

# Introduction to Modern Compilers (part of E0255)

Uday Kumar Reddy B

*udayb@iisc.ac.in*

**Dept of CSA**  
**Indian Institute of Science**

- **Current:**
  - C, C++, Rust, Java, Python, MATLAB, R, ...
- What will the new and disruptive programming technologies of the 21st century be?

- **Current:**
  - C, C++, Rust, Java, Python, MATLAB, R, ...
- **What will the new and disruptive programming technologies of the 21st century be?**

- ❶ What do programmers want?
- ❷ How are architectures evolving?
  - Multiple cores and many cores on a chip
  - GPUs, accelerators, and heterogeneous parallel architectures
  - Wider vector processing units
  - Deep memory hierarchies
  - Reduced precision

# HIGH-PERFORMANCE COMPILATION: WHAT DO YOU WANT TO PROGRAM?

- Scientific and engineering simulations
  - Eg: Solving partial differential equations numerically
- Embedded vision (Eg: Autonomous/self-driving cars)
- Smartphones — HPC in data centers and cloud drives a number of smartphone technologies
- **Scientific and Engineering simulations**
- **Data Analytics**
- **Deep Learning**
- **Generative AI, LLMs**

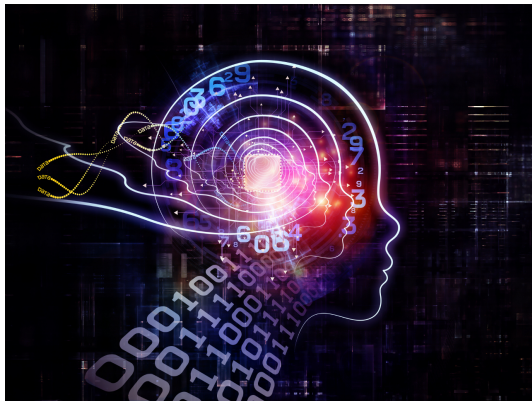
# QUESTIONS TO THINK ABOUT

- What will the new programming technologies for the emerging domains be?
  - **Current:** C, C++, Rust with OpenMP, MPI, CUDA, OpenCL
  - **Future:** New languages, compilers, libraries, and DSLs



# QUESTIONS TO THINK ABOUT

- **What will the new programming technologies for AI be?**
  - PyTorch is dominant today; JAX is another high-level one. OpenAI Triton is mid-level.
  - **Just scratches the surface**



# THE NEED FOR HIGH PERFORMANCE

- **More/Larger Data**
  - Instagram — 60 million photos / day
  - YouTube — 100 hours of video uploaded every minute
- **Need for a fast/real-time response in some domains**
- **More complex algorithms**
- **Science/Engineering simulations/modeling: Time to solution**



# PROGRAMMING MODERN HARDWARE EFFECTIVELY

- Compute speed: Eg.: 16 multiply-adds per cycle (AVX-512 unit for fp32)
- Synchronization (2 cores  $0.25\ \mu\text{s}$ , 8 cores  $1.25\ \mu\text{s}$ , 2x8 cores  $1.54\ \mu\text{s}$ )
- Memory bandwidth ( 10 GB/s per core, 500 GB/s per socket)
- High-Performance Programming and Compilation
  - Exploiting locality (caches, registers)
  - Exploit single core hardware well (vectorization, ...)
  - Multi-core parallelism
  - Reduce synchronization and communication as much as possible
- Good scaling without good single thread performance is a great waste of resources (power, equipment cost)

# PROGRAMMING MODERN HARDWARE EFFECTIVELY

- Compute speed: Eg.: 16 multiply-adds per cycle (AVX-512 unit for fp32)
- Synchronization (2 cores  $0.25\ \mu s$ , 8 cores  $1.25\ \mu s$ , 2x8 cores  $1.54\ \mu s$ )
- Memory bandwidth ( 10 GB/s per core, 500 GB/s per socket)
- High-Performance Programming and Compilation
  - Exploiting locality (caches, registers)
  - Exploit single core hardware well (vectorization, ...)
  - Multi-core parallelism
  - Reduce synchronization and communication as much as possible
- Good scaling without good single thread performance is a great waste of resources (power, equipment cost)

# A CLASSIFICATION OF VARIOUS APPROACHES

- ➊ Manual low-level (C, C++) with parallel programming models (OpenMP, CUDA, MPI) with the best optimizing compilers
- ➋ Library-based: C, C++, Python with libraries/packages: MKL, ScaLAPACK, CuBLAS, CuDNN, Cutlass, CuB
- ➌ Mid-level: Triton, CuTile, Pallas
- ➍ Ultra-high level languages and models including embedded DSLs: Tensorflow, PyTorch, JAX, R, MATLAB, Halide, Spiral
- ➎ General goal: Obtain productivity of the last class and the performance of the first

# A CLASSIFICATION OF VARIOUS APPROACHES

- ➊ Manual low-level (C, C++) with parallel programming models (OpenMP, CUDA, MPI) with the best optimizing compilers
- ➋ Library-based: C, C++, Python with libraries/packages: MKL, ScaLAPACK, CuBLAS, CuDNN, Cutlass, CuB
- ➌ Mid-level: Triton, CuTile, Pallas
- ➍ Ultra-high level languages and models including embedded DSLs: Tensorflow, PyTorch, JAX, R, MATLAB, Halide, Spiral
- ➎ General goal: Obtain productivity of the last class and the performance of the first

# A CLASSIFICATION OF VARIOUS APPROACHES

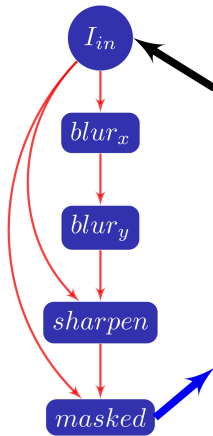
- ① Manual low-level (C, C++) with parallel programming models (OpenMP, CUDA, MPI) with the best optimizing compilers
  - ② Library-based: C, C++, Python with libraries/packages: MKL, ScaLAPACK, CuBLAS, CuDNN, Cutlass, CuB
  - ③ Mid-level: Triton, CuTile, Pallas
  - ④ Ultra-high level languages and models including embedded DSLs: Tensorflow, PyTorch, JAX, R, MATLAB, Halide, Spiral
- **General goal:** Obtain productivity of the last class and the performance of the first

# A CLASSIFICATION OF VARIOUS APPROACHES

- ① Manual low-level (C, C++) with parallel programming models (OpenMP, CUDA, MPI) with the best optimizing compilers
  - ② Library-based: C, C++, Python with libraries/packages: MKL, ScaLAPACK, CuBLAS, CuDNN, Cutlass, CuB
  - ③ Mid-level: Triton, CuTile, Pallas
  - ④ Ultra-high level languages and models including embedded DSLs: Tensorflow, PyTorch, JAX, R, MATLAB, Halide, Spiral
- **General goal:** Obtain productivity of the last class and the performance of the first

# EXAMPLE: UNSHARP MASK – AN IMAGE PROCESSING PIPELINE

(C) Bernie Saunders, CC BY-NC-ND 3.0



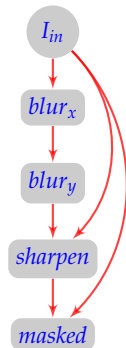
# UNSHARP MASK: COMPUTATION

```
for (i = 0; i <= 2; i++)
  for (j = 2; j <= (R + 1); j++)
    for (k = 0; (k <= (C + 3)); k++)
      blurx[i][j-2][k] = img[i][j-2][k]*0.0625f + img[i][j-1][k]*0.25f
        + img[i][j][k]*0.375f + img[i][j+1][k]*0.25f + img[i][j+2][k]*0.0625f;

for (i = 0; (i <= 2); i++)
  for (j = 2; (j <= (R + 1)); j++)
    for (k = 2; (k <= (C + 1)); k++)
      blurry[i][j][k-2] = blurx[i][j-2][k-2]*0.0625f + blurx[i][j-2][k-1]*0.25f
        + blurx[i][j-2][k]*0.375f + blurx[i][j-2][k+1]*0.25f + blurx[i][j-2][k+2]*0.0625f;

for (i = 0; (i <= 2); i++)
  for (j = 2; (j <= (R + 1)); j++)
    for (k = 2; (k <= (C + 1)); k++)
      sharpen[i][j][k-2] = img[i][j][k]*(1 + weight) + blurry[i][j-2][k-2]*(-weight);

for (i = 0; i <= 2; i++)
  for (j = 2; j <= R + 1; j++)
    for (k = 2; k <= C + 1; k++) {
      _ct0 = img[i][j][k];
      _ct1 = sharpen[i][j-2][k-2];
      _ct2 = (std::abs((img[i][j][k] - blurry[i][j-2][k-2])) < threshold)? _ct0: _ct1;
      mask[i][j-2][k-2] = _ct2;
    }
```



A sequential version in C: **18.6 ms** / frame  
(using GCC with opts, quad-core Nehalem, 720p video)



# UNSHARP MASK - A NAIVE OPENMP VERSION

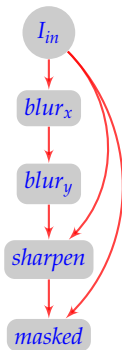
```
for (i = 0; i <= 2; i++)
#pragma omp parallel for
  for (j = 2; j <= (R + 1); j++)
#pragma omp ivdep
    for (k = 0; k <= C + 3; k++)
      blurx[i][j-2][k] = img[i][j-2][k]*0.0625f + img[i][j-1][k]*0.25f
        + img[i][j][k]*0.375f + img[i][j+1][k]*0.25f + img[i][j+2][k]*0.0625f;

for (i = 0; i <= 2; i++)
#pragma omp parallel for
  for (j = 2; j <= R + 1; j++)
#pragma omp ivdep
    for (k = 2; k <= C + 1; k++)
      blury[i][j][k-2] = blurx[i][j-2][k-2]*0.0625f + blurx[i][j-2][k-1]*0.25f
        + blurx[i][j-2][k]*0.375f + blurx[i][j-2][k+1]*0.25f + blurx[i][j-2][k+2]*0.0625f;

for (i = 0; i <= 2; i++)
#pragma omp parallel for
  for (j = 2; j <= R + 1; j++)
#pragma omp ivdep
    for (k = 2; k <= C + 1; k++)
      sharpen[i][j][k-2] = img[i][j][k]*(1 + weight) + blury[i][j-2][k-2]*(-weight);

for (i = 0; i <= 2; i++)
#pragma omp parallel for private(_ct0,_ct1,_ct2)
  for (j = 2; j <= R + 1; j++)
#pragma omp ivdep
    for (k = 2; k <= C + 1; k++) {
      _ct0 = img[i][j][k];
      _ct1 = sharpen[i][j-2][k-2];
      _ct2 = (std::abs((img[i][j][k] - blury[i][j-2][k-2])) < threshold)? _ct0: _ct1;
      mask[i][j-2][k-2] = _ct2;
    }
}
```

20.2 ms / frame on 1 thread, 18.02 ms / frame on 4 threads



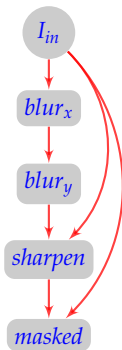
# UNSHARP MASK - A BETTER OPENMP VERSION

```
#pragma omp parallel for
for (j = 2; j <= (R + 1); j++)
    for (i = 0; i <= 2; i++)
        #pragma ivdep
        for (k = 0; (k <= (C + 3)); k++)
            blurx[i][j-2][k] = img[i][j-2][k]*0.0625f + img[i][j-1][k]*0.25f
            + img[i][j][k]*0.375f + img[i][j+1][k]*0.25f + img[i][j+2][k]*0.0625f;

#pragma omp parallel for
for (j = 2; (j <= (R + 1)); j++)
    for (i = 0; i <= 2; i++)
        #pragma ivdep
        for (k = 2; (k <= (C + 1)); k++)
            blurry[i][j][k-2] = blurx[i][j-2][k-2]*0.0625f + blurx[i][j-2][k-1]*0.25f
            + blurx[i][j-2][k]*0.375f + blurx[i][j-2][k+1]*0.25f + blurx[i][j-2][k+2]*0.0625f;

#pragma omp parallel for
for (j = 2; (j <= (R + 1)); j++)
    for (i = 0; i <= 2; i++)
        #pragma ivdep
        for (k = 2; (k <= (C + 1)); k++)
            sharpen[i][j][k-2] = img[i][j][k]*(1 + weight) + blurry[i][j-2][k-2]*(-weight);

#pragma omp parallel for private(_ct0,_ct1,_ct2)
for (j = 2; j <= R + 1; j++)
    for (i = 0; i <= 2; i++)
        #pragma ivdep
        for (k = 2; k <= C + 1; k++) {
            _ct0 = img[i][j][k];
            _ct1 = sharpen[i][j-2][k-2];
            _ct2 = (std::abs((img[i][j][k] - blurry[i][j-2][k-2])) < threshold)? _ct0: _ct1;
            mask[i][j-2][k-2] = _ct2;
        }
```



**18.6 ms / frame on 1 thread, 15.03 ms / frame on 4 threads**

# OPTIMIZING UNSHARP MASK

## ❶ Write with OpenCV library (with Python bindings)

```
@jit("float32[:] (uint8[:], int64)", cache = True, nogil = True)
def unsharp_cv(frame, lib_func):
    frame_f = np.float32(frame) / 255.0
    res = frame_f
    kernelx = np.array([1, 4, 6, 4, 1], np.float32) / 16
    kernely = np.array([[1], [4], [6], [4], [1]], np.float32) / 16
    blur = sepFilter2D(frame_f, -1, kernelx, kernely)
    sharpen = addWeighted(frame_f, (1 + weight), blur, (-weight), 0)
    th, choose = threshold(absdiff(frame_f, blur), thresh, 1, THRESH_BINARY)
    choose = choose.astype(bool)
    np.copyto(res, sharpen, 'same_kind', choose)
    return res
```

Performance: **35.9 ms** / frame

- ❷ Write in a dynamic language like Python and use a JIT (Numba) — performance: **79 ms** / frame
- ❸ A naive C version parallelized with OpenMP: **18.02 ms** / frame
- ❹ A version with sophisticated optimizations (fusion + overlapped tiling): **8.97 ms** / frame (in this course, we will study how to get to this, and build compilers/code generators that can achieve this automatically)

## • Video demo

# OPTIMIZING UNSHARP MASK

## ❶ Write with OpenCV library (with Python bindings)

```
@jit("float32[:] (uint8[:], uint64)", cache = True, nogil = True)
def unsharp_cv(frame, lib_func):
    frame_f = np.float32(frame) / 255.0
    res = frame_f
    kernelx = np.array([1, 4, 6, 4, 1], np.float32) / 16
    kernely = np.array([[1], [4], [6], [4], [1]], np.float32) / 16
    blur = sepFilter2D(frame_f, -1, kernelx, kernely)
    sharpen = addWeighted(frame_f, (1 + weight), blur, (-weight), 0)
    th, choose = threshold(absdiff(frame_f, blur), thresh, 1, THRESH_BINARY)
    choose = choose.astype(bool)
    np.copyto(res, sharpen, 'same_kind', choose)
    return res
```

Performance: **35.9 ms** / frame

## ❷ Write in a dynamic language like Python and use a JIT (Numba) — performance: **79 ms** / frame

## ❸ A naive C version parallelized with OpenMP: **18.02 ms** / frame

## ❹ A version with sophisticated optimizations (fusion + overlapped tiling): **8.97 ms** / frame (in this course, we will study how to get to this, and build compilers/code generators that can achieve this automatically)

## • Video demo

# OPTIMIZING UNSHARP MASK

## ❶ Write with OpenCV library (with Python bindings)

```
@jit("float32[:] (uint8[:], uint64)", cache = True, nogil = True)
def unsharp_cv(frame, lib_func):
    frame_f = np.float32(frame) / 255.0
    res = frame_f
    kernelx = np.array([1, 4, 6, 4, 1], np.float32) / 16
    kernely = np.array([[1], [4], [6], [4], [1]], np.float32) / 16
    blur = sepFilter2D(frame_f, -1, kernelx, kernely)
    sharpen = addWeighted(frame_f, (1 + weight), blur, (-weight), 0)
    th, choose = threshold(absdiff(frame_f, blur), thresh, 1, THRESH_BINARY)
    choose = choose.astype(bool)
    np.copyto(res, sharpen, 'same_kind', choose)
    return res
```

Performance: **35.9 ms** / frame

## ❷ Write in a dynamic language like Python and use a JIT (Numba) — performance: **79 ms** / frame

## ❸ A naive C version parallelized with OpenMP: **18.02 ms** / frame

## ❹ A version with sophisticated optimizations (fusion + overlapped tiling): **8.97 ms** / frame (in this course, we will study how to get to this, and build compilers/code generators that can achieve this automatically)

## • Video demo

# OPTIMIZING UNSHARP MASK

## ❶ Write with OpenCV library (with Python bindings)

```
@jit("float32[:,:](uint8[:,:],_int64)", cache = True, nogil = True)
def unsharp_cv(frame, lib_func):
    frame_f = np.float32(frame) / 255.0
    res = frame_f
    kernelx = np.array([1, 4, 6, 4, 1], np.float32) / 16
    kernely = np.array([[1], [4], [6], [4], [1]], np.float32) / 16
    blur = sepFilter2D(frame_f, -1, kernelx, kernely)
    sharpen = addWeighted(frame_f, (1 + weight), blur, (-weight), 0)
    th, choose = threshold(absdiff(frame_f, blur), thresh, 1, THRESH_BINARY)
    choose = choose.astype(bool)
    np.copyto(res, sharpen, 'same_kind', choose)
    return res
```

Performance: **35.9 ms** / frame

- ❷ Write in a dynamic language like Python and use a JIT (Numba) — performance: **79 ms** / frame
- ❸ A naive C version parallelized with OpenMP: **18.02 ms** / frame
- ❹ A version with sophisticated optimizations (fusion + overlapped tiling): **8.97 ms** / frame (in this course, we will study how to get to this, and build compilers/code generators that can achieve this automatically)

## ● Video demo

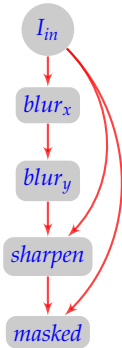
## UNSHARP MASK - A HIGHLY OPTIMIZED VERSION

**Note:** Code below is indicative and not meant for reading. Zoom into soft copy or browse source code repo listed in references.

```

#pragma omp parallel for schedule(static)
for (int _i1 = 0; (_i1 <= (R + 1) / 32); _i1 = (_i1 + 1))
{
    int _ct0 = (((R + 1) < ((32 + _T.i1) + 31)) ? (R + 1) : ((32 + _T.i1) + 31));
    int _ct1 = ((2 > (32 + _T.i1)) ? 2 : (32 + _T.i1));
    int _ct4 = (((R + 1) < ((32 + _T.i1) + 31)) ? (R + 1) : ((32 + _T.i1) + 31));
    int _ct5 = ((2 > (32 + _T.i1)) ? 2 : (32 + _T.i1));
    int _ct8 = (((R + 1) < ((32 + _T.i1) + 31)) ? (R + 1) : ((32 + _T.i1) + 31));
    int _ct9 = ((2 > (32 + _T.i1)) ? 2 : (32 + _T.i1));
    int _ct12 = (((R + 1) < ((32 + _T.i1) + 31)) ? (R + 1) : ((32 + _T.i1) + 31));
    int _ct13 = ((2 > (32 + _T.i1)) ? 2 : (32 + _T.i1));
    for (int _T.i2 = 1; _T.i2 <= ((C + 2) / 256); _T.i2 = (_T.i2 + 1))
    {
        int _ct2 = (((C + 2) < ((256 + _T.i2) + 261)) ? (C + 2) : ((256 + _T.i2) + 261));
        int _ct3 = ((8 > (256 + _T.i2)) ? 8 : (256 + _T.i2));
        int _ct4 = (((C + 2) < ((256 + _T.i2) + 260)) ? (C + 1) : ((256 + _T.i2) + 260));
        int _ct7 = ((2 > ((256 + _T.i2) + 1)) ? 2 : ((256 + _T.i2) + 1));
        int _ct10 = (((C + 1) < ((256 + _T.i2) + 259)) ? (C + 1) : ((256 + _T.i2) + 259));
        int _ct11 = ((2 > ((256 + _T.i2) + 2)) ? 2 : ((256 + _T.i2) + 2));
        int _ct14 = (((C + 1) < ((256 + _T.i2) + 258)) ? (C + 1) : ((256 + _T.i2) + 258));
        int _ct15 = ((2 > ((256 + _T.i2) + 3)) ? 2 : ((256 + _T.i2) + 3));
        for (int _i8 = 0; _i8 <= 2; _i8 = (_i8 + 1))
        {
            for (int _i1 = _ct1; _i1 <= _ct0; _i1 = (_i1 + 1))
            {
#pragma loopvec
                for (int _i2 = _ct3; _i2 <= _ct12; _i2 = (_i2 + 1))
                {
                    blury[_i0][((-32 + _T.i1) + _i2 <= ((-256 + _T.i2) + 12)) ? (((((blury[_i0][((-32 + _T.i1) + (C + 4)) + (-2 + _i1) - (C + 4)) + 12)) + 0.8625f) + (img[((-10 + ((R + 4) + (C + 4)) + (-1 + _i1) - (C + 4)) + (-256 + _T.i2) + 12)) * 0.25f) + (img[((-10 + ((R + 4) + (C + 4)) + (-1 + _i1) - (C + 4)) + 12]) * 0.25f)) : (img[((-10 + ((R + 4) + (C + 4)) + (-1 + _i1) - (C + 4)) + 12]) * 0.25f)];
                }
            }
            for (int _i8 = 0; _i8 <= 2; _i8 = (_i8 + 1))
            {
                for (int _i1 = _ct5; _i1 <= _ct4; _i1 = (_i1 + 1))
                {
#pragma loopvec
                    for (int _i2 = _ct7; _i2 <= _ct6; _i2 = (_i2 + 1))
                    {
                        blury[_i0][((-32 + _T.i1) + _i2 <= ((-256 + _T.i2) + 12)) ? (((((blury[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12)) + 0.8625f) + (blury[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12)) + 0.25f)) + (blury[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12)) + 0.25f)) : (blury[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12)) + 0.25f)];
                    }
                }
                for (int _i8 = 0; _i8 <= 2; _i8 = (_i8 + 1))
                {
                    for (int _i1 = _ct9; _i1 <= _ct8; _i1 = (_i1 + 1))
                    {
#pragma loopvec
                        for (int _i2 = _ct11; _i2 <= _ct10; _i2 = (_i2 + 1))
                        {
                            sharpen[_i0][((-32 + _T.i1) + _i2 <= ((-256 + _T.i2) + 12)) ? (((img[((-10 + ((R + 4) + (C + 4)) + (-1 + (C + 4)) + (-2) + (1 + weight)) + (blury[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12)) + 0.25f)]) - (weight));
                        }
                    }
                }
                for (int _i8 = 0; _i8 <= 2; _i8 = (_i8 + 1))
                {
                    for (int _i1 = _ct13; _i1 <= _ct12; _i1 = (_i1 + 1))
                    {
#pragma loopvec
                        for (int _i2 = _ct15; _i2 <= _ct14; _i2 = (_i2 + 1))
                        {
                            float _ct16 = img[((-10 + ((R + 4) + (C + 4)) + (-1 + (C + 4)) + (-2))];
                            float _ct17 = sharpen[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12))];
                            float _ct18 = ((ct16 <= (img[((-10 + ((R + 4) + (C + 4)) + (-1 + (C + 4)) + (-2)) + (1 + weight)) + (blury[_i0][((-32 + _T.i1) + _i1) + (-1 + _i2 - ((-256 + _T.i2) + 12)) + 0.25f)]) < threshold) ? _ct16 : _ct17);
                            mask.flip[(((11 - 2) + (3 + 2) + ((12 - 2) + 3)) + _i8)] = _ct18;
                        }
                    }
                }
            }
        }
    }
}

```



**15.5 ms / frame on 1 threads, 8.97 ms / frame on 4 threads**

# DOMAIN-SPECIFIC LANGUAGES (DSL)

- **The example motivates a domain-specific language + compiler approach**
- **High-performance domain-specific language + compiler:** productivity similar to ultra high-level or high-level but performance similar to manual or even better!



# DOMAIN-SPECIFIC LANGUAGES (DSL)

- **The example motivates a domain-specific language + compiler approach**
- **High-performance domain-specific language + compiler:** productivity similar to ultra high-level or high-level but performance similar to manual or even better!

# DOMAIN-SPECIFIC LANGUAGES (DSL)

## DSLs

- Exploit domain information to improve programmability, performance, and portability
- Expose greater information to the compiler and programmer specifies less
- abstract away many things from programmers (parallelism, memory)

## DSL compilers

- can “see” **across** routines – allow whole program optimization
- generate optimized code for multiple targets
- Programmers say **what** to execute and not **how** to execute

# DOMAIN-SPECIFIC LANGUAGES (DSL)

## DSLs

- Exploit domain information to improve programmability, performance, and portability
- Expose greater information to the compiler and programmer specifies less
- abstract away many things from programmers (parallelism, memory)

## DSL compilers

- can “see” **across** routines – allow whole program optimization
- generate optimized code for multiple targets
- Programmers say **what** to execute and not **how** to execute

# BIG PICTURE: ROLE OF COMPILERS

## General-Purpose

- Improve existing **general-purpose** compilers (for C, C++, Python, ...)
- Programmers say a **LOT**
- LLVM/Polly, GCC/Graphite

## Domain-Specific

- Build new **domain-specific languages and compilers**
- Programmers say **WHAT** they execute and not **HOW** they execute
- SPIRAL, Halide, Tensorflow, Pytorch, ...

# BIG PICTURE: ROLE OF COMPILERS

## General-Purpose

- Improve existing **general-purpose** compilers (for C, C++, Python, ...)
- Programmers say a **LOT**
- LLVM/Polly, GCC/Graphite
- Limited improvements, not everything is possible
- Broad impact

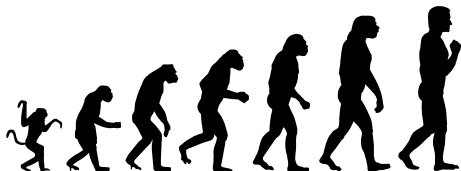
## Domain-Specific

- Build new **domain-specific languages and compilers**
- Programmers say **WHAT** they execute and not **HOW** they execute
- SPIRAL, Halide, Tensorflow, Pytorch, ...
- Dramatic speedups, Automatic parallelization
- Narrower impact and adoption

# BIG PICTURE: ROLE OF COMPILERS

## EVOLUTIONARY approach

- Improve existing **general-purpose** compilers (for C, C++, Python, ...)
- Programmers say a **LOT**
- LLVM/Polly, GCC/Graphite



## REVOLUTIONARY approach

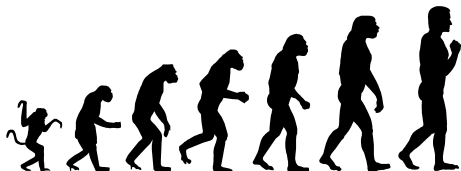
- Build new **domain-specific languages and compilers**
- Programmers say **WHAT** they execute and not **HOW** they execute
- SPIRAL, Halide, Tensorflow, Pytorch, ...



# BIG PICTURE: ROLE OF COMPILERS

## EVOLUTIONARY approach

- Improve existing **general-purpose** compilers (for C, C++, Python, ...)
- Programmers say a **LOT**
- LLVM/Polly, GCC/Graphite



- Both approaches share infrastructure
- Important to pursue both

## REVOLUTIONARY approach

- Build new **domain-specific languages and compilers**
- Programmers say **WHAT** they execute and not **HOW** they execute
- SPIRAL, Halide, Tensorflow, Pytorch, ...

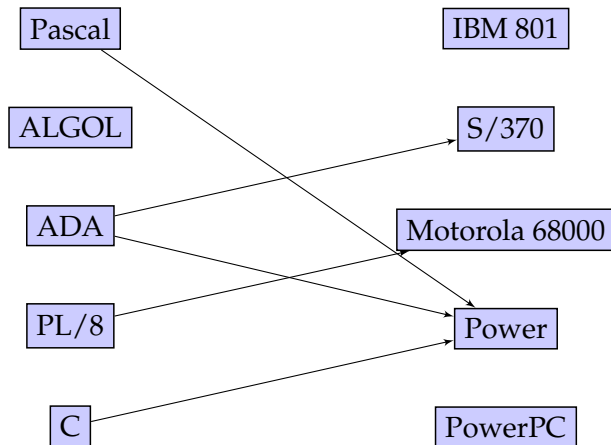


- Compilers for AI

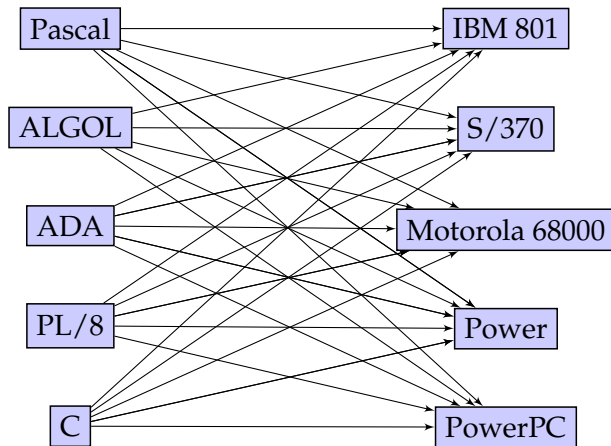


- **Compilers** are language translators: they translate programming languages to instructions hardware can execute
- One of the pillars of Computer Systems

# COMPILERS - THE EARLY DAYS

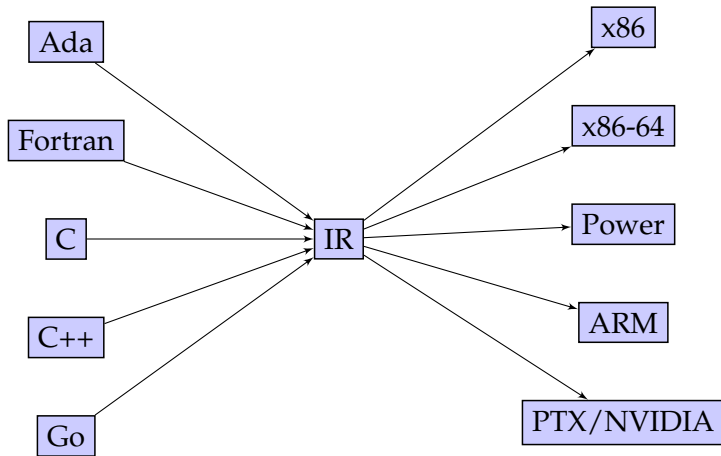


# COMPILERS - THE EARLY DAYS



- $M$  languages,  $N$  targets  $\Rightarrow M * N$  compilers! Not scalable!

# COMPILERS EVOLUTION - $M + N$



- With an common IR, we have  $M + N + 1$  compilers!

# WHAT DOES AN IR LOOK LIKE?

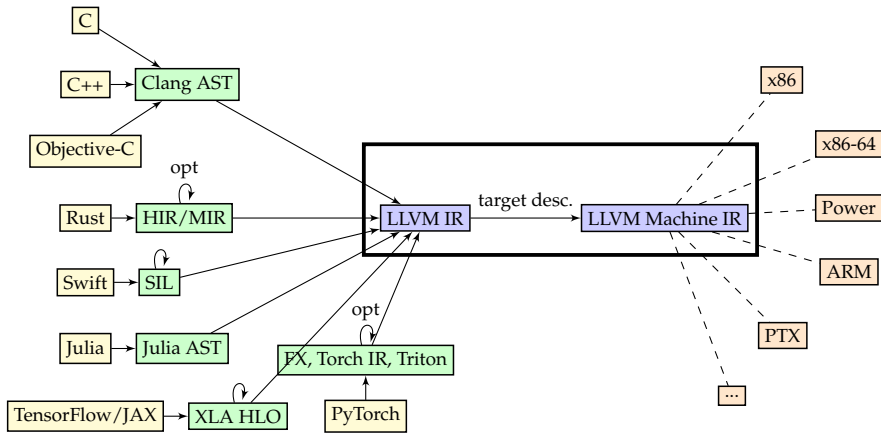
- A representation convenient to analyze and transform
- Round-trippable form that you can parse and print
- Low-level IRs are three-address code-like
- IRs have used expressions trees, 3-address code, graphs.
- Static Single Assignment: a property of IRs that makes it convenient; most IRs now use SSA

```
define void @foo(ptr nocapture %a) {
entry:
  br label %for.body

for.body:  ; preds = %for.body, %entry
  %indvars.iv = phi i64 [ 0, %entry ], [ %indvars.iv.next, %for.body ]
  %arrayidx = getelementptr inbounds i32, ptr %a, i64 %indvars.iv
  %0 = load i32, ptr %arrayidx, align 4
  %1 = add i32 %0, 2
  %inc = add nsw i32 %0, 1
  store i32 %inc, ptr %arrayidx, align 4
  %indvars.iv.next = add nuw nsw i64 %indvars.iv, 1
  %exitcond = icmp eq i64 %indvars.iv.next, 64
  br i1 %exitcond, label %for.end, label %for.body

for.end:      ; preds = %for.body
  ret void
}
```

# MODERN COMPILERS - LLVM IR-BASED



- LLVM: modular, reusable, open-source, but too low-level, not extensible for higher-order languages.

## AI programming frameworks



... ?

## Compiler infrastructure?

## Explosion of AI chips and accelerators



- Space in between is ruled by hand-written libraries. Not scalable.
- The right tools weren't available until recently.

# HOW DO YOU PROGRAM AI HARDWARE?

- High, mid, and low-level abstractions
  - High: PyTorch, JAX, ...
  - Mid: OpenAI Triton, cuTile
  - Low: CUDA, C/C++, CUTLASS, ...
- All three approaches need/use compilers in different ways
- They also share/rest on the same underlying infrastructure
  - Eg: Triton, MLIR, LLVM, PTX

```
# Triton JIT 'ed functions can be auto-tuned by using the 'triton.autotune' decorator, which
# - a list of 'triton.config' objects that define different configurations of
#   meta-parameters (e.g., BLOCK_SIZE, M) and compilation options (e.g., 'num_warps') to
#   an auto-tuning 'loop' whose shape in values will trigger evaluation of all the
#   provided configs
@triton.autotune(
    config_keys=['BLOCK_SIZE', 'M'],
    key=['M', 'N', 'K'],
)
@triton.jit
def torch_matmul(a_ptr, b_ptr, c_ptr,
                 M, N, K,
                 # The stride variables represent how much to increase the ptr by when moving by 1
                 # element in a particular dimension. E.g. 'stride_m' is how much to increase 'm'
                 # by to get the element one row down (i.e. has m rows).
                 stride_m, stride_n,
                 stride_a, stride_b,
                 stride_c,
                 # Meta-parameters
                 BLOCK_SIZE_M: tl.constexpr, BLOCK_SIZE_N: tl.constexpr,
                 GROUP_SIZE_M: tl.constexpr,
                 ACTIVATION: tl.constexpr,
):
    """Kernel for computing the matmul C = A x B.
    A has shape [M, N], B has shape [N, K] and C has shape [M, K]."""
    # For program id 0, just do the block of C to avoid compute.
    # This is done in a grouped ordering to promote L2 data reuse.
    # See above 'L2 Cache Optimizations' section for details.
    pid_m = tl.program_id(0)
    pid_n = tl.program_id(1)
    num_pid_n = GROUP_SIZE_M * num_pid_m
    group_id = pid_m // num_pid_n
    first_pid_m = group_id * GROUP_SIZE_M
    group_size_m = min(num_pid_m - first_pid_m, GROUP_SIZE_M)
    pid_m = first_pid_m + (pid_m % num_pid_n) * group_size_m
    pid_n = (pid_m // num_pid_n) // group_size_m

    # Create pointers for the first blocks of A and B.
    # We will advance this pointer as we move in the K direction
    # and accumulator
    # 'a_ptrs' is a block of [BLOCK_SIZE_M, BLOCK_SIZE_N] pointers
    # 'b_ptrs' is a block of [BLOCK_SIZE_M, BLOCK_SIZE_N] pointers
    # See above pointer arithmetic section for details.
    offs_m = (pid_m * BLOCK_SIZE_M + tl.arange(0, BLOCK_SIZE_M)) % M
    offs_n = (pid_n * BLOCK_SIZE_N + tl.arange(0, BLOCK_SIZE_N)) % N
    a_ptrs = tl.make_block_ptr(a_ptr, [M, N], [0, 0], [BLOCK_SIZE_M, BLOCK_SIZE_N],
                               [0, 0], [1, 0], [0, 0], [0, 0])
    b_ptrs = tl.make_block_ptr(b_ptr, [N, K], [0, 0], [BLOCK_SIZE_M, BLOCK_SIZE_N],
                               [0, 0], [0, 1], [0, 0], [0, 0])

    # Iterate to compute a block of the C matrix.
    # We accumulate into a [BLOCK_SIZE_M, BLOCK_SIZE_N] block
    # of float values for higher accuracy.
    # accumulator will be converted back to float after the loop.
    accumulator = tl.zeros([BLOCK_SIZE_M, BLOCK_SIZE_N], dtype=tl.float32)
    for k in range(0, K, BLOCK_SIZE_N):
        # Load the next block of A and B, generate a mask by checking the k dimension.
        # If it is out of bounds, set it to 0.
        a = tl.load(a_ptrs, mask=(offs_m + k * BLOCK_SIZE_M < M), other=0)
        b = tl.load(b_ptrs, mask=(offs_n + k * BLOCK_SIZE_N < N), other=0)
        # We accumulate along the k dimension.
        accumulator = tl.dot(a, b, accumulator)
        # Advance the ptrs in the next k block.
        a_ptrs += BLOCK_SIZE_M * stride_a
        b_ptrs += BLOCK_SIZE_M * stride_b

    # You can fuse arbitrary activation functions here
    # while the accumulator is still in FP32!
    if ACTIVATION == "leaky_relu":
        accumulator = tl.where(accumulator > 0, accumulator, 0)

    # Write back the block of the output matrix C with more.
    offs_c = (pid_m * BLOCK_SIZE_M + tl.arange(0, BLOCK_SIZE_M)) % M
    offs_n = (pid_n * BLOCK_SIZE_N + tl.arange(0, BLOCK_SIZE_N)) % N
    c_ptrs = tl.make_block_ptr(c_ptr, [M, K], [0, 0], [BLOCK_SIZE_M, BLOCK_SIZE_N],
                               [0, 0], [1, 0], [0, 0], [0, 0])
    c_mask = (offs_c + k * BLOCK_SIZE_M < M) & (offs_n + k * BLOCK_SIZE_N < N)
    tl.store(c_ptrs, accumulator, c_mask)
```

```
# Triton JIT 'ed functions can be auto-tuned by using the 'triton.autotune' decorator, which
# - a list of 'triton.config' objects that define different configurations of
#   meta-parameters (e.g., BLOCK_SIZE, M) and compilation options (e.g., 'num_warps') to
#   an auto-tuning 'loop' whose shape in values will trigger evaluation of all the
#   provided configs
@triton.autotune(
    config_keys=['BLOCK_SIZE', 'M'],
    key=['M', 'N', 'K'],
)
@triton.jit
def torch_matmul(a_ptr, b_ptr, c_ptr,
                 M, N, K,
                 # The stride variables represent how much to increase the ptr by when moving by 1
                 # element in a particular dimension. E.g. 'stride_m' is how much to increase 'm'
                 # by to get the element one row down (i.e. has m rows).
                 stride_m, stride_n,
                 stride_a, stride_b,
                 stride_c,
                 # Meta-parameters
                 BLOCK_SIZE_M: tl.constexpr, BLOCK_SIZE_N: tl.constexpr,
                 GROUP_SIZE_M: tl.constexpr,
                 ACTIVATION: tl.constexpr,
):
    """Kernel for computing the matmul C = A x B.
    A has shape [M, N], B has shape [N, K] and C has shape [M, K]."""
    # For program id 0, just do the block of C to avoid compute.
    # This is done in a grouped ordering to promote L2 data reuse.
    # See above 'L2 Cache Optimizations' section for details.
    pid_m = tl.program_id(0)
    pid_n = tl.program_id(1)
    num_pid_n = GROUP_SIZE_M * num_pid_m
    group_id = pid_m // num_pid_n
    first_pid_m = group_id * GROUP_SIZE_M
    group_size_m = min(num_pid_m - first_pid_m, GROUP_SIZE_M)
    pid_m = first_pid_m + (pid_m % num_pid_n) * group_size_m
    pid_n = (pid_m // num_pid_n) // group_size_m

    # Create pointers for the first blocks of A and B.
    # We will advance this pointer as we move in the K direction
    # and accumulator
    # 'a_ptrs' is a block of [BLOCK_SIZE_M, BLOCK_SIZE_N] pointers
    # 'b_ptrs' is a block of [BLOCK_SIZE_M, BLOCK_SIZE_N] pointers
    # See above pointer arithmetic section for details.
    offs_m = (pid_m * BLOCK_SIZE_M + tl.arange(0, BLOCK_SIZE_M)) % M
    offs_n = (pid_n * BLOCK_SIZE_N + tl.arange(0, BLOCK_SIZE_N)) % N
    a_ptrs = tl.make_block_ptr(a_ptr, [M, N], [0, 0], [BLOCK_SIZE_M, BLOCK_SIZE_N],
                               [0, 0], [1, 0], [0, 0], [0, 0])
    b_ptrs = tl.make_block_ptr(b_ptr, [N, K], [0, 0], [BLOCK_SIZE_M, BLOCK_SIZE_N],
                               [0, 0], [0, 1], [0, 0], [0, 0])

    # Iterate to compute a block of the C matrix.
    # We accumulate into a [BLOCK_SIZE_M, BLOCK_SIZE_N] block
    # of float values for higher accuracy.
    # accumulator will be converted back to float after the loop.
    accumulator = tl.zeros([BLOCK_SIZE_M, BLOCK_SIZE_N], dtype=tl.float32)
    for k in range(0, K, BLOCK_SIZE_N):
        # Load the next block of A and B, generate a mask by checking the k dimension.
        # If it is out of bounds, set it to 0.
        a = tl.load(a_ptrs, mask=(offs_m + k * BLOCK_SIZE_M < M), other=0)
        b = tl.load(b_ptrs, mask=(offs_n + k * BLOCK_SIZE_N < N), other=0)
        # We accumulate along the k dimension.
        accumulator = tl.dot(a, b, accumulator)
        # Advance the ptrs in the next k block.
        a_ptrs += BLOCK_SIZE_M * stride_a
        b_ptrs += BLOCK_SIZE_M * stride_b

    # You can fuse arbitrary activation functions here
    # while the accumulator is still in FP32!
    if ACTIVATION == "leaky_relu":
        accumulator = tl.where(accumulator > 0, accumulator, 0)

    # Write back the block of the output matrix C with more.
    offs_c = (pid_m * BLOCK_SIZE_M + tl.arange(0, BLOCK_SIZE_M)) % M
    offs_n = (pid_n * BLOCK_SIZE_N + tl.arange(0, BLOCK_SIZE_N)) % N
    c_ptrs = tl.make_block_ptr(c_ptr, [M, K], [0, 0], [BLOCK_SIZE_M, BLOCK_SIZE_N],
                               [0, 0], [1, 0], [0, 0], [0, 0])
    c_mask = (offs_c + k * BLOCK_SIZE_M < M) & (offs_n + k * BLOCK_SIZE_N < N)
    tl.store(c_ptrs, accumulator, c_mask)
```



# HOW DO YOU PROGRAM AI HARDWARE?

[illegible][illegible]

- High, mid, and low-level abstractions
- High: PyTorch, JAX, ...
- Mid: OpenAI Triton, cuTile
- Low: CUDA, C/C++, CUTLASS, ...
- All three approaches need/use compilers in different ways
- They also share/rest on the same underlying infrastructure
- Eg: Triton, MLIR, LLVM, PTX

# PYTHON-BASED PROGRAMMING FRAMEWORKS: TODAY

- Where does performance in Python-based frameworks come from?
- Largely from libraries written in C, C++, CUDA, and even assembly
- Compilers exist: XLA, TorchInductor, TensorRT
  - Limited in many ways: “semi-compilers”, fragmented infra, performance
  - Still evolving

```
class SelfAttentionLayer(nn.Module):
    def __init__(self, feature_size):
        super(SelfAttentionLayer, self).__init__()
        self.feature_size = feature_size

        # Linear transformations for Q, K, V from the same source.
        self.key = nn.Linear(feature_size, feature_size)
        self.query = nn.Linear(feature_size, feature_size)
        self.value = nn.Linear(feature_size, feature_size)

    def forward(self, x, mask=None):
        # Apply linear transformations.
        keys = self.key(x)
        queries = self.query(x)
        values = self.value(x)

        # Scaled dot-product attention.
        scores = torch.matmul(queries, keys.transpose(-2, -1))
            / torch.sqrt(torch.tensor(self.feature_size, dtype=torch.float32))

        # Apply mask (if provided).
        if mask is not None:
            scores = scores.masked_fill(mask == 0, -1e9)

        # Apply softmax.
        attention_weights = F.softmax(scores, dim=-1)

        # Multiply weights with values.
        output = torch.matmul(attention_weights, values)

        return output, attention_weights
```

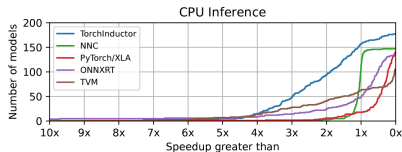
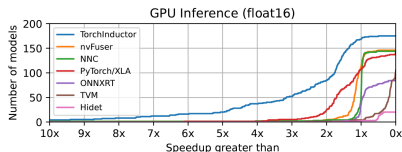
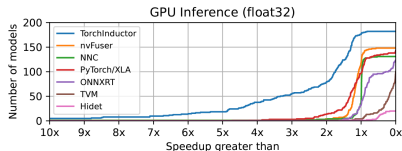
# OPENING MULTI-LEVEL DOORS TO PROGRAMMING AI HARDWARE

- ❶ High-level Python-based programming frameworks (e.g. PyTorch, JAX),
- ❷ Compiler support for (1) that could be turned off/on (e.g. `torch.compile`),
- ❸ Mid/low-level programming support (e.g., CUDA, CUTLASS, CuTile, Triton)
- ❹ Low-level MLIR dialects that expose their hardware intrinsics/virtual ISA on top of which both (2) compilers and (3) low-level frameworks rest,
- ❺ Ability to use inline virtual ISA.



# COMPILER AUTO-PARALLELIZATION IS ALREADY HERE!

- Recent PyTorch 2 ASPLOS paper  
*PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation*, Ansel et al. (Meta), ASPLOS 2024.



# HOW IS HARDWARE EVOLVING? (1/2)

- Multiple cores (early 2000s)
- Wider SIMD (early 2000s)
- Many cores (late 2000s)
- Heterogeneity (2000s/2010s)
- Tensor/matmul cores (mid 2010s)
- Low-precision compute instructions (late 2010s/2020s)

## From 2000s to now



## HOW IS HARDWARE EVOLVING? (2/2)

Example: NVIDIA H100 chip

- 80 GB of GPU DRAM
- 3.35 TB/s of memory bandwidth (HBM3).
- 990 TFLOPS for fp16 tensor operations, 1.98 PFLOPS for int8.
- 50 MB of L2 cache.

# HOW ARE PROGRAMMING FRAMEWORKS EVOLVING?

## AI programming frameworks

- Programmer productivity
- Write less and do more
- Hardware usability
- Deliver performance
- Deliver portability



## AI programming frameworks



... ?

## Compiler infrastructure?

## Explosion of AI chips and accelerators



- MLIR infrastructure: open-sourced by Google in 2019



## AI programming frameworks



... ?



## Explosion of AI chips and accelerators



- MLIR infrastructure: open-sourced by Google in 2019

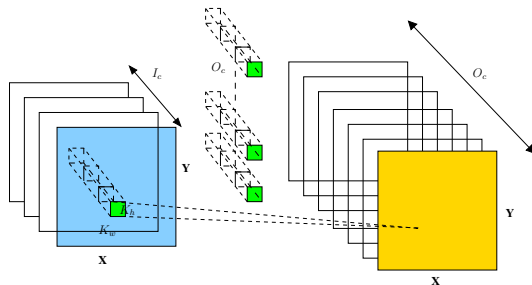


## MLIR

- ML in MLIR: Multi-level
- Characteristics
  - Loops and multi-dimensional arrays (tensors) had to be first class citizens
  - Had to be extensible (types, operations, attributes)
  - Had to enable building both general-purpose and domain-specific compilers and even more.
  - Had to be open-source with a permissive license

# MULTI-DIMENSIONALITY EVERYWHERE: CNN CONVOLUTION

```
for (n = 0; n < N; n++) // Samples in a batch.  
  for (o = 0; o < Oc; o++) // Output feature channels.  
    for (i = 0; i < Ic; i++) // Input feature channels.  
      for (y = 0; y < Y; y++) // Layer height.  
        for (x = 0; x < X; x++) // Layer width.  
          for (kh = 0; kh < Kh; kh++) // Convolution kernel height.  
            for (kw = 0; kw < Kw; kw++) // Convolution kernel width.  
              output[n, o, y, x] += input[n, i, y+kh, x+kw] * weights[o, i, kh, kw];
```

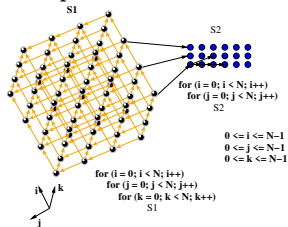


# MLIR: MULTI-LEVEL INTERMEDIATE REPRESENTATION

## 1. Ops (general purpose to domain specific) on tensor types / graph form

```
%patches = "tf.reshape"(%patches, %minus_one, %minor_dim_size)
           : (tensor<? x ? x ? x f32>, index, index) -> tensor<? x ? x f32>
%mat_out = "tf.matmul"(%patches_flat, %patches_flat)(transpose_a : true)
           : (tensor<? x ? x f32>, tensor<? x ? x f32>) -> tensor<? x ? x f32>
%vec_out = "tf.reduce_sum"(%patches_flat) [axis: 0] : (tensor<? x ? x f32>) -> tensor<? x f32>
```

## 2. Loop-level / mid-level form



```
affine.for %i = 0 to 8 step 4 {
  affine.for %j = 0 to 8 step 4 {
    affine.for %k = 0 to 8 step 4 {
      affine.for %ii = #map0(%i) to #map1(%i) {
        affine.for %jj = #map0(%j) to #map1(%j) {
          affine.for %kk = #map0(%k) to #map1(%k) {
            %5 = affine.load %arg0[%ii, %kk] : memref<8 x 8 x vector<64 x f32>>
            %6 = affine.load %arg1[%kk, %jj] : memref<8 x 8 x vector<64 x f32>>
            %7 = affine.load %arg2[%ii, %jj] : memref<8 x 8 x vector<64 x f32>>
            %8 = arith.mulf %5, %6 : vector<64xf32>
            %9 = arith.addf %7, %8 : vector<64xf32>
            affine.store %9, %arg2[%ii, %jj] : memref<8 x 8 x vector<64xf32>>
          }
        }
      }
    }
  }
}
```

## 3. Low-level form: closer to hardware

```
%v1 = memref.load %a[%i2, %i3] : memref<256 x 64 x vector<16 x f32>>
%v2 = memref.load %b[%i2, %i3] : memref<256 x 64 x vector<16 x f32>>
%v3 = addf %v1, %v2 : vector<16 x f32>
memref.store %v3, %d[%i2, %i3] : memref<256 x 64 x vector<16 x f32>>
```

# MODERN COMPILER TOPICS IN THIS COURSE

- ① Foundations: SSA, Dominance, Basic concepts for control flow analysis and data flow analysis
- ② Compiler optimizations for parallelism and locality
- ③ Affine abstraction/Polyhedral framework (only the basics)
- ④ MLIR
- ⑤ Practice: Building compilers using MLIR
- ⑥ Practice: Building compilers and optimizers for AI frameworks (basic overview, pointers)