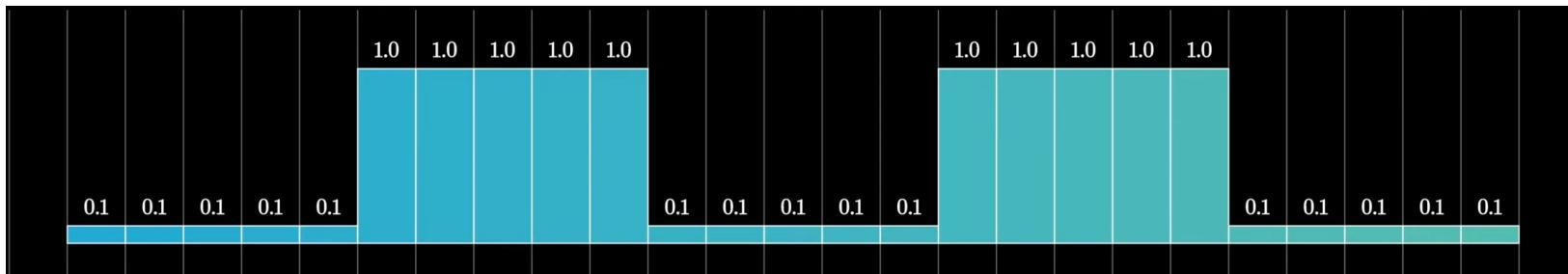


Deep Learning applied to EEG

Dung “Young” Truong & Arnaud Delorme



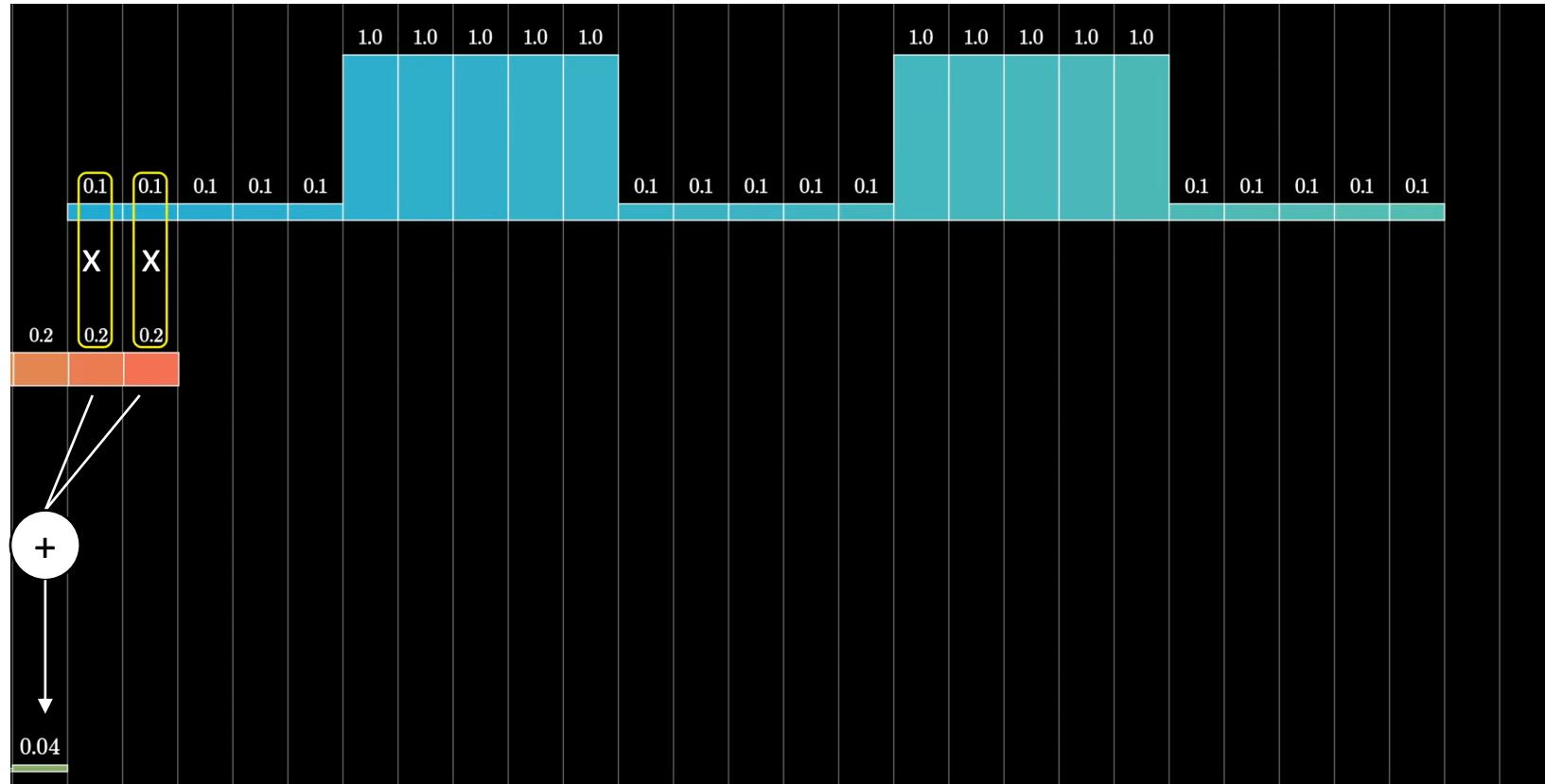
What kind of operation?



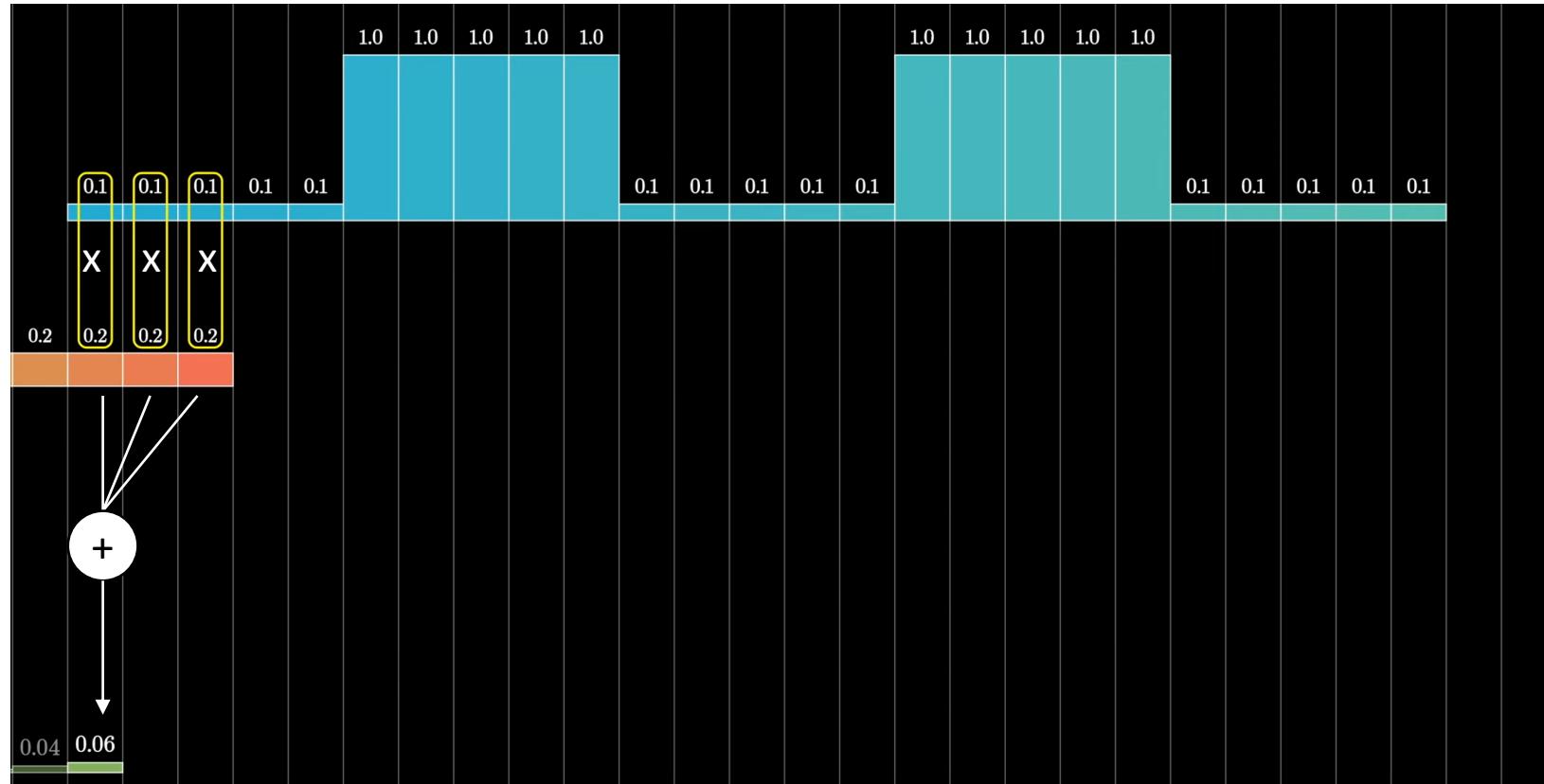
A diagram illustrating a convolution operation. The input layer consists of 10 units, each with a value of 0.1. A 5x5 kernel, also consisting of 25 units, each with a value of 0.2, is applied across the input. The output layer consists of 10 units, each with a value of 0.1. This visualizes how the kernel's value of 0.2 is distributed across the input units to produce the output.

$$\sum_i y_i = 1$$

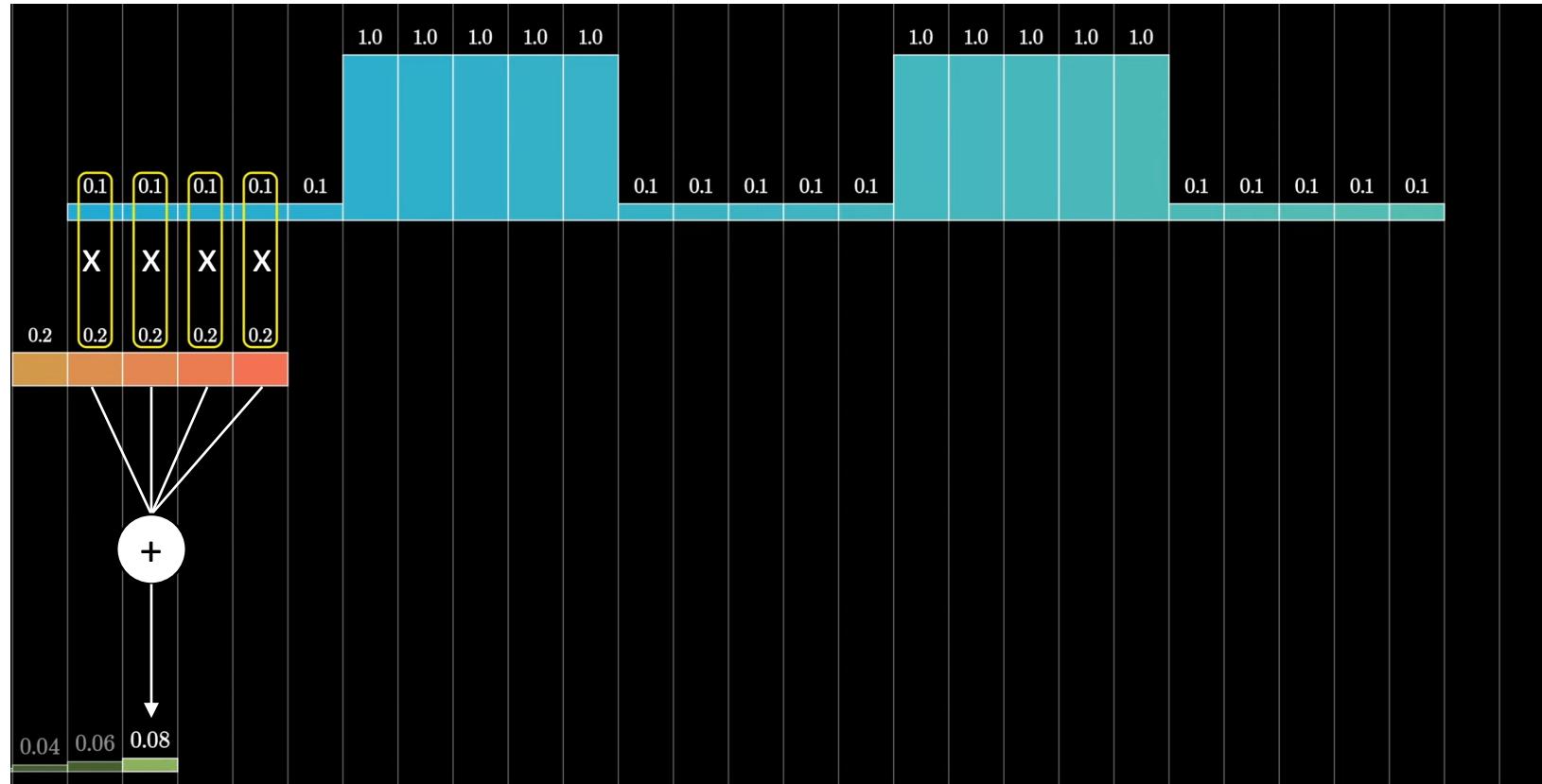
What kind of operation?



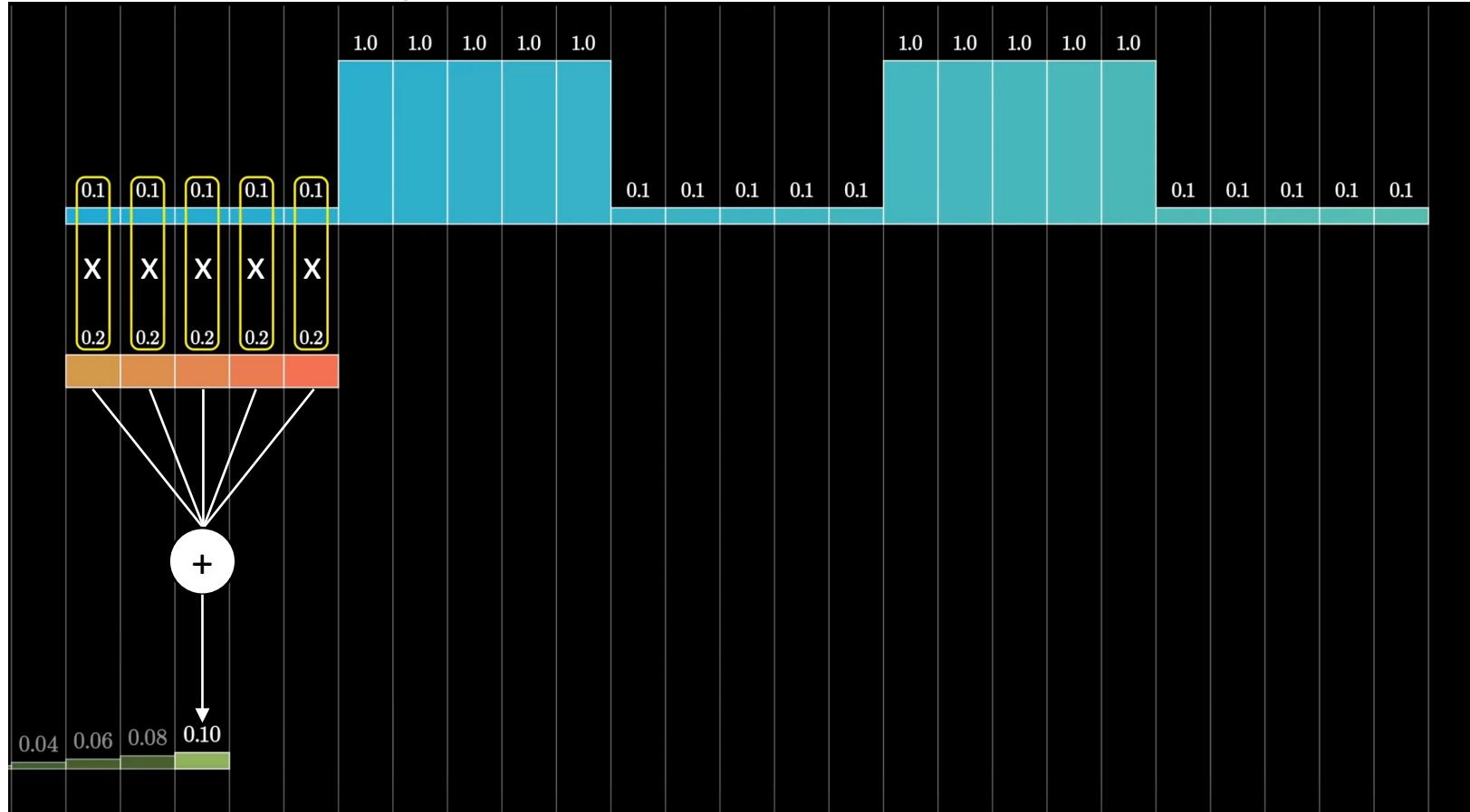
What kind of operation?



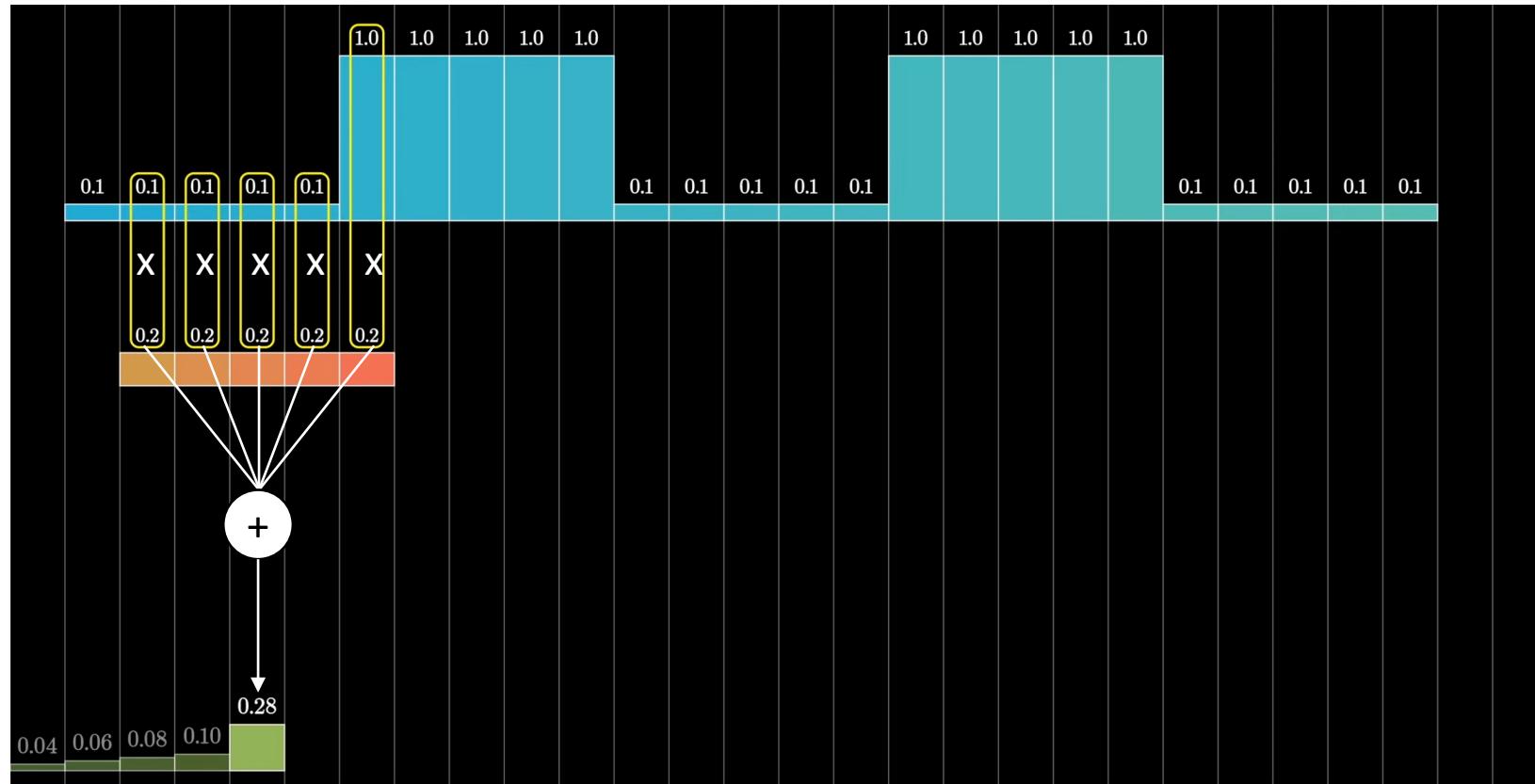
What kind of operation?



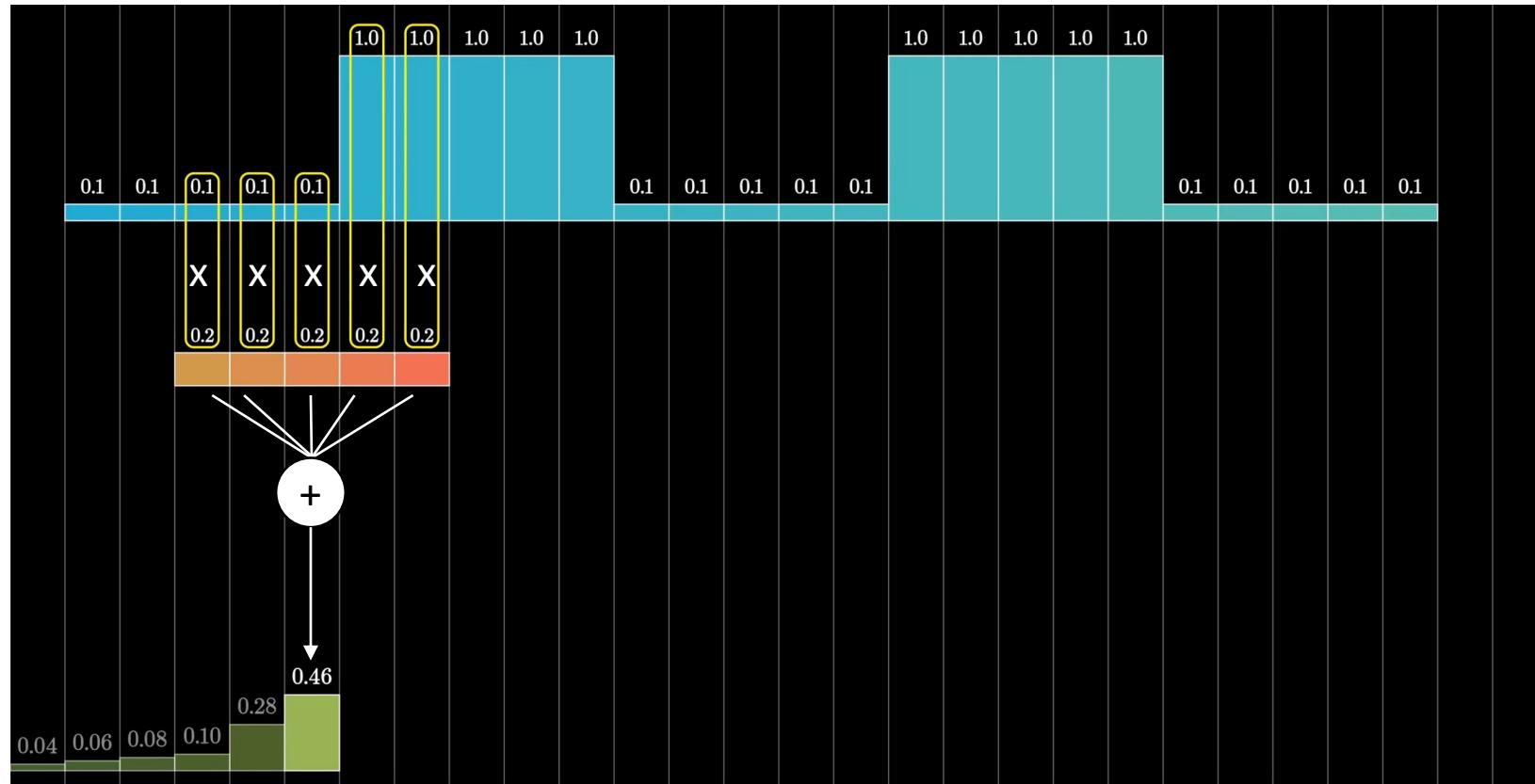
What kind of operation?



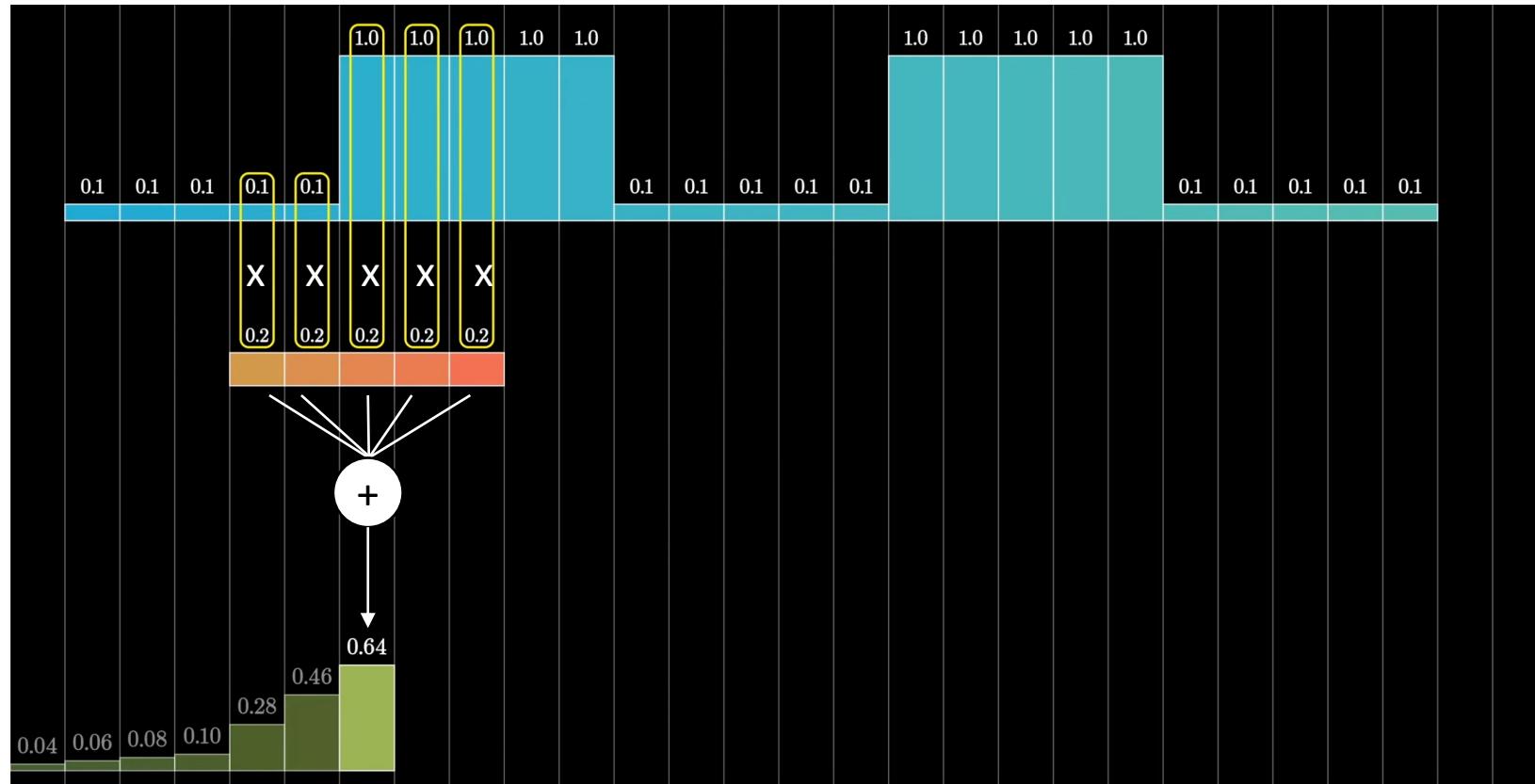
What kind of operation?



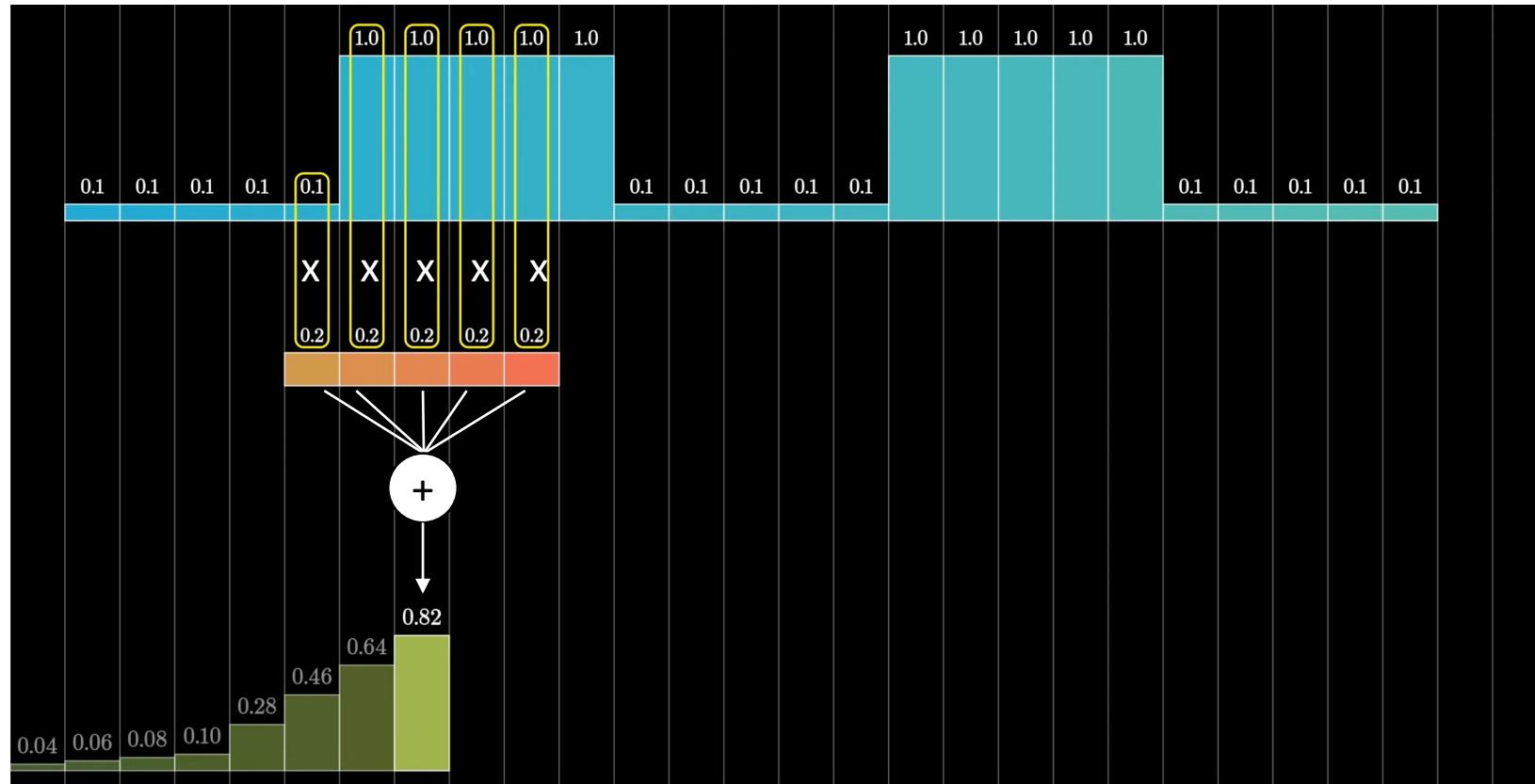
What kind of operation?



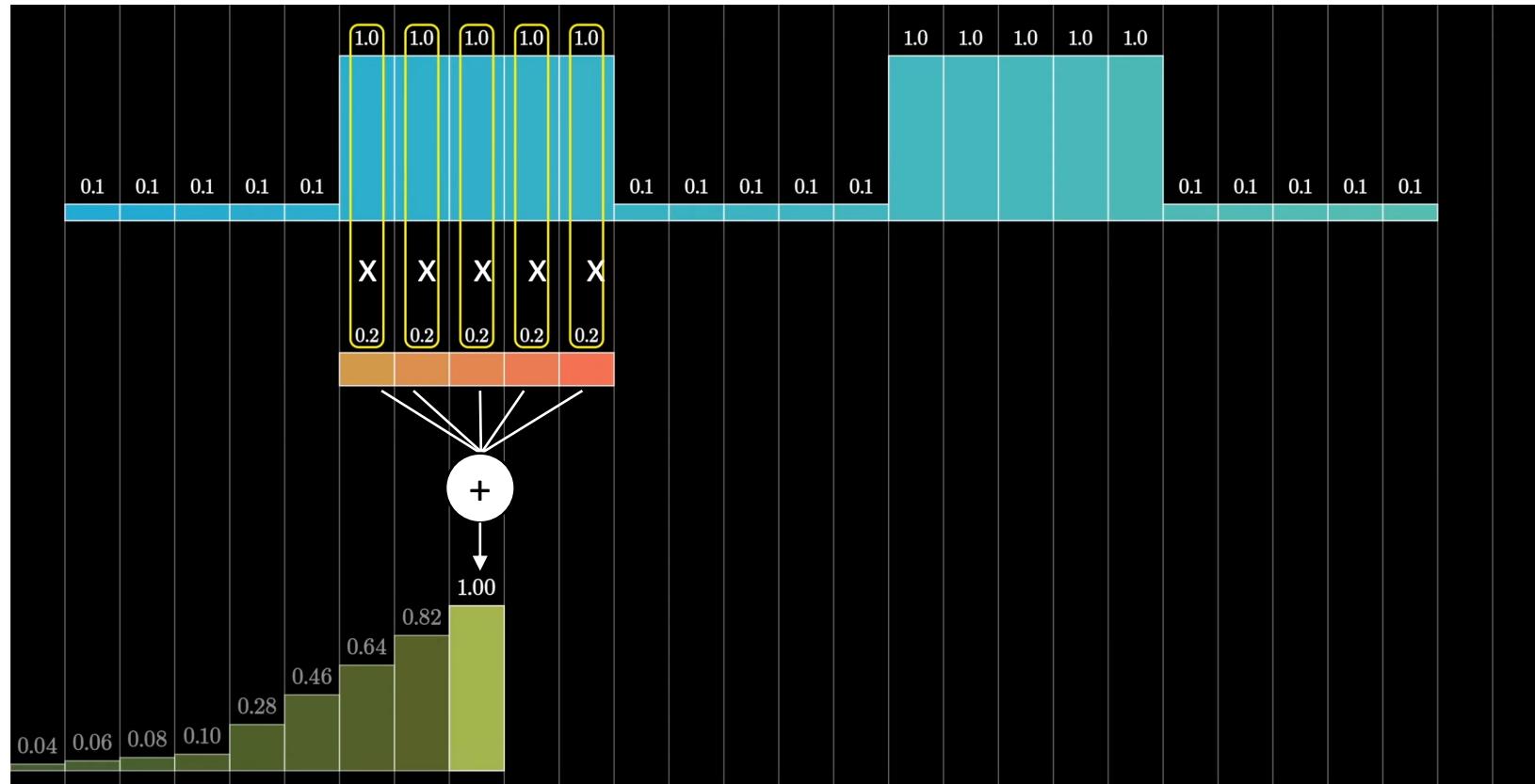
What kind of operation?



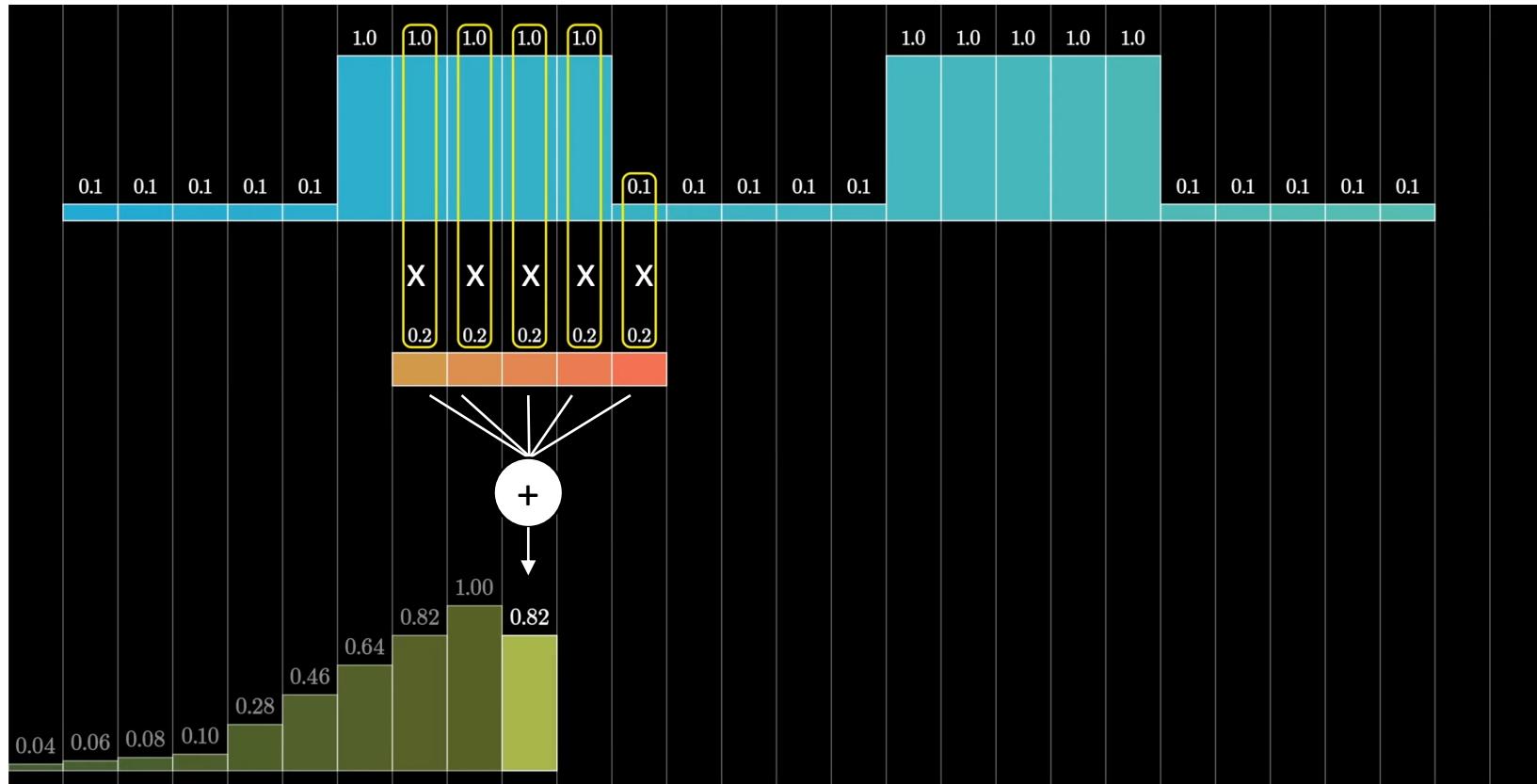
What kind of operation?



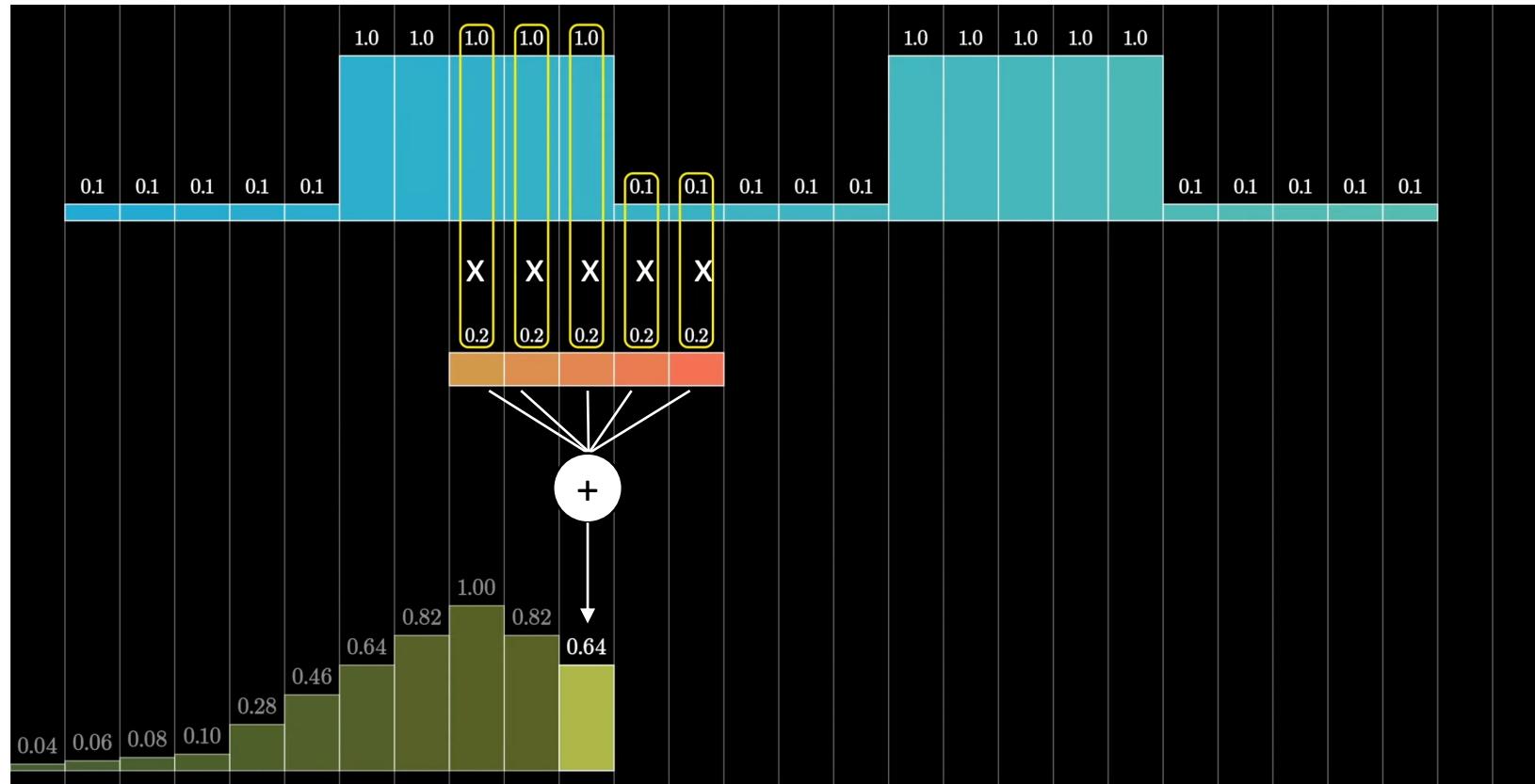
What kind of operation?



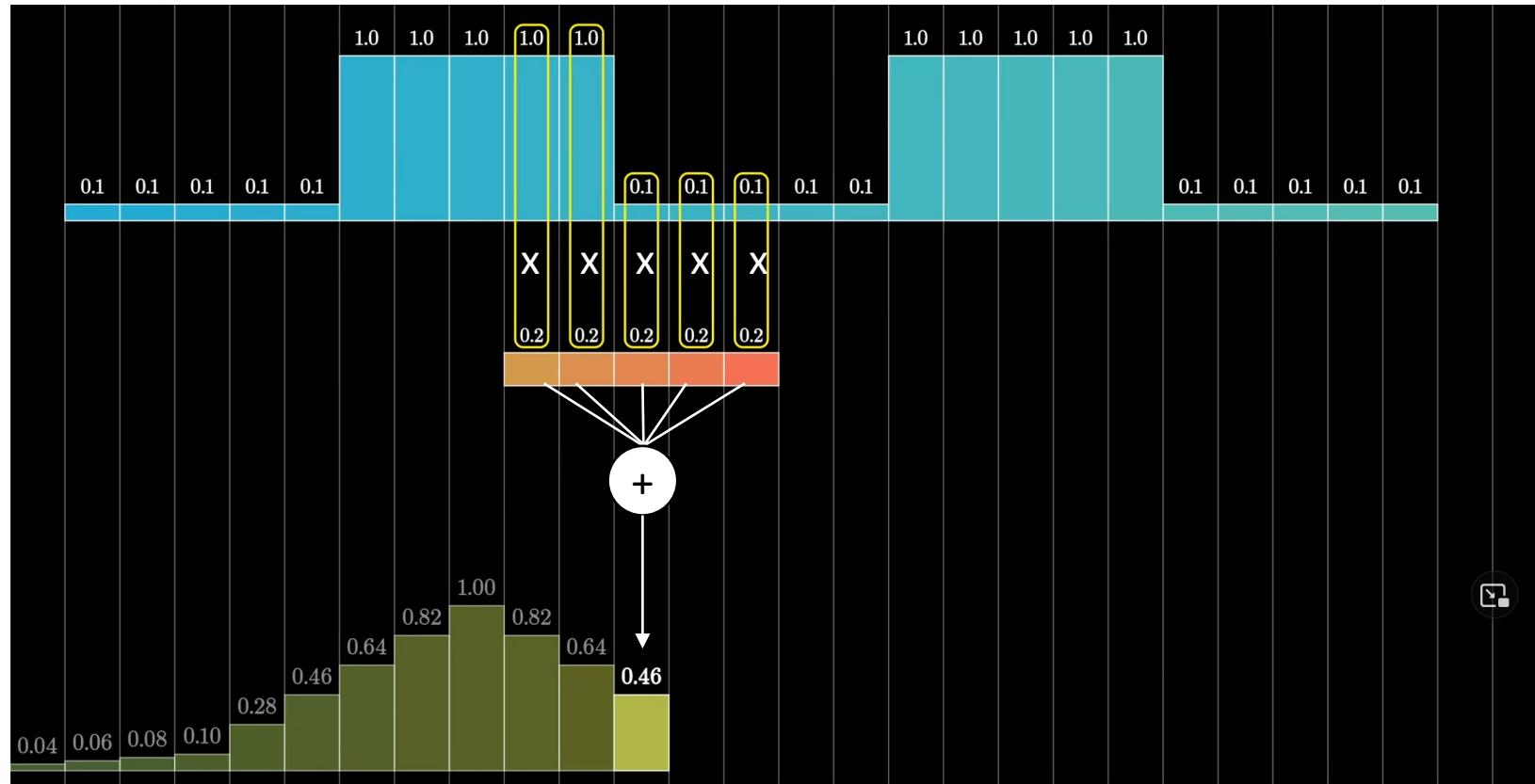
What kind of operation?



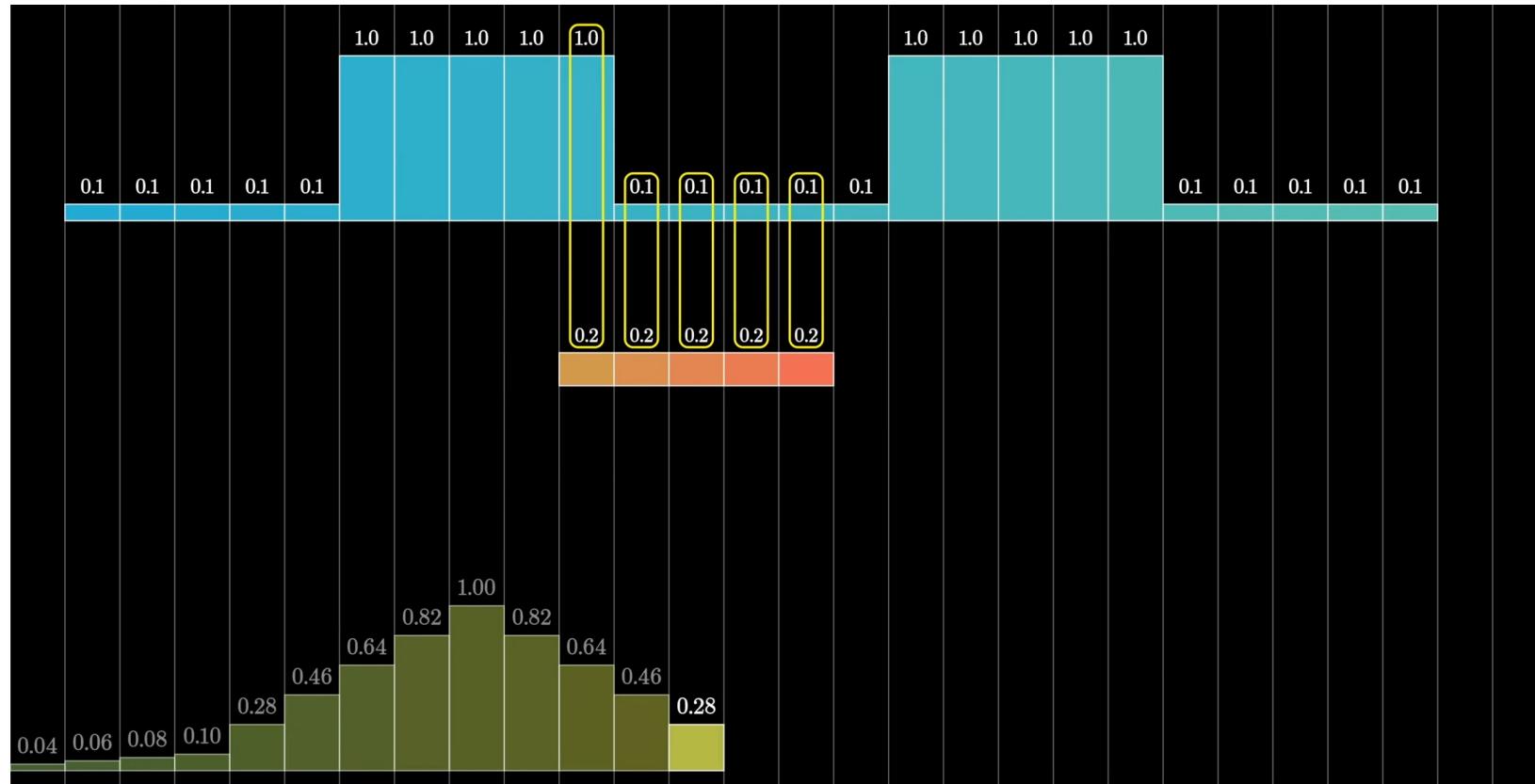
What kind of operation?



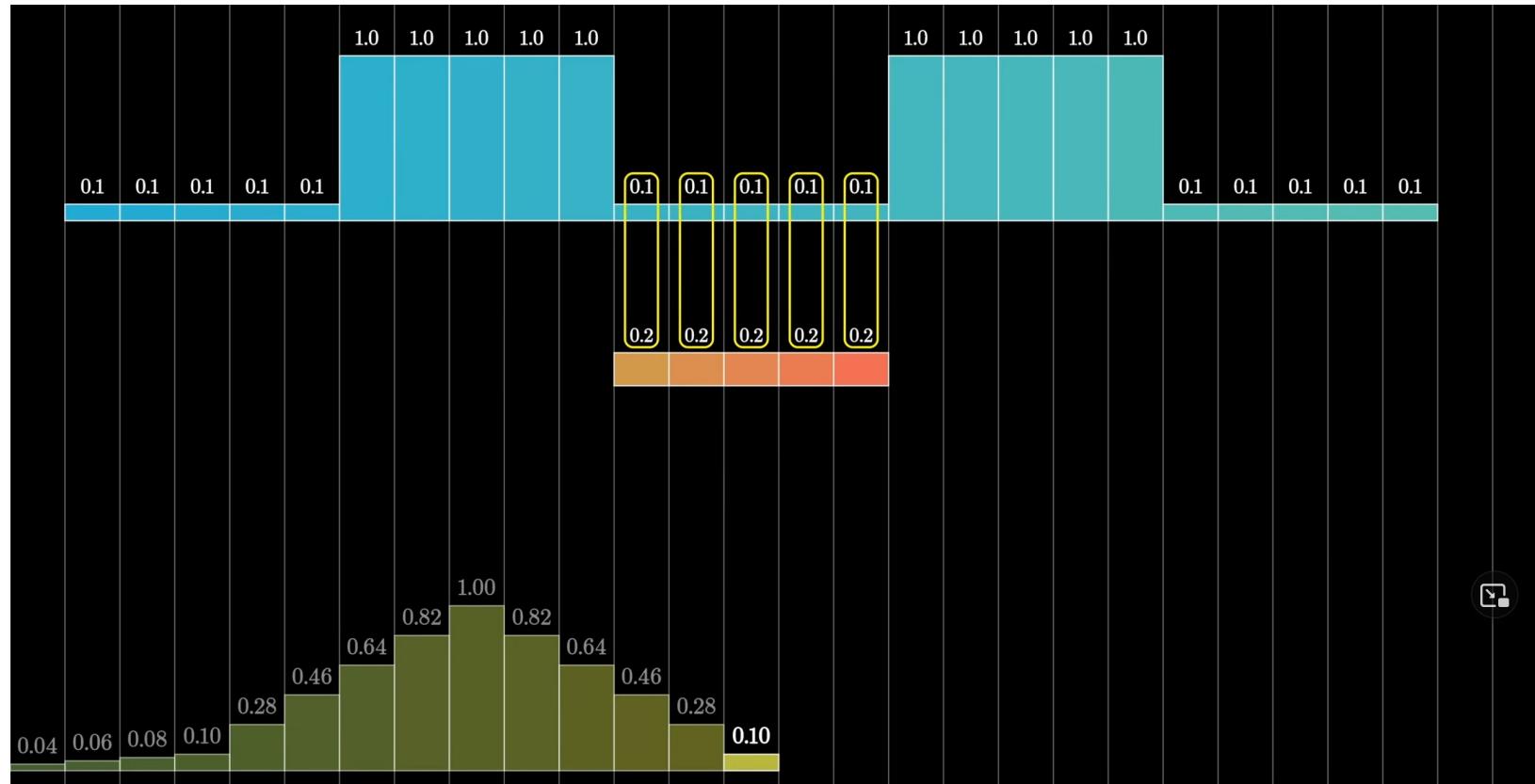
What kind of operation?



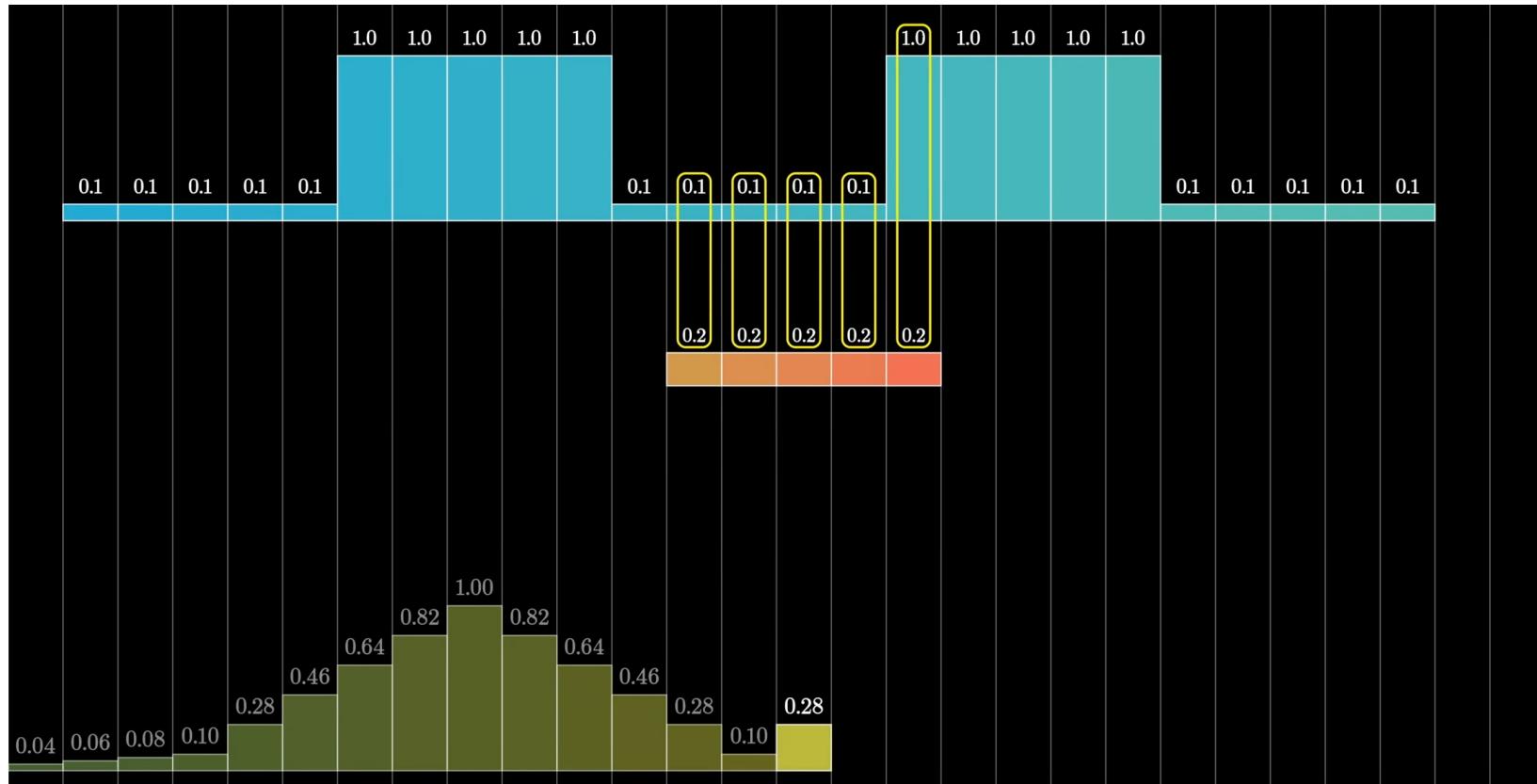
What kind of operation?



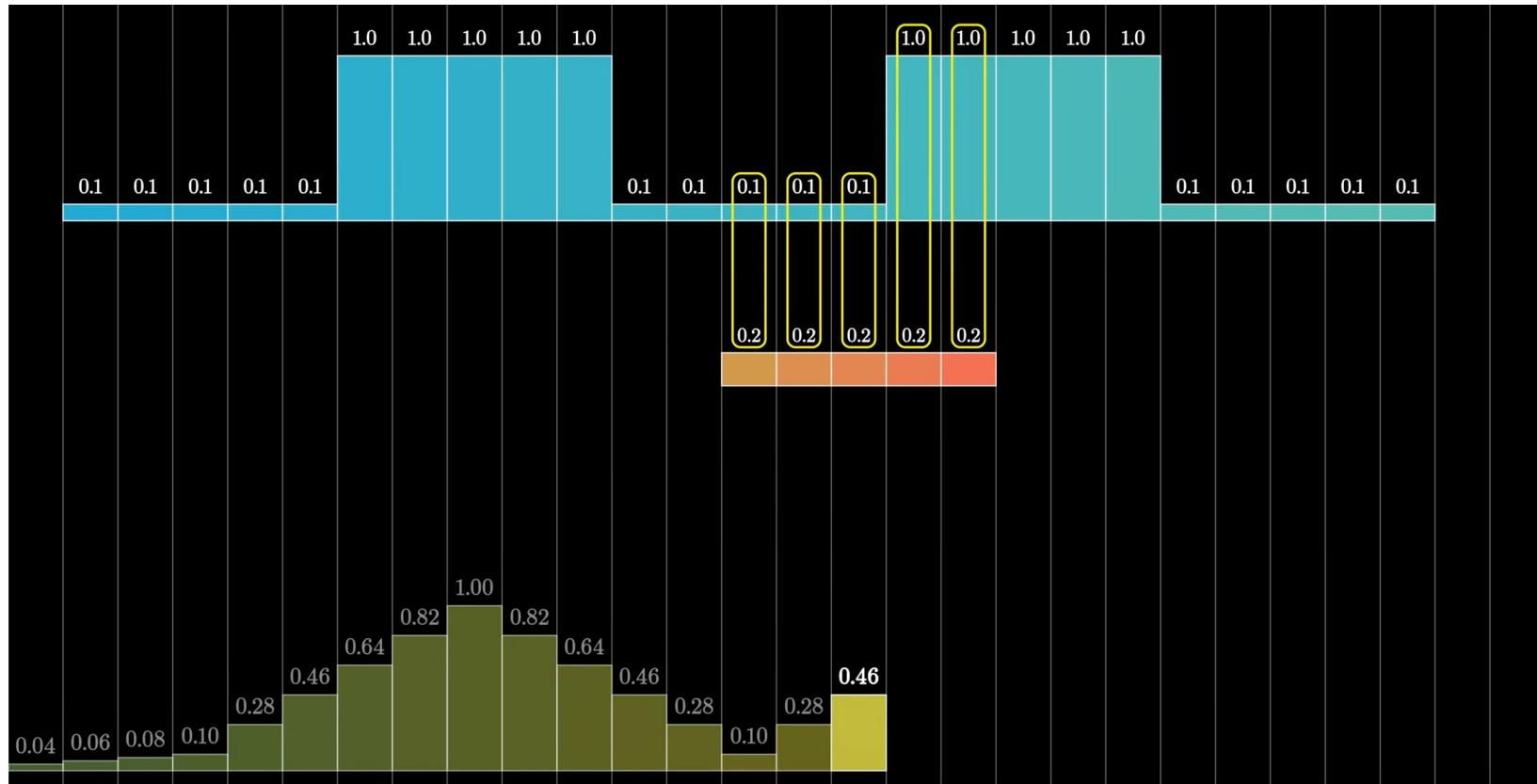
What kind of operation?



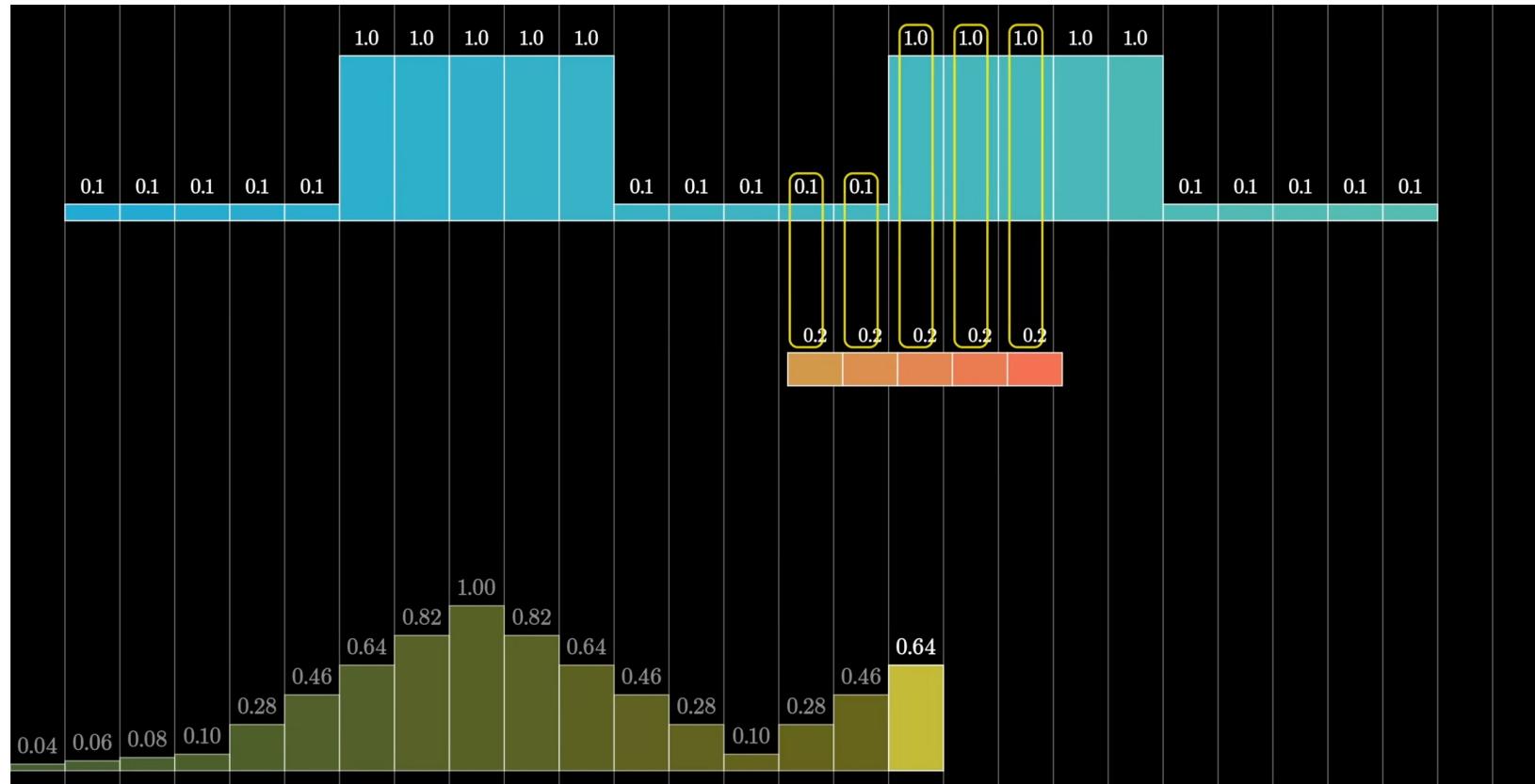
What kind of operation?



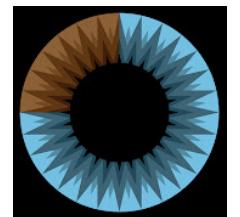
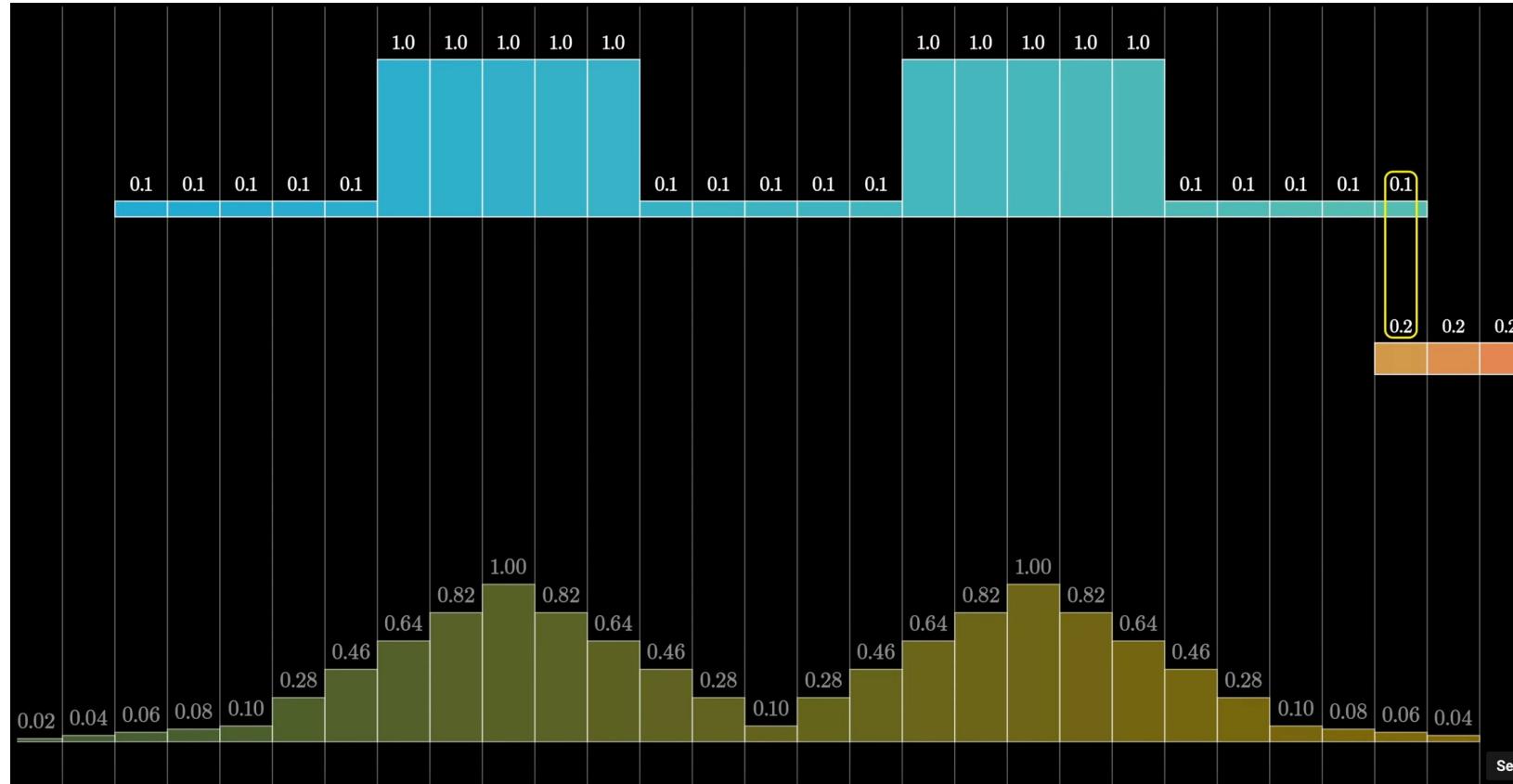
What kind of operation?



What kind of operation?



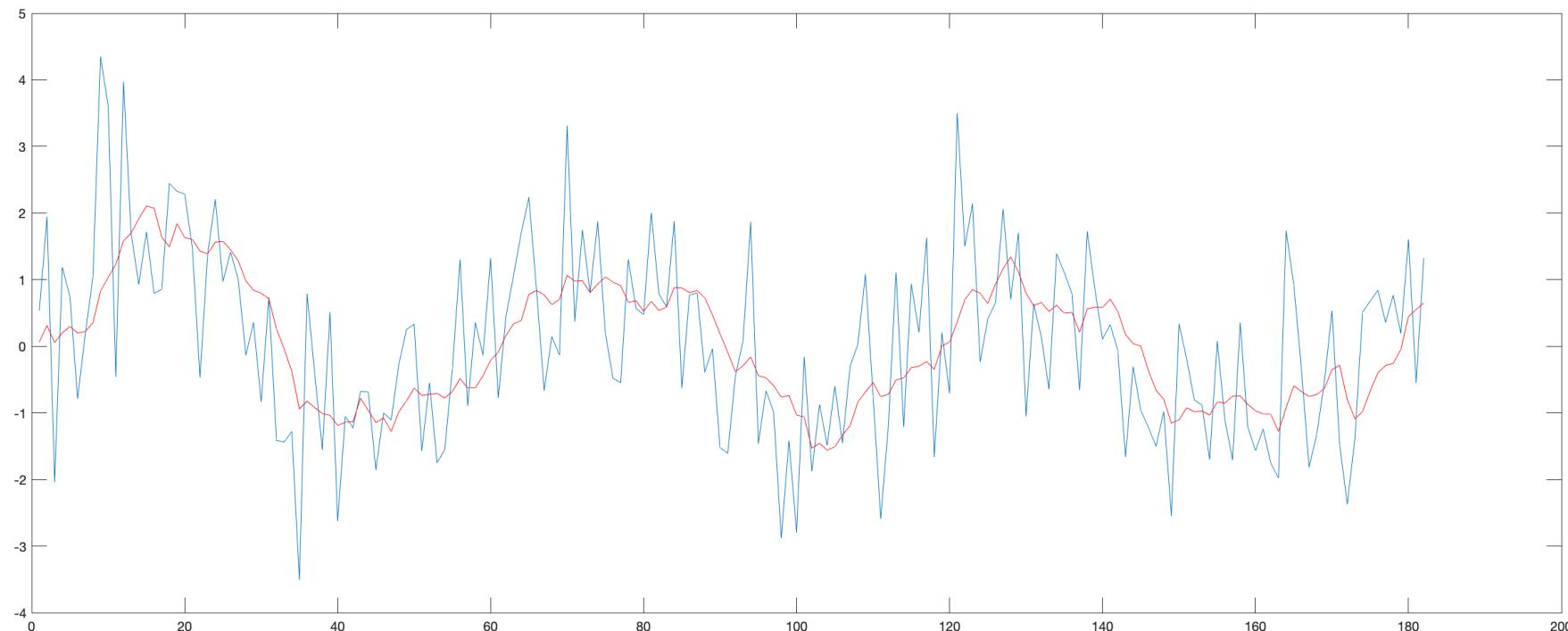
What kind of operation?



3Blue1Brown

<https://youtu.be/KuXjwB4LzSA>

Moving average – 1D Convolution



`filter = 1/8 * [1 1 1 1 1 1 1 1];`

2D Convolution

Input

7	2	3	3	8
4	5	3	8	4
3	3	2	8	4
2	8	7	2	7
5	4	4	5	4

*

1	0	-1
1	0	-1
1	0	-1

=

Feature/Activation map

6		

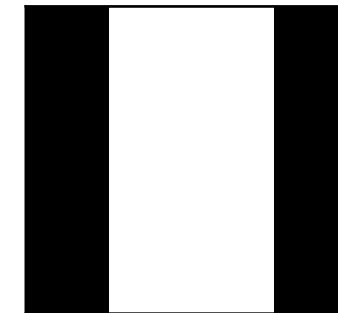
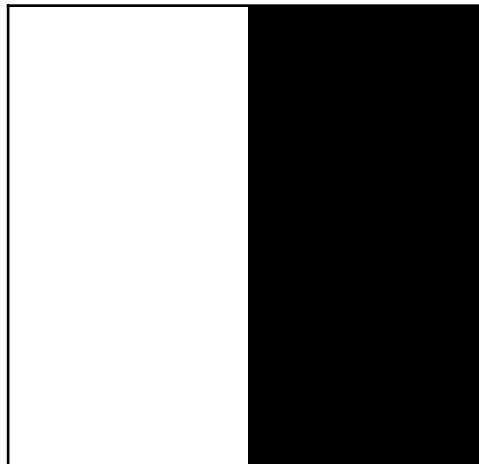
$$\begin{aligned} & 7 \times 1 + 4 \times 1 + 3 \times 1 + \\ & 2 \times 0 + 5 \times 0 + 3 \times 0 + \\ & 3 \times -1 + 3 \times -1 + 2 \times -1 \\ & = 6 \end{aligned}$$

Handcrafted filter for specific task

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

$$\begin{matrix} * & \begin{matrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{matrix} & = & \begin{matrix} 0 & 3 & 3 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 3 & 3 & 0 \\ 0 & 3 & 3 & 0 \end{matrix} \end{matrix}$$

Vertical edge detection filter



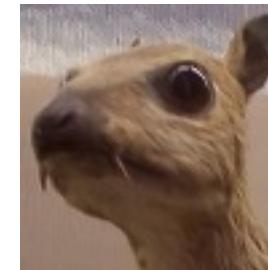
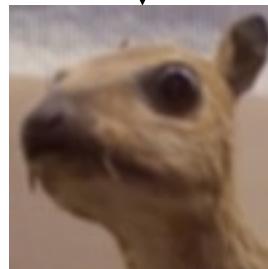
There are many specialized filters



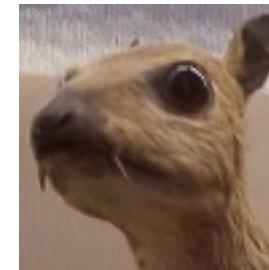
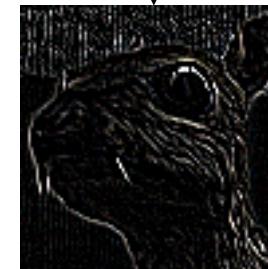
$$\text{Identity} \quad * \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$



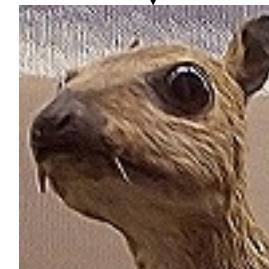
$$\text{Blur} \quad * \begin{bmatrix} 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \end{bmatrix}$$



$$\text{Edge detection} \quad * \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



$$\text{Sharpen} \quad * \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$



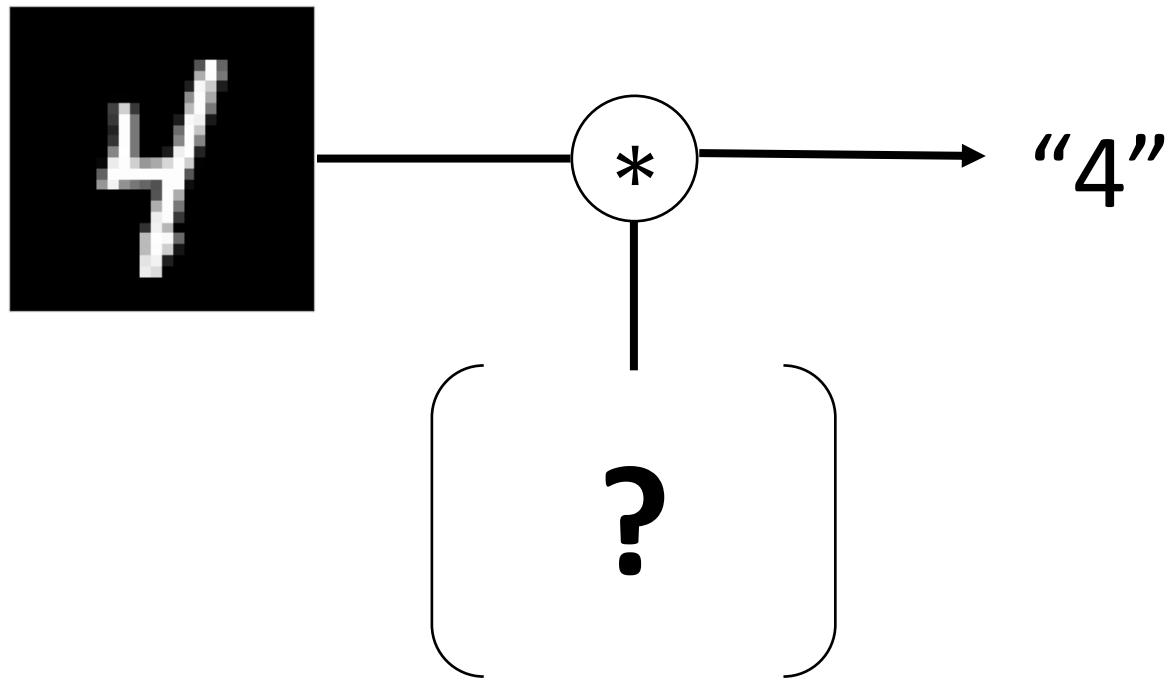
Identity

Blur

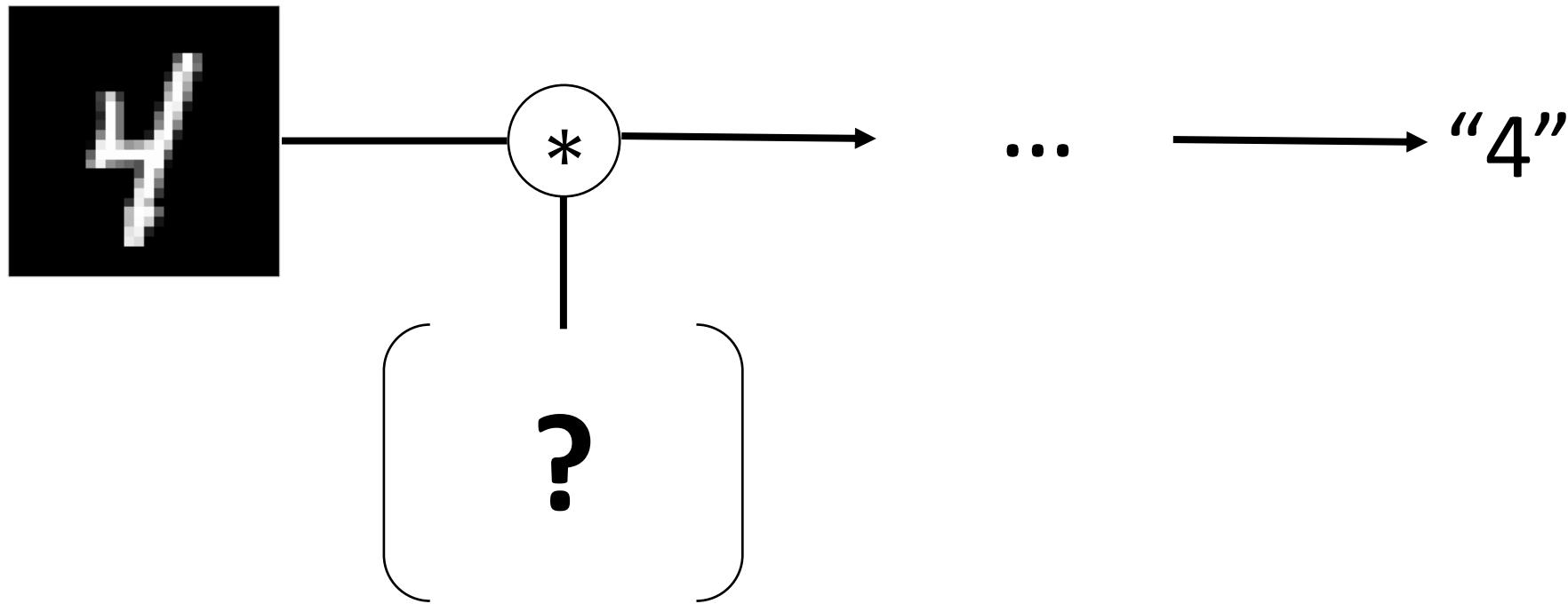
Edge detection

Sharpen

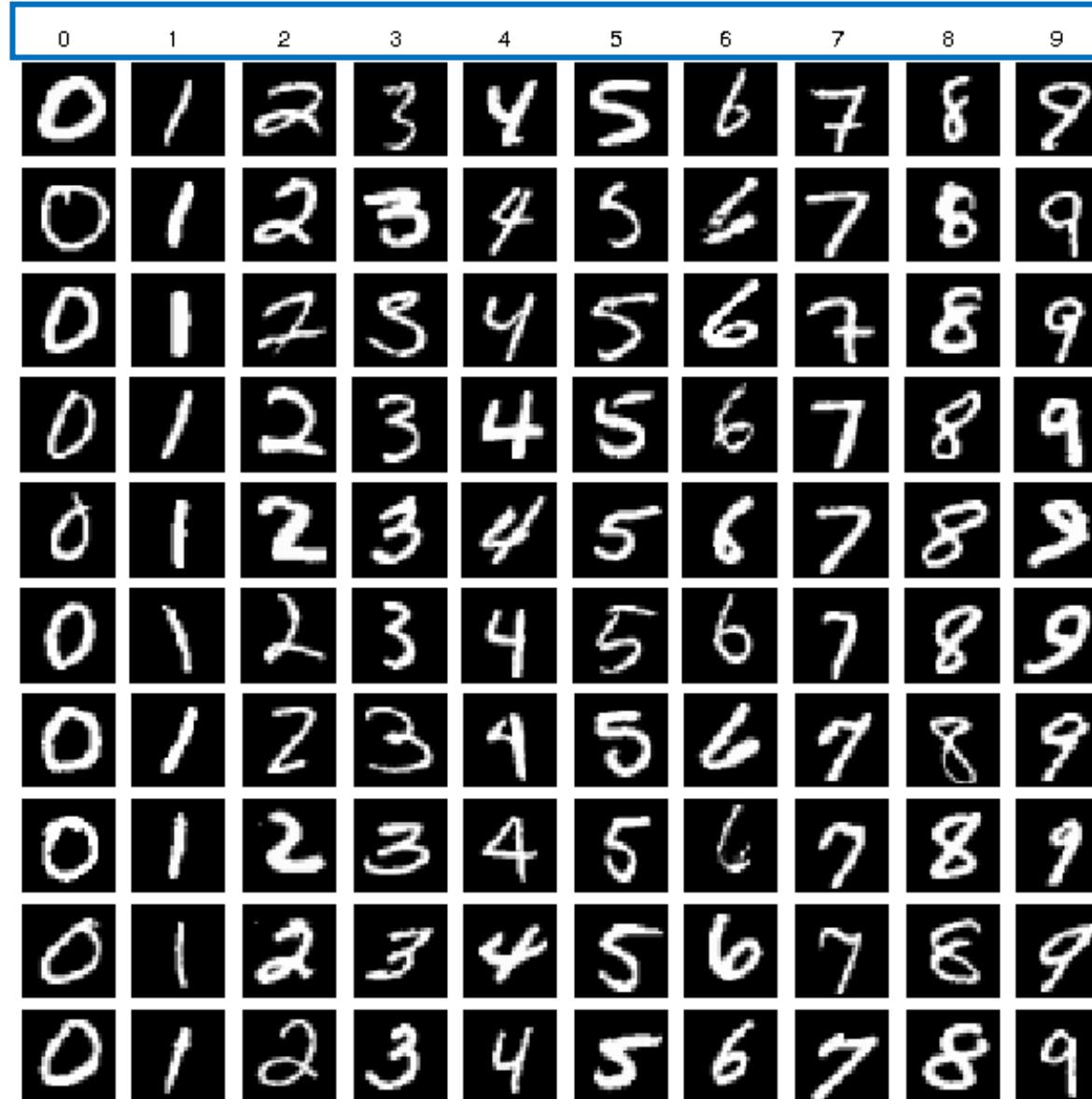
Filter for digit detection?



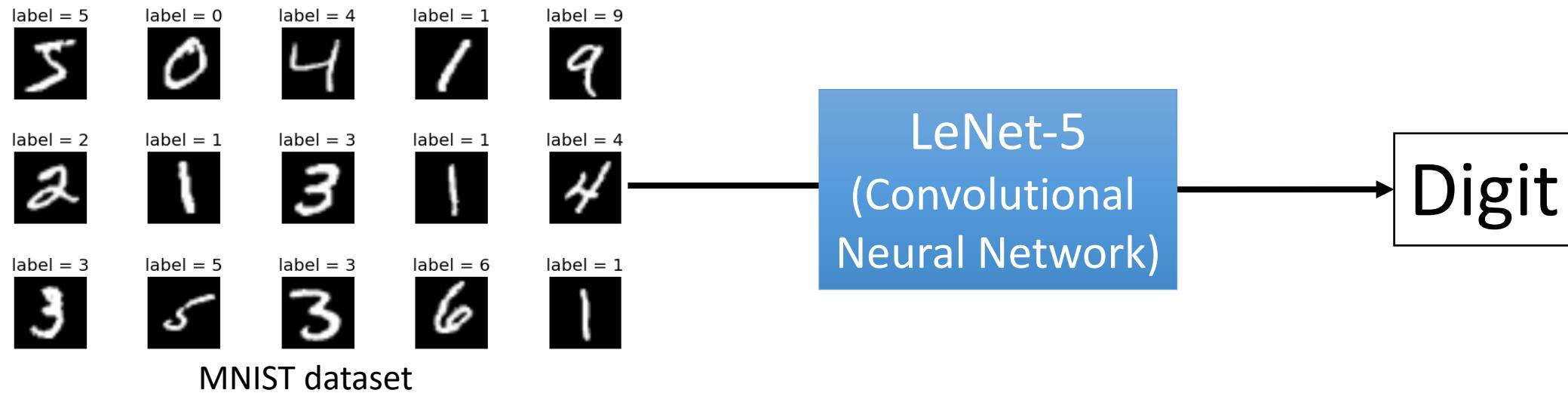
Filter for digit detection?



Let machine learn

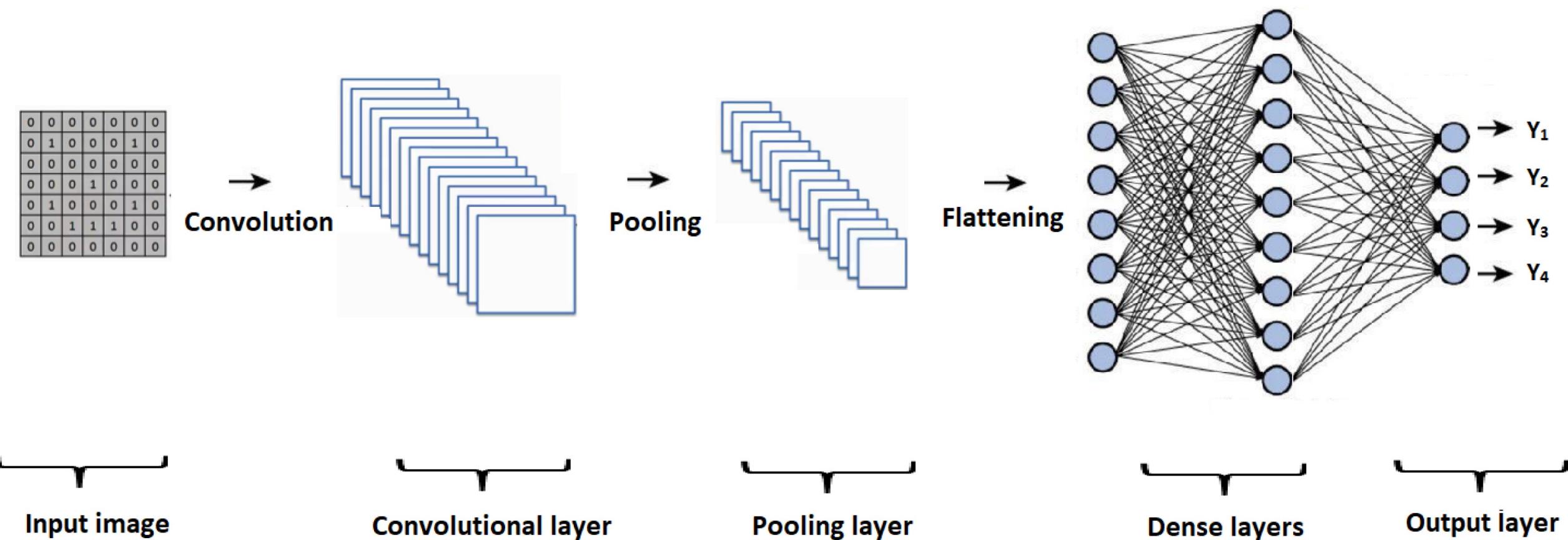


Algorithm meets Data meets Compute

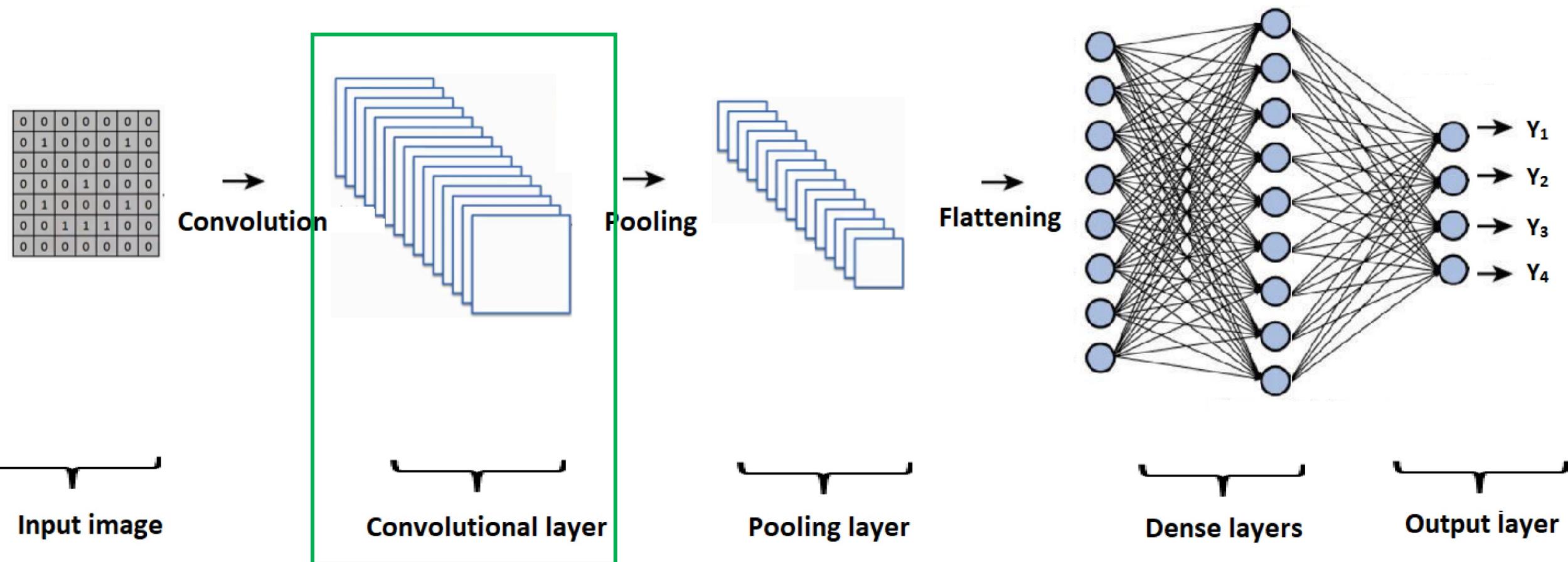


Y. Lecun, L. Bottou, Y. Bengio and P. Haffner, "Gradient-based learning applied to document recognition," in *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278-2324, Nov. 1998, doi: 10.1109/5.726791.

Convolutional Neural Network Architecture



Convolutional Neural Network Architecture



Convolutional layer

- Convolution operation

Input

7	2	3	3	8
4	5	3	8	4
3	3	2	8	4
2	8	7	2	7
5	4	4	5	4

* Kernel/Filter

1	0	-1
1	0	-1
1	0	-1

= Feature/Activation map

6		

$$7 \times 1 + 4 \times 1 + 3 \times 1 + \\ 2 \times 0 + 5 \times 0 + 3 \times 0 + \\ 3 \times -1 + 3 \times -1 + 2 \times -1 \\ = 6$$

Convolutional layer

- Padding

0	0	0	0	0	0	0
0	60	113	56	139	85	0
0	73	121	54	84	128	0
0	131	99	70	129	127	0
0	80	57	115	69	134	0
0	104	126	123	95	130	0
0	0	0	0	0	0	0

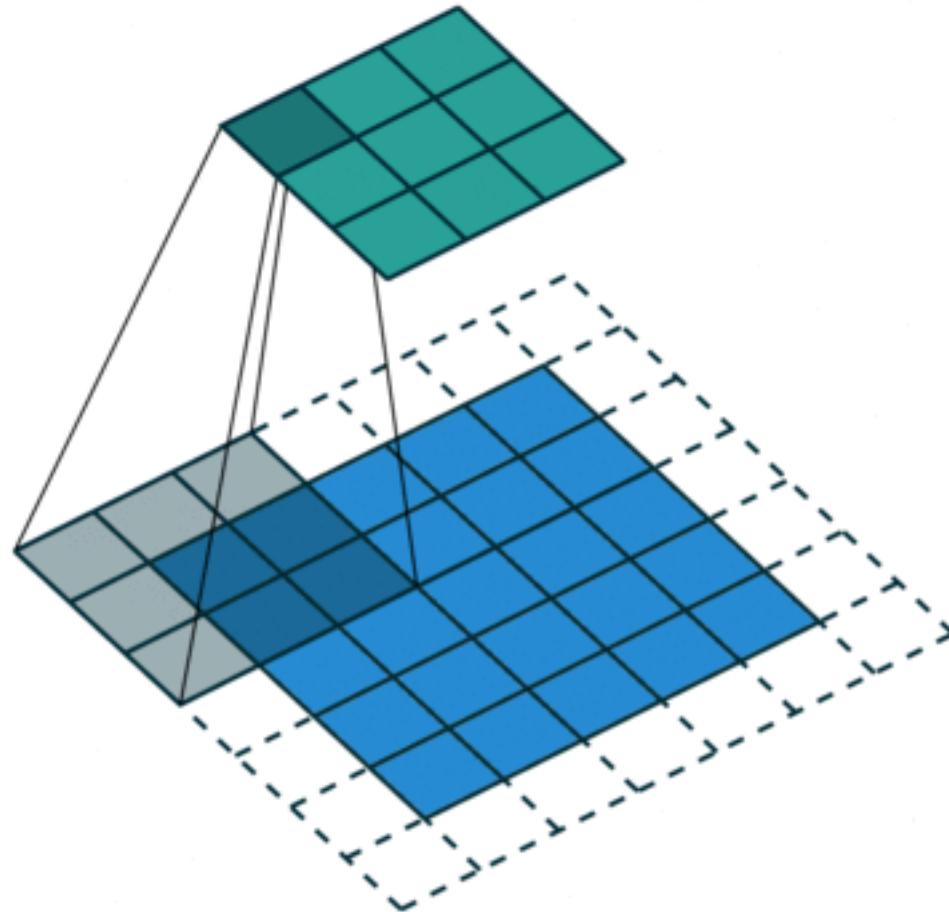
Kernel

0	-1	0
-1	5	-1
0	-1	0

114				

Convolutional layer

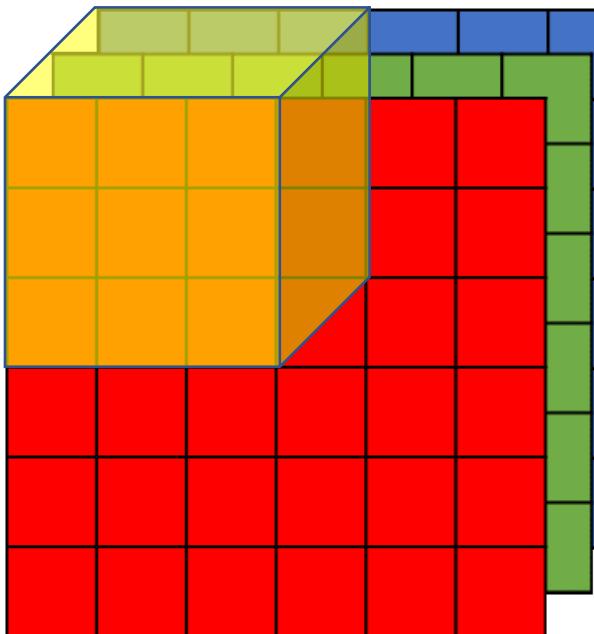
- Stride



Example: Padding = 1, Stride = 2

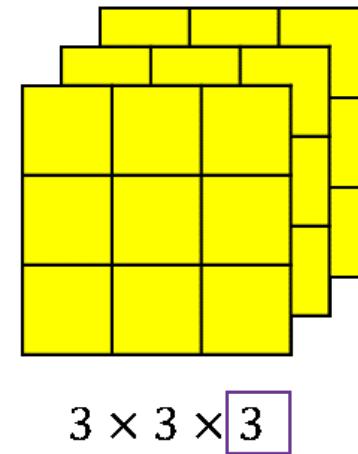
Convolutional layer

- Convolution over volume



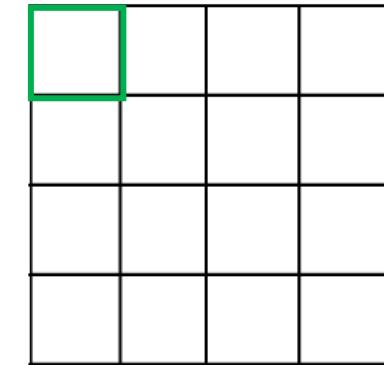
$6 \times 6 \times 3$

*



$3 \times 3 \times 3$

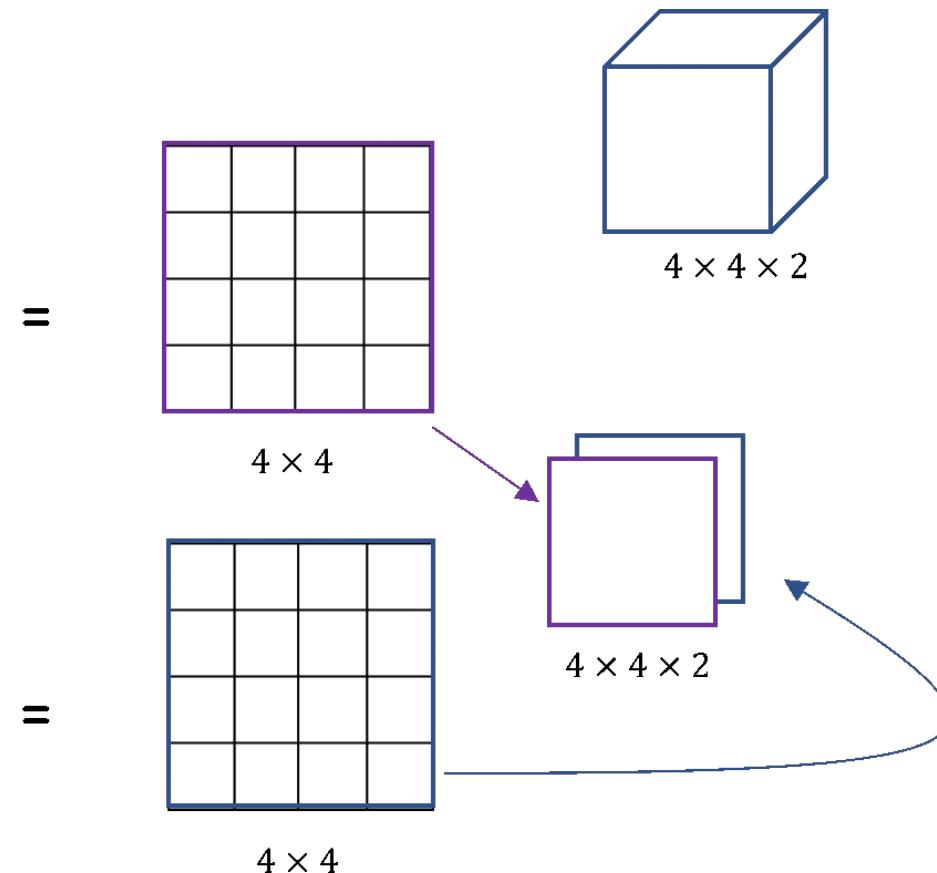
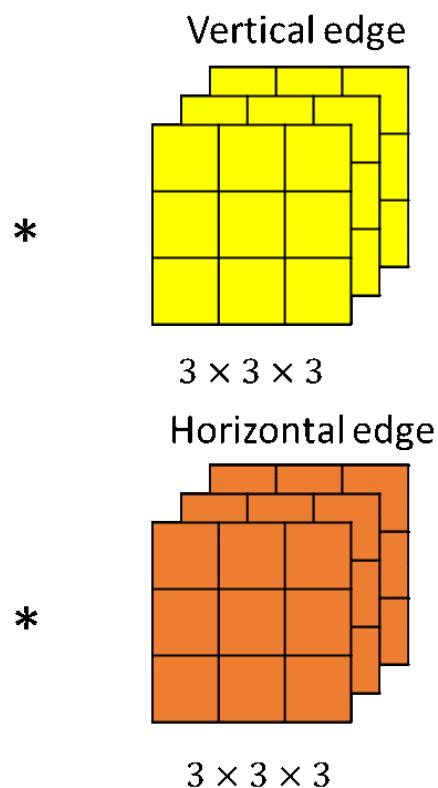
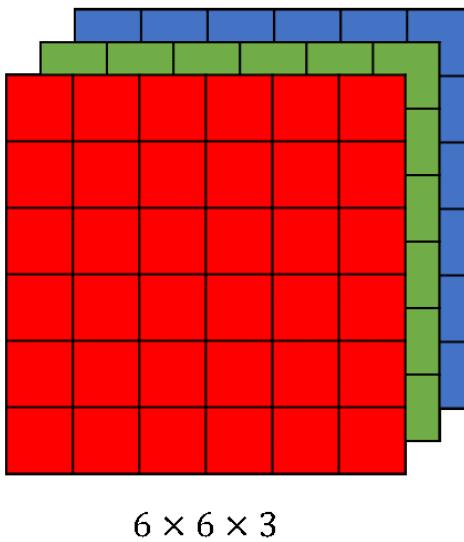
=



4×4

Convolutional layer

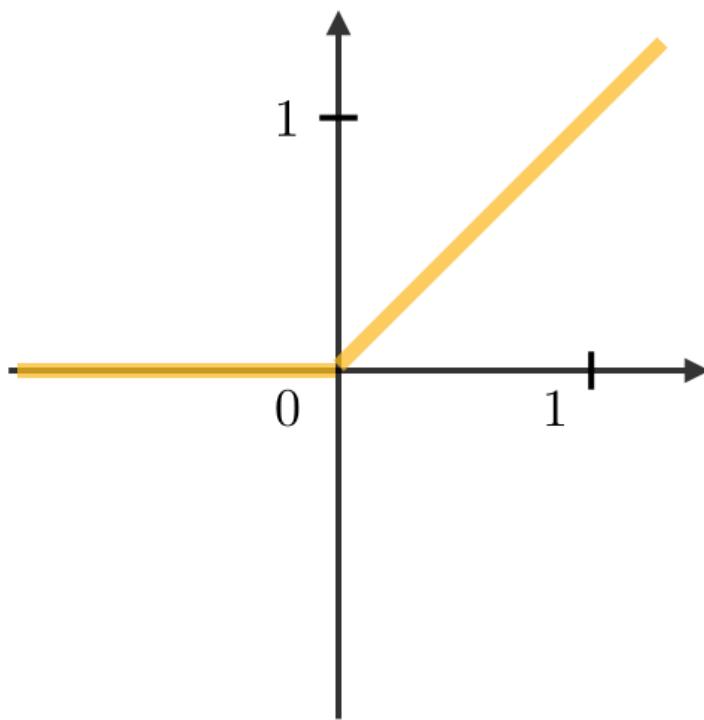
- Stacking activation maps



Activation functions

- Rectified Linear Unit (ReLU)

$$g(z) = \max(0, z)$$



Filter 1 Feature Map

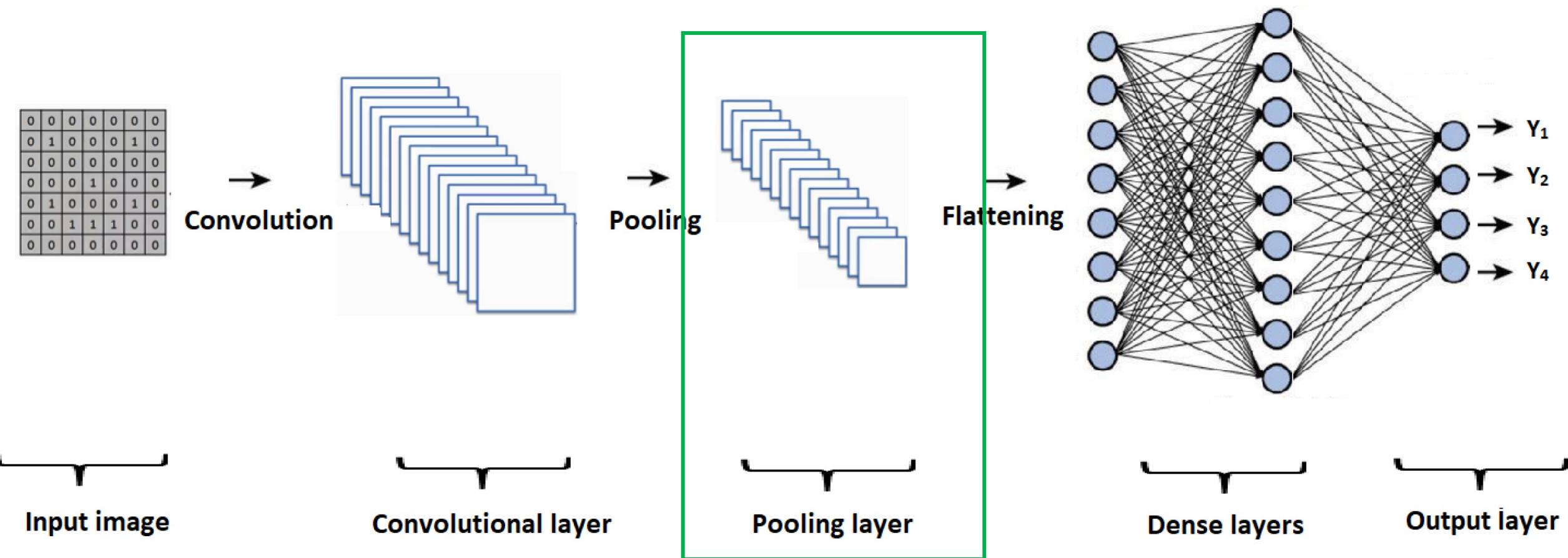
9	3	5	-8
-6	2	-3	1
1	3	4	1
3	-4	5	1

ReLU Layer

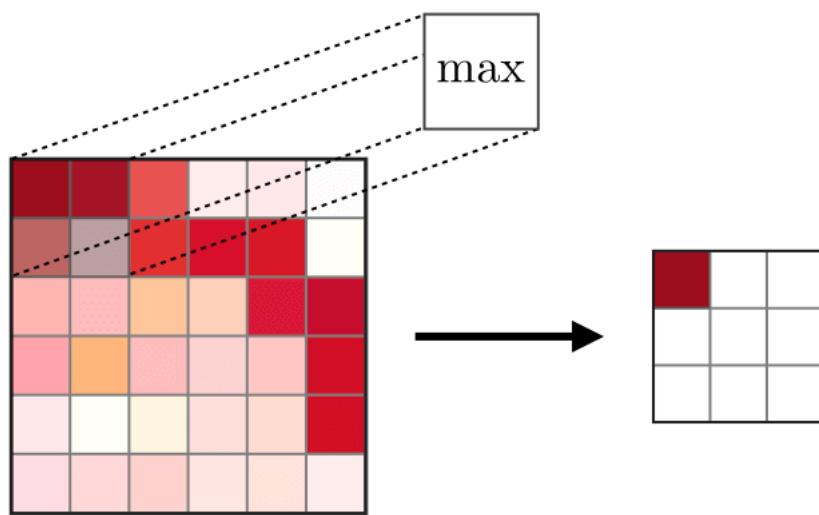


9	3	5	0
0	2	0	1
1	3	4	1
3	0	5	1

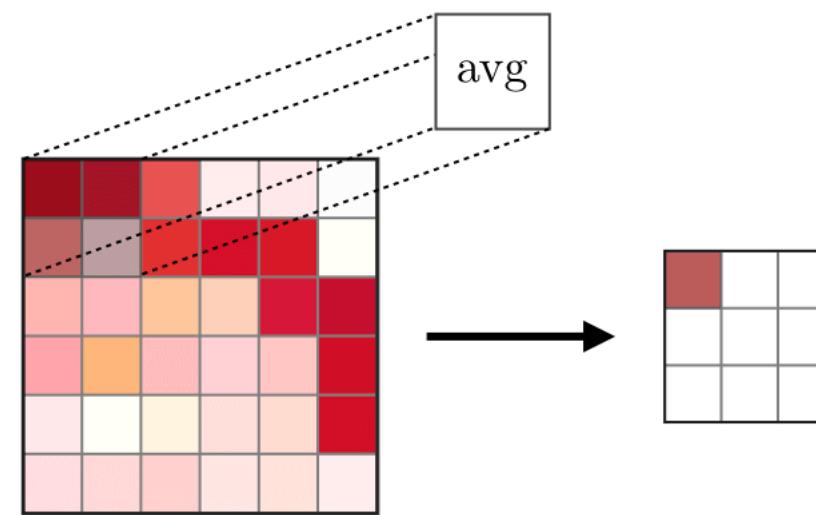
Convolutional Neural Network Architecture



Pooling layer

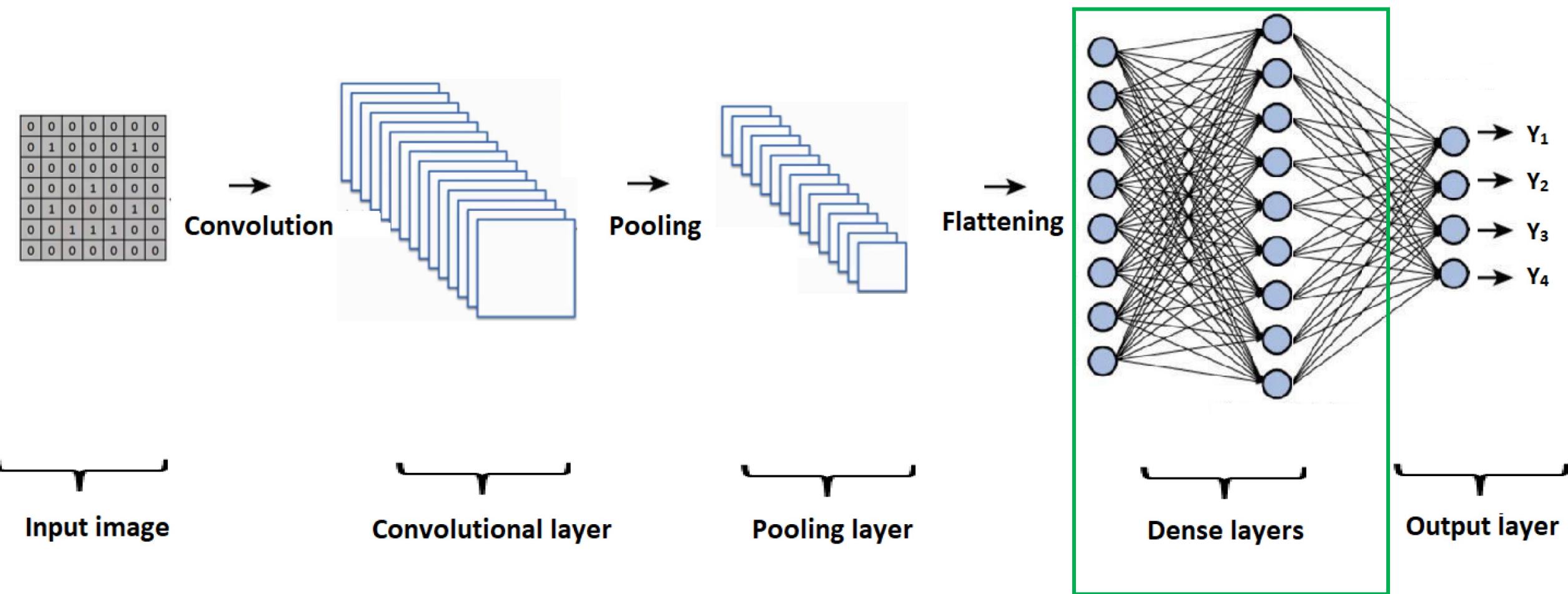


max



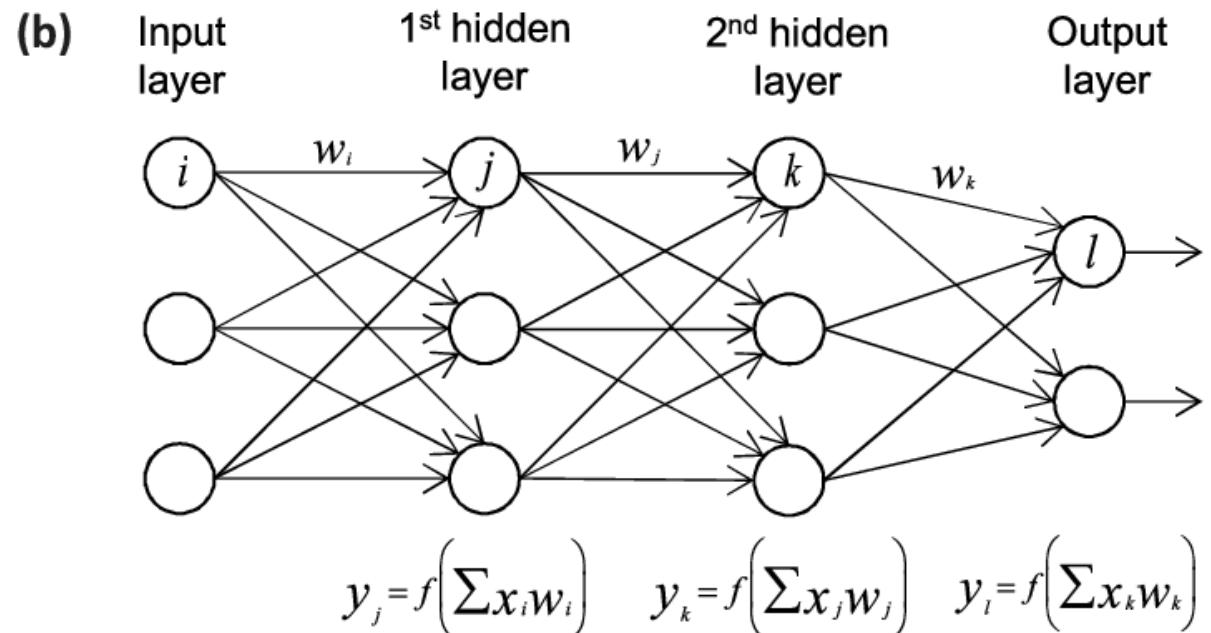
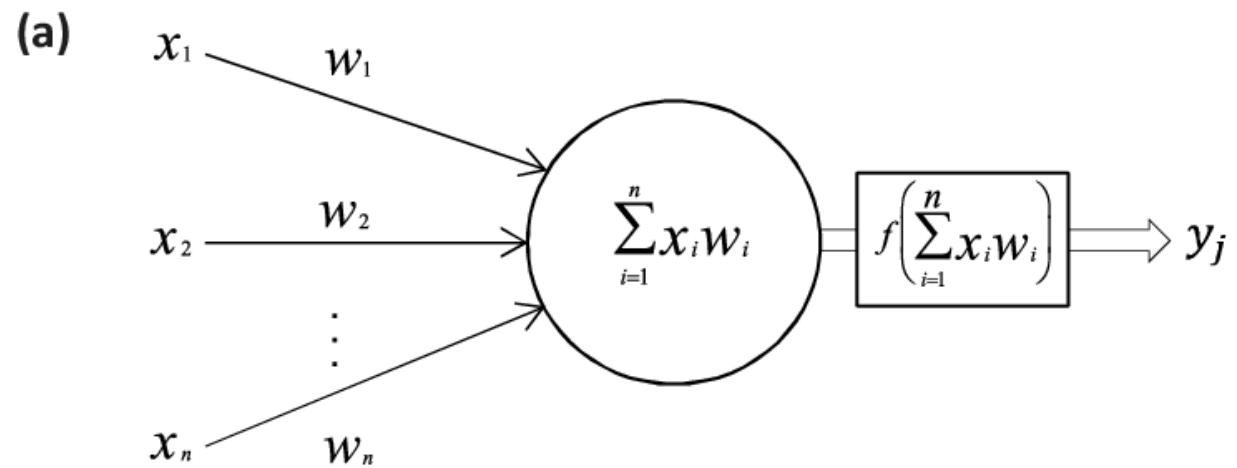
avg

Convolutional Neural Network Architecture



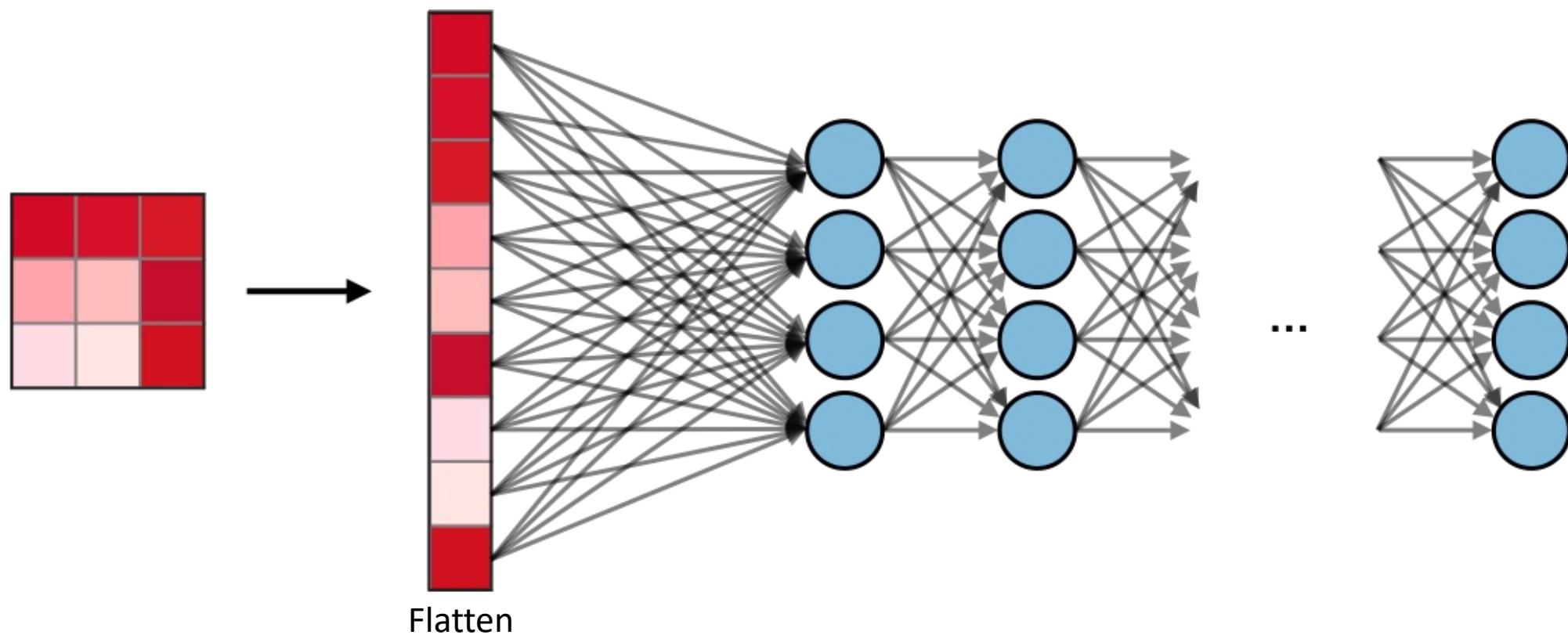
Fully connected layer

- Also known as dense layer
- Each input is connected to all hidden units

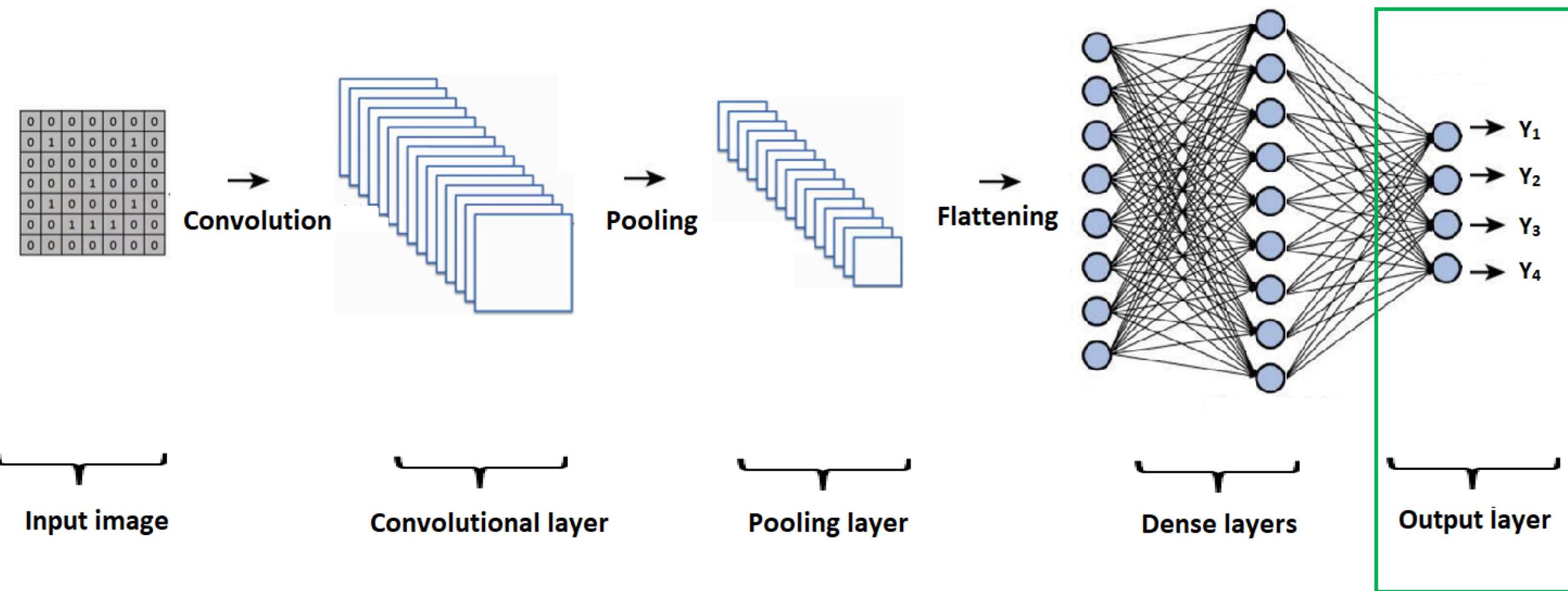


Fully connected layer

- Flattening



Convolutional Neural Network Architecture



Activation function for classification

- Softmax

$$p = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix} \quad \text{where} \quad p_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

Loss function

- Cross-entropy loss

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \cdot \log(\hat{y}_i)$$

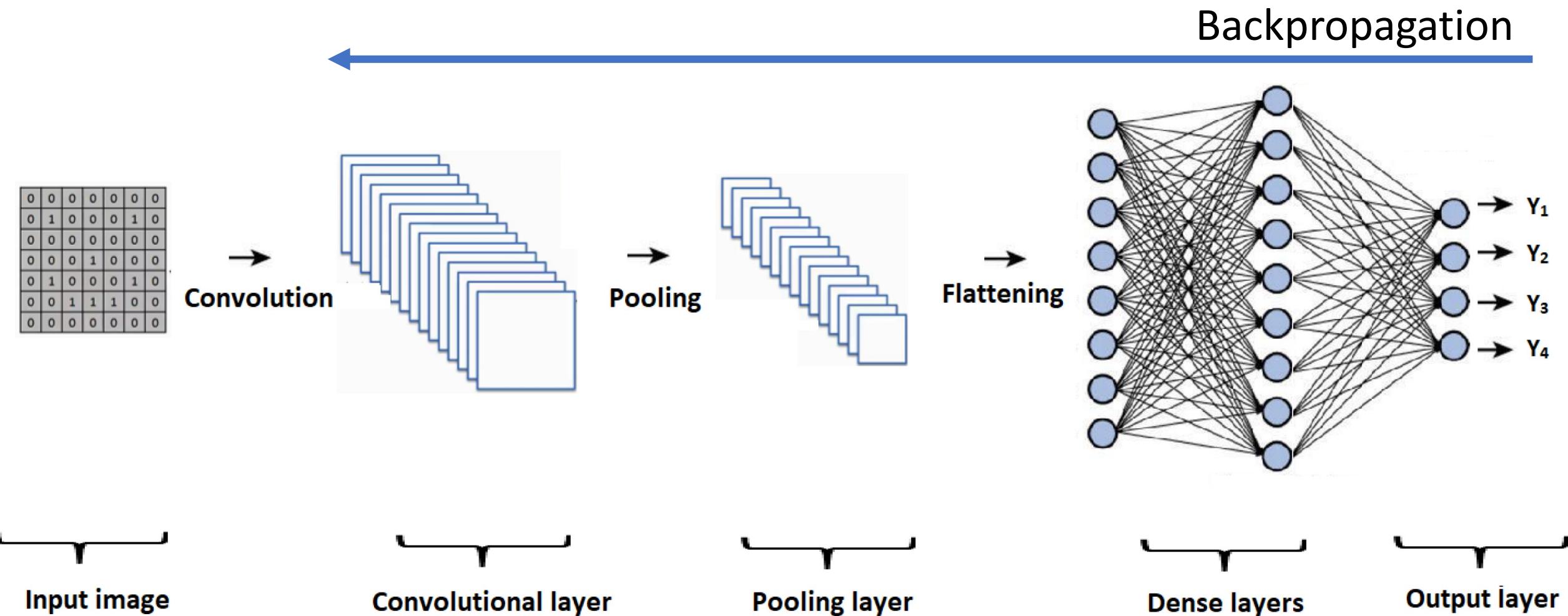
y is a one-hot vector ([0,0,1,...,0])

Gradient descent

$$\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} L$$

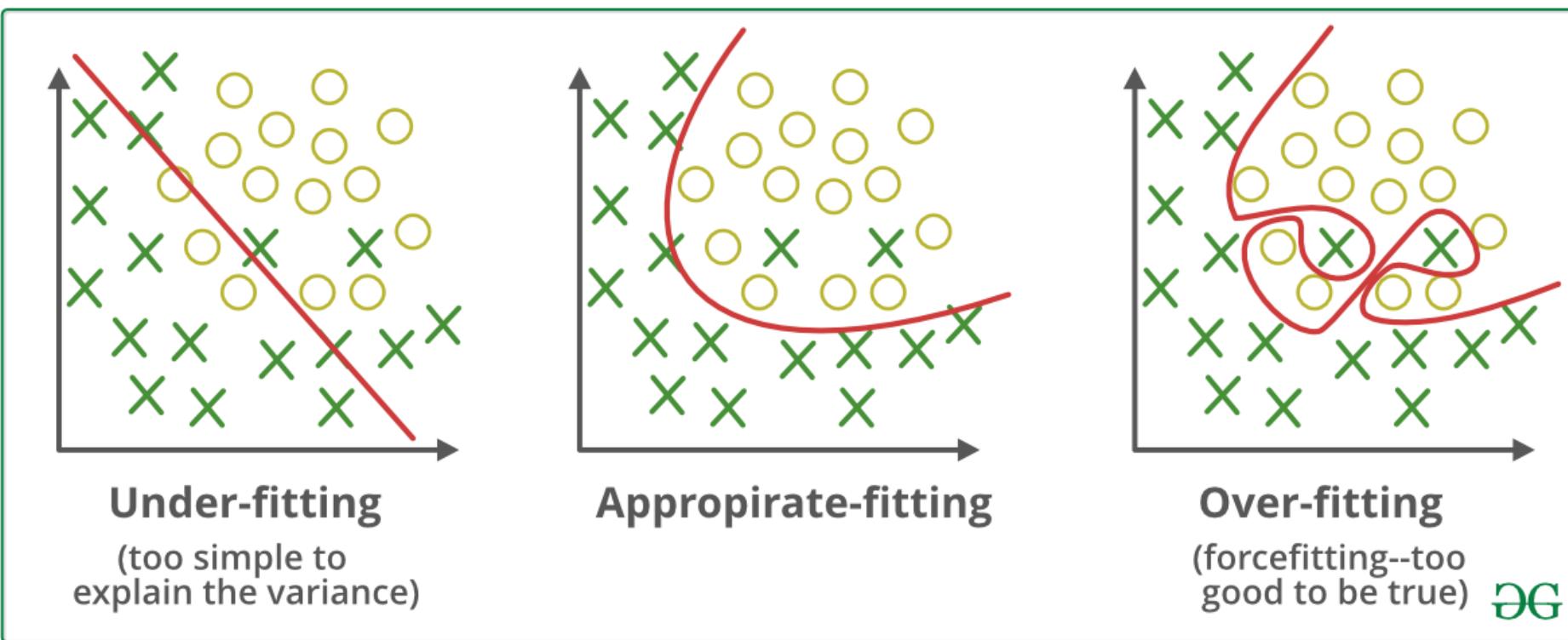
The diagram illustrates the gradient descent update rule. It features a central equation $\theta_j = \theta_j - \alpha \frac{\partial}{\partial \theta_j} L$. To the left of the first θ_j , there is an upward-pointing arrow labeled "Parameters". To the right of the α term, there is another upward-pointing arrow labeled "Learning rate".

Training the network

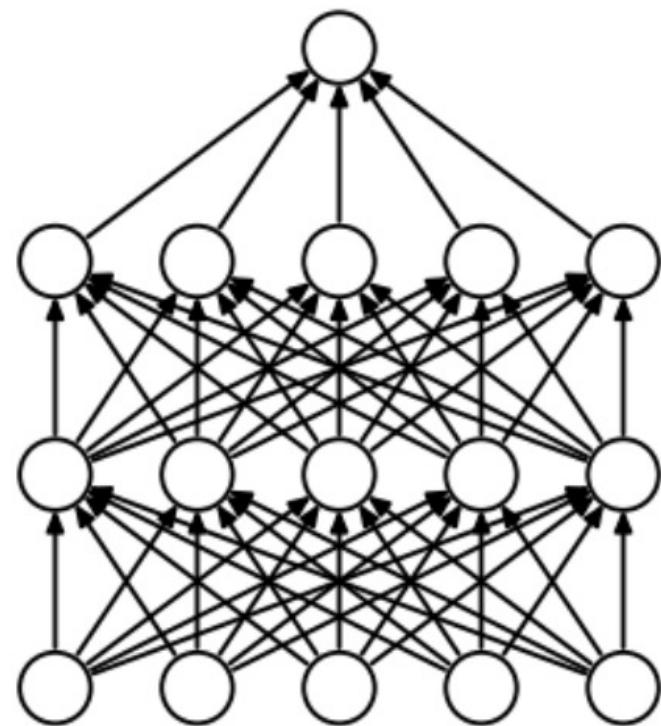


Dropout

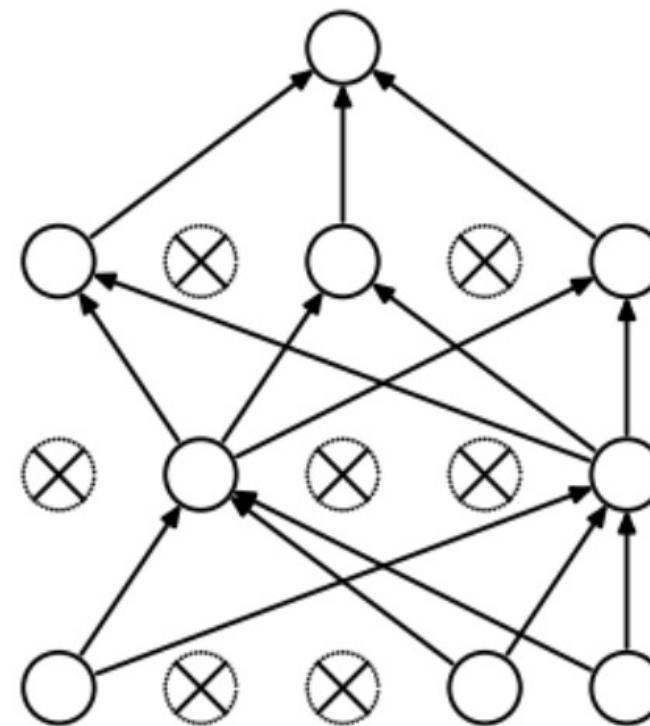
- Overfitting



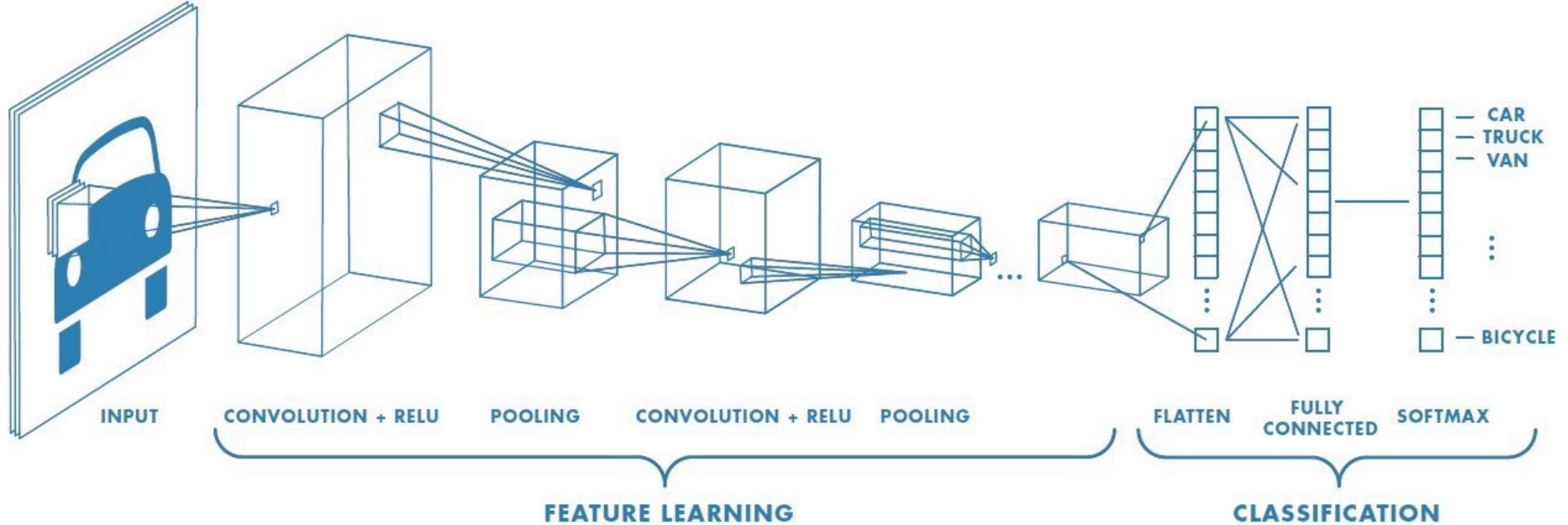
Dropout



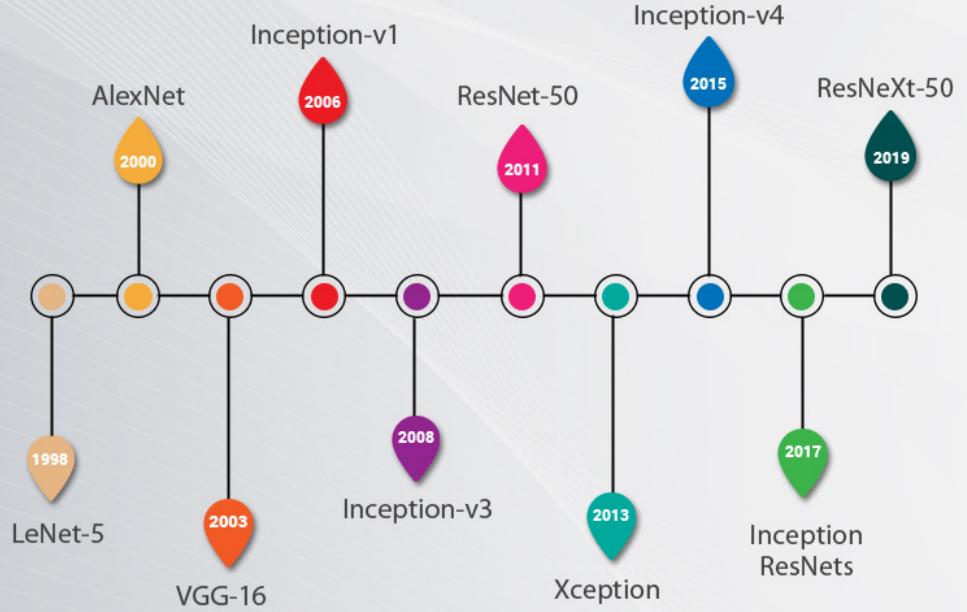
(a) Standard Neural Net



(b) After applying dropout.



CNN architectures over a timeline(1998-2019)

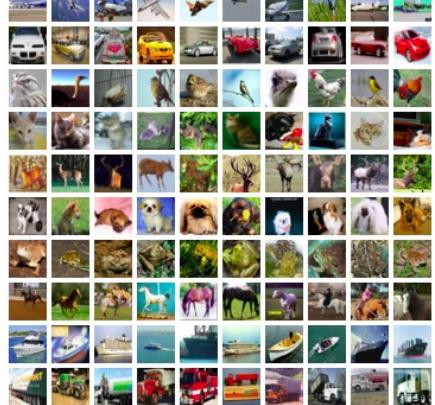


460,723 images

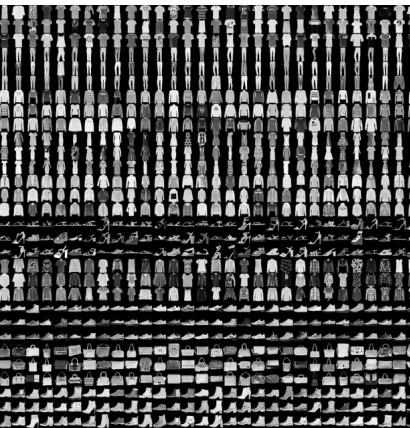


62,328 images

COCO 2020 Panoptic Segmentation Task



CIFAR-10

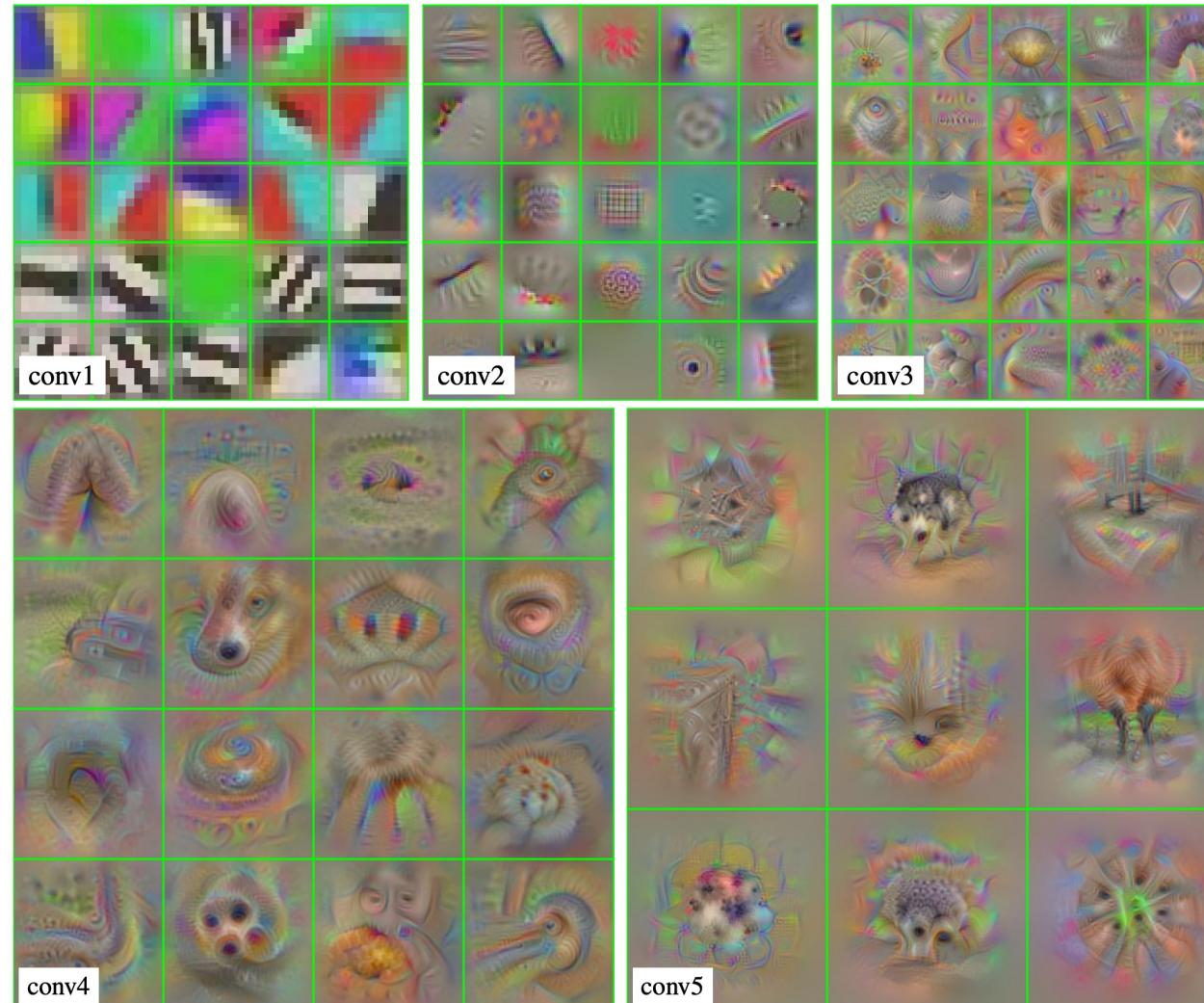


Fashion MNIST

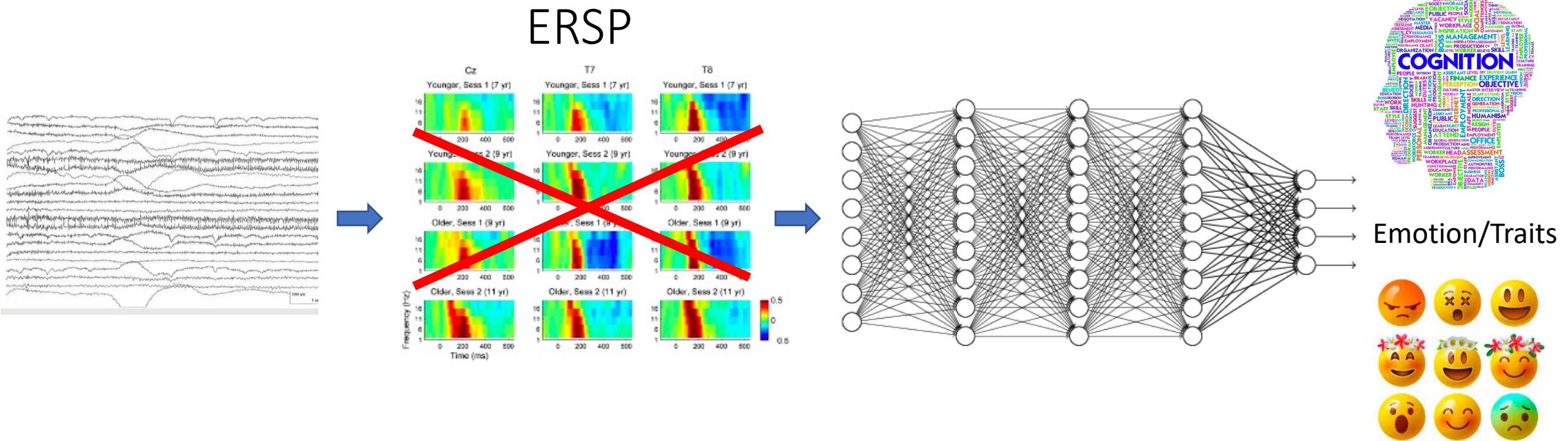
Visualizing learned filters



Visualizing learned filters



Applying to EEG



Raw EEG vs. Frequency-domain

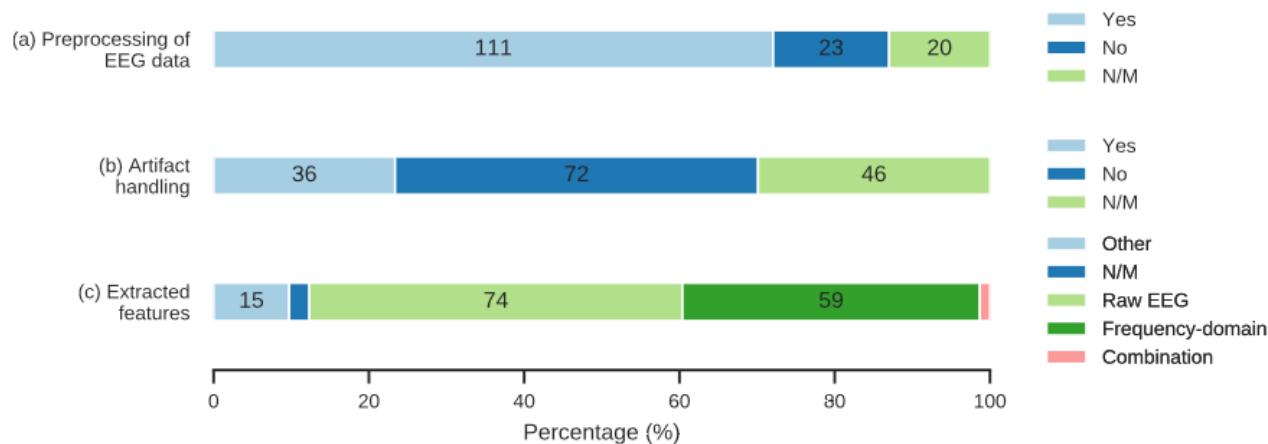
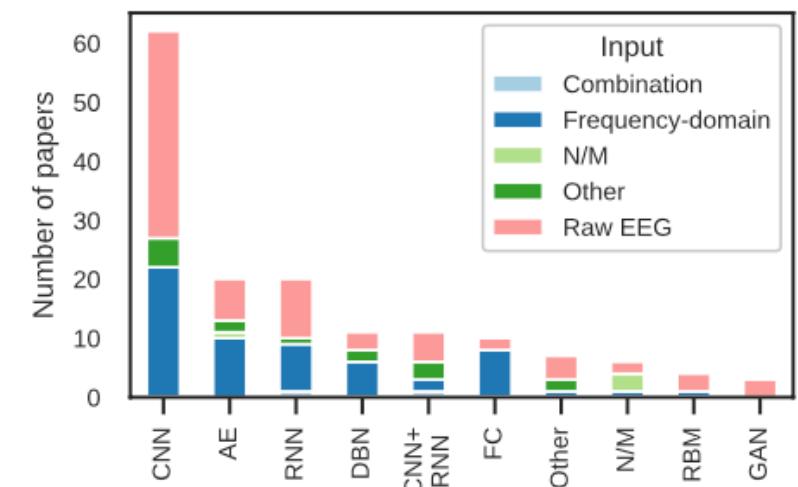
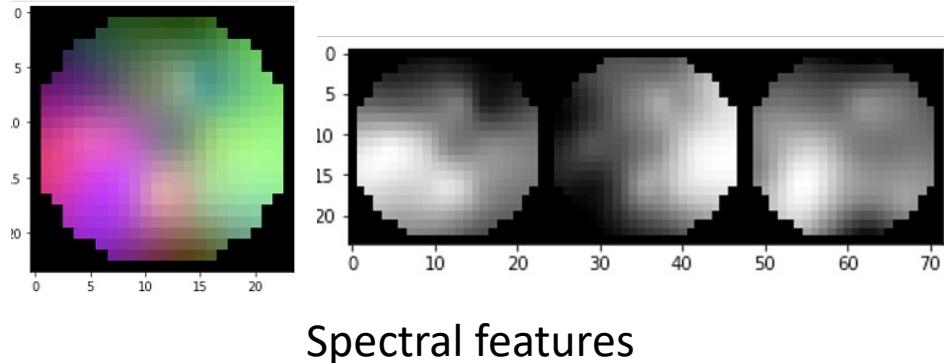


Figure 9. EEG processing choices. (a) Number of studies that used preprocessing steps, such as filtering, (b) number of studies that included, rejected or corrected artifacts in their data and (c) types of features that were used as input to the proposed models.



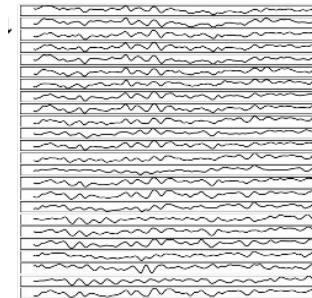
Distribution of input type according to the architecture category.

Our approach

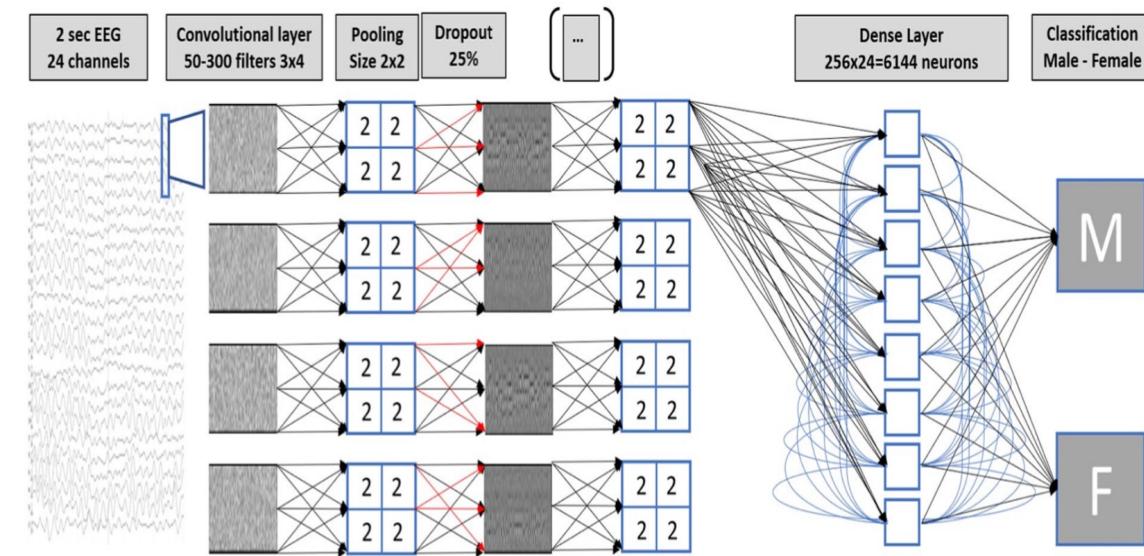


Spectral features

VS.

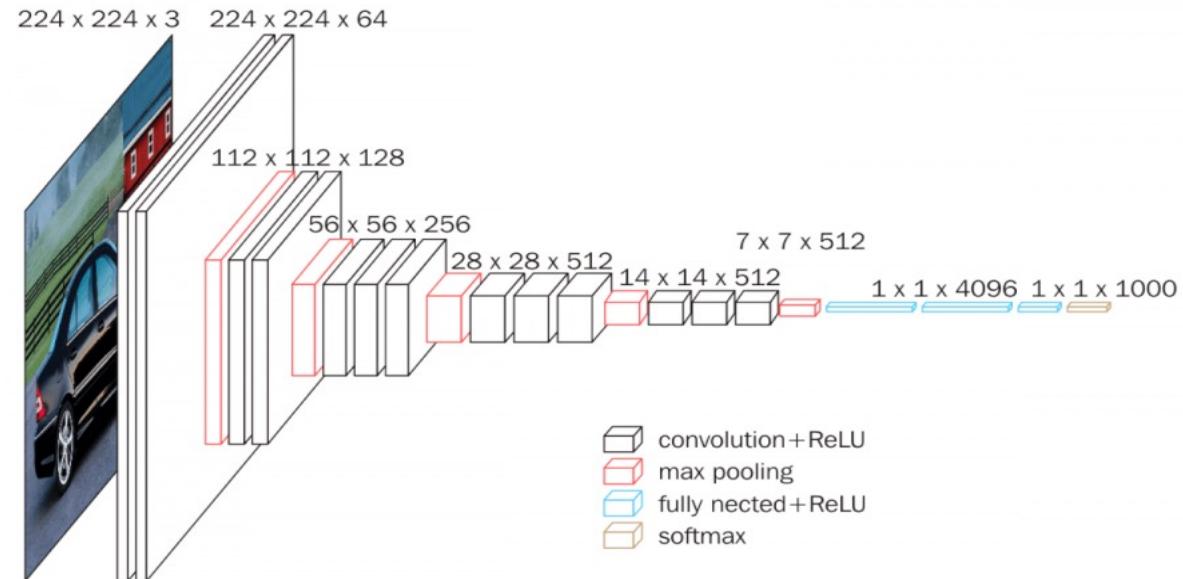


Raw EEG



van Putten et al. (2018)

VS.



Simonyan, K. and Zisserman A. (2014)

Data Data Data Data Data Data Data Data



Michael Milham, PI

- 128 channels EEG
- ~2,000 EEG datasets (planed 10,000)
- Tasks involving emotions (The Gift movie)
- Rest (eyes open and eyes closed)



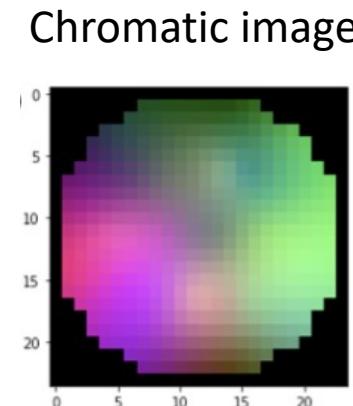
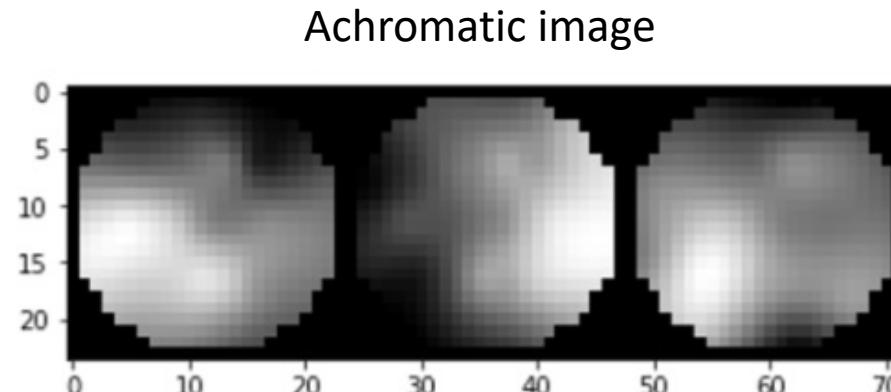
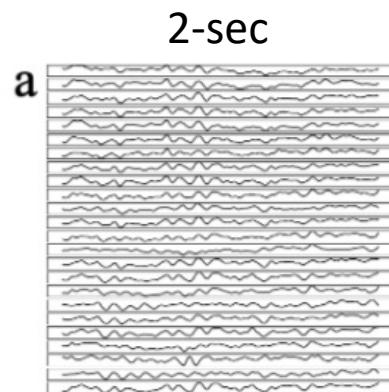
Data

- 2224 subjects:
 - 787 female (35%)
 - Ages: min 5, max 22, mean 10
- Pre-processing following the paper:
 - Remove baseline
 - Filter 0.25-25Hz
 - Resample 128Hz
 - Re-reference to average mastoids
 - Epoching: eye-closed, 3 40-second blocks. Ignored first and last 3 seconds of each block
 - clean_rawdata – ASR (our)
 - Sub-select 24 channels
 - Fp1, Fp2, F7, F3, Fz, F4, F8, FC3, FCz, FC4, T3, C3, C4, T4, CP3, CPz, CP4, T5, P3, Pz, P4, T6, O1, Cz
 - Segment 2-second non-overlapping windows
 - → ~ 81 samples per subject

Data

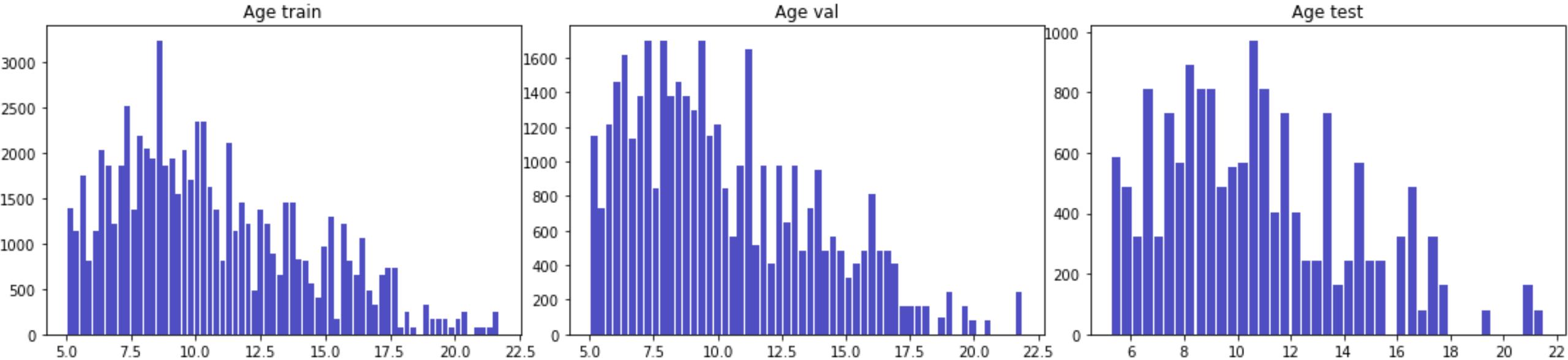
- Sub-select 1574 participants (50% female) 24-channels
- 2-second extracted epochs eyes open and eyes closed
- 5 predictors: **sex**, **handedness**, **eyes open/closed**, **age**, **segment count**

↑
Categorical
↑
Continuous

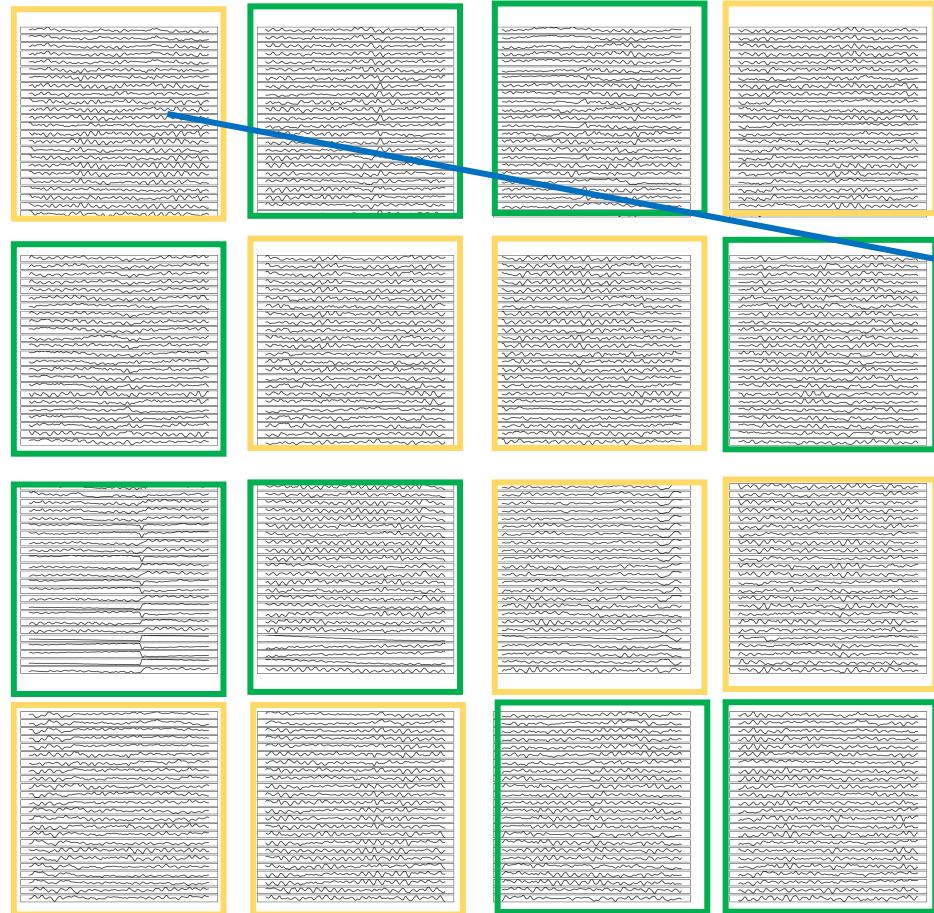


Data

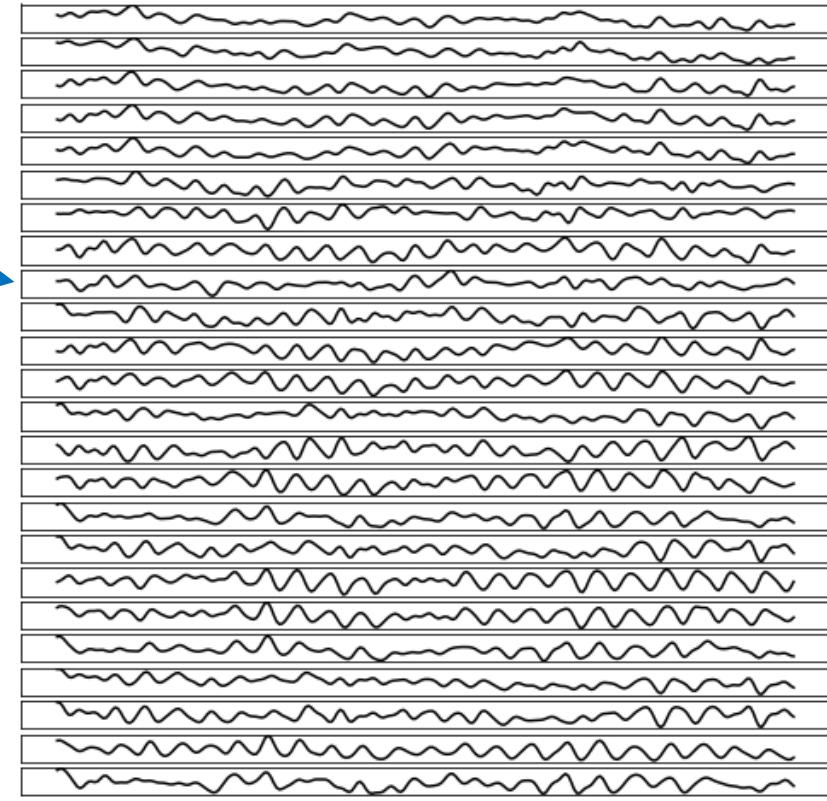
- 10-30-60 split
 - 885 subjects for training -> 71,381 samples
 - 492 subjects for validation -> 39,868 samples
 - 197 subjects for testing -> 15,925 samples



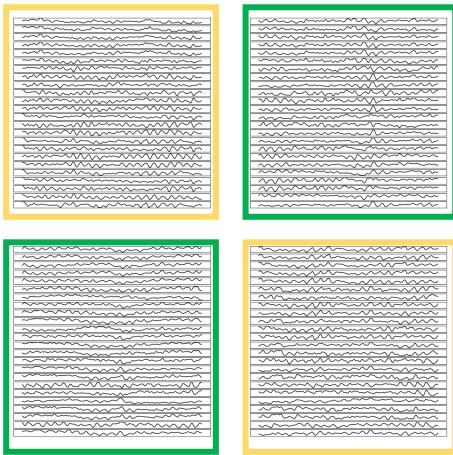
Input data $24 \times 256 \times n$



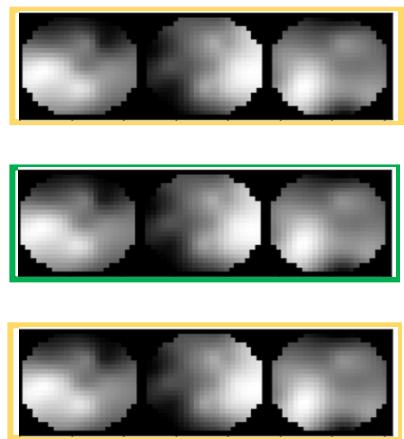
0 for male, 1 for female



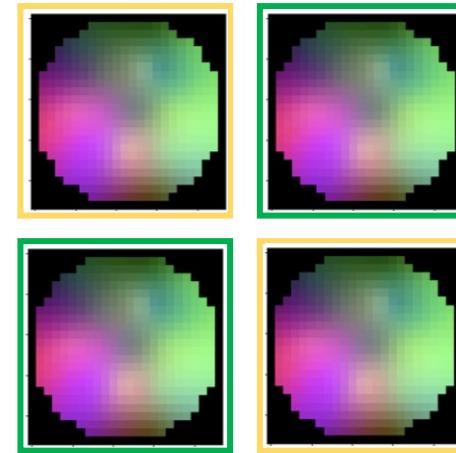
Raw input data
24 x 256 x n



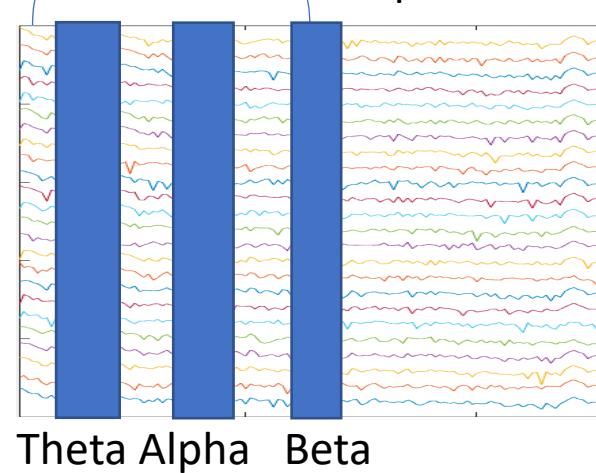
Spectral data
24 x 72 x n



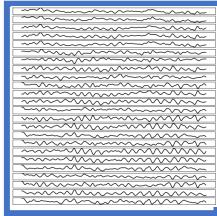
Spectral data
24 x 24 x n



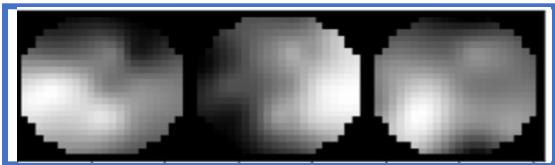
Tapered FFT
24 x 256 complex



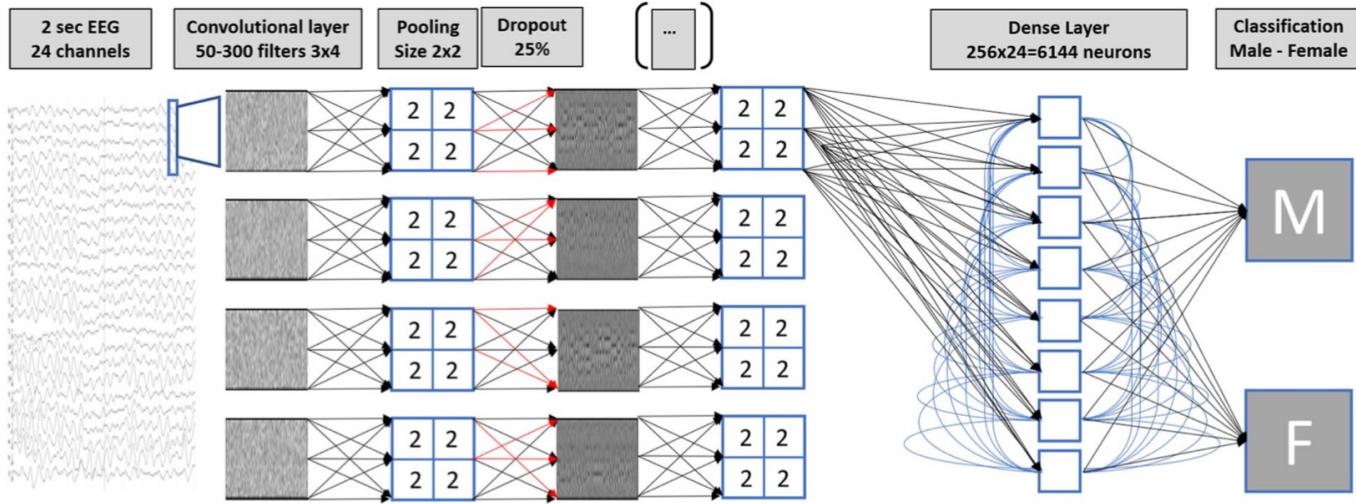
R-SCNN



S-SCNN



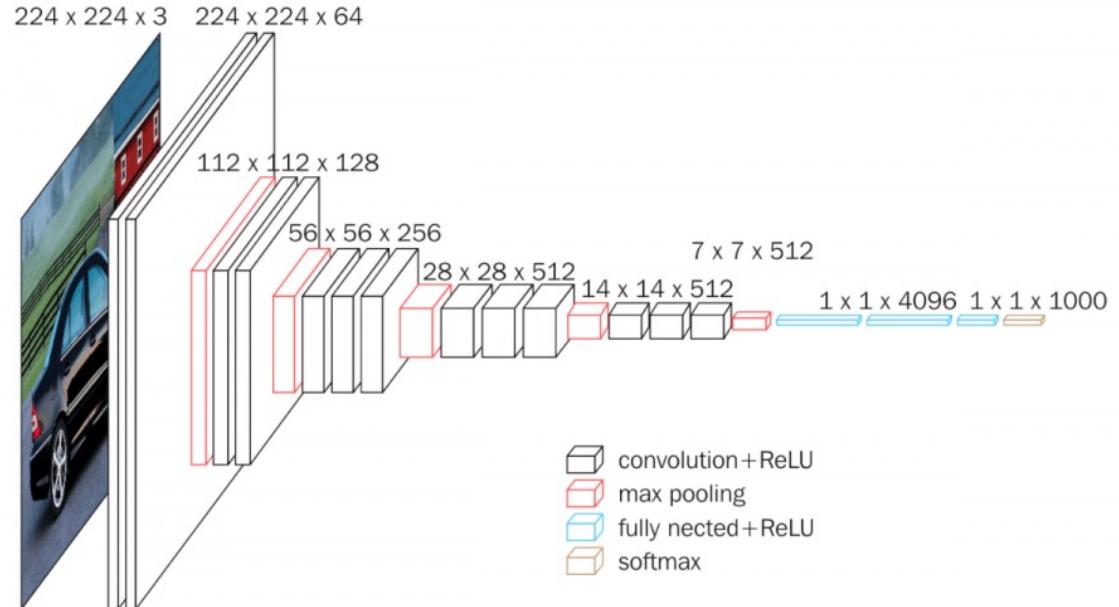
Or



Layer	Filter size	# of filters/hidden units
Convolutional	3x3	100
MaxPooling		
Dropout (25%)		
Convolutional	3x3	100
MaxPooling		
Dropout (25%)		
Convolutional	2x3	300
MaxPooling		
Dropout (25%)		
Convolutional†	1x7	300
MaxPooling*		
Dropout (25%)		
Convolutional†	1x3	100
Convolutional†	1x3	100
Fully connected		6144
Fully connected		2
Softmax		

Modified VGG16

- Changed input size
- Keep the same scaling between layers
- Number of convolutions divided by 8 for each layer
- Dropped layers 19 to 32
- Change number of output classes to 2
- Retrain the whole network



R-VGG

S-VGG

Or

Layer	Filter size	# of filters/hidden units
Convolutional†	3x3	16
Convolutional	3x3	16
MaxPooling		
Convolutional	3x3	32
Convolutional	3x3	32
MaxPooling		
Convolutional	3x3	64
Convolutional	3x3	64
Convolutional	3x3	64
MaxPooling		
Fully connected		1024
Dropout (50%)		
Fully connected		1024
Dropout (50%)		
Fully connected		2
Softmax		

Training the network

- Computing environment



Training the network

- Programming environment



v3.7.10



v1.3.1

A screenshot of a Jupyter Notebook interface. The top bar shows "jupyter SexPrediction-Final Last Checkpoint: Last Monday at 5:40 PM (unsaved changes)" and "Trusted | Python (ML) O Logout". The notebook has two cells. Cell In [41] contains a Python function "train" for training a PyTorch model. Cell In [82] contains code to define a neural network model and train it using Adamax optimizer. The output of cell In [82] shows a message indicating successful key matching.

```
In [41]: def train(model, optimizer, epochs=1):
    """
    Inputs:
    - model: A PyTorch Module giving the model to train.
    - optimizer: An Optimizer object we will use to train the model
    - epochs: (Optional) A Python integer giving the number of epochs to train for

    Returns: Nothing, but prints model accuracies during training.
    """
    model = model.to(device=device) # move the model parameters to CPU/GPU
    for e in range(epochs):
        for t, (x, y) in enumerate(loader_train):
            model.train() # put model to training mode
            x = x.to(device=device, dtype=dtype) # move to device, e.g. GPU
            y = y.to(device=device, dtype=torch.long)

            scores = model(x)
            loss = F.cross_entropy(scores, y)

            # Zero out all of the gradients for the variables which the optimizer
            # will update.
            optimizer.zero_grad()

            # This is the backwards pass: compute the gradient of the loss with
            # respect to each parameter of the model.
            loss.backward()

            # Actually update the parameters of the model using the gradients
            # computed by the backwards pass.
            optimizer.step()

            if t % print_every == 0:
                print('Epoch %d, Iteration %d, loss = %.4f' % (e, t, loss.item()))
                check_accuracy(loader_val, model)
                print()

In [82]: lr = 0.002
batch_size = 70
loader_train = DataLoader(train_data, batch_size=batch_size, shuffle=True)
loader_val = DataLoader(val_data, batch_size=batch_size)
model = nn.Sequential(
    nn.Conv2d(1,100,3),
    nn.ReLU(),
    nn.MaxPool2d(2, 2),
    nn.Dropout(0.25),
    nn.Conv2d(100,100,3),
    nn.ReLU(),
    nn.MaxPool2d(2, 2),
    nn.Dropout(0.25),
    nn.Conv2d(100,300,(2,3)),
    nn.ReLU(),
    nn.MaxPool2d(2, 2),
    nn.Dropout(0.25),
    nn.Conv2d(300,300,(1,7)),
    nn.ReLU(),
    nn.MaxPool2d(1,2), stride=1,
    nn.Dropout(0.25),
    nn.Conv2d(300,100,(1,3)),
    nn.Conv2d(100,100,(1,3)),
    nn.Flatten(),
    nn.Linear(1900,6144),
    nn.Linear(6144,2),
)
optimizer = torch.optim.Adamax(model.parameters(), lr=lr)
train(model, optimizer, epochs=70)
# model.load_state_dict(torch.load('logs/model_saved'))
```

Out[82]: <All keys matched successfully>

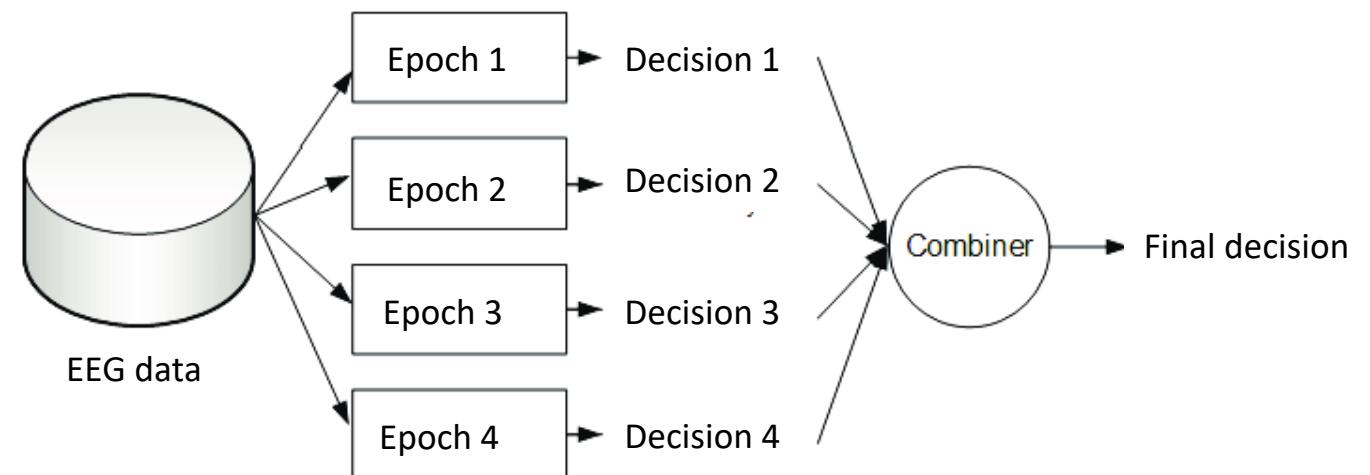
Training the network

- Adamax optimizer
 - $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 108$ and decay = 0.0
- Learning rate = 0.002
- Batch size = 70 with random shuffle
- 70 epochs

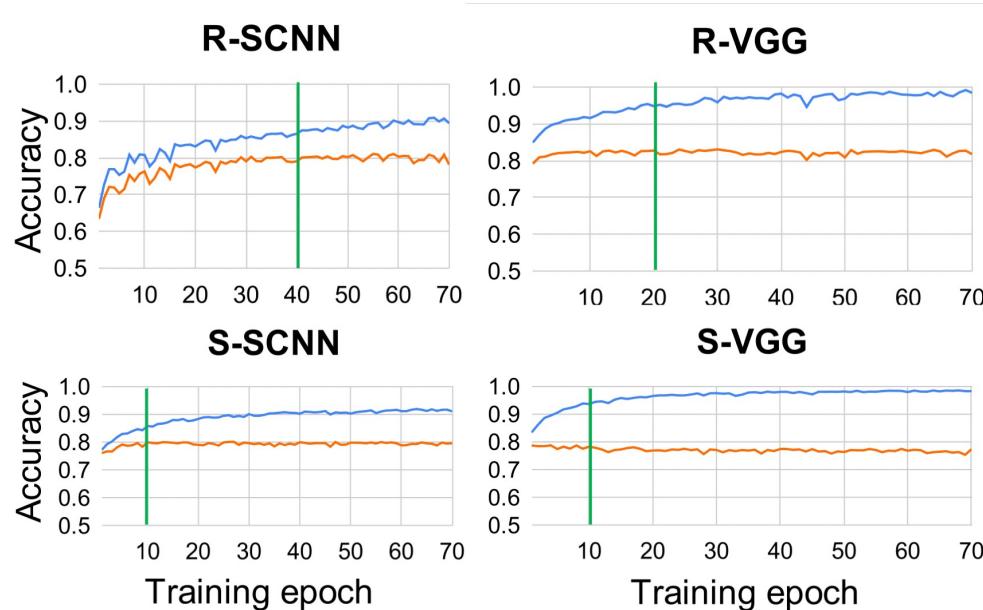
Testing scheme: Majority voting



- Estimates are on 2-second epochs
- We have about 81 samples per subject
- Boost performance by about 4% (see next slide)



Results



Model	2-sec	2-sec x n with vote
	Per-sample	Per-subject
R-SCNN	80.6 (79.7 to 81.5)	85.1 (84.3 to 85.9)
R-VGG	83.1 (82.7 to 83.4)	87.0 (86.6 to 87.4)
S-SCNN	79.0 (78.7 to 79.3)	83.2 (82.1 to 84.3)
S-VGG	77.1 (76.8 to 77.4)	81.3 (80.0 to 82.6)

Table 3. Models classification accuracy. 95% confidence interval is indicated in parenthesis. Bolded values indicate best performance.

Assessing features learned by the network



Krizhevsky et al., 2012

3/2/add_5

Type: Add
Channels: 1,280
Convolution: [1,1]

Technique

- Feature Visualization
- DeepDream
- Dataset Samples
- Caricature
- Text Feature Visualization

An artificial, optimized image that maximizes activations of the given unit. [Read more.](#)

Params

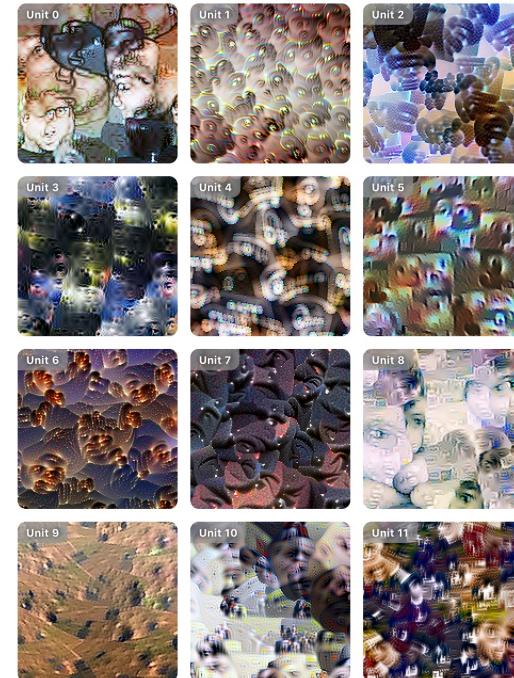
Optimization Objective
 channel
 neuron

View

Image Size



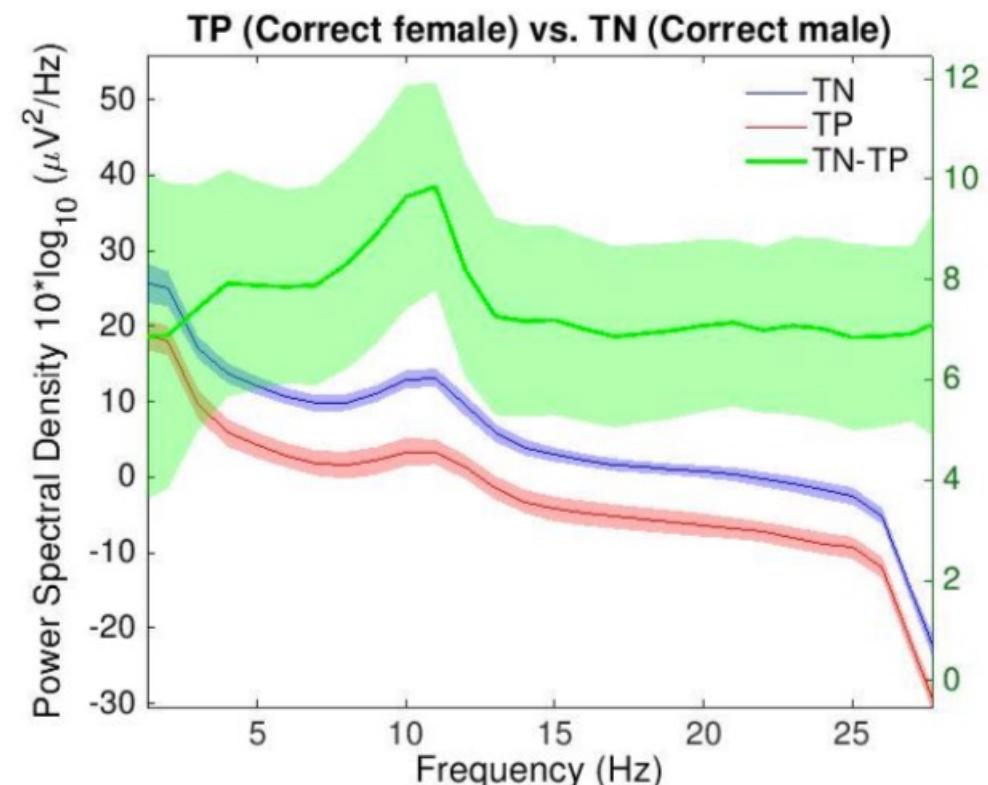
Resize Behavior
 Crop image
 Scale image



Best samples

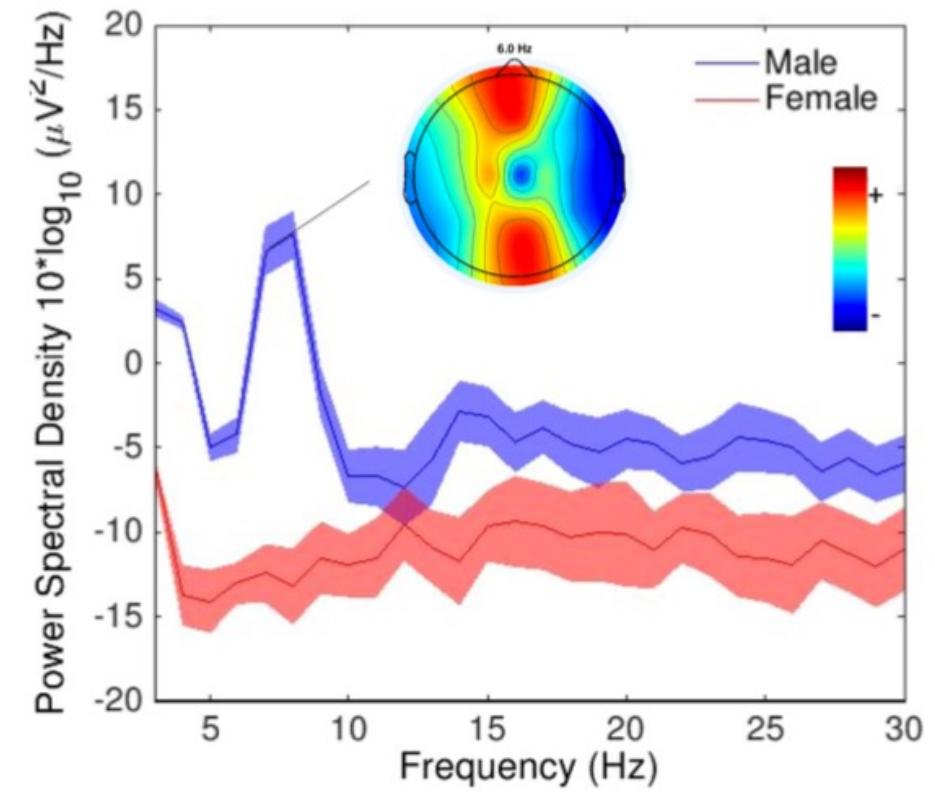
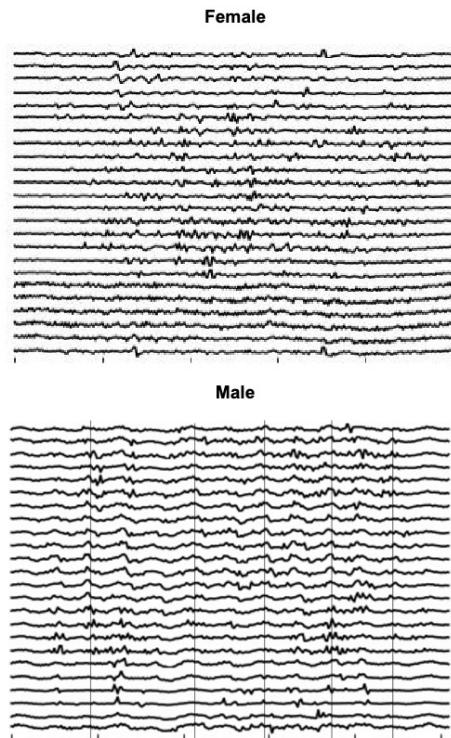
- “Best” = gives highest activation in classification neurons
- Get samples of the top 20 subjects in validation set
- Best male samples show higher spectral power across frequencies, most notably the alpha band near 10Hz

	Classified as female	Classified as male
Female sample	True Positive (TP)	False Negative (FN)
Male sample	False Positive (FP)	True Negative (TN)



Activation Maximization

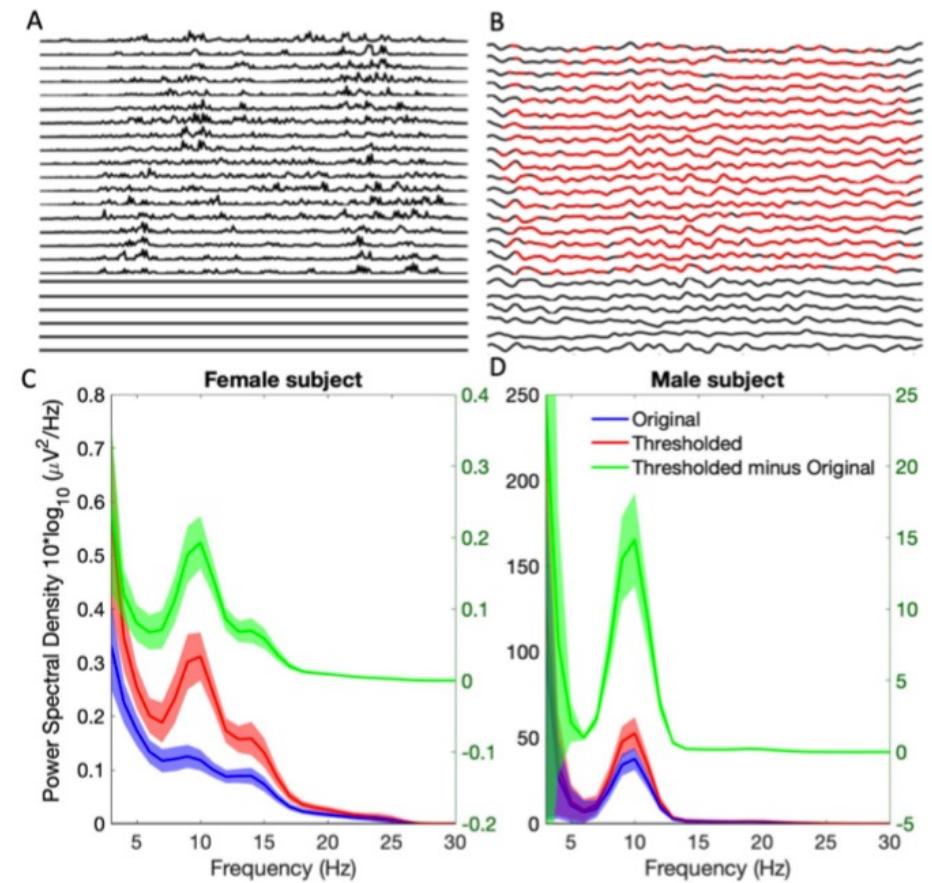
- Synthesize the input that maximize the activation of the classification neurons
- Get 20 samples for each sex
- Best male samples show distinctly higher high theta power for male samples (6-8 Hz)



Truong, D., Makeig, S., & Delorme, A. (2021). Assessing learned features of Deep Learning applied to EEG. *arXiv preprint arXiv:2111.04309*.

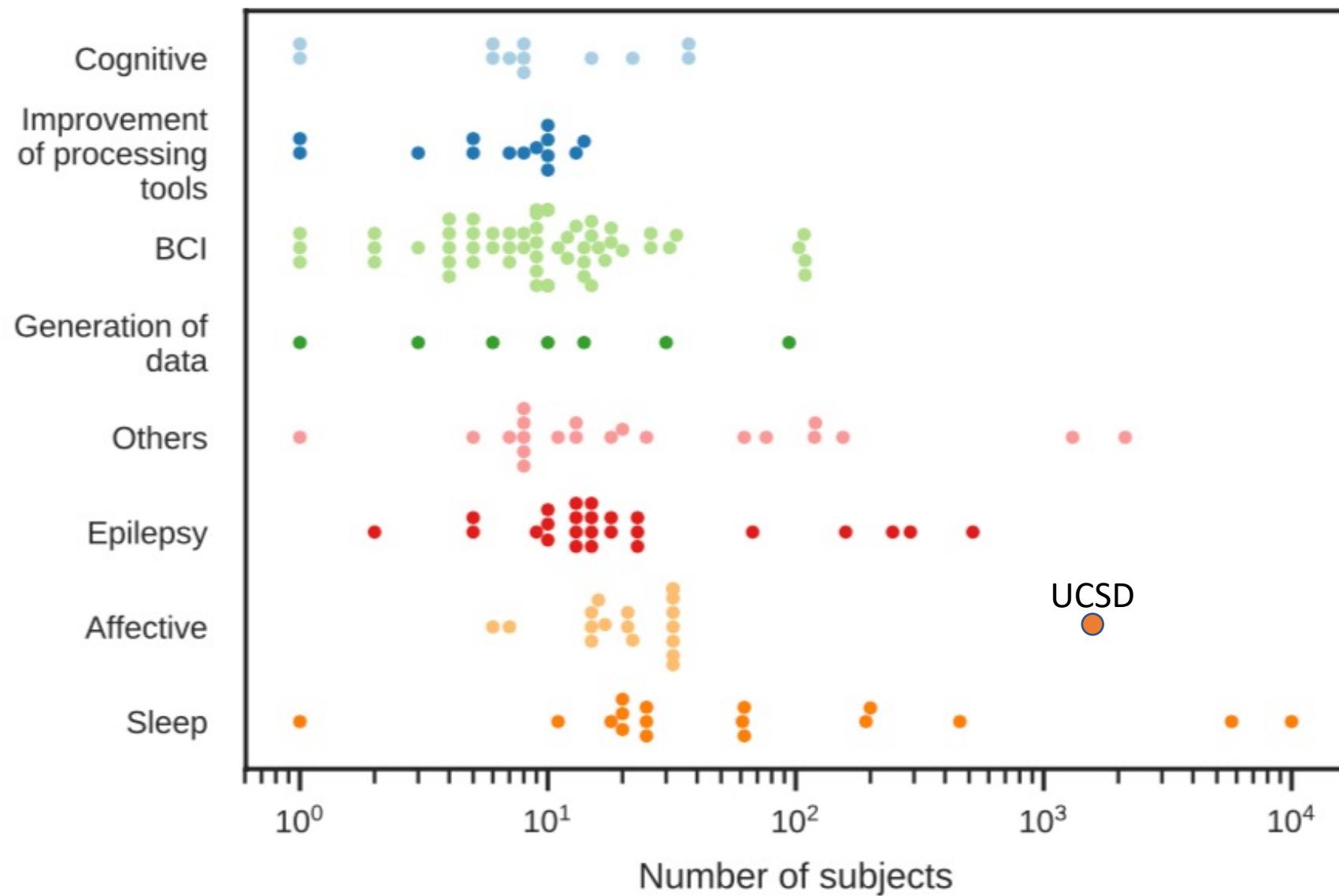
Saliency map

- Back-project the gradient of the classification neuron to the input
- Magnitude of the gradients for each input value indicate the importance of that value to the neuron's activation
- Raw EEG samples contributing the least to the classification (gradients fell below the 30% quantile threshold) were removed then linearly interpolated using the remaining samples
- Thresholded samples showed higher alpha power, most notably near 10Hz



Discussion

- Most EEG datasets are order of magnitude smaller



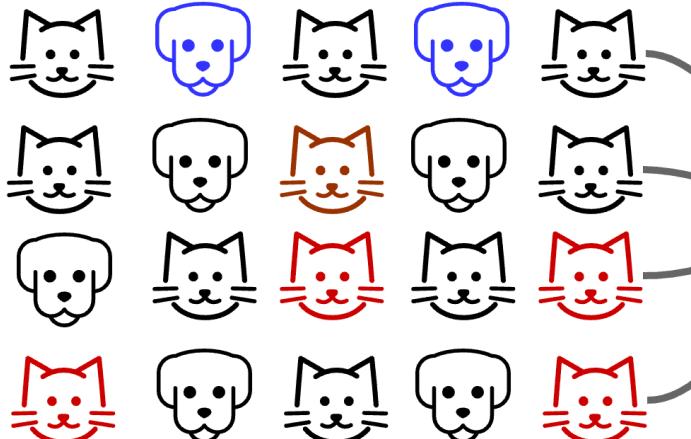
Discussion

- Most EEG datasets are order of magnitude smaller
- Not generalizable to different number of channels and sampling rate
- Problem chosen could be simple: boys move more than girls, known sex difference in skull thickness
- VGG-16 is no longer state-of-the-art model nor designed specifically for time-series/EEG data
- Scientific values of learned features visualization remain unclear

Future work: representation learning

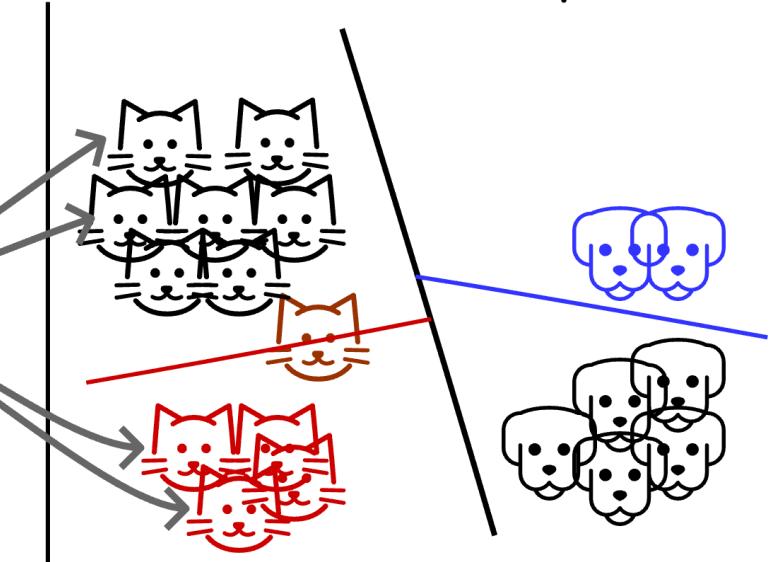
- DL models take data from the original space and map it to a “meaningful” representation

Default Representation



Deep Neural Network

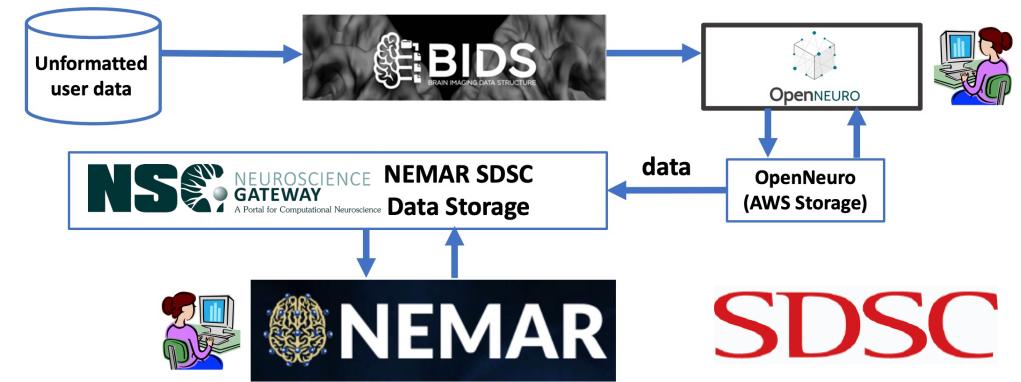
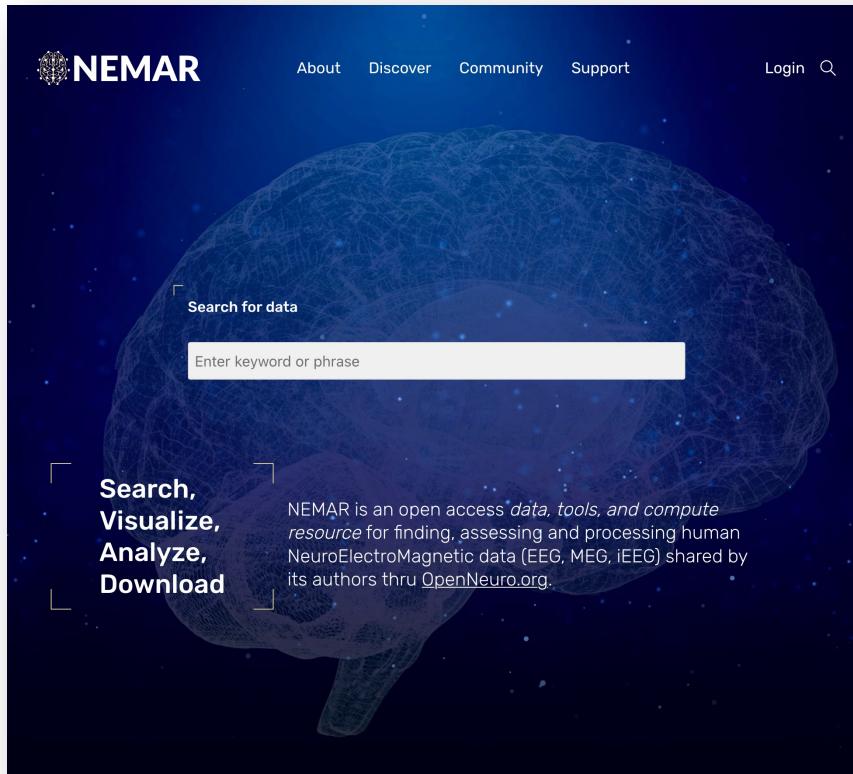
"Good" Semantic Representation



Future work: representation learning

- DL models take data from the original space and map it to a “meaningful” representation
- Can we learn embeddings that are useful across tasks?
- Can we leverage embeddings across modalities to study naturalistic cognition?
- What network architectures and training regimes can generalize across datasets?

Data, tools, and compute resources (DATCOR)



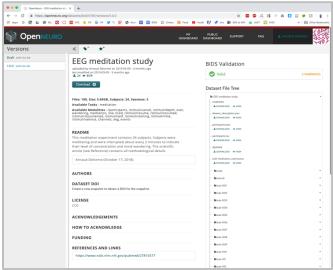
Total dataset number	128 (5TB)
Number of EEG dataset	95
Number of MEG dataset	20
Number of iEEG dataset	13
Hours of recording	3,000+
Number of participants	6,827

Task names from NEMAR datasets

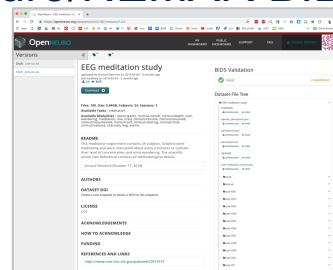


Standardized event description enables analysis across datasets

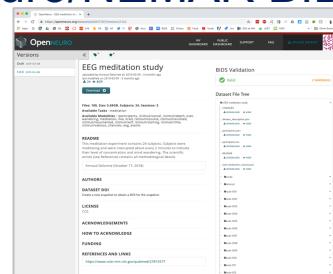
OpenNeuro/NEMAR BIDS #1



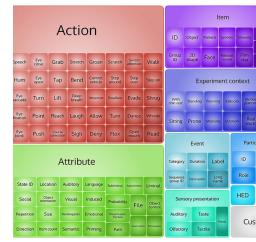
OpenNeuro/NEMAR BIDS #2



OpenNeuro/NEMAR BIDS #3



HED event selection

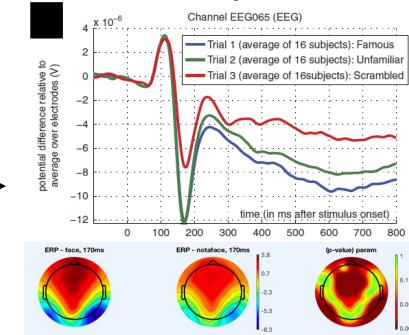


<https://hedtags.org>

EEGLAB automated pipeline



Output



Future development: *BIDS-DL* plug-in

- Automated pipeline to convert HED-identified data segments from BIDS datasets to DL-ready dataset
- Host data on the cloud and make it streamable so that no data download/upload required

Truong, D., Sinha, M., Venkataraju, K. U., Milham, M., & Delorme, A. (2022) A streamable large-scale clinical EEG dataset for Deep Learning. The 44th International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC); July 11-15, 2022.

Thank you