

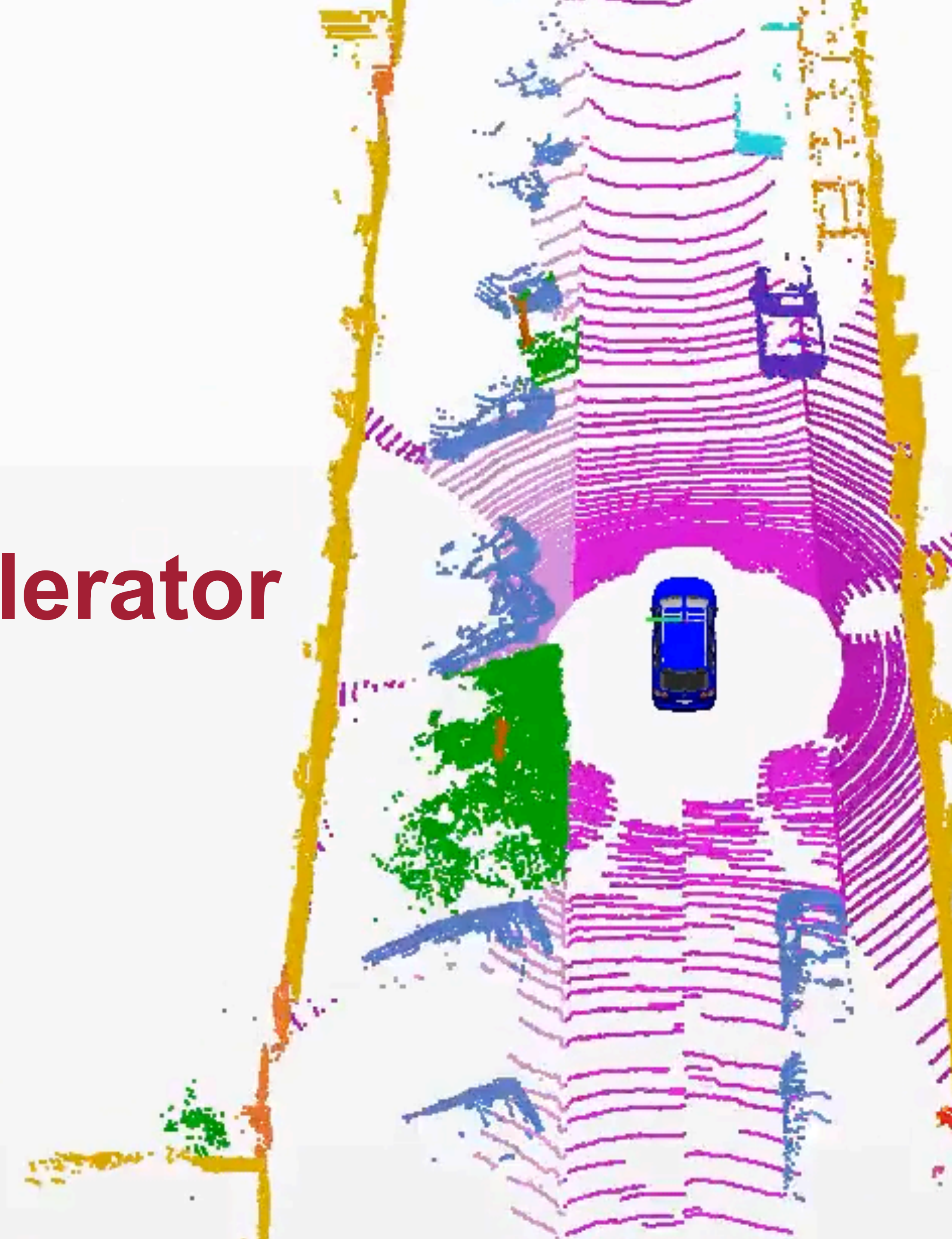
PointAcc

Efficient Point Cloud Accelerator

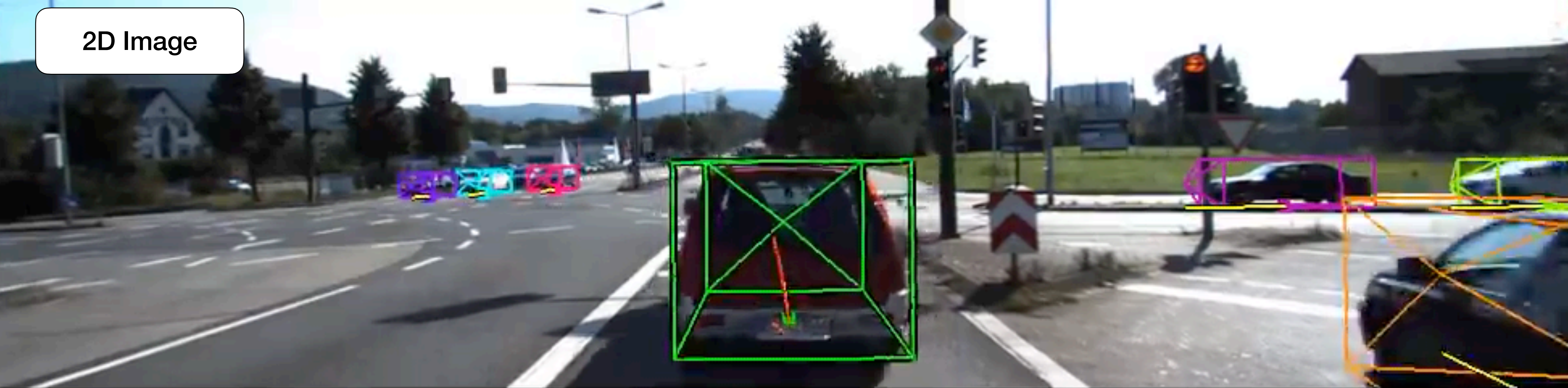
Yujun Lin, Zhekai Zhang, Haotian Tang,
Hanrui Wang, and Song Han

<https://hanlab.mit.edu/projects/pointacc>

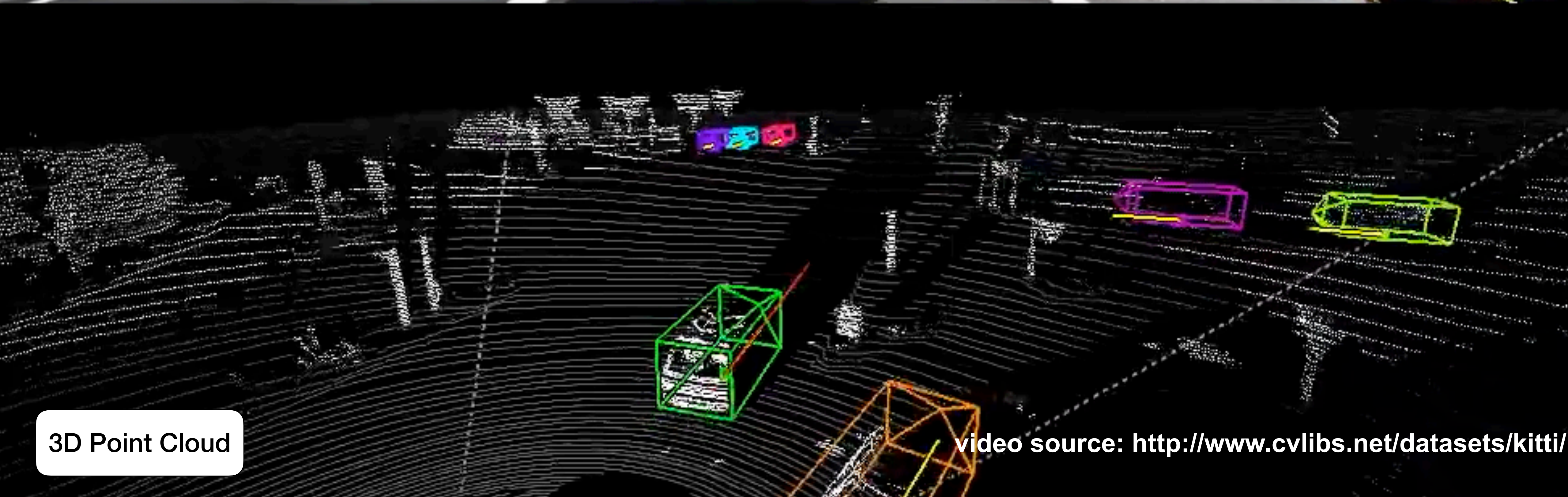
[video source: http://www.semantic-kitti.org/](http://www.semantic-kitti.org/)



2D Image



3D Point Cloud



video source: <http://www.cvlibs.net/datasets/kitti/>

Point Cloud Deep Learning are Everywhere

VR Glasses



Autonomous Driving Vehicles



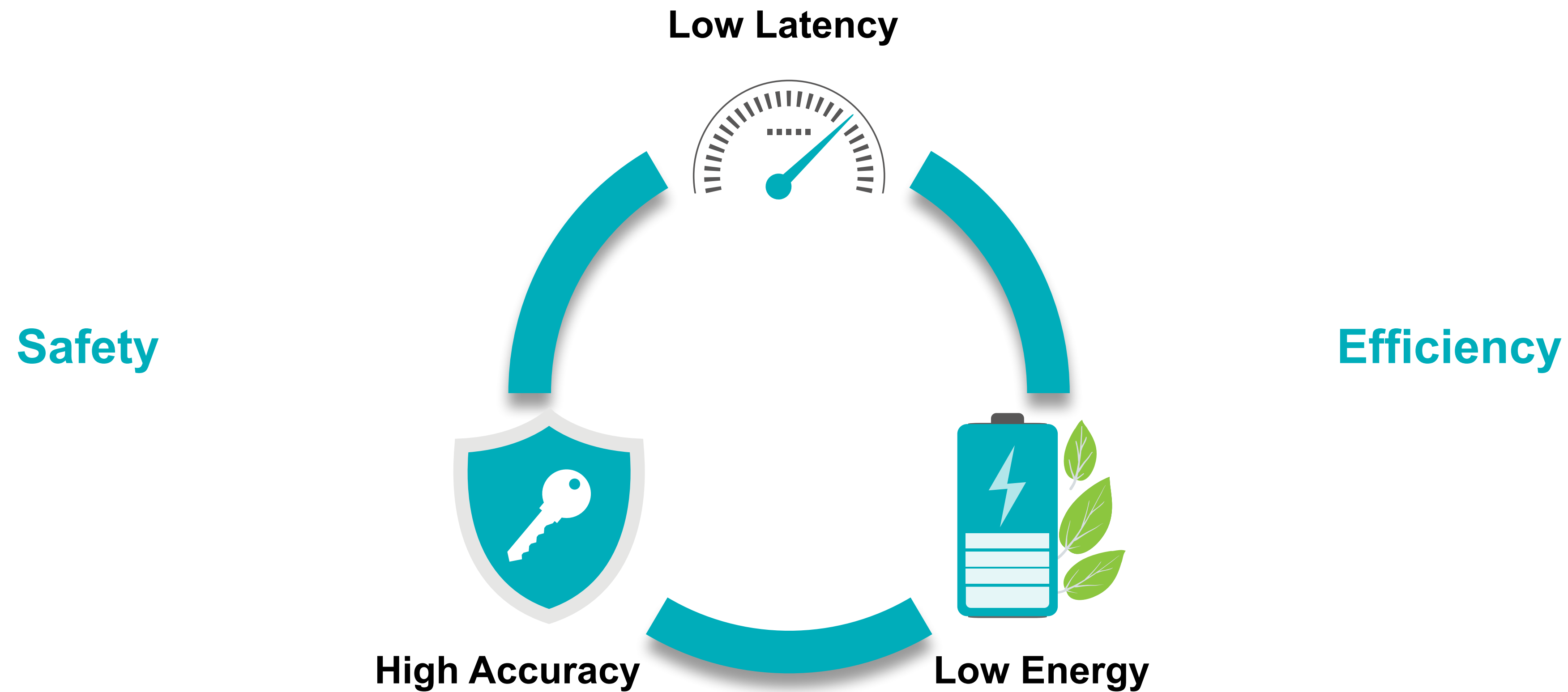
AR iPhones and iPad



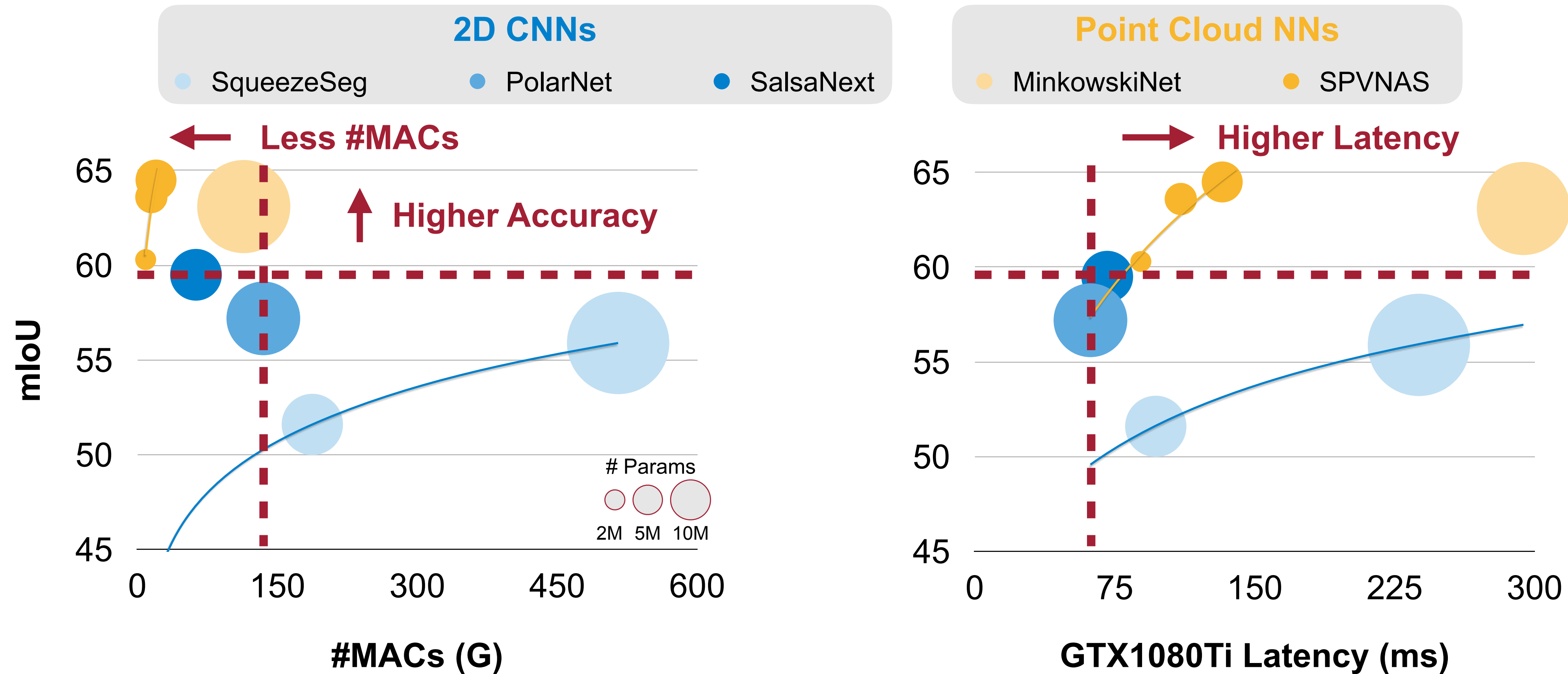
LiDAR Mapping Drones



Efficiency and Safety are Important

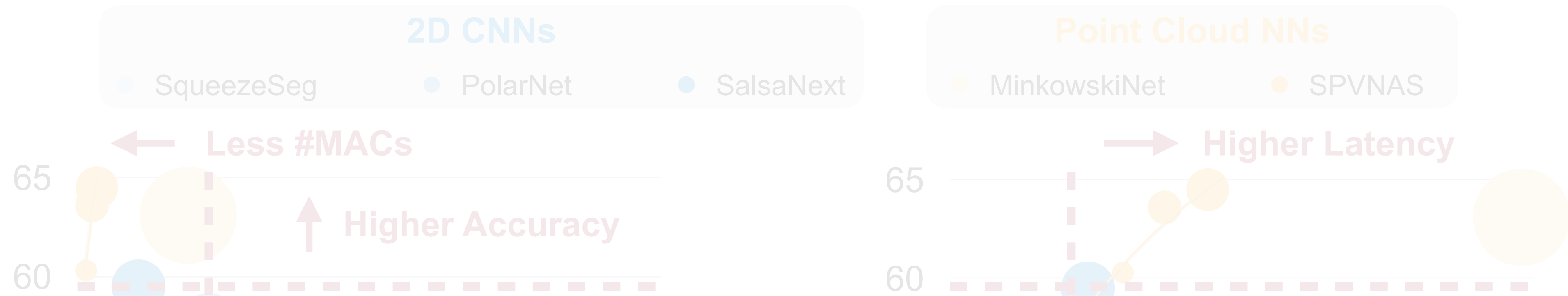


Efficiency and Safety are Important



Point Cloud Networks have **higher accuracy**, but **cannot run efficiently** on today's GPUs (**7X less #MACs but 1.3X slower**).

Efficiency and Safety are Important



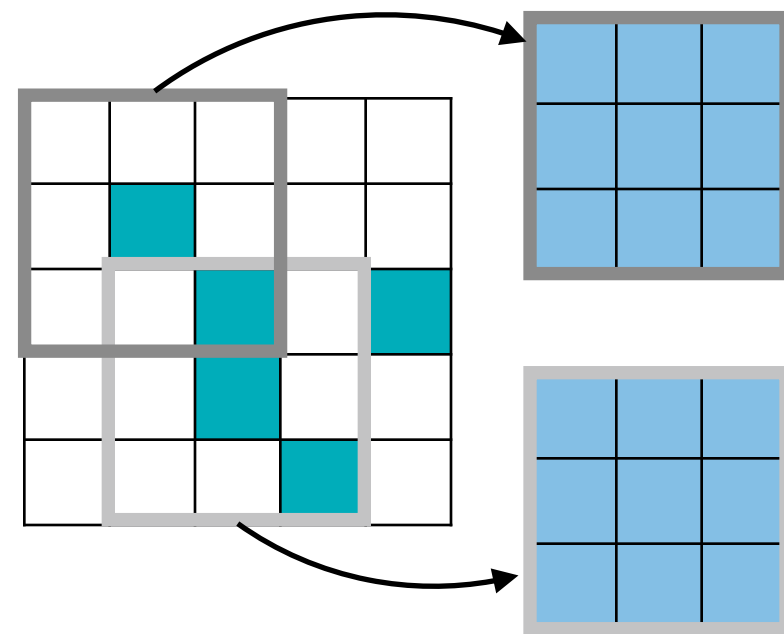
Can existing neural network accelerators solve the point cloud challenge?



Point Cloud Networks have **higher accuracy**, but **cannot run efficiently** on today's GPUs (7X less #MACs but 1.3X slower).

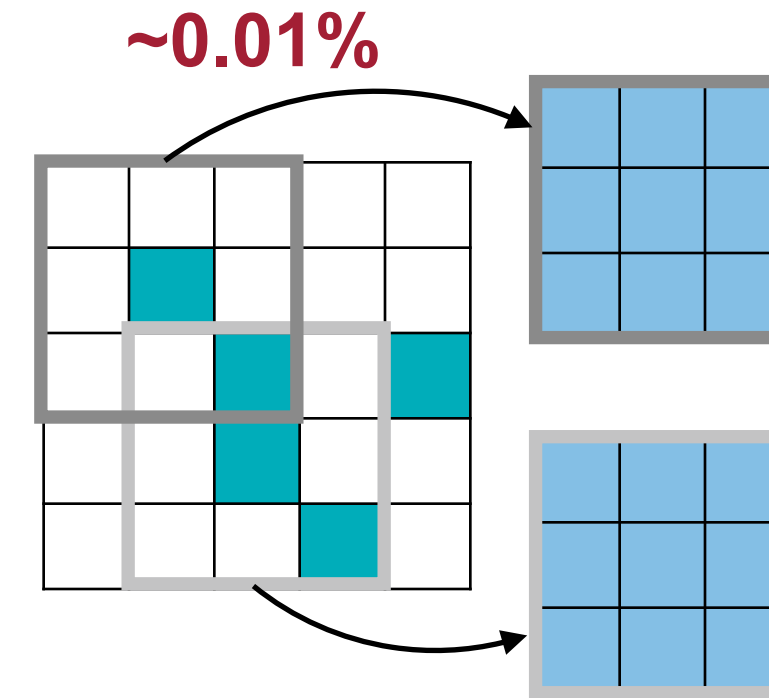
Point Cloud Convolution is Different from Conventional

Conventional Convolution

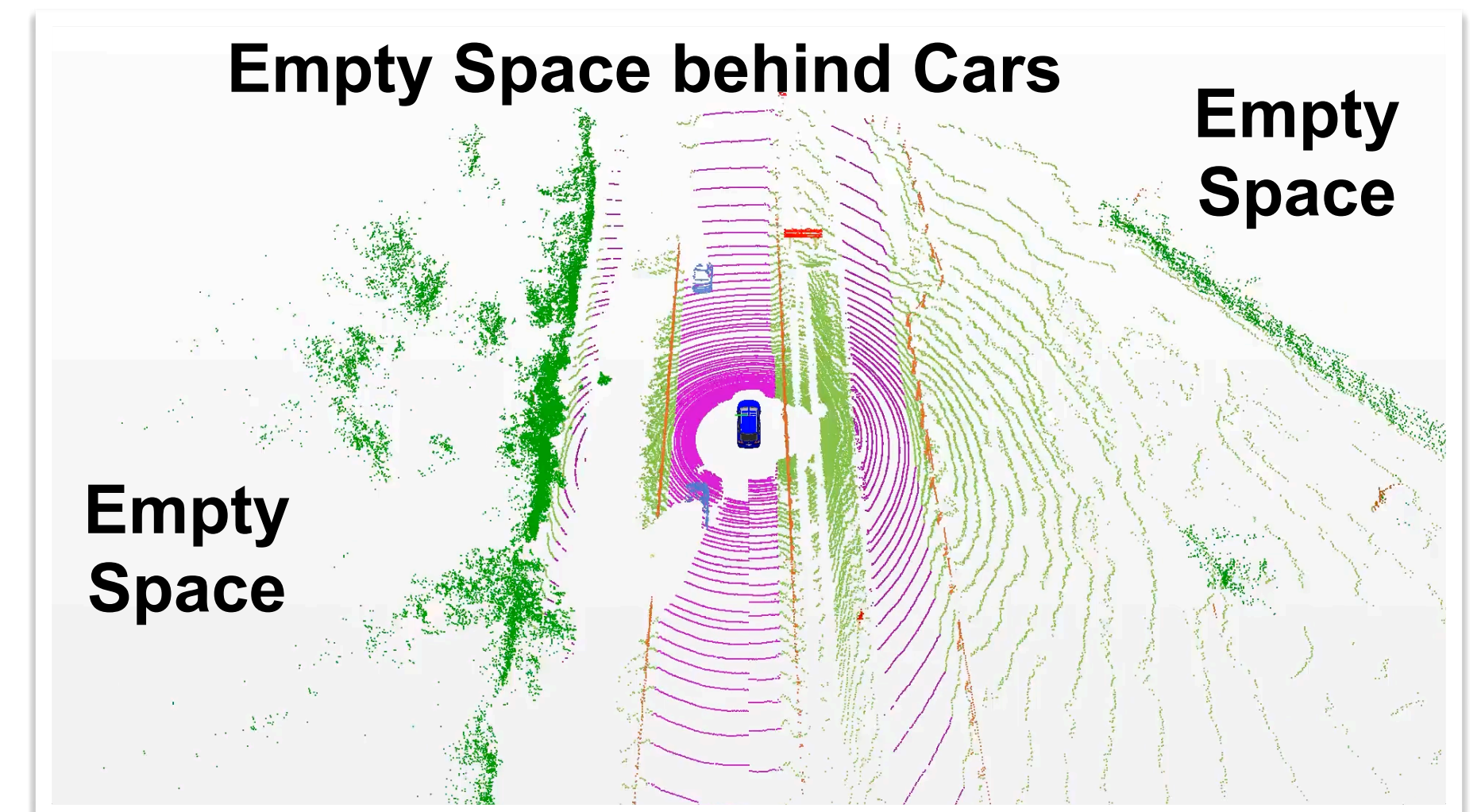
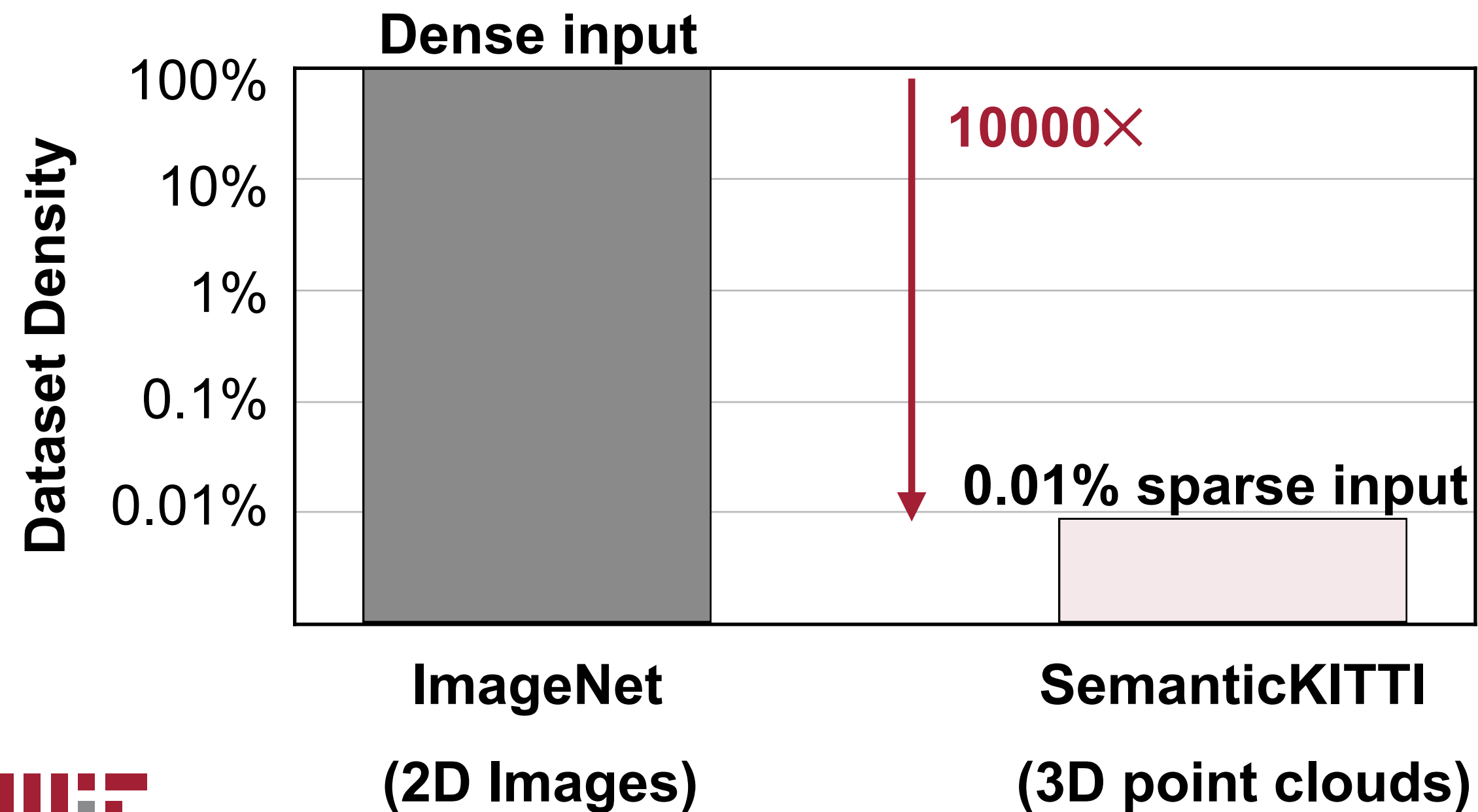


Input sparsity is from ReLU

Point Cloud Convolution



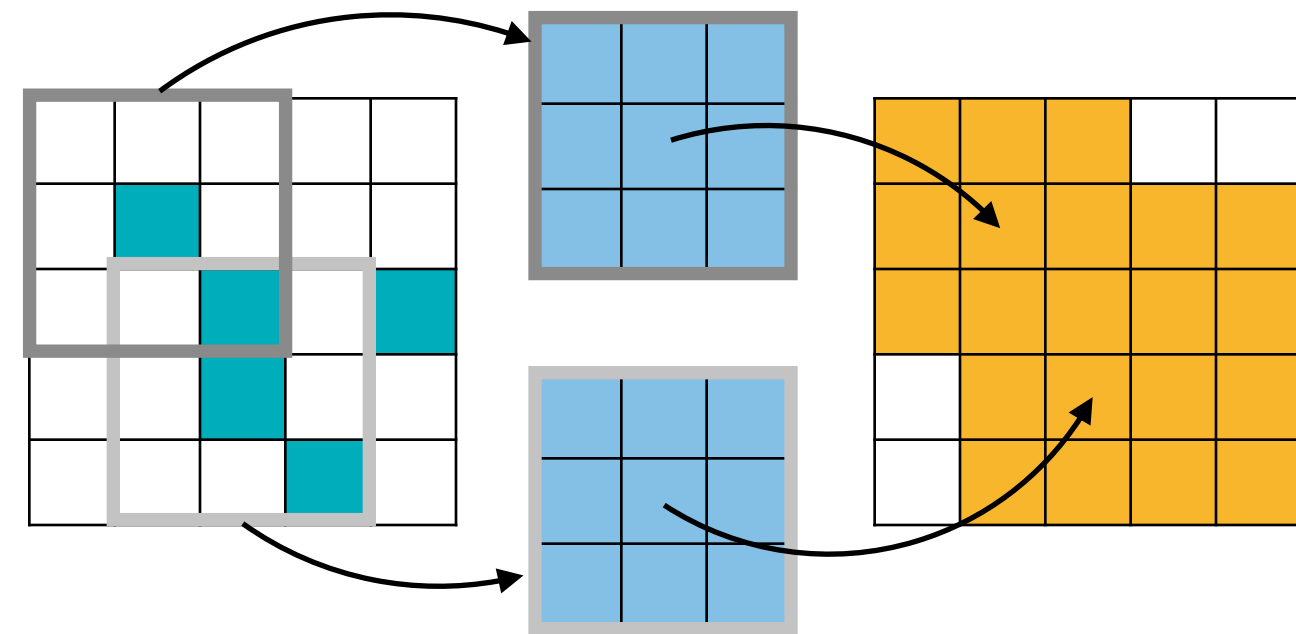
Input sparsity is from the distribution in physical space



video source: <http://www.semantic-kitti.org/>

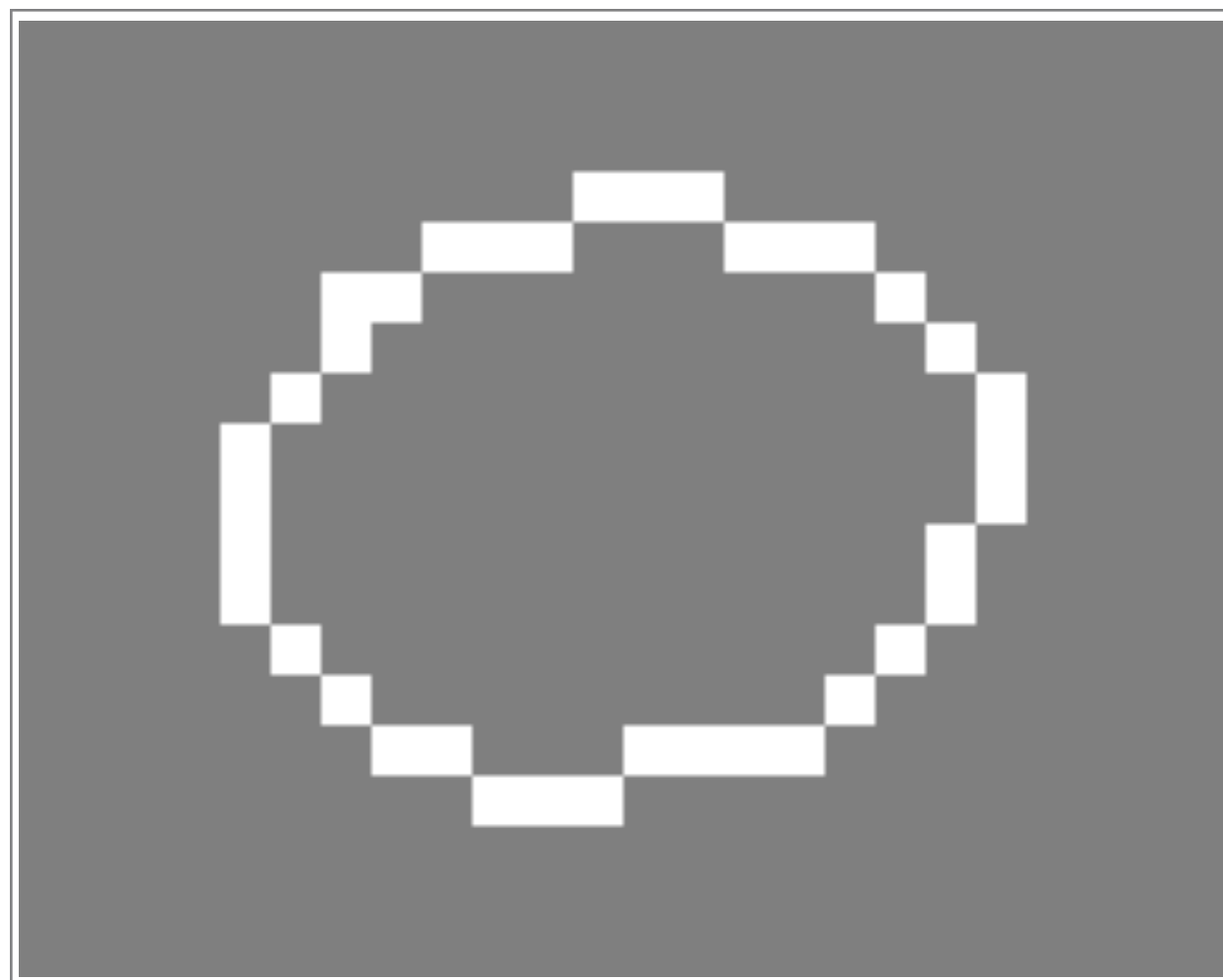
Point Cloud Convolution is Different from Conventional

Conventional Convolution



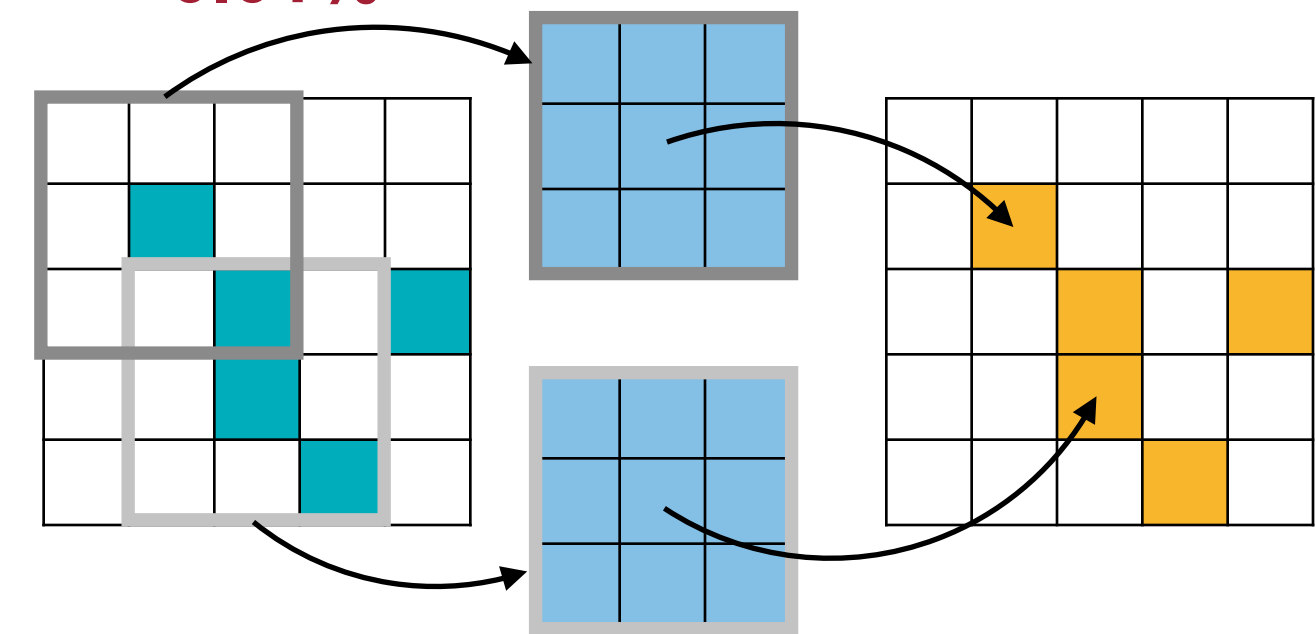
Input sparsity is from ReLU

Nonzeros will dilate



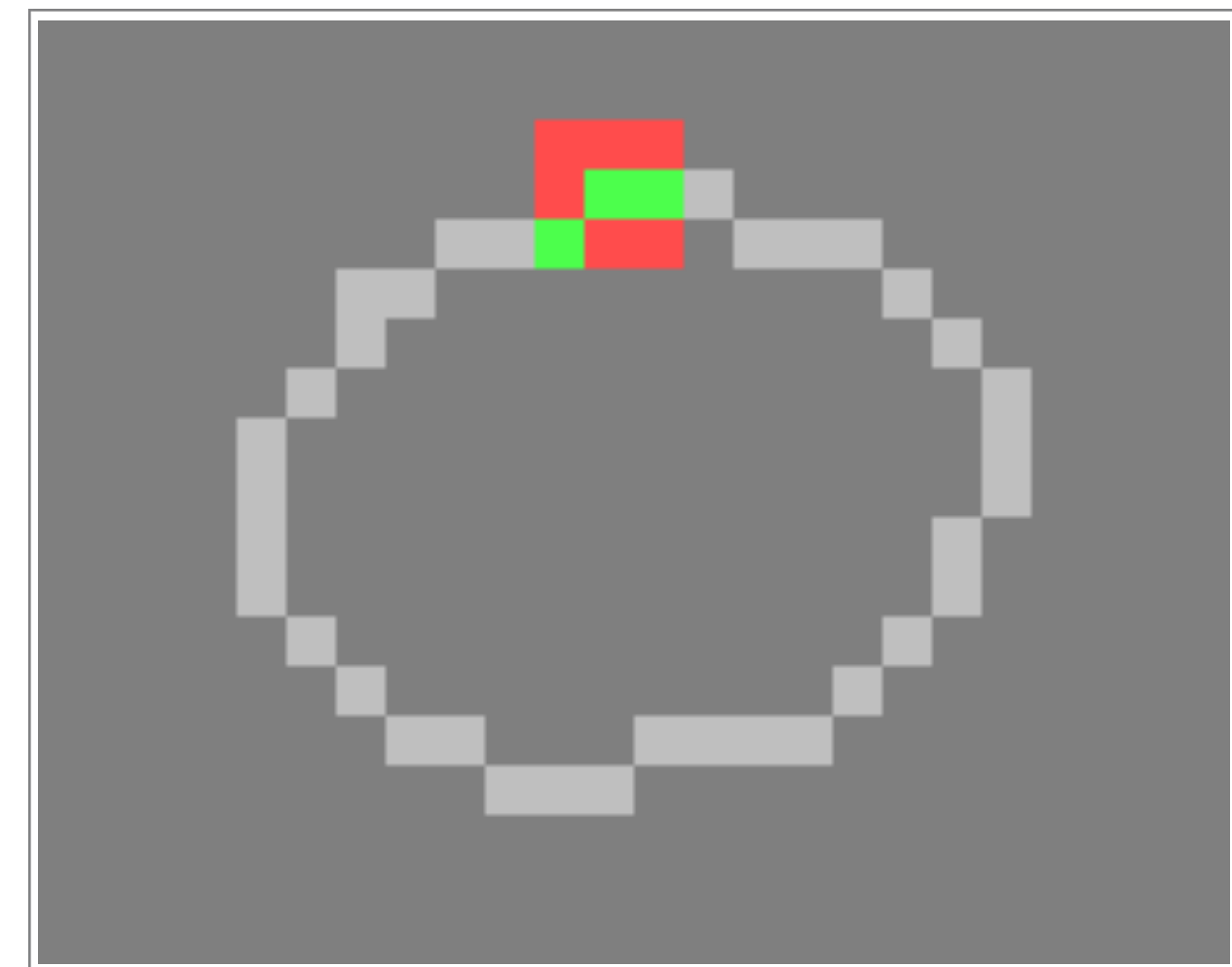
Point Cloud Convolution

~0.01%



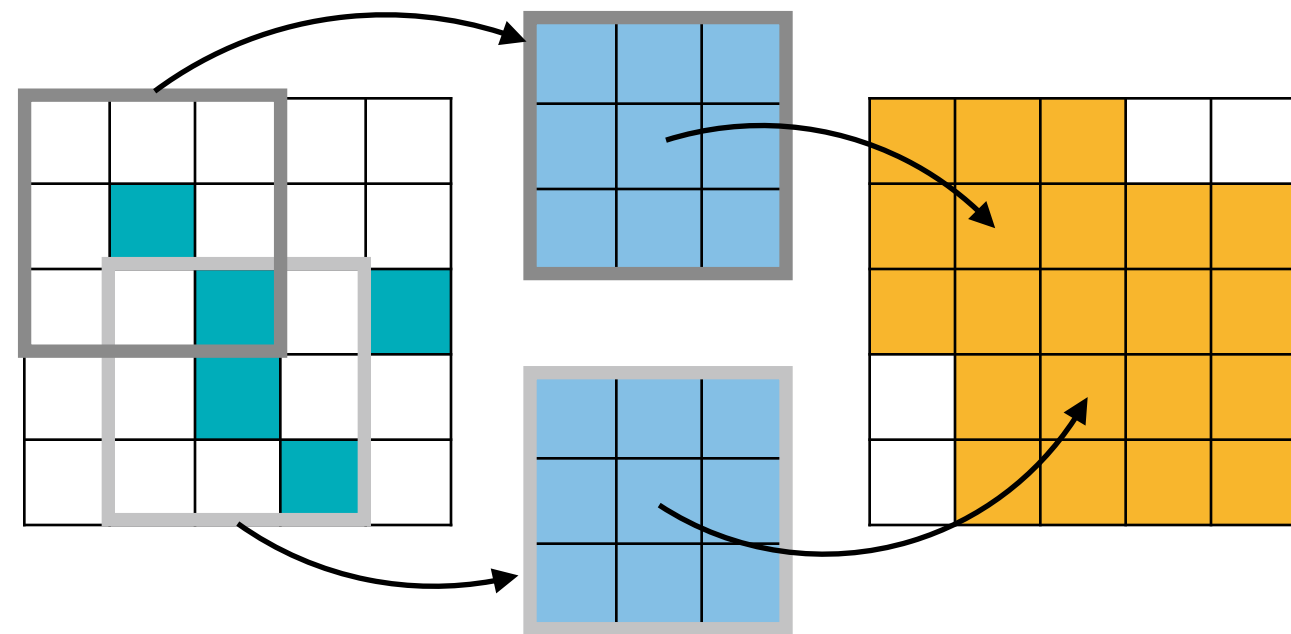
Input sparsity is from the distribution in physical space

Nonzeros will not dilate



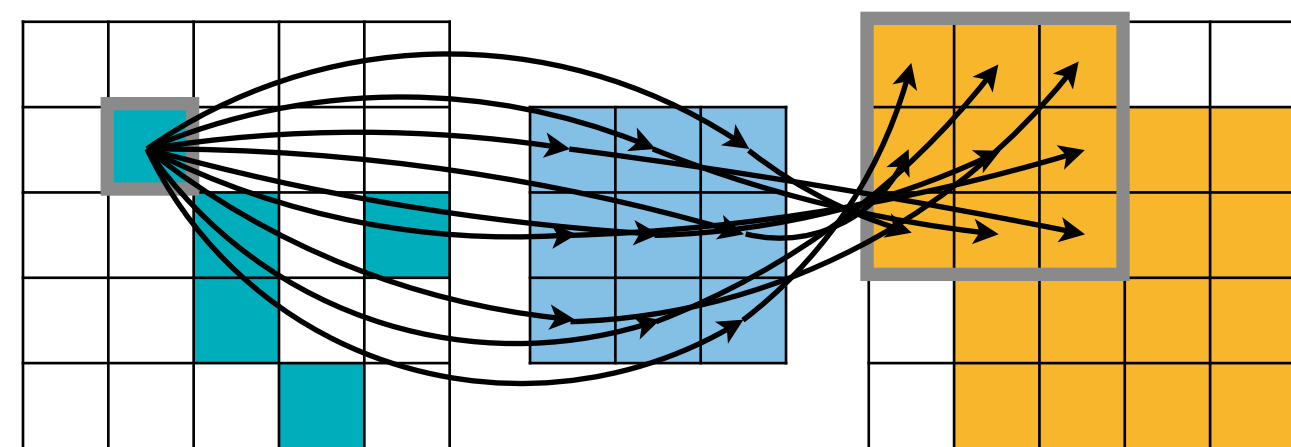
Point Cloud Convolution is Different from Conventional

Conventional Convolution



Input sparsity is from ReLU

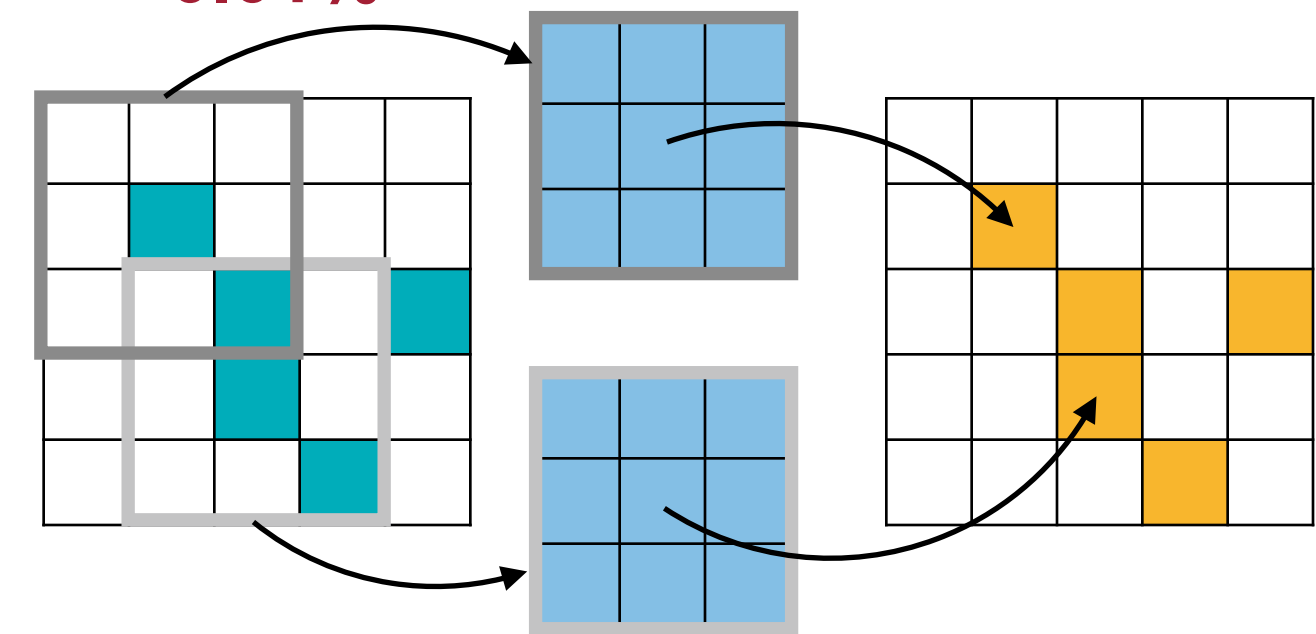
Nonzeros will dilate



Each nonzero input is multiplied with all nonzero weights

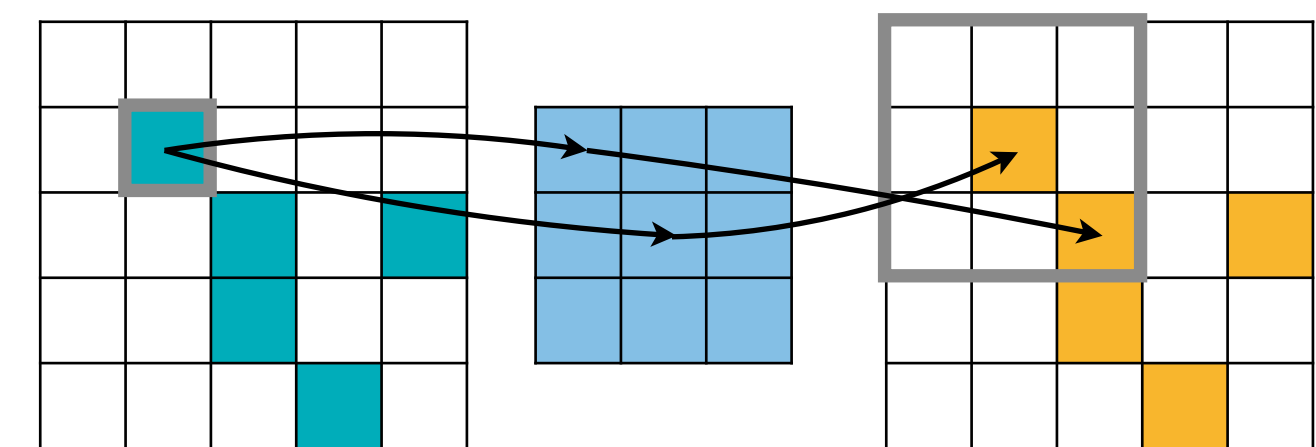
Point Cloud Convolution

~0.01%



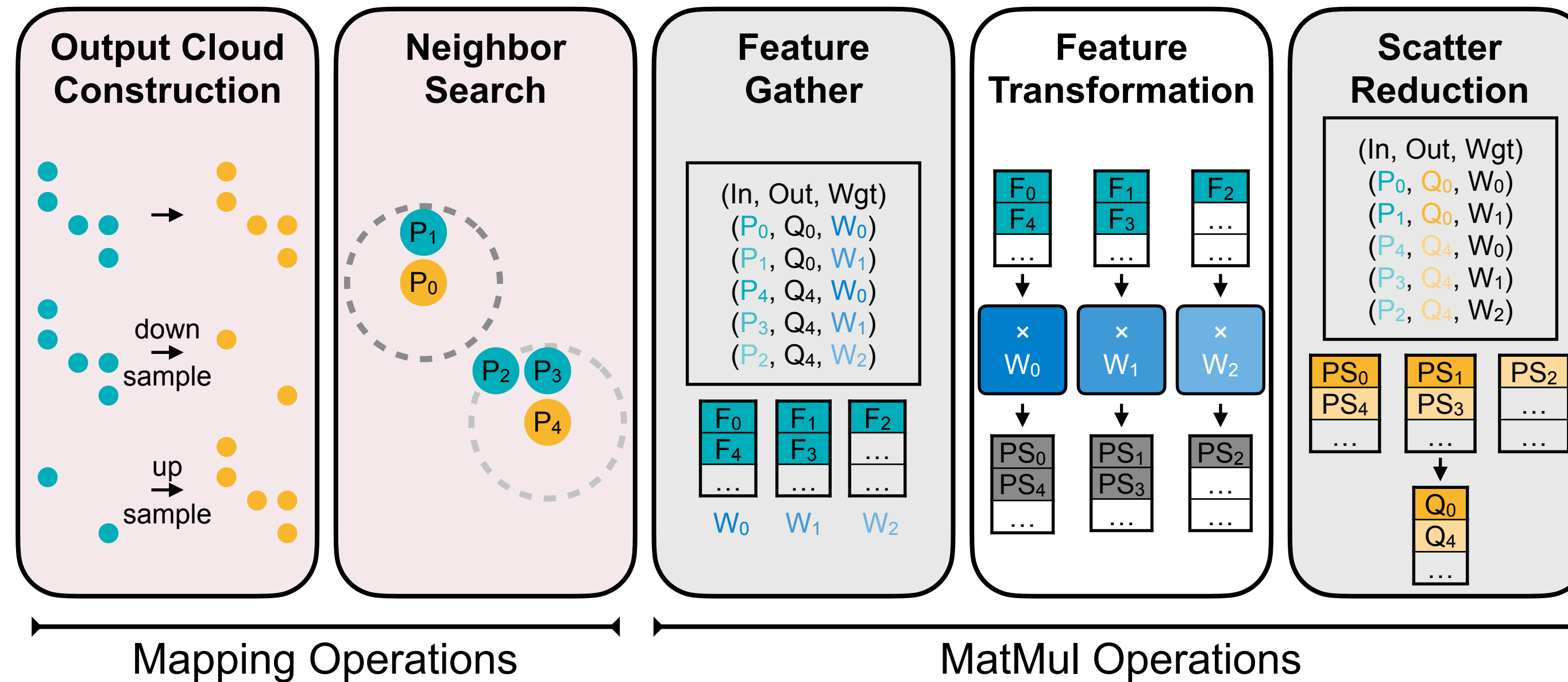
Input sparsity is from the distribution in physical space

Nonzeros will not dilate



Each nonzero input is not multiplied with all nonzero weights

Point Cloud Convolution is Different from Conventional



Sparsity in Physical Space
 Input sparsity is from the point distribution in physical space

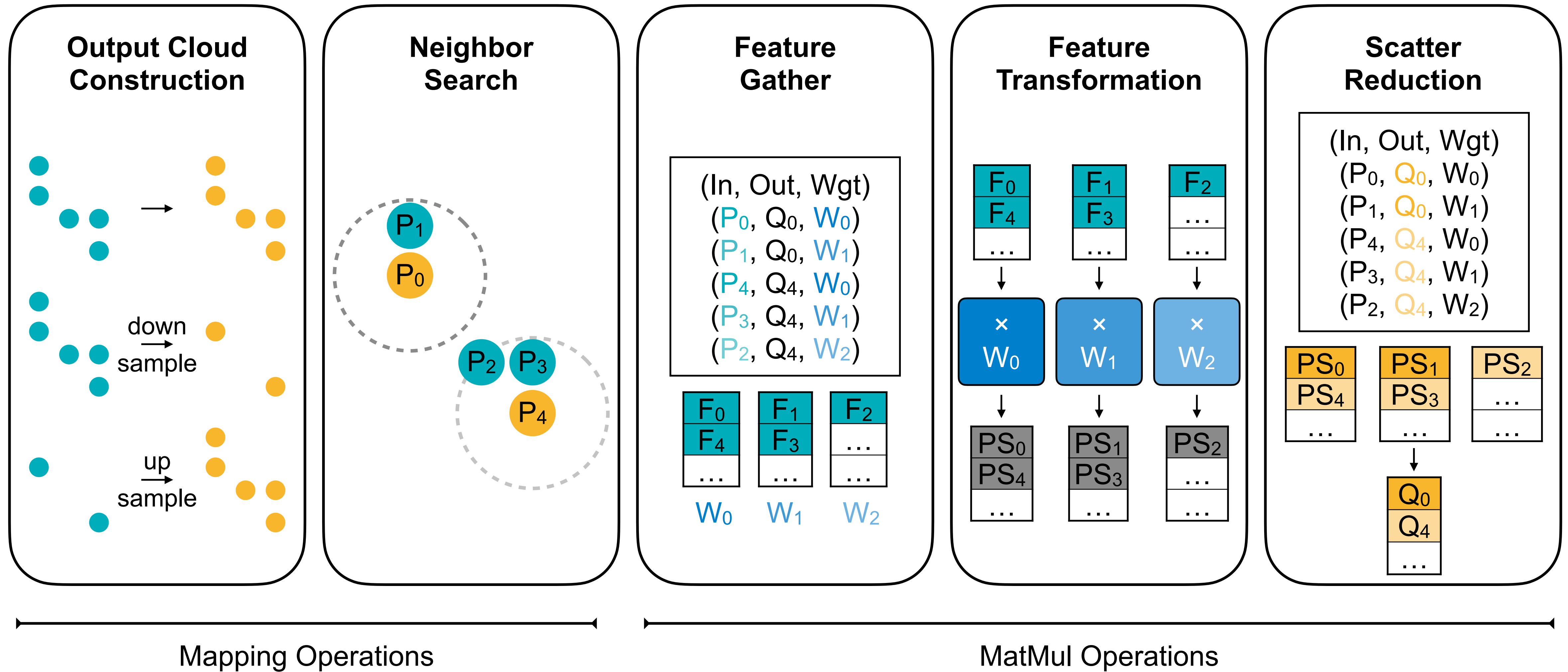
Nonzeros will not dilute

Irregular and sparse computation pattern
 Each nonzero input are not multiplied with all nonzero weights

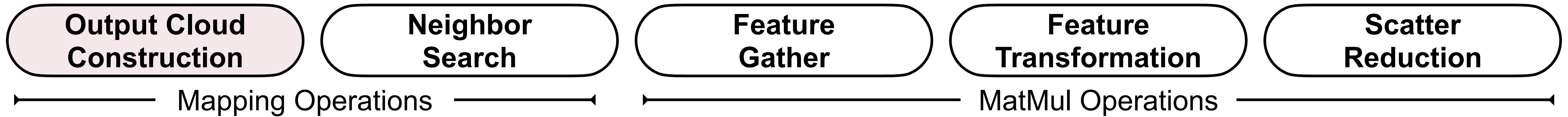
New operations to find neighbors
 Mapping operations

Explicit data movement overhead
 Gather features and Scatter psums

Introduction to Point Cloud Convolution



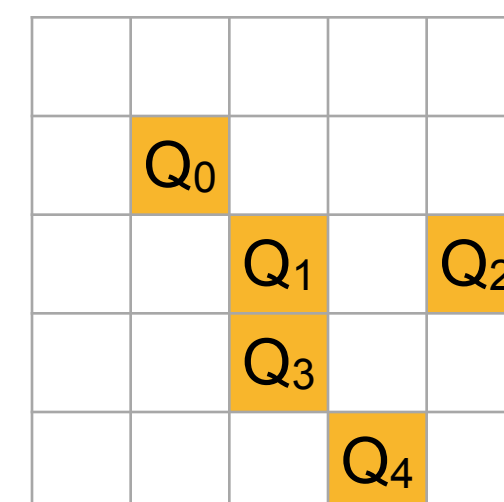
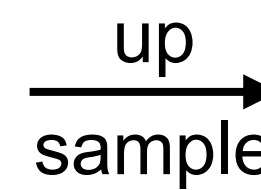
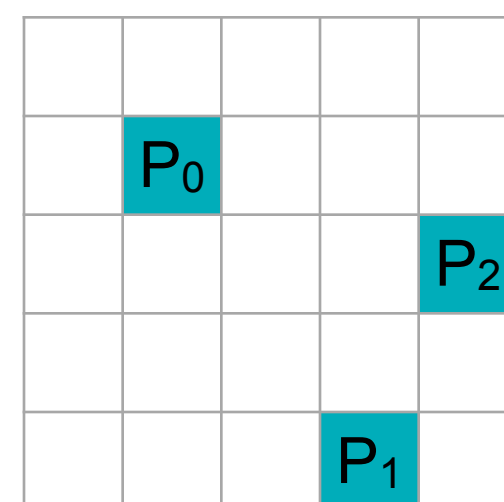
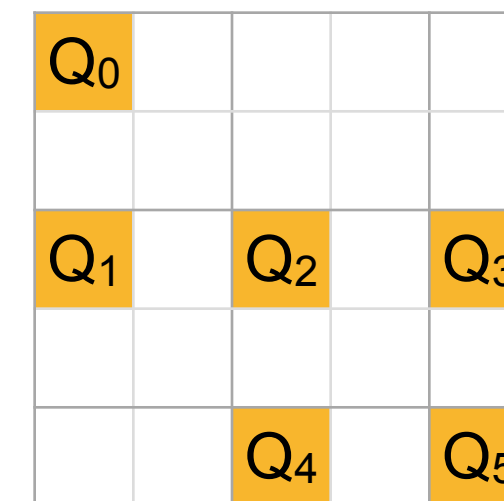
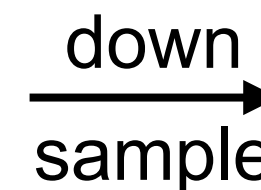
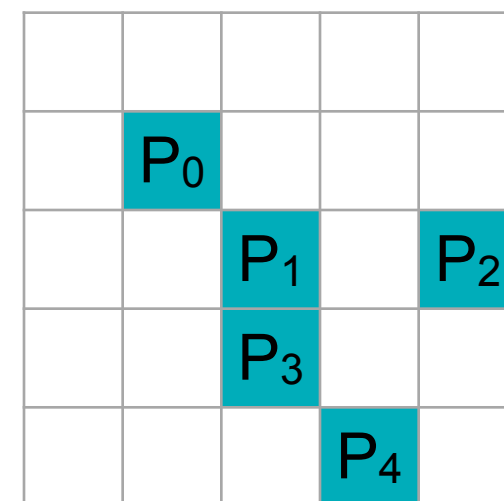
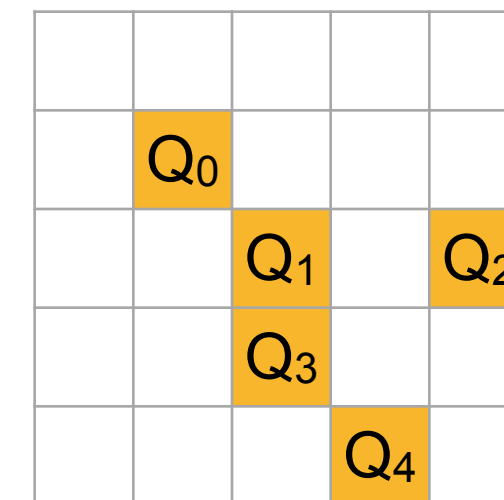
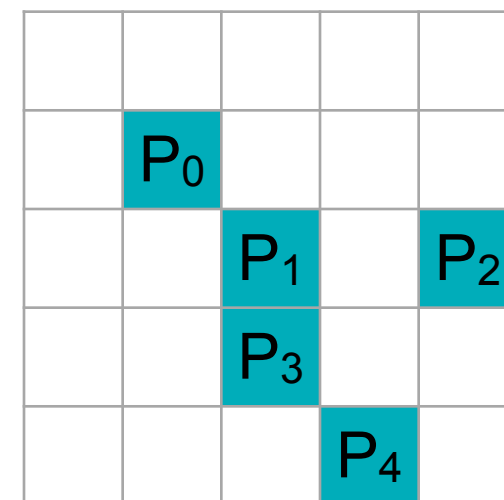
Step 1: Build Output Point Cloud



Input Point Cloud

(P_0, F_0)
(P_1, F_1)
(P_2, F_2)
(P_3, F_3)
(P_4, F_4)

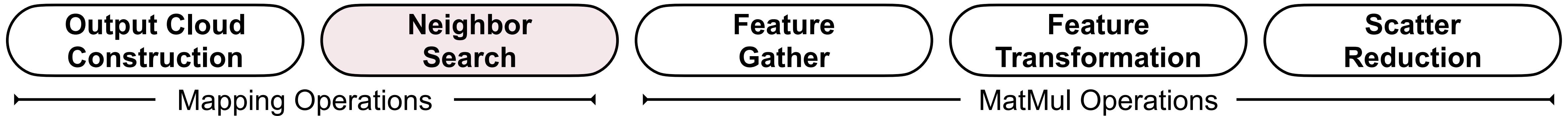
(Coords, Feature Vector)



Output Point Cloud

Q_0
Q_1
Q_2
Q_3
Q_4

Step 2: Search Neighbors

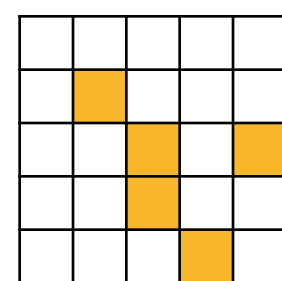
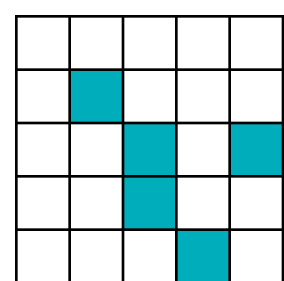


For each *output* point, we find the *input* points in its neighborhood and generate **(Input Point, Output Point, Weight Index)** maps

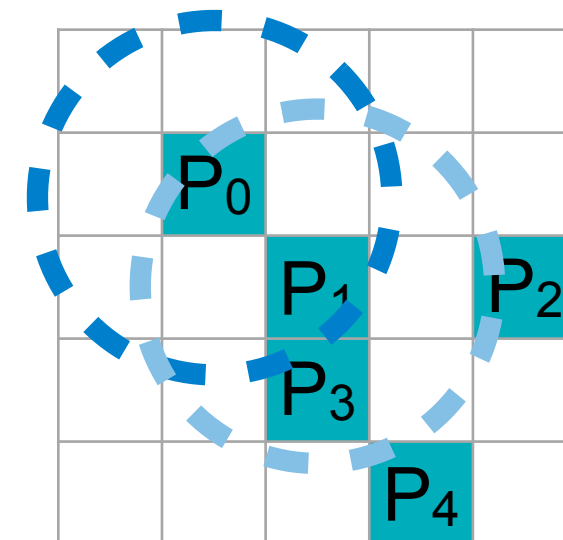
Input Point Cloud

(P_0, F_0)
(P_1, F_1)
(P_2, F_2)
(P_3, F_3)
(P_4, F_4)

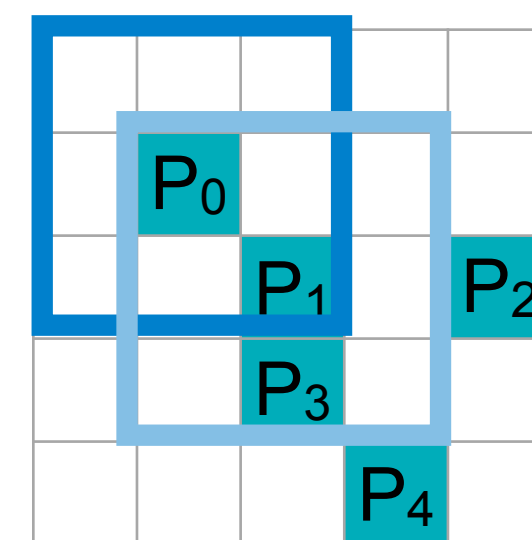
(Coords, Feature Vector)



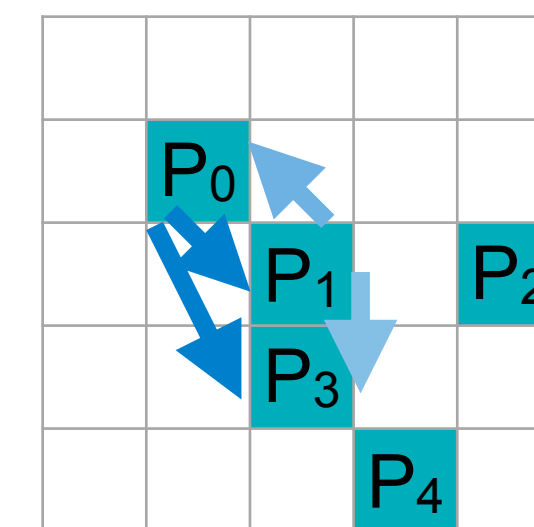
Neighborhood Shape		
Ball	Cube	Nearest Neighbor
Ball Query	Kernel Mapping	k-Nearest-Neighbors



(In, Out, Wgt)
(P_0, Q_0, W_0)
(P_1, Q_0, W_1)
(P_0, Q_1, W_0)
(P_3, Q_1, W_1)
...



(In, Out, Wgt)
$(P_0, Q_0, W_{0,0})$
$(P_1, Q_0, W_{1,1})$
$(P_0, Q_1, W_{-1,-1})$
$(P_1, Q_1, W_{0,0})$
$(P_3, Q_1, W_{1,0})$

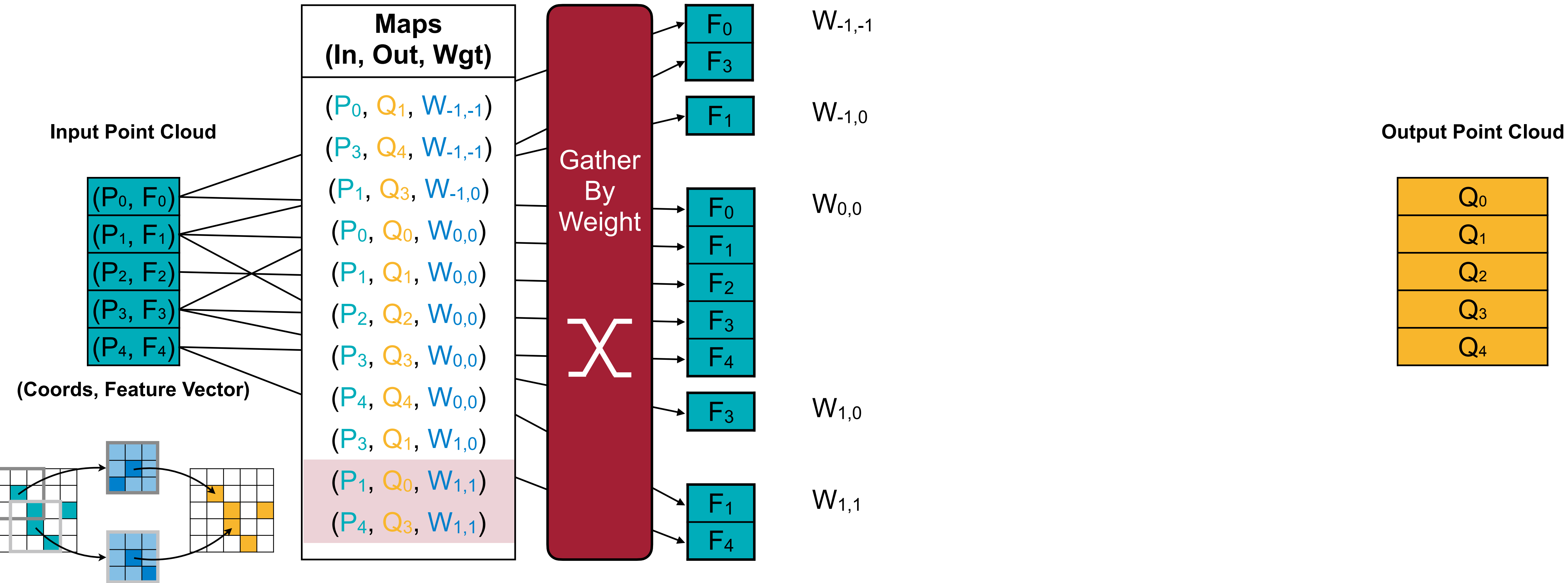
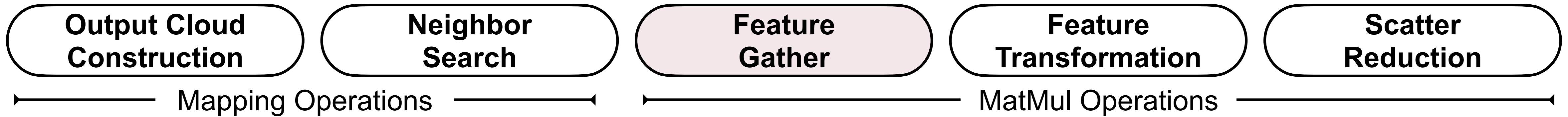


(In, Out, Wgt)
(P_0, Q_0, W_0)
(P_1, Q_0, W_1)
(P_3, Q_0, W_2)
(P_1, Q_1, W_0)
(P_3, Q_1, W_1)

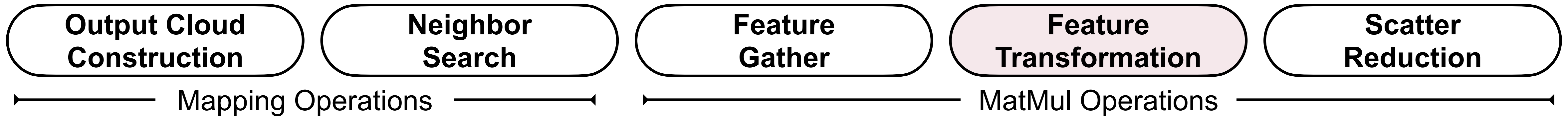
Output Point Cloud

Q_0
Q_1
Q_2
Q_3
Q_4

Step 3: Gather Input Features



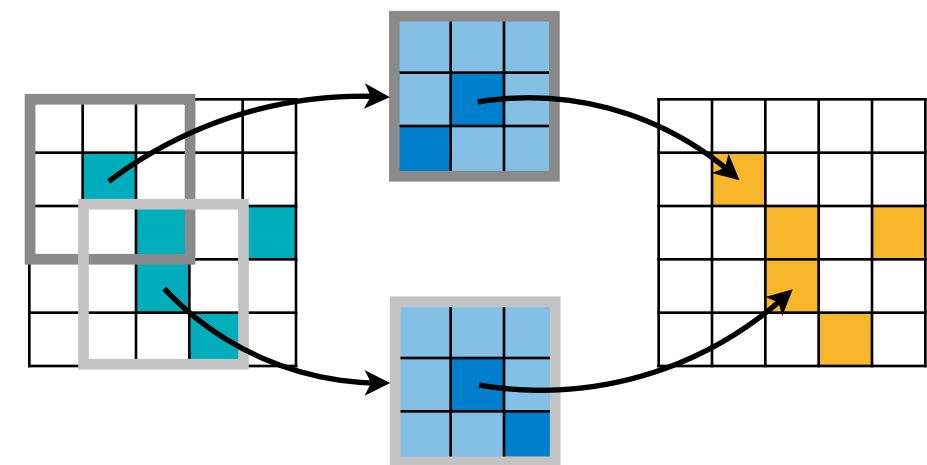
Step 4: Transform Input Features



Input Point Cloud

(P ₀ , F ₀)
(P ₁ , F ₁)
(P ₂ , F ₂)
(P ₃ , F ₃)
(P ₄ , F ₄)

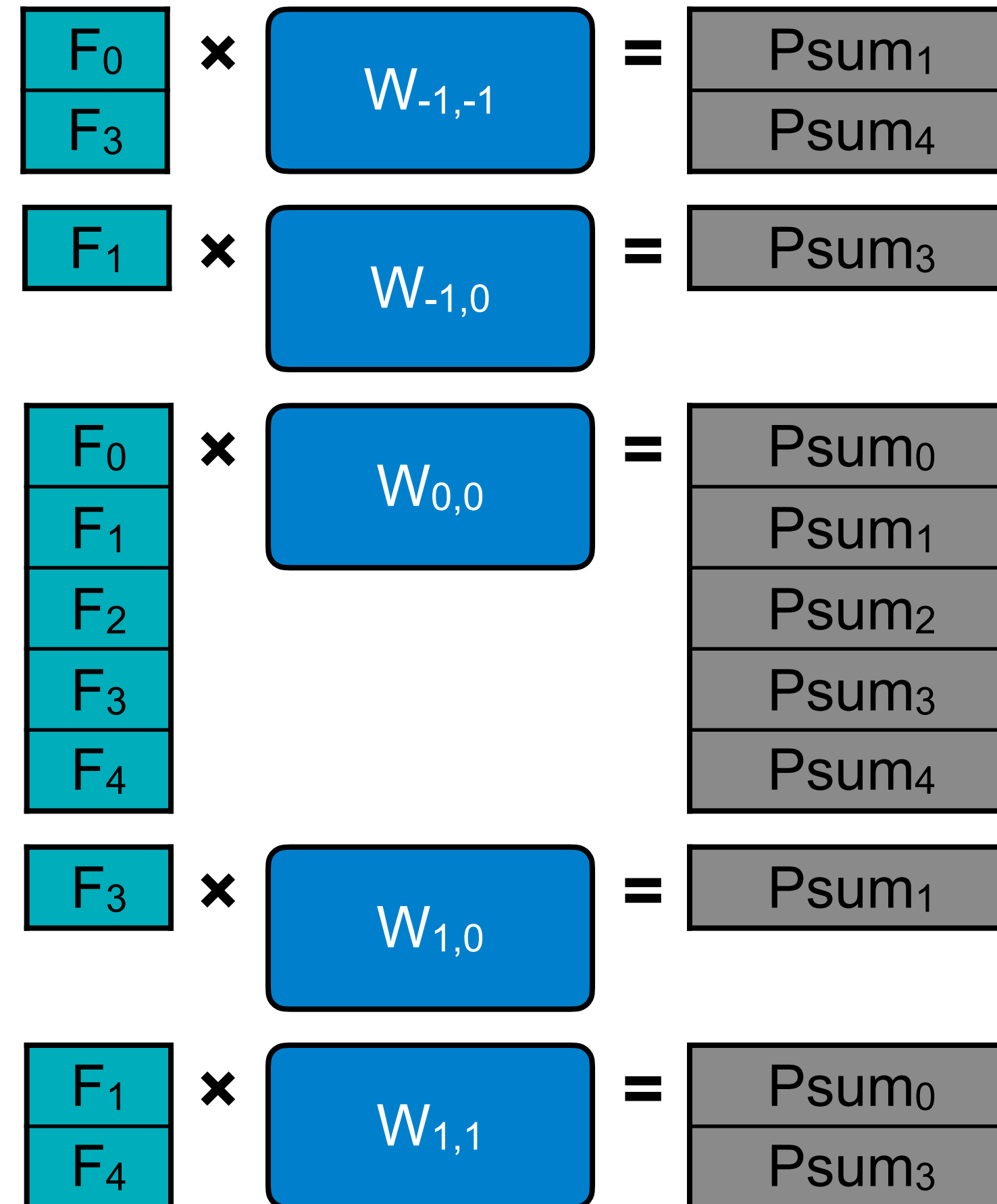
(Coords, Feature Vector)



Maps (In, Out, Wgt)
(P ₀ , Q ₁ , W _{-1,-1})
(P ₃ , Q ₄ , W _{-1,-1})
(P ₁ , Q ₃ , W _{-1,0})
(P ₀ , Q ₀ , W _{0,0})
(P ₁ , Q ₁ , W _{0,0})
(P ₂ , Q ₂ , W _{0,0})
(P ₃ , Q ₃ , W _{0,0})
(P ₄ , Q ₄ , W _{0,0})
(P ₃ , Q ₁ , W _{1,0})
(P ₁ , Q ₀ , W _{1,1})
(P ₄ , Q ₃ , W _{1,1})

Gather
By
Weight

X



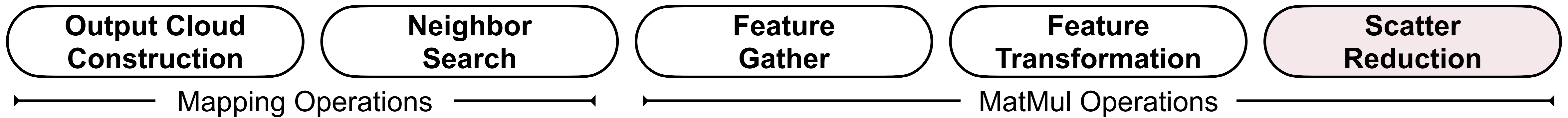
p -by- ic ic -by- oc p -by- oc × n

Matrix-Matrix Multiplication

Output Point Cloud

Q ₀
Q ₁
Q ₂
Q ₃
Q ₄

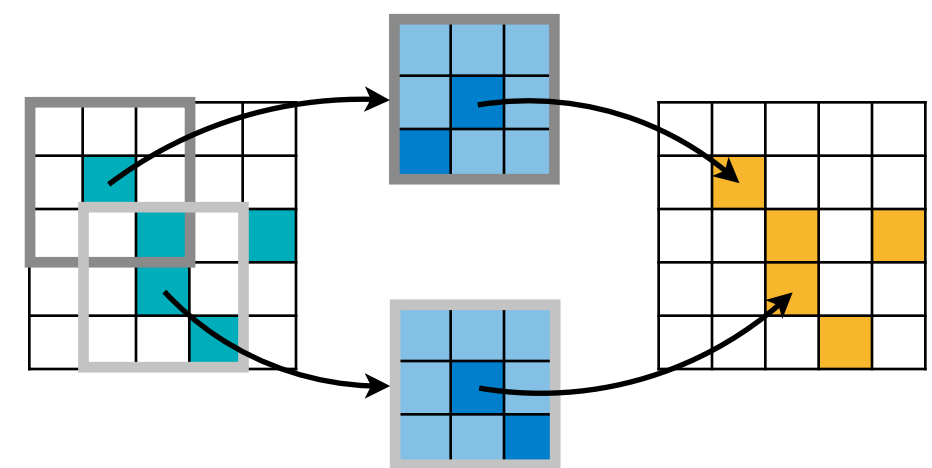
Step 5: Scatter and Reduce Partial Sums



Input Point Cloud

(P_0, F_0)
(P_1, F_1)
(P_2, F_2)
(P_3, F_3)
(P_4, F_4)

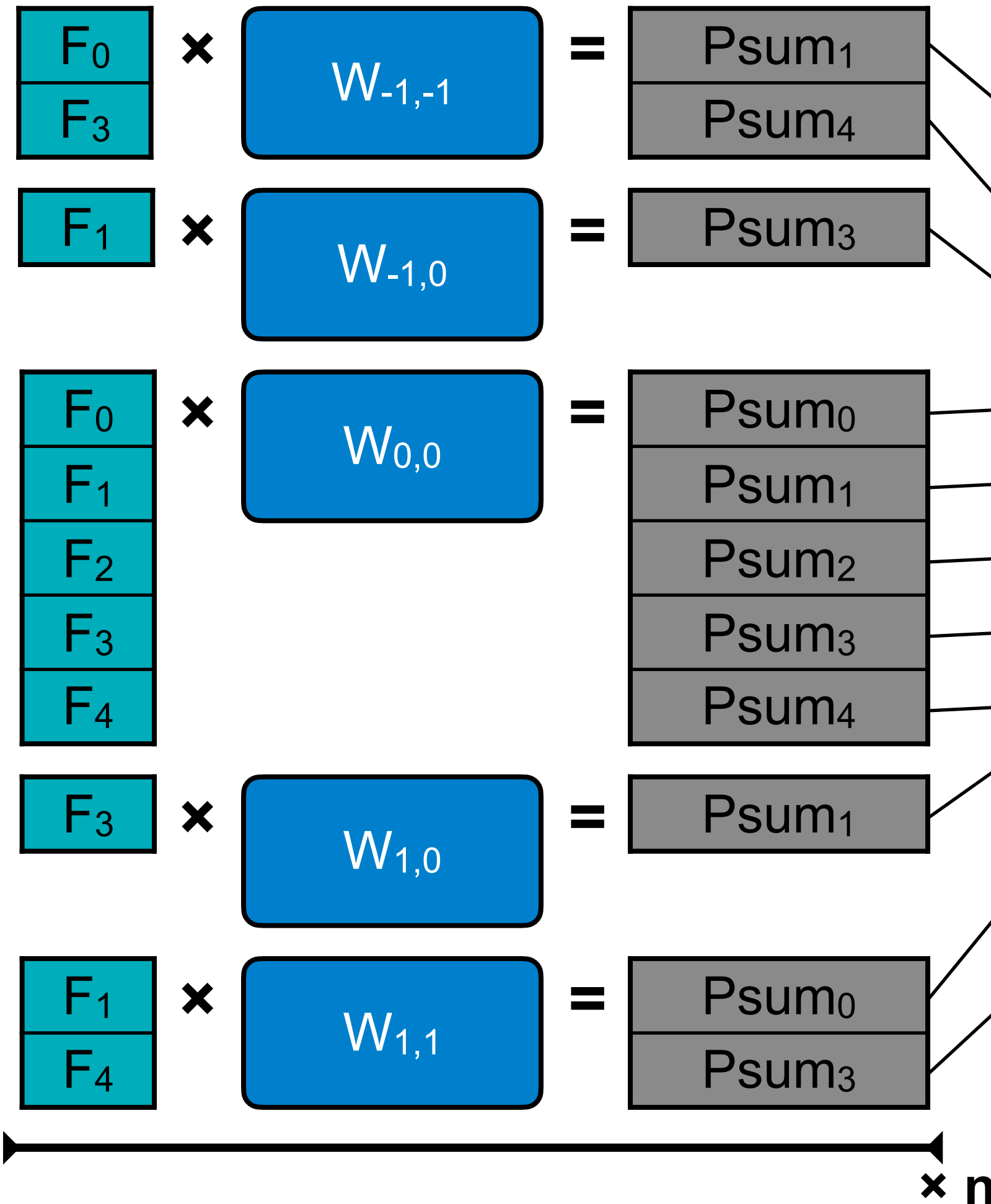
(Coords, Feature Vector)



Maps (In, Out, Wgt)		
$(P_0, Q_1, W_{-1,-1})$		
$(P_3, Q_4, W_{-1,-1})$		
$(P_1, Q_3, W_{-1,0})$		
$(P_0, Q_0, W_{0,0})$		
$(P_1, Q_1, W_{0,0})$		
$(P_2, Q_2, W_{0,0})$		
$(P_3, Q_3, W_{0,0})$		
$(P_4, Q_4, W_{0,0})$		
$(P_3, Q_1, W_{1,0})$		
$(P_1, Q_0, W_{1,1})$		
$(P_4, Q_3, W_{1,1})$		

Gather
By
Weight

X



Scatter
By
Output

X

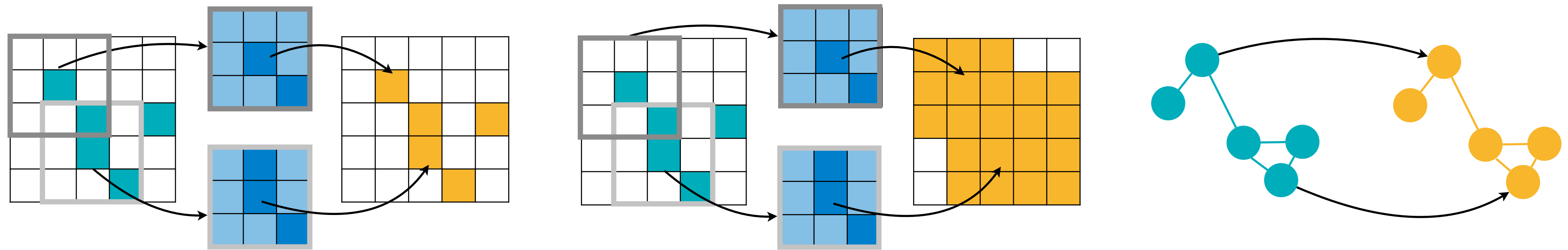
Reduction

Output Point Cloud

(Q_0, F_0)
(Q_1, F_1)
(Q_2, F_2)
(Q_3, F_3)
(Q_4, F_4)

MaxPool,
Accumulation, ...

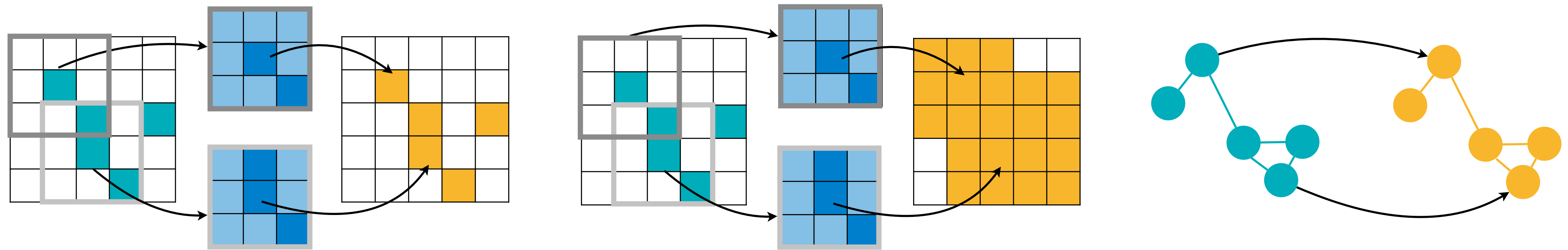
Point Cloud Deep Learning is Different



	Point Cloud NN	CNN	Graph CNN
source of sparsity	the distribution of points in physical space	ReLU and weight pruning	-
input sparsity	nonzeros do not dilate; input sparsity keeps low	nonzeros dilate; input sparsity quickly reduces	-

► Extreme low utilization on existing sparse CNN accelerators

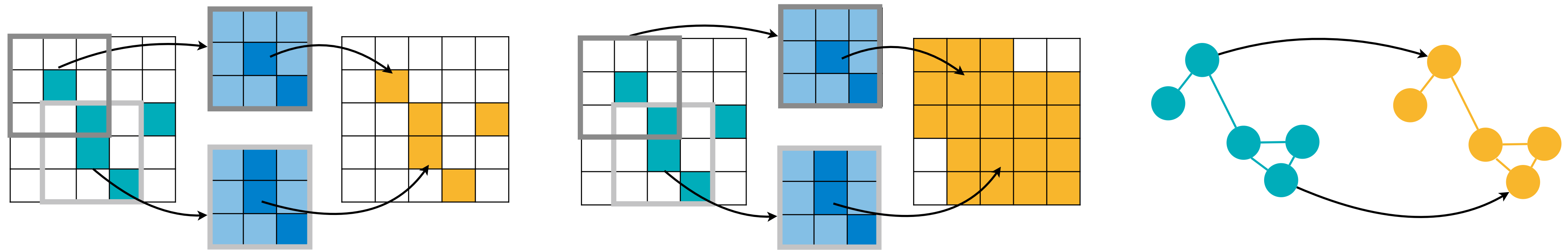
Point Cloud Deep Learning is Different



	Point Cloud NN	CNN	Graph CNN
source of sparsity	the distribution of points in physical space	ReLU and weight pruning	-
input sparsity	nonzeros do not dilate; input sparsity keeps low	nonzeros dilate; input sparsity quickly reduces	-
receptive field	both outputs' coords and neighbors need to be explicitly calculated	outputs' and neighbors' coords are inferred by pointer arithmetic	neighbors are given from the adjacency matrix

- ▶ Existing NN accelerators do not support mapping ops

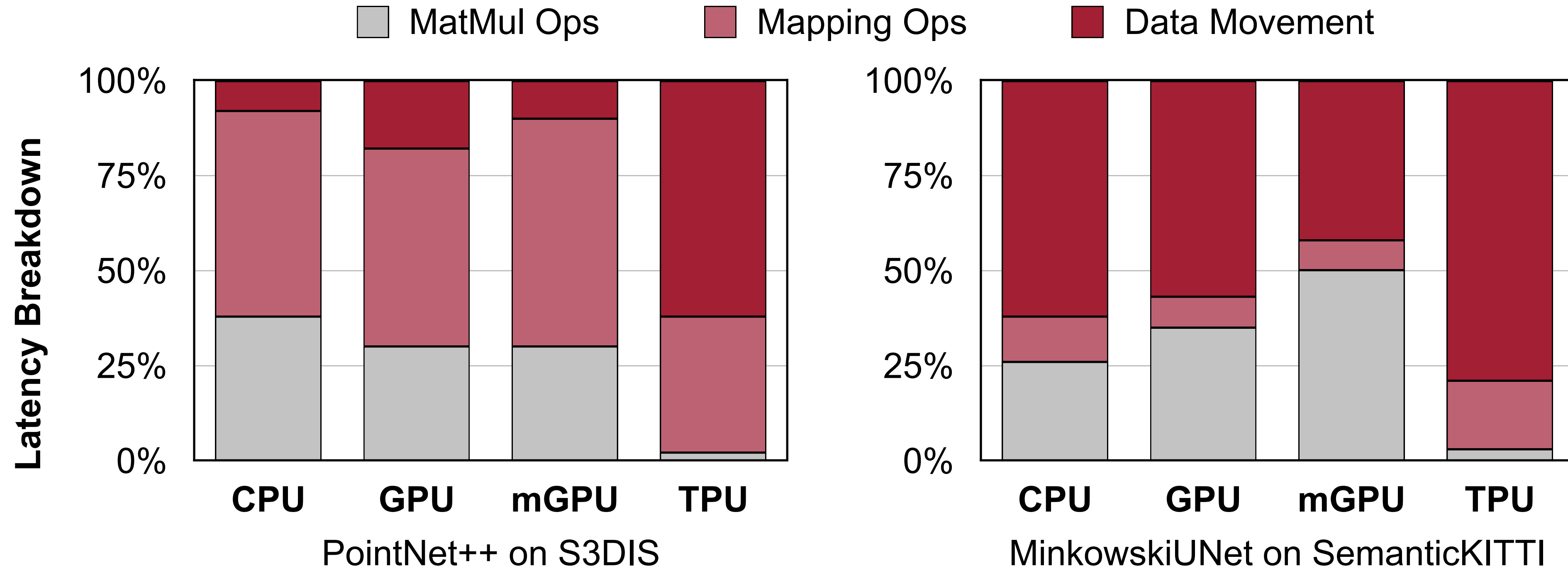
Point Cloud Deep Learning is Different



	Point Cloud NN	CNN	Graph CNN
source of sparsity	the distribution of points in physical space	ReLU and weight pruning	-
input sparsity	nonzeros do not dilate; input sparsity keeps low	nonzeros dilate; input sparsity quickly reduces	-
receptive field	both outputs' coords and neighbors need to be explicitly calculated	outputs' and neighbors' coords are inferred by pointer arithmetic	neighbors are given from the adjacency matrix
weights	weights can be shared or different for neighbors	weight are different for neighbors	weights are shared among neighbors
data movement	require both gather and scatter	no explicit gather or scatter	require either gather or scatter

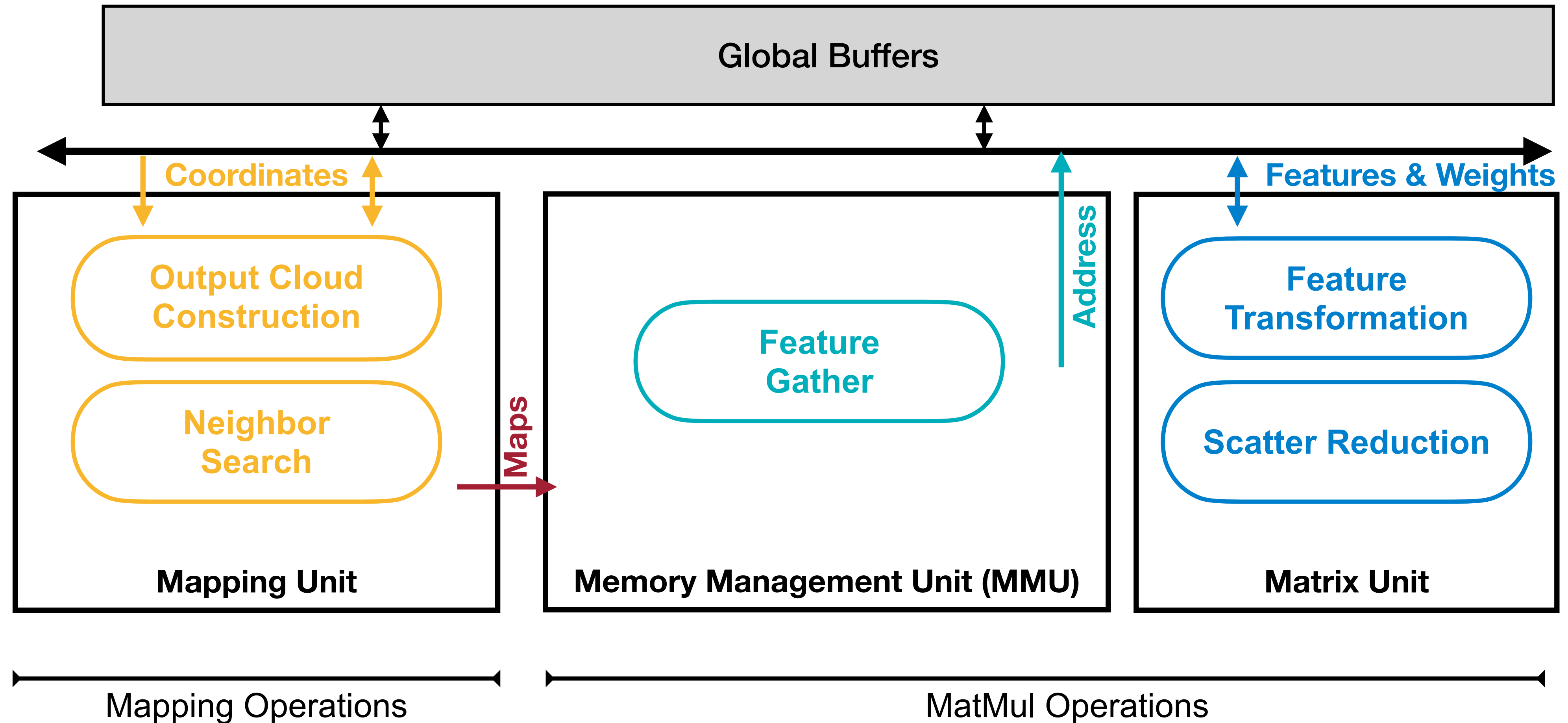
- ▶ GCN accelerators and previous PointNet accelerator Mesorasi do not support

Challenge of Point Cloud Deep Learning

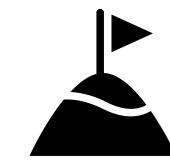
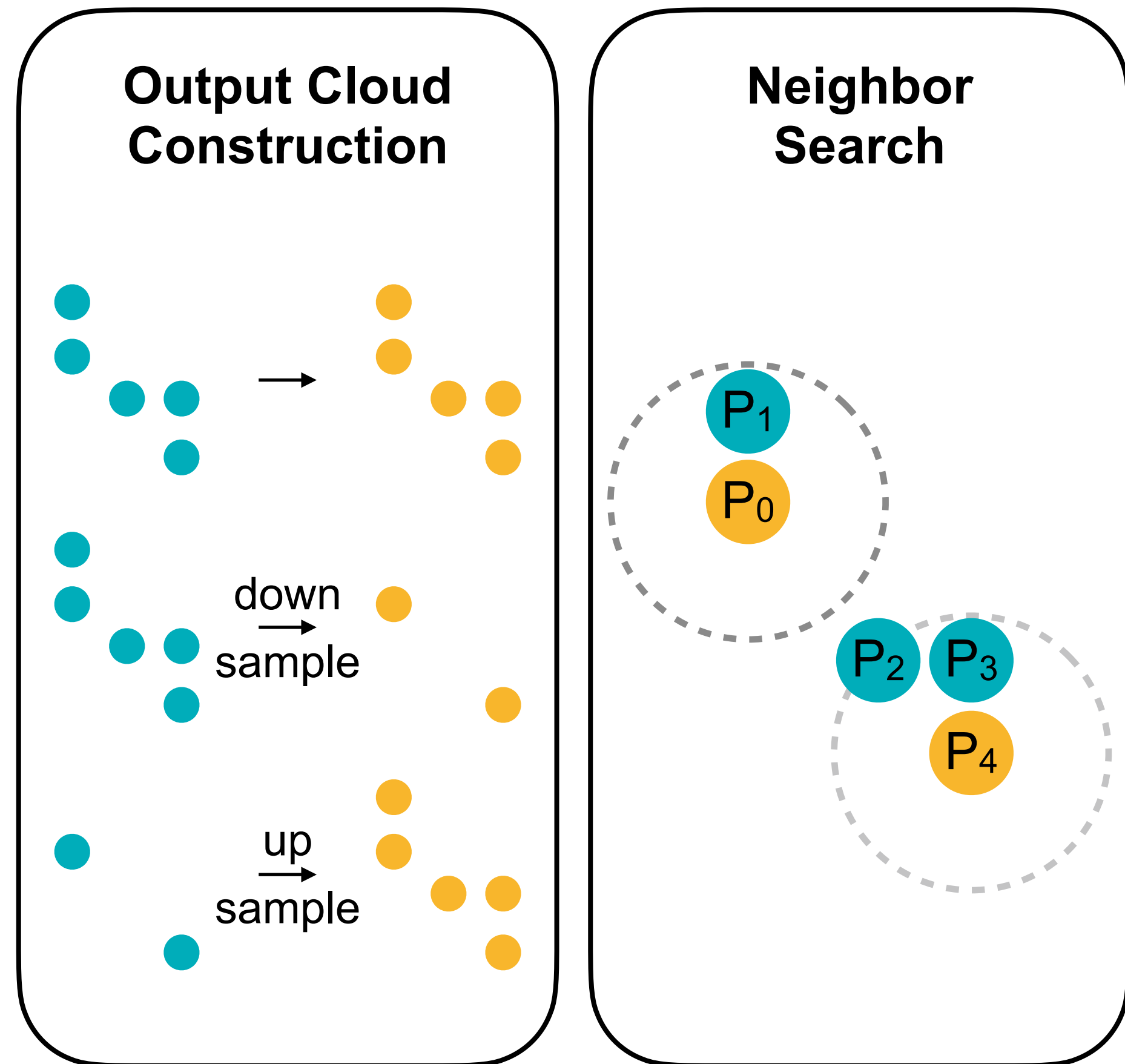


- Due to extreme sparsity, mapping operations and data movement overhead together take up **>50%** of the total runtime latency
- Due to unsupported mapping ops, data movement between co-processors (CPU and TPU) worsens the bottleneck

PointAcc: Efficient Point Cloud Accelerator



Mapping Unit: Diverse Mapping Ops in One Versatile Arch



Goal of Mapping Ops:

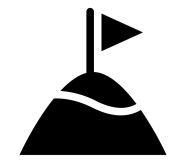
Generate Map (input point, output point, weight index)



Key Observation:

Maps are constructed based on the **comparison** among distances

Mapping Unit: Diverse Mapping Ops in One Versatile Arch



Goal of Mapping Ops: Generate Map (input point, output point, weight index)

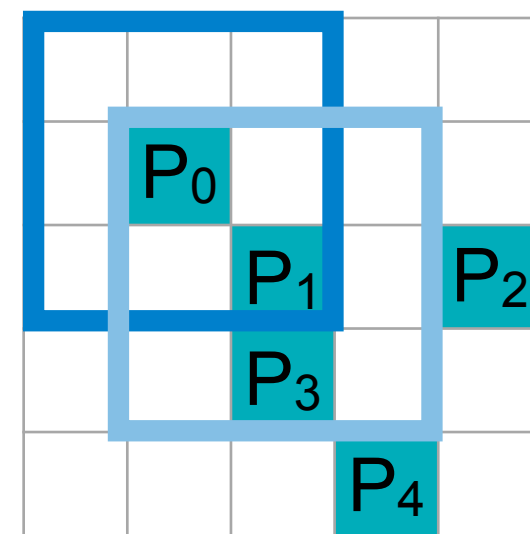


Key Observation: Maps are constructed based on the *comparison* among distances

Kernel Mapping

- ▶ comparison op: not greater than, $|\Delta x| \leq 1, |\Delta y| \leq 1, |\Delta z| \leq 1$
- ▶ query the hash table of input point cloud for each output point

```
For Q in O={Q0, Q1, Q2, ...}:  
  Find P in I, s.t. |Px-Qx| <= 1 && |Py-Qy| <= 1 && |Pz-Qz| <= 1
```

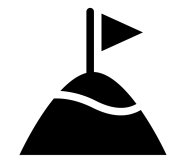


```
(In, Out, Wgt)  
(P0, Q0, W0,0)  
(P1, Q0, W1,1)  
(P0, Q1, W-1,-1)  
(P1, Q1, W0,0)  
...
```

```
For W in {W-1,-1, W-1,0, W1,0, ...}:  
  For Q in O={Q0, Q1, Q2, ...}:  
    Query Q+δw in HashTable(I)
```

- require on-chip memory as large as 160MB
- cannot be parallelized efficiently

Mapping Unit: Diverse Mapping Ops in One Versatile Arch



Goal of Mapping Ops: Generate Map (input point, output point, weight index)

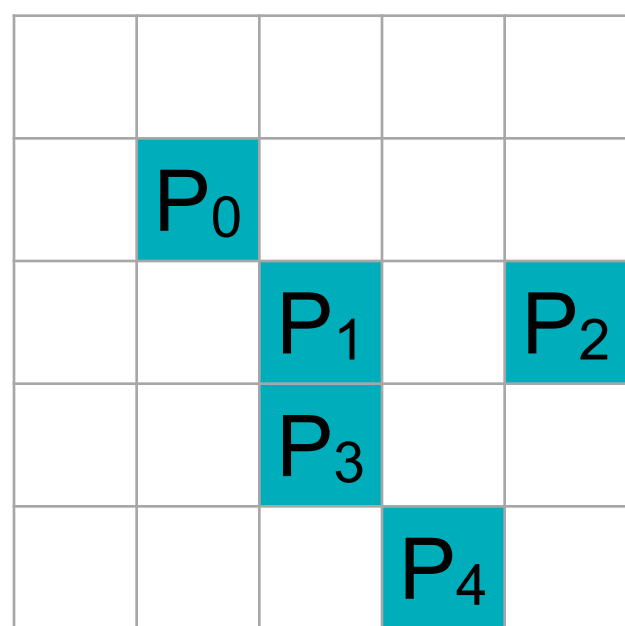


Key Observation: Maps are constructed based on the *comparison* among distances

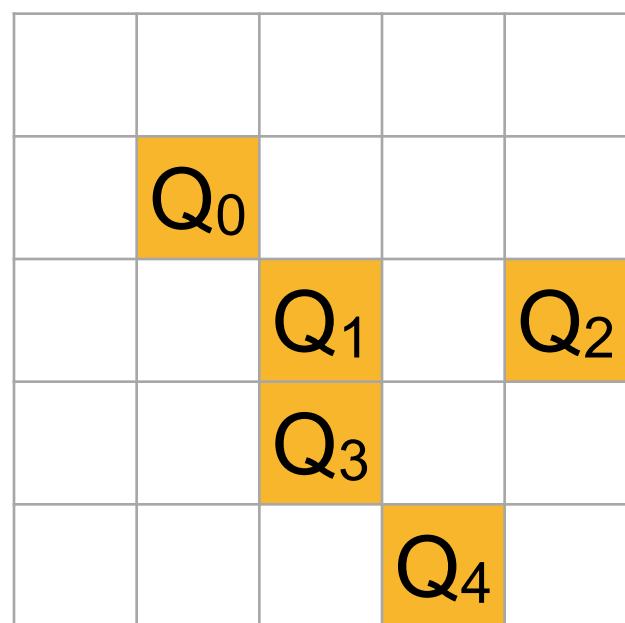
Kernel Mapping

- ▶ comparison op: not greater than, $|\Delta x| \leq 1, |\Delta y| \leq 1, |\Delta z| \leq 1$
- ▶ ~~query the hash table of input point cloud for each output point~~
- ▶ coordinates intersection for each neighbor position
- ▶ → parallelizable merge-sort and equal comparison

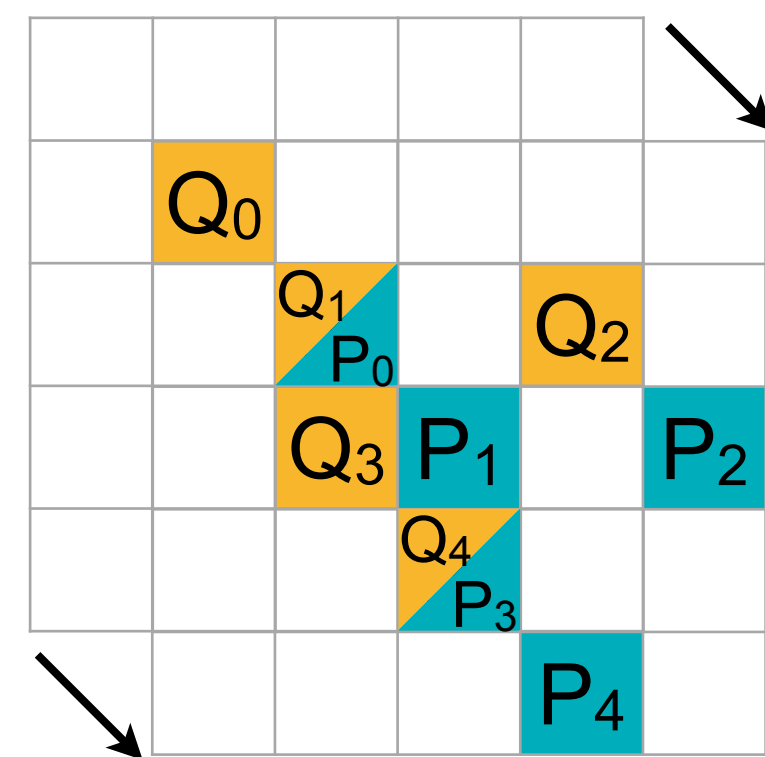
Input Point Cloud



↓ stride = 1



Output Point Cloud

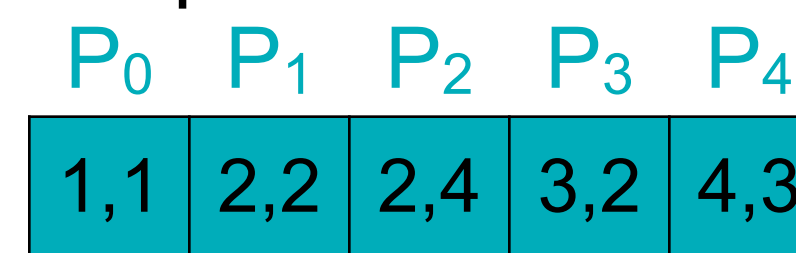


Shift Input for $W_{-1,-1}$

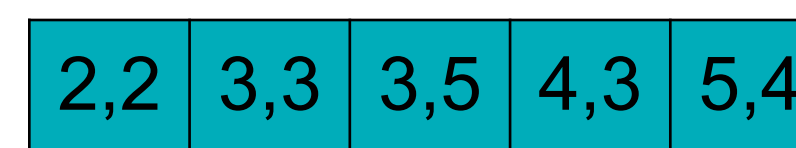
$W_{-1,-1}$	$W_{-1,0}$	$W_{-1,1}$
$W_{0,-1}$	$W_{0,0}$	$W_{0,1}$
$W_{1,-1}$	$W_{1,0}$	$W_{1,1}$

(In, Out, Wgt)
 (P₀, Q₁, W_{-1,-1})
 (P₃, Q₄, W_{-1,-1})

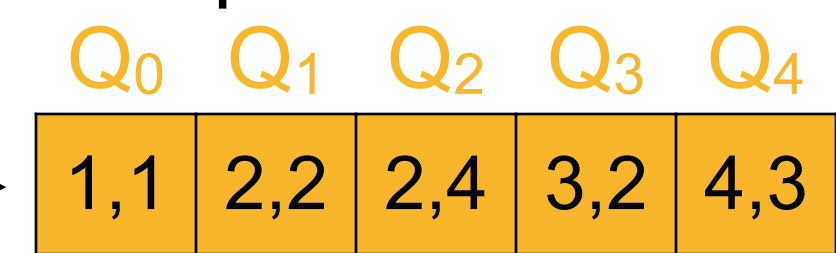
Input Point Cloud



+ (1, 1) ↓ for $w_{-1,-1}$



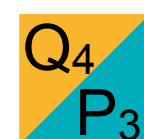
Output Point Cloud



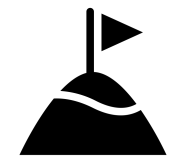
Merge Sort



Intersection



Mapping Unit: Diverse Mapping Ops in One Versatile Arch



Goal of Mapping Ops: Generate Map (input point, output point, weight index)

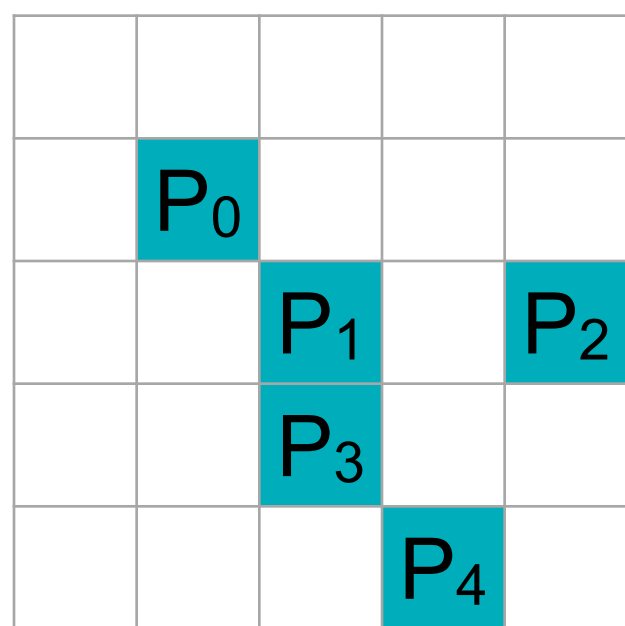


Key Observation: Maps are constructed based on the *comparison* among distances

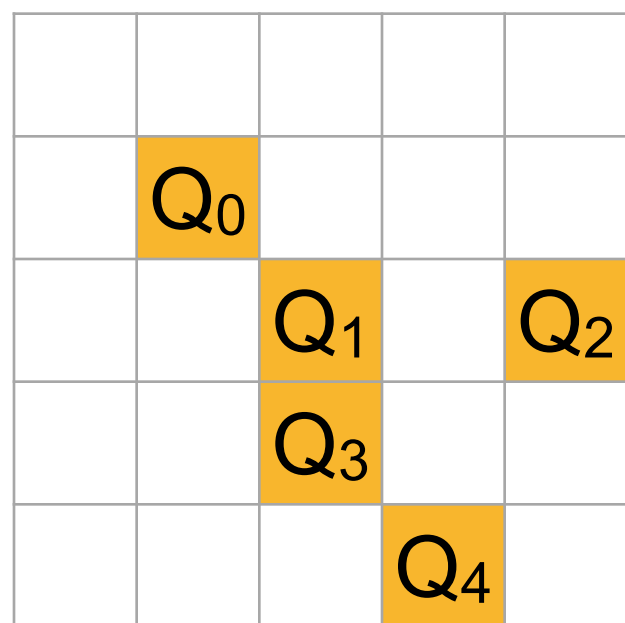
Kernel Mapping

- ▶ comparison op: not greater than, $|\Delta x| \leq 1, |\Delta y| \leq 1, |\Delta z| \leq 1$
- ▶ ~~query the hash table of input point cloud for each output point~~
- ▶ coordinates intersection for each neighbor position
- ▶ → parallelizable merge-sort and equal comparison

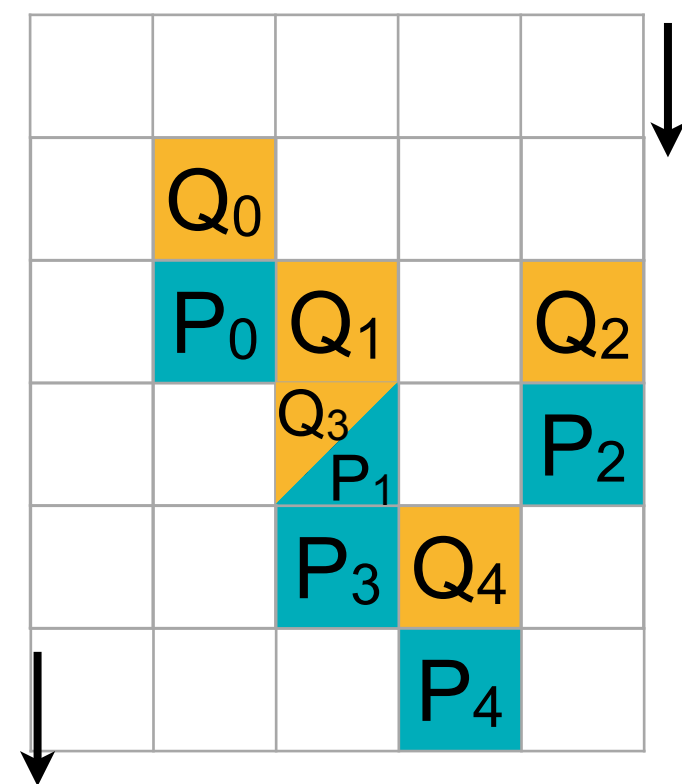
Input Point Cloud



↓ stride = 1



Output Point Cloud

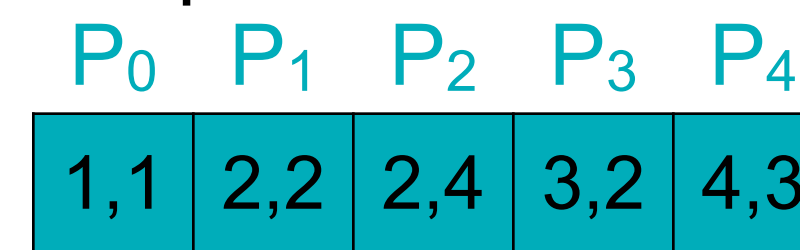


Shift Input for $W_{-1,0}$

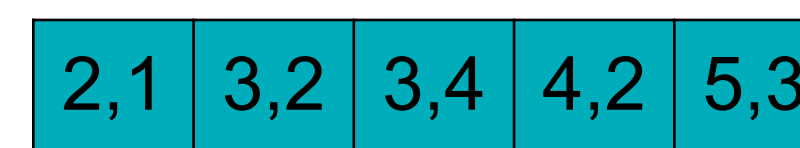
$W_{-1,-1}$	$W_{-1,0}$	$W_{-1,1}$
$W_{0,-1}$	$W_{0,0}$	$W_{0,1}$
$W_{1,-1}$	$W_{1,0}$	$W_{1,1}$

(In, Out, Wgt)
(P ₀ , Q ₁ , $W_{-1,-1}$)
(P ₃ , Q ₄ , $W_{-1,-1}$)
(P ₁ , Q ₃ , $W_{-1,0}$)

Input Point Cloud



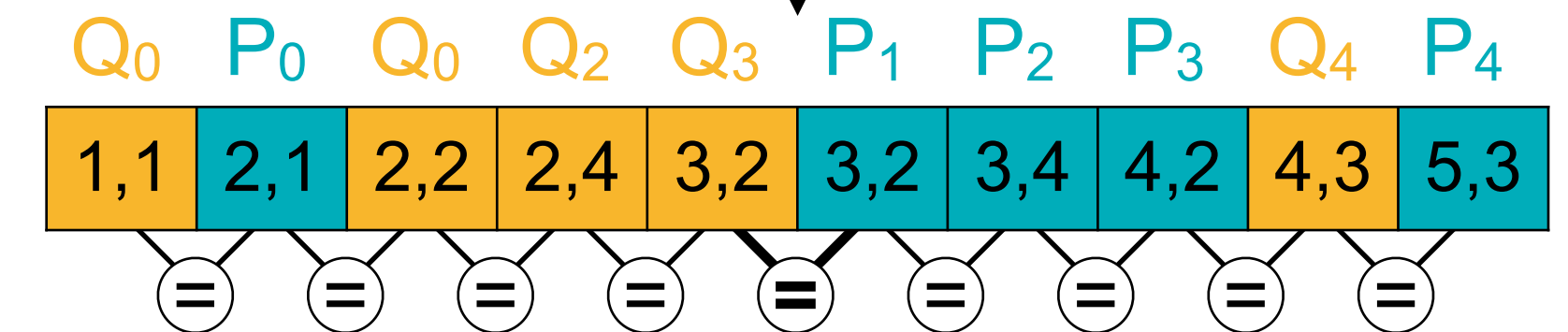
+ (1, 0) ↓ for $w_{-1,0}$



Output Point Cloud



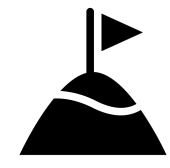
Merge Sort



= = = = = = = = = =



Mapping Unit: Diverse Mapping Ops in One Versatile Arch



Goal of Mapping Ops: Generate Map (input point, output point, weight index)

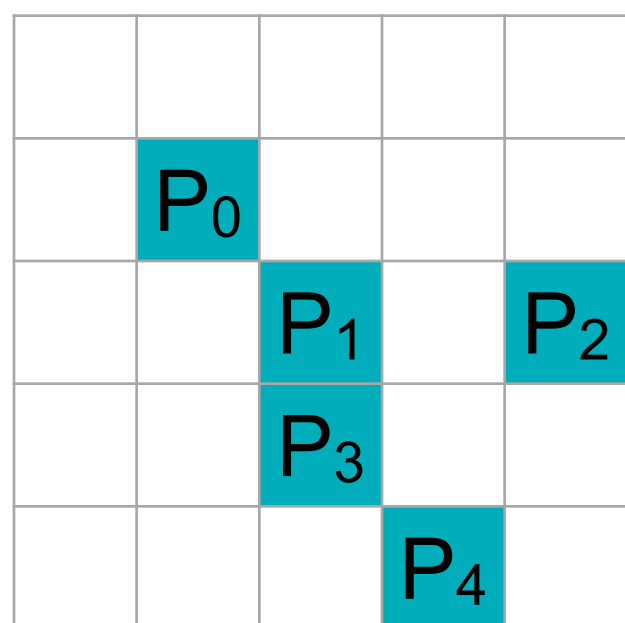


Key Observation: Maps are constructed based on the *comparison* among distances

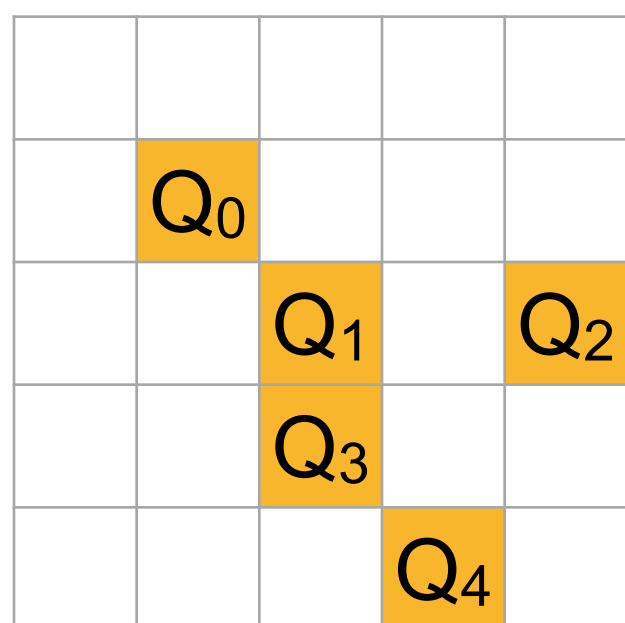
Kernel Mapping

- comparison op: not greater than, $|\Delta x| \leq 1, |\Delta y| \leq 1, |\Delta z| \leq 1$
- ~~query the hash table of input point cloud for each output point~~
- coordinates intersection for each neighbor position
- parallelizable merge-sort and equal comparison

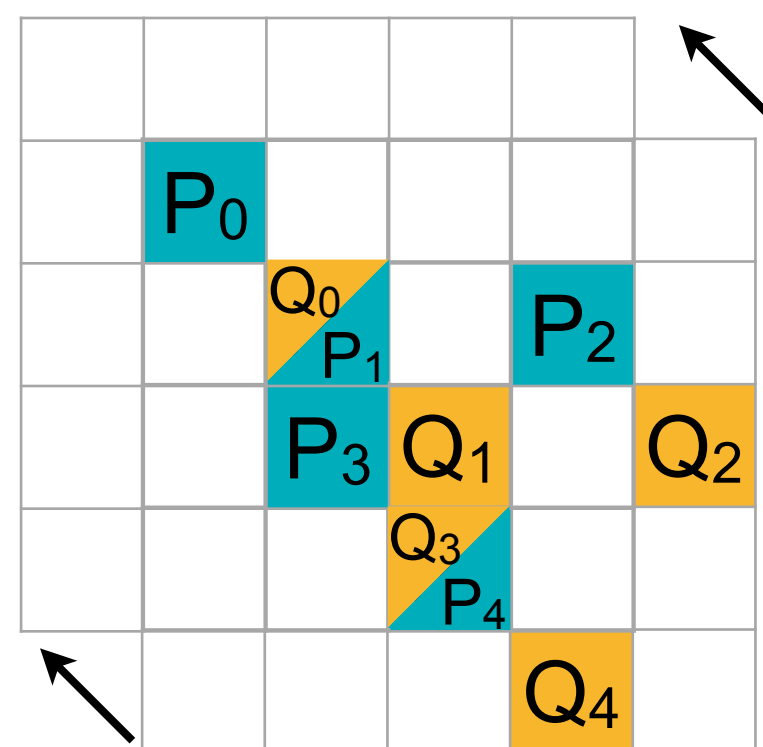
Input Point Cloud



↓ stride = 1



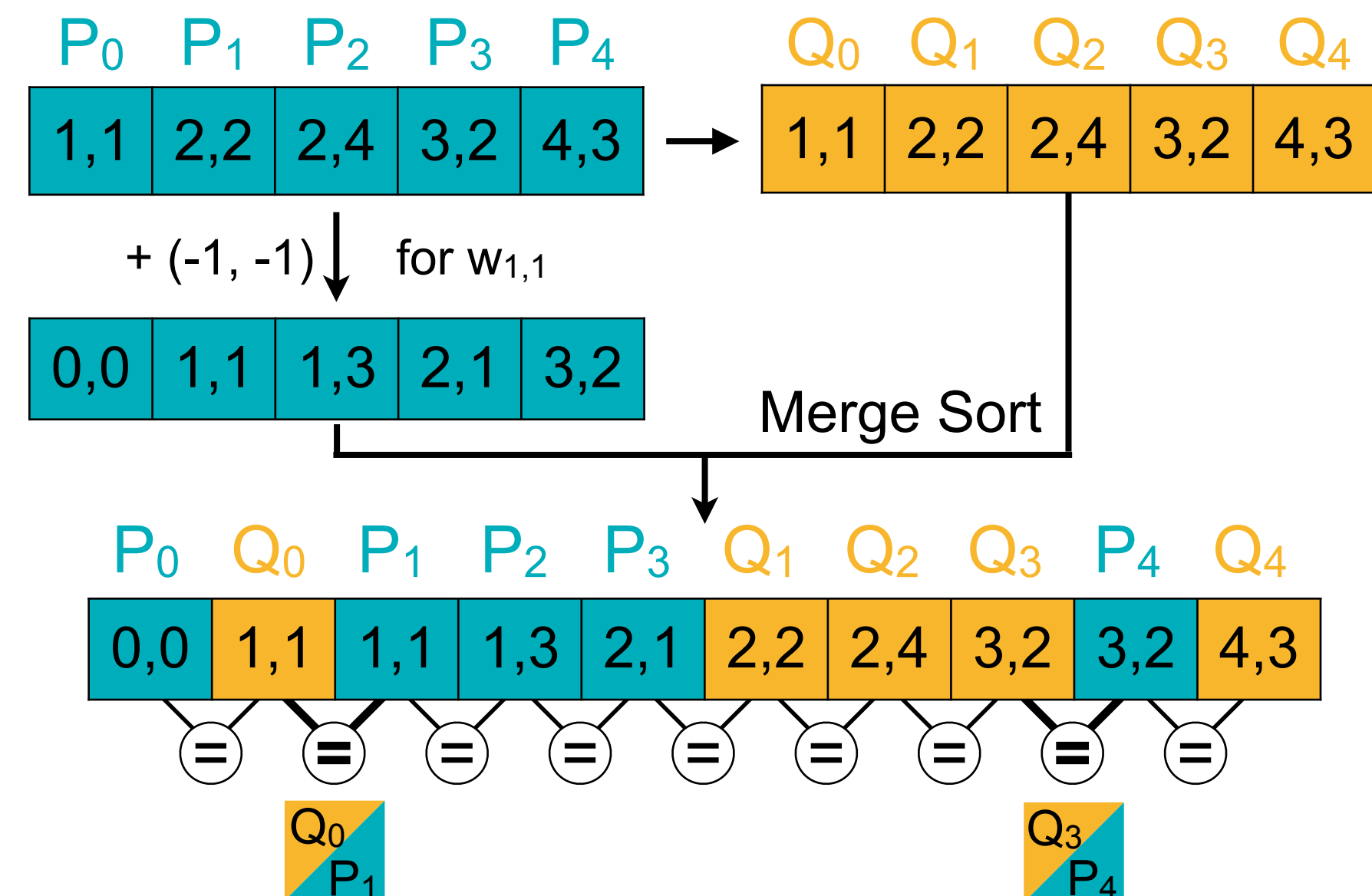
Output Point Cloud



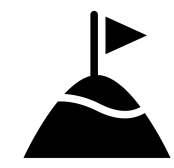
Shift Input for $W_{1,1}$

$W_{-1,-1}$	$W_{-1,0}$	$W_{-1,1}$
$W_{0,-1}$	$W_{0,0}$	$W_{0,1}$
$W_{1,-1}$	$W_{1,0}$	$W_{1,1}$

(In, Out, Wgt)
($P_0, Q_1, W_{-1,-1}$)
($P_3, Q_4, W_{-1,-1}$)
...
($P_1, Q_0, W_{1,1}$)
($P_4, Q_3, W_{1,1}$)



Mapping Unit: Diverse Mapping Ops in One Versatile Arch



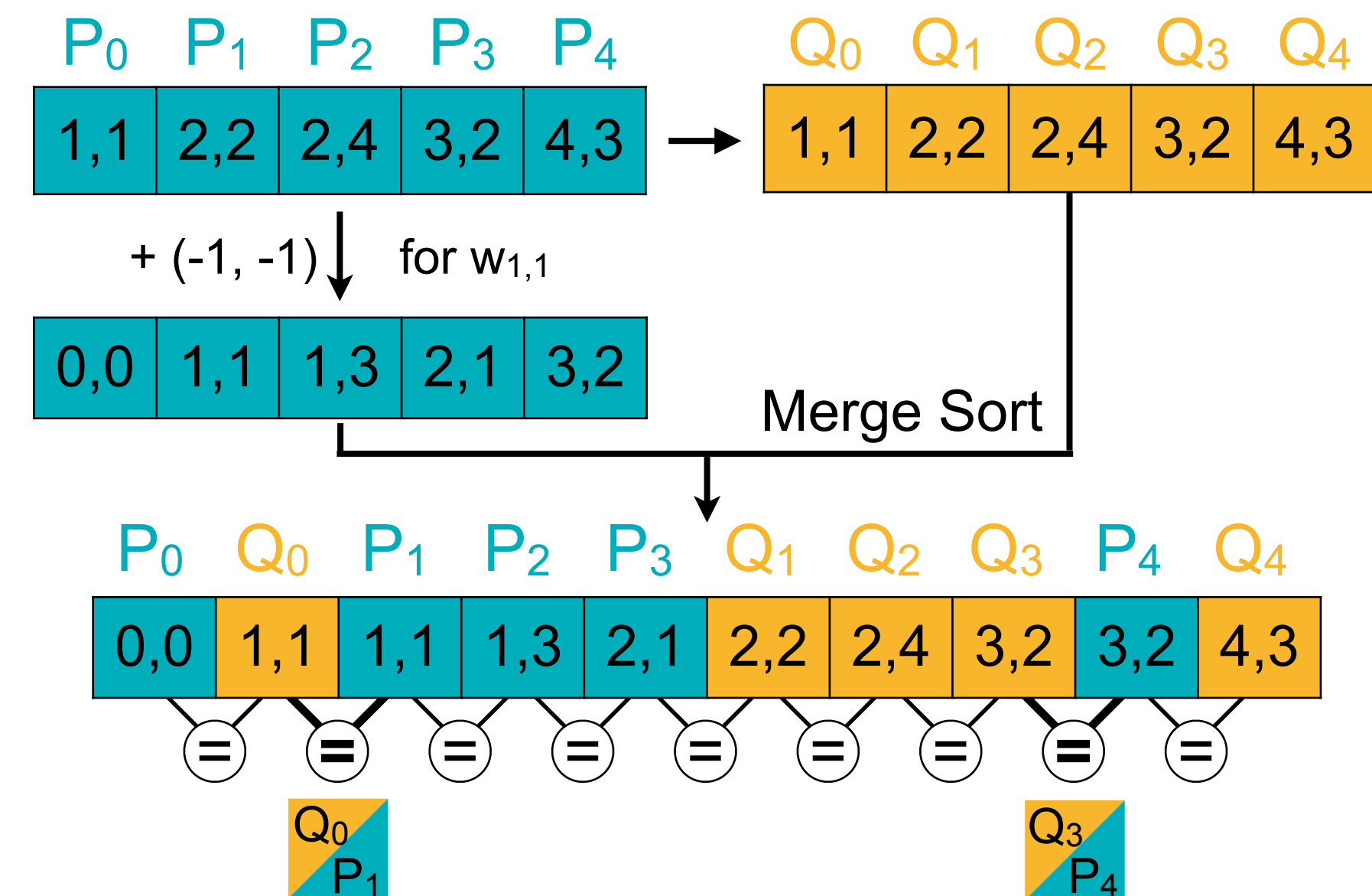
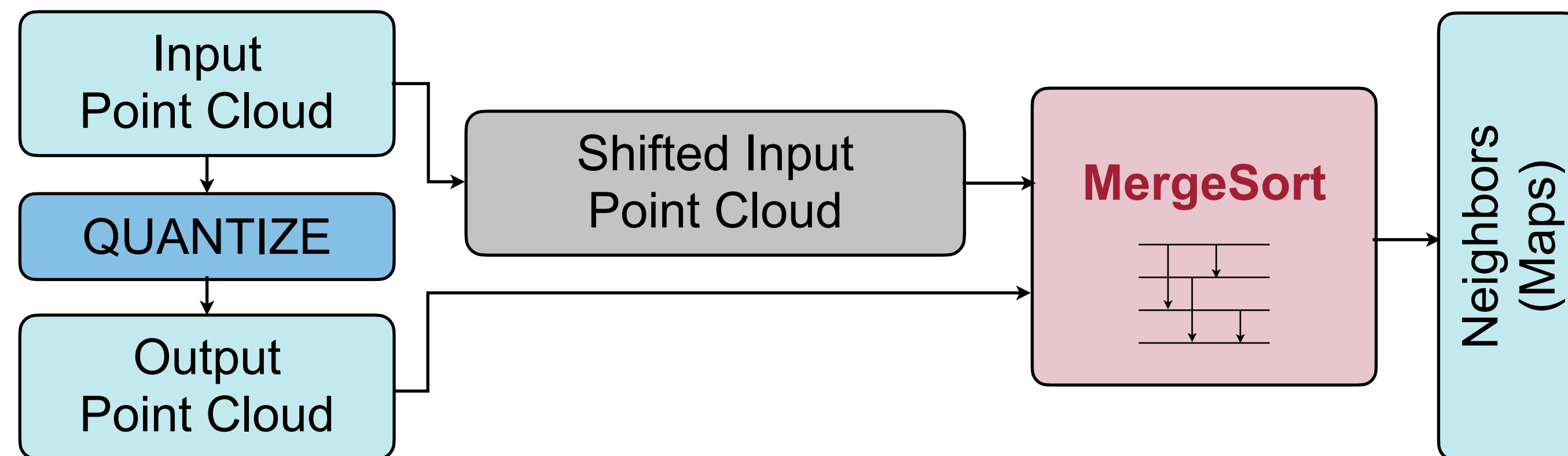
Goal of Mapping Ops: Generate Map (input point, output point, weight index)



Key Observation: Maps are constructed based on the *comparison* among distances

Kernel Mapping

- ▶ comparison op: not greater than, $|\Delta x| \leq 1, |\Delta y| \leq 1, |\Delta z| \leq 1$
- ▶ ~~query the hash table of input point cloud for each output point~~
- ▶ coordinates intersection for each neighbor position
- ▶ → parallelizable merge-sort and equal comparison



Mapping Unit: Diverse Mapping Ops in One Versatile Arch



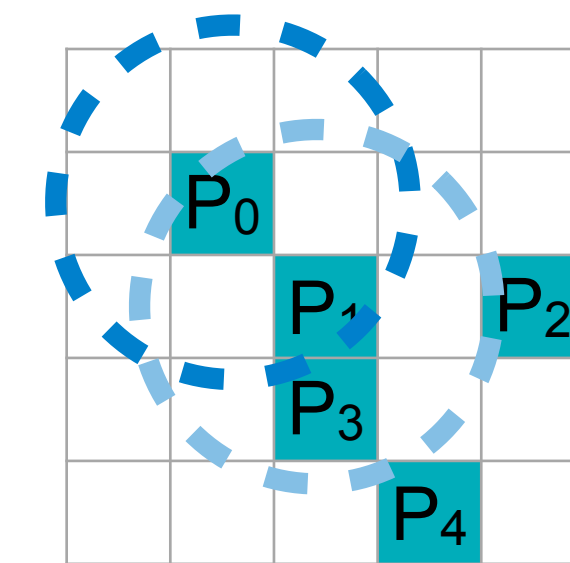
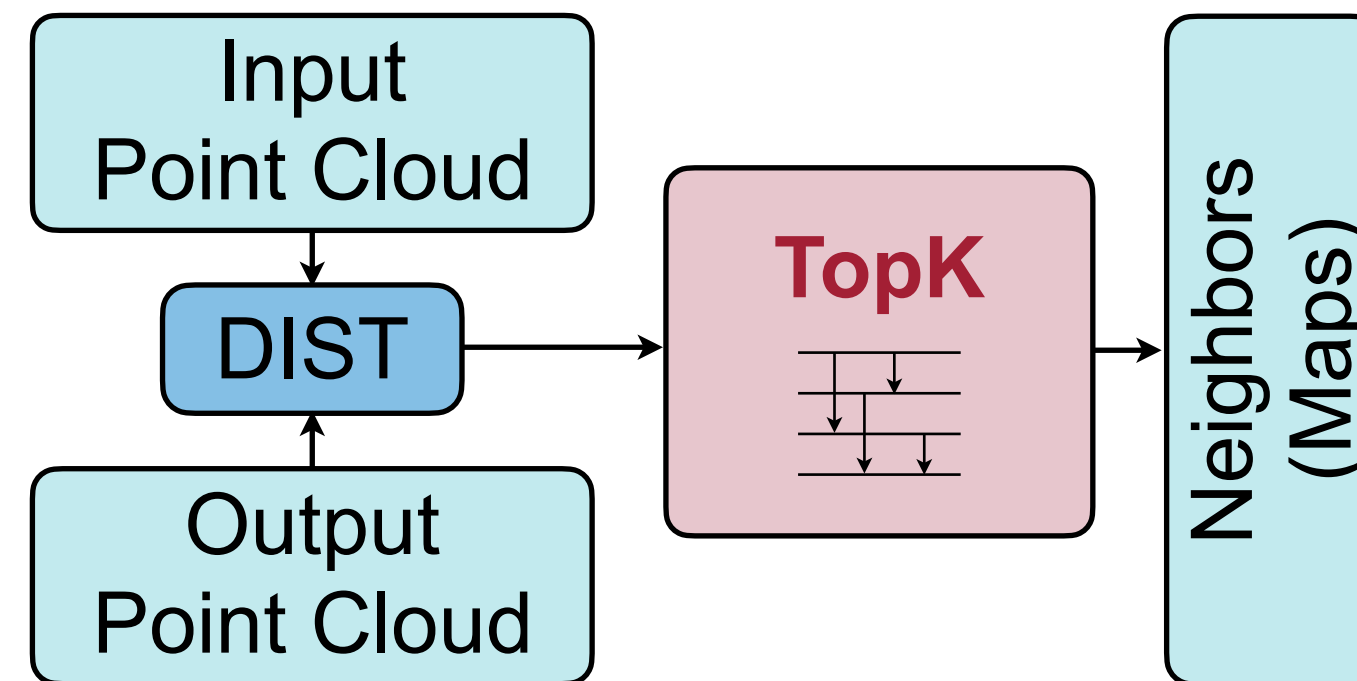
Goal of Mapping Ops: Generate Map (input point, output point, weight index)



Key Observation: Maps are constructed based on the *comparison* among distances

k-Nearest-Neighbor / Ball Query

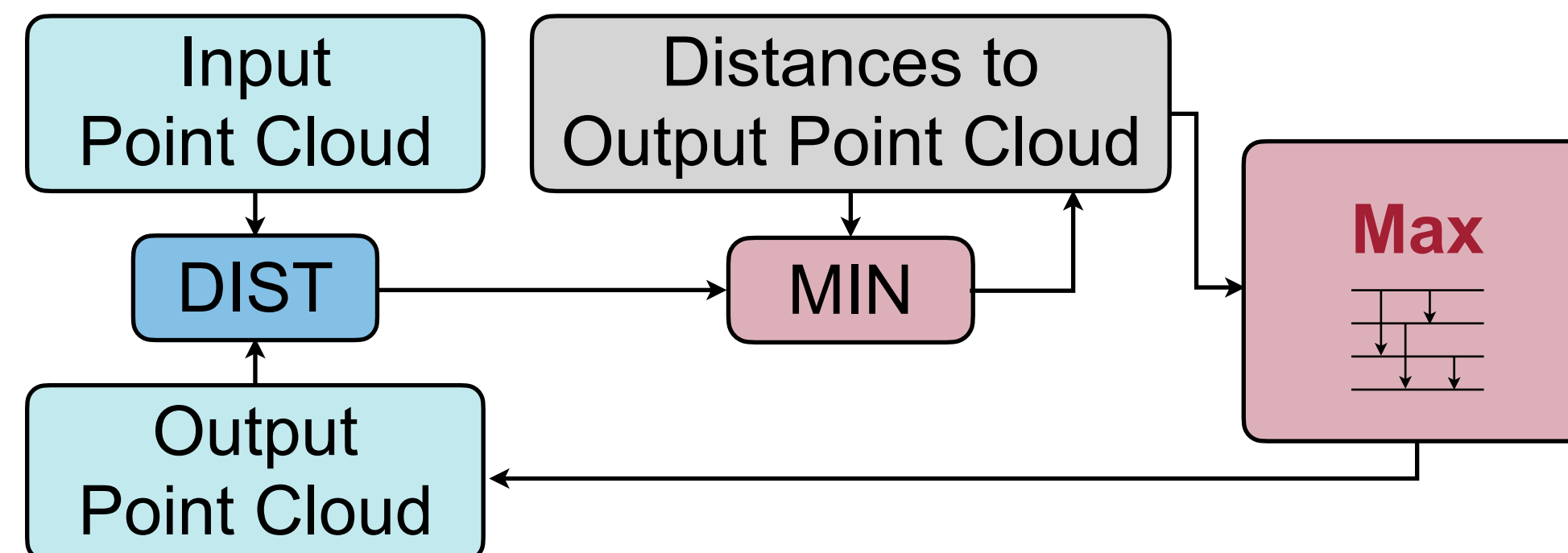
- comparison op: TopK
- ball query filters out the outsider in the nearest neighbors



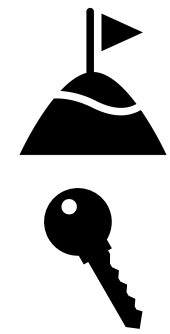
(In, Out, Wgt)		
P ₀	Q ₀	W ₀
P ₁	Q ₀	W ₁
P ₃	Q ₀	W ₂
P ₁	Q ₁	W ₀
P ₃	Q ₁	W ₁

Farthest Point Sampling

- comparison op: ArgMax / Max

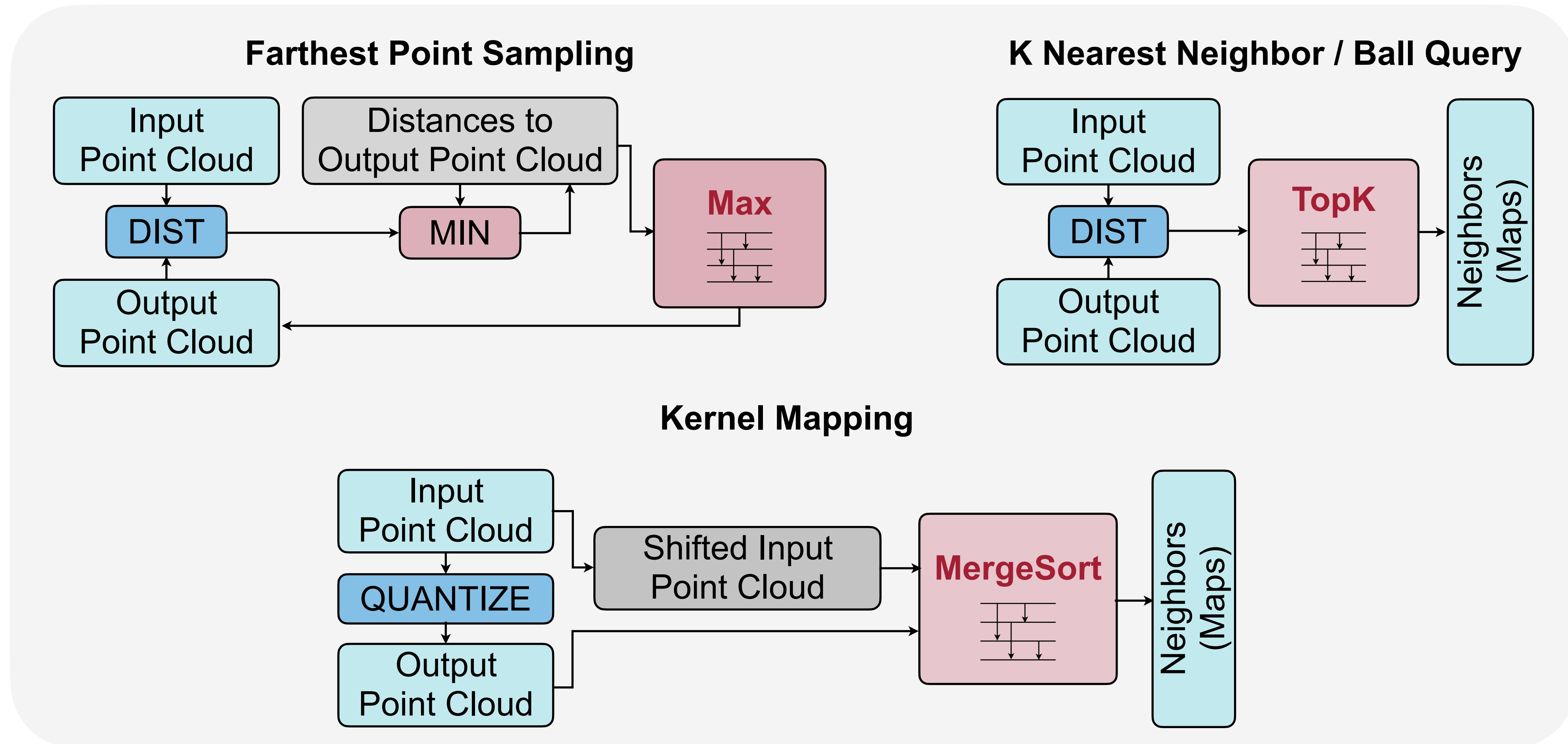


Mapping Unit: Diverse Mapping Ops in One Versatile Arch

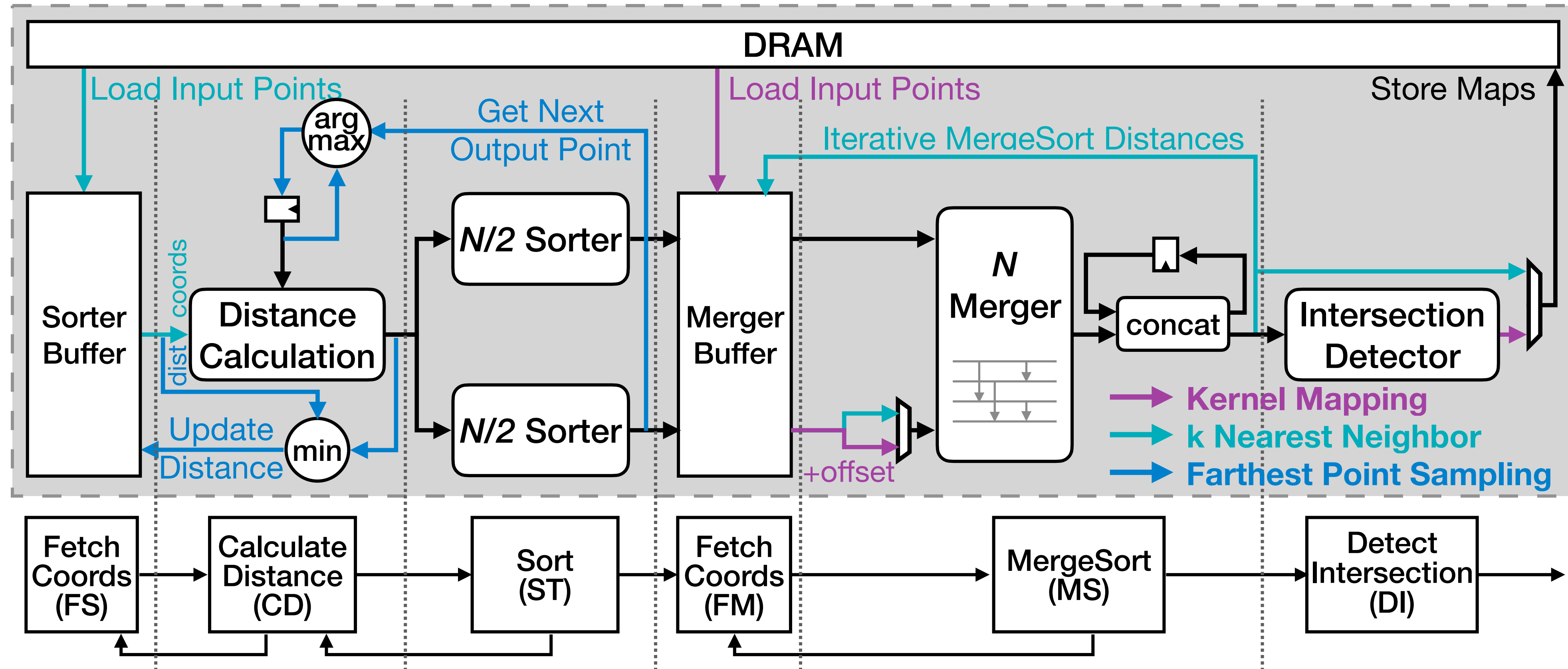


Goal of Mapping Ops: Generate Map (input point, output point, weight index)

Key Observation: Maps are constructed based on the *comparison* among distances

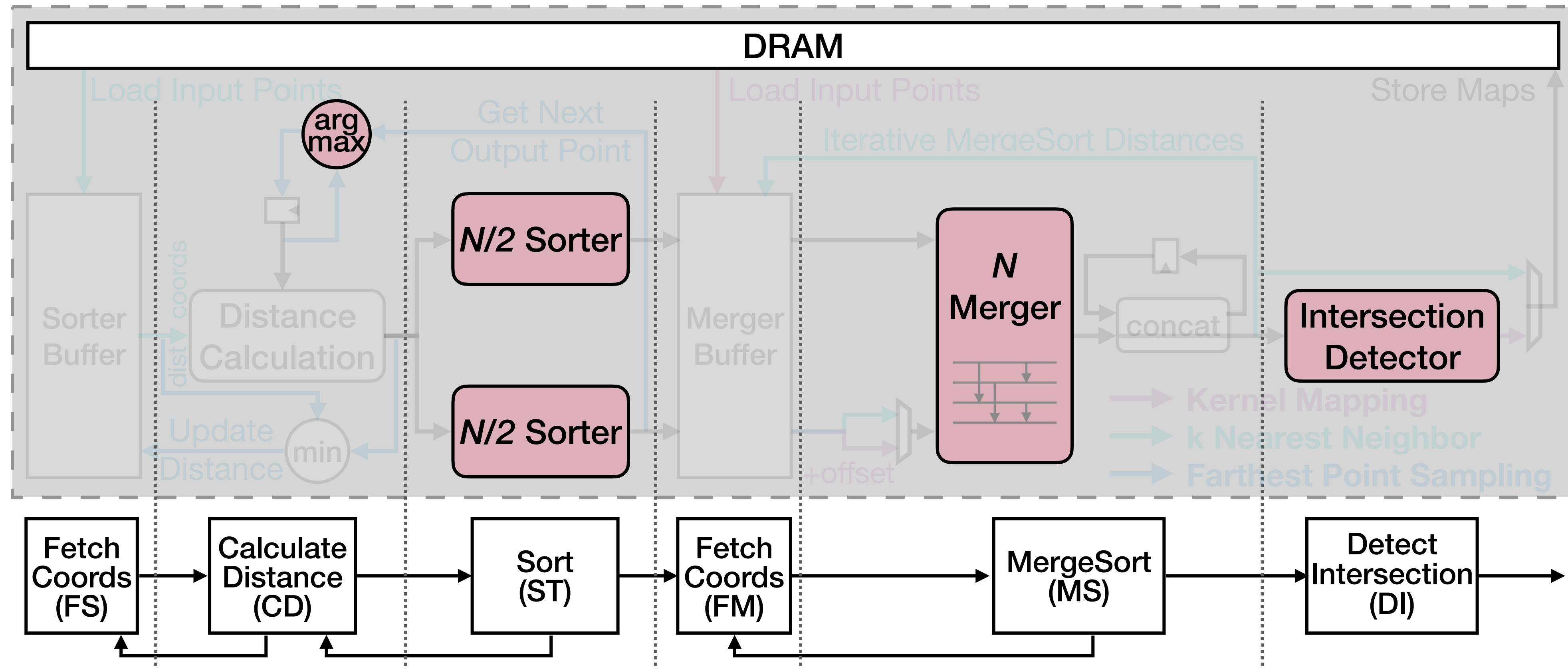


Mapping Unit: Diverse Mapping Ops in One Versatile Arch



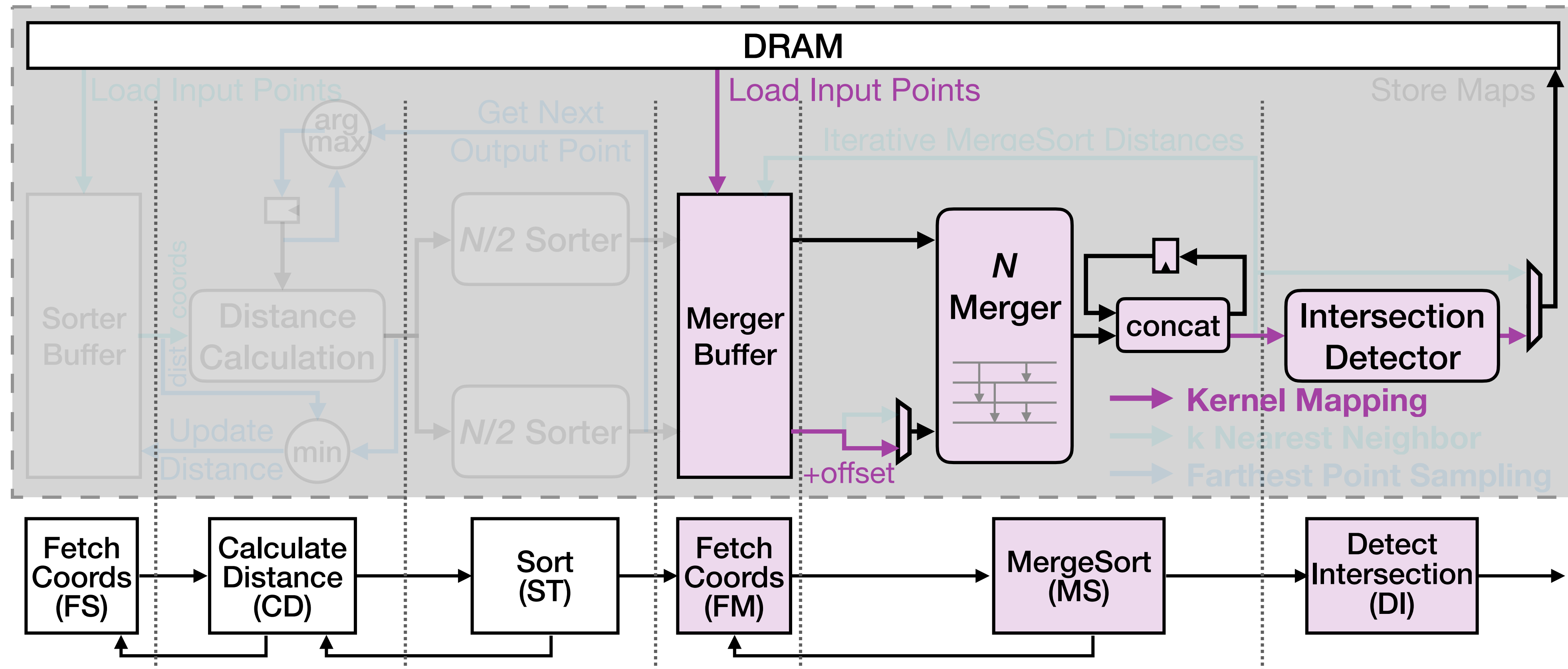
Mapping Unit: Diverse Mapping Ops in One Versatile Arch

- ▶ Main component for parallel comparison: sorters, merger



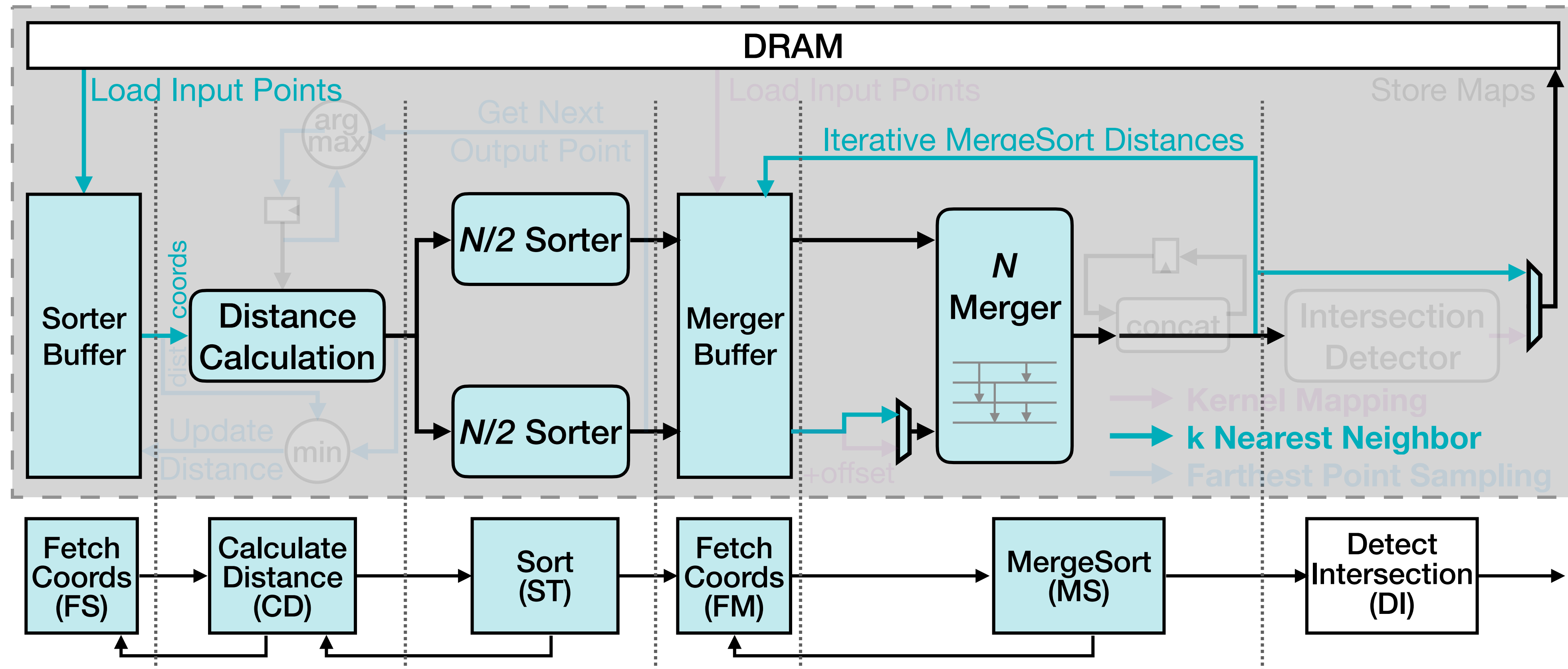
Mapping Unit: Diverse Mapping Ops in One Versatile Arch

- ▶ Data flow when running kernel mapping (*i.e.*, the neighborhood shape is cube)



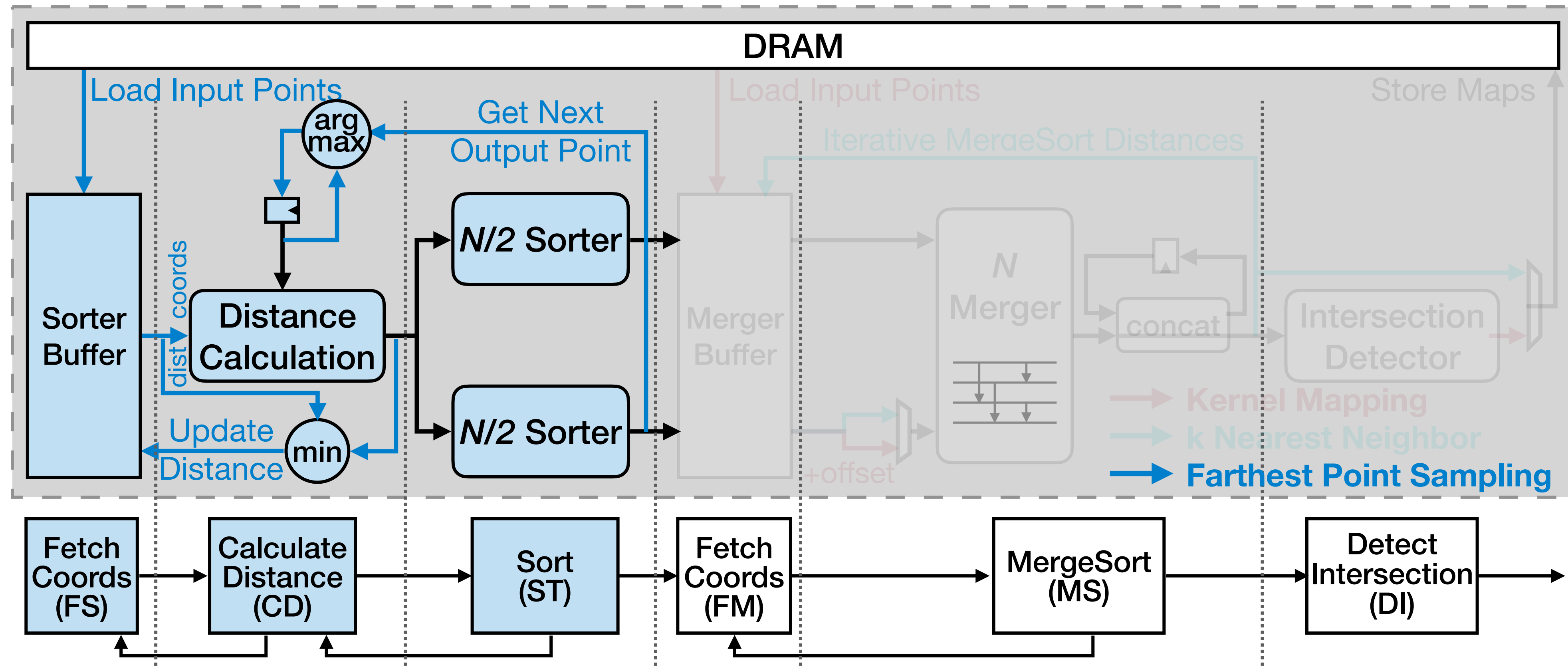
Mapping Unit: Diverse Mapping Ops in One Versatile Arch

- ▶ Data flow when running k-nearest-neighbor / ball query (i.e., the neighborhood shape is ball)



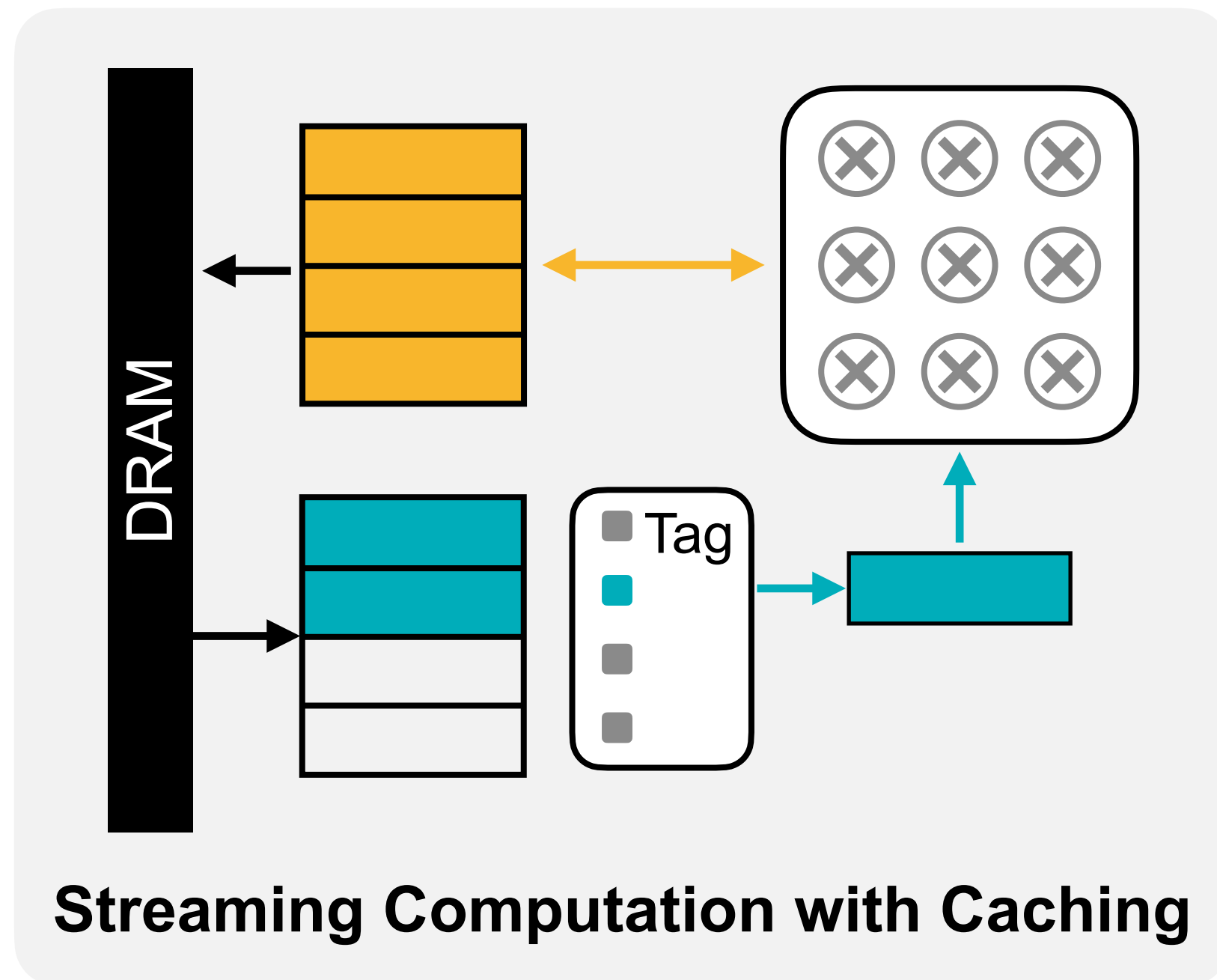
Mapping Unit: Diverse Mapping Ops in One Versatile Arch

- ▶ Data flow when running farthest point sampling

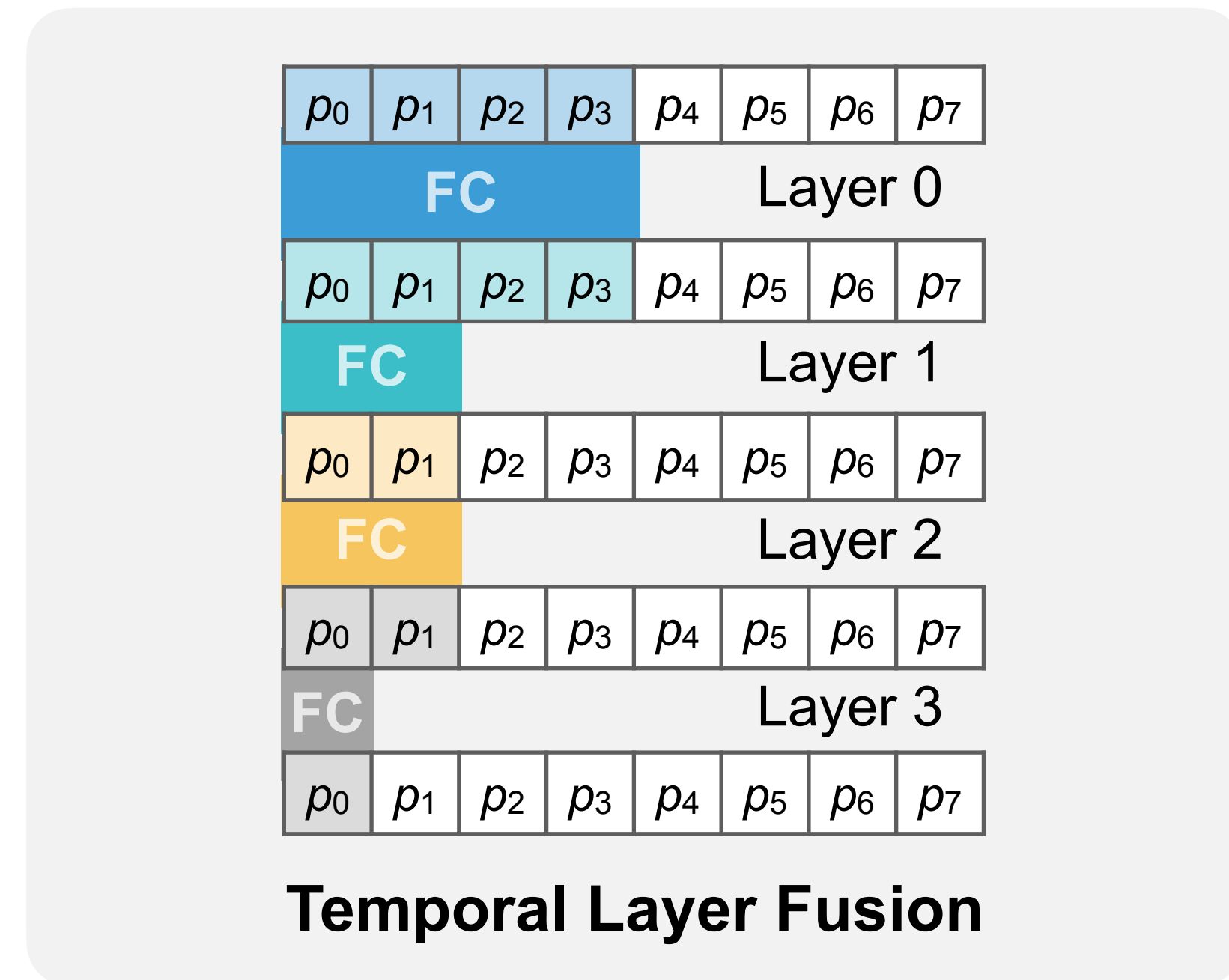


Flexible Memory Management

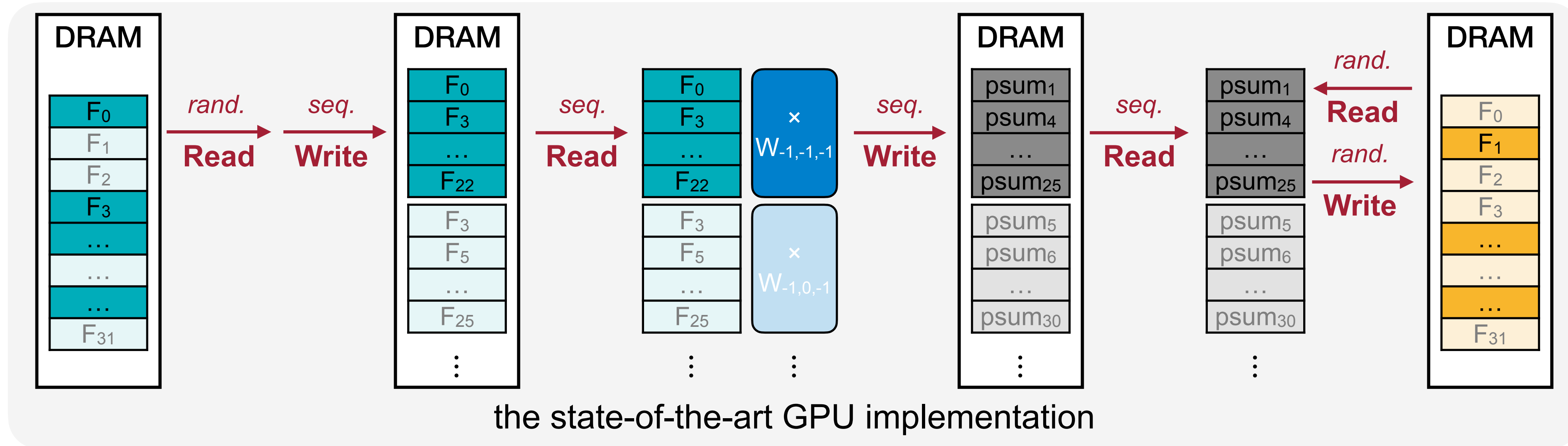
Sparse Computation



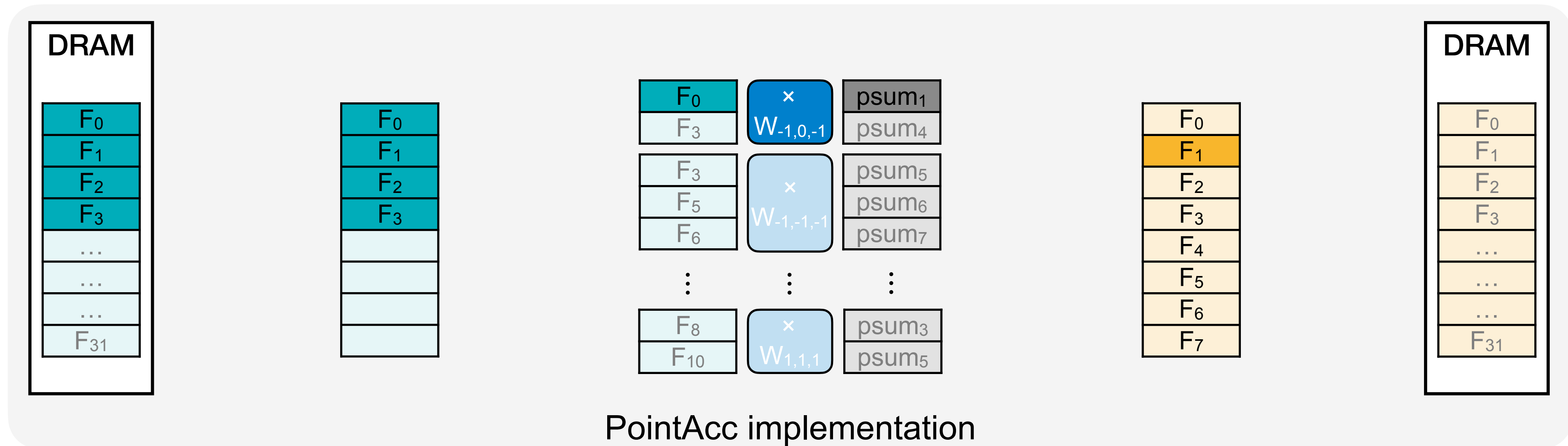
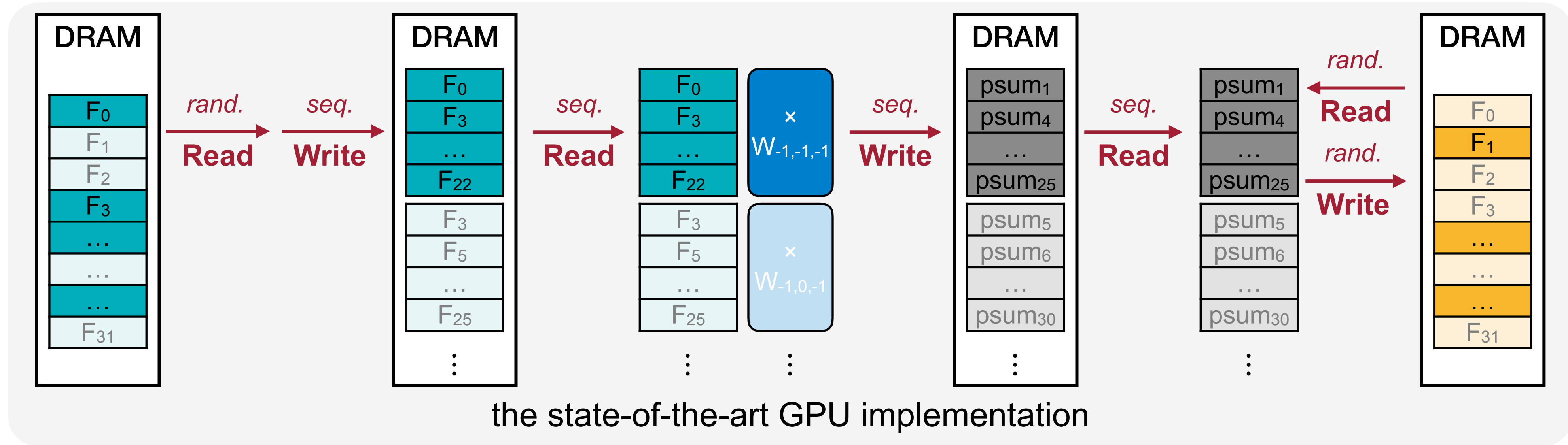
Dense Computation



Sparse MatMul Operations

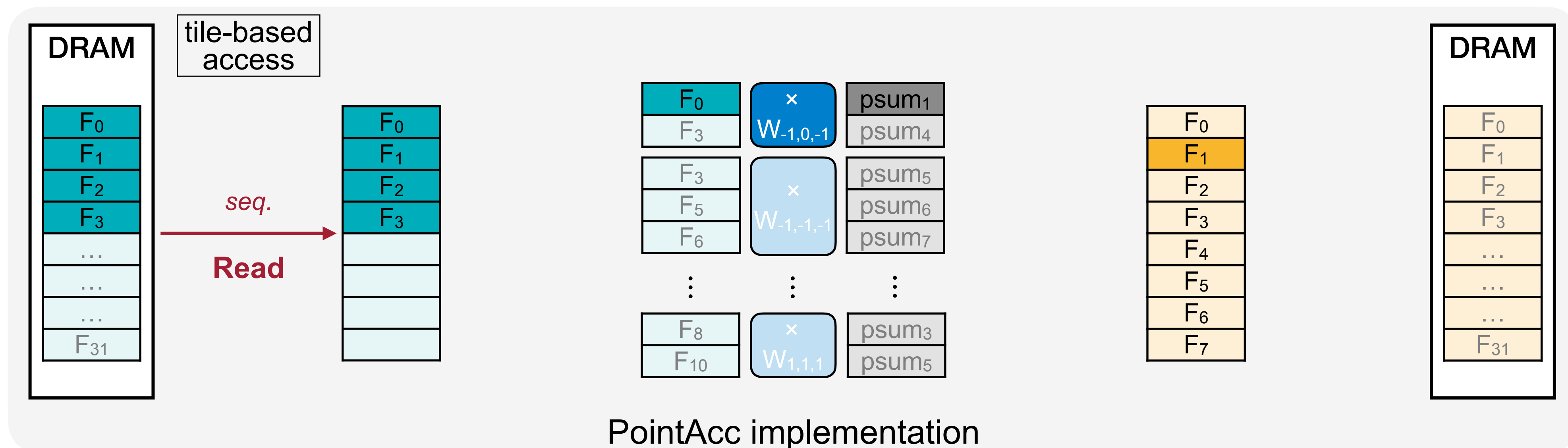


Streaming Sparse MatMul Operations



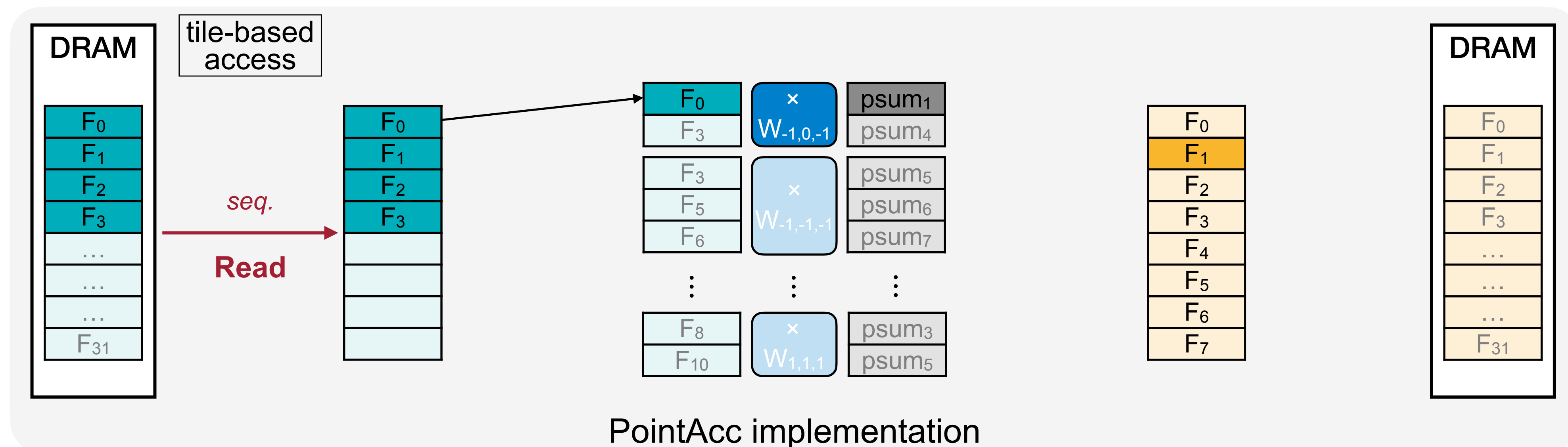
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
 - No more random access for gathering features



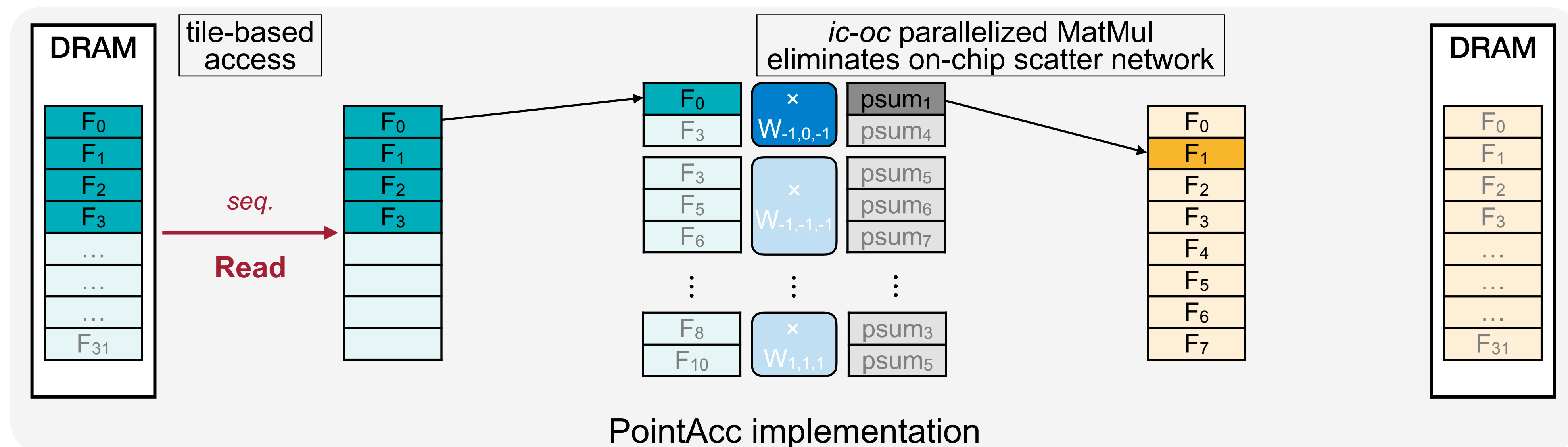
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
 - No more random access for gathering features
- ▶ MatMul computing parallelizes input channel (*ic*) and output channel (*oc*) dimension



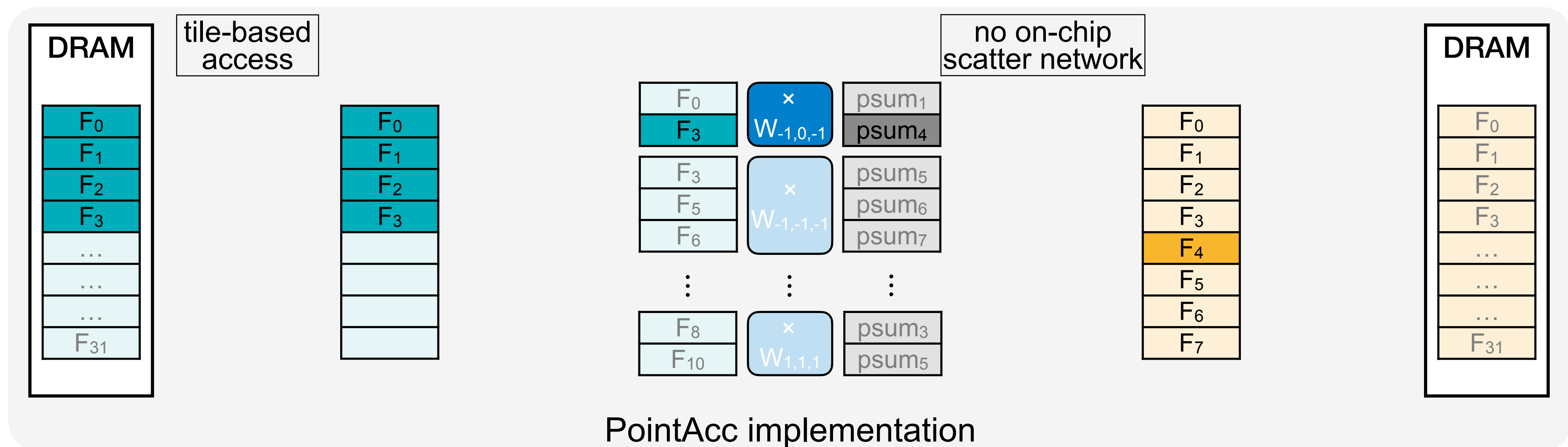
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
 - No more random access for gathering features
- ▶ MatMul computing parallelizes input channel (*ic*) and output channel (*oc*) dimension
 - No need for on-chip scatter network for scattering psums of different points simultaneously



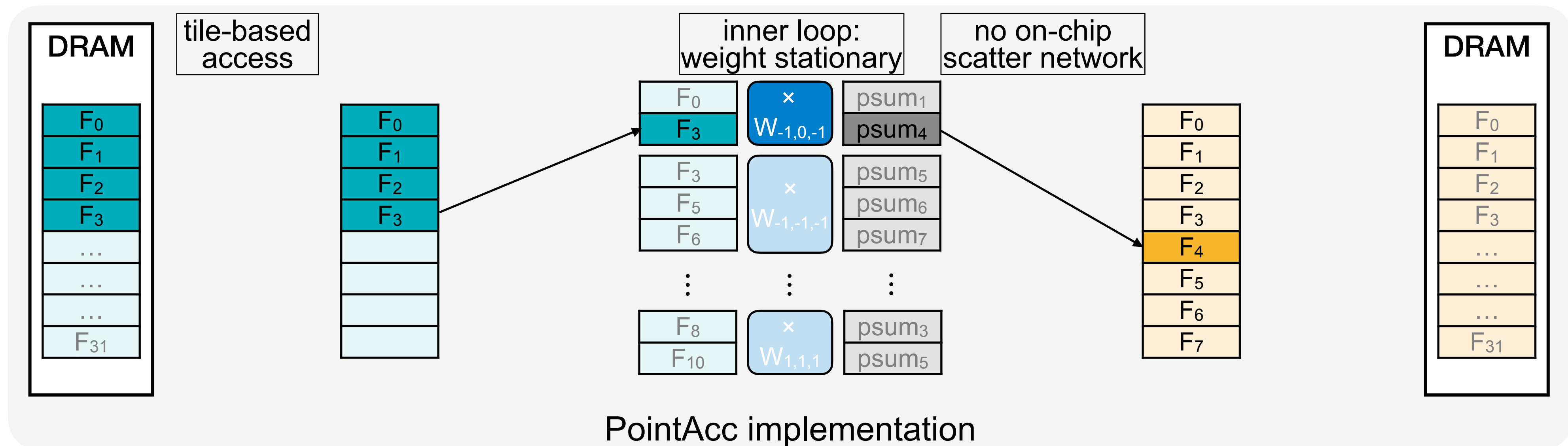
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network



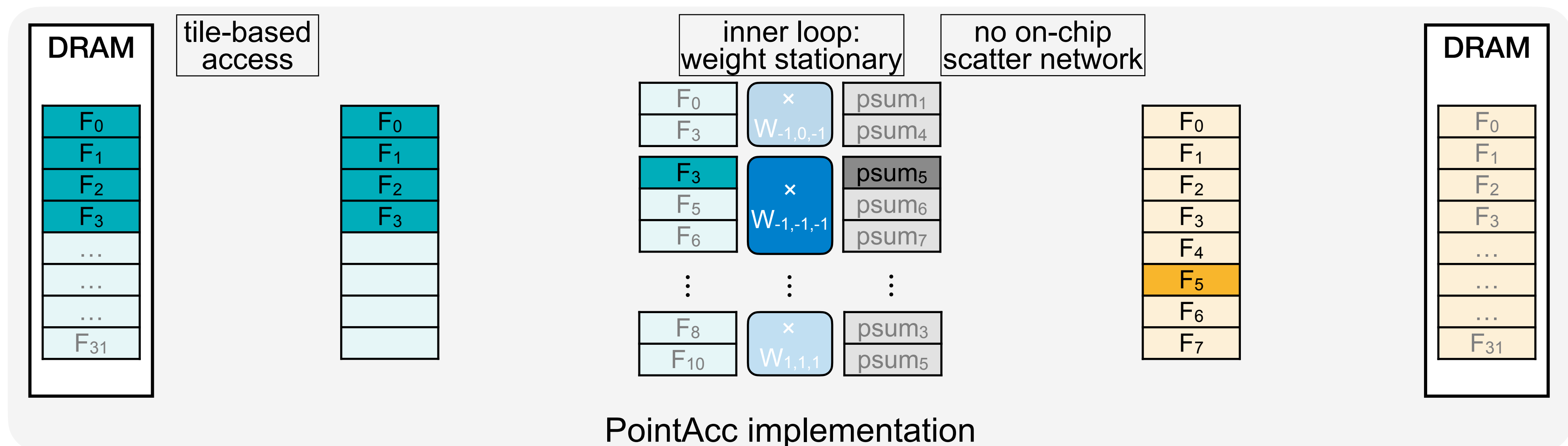
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network
- ▶ Weight stationary saves the on-chip memory footprint
 - Key: #points ($10^3 \sim 10^5$) \gg #channels ($10 \sim 10^3$)



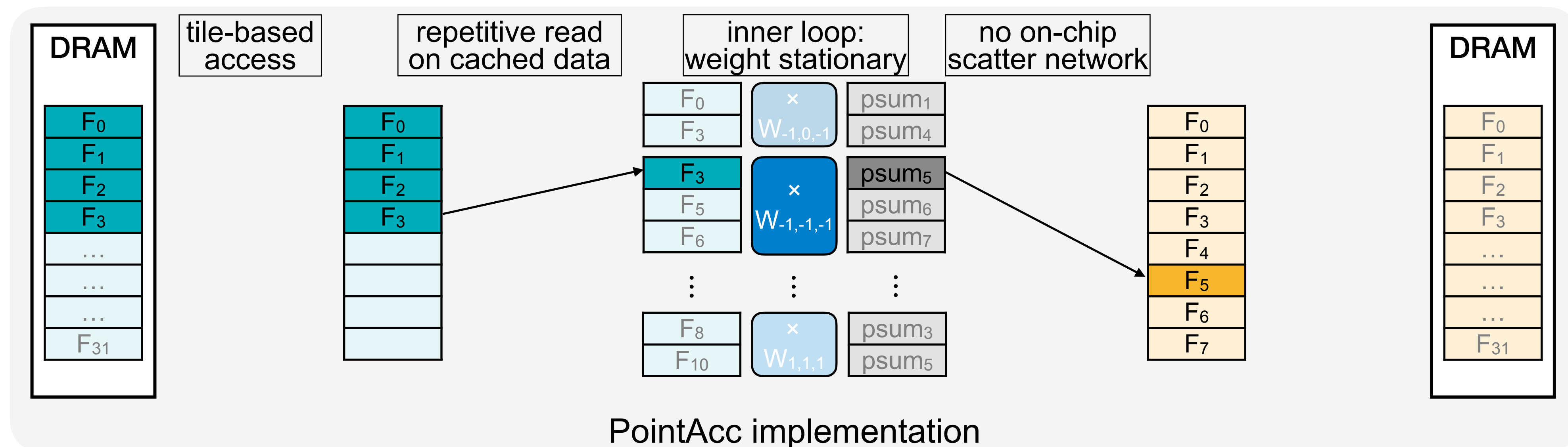
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network
- ▶ Weight stationary saves the on-chip memory footprint



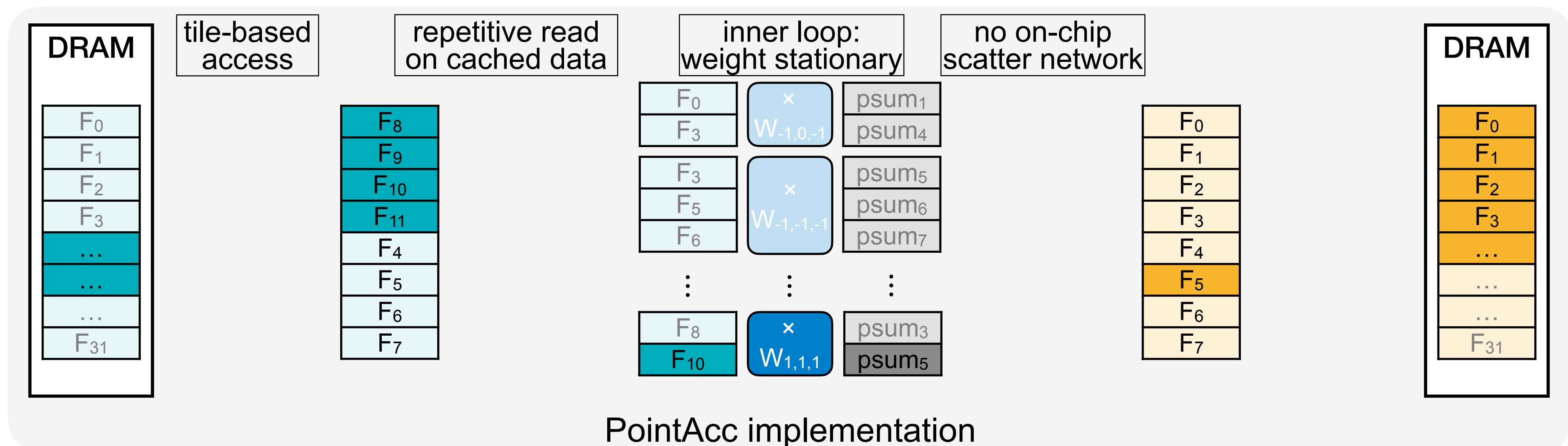
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network
- ▶ Weight stationary saves the on-chip memory footprint
- ▶ Caching reduces the #off-chip read of reused data to nearly 1 access



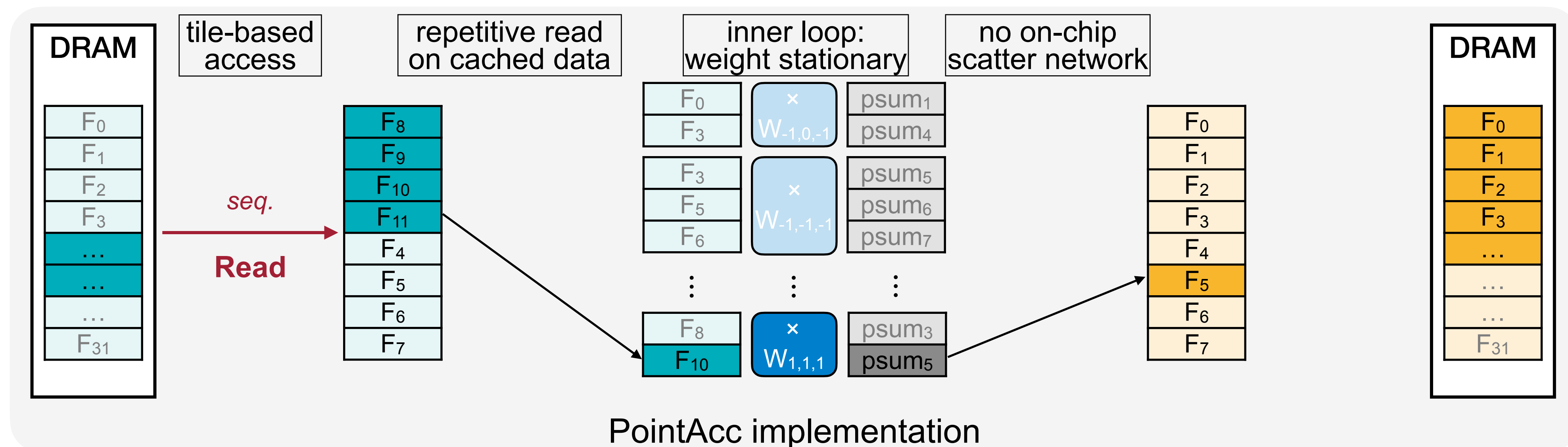
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network
- ▶ Weight stationary saves the on-chip memory footprint
- ▶ Caching reduces the #off-chip read of reused data to nearly 1 access



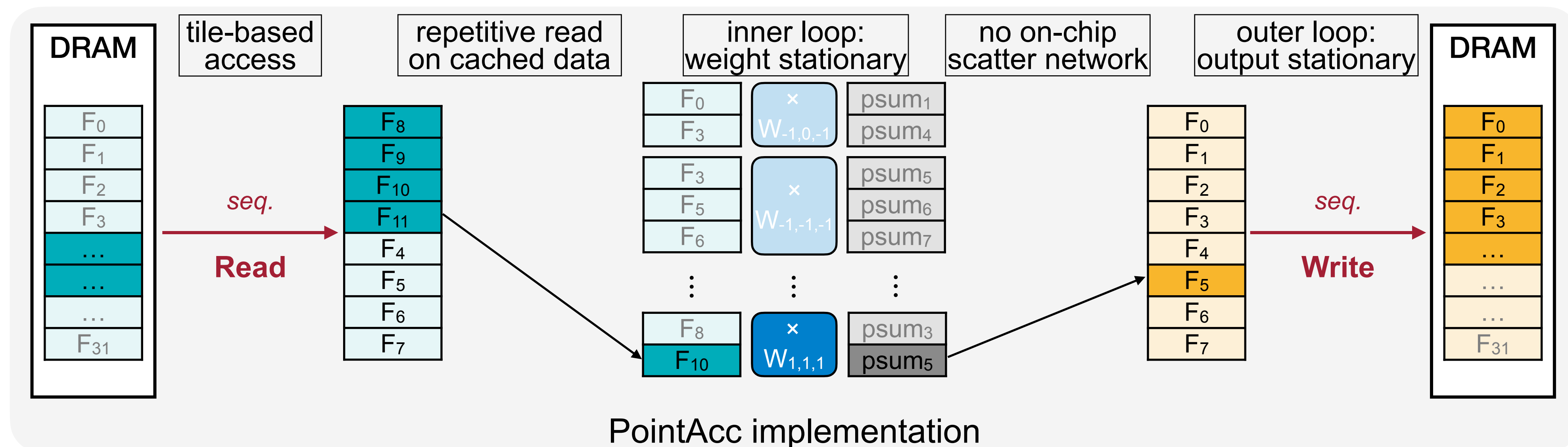
Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network
- ▶ Weight stationary saves the on-chip memory footprint
- ▶ Caching reduces the #off-chip read of reused data to nearly 1 access

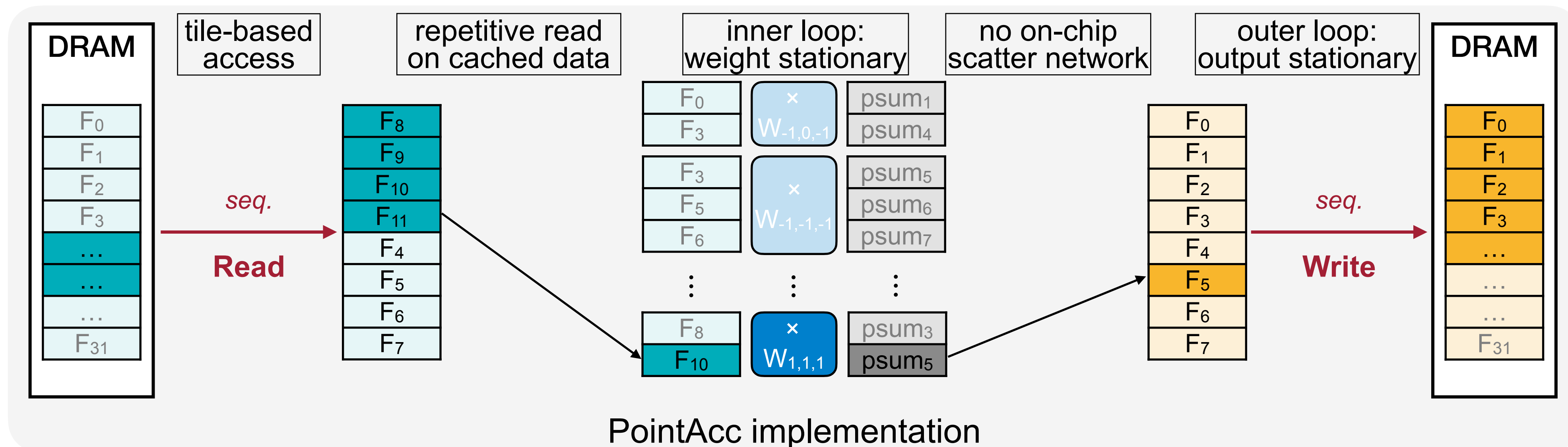
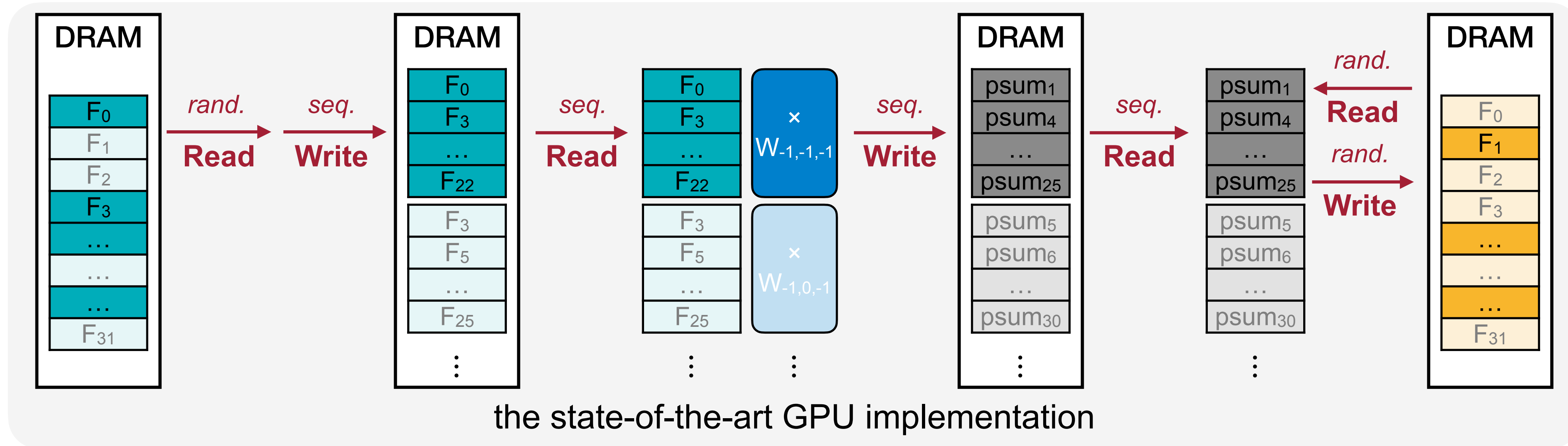


Streaming Sparse MatMul Operations

- ▶ Sequential fetch input features on demand in the granularity of tile (*i.e.*, “cache block”)
- ▶ MatMul computing parallelizes *ic* and *oc* dimension → no need for on-chip scatter network
- ▶ Weight stationary saves the on-chip memory footprint
- ▶ Caching reduces the #off-chip read of reused data to nearly 1 access
- ▶ Output stationary eliminates the off-chip scattering of partial sums

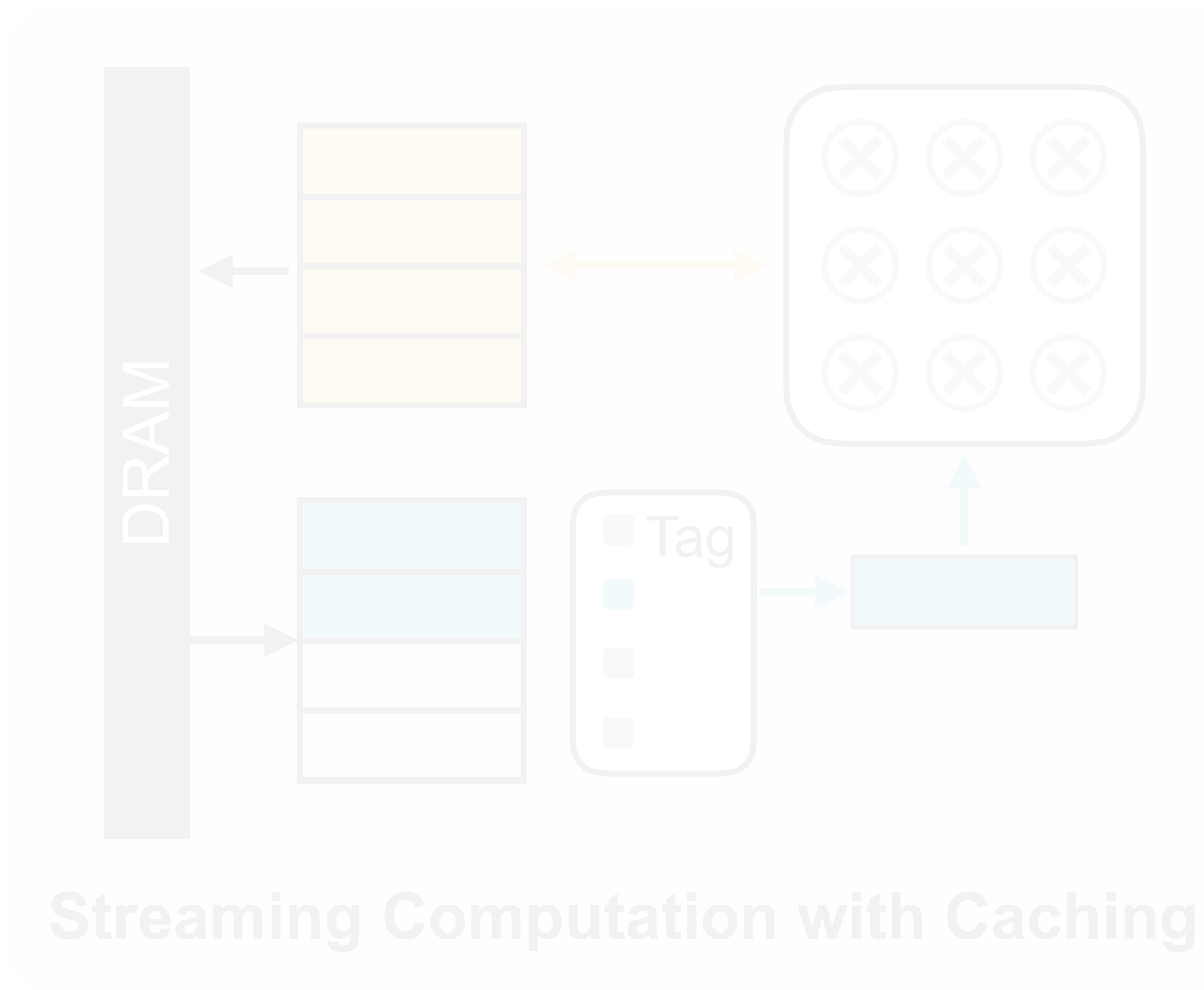


Streaming Sparse MatMul Operations

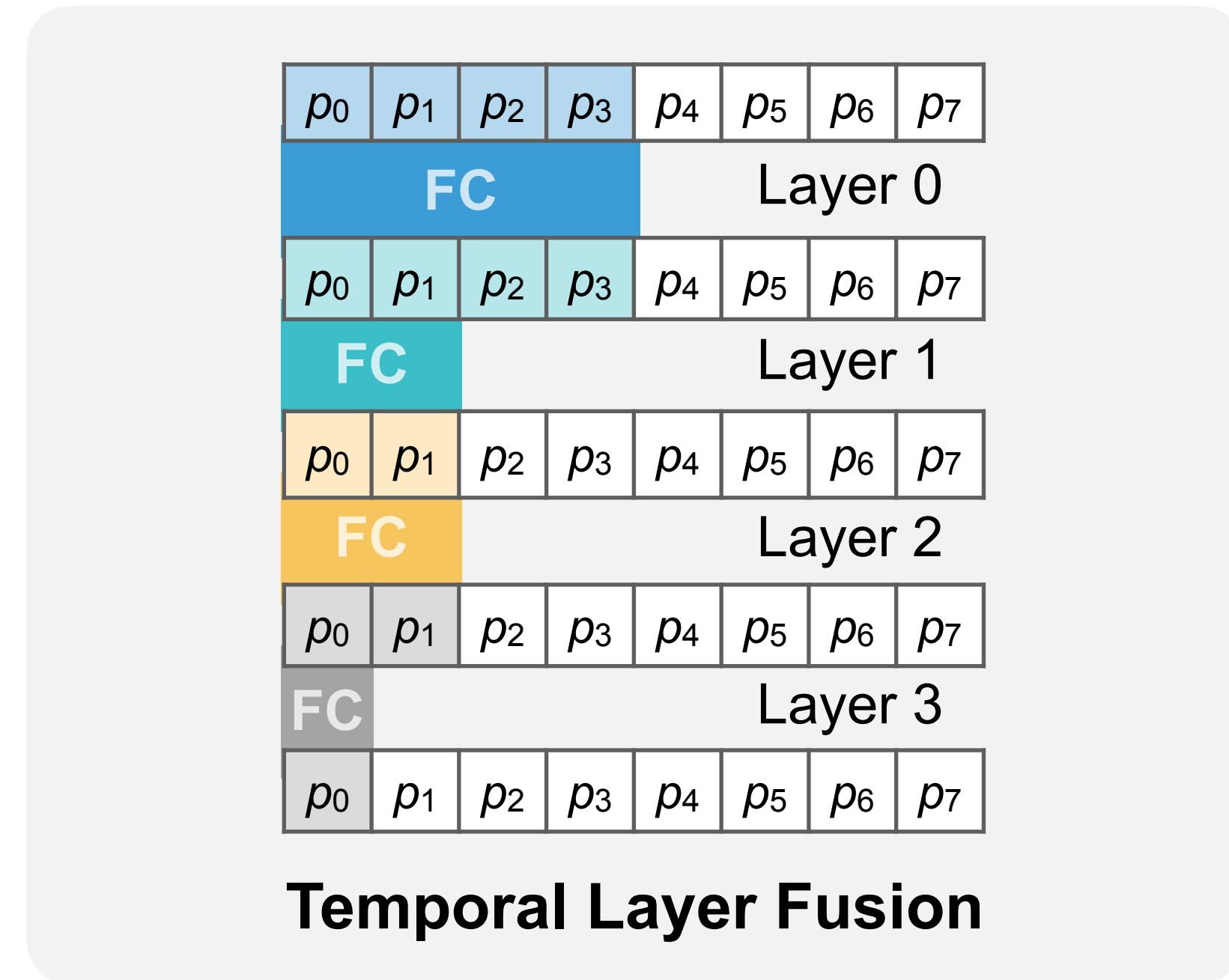


Flexible Memory Management

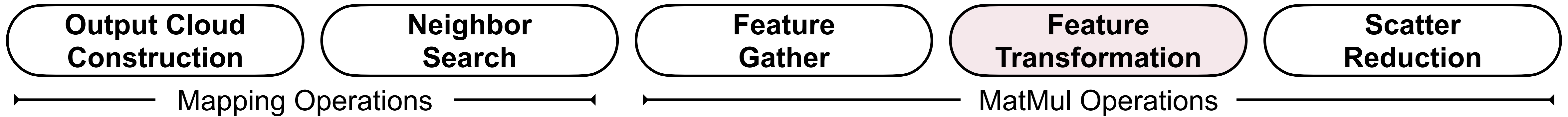
Sparse Computation



Dense Computation



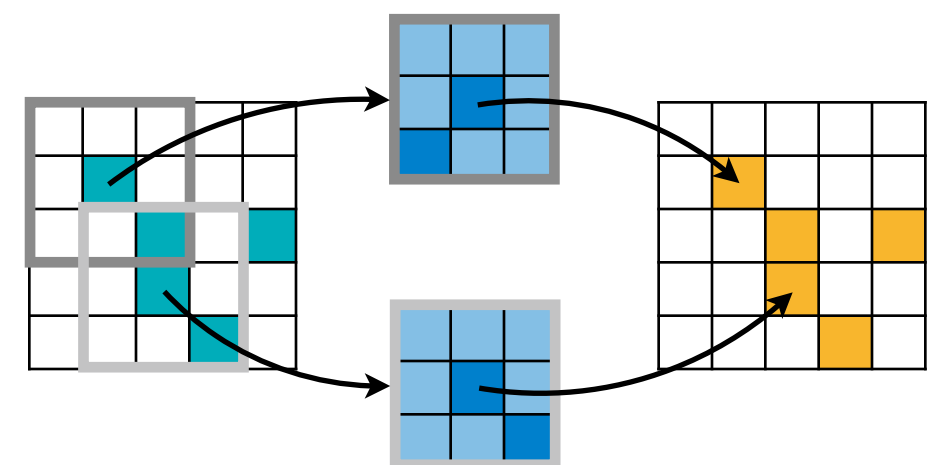
Step 4: Transform Input Features



Input Point Cloud

(P_0, F_0)
(P_1, F_1)
(P_2, F_2)
(P_3, F_3)
(P_4, F_4)

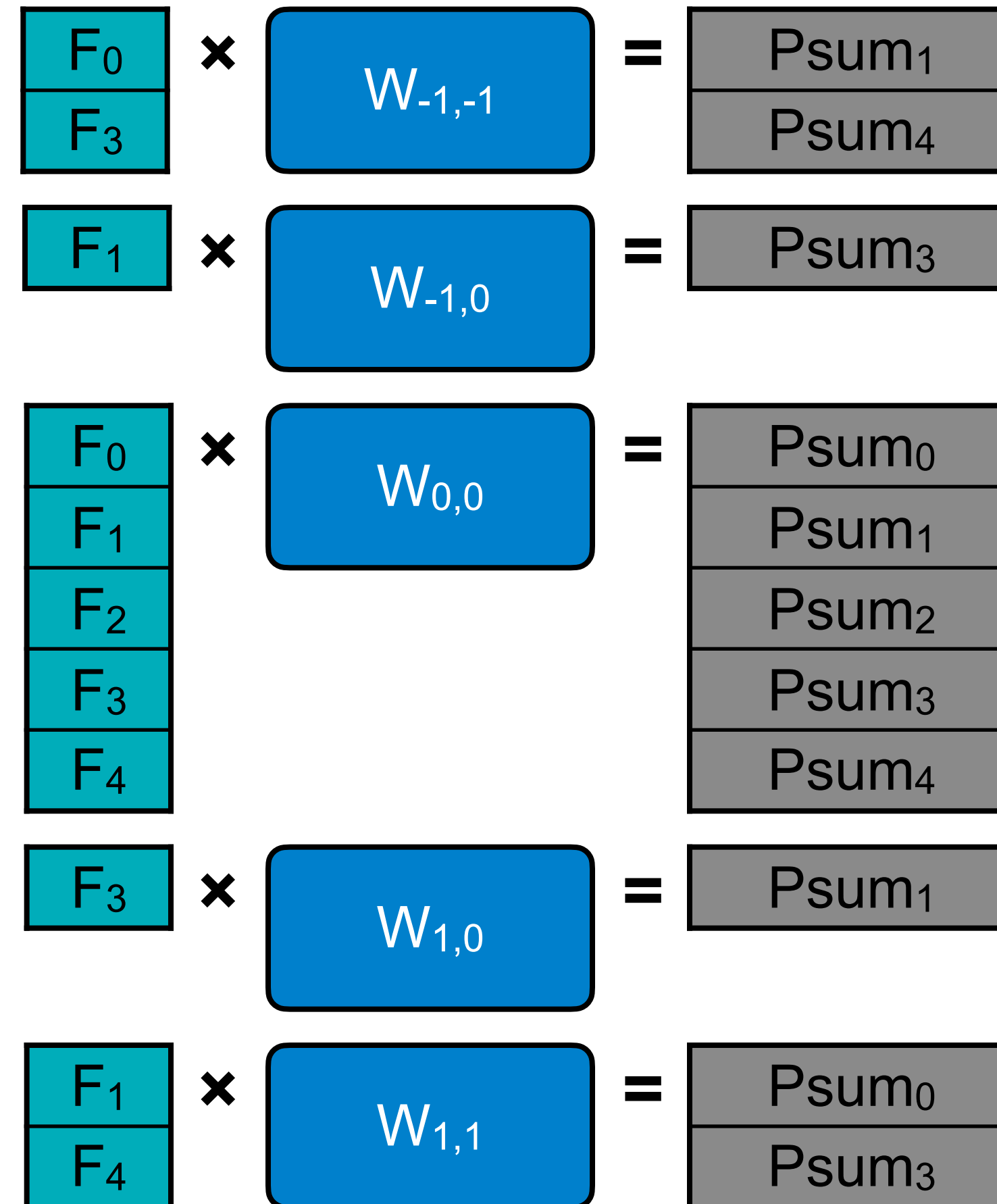
(Coords, Feature Vector)



Maps (In, Out, Wgt)		
$(P_0, Q_1, W_{-1,-1})$		
$(P_3, Q_4, W_{-1,-1})$		
$(P_1, Q_3, W_{-1,0})$		
$(P_0, Q_0, W_{0,0})$		
$(P_1, Q_1, W_{0,0})$		
$(P_2, Q_2, W_{0,0})$		
$(P_3, Q_3, W_{0,0})$		
$(P_4, Q_4, W_{0,0})$		
$(P_3, Q_1, W_{1,0})$		
$(P_1, Q_0, W_{1,1})$		
$(P_4, Q_3, W_{1,1})$		

Gather
By
Weight

X



p -by- ic ic -by- oc p -by- oc \times n

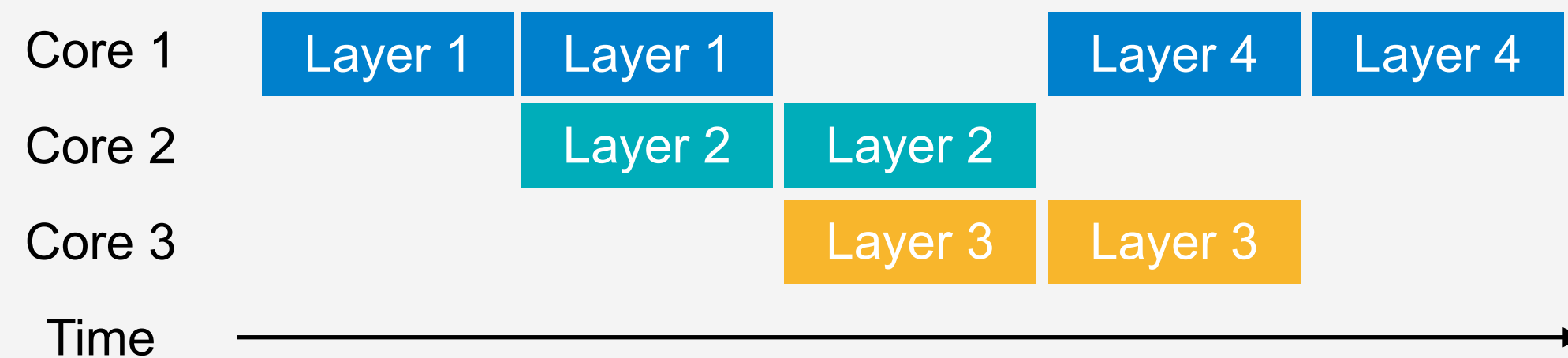
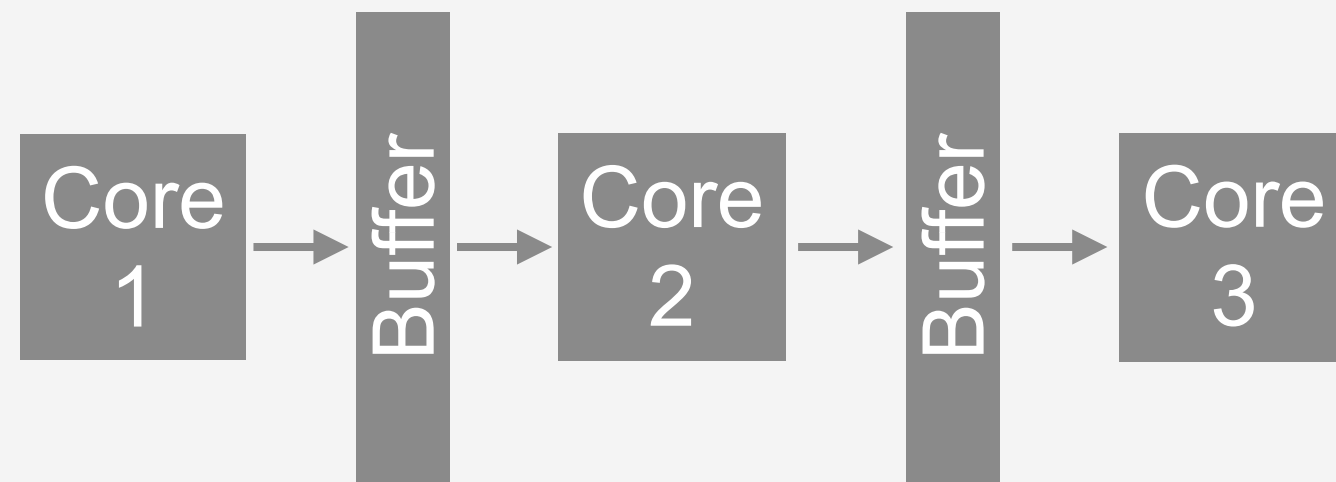
Matrix-Matrix Multiplication

Output Point Cloud

Q_0
Q_1
Q_2
Q_3
Q_4

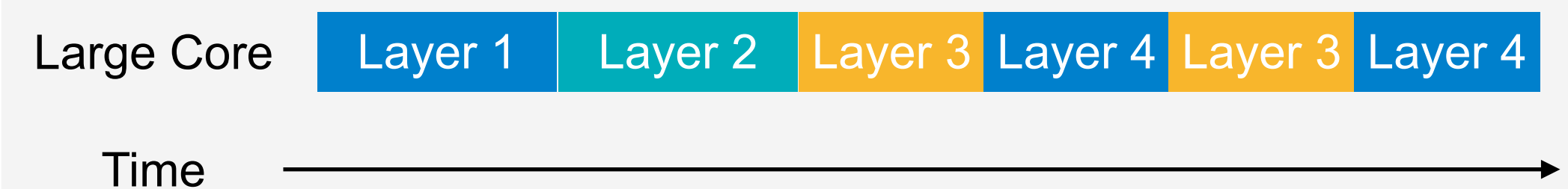
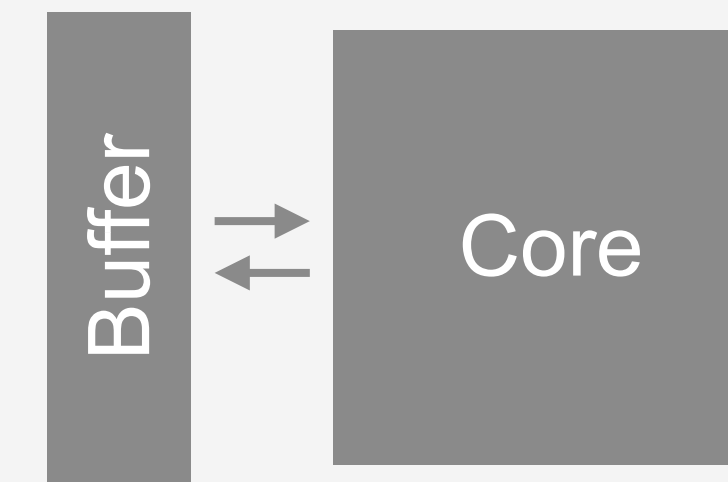
Layer Fusion

Spatial Layer Fusion



- ✓ Reduced off-chip memory access
- ✗ Complicated logic overhead
- ✗ Inflexible intermediate buffer use
- ✗ Fixed #layers for fusion

Temporal Layer Fusion



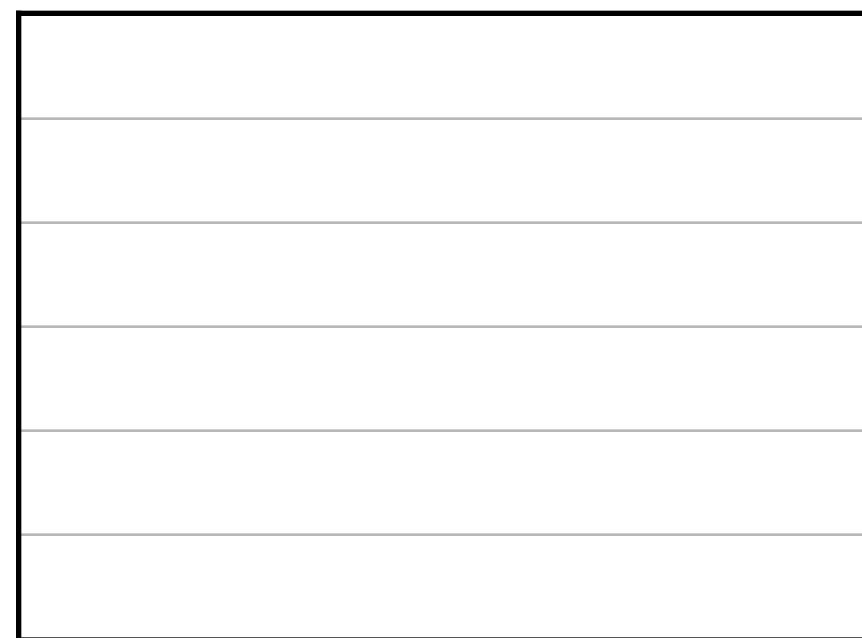
- ✓ Reduced off-chip memory access
- ✓ Simple logic modification
- ✓ Flexible buffer use for better layer tiling
- ✓ Flexible #layers for fusion

consecutive FCs varies among different point cloud NNs, and even among different blocks of the same model

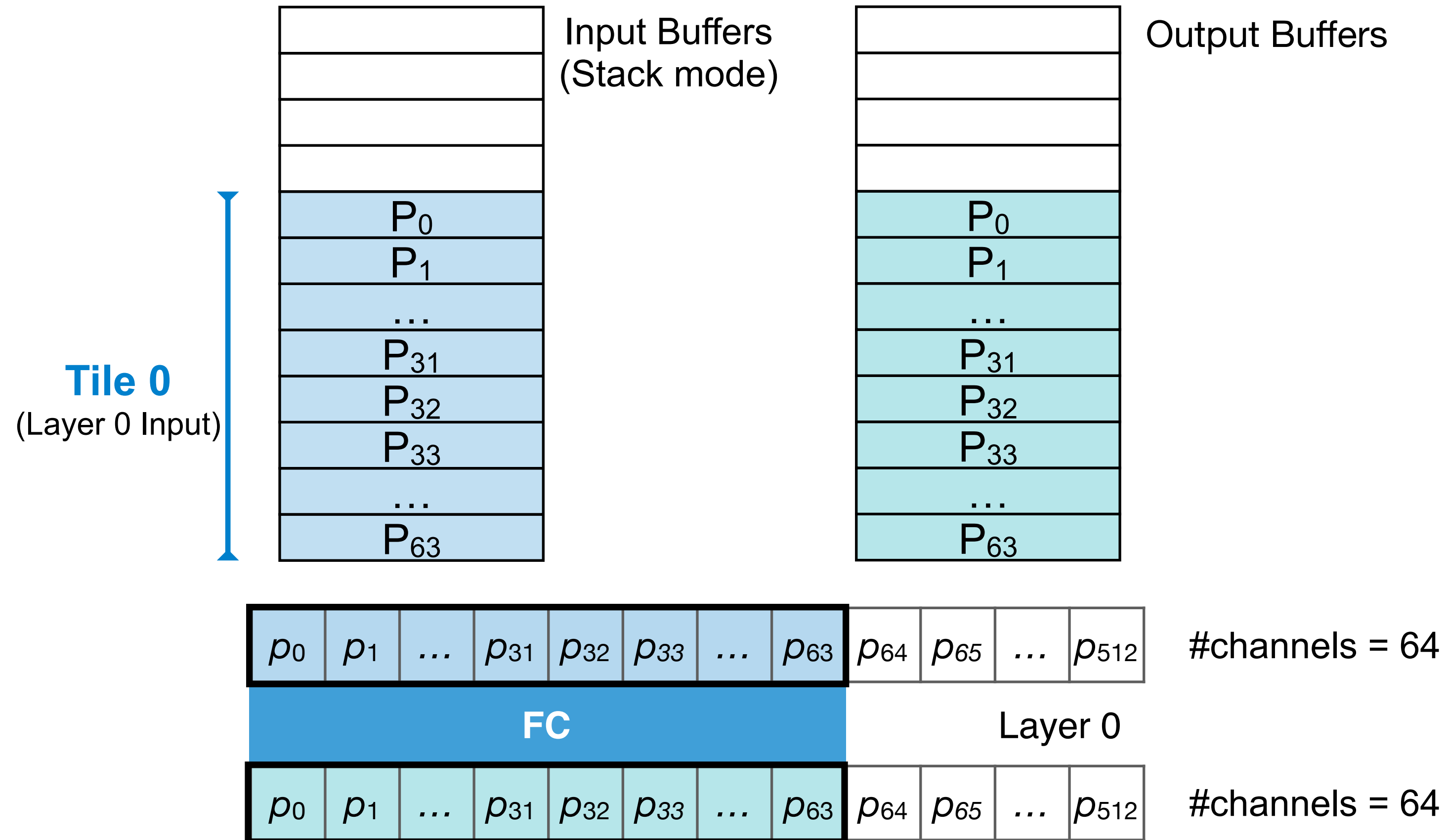
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read
 write



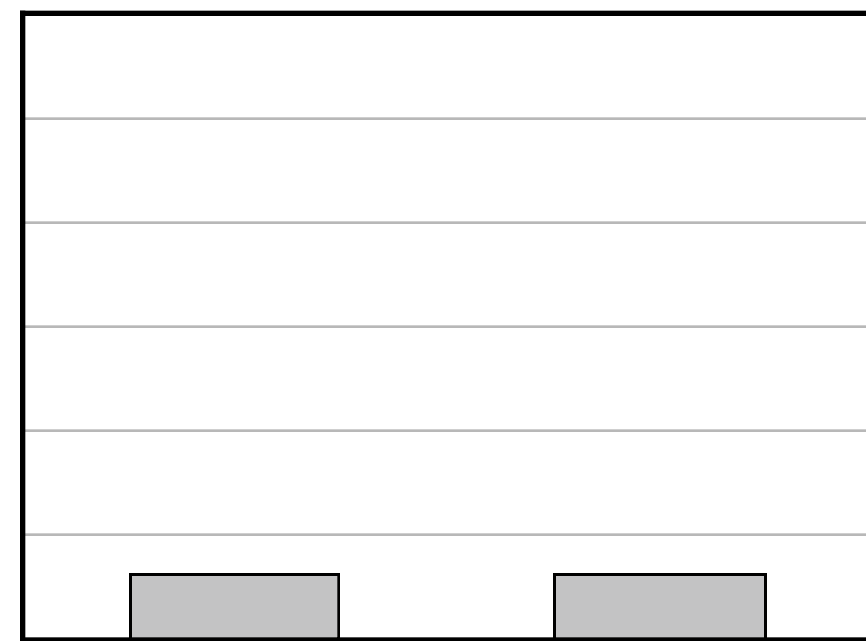
without layer fusion with layer fusion



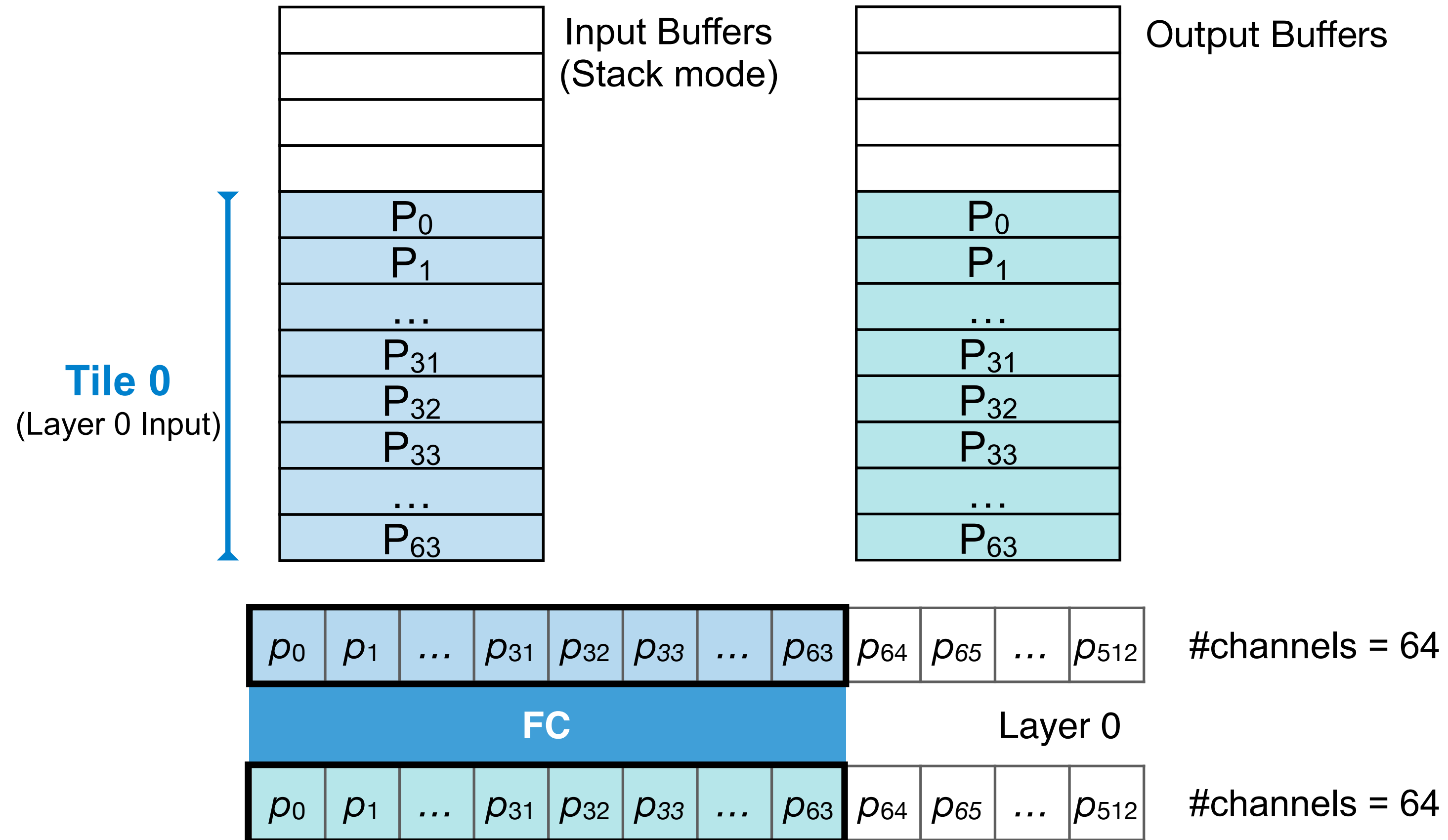
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



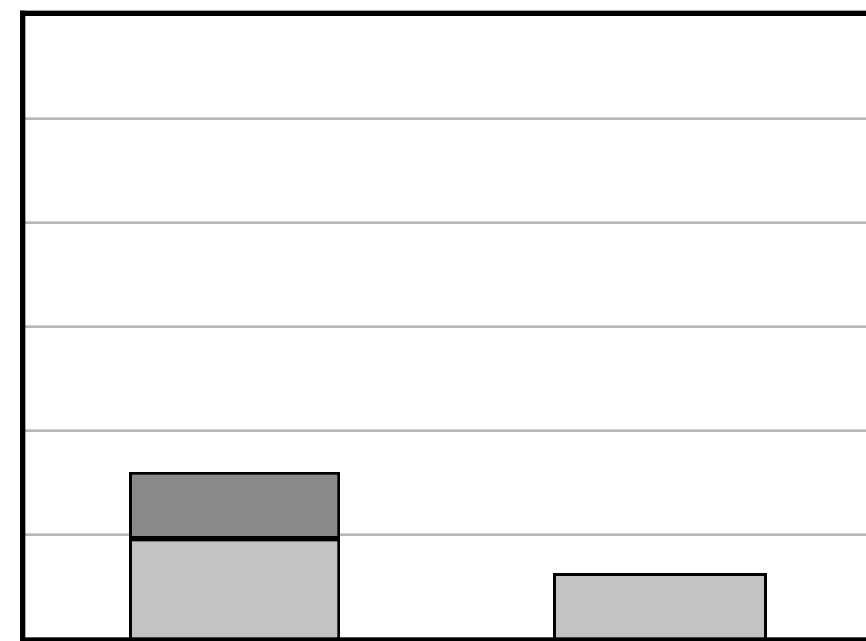
without layer fusion with layer fusion



Temporal Layer Fusion on consecutive FCs

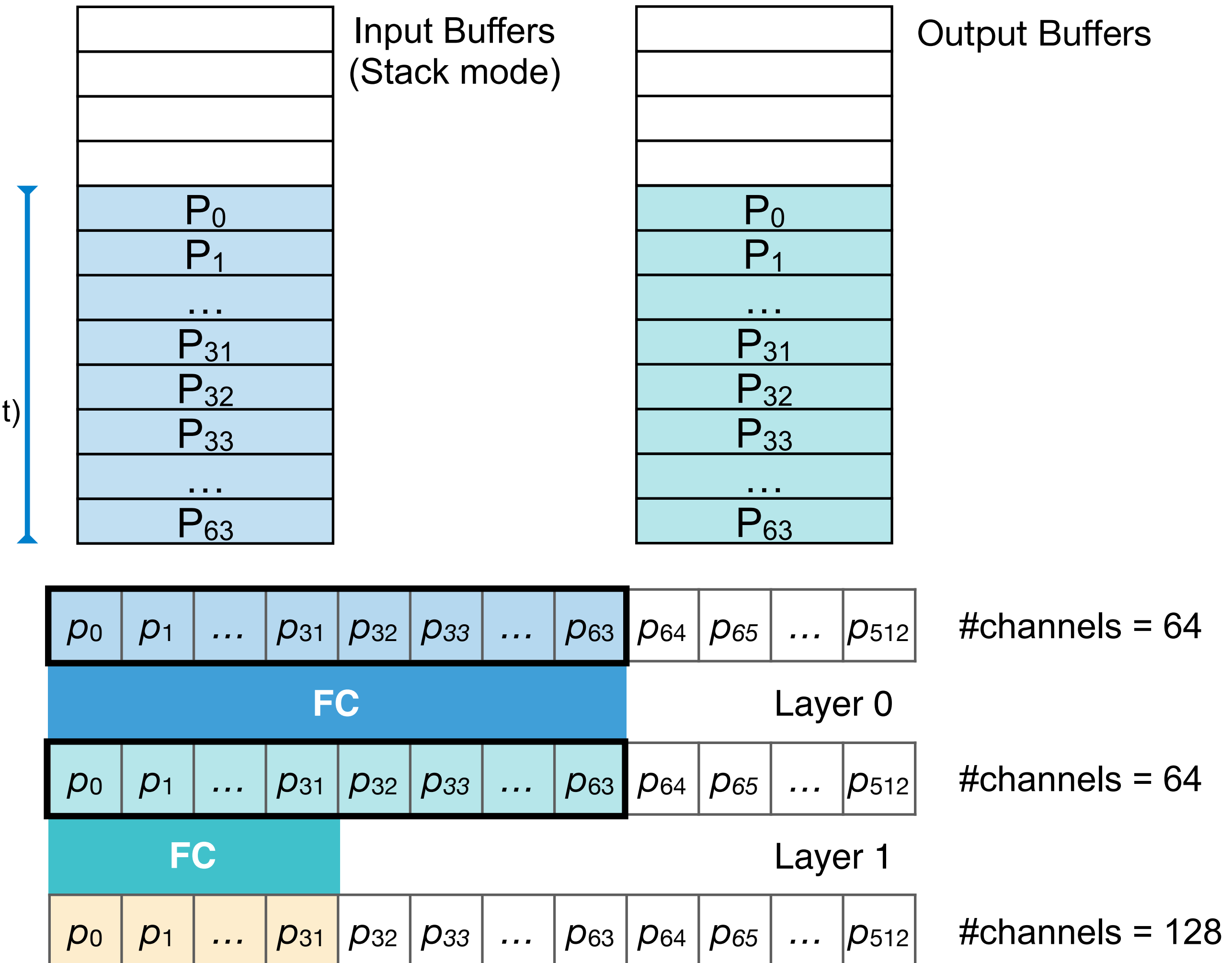
DRAM access per point

read write



without layer fusion with layer fusion

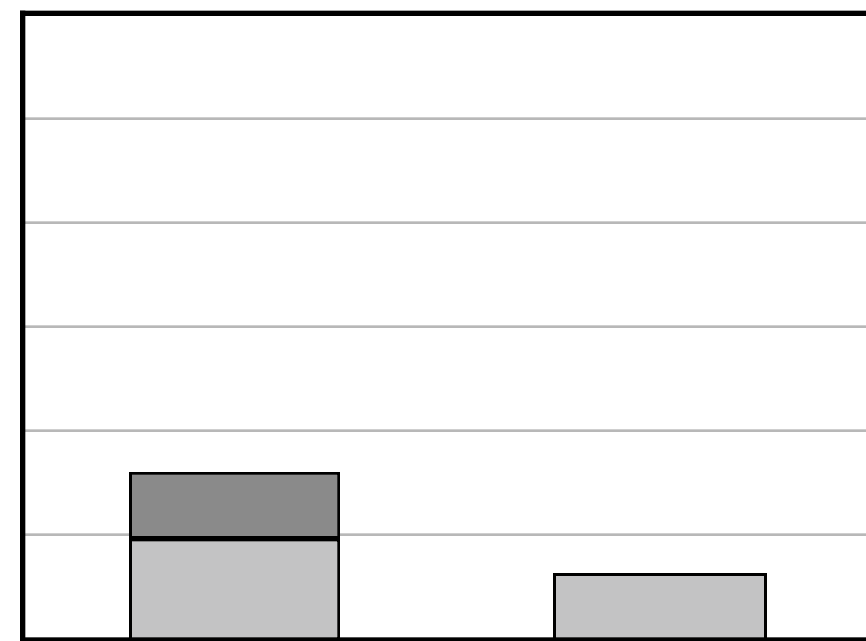
Tile 0
(Layer 0 Input)



Temporal Layer Fusion on consecutive FCs

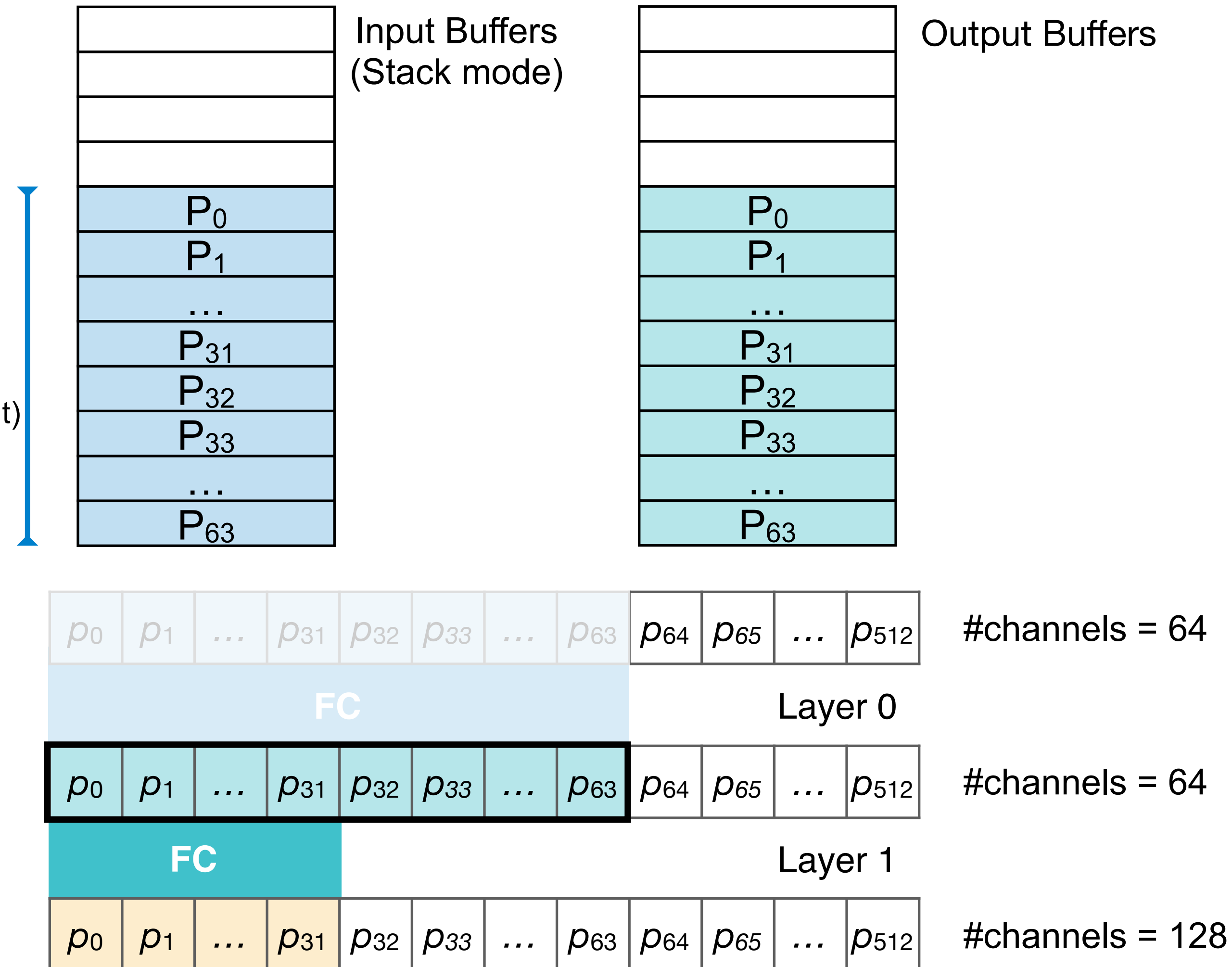
DRAM access per point

read write



without layer fusion with layer fusion

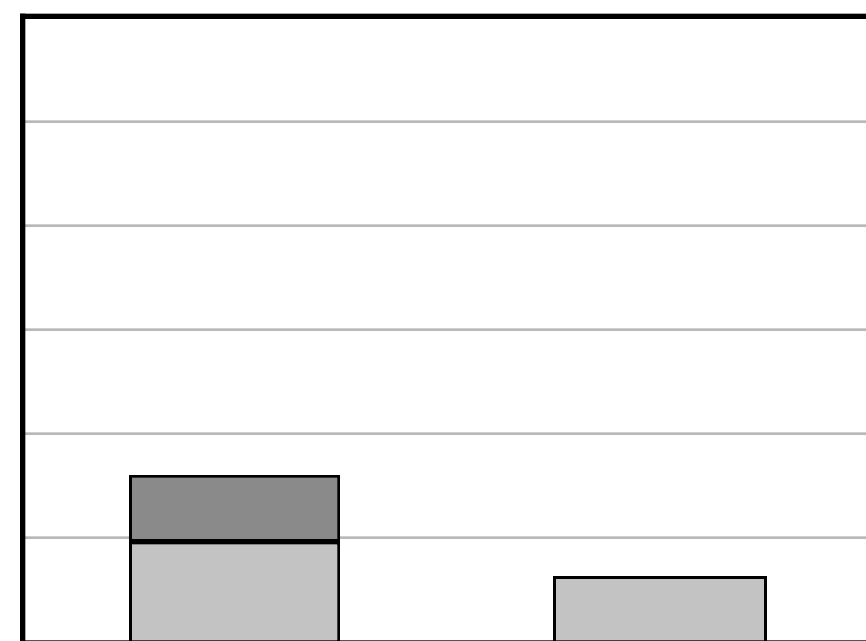
Tile 0
(Layer 0 Input)



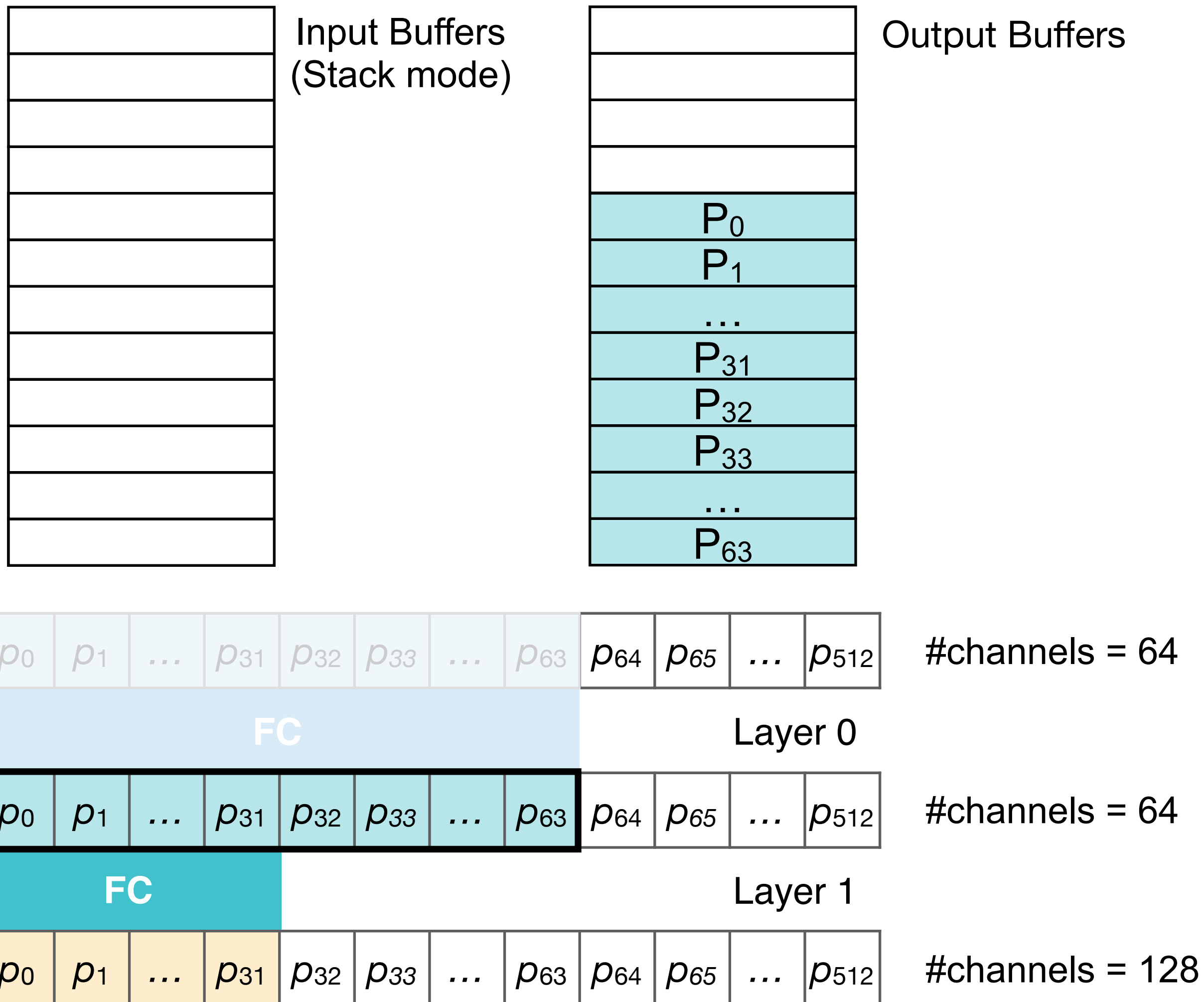
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



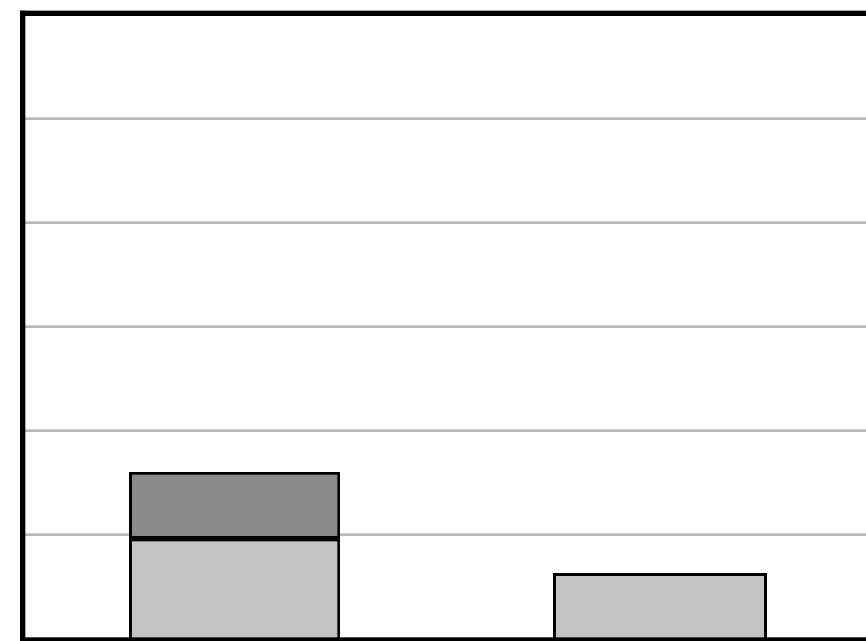
without layer fusion with layer fusion



Temporal Layer Fusion on consecutive FCs

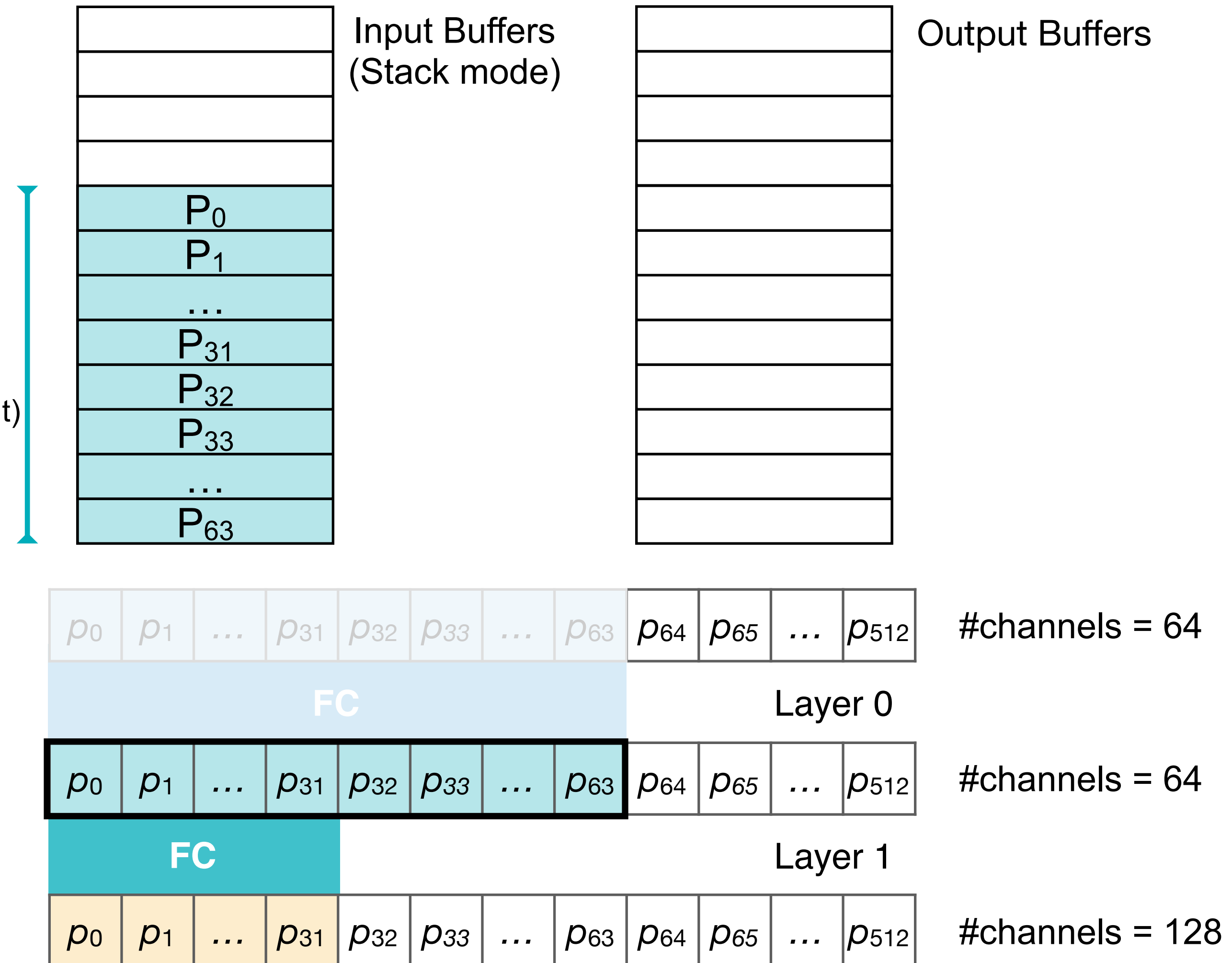
DRAM access per point

read write



without layer fusion with layer fusion

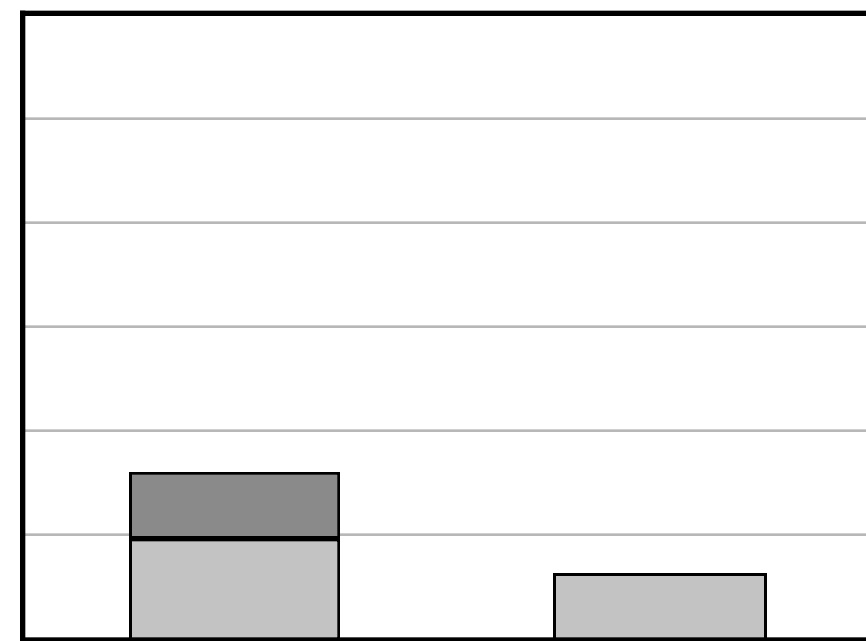
Tile 1
(Layer 1 Input)



Temporal Layer Fusion on consecutive FCs

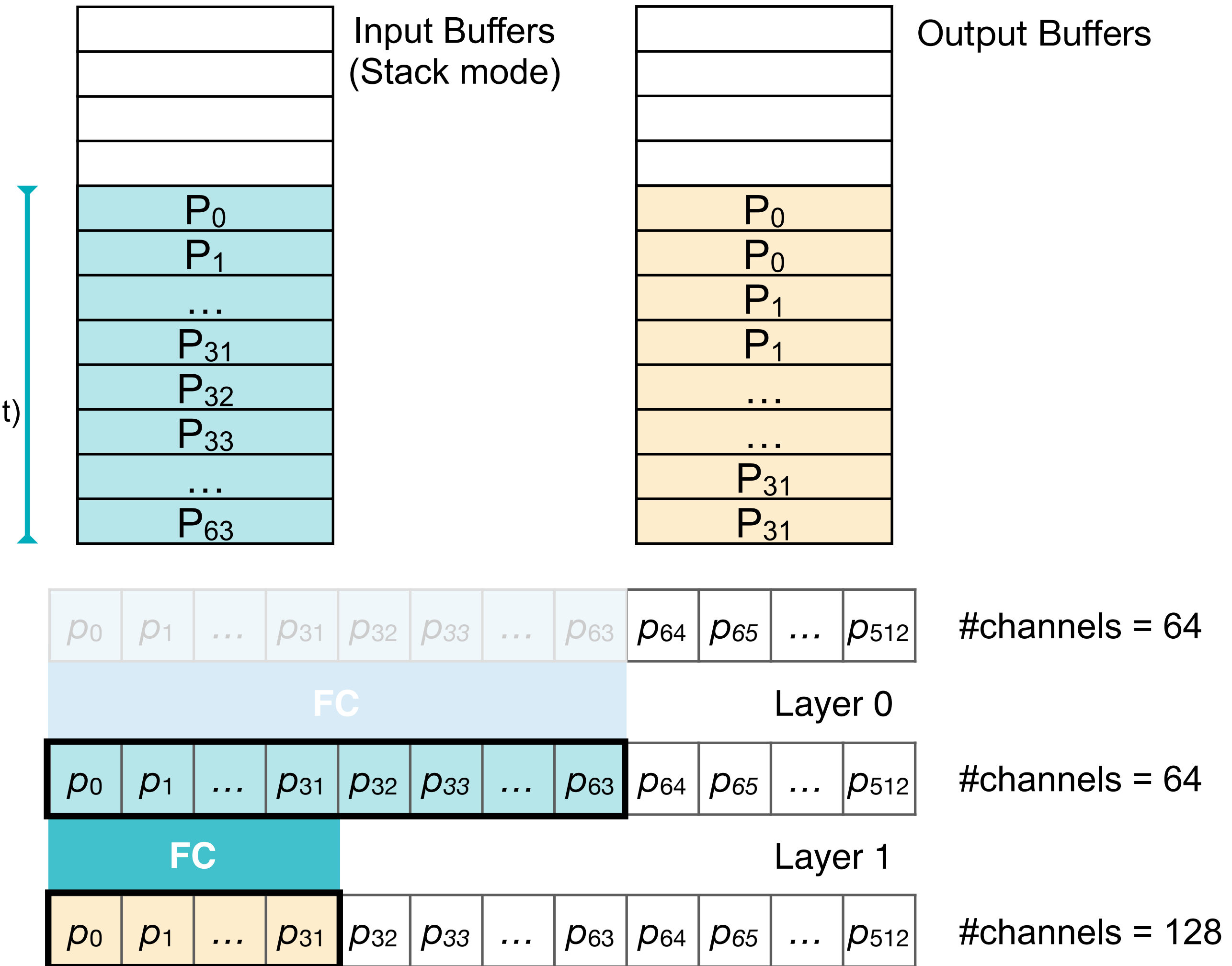
DRAM access per point

read write



without layer fusion with layer fusion

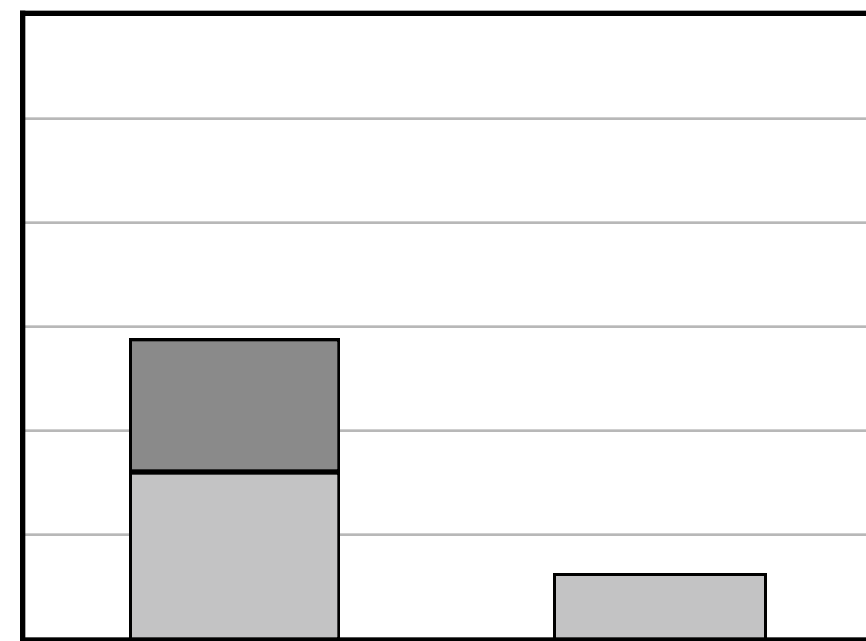
Tile 1
(Layer 1 Input)



Temporal Layer Fusion on consecutive FCs

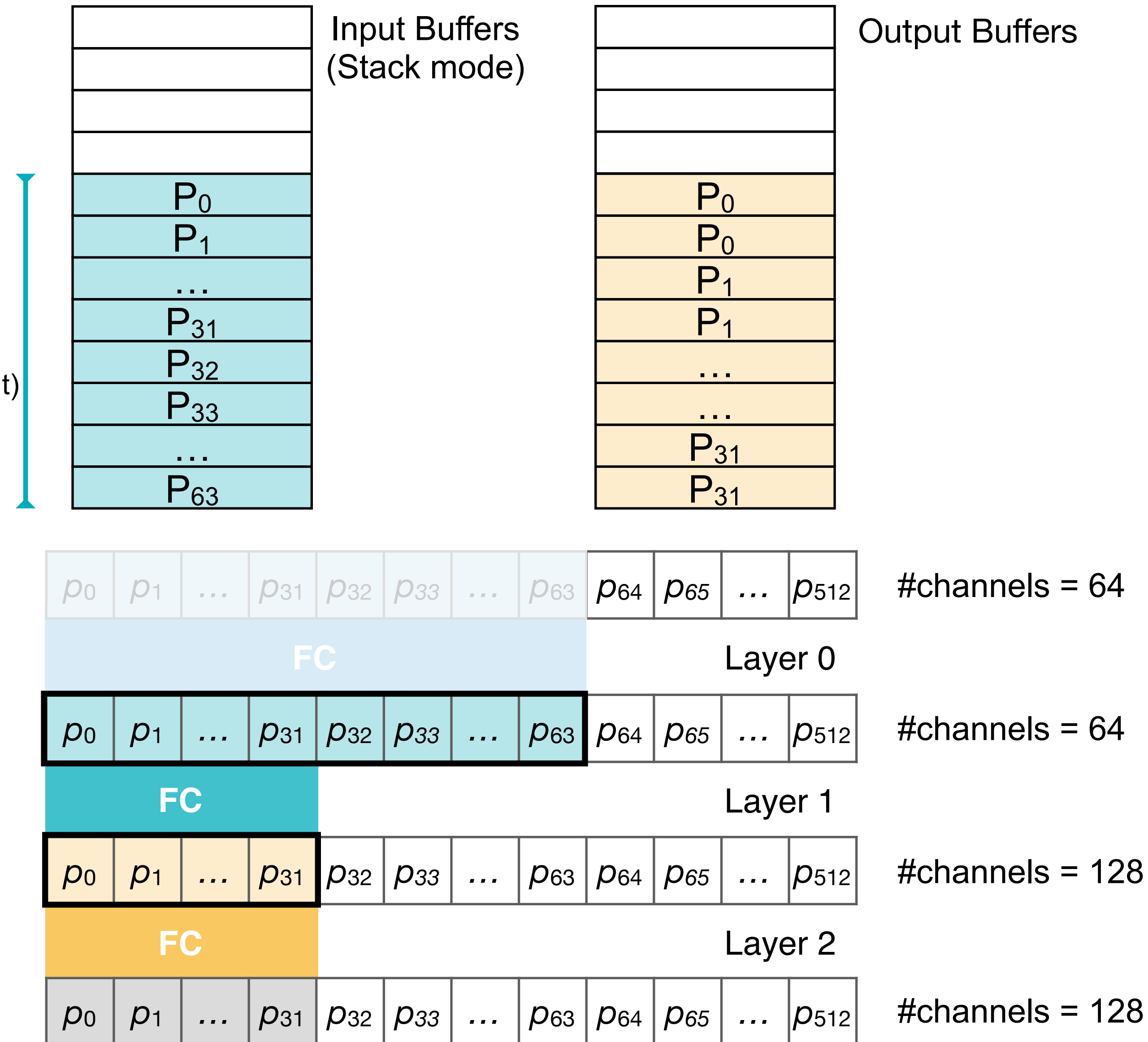
DRAM access per point

read write



without layer fusion with layer fusion

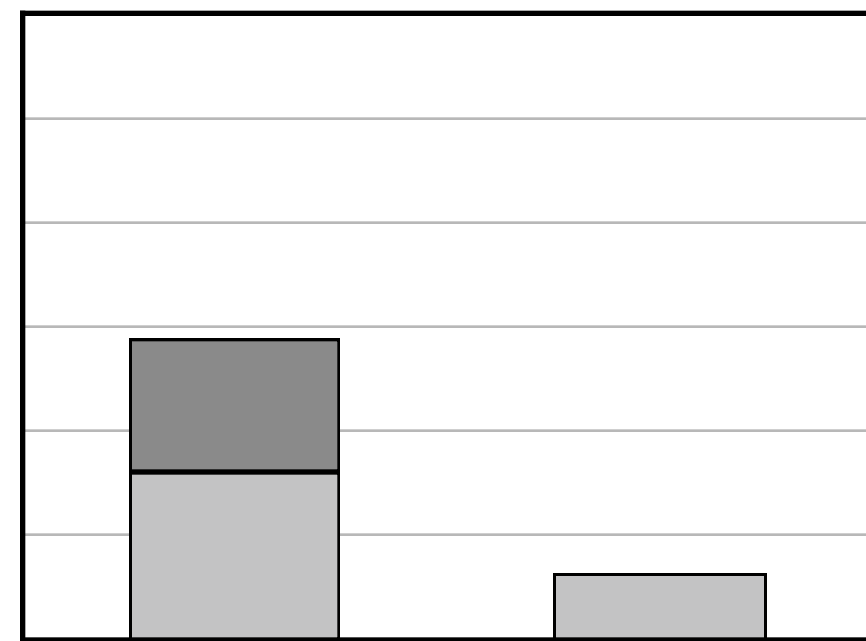
Tile 1
(Layer 1 Input)



Temporal Layer Fusion on consecutive FCs

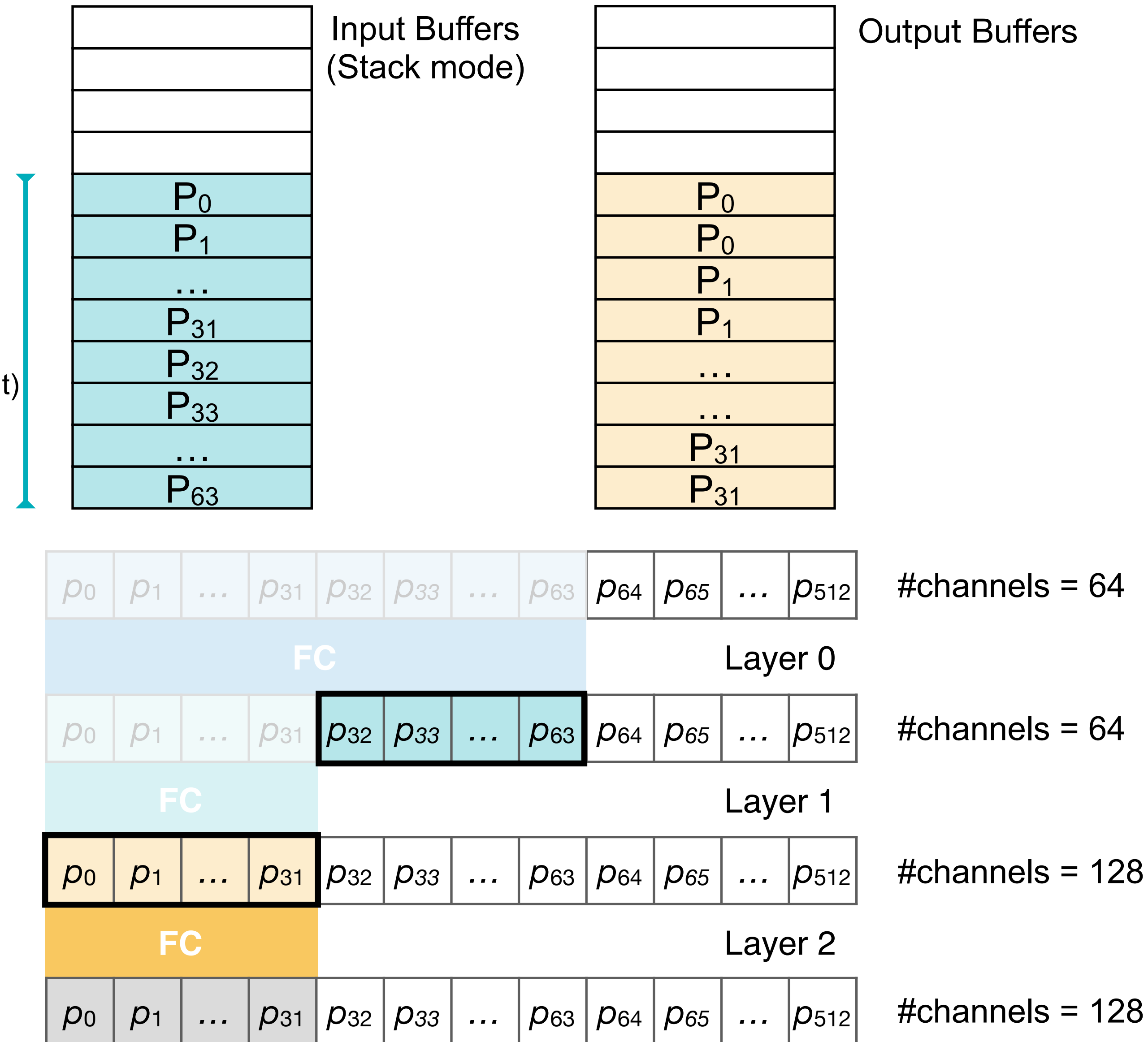
DRAM access per point

read write



without layer fusion with layer fusion

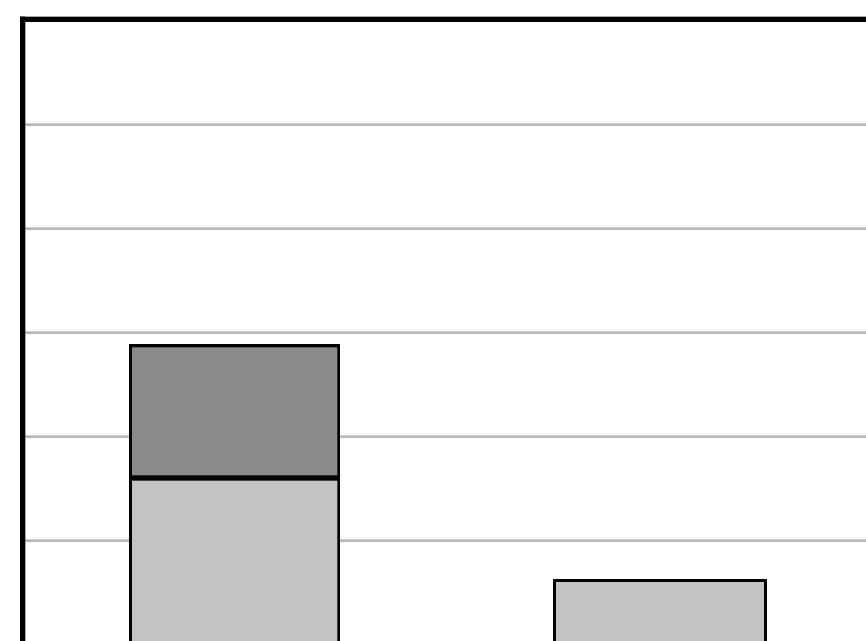
Tile 1
(Layer 1 Input)



Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write

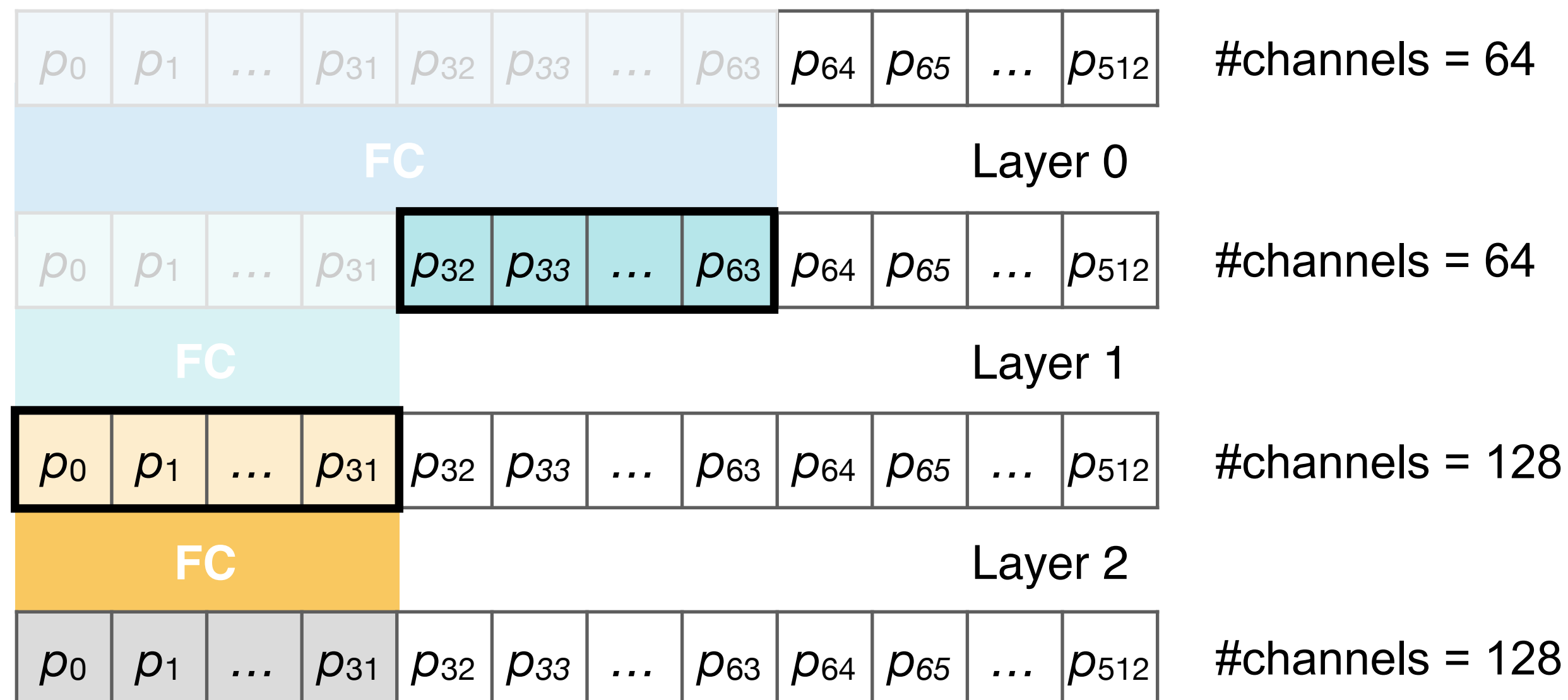
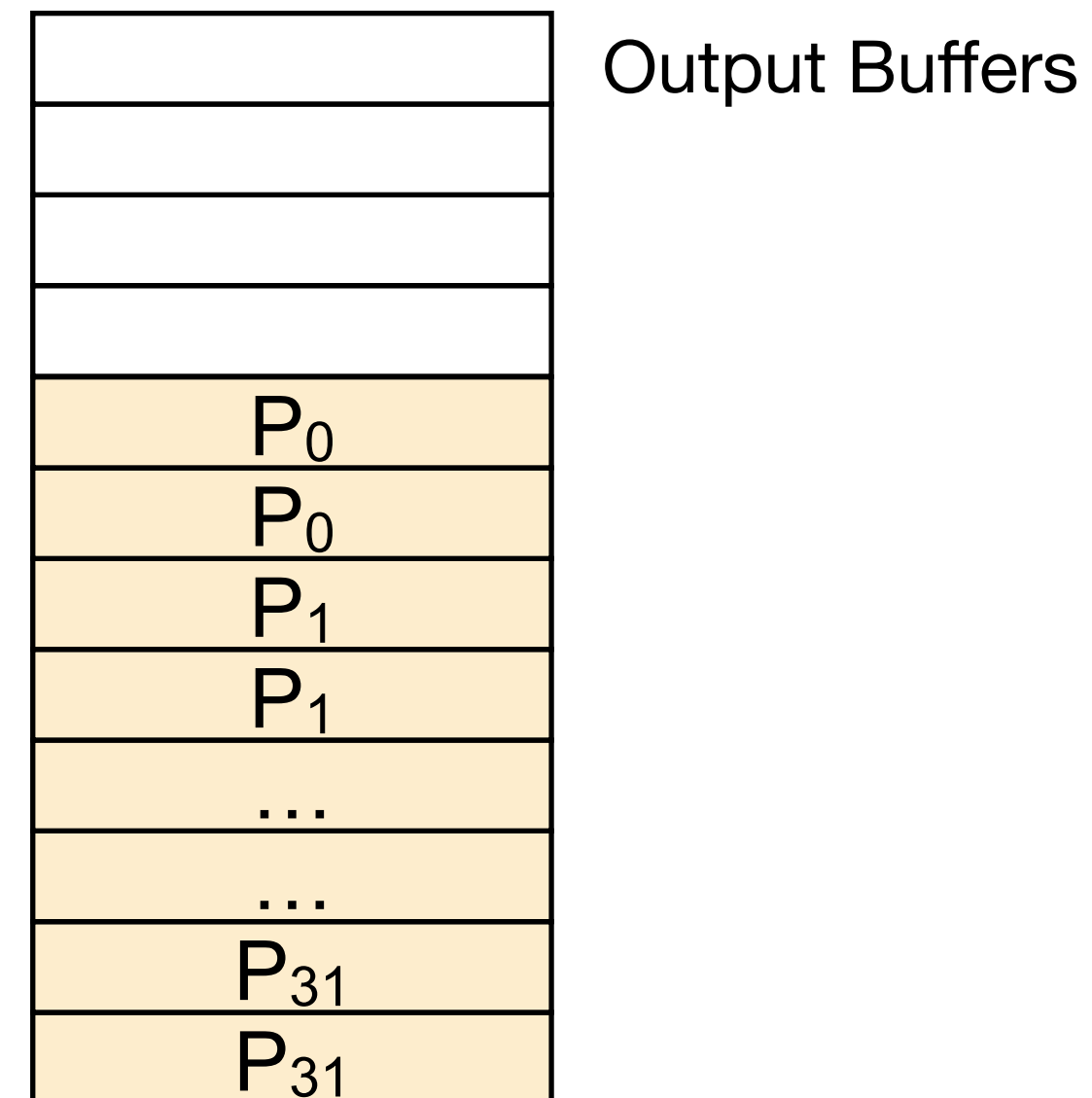
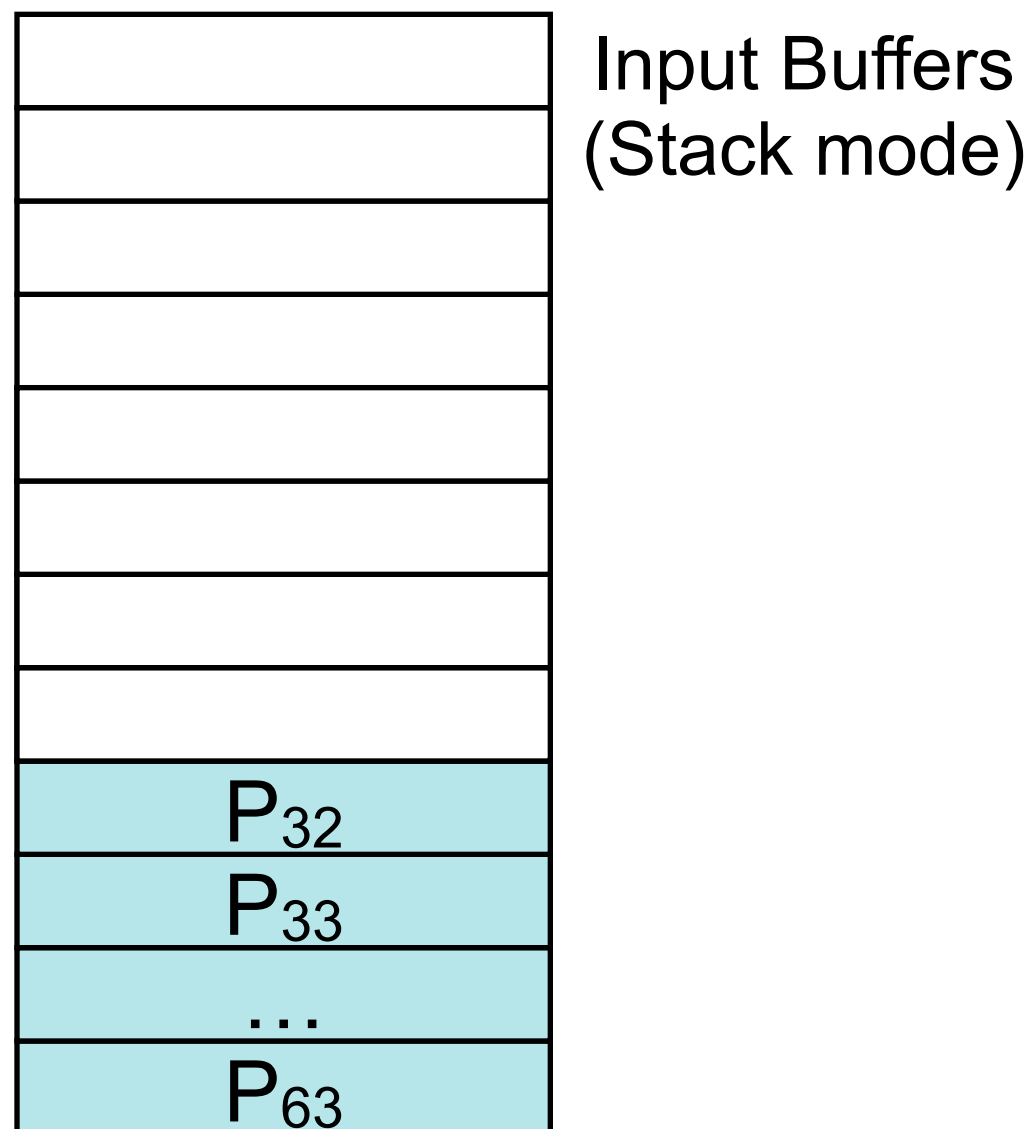


without layer fusion

with layer fusion

without layer fusion with layer fusion

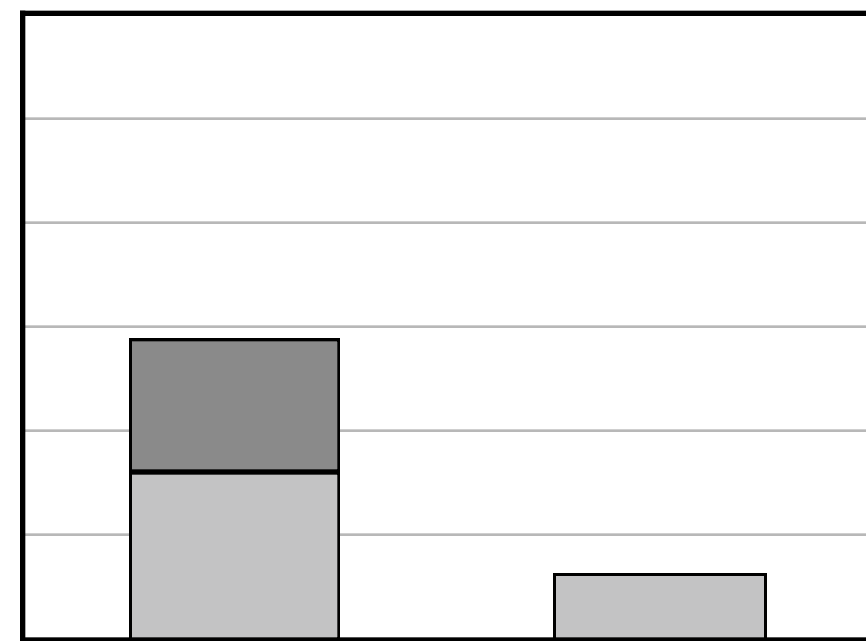
Tile 1
(Layer 1 Input)



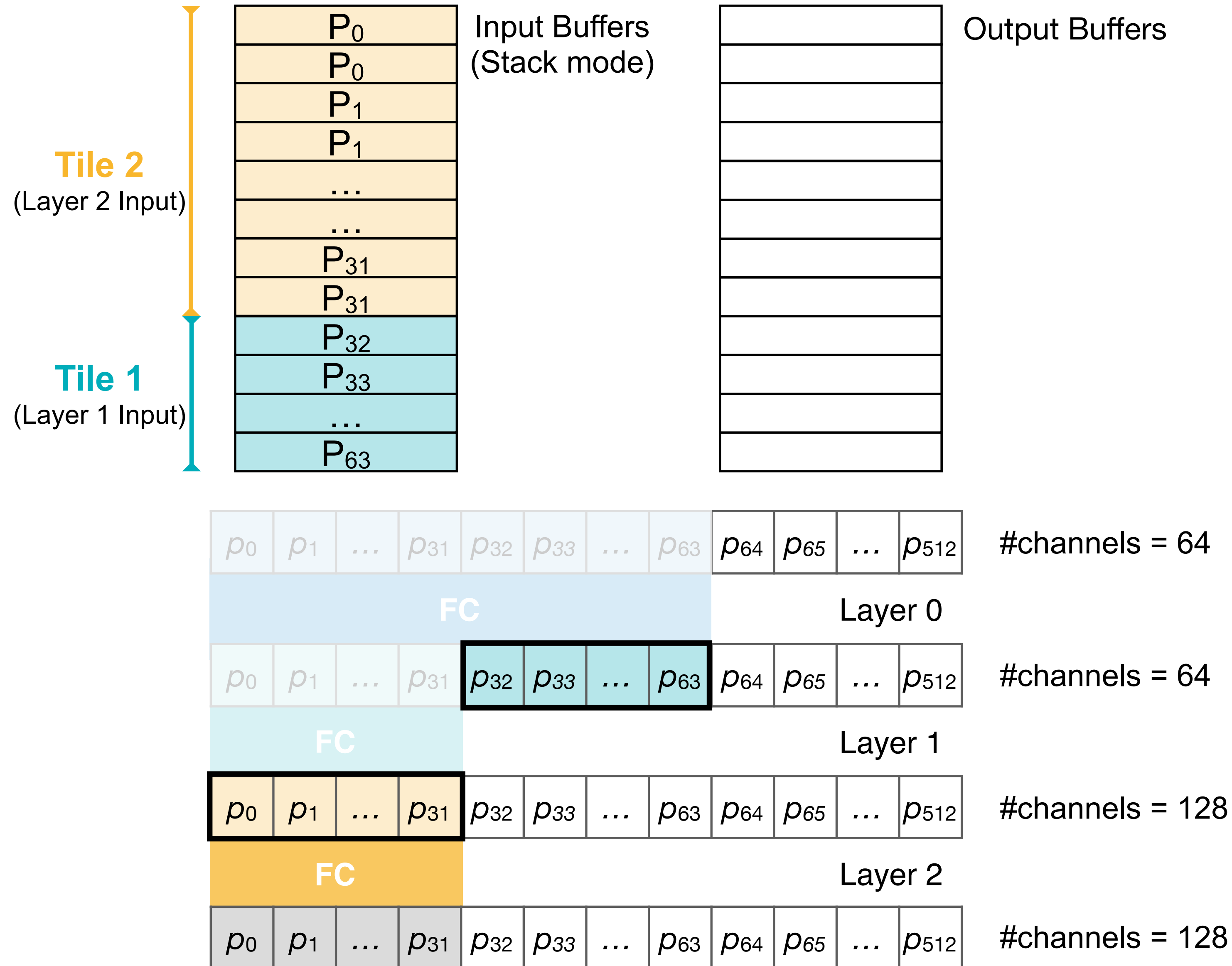
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



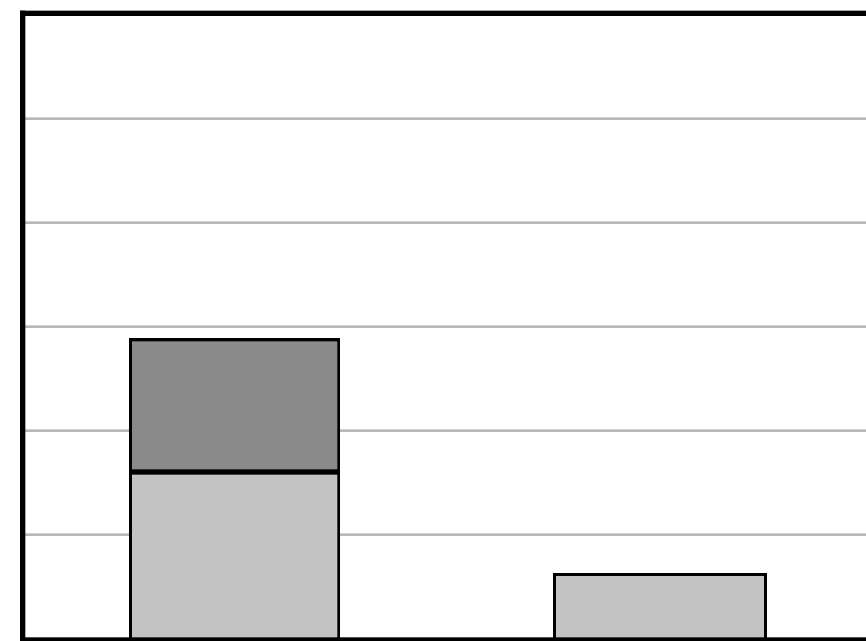
without layer fusion with layer fusion



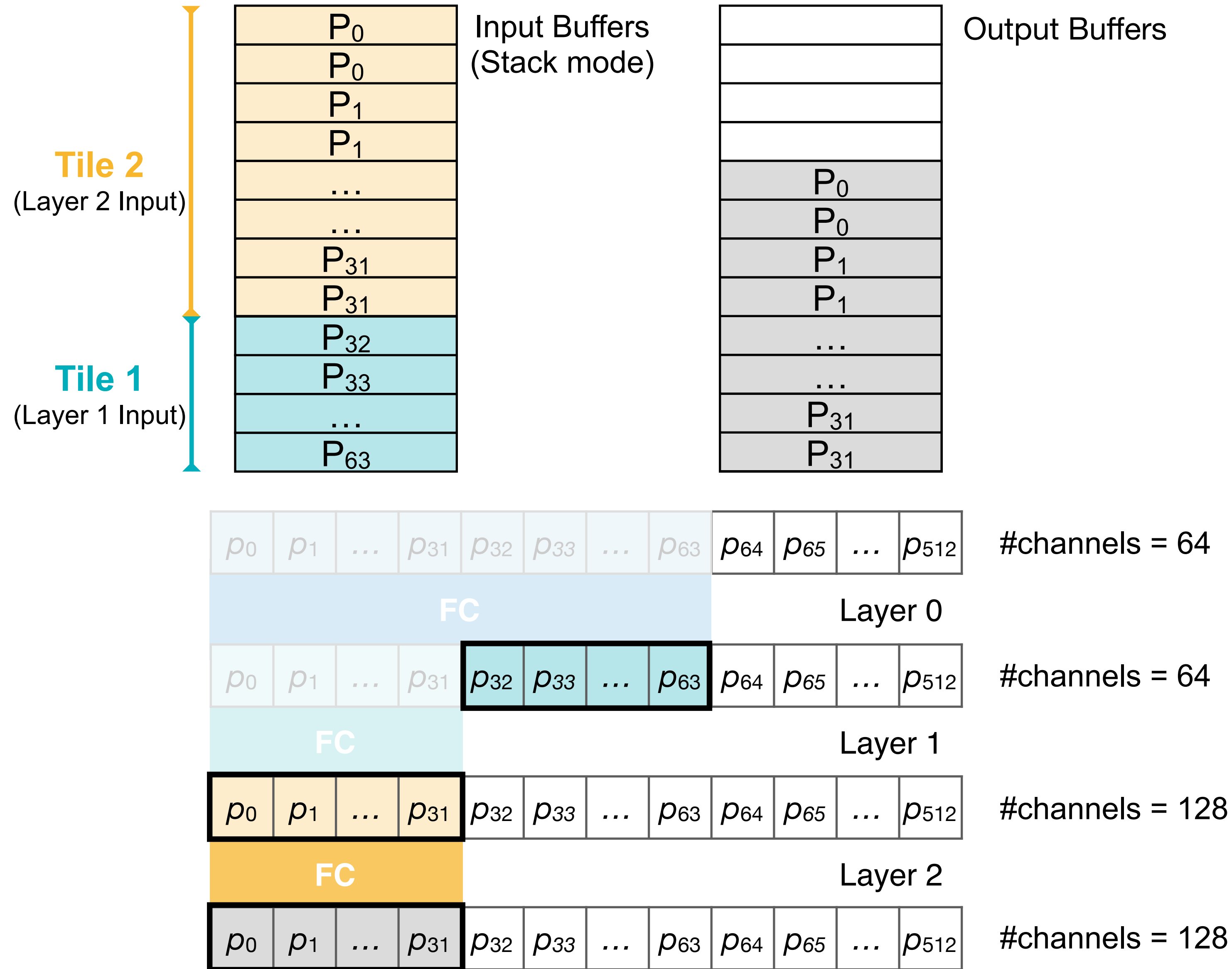
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



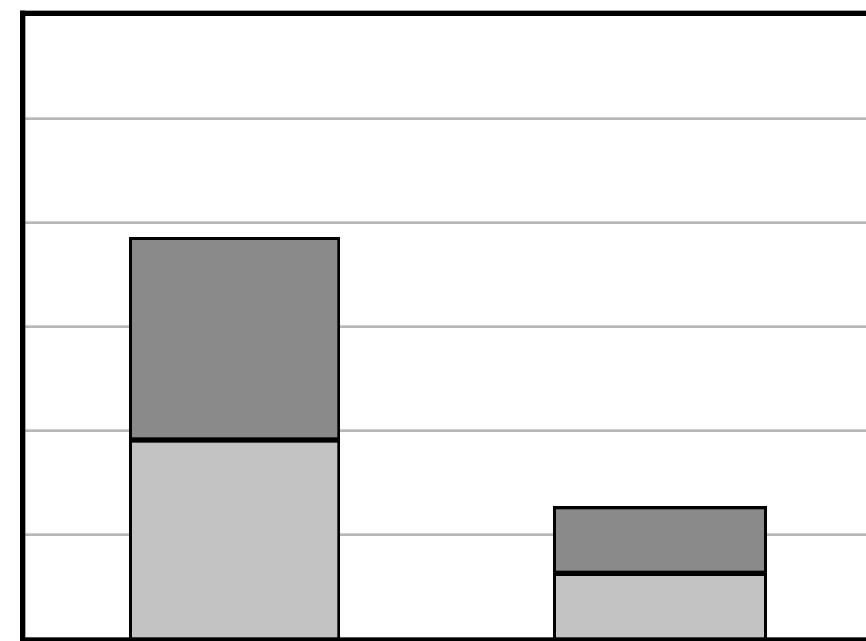
without layer fusion with layer fusion



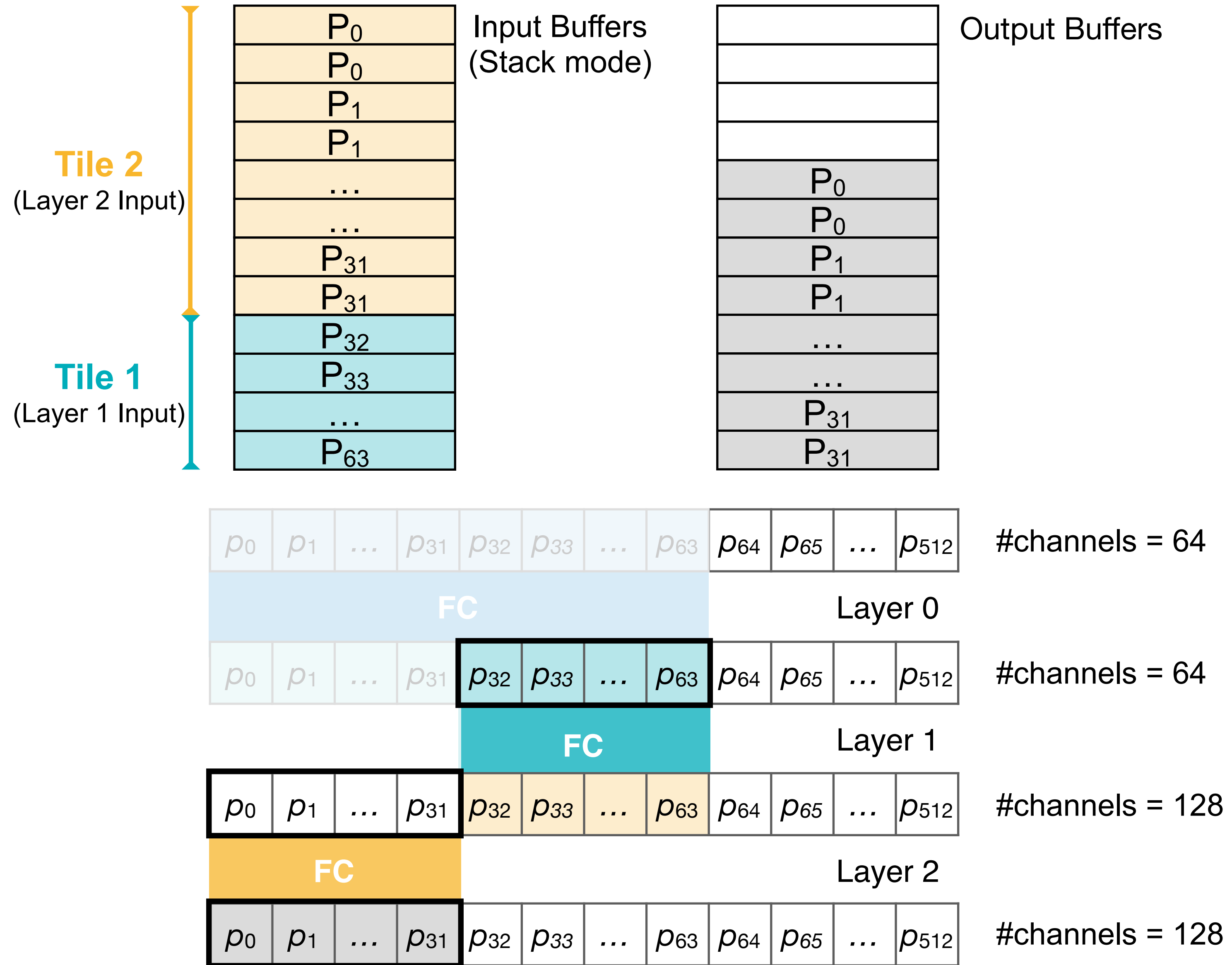
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



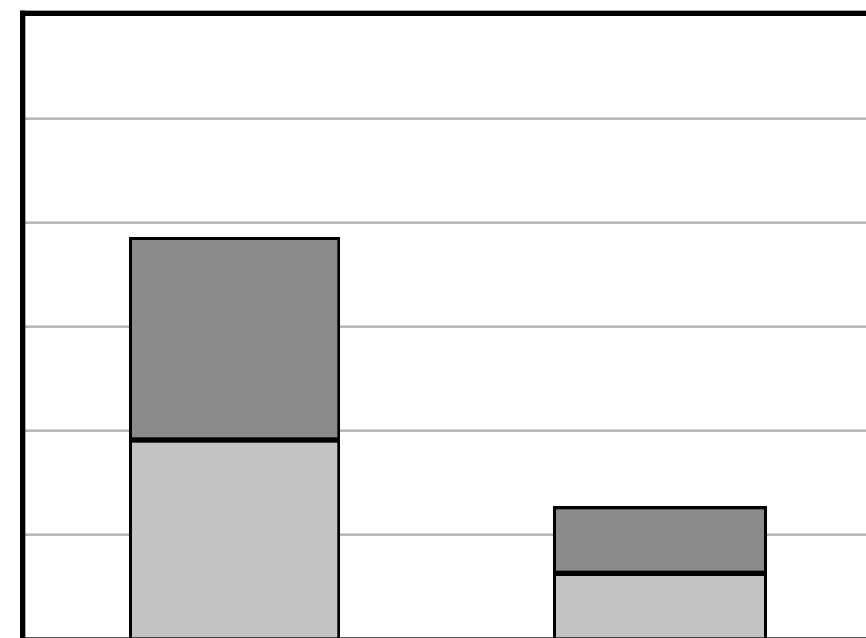
without layer fusion with layer fusion



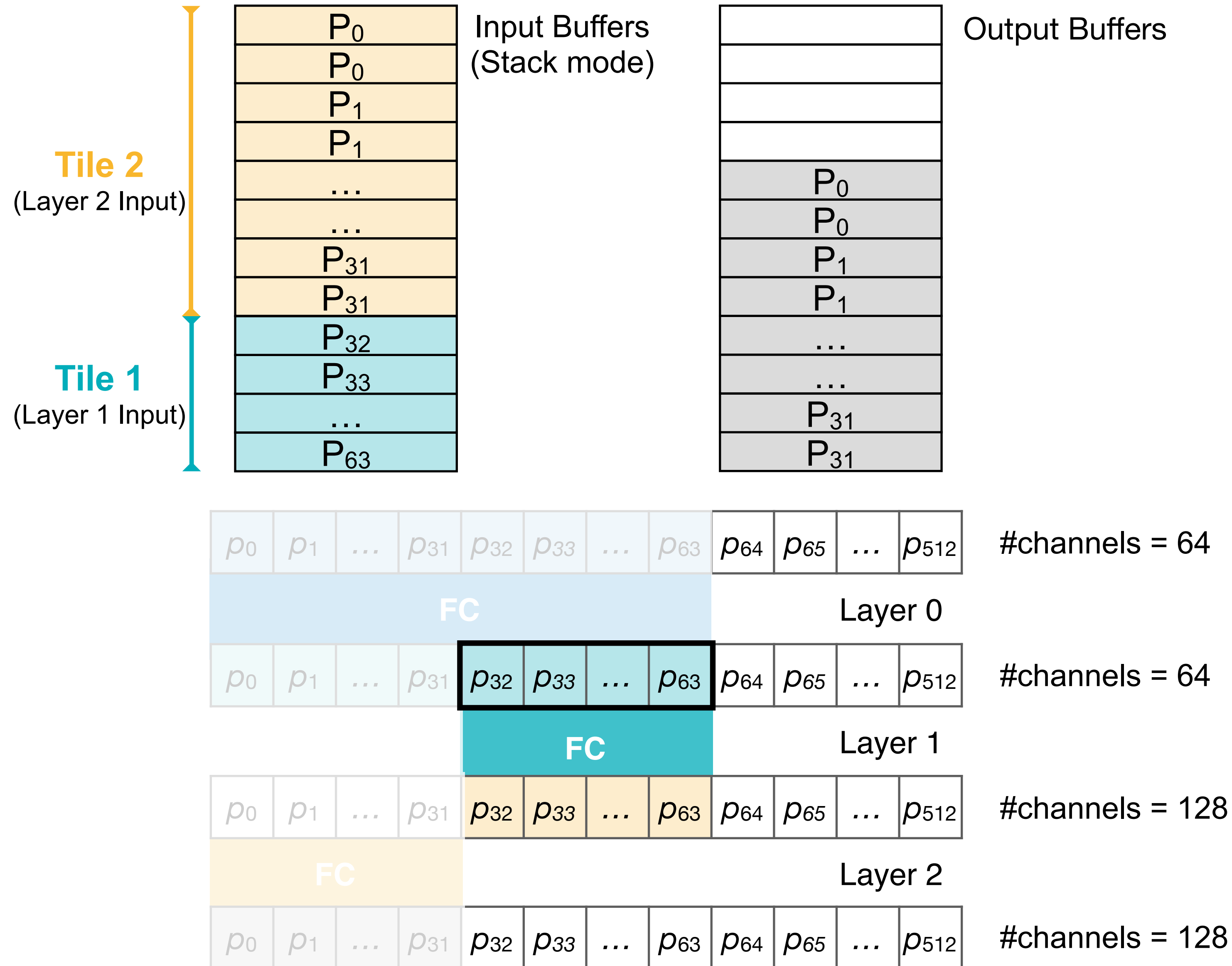
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



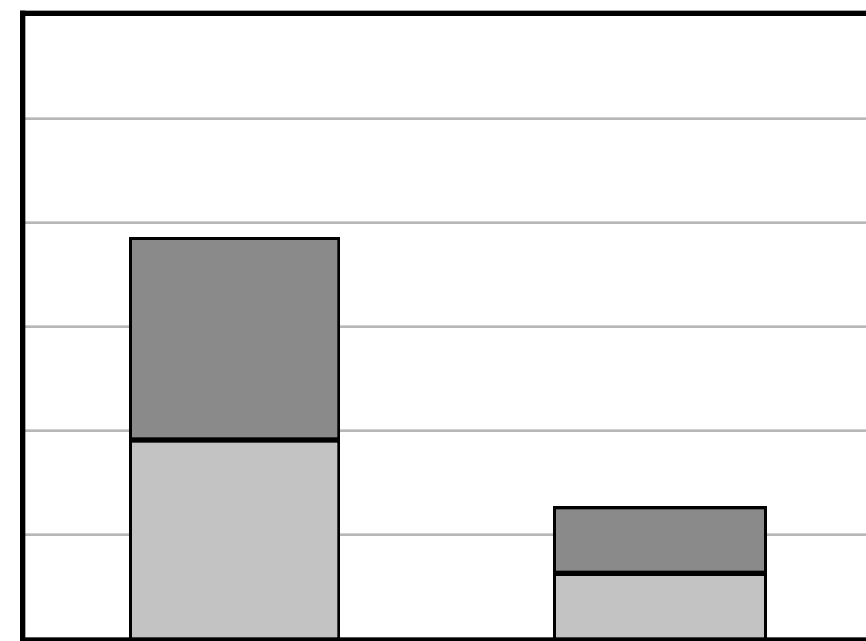
without layer fusion with layer fusion



Temporal Layer Fusion on consecutive FCs

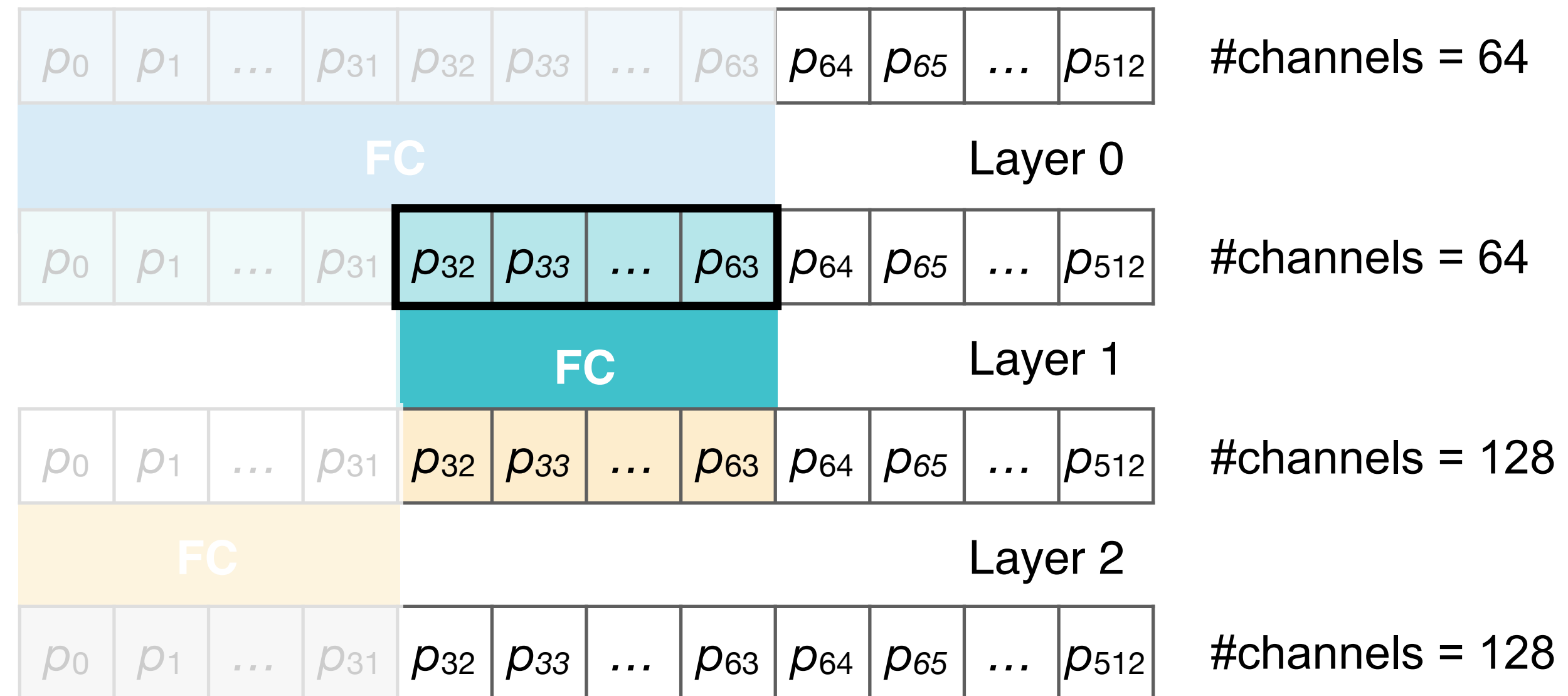
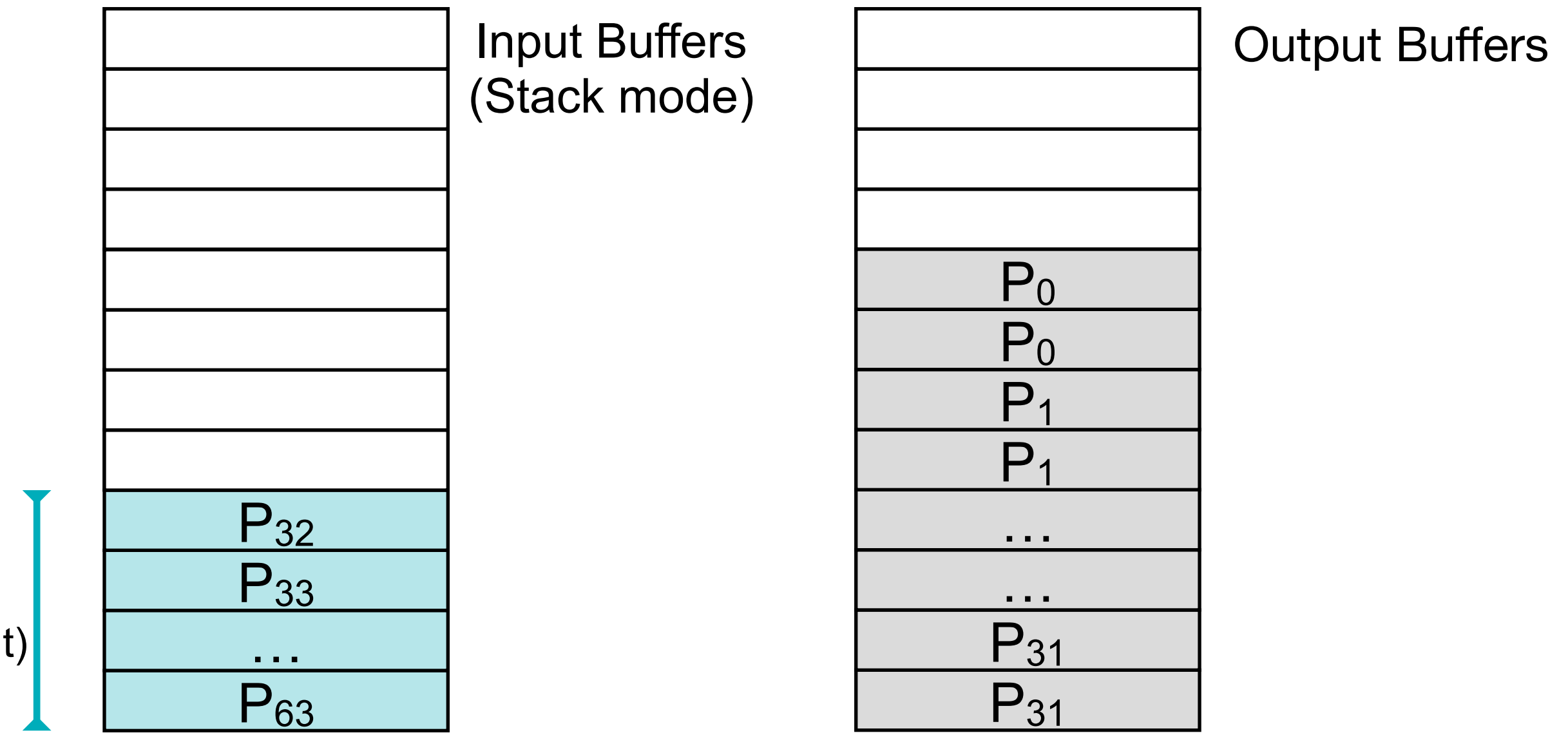
DRAM access per point

read write



without layer fusion with layer fusion

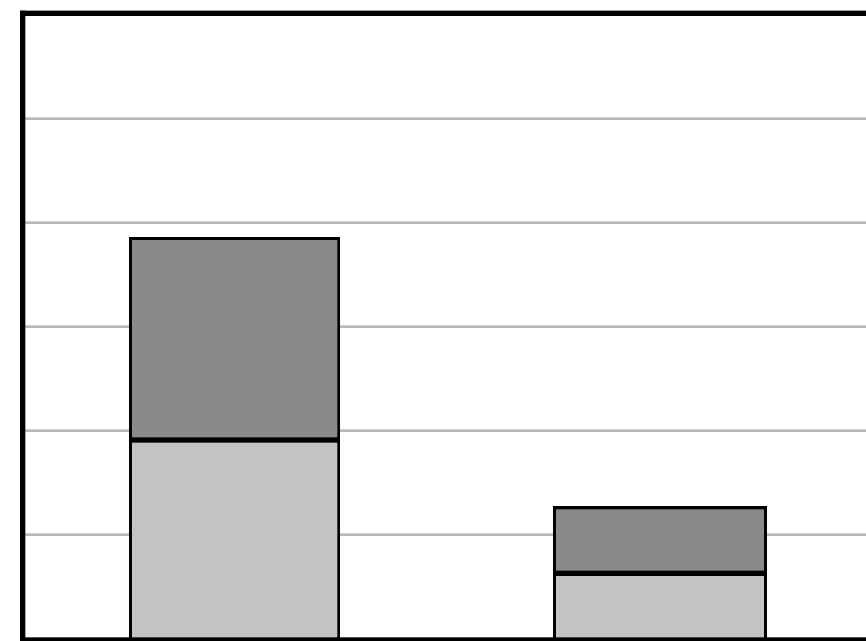
Tile 1
(Layer 1 Input)



Temporal Layer Fusion on consecutive FCs

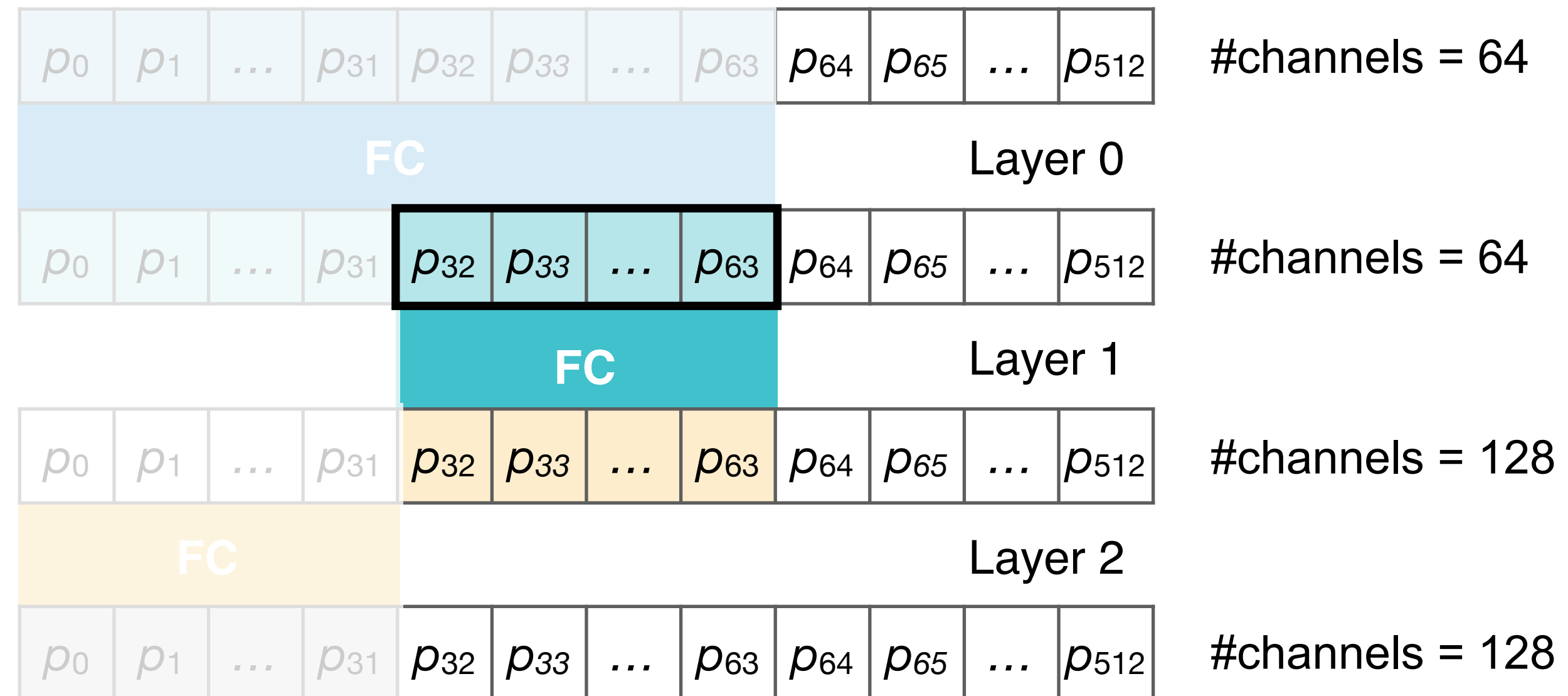
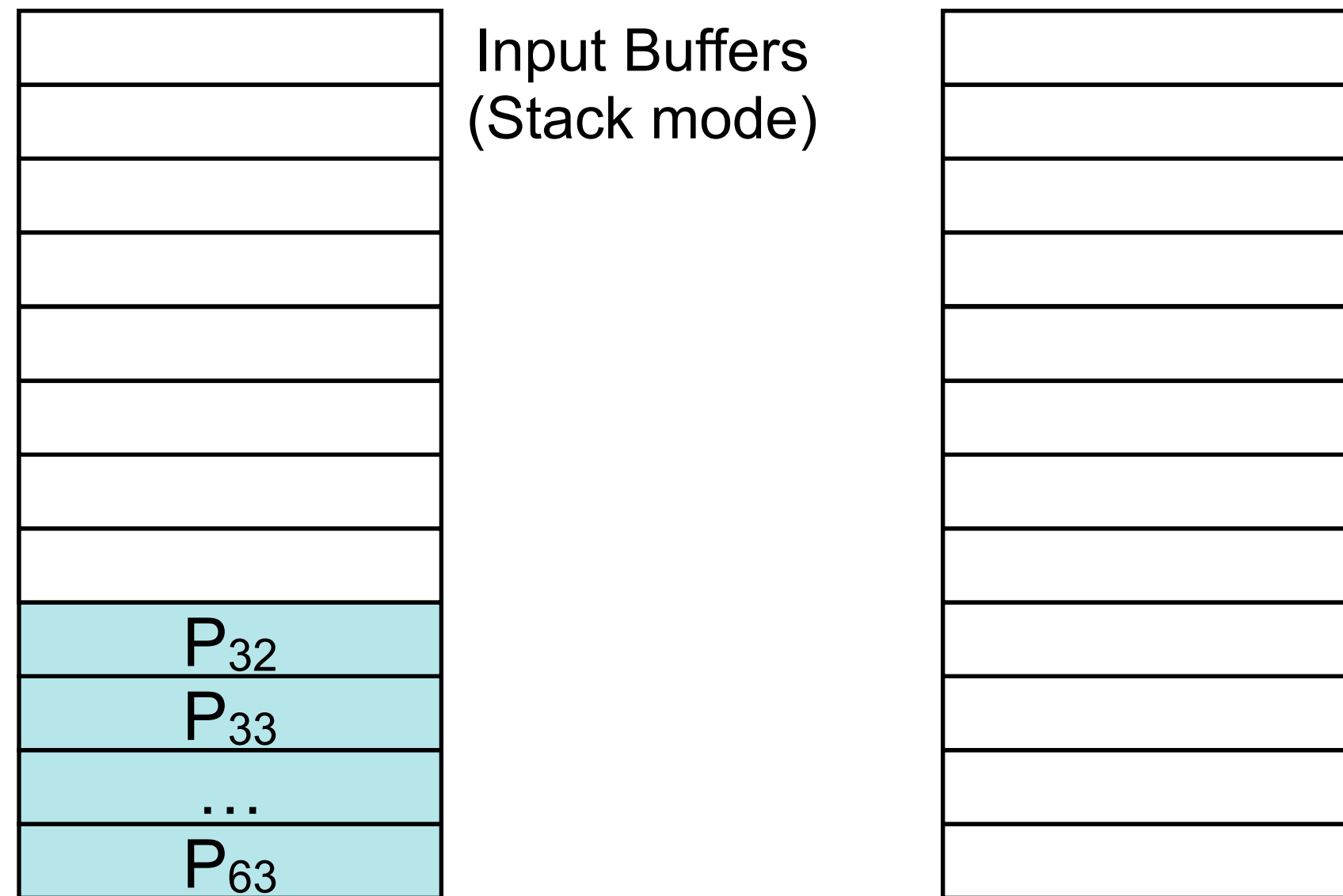
DRAM access per point

read write



without layer fusion with layer fusion

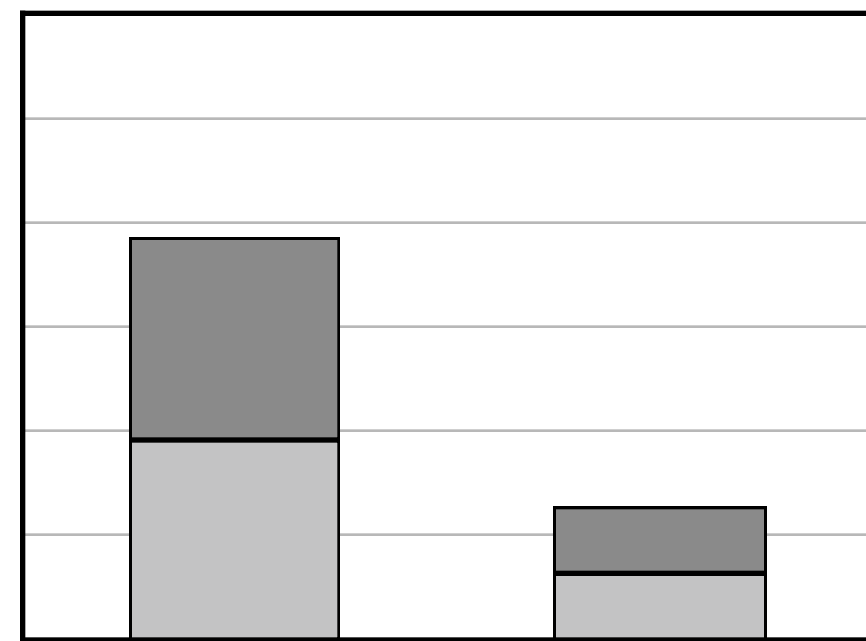
Tile 1
(Layer 1 Input)



Temporal Layer Fusion on consecutive FCs

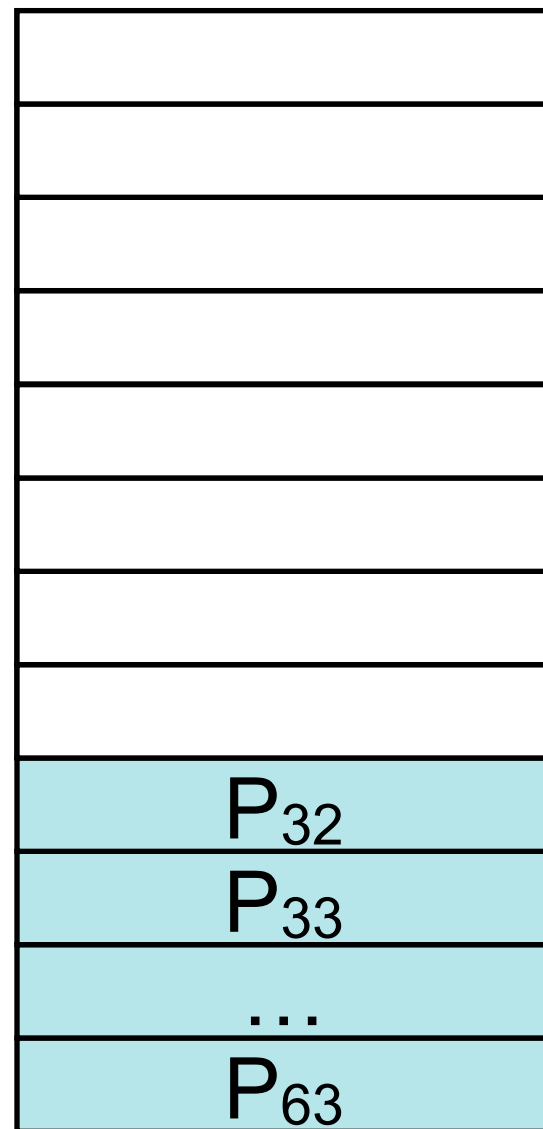
DRAM access per point

read write

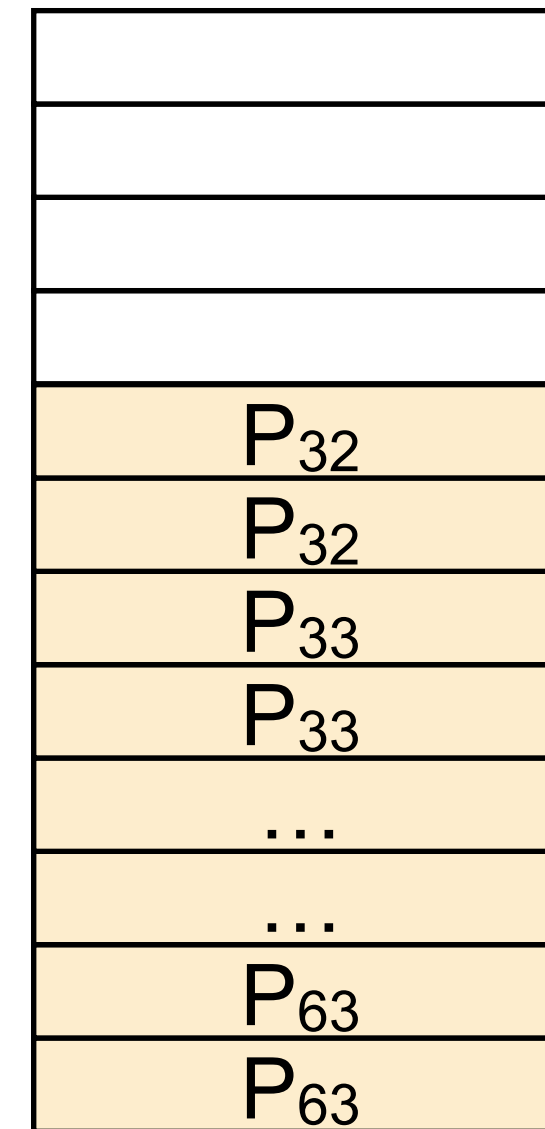


without layer fusion with layer fusion

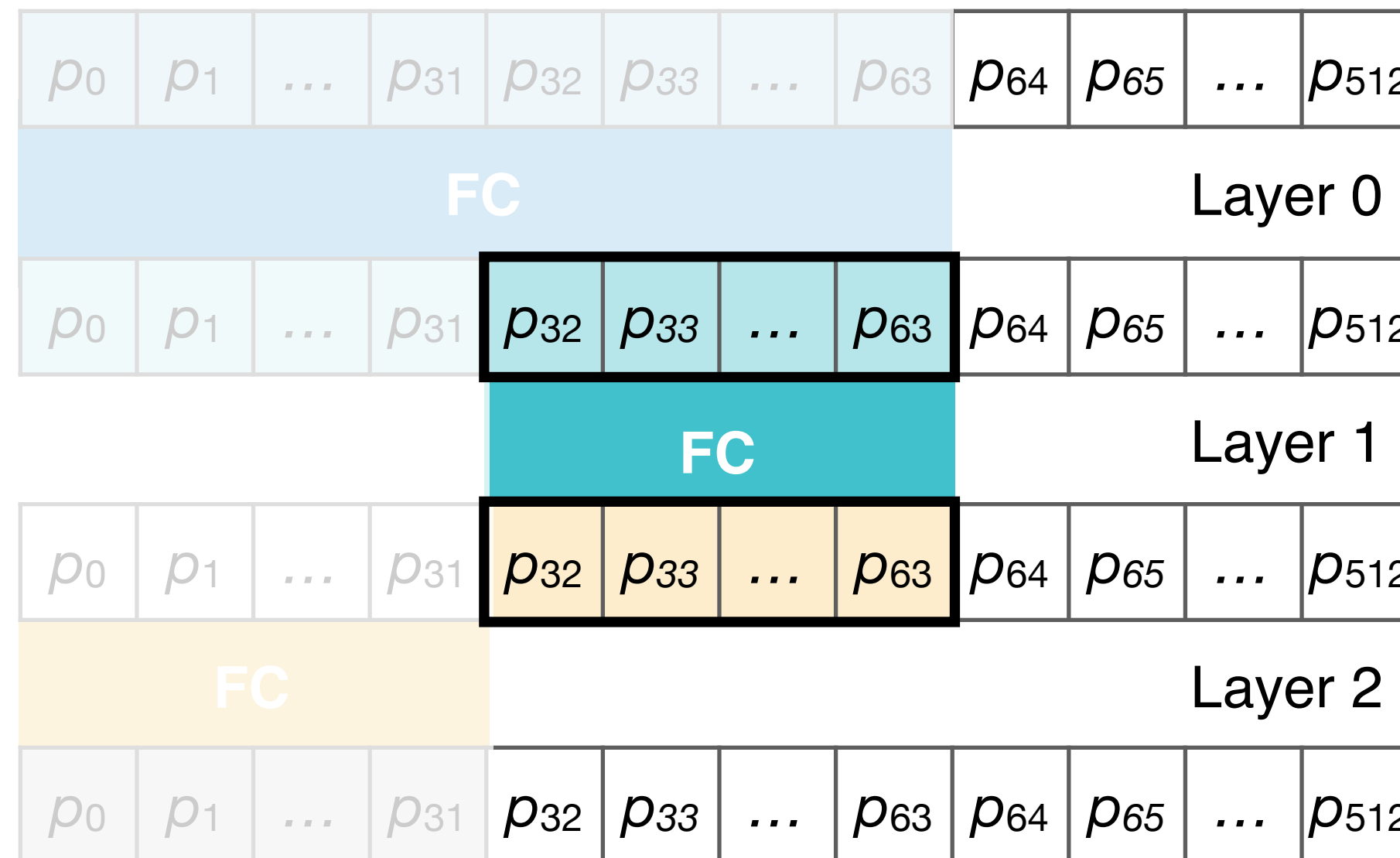
Tile 1
(Layer 1 Input)



Input Buffers
(Stack mode)



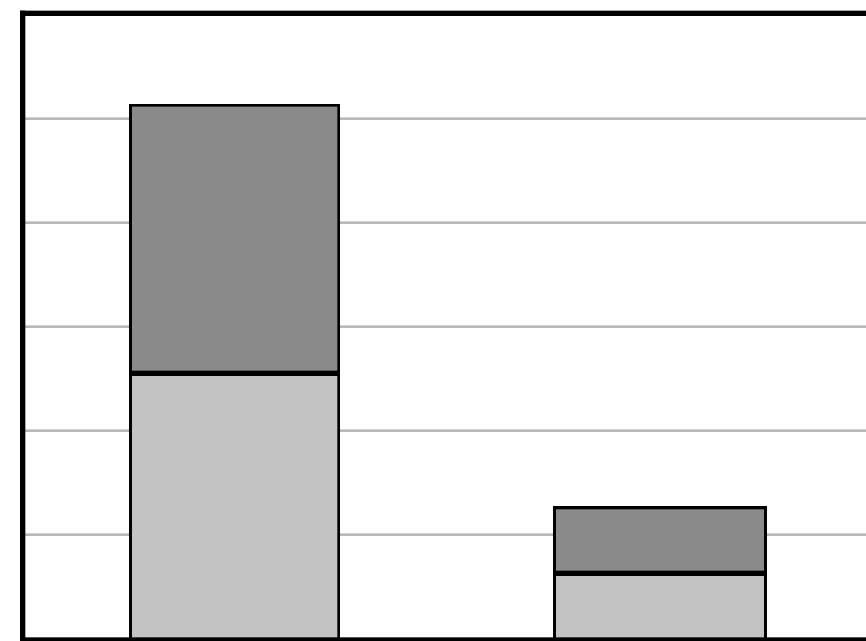
Output Buffers



Temporal Layer Fusion on consecutive FCs

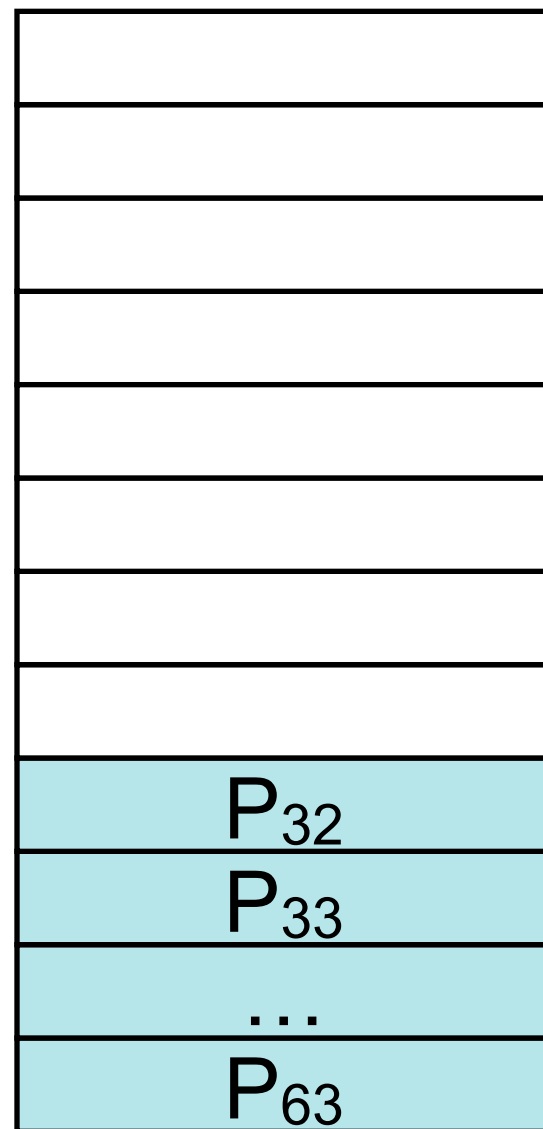
DRAM access per point

read write

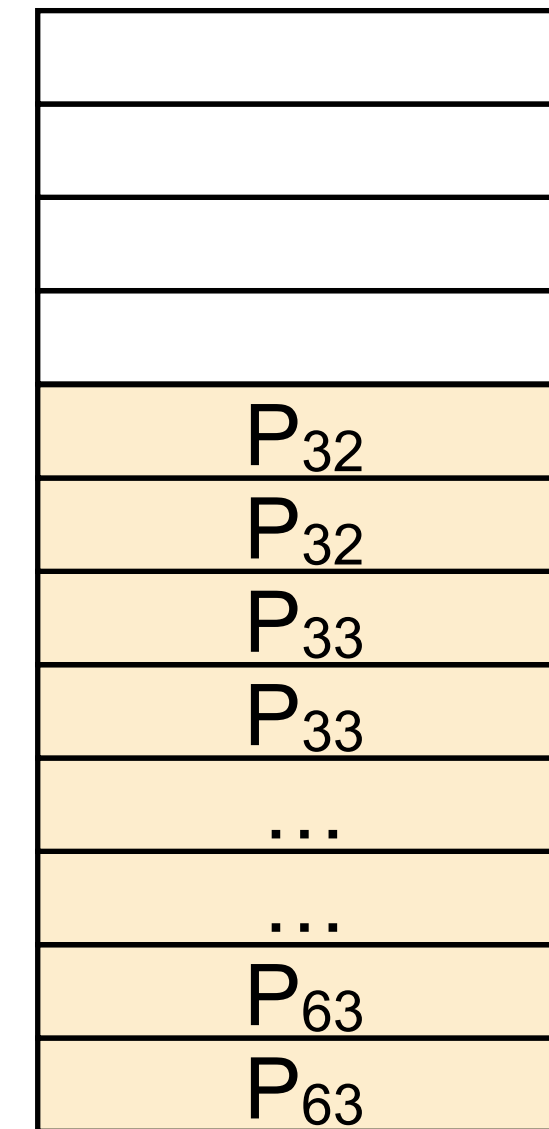


without layer fusion with layer fusion

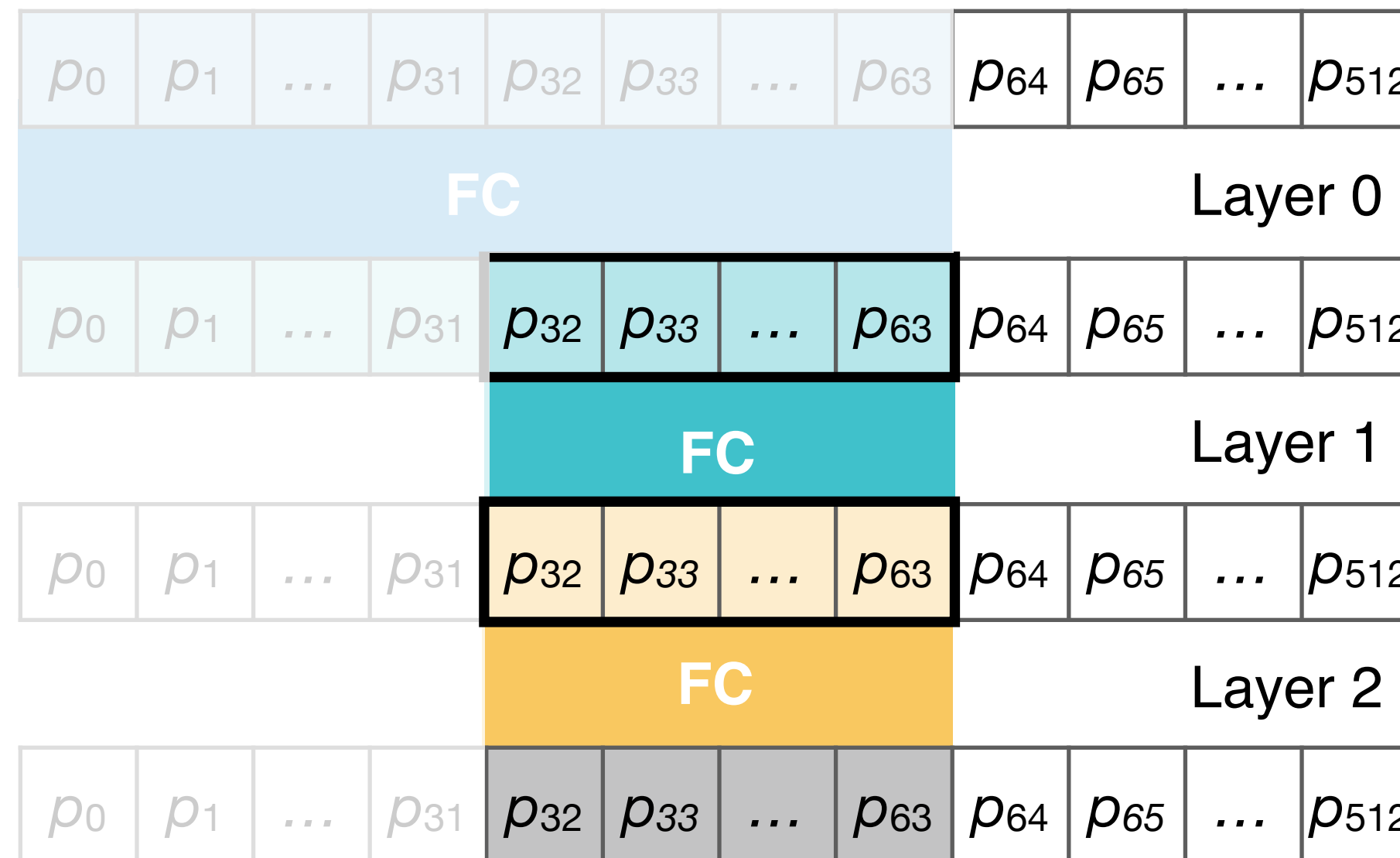
Tile 1
(Layer 1 Input)



Input Buffers
(Stack mode)



Output Buffers



#channels = 64

#channels = 64

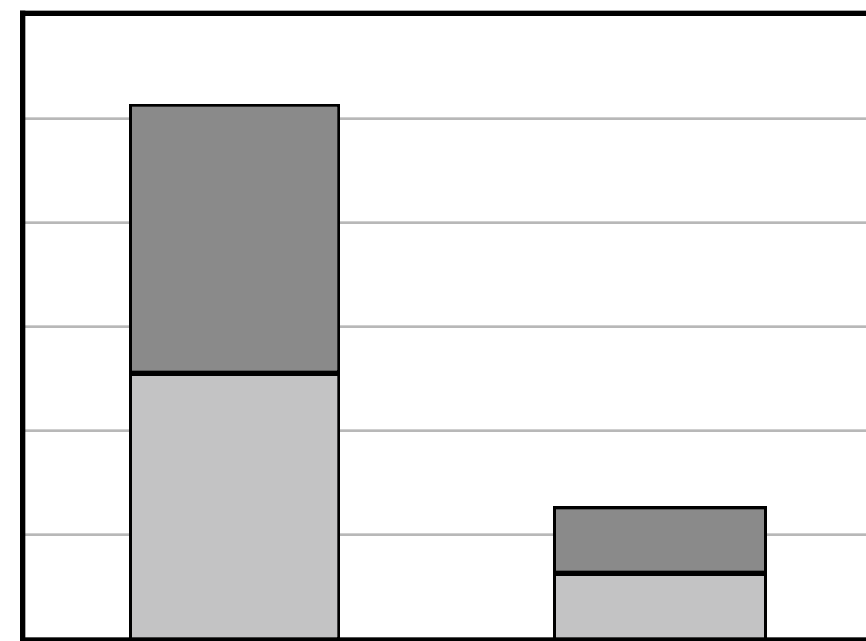
#channels = 128

#channels = 128

Temporal Layer Fusion on consecutive FCs

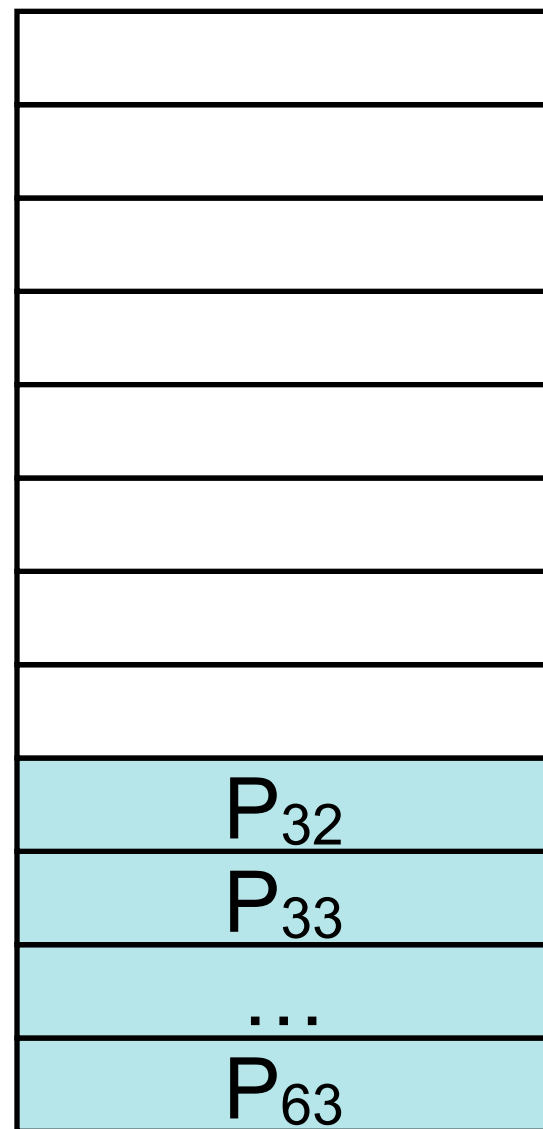
DRAM access per point

read write

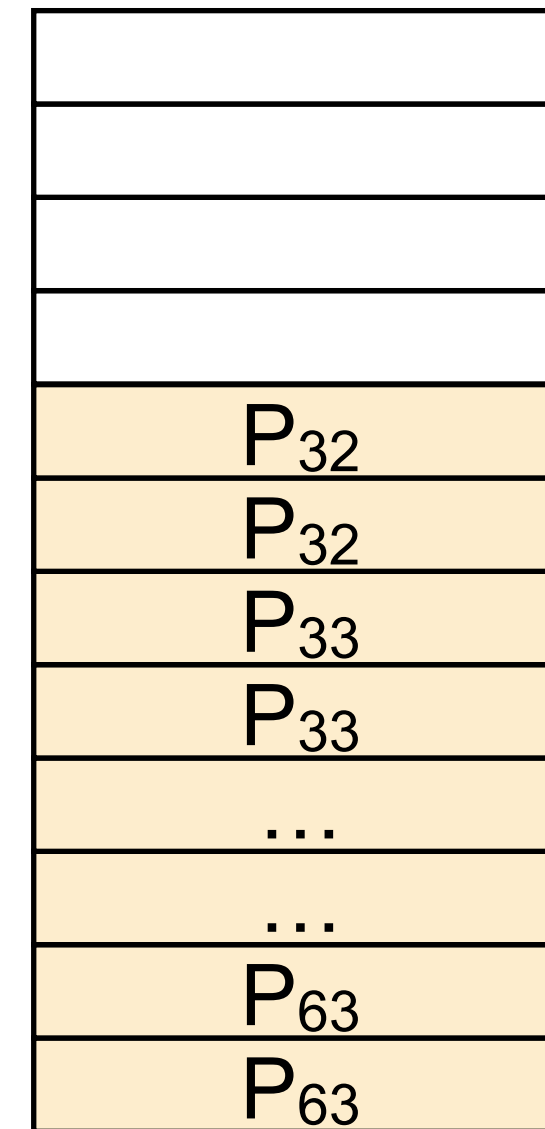


without layer fusion with layer fusion

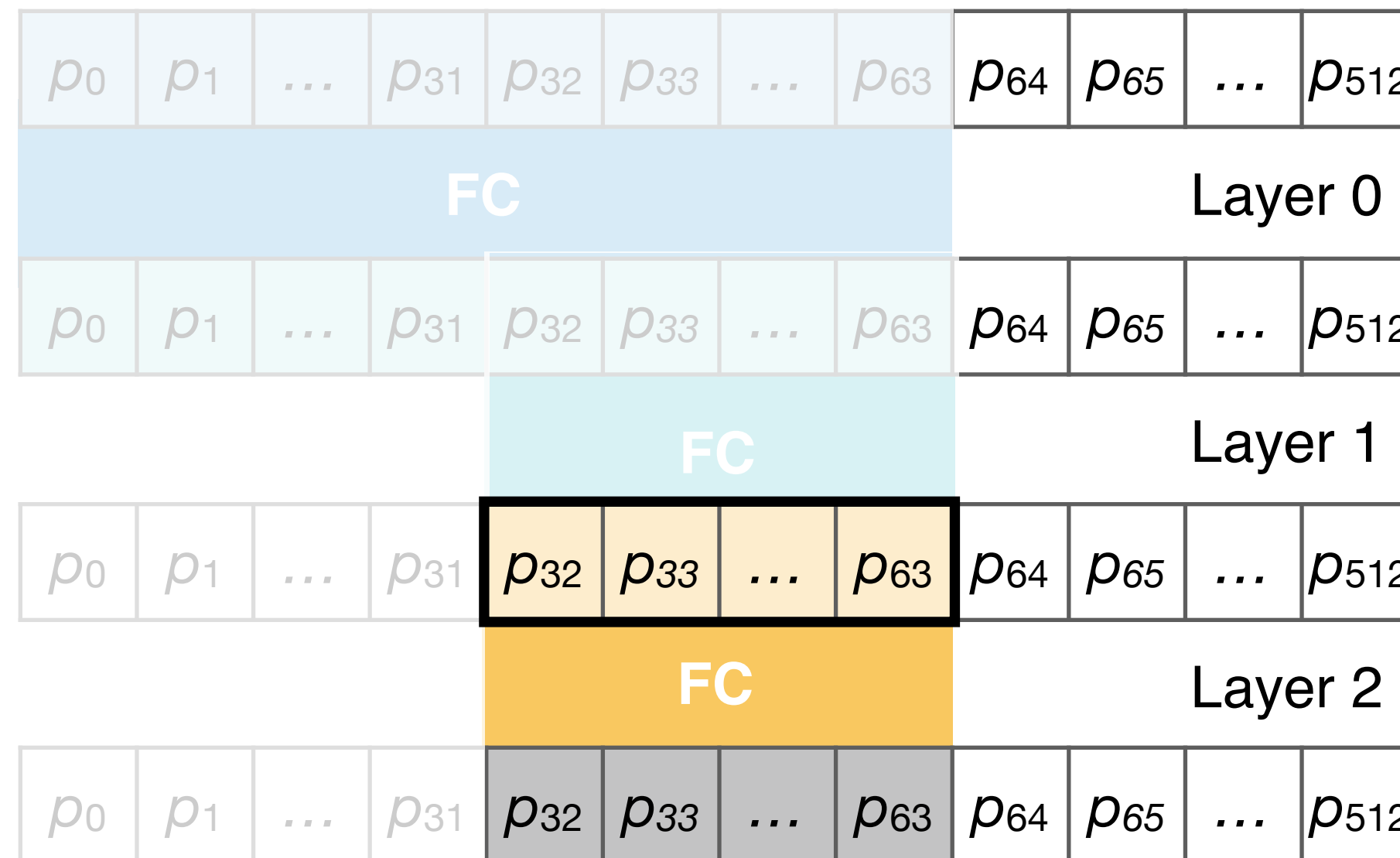
Tile 1
(Layer 1 Input)



Input Buffers
(Stack mode)



Output Buffers



#channels = 64

#channels = 64

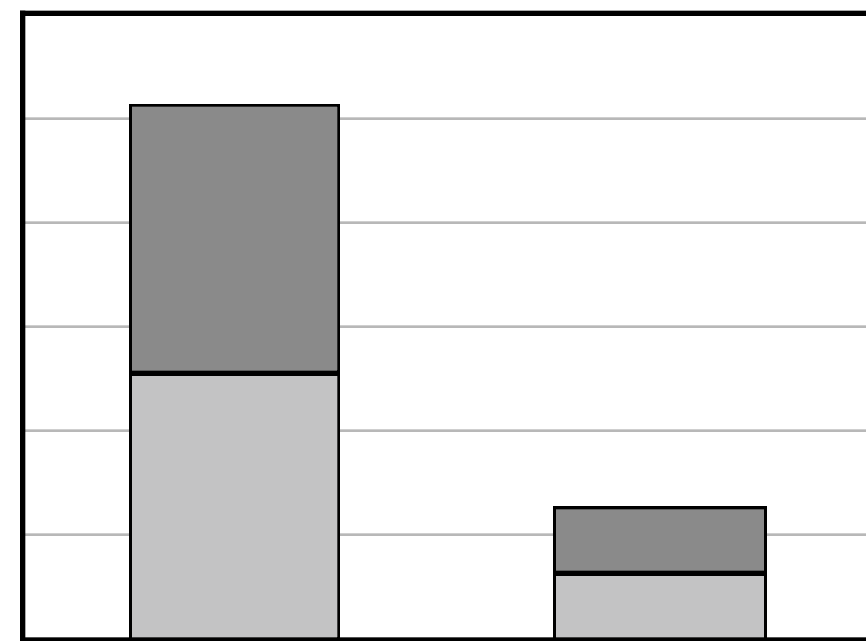
#channels = 128

#channels = 128

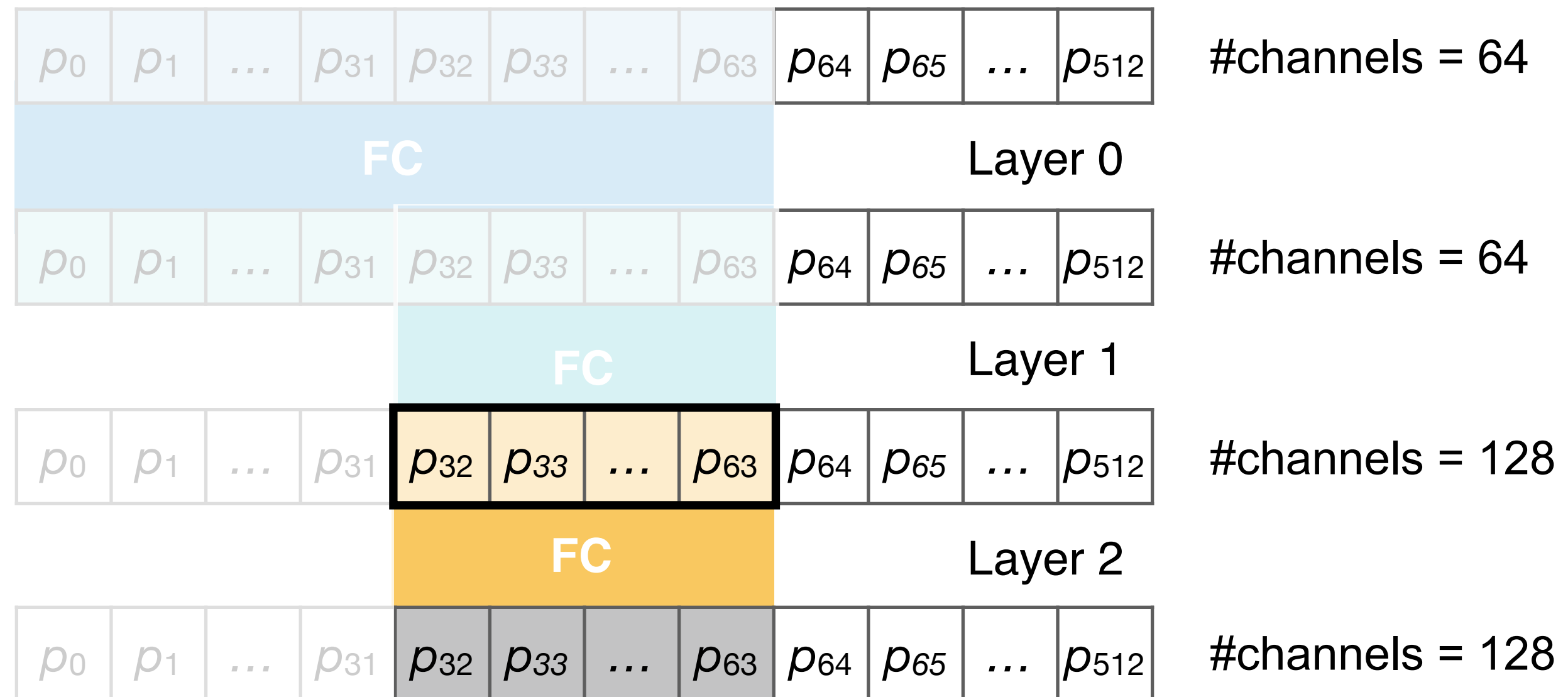
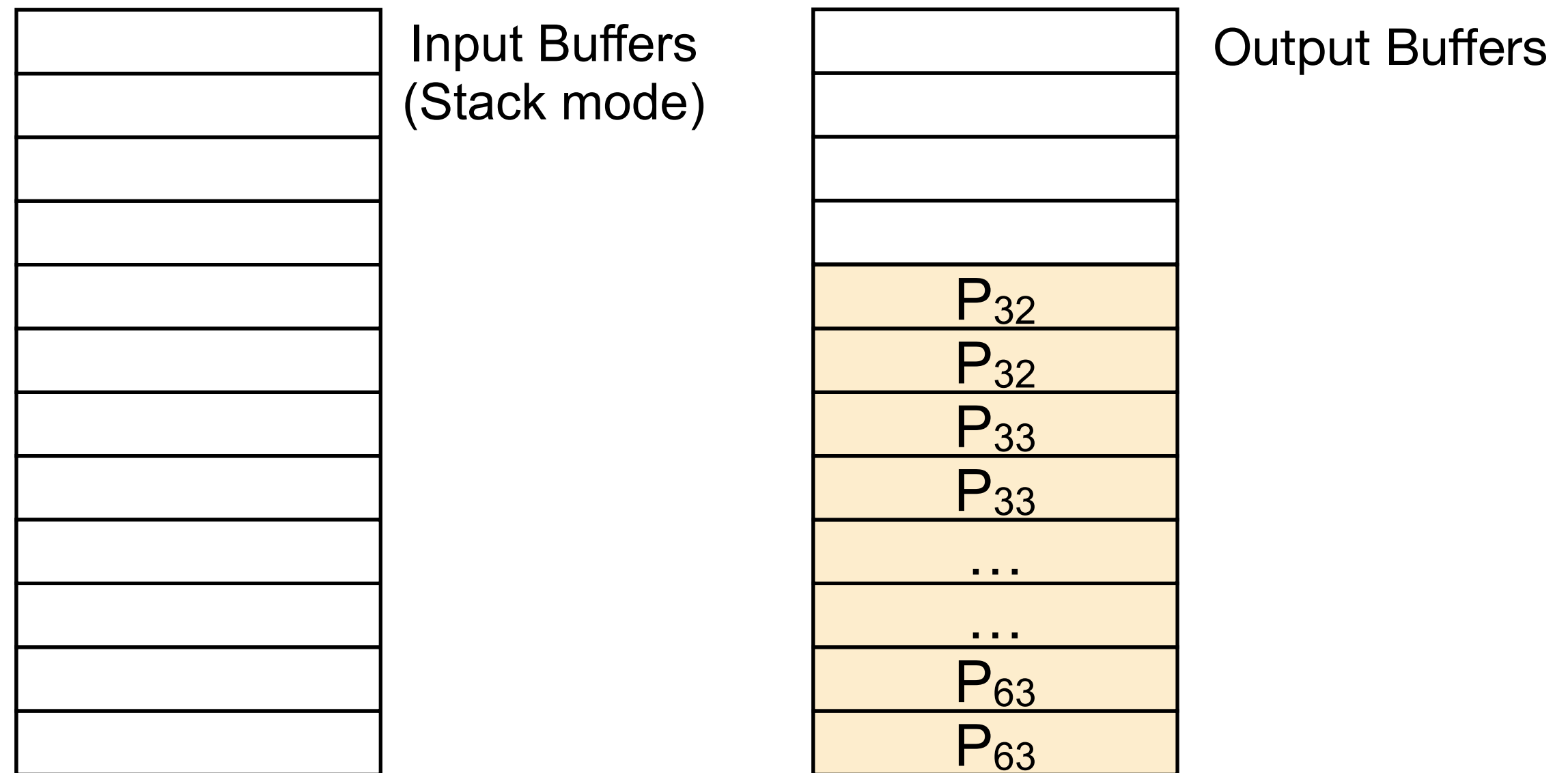
Temporal Layer Fusion on consecutive FCs

DRAM access per point

read write



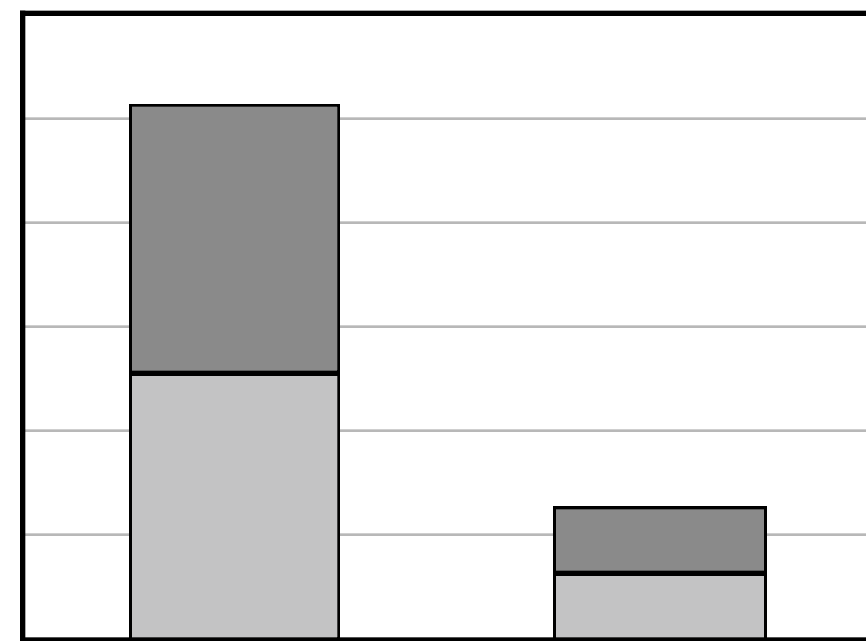
without layer fusion with layer fusion



Temporal Layer Fusion on consecutive FCs

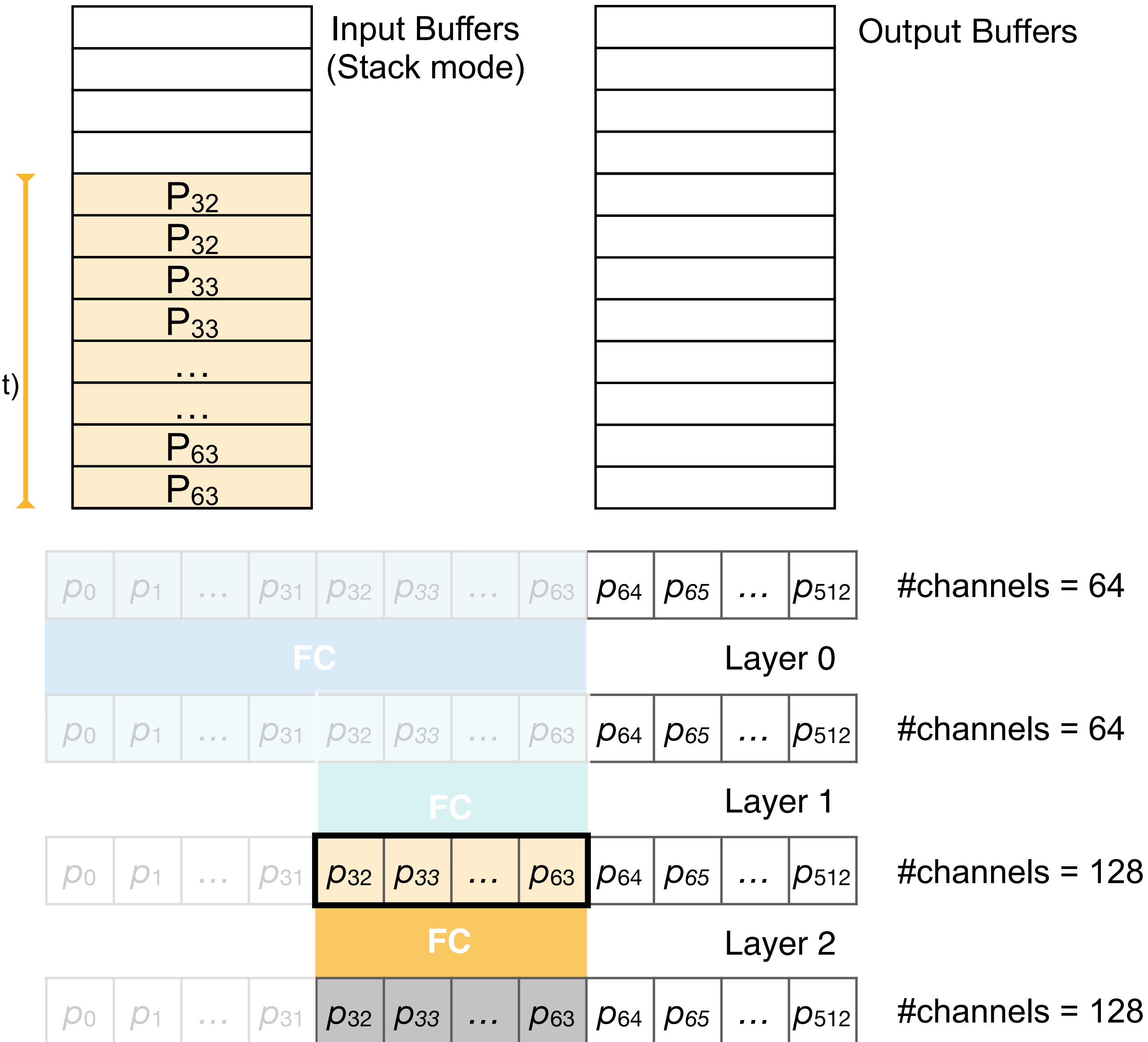
DRAM access per point

read write



without layer fusion with layer fusion

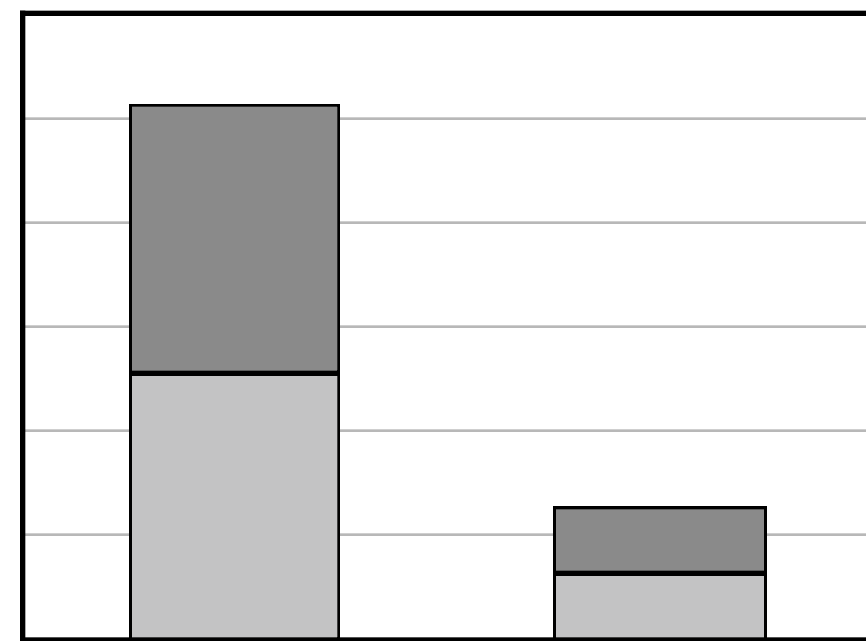
Tile 2
(Layer 2 Input)



Temporal Layer Fusion on consecutive FCs

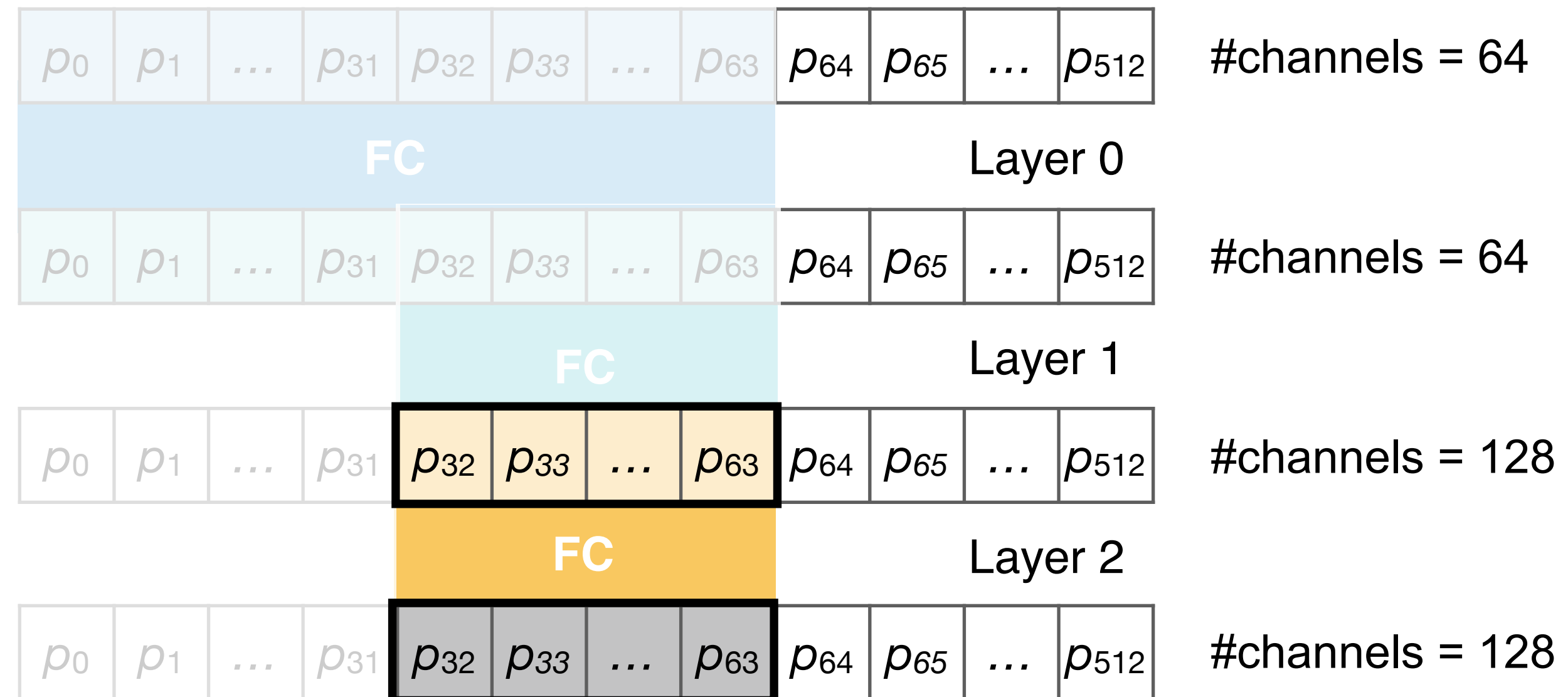
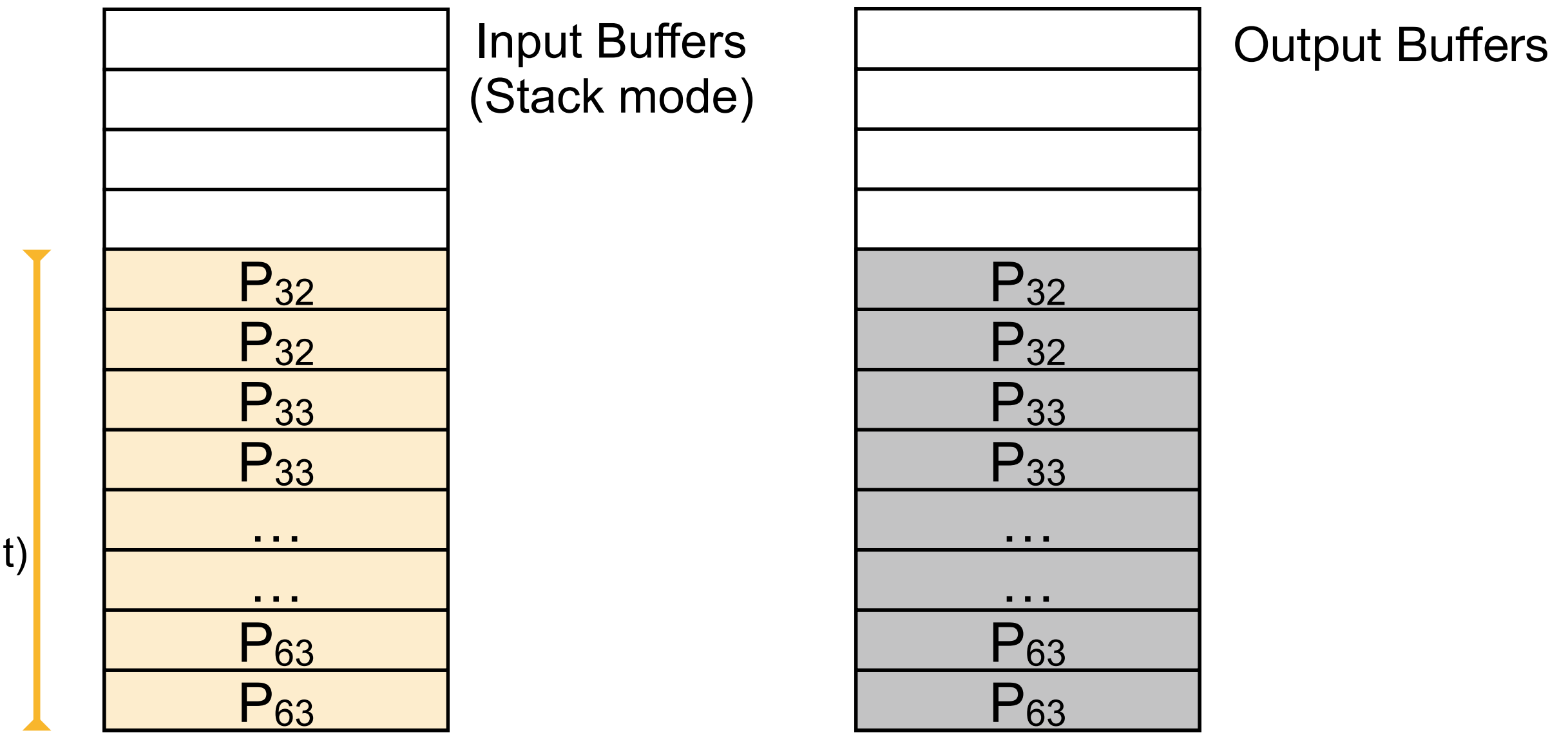
DRAM access per point

read write



without layer fusion with layer fusion

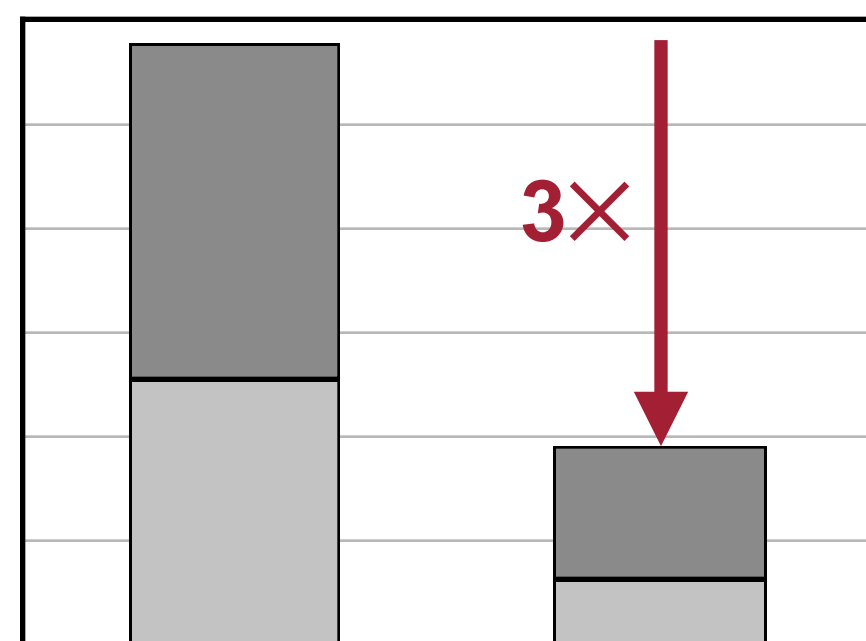
Tile 2
(Layer 2 Input)



Temporal Layer Fusion on consecutive FCs

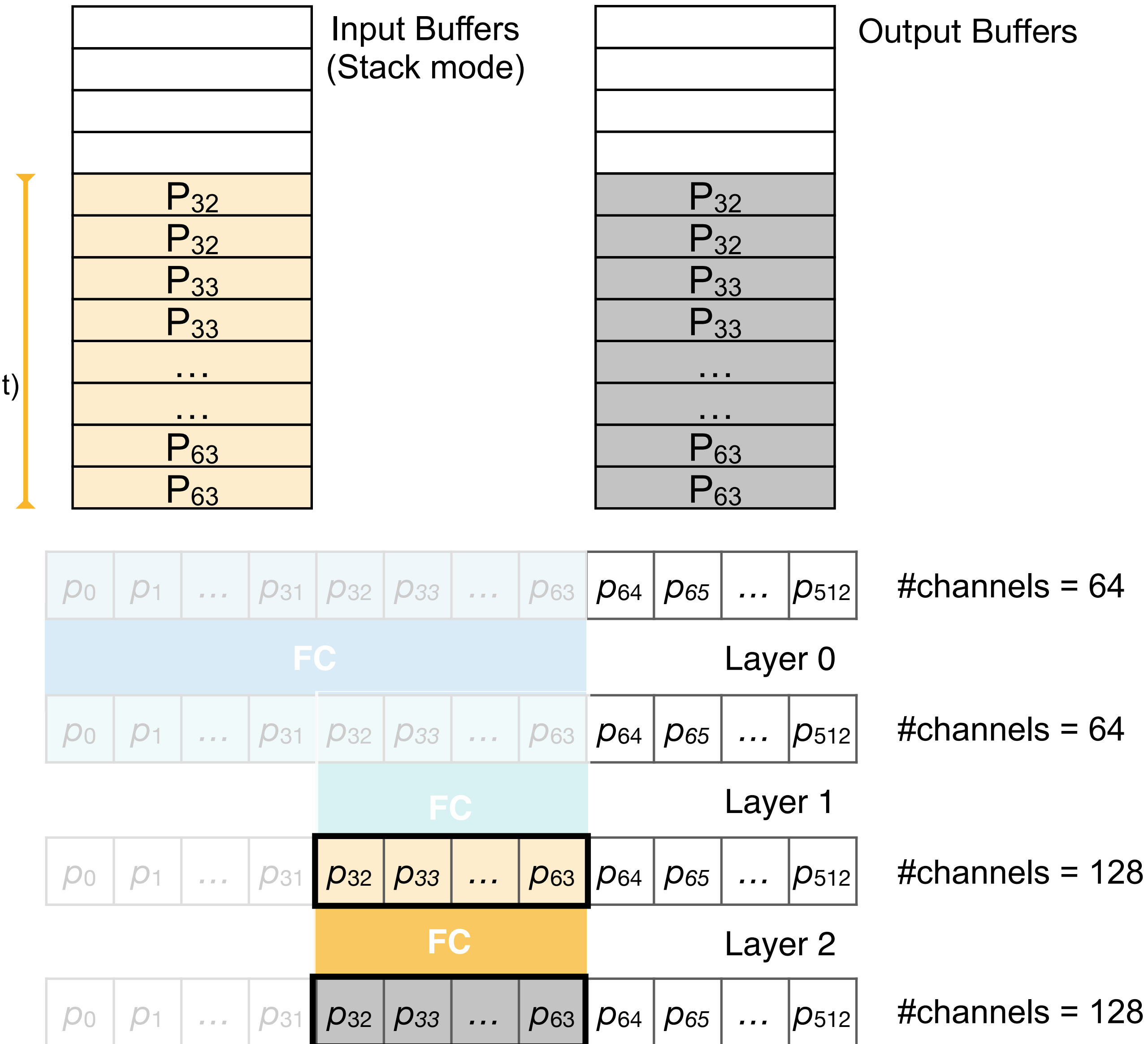
DRAM access per point

read write



without layer fusion with layer fusion

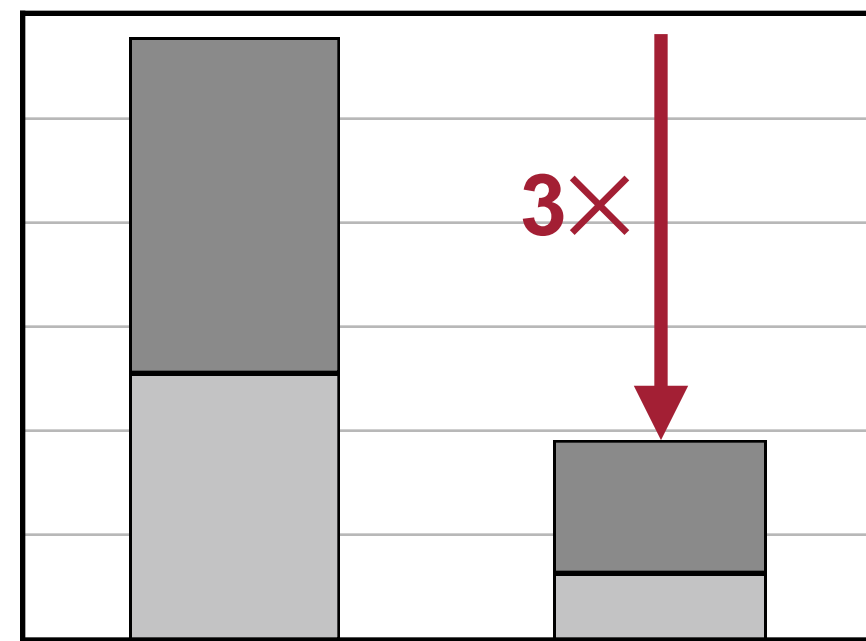
Tile 2
(Layer 2 Input)



Temporal Layer Fusion on consecutive FCs

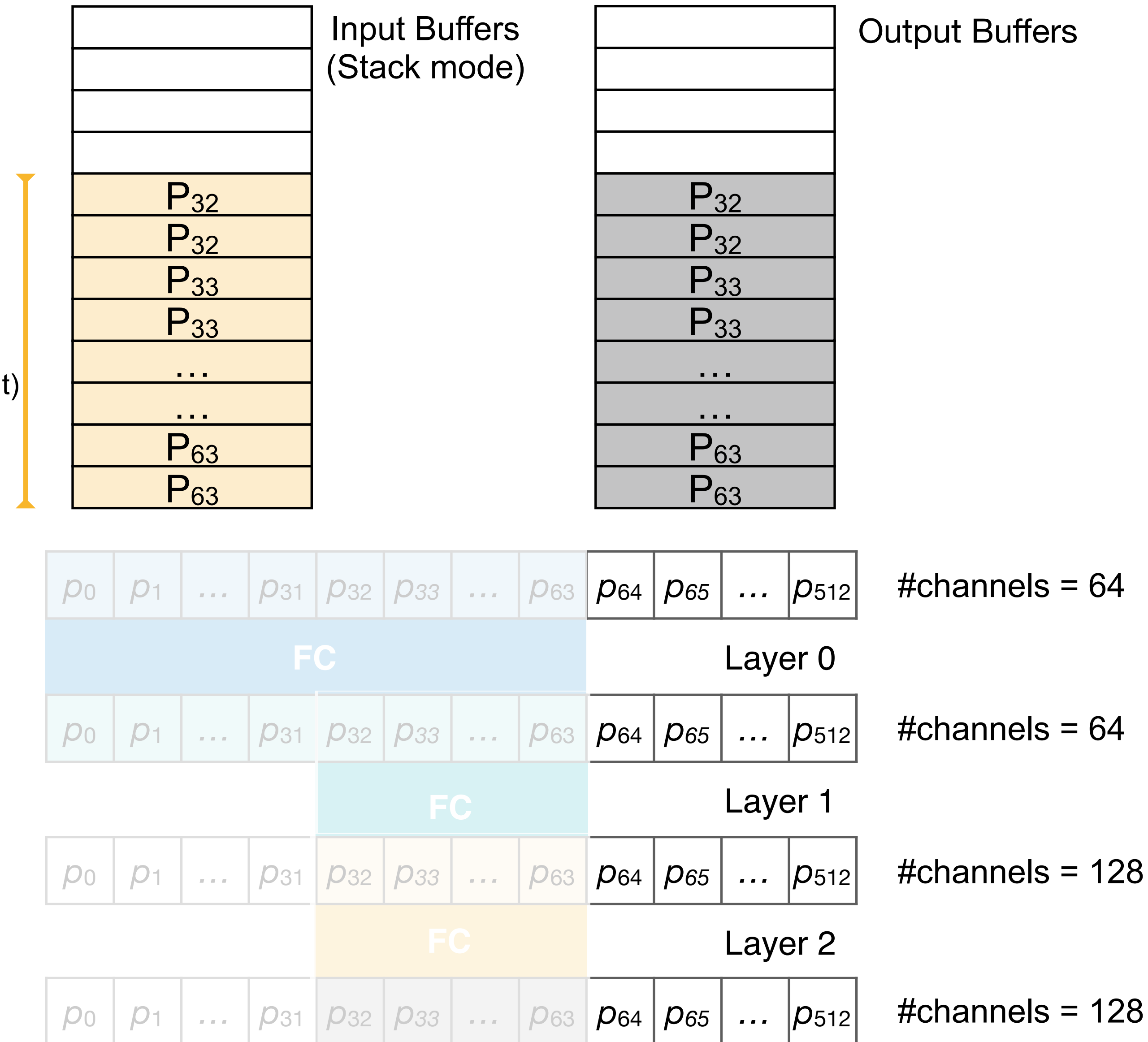
DRAM access per point

read write



without layer fusion with layer fusion

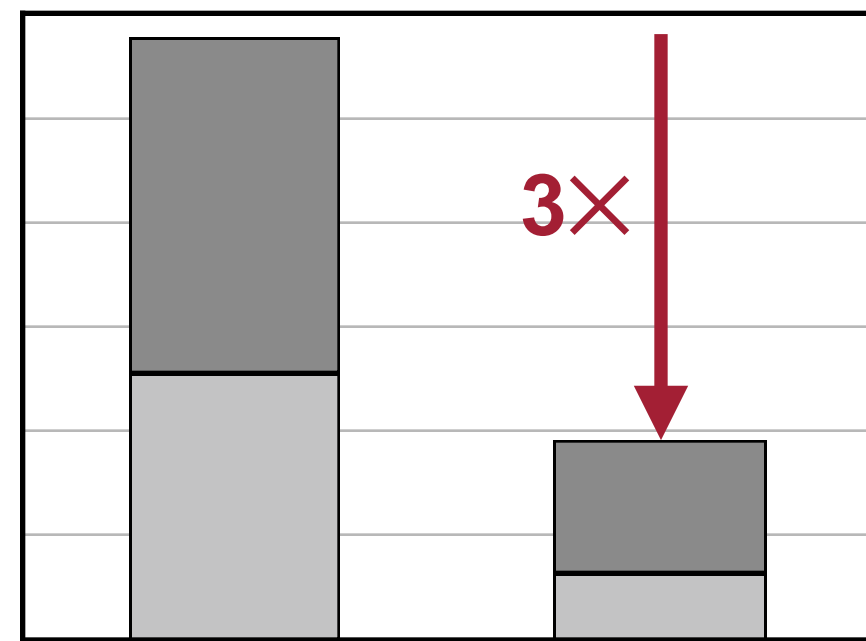
Tile 2
(Layer 2 Input)



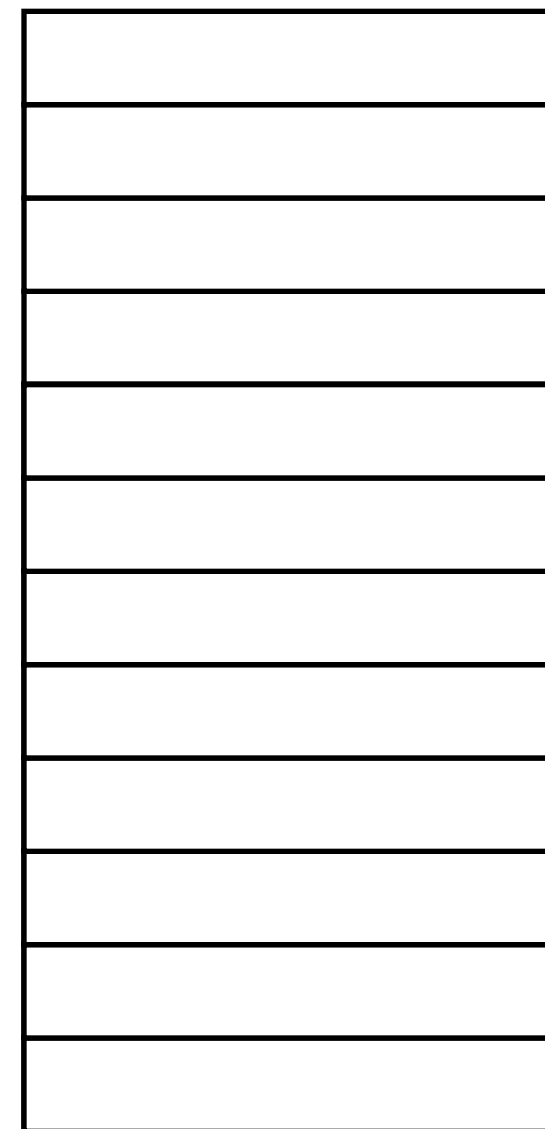
Temporal Layer Fusion on consecutive FCs

DRAM access per point

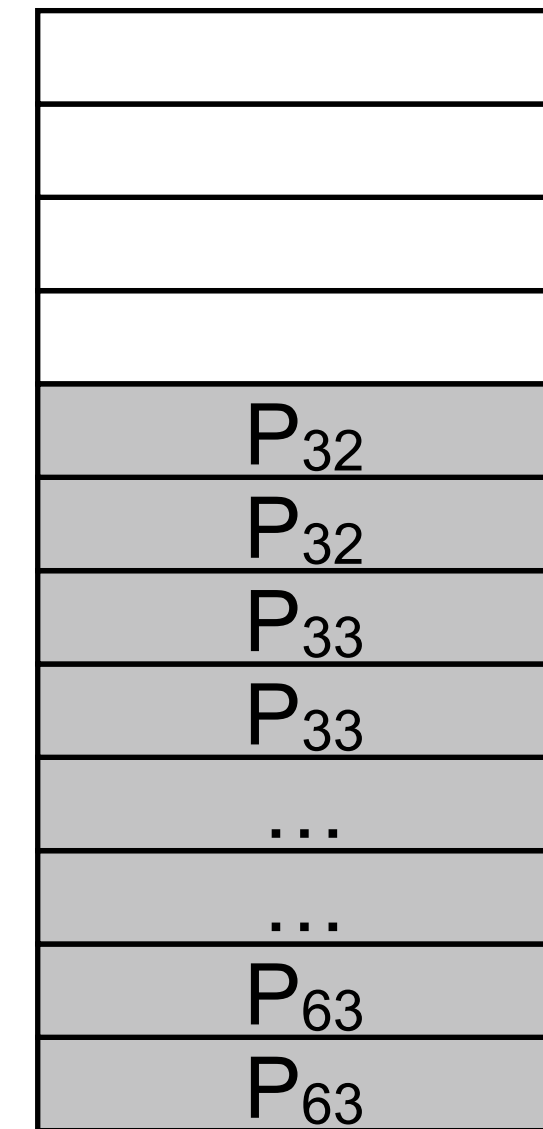
read write



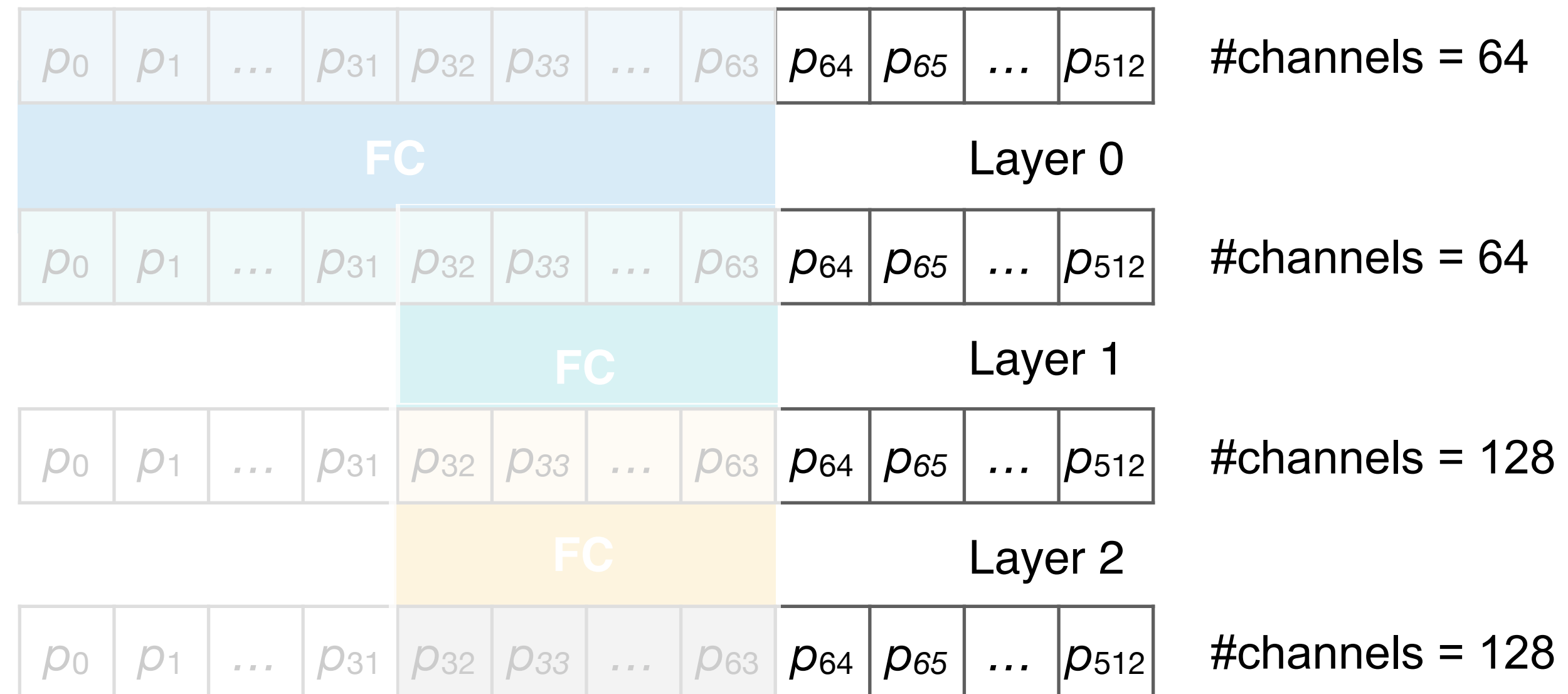
without layer fusion with layer fusion



Input Buffers (Stack mode)



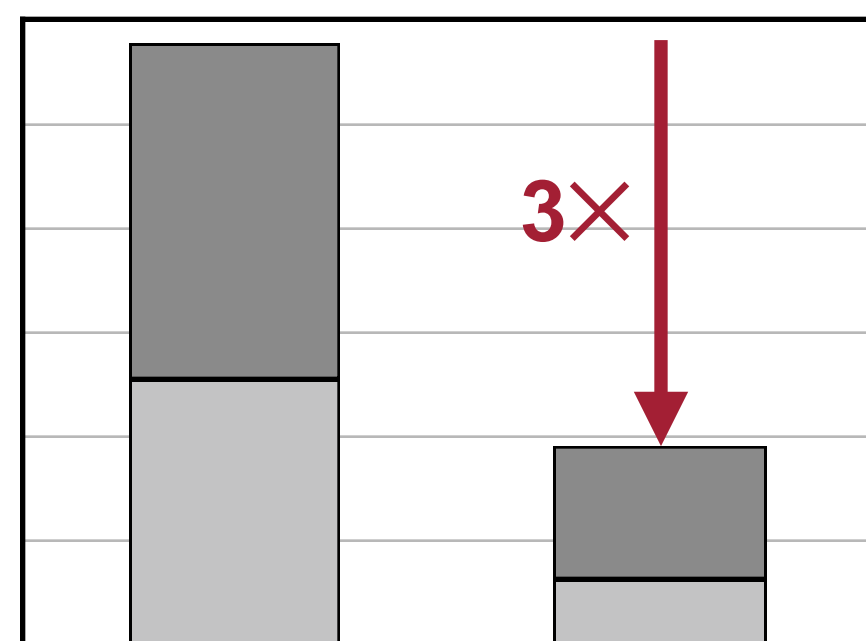
Output Buffers



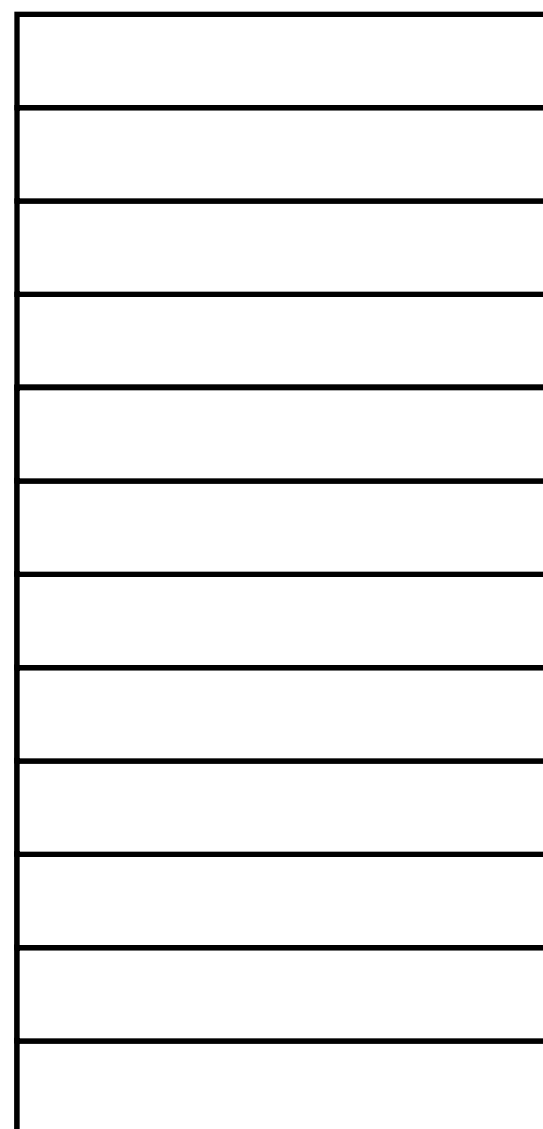
Temporal Layer Fusion on consecutive FCs

DRAM access per point

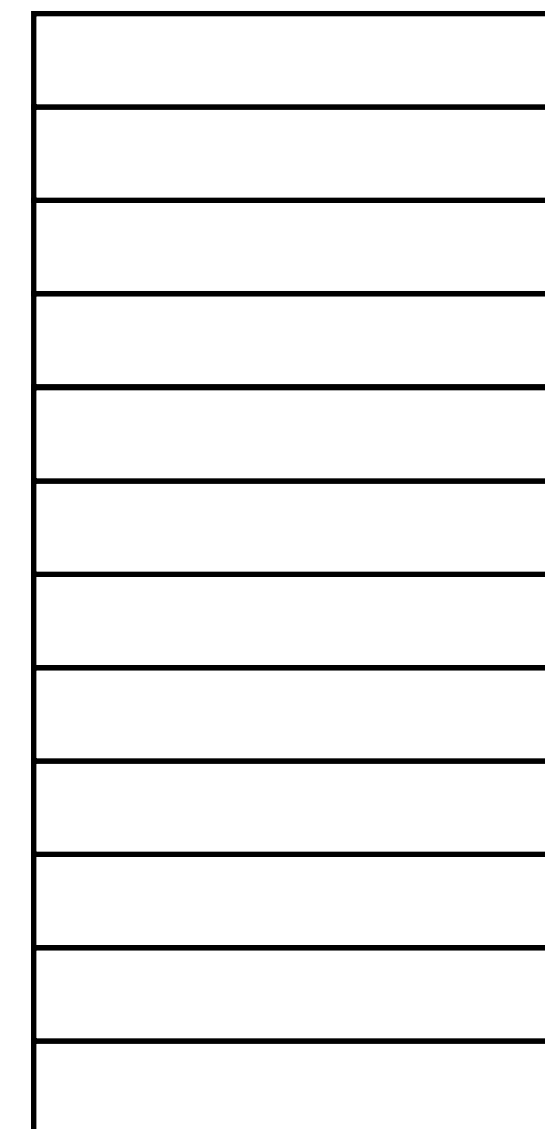
read write



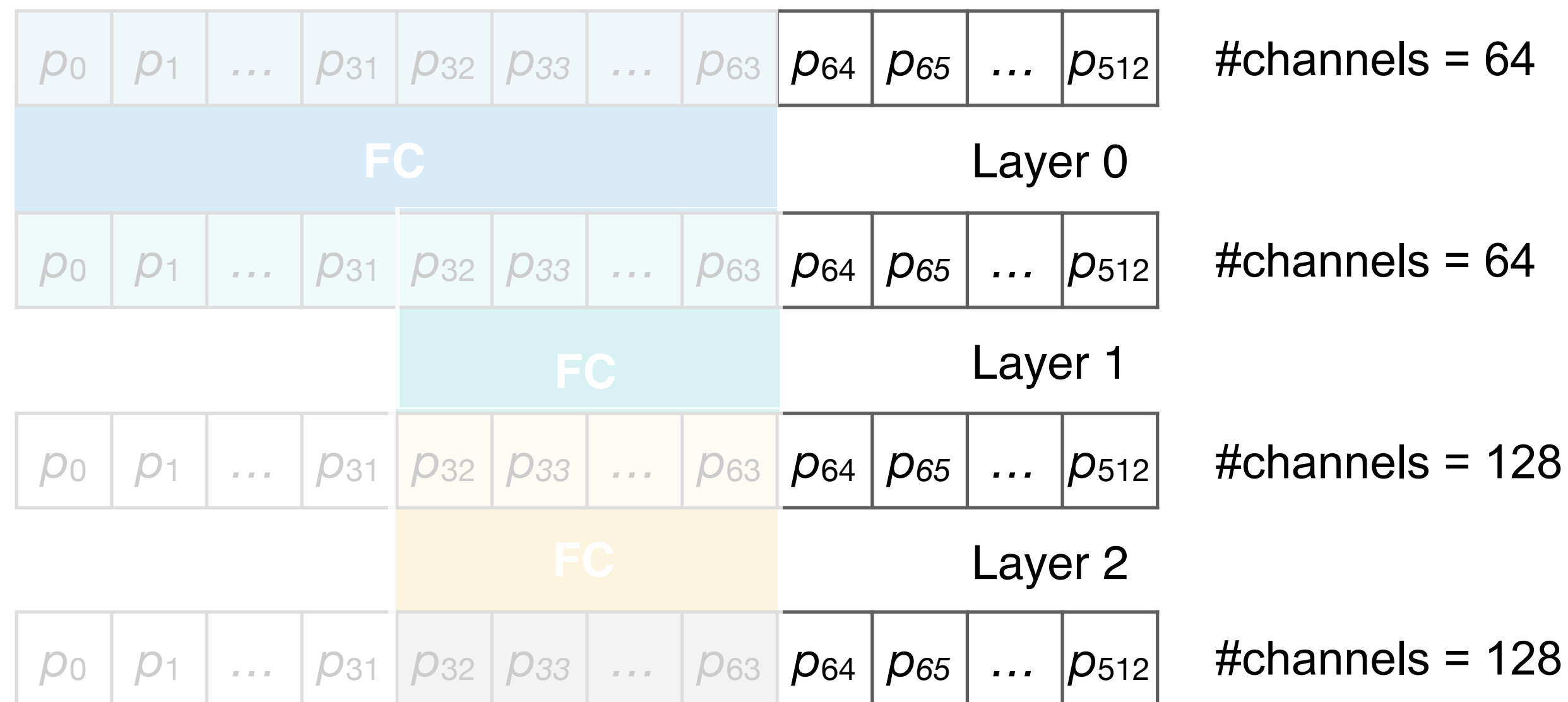
without layer fusion with layer fusion



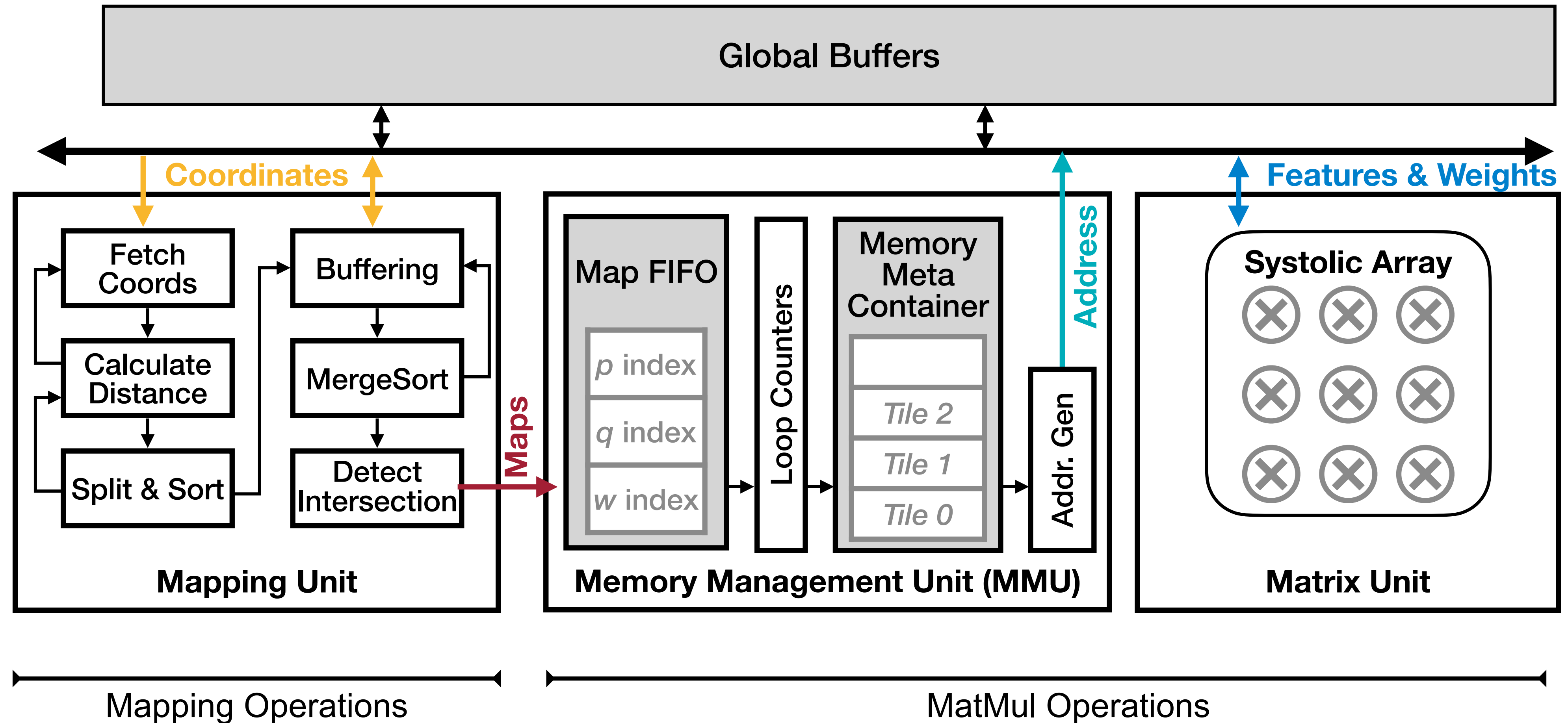
Input Buffers (Stack mode)



Output Buffers



PointAcc: Efficient Point Cloud Accelerator



Evaluation Setup

Evaluation benchmarks

- ▶ 4 different applications
- ▶ 5 different datasets
- ▶ 8 different point cloud models

Application	Dataset	Scene	Model
Classification	ModelNet40	Object	PointNet
			PointNet++ (c)
			PointNet++ (ps)
Part Segmentation	ShapeNet		DGCNN
Detection	KITTI	Outdoor	MinkNet(o)
Semantic Segmentation	S3DIS	Indoor	F-PointNet++
			PointNet++(s)
	SemanticKITT	Outdoor	MinkNet(i)

Evaluation Setup

Evaluation benchmarks

- ▶ 4 different applications: classification, part segmentation, detection, semantic segmentation
- ▶ 5 different datasets, ranging from single object to indoor scenes to outdoor scenes
- ▶ 8 different point cloud models, including classical and the state-of-the-art ones

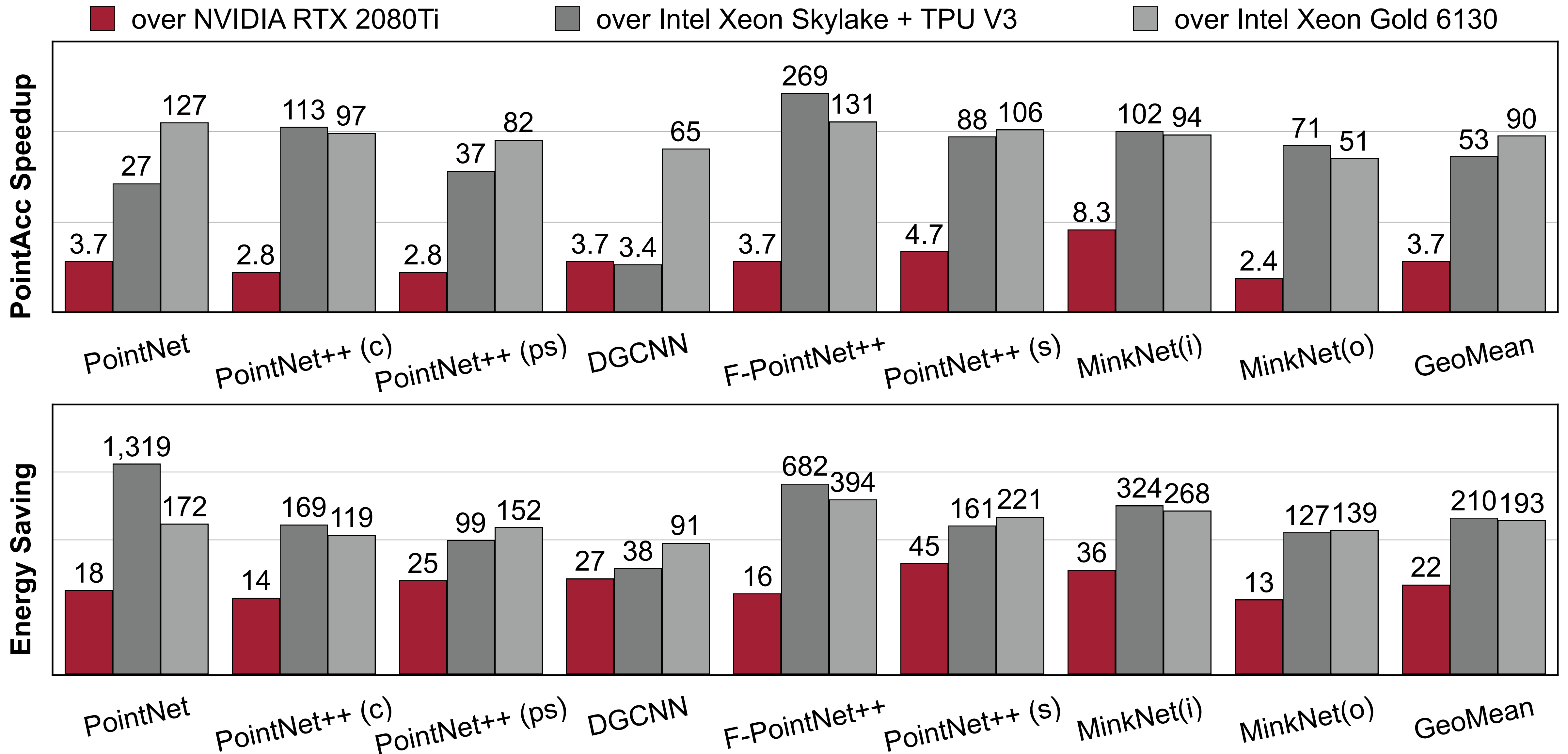
Hardware Baselines

- ▶ server-level products: Intel Xeon CPU, RTX 2080Ti GPU, TPU v3
- ▶ edge devices: Jetson Xavier NX, Jetson Nano, Raspberry Pi
- ▶ specialized point cloud NN ASIC: Mesorasi

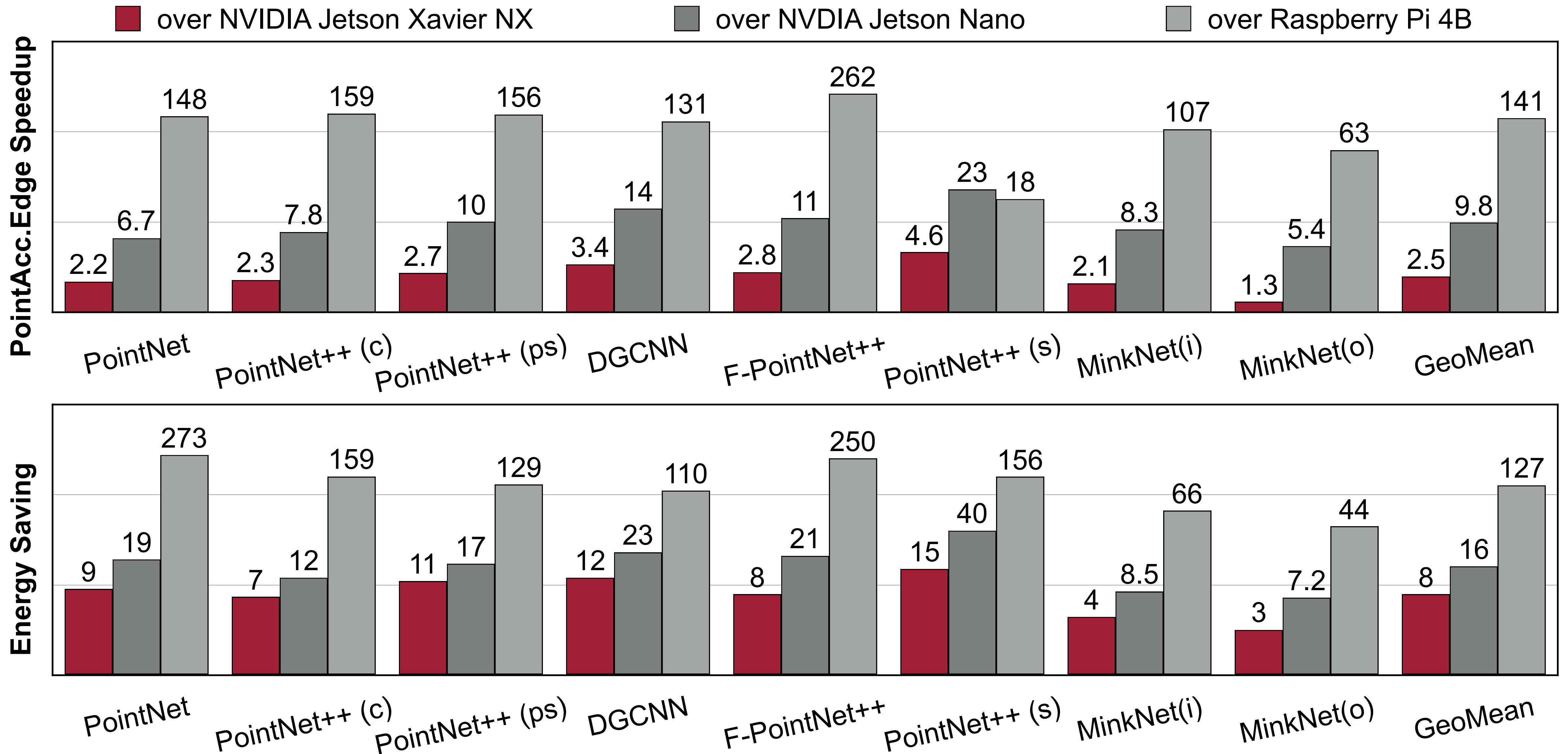
Variants

- ▶ PointAcc: 64×64 systolic array with 776 KB on-chip memory
- ▶ PointAcc.Edge: 16×16 systolic array with 274 KB on-chip memory

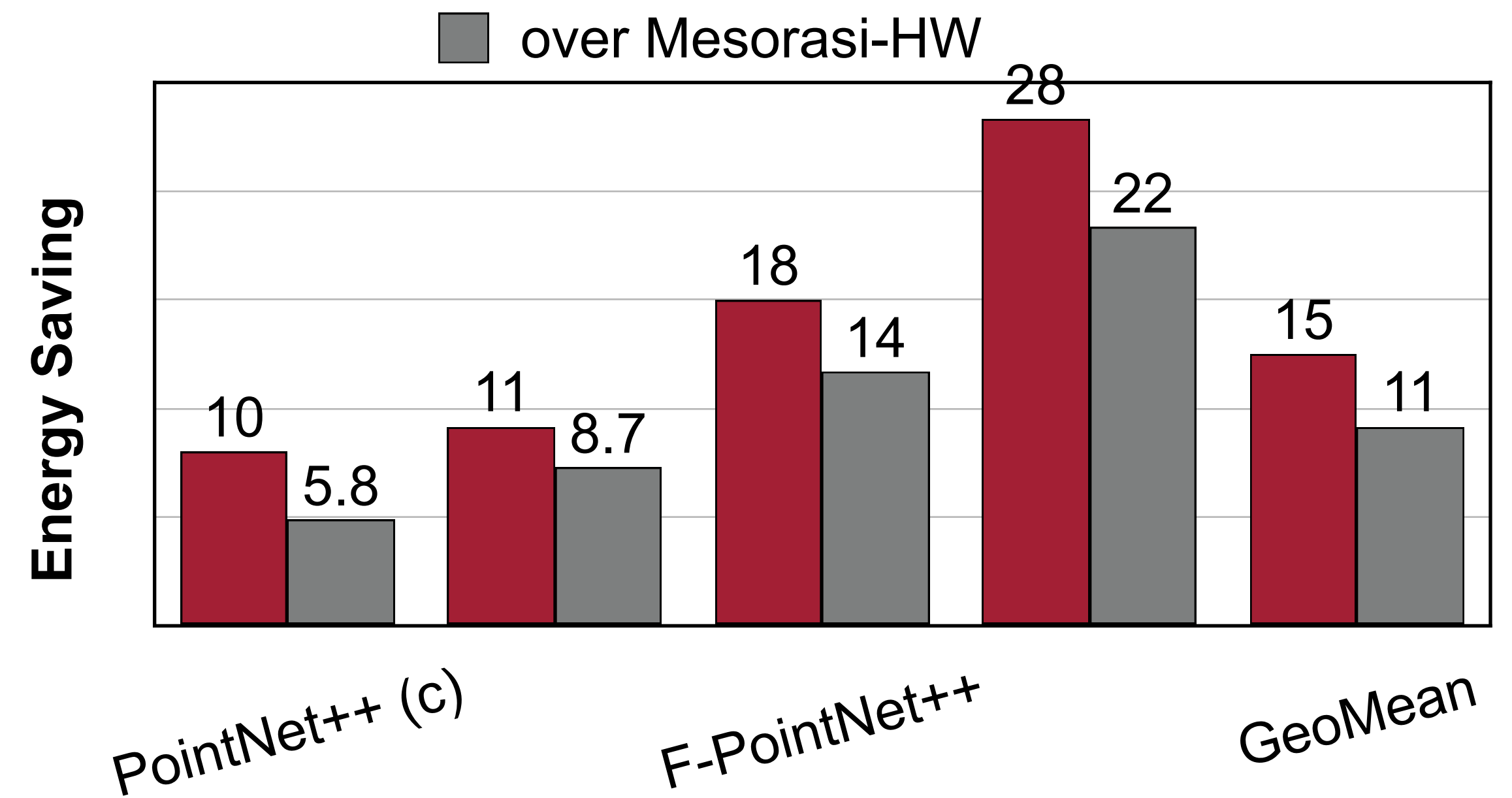
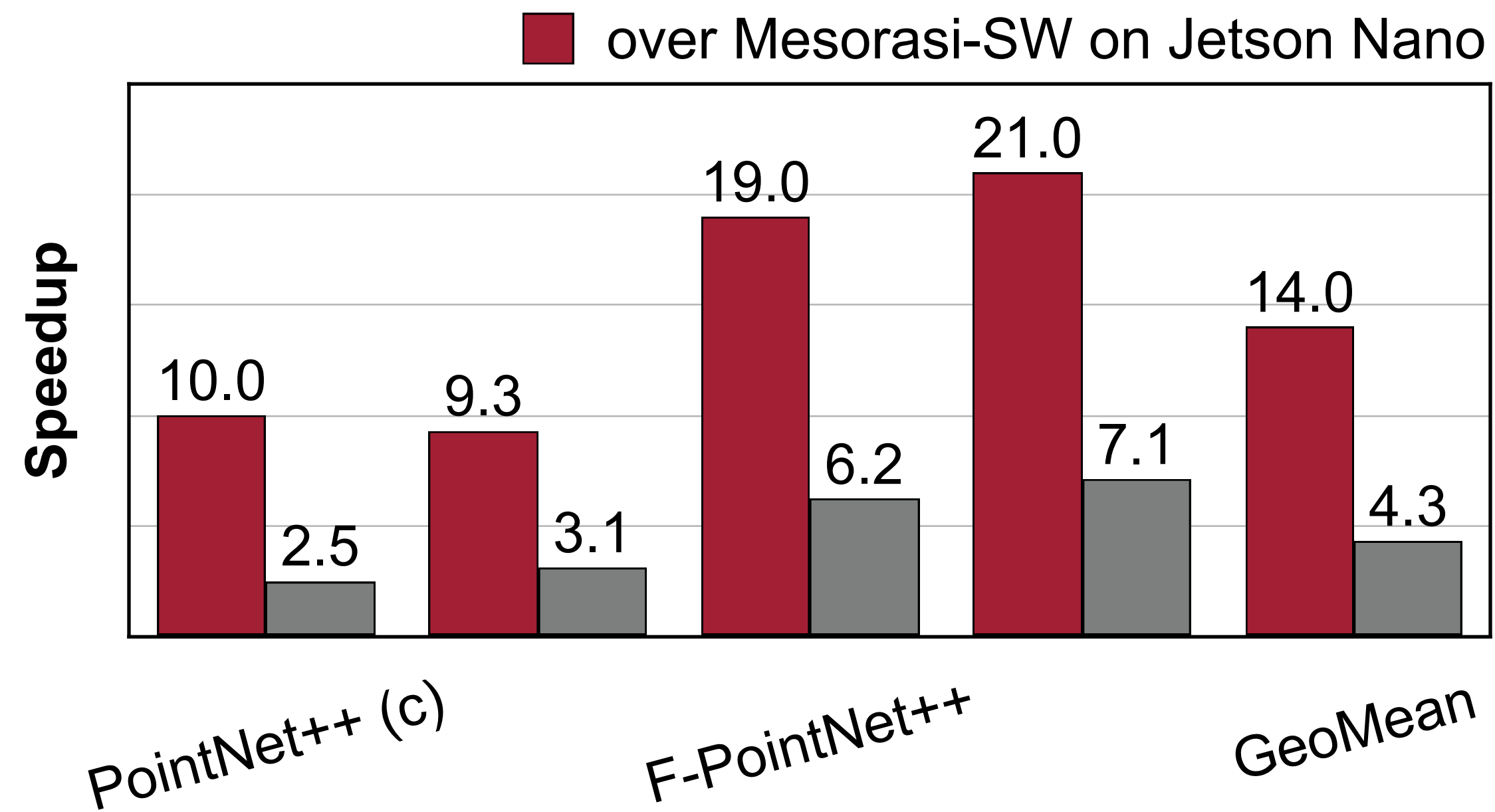
Performance Gain over the Server Products



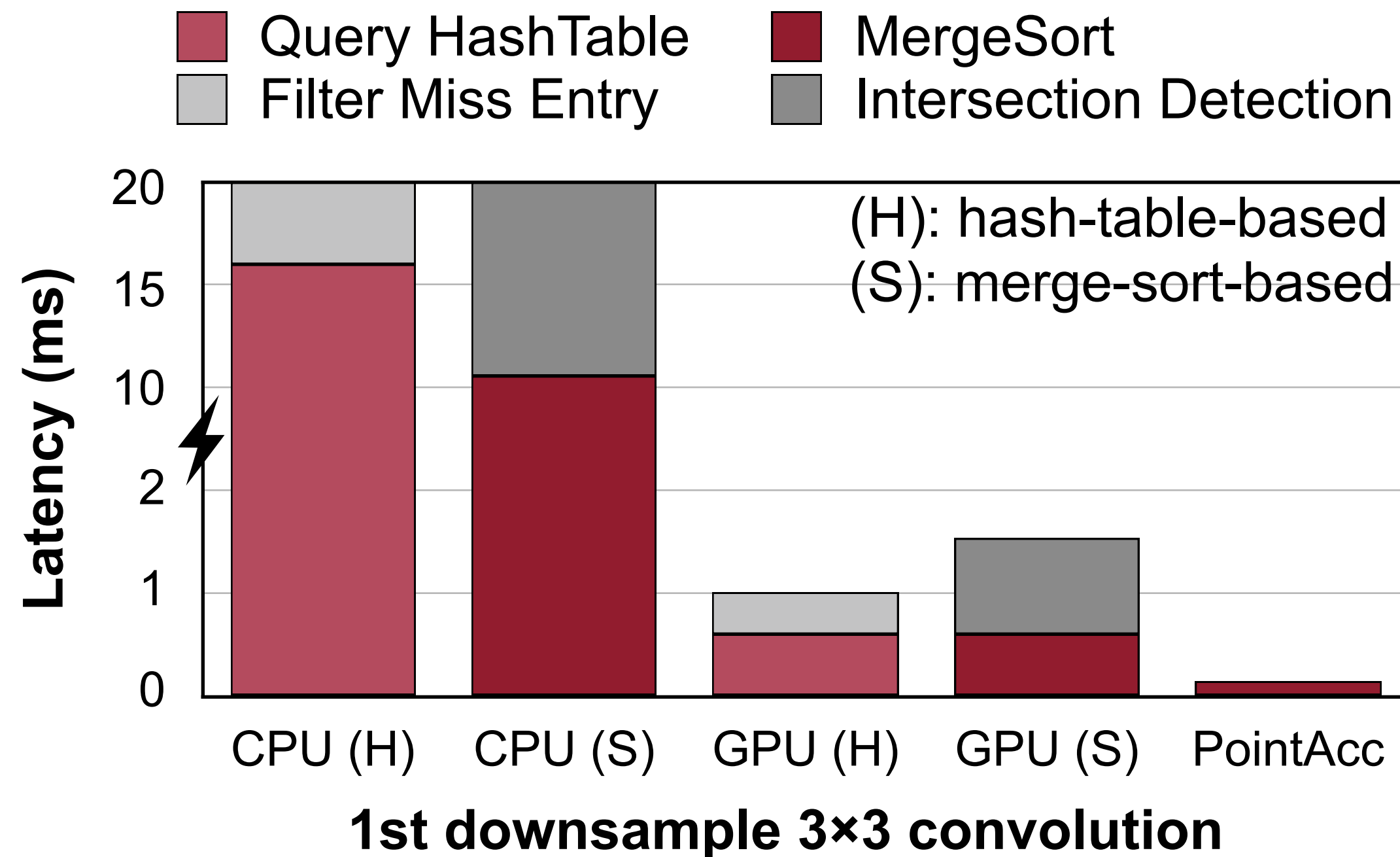
Performance Gain over the Edge Devices



Performance Gain of PointAcc.Edge over Mesorasi

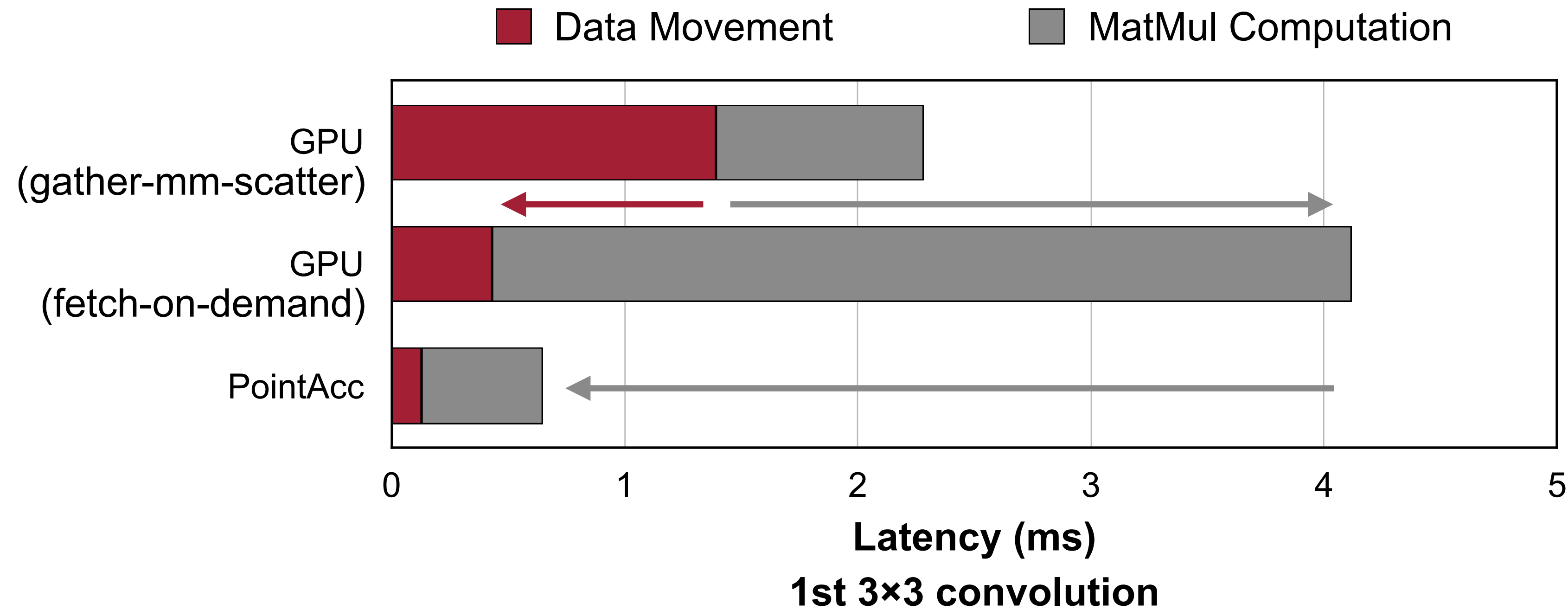


Source of Performance Gain



- ▶ Merge-sort-based implementation worsens the performance on CPU/GPU:
 - Excessive off-chip memory access between each stage of merge-sort.
 - Doubled #points in post-processing (2X #points after merge sort).
- ▶ PointAcc spatially pipelines the stages of merge sort and intersection detection.
- ▶ Using one versatile architecture for different mapping ops does not hinder the performances.

Source of Performance Gain



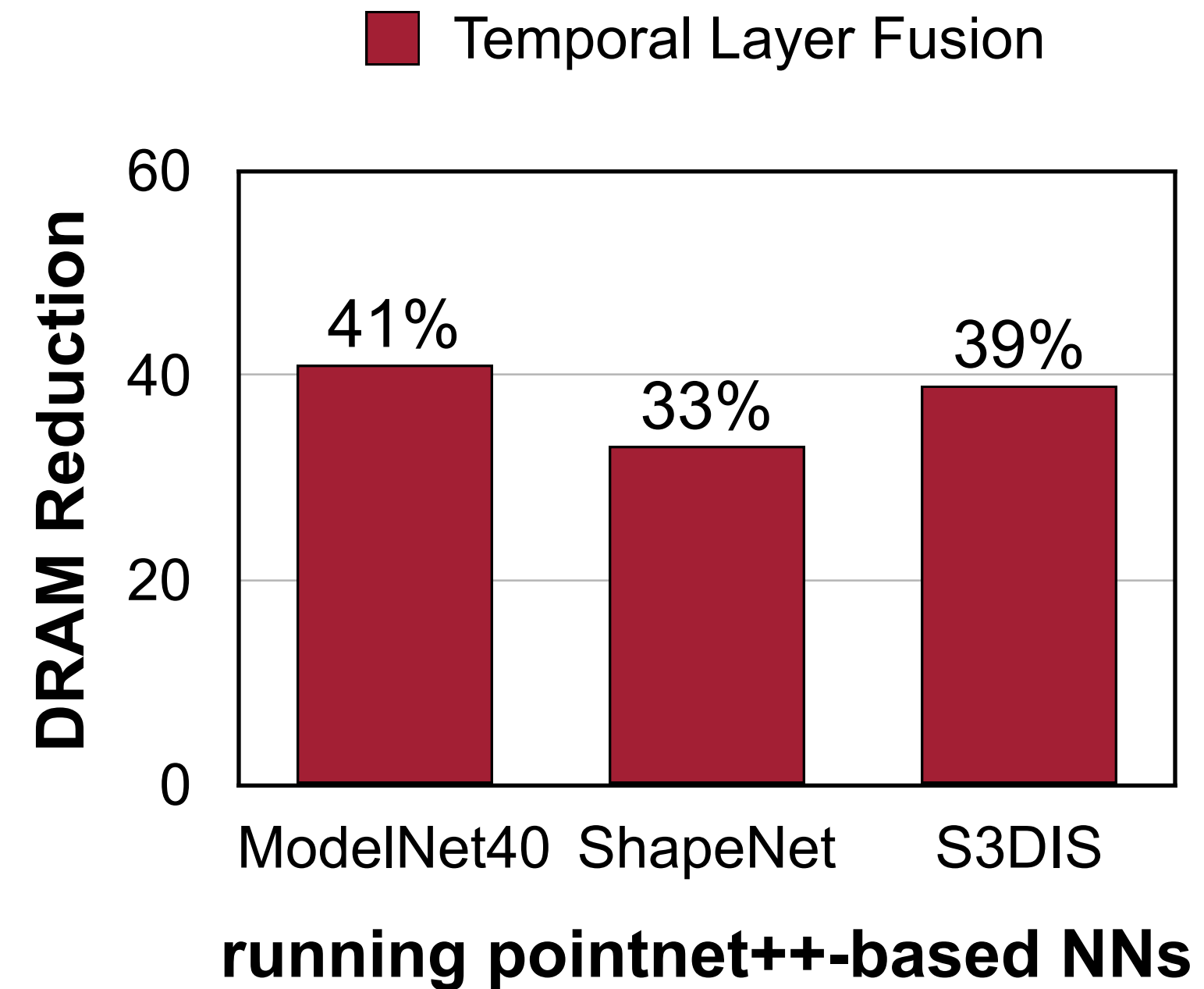
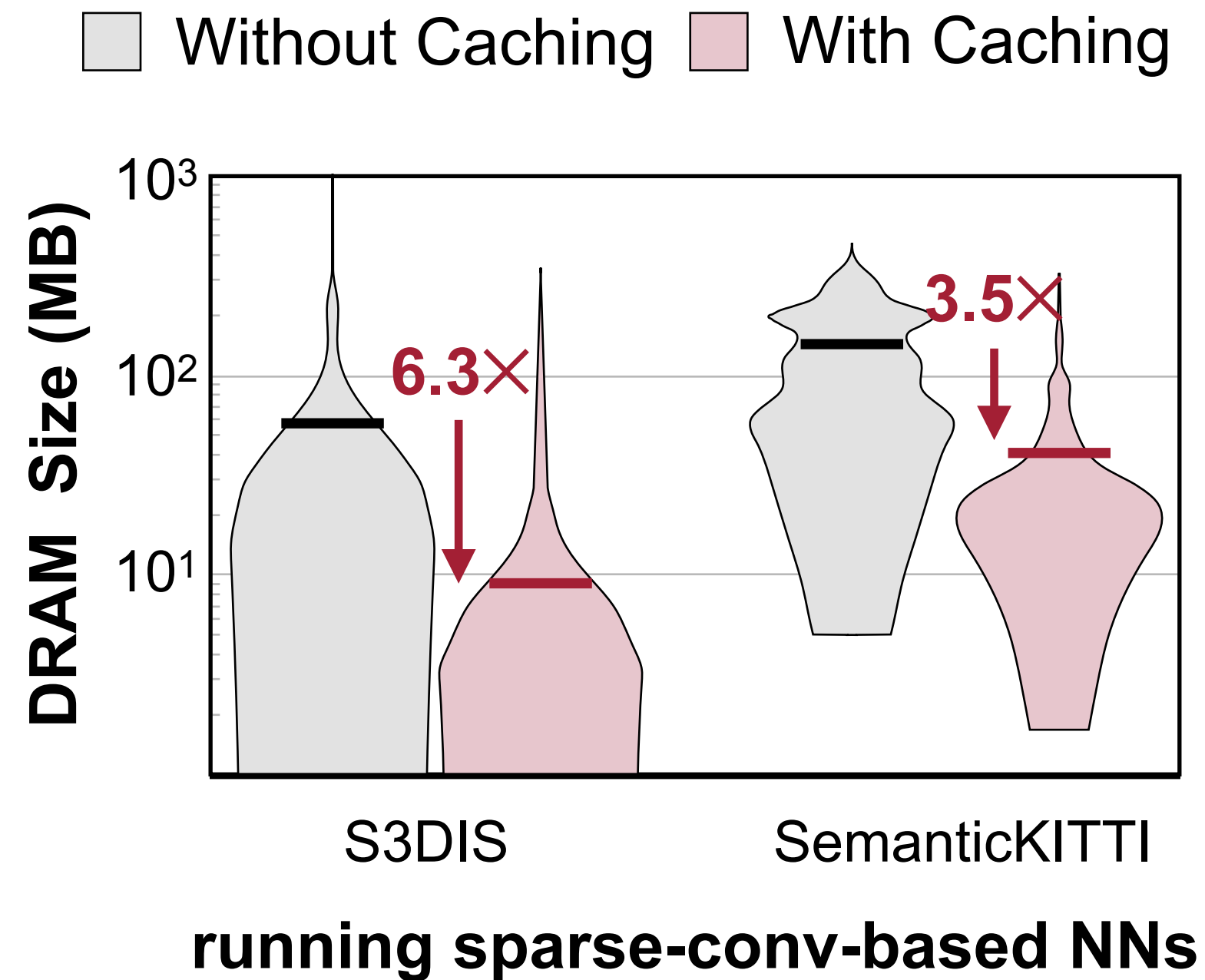
On GPU

- ▶ Fetch-on-demand flow saves the data movement cost by 3X.
- ▶ Decomposing the MM into MV multiplication significantly increases the computation overhead.

PointAcc

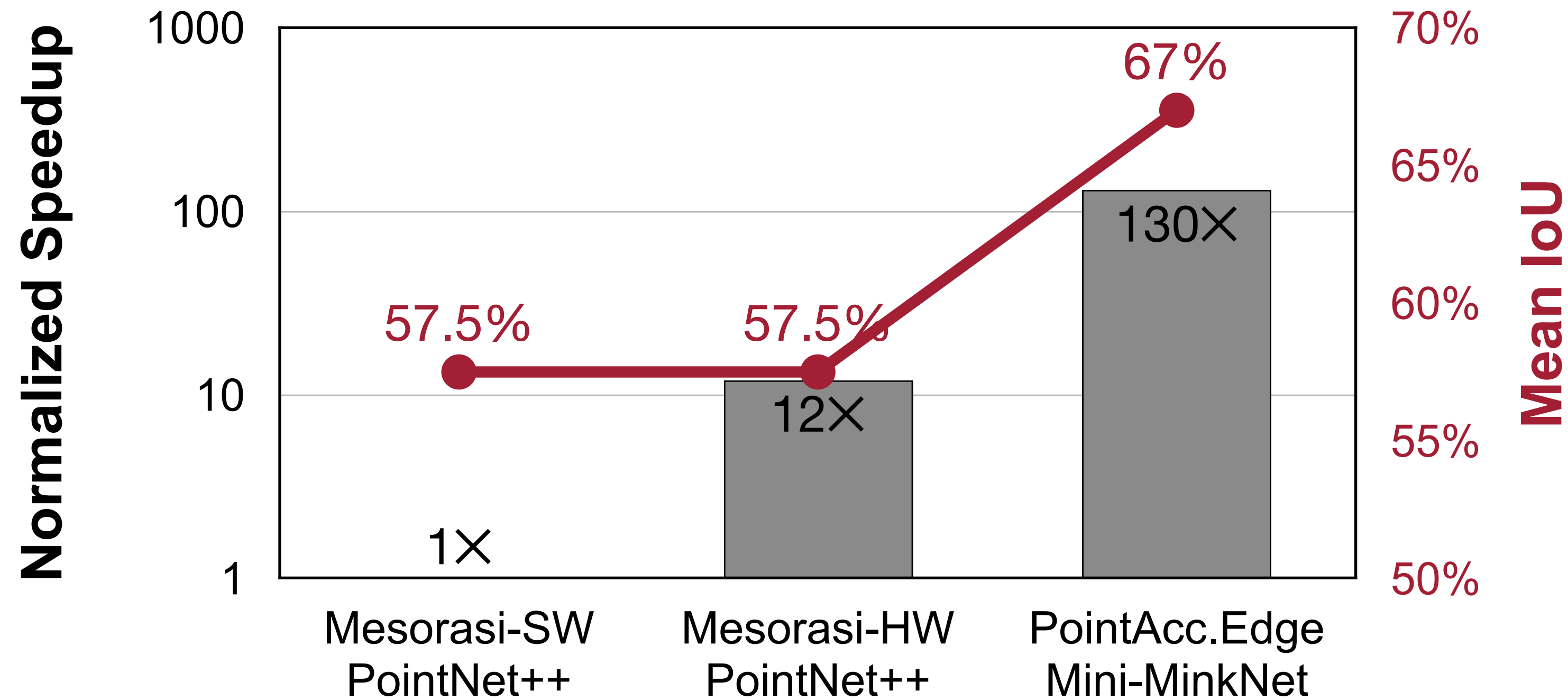
- ▶ Decoupled data orchestration and powerful systolic array cancel the MV overhead.

Source of Performance Gain



- ▶ The shape of distribution of the off-chip access data size per layer is nearly the same with and without caching → caching works consistently on different layers and different datasets.
- ▶ Caching reduces the off-chip memory footprint by 3.5X to 6.3X.
- ▶ Temporal layer fusion cuts the off-chip memory footprint from 33% to 41%.

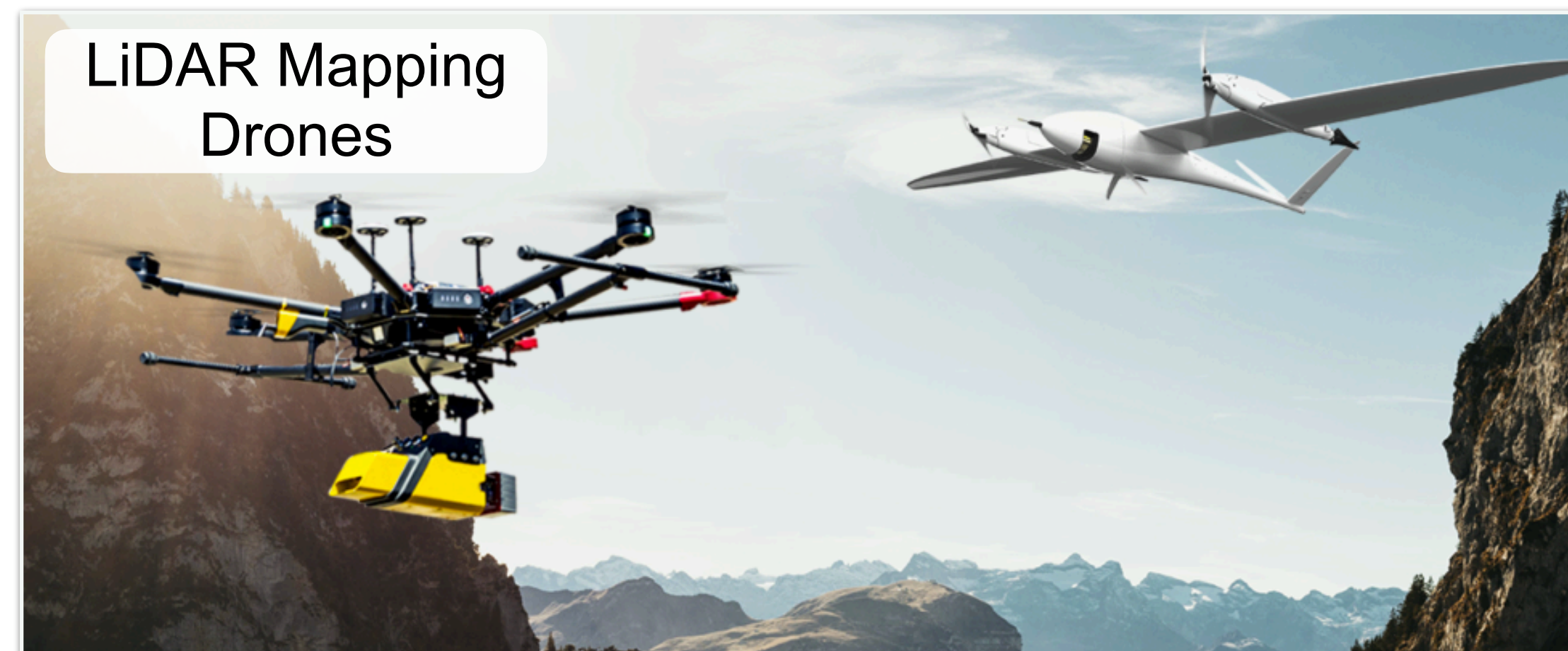
PointAcc.Edge v.s. Mesorasi



- ▶ Mesorasi does not support independent weights for different neighbors
- ▶ The state-of-the-art models tend to use independent weights for different neighbors
- ▶ Running the same segmentation task on S3DIS dataset, PointAcc.Edge is **9.1% higher accuracy with 130X lower latency.**

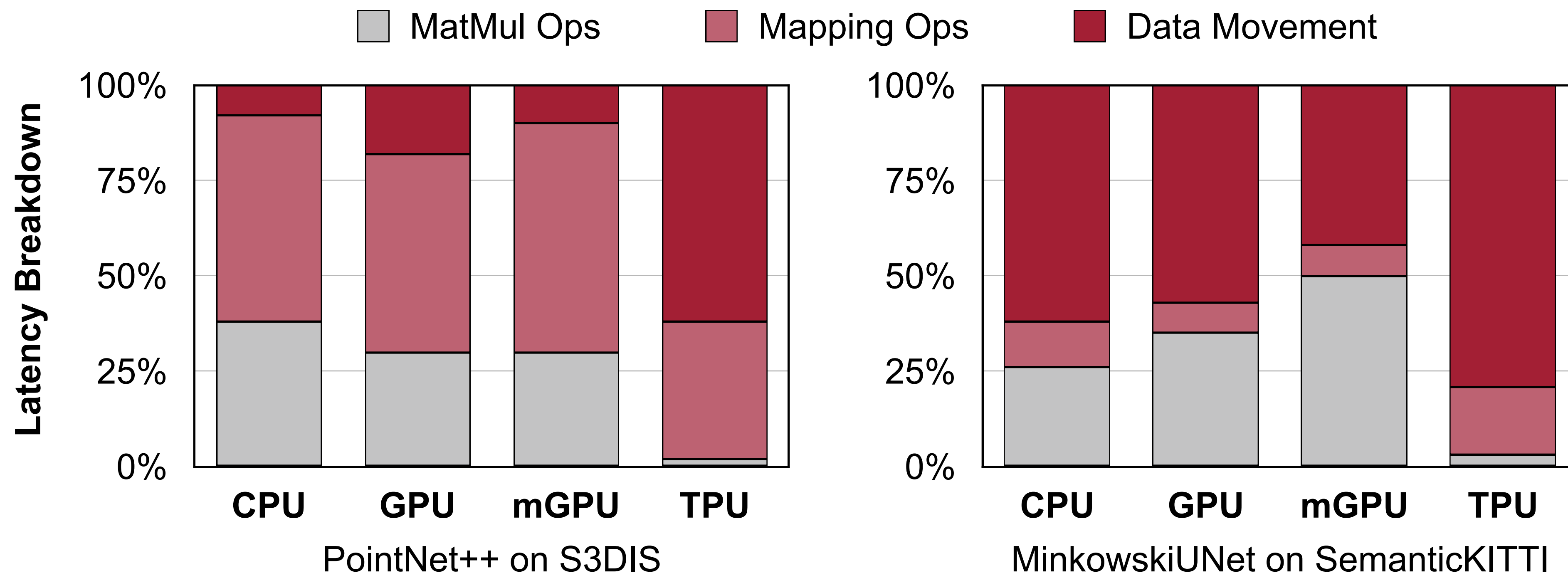
Conclusion

- ▶ With the rise of point cloud modality, the rapid development of point cloud deep learning brings new challenges and exciting opportunities for intelligent hardware design.



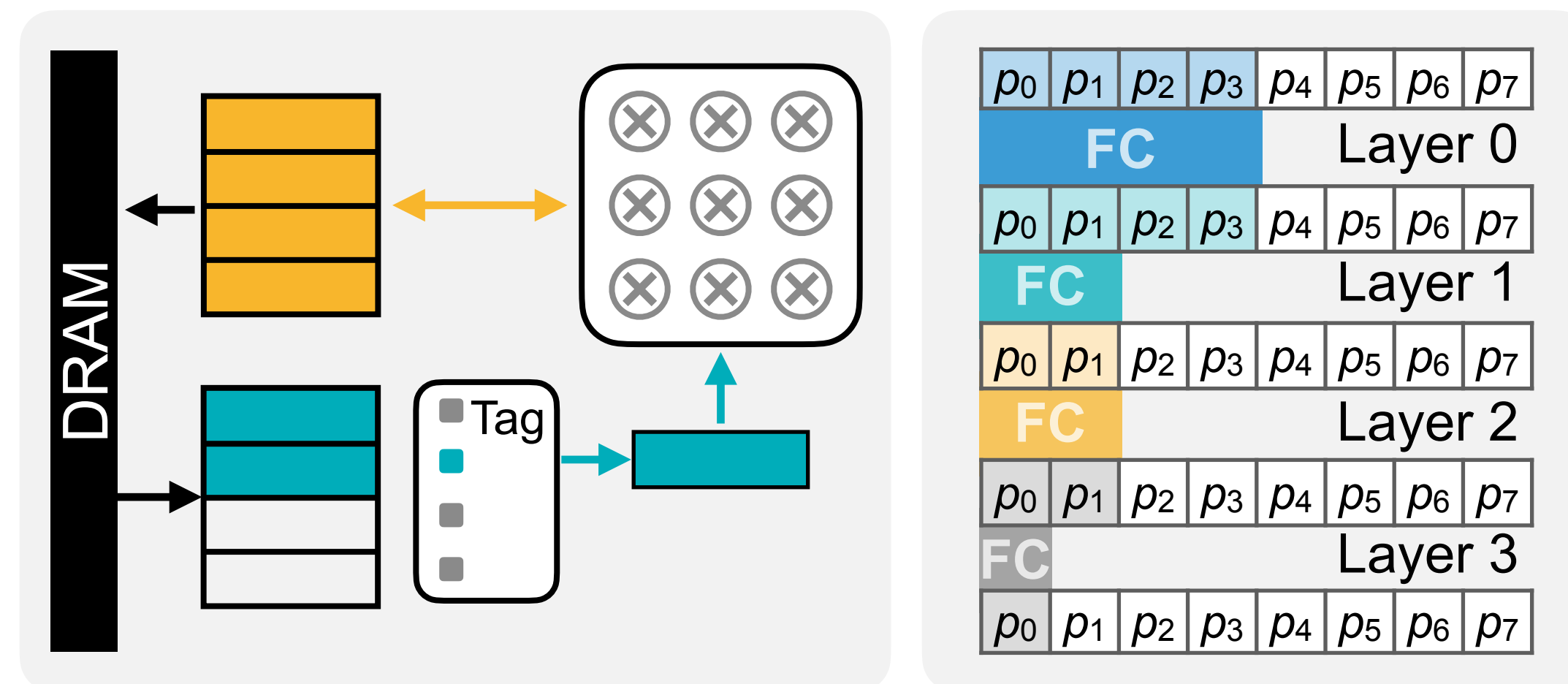
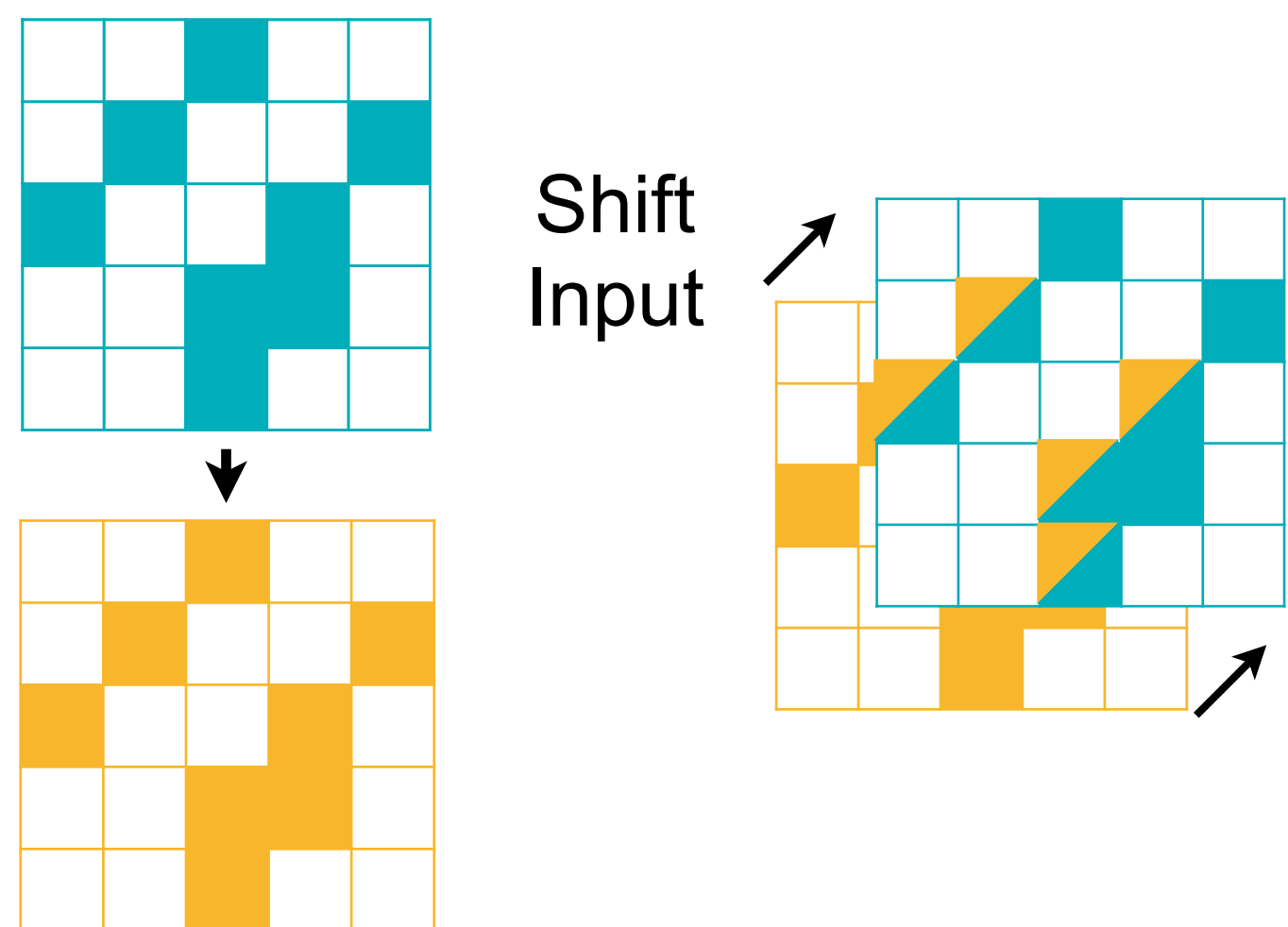
Conclusion

- ▶ With the rise of LiDAR sensors, the rapid development of point cloud deep learning brings new challenges and exciting opportunities for intelligent hardware design.
- ▶ The sparsity of point clouds in the physical space leads to two execution bottlenecks:
 - ▶ Newly introduced mapping operations for searching (input, output, weight) maps
 - ▶ Data movement overhead from gather and scatter the sparsely distributed features



Conclusion

- ▶ With the rise of LiDAR sensors, the rapid development of point cloud deep learning brings new challenges and exciting opportunities for intelligent hardware design.
- ▶ The sparsity of point clouds in the physical space leads to two execution bottlenecks:
 - ▶ Newly introduced mapping operations for searching (input, output, weight) maps
 - ▶ Data movement overhead from gather and scatter the sparsely distributed features
- ▶ PointAcc maps diverse mapping ops into sort-based computation with one versatile architecture.
- ▶ PointAcc reduces off-chip memory access and minimize the overhead of gather and scatter by flexible caching and layer fusion.

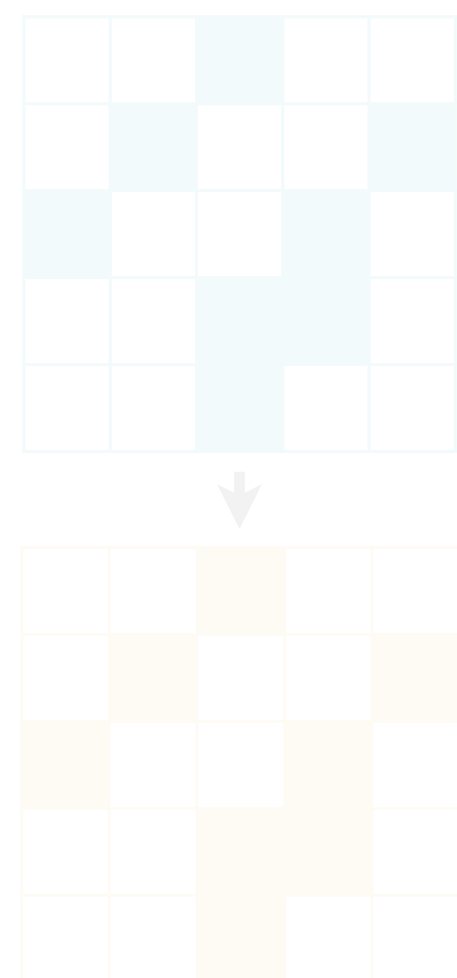


Conclusion

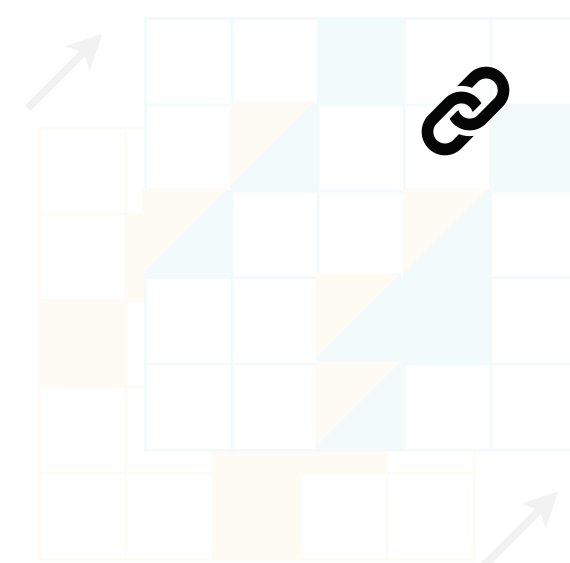
- ▶ With the rise of LiDAR sensors, the rapid development of point cloud deep learning brings new challenges and exciting opportunities for intelligent hardware design.
- ▶ The sparsity of point clouds in the physical space leads to two execution bottlenecks:
 - ▶ Newly introduced mapping operations for searching (input, output, weight) maps
 - ▶ Data movement overhead from gather and scatter the sparsely distributed features

Thank You

architecture.



Shift
Input



✉ **{yujunlin, songhan}@mit.edu**

🔗 <https://hanlab.mit.edu/projects/pointacc>

scatter by flexible caching and layer fusion.

