

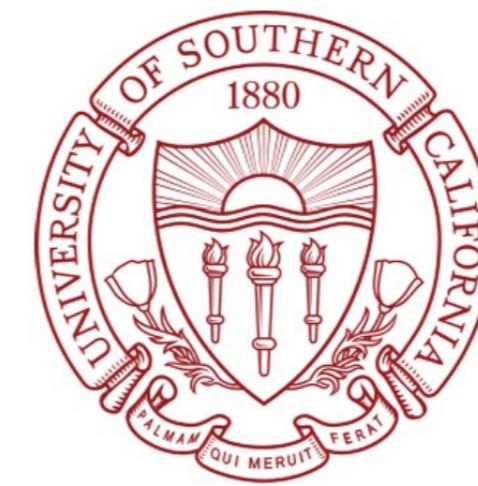
# HyPar: Towards Hybrid Parallelism for Deep Learning Accelerator Array

Linghao Song\*, Jiachen Mao\*, Youwei Zhuo<sup>#</sup>,  
Xuehai Qian<sup>#</sup>, Hai Li\*, Yiran Chen\*

*\*Duke University*

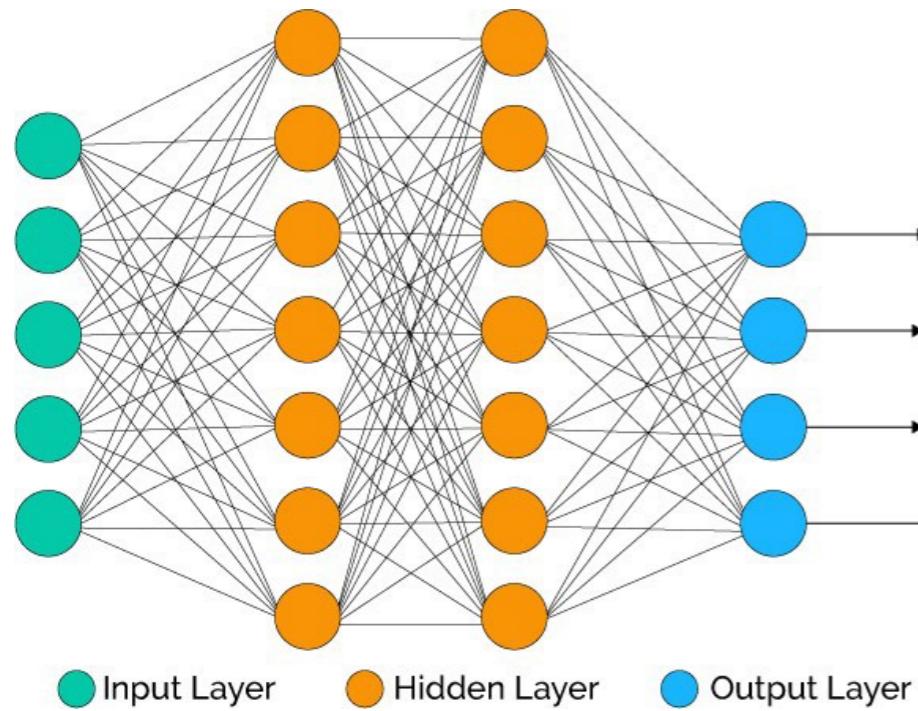
*# University of Southern California*

**CEI**  
[cei.pratt.duke.edu](http://cei.pratt.duke.edu)



**ALCHEM**  
[alchem.usc.edu](http://alchem.usc.edu)

# As a computer architect we know deep learning applications in various domains...

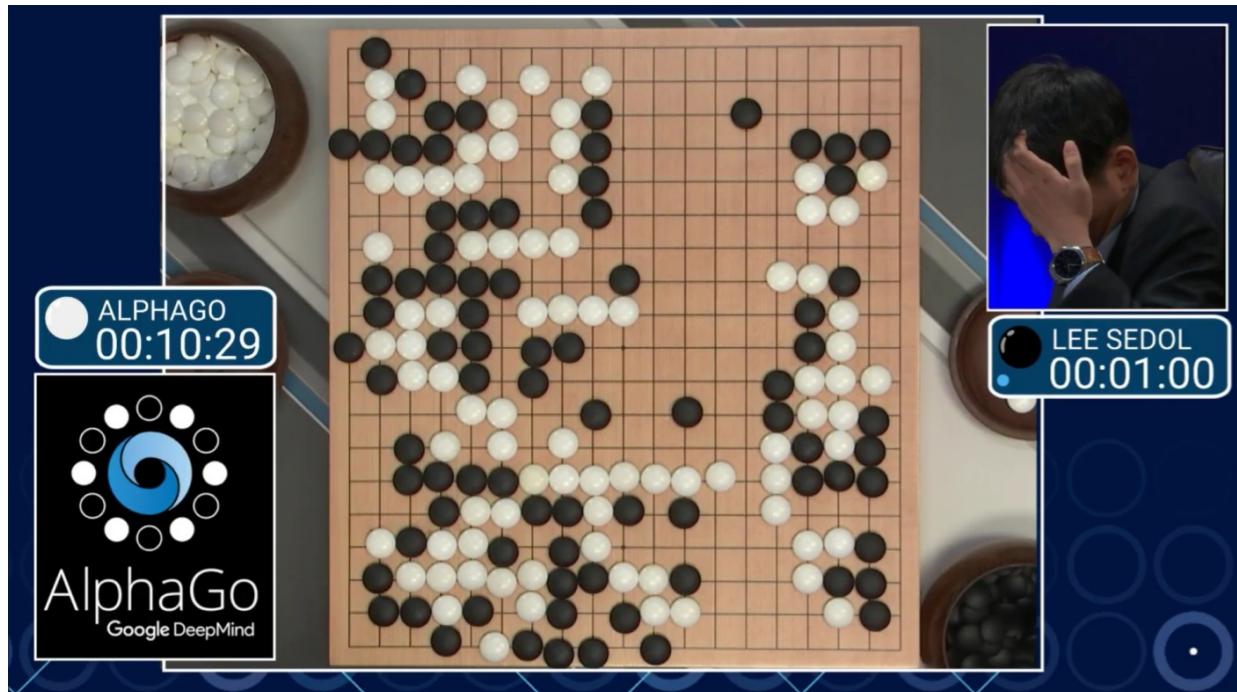
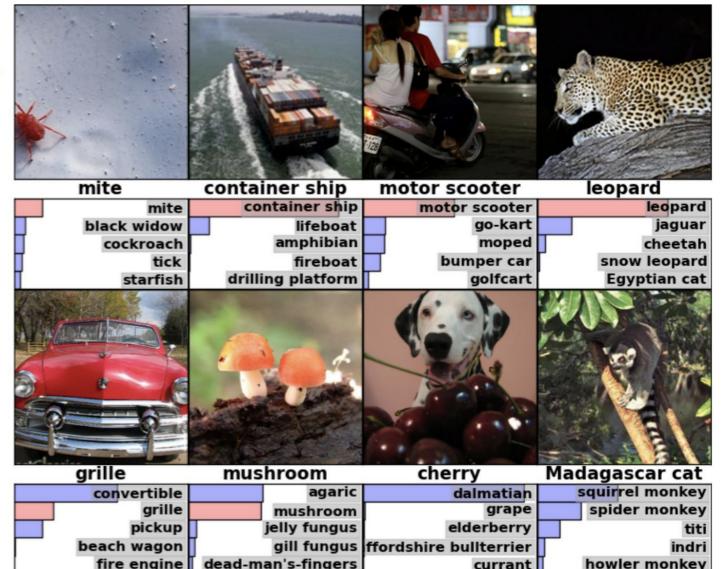


## ImageNet Challenge

(IJCV'15)

IMAGENET

- 1,000 object classes (categories).
- Images:
  - 1.2 M train
  - 100k test.



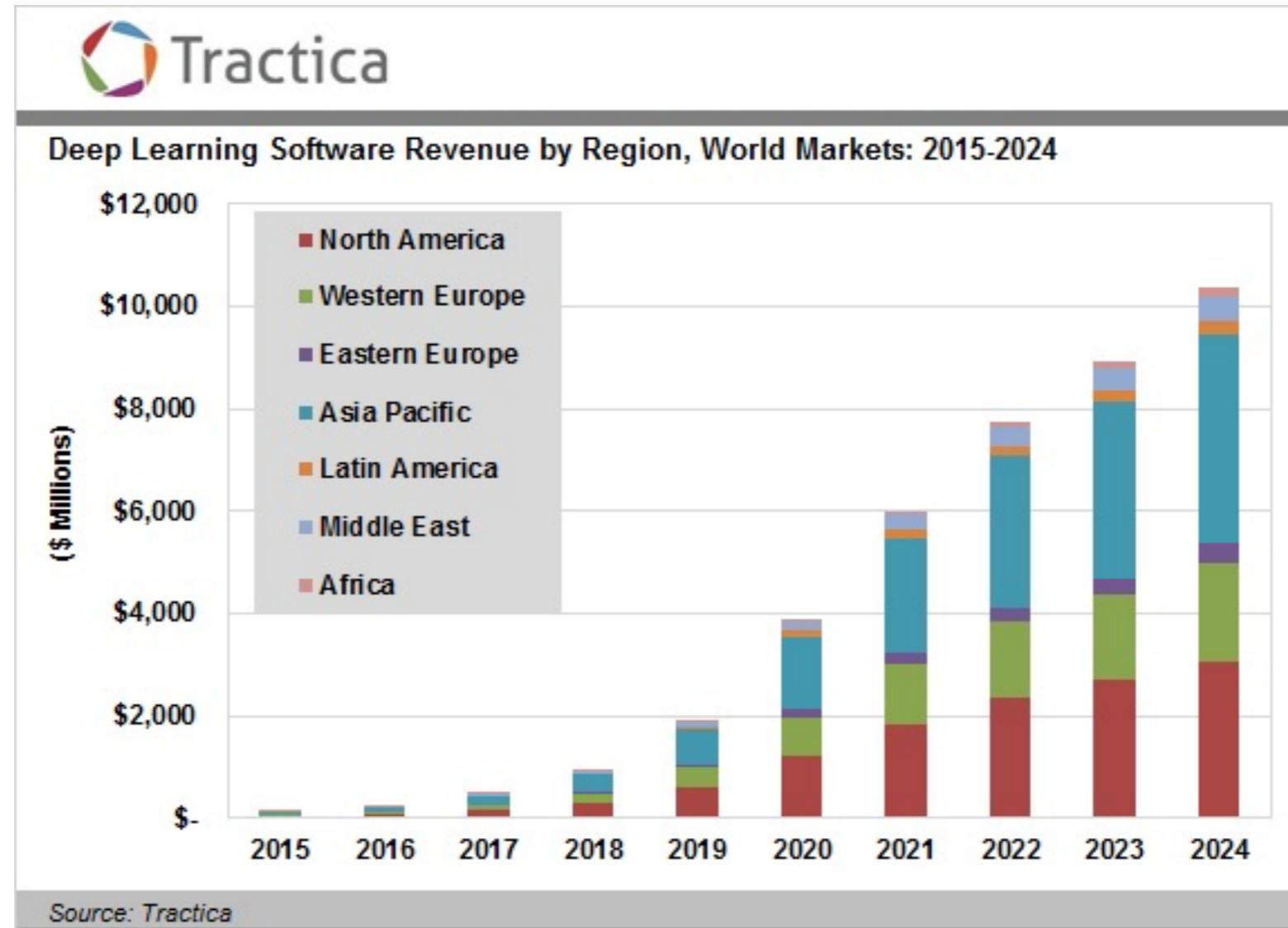
(AlphaGo, DeepMind)



(Biggan, ICLR'19)

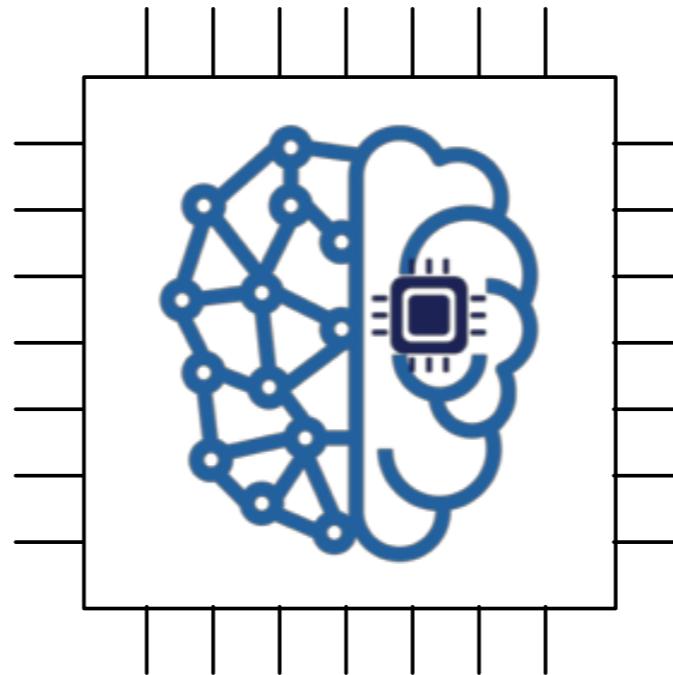
# We also know deep learning is rapidly growing..

---

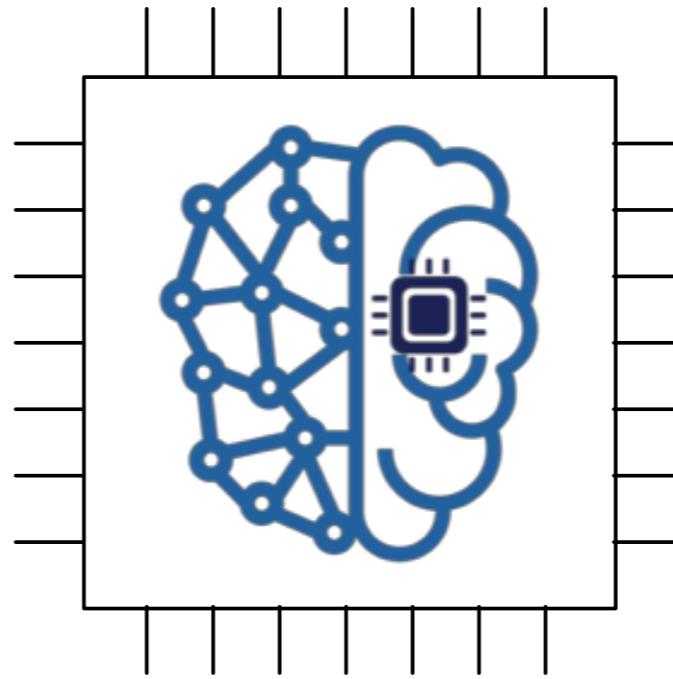


So we agree it is an application domain for specialized architecture.

# Current approach: NN->Chip, Mapping, Dataflow, etc.



Current approach: NN->Chip, Mapping, Dataflow, etc.



Current approach: NN->Chip, Mapping, Dataflow, etc.

But...

- That is not new, even boring...
- Morning Session Warning: Accelerator Wall
- So many DL accelerators...
- Why HyPar? What's next for DL accelerators?

# Outline

---

- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

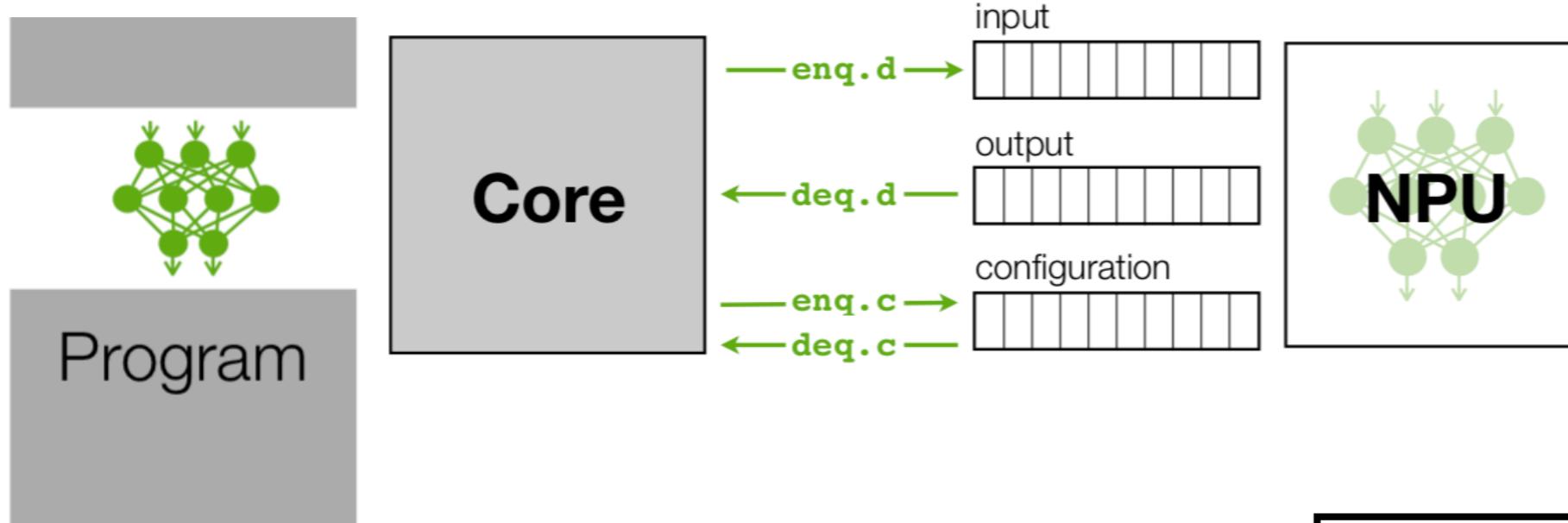
# Outline

---

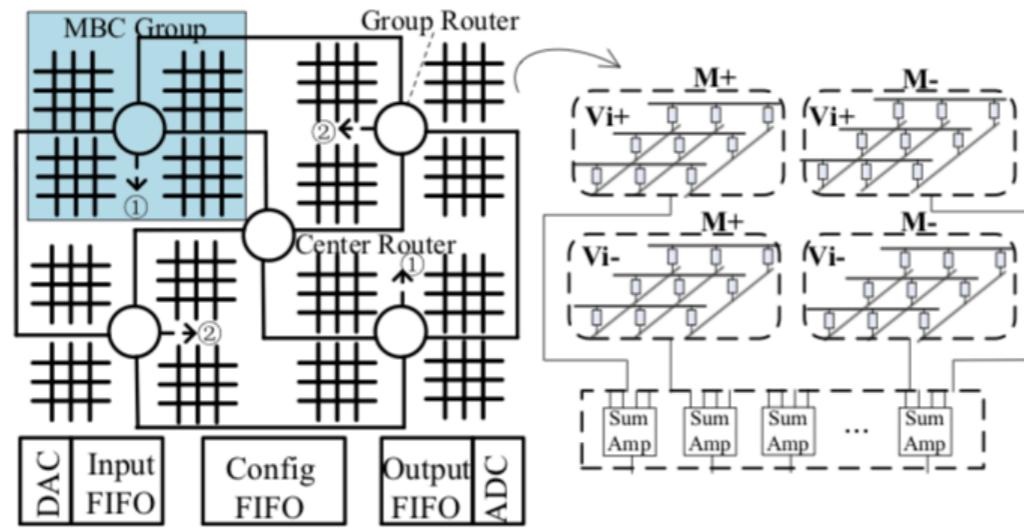
- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

# DLA early stage: on-chip design

- NPU: To run **some program segment** on NPU (neural processing unit) instead of CPU (MICRO'12)



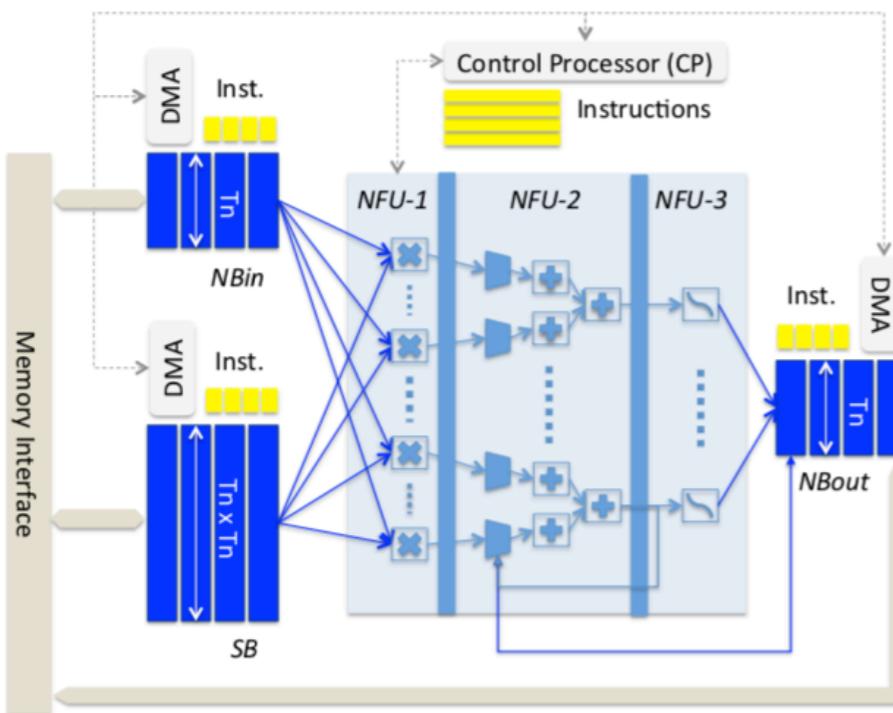
- Reno  
(DAC'15)



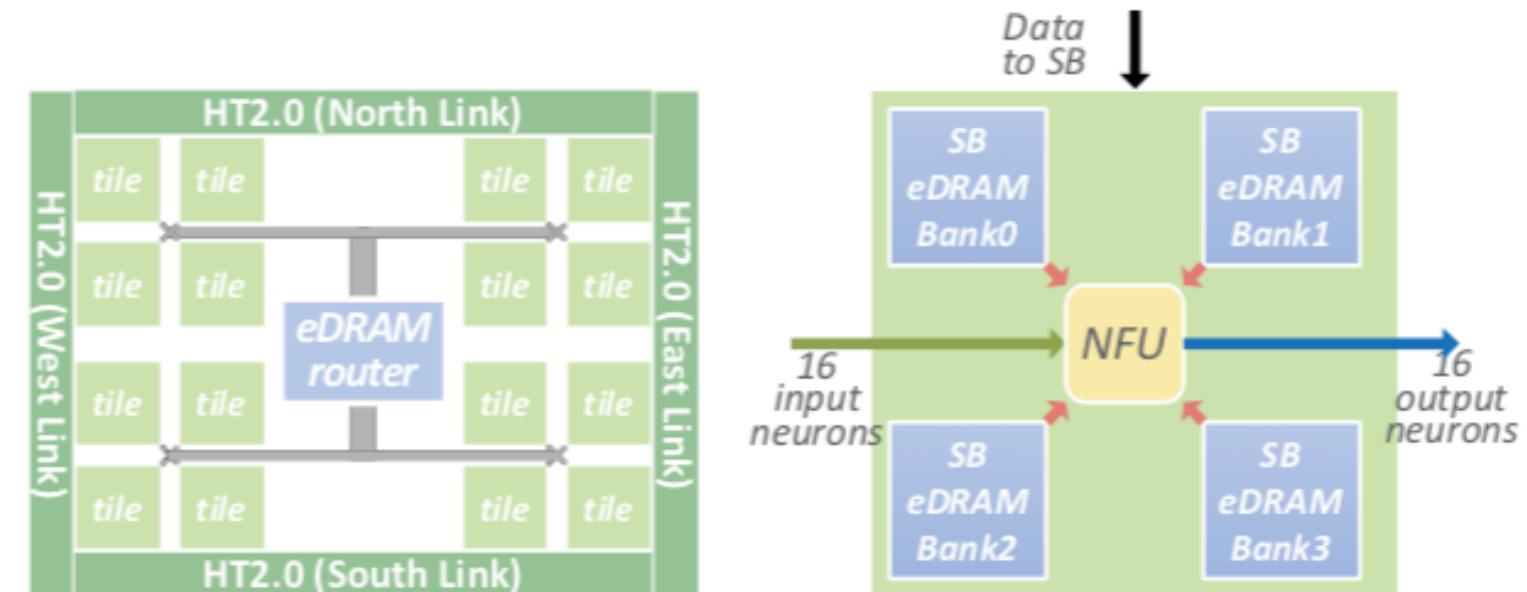
Limited application  
and performance

# DLA current stage: stand-alone accelerators

- Diannao (ASPLOS'14)



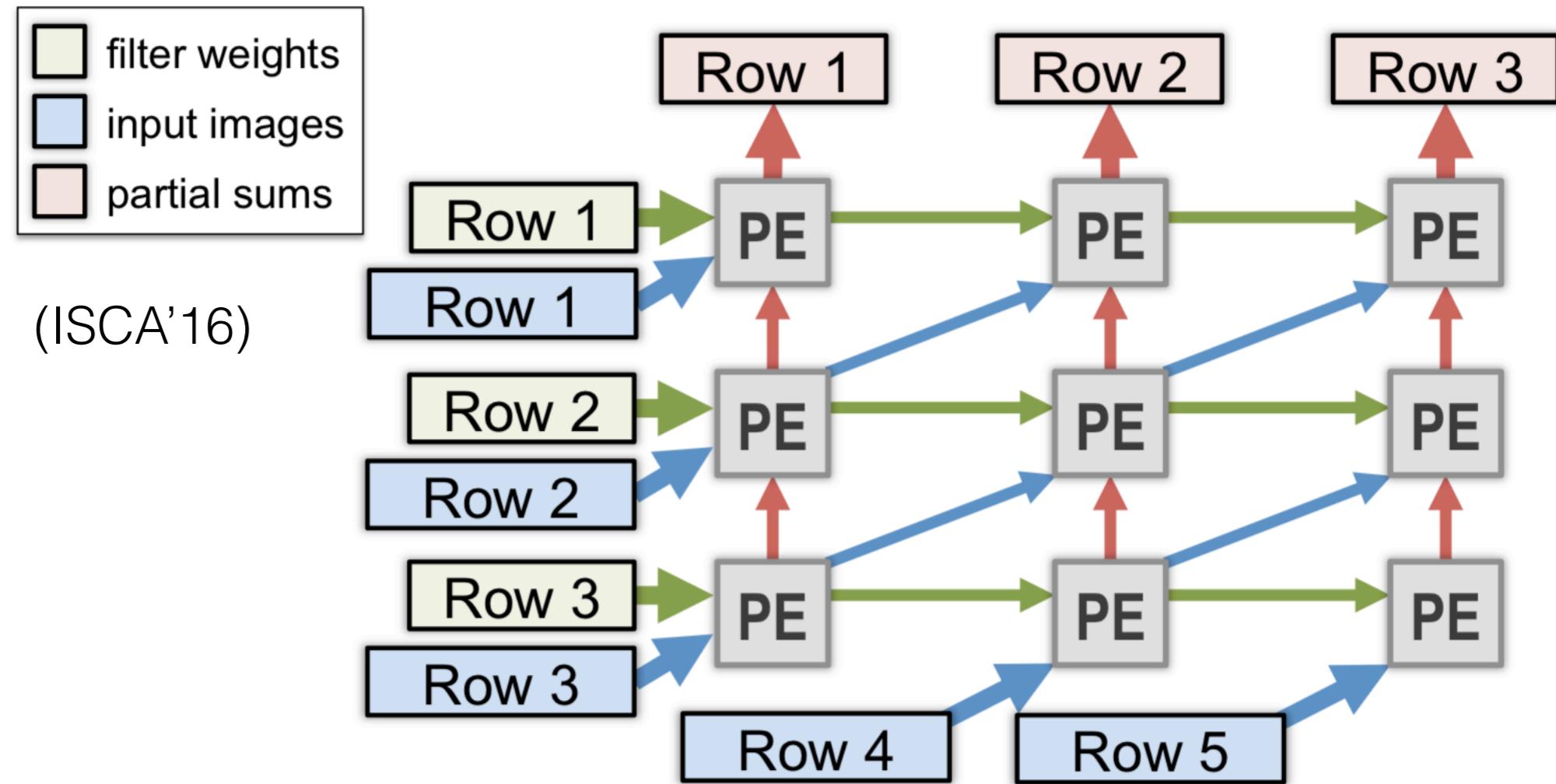
- DaDiannao (MICRO'14)



- ShiDiannao (ISCA'15), PuDiannao (ASPLOS'15), Cambricon-X (MICRO'16)

Support a whole application (NN) on an accelerator.

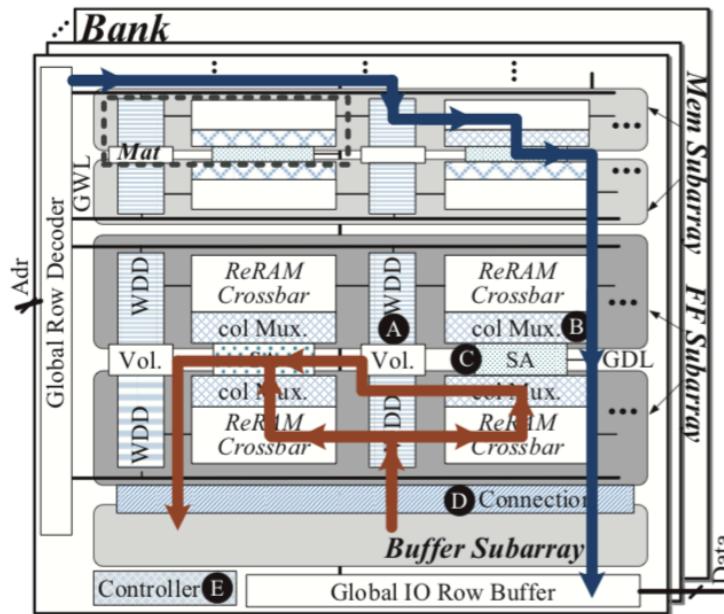
# DLA current stage: data flow architecture(Eyeriss)



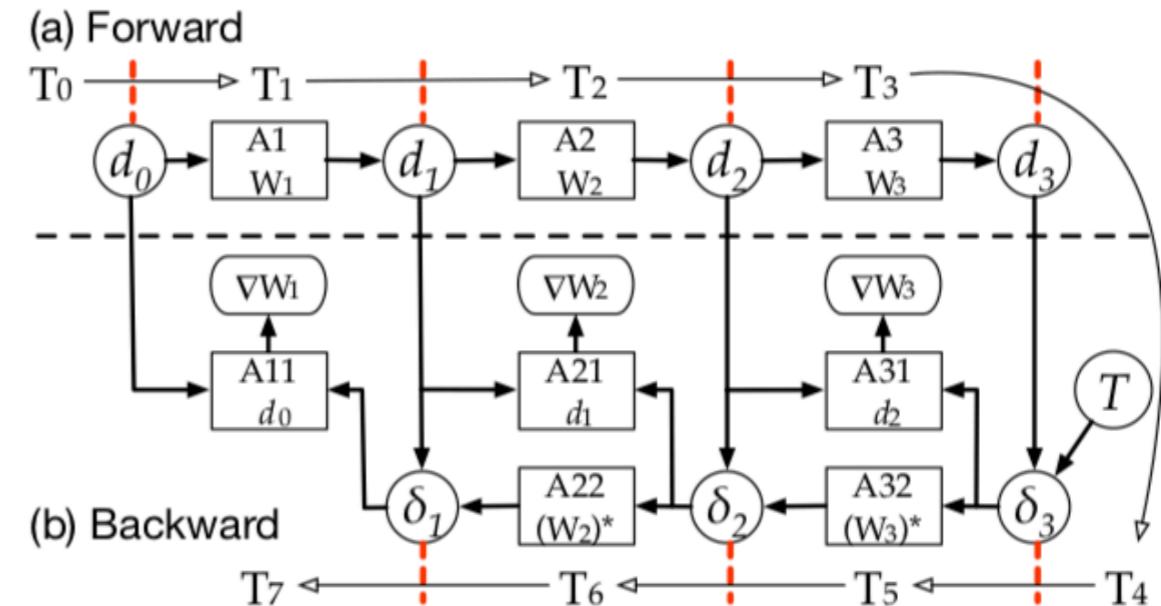
Fine exploration on data/resource mapping, sharing, reusing.

# DLA current stage: emerging memory(ReRAM)

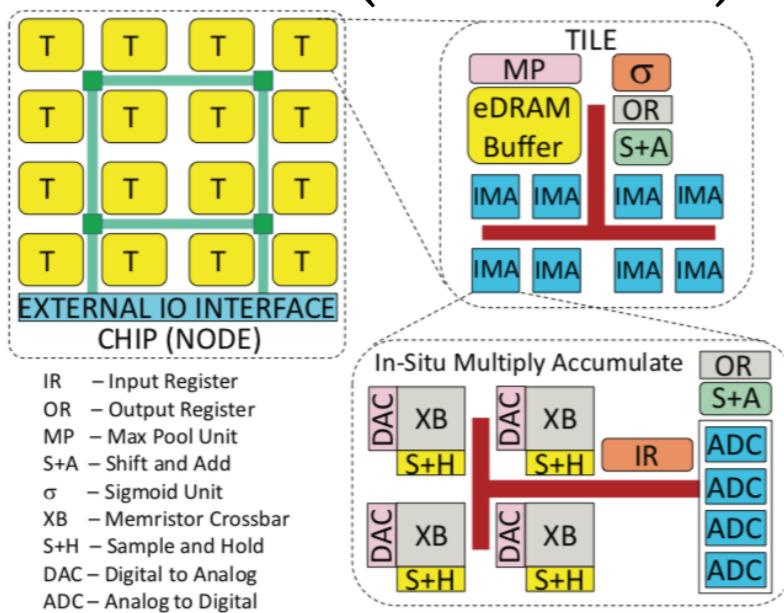
PRIME (ISCA'16)



PipeLayer (HPCA'17)



ISAAC (ISCA'16)



Take the benefit from emerging technologies beyond CMOS.

# Where are we in the development of DLAs?

---



# Where are we in the development of DLAs?

---



**CEI**

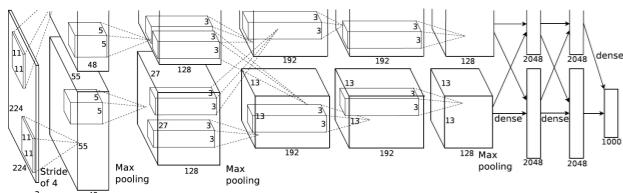
[cei.pratt.duke.edu](http://cei.pratt.duke.edu)

**ALCHEM**

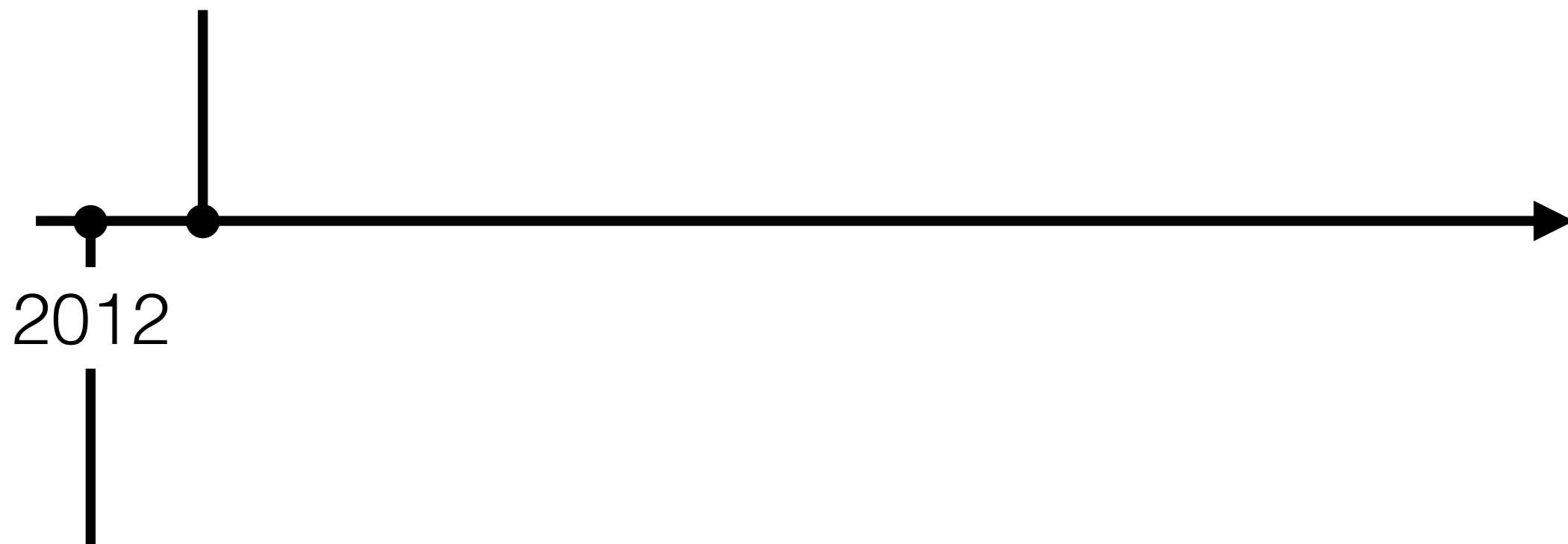
[alchem.usc.edu](http://alchem.usc.edu)

# Where are we in the development of DLAs?

---



AlexNet



NPU

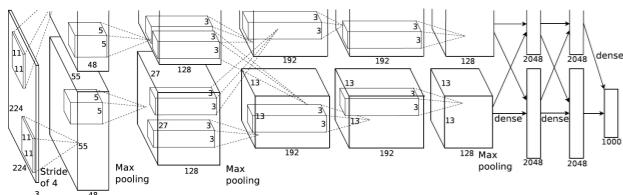
CEI

[cei.pratt.duke.edu](http://cei.pratt.duke.edu)

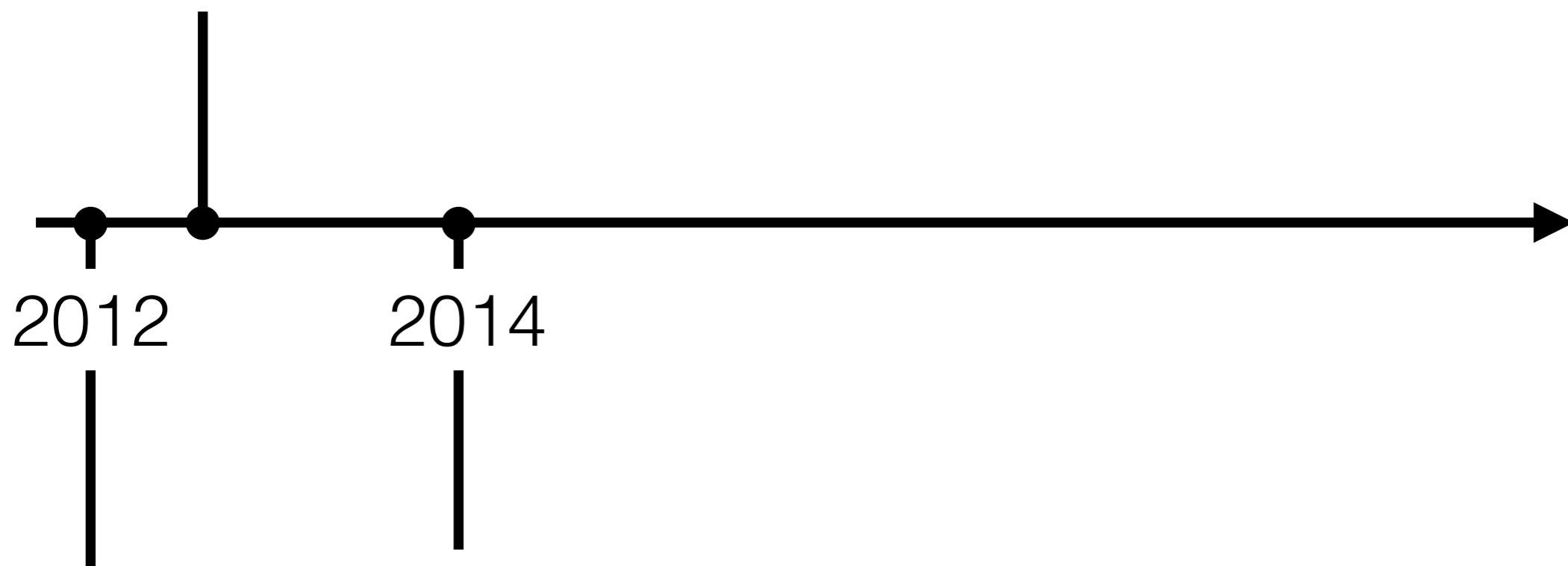
ALCHEM

[alchem.usc.edu](http://alchem.usc.edu)

# Where are we in the development of DLAs?

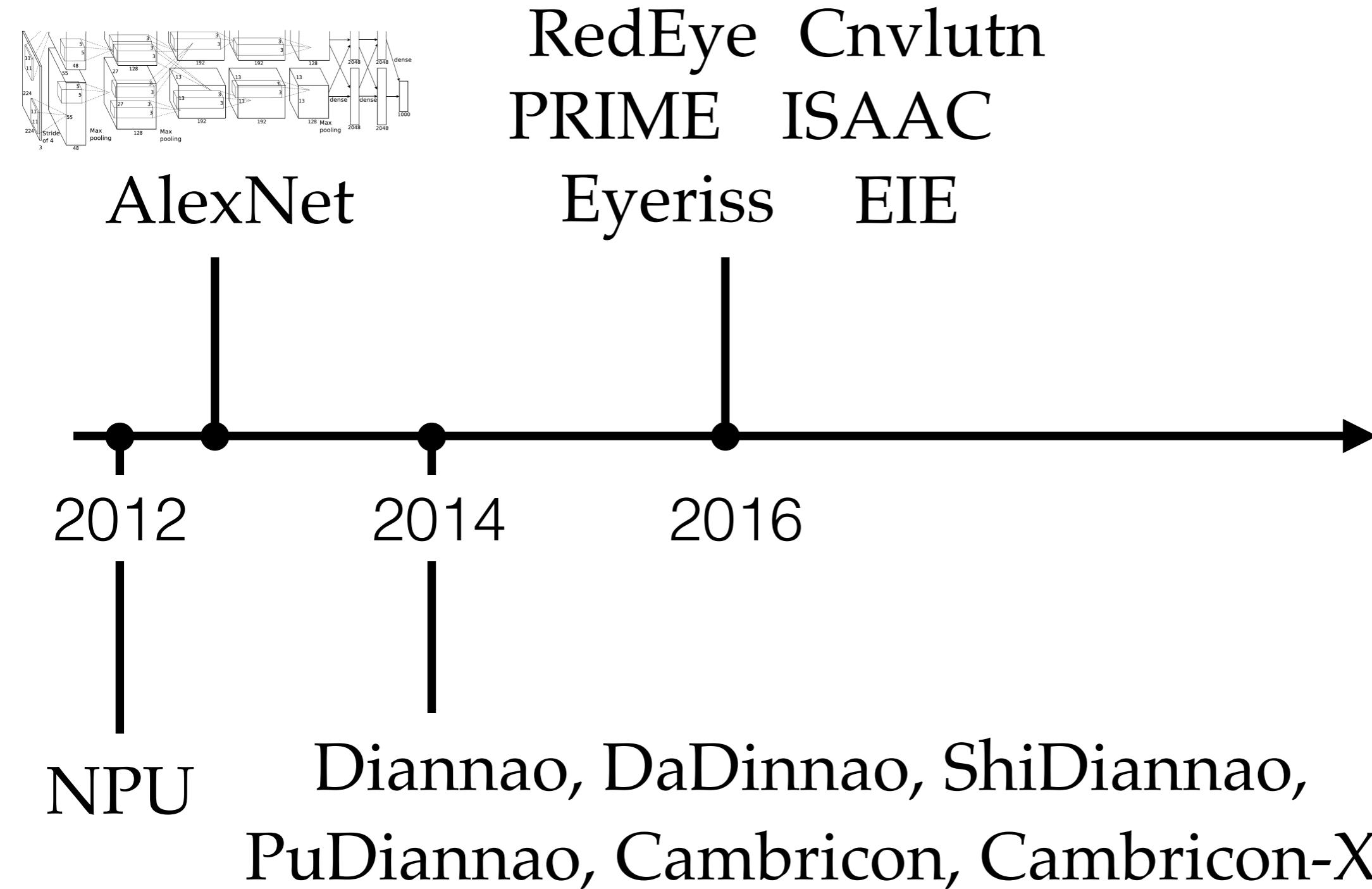


AlexNet

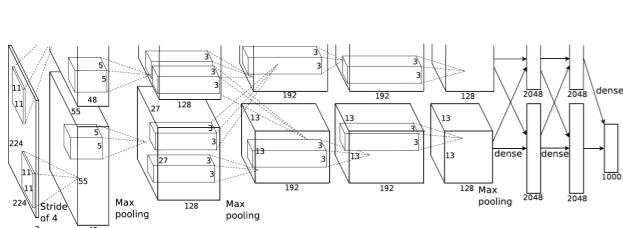


Diannao, DaDinnao, ShiDiannao,  
PuDiannao, Cambricon, Cambricon-X

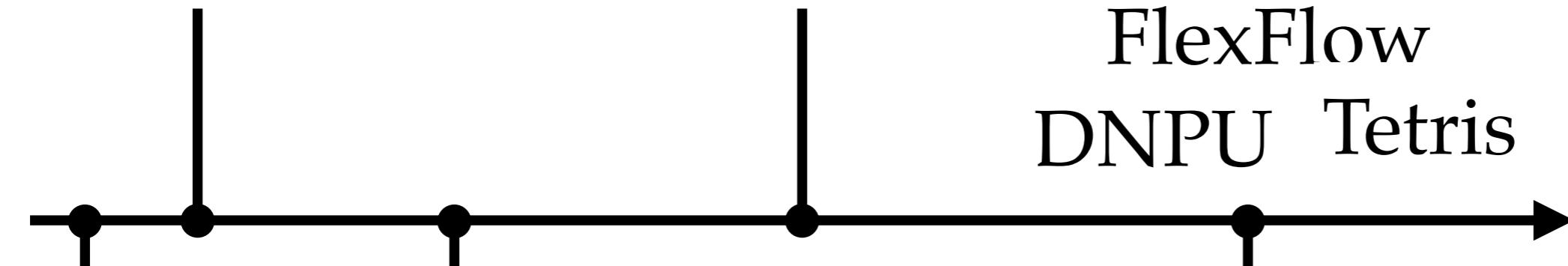
# Where are we in the development of DLAs?



# Where are we in the development of DLAs?



AlexNet



2012

NPU

2014

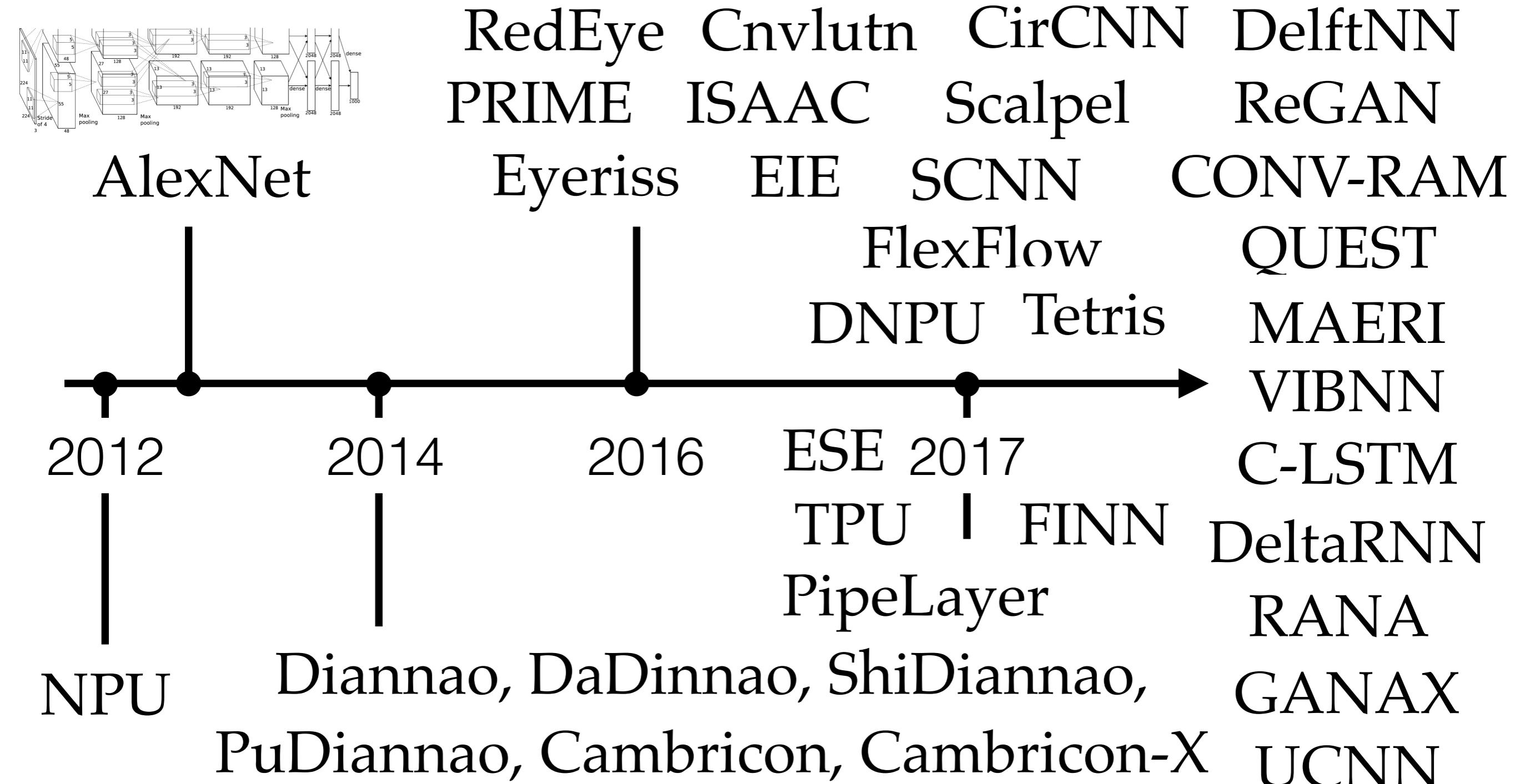
2016

ESE 2017

TPU | FINN

PipeLayer

# Where are we in the development of DLAs?



# Why HyPar

---

As a deep learning accelerator architect,  
we have well-designed DL accelerators,  
and we are approaching the accelerator  
optimization limitation (Accelerator Wall).

Why not use an accelerator array?

How can we parallelize them?

# Outline

---

- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

# Outline

---

- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

# HyPar: communication model

Accelerator Design is Guided by Cost

Arithmetic is Free  
(particularly low-precision)

Memory is expensive

Communication is prohibitively expensive



Bill Dally@AACBB(Sat.)

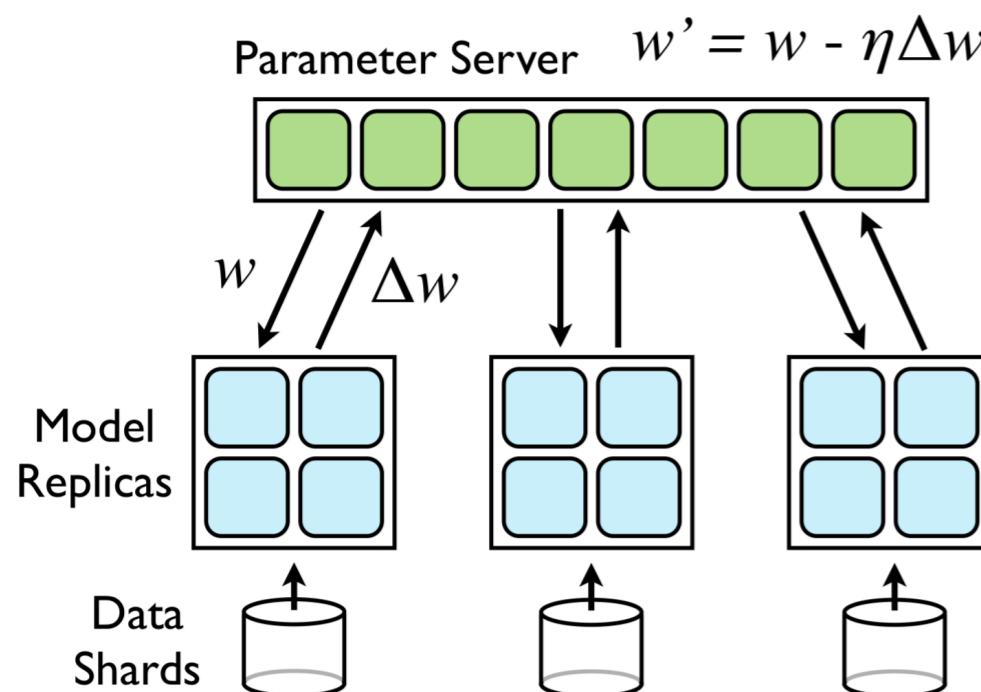
=> Let's reduce communication!

# HyPar: communication model

Two default parallelism:

- Data parallelism:

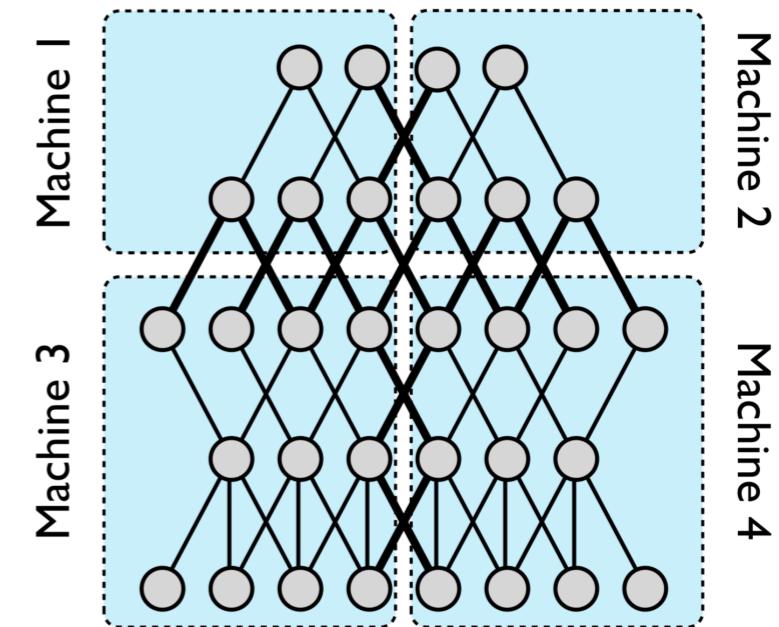
data partitioned



(NIPS'12)

- Model parallelism:

model partitioned



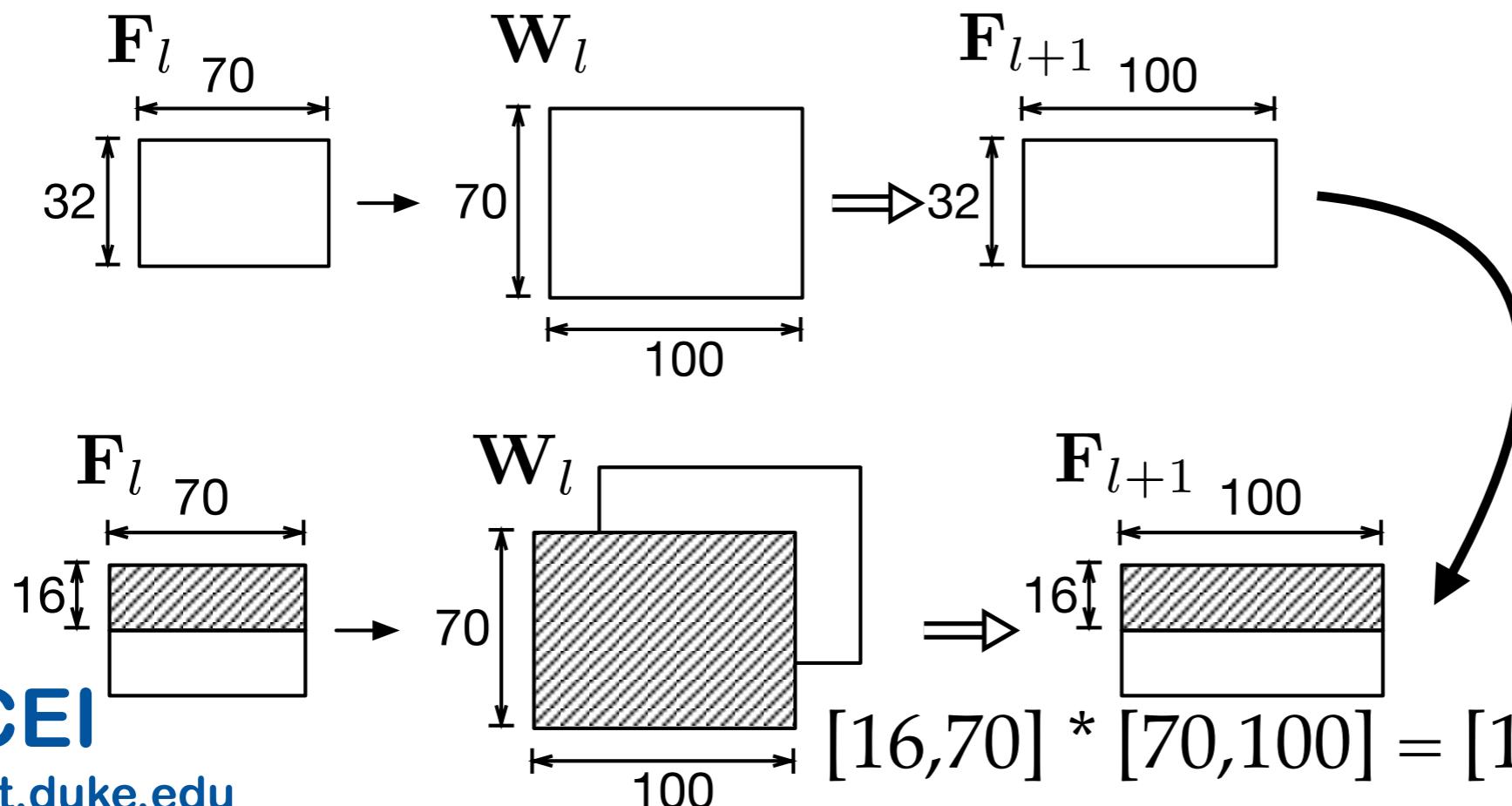
# HyPar: communication model

Where is the communication from? What's the amount?

- data parallelism:

Assume we are computing matrix multiplication with two accelerators.

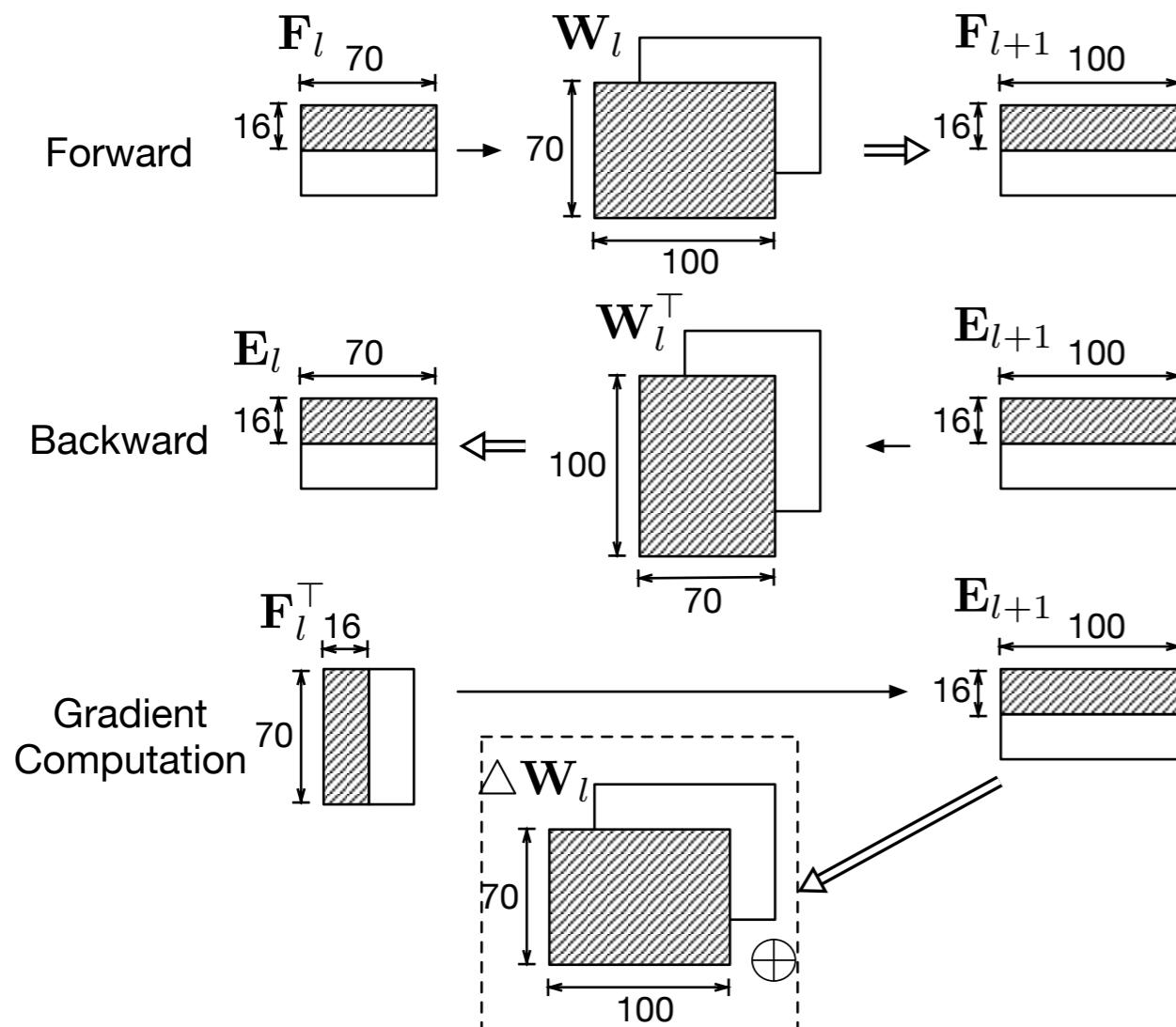
The tensor sizes are:  $[32,70] * [70,100] = [32,100]$



Tensors are partitioned and assigned to two accelerators.

# HyPar: communication model

- data parallelism:



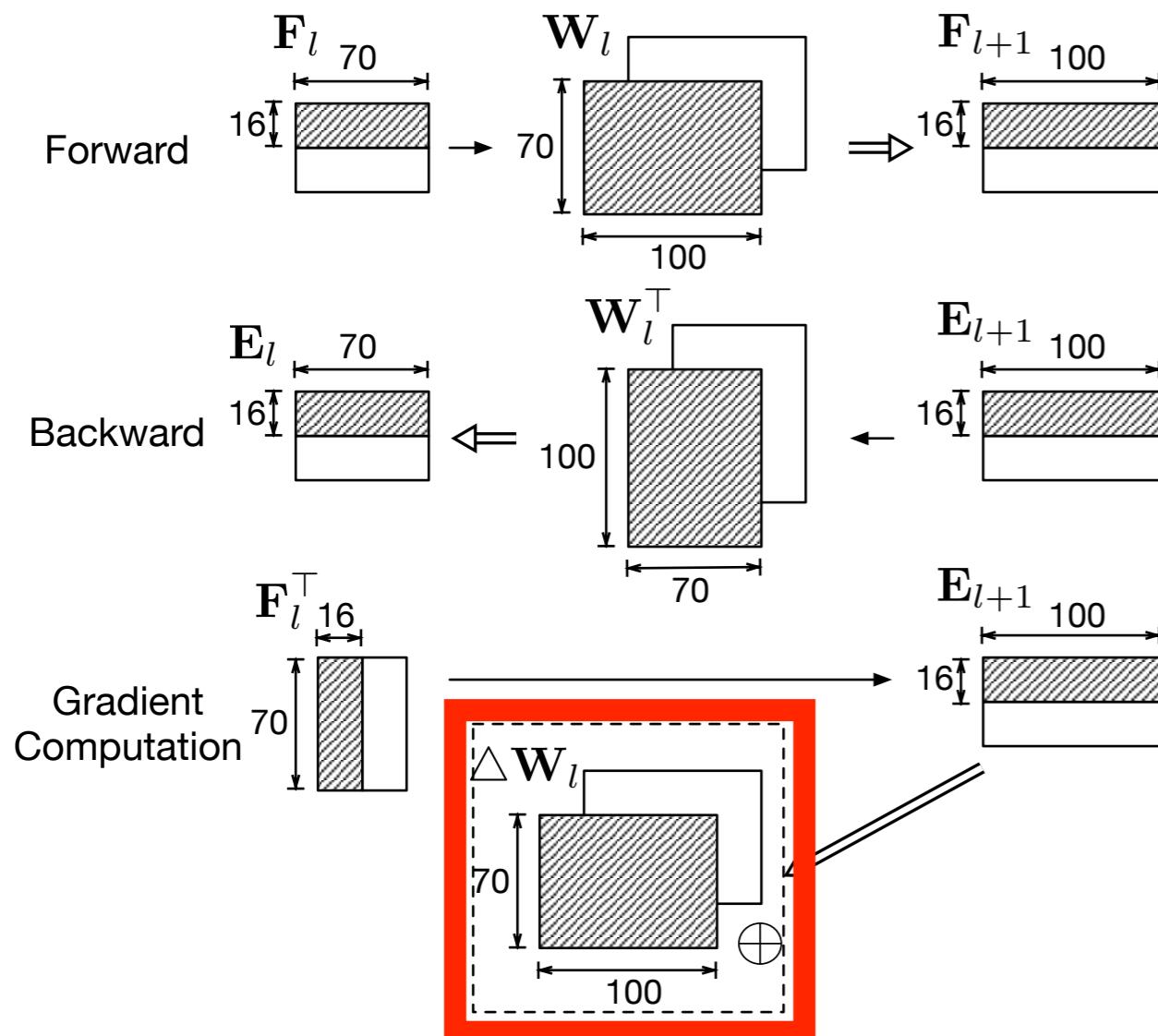
$$[16, 70] * [70, 100] = [16, 100]$$

$$[16, 100] * [100, 70] = [16, 70]$$

cei.prai [70, 16] \* [16, 100] = ?[70, 100]

# HyPar: communication model

- data parallelism:



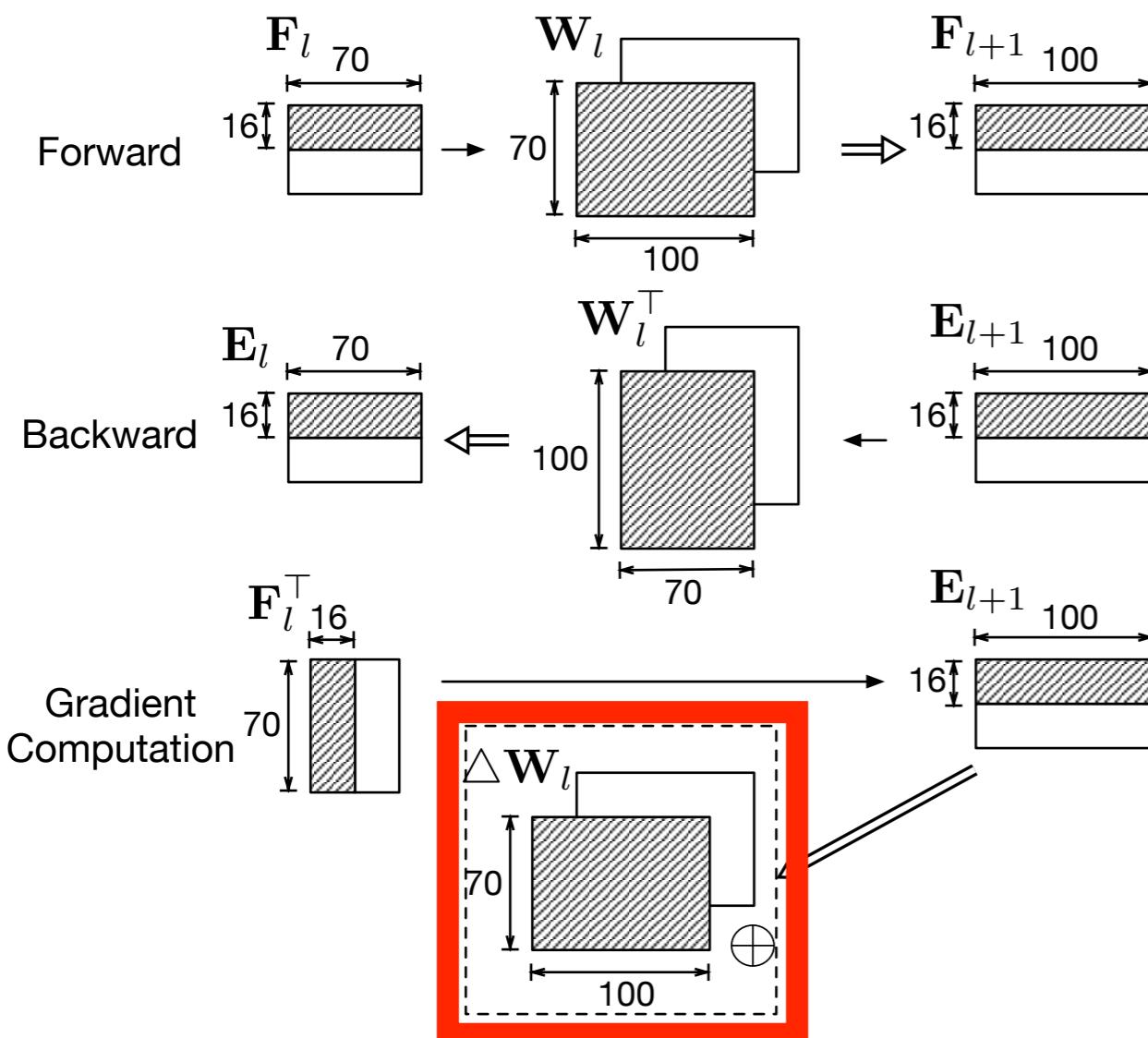
$$[16, 70] * [70, 100] = [16, 100]$$

$$[16, 100] * [100, 70] = [16, 70]$$

$[70, 16] * [16, 100] = ?[70, 100]$

# HyPar: communication model

- data parallelism:

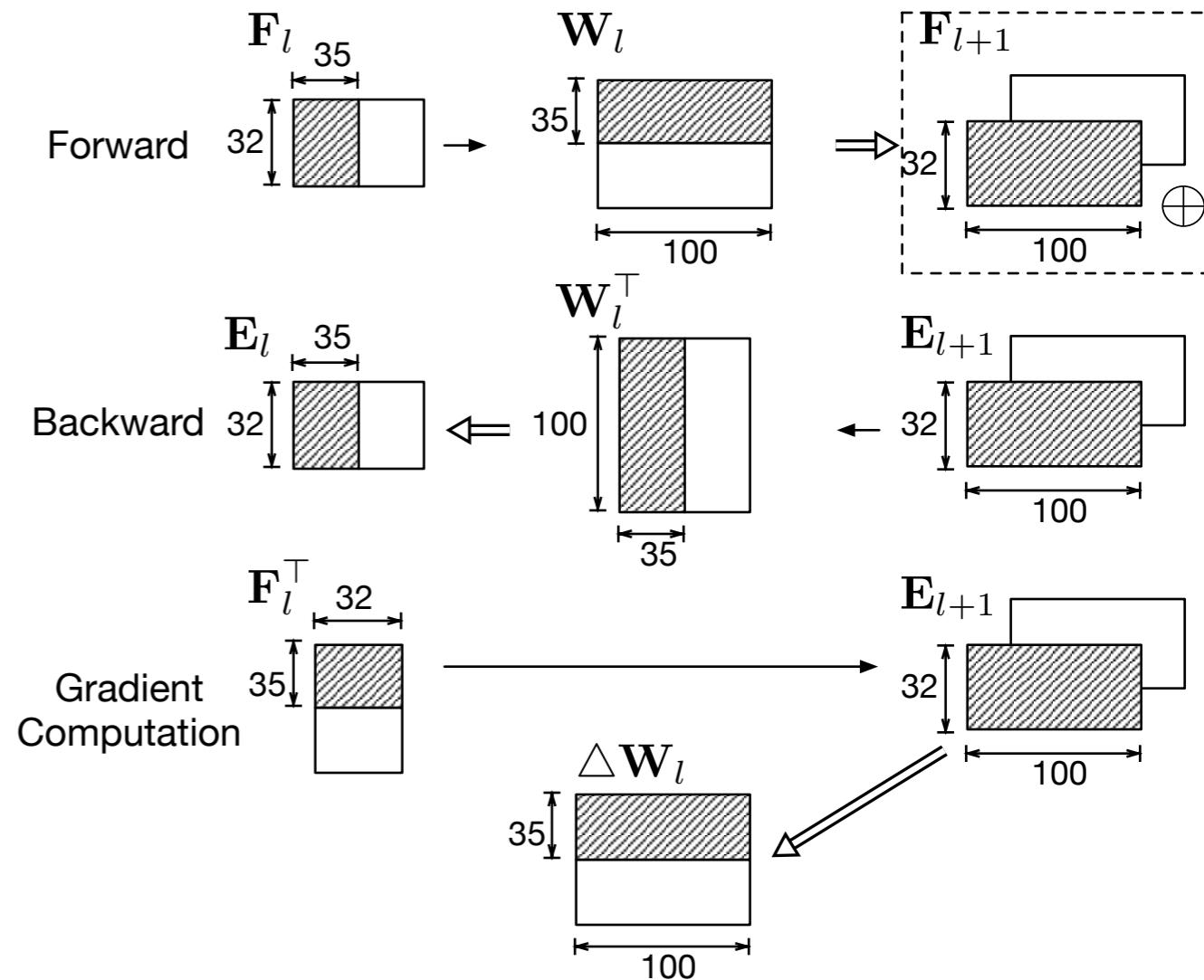


$$[16,70] * [70,100] = [16,100]$$

$$[16,100] * [100,70] = [16,70]$$

$$[70,16] * [16,100] = ?[70,100]$$

- model parallelism:



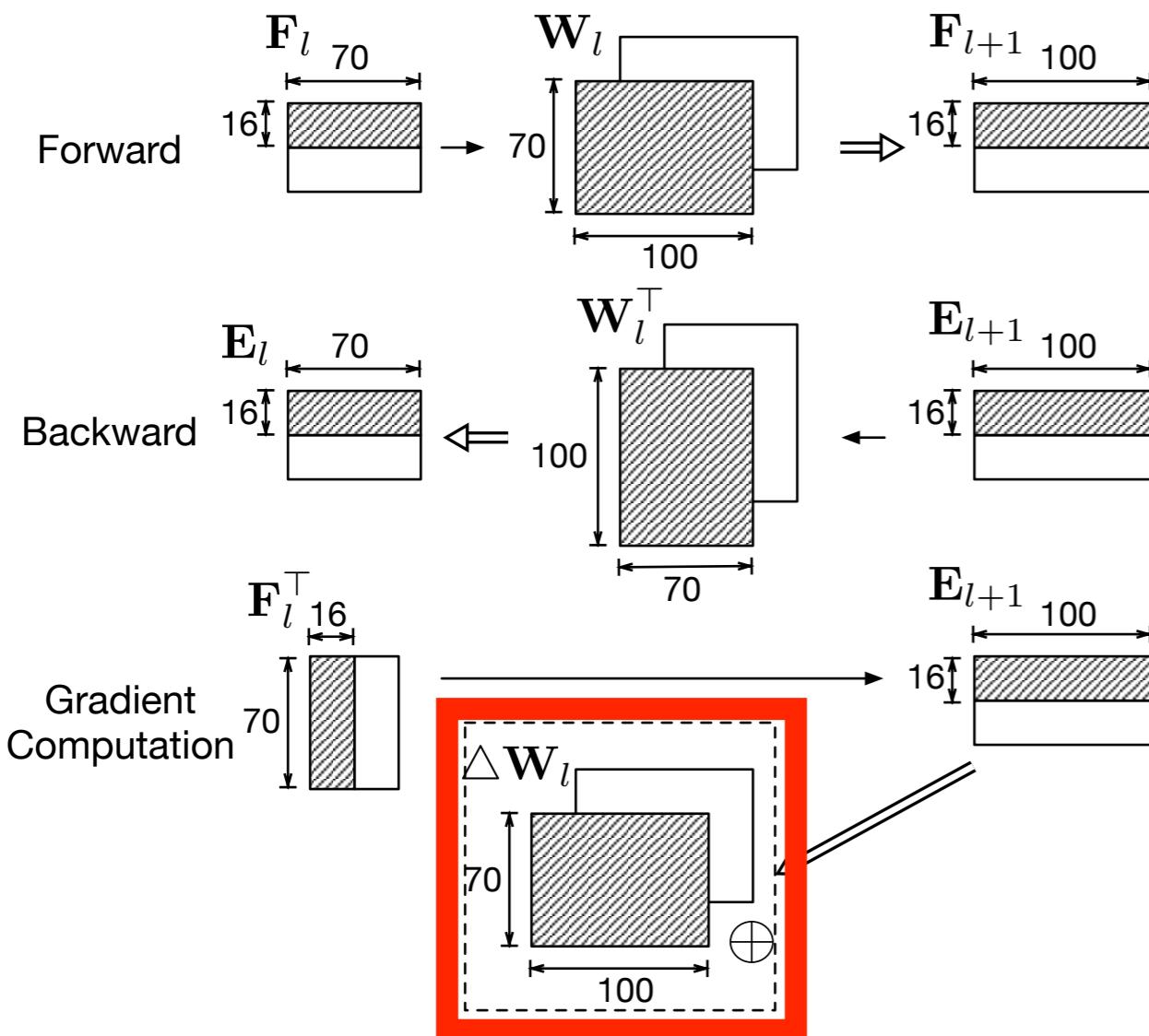
$$[32,35] * [35,100] = ?[32,100]$$

$$[32,100] * [100,35] = [32,35]$$

$$[35,32] * [32,100] = [35,100]$$

# HyPar: communication model

- data parallelism:

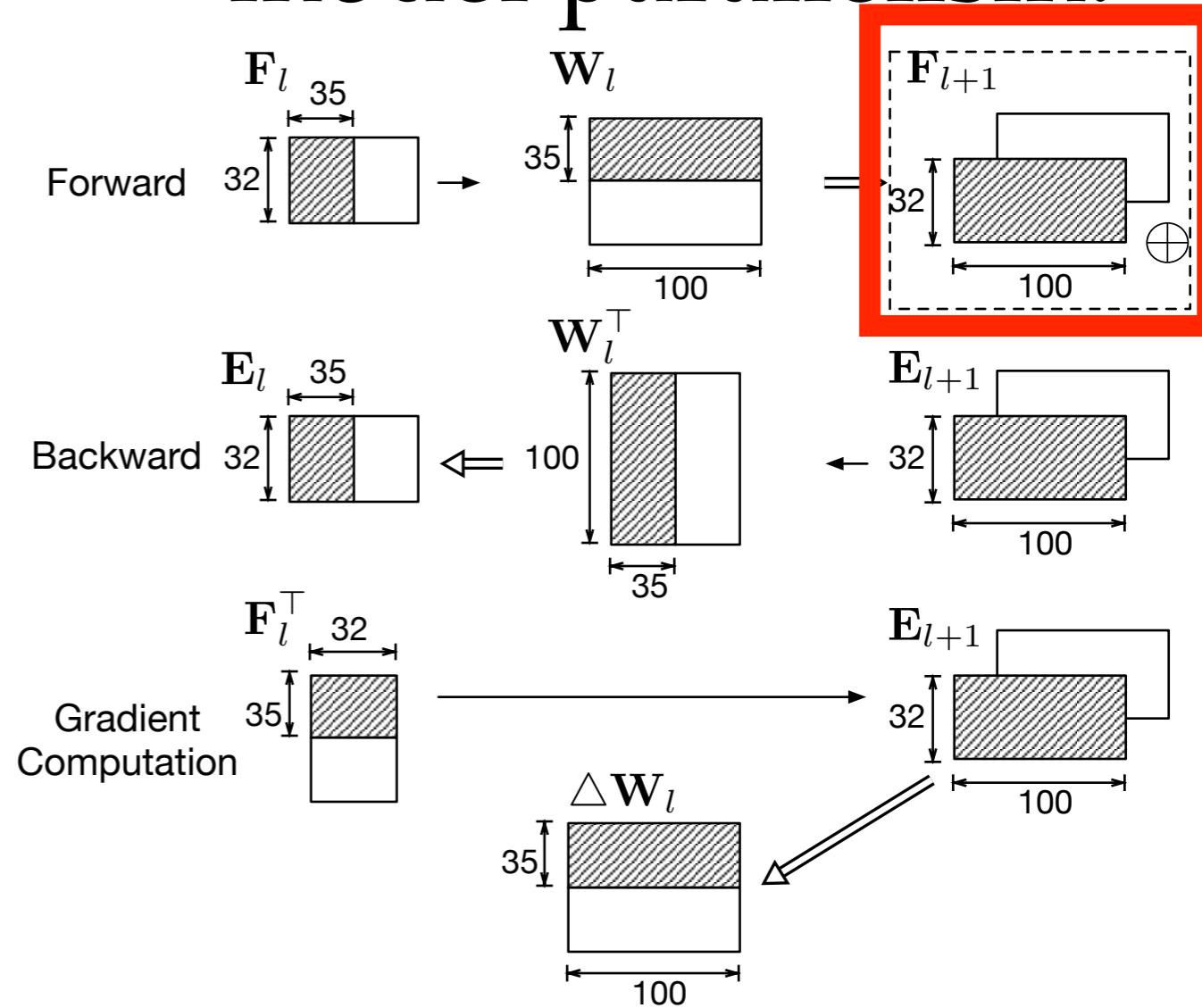


$$[16,70] * [70,100] = [16,100]$$

$$[16,100] * [100,70] = [16,70]$$

$$[70,16] * [16,100] = ?[70,100]$$

- model parallelism:



$$[32,35] * [35,100] = ?[32,100]$$

$$[32,100] * [100,35] = [32,35]$$

$$[35,32] * [32,100] = [35,100]$$

# HyPar: intra-layer com. model

---

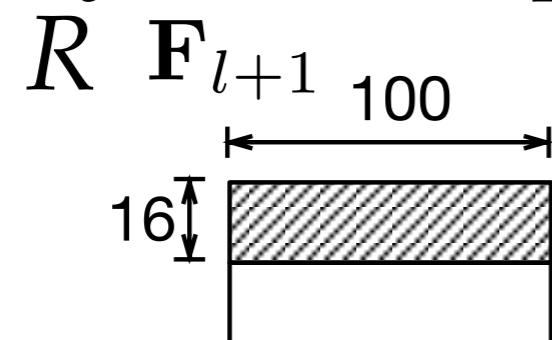
- Where: partial sum accumulation.
- Amount?

data parallelism	model parallelism
$\mathbb{A}(\Delta \mathbf{W}_l)$	$\mathbb{A}(\mathbf{F}_{l+1})$

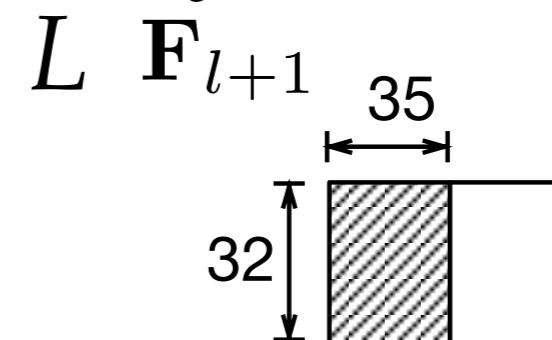
# HyPar: inter-layer com. model

- Where? Amount?

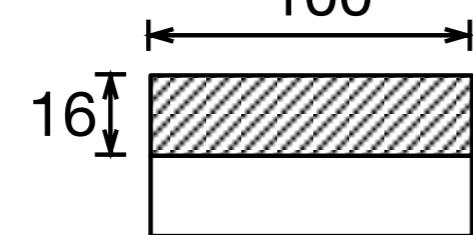
Layer L Output



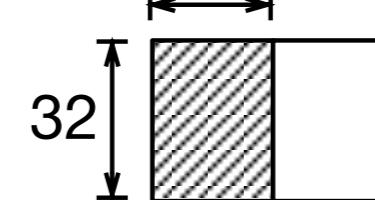
Layer L+1 Input



$R \ E_{l+1}$



$L \ E_{l+1}$

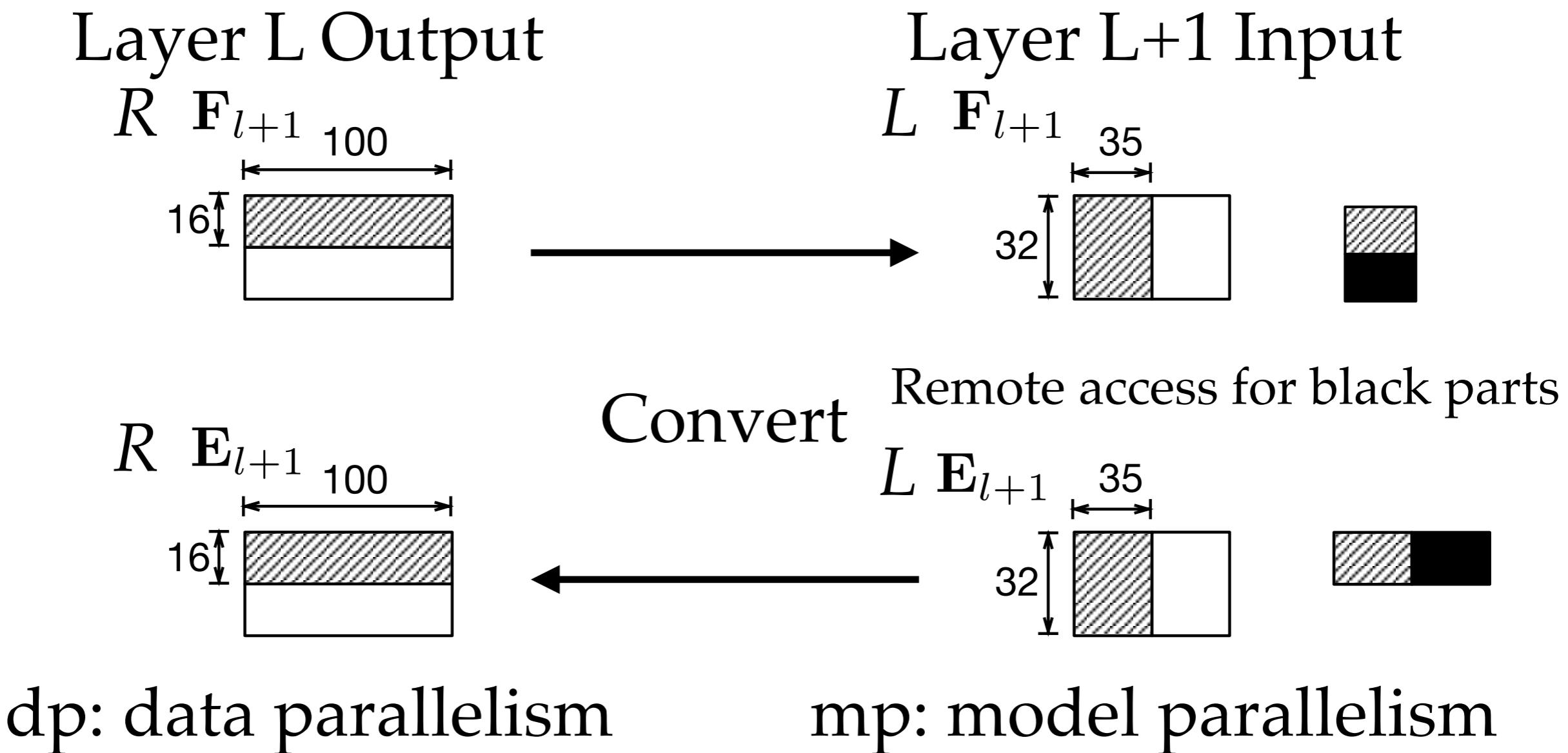


dp: data parallelism

mp: model parallelism

# HyPar: inter-layer com. model

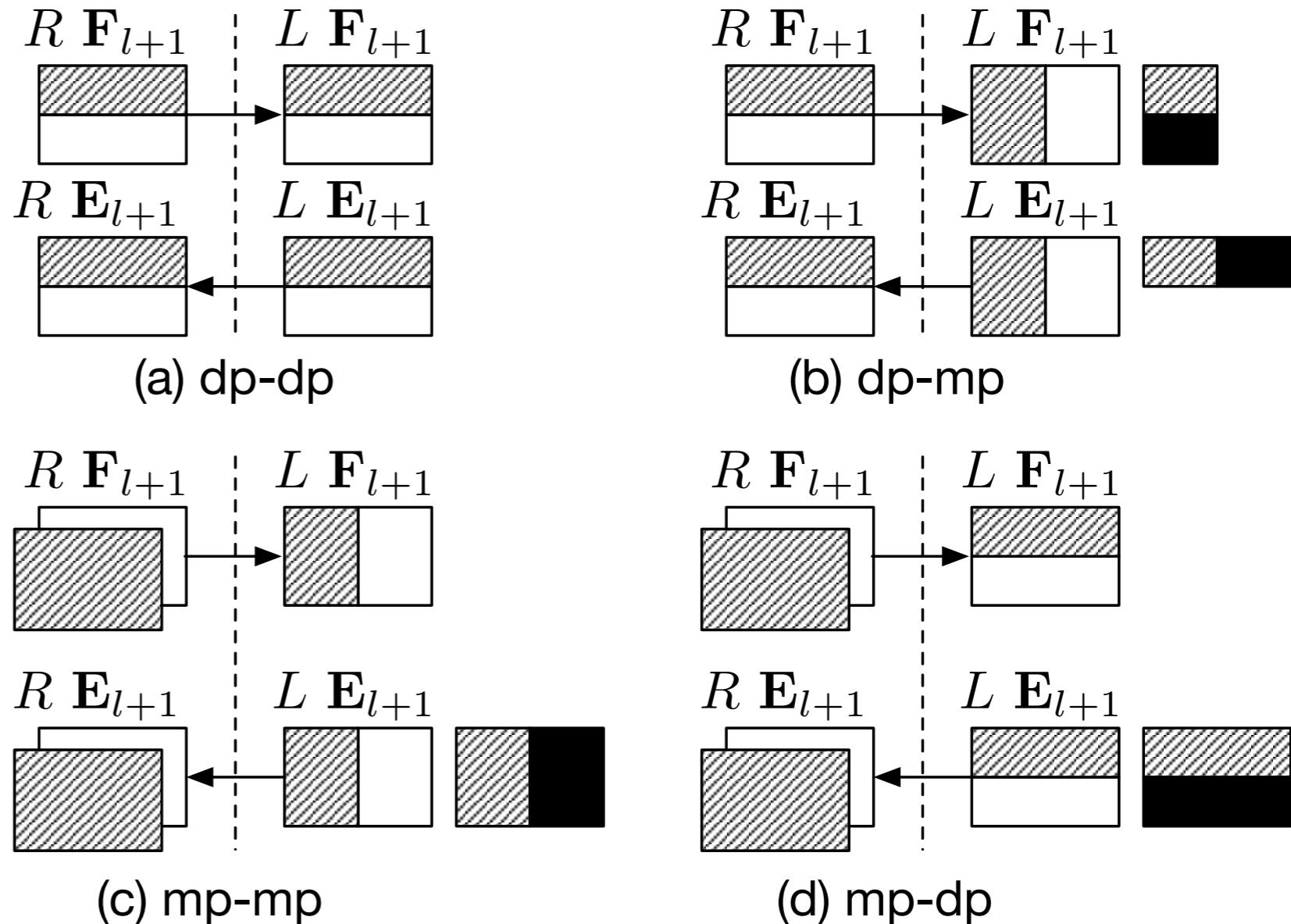
- Where? Amount?



# HyPar: inter-layer com. model

- Where: convert tensors between two layers.

- Amount?



# A slide added after the presentation

---

- For inference, we only need to use data parallelism for all layers, because
  - only forward is performed in inference, and no intra-layer communication in dp forward.
  - no inter-layer communication in dp forward

# Outline

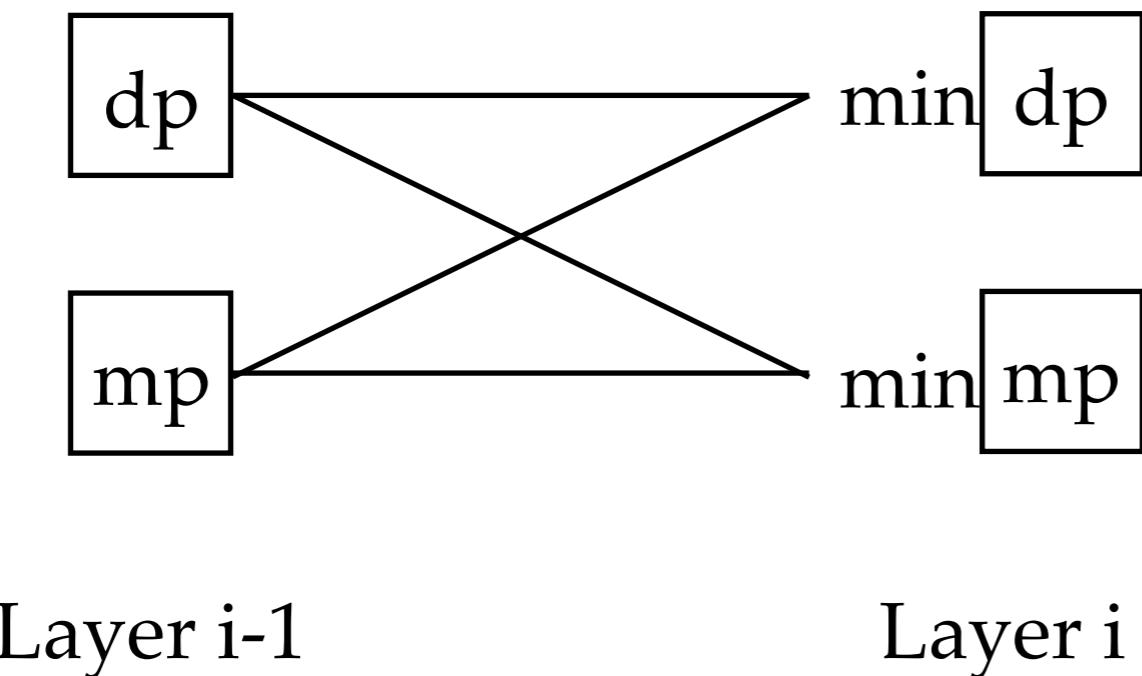
---

- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

# HyPar tensor partition: which parallelism to use?

---

- To find the partitioning which incurs less communication.
- How to find that? Brute Force?? Dynamic programming!



compute intra and inter layer communication using the communication model.

Layer  $i-1$

Layer  $i$

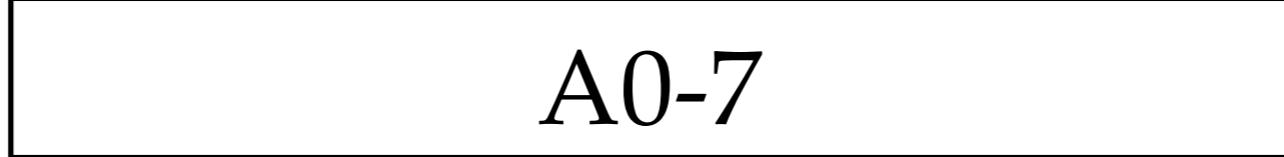
dp: data parallelism

mp: model parallelism

# HyPar: hierarchical partition

---

- We know that for 2 accelerators, how about 4, 8, 16, 32?
  - Hierarchical partition.

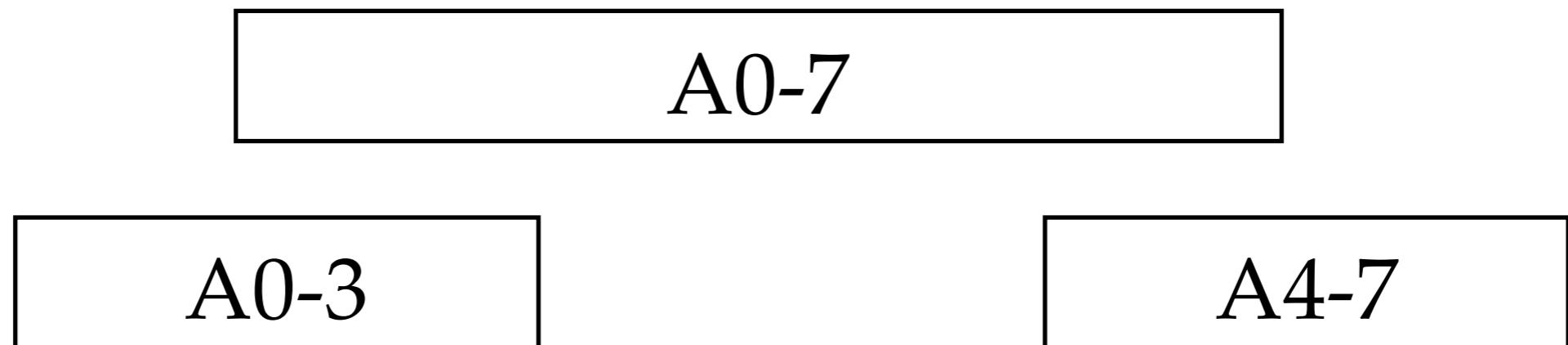


A0-7

# HyPar: hierarchical partition

---

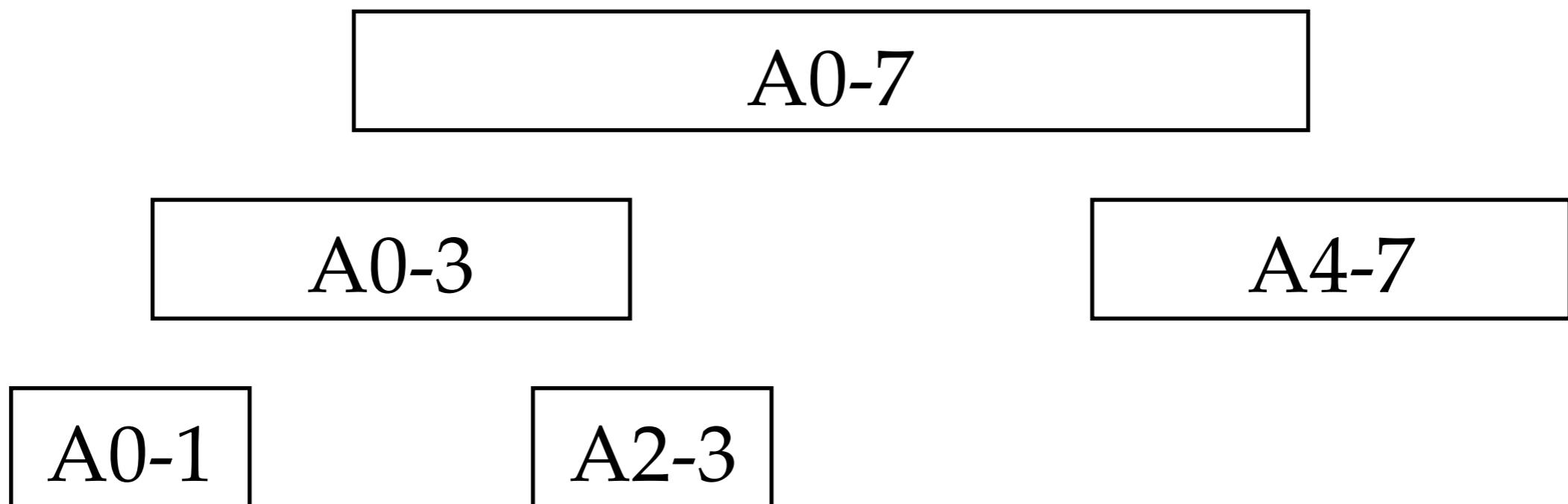
- We know that for 2 accelerators, how about 4, 8, 16, 32?
  - Hierarchical partition.



# HyPar: hierarchical partition

---

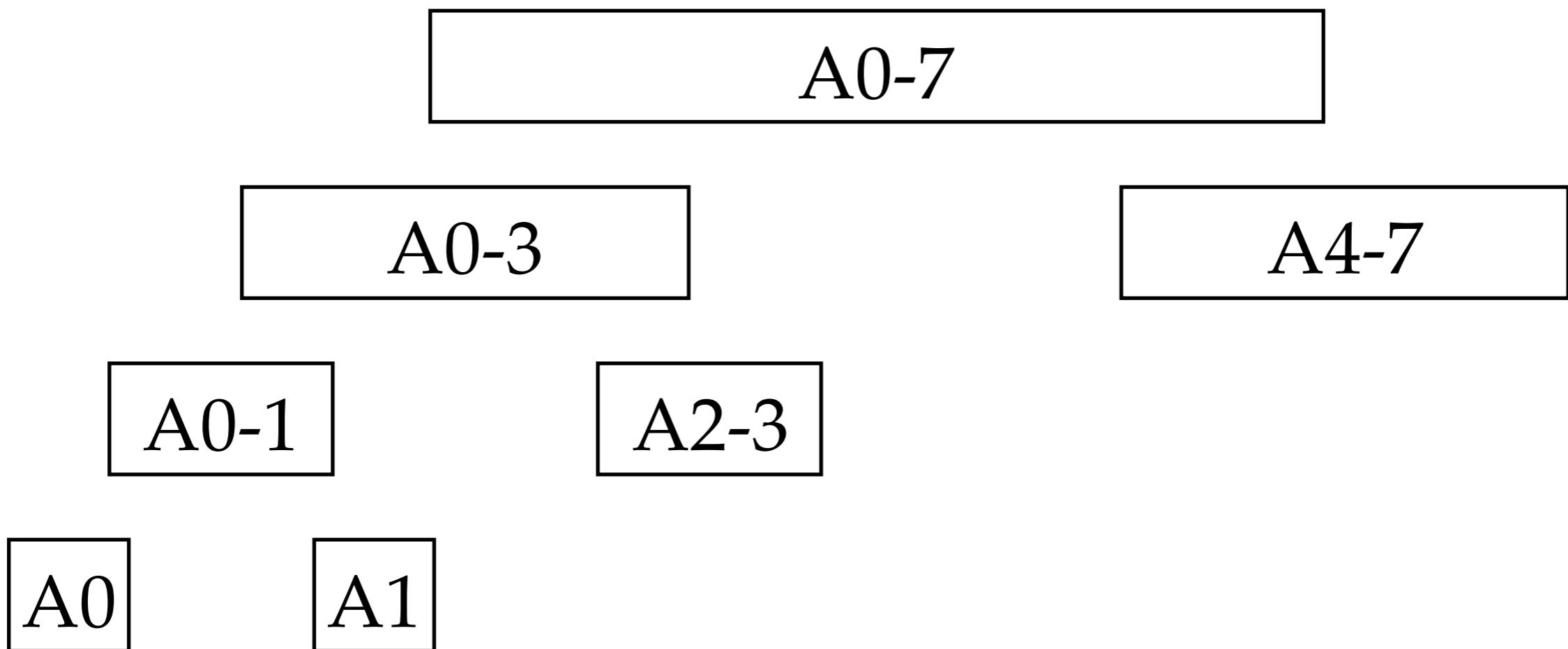
- We know that for 2 accelerators, how about 4, 8, 16, 32?
  - Hierarchical partition.



# HyPar: hierarchical partition

---

- We know that for 2 accelerators, how about 4, 8, 16, 32?
  - Hierarchical partition.

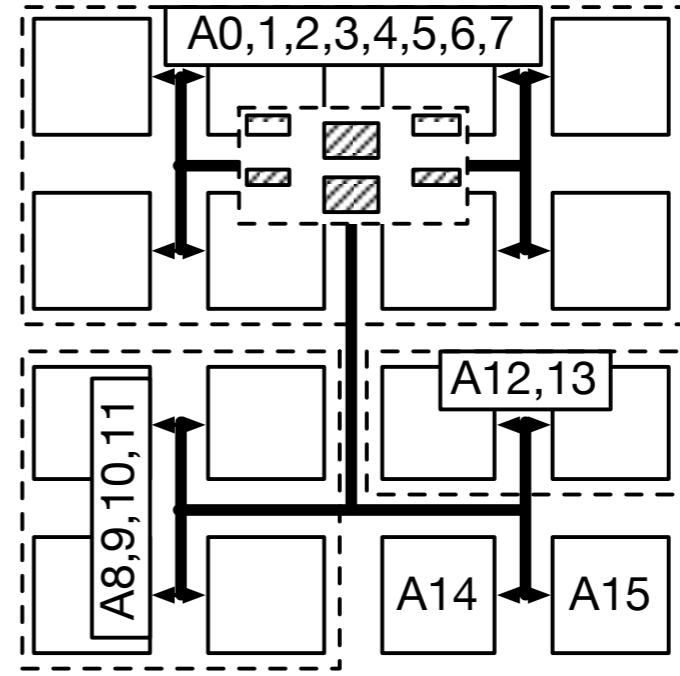
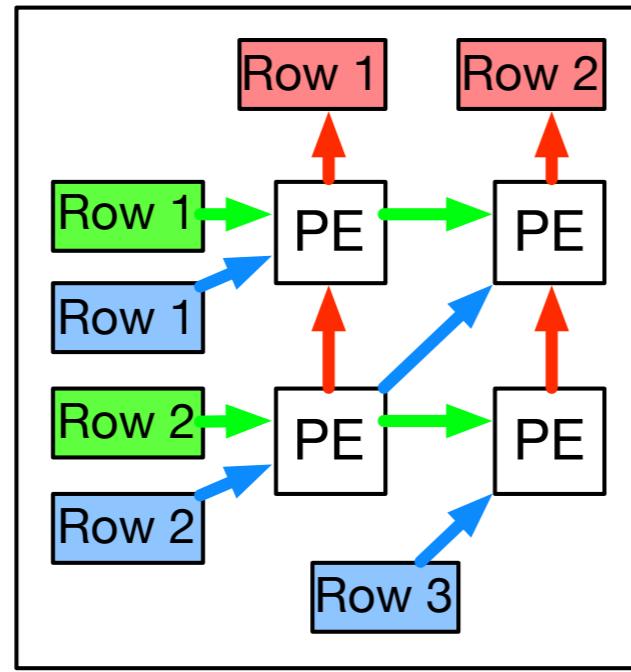


# Outline

---

- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

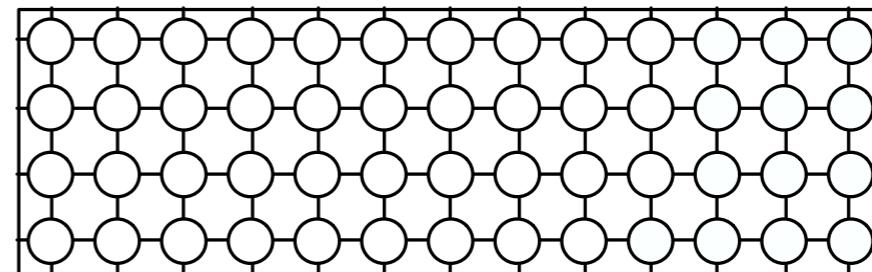
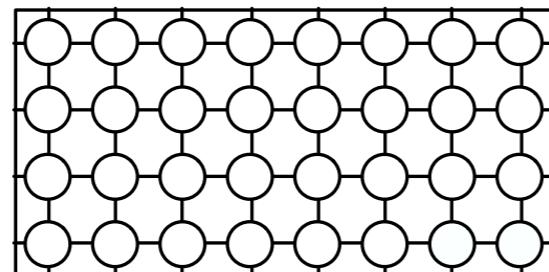
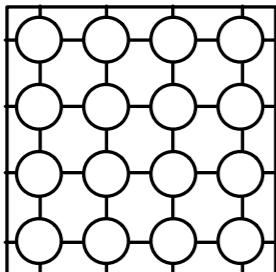
# HyPar evaluation setup



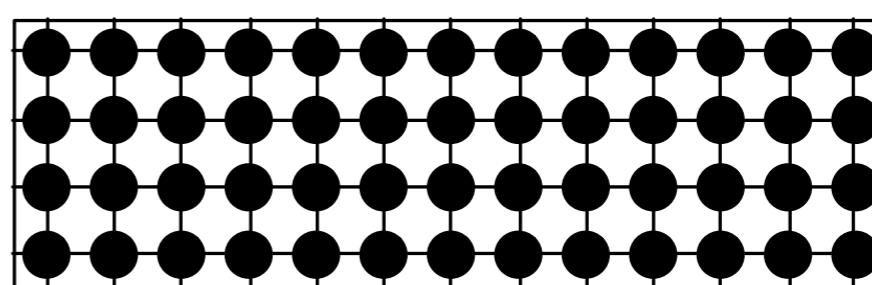
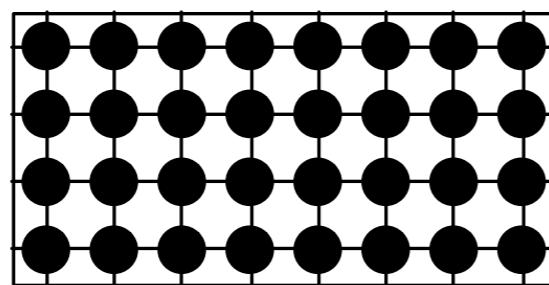
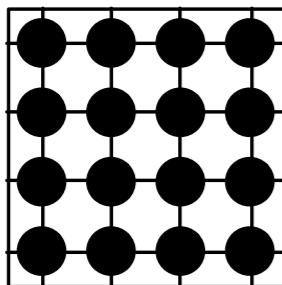
- 16 Eyeriss-like accelerators(4 hierarchies)
  - Total computation density: 1344 (=84\*16) GOPS / s
  - Total link bandwidth: 25.6 (=1.6\*8)Gb / s
- Ten benchmarks: from Lenet to VGGs.

# HyPar parallelism

DP

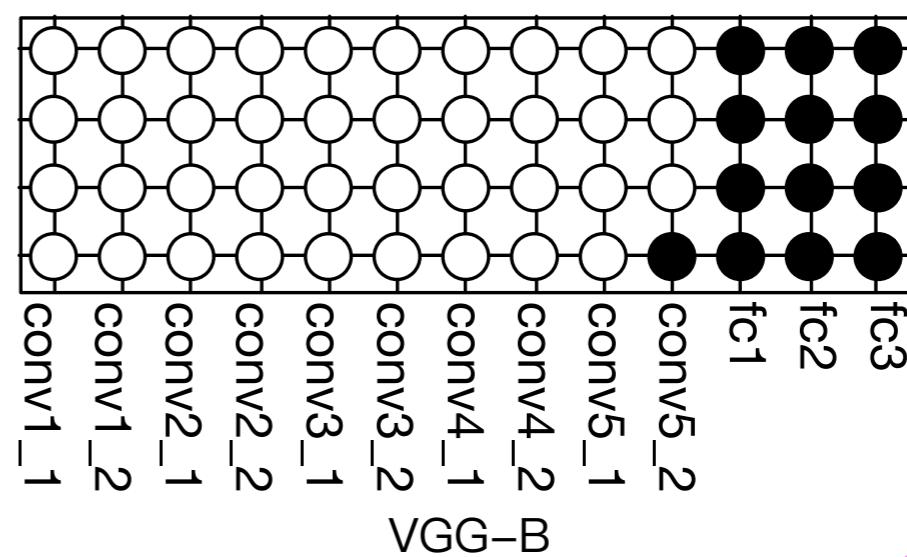
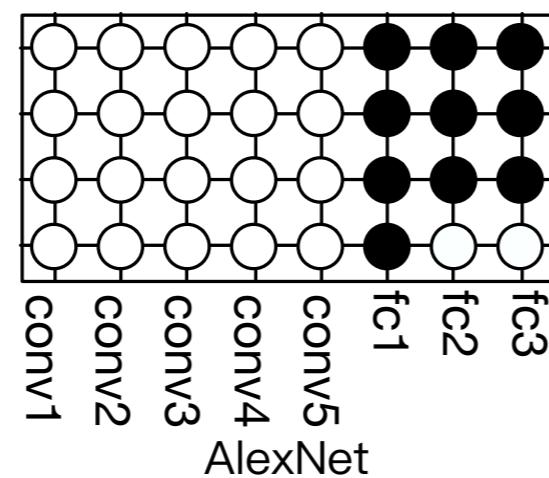
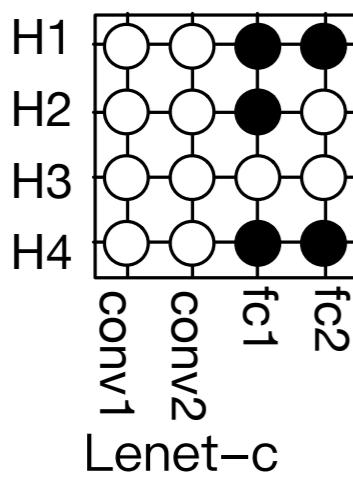


MP



- data parallelism
- model parallelism

HyPar



CEI

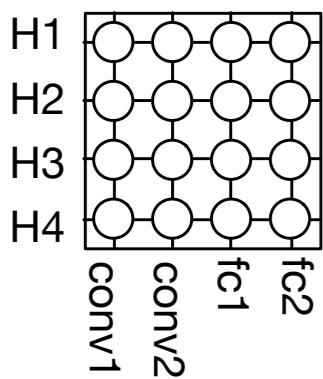
[cei.pratt.duke.edu](http://cei.pratt.duke.edu)

ALCHEM  
[alchem.usc.edu](http://alchem.usc.edu)

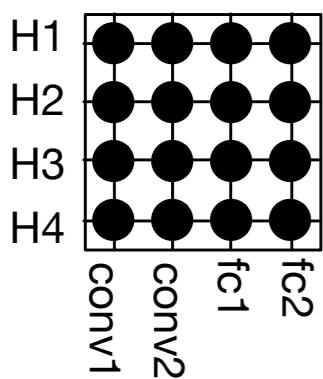
# HyPar parallelism

---

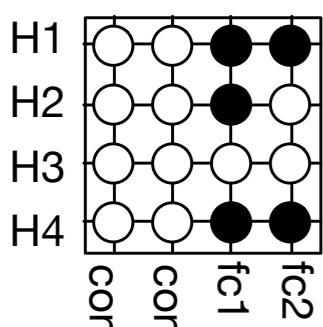
dp



mp



HyPar

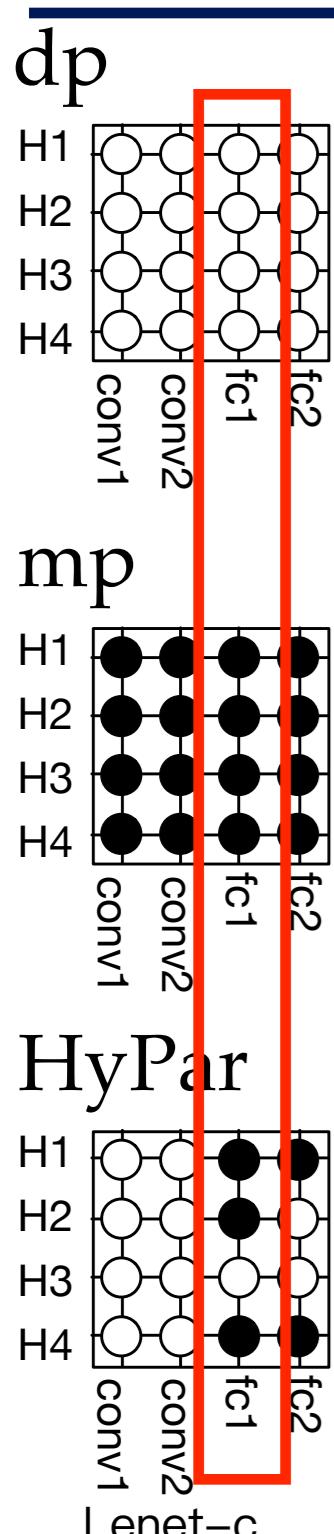


**CEI**

[cei.pratt.duke.edu](http://cei.pratt.duke.edu)

**ALCHEM**  
[alchem.usc.edu](http://alchem.usc.edu)

# HyPar parallelism

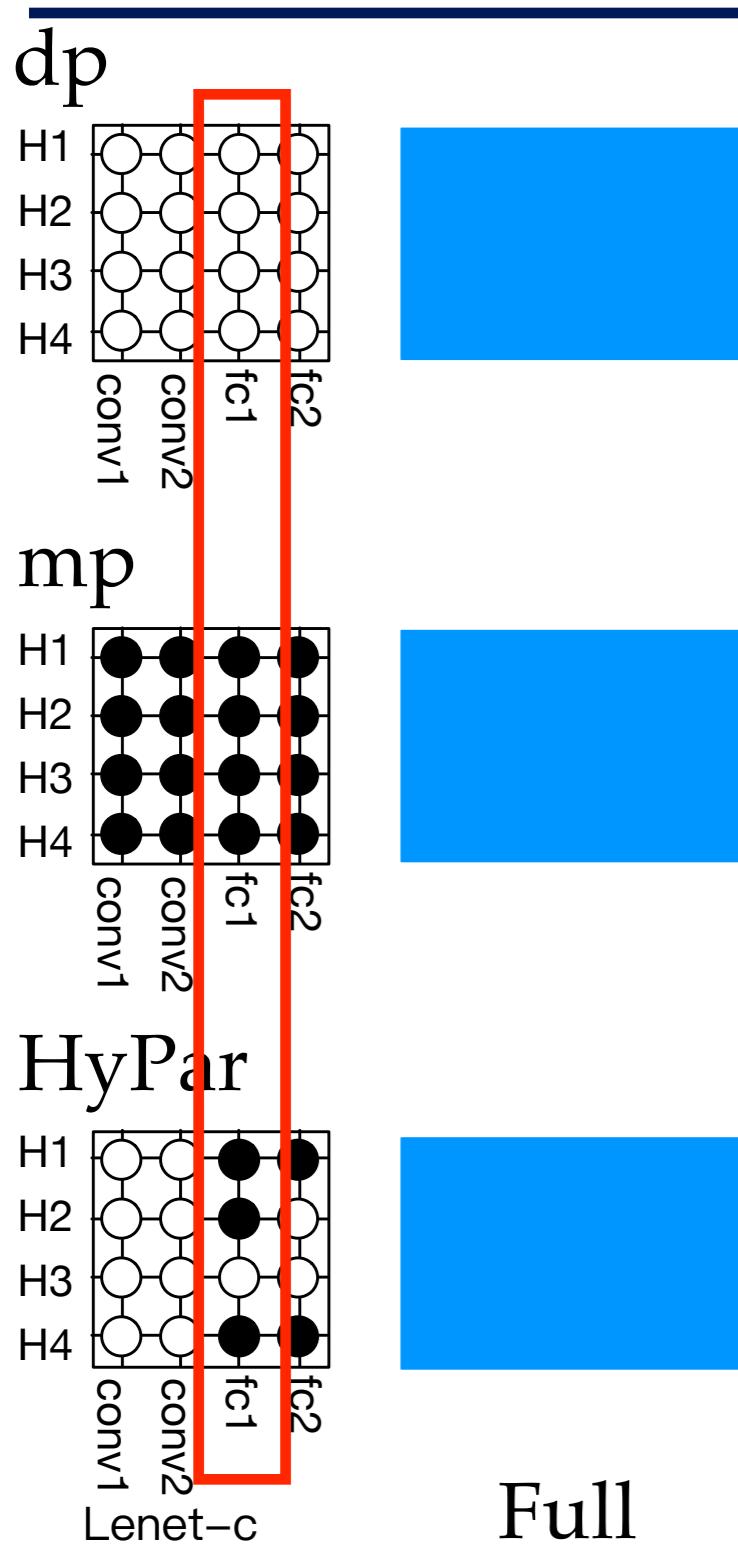


**CEI**

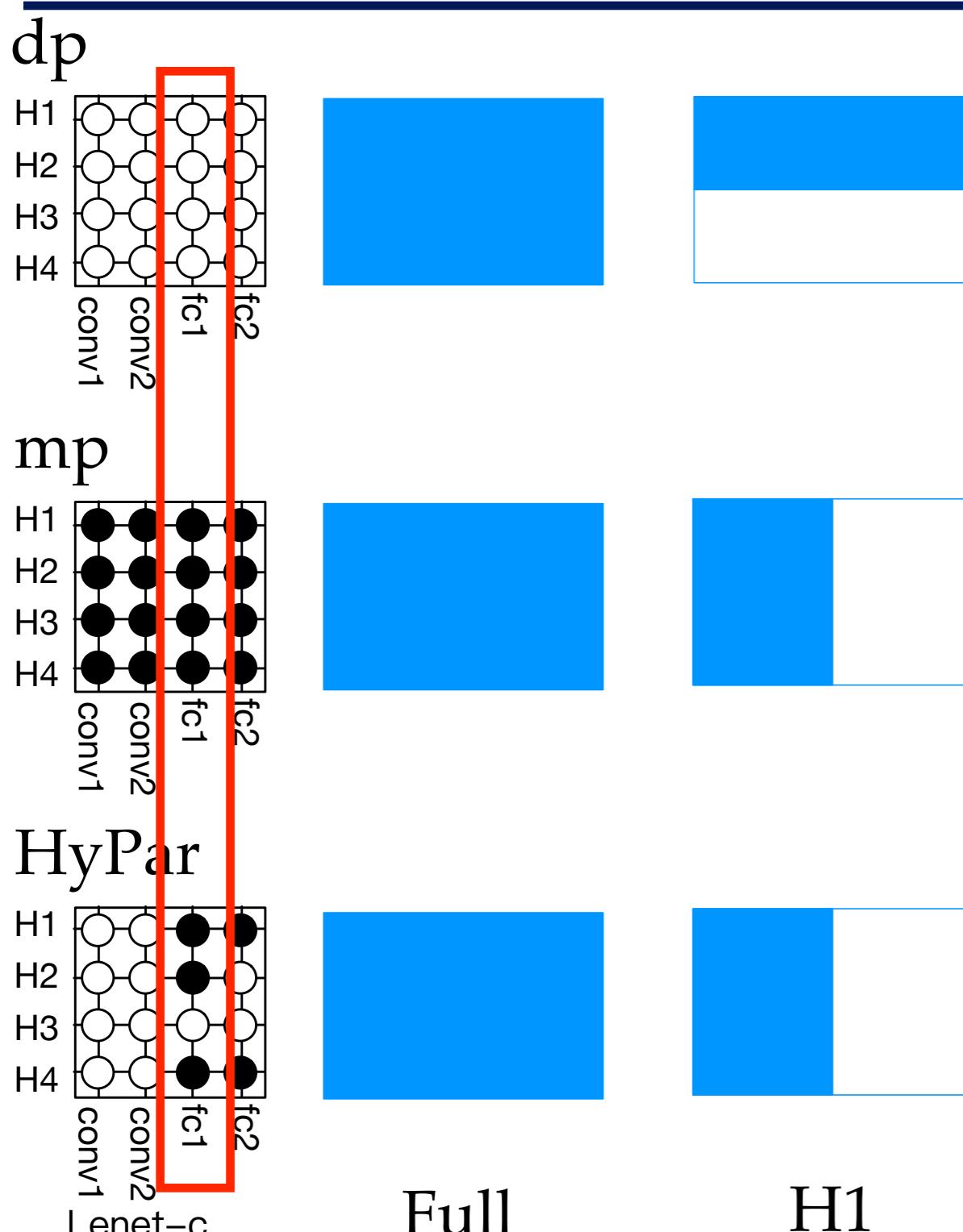
[cei.pratt.duke.edu](http://cei.pratt.duke.edu)

**ALCHEM**  
[alchem.usc.edu](http://alchem.usc.edu)

# HyPar parallelism



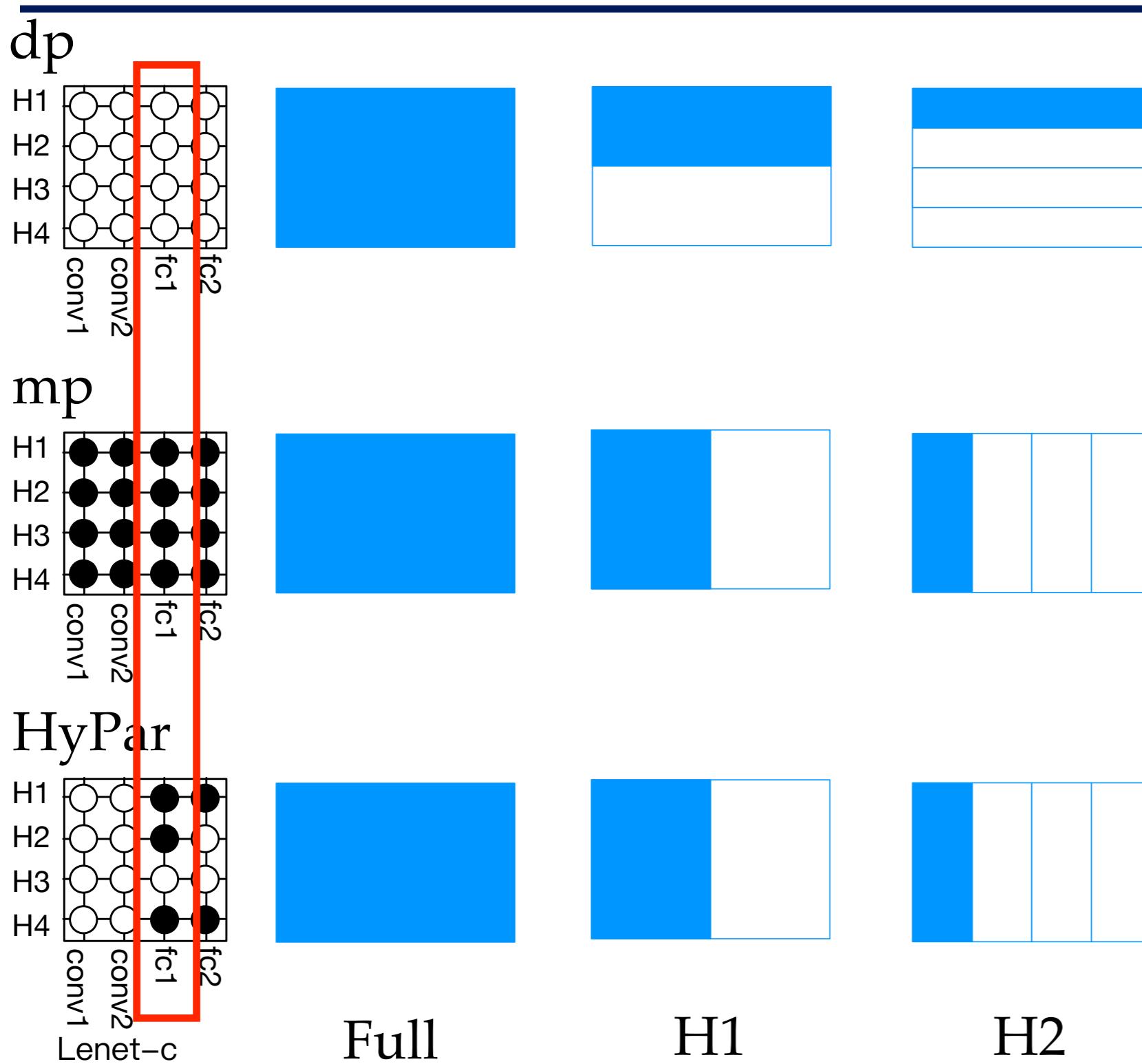
# HyPar parallelism



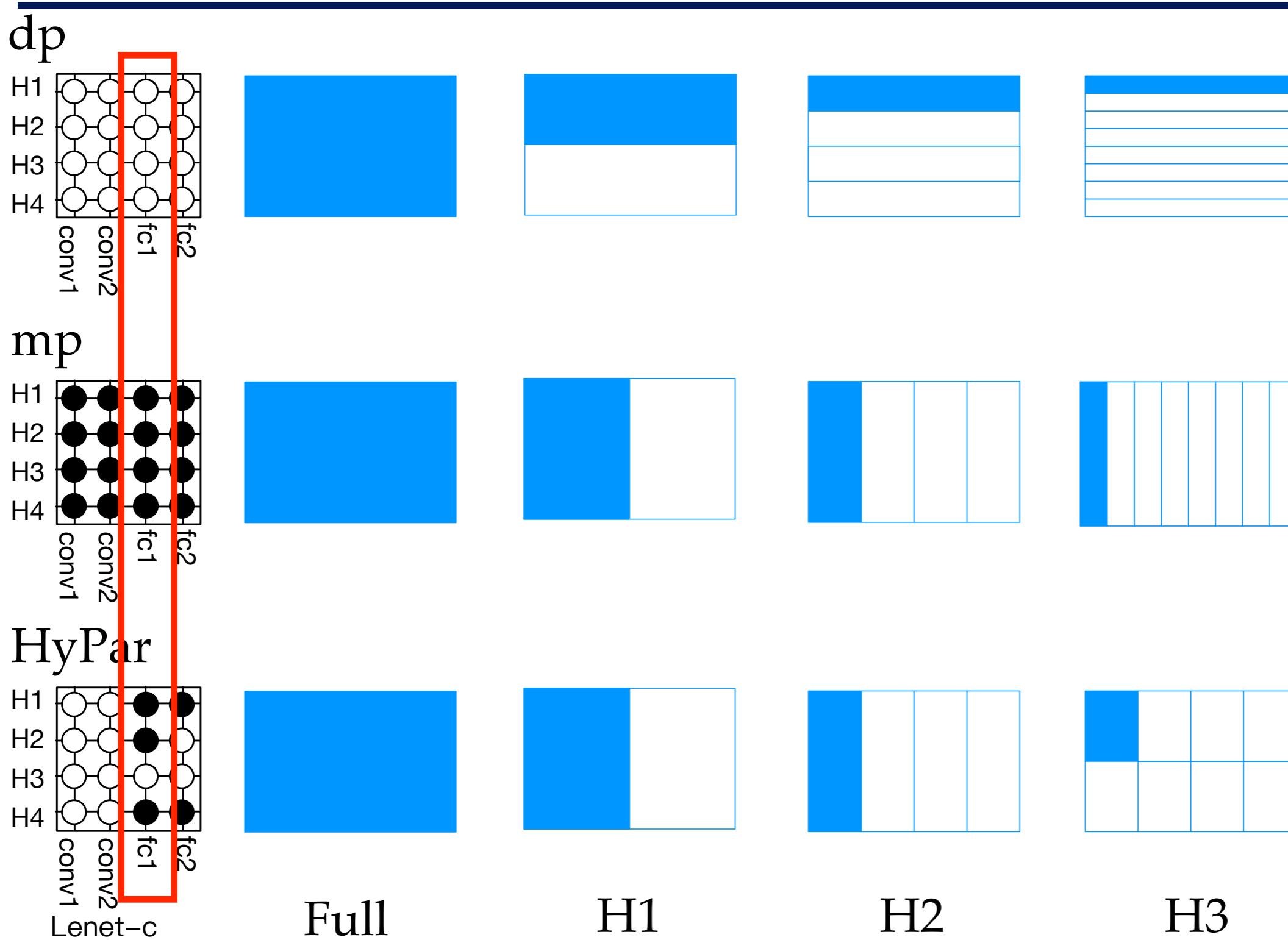
Full

H1

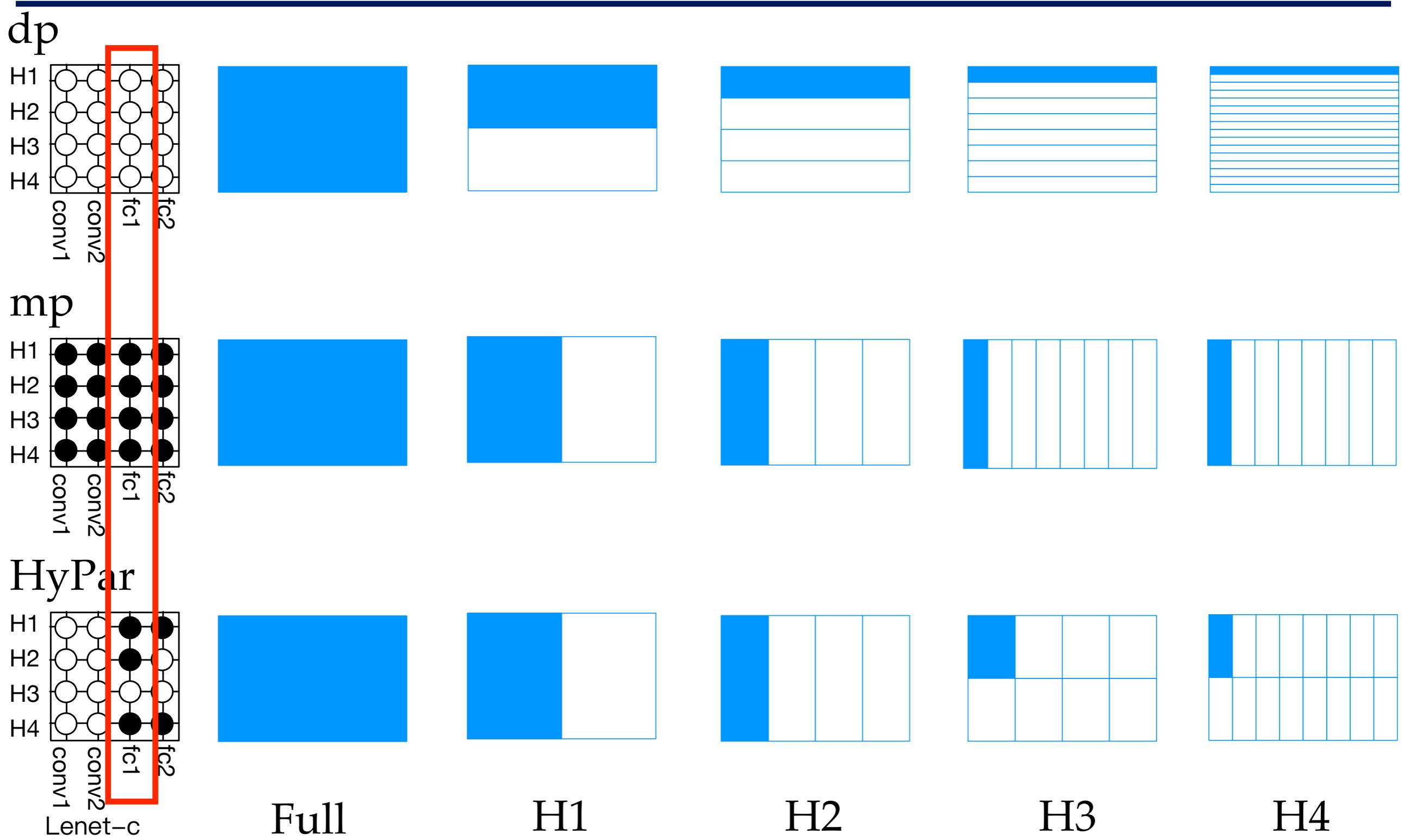
# HyPar parallelism



# HyPar parallelism



# HyPar parallelism

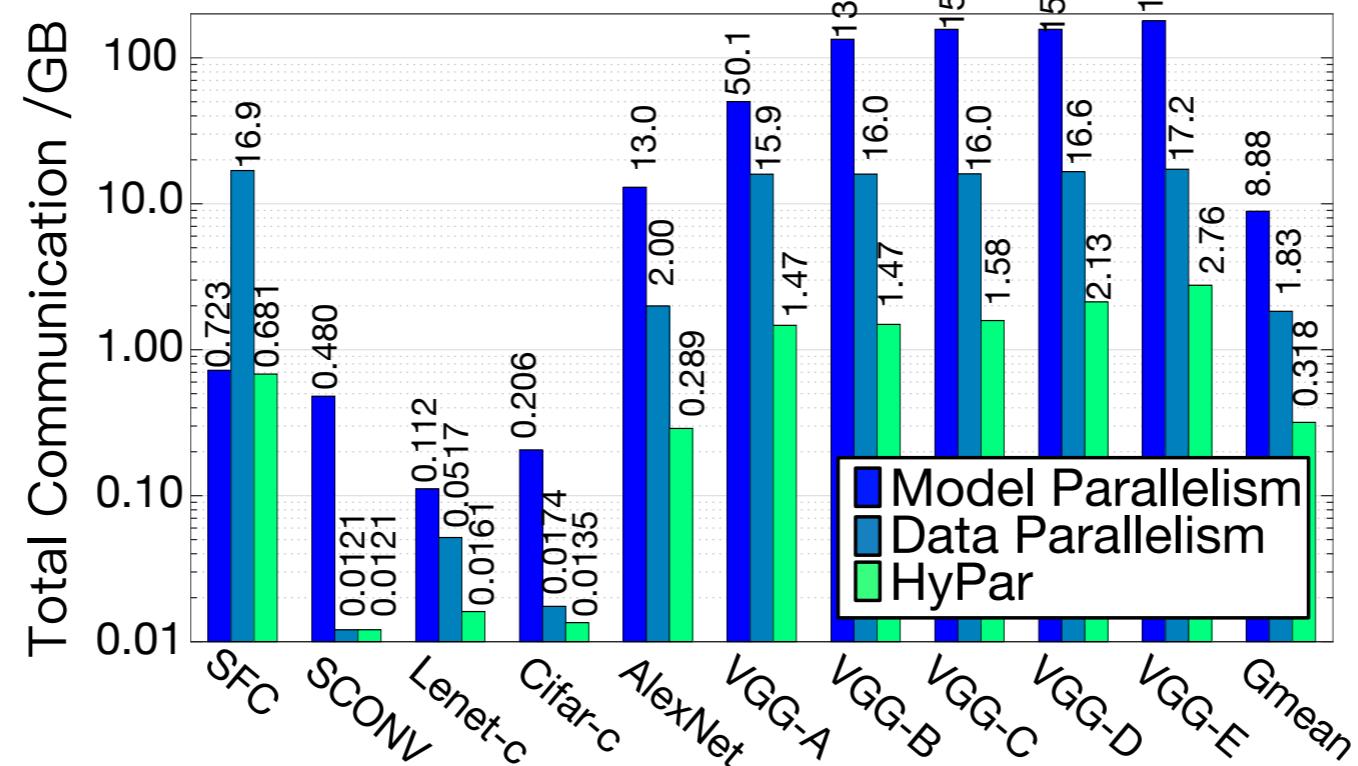
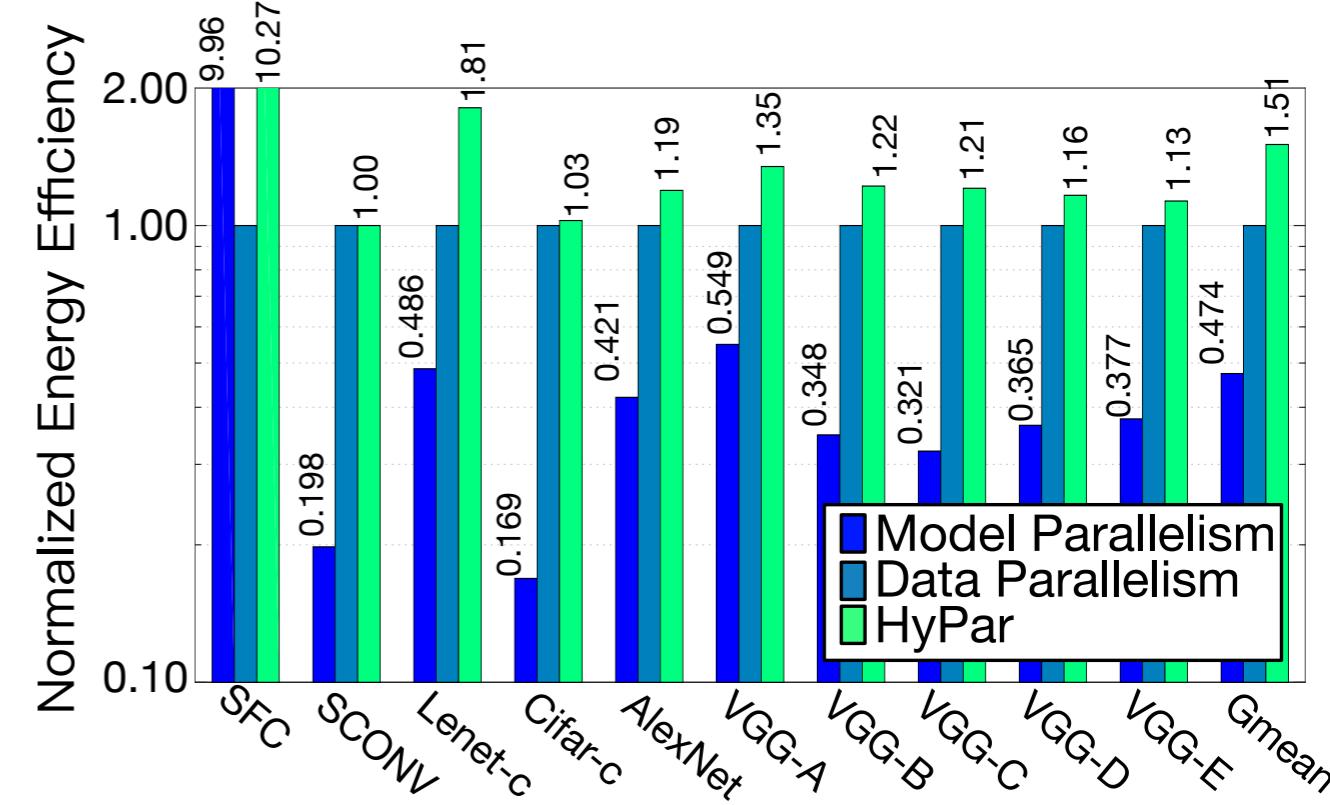
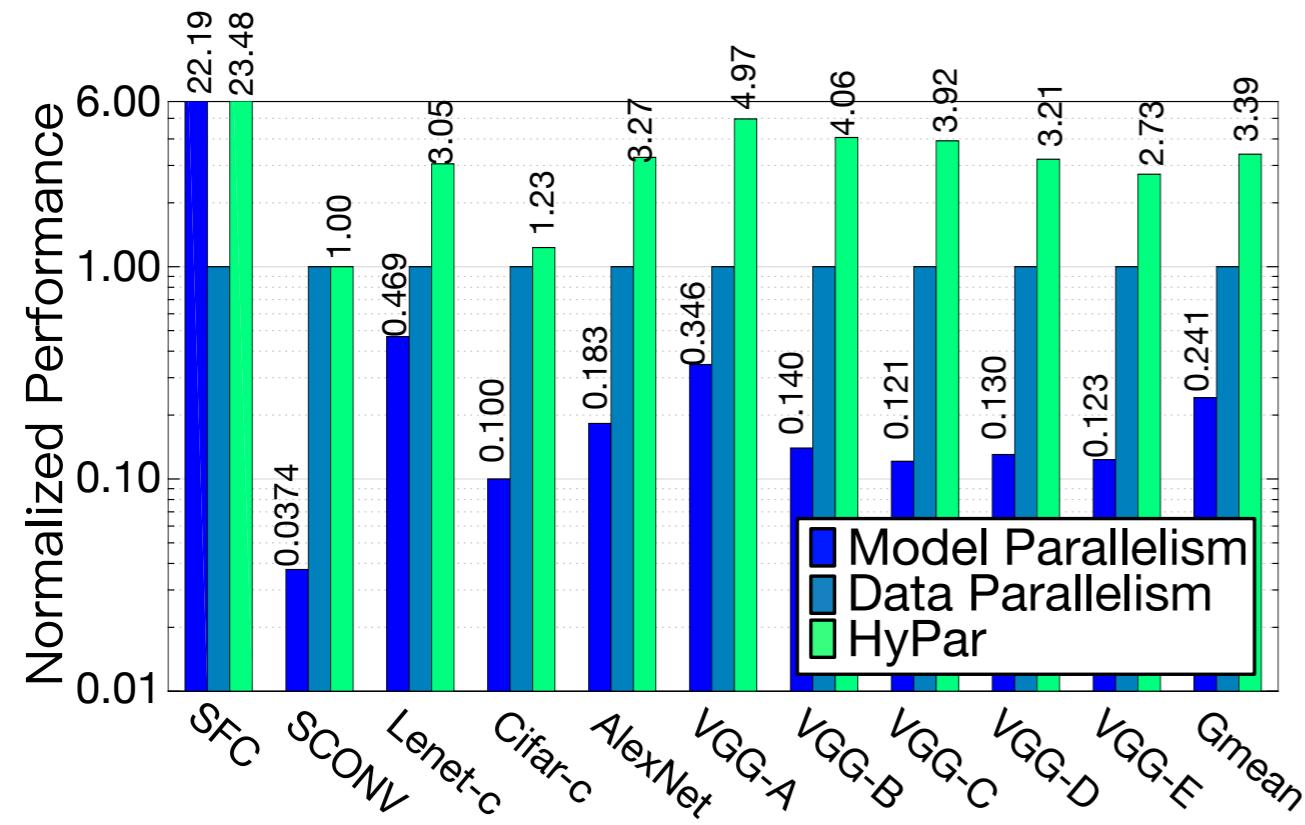


**CEI**

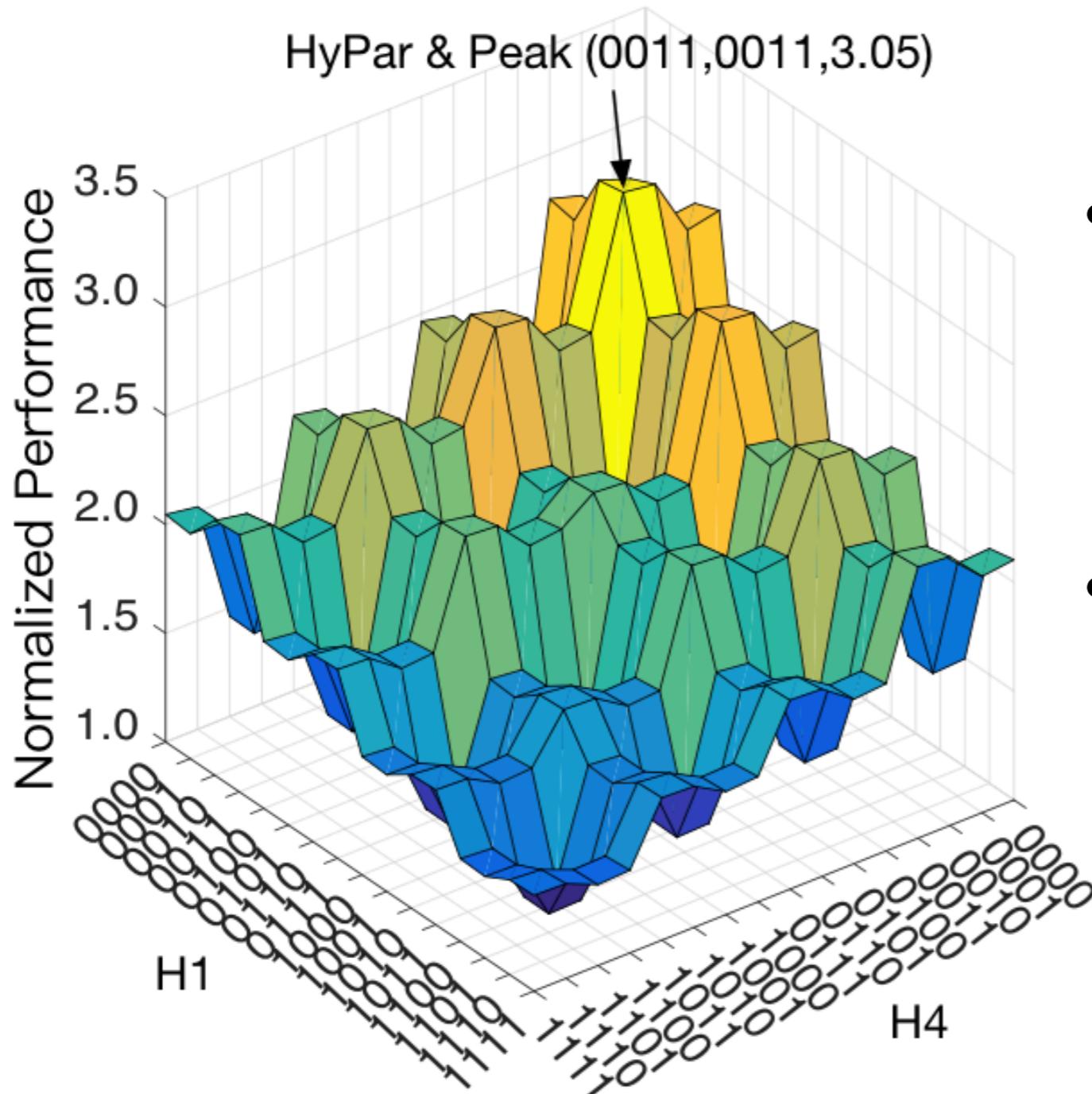
[cei.pratt.duke.edu](http://cei.pratt.duke.edu)

**ALCHEM**  
[alchem.usc.edu](http://alchem.usc.edu)

# HyPar: performance, energy efficiency & communication



# HyPar: parallelism space exploration for Lenet



- We enumerated all 256 possible combinations for H1 and H4.
- HyPar got the maximum as brute force search.

# Outline

---

- The development of deep learning accelerators
- Why HyPar
- HyPar: hybrid parallelism for deep learning accelerator array
  - Communication model
  - Tensor partition
  - Evaluation
- Conclusion

# Conclusion

---

- Optimization within an accelerator is reaching its ceiling.
- Communication cost is high, worth our attention.
- Models to understand communication sources in accelerator array:
  - Intra-layer communication
  - Inter-layer communication
- Dynamic programming and hierarchical partition to minimize the total communication.
- Hybrid parallelism reduces the total communication.

# Discussion: What's next for deep learning accelerators?

---

**Accelerator Design is Guided by Cost**

**Arithmetic is Free  
(particularly low-precision)**

**Memory is expensive**

**Communication is prohibitively expensive**

Bill Dally@AACBB(Sat.)

# Discussion: What's next for deep learning accelerators?

**Accelerator Design is Guided by Cost**

**Arithmetic is Free  
(particularly low-precision)**

**Memory is expensive**

**Communication is prohibitively expensive**

→ NPU, Diannao, etc

Bill Dally@AACBB(Sat.)

# Discussion: What's next for deep learning accelerators?

**Accelerator Design is Guided by Cost**

**Arithmetic is Free  
(particularly low-precision)**

**Memory is expensive**

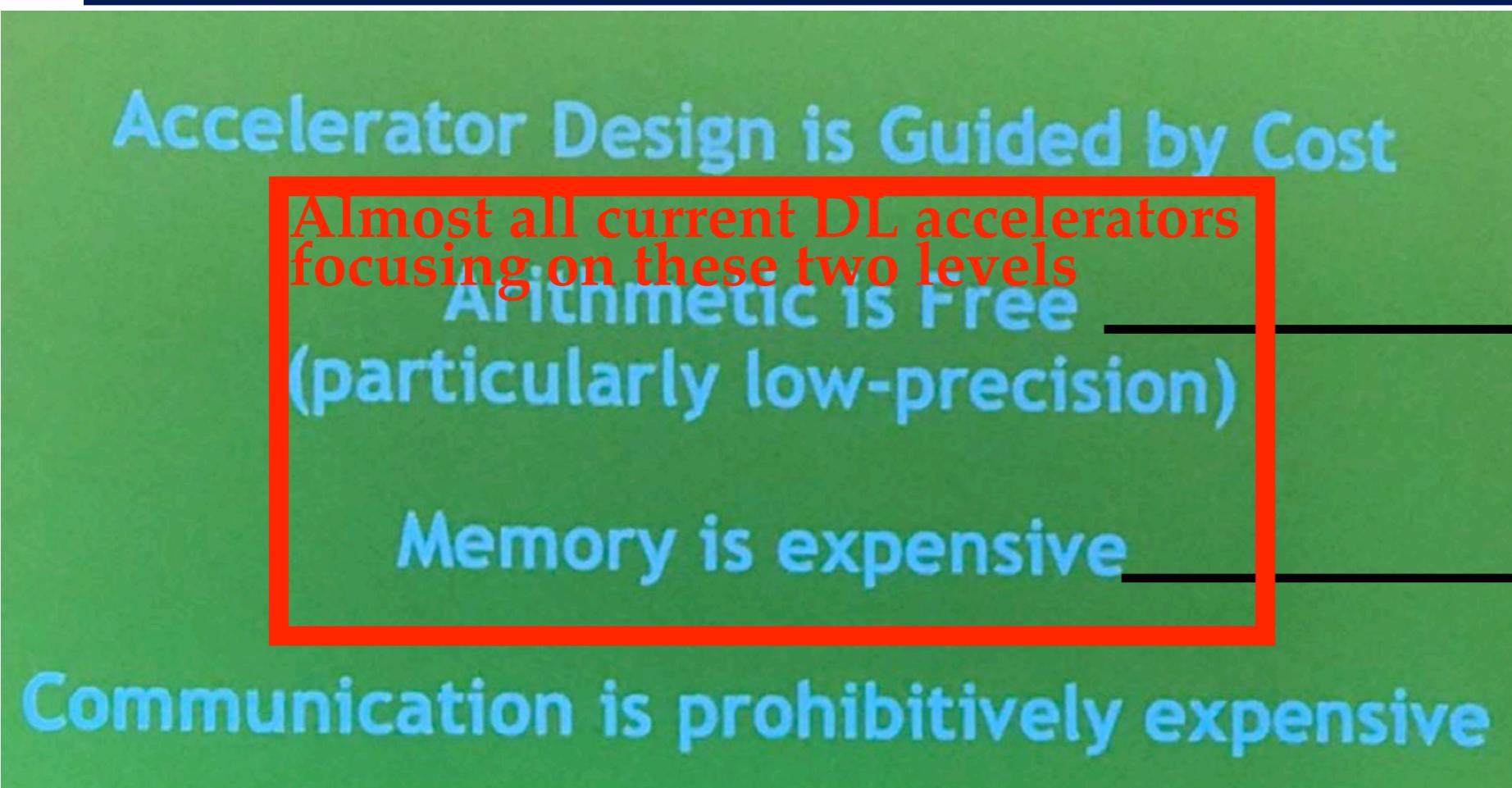
**Communication is prohibitively expensive**

→ NPU, Diannao, etc

→ DaDianao, Eyeriss,  
PipeLayer, etc

Bill Dally@AACBB(Sat.)

# Discussion: What's next for deep learning accelerators?



→ NPU, Diannao, etc

→ DaDianao, Eyeriss, PipeLayer, etc

Bill Dally@AACBB(Sat.)

# Discussion: What's next for deep learning accelerators?

## Accelerator Design is Guided by Cost

Almost all current DL accelerators focusing on these two levels

Arithmetic is Free  
(particularly low-precision)

Memory is expensive

Communication is prohibitively expensive

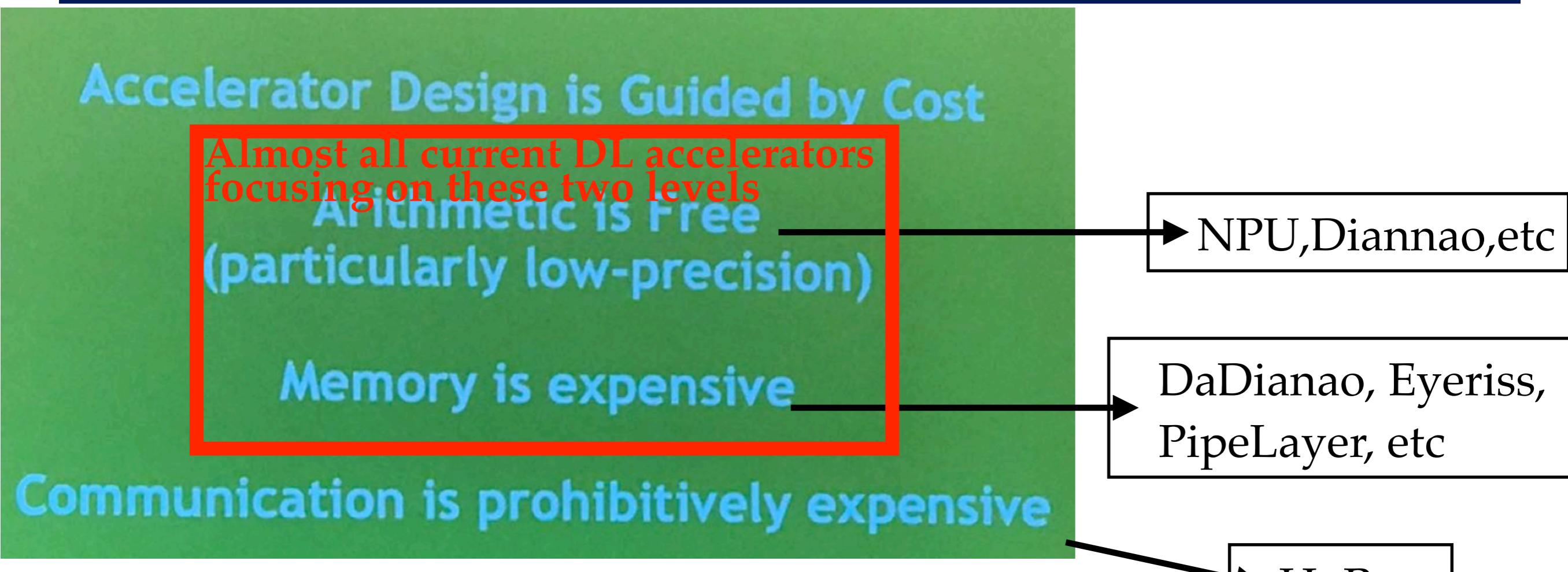
Bill Dally@AACBB(Sat.)

NPU, Diannao, etc

DaDianao, Eyeriss,  
PipeLayer, etc

HyPar

# Discussion: What's next for deep learning accelerators?



Bill Dally@AACBB(Sat.)

- Accelerator array(multi-accelerators)
- Optimization outside of chip, eg., communication

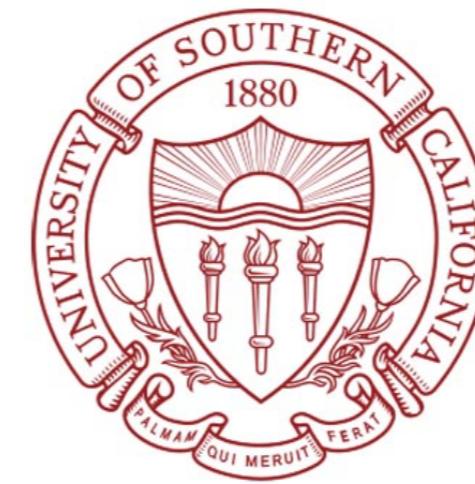
# HyPar: Towards Hybrid Parallelism for Deep Learning Accelerator Array

Linghao Song\*, Jiachen Mao\*, Youwei Zhuo<sup>#</sup>,  
Xuehai Qian<sup>#</sup>, Hai Li\*, Yiran Chen\*

*\*Duke University*

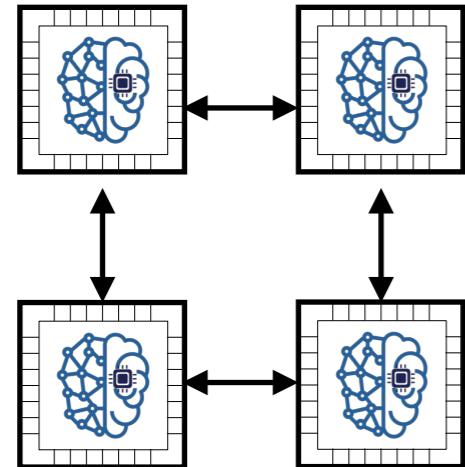
*#University of Southern California*

**CEI**  
[cei.pratt.duke.edu](http://cei.pratt.duke.edu)



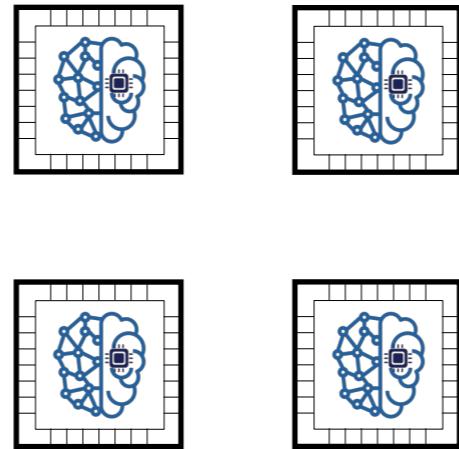
**ALCHEM**  
[alchem.usc.edu](http://alchem.usc.edu)

# Backup1:Layer wise execution model



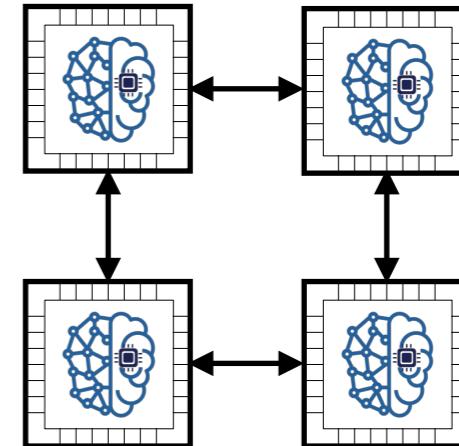
Communicate

Layer (i-1)->i



Compute

Layer i



Communicate

Layer i->(i+1)

Communication dominates!

# Backup2: three tensor computation phases in training

---

Data Forward:

$$\mathbf{F}_{l+1} = f(\mathbf{F}_l \otimes \mathbf{W}_l)$$

Error Backward:

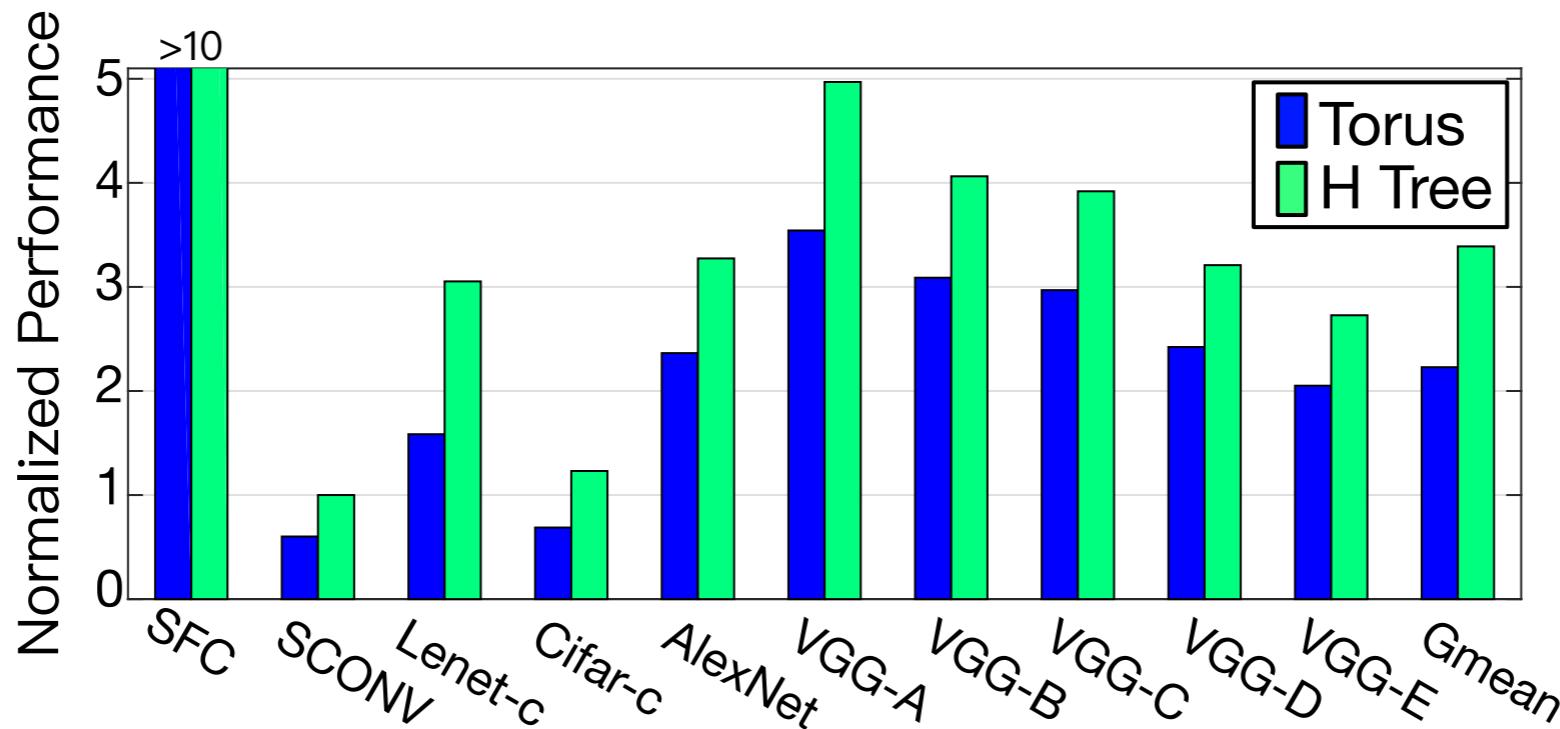
$$\mathbf{E}_l = (\mathbf{E}_{l+1} \otimes \mathbf{W}_l^*) \odot f'(\mathbf{F}_l)$$

Gradient Computation:

$$\Delta \mathbf{W}_l = \mathbf{F}_l^* \otimes \mathbf{E}_{l+1}$$

# Backup3: topologies

HyPar compared to torus:



Hierarchical partition favors H-Tree, and topology-aware partition is required.