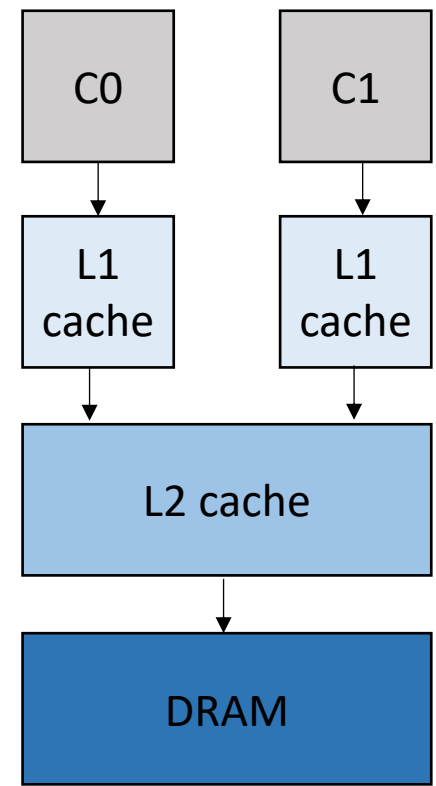


CSE211: Compiler Design

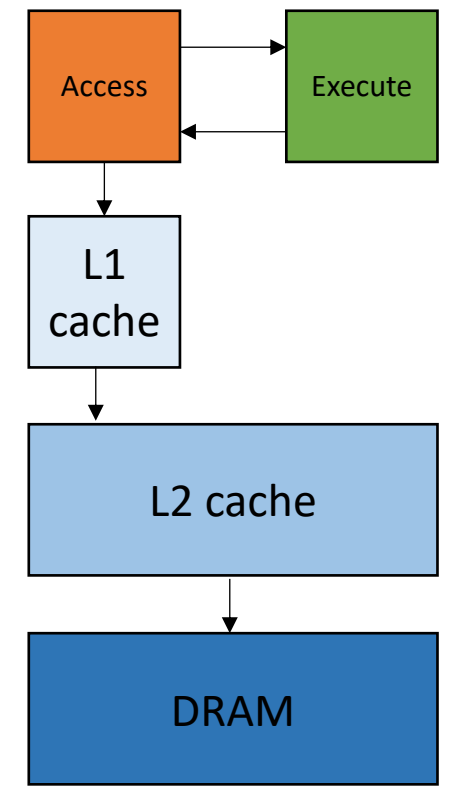
Nov. 22, 2022

- **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

- Homework 4 is released
 - pair up ASAP please
 - Exploring a DSL
 - pair assignment
 - due Dec. 5 (two weeks)
 - 6 page report
 - Discussion
- You can propose a DSL to explore, but you must do it within 1 week and there is no guarantee that I will approve it.

Announcements

- Working on grading HW 2 still
 - Sorry for the delay but it will be out soon!

Paper and project proposals

- Remember
 - Reports are due the day of the final: Dec 8.
 - I highly suggest not saving these until the last minute. They have a late deadline to give you flexibility, not to enable procrastination.
- Project presentation:
 - Everyone doing a project needs to prepare a 12 minute presentation
 - If your project isn't finished yet, then present background and as much as you have
 - I will randomly select 6 people to present in class on Dec. 1
 - Everyone else needs to submit a screen cast of the presentation.

Paper and project proposals

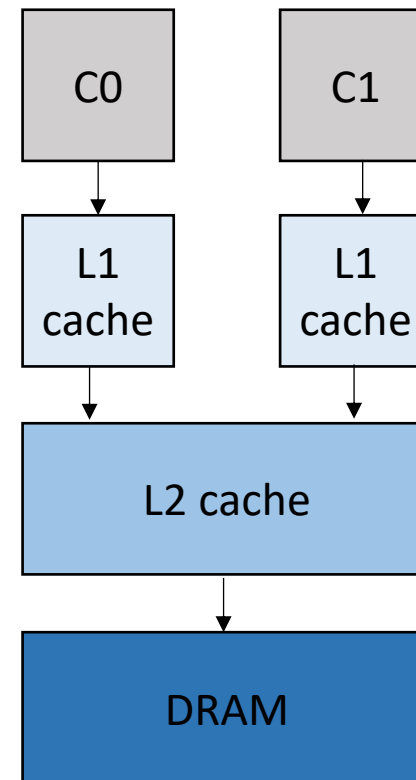
- Final:
 - Dec. 8
 - Take home (1 day instead of 1 week)
 - Similar to the midterm in length and question style
 - Designed to take ~2.5 hours not including studying
 - If you want to set aside time, 12 - 3 pm
- Same rules as midterm
 - Open note, open internet
 - Do not talk to each other about the test while it is out. Do not google for specific questions
 - ask questions as private posts on Piazza
 - No guarantees on answering questions past business hours

CSE211: Compiler Design

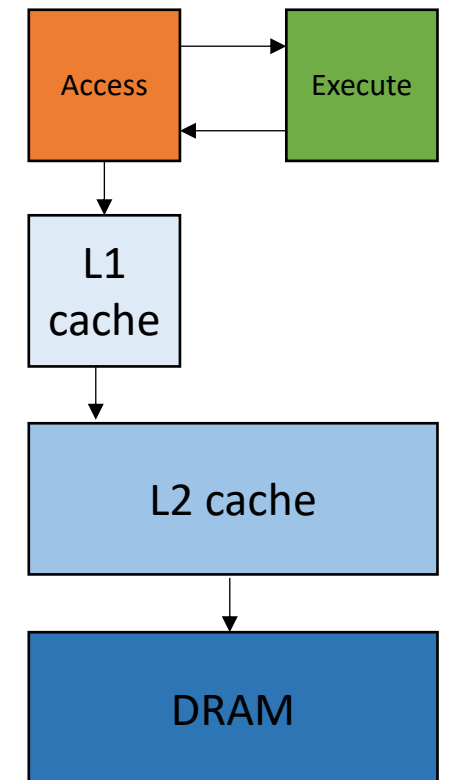
Nov. 22, 2022

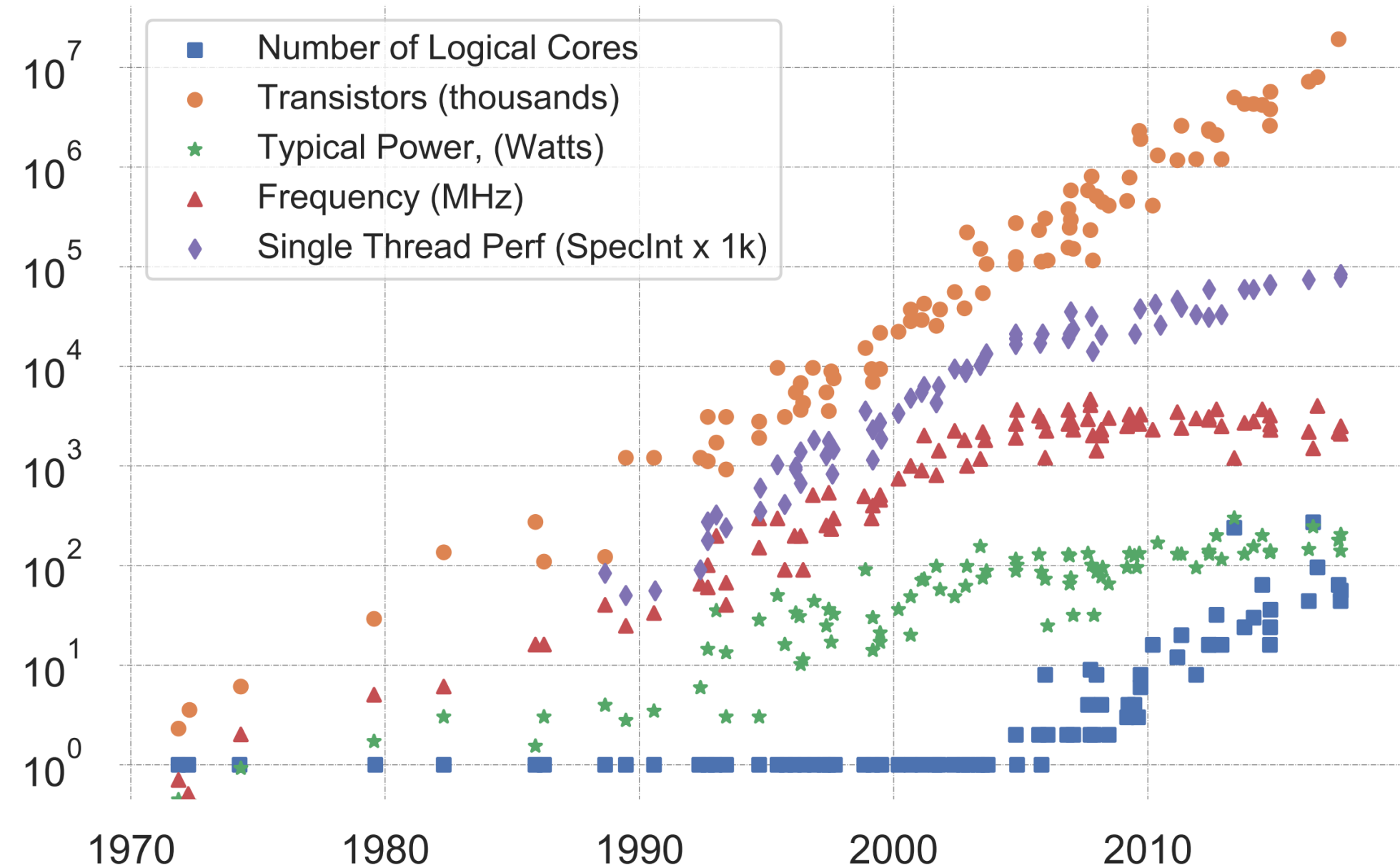
- **Topic:** Decoupled Access Execute (DAE)
- **Discussion questions:**
 - What does it mean for an application to be memory bound?
 - What are some techniques for dealing with memory bottlenecks

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.

Specialization discussion

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

Specialization discussion

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 discussion

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 discussion

- 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 discussion

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 discussion

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)

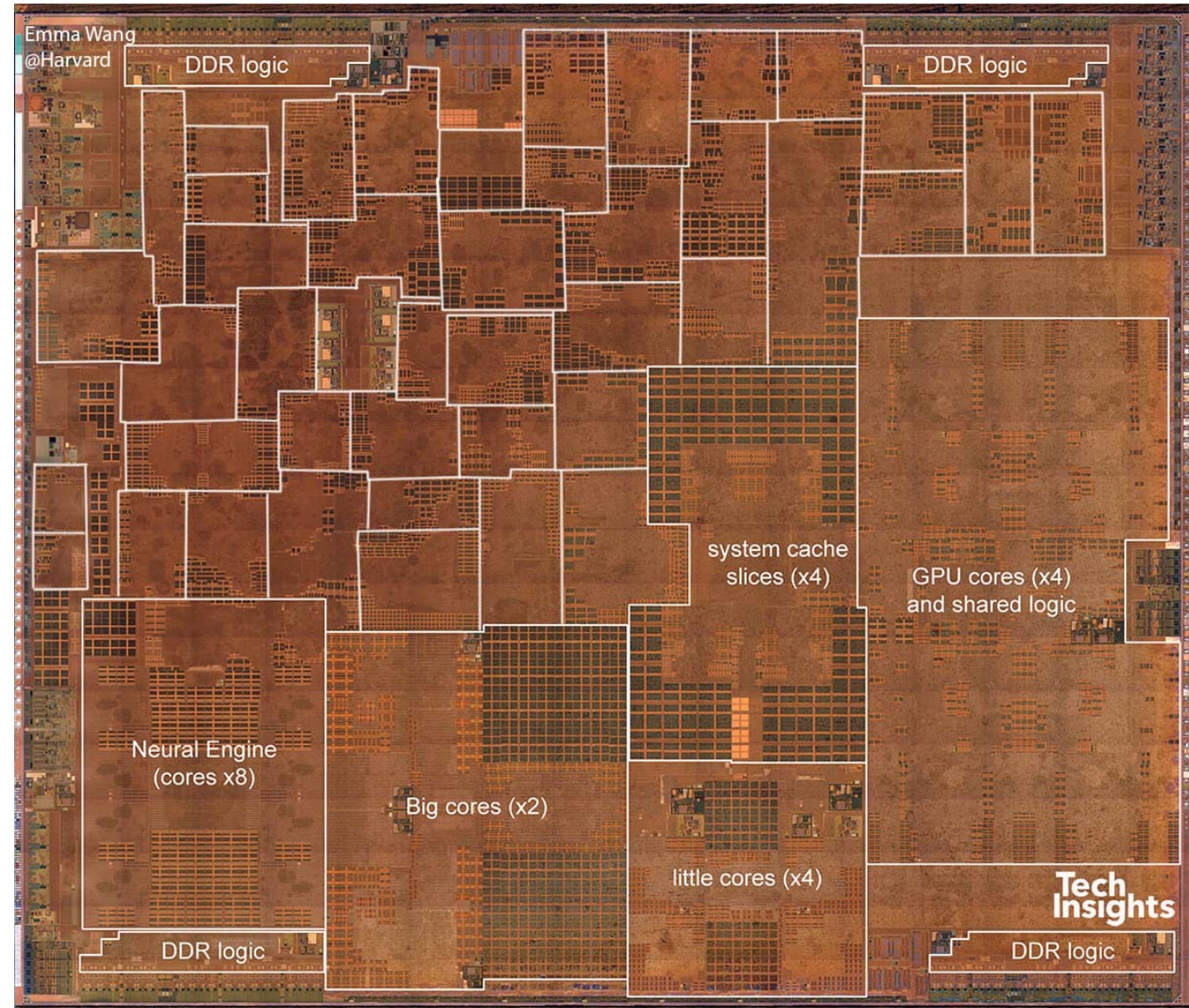
- TPUs:

- 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?

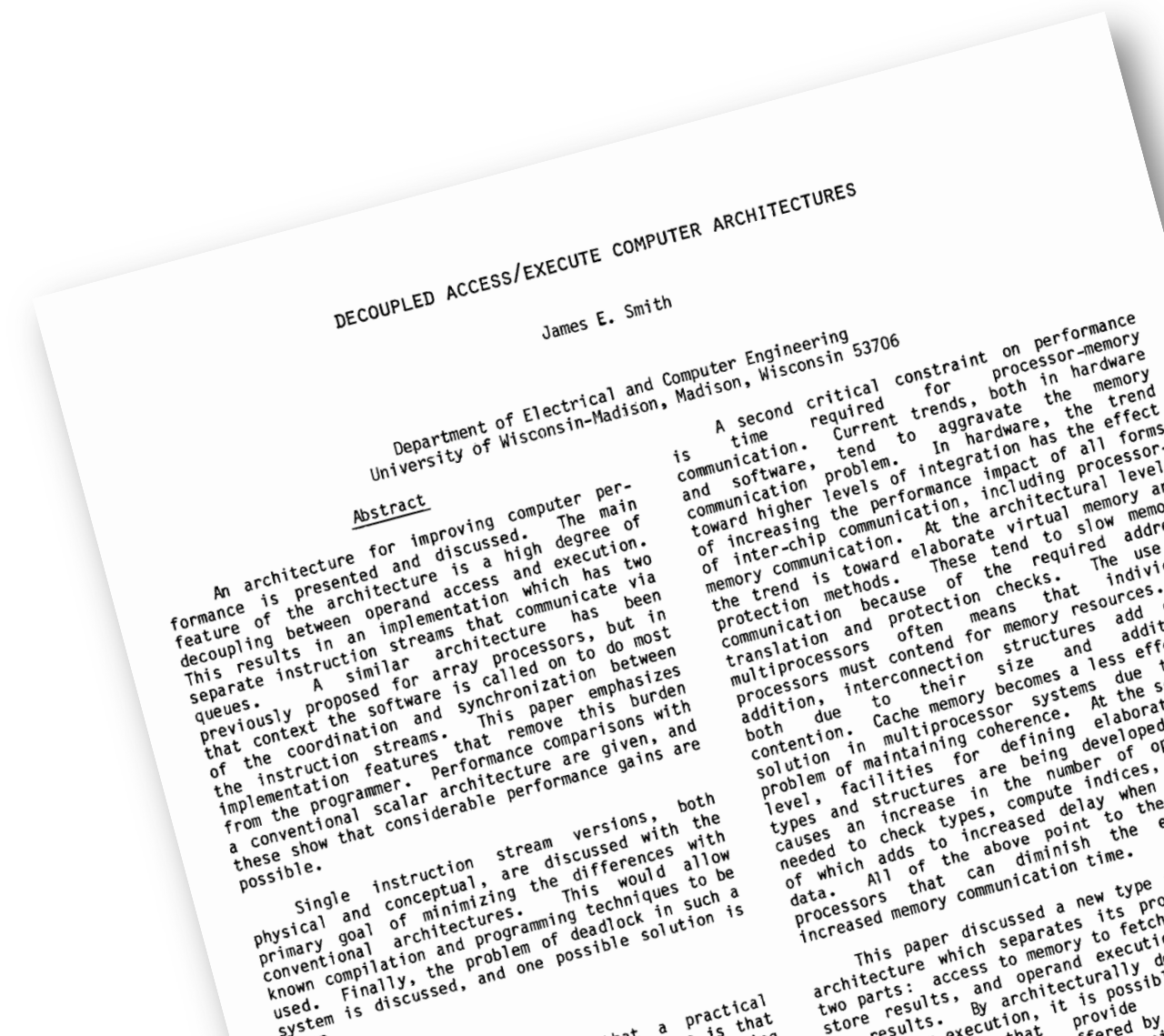
- **Programmer-centric:**
 - Programmers write specialized code that targets specific specialized processors
- **API-centric**
 - e.g. Tensorflow targets CPU, GPU, TPU
- **Hardware-centric:**
 - Hardware optimizes programs
 - Pipelining, super scalar, caches, etc. (what our traditional systems already 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

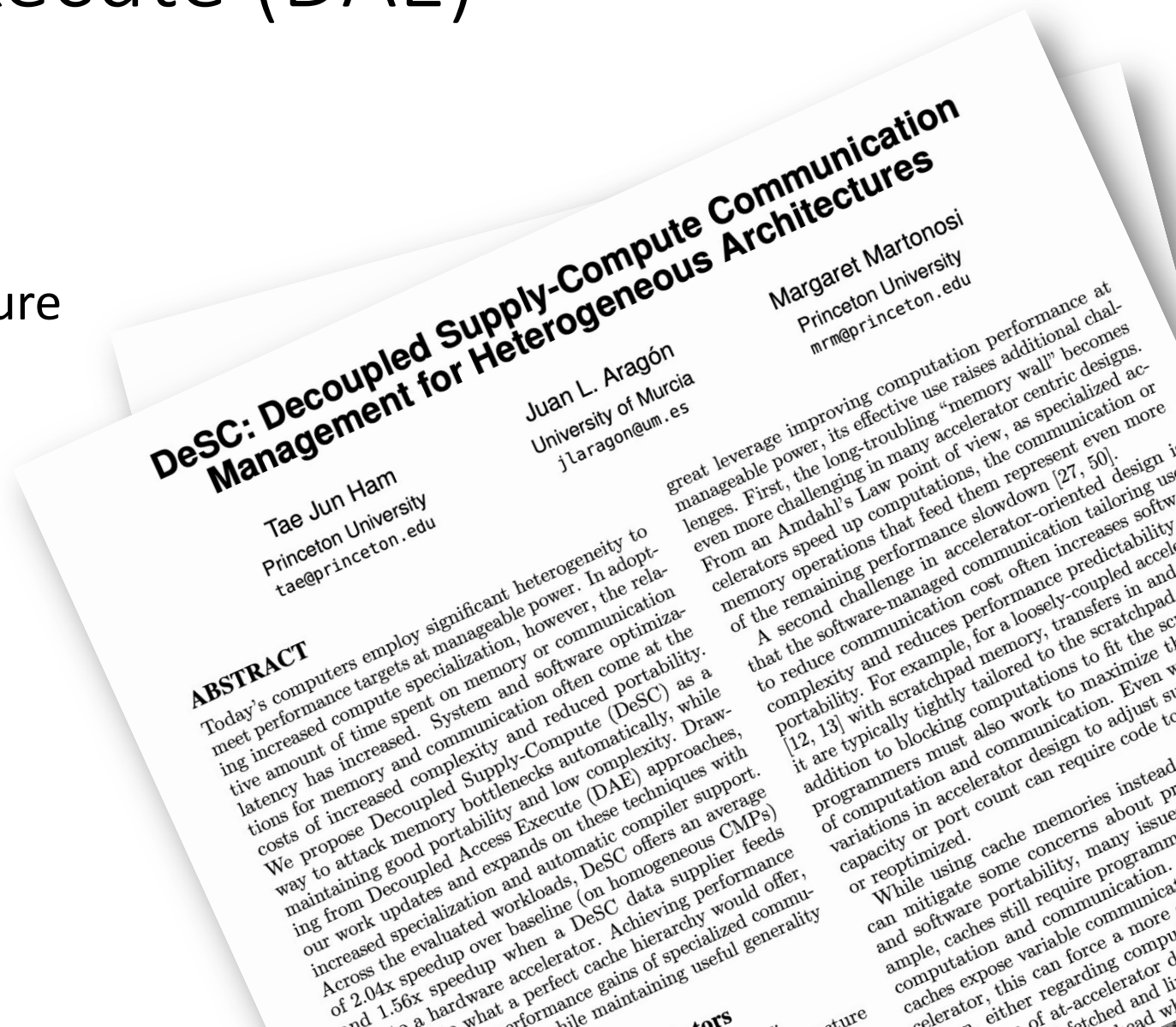
Decoupled Access/Execute (DAE)

- 1982: James E. Smith
 - Lives in Montana now and gives interesting keynotes at architecture conferences
 - “Reverse-Engineering the Brain: A Computer Architecture Grand Challenge” ISCA 2018*



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



DAE - motivation

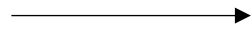
simple example program

```
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;  
}
```



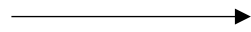
pseudo 3-address code

```
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

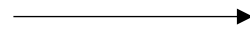


time

DAE - motivation

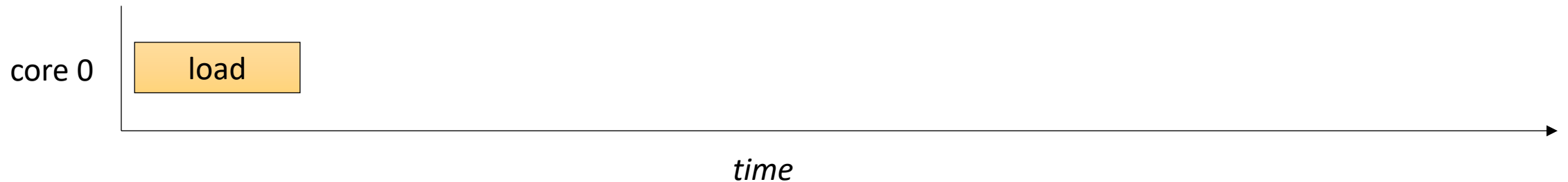
simple example program

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



pseudo 3-address code

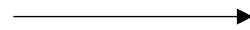
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



DAE - motivation

simple example program

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



pseudo 3-address code

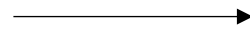
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



DAE - motivation

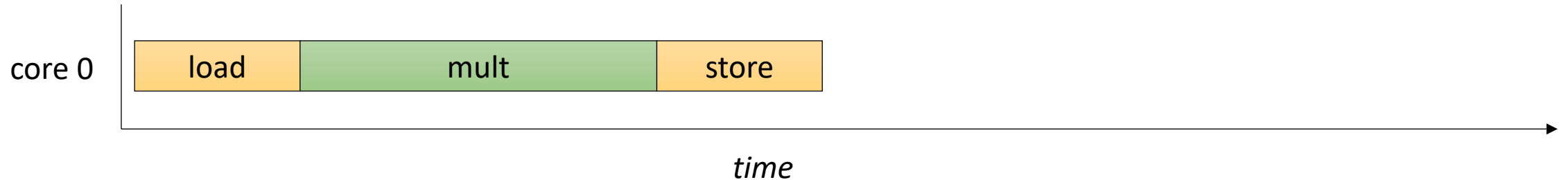
simple example program

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



pseudo 3-address code

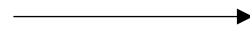
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



DAE - motivation

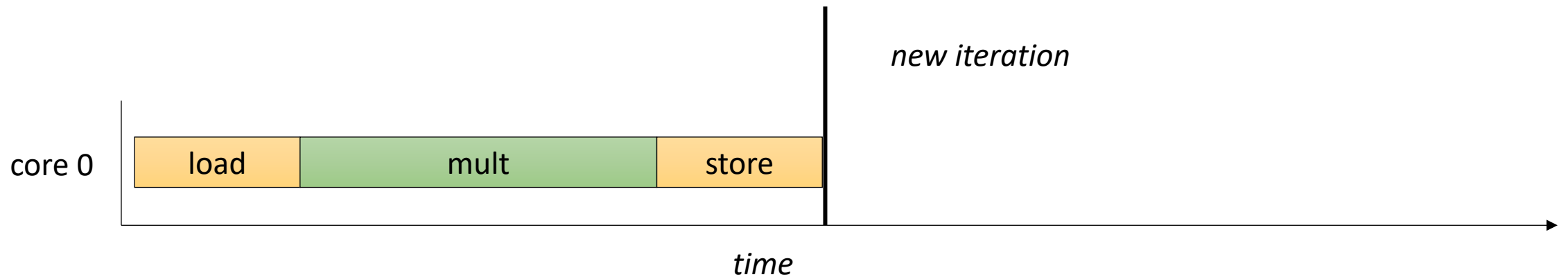
simple example program

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



pseudo 3-address code

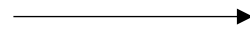
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



DAE - motivation

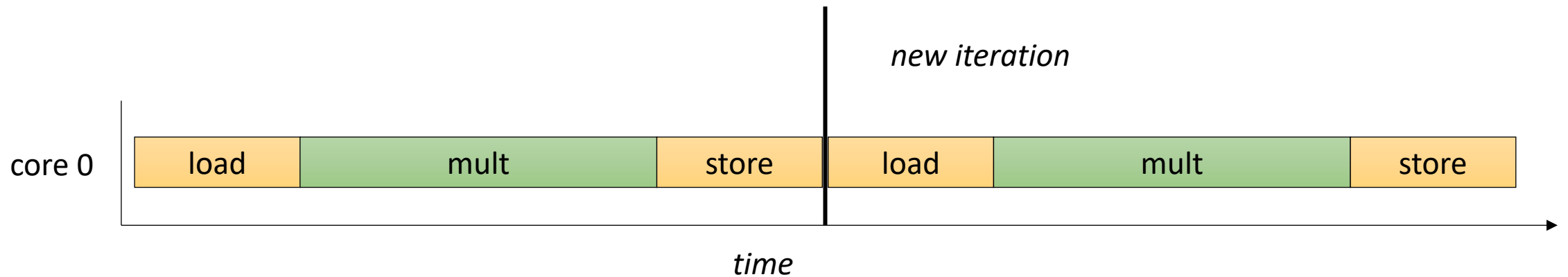
simple example program

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



pseudo 3-address code

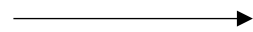
```
for (int i = 0; i < SIZE; i++) {  
    float r0 = load(b + i);  
    float r1 = r0 * 3.14;  
    store(a + i, r1);  
}
```



DAE - motivation

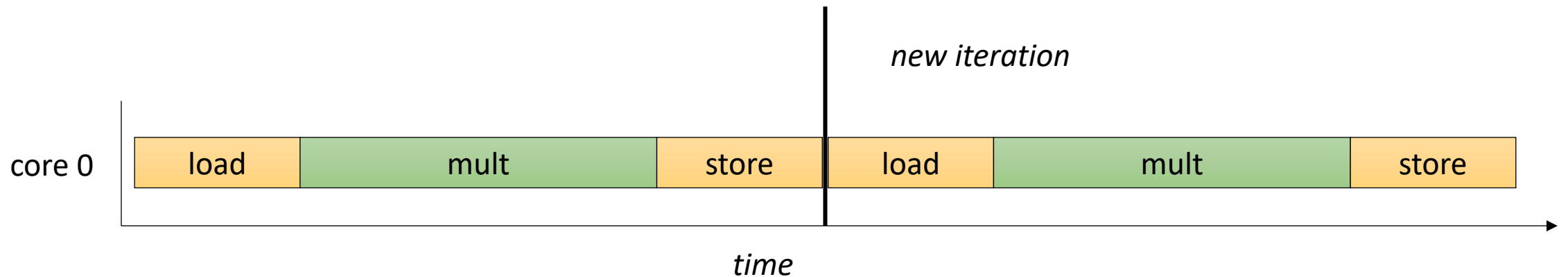
simple example program

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

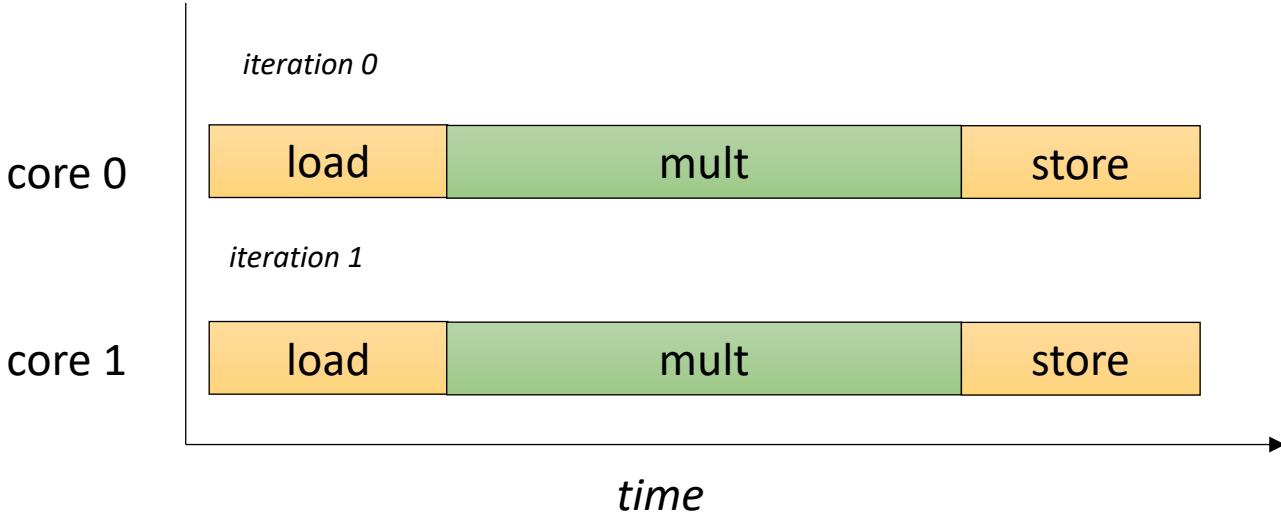


pseudo 3-address code

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



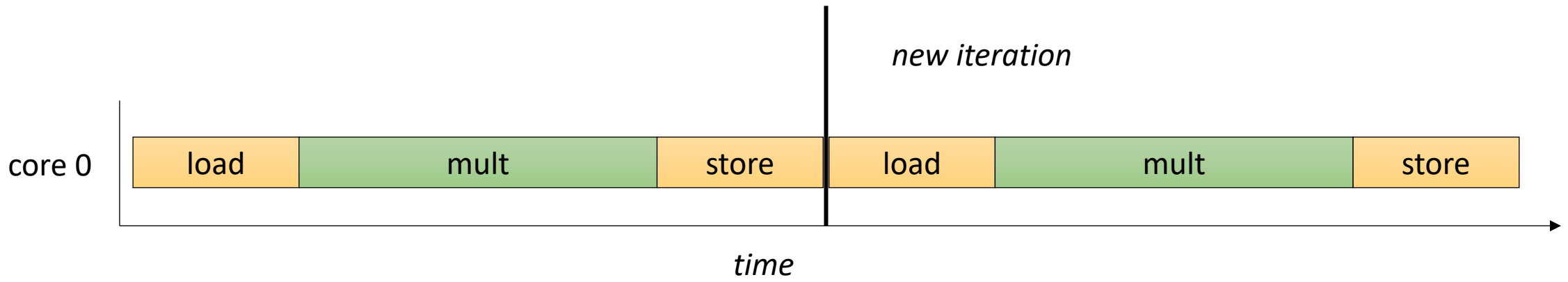
DAE - motivation



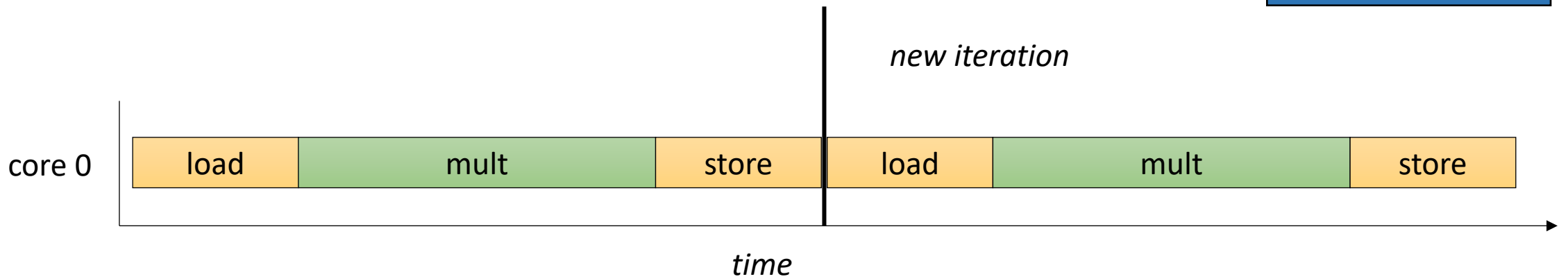
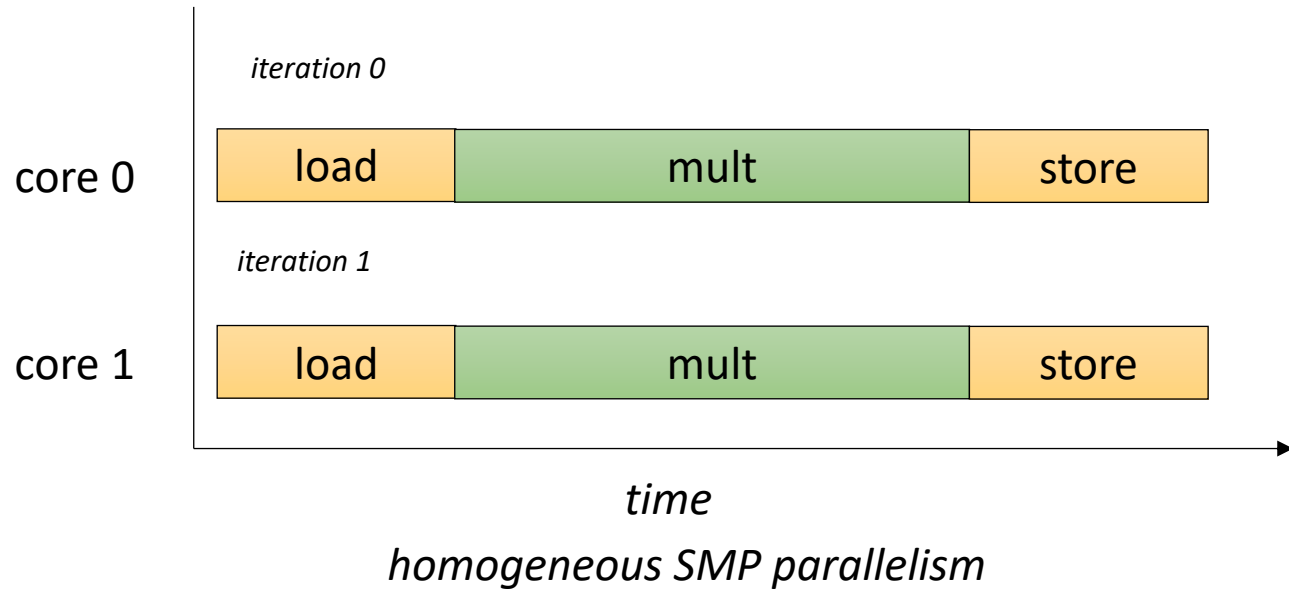
homogeneous SMP parallelism

pseudo 3-address code

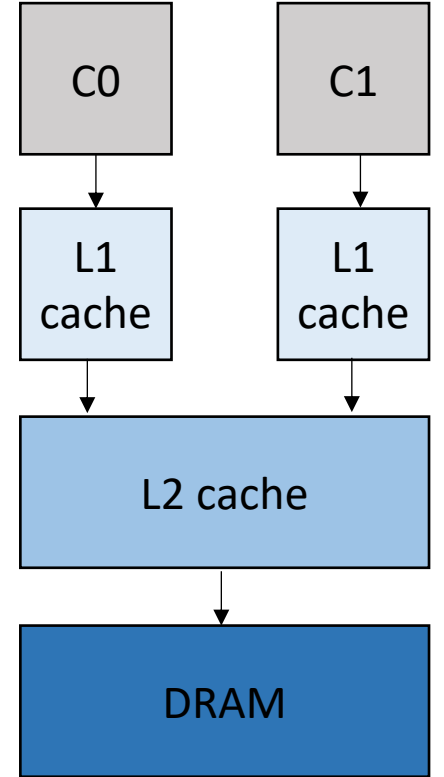
```
#pragma parallel
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
```



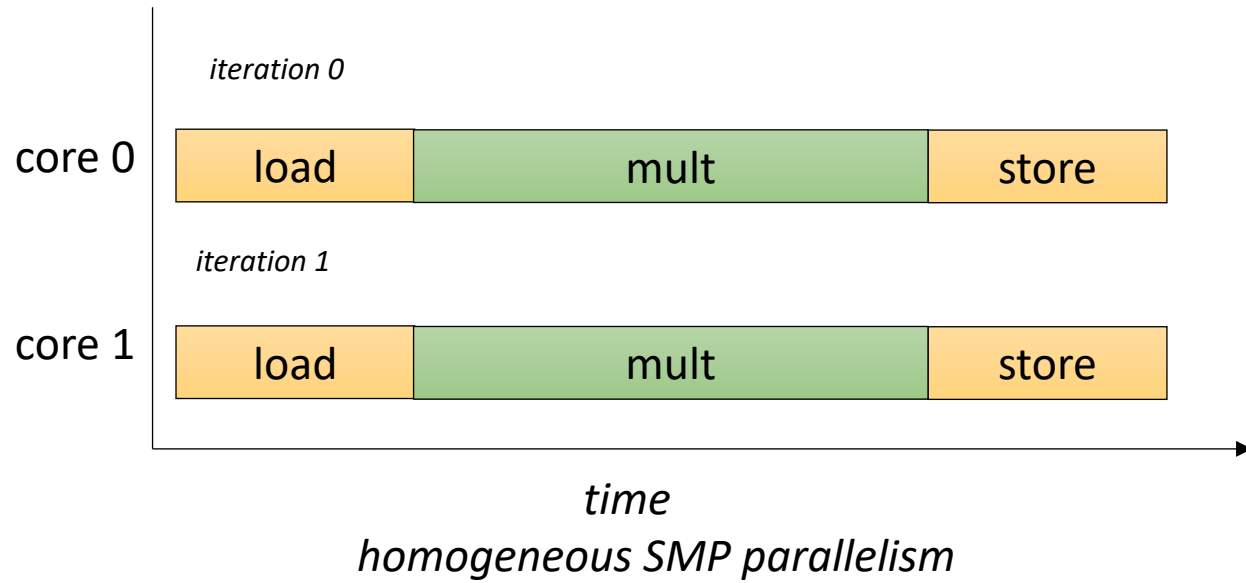
DAE - motivation



Traditional SMP System

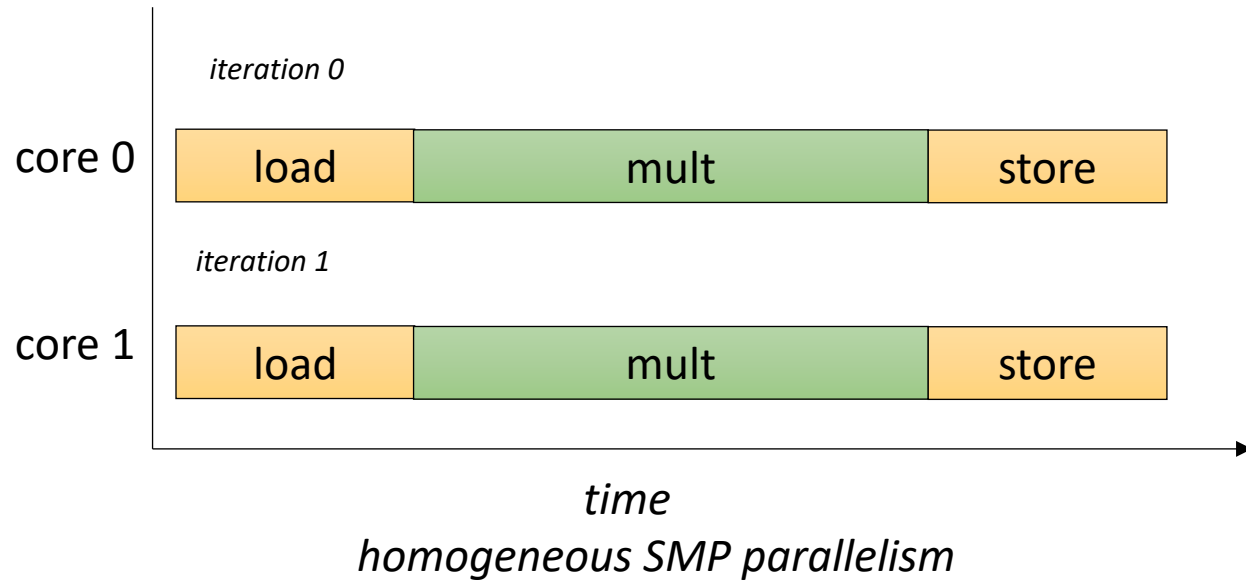


DAE - illustration

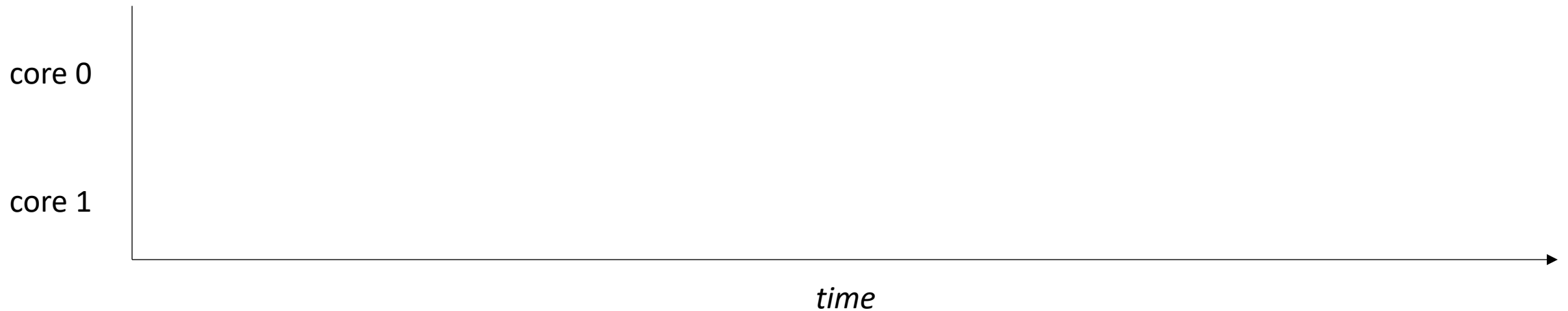


DAE: split into heterogeneous parallelism: one core does memory and one does computation

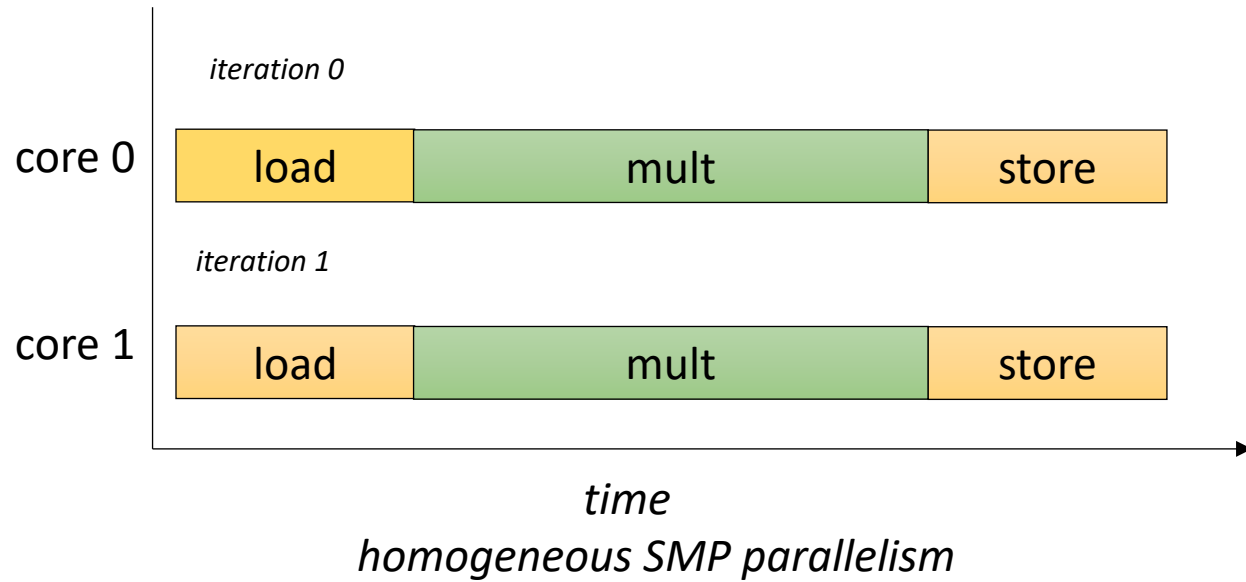
DAE - illustration



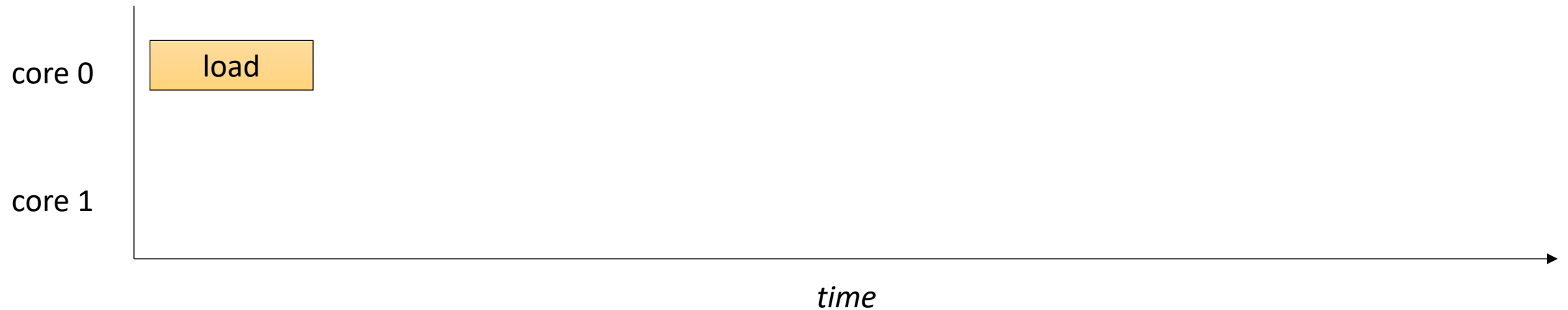
DAE: split into heterogeneous parallelism: one core does memory and one does computation



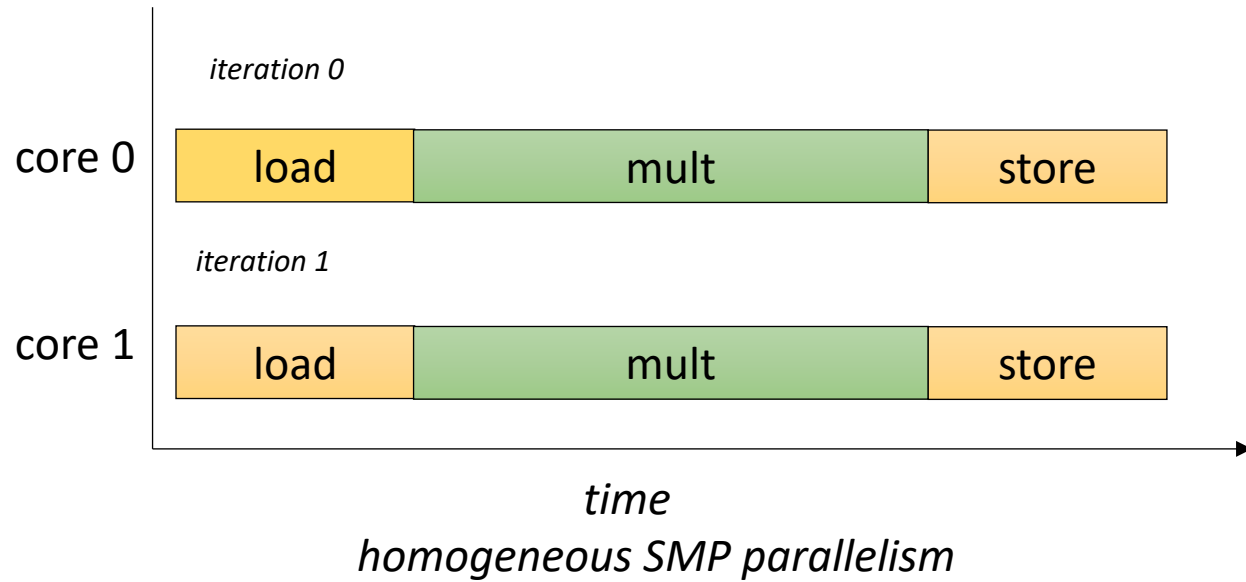
DAE - illustration



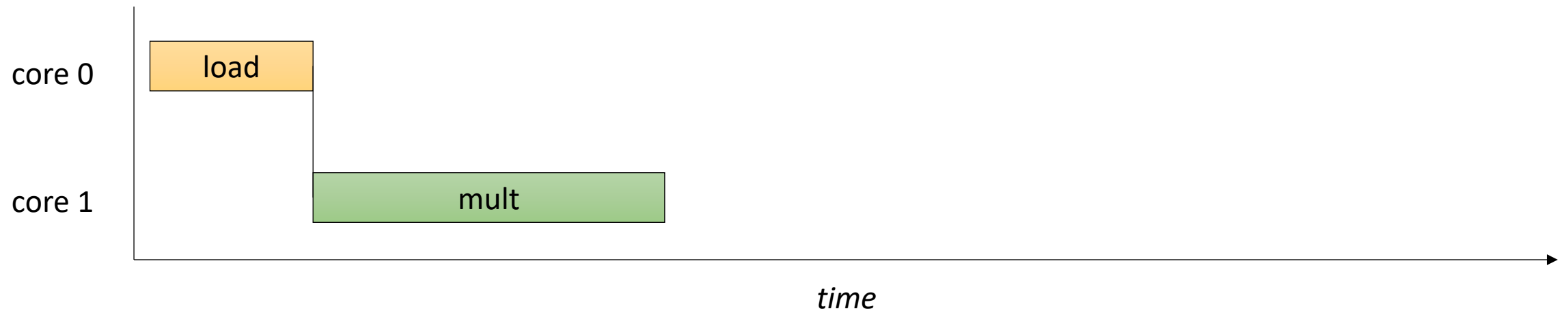
DAE: split into heterogeneous parallelism: one core does memory and one does computation



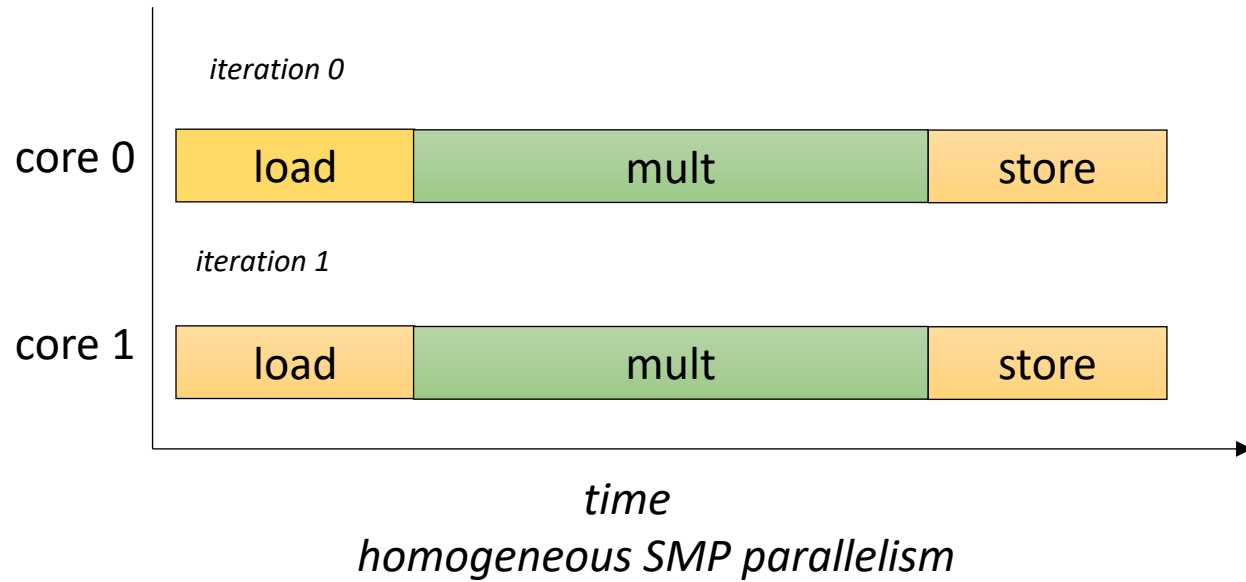
DAE - illustration



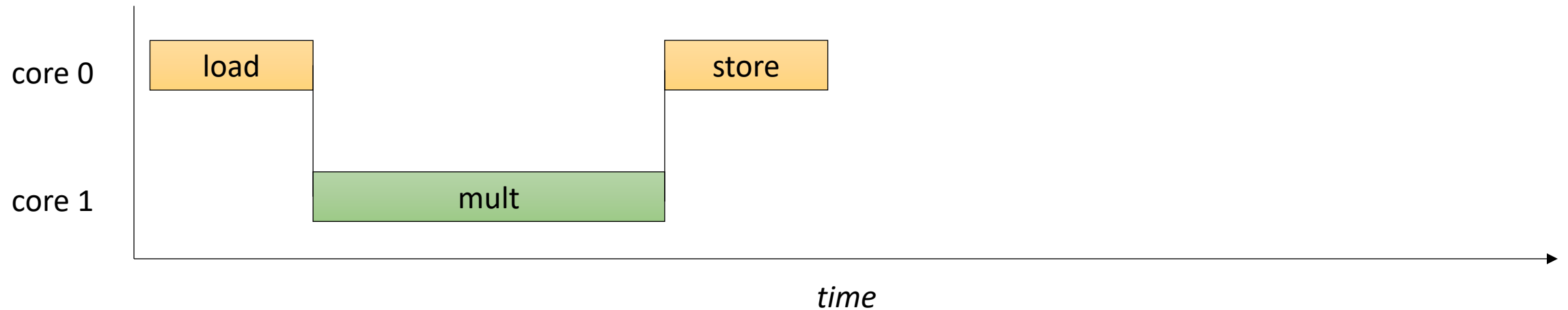
DAE: split into heterogeneous parallelism: one core does memory and one does computation



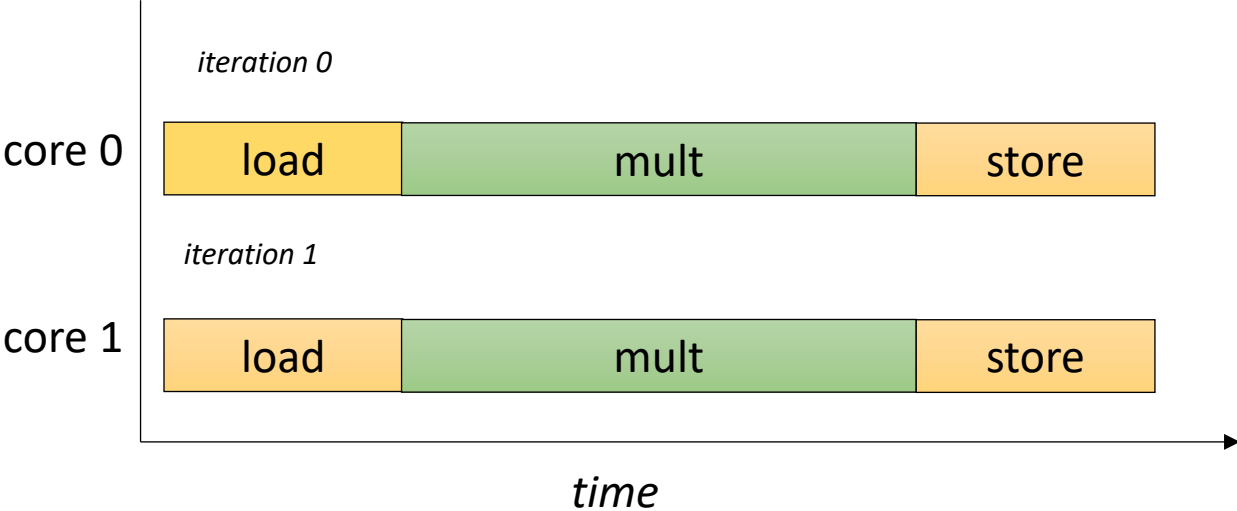
DAE - illustration



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

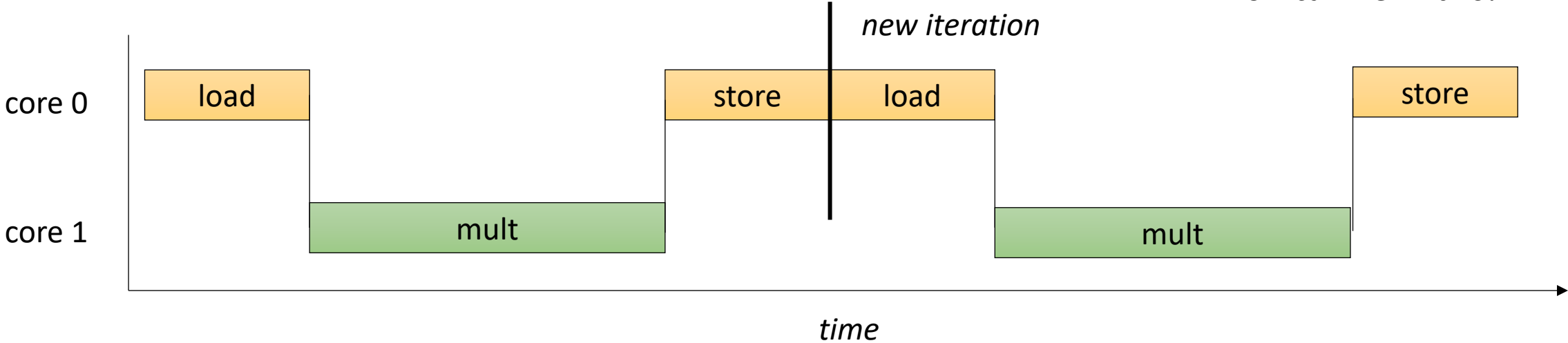


DAE - illustration



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

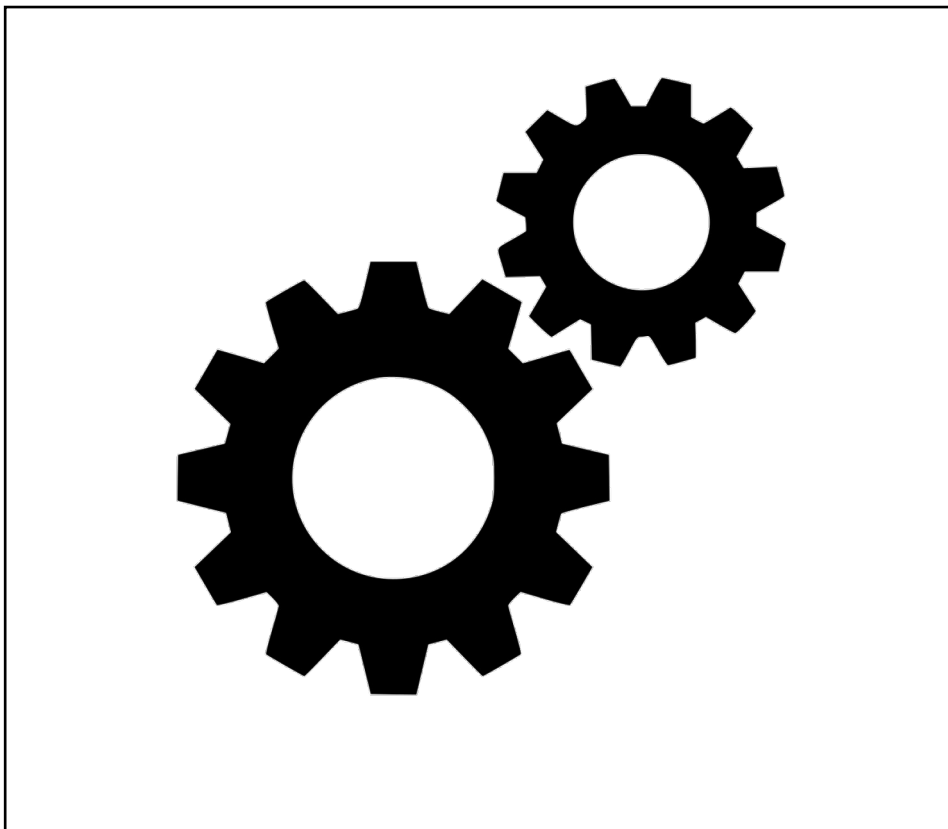
homogeneous SMP parallelism



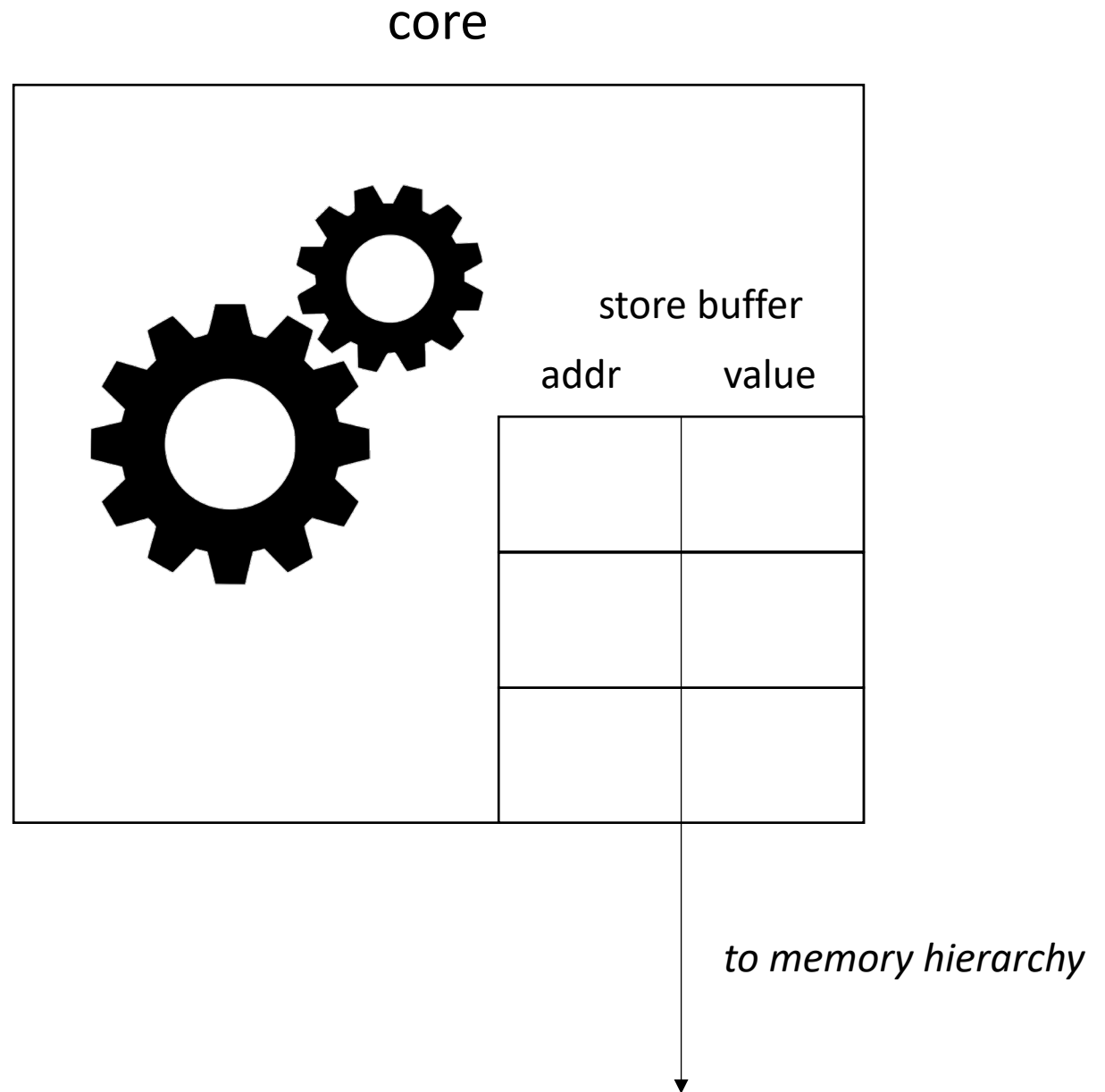
*This is sequentialized ☹️
How can we fix this?*

Store Buffers

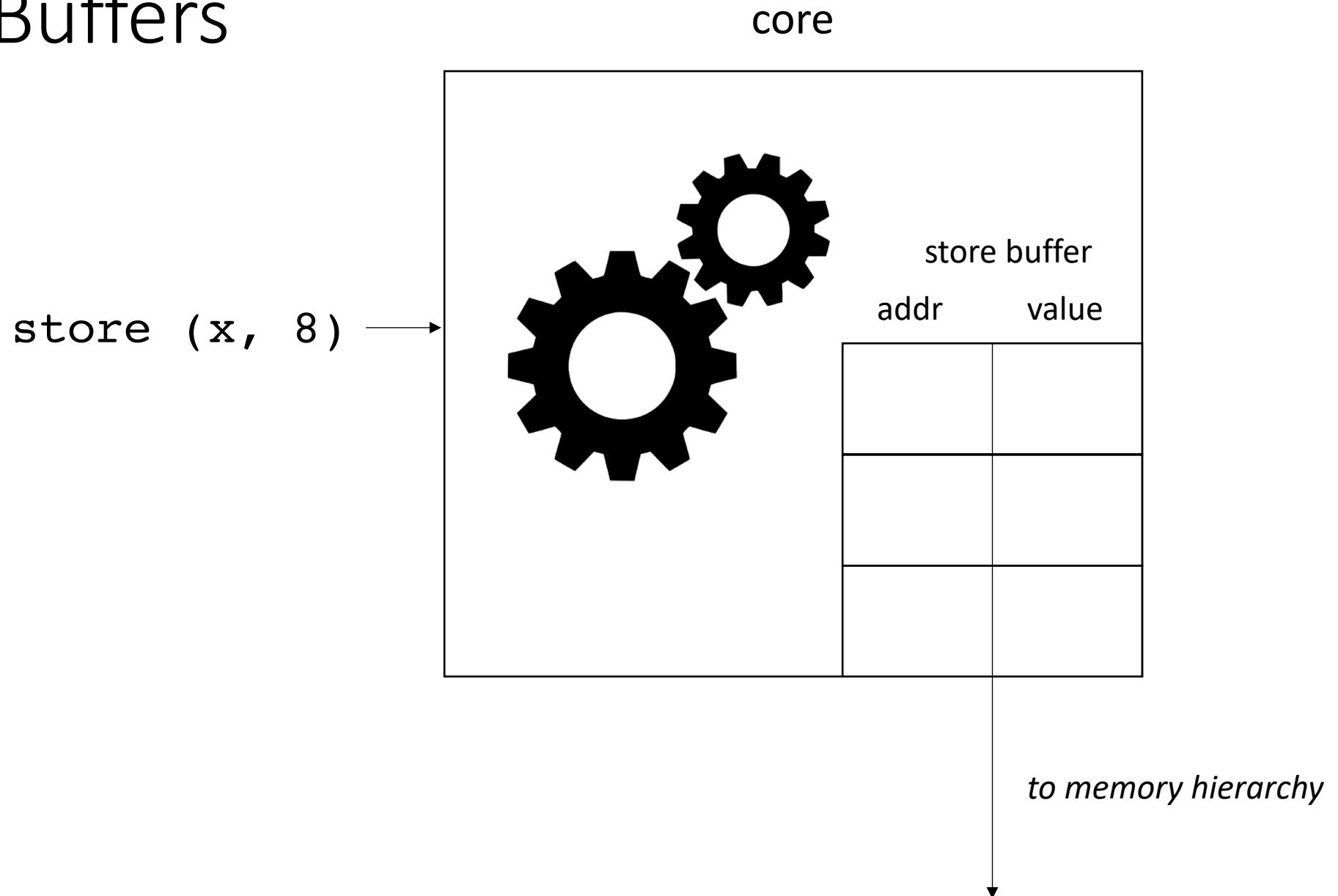
core



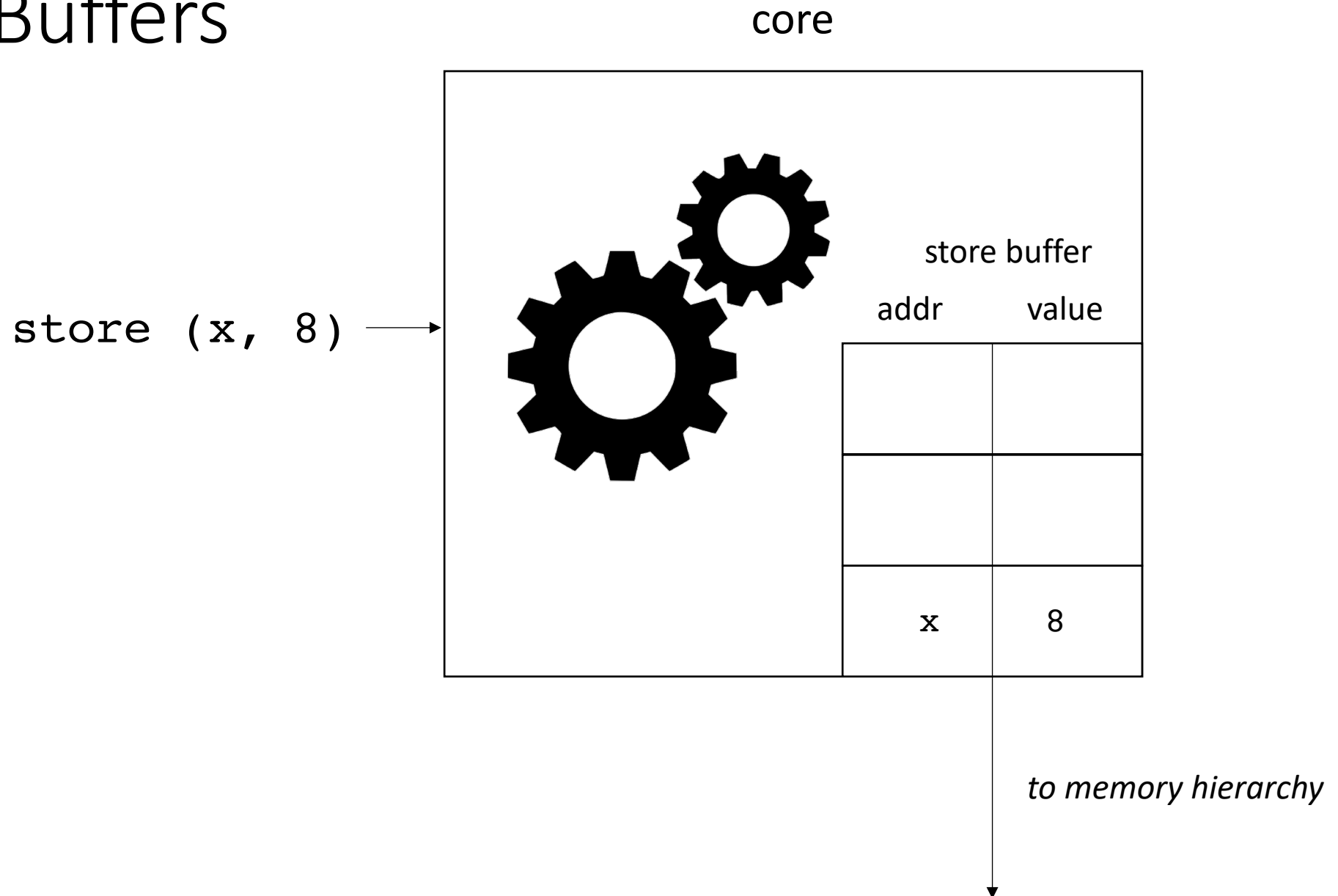
Store Buffers



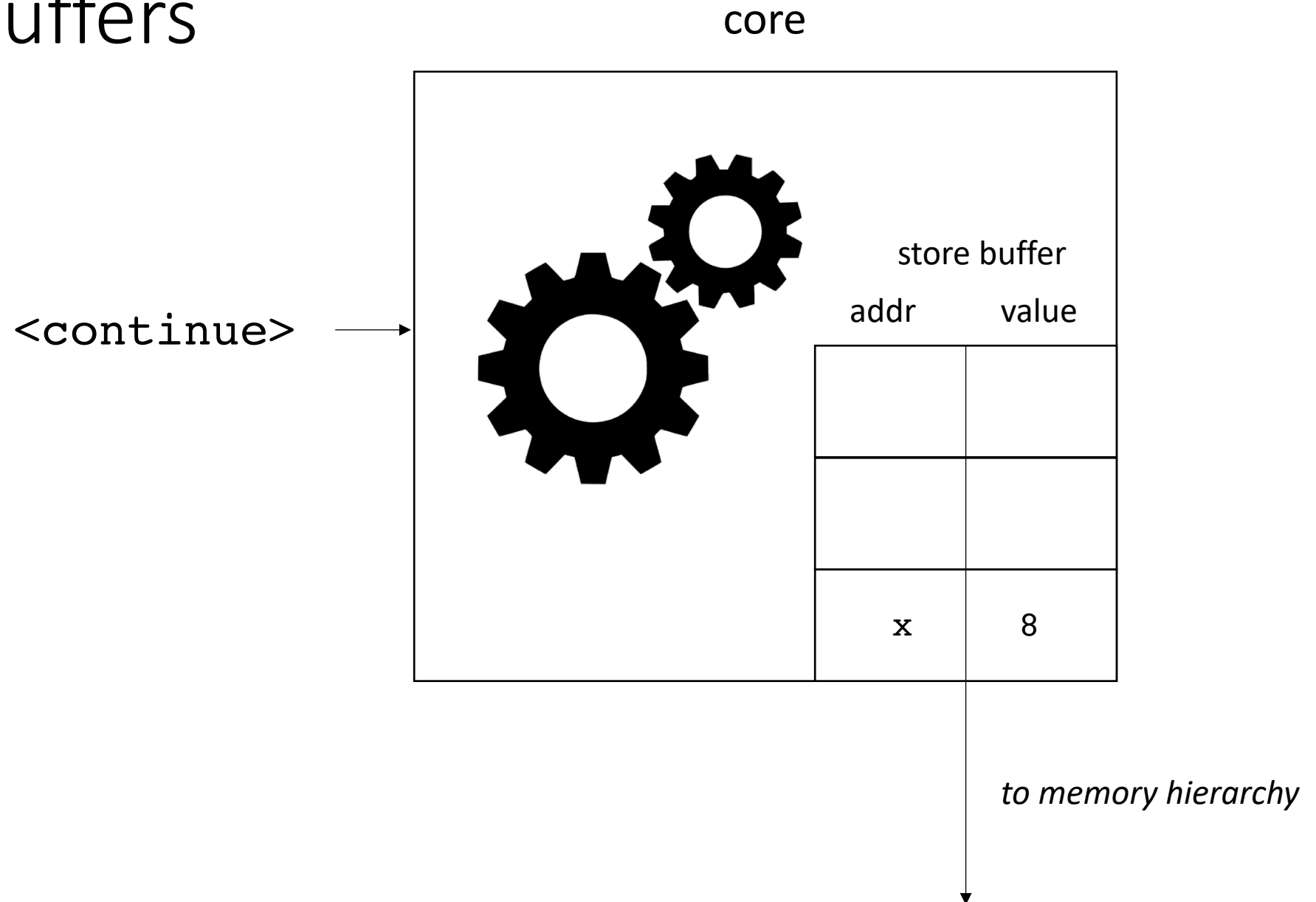
Store Buffers



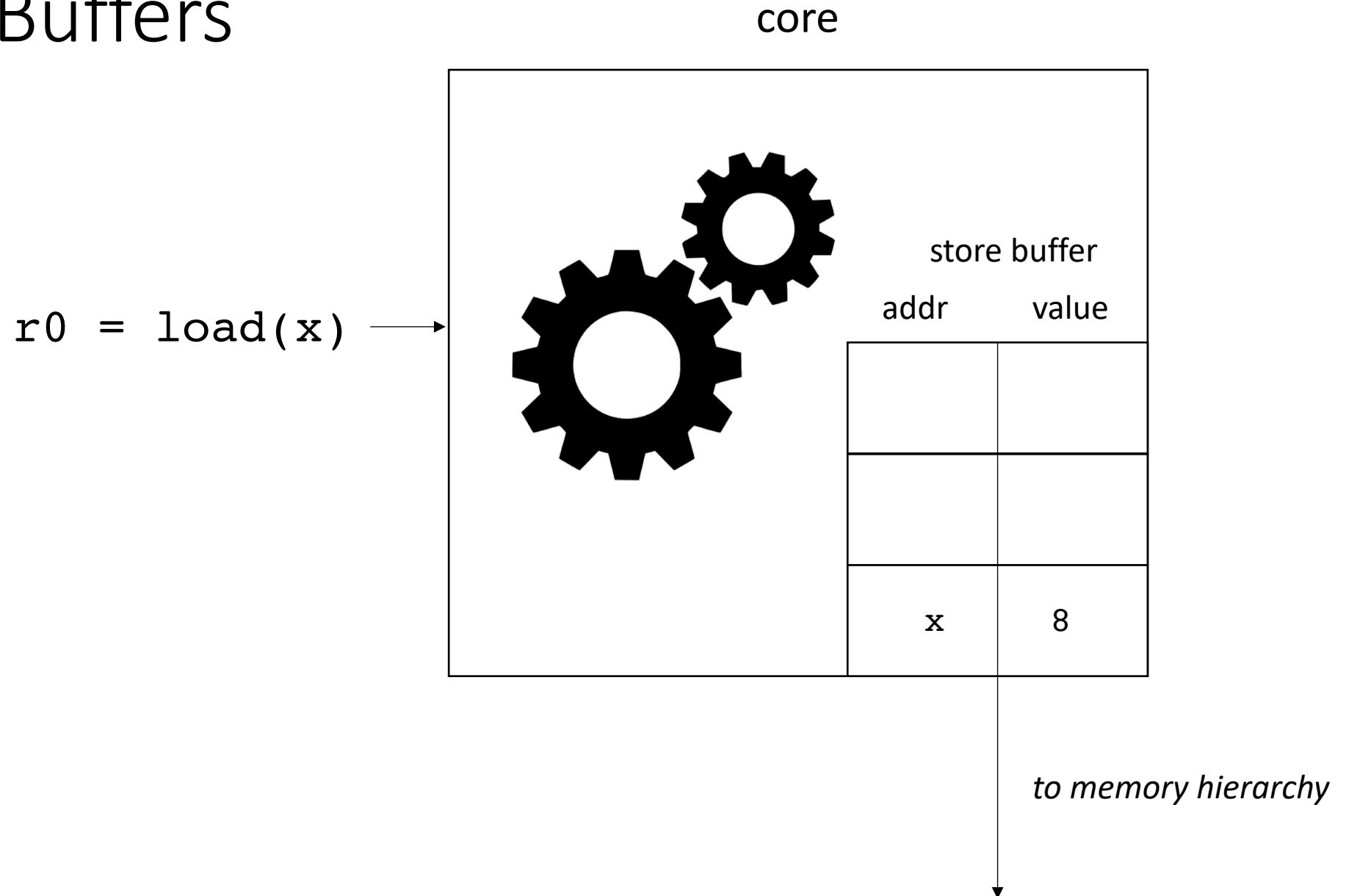
Store Buffers



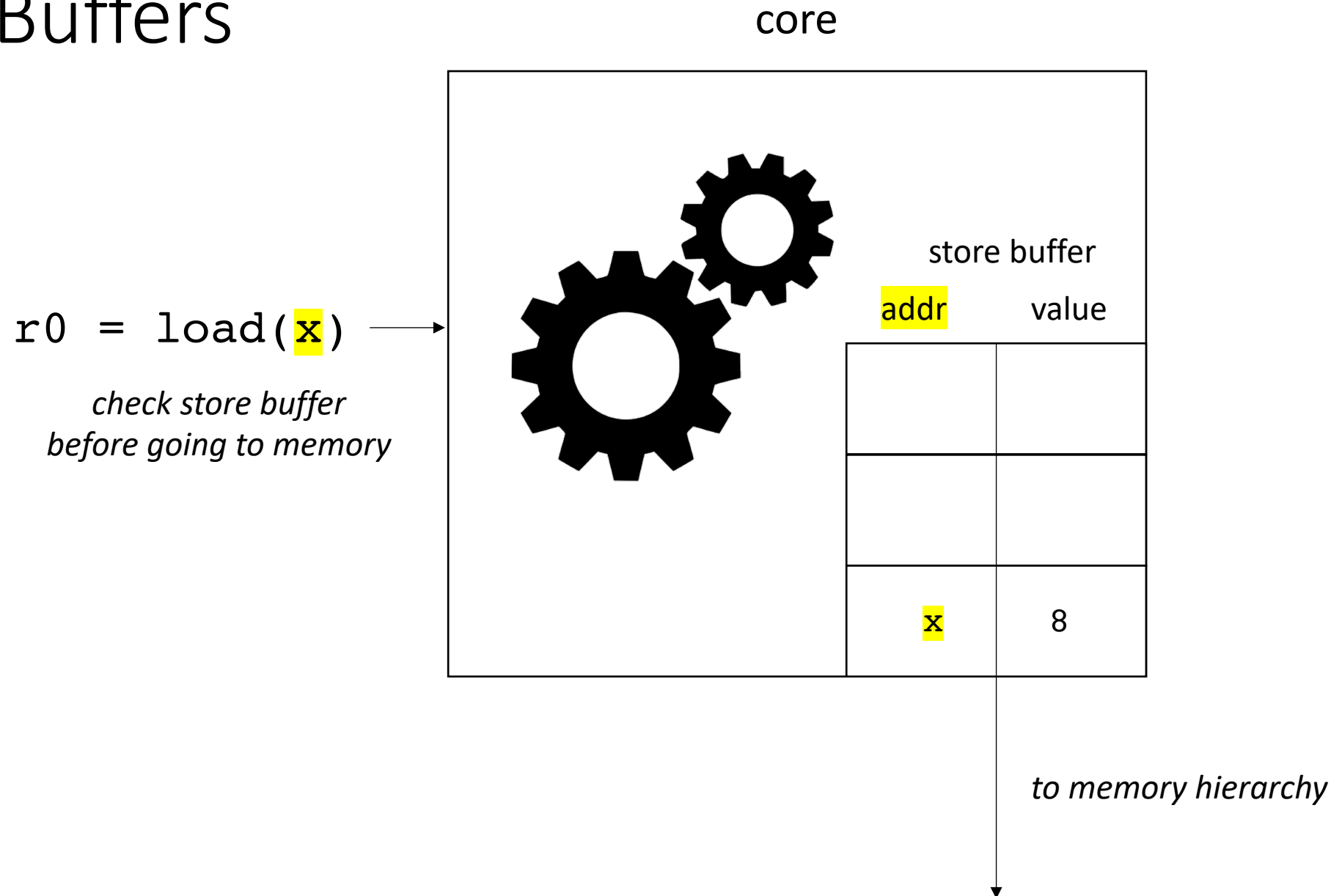
Store Buffers



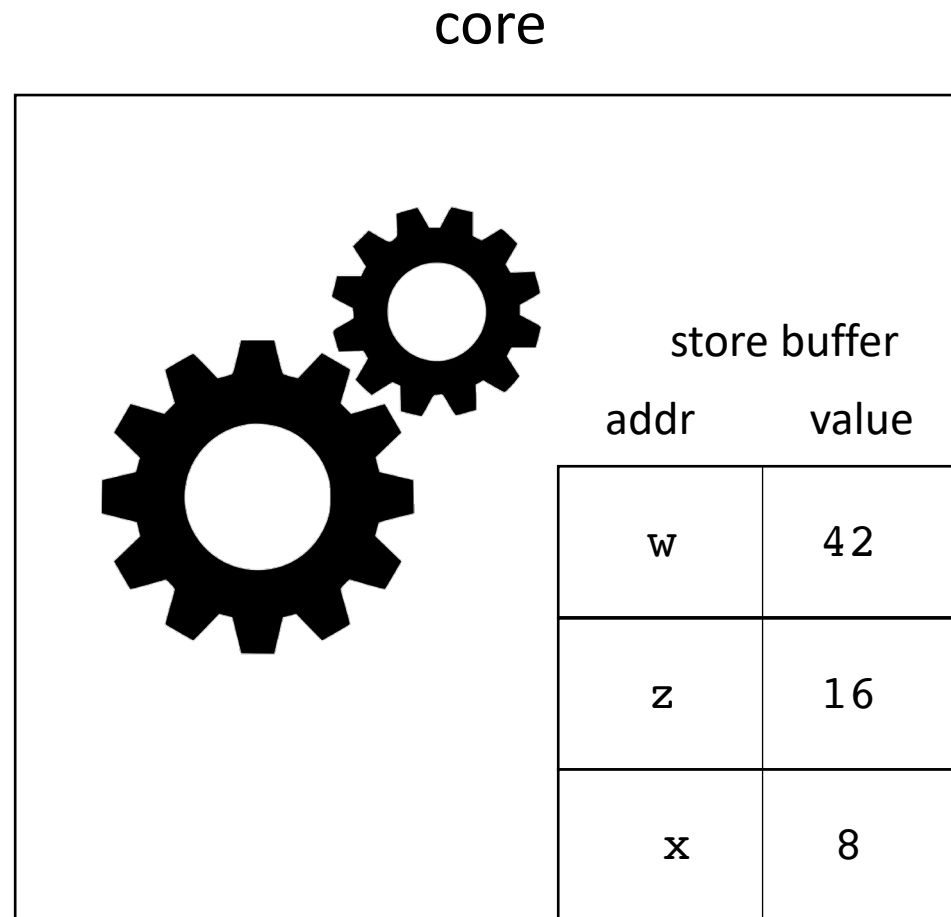
Store Buffers



Store Buffers

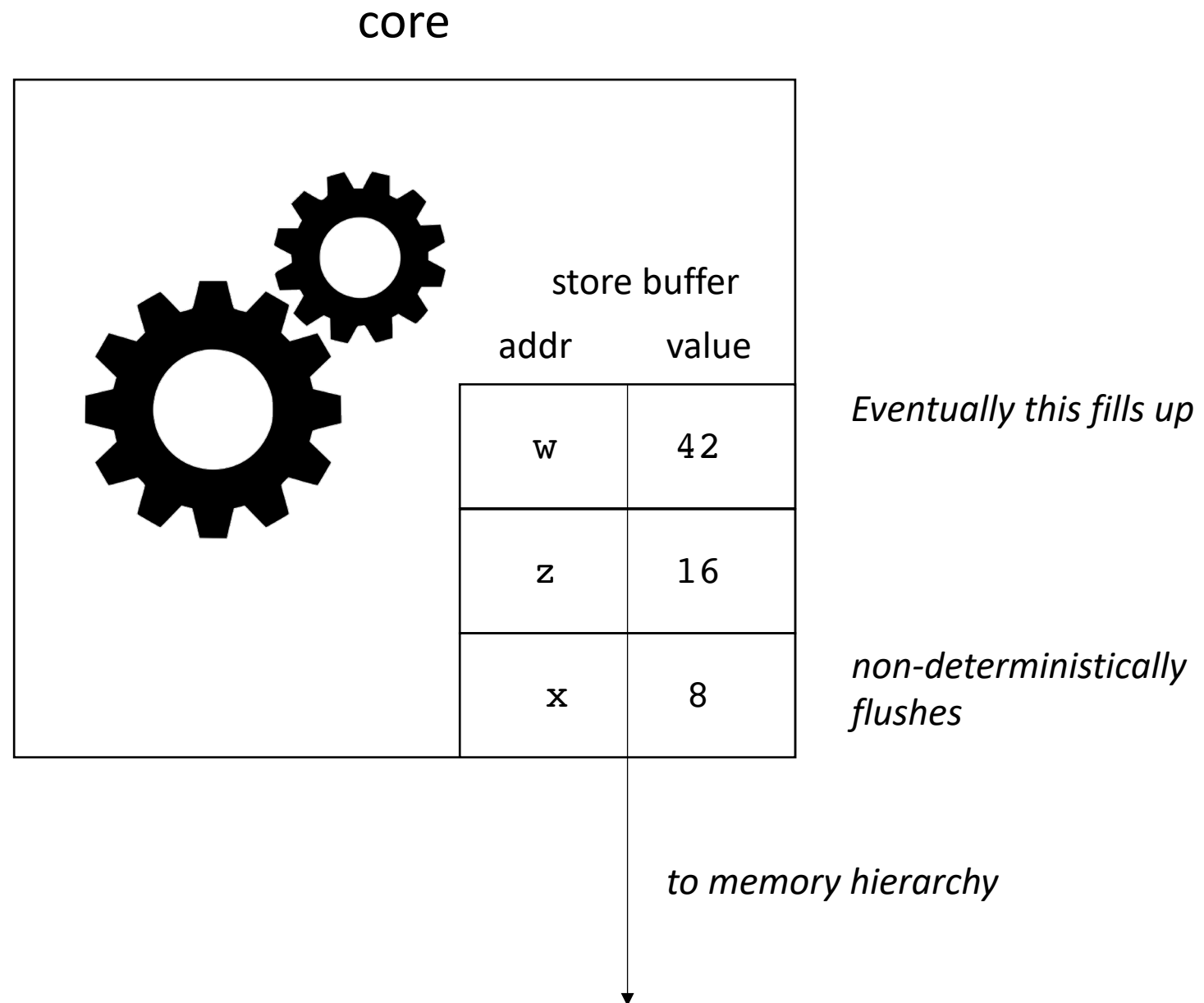


Store Buffers



Eventually this fills up

Store Buffers



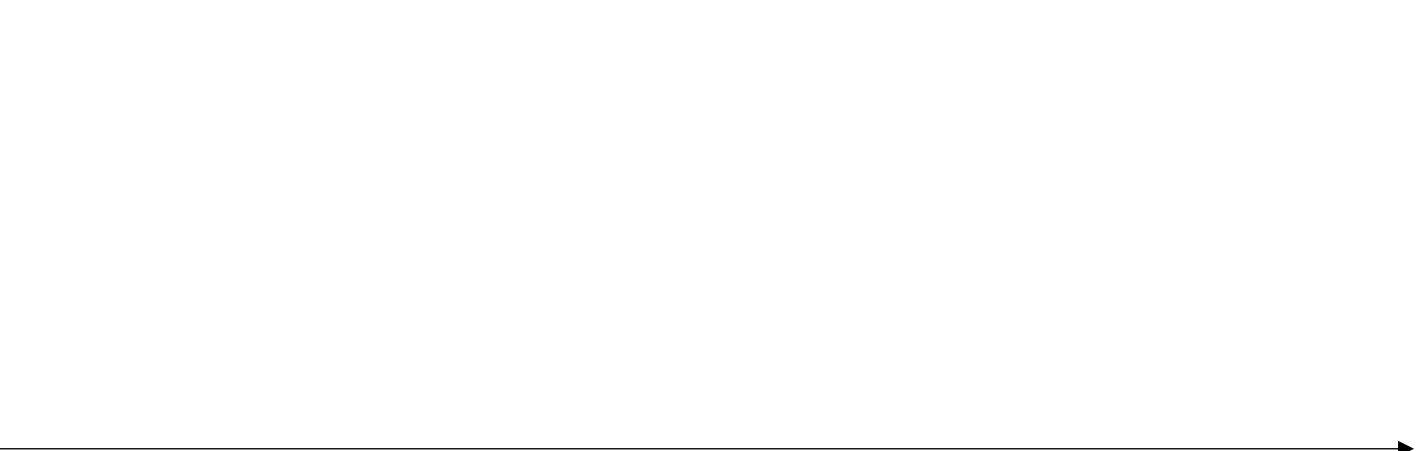
DAE Parallelism

store buffer

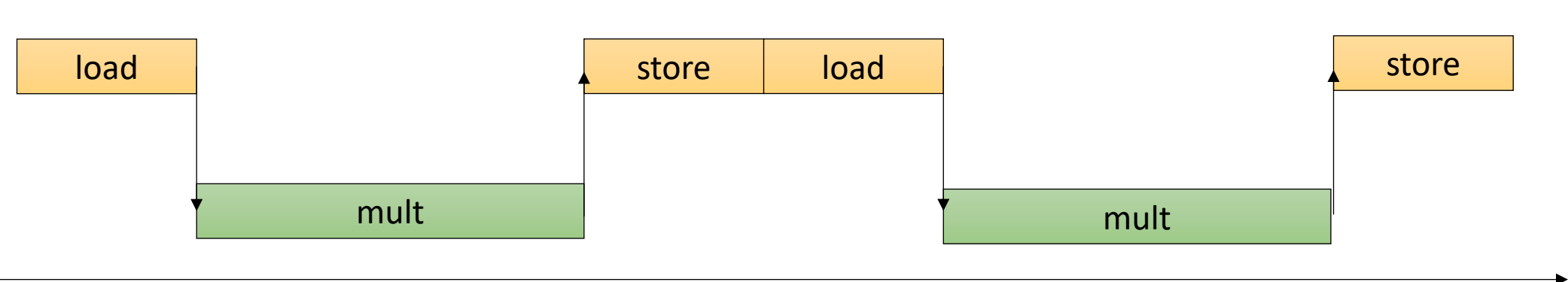
addr value


```
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);  
}
```



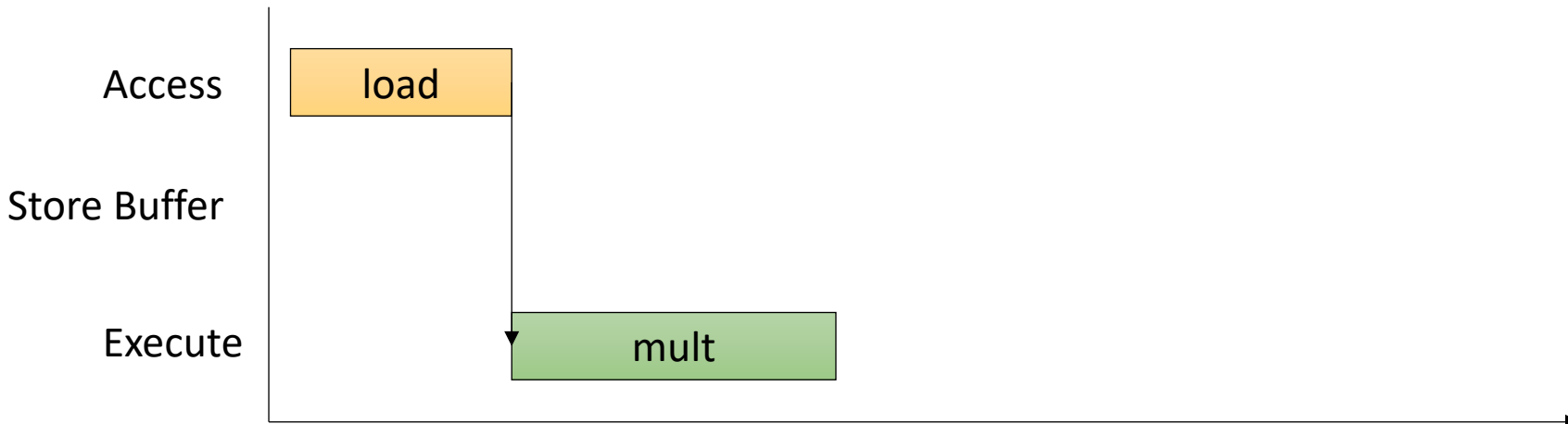
Store Buffer

Execute



time

DAE Parallelism

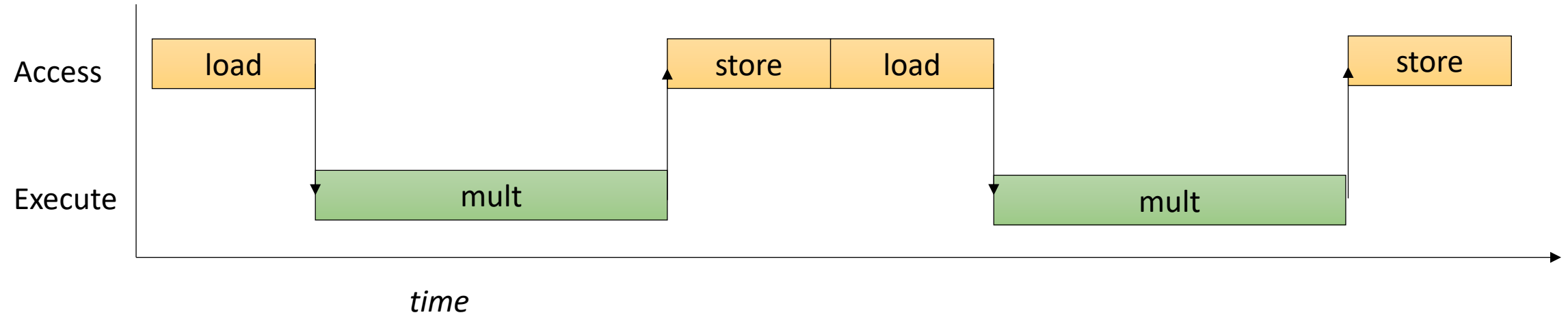


store buffer

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);  
}
```



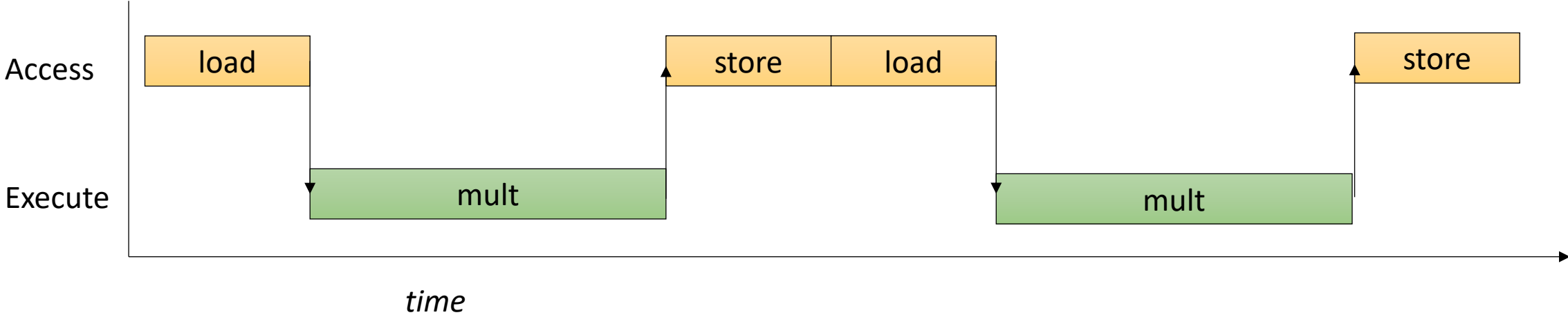
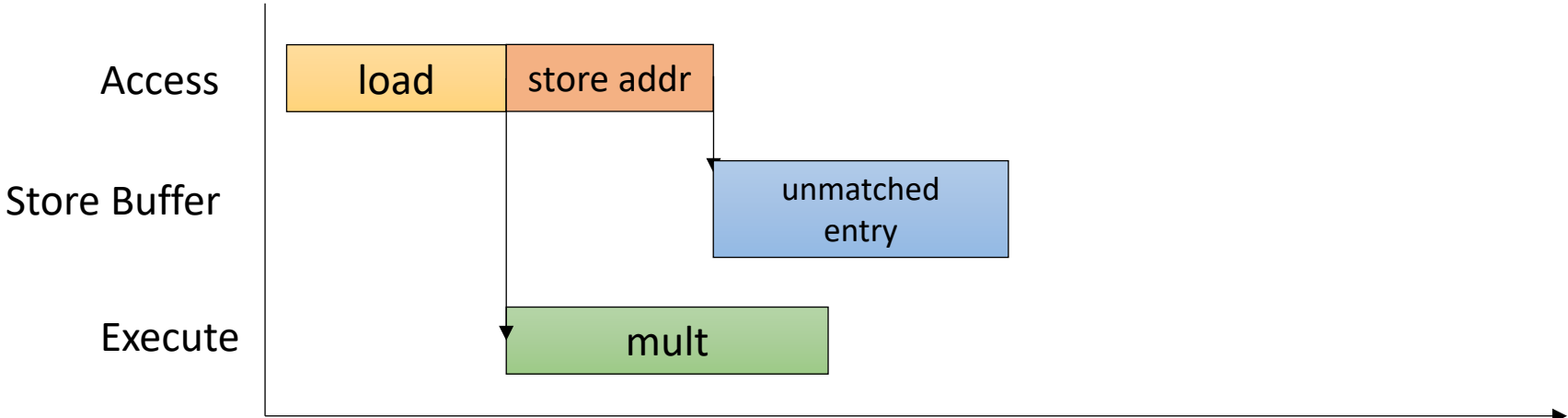
DAE Parallelism

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);
}
    
```



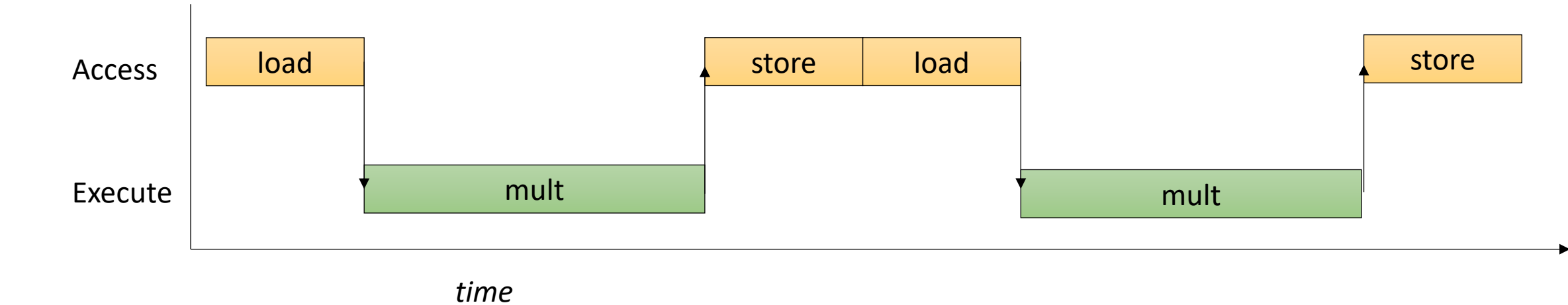
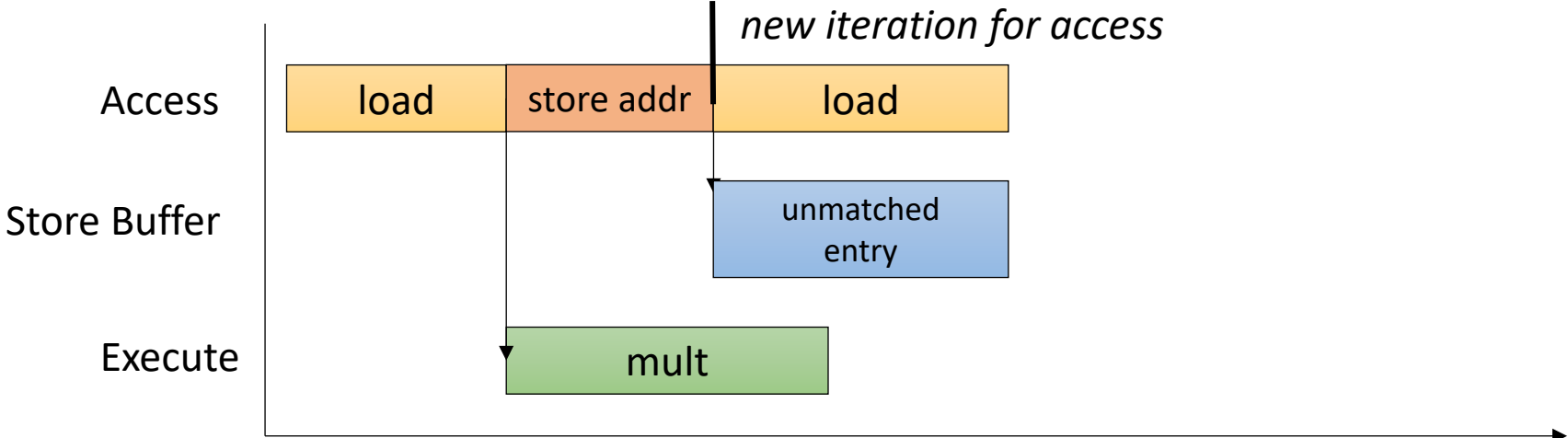
DAE Parallelism

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

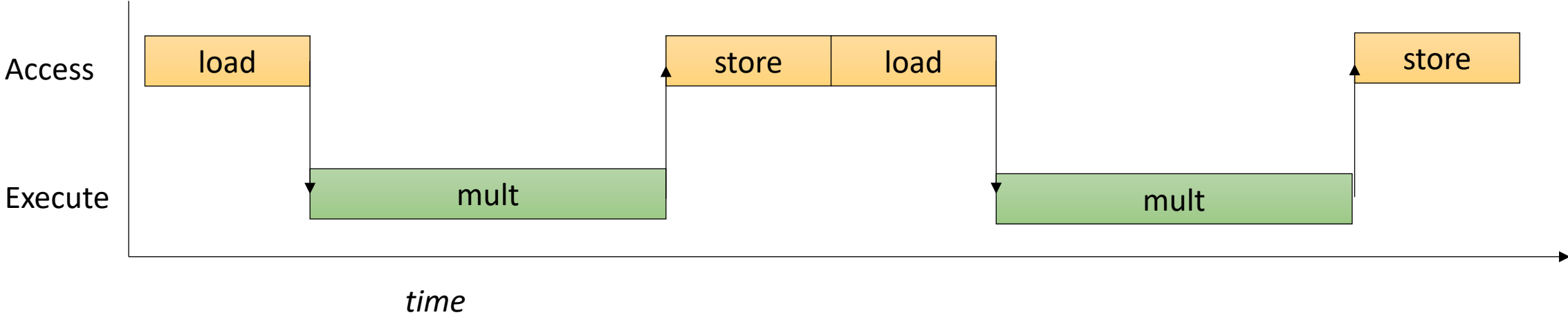
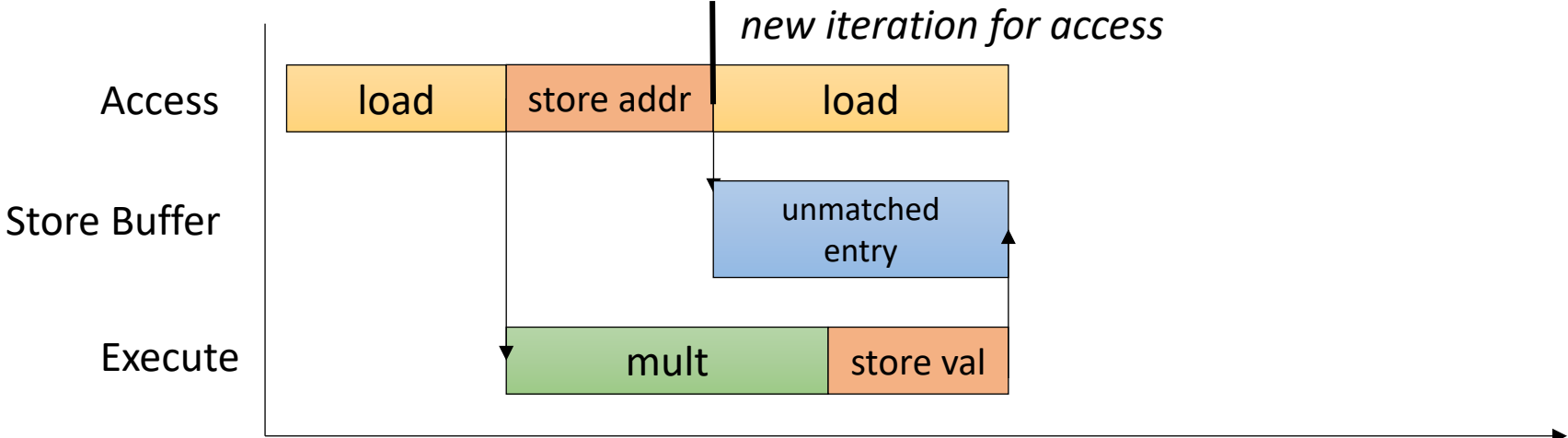
store buffer

addr value

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

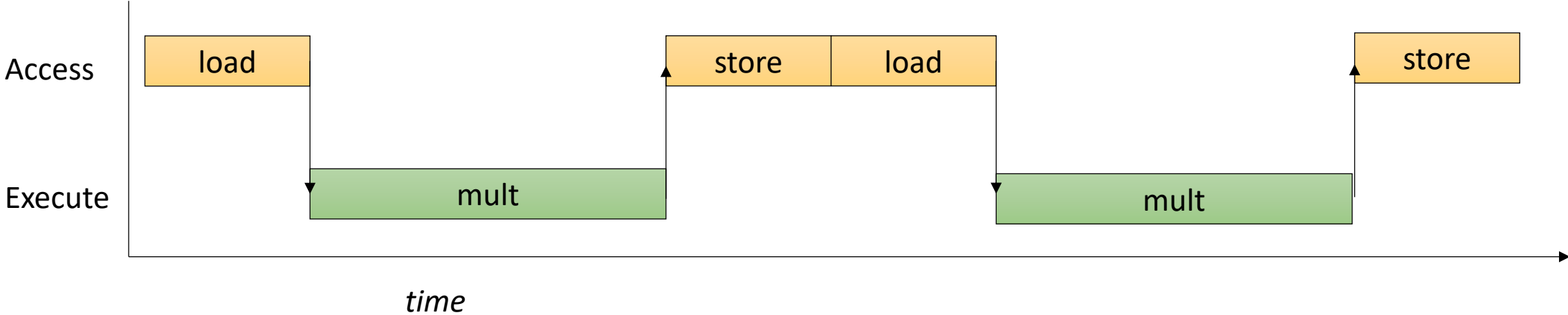
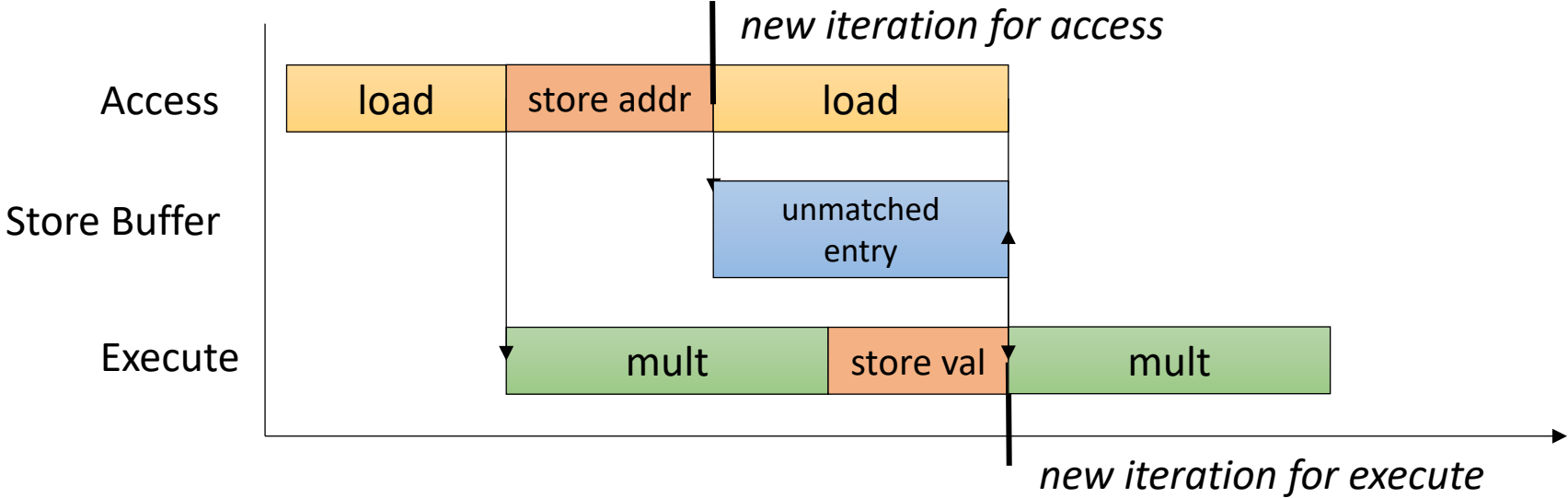
store buffer

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

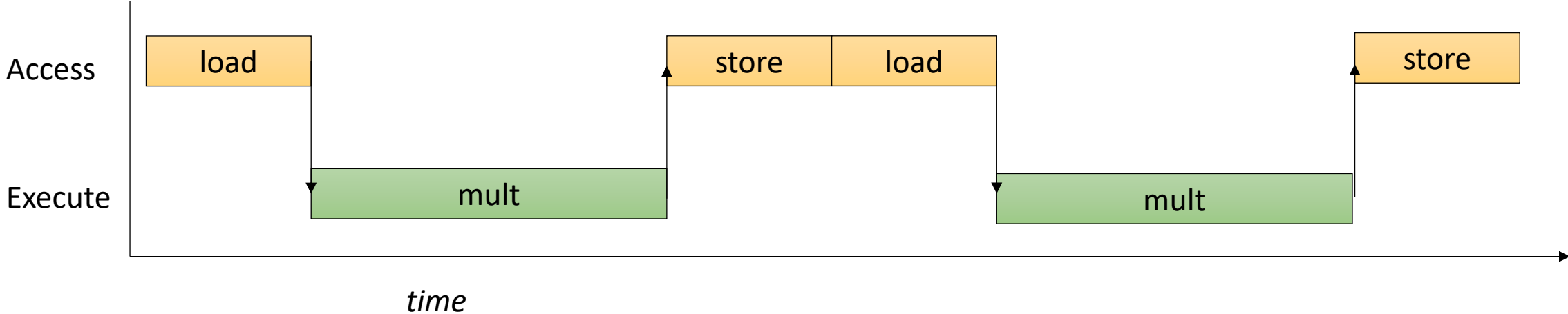
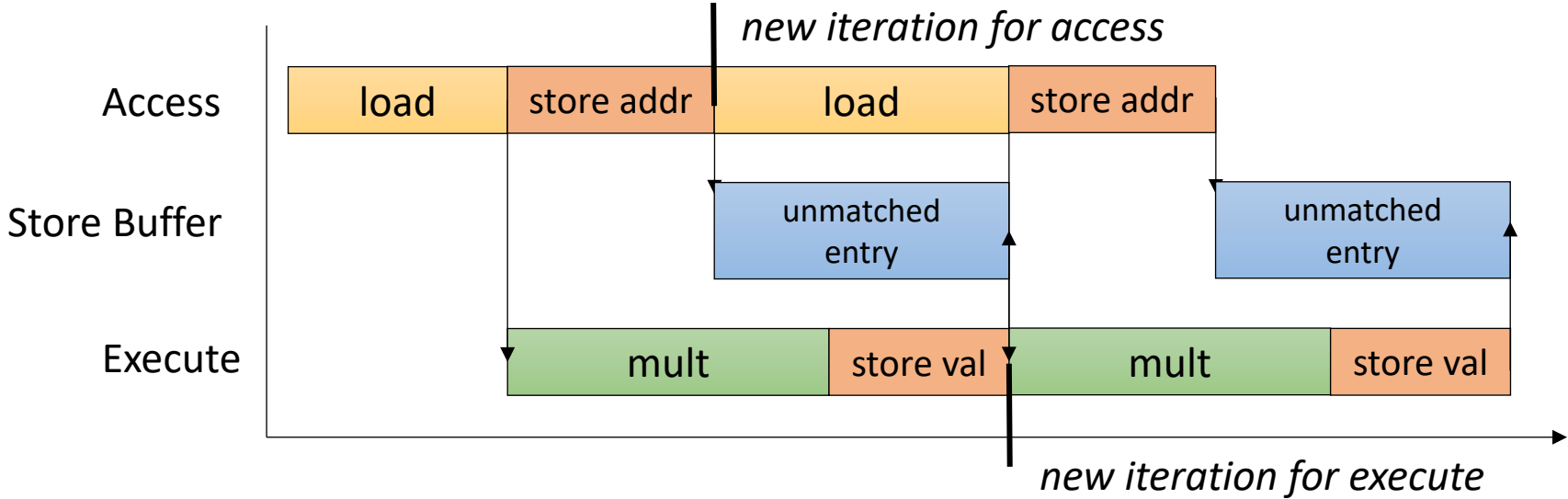
store buffer

addr value

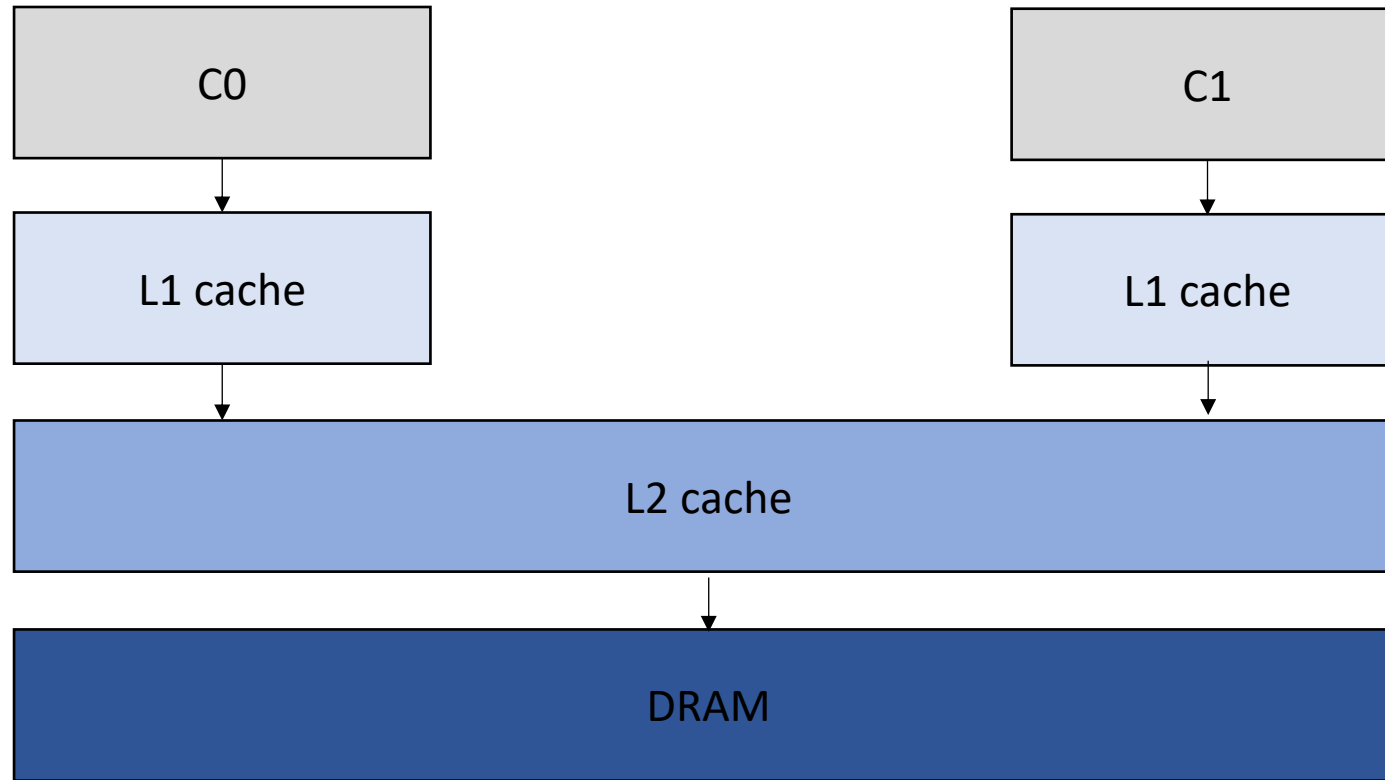
has 2 entries now

```

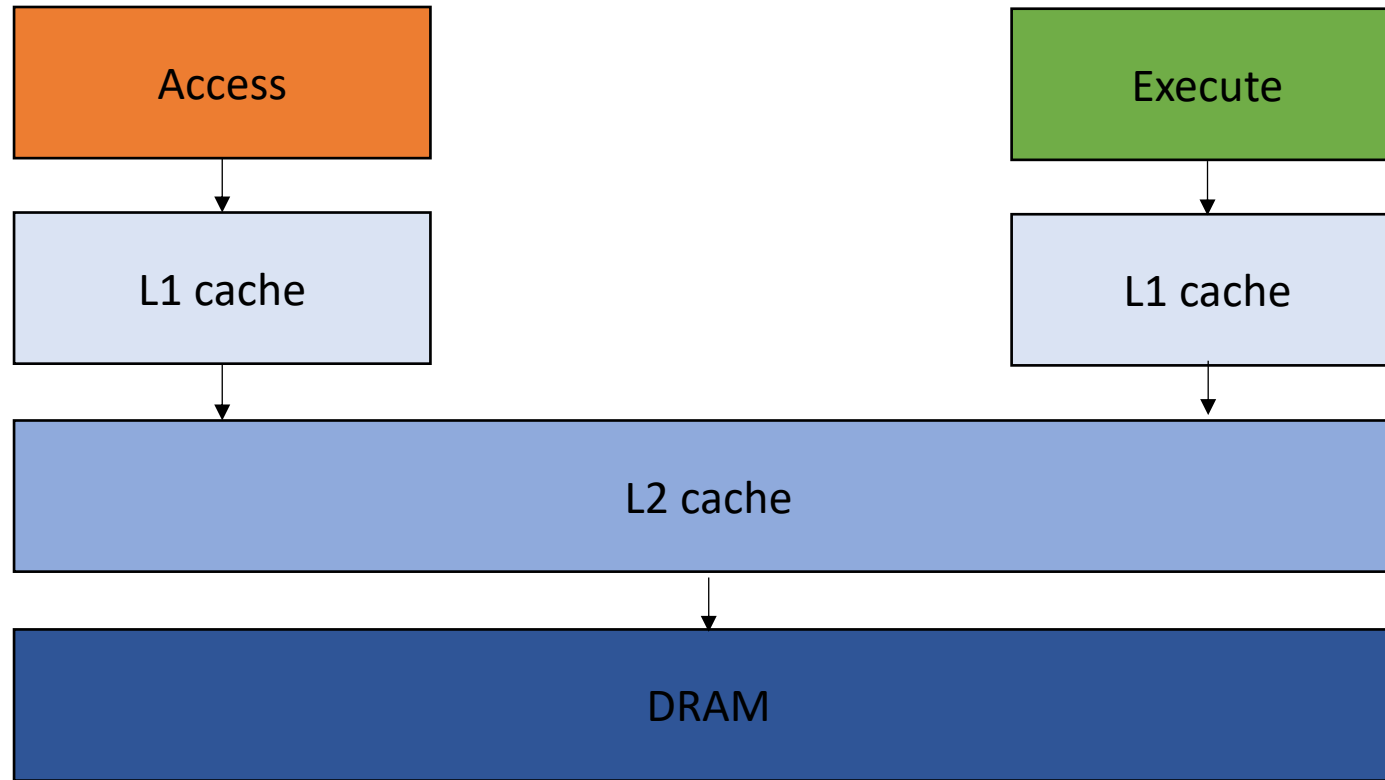
for (int i = 0; i < SIZE; i++) {
    float r0 = load(b + i);
    float r1 = r0 * 3.14;
    store(a + i, r1);
}
    
```



Specializing a DAE architecture

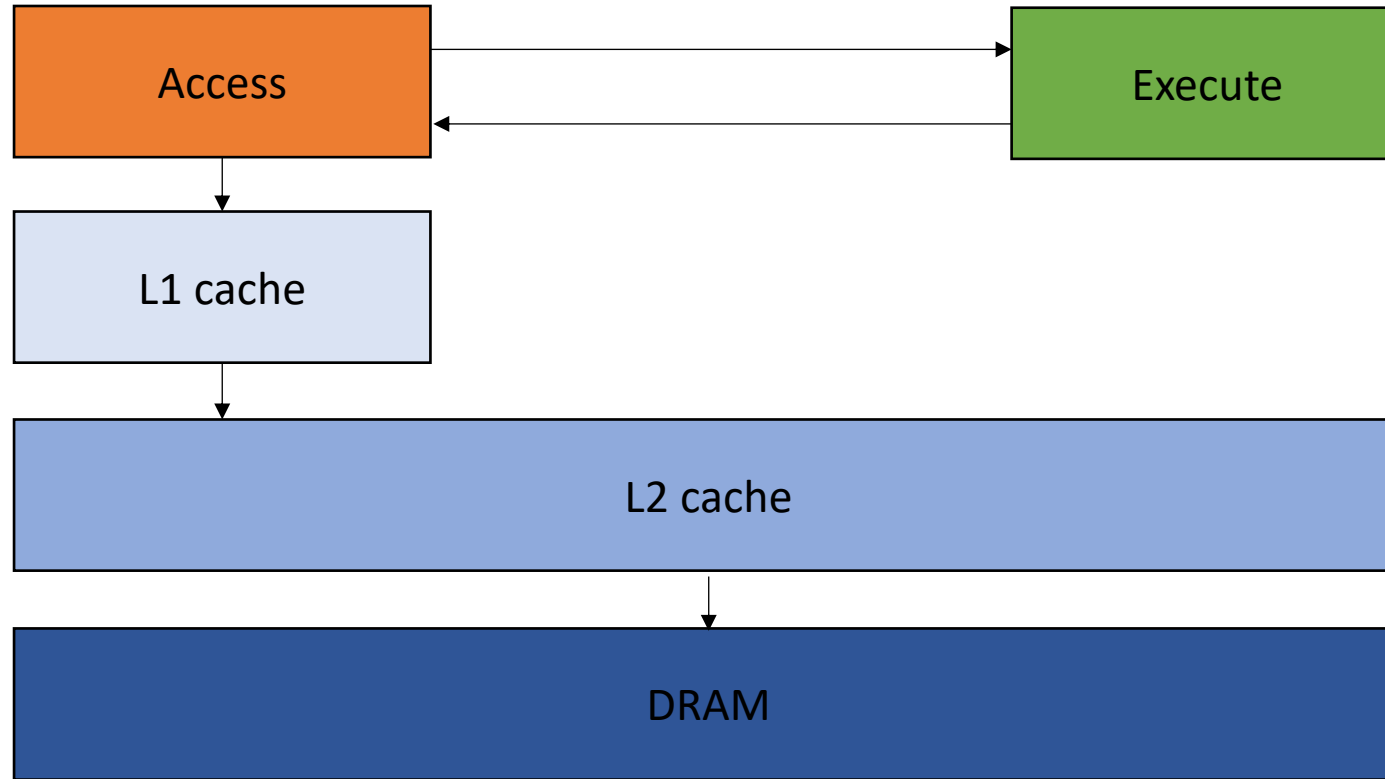


Specializing a DAE architecture



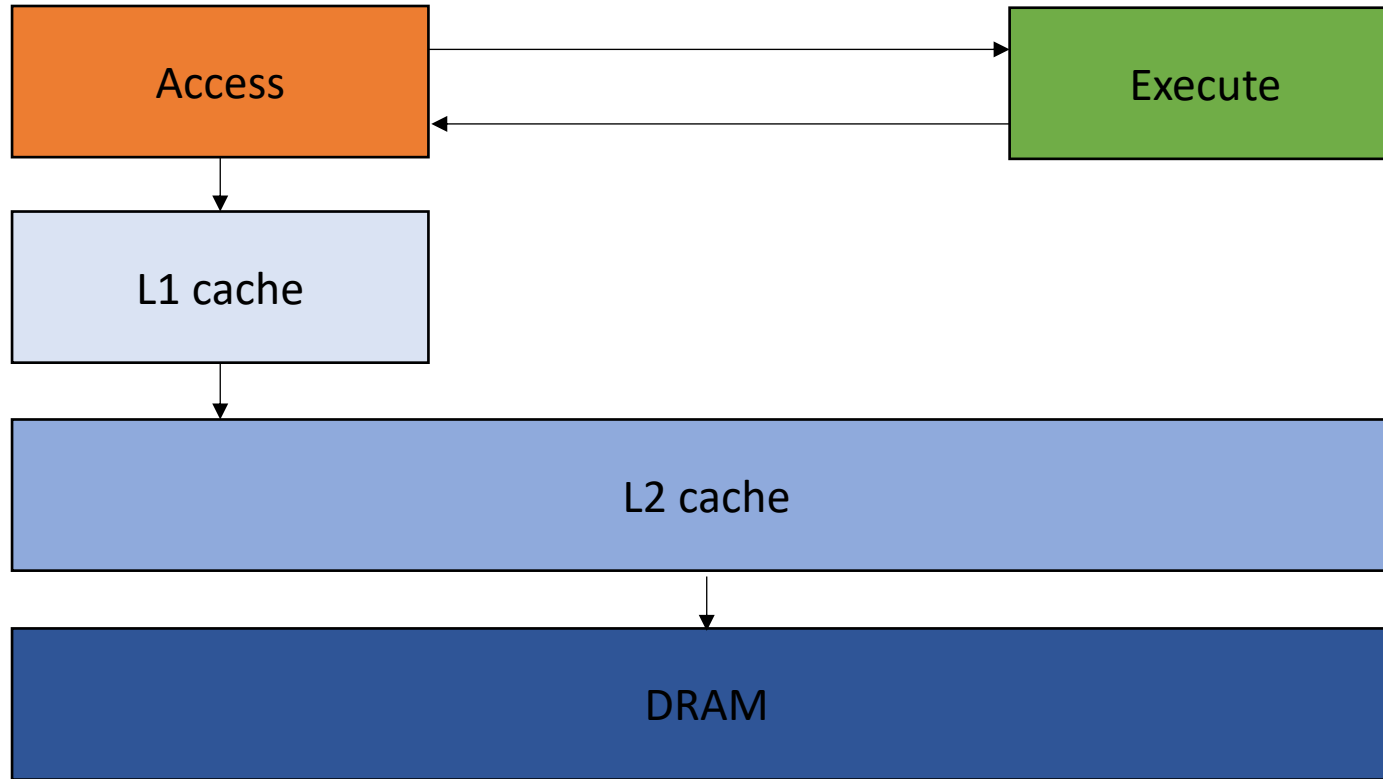
Specializing a DAE architecture

*Less contention
on memory hierarchy*



Specializing a DAE architecture

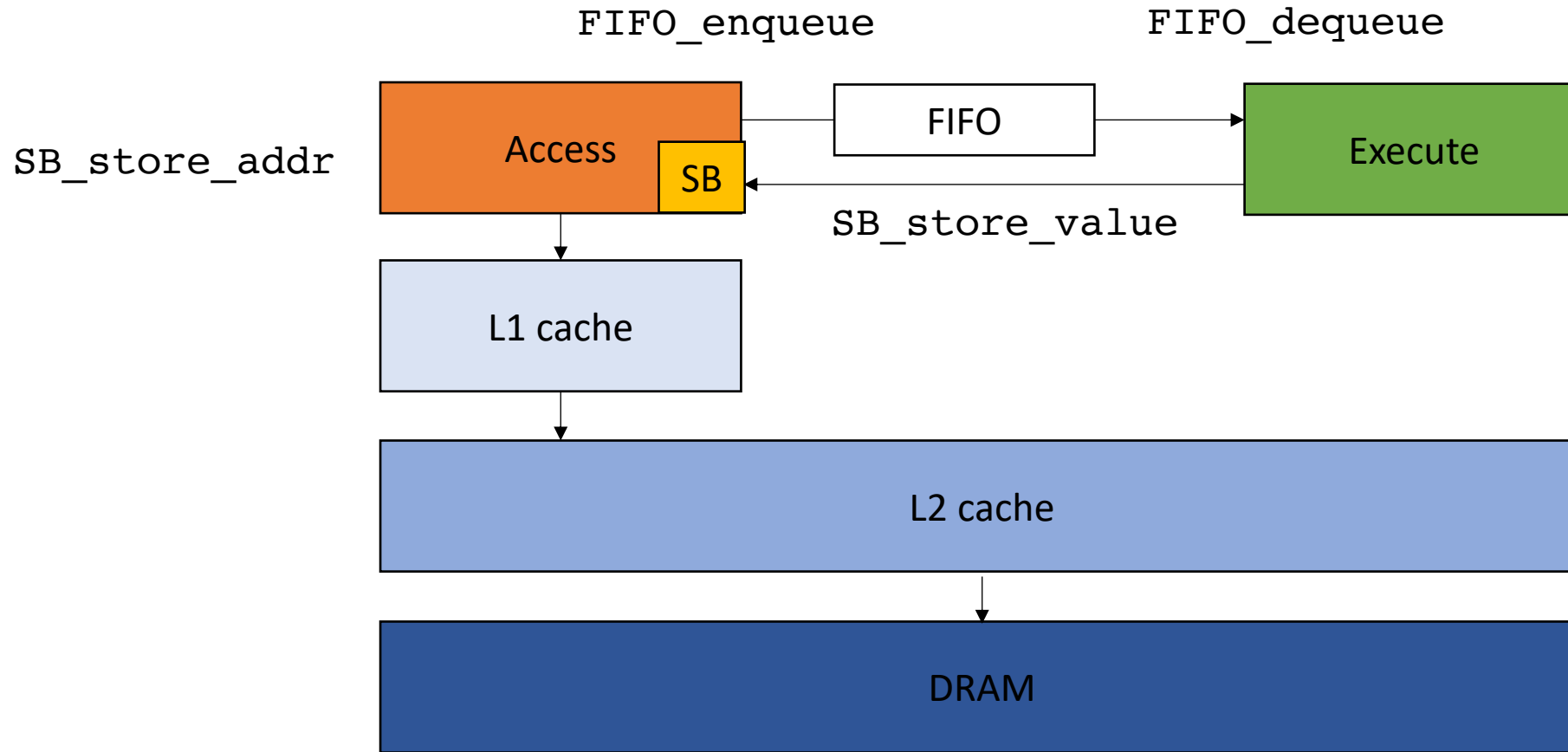
optimizations?
FP unit
Vector units



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 - 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)"

]

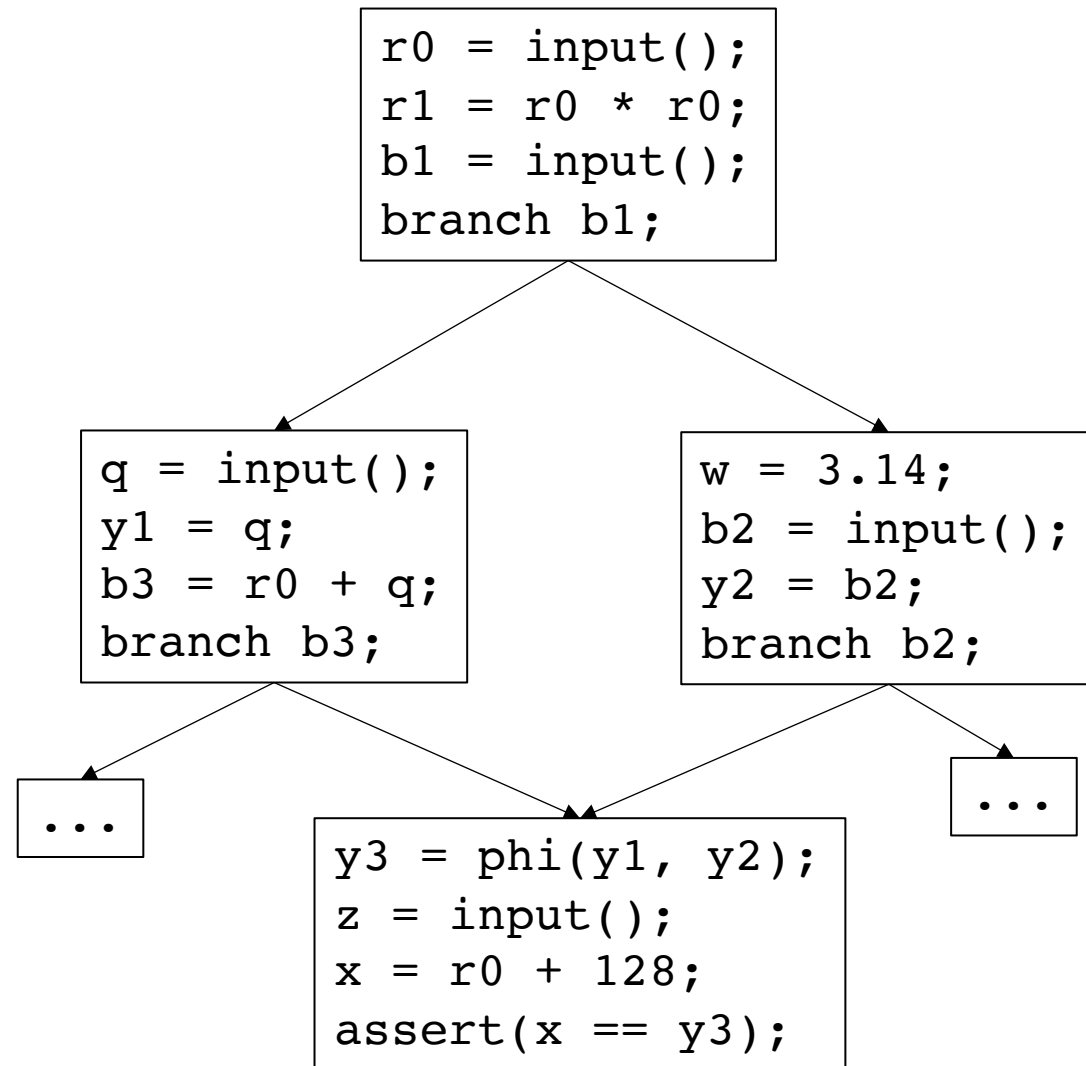
start with the statement and work backwards until there are no more dependencies

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

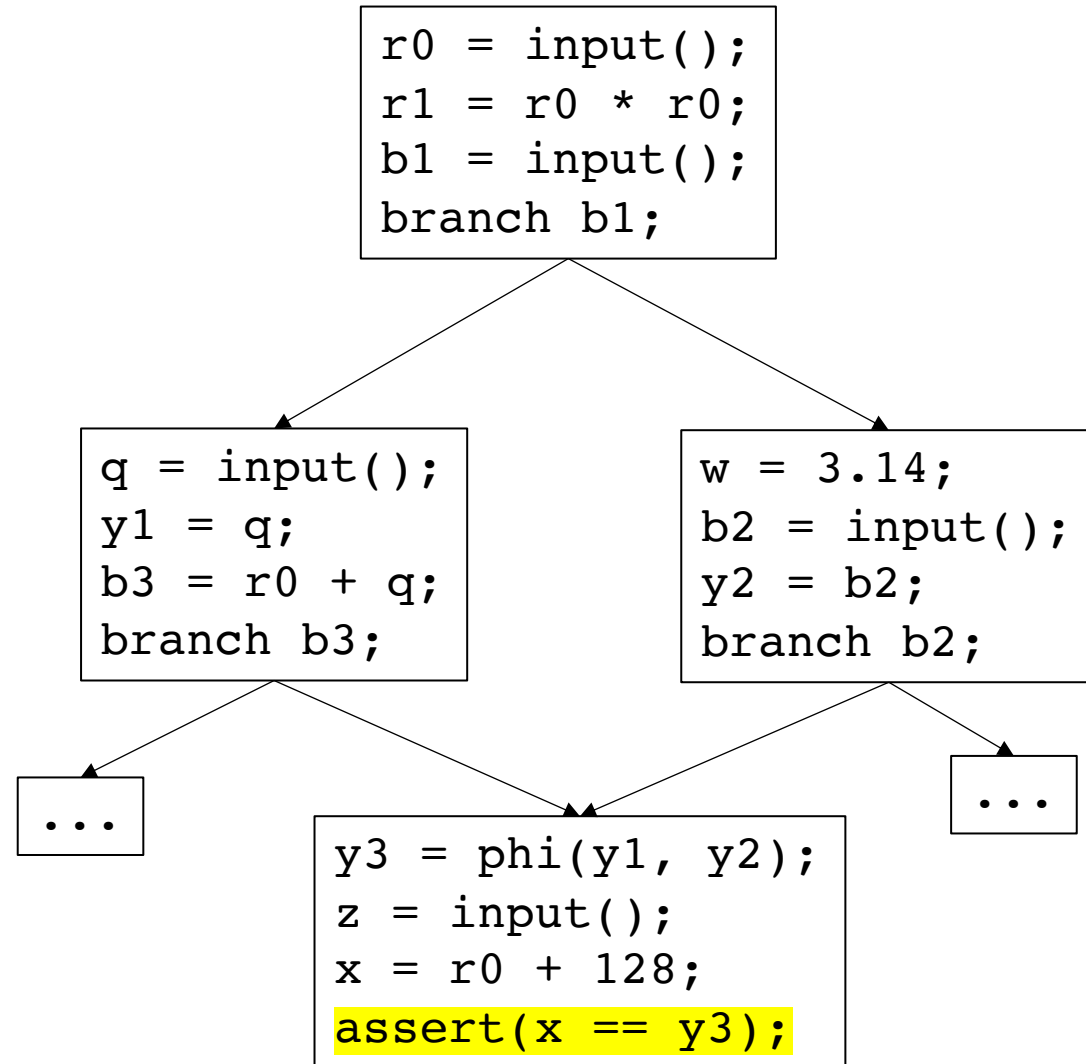
```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
    stmt = Worklist.pop();
    if (is_marked(stmt)) {
        continue;
    }
    mark(stmt);
    for a in stmt.args() {
        worklist.append(a);
    }
    for p in cfg[stmt].predecessors() {
        worklist.append(p.branch_stmt());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



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());
    }
}
```

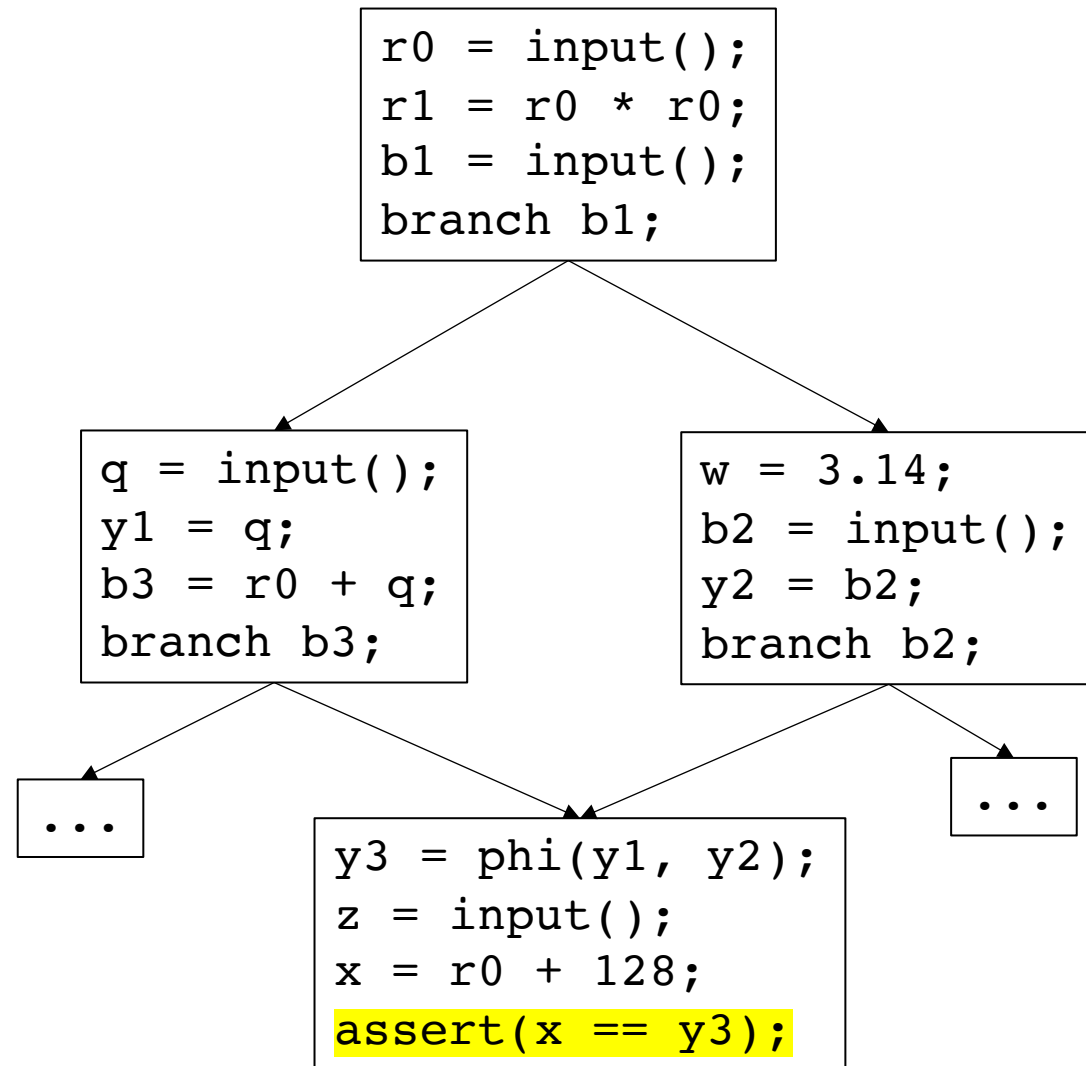
marked: worklist:
 assert()



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());
    }
}
```

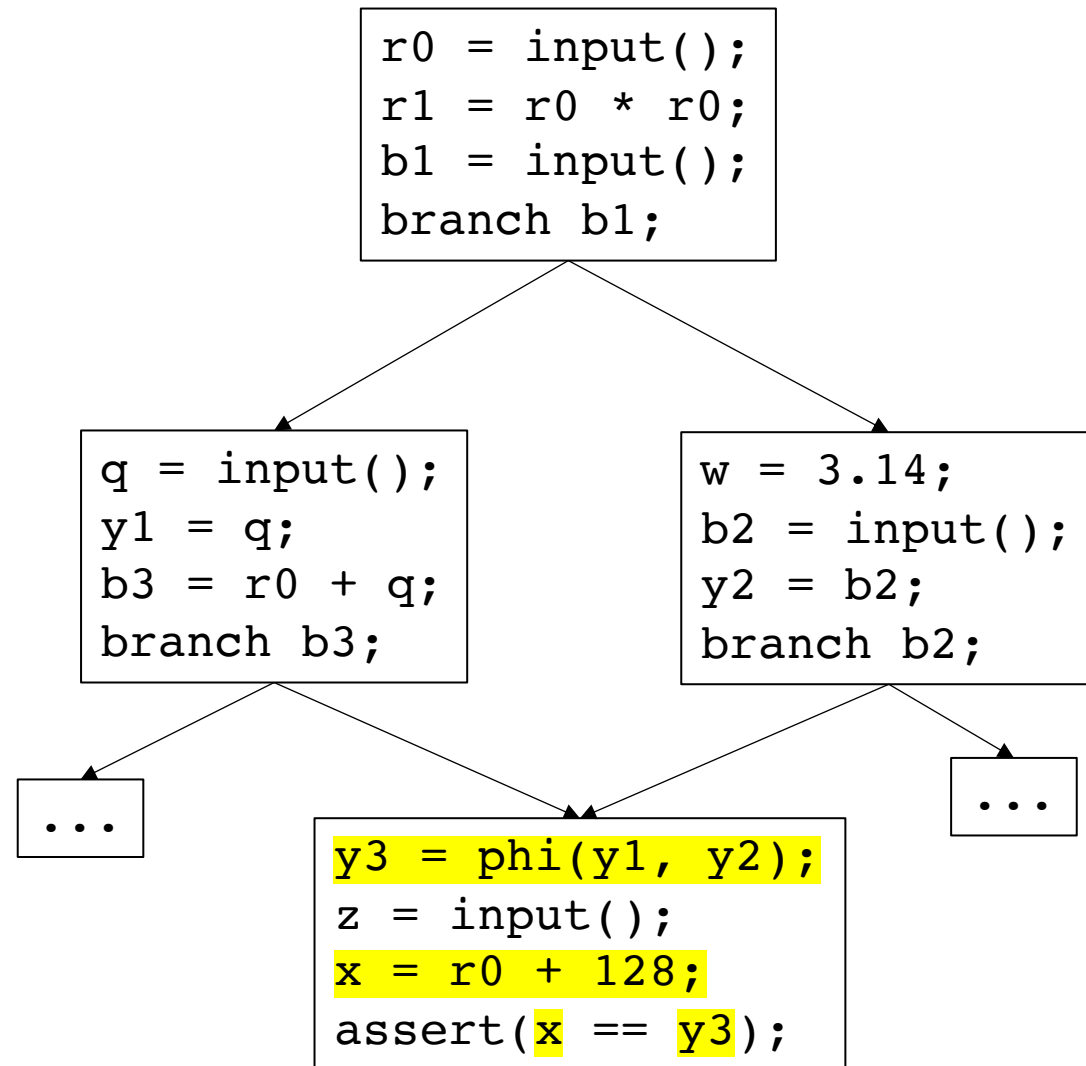
marked: worklist:
assert()



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());
    }
}
```

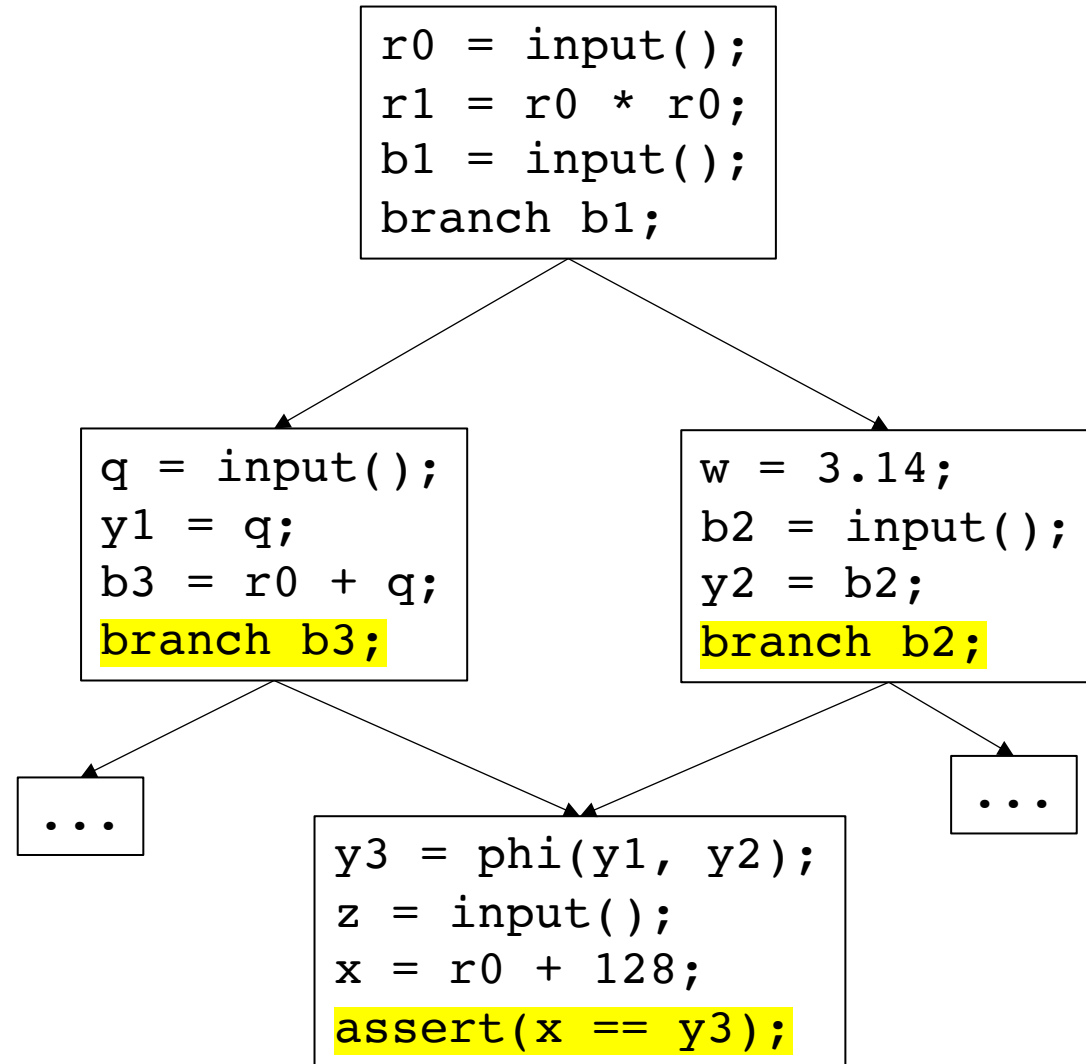
marked:	worklist:
assert()	x
	y3



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());
  }
}
```

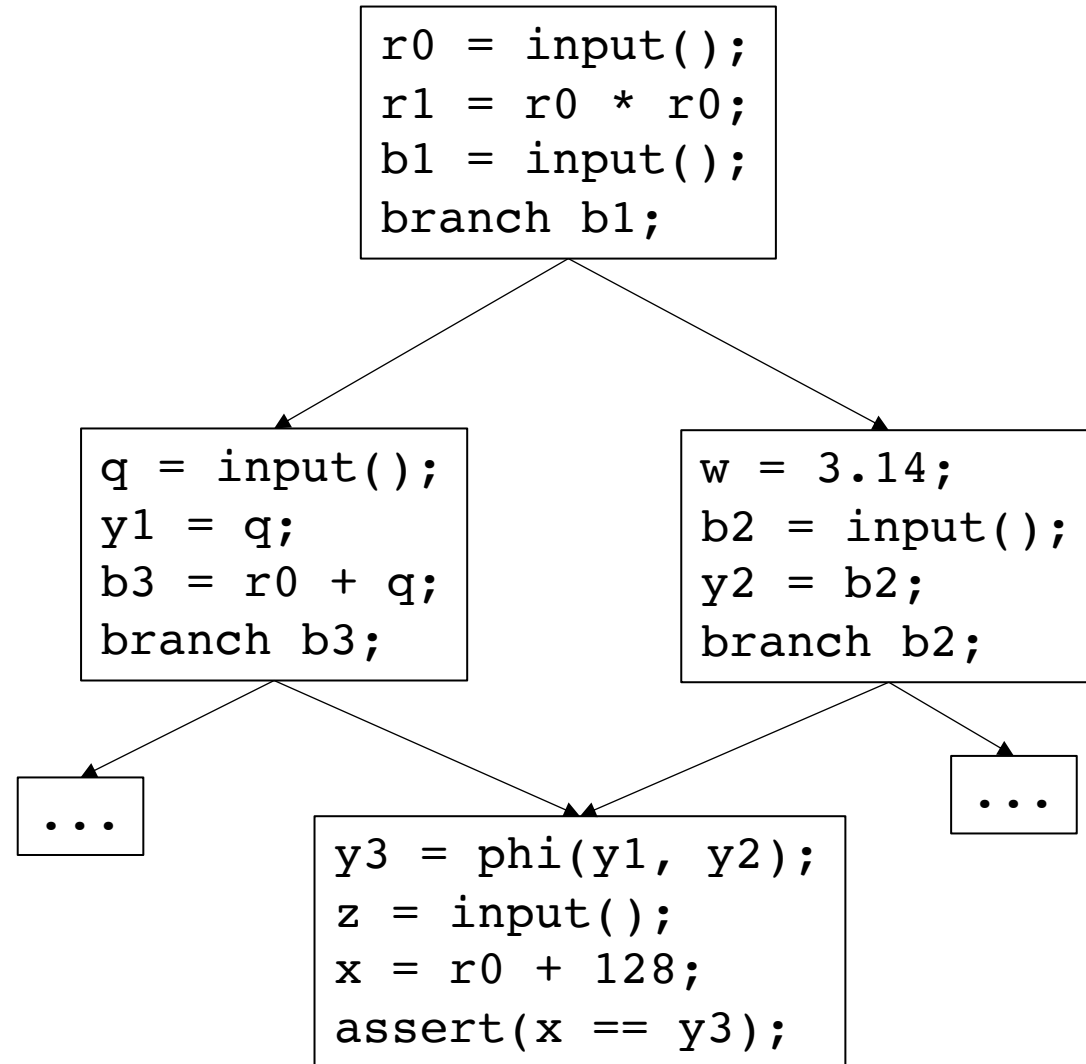
marked:	worklist:
assert()	x
	y3
	branch b3
	branch b2



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());
    }
}
```

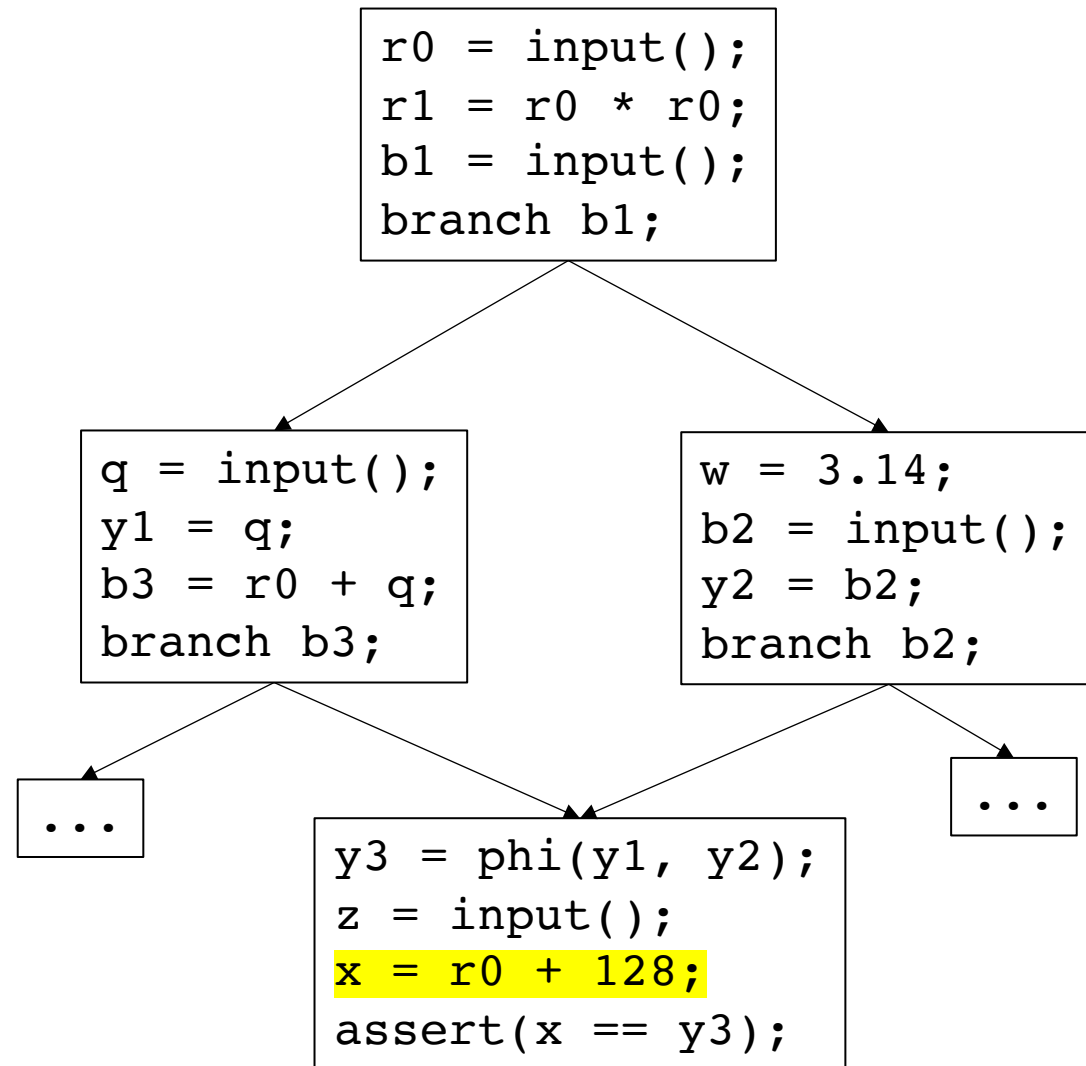
marked:	worklist:
assert()	x
	y3
	branch b3
	branch b2



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());
    }
}
```

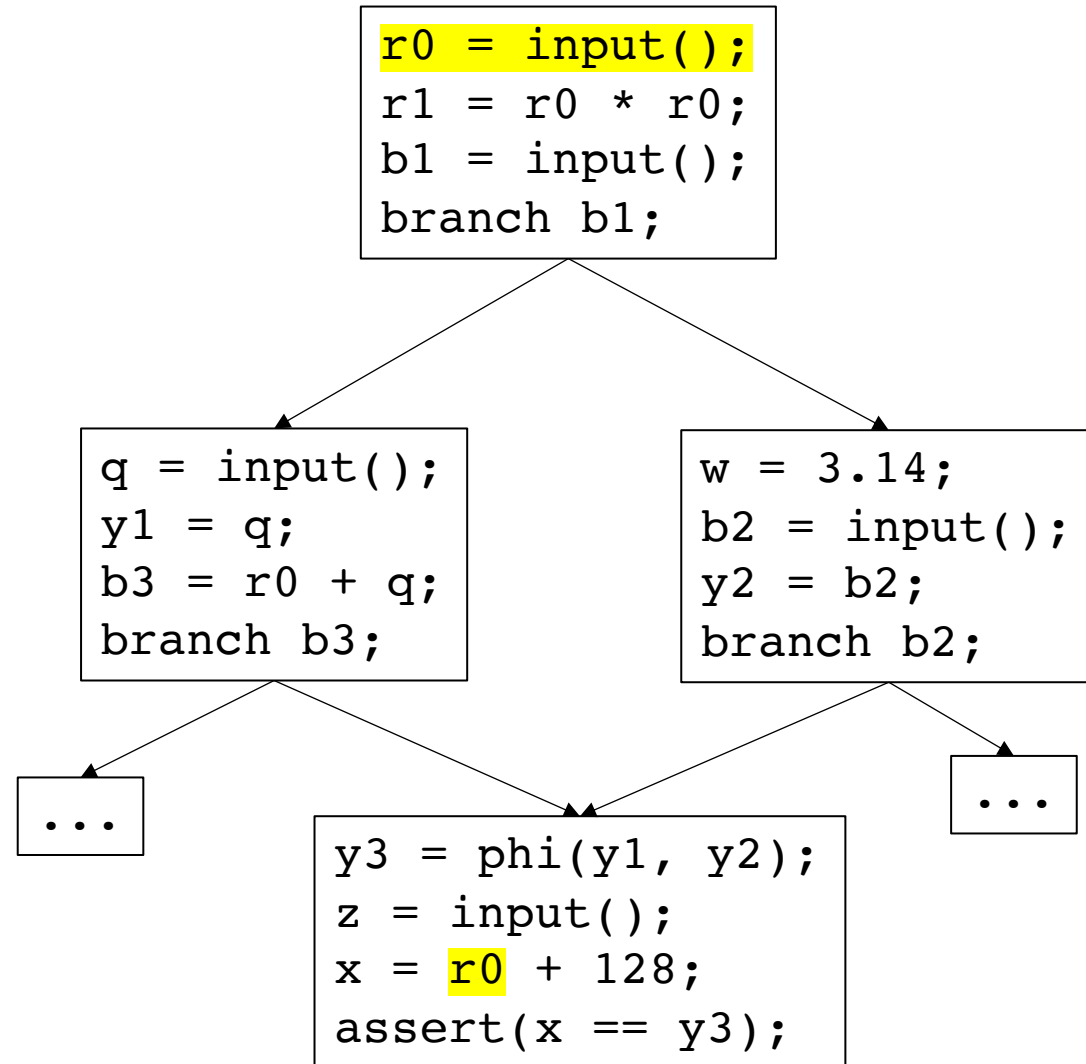
marked:	worklist:
assert()	y3
x	branch b3
	branch b2



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());
    }
}
```

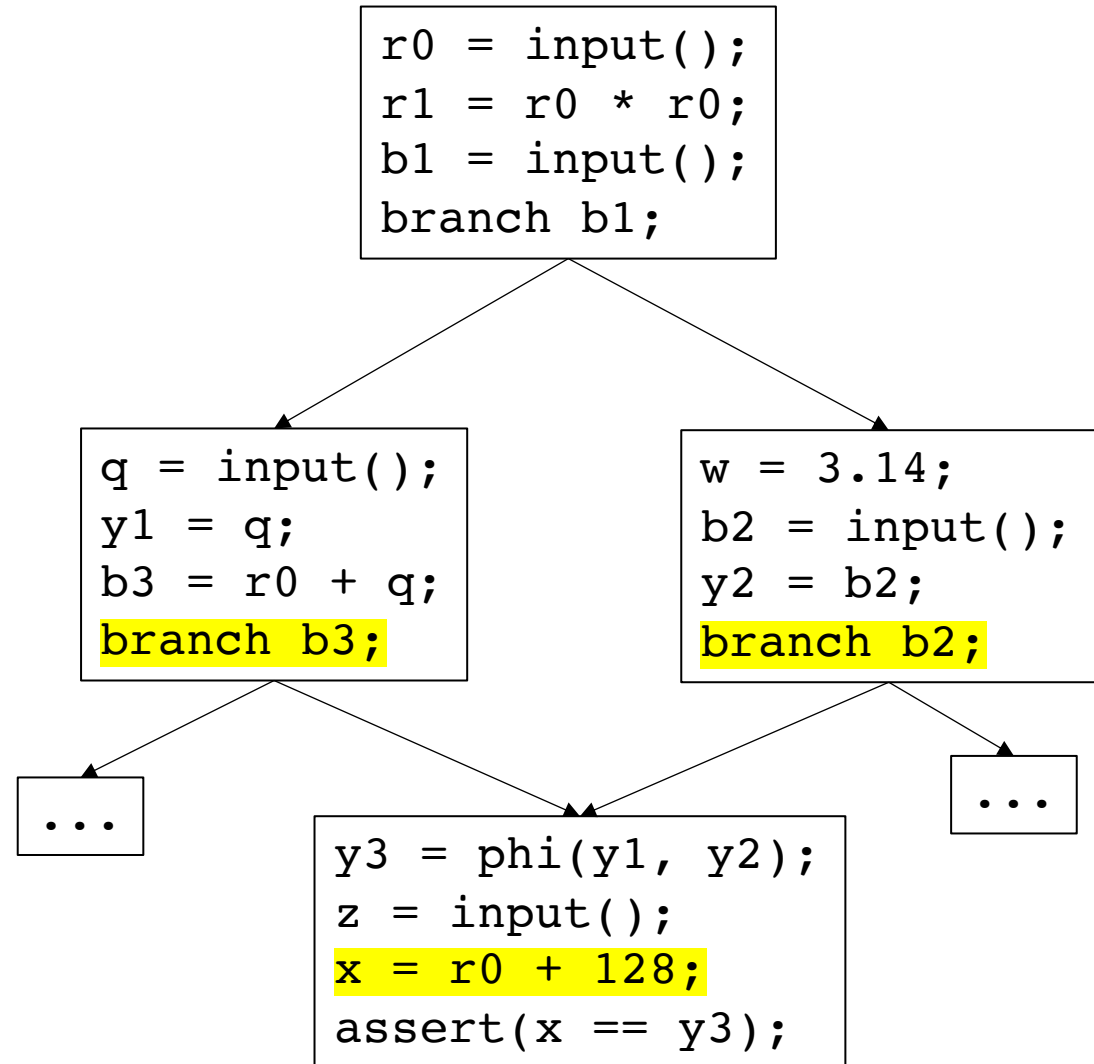
marked:	worklist:
assert()	y3
x	branch b3
	branch b2
	r0



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());
    }
}
```

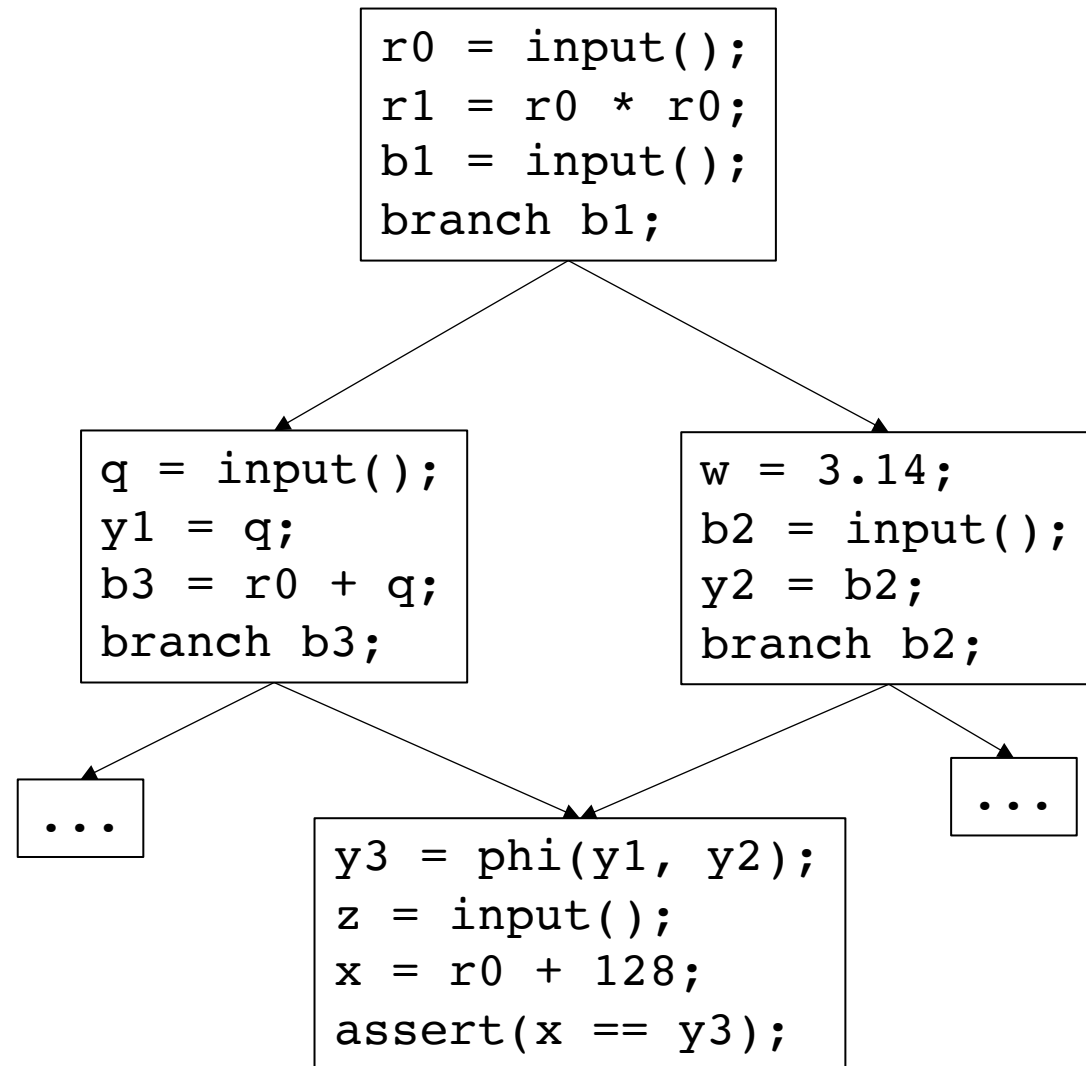
marked:	worklist:
assert()	y3
x	branch b3
	branch b2
	r0



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());
    }
}
```

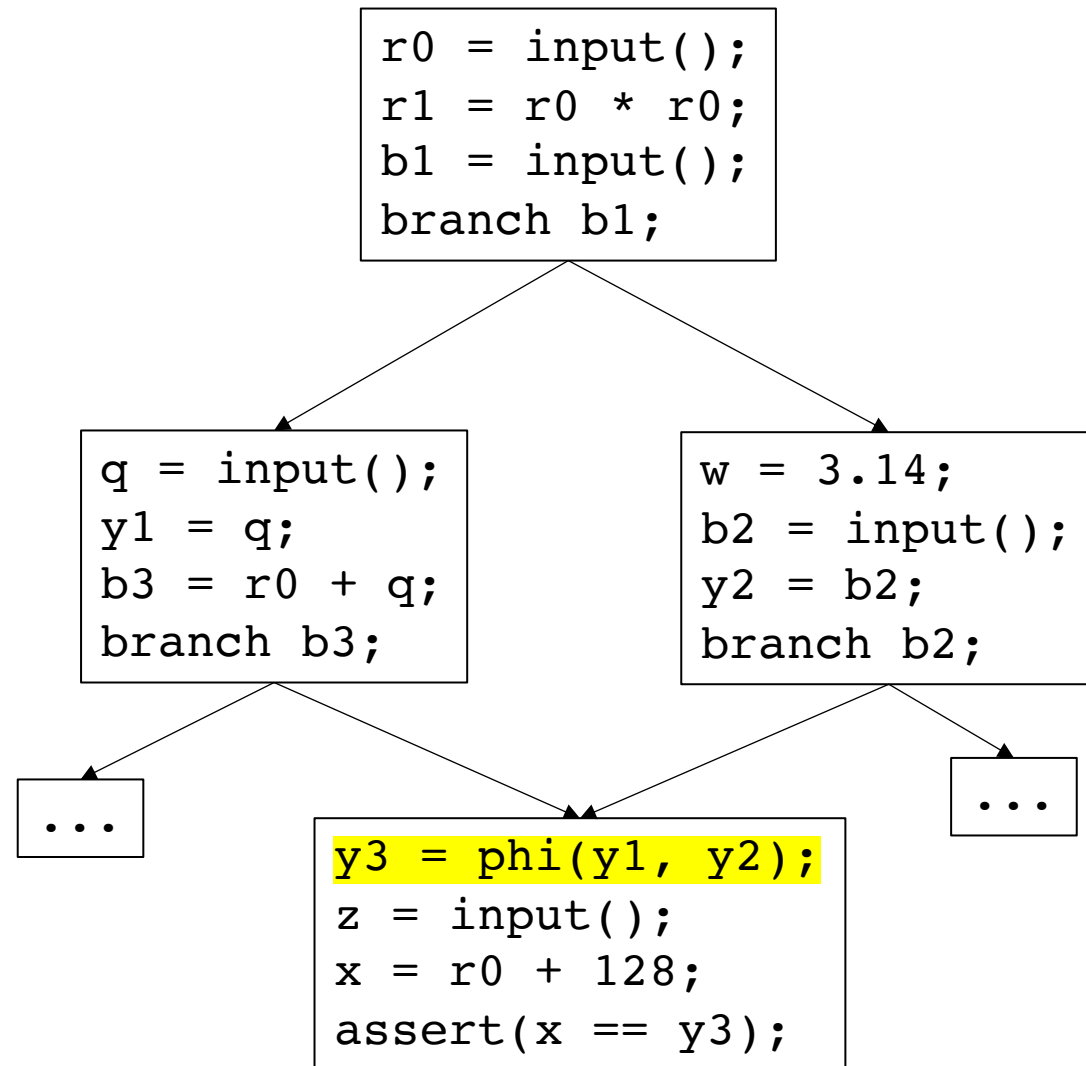
marked:	worklist:
assert()	y3
x	branch b3
	branch b2
	r0



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());
    }
}
```

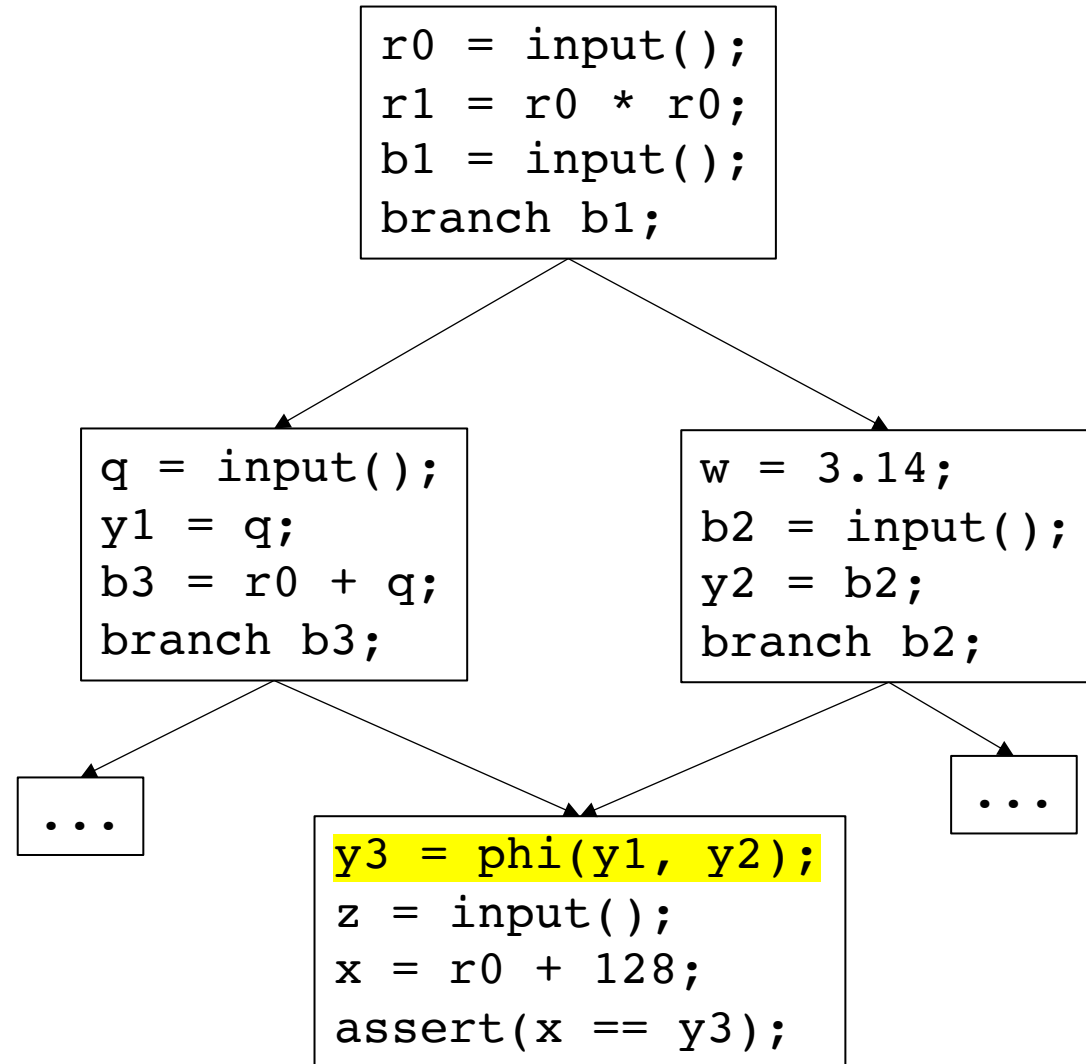
marked:	worklist:
assert()	y3
x	branch b3
	branch b2
	r0



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());
  }
}
```

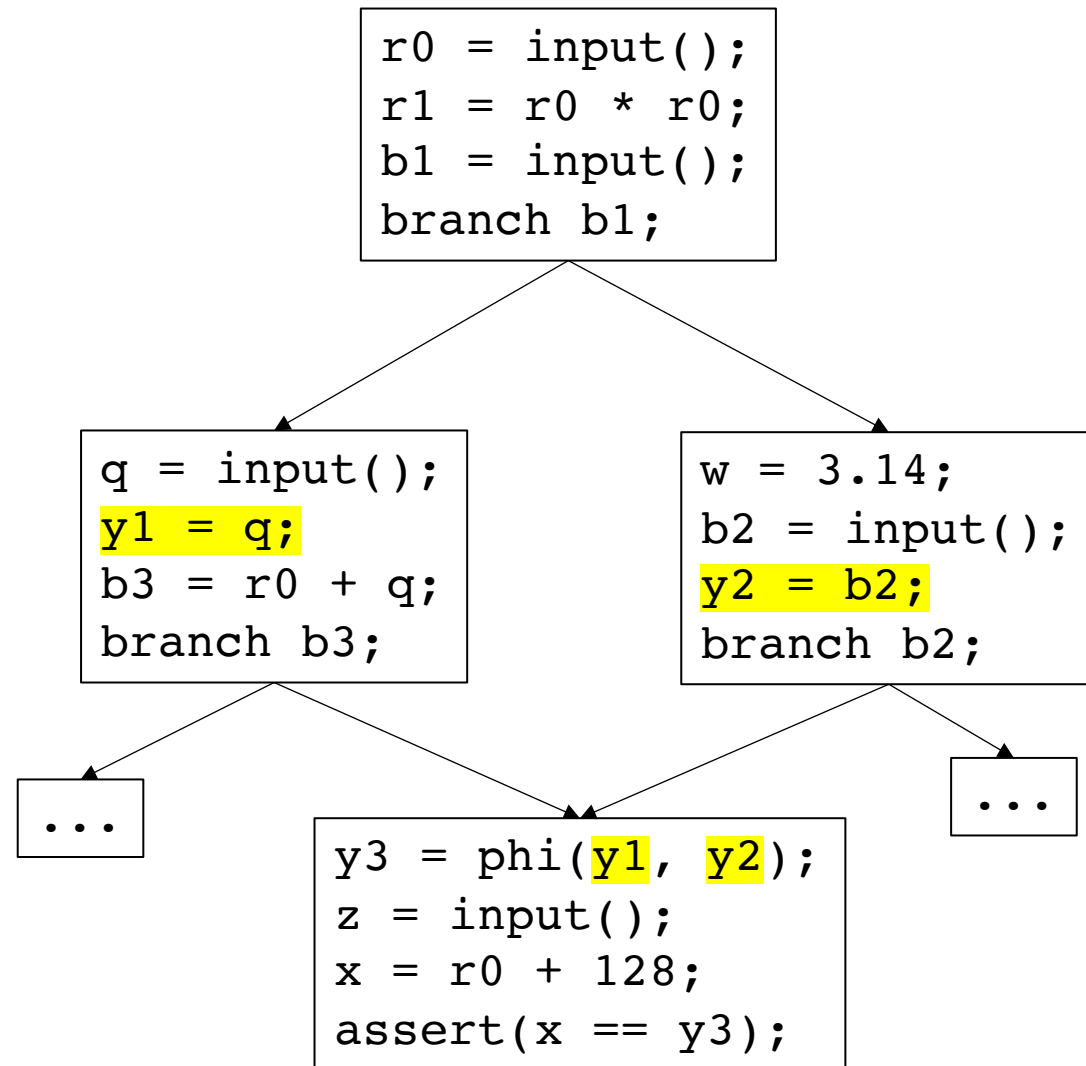
marked:	worklist:
assert()	branch b3
x	branch b2
y3	r0



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());
    }
}
```

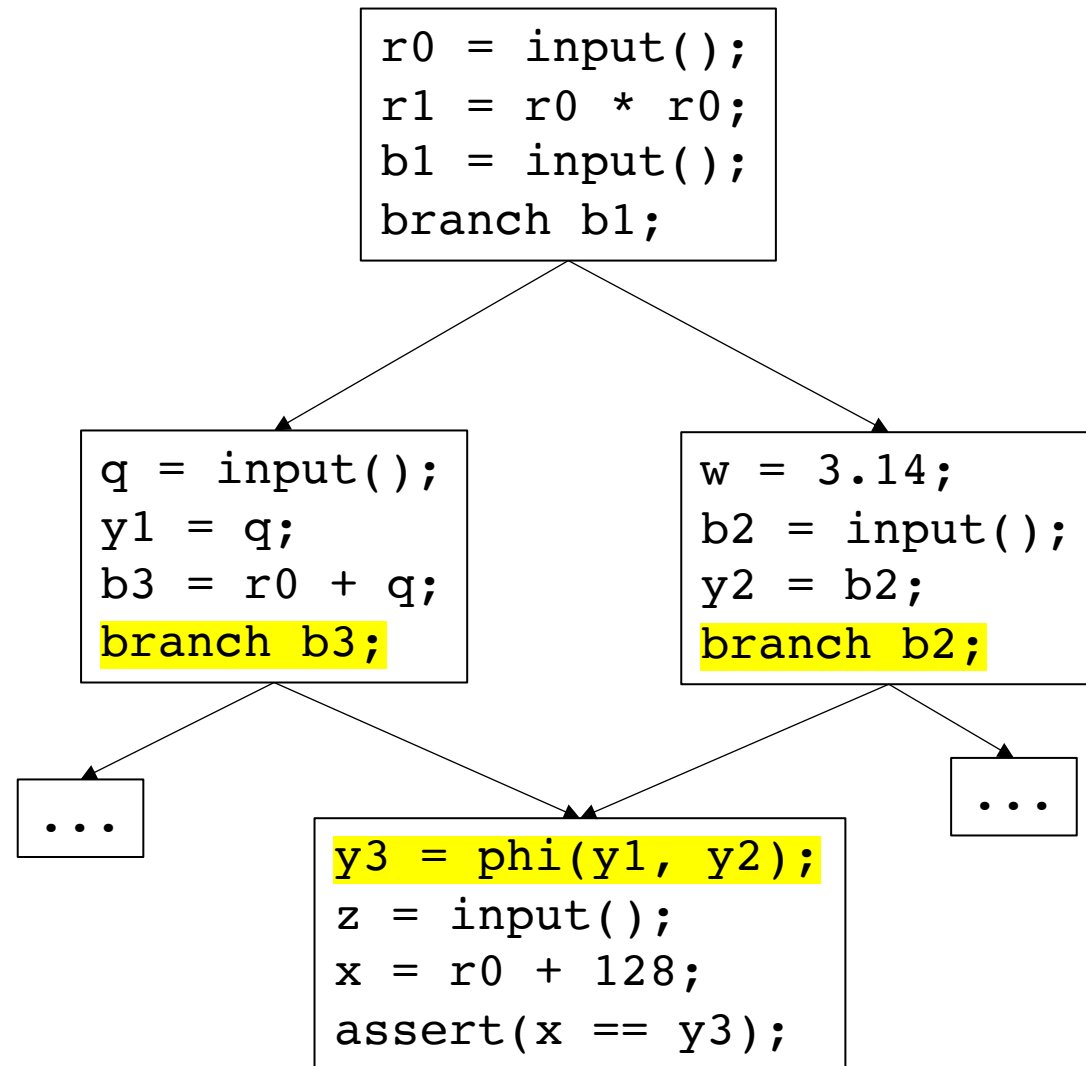
marked:	worklist:
assert()	branch b3
x	branch b2
y3	r0
	y1
	y2



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());
    }
}
```

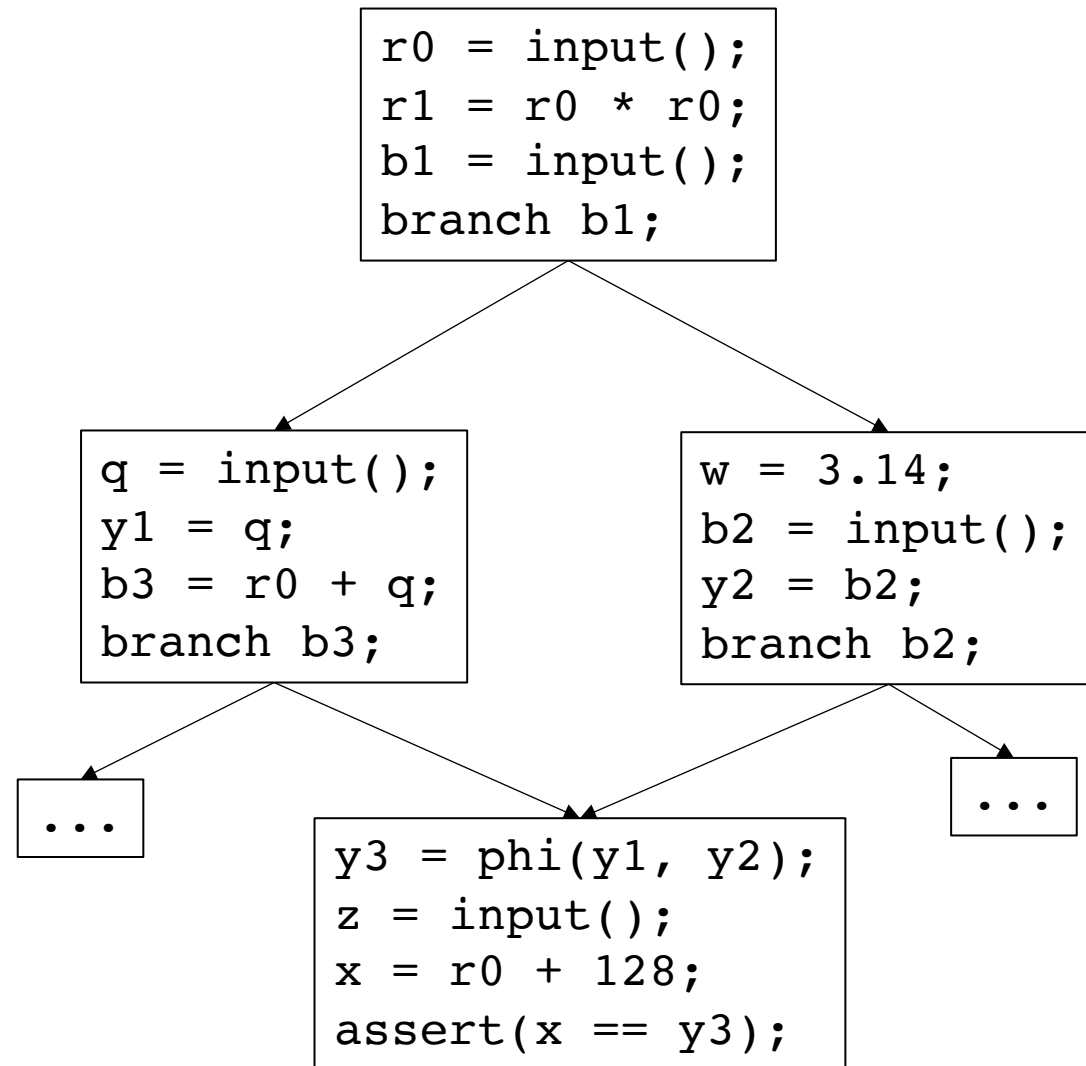
marked:	worklist:
assert()	branch b3
x	branch b2
y3	r0
	y1
	y2



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());
    }
}
```

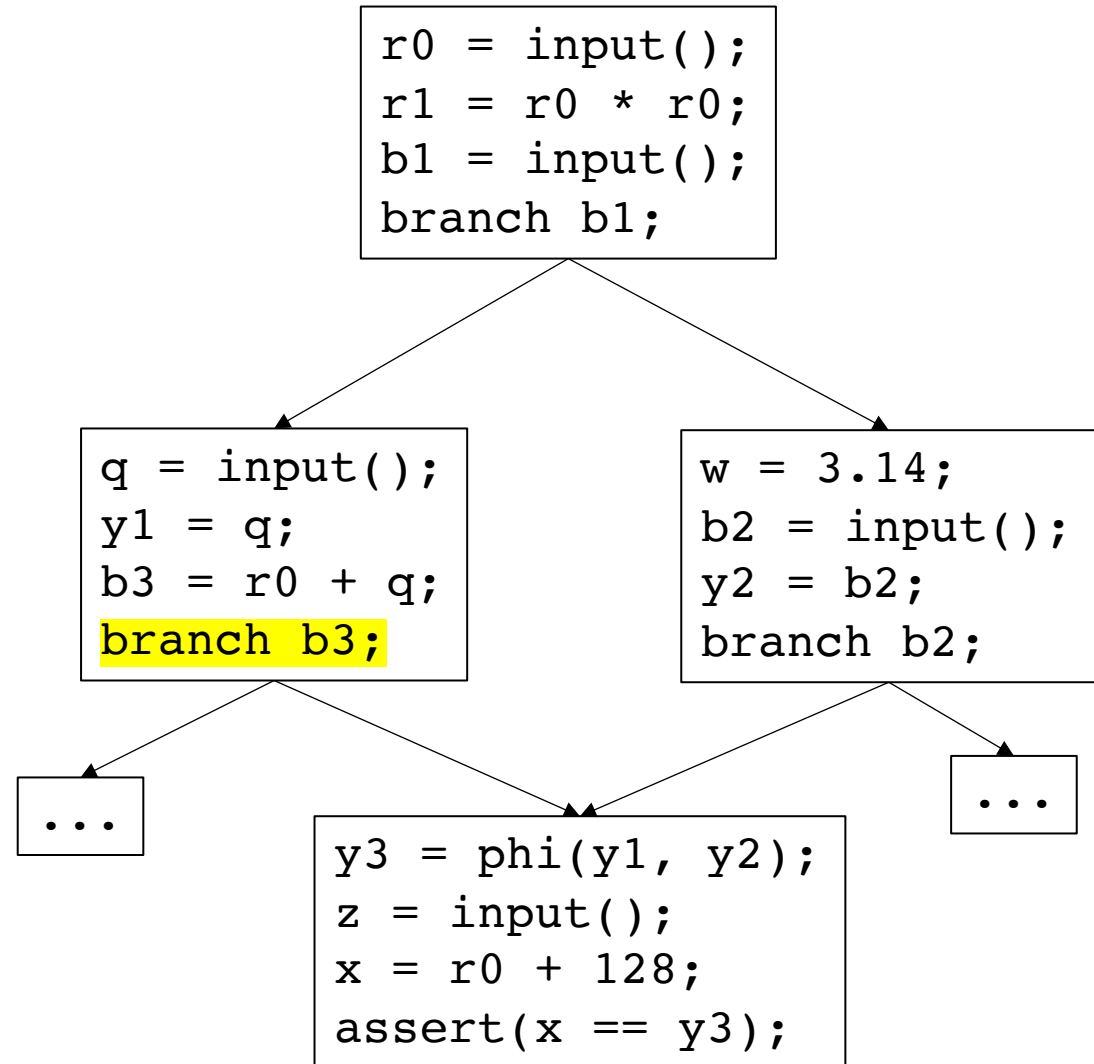
marked:	worklist:
assert()	branch b3
x	branch b2
y3	r0
	y1
	y2



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());
    }
}
```

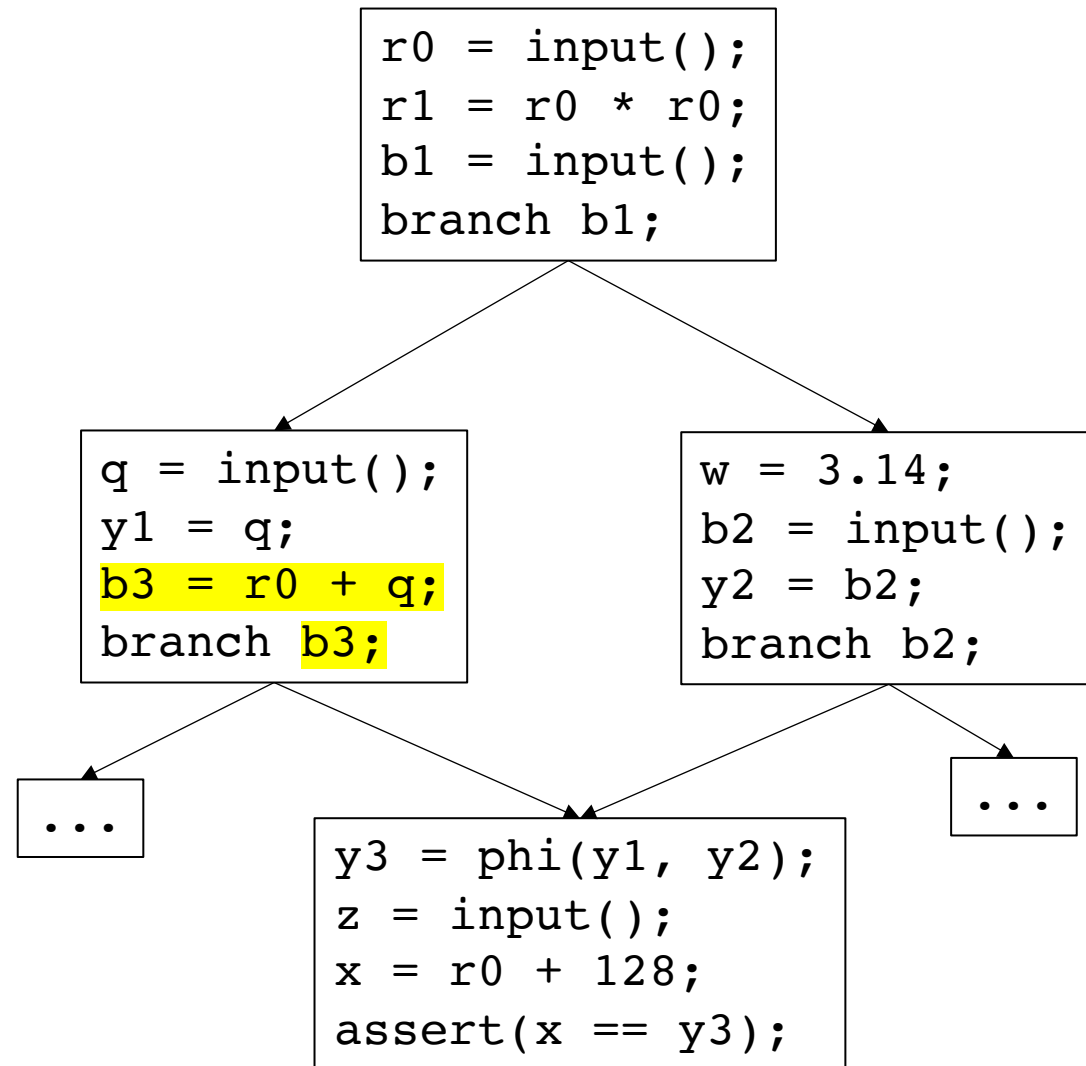
marked:	worklist:
assert()	branch b2
x	r0
y3	y1
branch b3	y2



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());
    }
}
```

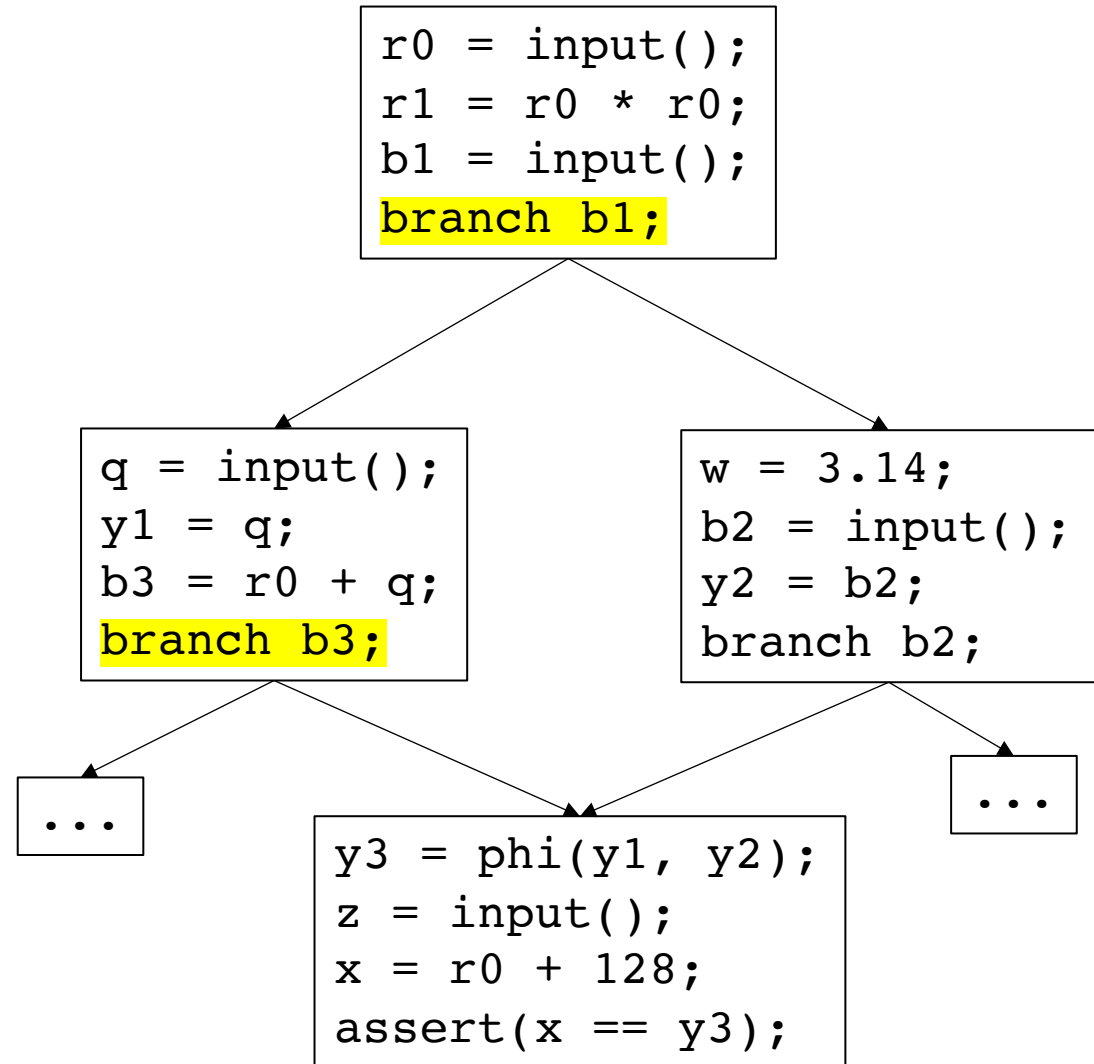
marked:	worklist:
assert()	branch b2
x	r0
y3	y1
branch b3	y2
	b3



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
marked: r0
assert() y1
x y2
y3 b3
branch b3 branch b1

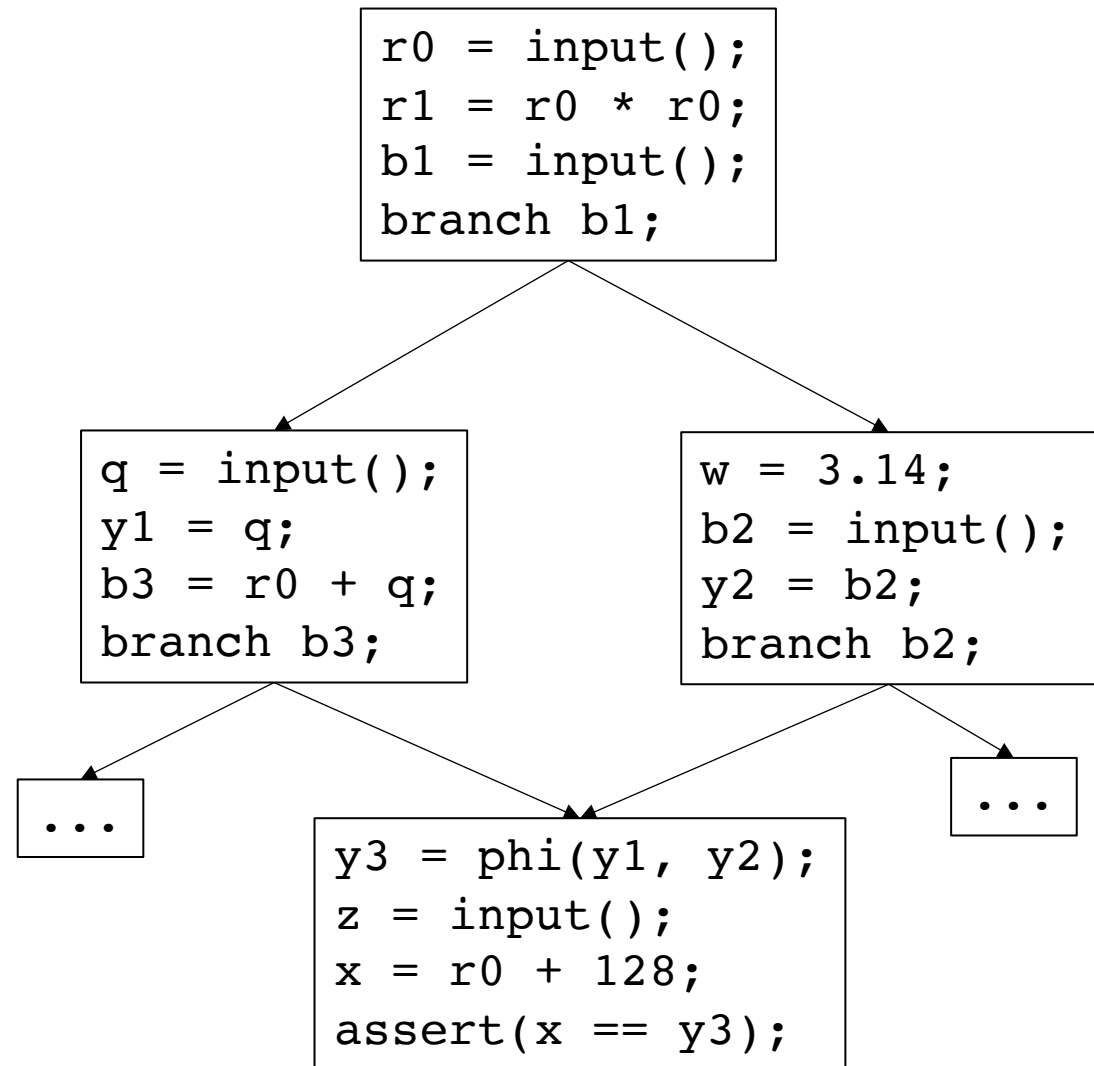


Backwards slicing algorithm

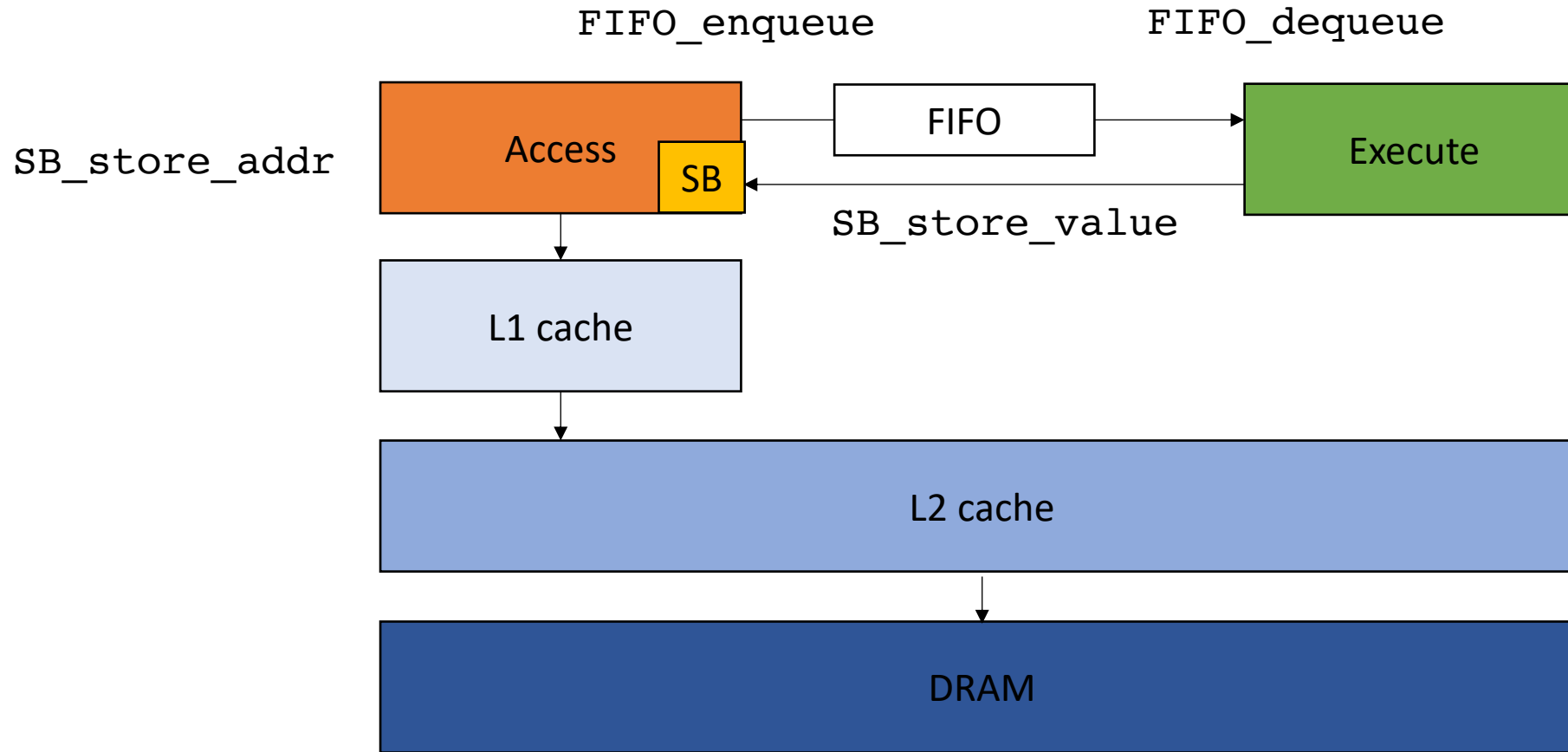
```
Worklist = S; // slicing criteria
while (!Worklist.empty()) {
  stmt = Worklist.pop();
  if (is_marked(stmt)) {
    continue;
  }
  mark(stmt);
  for a in stmt.args() {
    worklist.append(a);
  }
  for p in cfg[stmt].predecessors() {
    worklist.append(p.branch_stmt());
  }
}
```

worklist:
branch b2
marked: r0
assert() y1
x y2
y3 b3
branch b3 branch b1

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

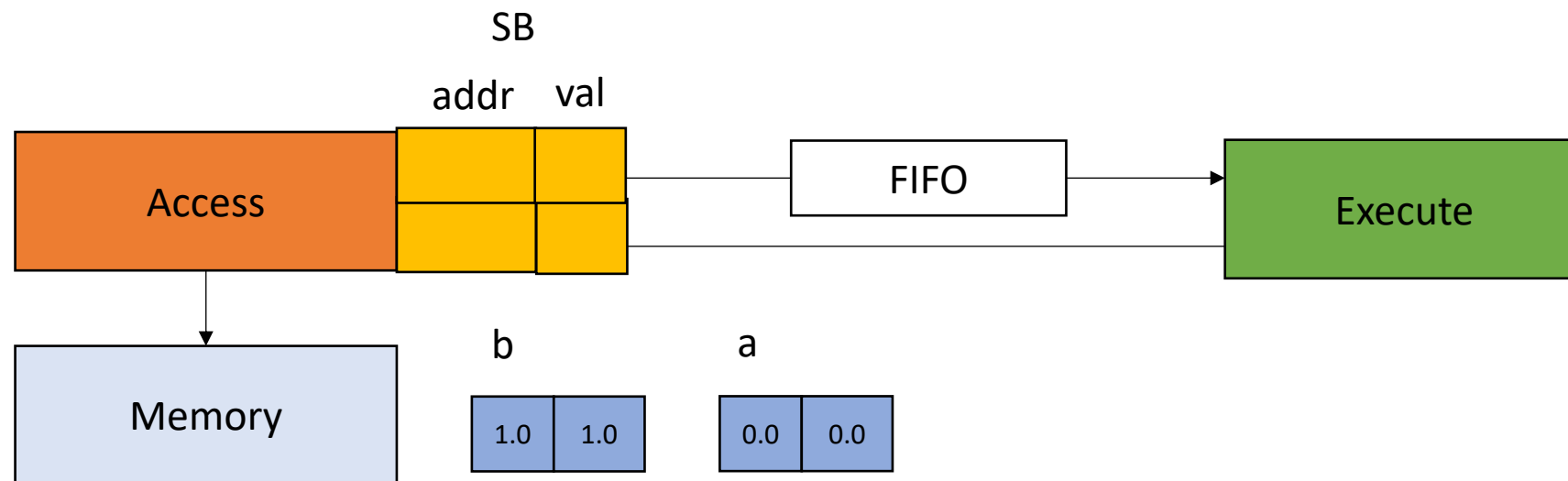
```
// SSA pseudo code
for (int i = 0; i < SIZE; i++) {
    float r0 = FIFO_dequeue()
    float r1 = r0 * 3.14;
    SB_store_value(r1);
}
```

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);
}
```



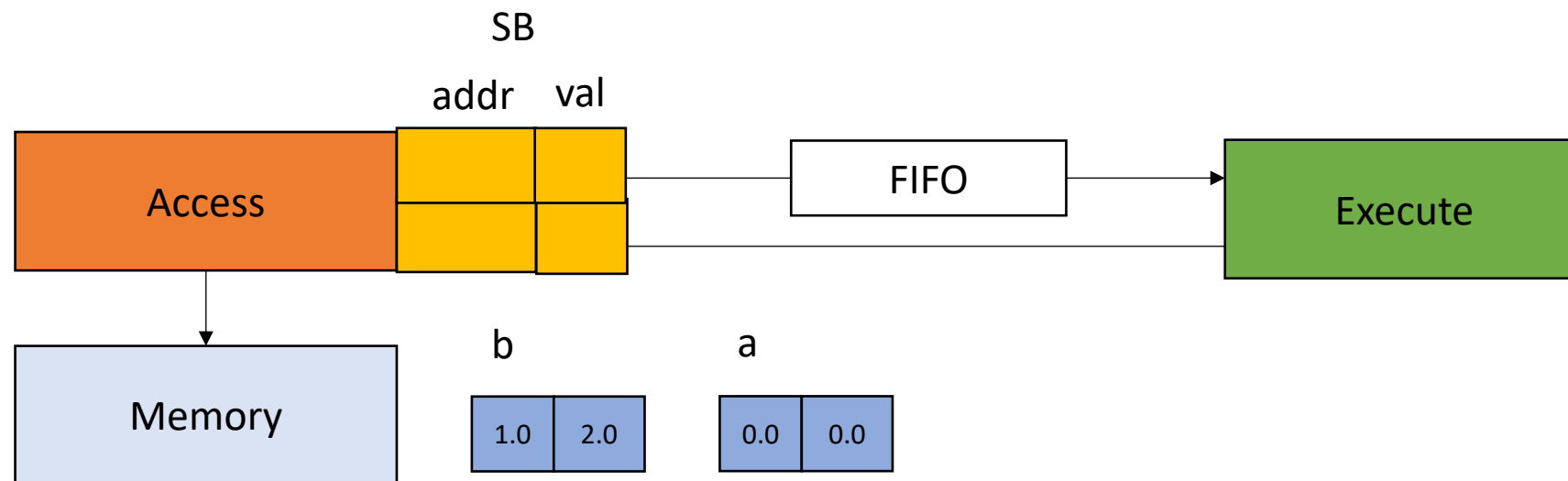
Access

```
// 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*



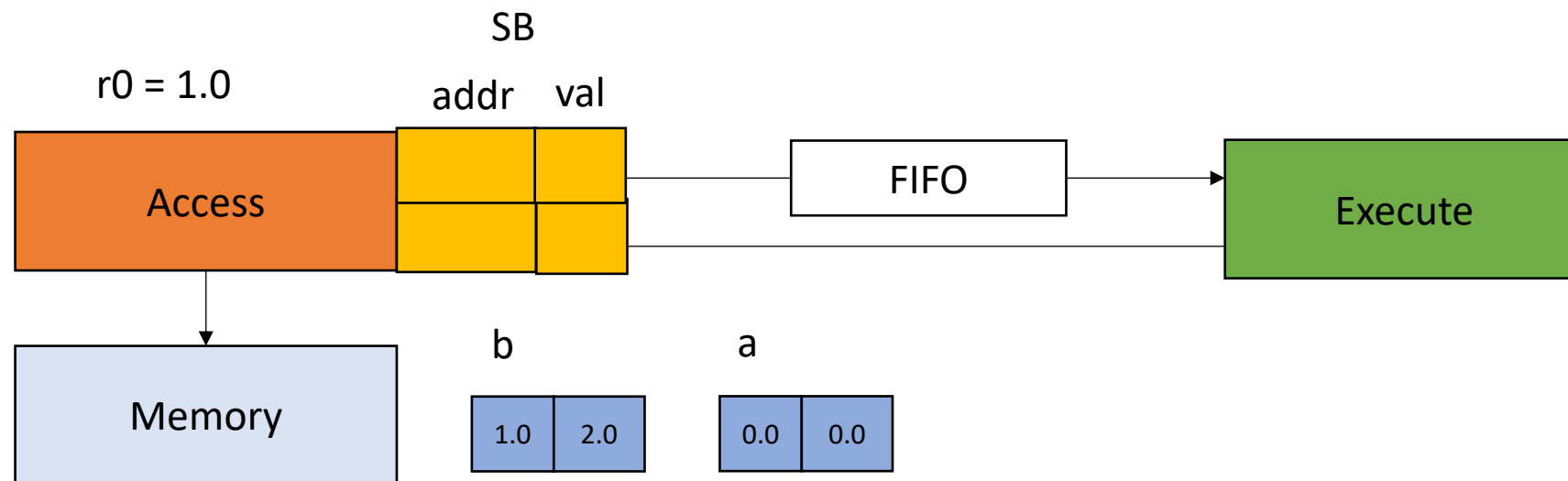
Access

```
// 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*



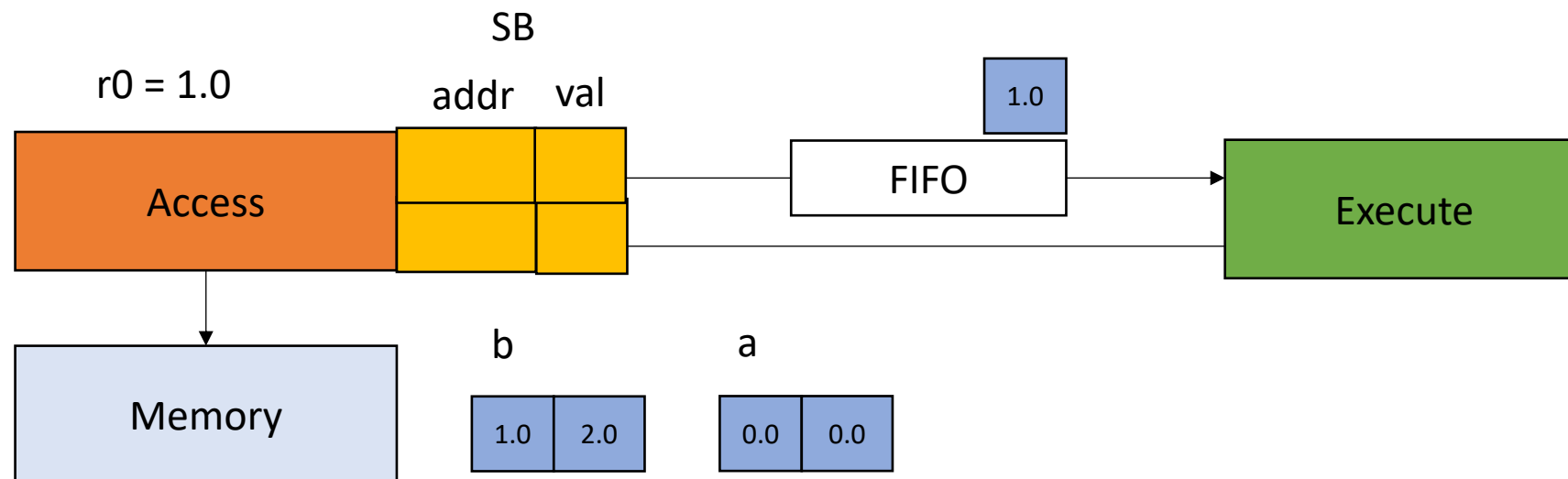
Access

```
// 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*



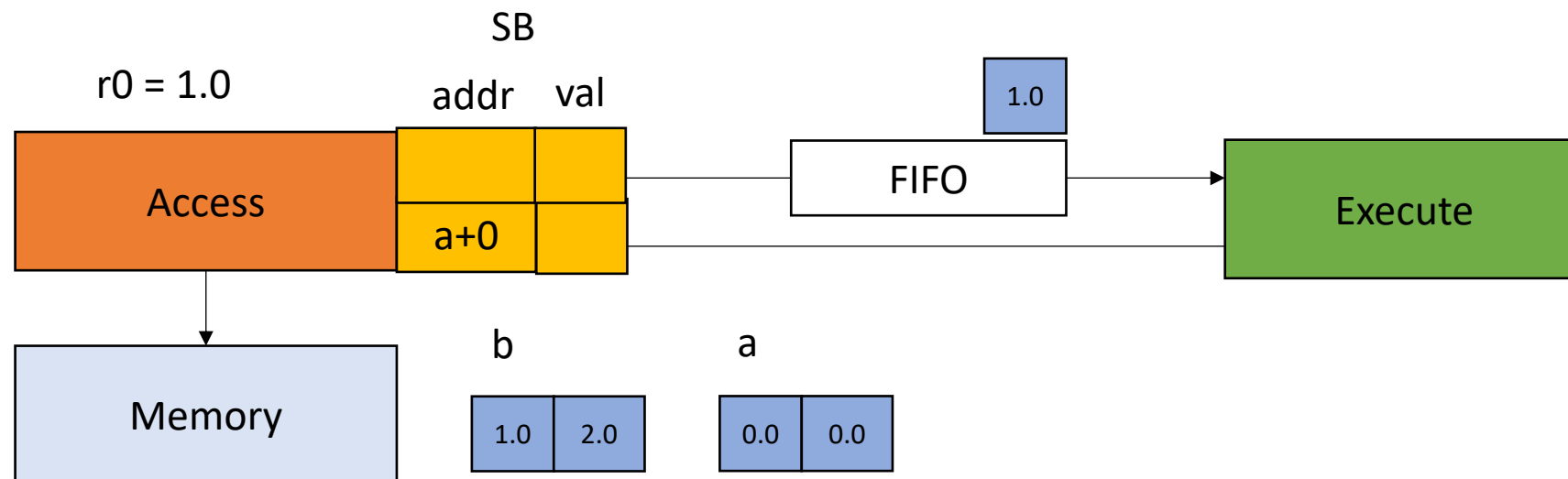
Access

```
// 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*



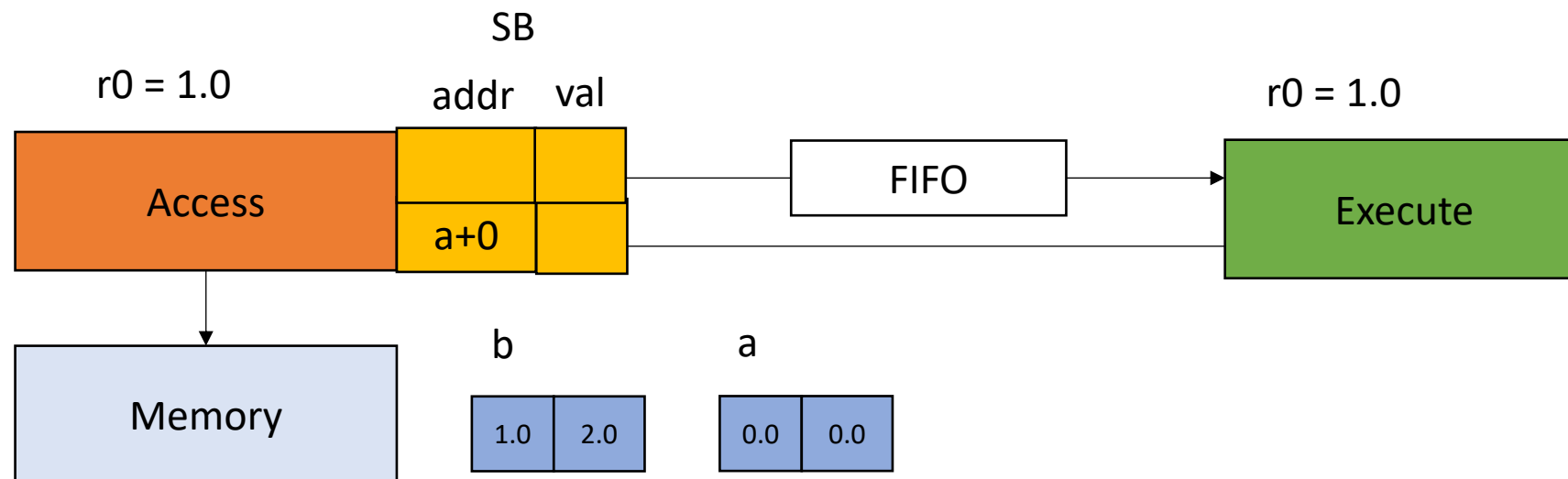
Access

```
// 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*



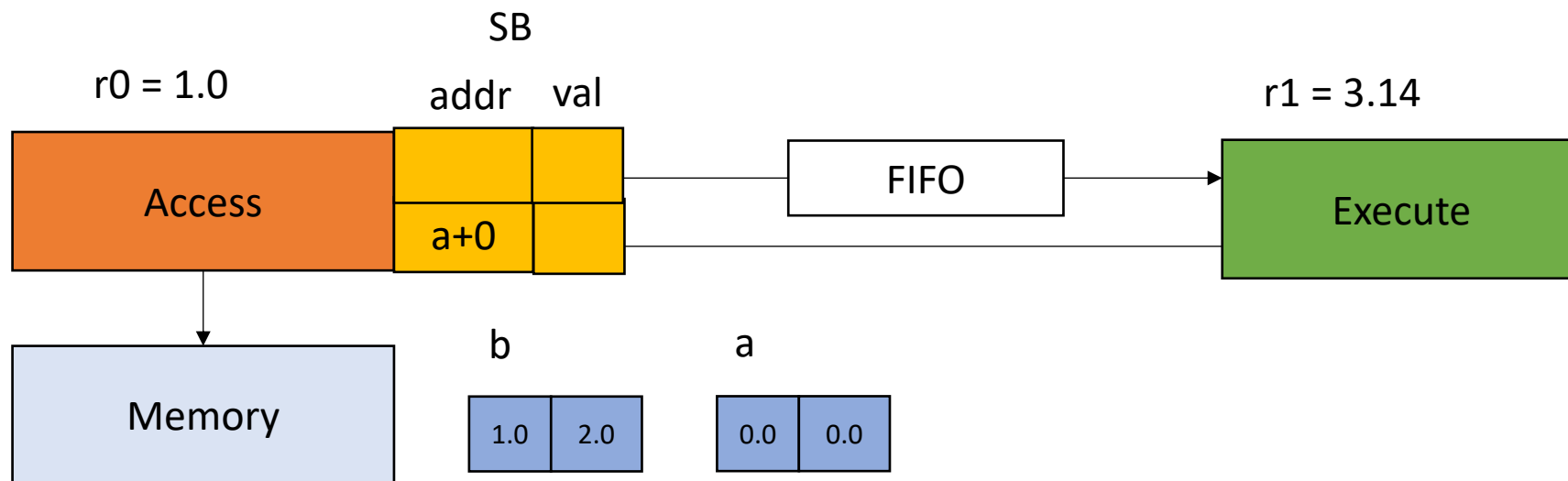
Access

```
// 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*



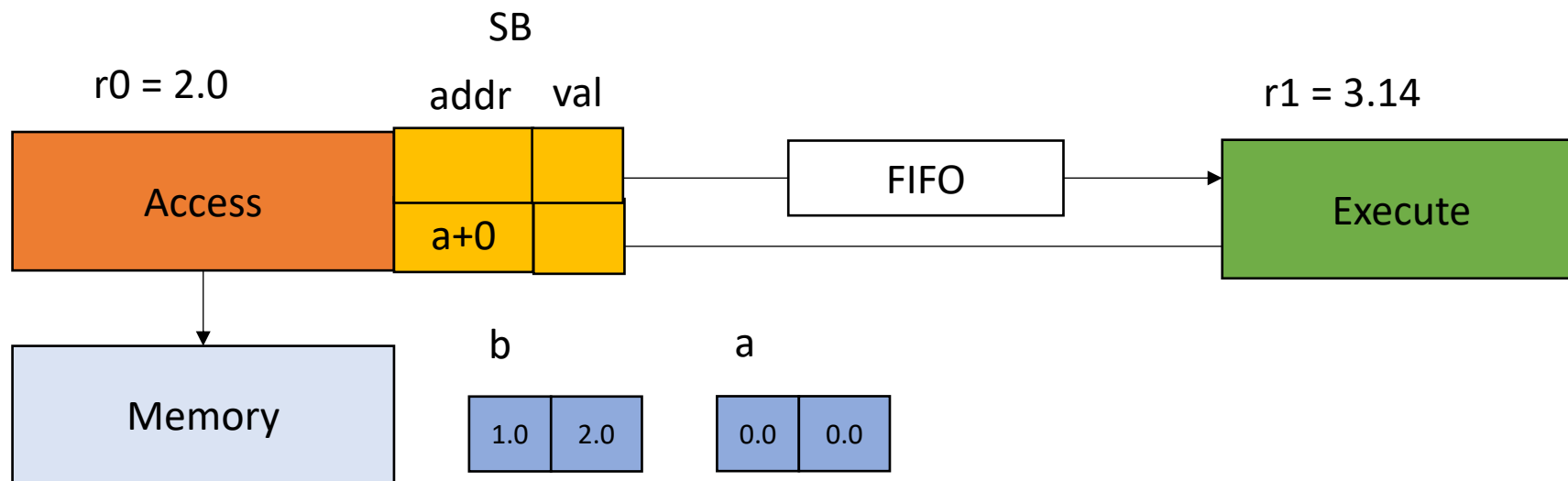
Access

```
// 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*



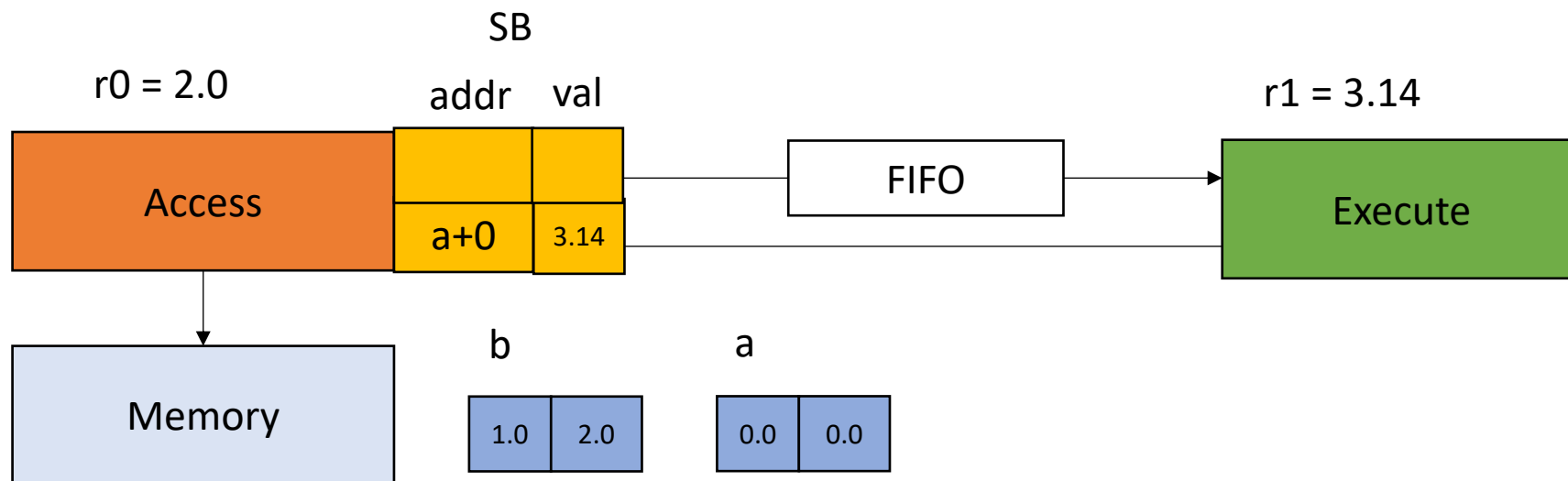
Access

```
// 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*



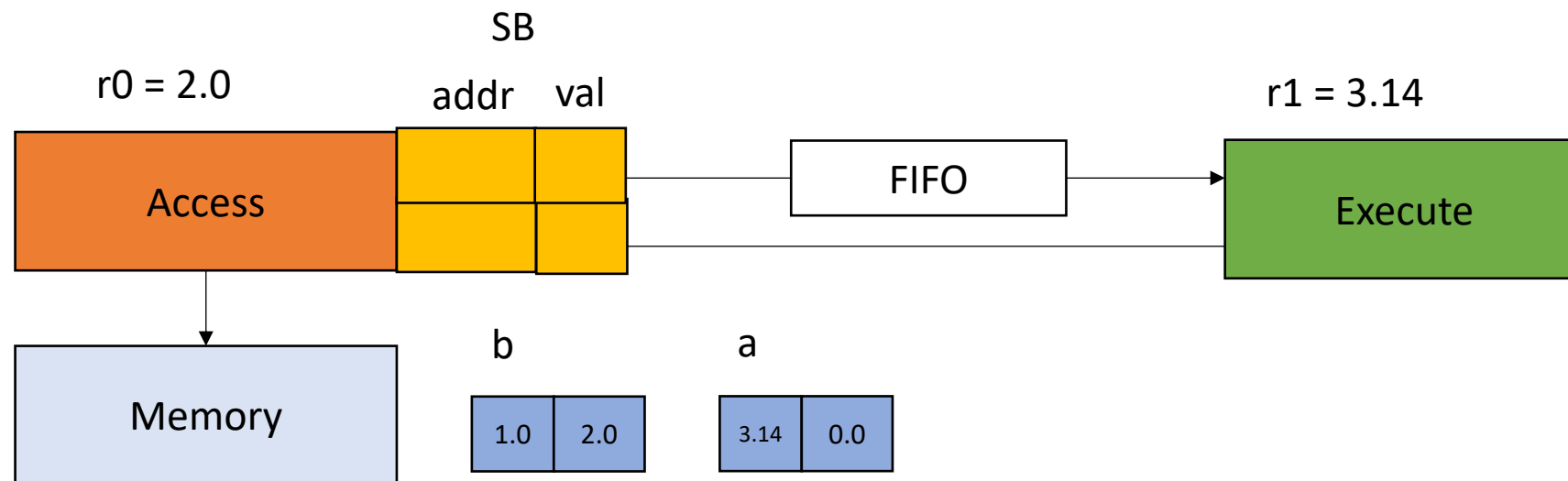
Access

```
// 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*



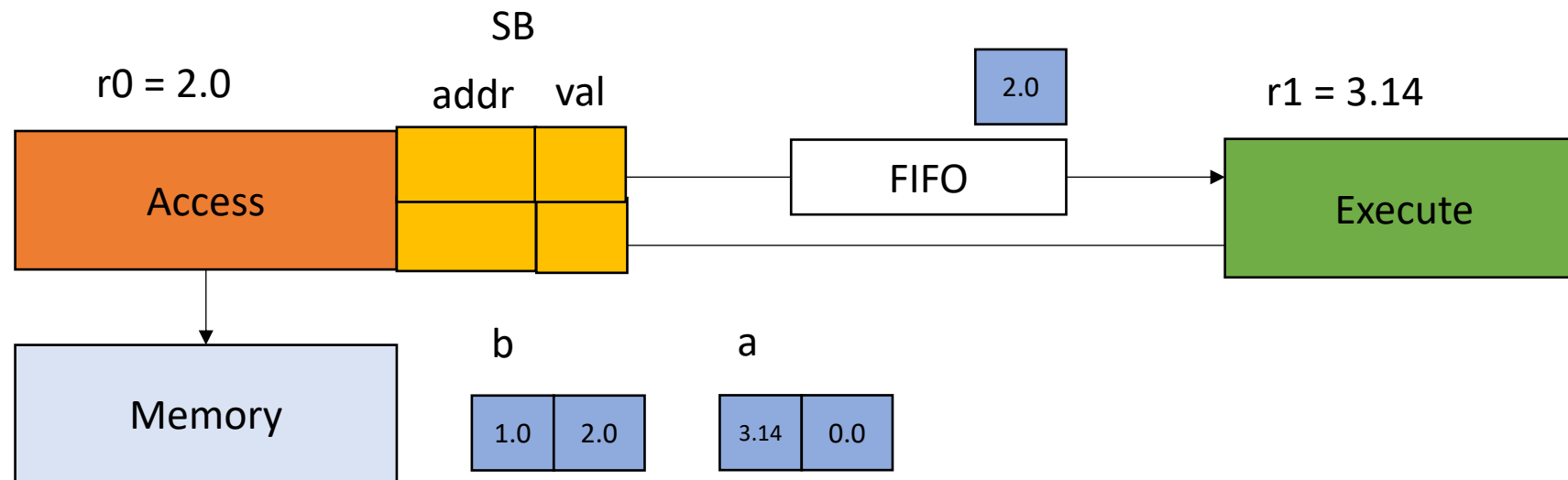
Access

```
// 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*



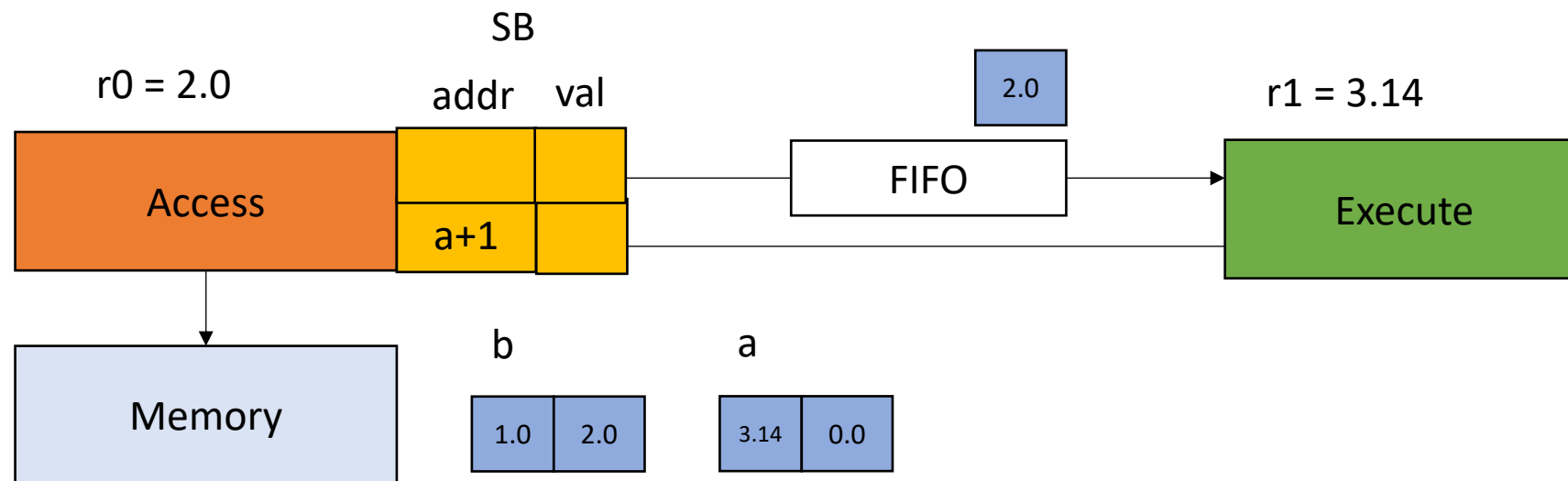
Access

```
// 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*



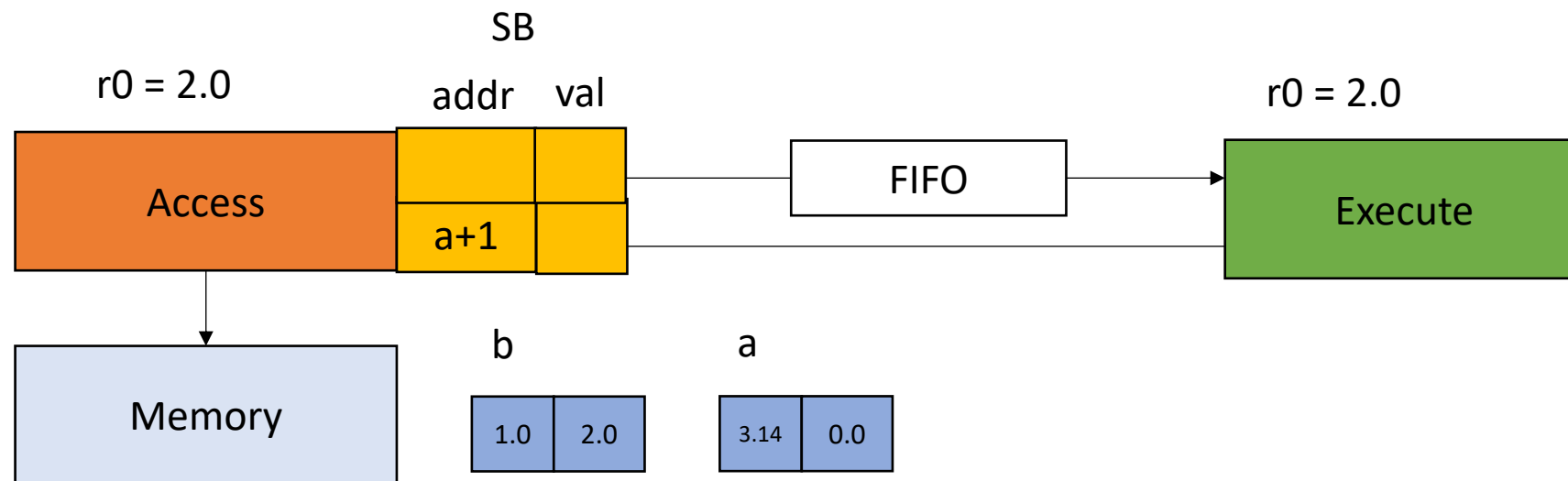
Access

```
// 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*



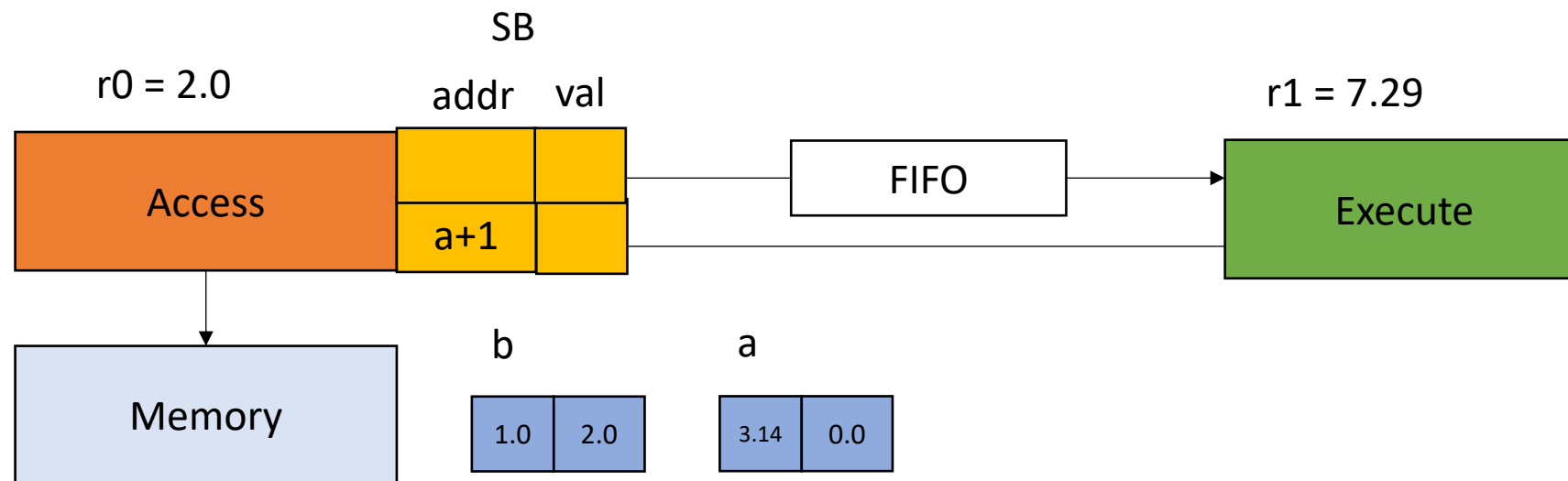
Access

```
// 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*



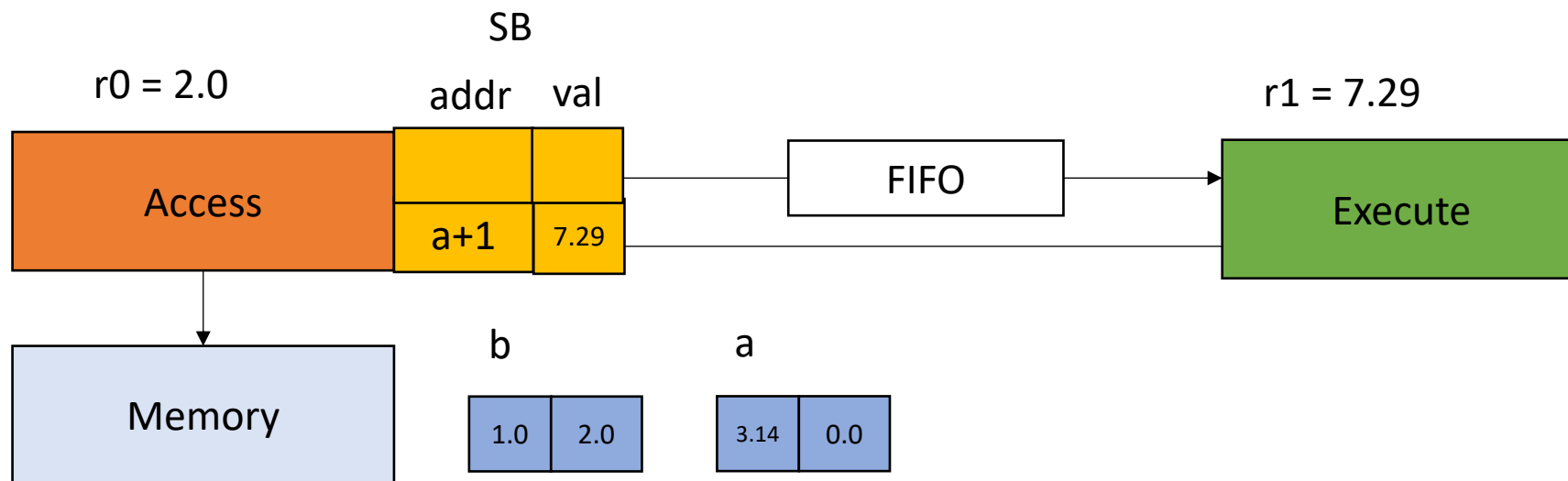
Access

```
// 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*



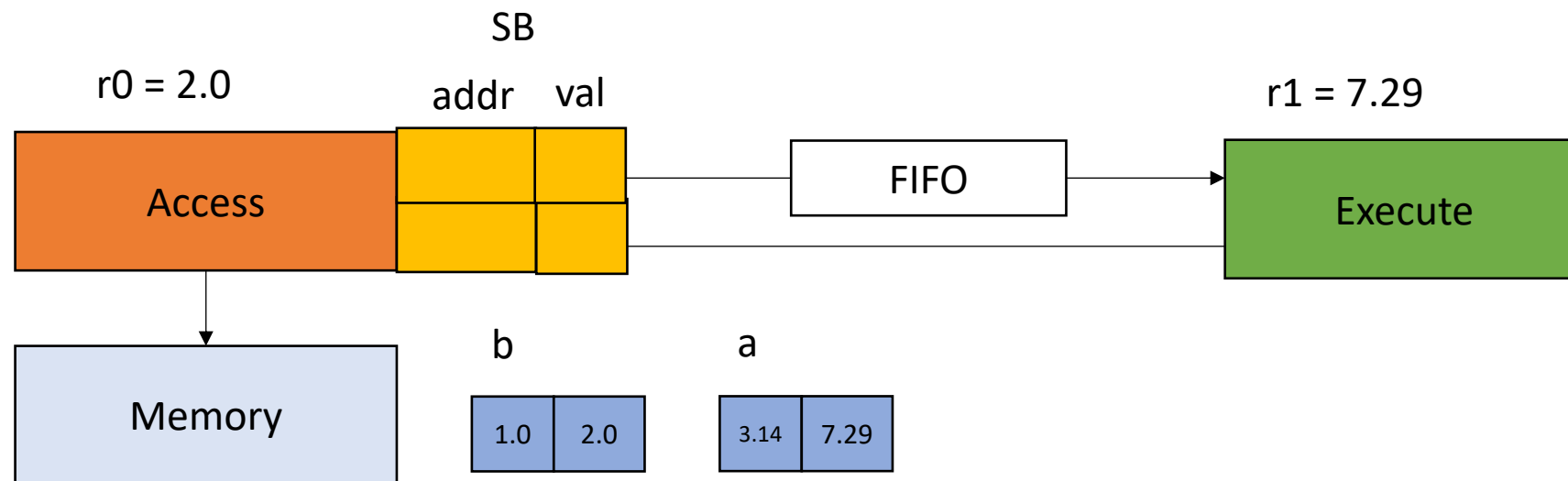
Access

```
// 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*

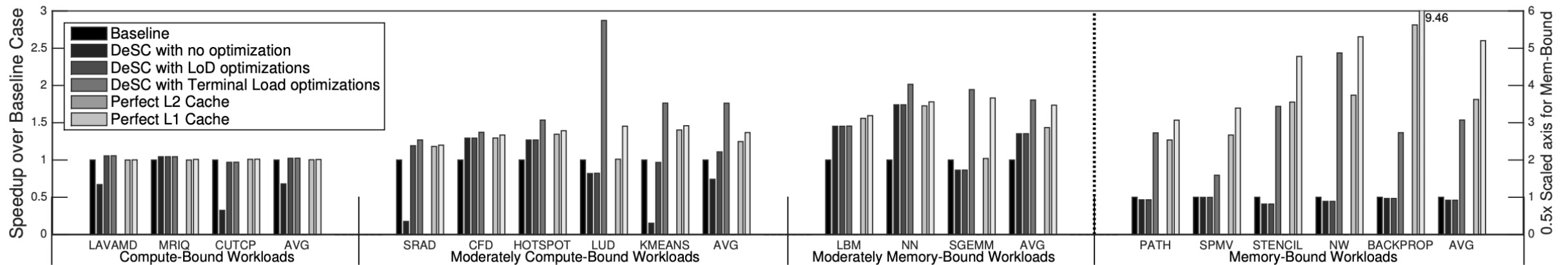


Performance bounds

- 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 simple core models, we can approximate: $E(p) = C(p) + M(p)$
 - Why might this not be completely accurate for more complex cores?
- 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?

Performance results in DeSC (micro '15)

Access and Execute are OoO cores



Final considerations

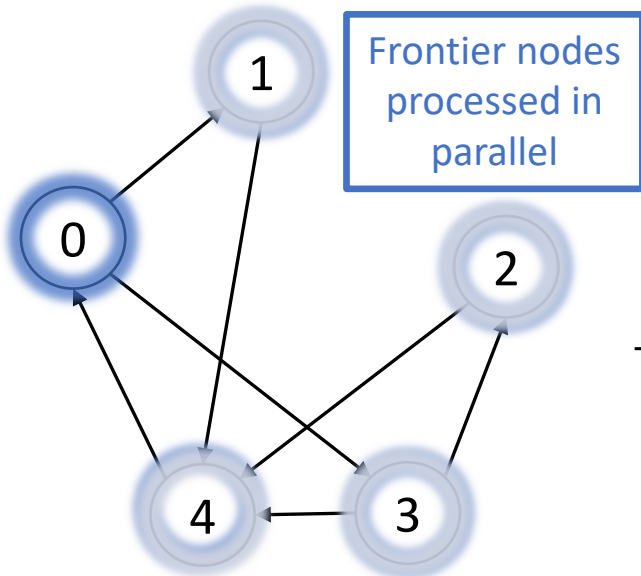
- Dependencies:
 - If Access depends on a value from Execute, performance can suffer
 - Also called LoD (loss of decoupling events)
- Coherence:
 - Access must read up-to-date values, even when waiting on Execute
- More optimizations:
 - Similar to asynchronous stores, some loads can be done asynchronously as well (if the value is not needed by the Access).

GraphAttack (taco '21)

- Follow up work led by Aninda Manocha at Princeton
 - She made some great slides!
- Further specializing DAE
 - Graph applications
 - in order cores

Graph Applications: Access Patterns are Irregular

- Iterative, frontier-based graph applications
 - Describes many graph processing workloads (e.g. BFS, SSSP, PR)
- **Indirect** accesses to neighbor data
 - Conditionally populate next frontier



	neighbors				
	0	1	2	3	4
0	0	1	0	1	0
1	0	0	1	0	0
2	0	0	0	0	1
3	0	0	1	0	1
4	1	0	0	0	0

```

for node in frontier:
    val = process_node(node)
    for neib in G.neighbors(node):
        update = update_neib(node_vals, val, neib)
        if (add_to_frontier(update)):
            new_frontier.push(neib)
    
```

Indirect memory access due to neighbor locations

Iterative, frontier-based graph application

Stores IDs of nodes to process

frontier

0	1	2	3	4
0	3			

Stores node property data

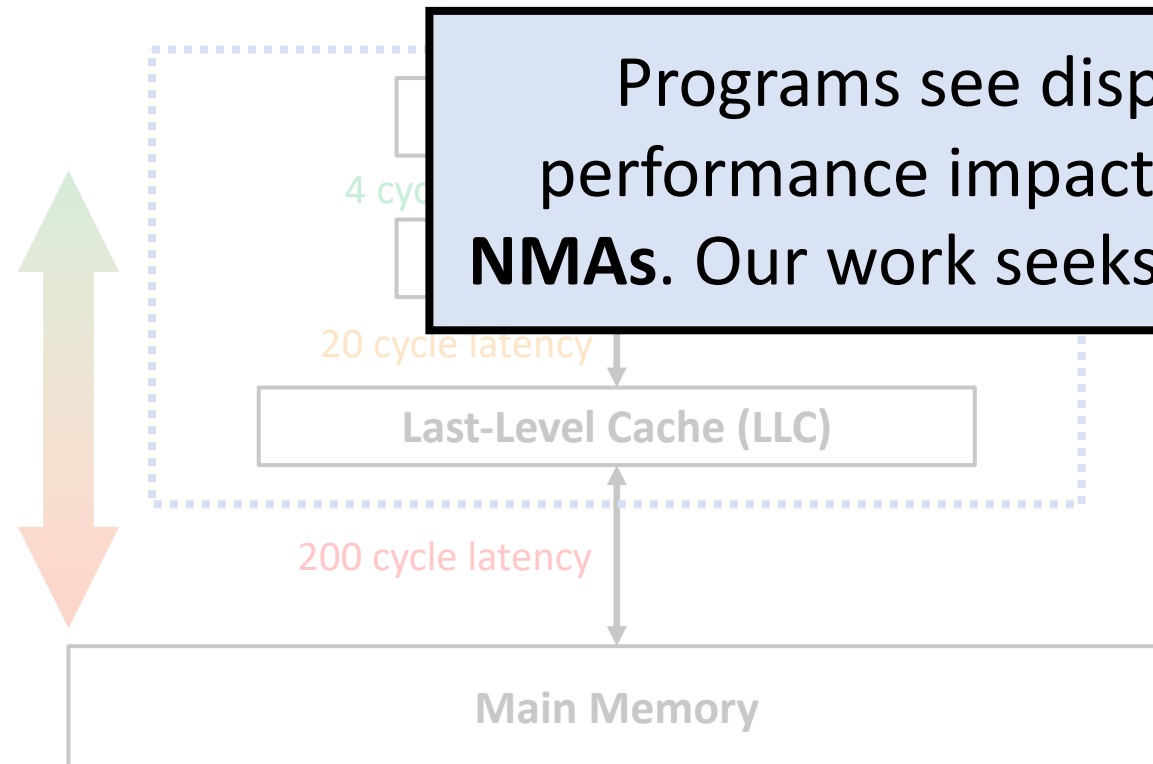
node_vals
(hops from 0)

0	1	2	3	4

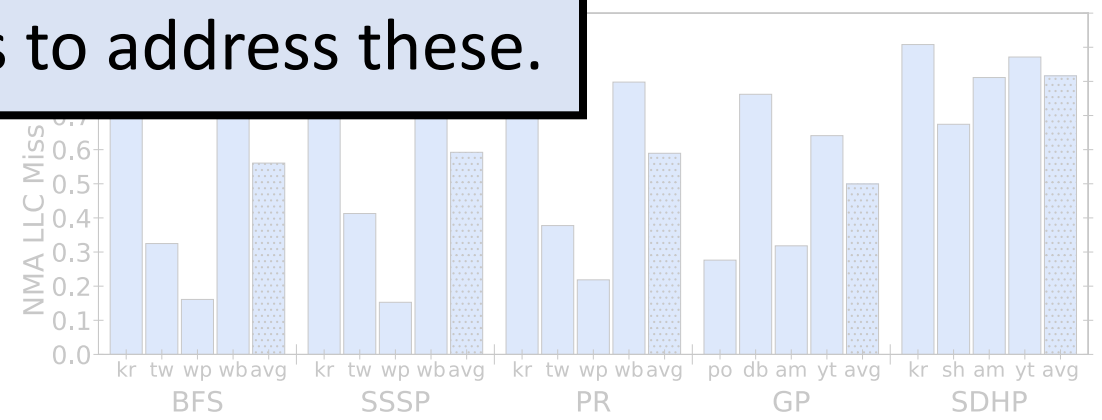
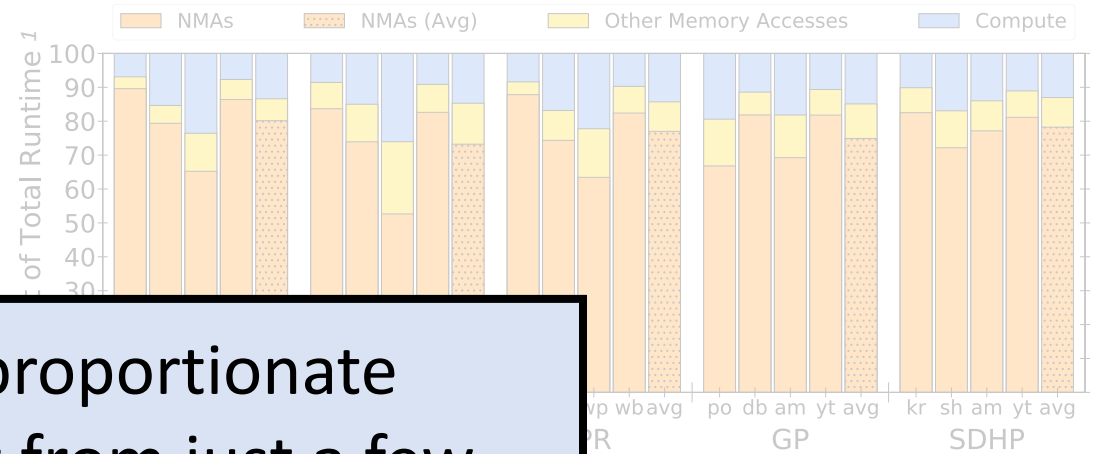
Updates are irregular!

Neighbor Memory Accesses: The Problem

- Irregular accesses experience cache misses
- **Neighbor Memory Accesses (NMAs)**: irregular memory accesses in critical path



Programs see disproportionate performance impact from just a few **NMAs**. Our work seeks to address these.



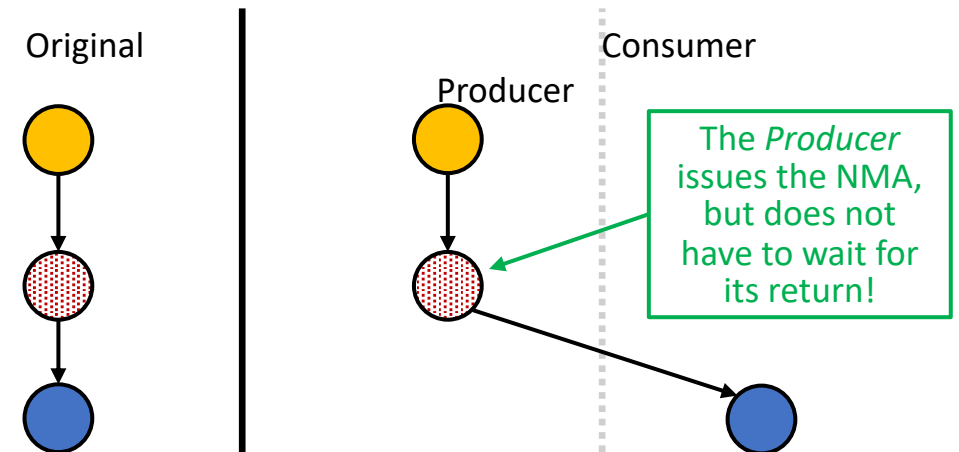
¹ Runtimes measured on a simulated in-order core.

Making NMAs Asynchronous Through Dependence Graph Analysis

- Graph application kernels can be sliced along their NMAs to create **asynchronous accesses**
 - Producer handles all instructions necessary to determine address of NMA and issues access
 - Consumer handles all instructions that depend on data returned from NMA
- **Dependence graph analysis** can identify NMAs as pointer indirect accesses in innermost loop
 - Competitive parallel implementations use RMWs, so NMAs can be loads or RMWs
 - Search for accesses with addresses originating from earlier load

```
for node in frontier:  
    val = process_node(node)  
    for neib in G.neighbors(node):  
        update = update_neib(node_vals, val, neib)  
        if(add_to_frontier(update)):  
            new frontier.push(neib)
```

Iterative, frontier-based graph application



Dependence graph analysis for Producer/Consumer slicing

Compilation Procedure

1

Find NMAs

Identify pointer indirect accesses in innermost kernel loop

Search for loads and RMWs where address originates from earlier load

2

Slice Program:

Create Producer/Consumer slices based on NMA locations

Producer issues NMAs, Consumer uses their data

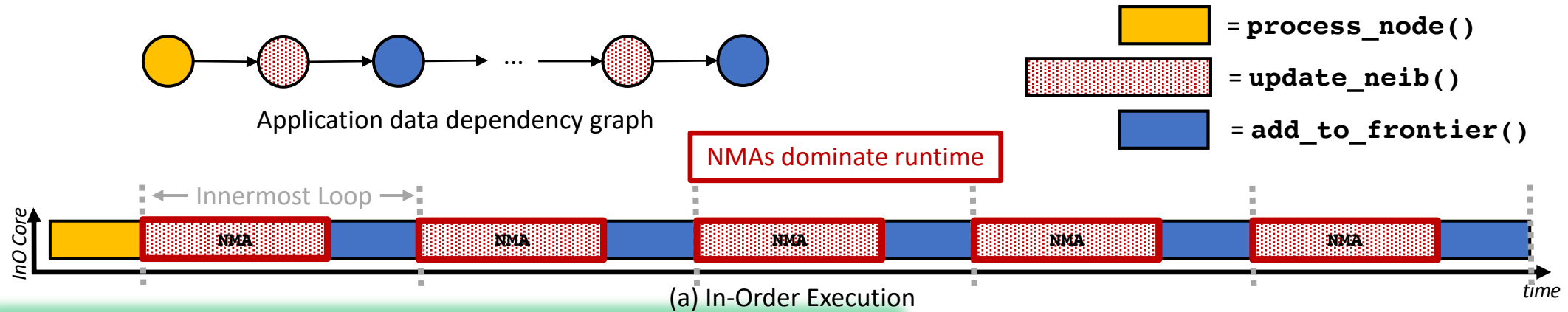
3

Check Memory Safety:

Ensure Producer, Consumer, and AAB do not have conflicts

Writes performed in disjoint memory regions

GraphAttack! Tolerates Latency in Graph Applications by Making NMAs Asynchronous

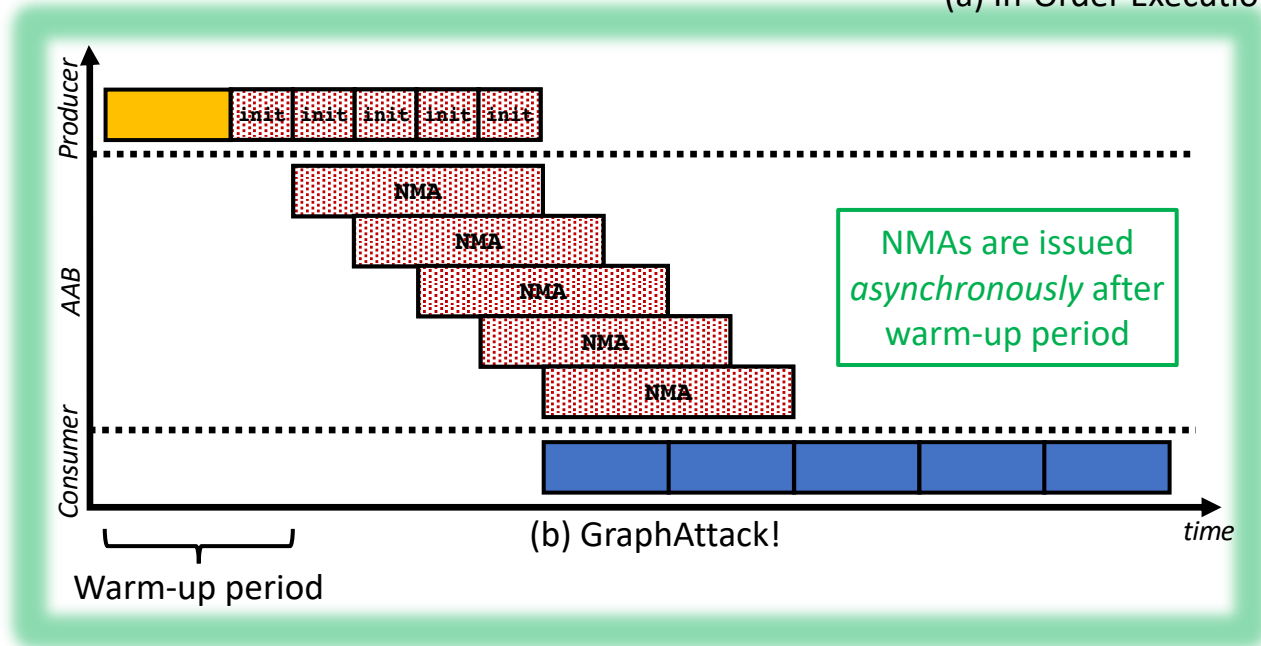


GraphAttack! eliminates dependencies on NMAs, so decoupling achieves latency tolerance on graph applications!

```

for node in frontier:
    val = process_node(node)
    for neib in G.neighbors(node):
        update = update_neib(node_vals, val, neib)
        if add_to_frontier(update):
            new_frontier.push(neib)
    
```

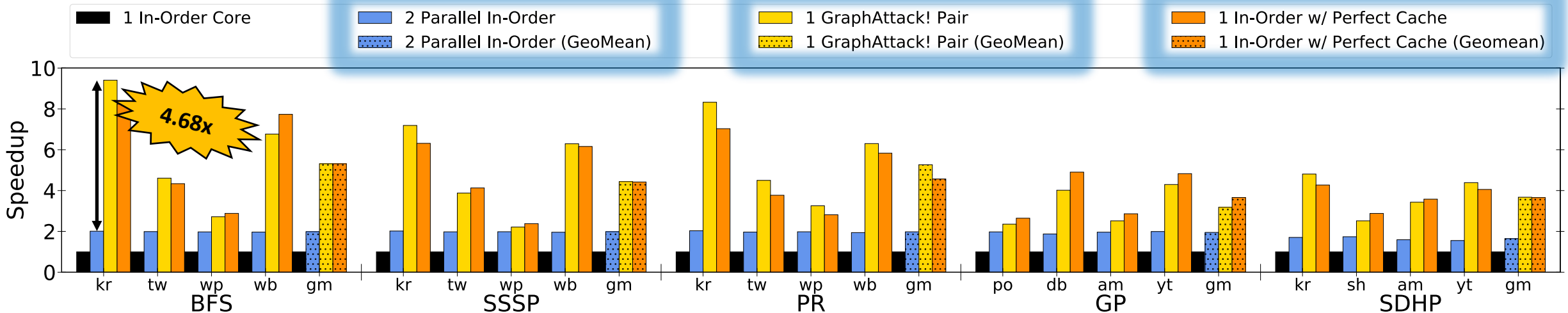
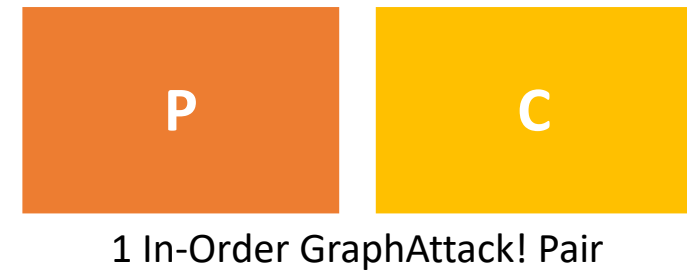
Iterative, frontier-based graph application template



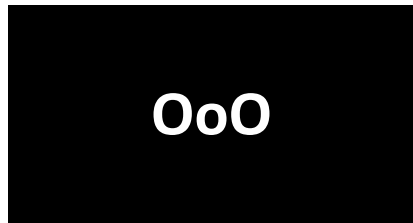
GraphAttack! Tolerates Latency and Better Utilizes In-Order Cores



vs.



GraphAttack! is Scalable and Energy-Efficient



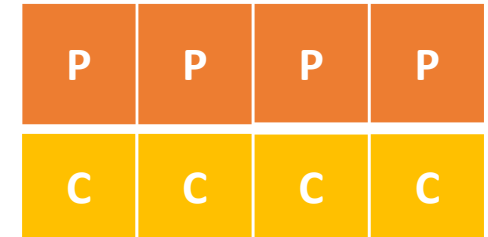
1 OoO Core

vs.

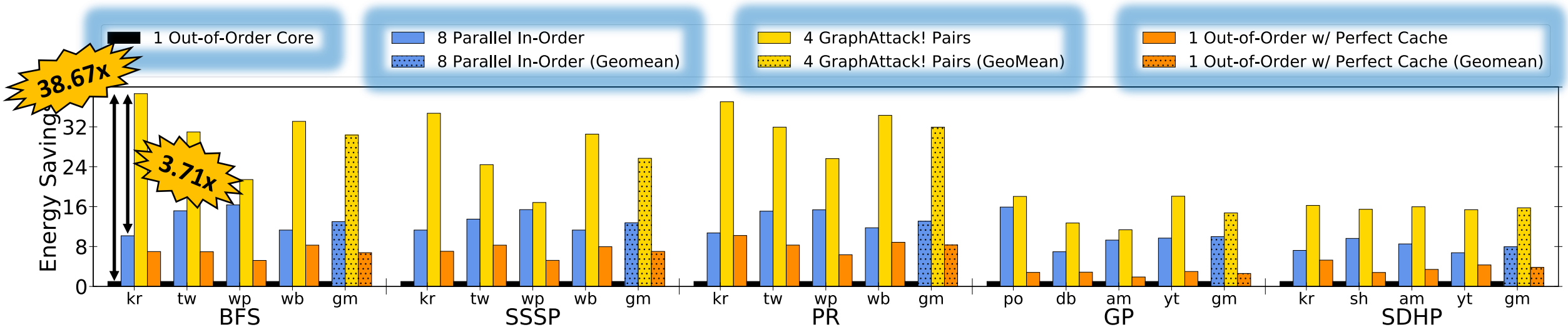


8 InO cores

vs.

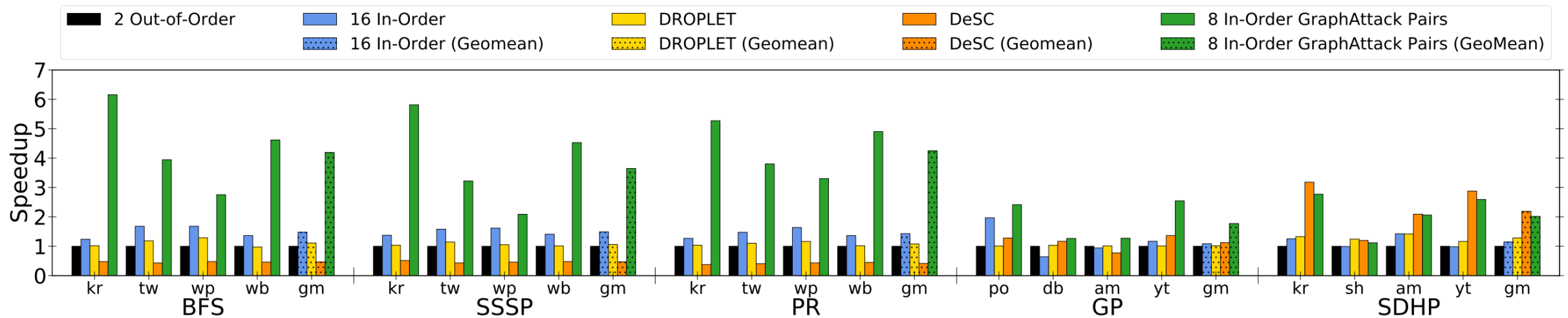


4 InO GraphAttack Pairs



Performance results in GraphAttack (taco '21)

Access and Execute are in order cores



Enjoy the holiday break!

- Get paired up for the assignment and start on it early
- Start working on paper review!