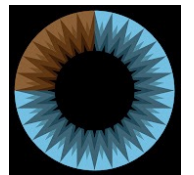# Deep Learning applied to EEG
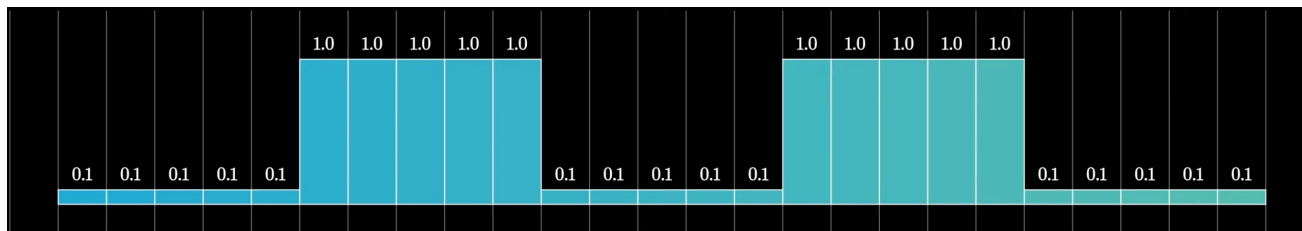
Dung "Young" Truong & Arnaud Delorme
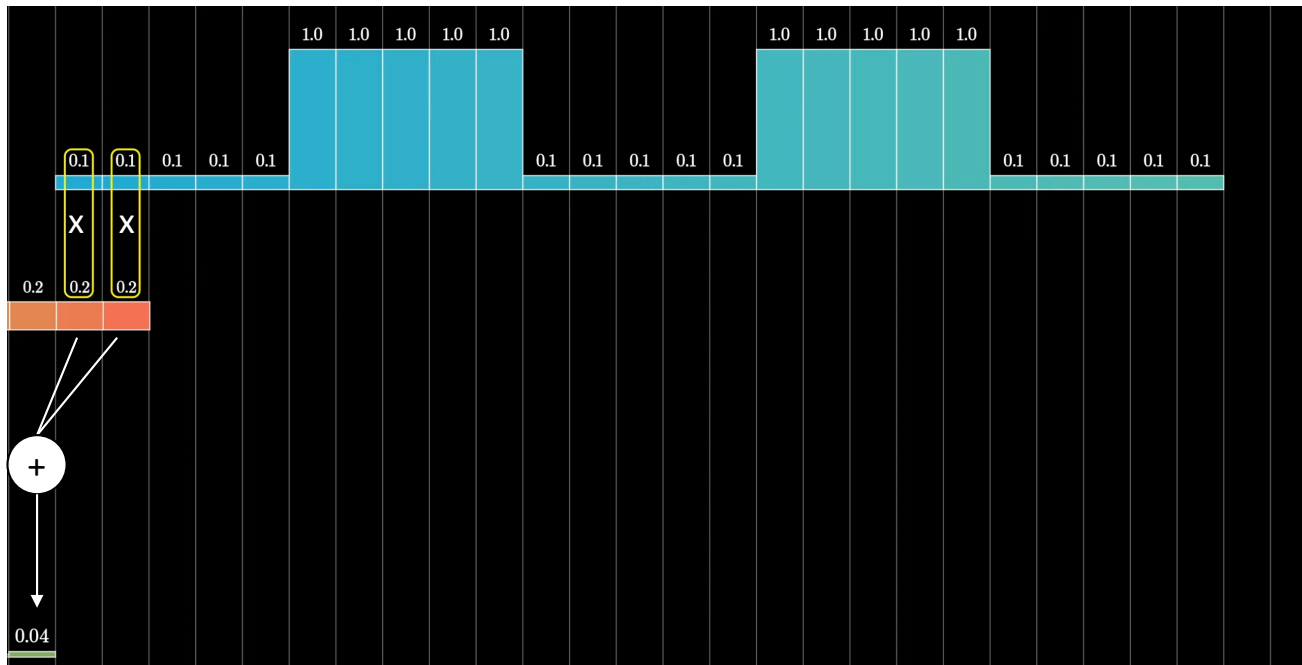
# What kind of operation?

# 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?
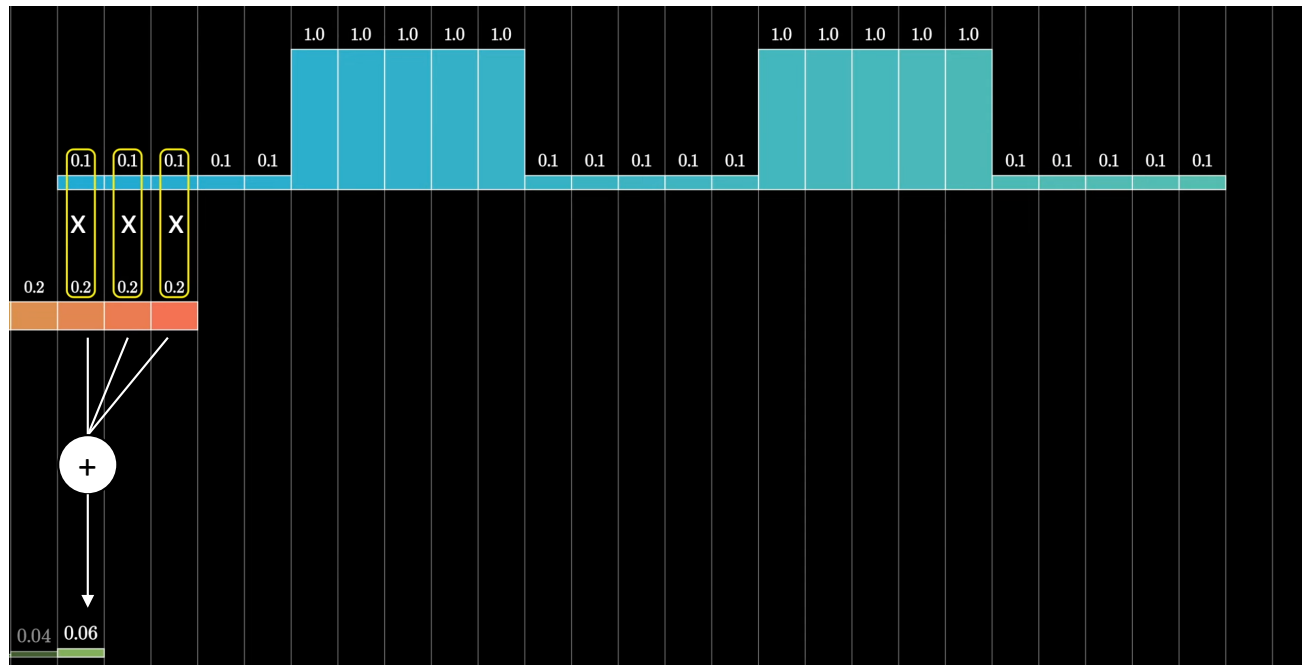
# What kind of operation?

# What kind of operation?

# What kind of operation?

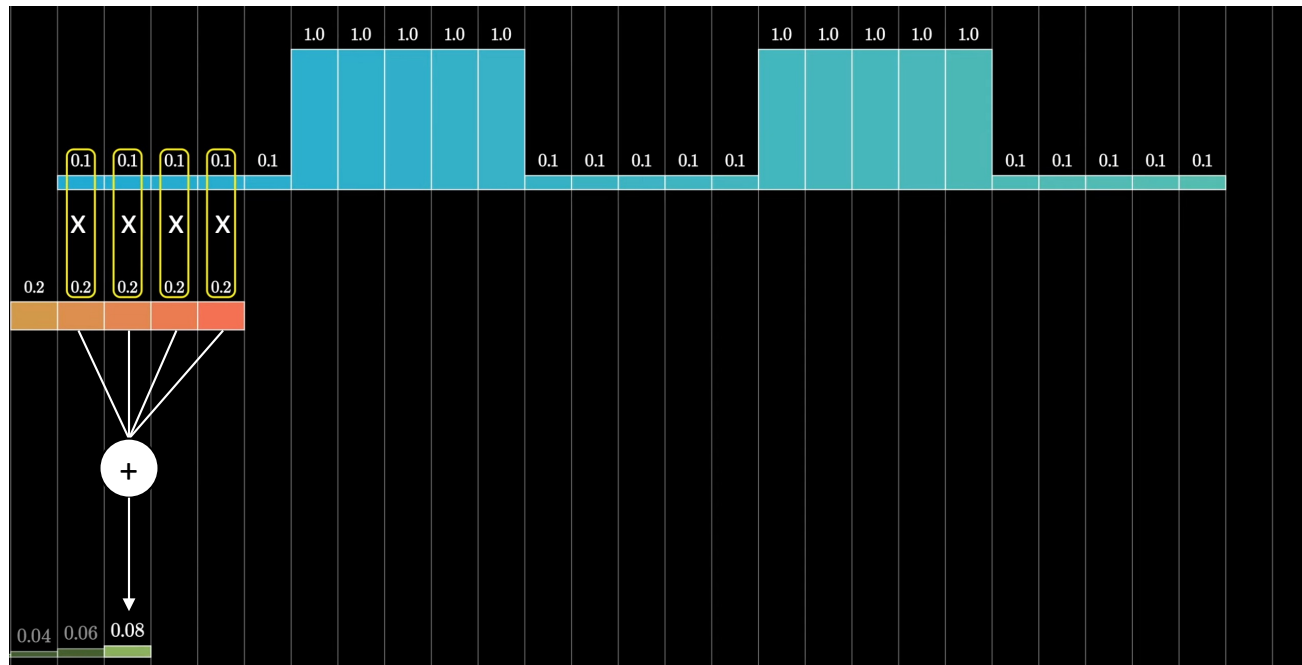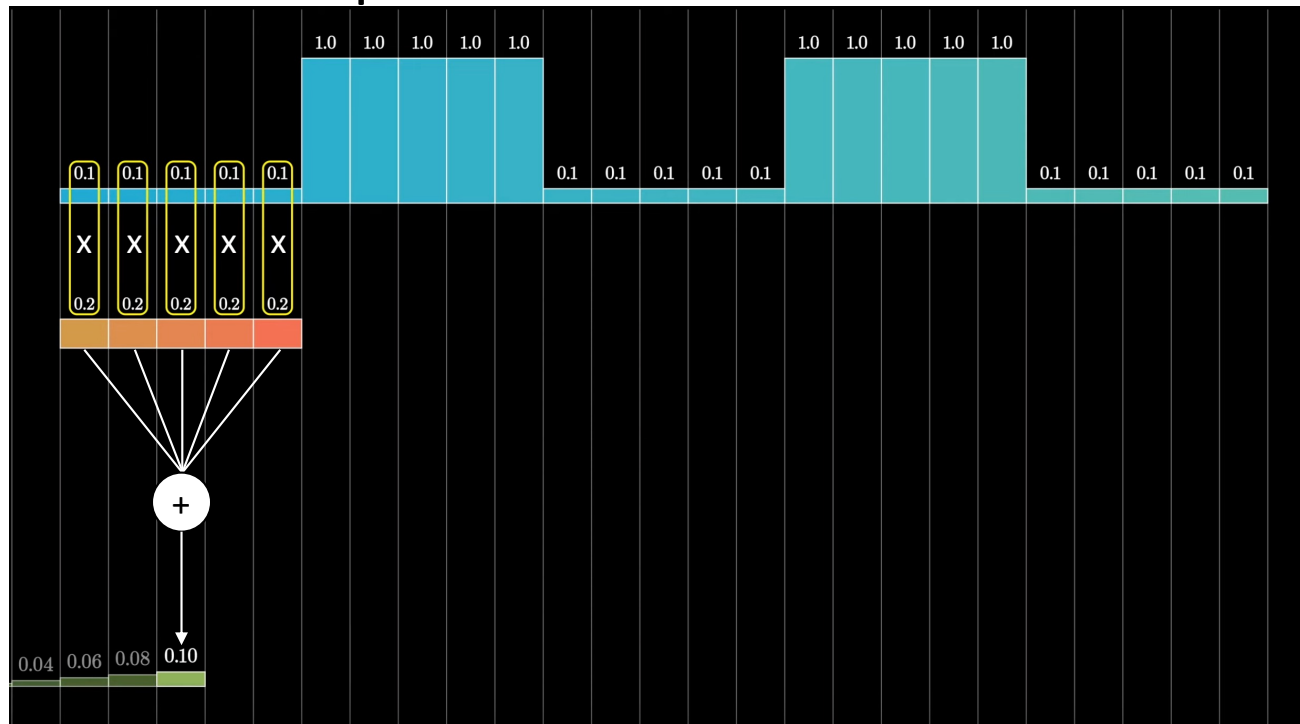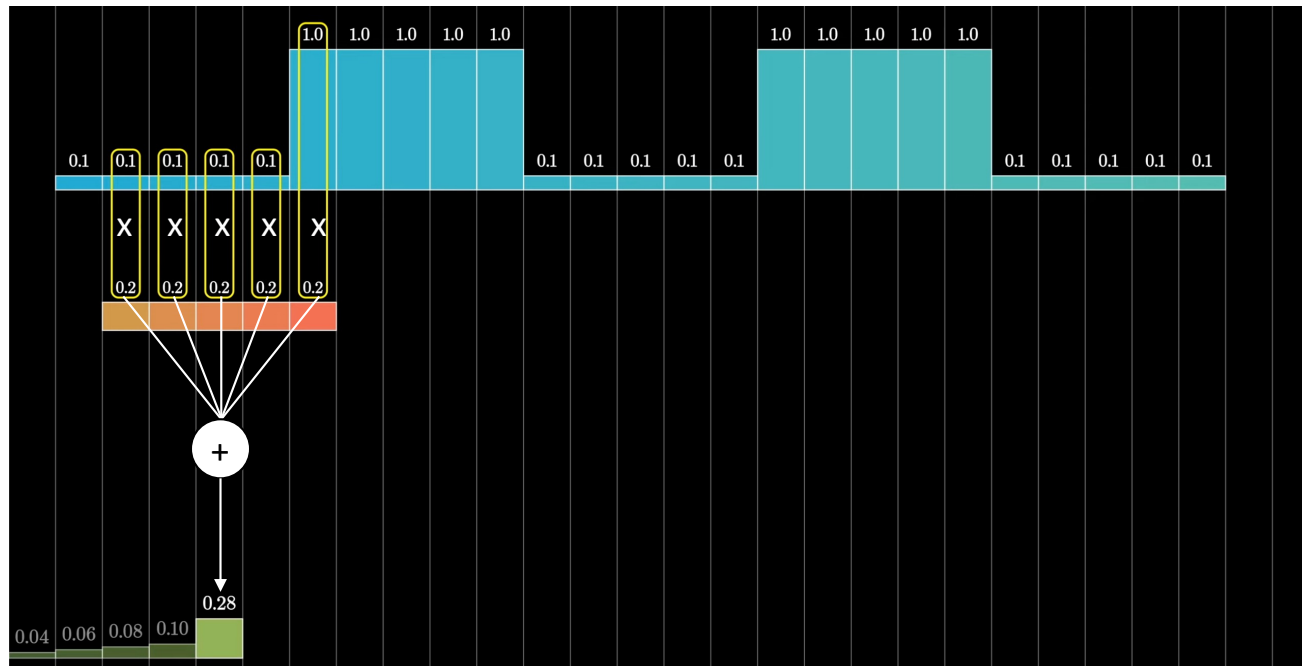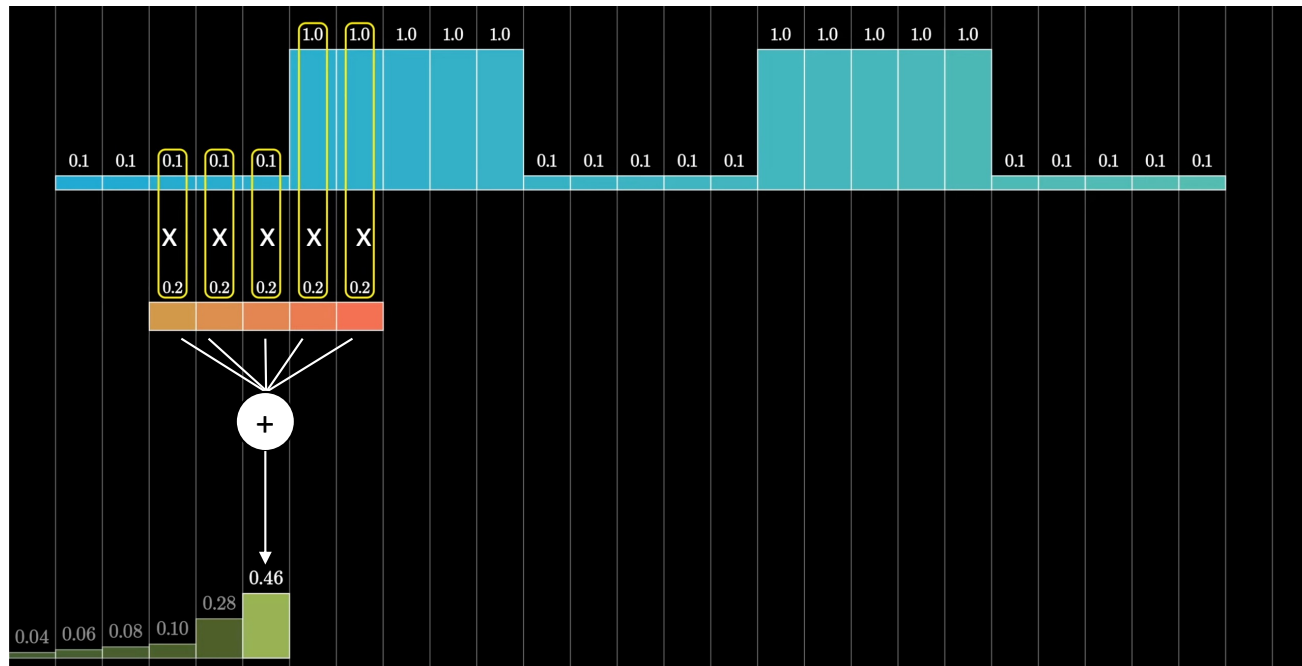# What kind of operation?
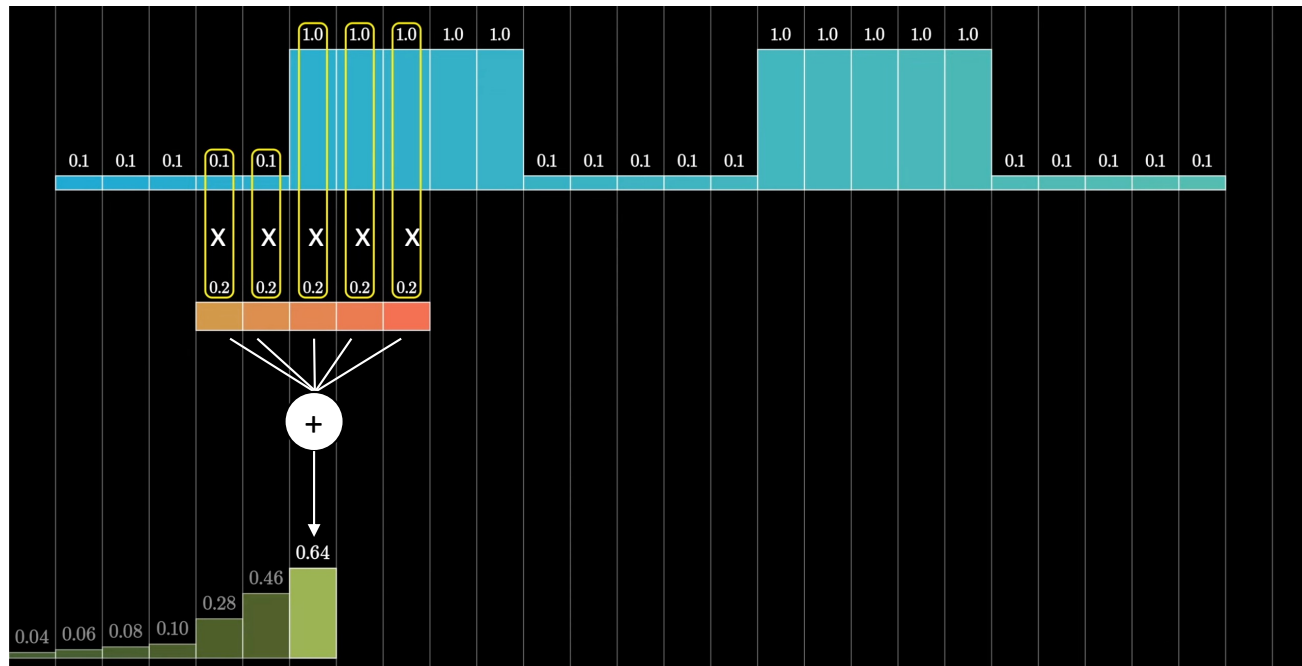
# What kind of operation?

# What kind of operation?

# What kind of operation?
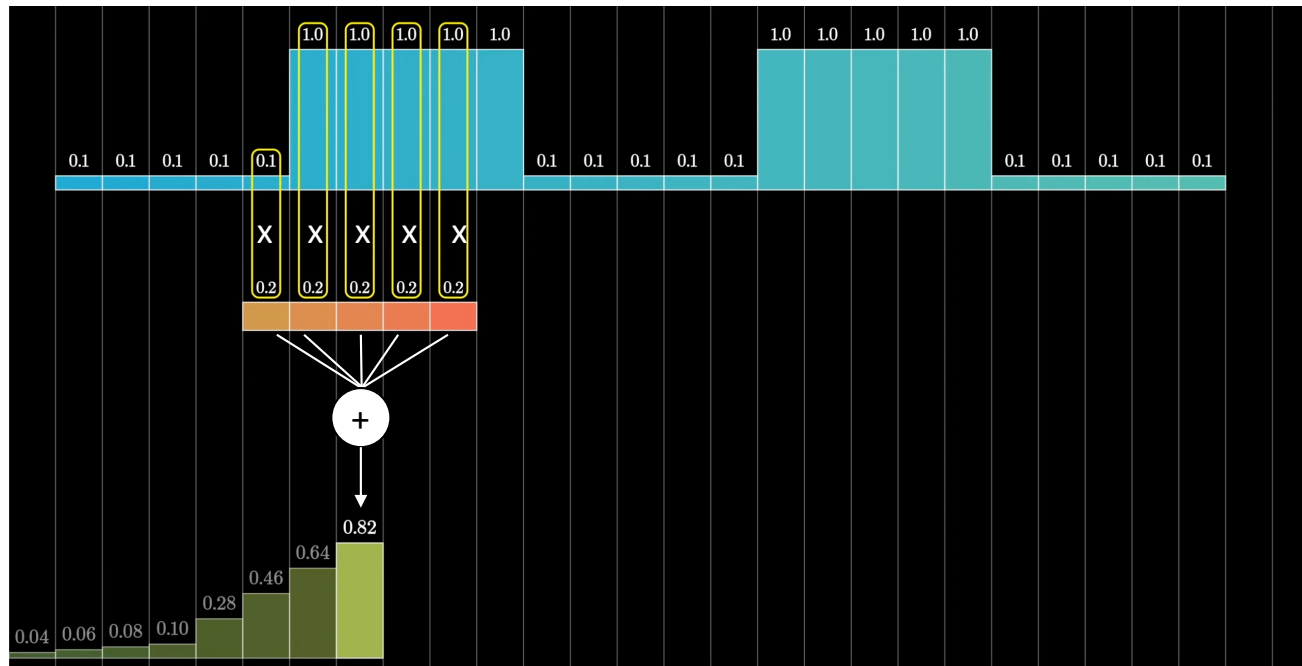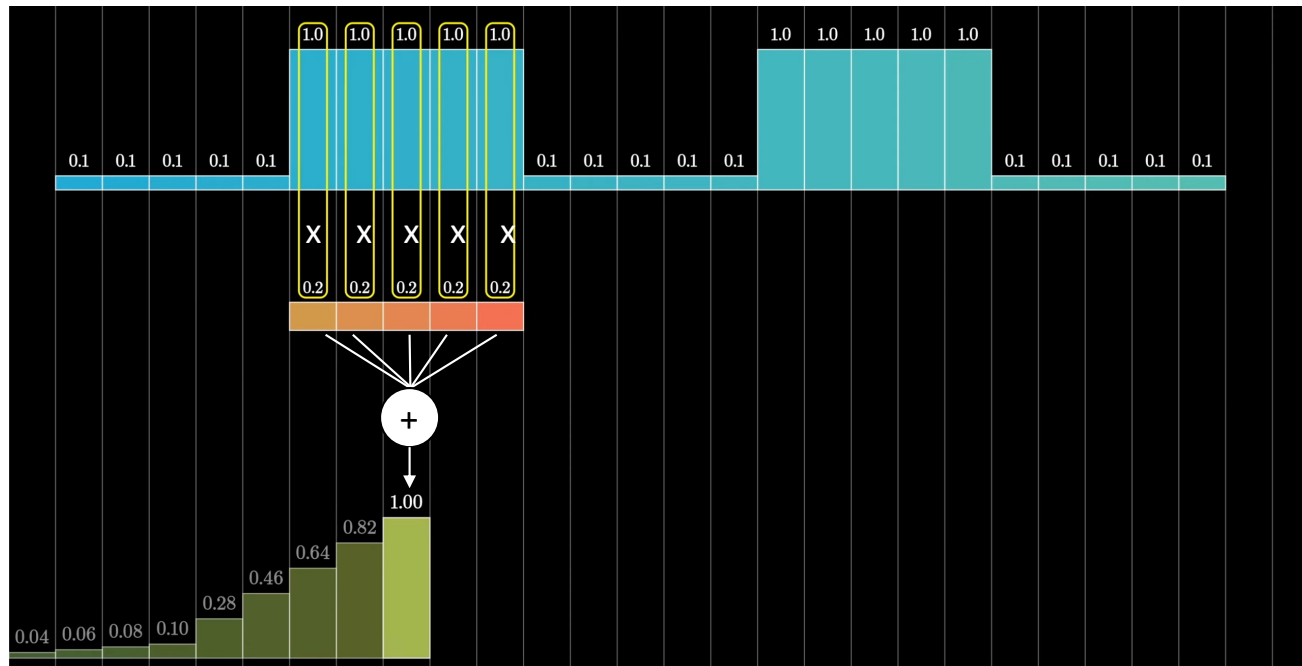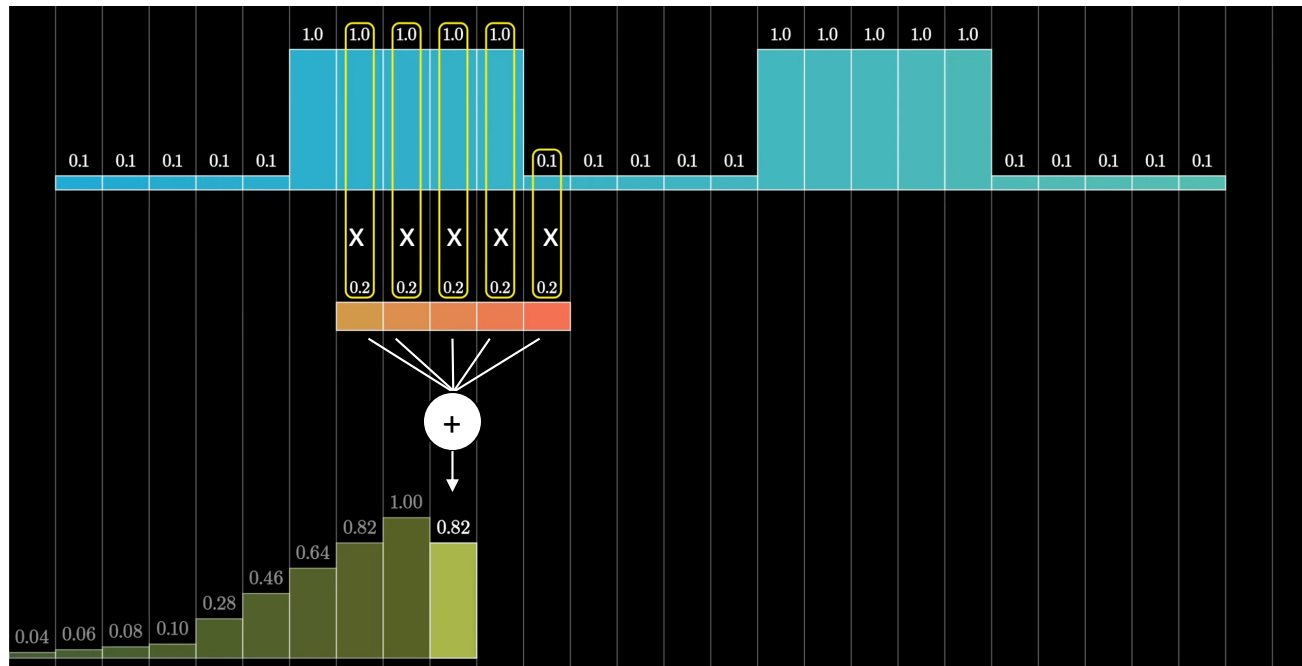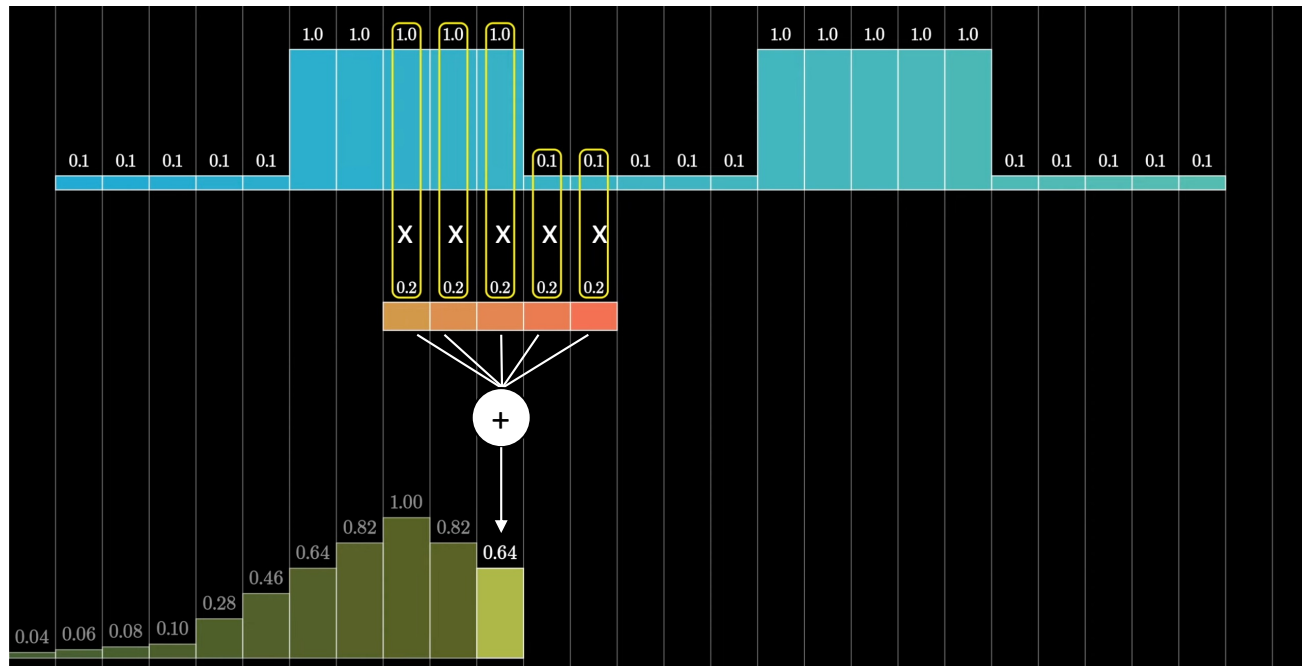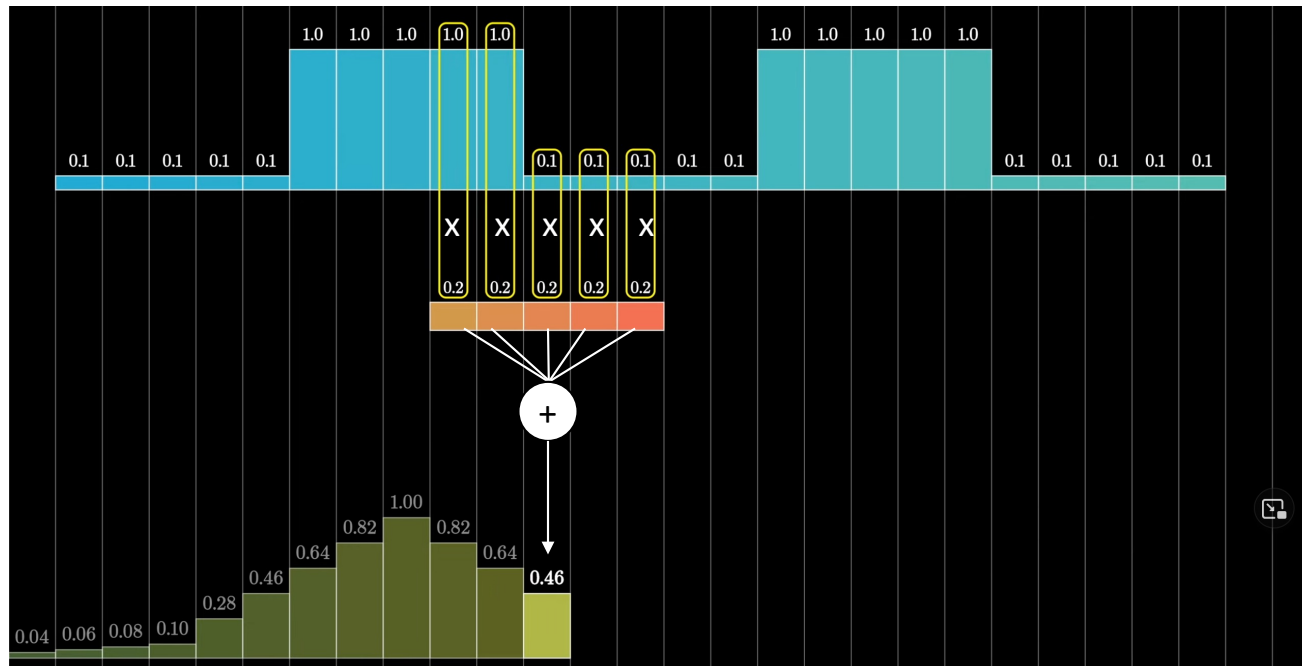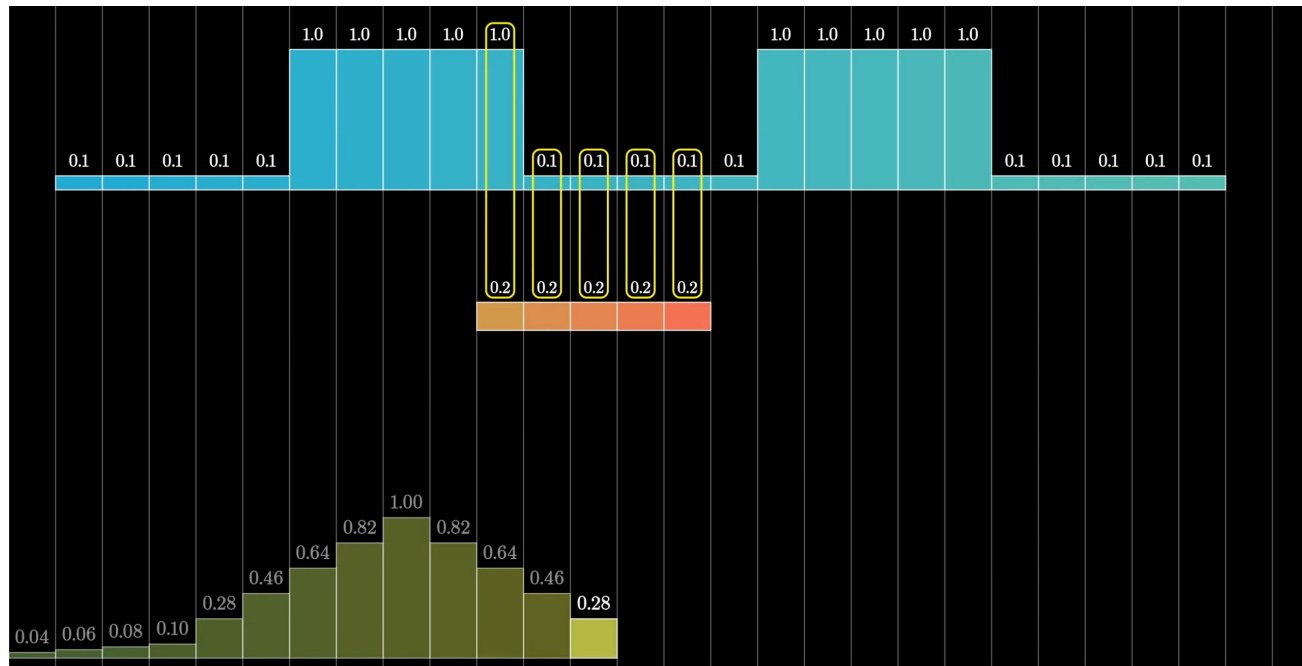
# What kind of operation?

# Moving average – 1D Convolution



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

# 2D Convolution

Input



Kernel/Filter

Feature/Activation map

7x1+4x1+3x1+
2x0+5x0+3x0+
3x-1+3x-1+2x-1
= 6

# Edge detection filter

| 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 |

\*

| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

=

| 0 | 3 | 3 | 0 |
|---|---|---|---|
| 0 | 3 | 3 | 0 |
| 0 | 3 | 3 | 0 |
| 0 | 3 | 3 | 0 |

# There are many types of filters



$$\circledast \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\circledast \begin{bmatrix} 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \\ 1/9 & 1/9 & 1/9 \end{bmatrix}$$

$$\circledast \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\circledast \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

Identity          Blur          Edge detection          Sharpen

https://en.wikipedia.org/wiki/Kernel_(image_processing)

# Filter for digit detection?



"4"

?

# Filter for digit detection?



$\ast$

?

... $\longrightarrow$ "4"

# Let machine learn

# Image processing + Data + Neural network



label = 5 label = 0 label = 4 label = 1 label = 9
label = 2 label = 1 label = 3 label = 1 label = 4
label = 3 label = 5 label = 3 label = 6 label = 1

MNIST dataset

LeNet-5
(Convolutional
Neural Network)

Digit

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



Input image → **Convolution** → Convolutional layer → **Pooling** → Pooling layer → **Flattening** → Dense layers → Output layer ($Y_1$, $Y_2$, $Y_3$, $Y_4$)

# Convolutional layer

- Convolution operation

Input



Kernel/Filter

Feature/Activation map

7x1+4x1+3x1+
2x0+5x0+3x0+
3x-1+3x-1+2x-1
= 6

# Convolutional layer

- Padding



Kernel

# Convolutional layer

- Stride

Example: Padding = 1, Stride = 2

# Convolutional layer

- Convolution over volume



$6 \times 6 \times \boxed{3}$   $*$   $3 \times 3 \times \boxed{3}$   $=$   $4 \times 4$

# Convolutional layer

- Stacking activation maps



$6 \times 6 \times 3$

Vertical edge

$3 \times 3 \times 3$

Horizontal edge

$3 \times 3 \times 3$

$4 \times 4$

$4 \times 4$

$4 \times 4 \times 2$

$4 \times 4 \times 2$

Andrew Ng

# Activation functions

- Rectified Linear Unit (ReLU)

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



**ReLU Layer**

**Filter 1 Feature Map**

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

→

| 9 | 3 | 5 | 0 |
|---|---|---|---|
| 0 | 2 | 0 | 1 |
| 1 | 3 | 4 | 1 |
| 3 | 0 | 5 | 1 |

# Convolutional Neural Network Architecture

# Pooling layer

# Convolutional Neural Network Architecture



Input image → **Convolution** → Convolutional layer → **Pooling** → Pooling layer → **Flattening** → Dense layers → Output layer ($Y_1$, $Y_2$, $Y_3$, $Y_4$)

# Fully connected layer



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

# Fully connected layer

- Flattening



Flatten

# Convolutional Neural Network Architecture



Input image → Convolution → Convolutional layer → Pooling → Pooling layer → Flattening → Dense layers → Output layer

$Y_1$
$Y_2$
$Y_3$
$Y_4$

# Activation function for classification

- Softmax

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

# 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

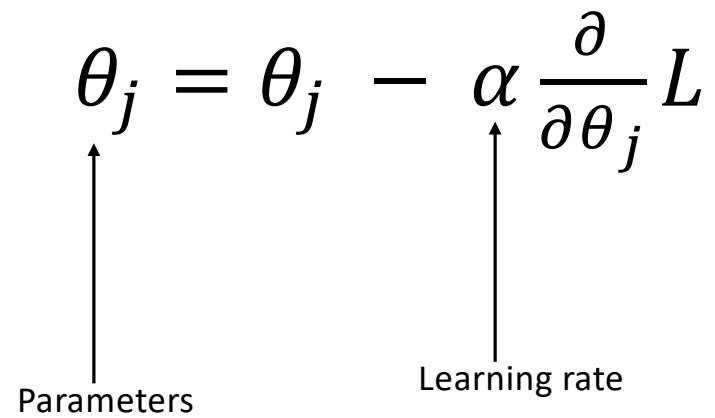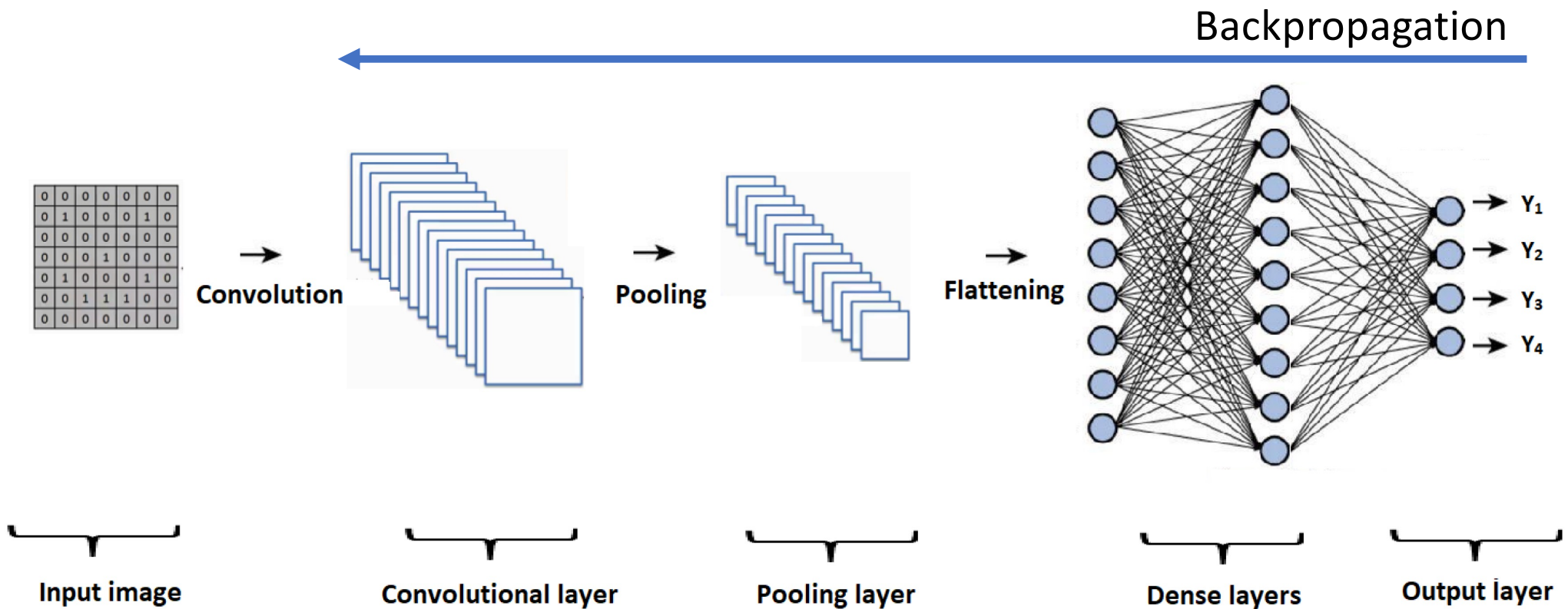$$\theta_j = \theta_j \; - \; \alpha \, \frac{\partial}{\partial \theta_j} L$$

Parameters
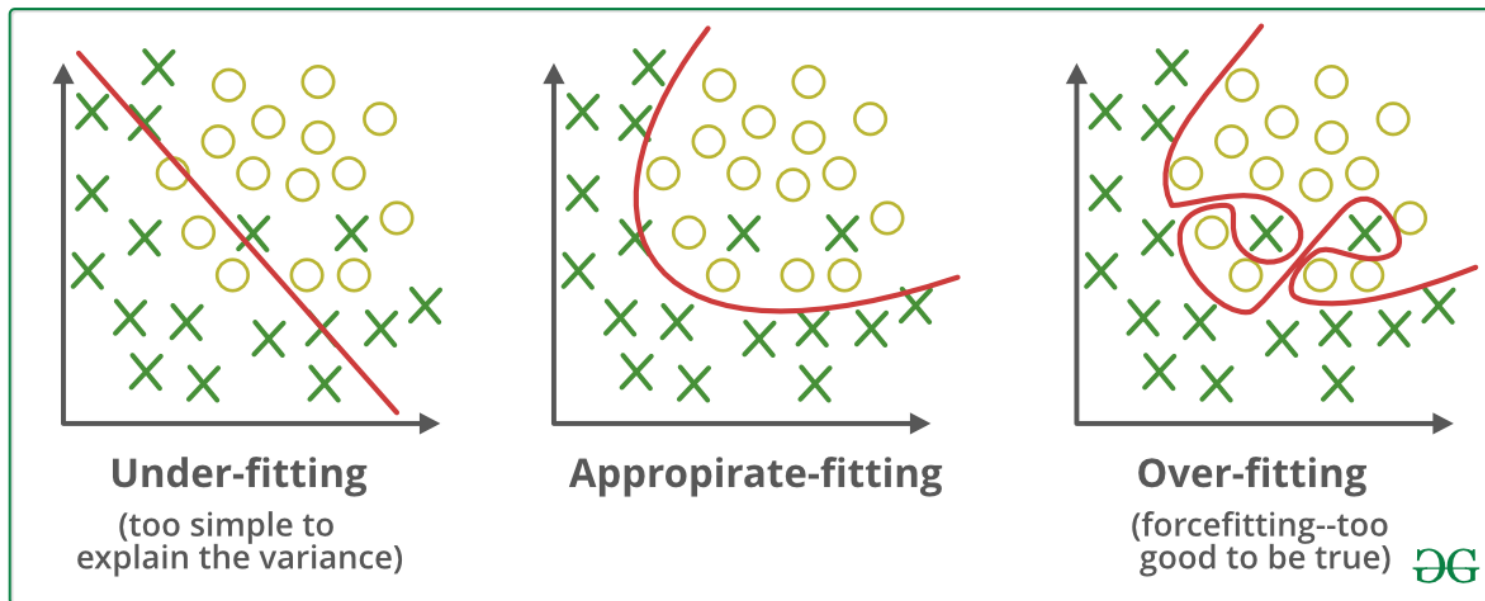
Learning rate

# Training the network



Backpropagation

Convolution → Pooling → Flattening → $Y_1$ $Y_2$ $Y_3$ $Y_4$

Input image • Convolutional layer • Pooling layer • Dense layers • Output layer

https://www.mdpi.com/2076-3417/10/4/1245/htm

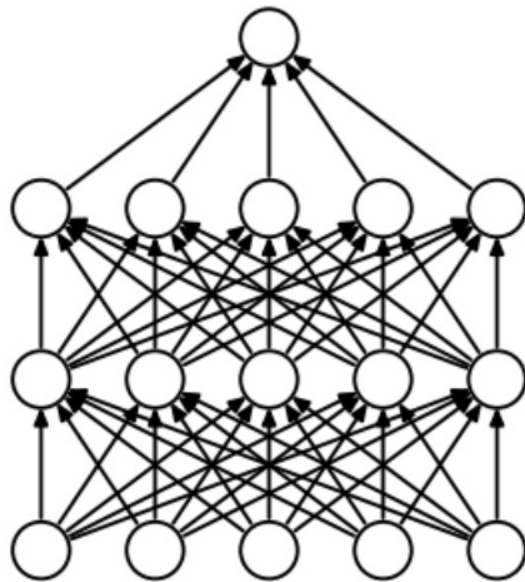# Dropout

- Overfitting



Under-fitting
(too simple to
explain the variance)

Appropirate-fitting

Over-fitting
(forcefitting--too
good to be true)
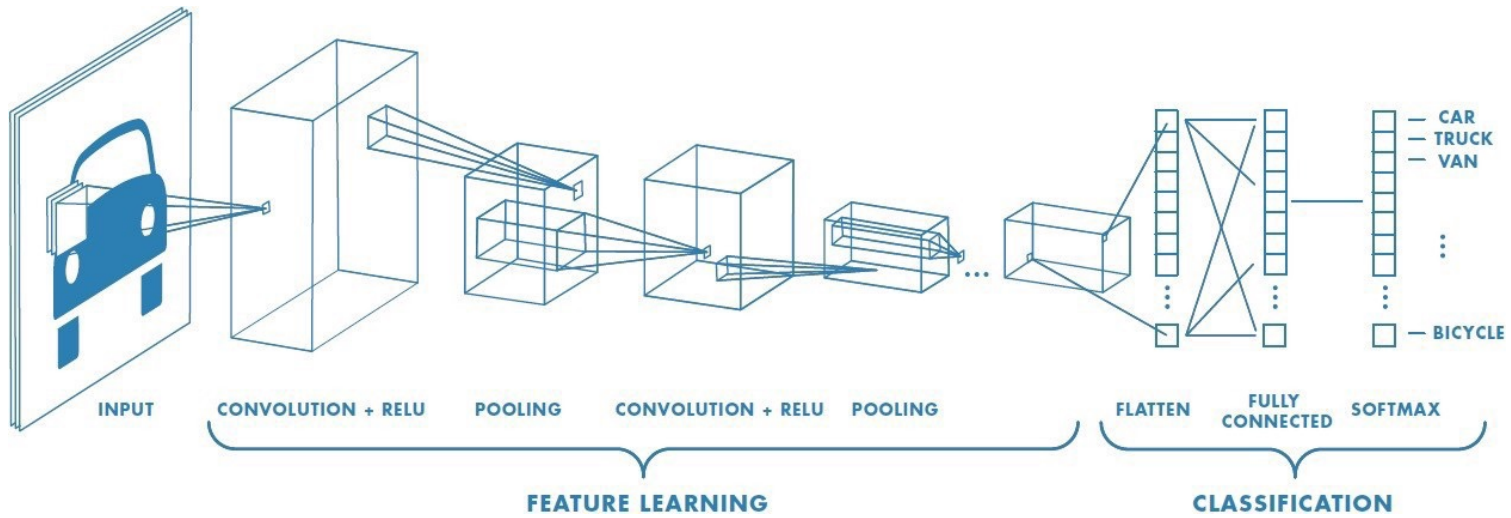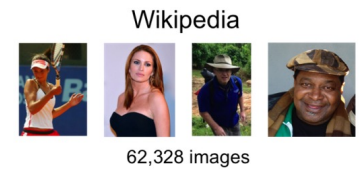
# Dropout



(a) Standard Neural Net

(b) After applying dropout.

INPUT   CONVOLUTION + RELU   POOLING   CONVOLUTION + RELU   POOLING   FLATTEN   FULLY CONNECTED   SOFTMAX

CAR
TRUCK
VAN
BICYCLE

FEATURE LEARNING

CLASSIFICATION

**CNN architectures** over a timeline(1998-2019)

AlexNet
2000

Inception-v1
2006

ResNet-50
2011

Inception-v4
2015

ResNeXt-50
2019

1998
LeNet-5

2003
VGG-16

2008
Inception-v3

2013
Xception

2017
Inception ResNets

IMAGENET

COCO 2020 Panoptic Segmentation Task

IMDb
460,723 images

Wikipedia
62,328 images

CIFAR-10

Fashion MNIST

# Applying to EEG

ERSP

Emotion/Traits

# Raw EEG vs. Frequency-domain



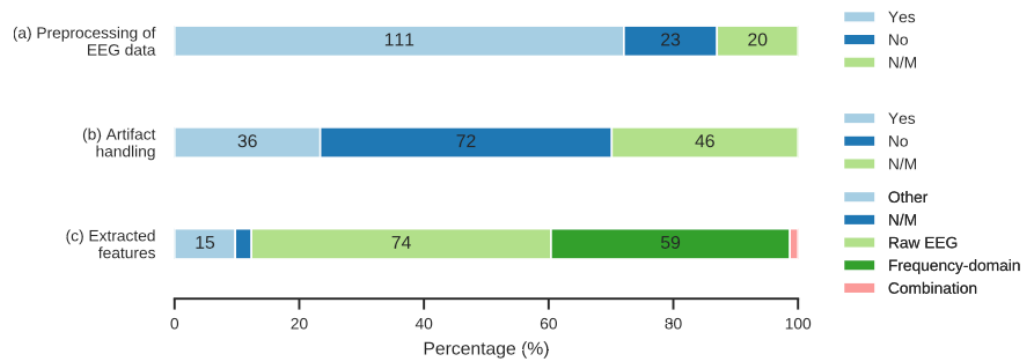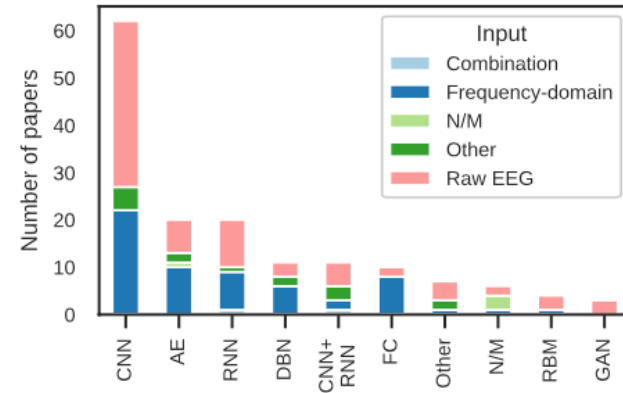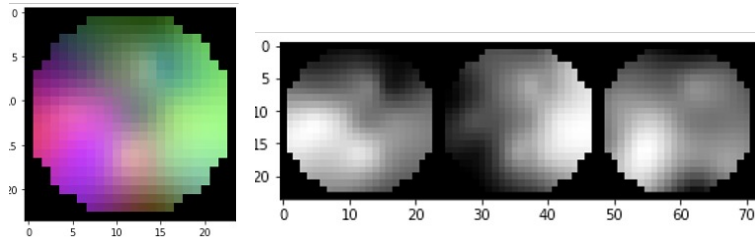Figure 9. EEG processing choices. (a) Number of studies that used preprocessing steps, such as filtering, (b) number of studies that included, rejected or corrected artifacts in their data and (c) types of features that were used as input to the proposed models.
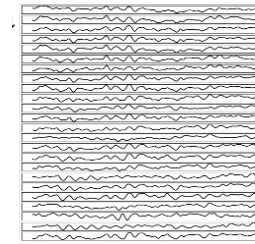


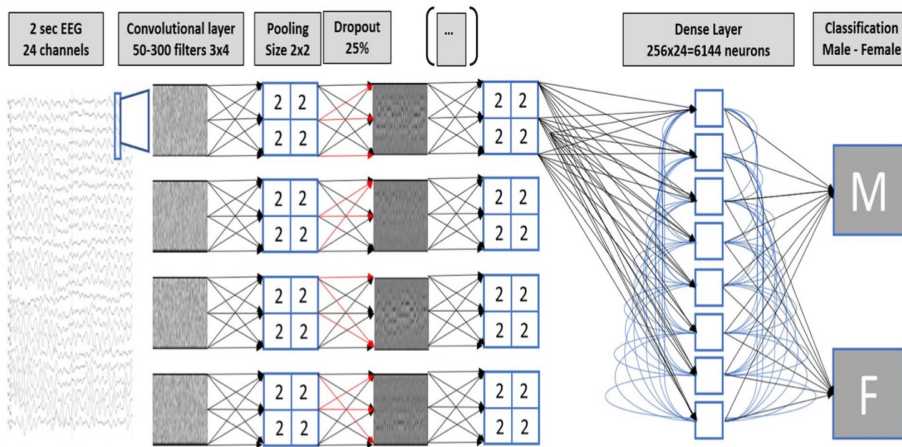Distribution of input type according to the architecture category.

Roy Y. et al (2019)
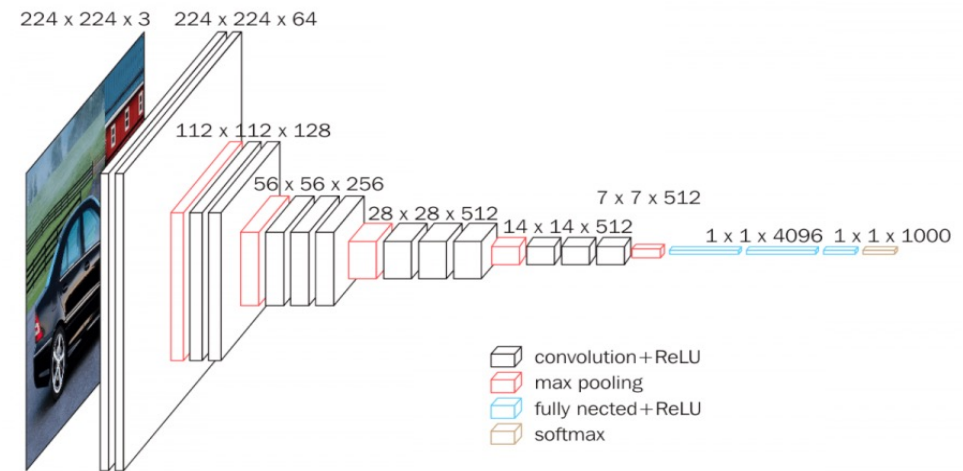
# Our approach



Spectral features



Raw EEG

vs.



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
- ~3,000 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
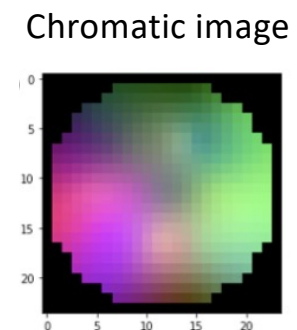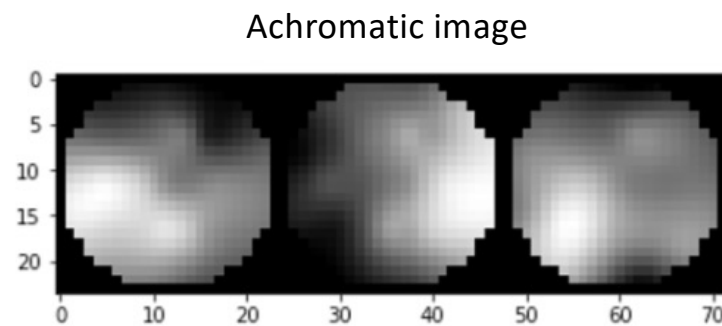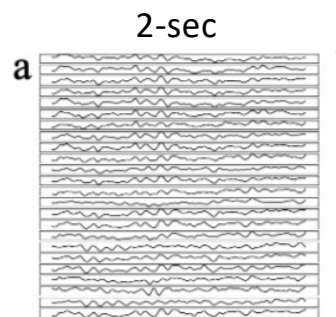    - → ~ 81 samples per subject

# Data

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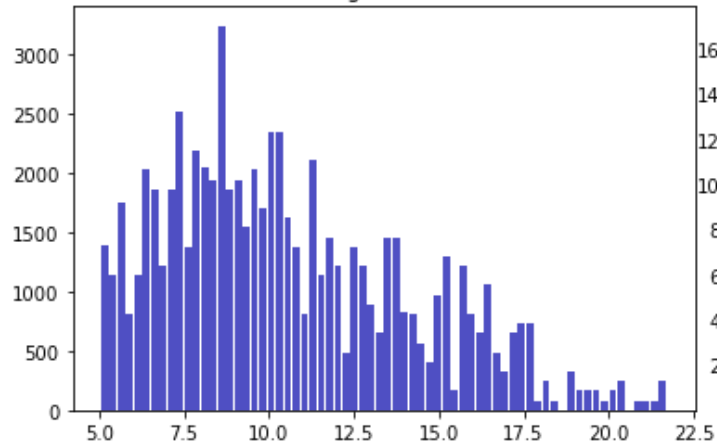- 5 predictors: sex, handedness, eyes open/closed, age, segment count

Categorical

Continuous



2-sec

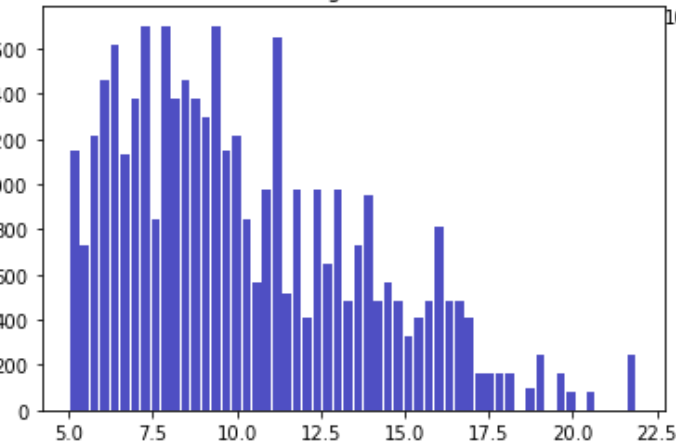Achromatic image

Chromatic image

# Data

- 10-30-60 split
  - 885 subjects for training     -> 71,381 samples
  - 492 subjects for validation -> 39,868 samples
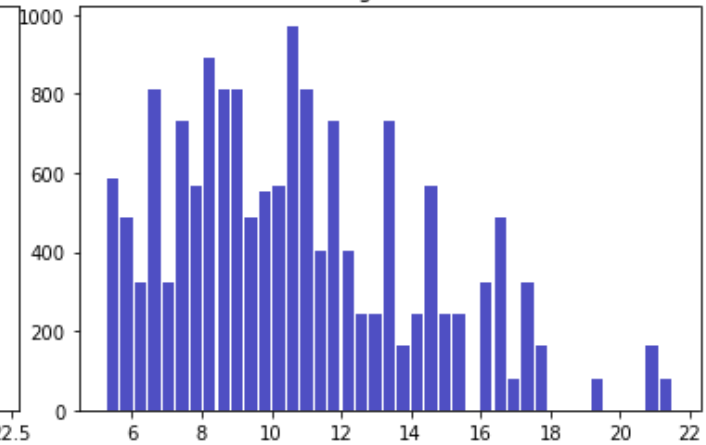  - 197 subjects for testing      -> 15,925 samples
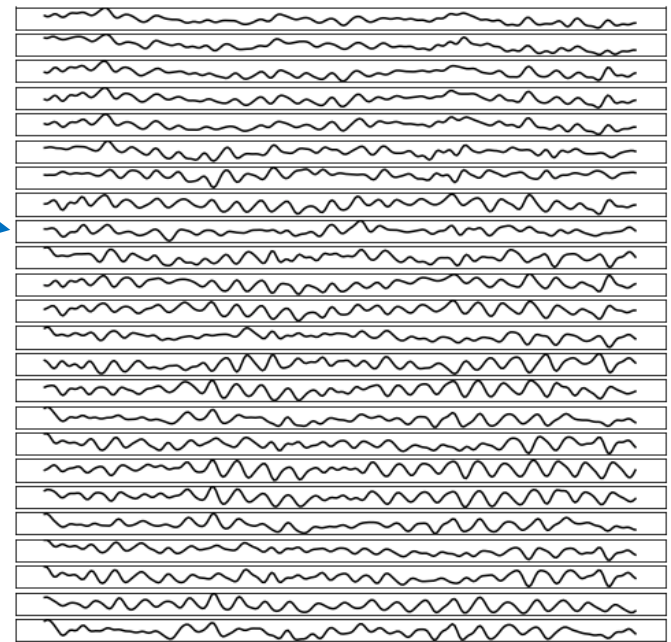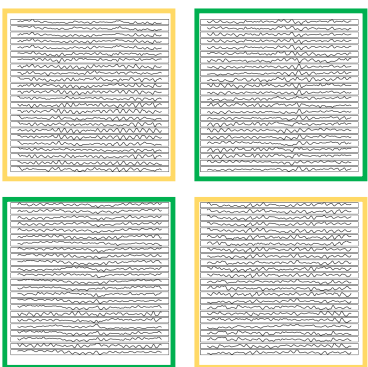
**Input data 24 x 256 x n**



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**

theta    alpha    beta

Spectral
amplitude

Tapered FFT
24 x 256 complex

Theta Alpha    Beta

| 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 l
- Dropped layers 19 to 32
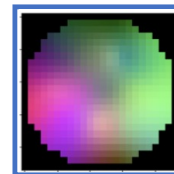- Change number of output classes to 2
- Retrain the whole network



| 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 | | |

R-VGG

Or

S-VGG

# Training the network

- Computing environment

# Training the network

- Programming environment



v3.7.10



v1.3.1

# Results



| Model | Per-sample |
|-------|-----------|
| R-SCNN | 80.6 (79.7 to 81.5) |
| R-VGG | **83.1** (82.7 to 83.4) |
| S-SCNN | 79.0 (78.7 to 79.3) |
| S-VGG | 77.1 (76.8 to 77.4) |

2-sec

Truong, D., Milham, M., Makeig, S., and Delorme, A. Deep Convolutional Neural Network Applied to Electroencephalography: Raw Data vs Spectral Feature. Annu Int Conf IEEE Eng Med Biol Soc, 2021, pp. 1039-1042, doi: 10.1109/EMBC46164.2021.9630708.

# Assessing features learned by the network



Krizhevsky et al., 2012

# 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)** |



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

# 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



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

# Discussion

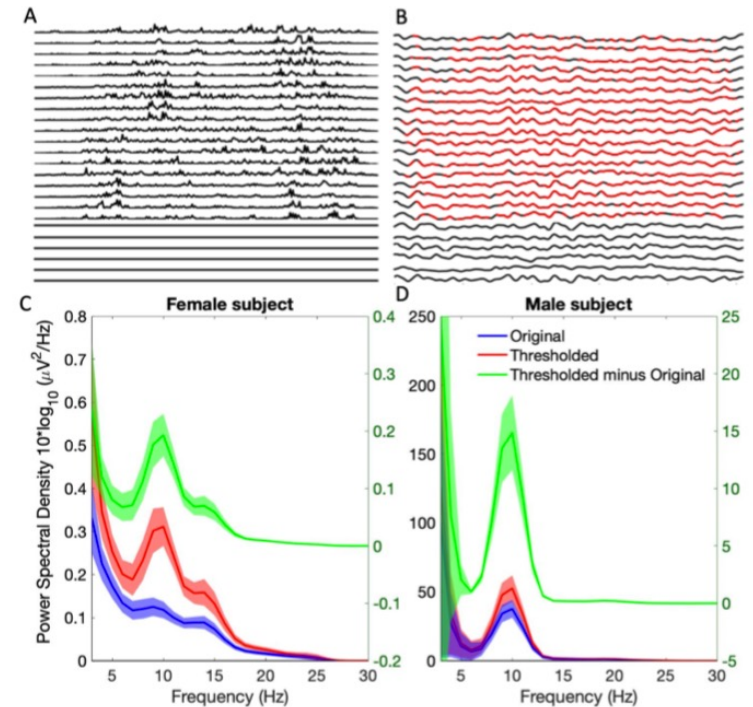- Architecture design is not physiologically driven
- Not generalizable for different number of channels and sampling rate
- Could be a simple problem: boys move more → Train on normalized data

# Decoding speech from non-invasive brain recordings

Alexandre Défossez[1,*], Charlotte Caucheteux[1,2], Jérémy Rapin[1], Ori Kabeli[1], and Jean-Rémi King[1,*]

# Locating and Editing Factual Associations in GPT

**Kevin Meng***
MIT CSAIL

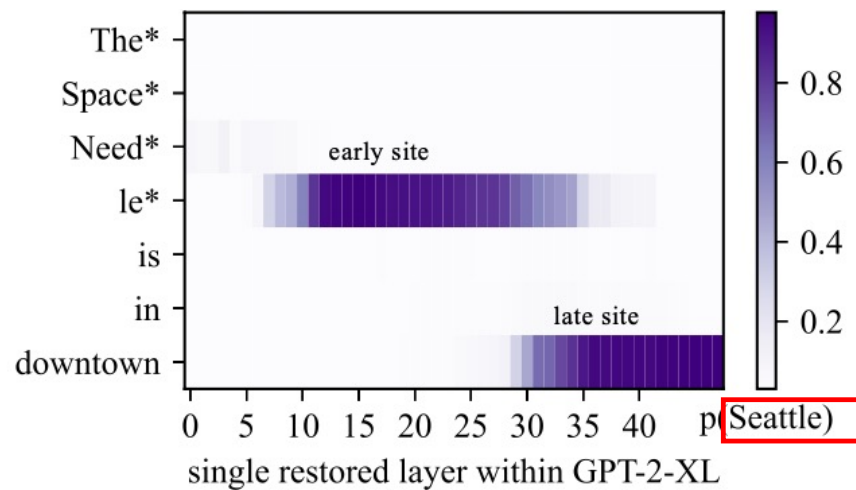**David Bau***
Northeastern University

**Alex Andonian**
MIT CSAIL

**Yonatan Belinkov**[†]
Technion – IIT

single restored layer within GPT-2-XL

# Locating and Editing Factual Associations in GPT

**Kevin Meng**[*]
MIT CSAIL

**David Bau**[*]
Northeastern University

**Alex Andonian**
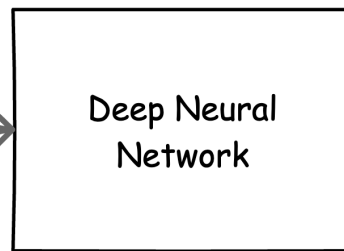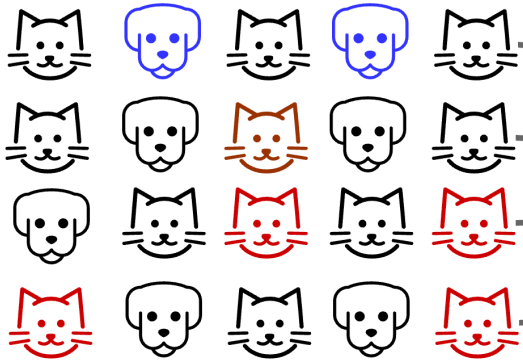MIT CSAIL

**Yonatan Belinkov**[†]
Technion – IIT

# Future work: representation learning

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



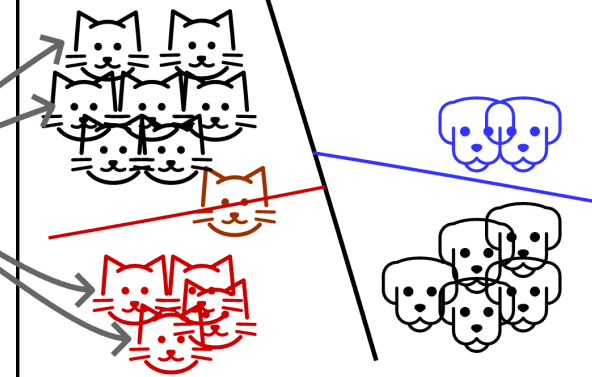Cat by Martin LEBRETON, Dog by Serhii Smirnov from the Noun Project

# Future work: representation learning

- DL models take data from the original space and map it to a "meaningful" representation
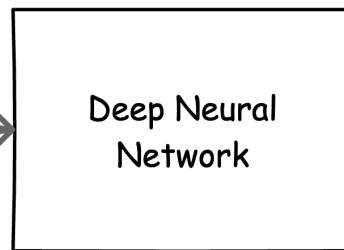
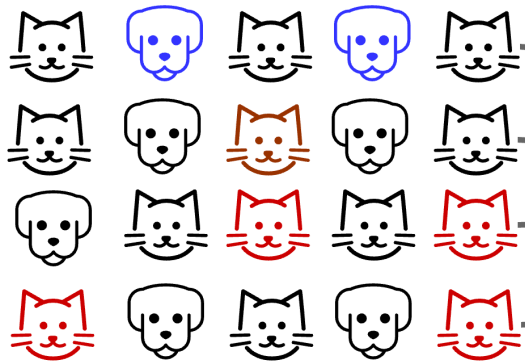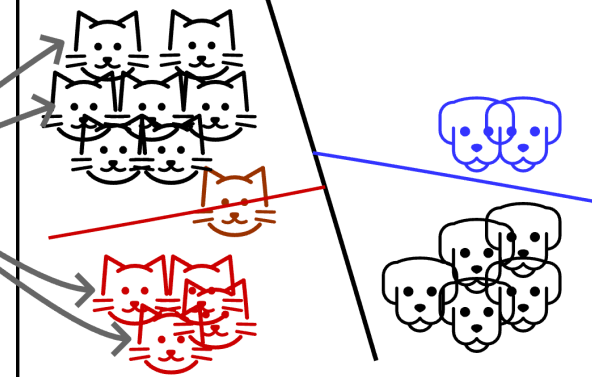- Can we learn meaningful embeddings that are also generalizable?



Default Representation

"Good" Semantic Representation

Deep Neural Network

Cat by Martin LEBRETON, Dog by Serhii Smirnov from the Noun Project

# 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