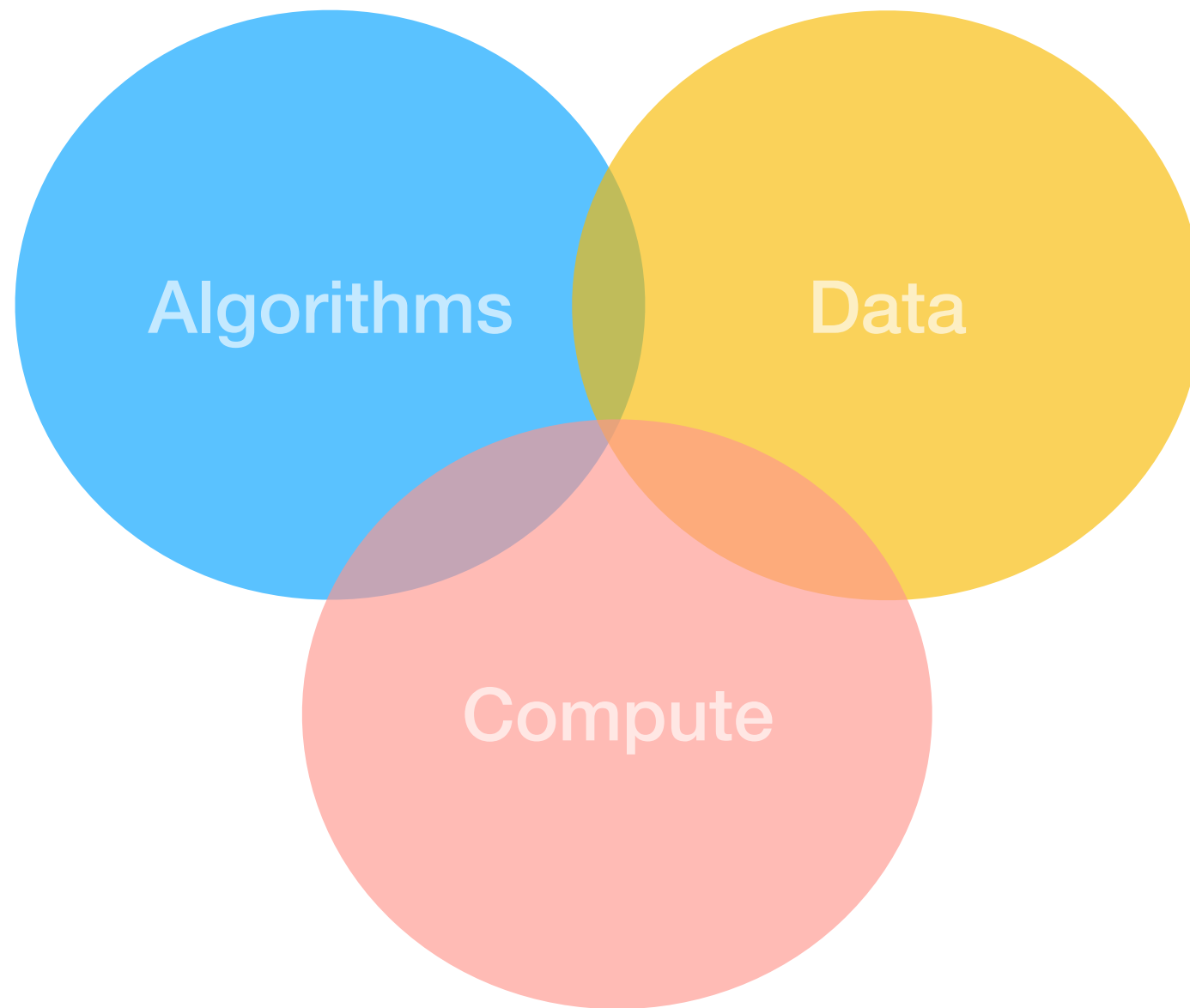


Compiling Dense and Sparse Neural Networks using Tiramisu

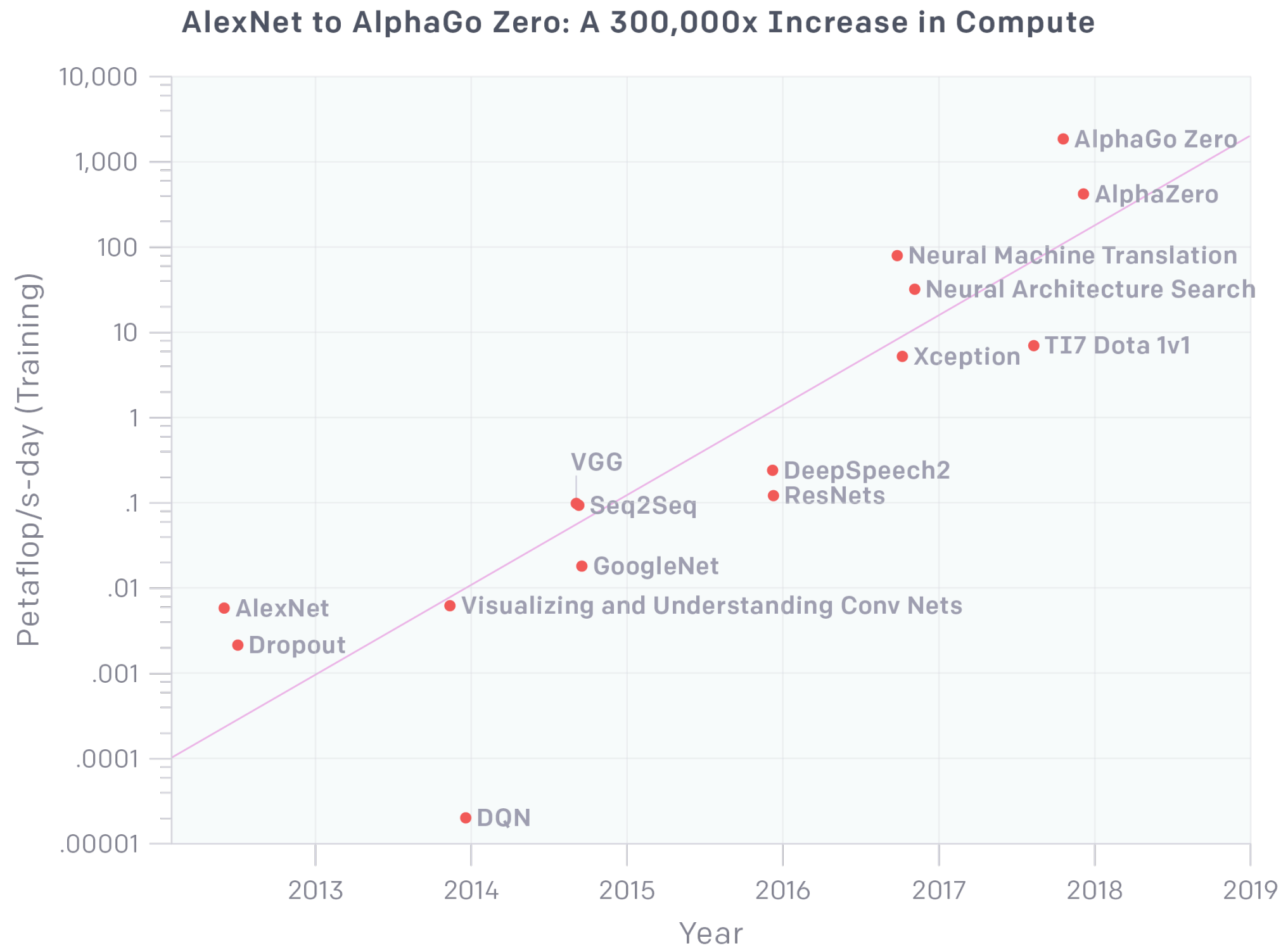
Riyadh Baghdadi, J. Ray, E. Del Sozzo, Y. Zhang, P. Suriana, S. Kamil, M. Ben Romdhane, A. Akkas, A. Renda, J. Elliott Frankle, K. Abdous, A. N. Debbagh, F. Z. Benhamida, T. Arbaoui, B. Karima, Michael Carbin, **Saman Amarasinghe**



Deep Learning Advancements

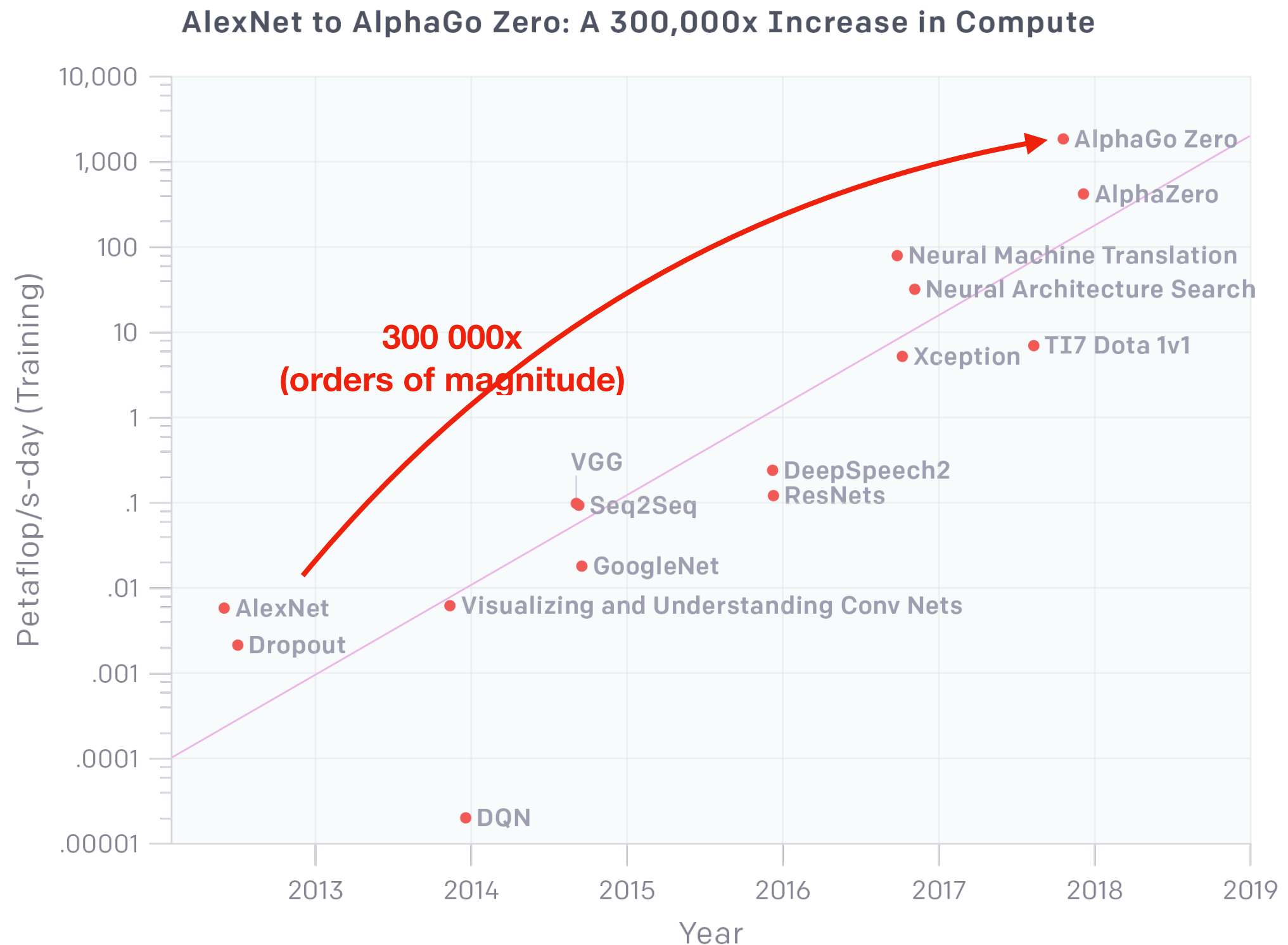


Increasing Need for Computing Power



[source: <https://openai.com/blog/ai-and-compute/>]

Increasing Need for Computing Power



[source: <https://openai.com/blog/ai-and-compute/>]

Diverse Hardware Architectures



Diverse Hardware Architectures

AI Chip Landscape

V0.5 August, 2019

S.T.

Tech Giants/System



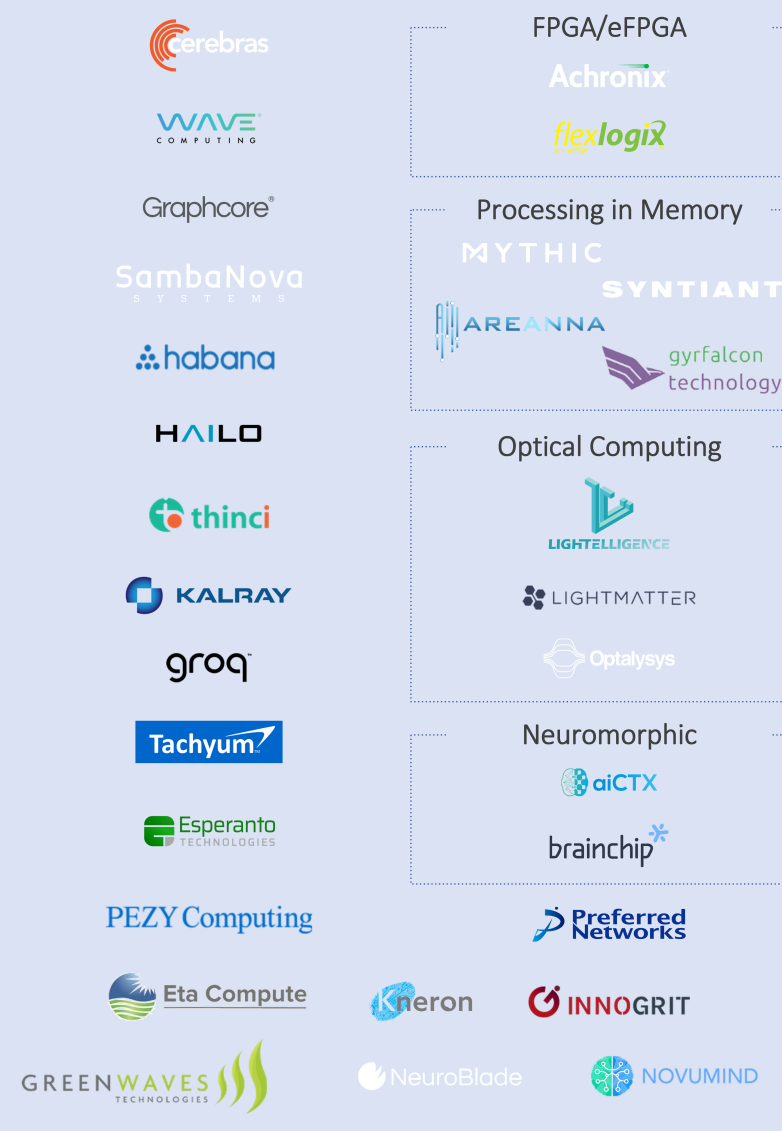
IC Vender/Fabless



Startup in China



Startup Worldwide



More at <https://basicmi.github.io/AI-Chip/>

IP/Design Service



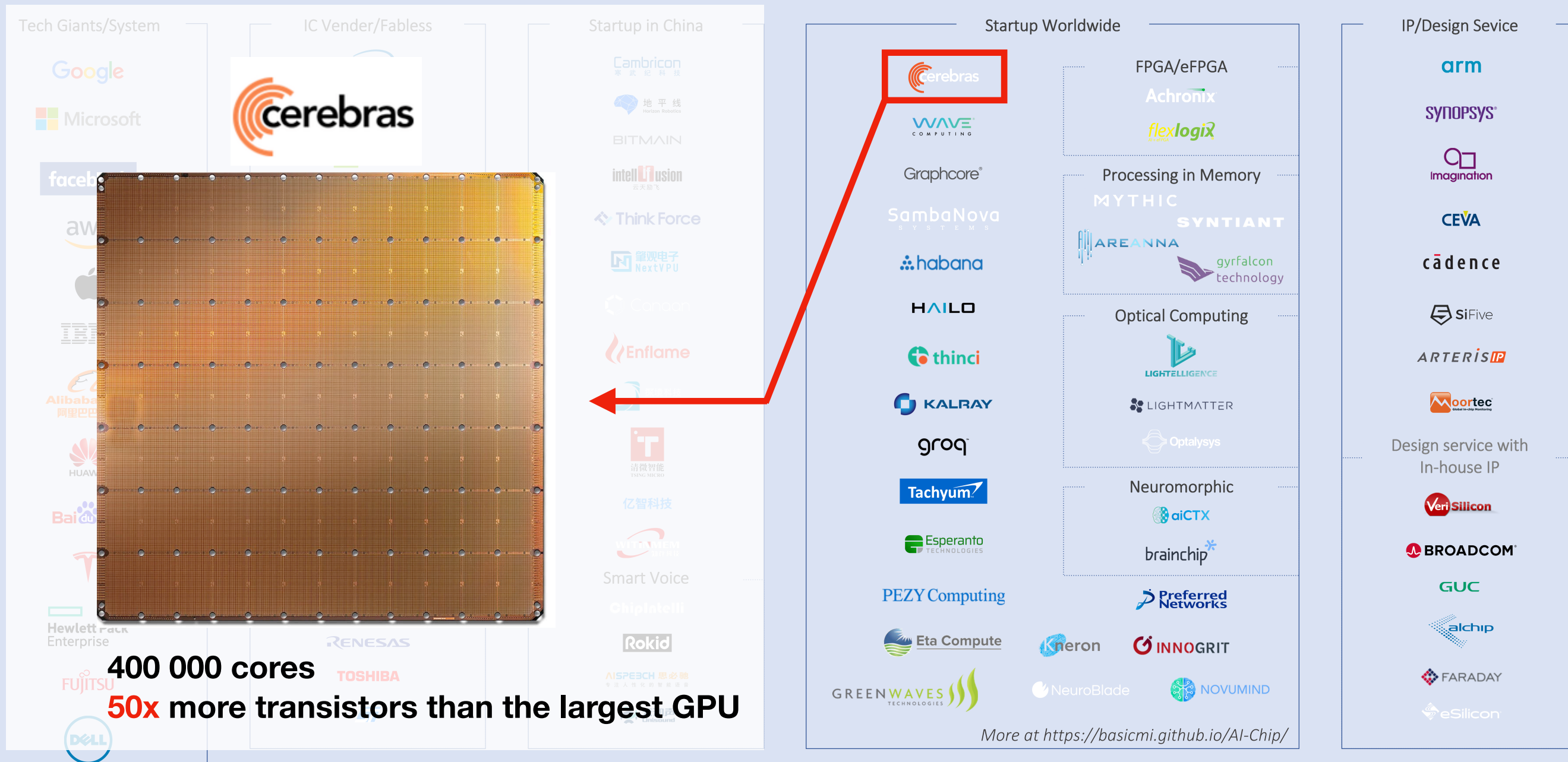
[source: <http://github.com/basicmi/AI-Chip>]

Diverse Hardware Architectures

AI Chip Landscape

V0.5 August, 2019

S.T.



[source: <http://github.com/basicmi/AI-Chip>]

We need compilers to
target these hardware architectures

State-of-the-art DNN Compilers

Support

- Feed forward NN
- Multiple hardware architectures

State-of-the-art DNN Compilers

Support

- Feed forward NN
- Multiple hardware architectures

Limitations in supporting

- Sparse NN
- Optimizing RNNs
- Distributed architectures

State-of-the-art DNN Compilers

Support

- Feed forward NN
- Multiple hardware architectures

Limitations in supporting

- Sparse NN
- Optimizing RNNs
- Distributed architectures

} **Tiramisu**

State-of-the-art DNN Compilers

Support

- Feed forward NN
- Multiple hardware architectures

5x
over industrial
libraries!

Limitations in supporting

- Sparse NN
- Optimizing RNNs
- Distributed architectures

Tiramisu

State-of-the-art DNN Compilers

Support

- Feed forward NN
- Multiple hardware architectures

5x
over industrial
libraries!

Limitations in supporting

- Sparse NN
- Optimizing RNNs
- Distributed architectures

Tiramisu

2x
over state-of-the-art
compilers!

Tiramisu

Tiramisu

Phase I:

Tiramisu

Phase I:

Manual Optimization
(Commands)

Tiramisu

Phase I:

Manual Optimization
(Commands)

Algorithm

Tiramisu

Phase I:

Manual Optimization
(Commands)

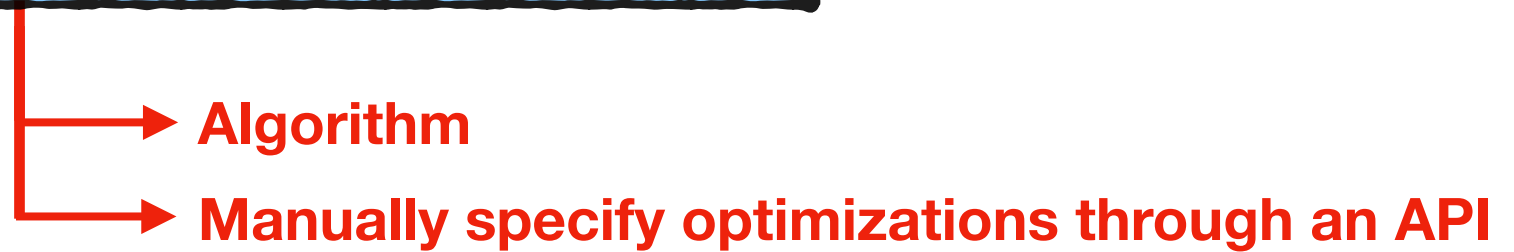
→ Algorithm

→ Manually specify optimizations through an API

Tiramisu

Phase I:

Manual Optimization
(Commands)



Phase II:

Tiramisu

Phase I:

Manual Optimization
(Commands)

- Algorithm
- Manually specify optimizations through an API

Phase II:

Automatic
Optimization

Tiramisu

Phase I:

Manual Optimization
(Commands)

Polyhedral Model
(Sparse)

- Algorithm
- Manually specify optimizations through an API

Phase II:

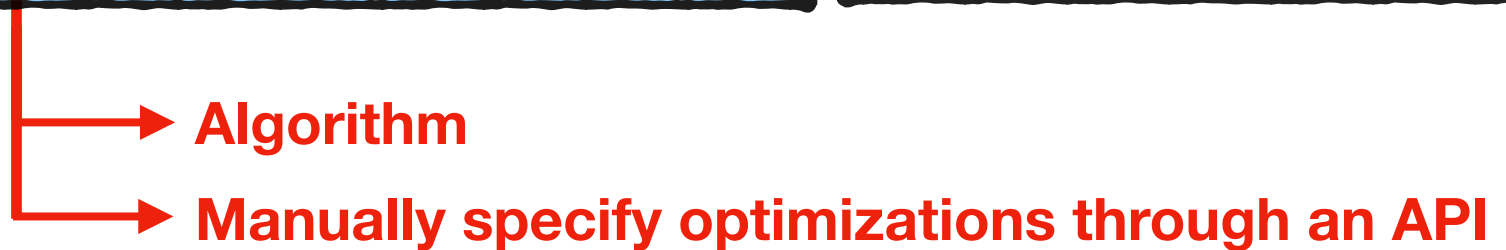
Automatic
Optimization

Tiramisu

Phase I:

Manual Optimization
(Commands)

Polyhedral Model
(Sparse)



Phase II:

Automatic
Optimization

Deep Learning
Cost Model

Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Outline

- **Phase I**

- Tiramisu Overview

- Example

- RNNs

- Sparse DNNs

- **Phase II**

- Automatic Optimization

Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Tiramisu

Tiramisu

- A polyhedral compiler

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```


Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```

- **Optimizations Commands**

Tiramisu

- A polyhedral compiler
- **C++ API** for expressing algorithms and optimizations
- Designed for optimizing **loop nests manipulating arrays**
- Separates **Algorithm** from **Optimization Commands**
- **Algorithm**
 - Declare pure computations (unoptimized)

Pseudocode

```
for x in 0 .. N
  for y in 0 .. N
    out[x, y] = 0;
```

Tiramisu

```
var x(0, N), y(0, N);
computation out(x, y);
out(x, y) = 0;
```

- **Optimizations Commands**

```
out.parallelize(x);
out.vectorize(y, 4);
```

Advantages of the Polyhedral Model

Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)

Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

```
//Halide github issue #2373
for x in 0 ... N
  for y in 0 ... N
    for r in 0 ... 32
      if (x >= r)
        out(x, y, r) = in(x - r, y);
```

Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

```
//Halide github issue #2373
for x in 0 ... N
  for y in 0 ... N
    for r in 0 ... 32
      if (x >= r)
        out(x, y, r) = in(x - r, y);
```


Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

```
//Halide github issue #2373
for x in 0 ... N
  for y in 0 ... N
    for r in 0 ... 32
      if (x >= r)
        out(x, y, r) = in(x - r, y);
```

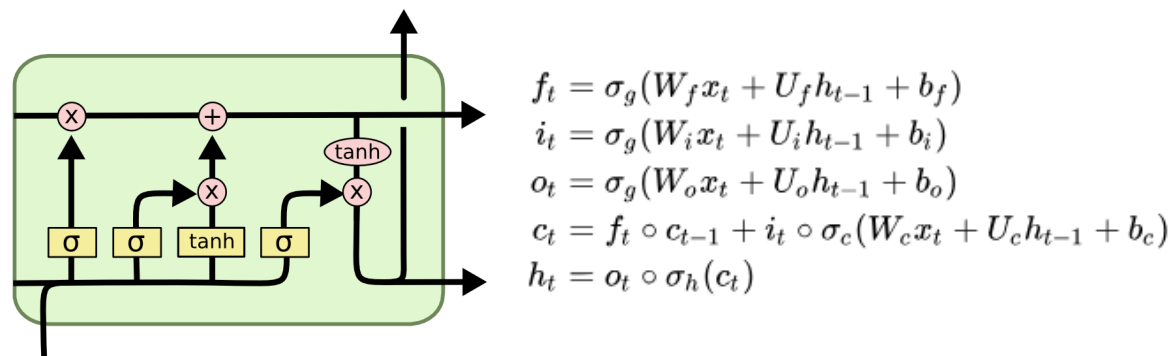
- Express algorithms with cyclic data flow graphs

Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

```
//Halide github issue #2373
for x in 0 ... N
  for y in 0 ... N
    for r in 0 ... 32
      if (x >= r)
        out(x, y, r) = in(x - r, y);
```

- Express algorithms with cyclic data flow graphs

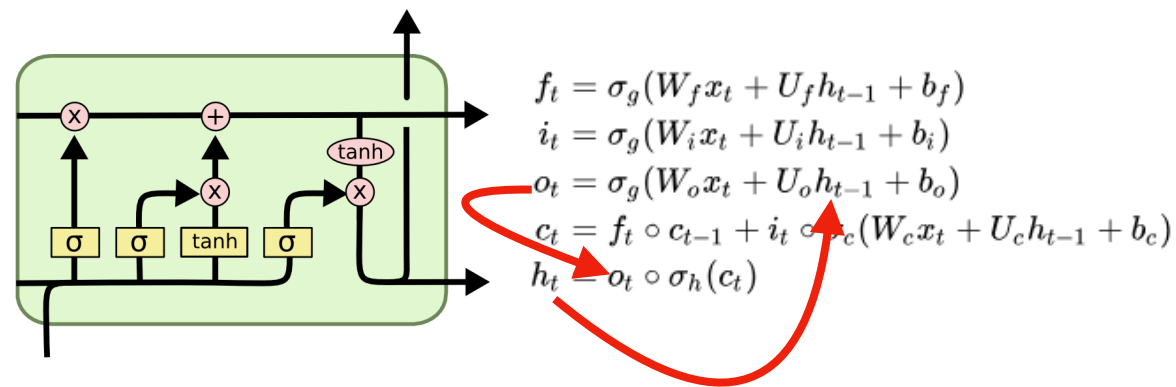


Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

```
//Halide github issue #2373
for x in 0 ... N
  for y in 0 ... N
    for r in 0 ... 32
      if (x >= r)
        out(x, y, r) = in(x - r, y);
```

- Express algorithms with cyclic data flow graphs

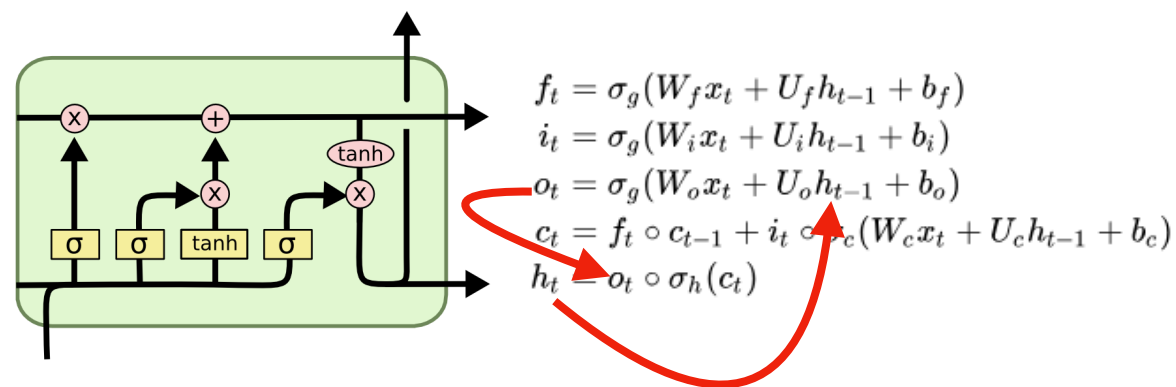


Advantages of the Polyhedral Model

- Complex transformations (all affine transformations)
- Naturally support non-rectangular loops

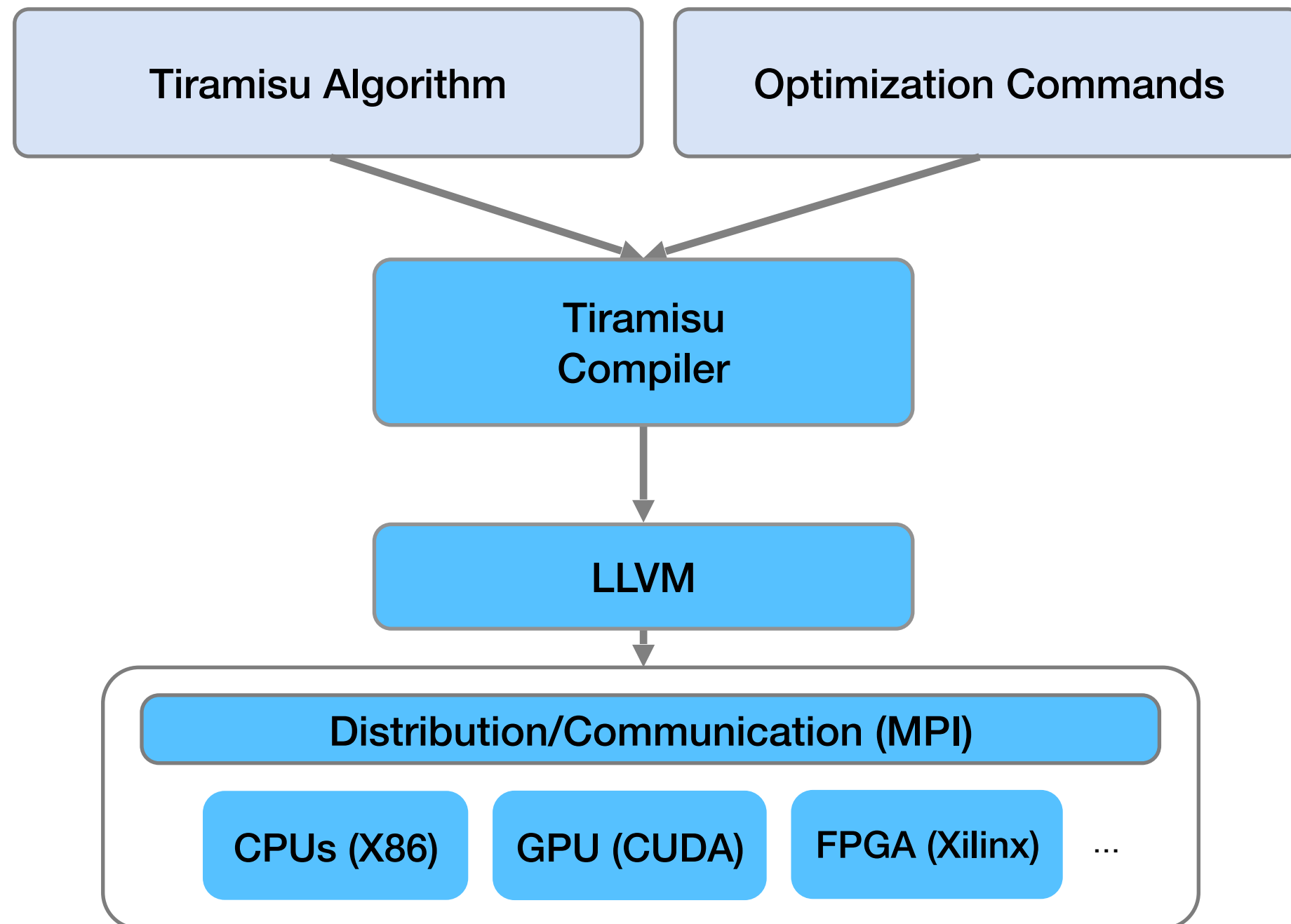
```
//Halide github issue #2373
for x in 0 ... N
  for y in 0 ... N
    for r in 0 ... 32
      if (x >= r)
        out(x, y, r) = in(x - r, y);
```

- Express algorithms with cyclic data flow graphs

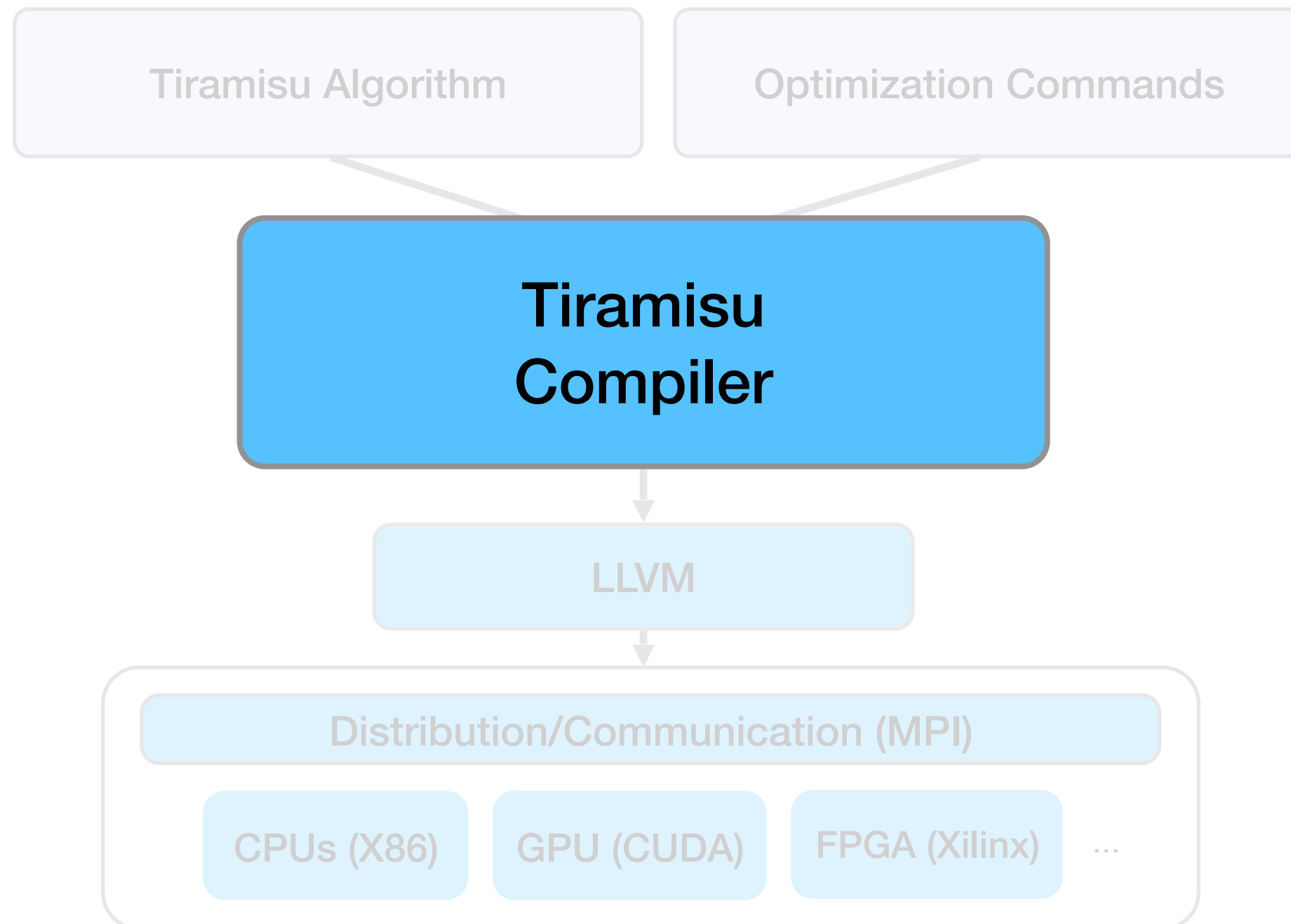


- Adequate for distributed code generation

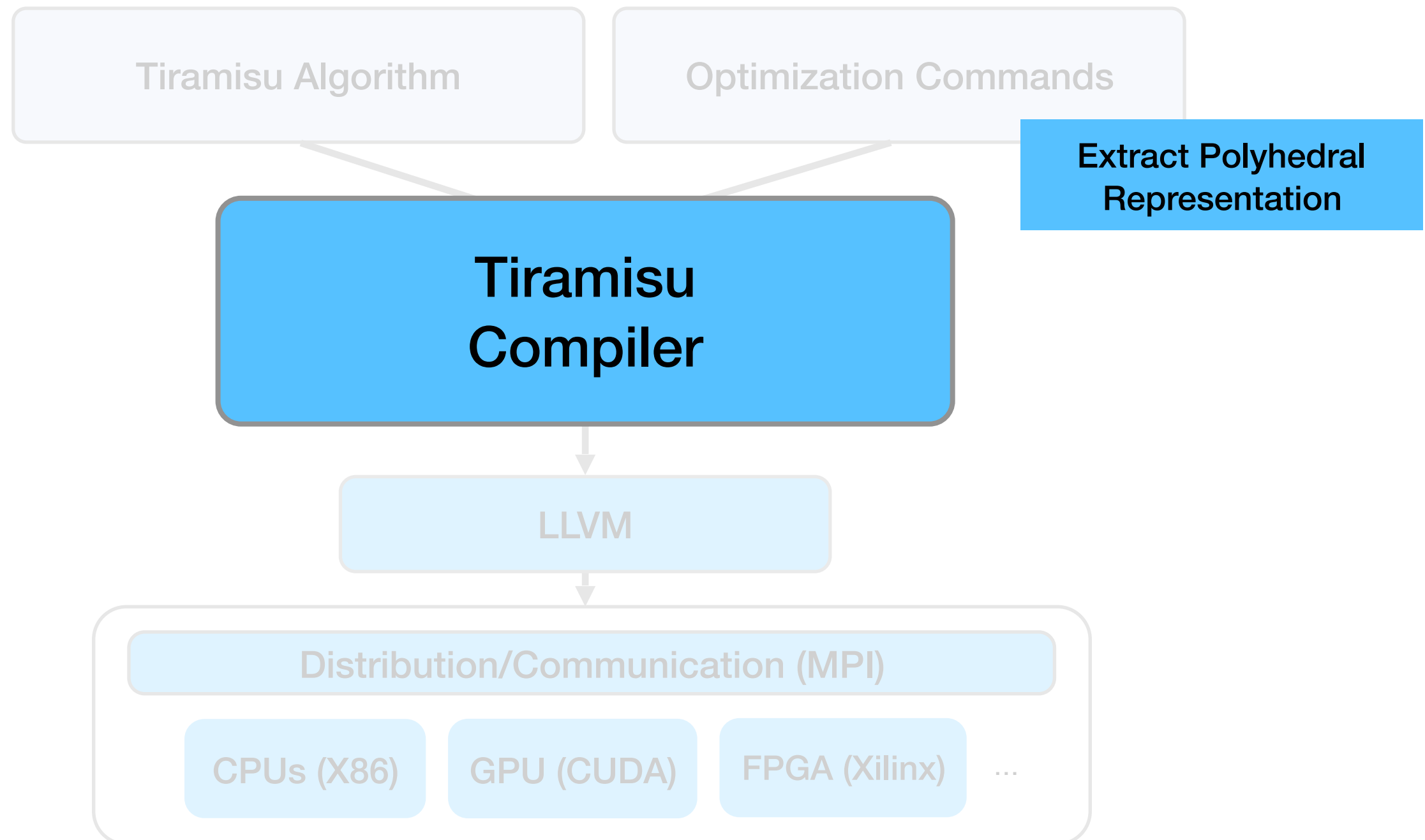
Compilation Flow



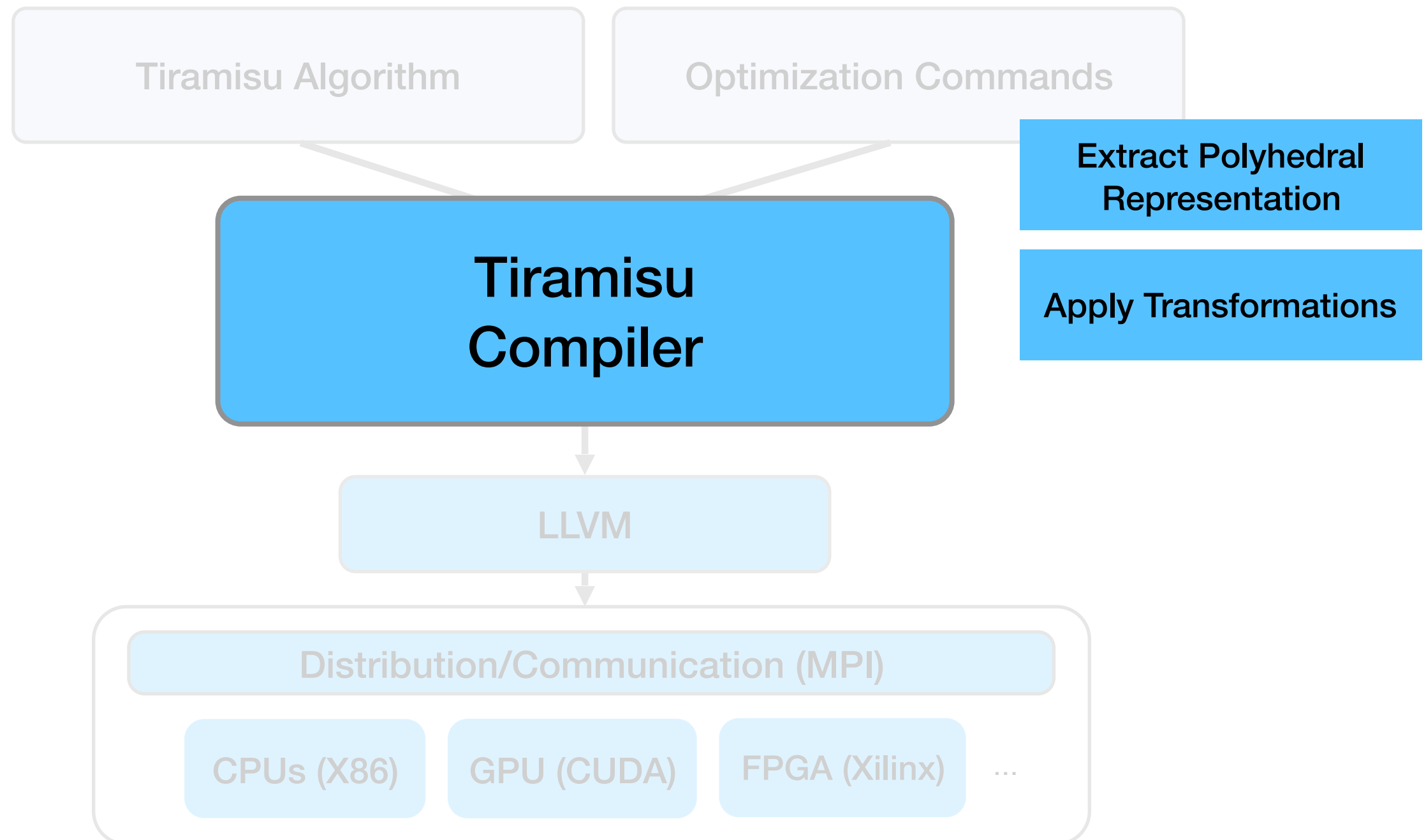
Compilation Flow



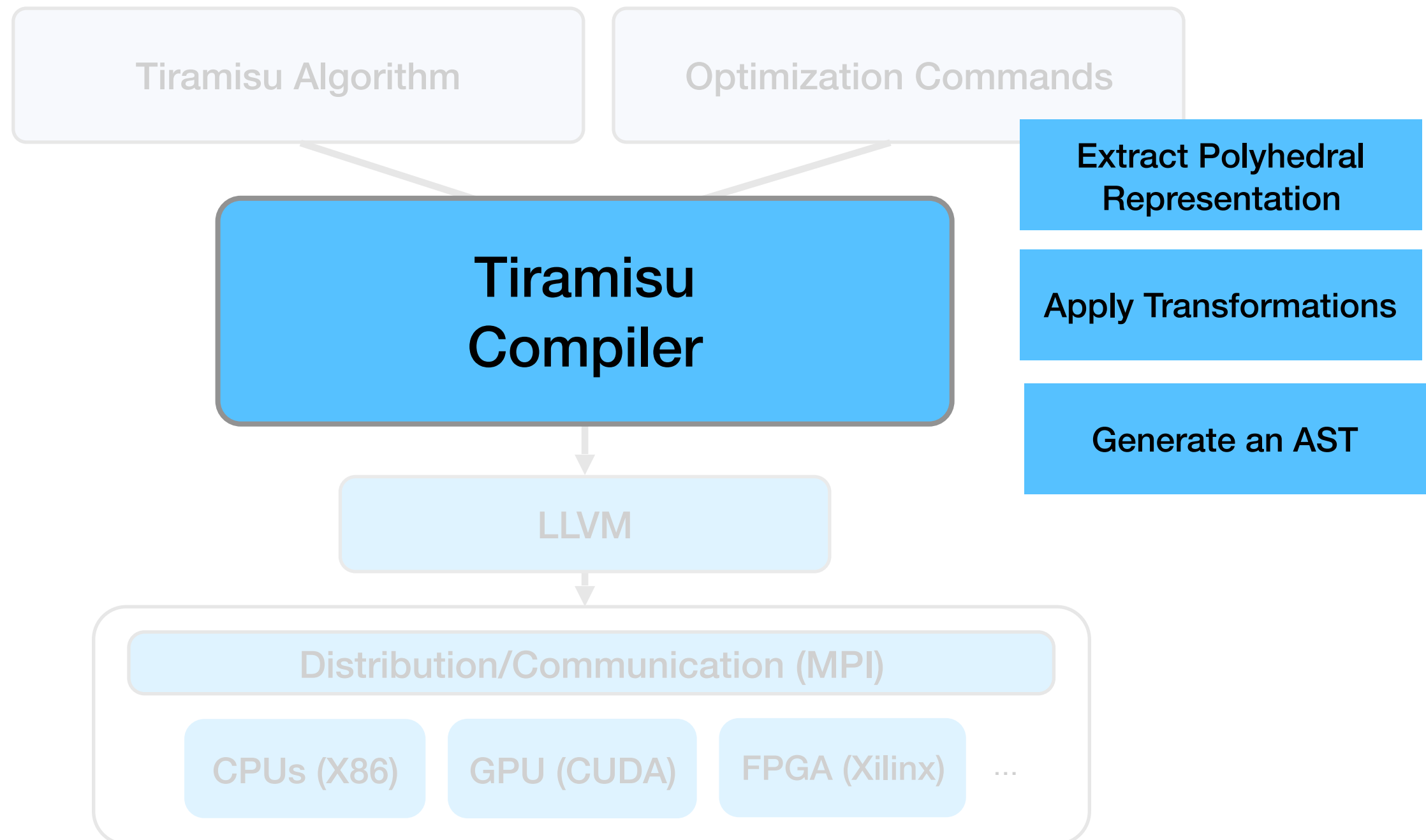
Compilation Flow



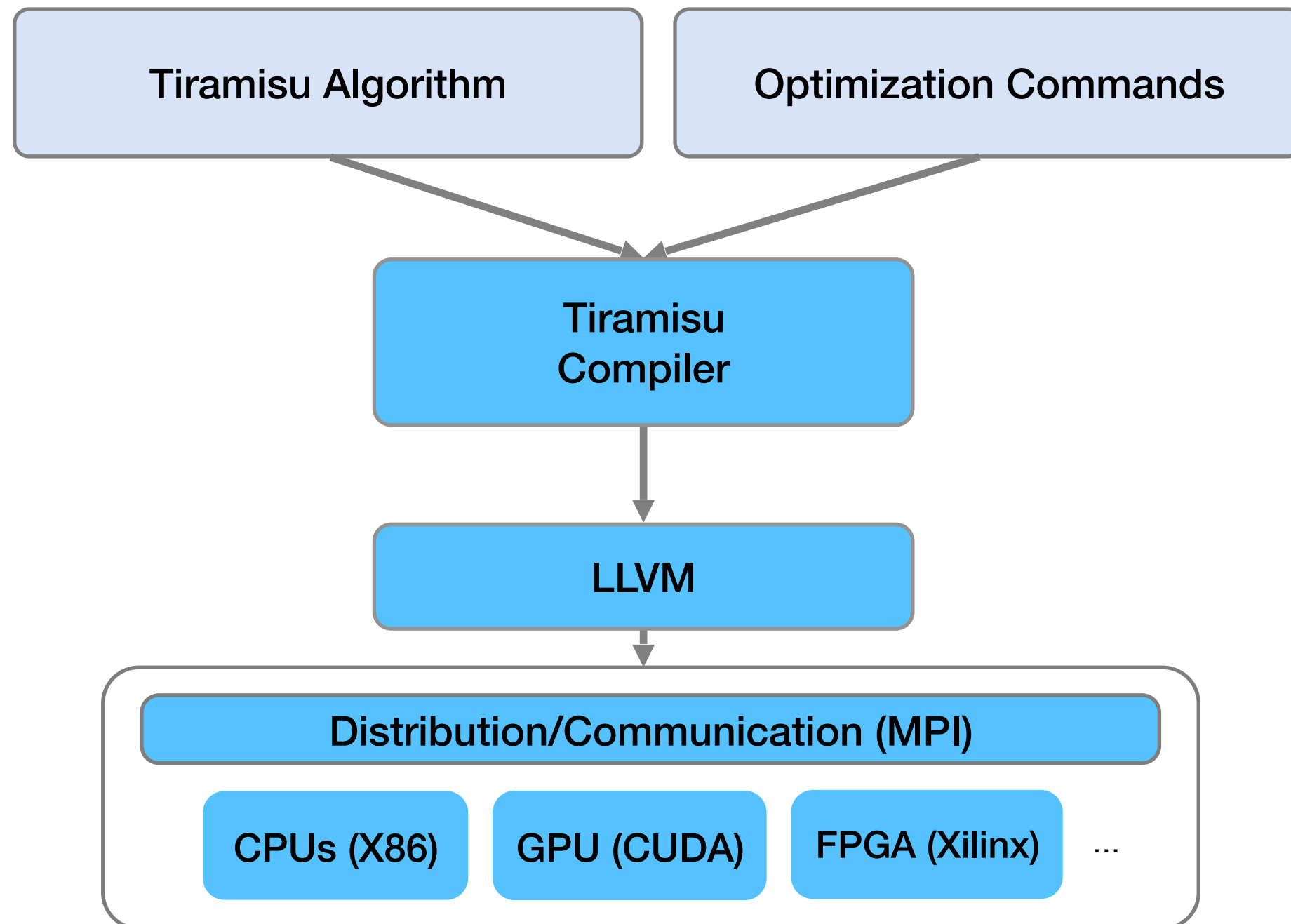
Compilation Flow



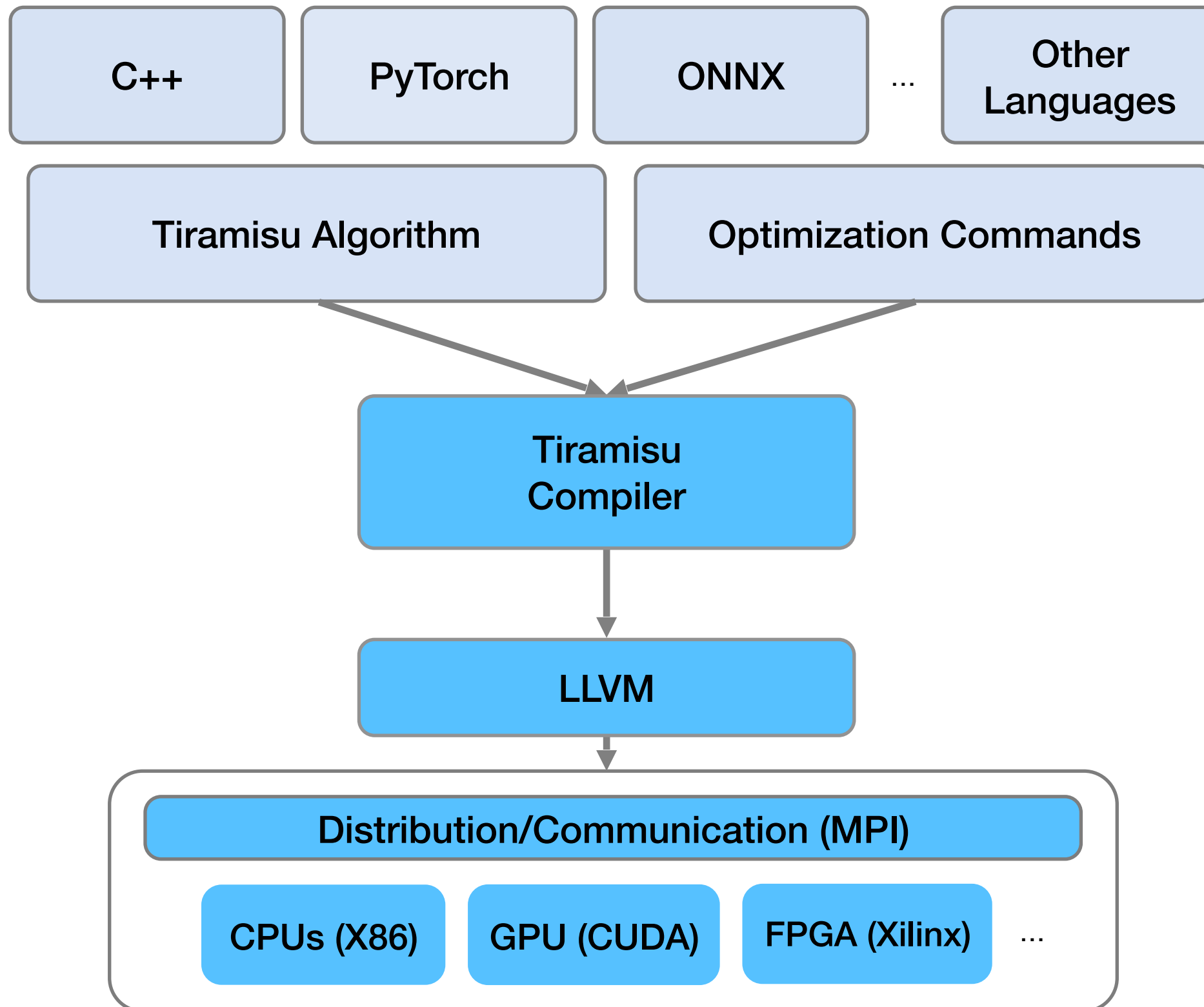
Compilation Flow



Compilation Flow



Compilation Flow



Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Scheduling

Outline

- **Phase I**

- Tiramisu Overview

- Example

- RNNs

- Sparse DNNs

- **Phase II**

- Automatic Scheduling

Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Scheduling

Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);
```

Matrix Multiplication

```
// Algorithm
```

```
var i(0, N), j(0, N), k(0, N);
```

```
input A(i,j), B(i,j);
```

```
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```


Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);
```

Matrix Multiplication

```
// Algorithm
```

```
var i(0, N), j(0, N), k(0, N);
```

```
input A(i,j), B(i,j);
```

```
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);
```

Matrix Multiplication

```
// Algorithm  
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);  
  
C(i,j) = 0;  
C(i,j) += A(i, k) * B(k, j);
```

Matrix Multiplication

```
// Algorithm  
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);  
  
C(i,j) = 0;  
C(i,j) += A(i, k) * B(k, j);
```

Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);
```

```
For (int i = 0; i<N; i++)
  For (int j = 0; j<N; j++)
    C[i, j] = 0;
    For (int k = 0; k<N; k++)
      C[i, j] += A[i,k]*B[k,j];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;  
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
C.parallel(i);
```

```
Parallel For (int i = 0; i<N; i++)  
  For (int j = 0; j<N; j++)  
    C[i, j] = 0;  
    For (int k = 0; k<N; k++)  
      C[i, j] += A[i,k]*B[k,j];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;
```

```
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,      i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)  
  For (int j0 = 0; j0<N/32; j0+=32)  
    For (int i1 = i0; i1<min(i0+64, N); i1++)  
      For (int j1 = j0; j1<min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0<N/32; k0++)  
      For (int i1 = i0; i1<min(i0+64, N); i1++)  
        For (int j1 = j0; j1<min(i0+32, N); j1++)  
          For (int k1 = 0; k1<min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```


Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;  
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,     i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0 < N/64; i0 += 64)  
  For (int j0 = 0; j0 < N/32; j0 += 32)  
    For (int i1 = i0; i1 < min(i0+64, N); i1++)  
      For (int j1 = j0; j1 < min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0 < N/32; k0++)  
      For (int i1 = i0; i1 < min(i0+64, N); i1++)  
        For (int j1 = j0; j1 < min(i0+32, N); j1++)  
          For (int k1 = 0; k1 < min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;
```

```
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,      i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)  
  For (int j0 = 0; j0<N/32; j0+=32)  
    For (int i1 = i0; i1<min(i0+64, N); i1++)  
      For (int j1 = j0; j1<min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0<N/32; k0++)  
      For (int i1 = i0; i1<min(i0+64, N); i1++)  
        For (int j1 = j0; j1<min(i0+32, N); j1++)  
          For (int k1 = 0; k1<min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;  
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,     i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0 < N/64; i0 += 64)  
  For (int j0 = 0; j0 < N/32; j0 += 32)  
    For (int i1 = i0; i1 < min(i0+64, N); i1++)  
      For (int j1 = j0; j1 < min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0 < N/32; k0++)  
      For (int i1 = i0; i1 < min(i0+64, N); i1++)  
        For (int j1 = j0; j1 < min(i0+32, N); j1++)  
          For (int k1 = 0; k1 < min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;
```

```
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,      i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)  
  For (int j0 = 0; j0<N/32; j0+=32)  
    For (int i1 = i0; i1<min(i0+64, N); i1++)  
      For (int j1 = j0; j1<min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0<N/32; k0++)  
      For (int i1 = i0; i1<min(i0+64, N); i1++)  
        For (int j1 = j0; j1<min(i0+32, N); j1++)  
          For (int k1 = 0; k1<min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;  
C.parallel(i);  
C.up(0).tile(i,j,          64,32,  
              i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
              i0,j0,k0,    i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0 < N/64; i0 += 64)  
  For (int j0 = 0; j0 < N/32; j0 += 32)  
    For (int i1 = i0; i1 < min(i0+64, N); i1++)  
      For (int j1 = j0; j1 < min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0 < N/32; k0++)  
      For (int i1 = i0; i1 < min(i0+64, N); i1++)  
        For (int j1 = j0; j1 < min(i0+32, N); j1++)  
          For (int k1 = 0; k1 < min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

// Algorithm

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

// Schedule

```
var i0, j0, k0, i1, j1, k1;
```

```
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,      i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)  
  For (int j0 = 0; j0<N/32; j0+=32)  
    For (int i1 = i0; i1<min(i0+64, N); i1++)  
      For (int j1 = j0; j1<min(i0+32, N); j1++)  
        C[i1, j1] = 0;  
    For (int k0 = 0; k0<N/32; k0++)  
      For (int i1 = i0; i1<min(i0+64, N); i1++)  
        For (int j1 = j0; j1<min(i0+32, N); j1++)  
          For (int k1 = 0; k1<min(k0+32, N); k1++)  
            C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

```
// Algorithm
```

```
var i(0, N), j(0, N), k(0, N);
```

```
input A(i,j), B(i,j);
```

```
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

```
// Schedule
```

```
var i0, j0, k0, i1, j1, k1;
```

```
C.parallel(i);
```

```
C.up(0).tile(i,j,          64,32,  
              i0,j0,        i1,j1)
```

```
C.up(1).tile(i,j,k,        64,32,32,  
              i0,j0,k0,     i1,j1,k1)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)
```

```
For (int j0 = 0; j0<N/32; j0+=32)
```

```
For (int i1 = i0; i1<min(i0+64, N); i1++)
```

```
For (int j1 = j0; j1<min(i0+32, N); j1++)
```

```
C[i1, j1] = 0;
```

```
For (int k0 = 0; k0<N/32; k0++)
```

```
For (int i1 = i0; i1<min(i0+64, N); i1++)
```

```
For (int j1 = j0; j1<min(i0+32, N); j1++)
```

```
For (int k1 = 0; k1<min(k0+32, N); k1++)
```

```
C[i1, j1] += A[i1, k1]*B[k1, j1];
```

Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);

// Schedule
var i0, j0, k0, i1, j1, k1;
C.parallel(i);
C.up(0).tile(i,j,          64,32,
             i0,j0,        i1,j1)
    .separate(i1, j1, 64, 32)

C.up(1).tile(i,j,k,        64,32,32,
             i0,j0,k0,      i1,j1,k1)
    .separate(j1, k1, 32, 32)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)
    For (int j0 = 0; j0<N/32; j0+=32)
        For (int i1 = i0; i1<i0+64; i1++)
            For (int j1 = j0; j1<i0+32; j1++)
                C[i1, j1] = 0;
            ...
        ...
    For (int k0 = 0; k0<N/32; k0++)
        For (int i1 = i0; i1<i0+64; i1++)
            For (int j1 = j0; j1<i0+32; j1++)
                For (int k1 = 0; k1<k0+32; k1++)
                    C[i1, j1] += A[i1, k1]*B[k1, j1];
            ...
        ...
    ...
```


Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);

// Schedule
var i0, j0, k0, i1, j1, k1;
C.parallel(i);
C.up(0).tile(i,j,          64,32,
             i0,j0,        i1,j1)
    .separate(i1, j1, 64, 32)
    .vectorize(j1,32)
C.up(1).tile(i,j,k,        64,32,32,
             i0,j0,k0,      i1,j1,k1)
    .separate(j1, k1, 32, 32)
    .vectorize(j1,32)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)
    For (int j0 = 0; j0<N/32; j0+=32)
        For (int i1 = i0; i1<i0+64; i1++)
            Vectorized For (int j1 = j0; j1<i0+32; j1++)
                C[i1, j1] = 0;
            ...
            ...
        For (int k0 = 0; k0<N/32; k0++)
            For (int i1 = i0; i1<i0+64; i1++)
                Vectorized For (int j1 = j0; j1<i0+32; j1++)
                    For (int k1 = 0; k1<k0+32; k1++)
                        C[i1, j1] += A[i1, k1]*B[k1, j1];
                    ...
                    ...
                    ...
```

Matrix Multiplication

```
// Algorithm
```

```
var i(0, N), j(0, N), k(0, N);  
input A(i,j), B(i,j);  
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

```
// Schedule
```

```
var i0, j0, k0, i1, j1, k1;  
C.parallel(i);  
C.up(0).tile(i,j,          64,32,  
             i0,j0,        i1,j1)  
    .separate(i1, j1, 64, 32)  
    .vectorize(j1,32).unroll(i1)  
C.up(1).tile(i,j,k,        64,32,32,  
             i0,j0,k0,      i1,j1,k1)  
    .separate(j1, k1, 32, 32)  
    .vectorize(j1,32).unroll(k1)
```

```
Parallel For (int i0 = 0; i0<N/64; i+=64)  
    For (int j0 = 0; j0<N/32; j0+=32)  
        Unroll For (int i1 = i0; i1<i0+64; i1++)  
        Vectorized For (int j1 = j0; j1<i0+32; j1++)  
            C[i1, j1] = 0;  
            ...  
            ...  
            For (int k0 = 0; k0<N/32; k0++)  
                For (int i1 = i0; i1<i0+64; i1++)  
                    Vectorized For (int j1 = j0; j1<i0+32; j1++)  
                        Unroll For (int k1 = 0; k1<k0+32; k1++)  
                            C[i1, j1] += A[i1, k1]*B[k1, j1];  
                            ...  
                            ...  
                            ...
```

Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);

// Schedule
var i0, j0, k0, i1, j1, k1;
C.parallel(i);
C.up(0).tile(i,j,          64,32,
             i0,j0,        i1,j1)
    .separate(i1, j1, 64, 32)
    .vectorize(j1,32).unroll(i1);
C.up(1).tile(i,j,k,        64,32,32,
             i0,j0,k0,      i1,j1,k1)
    .separate(j1, k1, 32, 32)
    .vectorize(j1,32).unroll(k1);

B.pack(64,32,32);
```

Matrix Multiplication

```
// Algorithm
var i(0, N), j(0, N), k(0, N);
input A(i,j), B(i,j);
computation C(i,j);

C(i,j) = 0;
C(i,j) += A(i, k) * B(k, j);

// Schedule
var i0, j0, k0, i1, j1, k1;
C.parallel(i);
C.up(0).tile(i,j,          64,32,
             i0,j0,        i1,j1)
    .separate(i1, j1, 64, 32)
    .vectorize(j1,32).unroll(i1);
C.up(1).tile(i,j,k,        64,32,32,
             i0,j0,k0,      i1,j1,k1)
    .separate(j1, k1, 32, 32)
    .vectorize(j1,32).unroll(k1);

B.pack(64,32,32);
```

Speedup: 350x
(4 Cores)

Matrix Multiplication (GPU)

```
// Algorithm
```

```
var i(0, N), j(0, N), k(0, N);
```

```
input A(i,j), B(i,j);
```

```
computation C(i,j);
```

```
C(i,j) = 0;
```

```
C(i,j) += A(i, k) * B(k, j);
```

```
// Schedule
```

```
var i0, j0, k0, i1, j1, k1;
```

```
C.gpu_tile(i,j,k, 64,32,32, i0,j0,k0, i1,j1,k1).separate();
```

```
A.cache_shared_at(C.update(1), j0);
```

```
A.cache_local_at(C.update(1), j0);
```

```
B.cache_shared_at(C.update(1), j0);
```

```
B.cache_local_at(C.update(1), j0);
```

Matrix Multiplication (CUDA)

```
// Host code
```

```
float *b_A_glb, *b_B_glb, *b_C_glb;
extern "C" void matmul(float *b_A, float *b_B, float *b_C)
{
    cudaMalloc(&b_A_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_B_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_C_glb, 3072 * 3072 * sizeof(float));
    cudaMemcpy(b_A_glb, b_A, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_B_glb, b_B, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_C_glb, b_C, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    {
        dim3 blocks((15 + 1), (11 + 1), 1);
        dim3 threads((15 + 1), (15 + 1), 1);
        _kernel_0<<<blocks, threads>>>(b_A_glb, b_B_glb, b_C_glb);
    };
    cudaMemcpy(b_C, b_C_glb, 3072 * 3072 * sizeof(float), cudaMemcpyDeviceToHost);
    cudaFree(b_A_glb);
    cudaFree(b_B_glb);
    cudaFree(b_C_glb);
}
```

```
// Device code
```

```
static __global__ void _kernel_0(float *b_A_glb, float *b_B_glb, float *b_C_glb)
{
    const int32_t __bx__ = (blockIdx.x + 0);
    const int32_t __by__ = (blockIdx.y + 0);
    const int32_t __tx__ = (threadIdx.x + 0);
    const int32_t __ty__ = (threadIdx.y + 0);
    __shared__ float b_A_shr[16 * 193];
    __shared__ float b_B_shr[16 * 256];
    float b_A_reg;
    float b_B_reg[16];
    float b_acc[12 * 16];
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] = 0;
        };
    };

    for (int32_t c9 = 0; (c9 <= 191); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 11); (c11 += 1))
        {
            b_A_shr[((0 + (((12 * __tx__) + c11) * 1)) + (__ty__ * (1 * ((int32_t) 193))))] =
                b_A_glb[((0 + ((__ty__ + (16 * c9)) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c11) * (1 * ((int32_t) 3072))))];
        }
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_B_shr[((0 + ((__tx__ + (16 * __ty__)) * 1)) + (c11 * (1 * ((int32_t) 256))))] =
                b_B_glb[((0 + (((256 * __by__) + __tx__) + (16 * __ty__)) * 1)) + (((16 * c9) + c11) * (1 * ((int32_t) 3072))))];
        }
        __syncthreads();
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            for (int32_t c13 = 0; (c13 <= 15); (c13 += 1))
            {
                b_B_reg[(0 + (c13 * 1))] =
                    b_B_shr[((0 + ((16 * __ty__) + c13) * 1)) + (c11 * (1 * ((int32_t) 256))))];
            };
            for (int32_t c13 = 0; (c13 <= 11); (c13 += 1))
            {
                b_A_reg = b_A_shr[((0 + (((12 * __tx__) + c13) * 1)) + (c11 * (1 * ((int32_t) 193))))];
                for (int32_t c15 = 0; (c15 <= 15); (c15 += 1))
                {
                    b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] =
                        (b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] + (b_A_reg * b_B_reg[(0 + (c15 * 1)))]);
                };
            };
        };
        __syncthreads();
    };
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072))))] =
                ((b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] * 3) +
                 (b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072))))] * 2));
        };
    };
};
```


Matrix Multiplication (CUDA)

```
// Host code
```

```
float *b_A_glb, *b_B_glb, *b_C_glb;
extern "C" void matmul(float *b_A, float *b_B, float *b_C)
{
    cudaMalloc(&b_A_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_B_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_C_glb, 3072 * 3072 * sizeof(float));
    cudaMemcpy(b_A_glb, b_A, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_B_glb, b_B, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_C_glb, b_C, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    {
        dim3 blocks((15 + 1), (11 + 1), 1);
        dim3 threads((15 + 1), (15 + 1), 1);
        _kernel_0<<<blocks, threads>>>(b_A_glb, b_B_glb, b_C_glb);
    };
    cudaMemcpy(b_C, b_C_glb, 3072 * 3072 * sizeof(float), cudaMemcpyDeviceToHost);
    cudaFree(b_A_glb);
    cudaFree(b_B_glb);
    cudaFree(b_C_glb);
}
```

```
// Device code
```

```
static __global__ void _kernel_0(float *b_A_glb, float *b_B_glb, float *b_C_glb)
{
    const int32_t __bx__ = (blockIdx.x + 0);
    const int32_t __by__ = (blockIdx.y + 0);
    const int32_t __tx__ = (threadIdx.x + 0);
    const int32_t __ty__ = (threadIdx.y + 0);
    __shared__ float b_A_shr[16 * 193];
    __shared__ float b_B_shr[16 * 256];
    float b_A_reg;
    float b_B_reg[16];
    float b_acc[12 * 16];
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] = 0;
        };
    };

    for (int32_t c9 = 0; (c9 <= 191); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 11); (c11 += 1))
        {
            b_A_shr[((0 + (((12 * __tx__) + c11) * 1)) + (__ty__ * (1 * ((int32_t) 193))))] =
                b_A_glb[((0 + ((__ty__ + (16 * c9)) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c11) * (1 * ((int32_t) 3072))))];
        }
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_B_shr[((0 + ((__tx__ + (16 * __ty__)) * 1)) + (c11 * (1 * ((int32_t) 256))))] =
                b_B_glb[((0 + (((256 * __by__) + __tx__) + (16 * __ty__)) * 1)) + (((16 * c9) + c11) * (1 * ((int32_t) 3072))))];
        }
        __syncthreads();
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            for (int32_t c13 = 0; (c13 <= 15); (c13 += 1))
            {
                b_B_reg[(0 + (c13 * 1))] =
                    b_B_shr[((0 + ((16 * __ty__) + c13) * 1)) + (c11 * (1 * ((int32_t) 256))))];
            };
            for (int32_t c13 = 0; (c13 <= 11); (c13 += 1))
            {
                b_A_reg = b_A_shr[((0 + (((12 * __tx__) + c13) * 1)) + (c11 * (1 * ((int32_t) 193))))];
                for (int32_t c15 = 0; (c15 <= 15); (c15 += 1))
                {
                    b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] =
                        (b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] + (b_A_reg * b_B_reg[(0 + (c15 * 1)))]);
                };
            };
        };
        __syncthreads();
    };
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] =
                ((b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] * 3) +
                 (b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] * 2));
        };
    };
};
```

Host/device data copies

Matrix Multiplication (CUDA)

// Host code

```
float *b_A_glb, *b_B_glb, *b_C_glb;
extern "C" void matmul(float *b_A, float *b_B, float *b_C)
{
    cudaMalloc(&b_A_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_B_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_C_glb, 3072 * 3072 * sizeof(float));
    cudaMemcpy(b_A_glb, b_A, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_B_glb, b_B, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_C_glb, b_C, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    {
        dim3 blocks((15 + 1), (11 + 1), 1);
        dim3 threads((15 + 1), (15 + 1), 1);
        _kernel_0<<<blocks, threads>>>(b_A_glb, b_B_glb, b_C_glb);
    };
    cudaMemcpy(b_C, b_C_glb, 3072 * 3072 * sizeof(float), cudaMemcpyDeviceToHost);
    cudaFree(b_A_glb);
    cudaFree(b_B_glb);
    cudaFree(b_C_glb);
}
```

// Device code

```
static __global__ void _kernel_0(float *b_A_glb, float *b_B_glb, float *b_C_glb)
{
    const int32_t __bx__ = (blockIdx.x + 0);
    const int32_t __by__ = (blockIdx.y + 0);
    const int32_t __tx__ = (threadIdx.x + 0);
    const int32_t __ty__ = (threadIdx.y + 0);
    __shared__ float b_A_shr[16 * 193];
    __shared__ float b_B_shr[16 * 256];
    float b_A_reg;
    float b_B_reg[16];
    float b_acc[12 * 16];
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] = 0;
        };
    };

    for (int32_t c9 = 0; (c9 <= 191); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 11); (c11 += 1))
        {
            b_A_shr[((0 + (((12 * __tx__) + c11) * 1)) + (__ty__ * (1 * ((int32_t) 193))))] =
                b_A_glb[((0 + ((__ty__ + (16 * c9)) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c11) * (1 * ((int32_t) 3072))))];
        }
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_B_shr[((0 + ((__tx__ + (16 * __ty__)) * 1)) + (c11 * (1 * ((int32_t) 256))))] =
                b_B_glb[((0 + (((256 * __by__) + __tx__) + (16 * __ty__)) * 1)) + (((16 * c9) + c11) * (1 * ((int32_t) 3072))))];
        }
        __syncthreads();
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            for (int32_t c13 = 0; (c13 <= 15); (c13 += 1))
            {
                b_B_reg[(0 + (c13 * 1))] =
                    b_B_shr[((0 + (((16 * __ty__) + c13) * 1)) + (c11 * (1 * ((int32_t) 256))))];
            };
            for (int32_t c13 = 0; (c13 <= 11); (c13 += 1))
            {
                b_A_reg = b_A_shr[((0 + (((12 * __tx__) + c13) * 1)) + (c11 * (1 * ((int32_t) 193))))];
                for (int32_t c15 = 0; (c15 <= 15); (c15 += 1))
                {
                    b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] =
                        (b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] + (b_A_reg * b_B_reg[(0 + (c15 * 1)))]);
                };
            };
        };
        __syncthreads();
    };
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] =
                ((b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] * 3) +
                 (b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] * 2));
        };
    };
};
```

Host/device data copies

global —> shared

Matrix Multiplication (CUDA)

// Host code

```
float *b_A_glb, *b_B_glb, *b_C_glb;
extern "C" void matmul(float *b_A, float *b_B, float *b_C)
{
    cudaMalloc(&b_A_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_B_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_C_glb, 3072 * 3072 * sizeof(float));
    cudaMemcpy(b_A_glb, b_A, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_B_glb, b_B, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_C_glb, b_C, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    {
        dim3 blocks((15 + 1), (11 + 1), 1);
        dim3 threads((15 + 1), (15 + 1), 1);
        _kernel_0<<<blocks, threads>>>(b_A_glb, b_B_glb, b_C_glb);
    };
    cudaMemcpy(b_C, b_C_glb, 3072 * 3072 * sizeof(float), cudaMemcpyDeviceToHost);
    cudaFree(b_A_glb);
    cudaFree(b_B_glb);
    cudaFree(b_C_glb);
}
```

// Device code

```
static __global__ void _kernel_0(float *b_A_glb, float *b_B_glb, float *b_C_glb)
{
    const int32_t __bx__ = (blockIdx.x + 0);
    const int32_t __by__ = (blockIdx.y + 0);
    const int32_t __tx__ = (threadIdx.x + 0);
    const int32_t __ty__ = (threadIdx.y + 0);
    __shared__ float b_A_shr[16 * 193];
    __shared__ float b_B_shr[16 * 256];
    float b_A_reg;
    float b_B_reg[16];
    float b_acc[12 * 16];
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] = 0;
        };
    };

    for (int32_t c9 = 0; (c9 <= 191); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 11); (c11 += 1))
        {
            b_A_shr[((0 + (((12 * __tx__) + c11) * 1)) + (__ty__ * (1 * ((int32_t) 193))))] =
                b_A_glb[((0 + ((__ty__ + (16 * c9)) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c11) * (1 * ((int32_t) 3072))))];
        }
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_B_shr[((0 + ((__tx__ + (16 * __ty__)) * 1)) + (c11 * (1 * ((int32_t) 256))))] =
                b_B_glb[((0 + (((256 * __by__) + __tx__) + (16 * __ty__)) * 1)) + (((16 * c9) + c11) * (1 * ((int32_t) 3072))))];
        }
        __syncthreads();
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            for (int32_t c13 = 0; (c13 <= 15); (c13 += 1))
            {
                b_B_reg[(0 + (c13 * 1))] =
                    b_B_shr[((0 + (((16 * __ty__) + c13) * 1)) + (c11 * (1 * ((int32_t) 256))))];
            };
            for (int32_t c13 = 0; (c13 <= 11); (c13 += 1))
            {
                b_A_reg = b_A_shr[((0 + (((12 * __tx__) + c13) * 1)) + (c11 * (1 * ((int32_t) 193))))];
                for (int32_t c15 = 0; (c15 <= 15); (c15 += 1))
                {
                    b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] =
                        (b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] + (b_A_reg * b_B_reg[(0 + (c15 * 1)))]);
                };
            };
        };
        __syncthreads();
    };
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] =
                ((b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] * 3) +
                 (b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] * 2));
        };
    };
};
```

Host/device data copies

global —> shared

shared —> local

Matrix Multiplication (CUDA)

```
// Host code

float *b_A_glb, *b_B_glb, *b_C_glb;
extern "C" void matmul(float *b_A, float *b_B, float *b_C)
{
    cudaMalloc(&b_A_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_B_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_C_glb, 3072 * 3072 * sizeof(float));
    cudaMemcpy(b_A_glb, b_A, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_B_glb, b_B, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_C_glb, b_C, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    {
        dim3 blocks((15 + 1), (11 + 1), 1);
        dim3 threads((15 + 1), (15 + 1), 1);
        _kernel_0<<<blocks, threads>>>(b_A_glb, b_B_glb, b_C_glb);
    };
    cudaMemcpy(b_C, b_C_glb, 3072 * 3072 * sizeof(float), cudaMemcpyDeviceToHost);
    cudaFree(b_A_glb);
    cudaFree(b_B_glb);
    cudaFree(b_C_glb);
}

// Device code
static __global__ void _kernel_0(float *b_A_glb, float *b_B_glb, float *b_C_glb)
{
    const int32_t __bx__ = (blockIdx.x + 0);
    const int32_t __by__ = (blockIdx.y + 0);
    const int32_t __tx__ = (threadIdx.x + 0);
    const int32_t __ty__ = (threadIdx.y + 0);
    __shared__ float b_A_shr[16 * 193];
    __shared__ float b_B_shr[16 * 256];
    float b_A_reg;
    float b_B_reg[16];
    float b_acc[12 * 16];
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] = 0;
        };
    };
    for (int32_t c9 = 0; (c9 <= 191); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 11); (c11 += 1))
        {
            b_A_shr[((0 + (((12 * __tx__) + c11) * 1)) + (__ty__ * (1 * ((int32_t) 193))))] =
                b_A_glb[((0 + ((__ty__ + (16 * c9)) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c11) * (1 * ((int32_t) 3072))))];
        }
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_B_shr[((0 + ((__tx__ + (16 * __ty__)) * 1)) + (c11 * (1 * ((int32_t) 256))))] =
                b_B_glb[((0 + (((256 * __by__) + __tx__) + (16 * __ty__)) * 1)) + (((16 * c9) + c11) * (1 * ((int32_t) 3072))))];
        }
        __syncthreads();
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            for (int32_t c13 = 0; (c13 <= 15); (c13 += 1))
            {
                b_B_reg[(0 + (c13 * 1))] =
                    b_B_shr[((0 + (((16 * __ty__) + c13) * 1)) + (c11 * (1 * ((int32_t) 256))))];
            };
            for (int32_t c13 = 0; (c13 <= 11); (c13 += 1))
            {
                b_A_reg = b_A_shr[((0 + (((12 * __tx__) + c13) * 1)) + (c11 * (1 * ((int32_t) 193))))];
                for (int32_t c15 = 0; (c15 <= 15); (c15 += 1))
                {
                    b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] =
                        (b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] + (b_A_reg * b_B_reg[(0 + (c15 * 1)))]);
                };
            };
        };
        __syncthreads();
    };
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] =
                ((b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] * 3) +
                 (b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c9) * (1 * ((int32_t) 3072)))] * 2));
        };
    };
};
```

Host/device data copies

**Compute shared/local
memory sizes**

global —> shared

shared —> local

Matrix Multiplication (CUDA)

```
// Host code

float *b_A_glb, *b_B_glb, *b_C_glb;
extern "C" void matmul(float *b_A, float *b_B, float *b_C)
{
    cudaMalloc(&b_A_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_B_glb, 3072 * 3072 * sizeof(float));
    cudaMalloc(&b_C_glb, 3072 * 3072 * sizeof(float));
    cudaMemcpy(b_A_glb, b_A, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_B_glb, b_B, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(b_C_glb, b_C, 3072 * 3072 * sizeof(float), cudaMemcpyHostToDevice);
    {
        dim3 blocks((15 + 1), (11 + 1), 1);
        dim3 threads((15 + 1), (15 + 1), 1);
        _kernel_0<<<blocks, threads>>>(b_A_glb, b_B_glb, b_C_glb);
    }
    cudaMemcpy(b_C, b_C_glb, 3072 * 3072 * sizeof(float), cudaMemcpyDeviceToHost);
    cudaFree(b_A_glb);
    cudaFree(b_B_glb);
    cudaFree(b_C_glb);
}

// Device code
static __global__ void _kernel_0(float *b_A_glb, float *b_B_glb, float *b_C_glb)
{
    const int32_t __bx__ = (blockIdx.x + 0);
    const int32_t __by__ = (blockIdx.y + 0);
    const int32_t __tx__ = (threadIdx.x + 0);
    const int32_t __ty__ = (threadIdx.y + 0);
    __shared__ float b_A_shr[16 * 193];
    __shared__ float b_B_shr[16 * 256];
    float b_A_reg;
    float b_B_reg[16];
    float b_acc[12 * 16];
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] = 0;
        }
    };
    for (int32_t c9 = 0; (c9 <= 191); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 11); (c11 += 1))
        {
            b_A_shr[((0 + (((12 * __tx__) + c11) * 1)) + (__ty__ * (1 * ((int32_t) 193))))] =
                b_A_glb[((0 + ((__ty__ + (16 * c9)) * 1)) + (((192 * __bx__) + (12 * __tx__)) + c11) * (1 * ((int32_t) 3072))))];
        }
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_B_shr[((0 + ((__tx__ + (16 * __ty__)) * 1)) + (c11 * (1 * ((int32_t) 256))))] =
                b_B_glb[((0 + (((256 * __by__) + __tx__) + (16 * __ty__)) * 1)) + (((16 * c9) + c11) * (1 * ((int32_t) 3072))))];
        }
    }
    __syncthreads();
    for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
    {
        for (int32_t c13 = 0; (c13 <= 15); (c13 += 1))
        {
            b_B_reg[(0 + (c13 * 1))] =
                b_B_shr[((0 + (((16 * __ty__) + c13) * 1)) + (c11 * (1 * ((int32_t) 256))))];
        }
    };
    for (int32_t c13 = 0; (c13 <= 11); (c13 += 1))
    {
        b_A_reg = b_A_shr[((0 + (((12 * __tx__) + c13) * 1)) + (c11 * (1 * ((int32_t) 193))))];
        for (int32_t c15 = 0; (c15 <= 15); (c15 += 1))
        {
            b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] =
                (b_acc[((0 + (c15 * 1)) + (c13 * (1 * ((int32_t) 16))))] + (b_A_reg * b_B_reg[(0 + (c15 * 1)))]));
        }
    };
    __syncthreads();
    for (int32_t c9 = 0; (c9 <= 11); (c9 += 1))
    {
        for (int32_t c11 = 0; (c11 <= 15); (c11 += 1))
        {
            b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__) + c9) * (1 * ((int32_t) 3072))))] =
                ((b_acc[((0 + (c11 * 1)) + (c9 * (1 * ((int32_t) 16))))] * 3) +
                 (b_C_glb[((0 + (((256 * __by__) + (16 * __ty__)) + c11) * 1)) + (((192 * __bx__) + (12 * __tx__) + c9) * (1 * ((int32_t) 3072))))] * 2));
        }
    };
};
```

Host/device data copies

Compute shared/local
memory sizes

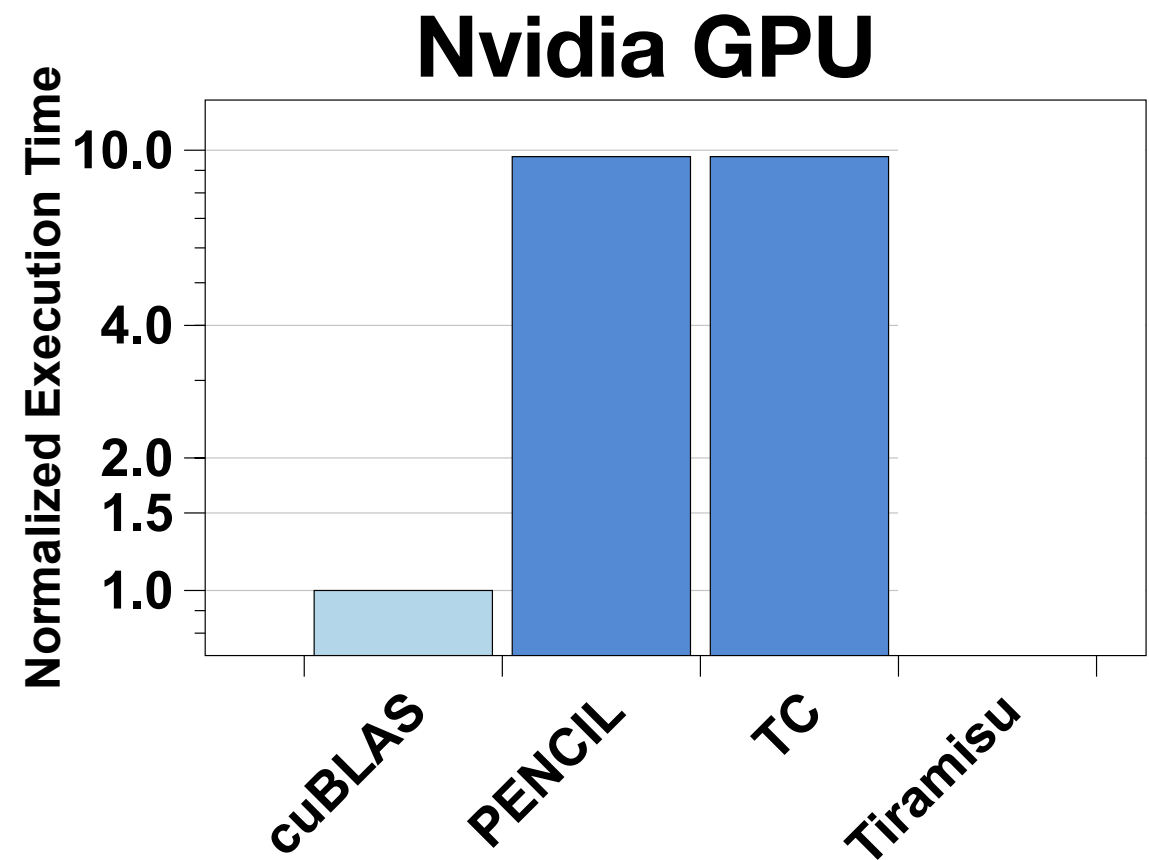
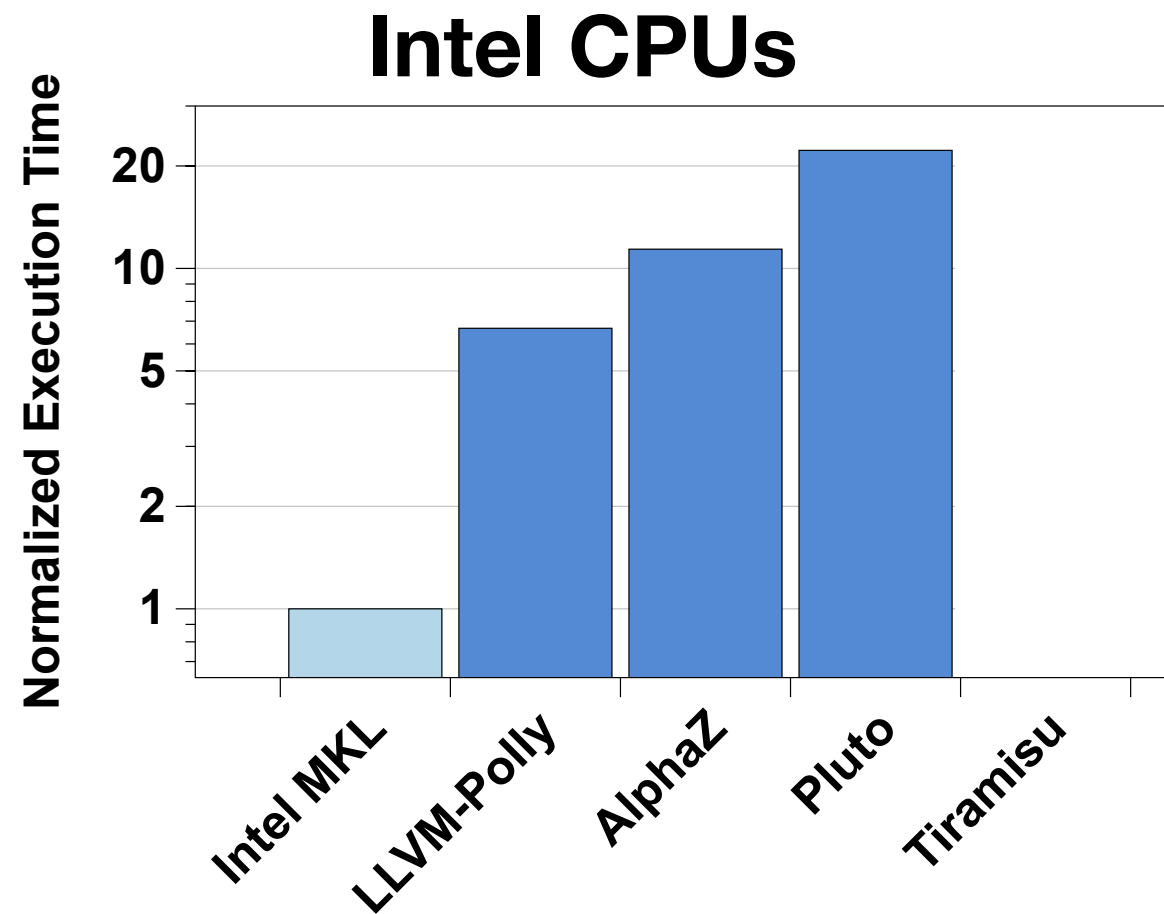
global —> shared

shared —> local

Synchronization

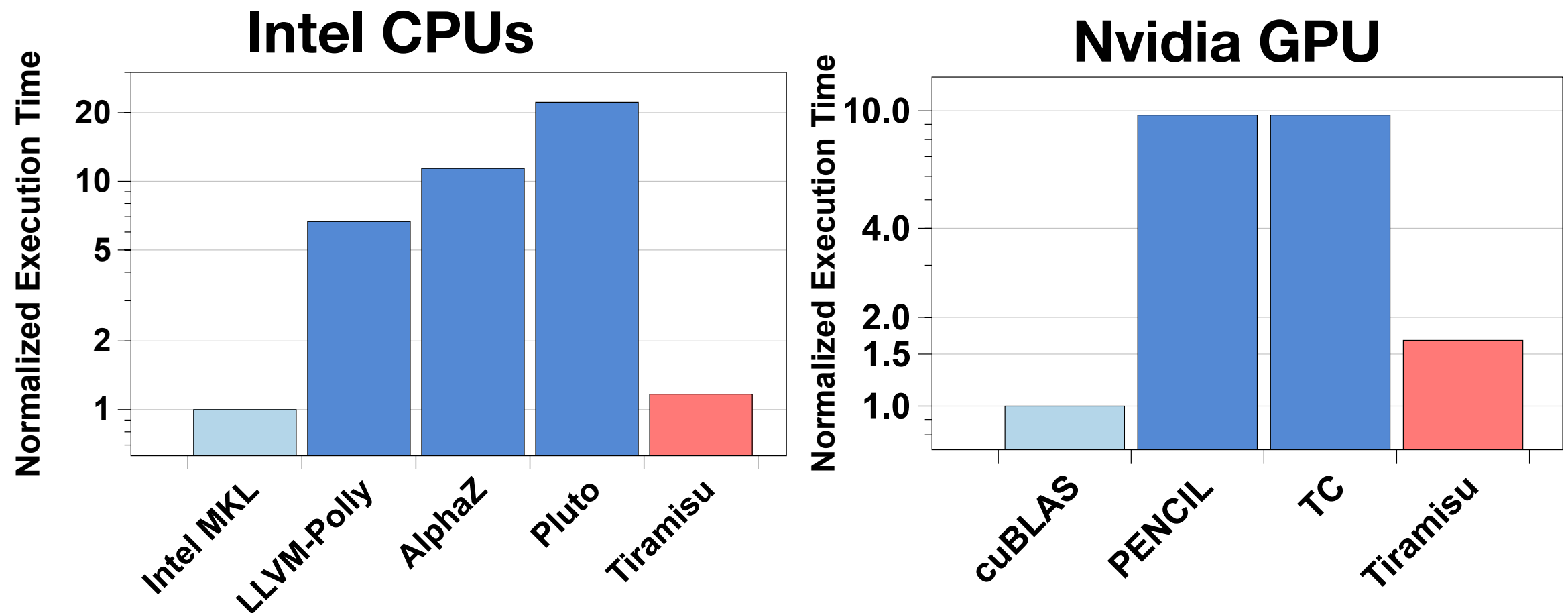
GEMM (Comparison with Polyhedral Compilers)

GEMM (Comparison with Polyhedral Compilers)



[Matrix size: 1060x1060]
[4-Cores CPU, Nvidia Pascal P4]

GEMM (Comparison with Polyhedral Compilers)



[Matrix size: 1060x1060]
[4-Cores CPU, Nvidia Pascal P4]

Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Outline

- **Phase I**

- Tiramisu Overview
- Example
- RNNs
- Sparse DNNs

- **Phase II**

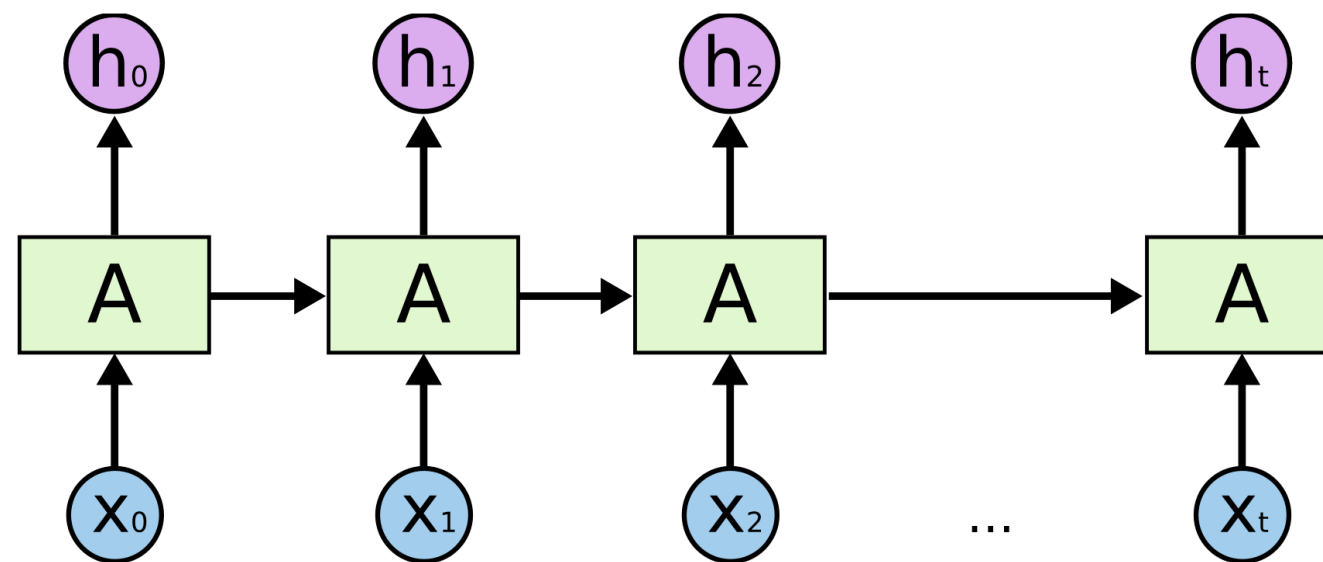
- Automatic Optimization

Outline

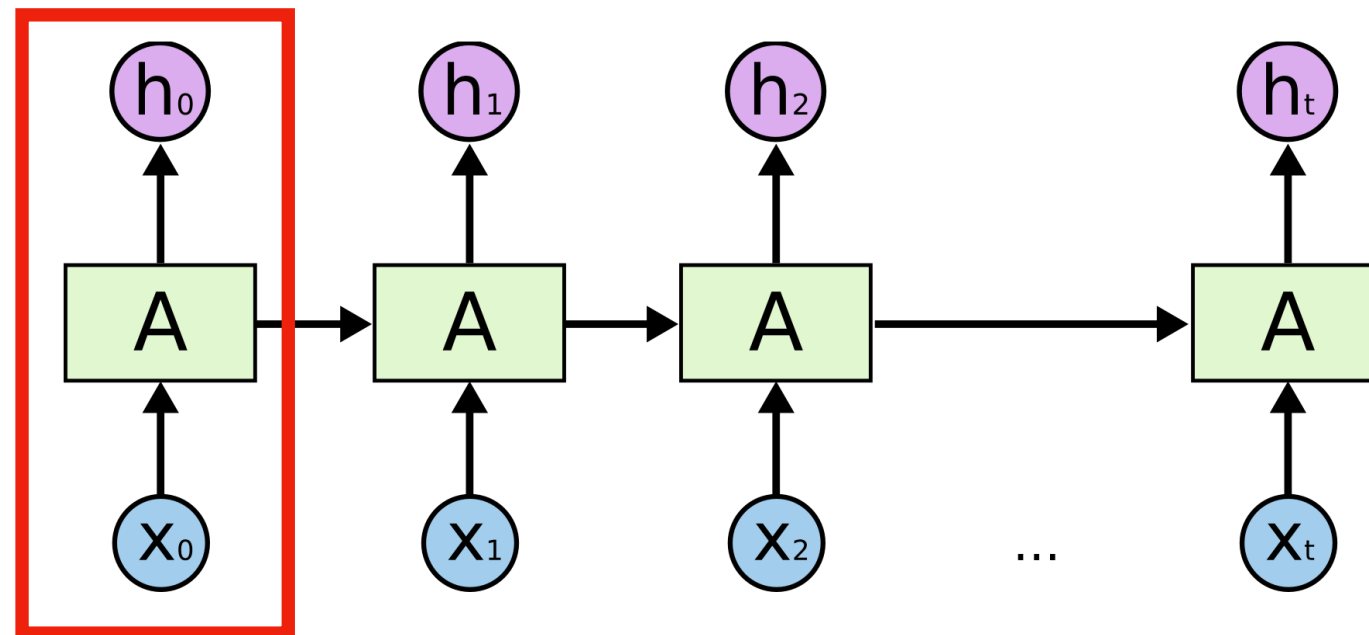
- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Recurrent Neural Networks

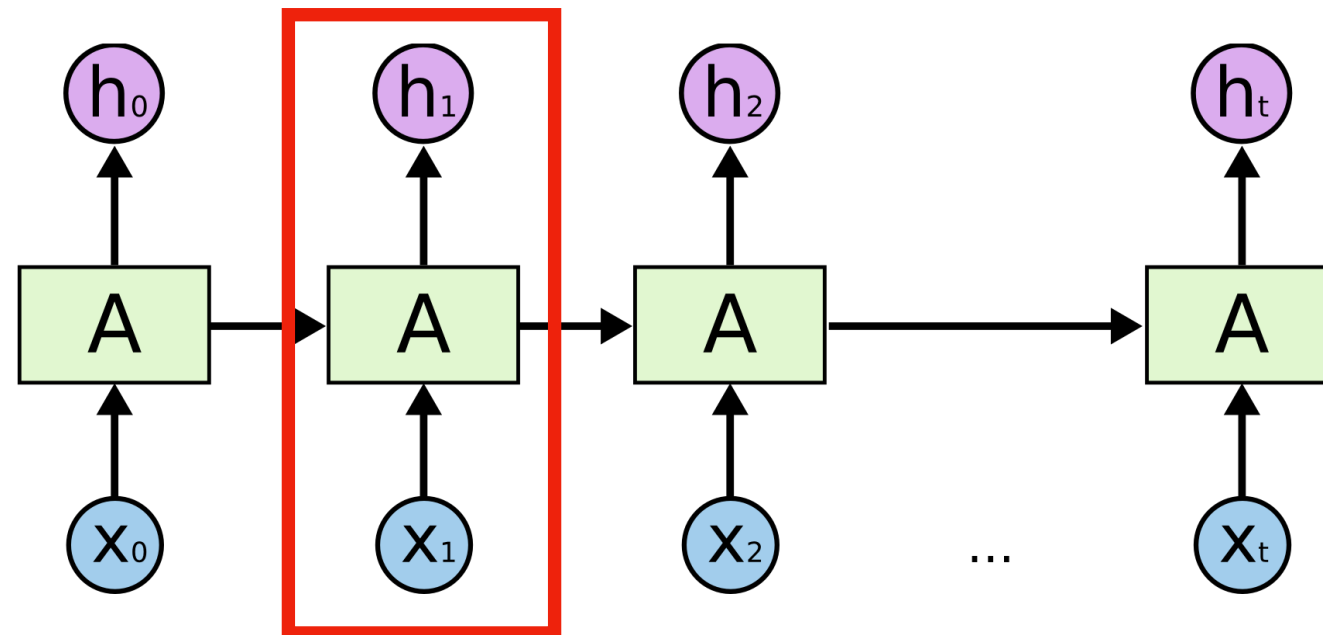
Recurrent Neural Networks



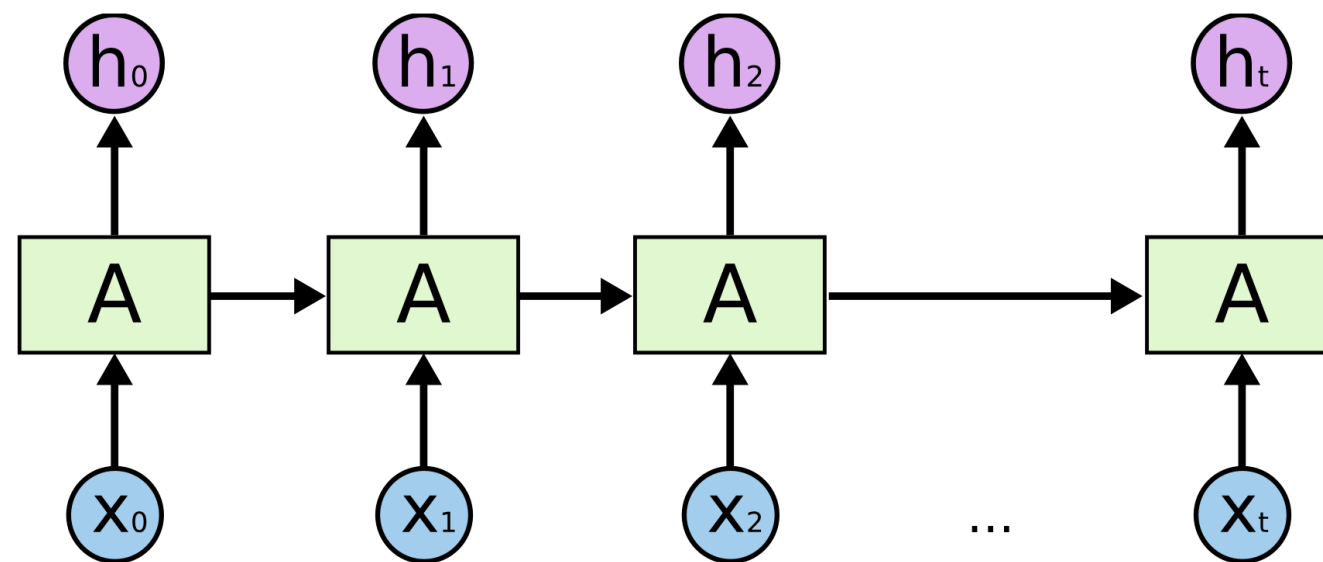
Recurrent Neural Networks



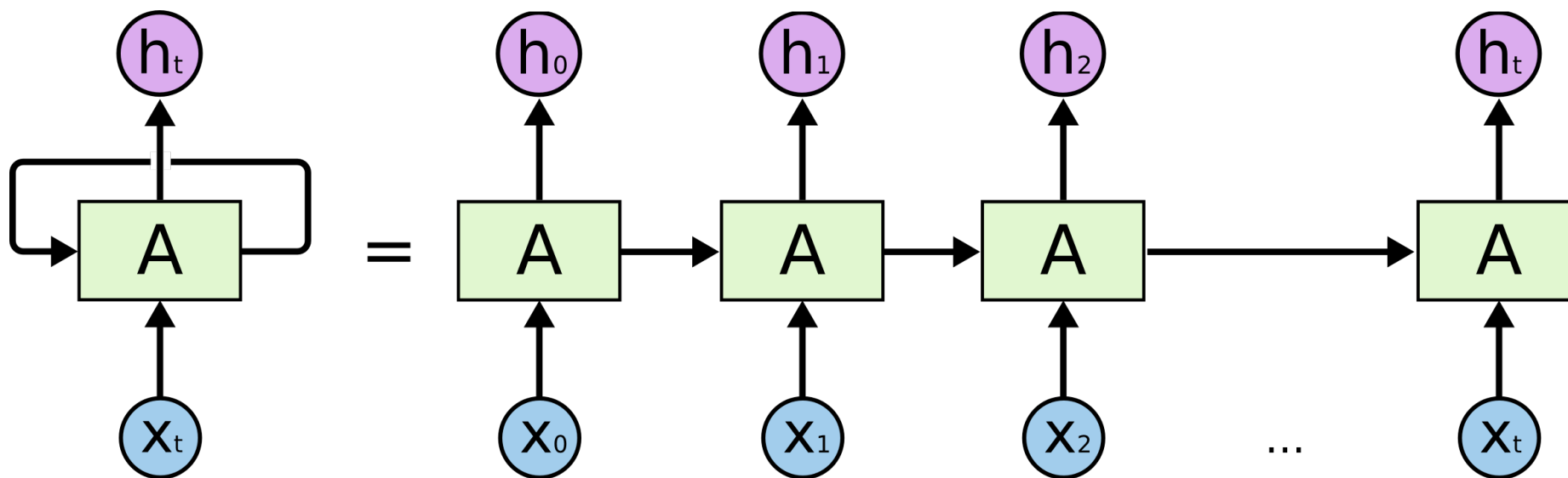
Recurrent Neural Networks



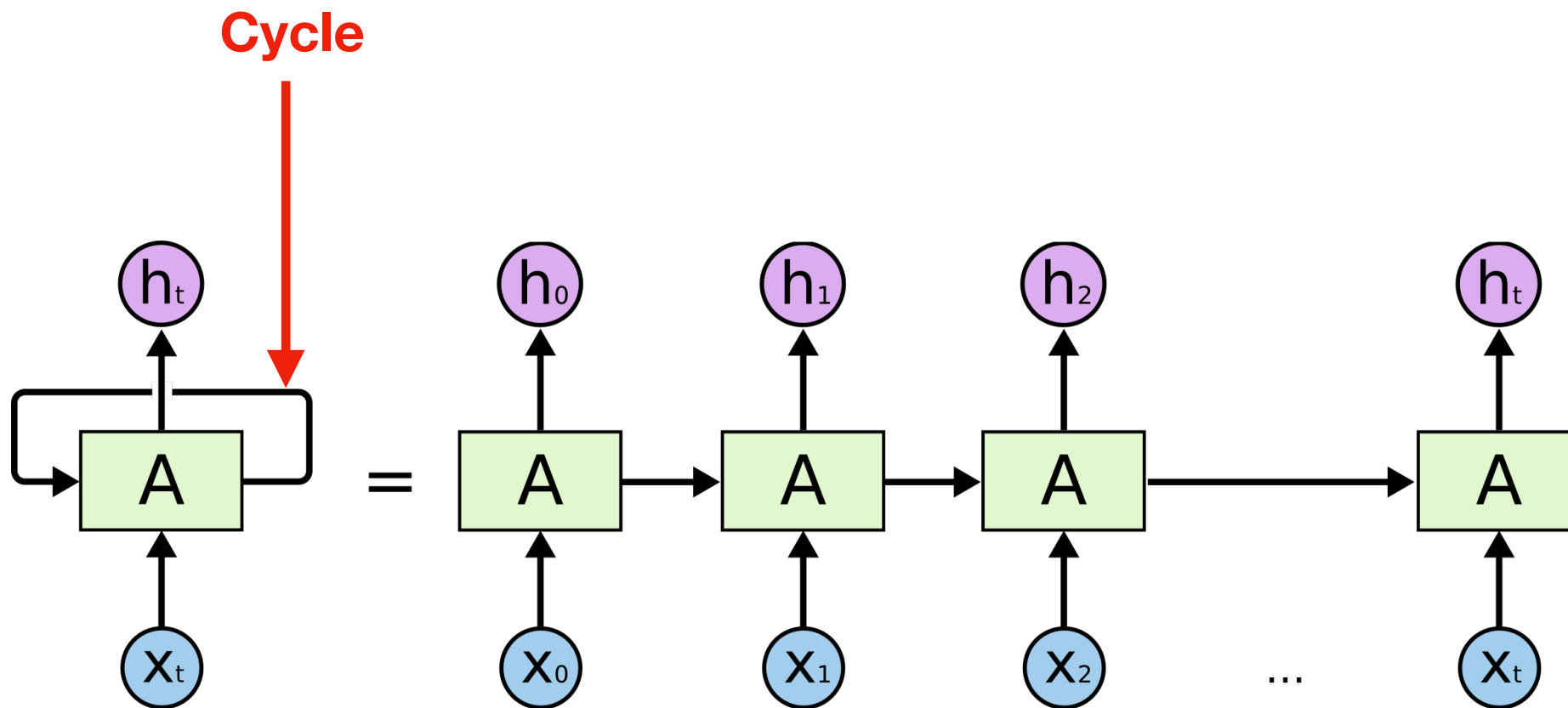
Recurrent Neural Networks



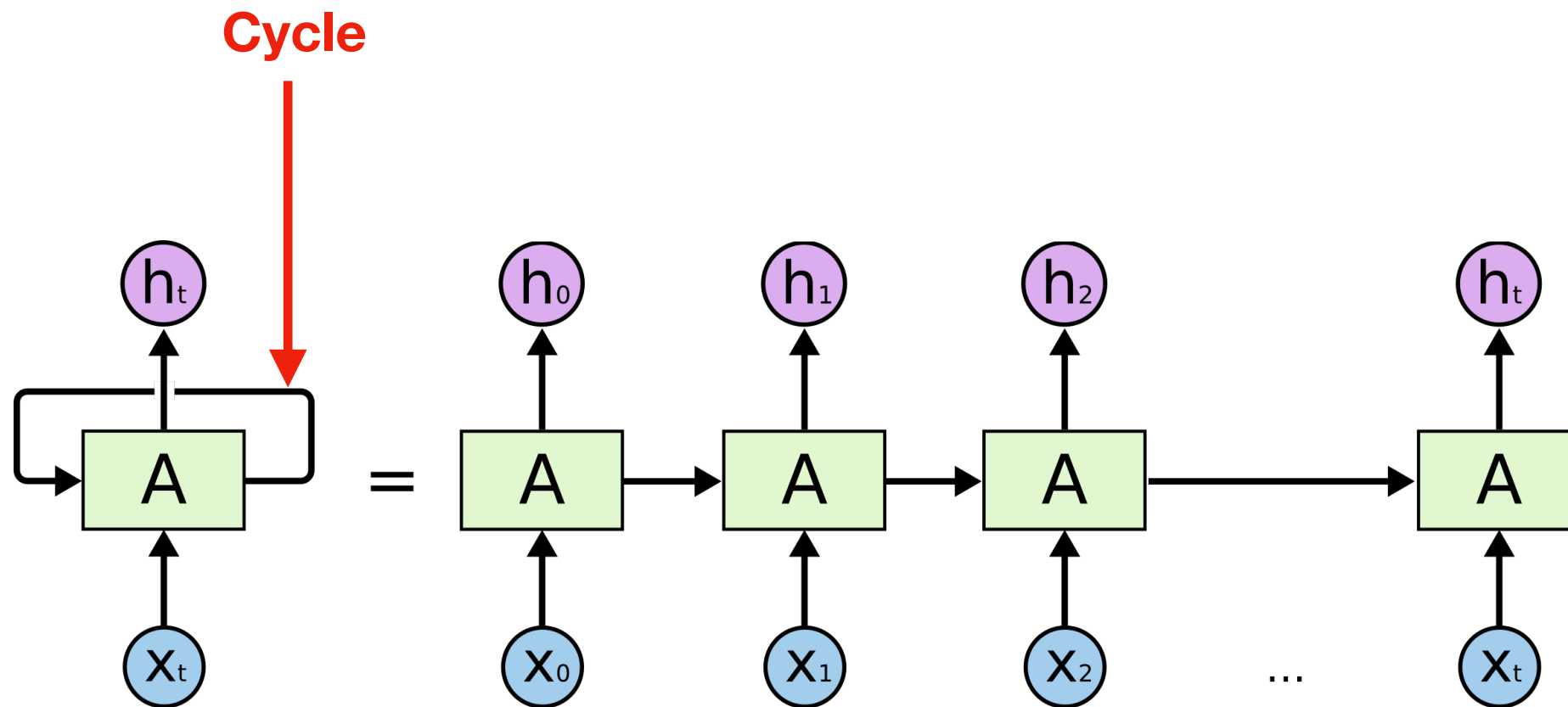
Recurrent Neural Networks



Recurrent Neural Networks

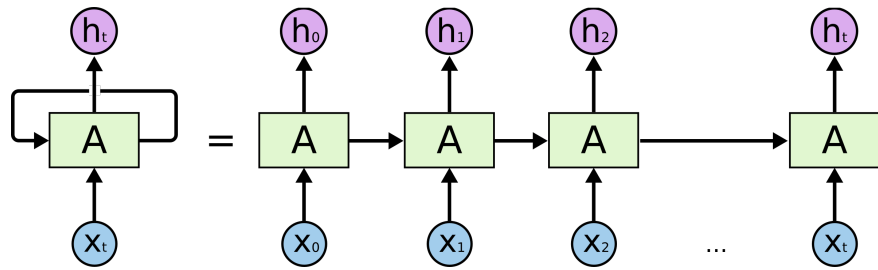


Recurrent Neural Networks

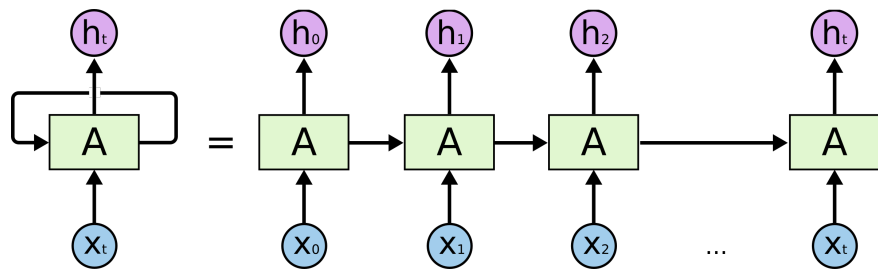


We use dependence analysis!

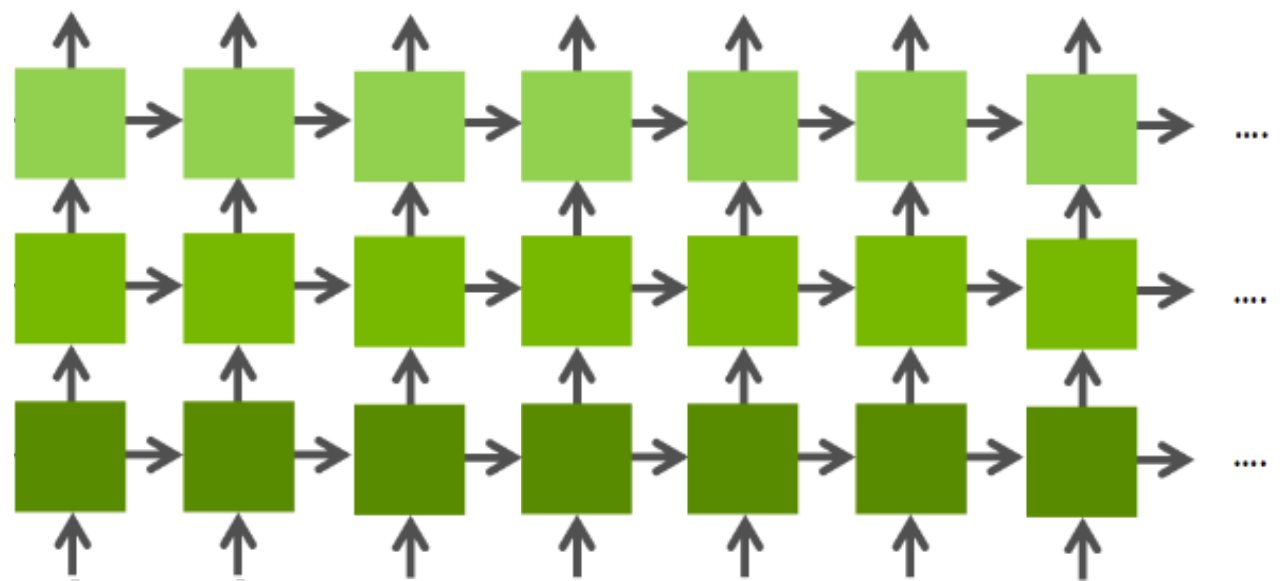
Recurrent Neural Networks



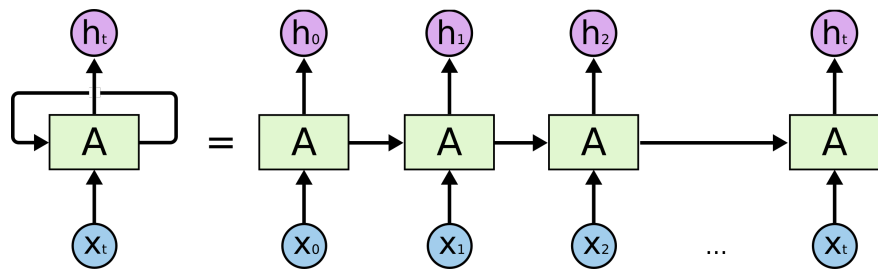
Recurrent Neural Networks



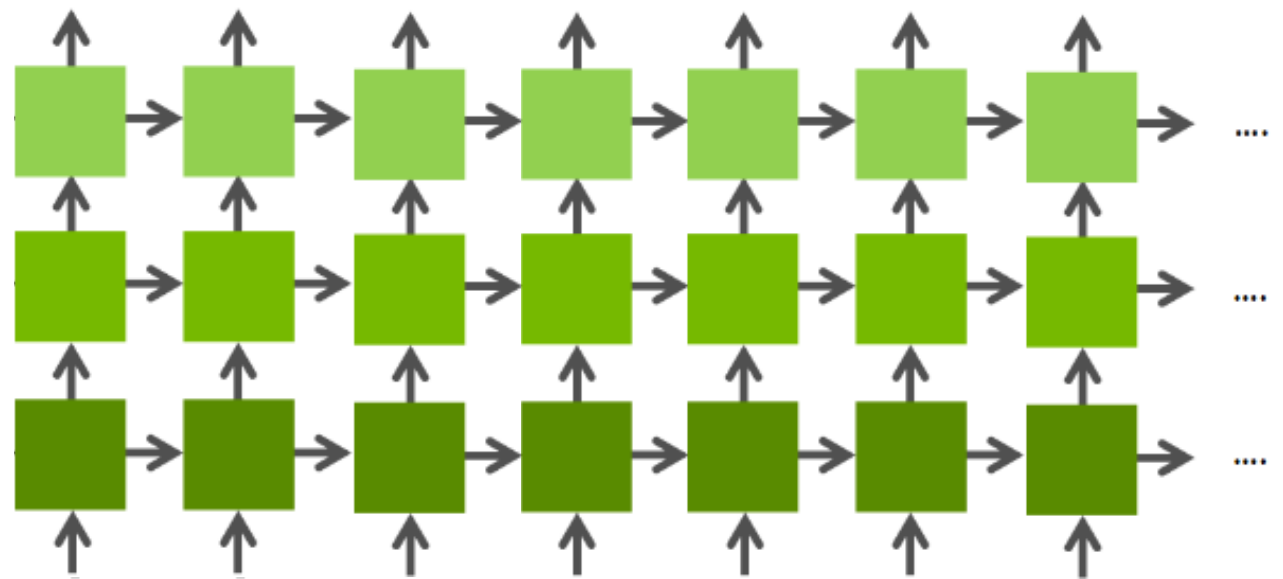
Multi-layer LSTM



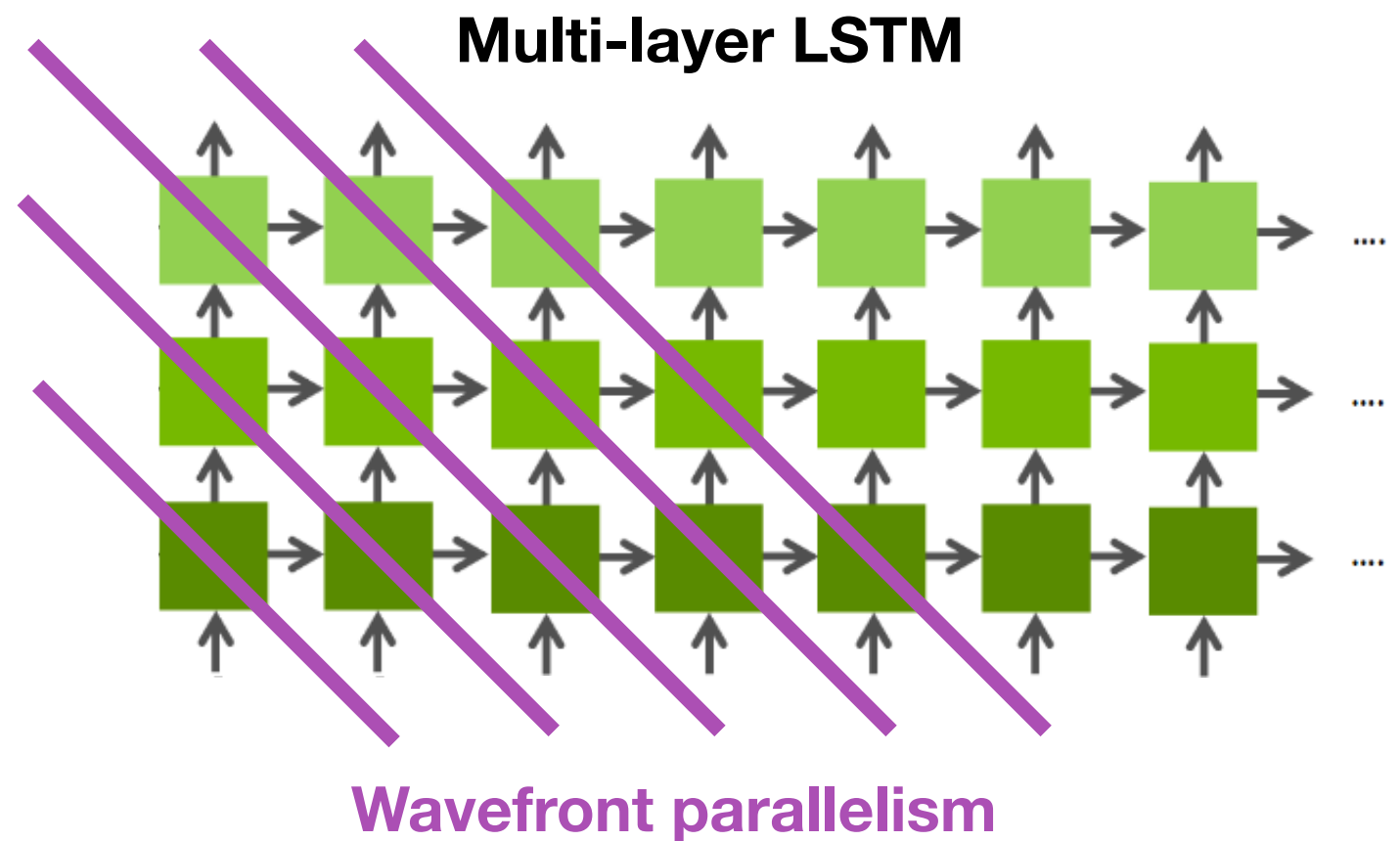
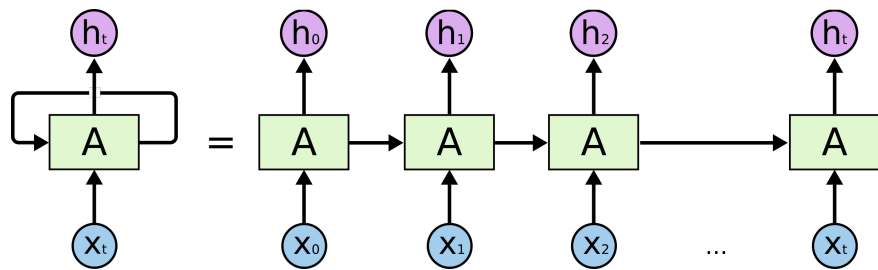
Recurrent Neural Networks



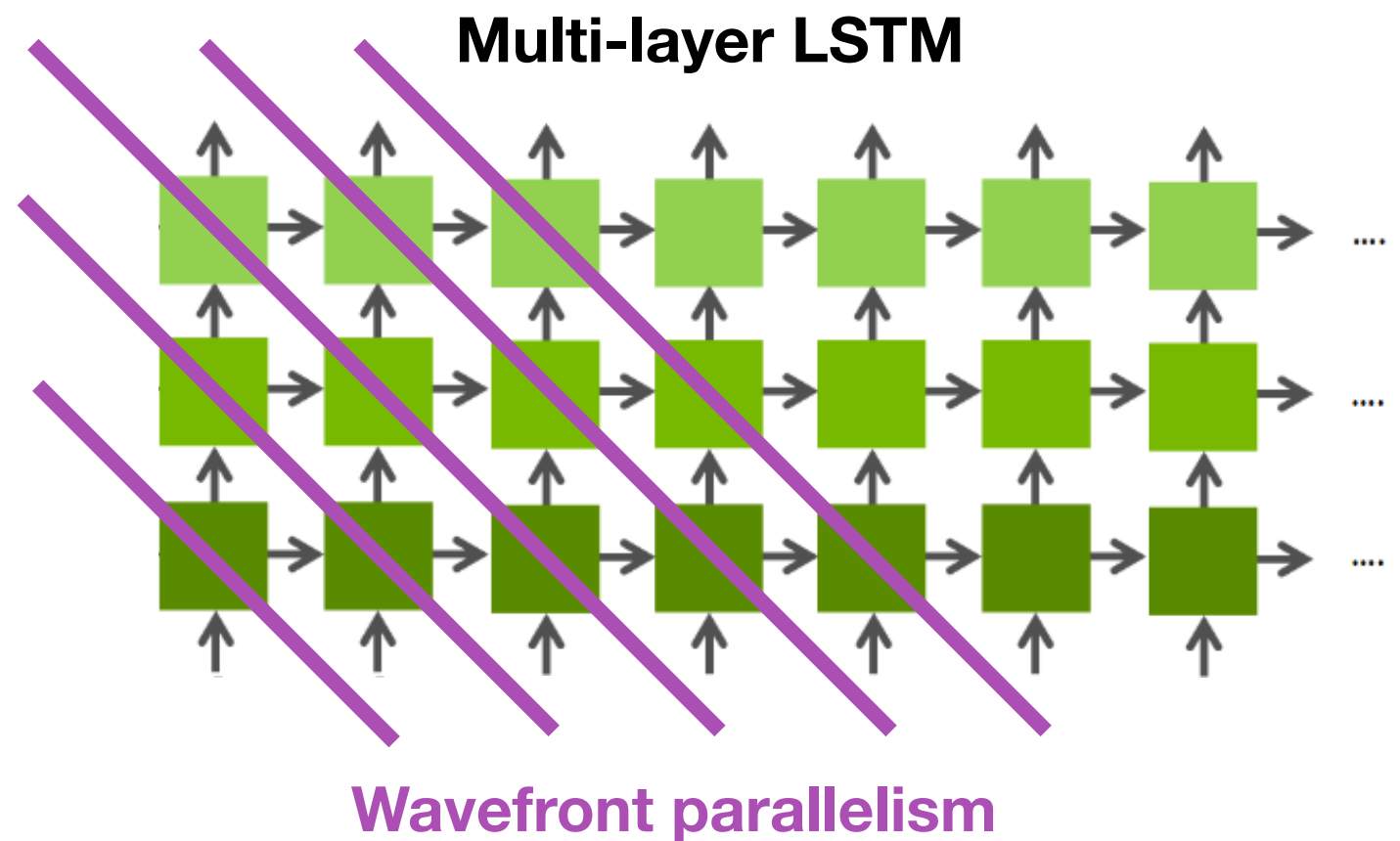
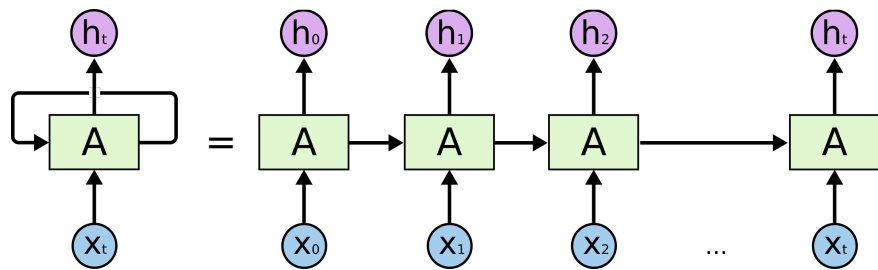
Multi-layer LSTM



Recurrent Neural Networks

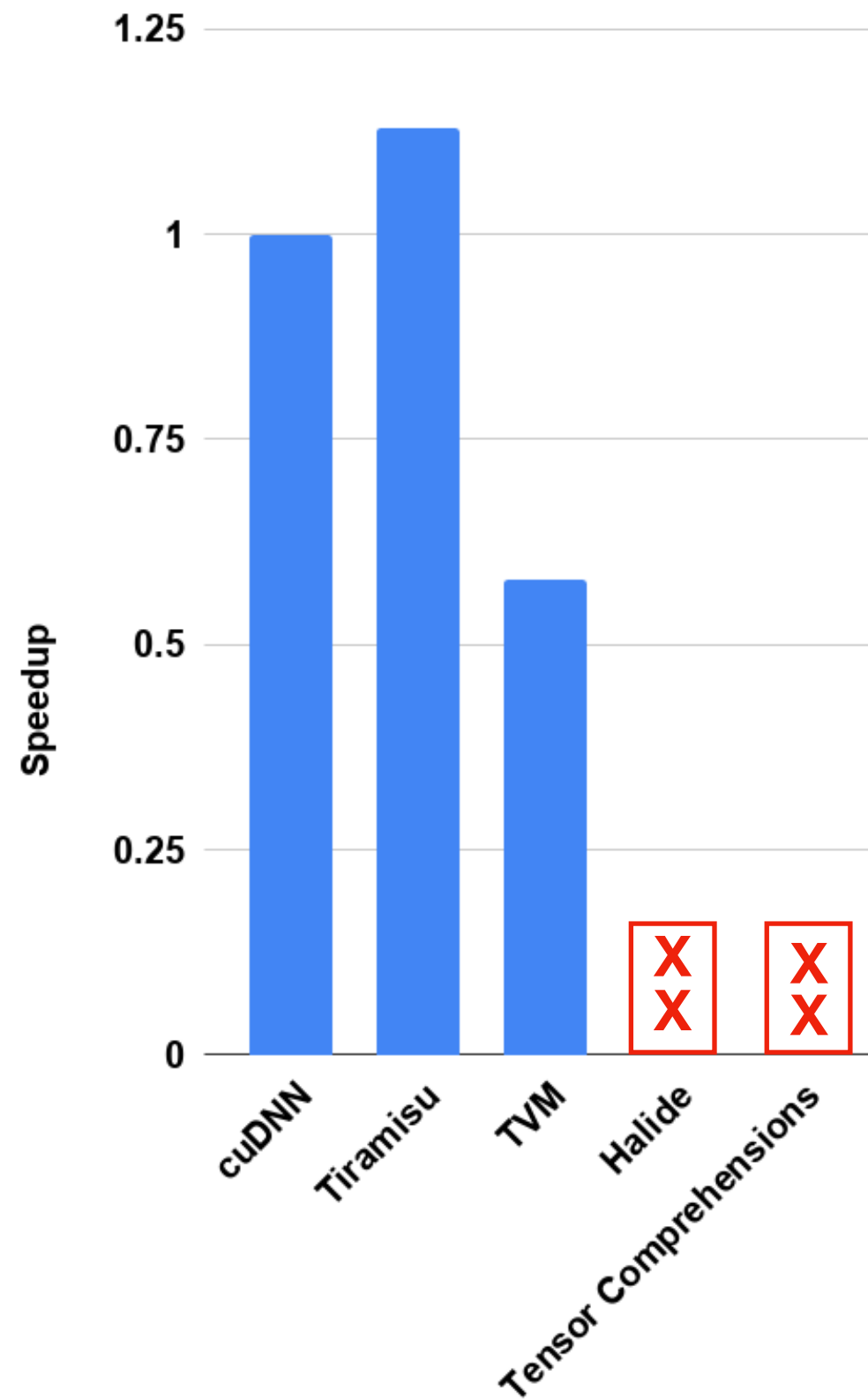


Recurrent Neural Networks

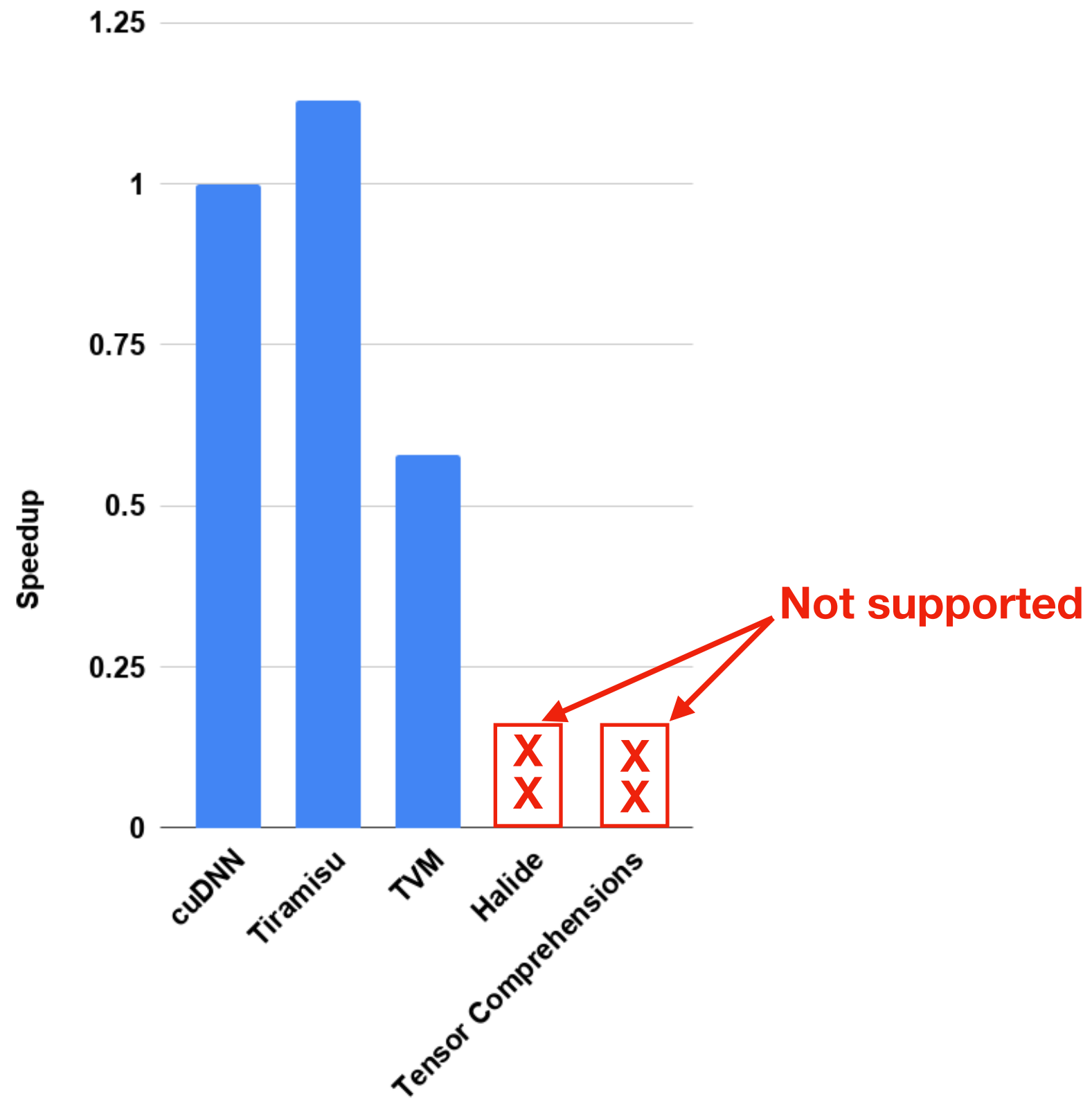


We use loop skewing

Multi-layer LSTM (GPU)



Multi-layer LSTM (GPU)



Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Outline

- **Phase I**

- Tiramisu Overview
- Example
- RNNs

- Sparse DNNs

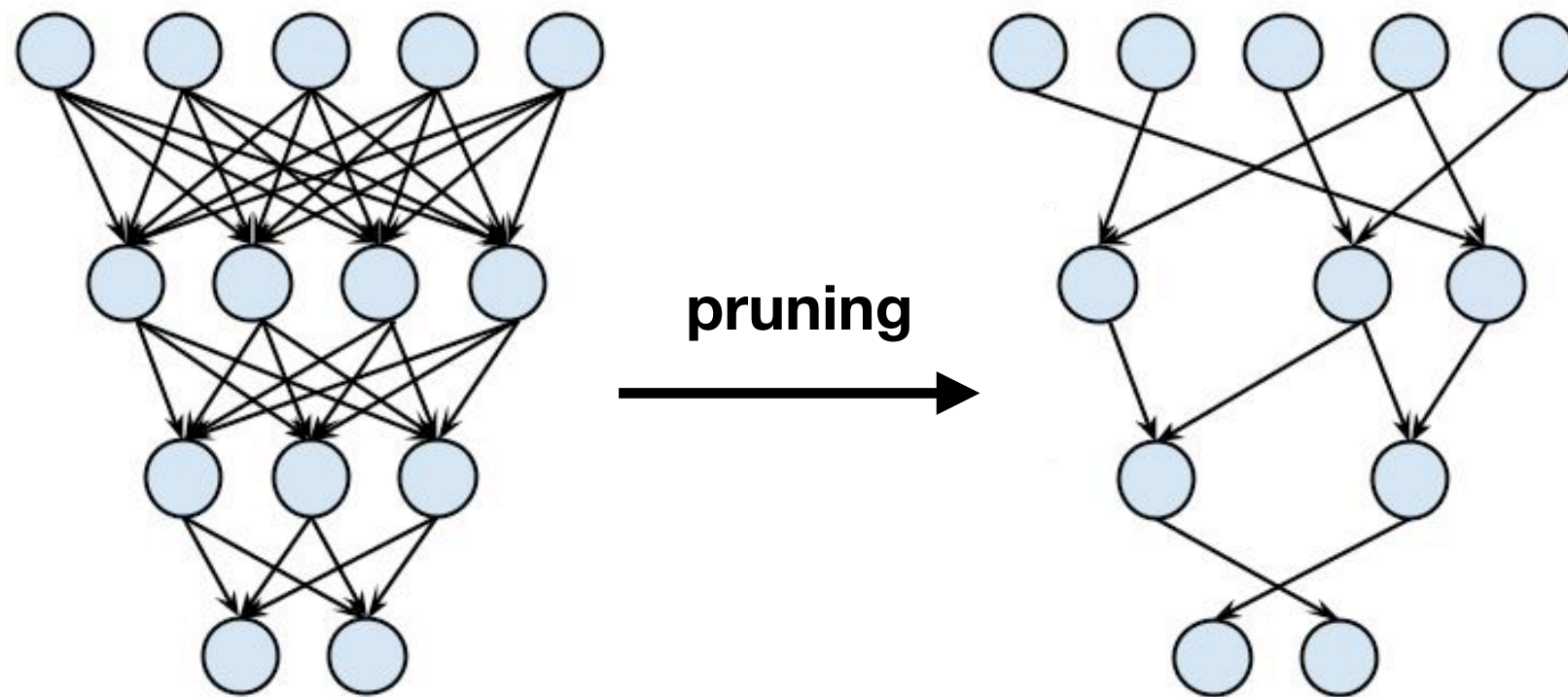
- **Phase II**

- Automatic Optimization

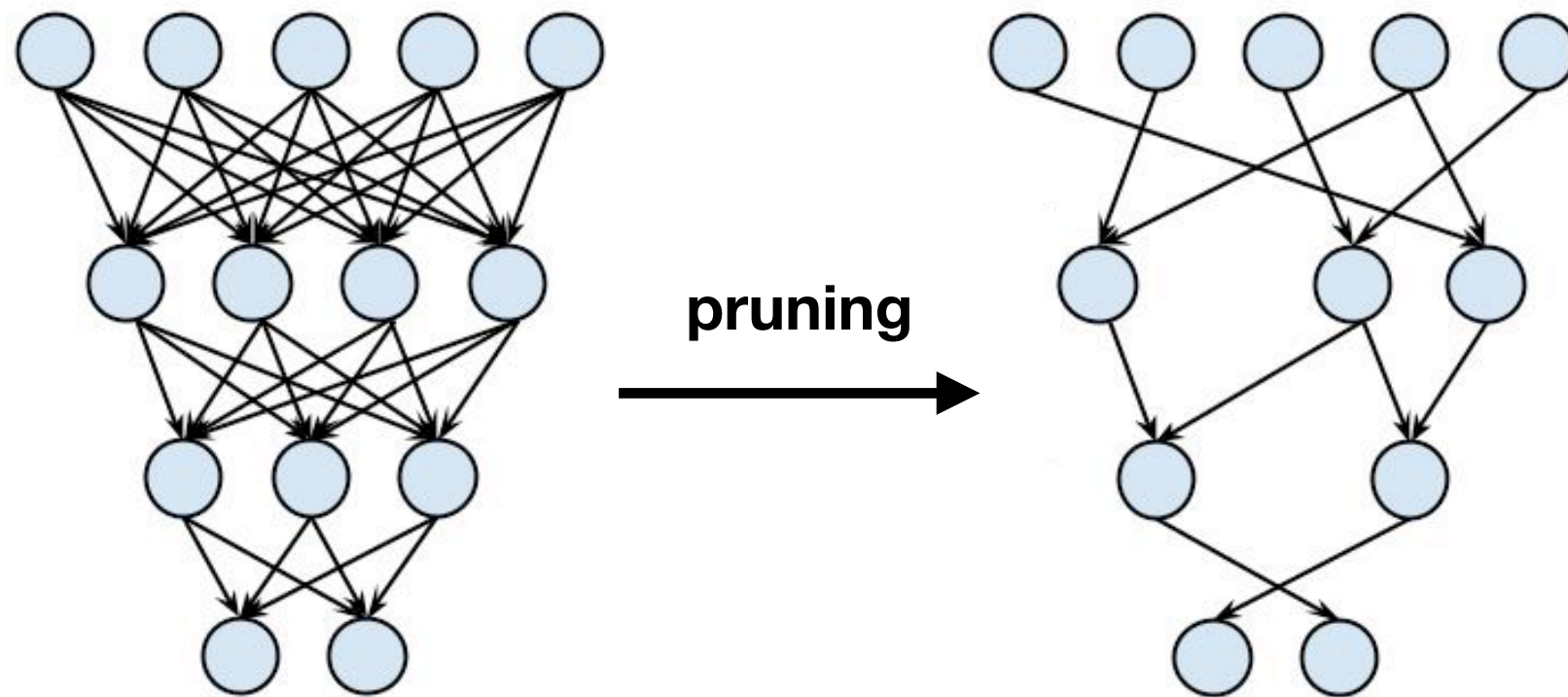
Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Exploiting Sparsity in DNN

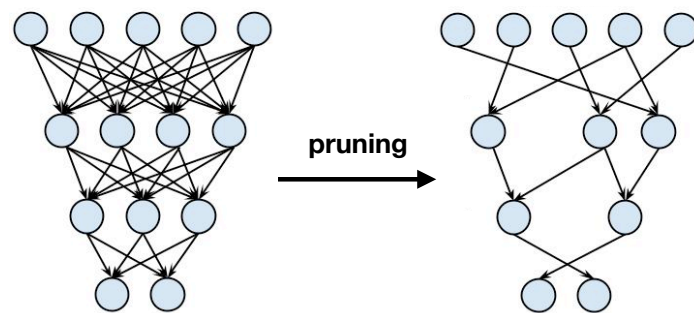


Exploiting Sparsity in DNN



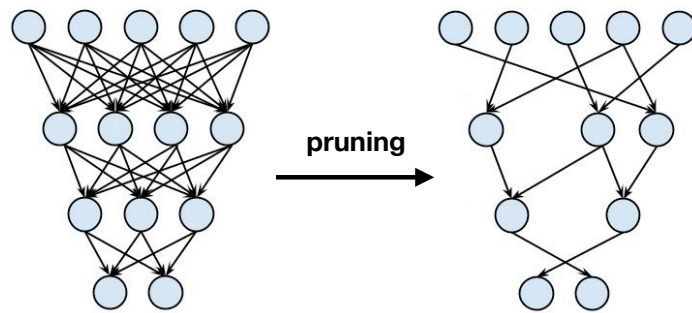
Reduce weights by **orders of magnitude**

Exploiting Sparsity in DNN



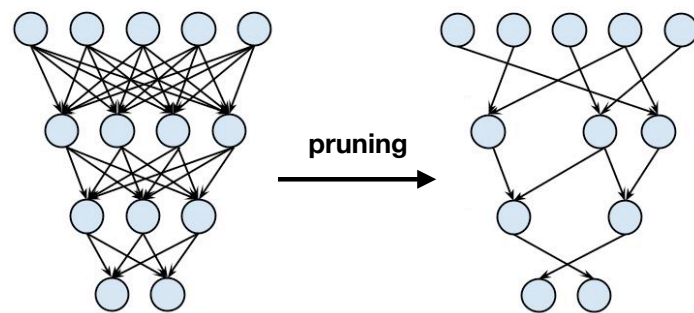
Exploiting Sparsity in DNN

Compiling sparse code is challenging!



Exploiting Sparsity in DNN

Compiling sparse code is challenging!

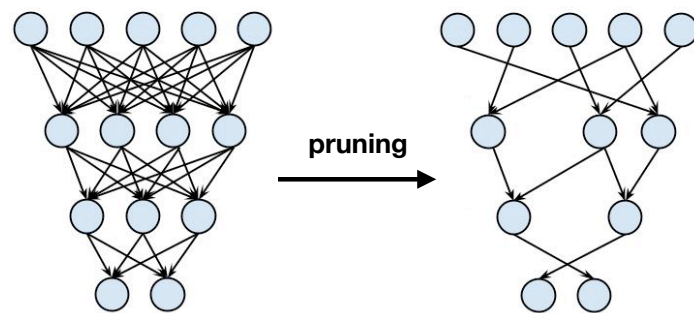


Irregular loop bounds

```
.....  
for j in (rowptr[n], rowptr[n+1])  
.....
```

Exploiting Sparsity in DNN

Compiling sparse code is challenging!

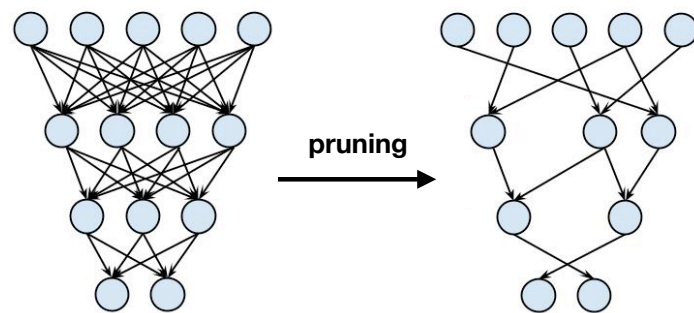


Irregular loop bounds

```
.....  
for j in (rowptr[n], rowptr[n+1])  
.....
```

Exploiting Sparsity in DNN

Compiling sparse code is challenging!



Irregular loop bounds

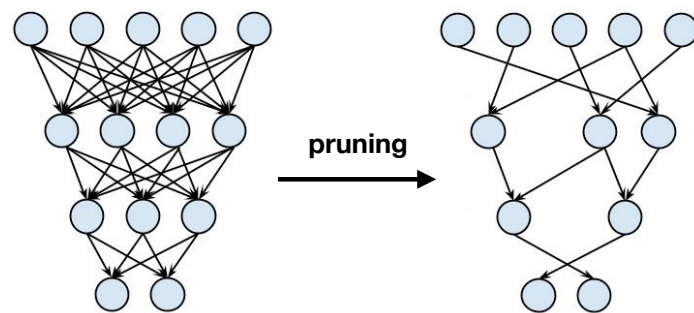
```
.....  
for j in (rowptr[n], rowptr[n+1])  
.....
```

Indirect array accesses

```
.....  
out[n][y][x] = ... + in[colidx[j]]  
.....
```

Exploiting Sparsity in DNN

Compiling sparse code is challenging!



Irregular loop bounds

```
.....  
for j in (rowptr[n], rowptr[n+1])  
.....
```

Indirect array accesses

```
.....  
out[n][y][x] = ... + in[colidx[j]]  
.....
```

Sparsity in DNNs

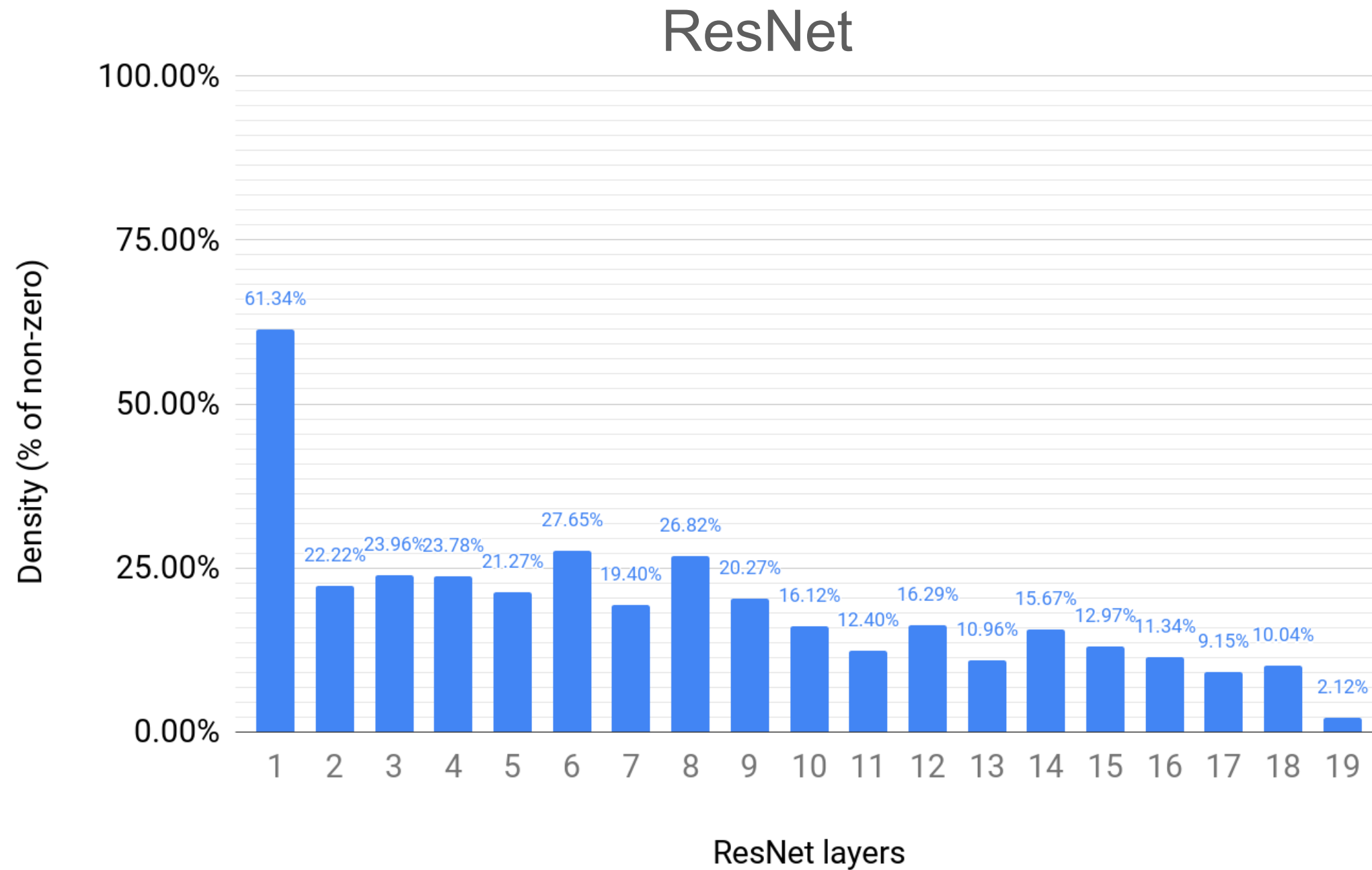
Median sparsity

ResNet: 20%

VGG: 2%

Sparsity in DNNs

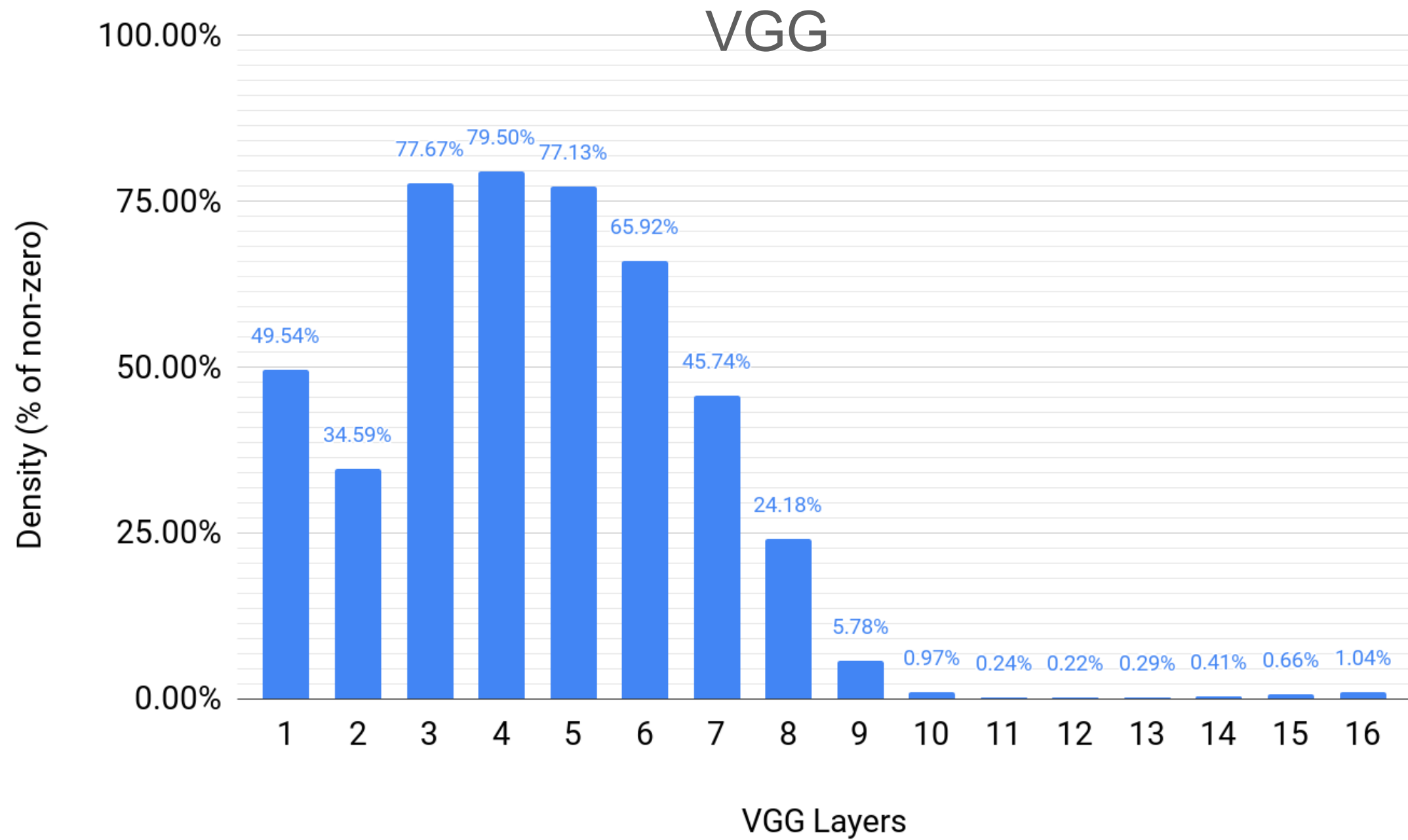
Sparsity in DNNs



[Frankle & Carbin, 2019]

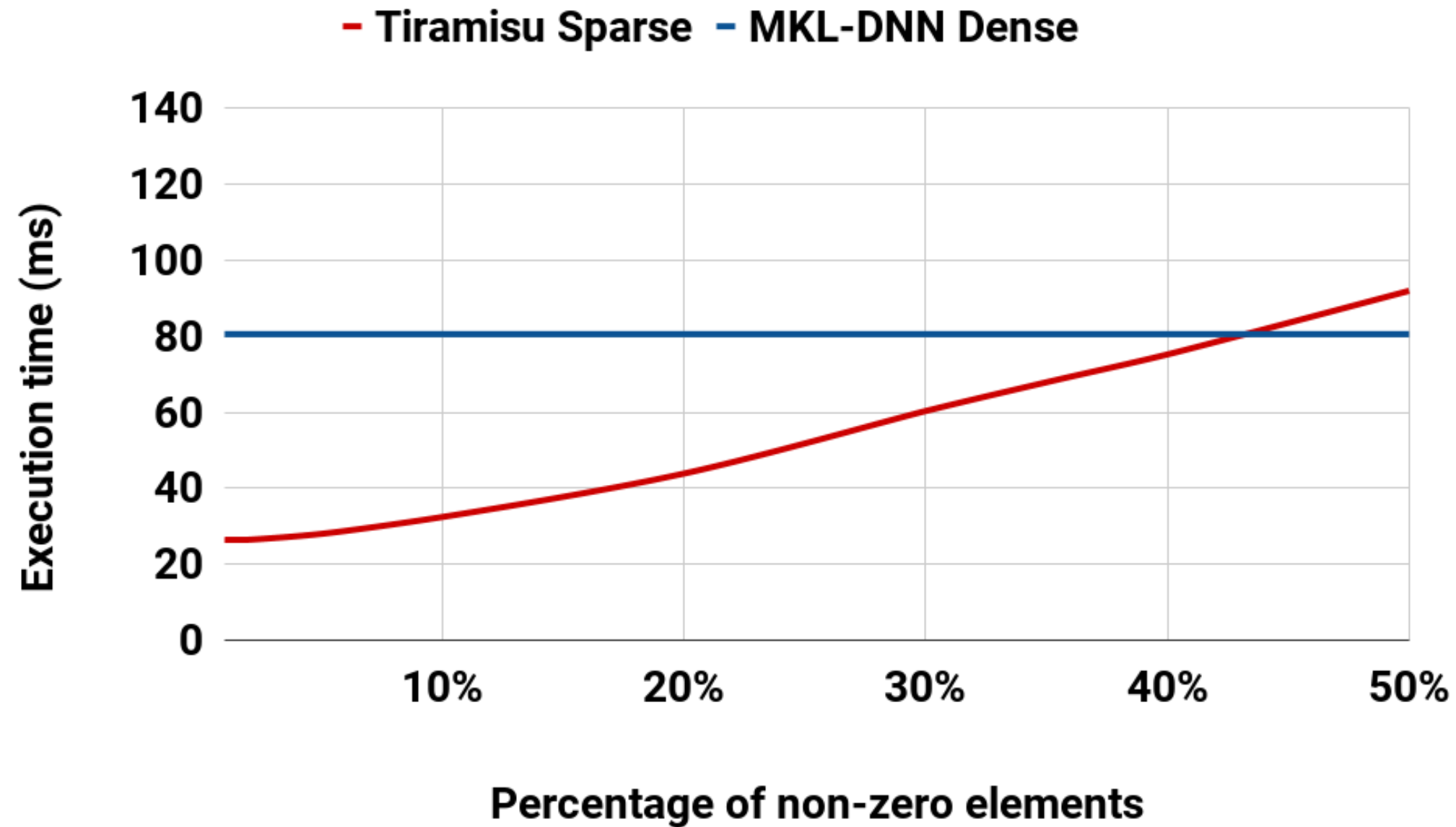
Sparsity in DNNs

Sparsity in DNNs

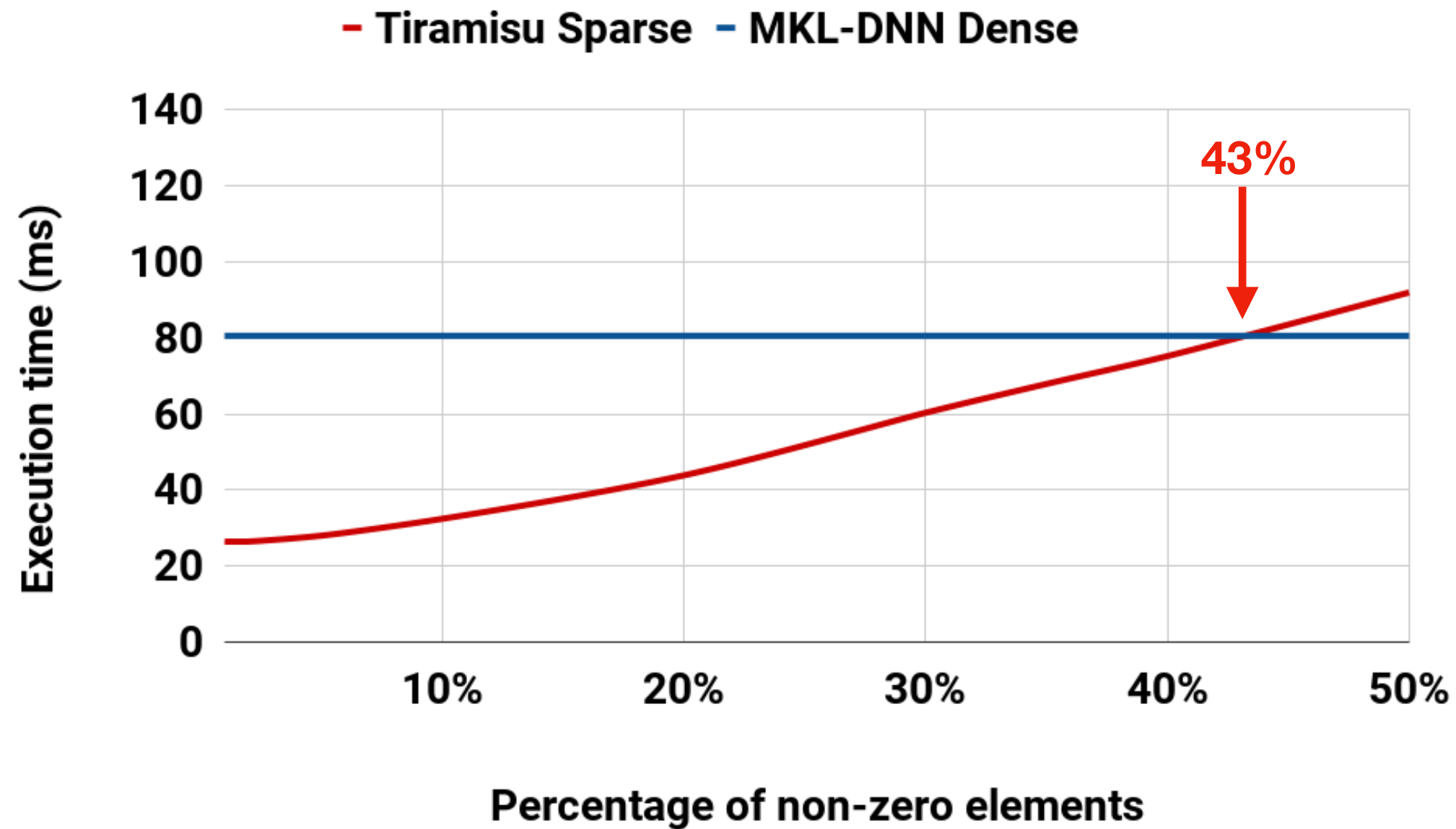


[Frankle & Carbin, 2019]

Break-even Density Level



Break-even Density Level



Optimizing Sparse DNNs

Optimizing Sparse DNNs

- Focus on weight sparsity

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights

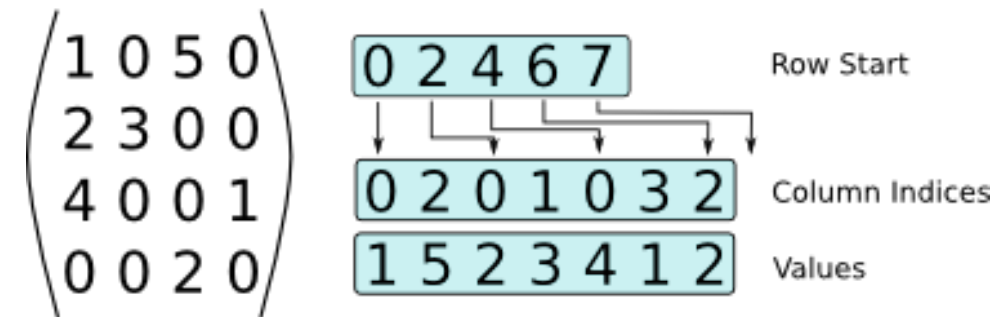
Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

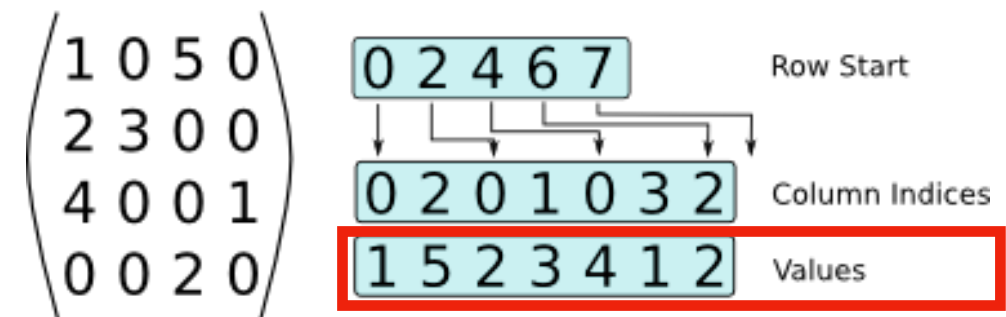
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

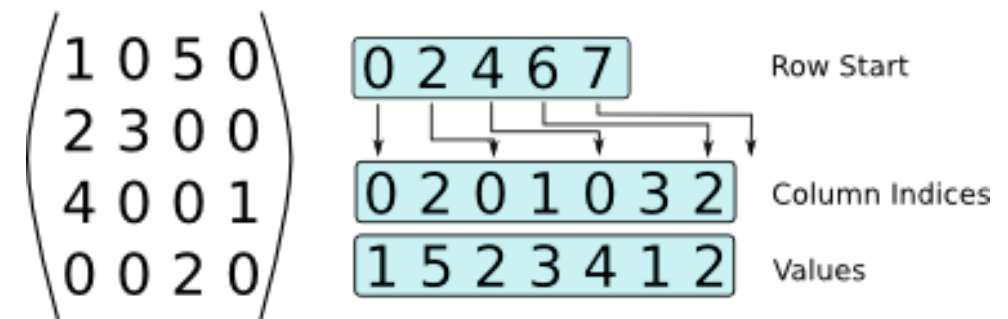
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

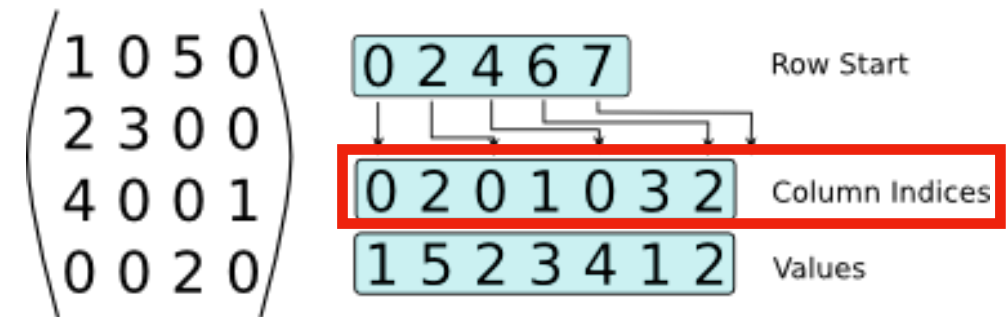
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

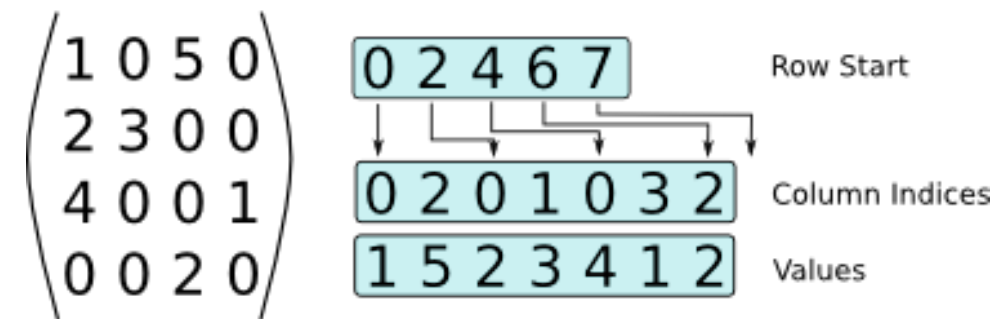
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

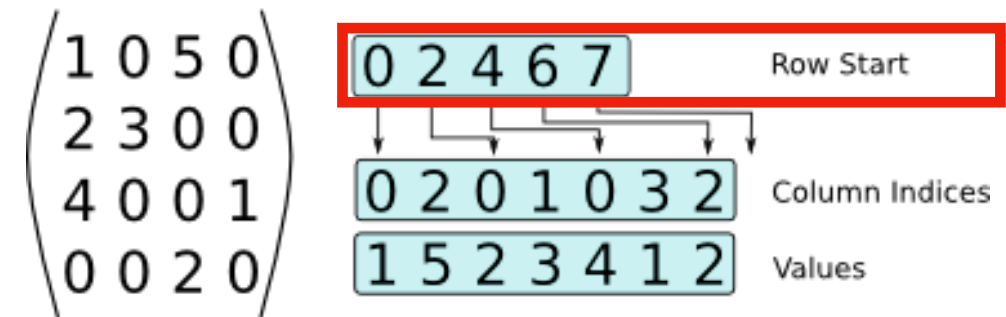
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

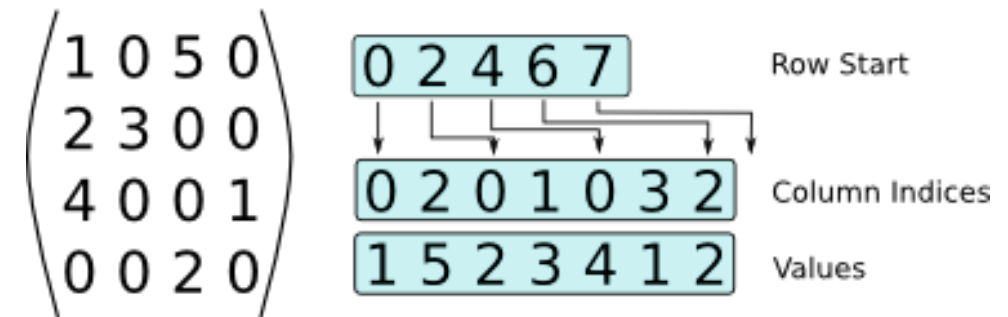
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

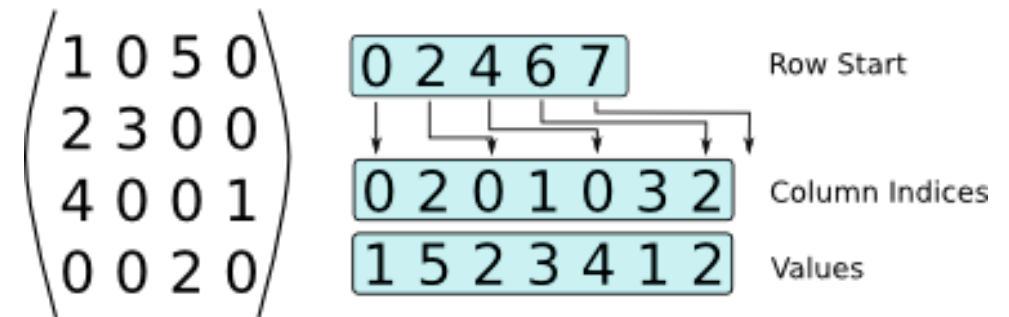
Example: CSR Storage



Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



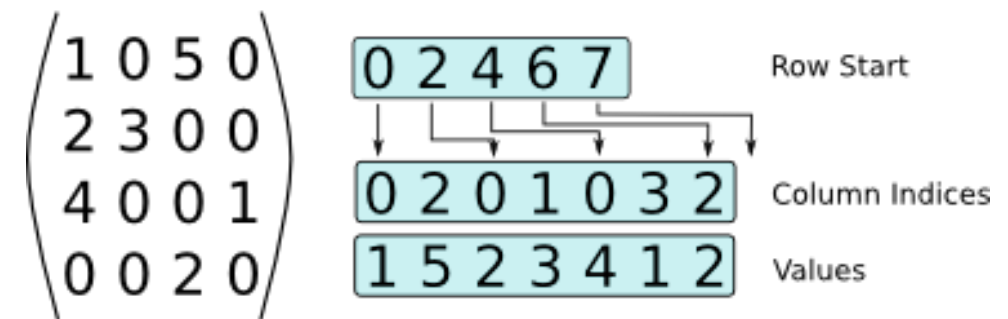
```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

[Park et al, 2016]

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



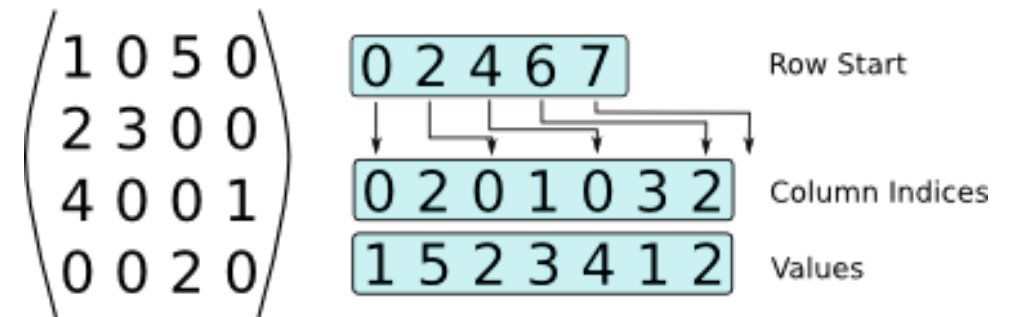
```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

[Park et al, 2016]

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



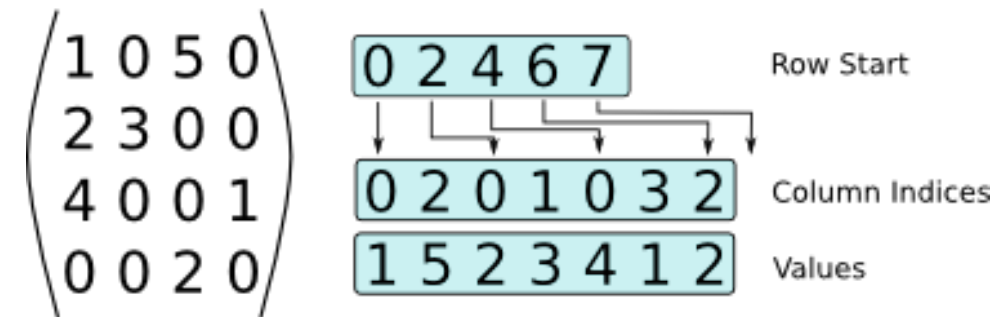
```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

[Park et al, 2016]

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

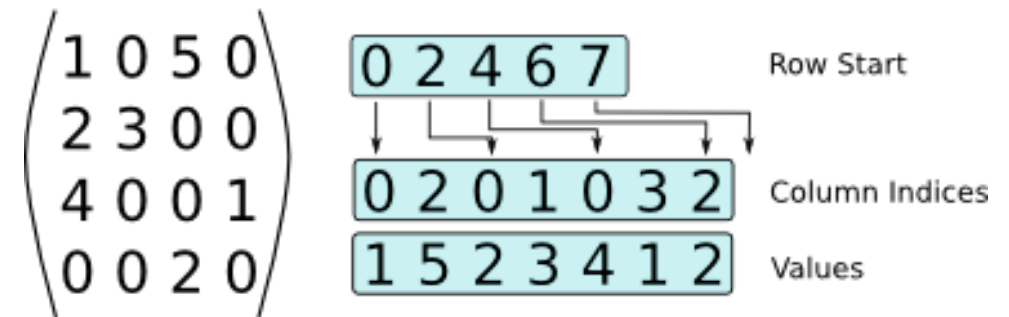
Get a non-zero weight

[Park et al, 2016]

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



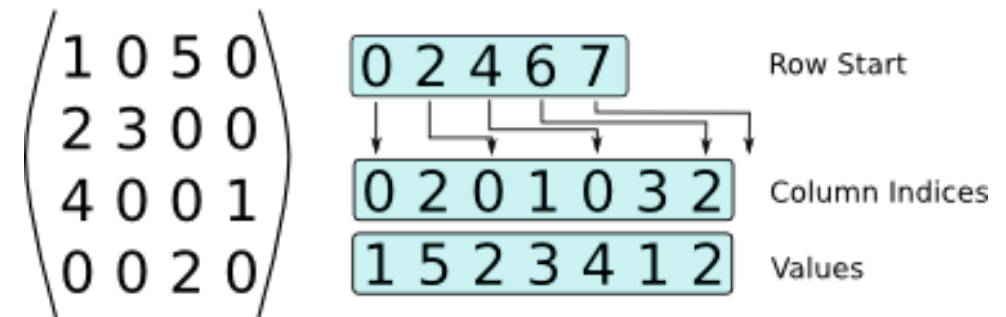
```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

[Park et al, 2016]

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



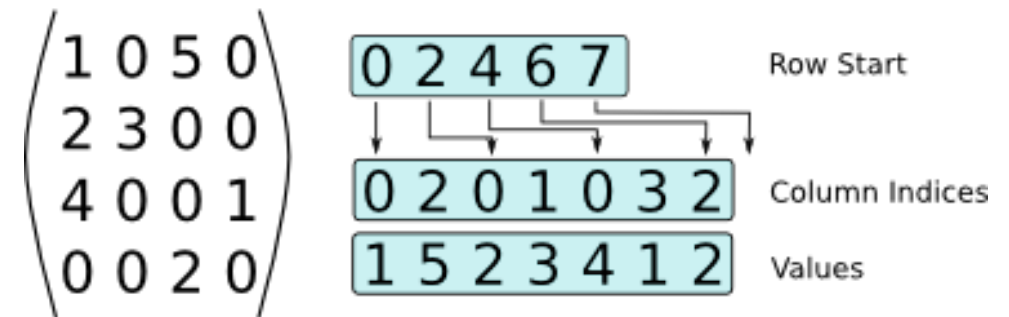
```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

[Park et al, 2016]

Optimizing Sparse DNNs

- Focus on weight sparsity
- CSR format for the weights
- Use direct convolution

Example: CSR Storage



```
for each output channel n
  for j in (W.rowptr[n], W.rowptr[n+1])
  {
    off = W.colidx[j];
    coeff = W.value[j];
    for (int y = 0; y < H_OUT; ++y)
      for (int x = 0; x < W_OUT; ++x)
        out[n][y][x] += coeff*in[y*W_OUT+x+off]
  }
```

[Park et al, 2016]

Extensions to the Polyhedral Model

Extensions to the Polyhedral Model

- Indirect array accesses

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```


Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```


Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

```
for (int y = 0; y < H_OUT; ++y)
  if (A[y])
    out[y] = ...
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

```
for (int y = 0; y < H_OUT; ++y)
  if (A[y])
    out[y] = ...
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

```
for (int y = 0; y < H_OUT; ++y)
  if (A[y])
    out[y] = ...
```


Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

```
for (int y = 0; y < H_OUT; ++y)
  if (A[y])
    out[y] = ...
```

```
for (int y = 0; y < H_OUT; ++y)
  p = A[y]
  if (p)
    out[y] = ...
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

```
for (int y = 0; y < H_OUT; ++y)
  if (A[y])
    out[y] = ...
```

```
for (int y = 0; y < H_OUT; ++y)
  p = A[y]
  if (p)
    out[y] = ...
```

Extensions to the Polyhedral Model

- Indirect array accesses

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    out[n][y][x] += ...*in[colidx[j]]
```

```
for (int y = 0; y < H_OUT; ++y)
  for (int x = 0; x < W_OUT; ++x)
    p = colidx[j]
    out[n][y][x] += ...*in[p]
```

- Irregular loop bounds

```
for j in (rowptr[n], rowptr[n+1])
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

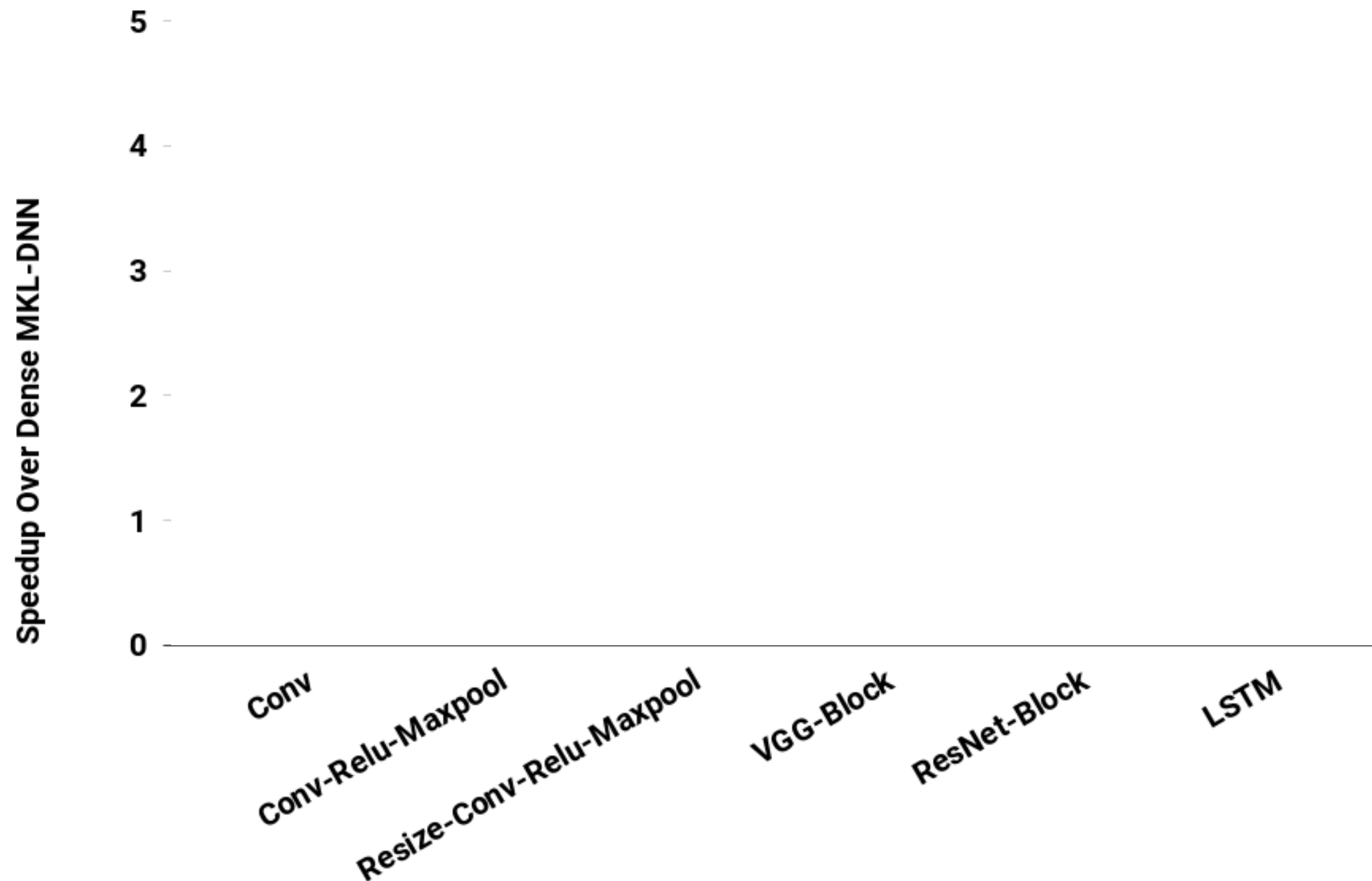
```
b0 = rowptr[n]; b1 = rowptr[n+1];
for j in (b0, b1)
  for (int y = 0; y < H_OUT; ++y)
    for (int x = 0; x < W_OUT; ++x)
      out[n][y][x] += ...*in[...]
```

- Irregular conditionals

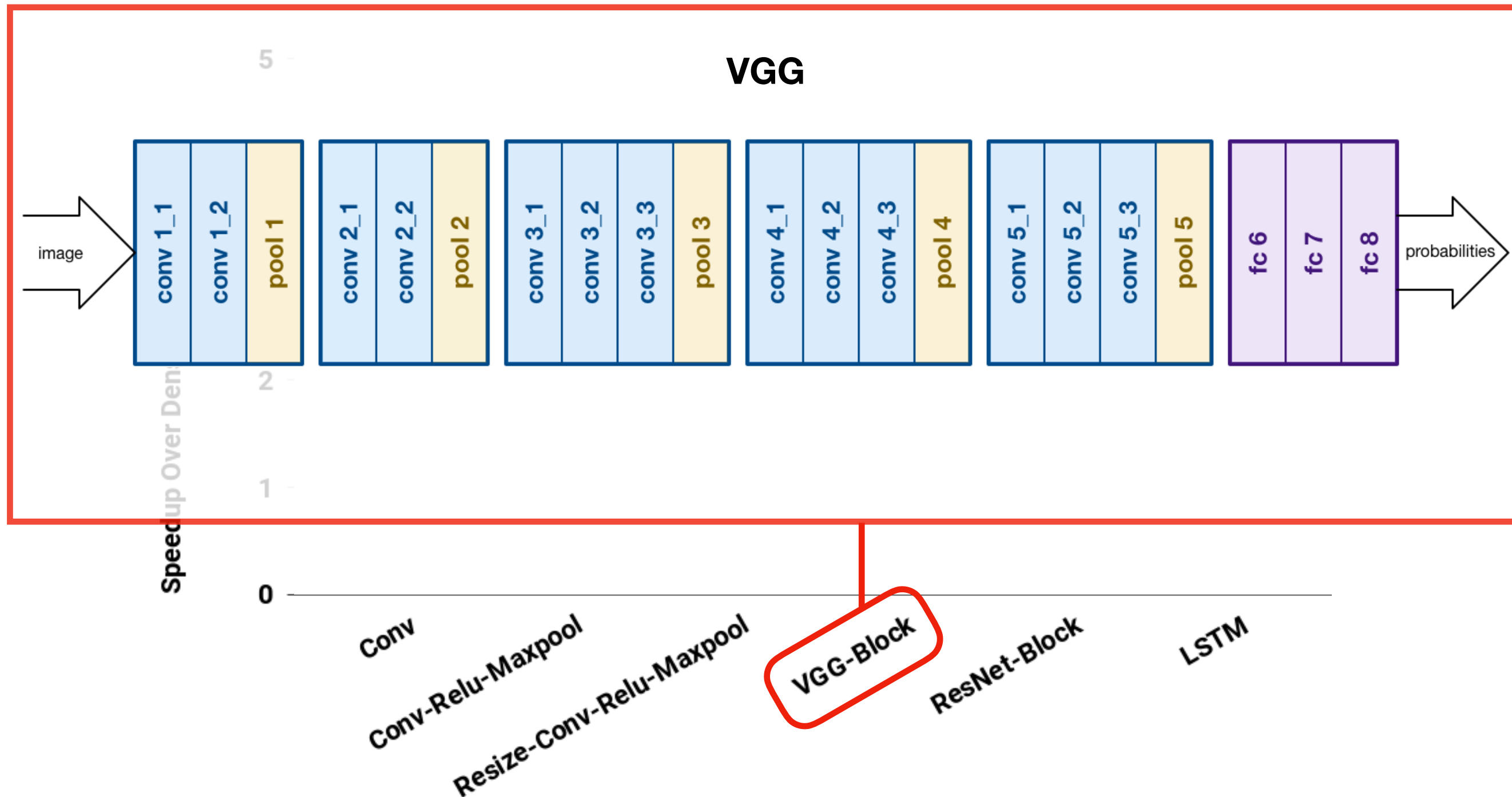
```
for (int y = 0; y < H_OUT; ++y)
  if (A[y])
    out[y] = ...
```

```
for (int y = 0; y < H_OUT; ++y)
  p = A[y]
  if (p)
    out[y] = ...
```

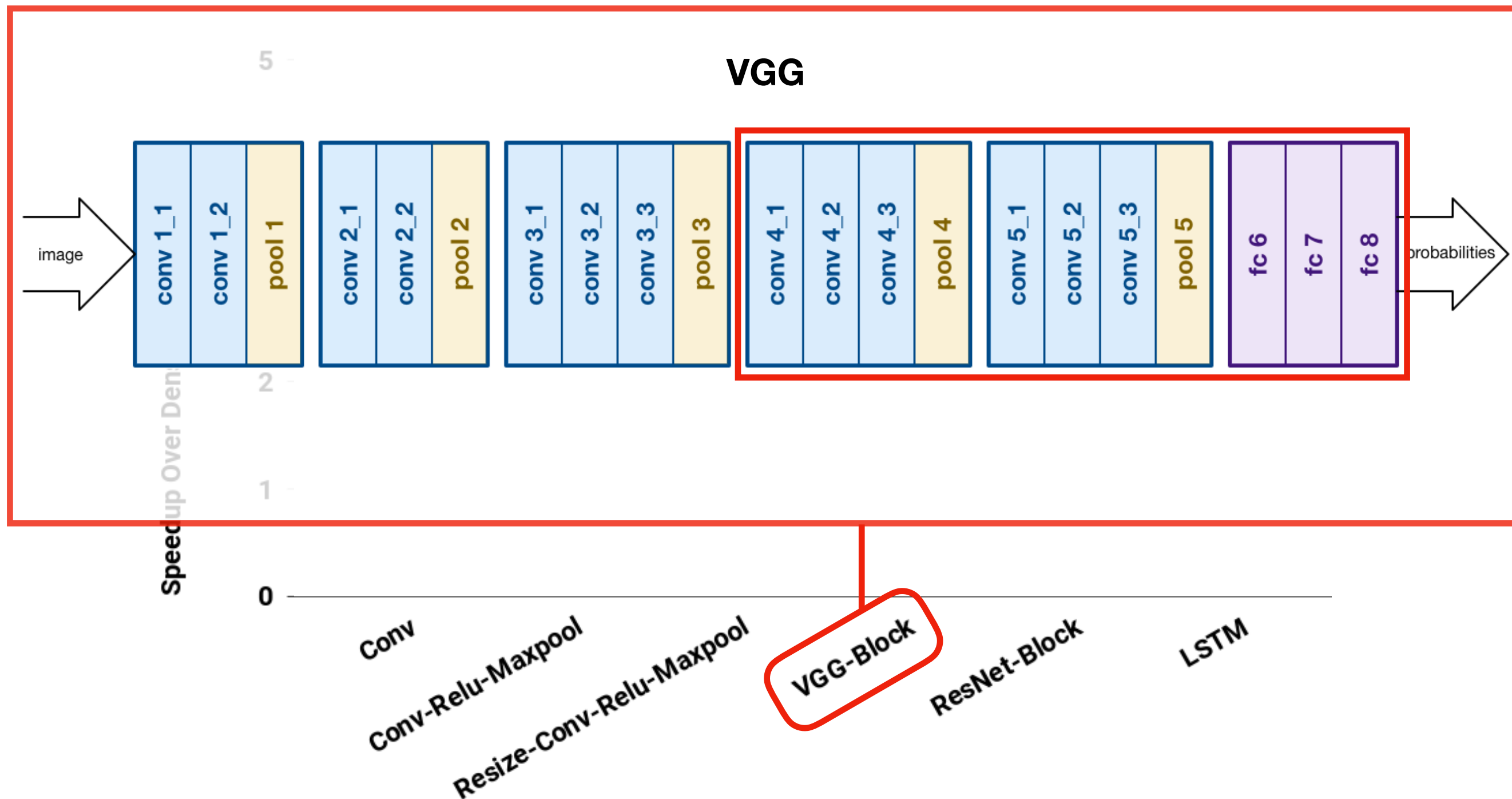
Tiramisu Dense and Sparse (CPU - 4 cores)



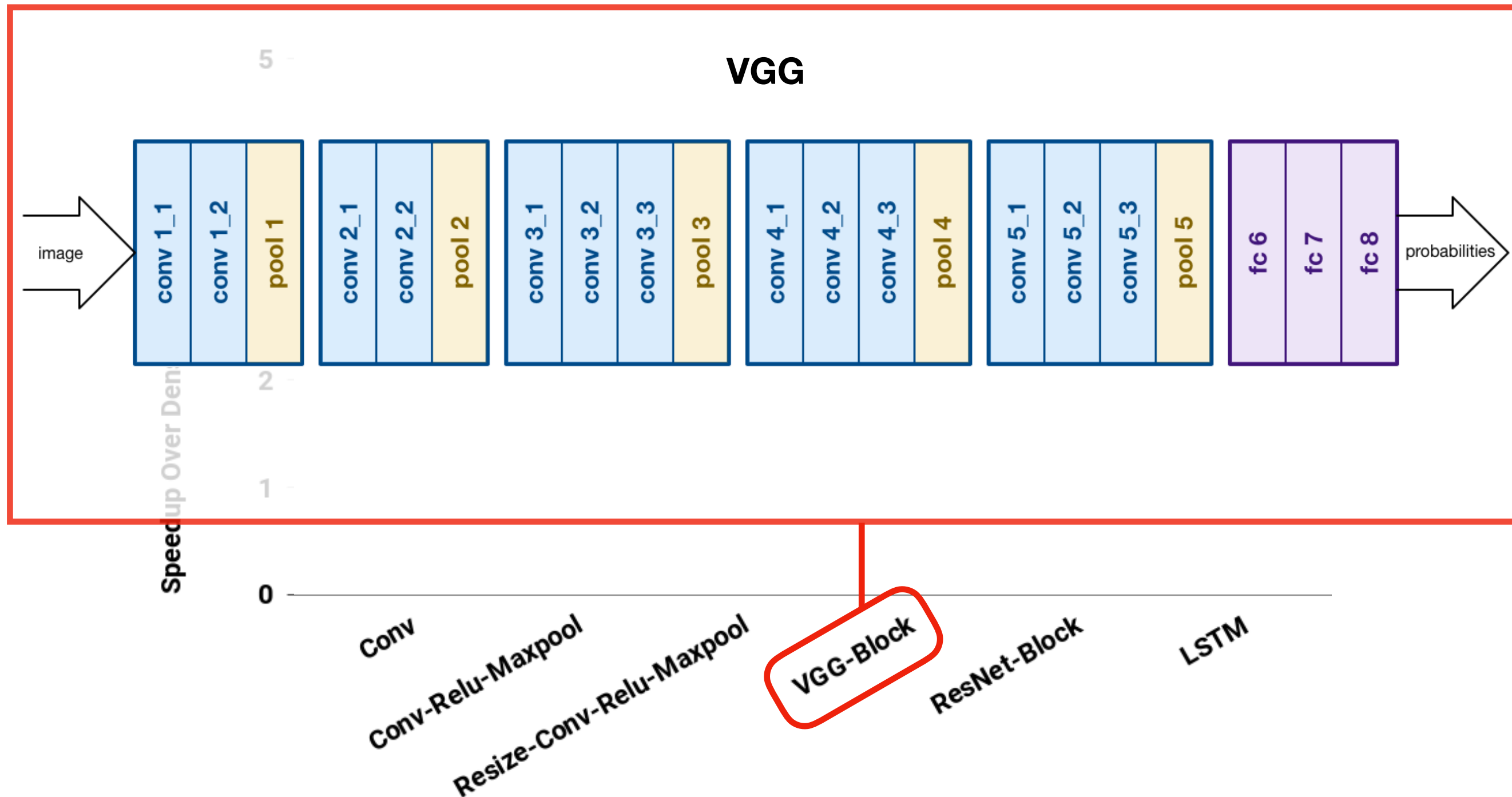
Tiramisu Dense and Sparse (CPU - 4 cores)



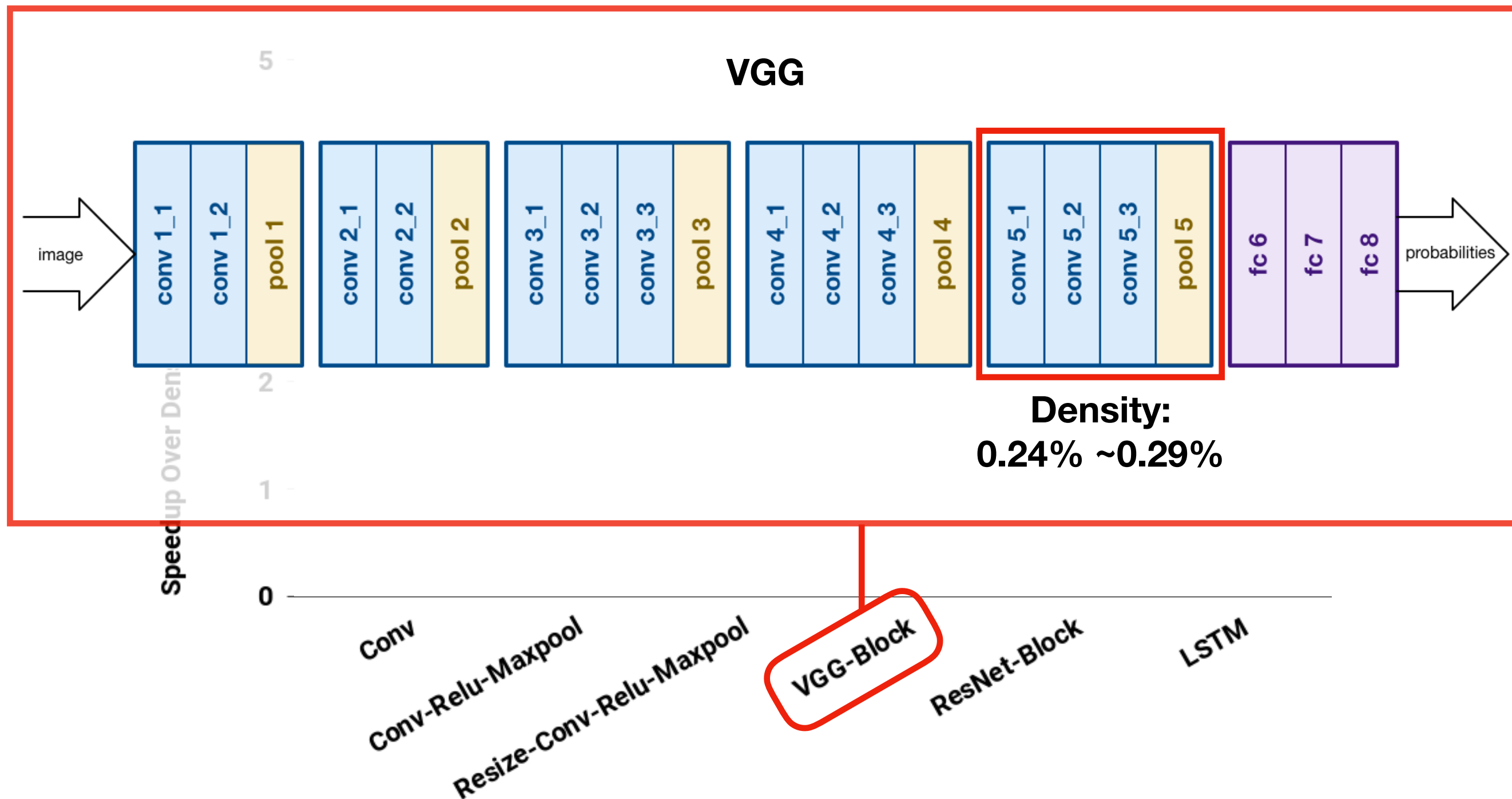
Tiramisu Dense and Sparse (CPU - 4 cores)



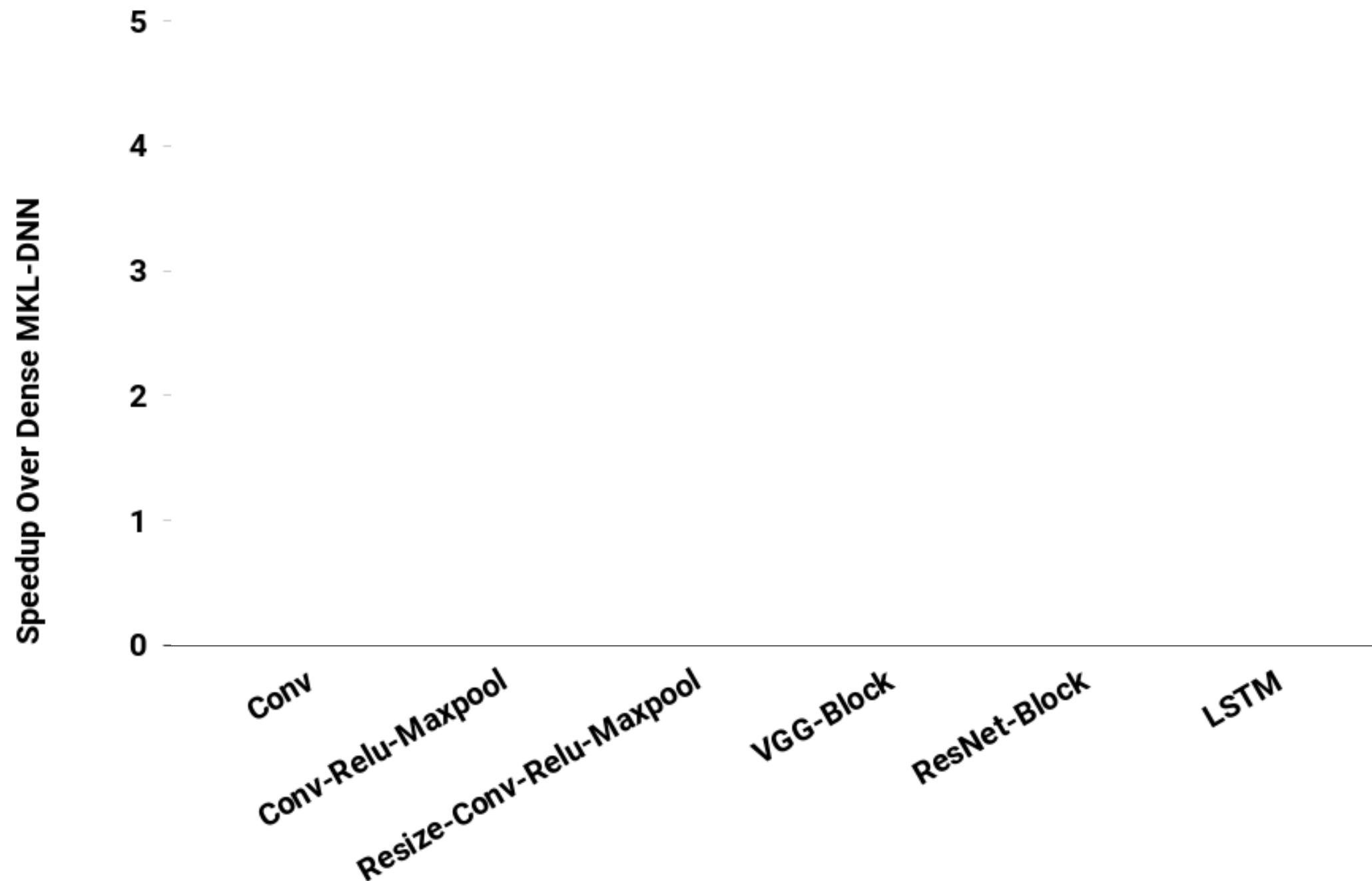
Tiramisu Dense and Sparse (CPU - 4 cores)



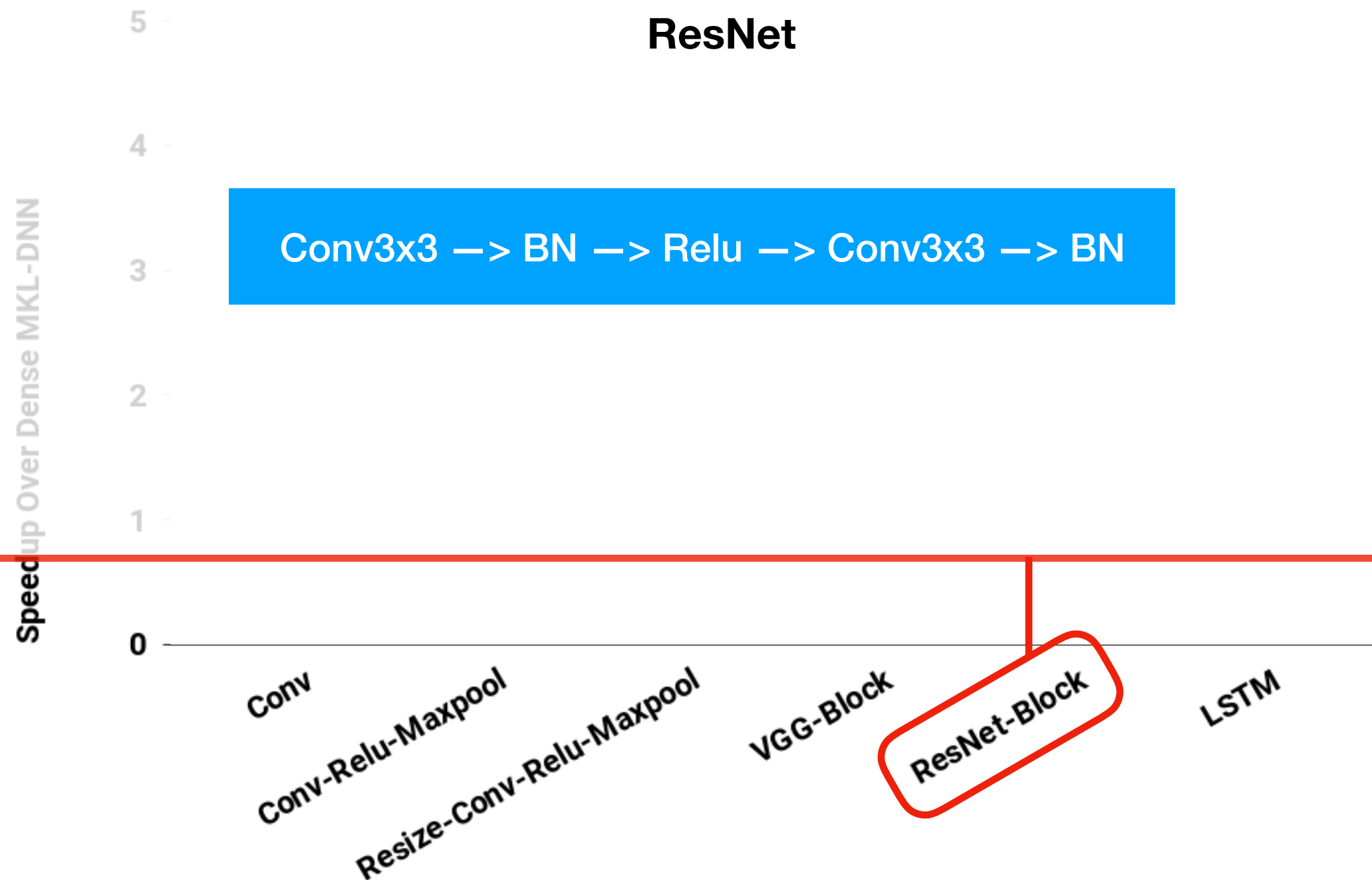
Tiramisu Dense and Sparse (CPU - 4 cores)



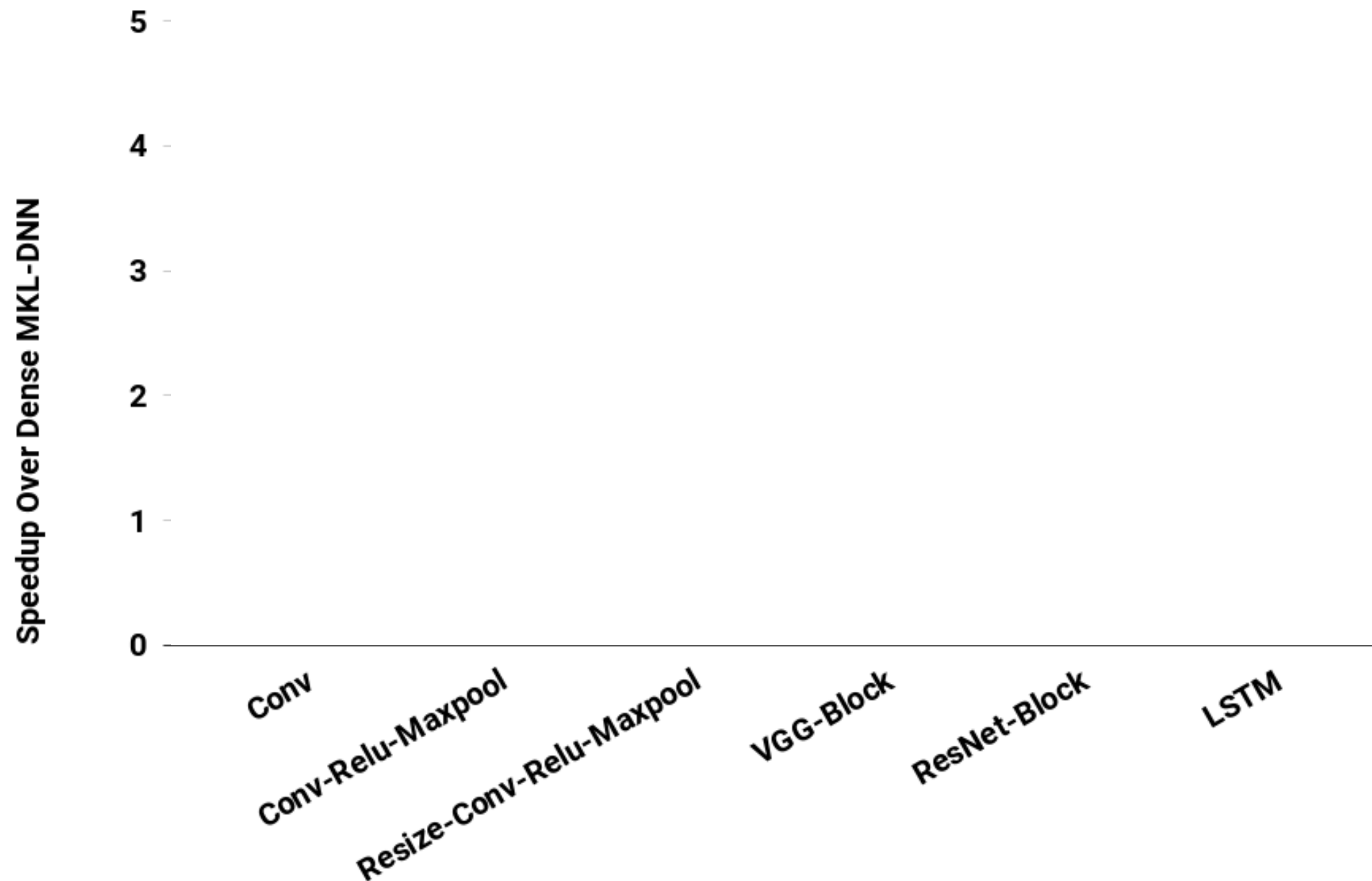
Tiramisu Dense and Sparse (CPU - 4 cores)



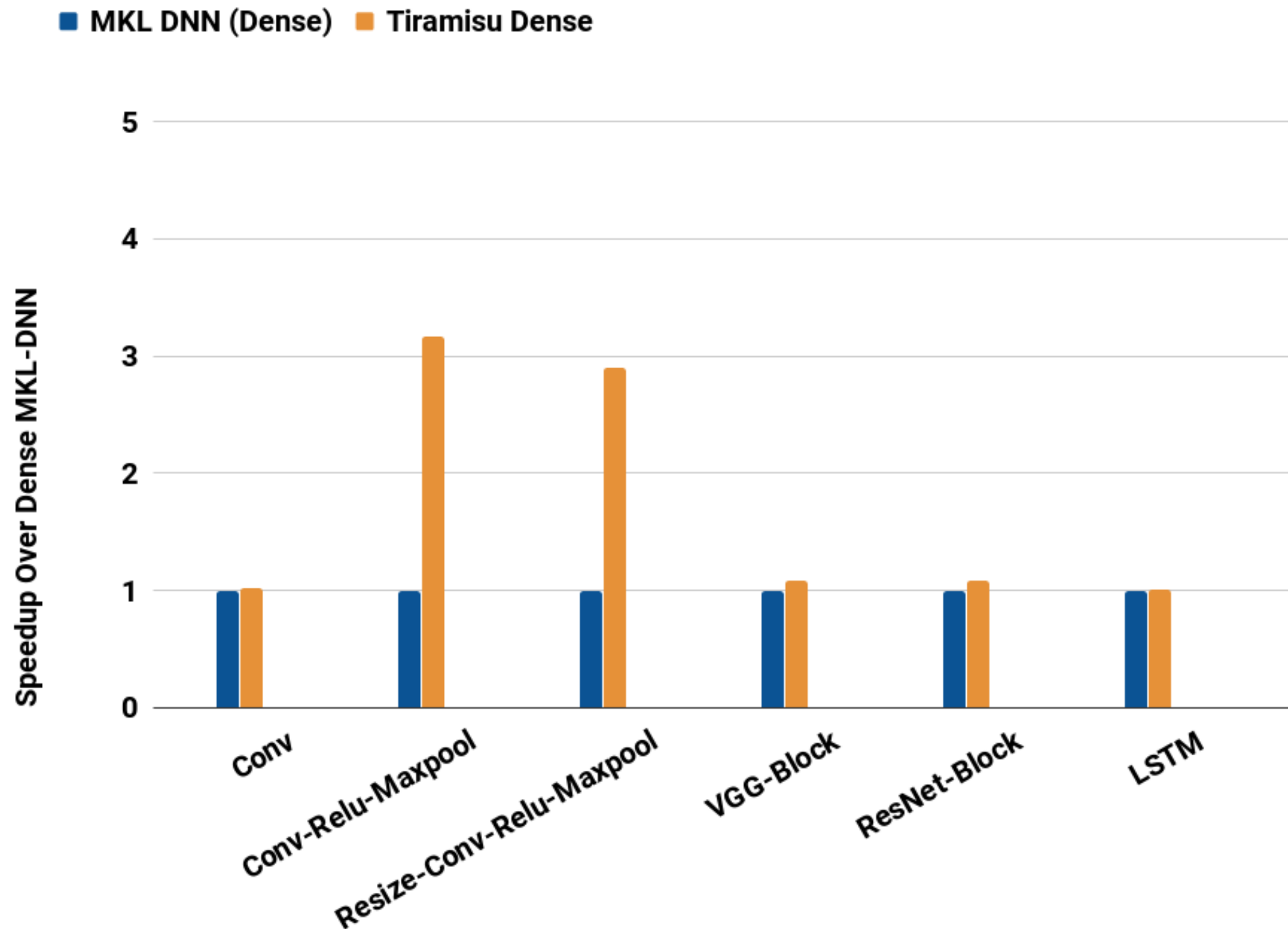
Tiramisu Dense and Sparse (CPU - 4 cores)



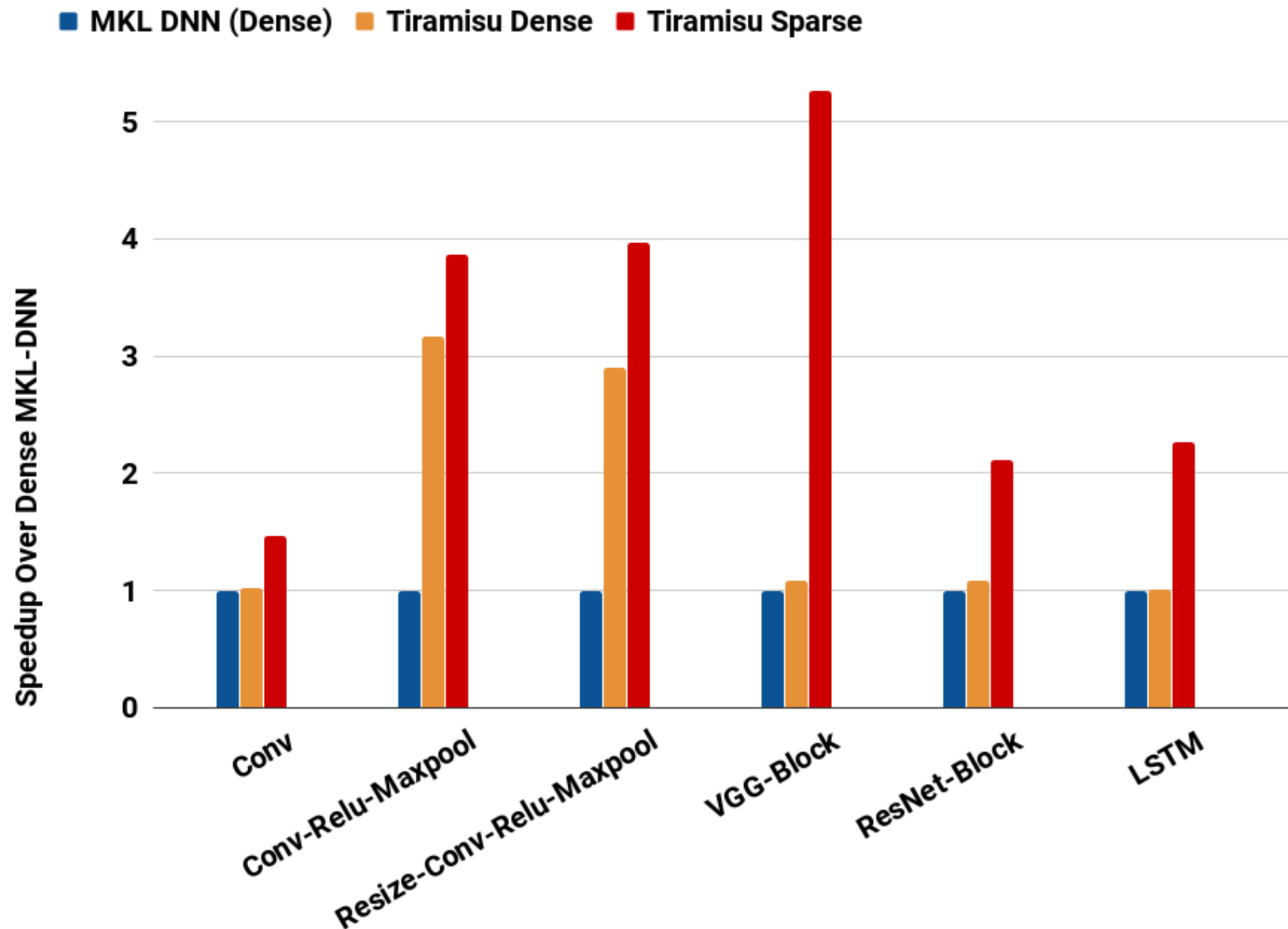
Tiramisu Dense and Sparse (CPU - 4 cores)



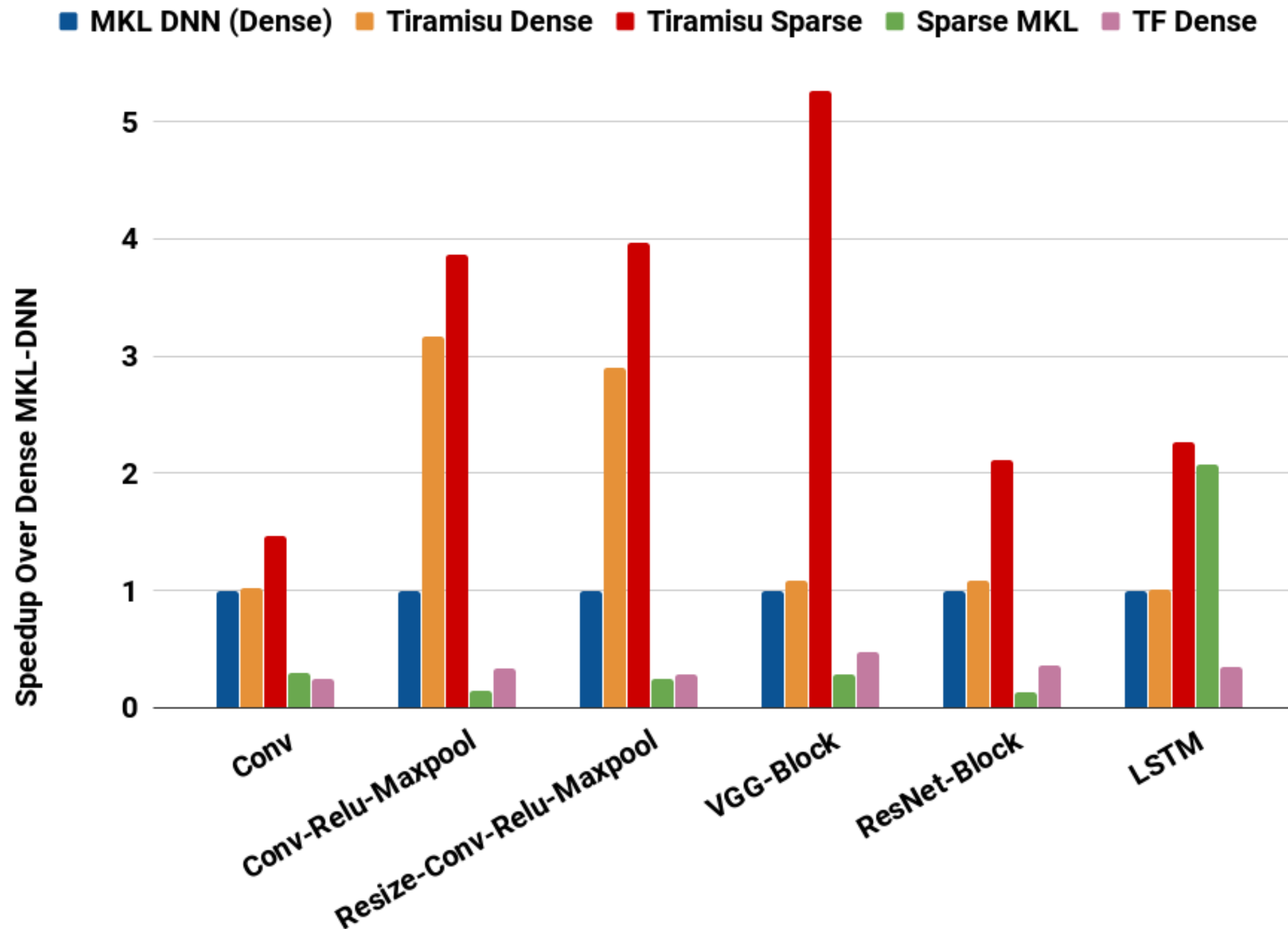
Tiramisu Dense and Sparse (CPU - 4 cores)



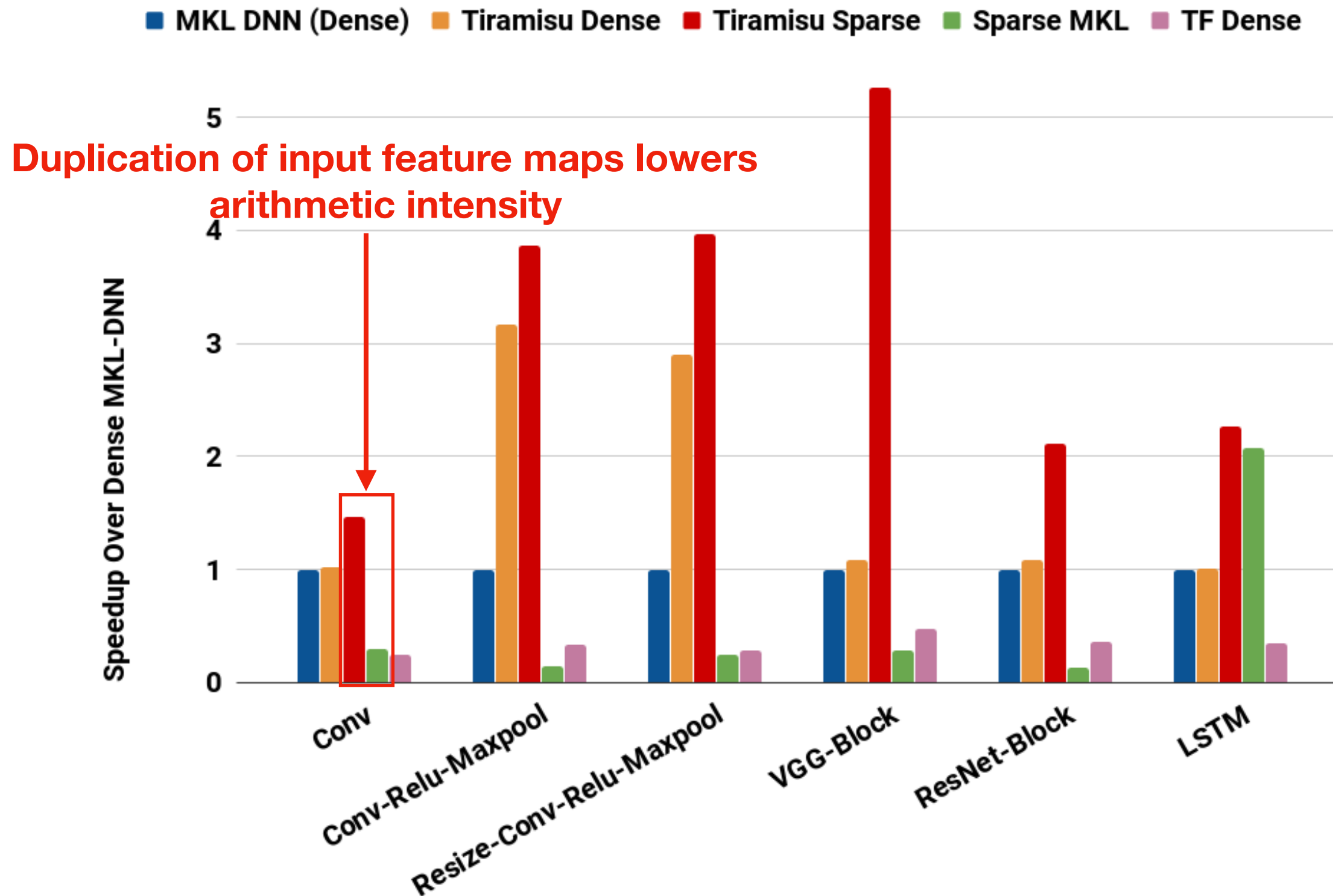
Tiramisu Dense and Sparse (CPU - 4 cores)



Tiramisu Dense and Sparse (CPU - 4 cores)



Tiramisu Dense and Sparse (CPU - 4 cores)



Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Outline

- **Phase I**

- Tiramisu Overview
- Example
- RNNs
- Sparse DNNs

- **Phase II**

- Automatic Optimization

Outline

- **Phase I**
 - Tiramisu Overview
 - Example
 - RNNs
 - Sparse DNNs
- **Phase II**
 - Automatic Optimization

Automatic Optimization

```
graph TD; A[Automatic Optimization] --> B[Selection]; A --> C[Correctness];
```

Selection

Correctness

Optimization Selection

C Explore: tiling, unrolling

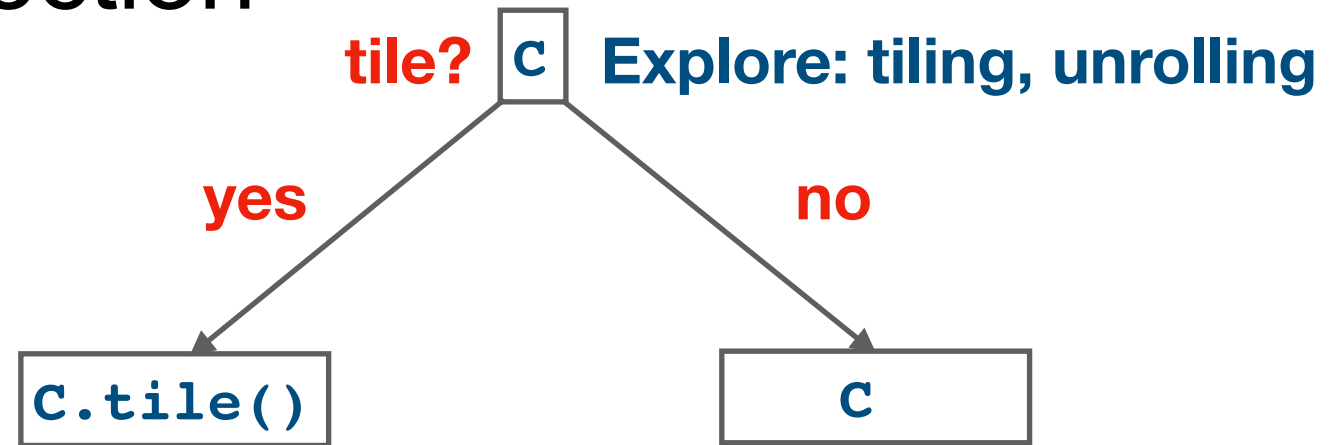
Optimization Selection

tile?

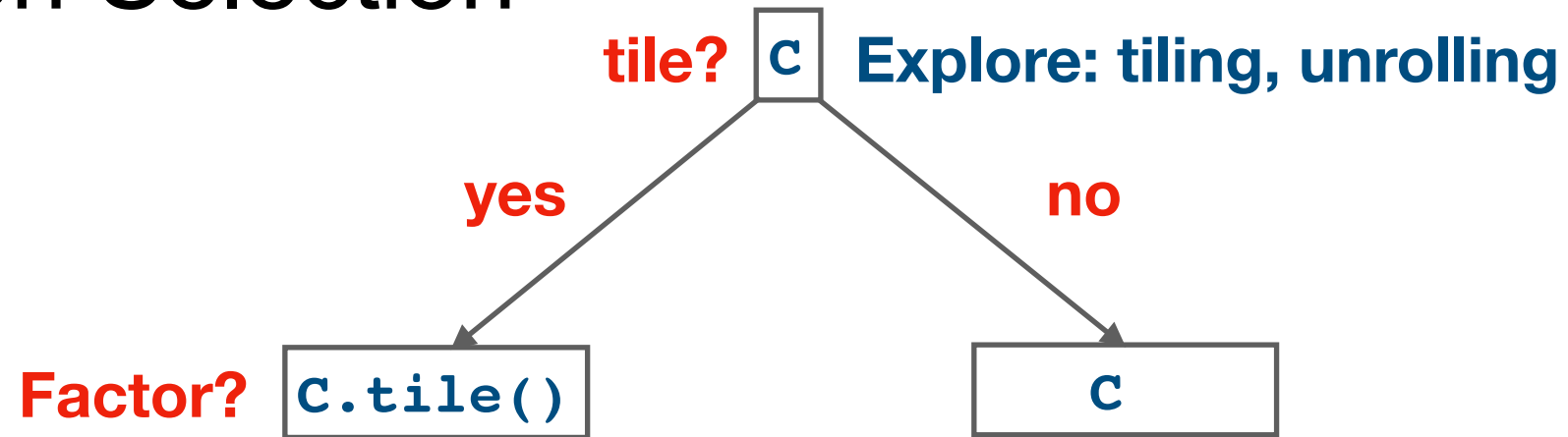


Explore: tiling, unrolling

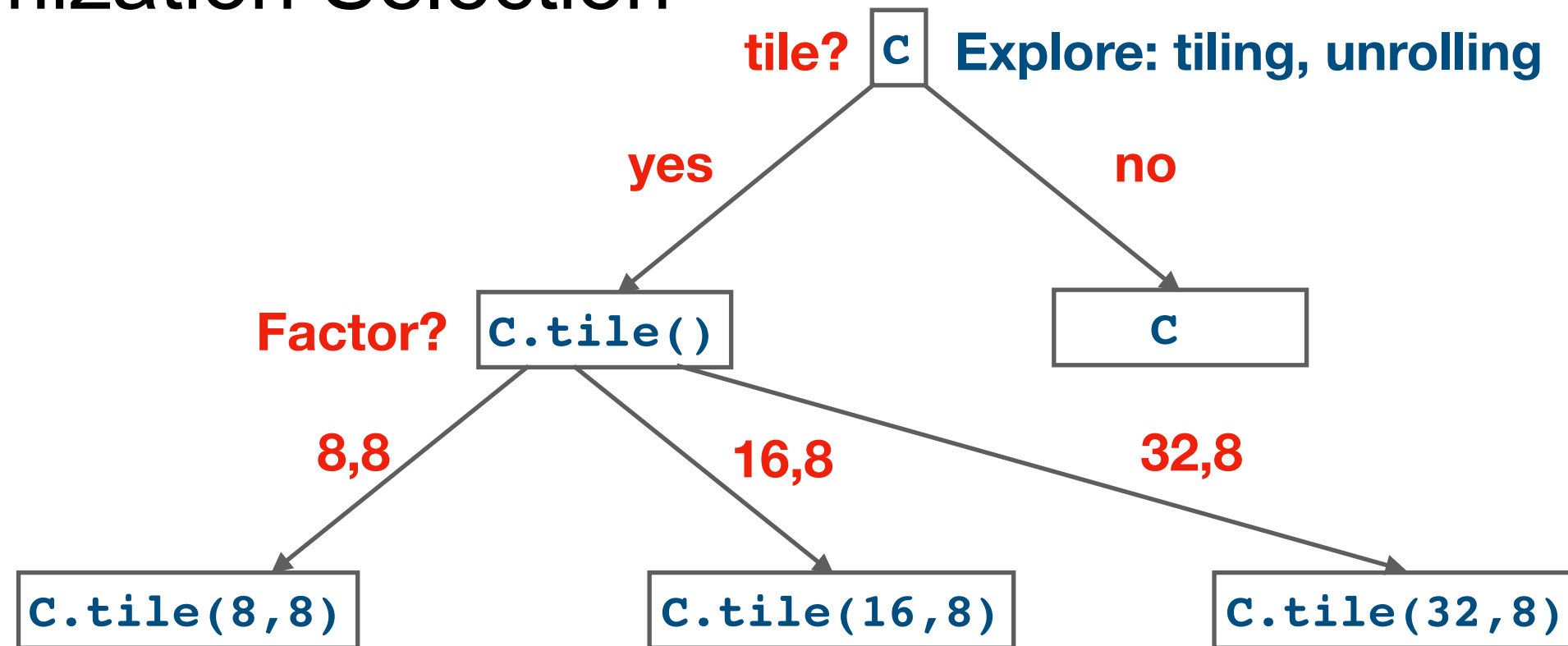
Optimization Selection



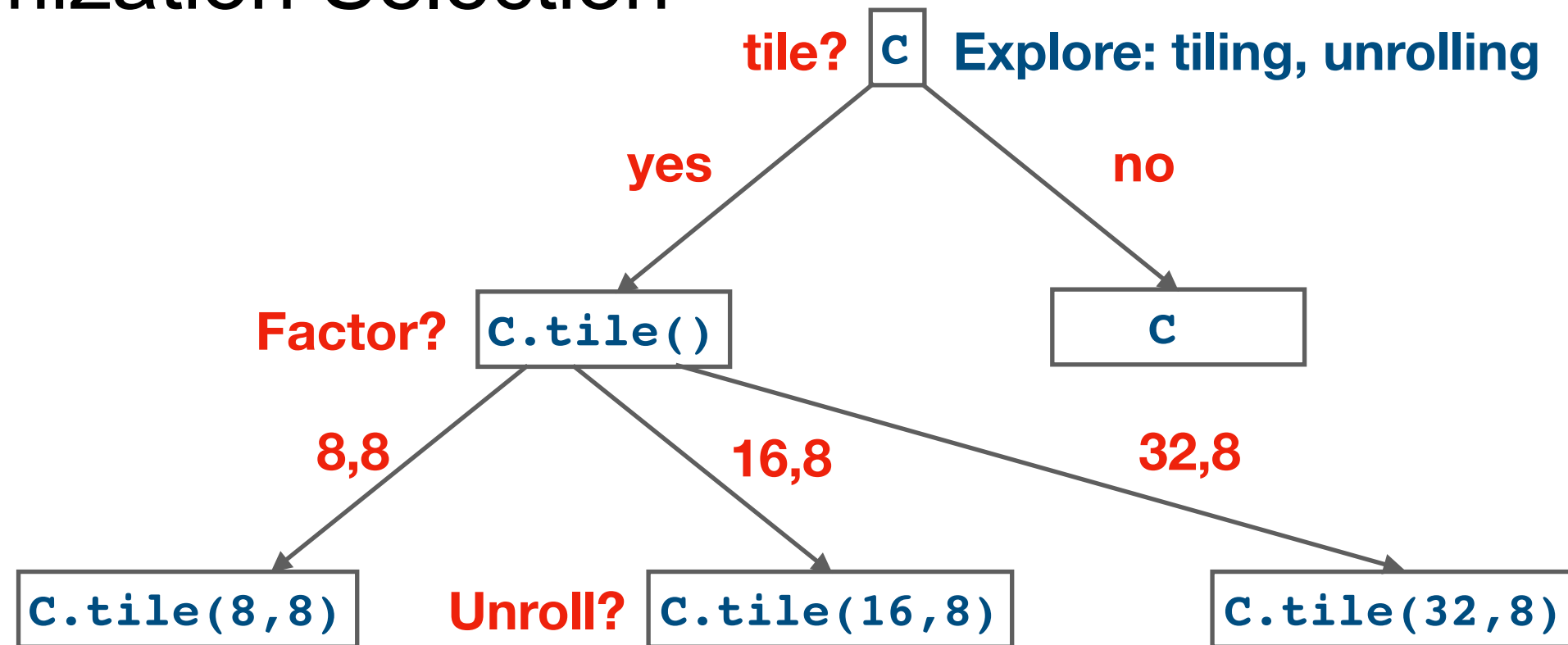
Optimization Selection



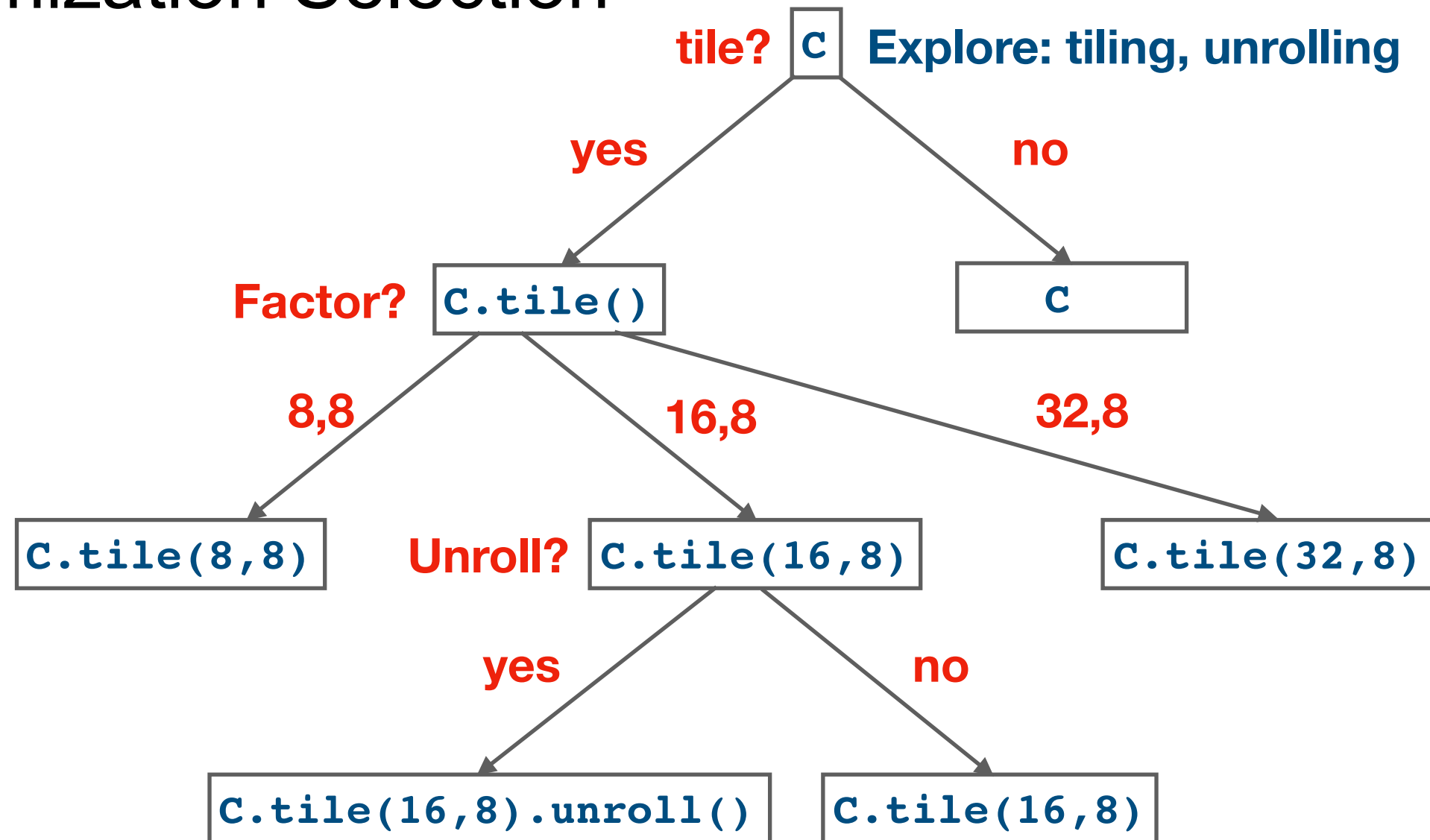
Optimization Selection



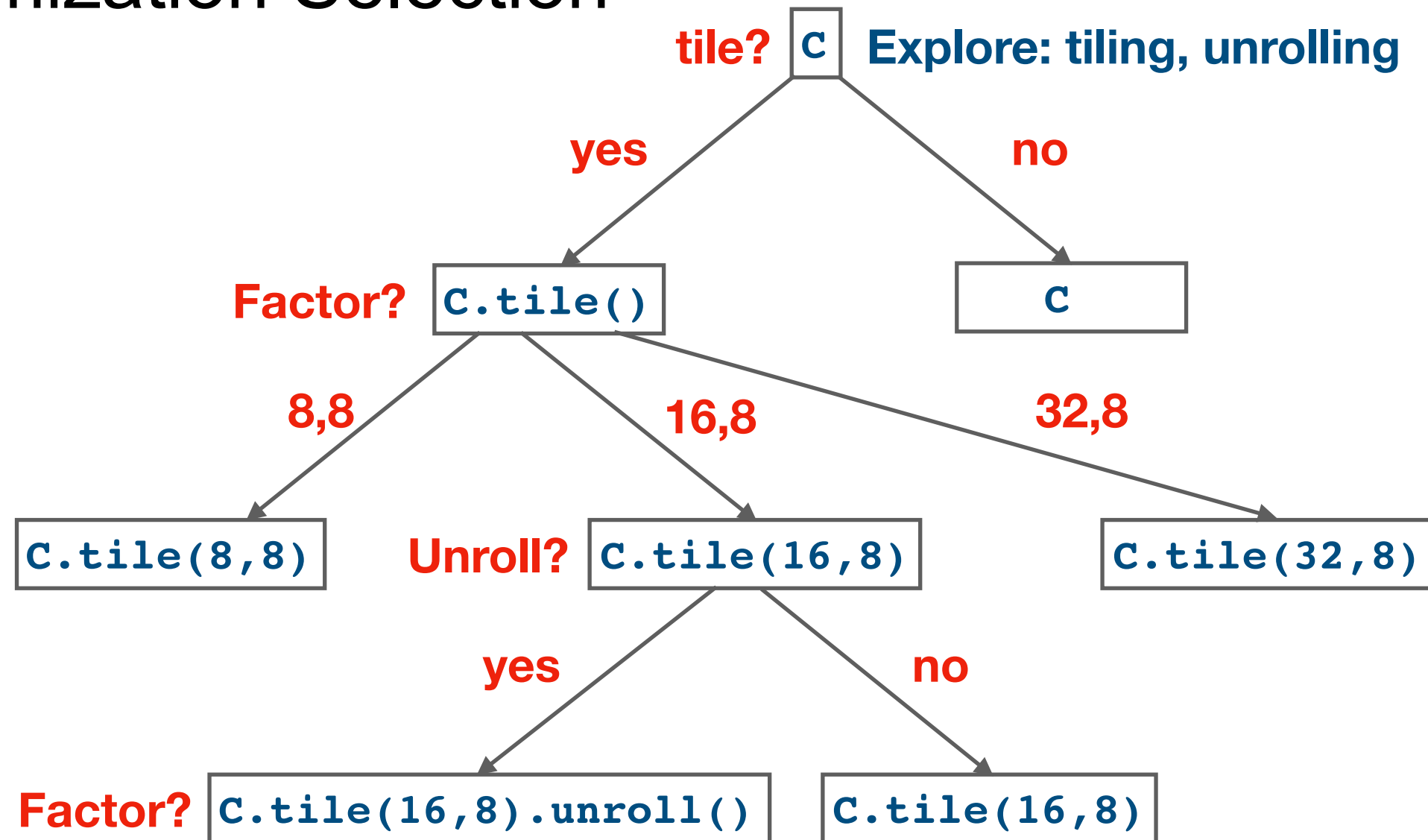
Optimization Selection



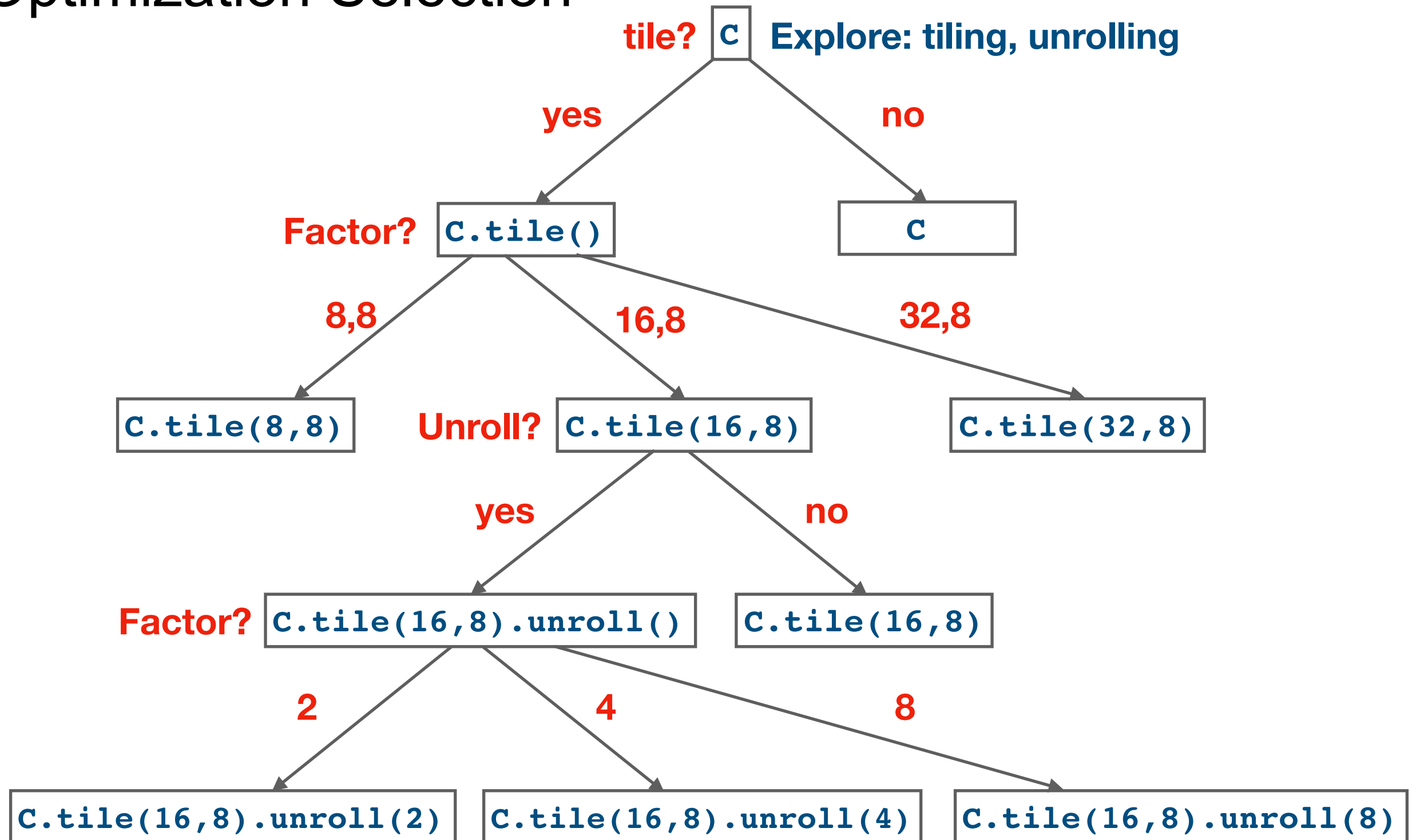
Optimization Selection



Optimization Selection



Optimization Selection



Automatic Optimization

```
graph TD; A[Automatic Optimization] --> B[Selection]; A --> C[Correctness];
```

Selection

Correctness

Automatic Optimization

```
graph TD; A[Automatic Optimization] --> B[Selection]; A --> C[Correctness]; B --> D[Cost model]; B --> E[Search technique];
```

Selection

Correctness

Cost model

Search
technique

Automatic Optimization

```
graph TD; A[Automatic Optimization] --> B[Selection]; A --> C[Correctness]; B --> D[Cost model]; B --> E[Search technique];
```

Selection

Correctness

Cost model

Search
technique

DNN-based
model

Automatic Optimization

```
graph TD; A[Automatic Optimization] --> B[Selection]; A --> C[Correctness]; B --> D[Cost model]; B --> E[Search technique]; D --- F[DNN-based model]; E --- G["Classical techniques (Beam search, ...)"]
```

Selection

Correctness

Cost model

Search
technique

DNN-based
model

Classical techniques
(Beam search, ...)

Automatic Optimization

```
graph TD; AO[Automatic Optimization] --> S[Selection]; AO --> C[Correctness]; S --> CM[Cost model]; S --> ST[Search technique]; CM --- CM_note[DNN-based model]; ST --- ST_note[Classical techniques (Beam search, ...)];
```

Selection

Correctness

Polyhedral dependence
analysis [Feautrier, 91]

Cost model

DNN-based
model

Search
technique

Classical techniques
(Beam search, ...)

Automatic Optimization

```
graph TD; A[Automatic Optimization] --> B[Selection]; A --> C[Correctness]; B --> D[Cost model]; B --> E[Search technique];
```

Selection

Correctness

Polyhedral dependence
analysis [Feautrier, 91]

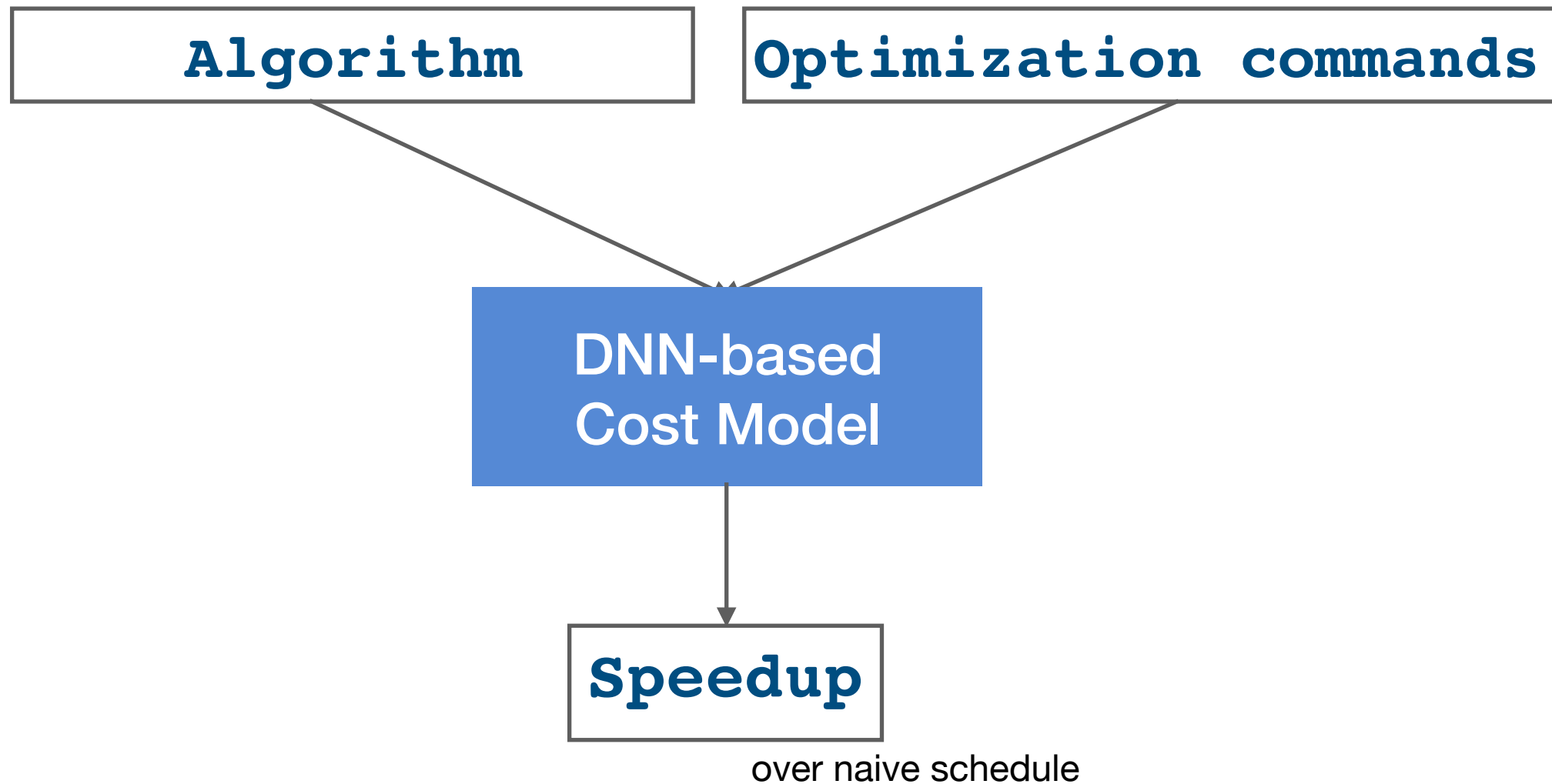
Cost model

DNN-based
model

Search
technique

Classical techniques
(Beam search, ...)

DNN-based Cost Model



DNN-based Cost Model

DNN-based Cost Model

- Training data:
 - Generate random algorithms/schedules
 - Use patterns found in real code
 - Generate 1M data points

DNN-based Cost Model

- Training data:
 - Generate random algorithms/schedules
 - Use patterns found in real code
 - Generate 1M data points
- No feature engineering

DNN-based Cost Model

- Training data:
 - Generate random algorithms/schedules
 - Use patterns found in real code
 - Generate 1M data points
- No feature engineering

DNN-based Cost Model

- Training data:
 - Generate random algorithms/schedules
 - Use patterns found in real code
 - Generate 1M data points
- No feature engineering
- Error rates: 24.95%

Summary

Summary

- Tiramisu has 3 unique features:
 - Optimize sparse DNNs
 - Optimize RNNs
 - Generate code for distributed architectures

Summary

- Tiramisu has 3 unique features:
 - Optimize sparse DNNs
 - Optimize RNNs
 - Generate code for distributed architectures
- Two main design choices
 - Polyhedral representation
 - Separates algorithm from optimization commands

Summary

- Tiramisu has 3 unique features:
 - Optimize sparse DNNs
 - Optimize RNNs
 - Generate code for distributed architectures
- Two main design choices
 - Polyhedral representation
 - Separates algorithm from optimization commands
- Uses deep learning to build a cost model for automatic optimization

Summary

- Tiramisu has 3 unique features:
 - Optimize sparse DNNs
 - Optimize RNNs
 - Generate code for distributed architectures
- Two main design choices
 - Polyhedral representation
 - Separates algorithm from optimization commands
- Uses deep learning to build a cost model for automatic optimization
- Web: <http://tiramisu-compiler.org>

Summary

- Tiramisu has 3 unique features:
 - Optimize sparse DNNs
 - Optimize RNNs
 - Generate code for distributed architectures
- Two main design choices
 - Polyhedral representation
 - Separates algorithm from optimization commands
- Uses deep learning to build a cost model for automatic optimization
- Web: <http://tiramisu-compiler.org>
- Interested in building a community around Tiramisu