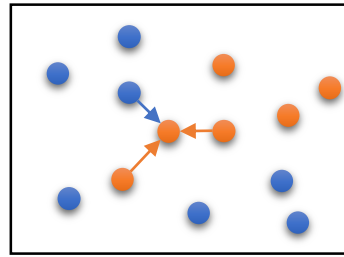


# FPGA BASED LOW LATENCY, LOW POWER STREAM PROCESSING AI

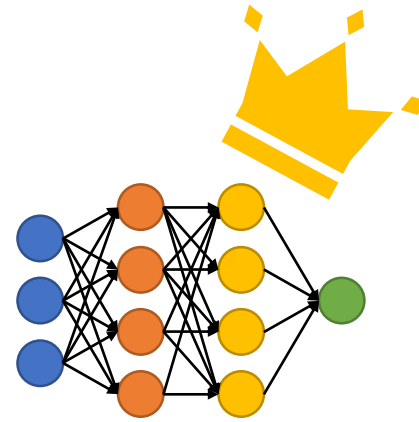
**Domenik Helms<sup>(1)</sup>, Mark Kettner<sup>(1)</sup>, Behnam Razi Perjikolaei<sup>(1)</sup>, Lukas Einhaus<sup>(2)</sup>, Christopher Ringhofer<sup>(2)</sup>, Chao Qian<sup>(2)</sup>, Gregor Schiele<sup>(2)</sup>**

**(1) OFFIS, (2) University of Duisburg-Essen**

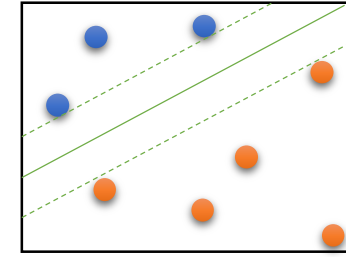
# Motivation



nearest neighbor

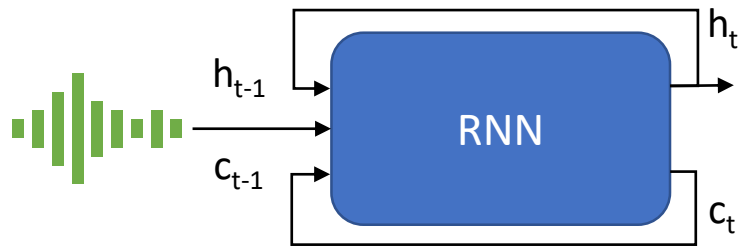


neural network

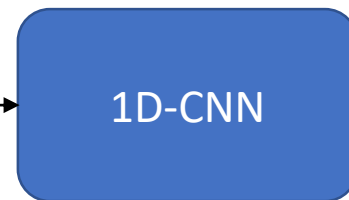
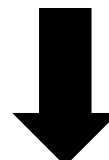


support vector machine

etc. ...

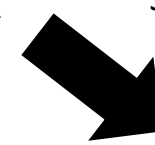


recurrent neural network



1D-CNN

“heart attack”



convolutional neural network

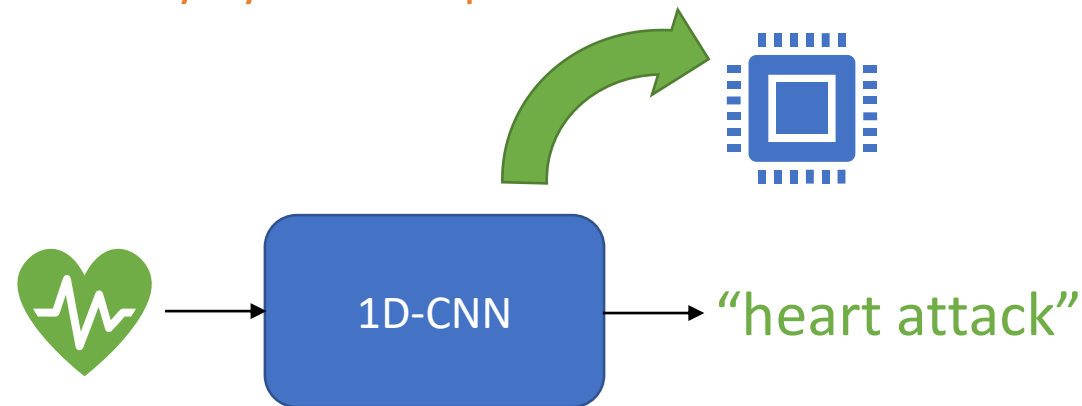
“cat”

# Motivation

- Wanted:
  - high throughput
    - ⚡ limited mainly by data access bandwidth
  - power efficiency
    - ⚡ limited mainly by data access cost
  - radiation hardness
    - ⚡ limited mainly by weak sequential cells

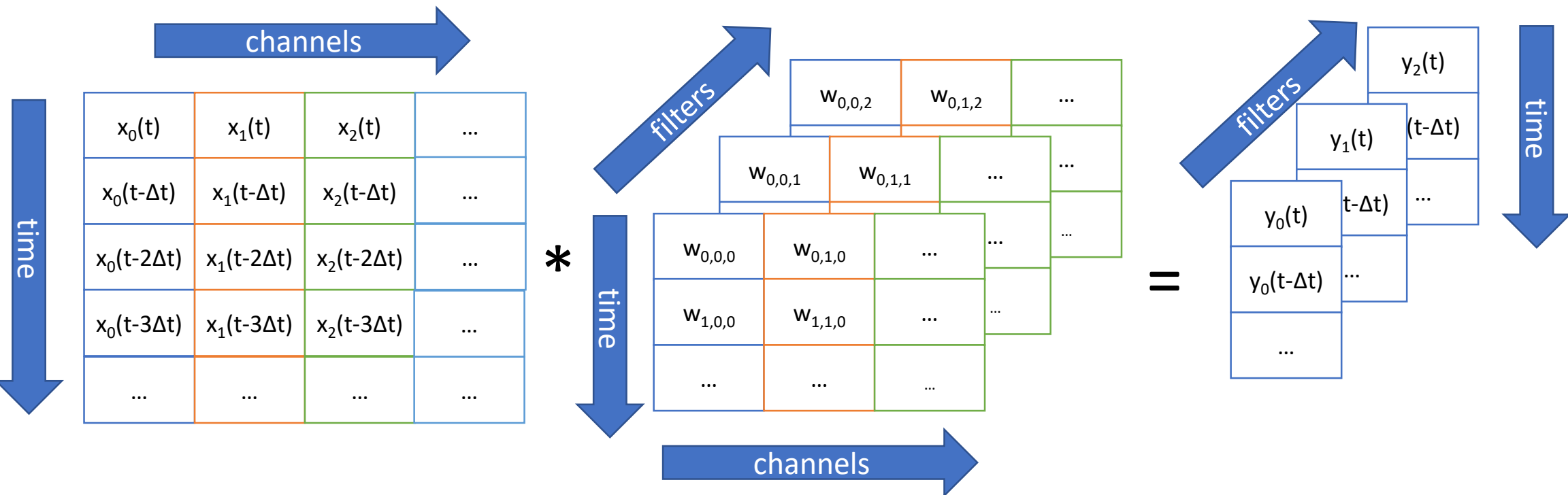
➔ limited mainly by memory!!!

do we need memory?

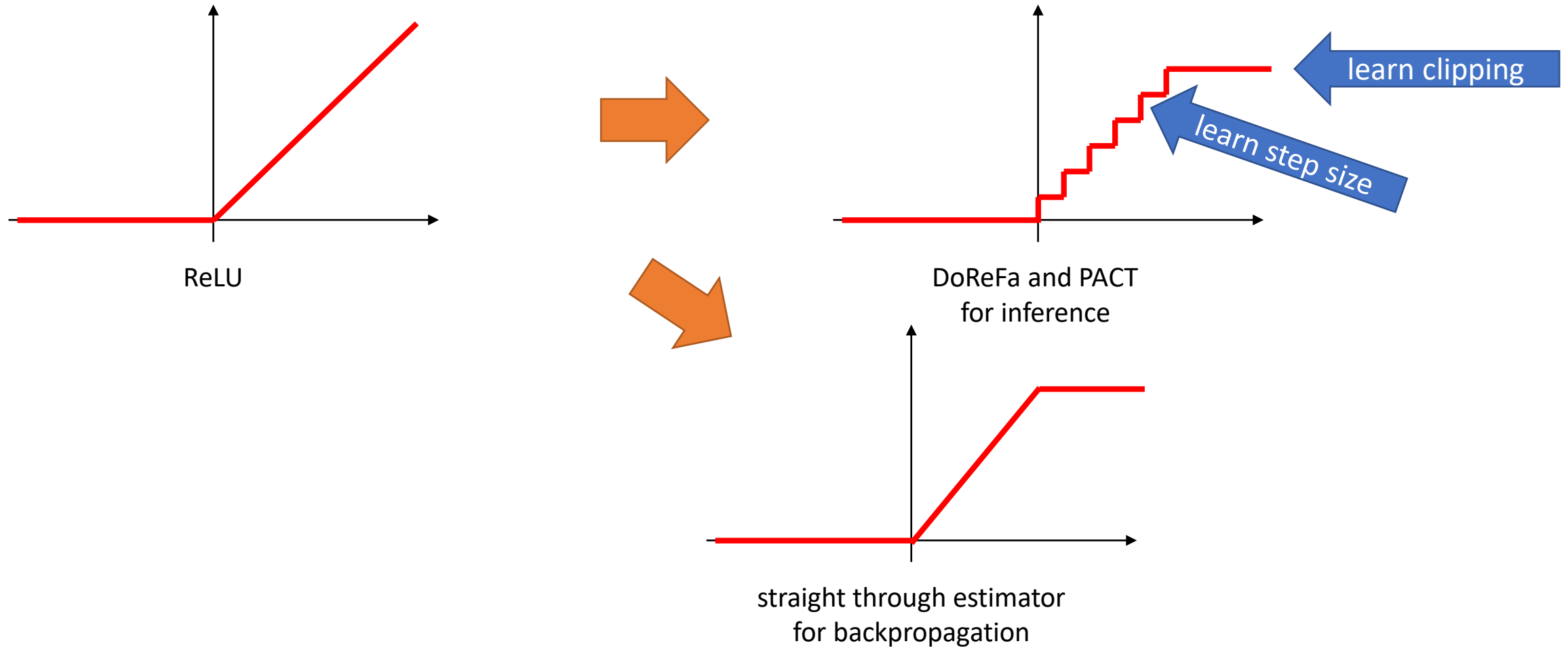


# How to build simple NN without memory?

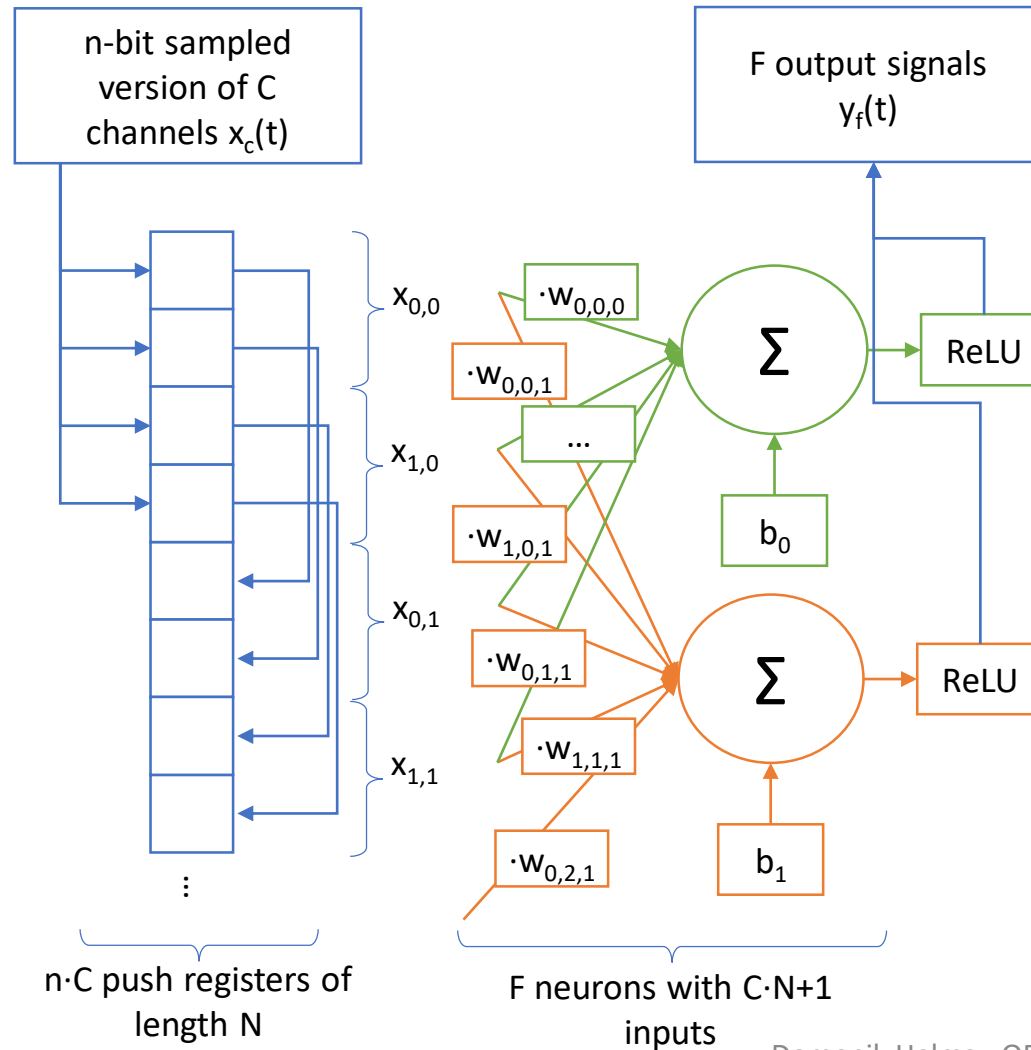
$$\text{1D-Convolution: } y_f(t) = \max \left( 0, \sum_{c=0}^{C-1} \sum_{i=0}^{N-1} w_{i,c,f} \cdot x_c(t - i \cdot \Delta t) + b_f \right)$$



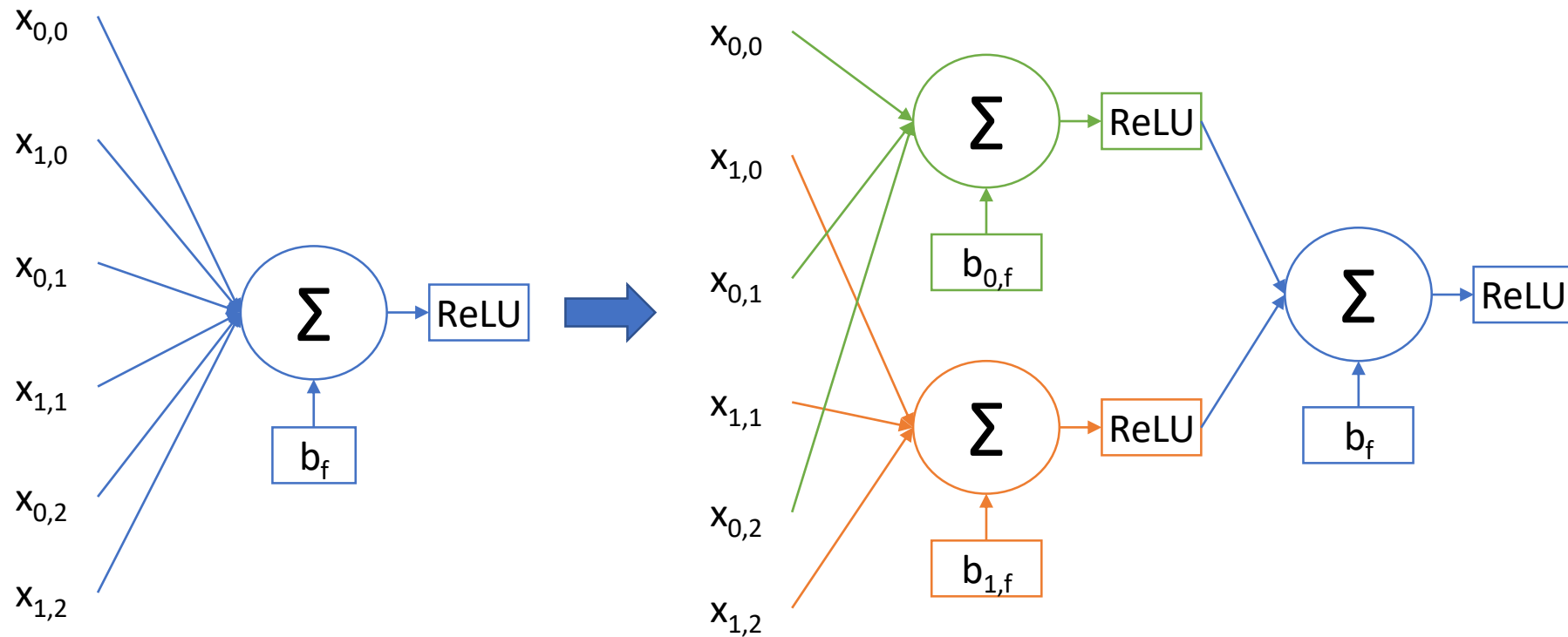
# Step 1: Quantization



# Step 2: Convolution



# Step 3: Reducing the input count



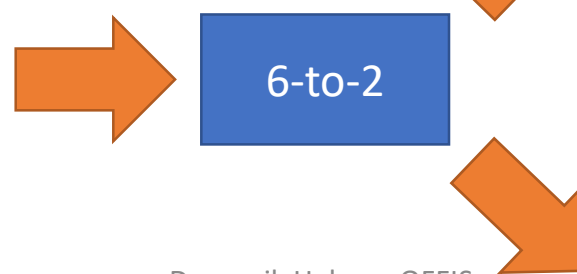
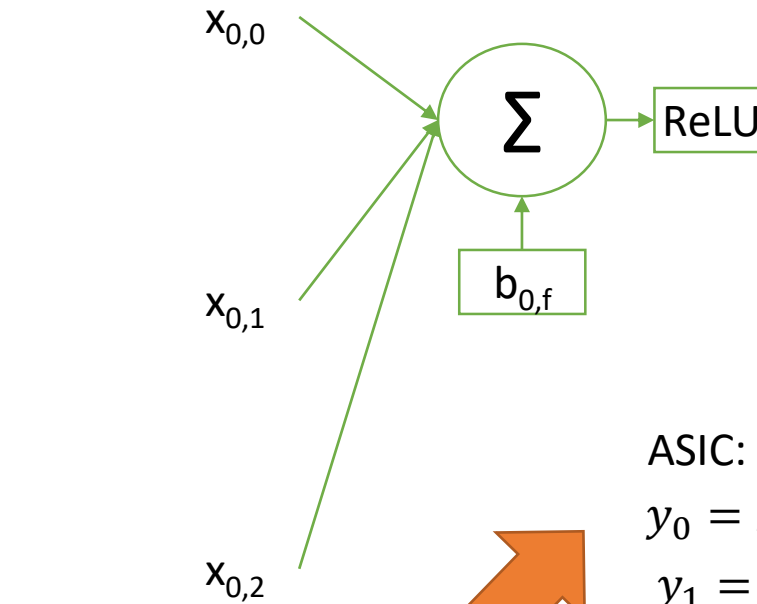
# Step 4: Precompute multiplication

$$y = \max\left(0, \sum_{i=0}^{N-1} w_i \cdot x_i + b\right)$$

$$w_0 = 0.43, w_1 = -0.24$$

$$w_2 = 0.65, b = 0.44$$

$x_0$	$x_1$	$x_2$	$y$	$ y $		$x$	$y$
0	0	0	0.44	0		0 0 0 0 0 0	0 0
0	0	1	1.09	1		0 0 0 0 0 1	0 1
0	0	2	1.74	2		0 0 0 0 1 0	1 0
0	0	3	2.39	2		0 0 0 0 1 1	1 0
0	1	0	0.2	0		0 0 0 1 0 0	0 0
0	1	1	0.85	1		0 0 0 1 0 1	0 1
⋮	⋮	⋮	⋮	⋮			
1	0	0	0.87	1		0 1 0 0 0 0	0 1
1	0	1	1.52	2		0 1 0 0 0 1	1 0
⋮	⋮	⋮	⋮	⋮			
3	3	2	2.31	2		1 1 1 1 1 0	1 0
3	3	3	2.96	3		1 1 1 1 1 1	1 1



ASIC:

$$y_0 = x_0 \cap (\overline{x_1} \cup x_4) \cup \overline{x_3} \cap (x_5 \cup x_6) \cup \dots$$

$$y_1 = \dots$$

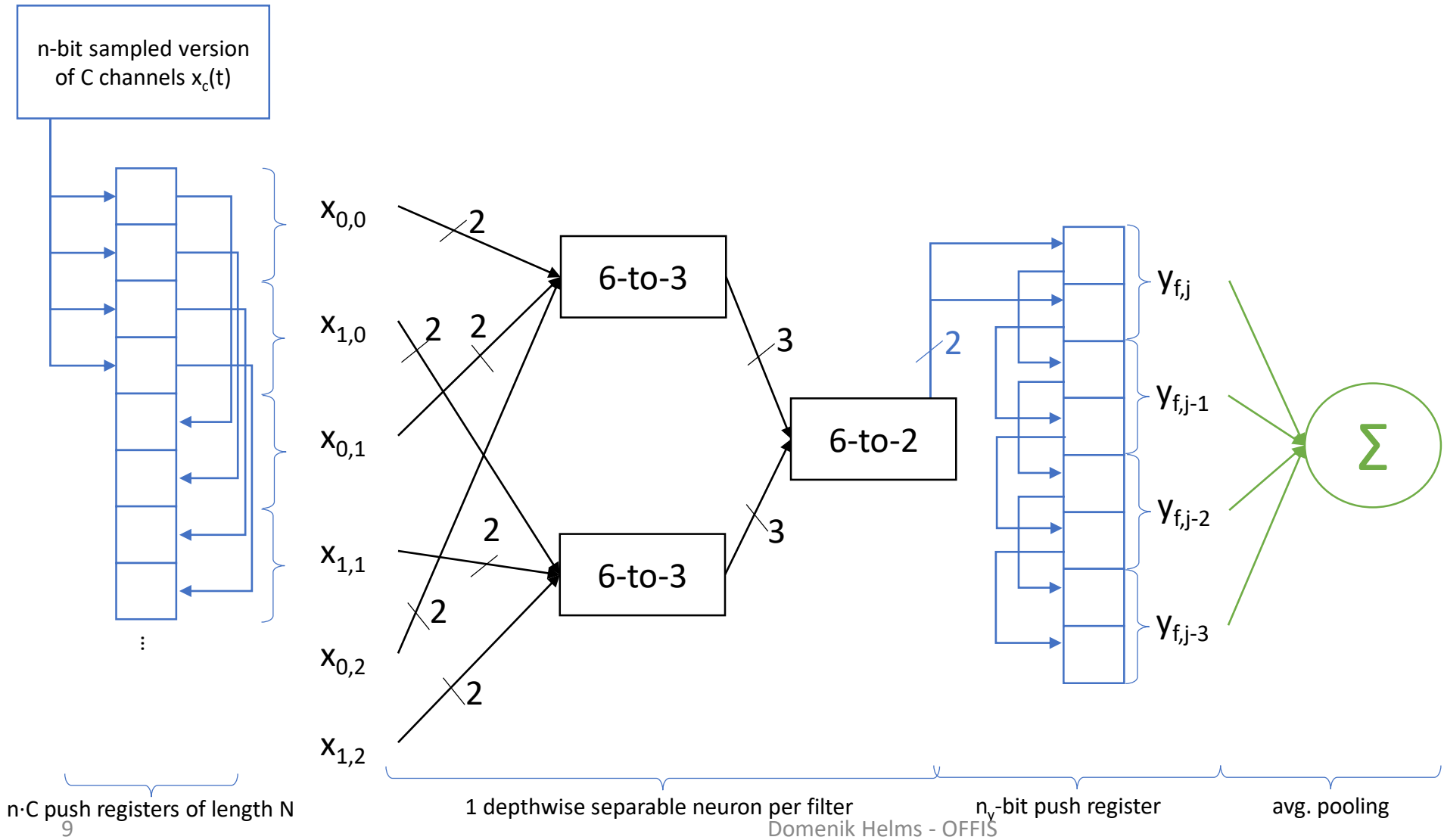
FPGA:

$$y_0 = LUT(x_0, x_1, \dots)$$

$$y_1 = LUT(x_0, x_1, \dots)$$



# Step 5: Pooling



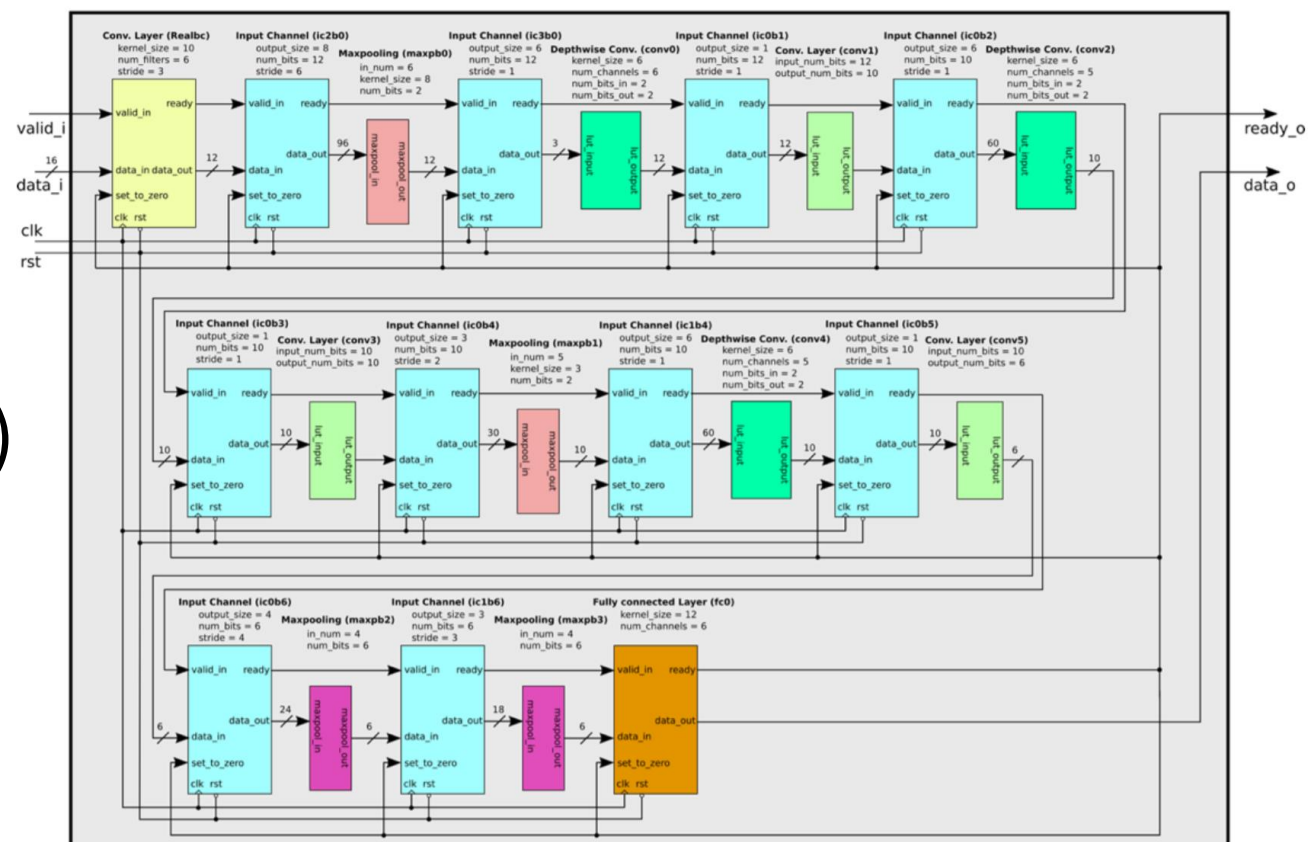
# Results

- All activations are quantized and stored in registers
  - quick, efficient and radiation hard
- All weights and biases are still floating values and implicitly coded into the Boolean functions
  - better training
- No RAM access needed at all
- Each neuron is explicitly and exclusively implemented
  - ~300 transistors per neuron in ASIC
  - ~10 LUTs per neuron in FPGA
- But:
  - it is not possible to retrain ASIC (needs reproduction → \$)
  - FPGA can be retrained, but it contains 5T SRAM cells → radiation)



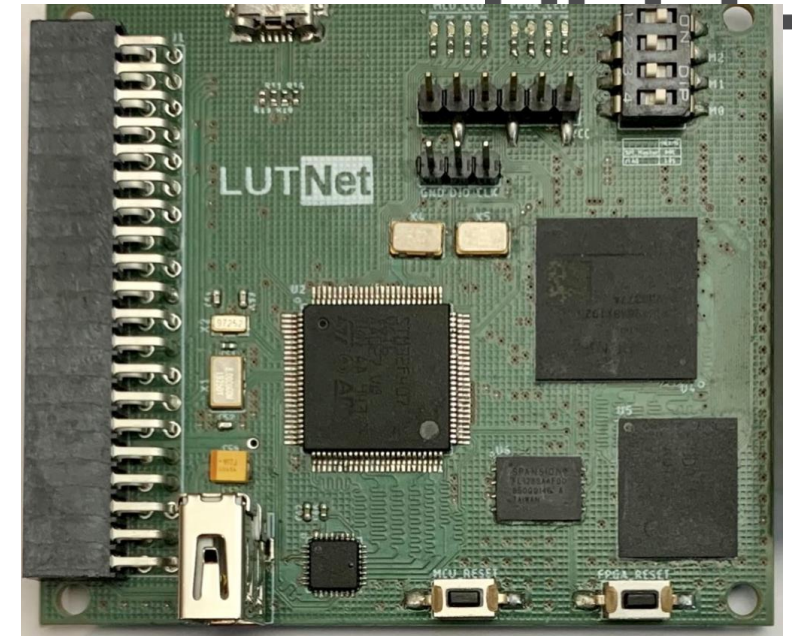
# Evaluation

- Spartan 7 S15
- 2155 LUT (26%)
- 26nJ per sample (full inference)
- 10MHz input sample rate



# Conclusion

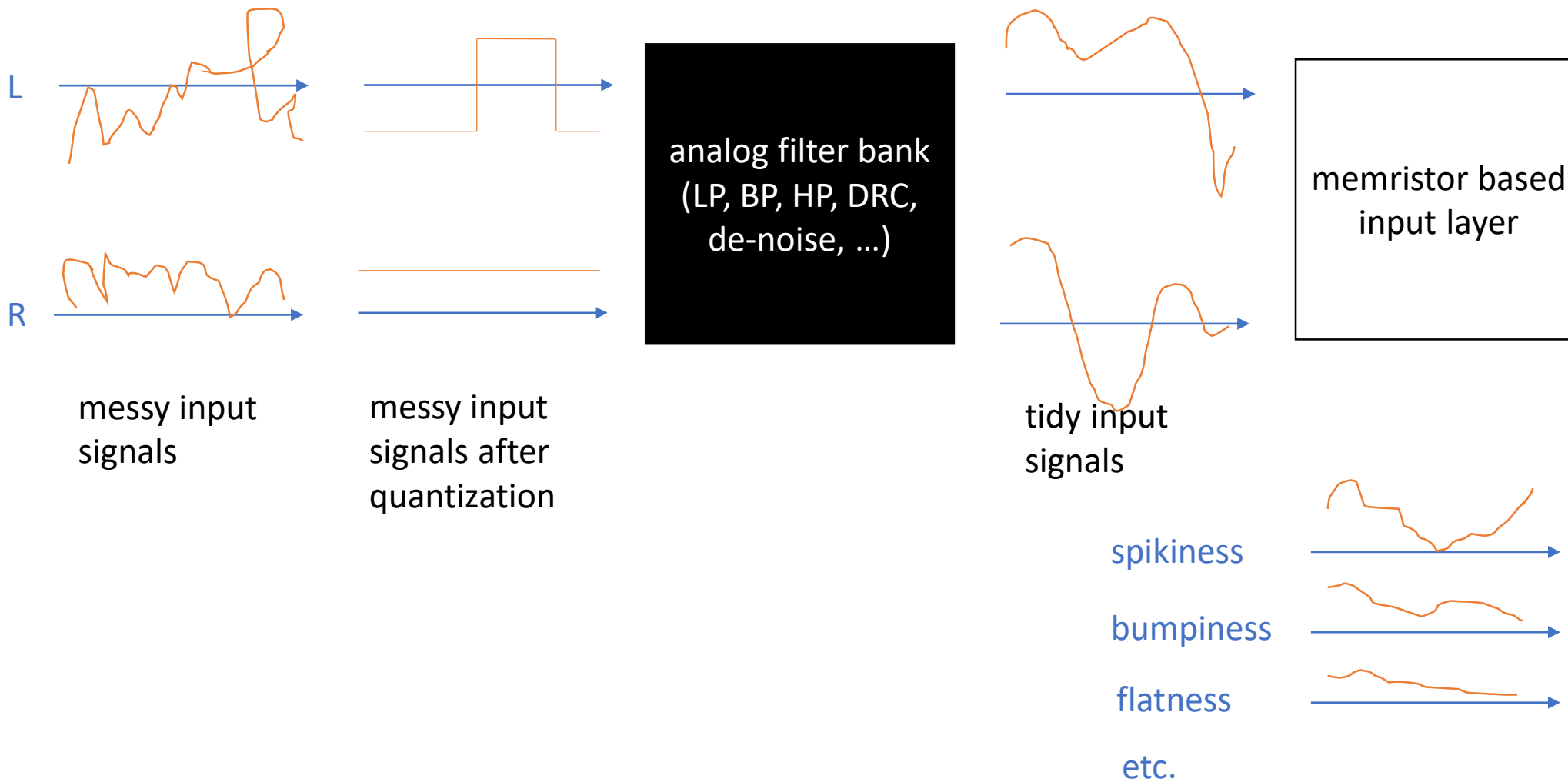
- The method has some fundamental limitations
  - low bitwidth
  - low synapse count
  - 1 dimensional convolution
- Within those limitations, it outperforms any other methodology
- An ASIC solution would improve performance and efficiency by at least an order of magnitude



# The OCTOPUS project idea

1. Develop an analogue frontend chip
2. Improve the methodology
3. Medical application
4. Industry application

# 1 Develop an analogue frontend chip



# 2 Extend the methodology

available

Conv1D	quantization while training
Max Pooling 1D	depthwise separation
Dense	
Global average pooling	ReLU

planned

Conv2D	width&height-wise separation
other std. activation	
nonlinear tensor functions	

$$y = \max \left( 0, \sum_{i=0}^{N-1} v_i \cdot x_i^2 + w_i \cdot x_i + b \right)$$



# 3 Health application

- Use the system as is to detect blood pressure
  - determine blood speed by following the pressure peaks

## 4 Railway application

- Measure the sound (in air or in metal) of a train passing
- use on sensor AI to detect
  - axle traversing
  - weight (load) per axle
  - speed of train
  - type of train
- carefully!!! design and verify a system which uses unreliable AI components in a safety critical system
  - decide if the track is free based on our system's "best guess"