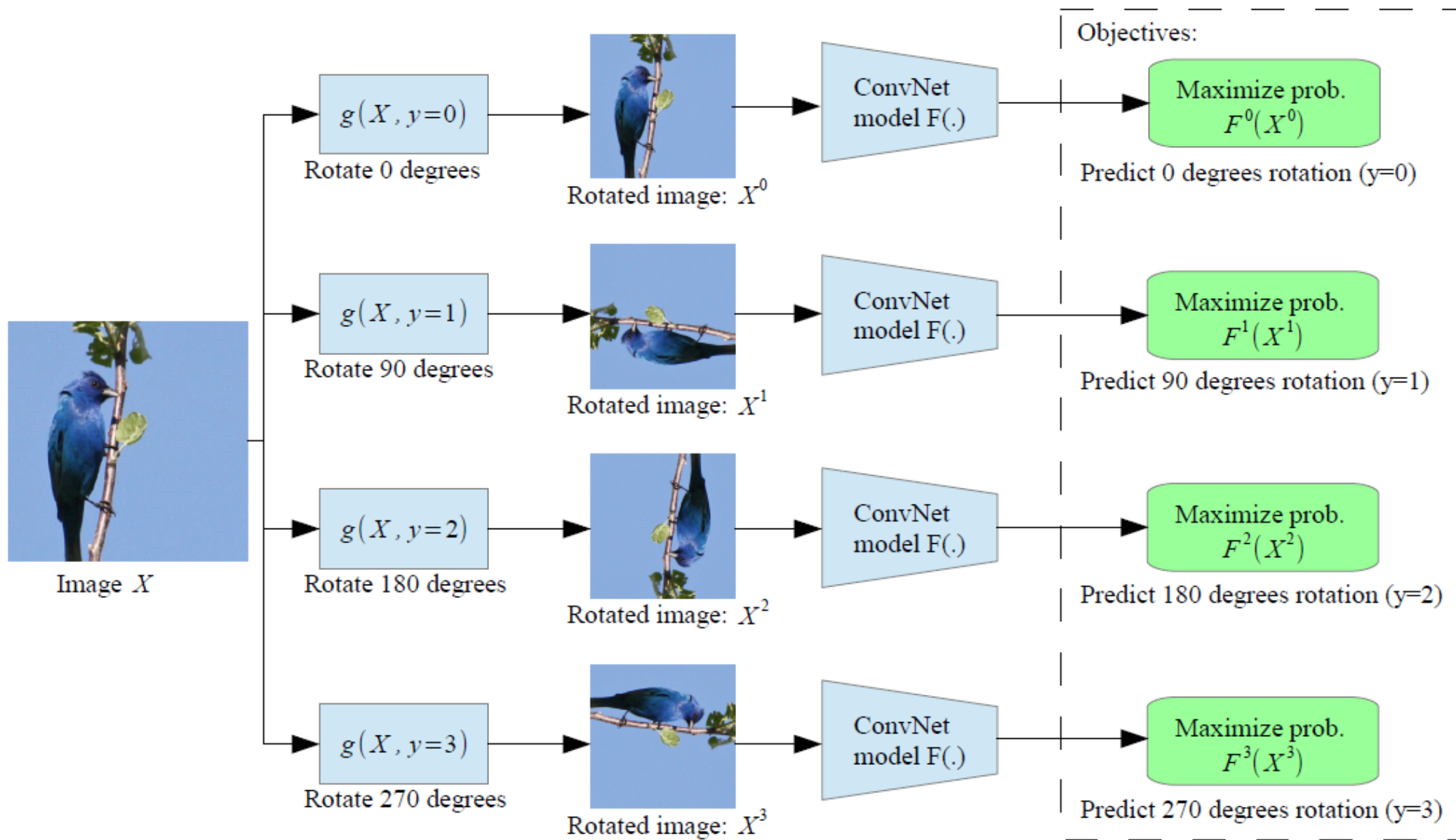


0° rotation



270° rotation

Figure 1: Images rotated by random multiples of 90 degrees (e.g., 0, 90, 180, or 270 degrees). The core intuition of our self-supervised feature learning approach is that if someone is not aware of the concepts of the objects depicted in the images, he cannot recognize the rotation that was applied to them.





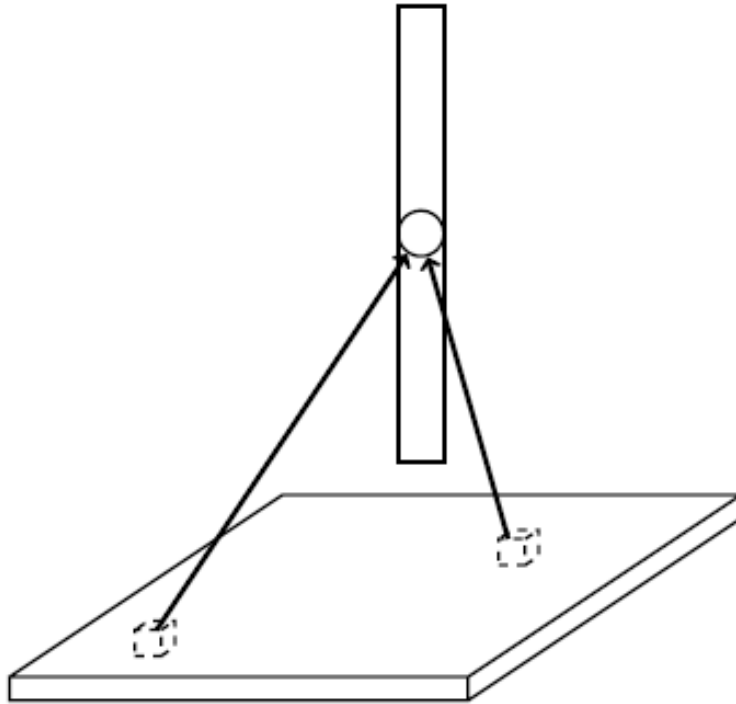
08	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	81	88
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	48	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	58	88	30	03	49	13	36	65
52	70	95	23	04	60	11	42	62	24	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	67	85	59	41	92	36	54	22	40	40	28	66	33	13	80
24	47	38	60	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
32	98	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	43	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
88	46	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	58	25	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	62	99	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	45	88	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	17	62	48

What the computer sees

image classification

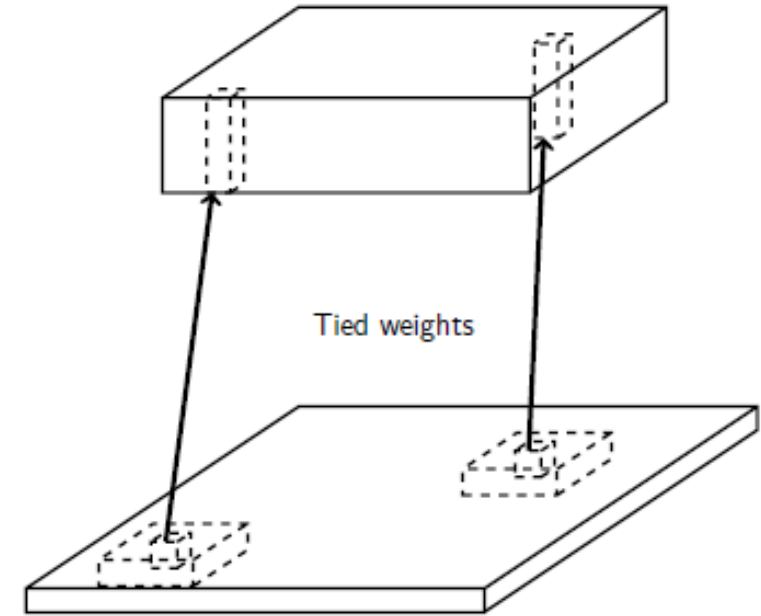
82% cat
15% dog
2% hat
1% mug

Fully connected



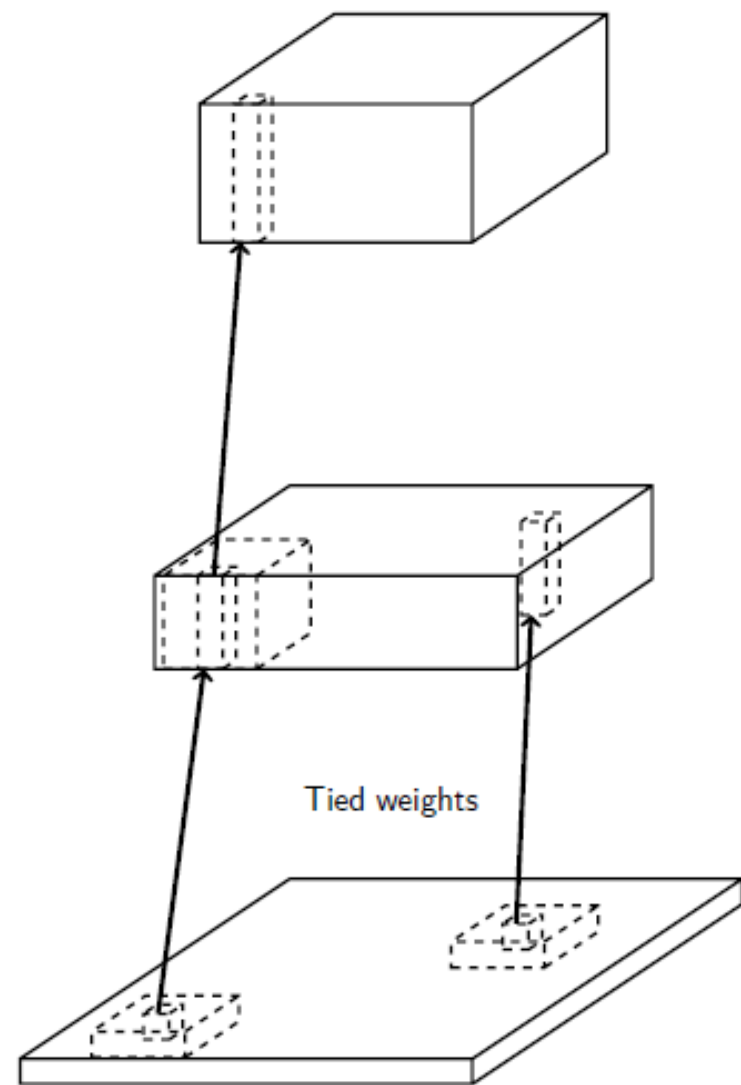
Each hidden unit looks at the entire image.

Convolutional



Each column of hidden units looks at a different patch of input.

Stack multiple layers of convolutions



2-D Convolution

What does this convolution kernel do?



*

0	1	0
1	4	1
0	1	0



What does this convolution kernel do?



*

0	-1	0
-1	8	-1
0	-1	0



What does this convolution kernel do?

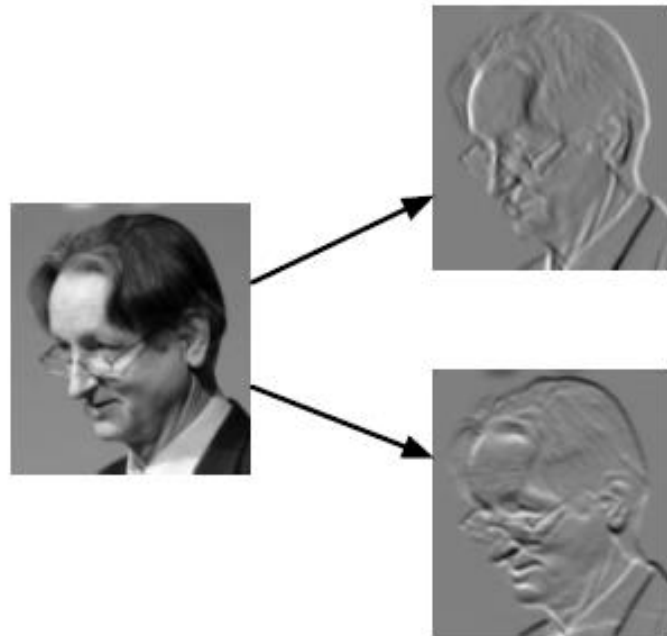


*

0	-1	0
-1	4	-1
0	-1	0

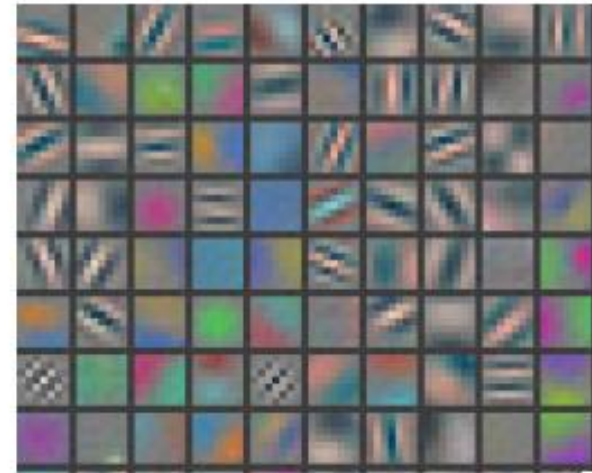


The convolution layer has a set of filters. Its output is a set of **feature maps**, each one obtained by convolving the image with a filter.

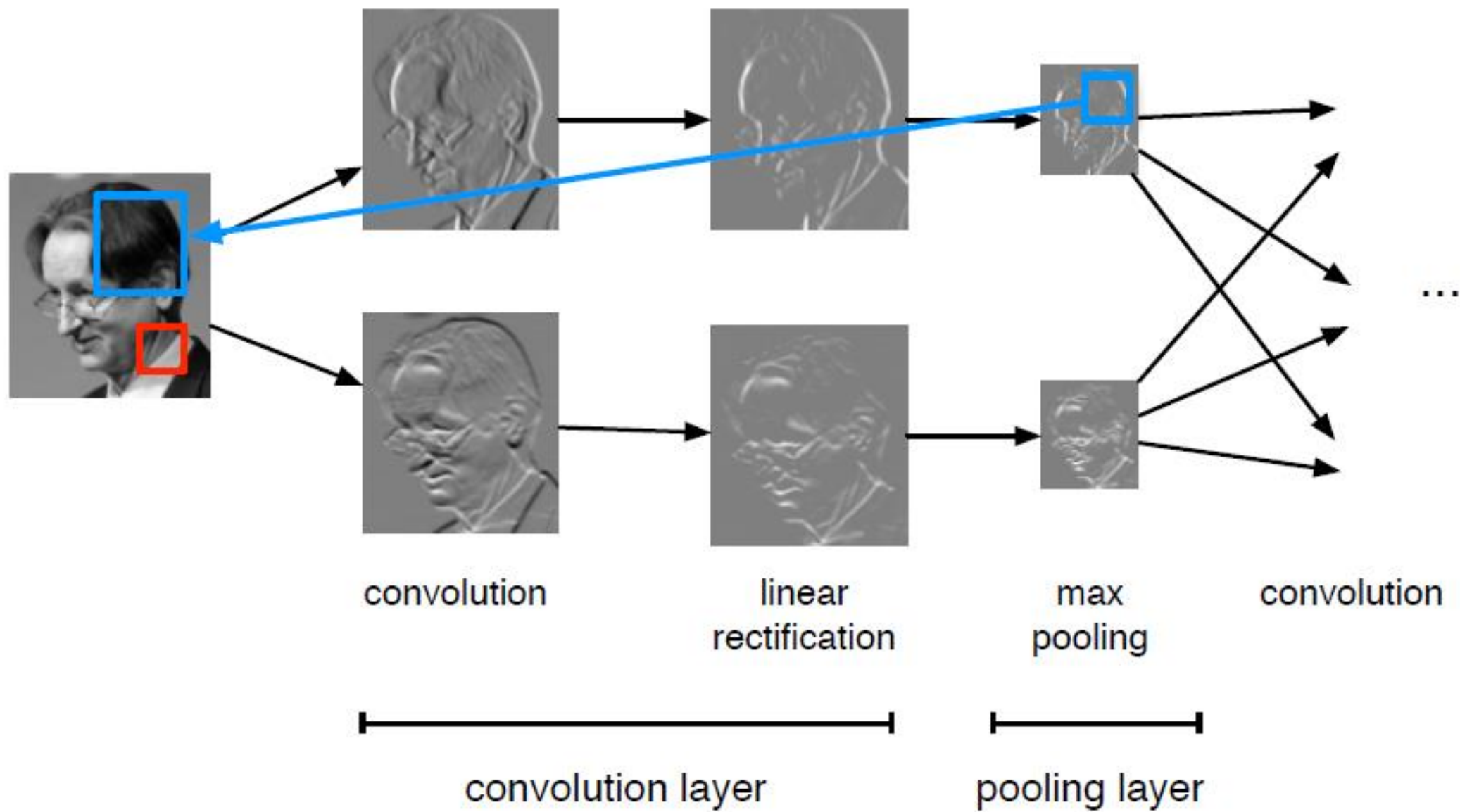


convolution

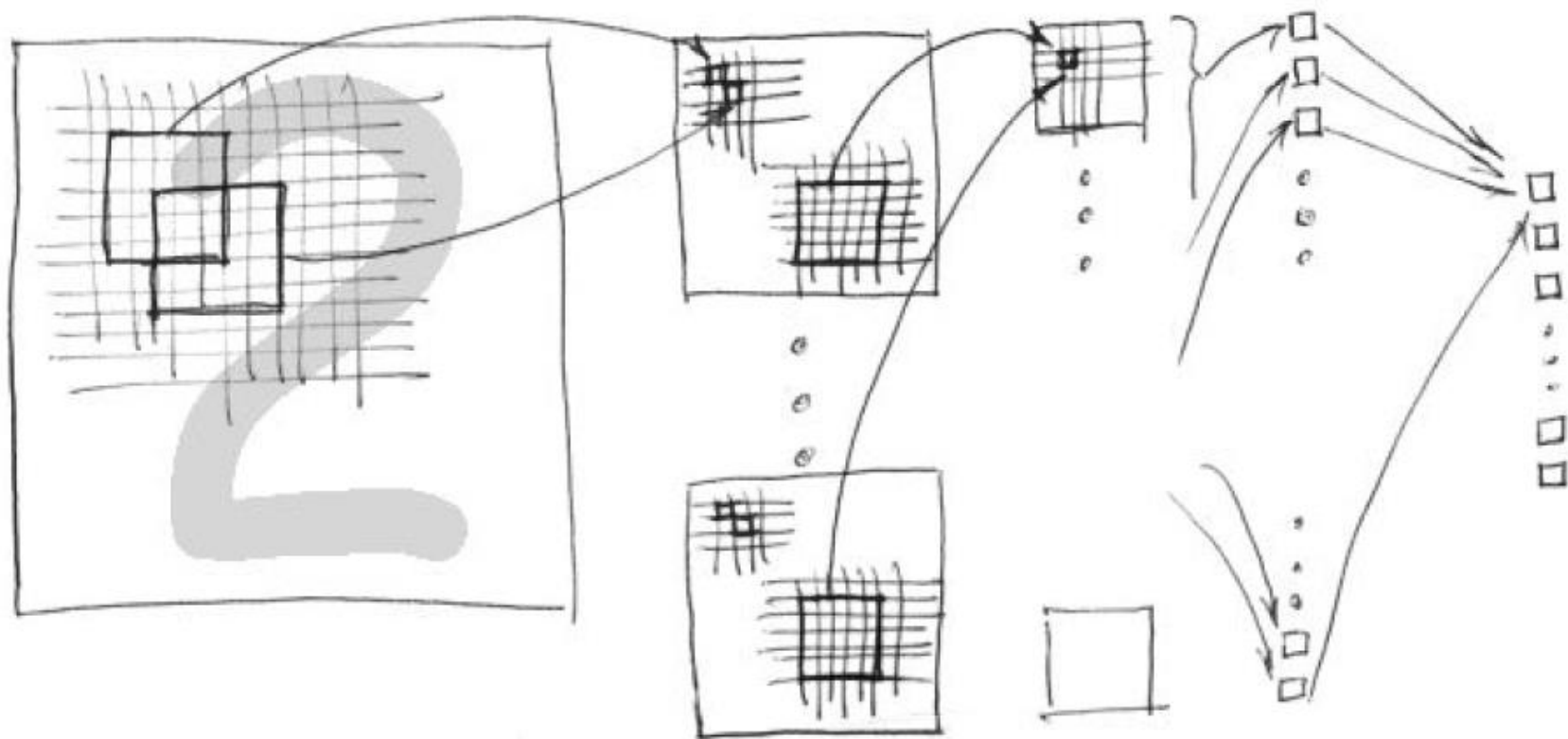
Example first-layer filters



(Zeiler and Fergus, 2013, Visualizing and understanding
convolutional networks)



Convolutional Neural Network



LeNet

LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W. and Jackel, L.D., 1989. **Backpropagation applied to handwritten zip code recognition.** *Neural computation*, 1(4), pp.541-551.

LeNet

10 output units

layer H3

30 hidden units

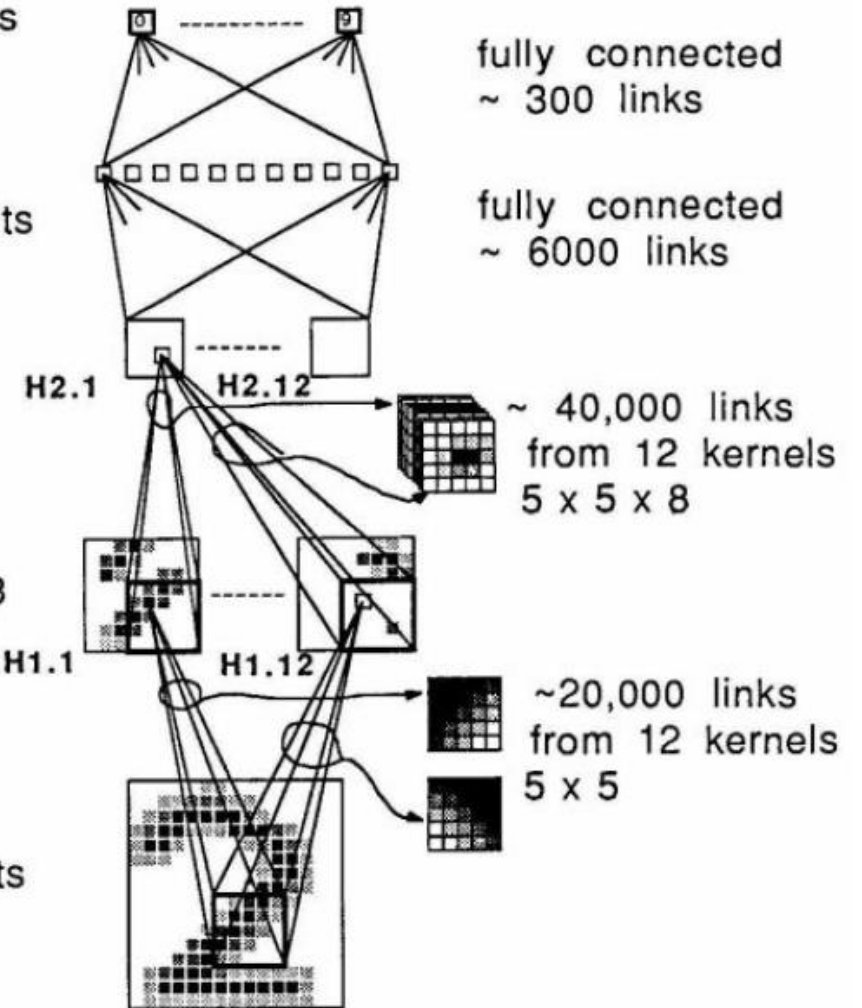
layer H2

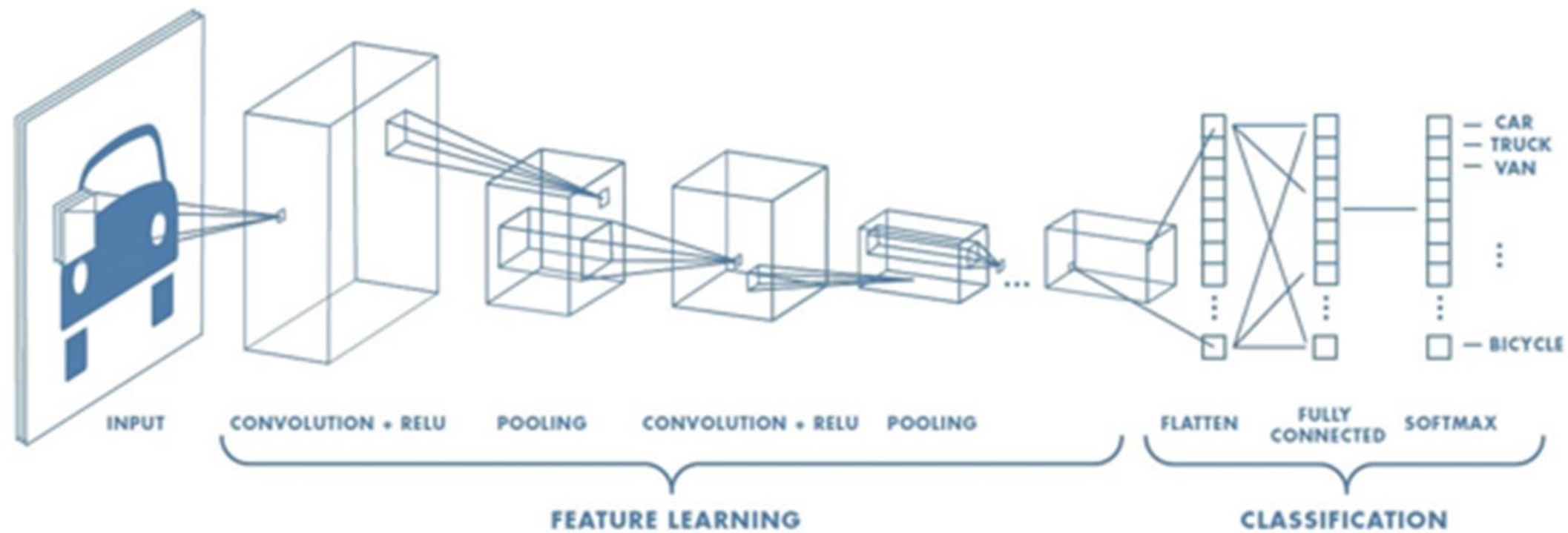
$12 \times 16 = 192$
hidden units

layer H1

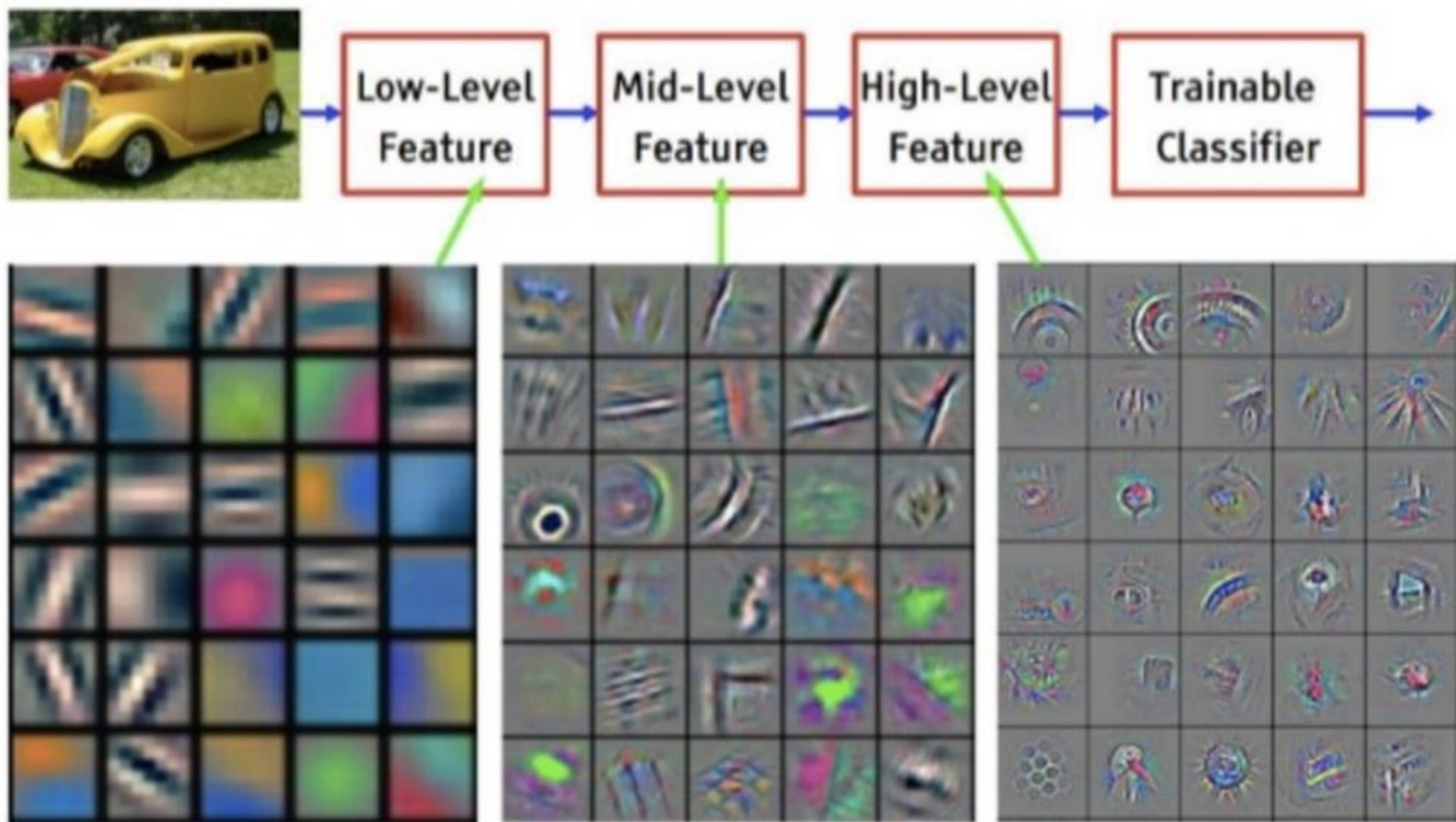
$12 \times 64 = 768$
hidden units

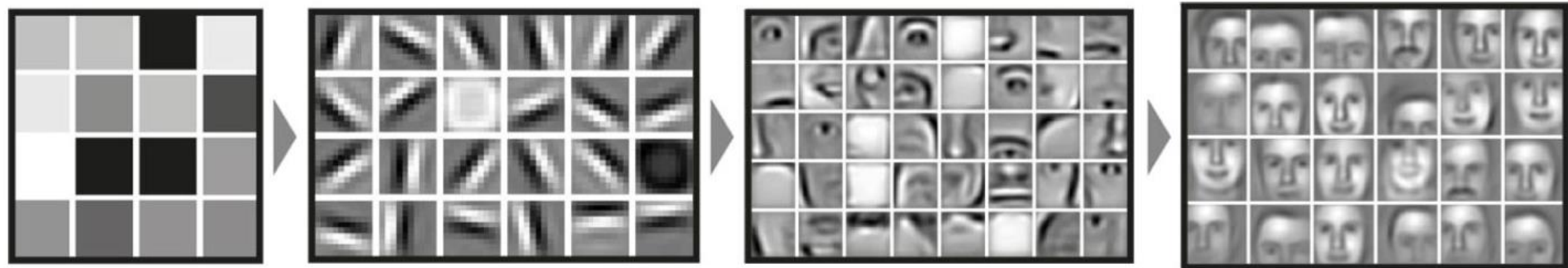
256 input units



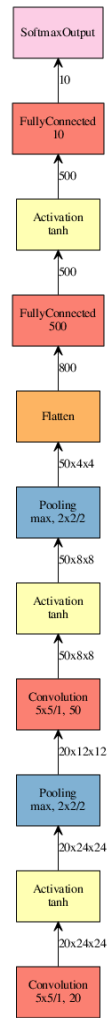


Standard architecture of a CNN





“LeNet”
1998:



“AlexNet”
2012:



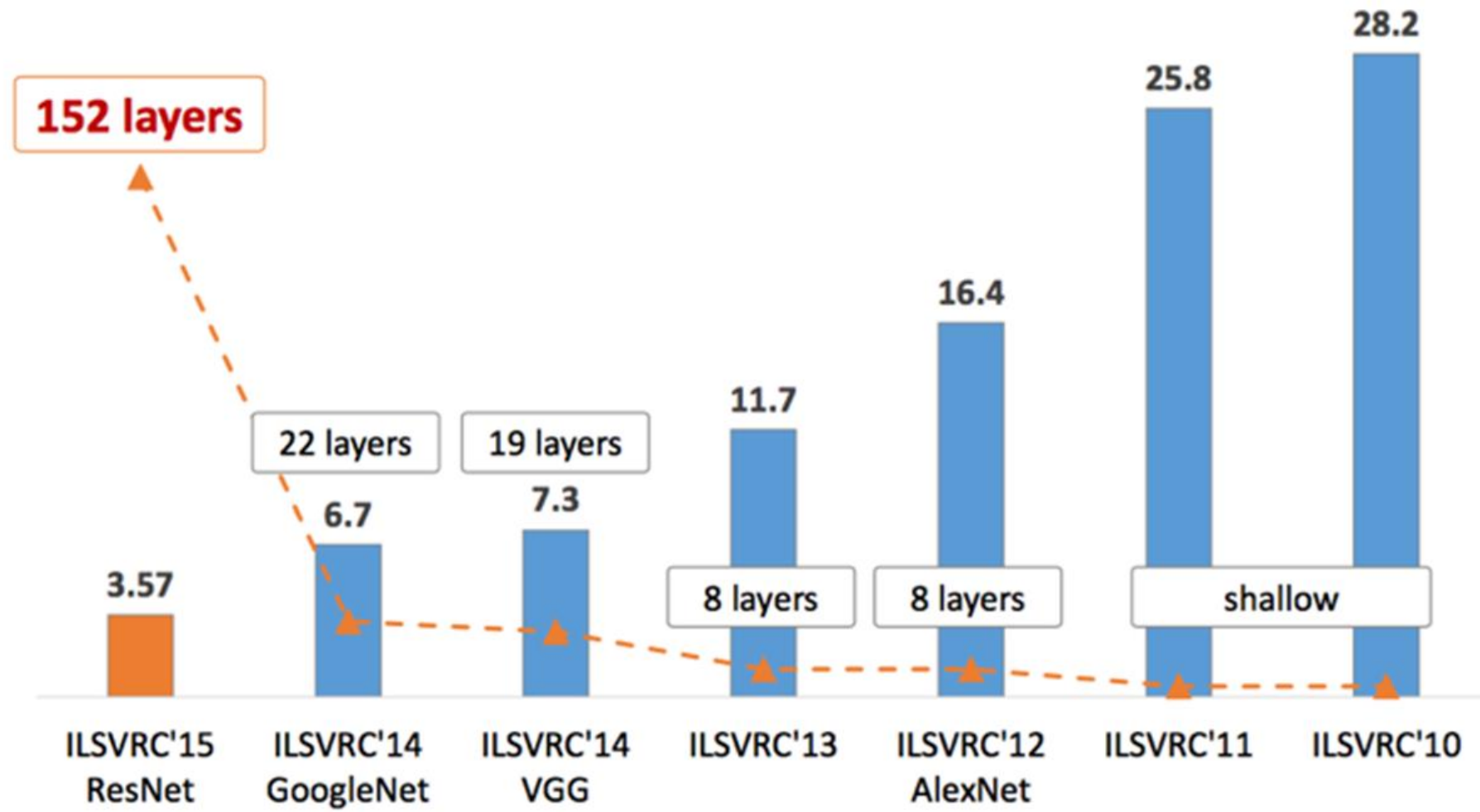
“GoogLeNet”
2014:



“Inception v3”
2015:



ImageNet Large Scale Visual Recognition Challenge (ILSVRC)



Back to...

Published as a conference paper at ICLR 2018

UNSUPERVISED REPRESENTATION LEARNING BY PREDICTING IMAGE ROTATIONS

Spyros Gidaris, Praveer Singh, Nikos Komodakis

University Paris-Est, LIGM

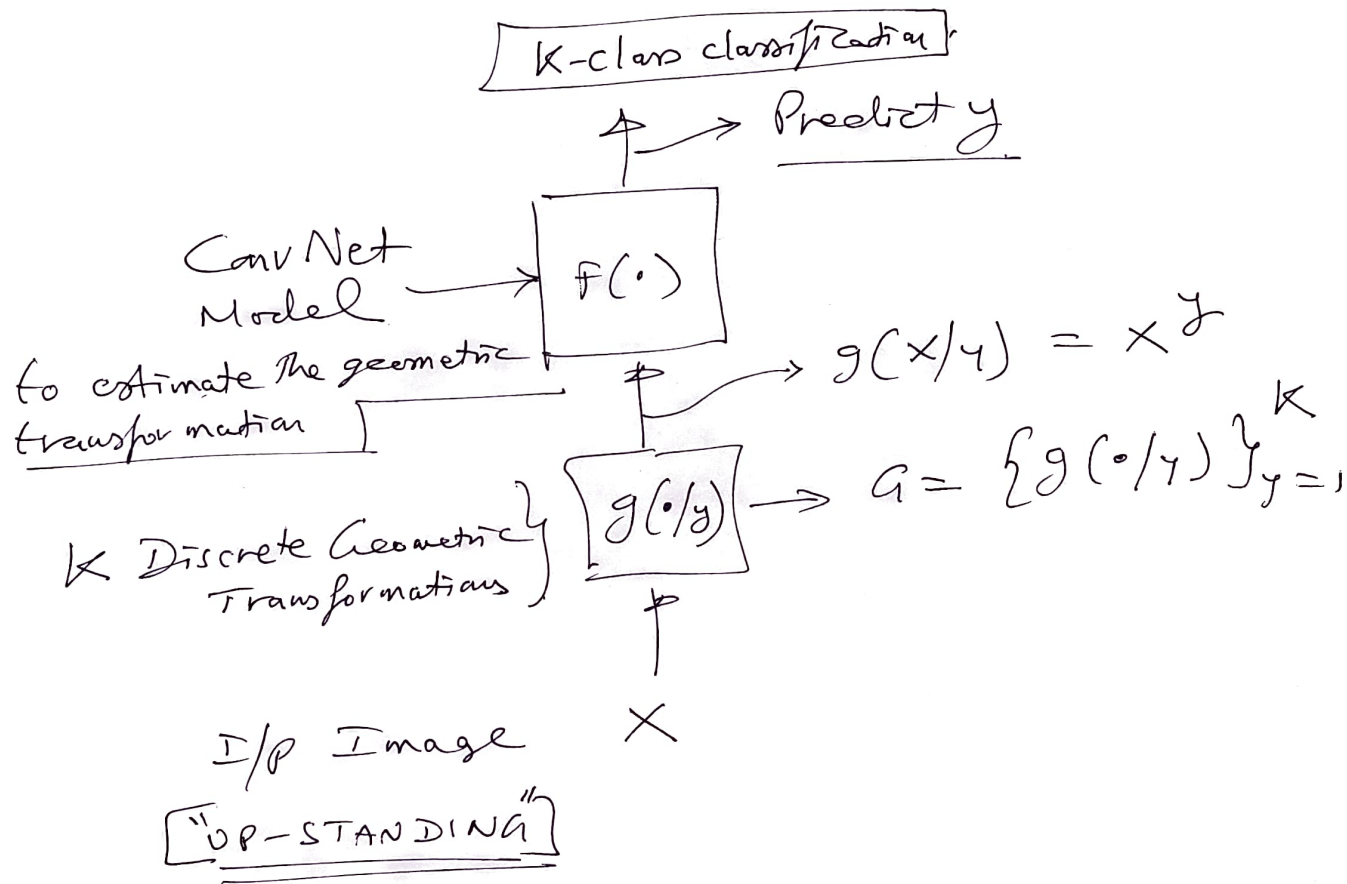
Ecole des Ponts ParisTech

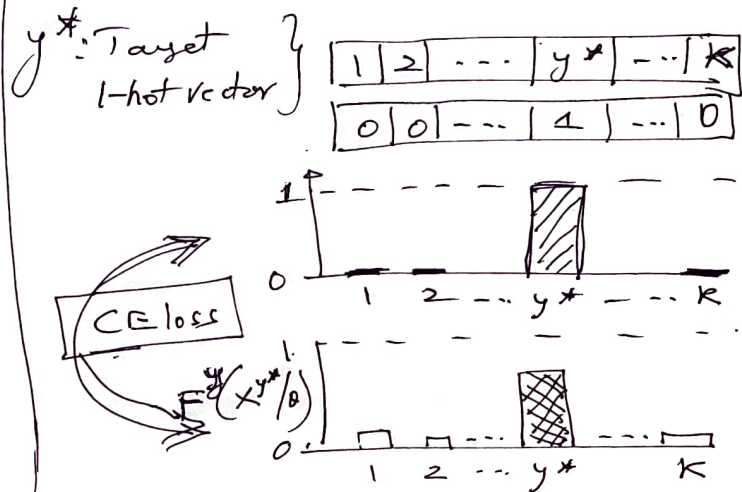
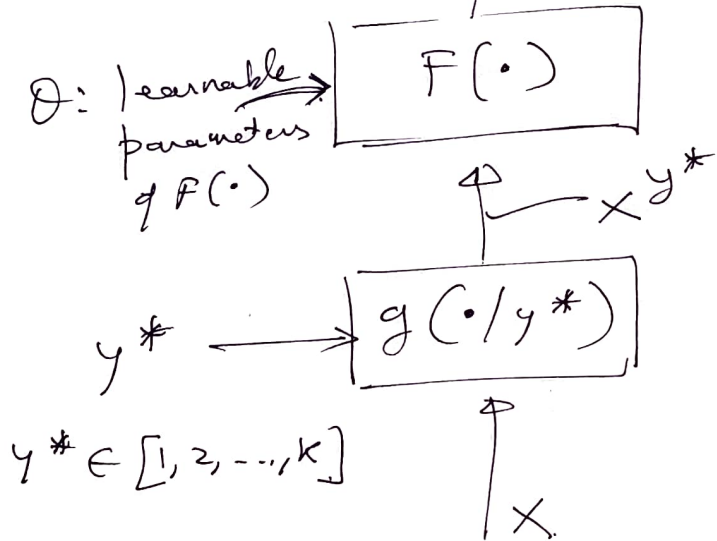
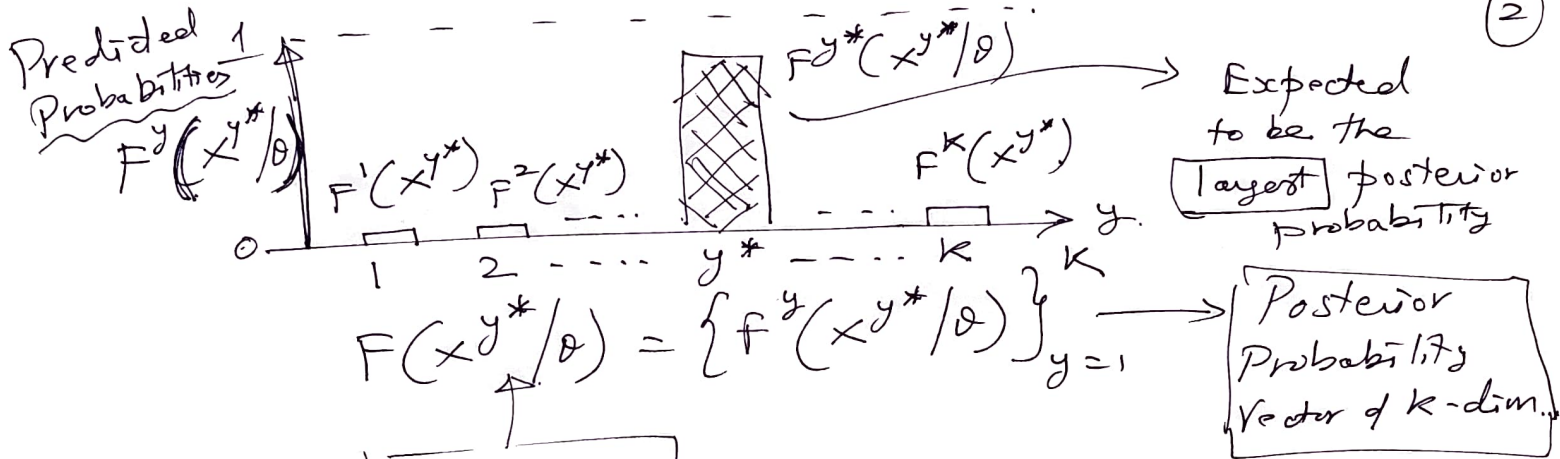
`{spyros.gidaris, praveer.singh, nikos.komodakis}@enpc.fr`

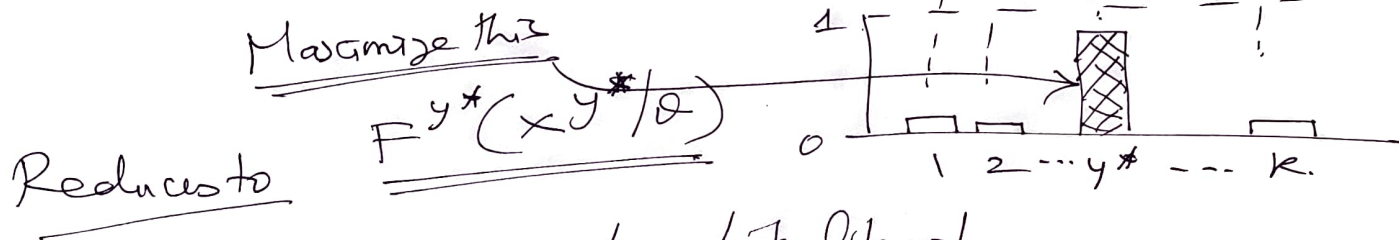
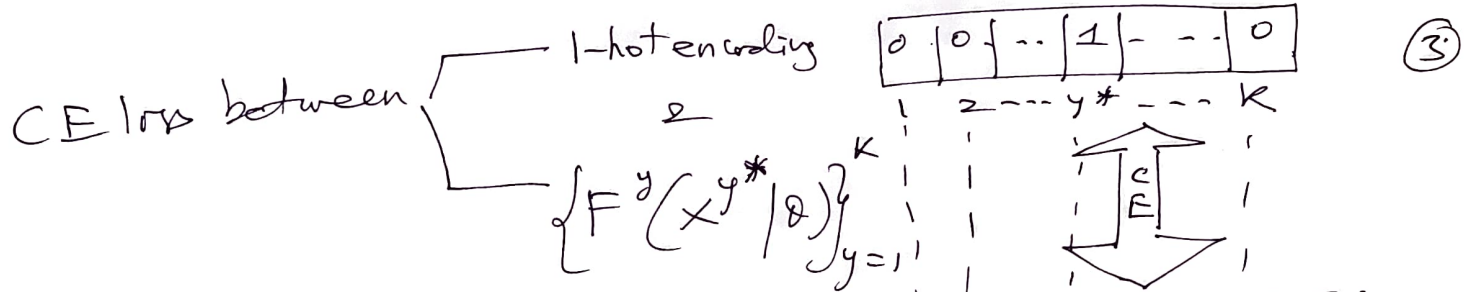
ICLR-2018

PREDICTING IMAGE ROTATIONS

①







NLL: Negative Log Likelihood

$$\text{Loss} = -\log(F^{y*}(g(x/y^*)/\theta))$$

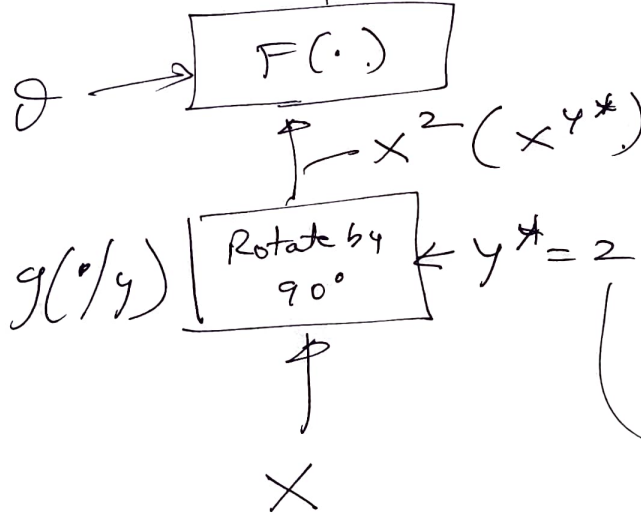
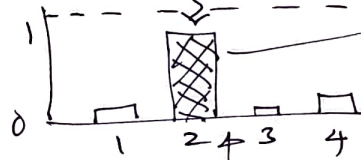
$K=4$ $y^*=2$

$\text{Rot}(x, \phi)$: operator that rotates image by $0, 90, 180, 270^\circ$

x ϕ

$$G = \{g(x/y)\}_{y=1}^4 \parallel g(x/y) = \text{Rot}(x, (y-1)90)$$

1	2	3	4
0	1	0	0



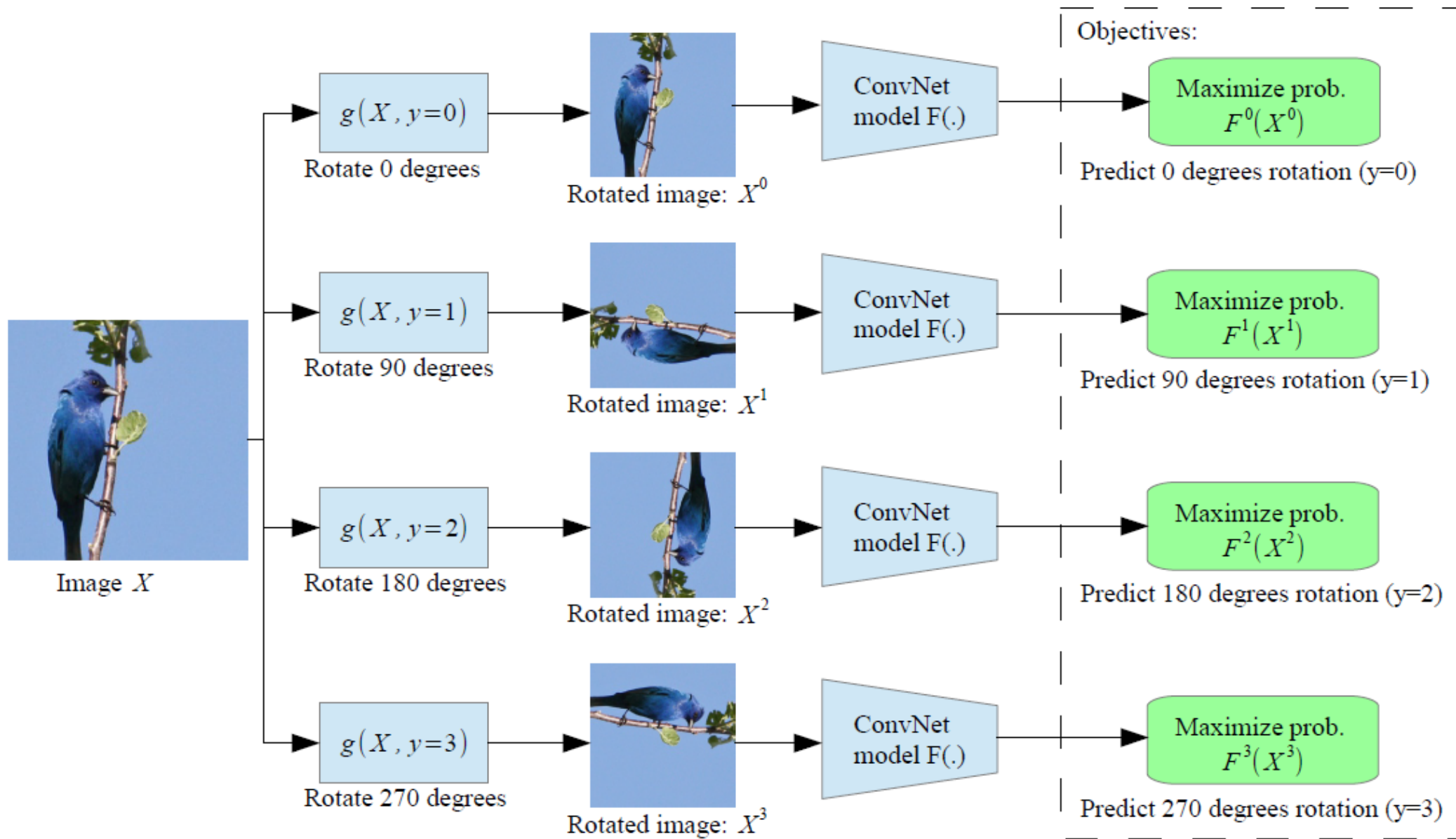
Maximize this
 $\log(F^2(x^2/\theta))$

Minimize loss [NLL]

$$-\log(F^2(x^2/\theta))$$

1	2	3	4
0	90	180	270

(4)



On a given unlabelled/unsupervised SSL T4 Set (5)

$$D = \{x_i\}_{i=1}^N \quad \text{1 } N \text{ samples}$$

Find ConvNet Model θ^*

$$\theta^* = \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \text{loss}(x_i, \theta)$$

with associated loss
objective $\min_{\theta} \text{Loss} = \min_{\theta} \frac{1}{N} \sum_{i=1}^N \text{loss}(x_i, \theta)$

fn

where

$$\text{Loss}(x_i, \theta) =$$

$$-\frac{1}{K} \sum_{y=1}^K \log(F^y(g(x_i/y)/\theta))$$

K=4

y	$y-1$	$g(x/y)$	$y=1$
$y=1$	0	$\text{Rot}(x, 0.90)$	$\rightarrow 0^\circ$
$y=2$	1	$\text{Rot}(x, 1.90)$	$\rightarrow 90^\circ$
$y=3$	2	$\text{Rot}(x, 2.90)$	$\rightarrow 180^\circ$
$y=4$	3	$\text{Rot}(x, 3.90)$	$\rightarrow 270^\circ$

⑥

$$\text{Loss}(x_i, \theta) =$$

$$-\frac{1}{K} \sum_{y=1}^K \log(F^y(g(x_i/\gamma) | \theta))$$

For a $x_i \in \mathcal{D} = \{x_i\}_{i=1}^N$

— apply all K geometric transformations [Rotations]
 $K=4 \Rightarrow 0, 90, 180, 270$ $\nearrow K$

— Create $\{x_i^y\} = g(x_i/\gamma) \quad \forall y=1, 2, 3, 4$

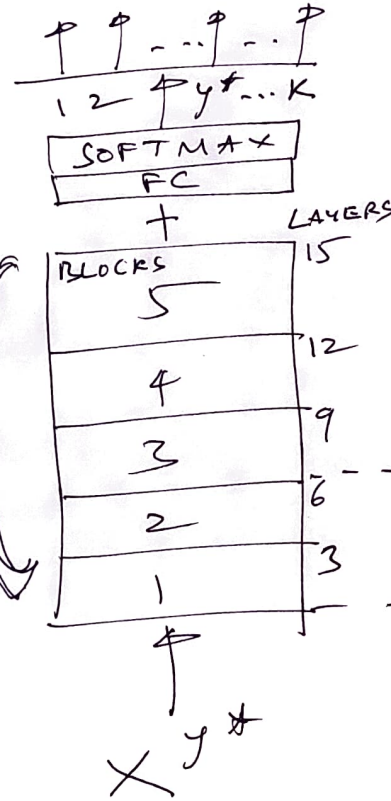
— Compute Loss $(x_i, \theta) \quad \forall \{x_i^y\}_{y=1}^K$

PRETRAINING

K-class
Pretext Task



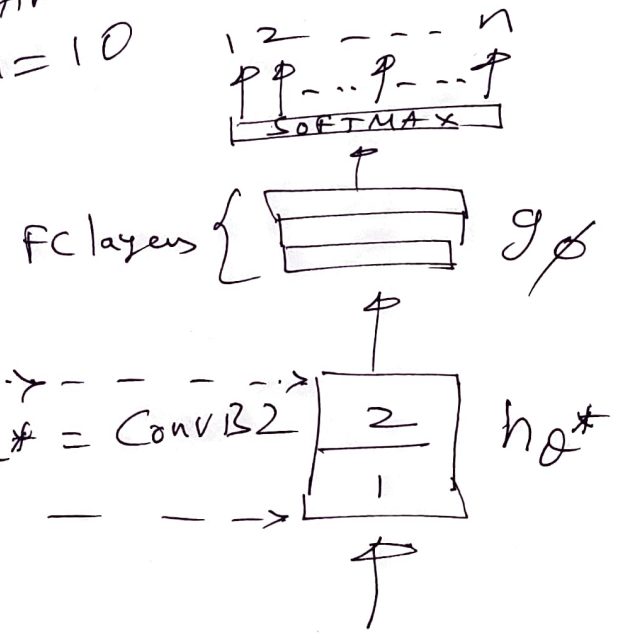
K-class classifi



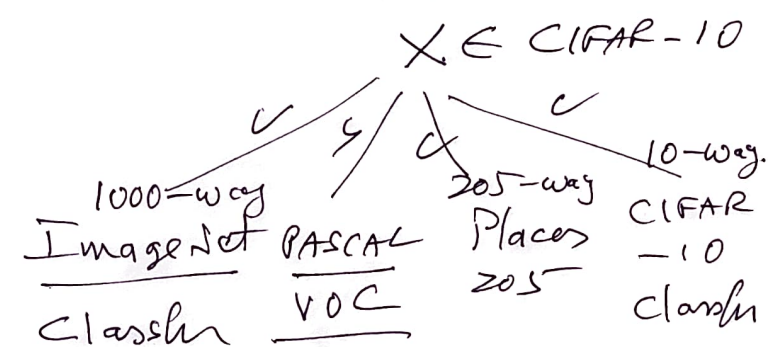
DOWNSTREAM

e.g
CIFAR-10
 $n=10$

n-class classifi



See TABLE-1





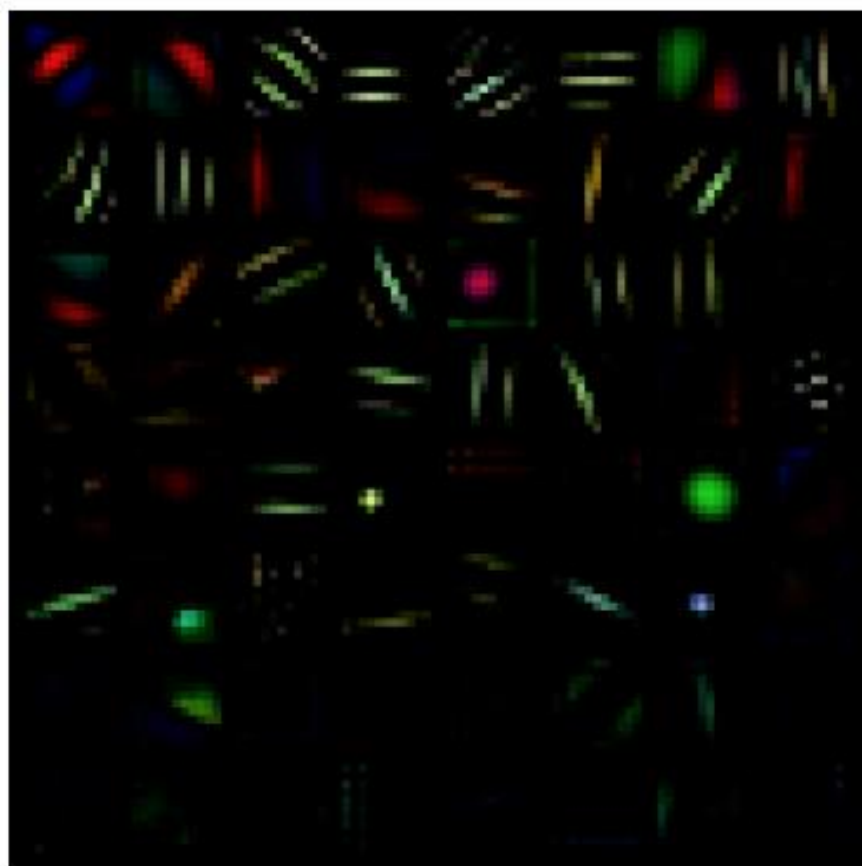
Input images on the models



(a) Attention maps of supervised model

(b) Attention maps of our self-supervised model

Figure 3: Attention maps generated by an AlexNet model trained (a) to recognize objects (supervised), and (b) to recognize image rotations (self-supervised). In order to generate the attention map of a conv. layer we first compute the feature maps of this layer, then we raise each feature activation on the power p , and finally we sum the activations at each location of the feature map. For the conv. layers 1, 2, and 3 we used the powers $p = 1$, $p = 2$, and $p = 4$ respectively. For visualization of our self-supervised model’s attention maps for all the rotated versions of the images see Figure 6 in appendix A.



(a) Supervised



(b) Self-supervised to recognize rotations

Figure 4: First layer filters learned by a AlexNet model trained on (a) the supervised object recognition task and (b) the self-supervised task of recognizing rotated images. We observe that the filters learned by the self-supervised task are mostly oriented edge filters on various frequencies and, remarkably, they seem to have more variety than those learned on the supervised task.

Table 1: Evaluation of the unsupervised learned features by measuring the classification accuracy that they achieve when we train a non-linear object classifier on top of them. The reported results are from CIFAR-10. The size of the ConvB1 feature maps is $96 \times 16 \times 16$ and the size of the rest feature maps is $192 \times 8 \times 8$.

Model	ConvB1	ConvB2	ConvB3	ConvB4	ConvB5
RotNet with 3 conv. blocks	85.45	88.26	62.09	-	-
RotNet with 4 conv. blocks	85.07	89.06	86.21	61.73	-
RotNet with 5 conv. blocks	85.04	89.76	86.82	74.50	50.37

Table 2: Exploring the quality of the self-supervised learned features w.r.t. the number of recognized rotations. For all the entries we trained a non-linear classifier with 3 fully connected layers (similar to Table 1) on top of the feature maps generated by the 2nd conv. block of a RotNet model with 4 conv. blocks in total. The reported results are from CIFAR-10.

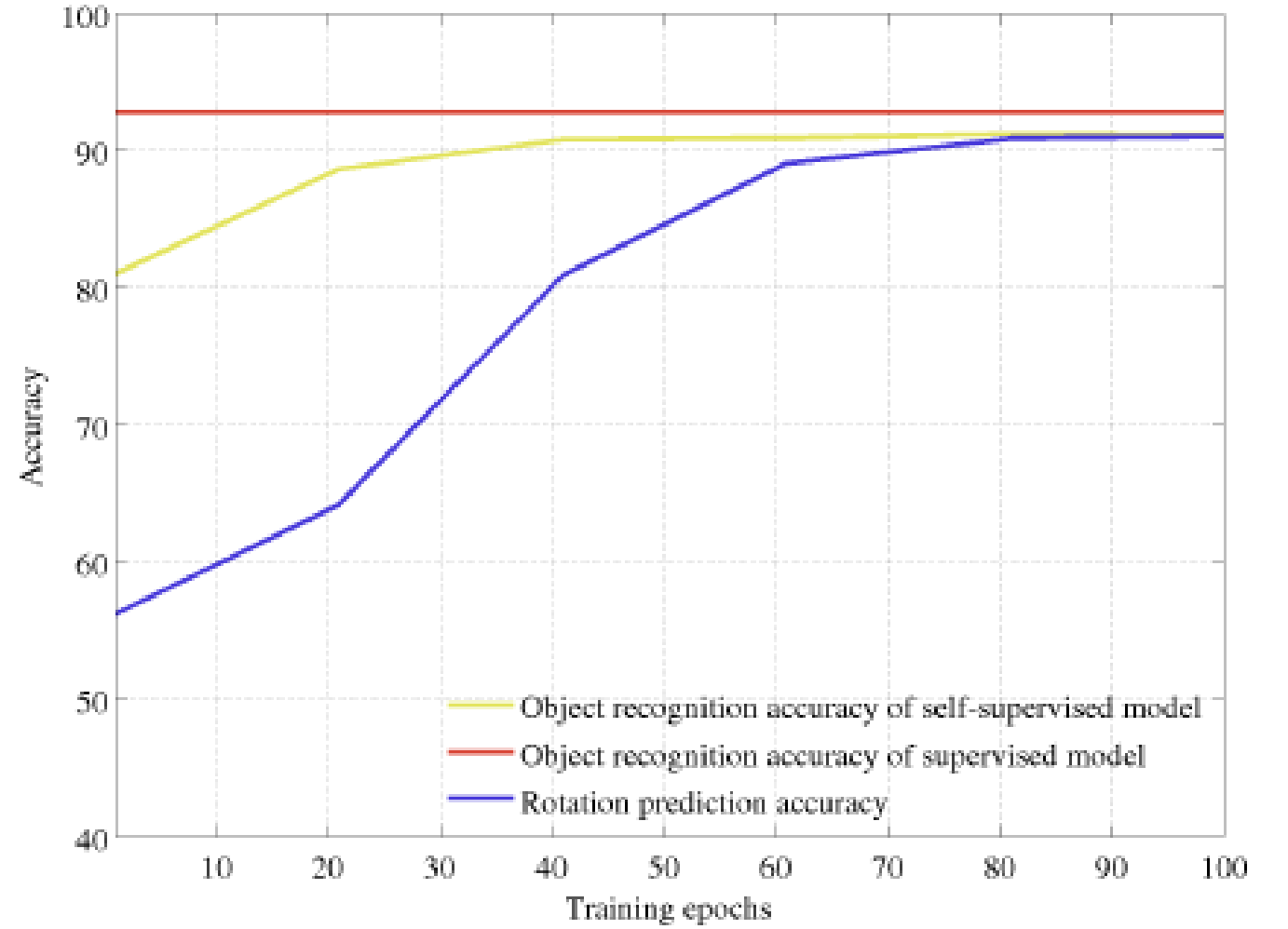
# Rotations	Rotations	CIFAR-10 Classification Accuracy
4	$0^\circ, 90^\circ, 180^\circ, 270^\circ$	89.06
8	$0^\circ, 45^\circ, 90^\circ, 135^\circ, 180^\circ, 225^\circ, 270^\circ, 315^\circ$	88.51
2	$0^\circ, 180^\circ$	87.46
2	$90^\circ, 270^\circ$	85.52

Table 3: Evaluation of unsupervised feature learning methods on CIFAR-10. The *Supervised NIN* and the *(Ours) RotNet + conv* entries have exactly the same architecture but the first was trained fully supervised while on the second the first 2 conv. blocks were trained unsupervised with our rotation prediction task and the 3rd block only was trained in a supervised manner. In the *Random Init. + conv* entry a conv. classifier (similar to that of *(Ours) RotNet + conv*) is trained on top of two NIN conv. blocks that are randomly initialized and stay frozen. Note that each of the prior approaches has a different ConvNet architecture and thus the comparison with them is just indicative.

Method	Accuracy
Supervised NIN	92.80
Random Init. + conv	72.50
(Ours) RotNet + non-linear	89.06
(Ours) RotNet + conv	91.16
(Ours) RotNet + non-linear (fine-tuned)	91.73
(Ours) RotNet + conv (fine-tuned)	92.17
Roto-Scat + SVM Oyallon & Mallat (2015)	82.3
ExemplarCNN Dosovitskiy et al. (2014)	84.3
DCGAN Radford et al. (2015)	82.8
Scattering Oyallon et al. (2017)	84.7

Table 9: **Per class CIFAR-10 classification accuracy.**

Classes	aero	car	bird	cat	deer	dog	frog	horse	ship	truck
Supervised	93.7	96.3	89.4	82.4	93.6	89.7	95.0	94.3	95.7	95.2
(Ours) RotNet	91.7	95.8	87.1	83.5	91.5	85.3	94.2	91.9	95.7	94.2



(a)

Figure 5: (a) Plot with the rotation prediction accuracy and object recognition accuracy as a function of the training epochs used for solving the rotation prediction task. The red curve is the object recognition accuracy of a fully supervised model (a NIN model), which is independent from the training epochs on the rotation prediction task. The yellow curve is the object recognition accuracy of an object classifier trained on top of feature maps learned by a *RotNet* model at different snapshots of the training procedure. (b) Accuracy as a function of the number of training examples per category

VISUALIZING ATTENTION MAPS OF ROTATED IMAGES

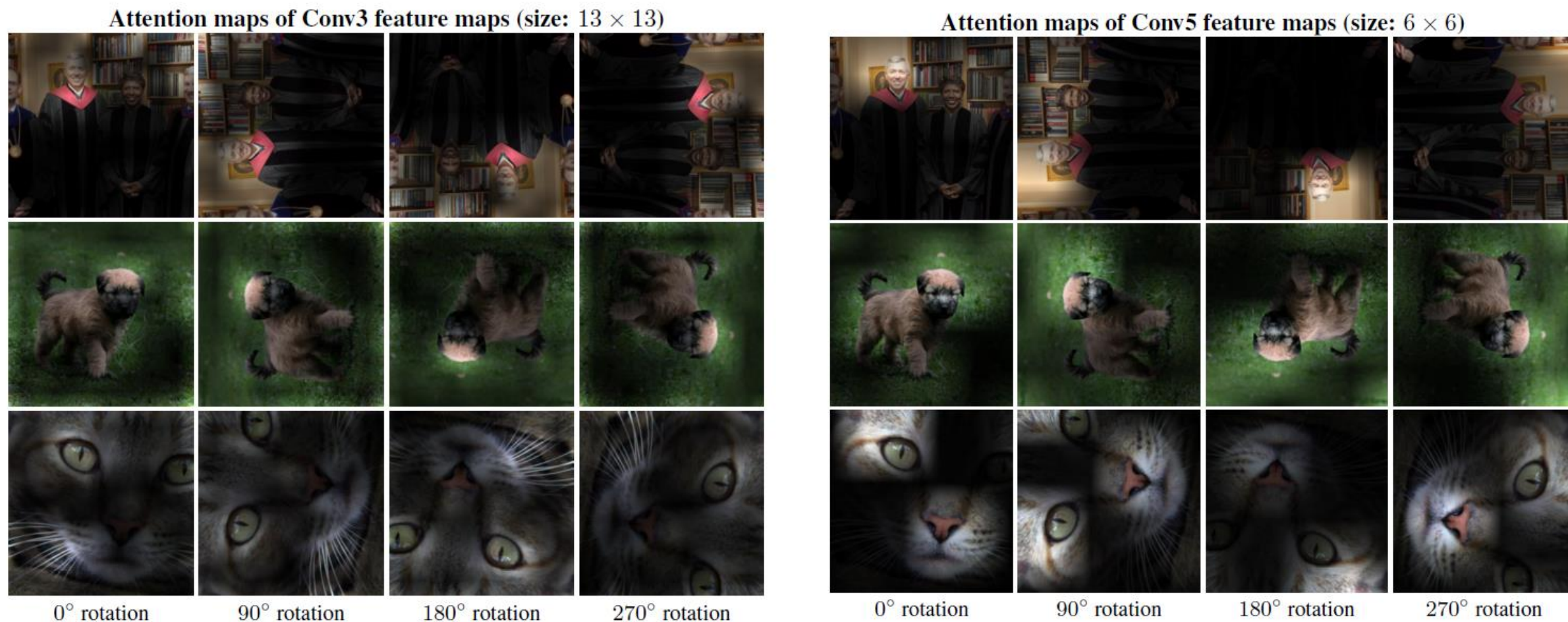


Figure 6: Attention maps of the Conv3 and Conv5 feature maps generated by an AlexNet model trained on the self-supervised task of recognizing image rotations. Here we present the attention maps generated for all the 4 rotated copies of an image.

Thank you !!