

# IMPLEMENTATION OF TRACKING ALGORITHMS FOR LIVE RECONSTRUCTION USING AI PROCESSORS

Klaus Desch

Jochen Kaminski

Michael Lupberger

Stephen Neuendorffer

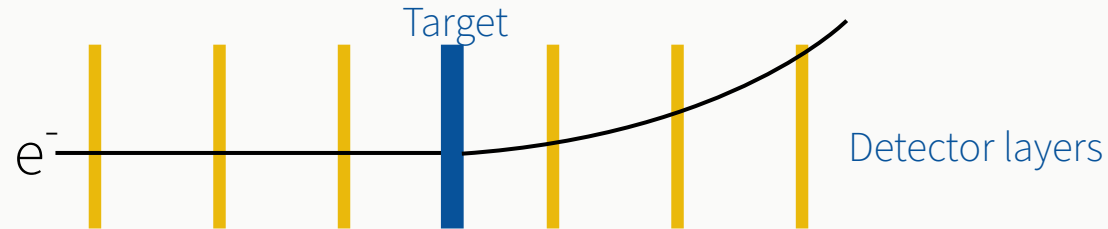
**Patrick Schwäbig**

2022 DPG Spring Meeting Heidelberg, 3/23/2022



## MOTIVATION

- Dark photon experiment, electron on target
- To be built at the ELSA accelerator at University of Bonn



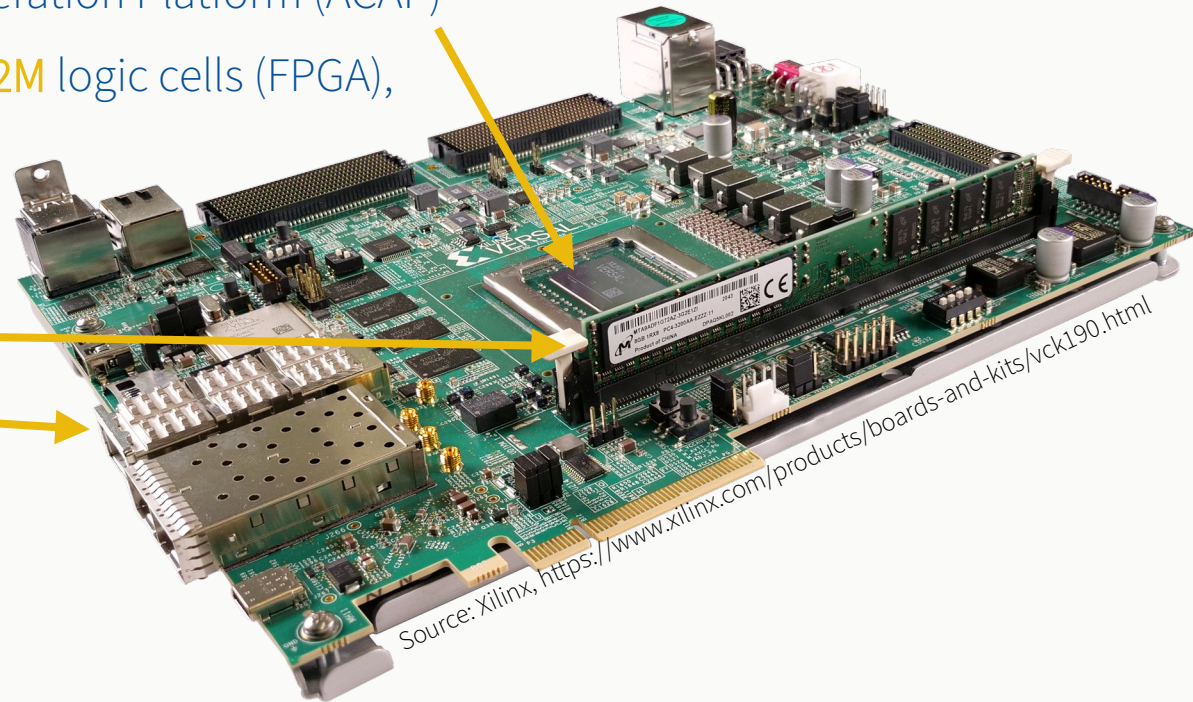
- Low signal expected, want to measure  $4 \times 10^{14}$  electrons
- For more information see **T 64.8**, directly after this talk, but different sessions
- Want sophisticated, fast trigger
  - Use hardware accelerated algorithms for the trigger → live tracking
- This talk: **pattern recognition with hardware-accelerated neural network**

## HARDWARE: XILINX VERSAL / VCK190

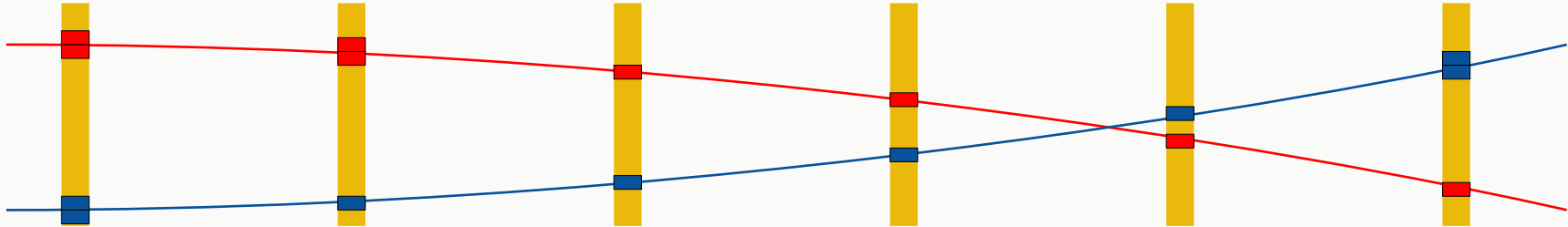
- Evaluation board **VCK190**
- Versal **VC1902** Adaptive Compute Acceleration Platform (ACAP)
- **400** AI processors (“AI engines”), nearly **2M** logic cells (FPGA),  
**2k** DSPs, Arm CPU, Arm RPU

Board with

- 8 GB DDR4 RAM
- QSFP28 for 100 Gbit Ethernet

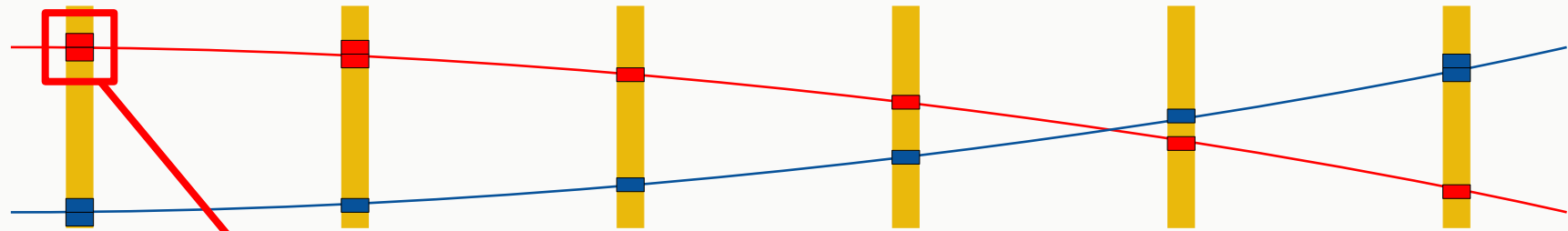


## BUILDING NEURAL NETWORK INPUT

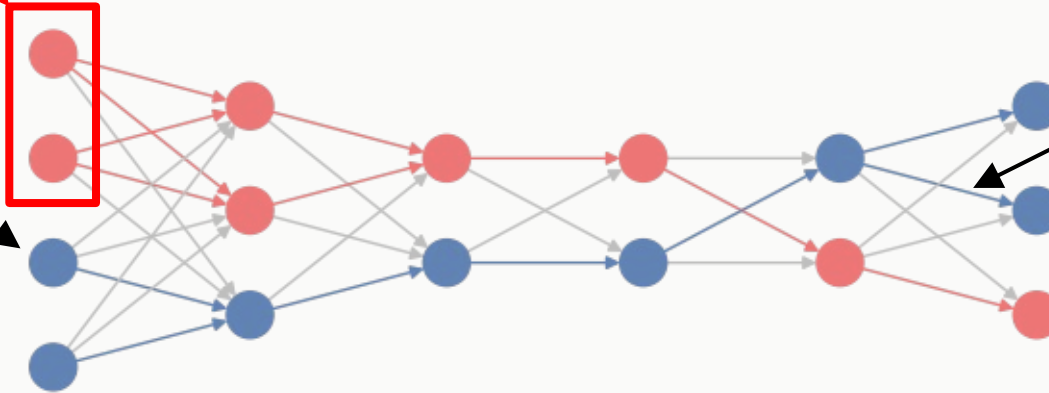


- Using **Timepix3** chip → ToT, ToA, Position for each hit → read out with ACAP
- Build a graph...

# BUILDING NEURAL NETWORK INPUT

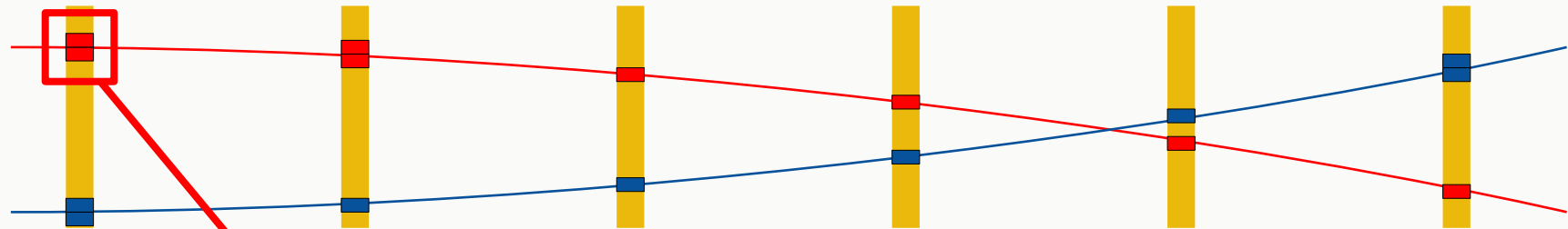


Node = Hit  
with features:  
ToA, ToT, Pos.

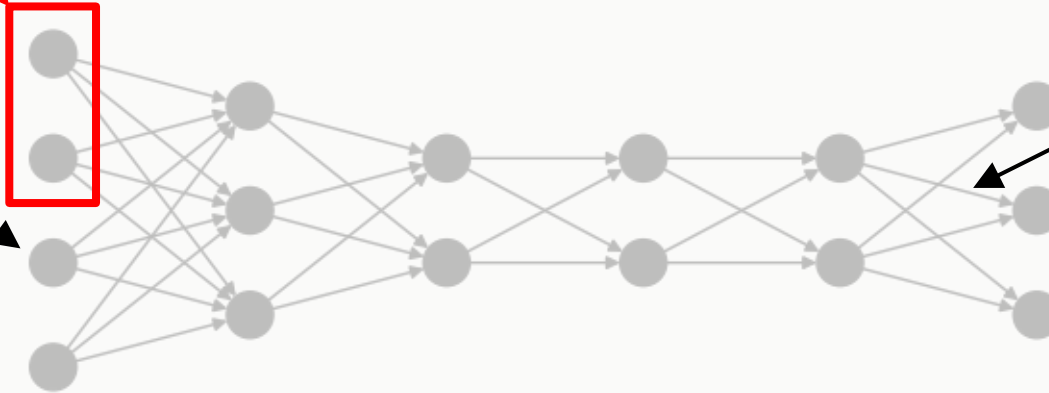


Edges = possible connections (hits between two adjacent layers)

## BUILDING NEURAL NETWORK INPUT

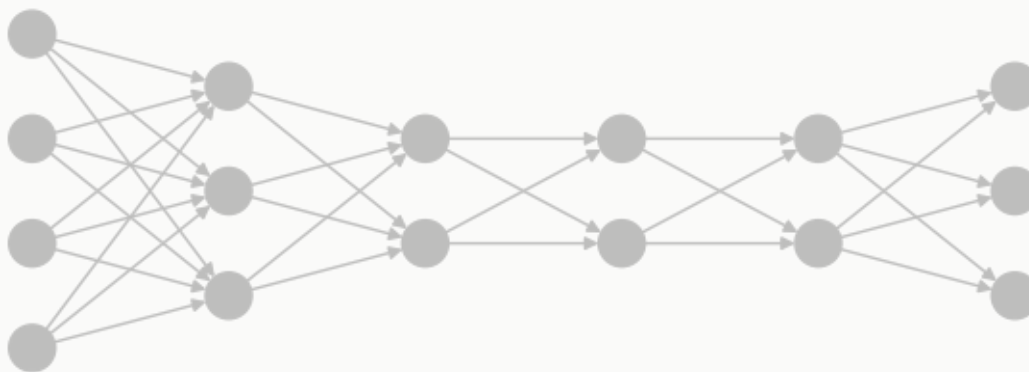


Node = Hit  
with features:  
ToA, ToT, Pos.



Edges =  
**possible**  
connections  
(hits between two  
adjacent layers)

## GNN FOR PATTERN RECOGNITION



- Classify edges: true edge or false edge  
→ Retrieve original track (pattern recognition)
- Use Interaction Network\* which is a type of **Graph Neural Network** (GNN)
- Uses graph (hits and possible connections) as input
- Input graph **is not** the neural network (which can also be shown as a graph)
- GNNs generalization of other NNs: data is not 2D or 1D anymore

\* described in DeZoort et al. (2021)

## STRUCTURE OF THE GNN

Detector/Timepix3 → Readout & build graph (already on ACAP)

Classically:  
Executed sequentially

- 1) Loop possible connections:
  - Run neural network “R1”
- 2) Loop hits:
  - Rearrange output of “R1”  
→ Run neural network “O”
- 3) Loop possible connections:
  - Rearrange output of “O”  
→ Run neural network “R2”  
& Sigmoid/apply threshold

→ Classified edges



## IMPLEMENTATION FOR ACAP — PIPELINING

Detector/Timepix3 → Readout & build graph (already on ACAP)

1) Loop possible connections:

- Run neural network “R1” → AIE 1

2) Loop hits:

- Rearrange output of “R1” → AIE 2
- Run neural network “O” → AIE 3

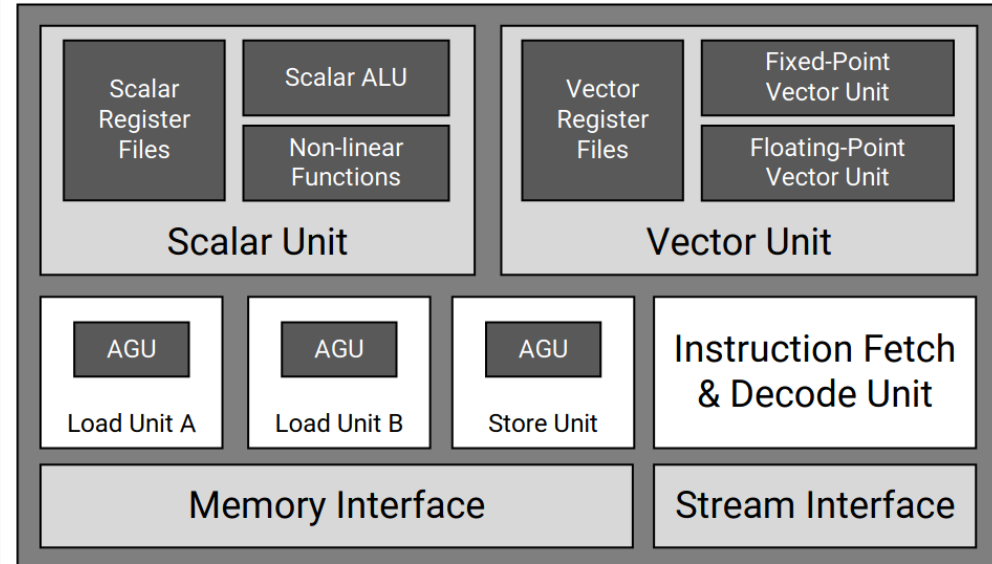
3) Loop possible connections:

- Rearrange output of “O” → AIE 4
- Run neural network “R2” → AIE 5
- & Sigmoid/apply threshold → Classified edges

## IMPLEMENTATION FOR ACAP — PARALLEL EXECUTION

### AI Engines of VCK190:

- Highly parallelized: SIMD VLIW



Source: Xilinx UG1079

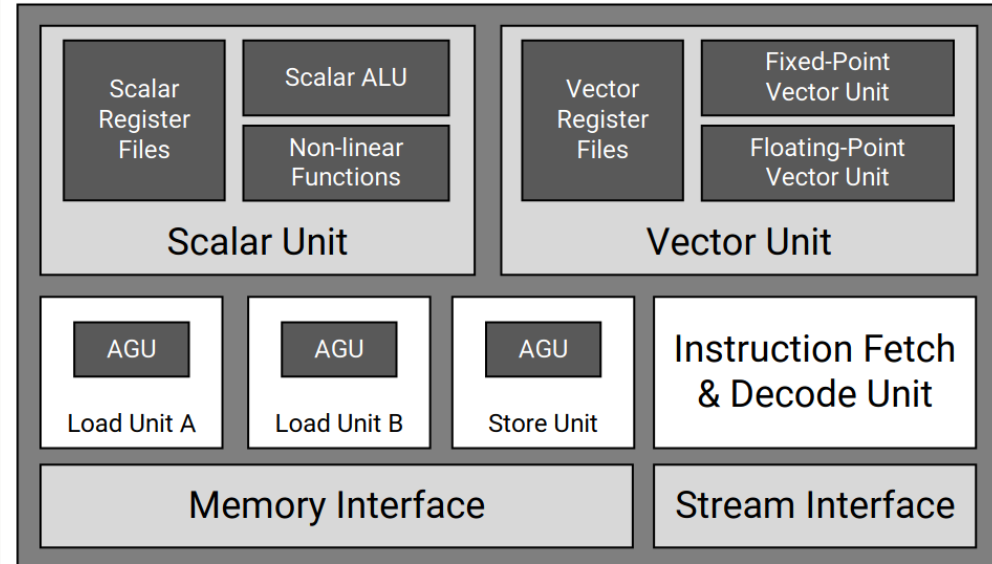
## IMPLEMENTATION FOR ACAP — PARALLEL EXECUTION

Single instruction, multiple data

AI Engines of VCK190:

- Highly parallelized: SIMD VLIW

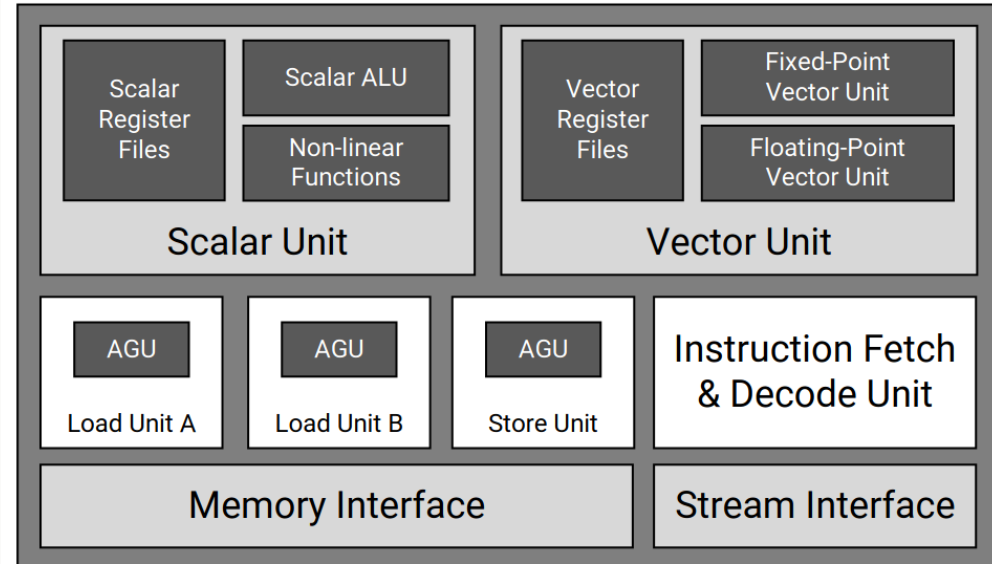
Very long instruction word



Source: Xilinx UG1079

### AI Engines of VCK190:

- Highly parallelized: SIMD VLIW
- VLIW → Simultaneous execution of:
  - Loading data
  - Storing data
  - Computations (scalar and **vector** → SIMD)
- Running at 1 GHz
- 8 Accumulators for 32b floating-point operations

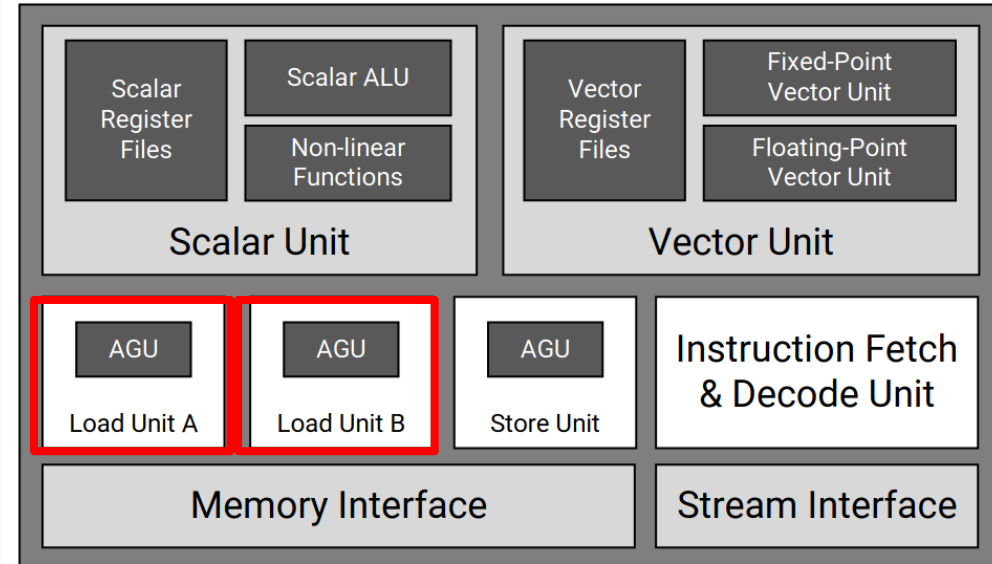


Source: Xilinx UG1079

→ Work on 8 possible connections simultaneously  
while loading data, moving data & incrementing pointers...

### AI Engines of VCK190:

- Highly parallelized: SIMD VLIW
- VLIW → Simultaneous execution of:
  - Loading data
  - Storing data
  - Computations (scalar and **vector** → SIMD)
- Running at 1 GHz
- 8 Accumulators for 32b floating-point operations



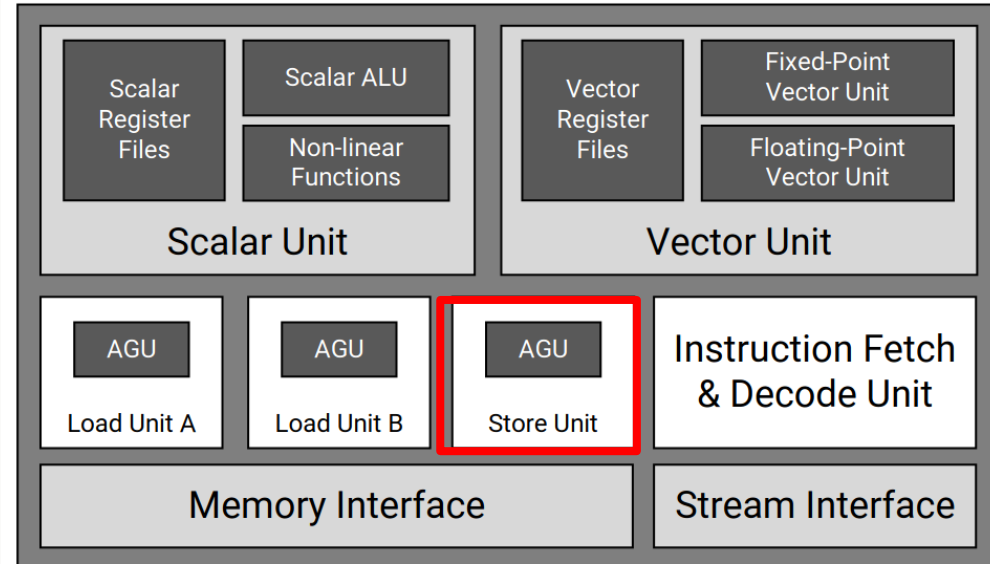
Source: Xilinx UG1079

→ Work on 8 possible connections simultaneously  
while loading data, moving data & incrementing pointers...

### AI Engines of VCK190:

- Highly parallelized: SIMD VLIW
- VLIW → Simultaneous execution of:
  - Loading data
  - Storing data
  - Computations (scalar and **vector** → SIMD)
- Running at 1 GHz
- 8 Accumulators for 32b floating-point operations

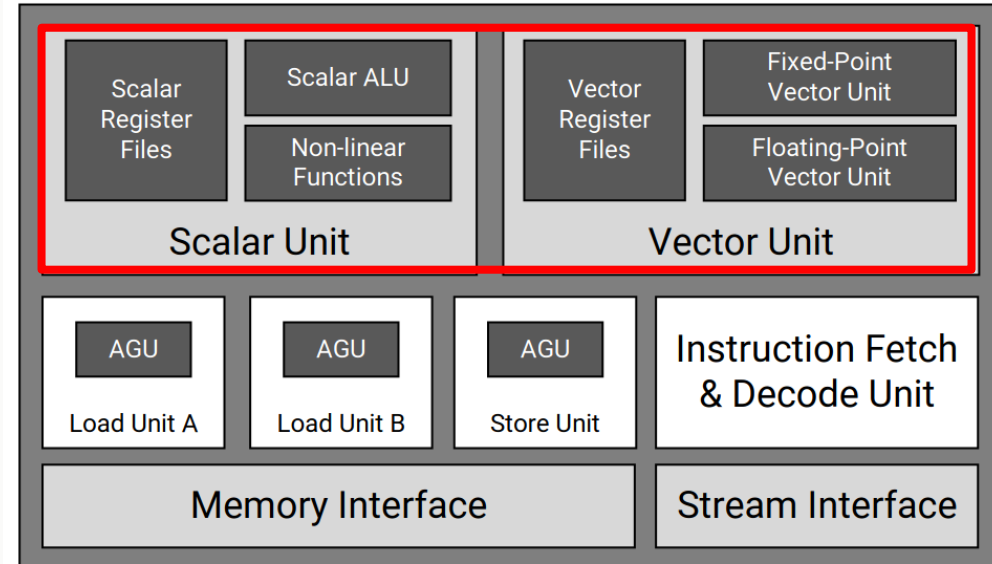
→ Work on 8 possible connections simultaneously  
while loading data, moving data & incrementing pointers...



Source: Xilinx UG1079

### AI Engines of VCK190:

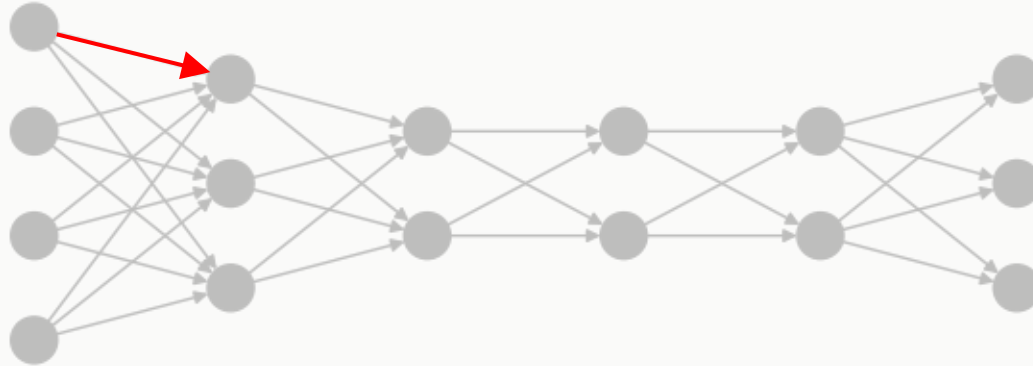
- Highly parallelized: SIMD VLIW
- VLIW → Simultaneous execution of:
  - Loading data
  - Storing data
  - Computations (scalar and **vector** → SIMD)
- Running at 1 GHz
- 8 Accumulators for 32b floating-point operations



Source: Xilinx UG1079

→ Work on 8 possible connections simultaneously  
while loading data, moving data & incrementing pointers...

## GNN FOR PATTERN RECOGNITION



1) Loop possible connections:

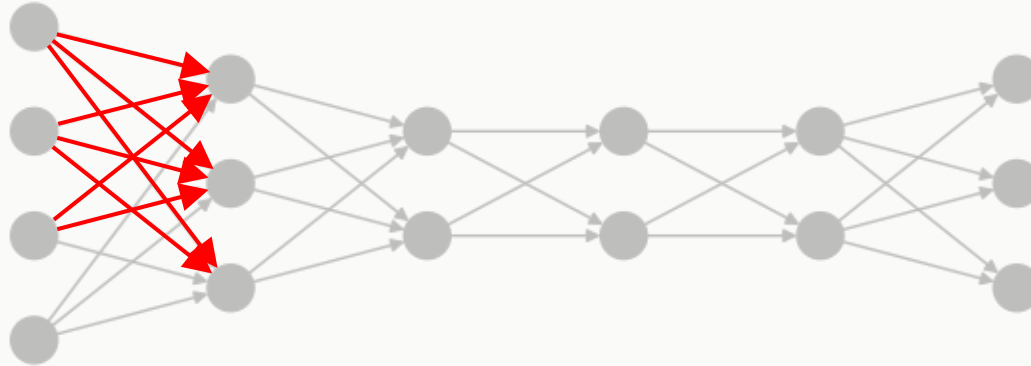
- Run neural network “R1” → AIE 1

Con. 1:

ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>



## GNN FOR PATTERN RECOGNITION



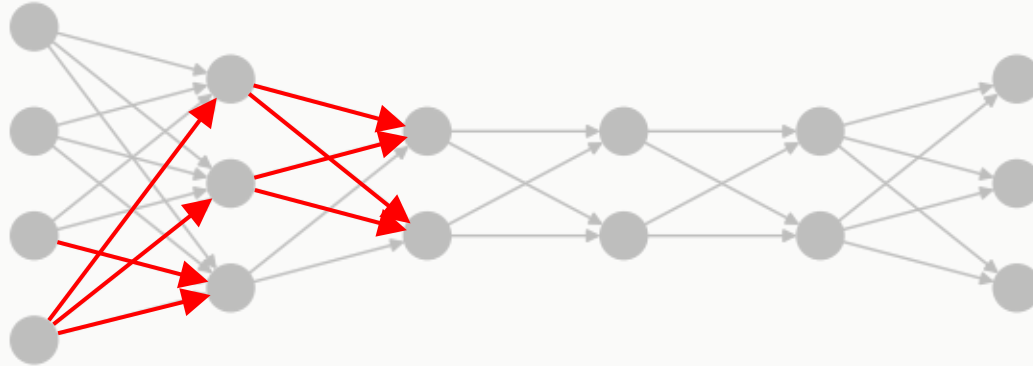
1) Loop possible connections/8:

- Run neural network “R1” → AIE 1

...with 8 accumulators!

Con. 1:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 2:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 3:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 4:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 5:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 6:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 7:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 8:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>

## GNN FOR PATTERN RECOGNITION



1) Loop possible connections/8:

- Run neural network “R1” → AIE 1

...with 8 accumulators!

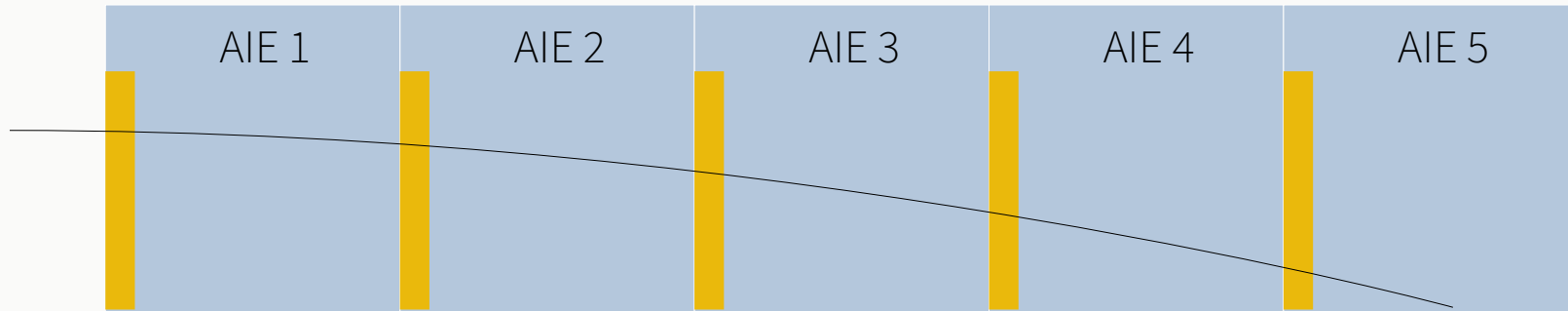
Con. 9:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con.10:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con.11:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 12:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con.13:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con.14:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 15:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>
Con. 16:	ToA <sub>Start</sub>	ToA <sub>End</sub>	ToT <sub>Start</sub>	ToT <sub>End</sub>

## CURRENT STATUS & OUTLOOK

- Status: Implemented, currently improving performance
- Targeting  $O(\mu\text{s})$  to classify each graph
- 32b floating-point  $\rightarrow$  Quantize the NN? Use 8b integer, x16 speedup
- Could also include FPGA-part of ACAP
- Only using 5 AI engines! Space to go faster or do more

## CURRENT STATUS & OUTLOOK

- Status: Implemented, currently improving performance
- Targeting  $O(\mu s)$  to classify each graph
- 32b floating-point  $\rightarrow$  Quantize the NN? Use 8b integer, x16 speedup
- Could also include FPGA-part of ACAP
- Only using 5 AI engines! Space to go faster or do more:
  - Build tracks
  - Fit tracks, e.g. using a Kalman filter:



Thanks for your attention!

Questions?