

# COMPILER CONFIDENTIAL

ERIC BRUMER

WHEN YOU THINK “COMPILER”...

```
c:\work\a.cpp(82): error C2059: syntax error : ')'  
c:\work\a.cpp(84): error C2015: too many characters in constant  
c:\work\a.cpp(104): error C2015: too many characters in constant  
c:\work\a.cpp(104): error C2001: newline in constant  
c:\work\a.cpp(116): error C2015: too many characters in constant  
c:\work\a.cpp(116): error C2001: newline in constant  
c:\work\a.cpp(122): error C2153: hex constants must have at least one hex digit  
c:\work\a.cpp(122): error C2001: newline in constant  
c:\work\a.cpp(122): error C2015: too many characters in constant  
c:\work\a.cpp(134): error C2015: too many characters in constant  
c:\work\a.cpp(134): error C2001: newline in constant  
c:\work\a.cpp(140): error C2015: too many characters in constant  
c:\work\a.cpp(140): error C2001: newline in constant  
c:\work\a.cpp(146): error C2001: newline in constant  
c:\work\a.cpp(146): error C2015: too many characters in constant  
c:\work\a.cpp(154): error C2146: syntax error : missing ';' before identifier 'modern'  
c:\work\a.cpp(154): error C4430: missing type specifier - int assumed. Note: C++ does not support default-int  
c:\work\a.cpp(154): error C2143: syntax error : missing ';' before '-'  
c:\work\a.cpp(154): error C2015: too many characters in constant  
c:\work\a.cpp(155): error C2059: syntax error : 'constant'  
c:\work\a.cpp(155): error C2059: syntax error : 'bad suffix on number'  
c:\work\a.cpp(158): error C2015: too many characters in constant  
c:\work\a.cpp(158): error C2059: syntax error : ')'  
c:\work\a.cpp(161): error C2001: newline in constant  
c:\work\a.cpp(161): error C2015: too many characters in constant  
c:\work\a.cpp(164): error C2059: syntax error : 'bad suffix on number'  
c:\work\a.cpp(164): error C2059: syntax error : 'constant'  
c:\work\a.cpp(168): error C2001: newline in constant  
c:\work\a.cpp(168): error C2015: too many characters in constant  
c:\work\a.cpp(178): error C2146: syntax error : missing ';' before identifier 'Examples'  
c:\work\a.cpp(178): error C4430: missing type specifier - int assumed. Note: C++ does not support default-int  
c:\work\a.cpp(178): error C2146: syntax error : missing ';' before identifier 'in'  
c:\work\a.cpp(178): error C2146: syntax error : missing ';' before identifier 'C'  
c:\work\a.cpp(178): error C2143: syntax error : missing ';' before '++'  
c:\work\a.cpp(181): error C2146: syntax error : missing ';' before identifier 'Examples'  
c:\work\a.cpp(181): error C4430: missing type specifier - int assumed. Note: C++ does not support default-int
```

```
void test(bool b) {  
    try {  
        if (b) {  
            MyType obj;  
            some_func(obj);  
            // ...  
        }  
    } catch (...) {  
        // ...  
    }  
}
```

**Destructor  
placement**



# CODE GENERATION & OPTIMIZATION

```
int binary_search2(std::vector<int> &arr, int value, int lowIndex, int highIndex) {  
    int midIndex = (highIndex + lowIndex) / 2;  
    int mid = arr[midIndex];  
    if (value == mid)  
        return midIndex;  
    if (lowIndex >= highIndex)  
        return -1;  
    if (value > mid) {  
        int l = midIndex+1;  
        int h = highIndex;  
        return binary_search2(arr, value, l, h);  
    } else { // value < mid  
        int l = lowIndex;  
        int h = midIndex-1;  
        return binary_search2(arr, value, l, h);  
    }  
}  
  
int binary_search(std::vector<int> &arr, int value) {  
    int l = 0;  
    int h = arr.size() - 1;  
    return binary_search2(arr, value, l, h);  
}
```



MAKE MY CODE RUN: CODE GENERATION

MAKE MY CODE RUN FAST: **OPTIMIZATION**

# MISSION: EXPOSE SOME OPTIMIZER GUTS

THERE WILL BE RAW LOOPS

THERE WILL BE ASSEMBLY CODE

THERE WILL BE MICROARCHITECTURE



**I sense much  
fear in you**

# AGENDA

CPU HARDWARE LANDSCAPE

VECTORIZING FOR MODERN CPUs

INDIRECT CALL OPTIMIZATIONS

# AGENDA

CPU HARDWARE LANDSCAPE

VECTORIZING FOR MODERN CPUs

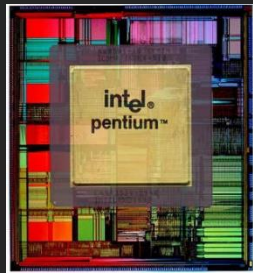
INDIRECT CALL OPTIMIZATIONS



# HARDWARE LANDSCAPE

“Yesterday”

3.1 million transistors

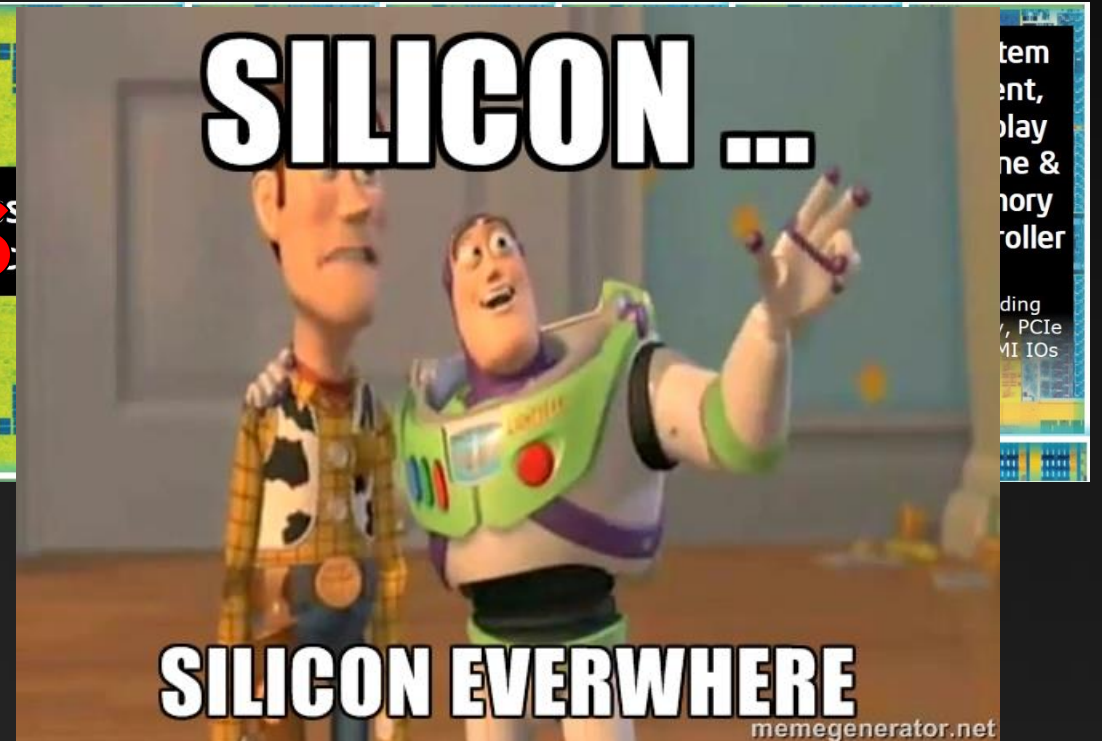


Not to s



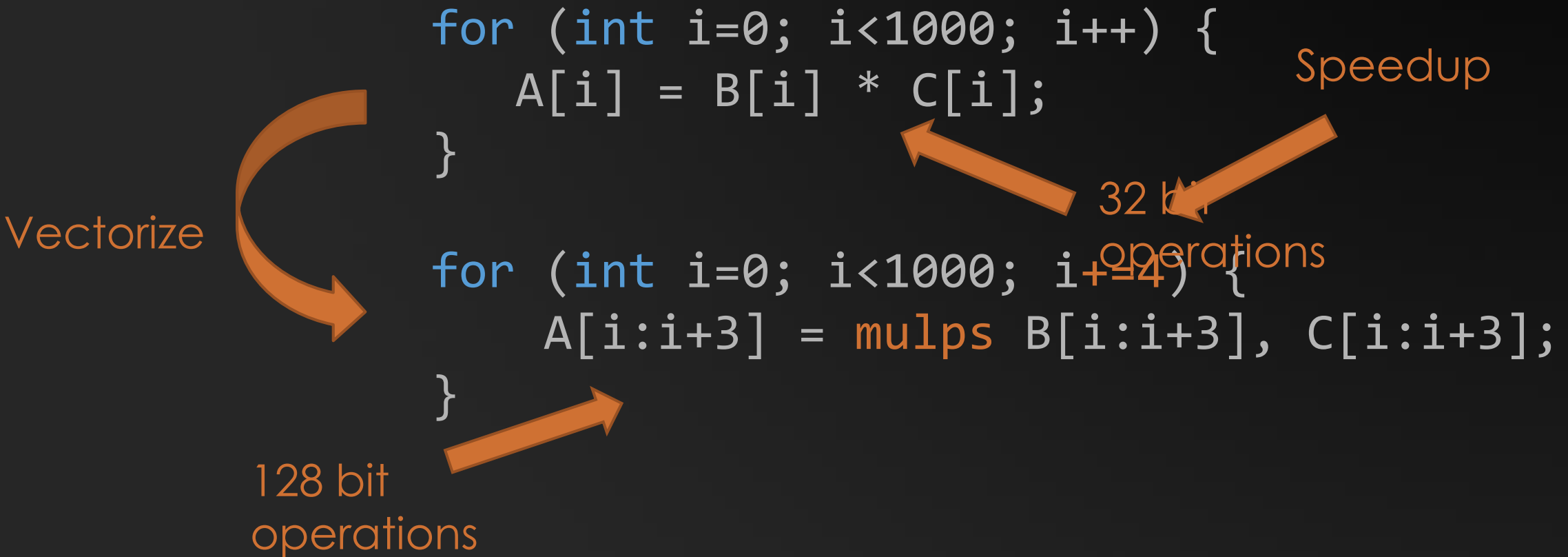
Today

1.4 billion transistors



# AUTOMATIC VECTORIZATION

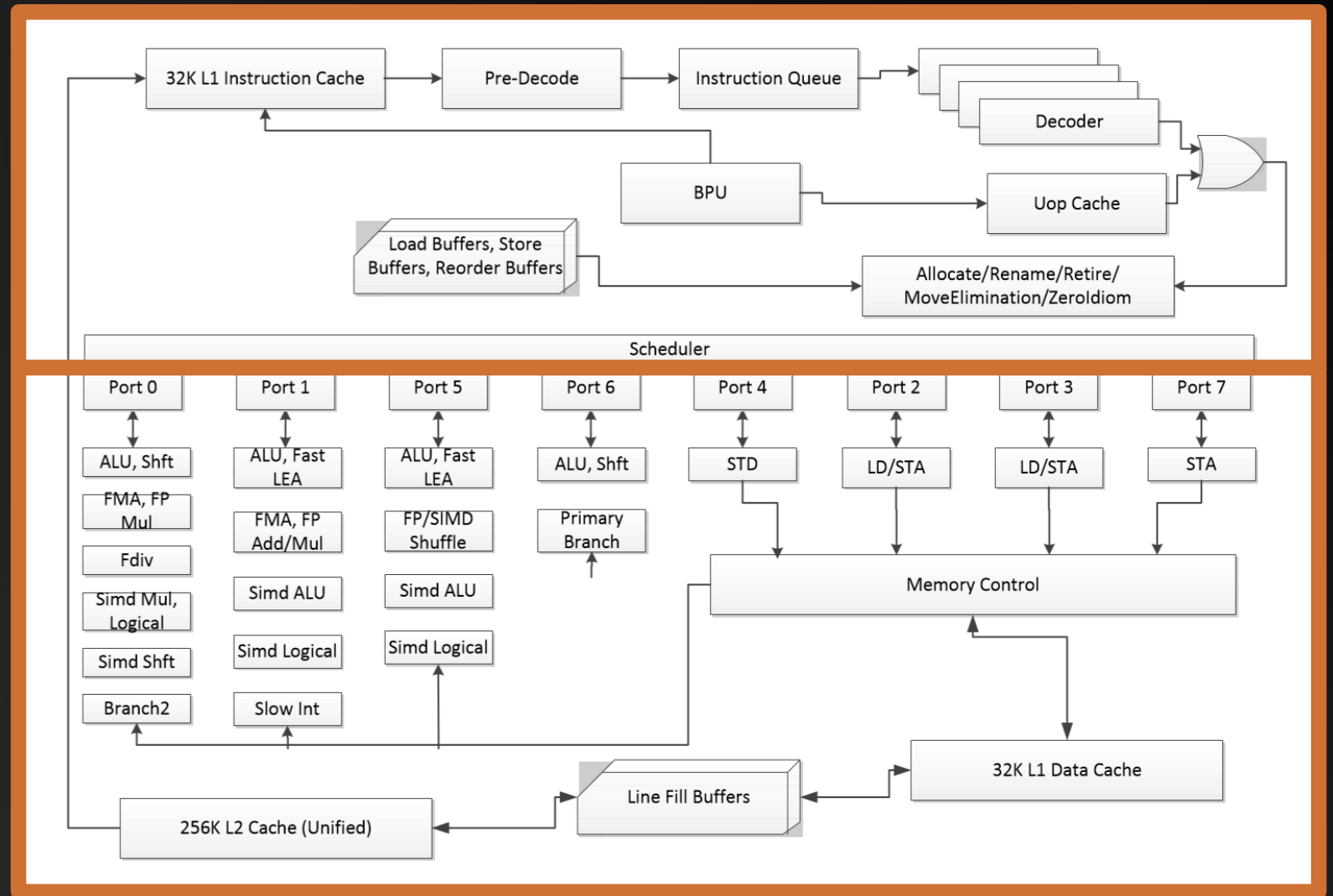
- TAKE ADVANTAGE OF (FAST) VECTOR HARDWARE
- EXECUTE MULTIPLE LOOP ITERATIONS IN PARALLEL



## Front-end

- Powerful branch predictor
- Ship instructions to backend as fast as possible

## Haswell core microarchitecture



## Back-end

- 8 wide super scalar
- Powerful vector units

# AGENDA

CPU HARDWARE LANDSCAPE

VECTORIZING FOR MODERN CPUs

INDIRECT CALL OPTIMIZATIONS

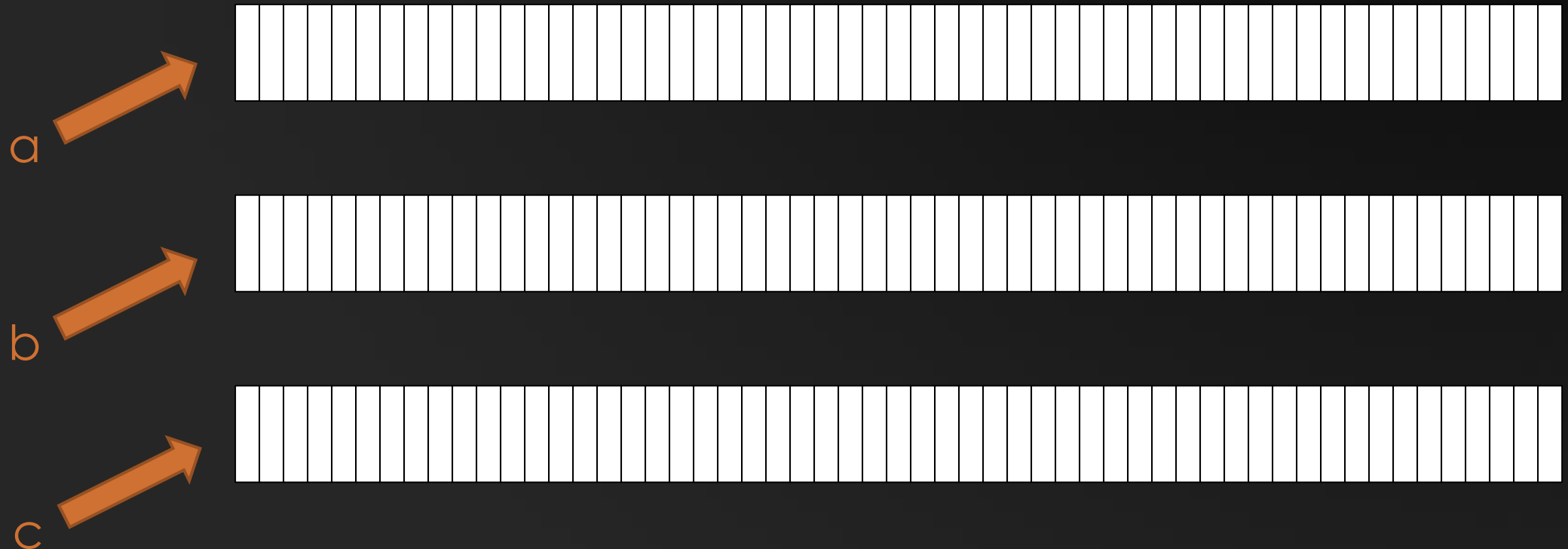
APPROACH TO VECTORIZING FOR MODERN CPUs:  
TAKE ADVANTAGE OF ALL THE EXTRA SILICON



KEY IDEA: CONDITIONAL VECTORIZATION

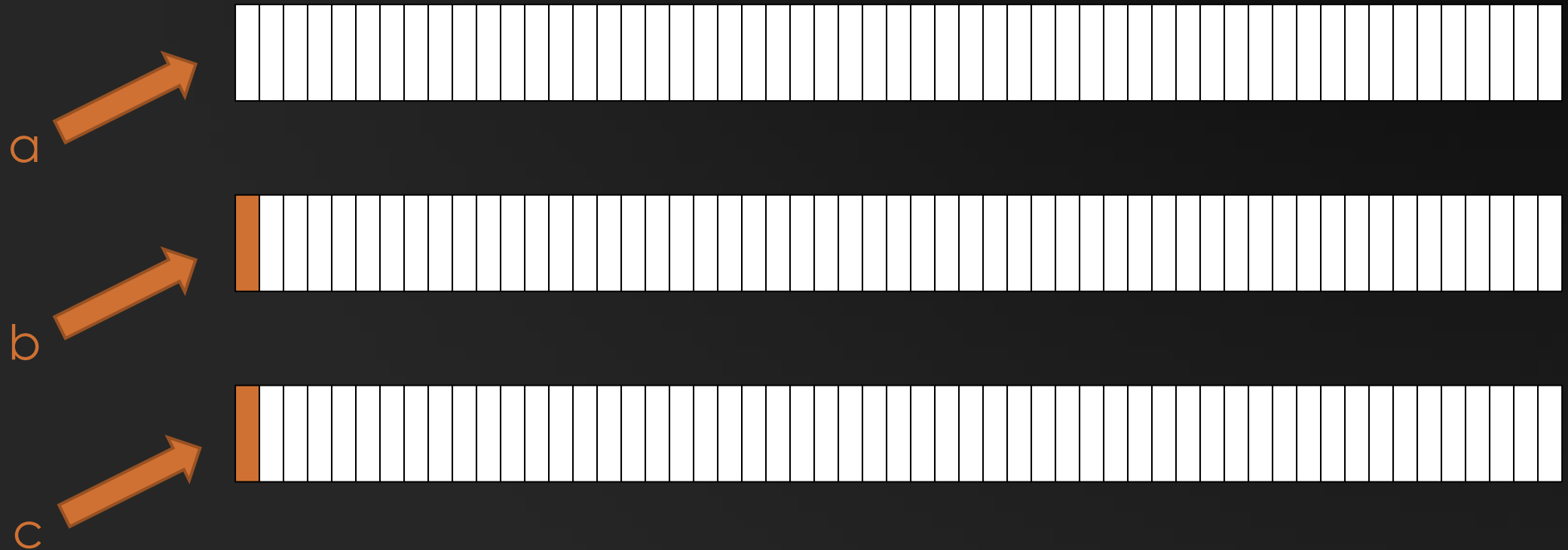
# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```



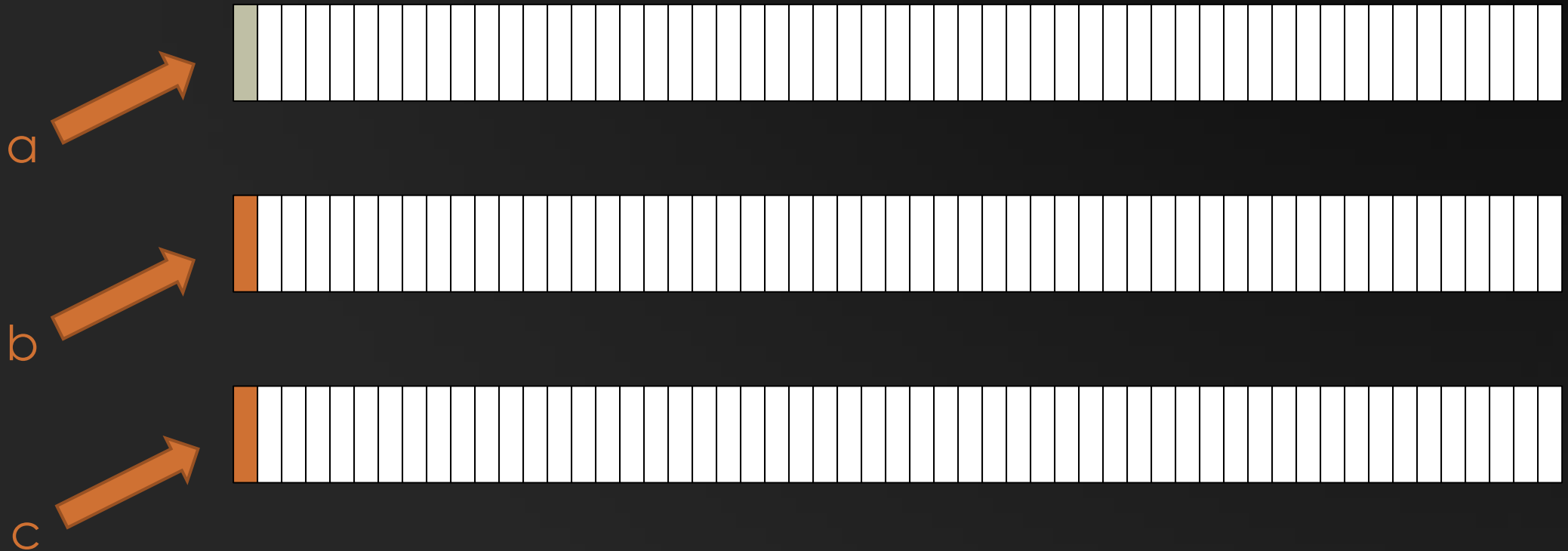
# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```



# MOTIVATING EXAMPLE

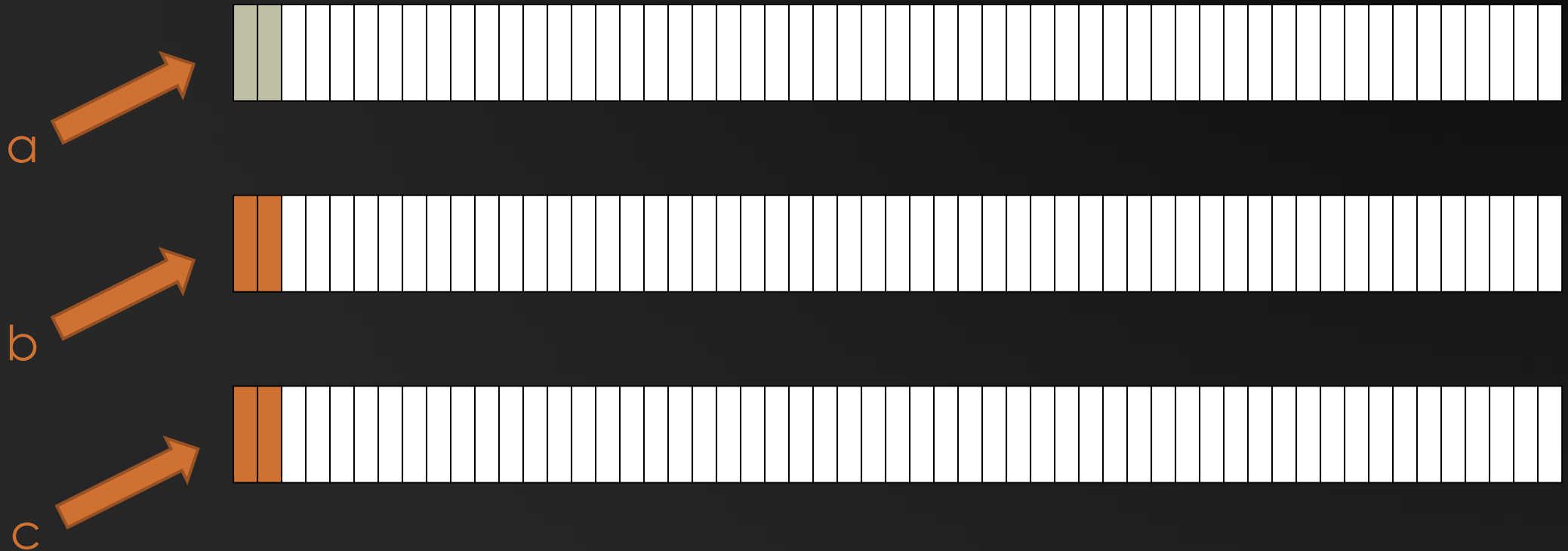
```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```





# MOTIVATING EXAMPLE

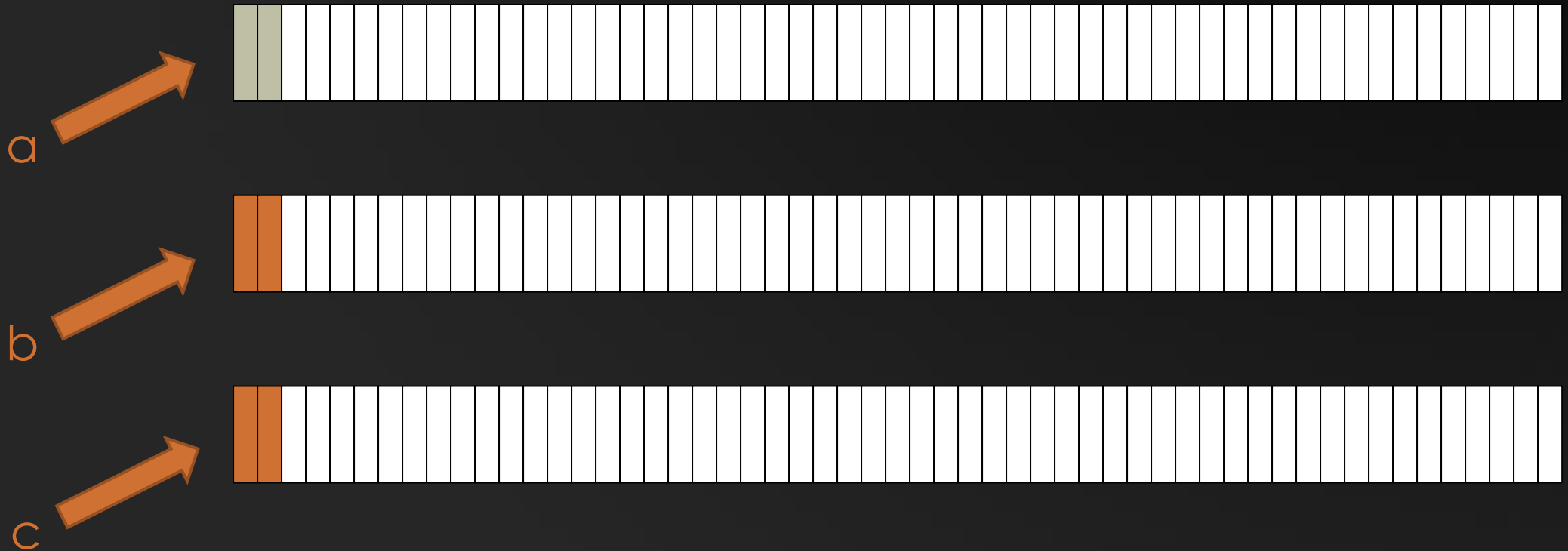
```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

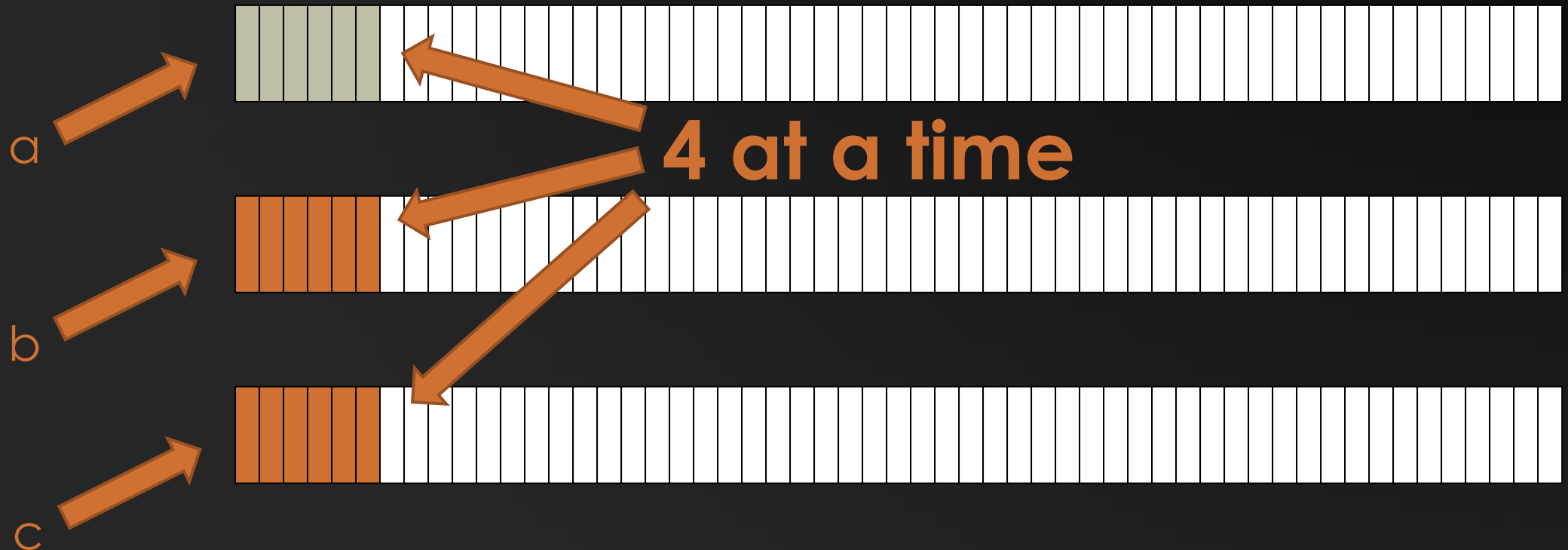
**Easy to vectorize**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

Easy to vectorize



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**

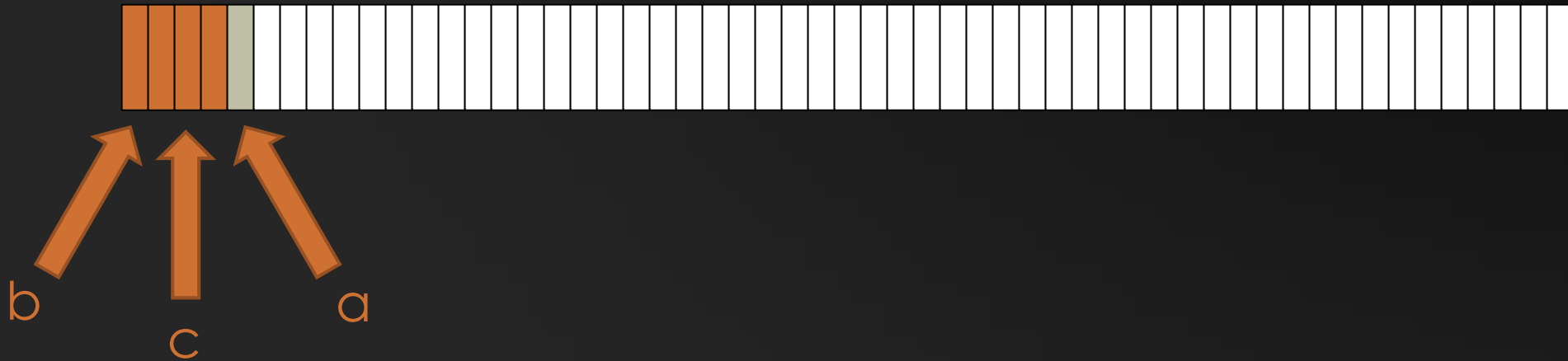




# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

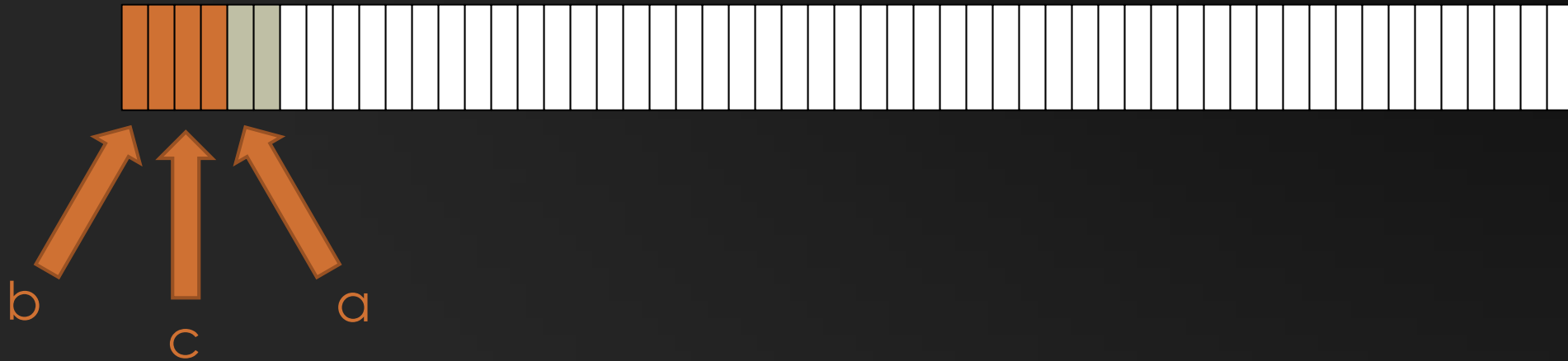
**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**What if there is  
overlap?**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

What if there is  
overlap?



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**Vectorization is not legal!**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**Vectorization is not legal!**





# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**Vectorization is not legal!**



# MOTIVATING EXAMPLE

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

**Vectorization is not legal!**



**WRONG! We are reading c[2] without having first stored to a[0]**

THE PRESENCE OF ~~OVERLAP PROHIBITS~~ VECTORIZATION

THE PRESENCE OF *POSSIBLE* OVERLAP PROHIBITS VECTORIZATION

THE COMPILER CAN STILL GENERATE FAST CODE

# CONDITIONAL VECTORIZATION #1

Source  
code:

```
void mul_flt(float *a, float *b, float *c) {  
    for (int i=0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

What we  
generate for  
you:

```
void mul_flt(float *a, float *b, float *c) {  
    if (a overlaps b) goto scalar_loop;  
    if (a overlaps c) goto scalar_loop;  
  
    for (int i = 0; i<1000; i+=4)  
        a[i:i+3] = mulps b[i:i+3], c[i:i+3];  
    return;  
scalar_loop:  
    for (int i = 0; i<1000; i++)  
        a[i] = b[i] * c[i];  
}
```

Runtime  
overlap  
checks



Vector  
loop



Scalar  
duplicate



# CONDITIONAL VECTORIZATION #1

```
for (int i=0; i<1000; i++)  
    a[i] = b[i] * c[i];
```

- 4 INSTRS OF RUNTIME CHECK, PLUS DUPLICATE LOOP
- `mul_float()` CODE SIZE INCREASES BY 7X

## 2.63X SPEEDUP

FOR REFERENCE, 2.64X SPEEDUP FOR VECT W/O RUNTIME CHECK AND THE DUPLICATE LOOP.

WHY?



# CONDITIONAL VECTORIZATION #2

Loop



```
for (k = 1; k <= M; k++) {
    mc[k] = mpp[k-1] + tpmm[k-1];
    if ((sc = ip[k-1] + tpim[k-1]) > mc[k]) mc[k] = sc;
    if ((sc = dpp[k-1] + tpdm[k-1]) > mc[k]) mc[k] = sc;
    if ((sc = xmb + bp[k]) > mc[k]) mc[k] = sc;
    mc[k] += ms[k];
    if (mc[k] < -INFTY) mc[k] = -INFTY;

    dc[k] = dc[k-1] + tpdd[k-1];
    if ((sc = mc[k-1] + tpmd[k-1]) > dc[k]) dc[k] = sc;
    if (dc[k] < -INFTY) dc[k] = -INFTY;

    if (k < M) {
        ic[k] = mpp[k] + tpmi[k];
        if ((sc = ip[k] + tpim[k]) > ic[k]) ic[k] = sc;
        ic[k] += is[k];
        if (ic[k] < -INFTY) ic[k] = -INFTY;
    }
}
```

# CONDITIONAL VECTORIZATION #2

```
for (k = 1; k <= M; k++) {
    mc[k] = mpp[k-1] + tpmm[k-1];
    if ((sc = ip[k-1] + tpim[k-1]) > mc[k]) mc[k] = sc;
    if ((sc = dpp[k-1] + tpdm[k-1]) > mc[k]) mc[k] = sc;
    if ((sc = xmb + bp[k]) > mc[k]) mc[k] = sc;
    mc[k] += ms[k];
    if (mc[k] < -INFTY) mc[k] = -INFTY;

    dc[k] = dc[k-1] + tpdd[k-1];
    if ((sc = mc[k-1] + tpmd[k-1]) > dc[k]) dc[k] = sc;
    if (dc[k] < -INFTY) dc[k] = -INFTY;

    if (k < M) {
        ic[k] = mpp[k] + tpmi[k];
        if ((sc = ip[k] + tpii[k]) > ic[k]) ic[k] = sc;
        ic[k] += is[k];
        if (ic[k] < -INFTY) ic[k] = -INFTY;
    }
}
```

# CONDITIONAL VECTORIZATION #2

- 42 RUNTIME CHECKS NEEDED
- 84 CMP/BR INSTRUCTIONS, DUPLICATE LOOP
- LOOP CODE SIZE INCREASES BY 4X

DOESN'T THIS SUCK?

**2X LOOP SPEEDUP**

**30% OVERALL BENCHMARK SPEEDUP**

FOR REFERENCE, 2.1X SPEEDUP FOR VECT W/O RUNTIME CHECK

```
for (k = 1; k <= M; k++) {
    mc[k] = mpp[k-1] + tpmm[k-1];
    if ((sc = ip[k-1] + tpim[k-1]) > mc[k]) mc[k] = sc;
    if ((sc = dpp[k-1] + tpdm[k-1]) > mc[k]) mc[k] = sc;
    if ((sc = xmb + bp[k]) > mc[k]) mc[k] = sc;
    mc[k] += ms[k];
    if (mc[k] < -INFTY) mc[k] = -INFTY;

    dc[k] = dc[k-1] + tpdd[k-1];
    if ((sc = mc[k-1] + tpmd[k-1]) > dc[k]) dc[k] = sc;
    if (dc[k] < -INFTY) dc[k] = -INFTY;

    if (k < M) {
        ic[k] = mpp[k] + tpmi[k];
        if ((sc = ip[k] + tpim[k]) > ic[k]) ic[k] = sc;
        ic[k] += is[k];
        if (ic[k] < -INFTY) ic[k] = -INFTY;
    }
}
```



# AGENDA

CPU HARDWARE LANDSCAPE

VECTORIZING FOR MODERN CPUs

INDIRECT CALL OPTIMIZATIONS

```
typedef int (PFUNC)(int);
```

```
int func1(int x) {  
    return x + 100;  
}
```

```
int func2(int x) {  
    return x + 200;  
}
```

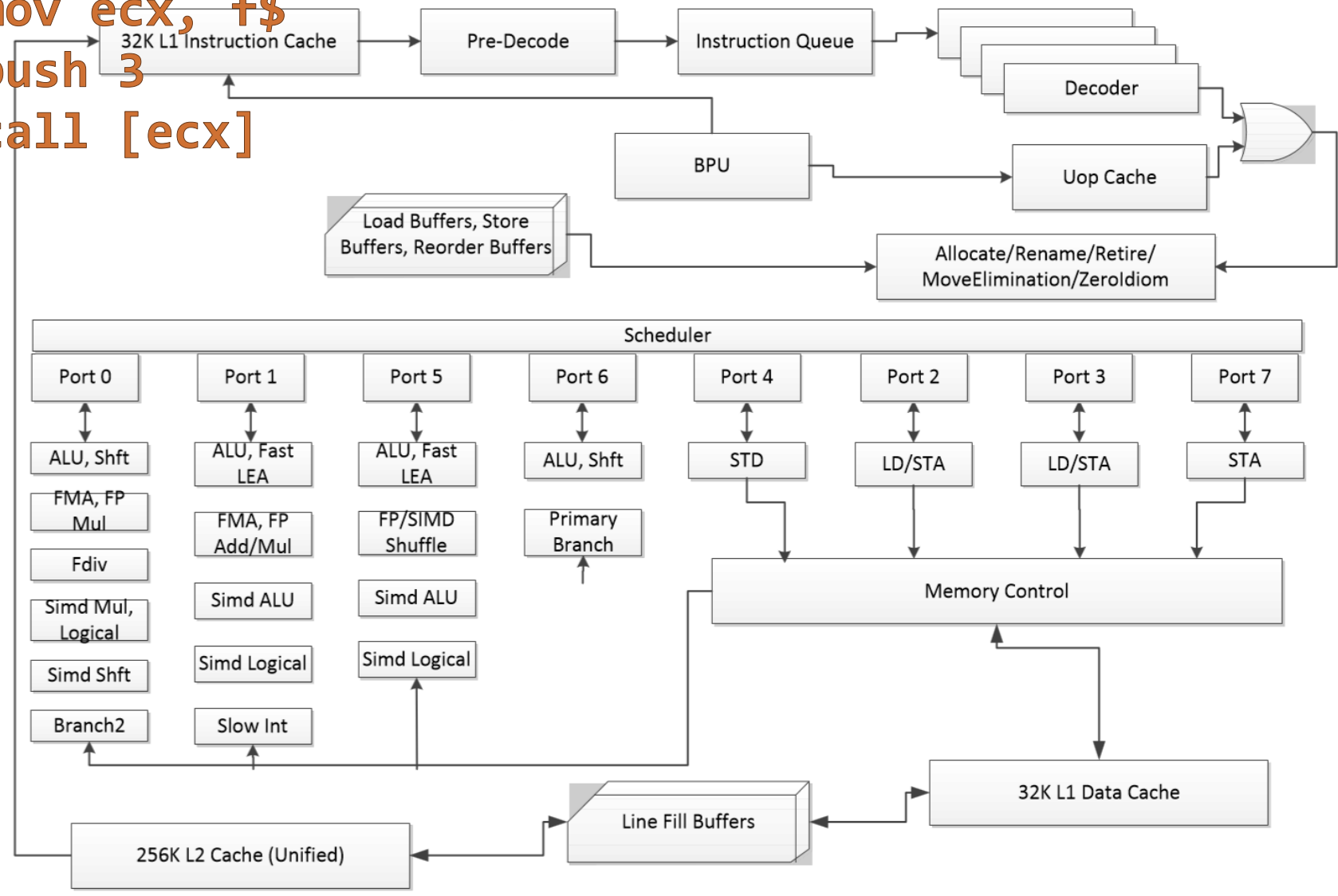
```
int test(PFUNC f) {  
    return f(3);  
}
```

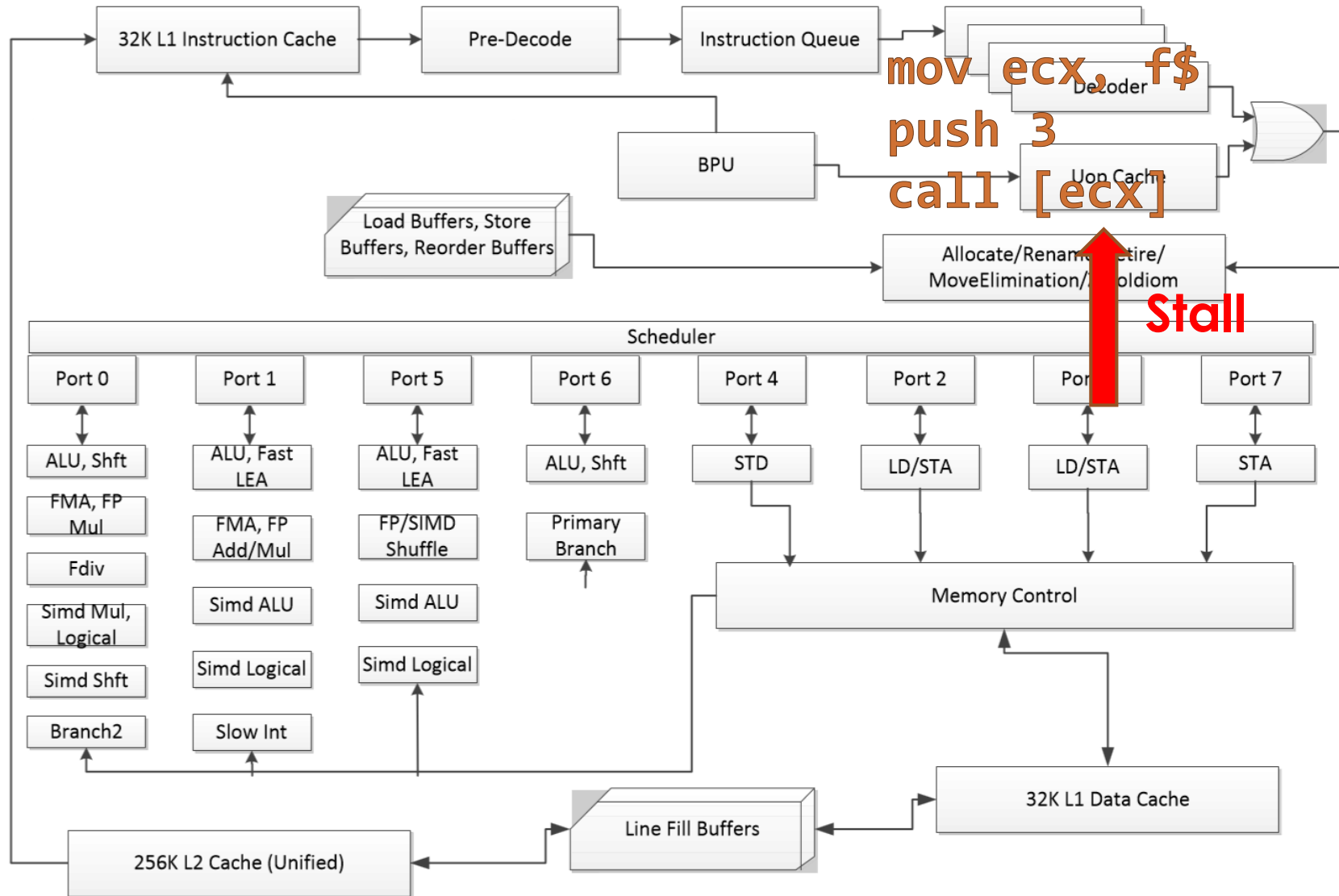
```
mov ecx, f$  
push 3  
call [ecx]
```

**This sucks**



mov ecx, f\$  
push 3  
call [ecx]





```
typedef int (PFUNC)(int);
```

```
int func1(int x) {  
    return x + 100;  
}
```

```
int func2(int x) {  
    return x + 200;  
}
```

```
int test(PFUNC f) {  
    return f(3);  
}
```

```
int test(PFUNC f) {  
    if (f == func1) return func1(3);  
    if (f == func2) return func2(3);  
    return f(3);  
}
```

```
mov ecx, f$  
push 3  
cmp ecx, &func1  
jne $LN1  
call func1  
ret
```

```
$LN1:  
    cmp ecx, &func2  
    jne $LN2  
    call func2  
    ret
```

```
$LN2:  
    call [ecx]
```

Leverage  
branch  
predictor



```

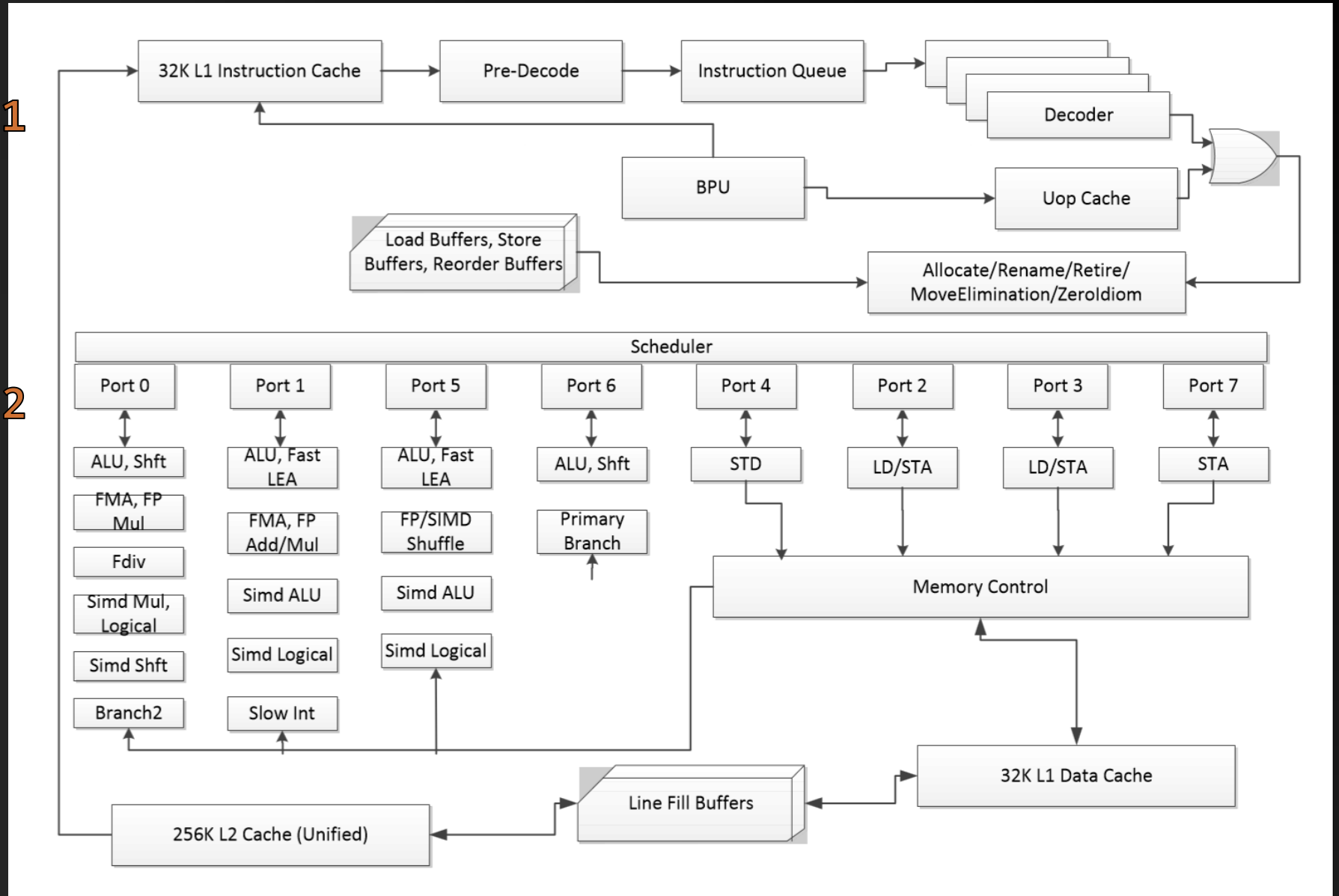
mov ecx, f$
push 3
cmp ecx, &func1
jne $LN1
call func1
ret

```

```

$LN1:
cmp ecx, &func2
jne $LN2
call func2
ret
$LN2:
call [ecx]

```



```

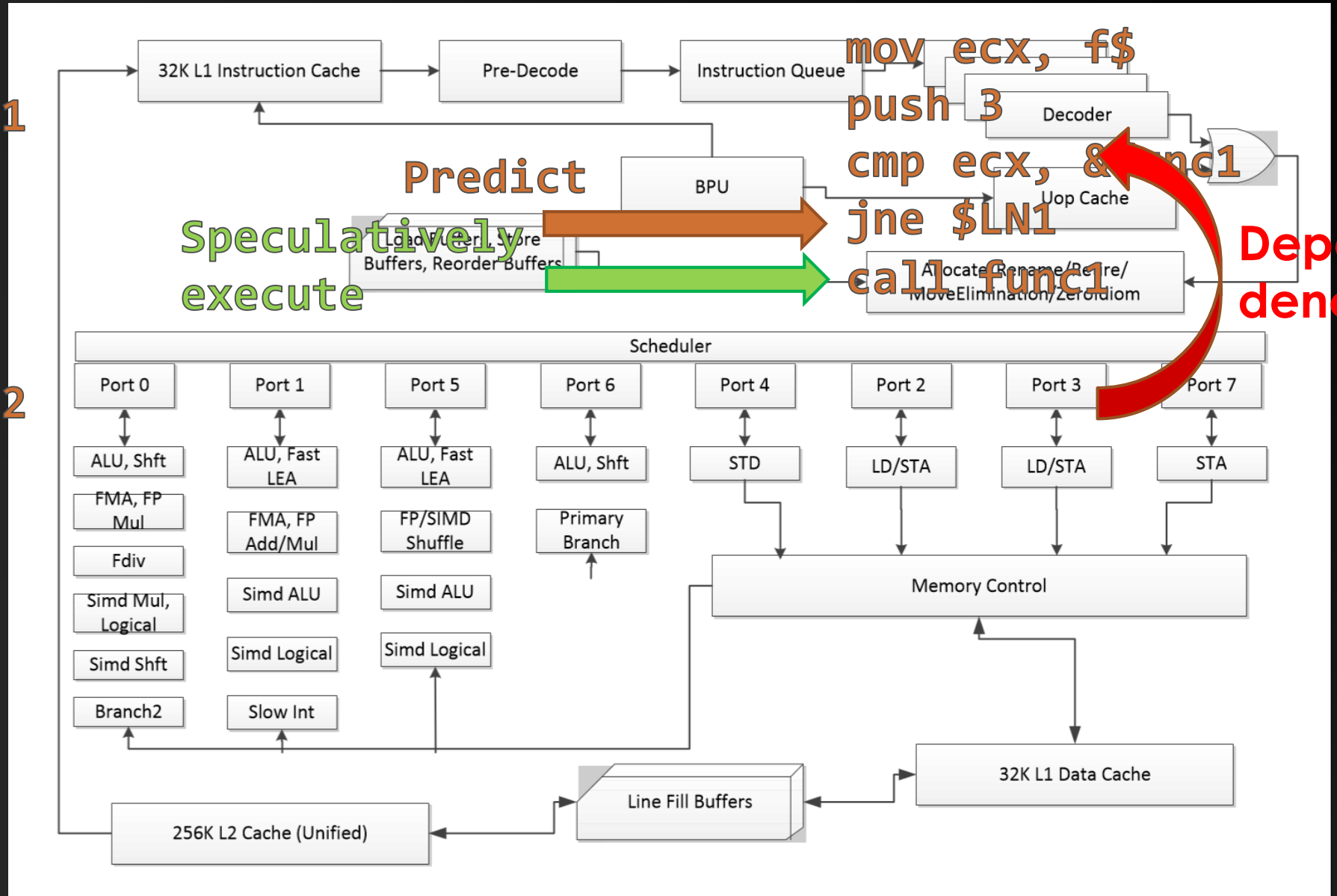
mov ecx, f$
push 3
cmp ecx, &func1
jne $LN1
call func1
ret

```

```

$LN1:
cmp ecx, &func2
jne $LN2
call func2
ret
$LN2:
call [ecx]

```



```
int test(PFUNC f) {
    return f(3);
}
```

```
int test(PFUNC f) {
    if (f == func1) return func1(3);
    if (f == func2) return func2(3);
    return f(3);
}
```

```
mov ecx, f$
push 3
call [ecx] ← Stall
```

```
mov ecx, f$
push 3
cmp ecx, &func1
jne $LN1 ← Not a stall
call func1
ret
$LN1:
cmp ecx, &func2
jne $LN2
call func2
ret
$LN2:
call [ecx]
```

Speedup due to if-statements + branch prediction

You could add if-statements by hand...

But with profile counts, the compiler does it for you.



## Source code:

```
typedef int (PFUNC)(int);
```

```
int func1(int x) {  
    return x + 100;  
}
```

```
int func2(int x) {  
    return x + 200;  
}
```

```
int test(PFUNC f) {  
    return f(3);  
}
```

If counts say test() calls func1() as often as func2():

- Compiler inserts two if-checks
- test() code size increases **5.4x**
- **10%** performance win

If counts say test() calls func1() way more than func2():

- Compiler inserts one if-check
- test() code size increases **3.4x**
- **15%** performance win

If counts say test() calls func1() way more than func2(), and we decide to inline func1():

- Compiler inserts one if-check
- test() code size increases **2.7x**
- **30%** performance win

```
if (f == func1)  
    return func1(3);  
if (f == func2)  
    return func2(3);  
return f(3);
```

```
if (f == func1)  
    return func1(3);  
return f(3);
```

```
if (f == func1)  
    return 103;  
return f(3);
```

**All compiler driven – no code changes!**

THAT'S NICE, BUT I DON'T USE FUNCTION POINTERS

```

class Base {
public:
    virtual int func(int x) = 0;
};

class A : public Base {
    int func(int x) { return x + 100; };
};

class B : public Base {
    int func(int x) { return x + 200; };
};

class C : public Base {
    int func(int x) { return x + 300; };
};

```

```

int test(Base *x) {
    return x->foo(3);
}

```

Load  
vtable

Load right  
'func'

Push  
argument

Indirect call

```

mov    ecx, x$
mov    eax, [ecx]
push   3
call   [eax]

```

**Compiler-driven speculative  
devirtualization & inlining**

# RECAP & OTHER RESOURCES

COMPILER HAS TO TAKE ADVANTAGE OF SILICON

GUARD OPTIMIZATIONS WITH RUNTIME CHECKING

`/Qvec-report:2` MESSAGES (15XX CODES ~ RUNTIME CHECKS)

PROFILE COUNTS: PROFILE GUIDED OPTIMIZATIONS

PROFILING TOOLS

VISUAL STUDIO PERFORMANCE ANALYSIS

INTEL VTUNE AMPLIFIER XE

AMD CODEXL

COMPILER SWITCHES

<http://msdn.microsoft.com>

AUTOMATIC VECTORIZATION BLOG & COOKBOOK

<http://blogs.msdn.com/b/nativeconcurrency>

VISUAL C++ BLOG

<http://blogs.msdn.com/b/vcblog/>

CHANNEL 9 GOING NATIVE

<http://channel9.msdn.com/Shows/C9-GoingNative>

Q&A

BACKUP SLIDES

# WILD AND CRAZY RUNTIME CHECKS

```
for (int i=0; i<1000; i++)  
    a[i] = b[i] * 2.0f;
```

Range of a: `&a[0]` to `&a[999]`

Range of b: `&b[0]` to `&b[999]`



# WILD AND CRAZY RUNTIME CHECKS

```
for (int i=0; i<1000; i++)  
    a[i] = b[i+1] * 2.0f;
```

Range of a: `&a[0]` to `&a[999]`

Range of b: `&b[1]` to `&b[1000]`

# WILD AND CRAZY RUNTIME CHECKS

```
for (int i=0; i<1000; i++)  
    a[i] = b[i+1] + b[i+5];
```

Range of a: `&a[0]` to `&a[999]`  
Range of b: `&b[1]` to `&b[1004]`

Messup in the presentation slides. B ends at `b[1004]`.

Another reason why the compiler should do this for you!



# WILD AND CRAZY RUNTIME CHECKS

```
for (int i=0; i<1000; i++)  
    a[i] = b[i+1] + b[i+x];
```

Range of a: `&a[0]` to `&a[999]`

Range of b: `&b[?]` to `&b[?]`

# WILD AND CRAZY RUNTIME CHECKS

```
for (int i=lb; i<ub; i++)  
    a[i] = b[i*i];
```

Range of a: `&a[lb]` to `&a[ub]`

Range of b: `&b[?]` to `&b[?]`