



# Distributed DNN Serving in the Network Data Plane

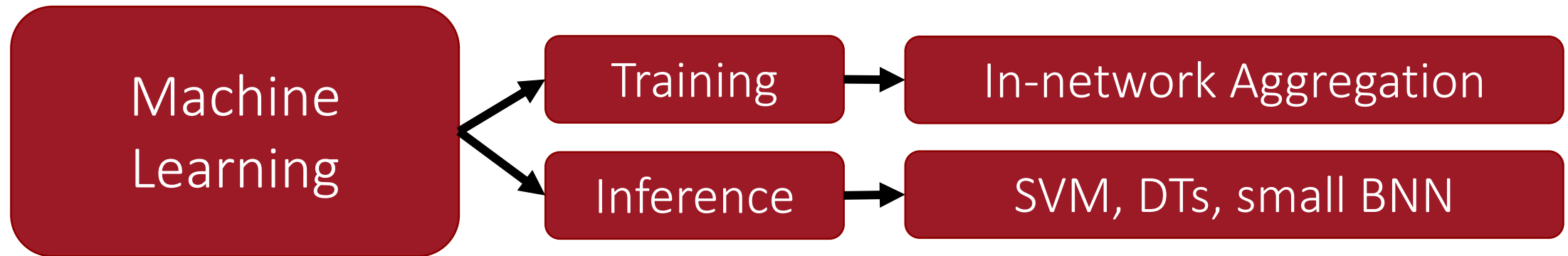
Kamran Razavi<sup>†</sup>, George Karlos<sup>‡</sup>, Vinod Nigade<sup>‡</sup>, Max Mühlhäuser<sup>†</sup>, Lin Wang<sup>†‡</sup>

<sup>†</sup>University of Darmstadt

<sup>‡</sup>Vrije Universiteit Amsterdam

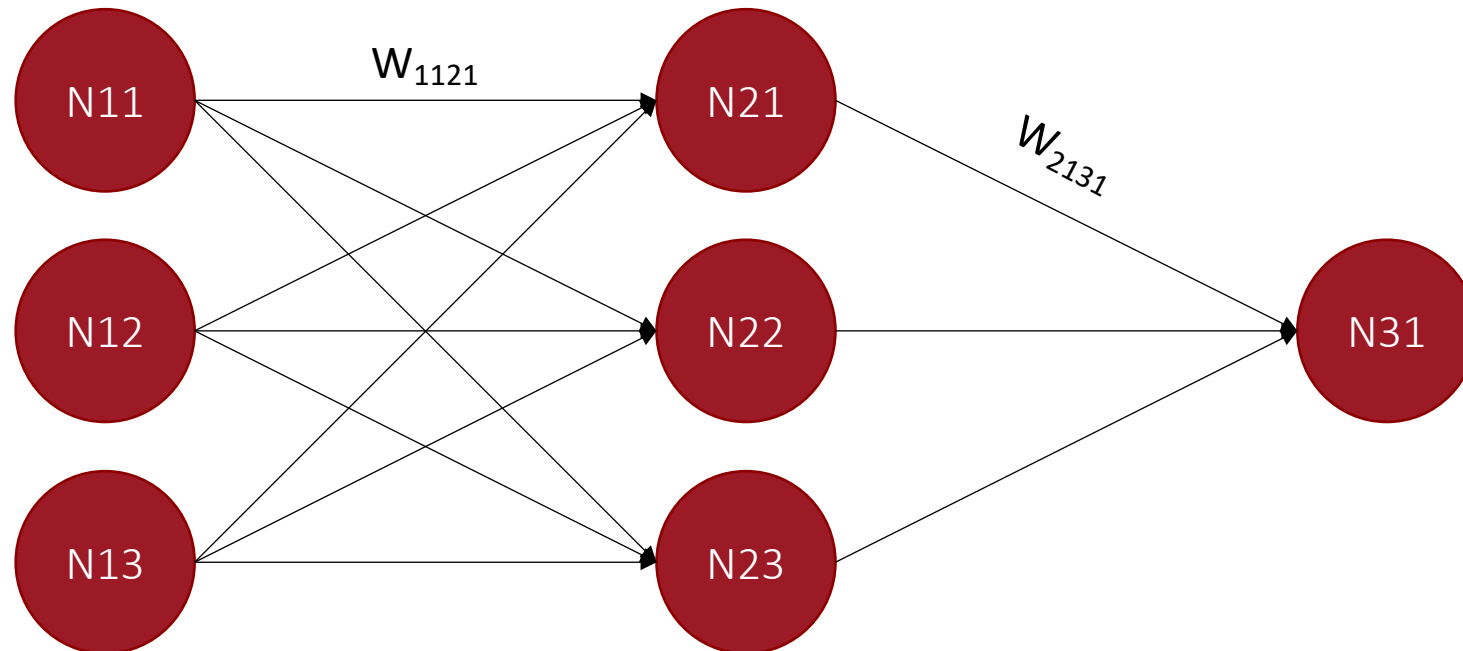


# Programmable Network

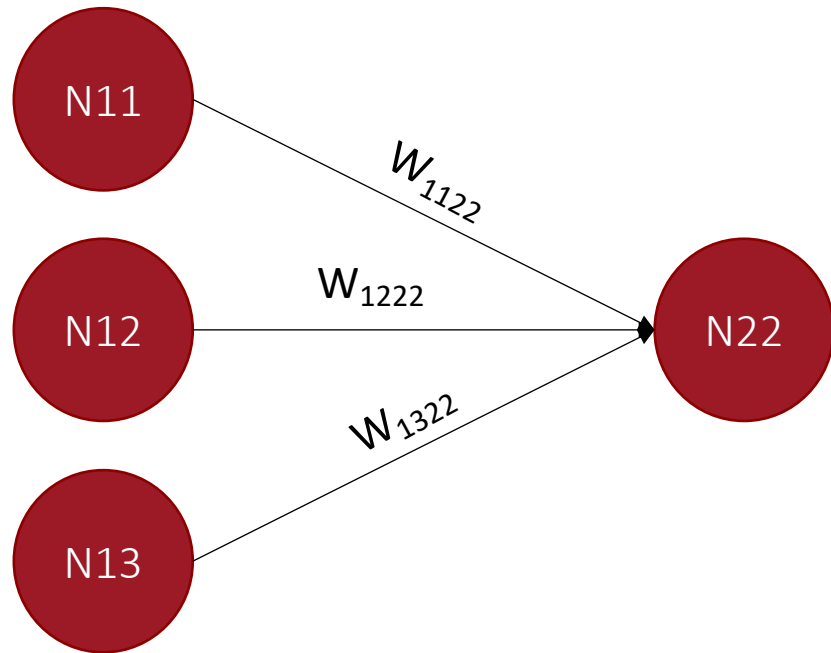


*Can we leverage a programmable network to perform DNN serving?*

# Deep Neural Network Inference



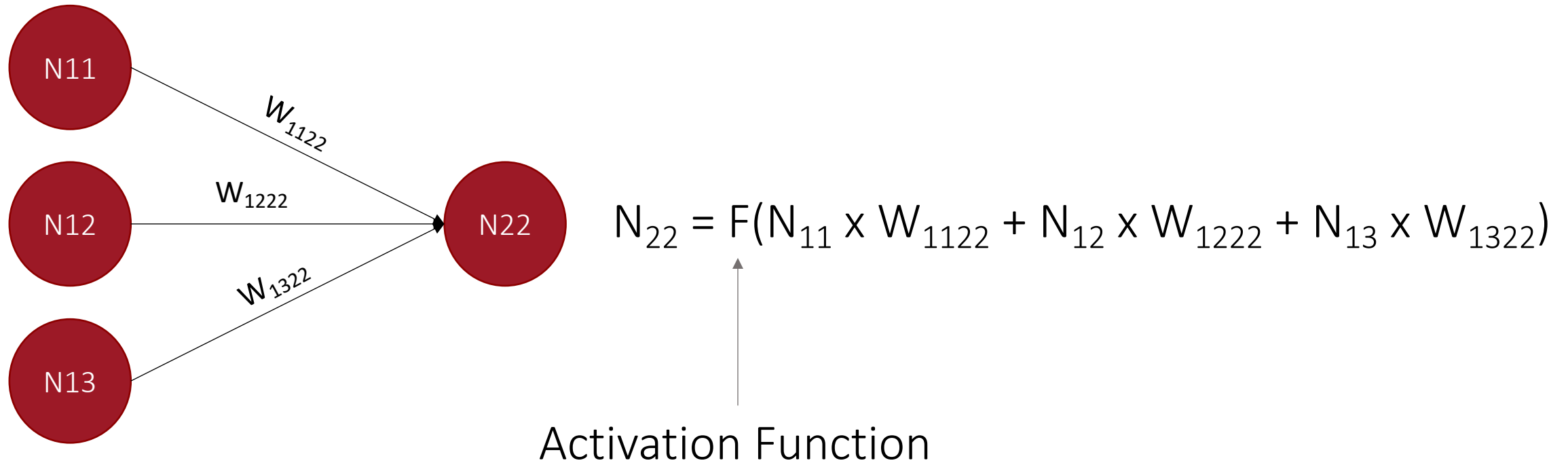
# Deep Neural Network Inference



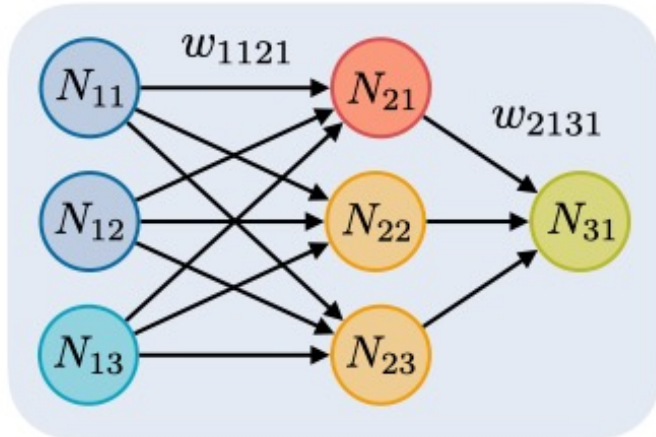
$$N_{22} = N_{11} \times W_{1122} + N_{12} \times W_{1222} + N_{13} \times W_{1322}$$



# Deep Neural Network Inference



# DNN Serving in the Network



# On-Switch Execution

No FPU:  $1.1 * 2$  is time consuming!

- Int-8 (first bit is the sign bit)
- Parallel multiplication

Multiplicand = 13 

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

Multiplier = 5 

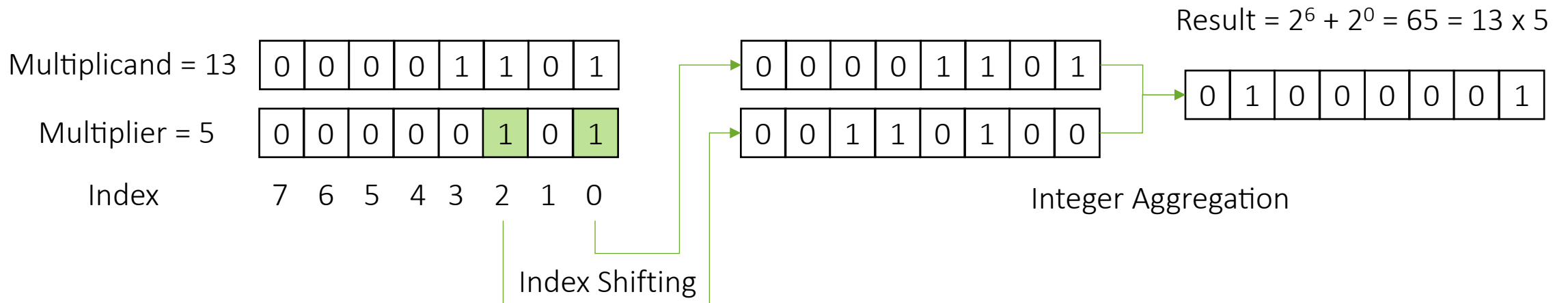
0	0	0	0	0	1	0	1
---	---	---	---	---	---	---	---



# On-Switch Execution

No FPU:  $1.1 * 2$  is time consuming!

- Int-8 (first bit is the sign bit)
- Parallel multiplication



8 ALUs per multiplication -> # of parallel multiplication  $\sim N/8$





# On-Switch Execution

No FPU:  $1.1 * 2$  is time consuming!

- Int-8 (first bit is the sign bit)
- Parallel multiplication

## Activation functions

For example, ReLU:

- If the sign bit == 1, put 0
- Else keep the number

## Layer specific requirements: e.g. MaxPooling

- State storage
- Order of packets



# Case Study: Mini-AlexNet



---

## mini-AlexNet (CIFAR-10)

---

<b>Layer 1</b>	Input: $32 \times 32 \times 3$
<b>Layer 2</b>	Conv1: $3 \times 3 \times 3 \times 64$
<b>Layer 3</b>	Conv2: $3 \times 3 \times 64 \times 192$
<b>Layer 4</b>	Conv3: $3 \times 3 \times 192 \times 384$
<b>Layer 5</b>	FC1: 4096
<b>Layer 6</b>	FC2: 2048
<b>Layer 7</b>	FC3 (Output): 10

---

