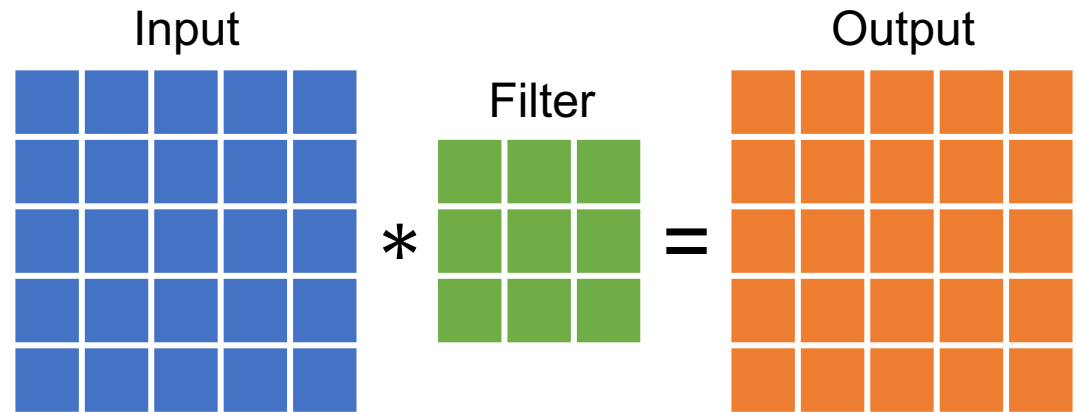




Unified Convolution Framework

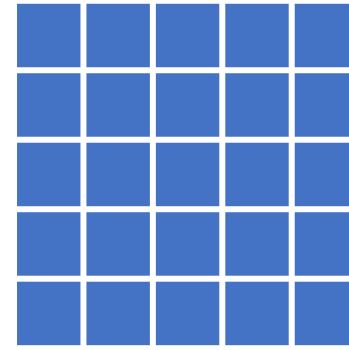
Jaeyeon Won, Changwan Hong, Charith Mendis, Joel Emer, Saman Amarasinghe

Convolution

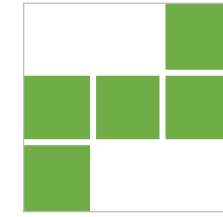


Sparsity in Convolution

Filter Sparse Convolution

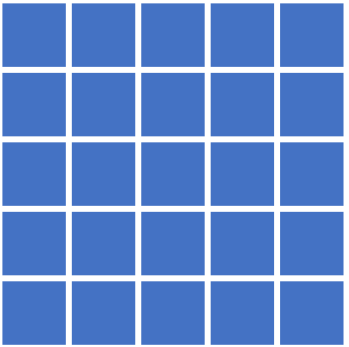


Sparse Filter



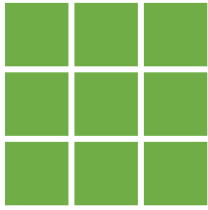
*

Input



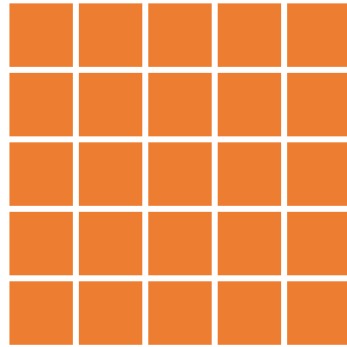
*

Filter

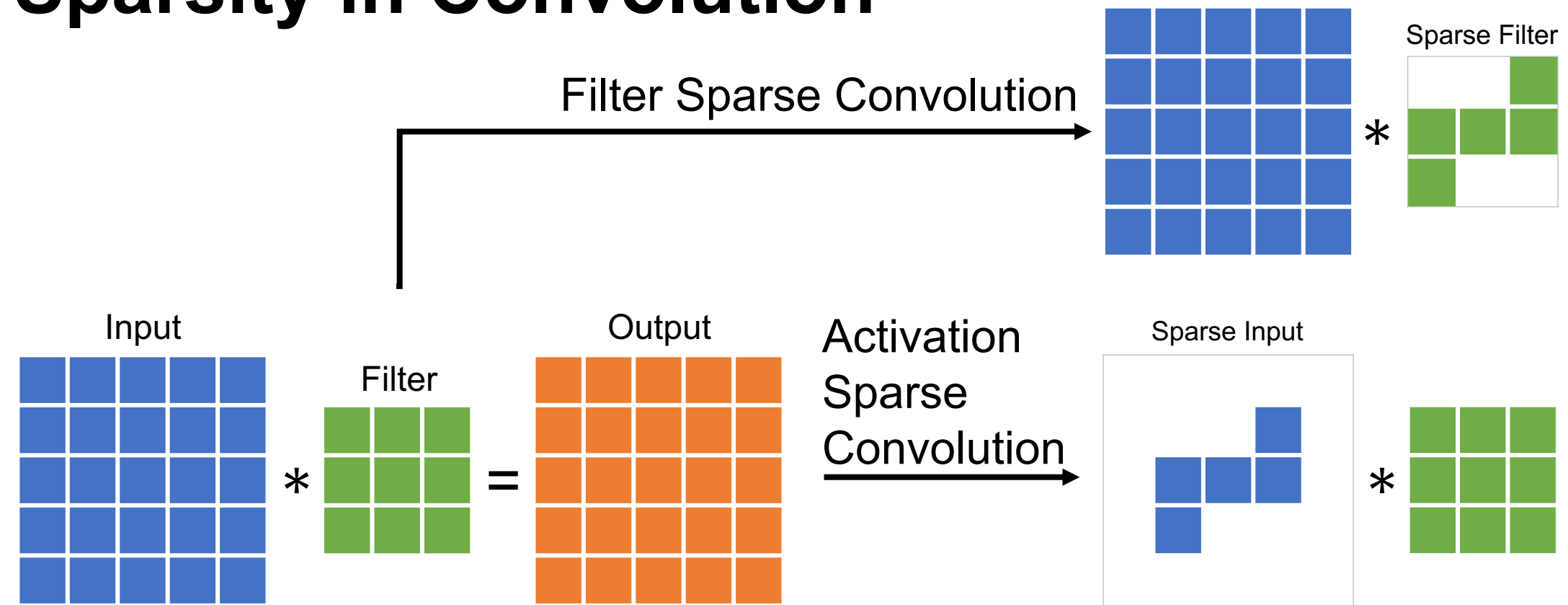


=

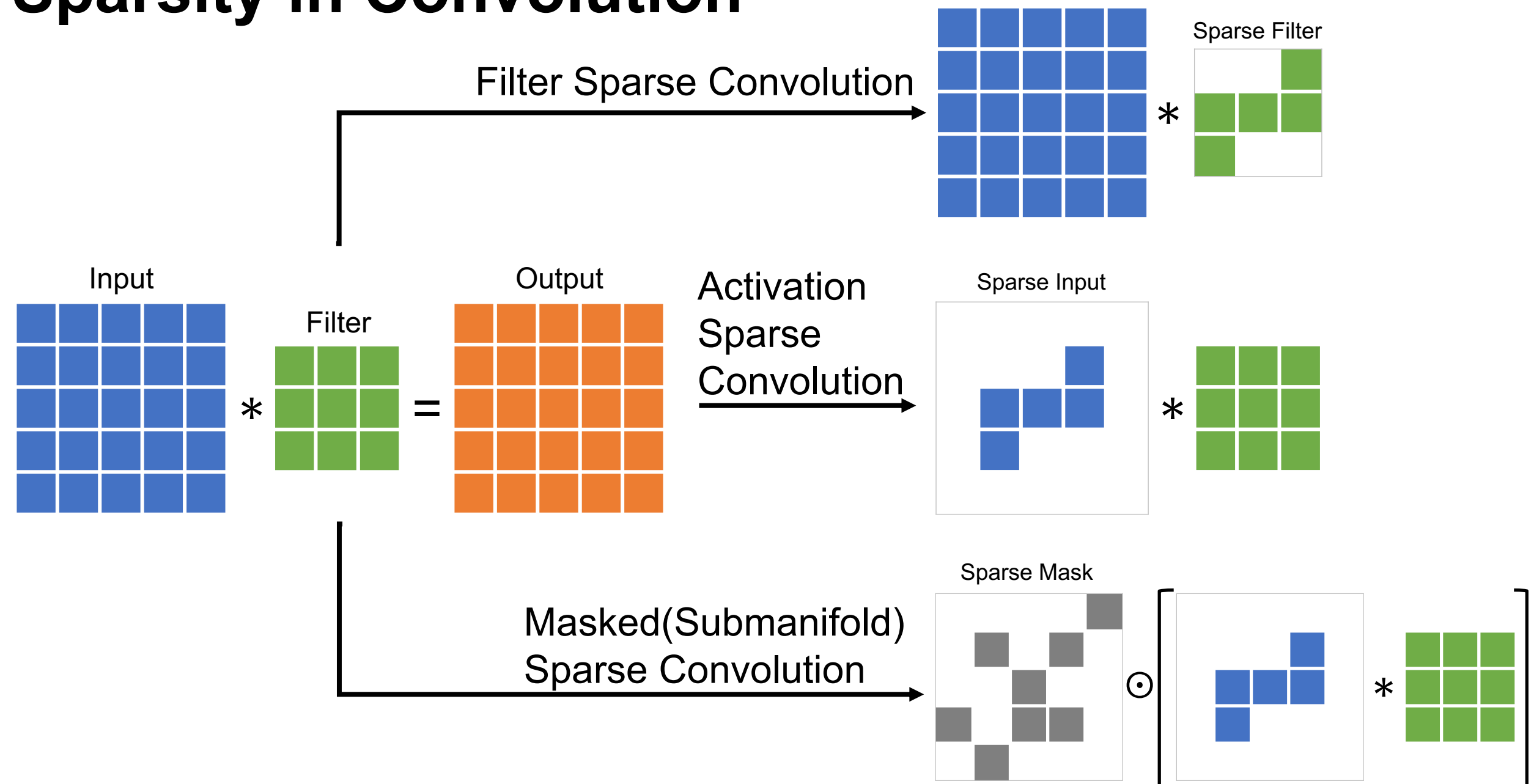
Output



Sparsity in Convolution



Sparsity in Convolution



Sparsity in Convolution

Filter Sparse Convolution



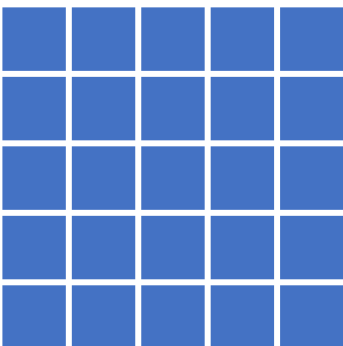
IntelLabs/
SkimCaffe

intel labs

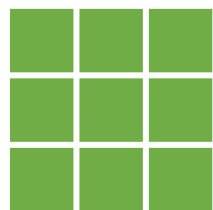
Caffe for Sparse Convolutional Neural Network

223 Contributors 16 Issues 236 Stars 64 Forks

Input



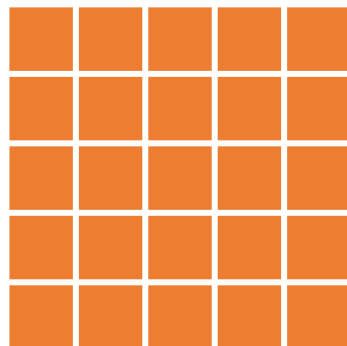
Filter



*

=

Output



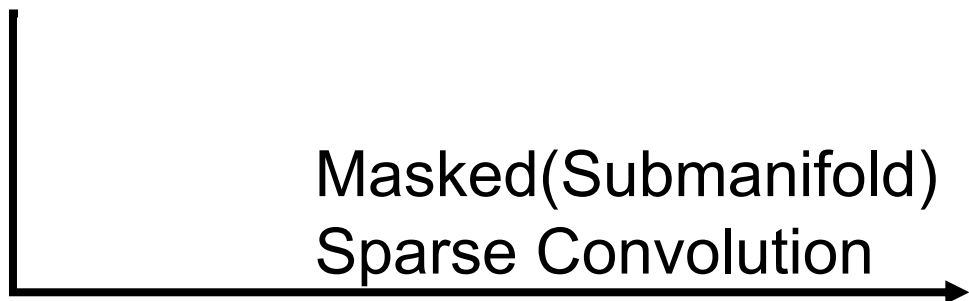
Activation
Sparse
Convolution



DeepSparse Engine

A sparsity-aware neural network inference engine that delivers GPU-class performance on commodity CPUs, anywhere.

Masked(Submanifold)
Sparse Convolution



mit-han-lab/
torchsparse

MIT HAN

[MLSys'22] TorchSparse: Efficient Point Cloud Inference Engine

13 Contributors 6 Used by 1 Discussion 770 Stars 102 Forks

Limitation of Library-based Approach

Name	Sparse Convolutions				Formats	Backends	
	Filter SpConv	Activation SpConv	Submanifold SpConv	Dual SpConv		CPU	GPU
SkimCaffe	✓	✗	✗	✗	1	✓	✗
TorchSparse	✗	✗	✓	✗	1	✓	✓
DeepSparse	✓	✓	✗	✗	1	✓	✗

1. Unoptimized for new sparse convolutions
2. Unoptimized for various formats and backends

Limitation of Library-based Approach

Name	Sparse Convolutions				Formats	Backends	
	Filter SpConv	Activation SpConv	Submanifold SpConv	Dual SpConv		CPU	GPU
SkimCaffe	✓	✗	✗	✗	1	✓	✗
TorchSparse	✗	✗	✓	✗	1	✓	✓
DeepSparse	✓	✓	✗	✗	1	✓	✗

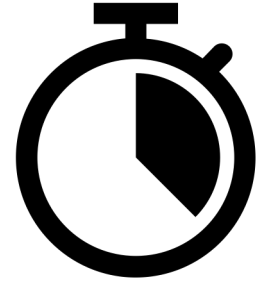
1. Unoptimized for new sparse convolutions
2. Unoptimized for various formats and backends

Limitation of Library-based Approach

Name	Sparse Convolutions				Formats	Backends	
	Filter SpConv	Activation SpConv	Submanifold SpConv	Dual SpConv		CPU	GPU
SkimCaffe	✓	✗	✗	✗	1	✓	✗
TorchSparse	✗	✗	✓	✗	1	✓	✓
DeepSparse	✓	✓	✗	✗	1	✓	✗

1. Unoptimized for new sparse convolutions
2. Unoptimized for various formats and backends

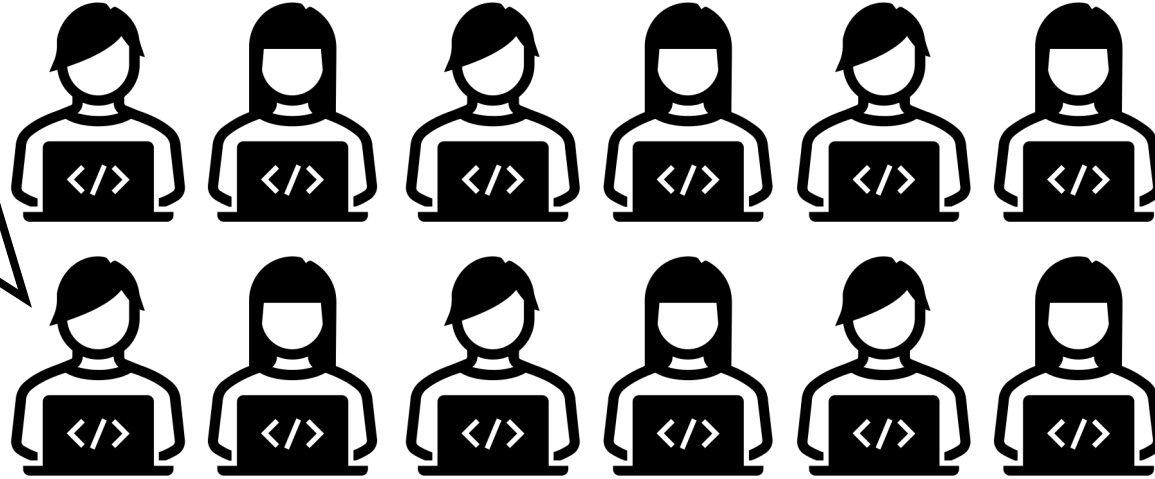
Limitation of Library-based Approach



<TODO list>

1. Optimize on Edge Device
2. Optimize on new GPU
3. Add New Features

...



1. Unoptimized for new sparse convolutions
2. Unoptimized for various formats and backends

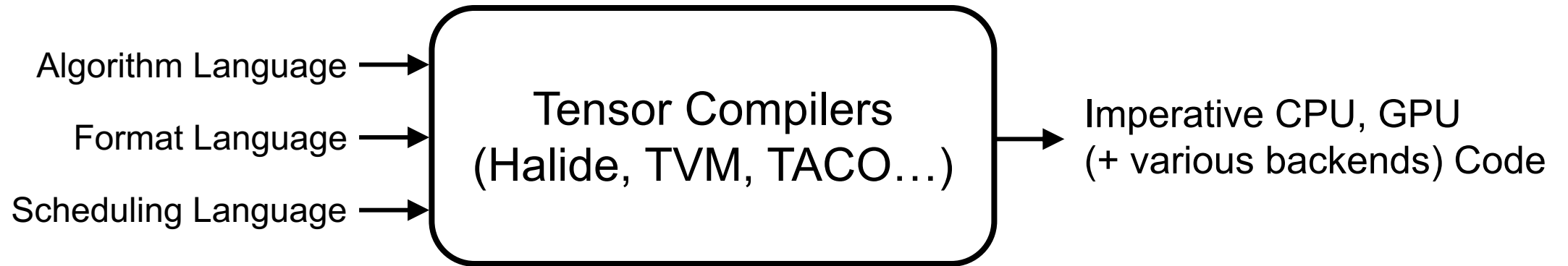
Unified Convolution Framework (UCF)

Name	Sparse Convolutions				Formats	Backends	
	Filter SpConv	Activation SpConv	Submanifold SpConv	Dual SpConv		CPU	GPU
SkimCaffe	✓	✗	✗	✗	1	✓	✗
TorchSparse	✗	✗	✓	✗	1	✓	✓
DeepSparse	✓	✓	✗	✗	1	✓	✗
Our Work (TACO-UCF)	✓	✓	✓	✓	> 100	✓	✓

Unified Convolution Framework is a compiler that supports

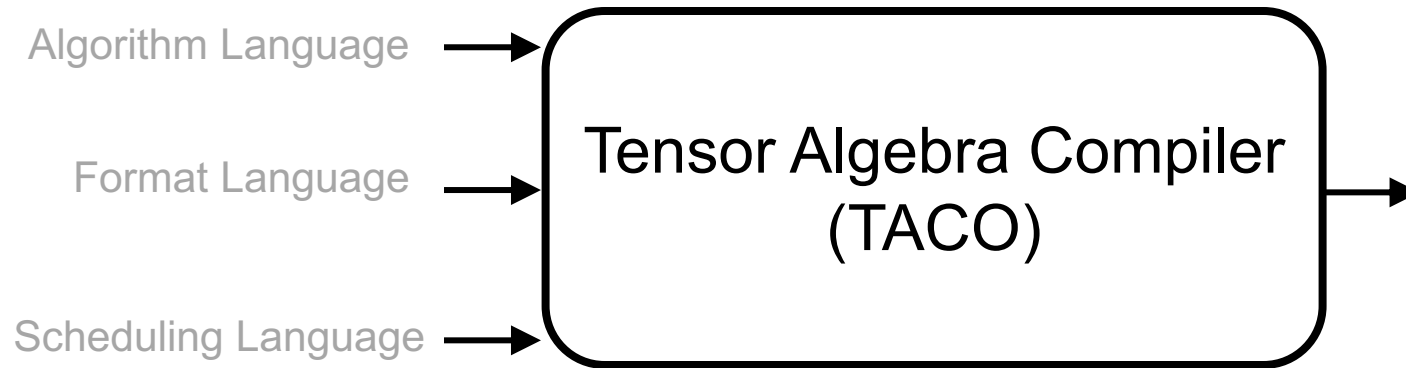
- (1) all sparse convolutions
- (2) on various formats and backends

Background : Tensor Compiler



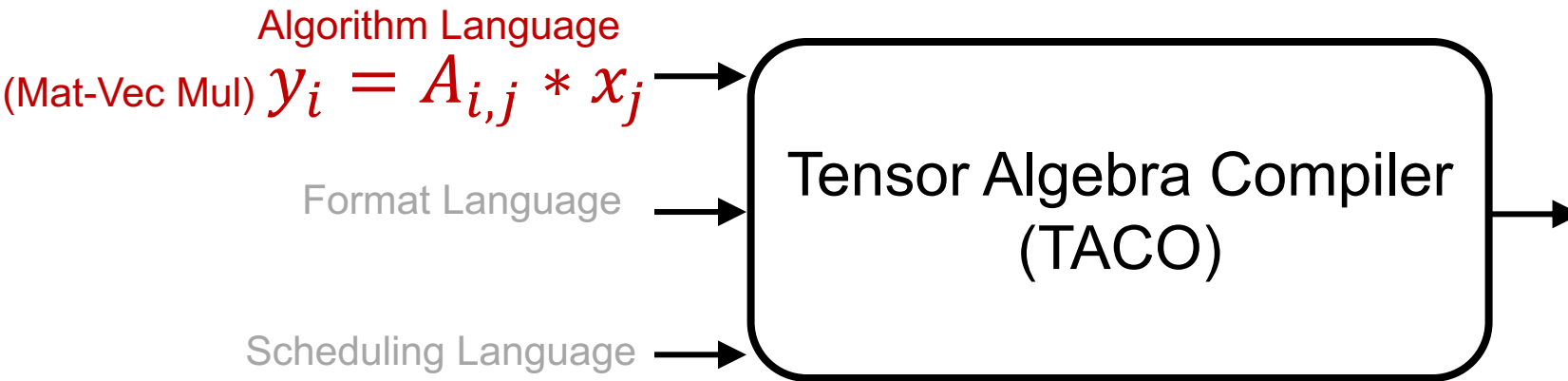
Background : Tensor Compiler (TACO)

What we want : Sparse Matrix – Dense Vector Multiplication Kernel

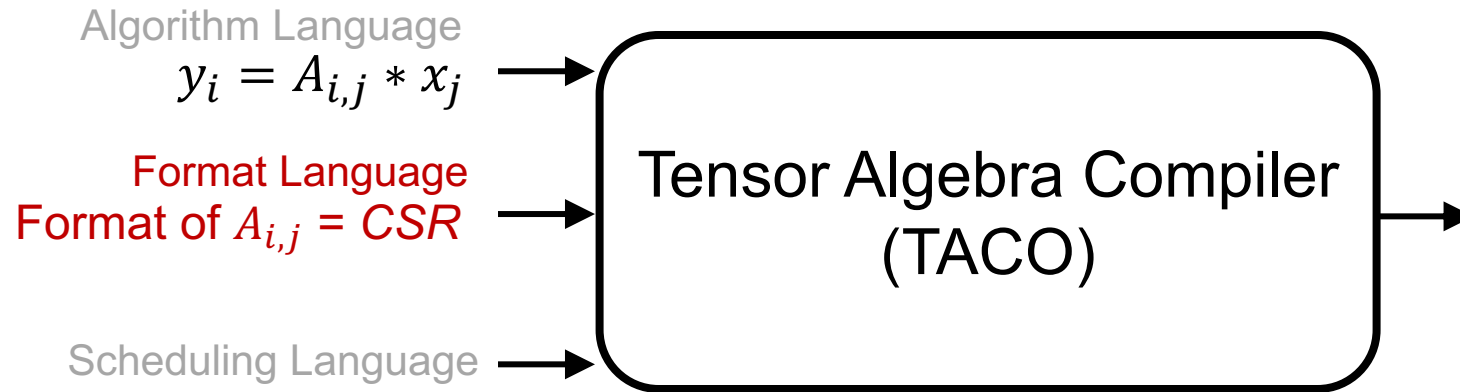


Background : Tensor Compiler (TACO)

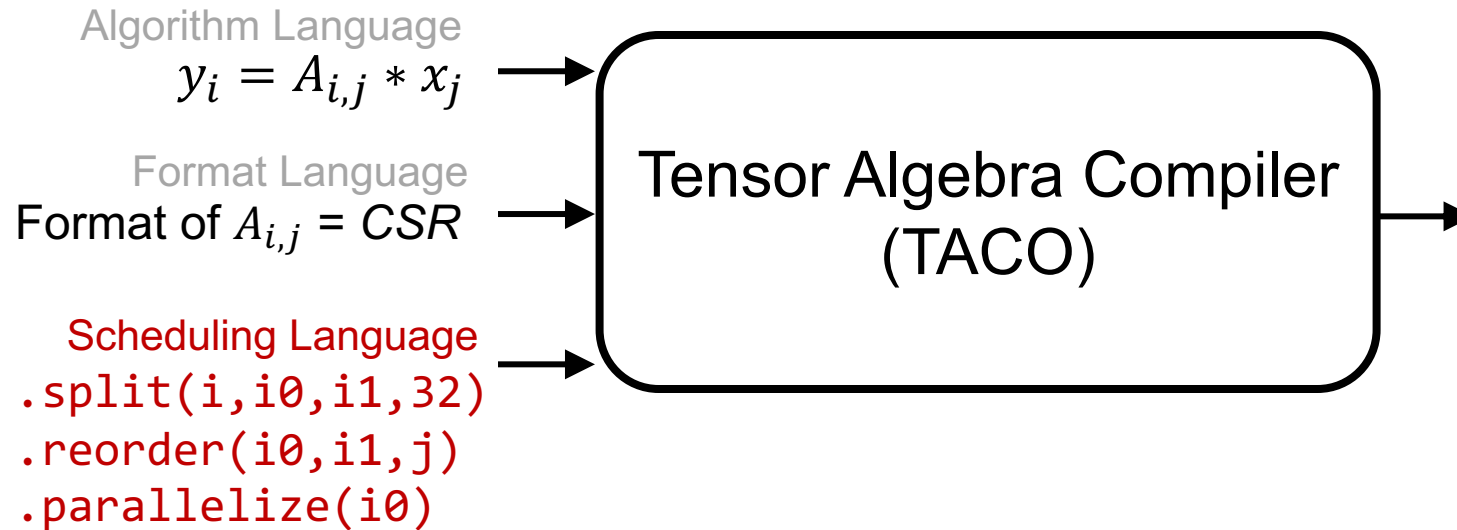
What we want : Sparse Matrix – Dense Vector Multiplication Kernel



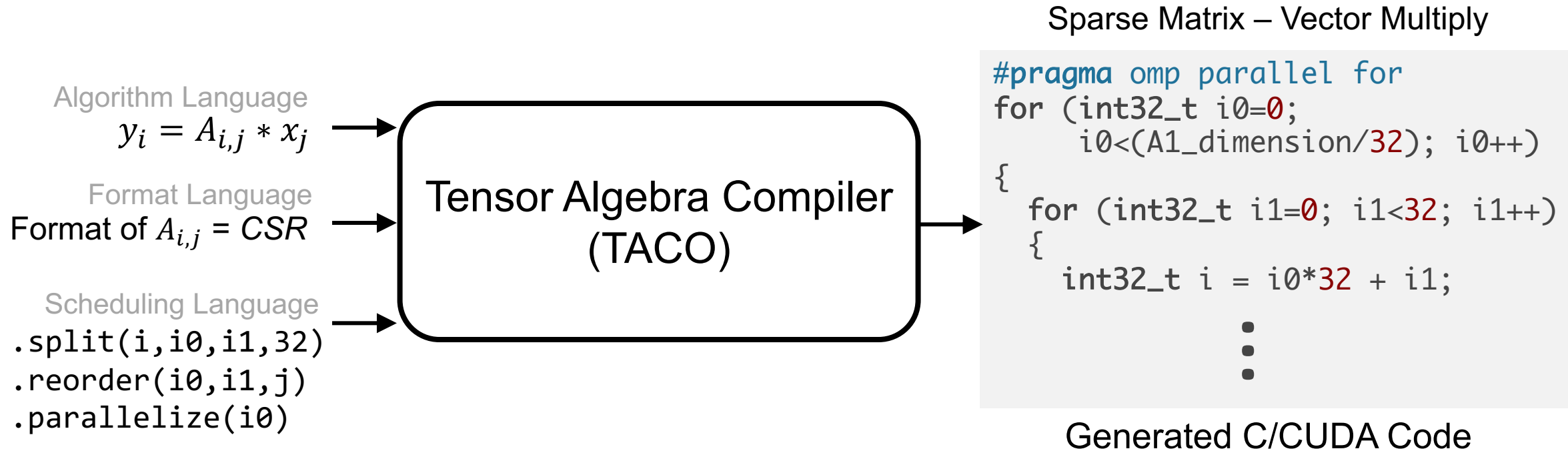
Background : Tensor Compiler (TACO)



Background : Tensor Compiler (TACO)

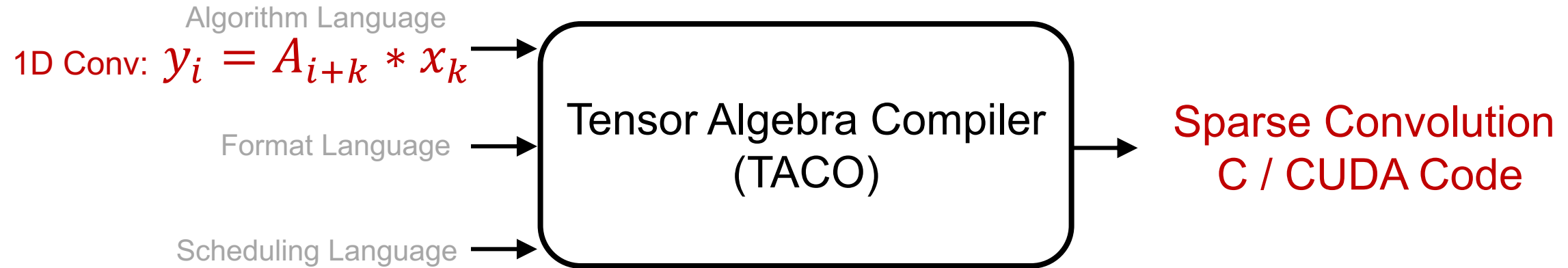


Background : Tensor Compiler (TACO)

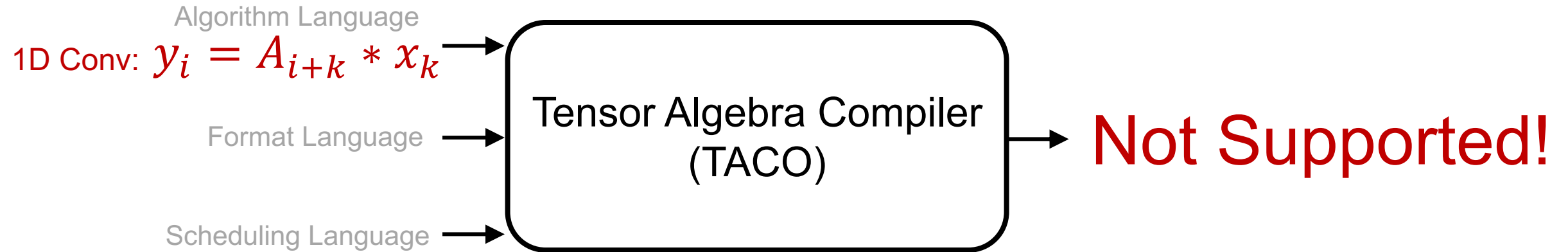


Background : Tensor Compiler (TACO)

What we want : **Sparse Convolution Kernel**

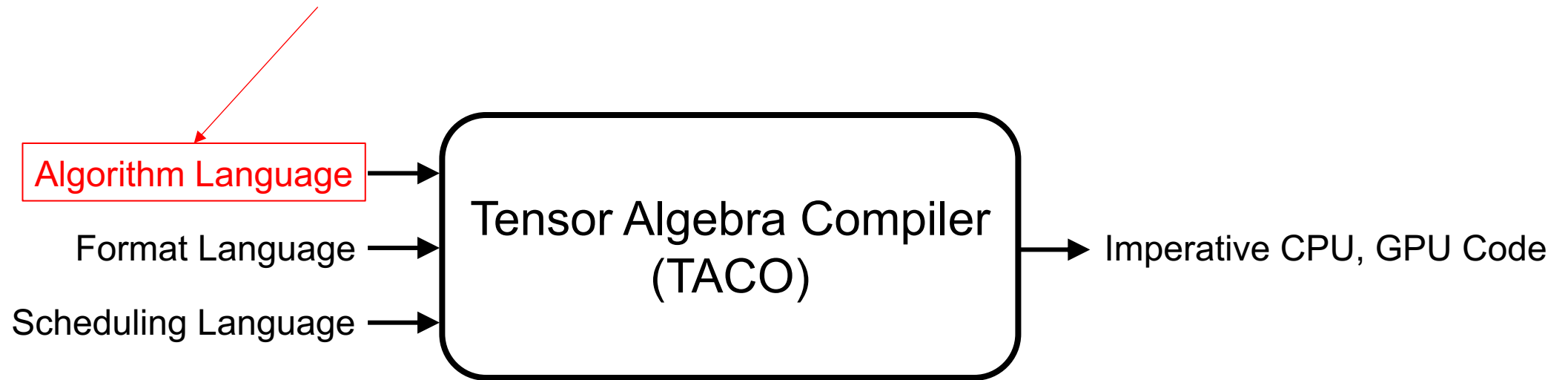


Background : Tensor Compiler (TACO)

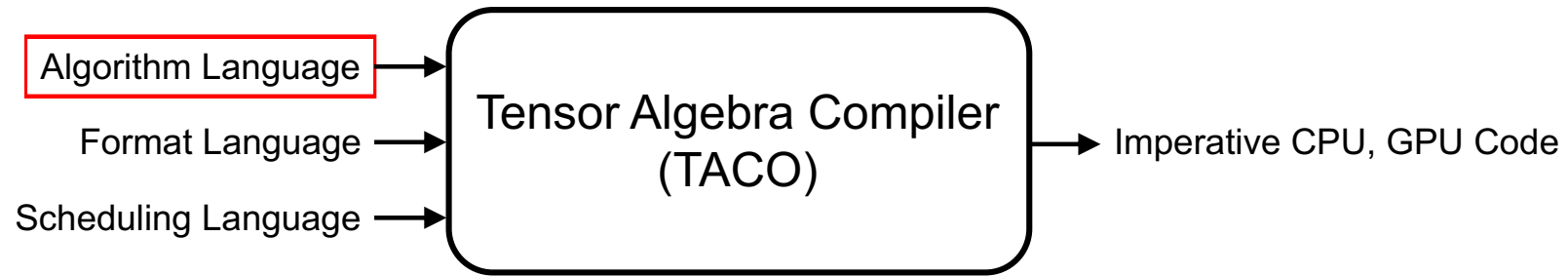


Why can't TACO support convolutions?

TACO's Algorithm Language does not accept "Affine Index"



Why can't TACO support convolutions?

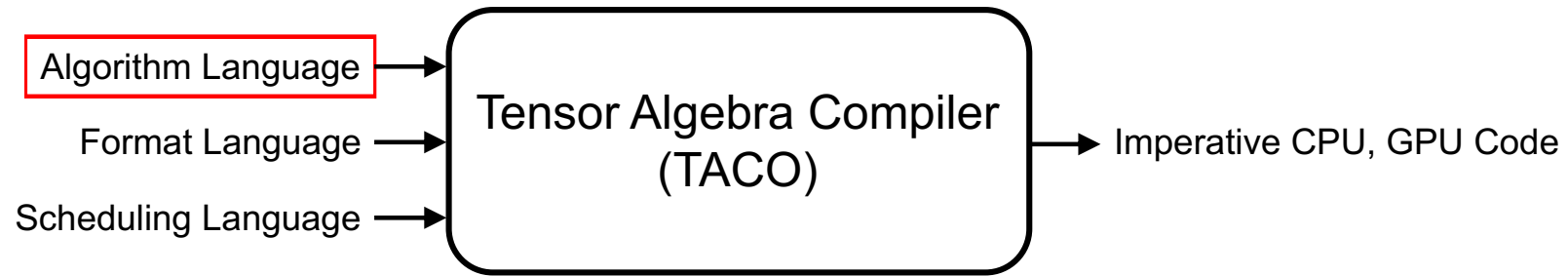


What TACO
can support



What TACO
cannot support

Why can't TACO support convolutions?



What TACO
can support

$$C_i = A_i * B_i \text{ (Element-wise Mul)}$$

$$C_i = A_{2i+1} * B_{3i-1}$$

$$C_{i,j} = A_{i,k} * B_{j,k} \text{ (MatMul)}$$

**Single Variable Affine Expression
(SVAE)**



What TACO
cannot support

Why can't TACO support convolutions?



What TACO can support

$$C_i = A_i * B_i \text{ (Element-wise Mul)}$$

$$C_i = A_{2i+1} * B_{3i-1}$$

Single Variable Affine Expression (SVAE)

$$C_{i,j} = A_{i,k} * B_{j,k} \text{ (MatMul)}$$



What TACO cannot support

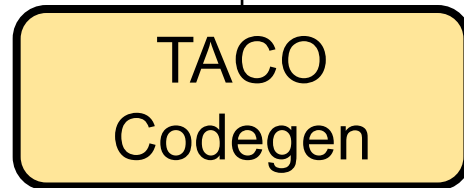
$$C_i = A_{\underline{i+k}} * B_k \text{ (1DConv)}$$

Multiple Variable Affine Expression (MVAE)

multiple variables (i and k)

Lowering Technique

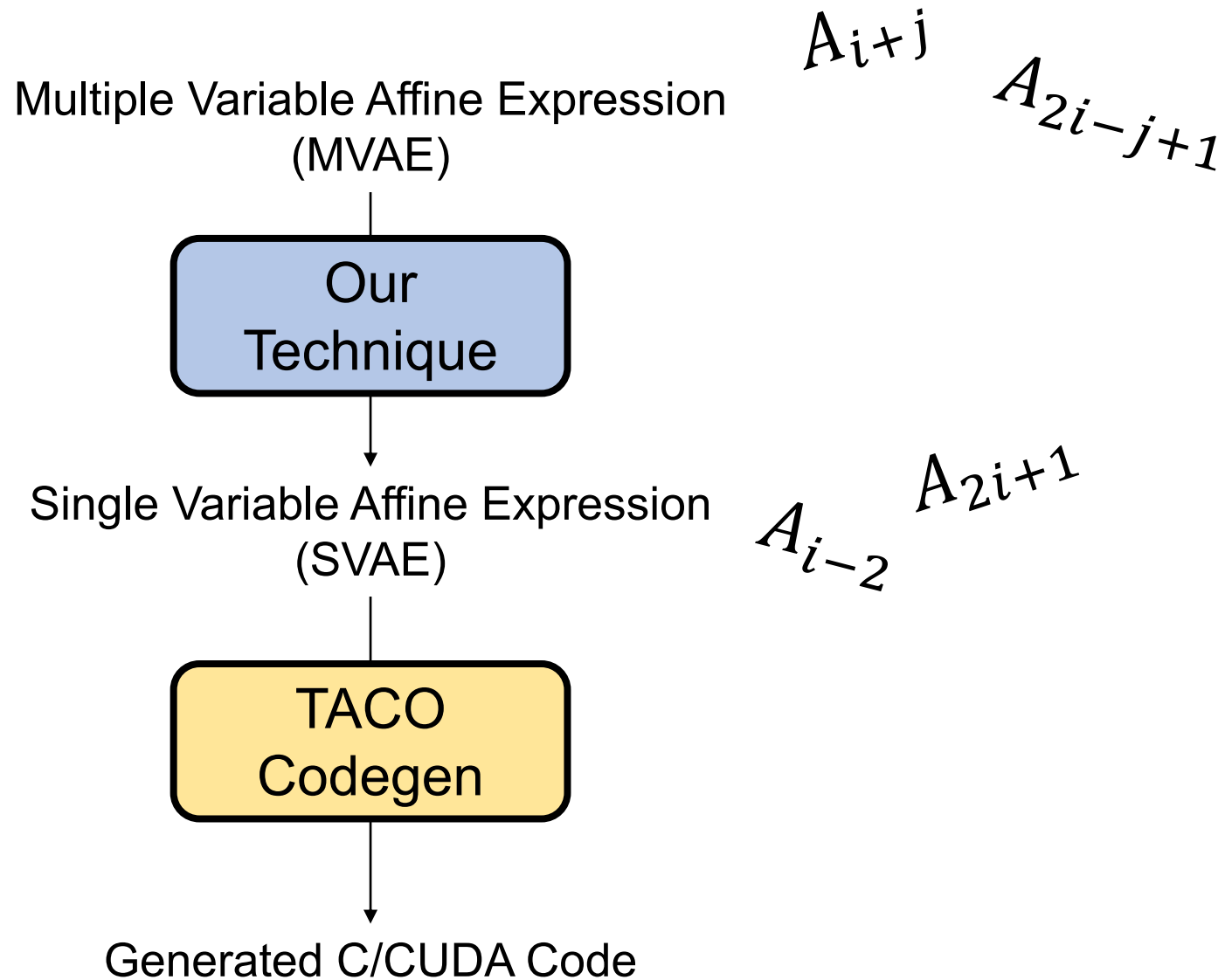
Single Variable Affine Expression
(SVAE)



Generated C/CUDA Code

A_{i-2} A_{2i+1}

Lowering Technique



Lowering Technique

Multiple Variable Affine Expression
(MVAE)

Our
Technique

Single Variable Affine Expression
(SVAE)

TACO
Codegen

Generated C/CUDA Code

MVAE = Stride * Base-variable + Offset

$$A_{3i+2j+5} = A_{2*j+(3i+5)}$$

Stride Base-variable Offset

Lowering Technique

Multiple Variable Affine Expression
(MVAE)

Our
Technique

Single Variable Affine Expression
(SVAE)

TACO
Codegen

Generated C/CUDA Code

MVAE = Stride * Base-variable + Offset

$$A_{3i+2j+5} = A_{2*j+(3i+5)}$$

Stride Base-variable Offset

```
0: For i:  
1:   offset = 3*i+5  
2:   For j: //Base-variable j  
3:     access A[2*j + offset]
```

Lowering Technique

Multiple Variable Affine Expression
(MVAE)

Our
Technique

Single Variable Affine Expression
(SVAE)

TACO
Codegen

Generated C/CUDA Code

MVAE = Stride * Base-variable + Offset

$$A_{3i+2j+5} = A_{2*j+(3i+5)}$$

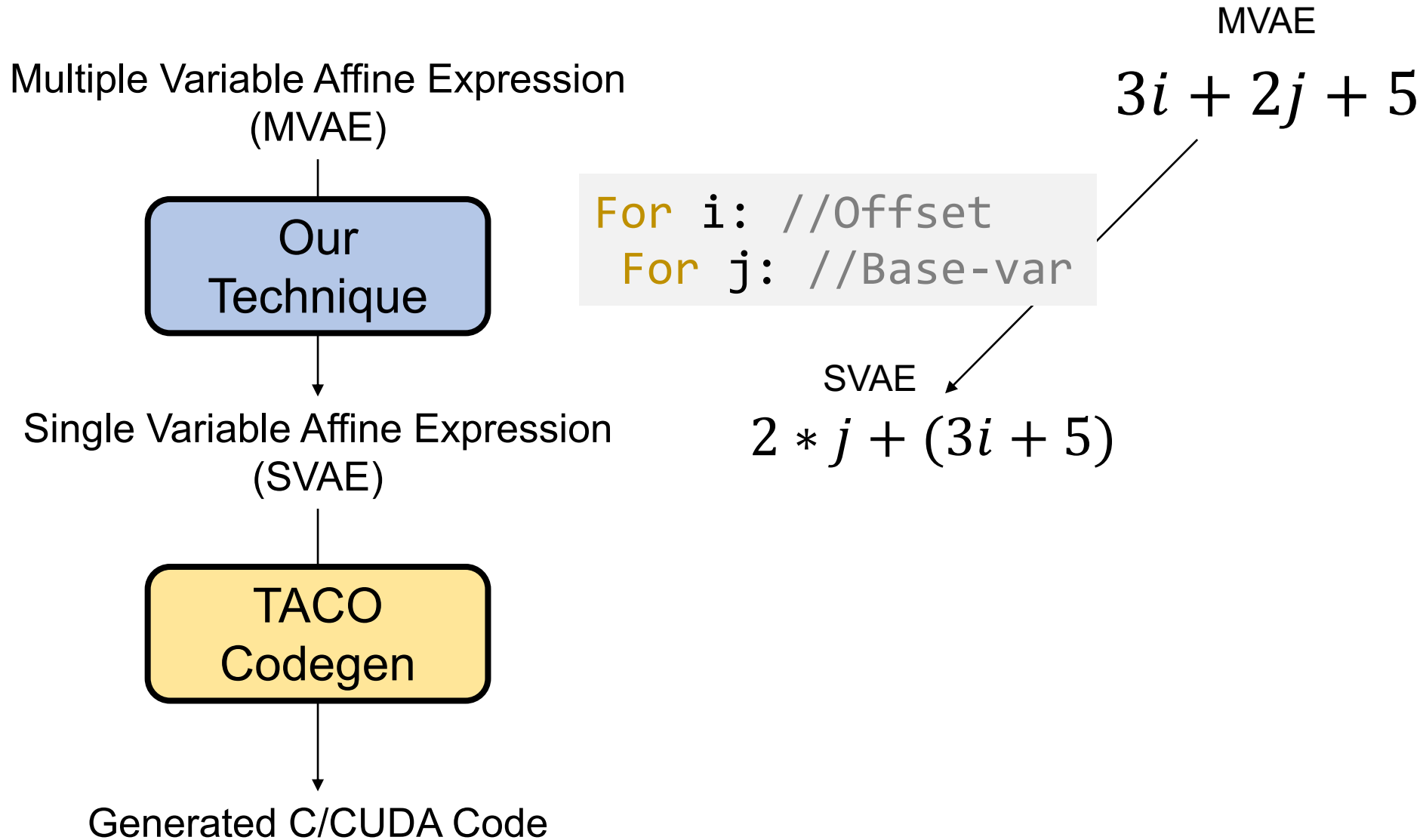
Stride Base-variable Offset

||

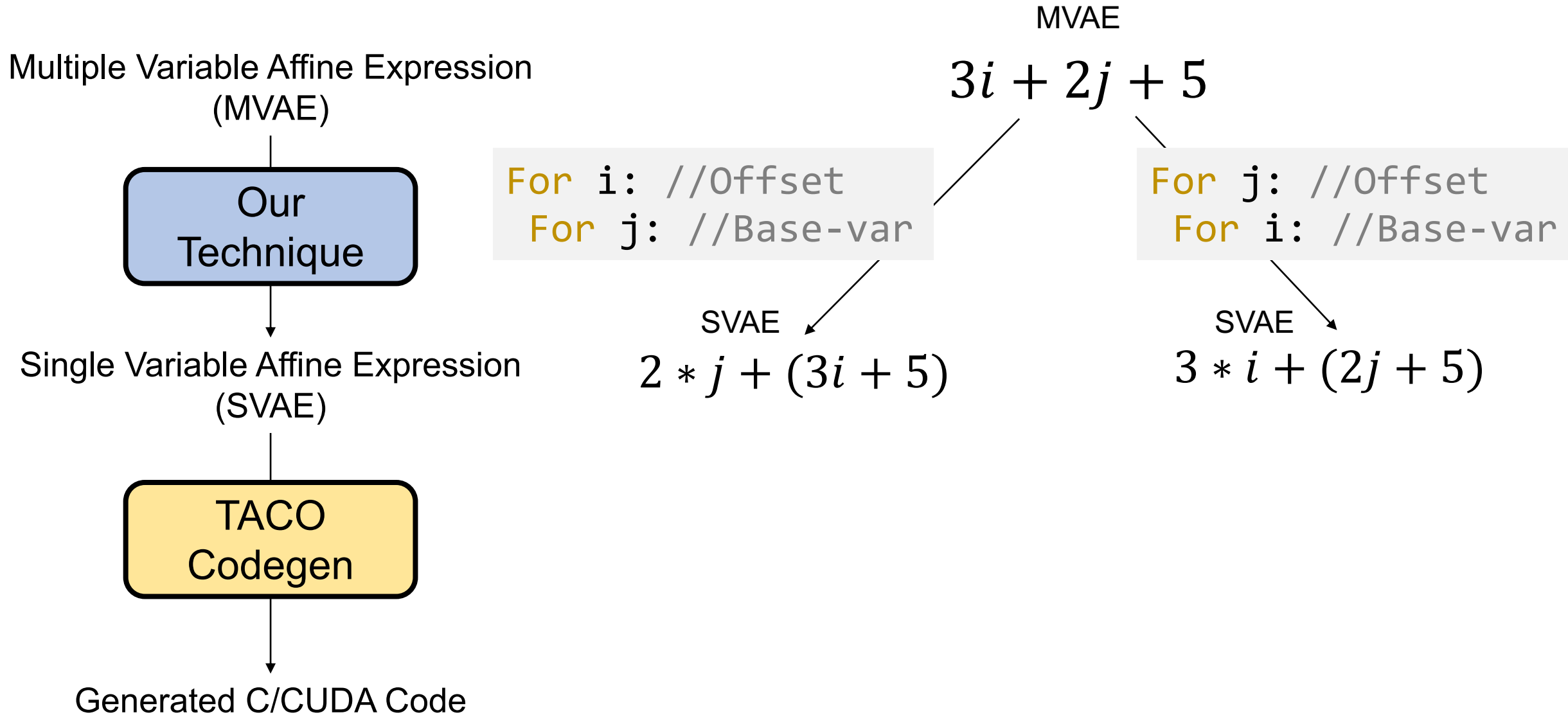
$$A_{3*i+(2j+5)}$$

Stride Base-variable Offset

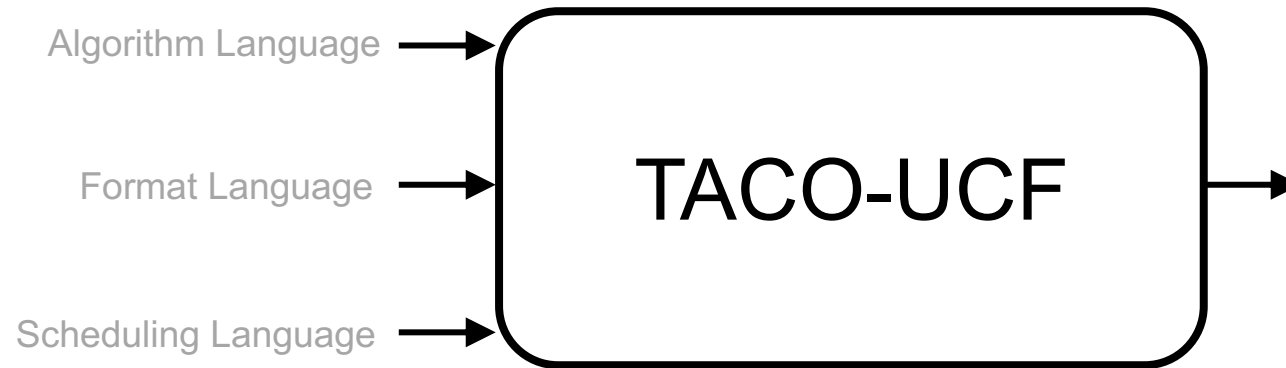
Lowering Technique



Lowering Technique



Sparse Convolution w/ UCF



Sparse Convolution w/ UCF

2D Conv

$$O_{p,q,m} = I_{p+r,q+s,c} * F_{m,r,s,c}$$

Algorithm Language

Format Language

Scheduling Language

TACO-UCF

```
158
159 #pragma omp parallel for schedule(static)
160 for (int32_t pOut = 0; pOut < ((Out2_dimension * Out3_dimens
161     Out_vals[pOut] = 0.0;
162 }
163
164 for (int32_t n = 0; n < 1; n++) {
165     #pragma omp parallel for schedule(runtime)
166     for (int32_t p = 0; p < 56; p++) {
167         for (int32_t r = 0; r < 1; r++) {
168             for (int32_t sF = F2_pos[r]; sF < F2_pos[(r + 1)]; sF+
169                 int32_t s = F2_crd[sF];
170                 for (int32_t c = 0; c < 64; c++) {
171                     int32_t cIn = n * In2_dimension + c;
172                     int32_t rIn = cIn * In3_dimension + (r + 2 * p);
173                     int32_t cF = sF * F2_dimension + c;
174                     for (int32_t mF = F4_pos[cF]; mF < F4_pos[(cF + 1)
175                         int32_t m = F4_crd[mF];
176                         int32_t mOut = n * Out2_dimension + m;
177                         int32_t pOut0 = mOut * Out3_dimension + p;
178
```

Sparse CPU/GPU Convolution Kernel

Evaluation

CPU : Intel Xeon E5-2680 v3 (24 threads)

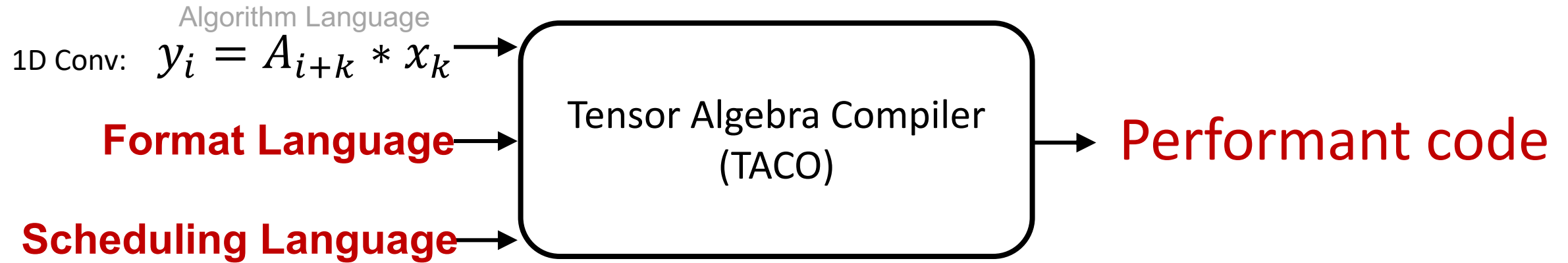
GPU : Nvidia V100

1. Importance of Format and Schedule.

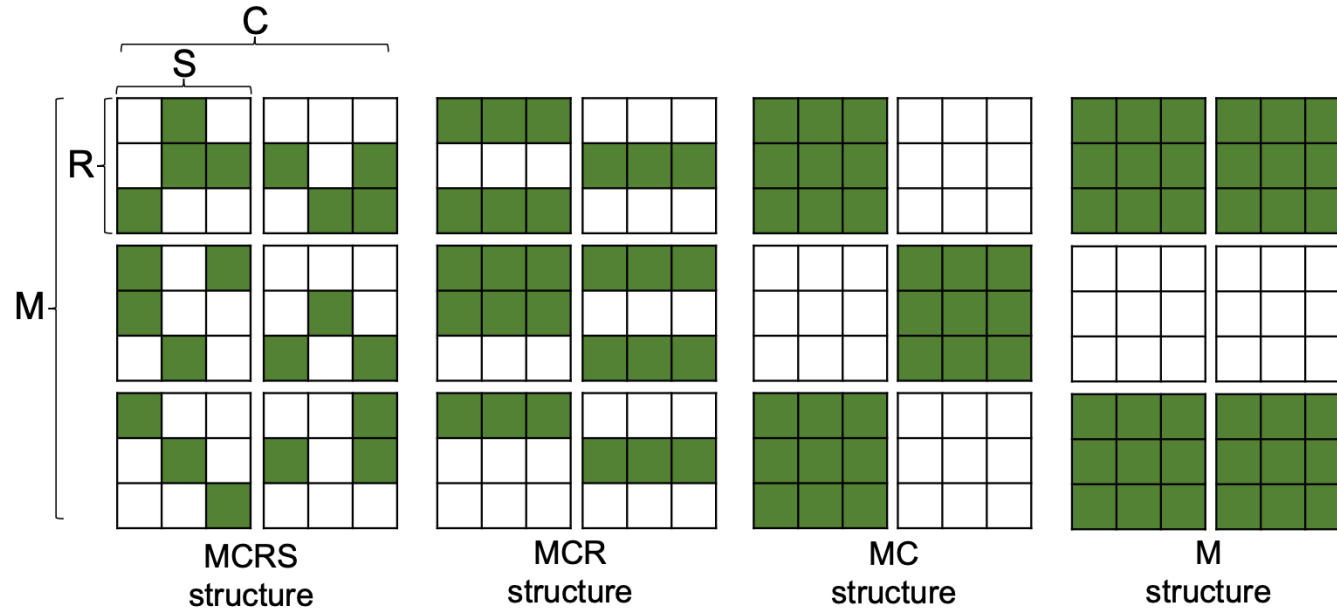
2. Performance Comparison

- Filter Sparse Convolution
- Submanifold Sparse Convolution

Format and Schedule Matters.

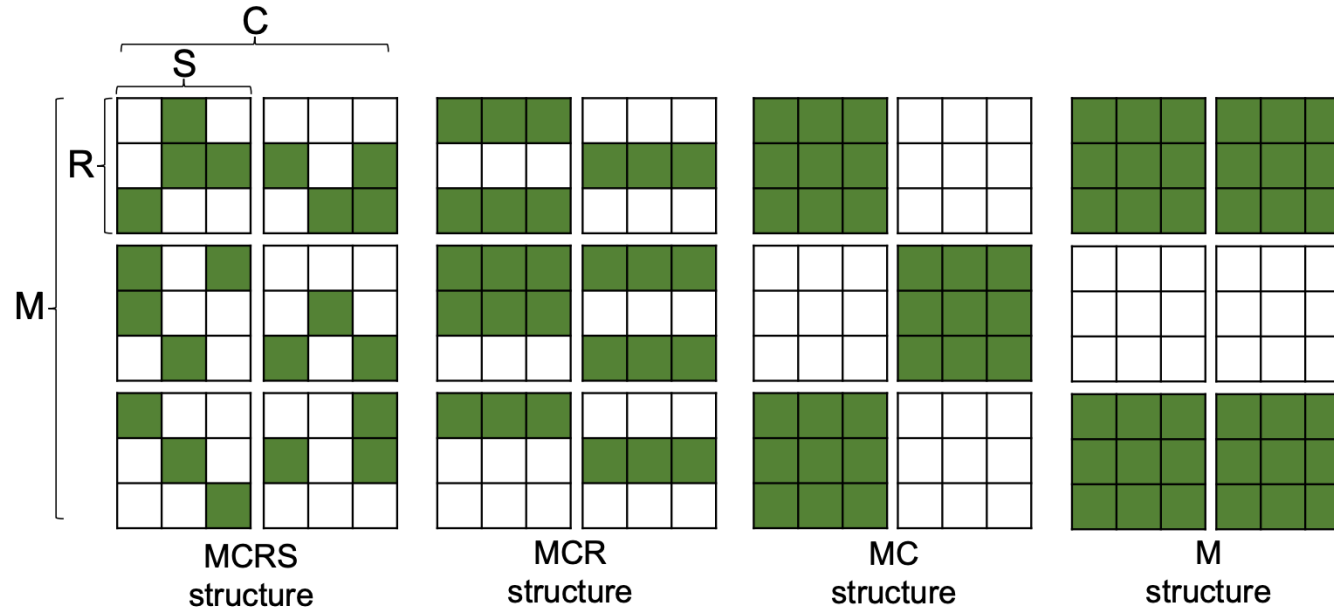


Format and Schedule Matters.



More Structured Pattern

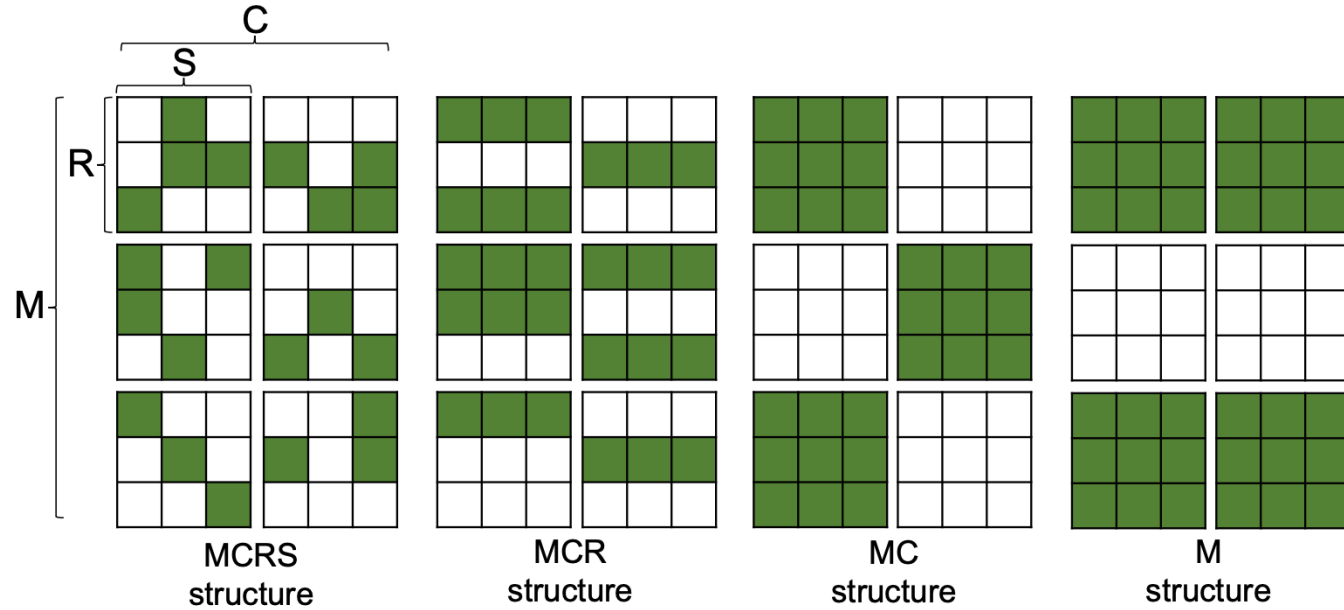
Format and Schedule Matters.



Memory Saving over Uncompressed(Dense) Representation

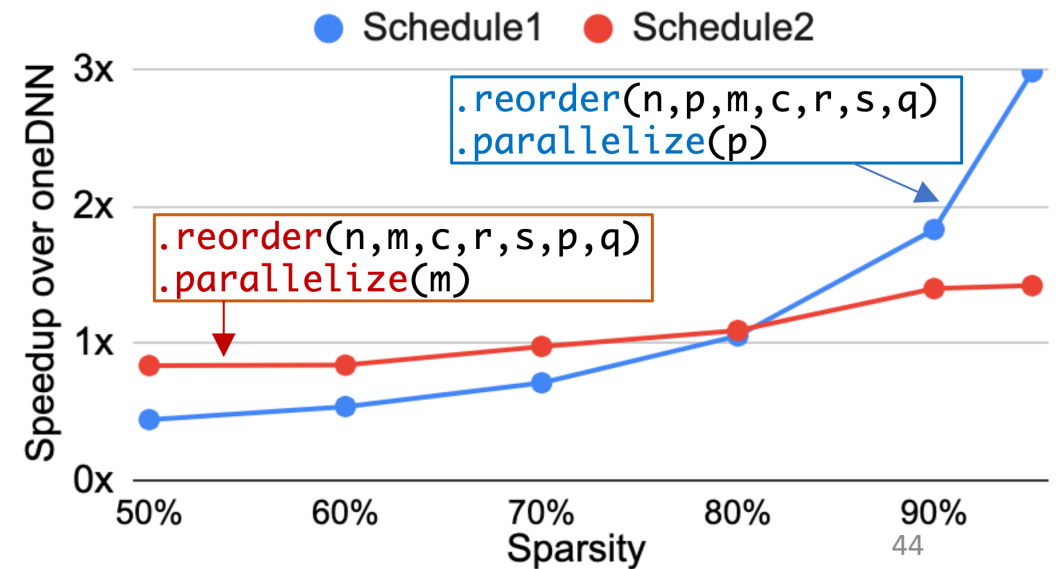
80% Sparsity	MCRS Structure (Unstructured)	MCR Structure	MC Structure	M Structure
Format1	1.08x	1.55x	1.72x	1.84x
Format2	1.17x	2.64x	3.17x	3.41x
Format3	1.03x	1.83x	4.42x	4.78x
Format4	0.99x	0.99x	0.99x	5x

Format and Schedule Matters.



Memory Saving over Uncompressed(Dense) Representation

80% Sparsity	MCRS Structure (Unstructured)	MCR Structure	MC Structure	M Structure
Format1	1.08x	1.55x	1.72x	1.84x
Format2	1.17x	2.64x	3.17x	3.41x
Format3	1.03x	1.83x	4.42x	4.78x
Format4	0.99x	0.99x	0.99x	5x



Evaluation

CPU : Intel Xeon E5-2680 v3 (24 threads)

GPU : Nvidia V100

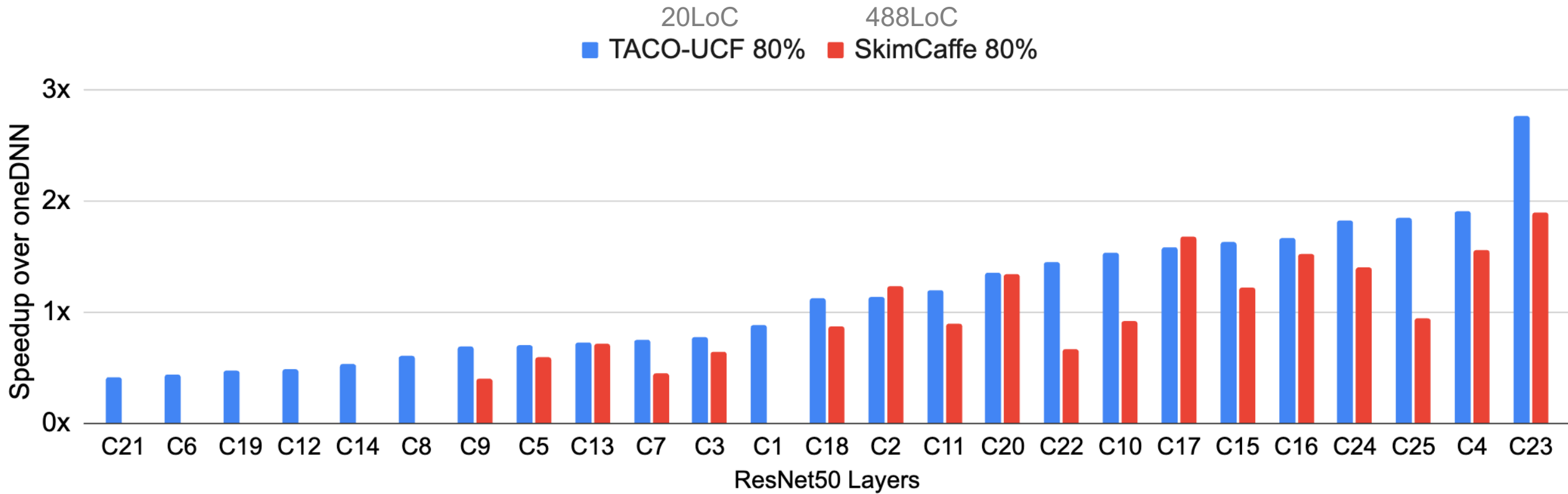
1. Importance of Format and Schedule.

2. Performance Comparison

- **Filter Sparse Convolution**
- Submanifold Sparse Convolution

Evaluation – Filter Sparse Convolution

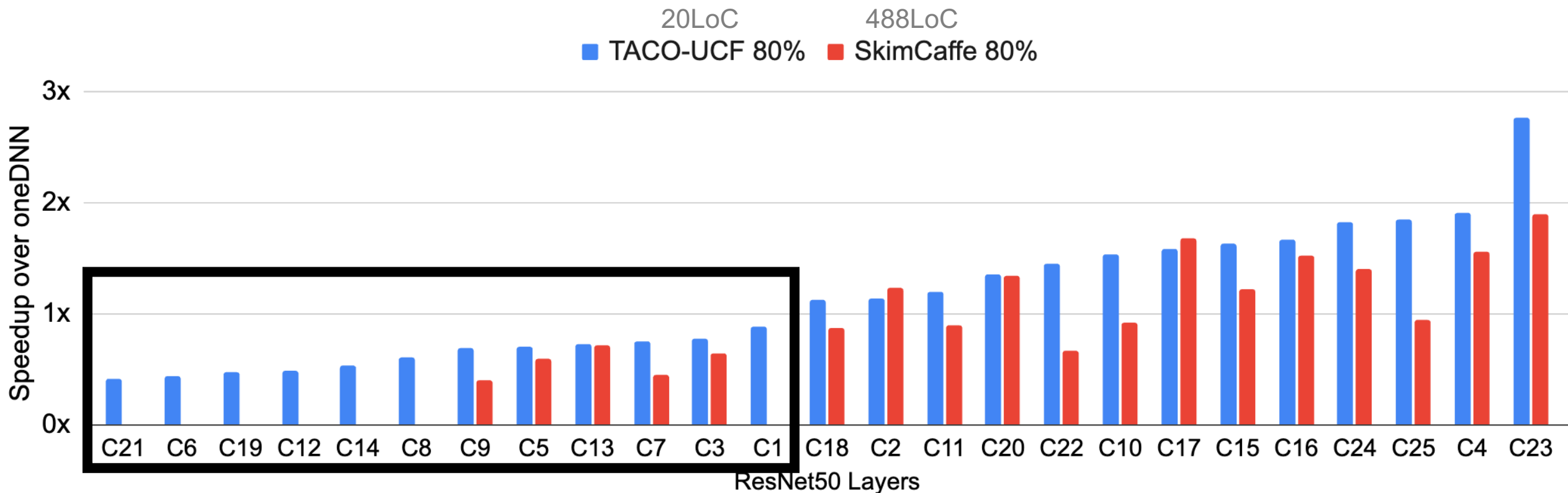
80% sparsity - pruned ResNet50
CPU : 24-core Intel Xeon



Evaluation – Filter Sparse Convolution

80% sparsity - pruned ResNet50

CPU : 24-core Intel Xeon



Not all layers can benefit from pruning!

Evaluation – Filter Sparse Convolution

ResNet50 on Nvidia V100 GPU

Pruning Sparsity	80%	91%	96%	98%
cuDNN	1.0×	1.0×	1.0×	1.0×
Escort	0.78×	1.09×	1.35×	1.49×
TACO-UCF	1.08×	1.61×	2.15×	2.57×

TACO-UCF > cuDNN at 80% Sparsity

Escort > cuDNN at 91% Sparsity

Evaluation

CPU : Intel Xeon E5-2680 v3 (24 threads)

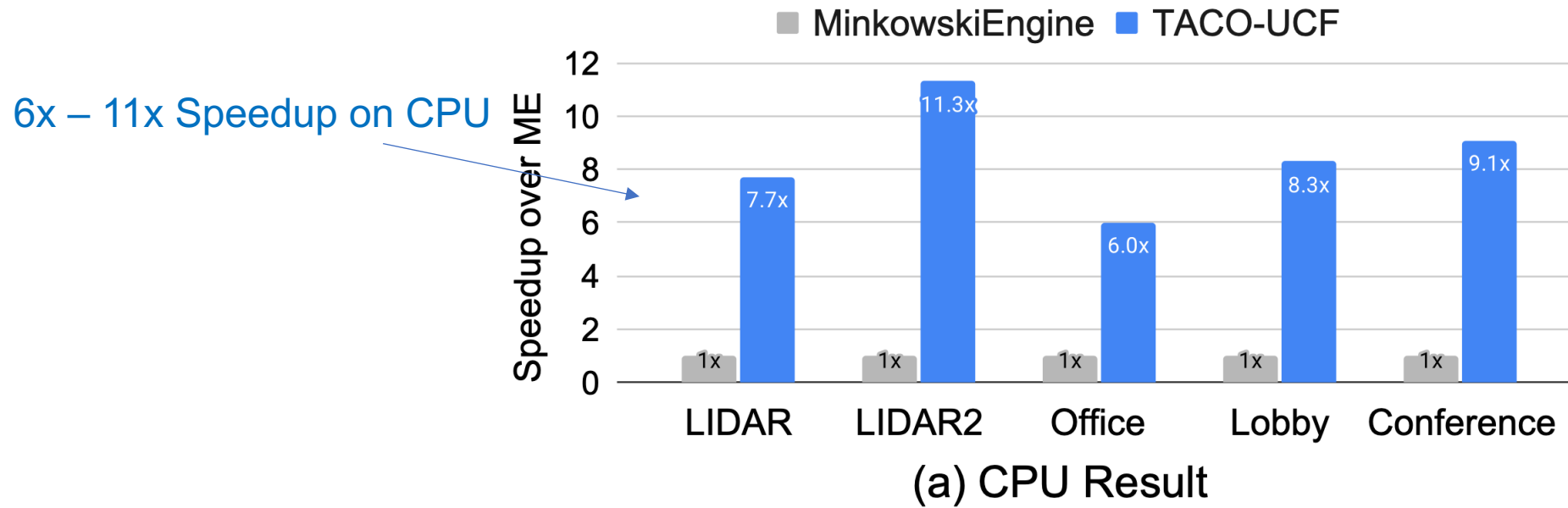
GPU : Nvidia V100

1. Importance of Format and Schedule.

2. Performance Comparison

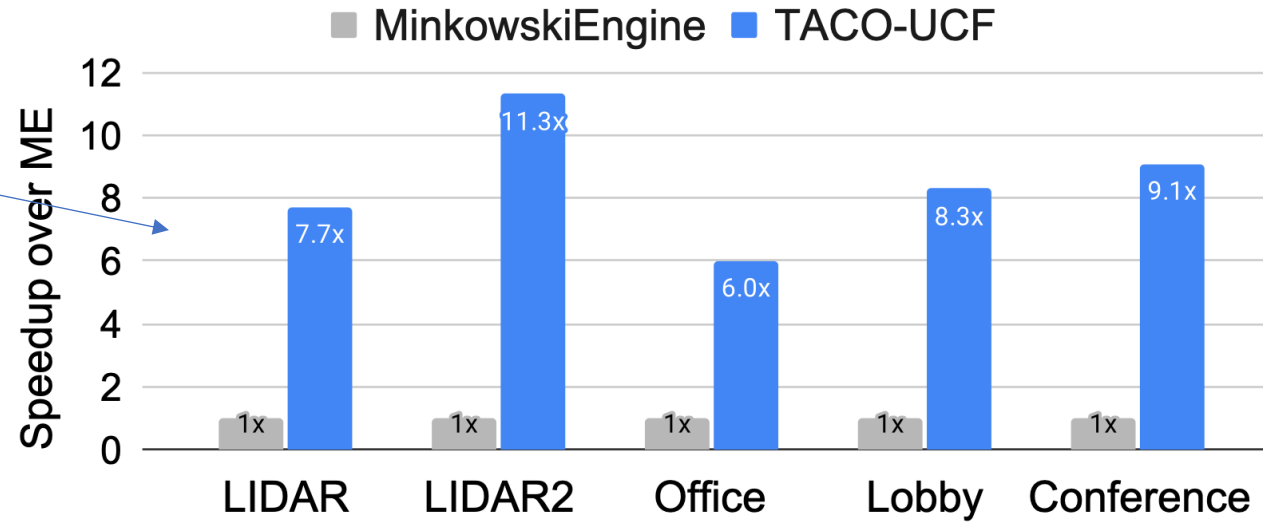
- Filter Sparse Convolution
- **Submanifold Sparse Convolution**

Evaluation – Submanifold Sparse Convolution

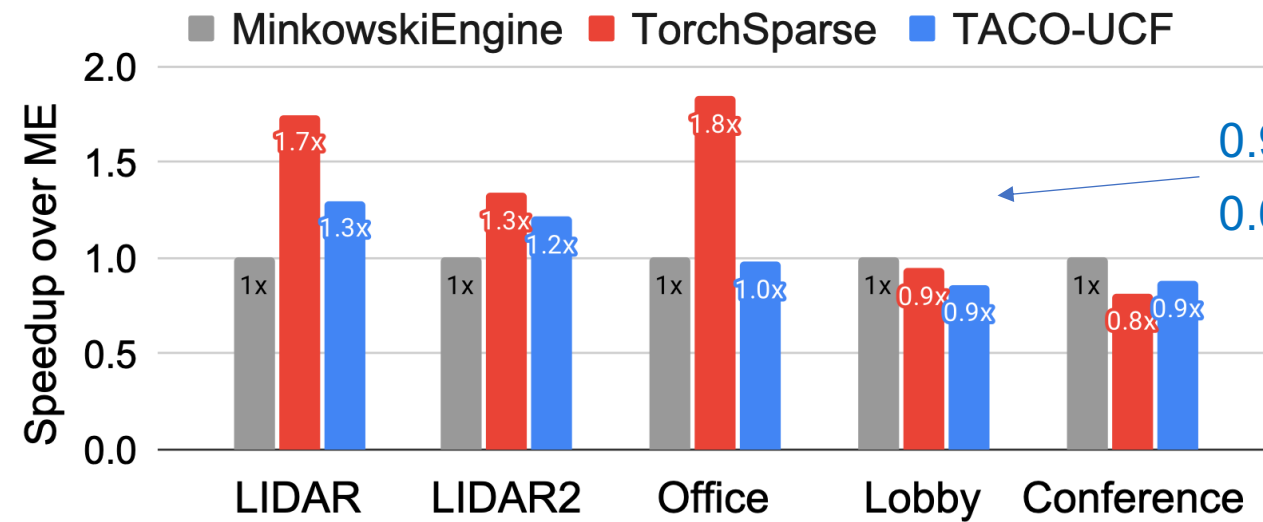


Evaluation – Submanifold Sparse Convolution

6x – 11x Speedup on CPU



(a) CPU Result



0.9x – 1.3x over MinkowskiEngine

0.6x – 1.1x over TorchSparse

(b) GPU Result

Evaluation

CPU : Intel Xeon E5-2680 v3 (24 threads)

GPU : Nvidia V100

1. Importance of Format and Schedule.

2. Performance Comparison

Library Name	Filter Sparse Conv		Submanifold Sparse Conv	
	CPU	GPU	CPU	GPU
SkimCaffe				
Escort				
MinkowskiEngine				
TorchSparse				
Ours (Normalized)				

Evaluation

CPU : Intel Xeon E5-2680 v3 (24 threads)

GPU : Nvidia V100

1. Importance of Format and Schedule.

2. Performance Comparison

Library Name	Filter Sparse Conv		Submanifold Sparse Conv	
	CPU	GPU	CPU	GPU
SkimCaffe	76%	-	-	-
Escort	-	67%	-	-
MinkowskiEngine	-	-	12%	97%
TorchSparse	-	-	< 5%	123%
Ours (Normalized)	100%	100%	100%	100%

1. Better Performance
2. Versatile convolution support
3. Flexible Hardware
4. Less lines of code!

Evaluation

CPU : Intel Xeon E5-2680 v3 (24 threads)

GPU : Nvidia V100

1. Importance of Format and Schedule.

2. Performance Comparison

Library Name	Filter Sparse Conv		Submanifold Sparse Conv		Dual Submanifold Sparse Conv	
	CPU	GPU	CPU	GPU	CPU	GPU
SkimCaffe	76%	-	-	-	Details In Paper!	
Escort	-	67%	-	-		
MinkowskiEngine	-	-	12%	97%		
TorchSparse	-	-	< 5%	123%		
Ours (Normalized)	100%	100%	100%	100%		

Thanks!