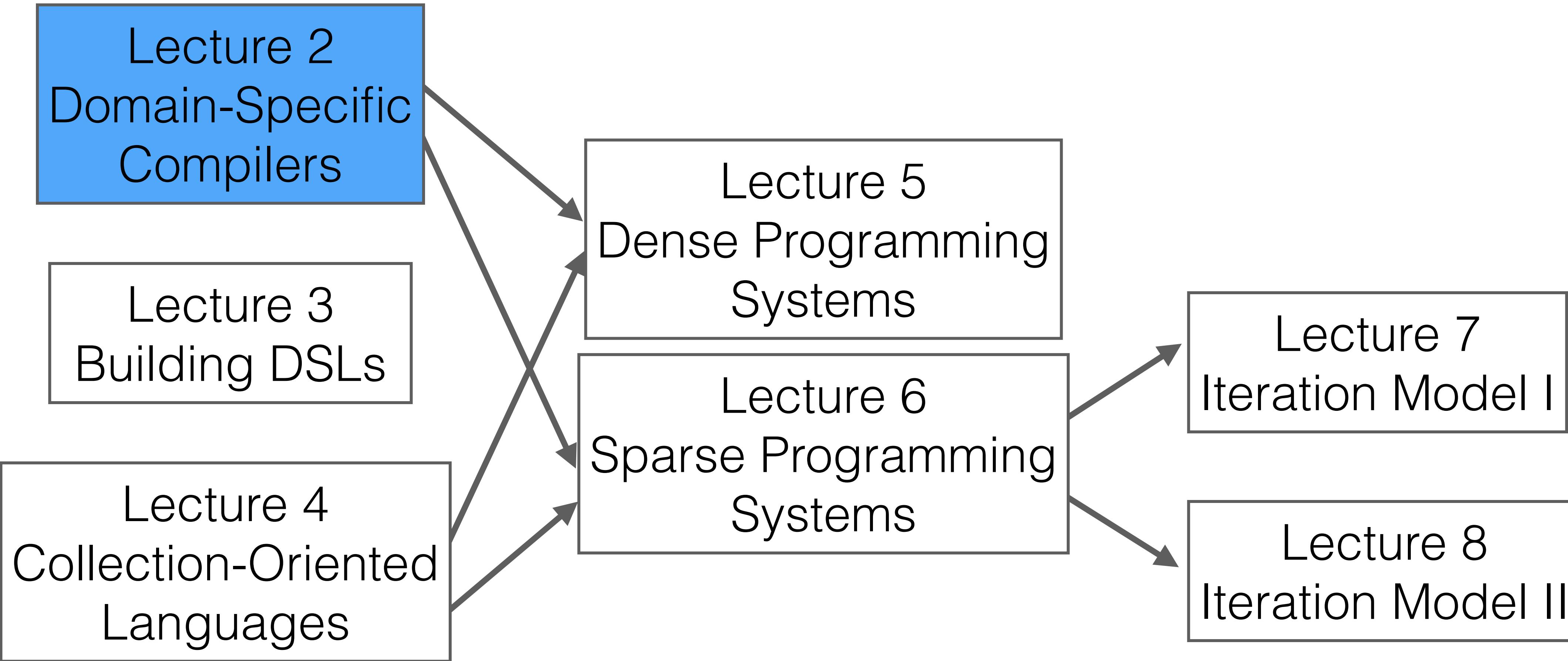
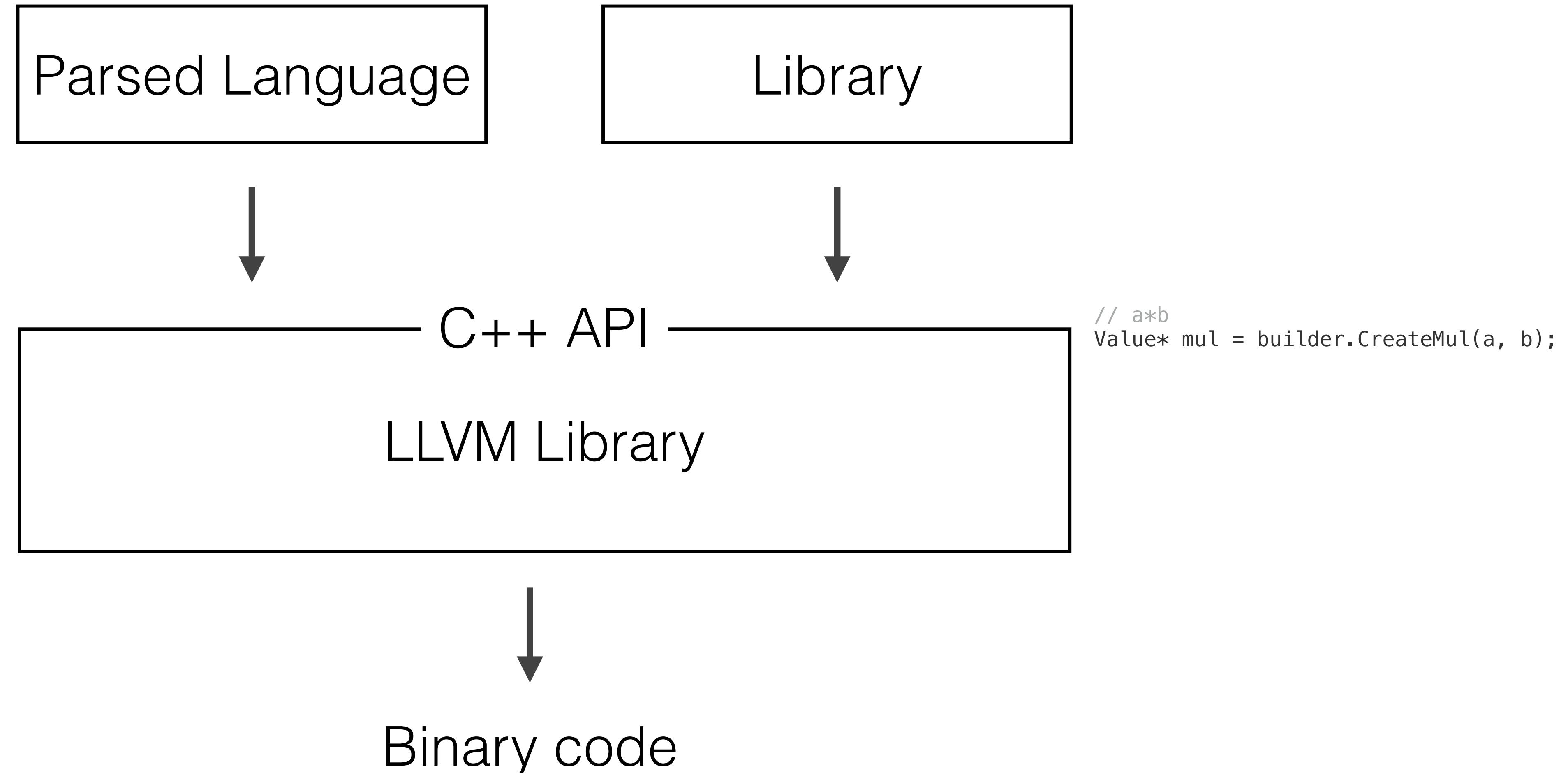


Lecture 2 – Why Domain-Specific Compilers

Stanford CS343D (Winter 2025)
Fred Kjolstad



Languages vs libraries: LLVM is a compiler for general languages, yet it is a library that does not require a parser



Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

1960s - Simple expression optimizations

- Common subexpression elimination
- Operator strength reduction
- Program optimization (Allen 1969)

Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

1960s - Simple expression optimizations

- Common subexpression elimination
- Operator strength reduction
- Program optimization (Allen 1969)

1970s - Loop transformations

- Strip-mining and tiling
- Simple loop parallelization
- Dependence testing
- Lots of work at Illinois

Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

1980s - Loop analysis and parallelization

- Auto-vectorization
- Auto-parallelization

1960s - Simple expression optimizations

- Common subexpression elimination
- Operator strength reduction
- Program optimization (Allen 1969)

1970s - Loop transformations

- Strip-mining and tiling
- Simple loop parallelization
- Dependence testing
- Lots of work at Illinois

Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

1960s - Simple expression optimizations

- Common subexpression elimination
- Operator strength reduction
- Program optimization (Allen 1969)

1970s - Loop transformations

- Strip-mining and tiling
- Simple loop parallelization
- Dependence testing
- Lots of work at Illinois

1980s - Loop analysis and parallelization

- Auto-vectorization
- Auto-parallelization

1990s - Polyhedral model

- Auto-scheduling
- Code generation through scanning
- Remove cost of object-orientation

Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

1960s - Simple expression optimizations

- Common subexpression elimination
- Operator strength reduction
- Program optimization (Allen 1969)

1970s - Loop transformations

- Strip-mining and tiling
- Simple loop parallelization
- Dependence testing
- Lots of work at Illinois

1980s - Loop analysis and parallelization

- Auto-vectorization
- Auto-parallelization

1990s - Polyhedral model

- Auto-scheduling
- Code generation through scanning
- Remove cost of object-orientation

2000s - SSA and LLVM

- SSA becomes widely used
- LLVM makes compilers accessible
- Remove cost of dynamic languages

Selective highlights from the history of compilers

1950s - Fortran Compiler

- The History of Fortran (Backus 1982)

1960s - Simple expression optimizations

- Common subexpression elimination
- Operator strength reduction
- Program optimization (Allen 1969)

1970s - Loop transformations

- Strip-mining and tiling
- Simple loop parallelization
- Dependence testing
- Lots of work at Illinois

1980s - Loop analysis and parallelization

- Auto-vectorization
- Auto-parallelization

1990s - Polyhedral model

- Auto-scheduling
- Code generation through scanning
- Remove cost of object-orientation

2000s - SSA and LLVM

- SSA becomes widely used
- LLVM makes compilers accessible
- Remove cost of dynamic languages

2010-2020s - DSLs and Synthesis

- Halide, TensorFlow/XLA, Taco
- Code generation for SQL

Automatic programming

The compiler as an optimizer

```
for (int i = 0; i < M; i++) {  
    double t = 0.0;  
    for (int j = 0; j < N; j++) {  
        int pB2 = i*N + j;  
        t += B[pB2] * c[j];  
    }  
    a[i] = t;  
}
```

optimize →

```
for (int i = 0; i < M; i++) {  
    double t = 0.0;  
    for (int p = B_pos[i]; p < B_pos[i+1]; p++) {  
        int j = B_crd[p];  
        t += B[p] * c[j];  
    }  
    a[i] = t;  
}
```

The compiler as generator

$$a = Bc$$

↓
lower

```
for (int i = 0; i < M; i++) {  
    double t = 0.0;  
    for (int p = B_pos[i]; p < B_pos[i+1]; p++) {  
        int j = B_crd[p];  
        t += B[p] * c[j];  
    }  
    a[i] = t;  
}
```

Automatic programming

The compiler as an optimizer

```
for (int i = 0; i < M; i++) {  
    double t = 0.0;  
    for (int j = 0; j < N; j++) {  
        int pB2 = i*N + j;  
        t += B[pB2] * c[j];  
    }  
    a[i] = t;  
}
```

optimize →

```
for (int i = 0; i < M; i++) {  
    double t = 0.0;  
    for (int p = B_pos[i]; p < B_pos[i+1]; p++) {  
        int j = B_crd[p];  
        t += B[p] * c[j];  
    }  
    a[i] = t;  
}
```

The compiler as generator

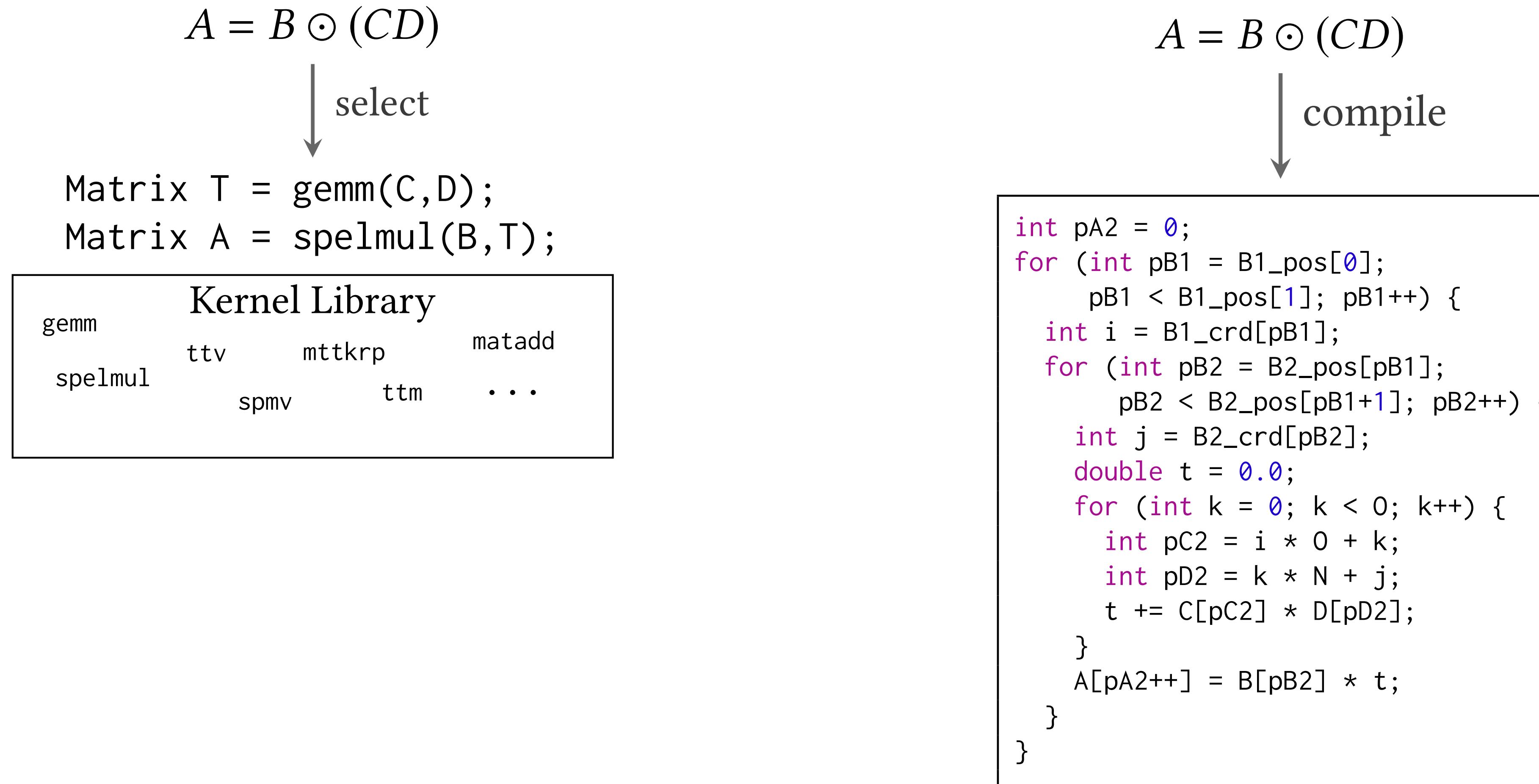
$$a = Bc$$

↓
lower

```
for (int i = 0; i < M; i++) {  
    double t = 0.0;  
    for (int p = B_pos[i]; p < B_pos[i+1]; p++) {  
        int j = B_crd[p];  
        t += B[p] * c[j];  
    }  
    a[i] = t;  
}
```

“In short, automatic programming always has been a euphemism for programming with a higher-level language than was then available to the programmer. Research in automatic programming is simply research in the implementation of higher-level programming languages.”
- David Parnas

Granularity of generated code



What does a compiler do for you?

1. Lets you program a different machine than the one you actually have
 - A high-level language is an imaginary machine (virtual machine)
 - The compiler automatically programs the actual machine for you
2. Lets you know if you are using the language incorrectly
3. Optimizes the performance of your program

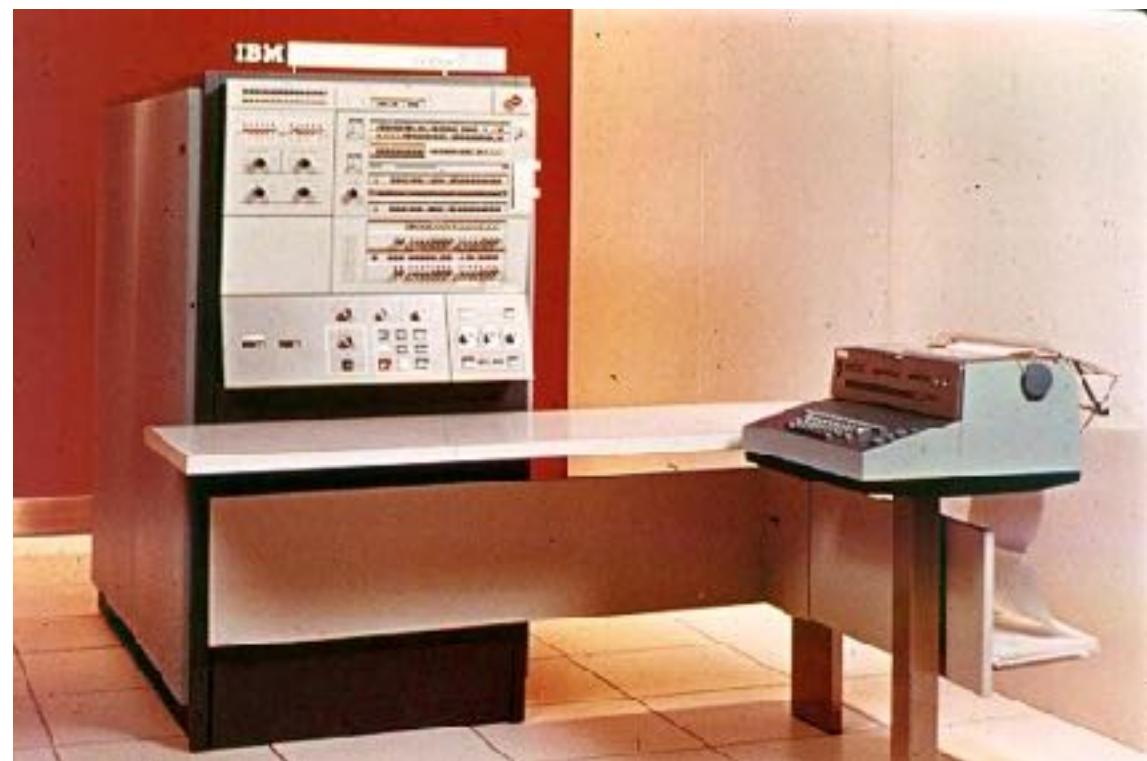
In the Golden Era of Computing,
Performance Engineering Ruled the World

Every programmer was a Performance Engineer

In the Golden Era of Computing, Performance Engineering Ruled the World

Every programmer was a Performance Engineer

IBM System/360

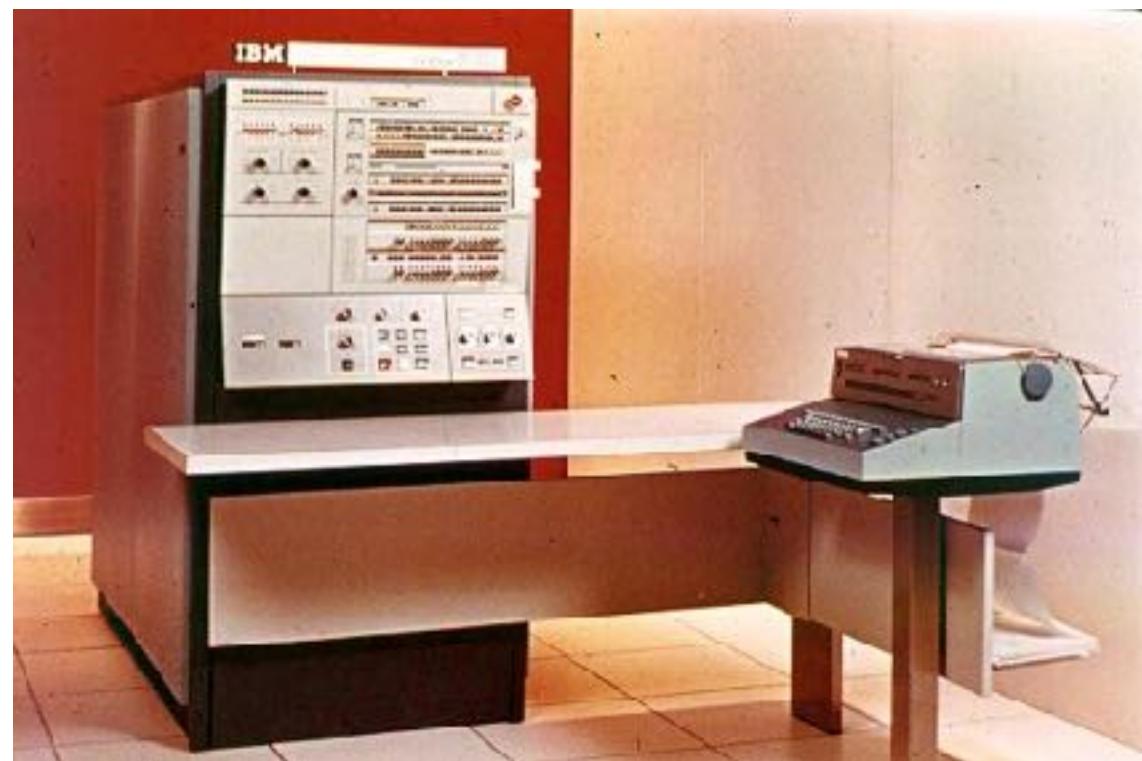


Launched: 1964
Clock rate: 33 KHz
Data path: 32bits
Memory: 524 Kbytes
Cost: \$5,000 per month

In the Golden Era of Computing, Performance Engineering Ruled the World

Every programmer was a Performance Engineer

IBM System/360



Launched: 1964
Clock rate: 33 KHz
Data path: 32bits
Memory: 524 Kbytes
Cost: \$5,000 per month

Apple II

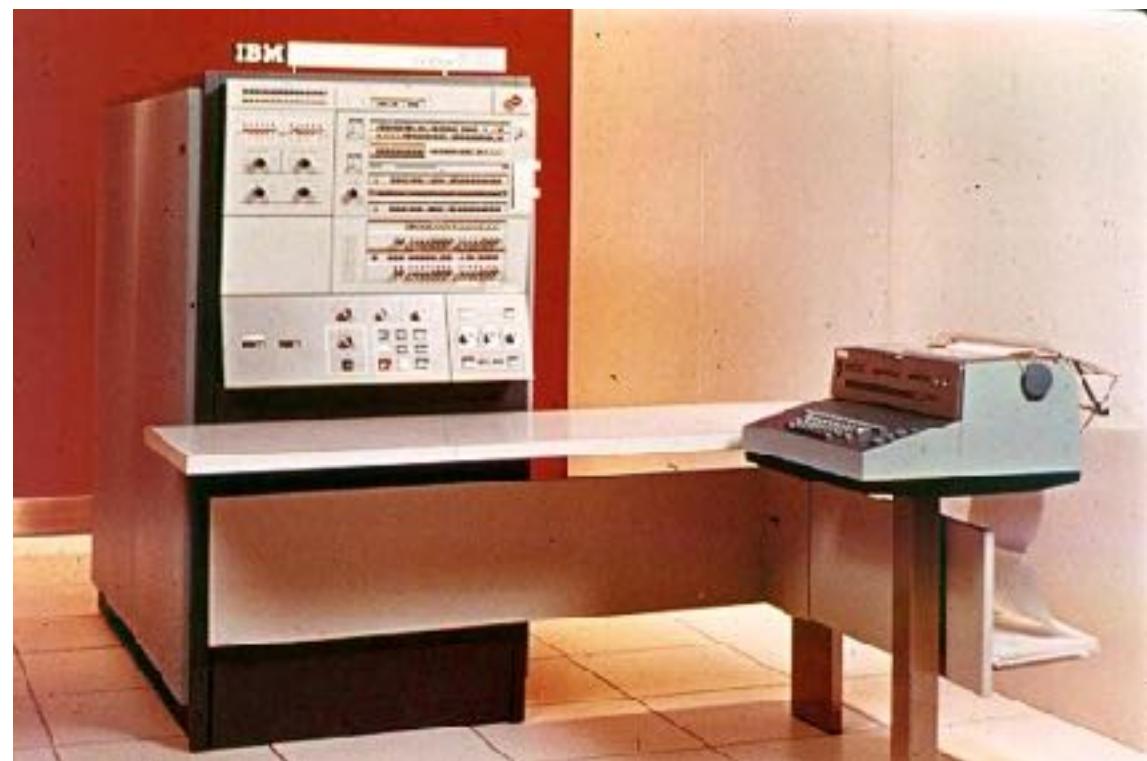


Launched: 1977
Clock rate: 1 MHz
Data path: 8 bits
Memory: 48 Kbytes
Cost: \$1,395

In the Golden Era of Computing, Performance Engineering Ruled the World

Every programmer was a Performance Engineer

IBM System/360



Launched: 1964
Clock rate: 33 KHz
Data path: 32bits
Memory: 524 Kbytes
Cost: \$5,000 per month

Apple II



Launched: 1977
Clock rate: 1 MHz
Data path: 8 bits
Memory: 48 Kbytes
Cost: \$1,395

Any useful program would stretch the machine resources
Program had to be planned around the machine
Many would not ‘fit’ without intense performance hacks

Software Properties

What do programmers want to add?

Software Properties

What do programmers want to add?

- Functionality

Software Properties

What do programmers want to add?

- Functionality
 - ... and...

Software Properties

What do programmers want to add?

- Functionality

... and...

- Scalability
- Compatibility
- Correctness
- Clarity

- Low Power
- Maintainability
- Modularity
- Portability

- Reliability
- Robustness
- Testability
- Usability

... and more.

Software Properties

What do programmers want to add?

- Functionality

... and...

- Scalability
- Compatibility
- Correctness
- Clarity

- Low Power
- Maintainability
- Modularity
- Portability

- Reliability
- Robustness
- Testability
- Usability

... and more.

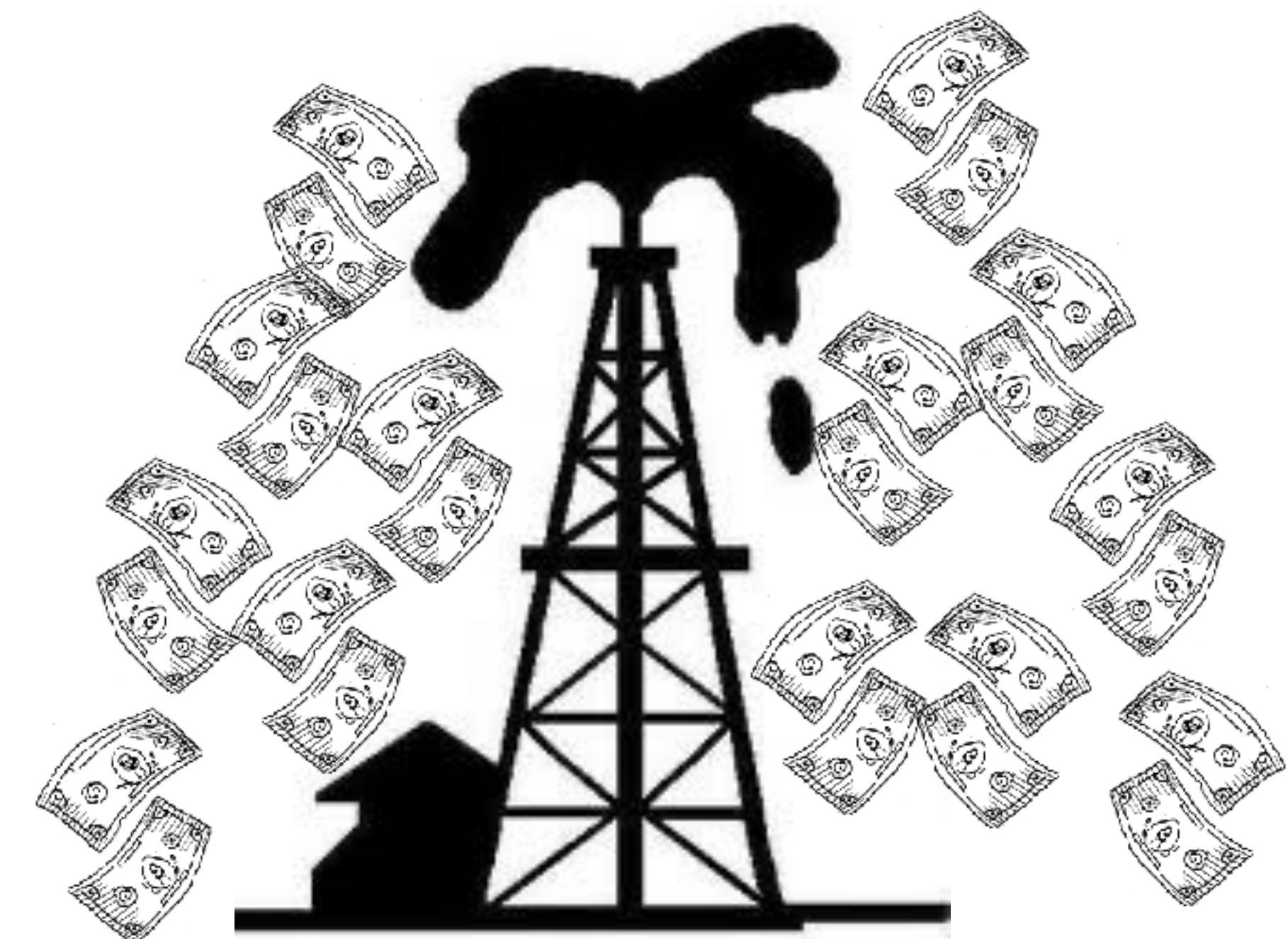
Performance is the **currency** of computing. You can often "buy" needed properties with performance.

In the Dominant Era of Computing, Performance became Free

The currency was free

Only need to wait a few months

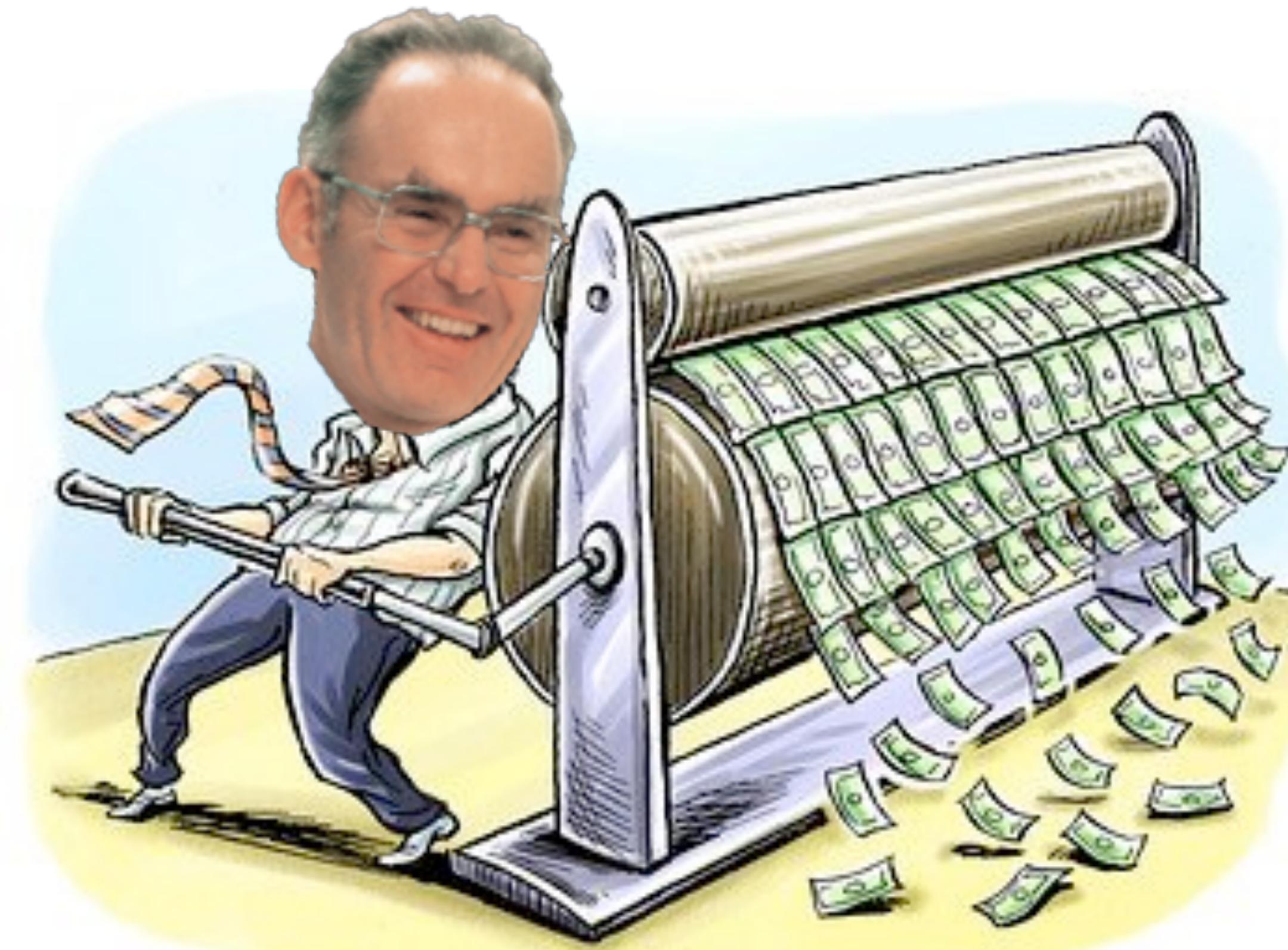
Performance doubled every 2 years



Performance is the **currency** of computing. You can often "buy" needed properties with performance.

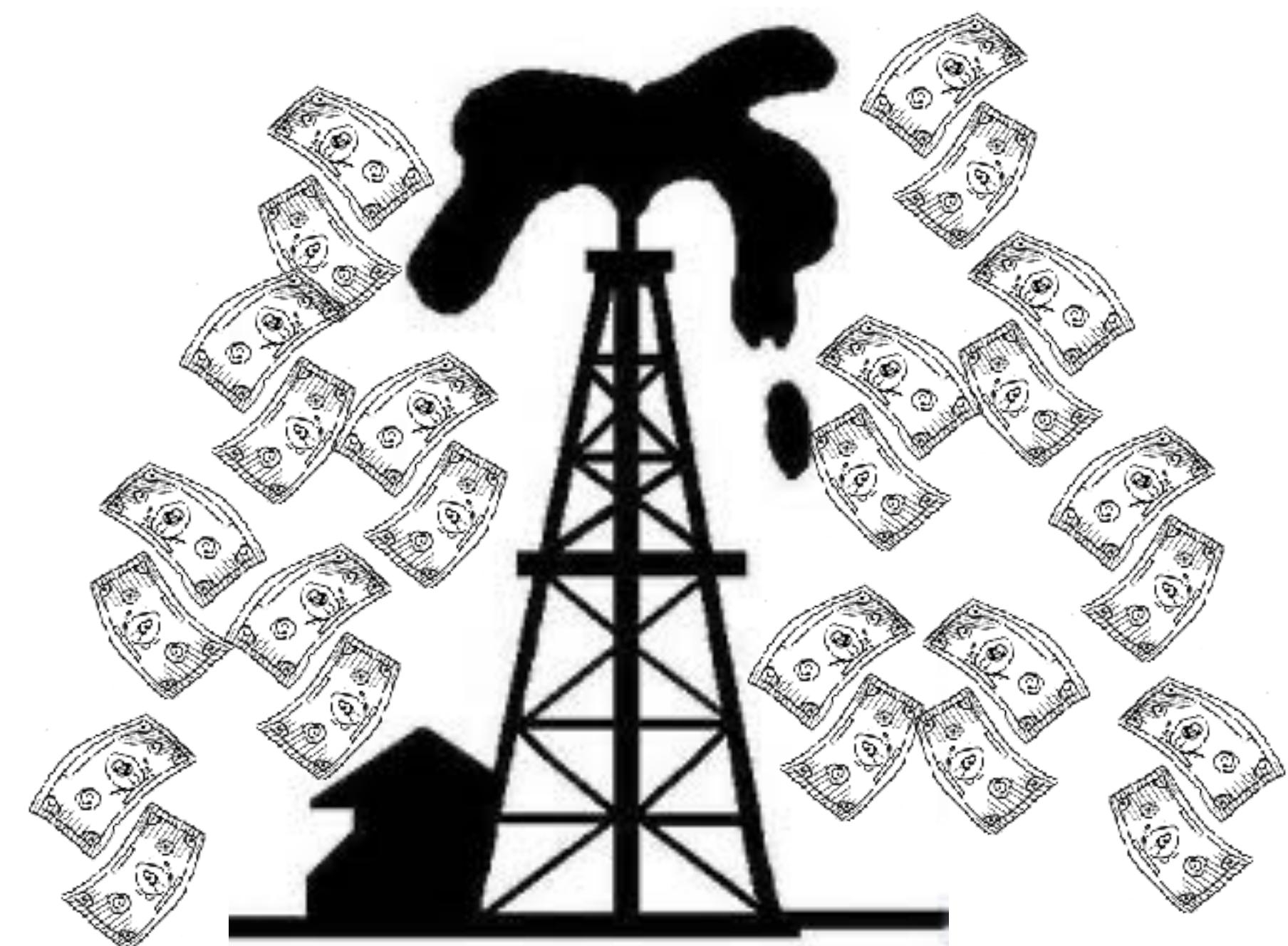
In the Dominant Era, Performance was Free

Moore's Law and the scaling of clock frequency
= printing press for the currency of performance



In the Dominant Era, Performance was Free

Performance engineering was
'optional' at best and
'irrelevant' for most programmers



Performance is the **currency** of computing. You can often "buy" needed properties with performance.

The Age of Free Performance is Over

Moore's law is not giving free performance any more



Performance is the **currency** of computing. You can often "buy" needed properties with performance.

The Age of Free Performance is Over

Two ways to get better performance:

1. Remove software abstractions costs
2. Build domain-specific hardware



Both requires specialization. A compiler is
a generator of specialized code.

Inefficient abstractions mechanisms in software and inefficient use of hardware

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96
7	plus AVX intrinsics	0.41	337.812	62,806	2.7	40.45

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96
7	plus AVX intrinsics	0.41	337.812	62,806	2.7	40.45

17h vs 1s!

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96
7	plus AVX intrinsics	0.41	337.812	62,806	2.7	40.45

17h vs 1s!

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Parallelism

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96
7	plus AVX intrinsics	0.41	337.812	62,806	2.7	40.45

17h vs 1s!

Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Parallelism

Locality

Inefficient abstractions mechanisms in software and inefficient use of hardware

“Abstraction is selective ignorance”

- Andrew Koenig

“All problems in computer science can be solved by another level of abstraction, except the problems of too many levels of abstraction.”

- David Wheeler

Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices. Each version represents a successive refinement of the original Python code. “Running time” is the running time of the version. “GFLOPS” is the billions of 64-bit floating-point operations per second that the version executes. “Absolute speedup” is time relative to Python, and “relative speedup,” which we show with an additional digit of precision, is time relative to the preceding line. “Fraction of peak” is GFLOPS relative to the computer’s peak 835 GFLOPS. See Methods for more details.

Version	Implementation	Running time (s)	GFLOPS	Absolute speedup	Relative speedup	Fraction of peak (%)
1	Python	25,552.48	0.005	1	–	0.00
2	Java	2,372.68	0.058	11	10.8	0.01
3	C	542.67	0.253	47	4.4	0.03
4	Parallel loops	69.80	1.969	366	7.8	0.24
5	Parallel divide and conquer	3.80	36.180	6,727	18.4	4.33
6	plus vectorization	1.10	124.914	23,224	3.5	14.96
7	plus AVX intrinsics	0.41	337.812	62,806	2.7	40.45

17h vs 1s!

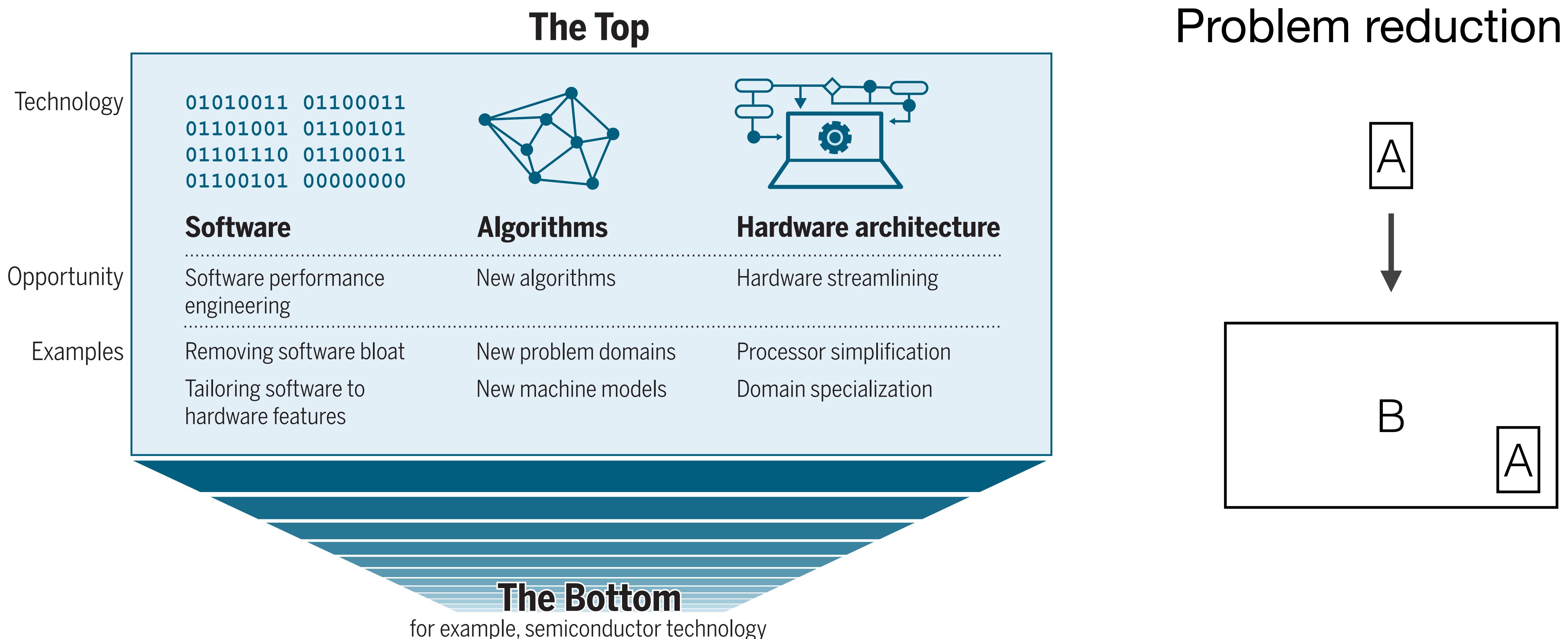
Figure from “There’s plenty of room at the Top: What will drive computer performance after Moore’s law?” Leiserson et al., Science 368

Parallelism

Locality

Specialization

There's plenty of room at the top



Abstraction with friction from traditional library composition

$$A = B \odot (CD)$$

Traditional Library Composition

```
T = matmul(C, D);  
A = elmul(B, T);
```

Three pitfalls:

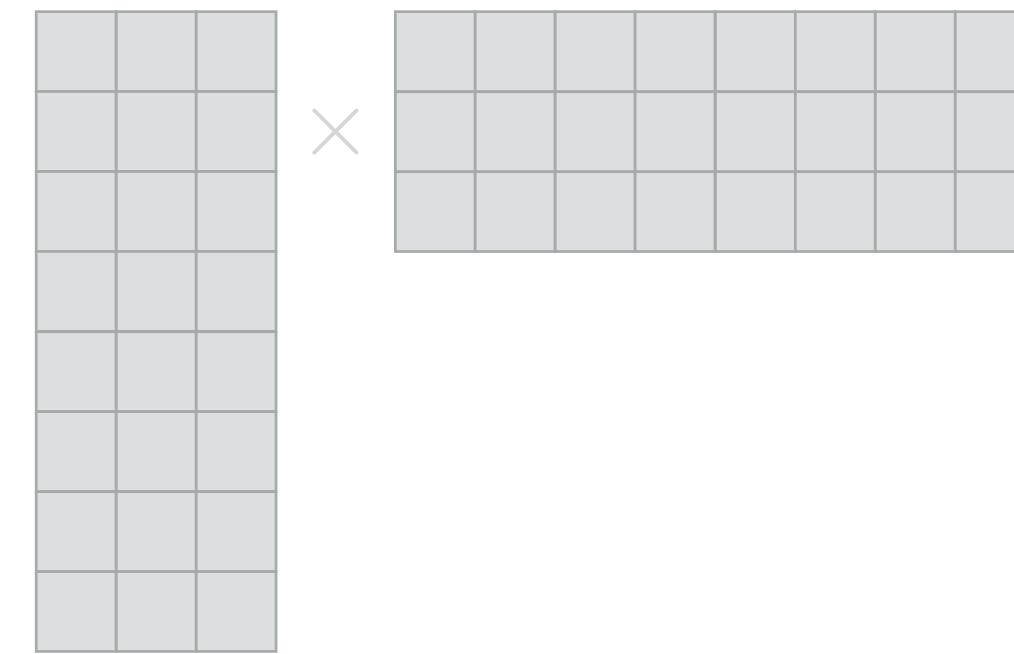
1. Loose temporal locality
2. Data structures must match what functions expect
3. May cause asymptotic slow-down

Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$

Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$


$$\begin{matrix} & \times & \end{matrix}$$

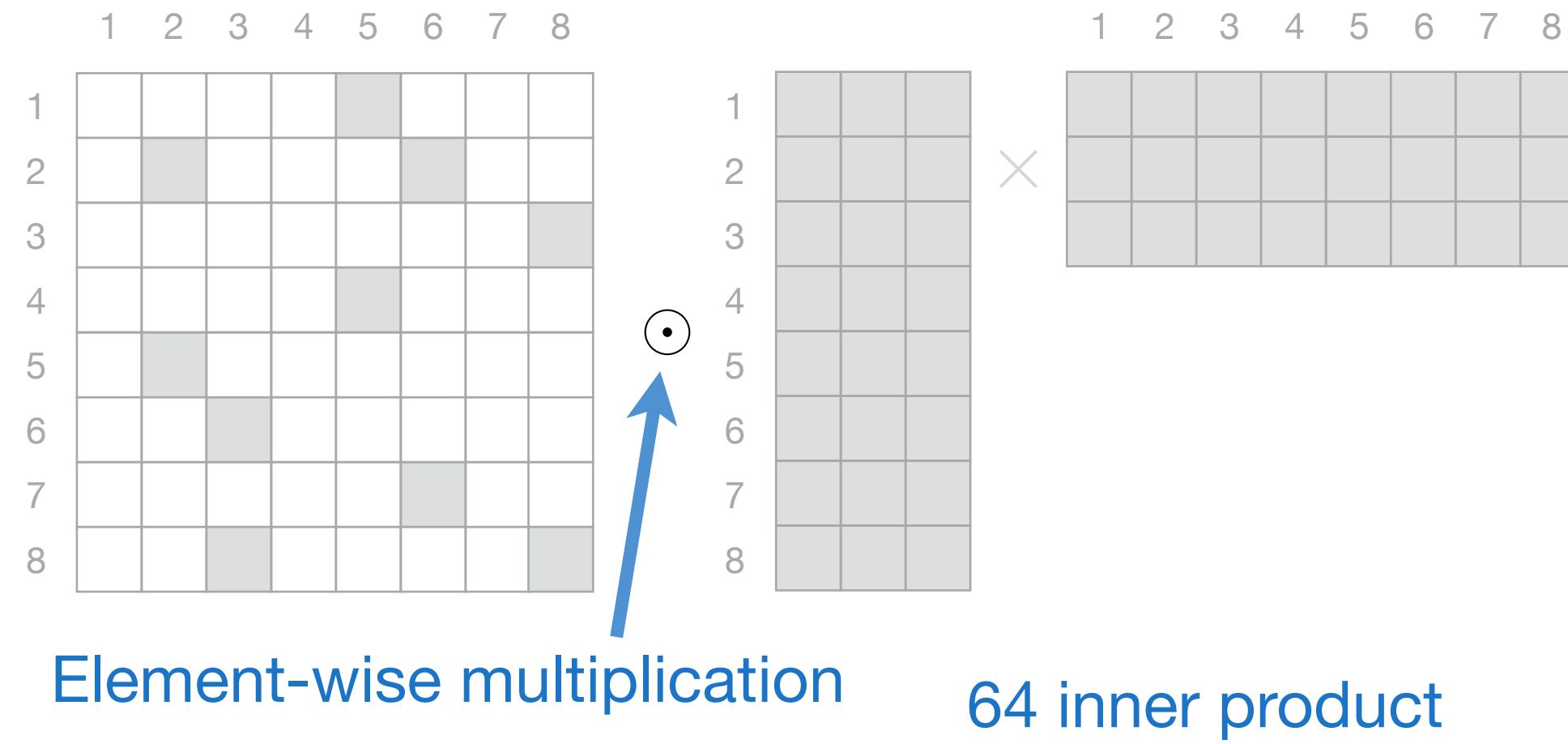
Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$

64 inner product

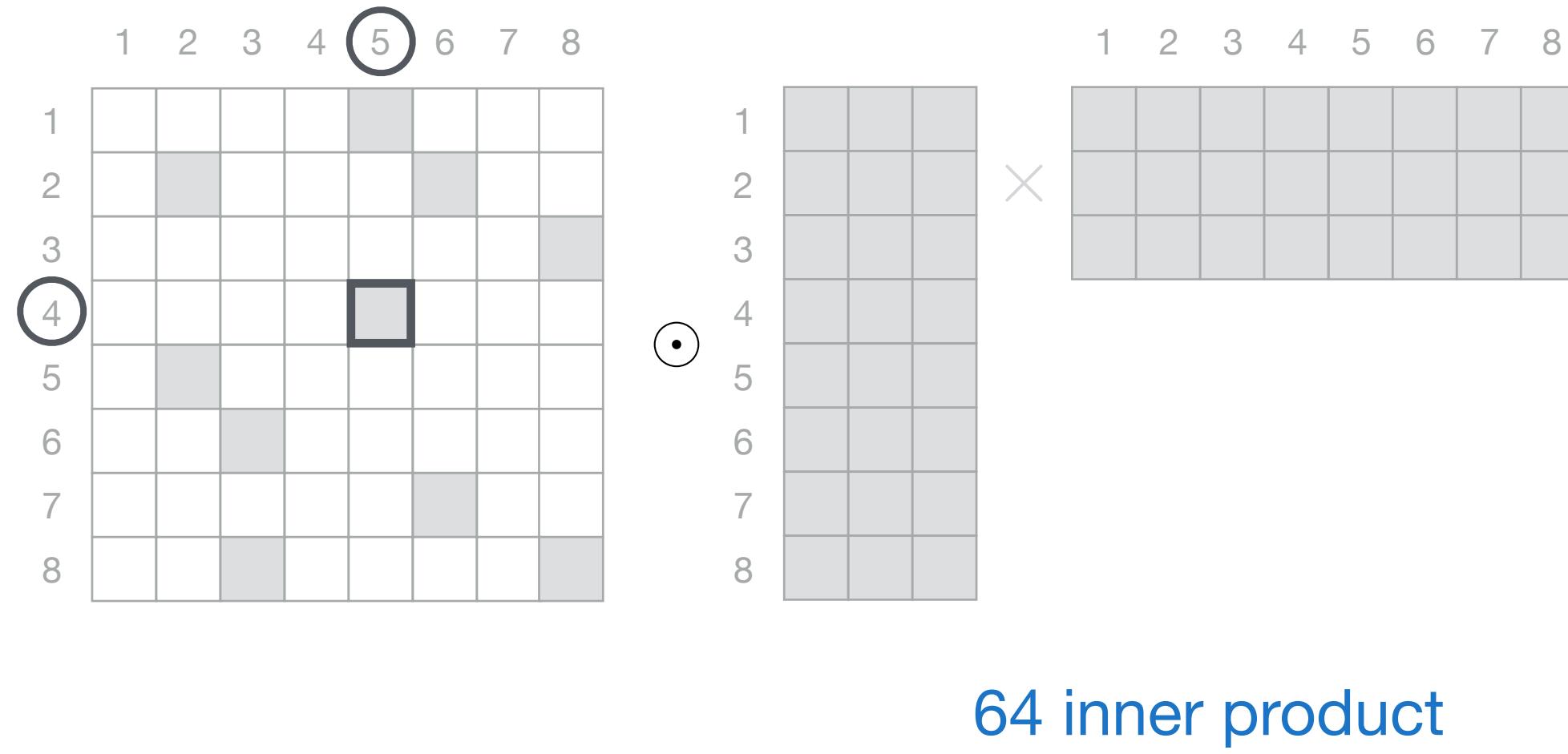
Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$



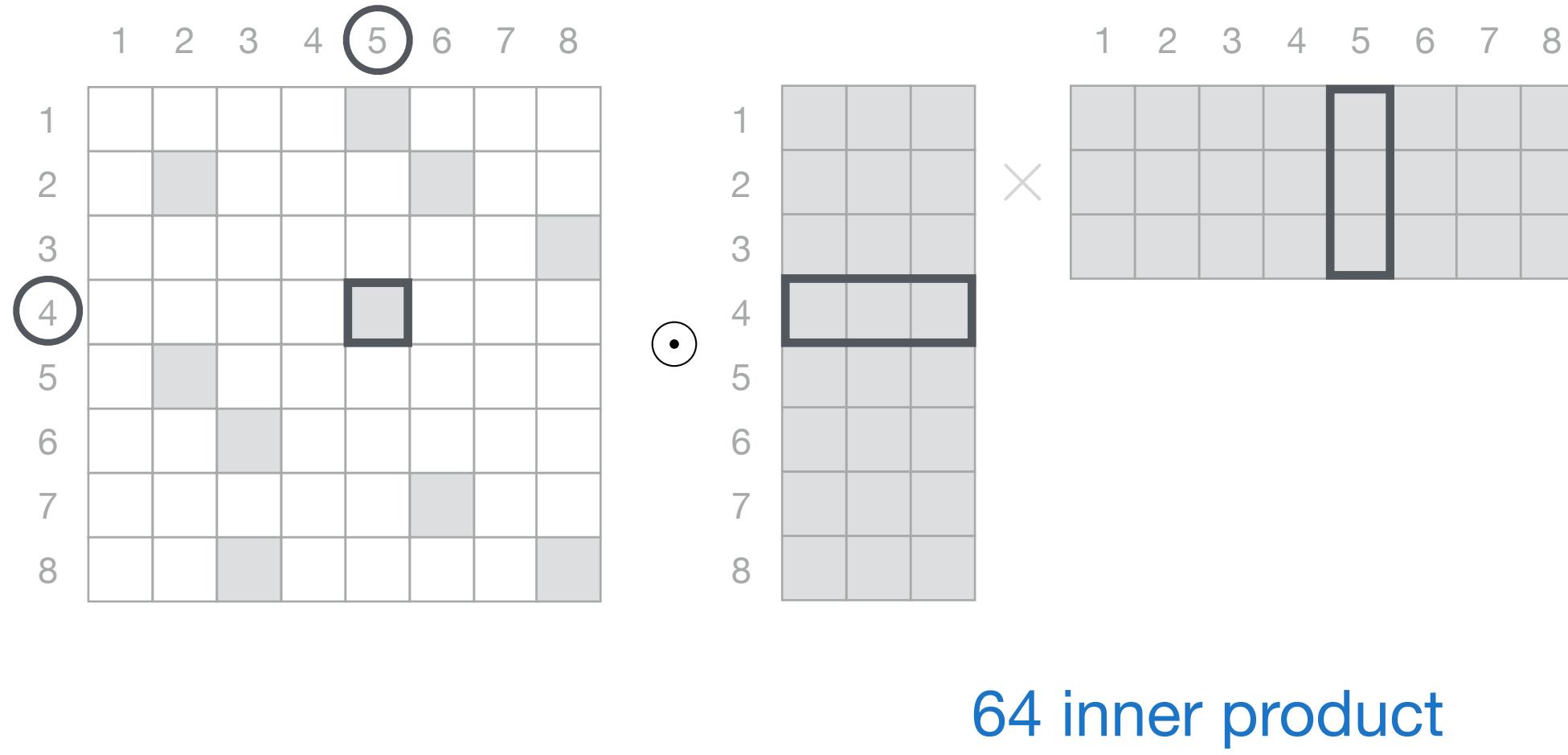
Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$



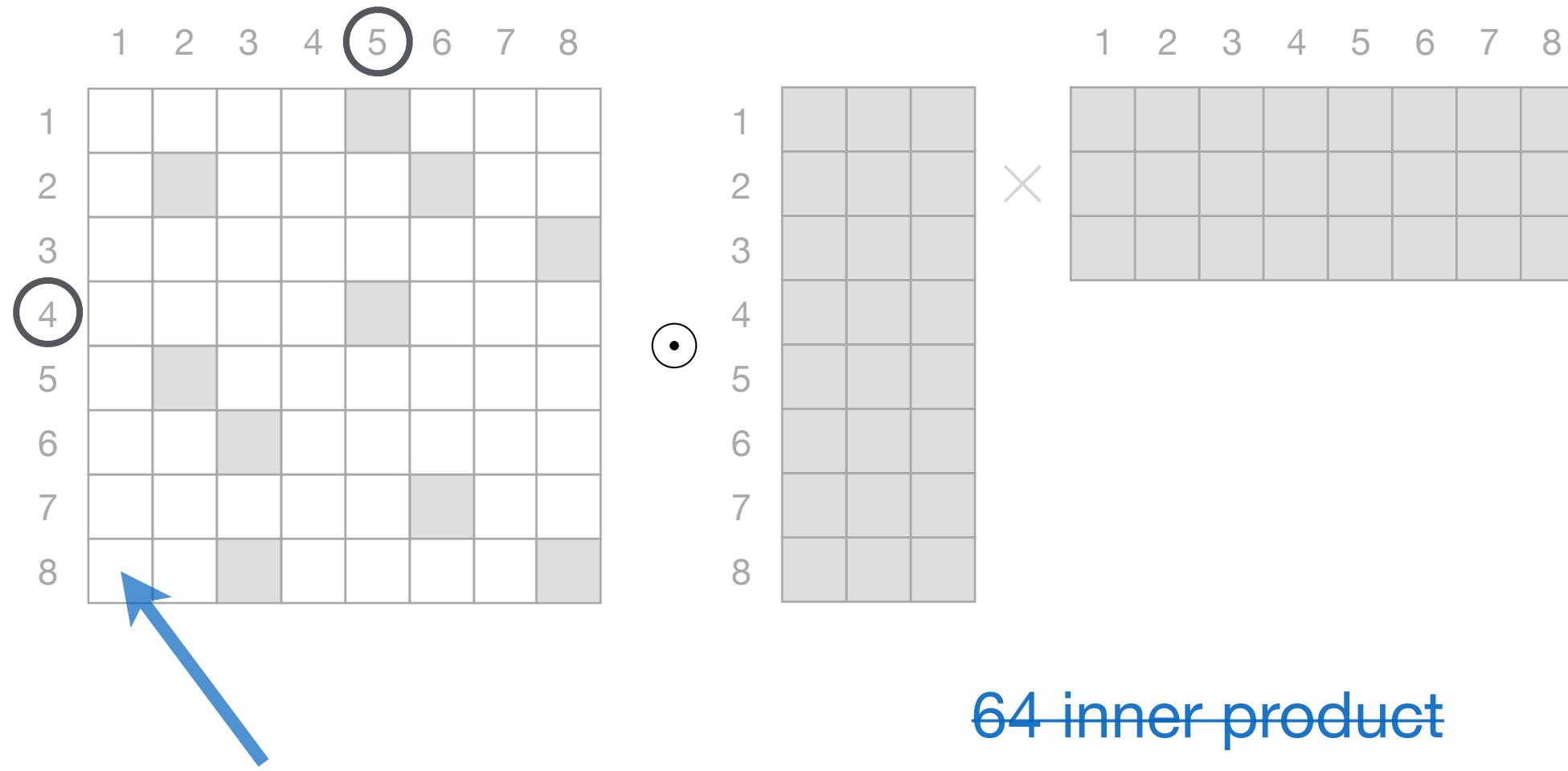
Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$



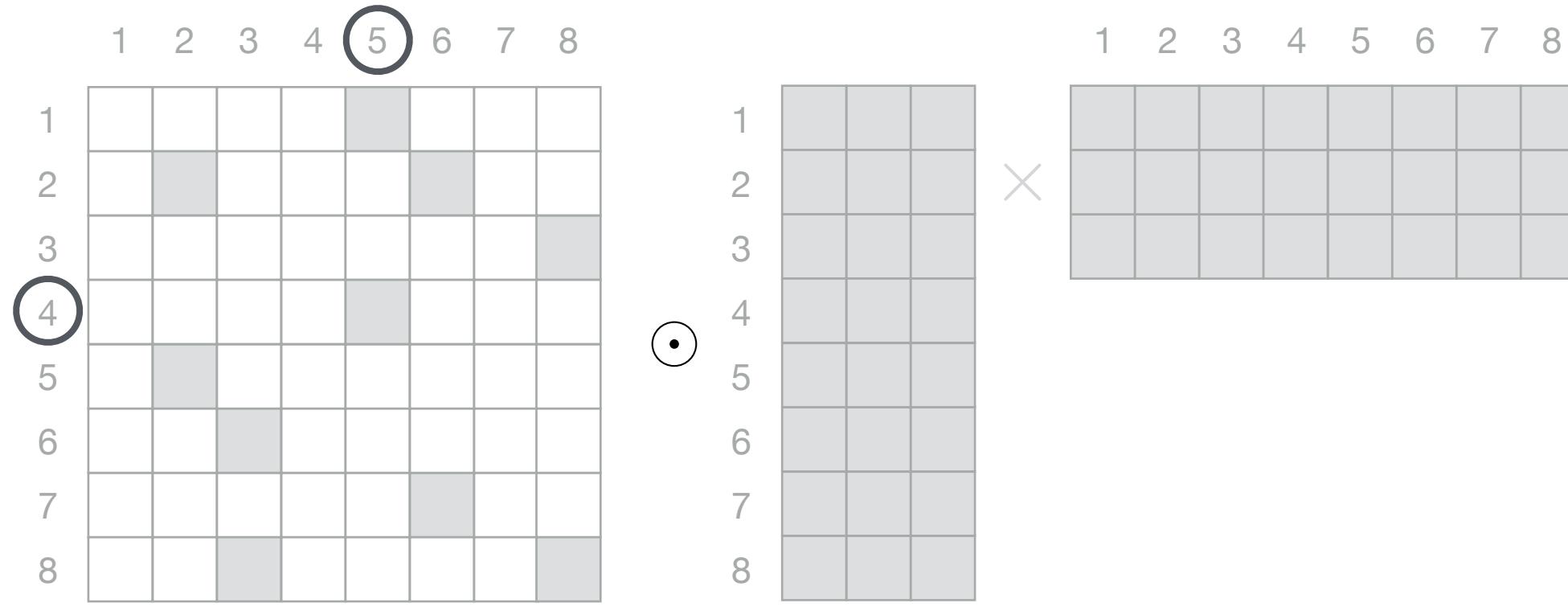
Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$



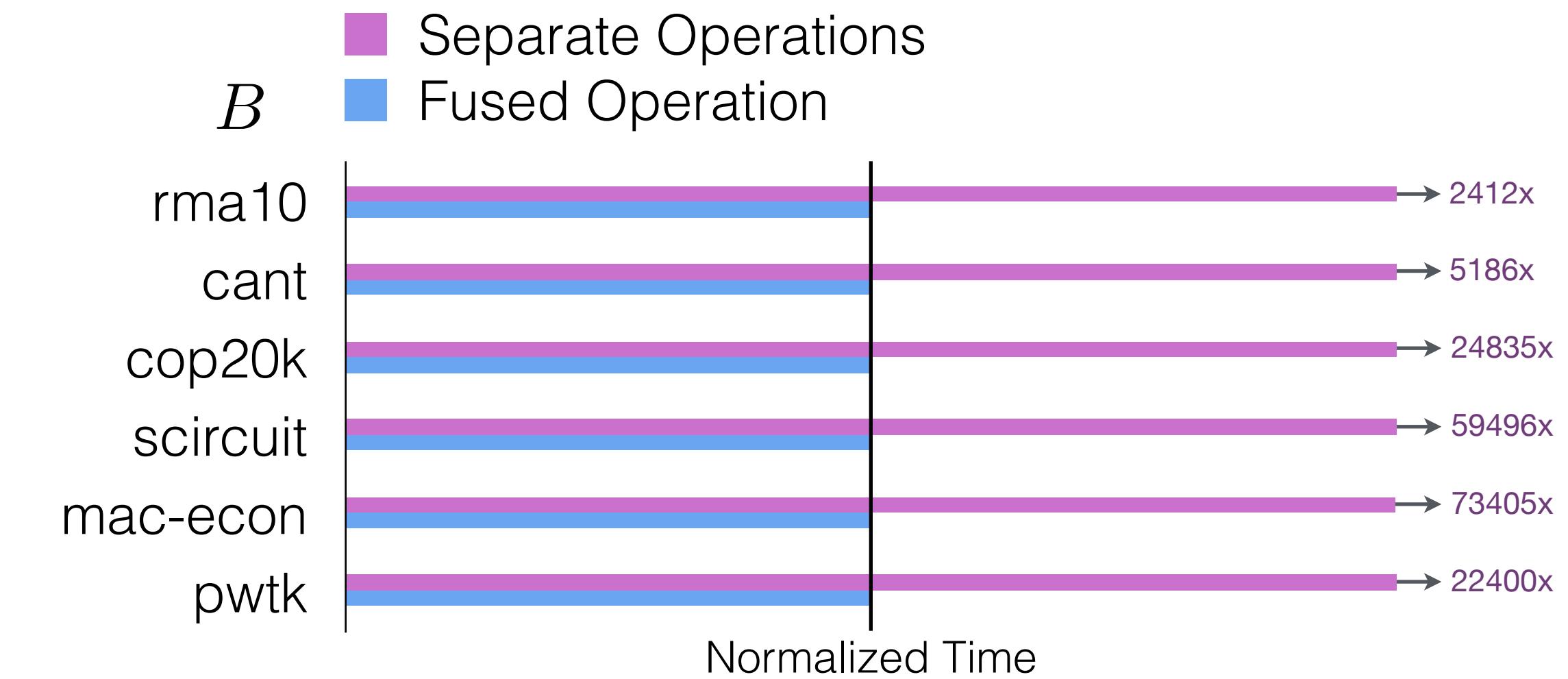
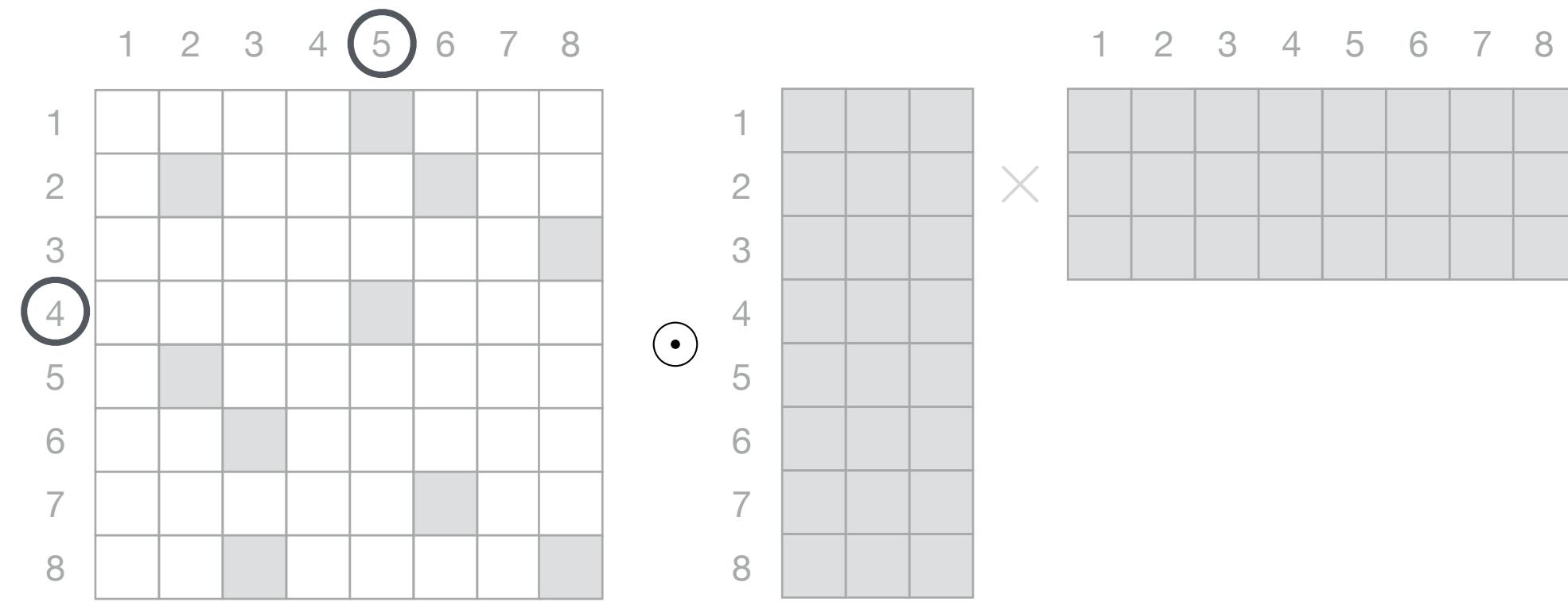
Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$



Example 1: Sampled Dense-Dense Matrix Multiplication with Linear Algebra

$$A = B \odot (CD)$$



Example 2: Triangle Query with Relational Algebra

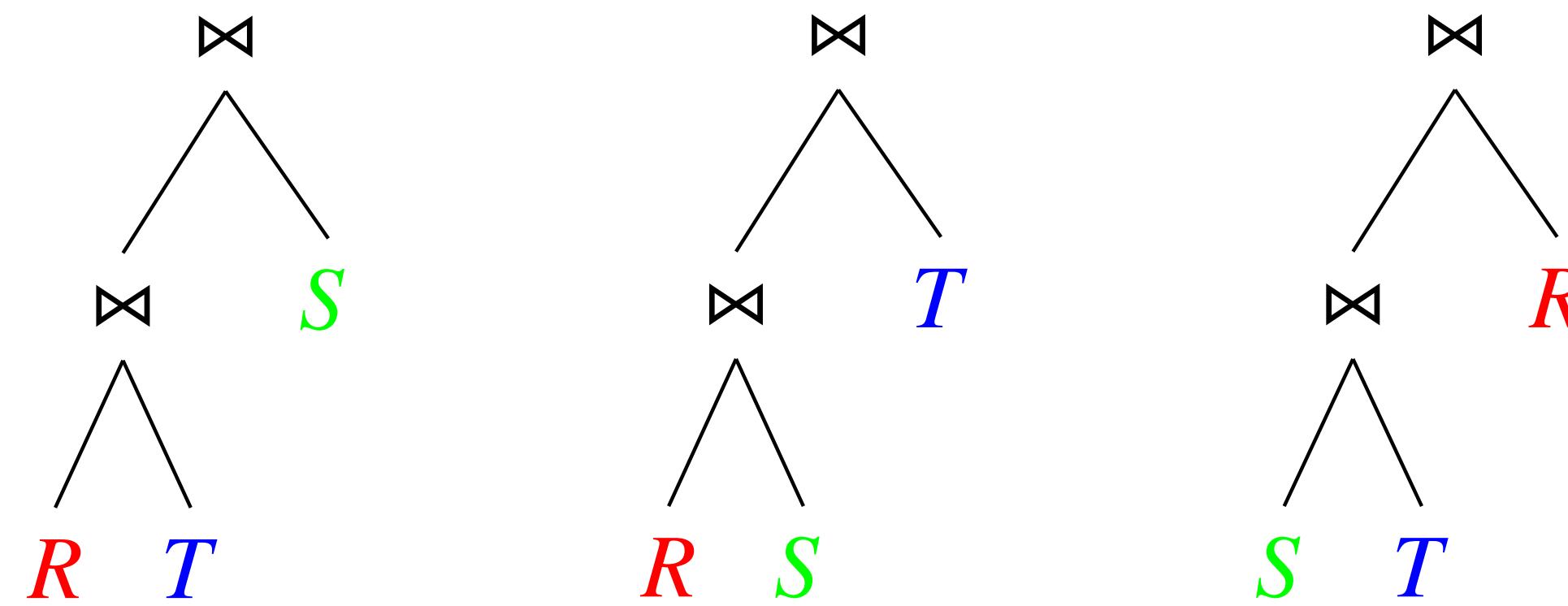
$$Q_{\Delta} = R(A, B) \bowtie S(B, C) \bowtie T(A, C)$$

Example 2: Triangle Query with Relational Algebra

$$Q_{\Delta} = R(A, B) \bowtie S(B, C) \bowtie T(A, C) \quad O(N^{3/2})$$

Example 2: Triangle Query with Relational Algebra

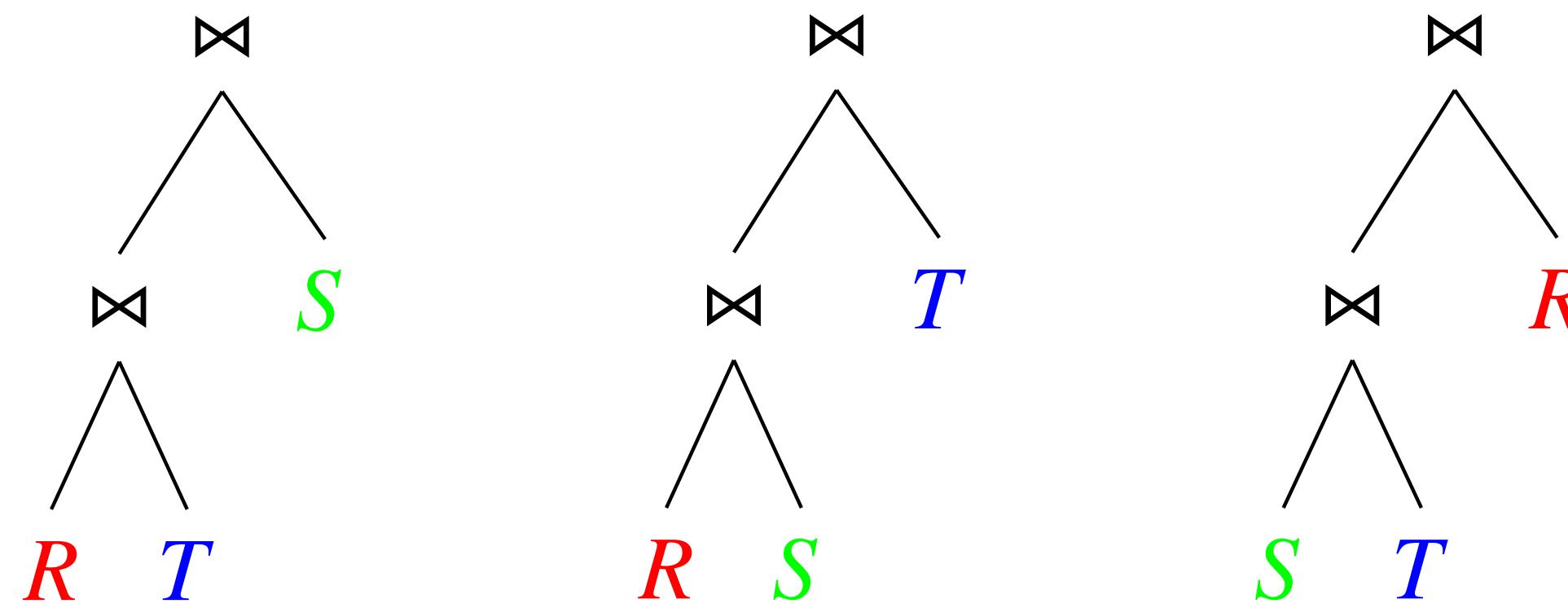
$$Q_{\Delta} = R(A, B) \bowtie S(B, C) \bowtie T(A, C) \quad O(N^{3/2})$$



$O(N^2)$

Example 2: Triangle Query with Relational Algebra

$$Q_{\Delta} = R(A, B) \bowtie S(B, C) \bowtie T(A, C) \quad O(N^{3/2})$$



Algorithm 2 Computing Q_{Δ} by delaying computation.

Input: $R(A, B), S(B, C), T(A, C)$ in sorted order

- 1: $Q \leftarrow \emptyset$
 - 2: $L_A \leftarrow \pi_A(R) \cap \pi_A(T)$
 - 3: **For** each $a \in L_A$ **do**
 - 4: $L_B^a \leftarrow \pi_B(\sigma_{A=a}(R)) \cap \pi_B(S)$
 - 5: **For** each $b \in L_B^a$ **do**
 - 6: $L_C^{a,b} \leftarrow \pi_C(\sigma_{B=b}(S)) \cap \pi_C(\sigma_{A=a}(T))$
 - 7: **For** each $c \in L_C^{a,b}$ **do**
 - 8: Add (a, b, c) to Q
 - 9: **Return** Q
-

Figures from Ngo, Ré and Rudra (2013),
with algorithm from Veldhuizen (2014)

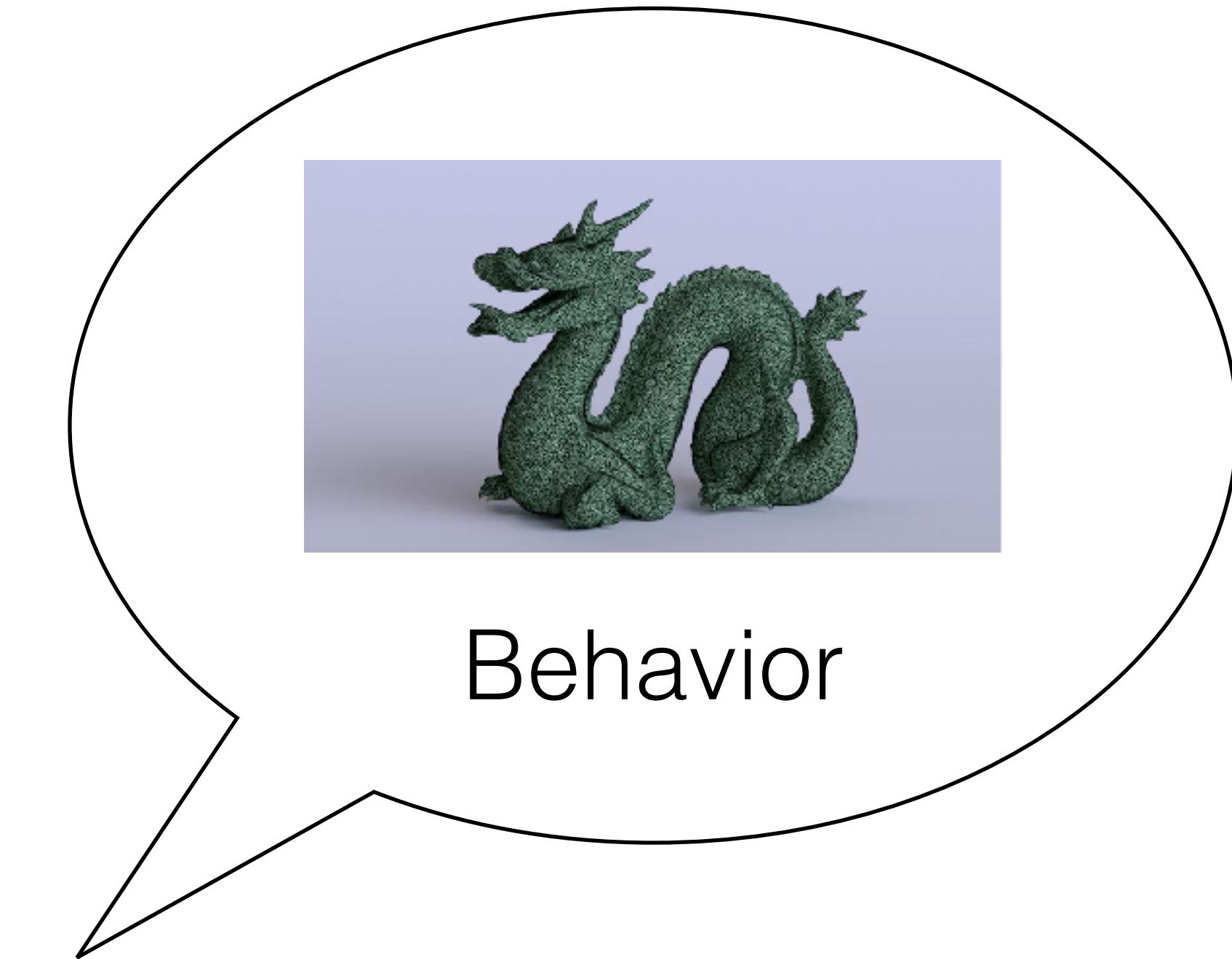
$O(N^2)$

$O(N^{3/2})$

Example 3: Simulation with Meshes and Linear Algebra

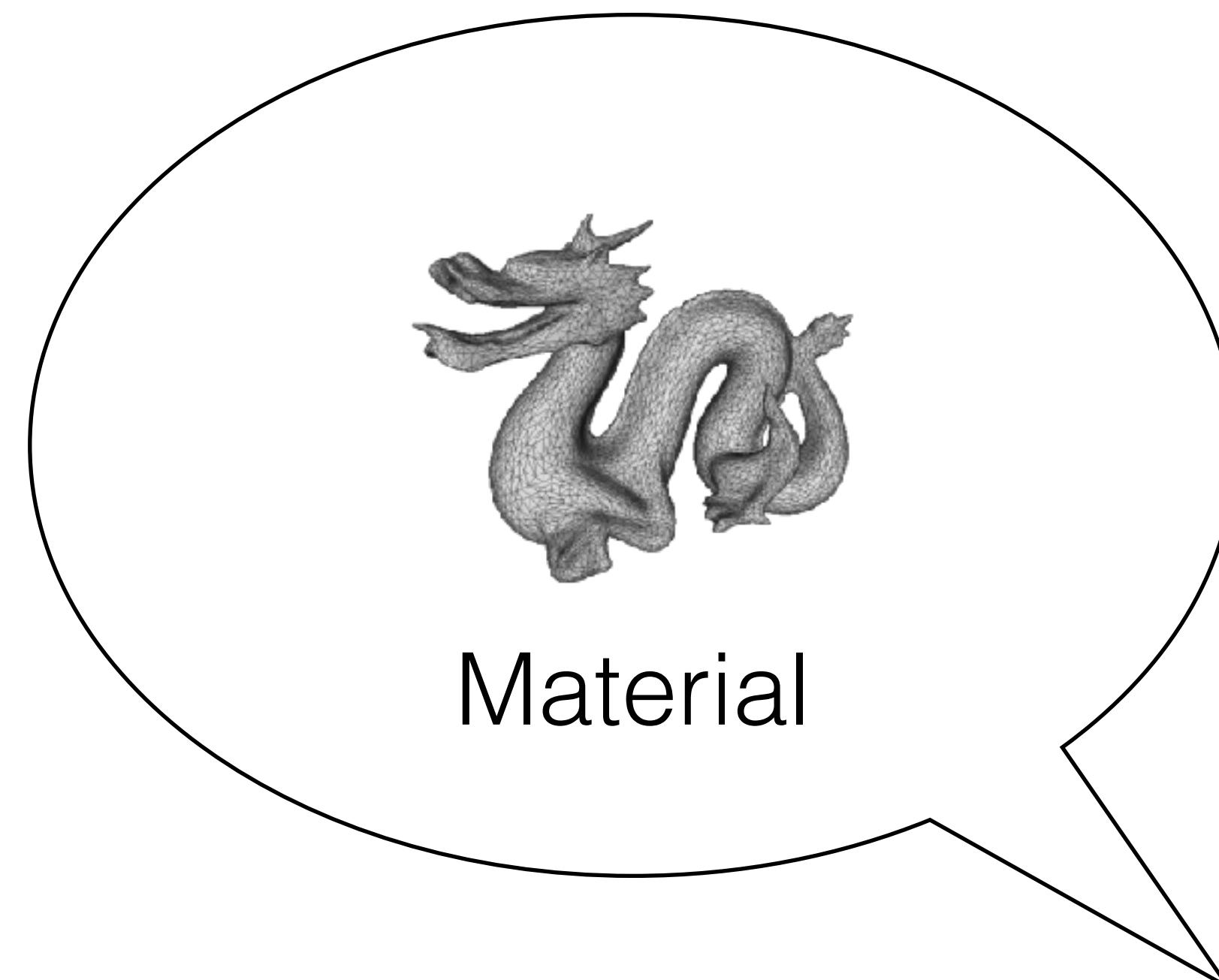


Example 3: Simulation with Meshes and Linear Algebra

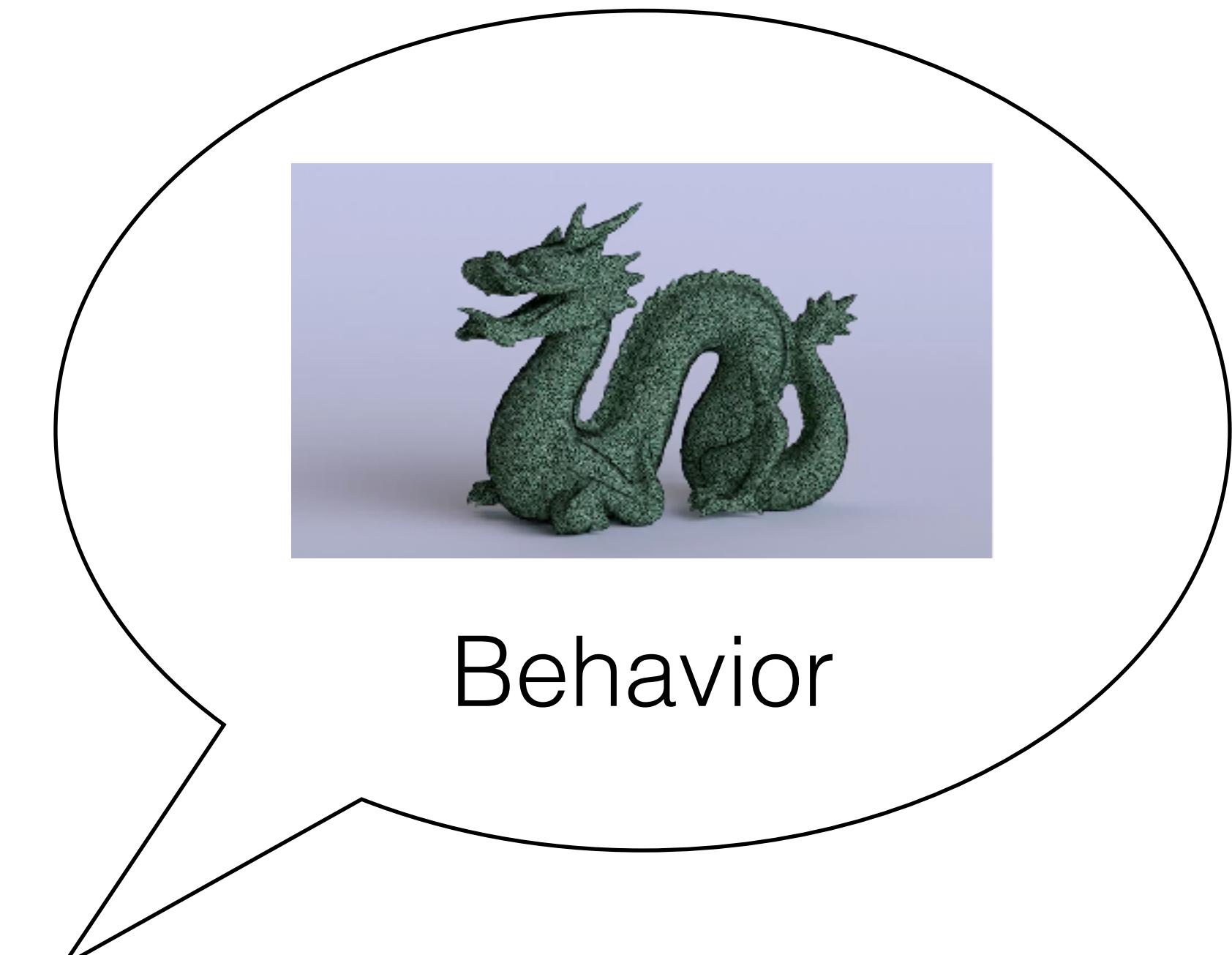


Behavior

Example 3: Simulation with Meshes and Linear Algebra



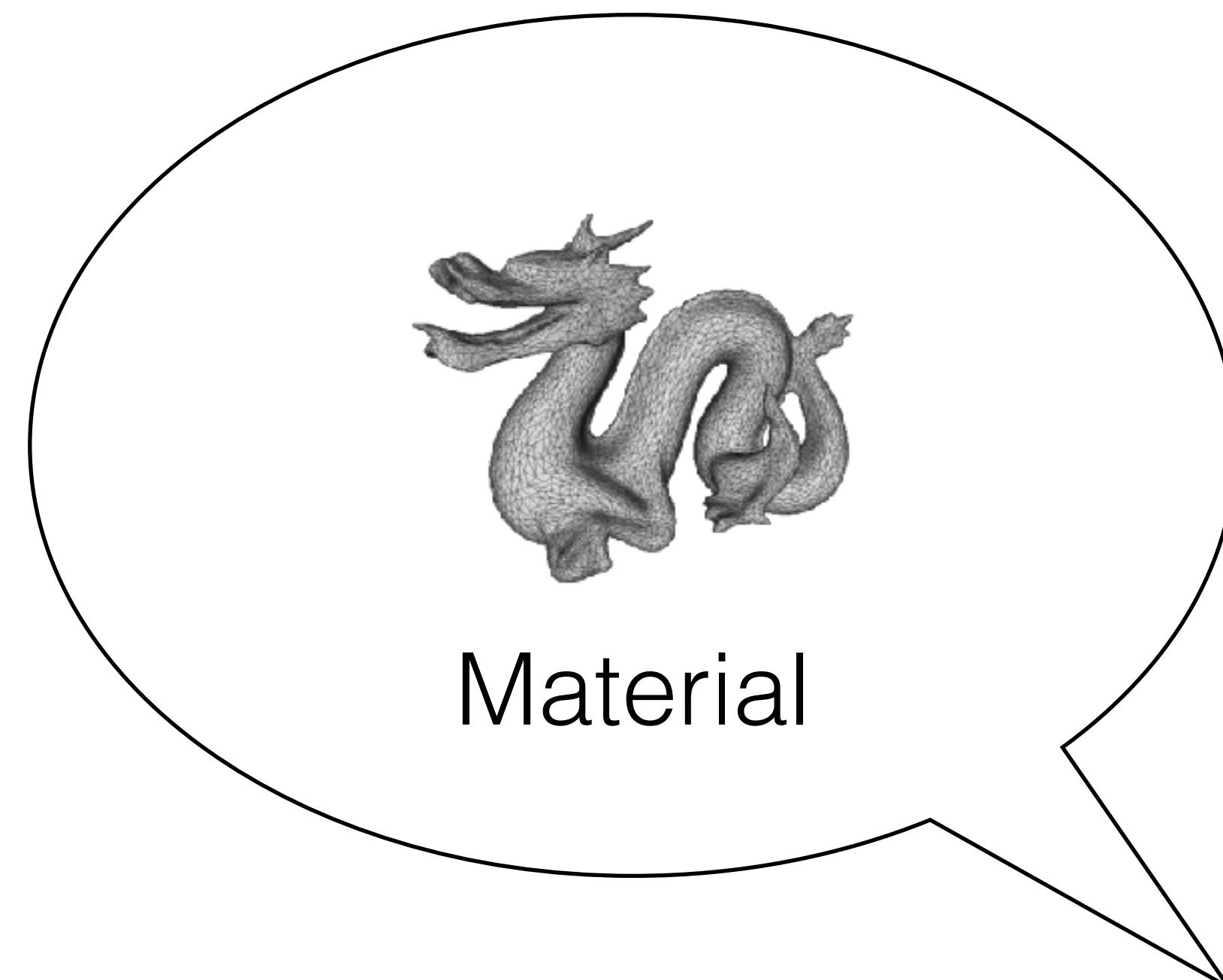
Material



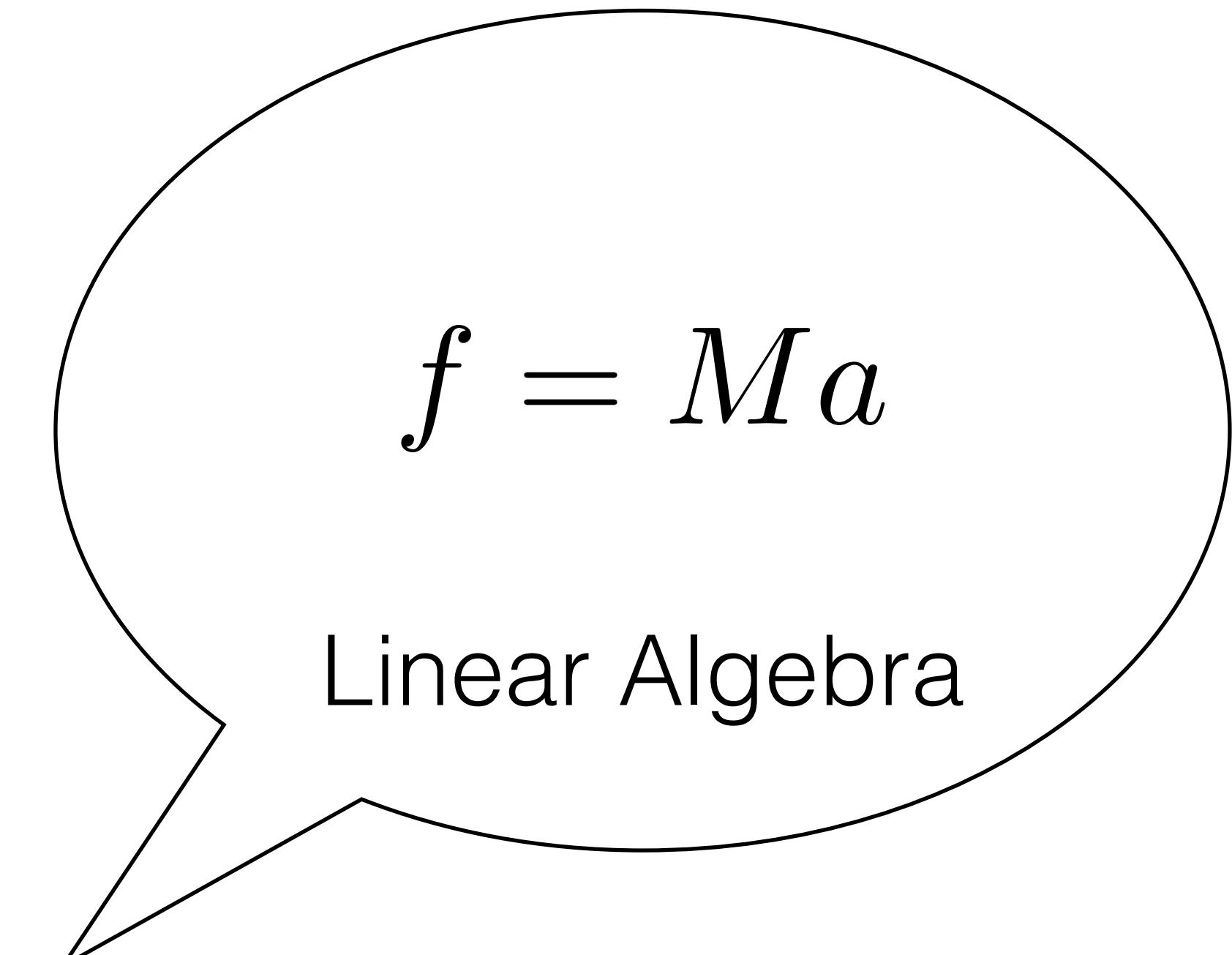
Behavior



Example 3: Simulation with Meshes and Linear Algebra



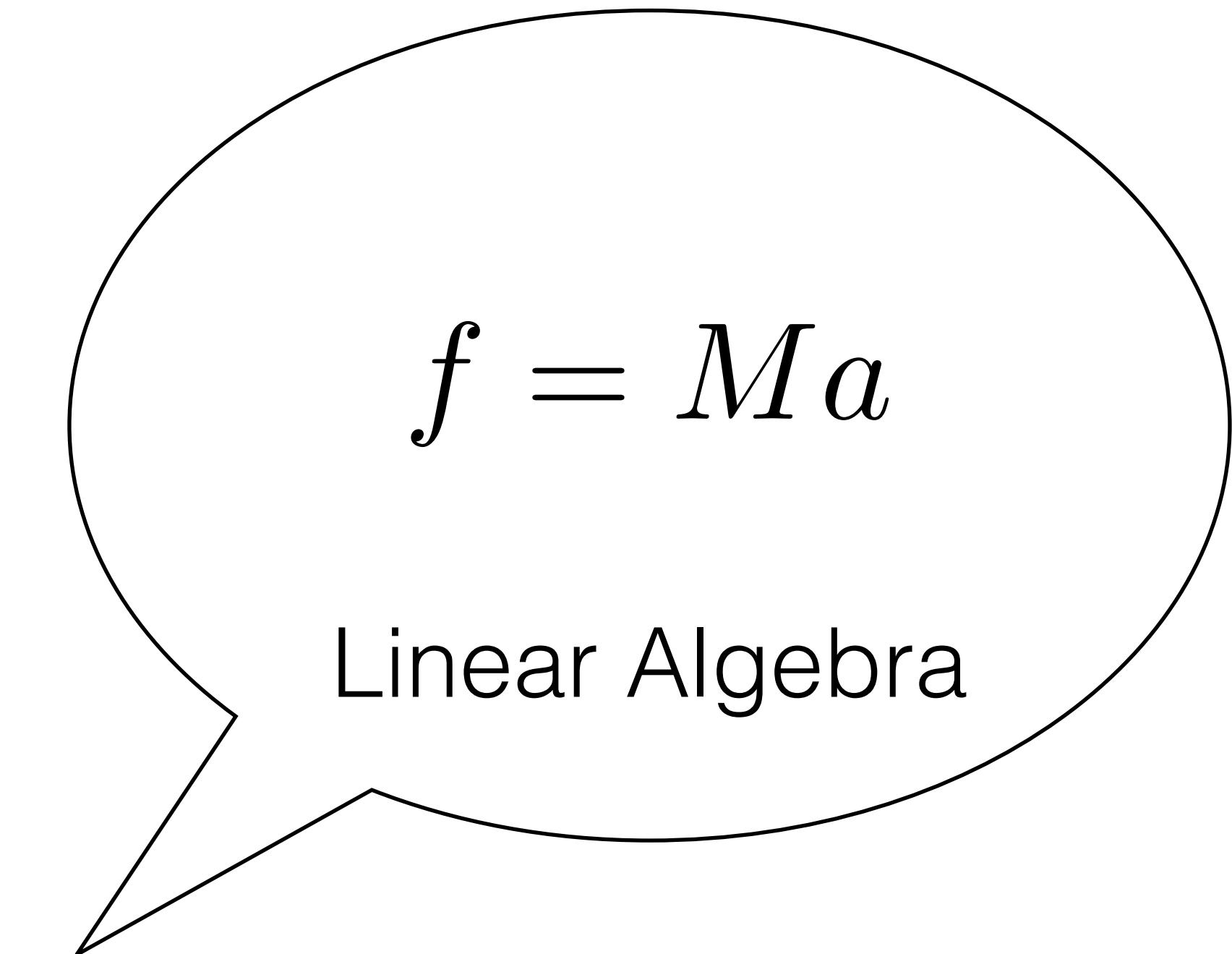
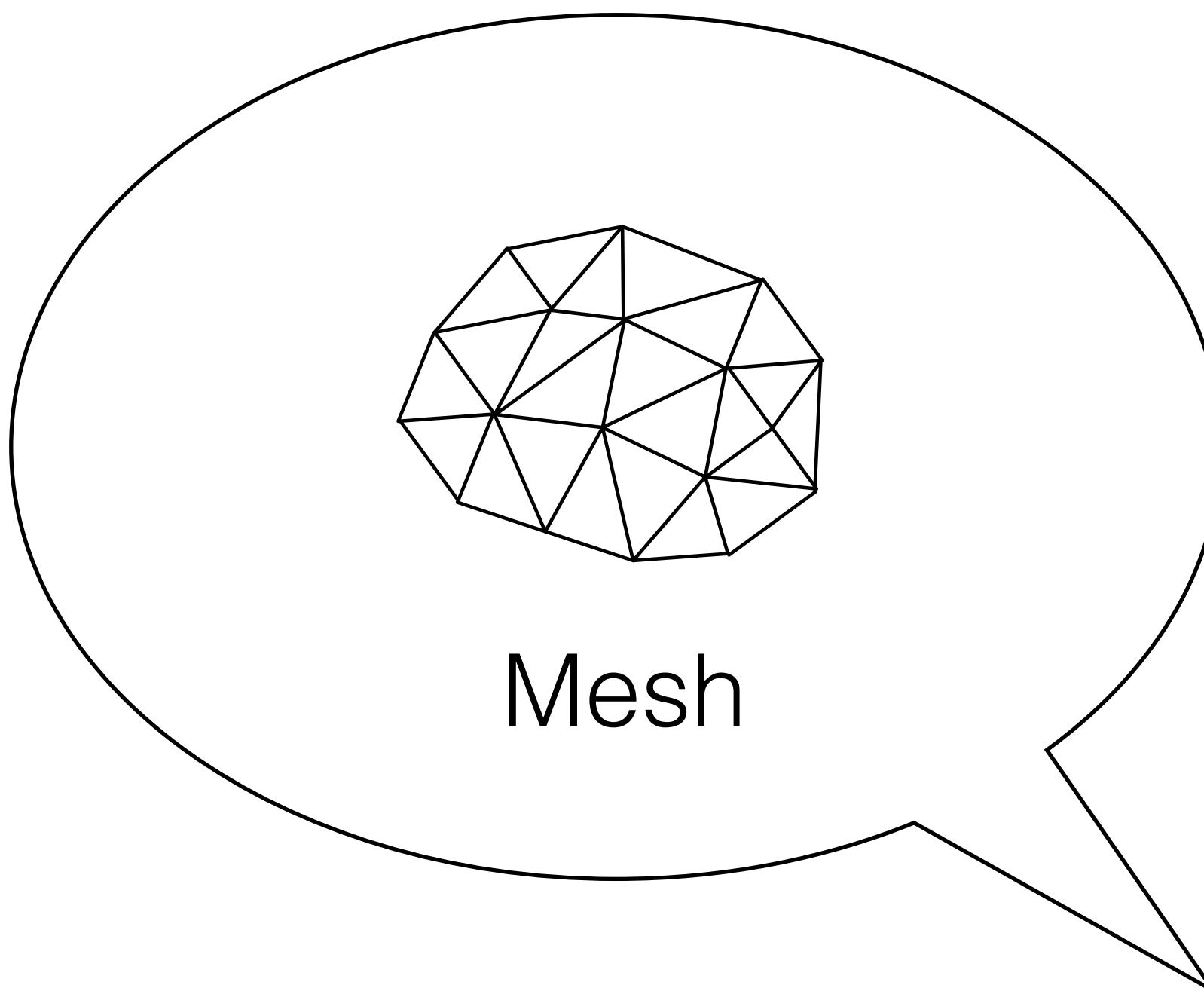
Material



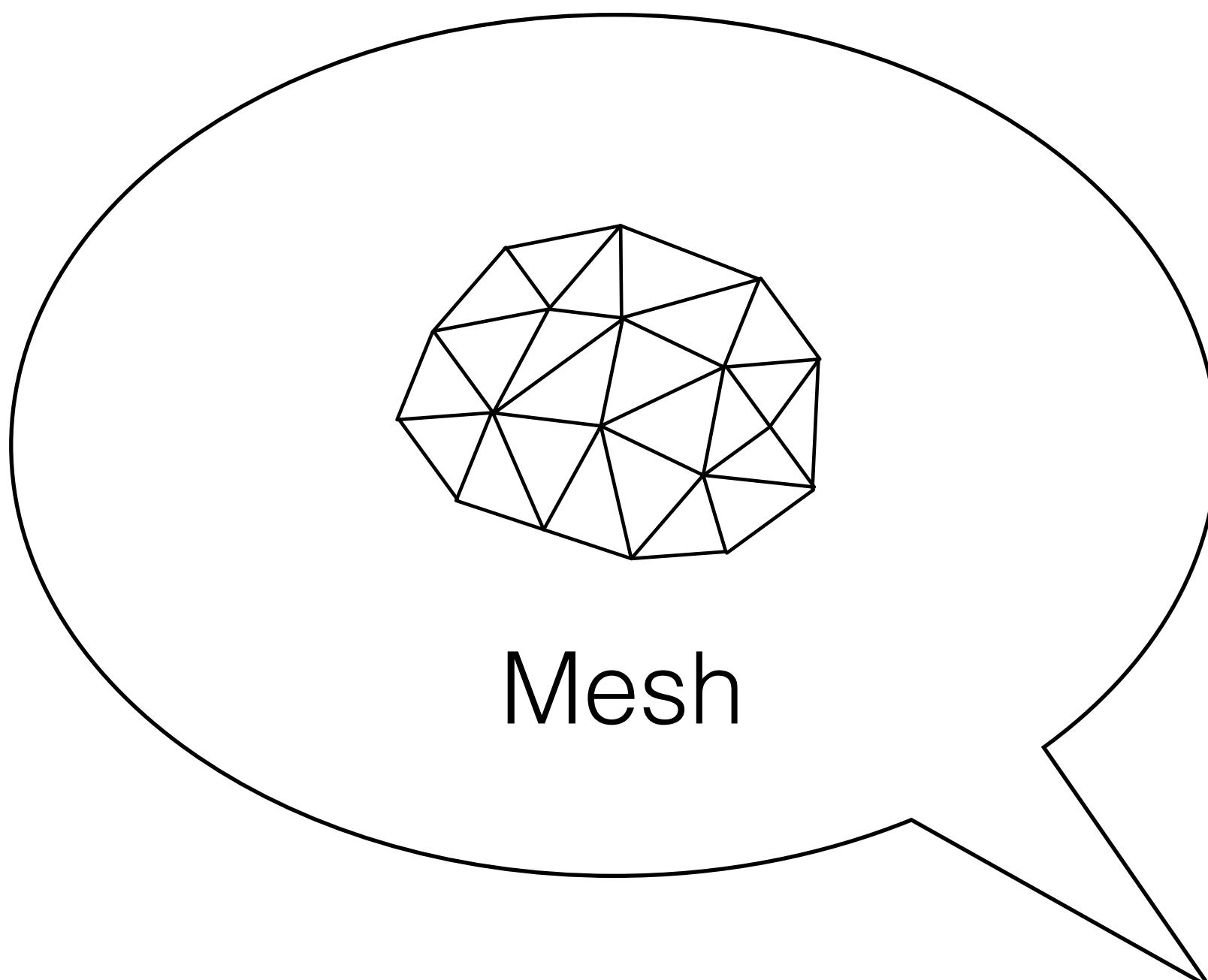
Linear Algebra



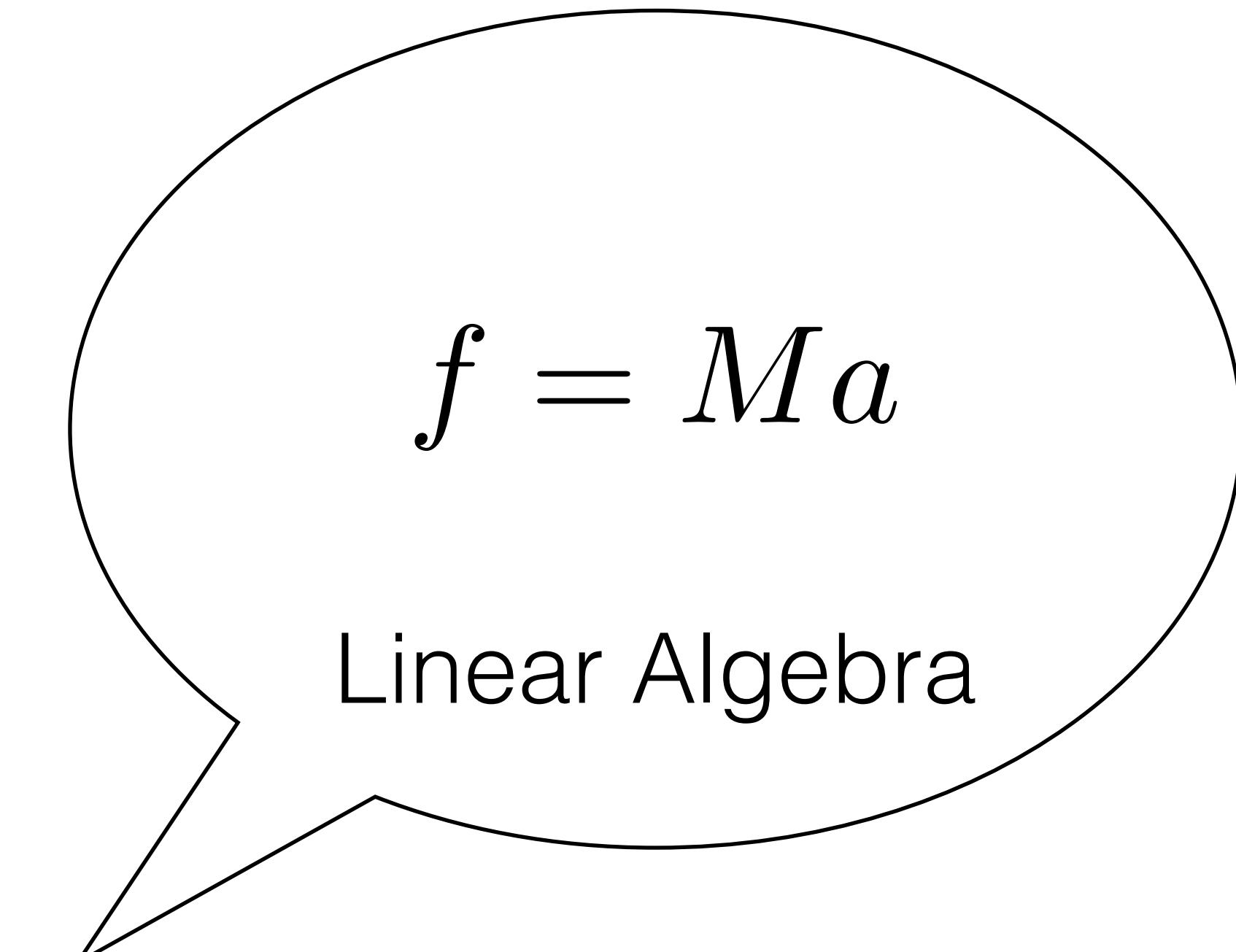
Example 3: Simulation with Meshes and Linear Algebra



Example 3: Simulation with Meshes and Linear Algebra



Mesh

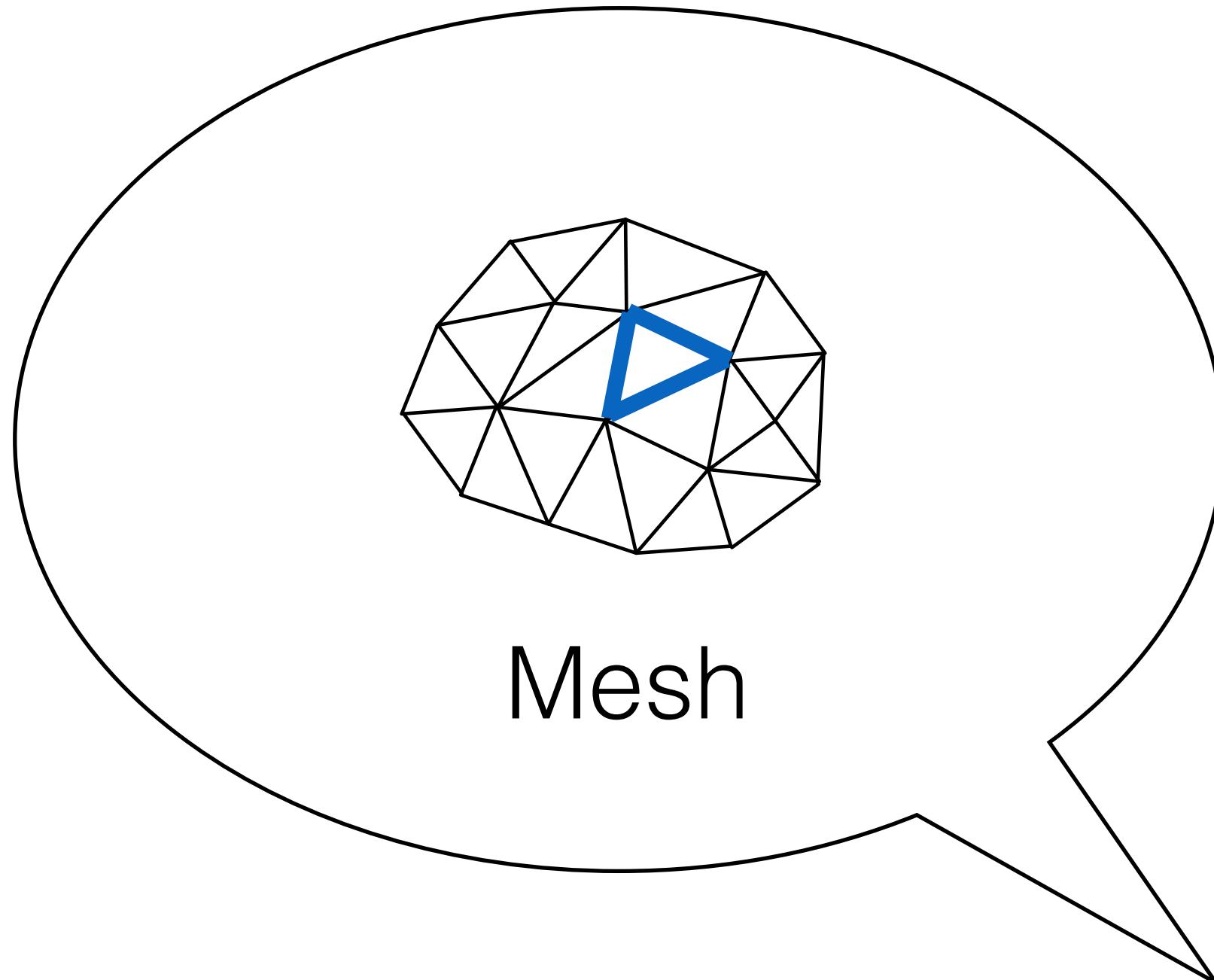


Linear Algebra



Example 3: Simulation with Meshes and Linear Algebra

Matrix-free in-place stencil computation

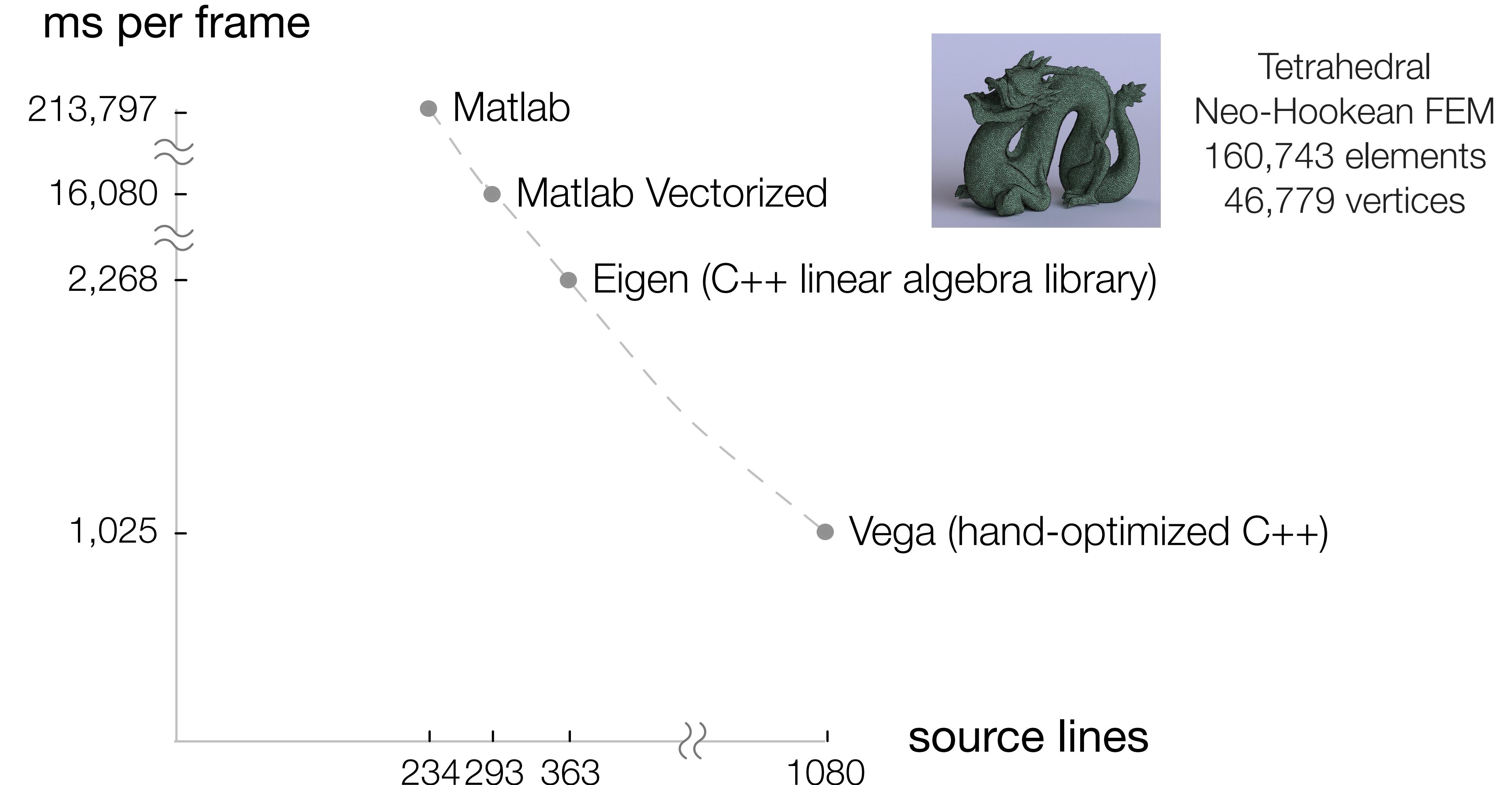


Mesh

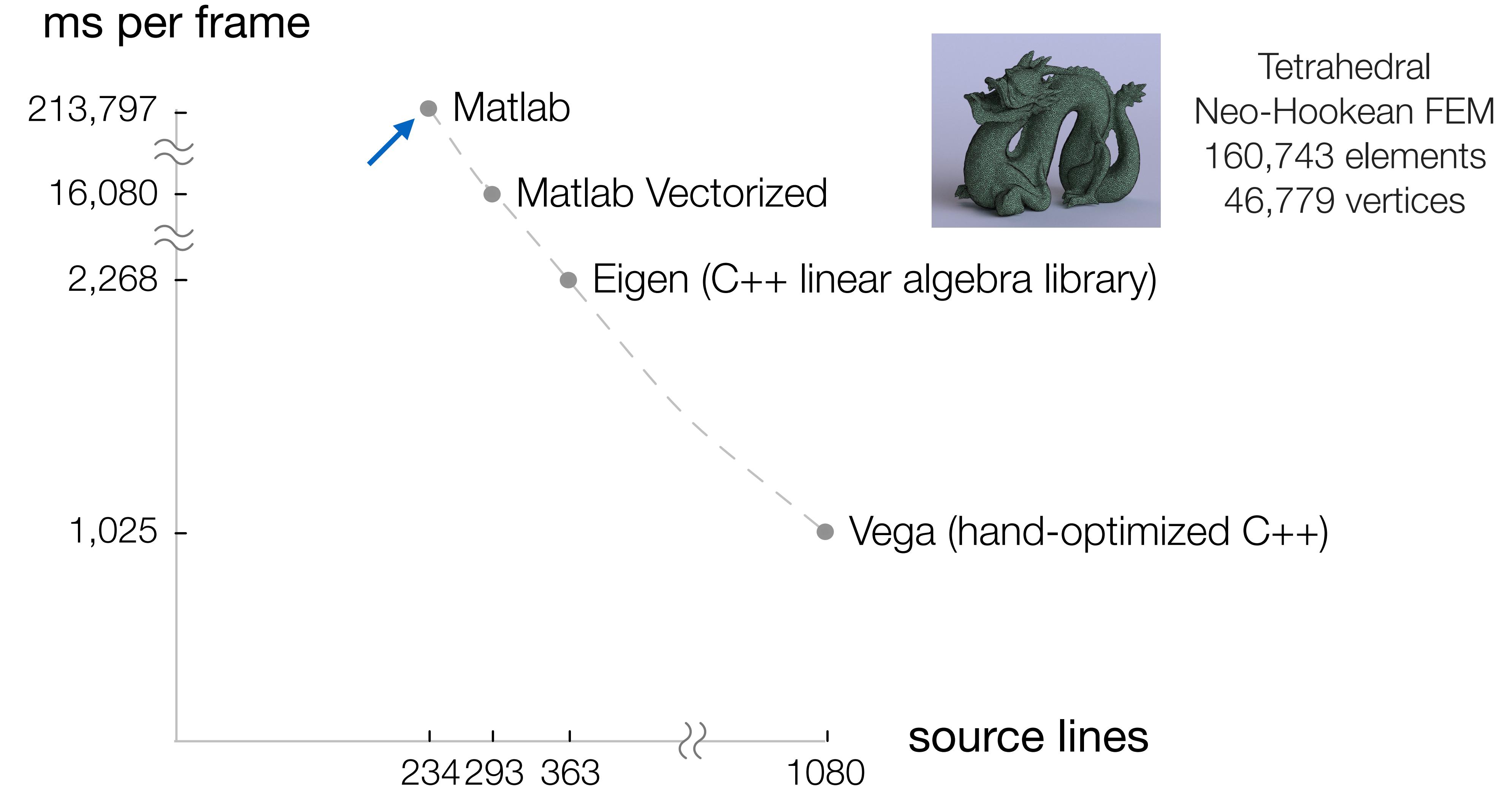


A detailed anatomical illustration of a human brain, showing the cerebral cortex with its characteristic gyri (ridges) and sulci (grooves). The brain is colored in shades of pink and red, set against a white background. To the right of the brain, there are four vertical columns of dense, multi-colored text, likely representing neural activity or connectivity data. The text is color-coded with various shades of green, blue, purple, and red, possibly indicating different functional regions or activity levels. The overall composition suggests a scientific or medical context, such as a neuroimaging study or a brain mapping project.

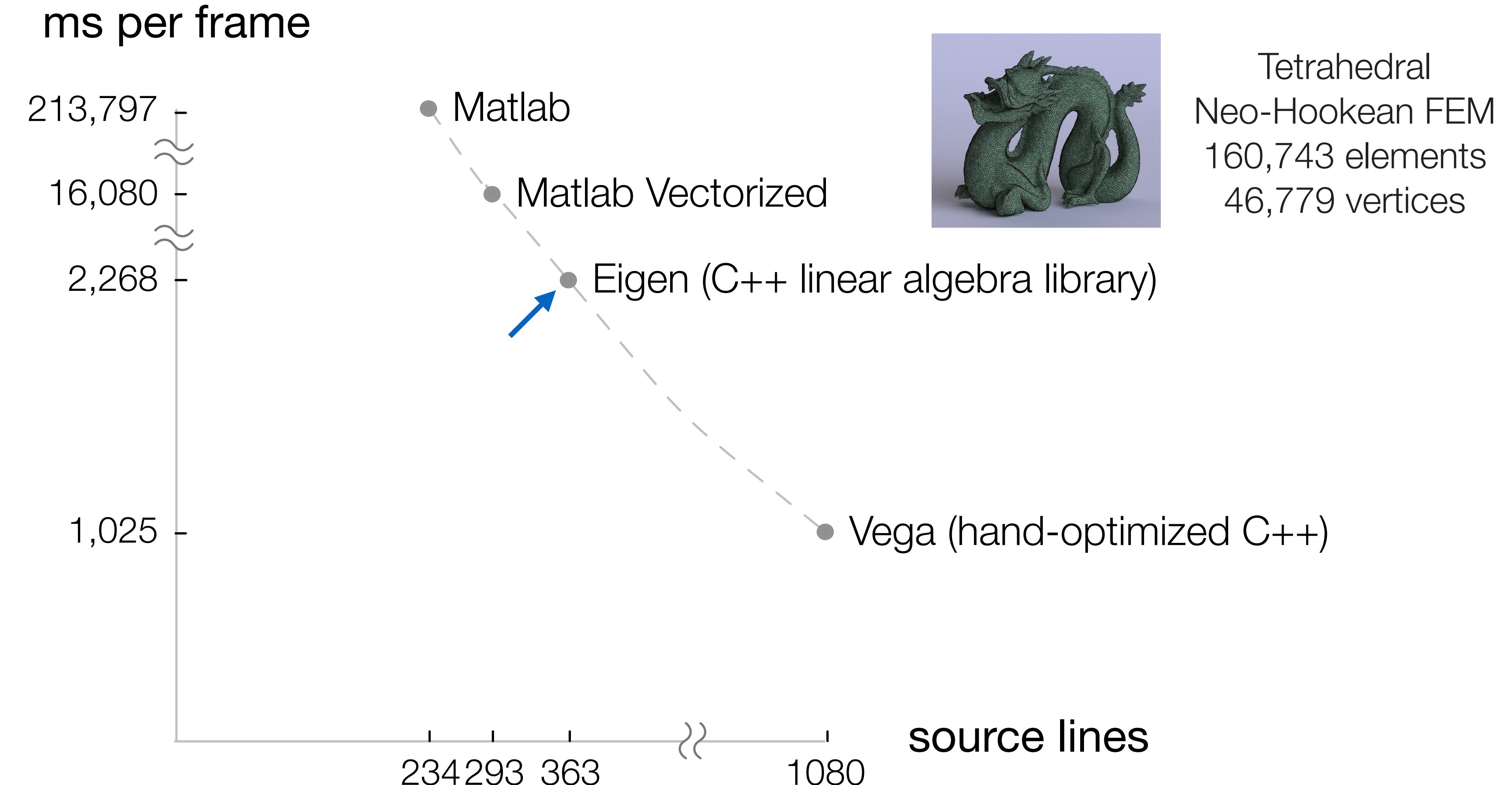
Example 3: Simulation with Graphs and Linear Algebra



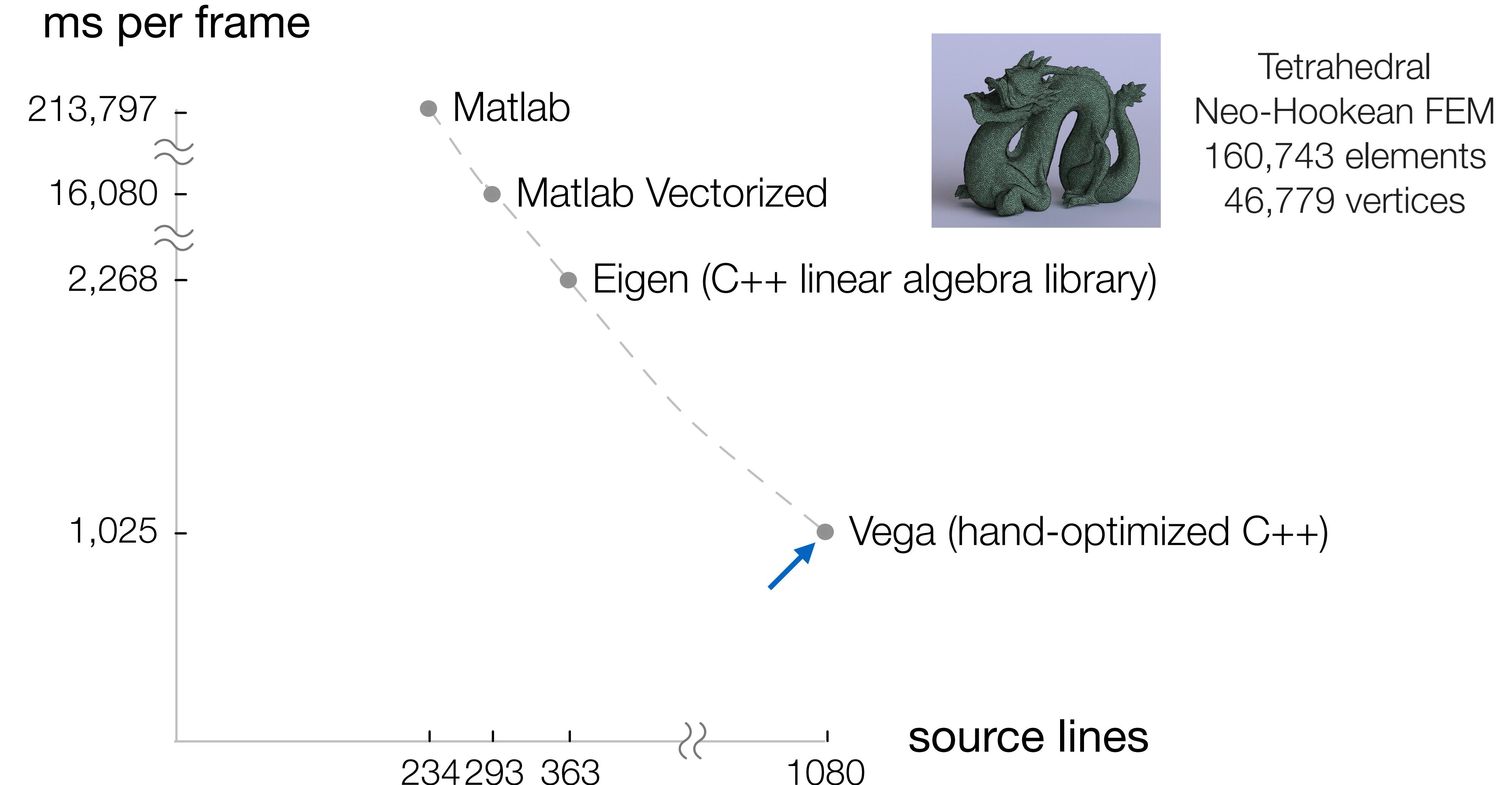
Example 3: Simulation with Graphs and Linear Algebra



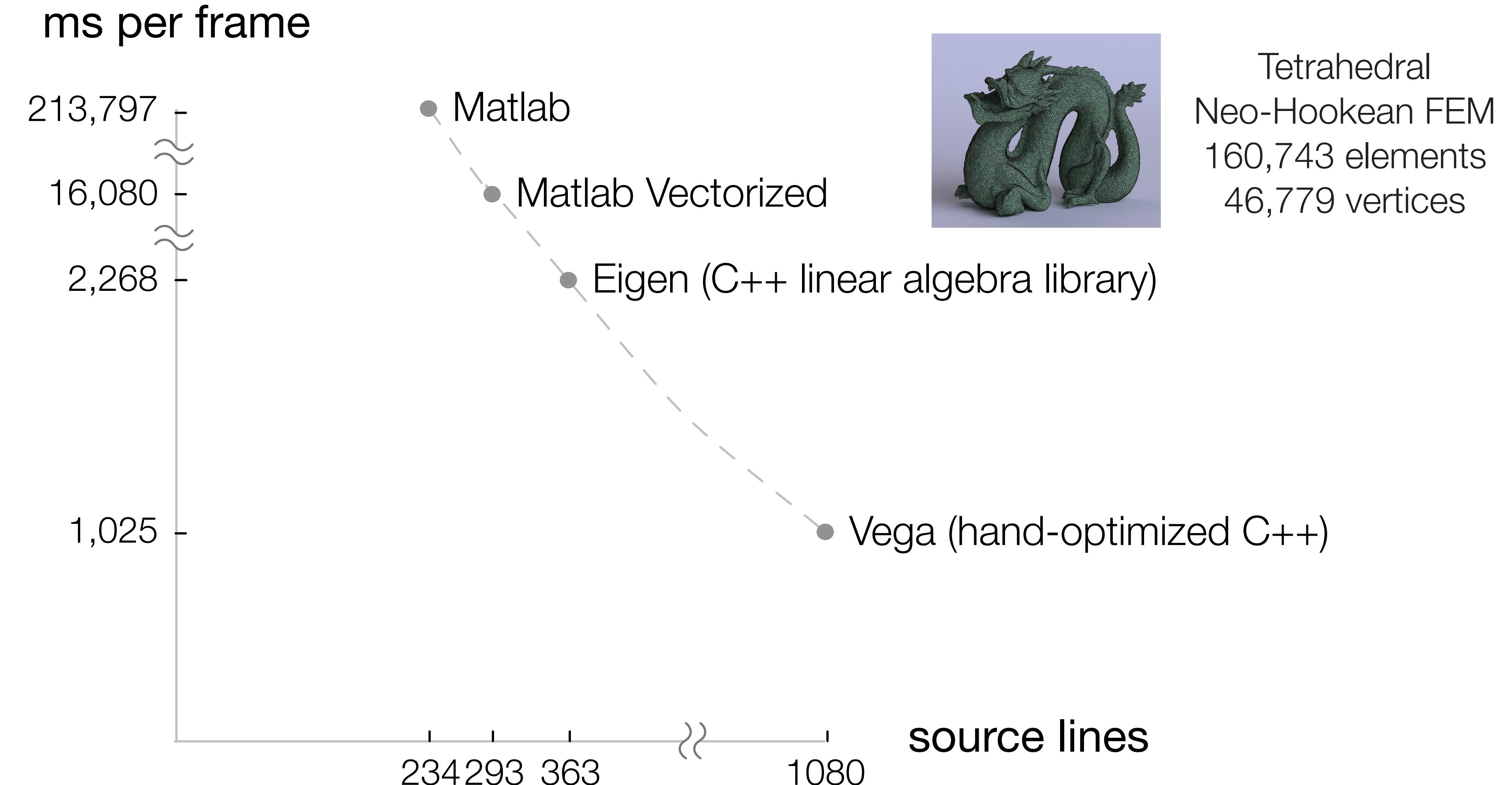
Example 3: Simulation with Graphs and Linear Algebra



Example 3: Simulation with Graphs and Linear Algebra



Example 3: Simulation with Graphs and Linear Algebra



Too many combinations for a fixed-function library

$$a = Bc$$

$$a = Bc + a$$

$$a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a$$

$$a = B^T c \quad A = \alpha B \quad a = B(c + d)$$

$$a = B^T c + d \quad A = B + C + D \quad A = BC$$

$$A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD)$$

$$A = BCd \quad A = B^T \quad a = B^T Bc$$

$$a = b + c \quad A = B \quad K = A^T CA$$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k$$

$$A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j$$

$$A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj}$$

$$\tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$$

$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}}$$

$$a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$$

Too many combinations for a fixed-function library

$$\begin{aligned} & a = Bc \\ & a = Bc + a \\ & a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\ & \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\ & a = B^T c + d \quad A = B + C + D \quad A = BC \\ & A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\ & \quad A = BCd \quad A = B^T \quad a = B^T Bc \\ & a = b + c \quad A = B \quad K = A^T CA \end{aligned}$$

Linear Algebra

$$\begin{aligned} A_{ij} &= \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\ A_{lj} &= \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\ A_{ijk} &= \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\ A_{jk} &= \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\ C &= \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\ a &= \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}} \end{aligned}$$

Too many combinations for a fixed-function library

$$a = Bc$$

$$a = Bc + a$$

$$a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a$$

$$a = B^T c \quad A = \alpha B \quad a = B(c + d)$$

$$a = B^T c + d \quad A = B + C + D \quad A = BC$$

$$A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD)$$

$$A = BCd \quad A = B^T \quad a = B^T Bc$$

$$a = b + c \quad A = B \quad K = A^T CA$$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k$$

$$A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j$$

$$A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj}$$

$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$$

$$a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$$

Data analytics
(tensor factorization)

Too many combinations for a fixed-function library

$$a = Bc$$

$$a = Bc + a$$

$$a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a$$

$$a = B^T c \quad A = \alpha B \quad a = B(c + d)$$

$$a = B^T c + d \quad A = B + C + D \quad A = BC$$

$$A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD)$$

$$A = BCd \quad A = B^T \quad a = B^T Bc$$

$$a = b + c \quad A = B \quad K = A^T CA$$

$$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$$

$$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k$$

$$A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j$$

$$A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj}$$

$$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$$

$$a = \sum_{ijklmno} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$$

Quantum Chromodynamics

Too many combinations for a fixed-function library

CSpars	Eigen (SpMV)			
	$a = Bc$			
$a = Bc + a$				
$a = Bc + b$	$A = B + C$	$a = \alpha Bc + \beta a$		
PETSc	$a = B^T c$	$A = \alpha B$	$a = B(c + d)$	
	$a = B^T c + d$	$A = B + C + D$	$A = BC$	
	$A = B \odot C$	$a = b \odot c$	$A = 0$	$A = B \odot (CD)$
	$A = BCd$	$A = B^T$	$a = B^T Bc$	
	$a = b + c$	$A = B$	$K = A^T CA$	
	$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj}$	$A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$		
	$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj}$	$A_{ij} = \sum_k B_{ijk} c_k$		
	$A_{ijk} = \sum_l B_{ikl} C_{lj}$	$A_{ik} = \sum_j B_{ijk} c_j$		
	$A_{jk} = \sum_i B_{ijk} c_i$	$A_{ijl} = \sum_k B_{ikl} C_{kj}$		
	$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}}$	$\tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$		
	$a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$			

Too many combinations for a fixed-function library

CSpars	Eigen (SpMV)			
	$a = Bc$			
$a = Bc + a$				
$a = Bc + b$	$A = B + C$	$a = \alpha Bc + \beta a$		
PETSc	$a = B^T c$	$A = \alpha B$	$a = B(c + d)$	
	$a = B^T c + d$	$A = B + C + D$	$A = BC$	
	$A = B \odot C$	$a = b \odot c$	$A = 0$	$A = B \odot (CD)$
		$A = BCd$	$A = B^T$	$a = B^T Bc$
	$a = b + c$	$A = B$	$K = A^T CA$	
		$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj}$	$A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$	
		$A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj}$	$A_{ij} = \sum_k B_{ijk} c_k$	
		$A_{ijk} = \sum_l B_{ikl} C_{lj}$	$A_{ik} = \sum_j B_{ijk} c_j$	
		$A_{jk} = \sum_i B_{ijk} c_i$	$A_{ijl} = \sum_k B_{ikl} C_{kj}$	
		$C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}}$	$\tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$	
		$a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}$		

OSKI has 282 specialized variants of this expression

Too many combinations for a fixed-function library

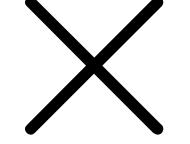
$$\begin{aligned}
& a = Bc \\
& a = Bc + a \\
& a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\
& \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\
& a = B^T c + d \quad A = B + C + D \quad A = BC \\
& A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\
& \quad A = BCd \quad A = B^T \quad a = B^T Bc \\
& a = b + c \quad A = B \quad K = A^T CA
\end{aligned}$$

×

$$\begin{aligned}
A_{ij} &= \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\
A_{lj} &= \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\
A_{ijk} &= \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\
A_{jk} &= \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\
C &= \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\
a &= \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} \overline{P_{no}} \overline{M_{po}} \overline{P_{ip}}
\end{aligned}$$

Dense Matrix		
CSR	DCSR	BCSR
COO	ELLPACK	CSB
Blocked COO		CSC
DIA	Blocked DIA	DCSC

Too many combinations for a fixed-function library

$a = Bc$ $a = Bc + a$ $a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a$ $a = B^T c \quad A = \alpha B \quad a = B(c + d)$ $a = B^T c + d \quad A = B + C + D \quad A = BC$ $A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD)$ $A = BCd \quad A = B^T \quad a = B^T Bc$ $a = b + c \quad A = B \quad K = A^T CA$		<p style="margin: 0;">Dense Matrix</p> <table style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 33%;">CSR</th> <th style="width: 33%;">DCSR</th> <th style="width: 33%;">BCSR</th> </tr> </thead> <tbody> <tr> <td style="text-align: center;">COO</td> <td style="text-align: center;">ELLPACK</td> <td style="text-align: center;">CSB</td> </tr> <tr> <td style="text-align: center;">Blocked COO</td> <td></td> <td style="text-align: center;">CSC</td> </tr> <tr> <td style="text-align: center;">DIA</td> <td style="text-align: center;">Blocked DIA</td> <td style="text-align: center;">DCSC</td> </tr> </tbody> </table> <p style="margin: 0;">Thermal Simulation</p>	CSR	DCSR	BCSR	COO	ELLPACK	CSB	Blocked COO		CSC	DIA	Blocked DIA	DCSC
CSR	DCSR	BCSR												
COO	ELLPACK	CSB												
Blocked COO		CSC												
DIA	Blocked DIA	DCSC												

$A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij}$
 $A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k$
 $A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j$
 $A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj}$
 $C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik})$
 $a = \sum_{ijklmnop} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}}$

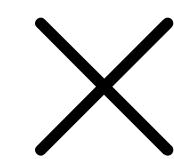
Too many combinations for a fixed-function library

$$\begin{aligned}
 & a = Bc \\
 & a = Bc + a \\
 & a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\
 & \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\
 & a = B^T c + d \quad A = B + C + D \quad A = BC \\
 & A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\
 & \quad A = BCd \quad A = B^T \quad a = B^T Bc \\
 & a = b + c \quad A = B \quad K = A^T CA \\
 & A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\
 & A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\
 & A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\
 & A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\
 & C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\
 & a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}}
 \end{aligned}$$



Too many combinations for a fixed-function library

$$\begin{aligned} & a = Bc \\ & a = Bc + a \\ & a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\ & \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\ & a = B^T c + d \quad A = B + C + D \quad A = BC \\ & A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\ & \quad A = BCd \quad A = B^T \quad a = B^T Bc \\ & a = b + c \quad A = B \quad K = A^T CA \\ & \\ & A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\ & \quad A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\ & A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\ & \quad A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\ & \\ & C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\ & \quad a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}} \end{aligned}$$



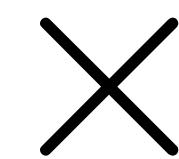
Dense Matrix

CSR	DCSR	BCSR
COO	ELLPACK	CSB
Blocked COO	CSC	
DIA	Blocked DIA	DCSC

Sparse vector Hash Maps

Too many combinations for a fixed-function library

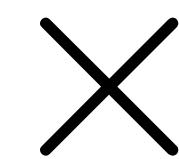
$$\begin{aligned} & a = Bc \\ & a = Bc + a \\ & a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\ & \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\ & a = B^T c + d \quad A = B + C + D \quad A = BC \\ & A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\ & \quad A = BCd \quad A = B^T \quad a = B^T Bc \\ & a = b + c \quad A = B \quad K = A^T CA \\ & \\ & A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\ & \quad A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\ & A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\ & \quad A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\ & \\ & C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\ & \quad a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}} \end{aligned}$$



Dense Matrix
CSR DCSR BCSR
COO ELLPACK CSB
Blocked COO CSC
DIA Blocked DIA DCSC
Sparse vector Hash Maps
Coordinates
CSF Dense Tensors
 Blocked Tensors

Too many combinations for a fixed-function library

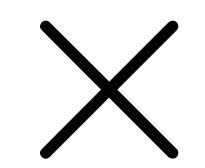
$$\begin{aligned} & a = Bc \\ & a = Bc + a \\ & a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\ & \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\ & a = B^T c + d \quad A = B + C + D \quad A = BC \\ & A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\ & \quad A = BCd \quad A = B^T \quad a = B^T Bc \\ & a = b + c \quad A = B \quad K = A^T CA \\ & \\ & A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\ & \quad A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\ & A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\ & \quad A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\ & \\ & C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\ & \quad a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}} \end{aligned}$$



Dense Matrix
CSR DCSR BCSR
COO ELLPACK CSB
Blocked COO CSC
DIA Blocked DIA DCSC
Sparse vector Hash Maps
Coordinates
CSF Dense Tensors
 Blocked Tensors
Linked Lists Database
Compression Schemes
Cloud Storage

Too many combinations for a fixed-function library

$$\begin{aligned}
 & a = Bc \\
 & a = Bc + a \\
 & a = Bc + b \quad A = B + C \quad a = \alpha Bc + \beta a \\
 & \quad a = B^T c \quad A = \alpha B \quad a = B(c + d) \\
 & a = B^T c + d \quad A = B + C + D \quad A = BC \\
 & A = B \odot C \quad a = b \odot c \quad A = 0 \quad A = B \odot (CD) \\
 & \quad A = BCd \quad A = B^T \quad a = B^T Bc \\
 & a = b + c \quad A = B \quad K = A^T CA \\
 & A_{ij} = \sum_{kl} B_{ikl} C_{lj} D_{kj} \quad A_{kj} = \sum_{il} B_{ikl} C_{lj} D_{ij} \\
 & \quad A_{lj} = \sum_{ik} B_{ikl} C_{ij} D_{kj} \quad A_{ij} = \sum_k B_{ijk} c_k \\
 & A_{ijk} = \sum_l B_{ikl} C_{lj} \quad A_{ik} = \sum_j B_{ijk} c_j \\
 & \quad A_{jk} = \sum_i B_{ijk} c_i \quad A_{ijl} = \sum_k B_{ikl} C_{kj} \\
 & C = \sum_{ijkl} M_{ij} P_{jk} \overline{M_{lk}} \overline{P_{il}} \quad \tau = \sum_i z_i (\sum_j z_j \theta_{ij}) (\sum_k z_k \theta_{ik}) \\
 & \quad a = \sum_{ijklmноп} M_{ij} P_{jk} M_{kl} P_{lm} \overline{M_{nm}} P_{no} \overline{M_{po}} \overline{P_{ip}}
 \end{aligned}$$



Dense Matrix				
CSR	DCSR	BCSR		
COO	ELLPACK	CSB	CPU	
Blocked COO		CSC	GPUs	TPUs
DIA	Blocked DIA	DCSC	FPGA	
Sparse vector		Hash Maps		Sparse Tensor Hardware
Coordinates				
CSF		Dense Tensors		Cloud Computers
		Blocked Tensors		
Linked Lists		Database		
				Supercomputers
Compression Schemes				
		Cloud Storage		

Optimized code is often complex, especially when it iterates over irregular data structures

$$A_{ij} = \sum_k B_{ijk} c_k$$

Optimized code is often complex, especially when it iterates over irregular data structures

$$A_{ij} = \sum_k B_{ijk} c_k$$


The diagram illustrates a summation operation in a formula. The term B_{ijk} is highlighted with a blue arrow pointing to its subscript k , labeled "dense". Similarly, the term c_k is highlighted with a blue arrow pointing to its subscript k , also labeled "dense". This visual emphasizes that both the matrix B and the vector c are dense, which can lead to complex memory access patterns in optimized code.

Optimized code is often complex, especially when it iterates over irregular data structures

$$A_{ij} = \sum_k B_{ijk} c_k$$


dense dense

```
for (int i = 0; i < m; i++) {  
  
    for (int j = 0; j < n; j++) {  
        int pB2 = i*n + j;  
        int pA2 = i*n + j;  
        double t = 0.0;  
        for (int k = 0; k < o; k++) {  
            int pB3 = pB2*o + k;  
            t += B[pB3] * c[k];  
        }  
        A[pA2] = t;  
    }  
}
```

Optimized code is often complex, especially when it iterates over irregular data structures

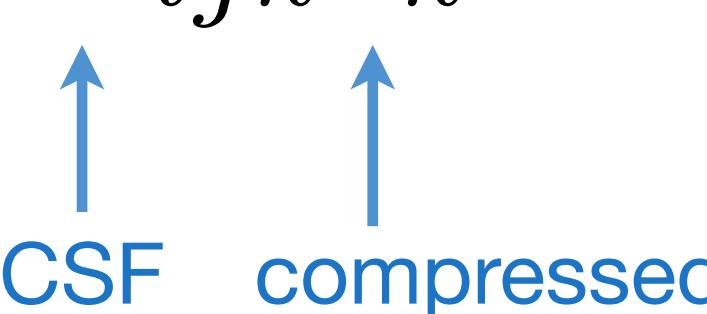
$$A_{ij} = \sum_k B_{ijk} c_k$$


CSF dense

```
for (int pA = 0; pA < m*n; pA++) {
    A[pA] = 0.0;
}
for (int pB1 = B1_pos[0]; pB1 < B1_pos[1]; pB1++) {
    int i = B1_crd[pB1];
    for (int pB2 = B2_pos[pB1]; pB2 < B2_pos[pB1+1]; pB2++) {
        int j = B2_crd[pB2];
        int pA2 = i*n + j;
        double t = 0.0;
        for (int pB3 = B3_pos[pB2]; pB3 < B3_pos[pB2+1]; pB3++) {
            int k = B3_crd[pB3];
            t += B[pB3] * c[k];
        }
        A[pA2] = t;
    }
}
```

Optimized code is often complex, especially when it iterates over irregular data structures

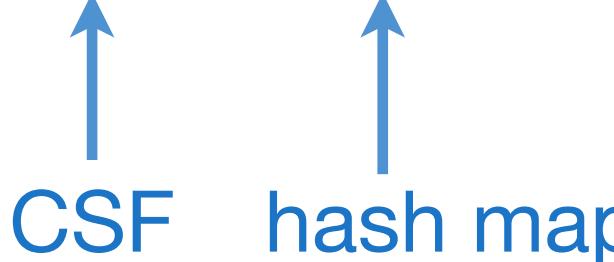
$$A_{ij} = \sum_k B_{ijk} c_k$$

CSF compressed

```
for (int pA = 0; pA < m*n; pA++) {
    A[pA] = 0.0;
}
for (int pB1 = B1_pos[0]; pB1 < B1_pos[1]; pB1++) {
    int i = B1_crd[pB1];
    for (int pB2 = B2_pos[pB1]; pB2 < B2_pos[pB1+1]; pB2++) {
        int j = B2_crd[pB2];
        int pA2 = i*n + j;
        double t = 0.0;
        int pB3 = B3_pos[pB2];
        int pc1 = c1_pos[0];
        while (pB3 < B3_pos[pB2+1] && pc1 < c1_pos[1]) {
            int kB = B3_crd[pB3];
            int kc = c1_crd[pc1];
            int k = min(kB, kc);
            if (kB == k && kc == k) {
                t += B[pB3] * c[pc1];
            }
            pB3 += (int)(kB == k);
            pc1 += (int)(kc == k);
        }
        A[pA2] = t;
    }
}
```

Optimized code is often complex, especially when it iterates over irregular data structures

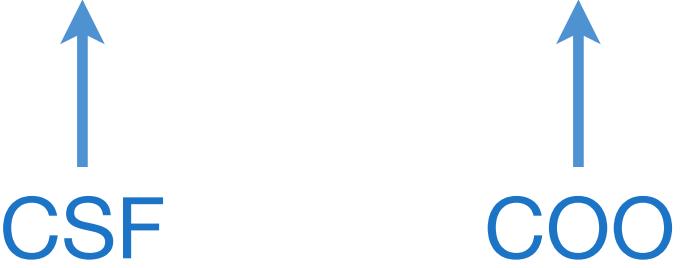
$$A_{ij} = \sum_k B_{ijk} c_k$$


CSF hash map

```
for (int pA = 0; pA < m*n; pA++) {
    A[pA] = 0.0;
}
for (int pB1 = B1_pos[0]; pB1 < B1_pos[1]; pB1++) {
    int i = B1_crd[pB1];
    for (int pB2 = B2_pos[pB1]; pB2 < B2_pos[pB1+1]; pB2++) {
        int j = B2_crd[pB2];
        int pA2 = i*n + j;
        double t = 0.0;
        for (int pB3 = B3_pos[pB2]; pB3 < B3_pos[pB2+1]; pB3++) {
            int k = B3_crd[pB3];
            int pc1 = k % c_size;
            if (c_crd[pc1] != k && c_crd[pc1] != -1) {
                int end = pc;
                do {
                    pc = (pc+1) % c_size;
                } while (c_crd[pc1] != k &&
                         c_crd[pc1] != -1 && pc1 != end);
            }
            if (c_crd[pc1] == k) {
                t += B[pB3] * c[pc1];
            }
        }
        A[pA2] = t;
    }
}
```

Optimized code is often complex, especially when it iterates over irregular data structures

$$A_{ijk} = B_{ijk} + C_{ijk}$$



```

int iB = 0;
int C0_pos = C0_pos[0];
while (C0_pos < C0_pos[1]) {
    int iC = C0_crd[C0_pos];
    int C0_end = C0_pos + 1;
    if (iC == iB)
        while ((C0_end < C0_pos[1]) && (C0_crd[C0_end] == iB)) {
            C0_end++;
        }
    if (iC == iB) {
        int B1_pos = B1_pos[iB];
        int C1_pos = C0_pos;
        while ((B1_pos < B1_pos[iB + 1]) && (C1_pos < C0_end)) {
            int jB = B1_crd[B1_pos];
            int jC = C1_crd[C1_pos];
            int j = min(jB, jC);
            int A1_pos = (iB * A1_size) + j;
            int C1_end = C1_pos + 1;
            if (jC == j)
                while ((C1_end < C0_end) && (C1_crd[C1_end] == j)) {
                    C1_end++;
                }
            if ((jB == j) && (jC == j)) {
                int B2_pos = B2_pos[B1_pos];
                int C2_pos = C1_pos;
                while ((B2_pos < B2_pos[B1_pos + 1]) && (C2_pos < C1_end)) {
                    int kB = B2_crd[B2_pos];
                    int kC = C2_crd[C2_pos];
                    int k = min(kB, kC);
                    int A2_pos = (A1_pos * A2_size) + k;
                    if ((kB == k) && (kC == k)) {
                        A[A2_pos] = B[B2_pos] + C[C2_pos];
                    } else if (kB == k) {
                        A[A2_pos] = B[B2_pos];
                    } else {
                        A[A2_pos] = C[C2_pos];
                    }
                    if (kB == k) B2_pos++;
                    if (kC == k) C2_pos++;
                }
                while (B2_pos < B2_pos[B1_pos + 1]) {
                    int kB0 = B2_crd[B2_pos];
                    int A2_pos0 = (A1_pos * A2_size) + kB0;
                    A[A2_pos0] = B[B2_pos];
                    B2_pos++;
                }
                while (C2_pos < C1_end) {
                    int kC0 = C2_crd[C2_pos];
                    int A2_pos1 = (A1_pos * A2_size) + kC0;
                    A[A2_pos1] = C[C2_pos];
                    C2_pos++;
                }
            } else if (jB == j) {
                for (int B2_pos0 = B2_pos[B1_pos]; B2_pos0 < B2_pos[B1_pos + 1]; B2_pos0++) {
                    int kB1 = B2_crd[B2_pos0];
                    int A2_pos2 = (A1_pos * A2_size) + kB1;
                    A[A2_pos2] = B[B2_pos0];
                }
            } else {
                for (int C2_pos0 = C1_pos; C2_pos0 < C1_end; C2_pos0++) {
                    int kC1 = C2_crd[C2_pos0];
                    int A2_pos3 = (A1_pos * A2_size) + kC1;
                    A[A2_pos3] = C[C2_pos0];
                }
            }
            if (jB == j) B1_pos++;
            if (jC == j) C1_pos = C1_end;
        }
    }
}

```

```

while (B1_pos < B1_pos[iB + 1]) {
    int jB0 = B1_crd[B1_pos];
    int A1_pos0 = (iB * A1_size) + jB0;
    for (int B2_pos1 = B2_pos[B1_pos]; B2_pos1 < B2_pos[B1_pos + 1]; B2_pos1++) {
        int kB2 = B2_crd[B2_pos1];
        int A2_pos4 = (A1_pos0 * A2_size) + kB2;
        A[A2_pos4] = B[B2_pos1];
    }
    B1_pos++;
}
while (C1_pos < C0_end) {
    int jC0 = C1_crd[C1_pos];
    int A1_pos1 = (iB * A1_size) + jC0;
    int C1_end0 = C1_pos + 1;
    while ((C1_end0 < C0_end) && (C1_crd[C1_end0] == jC0)) {
        C1_end0++;
    }
    for (int C2_pos1 = C1_pos; C2_pos1 < C1_end0; C2_pos1++) {
        int kB2 = C2_crd[C2_pos1];
        int A2_pos5 = (A1_pos1 * A2_size) + kB2;
        A[A2_pos5] = C[C2_pos1];
    }
    C1_pos = C1_end0;
}
else {
    for (int B1_pos0 = B1_pos[iB]; B1_pos0 < B1_pos[iB + 1]; B1_pos0++) {
        int jB1 = B1_crd[B1_pos0];
        int A1_pos2 = (iB * A1_size) + jB1;
        for (int B2_pos2 = B2_pos[B1_pos0]; B2_pos2 < B2_pos[B1_pos0 + 1]; B2_pos2++) {
            int kB3 = B2_crd[B2_pos2];
            int A2_pos6 = (A1_pos2 * A2_size) + kB3;
            A[A2_pos6] = B[B2_pos2];
        }
    }
    if (iC == iB) C0_pos = C0_end;
    iB++;
}
while (iB < B0_size) {
    for (int B1_pos1 = B1_pos[iB]; B1_pos1 < B1_pos[iB + 1]; B1_pos1++) {
        int jB2 = B1_crd[B1_pos1];
        int A1_pos3 = (iB * A1_size) + jB2;
        for (int B2_pos3 = B2_pos[B1_pos1]; B2_pos3 < B2_pos[B1_pos1 + 1]; B2_pos3++) {
            int kB4 = B2_crd[B2_pos3];
            int A2_pos7 = (A1_pos3 * A2_size) + kB4;
            A[A2_pos7] = B[B2_pos3];
        }
    }
    iB++;
}

```

Can we get abstractions *without* friction by moving the abstractions into the compiler?

Can we get abstractions *without* friction by moving the abstractions into the compiler?

Domain-Specific Language Constructs

Can we get abstractions *without friction* by moving the abstractions into the compiler?

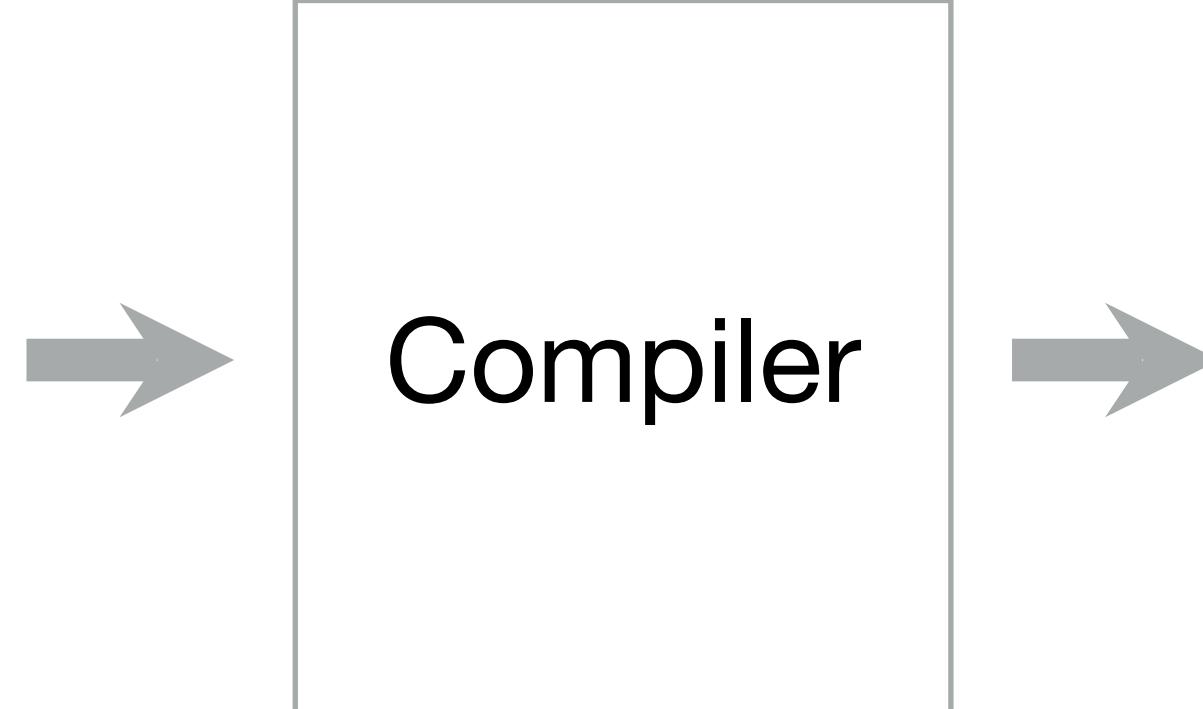
Domain-Specific Language Constructs

Generated Code

```
int iB = 0;
int C0_pos = C0_pos_arr[0];
while (C0_pos < C0_pos_arr[1]) {
    int iC = C0_idx_arr[C0_pos];
    int C0_end = C0_pos + 1;
    if (iC == iB)
        while ((C0_end < C0_pos_arr[1]) && (C0_idx_arr[C0_end] == iB)) {
            C0_end++;
        }
    if (iC == iB) {
        int B1_pos = B1_pos_arr[iB];
        int C1_pos = C0_pos;
        while ((B1_pos < B1_pos_arr[iB + 1]) && (C1_pos < C0_end)) {
            int jB = B1_idx_arr[B1_pos];
            int jC = C1_idx_arr[C1_pos];
            int j = min(jB, jC);
            int A1_pos = (iB * A1_size) + j;
            int C1_end = C1_pos + 1;
            if (jC == j)
                while ((C1_end < C0_end) && (C1_idx_arr[C1_end] == j)) {
                    C1_end++;
                }
            if ((jB == j) && (jC == j)) {
                int B2_pos = B2_pos_arr[B1_pos];
                int C2_pos = C1_pos;
                while ((B2_pos < B2_pos_arr[B1_pos + 1]) && (C2_pos < C1_end)) {
                    int kB = B2_idx_arr[B2_pos];
                    int kC = C2_idx_arr[C2_pos];
                    int k = min(kB, kC);
                    int A2_pos = (A1_pos * A2_size) + k;
                    if ((kB == k) && (kC == k)) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos] + C_val_arr[C2_pos];
                    } else if (kB == k) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos];
                    } else {
                        A_val_arr[A2_pos] = C_val_arr[C2_pos];
                    }
                    if (kB == k) B2_pos++;
                    if (kC == k) C2_pos++;
                }
                while (B2_pos < B2_pos_arr[B1_pos + 1]) {
                    int kB0 = B2_idx_arr[B2_pos];
                    int A2_pos0 = (A1_pos * A2_size) + kB0;
                    A_val_arr[A2_pos0] = B_val_arr[B2_pos];
                    B2_pos++;
                }
                while (C2_pos < C1_end) {
                    int kC0 = C2_idx_arr[C2_pos];
                    int A2_pos1 = (A1_pos * A2_size) + kC0;
                    A_val_arr[A2_pos1] = C_val_arr[C2_pos];
                    C2_pos++;
                }
            } else if (jB == j) {
                for (int B2_pos0 = B2_pos_arr[B1_pos];
                     B2_pos0 < B2_pos_arr[B1_pos + 1]; B2_pos0++) {
                    int kB1 = B2_idx_arr[B2_pos0];
                    int A2_pos2 = (A1_pos * A2_size) + kB1;
                    A_val_arr[A2_pos2] = B_val_arr[B2_pos0];
                }
            } else {
                for (int C2_pos0 = C1_pos;
                     C2_pos0 < C1_end; C2_pos0++) {
                    int kC1 = C2_idx_arr[C2_pos0];
                    int A2_pos3 = (A1_pos * A2_size) + kC1;
                    A_val_arr[A2_pos3] = C_val_arr[C2_pos0];
                }
            }
            if (jB == j) B1_pos++;
            if (jC == j) C1_pos = C1_end;
        }
    }
    while (B1_pos < B1_pos_arr[iB + 1]) {
        int jB0 = B1_idx_arr[B1_pos];
        int A1_pos0 = (iB * A1_size) + jB0;
        for (int B2_pos1 = B2_pos_arr[B1_pos];
             B2_pos1 < B2_pos_arr[B1_pos + 1]; B2_pos1++) {
            int kB2 = B2_idx_arr[B2_pos1];
            int A2_pos4 = (A1_pos0 * A2_size) + kB2;
            A_val_arr[A2_pos4] = B_val_arr[B2_pos1];
        }
        B1_pos++;
    }
    while (C1_pos < C0_end) {
        int jC0 = C1_idx_arr[C1_pos];
        int A1_pos1 = (iB * A1_size) + jC0;
        int C1_end0 = C1_pos + 1;
        while ((C1_end0 < C0_end) && (C1_idx_arr[C1_end0] == jC0)) {
            C1_end0++;
        }
        for (int C2_pos1 = C1_pos;
             C2_pos1 < C1_end0; C2_pos1++) {
            int kB2 = C2_idx_arr[C2_pos1];
            int A2_pos5 = (A1_pos1 * A2_size) + kB2;
            A_val_arr[A2_pos5] = C_val_arr[C2_pos1];
        }
        C1_pos = C1_end0;
    }
} else {
    for (int B1_pos0 = B1_pos_arr[iB];
         B1_pos0 < B1_pos_arr[iB + 1]; B1_pos0++) {
        int jB1 = B1_idx_arr[B1_pos0];
        int A1_pos2 = (iB * A1_size) + jB1;
        for (int B2_pos2 = B2_pos_arr[B1_pos0];
             B2_pos2 < B2_pos_arr[B1_pos0 + 1]; B2_pos2++) {
            int kB3 = B2_idx_arr[B2_pos2];
            int A2_pos6 = (A1_pos2 * A2_size) + kB3;
            A_val_arr[A2_pos6] = B_val_arr[B2_pos2];
        }
    }
}
if (iC == iB) C0_pos = C0_end;
iB++;
}
while (iB < B0_size) {
    for (int B1_pos1 = B1_pos_arr[iB];
         B1_pos1 < B1_pos_arr[iB + 1]; B1_pos1++) {
        int jB2 = B1_idx_arr[B1_pos1];
        int A1_pos3 = (iB * A1_size) + jB2;
        for (int B2_pos3 = B2_pos_arr[B1_pos1];
             B2_pos3 < B2_pos_arr[B1_pos1 + 1]; B2_pos3++) {
            int kB4 = B2_idx_arr[B2_pos3];
            int A2_pos7 = (A1_pos3 * A2_size) + kB4;
            A_val_arr[A2_pos7] = B_val_arr[B2_pos3];
        }
    }
    iB++;
}
```

Can we get abstractions *without friction* by moving the abstractions into the compiler?

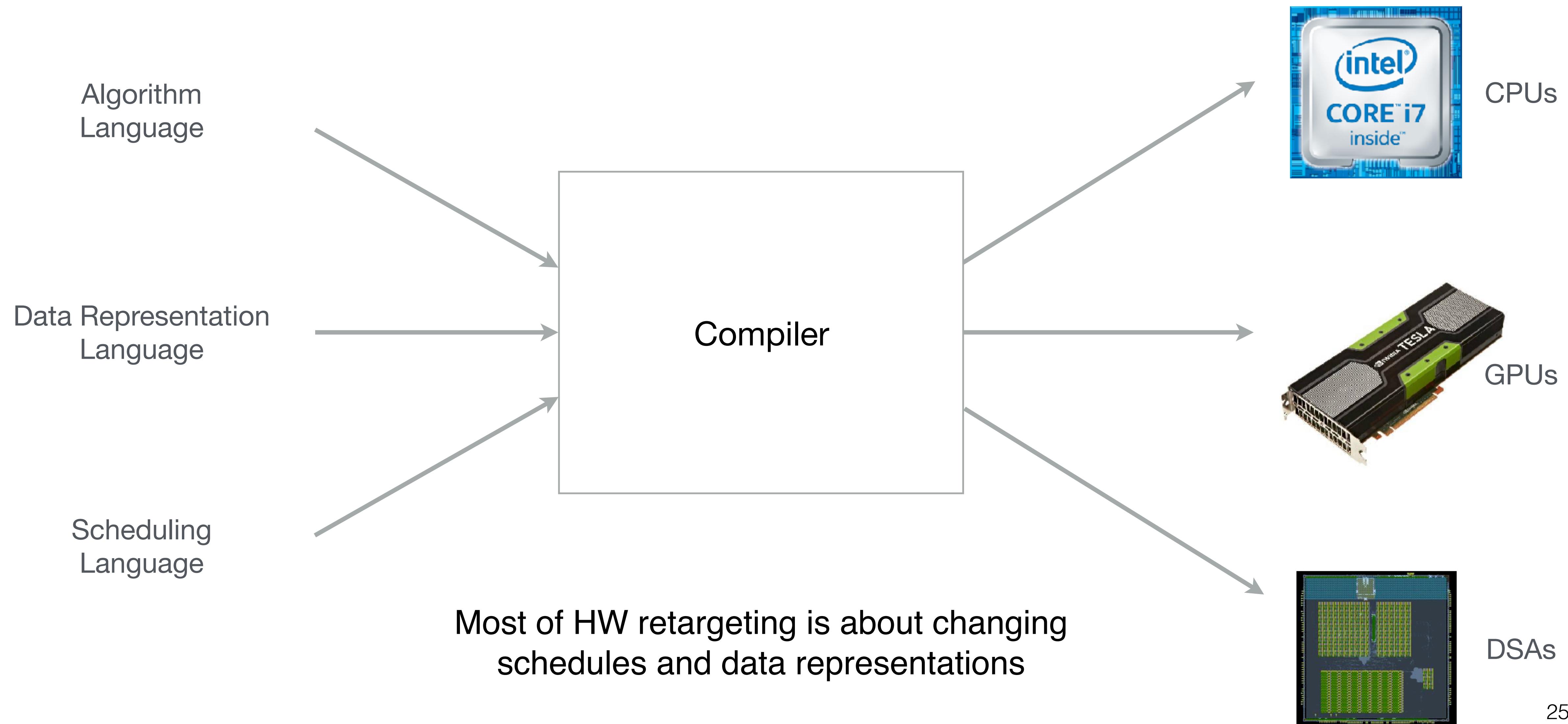
Domain-Specific
Language
Constructs



Generated Code

```
int iB = 0;
int C0_pos = C0_pos_arr[0];
while (C0_pos < C0_pos_arr[1]) {
    int iC = C0_idx_arr[C0_pos];
    int C0_end = C0_pos + 1;
    if (iC == iB)
        while ((C0_end < C0_pos_arr[1]) && (C0_idx_arr[C0_end] == iB)) {
            C0_end++;
        }
    if (iC == iB) {
        int B1_pos = B1_pos_arr[iB];
        int C1_pos = C0_pos;
        while ((B1_pos < B1_pos_arr[iB + 1]) && (C1_pos < C0_end)) {
            int jB = B1_idx_arr[B1_pos];
            int jC = C1_idx_arr[C1_pos];
            int j = min(jB, jC);
            int A1_pos = (iB * A1_size) + j;
            int C1_end = C1_pos + 1;
            if (jC == j)
                while ((C1_end < C0_end) && (C1_idx_arr[C1_end] == j)) {
                    C1_end++;
                }
            if ((jB == j) && (jC == j)) {
                int B2_pos = B2_pos_arr[B1_pos];
                int C2_pos = C1_pos;
                while ((B2_pos < B2_pos_arr[B1_pos + 1]) && (C2_pos < C1_end)) {
                    int kB = B2_idx_arr[B2_pos];
                    int kC = C2_idx_arr[C2_pos];
                    int k = min(kB, kC);
                    int A2_pos = (A1_pos * A2_size) + k;
                    if ((kB == k) && (kC == k)) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos] + C_val_arr[C2_pos];
                    } else if (kB == k) {
                        A_val_arr[A2_pos] = B_val_arr[B2_pos];
                    } else {
                        A_val_arr[A2_pos] = C_val_arr[C2_pos];
                    }
                    if (kB == k) B2_pos++;
                    if (kC == k) C2_pos++;
                }
                while (B2_pos < B2_pos_arr[B1_pos + 1]) {
                    int kB0 = B2_idx_arr[B2_pos];
                    int A2_pos0 = (A1_pos * A2_size) + kB0;
                    A_val_arr[A2_pos0] = B_val_arr[B2_pos];
                    B2_pos++;
                }
                while (C2_pos < C1_end) {
                    int kC0 = C2_idx_arr[C2_pos];
                    int A2_pos1 = (A1_pos * A2_size) + kC0;
                    A_val_arr[A2_pos1] = C_val_arr[C2_pos];
                    C2_pos++;
                }
            } else if (jB == j) {
                for (int B2_pos0 = B2_pos_arr[B1_pos];
                     B2_pos0 < B2_pos_arr[B1_pos + 1]; B2_pos0++) {
                    int kB1 = B2_idx_arr[B2_pos0];
                    int A2_pos2 = (A1_pos * A2_size) + kB1;
                    A_val_arr[A2_pos2] = B_val_arr[B2_pos0];
                }
            } else {
                for (int C2_pos0 = C1_pos;
                     C2_pos0 < C1_end; C2_pos0++) {
                    int kC1 = C2_idx_arr[C2_pos0];
                    int A2_pos3 = (A1_pos * A2_size) + kC1;
                    A_val_arr[A2_pos3] = C_val_arr[C2_pos0];
                }
            }
            if (jB == j) B1_pos++;
            if (jC == j) C1_pos = C1_end;
        }
    }
}
```

Separation of Algorithm, Data Representation, and Schedule



How do you develop new language and compiler abstractions

How do you develop new language and compiler abstractions

“Hitching our research to someone else’s driving problem,
and solving those problems on the owners’ terms
leads us to richer computer science research.”

— Fred Brooks

How do you develop new language and compiler abstractions

“Hitching our research to someone else’s driving problem, and solving those problems on the owners’ terms leads us to richer computer science research.”

— Fred Brooks

“Like other great software, great little languages are grown, not built. Start with a solid simple design, expressed in a notation like Backus-Naur form. Before implementing the language, test your design by describing a wide variety of objects in the proposed language. After the language is up and running, iterate designs to add new features as dictated by real use.”

— Jon Bentley (Little Languages)

A process for developing DSLs

A process for developing DSLs

- Study applications to find patterns in their computations
 - Best to work closely with application people
 - Empirical at first: you must see many examples

A process for developing DSLs

- Study applications to find patterns in their computations
 - Best to work closely with application people
 - Empirical at first: you must see many examples
- Then you generalize
 - Deductively from examples to a natural coherent class
 - Inductively from new examples and observed patterns
 - Generalization must work for observed cases
 - Generalization must work for something new

A process for developing DSLs

- Study applications to find patterns in their computations
 - Best to work closely with application people
 - Empirical at first: you must see many examples
- Then you generalize
 - Deductively from examples to a natural coherent class
 - Inductively from new examples and observed patterns
 - Generalization must work for observed cases
 - Generalization must work for something new
- Look for ways to build abstractions into the compiler
 - Lets you separately describe different concerns
 - E.g., describe data structures independently of program

For Discussion

c := 0

for i := 1 step 1 until n do

c := c + a[i]×b[i]

```
c := 0  
for i := 1 step 1 until n do  
  c := c + a[i]×b[i]
```

Def Innerproduct

= (Insert +)°(ApplyToAll ×)°Transpose

```

c := 0
for i := 1 step 1 until n do
  c := c + a[i]×b[i]

```

Def Innerproduct

$$= (\text{Insert } +) \circ (\text{ApplyToAll } \times) \circ \text{Transpose}$$

IP:<<1,2,3>, <6,5,4>> =	
Definition of IP	$\Rightarrow (/+) \circ (\alpha \times) \circ \text{Trans: } <<1,2,3>, <6,5,4>>$
Effect of composition, \circ	$\Rightarrow (/+):((\alpha \times):(\text{Trans: } <<1,2,3>, <6,5,4>>))$
Applying Transpose	$\Rightarrow (/+):((\alpha \times): <<1,6>, <2,5>, <3,4>>)$
Effect of ApplyToAll, α	$\Rightarrow (/+): <\times: <1,6>, \times: <2,5>, \times: <3,4>>$
Applying \times	$\Rightarrow (/+): <6,10,12>$
Effect of Insert, /	$\Rightarrow +: <6, +: <10,12>>$
Applying +	$\Rightarrow +: <6,22>$
Applying + again	$\Rightarrow 28$