



Classification of ncRNAs using machine learning methods

Stefan Simm

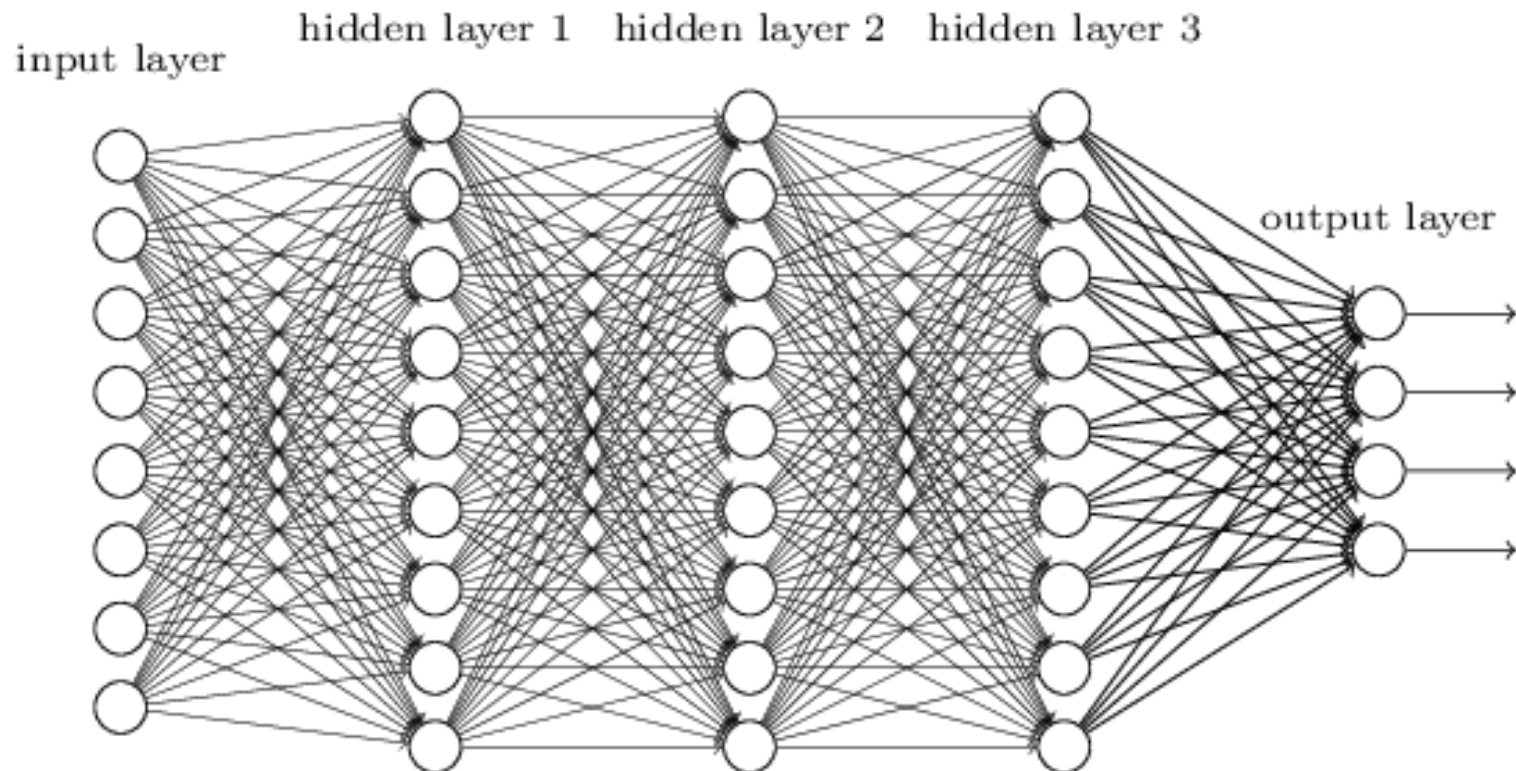
31. March 2023



CONVOLUTIONAL NEURAL NETWORK

ANN or CNN?

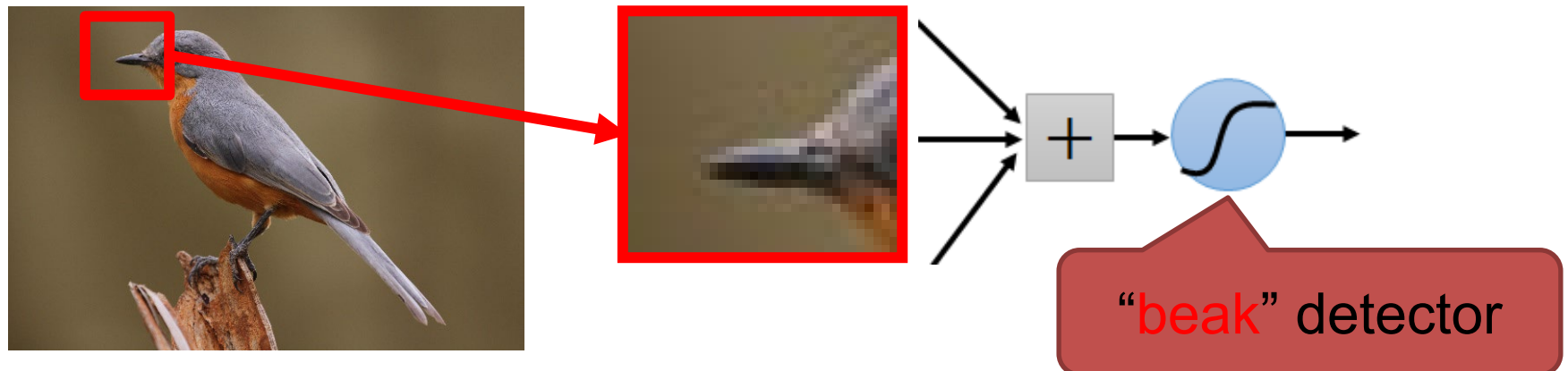
- We know it is better to train on small models if possible.
- From this fully connected model, do we really need all the edges?
- Can some of these be shared?



Consider learning from images:

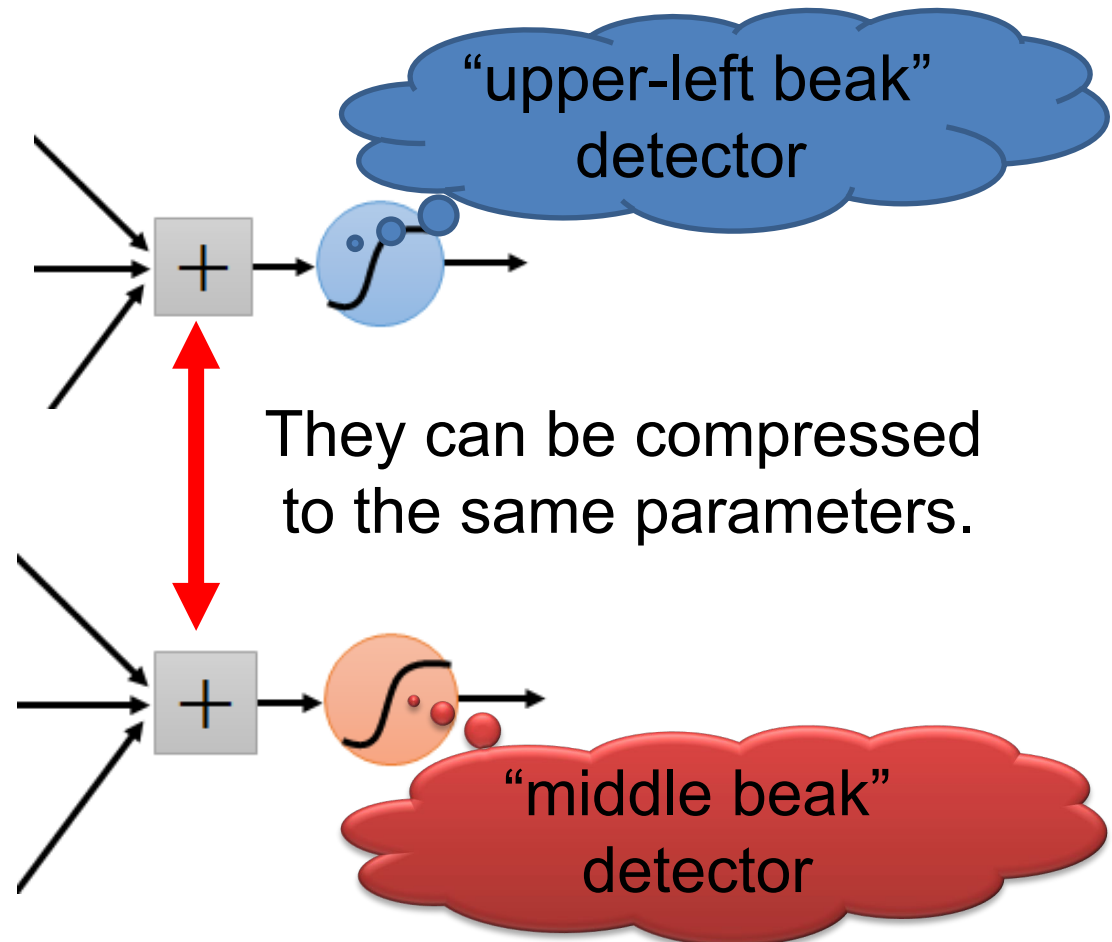
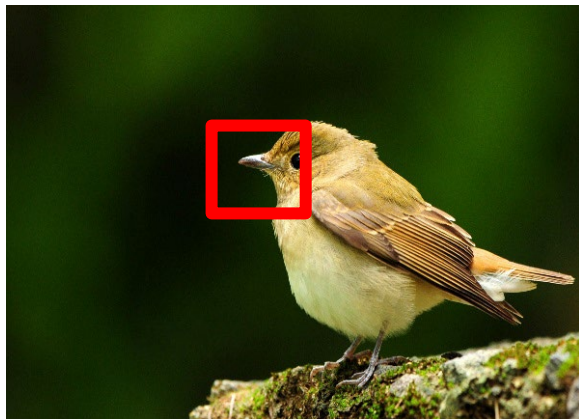
- Some patterns are much smaller than the whole image

Can represent a small region with fewer parameters



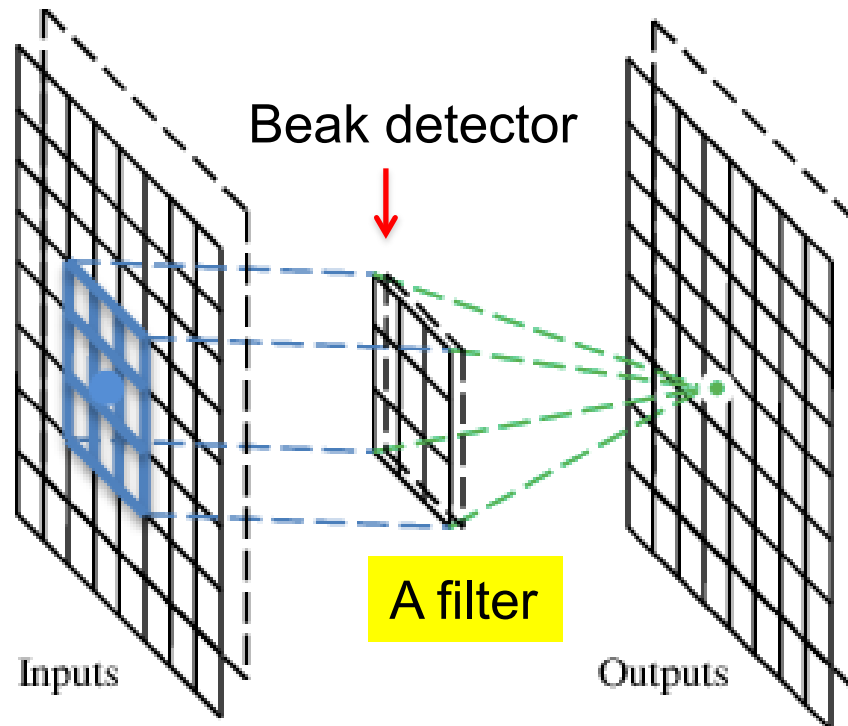
Same pattern appears in different places:
They can be compressed!

What about training a lot of such “small”
detectors
and each detector must “move around”.



A convolutional layer

A CNN is a neural network with some convolutional layers (and some other layers). A convolutional layer has a number of filters that does convolutional operation.



Convolution

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

These are the network parameters to be learned.

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

⋮ ⋮

Each filter detects a small pattern (3 x 3).

Convolution

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

Dot
product



3

-1

.Convolution

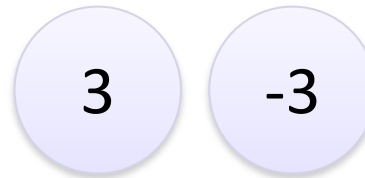
If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

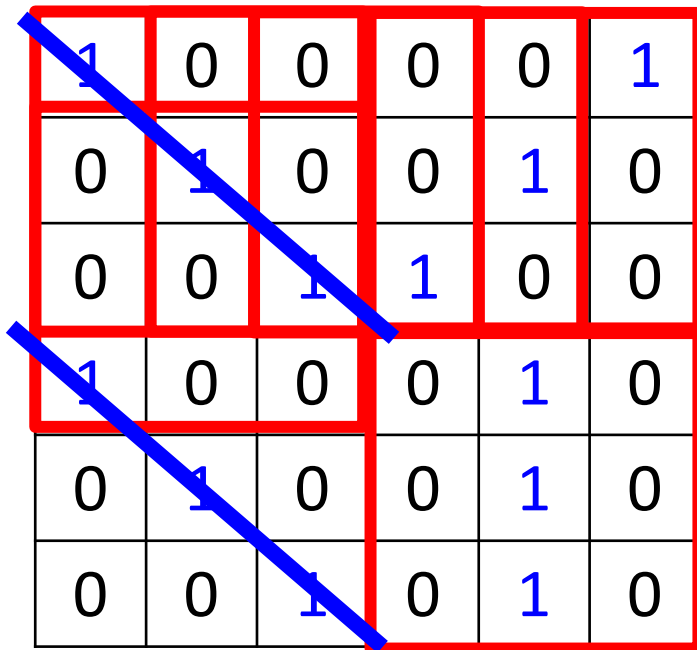
1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



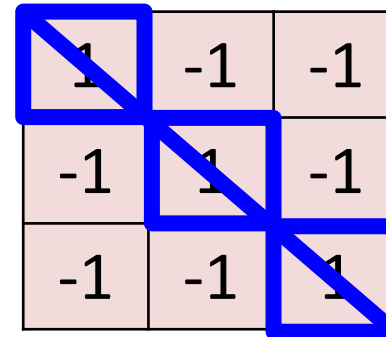
Convolution

stride=1



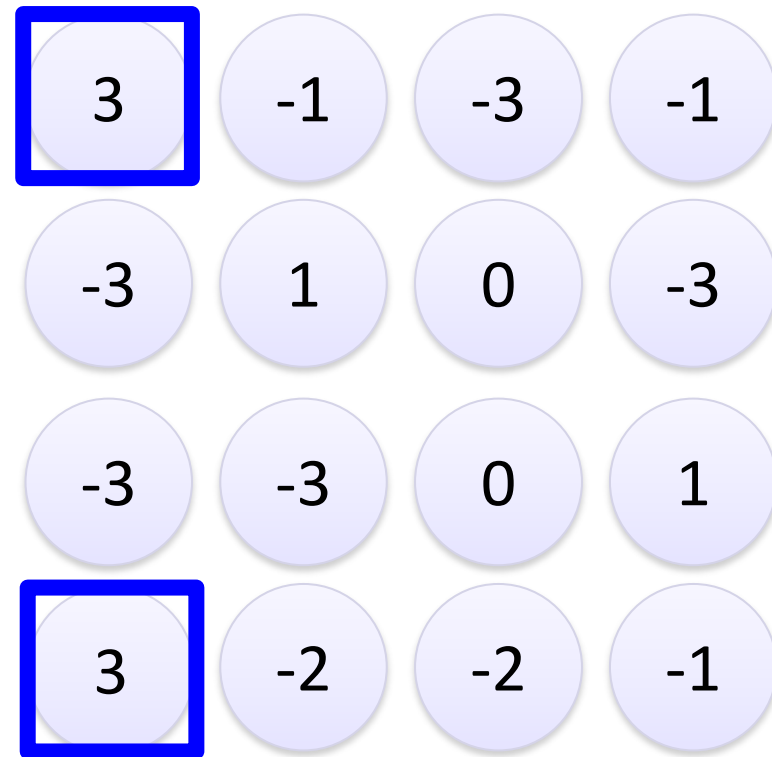
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image



1	-1	-1
-1	1	-1
-1	-1	1

Filter 1



3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

Convolution

stride=1

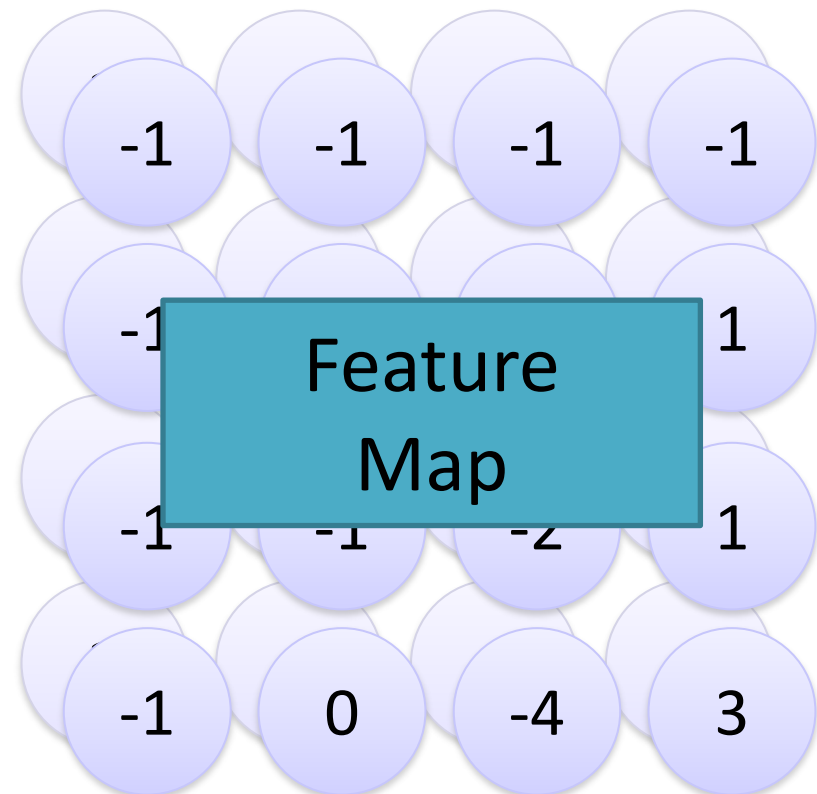
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

-1	1	-1
-1	1	-1
-1	1	-1

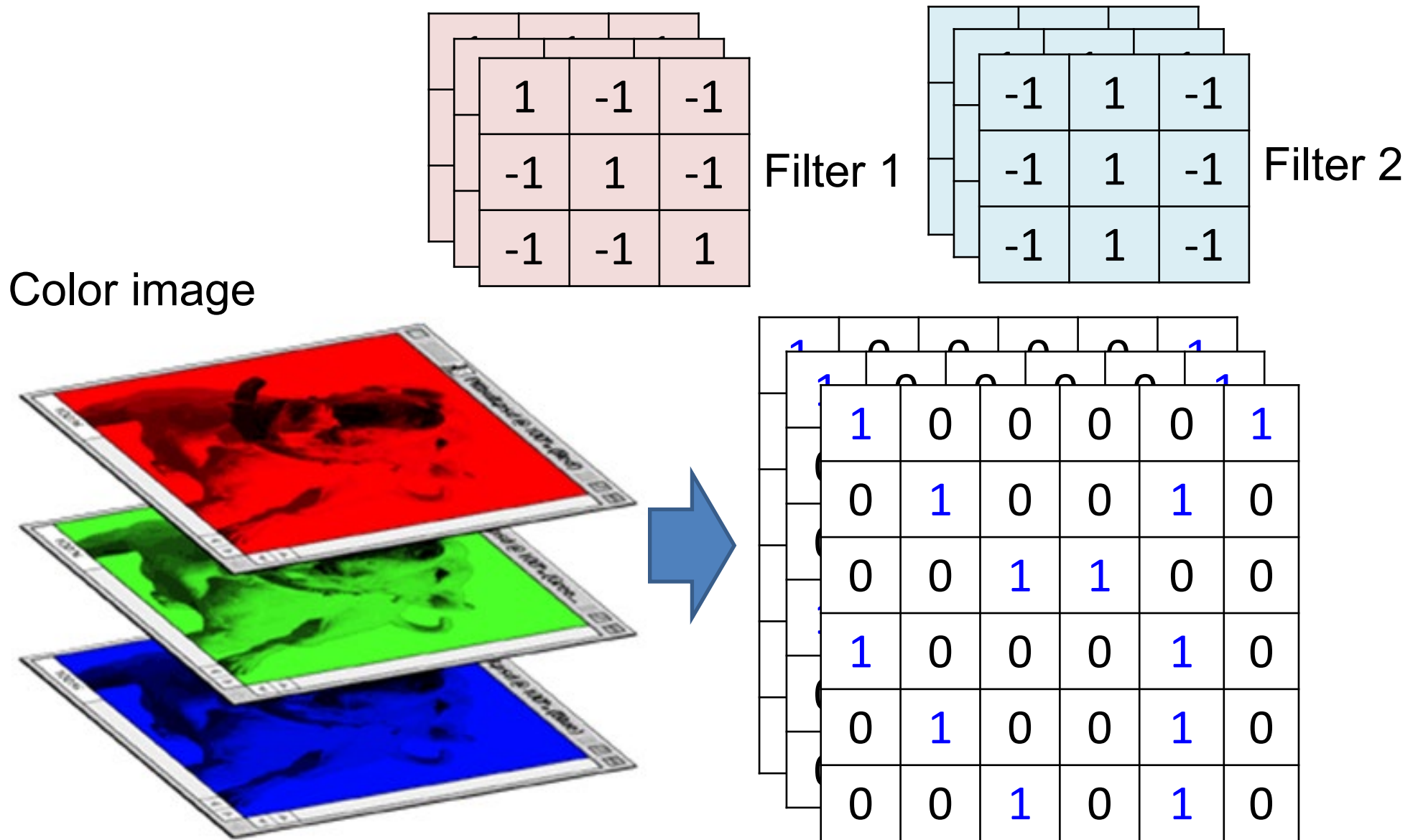
Filter 2

Repeat this for each filter



Two 4 x 4 images
Forming 2 x 4 x 4 matrix

Color image: RGB 3 channels



Convolution vs. Fully Connected

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

image

1	-1	-1
-1	1	-1
-1	-1	1

-1	1	-1
-1	1	-1
-1	1	-1

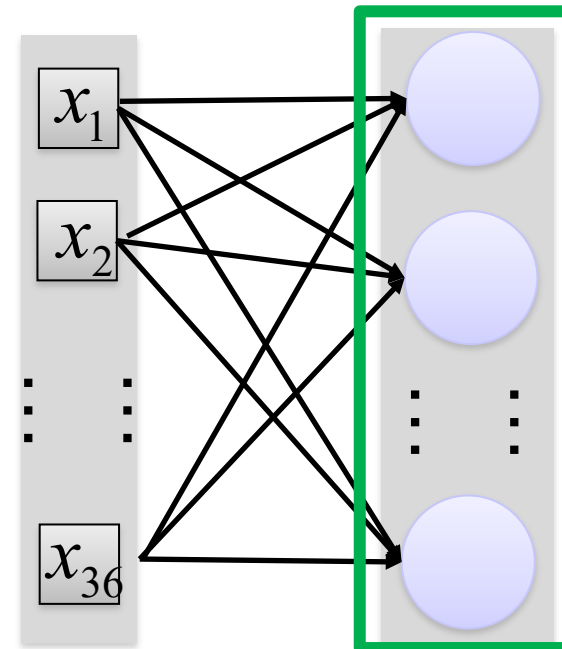


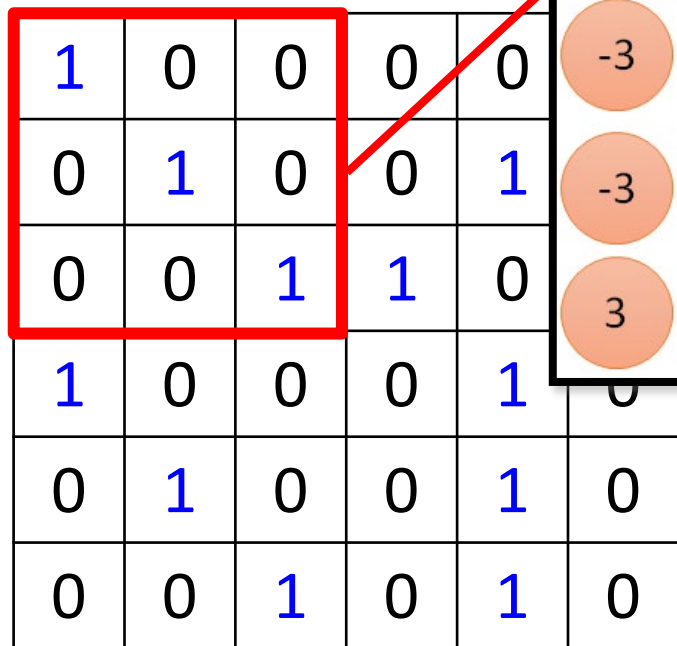
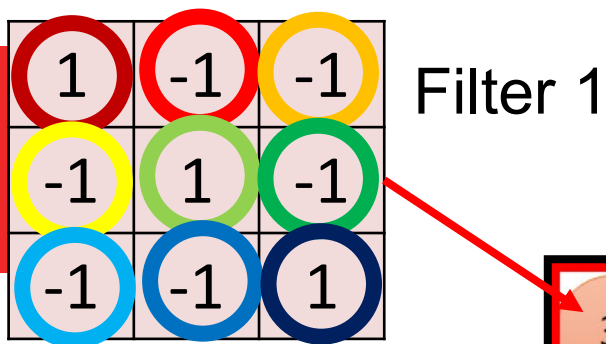
convolution

-1	-1	-1	-1
-1	-1	-2	1
-1	-1	-2	1
-1	0	-4	3

Fully-
connected

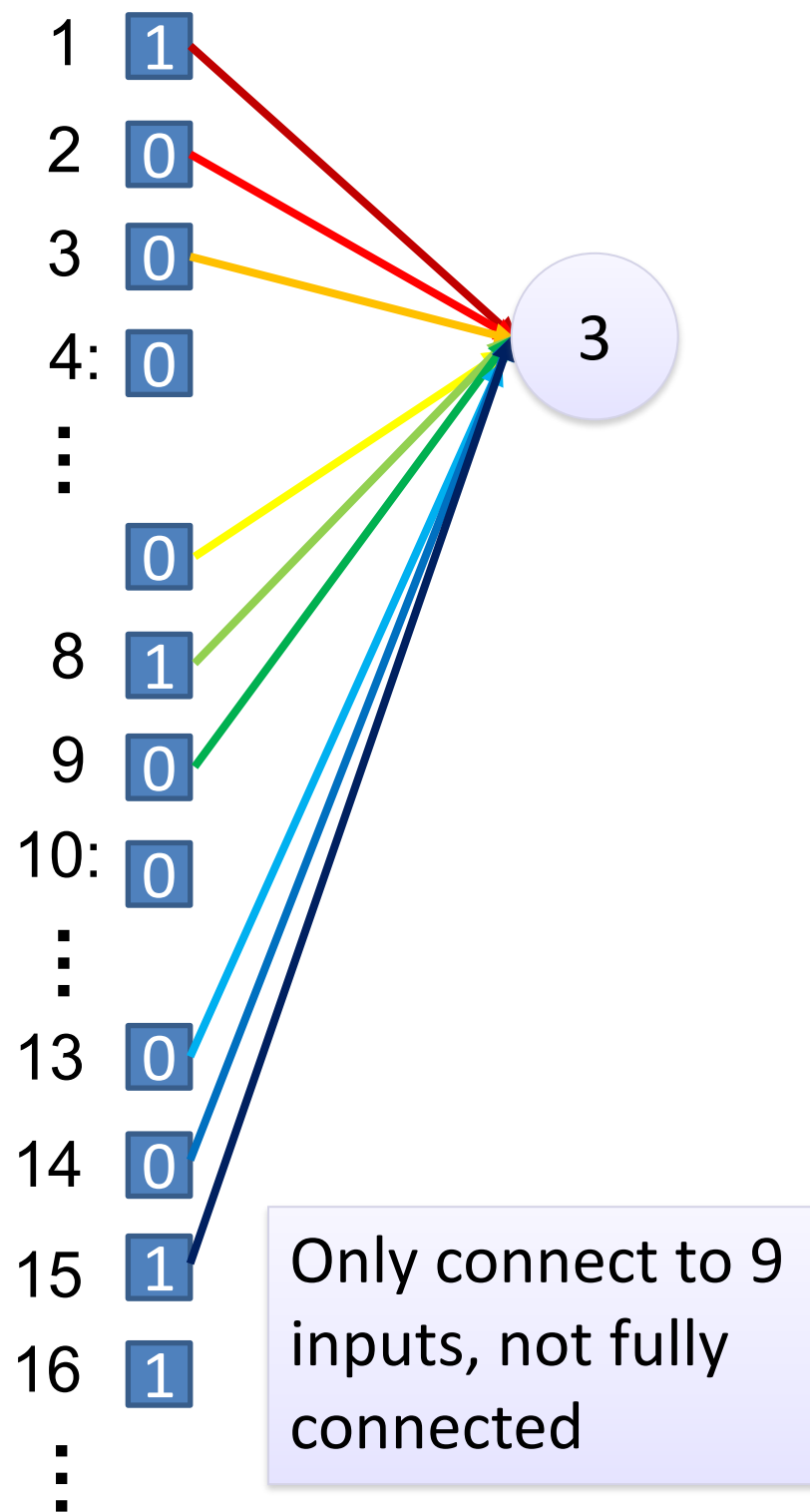
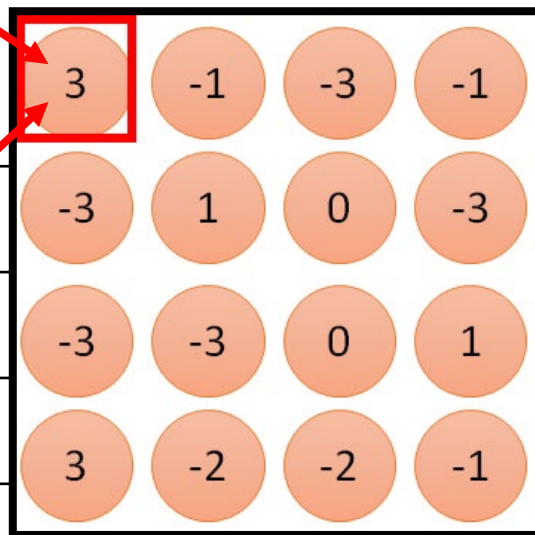
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

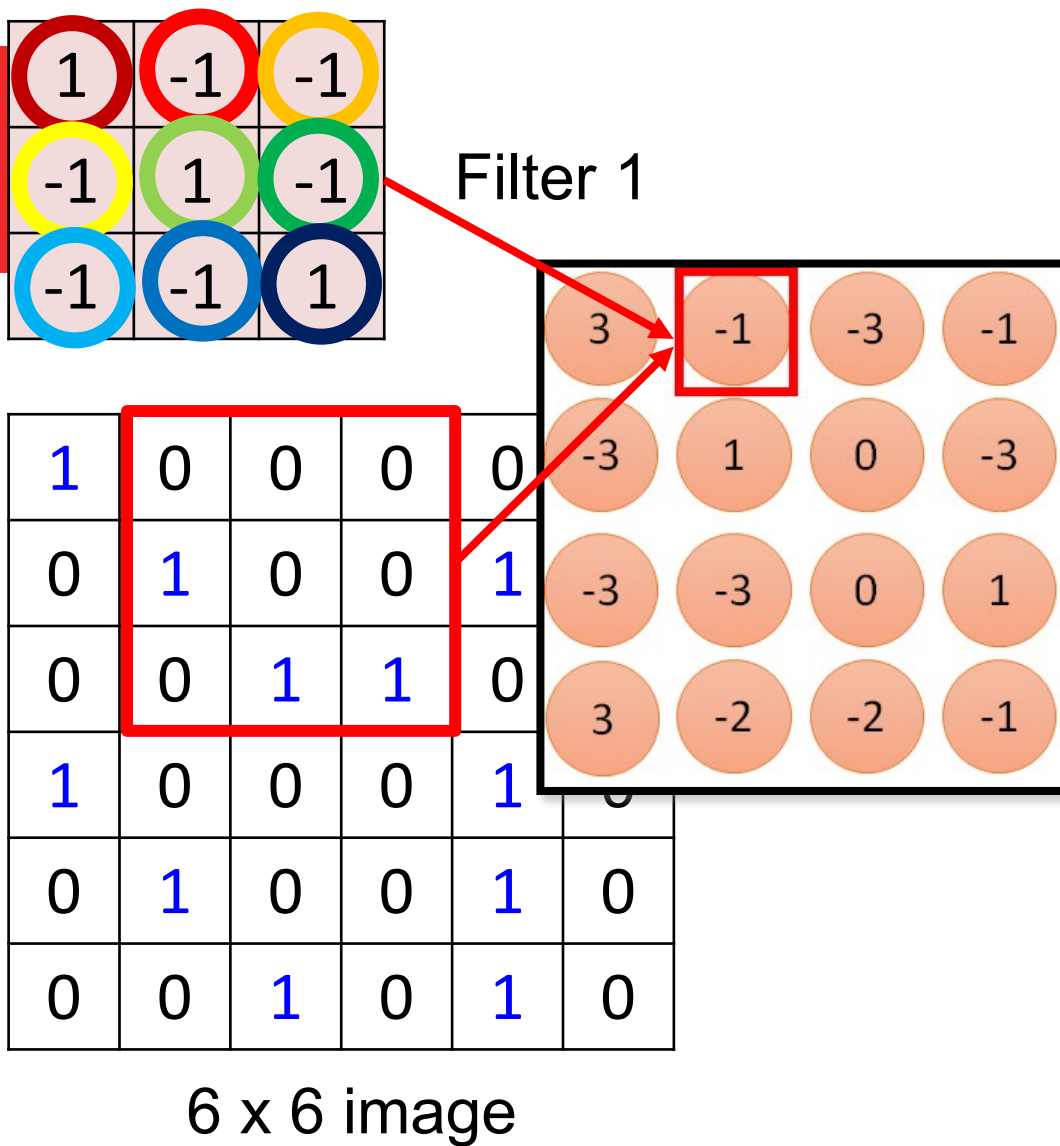




6 x 6 image

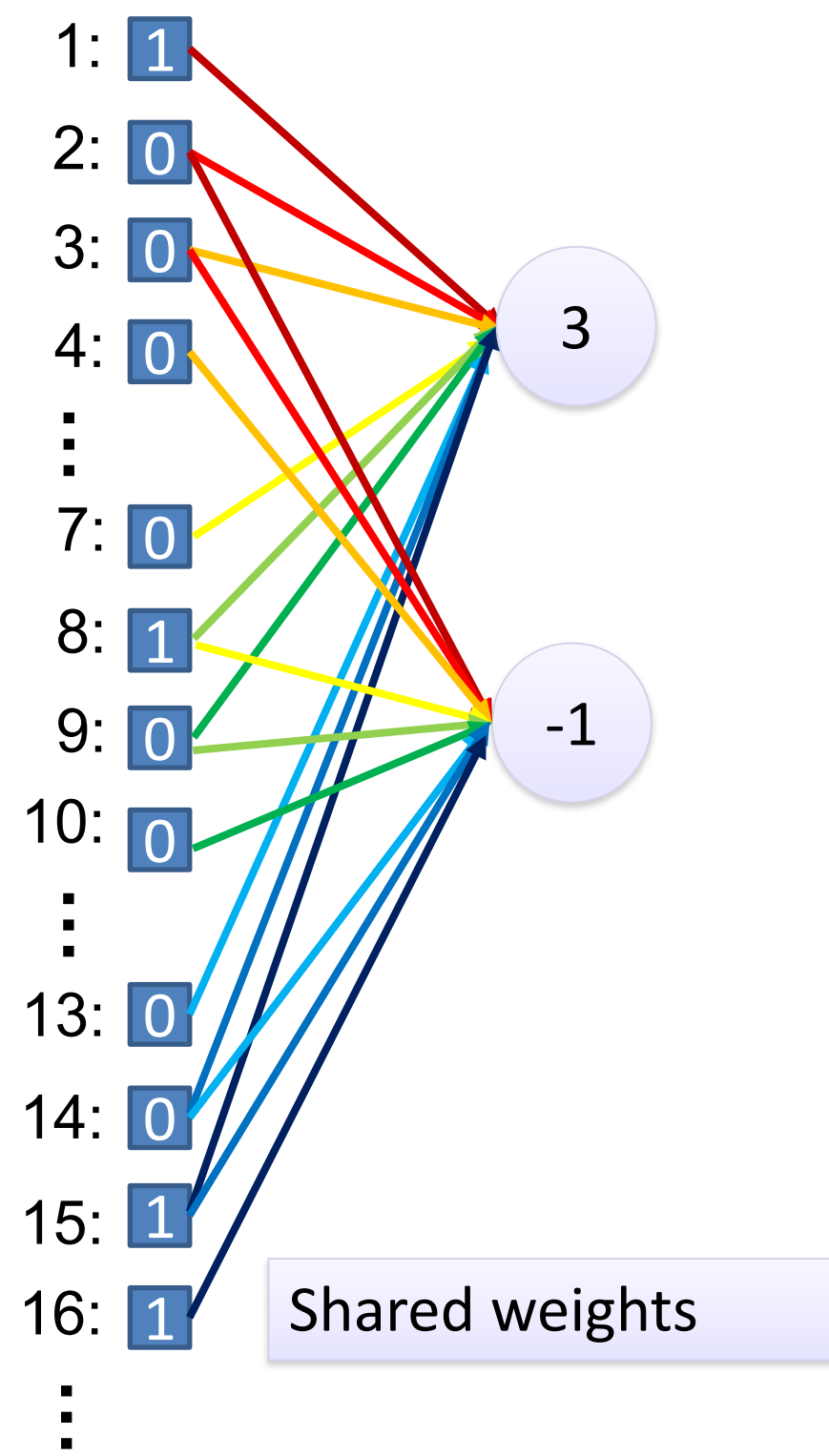
fewer parameters!



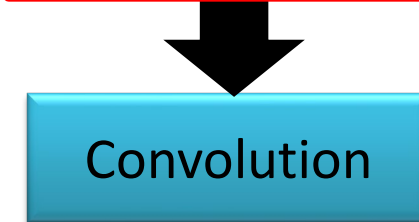
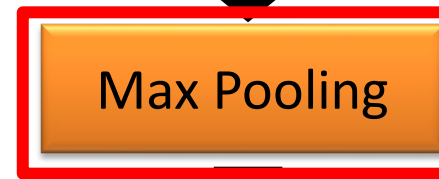
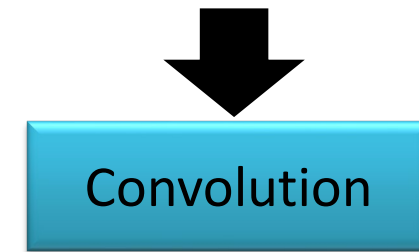
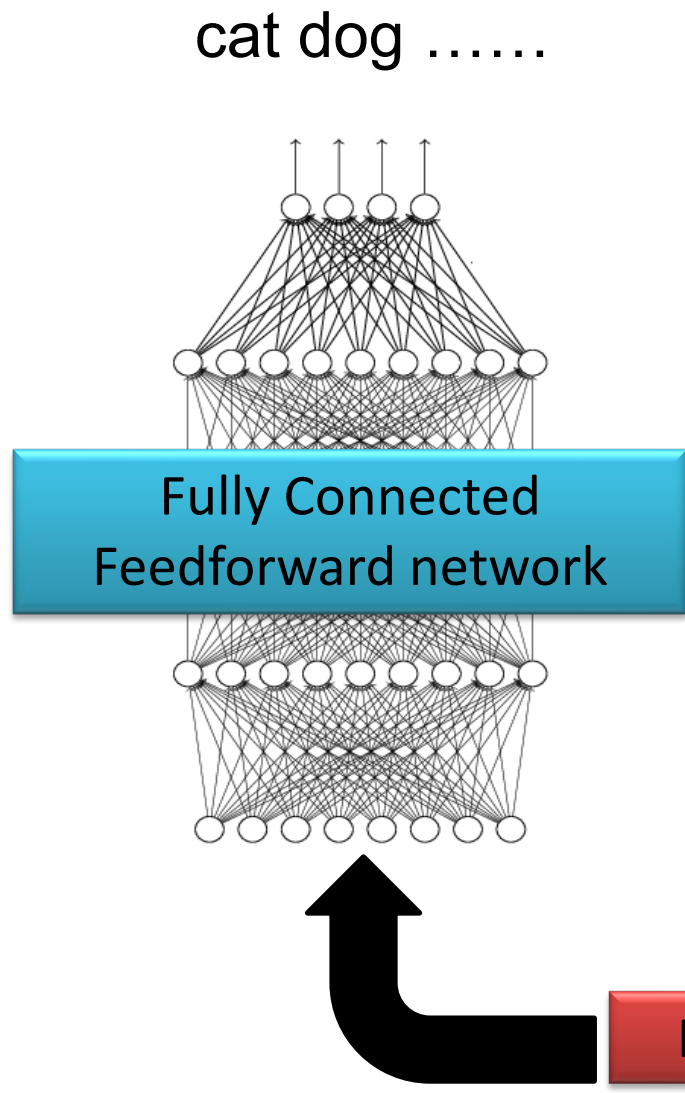


Fewer parameters

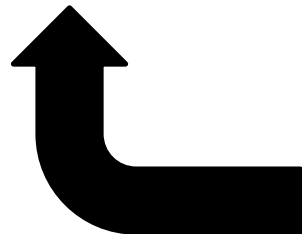
Even fewer parameters



The whole CNN



Can repeat many times



Max Pooling

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

-1	-1	-1	-1
-1	-1	-2	1
-1	-1	-2	1
-1	0	-4	3

Why Pooling

- Subsampling pixels will not change the object

bird

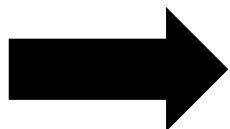


Subsampling

bird



We can subsample the pixels to make image smaller

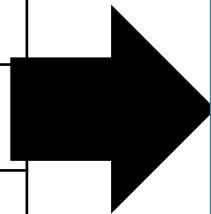


fewer parameters to characterize the image

Max Pooling

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

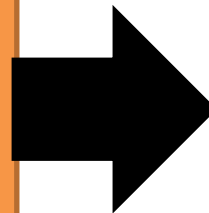
6 x 6 image



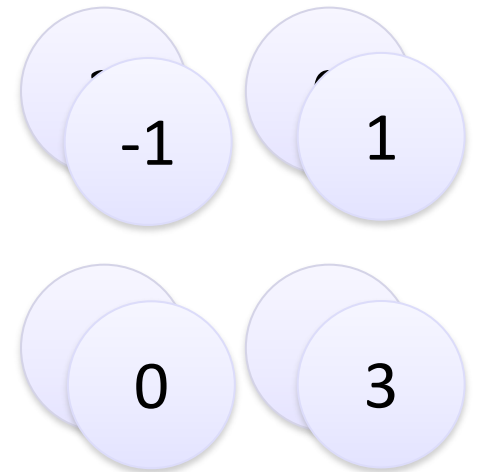
Conv



Max
Pooling



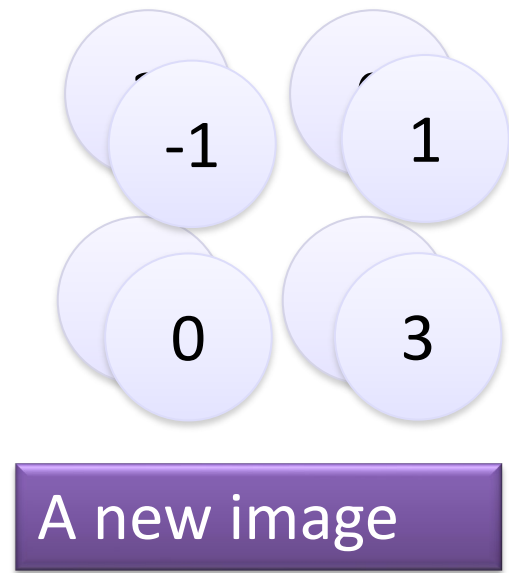
New image
but smaller



2 x 2 image

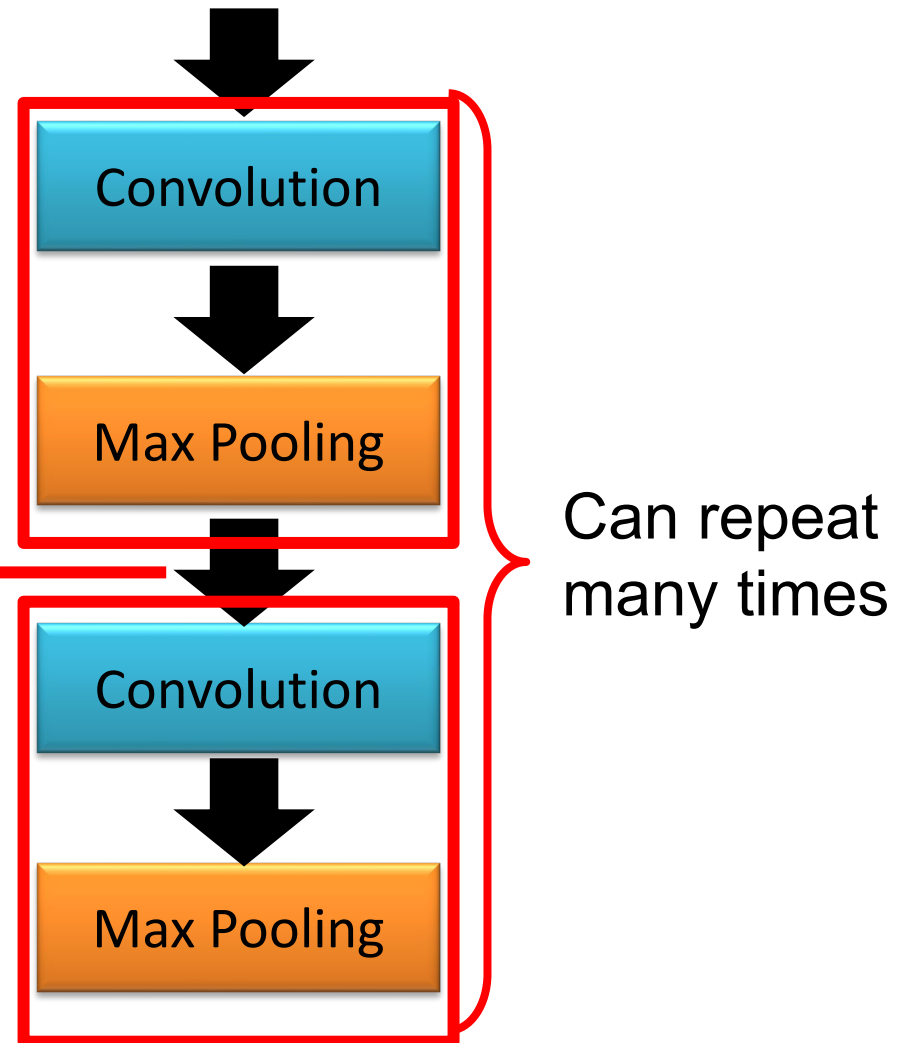
Each filter
is a channel

The whole CNN



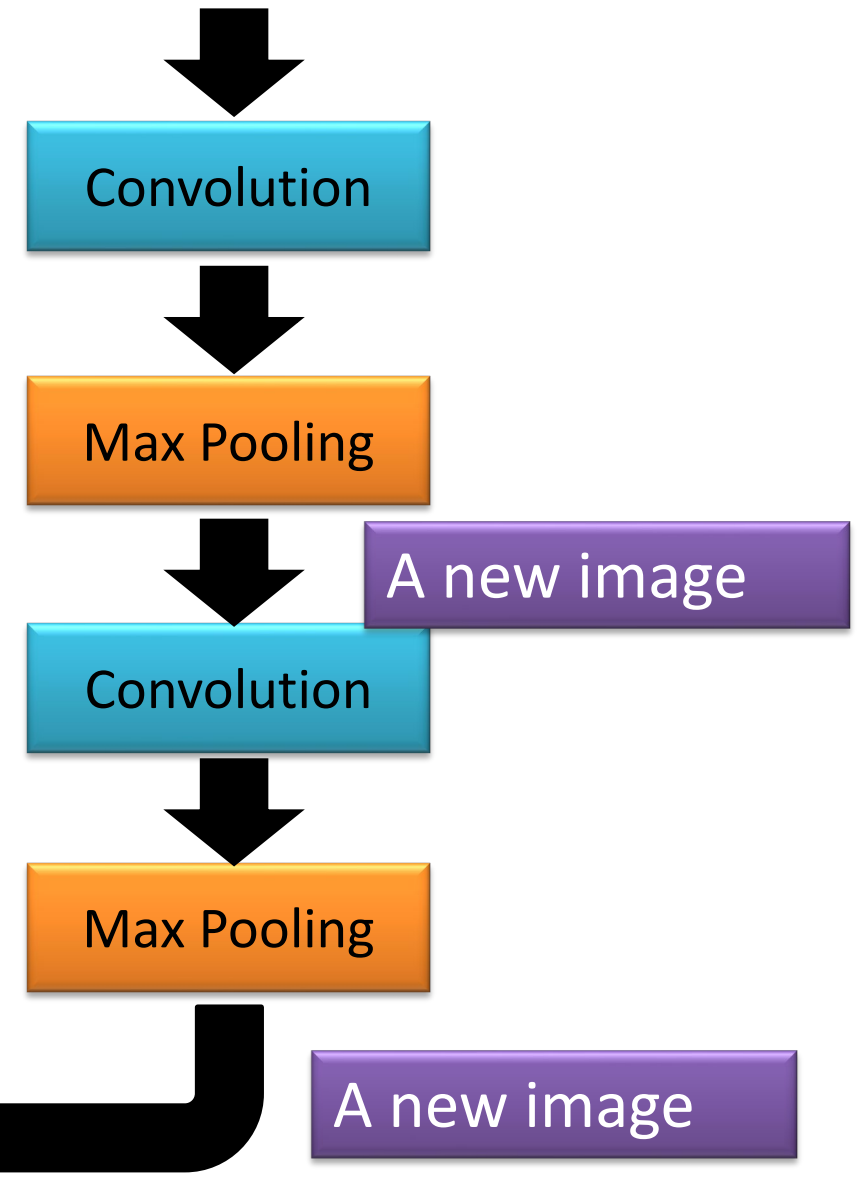
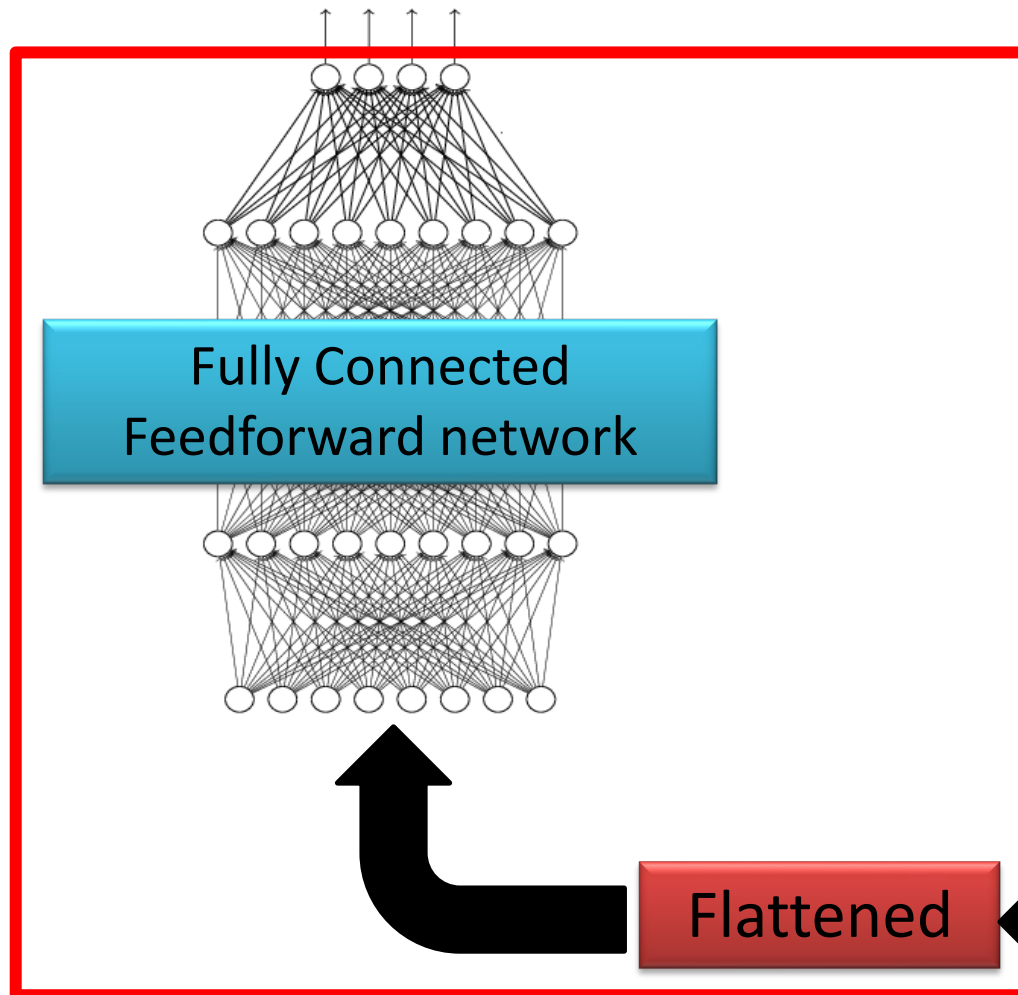
Smaller than the original image

The number of channels is the number of filters

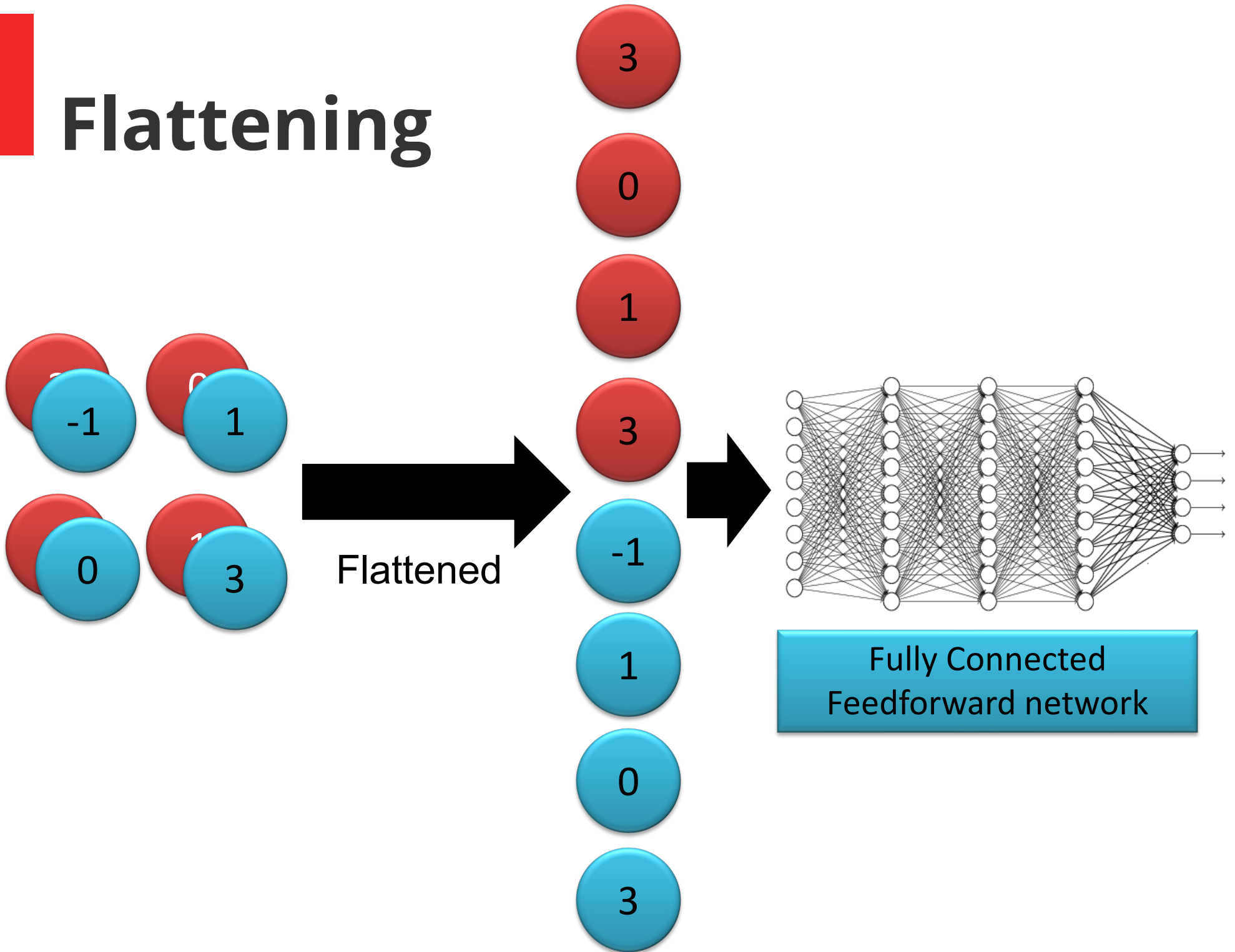


The whole CNN

cat dog



Flattening

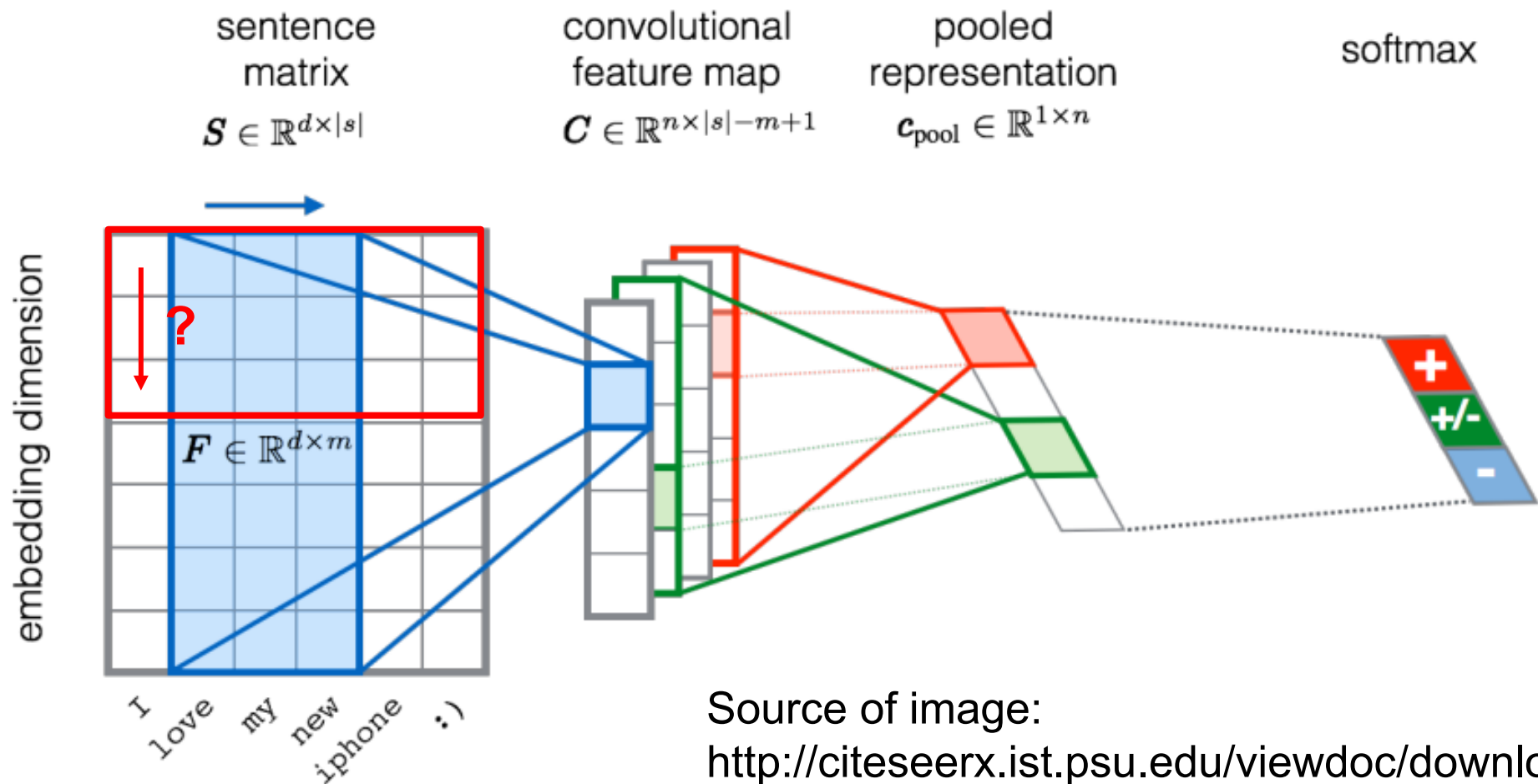




Transfer on other data input

- How to deal with Text as input?
 - Transform to image?
 - Is the input parameter independent?
- Same requirements fitting?
 - Letters transform to integers?
 - Is the sequence order important?
 - Why smaller portions?

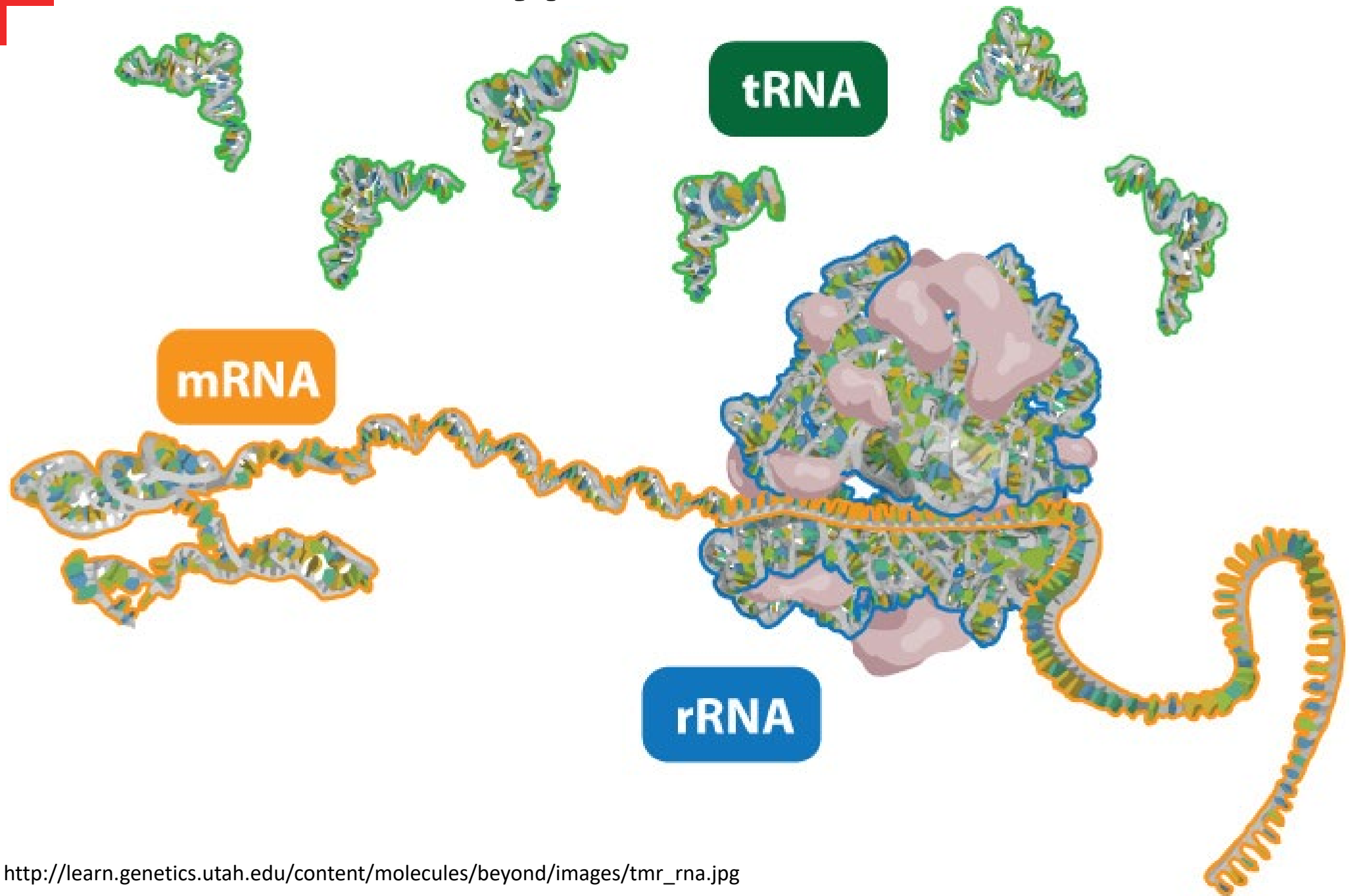
CNN in text classification





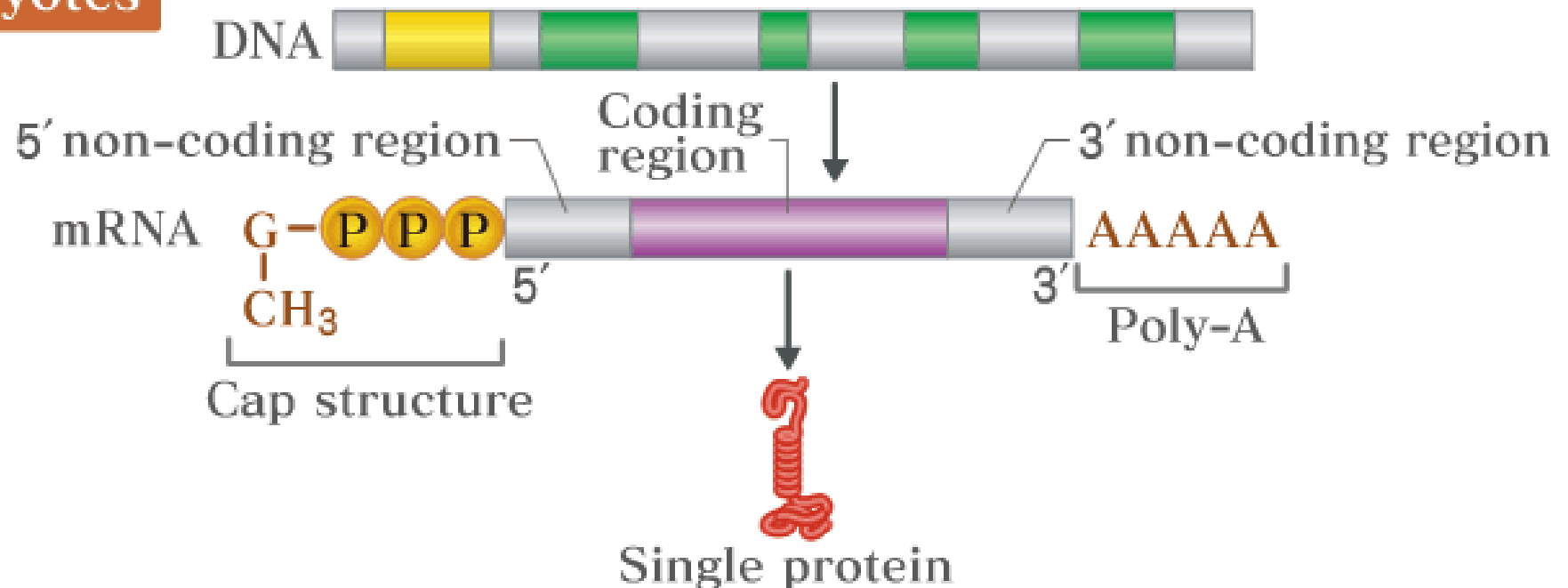
TYPES OF RNA AND THEIR SPECIFIC STRUCTURES

Types of RNA



mRNA features / structure

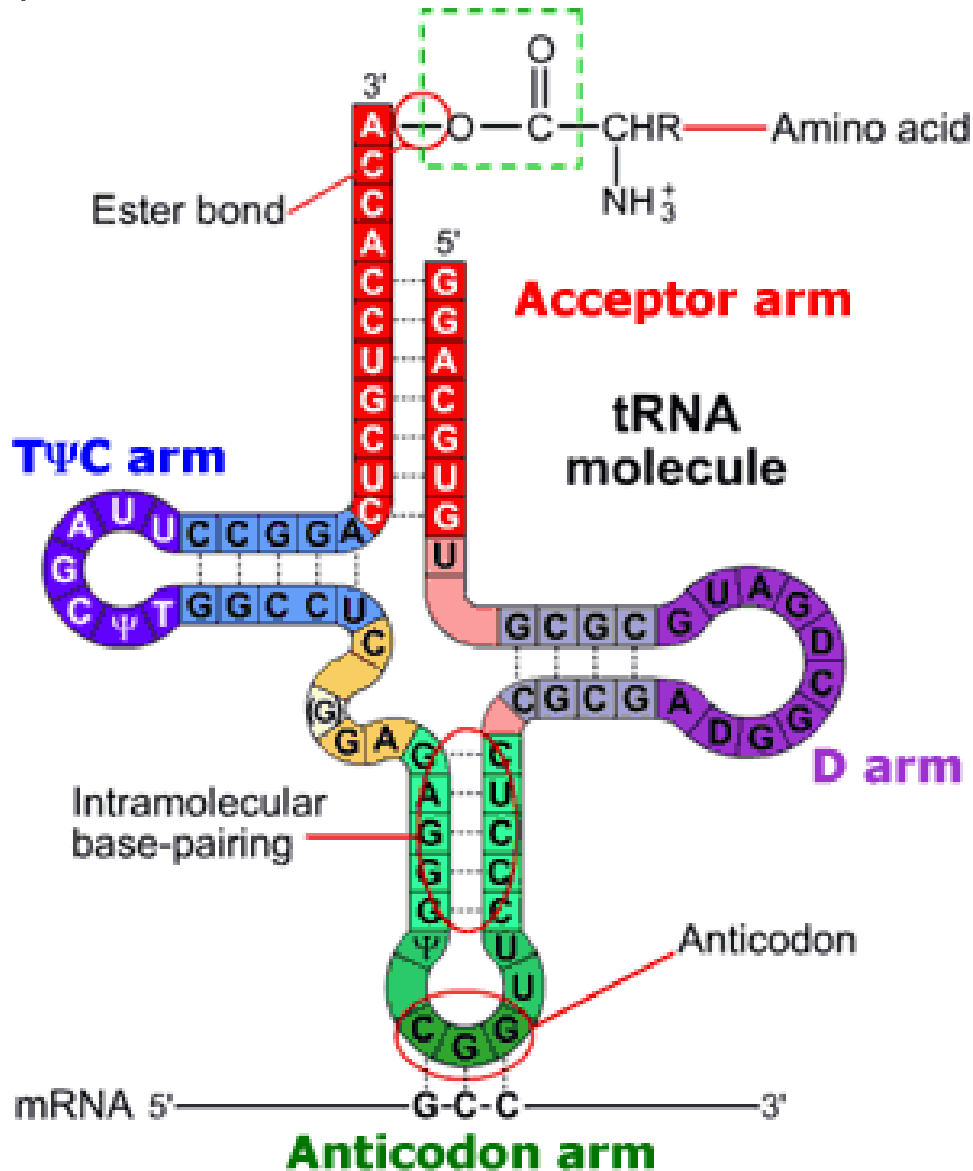
Eukaryotes



- ❑ Comprises only 5% of the RNA in the cell
- ❑ Most heterogeneous in size and base sequence
- ❑ All members of the class function as messengers carrying the information in a gene to the protein synthesizing machinery

tRNA features / structure

<http://www.alamy.com/search-results.asp?qt=Friedrich+Sauer+Einstein&pg=7>



- Small RNA species (74 – 95 nucleotides)
- Transfer amino acid to translation machinery (easily soluble RNA)
- At least 20 species for 20 amino acids required

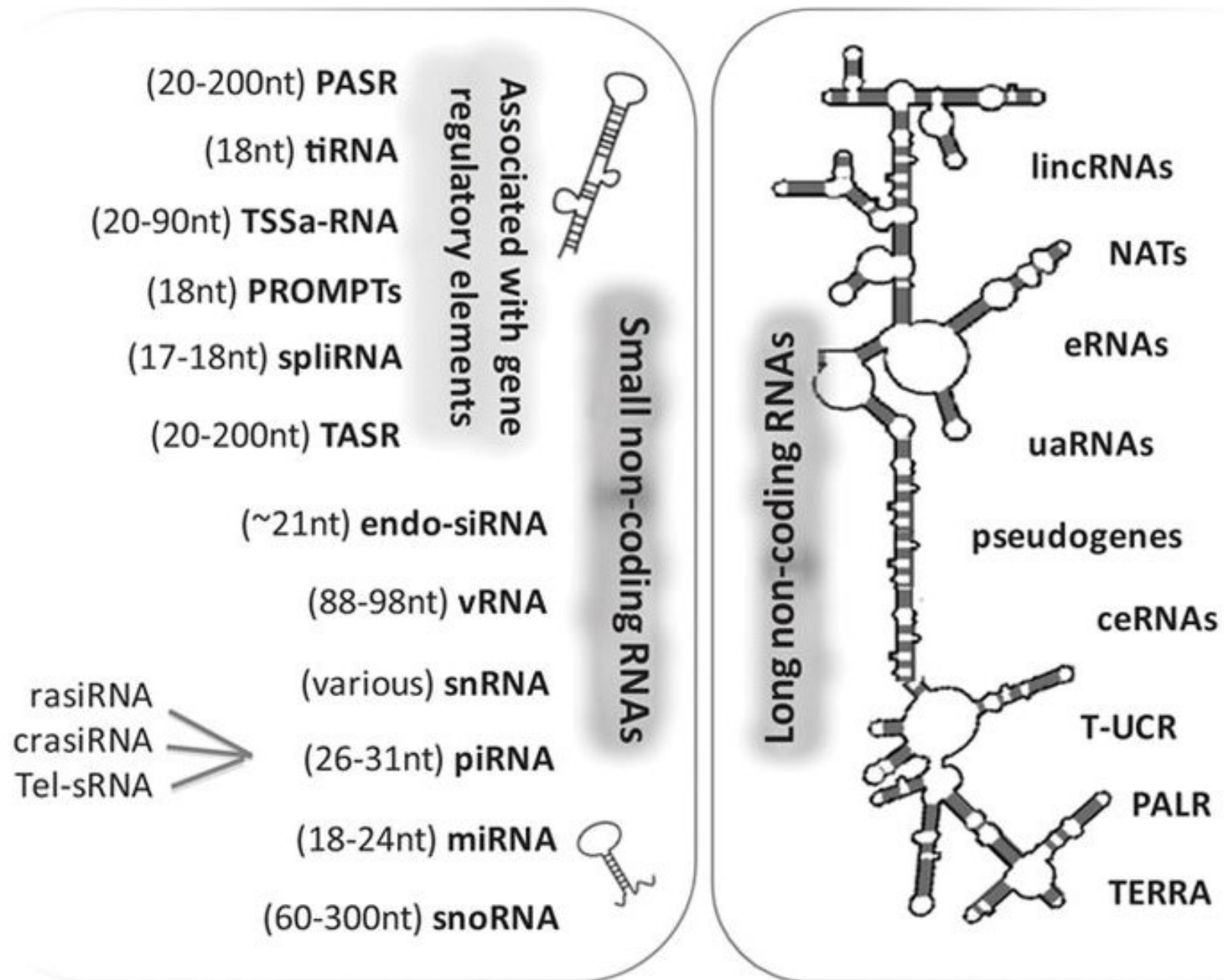
10



<http://www.pnas.org/content/107/46/19748/F3.large.jpg>

- ❖ Large and small ribosomal subunit:
- ❖ 60S → 5S, 5.8S and 28S rRNA
- ❖ 40S → 18S rRNA
- ❖ 5S is independently transcribed
- ❖ All other comes from 45S precursor

non-coding (nc) RNAs



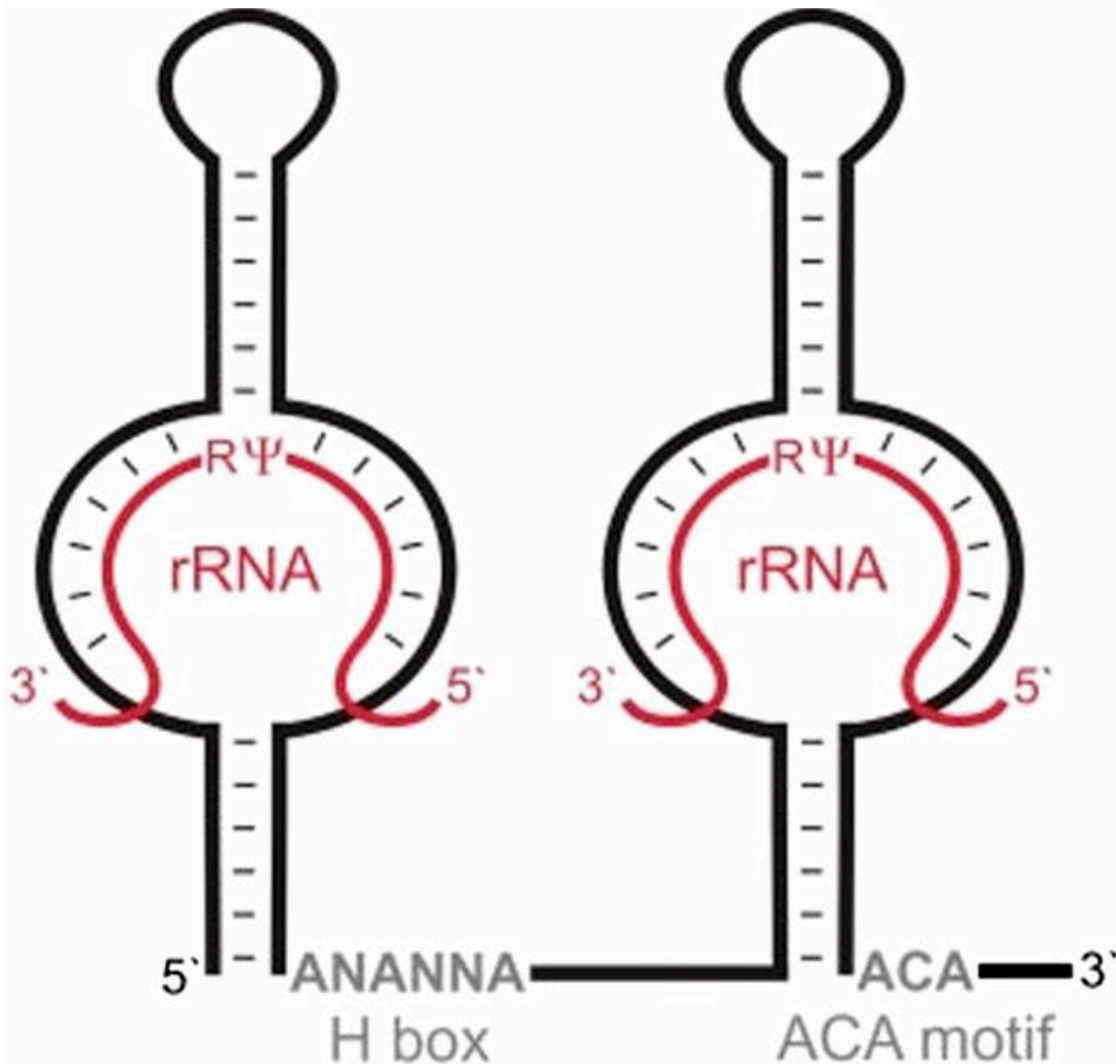


Aim for today's lesson

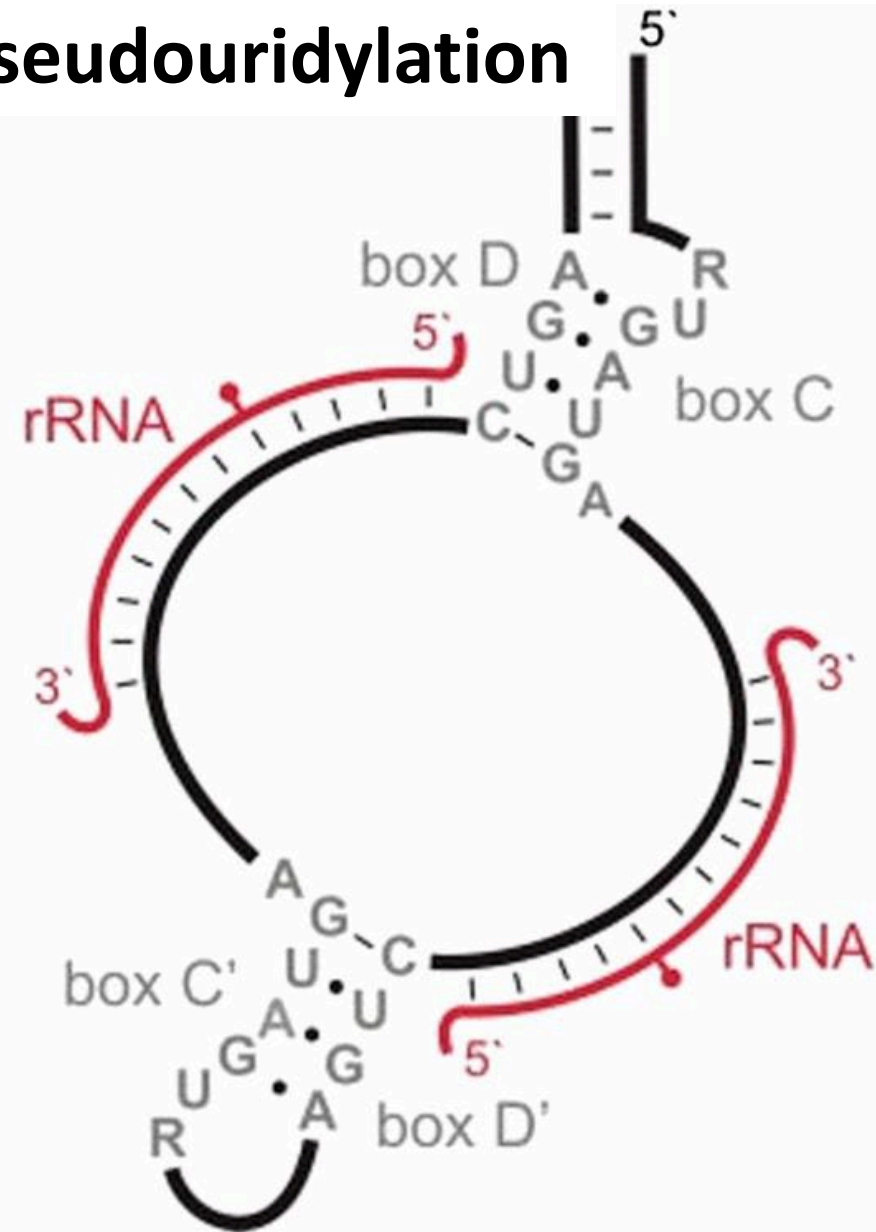
CNN FOR SEQUENCE CLASSIFICATION

snoRNAs

2'-O-ribose methylation

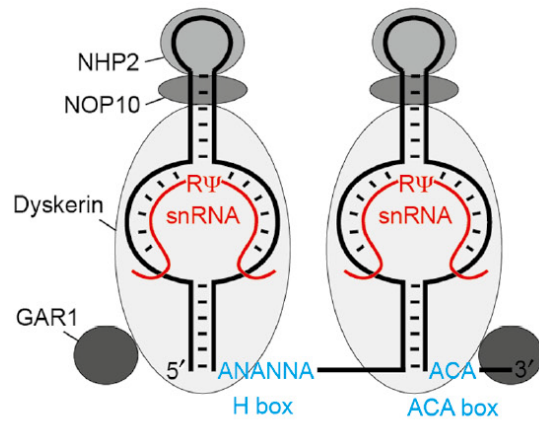


Pseudouridylation



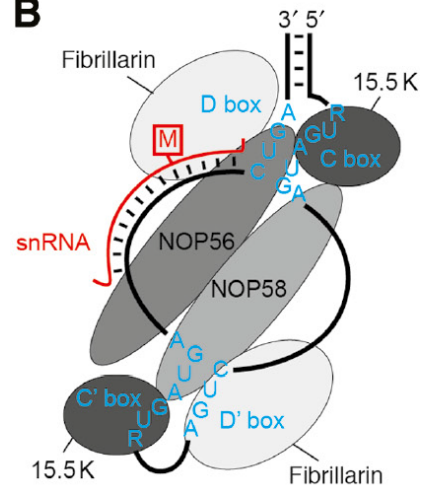
scaRNAs

A



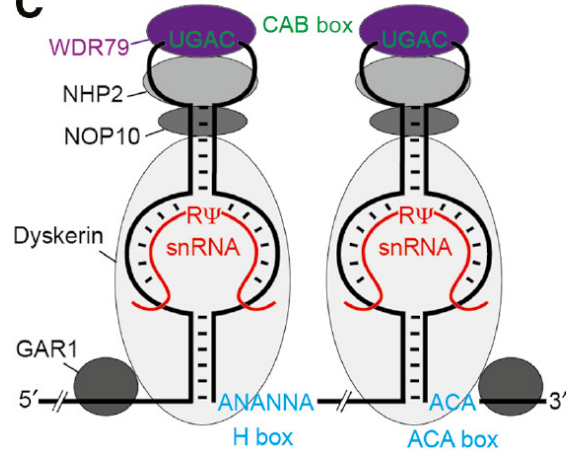
Canonical H/ACA box snoRNP

B



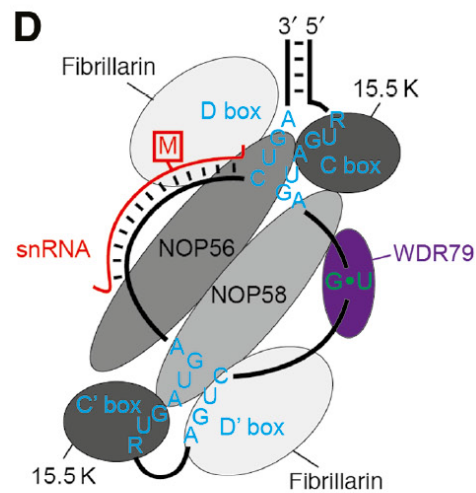
Canonical C/D box snoRNP

C



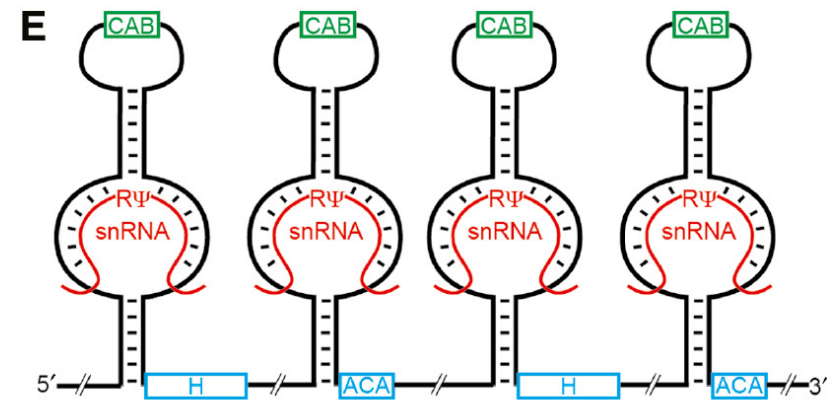
Canonical H/ACA box scaRNP

D



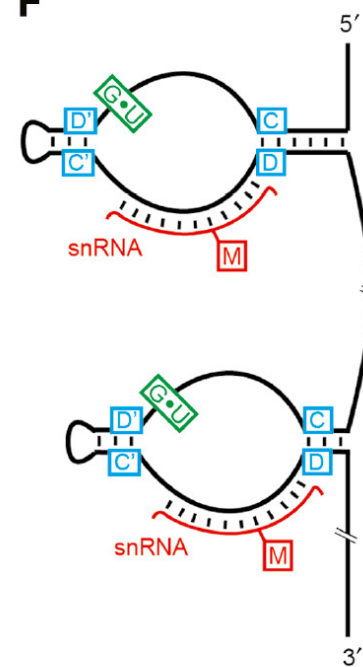
Canonical C/D box scaRNP

E



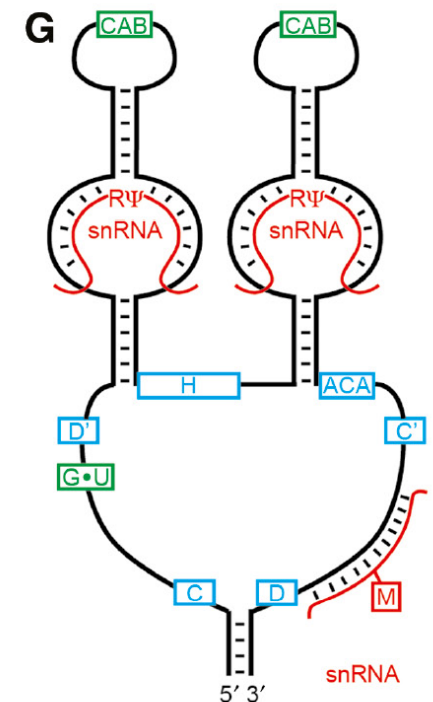
Tandem H/ACA box scaRNA

F



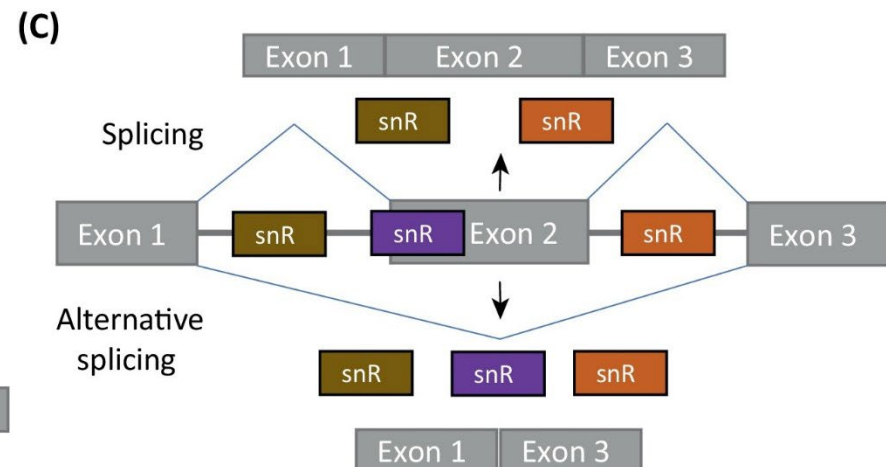
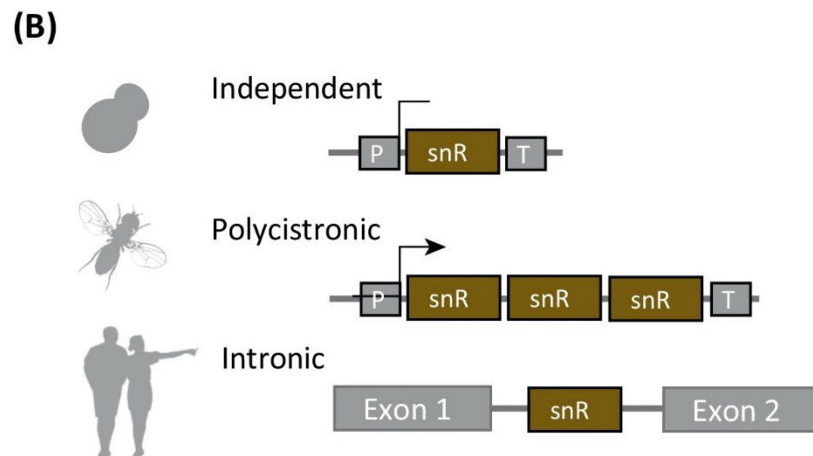
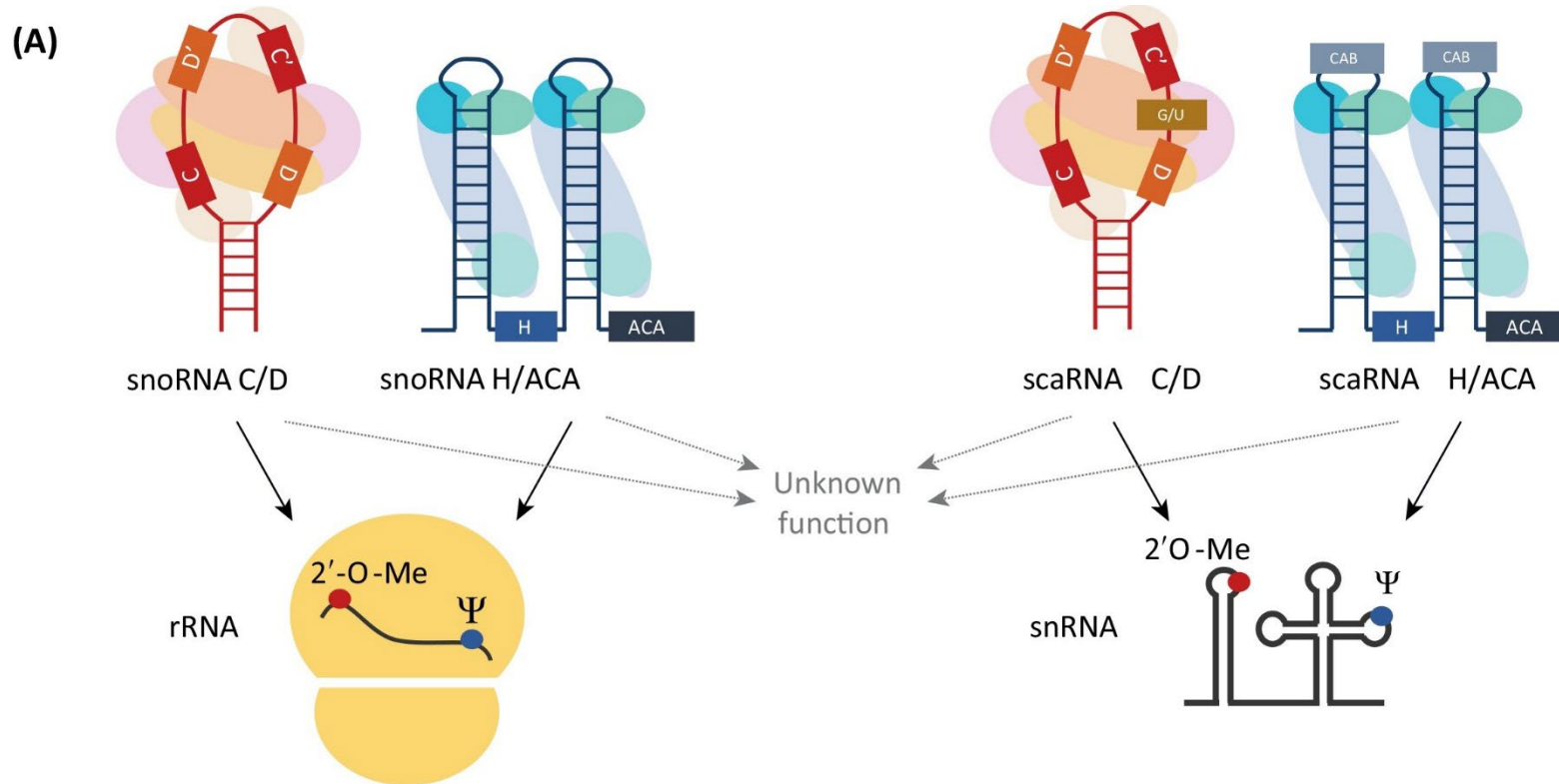
Tandem C/D box scaRNA

G



Hybrid H/ACA-C/D box scaRNA

scaRNAs vs snoRNAs



The dataset

- **RNAcentral**

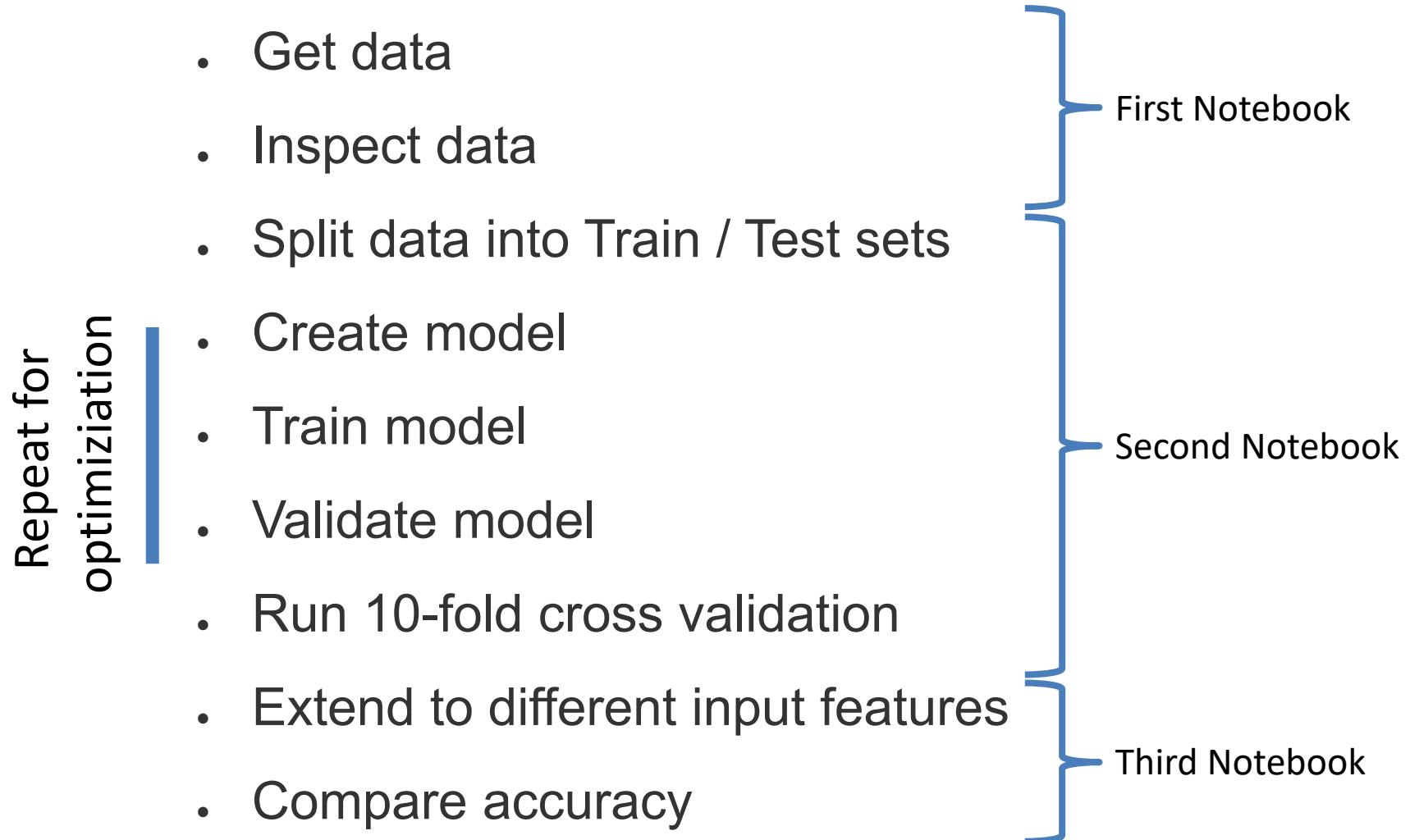
- snoRNA (C/D-box, H/ACA-box), scaRNA
- No hypotheticals, no partials

	snoRNA (H/ACA)	snoRNA (C/D)	scaRNA
Overall	89,522	152,235	6,747
Mammalia	57,640	34,286	4,094
Insecta	976	2,388	47
Fish (Actinopterygii)	7,320	11,530	752
Plants (Viridiplantae)	5,808	76,445	28
Fungi (Ascomycota)	5,375	7,835	1

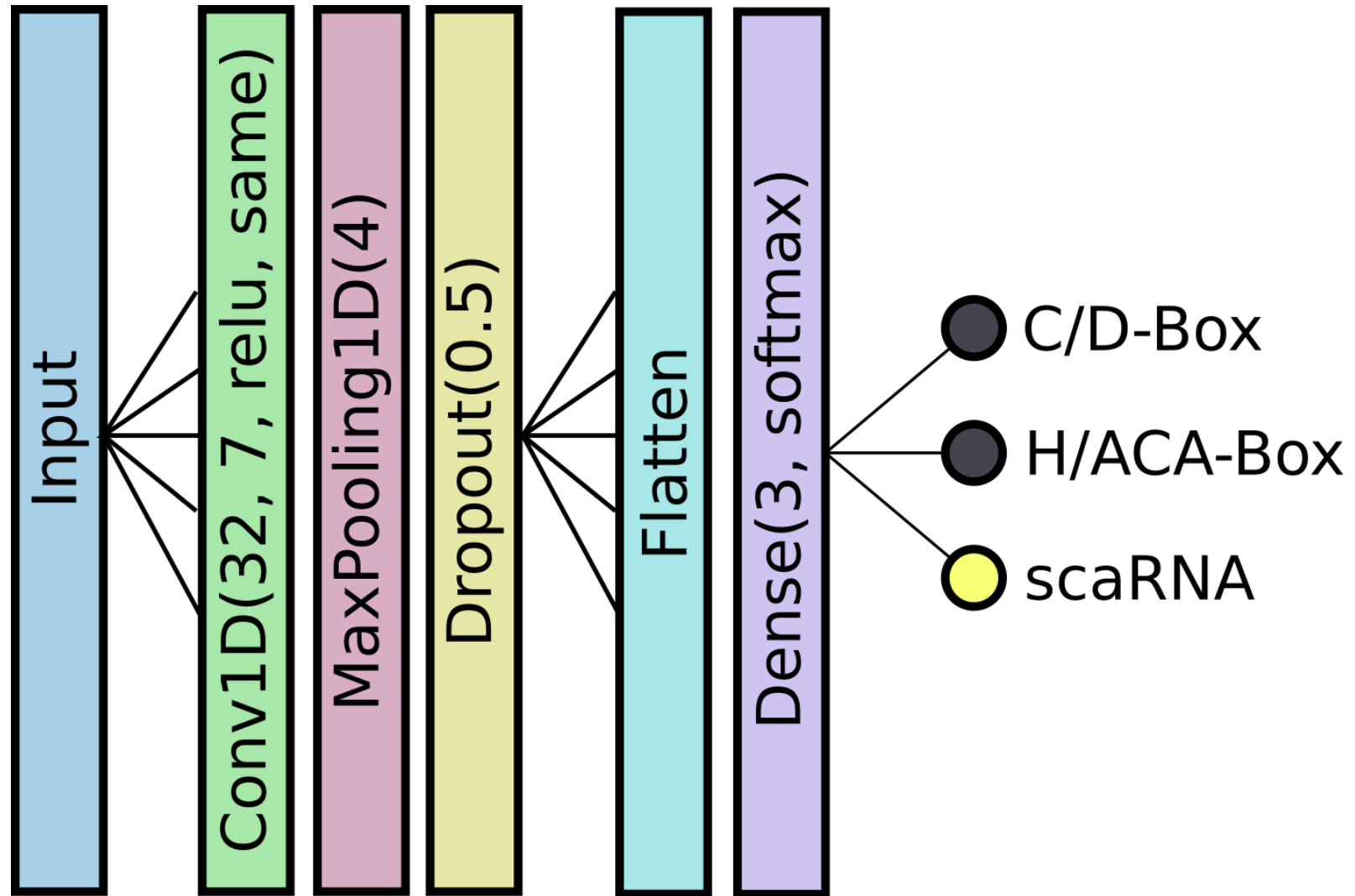
- **Preprocessing step**

- Clustered via CD-HIT to a max. of 0.9 (90%) similarity
- Balanced to 1913 sequences per class

Workflow



CNN our Standard Architecture



Let's Classify

