
Fusion of Operators of Computational Graphs via Greedy Clustering: The XNNC Experience

Michael Canesche canesche@cadence.com

Vanderson Rosário vrosario@cadence.com

Edson Borin edson@ic.unicamp.br

Fernando Pereira fernando@dcc.ufmg.br

Goal

1. Design and implement a new kernel fusion algorithm

Goal

1. Design and implement a new kernel fusion algorithm
2. Deploy this algorithm on the XNNC Compiler

Goal

1. Design and implement a new kernel fusion algorithm
2. Deploy this algorithm on the XNNC Compiler

Xtensa® Neural Network Compiler (XNNC)

1. Design and implement a new kernel fusion algorithm
2. Deploy this algorithm on the XNNC Compiler

Xtensa® Neural Network Compiler (XNNC)

- Floating-point Convolutional Neural Network → Optimized, fixed-point code
- Cadence Tensilica® processors.

XNNC

1. Automotive Audio DSPs
2. Light Detection and Ranging
3. Consumer Electronics
4. etc



Xtensa® Neural Network Compiler (XNNC)

- Floating-point Convolutional Neural Network → Optimized, fixed-point code
- Cadence Tensilica® processors.

XNNC

1. C++
2. 1.5 Million LoC
3. LLVM (MLIR)



Xtensa® Neural Network Compiler (XNNC)

- Floating-point Convolutional Neural Network → Optimized, fixed-point code
- Cadence Tensilica® processors.

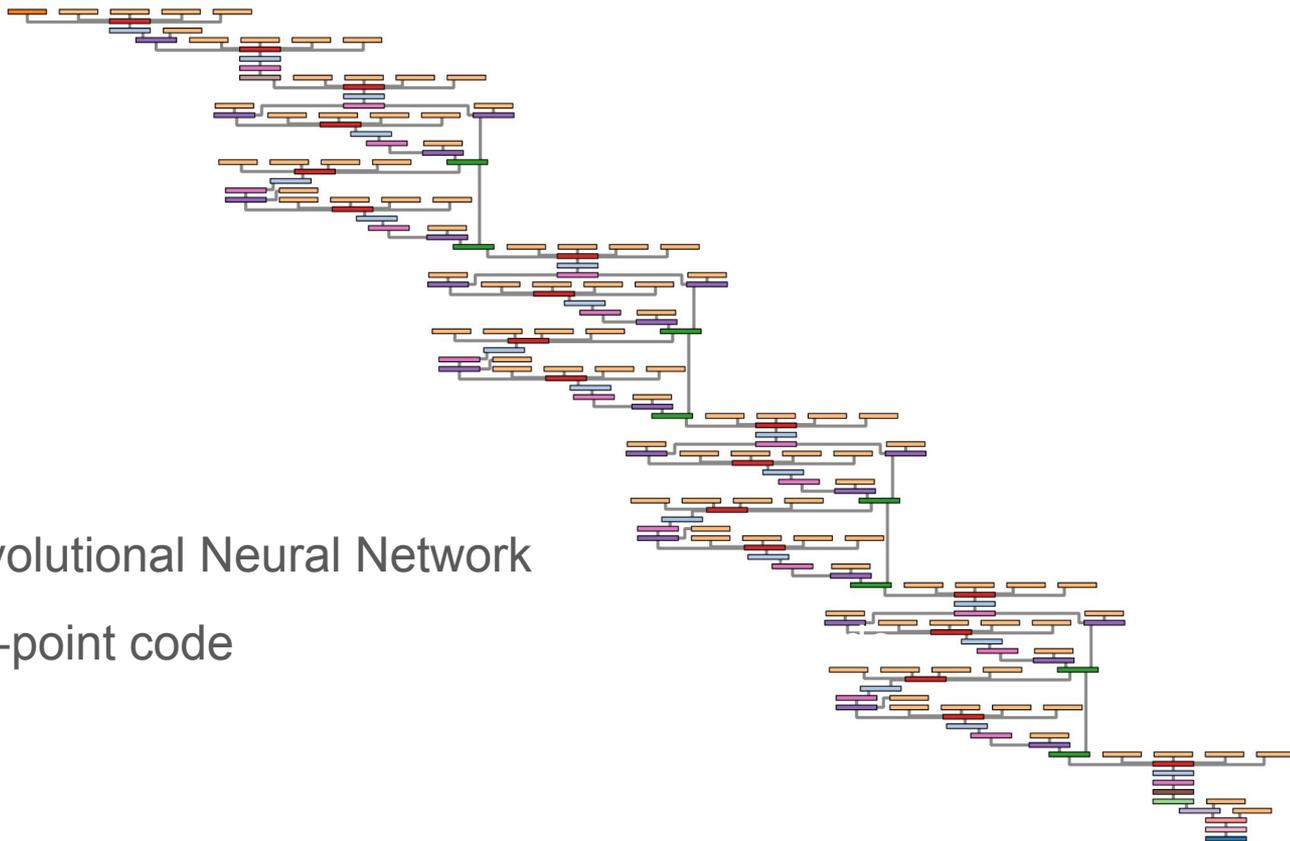
1. C++
2. 1.5 Million LoC
3. LLVM (MLIR)



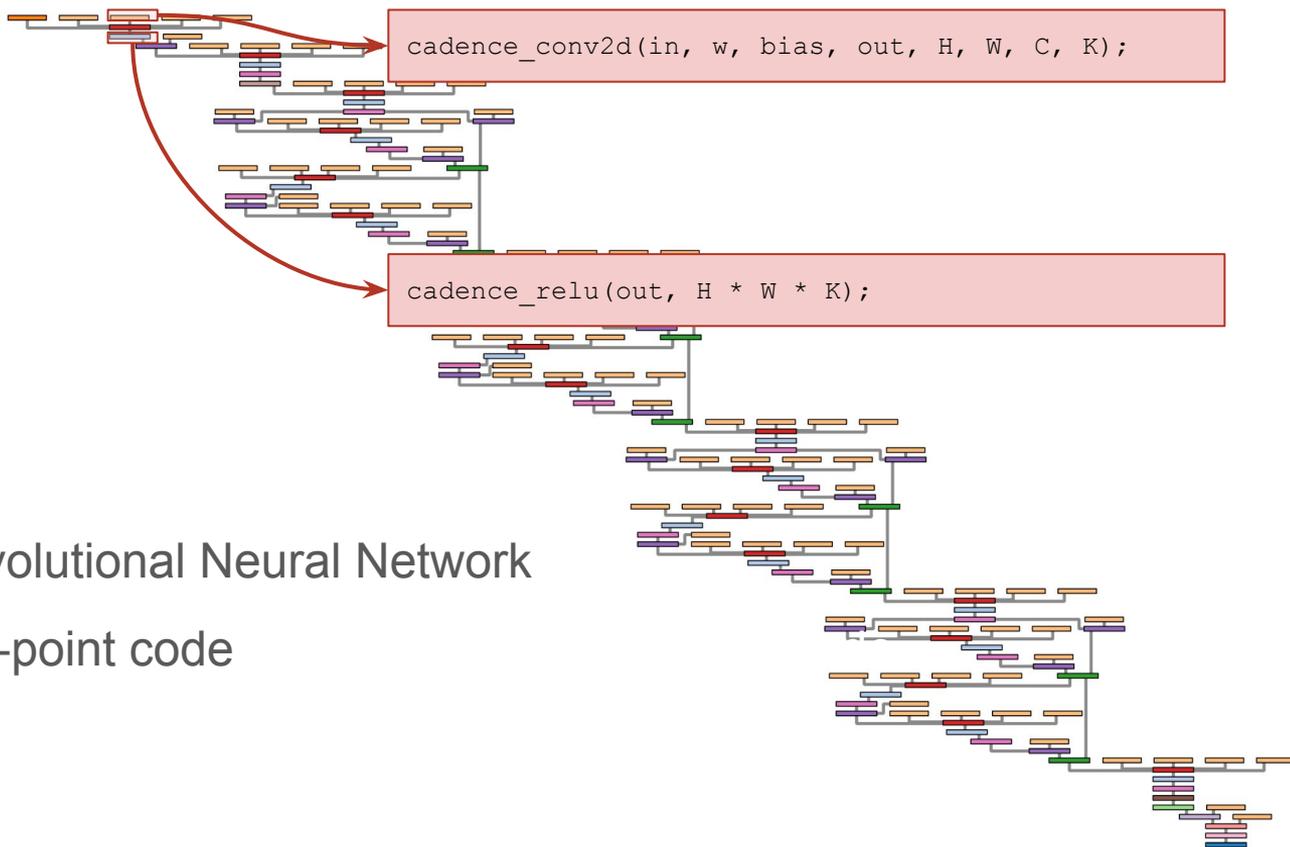
Xtensa® Neural Network Compiler (XNNC)

- Floating-point Convolutional Neural Network → Optimized, fixed-point code
- Cadence Tensilica® processors.

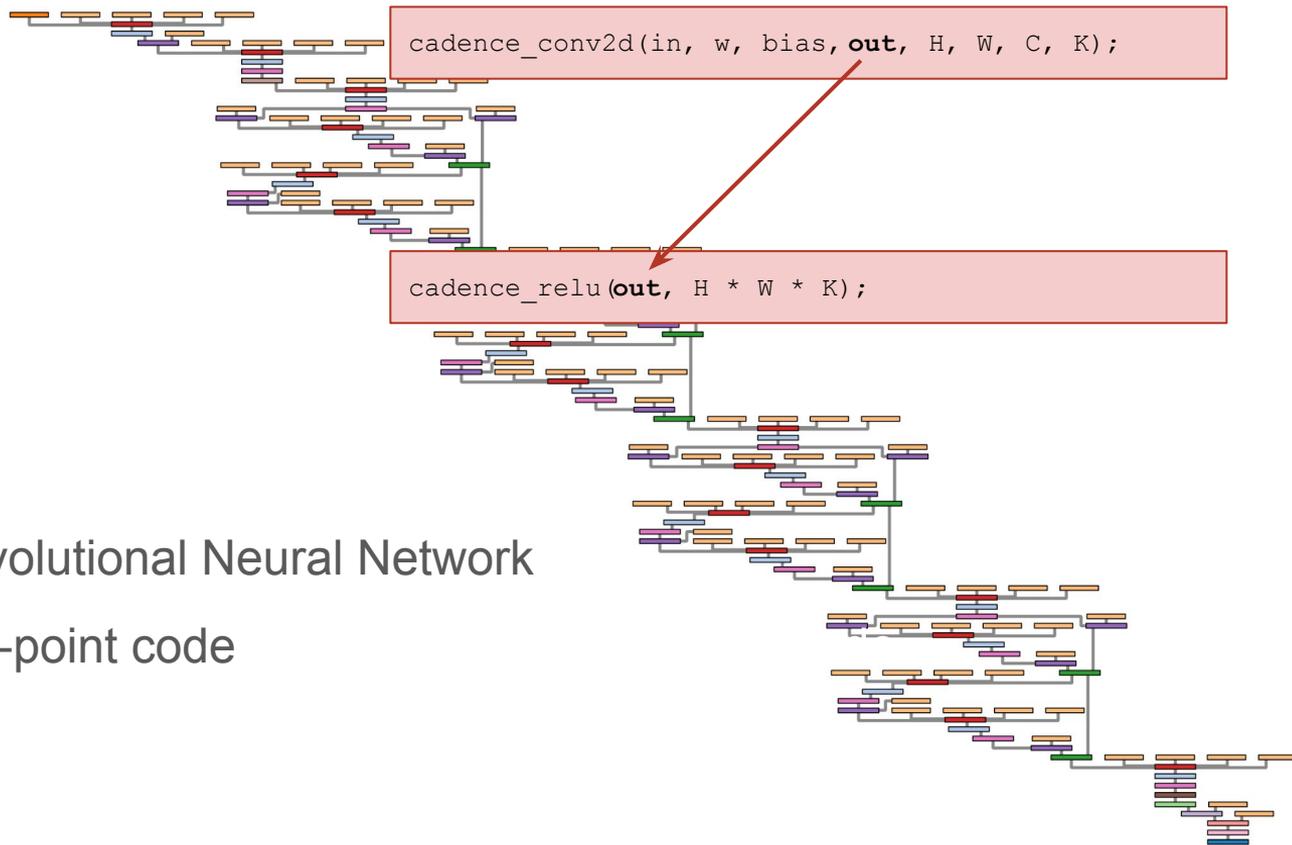
- Floating-point Convolutional Neural Network
 - Optimized, fixed-point code



- Floating-point Convolutional Neural Network
→ Optimized, fixed-point code

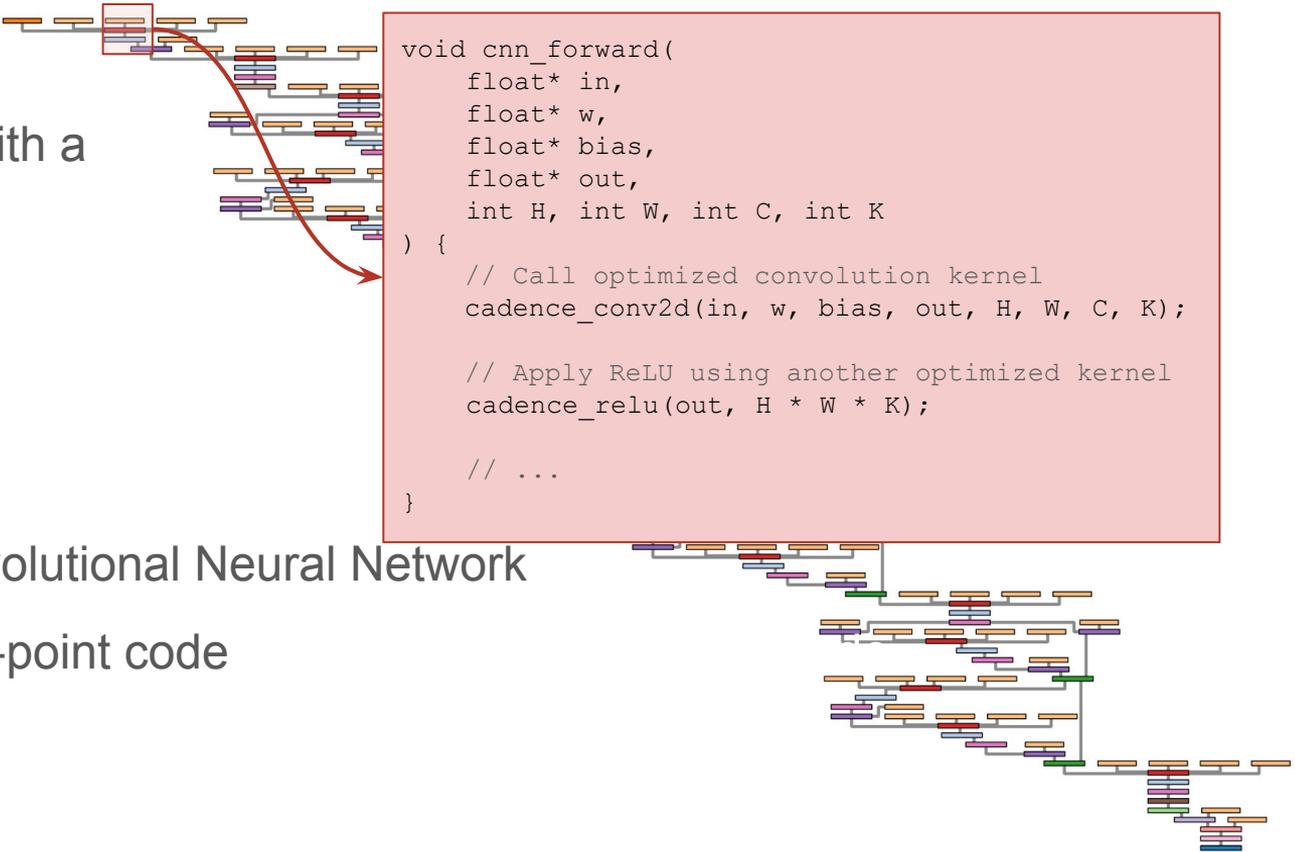


- Floating-point Convolutional Neural Network
→ Optimized, fixed-point code



- Floating-point Convolutional Neural Network
→ Optimized, fixed-point code

1. Pattern matching with a kernel library

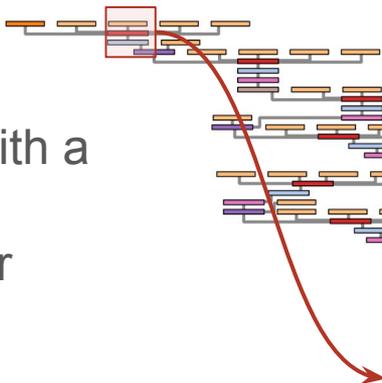


```
void cnn_forward(  
    float* in,  
    float* w,  
    float* bias,  
    float* out,  
    int H, int W, int C, int K  
) {  
    // Call optimized convolution kernel  
    cadence_conv2d(in, w, bias, out, H, W, C, K);  
  
    // Apply ReLU using another optimized kernel  
    cadence_relu(out, H * W * K);  
  
    // ...  
}
```

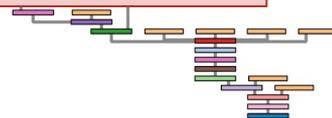
- Floating-point Convolutional Neural Network

→ Optimized, fixed-point code

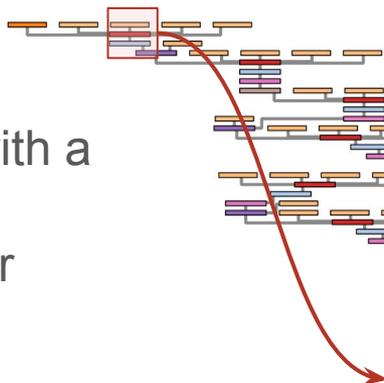
1. Pattern matching with a kernel library
2. Code generation for unmatched parts



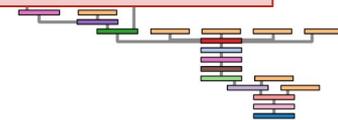
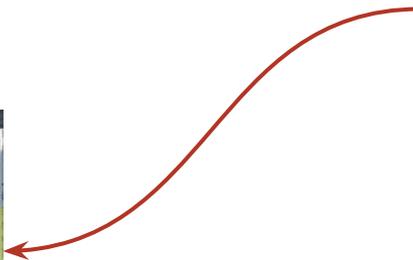
```
void cnn_forward(  
    float* in,  
    float* w,  
    float* bias,  
    float* out,  
    int H, int W, int C, int K  
) {  
    // Call optimized convolution kernel  
    cadence_conv2d(in, w, bias, out, H, W, C, K);  
  
    // Apply ReLU using another optimized kernel  
    cadence_relu(out, H * W * K);  
  
    // Handle any additional logic  
    // (e.g., padding, custom ops)  
    for (int i = 0; i < H * W * K; i++) {  
        if (output[i] < 0) {  
            output[i] = 0;  
        }  
    }  
}
```



1. Pattern matching with a kernel library
2. Code generation for unmatched parts
3. Final Translation to Hardware-Specific Instructions



```
void cnn_forward(  
    float* in,  
    float* w,  
    float* bias,  
    float* out,  
    int H, int W, int C, int K  
) {  
    // Call optimized convolution kernel  
    cadence_conv2d(in, w, bias, out, H, W, C, K);  
  
    // Apply ReLU using another optimized kernel  
    cadence_relu(out, H * W * K);  
  
    // Handle any additional logic  
    // (e.g., padding, custom ops)  
    for (int i = 0; i < H * W * K; i++) {  
        if (output[i] < 0) {  
            output[i] = 0;  
        }  
    }  
}
```



XNNC: "Instruction Selection"

└── Computational graph → kernels + loops

XNNC: "Instruction Selection"

└── Computational graph → kernels + loops

Traditional compiler: Instruction Selection

└── AST → Instructions

■ But, what is "Operator Fusion"?

But, what is "Operator Fusion"?

A CATALOGUE OF OPTIMIZING TRANSFORMATIONS



Frances E. Allen

John Cocke

IBM Thomas J. Watson Research Center
Yorktown Heights

Frances Allen and John Cocke's *"A Catalogue of Optimizing Transformations"* was written in 1971. More than fifty years later, the optimizations they listed in that work still form the backbone of optimizing compilers.



Frances E. Allen receiving the Erna Hamburger Distinguished Lecture Award at the EPFL



John Cocke in his garden. Computer History Museum archive, (c) Louis Fabian Bachrach

Example: Single-Head Attention

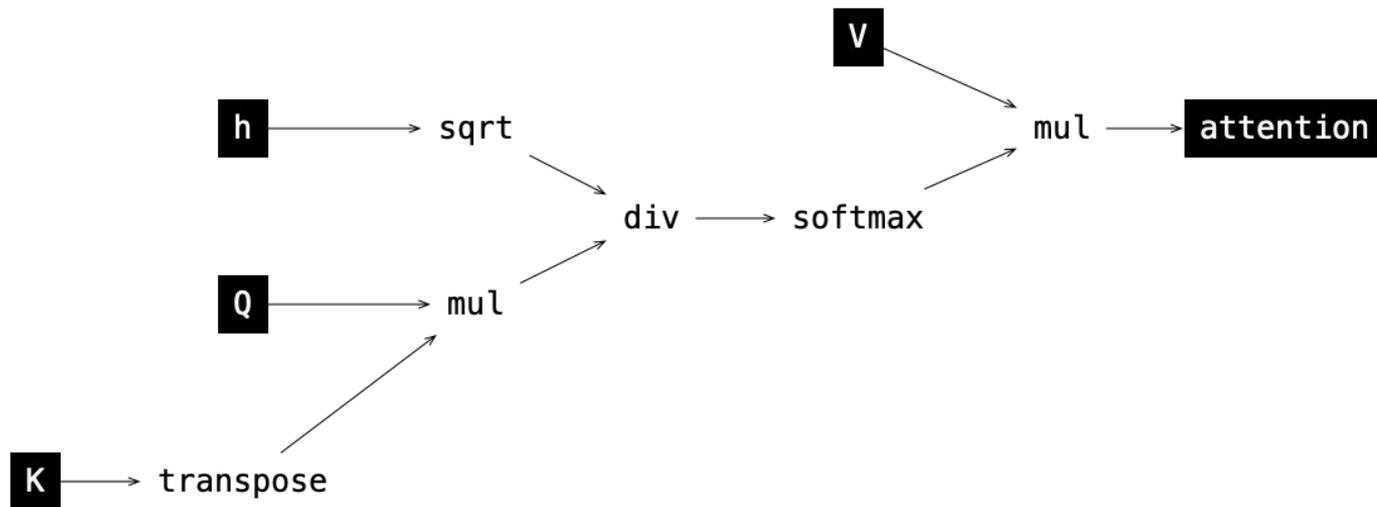
- Core of the Transformer Architecture

Example: Single-Head Attention

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^\top}{\sqrt{d_k}} \right) V$$

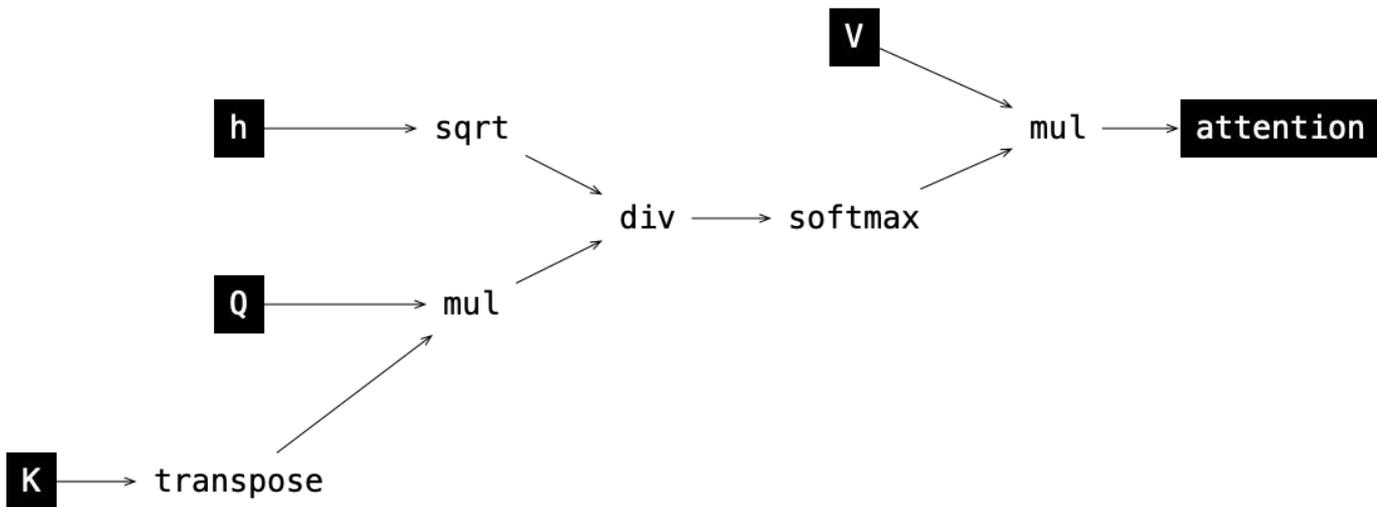
Example: Single-Head Attention

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^{\top}}{\sqrt{d_k}} \right) V$$



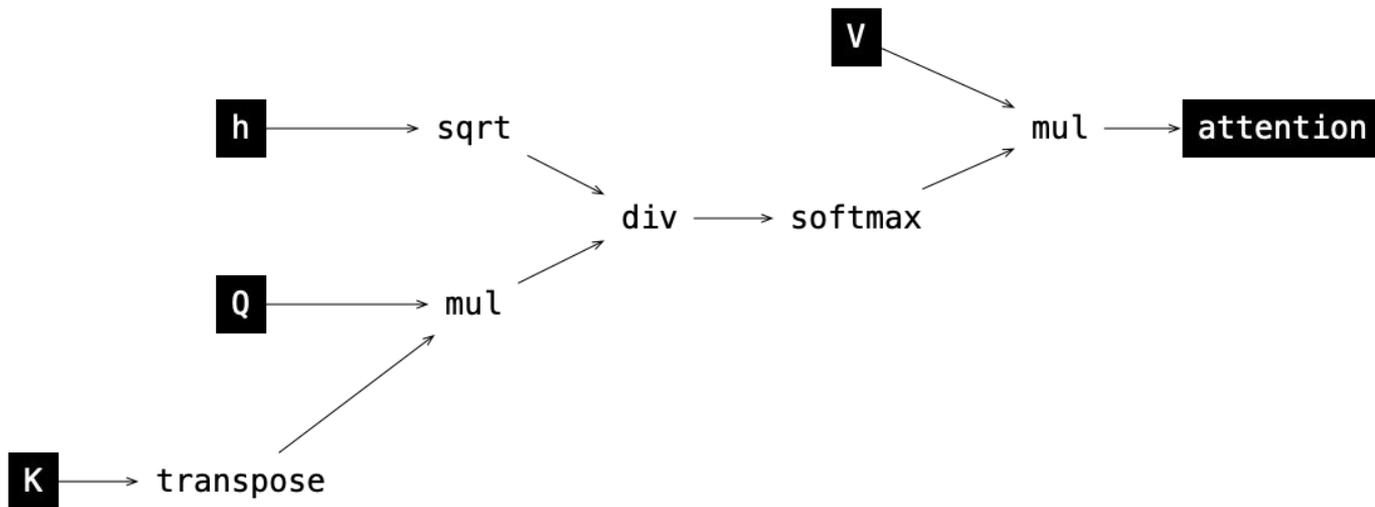
Example: Single-Head Attention

```
scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(h)
attention_weights = torch.softmax(scores, dim=-1)
attention_output = torch.matmul(attention_weights, V)
```



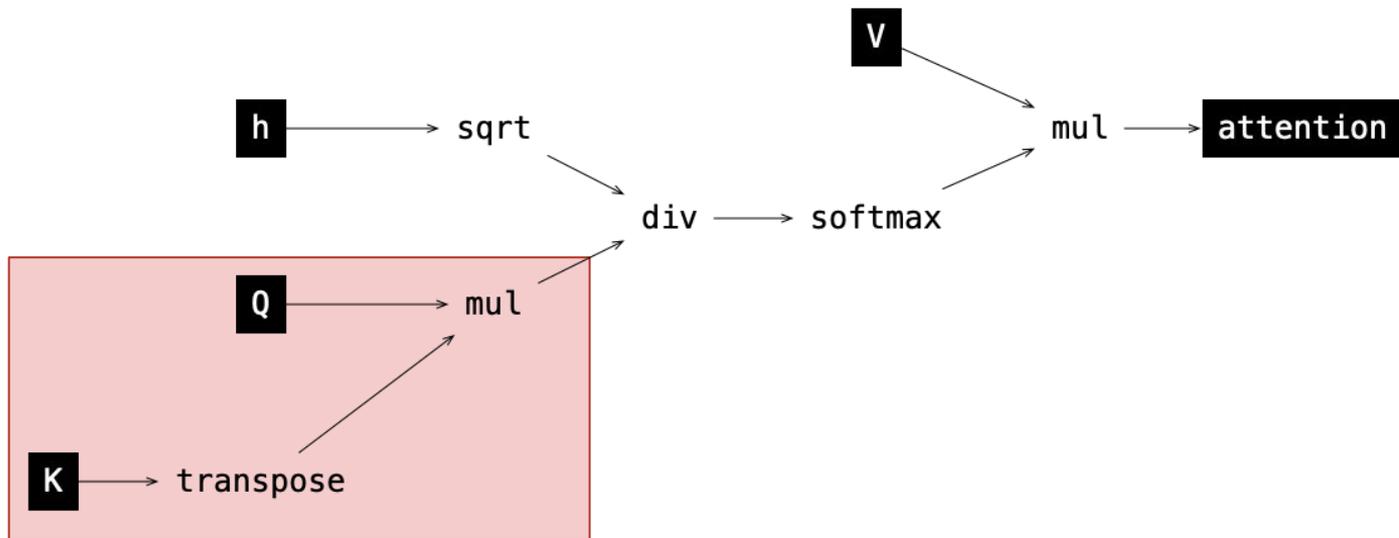
Example: Single-Head Attention

```
scores = torch.matmul(Q, K.transpose(-2, -1)) / math.sqrt(h)
attention_weights = torch.softmax(scores, dim=-1)
attention_output = torch.matmul(attention_weights, V)
```



Example: Single-Head Attention

```
scores = fused_matmul_transpose(Q, K, -2, -1) / math.sqrt(h)
attention_weights = torch.softmax(scores, dim=-1)
attention_output = torch.matmul(attention_weights, V)
```

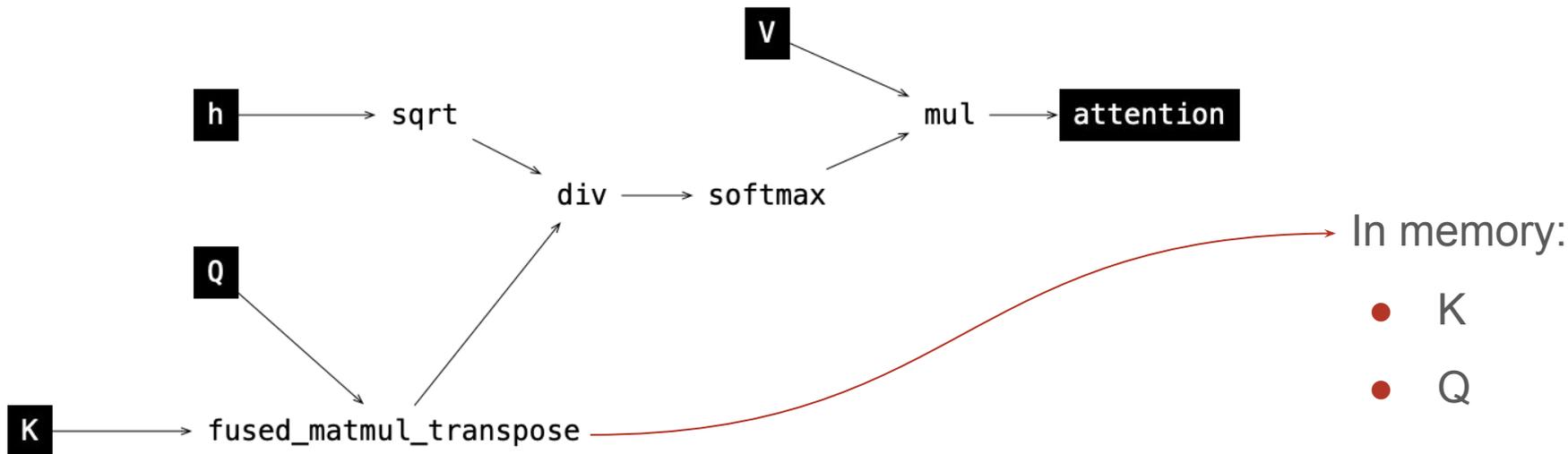


Example: Single-Head Attention

```
scores = fused_matmul_transpose(Q, K, -2, -1) / math.sqrt(h)
```

Compare with previous code:

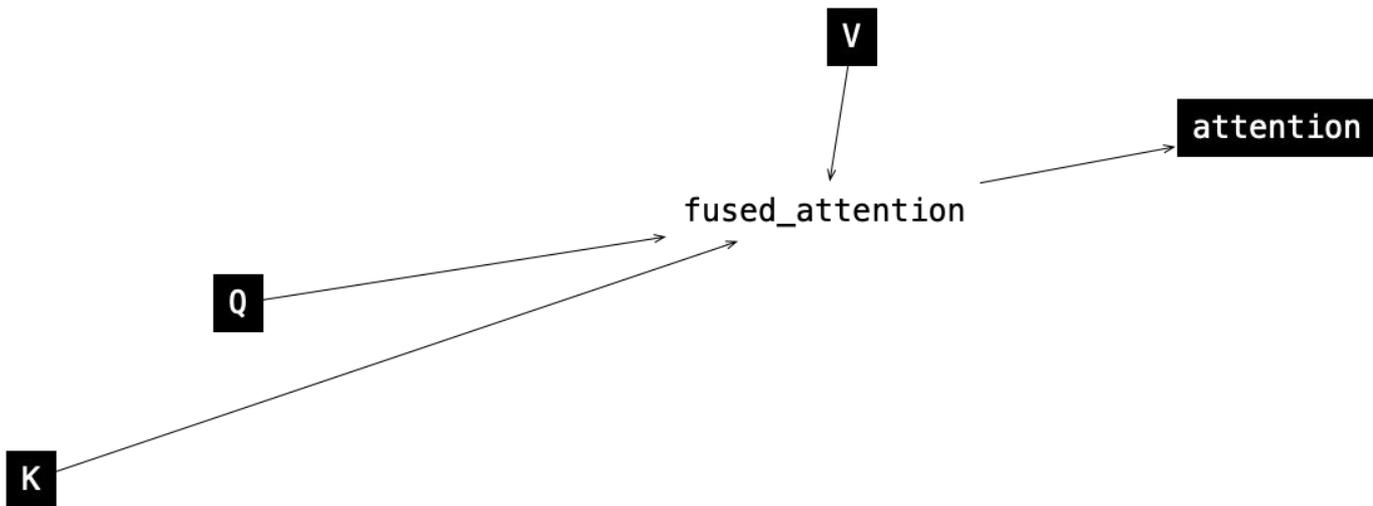
```
scores = torch.matmul(Q, K.transpose()) / math.sqrt(h)
```



Example: Single-Head Attention

```
from flash_attention import fused_attention
```

```
attention_output = fused_attention(Q, K, V, scale=1.0 / math.sqrt(h))
```



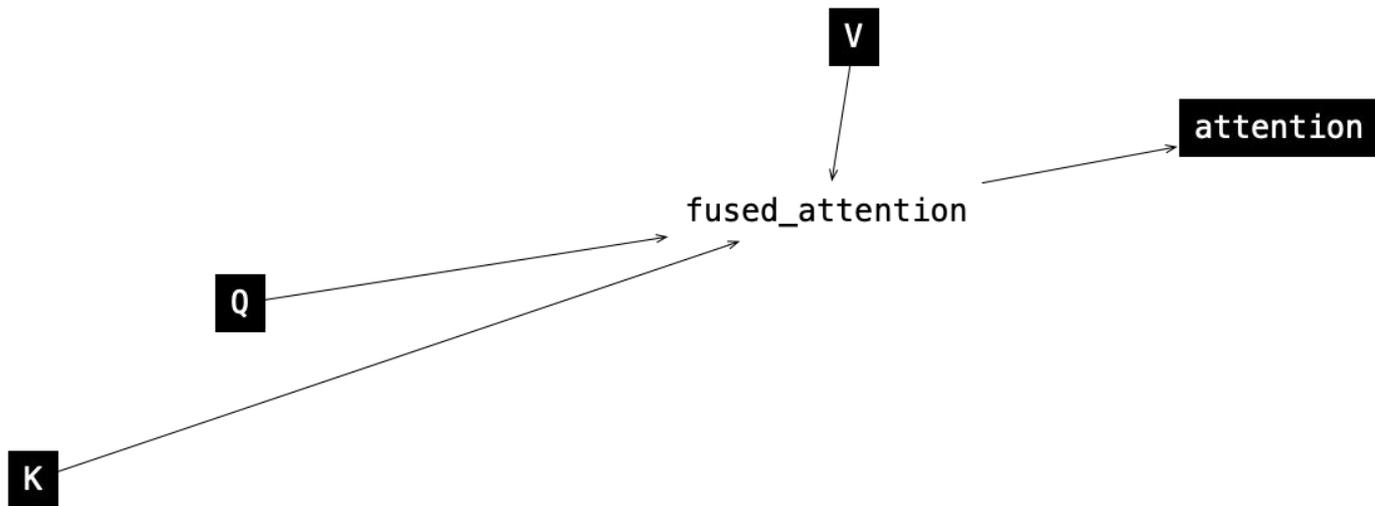
In memory:

- K
- Q

Example: Single-Head Attention

```
from flash_attention import fused_attention
```

```
attention_output = fused_attention(Q, K, V, scale=1.0 / math.sqrt(h))
```



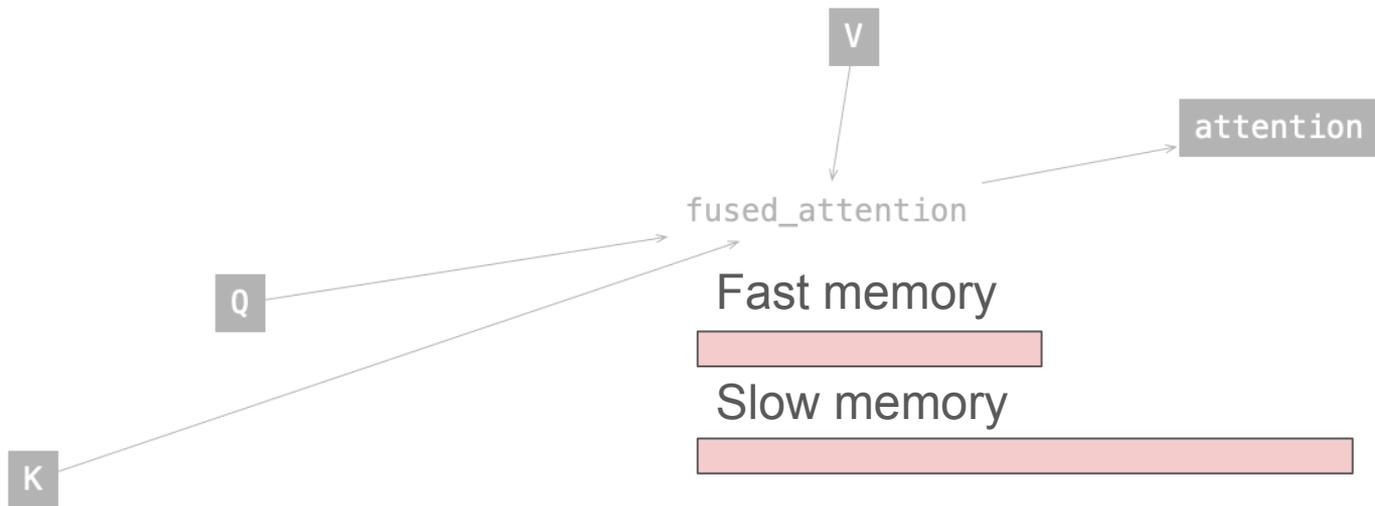
In memory:

- K
- Q
- V

Example: Single-Head Attention

```
from flash_attention import fused_attention
```

```
attention_output = fused_attention(Q, K, V, scale=1.0 / math.sqrt(h))
```



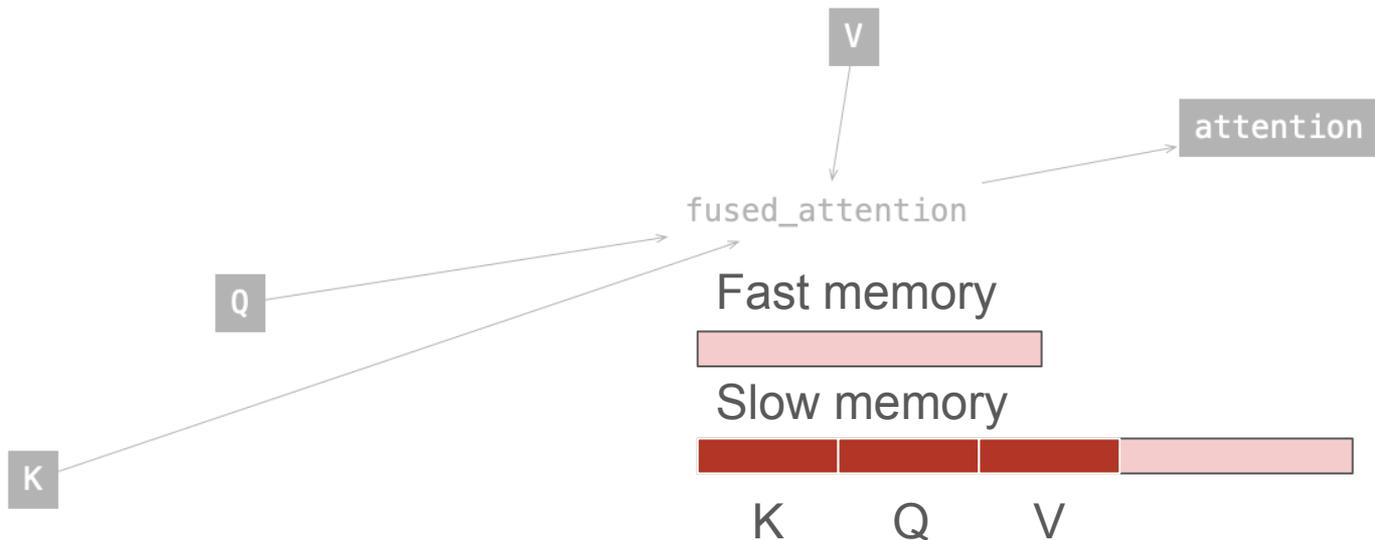
In memory:

- K
- Q
- V

Example: Single-Head Attention

```
from flash_attention import fused_attention
```

```
attention_output = fused_attention(Q, K, V, scale=1.0 / math.sqrt(h))
```



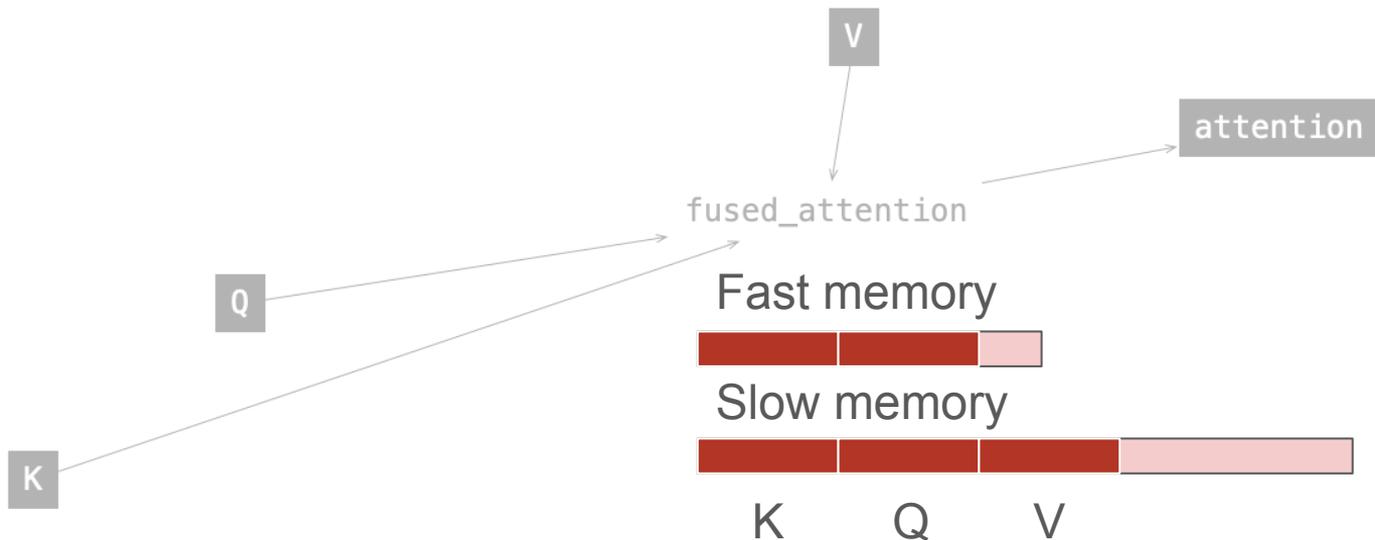
In memory:

- K
- Q
- V

Example: Single-Head Attention

```
from flash_attention import fused_attention
```

```
attention_output = fused_attention(Q, K, V, scale=1.0 / math.sqrt(h))
```



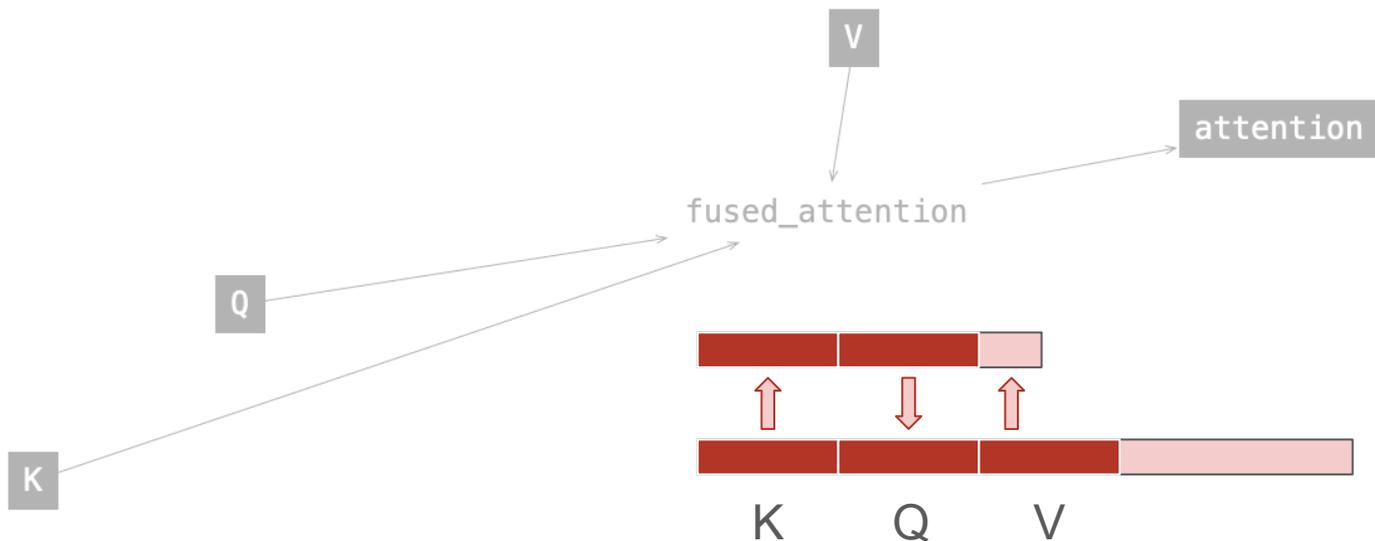
In memory:

- K
- Q

Example: Single-Head Attention

```
from flash_attention import fused_attention
```

```
attention_output = fused_attention(Q, K, V, scale=1.0 / math.sqrt(h))
```



In memory:

- K
- Q
- V

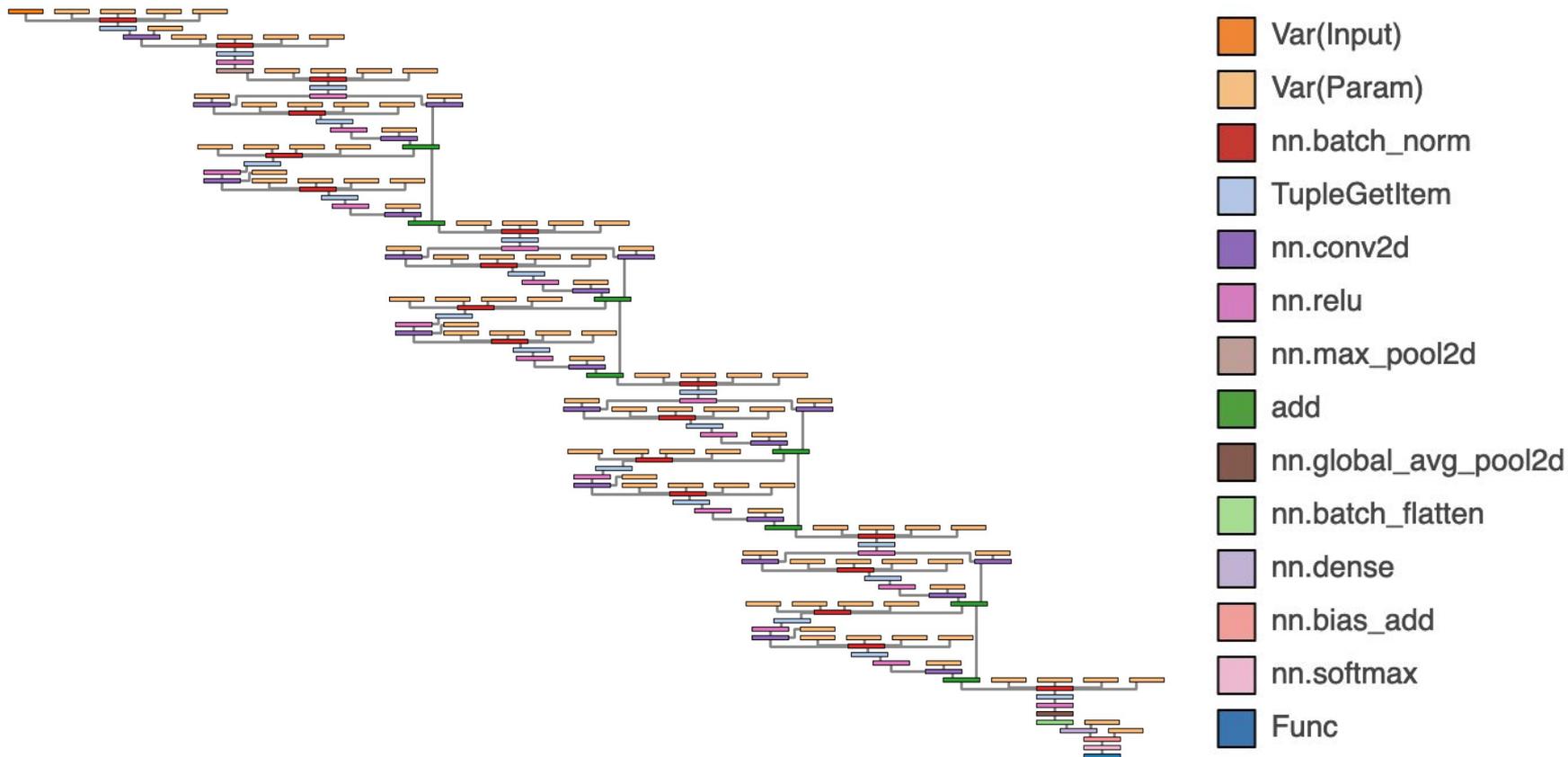
First Observation

1. Fusion improves performance, but not always

Second Observation

1. Fusion improves performance, but not always
2. Deciding which operators to fuse is a challenging problem

Example of a Computational Graph: ResNet18



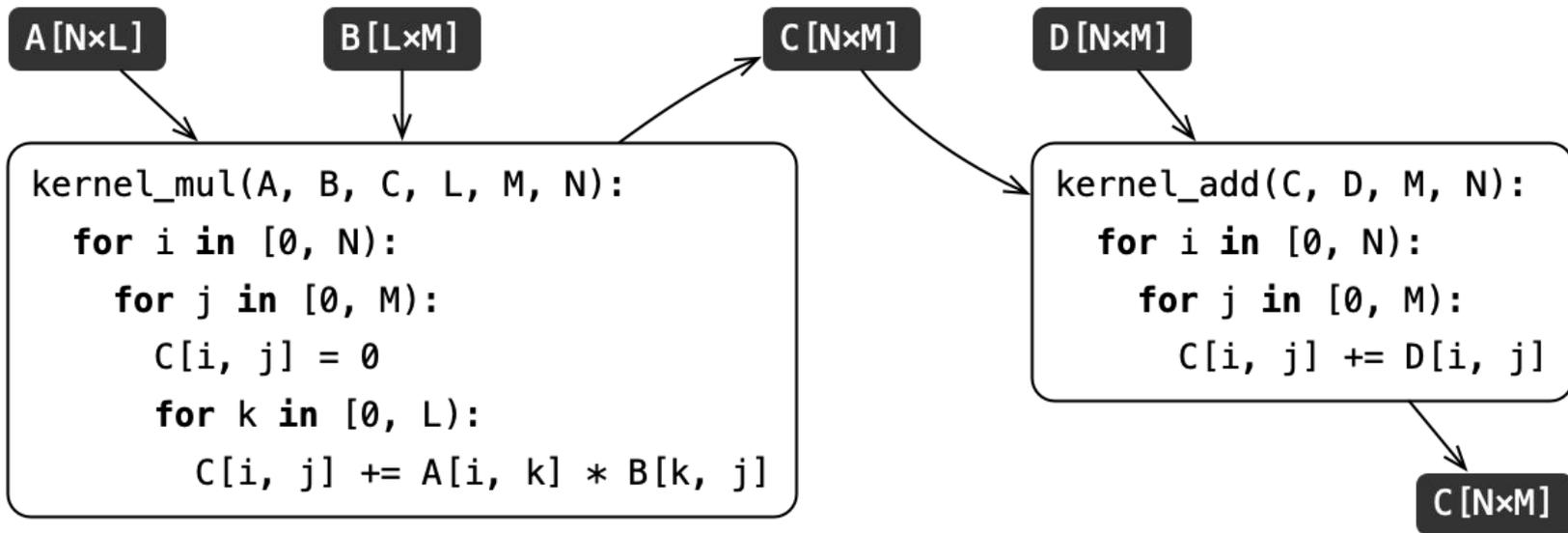
Which Operators does it Make Sense to Fuse?

Which Operators does it Make Sense to Fuse?

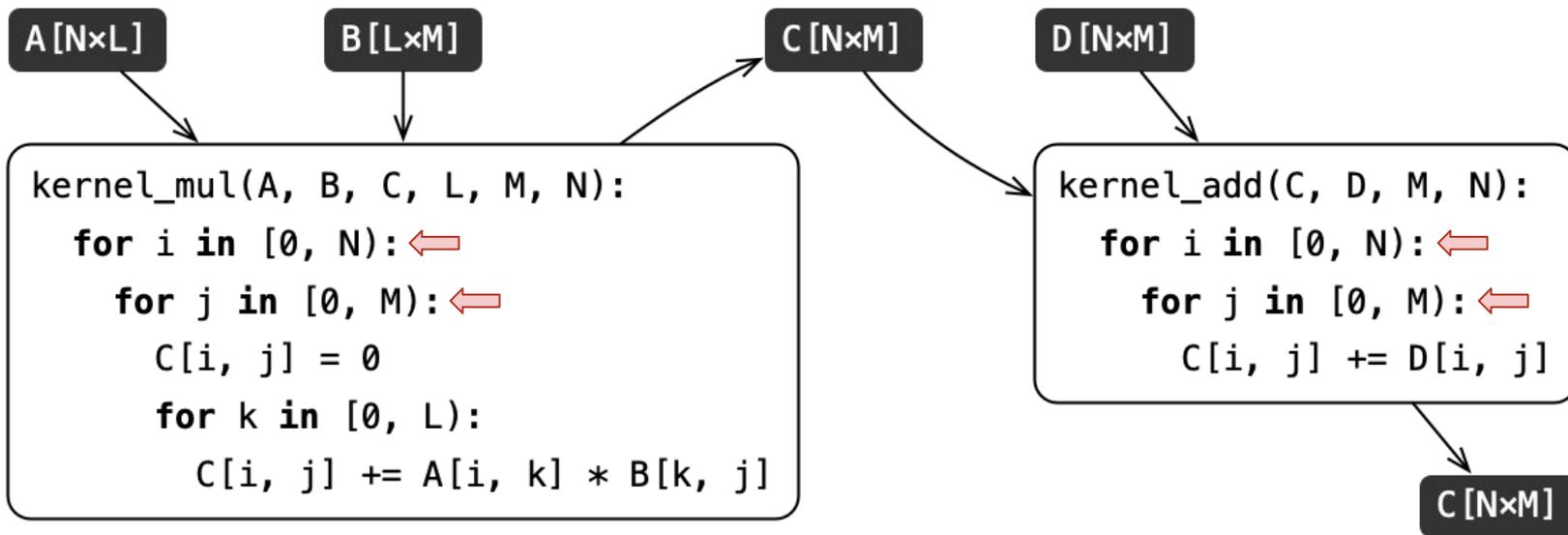
```
kernel_mul(A, B, C, L, M, N):  
    for i in [0, N):  
        for j in [0, M):  
            C[i, j] = 0  
            for k in [0, L):  
                C[i, j] += A[i, k] * B[k, j]
```

```
kernel_add(C, D, M, N):  
    for i in [0, N):  
        for j in [0, M):  
            C[i, j] += D[i, j]
```

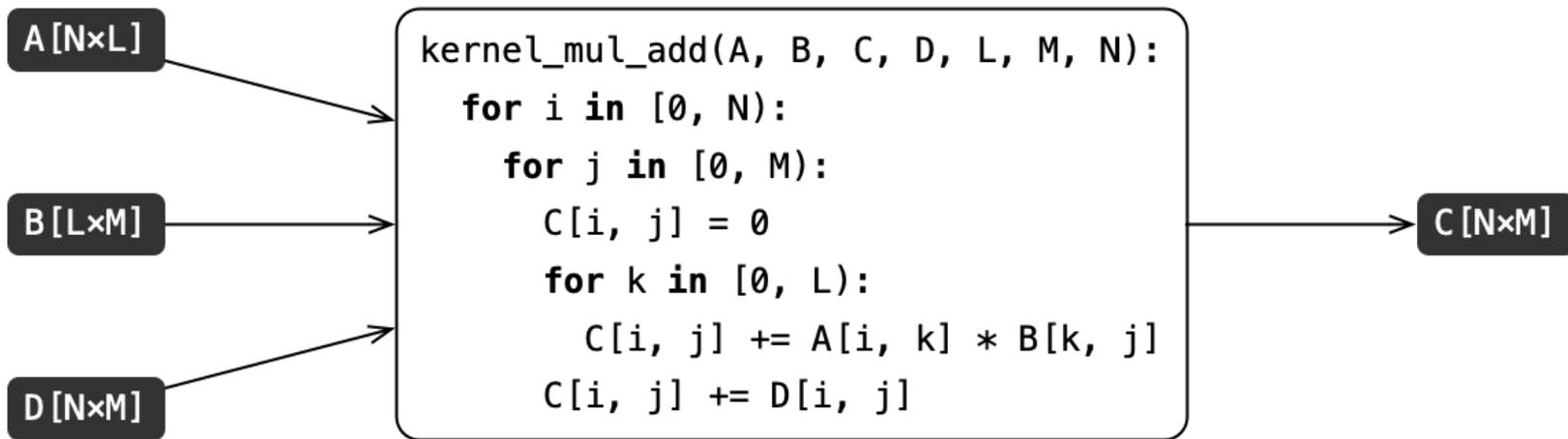
Which Operators does it Make Sense to Fuse?



Which Operators does it Make Sense to Fuse?



Which Operators does it Make Sense to Fuse?



Third Observation

1. Fusion improves performance, but not always
2. Deciding which operators to fuse is a challenging problem
3. Fusion, when seen as an optimization problem, is guided by a cost model

Third Observation

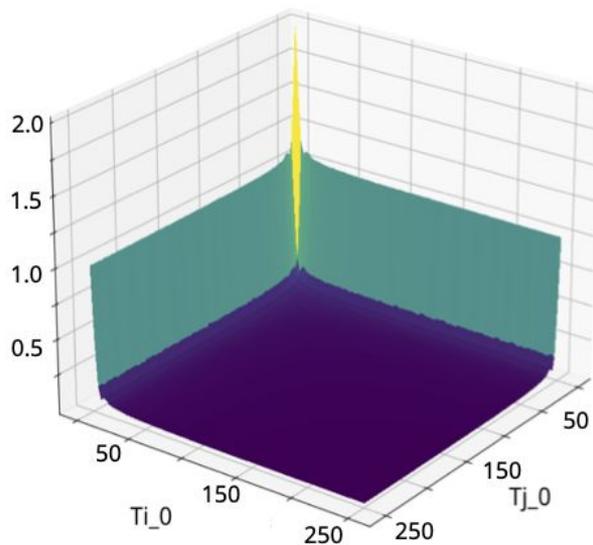
1. Fusion improves performance, but not always
2. Deciding which operators to fuse is a challenging problem
3. Fusion, when seen as an optimization problem, is guided by a cost model
 - a. **Dynamic:** profile and decide

Third Observation

1. Fusion improves performance, but not always
2. Deciding which operators to fuse is a challenging problem
3. Fusion, when seen as an optimization problem, is guided by a cost model
 - a. **Dynamic:** profile and decide
 - b. **Static:** analytical function

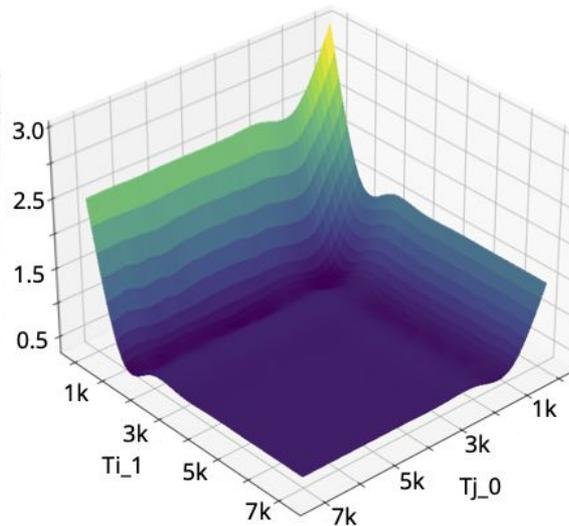
■ Static Cost Model: Program → Estimated Running Time

Static Cost Model: Program \rightarrow Estimated Running Time



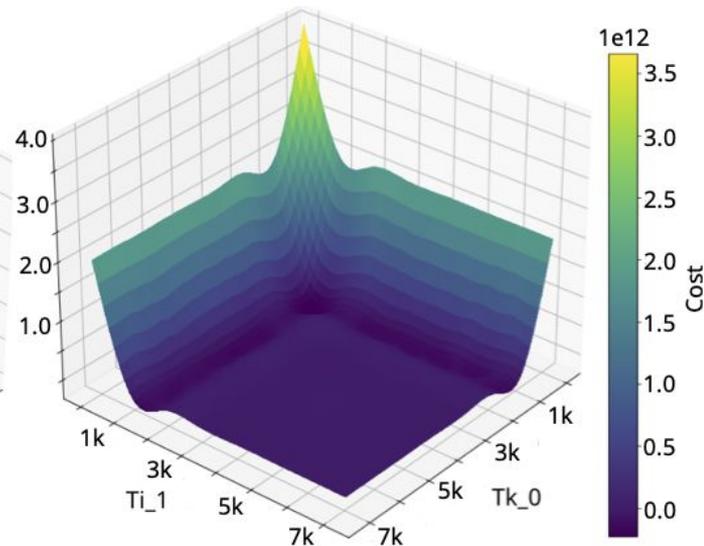
(a) Single cache L1.

$$\text{Cost} = N_i * N_j * N_k * (1/T_{j_0} + 1/T_{i_0} + 2/N_k)$$



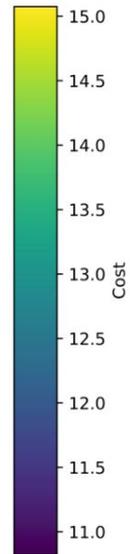
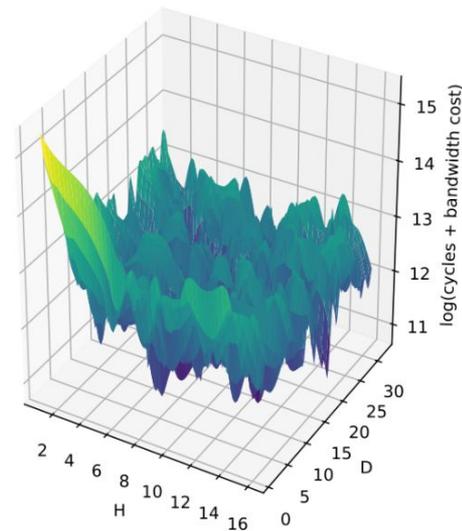
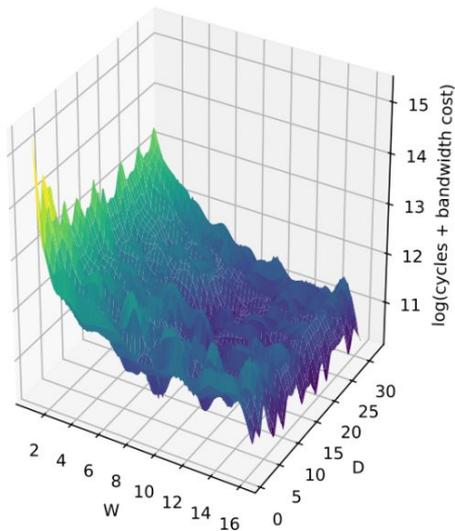
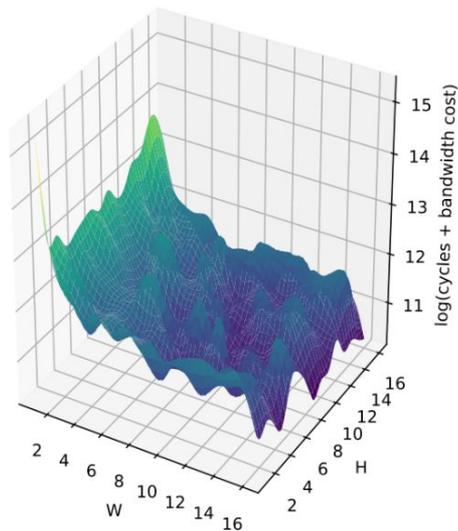
(b) L1 and L2, varying T_{i_1} and T_{j_0} .

$$\text{Cost} = N_i * N_j * N_k * (1/T_{j_1} + 1/T_{i_1} + 2/N_k) + N_i * N_j * N_k * (2/T_{k_0} + 1/T_{j_0} + 1/T_{i_1})$$



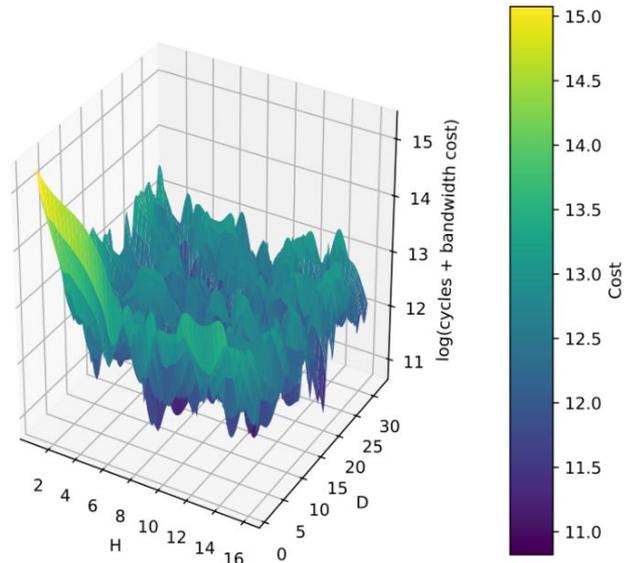
(c) L1 and L2, varying T_{i_1} and T_{k_0} .

Static Cost Model: Program \rightarrow Estimated Running Time

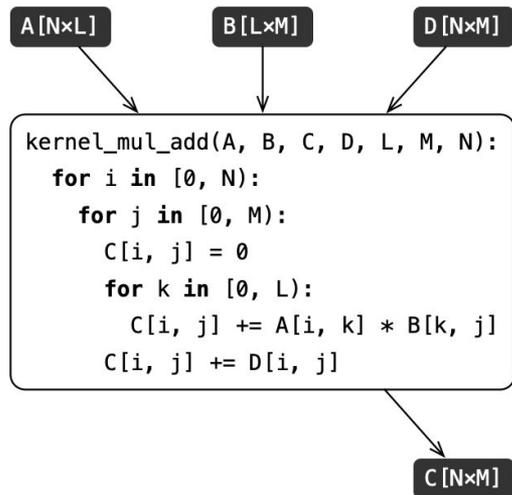


Static Cost Model: Program \rightarrow Estimated Running Time

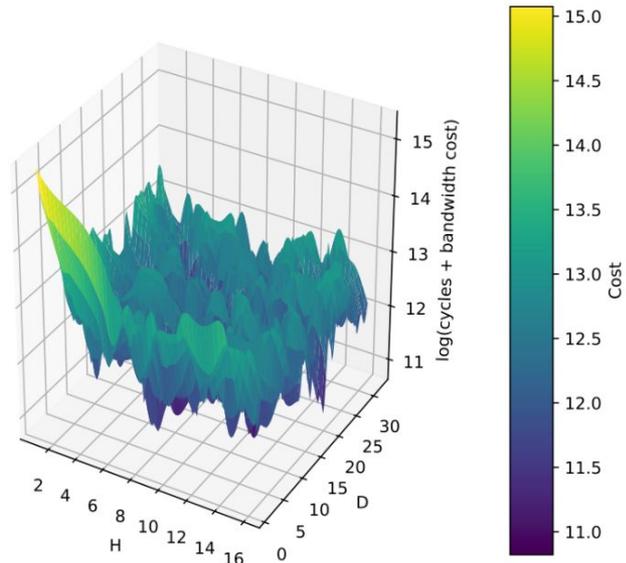
```
def perfModel(L, M, N, SIMD_WIDTH):
    oW, oH, oD = L, M, N
    vw = SIMD_WIDTH - 1
    return TreePerfModel(
        name,
        TPLoops(
            (oW.dim + vw - 1) // vw,
            TPSeq(
                TPLoops([oH, oD, oW]),
                TPCond(
                    cneq(oH % 2, 0),
                    TPLoops([(oD + 1, oH)])
                ),
            ),
        ),
    )
```



Static Cost Model: Program \rightarrow Estimated Running Time

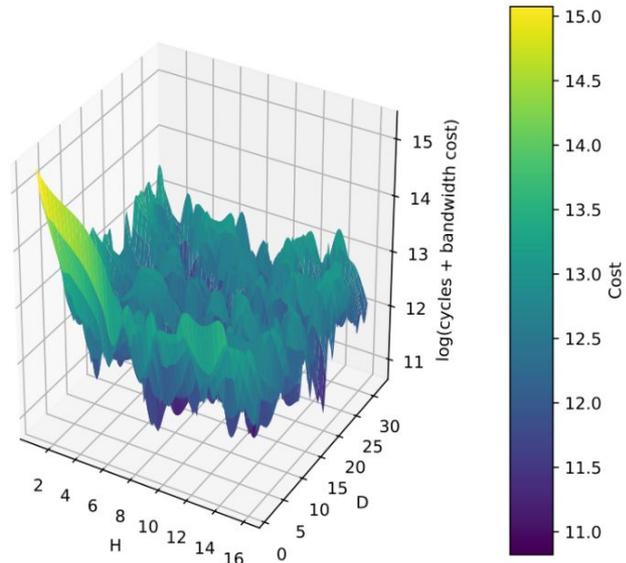
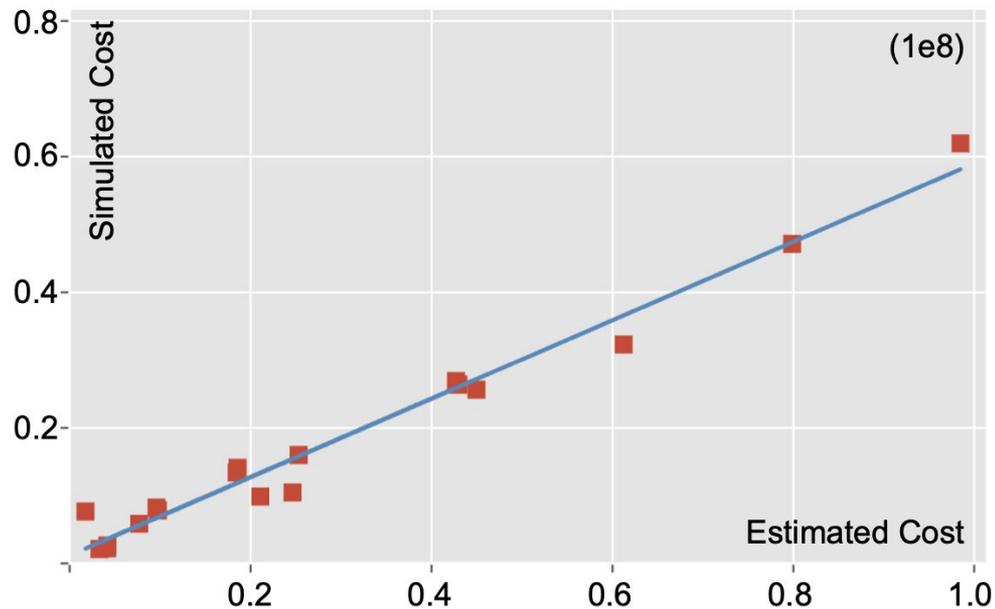


```
def perfModel(L, M, N, SIMD_WIDTH):  
    oW, oH, oD = L, M, N  
    vw = SIMD_WIDTH - 1  
    return TreePerfModel(  
        name,  
        TPLoops(  
            (oW.dim + vw - 1) // vw,  
            TPSeq(  
                TPLoops([oH, oD, oW]),  
                TPCond(  
                    cneq(oH % 2, 0),  
                    TPLoops([(oD + 1, oH])  
                ),  
            ),  
        ),  
    )
```

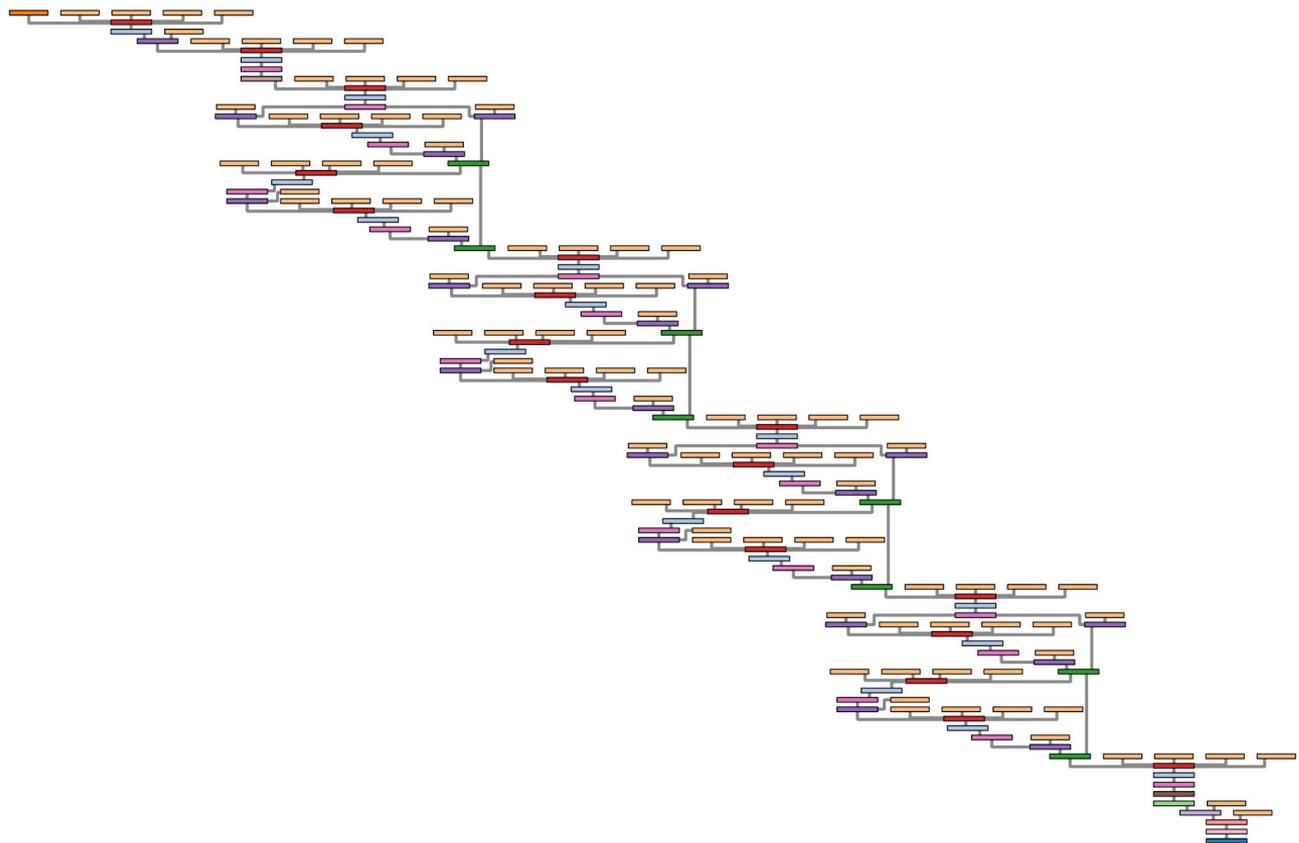


Static Cost Model: Program \rightarrow Estimated Running Time

DCC

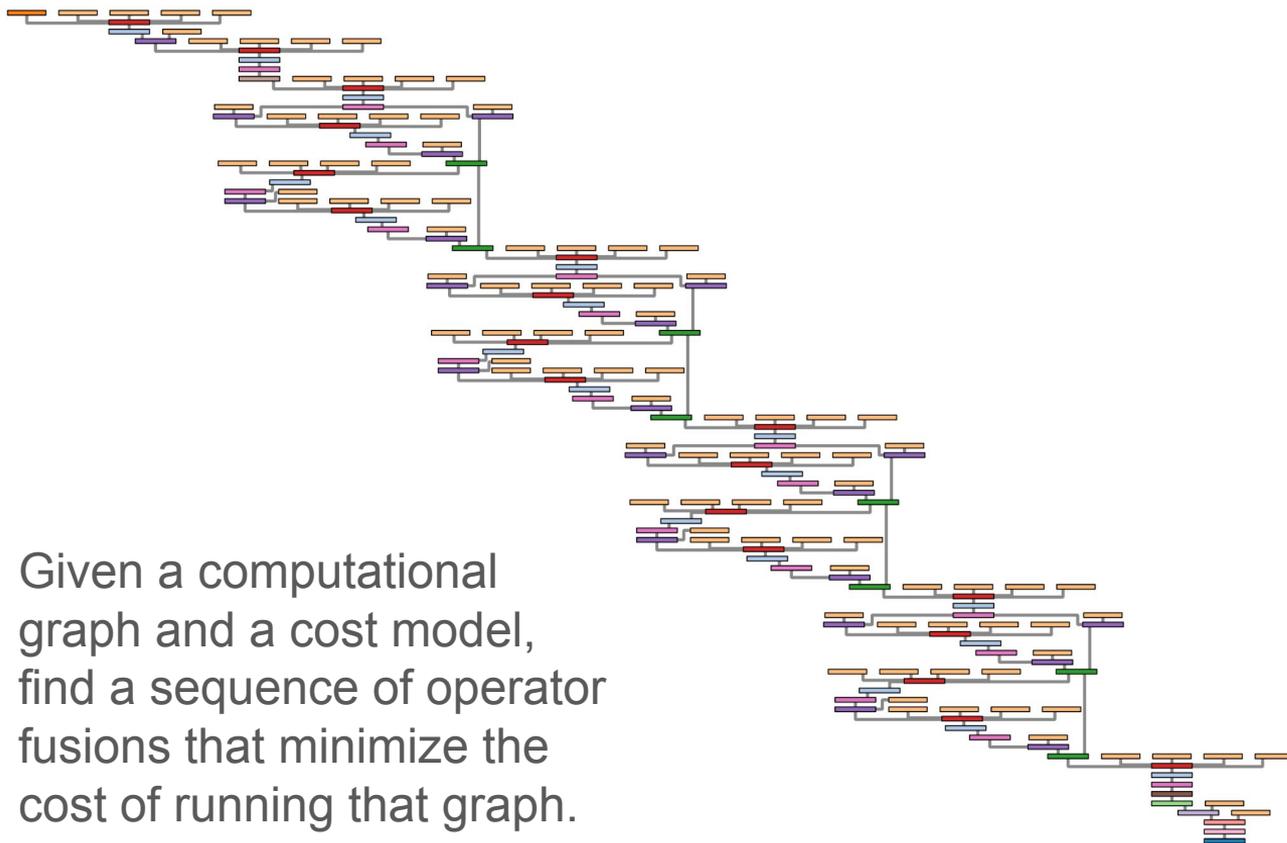


Fusion as an Optimization Problem



- Var(Input)
- Var(Param)
- nn.batch_norm
- TupleGetItem
- nn.conv2d
- nn.relu
- nn.max_pool2d
- add
- nn.global_avg_pool2d
- nn.batch_flatten
- nn.dense
- nn.bias_add
- nn.softmax
- Func

Fusion as an Optimization Problem



- Var(Input)
- Var(Param)
- nn.batch_norm
- TupleGetItem
- nn.conv2d
- nn.relu
- nn.max_pool2d
- add
- nn.global_avg_pool2d
- nn.batch_flatten
- nn.dense
- nn.bias_add
- nn.softmax
- Func

Fusion as an Optimization Problem

Exhaustive:

Xuyi Cai, Ying Wang, Lei Zhang: **Optimus: An Operator Fusion Framework for Deep Neural Networks**. ACM Trans. Embed. Comput. Syst. 22(1): 1:1-1:26 (2023)

Fusion as an Optimization Problem

Exhaustive:

Xuyi Cai, Ying Wang, Lei Zhang: **Optimus: An Operator Fusion Framework for Deep Neural Networks**. ACM Trans. Embed. Comput. Syst. 22(1): 1:1-1:26 (2023)

Evolutionary:

Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, Joseph E. Gonzalez, Ion Stoica: **Ansor: Generating High-Performance Tensor Programs for Deep Learning**. OSDI 2020: 863-879

Fusion as an Optimization Problem

Exhaustive:

Xuyi Cai, Ying Wang, Lei Zhang: **Optimus: An Operator Fusion Framework for Deep Neural Networks**. ACM Trans. Embed. Comput. Syst. 22(1): 1:1-1:26 (2023)

Evolutionary:

Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, Joseph E. Gonzalez, Ion Stoica: **Ansor: Generating High-Performance Tensor Programs for Deep Learning**. OSDI 2020: 863-879

Minimum Spanning Tree:

Michael Canesche, Vanderson Rosário, Edson Borin, Fernando Magno Quintão Pereira: **Fusion of Operators of Computational Graphs via Greedy Clustering: The XNNC Experience**. CC 2025: 1-13

The Minimum Spanning Tree Approach

The Minimum Spanning Tree Approach

- Greedy algorithm guided by cost model

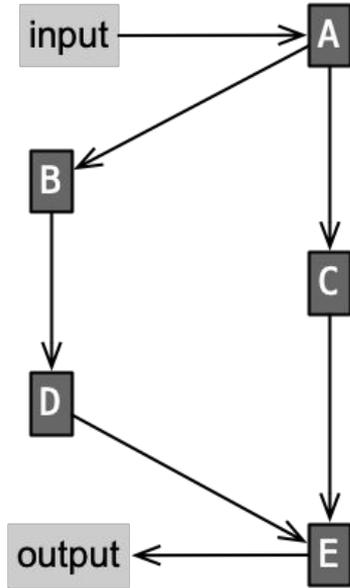
The Minimum Spanning Tree Approach

- Greedy algorithm guided by cost model
- Group nodes into clusters

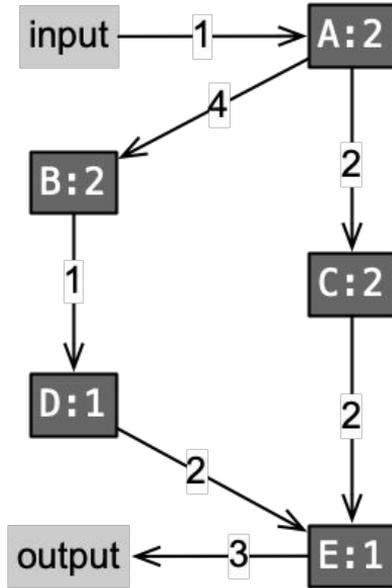
The Minimum Spanning Tree Approach

- Greedy algorithm guided by cost model
- Group nodes into clusters
- Fuse the two most profitable clusters (Until no more fusions are possible)

The Minimum Spanning Tree Approach



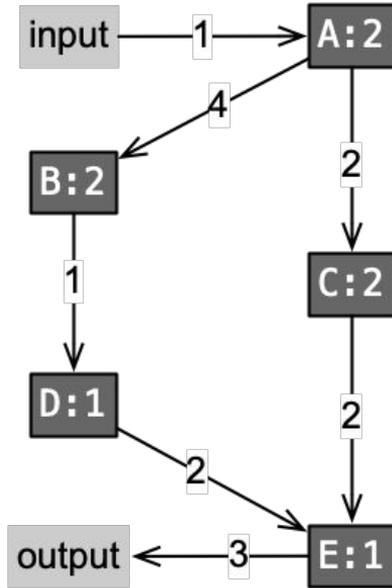
The Minimum Spanning Tree Approach



Cost model:

- Edges: cost of data transfer
- Vertices: cost of computation

The Minimum Spanning Tree Approach



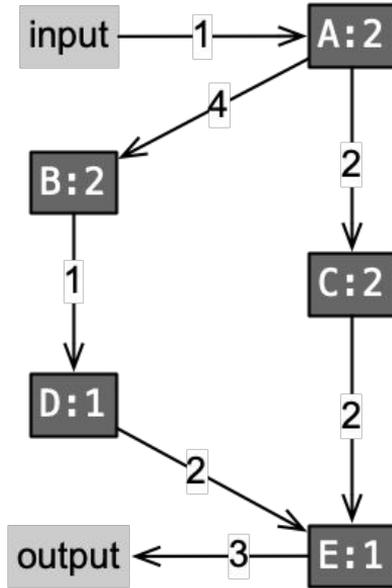
Cost model:

- Edges: cost of data transfer
- Vertices: cost of computation

(Non-) Properties:

- Non-associative
 - $\text{cost}(\text{fuse}(A, \text{fuse}(B, C))) \neq \text{cost}(\text{fuse}(\text{fuse}(A, B), C))$

The Minimum Spanning Tree Approach



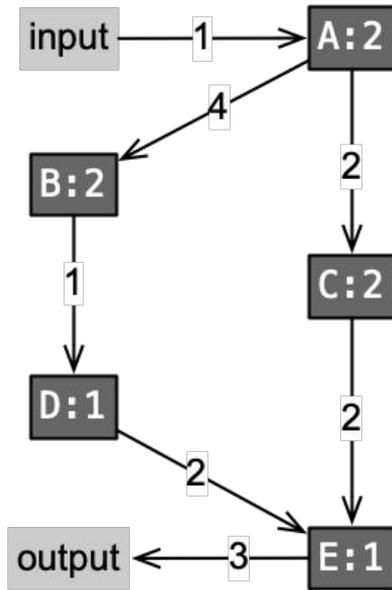
Cost model:

- Edges: cost of data transfer
- Vertices: cost of computation

(Non-) Properties:

- Non-associative
 - $\text{cost}(\text{fuse}(A, \text{fuse}(B, C))) \neq \text{cost}(\text{fuse}(\text{fuse}(A, B), C))$
- No Matroid Structure
 - Greedy choices are not guaranteed to be optimal

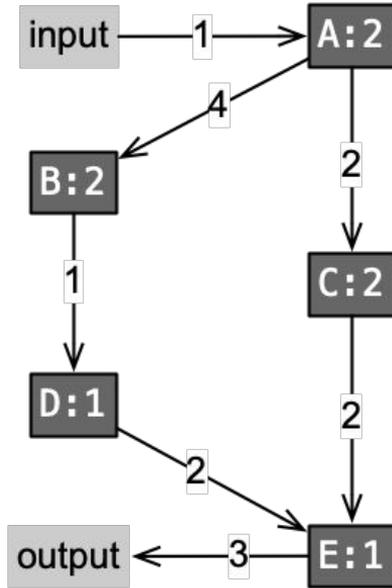
The Minimum Spanning Tree Approach



Assumption (example only!)

- $\text{Cost}(\text{edges}) < 5$

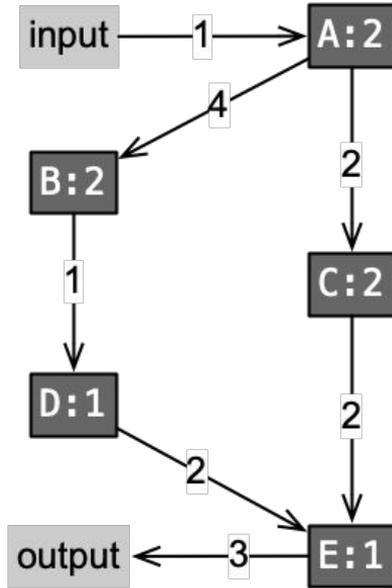
The Minimum Spanning Tree Approach



Assumption (example only!)

- Cost(edges) < 5
- Eg.: memory capacity == 4

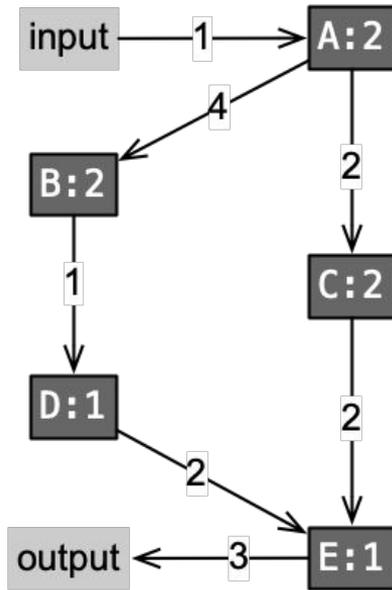
The Minimum Spanning Tree Approach



Assumption (example only!)

- $\text{Cost}(\text{edges}) < 5$
- Eg.: memory capacity == 4
- $\text{Cost}(\text{fuse}(V_A, V_B)) == \text{Cost}(V_A) + \text{Cost}(V_B)$

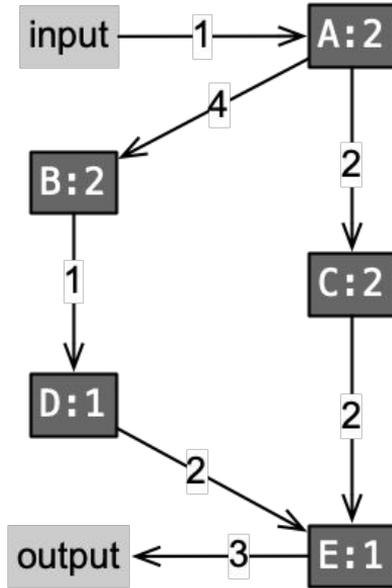
The Minimum Spanning Tree Approach



Assumption (example only!)

- $\text{Cost}(\text{edges}) < 5$
- Eg.: memory capacity == 4
- $\text{Cost}(\text{fuse}(V_A, V_B)) == \text{Cost}(V_A) + \text{Cost}(V_B)$
 - "cost is homomorphic with regards to fusion"

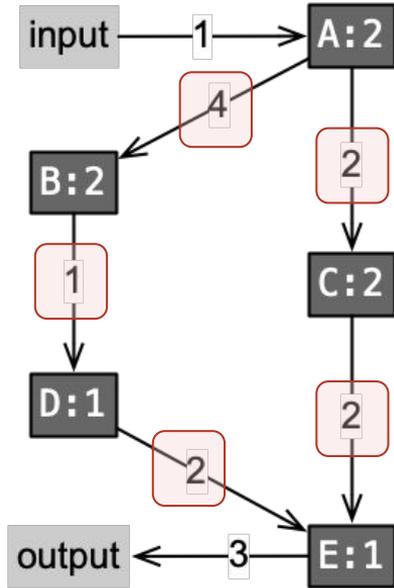
The Minimum Spanning Tree Approach



Assumption (example only!)

- $\text{Cost}(\text{edges}) < 5$
- Eg.: memory capacity == 4
- $\text{Cost}(\text{fuse}(V_A, V_B)) == \text{Cost}(V_A) + \text{Cost}(V_B)$
- $\text{Cost}(\text{fuse}(V_A \rightarrow V_B)) == 0$

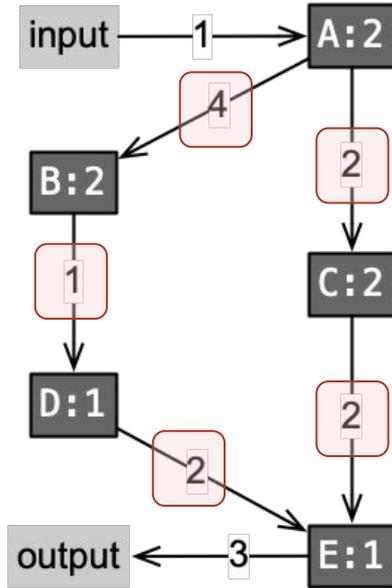
The Minimum Spanning Tree Approach



Edges: 15, Vertices: 8

- Five fusion candidates

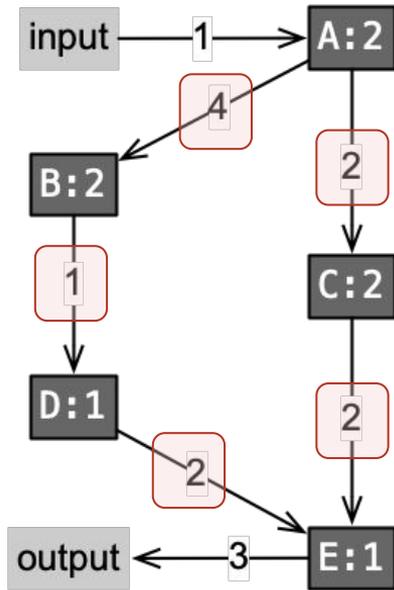
The Minimum Spanning Tree Approach



Edges: 15, Vertices: 8

- Five fusion candidates
- The most profitable?

The Minimum Spanning Tree Approach

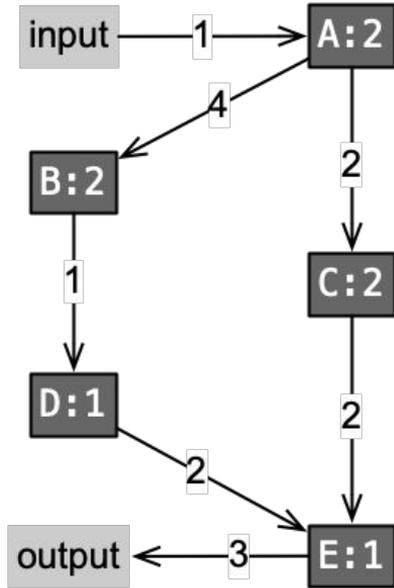


Reduces total cost the most,
e.g.: $C(\text{edges}) + C(\text{vertices})$

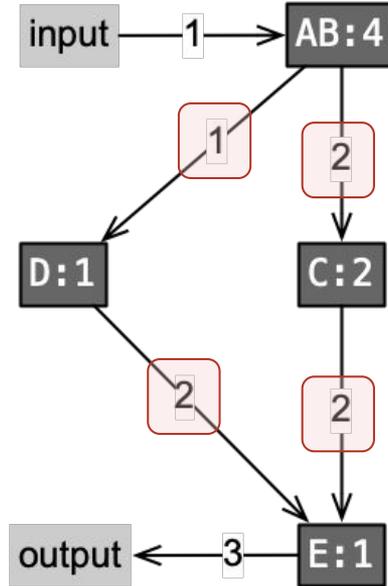
Edges: 15, Vertices: 8

- Five fusion candidates
- The most profitable?

The Minimum Spanning Tree Approach



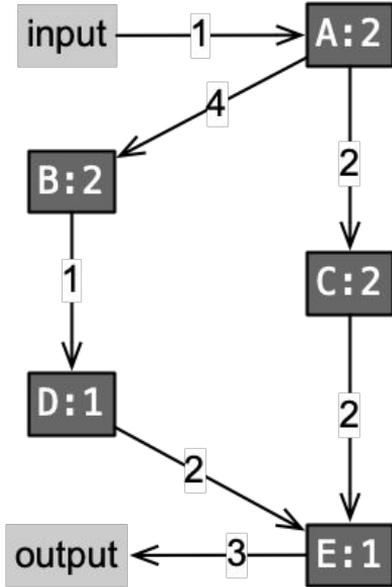
Edges: 15, Vertices: 8



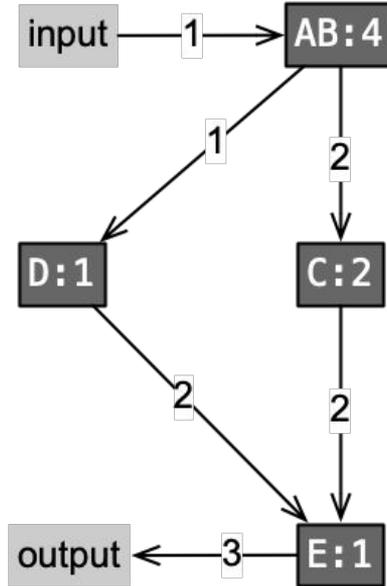
Edges: 11, Vertices: 8

- Four fusion candidates
- The most profitable?

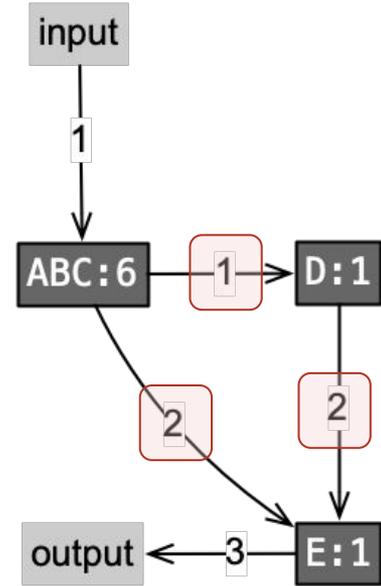
The Minimum Spanning Tree Approach



Edges: 15, Vertices: 8

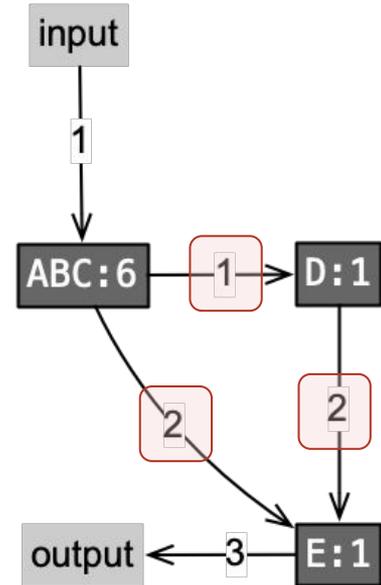


Edges: 11, Vertices: 8



Edges: 9, Vertices: 8

The Minimum Spanning Tree Approach



Edges: 9, Vertices: 8

When no more fusions are possible, we are done

How Good is the Cluster-Based Algorithm?

How Good is the Cluster-Based Algorithm?

- Benchmarks: 19 convolutional neural networks

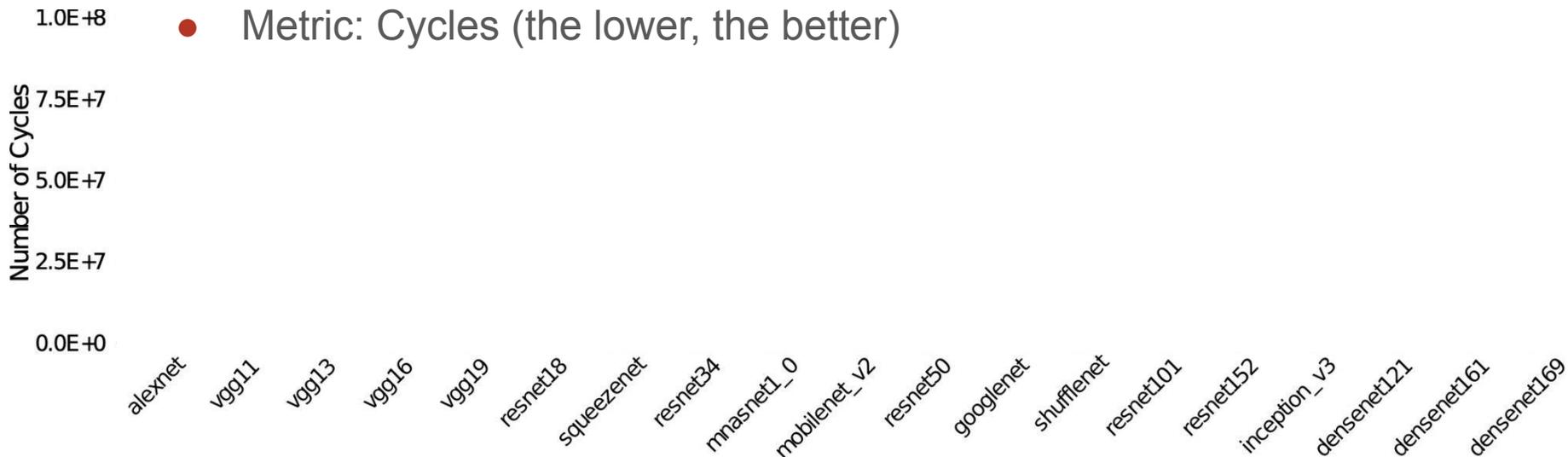
alexnet vgg11 vgg13 vgg16 vgg19 resnet18 squeezeenet resnet34 mnasnet1_0 mobilenet_v2 resnet50 googlenet shufflenet resnet101 resnet152 inception_v3 densenet121 densenet161 densenet169

How Good is the Cluster-Based Algorithm?

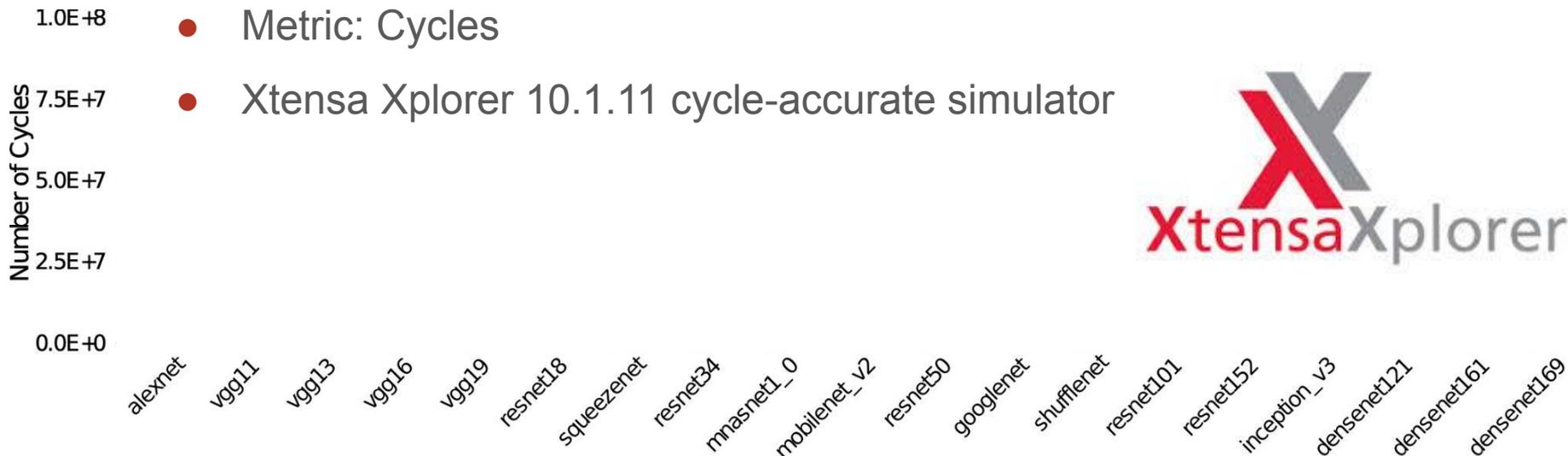
- Benchmarks: 19 convolutional neural networks
- Same image (batch=1, channels=3, height=224, width=224)

alexnet vgg11 vgg13 vgg16 vgg19 resnet18 squeezeenet resnet34 mmasnet1_0 mobilenet_v2 resnet50 googlenet shufflenet resnet101 resnet152 inception_v3 densenet121 densenet161 densenet169

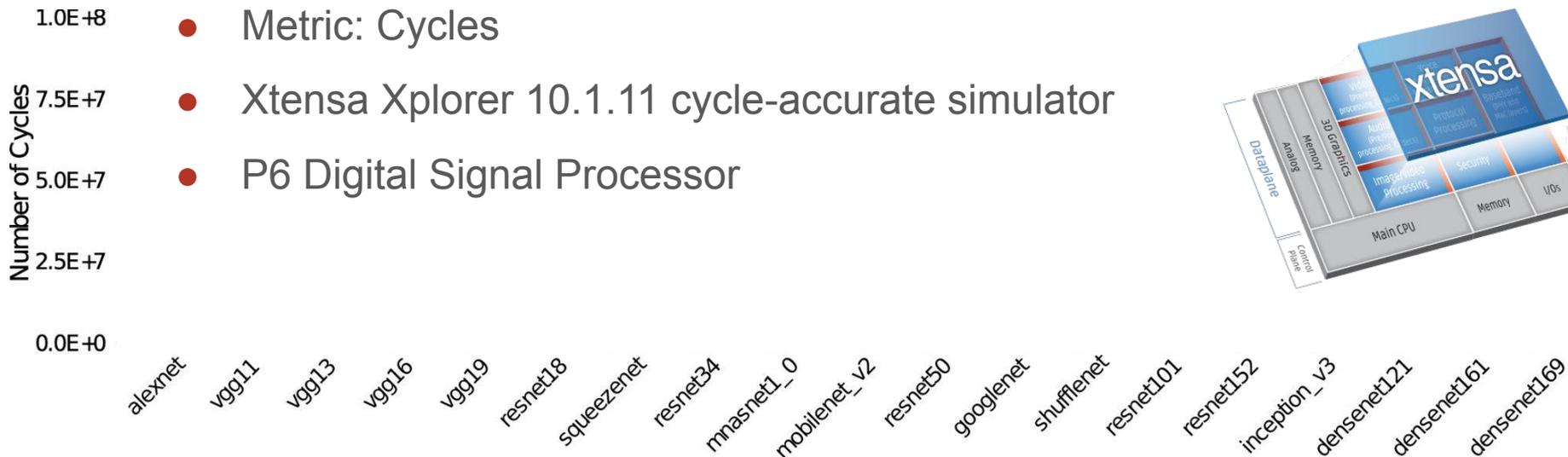
How Good is the Cluster-Based Algorithm?



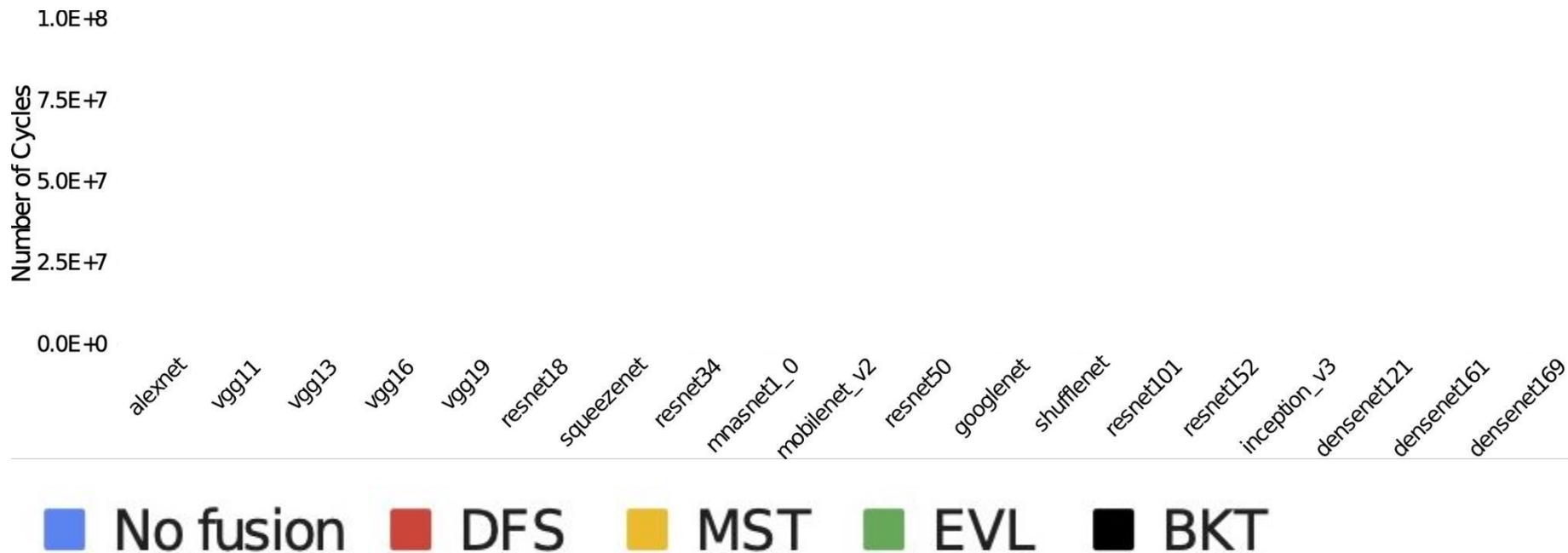
How Good is the Cluster-Based Algorithm?



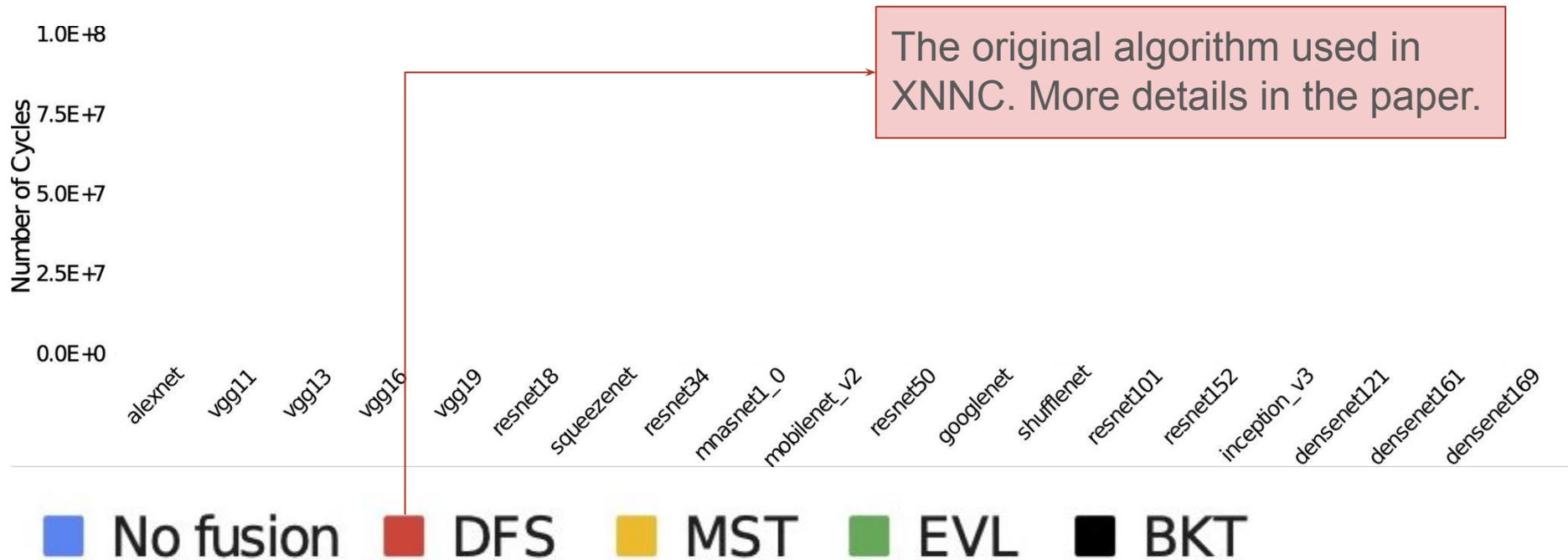
How Good is the Cluster-Based Algorithm?



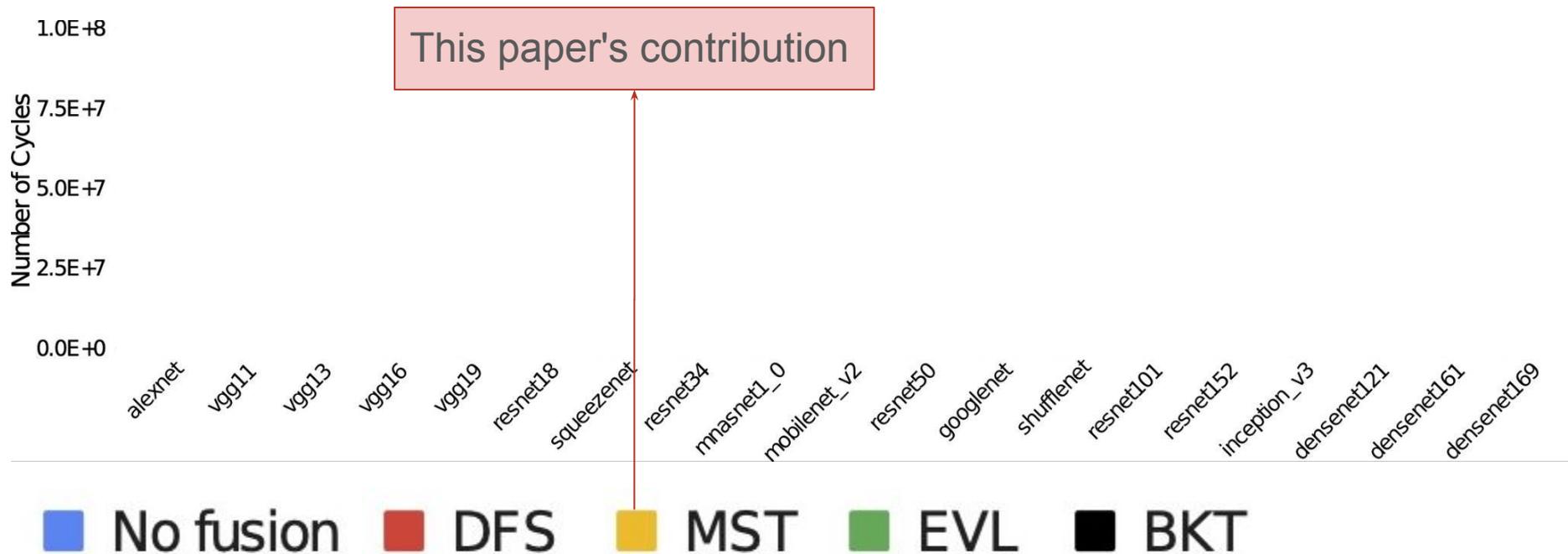
How Good is the Cluster-Based Algorithm?



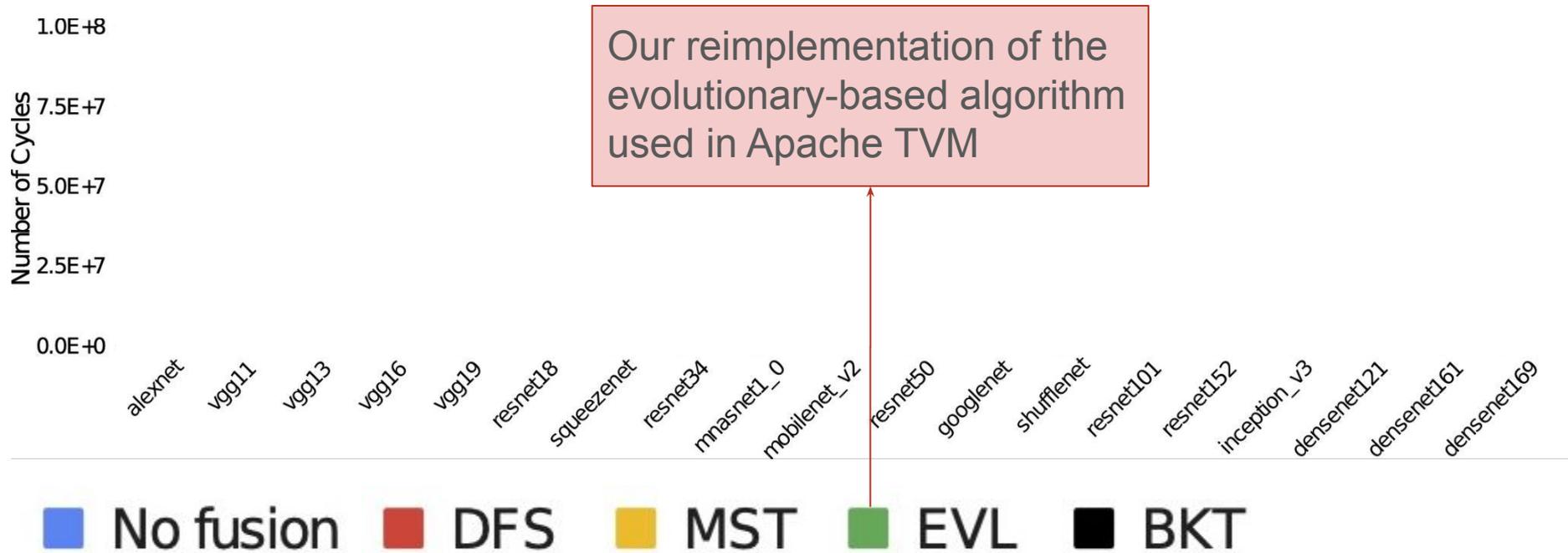
How Good is the Cluster-Based Algorithm?



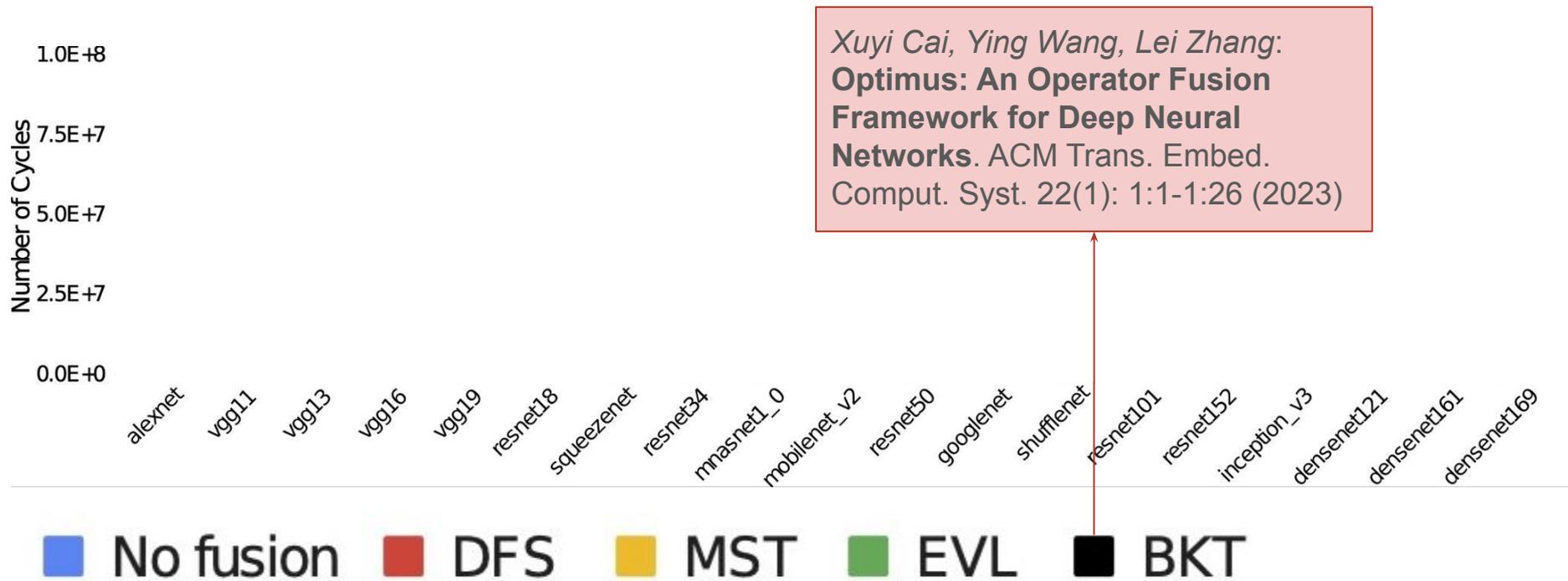
How Good is the Cluster-Based Algorithm?



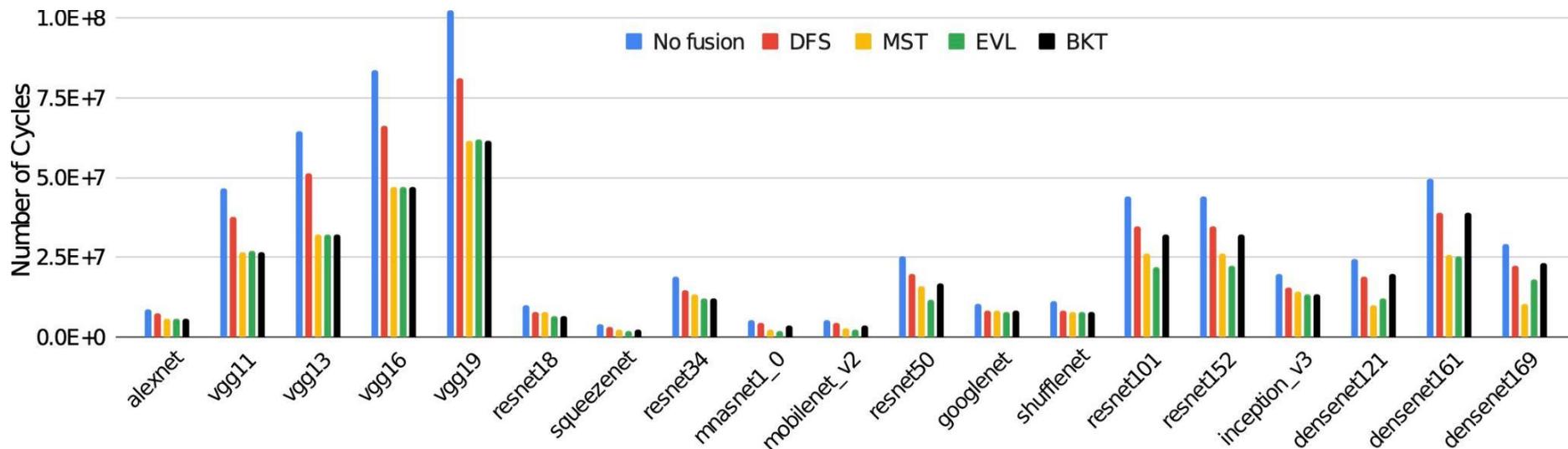
How Good is the Cluster-Based Algorithm?



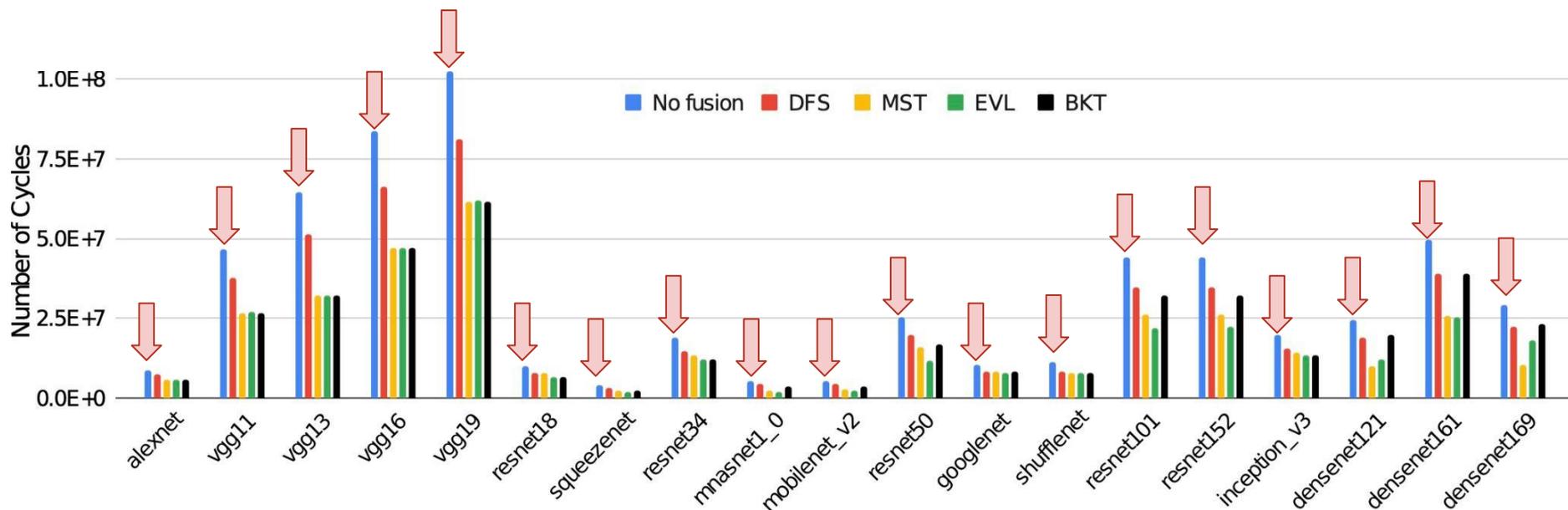
How Good is the Cluster-Based Algorithm?



How Good is the Cluster-Based Algorithm?

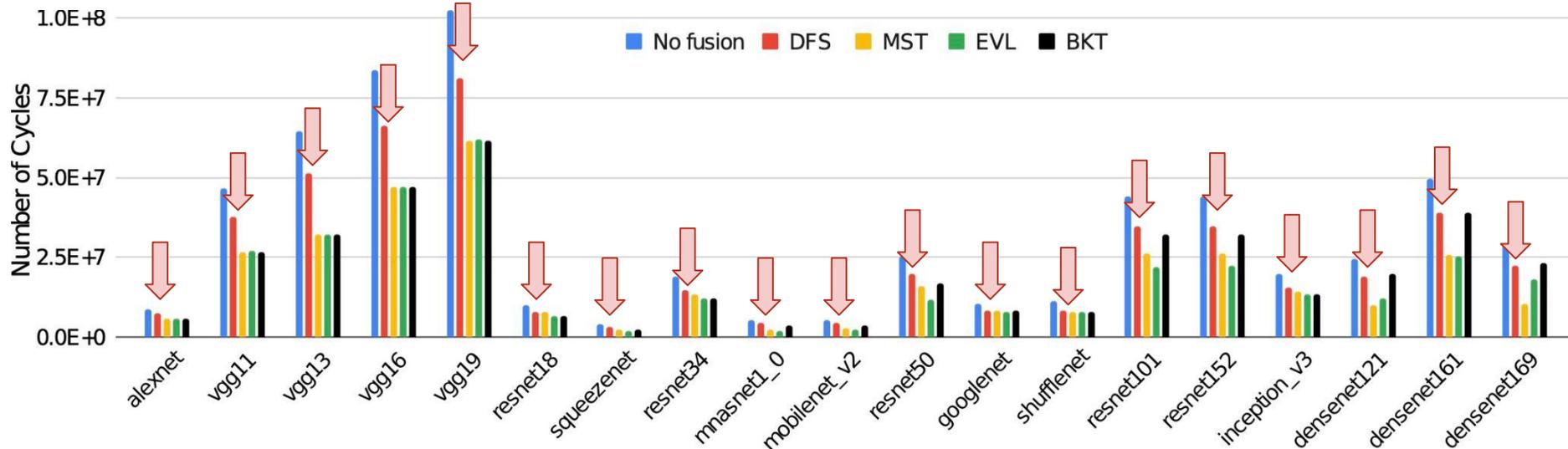


How Good is the Cluster-Based Algorithm?



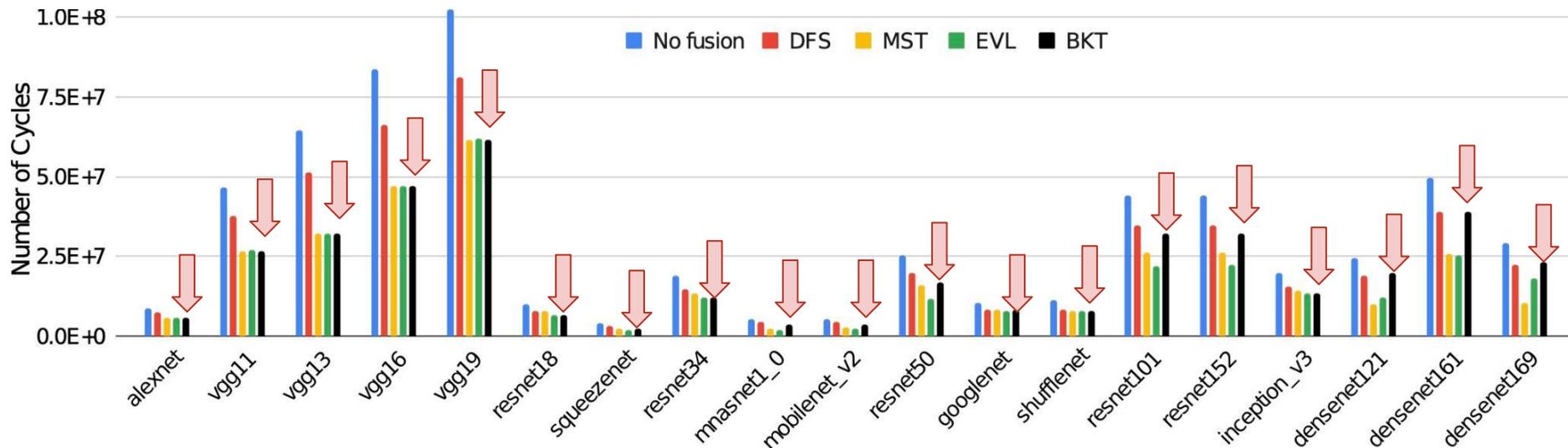
Operator fusion, regardless of the approach, is usually good

How Good is the Cluster-Based Algorithm?



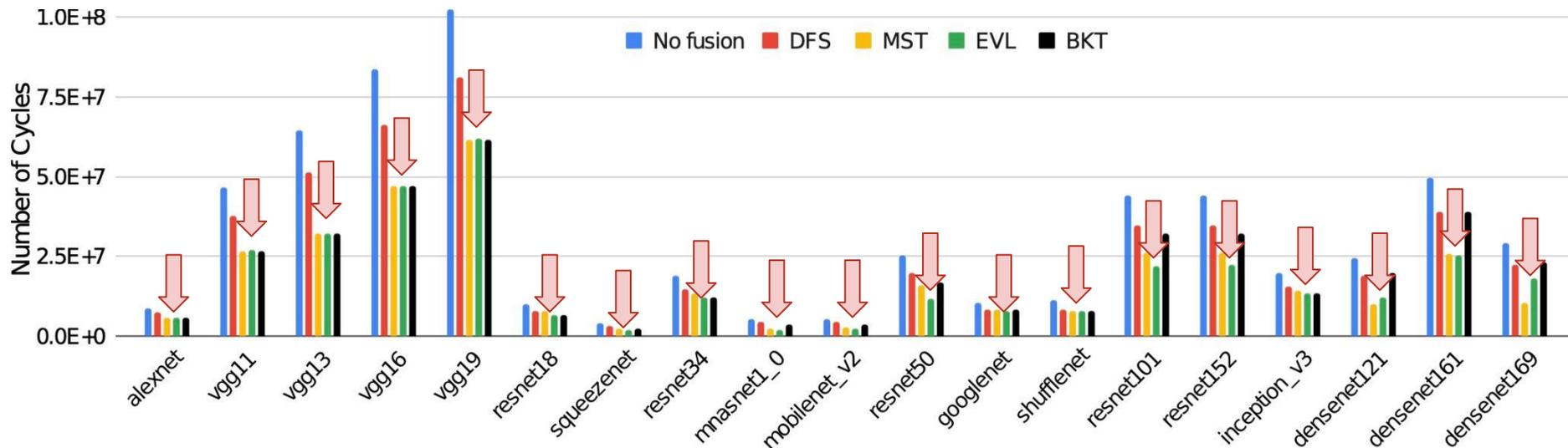
The new algorithm, MST, improves on XNNC's old algorithm (DFS) by 39%

How Good is the Cluster-Based Algorithm?



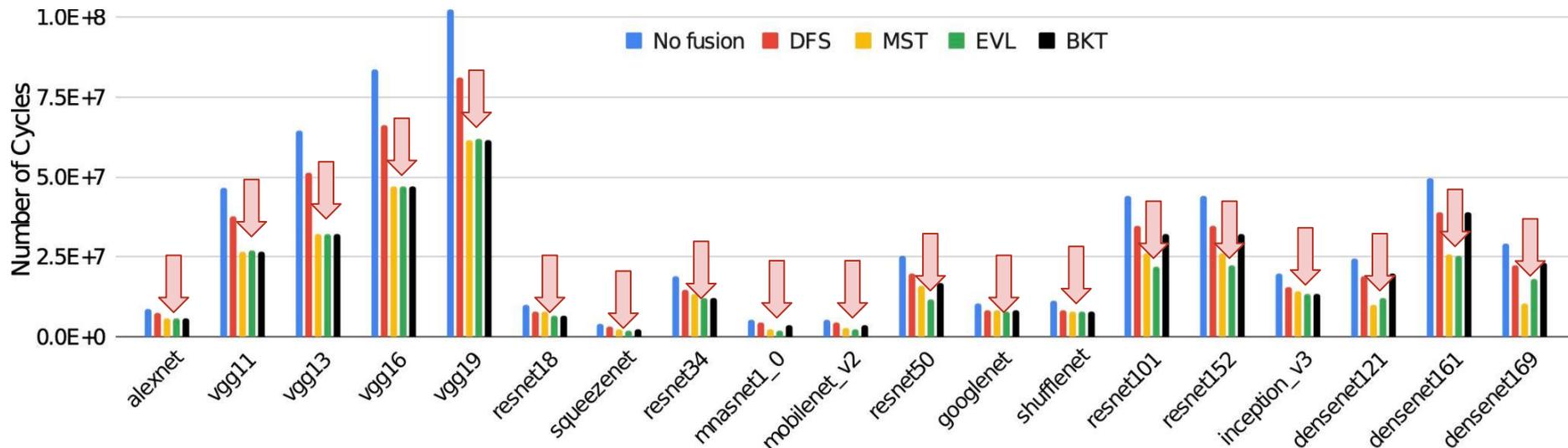
The exhaustive BKT did not yield the best performance (due to time-outs)

How Good is the Cluster-Based Algorithm?



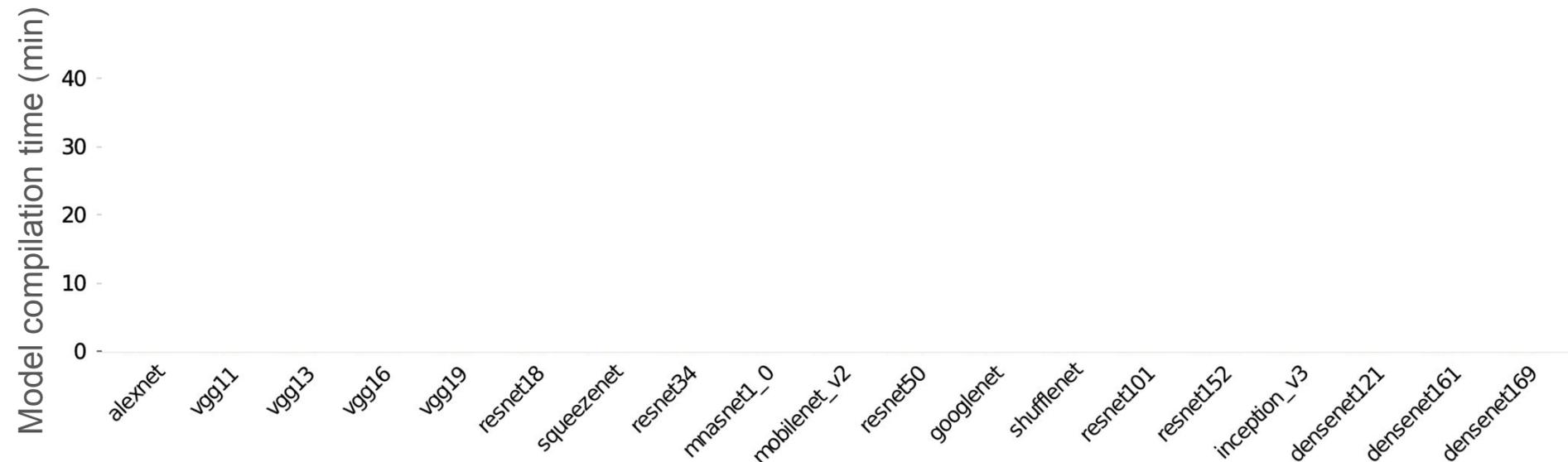
There was no statistical difference between MST and EVL

How Good is the Cluster-Based Algorithm?



There was no statistical difference between MST and EVL.
But MST is usually faster.

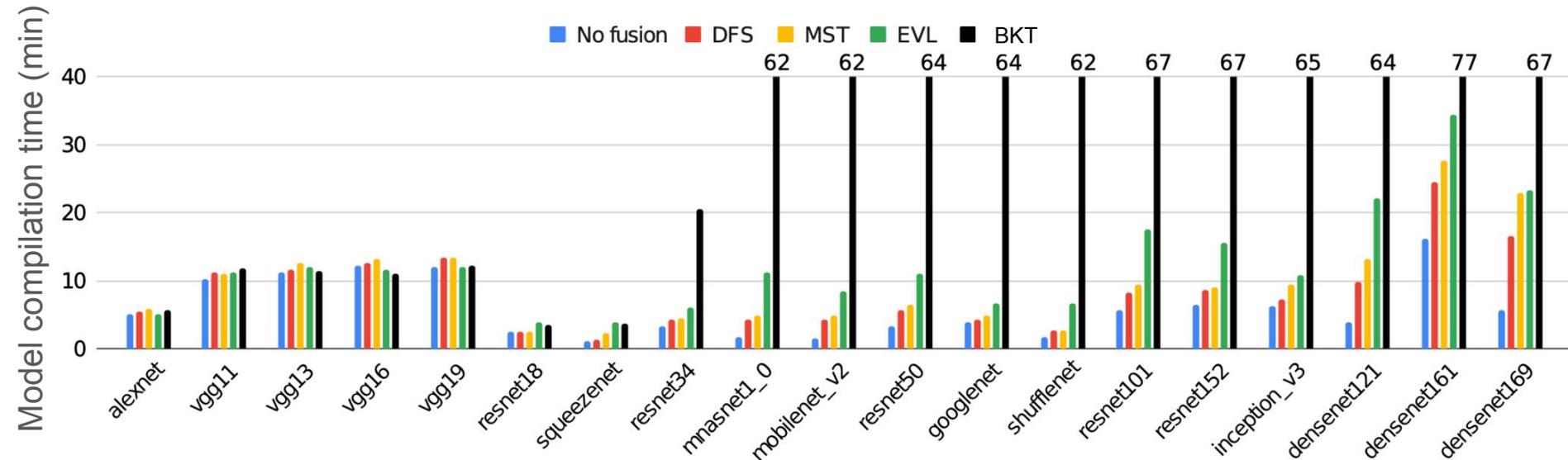
How Good is the Cluster-Based Algorithm?



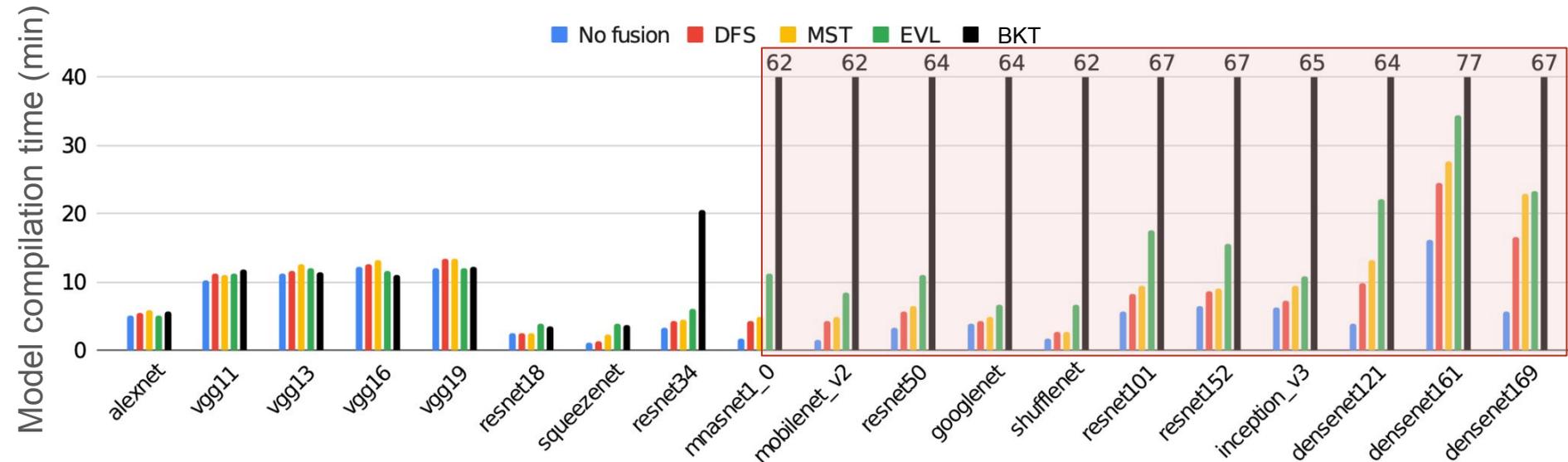
There was no statistical difference between MST and EVL.

But MST is usually faster.

How Good is the Cluster-Based Algorithm?

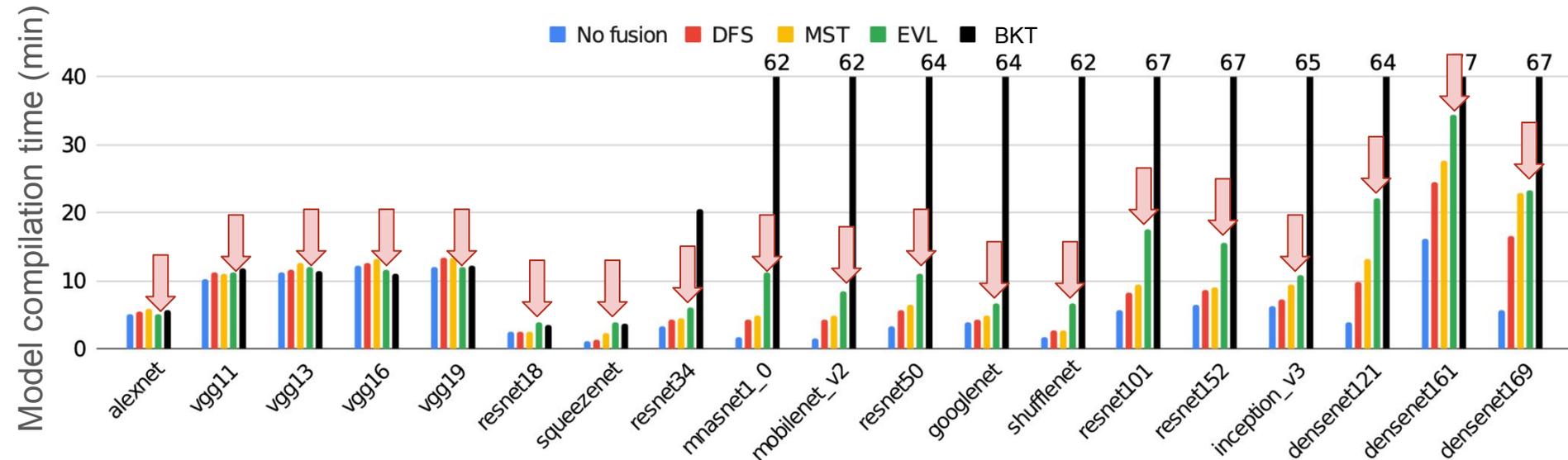


How Good is the Cluster-Based Algorithm?



The exhaustive approach experienced frequent timeouts (within one hour).

How Good is the Cluster-Based Algorithm?



The evolutionary algorithm (EVL) is 29% slower than MST

How Good is the Cluster-Based Algorithm?

- MST is 29% faster than EVL, but generates code equality as good

How Good is the Cluster-Based Algorithm?

- MST is 29% faster than EVL, but generates code equality as good
- MST is 11% slower than DFS, but generates code that is 39% faster

How Good is the Cluster-Based Algorithm?

- MST is 29% faster than EVL, but generates code equality as good
- MST is 11% slower than DFS, but generates code that is 39% faster
- MST's implementation is a bit larger than DFS's but much shorter than EVL's

How Good is the Cluster-Based Algorithm?

- MST is 29% faster than EVL, but generates code equality as good
- MST is 11% slower than DFS, but generates code that is 39% faster
- MST's implementation is a bit larger than DFS's but much shorter than EVL's
- The exhaustive approach seems unpractical in XNNC's environment

How Good is the Cluster-Based Algorithm?

- MST is 29% faster than EVL, but generates code equality as good
- MST is 11% slower than DFS, but generates code that is 39% faster
- MST's implementation is a bit larger than DFS's but much shorter than EVL's
- The exhaustive approach seems unpractical in XNNC's environment
 - Ex.: densenet169 has 1376 operators and 2231 edges

Fernando Pereira

fernando@dcc.ufmg.br

<http://lac.dcc.ufmg.br>



*Michael Canesche, Vanderson Rosário, Edson Borin, Fernando Magno Quintão Pereira: **Fusion of Operators of Computational Graphs via Greedy Clustering: The XNNC Experience.** CC 2025: 1-13*