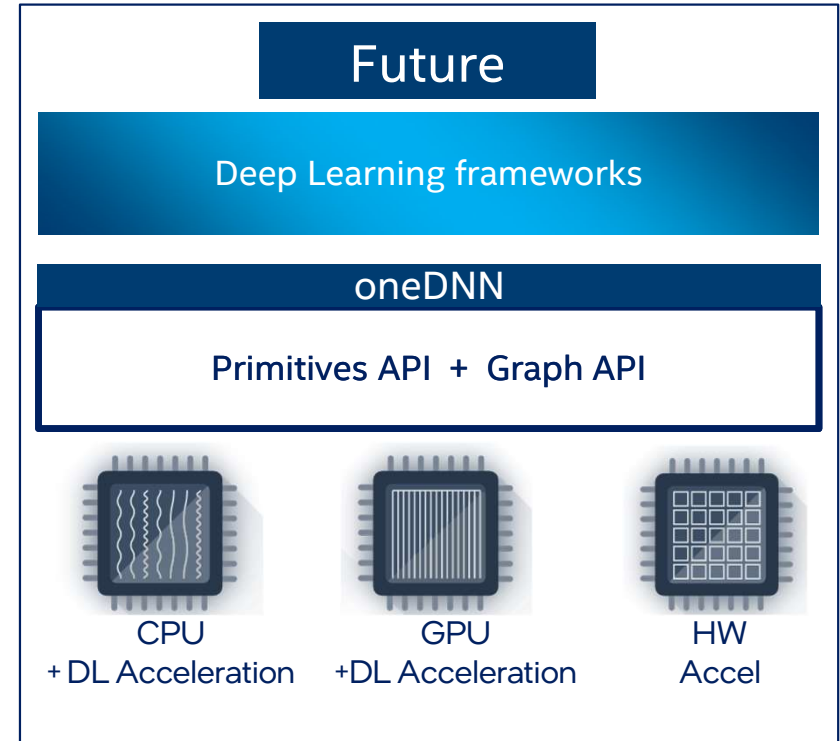
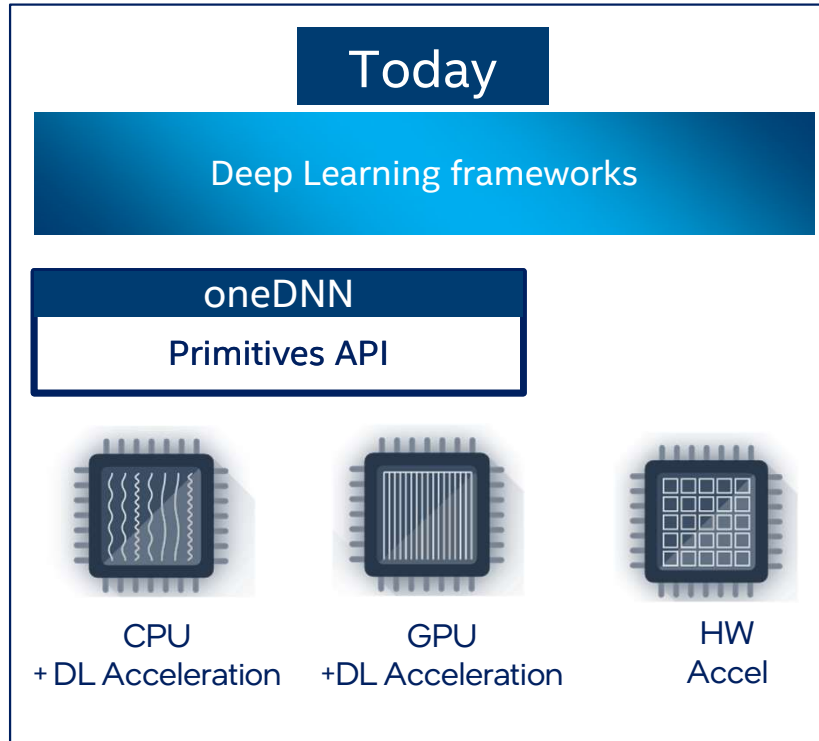


# Generating Efficient Code for Deep Neural Network

Jianhui Li  
Software and Advanced Technology Group

The AIA logo, featuring the letters "AIA" in a bold, dark blue, sans-serif font. The letters are surrounded by a grid of small, light blue squares.

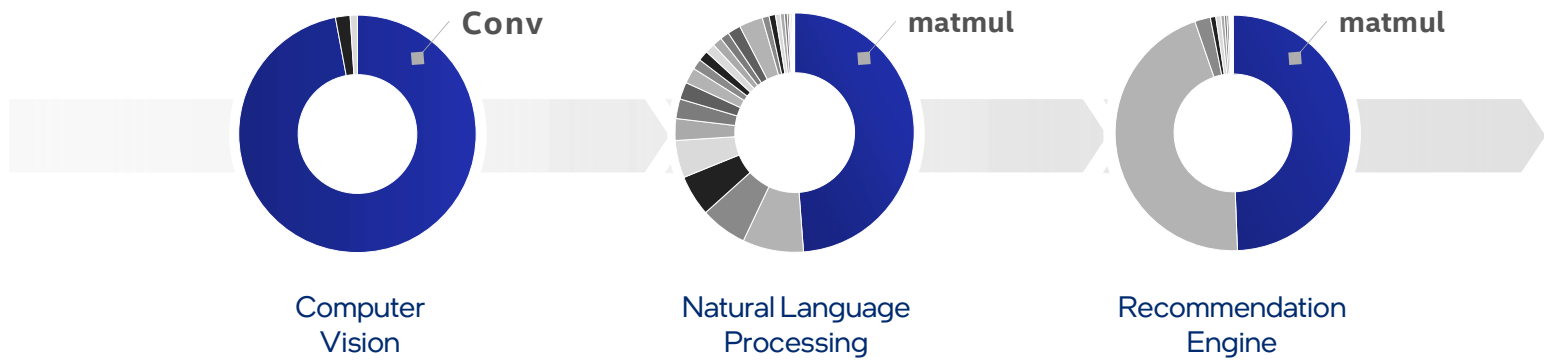
# oneDNN is evolving...



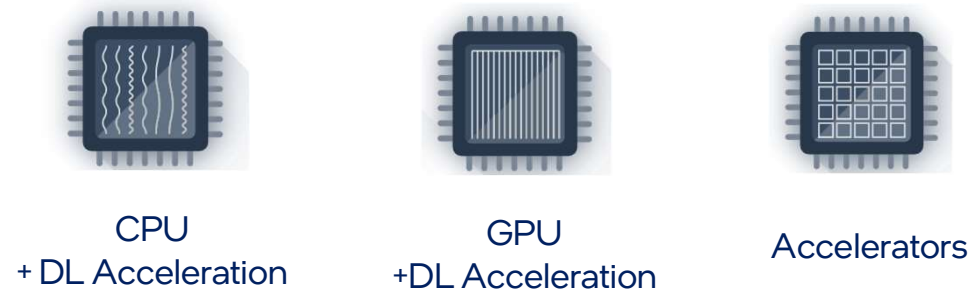
- Graph API allows HW backend to achieve high efficiency
- Same integration for multiple AI HW: CPU, GPU, and accelerators

# The driving forces of AI Optimization

Diversifying AI application



Hardware Acceleration for AI

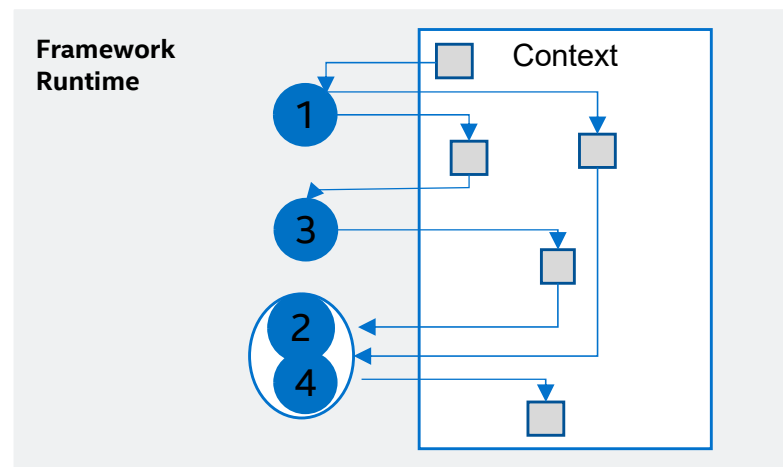
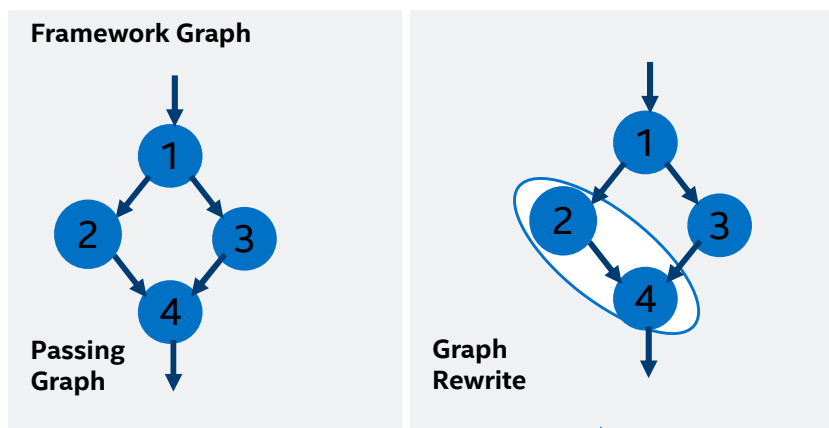


# oneDNN Graph API

DL Framework

oneDNN Graph API

oneDNN Graph Impl.

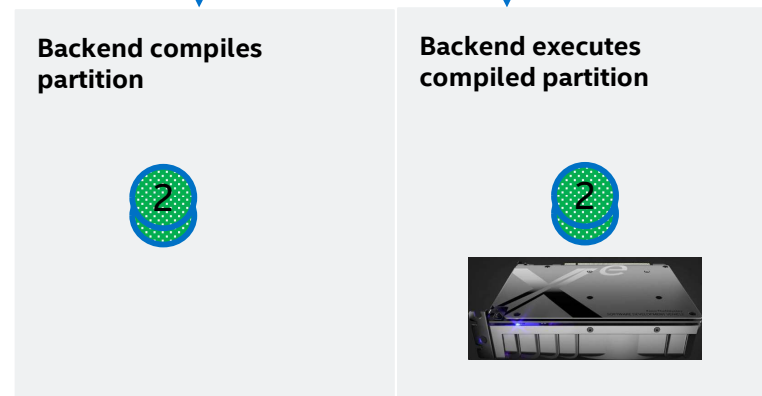
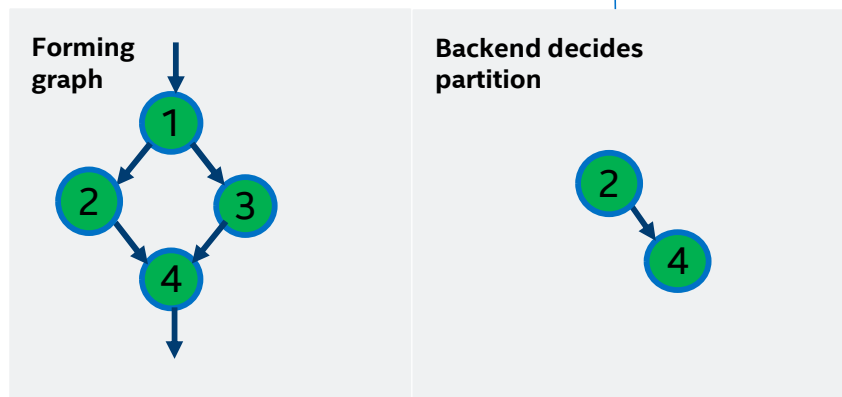


add\_op()

get\_partitions()

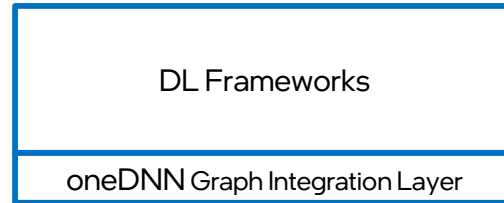
compile()

execute()



Designed to Interoperate and compose with framework components

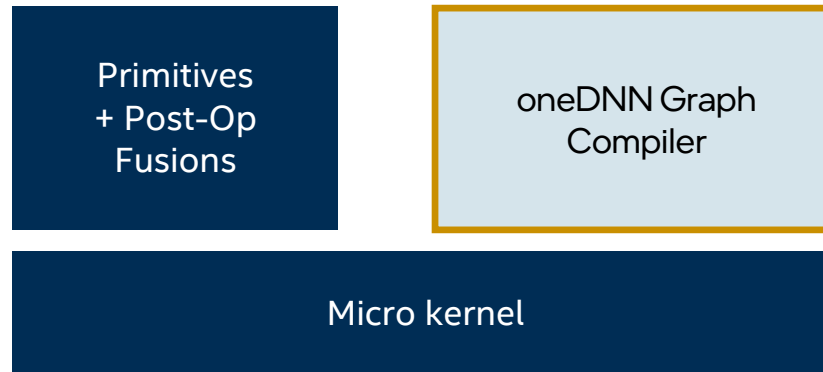
# Deep learning with oneDNN Graph



oneDNN Graph API

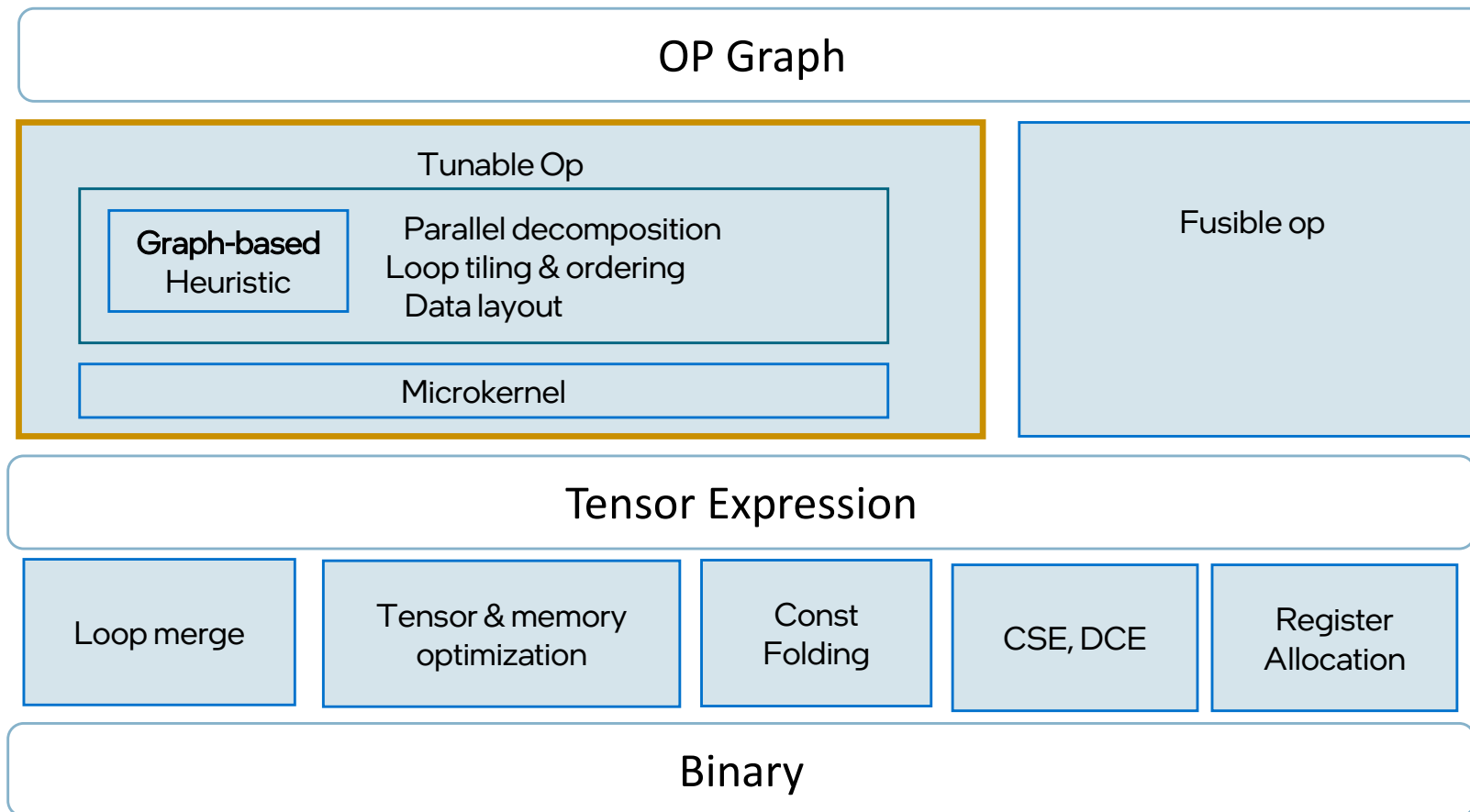


oneDNN Graph Implementation

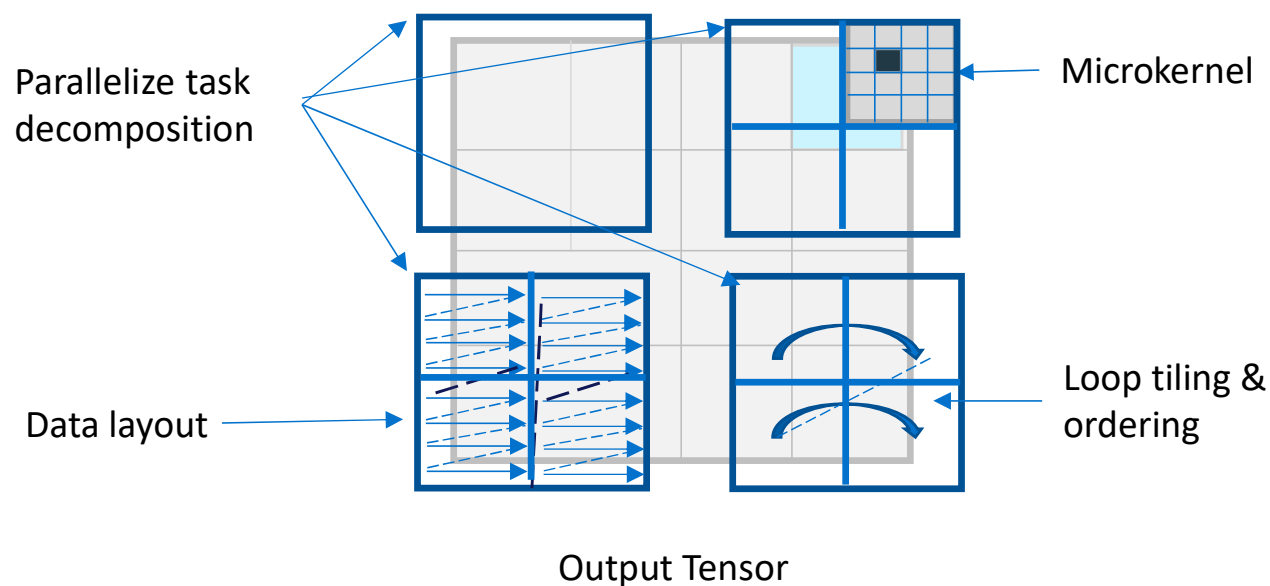


low-level graph compiler focusing on code generation for performance-critical DNN graph partition

# oneDNN Graph Compiler

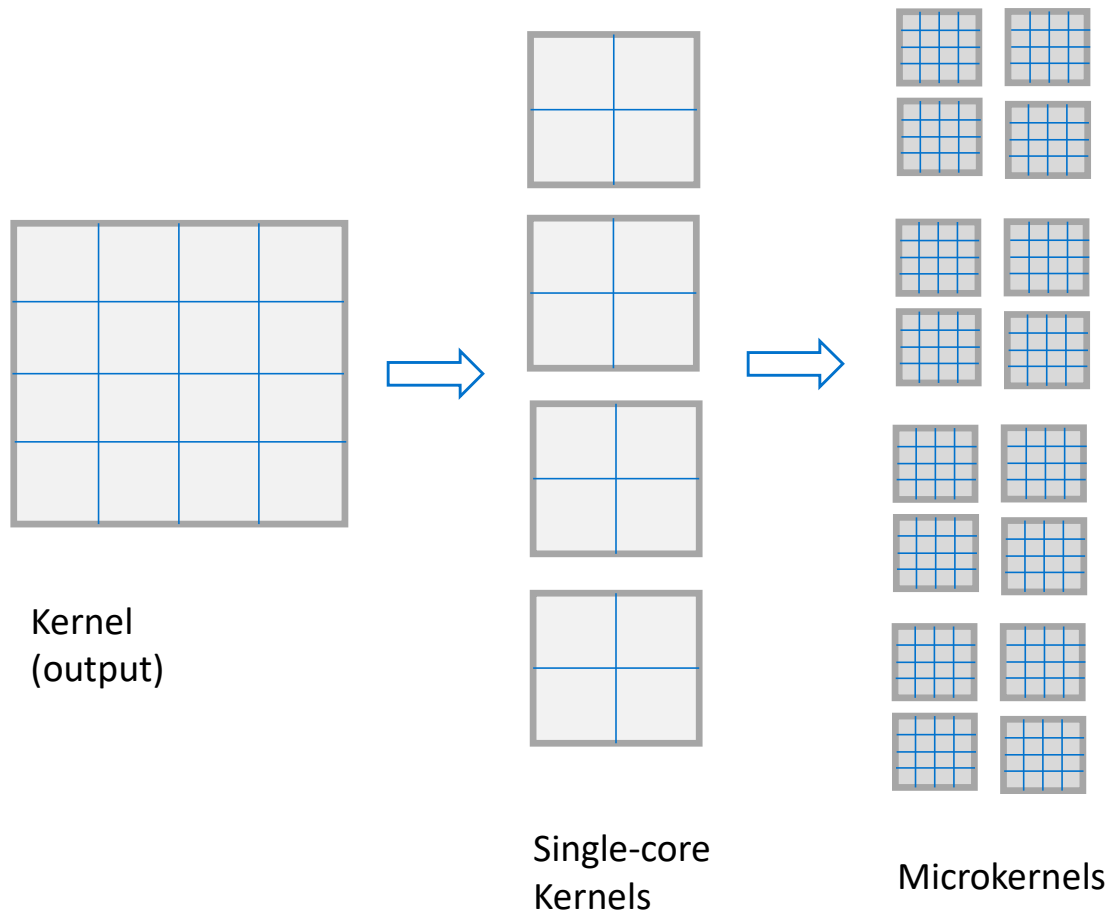


# Recipe for efficient kernel



Tunable hyperparameters	Hardware efficiency
Single-core task size	Multi-Core
Microkernel size	Vector/Matrix
Loop tiling size Loop order	Cache/TLB
Data blocking factor	Cache/TLB

# Parallel task decomposition



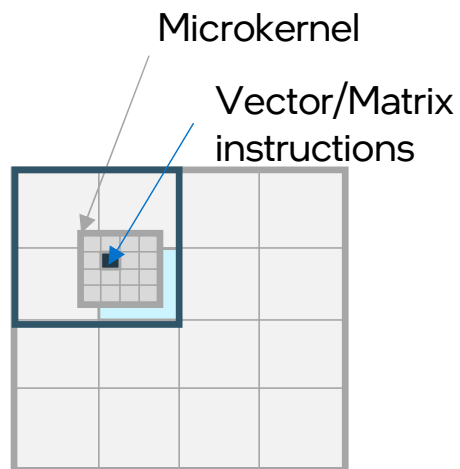
Propose single-core kernel size options, using all cores with balanced load

Propose options of microkernel sizes with high vector/matrix utilization

Search for a pair of single kernel size and microkernel size based on a cost model



# Microkernel



$$C[M, N] = \sum_{i=0}^{batch} A_i[M, K] * B_i[K, N]$$

Batch-Reduce gemm interface

High-Performance Deep Learning via a Single Building Block, Evangelos Georganas, Alexander Heinecke

## Microkernel code example

$$C[32,32] = \sum_{i=0}^{16} A_i[32,32] * B_i[32,32]$$

```
Batch_reduce_gemm(A[0:32, 0:32],
                  B[0:32, 0:32],
                  C[0:32, 0:32],
                  batch = 16)
{
  tileloaddt1(tmm0, C);
  tileloaddt1(tmm1, C + 64);
  tileloaddt1(tmm2, C + 2048);
  tileloaddt1(tmm3, C + 2112);
  for (int i = 0; i < batch; i++) {
    tileloaddt1(tmm4, A);
    tileloaddt1(tmm5, A + 1024);
    tileloaddt1(tmm6, B);
    tdbpf16ps(tmm0, tmm4, tmm6);
    tdbpf16ps(tmm2, tmm5, tmm6);
    tileloaddt1(tmm7, B + 64);
    tdbpf16ps(tmm1, tmm4, tmm7);
    tdbpf16ps(tmm3, tmm5, tmm7);
    A += 1024; B += 1024;
  }
  tilestored(C, tmm0);
  tilestored(C + 64, tmm1);
  tilestored(C + 2048, tmm2);
  tilestored(C + 2112, tmm3);
}
```

Specialized for specific tensor shape

Aligned with HW vector/matrix size

Unroll, software pipelining, prefetching

Vectorize/Tensorize

Careful and "manual" register blocking and allocation

Based on carefully hand-crafted code template

# Synthesized kernel for fused op

Parallelized for multiple core

Loop ordering & tiling

Microkernel

```

Parallel loop m_p = 0, M, MBP {
  Parallel loop n_p = 0, N, NBP {
    Loop m_o = m_p * MBP, (m_p + 1) * MBP, MB {
      Loop n_o = n_p * NBP, (n_p + 1) * NBP, NB {
        Loop k_o = 0, K/KB {
          A'[0:MB, 0:KB] = A[m_o, k_o, 0:MB, 0:KB];
          B'[0:KB, 0:NB] = B[k_o, n_o, 0:KB, 0:NB]
          Call Batch_reduce_gemm(A'[0:MB, 0:KB],
                                B'[0:NB, 0:KB],
                                C'[0:MB, 0:NB], Batch = 1);
        }
        C[m_o, n_o, 0:MB, 0:NB] = Relu(C'[0:MB, 0:NB])
      }
    }
  }
}
    
```

Blocked Layout

Potential place to apply fusion

Good heuristic required to decompose to single-core kernel and microkernel, loop order, and data layout

# Compile subgraph by composing kernels

matmul  
relu

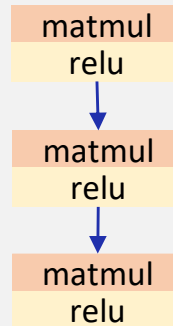
C - Output tensor, [M][N]  
A - Input tensor (activations) [M][K]  
B - Input tensor (weights) [K][N]

```
Parallel loop m_p = 0, M, MBP {
  Parallel loop n_p = 0, N, NBP {
    Loop m_o = m_p * MBP, (m_p + 1) * MBP, MB {
      Loop n_o = n_p * NBP, (n_p + 1) * NBP, NB {
        Loop k_o = 0, K/KB {
          A'[0:MB, 0:KB] = A[m_o, k_o, 0:MB, 0:KB];
          B'[0:KB, 0:NB] = B[k_o, n_o, 0:KB, 0:NB]
          Call Batch_reduce_gemm(A'[0:MB, 0:KB],
                                B'[0:NB, 0:KB],
                                C'[0:MB, 0:NB]);
        }
        C[m_o, n_o, 0:MB, 0:NB] = relu ( C'[0:MB, 0:NB])
      }
    }
  }
}
```

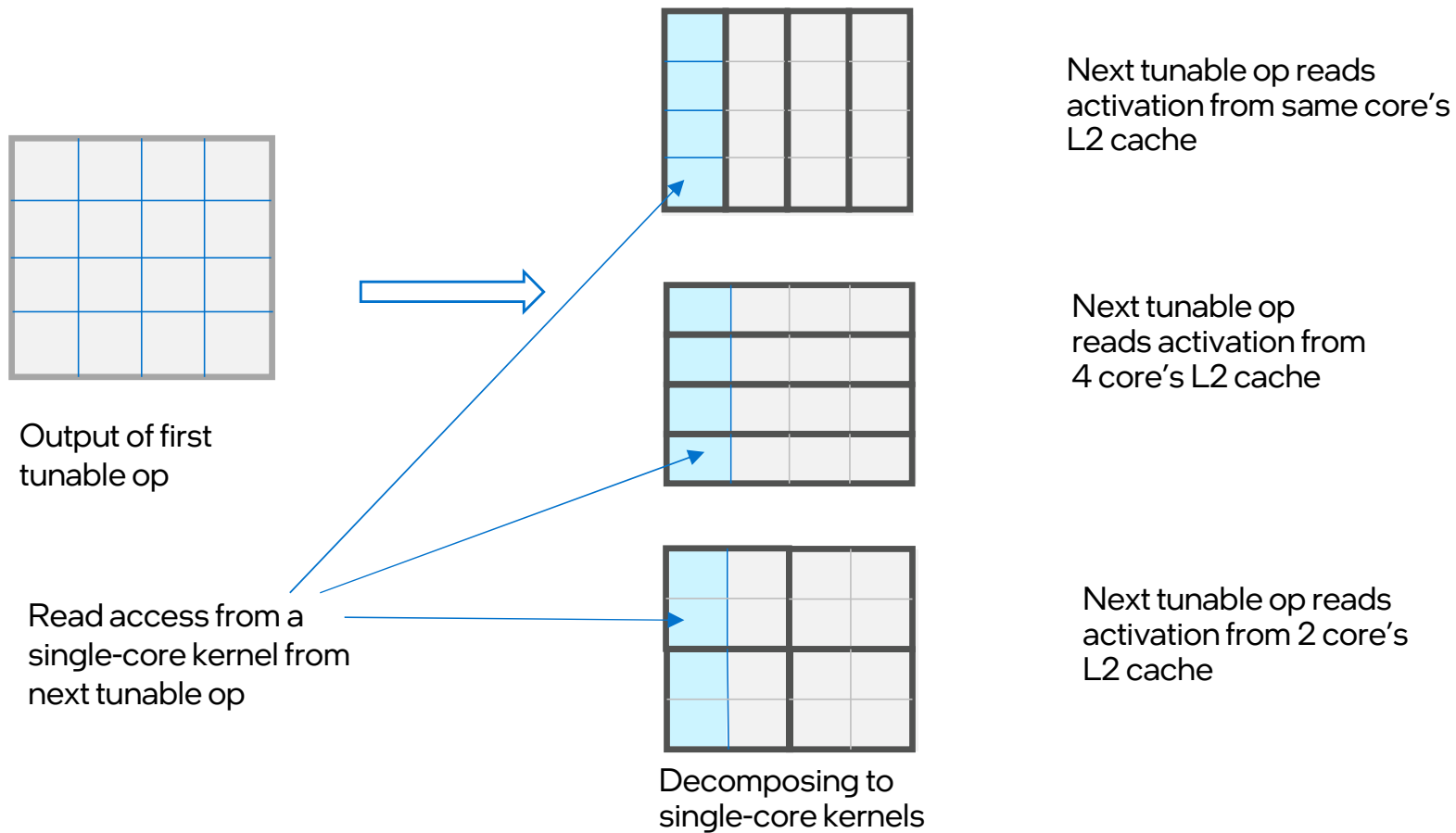
```
Parallel loop m_p = 0, M, MBP {
  Parallel loop n_p = 0, N, NBP {
    Loop m_o = m_p * MBP, (m_p + 1) * MBP, MB {
      Loop n_o = n_p * NBP, (n_p + 1) * NBP, NB {
        Loop k_o = 0, K/KB {
          A'[0:MB, 0:KB] = A[m_o, k_o, 0:MB, 0:KB];
          B'[0:KB, 0:NB] = B[k_o, n_o, 0:KB, 0:NB]
          Call Batch_reduce_gemm(A'[0:MB, 0:KB],
                                B'[0:NB, 0:KB],
                                C'[0:MB, 0:NB]);
        }
        C[m_o, n_o, 0:MB, 0:NB] = relu ( C'[0:MB, 0:NB])
      }
    }
  }
}
```

```
Parallel loop m_p = 0, M, MBP {
  Parallel loop n_p = 0, N, NBP {
    Loop m_o = m_p * MBP, (m_p + 1) * MBP, MB {
      Loop n_o = n_p * NBP, (n_p + 1) * NBP, NB {
        Loop k_o = 0, K/KB {
          A'[0:MB, 0:KB] = A[m_o, k_o, 0:MB, 0:KB];
          B'[0:KB, 0:NB] = B[k_o, n_o, 0:KB, 0:NB]
          Call Batch_reduce_gemm(A'[0:MB, 0:KB],
                                B'[0:NB, 0:KB],
                                C'[0:MB, 0:NB]);
        }
        C[m_o, n_o, 0:MB, 0:NB] = relu ( C'[0:MB, 0:NB])
      }
    }
  }
}
```

```
Parallel loop m_p = 0, M, MBP {
  Parallel loop n_p = 0, N, NBP {
    Loop m_o = m_p * MBP, (m_p + 1) * MBP, MB {
      Loop n_o = n_p * NBP, (n_p + 1) * NBP, NB {
        Loop k_o = 0, K/KB {
          A'[0:MB, 0:KB] = A[m_o, k_o, 0:MB, 0:KB];
          B'[0:KB, 0:NB] = B[k_o, n_o, 0:KB, 0:NB]
          Call Batch_reduce_gemm(A'[0:MB, 0:KB],
                                B'[0:NB, 0:KB],
                                C'[0:MB, 0:NB]);
        }
        C[m_o, n_o, 0:MB, 0:NB] = relu ( C'[0:MB, 0:NB])
      }
    }
  }
}
```

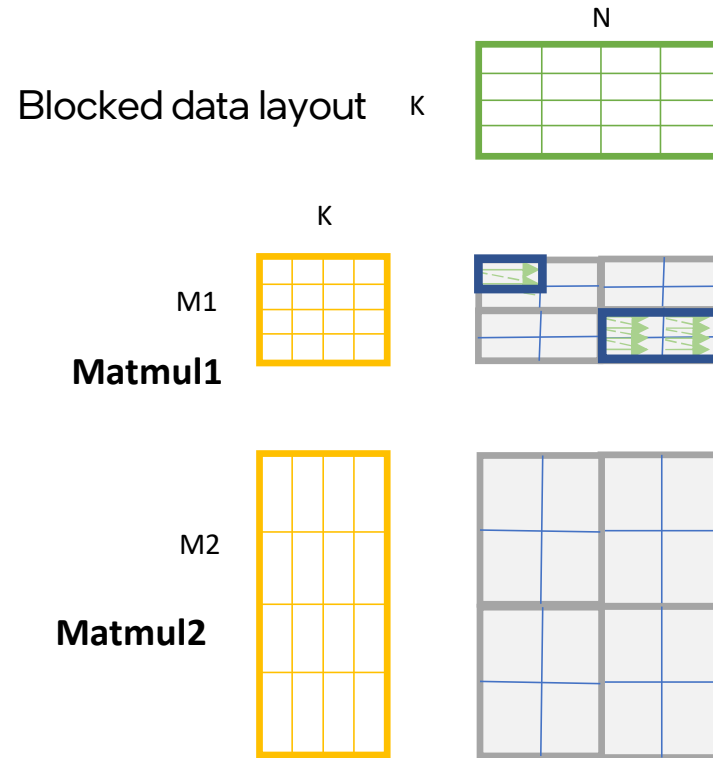
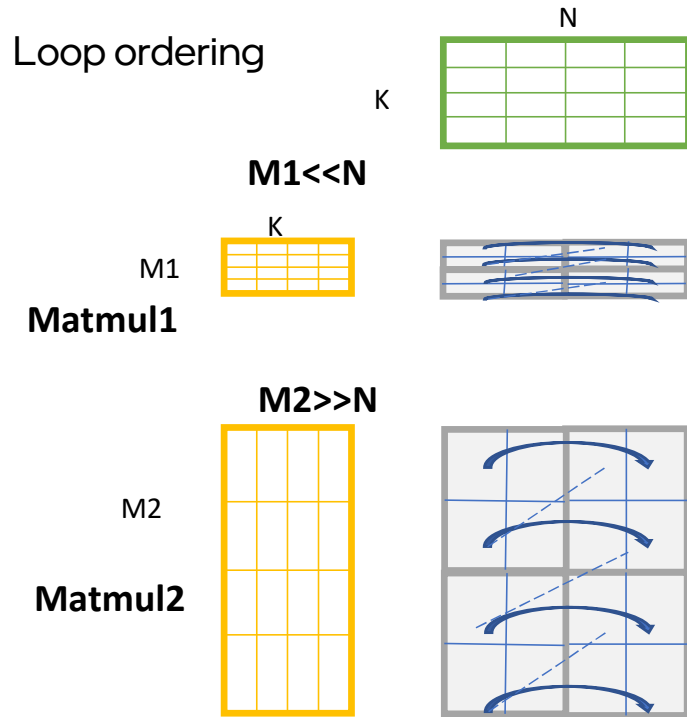


# Parallel task decomposition within Graph context



Graph-based heuristic considers the cost of reading activation data from last tunable op

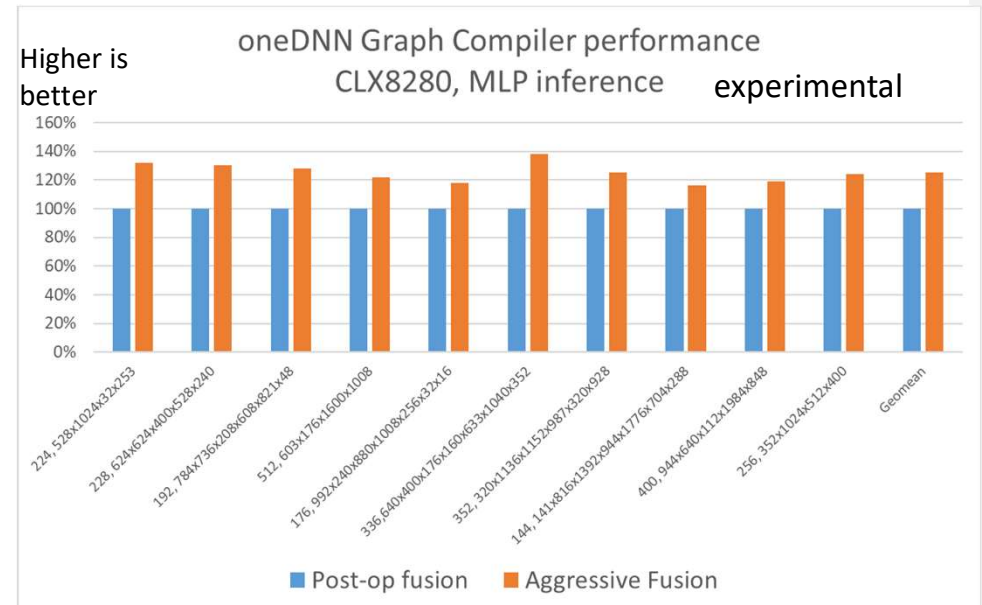
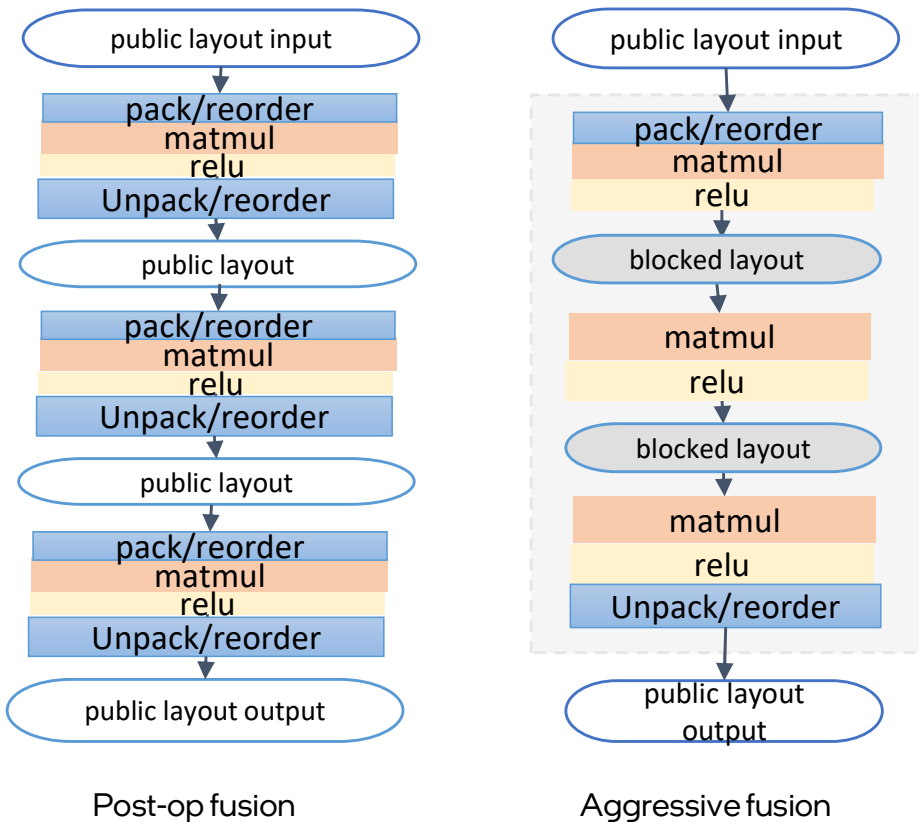
# Graph-based heuristic for Single-core Kernel



Graph knows input/weight is "hot" or "cold"  
 Input stationary preferred as weight is usually "cold"

Choose blocked layout works best  
 for overall performance

# Aggressive Fusion



Aggressive fusion covers MLP, MHA, Conv Block, both inference and training, fp32, int8, and bf16 on CPU

Substantial performance gains observed

# Summary

oneDNN Graph API offers flexible API to accelerate DNN Graph with partitioning

oneDNN Graph compiler merges nested loops of tuable op and fusible op, compose a large kernel for aggressive fusion, and Automate kernel generation

Recipes from hand-tuned kernel, micro-kernel, and blocked layout

Graph-based heuristic further improves kernel performance within graph context

# Notices and Disclaimers

Intel technologies may require enabled hardware, software or service activation.

No product or component can be absolutely secure.

Your costs and results may vary.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.

