

# Like an Open Book? Read Neural Network Architecture with Simple Power Analysis on 32-bit Microcontrollers

Raphaël Joud<sup>1,2</sup>, Pierre-Alain Moëllic<sup>1,2</sup>, Simon Pontié<sup>1,2</sup>, and Jean-Baptiste Rigaud<sup>3</sup>

1 CEA Tech, Centre CMP, Equipe Commune CEA Tech - Mines Saint-Etienne, F-13541 Gardanne, France, {name.surname}@cea.fr

2 Univ. Grenoble Alpes, CEA, Leti, F-38000 Grenoble, France

3 Mines Saint-Etienne, CEA, Leti, Centre CMP, F-13541 Gardanne France, surname@emse.fr



# Outline

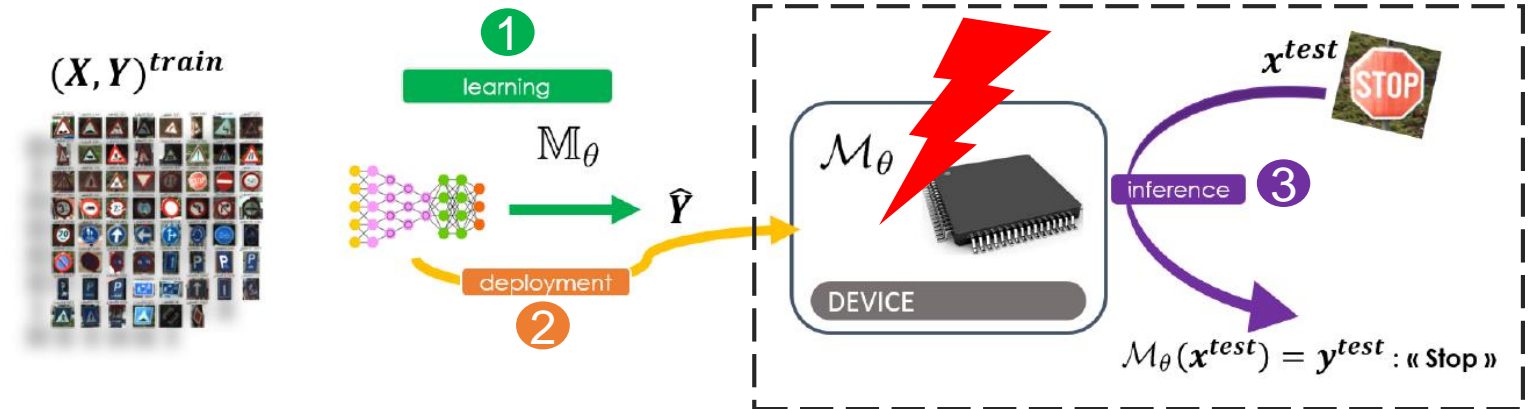
- Introduction
- Challenges
- Scope and threat model
- Model analysis
- Layer analysis
- Discussion & Perspectives



# Threats related to Machine Learning

- Large-scale Machine Learning (ML) model deployment
- Wide variety of applications and HW platforms involved
- Questions about ML security and related attack surface
- Attacking a ML model:

- Integrity
- Availability
- **Confidentiality**

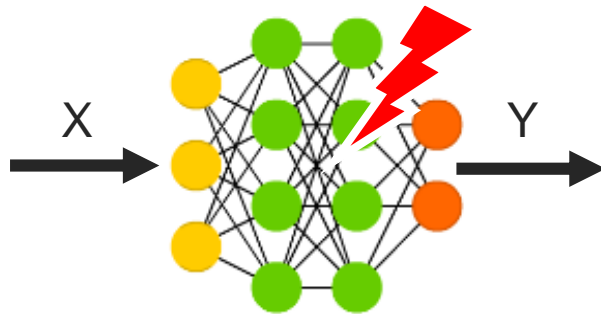


# Overall attack surface

## Attack Surface

Algorithmic / API-based attacks

Target theoretical flaws – Include adversarial examples, data poisoning, membership inference, model extraction...



# Overall attack surface

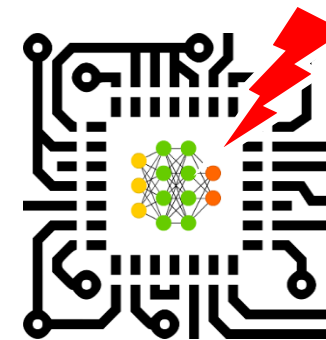
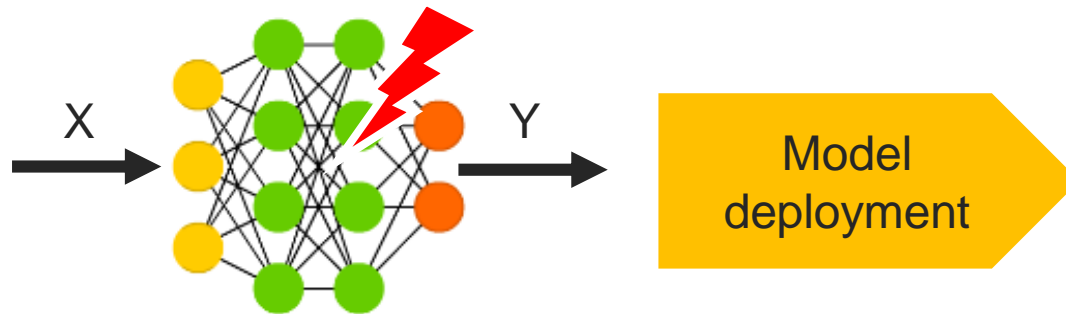
## Attack Surface

Algorithmic / API-based attacks

**Implementation-based attacks**

Target theoretical flaws – Include adversarial examples, data poisoning, membership inference, model extraction...

Target implementation flaws – Include physical attacks like Fault Injection Analysis (FIA) and **Side-Channel Analysis (SCA)**



# Outline

- Introduction
- **Challenges**
- Scope and threat model
- Model analysis
- Layer analysis
- Discussion & Perspectives



# Stakes related to model architecture

- Third-party vs Ad-hoc architectures
- Architecture knowledge is required to carry out various attacks
- Often required for parameters extraction whatever the attack surface considered [1 to 7]

Attack surface / method	Example of attacks requiring architecture knowledge
API-based	Cryptanalytic extraction [1,2] High Accuracy and High Fidelity [3]
FIA-based	Deepsteal [4] Fault Injection & Safe-error [5]
SCA-based	Practical Introduction [6] Reverse-engineering DNN [7]

# State of the art of architecture extraction

Attack	Physical target	Targeted models	Used techniques
Duddu <i>et al.</i> (2018) [8]	MLaaS	CNN	TA & Regression
Yu <i>et al.</i> (2020) [9] Yli <i>et al.</i> (2021) [10]	FPGA	BNN	SEMA
Luo <i>et al.</i> (2022) [11]	FPGA	CNN & ResNet	SEMA
Chmielewski <i>et al.</i> (2021) [12]	GPU	CNN	SEMA & TA
Batina <i>et al.</i> (2019) [13]	$\mu$ C	MLP & CNN	CEMA
Xiang <i>et al.</i> (2020) [14]	$\mu$ C	CNN	SPA & ML
<b>Ours</b>	$\mu$ C	<b>MLP &amp; CNN</b>	<b>SEMA only</b>

MLaaS: ML as a Service,  
TA: Timing Analysis,

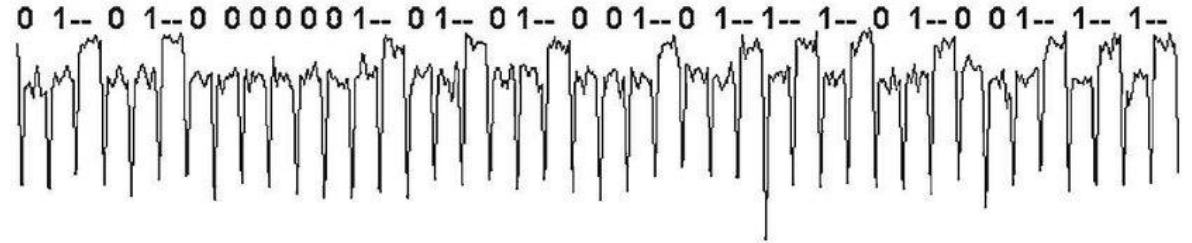
SEMA: Simple EM Analysis,  
SPA: Simple Power Analysis



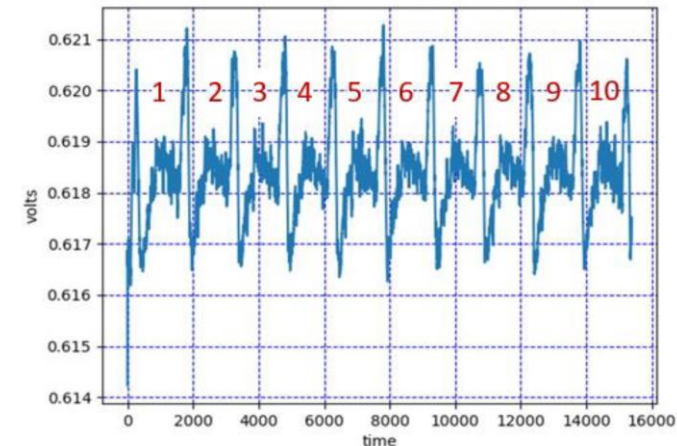
# What is a Simple Power / EM Analysis ?

SPA / SEMA:

- “Direct” interpretation of power consumption / EM measurements collected
- Can allow to identify parts of codes or even secret information
- Often requires knowledge of used implementation library



SPA on RSA emanations



SPA on AES emanations

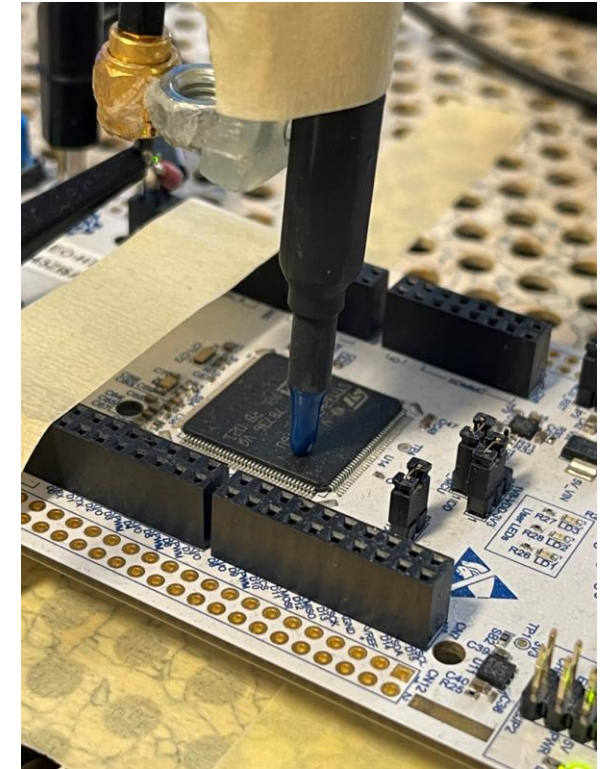
# Outline

- Introduction
- Challenges
- **Scope and threat model**
- Model analysis
- Layer analysis
- Discussion & Perspectives



# Scope and Threat model

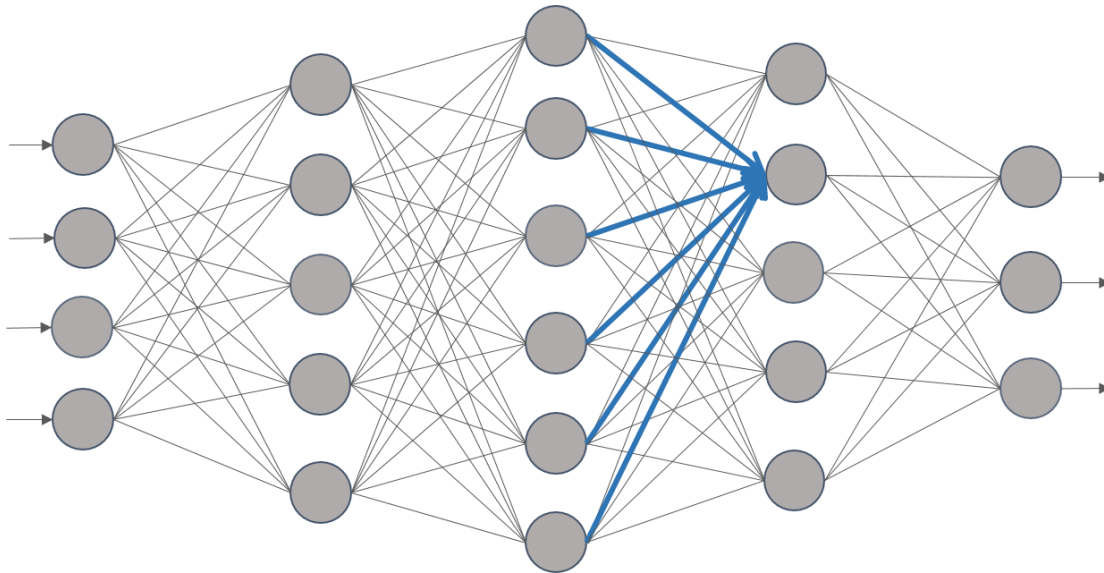
- Goal: extract models architecture as precisely as possible
- Targeted models: quantified (8-bits) and implemented with open-source libraries: **NNoM** [15] & **CMSIS-NN** by ARM [16]
  - ➔ Analyse layer C code to anticipate corresponding EM leakages (cache optimisations are disabled)
- Black-box context: no knowledge of model architecture nor parameters
- EM traces acquired with only 1 input: avoids desynchronisation problems & allows to average traces



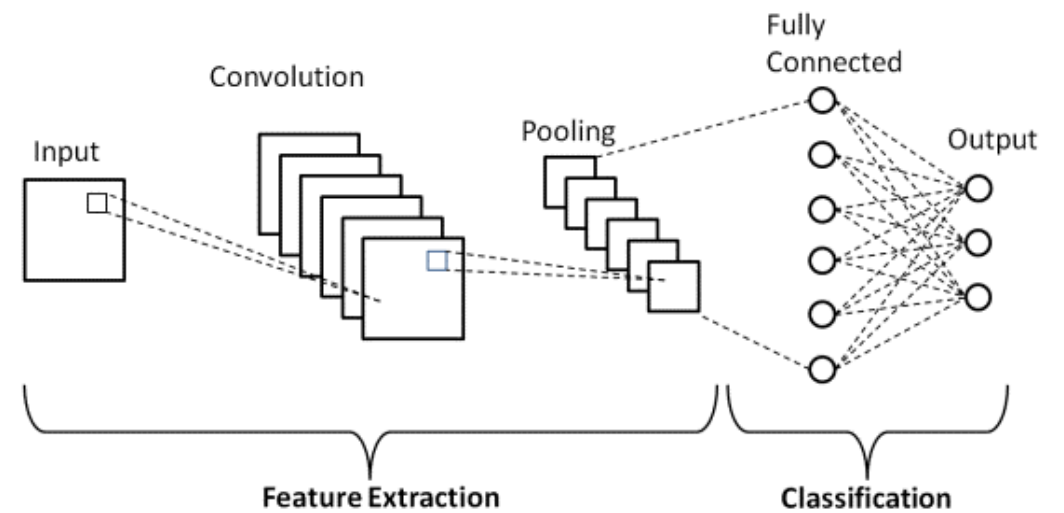
# Scope and Threat model

- Studied architectures:

## MultiLayer Perceptron (MLP)

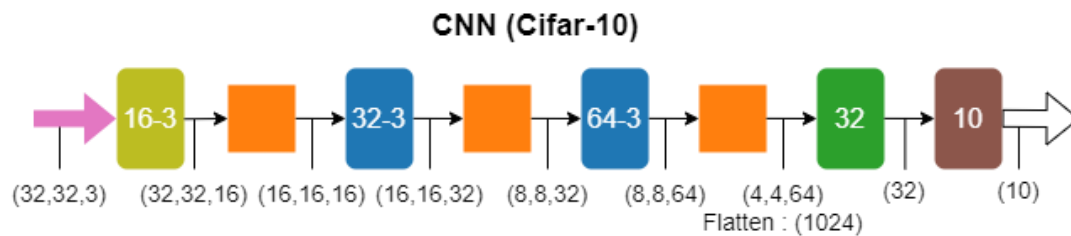
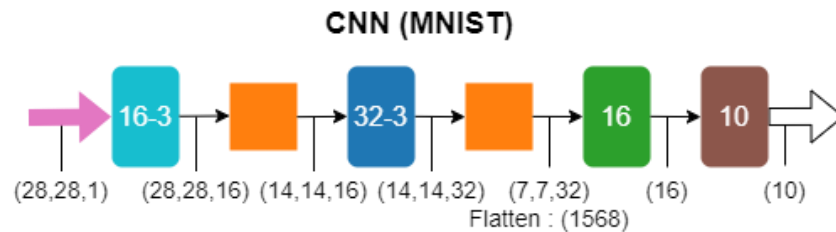
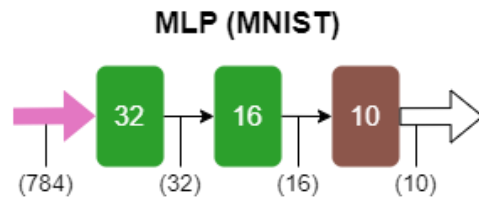


## Convolutional Neural Networks (CNN)



# Scope and Threat model

- 3 different targeted models:



Target	Hyper-parameter	Notation
Conv. layer	Output shape	$H_{out}$
	# Kernels	$K$
	Kernel size	$Z$
	Stride, Padding	$S, P$
MaxPool layer	Output shape	$H_{out}$
	Filter size	$Z_{pool}$
Dense layer	# Neurons	$N_e$
Activation layer	ReLU or not	$\emptyset$
Model	# Layers	$L$
	Layers nature	$\emptyset$

# Outline

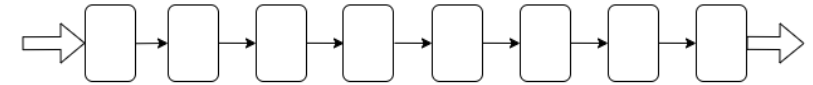
- Introduction
- Challenges
- Scope and threat model
- **Model analysis**
- Layer analysis
- Discussion & Perspectives



# Model Analysis

Analyse the overall shape of inference trace:

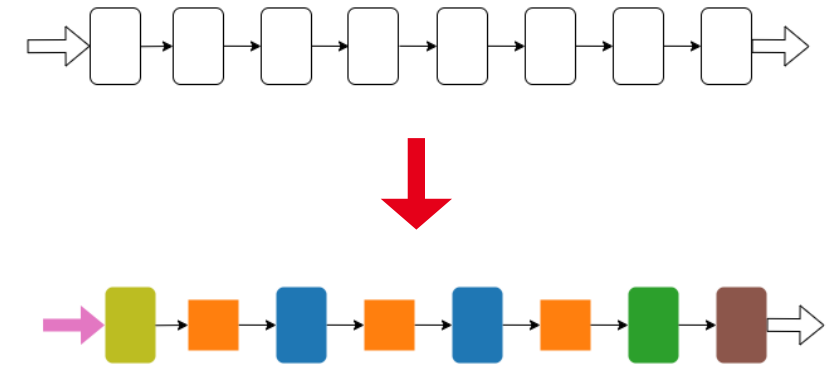
1. Find the number of layers:
  - Straightforward or can rely on frequency spectrum
2. Identify their nature:
  - Use layer execution time (complexity)
  - Use patterns of specific layer
3. Extract hyper-parameters of each of them



# Model Analysis

Analyse the overall shape of inference trace:

1. Find the number of layers:
  - Straightforward or can rely on frequency spectrum
2. Identify their nature:
  - Use layer execution time (complexity)
  - Use patterns of specific layer
3. Extract hyper-parameters of each of them

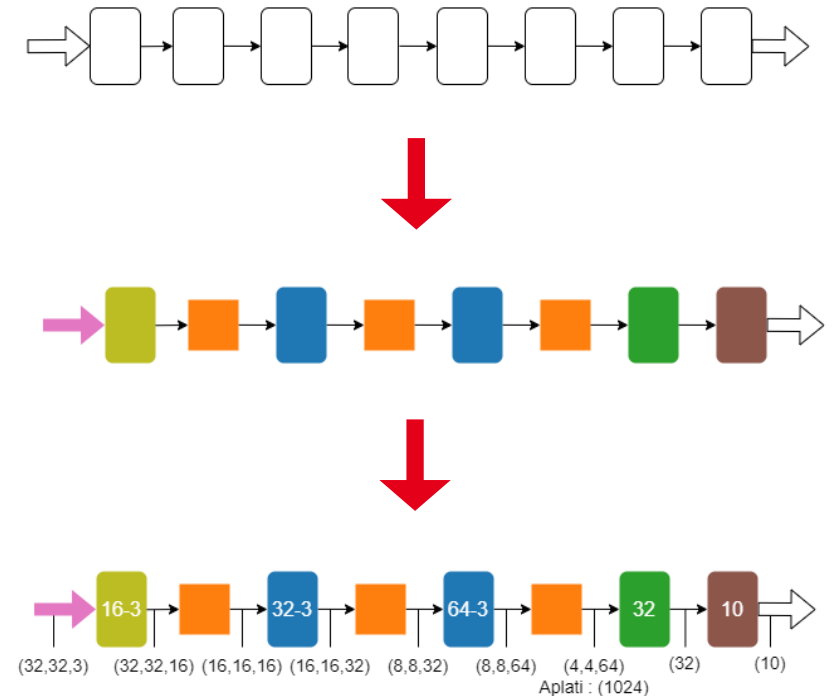




# Model Analysis

Analyse the overall shape of inference trace:

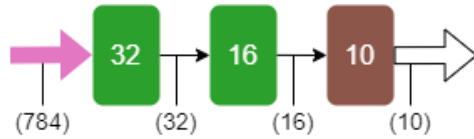
1. Find the number of layers:
  - Straightforward or can rely on frequency spectrum
2. Identify their nature:
  - Use layer execution time (complexity)
  - Use patterns of specific layer
3. Extract hyper-parameters of each of them



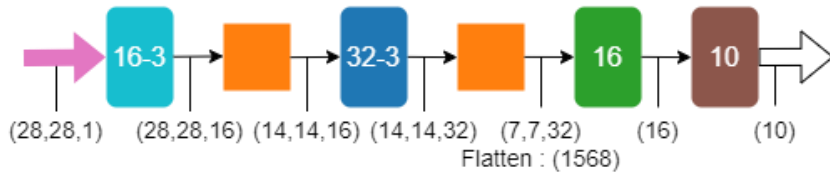
# Model Analysis



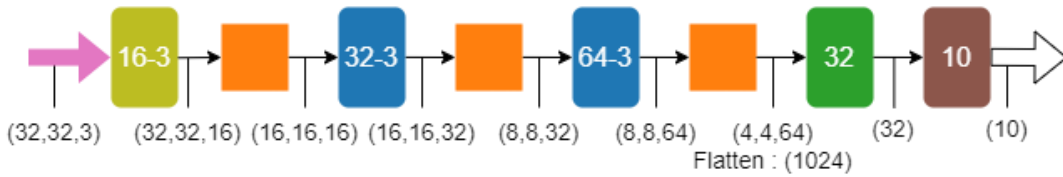
### MLP (MNIST)



### CNN (MNIST)



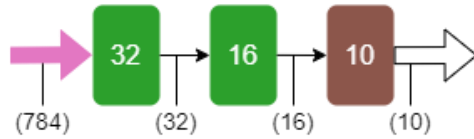
### CNN (Cifar-10)



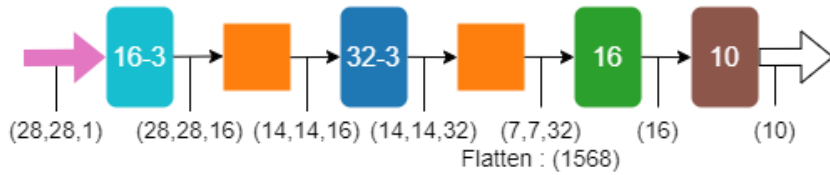
# Model Analysis



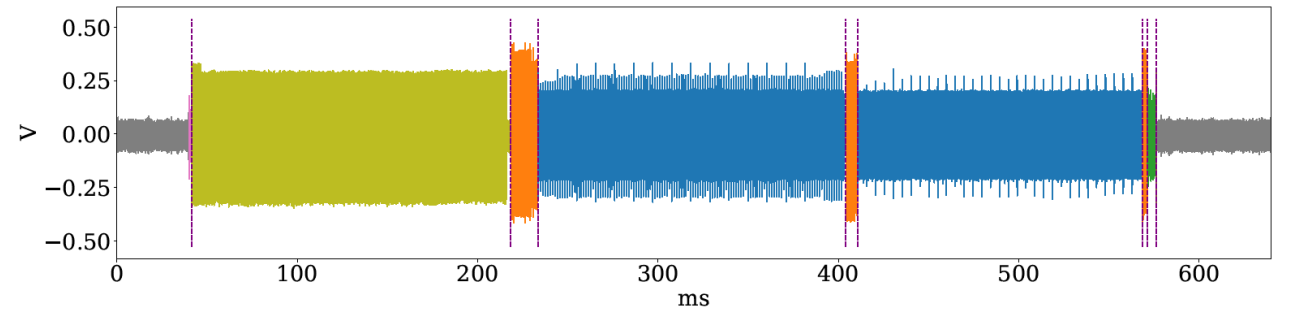
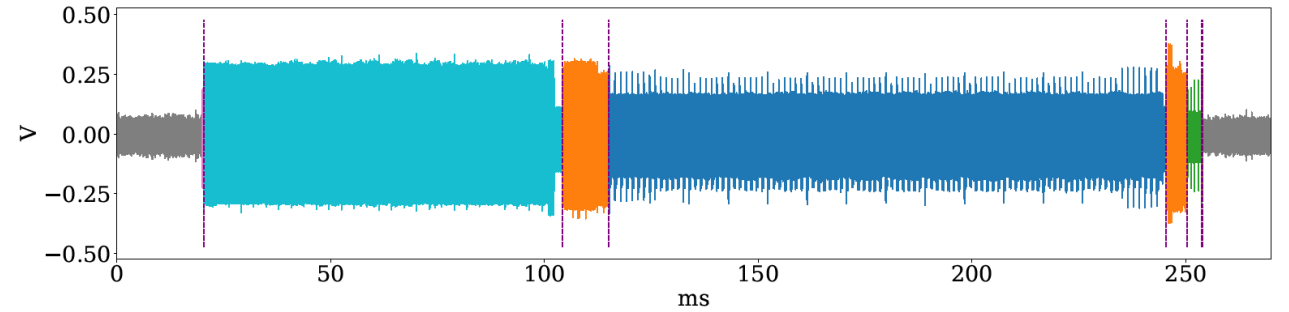
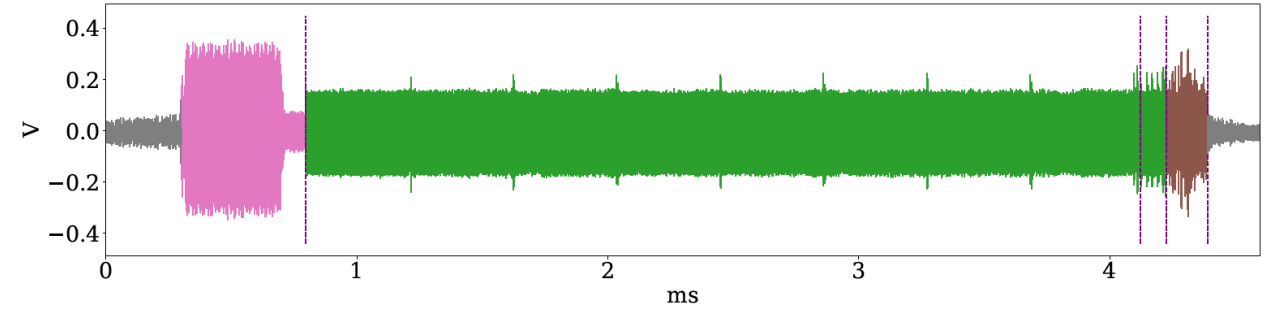
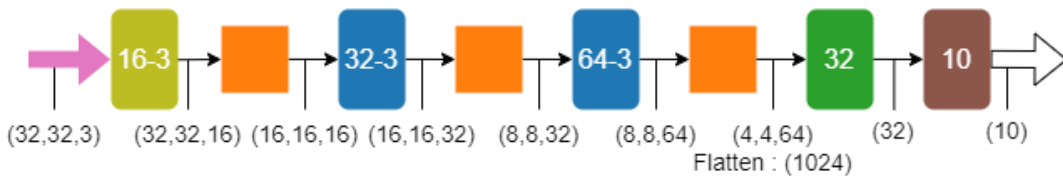
MLP (MNIST)



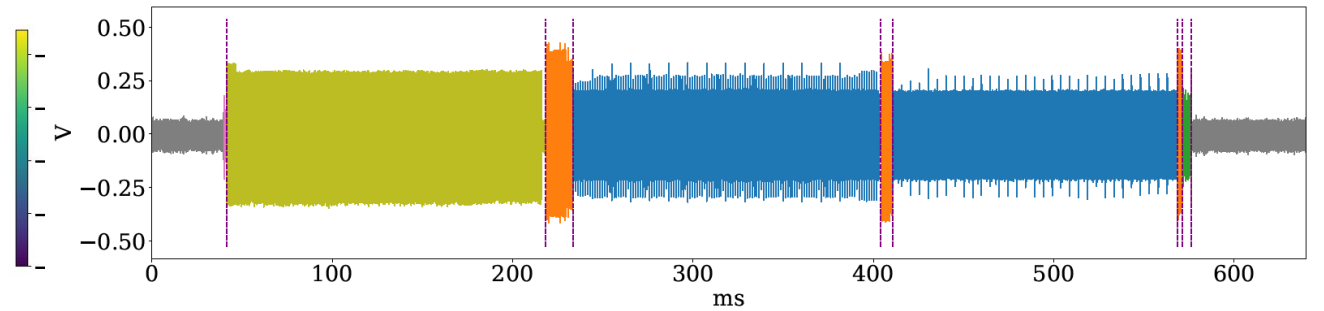
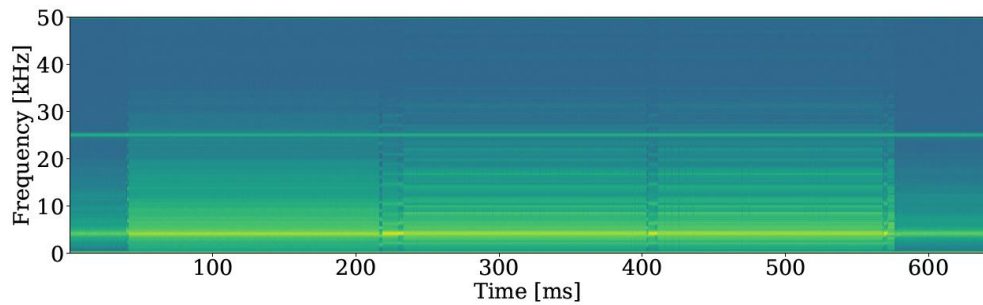
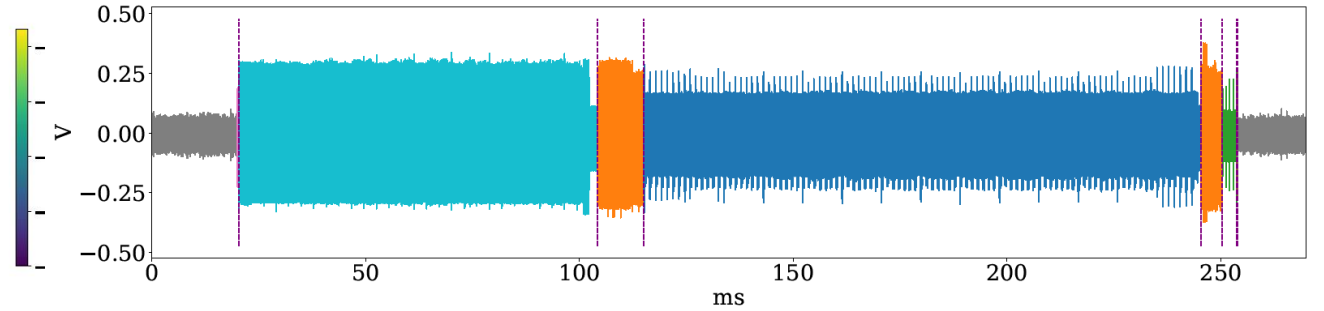
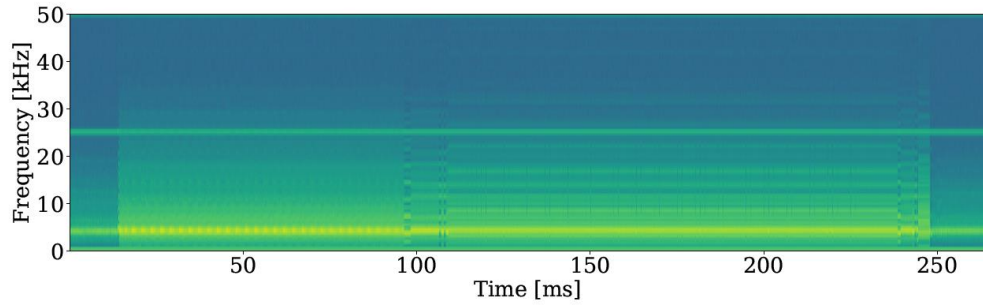
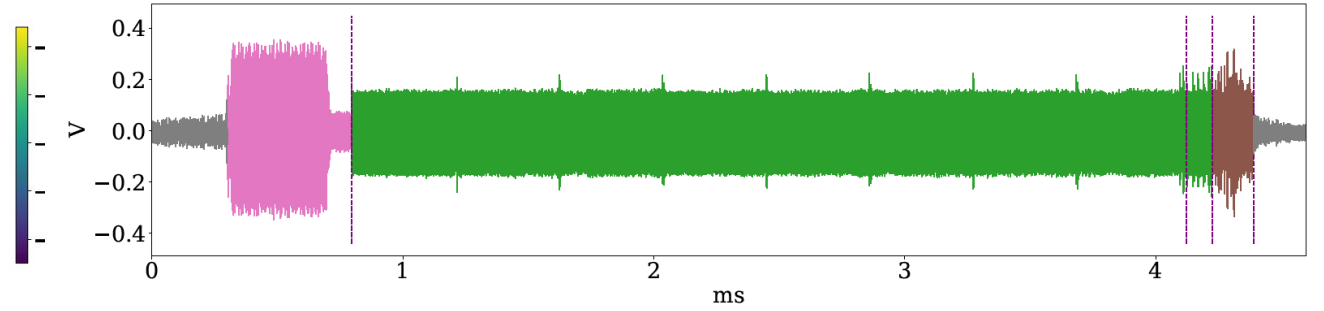
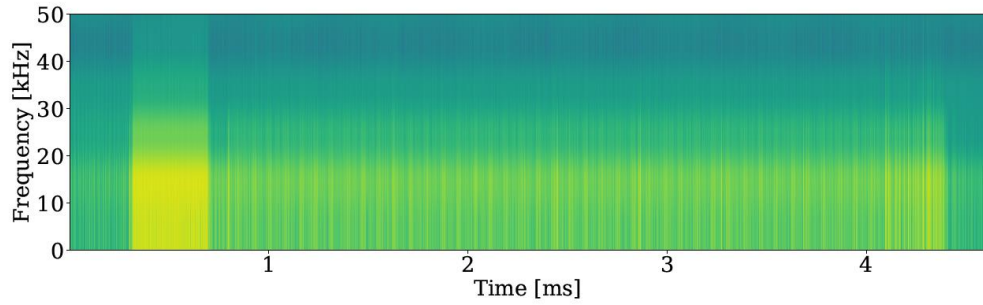
CNN (MNIST)



CNN (Cifar-10)



# Model Analysis



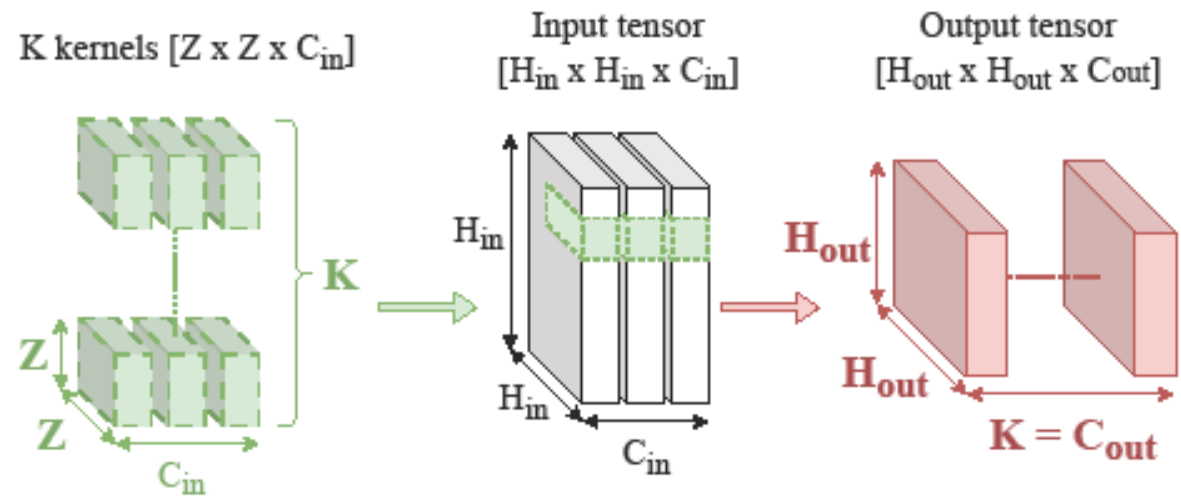
# Outline

- Introduction
- Challenges
- Scope and threat model
- Model analysis
- **Layer analysis**
- Discussion & Perspectives



# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$




---

## Algorithm 1 General convolution implementation

---

**Input:** Input tensor  $I_{in}$  of size  $H_{in}^2 \cdot C_{in}$ ,  $ker$  (Kernel tensor),  $S$ ,  $P$ , Output size  $I_{out}$  of size  $H_{out}^2 \cdot C_{out}$

**Output:** Filled  $I_{out}$

```

1: for  $i_y \leftarrow 0, i_y < H_{out}, i_y ++1$  do
2:   for  $i_x \leftarrow 0, i_x < H_{out}, i_x ++1$  do
3:      $buff_{in} \leftarrow im2col(I_{in}, i_y, i_x, S, P, H_{in}, C_{in})$ 
4:     if  $len(buff_{in}) == 2 \times C_{in} \times Z^2$  then
5:        $GeMM(buff_{in}, ker, C_{out}, C_{in} \times Z^2, I_{out})$ 
6:        $buff_{in} \leftarrow 0$ 
7:     end if
8:   end for
9: end for

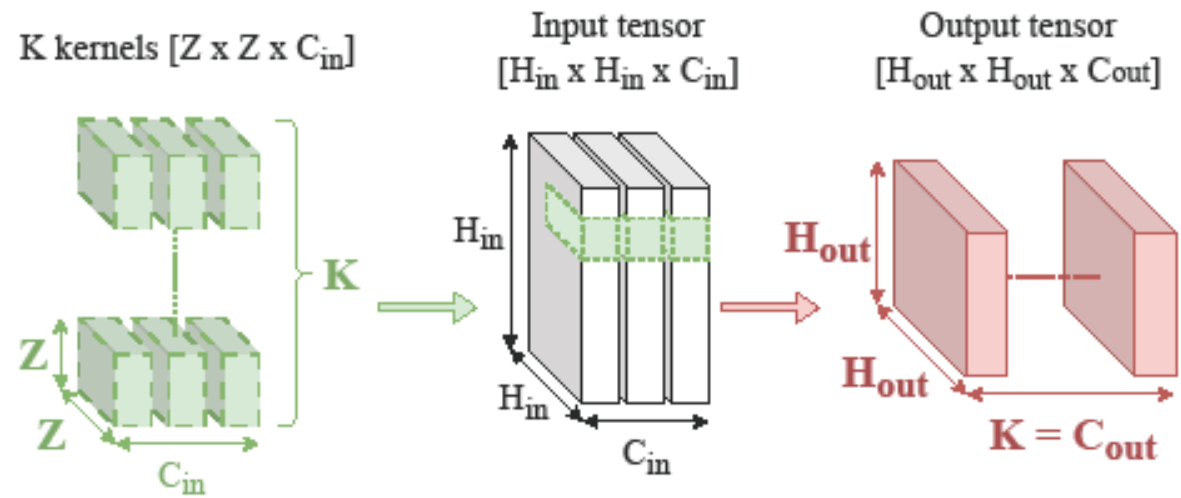
```

- ▷ Iterate over  $H_{out}$  ( $y$ -axis)
- ▷ idem ( $x$ -axis)
- ▷ Apply Im2col conversion
- ▷ Checks if 2 input columns are set
- ▷ Perform matrix-product
- ▷ Buffer reset

---

# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$




---

## Algorithm 1 General convolution implementation

---

**Input:** Input tensor  $I_{in}$  of size  $H_{in}^2 \cdot C_{in}$ ,  $ker$  (Kernel tensor),  $S$ ,  $P$ , Output size  $I_{out}$  of size  $H_{out}^2 \cdot C_{out}$

**Output:** Filled  $I_{out}$

```

1: for  $i_y \leftarrow 0, i_y < H_{out}, i_y + 1$  do           ▷ Iterate over  $H_{out}$  ( $y$ -axis)
2:   for  $i_x \leftarrow 0, i_x < H_{out}, i_x + 1$  do     ▷ idem ( $x$ -axis)
3:      $buff_{in} \leftarrow im2col(I_{in}, i_y, i_x, S, P, H_{in}, C_{in})$    ▷ Apply Im2col conversion
4:     if  $len(buff_{in}) == 2 \times C_{in} \times Z^2$  then   ▷ Checks if 2 input columns are set
5:        $GeMM(buff_{in}, ker, C_{out}, C_{in} \times Z^2, I_{out})$    ▷ Perform matrix-product
6:        $buff_{in} \leftarrow 0$                                ▷ Buffer reset
7:     end if
8:   end for
9: end for

```

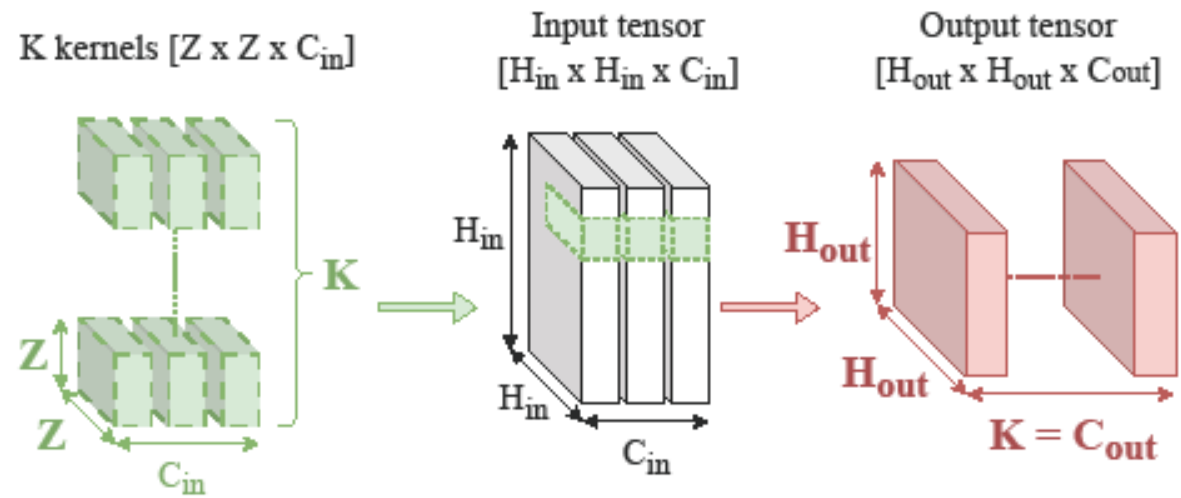
$H_{out} \times H_{out}$  iterations

GeMM function called  
 $H_{out} \times (H_{out}/2)$  times



# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$




---

## Algorithm 1 General convolution implementation

---

**Input:** Input tensor  $I_{in}$  of size  $H_{in}^2 \cdot C_{in}$ ,  $ker$  (Kernel tensor),  $S$ ,  $P$ , Output size  $I_{out}$  of size  $H_{out}^2 \cdot C_{out}$

**Output:** Filled  $I_{out}$

```

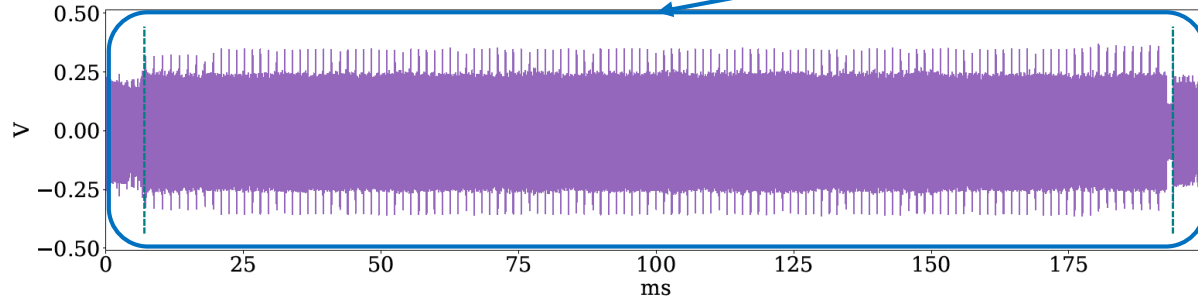
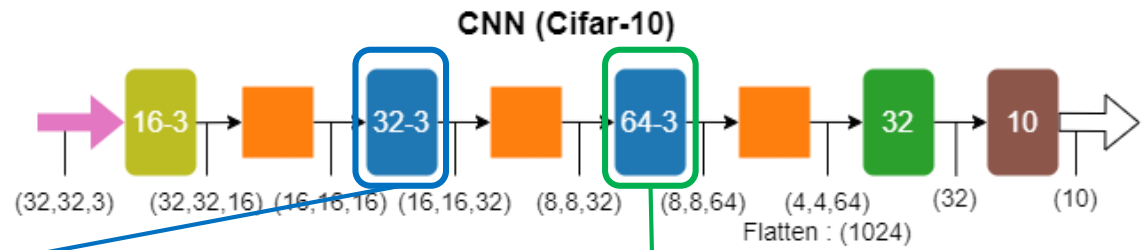
1: for  $i_y \leftarrow 0, i_y < H_{out}, i_y ++1$  do ▷ Iterate over  $H_{out}$  ( $y$ -axis)
2:   for  $i_x \leftarrow 0, i_x < H_{out}, i_x ++1$  do ▷ idem ( $x$ -axis)
3:      $buff_{in} \leftarrow im2col(I_{in}, i_y, i_x, S, P, H_{in}, C_{in})$  ▷ Apply Im2col conversion
4:     if  $len(buff_{in}) == 2 \times C_{in} \times Z^2$  then ▷ Checks if 2 input columns are set
5:        $GeMM(buff_{in}, ker, C_{out}, C_{in} \times Z^2, I_{out})$  ▷ Perform matrix-product
6:        $buff_{in} \leftarrow 0$  ▷ Buffer reset
7:     end if
8:   end for
9: end for
  
```

With  $N_p$  the number of pattern corresponding to GeMM function call:

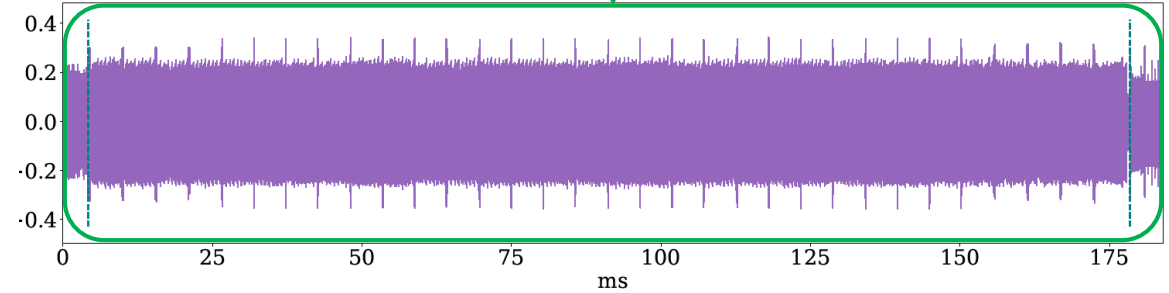
$$H_{out} = \sqrt{2 \times N_p}$$



# Convolution layers



$$N_p = 128 \text{ iterations} \Leftrightarrow H_{out} = 16$$

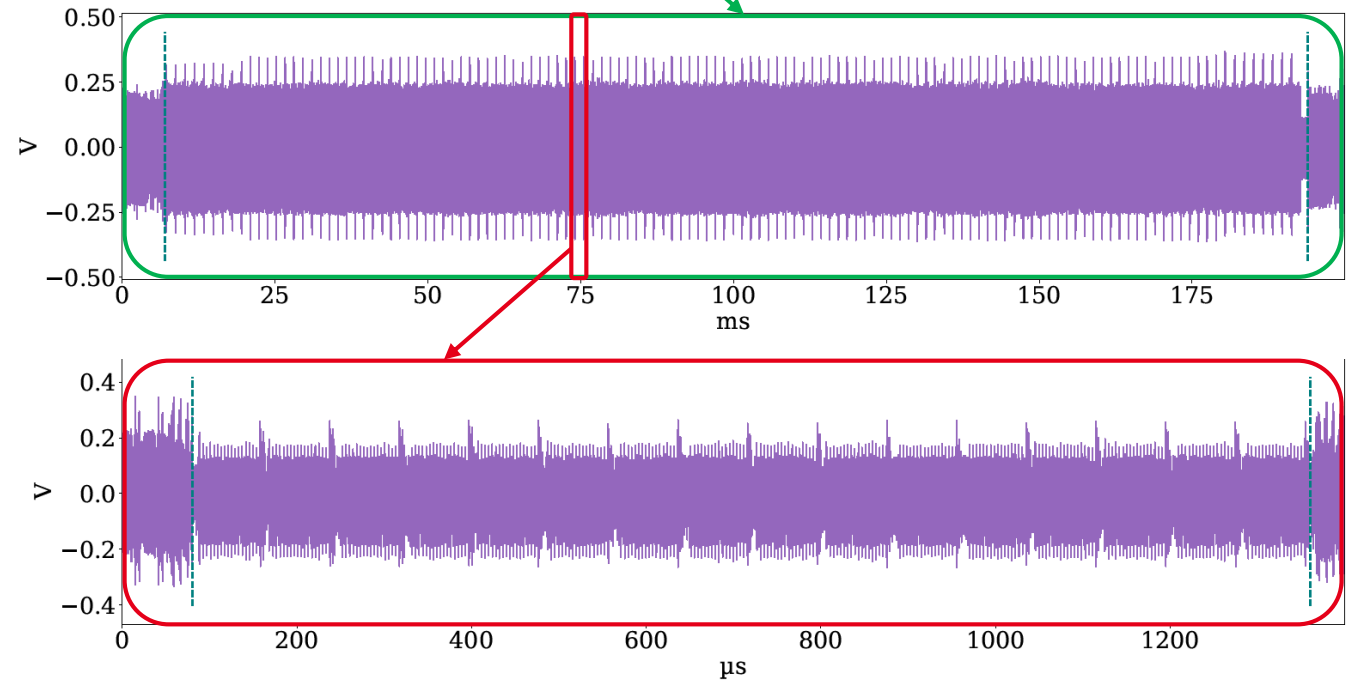
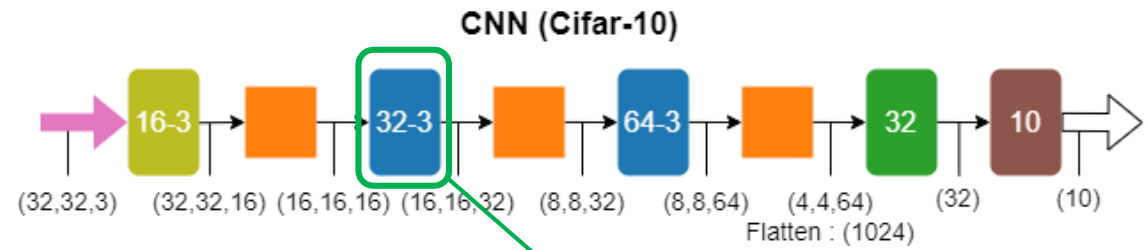


$$N_p = 32 \text{ iterations} \Leftrightarrow H_{out} = 8$$

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$
- $\rightarrow H_{out} = \sqrt{2 \times N_p}$

# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$
- Russian dolls effect to get  $K$



# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$
- Russian dolls effect to get  $K$

---

## Algorithm 1 Matrix-product for convolution (GeMM)

---

**Input:**  $buff_{in}$ ,  $ker$ ,  $C_{out}$ ,  $C_{in} \times Z^2$ ,  $I_{out}$ ,  $bias$  (bias tensor)

**Output:** Partly filled  $I_{out}$

```
1: rowCnt ←  $C_{out} \gg 1$  ▷ Set rowCnt =  $K/2$ 
2: for rowCnt > 0, rowCnt - -1 do ▷ Iterate over  $K/2$ 
3:   sum, sum1, sum2, sum3 = init_sum(bias,  $C_{in} \times Z^2$ )
4:   colCnt ←  $(C_{in} \times Z^2) \gg 2$  ▷ Set colCnt
5:   for colCnt > 0, colCnt - -1 do
6:     simd_mac(sum, sum1, sum2, sum3, buff_in, ker)
7:   end for
8:   if  $(C_{in} \times Z^2) \& 0x3$  then
9:     Manage_colCnt_remainder(sum, sum1, sum2, sum3, buff_in, bias,  $C_{in} \times Z^2$ )
10:  end if
11:  apply_mac(sum, sum1, sum2, sum3,  $C_{out}$ ,  $I_{out}$ )
12: end for
13: if  $C_{out} \& 0x3$  then ▷ Check presence of remainders
14:   sum, sum1 = init_sum(bias,  $C_{in} \times Z^2$ )
15:   colCnt ←  $(C_{in} \times Z^2) \gg 2$  ▷ Set colCnt
16:   for colCnt > 0, colCnt - -1 do
17:     simd_mac(sum, sum1, buff_in, ker)
18:   end for
19:   if  $(C_{in} \times Z^2) \& 0x3$  then
20:     Manage_colCnt_remainder(sum, sum1, buff_in, bias,  $C_{in} \times Z^2$ )
21:   end if
22:   apply_mac(sum, sum1,  $C_{out}$ ,  $I_{out}$ )
23: end if
```

---

# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$
- Russian dolls effect to get  $K$
- With  $N_p$  the number of visible patterns corresponding to iterations of for loop (3-11):

$$rowCnt: K = C_{out} = 2 \times N_p$$

---

## Algorithm 1 Matrix-product for convolution (GeMM)

---

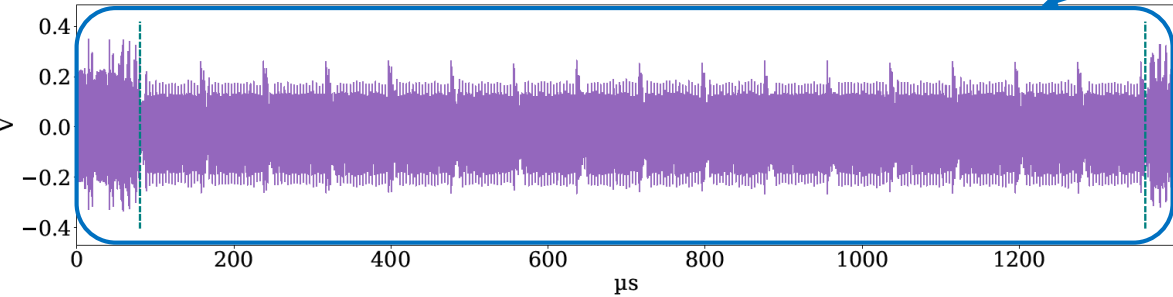
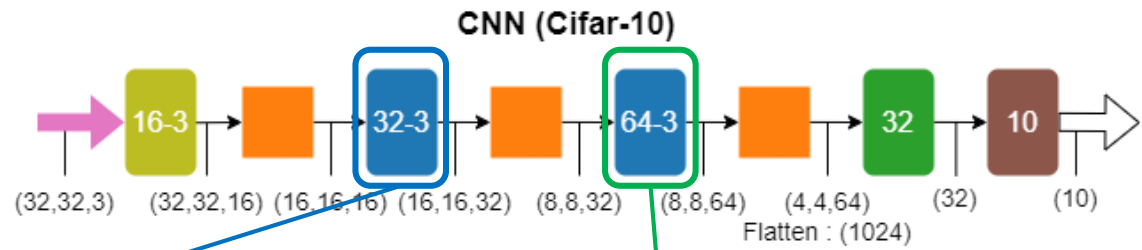
**Input:**  $buff_{in}$ ,  $ker$ ,  $C_{out}$ ,  $C_{in} \times Z^2$ ,  $I_{out}$ ,  $bias$  (bias tensor)

**Output:** Partly filled  $I_{out}$

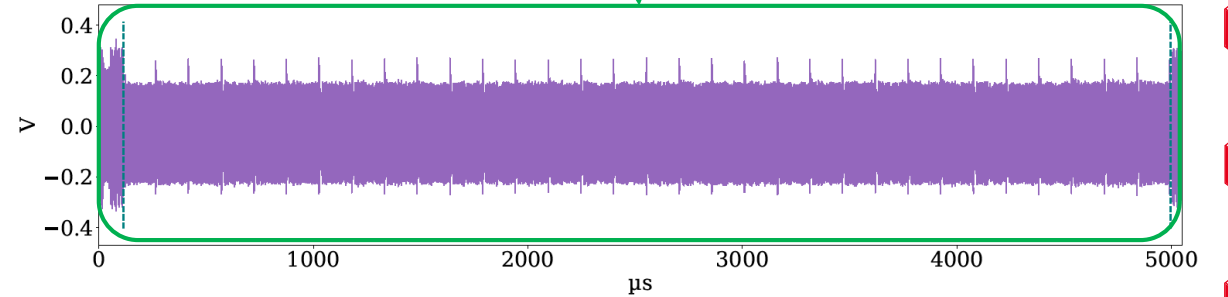
```
1: rowCnt ←  $C_{out} \gg 1$  ▷ Set rowCnt =  $K/2$ 
2: for rowCnt > 0, rowCnt - -1 do ▷ Iterate over  $K/2$ 
3:   sum, sum1, sum2, sum3 = init_sum(bias,  $C_{in} \times Z^2$ )
4:   colCnt ←  $(C_{in} \times Z^2) \gg 2$  ▷ Set colCnt
5:   for colCnt > 0, colCnt - -1 do
6:     simd_mac(sum, sum1, sum2, sum3, buff_in, ker)
7:   end for
8:   if  $(C_{in} \times Z^2) \& 0x3$  then
9:     Manage_colCnt_remainder(sum, sum1, sum2, sum3, buff_in, bias,  $C_{in} \times Z^2$ )
10:  end if
11:  apply_mac(sum, sum1, sum2, sum3,  $C_{out}$ ,  $I_{out}$ )
12: end for
13: if  $C_{out} \& 0x3$  then ▷ Check presence of remainders
14:   sum, sum1 = init_sum(bias,  $C_{in} \times Z^2$ )
15:   colCnt ←  $(C_{in} \times Z^2) \gg 2$  ▷ Set colCnt
16:   for colCnt > 0, colCnt - -1 do
17:     simd_mac(sum, sum1, buff_in, ker)
18:   end for
19:   if  $(C_{in} \times Z^2) \& 0x3$  then
20:     Manage_colCnt_remainder(sum, sum1, buff_in, bias,  $C_{in} \times Z^2$ )
21:   end if
22:   apply_mac(sum, sum1,  $C_{out}$ ,  $I_{out}$ )
23: end if
```

---

# Convolution layers



$$N_p = 16 \text{ iterations} \Leftrightarrow K = 32$$



$$N_p = 32 \text{ iterations} \Leftrightarrow K = 64$$

- With  $N_p$  the number of visible patterns corresponding to iterations of first for loop:

$$\text{rowCount: } K = C_{out} = 2 \times N_p$$

# Convolution layers

- Allows to extract:  $H_{out}$ ,  $K$  and  $Z$
- Allows to deduce:  $S$  and  $P$
- Russian dolls effect to get  $K$  and  $Z$
- $rowCnt$ :  $K = C_{out} = 2 \times N_p$
- $colCnt$ :  $Z = \sqrt{\frac{(4 \times N_p)}{C_{in}}}$

---

## Algorithm 1 Matrix-product for convolution (GeMM)

---

**Input:**  $buff_{in}$ ,  $ker$ ,  $C_{out}$ ,  $C_{in} \times Z^2$ ,  $I_{out}$ ,  $bias$  (bias tensor)

**Output:** Partly filled  $I_{out}$

```

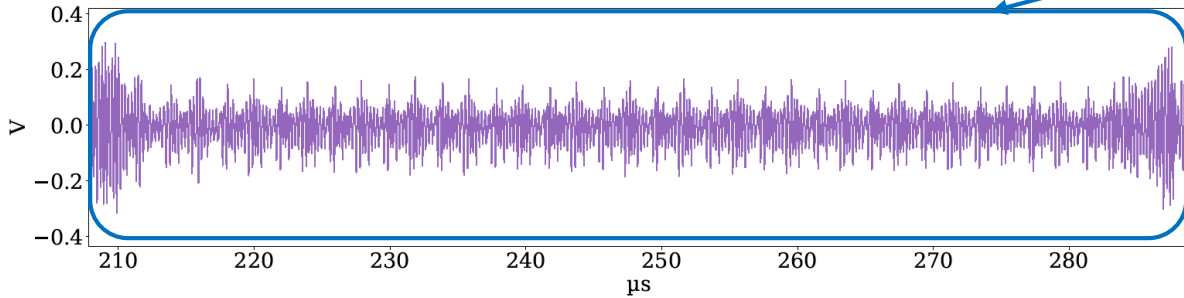
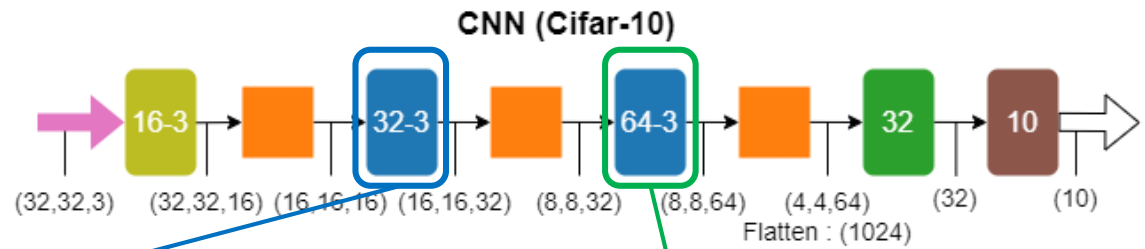
1: rowCnt ← Cout >> 1                                     ▷ Set rowCnt = K/2
2: for rowCnt > 0, rowCnt - -1 do                          ▷ Iterate over K/2
3:   sum, sum1, sum2, sum3 = init sum(bias, Cin × Z2)
4:   colCnt ← (Cin × Z2) >> 2                             ▷ Set colCnt
5:   for colCnt > 0, colCnt - -1 do
6:     simd_mac(sum, sum1, sum2, sum3, buffin, ker)
7:   end for
8:   if (Cin × Z2) & 0x3 then
9:     Manage_colCnt_remainder(sum, sum1, sum2, sum3, buffin, bias, Cin × Z2)
10:  end if
11:  apply_mac(sum, sum1, sum2, sum3, Cout, Iout)
12: end for
13: if Cout & 0x3 then                                     ▷ Check presence of remainders
14:   sum, sum1 = init sum(bias, Cin × Z2)
15:   colCnt ← (Cin × Z2) >> 2                             ▷ Set colCnt
16:   for colCnt > 0, colCnt - -1 do
17:     simd_mac(sum, sum1, buffin, ker)
18:   end for
19:   if (Cin × Z2) & 0x3 then
20:     Manage_colCnt_remainder(sum, sum1, buffin, bias, Cin × Z2)
21:   end if
22:   apply_mac(sum, sum1, Cout, Iout)
23: end if

```

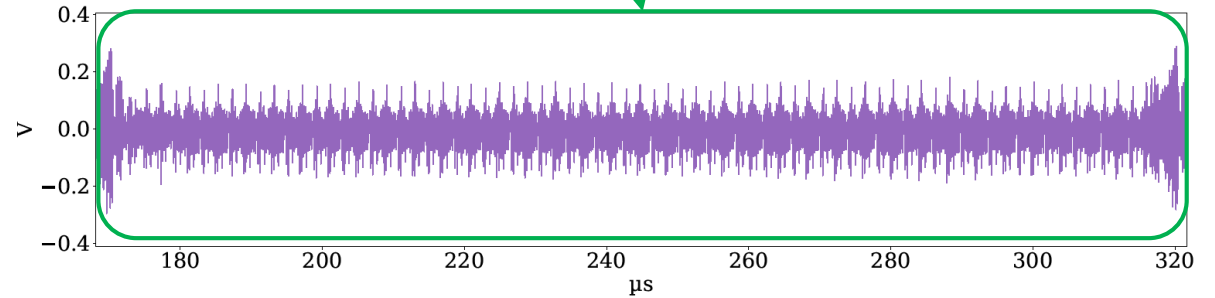
---



# Convolution layers



$N_p = 36$  iterations  $\Leftrightarrow Z = 3$  with  $C_{in} = 16$



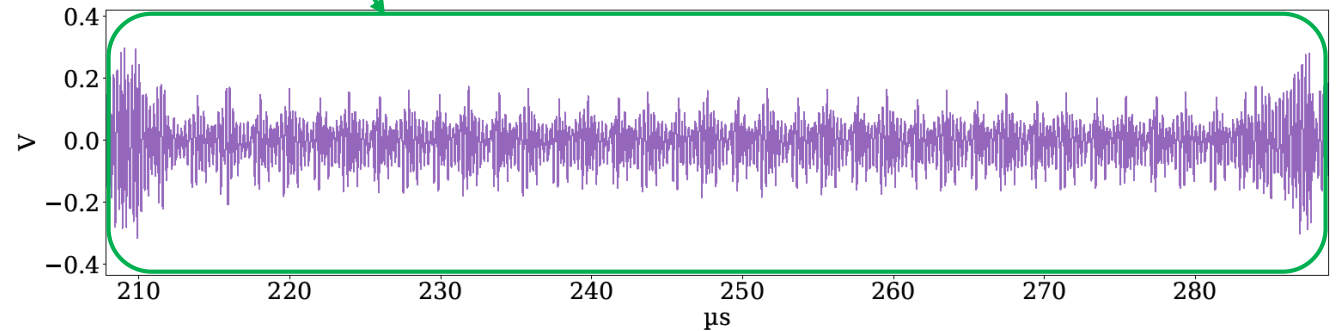
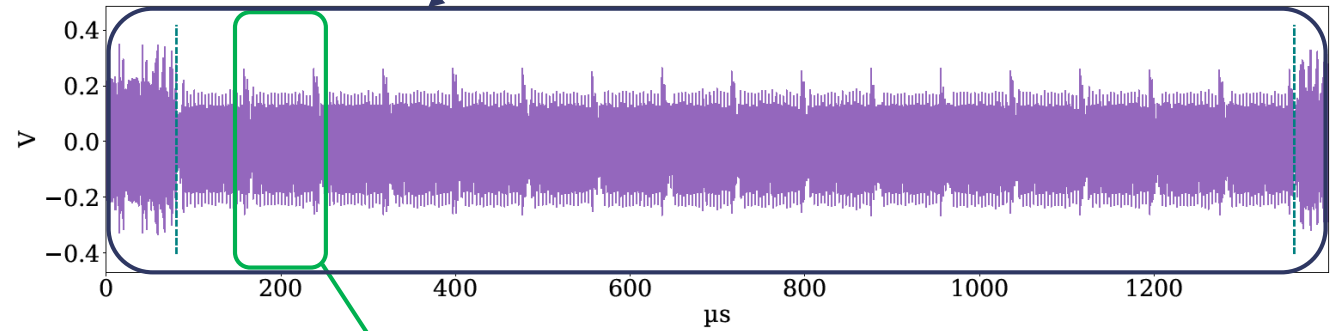
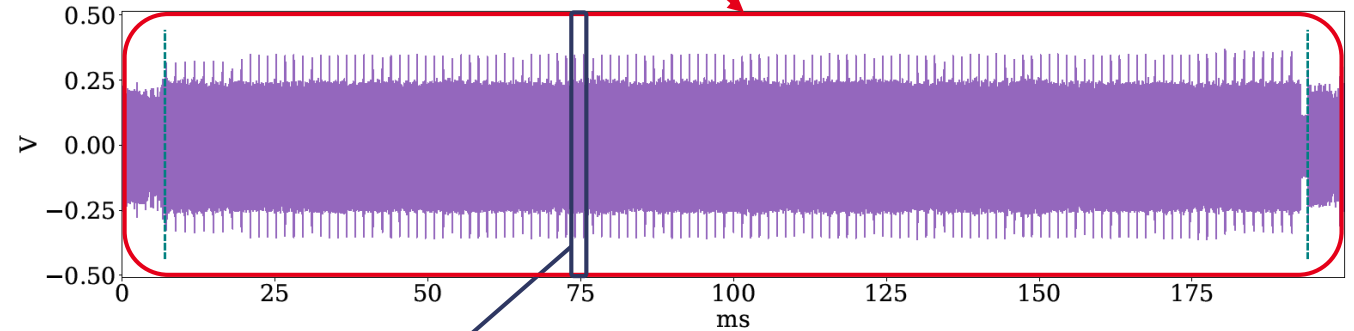
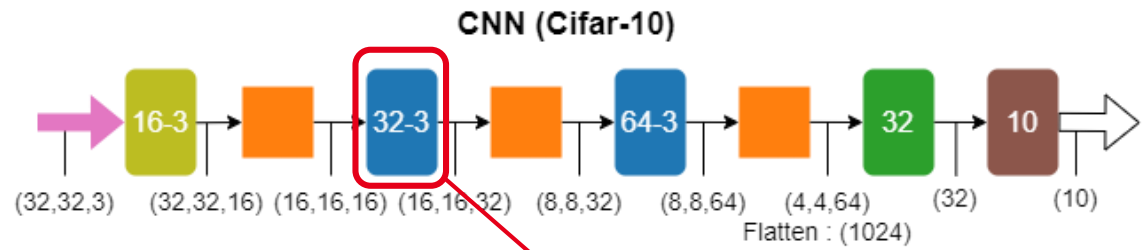
$N_p = 72$  iterations  $\Leftrightarrow Z = 3$  with  $C_{in} = 32$

- With  $N_p$  the number of visible patterns corresponding to iterations of second for loop:

$$colCnt: Z = \sqrt{\frac{(4 \times N_p)}{C_{in}}}$$

# Convolution layers

- Look at for loops iterations
- Russian Dolls effect
- Allows to extract:
  - $H_{out}$
  - $K$
  - $Z$





# Convolution layers

- Look at for loops iterations
- Russian Dolls effect
- Allows to extract:

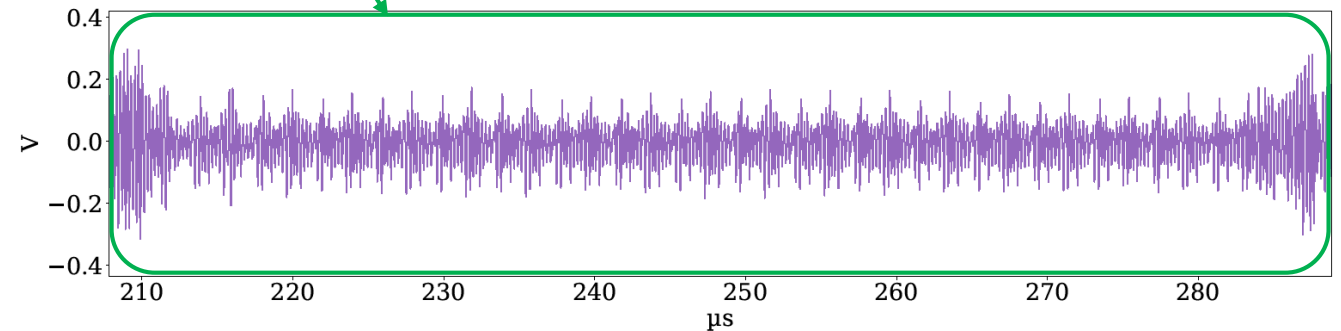
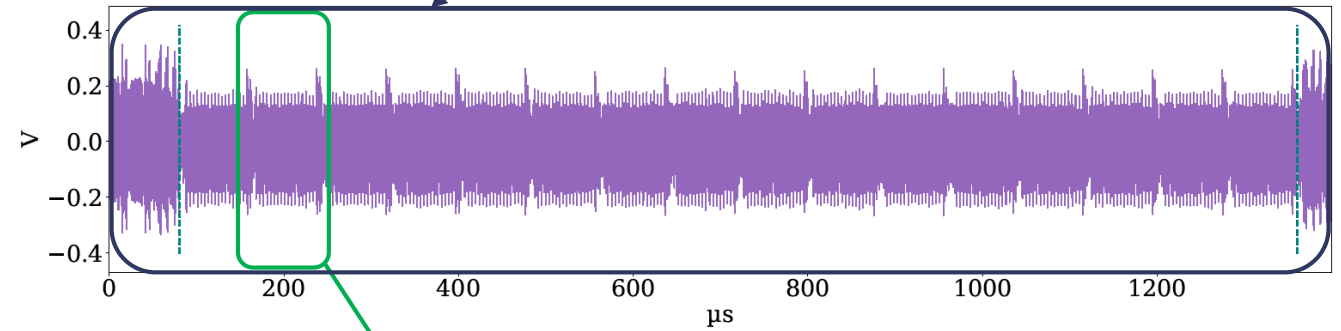
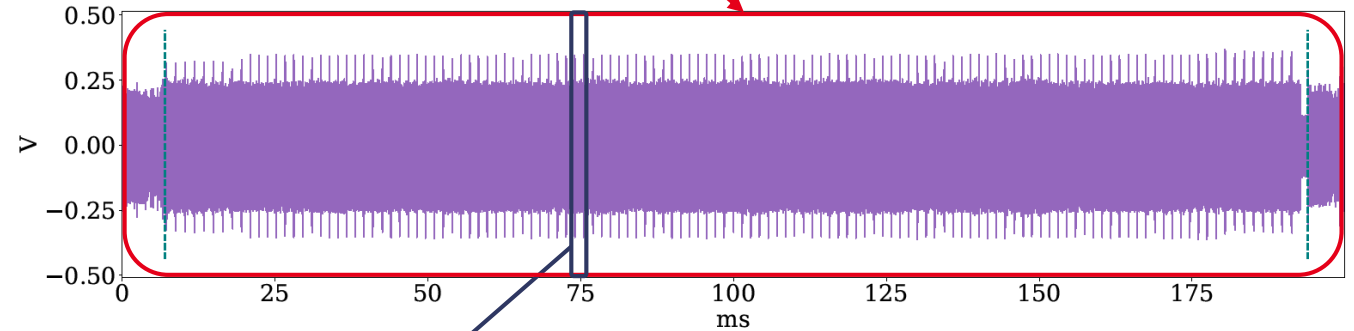
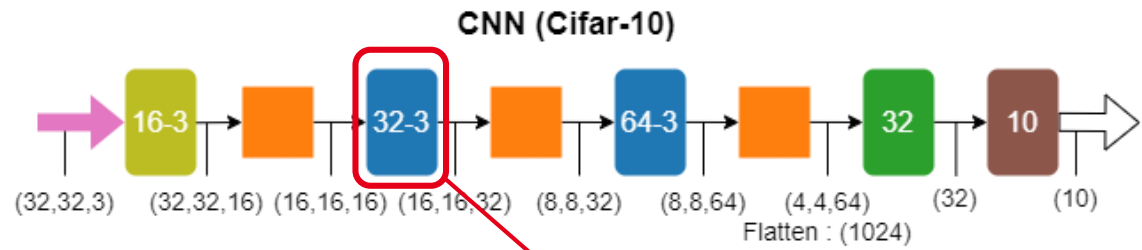
- $H_{out}$

- $K$

- $Z$

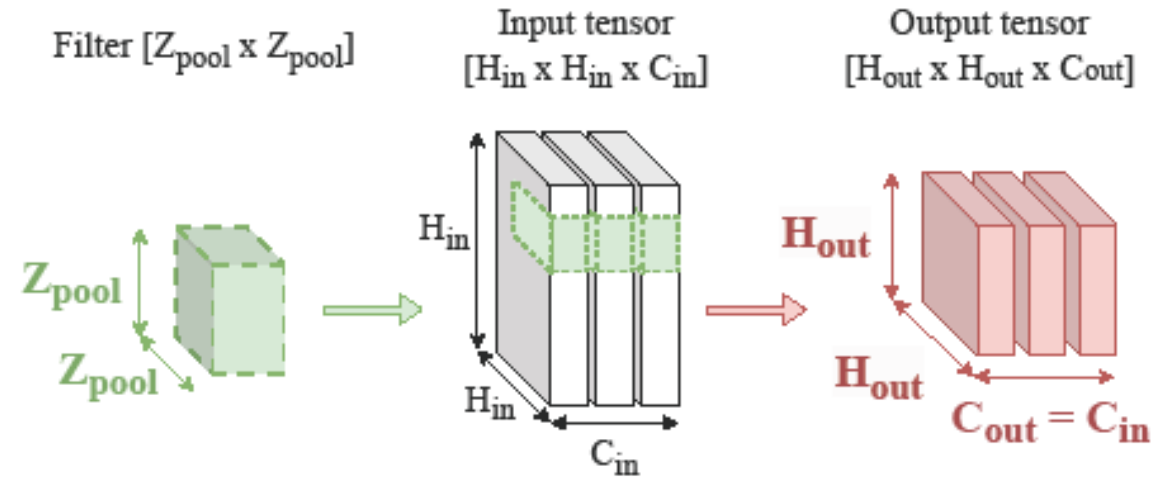
- Then deducing:

- $S$  and  $P$  as  $P \leq Z$



# MaxPool layers

- Allows to extract  $H_{out}$  and deduce  $Z_{pool}$
- 2 distinct parts



$$H_{out} = \frac{H_{in}}{Z_{pool}}$$

---

## Algorithm 1 MaxPool - arm\_maxpool\_q7\_HWC

---

**Input:** Input tensor  $I_{in}$  of size  $H_{in}^2 \cdot C_{in}$ , Output tensor  $I_{out}$  of size  $H_{out}^2 \cdot C_{out}$ ,  $P$ ,  $S$ ,  $H_{ker}$

**Output:** Filled  $I_{out}$

```

1: for  $i_y \leftarrow 0, i_y < H_{in}, i_y ++1$  do
2:   for  $i_x \leftarrow 0, i_x < H_{out}, i_x ++1$  do
3:      $win_{start}, win_{stop} \leftarrow set\_window(i_y, i_x, I_{in}, H_{ker}, P, S)$ 
4:      $compare\_and\_replace\_if\_larger(win_{start}, win_{stop}, i_y, i_x, I_{in})$ 
5:   end for
6: end for
7: for  $i_y \leftarrow 0, i_y < H_{out}, i_y ++1$  do
8:    $row_{start}, row_{stop} \leftarrow set\_rows(i_y, I_{in}, I_{out}, H_{ker}, P, S)$ 
9:    $compare\_replace\_then\_apply(row_{start}, row_{stop}, I_{in}, I_{out})$ 
10: end for

```

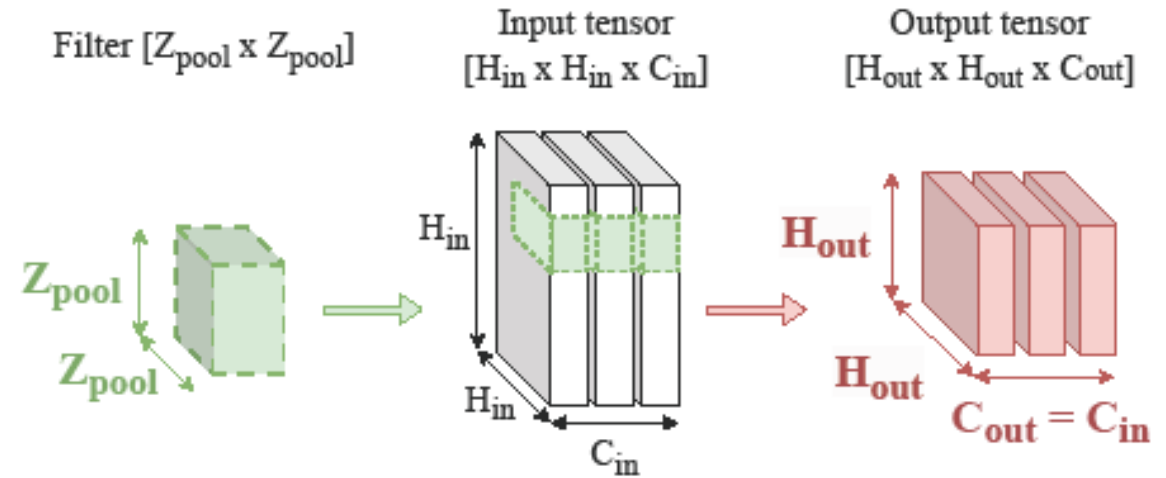
▷ Pooling along x-axis

▷ Pooling along y-axis, iterates directly over  $H_{out}$

---

# MaxPool layers

- Allows to extract  $H_{out}$  and deduce  $Z_{pool}$
- 2 distinct parts



$$H_{out} = \frac{H_{in}}{Z_{pool}}$$

---

## Algorithm 1 MaxPool - arm\_maxpool\_q7\_HWC

---

**Input:** Input tensor  $I_{in}$  of size  $H_{in}^2 \cdot C_{in}$ , Output tensor  $I_{out}$  of size  $H_{out}^2 \cdot C_{out}$ ,  $P$ ,  $S$ ,  $H_{ker}$

**Output:** Filled  $I_{out}$

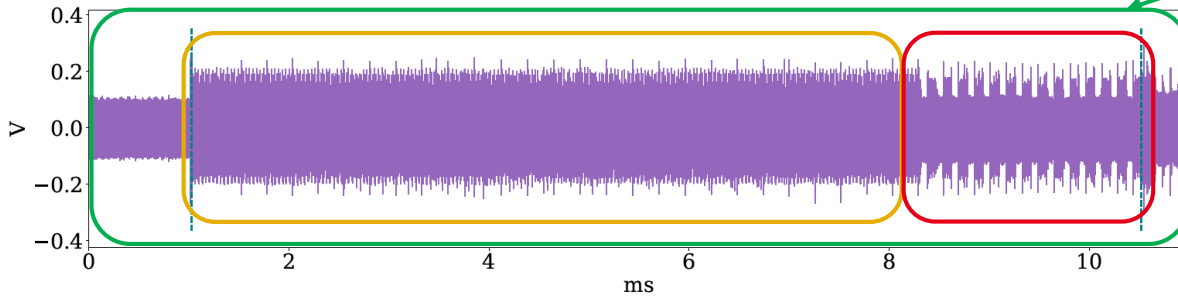
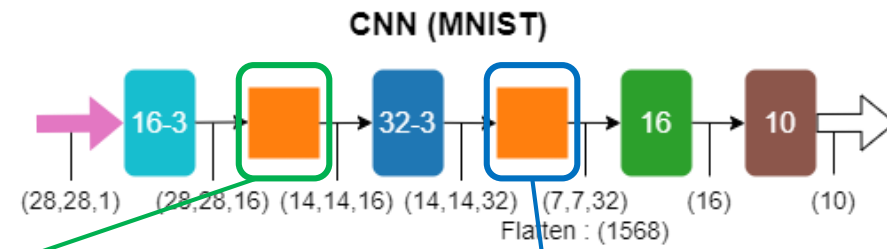
```

1: for  $i_y \leftarrow 0, i_y < H_{in}, i_y ++1$  do ▷ Pooling along x-axis
2:   for  $i_x \leftarrow 0, i_x < H_{out}, i_x ++1$  do
3:      $win_{start}, win_{stop} \leftarrow set\_window(i_y, i_x, I_{in}, H_{ker}, P, S)$ 
4:      $compare\_and\_replace\_if\_larger(win_{start}, win_{stop}, i_y, i_x, I_{in})$ 
5:   end for
6: end for
7: for  $i_y \leftarrow 0, i_y < H_{out}, i_y ++1$  do ▷ Pooling along y-axis, iterates directly over  $H_{out}$ 
8:    $row_{start}, row_{stop} \leftarrow set\_rows(i_y, I_{in}, I_{out}, H_{ker}, P, S)$ 
9:    $compare\_replace\_then\_apply(row_{start}, row_{stop}, I_{in}, I_{out})$ 
10: end for

```

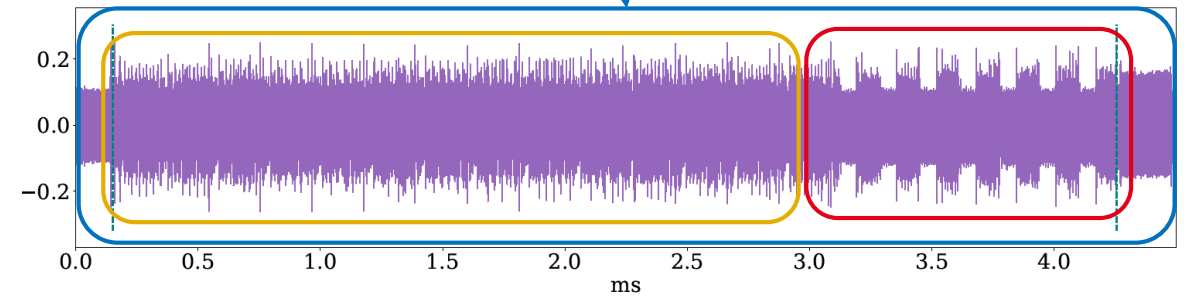
---

# MaxPool layers



$$N_p = 14 \text{ iterations} \Leftrightarrow H_{out} = 14$$

- Look at for loops iterations
- Two distinct parts: one easier to analyse



$$N_p = 7 \text{ iterations} \Leftrightarrow H_{out} = 7$$

- Allows to extract:  $H_{out}$
- Allows to deduce:  $Z_{pool}$

# Dense layers

- Neurons managed by groups of 4
- Patterns length depends on input shape

---

## Algorithm 1 Dense layer - `arm_fully_connected_q7_opt`

---

**Input:** Input vector  $I_{in}$  of size  $H_{in}$ , weight vector  $ker$  of size  $N_e$ , bias matrix  $bias$  of size  $N_e$ , output vector  $I_{out}$  of size  $H_{out}$ ,  $P$ ,  $S$ ,  $H_{ker}$

**Output:** Filled  $I_{out}$

```
1: rowCnt ←  $N_e \gg 2$  ▷ Nb. neurons divided by 4
2: for rowCnt > 0, rowCnt - -1 do ▷ Iterate directly over  $N_e/4$ 
3:   sum, sum1, sum2, sum3 = init_sum_with_bias(bias, rowCnt)
4:   colCnt ←  $H_{in} \gg 2$ 
5:   for colCnt > 0, colCnt - -1 do
6:     simd_mac(sum, sum1, sum2, sum3, ker, colCnt)
7:   end for
8:   apply_mac(sum, sum1, sum2, sum3, rowCnt, Iout)
9: end for
10: rowCnt ←  $N_e \& 0x3$  ▷ Manage remainders if any
11: for rowCnt > 0, rowCnt - -1 do
12:   sum = init_sum_with_bias(bias, rowCnt)
13:   colCnt ←  $H_{in} \gg 2$ 
14:   for colCnt > 0, colCnt - -1 do
15:     mac(sum, ker, colCnt)
16:   end for
17:   apply_mac(sum, rowCnt, Iout)
18: end for
```

---

# Dense layers

- Neurons managed by groups of 4
- Patterns length depends on input shape
- Allows to extract:

- $N_e = 4 \times N_p$

---

## Algorithm 1 Dense layer - `arm_fully_connected_q7_opt`

---

**Input:** Input vector  $I_{in}$  of size  $H_{in}$ , weight vector  $ker$  of size  $N_e$ , bias matrix  $bias$  of size  $N_e$ , output vector  $I_{out}$  of size  $H_{out}$ ,  $P$ ,  $S$ ,  $H_{ker}$

**Output:** Filled  $I_{out}$

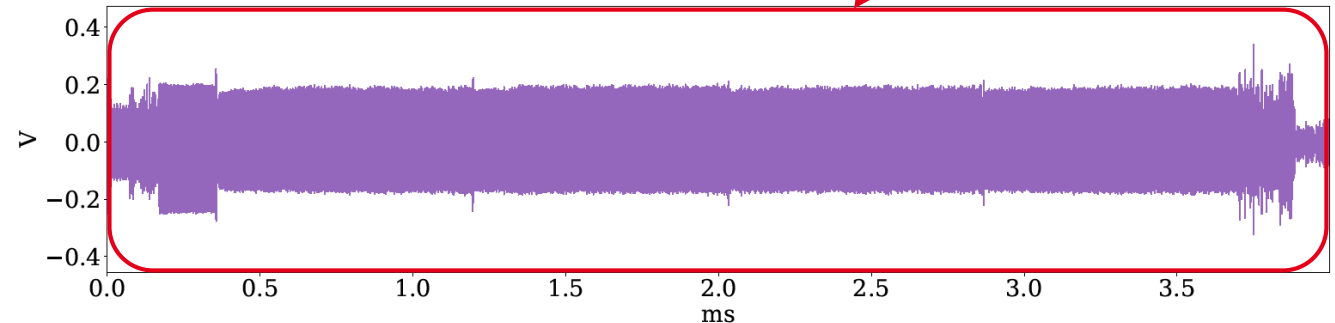
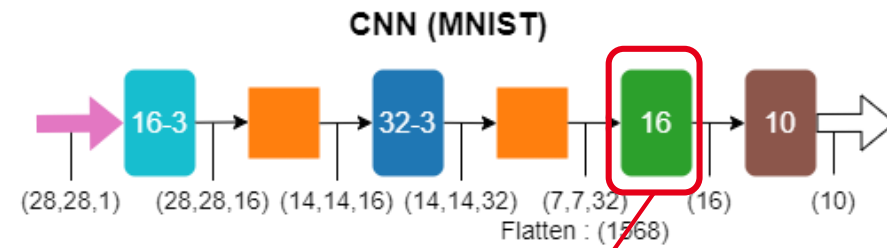
```
1: rowCnt ←  $N_e \gg 2$  ▷ Nb. neurons divided by 4
2: for rowCnt > 0, rowCnt - -1 do ▷ Iterate directly over  $N_e/4$ 
3:   sum, sum1, sum2, sum3 = init_sum_with_bias(bias, rowCnt)
4:   colCnt ←  $H_{in} \gg 2$ 
5:   for colCnt > 0, colCnt - -1 do
6:     simd_mac(sum, sum1, sum2, sum3, ker, colCnt)
7:   end for
8:   apply_mac(sum, sum1, sum2, sum3, rowCnt, Iout)
9: end for
10: rowCnt ←  $N_e \& 0x3$  ▷ Manage remainders if any
11: for rowCnt > 0, rowCnt - -1 do
12:   sum = init_sum_with_bias(bias, rowCnt)
13:   colCnt ←  $H_{in} \gg 2$ 
14:   for colCnt > 0, colCnt - -1 do
15:     mac(sum, ker, colCnt)
16:   end for
17:   apply_mac(sum, rowCnt, Iout)
18: end for
```

---



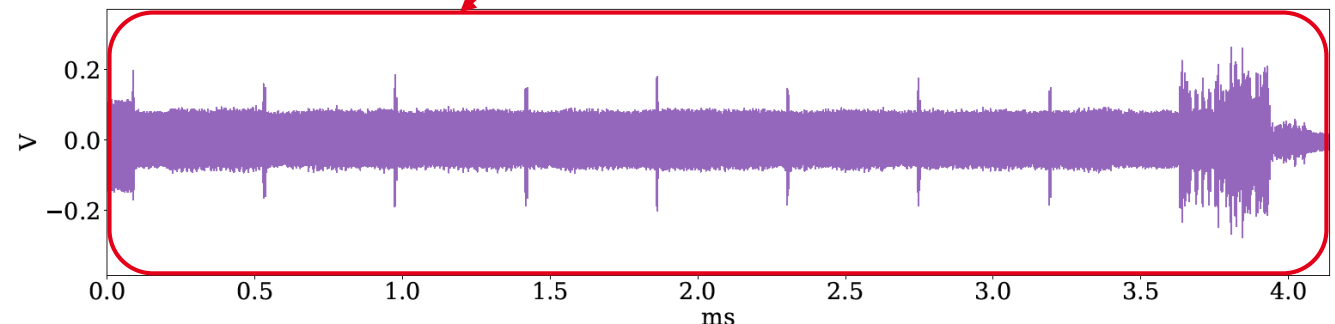
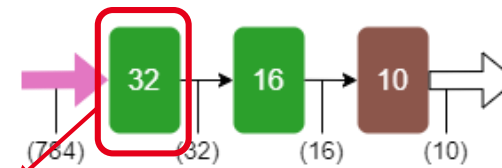
# Dense layers

- Neuron managed by groups of 4
- Patterns length depends on input shape
- Allows to extract:
  - $N_e = 4 \times N_p$
- Usual cases when  $N_e \% 4 = 0$



$$N_p = 4 \text{ iterations} \Leftrightarrow N_e = 16$$

**MLP (MNIST)**



$$N_p = 8 \text{ iterations} \Leftrightarrow N_e = 32$$

# Dense layers

- Neuron managed by groups of 4
- Patterns length depends on input shape
- Allows to extract:
  - $N_e = 4 \times N_p$
- Usual cases when  $N_e \% 4 = 0$
- Special cases when  $N_e \% 4 \neq 0$

---

## Algorithm 1 Dense layer - arm\_fully\_connected\_q7\_opt

---

**Input:** Input vector  $I_{in}$  of size  $H_{in}$ , weight vector  $ker$  of size  $N_e$ , bias matrix  $bias$  of size  $N_e$ , output vector  $I_{out}$  of size  $H_{out}$ ,  $P$ ,  $S$ ,  $H_{ker}$

**Output:** Filled  $I_{out}$

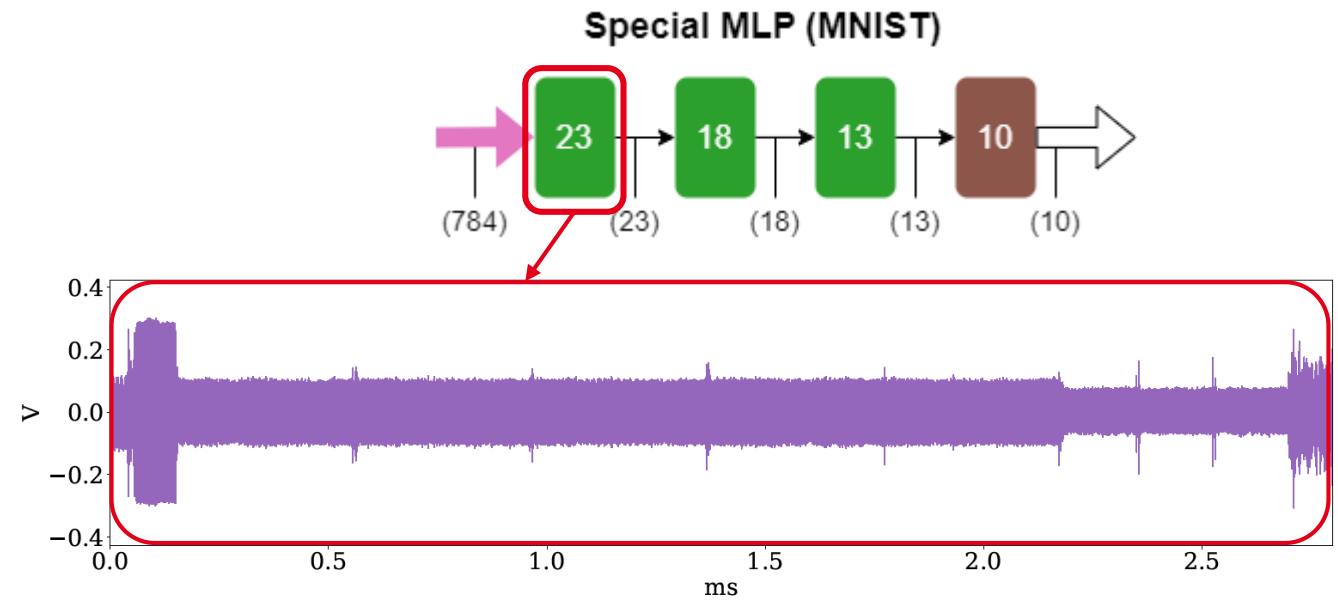
```
1: rowCnt ←  $N_e \gg 2$  ▷ Nb. neurons divided by 4  
2: for rowCnt > 0, rowCnt - -1 do ▷ Iterate directly over  $N_e/4$   
3:   sum, sum1, sum2, sum3 = init_sum_with_bias(bias, rowCnt)  
4:   colCnt ←  $H_{in} \gg 2$   
5:   for colCnt > 0, colCnt - -1 do  
6:     simd_mac(sum, sum1, sum2, sum3, ker, colCnt)  
7:   end for  
8:   apply_mac(sum, sum1, sum2, sum3, rowCnt, I_out)  
9: end for  
10: rowCnt ←  $N_e \& 0x3$  ▷ Manage remainders if any  
11: for rowCnt > 0, rowCnt - -1 do  
12:   sum = init_sum_with_bias(bias, rowCnt)  
13:   colCnt ←  $H_{in} \gg 2$   
14:   for colCnt > 0, colCnt - -1 do  
15:     mac(sum, ker, colCnt)  
16:   end for  
17:   apply_mac(sum, rowCnt, I_out)  
18: end for
```

---



# Dense layers

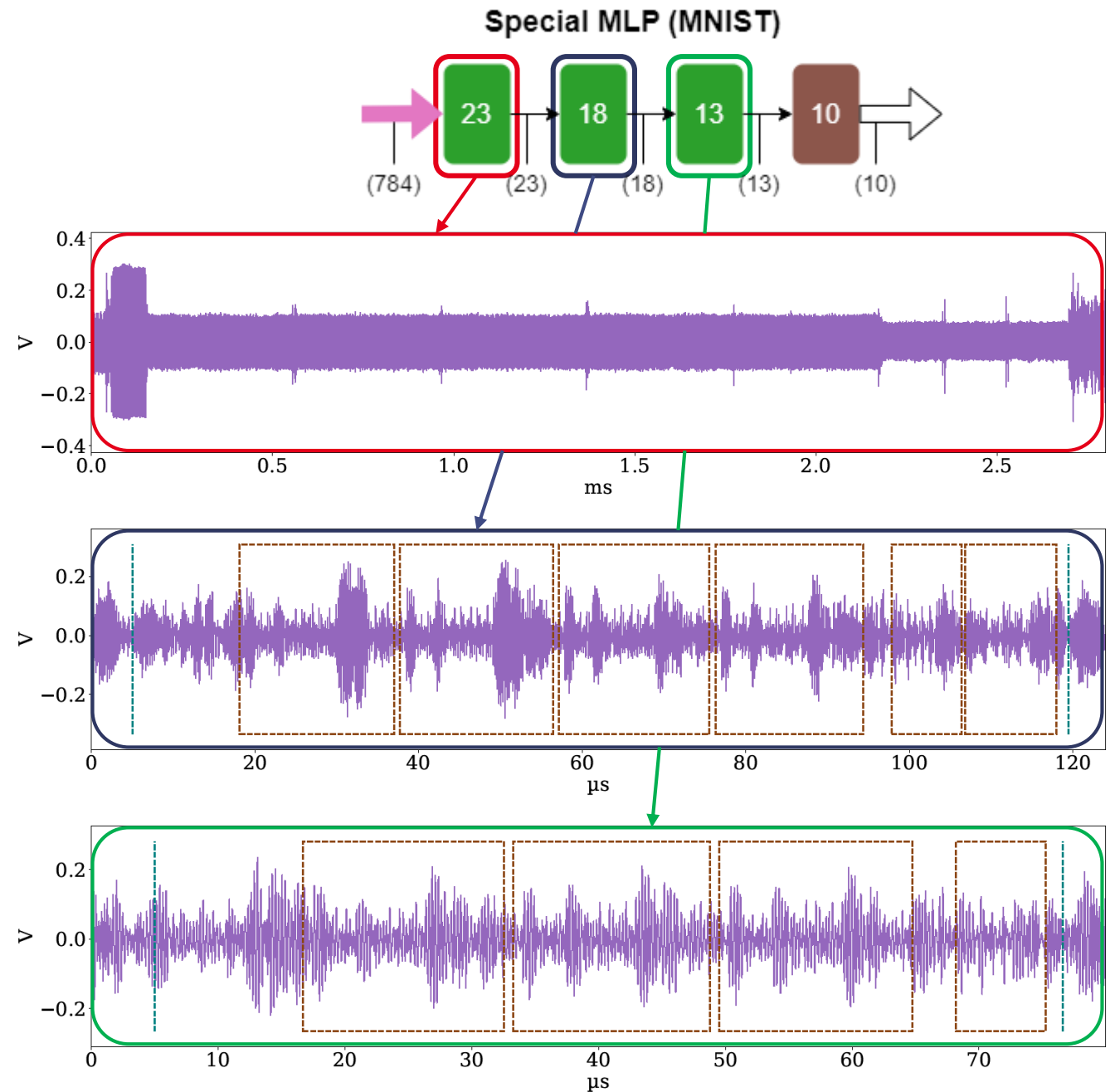
- Neuron managed by groups of 4
- Patterns length depends on input shape
- Allows to extract:
  - $N_e = 4 \times N_p$
- Usual cases when  $N_e \% 4 = 0$
- Special cases when  $N_e \% 4 \neq 0$



$$N_p = 5 + 3 \text{ iterations} \Leftrightarrow N_e = 4 \times 5 + 3 = 23$$

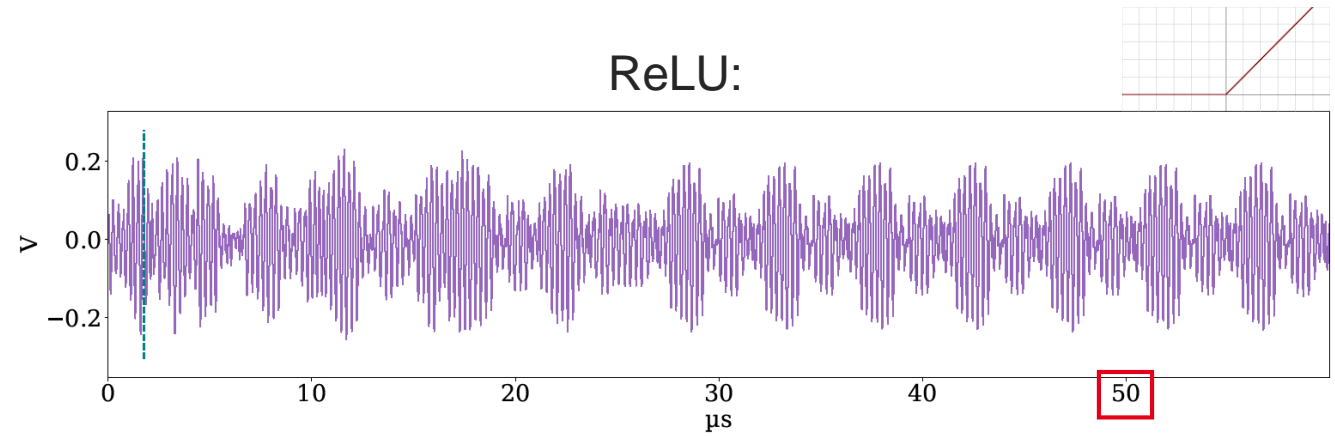
# Dense layers

- Neuron managed by groups of 4
- Patterns length depends on input shape
- Allows to extract:
  - $N_e = 4 \times N_p$
- Usual cases when  $N_e \% 4 = 0$
- Special cases when  $N_e \% 4 \neq 0$
- Illustration with additional triggers



# Activation layers

- Distinguish ReLU from Tanh and Sigmoid
- Can use duration of activation layers
- EM patterns give additional hints

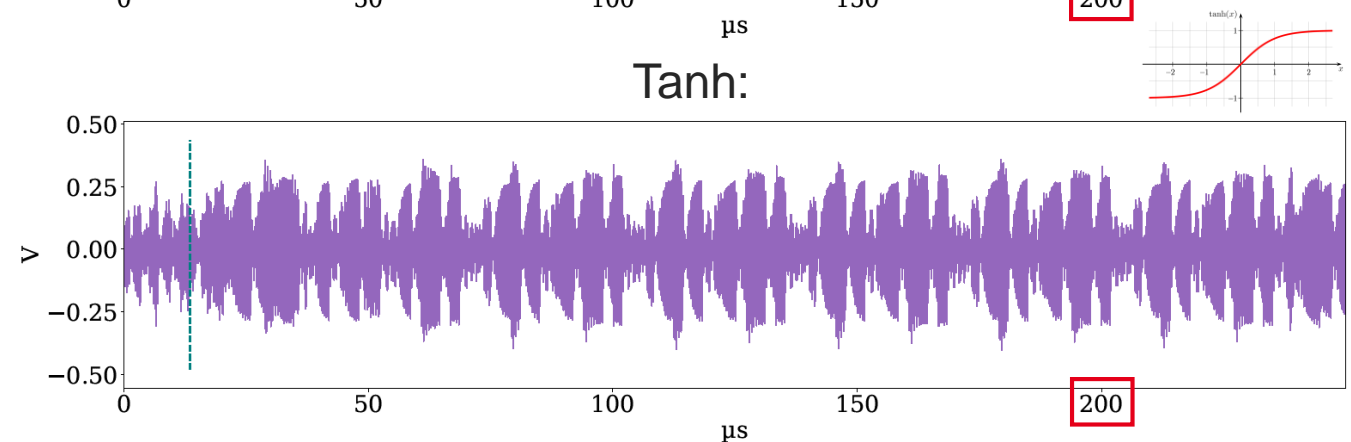
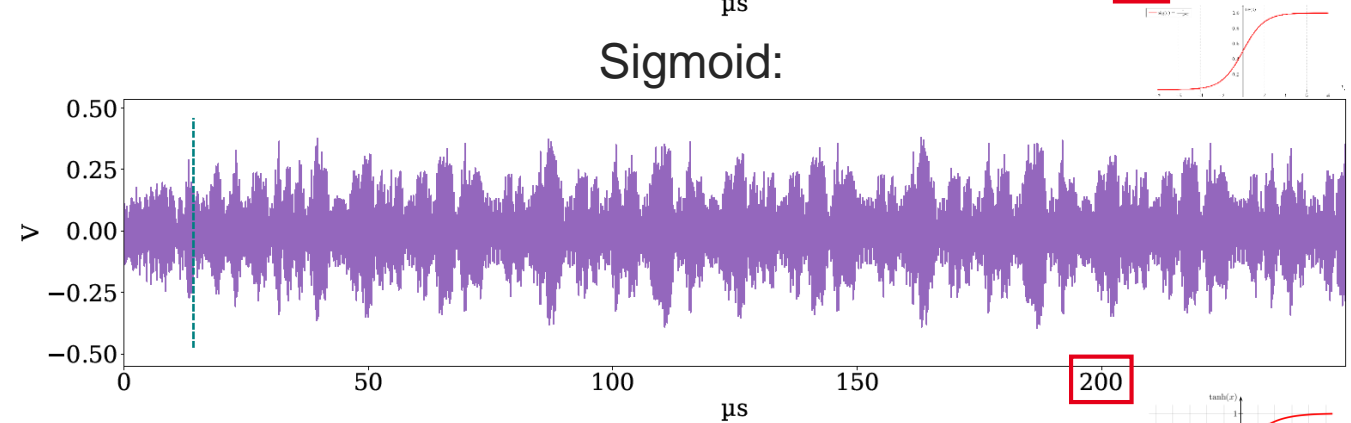
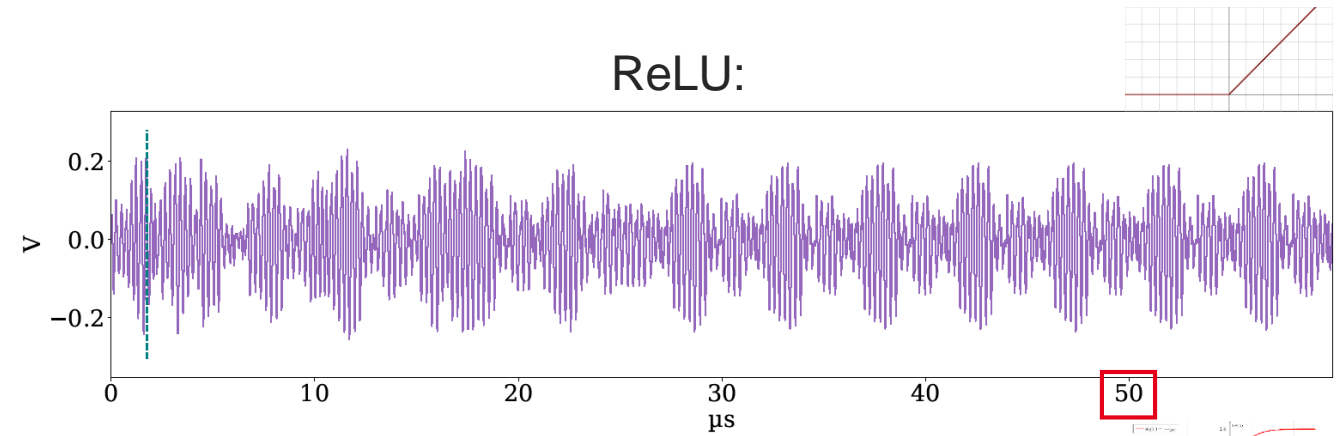


Layer	ReLU	Sigmoid	Tanh
0	1,84 ms	8,24 ms	8,24 ms
1	0,92 ms	3,21 ms	3,58 ms
2	7,60 μs	15,6 μs	20,64 μs

# Activation layers

- Distinguish ReLU from Tanh and Sigmoid
- Can use duration of activation layers
- EM patterns give additional hints

Layer	ReLU	Sigmoid	Tanh
0	1,84 ms	8,24 ms	8,24 ms
1	0,92 ms	3,21 ms	3,58 ms
2	7,60 $\mu$ s	15,6 $\mu$ s	20,64 $\mu$ s



# Outline

- Introduction
- Challenges
- Scope and threat model
- Model analysis
- Layer analysis
- **Discussion & Perspectives**



# Discussions and Perspectives

- Several works studying architecture extraction
- Attack surface for such a threat is significantly extended by SCA
- No need of complex exploitation methods (e.g., heavy supervised profiling step)
- Important step for parameters extraction [1 to 7]
- Some more complex cases allowing to approximate hyper-parameter (could be handled with pattern detection tools)
- Critical need to protect architecture from extraction [11, 17, 18]
- Hard challenges to develop efficient protections 32-bit microcontrollers



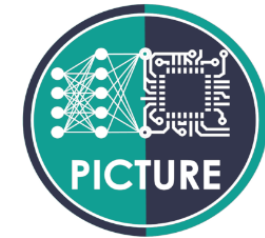
# References

1. Carlini, N. et al. "Cryptanalytic extraction of neural network models." Annual International Cryptology Conference. 2020.
2. Canales-Martínez, I., et al. "Polynomial Time Cryptanalytic Extraction of Neural Network Models." arXiv 2023.
3. Jagielski, M., et al. "High accuracy and high fidelity extraction of neural networks." 29th USENIX security symposium (USENIX Security 20). 2020.
4. Rakin, A. et al. "Deepsteal: Advanced model extractions leveraging efficient weight stealing in memories." IEEE SP, 2022.
5. Hector, K., et al. "Fault Injection and Safe-Error Attack for Extraction of Embedded Neural Network Models." SECAI (ESORICS Workshop) 2023.
6. Joud, R, et al. "A Practical Introduction to Side-Channel Extraction of Deep Neural Network Parameters." CARDIS 2022.
7. Gongye, C. et al. "Reverse-engineering deep neural networks using floating-point timing side-channels." ACM/IEEE DAC, 2020.
8. Duddu, V., et al. "Stealing neural networks via timing side channels." arXiv 2018.
9. Yu, H. et al. "Deepem: Deep neural networks model recovery through em side-channel information leakage." IEEE HOST 2020.
10. Yli-Mäyry, V., et al. "Extraction of binarized neural network architecture and secret parameters using side-channel information." IEEE ISCAS. 2021.
11. Luo, Y., et al. "NNReArch: A Tensor Program Scheduling Framework Against Neural Network Architecture Reverse Engineering." IEEE FCCM. 2022.
12. Chmielewski, L. et al. "On reverse engineering neural network implementation on gpu." ACNS 2021, Satellite Workshops, AIBlock, AIHWS, AIoTS, CIMSS, Cloud S&P, SCI, SecMT, and SiMLA, 2021.
13. Batina, L., et al. "CSI-NN: Reverse engineering of neural network architectures through electromagnetic side channel." USENIX Security. 2019.
14. Xiang, Y., et al. "Open dnn box by power side-channel attack." IEEE Transactions on Circuits and Systems II: Express Briefs 2020
15. Ma, J.: A higher-level Neural Network library on Microcontrollers (NNoM) (2020).
16. Lai, L. et al. "Cmsis-nn: Efficient neural network kernels for arm cortex-m cpus." arXiv. 2018.
17. Chabanne, H., et al. "Telepathic headache: Mitigating cache side-channel attacks on convolutional neural networks." ACNS, 2021.
18. Li, J., et al. "Neurofuscator: A full-stack obfuscation tool to mitigate neural architecture stealing." IEEE HOST. 2021.



# Acknowledgements

- This work is supported by (CEA-Leti) the European project InSecTT (ECSEL JU 876038) and by the French ANR in the framework of the Investissements d'avenir program (ANR-10-AIRT-05, irtnanoelec); and (Mines Saint-Etienne) by the French program ANR PICTURE (AAPG2020).
- This work benefited from the French Jean Zay supercomputer with the AI dynamic access program.



# Thank you for your attention

Contacts:

Pierre-Alain Moëllic  
[pierre-alain.moellic@cea.fr](mailto:pierre-alain.moellic@cea.fr)

Simon Pontié  
[Simon.pontie@cea.fr](mailto:Simon.pontie@cea.fr)

Raphaël Joud  
[raphael.joud@cea.fr](mailto:raphael.joud@cea.fr)

Jean-Baptiste Rigaud  
[rigaud@emse.fr](mailto:rigaud@emse.fr)