



Scientific and Technical Computing

Hardware and Code Optimization

Lars Koesterke

UT Austin, 9/1/22 & 9/8/22 & 9/15/22

Discussion

What are the primary components of a computer?



Discussion

What are the **primary components** of a computer?

Can you detect some limitations?

Discussion

What are the primary components of a computer?

- CPU
- Memory
- (Storage)
- (Motherboard, lots of 'wires')
- (Keyboard, screen)

Can you detect some limitations?

- Number of pins connecting CPU and Memory to motherboard

What **sciences and technologies** are
involved in designing and building a
CPU/Memory/Computer?

What sciences and technologies are involved in designing and building a CPU?

- Electrical engineering, physics, chemistry, math ... (all the things you study)
- Cutting edge research
- A lot of institutional knowledge by people in companies and research institutions
 - Not everything is an exact science; many tricks are applied (a bit like cooking)
 - Design decisions are made on incomplete facts (humans weigh the pros and cons)
- Certainly, many computers using current/previous generations of CPUs

This is **not** what we will talk about in this section of the class segment
(Hardware and Code Optimization)

Scope of this class

What are we going to explore in the next weeks?

High level overview of the architecture

- High level of abstraction
- Simplified implementation details
- Features that allow for a very high peak performance

How to write code that exploits the hardware features?

Scope of this class

Matching software to hardware: Why?

Much higher performance

- Orders of magnitude!

There is, of course, the idea to match the hardware to the software/purpose

This is done in other areas
In HPC the idea is not feasible (with one notable exception)

Conceptual understanding of hardware features guides software design

Assumption: You are in this class (and other TACC classes) to learn how to

- learn about high-performance computing, HPC (language implications)
- use a supercomputer (or any computer!) in an **efficient** and **effective** way
- write **fast** code
 - This class: exploiting parallelism of the hardware
 - Note: some/many bad code design decision cannot be reversed later
- write **parallel** code with OpenMP and MPI (PCSE in the spring semester)

Getting scientific calculations done
Better than the competition

Terminology

Today's supercomputers are clusters

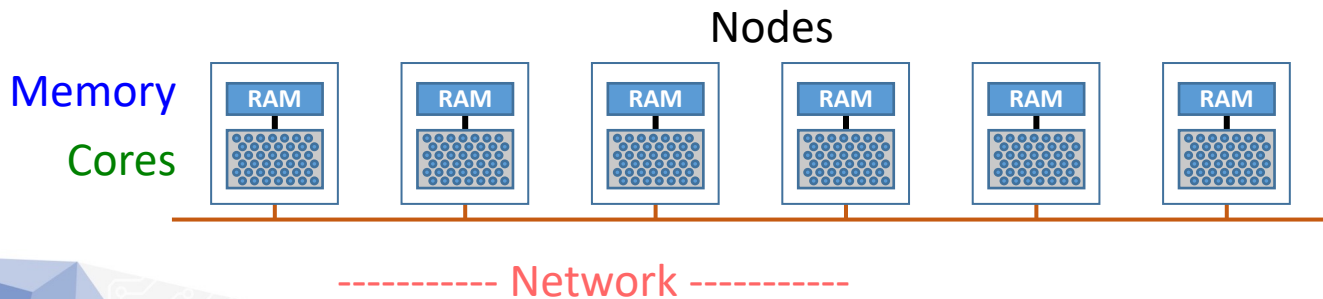
Components of a cluster

- Nodes
 - CPUs and memory
- Network
- (File system)

In this class we strongly focus on a single-core and single-node

A cluster is built from individual computers which are called nodes

The nodes are connected through an interconnect



Discussion

What is a clock tick?

What is a clock cycle?

Discussion

What is a clock tick?

What is a clock cycle?

What does this mean?

Clock frequency = 2.2 GHz

Discussion

What is a clock tick?

What is a clock cycle?

Think of an assembly line

- Smallest unit of time to 'do' something
- One or multiple instructions are executed
- An instruction may take several cycles

Examples of instructions are

- Multiply two numbers
- Load data into a register

Some answers from the web

Computers use an internal clock to synchronize all of their calculations. The clock ensures that the various circuits inside a computer work together at the same time.

Same as a cycle, the smallest unit of time recognized by a device. For personal computers, clock ticks generally refer to the main system clock, which runs at 66 MHz. This means that there are 66 million clock ticks (or cycles) per second. Since modern CPUs run much faster (up to 3 GHz), the CPU can execute several instructions in a single clock tick.

"The processor clock coordinates all CPU and memory operations by periodically generating a time reference signal called a *clock cycle* or *tick*. Clock frequency is specified in gigahertz (GHz), which specifies billions of ticks per second. Clock speed determines how fast instructions execute. Some instructions require one tick, others multiple ticks, and some processors execute multiple instructions during one tick."

Intermission

Let's talk about how to proceed

I'd like to organize this class having this in mind:

What and how to learn?

What will be on the slides?

How to participate?

Give me feedback!

How will participation affect your grade?

Teamwork

Let me know (at a later point) what you are interested in

Experiment: Prepare something at home

Your tasks

- Look up what the 'Horner scheme' is
 - Wikipedia entry (English Wikipedia site) is very good
 - To be specific, Horner's notation: $((((z+\dots)\times z+\dots)\times z \dots$
- Describe the 'Horner notation'
 - Three to four sentences
 - General context (Note: you don't have to explain what a polynomial is)
 - What is the 'trick'?
 - Why would you use it when writing code?
- 'Present' in class next week
 - Nothing dramatic, I'll explain
 - Write the 3 sentences down, if that helps you

Horner Scheme

1. The Horner scheme is an efficient way to evaluate polynomials
2. The number of multiplications is minimized
3. This is achieved by successively building x^{n+1} from x^n
4. $y = (... (a_n x + a_{n-1}) x + a_{n-2}) x + ...$
5. This scheme can be found in many HPC applications
6. Question for later: Is this really the fastest scheme for today's architecture

Let's make our computer do something

$$a = b + c$$

Let's make our computer do something

$$a = b + c$$

What needs to happen so that the computer can perform the calculation?

Where is the data coming from?

Where is the data going?

What part of the hardware performs the operation?

Let's make our computer do something

$$a = b + c$$

What needs to happen so that the CPU can calculate?

Where is the data coming from?

Memory

Where is the data going?

Memory

What part performs the operation?

Floating Point Unit (FPU) in the CPU (Central Processing Unit)

Let's make our computer do something

$$a = b + c$$

What needs to happen so that the CPU can calculate?

How long does it take?

(What units to use?)

Where is the data coming from?

Memory

Where is the data going?

Memory

What part performs the operation?

FPU in the CPU

Let's make our computer do something

$$a = b + c$$

What needs to happen so that the CPU can calculate?

How long does it take?

Where is the data coming from?

Memory

A very long time

Where is the data going?

Memory

A very long time

What part performs the operation?

FPU in the CPU

A very short time

Let's make our computer do something

$$a = b + c$$

What needs to happen so that the CPU can calculate?

How long does it take?

Where is the data coming from?

Memory

300 cycles

Where is the data going?

Memory

300 cycles

What part performs the operation?

FPU in the CPU

3 cycles

What a bummer! Actual work takes 0.5% of the total time
We may have built a very ineffective computer

Does this help us?

$$a(i) = b(i) + c(i)$$

Does this help us?

$$a(i) = b(i) + c(i)$$

Hint: where would you likely find such a statement?

Does this help us?

$$a(i) = b(i) + c(i)$$

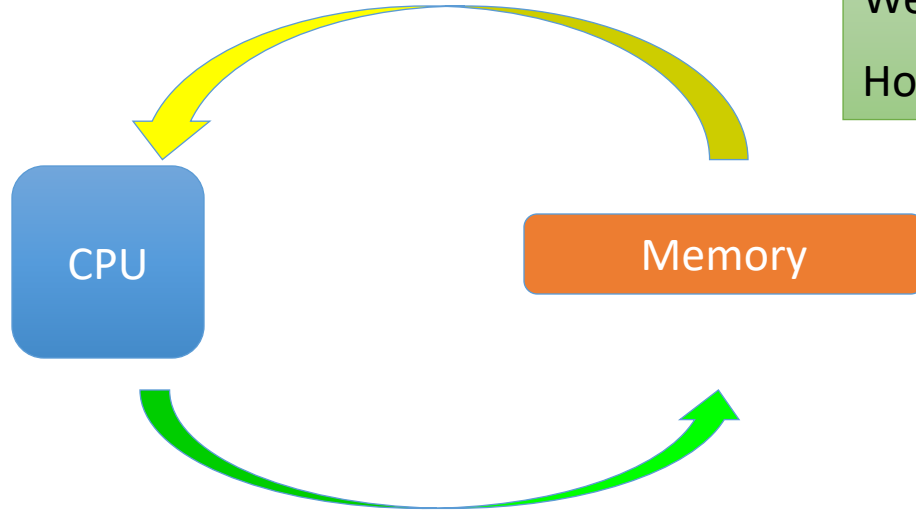
```
loop with index i  
  a(i) = b(i) + c(i)  
end loop
```


Data

$$a(i) = b(i) + c(i)$$

We have a lot of data to process

How can that help to **'getting more done'**?

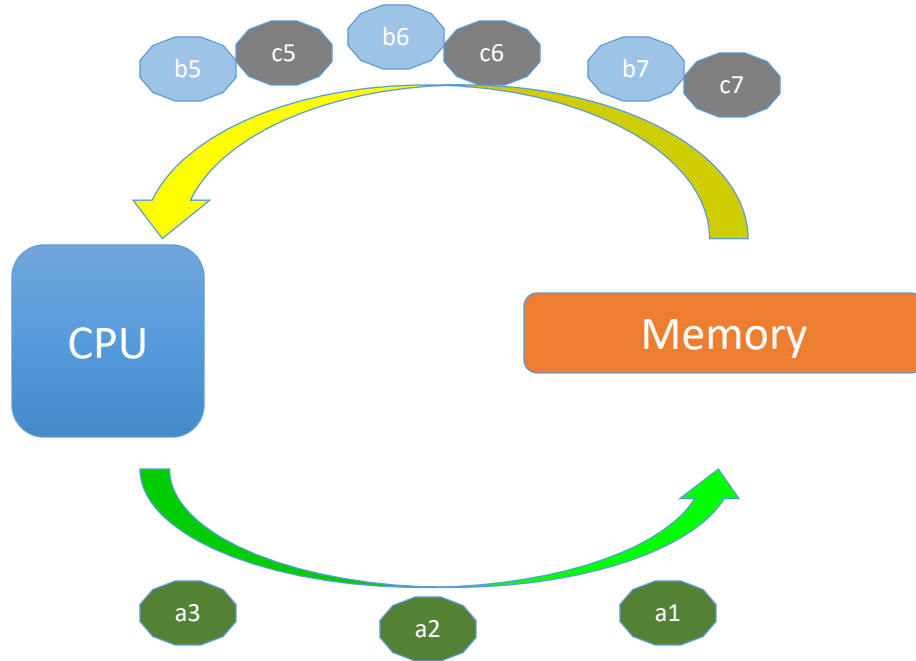


'getting more done'

We are mostly interested in floating point operations (flops) that move the calculation closer to the solution

Data Streams

$$a(i) = b(i) + c(i)$$



Some data is 'en route'

b and c: from memory to CPU

a: from CPU to memory

Some data is being processed

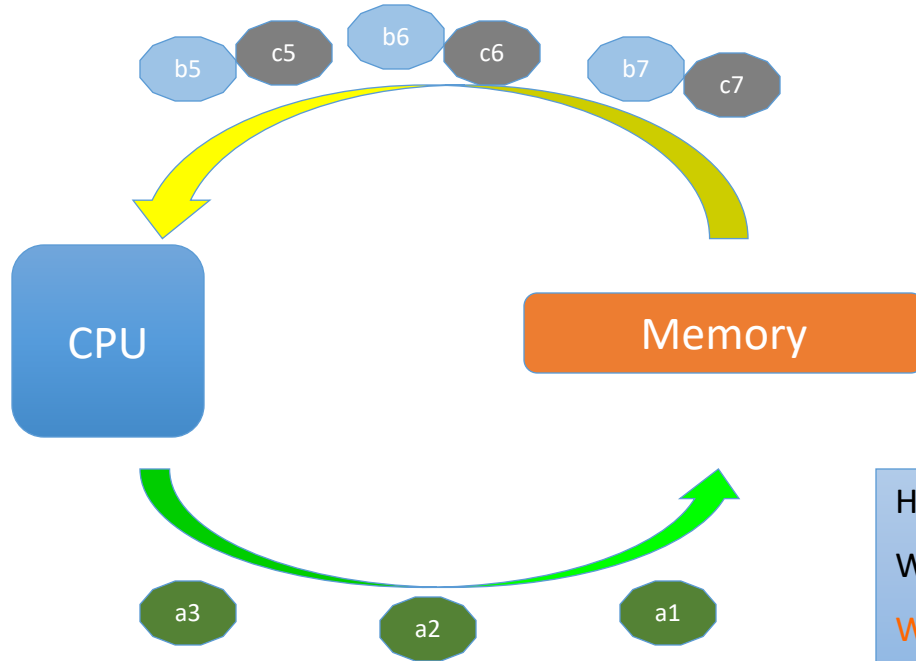
$$a(i) = b(i) + c(i)$$

How much data has to be 'en route'?

What are the main factors?

Data Streams

$$a(i) = b(i) + c(i)$$



Some data is 'en route'

b and c: from memory to CPU

a: from CPU to memory

Some data is being processed

$$a(i) = b(i) + c(i)$$

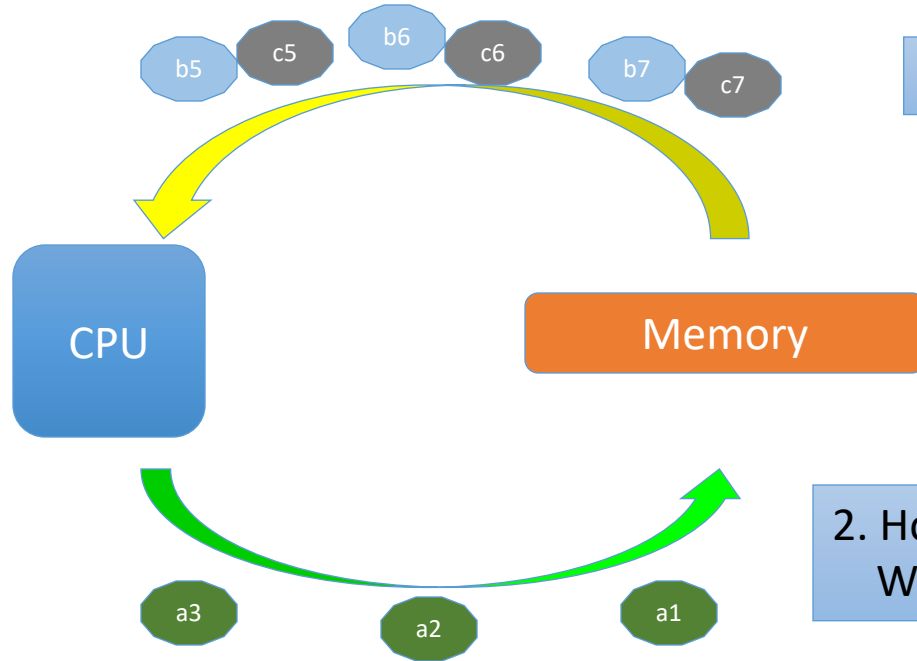
How much data has to be 'en route'?

What is one of the main factors?

What if I had drawn CPU and Memory further apart?

Data Streams

$$a(i) = b(i) + c(i)$$



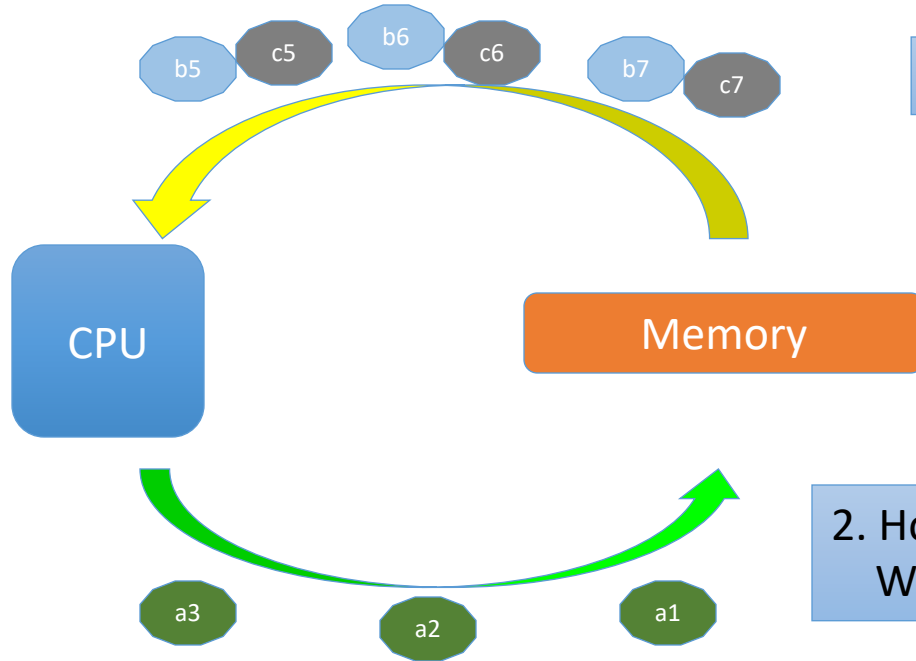
1. How much data has to be 'en route'?

Operation: 3 cycles
Data transfer: 300 cycles

2. How many results 'a' do we get per cycle?
What is the speedup?

Data Streams

$$a(i) = b(i) + c(i)$$



1. How much data has to be 'en route'?

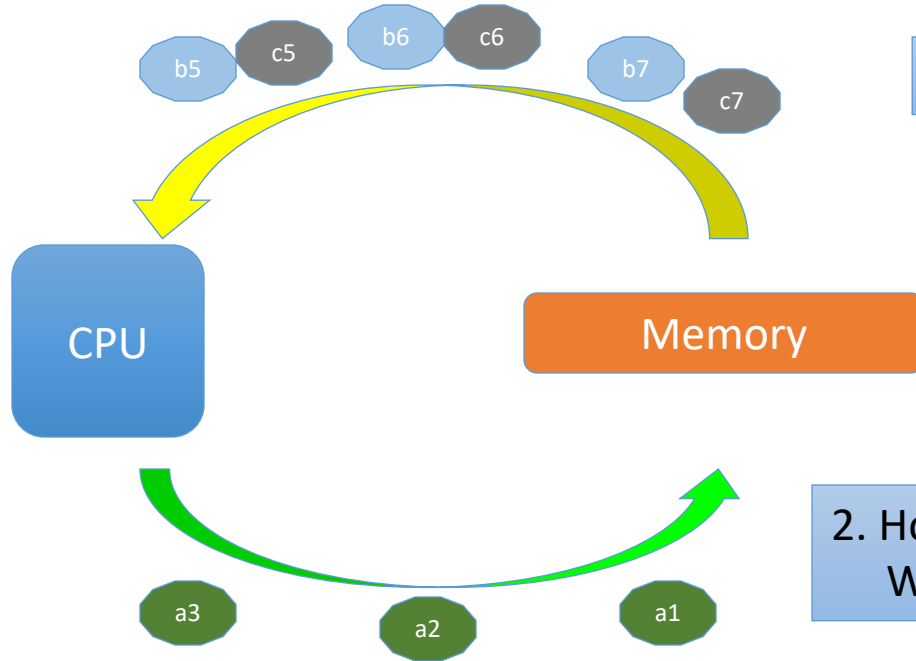
Operation: 3 cycles
Data transfer: 300 cycles

100 elements of a, b, and c

2. How many results (a) do we get per cycle?
What is the speedup?

Data Streams

$$a(i) = b(i) + c(i)$$



1. How much data has to be 'en route'?

Operation: 3 cycles
Data transfer: 300 cycles
Ratio: 1/100

100 elements of a, b, and c

2. How many results (a) do we get per cycle?
What is the speedup?

1 result every 3 cycles
Speedup is 100×

Hooray!

We have just discovered one of the most important hardware features in a computer: Data streams

- Data Streams
 - Long distance (in terms of cycles) between main memory and CPU
 - Short time to execute 'add' operation (few cycles)
 - Streaming data: Data 'en route' filling the stream between memory and CPU
- Questions for later
 - How do we or the CPU 'organize' the data stream?
 - How does this look in code?

There are 2 fundamental bottlenecks

1. The data supply to the CPU
2. The actual operations (flops)

Hooray!

We have just discovered one of the most important hardware features in a CPU: **Data streams**

- Data Streams
 - Long distance (in terms of cycles) between main memory and CPU
 - Short time to execute 'add' operation (few cycles)
 - Streaming data: Data 'en route' filling the pipeline between memory and CPU
- Questions for later
 - How do we or the CPU 'organize' the data stream?
 - How does this look in code?

There are 2 fundamental bottlenecks

1. The data supply to the CPU
2. The actual operations (flops)

3 major technologies are applied

(we can argue about the exact number)

1. Data streams
2. 'Improving CPU throughput'
3. 'Improving data movement and data use'

Discussion

Can we speed-up the actual numerical operation 'add'?

If so, how?

Assume

1. 'add' takes 3 cycles
2. The data streams supply data effortlessly, i.e. without delay and at infinite bandwidth

Discussion

How can we speed-up the actual operation 'add'?

Assume

1. 'add' takes 3 cycles
2. The data streams supply data effortlessly

Think of Henry Ford's moving assembly line

The assembly line predates H. Ford (see R. Olds in the automotive industry; but earlier assembly lines in other industries)

Discussion: Pipelining

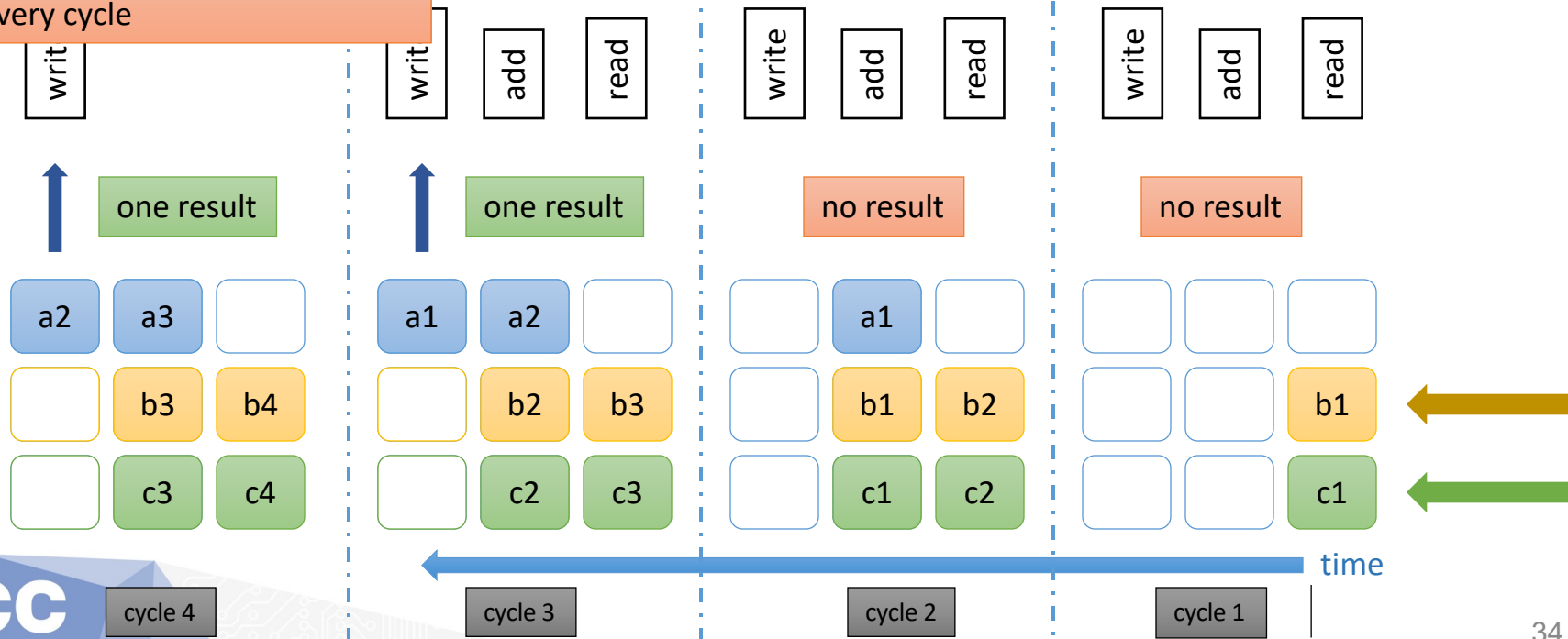
Simple toy model of the implementation of 'add'

'add' takes 3 cycles: For sake of argument, 'read', 'add', 'write'

After 2 cycles (ramp up phase) a result is produced every cycle

One result per cycle, once the pipeline is filled
3× performance increase

Operation: 1 cycle
Data transfer: 300 cycles
Ratio: 1/300
Ratio is now worse (no good deed goes unpunished)



Discussion: Pipelining

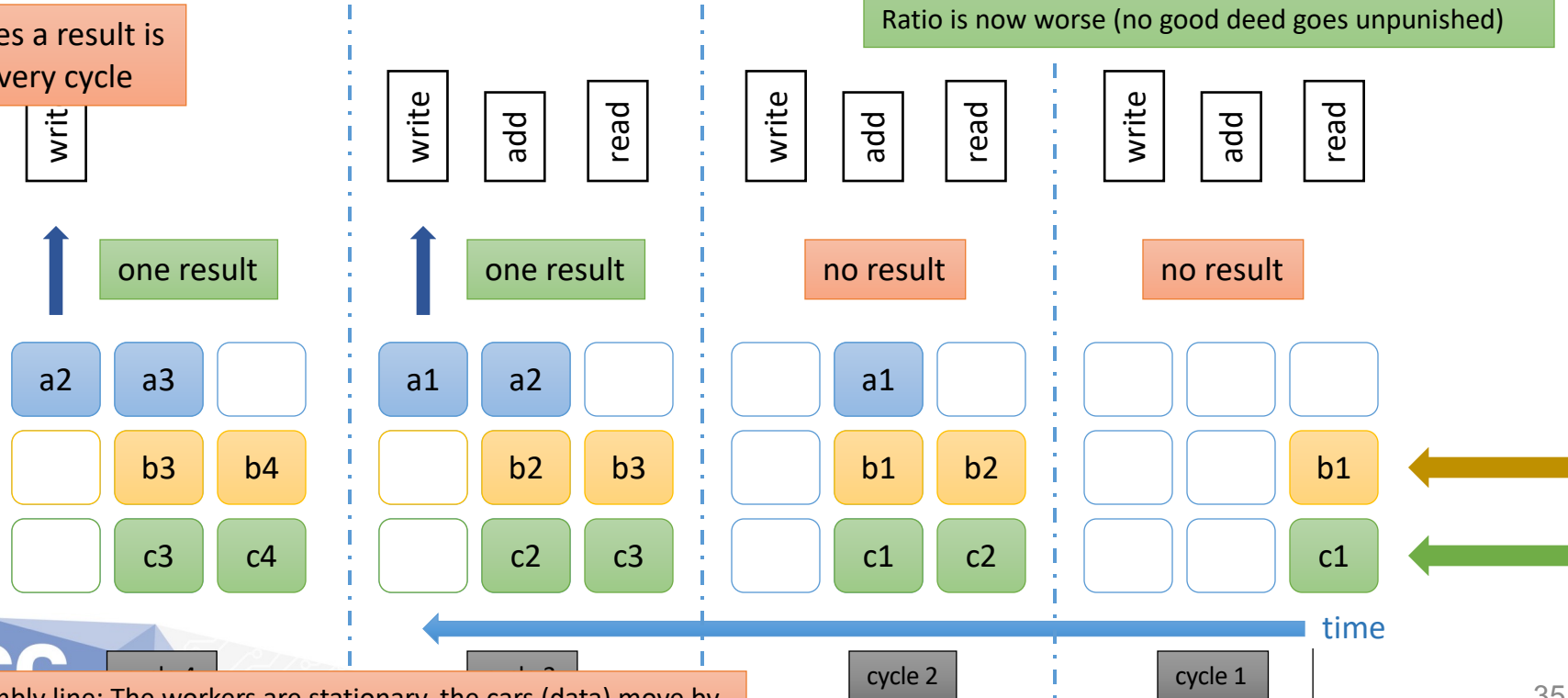
Simple toy model of the implementation of 'add'

'add' takes 3 cycles: 'read', 'add', write'

After 2 cycles a result is produced every cycle

One result per cycle, once the pipeline is filled
3× performance increase

Operation: 1 cycle
Data transfer: 300 cycles
Ratio: 1/300
Ratio is now worse (no good deed goes unpunished)



Moving assembly line: The workers are stationary, the cars (data) move by

We have just discovered another important hardware features in a CPU: **Pipelining**

- Pipelining
 - Single operation (add) takes more than a cycle
 - Pipelining (moving assembly line) allows to calculate **one result per cycle** once the pipeline is filled
- Questions for later
 - How do we or the CPU 'organize' pipelining?
 - How does this look in code?

2 fundamental bottlenecks

1. The data supply to the CPU
2. The actual operations

3 major technologies in CPU/memory

(we can argue about the exact number)

1. Data streams ($R=1/100$)
2. Pipelining ($R=1/300$ data supply even more important)
3. 'Improving data movement and data use'

Recap 'Data Streams'

Bandwidth & Latency

Bandwidth is the amount of data transferred per unit of time
Most convenient units for now based on **words** and **cycles**

So far we have talked about data streams and how to improve bandwidth

Recap 'Data streams'

- 300 cycles to move data between main memory to the CPU
- This is the latency (300 cycles to get the first data)

Assume no streams for the example: $a = b + c$

- Average bandwidth = 1 'word' per 200 cycles (0.005 wpc)
 - 'a', 'b', and 'c' are 4-byte or 8-byte words; single or double precision, respectively
- Let's pause here: Why exactly is the average bandwidth 1 word per 200 cycles?
- Let's discuss ...

I'm making this unit up
wpc: word per cycle

Recap 'Data Streams'

Bandwidth

So far we have talked about data streams and how to improve bandwidth

Recap 'Data streams'

- 300 cycles to move data between main memory to the CPU

Assume no streams for the example: $a = b + c$

- Average bandwidth = 1 'word' per 200 cycles (0.005 wpc)

- 'a', 'b', and 'c' are 4-byte or 8-byte words

wpc: our unit

- Moving 'b' and 'c' from memory to CPU: 2 words per 300 cycles
- Moving 'a' from CPU to memory: 1 word per 300 cycles
- Average: 3 words per 600 cycles = 0.005 wpc
round-trip: 600 cycles

Recap 'Data Streams'

Bandwidth

Same example, but now with streams and with pipelining: $a = b + c$

- Bandwidth = 3 wpc
- Again, let's pause here: Why exactly is the bandwidth 3 words per cycle?

Recap 'Data Streams'

Bandwidth

Bandwidth is the amount of data transferred per unit of time
Most convenient units for now: **words** and **cycles**

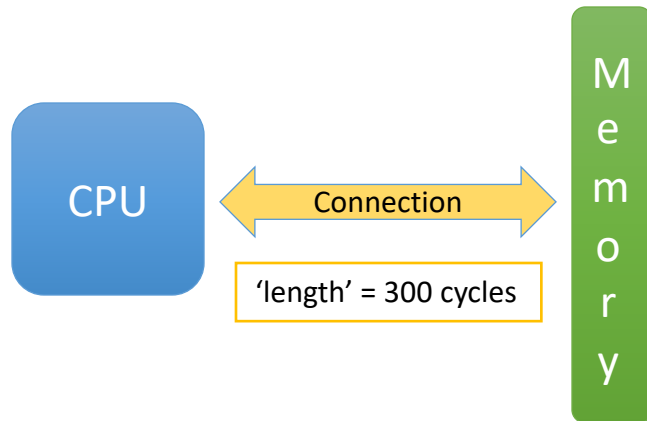
Same example, but now with streams and with pipelining: $a = b + c$

- Bandwidth = 3 wpc
- Every cycle one element of 'b' and 'c' are received, respectively
- Every cycle one element of 'a' is sent back
- In total, 3 words are received and sent every cycle
- Pipelining: one 'add' operation per cycle

Our first Computer: $a(i) = b(i) + c(i)$

Let's 'build' a computer and look at the requirements to achieve full performance

- CPU
 - Compute requirements: ...
 - Data movement: ...
- Memory
 - Data movement: ...
- Connection
 - Data movement: ...



Performance goal: 1 operation (add) per cycle
(we are going to ignore the ramp-up and ramp-down phase for now)

Our first Computer: $a(i) = b(i) + c(i)$

Let's 'build' a computer and look at the requirements to achieve full performance

- CPU

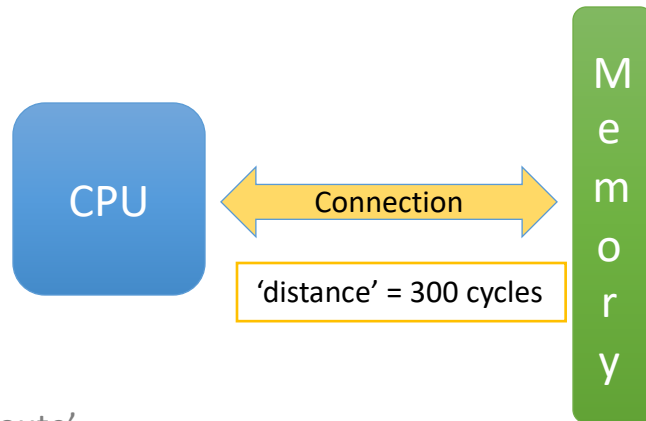
- Compute requirements: pipelined, 1 opc
- Data movement: 3 wpc

- Memory

- Data movement: 3 wpc

- Connection

- Data movement: 3 wpc and a total of 900 words 'en route'



Performance goal: 1 operation per cycle

(we are going to ignore the ramp-up and ramp-down phase for now)

But I'm starting to wonder how long these phases may be ...

Our first computer and our first source code

```
'Code kernel'  
n = 10000  
do i=1, n  
    a(i) = b(i) + c(i)  
enddo
```

Pretty simple code

- Data streams between memory and CPU
- CPU executes the 'add' operation

'Full code'

```
program add  
real,dimension(:),allocatable :: a,b,c  
n = 10000  
allocate(a(n),b(n),c(n))  
do i=1, n  
    a(i) = b(i) + c(i)  
enddo  
end program
```

Hints

- 'real' = 'float'
- 'allocate' = 'malloc'
- arrays start from 1, unless noted otherwise
- 'do' = 'for'

Technically, the average performance depends on the value of 'n', but we'll make it easy

1. Large 'n': bandwidth = 3 wpc → performance 1 opc (operation per cycle)
2. Small 'n': bandwidth and performance limited (down to fractions of a percent)

Recap: class period 1

Things we have discussed:

- Primary components of a computer
- Clock tick, clock frequency and its limitations
- Units: word, wpc, opc
- Data streams
- Memory latency & bandwidth
- Pipelining
- Compute latency

Second Source Code: Stencil Update

Let's beef-up our computer ... (just adding numbers is a bit dull)

We include 'mult' in our instruction set: multiply 2 numbers: $a = b \times c$

- 'add' 3 cycles, pipelined
- 'mult' 5 cycles, pipelined

3 and 5 cycles are realistic numbers for today's x86 hardware
We'll worry about the details later

Now we can implement a stencil update and can discover the next major hardware feature

Does anybody know:
What is a stencil update?
Where is a stencil update being used?

Second Source Code: Stencil Update

Example of a very simple case in 2d

'x' and 'y' are 2d arrays

For every position (i,j) an array element of 'y' is calculated as the average of its neighbors in the array 'x'

Each element:

$$y_{i,j} = 0.25 * (x_{i-1,j} + x_{i+1,j} + x_{i,j-1} + x_{i,j+1})$$
$$y = 0.25 * ('N' + 'S' + 'W' + 'E')$$

Second Source Code: Stencil Update

Most simple implementation in 2d

- Square arrays of identical size
- Only inner points are updated
- A little bit of space wasted, but no `if` statements for the boundary

```
'Code kernel'  
n = 10000  
m = n + 1  
allocate (x(0:m,0:m),y(0:m,0:m)) ! We allocate n+2 elements in both directions  
                                     ! Simplification: array boundaries start at 0 in C and Fortran  
do j=1, n  
  do i=1, n  
    y(i,j) = 0.25 * (x(i-1,j) + x(i+1,j) + x(i,j-1) + x(i,j+1))  
  enddo  
enddo
```

Only update the inner points: 1 to n

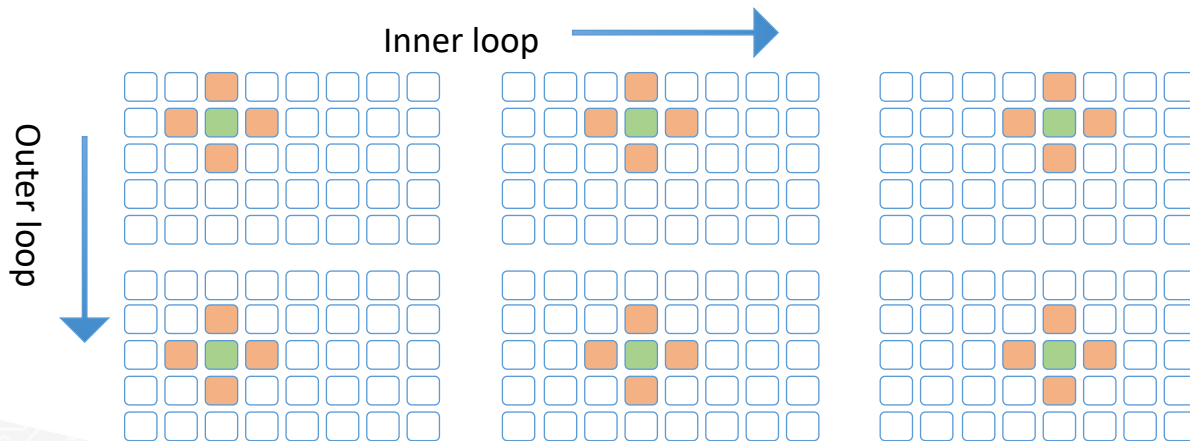
Second Source Code: Stencil Update

Let's discuss performance in terms of bandwidth

How much bandwidth (provided by the CPU-to-Memory connection) is needed for optimal/maximum performance?

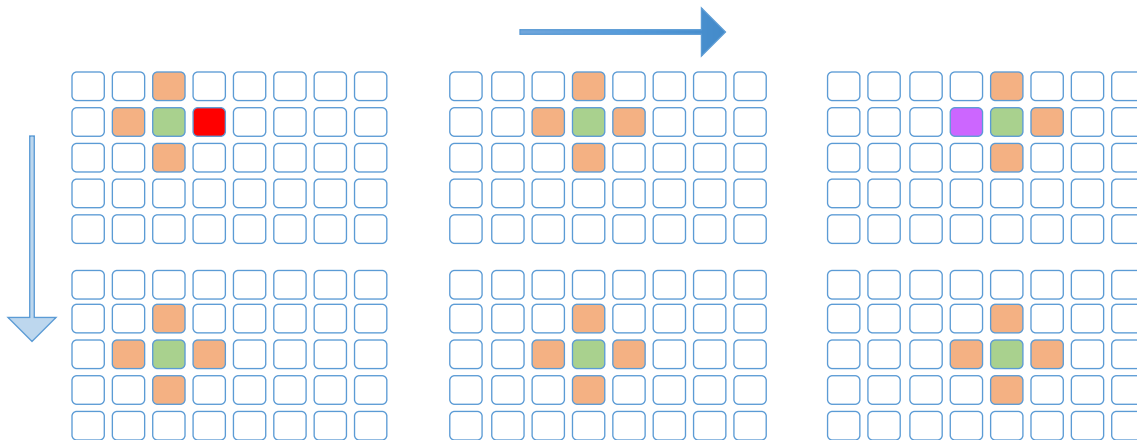
4 words input, 1 word output → 5 wpc

Can we lower the bandwidth requirement for the CPU-to-Memory connection?



Second Source Code: Stencil Update

Can we lower the bandwidth requirement?



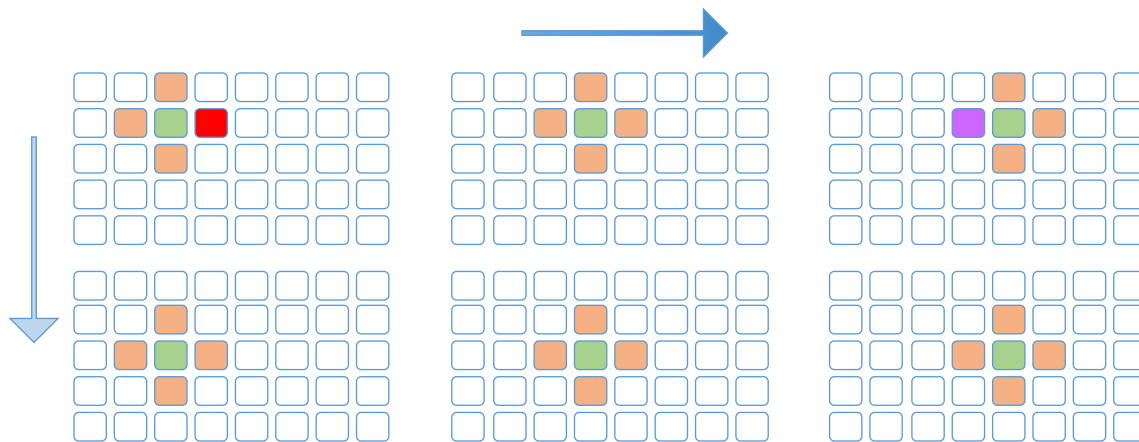
Second Source Code: Stencil Update

Can we lower the bandwidth requirement?

The element that we have loaded **here** is also loaded over **here**

If we go in x-direction first, this is just 2 loop iterations (inner loop) later

The data is being held in a register. No bandwidth is use to store the data between the 2 accesses

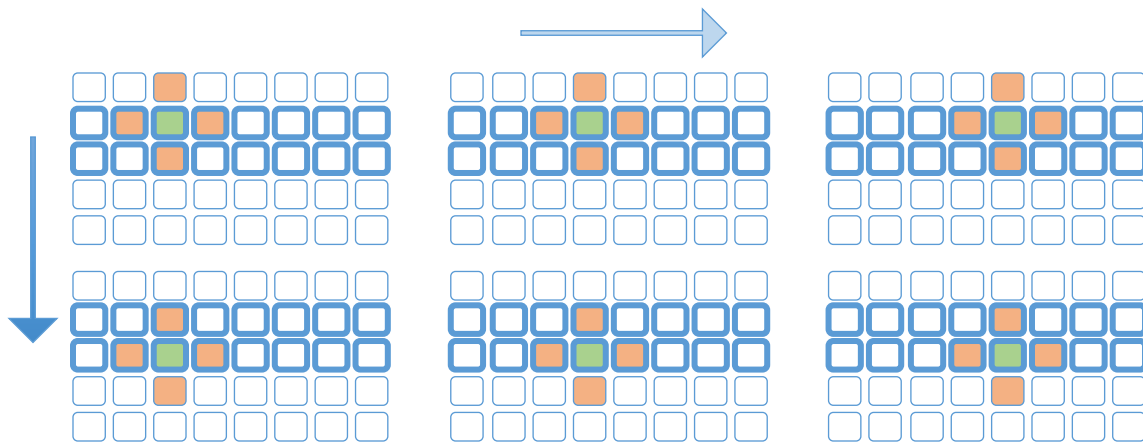


Second Source Code: Stencil Update

Can we lower the bandwidth requirement?

Similarly, two rows have been loaded in the previous y-loop

Can we use this to our advantage?



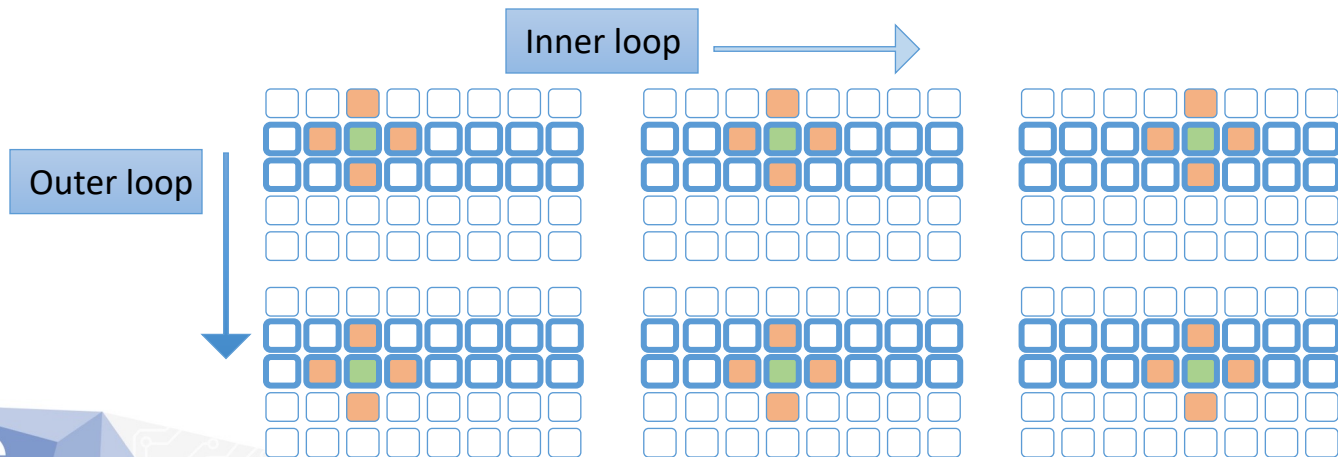
Second Source Code: Stencil Update

Can we lower the bandwidth requirement?

Two rows have been loaded in the previous y-loop

Can we use this to our advantage?

What if we could store two rows in a special buffer that provides higher bandwidth?



Second Source Code: Stencil Update

Can we lower the bandwidth requirement?

Similarly, two rows have been loaded in the previous y-loop

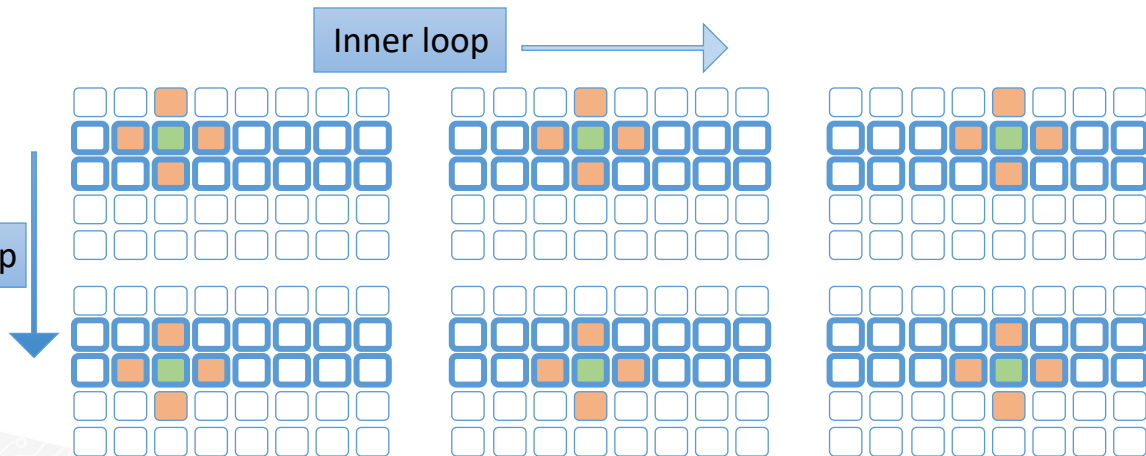
Can we use this to our advantage?

What if we could store two rows in a special buffer that provides higher bandwidth?

Such a buffer would be called a 'cache'

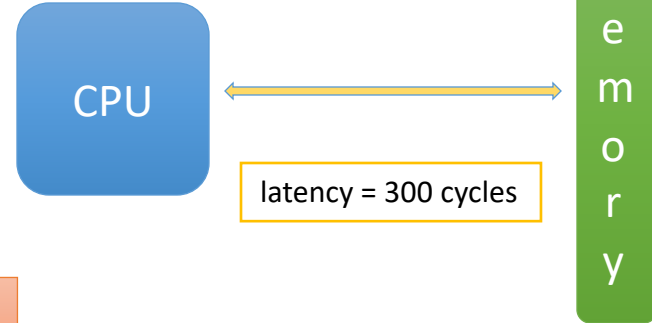
Outer loop

Inner loop



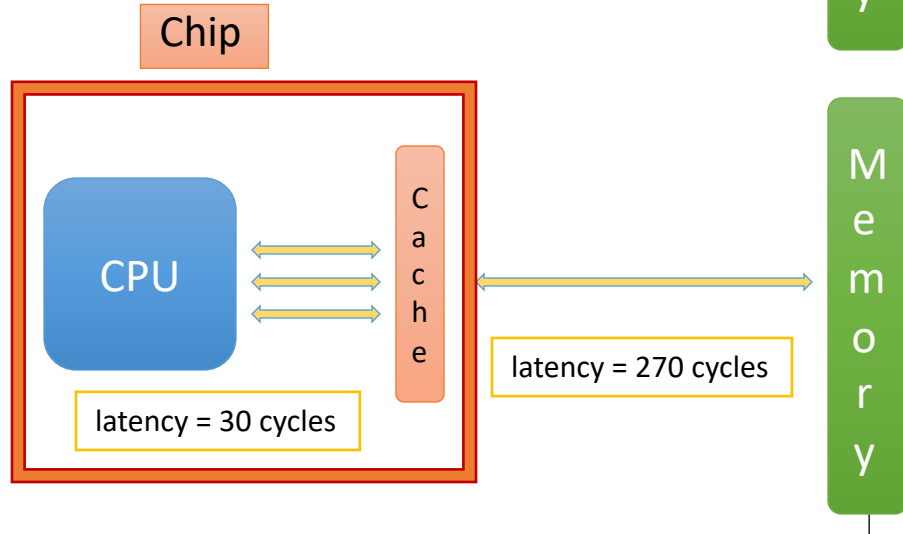
Computer with an added Cache

I hope you remember this diagram



Updated diagram with a cache

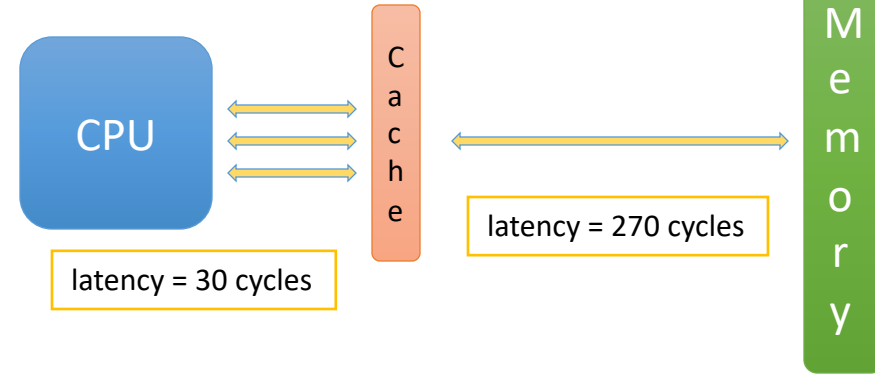
Cache between CPU and Memory
Cache is much(!) smaller than main memory
Bandwidth is higher
Cache closer to CPU (than to memory)
Caches usually/nowadays on chip



Computer with an added Cache

Some questions to ponder:

1. How big should the cache be for our specific example?
 - Answer in terms of 'n'
2. How does the cache provide 3× higher bandwidth?



```
'Code kernel'  
n = 10000  
m = n+1  
allocate (x(0:m,0:m), y(0:m,0:m))  
do j=1, n  
  do i=1, n  
    y(i,j) = 0.25 * (x(i-1,j) + ...  
  enddo  
enddo
```


Computer with an added Cache

Some questions to ponder:

How big should the cache be for our specific example?

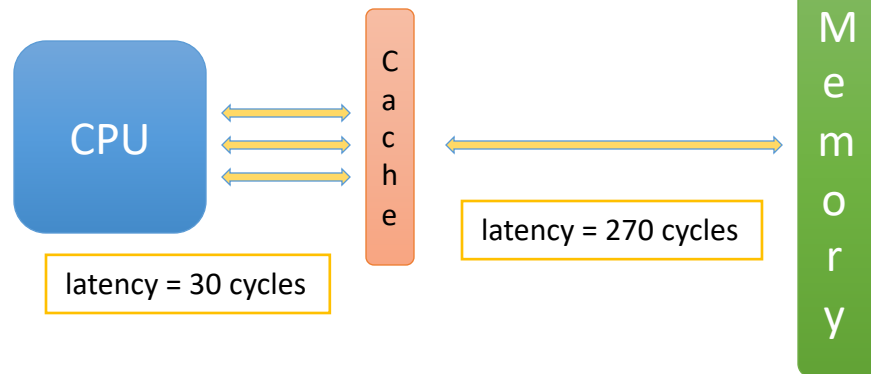
At least $2 \times n$ words

How does the cache provide $3 \times$ higher bandwidth?

Much closer

$3 \times$ as many wires

...



There is a deeper question here as well

How does the cache put ' $3 \times$ ' more data onto the connection?

Note that the main memory is fast enough to feed ' $1 \times$ ' data (the capabilities of the memory and the connection match)

```
'Code kernel'  
n = 10000  
m = n+1  
allocate (x(0:m,0:m), y(0:m,0:m))  
do j=1, n  
  do i=1, n  
    y(i,j) = 0.25 * (x(i-1,j) + ...  
  enddo  
enddo
```

Computer with an added Cache

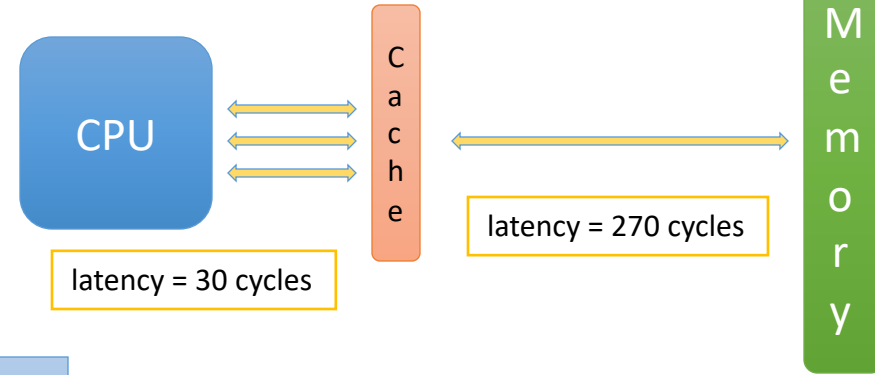
Some questions to ponder:

How big should the cache be for our specific example?

At least $2 \times n$ words

How does the cache provide $3 \times$ higher bandwidth?

More closely located, $3 \times$ as many wires, different hardware



There is a deeper question here as well

How does the cache put $3 \times$ more data onto the connection?

Assume that the main memory is fast enough to feed $1 \times$ data

Caches use a different storage technology

Main memory: DRAM --- Dynamic Random-Access Memory

Cache: SRAM --- Static Random-Access Memory

More later (maybe)

Random-access means that you can access data in any order

A different method would be a stack (of cards)
You can only add (write) to the top of the stack
You can only take (read) from the top of the stack

Computer with an added Cache

An additional fact about caches

- Caches are managed by the run-time
- User (user code) has no control over it

If a 'cache' were user-controlled, then it would (typically) have a different name

- Such user-controlled 'caches' are sometimes part of non-x86 architectures

This leaves us with a big question

What strategy can we devise to make the cache most useful *and* 'automatic'?

The problem:

The cache is much(!) smaller (at least 1000×) than the memory

What do we do when the cache has filled up?

Computer with an added Cache

An additional fact about caches

- Caches are managed by the run-time
- Users (user code) has no control over it

This leaves us with a big question

What strategy can we devise to make the cache most useful *and* 'automatic'?

The problem:

The cache is much(!) smaller (at least 1000×) then the memory

What do we do when the cache has filled up?

One solution: **FIFO**, **First-in first-out**

When the cache is full and new data is entered, then the oldest data is evicted

In contrast: A stack is a LIFO, last-in first-out
Can you explain why?

Let's return to a slide we have discussed before

Some questions to ponder:

1. How big should the cache be for our specific example?

At least $2 \times n$ words

2. How does the cache provide $3 \times$ higher bandwidth?

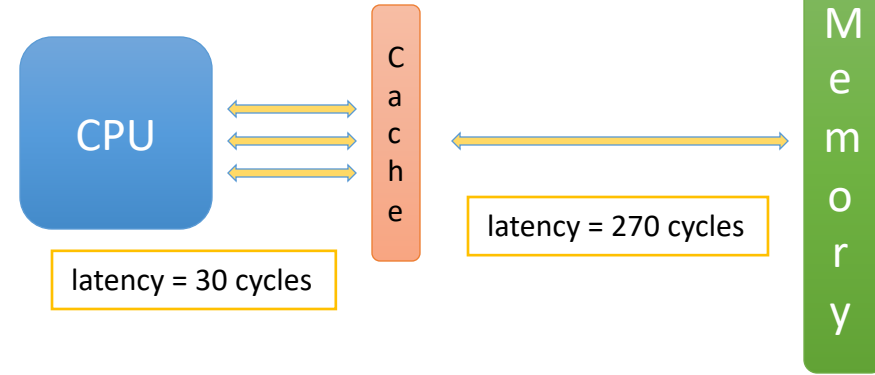
Closer, $3 \times$ as many wires, different hardware

3. How is the cache managed?

Cache is managed by the run-time

Oldest data is evicted (FIFO)

I added this new
information to the
slide



It seems to me that we must revise one of the statements above.

Which one, and why?

Reminder: Which rows are being re-used?

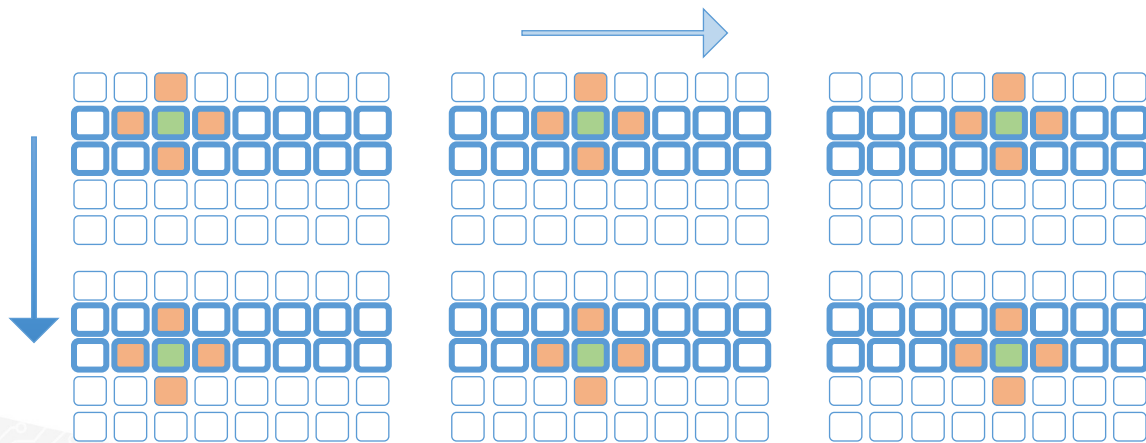
Can we lower the bandwidth requirement?

Similarly, two rows have been loaded in the previous y-loop

Can we use this to our advantage?

What if we could store two rows in a special buffer that provides higher bandwidth?

Such a buffer
would be called a
'cache'



Let's revisit this slide

Some questions to ponder:

1. How big should the cache be for our specific example?

At least $3 \times n$ words

2. How does the cache provide $3 \times$ higher bandwidth?

Closer, $3 \times$ as many wires, different hardware

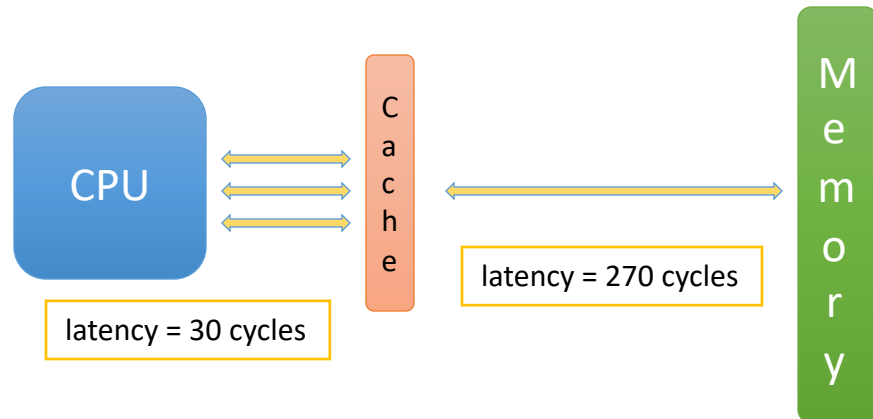
3. How is the cache managed?

Cache is managed by the run-time

Oldest data is evicted (FIFO)

It seems to me that we have to revise one of the statements above.

Which one, and why? Number of rows stored



The first re-use of cached data may happen when we start the second pass

During the second pass we evict data from the cache

The oldest elements are the first (leftmost) elements of the first 3 rows

In a user-managed setup we could just evict elements from the very first row

However, the cache will evict old elements from all rows. Hence we have to have room to store at least $3 \times n$ words, not just $2 \times n$

Summary

Part 1

Keep in mind that we have mostly discussed concepts

Many relevant details of the 3 major hardware features are still missing
We have not discussed the ramifications for code design

Particularly our discussion of caches is very much incomplete!
(and we will continue with caches in the next section)

Main bottlenecks

Transfer the data between memory and CPU

Actual computation

Major technologies to increase **concurrency**

Streaming of data between memory and CPU to hide memory latency and increase memory bandwidth

Pipelining computation (add/mult) to increase compute throughput

Caches to re-use data in order to decrease pressure on the Memory-to-CPU connection
and to also to hide memory latency and to increase memory bandwidth

Concurrency

A single action (flops, memory operation, etc.) can only be so fast

Clock speed, power consumption, speed of light

Increasing the '**concurrency**' is the best (maybe only) way to increase performance substantially