



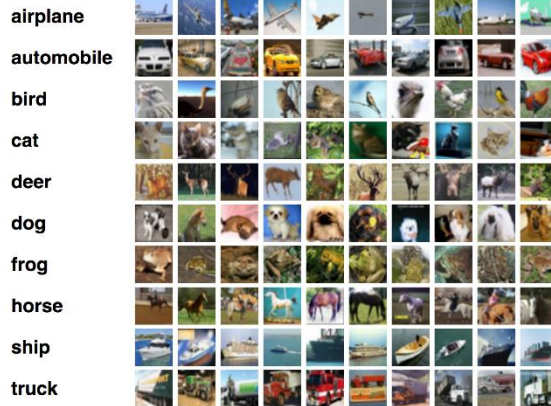
# TNPU: An Efficient Accelerator Architecture for Training Convolutional Neural Networks

Jiajun Li, Guihai Yan, Wenyan Lu, Shuhao Jiang, Shijun Gong,  
Jingya Wu, Junchao Yan, Xiaowei Li

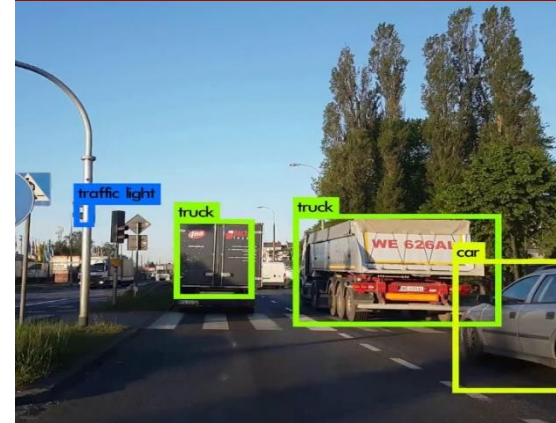
2019/1/23

# Convolutional neural networks

## Image classification



## Object detection



## Speech recognition

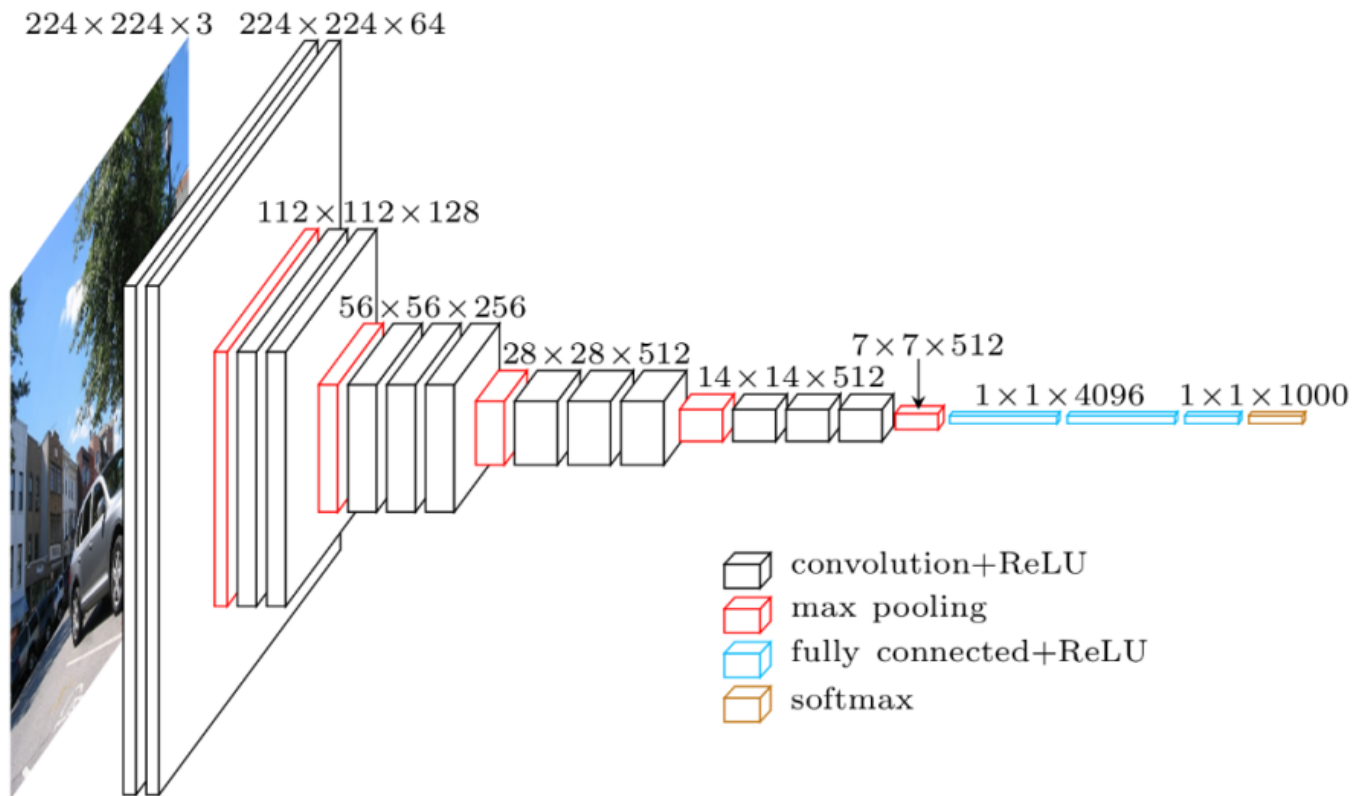


## Playing games

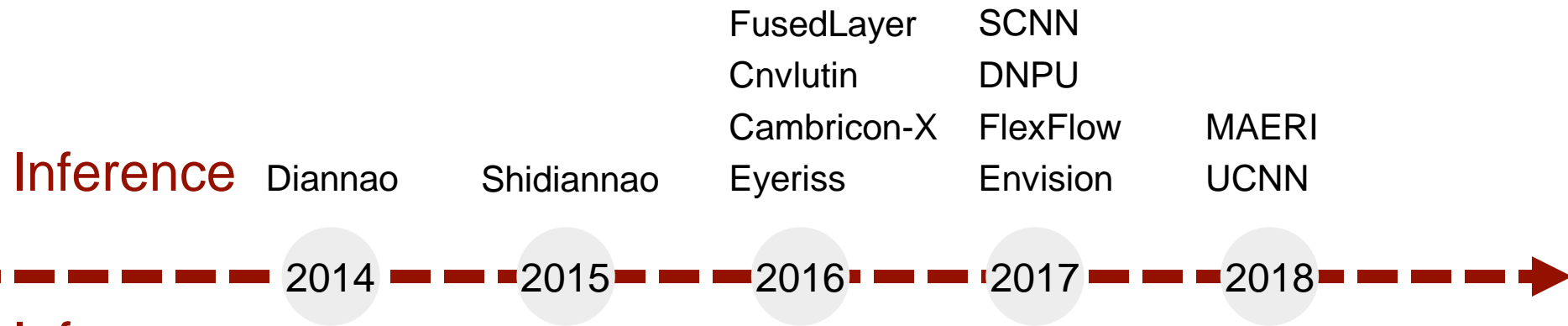


# Intensive operations in CNNs

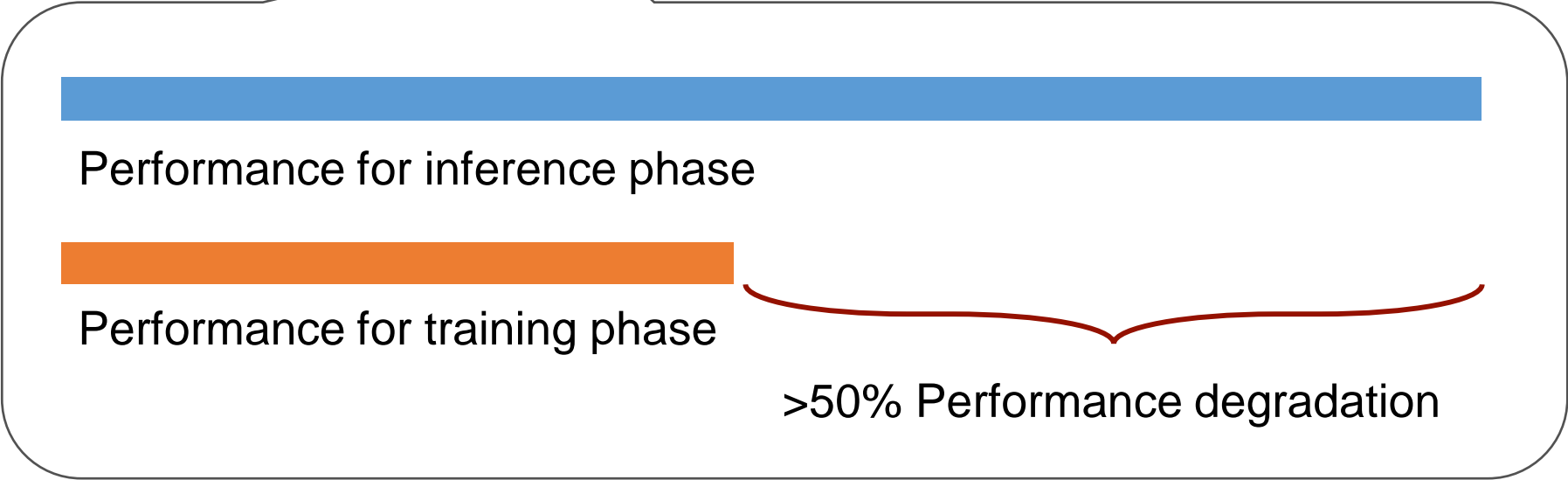
- Training AlexNet takes six days on two GPUs



# CNN accelerators

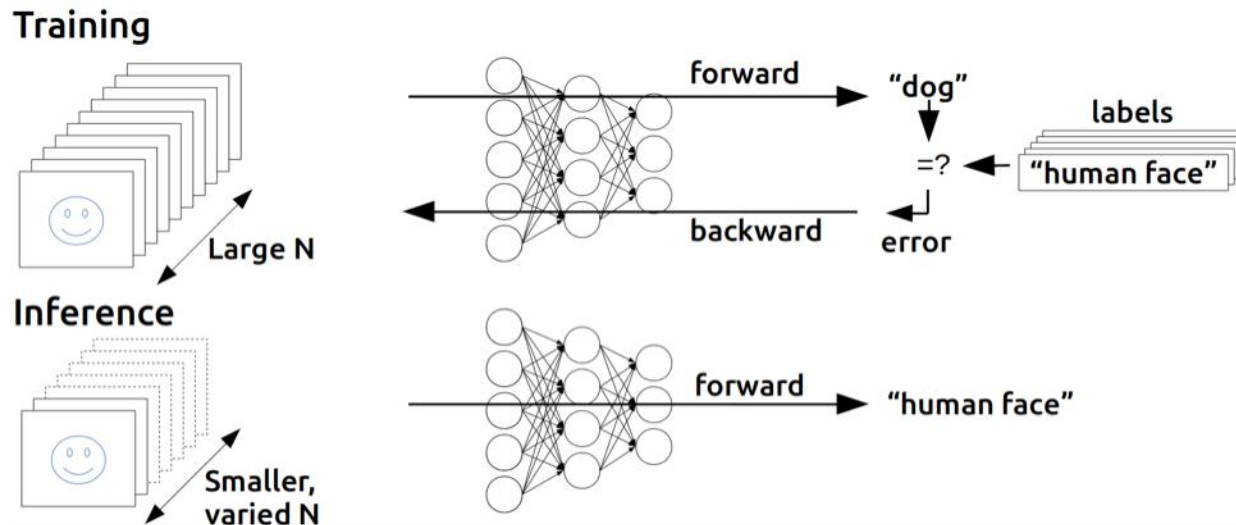


Inference  
& Training



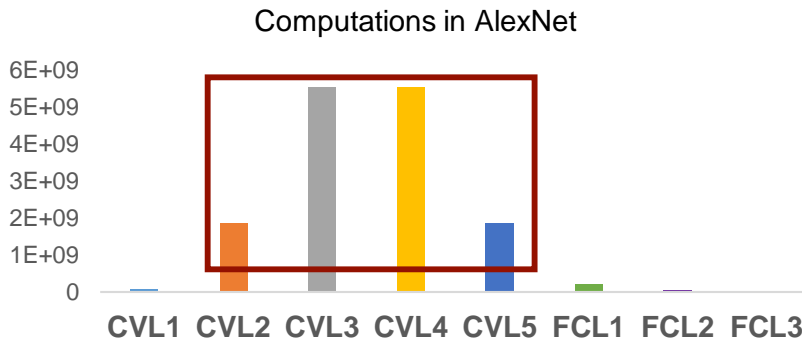
# Inference vs. Training

- New challenges in training
  - Optimization for fully-connected layers
  - Bidirectional data dependency
  - Wide-range of convolutional kernel sizes

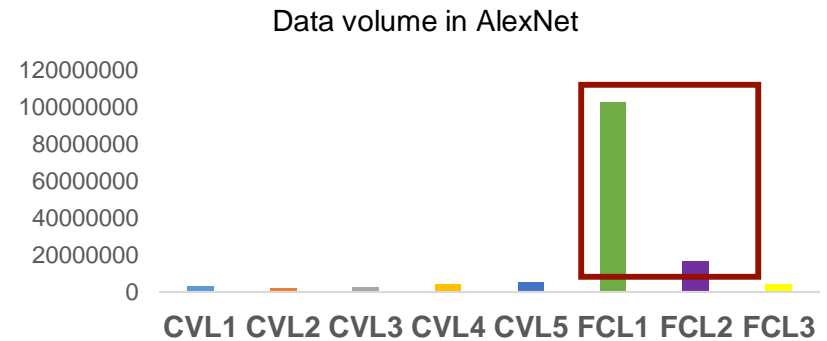


# New challenges in training (1/3)

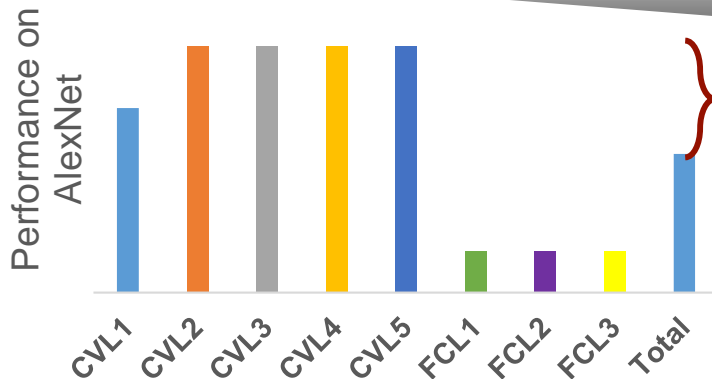
- Comparison between CVLs and FCLs



CVLs dominate computation volume



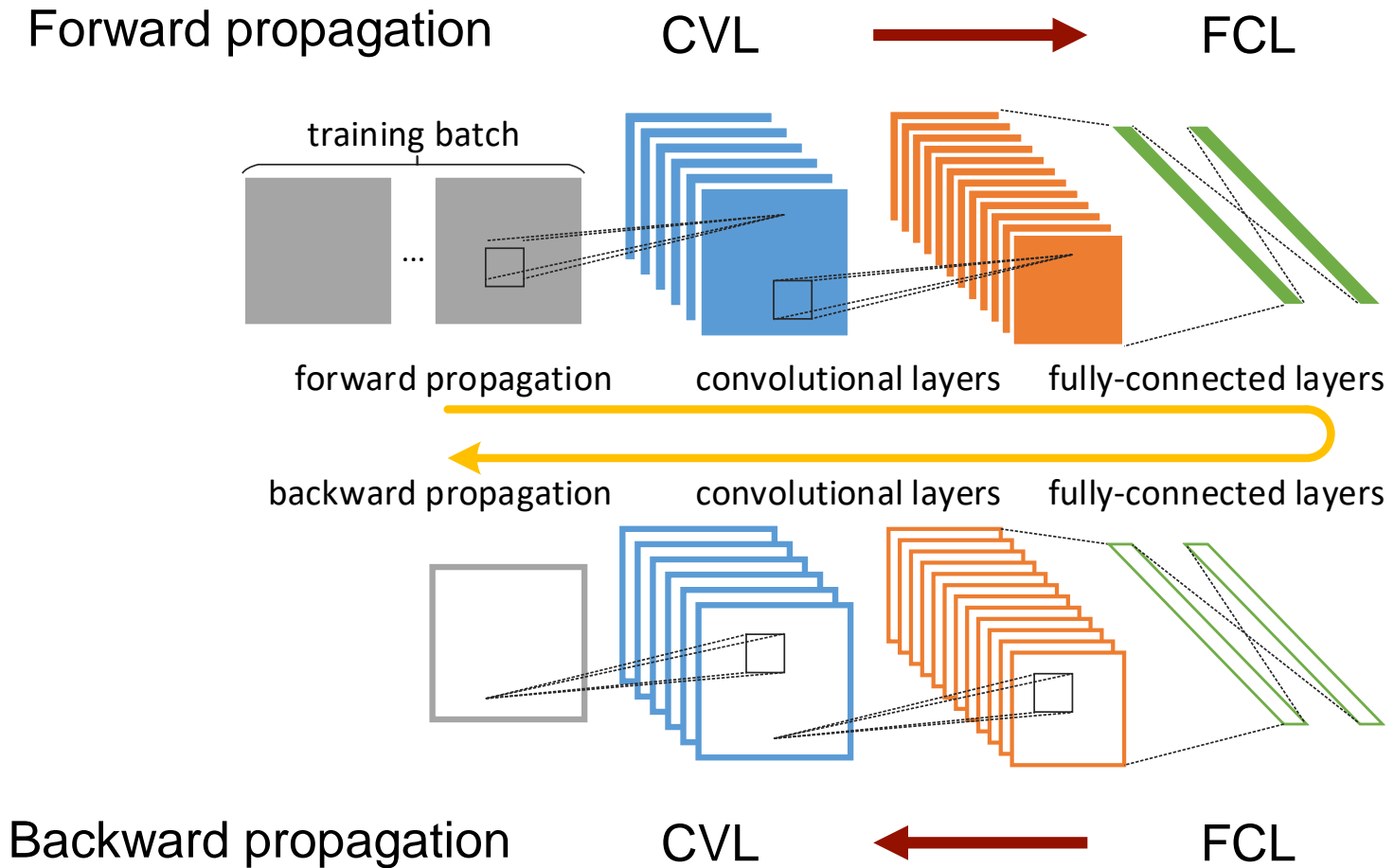
FCLs dominate data volume  
Cannot compress in training



Low performance on FCLs lead to significant performance degradation

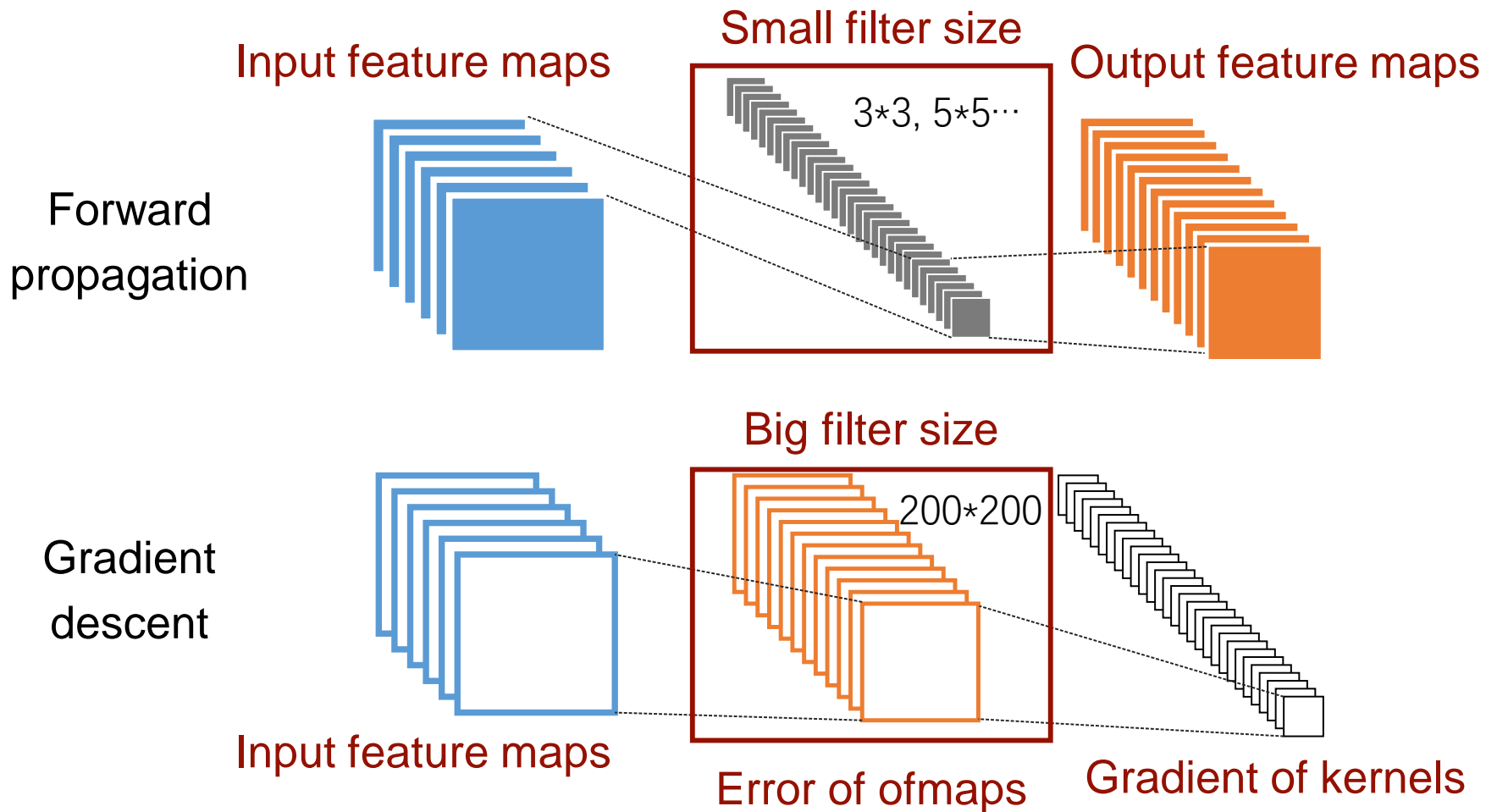
# New challenges in training (2/3)

- Bidirectional data dependency



# New challenges in training (3/3)

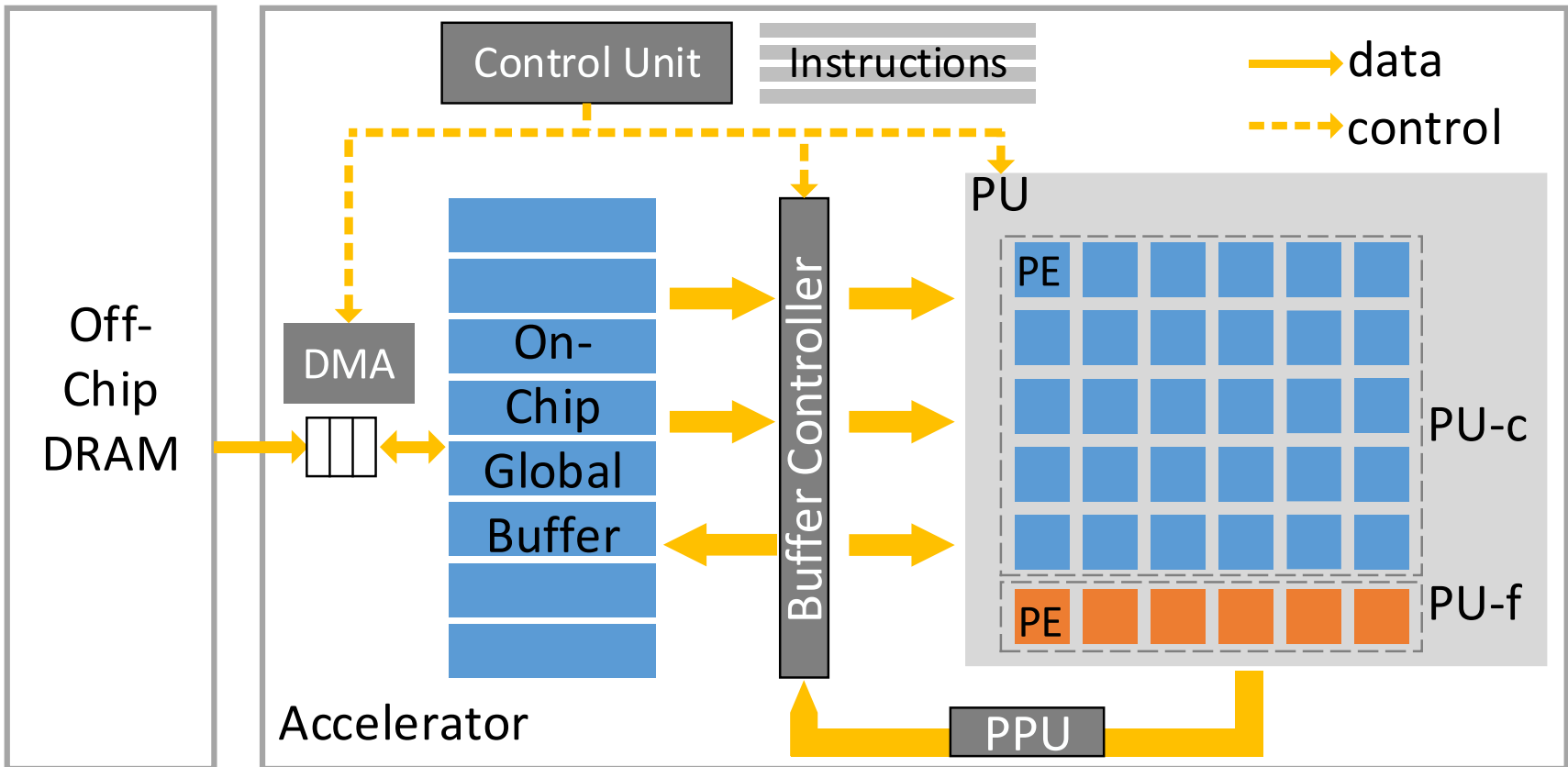
- Extremely large convolutional kernels





# Our solution: TNPU

- Training **N**eural **N**etwork **P**rocessing **U**nit



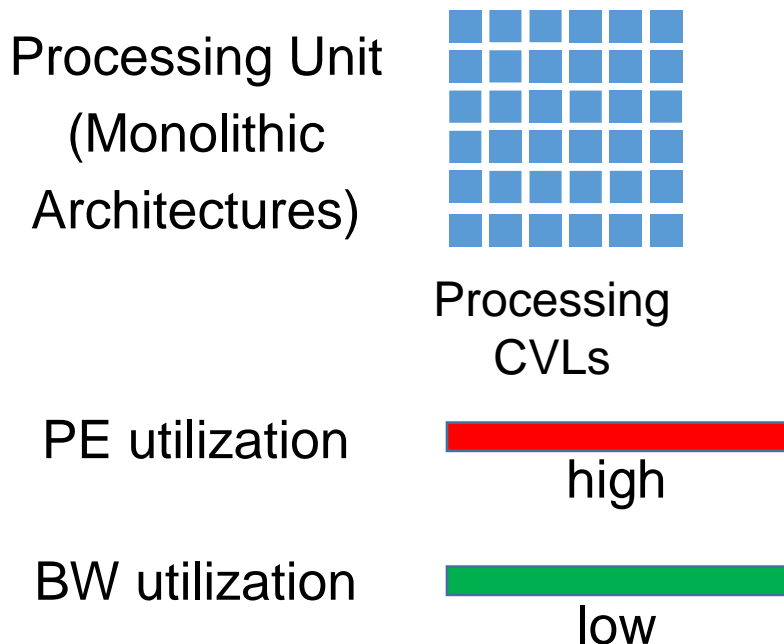
# CVLs vs. FCLs

---

- CVLs: computation  $\gg$  parameter
  - Computation hungry  $\rightarrow$  High PE utilization
- FCLs: computation  $\sim$  parameter
  - Memory bandwidth hungry  $\rightarrow$  High BW utilization

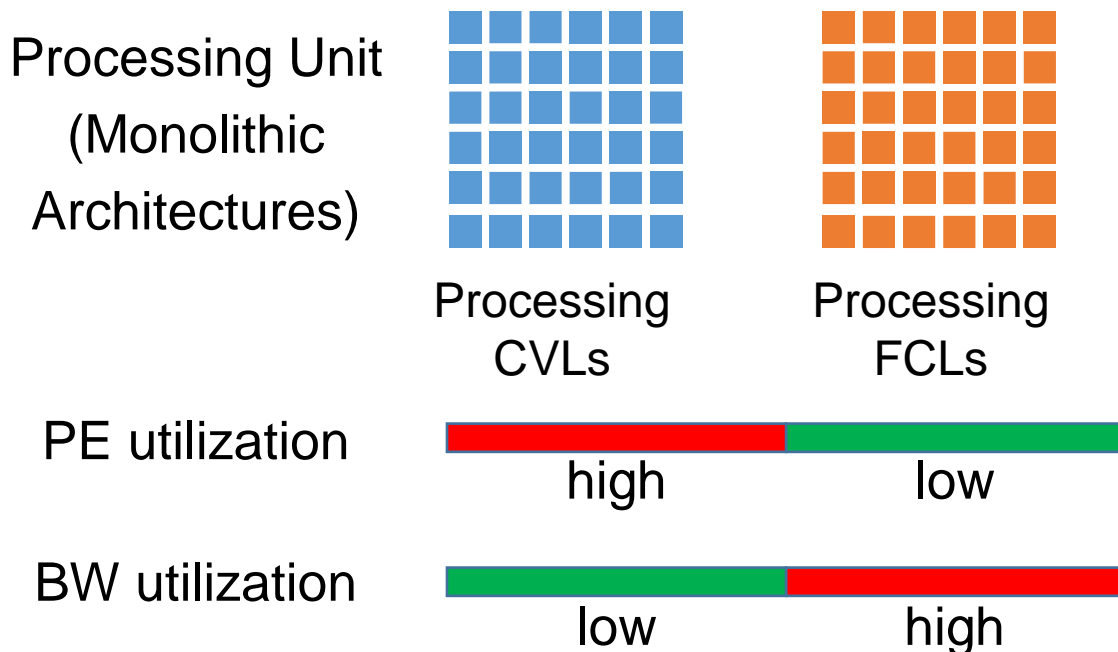
# CVLs vs. FCLs

- CVLs: computation  $\gg$  parameter
  - Computation hungry  $\rightarrow$  High PE utilization
- FCLs: computation  $\sim$  parameter
  - Memory bandwidth hungry  $\rightarrow$  High BW utilization



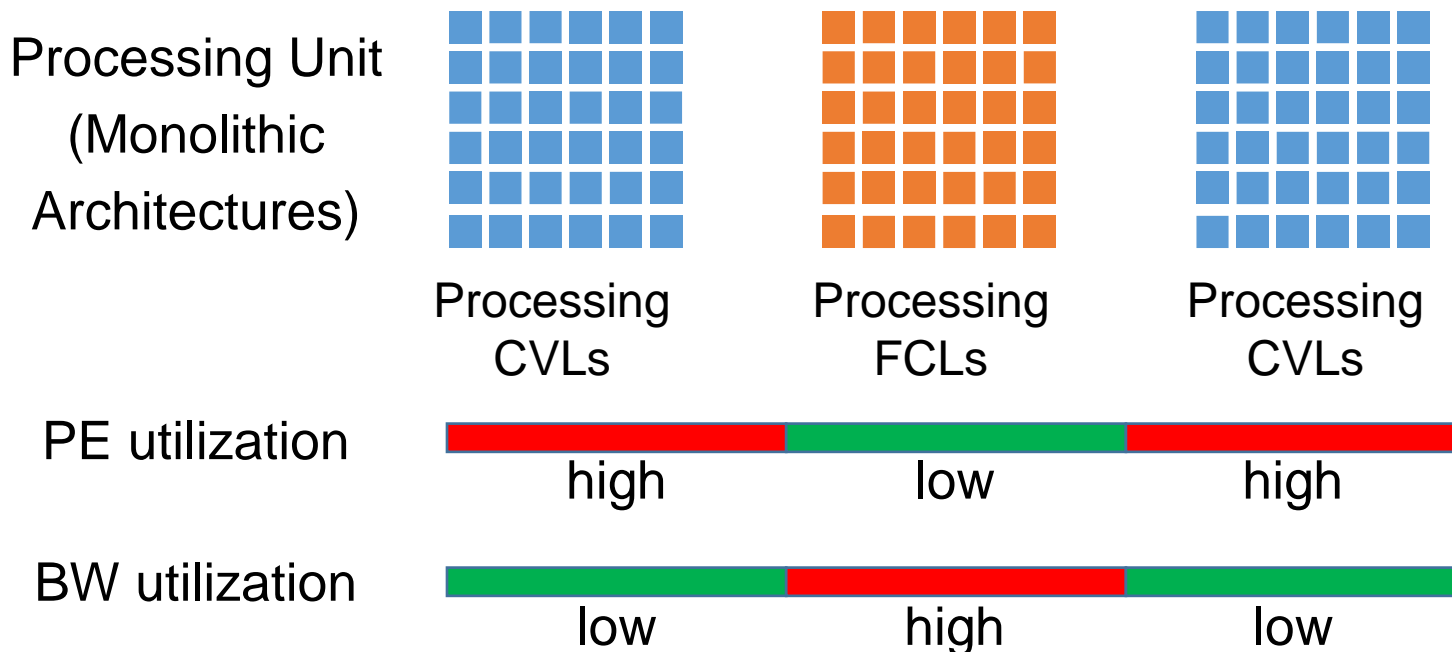
# CVLs vs. FCLs

- CVLs: computation  $\gg$  parameter
  - Computation hungry  $\rightarrow$  High PE utilization
- FCLs: computation  $\sim$  parameter
  - Memory bandwidth hungry  $\rightarrow$  High BW utilization



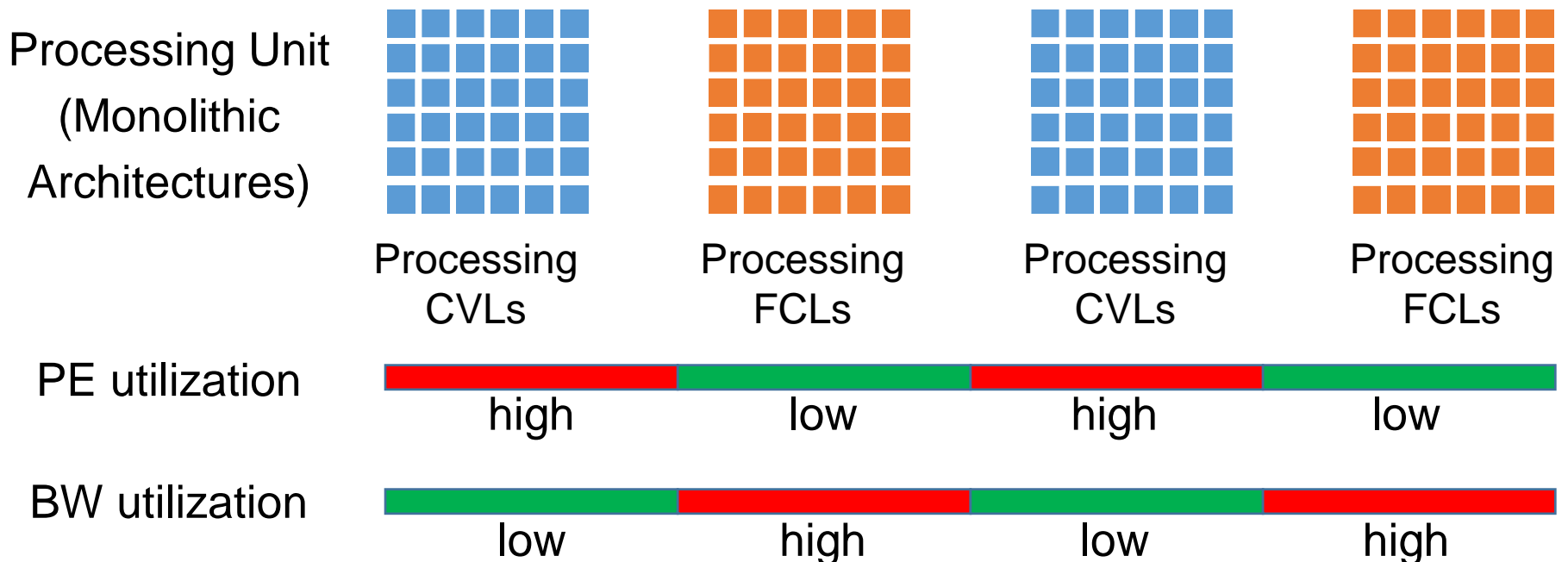
# CVLs vs. FCLs

- CVLs: computation  $\gg$  parameter
  - Computation hungry  $\rightarrow$  High PE utilization
- FCLs: computation  $\sim$  parameter
  - Memory bandwidth hungry  $\rightarrow$  High BW utilization



# CVLs vs. FCLs

- CVLs: computation  $\gg$  parameter
  - Computation hungry  $\rightarrow$  High PE utilization
- FCLs: computation  $\sim$  parameter
  - Memory bandwidth hungry  $\rightarrow$  High BW utilization



# CVLs vs. FCLs

- CVLs: computation >> parameter

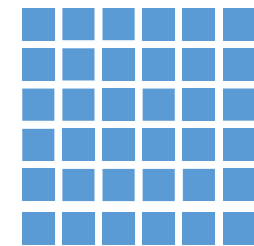
- Computation hungry

- FCLs: computation

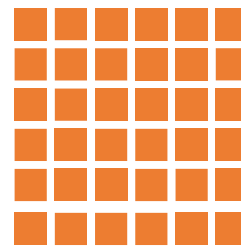
- Memory bandwidth h

Either PEs or memory bandwidth will be underutilized, thereby causing performance degradation

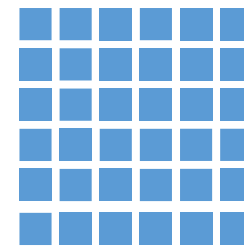
Processing Unit  
(Monolithic  
Architectures)



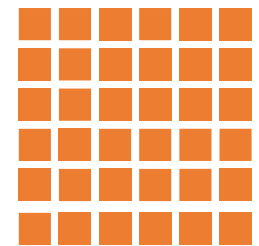
Processing  
CVLs



Processing  
FCLs



Processing  
CVLs



Processing  
FCLs

PE utilization



BW utilization



# Complementary effects

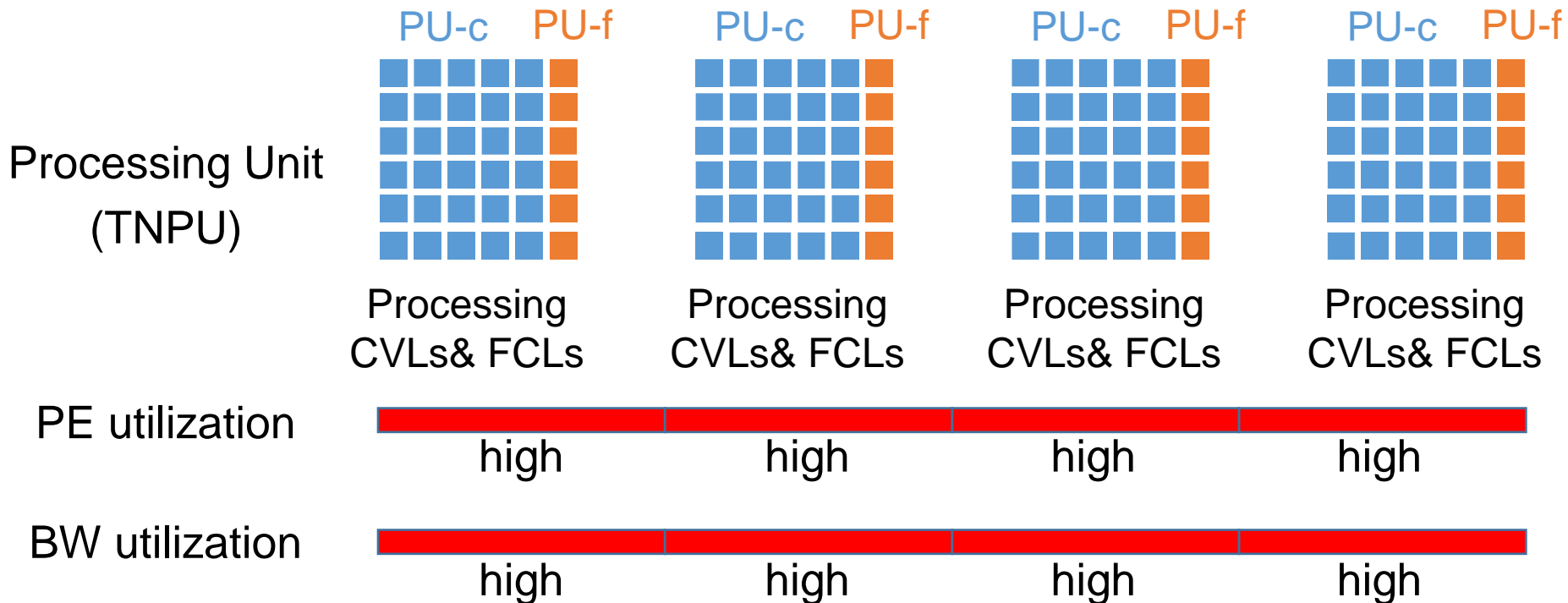
---

- Resource requirement diversity
  - CVLs: more PE, less BW
  - FCLs: more BW, less PE
- Orchestrate CVLs and FCLs to be processed concurrently
  - Underutilized PEs in FCLs can be leveraged by CVLs
  - Underutilized BW in CVLs can be leveraged by FCLs



# Asymmetric tandem architecture

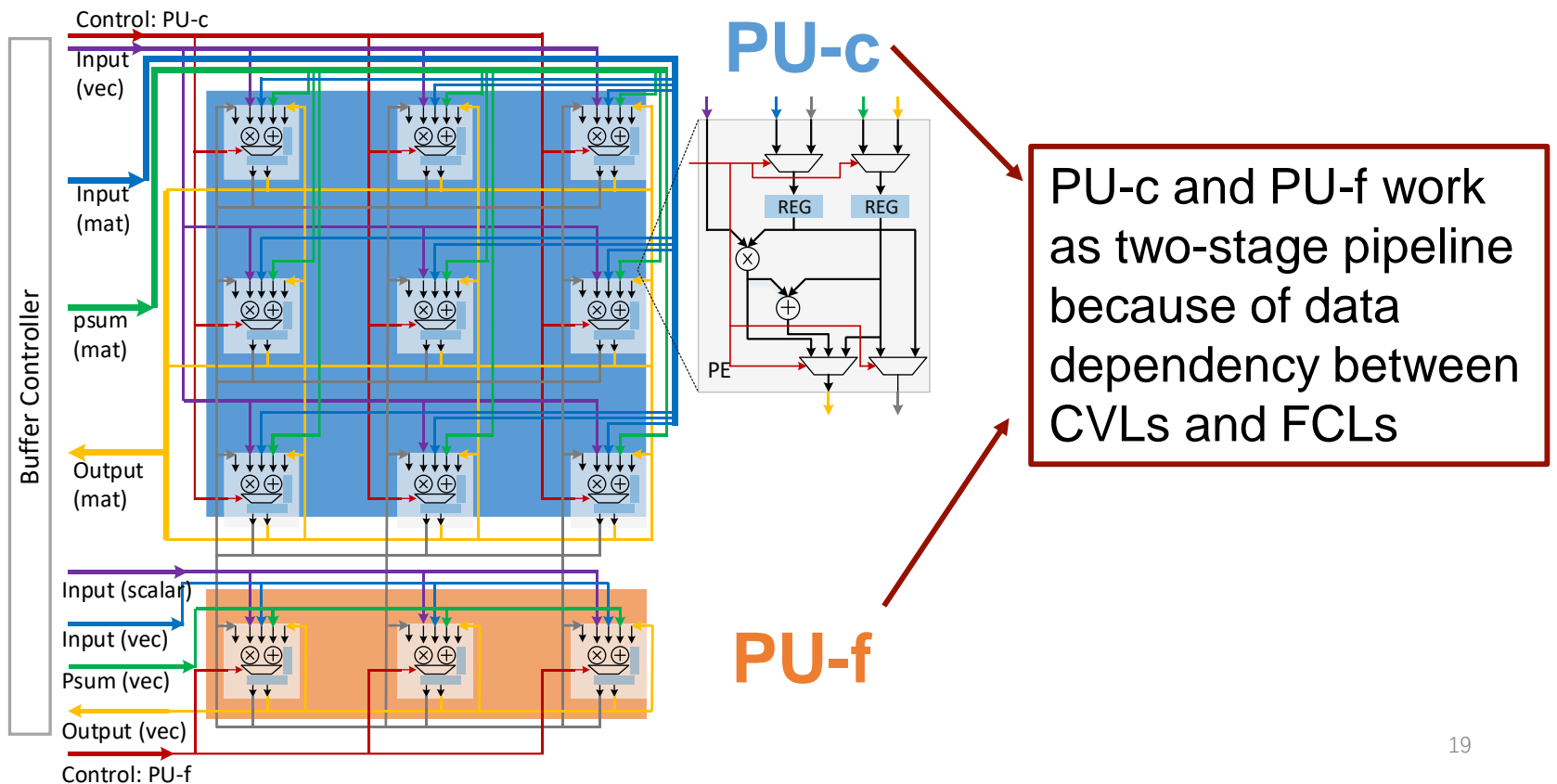
- PU: Space division multiplexing
  - Exploits complementary effect of resource requirement between CVLs and FCLs to boost utilization



# Asymmetric tandem architecture

- PE micro-architecture

- Virtualize PU-c and PU-f from a 2D mesh PE array
- PEs connect with each other via a NoC



# Bidirectional dataflow

- Data dependency:  $CVL \rightarrow FCL \rightarrow CVL$



Cycle	T0	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10
0-FP	0-FP-C	0-FP-F									
0-BP											
0-GD											
1-FP											
1-BP											
1-GD											
2-FP											
States	Idle	Idle									

FP-forward propagation, BP-backward propagation, GD-gradient descent

# Bidirectional dataflow

- Data dependency:  $CVL \rightarrow FCL \rightarrow CVL$



Cycle	T0	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10
0-FP	0-FP-C	0-FP-F									
0-BP			0-BP-F	0-BP-C							
0-GD											
1-FP											
1-BP											
1-GD											
2-FP											
States	Idle	Idle	Idle	Idle							

FP-forward propagation, BP-backward propagation, GD-gradient descent

# Bidirectional dataflow

- Data dependency:  $CVL \rightarrow FCL \rightarrow CVL$

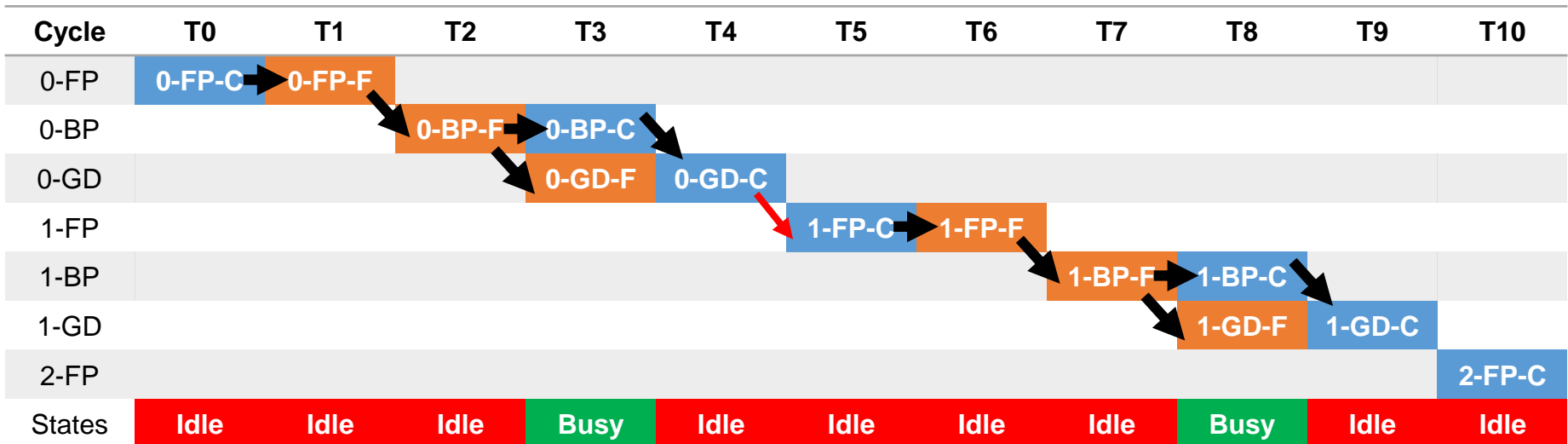
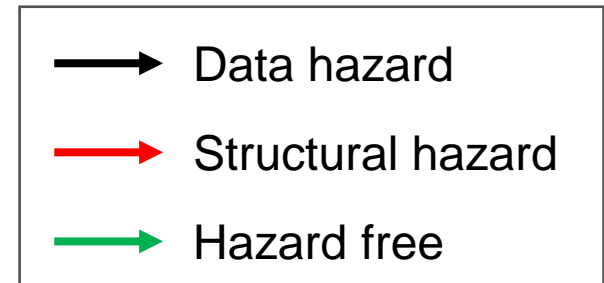
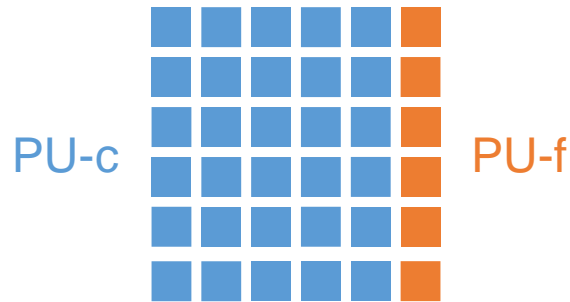


Cycle	T0	T1	T2	T3	T4	T5	T6	T7	T8	T9	T10
0-FP	0-FP-C	0-FP-F									
0-BP			0-BP-F	0-BP-C							
0-GD				0-GD-F	0-GD-C						
1-FP											
1-BP											
1-GD											
2-FP											
States	Idle	Idle	Idle	Busy	Idle						

FP-forward propagation, BP-backward propagation, GD-gradient descent

# Bidirectional dataflow

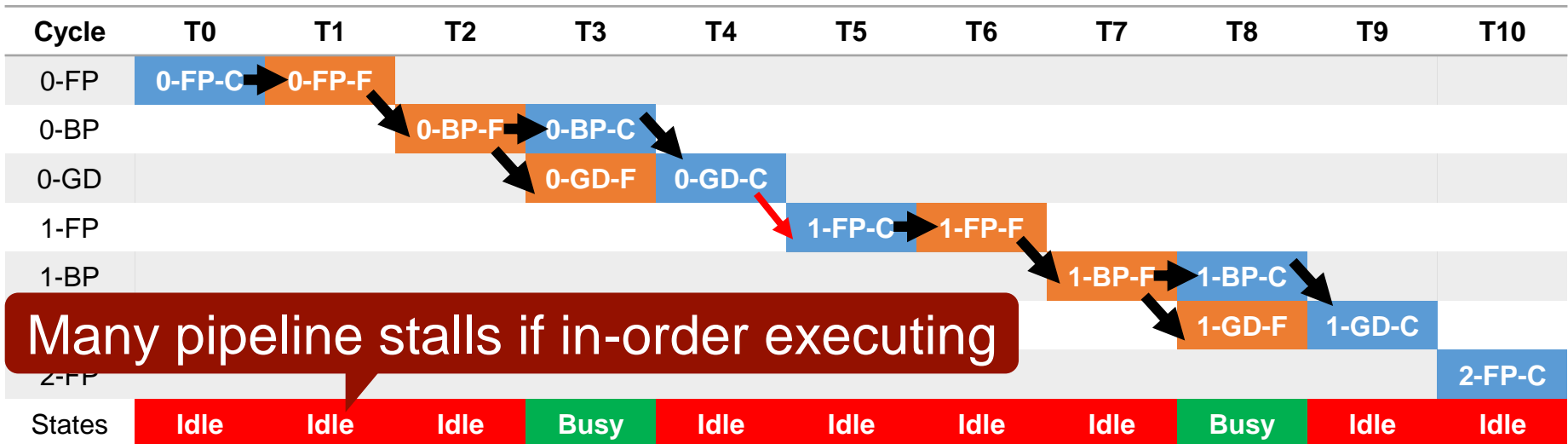
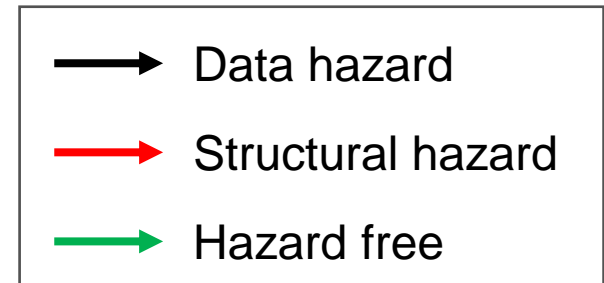
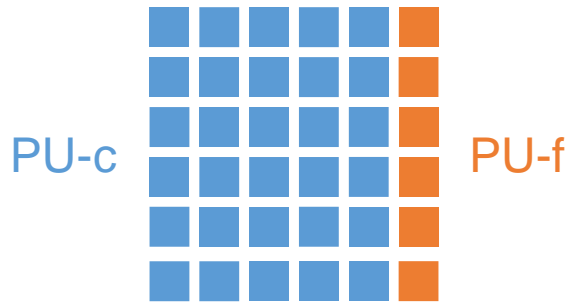
- Data dependency:  $CVL \rightarrow FCL \rightarrow CVL$



FP-forward propagation, BP-backward propagation, GD-gradient descent

# Bidirectional dataflow

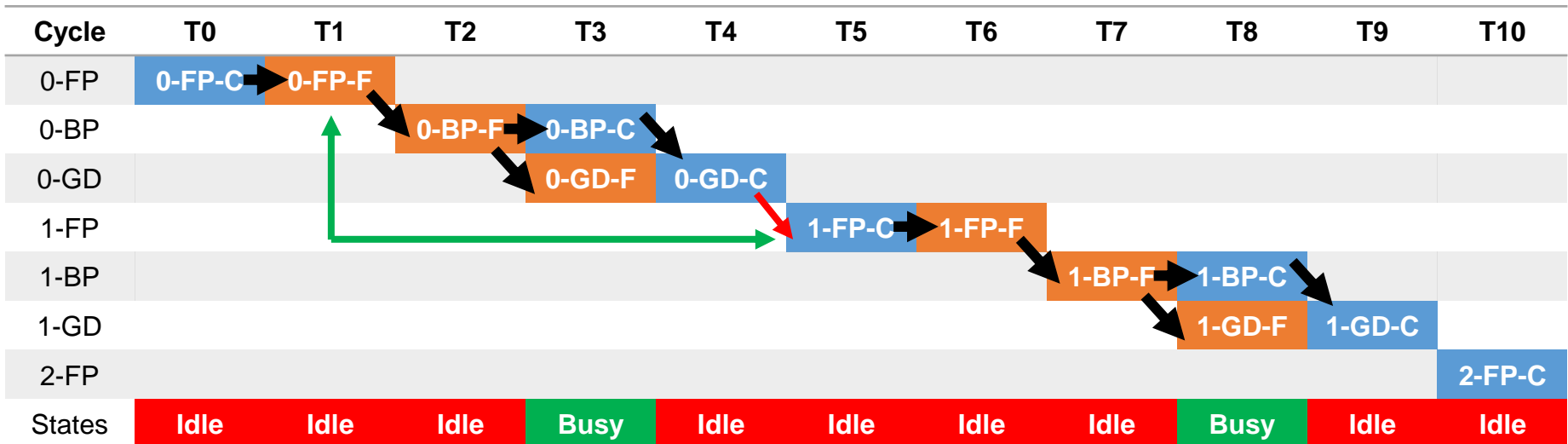
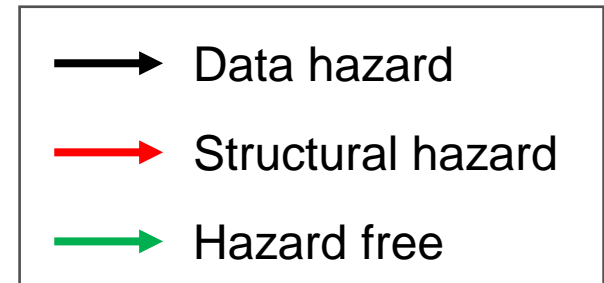
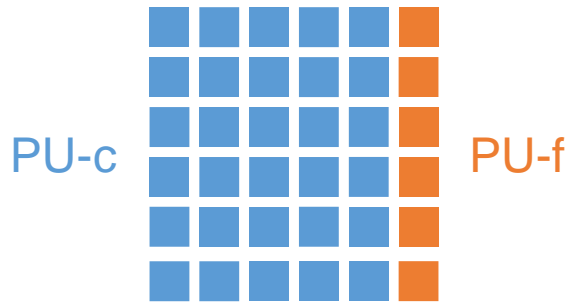
- Data dependency:  $CVL \rightarrow FCL \rightarrow CVL$



FP-forward propagation, BP-backward propagation, GD-gradient descent

# Out-of-order scheduling

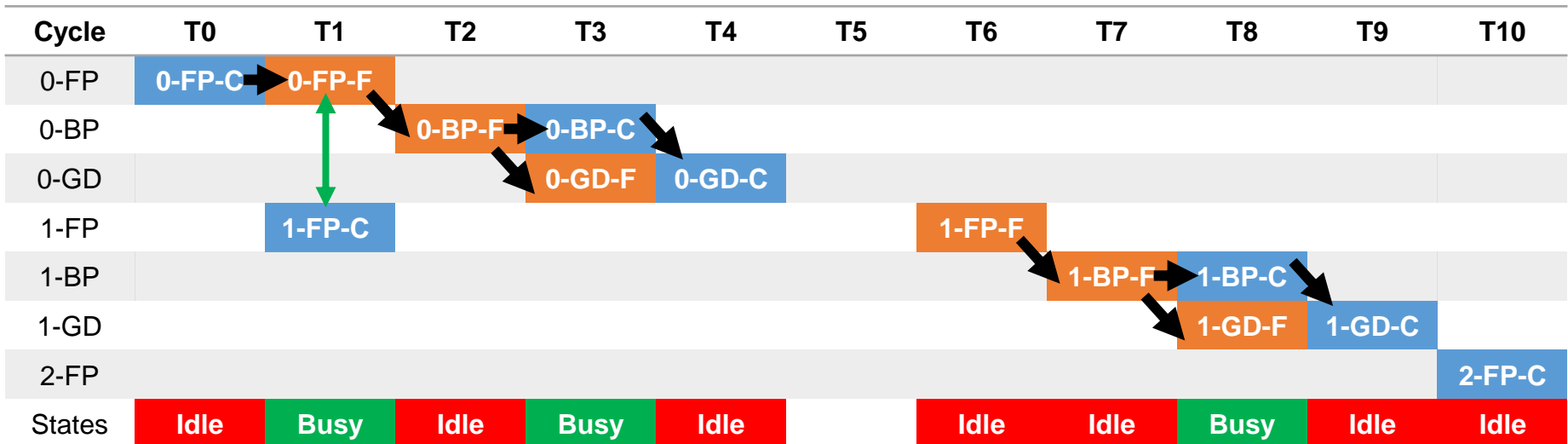
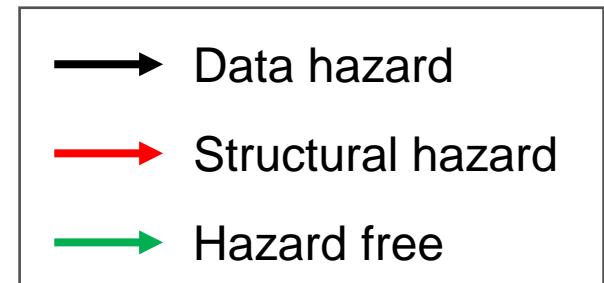
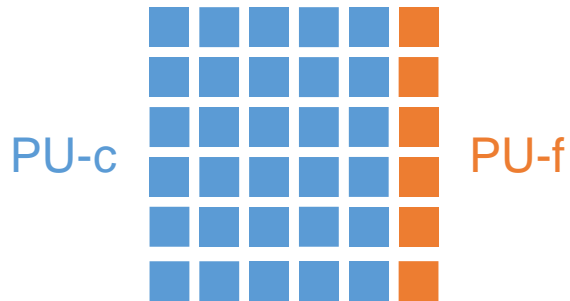
- Forwarding to eliminate hazard





# Out-of-order scheduling

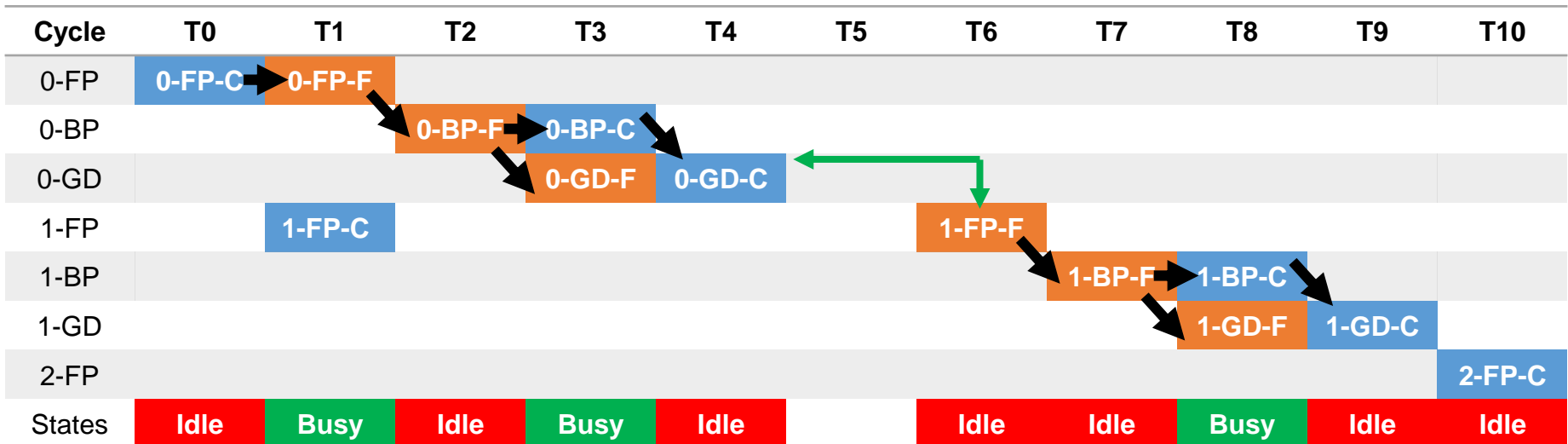
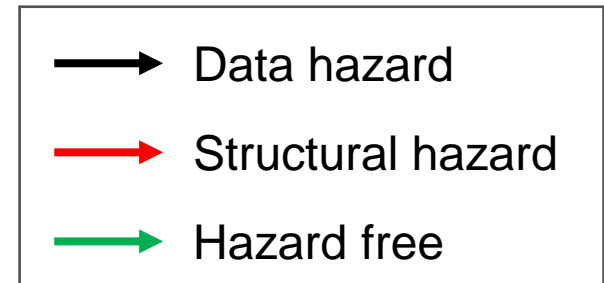
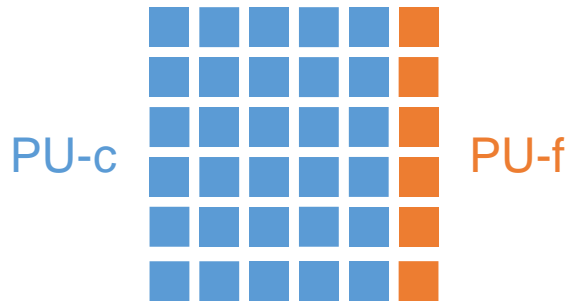
- Forwarding to eliminate hazard



FP-forward propagation, BP-backward propagation, GD-gradient descent

# Out-of-order scheduling

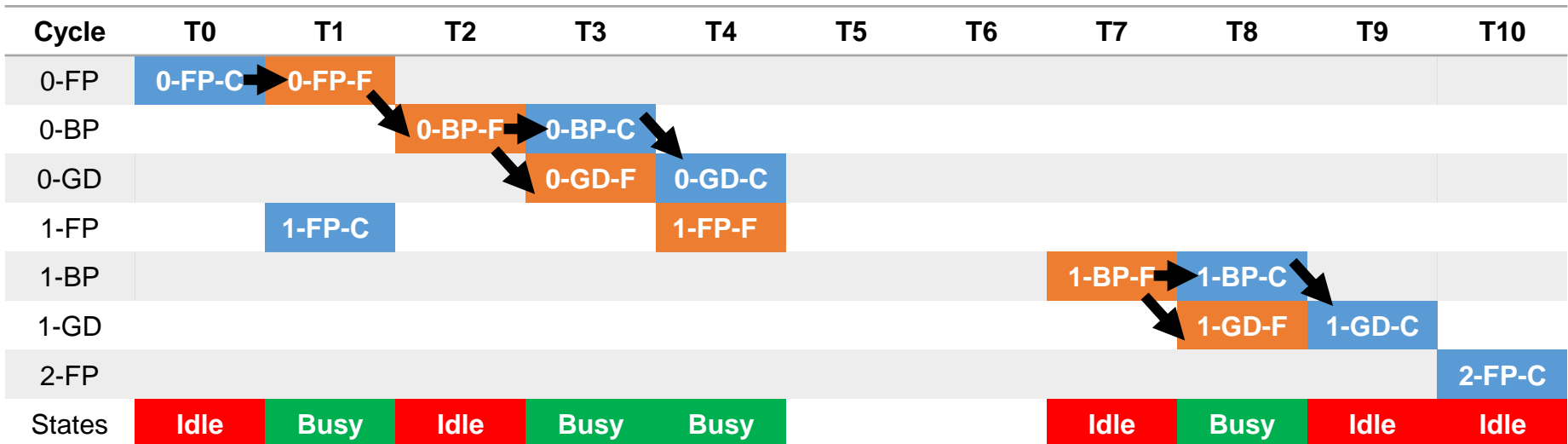
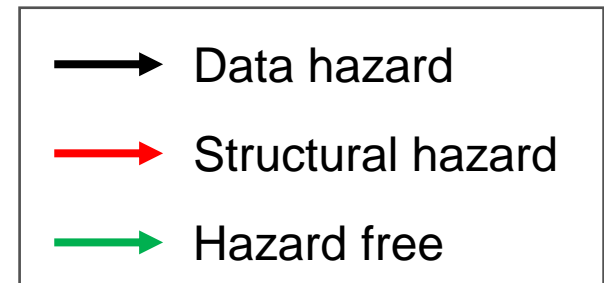
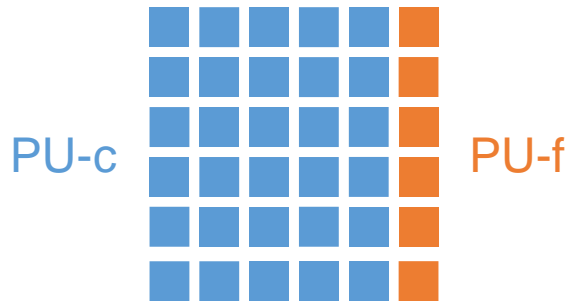
- Forwarding to eliminate hazard



FP-forward propagation, BP-backward propagation, GD-gradient descent

# Out-of-order scheduling

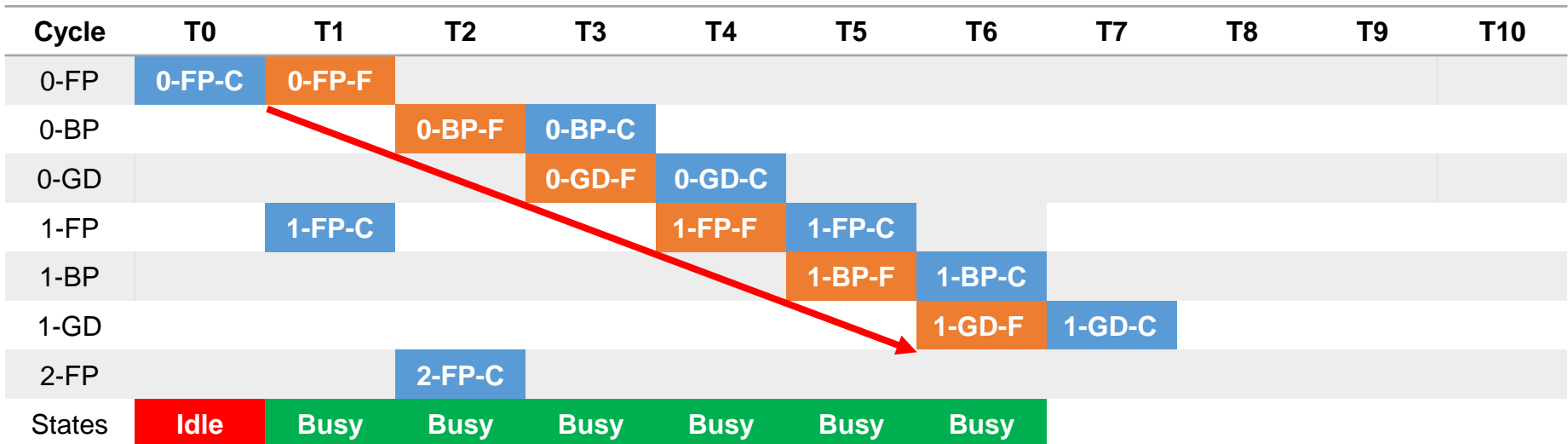
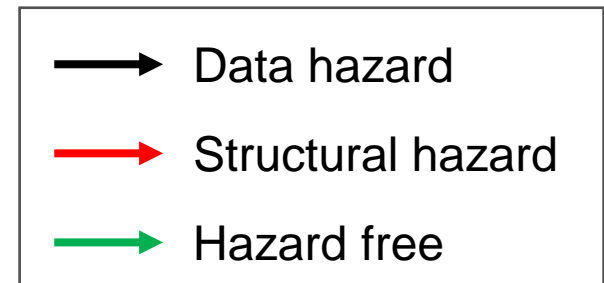
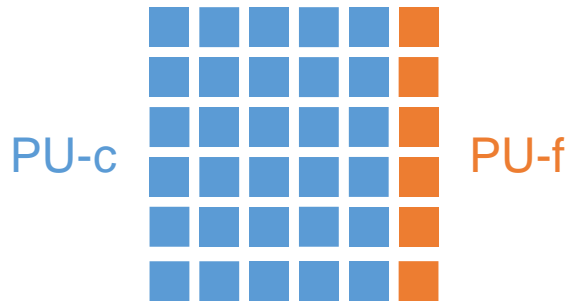
- Forwarding to eliminate hazard



FP-forward propagation, BP-backward propagation, GD-gradient descent

# Out-of-order scheduling

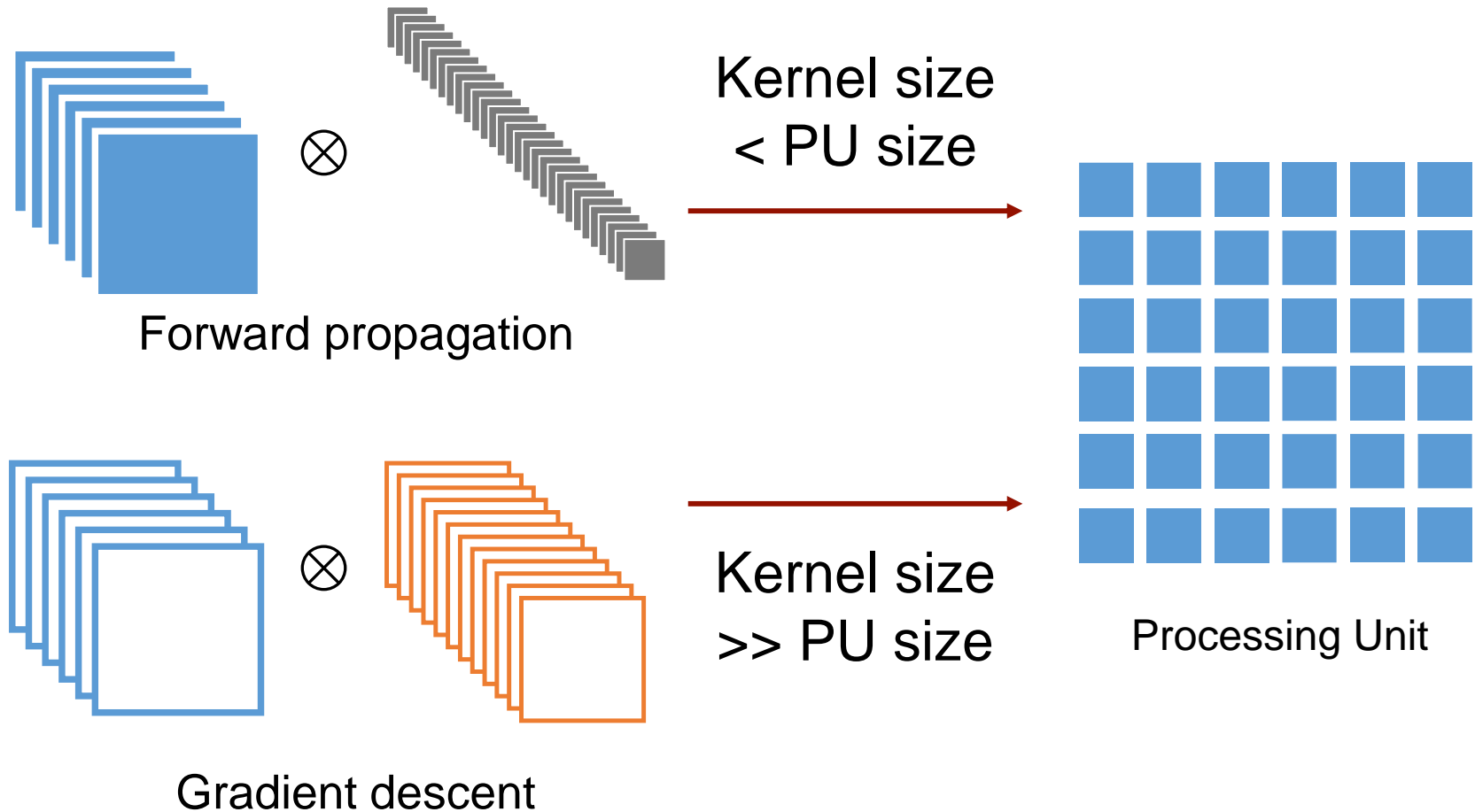
- Forwarding to eliminate hazard



FP-forward propagation, BP-backward propagation, GD-gradient descent

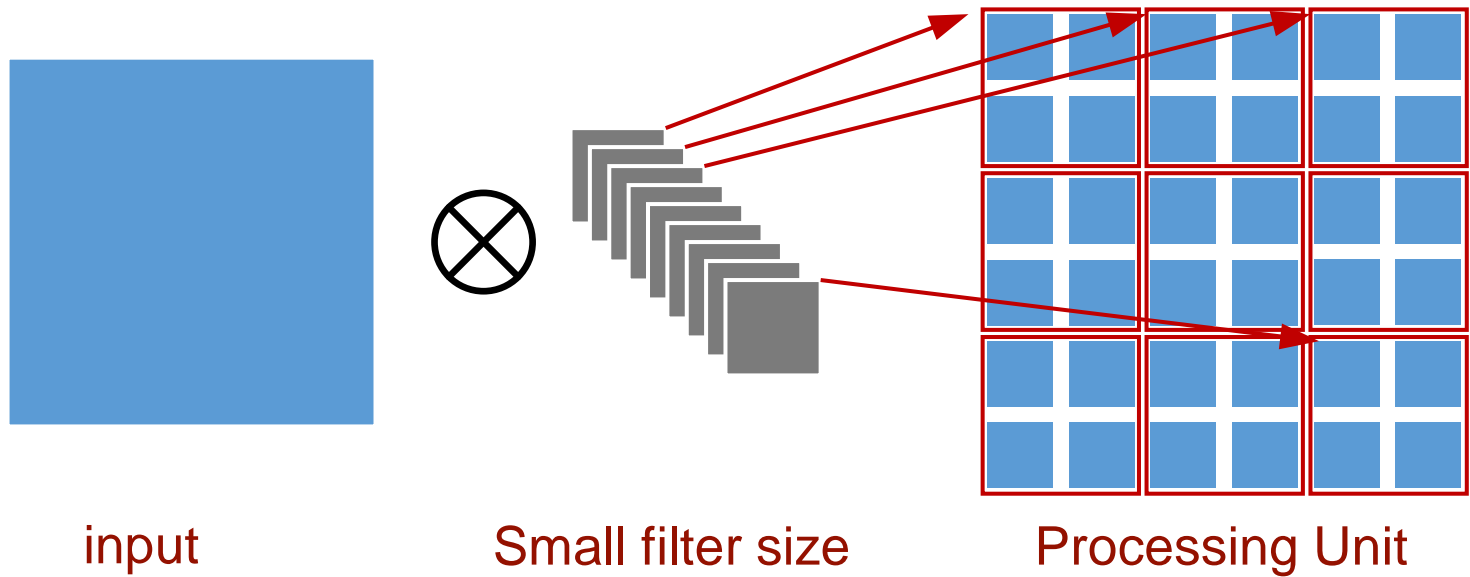
# Wide-range of convolution kernels

- Flexible convolution allocation



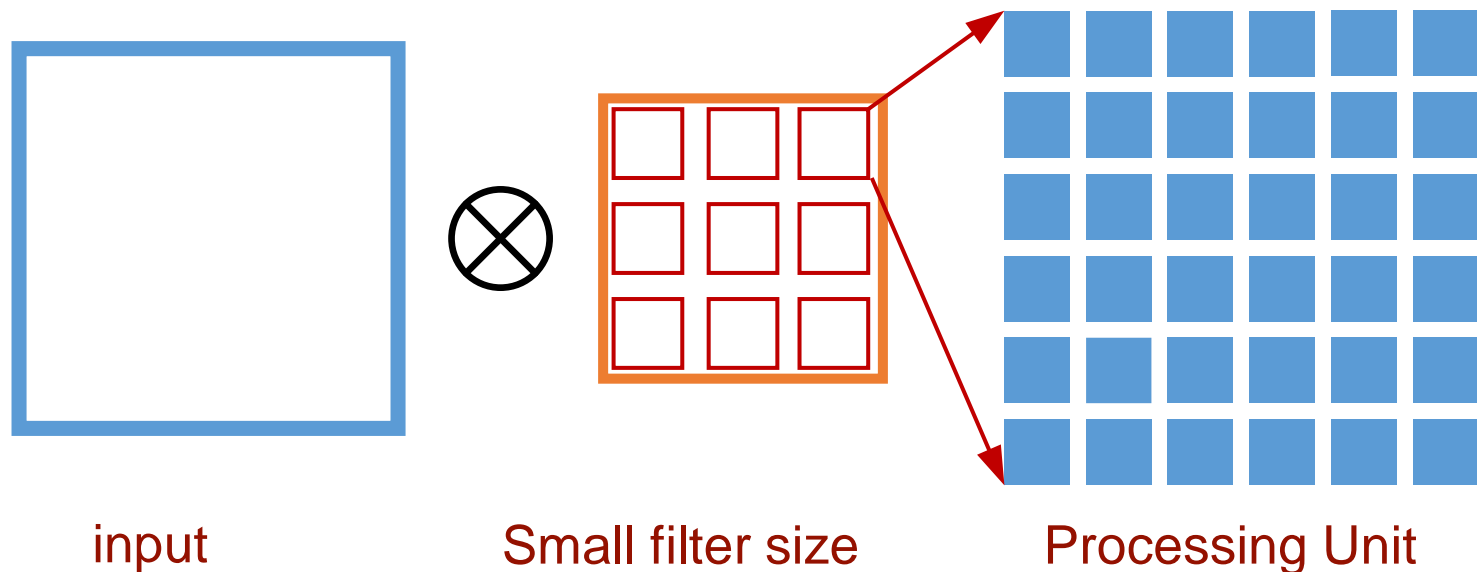
# Flexible convolution allocation

- Kernel size < PU size: kernel aggregation
  - Aggregate multiple kernels to generate a big one, then map to Processing Unit



# Flexible convolution allocation

- Kernel size > PU size: kernel partitioning
  - Partition large kernels into small sub-kernels [DATE'2018]



# Experimental setup

- Evaluation platform

- RTL-level cycle-accurate simulation for performance estimation
- System-level energy estimation, based on synthesis and CACTI

- Platform configurations

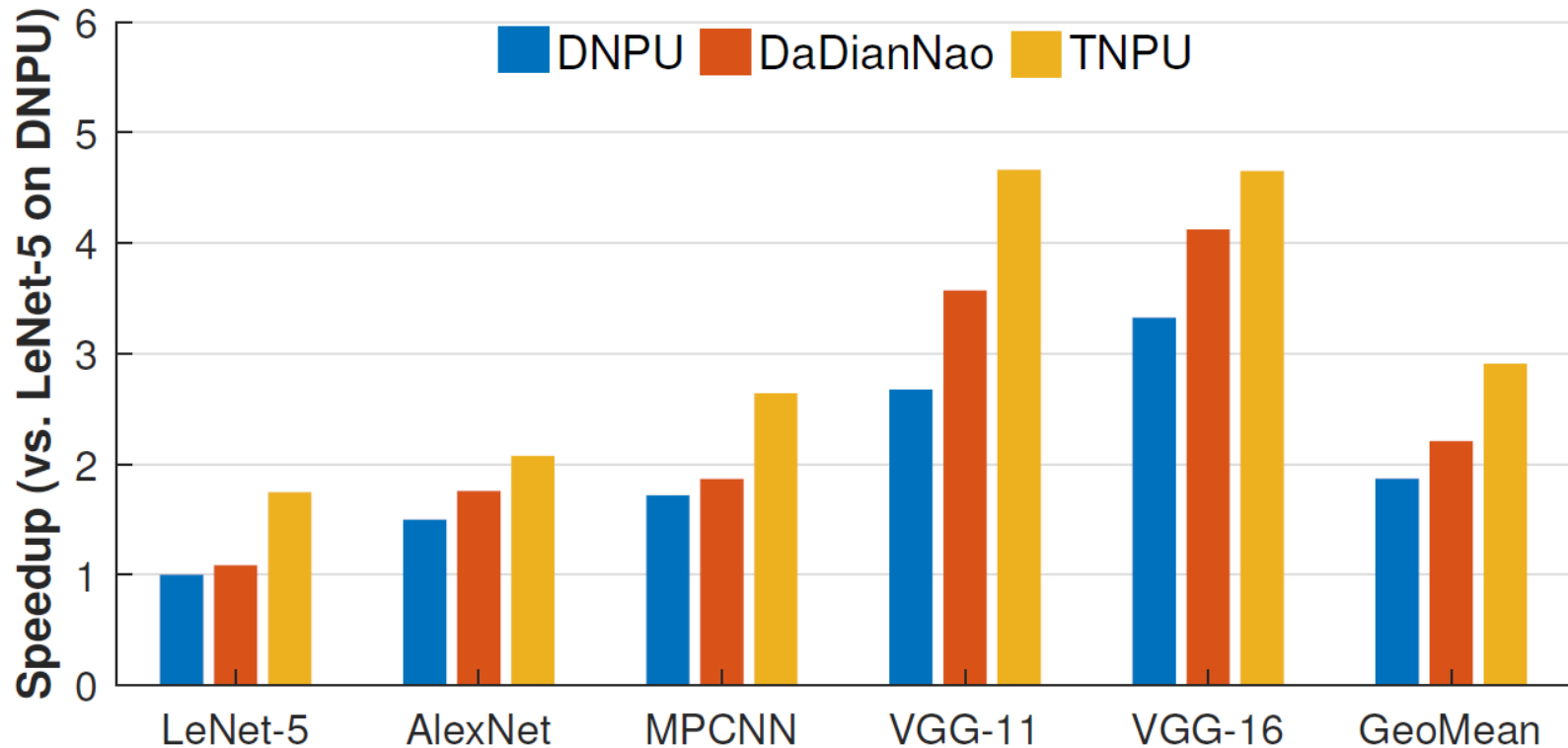
---

	<b>MACs</b>	<b>SRAM</b>	<b>Clock</b>	<b>DRAM bandwidth</b>
TNPU	272	256KB	400MHz	17.1GB/s
Dadiannao-re	272	256KB	400MHz	17.1GB/s
DNPU-re	272	256KB	400MHz	17.1GB/s



# Experimental results

- Performance

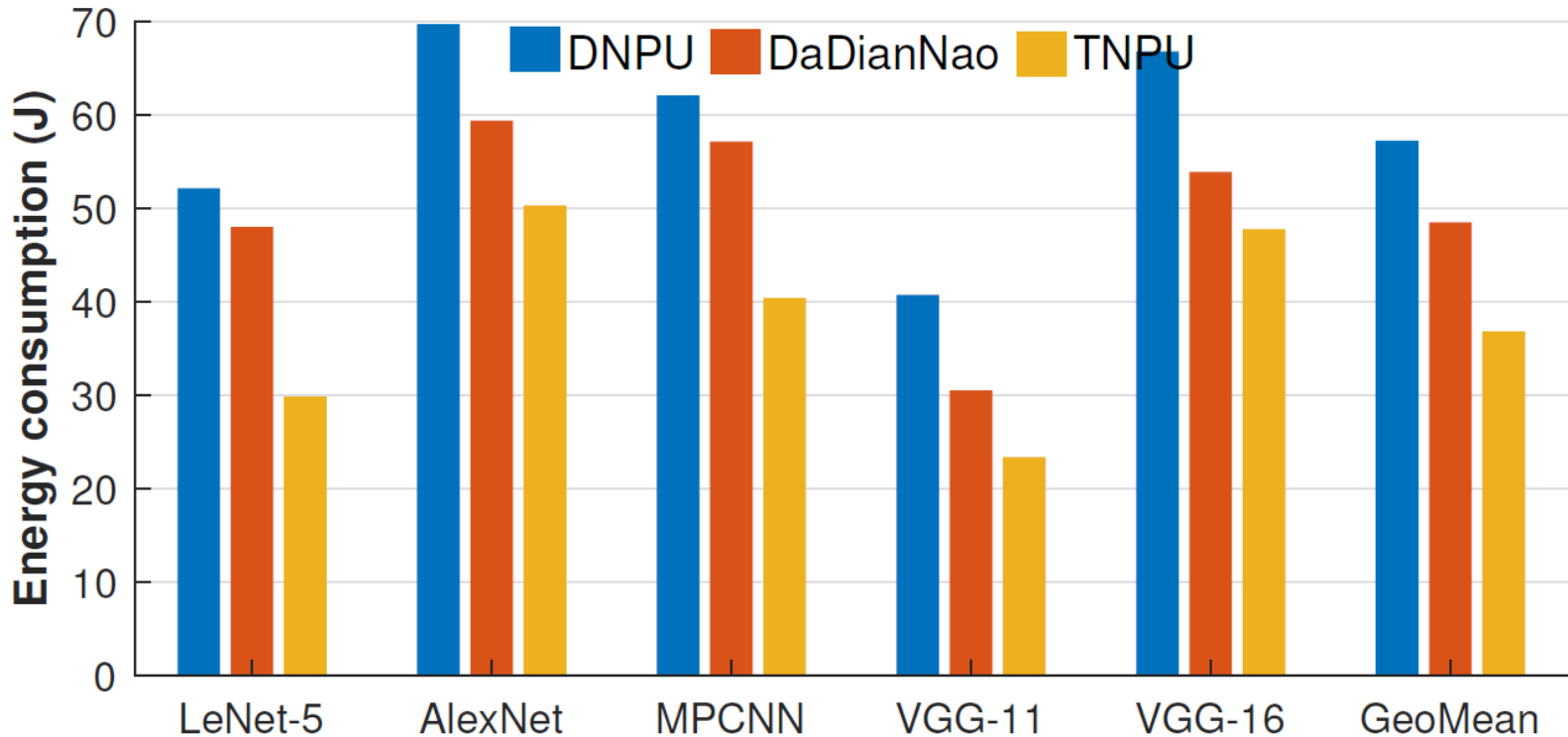


55% performance improvement over DNPU

32% performance improvement over Dadiannao

# Experimental results

- Energy consumption

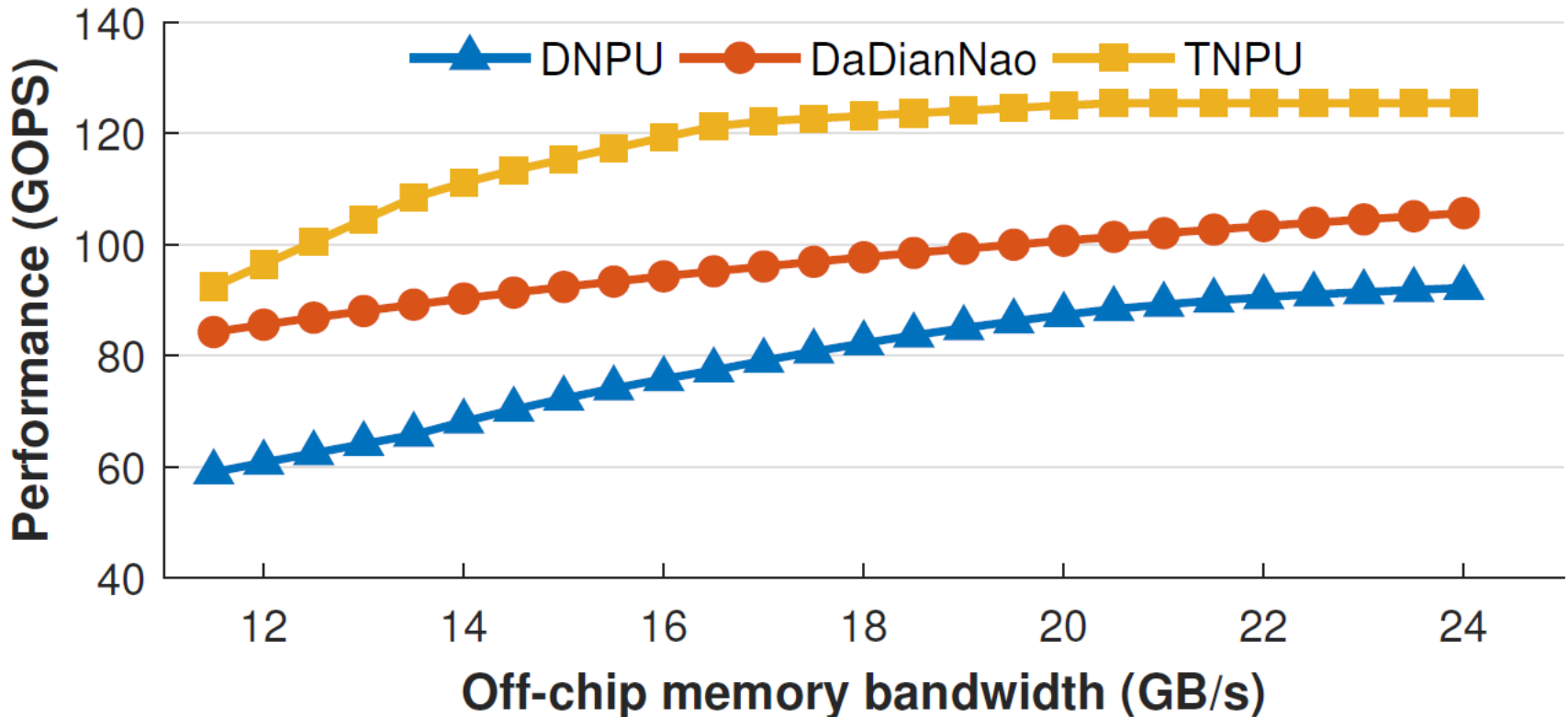


35% energy saving compared to DNPU

24% energy saving compared to Dadiannaao

# Experimental results

- Scalability to DRAM bandwidth



TNPU maintains a stable high performance and consistently outperforms DNPU and Dadiannaio

# Conclusion

---

- TNPU, an accelerator architecture for CNN training:
  - 55% and 32% performance improvement over DNPU and Dadiannao, respectively
- TNPU addresses the following challenges in CNN training:
  - Diversity between CVLs and FCLs
  - Bidirectional data dependency
  - Extremely large convolutional kernels



# Thanks for your attention!

Email: [lijiajun@ict.ac.cn](mailto:lijiajun@ict.ac.cn)

