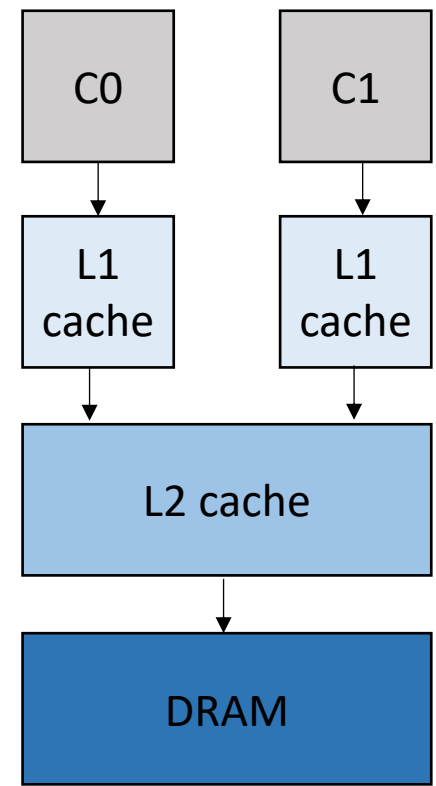


# CSE211: Compiler Design

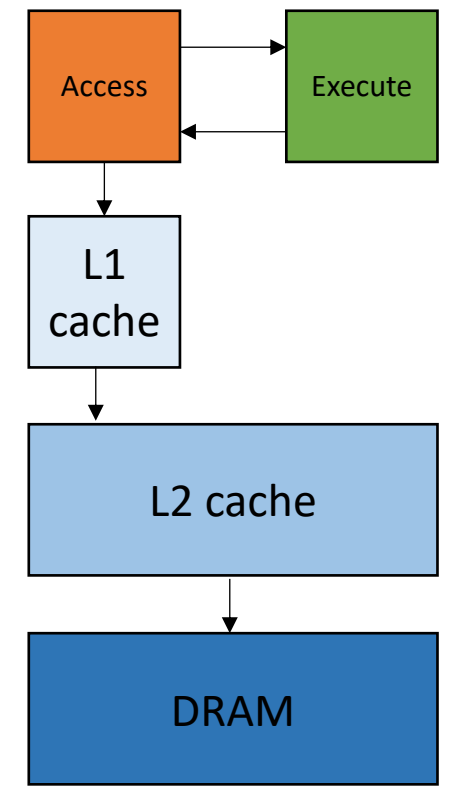
Dec. 1, 2020

- **Topic:** Decoupled Access Execute (DAE)
- **Discussion questions:**
  - What does it mean for an application to be memory bound?
  - What are some techniques for dealing with memory bottlenecks

*Traditional SMP System*



*Decoupled Access/Execute System*



# Announcements

- Midterm Graded; message me if you have any questions
- SETs are out. Please fill it out!
  - Personalized feedback is also welcome
- Projects set!
  - 10-15 minute presentations per paper
  - 1.5 days of presentations
  - papers listed in the announcements (please have a look before presentations)
  - Turn in report on Canvas
- Final will be released Dec. 13 night and due Dec. 14 midnight AoE

# Announcements

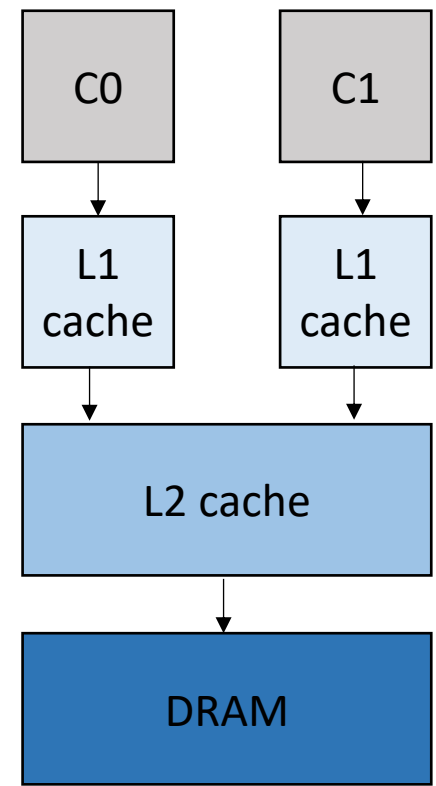
- Homework 3 due on Thursday
  - one more office hour on Wednesday
- Homework 4 released on Thursday. Due on day of Final: Dec. 14

# CSE211: Compiler Design

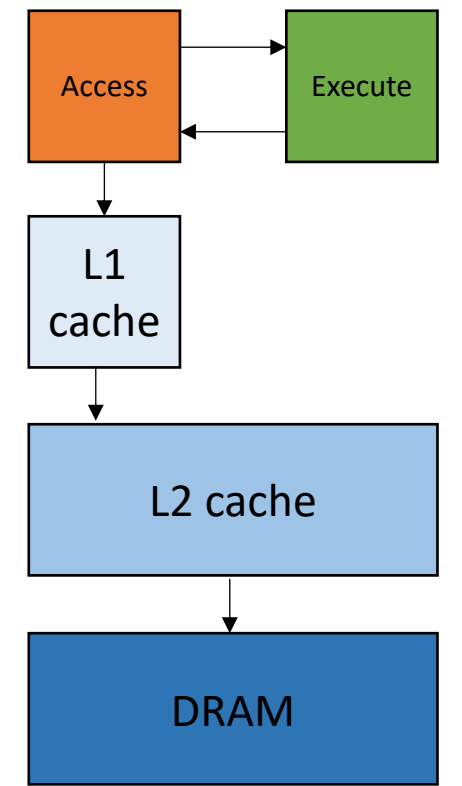
Dec. 1, 2020

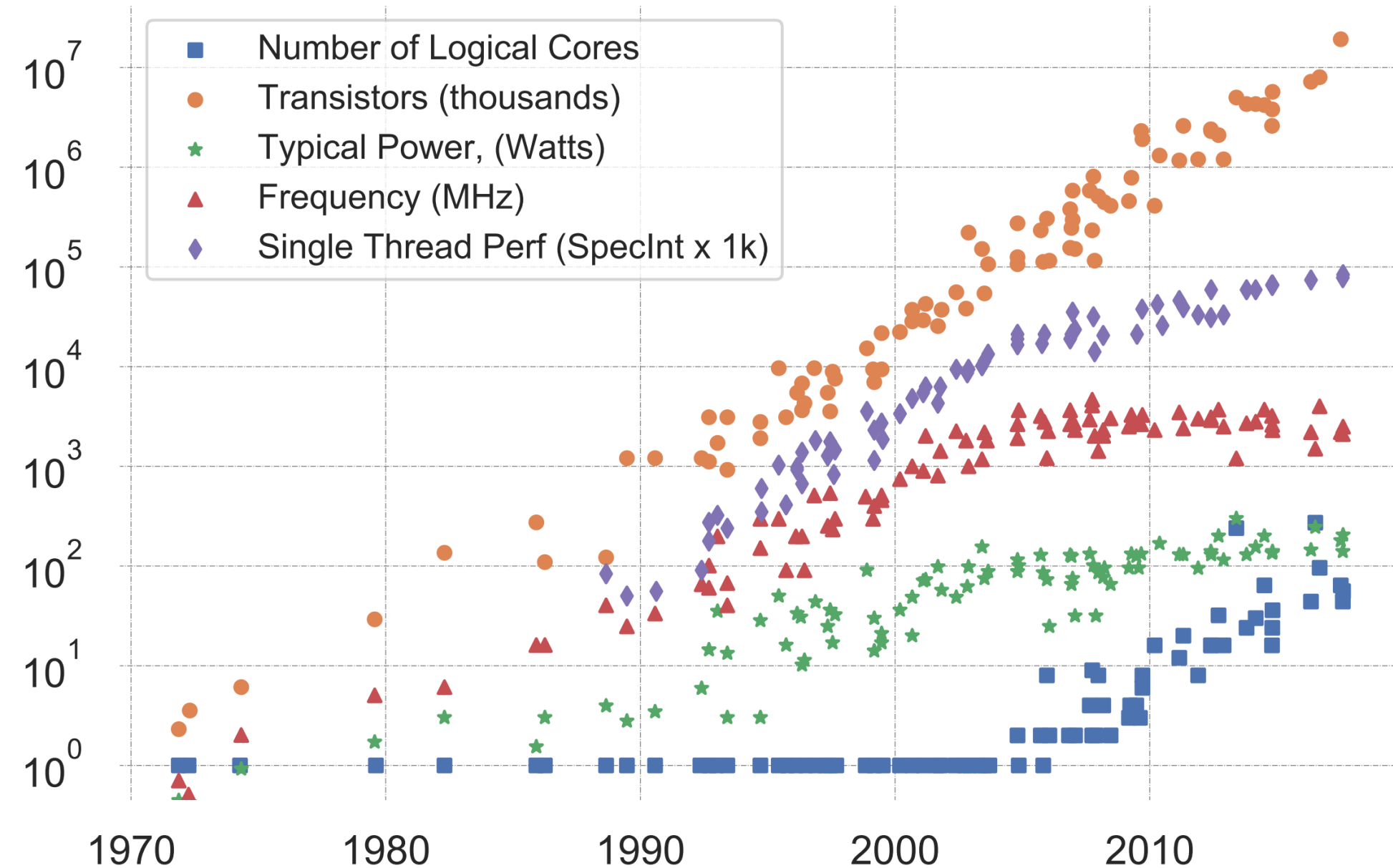
- **Topic:** Decoupled Access Execute (DAE)
- **Discussion questions:**
  - What does it mean for an application to be memory bound?
  - What are some techniques for dealing with memory bottlenecks

*Traditional SMP System*



*Decoupled Access/Execute System*





K. Rupp, "40 Years of Microprocessor Trend Data," <https://www.karlrupp.net/2015/06/40-years-of-microprocessor-trend-data>, 2015.

# Trends

- Frequency scaling: **Dennard's scaling**
  - Mostly agreed that this is over (2005 -2007)
- Number of transistors: **Moore's law**
  - On its last legs.
  - Intel delaying 7nm chips. Apple has a 5nm. Some roadmaps project up to 3nm
- *Chips are not increasing in raw frequency, and space is becoming more valuable*

# Specialization example

- CPUs:
  - Aim to be good at general tasks
  - poor area and energy utilization

# Specialization example

How many floating point operations per second (FLOPS) on matrix multiplication

2 TFLOPS

- CPUs:
  - Aim to be good at general tasks
  - poor area and energy utilization



# Specialization example

How many floating point operations per second (FLOPS) on matrix multiplication

2 TFLOPS

- CPUs:
  - Aim to be good at general tasks
  - poor area and energy utilization
- GPUs:
  - Good at regular, uniform parallelism
  - Bad at irregular parallelism and programs with control dependencies

# Specialization example

- CPUs:
  - Aim to be good at general tasks
  - poor area and energy utilization
- GPUs:
  - Good at regular, uniform parallelism
  - Bad at irregular parallelism and programs with control dependencies

How many floating point operations per second (FLOPS) on matrix multiplication

2 TFLOPS

125 TFLOPS  
(62x faster than CPU)

# Specialization example

How many floating point operations per second (FLOPS) on matrix multiplication

- CPUs:
  - Aim to be good at general tasks
  - poor area and energy utilization
- GPUs:
  - Good at regular, uniform parallelism
  - Bad at irregular parallelism and programs with control dependencies
- TPUs:
  - Good at matrix multiplication
  - Not good at much else (12 instructions)

2 TFLOPS

125 TFLOPS  
(62x faster than CPU)

# Specialization example

How many floating point operations per second (FLOPS) on matrix multiplication

- CPUs:

- Aim to be good at general tasks
- poor area and energy utilization

2 TFLOPS

- GPUs:

- Good at regular, uniform parallelism
- Bad at irregular parallelism and programs with control dependencies

125 TFLOPS  
(62x faster than CPU)

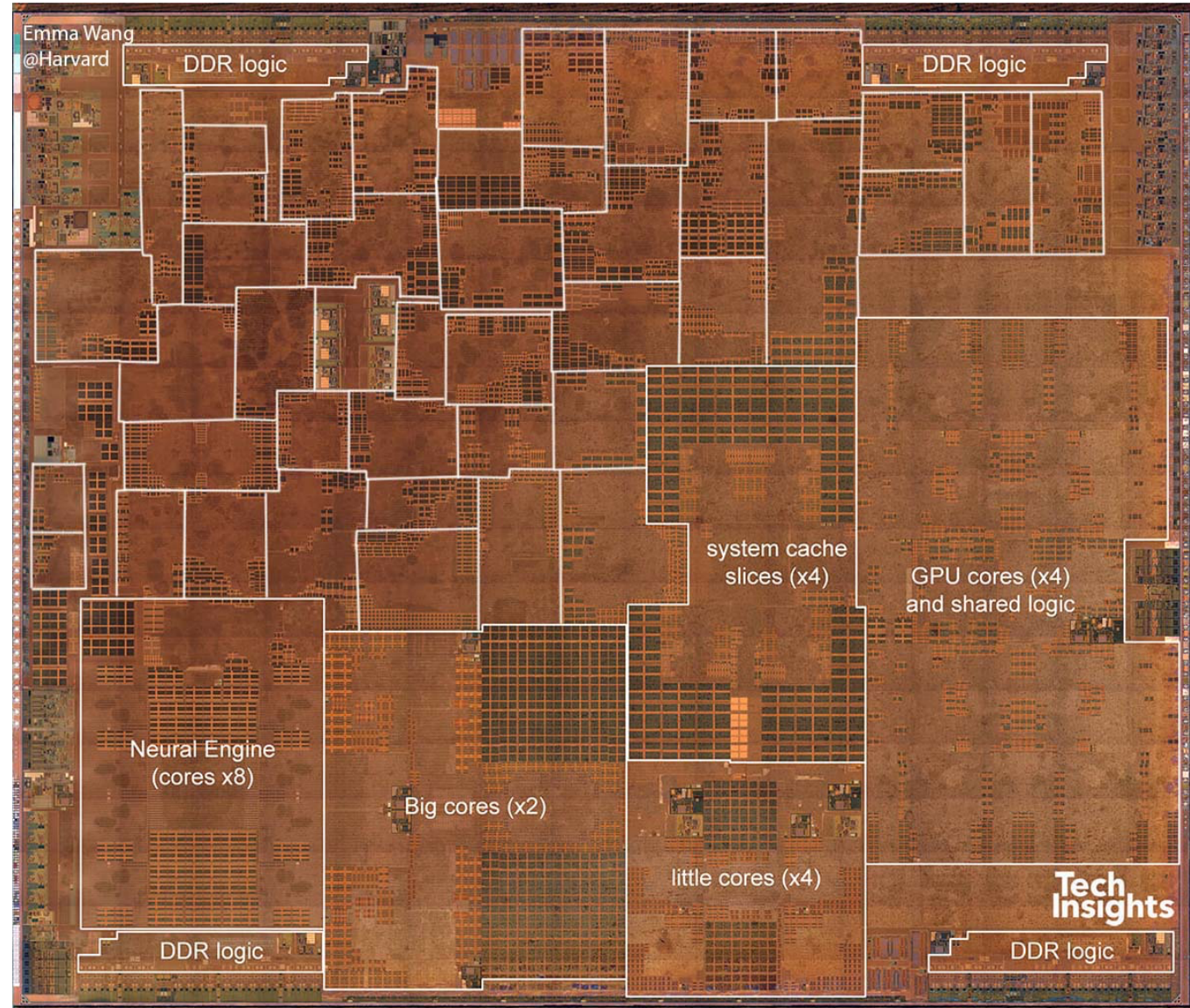
- TPUs:

- Good at matrix multiplication
- Not good at much else (12 instructions)

180 TFLOPS  
(much faster than CPU,  
1.4x faster than GPU)

# Specialization in modern SoCs

- From David Brooks lab at Harvard:  
<http://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis/>
- CPUs, GPUs, Neural Engine, IP blocks (cryptography, DSP, etc.)



# How do programs take advantage of specialization?

- Programmer-centric: Programmers write code using a specific API
- Hardware-centric: Hardware transparently targets optimized hardware
- Compiler-centric: Compiler performs non-trivial transformations to target specialized hardware

# How do programs take advantage of specialization?

- **Programmer-centric:**

- Programmers write code using a specific API
- e.g. Tensorflow targets CPU, GPU, TPU

- **Hardware-centric:**

- Hardware transparently targets optimized hardware
- Pipelining, super scalar, caches, etc. (what we already do!)

- **Compiler-centric:**

- Compiler performs non-trivial transformations to target specialized hardware

# Specialization is not new

- First GPU in 1951 (MIT flight simulator)
- Architecture academic work proposes many new designs
  - Evaluated on detailed simulators; rarely taped out
- Had a hard time breaking into the mainstream:
  - benefits had to outweigh eventual returns from Dennard's Scaling and Moore's Law
- But now...
  - Hennessy and Patterson's 2017 Turing award lecture: The New Golden Age of Computer Architecture



# Decoupled Access/Execute (DAE)

- 1982: James E. Smith
  - Lives in Montana now and gives interesting keynotes at architecture conferences

DECOUPLED ACCESS/EXECUTE COMPUTER ARCHITECTURES  
James E. Smith  
Department of Electrical and Computer Engineering  
University of Wisconsin-Madison, Madison, Wisconsin 53706

## Abstract

An architecture for improving computer performance is presented and discussed. The main feature of the architecture is a high degree of decoupling between operand access and execution. This results in an implementation which has two separate instruction streams that communicate via queues. A similar architecture has been previously proposed for array processors, but in that context the software is called on to do most of the coordination and synchronization between the instruction streams. This paper emphasizes implementation features that remove this burden from the programmer. Performance comparisons with a conventional scalar architecture are given, and these show that considerable performance gains are possible.

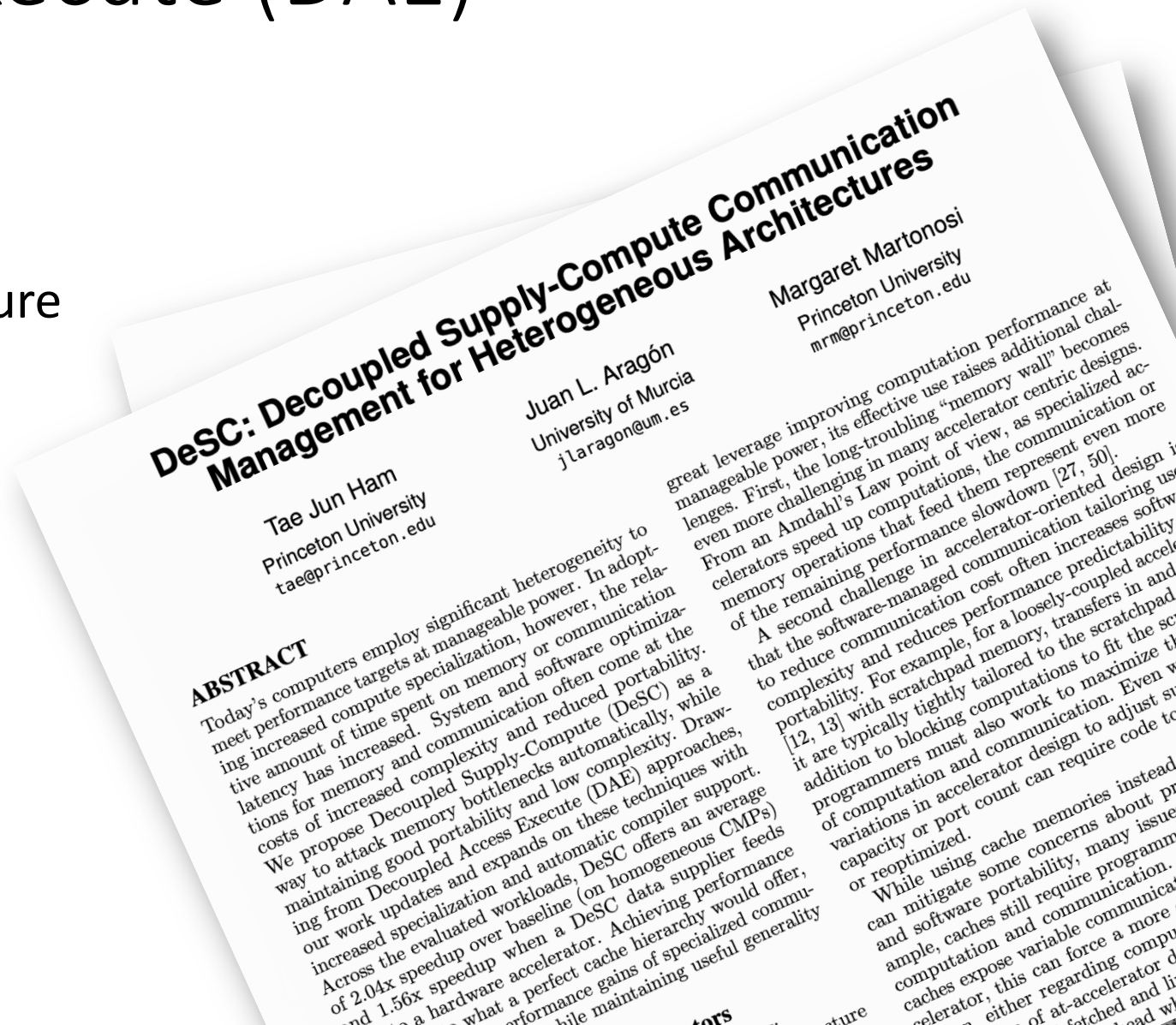
Single instruction stream versions, both physical and conceptual, are discussed with the primary goal of minimizing the differences with conventional architectures. This would allow known compilation and programming techniques to be used. Finally, the problem of deadlock in such a system is discussed, and one possible solution is

A second critical constraint on performance is time required for processor-memory communication. Current trends, both in hardware and software, tend to aggravate the memory communication problem. In hardware, the trend toward higher levels of integration has the effect of increasing the performance impact of all forms of inter-chip communication, including processor-memory communication. At the architectural level, the trend is toward elaborate virtual memory protection methods. These tend to slow address translation because of the required address communication and protection checks. The use of multiprocessors often means that individual processors must contend for memory resources, both due to their size and additional interconnection structures added. Cache memory becomes a less effective solution in multiprocessor systems due to contention in maintaining coherence. At the system level, facilities for defining elaborate types and structures are being developed, which causes an increase in the number of types needed to check types, compute indices, of which adds to increased delay when data. All of the above point to the processors that can diminish the effect of increased memory communication time.

This paper discussed a new type architecture which separates its processor into two parts: access to memory to fetch store results, and operand execution. By architecturally separating the execution, it is possible that execution that is provided by

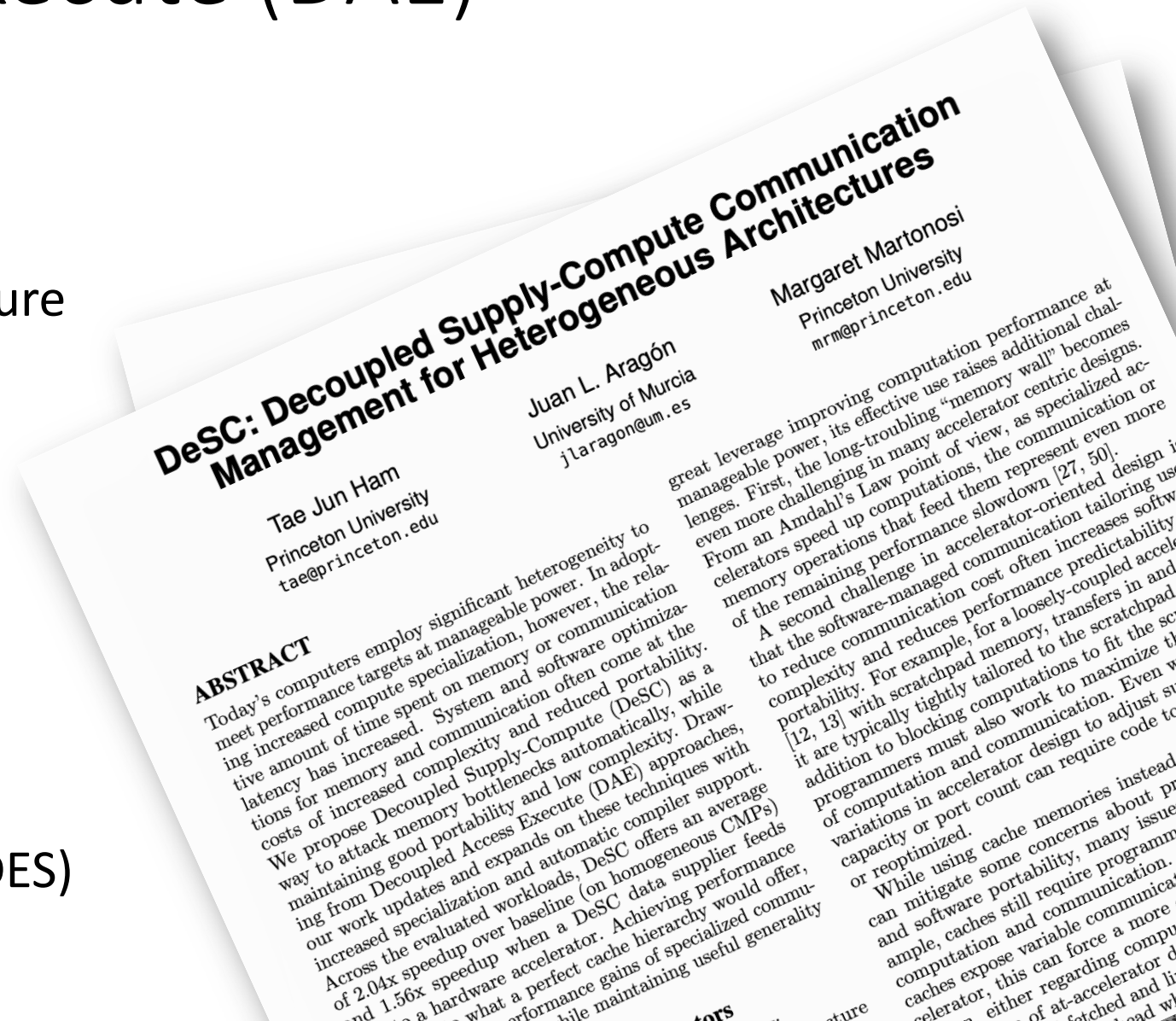
# Decoupled Access/Execute (DAE)

- 1982: James E. Smith
  - Lives in Montana now and gives interesting keynotes at architecture conferences
- 2015: DeSC by Ham et al.
  - More optimizations and practicalities



# Decoupled Access/Execute (DAE)

- 1982: James E. Smith
  - Lives in Montana now and gives interesting keynotes at architecture conferences
- 2015: DeSC by Ham et al.
  - More optimizations and practicalities
- Do chips exist yet?
  - Working on it! (Princeton DECADES)



# DAE - motivation

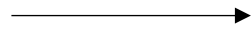
*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```

# DAE - motivation

*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



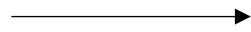
*pseudo 3-address code*

```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```

# DAE - motivation

*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



*pseudo 3-address code*

```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```

core 0

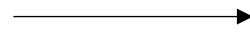


*time*

# DAE - motivation

*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



*pseudo 3-address code*

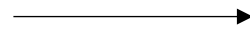
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



# DAE - motivation

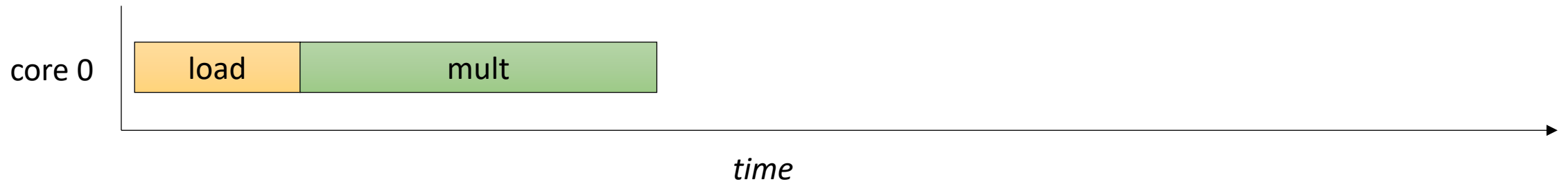
*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



*pseudo 3-address code*

```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```

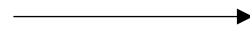




# DAE - motivation

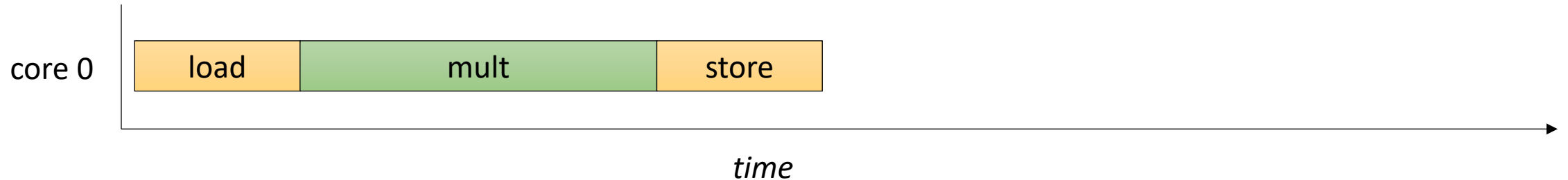
*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



*pseudo 3-address code*

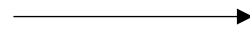
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



# DAE - motivation

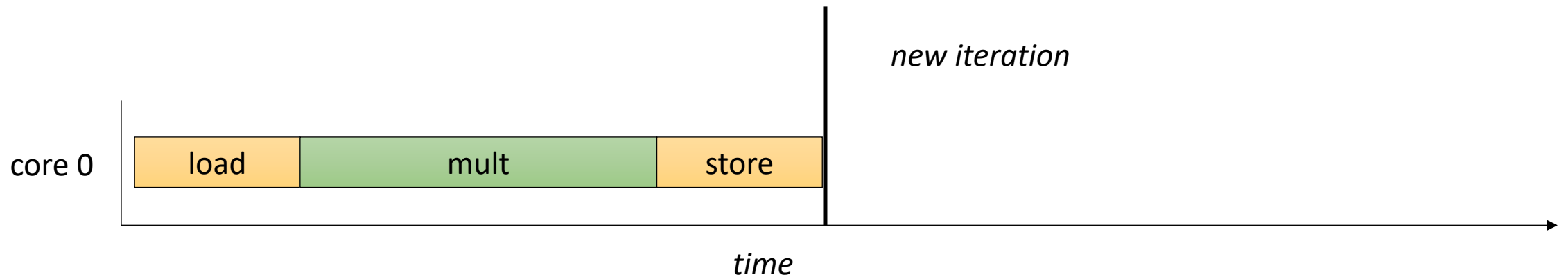
*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



*pseudo 3-address code*

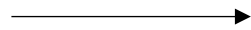
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



# DAE - motivation

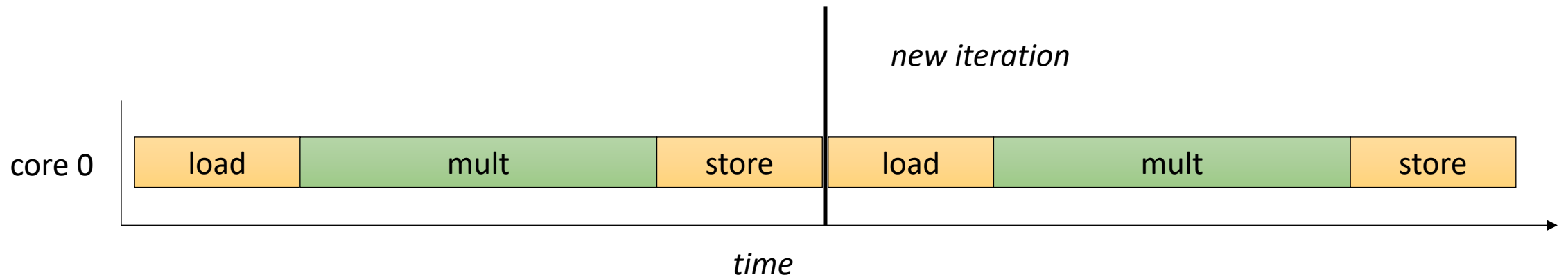
*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



*pseudo 3-address code*

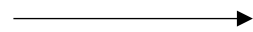
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



# DAE - motivation

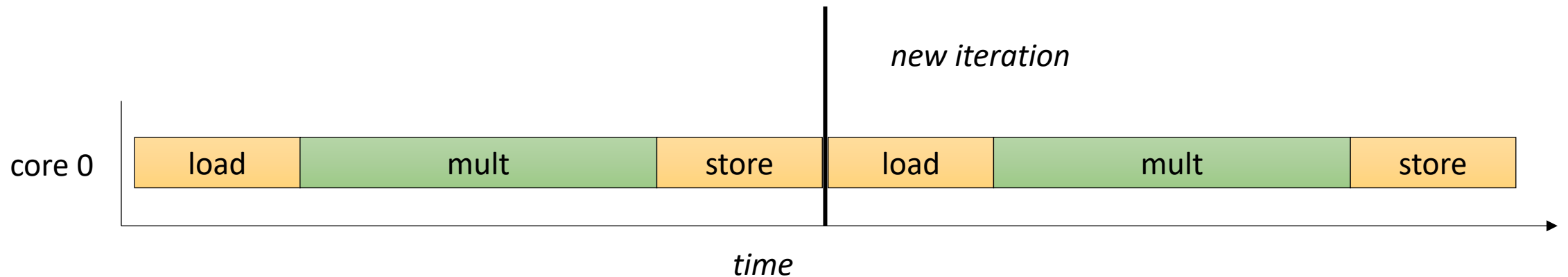
*simple example program*

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```

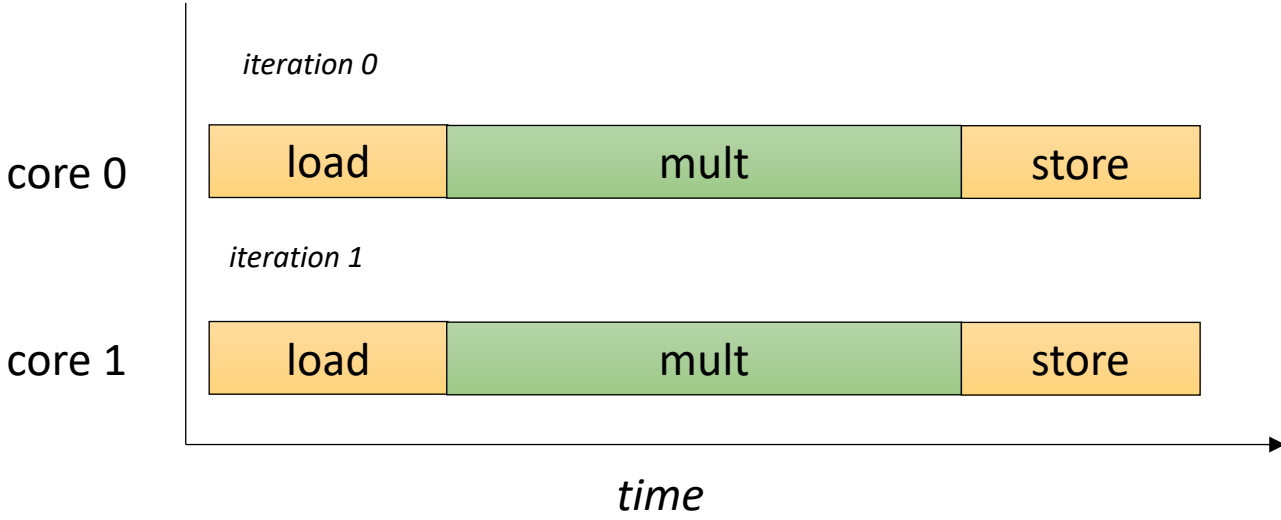


*pseudo 3-address code*

```
#pragma parallel  
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



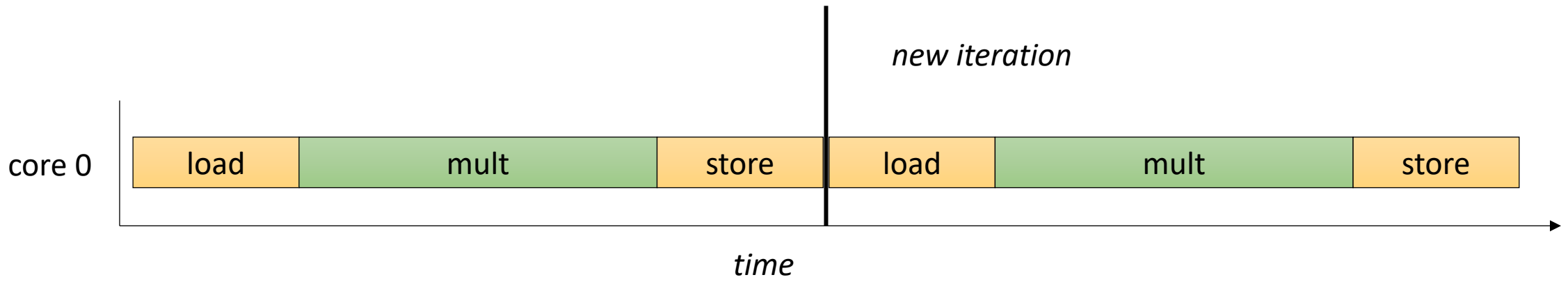
# DAE - motivation



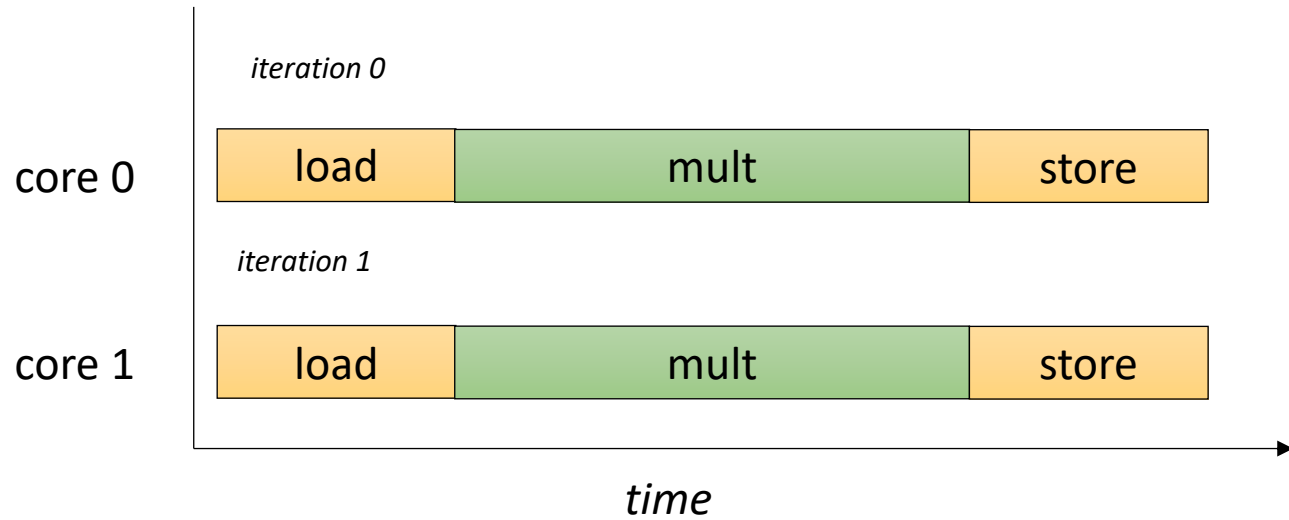
homogeneous SMP parallelism

*pseudo 3-address code*

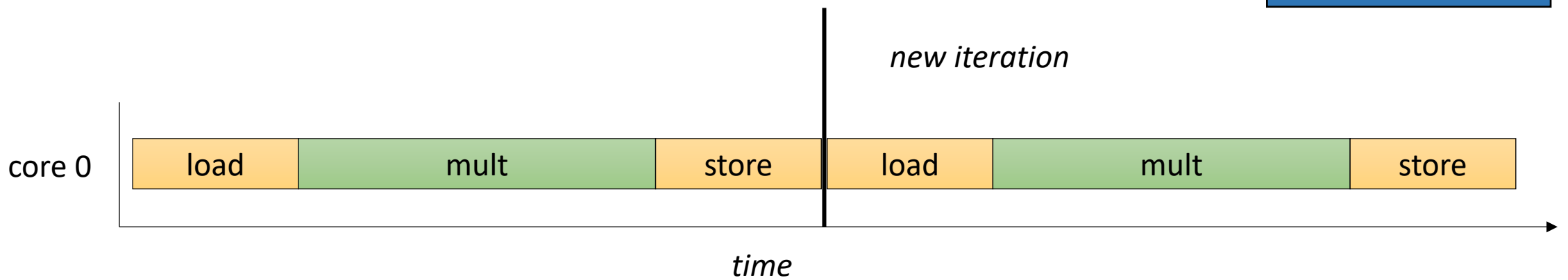
```
#pragma parallel
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```



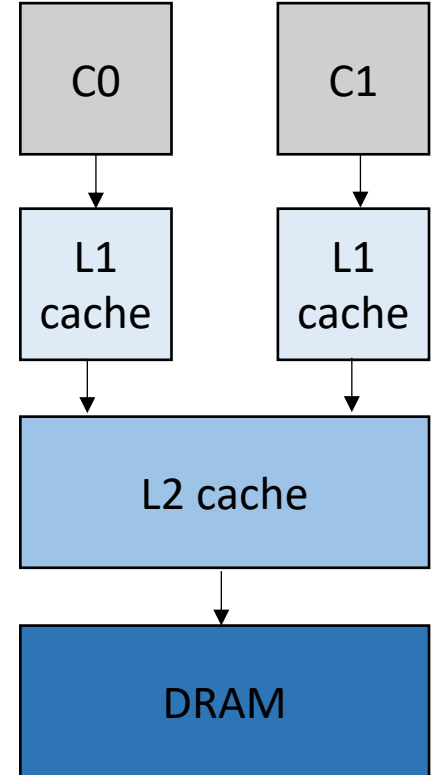
# DAE - motivation



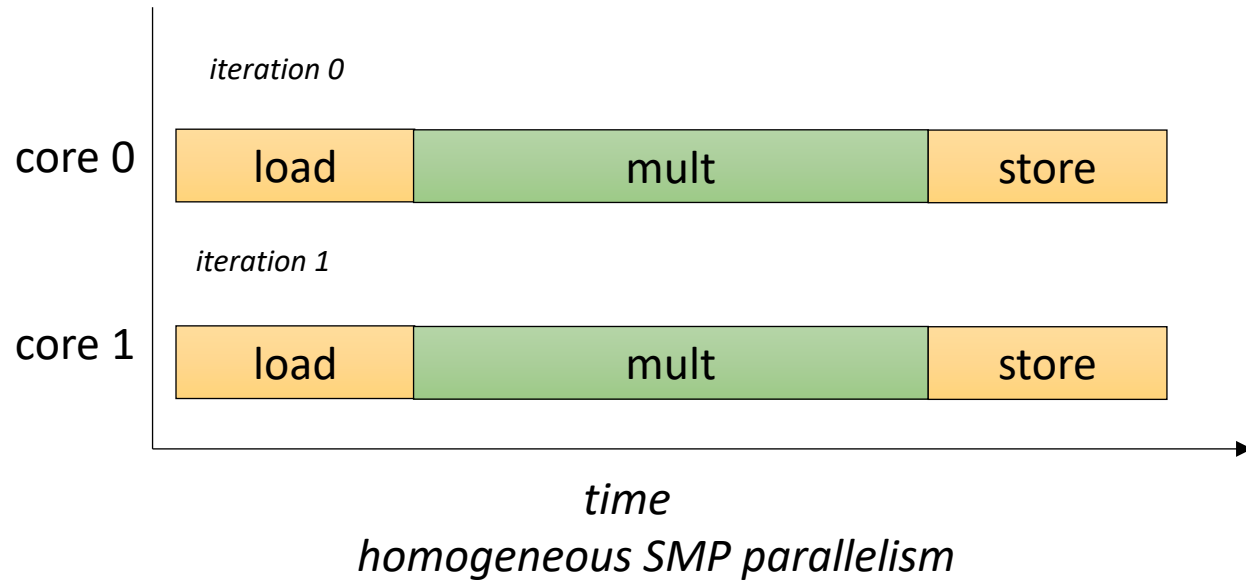
*homogeneous SMP parallelism*



*Traditional SMP System*

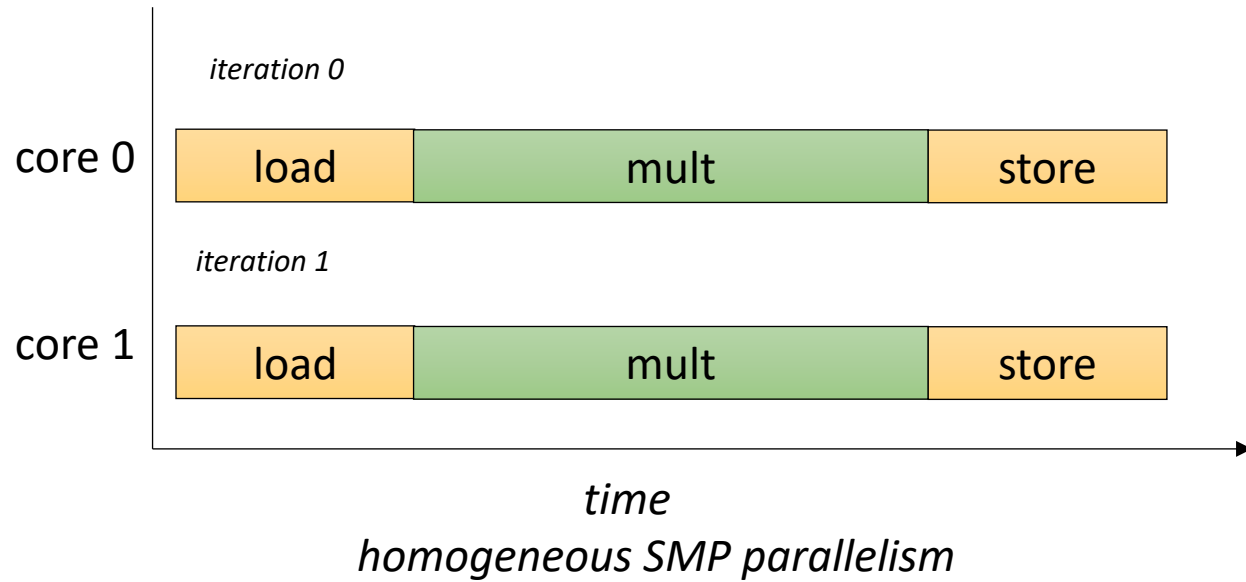


# DAE - illustration

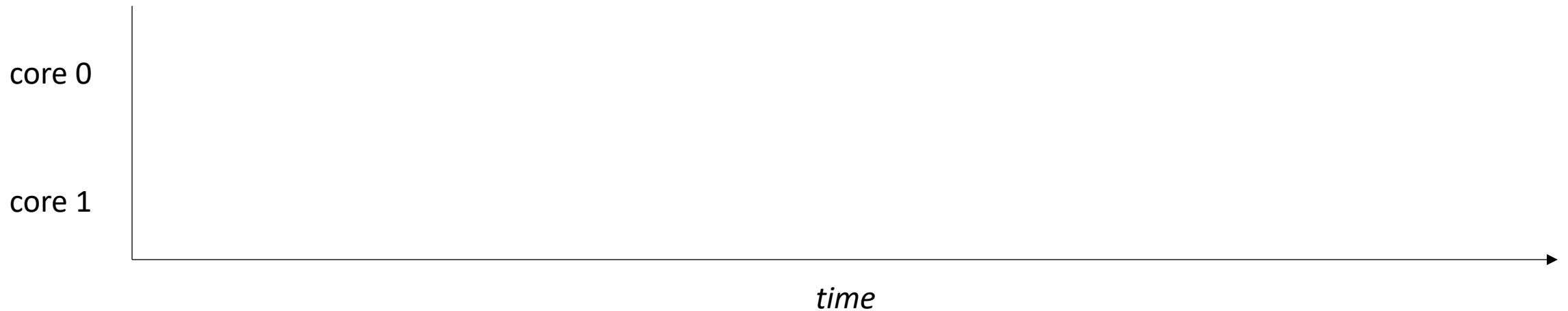


*DAE: split into heterogeneous parallelism: one core does memory and one does computation*

# DAE - illustration

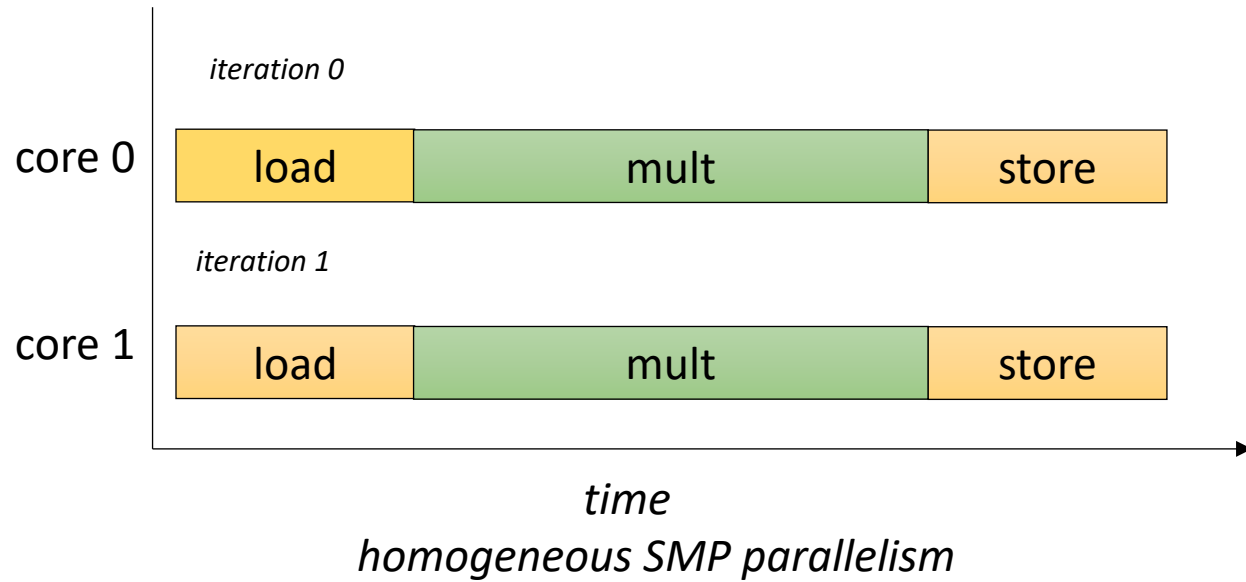


*DAE: split into heterogeneous parallelism: one core does memory and one does computation*

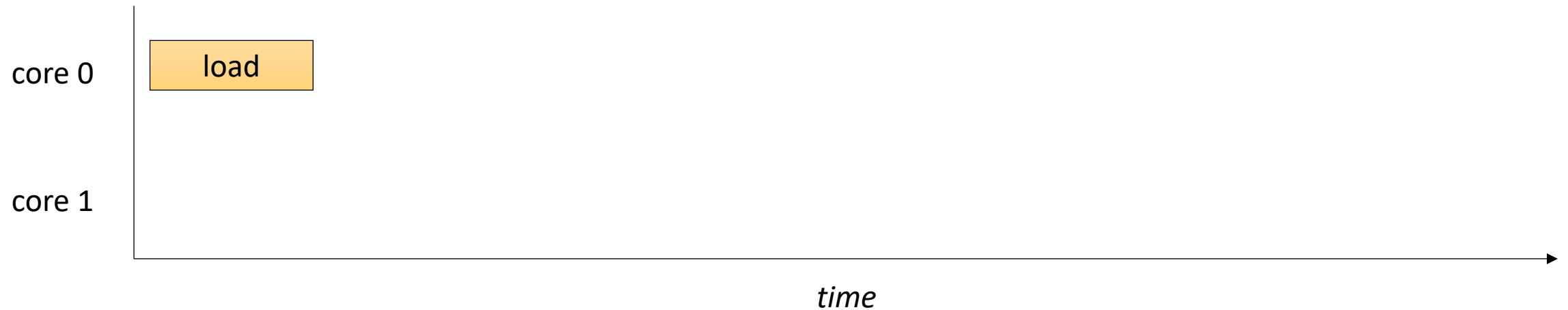




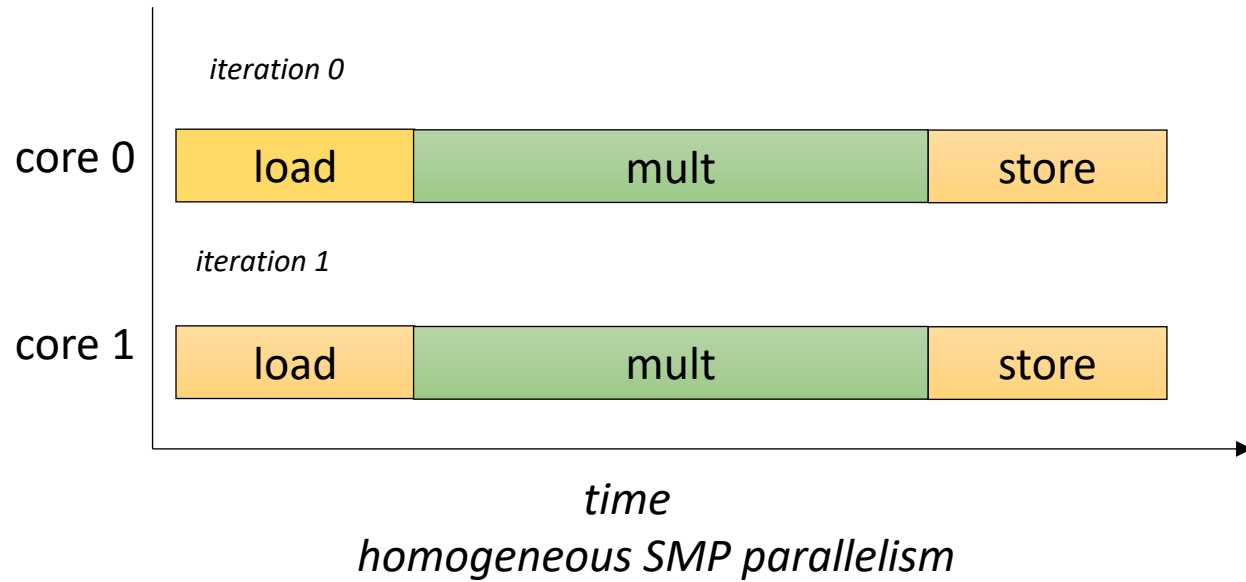
# DAE - illustration



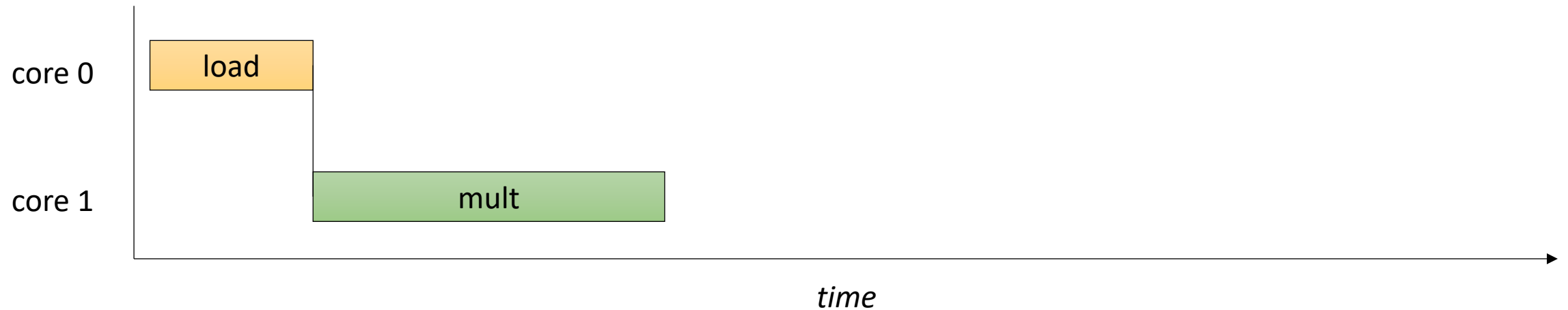
*DAE: split into heterogeneous parallelism: one core does memory and one does computation*



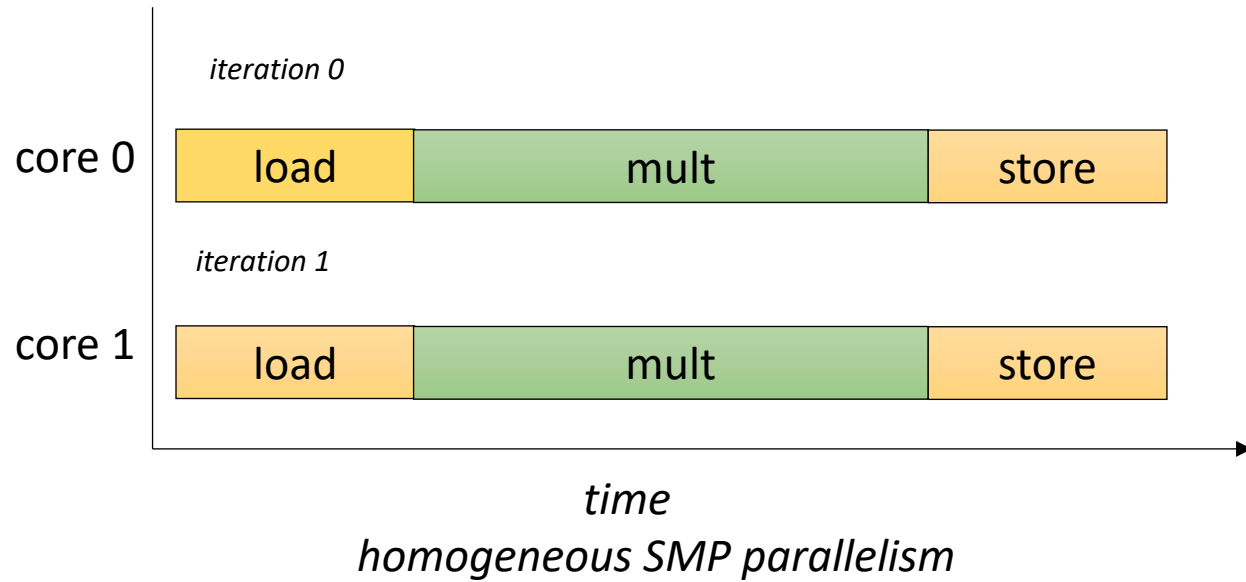
# DAE - illustration



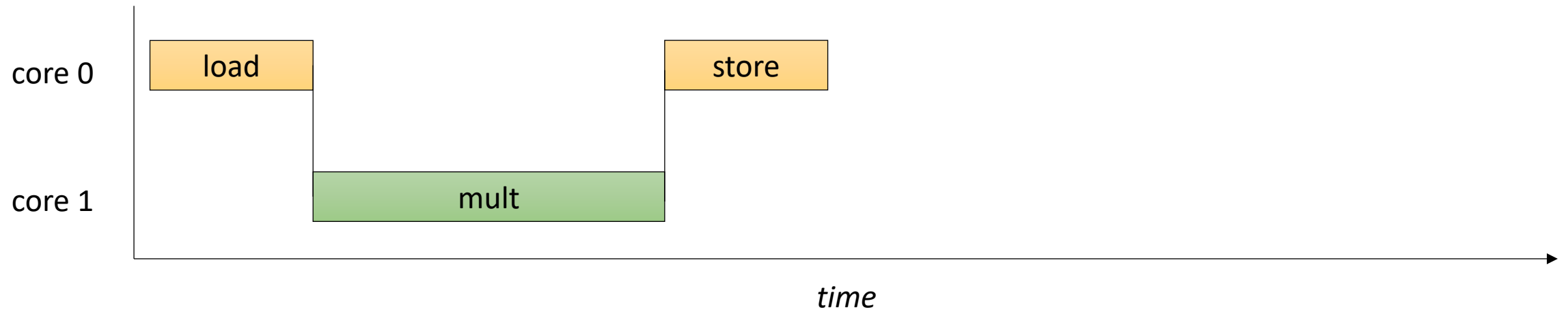
*DAE: split into heterogeneous parallelism: one core does memory and one does computation*



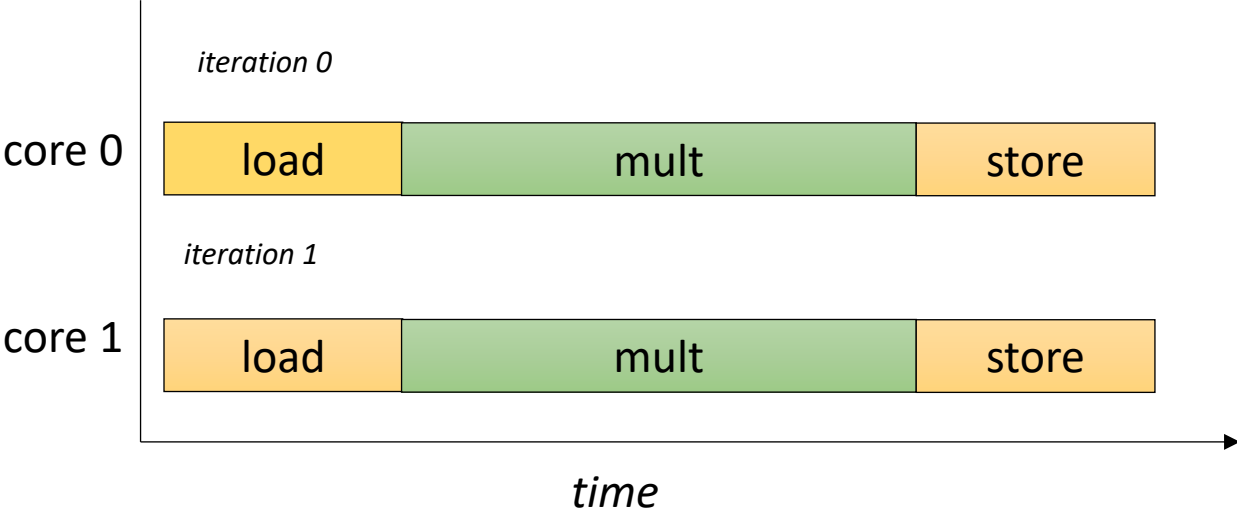
# DAE - illustration



*DAE: split into heterogeneous parallelism: one core does memory and one does computation*

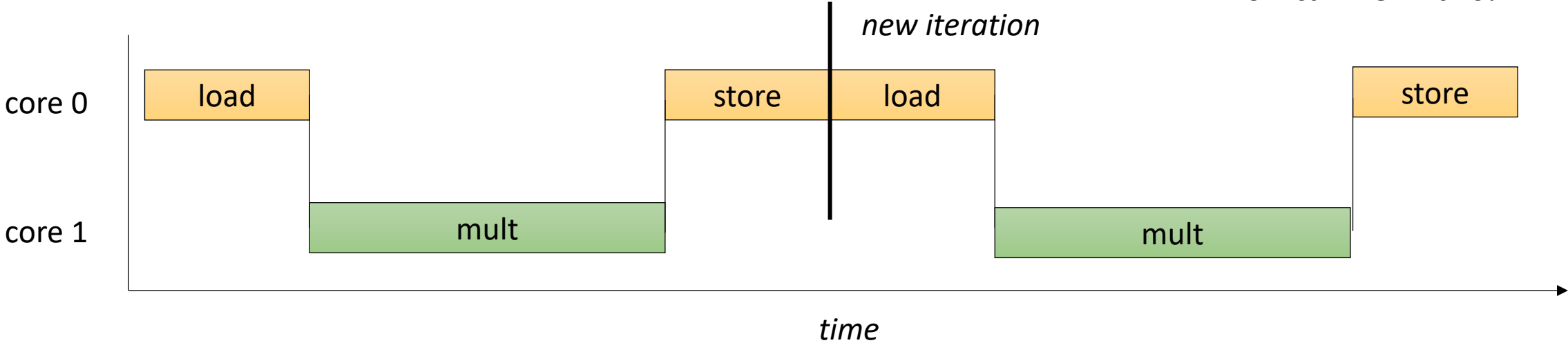


# DAE - illustration



*DAE: split into heterogeneous parallelism: one core does memory and one does computation*

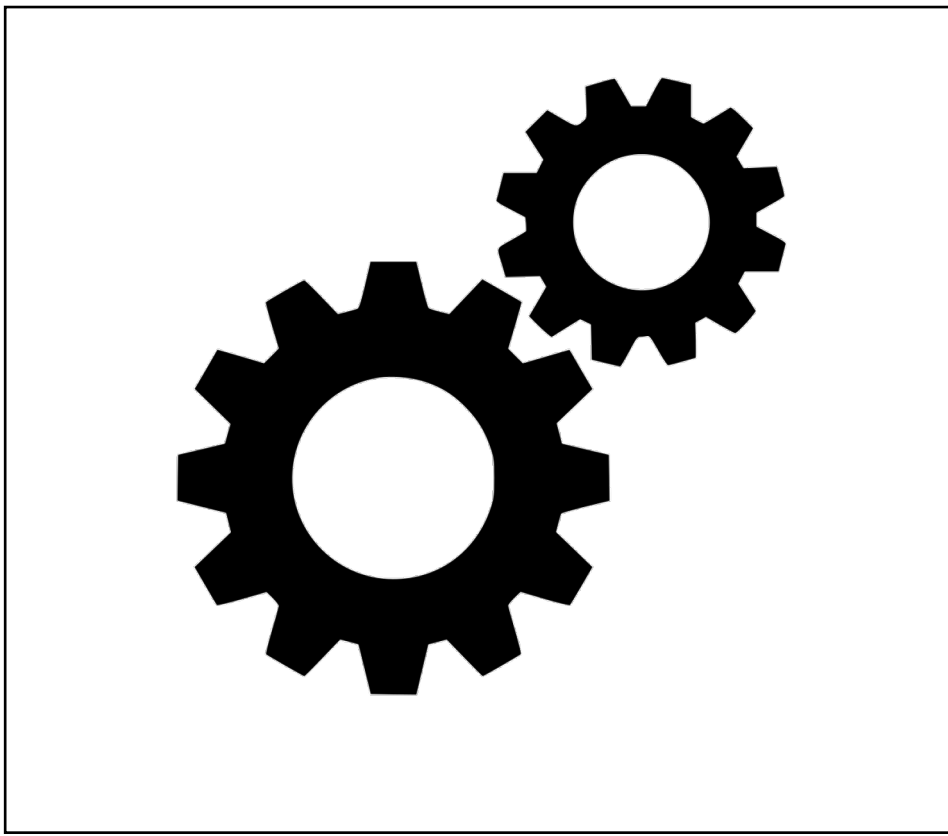
*homogeneous SMP parallelism*



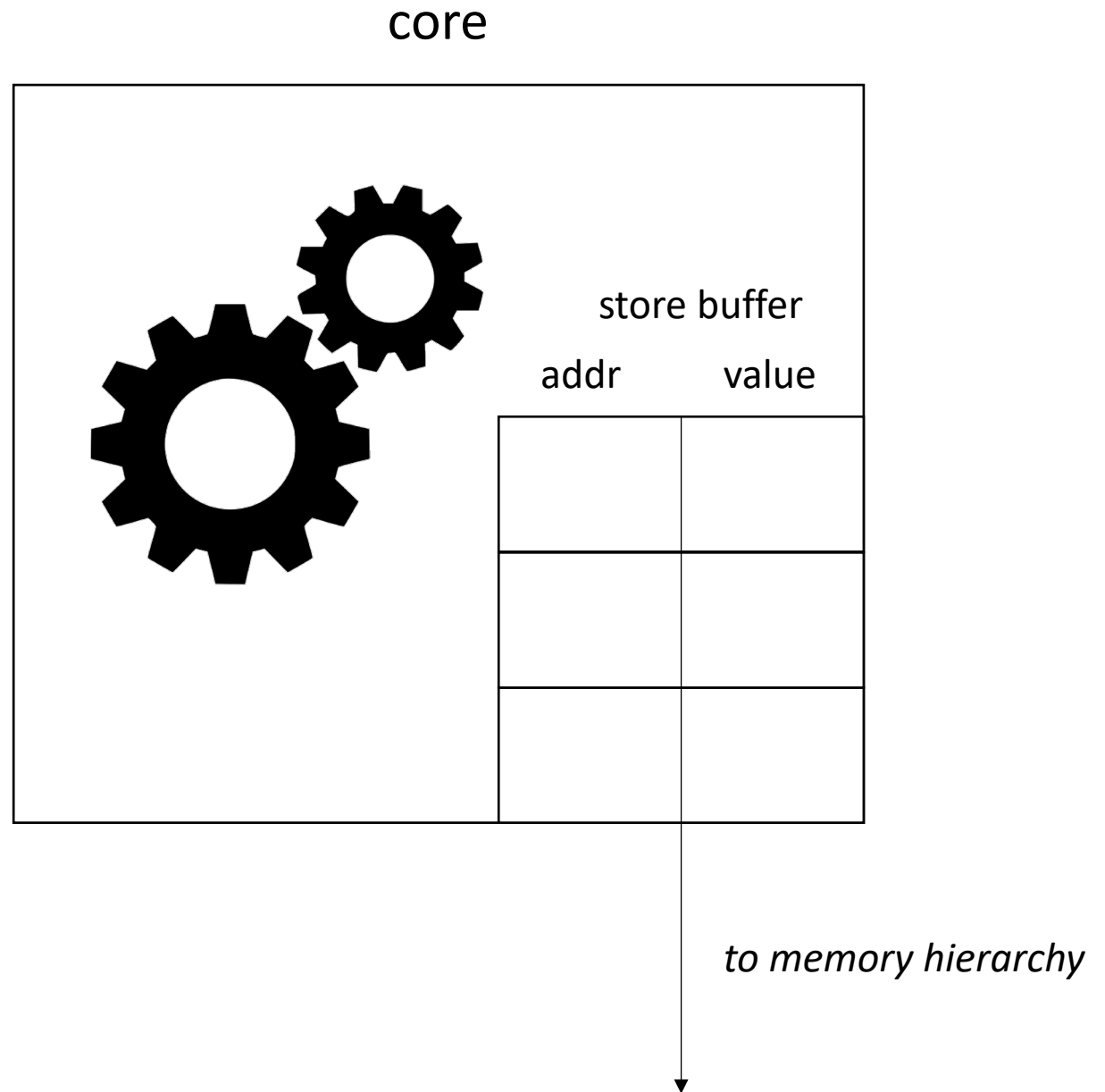
*This is sequentialized ☹️  
How can we fix this?*

# Store Buffers

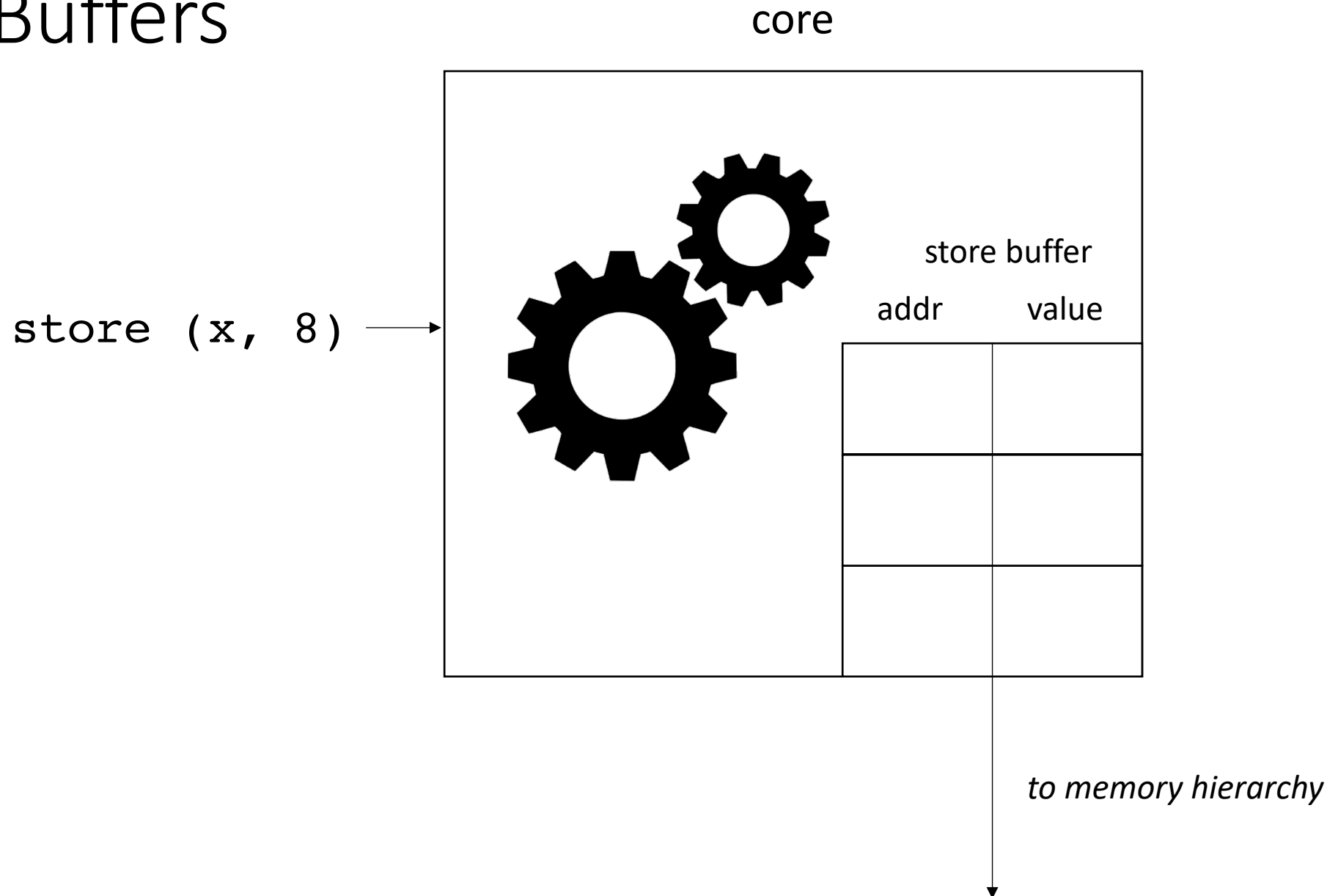
core



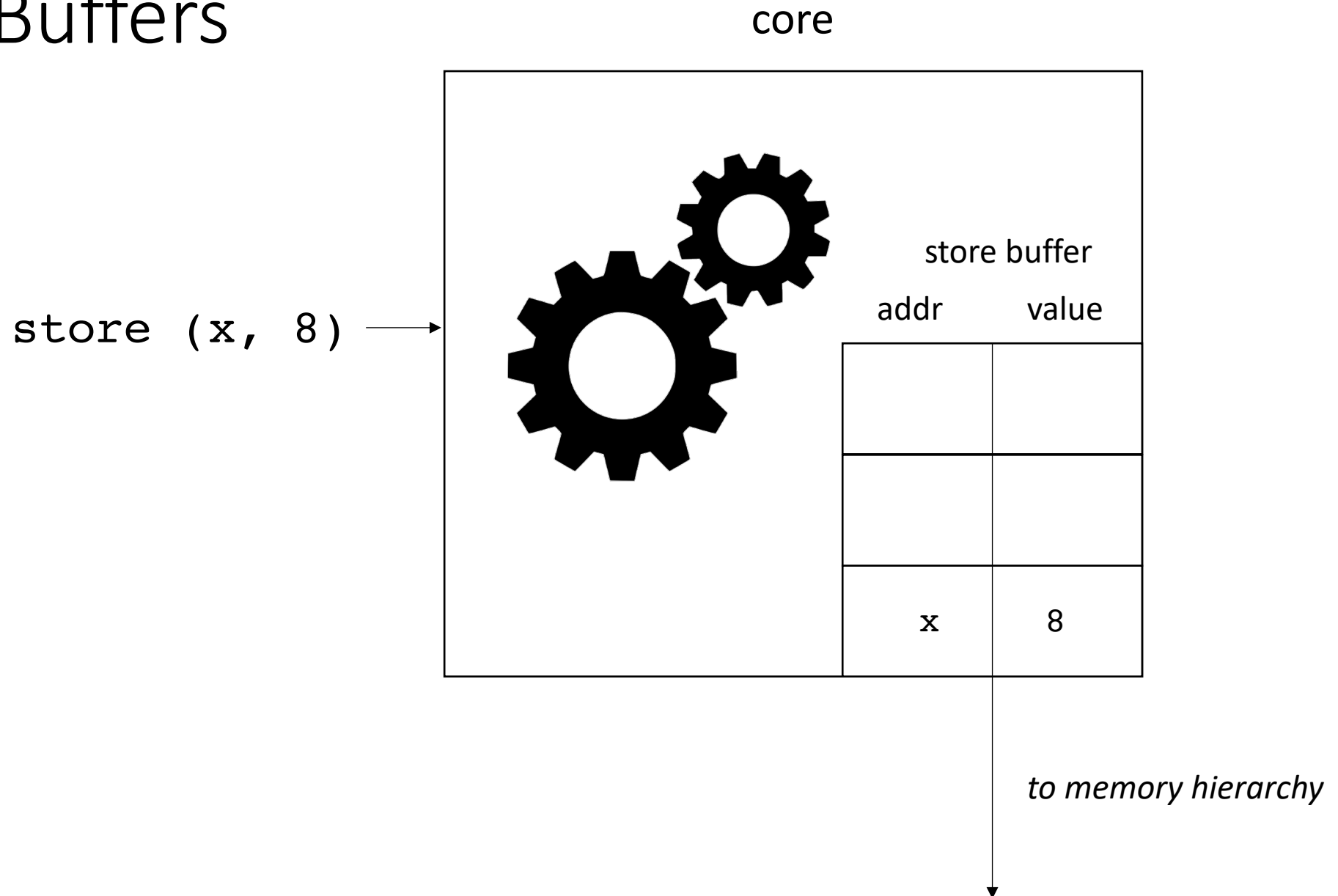
# Store Buffers



# Store Buffers

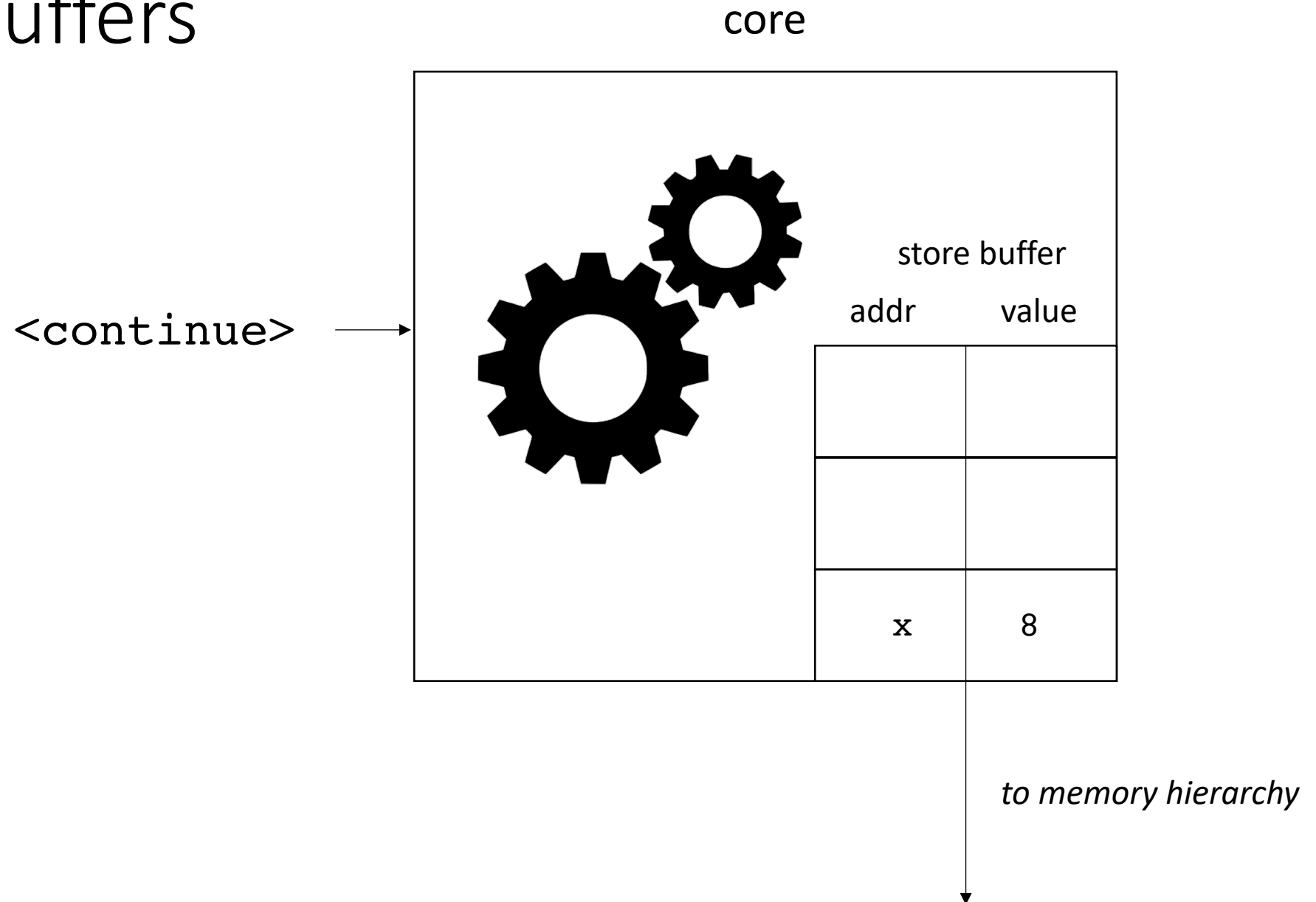


# Store Buffers

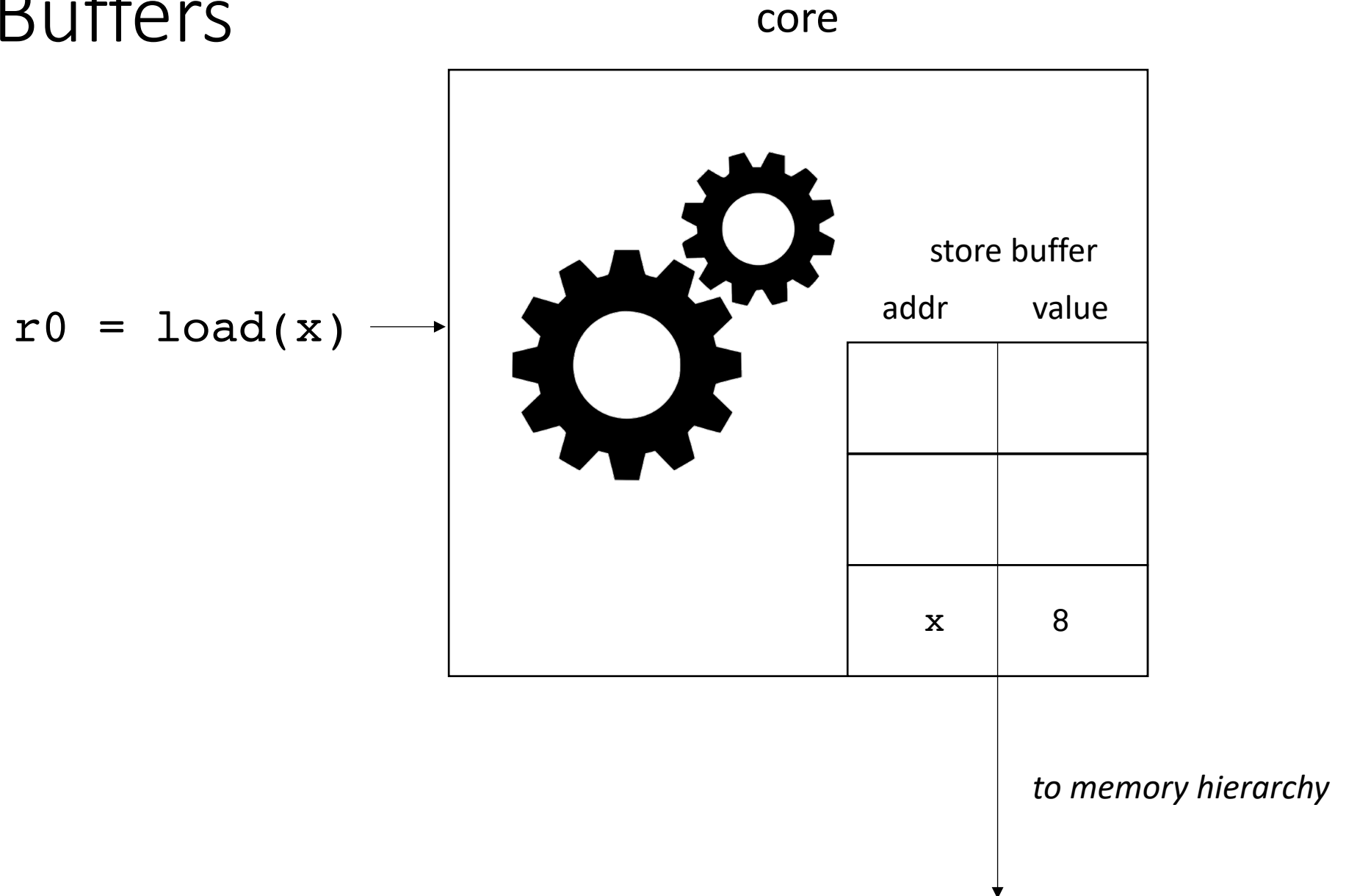




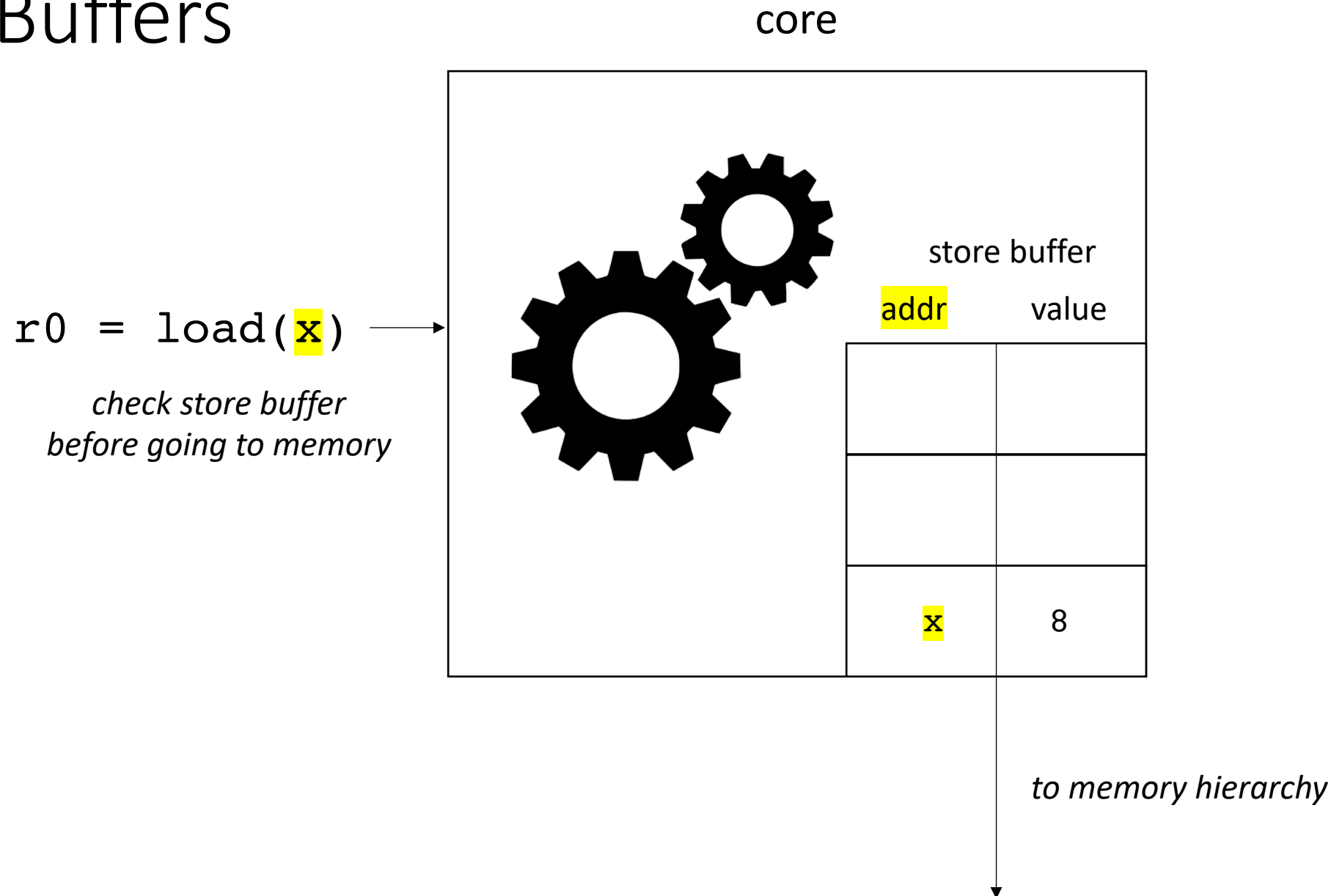
# Store Buffers



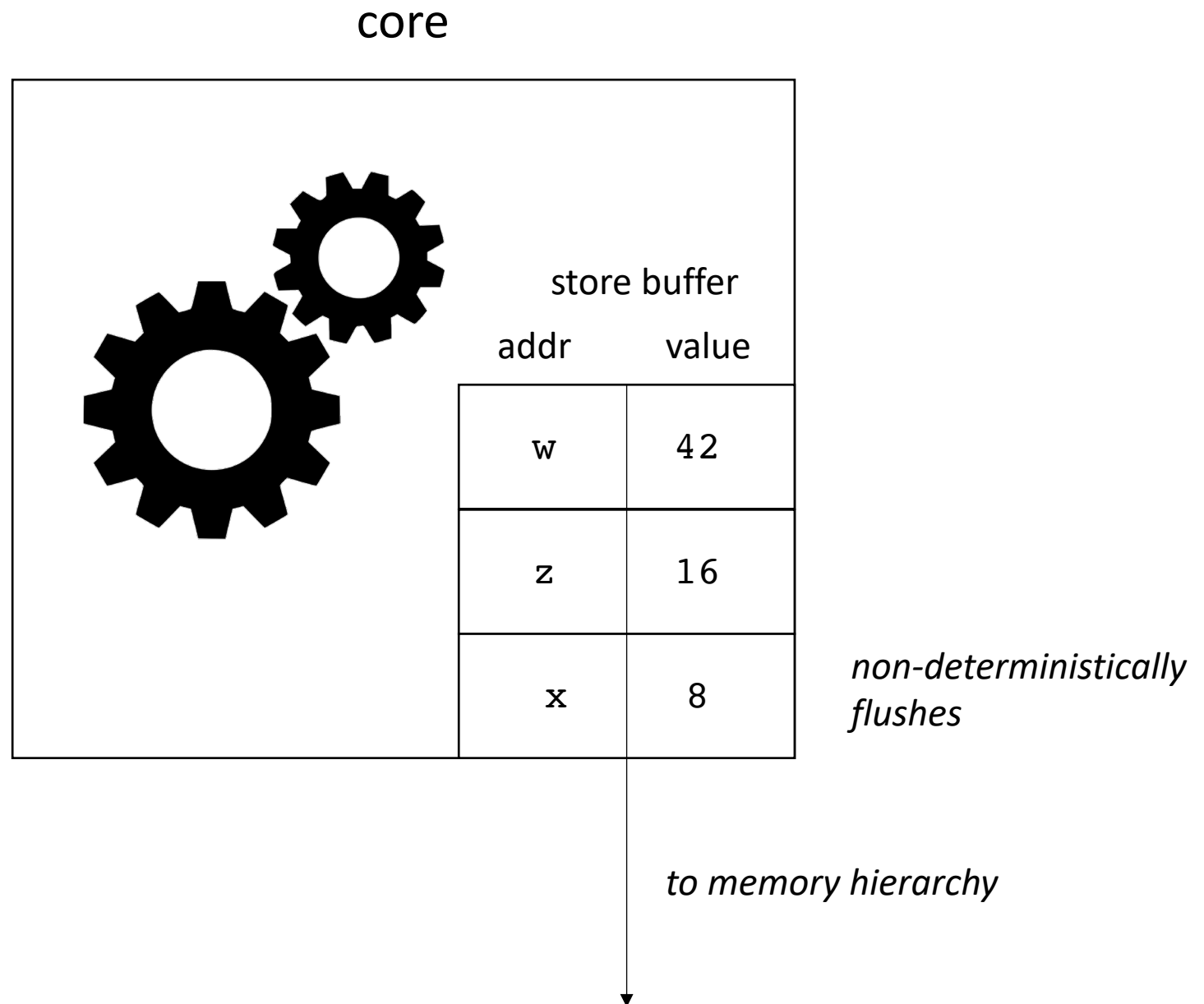
# Store Buffers



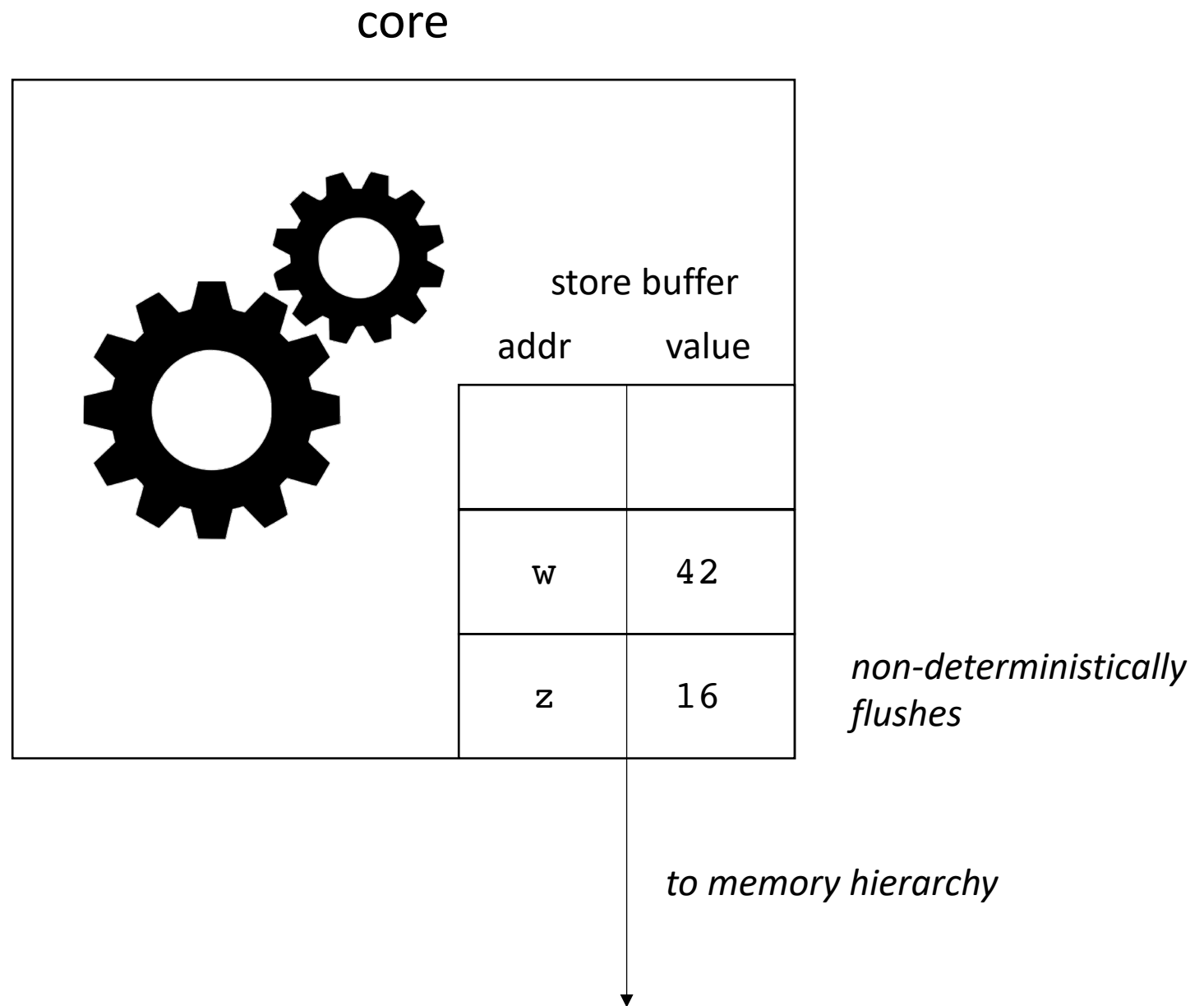
# Store Buffers



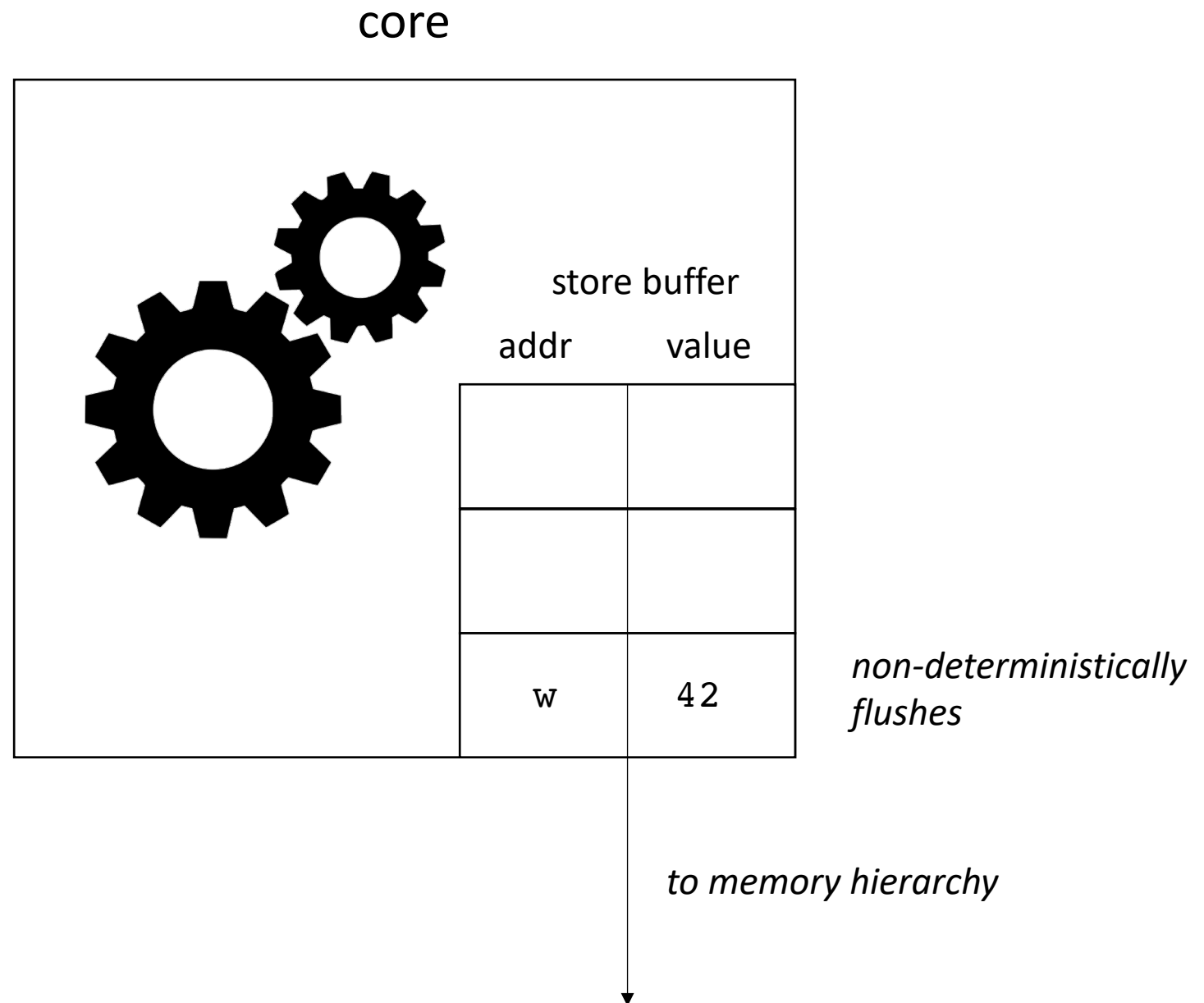
# Store Buffers



# Store Buffers

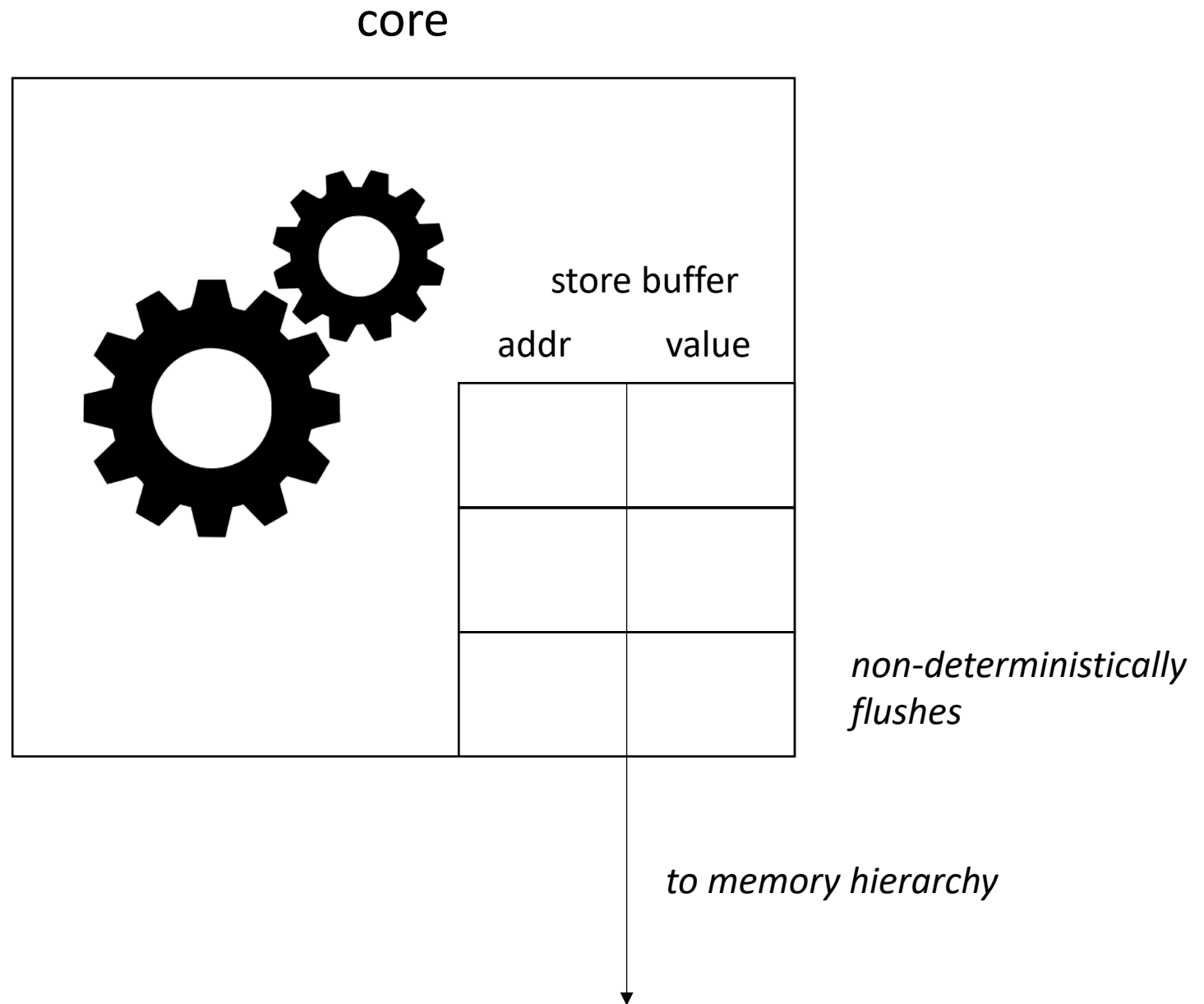


# Store Buffers



# Store Buffers

*key insight:  
Store buffers allow  
asynchronous stores!*



# DAE Parallelism

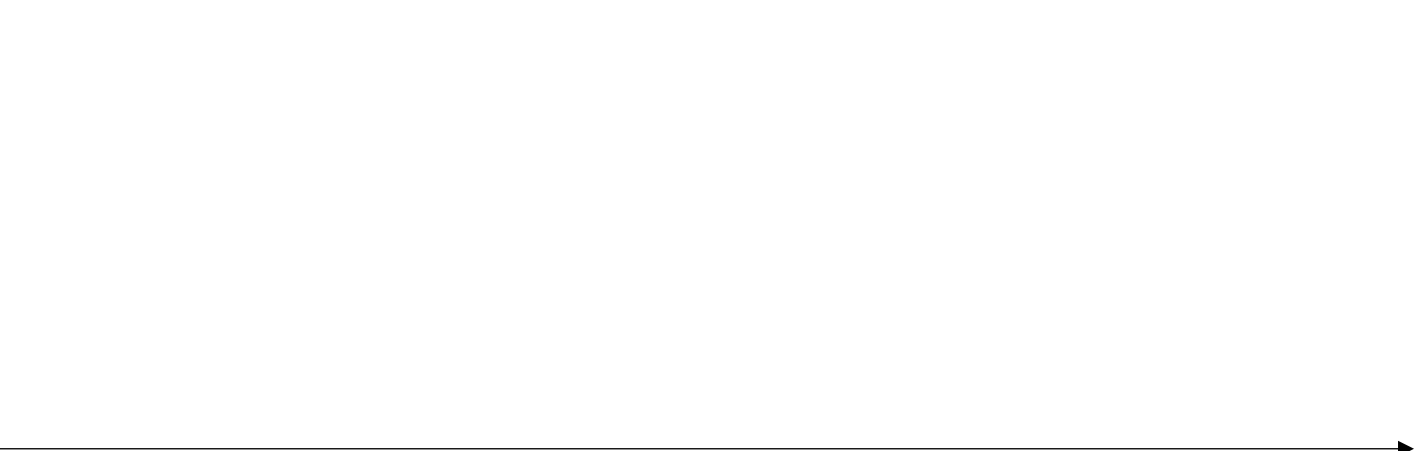
store buffer

addr value

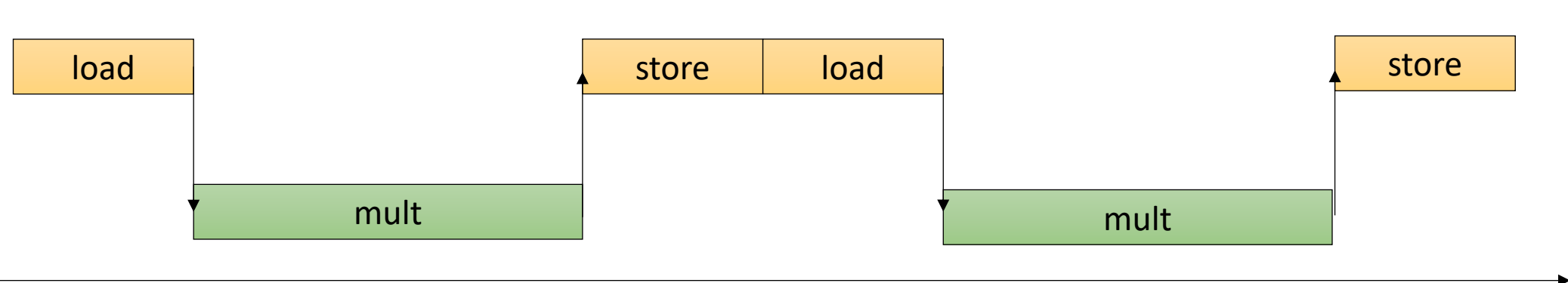
|  |  |
|--|--|
|  |  |
|  |  |
|  |  |

```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```

Access  
Store Buffer  
Execute



Access  
Execute



*time*



# DAE Parallelism

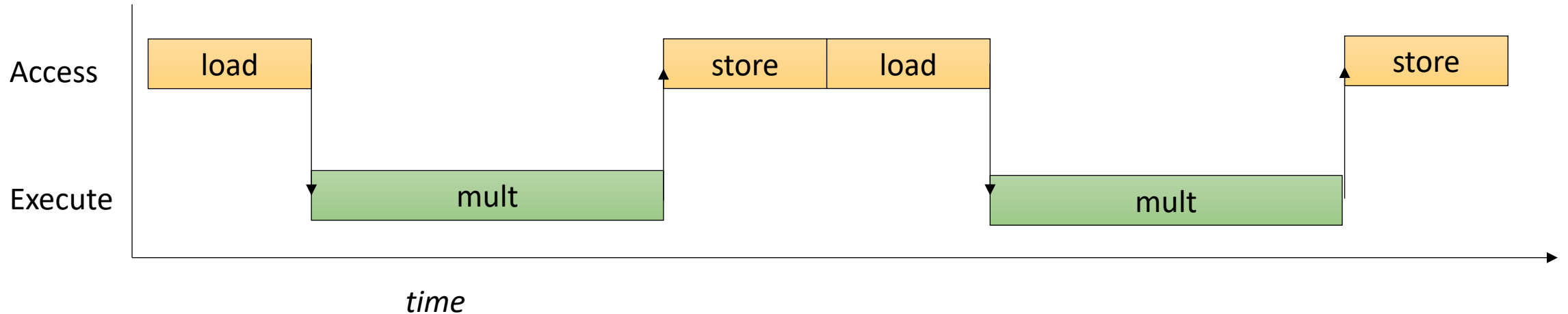
store buffer

addr value

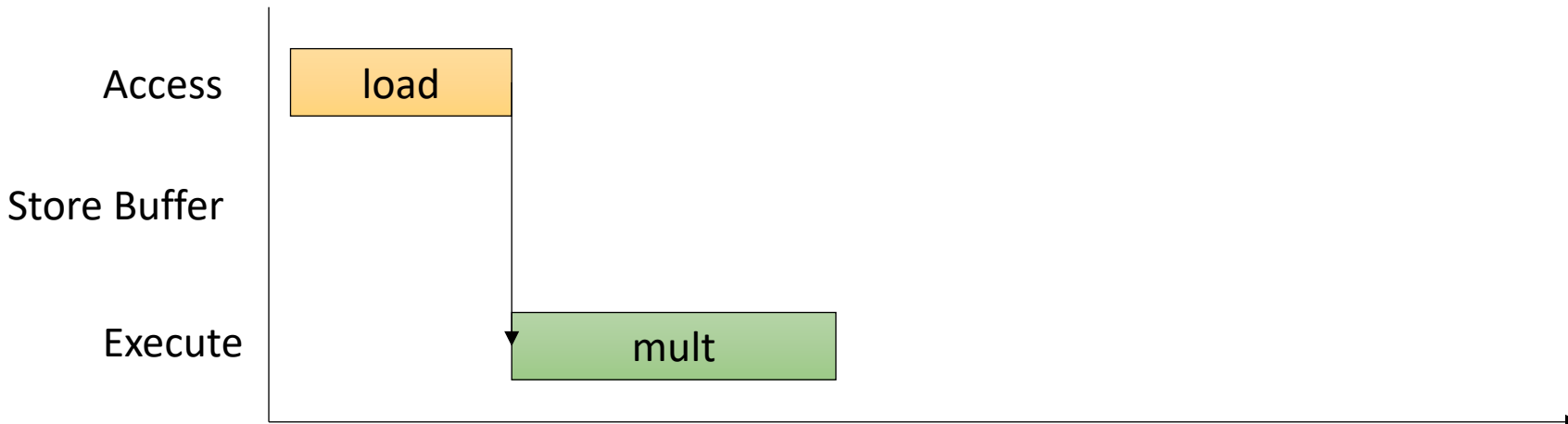
|  |  |
|--|--|
|  |  |
|  |  |
|  |  |



```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



# DAE Parallelism

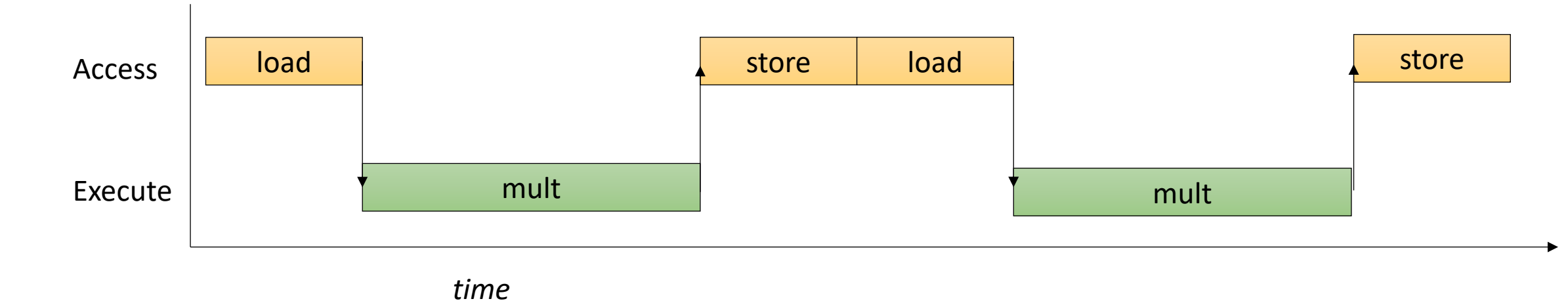


store buffer

addr value

| addr | value |
|------|-------|
|      |       |
|      |       |
|      |       |

```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



# DAE Parallelism

store buffer

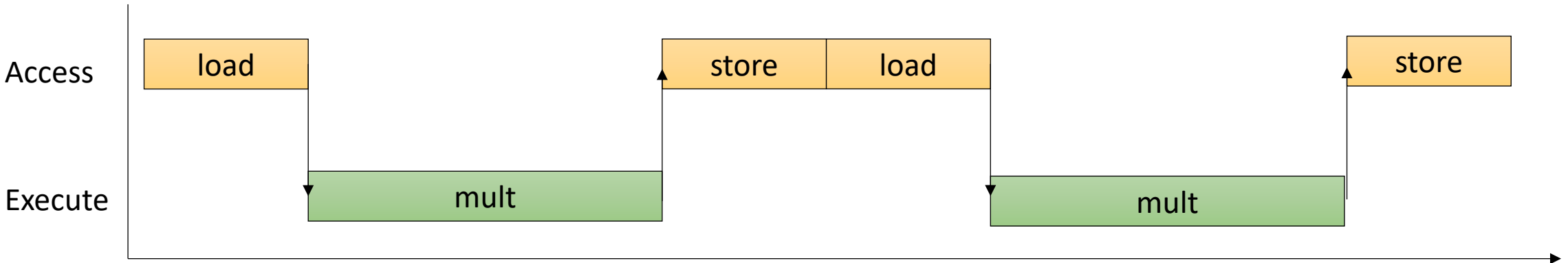
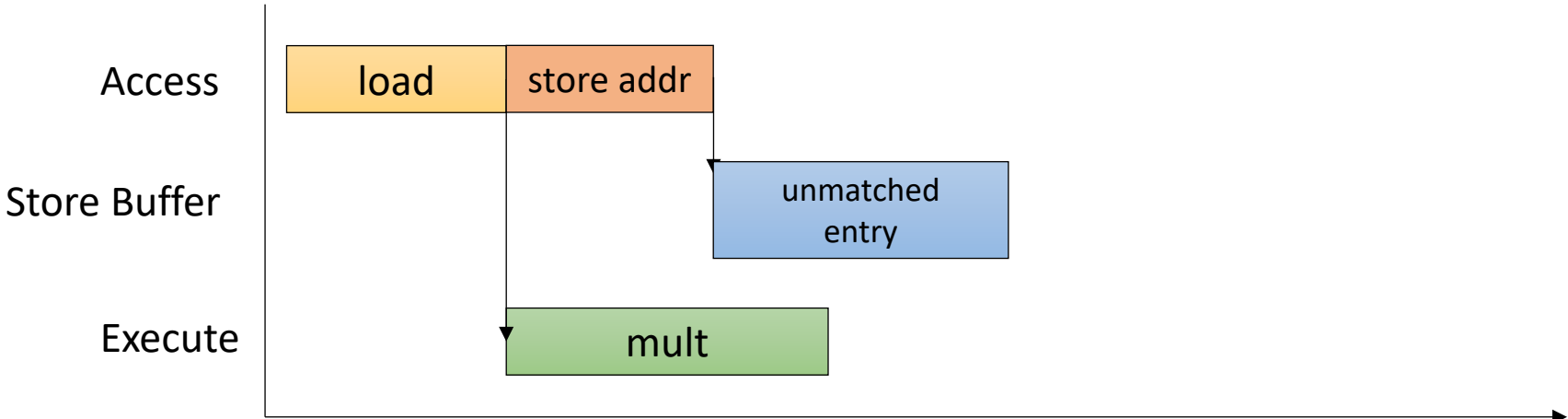
addr value

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |

```

for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}

```



*time*

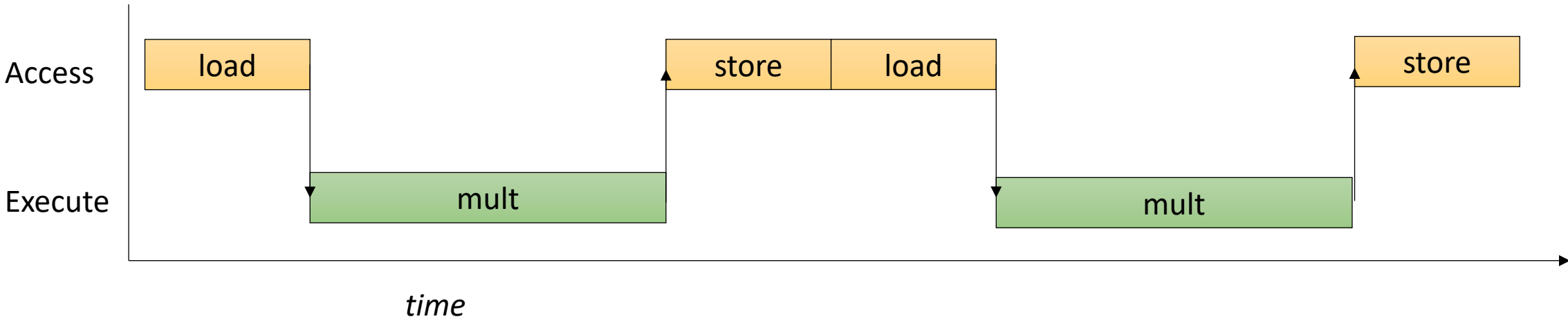
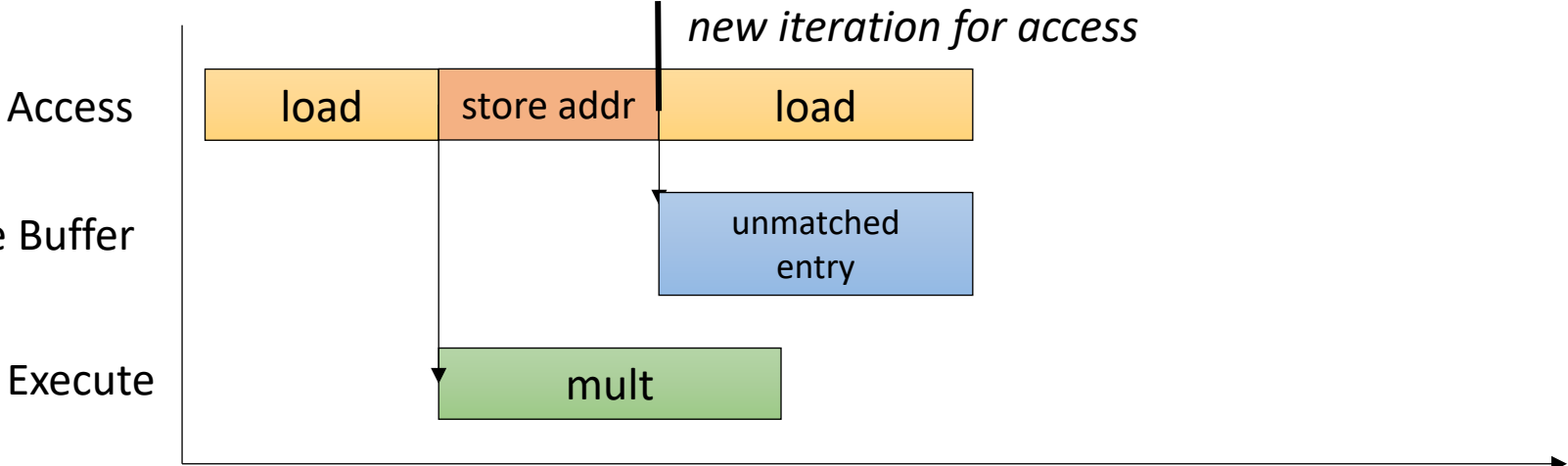
# DAE Parallelism

store buffer

addr value

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |

```
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```



time

# DAE Parallelism

store buffer

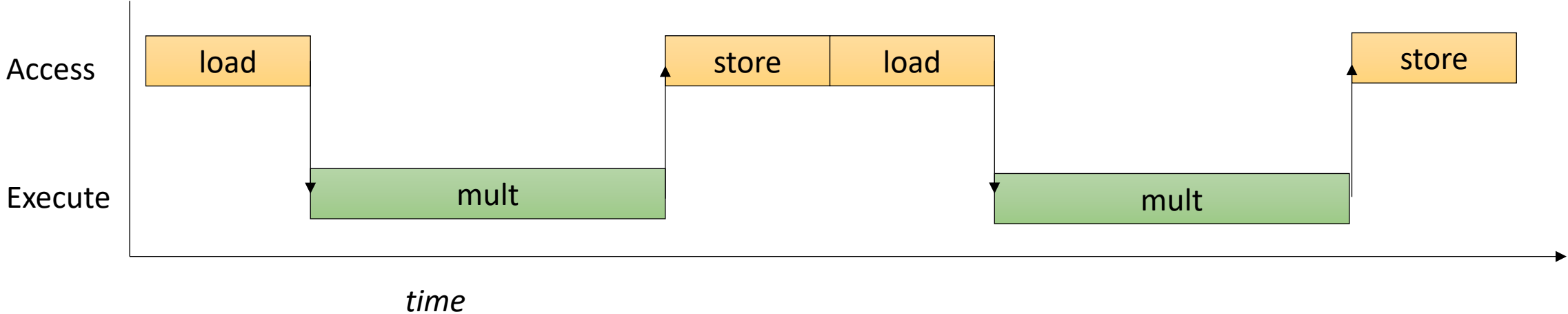
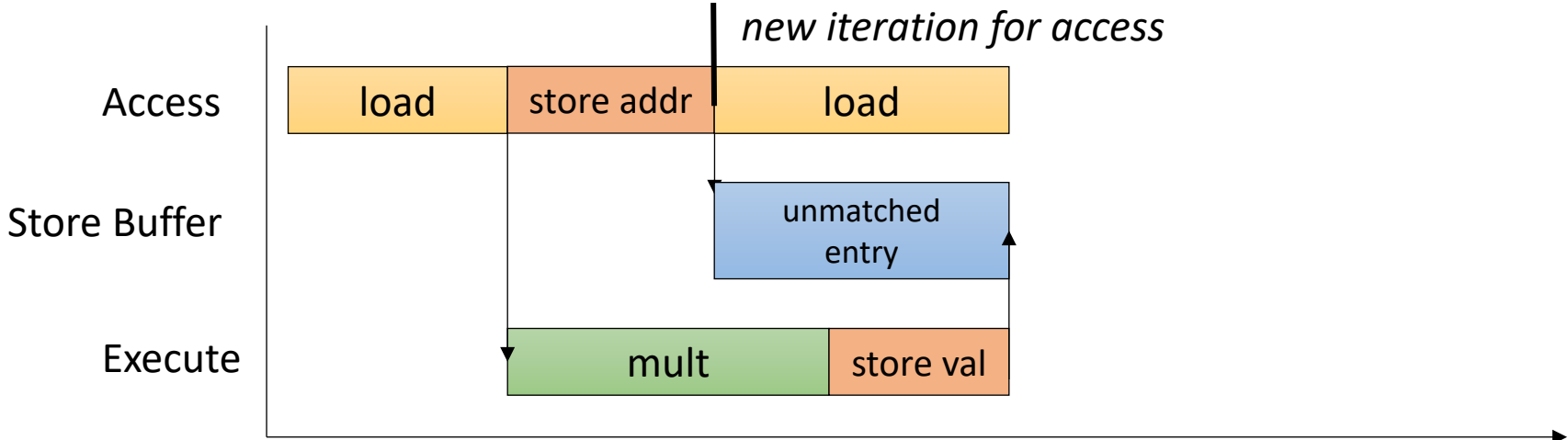
addr value

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |

*can flush now!*

```

for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
    
```



# DAE Parallelism

store buffer

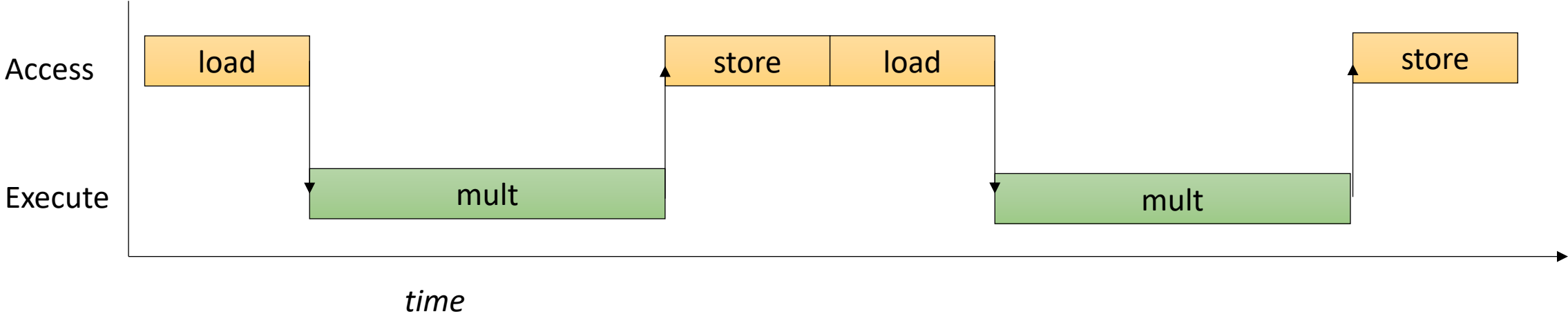
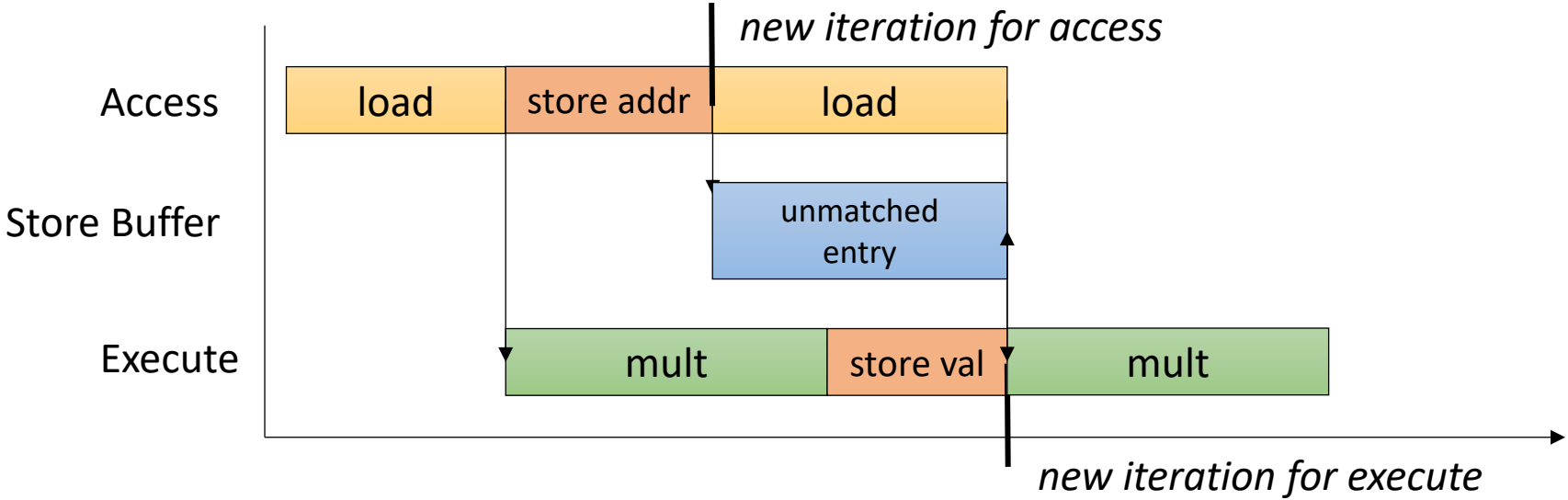
addr value

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |

*can flush now!*

```

for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
    
```



# DAE Parallelism

store buffer

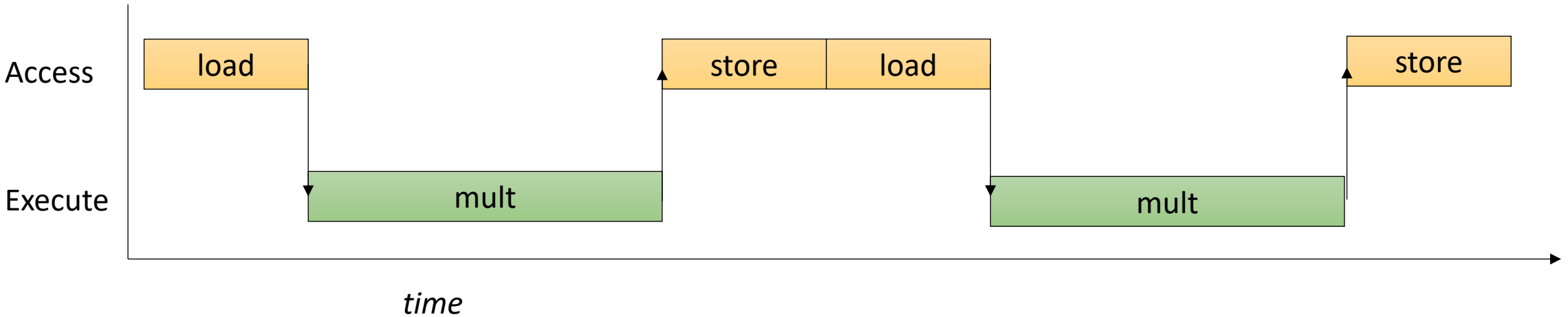
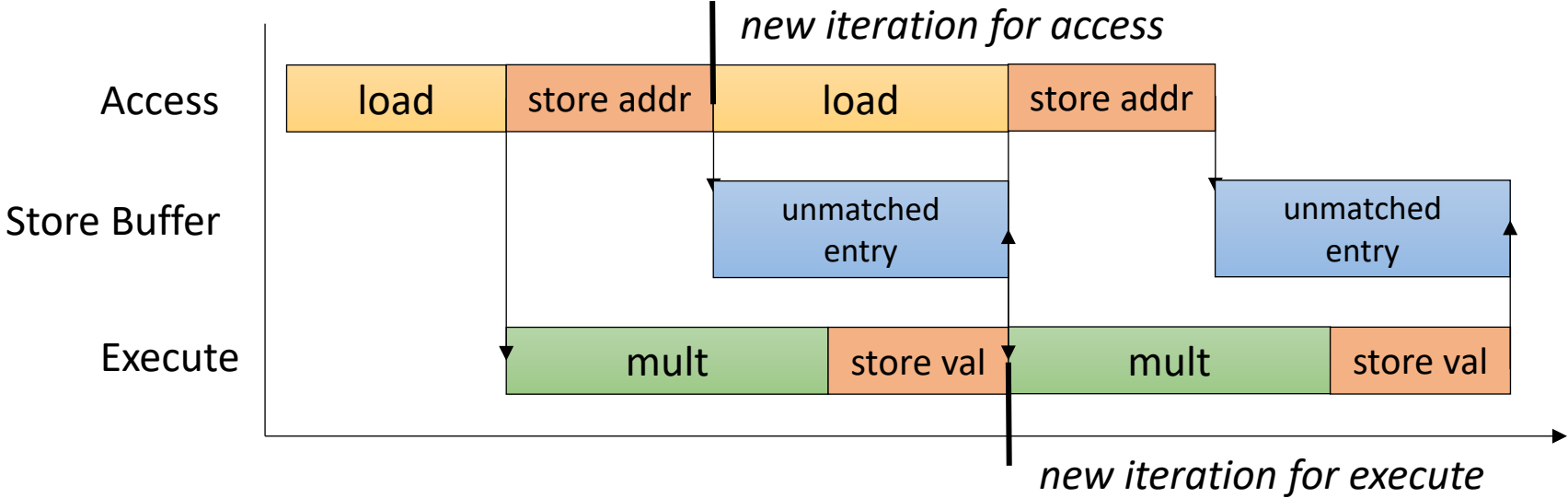
addr value

|  |  |
|--|--|
|  |  |
|  |  |
|  |  |

*has 2 entries now*

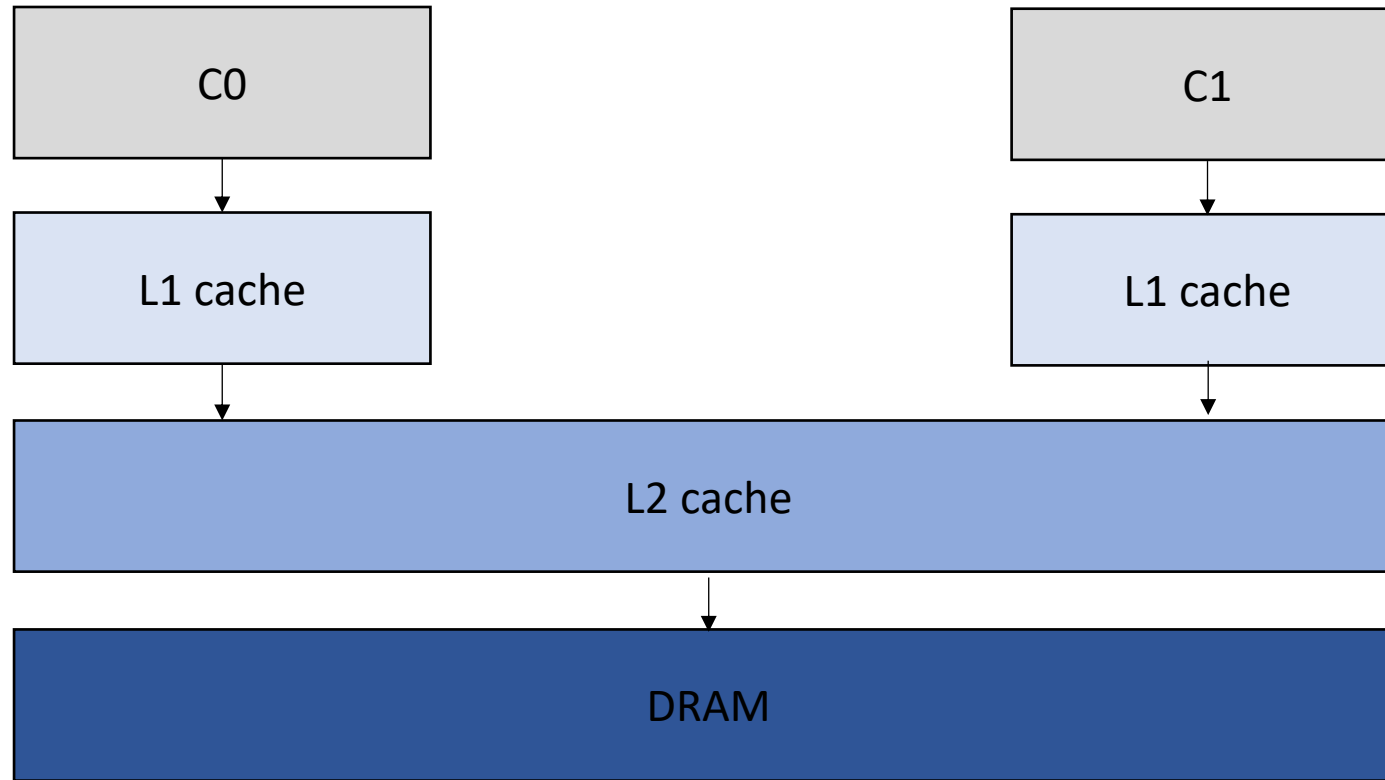
```

for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
    
```



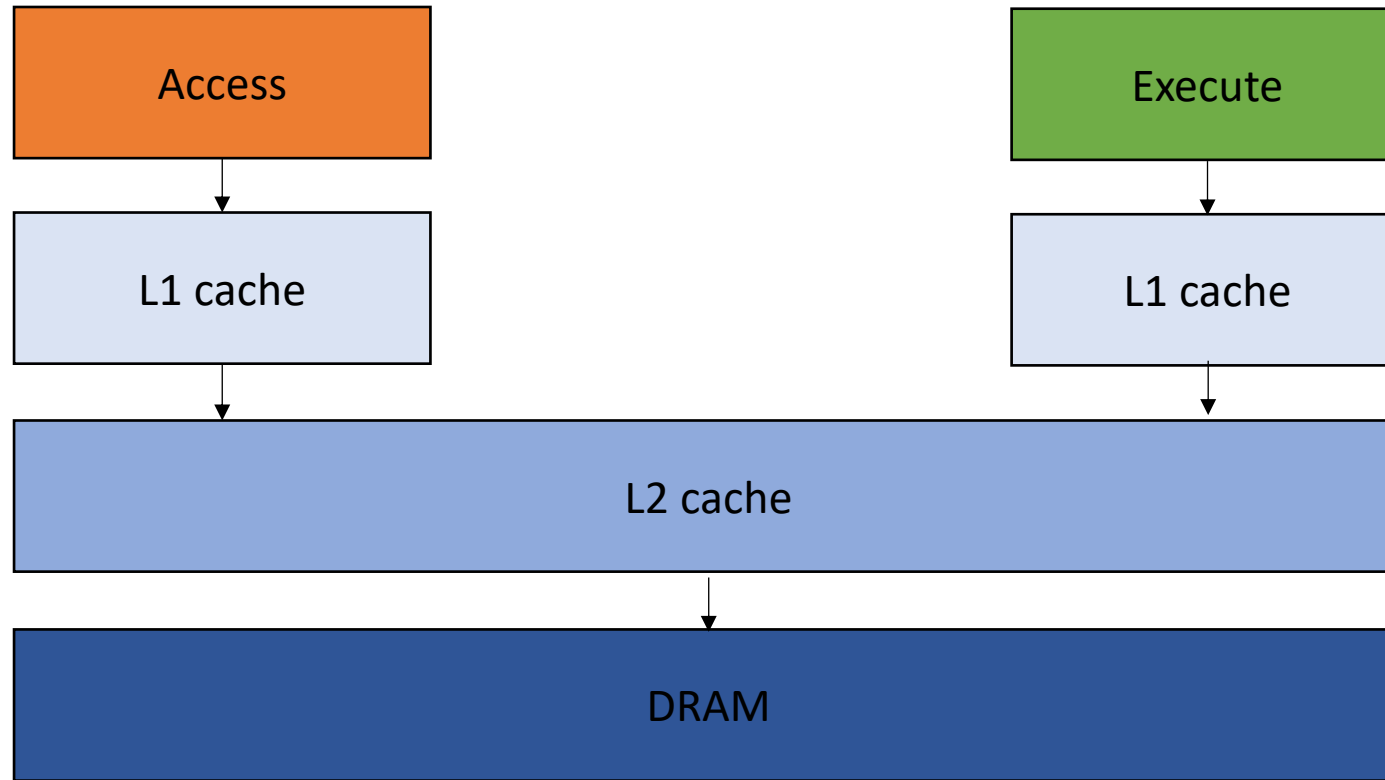
*time*

# Specializing a DAE architecture



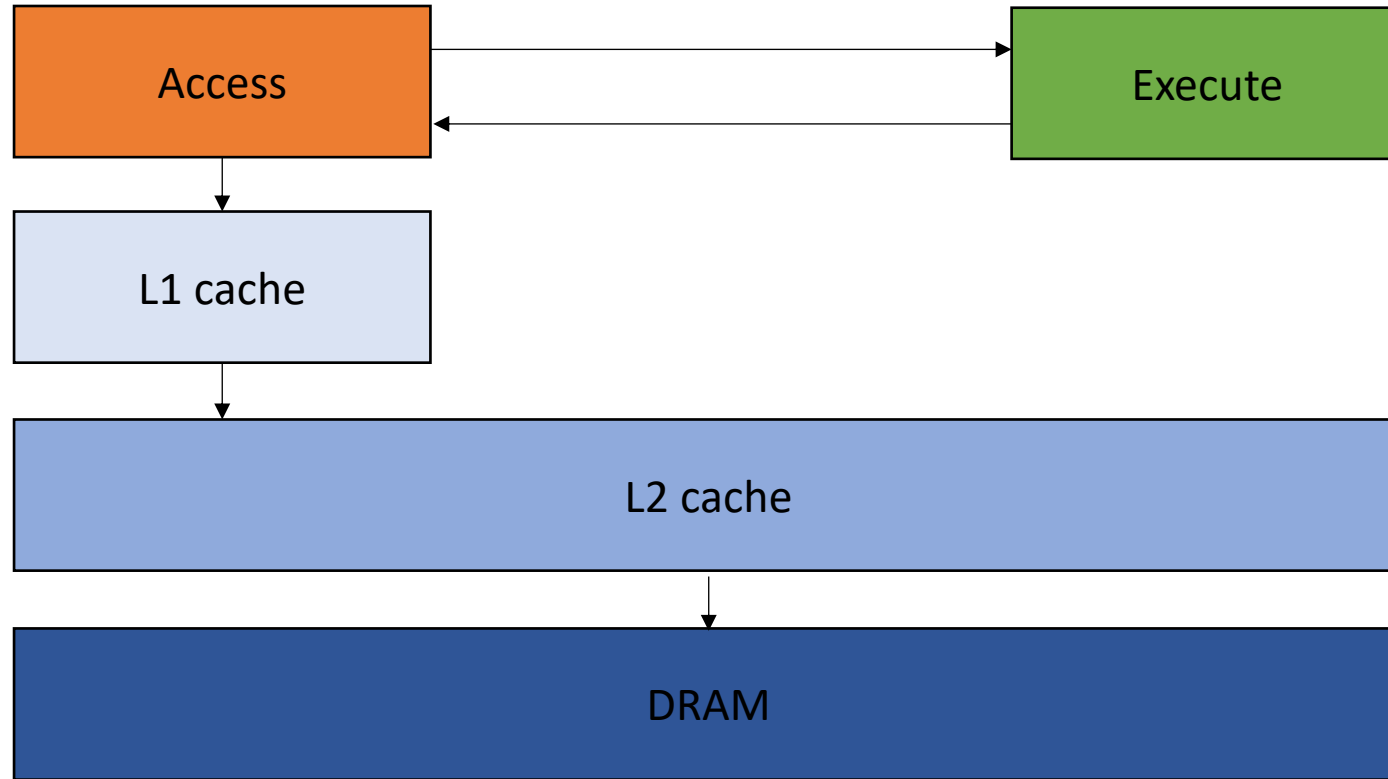


# Specializing a DAE architecture



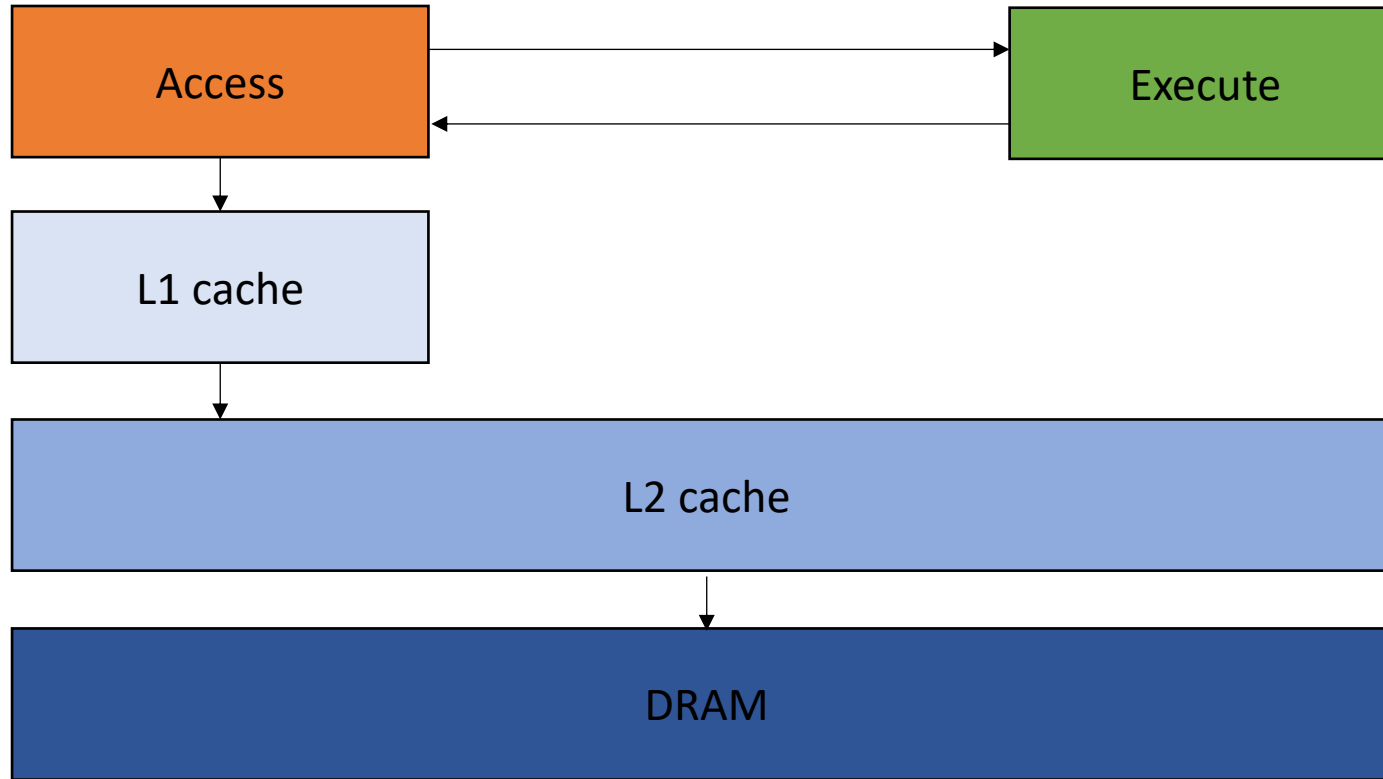
# Specializing a DAE architecture

*Less contention  
on memory hierarchy*



# Specializing a DAE architecture

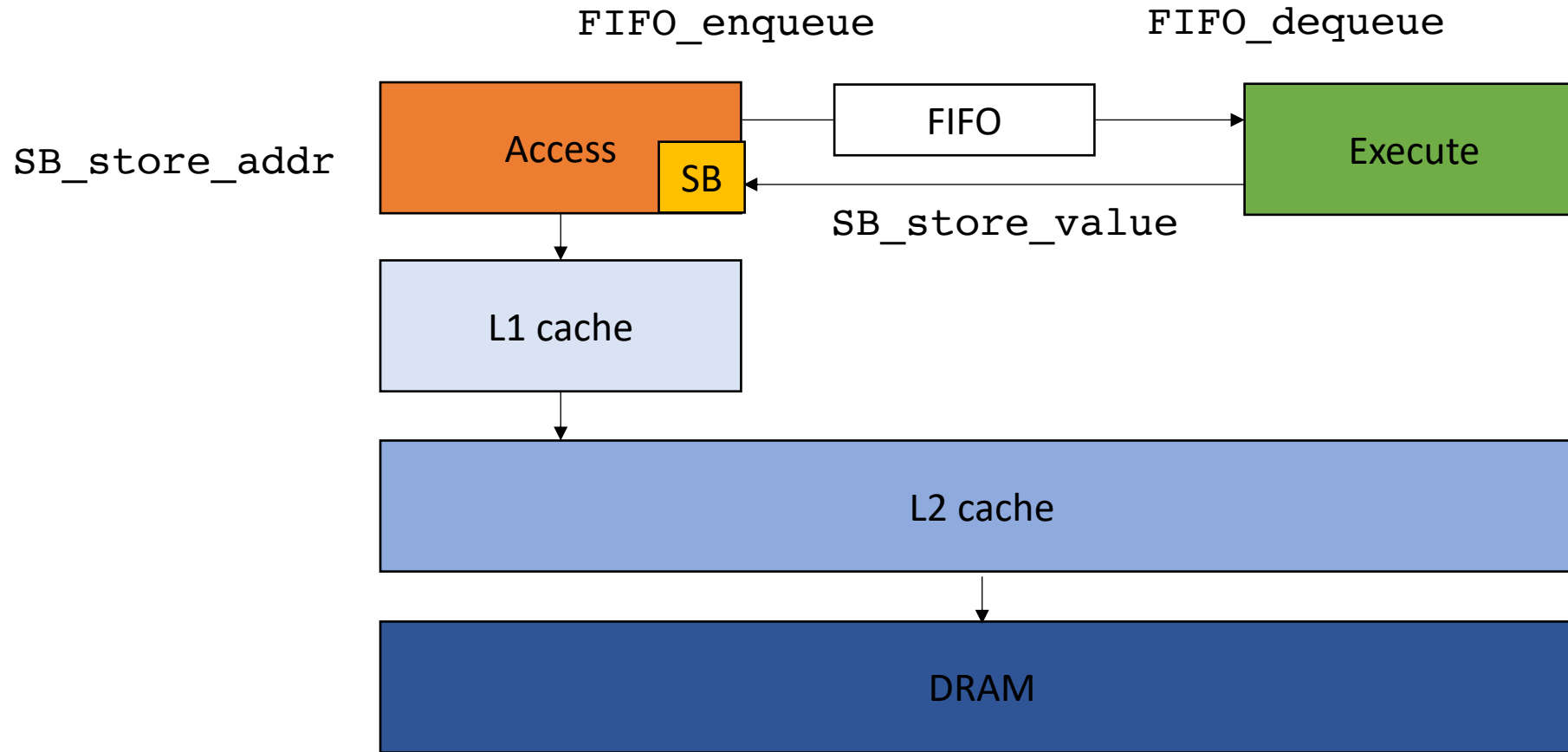
optimizations?  
FP unit  
Vector units



optimizations?  
Load/Store unit  
Storebuffer

*Less contention  
on memory hierarchy*

# DAE API



# Compiler

- Given a sequential program, how can we automatically target a DAE architecture?

# Program slicing

- Mark Weiser in 1981.
  - presented as a formalization of debugging



# Program slicing

Main idea:

# Program slicing

Main idea:

- **Forward Slicing:** given statements  $S$ , remove all statements except for those that depend (control or data) on  $s \in S$ 
  - Intuitively: get a minimal (heuristically) program where all actions depend on statements in  $S$



# Program slicing

Main idea:

- **Forward Slicing:** given statements  $S$ , remove all statements except for those that depend (control or data) on  $s \in S$ 
  - Intuitively: get a minimal (heuristically) program where all actions depend on statements in  $S$
- **Backward Slicing:** given statements  $S$ , remove all statements except for those that for which  $s \in S$  depend on.
  - Intuitively: get a minimal (heuristically) program where all actions influence statements in  $S$

# Program slicing

Main idea:

- **Forward Slicing:** given statements  $S$ , remove all statements except for those that depend (control or data) on  $s \in S$ 
  - Intuitively: get a minimal (heuristically) program where all actions depend on statements in  $S$

This is the one we will focus on

- **Backward Slicing:** given statements  $S$ , remove all statements except for those that for which  $s \in S$  depend on.
  - Intuitively: get a minimal (heuristically) program where all actions influence statements in  $S$

# Program slicing

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. r5 = r2 + r0;  
6. r6 = 128;  
7. assert(r5 == r6)
```

slicing criterion: [  
"7. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. r5 = r2 + r0;  
6. r6 = 128;  
7. assert(r5 == r6)
```

slicing criterion: [  
"7. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. r5 = r2 + r0;  
6. r6 = 128;  
7. assert(r5 == r6)
```

slicing criterion: [  
"7. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. r5 = r2 + r0;  
6. r6 = 128;  
7. assert(r5 == r6)
```

slicing criterion: [  
"7. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. r5 = r2 + r0;  
6. r6 = 128;  
7. assert(r5 == r6)
```

slicing criterion: [  
"7. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. r5 = r2 + r0;  
6. r6 = 128;  
7. assert(r5 == r6)
```

slicing criterion: [  
"7. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*



# Program slicing - Control dependence

```
1.  r0 = a + b;  
2.  r1 = b + c;  
3.  r2 = r0 * r0;  
4.  r4 = r1 + r0;  
5.  bne r4, 64, END  
6.  r5 = r2 + r0;  
7.  r6 = 128;  
8.  assert(r5 == r6)  
9.END:
```

slicing criterion: [  
"8. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing - Control dependence

```
1.  r0 = a + b;  
2.  r1 = b + c;  
3.  r2 = r0 * r0;  
4.  r4 = r1 + r0;  
5.  bne r4, 64, END
```

```
6.  r5 = r2 + r0;  
7.  r6 = 128;  
8.  assert(r5 == r6)
```

```
9. END:
```

slicing criterion: [

```
"8. assert(r5 == r6)"
```

]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing - Control dependence

```
1.  r0 = a + b;  
2.  r1 = b + c;  
3.  r2 = r0 * r0;  
4.  r4 = r1 + r0;  
5.  bne r4, 64, END
```

```
6.  r5 = r2 + r0;  
7.  r6 = 128;  
8.  assert(r5 == r6)
```

```
9. END:
```

slicing criterion: [

"8. assert(r5 == r6)"

]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing - Control dependence

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. bne r4, 64, END
```

*branch statement*

```
6. r5 = r2 + r0;  
7. r6 = 128;  
8. assert(r5 == r6)
```

```
9.END:
```

slicing criterion: [  
"8. assert(r5 == r6)"  
]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing - Control dependence

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. bne r4, 64, END
```

```
6. r5 = r2 + r0;  
7. r6 = 128;  
8. assert(r5 == r6)
```

```
9.END:
```

slicing criterion: [

"8. assert(r5 == r6)"

]

*start with the statement and work backwards until there are no more dependencies*

# Program slicing - Control dependence

```
1. r0 = a + b;  
2. r1 = b + c;  
3. r2 = r0 * r0;  
4. r4 = r1 + r0;  
5. bne r4, 64, END
```

```
6. r5 = r2 + r0;  
7. r6 = 128;  
8. assert(r5 == r6)
```

```
9. END:
```

slicing criterion: [

"8. assert(r5 == r6)"

]

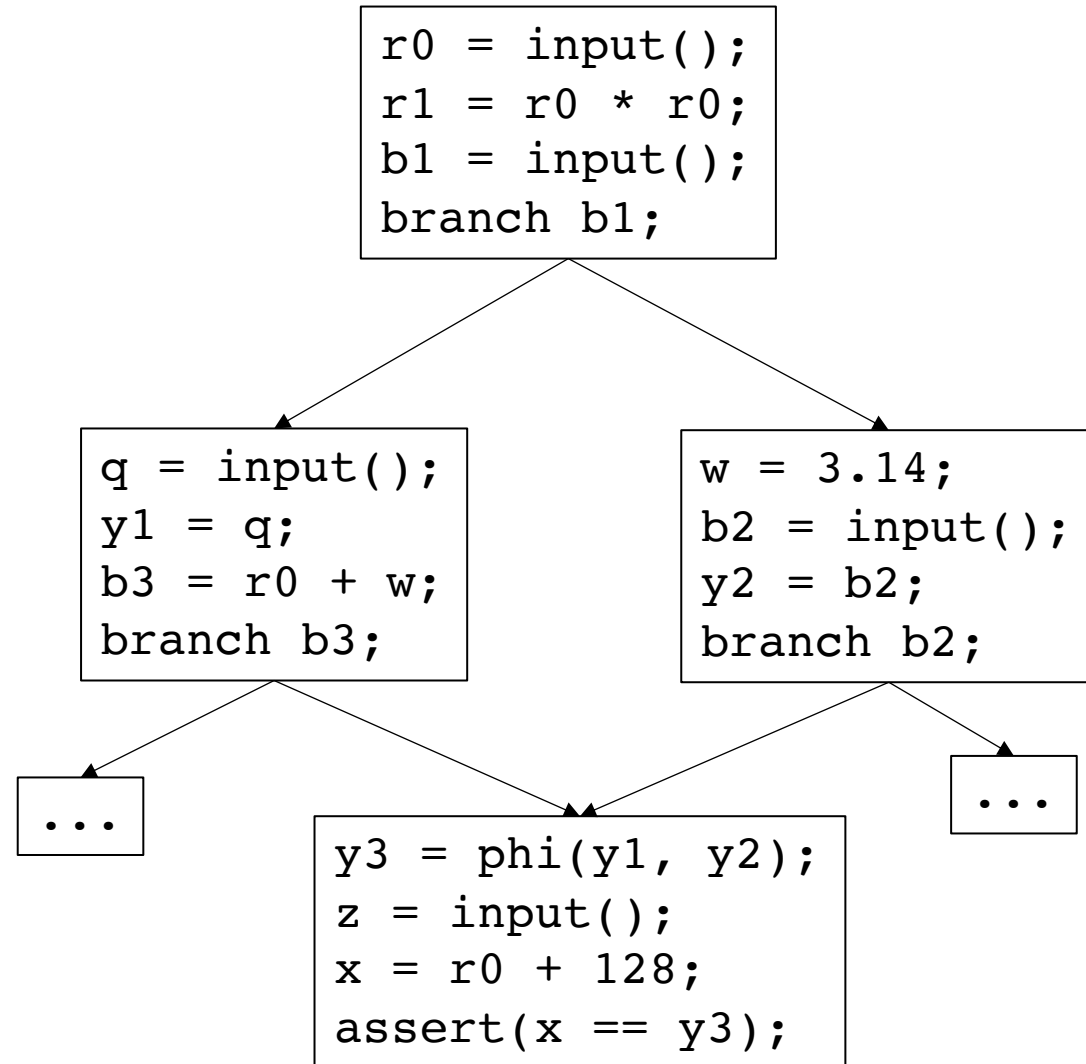
*start with the statement and work backwards until there are no more dependencies*

# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

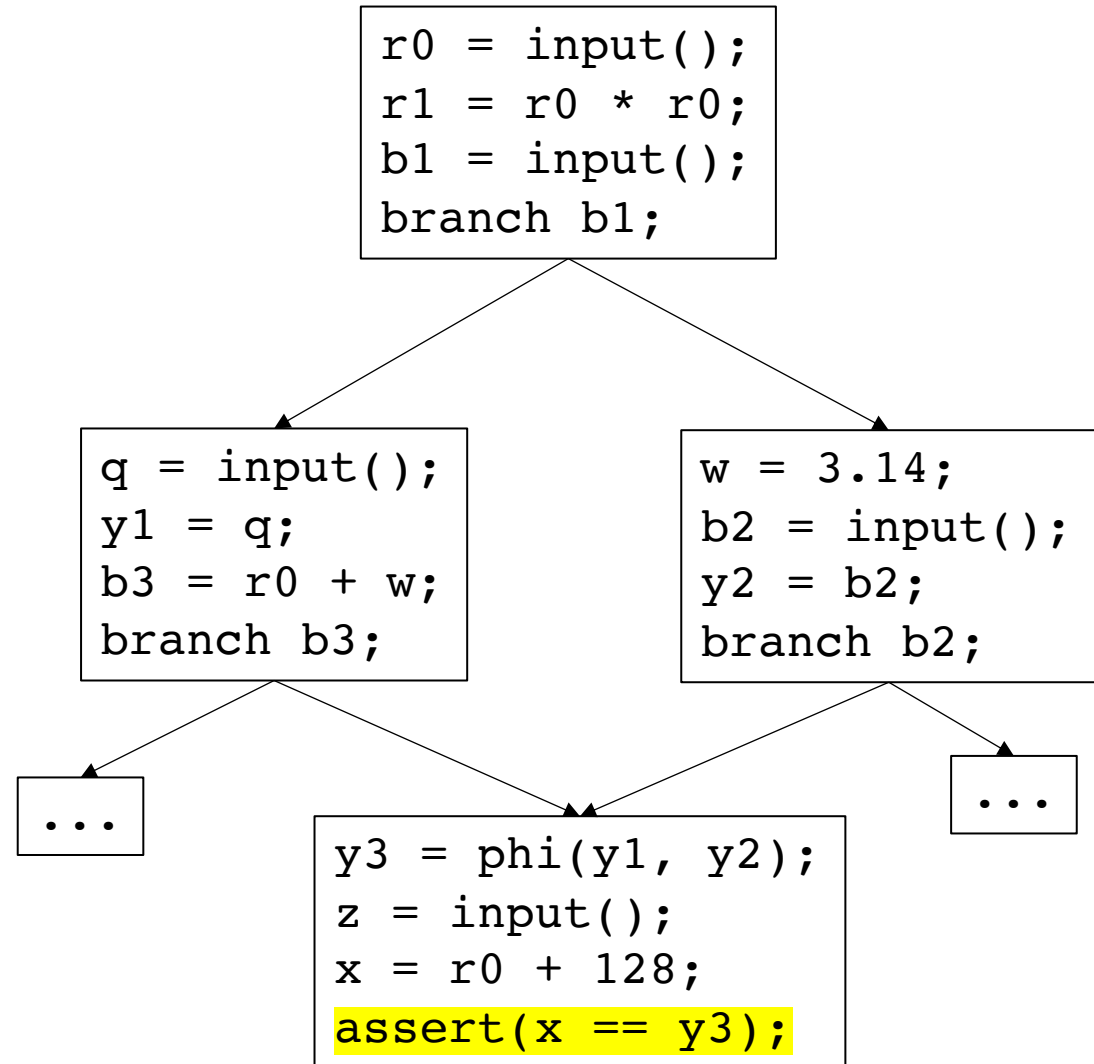




# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

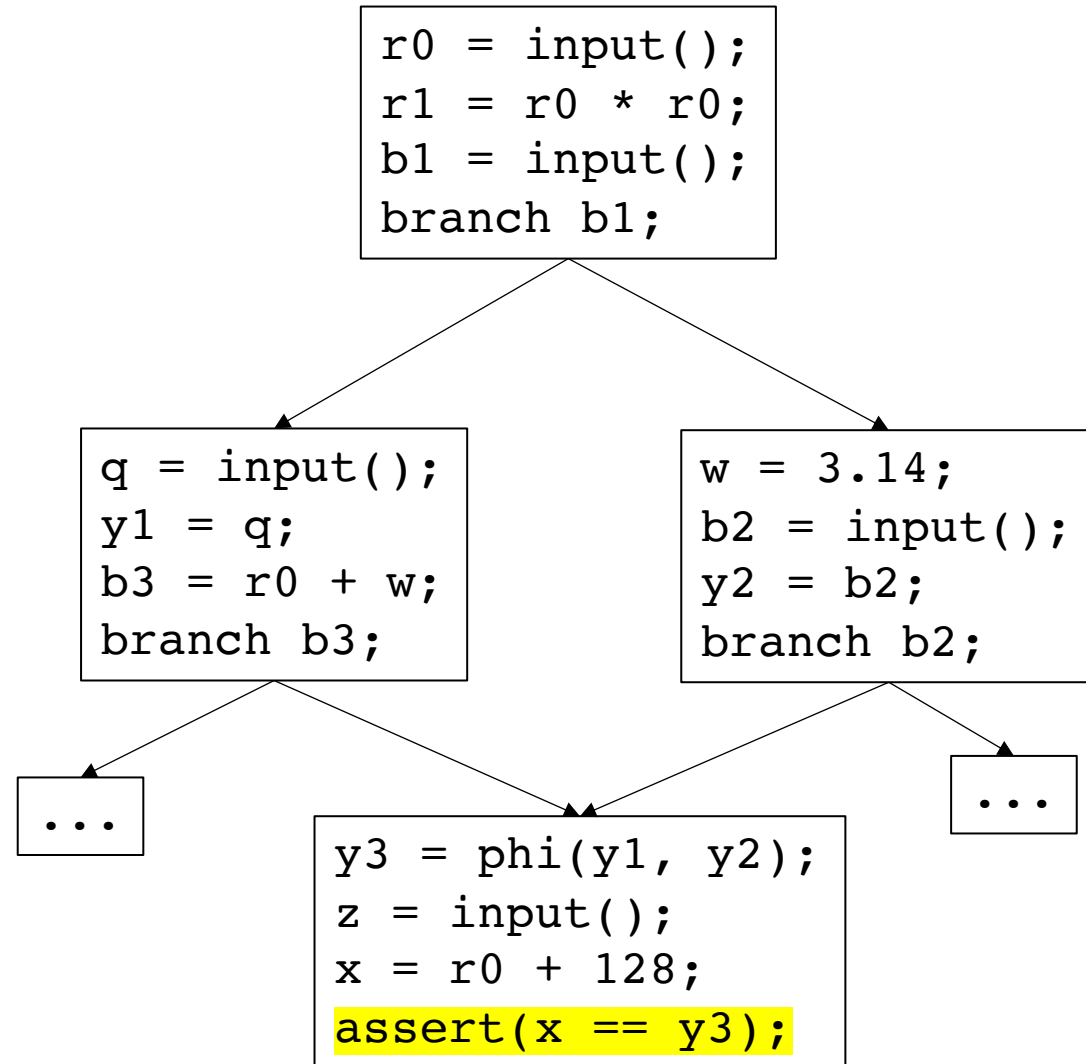
marked:           worklist:  
                  assert()



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

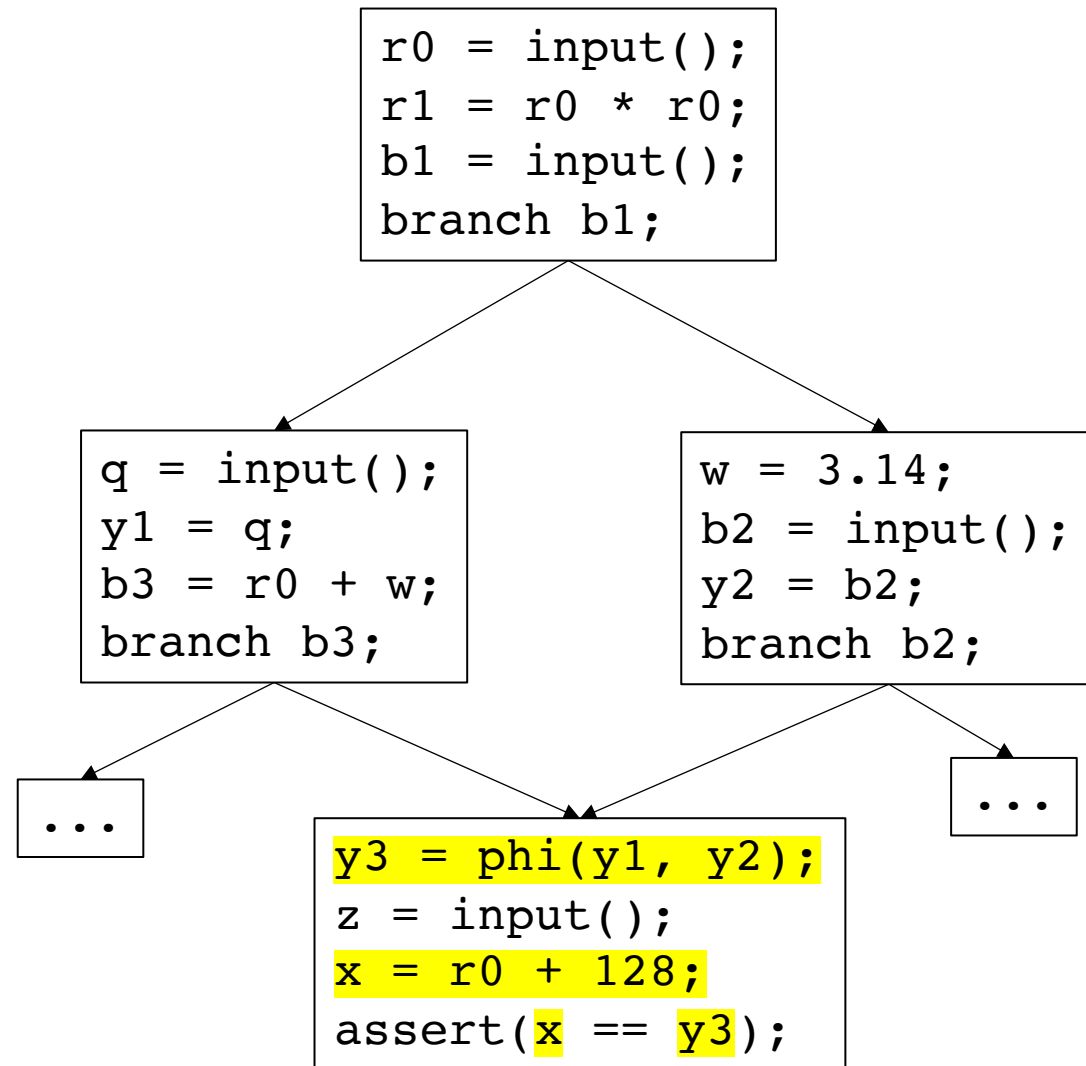
marked:           worklist:  
assert()



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

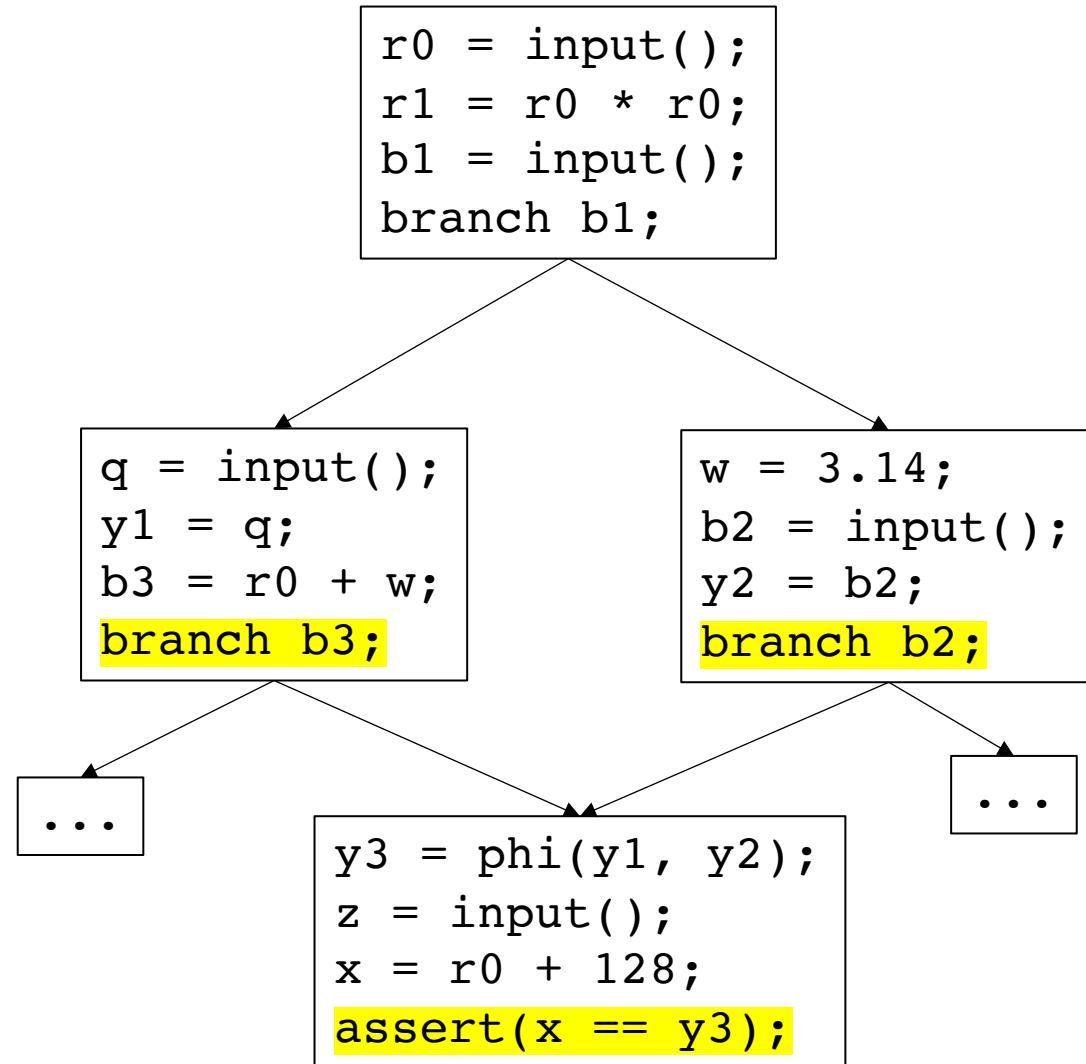
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | x         |
|          | y3        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

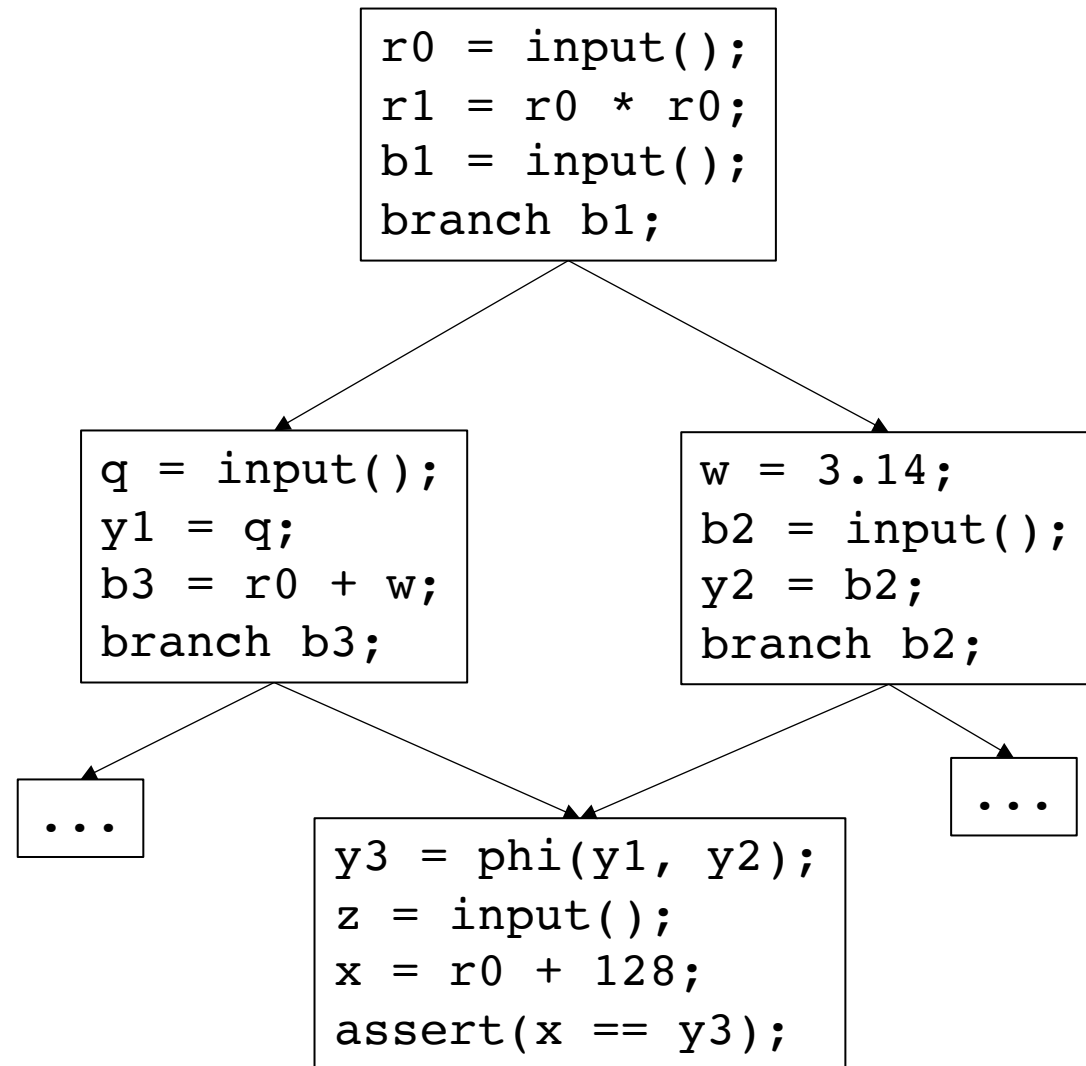
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | x         |
|          | y3        |
|          | branch b3 |
|          | branch b2 |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

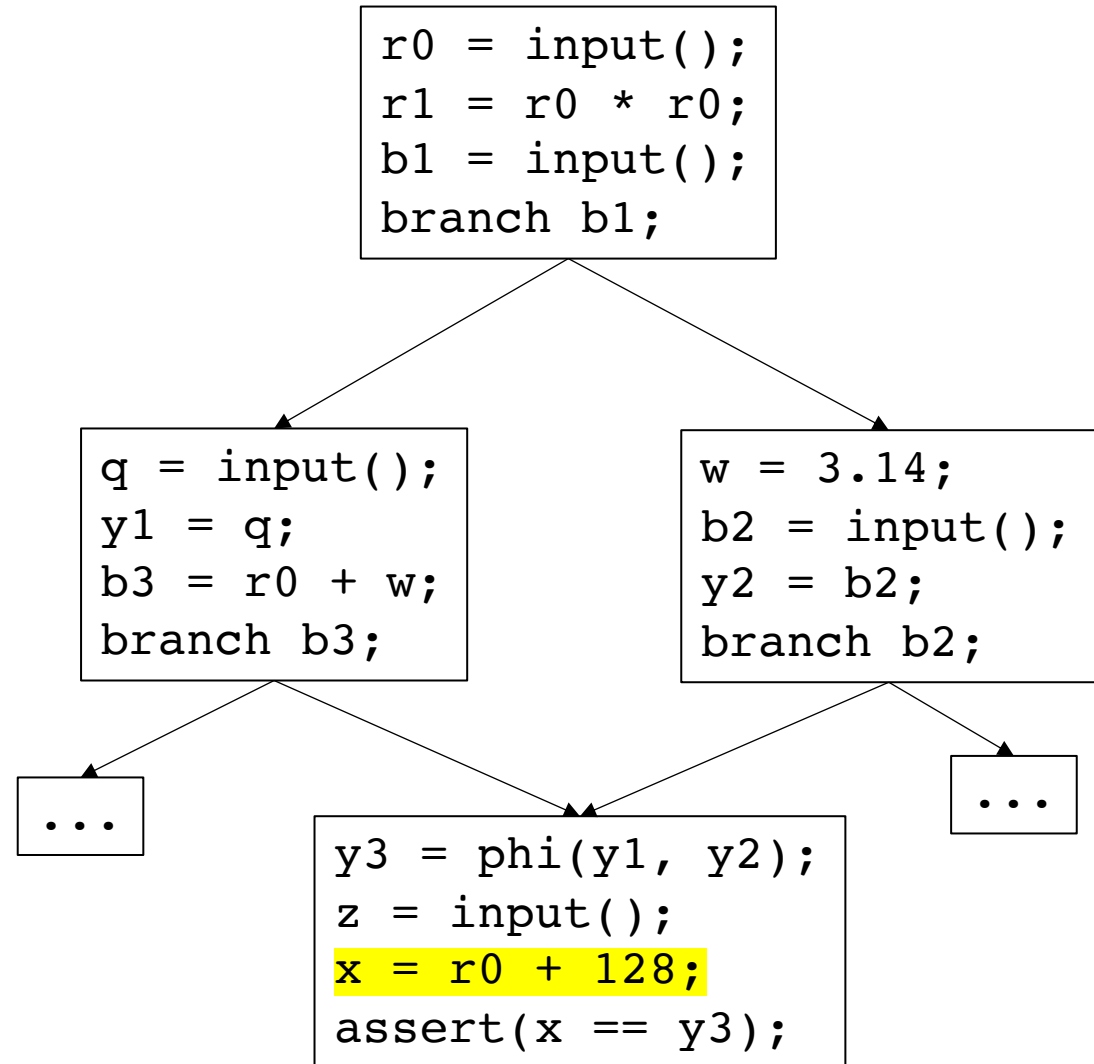
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | x         |
|          | y3        |
|          | branch b3 |
|          | branch b2 |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

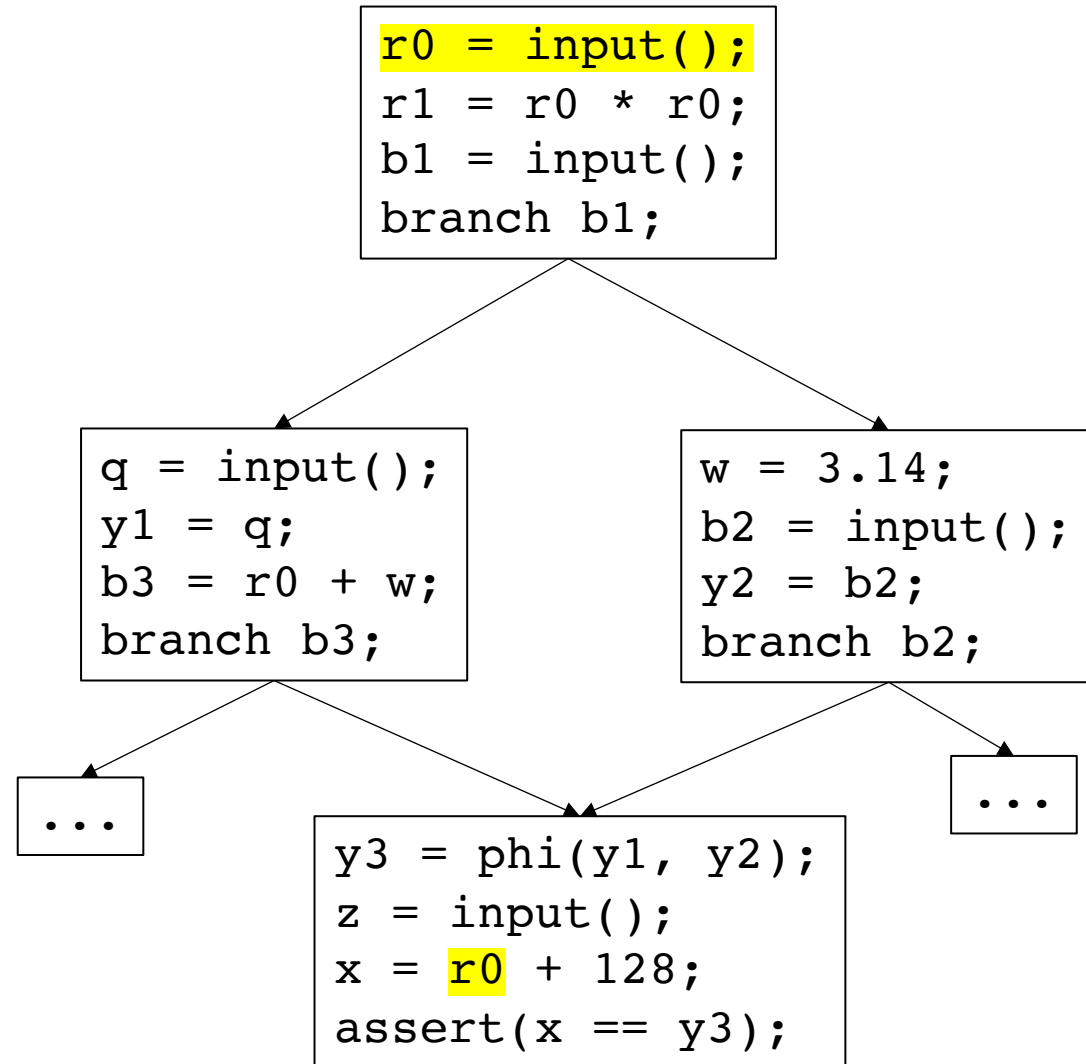
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | y3        |
| x        | branch b3 |
|          | branch b2 |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

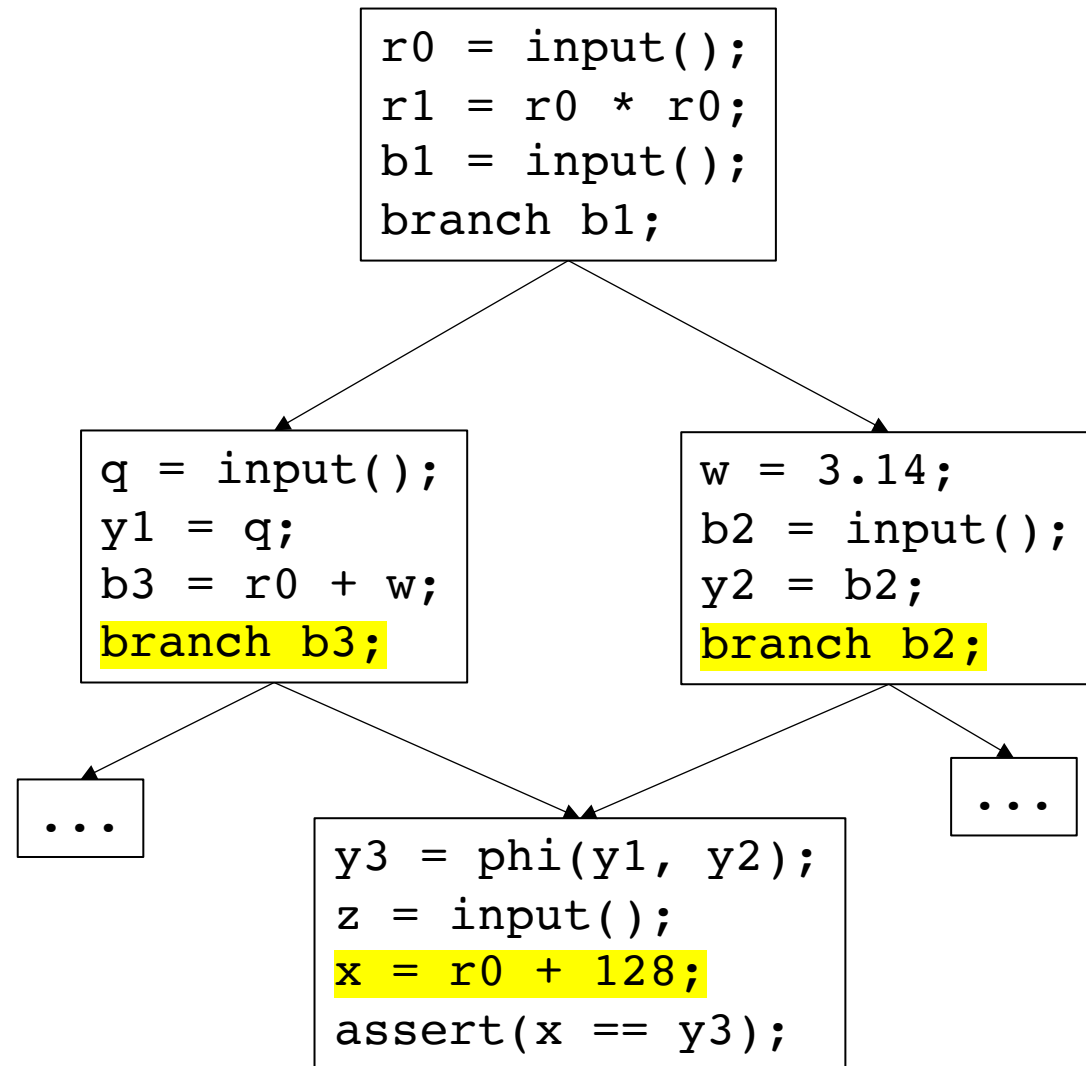
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | y3        |
| x        | branch b3 |
|          | branch b2 |
|          | r0        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | y3        |
| x        | branch b3 |
|          | branch b2 |
|          | r0        |

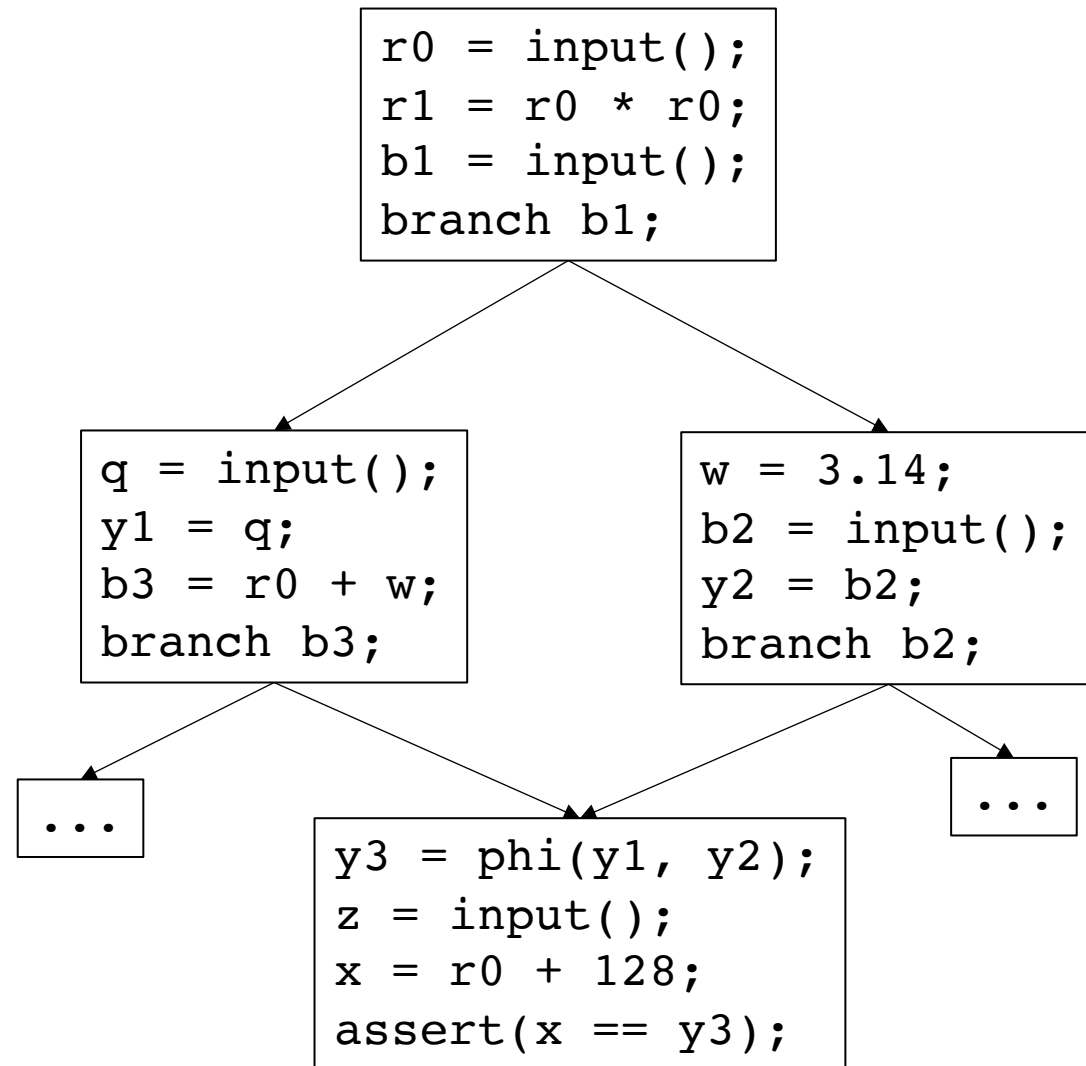




# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

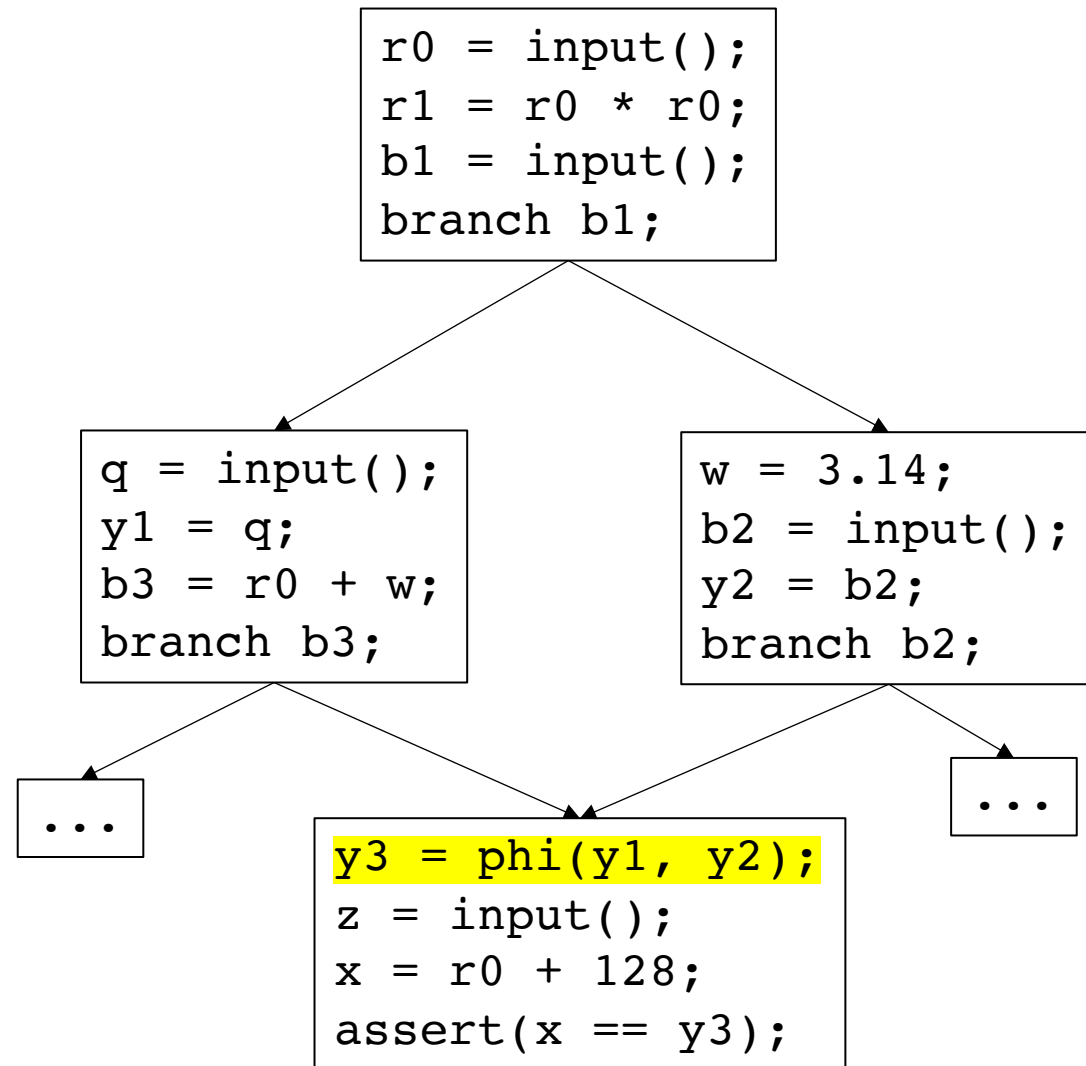
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | y3        |
| x        | branch b3 |
|          | branch b2 |
|          | r0        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

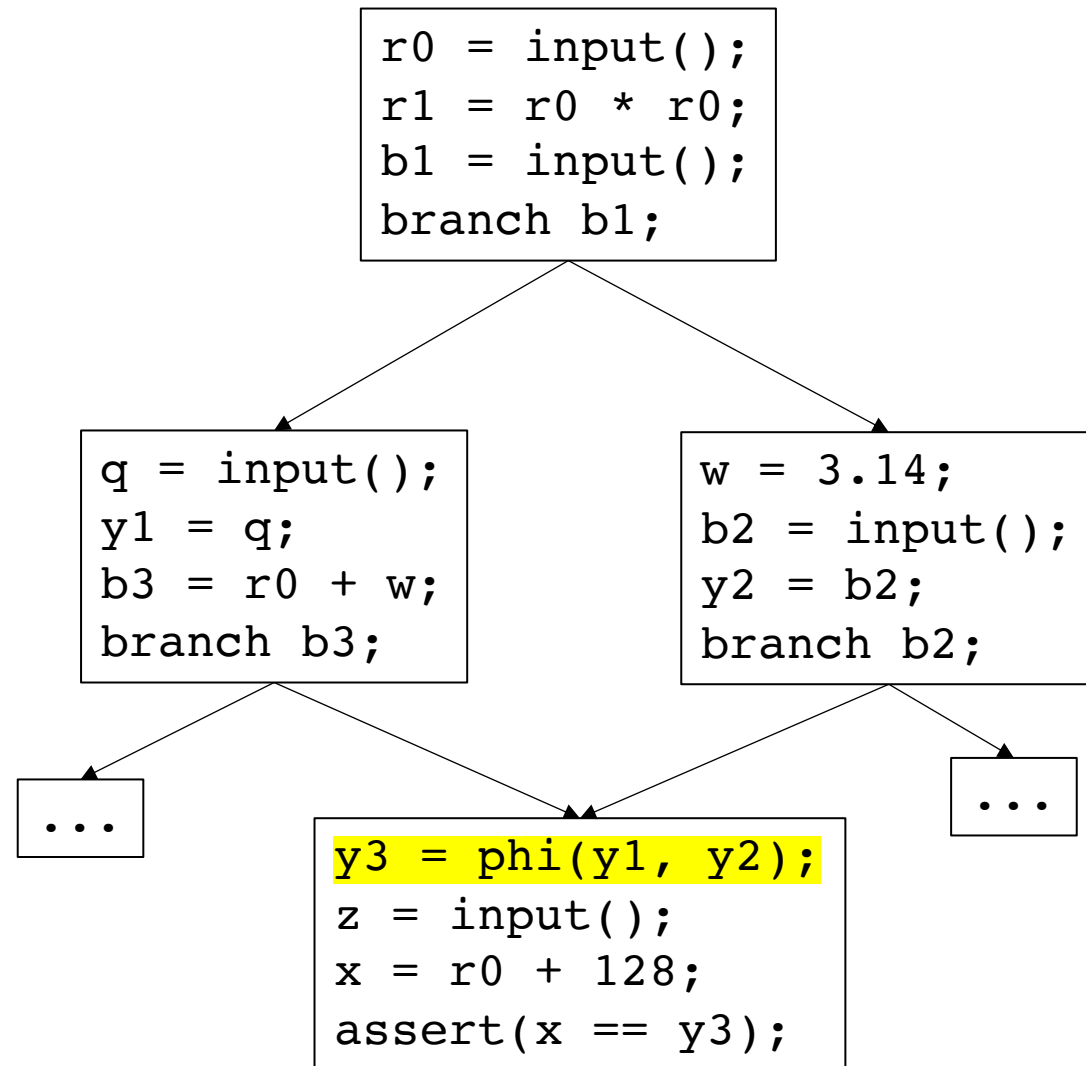
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | y3        |
| x        | branch b3 |
|          | branch b2 |
|          | r0        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
  stmt = Worklist.pop();
  if (is_marked(stmt)) {
    continue;
  }
  mark(stmt);
  for a in stmt.args() {
    worklist.append(a);
  }
  for p in cfg[stmt].predecessors() {
    worklist.append(p.branch_stmt());
  }
}
```

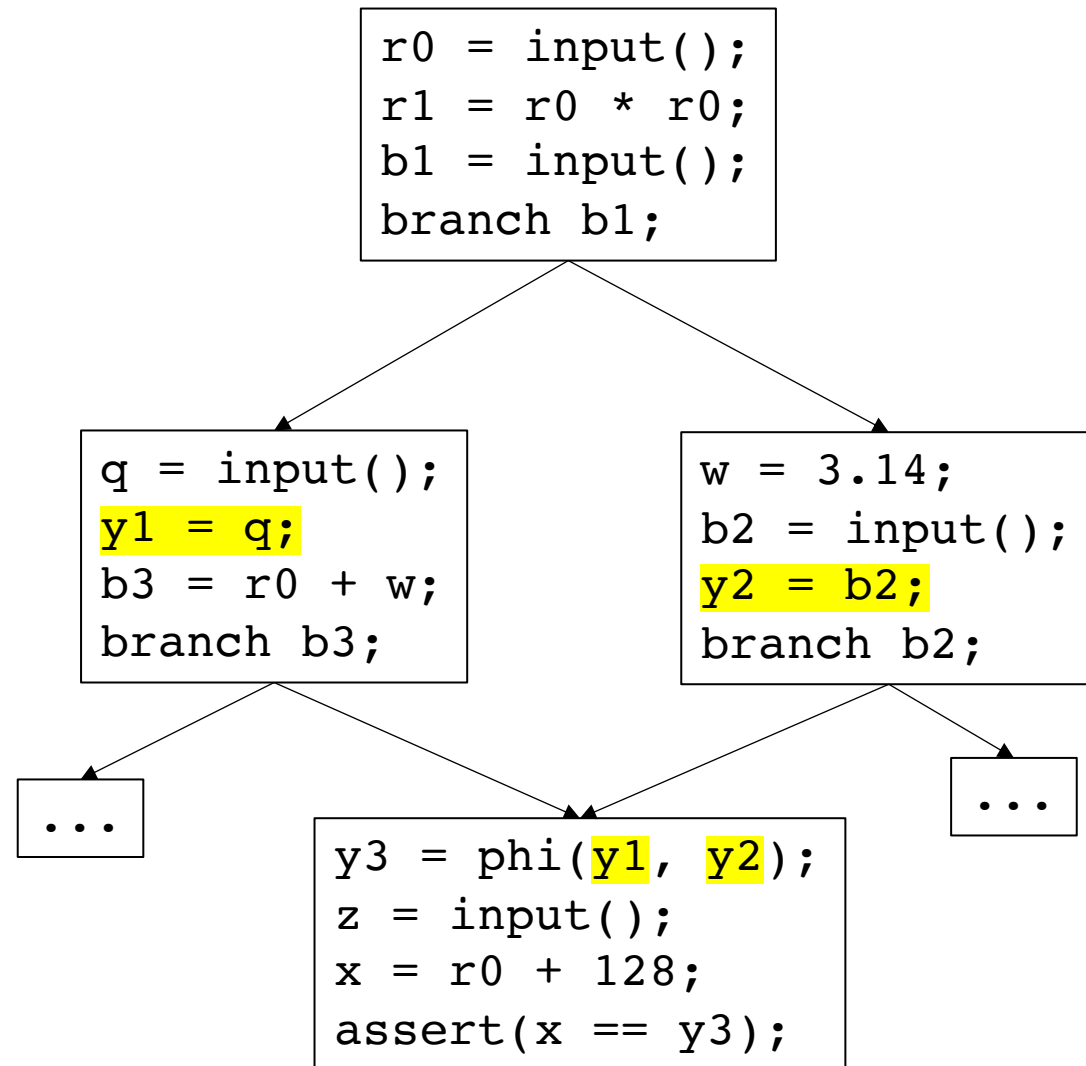
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | branch b3 |
| x        | branch b2 |
| y3       | r0        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

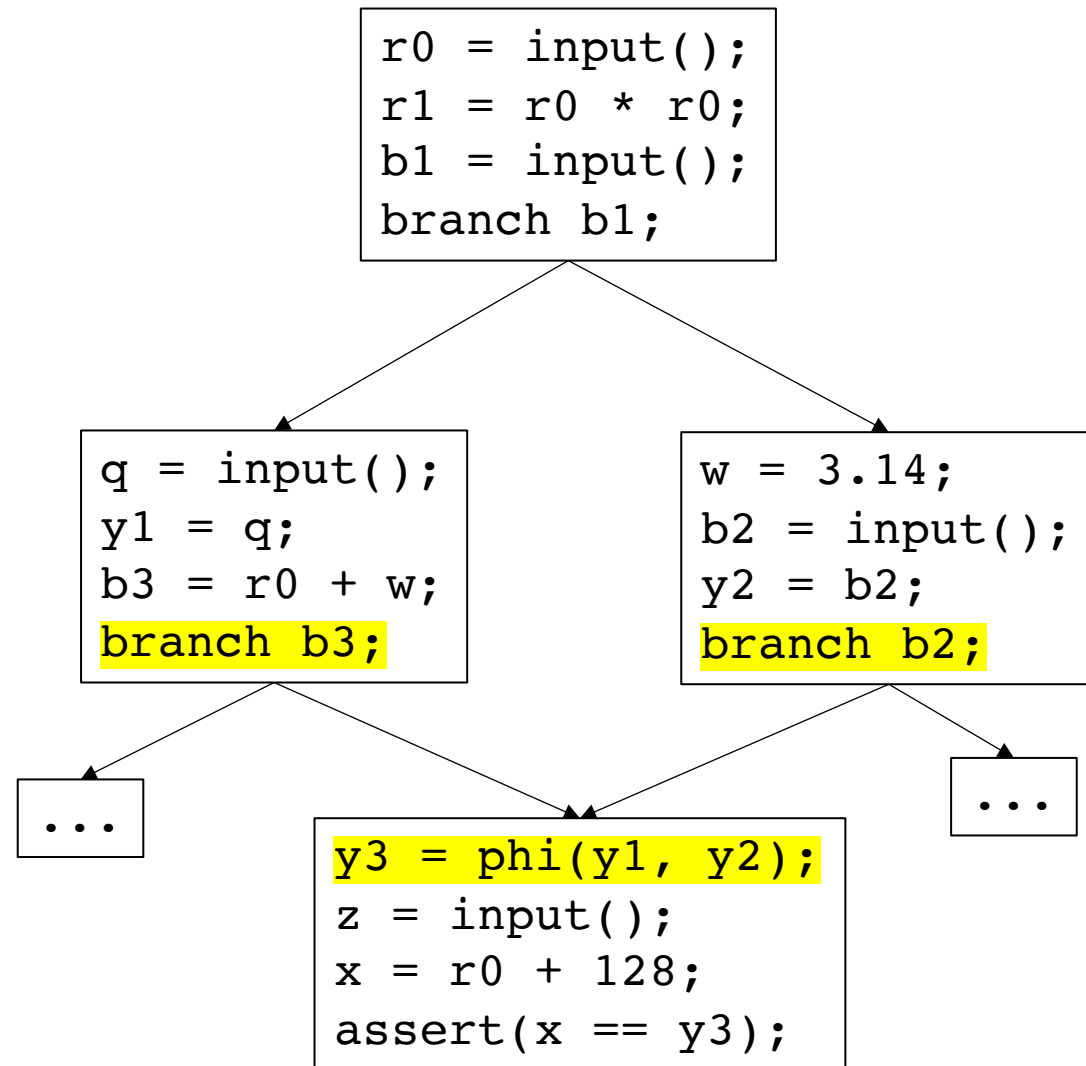
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | branch b3 |
| x        | branch b2 |
| y3       | r0        |
|          | y1        |
|          | y2        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

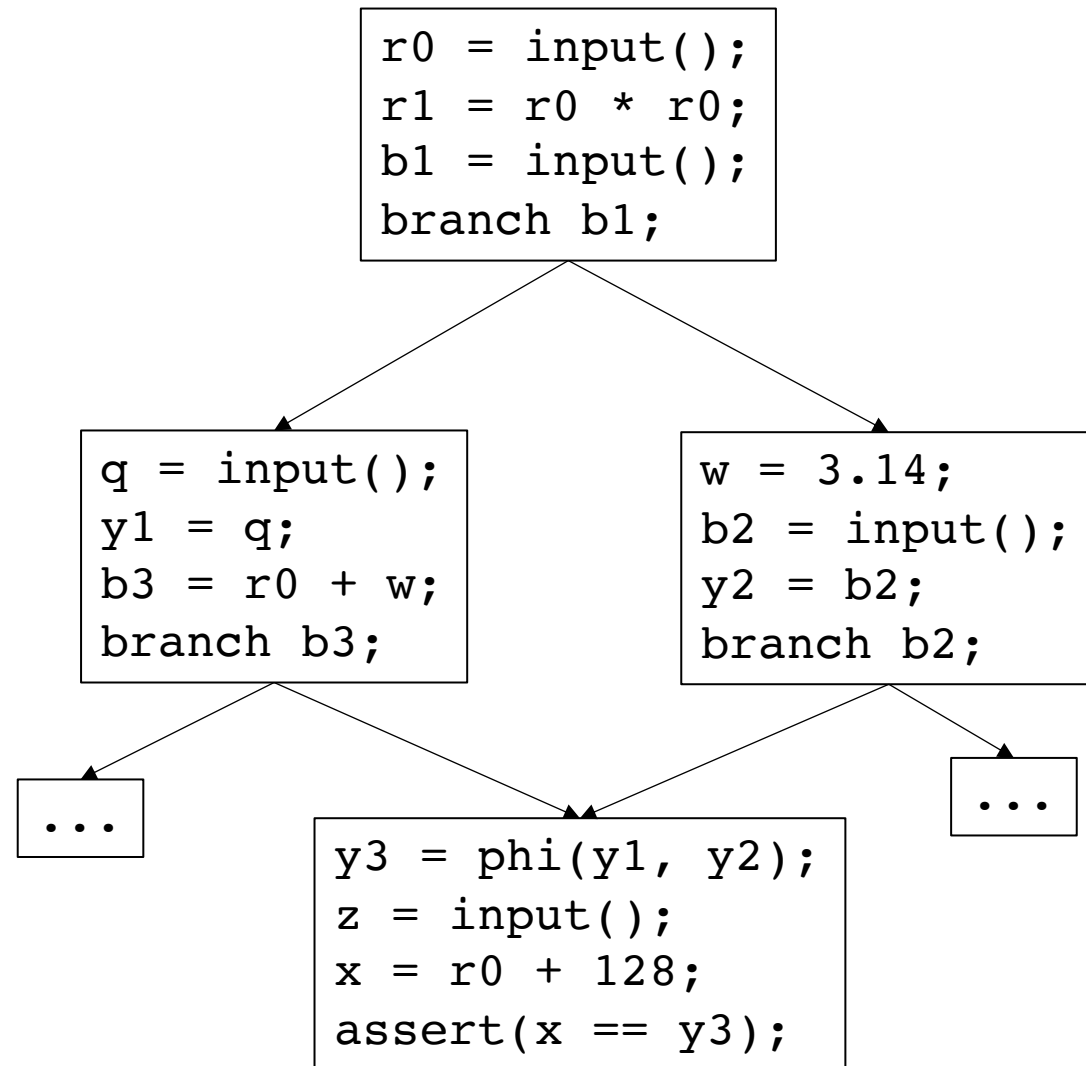
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | branch b3 |
| x        | branch b2 |
| y3       | r0        |
|          | y1        |
|          | y2        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

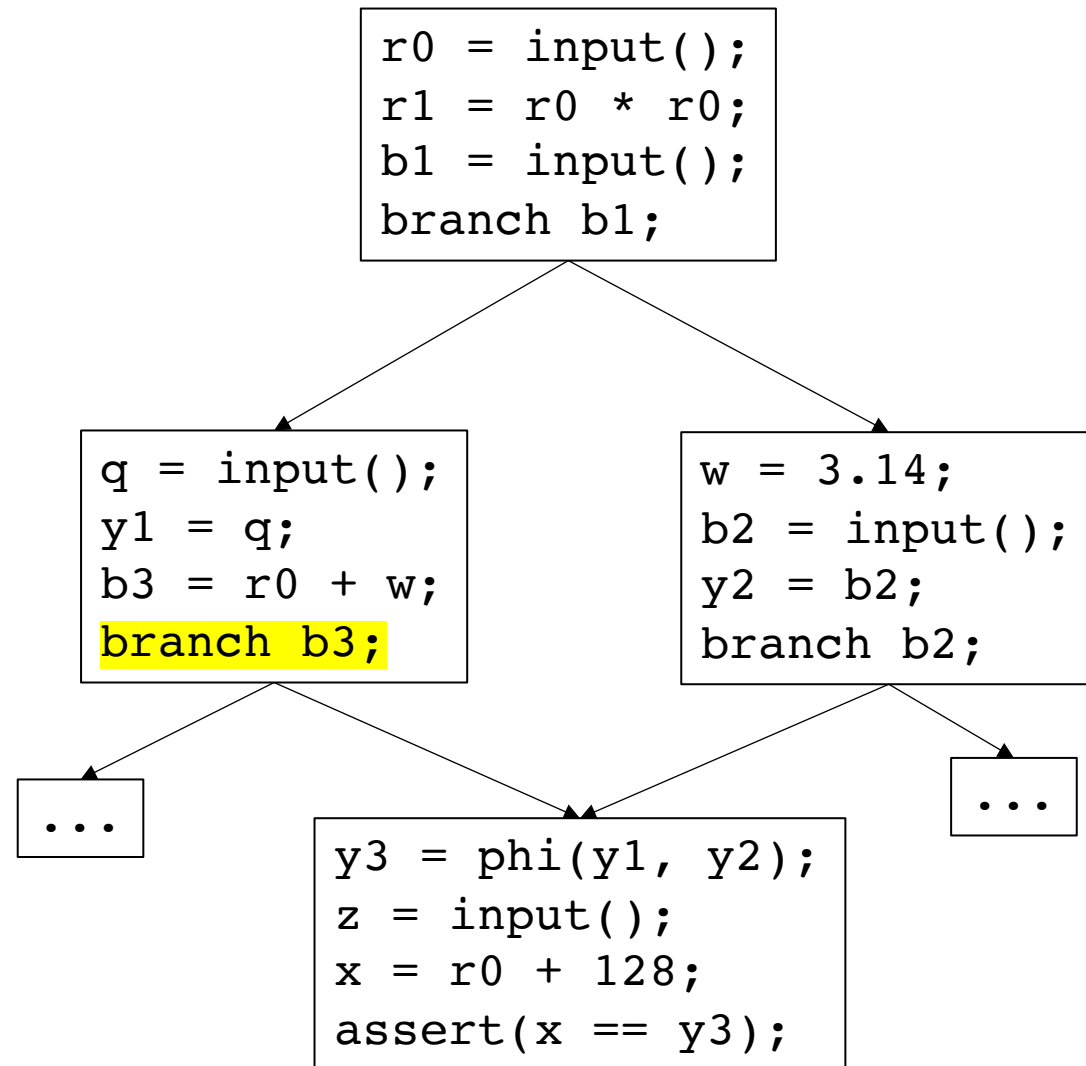
|          |           |
|----------|-----------|
| marked:  | worklist: |
| assert() | branch b3 |
| x        | branch b2 |
| y3       | r0        |
|          | y1        |
|          | y2        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

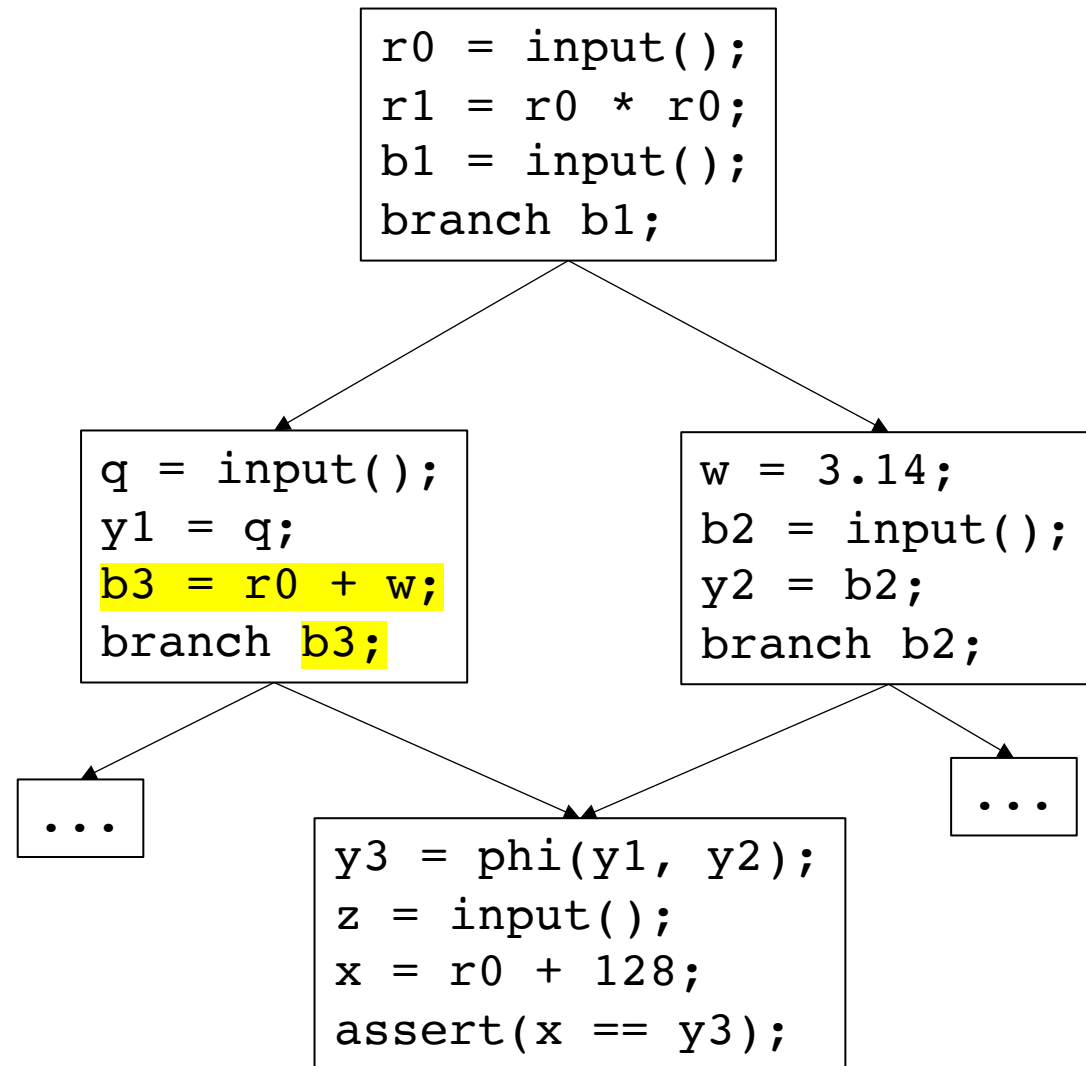
|           |           |
|-----------|-----------|
| marked:   | worklist: |
| assert()  | branch b2 |
| x         | r0        |
| y3        | y1        |
| branch b3 | y2        |



# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

|           |           |
|-----------|-----------|
| marked:   | worklist: |
| assert()  | branch b2 |
| x         | r0        |
| y3        | y1        |
| branch b3 | y2        |
|           | b3        |



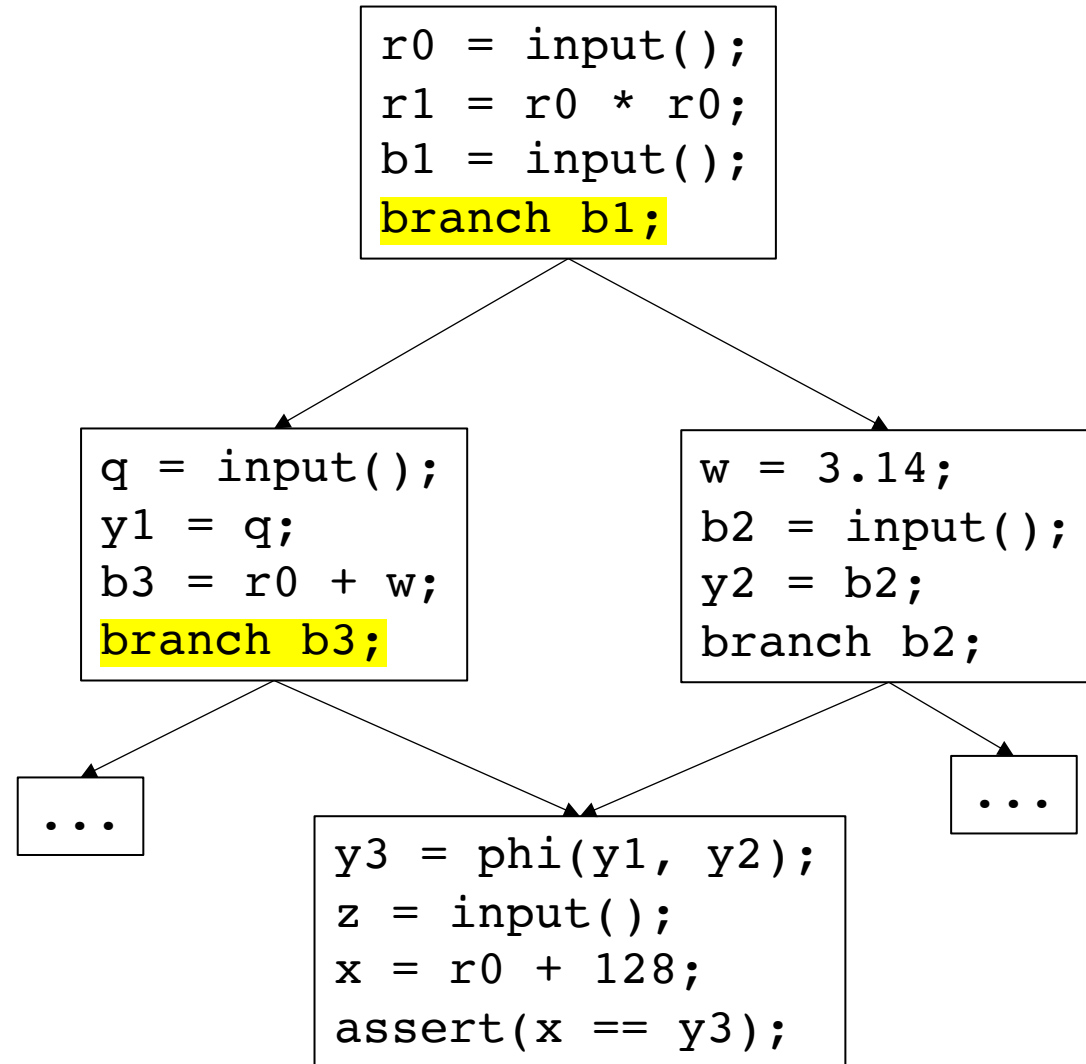


# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

worklist:  
branch b2  
r0  
y1  
y2  
b3  
branch b1

marked:  
assert()  
x  
y3  
branch b3



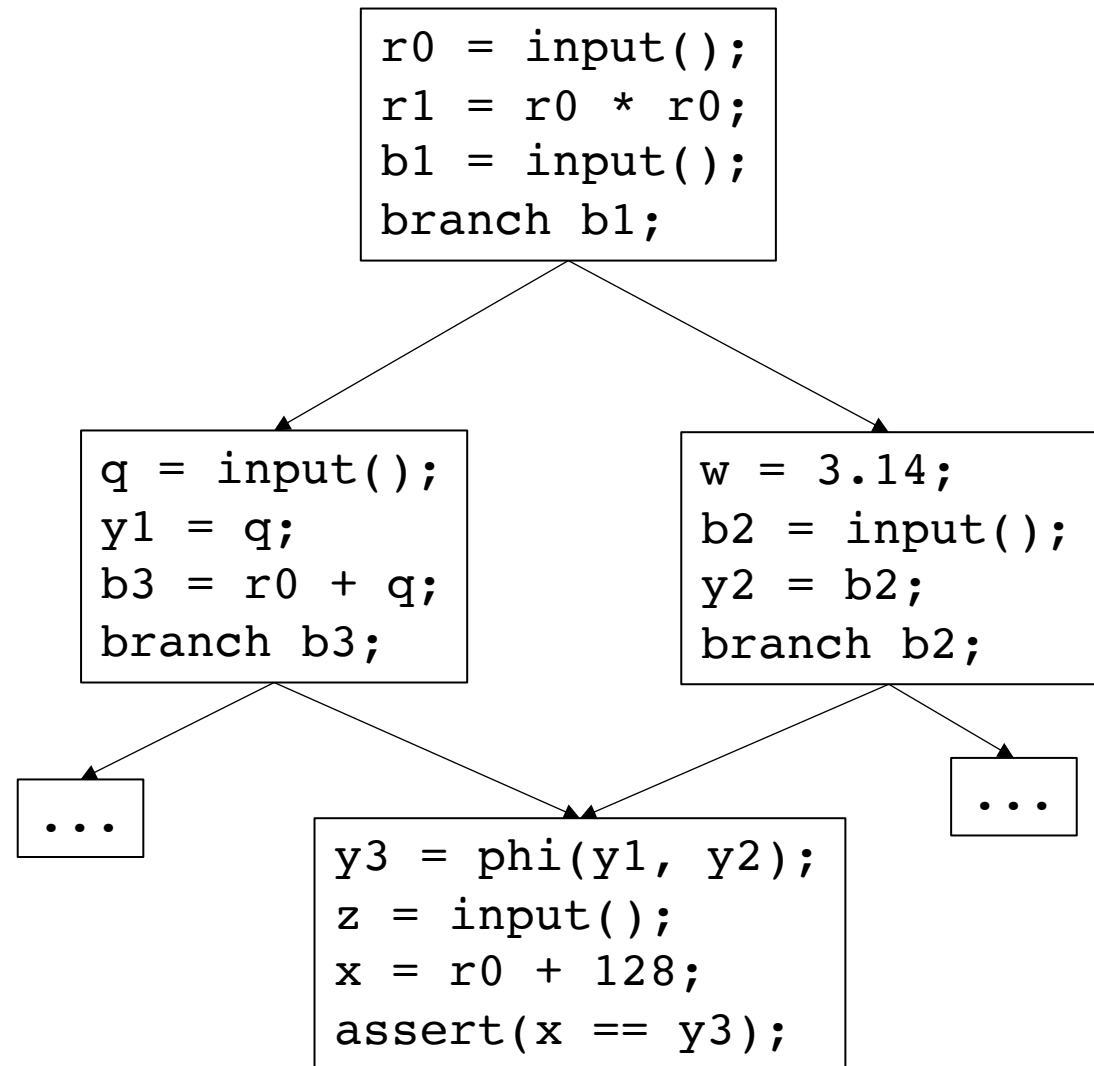
# Backwards slicing algorithm

```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());
    }
}
```

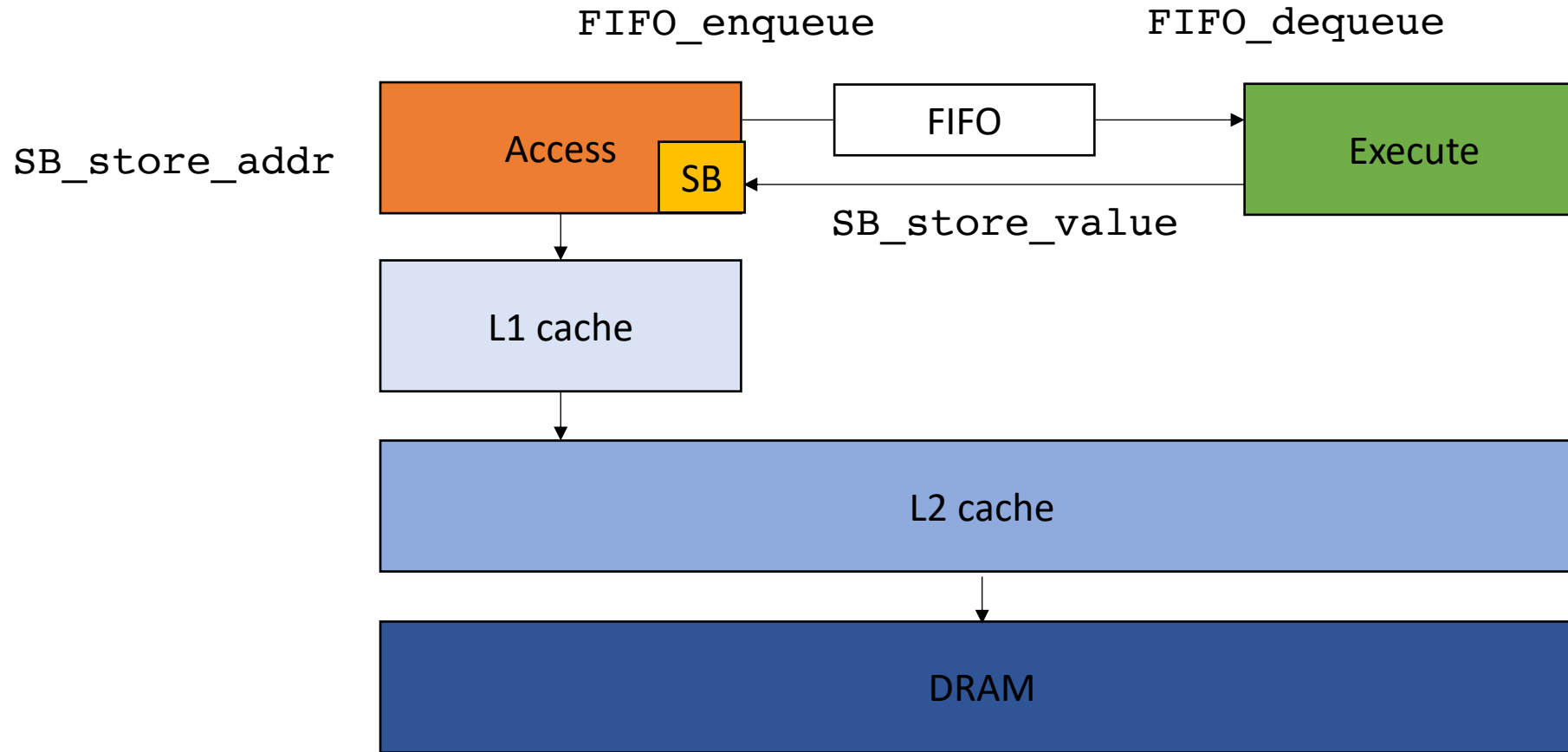
worklist:  
branch b2  
r0  
y1  
y2  
b3  
branch b1

marked:  
assert()  
x  
y3  
branch b3

rest of example  
is an exercise



# Back to DAE



# Compiler

## Step 1: compile to SSA

```
for (int i = 0; i < SIZE; i++) {  
    a[i] = b[i] * 3.14;  
}
```



```
// SSA pseudo code  
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```

# Compiler

**Step 2:** Create two copies, one for the access and one for the execute

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```

# Compiler

**Step 3:** Replace loads in Execute with FIFO reads, stores with SB\_store\_values

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```

# Compiler

## Step 3: Replace loads in Execute with FIFO reads

### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue();
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```

# Compiler

**Step 4:** Enqueue loaded values on the Access. Store addresses instead of values

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue()
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```



# Compiler

**Step 4:** Enqueue loaded values on the Access. Store addresses instead of values

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    FIFO_enqueue(r0);
    float r1 = r0 * 3.14;
    SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue()
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```

# Compiler

**Step 5:** Slice the Execute on all FIFO dequeue and SB store value calls

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    FIFO_enqueue(r0);
    float r1 = r0 * 3.14;
    SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue()
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```

# Compiler

**Step 6:** Slice the Access on all FIFO enqueue and SB store address calls

**Access**

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    FIFO_enqueue(r0);
    float r1 = r0 * 3.14;
    SB_store_addr(a + i);
}
```

**Execute**

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue()
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```

# Compiler

**Step 6:** Slice the Access on all FIFO enqueue and SB store address calls

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    FIFO_enqueue(r0);
float r1 = r0 * 3.14;
    SB_store_addr(a + i);
}
```

## Execute

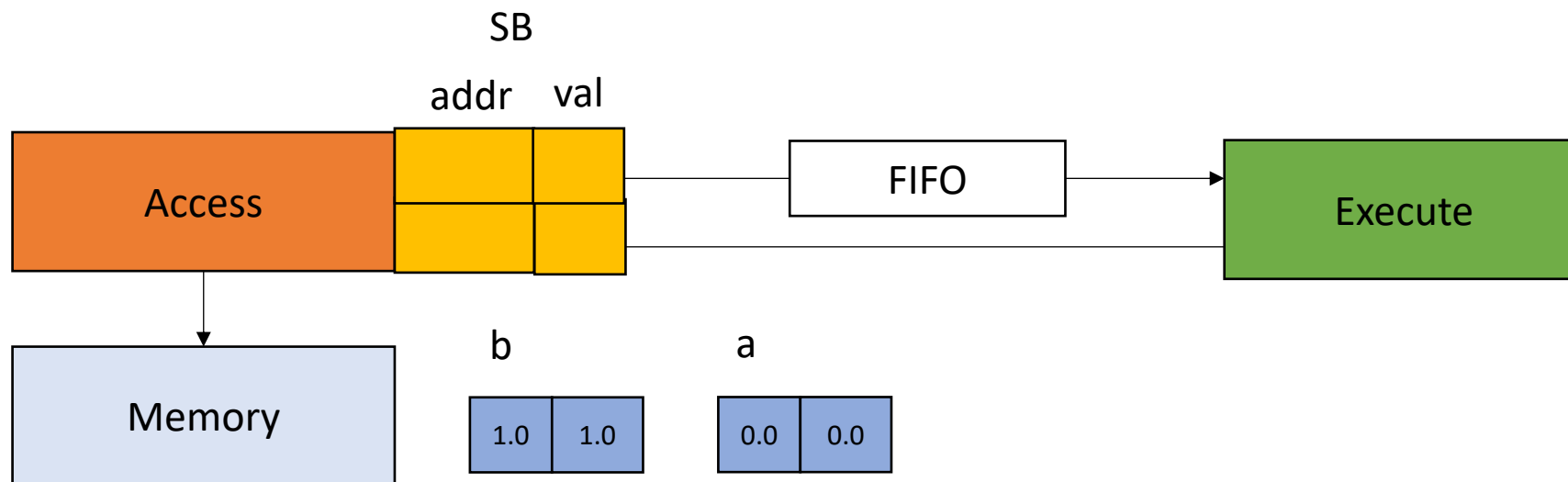
```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue()
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```

## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```



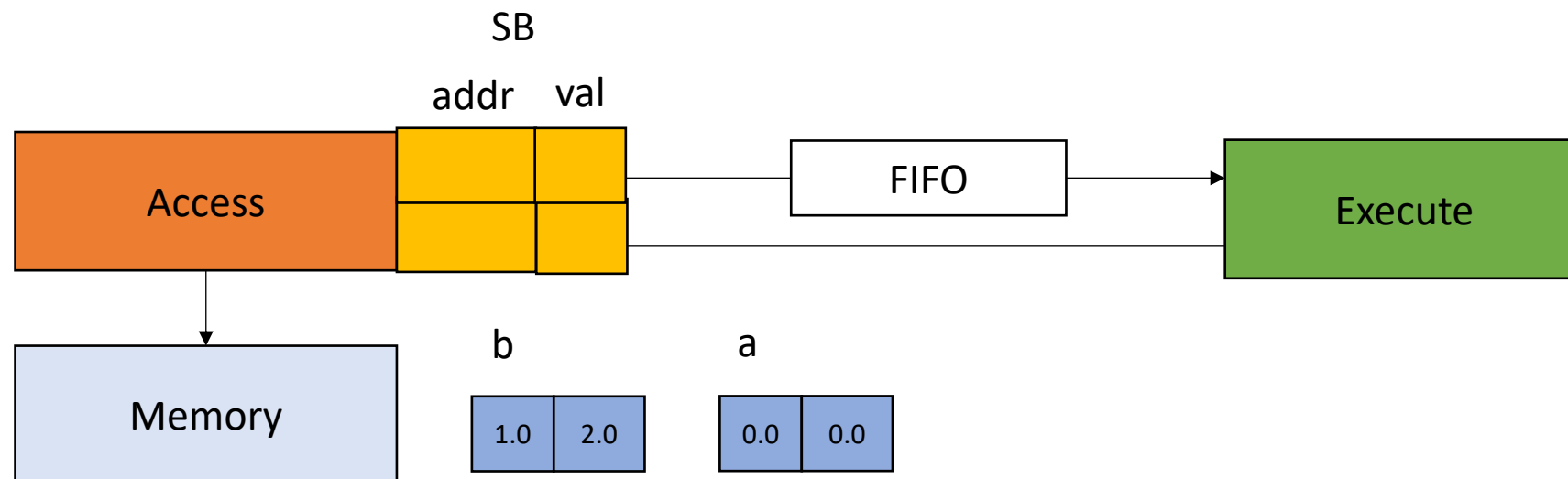
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



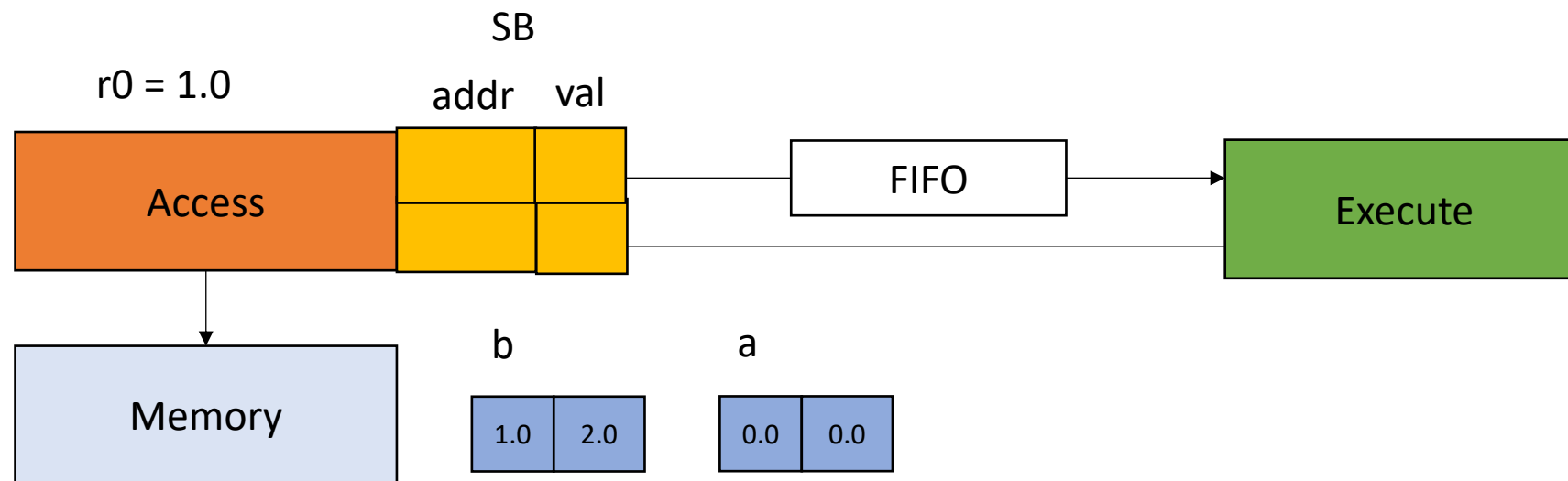
### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



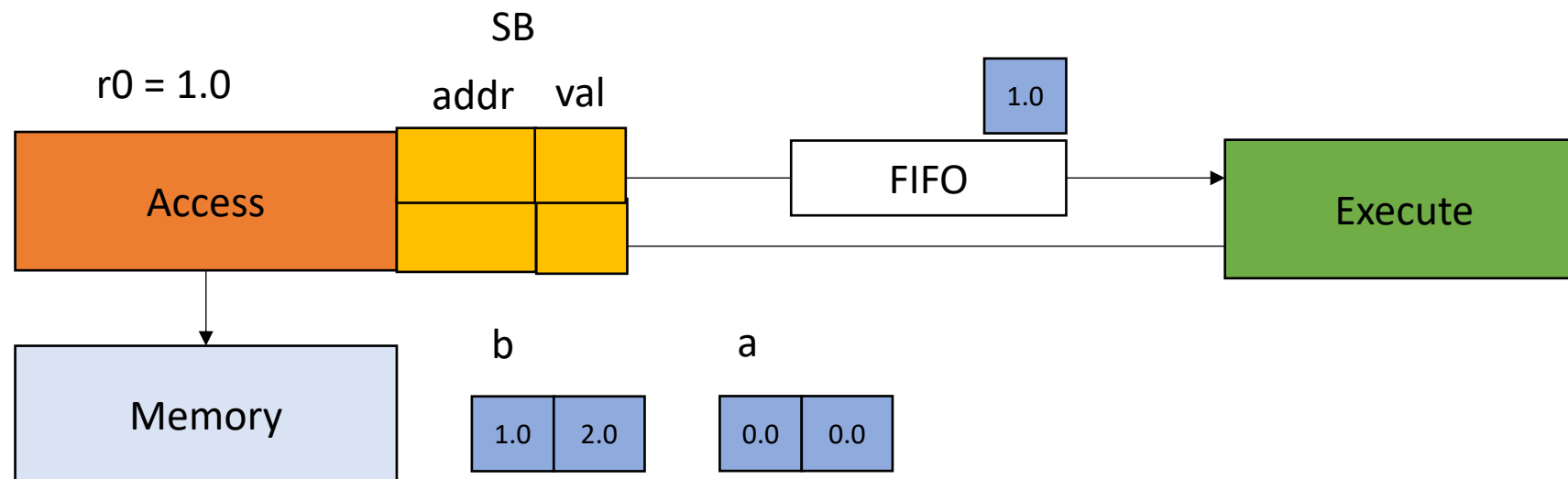
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue();
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*





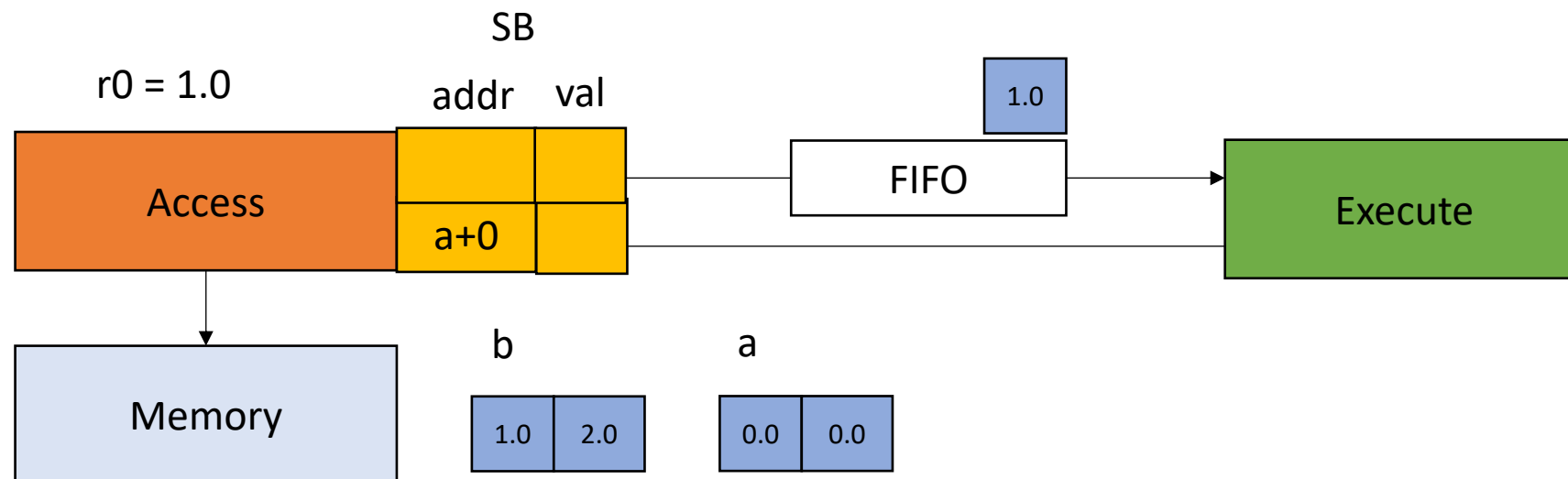
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



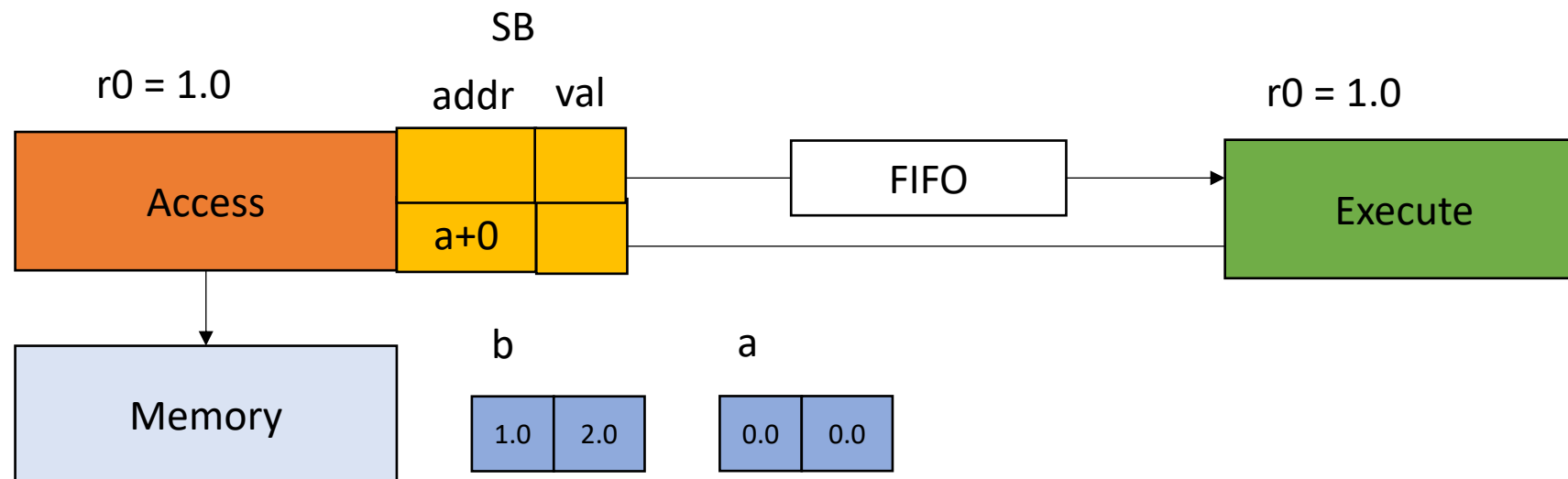
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



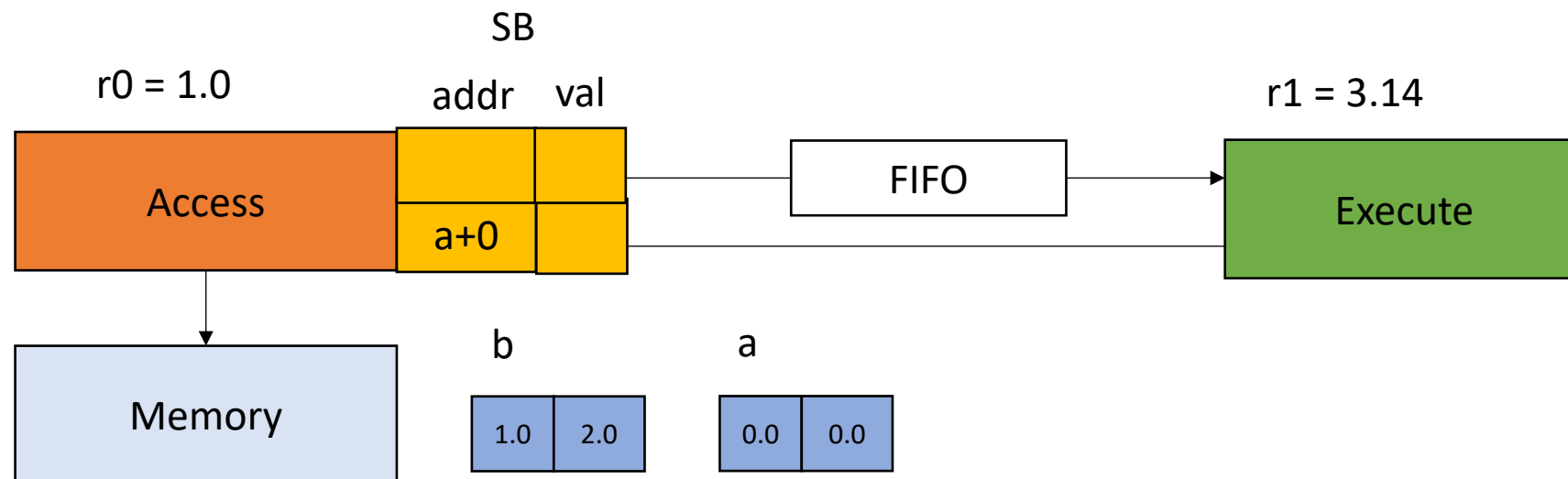
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



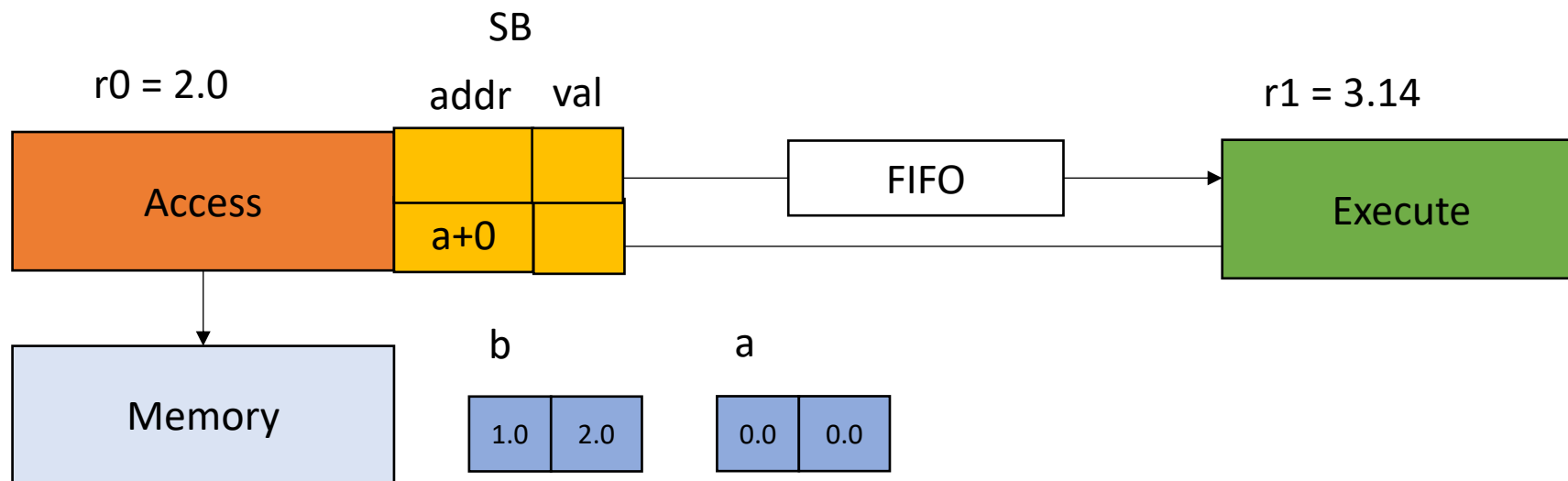
### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



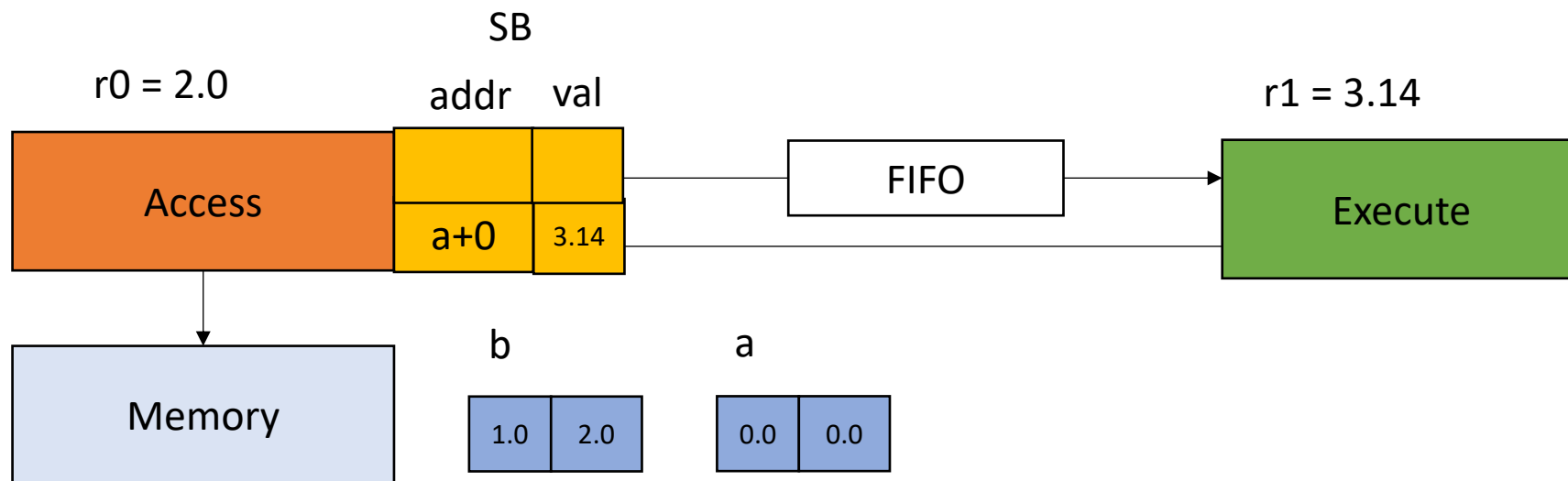
### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



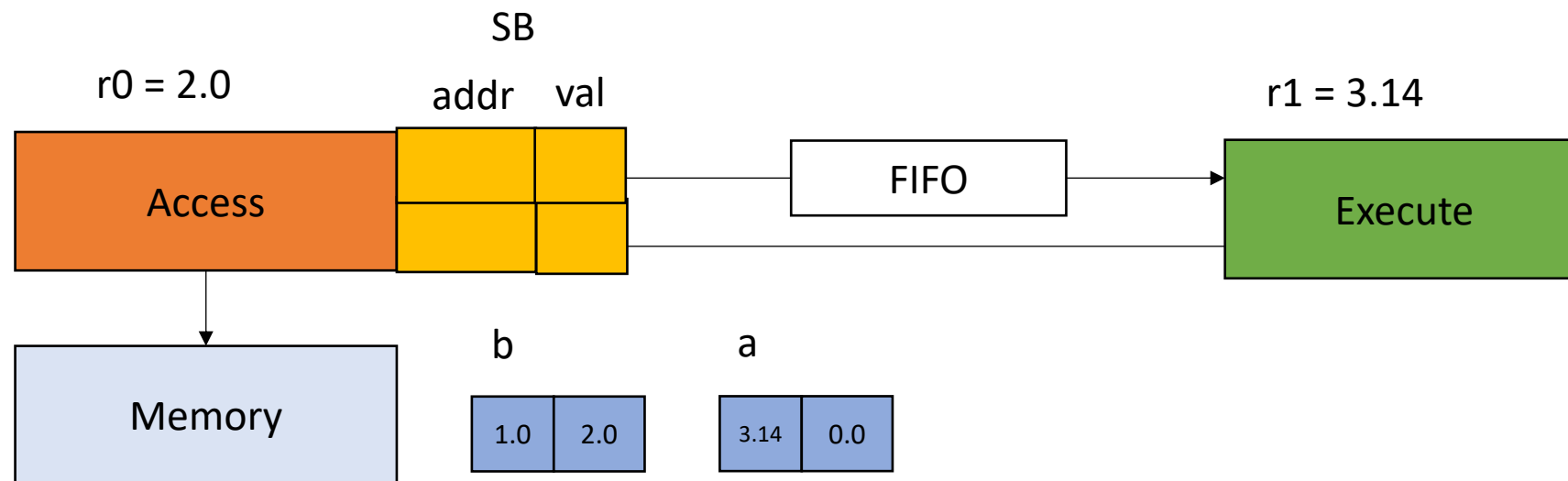
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



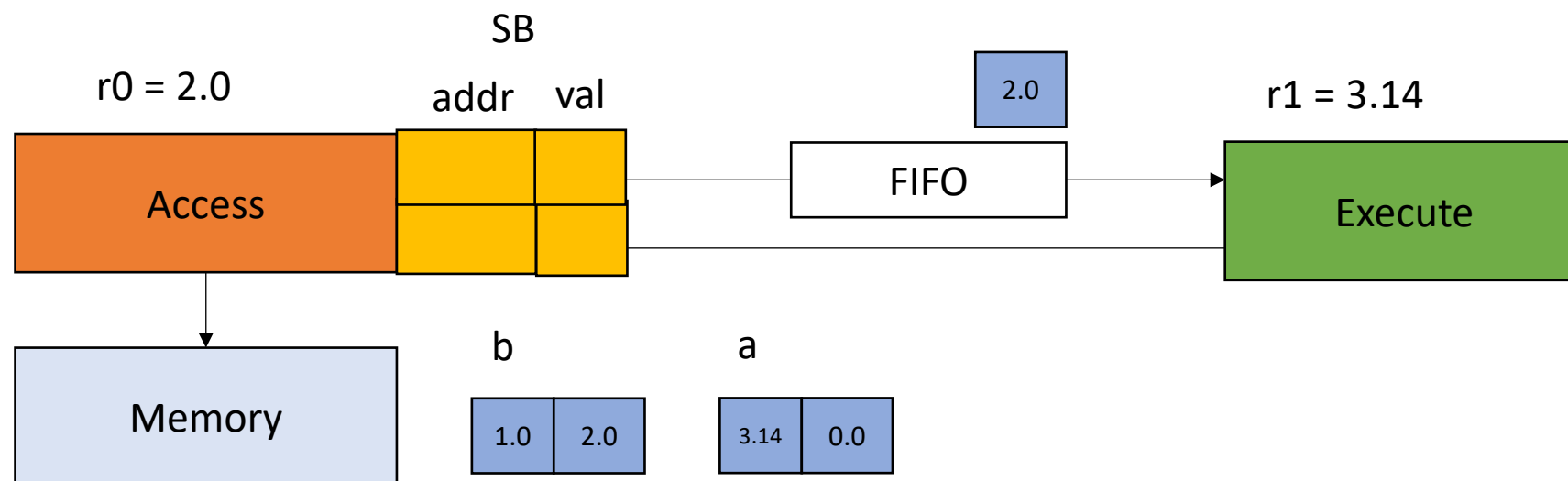
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



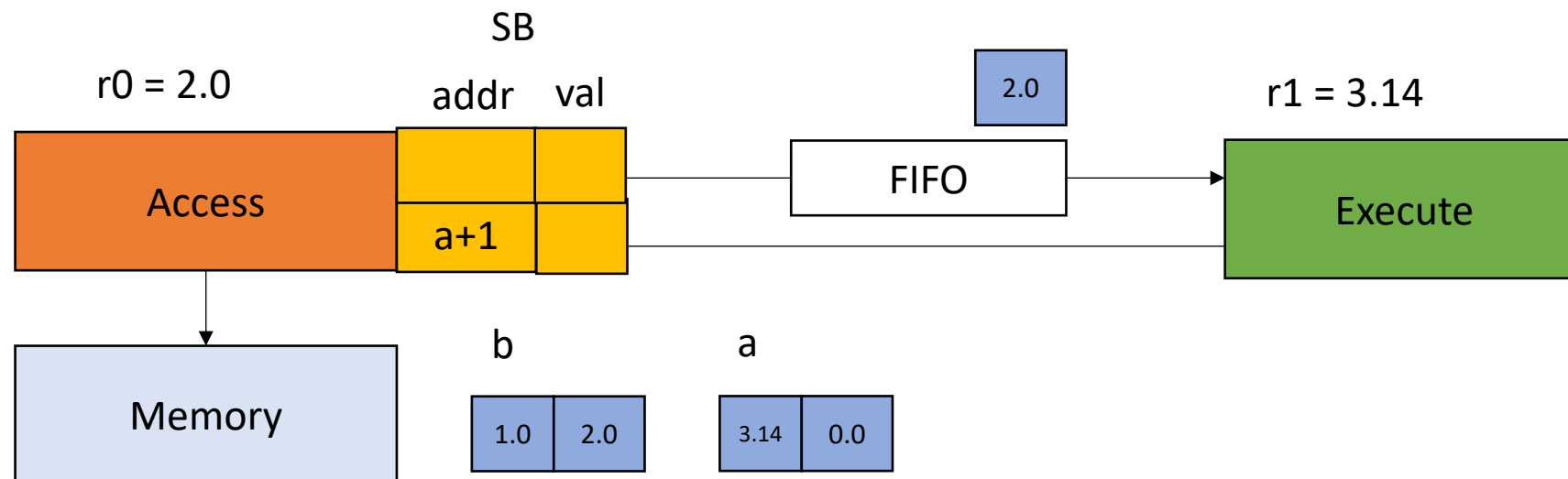
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*





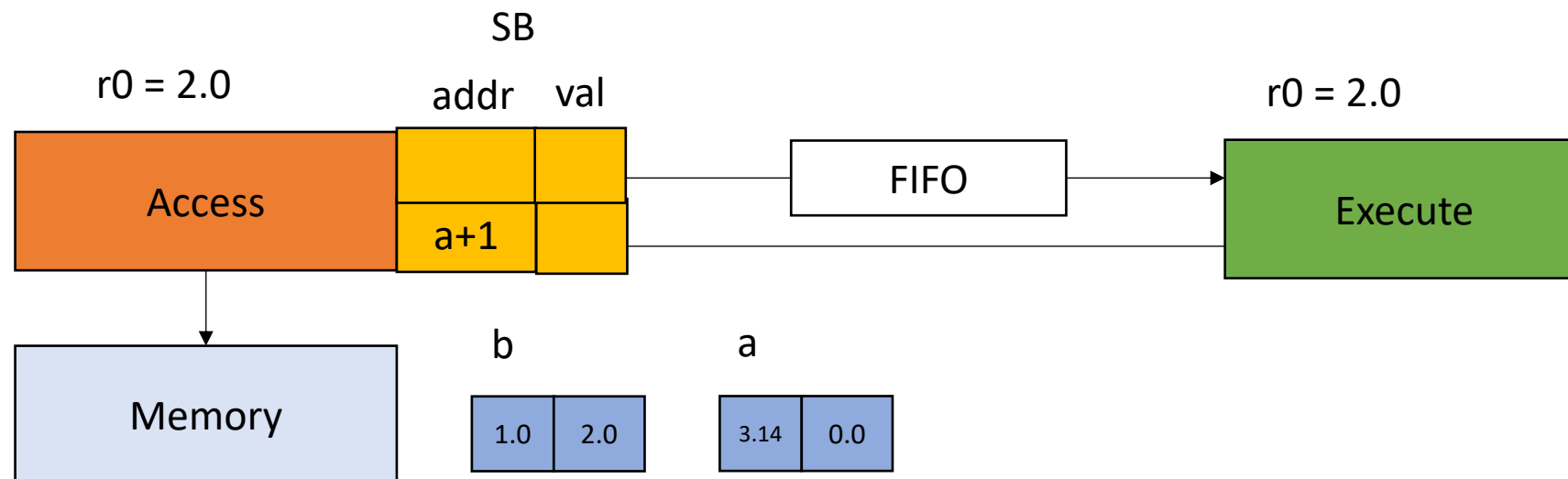
## Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

## Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue();
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



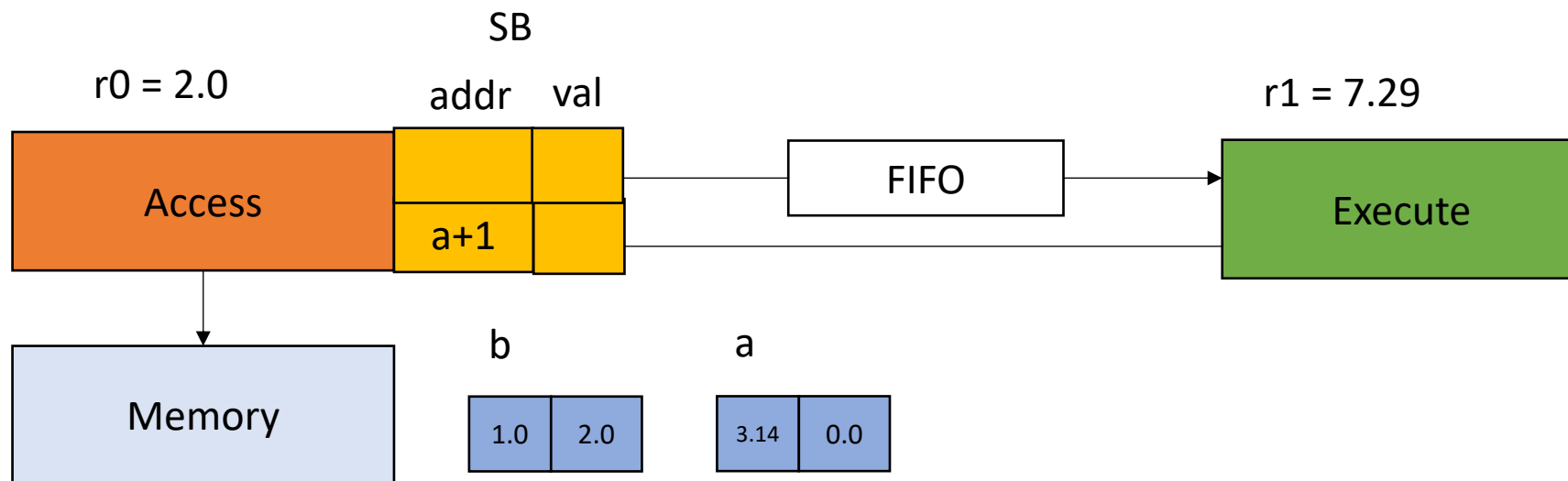
### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



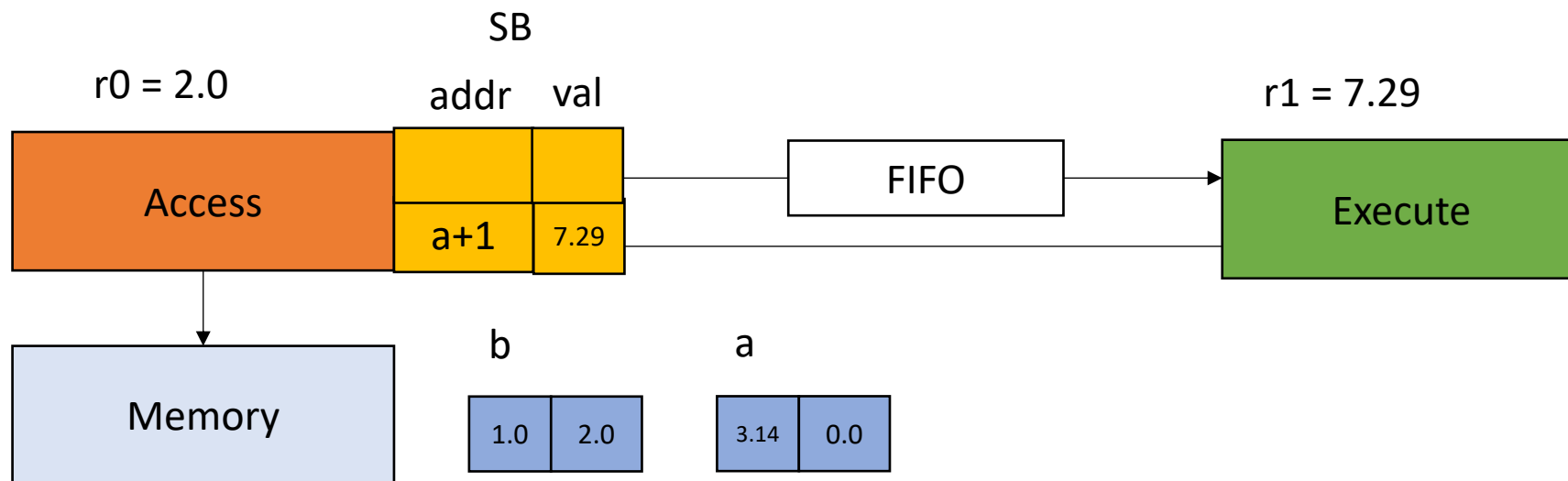
### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



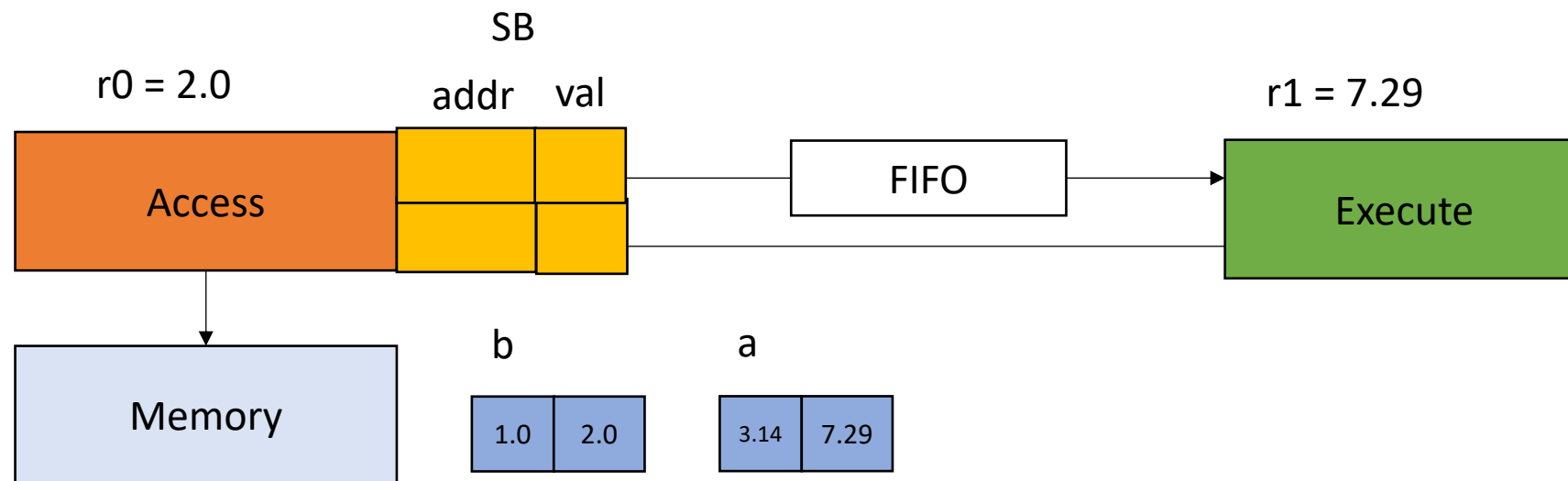
### Access

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = load(b + i);
  FIFO_enqueue(r0);
  SB_store_addr(a + i);
}
```

### Execute

```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
  float r0 = FIFO_dequeue()
  float r1 = r0 * 3.14;
  SB_store_value(r1);
}
```

*blocks until queue  
has an item to dequeue*



# Performance bounds

- A program  $p$  has execution time of  $E(p)$ . The time spent on compute (arithmetic) is  $C(p)$ . The time spent on memory latency is  $M(p)$ .
- For many scientific kernels we can approximate:  $E(p) = C(p) + M(p)$
- In DAE, the Execute time ideally is  $C(p)$ , and the Access ideally is  $M(p)$ .
- Optimistic estimates of DAE performance is
  - $\max(C(p), M(p))$
  - best case is when  $C(p) \sim M(p)$ , we get  $2x$  performance increase
- Pros/cons?

# Next week

- Coherence issues in DAE
- Optimizing loads in DAE