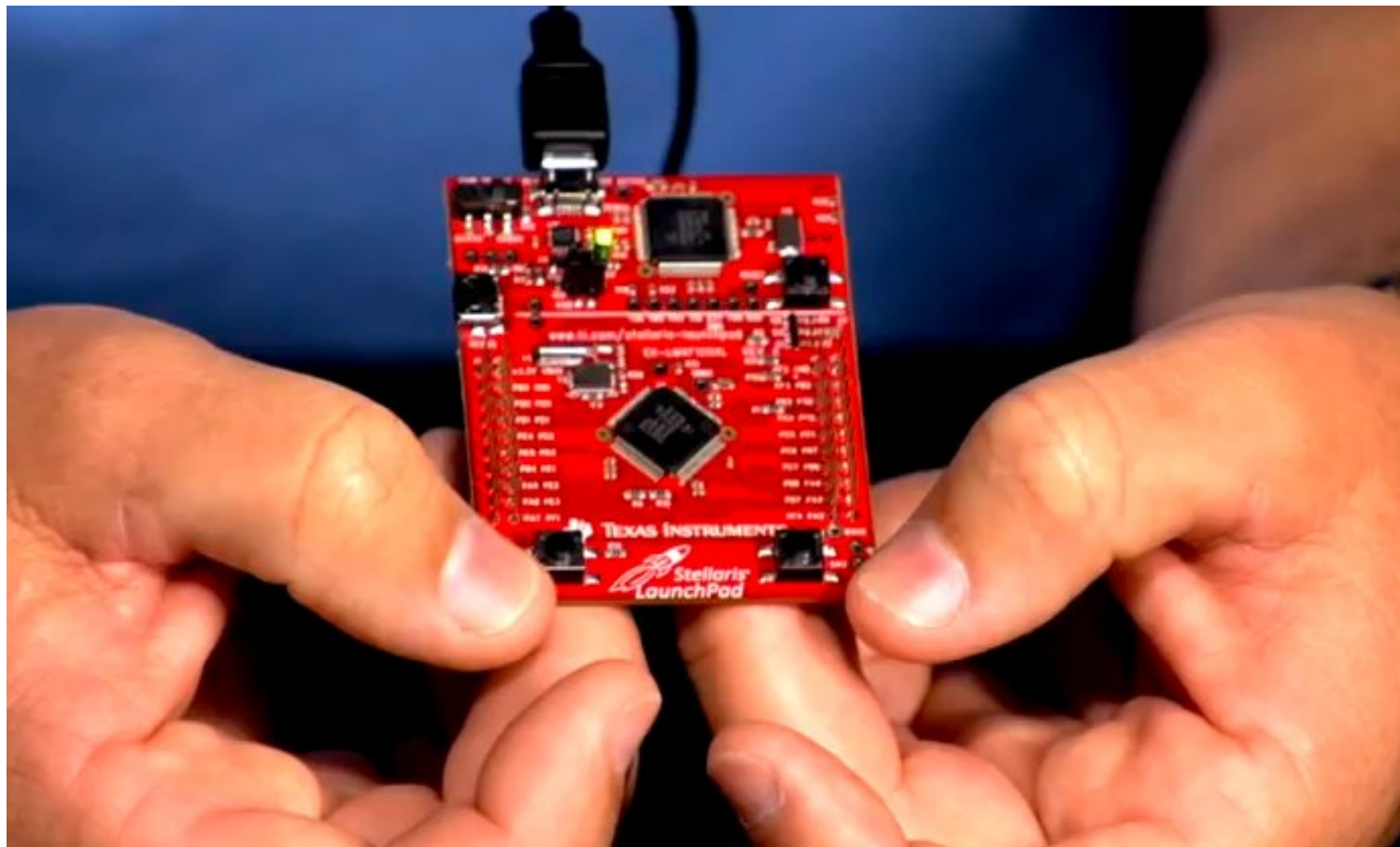


CPU-specific optimization

Example of a target CPU core:
ARM Cortex-M4F core inside
LM4F120H5QR microcontroller
in Stellaris LM4F120 Launchpad.



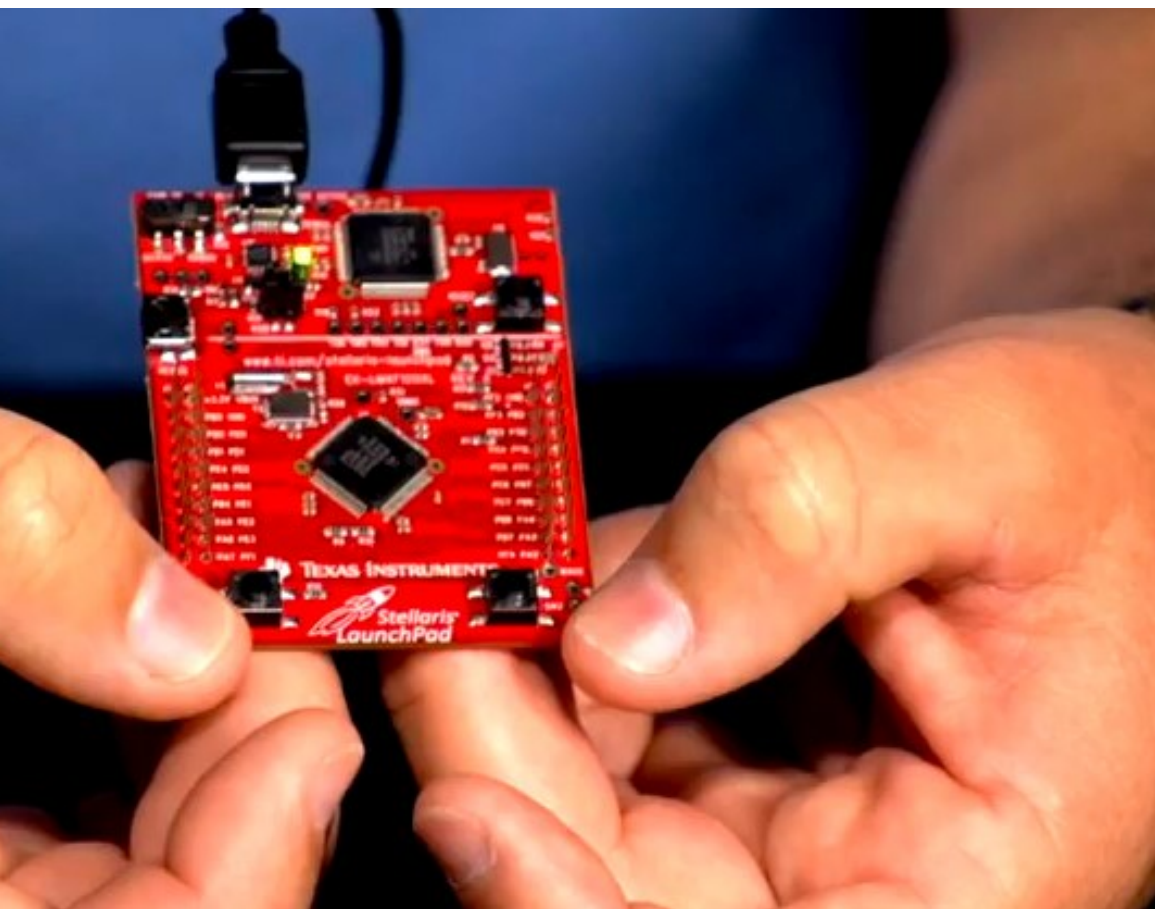
Example of a function
that we want to optimize:
adding 1000 integers mod 2^{32} .

Reference implementation:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += x[i];
    return result;
}
```

Specific optimization

of a target CPU core:
Cortex-M4F core inside
STM32F405QR microcontroller
on Texas Instruments LM4F120 Launchpad.



1

Example of a function
that we want to optimize:
adding 1000 integers mod 2^{32} .

Reference implementation:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += x[i];
    return result;
}
```

2

Counting

static v

*const

= (vo

...

int bef

int res

int aft

UARTpri

result

Output s

Change

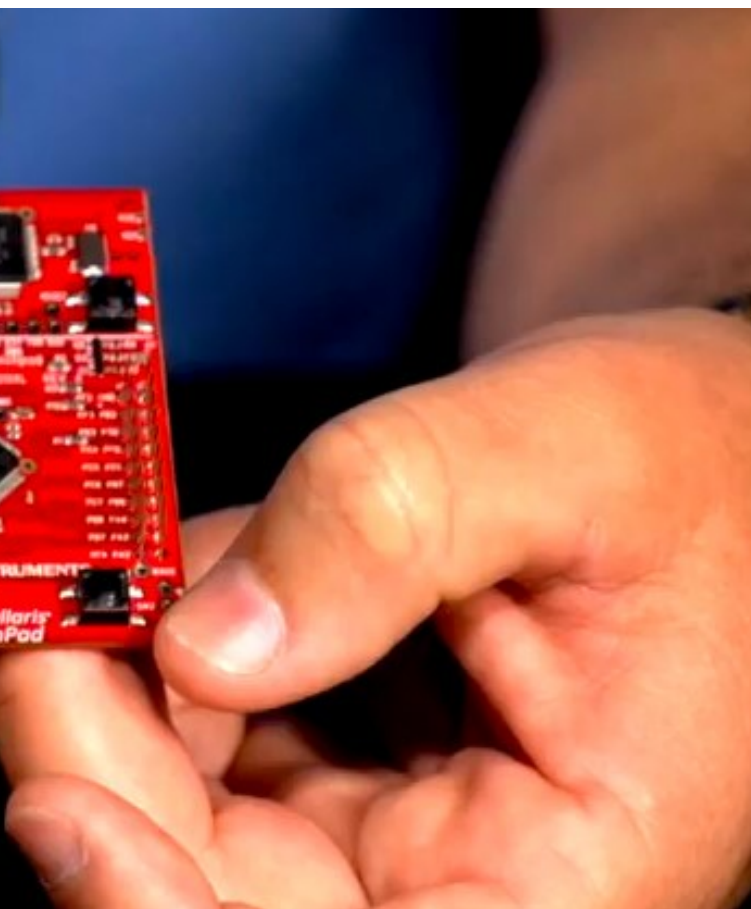
Optimization

Get CPU core:

core inside

microcontroller

20 Launchpad.



1

Example of a function

that we want to optimize:

adding 1000 integers mod 2^{32} .

Reference implementation:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += x[i];
    return result;
}
```

2

Counting cycles:

```
static volatile
    *const DWT_CYC
    = (void *) 0xE
...

```

```
int beforesum =
int result = sum
int aftersum = *
UARTprintf("sum
    result, aftersu

```

Output shows 801

Change 1000 to 50

1

Example of a function
that we want to optimize:
adding 1000 integers mod 2^{32} .

Reference implementation:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += x[i];
    return result;
}
```

2

Counting cycles:

```
static volatile unsigned
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...
```

```
int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum - beforesum);
```

Output shows 8012 cycles.
Change 1000 to 500: 4012.

Example of a function
that we want to optimize:
adding 1000 integers mod 2^{32} .

Reference implementation:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += x[i];
    return result;
}
```

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum - beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

e of a function
want to optimize:
1000 integers mod 2^{32} .

ce implementation:

```
(int *x)
```

```
result = 0;
```

```
;
```

```
i = 0; i < 1000; ++i)
```

```
ult += x[i];
```

```
n result;
```

2

Counting cycles:

```
static volatile unsigned int  
    *const DWT_CYCCNT  
    = (void *) 0xE0001004;
```

```
...
```

```
int beforesum = *DWT_CYCCNT;
```

```
int result = sum(x);
```

```
int aftersum = *DWT_CYCCNT;
```

```
UARTprintf("sum %d %d\n",
```

```
    result, aftersum - beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

3

“Okay, &
Um, are
really th

tion
optimize:
ers mod 2^{32} .

entation:

```
;
1000; ++i)
i];
```

2

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum - beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

3

“Okay, 8 cycles per
Um, are microcont
really this slow at

2

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

3

“Okay, 8 cycles per addition
Um, are microcontrollers
really this slow at addition?”

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -Os: 8012 cycles.

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -Os: 8012 cycles.

Try -O1: 8012 cycles.

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -0s: 8012 cycles.

Try -01: 8012 cycles.

Try -02: 8012 cycles.

Counting cycles:

```
static volatile unsigned int
    *const DWT_CYCCNT
    = (void *) 0xE0001004;
...

int beforesum = *DWT_CYCCNT;
int result = sum(x);
int aftersum = *DWT_CYCCNT;
UARTprintf("sum %d %d\n",
    result, aftersum-beforesum);
```

Output shows 8012 cycles.

Change 1000 to 500: 4012.

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -0s: 8012 cycles.

Try -01: 8012 cycles.

Try -02: 8012 cycles.

Try -03: 8012 cycles.

g cycles:

```
volatile unsigned int  
t DWT_CYCCNT
```

```
id *) 0xE0001004;
```

```
oresum = *DWT_CYCCNT;
```

```
ult = sum(x);
```

```
ersum = *DWT_CYCCNT;
```

```
ntf("sum %d %d\n",
```

```
t, aftersum-beforesum);
```

shows 8012 cycles.

1000 to 500: 4012.

3

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -0s: 8012 cycles.

Try -01: 8012 cycles.

Try -02: 8012 cycles.

Try -03: 8012 cycles.

4

Try mov

```
int sum
```

```
{
```

```
int r
```

```
int i
```

```
for (i
```

```
    res
```

```
return
```

```
}
```

```
unsigned int
CNT
0001004;

*DWT_CYCCNT;
(x);
DWT_CYCCNT;
%d %d\n",
m-beforesum);

2 cycles.
00: 4012.
```

3

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -0s: 8012 cycles.

Try -01: 8012 cycles.

Try -02: 8012 cycles.

Try -03: 8012 cycles.

4

Try moving the po

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i <
        result += *x;
    return result;
}
```


3

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -O0: 8012 cycles.

Try -O1: 8012 cycles.

Try -O2: 8012 cycles.

Try -O3: 8012 cycles.

4

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -0s: 8012 cycles.

Try -01: 8012 cycles.

Try -02: 8012 cycles.

Try -03: 8012 cycles.

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

“Okay, 8 cycles per addition.
Um, are microcontrollers
really this slow at addition?”

Bad approach:

Apply random “optimizations”
(and tweak compiler options)
until you get bored/frustrated.
Keep the fastest results.

Try -O0: 8012 cycles.

Try -O1: 8012 cycles.

Try -O2: 8012 cycles.

Try -O3: 8012 cycles.

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

8010 cycles.

3 cycles per addition.

microcontrollers

is slow at addition?"

approach:

random "optimizations"

(break compiler options)

you get bored/frustrated.

the fastest results.

: 8012 cycles.

: 8012 cycles.

: 8012 cycles.

: 8012 cycles.

4

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

8010 cycles.

5

Try counting

```
int sum
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

4

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

8010 cycles.

5

Try counting down

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x;
    return result;
}
```

4

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

8010 cycles.

5

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

8010 cycles.

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

Try moving the pointer:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; ++i)
        result += *x++;
    return result;
}
```

8010 cycles.

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

8010 cycles.

ing the pointer:

```
(int *x)

result = 0;

;

i = 0; i < 1000; ++i)

result += *x++;

n result;
```

cles.

5

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

8010 cycles.

6

Try using

```
int sum
{
    int r
    int *y
    while
        res
    return
}
```

printer:

```
;
for (i = 1000; ++i)
    ++;
```

5

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

8010 cycles.

6

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1;
    while (x != y)
        result += *x++;
    return result;
}
```

5

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

8010 cycles.

6

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

8010 cycles.

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

Try counting down:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 1000; i > 0; --i)
        result += *x++;
    return result;
}
```

8010 cycles.

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

8010 cycles.

Counting down:

```
(int *x)
```

```
result = 0;
```

```
;
```

```
i = 1000; i > 0; --i)
```

```
result += *x++;
```

```
return result;
```

cycles.

6

Try using an end pointer:

```
int sum(int *x)
```

```
{
```

```
    int result = 0;
```

```
    int *y = x + 1000;
```

```
    while (x != y)
```

```
        result += *x++;
```

```
    return result;
```

```
}
```

8010 cycles.

7

Back to

```
int sum
```

```
{
```

```
    int re
```

```
    int i
```

```
    for (i
```

```
        resu
```

```
        resu
```

```
}
```

```
return
```

```
}
```

6

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

8010 cycles.

7

Back to original.

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i++)
        result += x[i];
    return result;
}
```

6

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

8010 cycles.

7

Back to original. Try unrolli

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i++)
        result += x[i];
    return result;
}
```


Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

8010 cycles.

Back to original. Try unrolling:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 2) {
        result += x[i];
        result += x[i + 1];
    }
    return result;
}
```

Try using an end pointer:

```
int sum(int *x)
{
    int result = 0;
    int *y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

8010 cycles.

Back to original. Try unrolling:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 2) {
        result += x[i];
        result += x[i + 1];
    }
    return result;
}
```

5016 cycles.

g an end pointer:

```
(int *x)
{
    result = 0;
    y = x + 1000;
    while (x != y)
        result += *x++;
    return result;
}
```

cles.

7

Back to original. Try unrolling:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 2) {
        result += x[i];
        result += x[i + 1];
    }
    return result;
}
```

5016 cycles.

8

```
int sum
{
    int re
    int i
    for (
        resu
        resu
        resu
        resu
        resu
    }
    return
}
```

pointer:

;

000;

++;

7

Back to original. Try unrolling:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 2) {
        result += x[i];
        result += x[i + 1];
    }
    return result;
}
```

5016 cycles.

8

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i <
        result += x[
        result += x[
        result += x[
        result += x[
        result += x[
    }
    return result;
}
```

7

Back to original. Try unrolling:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 2) {
        result += x[i];
        result += x[i + 1];
    }
    return result;
}
```

5016 cycles.

8

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 1) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

Back to original. Try unrolling:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 2) {
        result += x[i];
        result += x[i + 1];
    }
    return result;
}
```

5016 cycles.

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

original. Try unrolling:

```
(int *x)
```

```
result = 0;
```

```
;
```

```
for (i = 0; i < 1000; i += 2) {
```

```
    result += x[i];
```

```
    result += x[i + 1];
```

```
return result;
```

cycles.

8

```
int sum(int *x)
```

```
{
```

```
    int result = 0;
```

```
    int i;
```

```
    for (i = 0; i < 1000; i += 5) {
```

```
        result += x[i];
```

```
        result += x[i + 1];
```

```
        result += x[i + 2];
```

```
        result += x[i + 3];
```

```
        result += x[i + 4];
```

```
    }
```

```
    return result;
```

```
}
```

9

4016 cycles

Try unrolling:

```
;
for (i = 0; i < 1000; i += 2) {
    result += x[i];
    result += x[i + 1];
}
```

8

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

9

4016 cycles. Are v

8

ng:

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

= 2) {

9

4016 cycles. Are we done no

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

4016 cycles. Are we done now?

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

4016 cycles. Are we done now?

Most random “optimizations”
that we tried seem useless.

Can spend time trying more.

Does frustration level tell us
that we’re close to optimal?

```
int sum(int *x)
{
    int result = 0;
    int i;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

4016 cycles. Are we done now?

Most random “optimizations”
that we tried seem useless.

Can spend time trying more.

Does frustration level tell us
that we’re close to optimal?

Good approach:

Figure out lower bound for
cycles spent on arithmetic etc.

Understand gap between
lower bound and observed time.

Let’s try this approach.

```
(int *x)
{
    result = 0;
    for (i = 0; i < 1000; i += 5) {
        result += x[i];
        result += x[i + 1];
        result += x[i + 2];
        result += x[i + 3];
        result += x[i + 4];
    }
    return result;
}
```

9

4016 cycles. Are we done now?

Most random “optimizations”
that we tried seem useless.

Can spend time trying more.

Does frustration level tell us
that we’re close to optimal?

Good approach:

Figure out lower bound for
cycles spent on arithmetic etc.

Understand gap between
lower bound and observed time.

Let’s try this approach.

10

Find “A
Technical
Rely on
M4F =
Manual
“implem
architect
Points to
Architec
which de
e.g., “A
First ma
ADD tal

```

;
1000; i += 5) {
i];
i + 1];
i + 2];
i + 3];
i + 4];

```

4016 cycles. Are we done now?

Most random “optimizations”
that we tried seem useless.

Can spend time trying more.

Does frustration level tell us
that we’re close to optimal?

Good approach:

Figure out lower bound for
cycles spent on arithmetic etc.

Understand gap between
lower bound and observed time.

Let’s try this approach.

Find “ARM Cortex

Technical Referenc

Rely on Wikipedia

M4F = M4 + float

Manual says that

“implements the A

architecture profile

Points to the “AR

Architecture Refer

which defines instr

e.g., “ADD” for 3

First manual says

ADD takes just 1

4016 cycles. Are we done now?

Most random “optimizations”
that we tried seem useless.

Can spend time trying more.

Does frustration level tell us
that we’re close to optimal?

Good approach:

Figure out lower bound for
cycles spent on arithmetic etc.

Understand gap between
lower bound and observed time.

Let’s try this approach.

Find “ARM Cortex-M4 Proc
Technical Reference Manual
Rely on Wikipedia comment
M4F = M4 + floating-point
Manual says that Cortex-M4
“implements the ARMv7E-M
architecture profile” .

Points to the “ARMv7-M
Architecture Reference Man
which defines instructions:
e.g., “ADD” for 32-bit addit

First manual says that
ADD takes just 1 cycle.

4016 cycles. Are we done now?

Most random “optimizations” that we tried seem useless.

Can spend time trying more.

Does frustration level tell us that we’re close to optimal?

Good approach:

Figure out lower bound for cycles spent on arithmetic etc.

Understand gap between lower bound and observed time.

Let’s try this approach.

Find “ARM Cortex-M4 Processor Technical Reference Manual” .

Rely on Wikipedia comment that $M4F = M4 + \text{floating-point unit}$.

Manual says that Cortex-M4 “implements the ARMv7E-M architecture profile” .

Points to the “ARMv7-M Architecture Reference Manual” , which defines instructions: e.g., “ADD” for 32-bit addition.

First manual says that ADD takes just 1 cycle.

cles. Are we done now?

andom “optimizations”
tried seem useless.

nd time trying more.

ustration level tell us
re close to optimal?

pproach:

ut lower bound for
pent on arithmetic etc.

and gap between

ound and observed time.

y this approach.

10

Find “ARM Cortex-M4 Processor
Technical Reference Manual” .

Rely on Wikipedia comment that
 $M4F = M4 + \text{floating-point unit}$.

Manual says that Cortex-M4
“implements the ARMv7E-M
architecture profile” .

Points to the “ARMv7-M
Architecture Reference Manual” ,
which defines instructions:
e.g., “ADD” for 32-bit addition.

First manual says that
ADD takes just 1 cycle.

11

Inputs a
“integer
has 16 i
special-p
and “pro

Each ele
be “load

Basic loa
Manual

a note a
Then mo
instructi
address
then it s

we done now?

optimizations”

useless.

ying more.

level tell us

o optimal?

ound for

ithmetic etc.

etween

observed time.

oach.

Find “ARM Cortex-M4 Processor Technical Reference Manual” .

Rely on Wikipedia comment that $M4F = M4 + \text{floating-point unit}$.

Manual says that Cortex-M4 “implements the ARMv7E-M architecture profile” .

Points to the “ARMv7-M Architecture Reference Manual” , which defines instructions: e.g., “ADD” for 32-bit addition.

First manual says that ADD takes just 1 cycle.

Inputs and output “integer registers”

has 16 integer reg

special-purpose “s

and “program cou

Each element of x

be “loaded” into a

Basic load instruct

Manual says 2 cyc

a note about “pip

Then more explan

instruction is also

address not based

then it saves 1 cyc

Find “ARM Cortex-M4 Processor Technical Reference Manual” .

Rely on Wikipedia comment that $M4F = M4 + \text{floating-point unit}$.

Manual says that Cortex-M4 “implements the ARMv7E-M architecture profile” .

Points to the “ARMv7-M Architecture Reference Manual” , which defines instructions: e.g., “ADD” for 32-bit addition.

First manual says that ADD takes just 1 cycle.

Inputs and output of ADD are “integer registers” . ARMv7-M has 16 integer registers, including special-purpose “stack pointer” and “program counter” .

Each element of x array needs to be “loaded” into a register.

Basic load instruction: LDR. Manual says 2 cycles but adds a note about “pipelining” .

Then more explanation: if next instruction is also LDR (with address not based on first LDR) then it saves 1 cycle.

Find “ARM Cortex-M4 Processor Technical Reference Manual” .

Rely on Wikipedia comment that $M4F = M4 + \text{floating-point unit}$.

Manual says that Cortex-M4

“implements the ARMv7E-M architecture profile” .

Points to the “ARMv7-M Architecture Reference Manual” , which defines instructions:

e.g., “ADD” for 32-bit addition.

First manual says that ADD takes just 1 cycle.

Inputs and output of ADD are “integer registers” . ARMv7-M has 16 integer registers, including special-purpose “stack pointer” and “program counter” .

Each element of x array needs to be “loaded” into a register.

Basic load instruction: LDR.

Manual says 2 cycles but adds a note about “pipelining” .

Then more explanation: if next instruction is also LDR (with address not based on first LDR) then it saves 1 cycle.

ARM Cortex-M4 Processor
Reference Manual”.

Wikipedia comment that
M4 + floating-point unit.

says that Cortex-M4
implements the ARMv7E-M
architecture profile”.

to the “ARMv7-M
Architecture Reference Manual”,
defines instructions:

“ADD” for 32-bit addition.

Manual says that
takes just 1 cycle.

Inputs and output of ADD are
“integer registers”. ARMv7-M
has 16 integer registers, including
special-purpose “stack pointer”
and “program counter”.

Each element of x array needs to
be “loaded” into a register.

Basic load instruction: LDR.

Manual says 2 cycles but adds
a note about “pipelining”.

Then more explanation: if next
instruction is also LDR (with
address not based on first LDR)
then it saves 1 cycle.

n consec
takes on
(“more
pipelined

Can ach
in other
but noth

Lower bo
 $2n + 1$ c
 n cycles

Why obs
non-cons
costs of

Cortex-M4 Processor
Reference Manual”.

comment that
the floating-point unit.

Cortex-M4

ARMv7E-M

”.

ARMv7-M

Reference Manual”,

instructions:

2-bit addition.

that

cycle.

Inputs and output of ADD are
“integer registers”. ARMv7-M
has 16 integer registers, including
special-purpose “stack pointer”
and “program counter”.

Each element of x array needs to
be “loaded” into a register.

Basic load instruction: LDR.

Manual says 2 cycles but adds
a note about “pipelining”.

Then more explanation: if next
instruction is also LDR (with
address not based on first LDR)
then it saves 1 cycle.

n consecutive LDR
takes only $n + 1$ cycles
(“more multiple LDR
pipelined together”)

Can achieve this savings
in other ways (LDR
but nothing seems

Lower bound for n
 $2n + 1$ cycles, including
 n cycles of arithmetic

Why observed time for
non-consecutive LDR
costs of manipulating

processor
".
that
unit.

Inputs and output of ADD are
"integer registers". ARMv7-M
has 16 integer registers, including
special-purpose "stack pointer"
and "program counter".

Each element of x array needs to
be "loaded" into a register.

Basic load instruction: LDR.

Manual says 2 cycles but adds
a note about "pipelining".

Then more explanation: if next
instruction is also LDR (with
address not based on first LDR)
then it saves 1 cycle.

n consecutive LDRs
takes only $n + 1$ cycles
("more multiple LDRs can be
pipelined together").

Can achieve this speed
in other ways (LDRD, LDM),
but nothing seems faster.

Lower bound for n LDR + n
 $2n + 1$ cycles, including
 n cycles of arithmetic.

Why observed time is higher
non-consecutive LDRs;
costs of manipulating i .

Inputs and output of ADD are “integer registers”. ARMv7-M has 16 integer registers, including special-purpose “stack pointer” and “program counter”.

Each element of x array needs to be “loaded” into a register.

Basic load instruction: LDR.

Manual says 2 cycles but adds a note about “pipelining”.

Then more explanation: if next instruction is also LDR (with address not based on first LDR) then it saves 1 cycle.

n consecutive LDRs takes only $n + 1$ cycles (“more multiple LDRs can be pipelined together”).

Can achieve this speed in other ways (LDRD, LDM) but nothing seems faster.

Lower bound for n LDR + n ADD: $2n + 1$ cycles, including n cycles of arithmetic.

Why observed time is higher: non-consecutive LDRs; costs of manipulating i .

and output of ADD are registers". ARMv7-M integer registers, including purpose "stack pointer" program counter".

element of x array needs to be loaded" into a register.

load instruction: LDR.

says 2 cycles but adds about "pipelining".

more explanation: if next instruction is also LDR (with address not based on first LDR) it saves 1 cycle.

n consecutive LDRs takes only $n + 1$ cycles ("more multiple LDRs can be pipelined together").

Can achieve this speed in other ways (LDRD, LDM) but nothing seems faster.

Lower bound for n LDR + n ADD: $2n + 1$ cycles, including n cycles of arithmetic.

Why observed time is higher: non-consecutive LDRs; costs of manipulating i .

2281 cycles

$y = x +$

$p = x$

result =

loop:

$x_{i9} =$

$x_{i8} =$

$x_{i7} =$

$x_{i6} =$

$x_{i5} =$

$x_{i4} =$

$x_{i3} =$

$x_{i2} =$

of ADD are
 . ARMv7-M
 isters, including
 tack pointer”
 nter”.

array needs to
 a register.

tion: LDR.

les but adds
 elining”.

ation: if next

LDR (with

on first LDR)

cle.

n consecutive LDRs
 takes only $n + 1$ cycles
 (“more multiple LDRs can be
 pipelined together”).

Can achieve this speed
 in other ways (LDRD, LDM)
 but nothing seems faster.

Lower bound for n LDR + n ADD:
 $2n + 1$ cycles, including
 n cycles of arithmetic.

Why observed time is higher:
 non-consecutive LDRs;
 costs of manipulating i .

2281 cycles using

```
y = x + 4000
```

```
p = x
```

```
result = 0
```

```
loop:
```

```
xi9 = *(uint32
```

```
xi8 = *(uint32
```

```
xi7 = *(uint32
```

```
xi6 = *(uint32
```

```
xi5 = *(uint32
```

```
xi4 = *(uint32
```

```
xi3 = *(uint32
```

```
xi2 = *(uint32
```

n consecutive LDRs
 takes only $n + 1$ cycles
 (“more multiple LDRs can be
 pipelined together”).

Can achieve this speed
 in other ways (LDRD, LDM)
 but nothing seems faster.

Lower bound for n LDR + n ADD:
 $2n + 1$ cycles, including
 n cycles of arithmetic.

Why observed time is higher:
 non-consecutive LDRs;
 costs of manipulating i .

2281 cycles using `ldr.w`:

```
y = x + 4000
```

```
p = x
```

```
result = 0
```

```
loop:
```

```
xi9 = *(uint32 *) (p +
```

```
xi8 = *(uint32 *) (p +
```

```
xi7 = *(uint32 *) (p +
```

```
xi6 = *(uint32 *) (p +
```

```
xi5 = *(uint32 *) (p +
```

```
xi4 = *(uint32 *) (p +
```

```
xi3 = *(uint32 *) (p +
```

```
xi2 = *(uint32 *) (p +
```

n consecutive LDRs
takes only $n + 1$ cycles
(“more multiple LDRs can be
pipelined together”).

Can achieve this speed
in other ways (LDRD, LDM)
but nothing seems faster.

Lower bound for n LDR + n ADD:
 $2n + 1$ cycles, including
 n cycles of arithmetic.

Why observed time is higher:
non-consecutive LDRs;
costs of manipulating i .

2281 cycles using `ldr.w`:

```
y = x + 4000
```

```
p = x
```

```
result = 0
```

```
loop:
```

```
xi9 = *(uint32 *) (p + 76)
```

```
xi8 = *(uint32 *) (p + 72)
```

```
xi7 = *(uint32 *) (p + 68)
```

```
xi6 = *(uint32 *) (p + 64)
```

```
xi5 = *(uint32 *) (p + 60)
```

```
xi4 = *(uint32 *) (p + 56)
```

```
xi3 = *(uint32 *) (p + 52)
```

```
xi2 = *(uint32 *) (p + 48)
```

secutive LDRs

only $n + 1$ cycles

multiple LDRs can be
"packed together").

to achieve this speed

several ways (LDRD, LDM)

streaming seems faster.

bound for n LDR + n ADD:

$n + 1$ cycles, including
one cycle of arithmetic.

observed time is higher:

for consecutive LDRs;

and for manipulating i .

2281 cycles using `ldr.w`:

```
y = x + 4000
```

```
p = x
```

```
result = 0
```

```
loop:
```

```
    xi9 = *(uint32 *) (p + 76)
```

```
    xi8 = *(uint32 *) (p + 72)
```

```
    xi7 = *(uint32 *) (p + 68)
```

```
    xi6 = *(uint32 *) (p + 64)
```

```
    xi5 = *(uint32 *) (p + 60)
```

```
    xi4 = *(uint32 *) (p + 56)
```

```
    xi3 = *(uint32 *) (p + 52)
```

```
    xi2 = *(uint32 *) (p + 48)
```

```
xi1 =
```

```
xi0 =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
result =
```

```
xi9 =
```

```
xi8 =
```

```
xi7 =
```

2281 cycles using `ldr.w`:

```
y = x + 4000
```

```
p = x
```

```
result = 0
```

```
loop:
```

```
    xi9 = *(uint32 *) (p + 76)
```

```
    xi8 = *(uint32 *) (p + 72)
```

```
    xi7 = *(uint32 *) (p + 68)
```

```
    xi6 = *(uint32 *) (p + 64)
```

```
    xi5 = *(uint32 *) (p + 60)
```

```
    xi4 = *(uint32 *) (p + 56)
```

```
    xi3 = *(uint32 *) (p + 52)
```

```
    xi2 = *(uint32 *) (p + 48)
```

```
xi1 = *(uint32
```

```
xi0 = *(uint32
```

```
result += xi9
```

```
result += xi8
```

```
result += xi7
```

```
result += xi6
```

```
result += xi5
```

```
result += xi4
```

```
result += xi3
```

```
result += xi2
```

```
result += xi1
```

```
result += xi0
```

```
xi9 = *(uint32
```

```
xi8 = *(uint32
```

```
xi7 = *(uint32
```

2281 cycles using ldr.w:

y = x + 4000

p = x

result = 0

loop:

xi9 = *(uint32 *) (p + 76)

xi8 = *(uint32 *) (p + 72)

xi7 = *(uint32 *) (p + 68)

xi6 = *(uint32 *) (p + 64)

xi5 = *(uint32 *) (p + 60)

xi4 = *(uint32 *) (p + 56)

xi3 = *(uint32 *) (p + 52)

xi2 = *(uint32 *) (p + 48)

xi1 = *(uint32 *) (p +

xi0 = *(uint32 *) (p +

result += xi9

result += xi8

result += xi7

result += xi6

result += xi5

result += xi4

result += xi3

result += xi2

result += xi1

result += xi0

xi9 = *(uint32 *) (p +

xi8 = *(uint32 *) (p +

xi7 = *(uint32 *) (p +

2281 cycles using ldr.w:

```
y = x + 4000
```

```
p = x
```

```
result = 0
```

```
loop:
```

```
    xi9 = *(uint32 *) (p + 76)
```

```
    xi8 = *(uint32 *) (p + 72)
```

```
    xi7 = *(uint32 *) (p + 68)
```

```
    xi6 = *(uint32 *) (p + 64)
```

```
    xi5 = *(uint32 *) (p + 60)
```

```
    xi4 = *(uint32 *) (p + 56)
```

```
    xi3 = *(uint32 *) (p + 52)
```

```
    xi2 = *(uint32 *) (p + 48)
```

```
xi1 = *(uint32 *) (p + 44)
```

```
xi0 = *(uint32 *) (p + 40)
```

```
result += xi9
```

```
result += xi8
```

```
result += xi7
```

```
result += xi6
```

```
result += xi5
```

```
result += xi4
```

```
result += xi3
```

```
result += xi2
```

```
result += xi1
```

```
result += xi0
```

```
xi9 = *(uint32 *) (p + 36)
```

```
xi8 = *(uint32 *) (p + 32)
```

```
xi7 = *(uint32 *) (p + 28)
```


cles using ldr.w:

4000

= 0

```
*(uint32 *) (p + 76)
*(uint32 *) (p + 72)
*(uint32 *) (p + 68)
*(uint32 *) (p + 64)
*(uint32 *) (p + 60)
*(uint32 *) (p + 56)
*(uint32 *) (p + 52)
*(uint32 *) (p + 48)
```

14

```
xi1 = *(uint32 *) (p + 44)
xi0 = *(uint32 *) (p + 40)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
xi9 = *(uint32 *) (p + 36)
xi8 = *(uint32 *) (p + 32)
xi7 = *(uint32 *) (p + 28)
```

15

```
xi6 =
xi5 =
xi4 =
xi3 =
xi2 =
xi1 =
xi0 =
result-
result-
result-
result-
result-
result-
result-
result-
```

ldr.w:

```

*) (p + 76)
*) (p + 72)
*) (p + 68)
*) (p + 64)
*) (p + 60)
*) (p + 56)
*) (p + 52)
*) (p + 48)

```

```

xi1 = *(uint32 *) (p + 44)
xi0 = *(uint32 *) (p + 40)

result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0

xi9 = *(uint32 *) (p + 36)
xi8 = *(uint32 *) (p + 32)
xi7 = *(uint32 *) (p + 28)

```

```

xi6 = *(uint32 *) (p + 24)
xi5 = *(uint32 *) (p + 20)
xi4 = *(uint32 *) (p + 16)
xi3 = *(uint32 *) (p + 12)
xi2 = *(uint32 *) (p + 8)
xi1 = *(uint32 *) (p + 4)
xi0 = *(uint32 *) (p)

result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2

```

14

```
xi1 = *(uint32 *) (p + 44)
xi0 = *(uint32 *) (p + 40)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
xi9 = *(uint32 *) (p + 36)
xi8 = *(uint32 *) (p + 32)
xi7 = *(uint32 *) (p + 28)
```

15

```
xi6 = *(uint32 *) (p +
xi5 = *(uint32 *) (p +
xi4 = *(uint32 *) (p +
xi3 = *(uint32 *) (p +
xi2 = *(uint32 *) (p +
xi1 = *(uint32 *) (p +
xi0 = *(uint32 *) p; p
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
```

76)

72)

68)

64)

60)

56)

52)

48)

```
xi1 = *(uint32 *) (p + 44)
xi0 = *(uint32 *) (p + 40)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
xi9 = *(uint32 *) (p + 36)
xi8 = *(uint32 *) (p + 32)
xi7 = *(uint32 *) (p + 28)
```

```
xi6 = *(uint32 *) (p + 24)
xi5 = *(uint32 *) (p + 20)
xi4 = *(uint32 *) (p + 16)
xi3 = *(uint32 *) (p + 12)
xi2 = *(uint32 *) (p + 8)
xi1 = *(uint32 *) (p + 4)
xi0 = *(uint32 *) p; p += 160
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
```

15

```
*(uint32 *) (p + 44)
*(uint32 *) (p + 40)
t += xi9
t += xi8
t += xi7
t += xi6
t += xi5
t += xi4
t += xi3
t += xi2
t += xi1
t += xi0
*(uint32 *) (p + 36)
*(uint32 *) (p + 32)
*(uint32 *) (p + 28)
```

16

```
xi6 = *(uint32 *) (p + 24)
xi5 = *(uint32 *) (p + 20)
xi4 = *(uint32 *) (p + 16)
xi3 = *(uint32 *) (p + 12)
xi2 = *(uint32 *) (p + 8)
xi1 = *(uint32 *) (p + 4)
xi0 = *(uint32 *) p; p += 160
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result
```

15

```

*) (p + 44)
*) (p + 40)

*) (p + 36)
*) (p + 32)
*) (p + 28)

```

```

xi6 = *(uint32 *) (p + 24)
xi5 = *(uint32 *) (p + 20)
xi4 = *(uint32 *) (p + 16)
xi3 = *(uint32 *) (p + 12)
xi2 = *(uint32 *) (p + 8)
xi1 = *(uint32 *) (p + 4)
xi0 = *(uint32 *) p; p += 160

result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2

```

16

```

result += xi1
result += xi0
xi9 = *(uint32 *)
xi8 = *(uint32 *)
xi7 = *(uint32 *)
xi6 = *(uint32 *)
xi5 = *(uint32 *)
xi4 = *(uint32 *)
xi3 = *(uint32 *)
xi2 = *(uint32 *)
xi1 = *(uint32 *)
xi0 = *(uint32 *)

result += xi9
result += xi8
result += xi7

```

15

```
44) xi6 = *(uint32 *) (p + 24)
40) xi5 = *(uint32 *) (p + 20)
xi4 = *(uint32 *) (p + 16)
xi3 = *(uint32 *) (p + 12)
xi2 = *(uint32 *) (p + 8)
xi1 = *(uint32 *) (p + 4)
xi0 = *(uint32 *) p; p += 160
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
36) result += xi4
32) result += xi3
28) result += xi2
```

16

```
result += xi1
result += xi0
xi9 = *(uint32 *) (p -
xi8 = *(uint32 *) (p -
xi7 = *(uint32 *) (p -
xi6 = *(uint32 *) (p -
xi5 = *(uint32 *) (p -
xi4 = *(uint32 *) (p -
xi3 = *(uint32 *) (p -
xi2 = *(uint32 *) (p -
xi1 = *(uint32 *) (p -
xi0 = *(uint32 *) (p -
result += xi9
result += xi8
result += xi7
```

```
xi6 = *(uint32 *) (p + 24)
xi5 = *(uint32 *) (p + 20)
xi4 = *(uint32 *) (p + 16)
xi3 = *(uint32 *) (p + 12)
xi2 = *(uint32 *) (p + 8)
xi1 = *(uint32 *) (p + 4)
xi0 = *(uint32 *) p; p += 160

result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
```

```
result += xi1
result += xi0
xi9 = *(uint32 *) (p - 4)
xi8 = *(uint32 *) (p - 8)
xi7 = *(uint32 *) (p - 12)
xi6 = *(uint32 *) (p - 16)
xi5 = *(uint32 *) (p - 20)
xi4 = *(uint32 *) (p - 24)
xi3 = *(uint32 *) (p - 28)
xi2 = *(uint32 *) (p - 32)
xi1 = *(uint32 *) (p - 36)
xi0 = *(uint32 *) (p - 40)

result += xi9
result += xi8
result += xi7
```


16

```

*(uint32 *) (p + 24)
*(uint32 *) (p + 20)
*(uint32 *) (p + 16)
*(uint32 *) (p + 12)
*(uint32 *) (p + 8)
*(uint32 *) (p + 4)
*(uint32 *) p; p += 160

t += xi9
t += xi8
t += xi7
t += xi6
t += xi5
t += xi4
t += xi3
t += xi2

```

17

```

result += xi1
result += xi0
xi9 = *(uint32 *) (p - 4)
xi8 = *(uint32 *) (p - 8)
xi7 = *(uint32 *) (p - 12)
xi6 = *(uint32 *) (p - 16)
xi5 = *(uint32 *) (p - 20)
xi4 = *(uint32 *) (p - 24)
xi3 = *(uint32 *) (p - 28)
xi2 = *(uint32 *) (p - 32)
xi1 = *(uint32 *) (p - 36)
xi0 = *(uint32 *) (p - 40)

result += xi9
result += xi8
result += xi7

```

16

```

*) (p + 24)
*) (p + 20)
*) (p + 16)
*) (p + 12)
*) (p + 8)
*) (p + 4)
*) p; p += 160

```

```

result += xi1
result += xi0
xi9 = *(uint32 *) (p - 4)
xi8 = *(uint32 *) (p - 8)
xi7 = *(uint32 *) (p - 12)
xi6 = *(uint32 *) (p - 16)
xi5 = *(uint32 *) (p - 20)
xi4 = *(uint32 *) (p - 24)
xi3 = *(uint32 *) (p - 28)
xi2 = *(uint32 *) (p - 32)
xi1 = *(uint32 *) (p - 36)
xi0 = *(uint32 *) (p - 40)

result += xi9
result += xi8
result += xi7

```

17

```

result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
xi9 = *(uint32 *) (p - 4)
xi8 = *(uint32 *) (p - 8)
xi7 = *(uint32 *) (p - 12)
xi6 = *(uint32 *) (p - 16)
xi5 = *(uint32 *) (p - 20)
xi4 = *(uint32 *) (p - 24)
xi3 = *(uint32 *) (p - 28)
xi2 = *(uint32 *) (p - 32)

```

```

24) result += xi1
20) result += xi0
16) xi9 = *(uint32 *) (p - 4)
12) xi8 = *(uint32 *) (p - 8)
8) xi7 = *(uint32 *) (p - 12)
4) xi6 = *(uint32 *) (p - 16)
+= 160 xi5 = *(uint32 *) (p - 20)
xi4 = *(uint32 *) (p - 24)
xi3 = *(uint32 *) (p - 28)
xi2 = *(uint32 *) (p - 32)
xi1 = *(uint32 *) (p - 36)
xi0 = *(uint32 *) (p - 40)

result += xi9
result += xi8
result += xi7

```

```

result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
xi9 = *(uint32 *) (p -
xi8 = *(uint32 *) (p -
xi7 = *(uint32 *) (p -
xi6 = *(uint32 *) (p -
xi5 = *(uint32 *) (p -
xi4 = *(uint32 *) (p -
xi3 = *(uint32 *) (p -
xi2 = *(uint32 *) (p -

```

```
result += xi1
result += xi0
xi9 = *(uint32 *) (p - 4)
xi8 = *(uint32 *) (p - 8)
xi7 = *(uint32 *) (p - 12)
xi6 = *(uint32 *) (p - 16)
xi5 = *(uint32 *) (p - 20)
xi4 = *(uint32 *) (p - 24)
xi3 = *(uint32 *) (p - 28)
xi2 = *(uint32 *) (p - 32)
xi1 = *(uint32 *) (p - 36)
xi0 = *(uint32 *) (p - 40)

result += xi9
result += xi8
result += xi7
```

```
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
xi9 = *(uint32 *) (p - 44)
xi8 = *(uint32 *) (p - 48)
xi7 = *(uint32 *) (p - 52)
xi6 = *(uint32 *) (p - 56)
xi5 = *(uint32 *) (p - 60)
xi4 = *(uint32 *) (p - 64)
xi3 = *(uint32 *) (p - 68)
xi2 = *(uint32 *) (p - 72)
```


17

```

result += xi6
result += xi5
*) (p - 4)
result += xi4
*) (p - 8)
result += xi3
*) (p - 12)
result += xi2
*) (p - 16)
result += xi1
*) (p - 20)
result += xi0
*) (p - 24)
xi9 = *(uint32 *) (p - 44)
*) (p - 28)
xi8 = *(uint32 *) (p - 48)
*) (p - 32)
xi7 = *(uint32 *) (p - 52)
*) (p - 36)
xi6 = *(uint32 *) (p - 56)
*) (p - 40)
xi5 = *(uint32 *) (p - 60)
xi4 = *(uint32 *) (p - 64)
xi3 = *(uint32 *) (p - 68)
xi2 = *(uint32 *) (p - 72)

```

18

```

xi1 = *(uint32 *) (p - 76)
xi0 = *(uint32 *) (p - 80)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
= ?
goto loop if !=

```

```

result += xi6
result += xi5
4) result += xi4
8) result += xi3
12) result += xi2
16) result += xi1
20) result += xi0
24) xi9 = *(uint32 *) (p - 44)
28) xi8 = *(uint32 *) (p - 48)
32) xi7 = *(uint32 *) (p - 52)
36) xi6 = *(uint32 *) (p - 56)
40) xi5 = *(uint32 *) (p - 60)
xi4 = *(uint32 *) (p - 64)
xi3 = *(uint32 *) (p - 68)
xi2 = *(uint32 *) (p - 72)

```

```

xi1 = *(uint32 *) (p -
xi0 = *(uint32 *) (p -
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
=? p - y
goto loop if !=

```


18

```

t += xi6
t += xi5
t += xi4
t += xi3
t += xi2
t += xi1
t += xi0
*(uint32 *) (p - 44)
*(uint32 *) (p - 48)
*(uint32 *) (p - 52)
*(uint32 *) (p - 56)
*(uint32 *) (p - 60)
*(uint32 *) (p - 64)
*(uint32 *) (p - 68)
*(uint32 *) (p - 72)

```

19

```

xi1 = *(uint32 *) (p - 76)
xi0 = *(uint32 *) (p - 80)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0
                                =? p - y
goto loop if !=

```

Wikiped

even per

optimizi

performa

```
xi1 = *(uint32 *) (p - 76)
```

```
xi0 = *(uint32 *) (p - 80)
```

```
result += xi9
```

```
result += xi8
```

```
result += xi7
```

```
result += xi6
```

```
result += xi5
```

```
result += xi4
```

```
result += xi3
```

```
result += xi2
```

```
result += xi1
```

```
result += xi0
```

```
=? p - y
```

```
goto loop if !=
```

Wikipedia: “By th
 even performance
 optimizing compile
 performance of hu

```
*) (p - 44)
```

```
*) (p - 48)
```

```
*) (p - 52)
```

```
*) (p - 56)
```

```
*) (p - 60)
```

```
*) (p - 64)
```

```
*) (p - 68)
```

```
*) (p - 72)
```

18

```
xi1 = *(uint32 *) (p - 76)
```

```
xi0 = *(uint32 *) (p - 80)
```

```
result += xi9
```

```
result += xi8
```

```
result += xi7
```

```
result += xi6
```

```
result += xi5
```

```
44) result += xi4
```

```
48) result += xi3
```

```
52) result += xi2
```

```
56) result += xi1
```

```
60) result += xi0
```

```
64)
68)          =? p - y
```

```
72) goto loop if !=
```

19

Wikipedia: “By the late 1990s, even performance sensitive code optimizing compilers exceeded the performance of human experts.”

```

xi1 = *(uint32 *) (p - 76)
xi0 = *(uint32 *) (p - 80)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0

```

```

      =? p - y

```

```

goto loop if !=

```

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

```

xi1 = *(uint32 *) (p - 76)
xi0 = *(uint32 *) (p - 80)
result += xi9
result += xi8
result += xi7
result += xi6
result += xi5
result += xi4
result += xi3
result += xi2
result += xi1
result += xi0

                =? p - y
goto loop if !=

```

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

Reality: The fastest software today relies on human experts understanding the CPU.

Cannot trust compiler to optimize instruction selection.

Cannot trust compiler to optimize instruction scheduling.

Cannot trust compiler to optimize register allocation.

```

*(uint32 *) (p - 76)
*(uint32 *) (p - 80)
t += xi9
t += xi8
t += xi7
t += xi6
t += xi5
t += xi4
t += xi3
t += xi2
t += xi1
t += xi0

      =? p - y
op if !=

```

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

Reality: The fastest software today relies on human experts understanding the CPU.

Cannot trust compiler to optimize instruction selection.

Cannot trust compiler to optimize instruction scheduling.

Cannot trust compiler to optimize register allocation.

The big

CPUs are
farther a
from naï

*) (p - 76)

*) (p - 80)

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

Reality: The fastest software today relies on human experts understanding the CPU.

Cannot trust compiler to optimize instruction selection.

Cannot trust compiler to optimize instruction scheduling.

Cannot trust compiler to optimize register allocation.

The big picture

CPUs are evolving farther and farther from naive models

76)

80)

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

Reality: The fastest software today relies on human experts understanding the CPU.

Cannot trust compiler to optimize instruction selection.

Cannot trust compiler to optimize instruction scheduling.

Cannot trust compiler to optimize register allocation.

The big picture

CPUs are evolving farther and farther away from naive models of CPUs.

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

Reality: The fastest software today relies on human experts understanding the CPU.

Cannot trust compiler to optimize instruction selection.

Cannot trust compiler to optimize instruction scheduling.

Cannot trust compiler to optimize register allocation.

The big picture

CPUs are evolving farther and farther away from naive models of CPUs.

Wikipedia: “By the late 1990s for even performance sensitive code, optimizing compilers exceeded the performance of human experts.”

Reality: The fastest software today relies on human experts understanding the CPU.

Cannot trust compiler to optimize instruction selection.

Cannot trust compiler to optimize instruction scheduling.

Cannot trust compiler to optimize register allocation.

The big picture

CPUs are evolving farther and farther away from naive models of CPUs.

Minor optimization challenges:

- Pipelining.
- Superscalar processing.

Major optimization challenges:

- Vectorization.
- Many threads; many cores.
- The memory hierarchy; the ring; the mesh.
- Larger-scale parallelism.
- Larger-scale networking.

ia: “By the late 1990s for performance sensitive code, writing compilers exceeded the performance of human experts.”

The fastest software relies on human experts understanding the CPU.

trust compiler to do instruction selection.

trust compiler to do instruction scheduling.

trust compiler to do register allocation.

The big picture

CPUs are evolving farther and farther away from naive models of CPUs.

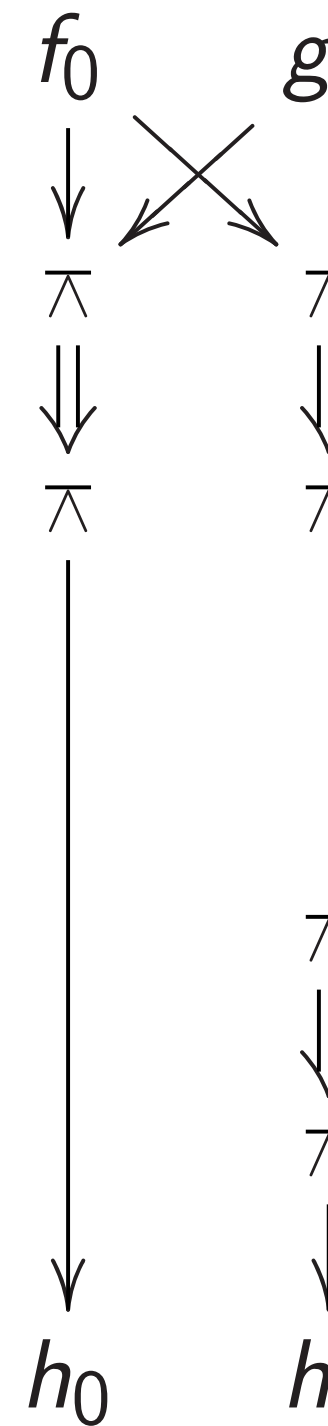
Minor optimization challenges:

- Pipelining.
- Superscalar processing.

Major optimization challenges:

- Vectorization.
- Many threads; many cores.
- The memory hierarchy; the ring; the mesh.
- Larger-scale parallelism.
- Larger-scale networking.

CPU des



Gates π
product
of integers

the late 1990s for
sensitive code,
ers exceeded the
man experts.”

st software
man experts
CPU.

piler to
on selection.

piler to
on scheduling.

piler to
allocation.

The big picture

CPUs are evolving
farther and farther away
from naive models of CPUs.

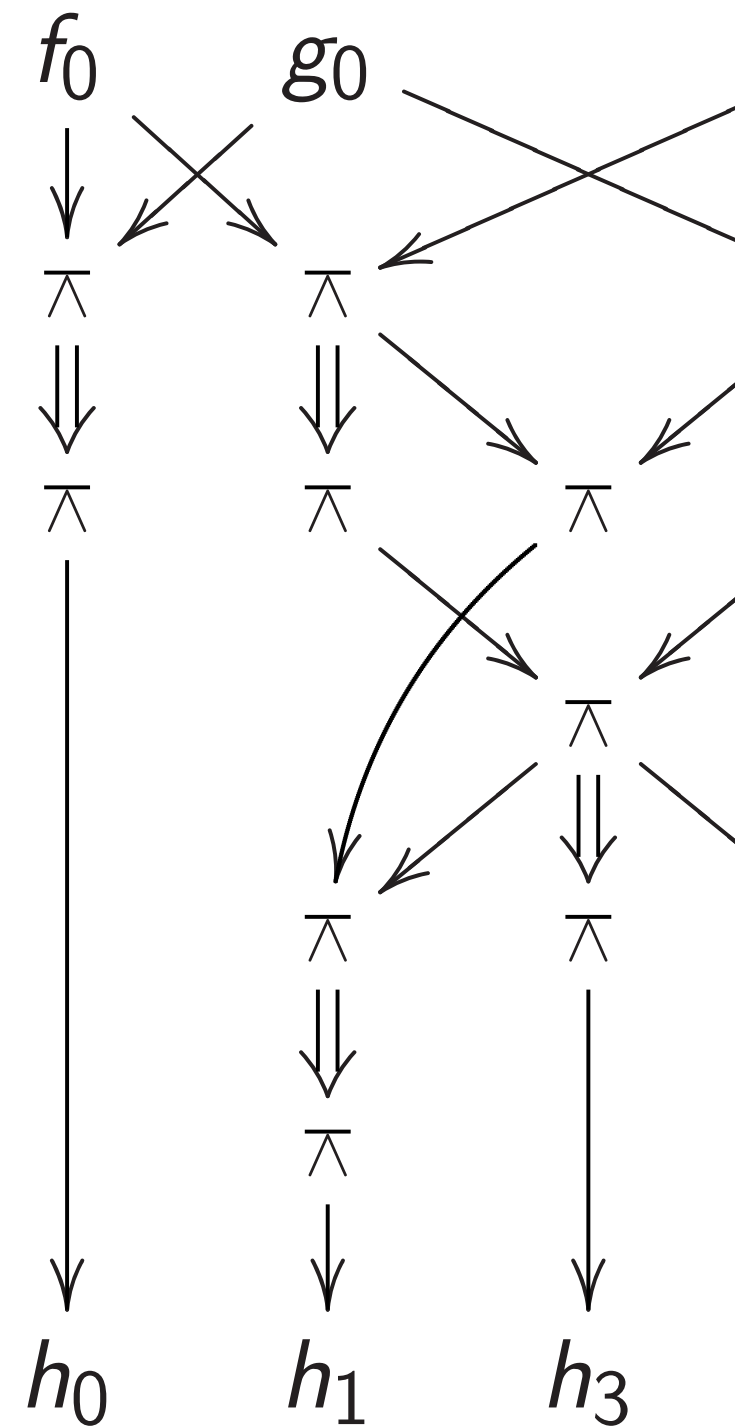
Minor optimization challenges:

- Pipelining.
- Superscalar processing.

Major optimization challenges:

- Vectorization.
- Many threads; many cores.
- The memory hierarchy;
the ring; the mesh.
- Larger-scale parallelism.
- Larger-scale networking.

CPU design in a n



Gates $\pi: a, b \mapsto 1$
product $h_0 + 2h_1 + h_3$
of integers $f_0 + 2f_1 + h_3$

00s for
code,
ed the
rts.”

e
rts

n.

ing.

The big picture

CPUs are evolving
farther and farther away
from naive models of CPUs.

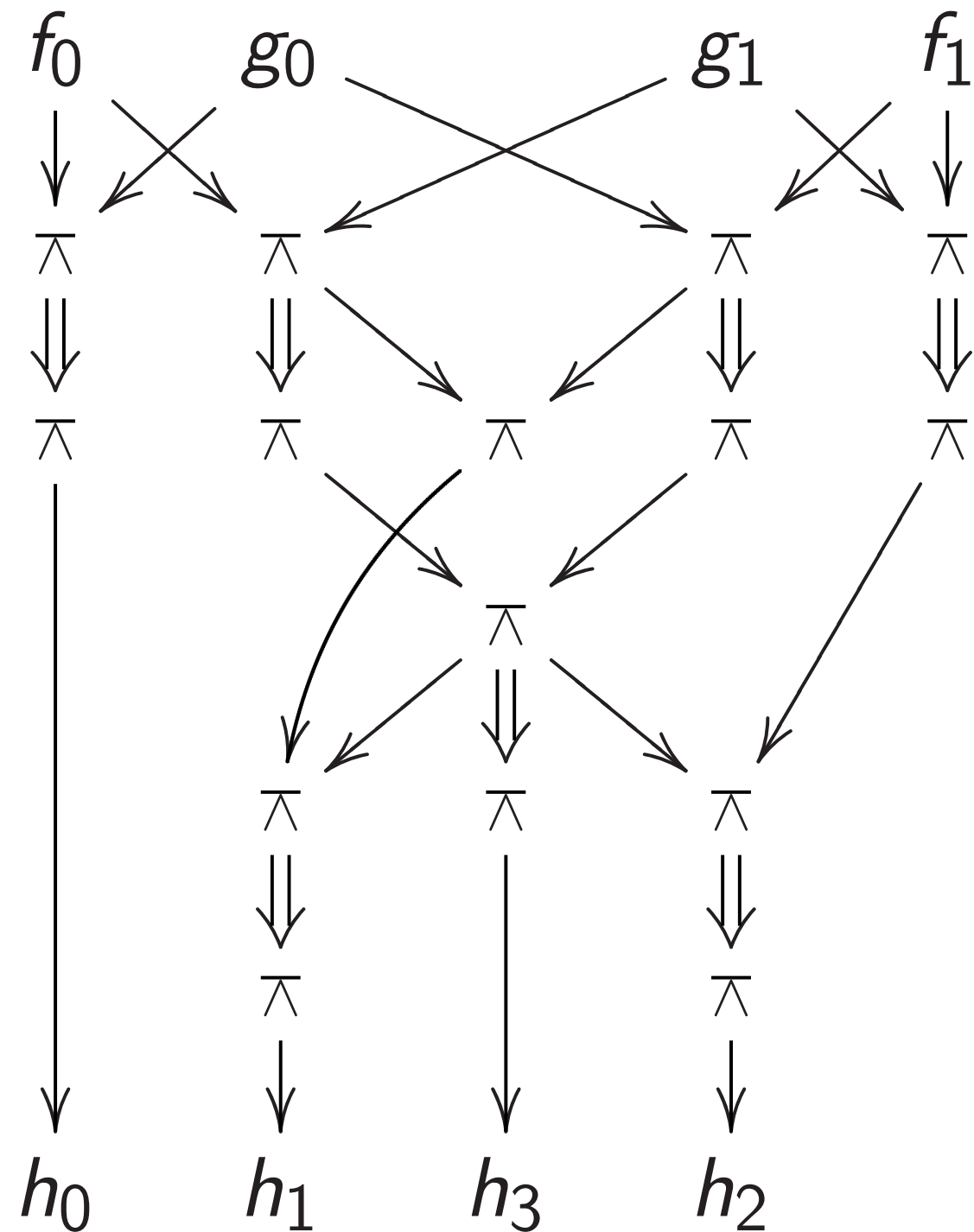
Minor optimization challenges:

- Pipelining.
- Superscalar processing.

Major optimization challenges:

- Vectorization.
- Many threads; many cores.
- The memory hierarchy;
the ring; the mesh.
- Larger-scale parallelism.
- Larger-scale networking.

CPU design in a nutshell



Gates $\pi : a, b \mapsto 1 - ab$ compute
product $h_0 + 2h_1 + 4h_2 + 8h_3$
of integers $f_0 + 2f_1, g_0 + 2g_1$.

The big picture

CPUs are evolving farther and farther away from naive models of CPUs.

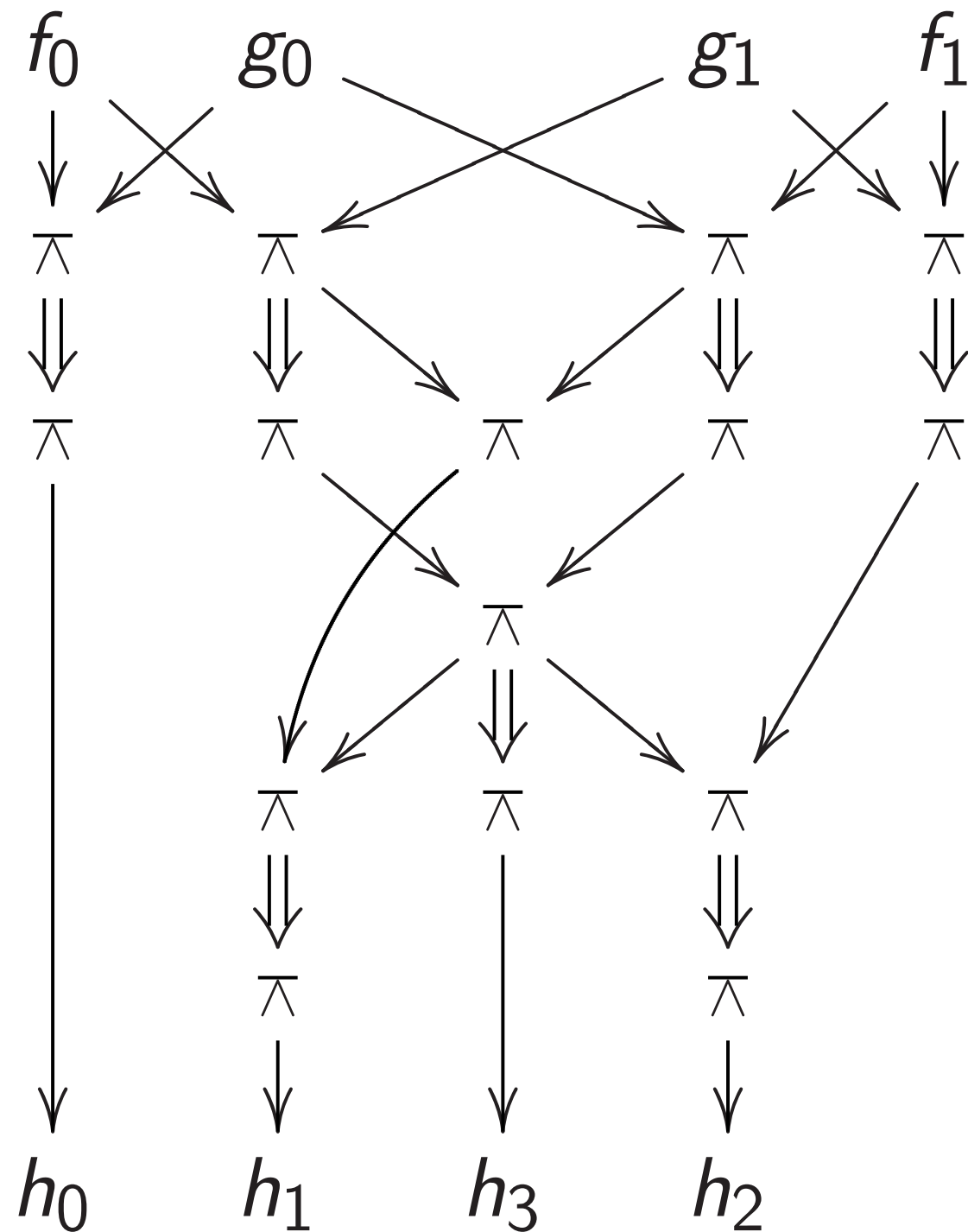
Minor optimization challenges:

- Pipelining.
- Superscalar processing.

Major optimization challenges:

- Vectorization.
- Many threads; many cores.
- The memory hierarchy; the ring; the mesh.
- Larger-scale parallelism.
- Larger-scale networking.

CPU design in a nutshell



Gates $\pi : a, b \mapsto 1 - ab$ computing product $h_0 + 2h_1 + 4h_2 + 8h_3$ of integers $f_0 + 2f_1, g_0 + 2g_1$.

picture

the evolving
and farther away
ve models of CPUs.

optimization challenges:
ning.

scalar processing.

optimization challenges:
rization.

threads; many cores.

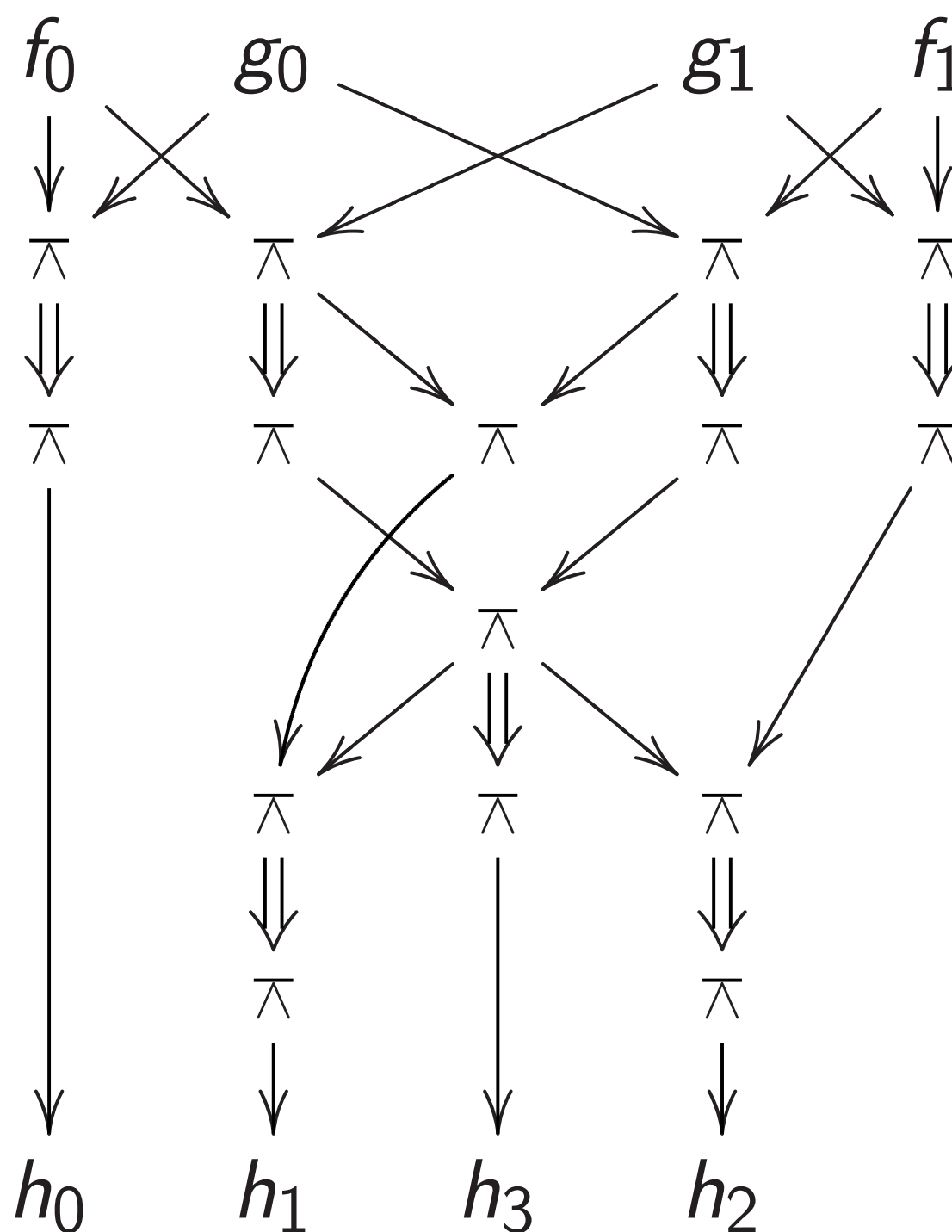
memory hierarchy;

g; the mesh.

-scale parallelism.

-scale networking.

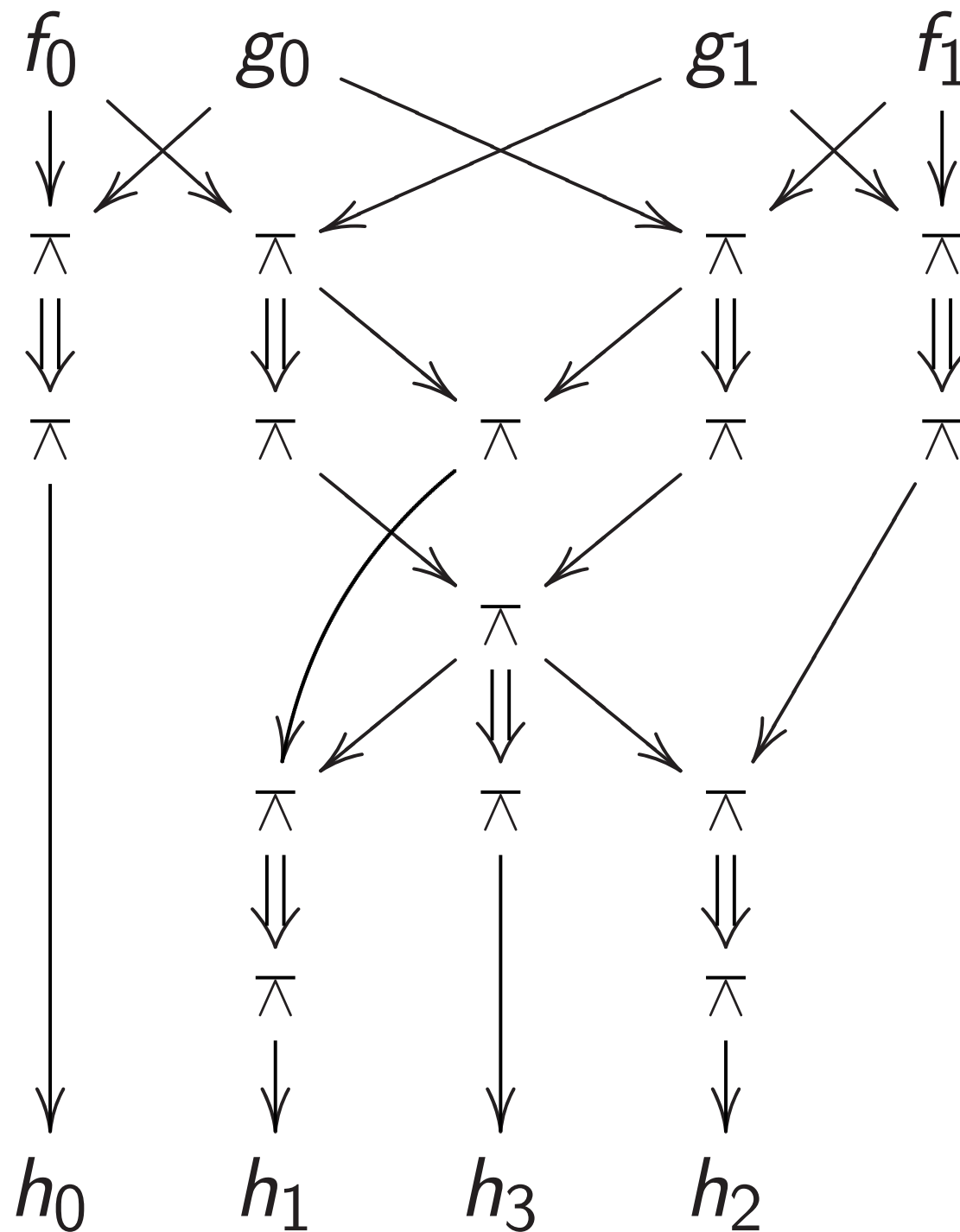
CPU design in a nutshell



Gates $\neg : a, b \mapsto 1 - ab$ computing
product $h_0 + 2h_1 + 4h_2 + 8h_3$
of integers $f_0 + 2f_1, g_0 + 2g_1$.

Electricity
percolates
If f_0, f_1, \dots
then h_0, \dots
a few m

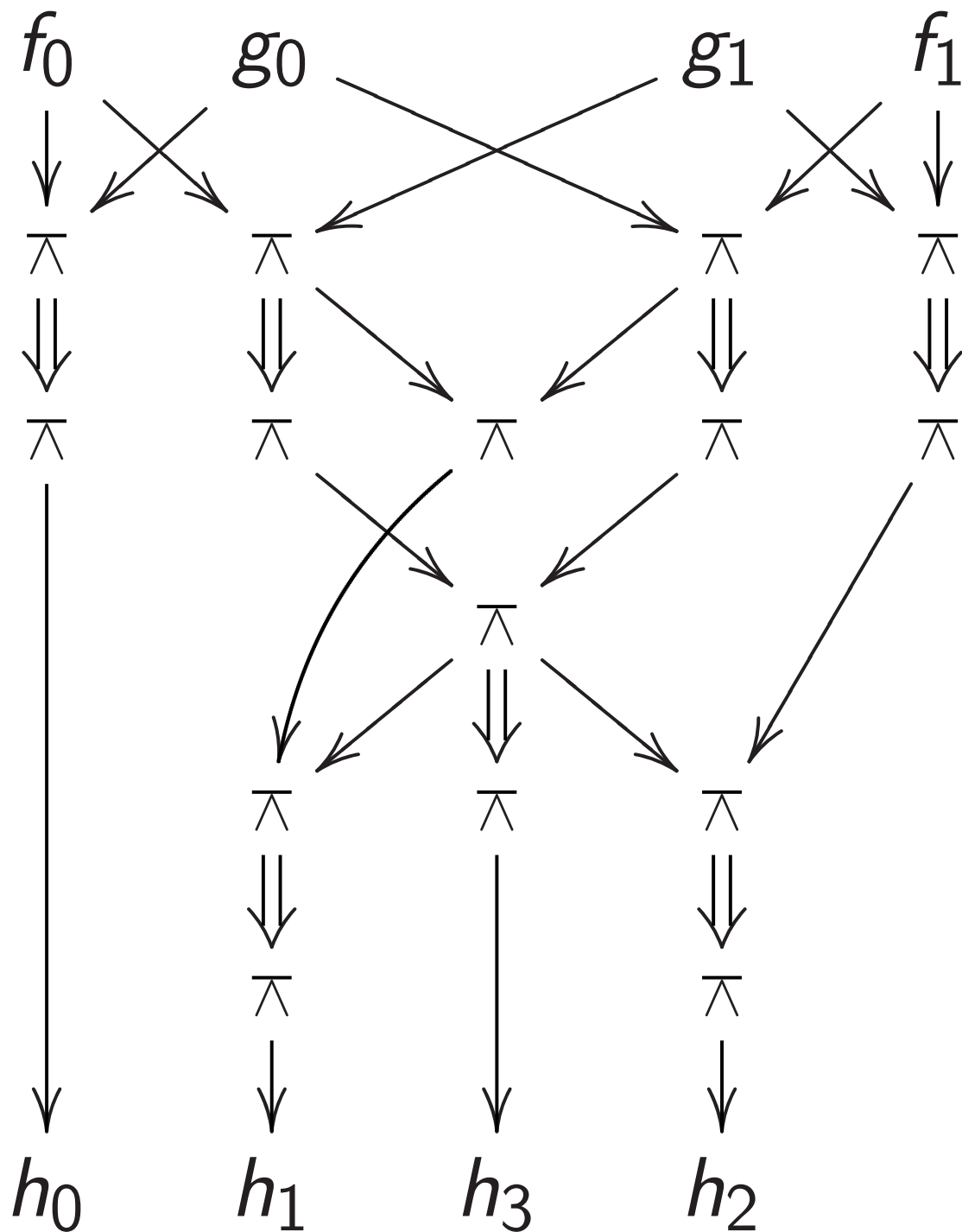
CPU design in a nutshell



Gates $\pi : a, b \mapsto 1 - ab$ computing product $h_0 + 2h_1 + 4h_2 + 8h_3$ of integers $f_0 + 2f_1, g_0 + 2g_1$.

Electricity takes time to percolate through the system. If f_0, f_1, g_0, g_1 are inputs, then h_0, h_1, h_2, h_3 are outputs a few moments later.

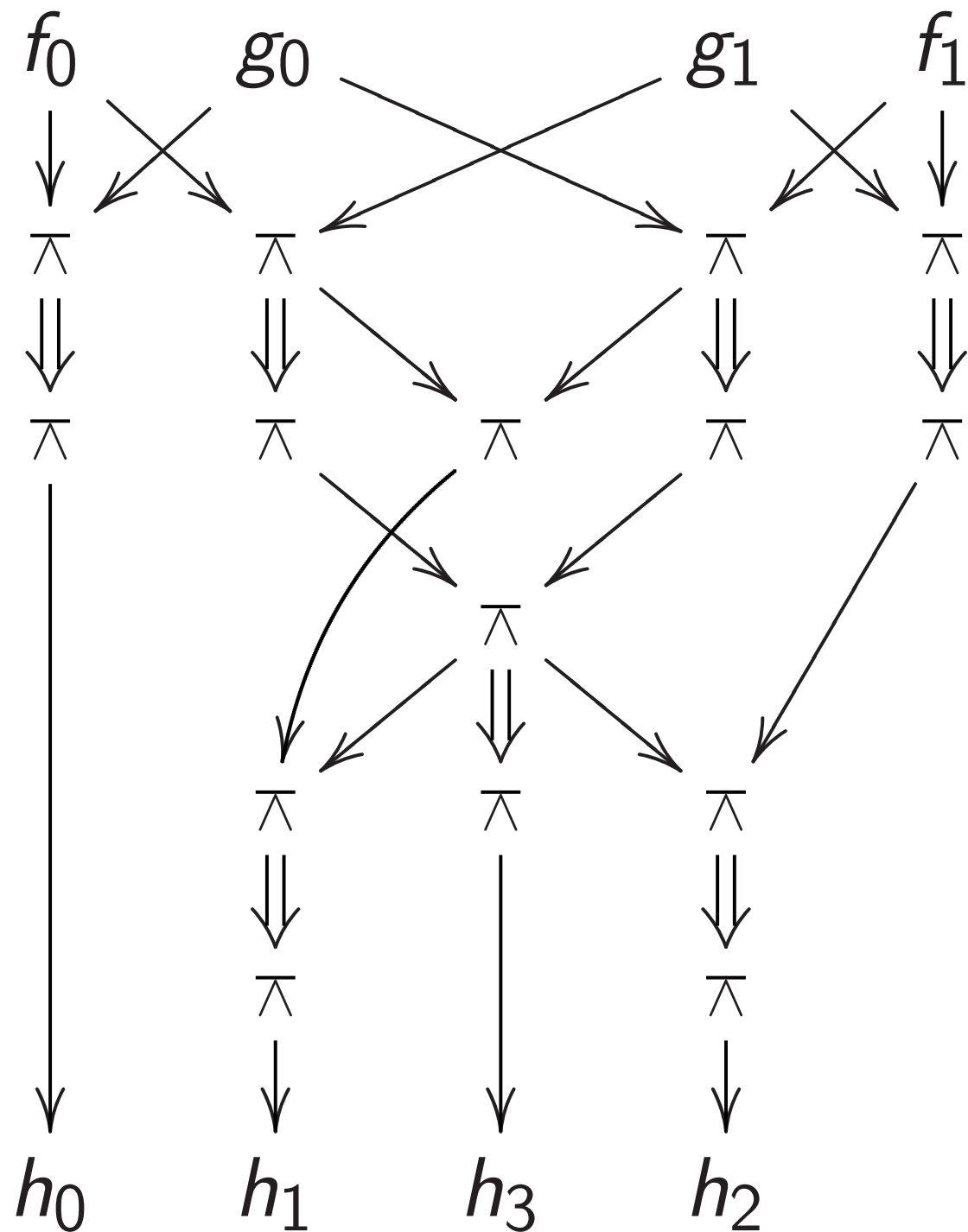
CPU design in a nutshell



Gates $\bar{\wedge} : a, b \mapsto 1 - ab$ computing
 product $h_0 + 2h_1 + 4h_2 + 8h_3$
 of integers $f_0 + 2f_1, g_0 + 2g_1$.

Electricity takes time to
 percolate through wires and
 If f_0, f_1, g_0, g_1 are stable
 then h_0, h_1, h_2, h_3 are stable
 a few moments later.

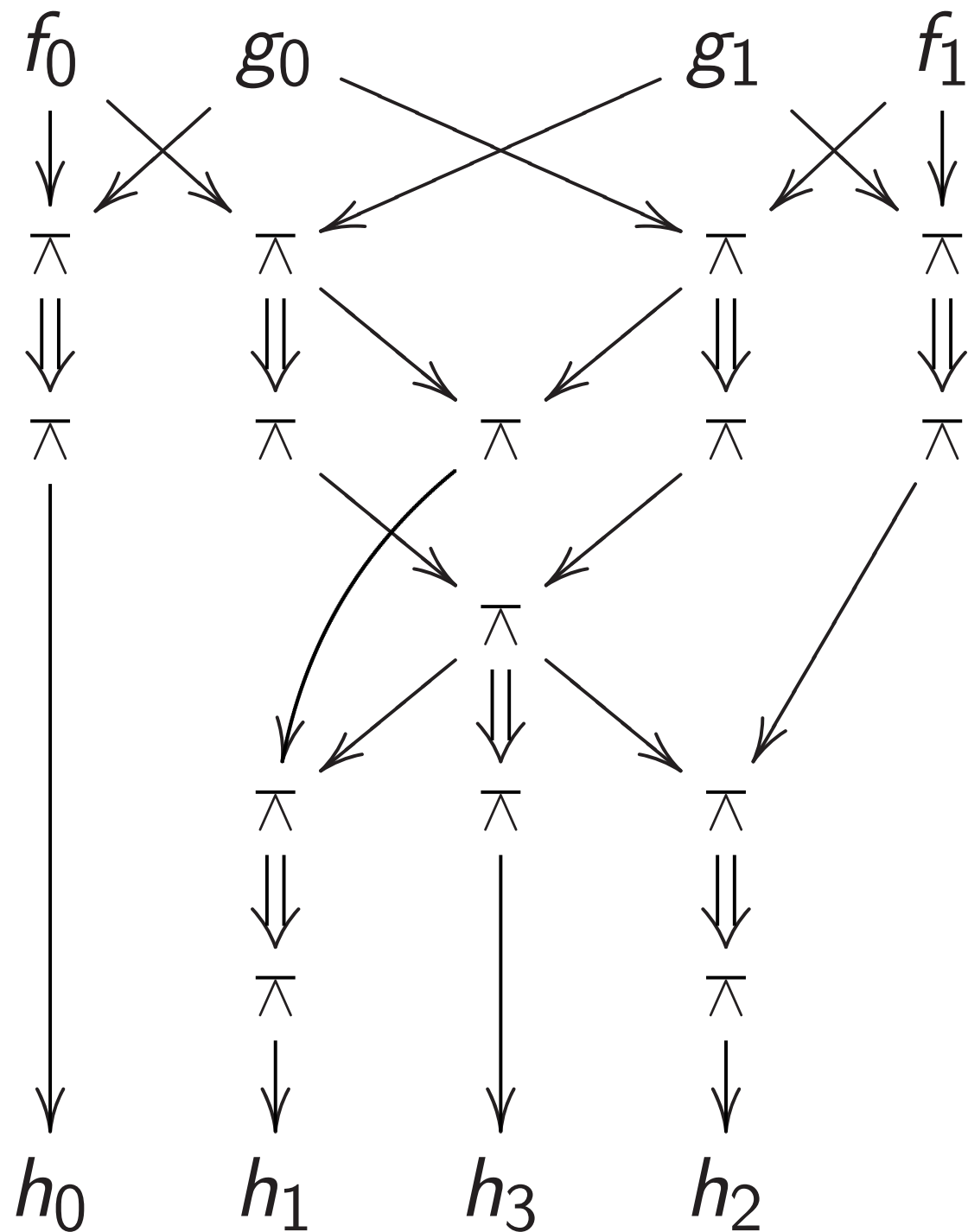
CPU design in a nutshell



Gates $\neg : a, b \mapsto 1 - ab$ computing product $h_0 + 2h_1 + 4h_2 + 8h_3$ of integers $f_0 + 2f_1, g_0 + 2g_1$.

Electricity takes time to percolate through wires and gates. If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

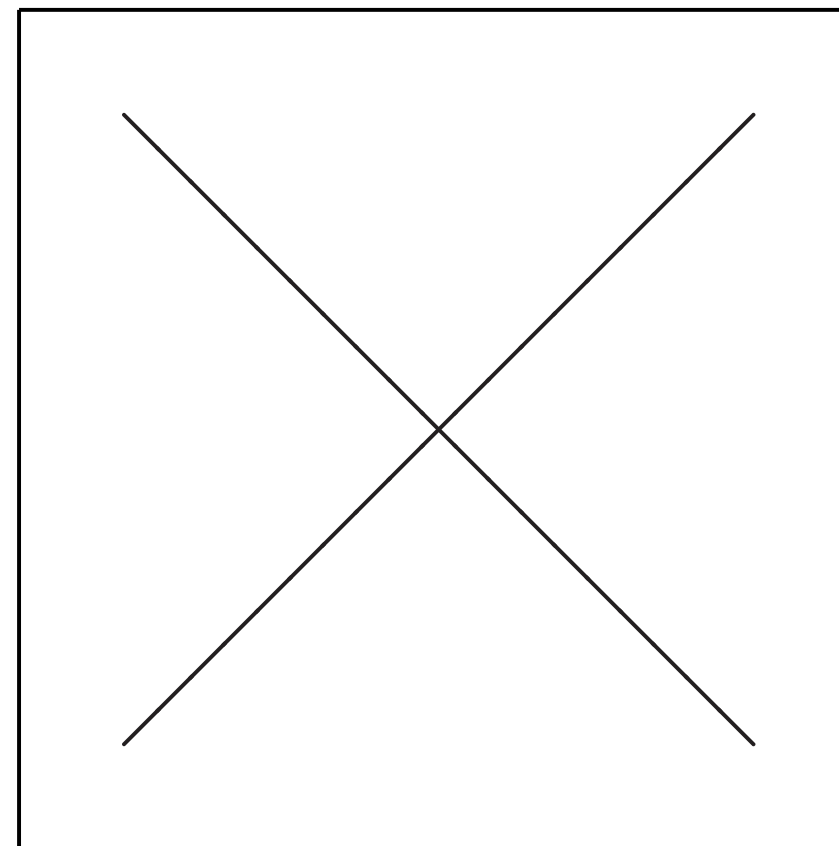
CPU design in a nutshell



Gates $\bar{\wedge} : a, b \mapsto 1 - ab$ computing product $h_0 + 2h_1 + 4h_2 + 8h_3$ of integers $f_0 + 2f_1, g_0 + 2g_1$.

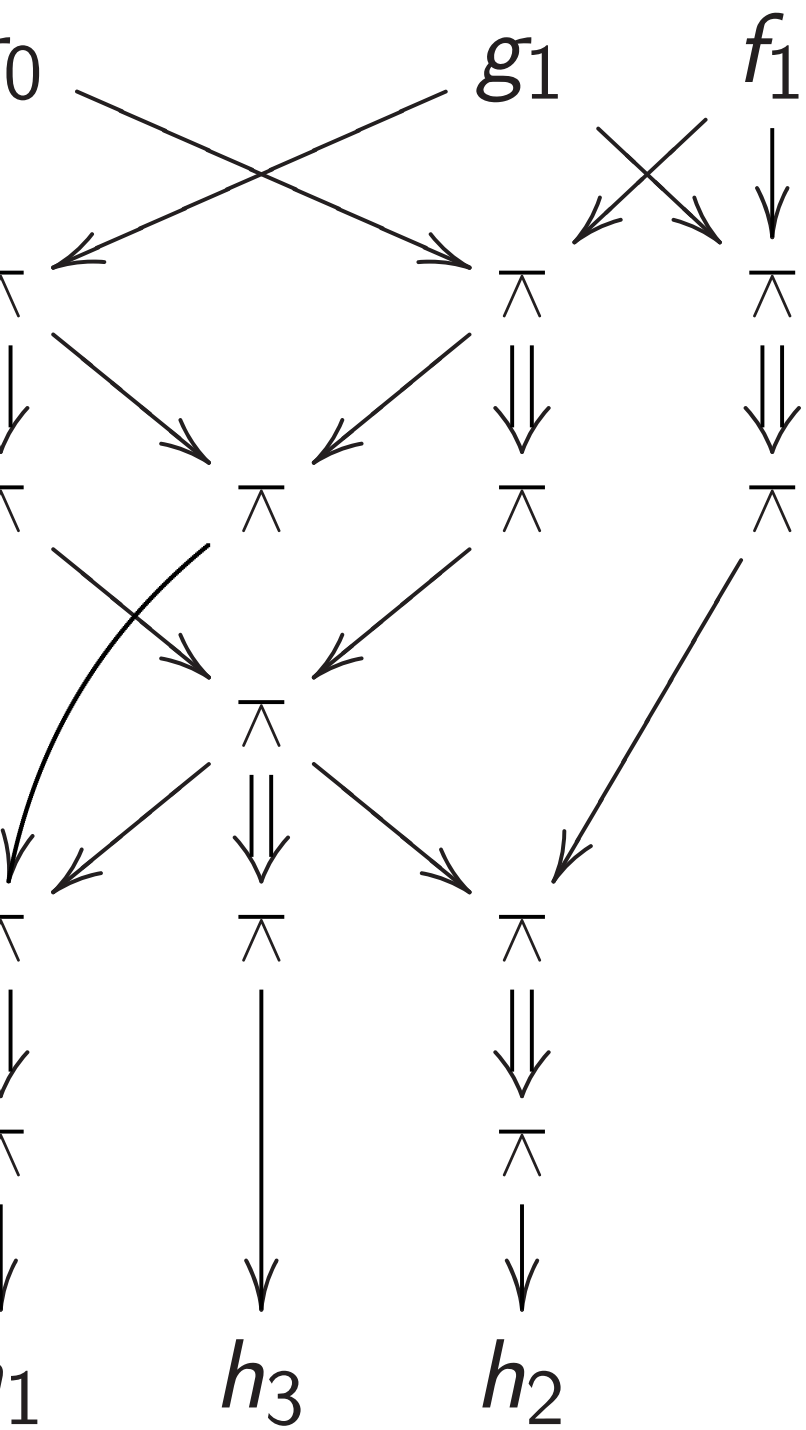
Electricity takes time to percolate through wires and gates. If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:



(Details omitted.)

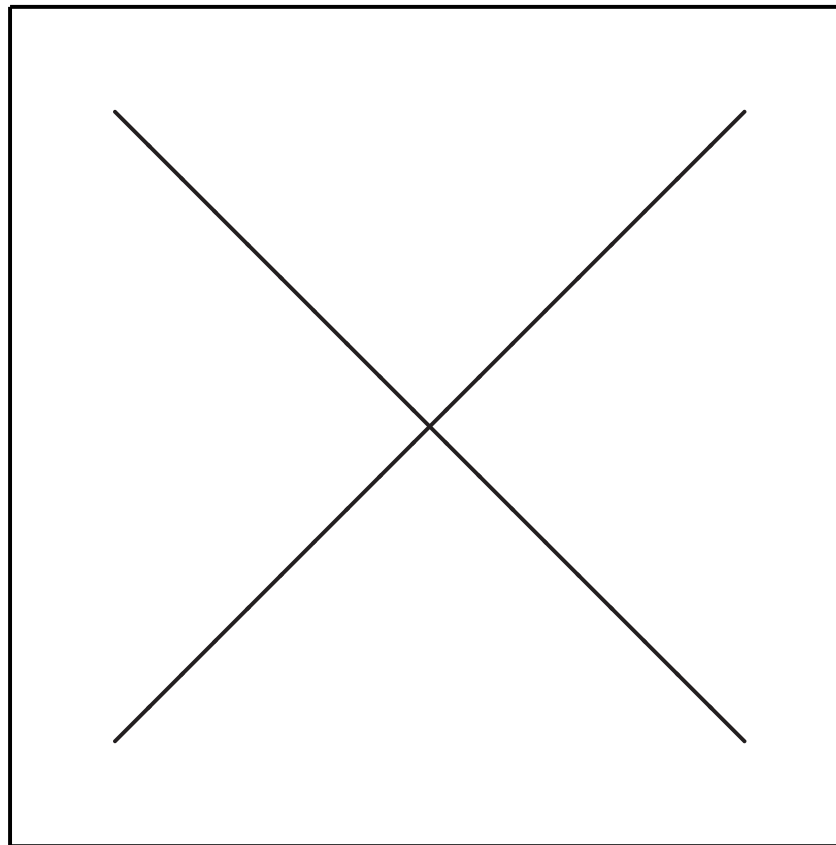
Design in a nutshell



$f_0, g_1 \mapsto 1 - ab$ computing
 $h_0 + 2h_1 + 4h_2 + 8h_3$
 using $f_0 + 2f_1, g_0 + 2g_1$.

Electricity takes time to percolate through wires and gates. If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:

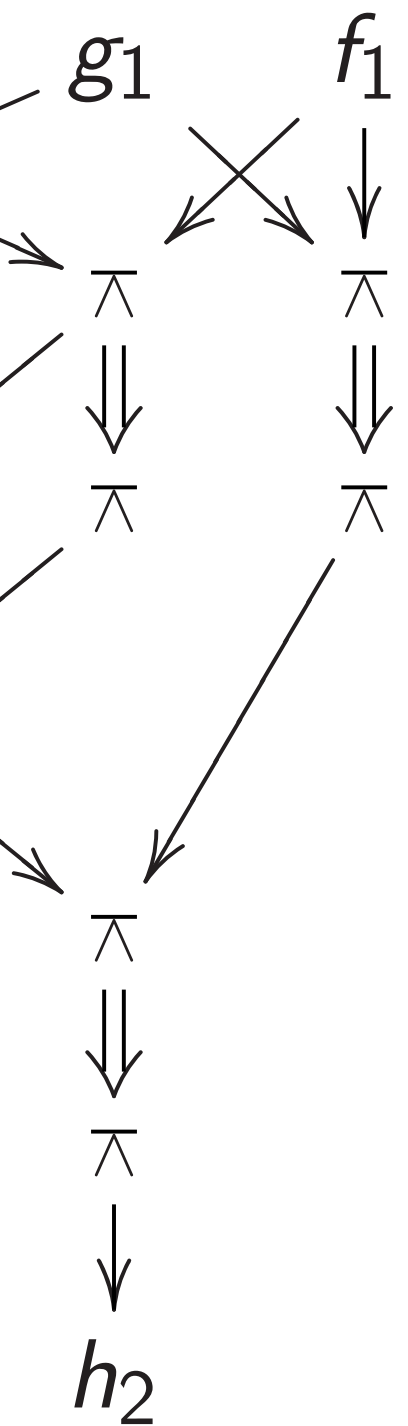


(Details omitted.)

Build circuit to multiply 32-bit integers using given 4-bit and 32-bit gates.



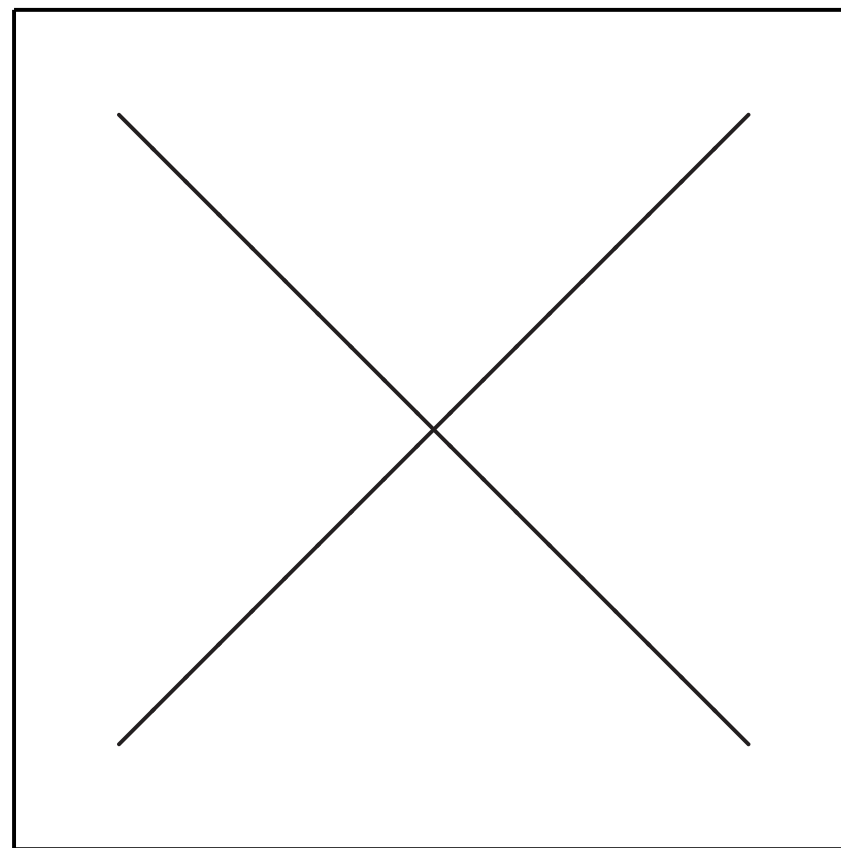
utshell



– ab computing
 $+ 4h_2 + 8h_3$
 $+ g_1, g_0 + 2g_1.$

Electricity takes time to percolate through wires and gates. If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:



(Details omitted.)

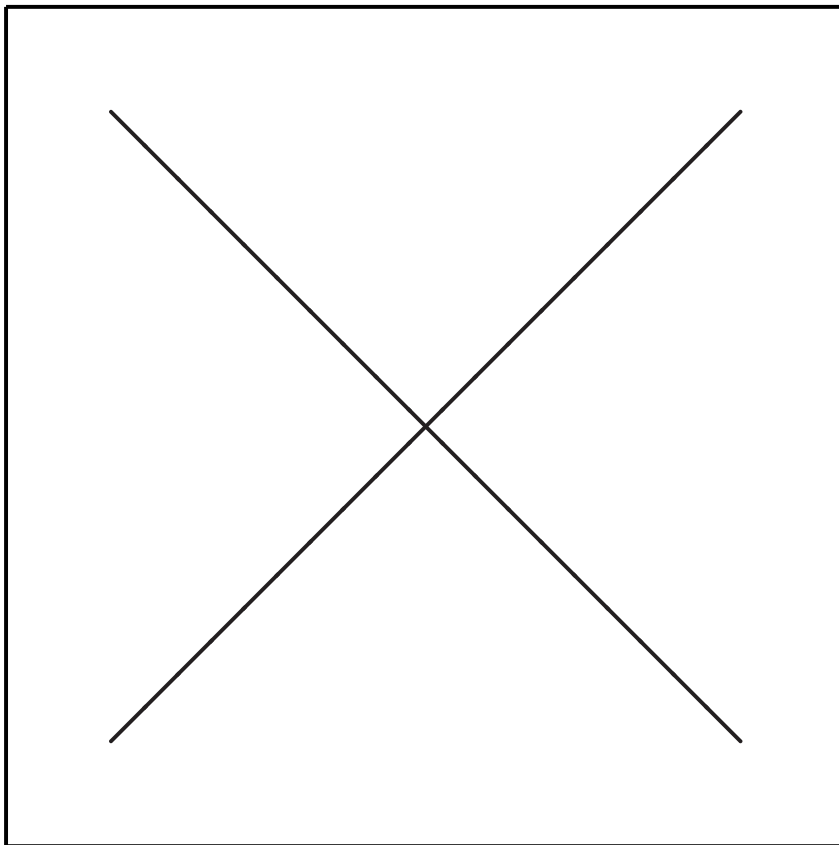
Build circuit to compute 32-bit integer r_i given 4-bit integers and 32-bit integers

register
read

Electricity takes time to percolate through wires and gates.

If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:



(Details omitted.)

Build circuit to compute 32-bit integer r_i

given 4-bit integer i

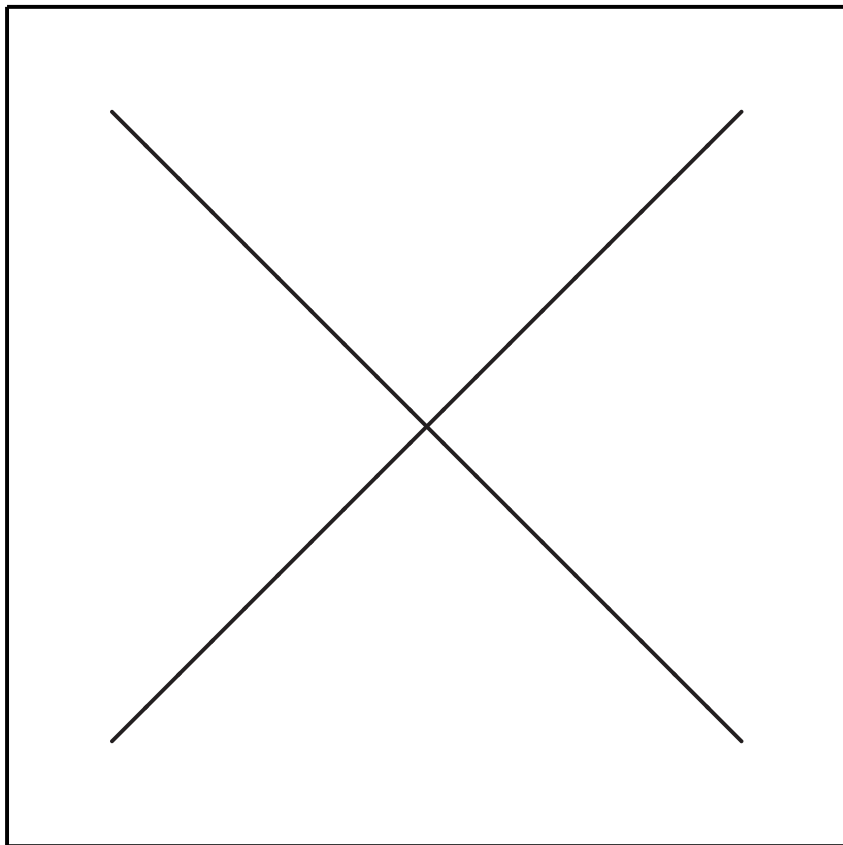
and 32-bit integers r_0, r_1, \dots

register
read

Electricity takes time to percolate through wires and gates.

If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:



(Details omitted.)

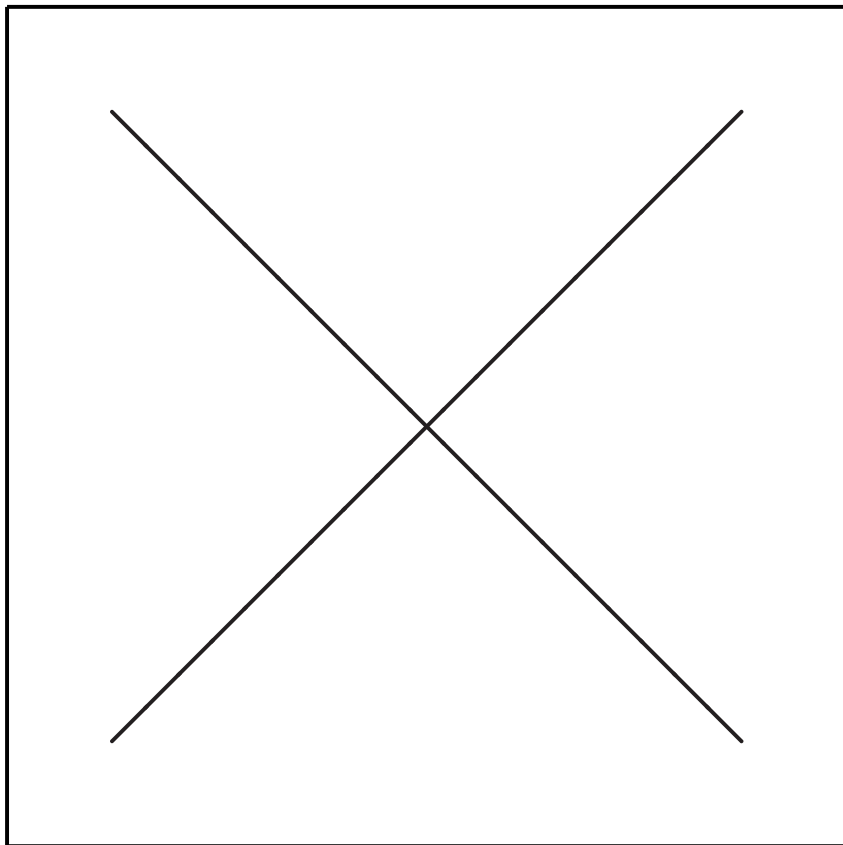
Build circuit to compute 32-bit integer r_i given 4-bit integer i and 32-bit integers r_0, r_1, \dots, r_{15} :

register
read

Electricity takes time to percolate through wires and gates.

If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:



(Details omitted.)

Build circuit to compute 32-bit integer r_i given 4-bit integer i and 32-bit integers r_0, r_1, \dots, r_{15} :

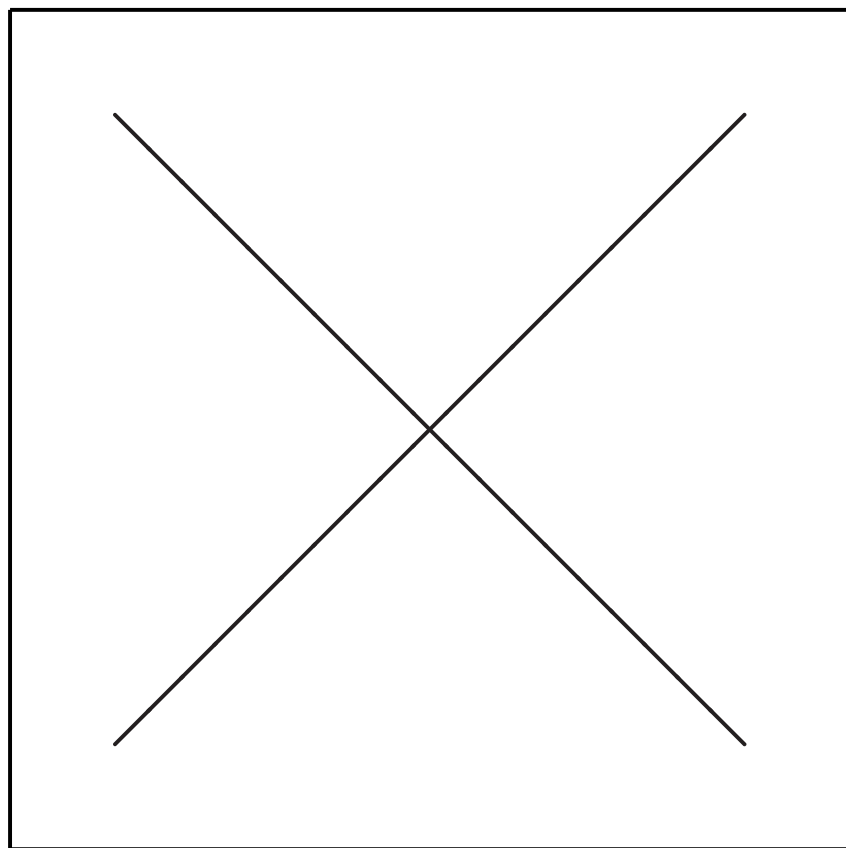
register
read

Build circuit for “register write”:
 $r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$
where $r'_j = r_j$ except $r'_i = s$.

Electricity takes time to percolate through wires and gates.

If f_0, f_1, g_0, g_1 are stable then h_0, h_1, h_2, h_3 are stable a few moments later.

Build circuit with more gates to multiply (e.g.) 32-bit integers:



(Details omitted.)

Build circuit to compute 32-bit integer r_i given 4-bit integer i and 32-bit integers r_0, r_1, \dots, r_{15} :

register
read

Build circuit for “register write”:

$r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$

where $r'_j = r_j$ except $r'_i = s$.

Build circuit for addition. Etc.

ity takes time to
e through wires and gates.

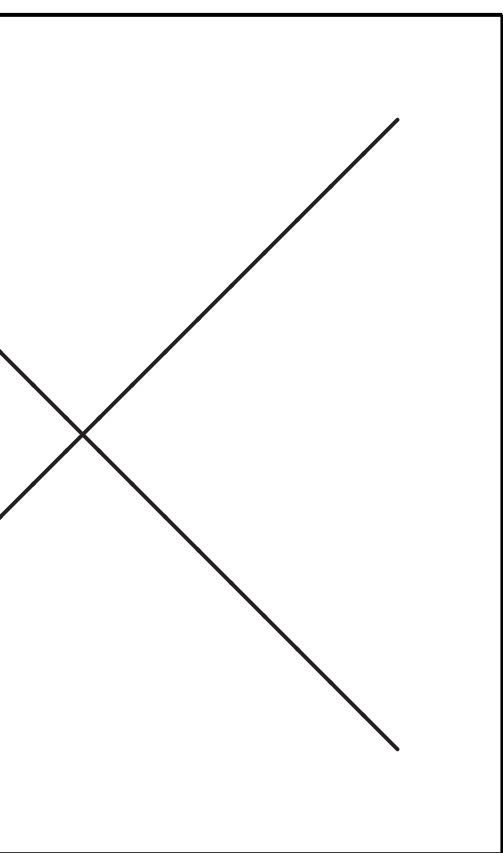
g_0, g_1 are stable

h_1, h_2, h_3 are stable

oments later.

circuit with more gates

ply (e.g.) 32-bit integers:



omitted.)

Build circuit to compute
32-bit integer r_i
given 4-bit integer i
and 32-bit integers r_0, r_1, \dots, r_{15} :

register
read

Build circuit for “register write”:

$r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$

where $r'_j = r_j$ except $r'_i = s$.

Build circuit for addition. Etc.

r_0, \dots, r_{15}
where r'_ℓ

regist
rea

r

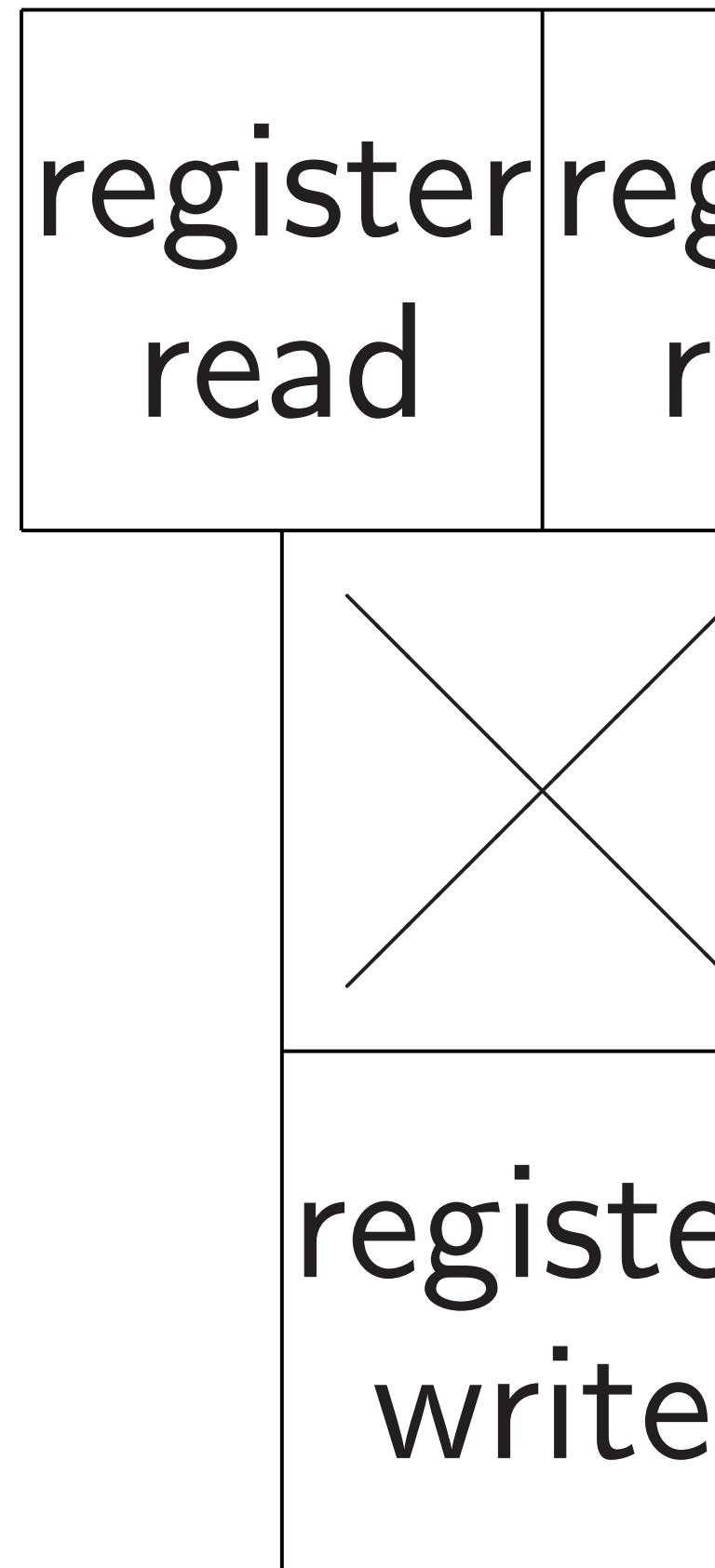
me to
wires and gates.
stable
are stable
ter.
more gates
32-bit integers:

Build circuit to compute
32-bit integer r_i
given 4-bit integer i
and 32-bit integers r_0, r_1, \dots, r_{15} :

register
read

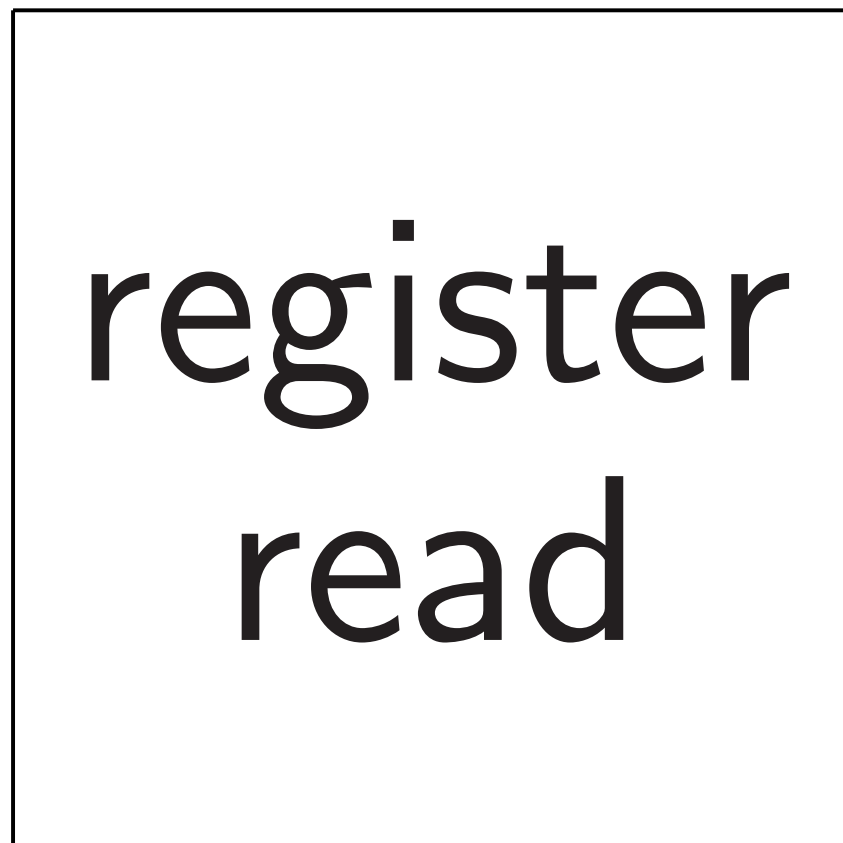
Build circuit for “register write”:
 $r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$
where $r'_j = r_j$ except $r'_i = s$.
Build circuit for addition. Etc.

$r_0, \dots, r_{15}, i, j, k \mapsto$
where $r'_\ell = r_\ell$ except



gates.

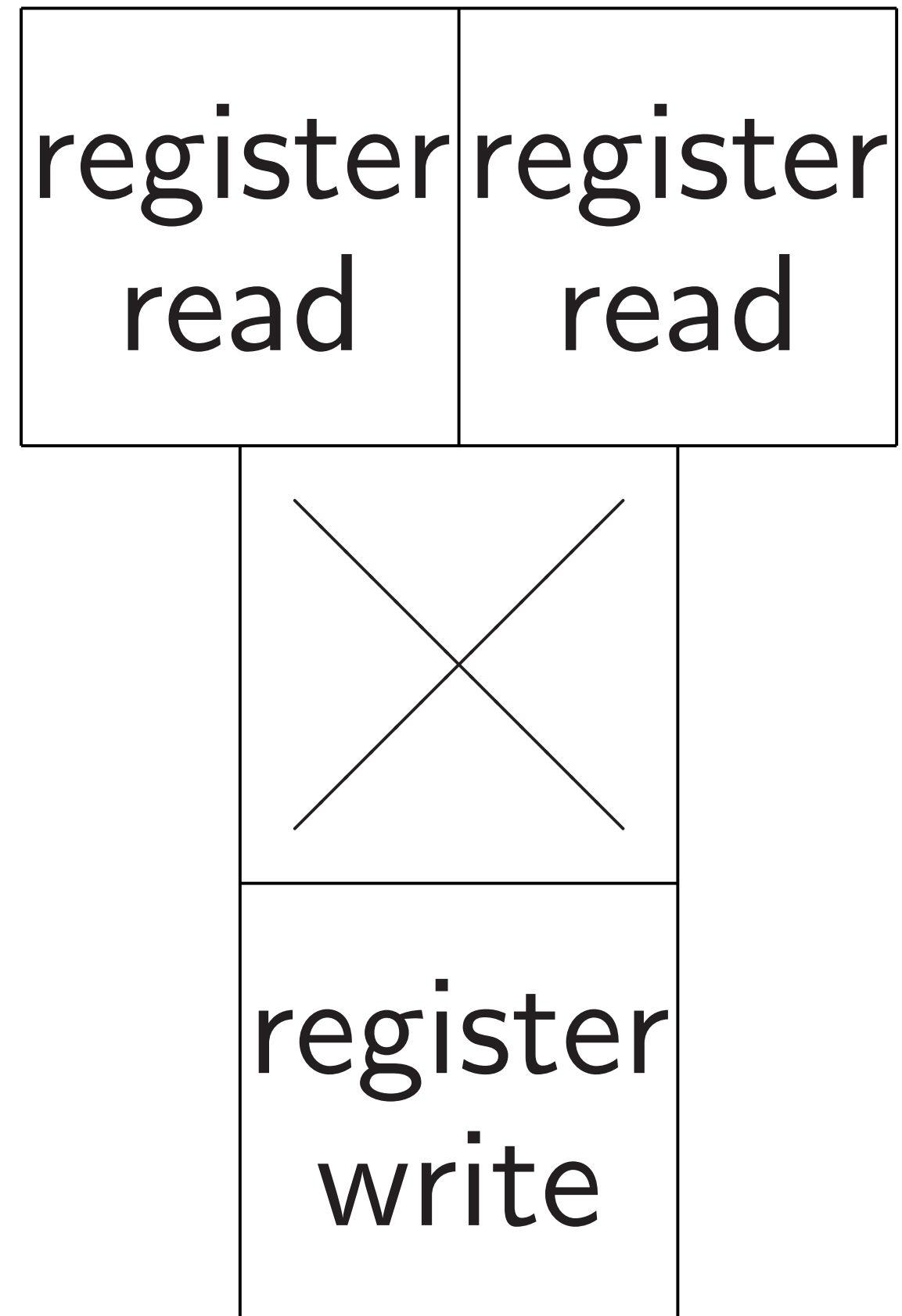
Build circuit to compute
32-bit integer r_i
given 4-bit integer i
and 32-bit integers r_0, r_1, \dots, r_{15} :



s
egers:

Build circuit for “register write”:
 $r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$
where $r'_j = r_j$ except $r'_i = s$.
Build circuit for addition. Etc.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
where $r'_\ell = r_\ell$ except $r'_i = r_j$

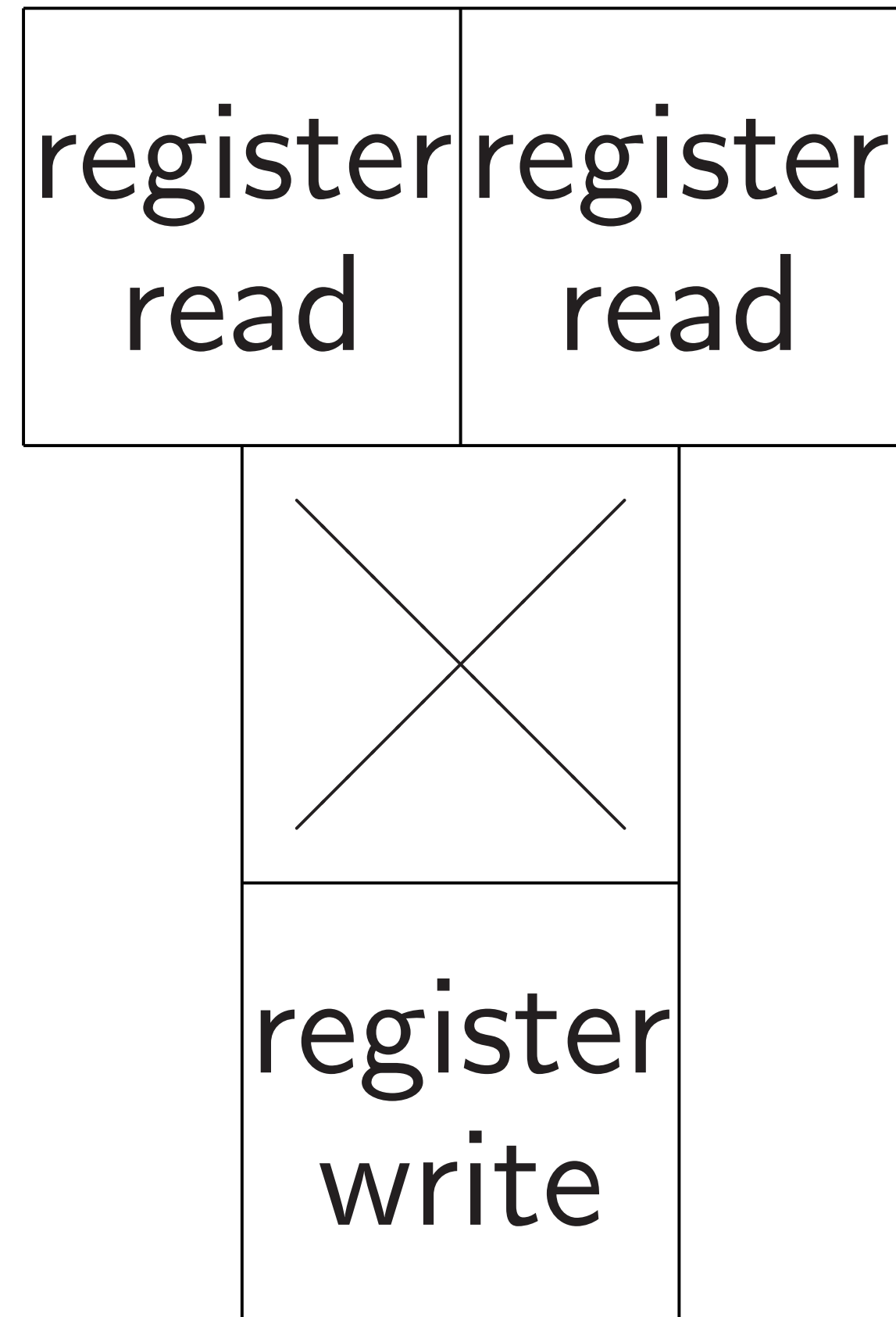


Build circuit to compute
 32-bit integer r_i
 given 4-bit integer i
 and 32-bit integers r_0, r_1, \dots, r_{15} :

register
read

Build circuit for “register write”:
 $r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$
 where $r'_j = r_j$ except $r'_i = s$.
 Build circuit for addition. Etc.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



circuit to compute

integer r_i

bit integer i

bit integers r_0, r_1, \dots, r_{15} :

register
read

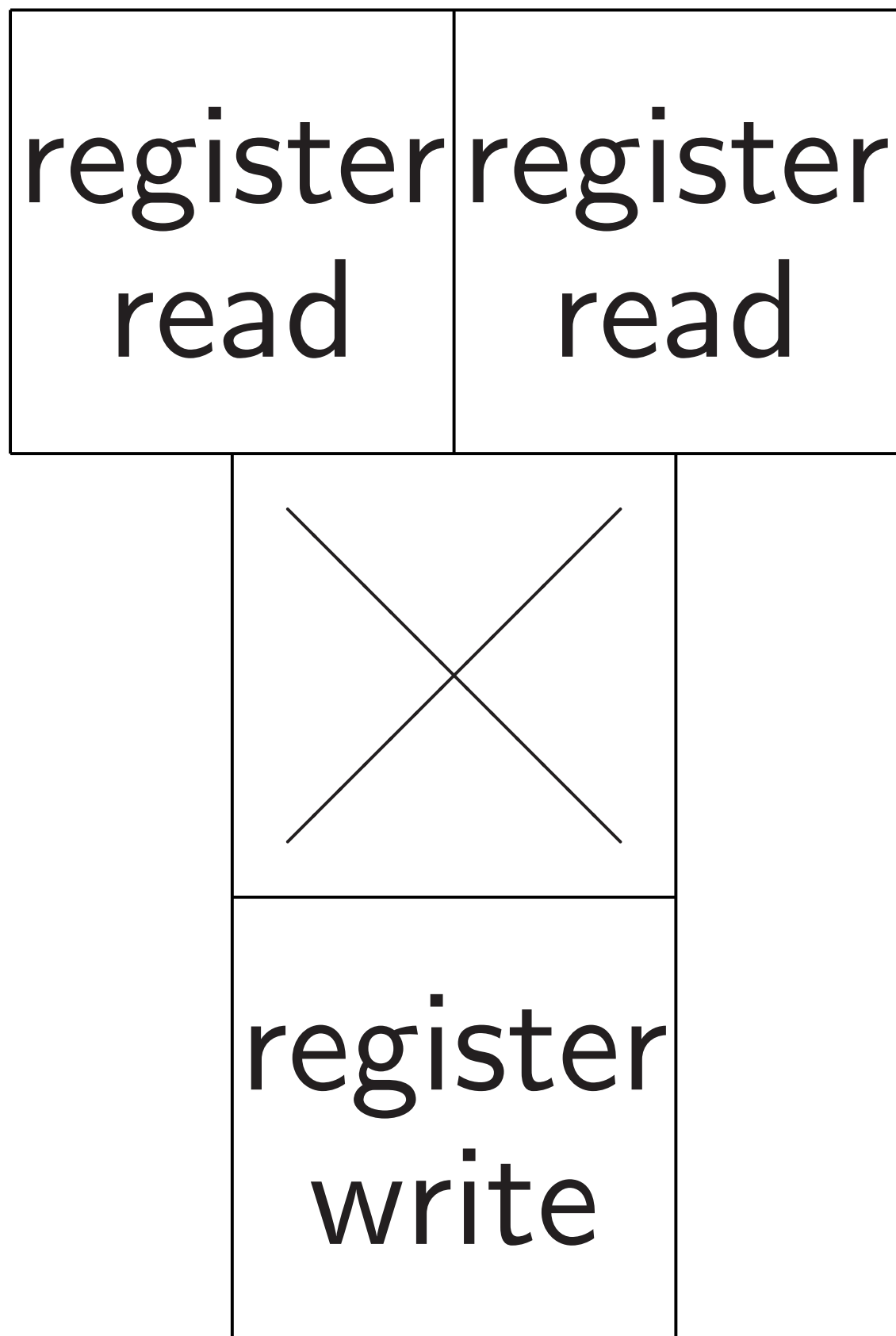
circuit for “register write”:

$r_0, \dots, r_{15}, s, i \mapsto r'_0, \dots, r'_{15}$

$r'_\ell = r_\ell$ except $r'_i = s$.

circuit for addition. Etc.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more

More arith

replace (

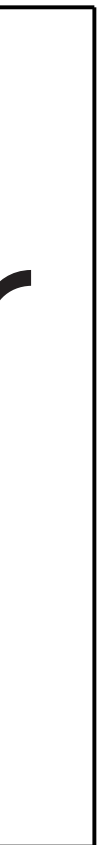
(“×”, i, j)

(“+”, i, j)

compute

i

registers r_0, r_1, \dots, r_{15} :



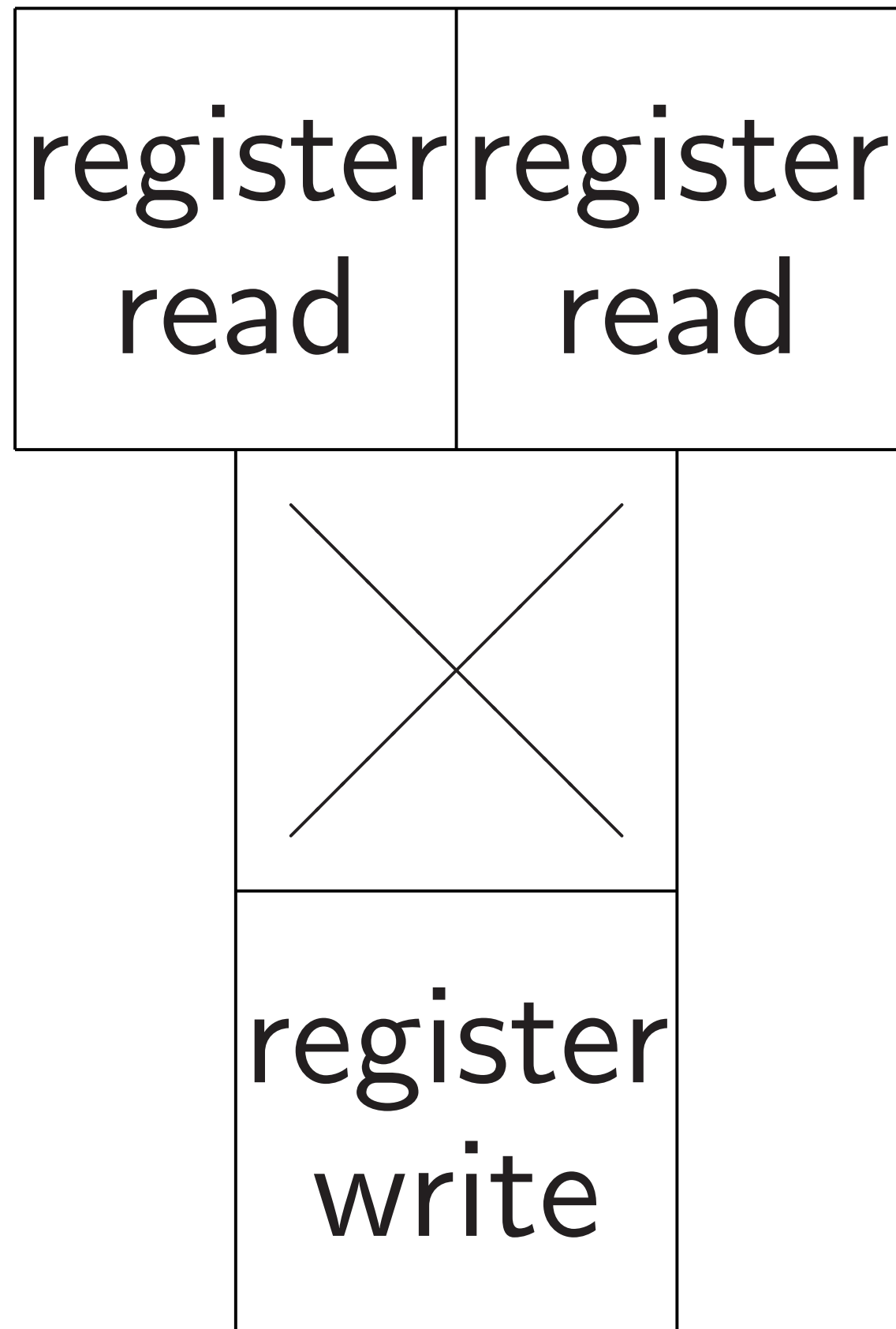
“register write”:

r'_0, \dots, r'_{15}

except $r'_i = s$.

addition. Etc.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more flexibility

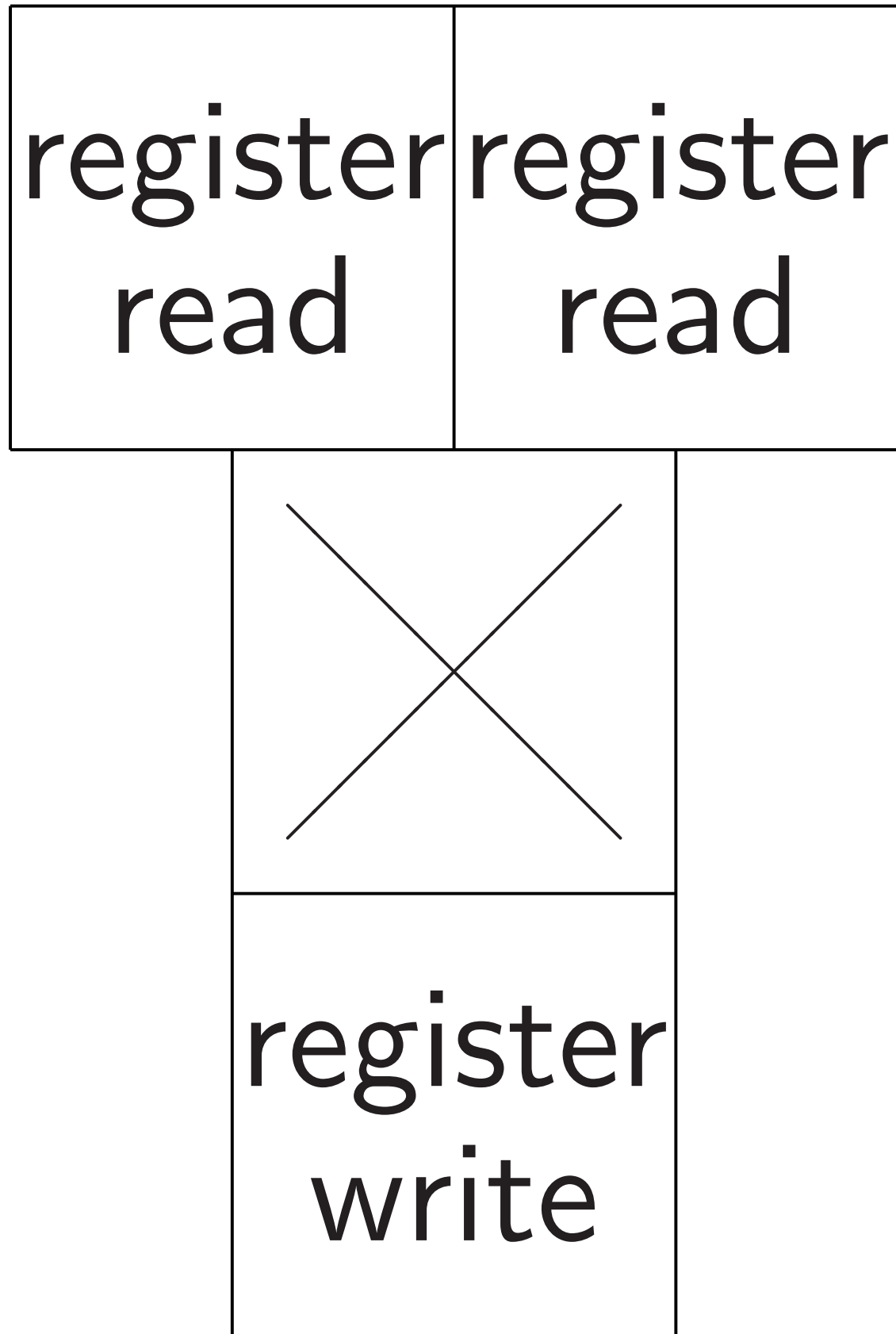
More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more flexibility.

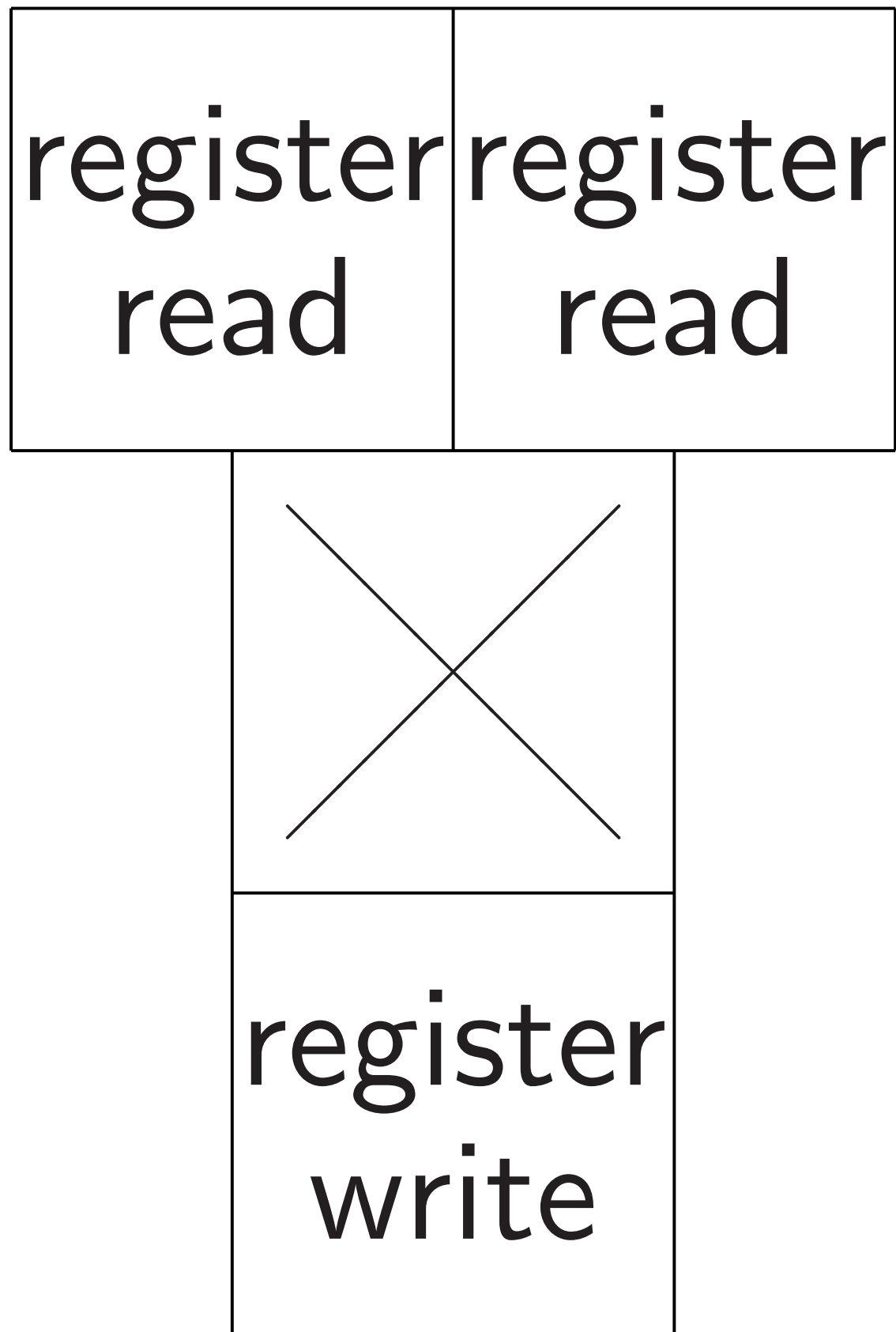
More arithmetic:

replace (i, j, k) with

$(\text{"\times"}, i, j, k)$ and

$(\text{"+"}, i, j, k)$ and more options

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more flexibility.

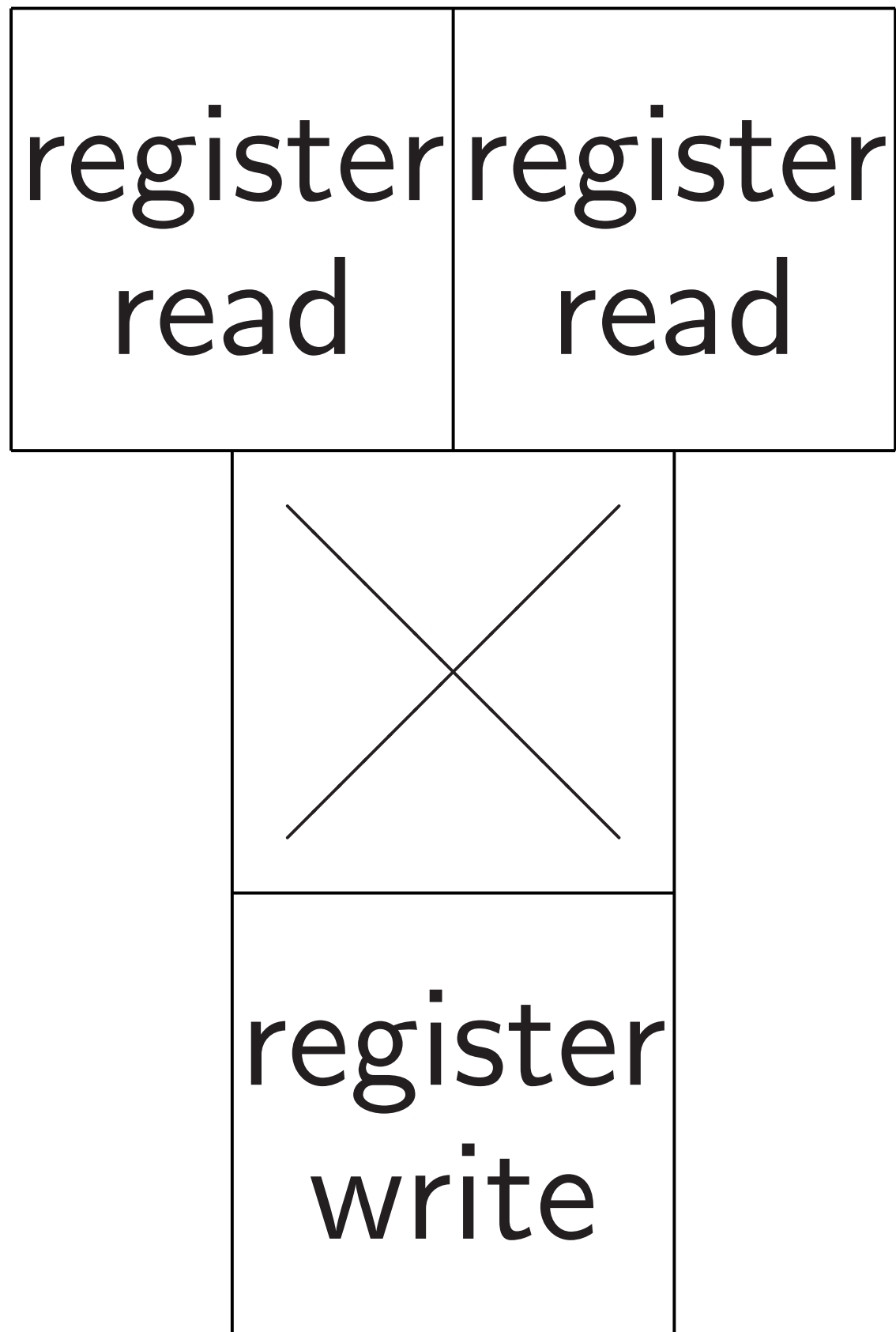
More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

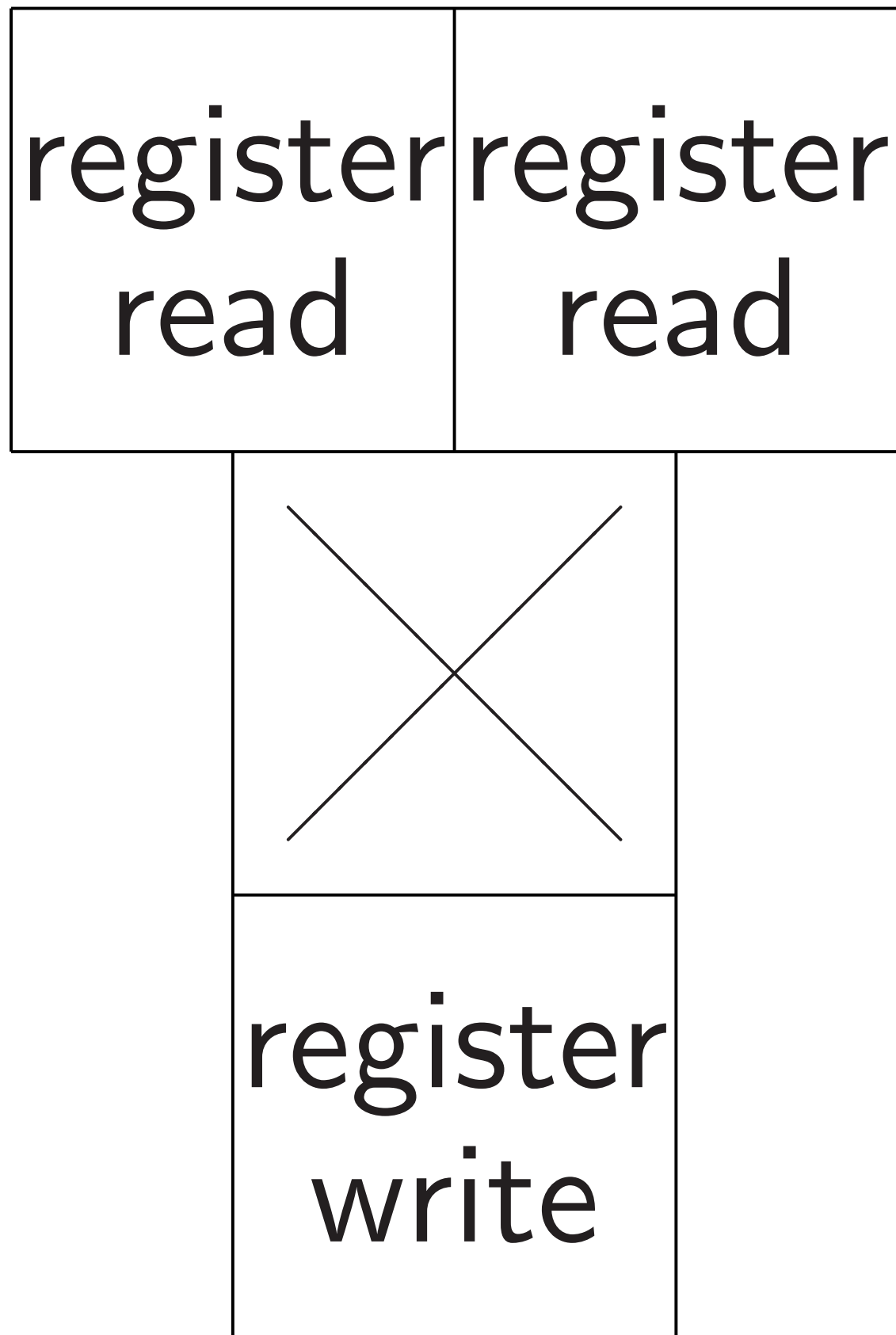
$(“+”, i, j, k)$ and more options.

More (but slower) storage:

“load” from and “store” to

larger “RAM” arrays.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

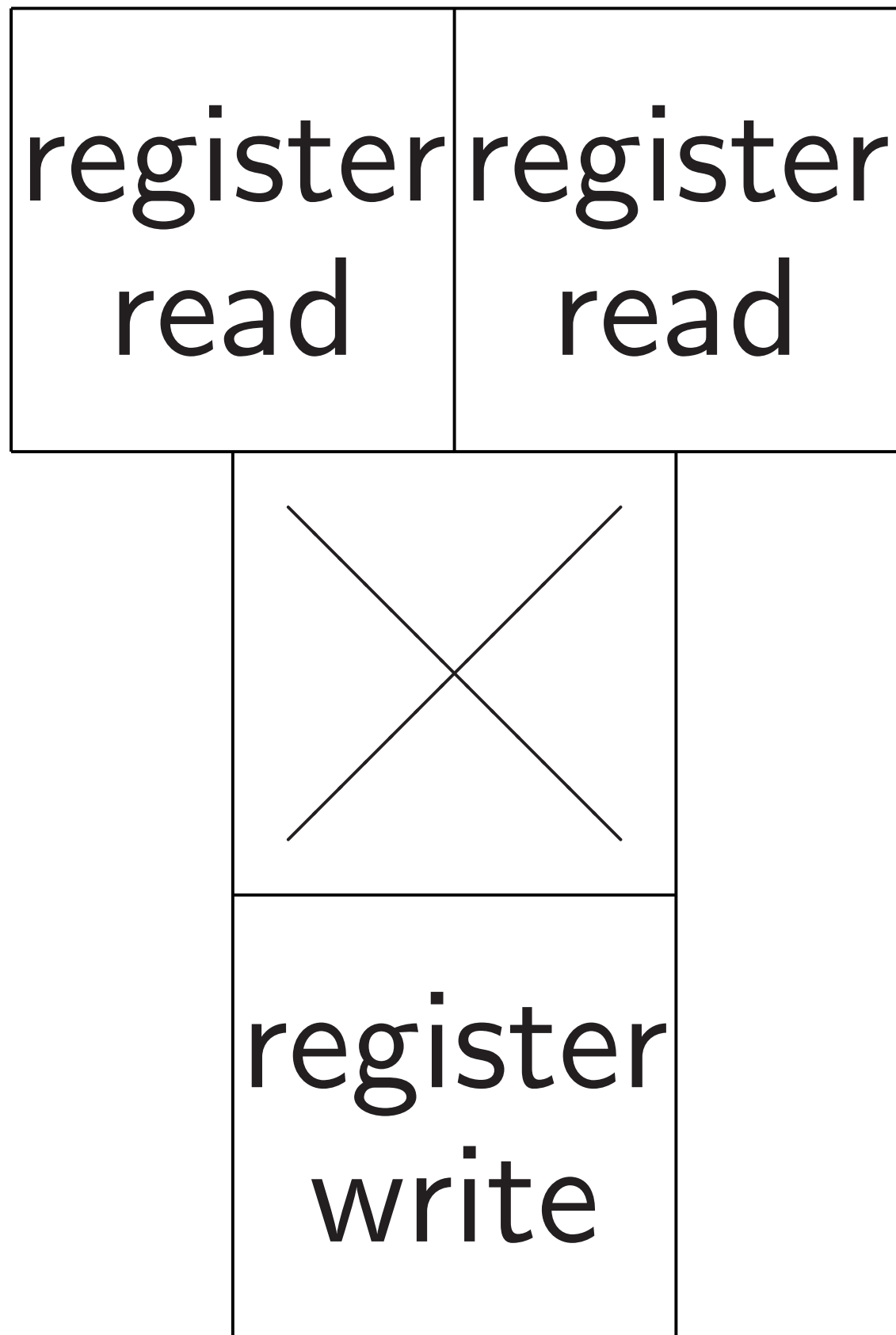
More (but slower) storage:

“load” from and “store” to
 larger “RAM” arrays.

“Instruction fetch”:

$p \mapsto o_p, i_p, j_p, k_p, p'$.

$r_0, \dots, r_{15}, i, j, k \mapsto r'_0, \dots, r'_{15}$
 where $r'_\ell = r_\ell$ except $r'_i = r_j r_k$:



Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

More (but slower) storage:

“load” from and “store” to
 larger “RAM” arrays.

“Instruction fetch”:

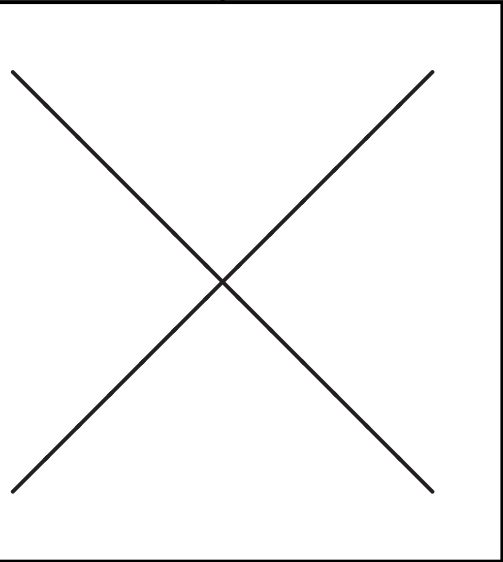
$p \mapsto o_p, i_p, j_p, k_p, p'$.

“Instruction decode”:

decompression of compressed
 format for o_p, i_p, j_p, k_p, p' .

$15, i, j, k \mapsto r'_0, \dots, r'_{15}$
 $= r_\ell$ except $r'_i = r_j r_k$:

store	register
d	read



register
write

Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

More (but slower) storage:

“load” from and “store” to
 larger “RAM” arrays.

“Instruction fetch”:

$p \mapsto o_p, i_p, j_p, k_p, p'$.

“Instruction decode”:

decompression of compressed
 format for o_p, i_p, j_p, k_p, p' .

Build “fl
 storing (

Hook (p
 flip-flops

Hook ou
 into the

At each
 flip-flops

with the

Clock ne
 for elect

all the w
 from flip

$\rightarrow r'_0, \dots, r'_{15}$
except $r'_i = r_j r_k$:

register
read

er

25

Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

More (but slower) storage:

“load” from and “store” to
larger “RAM” arrays.

“Instruction fetch”:

$p \mapsto o_p, i_p, j_p, k_p, p'$.

“Instruction decode”:

decompression of compressed
format for o_p, i_p, j_p, k_p, p' .

26

Build “flip-flops”
storing $(p, r_0, \dots,$

Hook (p, r_0, \dots, r_{15})
flip-flops into circuit.

Hook outputs $(p',$
into the same flip-

At each “clock tick”
flip-flops are overw-
with the outputs.

Clock needs to be
for electricity to p-
all the way throug-
from flip-flops to f-

Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

More (but slower) storage:

“load” from and “store” to
larger “RAM” arrays.

“Instruction fetch”:

$p \mapsto o_p, i_p, j_p, k_p, p'$.

“Instruction decode”:

decompression of compressed
format for o_p, i_p, j_p, k_p, p' .

Build “flip-flops”

storing (p, r_0, \dots, r_{15}) .

Hook (p, r_0, \dots, r_{15})

flip-flops into circuit inputs.

Hook outputs $(p', r'_0, \dots, r'_{15})$

into the same flip-flops.

At each “clock tick”,
flip-flops are overwritten
with the outputs.

Clock needs to be slow enough
for electricity to percolate
all the way through the circuit
from flip-flops to flip-flops.

Add more flexibility.

More arithmetic:

replace (i, j, k) with

$(“\times”, i, j, k)$ and

$(“+”, i, j, k)$ and more options.

More (but slower) storage:

“load” from and “store” to
larger “RAM” arrays.

“Instruction fetch”:

$p \mapsto o_p, i_p, j_p, k_p, p'$.

“Instruction decode”:

decompression of compressed
format for o_p, i_p, j_p, k_p, p' .

Build “flip-flops”

storing (p, r_0, \dots, r_{15}) .

Hook (p, r_0, \dots, r_{15})

flip-flops into circuit inputs.

Hook outputs $(p', r'_0, \dots, r'_{15})$

into the same flip-flops.

At each “clock tick”,

flip-flops are overwritten

with the outputs.

Clock needs to be slow enough

for electricity to percolate

all the way through the circuit,

from flip-flops to flip-flops.

re flexibility.

ithmetic:

(i, j, k) with

(j, k) and

(j, k) and more options.

(but slower) storage:

rom and “store” to

RAM” arrays.

tion fetch”:

i_p, j_p, k_p, p' .

tion decode”:

ression of compressed

or o_p, i_p, j_p, k_p, p' .

Build “flip-flops”

storing (p, r_0, \dots, r_{15}) .

Hook (p, r_0, \dots, r_{15})

flip-flops into circuit inputs.

Hook outputs $(p', r'_0, \dots, r'_{15})$

into the same flip-flops.

At each “clock tick”,

flip-flops are overwritten

with the outputs.

Clock needs to be slow enough

for electricity to percolate

all the way through the circuit,

from flip-flops to flip-flops.

Now have



Further

e.g., log

Build “flip-flops”
storing (p, r_0, \dots, r_{15}) .

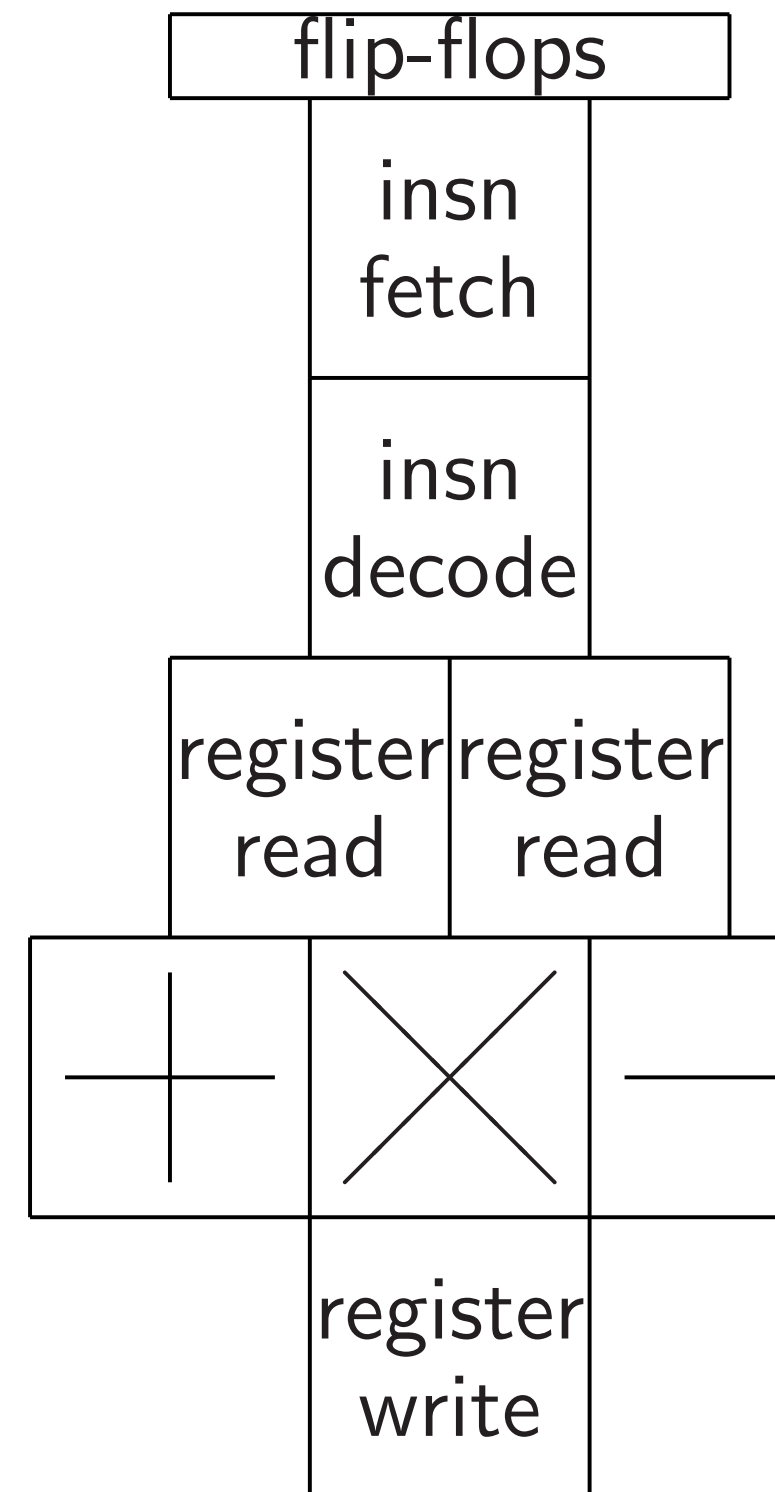
Hook (p, r_0, \dots, r_{15})
flip-flops into circuit inputs.

Hook outputs $(p', r'_0, \dots, r'_{15})$
into the same flip-flops.

At each “clock tick”,
flip-flops are overwritten
with the outputs.

Clock needs to be slow enough
for electricity to percolate
all the way through the circuit,
from flip-flops to flip-flops.

Now have semi-flip-flops



Further flexibility in
e.g., logic instructions

Build “flip-flops”
storing (p, r_0, \dots, r_{15}) .

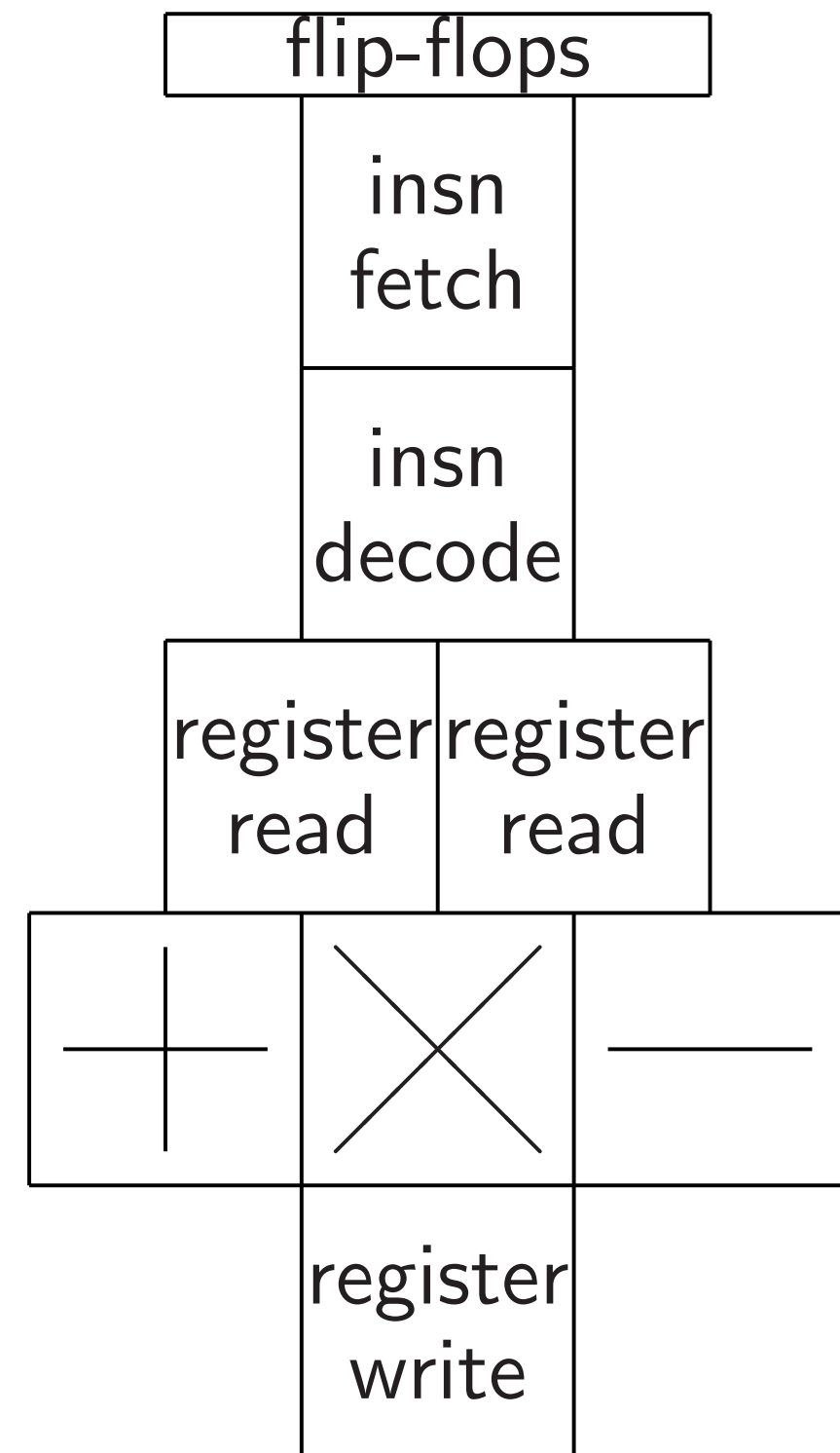
Hook (p, r_0, \dots, r_{15})
flip-flops into circuit inputs.

Hook outputs $(p', r'_0, \dots, r'_{15})$
into the same flip-flops.

At each “clock tick”,
flip-flops are overwritten
with the outputs.

Clock needs to be slow enough
for electricity to percolate
all the way through the circuit,
from flip-flops to flip-flops.

Now have semi-flexible CPU



Further flexibility is useful:
e.g., logic instructions.

Build “flip-flops”

storing (p, r_0, \dots, r_{15}) .

Hook (p, r_0, \dots, r_{15})

flip-flops into circuit inputs.

Hook outputs $(p', r'_0, \dots, r'_{15})$

into the same flip-flops.

At each “clock tick”,

flip-flops are overwritten

with the outputs.

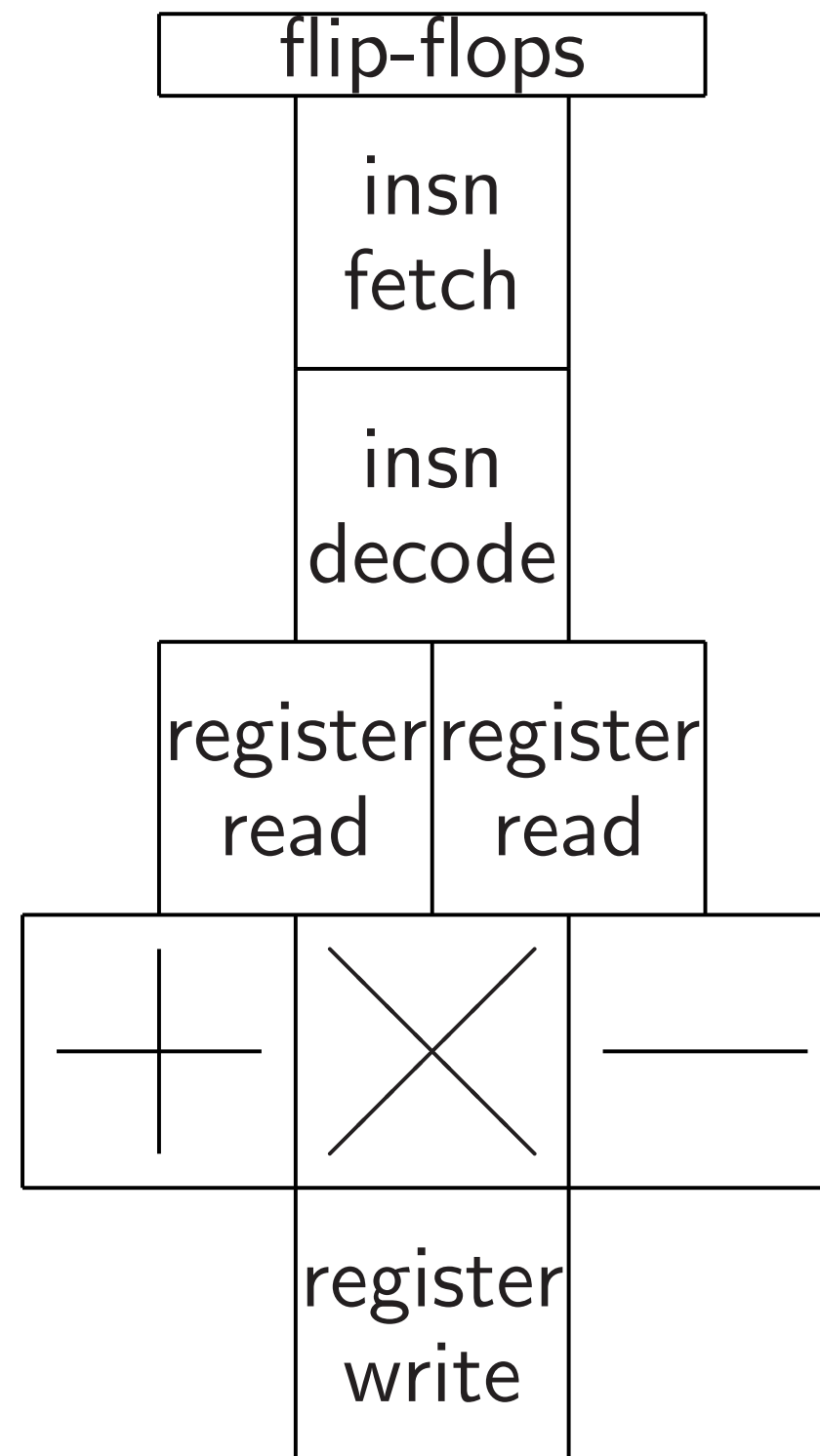
Clock needs to be slow enough

for electricity to percolate

all the way through the circuit,

from flip-flops to flip-flops.

Now have semi-flexible CPU:



Further flexibility is useful:

e.g., logic instructions.

flip-flops”

(p, r_0, \dots, r_{15}) .

(p, r_0, \dots, r_{15})

into circuit inputs.

outputs $(p', r'_0, \dots, r'_{15})$

same flip-flops.

“clock tick”,

are overwritten

outputs.

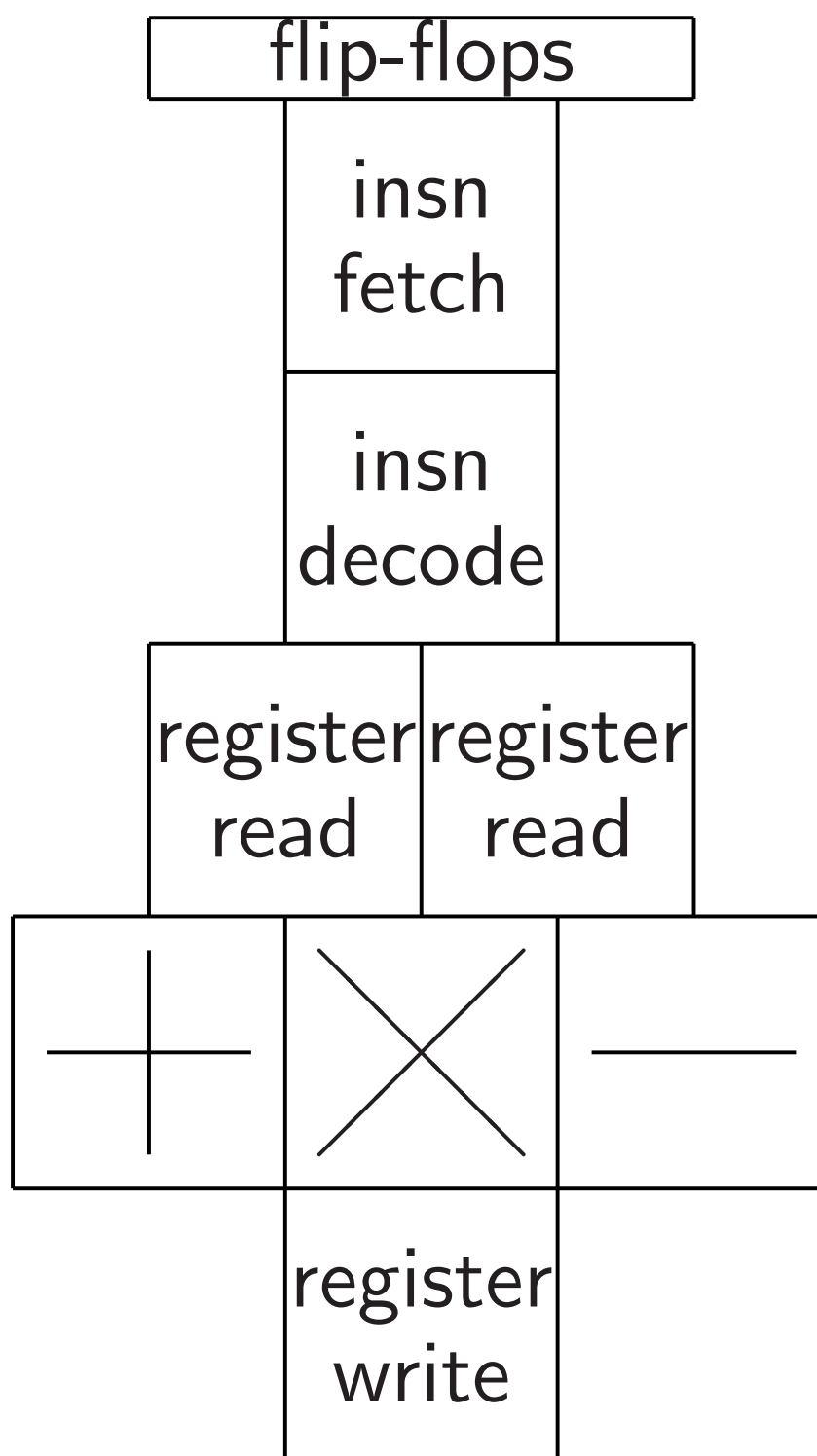
needs to be slow enough

ricity to percolate

way through the circuit,

o-flops to flip-flops.

Now have semi-flexible CPU:

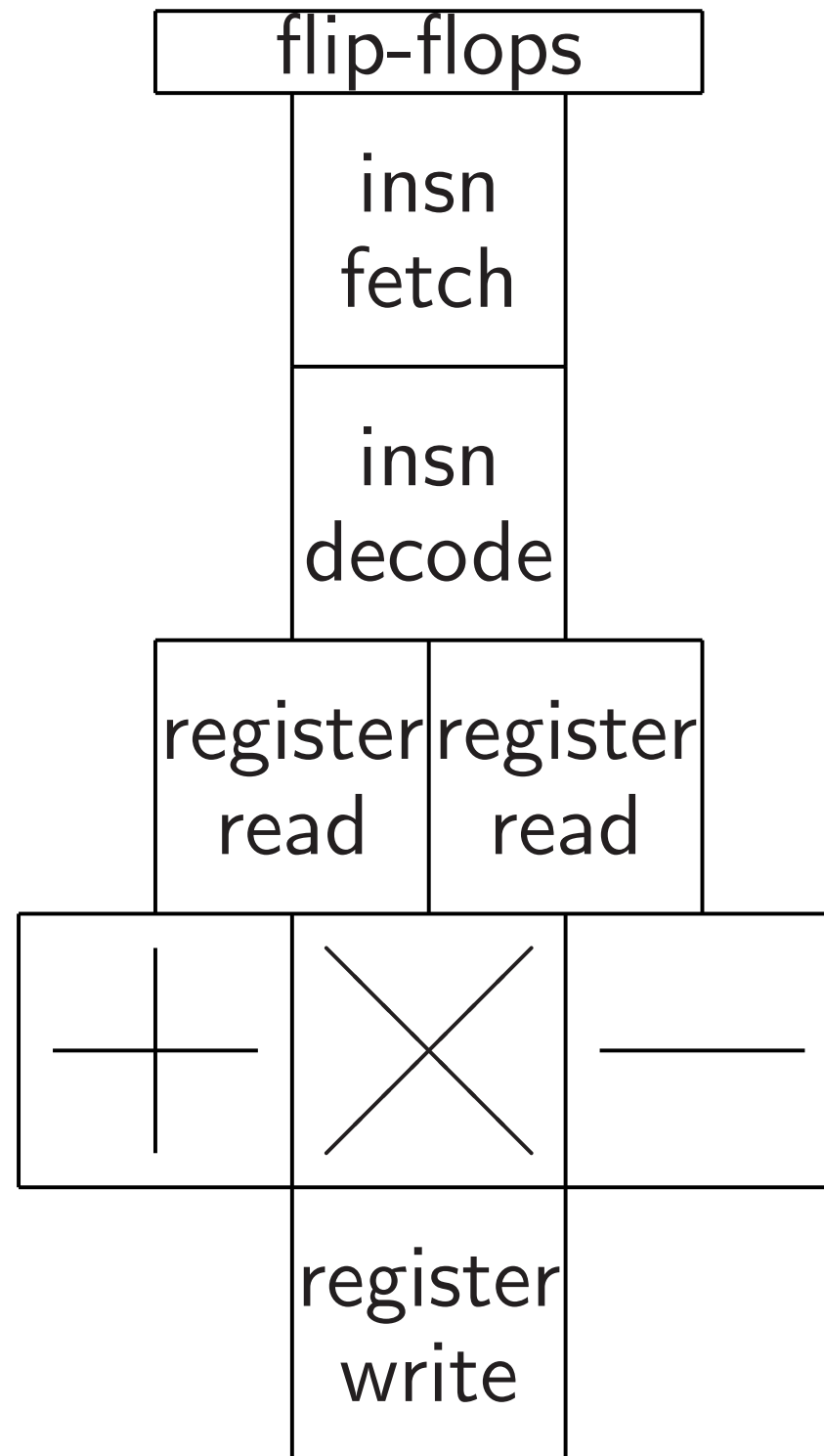


Further flexibility is useful:
e.g., logic instructions.

“Pipelined

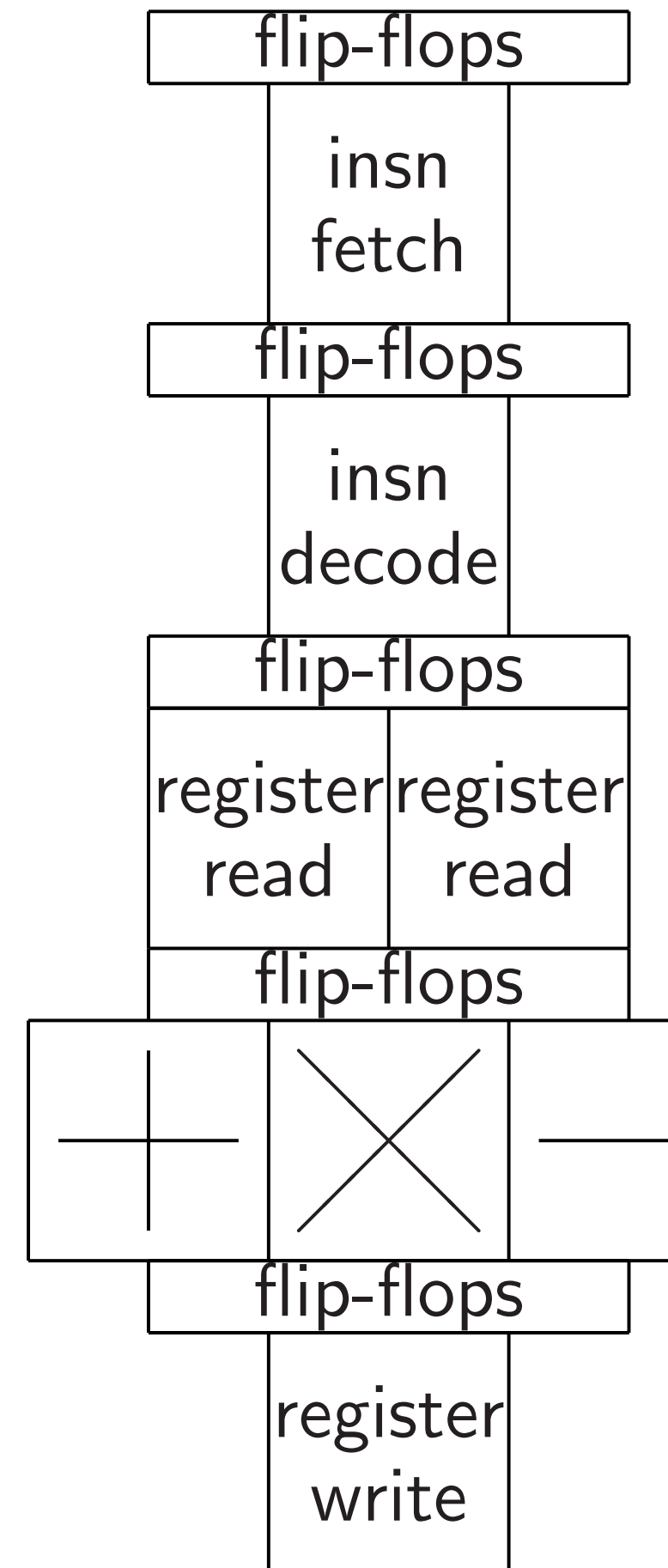


Now have semi-flexible CPU:

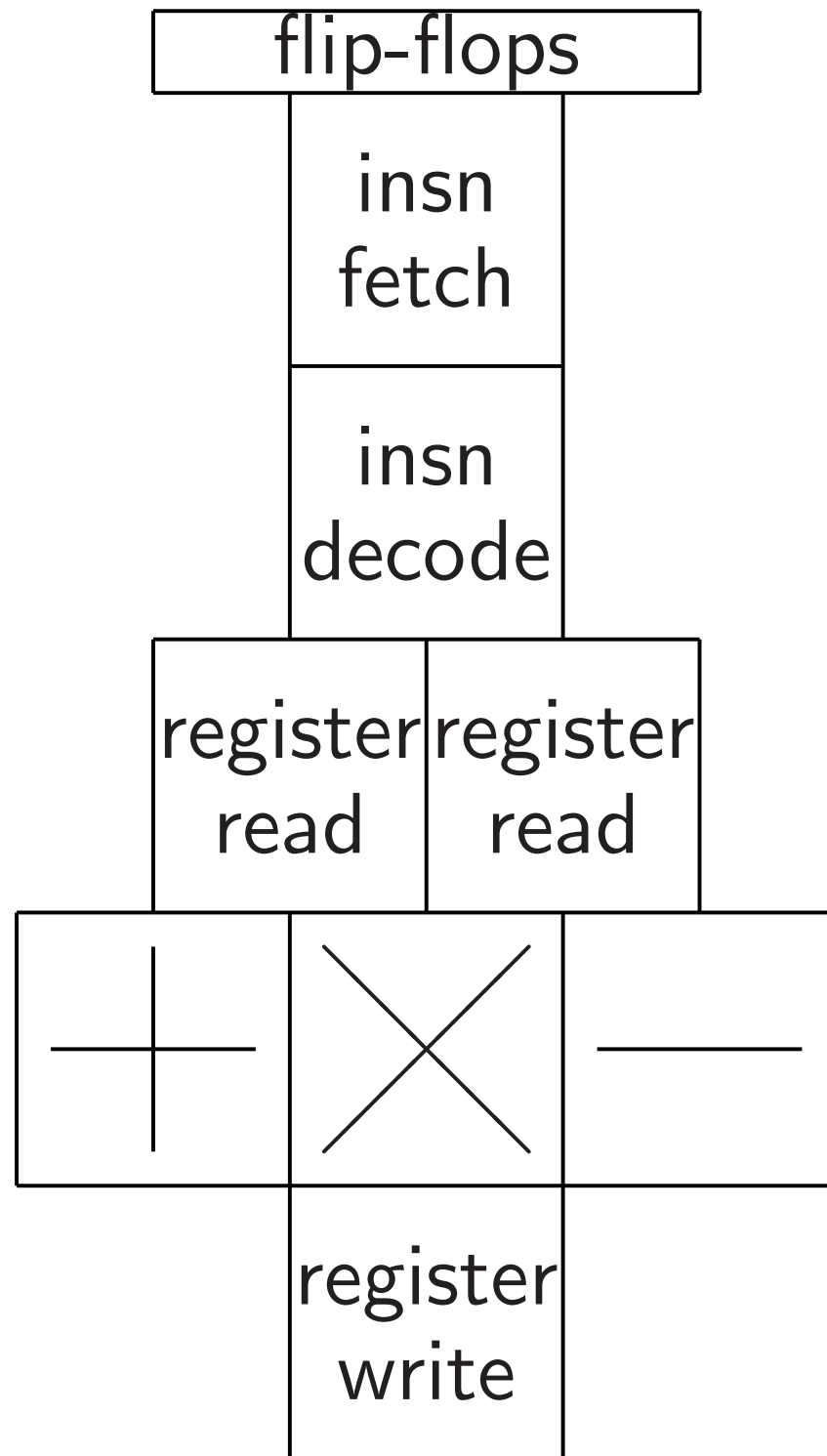


Further flexibility is useful:
e.g., logic instructions.

“Pipelining” allow

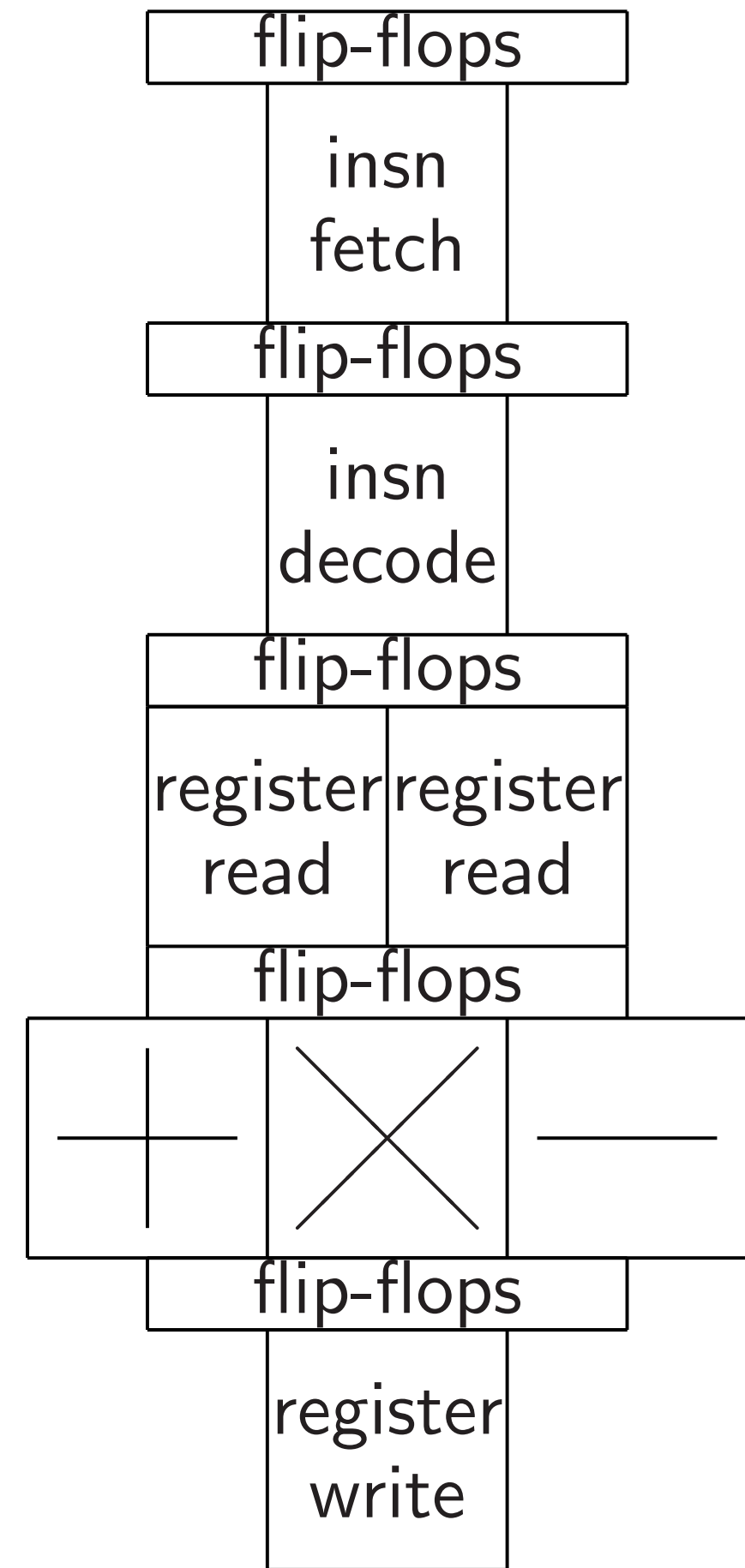


Now have semi-flexible CPU:



Further flexibility is useful:
e.g., logic instructions.

“Pipelining” allows faster clock



stage

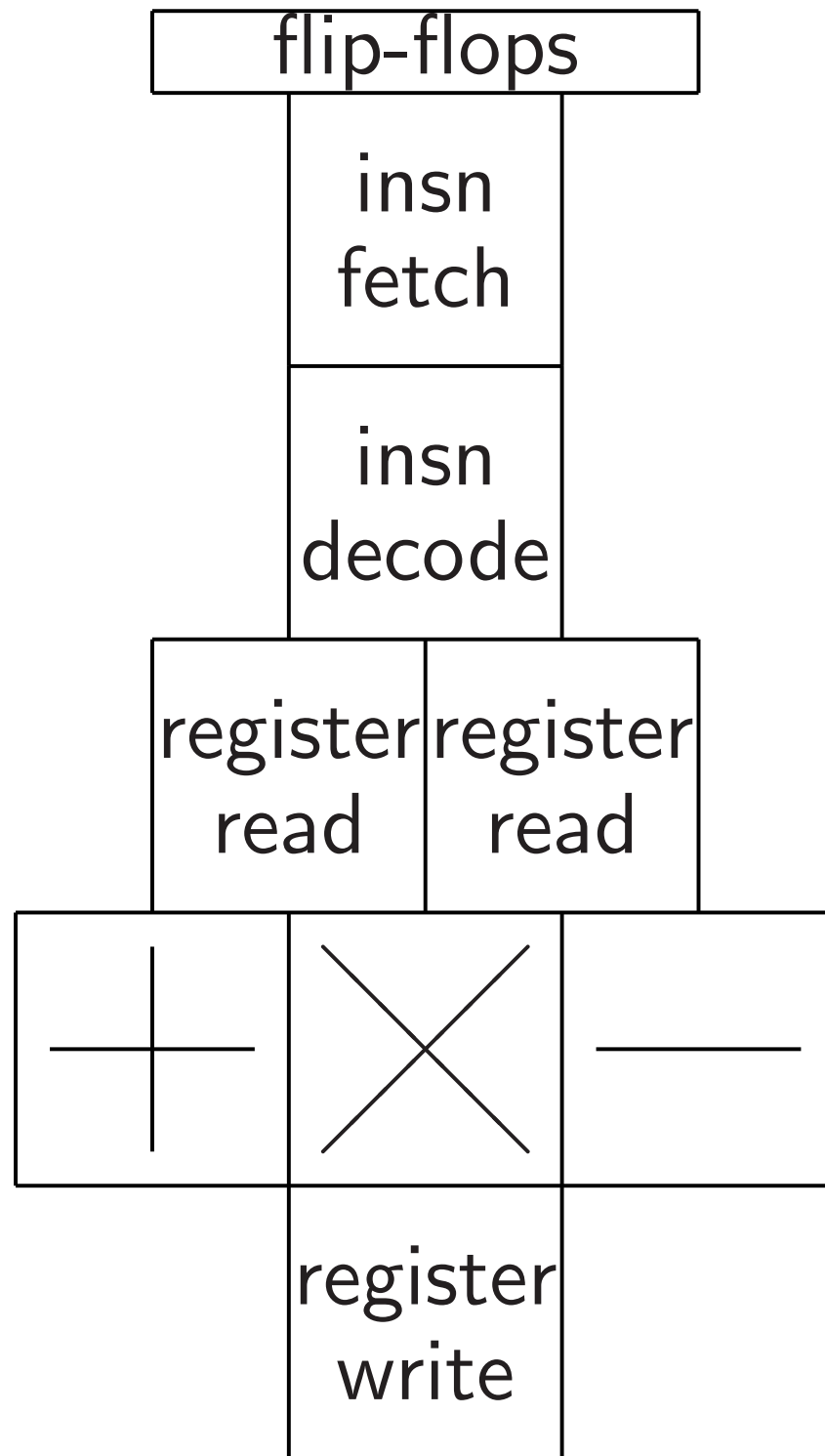
stage

stage

stage

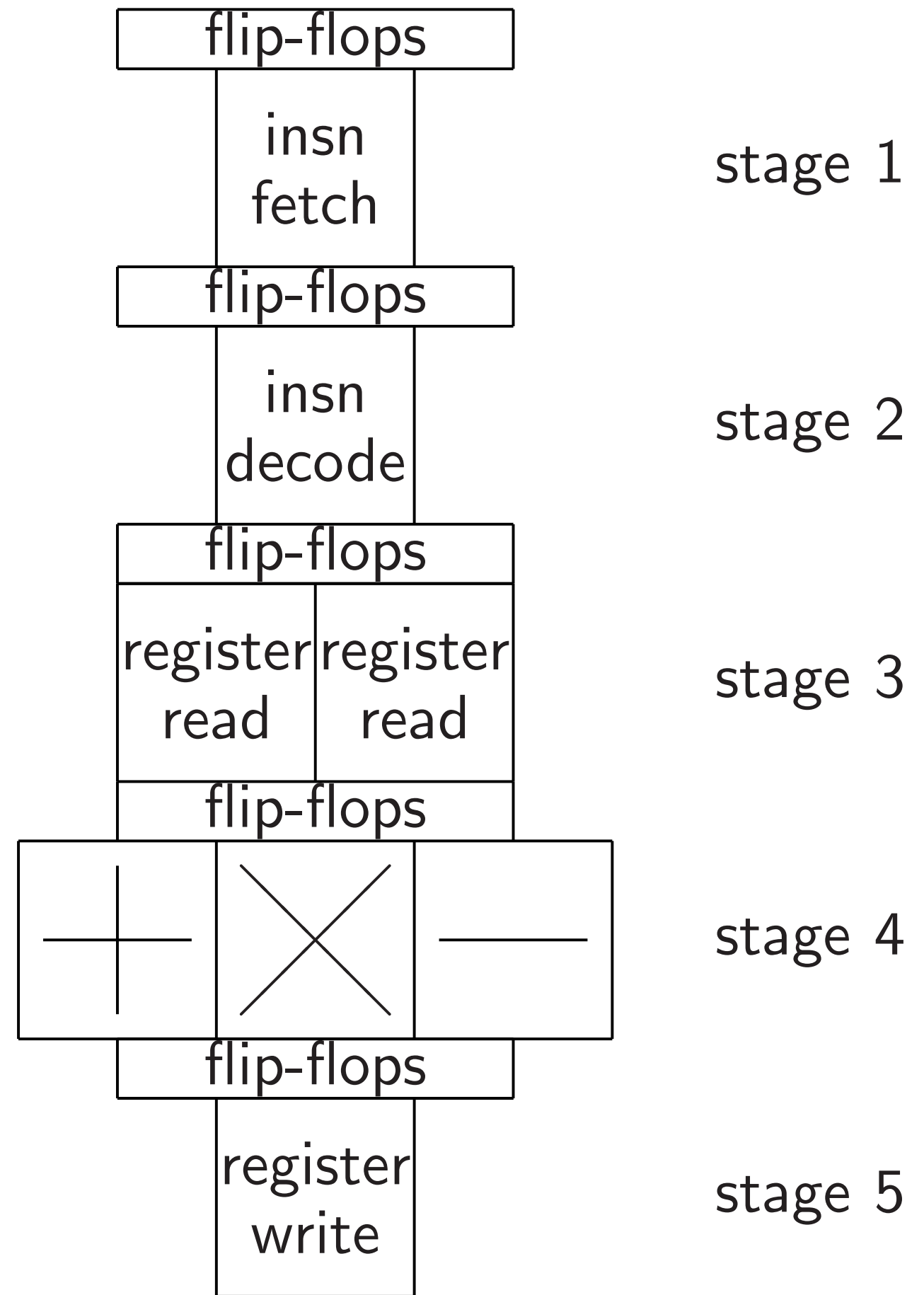
stage

Now have semi-flexible CPU:

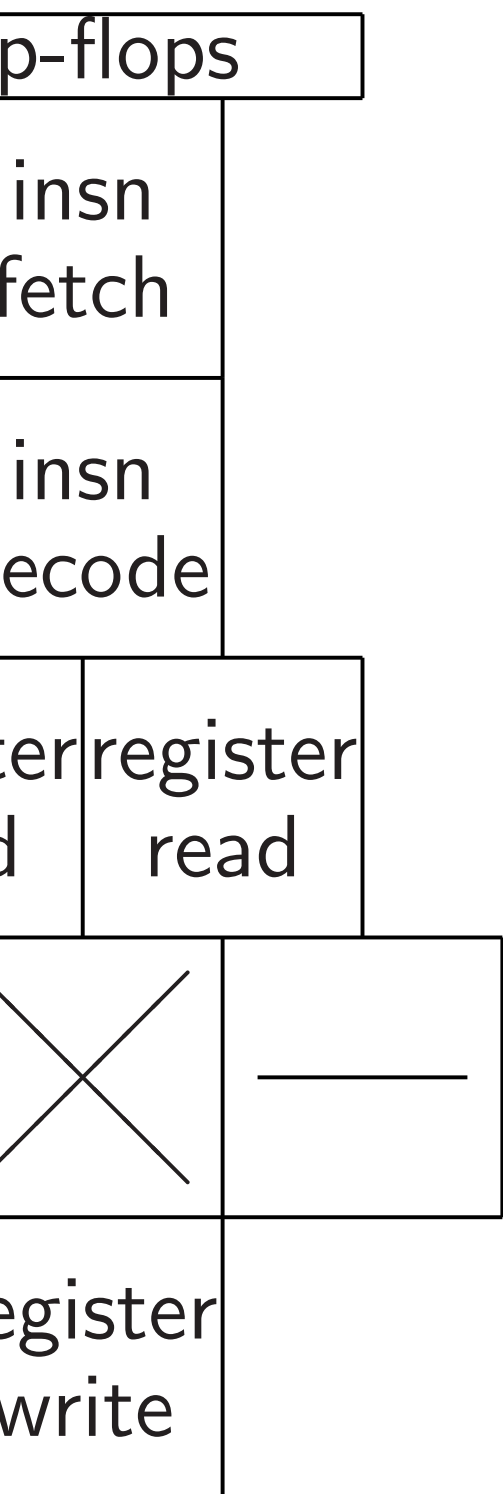


Further flexibility is useful:
e.g., logic instructions.

“Pipelining” allows faster clock:

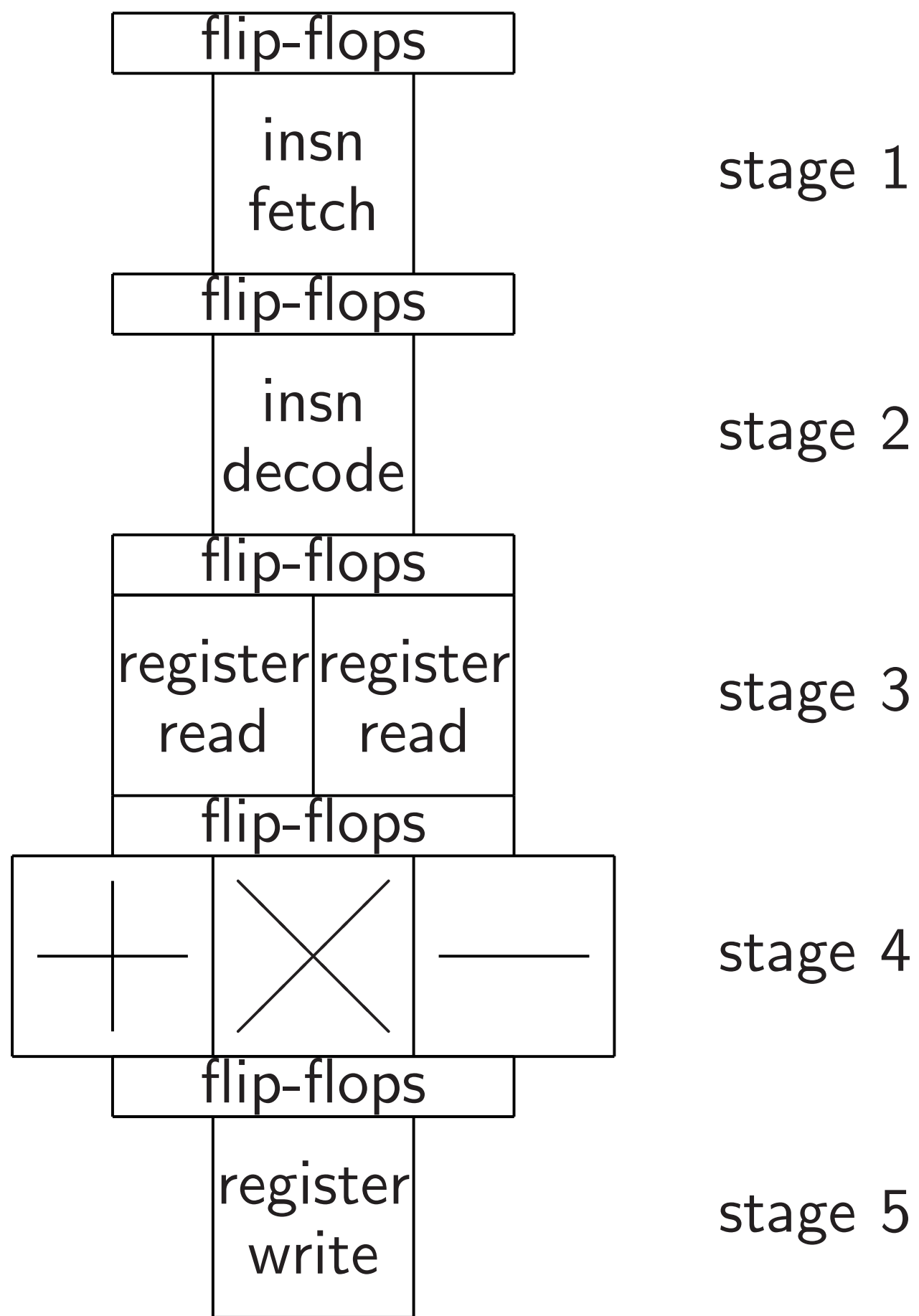


ve semi-flexible CPU:



flexibility is useful:
ic instructions.

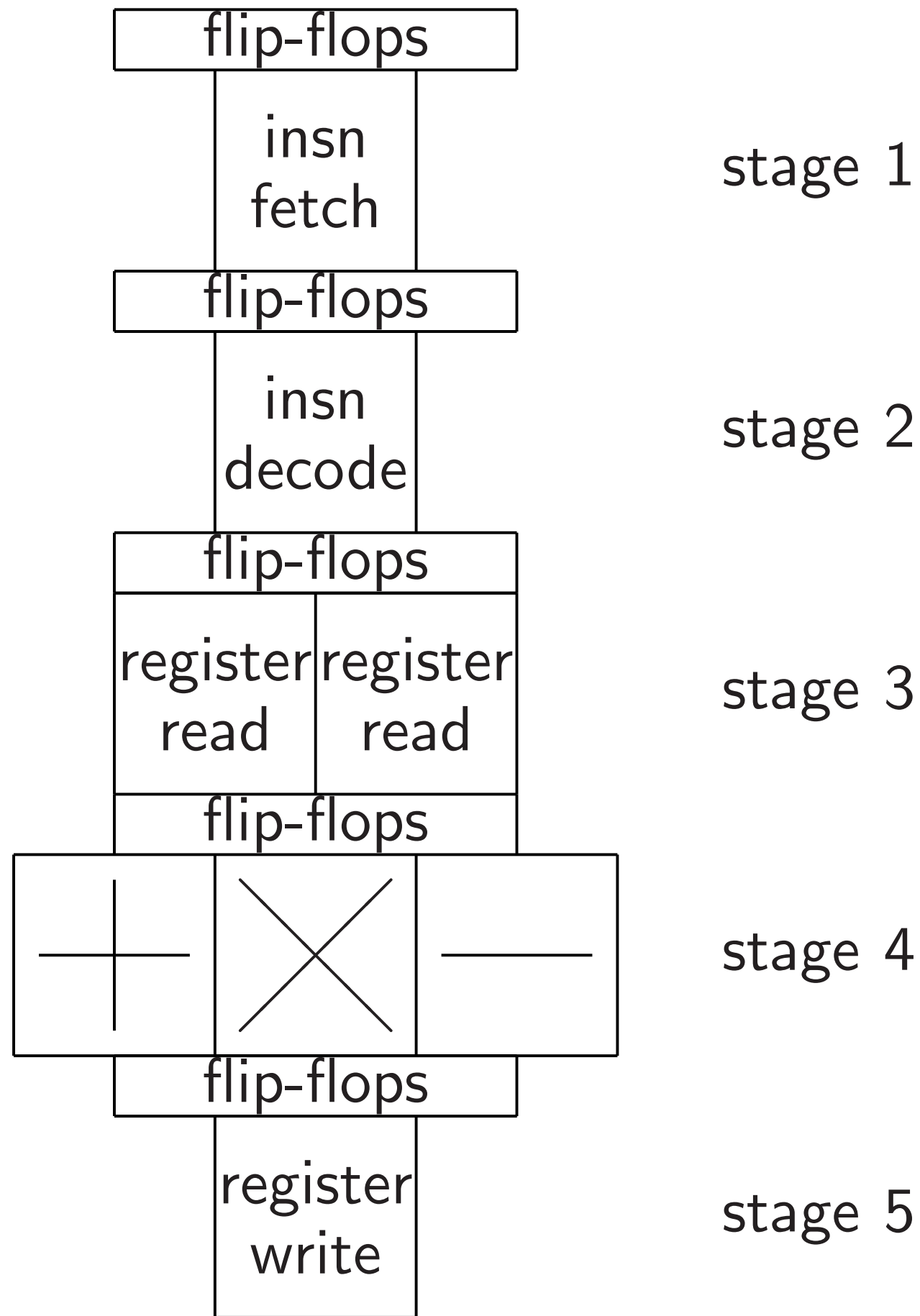
“Pipelining” allows faster clock:



Goal: St
one tick
Instructi
reads ne
feeds p'
After ne
instructi
uncompr
while ins
reads an
Some ex
Also ext
preserve
e.g., sta

flexible CPU:

“Pipelining” allows faster clock:



Goal: Stage n handles one tick after stage

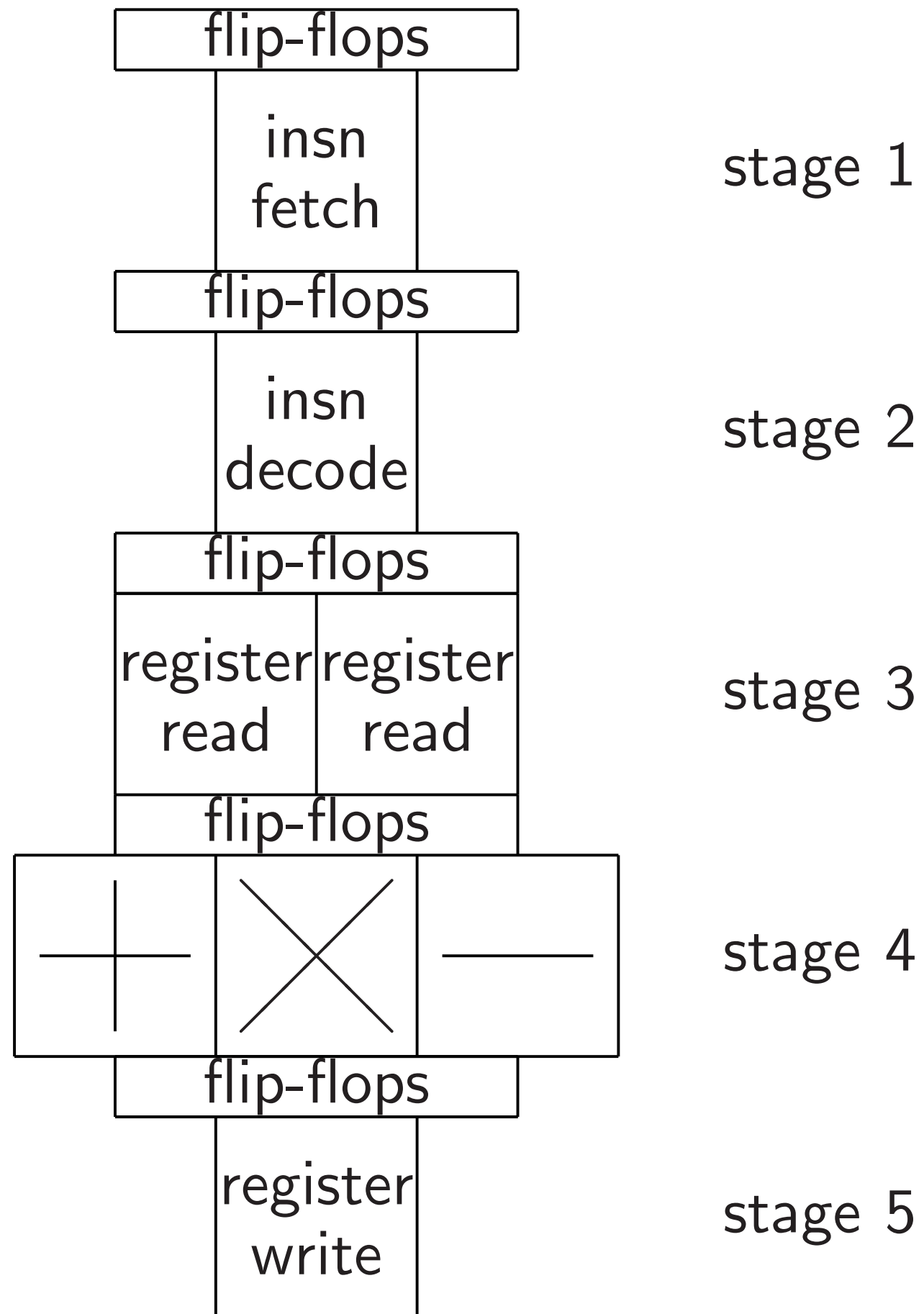
Instruction fetch reads next instruction feeds p' back, sends

After next clock tick instruction decode uncompresses this while instruction fetch reads another instruction

Some extra flip-flops Also extra area to preserve instructions e.g., stall on read-

is useful:
ions.

“Pipelining” allows faster clock:



Goal: Stage n handles instruction one tick after stage $n - 1$.

Instruction fetch

reads next instruction, feeds p' back, sends instruction

After next clock tick,

instruction decode

uncompresses this instruction

while instruction fetch

reads another instruction.

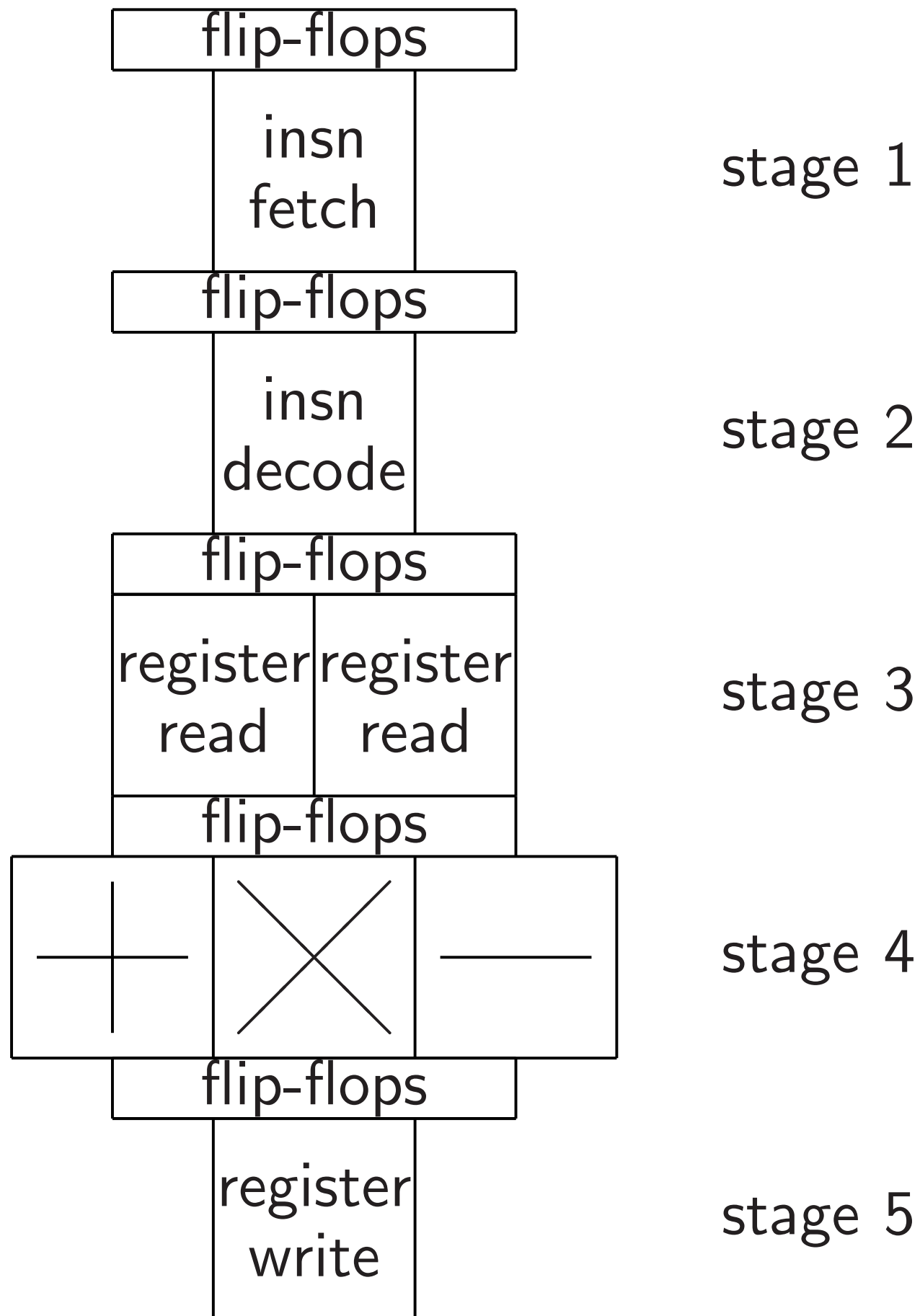
Some extra flip-flop area.

Also extra area to

preserve instruction semantics

e.g., stall on read-after-write

“Pipelining” allows faster clock:



Goal: Stage n handles instruction one tick after stage $n - 1$.

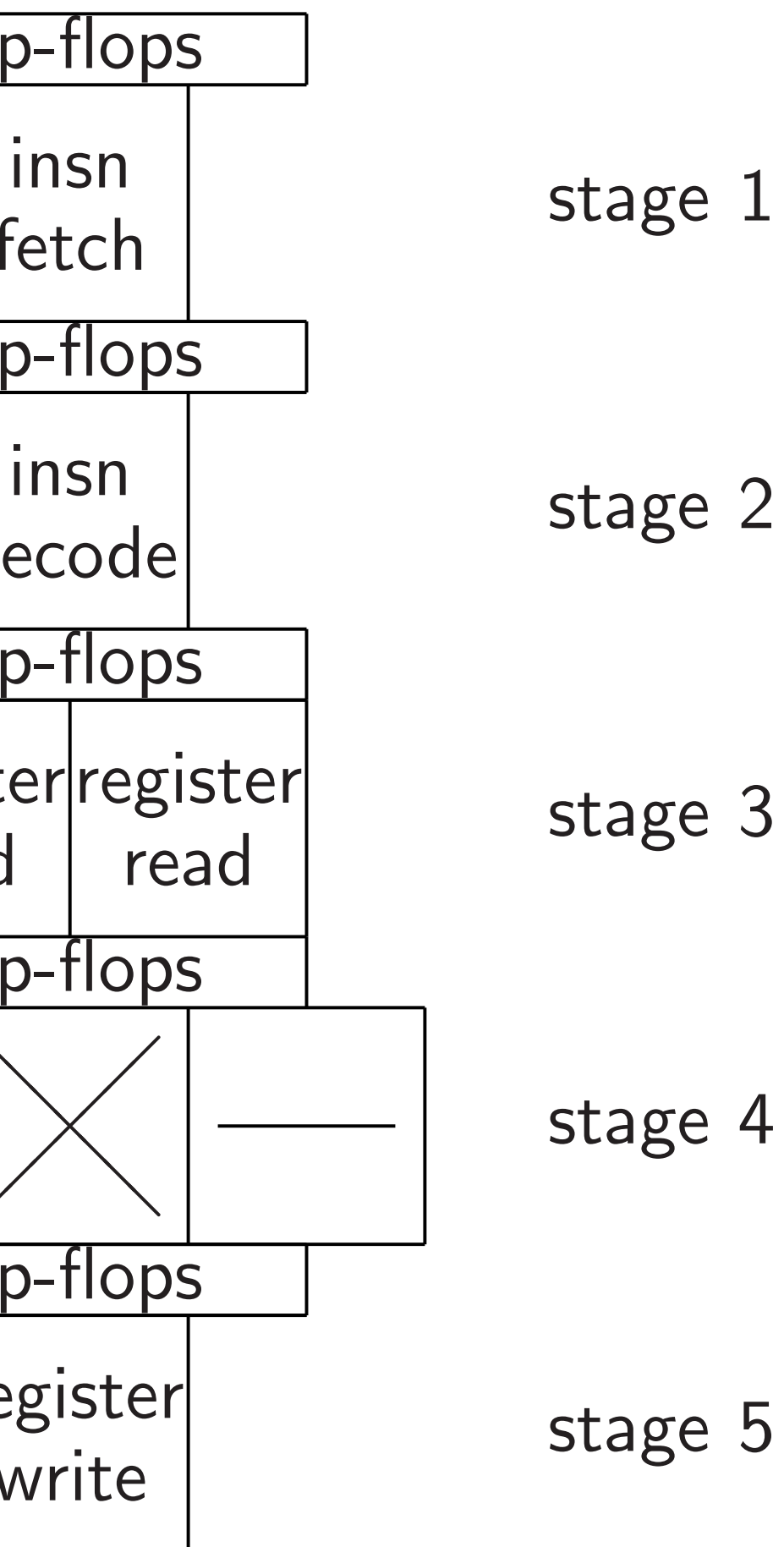
Instruction fetch reads next instruction, feeds p' back, sends instruction.

After next clock tick, instruction decode uncompresses this instruction, while instruction fetch reads another instruction.

Some extra flip-flop area.

Also extra area to preserve instruction semantics: e.g., stall on read-after-write.

“pipelining” allows faster clock:



Goal: Stage n handles instruction one tick after stage $n - 1$.

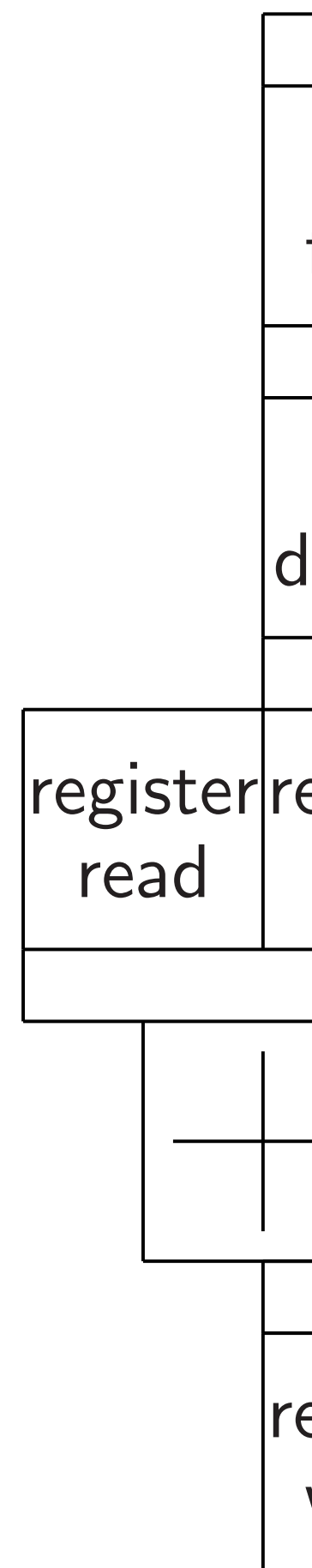
Instruction fetch reads next instruction, feeds p' back, sends instruction.

After next clock tick, instruction decode uncompresses this instruction, while instruction fetch reads another instruction.

Some extra flip-flop area.

Also extra area to preserve instruction semantics: e.g., stall on read-after-write.

“Superscalar”



s faster clock:

stage 1

Goal: Stage n handles instruction one tick after stage $n - 1$.

stage 2

Instruction fetch
reads next instruction,
feeds p' back, sends instruction.

stage 3

After next clock tick,
instruction decode
uncompresses this instruction,
while instruction fetch
reads another instruction.

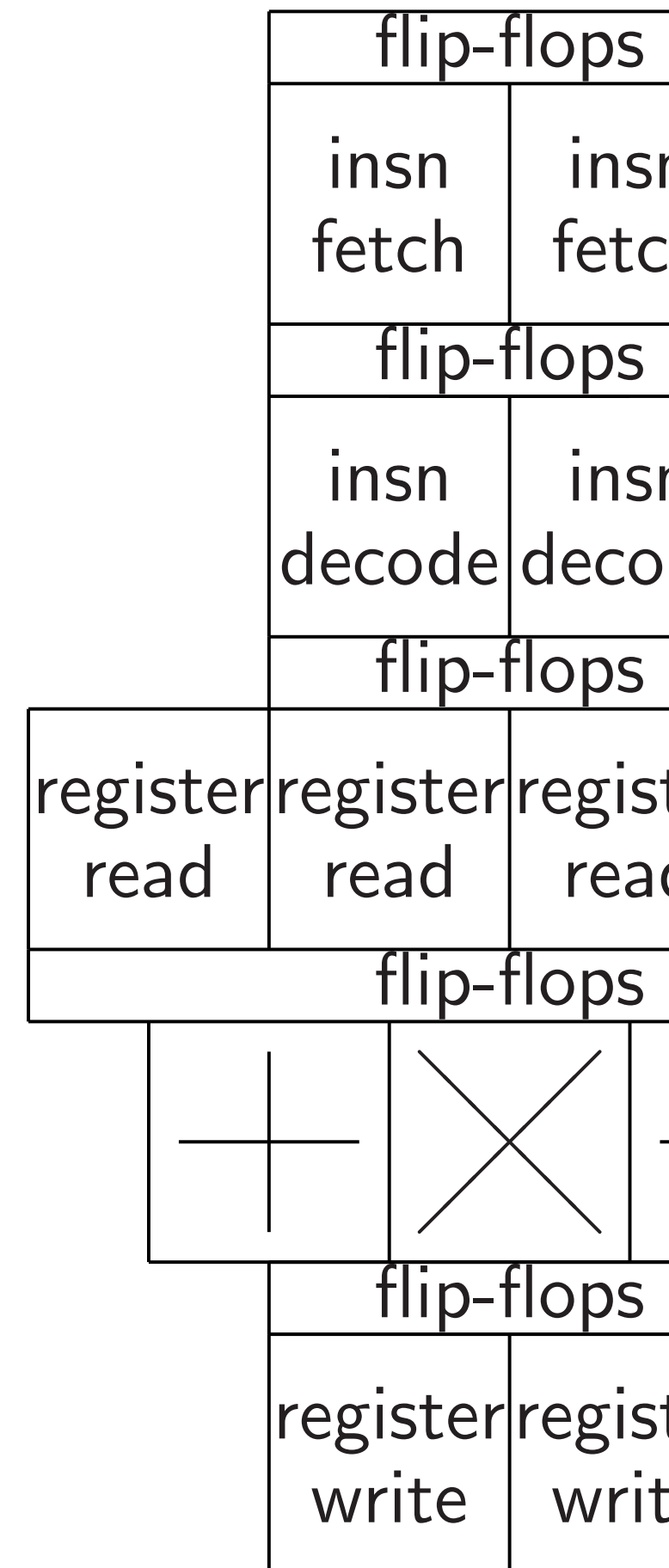
stage 4

Some extra flip-flop area.

stage 5

Also extra area to
preserve instruction semantics:
e.g., stall on read-after-write.

“Superscalar” proc



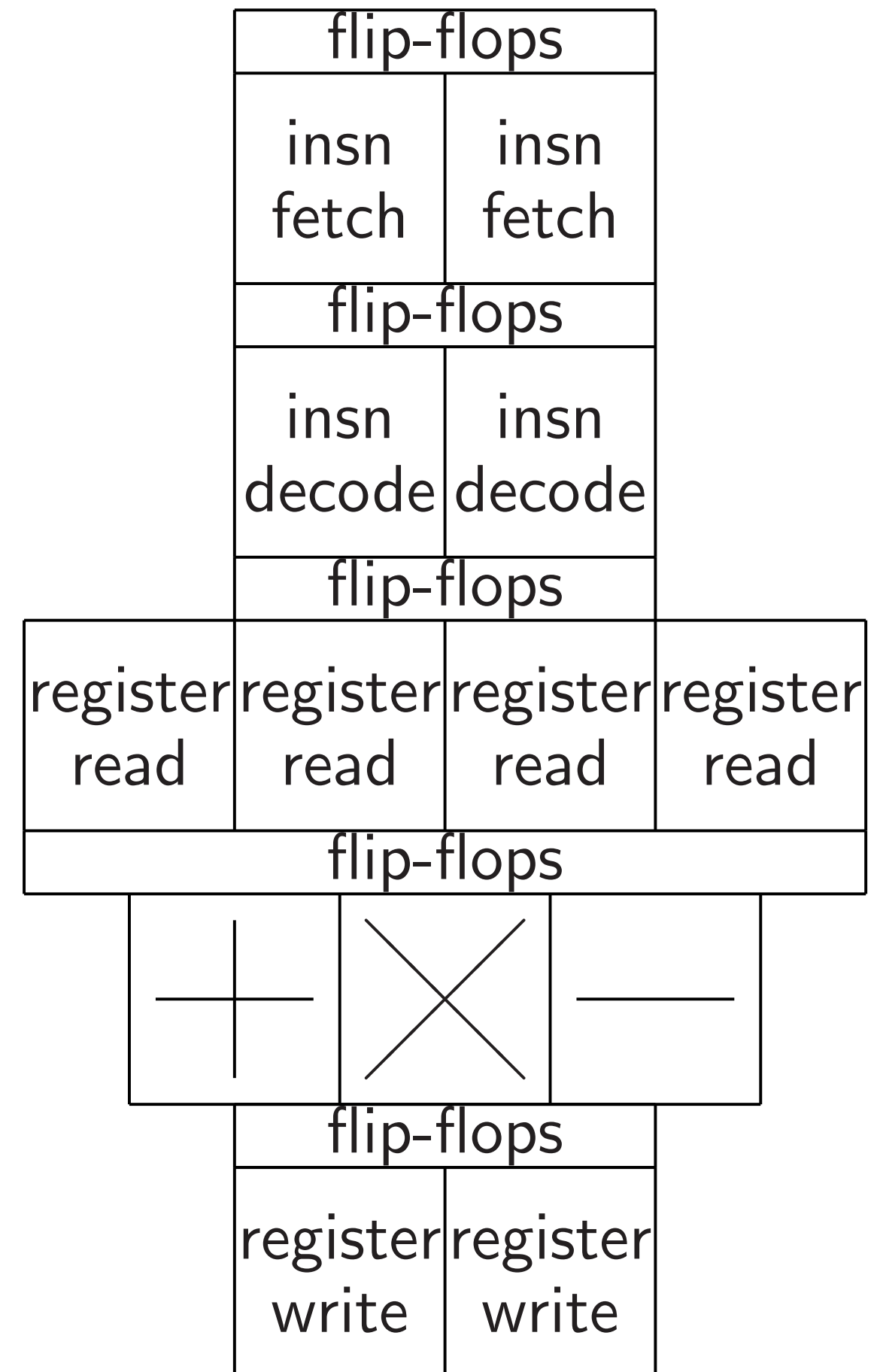
clock: Goal: Stage n handles instruction one tick after stage $n - 1$.

1 Instruction fetch
reads next instruction,
feeds p' back, sends instruction.

2 After next clock tick,
instruction decode
uncompresses this instruction,
while instruction fetch
reads another instruction.

3 Some extra flip-flop area.
Also extra area to
preserve instruction semantics:
e.g., stall on read-after-write.

“Superscalar” processing:



Goal: Stage n handles instruction one tick after stage $n - 1$.

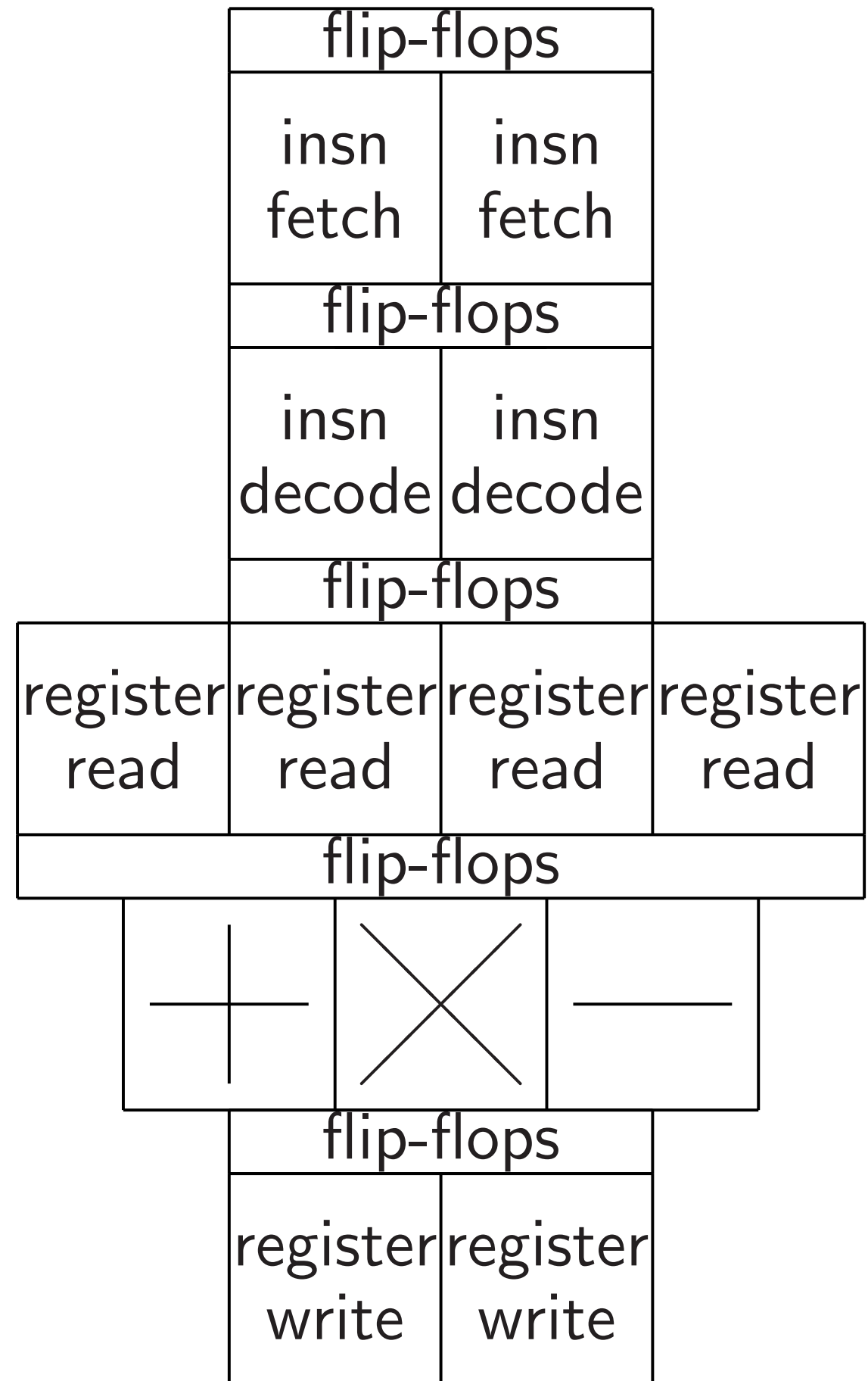
Instruction fetch
reads next instruction,
feeds p' back, sends instruction.

After next clock tick,
instruction decode
uncompresses this instruction,
while instruction fetch
reads another instruction.

Some extra flip-flop area.

Also extra area to
preserve instruction semantics:
e.g., stall on read-after-write.

“Superscalar” processing:



stage n handles instruction
after stage $n - 1$.

on fetch

xt instruction,

back, sends instruction.

xt clock tick,

on decode

resses this instruction,

struction fetch

other instruction.

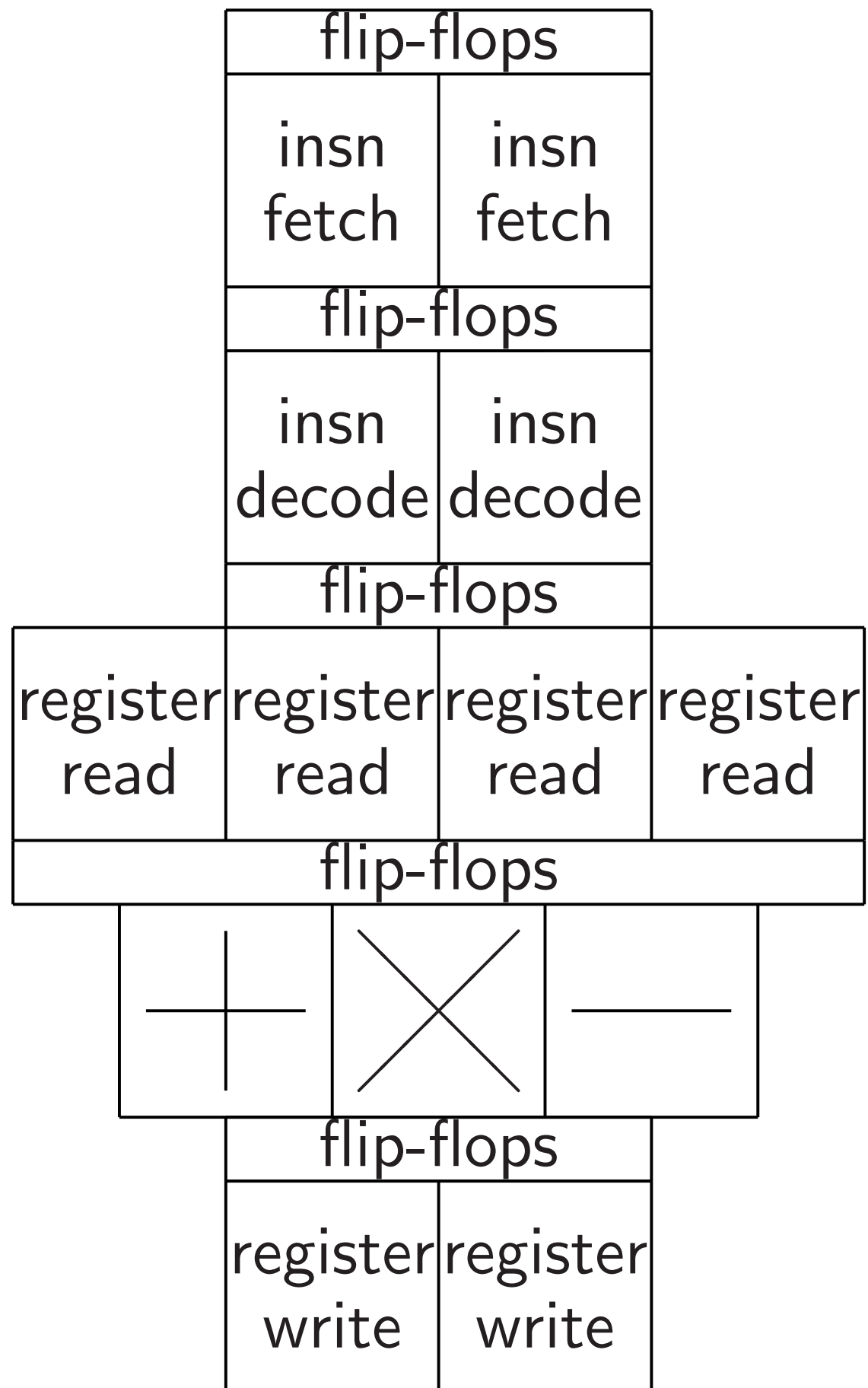
tra flip-flop area.

ra area to

instruction semantics:

ll on read-after-write.

“Superscalar” processing:



“Vector”

Expand

into n -ve

ARM “M

Intel “AV

Intel “AV

GPUs ha

handles instruction
stage $n - 1$.

tion,

ds instruction.

ck,

instruction,

etch

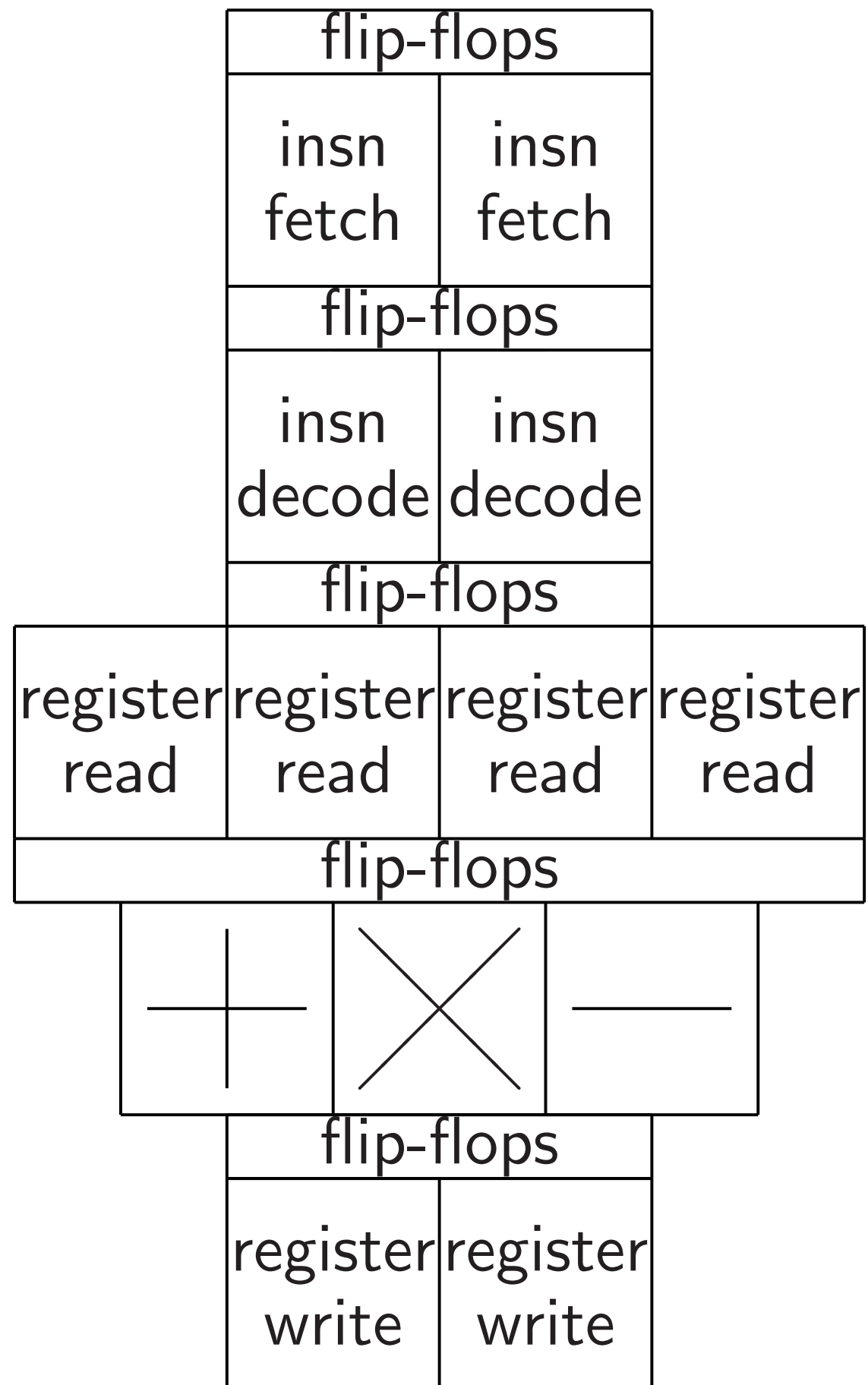
ruction.

op area.

n semantics:

after-write.

“Superscalar” processing:



“Vector” processing

Expand each 32-bit

into n -vector of 32

ARM “NEON” has

Intel “AVX2” has

Intel “AVX-512” h

GPUs have larger

action

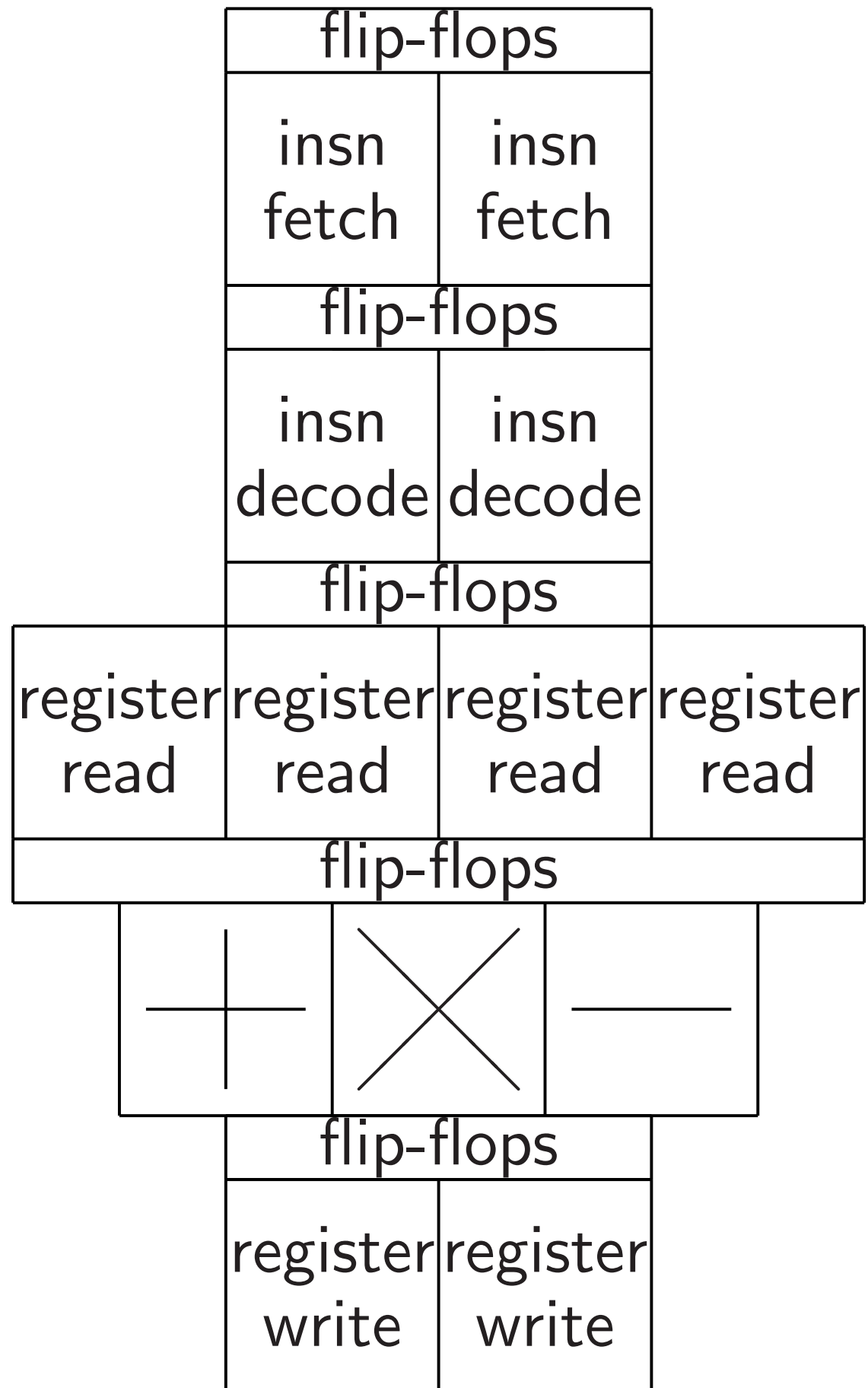
tion.

n,

CS:

e.

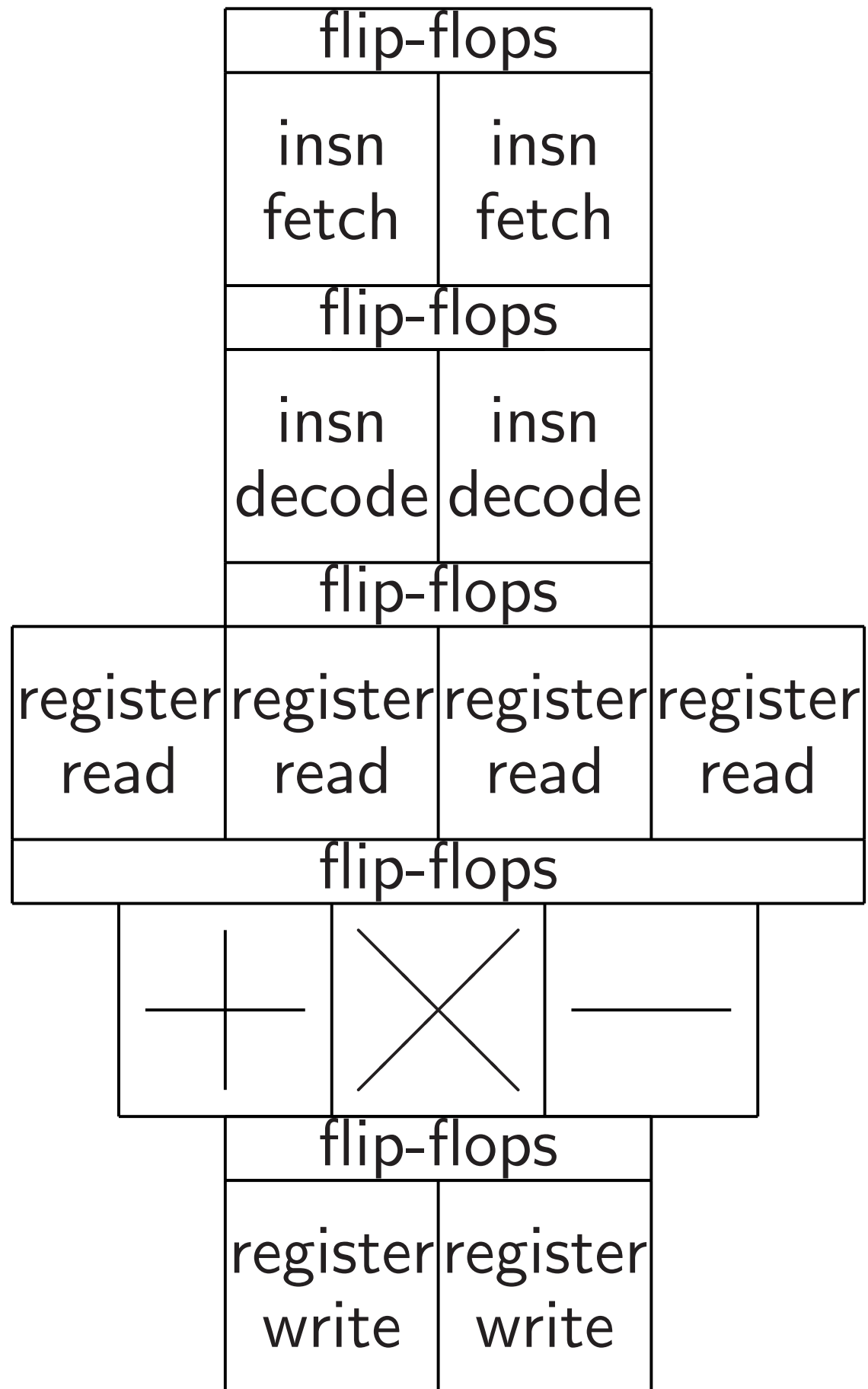
“Superscalar” processing:



“Vector” processing:

Expand each 32-bit integer into n -vector of 32-bit integers
 ARM “NEON” has $n = 4$;
 Intel “AVX2” has $n = 8$;
 Intel “AVX-512” has $n = 16$
 GPUs have larger n .

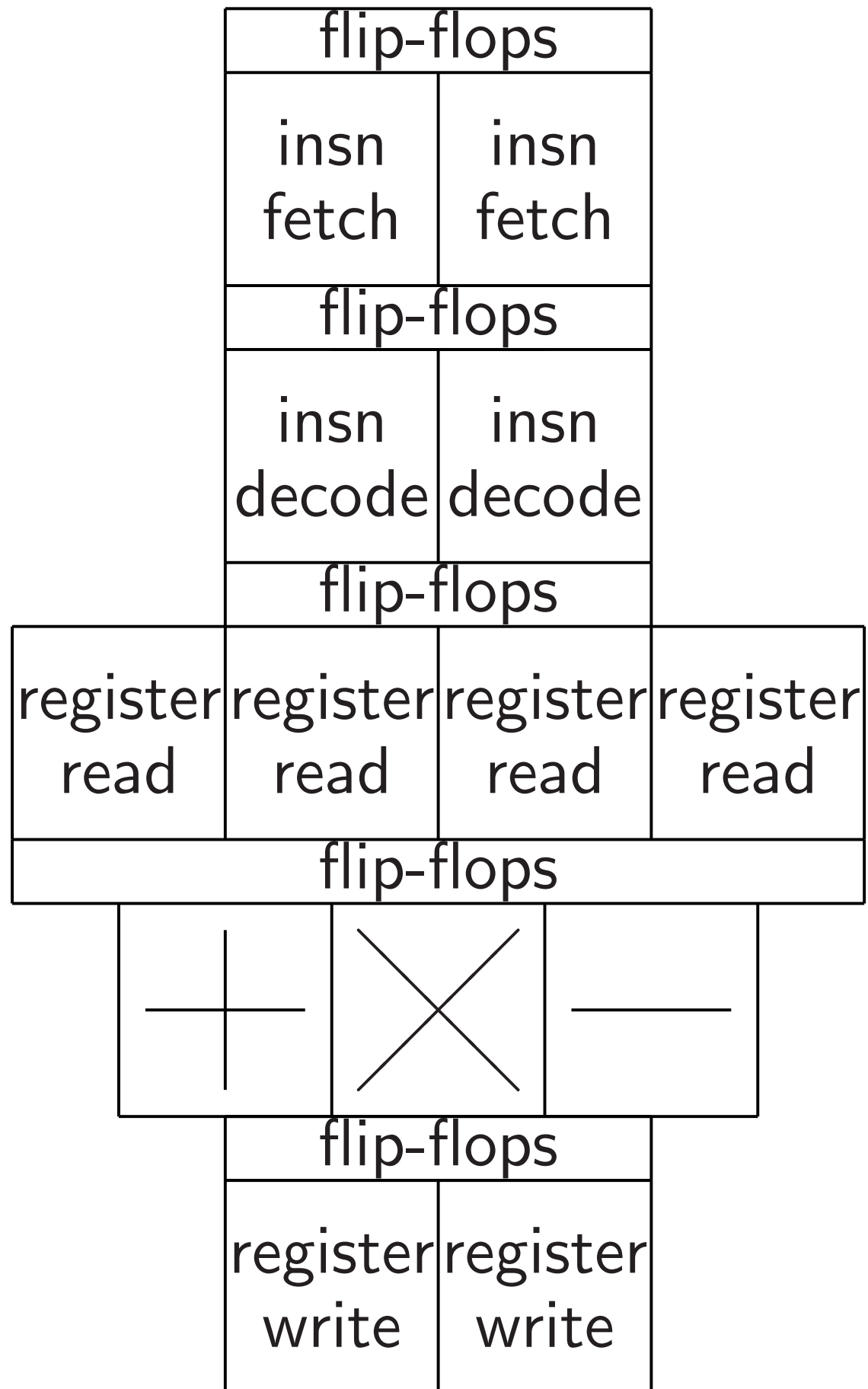
“Superscalar” processing:



“Vector” processing:

Expand each 32-bit integer into n -vector of 32-bit integers.
 ARM “NEON” has $n = 4$;
 Intel “AVX2” has $n = 8$;
 Intel “AVX-512” has $n = 16$;
 GPUs have larger n .

“Superscalar” processing:



“Vector” processing:

Expand each 32-bit integer into n -vector of 32-bit integers.
 ARM “NEON” has $n = 4$;
 Intel “AVX2” has $n = 8$;
 Intel “AVX-512” has $n = 16$;
 GPUs have larger n .

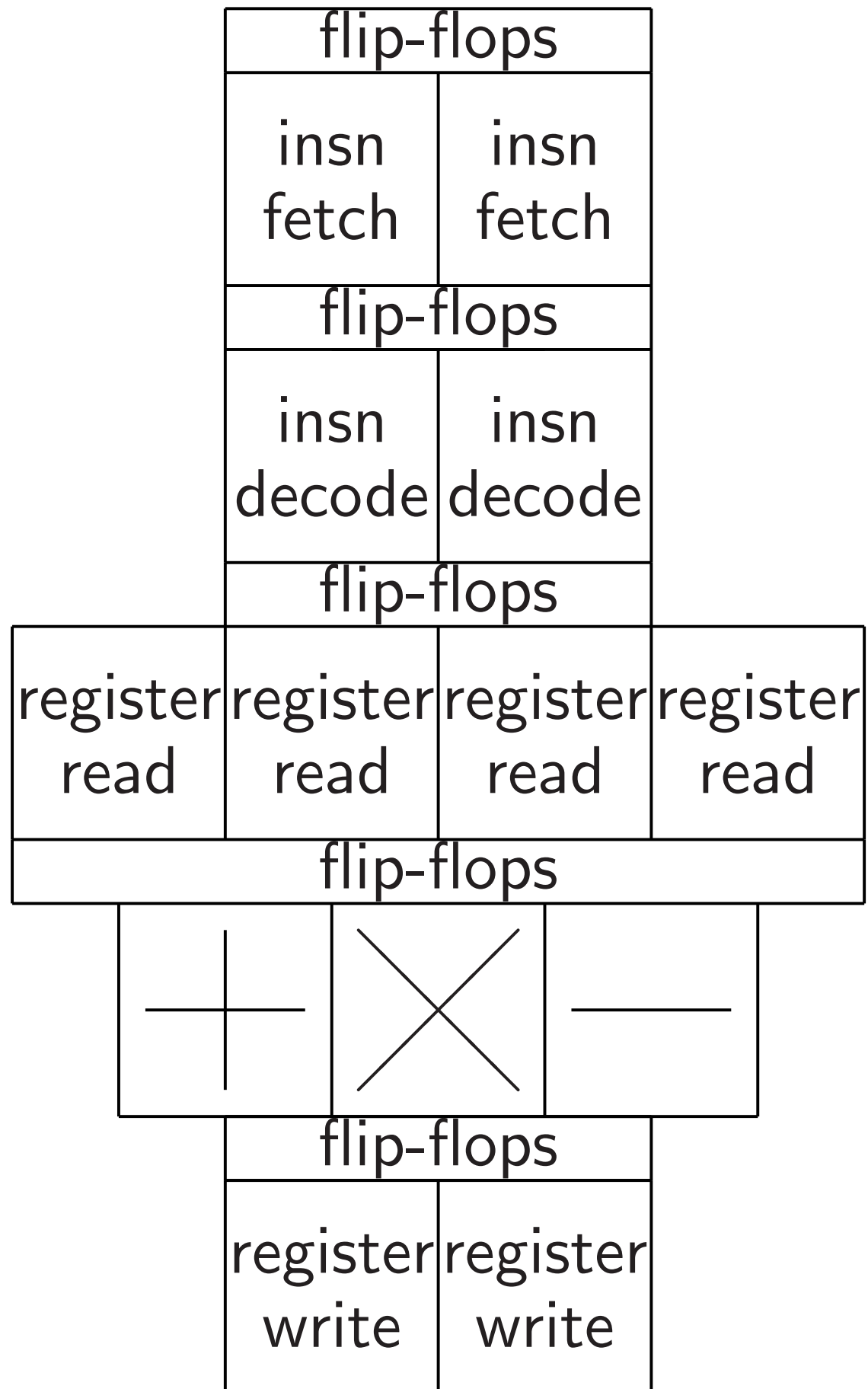
$n \times$ speedup if

$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

“Superscalar” processing:



“Vector” processing:

Expand each 32-bit integer into n -vector of 32-bit integers.
 ARM “NEON” has $n = 4$;
 Intel “AVX2” has $n = 8$;
 Intel “AVX-512” has $n = 16$;
 GPUs have larger n .

$n \times$ speedup if

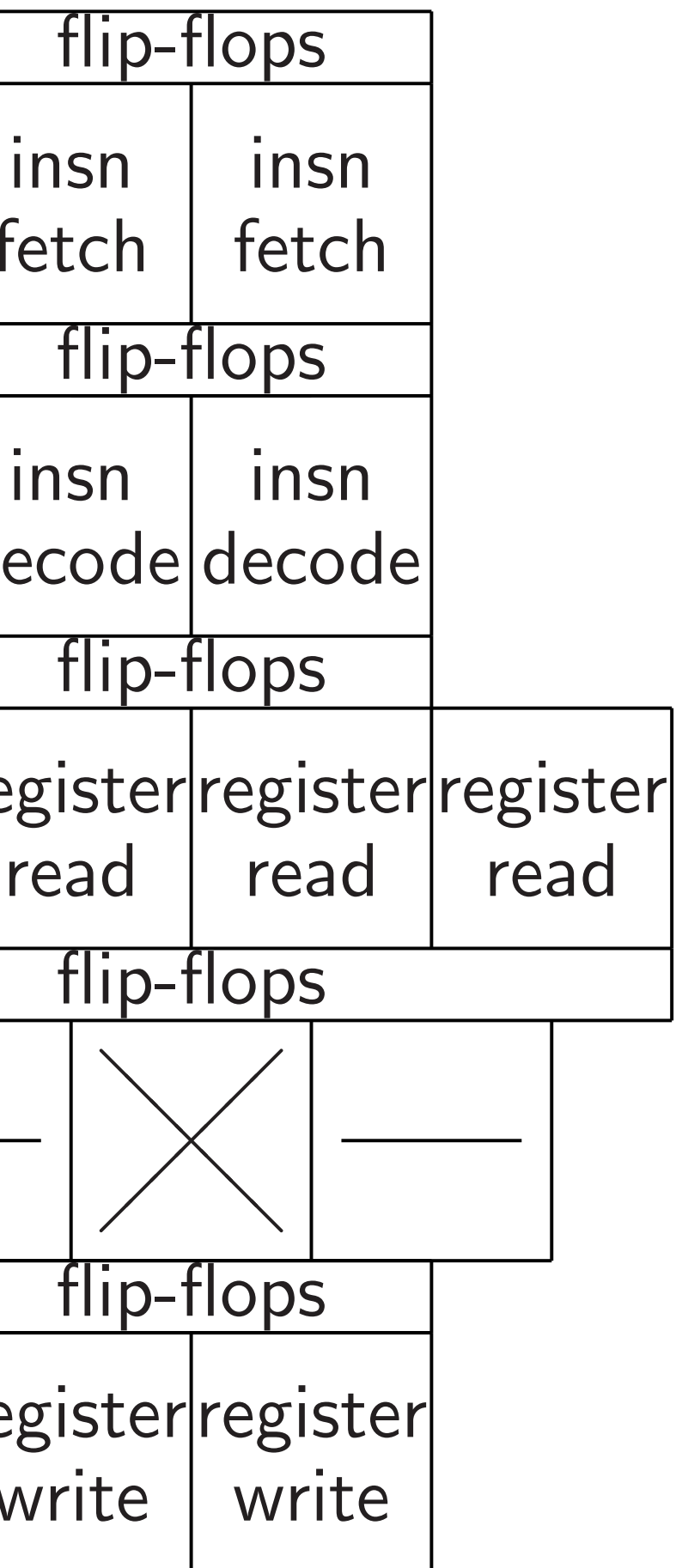
$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

Huge effect on higher-level algorithms and data structures.

“Scalar” processing:



“Vector” processing:

Expand each 32-bit integer into n -vector of 32-bit integers.

ARM “NEON” has $n = 4$;

Intel “AVX2” has $n = 8$;

Intel “AVX-512” has $n = 16$;

GPUs have larger n .

$n \times$ speedup if

$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

Huge effect on higher-level algorithms and data structures.

Network

How exp

Input: a

Each nu

represen

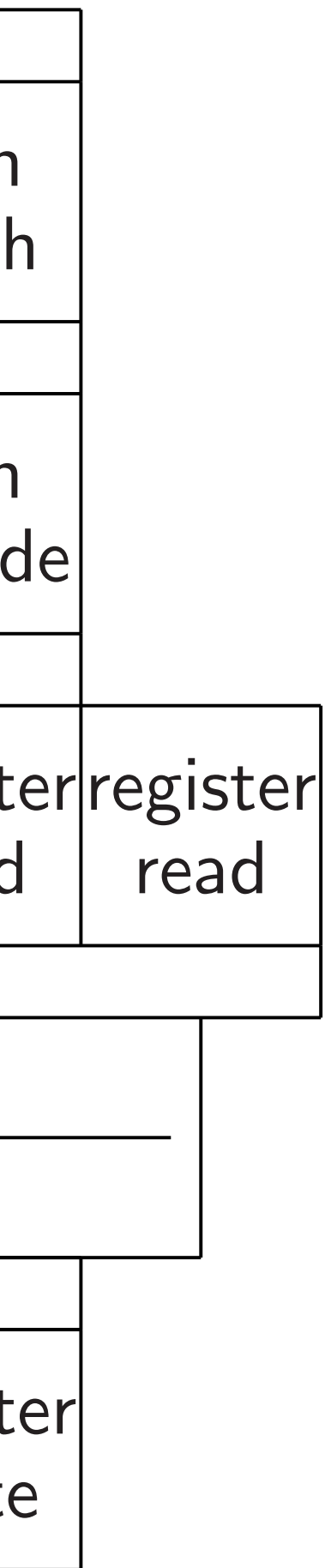
Output:

in increa

represen

same mu

processing:



“Vector” processing:

Expand each 32-bit integer into n -vector of 32-bit integers.

ARM “NEON” has $n = 4$;

Intel “AVX2” has $n = 8$;

Intel “AVX-512” has $n = 16$;

GPUs have larger n .

$n \times$ speedup if

$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

Huge effect on higher-level algorithms and data structures.

Network on chip:

How expensive is s

Input: array of n r

Each number in $\{$
represented in bina

Output: array of m

in increasing order

represented in bina

same multiset as i

“Vector” processing:

Expand each 32-bit integer
into n -vector of 32-bit integers.

ARM “NEON” has $n = 4$;

Intel “AVX2” has $n = 8$;

Intel “AVX-512” has $n = 16$;

GPUs have larger n .

$n \times$ speedup if

$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

Huge effect on higher-level
algorithms and data structures.

Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

Each number in $\{1, 2, \dots, n\}$
represented in binary.

Output: array of n numbers
in increasing order,

represented in binary;
same multiset as input.

“Vector” processing:

Expand each 32-bit integer
into n -vector of 32-bit integers.

ARM “NEON” has $n = 4$;

Intel “AVX2” has $n = 8$;

Intel “AVX-512” has $n = 16$;

GPUs have larger n .

$n \times$ speedup if

$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

Huge effect on higher-level
algorithms and data structures.

Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

Each number in $\{1, 2, \dots, n^2\}$,
represented in binary.

Output: array of n numbers,
in increasing order,
represented in binary;
same multiset as input.

“Vector” processing:

Expand each 32-bit integer
into n -vector of 32-bit integers.

ARM “NEON” has $n = 4$;

Intel “AVX2” has $n = 8$;

Intel “AVX-512” has $n = 16$;

GPUs have larger n .

$n \times$ speedup if

$n \times$ arithmetic circuits,

$n \times$ read/write circuits.

Benefit: Amortizes insn circuits.

Huge effect on higher-level
algorithms and data structures.

Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

Each number in $\{1, 2, \dots, n^2\}$,
represented in binary.

Output: array of n numbers,
in increasing order,
represented in binary;
same multiset as input.

Metric: seconds used by
circuit of area $n^{1+o(1)}$.

For simplicity assume $n = 4^k$.

' processing:

each 32-bit integer
ector of 32-bit integers.

"NEON" has $n = 4$;

"VX2" has $n = 8$;

"VX-512" has $n = 16$;

ave larger n .

dup if

arithmetic circuits,

/write circuits.

Amortizes insn circuits.

fect on higher-level

ms and data structures.

Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

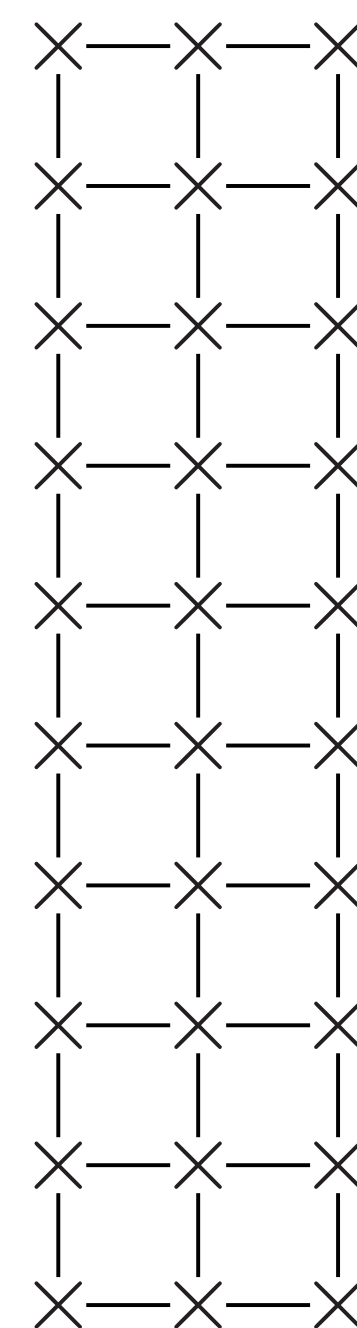
Each number in $\{1, 2, \dots, n^2\}$,
represented in binary.

Output: array of n numbers,
in increasing order,
represented in binary;
same multiset as input.

Metric: seconds used by
circuit of area $n^{1+o(1)}$.

For simplicity assume $n = 4^k$.

Spread a
square m
each of
with nea



Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

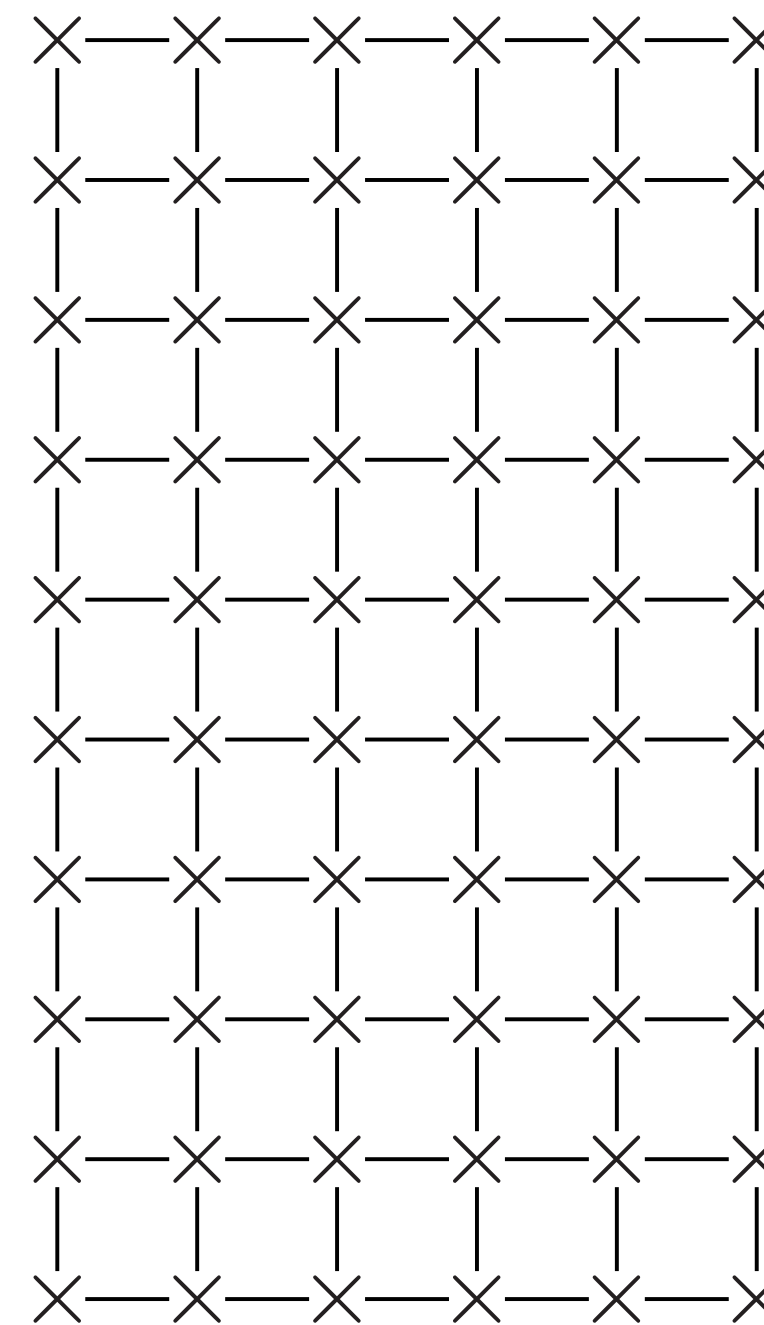
Each number in $\{1, 2, \dots, n^2\}$,
represented in binary.

Output: array of n numbers,
in increasing order,
represented in binary;
same multiset as input.

Metric: seconds used by
circuit of area $n^{1+o(1)}$.

For simplicity assume $n = 4^k$.

Spread array across
square mesh of n sites,
each of area $n^{o(1)}$
with near-neighborhood



Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

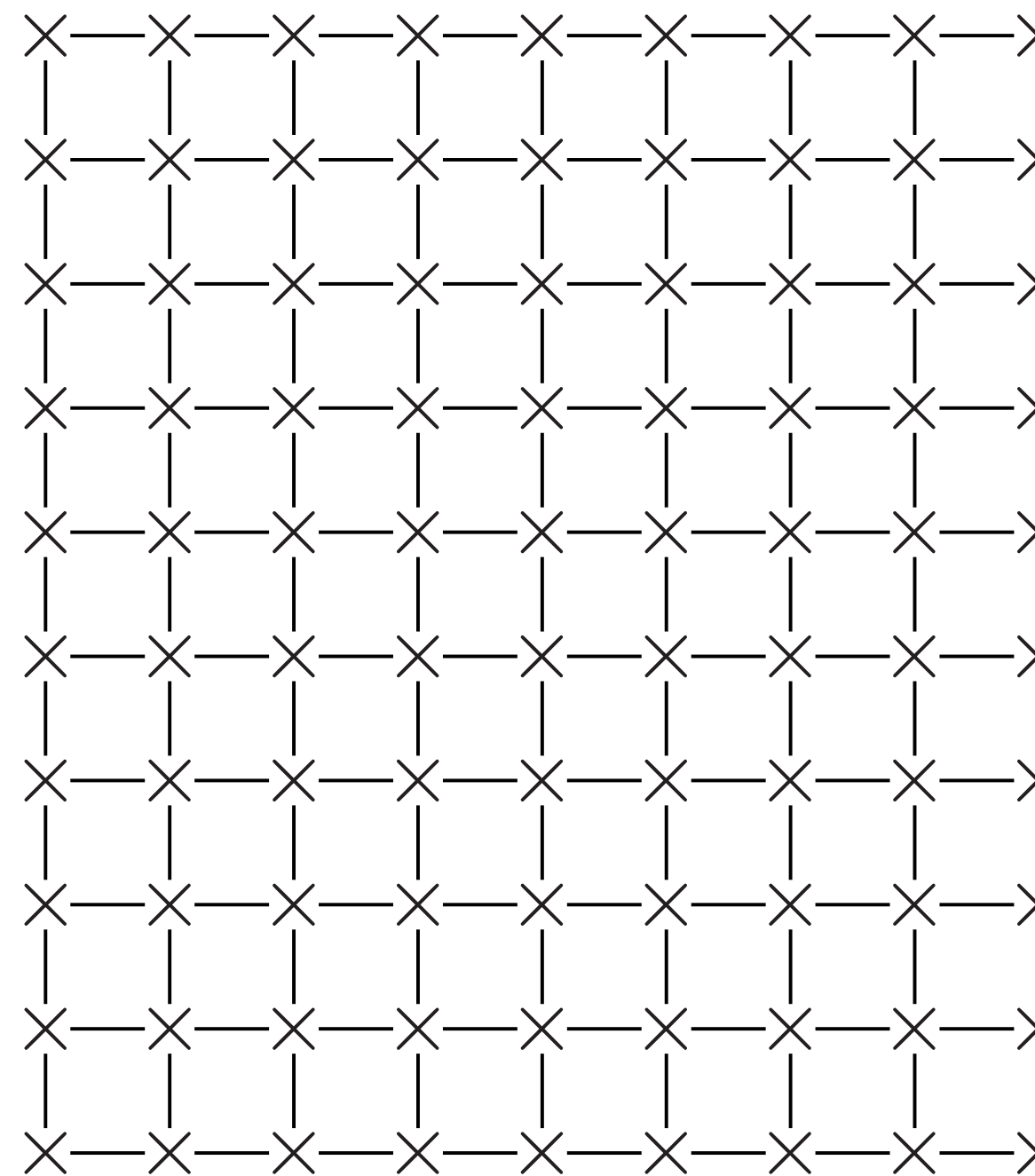
Each number in $\{1, 2, \dots, n^2\}$,
represented in binary.

Output: array of n numbers,
in increasing order,
represented in binary;
same multiset as input.

Metric: seconds used by
circuit of area $n^{1+o(1)}$.

For simplicity assume $n = 4^k$.

Spread array across
square mesh of n small cells
each of area $n^{o(1)}$,
with near-neighbor wiring:



Network on chip: the mesh

How expensive is sorting?

Input: array of n numbers.

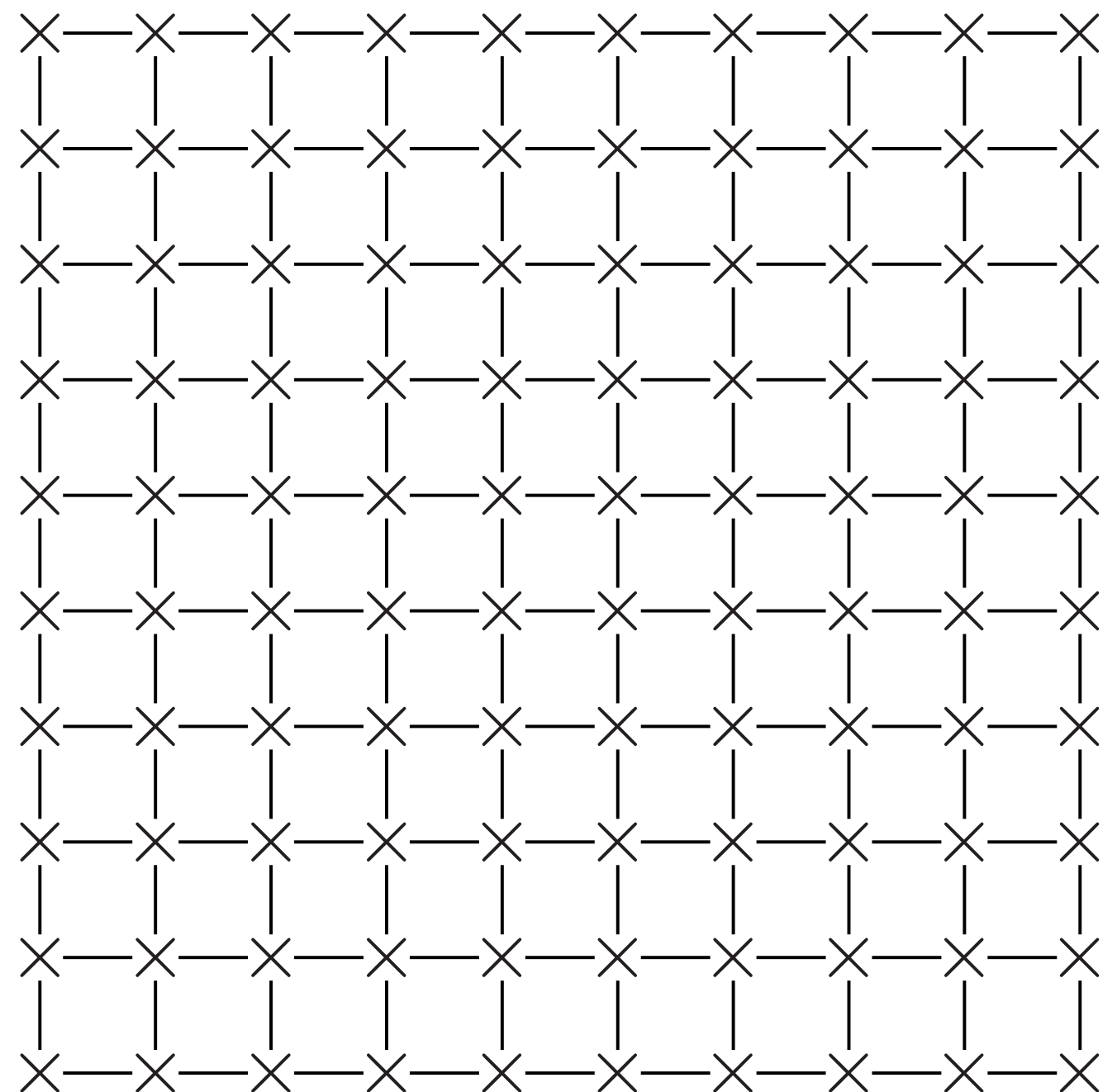
Each number in $\{1, 2, \dots, n^2\}$,
represented in binary.

Output: array of n numbers,
in increasing order,
represented in binary;
same multiset as input.

Metric: seconds used by
circuit of area $n^{1+o(1)}$.

For simplicity assume $n = 4^k$.

Spread array across
square mesh of n small cells,
each of area $n^{o(1)}$,
with near-neighbor wiring:



on chip: the mesh

ensive is sorting?

array of n numbers.

number in $\{1, 2, \dots, n^2\}$,

ted in binary.

array of n numbers,

asing order,

ted in binary;

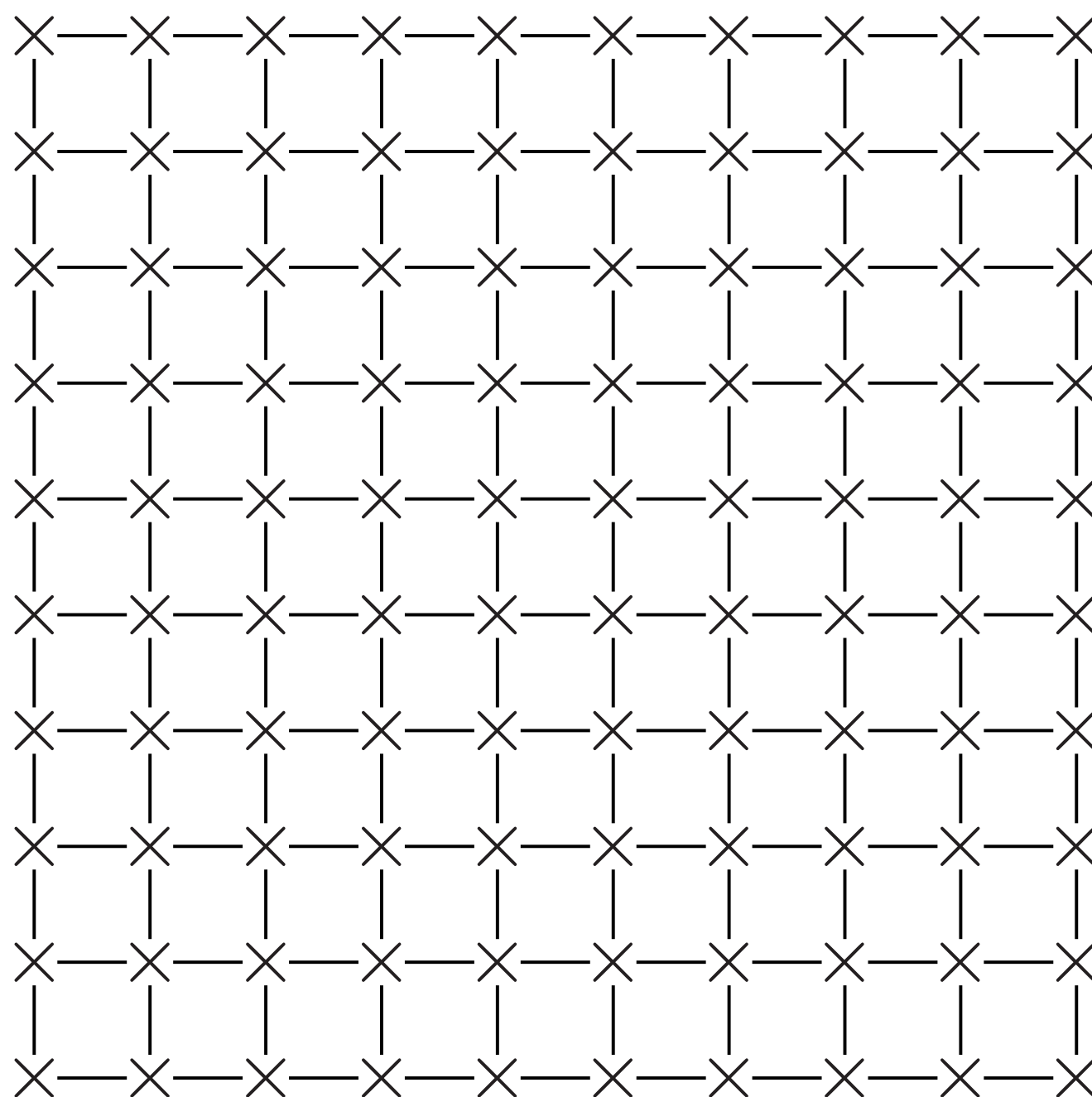
multiset as input.

seconds used by

f area $n^{1+o(1)}$.

licity assume $n = 4^k$.

Spread array across
square mesh of n small cells,
each of area $n^{o(1)}$,
with near-neighbor wiring:



Sort row
in $n^{0.5+o(1)}$

- Sort e

3 1 4

1 3 1 4

- Sort a

1 3 1 4

1 1 3 4

- Repea

equals

the mesh

sorting?

numbers.

$\{1, 2, \dots, n^2\}$,

ary.

n numbers,

,

ary;

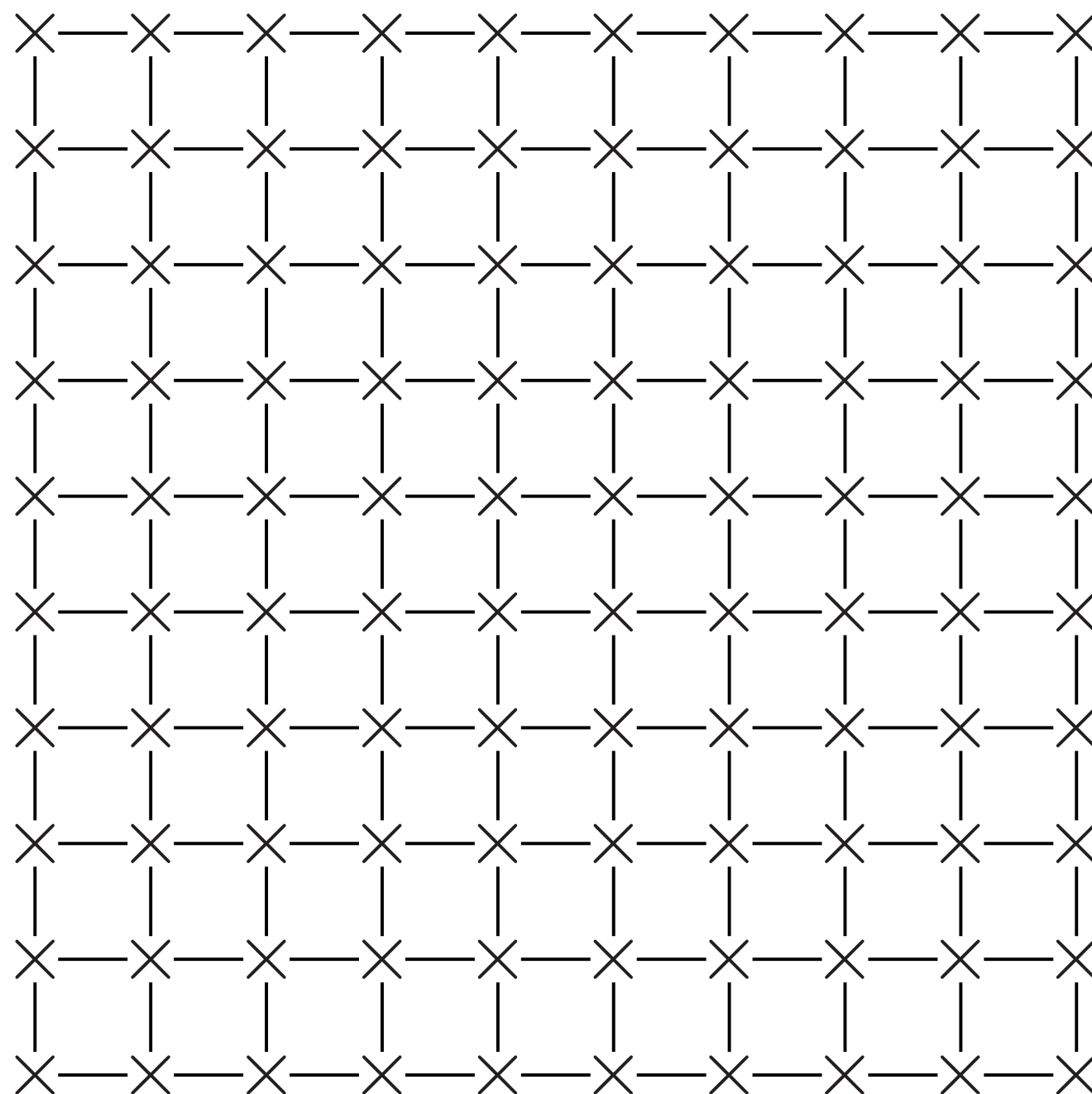
input.

sed by

$o(1)$.

me $n = 4^k$.

Spread array across
square mesh of n small cells,
each of area $n^{o(1)}$,
with near-neighbor wiring:



Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds

- Sort each pair in

3 1 4 1 5 9 2 6

1 3 1 4 5 9 2 6

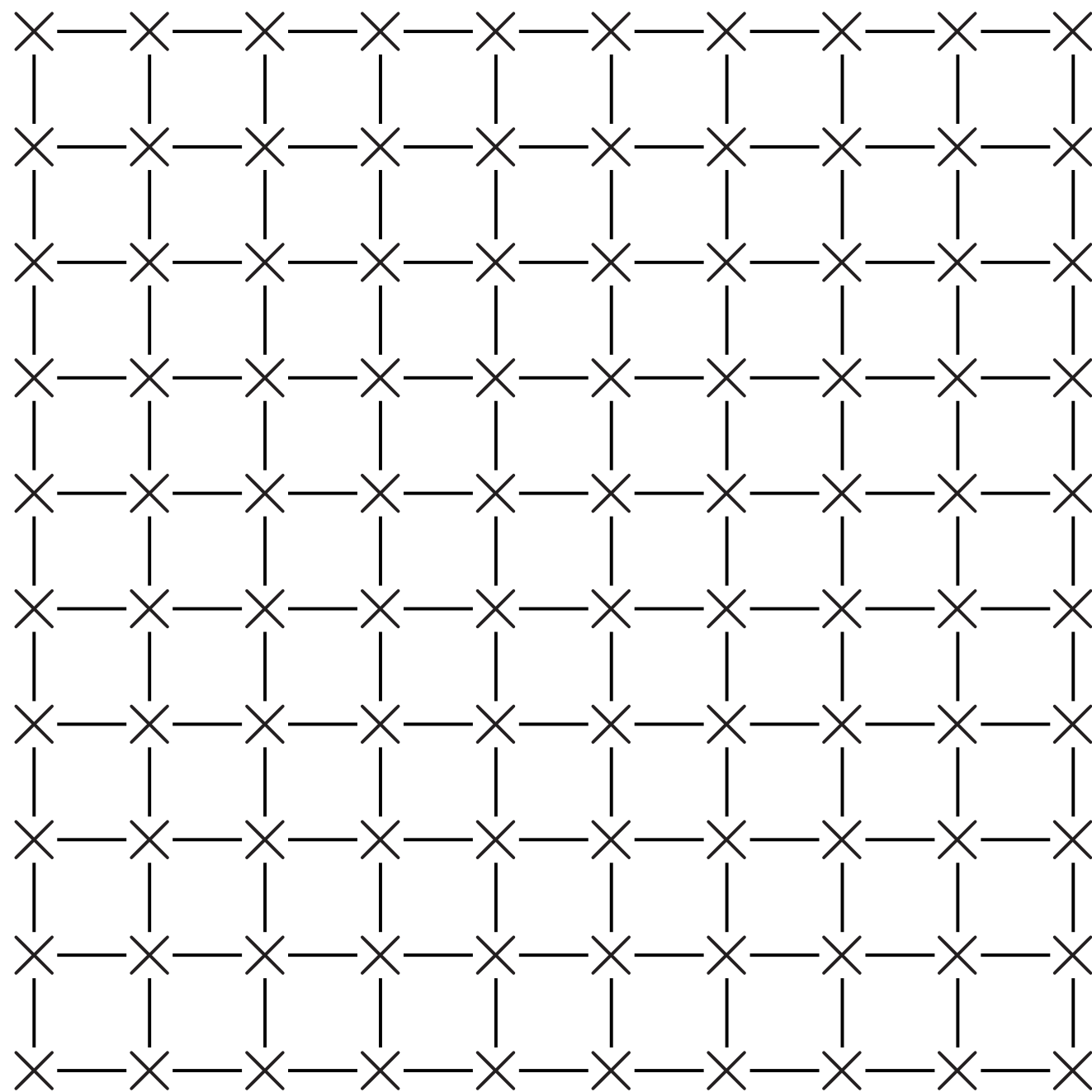
- Sort alternate pairs

1 3 1 4 5 9 2 6

1 1 3 4 5 2 9 6

- Repeat until number of rows equals row length

Spread array across
square mesh of n small cells,
each of area $n^{o(1)}$,
with near-neighbor wiring:



Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

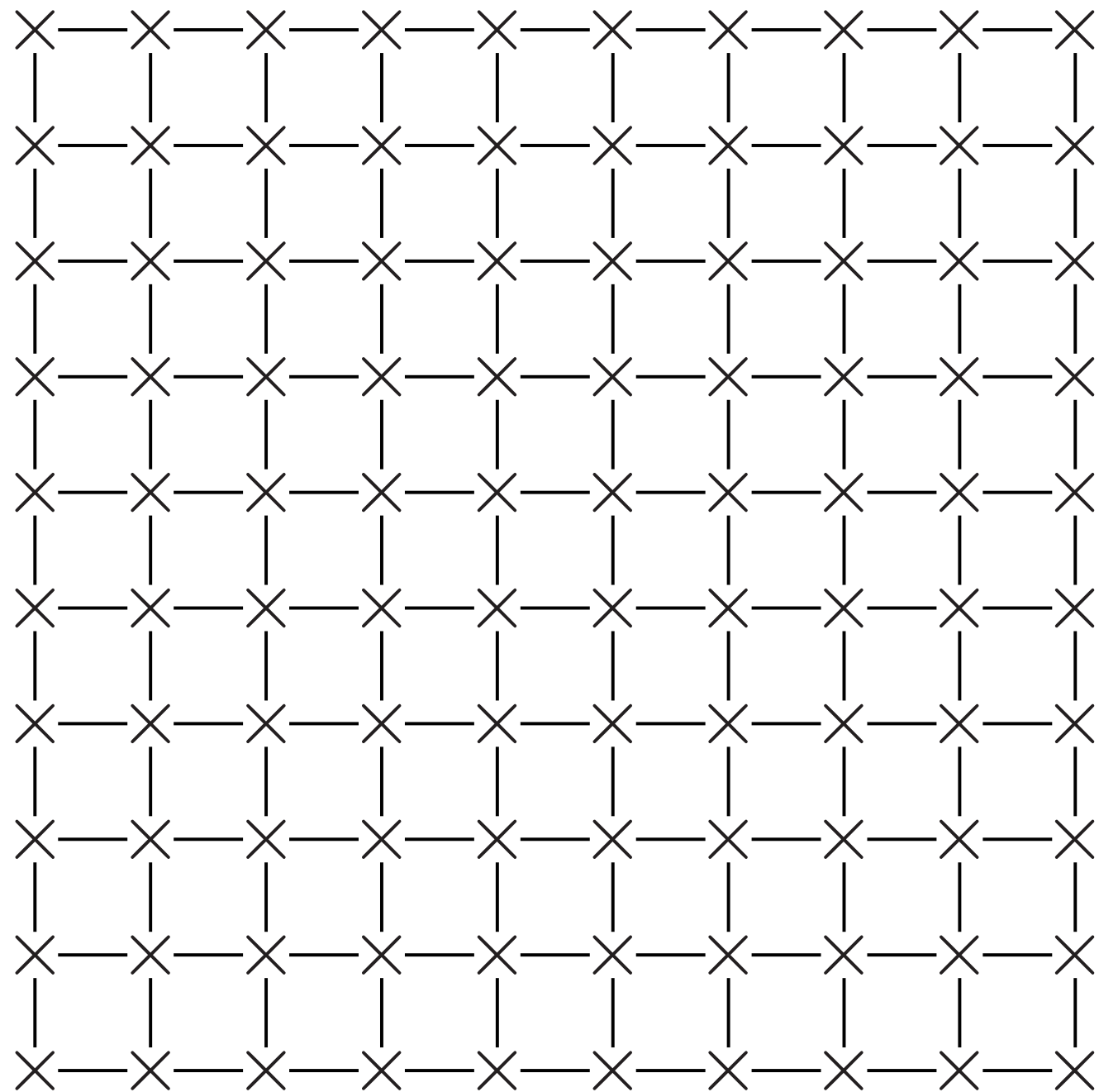
- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

- Repeat until number of steps equals row length.

Spread array across
square mesh of n small cells,
each of area $n^{o(1)}$,
with near-neighbor wiring:



Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

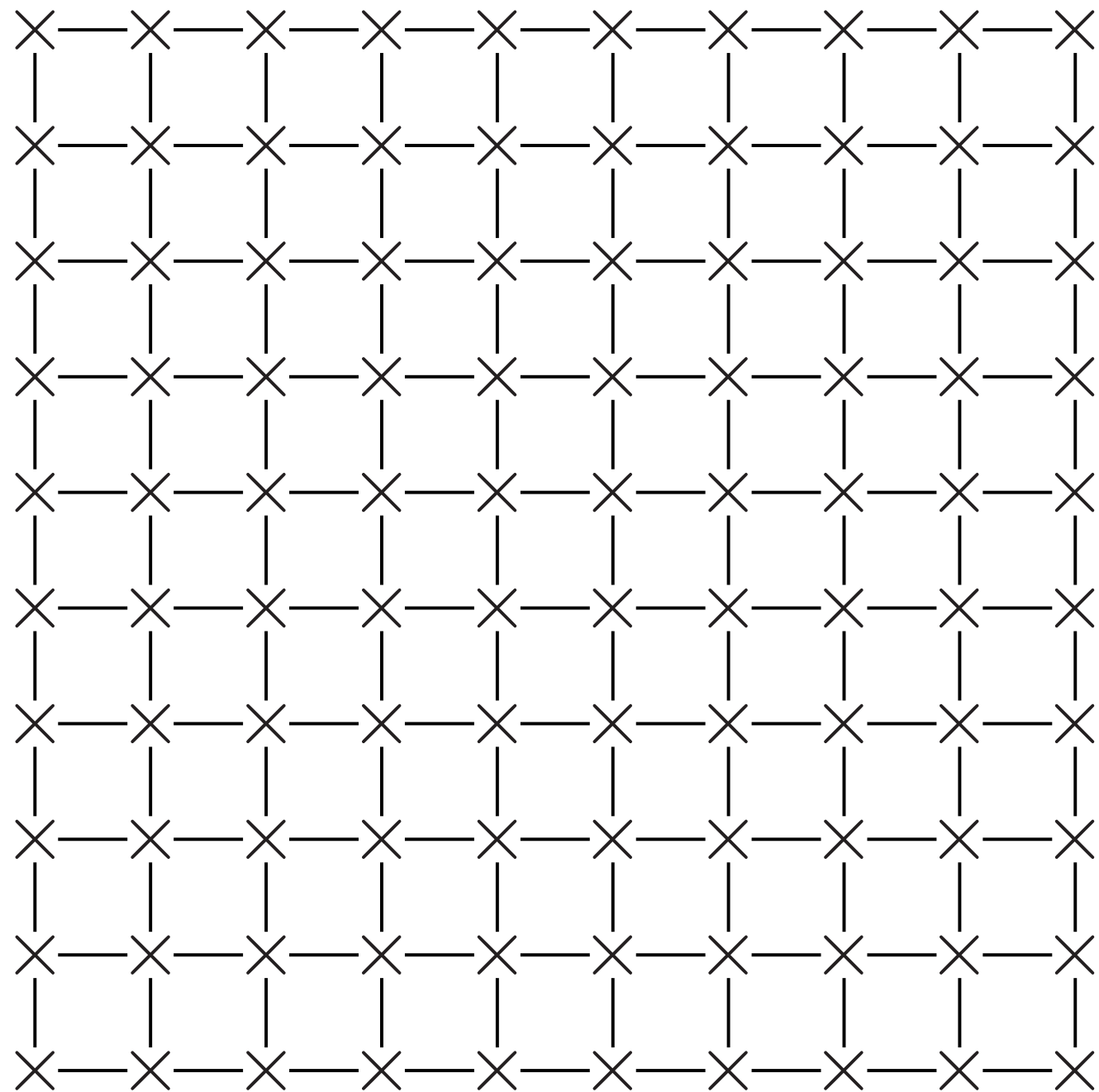
- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

- Repeat until number of steps equals row length.

Spread array across
square mesh of n small cells,
each of area $n^{o(1)}$,
with near-neighbor wiring:



Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

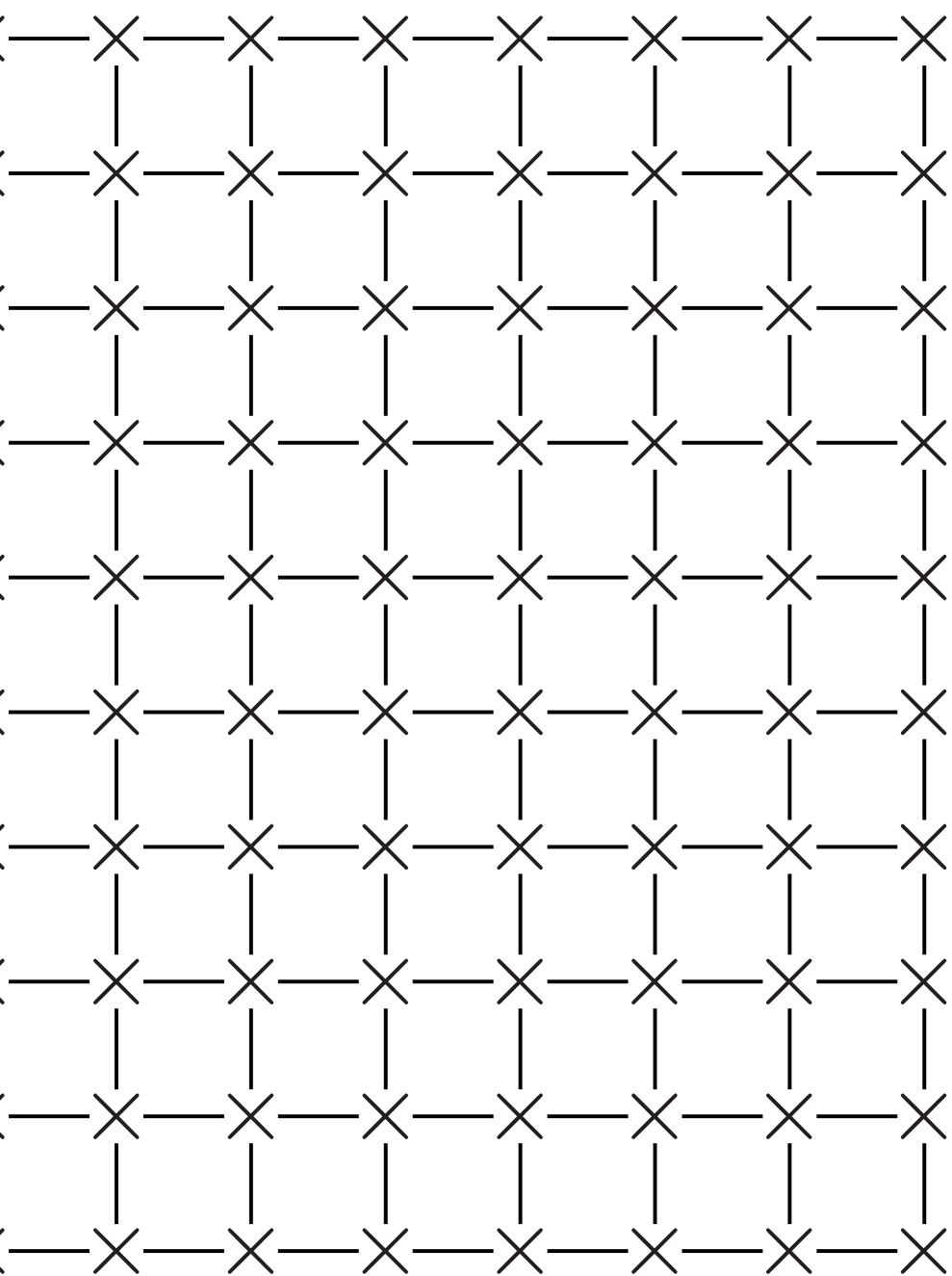
- Repeat until number of steps equals row length.

Sort *each* row, in parallel,
in a *total* of $n^{0.5+o(1)}$ seconds.

array across

mesh of n small cells,
area $n^{o(1)}$,

near-neighbor wiring:



Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

- Repeat until number of steps equals row length.

Sort *each* row, in parallel,
in a *total* of $n^{0.5+o(1)}$ seconds.

Sort all
in $n^{0.5+o(1)}$

- Recurs
in para

- Sort e

- Sort e

- Sort e

- Sort e

With pro

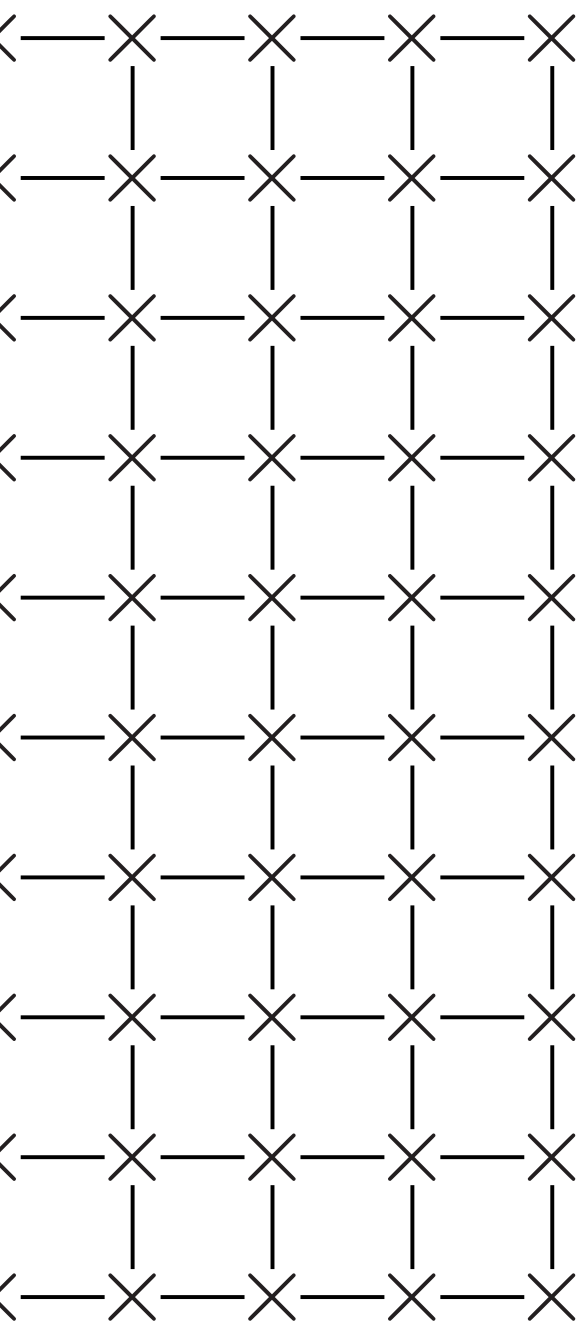
left-to-ri

for each

that this

ss
small cells,

wiring:



Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

- Repeat until number of steps equals row length.

Sort *each* row, in parallel,
in a *total* of $n^{0.5+o(1)}$ seconds.

Sort all n cells
in $n^{0.5+o(1)}$ seconds

- Recursively sort $n^{0.5}$ cells in parallel, if $n > n^{0.5}$
- Sort each column in parallel
- Sort each row in parallel
- Sort each column in parallel
- Sort each row in parallel

With proper choice of wiring,
left-to-right/right-to-left,
for each row, can be done
that this sorts whole array

Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

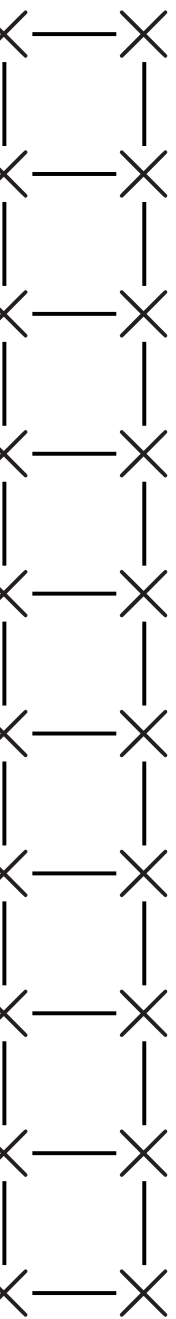
- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

- Repeat until number of steps equals row length.

Sort *each* row, in parallel,
in a *total* of $n^{0.5+o(1)}$ seconds.



Sort all n cells
in $n^{0.5+o(1)}$ seconds:

- Recursively sort quadrants in parallel, if $n > 1$.
- Sort each column in parallel.
- Sort each row in parallel.
- Sort each column in parallel.
- Sort each row in parallel.

With proper choice of left-to-right/right-to-left for each row, can prove that this sorts whole array.

Sort row of $n^{0.5}$ cells
in $n^{0.5+o(1)}$ seconds:

- Sort each pair in parallel.

3 1 4 1 5 9 2 6 \mapsto

1 3 1 4 5 9 2 6

- Sort alternate pairs in parallel.

1 3 1 4 5 9 2 6 \mapsto

1 1 3 4 5 2 9 6

- Repeat until number of steps equals row length.

Sort *each* row, in parallel,
in a *total* of $n^{0.5+o(1)}$ seconds.

Sort all n cells
in $n^{0.5+o(1)}$ seconds:

- Recursively sort quadrants in parallel, if $n > 1$.
- Sort each column in parallel.
- Sort each row in parallel.
- Sort each column in parallel.
- Sort each row in parallel.

With proper choice of left-to-right/right-to-left for each row, can prove that this sorts whole array.

of $n^{0.5}$ cells
 in $\Theta(n^{0.5})$ seconds:

Sort each pair in parallel.

1 5 9 2 6 \mapsto

4 5 9 2 6

Sort alternate pairs in parallel.

4 5 9 2 6 \mapsto

4 5 2 9 6

Continue until number of steps
 is proportional to row length.

Sort each row, in parallel,

in total of $n^{0.5+o(1)}$ seconds.

Sort all n cells
 in $n^{0.5+o(1)}$ seconds:

- Recursively sort quadrants in parallel, if $n > 1$.
- Sort each column in parallel.
- Sort each row in parallel.
- Sort each column in parallel.
- Sort each row in parallel.

With proper choice of
 left-to-right/right-to-left
 for each row, can prove
 that this sorts whole array.

For example,
 this 8×3 array

3	1	4
5	3	5
2	3	8
3	3	8
0	2	8
1	6	9
5	1	0
7	4	9

Sort all n cells
in $n^{0.5+o(1)}$ seconds:

- Recursively sort quadrants in parallel, if $n > 1$.
- Sort each column in parallel.
- Sort each row in parallel.
- Sort each column in parallel.
- Sort each row in parallel.

With proper choice of left-to-right/right-to-left for each row, can prove that this sorts whole array.

For example, assume
this 8×8 array is

3	1	4	1	5	9		
5	3	5	8	9	7		
2	3	8	4	6	2		
3	3	8	3	2	7		
0	2	8	8	4	1		
1	6	9	3	9	9		
5	1	0	5	8	2		
7	4	9	4	4	5		

Sort all n cells
in $n^{0.5+o(1)}$ seconds:

- Recursively sort quadrants in parallel, if $n > 1$.
- Sort each column in parallel.
- Sort each row in parallel.
- Sort each column in parallel.
- Sort each row in parallel.

With proper choice of left-to-right/right-to-left for each row, can prove that this sorts whole array.

For example, assume that this 8×8 array is in cells:

3	1	4	1	5	9	2	6
5	3	5	8	9	7	9	3
2	3	8	4	6	2	6	4
3	3	8	3	2	7	9	5
0	2	8	8	4	1	9	7
1	6	9	3	9	9	3	7
5	1	0	5	8	2	0	9
7	4	9	4	4	5	9	2

Sort all n cells
in $n^{0.5+o(1)}$ seconds:

- Recursively sort quadrants in parallel, if $n > 1$.
- Sort each column in parallel.
- Sort each row in parallel.
- Sort each column in parallel.
- Sort each row in parallel.

With proper choice of left-to-right/right-to-left for each row, can prove that this sorts whole array.

For example, assume that this 8×8 array is in cells:

3	1	4	1	5	9	2	6
5	3	5	8	9	7	9	3
2	3	8	4	6	2	6	4
3	3	8	3	2	7	9	5
0	2	8	8	4	1	9	7
1	6	9	3	9	9	3	7
5	1	0	5	8	2	0	9
7	4	9	4	4	5	9	2

n cells
 $p(1)$ seconds:

recursively sort quadrants

in parallel, if $n > 1$.

sort each column in parallel.

sort each row in parallel.

sort each column in parallel.

sort each row in parallel.

proper choice of

left/right/right-to-left

row, can prove

algorithm sorts whole array.

For example, assume that
 this 8×8 array is in cells:

3	1	4	1	5	9	2	6
5	3	5	8	9	7	9	3
2	3	8	4	6	2	6	4
3	3	8	3	2	7	9	5
0	2	8	8	4	1	9	7
1	6	9	3	9	9	3	7
5	1	0	5	8	2	0	9
7	4	9	4	4	5	9	2

Recursive

top \rightarrow , l

1	1	2
3	3	3
3	4	4
5	8	8
1	1	0
4	4	3
7	6	5
9	9	8

For example, assume that this 8×8 array is in cells:

3	1	4	1	5	9	2	6
5	3	5	8	9	7	9	3
2	3	8	4	6	2	6	4
3	3	8	3	2	7	9	5
0	2	8	8	4	1	9	7
1	6	9	3	9	9	3	7
5	1	0	5	8	2	0	9
7	4	9	4	4	5	9	2

Recursively sort qu
top \rightarrow , bottom \leftarrow

1	1	2	3	2	2
3	3	3	3	4	5
3	4	4	5	6	6
5	8	8	8	9	9
1	1	0	0	2	2
4	4	3	2	5	4
7	6	5	5	9	8
9	9	8	8	9	9

For example, assume that this 8×8 array is in cells:

3	1	4	1	5	9	2	6
5	3	5	8	9	7	9	3
2	3	8	4	6	2	6	4
3	3	8	3	2	7	9	5
0	2	8	8	4	1	9	7
1	6	9	3	9	9	3	7
5	1	0	5	8	2	0	9
7	4	9	4	4	5	9	2

el.

el.

Recursively sort quadrants, top \rightarrow , bottom \leftarrow :

1	1	2	3	2	2	2	3
3	3	3	3	4	5	5	6
3	4	4	5	6	6	7	7
5	8	8	8	9	9	9	9
1	1	0	0	2	2	1	0
4	4	3	2	5	4	4	3
7	6	5	5	9	8	7	7
9	9	8	8	9	9	9	9

For example, assume that this 8×8 array is in cells:

3	1	4	1	5	9	2	6
5	3	5	8	9	7	9	3
2	3	8	4	6	2	6	4
3	3	8	3	2	7	9	5
0	2	8	8	4	1	9	7
1	6	9	3	9	9	3	7
5	1	0	5	8	2	0	9
7	4	9	4	4	5	9	2

Recursively sort quadrants,
top \rightarrow , bottom \leftarrow :

1	1	2	3	2	2	2	3
3	3	3	3	4	5	5	6
3	4	4	5	6	6	7	7
5	8	8	8	9	9	9	9
1	1	0	0	2	2	1	0
4	4	3	2	5	4	4	3
7	6	5	5	9	8	7	7
9	9	8	8	9	9	9	9

Example, assume that
8 array is in cells:

1	5	9	2	6
8	9	7	9	3
4	6	2	6	4
3	2	7	9	5
8	4	1	9	7
3	9	9	3	7
5	8	2	0	9
4	4	5	9	2

Recursively sort quadrants,
top \rightarrow , bottom \leftarrow :

1	1	2	3	2	2	2	3
3	3	3	3	4	5	5	6
3	4	4	5	6	6	7	7
5	8	8	8	9	9	9	9
1	1	0	0	2	2	1	0
4	4	3	2	5	4	4	3
7	6	5	5	9	8	7	7
9	9	8	8	9	9	9	9

Sort each
in parallel

1	1	0
1	1	2
3	3	3
3	4	3
4	4	4
5	6	5
7	8	8
9	9	8

me that
in cells:

2	6
9	3
6	4
9	5
9	7
3	7
0	9
9	2

Recursively sort quadrants,
top \rightarrow , bottom \leftarrow :

1	1	2	3	2	2	2	3
3	3	3	3	4	5	5	6
3	4	4	5	6	6	7	7
5	8	8	8	9	9	9	9
1	1	0	0	2	2	1	0
4	4	3	2	5	4	4	3
7	6	5	5	9	8	7	7
9	9	8	8	9	9	9	9

Sort each column
in parallel:

1	1	0	0	2	2
1	1	2	2	2	2
3	3	3	3	4	4
3	4	3	3	5	5
4	4	4	5	6	6
5	6	5	5	9	8
7	8	8	8	9	9
9	9	8	8	9	9

Recursively sort quadrants,
top \rightarrow , bottom \leftarrow :

1	1	2	3	2	2	2	3
3	3	3	3	4	5	5	6
3	4	4	5	6	6	7	7
5	8	8	8	9	9	9	9
1	1	0	0	2	2	1	0
4	4	3	2	5	4	4	3
7	6	5	5	9	8	7	7
9	9	8	8	9	9	9	9

Sort each column
in parallel:

1	1	0	0	2	2	1	0
1	1	2	2	2	2	2	3
3	3	3	3	4	4	4	3
3	4	3	3	5	5	5	6
4	4	4	5	6	6	7	7
5	6	5	5	9	8	7	7
7	8	8	8	9	9	9	9
9	9	8	8	9	9	9	9

Recursively sort quadrants,
top \rightarrow , bottom \leftarrow :

1	1	2	3	2	2	2	3
3	3	3	3	4	5	5	6
3	4	4	5	6	6	7	7
5	8	8	8	9	9	9	9
1	1	0	0	2	2	1	0
4	4	3	2	5	4	4	3
7	6	5	5	9	8	7	7
9	9	8	8	9	9	9	9

Sort each column
in parallel:

1	1	0	0	2	2	1	0
1	1	2	2	2	2	2	3
3	3	3	3	4	4	4	3
3	4	3	3	5	5	5	6
4	4	4	5	6	6	7	7
5	6	5	5	9	8	7	7
7	8	8	8	9	9	9	9
9	9	8	8	9	9	9	9

ely sort quadrants,
bottom \leftarrow :

2	3	2	2	2	3
3	3	4	5	5	6
4	5	6	6	7	7
8	8	9	9	9	9
0	0	2	2	1	0
3	2	5	4	4	3
5	5	9	8	7	7
3	8	9	9	9	9

Sort each column
in parallel:

1	1	0	0	2	2	1	0
1	1	2	2	2	2	2	3
3	3	3	3	4	4	4	3
3	4	3	3	5	5	5	6
4	4	4	5	6	6	7	7
5	6	5	5	9	8	7	7
7	8	8	8	9	9	9	9
9	9	8	8	9	9	9	9

Sort each
alternate

0	0	0
3	2	2
3	3	3
6	5	5
4	4	4
9	8	7
7	8	8
9	9	9

quadrants,
:

2	3
5	6
7	7
9	9
1	0
4	3
7	7
9	9

Sort each column
in parallel:

1	1	0	0	2	2	1	0
1	1	2	2	2	2	2	3
3	3	3	3	4	4	4	3
3	4	3	3	5	5	5	6
4	4	4	5	6	6	7	7
5	6	5	5	9	8	7	7
7	8	8	8	9	9	9	9
9	9	8	8	9	9	9	9

Sort each row in parallel
alternately \leftarrow , \rightarrow :

0	0	0	1	1	1
3	2	2	2	2	2
3	3	3	3	3	4
6	5	5	5	4	3
4	4	4	5	6	6
9	8	7	7	6	5
7	8	8	8	9	9
9	9	9	9	9	9

h column
el:

0	0	2	2	1	0
2	2	2	2	2	3
3	3	4	4	4	3
3	3	5	5	5	6
4	5	6	6	7	7
5	5	9	8	7	7
3	8	9	9	9	9
3	8	9	9	9	9

39

Sort each row in parallel,
alternately \leftarrow , \rightarrow :

0	0	0	1	1	1	2	2
3	2	2	2	2	2	1	1
3	3	3	3	3	4	4	4
6	5	5	5	4	3	3	3
4	4	4	5	6	6	7	7
9	8	7	7	6	5	5	5
7	8	8	8	9	9	9	9
9	9	9	9	9	9	8	8

40

Sort each
in parallel

0	0	0
3	2	2
3	3	3
4	4	4
6	5	5
7	8	7
9	8	8
9	9	9

Sort each row in parallel,
alternately \leftarrow , \rightarrow :

1	0
2	3
4	3
5	6
7	7
7	7
9	9
9	9

0	0	0	1	1	1	2	2
3	2	2	2	2	2	1	1
3	3	3	3	3	4	4	4
6	5	5	5	4	3	3	3
4	4	4	5	6	6	7	7
9	8	7	7	6	5	5	5
7	8	8	8	9	9	9	9
9	9	9	9	9	9	8	8

Sort each column
in parallel:

0	0	0	1	1	1
3	2	2	2	2	2
3	3	3	3	3	3
4	4	4	5	4	4
6	5	5	5	6	5
7	8	7	7	6	6
9	8	8	8	9	9
9	9	9	9	9	9

Sort each row in parallel,
 sequentially \leftarrow , \rightarrow :

0	1	1	1	2	2
2	2	2	2	1	1
3	3	3	4	4	4
5	5	4	3	3	3
4	5	6	6	7	7
7	7	6	5	5	5
8	8	9	9	9	9
9	9	9	9	8	8

Sort each column
 in parallel:

0	0	0	1	1	1	1	1
3	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3
4	4	4	5	4	4	4	4
6	5	5	5	6	5	5	5
7	8	7	7	6	6	7	7
9	8	8	8	9	9	8	8
9	9	9	9	9	9	9	9

Sort each
 \leftarrow or \rightarrow

0	0	0
2	2	2
3	3	3
4	4	4
5	5	5
6	6	7
8	8	8
9	9	9

parallel,

2	2
1	1
4	4
3	3
7	7
5	5
9	9
8	8

Sort each column
in parallel:

0	0	0	1	1	1	1	1
3	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3
4	4	4	5	4	4	4	4
6	5	5	5	6	5	5	5
7	8	7	7	6	6	7	7
9	8	8	8	9	9	8	8
9	9	9	9	9	9	9	9

Sort each row in p

← or → as desired

0	0	0	1	1	1
2	2	2	2	2	2
3	3	3	3	3	3
4	4	4	4	4	4
5	5	5	5	5	5
6	6	7	7	7	7
8	8	8	8	8	9
9	9	9	9	9	9

h column
el:

0	1	1	1	1	1
2	2	2	2	2	2
3	3	3	3	3	3
4	5	4	4	4	4
5	5	6	5	5	5
7	7	6	6	7	7
8	8	9	9	8	8
9	9	9	9	9	9

Sort each row in parallel,
← or → as desired:

0	0	0	1	1	1	1	1
2	2	2	2	2	2	2	3
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	5
5	5	5	5	5	5	6	6
6	6	7	7	7	7	7	8
8	8	8	8	8	9	9	9
9	9	9	9	9	9	9	9

Chips ar
towards
parallelis
GPUs: p
Old Xeo
New Xeo

Sort each row in parallel,
 \leftarrow or \rightarrow as desired:

1	1
2	2
3	3
4	4
5	5
7	7
8	8
9	9

0	0	0	1	1	1	1	1
2	2	2	2	2	2	2	3
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	5
5	5	5	5	5	5	6	6
6	6	7	7	7	7	7	8
8	8	8	8	8	9	9	9
9	9	9	9	9	9	9	9

Chips are in fact e
 towards having thi
 parallelism and co

GPUs: parallel +

Old Xeon Phi: pa

New Xeon Phi: pa

Sort each row in parallel,
← or → as desired:

0	0	0	1	1	1	1	1
2	2	2	2	2	2	2	3
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	5
5	5	5	5	5	5	6	6
6	6	7	7	7	7	7	8
8	8	8	8	8	9	9	9
9	9	9	9	9	9	9	9

Chips are in fact evolving
towards having this much
parallelism and communicat

GPUs: parallel + global RA

Old Xeon Phi: parallel + rin

New Xeon Phi: parallel + m

Sort each row in parallel,
 \leftarrow or \rightarrow as desired:

0	0	0	1	1	1	1	1
2	2	2	2	2	2	2	3
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	5
5	5	5	5	5	5	6	6
6	6	7	7	7	7	7	8
8	8	8	8	8	9	9	9
9	9	9	9	9	9	9	9

Chips are in fact evolving
 towards having this much
 parallelism and communication.

GPUs: parallel + global RAM.

Old Xeon Phi: parallel + ring.

New Xeon Phi: parallel + mesh.

Sort each row in parallel,
 \leftarrow or \rightarrow as desired:

0	0	0	1	1	1	1	1
2	2	2	2	2	2	2	3
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	5
5	5	5	5	5	5	6	6
6	6	7	7	7	7	7	8
8	8	8	8	8	9	9	9
9	9	9	9	9	9	9	9

Chips are in fact evolving
 towards having this much
 parallelism and communication.

GPUs: parallel + global RAM.

Old Xeon Phi: parallel + ring.

New Xeon Phi: parallel + mesh.

Algorithm designers

don't even get the right exponent
 without taking this into account.

Sort each row in parallel,
 \leftarrow or \rightarrow as desired:

0	0	0	1	1	1	1	1
2	2	2	2	2	2	2	3
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	5
5	5	5	5	5	5	6	6
6	6	7	7	7	7	7	8
8	8	8	8	8	9	9	9
9	9	9	9	9	9	9	9

Chips are in fact evolving
 towards having this much
 parallelism and communication.

GPUs: parallel + global RAM.

Old Xeon Phi: parallel + ring.

New Xeon Phi: parallel + mesh.

Algorithm designers

don't even get the right exponent
 without taking this into account.

Shock waves from subroutines
 into high-level algorithm design.