

Optimization: Misc. + SIMD

changelog

Changes since first lecture:

11 November 2021: cleanup code in vector exercise 2

11 November 2021: add explanation slide for vector exercise 2

17 November 2021: post updated vector exercise 2 explanation slide to PDF

last time

compiler limitations

- match behavior on corner cases
- working across files/functions
- code size tradeoffs

aliasing: two names for one value

function inlining

- replace function call with function body
- saves function call/return/argument setup instructions
- probably extra copies of function body

logistics note

my office hours next week

moved+reduced on Wednesday (10:30a-11:30a instead of 11a-1p)

will help out with Tuesday afternoon

compiler limitations

needs to generate code that does the same thing...

...even in corner cases that “obviously don’t matter”

often doesn’t ‘look into’ a method

needs to assume it might do anything

can’t predict what inputs/values will be

e.g. lots of loop iterations or few?

can’t understand code size versus speed tradeoffs

remove redundant operations (1)

```
int number_of_As(const char *str) {  
    int count = 0;  
    for (int i = 0; i < strlen(str); ++i) {  
        if (str[i] == 'a')  
            count++;  
    }  
    return count;  
}
```

remove redundant operations (1, fix)

```
int number_of_As(const char *str) {  
    int count = 0;  
    int length = strlen(str);  
    for (int i = 0; i < length; ++i) {  
        if (str[i] == 'a')  
            count++;  
    }  
    return count;  
}
```

call strlen once, not once per character!

Big-Oh improvement!

remove redundant operations (1, fix)

```
int number_of_As(const char *str) {  
    int count = 0;  
    int length = strlen(str);  
    for (int i = 0; i < length; ++i) {  
        if (str[i] == 'a')  
            count++;  
    }  
    return count;  
}
```

call strlen once, not once per character!

Big-Oh improvement!

remove redundant operations (2)

```
int shiftArray(int *source, int *dest, int N, int amount) {  
    for (int i = 0; i < N; ++i) {  
        if (i + amount < N)  
            dest[i] = source[i + amount];  
        else  
            dest[i] = source[N - 1];  
    }  
}
```

compare $i + \text{amount}$ to N many times

remove redundant operations (2, fix)

```
int shiftArray(int *source, int *dest, int N, int amount) {  
    int i;  
    for (i = 0; i + amount < N; ++i) {  
        dest[i] = source[i + amount];  
    }  
    for (; i < N; ++i) {  
        dest[i] = source[N - 1];  
    }  
}
```

eliminate comparisons

compiler limitations

needs to generate code that does the same thing...
...even in corner cases that “obviously don’t matter”

often doesn’t ‘look into’ a method

needs to assume it might do anything

can’t predict what inputs/values will be

e.g. lots of loop iterations or few?

can’t understand code size versus speed tradeoffs

exercise: when optimizations backfire...

Which of these optimizations are likely to **increase** machine code size?
(**Select all that apply.**)

Which of these optimizations are likely to **increase** number of instructions executed? (**Select all that apply.**)

- A. cache blocking
- B. function inlining
- C. loop unrolling
- D. moving a calculation outside a loop
- E. multiple accumulators (after loop unrolling)

looplab speeds on my desktop

original assembly: 2.0 cycles/element

unrolled x2: 1.0 cycles element

unrolled x4: 1.0 cycles element

unrolled x8: 1.0 cycles element

unrolled x8, 4 accumulators: 0.5 cycles element

looplab speeds on my desktop

original assembly: 2.0 cycles/element

unrolled x2: 1.0 cycles element

unrolled x4: 1.0 cycles element

unrolled x8: 1.0 cycles element

unrolled x8, 4 accumulators: 0.5 cycles element

Clang 6 optimized code: 0.13 cycles/element

GCC optimized code: 0.14 cycles/element

looplab speeds on my desktop

original assembly: 2.0 cycles/element

unrolled x2: 1.0 cycles element

unrolled x4: 1.0 cycles element

unrolled x8: 1.0 cycles element

unrolled x8, 4 accumulators: 0.5 cycles element

Clang 6 optimized code: 0.13 cycles/element

GCC optimized code: 0.14 cycles/element

how? instructions that add *16 pairs of shorts* at once!

“vector” or “SIMD” (single instruction multiple data) instruction

unvectorized add (original)

```
unsigned int A[512], B[512];  
...  
for (int i = 0; i < N; i += 1) {  
    A[i] = A[i] + B[i];  
}
```


unvectorized add (unrolled)

```
unsigned int A[512], B[512];  
...  
for (int i = 0; i < 512; i += 8) {  
    A[i+0] = A[i+0] + B[i+0];  
    A[i+1] = A[i+1] + B[i+1];  
    A[i+2] = A[i+2] + B[i+2];  
    A[i+3] = A[i+3] + B[i+3];  
    A[i+4] = A[i+4] + B[i+4];  
    A[i+5] = A[i+5] + B[i+5];  
    A[i+6] = A[i+6] + B[i+6];  
    A[i+7] = A[i+7] + B[i+7];  
}
```

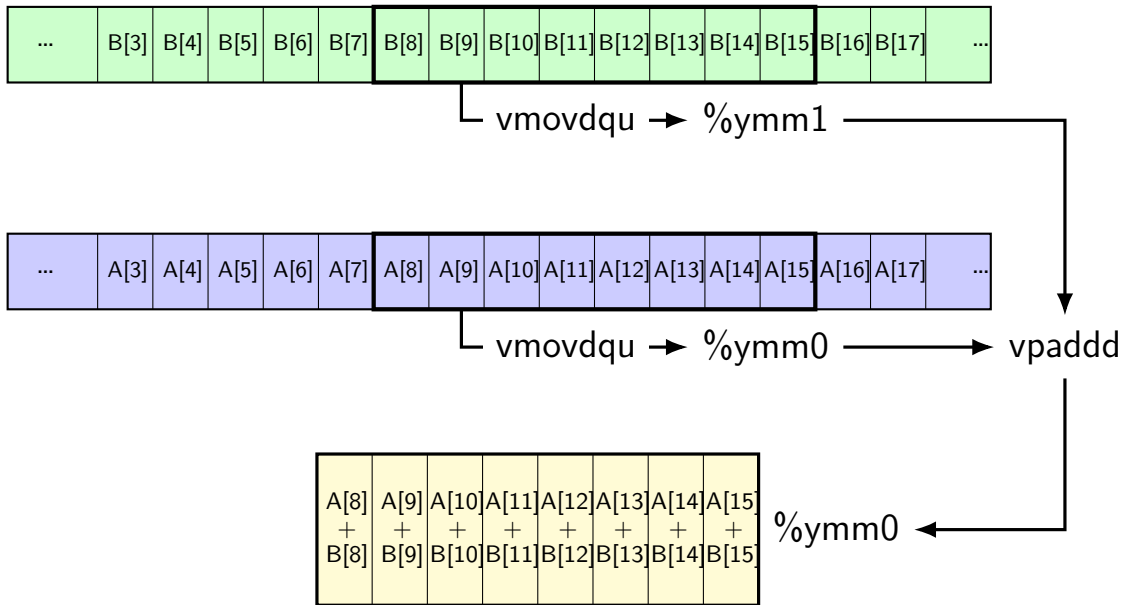
goal: use SIMD add instruction to do all 8 adds above
SIMD = single instruction, multiple data

desired assembly

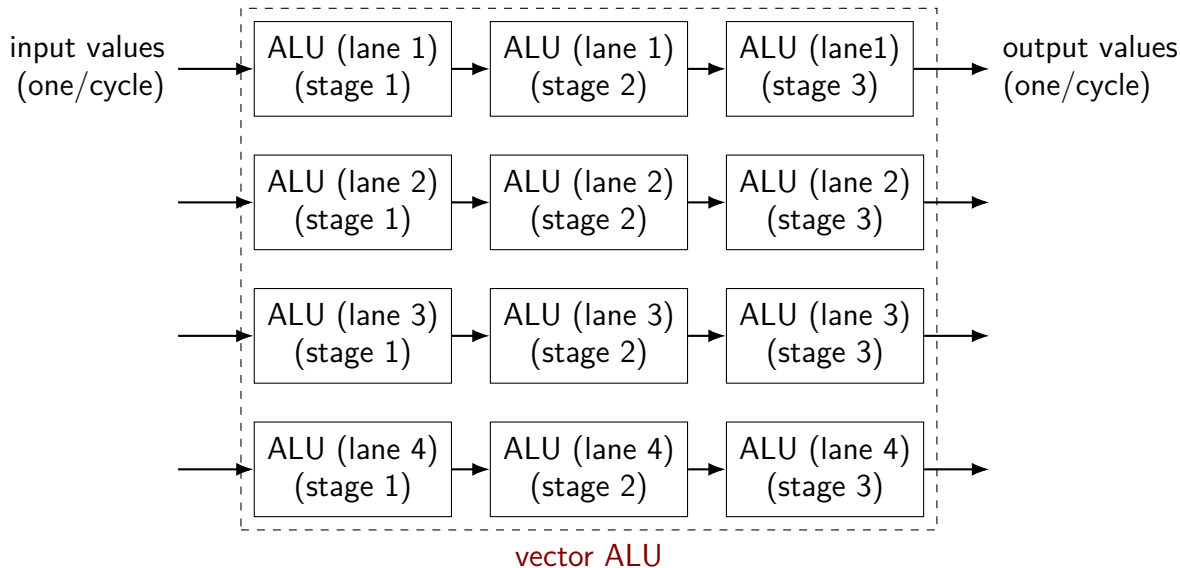
```
xor %rax, %rax
the_loop:
    vmovdqu A(%rax), %ymm0
    vmovdqu B(%rax), %ymm1
    vpaddq %ymm1, %ymm0, %ymm0
    vmovdqu %ymm0, A(%rax)
    addq $32, %rax
    cmpq $2048, %rax
    jne the_loop
```

```
/* load 256 bits of A into ymm0 */
/* load 256 bits of B into ymm1 */
/* ymm1 + ymm0 -> ymm0 */
/* store ymm0 into A */
/* increment index by 32 bytes */
/* offset < 2048 (= 512 * 4) bytes */
```

vector add picture



one view of vector functional units



why vector instructions?

lots of logic not dedicated to computation

- instruction queue

- reorder buffer

- instruction fetch

- branch prediction

- ...

adding vector instructions — little extra control logic

...but a lot more computational capacity

vector instructions and compilers

compilers can sometimes figure out how to use vector instructions
(and have gotten much, much better at it over the past decade)

but easily messed up:

- by aliasing

- by conditionals

- by some operation with no vector instruction

- ...

fickle compiler vectorization (1)

GCC 8.2 and Clang 7.0 generate vector instructions for this:

```
#define N 1024
void foo(unsigned int *A, unsigned int *B) {
    for (int k = 0; k < N; ++k)
        for (int i = 0; i < N; ++i)
            for (int j = 0; j < N; ++j)
                B[i * N + j] += A[i * N + k] * A[k * N + j];
}
```

but not:

```
#define N 1024
void foo(unsigned int *A, unsigned int *B) {
    for (int i = 0; i < N; ++i)
        for (int j = 0; j < N; ++j)
            for (int k = 0; k < N; ++k)
                B[i * N + j] += A[i * N + k] * A[j * N + k];
}
```

fickle compiler vectorization (2)

Clang 5.0.0 generates vector instructions for this:

```
void foo(int N, unsigned int *A, unsigned int *B) {  
    for (int k = 0; k < N; ++k)  
        for (int i = 0; i < N; ++i)  
            for (int j = 0; j < N; ++j)  
                B[i * N + j] += A[i * N + k] * A[k * N + j];  
}
```

but not: (fixed in later versions)

```
void foo(long N, unsigned int *A, unsigned int *B) {  
    for (long k = 0; k < N; ++k)  
        for (long i = 0; i < N; ++i)  
            for (long j = 0; j < N; ++j)  
                B[i * N + j] += A[i * N + k] * A[k * N + j];  
}
```


vector intrinsics

if compiler doesn't work...

could write vector instruction assembly by hand

second option: “intrinsic functions”

C functions that compile to particular instructions

vector intrinsics: add example

```
int A[512], B[512];
```

```
for (int i = 0; i < 512; i += 8) {  
    // "si256" --> 256 bit integer  
    // a_values = {A[i], A[i+1], ..., A[i+7]} (8 x 32 bits)  
    __m256i a_values = _mm256_loadu_si256((__m256i*) &A[i]);  
    // b_values = {B[i], B[i+1] ..., A[i+7]} (8 x 32 bits)  
    __m256i b_values = _mm256_loadu_si256((__m256i*) &B[i]);  
  
    // add eight 32-bit integers  
    // sums = {A[i] + B[i], A[i+1] + B[i+1], ....., A[i+7] + B[i+7]}  
    __m256i sums = _mm256_add_epi32(a_values, b_values);  
  
    // {A[i], A[i+1], A[i+2], A[i+3], ..., A[i+7]} = sums  
    _mm256_storeu_si256((__m256i*) &A[i], sums);  
}
```

vector intrinsics: add example

special type `__m256i` — “256 bits of integers”
other types: `__m256` (floats), `__m128d` (doubles)

```
int A[512]
for (int i = 0; i < 512; i += 8) {
    // "si256" --> 256 bit integer
    // a_values = {A[i], A[i+1], ..., A[i+7]} (8 x 32 bits)
    __m256i a_values = _mm256_loadu_si256((__m256i*) &A[i]);
    // b_values = {B[i], B[i+1] ..., A[i+7]} (8 x 32 bits)
    __m256i b_values = _mm256_loadu_si256((__m256i*) &B[i]);

    // add eight 32-bit integers
    // sums = {A[i] + B[i], A[i+1] + B[i+1], ..., A[i+7] + B[i+7]}
    __m256i sums = _mm256_add_epi32(a_values, b_values);

    // {A[i], A[i+1], A[i+2], A[i+3], ..., A[i+7]} = sums
    _mm256_storeu_si256((__m256i*) &A[i], sums);
}
```

vector intrinsics: add example

i functions to store/load
s i256 means “256-bit integer value”
f u for “unaligned” (otherwise, pointer address must be multiple of 32)

```
// "si256" --> 256 bit integer
// a_values = {A[i], A[i+1], ..., A[i+7]} (8 x 32 bits)
__m256i a_values = _mm256_loadu_si256((__m256i*) &A[i]);
// b_values = {B[i], B[i+1] ..., A[i+7]} (8 x 32 bits)
__m256i b_values = _mm256_loadu_si256((__m256i*) &B[i]);

// add eight 32-bit integers
// sums = {A[i] + B[i], A[i+1] + B[i+1], ..., A[i+7] + B[i+7]}
__m256i sums = _mm256_add_epi32(a_values, b_values);

// {A[i], A[i+1], A[i+2], A[i+3], ..., A[i+7]} = sums
_mm256_storeu_si256((__m256i*) &A[i], sums);
}
```

vector intrinsics: add example

```
int A[512], B[512];
```

```
for (int i = 0; i < 512; i += 8) {
```

```
    // "si256" --> function to add  
    // a_values = a_values (8 x 32 bits)  
    __m256i a_values = _mm256_loadu_si256((__m256i*) &A[i]);
```

```
    // b_values = {B[i], B[i+1] ..., A[i+7]} (8 x 32 bits)  
    __m256i b_values = _mm256_loadu_si256((__m256i*) &B[i]);
```

```
    // add eight 32-bit integers
```

```
    // sums = {A[i] + B[i], A[i+1] + B[i+1], ..., A[i+7] + B[i+7]}
```

```
    __m256i sums = _mm256_add_epi32(a_values, b_values);
```

```
    // {A[i], A[i+1], A[i+2], A[i+3], ..., A[i+7]} = sums
```

```
    _mm256_storeu_si256((__m256i*) &A[i], sums);
```

```
}
```

vector intrinsics: different size

```
long A[512], B[512]; /* instead of int */
...
for (int i = 0; i < 512; i += 4) {
    // a_values = {A[i], A[i+1], A[i+2], A[i+3]} (4 x 64 bits)
    __m256i a_values = _mm256_loadu_si256((__m256i*) &A[i]);
    // b_values = {B[i], B[i+1], B[i+2], B[i+3]} (4 x 64 bits)
    __m256i b_values = _mm256_loadu_si256((__m256i*) &B[i]);
    // add four 64-bit integers: vpaddq %ymm0, %ymm1
    // sums = {A[i] + B[i], A[i+1] + B[i+1], ...}
    __m256i sums = _mm256_add_epi64(a_values, b_values);
    // {A[i], A[i+1], A[i+2], A[i+3]} = sums
    _mm256_storeu_si256((__m256i*) &A[i], sums);
}
```

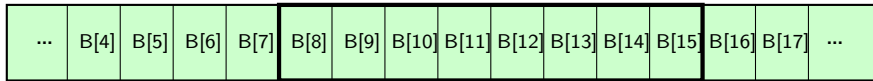
vector intrinsics: different size

```
long A[512], B[512]; /* instead of int */
```

```
...
```

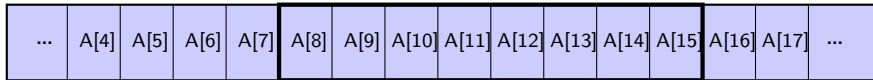
```
for (int i = 0; i < 512; i += 4) {  
    // a_values = {A[i], A[i+1], A[i+2], A[i+3]} (4 x 64 bits)  
    __m256i a_values = _mm256_loadu_si256((__m256i*) &A[i]);  
    // b_values = {B[i], B[i+1], B[i+2], B[i+3]} (4 x 64 bits)  
    __m256i b_values = _mm256_loadu_si256((__m256i*) &B[i]);  
    // add four 64-bit integers: vpaddq %ymm0, %ymm1  
    // sums = {A[i] + B[i], A[i+1] + B[i+1], ...}  
    __m256i sums = _mm256_add_epi64(a_values, b_values);  
    // {A[i], A[i+1], A[i+2], A[i+3]} = sums  
    _mm256_storeu_si256((__m256i*) &A[i], sums);  
}
```

vector add picture (intrinsics)



`_mm256_loadu_si256`
(asm: vmovdqu)

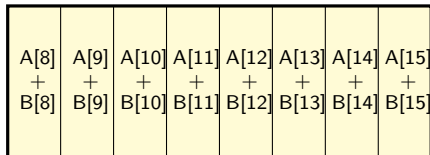
→ `b_values`
(`%ymm1?`)



`_mm256_loadu_si256`
(asm: vmovdqu)

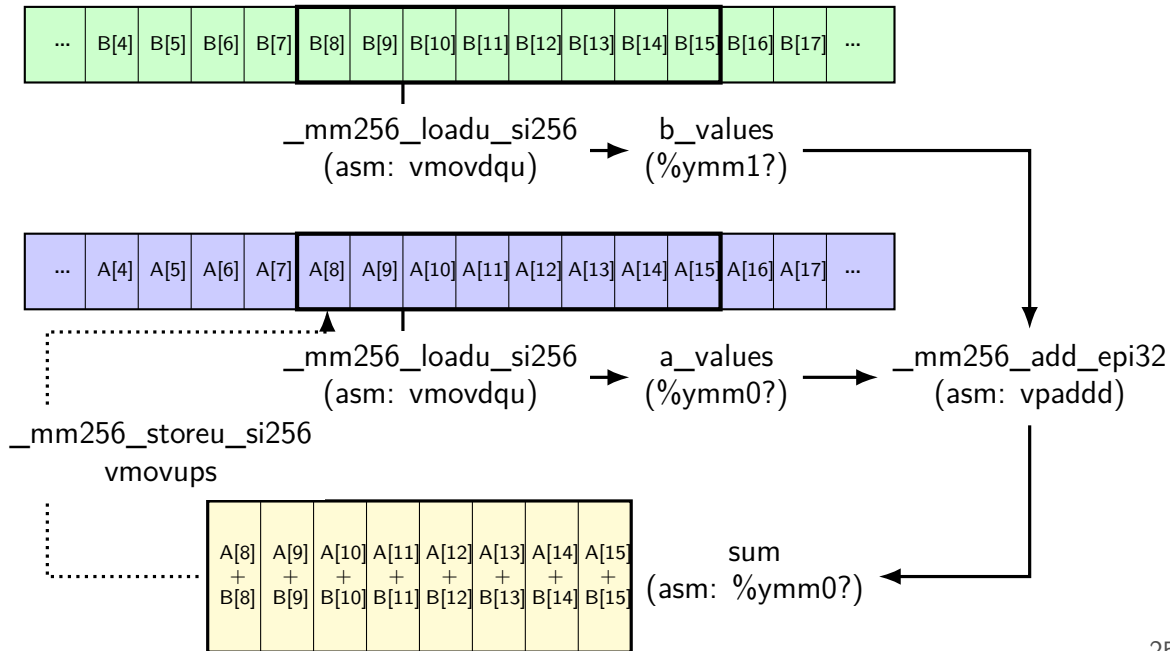
→ `a_values`
(`%ymm0?`)

→ `_mm256_add_epi32`
(asm: vpaddd)



`sum`
(asm: `%ymm0?`)

vector add picture (intrinsics)



exercise

```
long foo[8] = {1,1,2,2,3,3,4,4};
long bar[8] = {2,2,2,3,3,3,4,4};
__mm256i foo0_as_vector = _mm256_loadu_si256((__m256i*)&foo[0])
__mm256i foo4_as_vector = _mm256_loadu_si256((__m256i*)&foo[4])
__mm256i bar0_as_vector = _mm256_loadu_si256((__m256i*)&bar[0])

__mm256i result = _mm256_add_epi64(foo0_as_vector, foo4_as_vector);
result = _mm256_mullo_epi64(result, bar0_as_vector);
_mm256_storeu_si256((__m256i*) &bar[4], result);
```

Final value of bar array?

- A. {2,2,2,3,12,12,24,24} B. {2,2,2,3,15,15,28,28}
C. {2,2,2,3,10,10,20,20} D. {12,12,24,24,3,3,4,4}
E. {14,14,26,27,3,3,4,4} F. {14,14,26,27,12,12,24,24}
G. something else

128-bit version, too

history: 256-bit vectors added in extension called AVX (c. 2011)

before: 128-bit vectors added in extension called SSE (c. 1999)

128-bit intrinsics exist, too:

`__m256i` becomes `__m128i`

`_mm256_add_epi32` becomes `_mm_add_epi32`

`_mm256_loadu_si256` becomes `_mm_loadu_si128`

matrix multiply

```
void matmul(unsigned int *A, unsigned int *B, unsigned int *C)
    for (int k = 0; k < N; ++k)
        for (int i = 0; i < N; ++i)
            for (int j = 0; j < N; ++j)
                C[i * N + j] += A[i * N + k] * B[k * N + j];
}
```

(simple version, no cache blocking, no avoiding aliasing between C, B, A,...)

matmul unrolled

```
void matmul(unsigned int *A, unsigned int *B, unsigned int *C) {  
    for (int k = 0; k < N; ++k) {  
        for (int i = 0; i < N; ++i)  
            for (int j = 0; j < N; j += 8) {  
                /* goal: vectorize this */  
                C[i * N + j + 0] += A[i * N + k] * B[k * N + j + 0];  
                C[i * N + j + 1] += A[i * N + k] * B[k * N + j + 1];  
                C[i * N + j + 2] += A[i * N + k] * B[k * N + j + 2];  
                C[i * N + j + 3] += A[i * N + k] * B[k * N + j + 3];  
                C[i * N + j + 4] += A[i * N + k] * B[k * N + j + 4];  
                C[i * N + j + 5] += A[i * N + k] * B[k * N + j + 5];  
                C[i * N + j + 6] += A[i * N + k] * B[k * N + j + 6];  
                C[i * N + j + 7] += A[i * N + k] * B[k * N + j + 7];  
            }  
    }  
}
```

(NB: would probably also want to do cache blocking...)

handy intrinsic functions for matmul

`_mm256_set1_epi32` — load eight copies of a 32-bit value into a 256-bit value

instructions generated vary; one example: `vmovd + vpbroadcastd`

`_mm256_mullo_epi32` — multiply eight pairs of 32-bit values, give lowest 32-bits of results

generates `vpmulld`

vectorizing matmul

/ goal: vectorize this */*

`C[i * N + j + 0] += A[i * N + k] * B[k * N + j + 0];`

`C[i * N + j + 1] += A[i * N + k] * B[k * N + j + 1];`

`...`

`C[i * N + j + 6] += A[i * N + k] * B[k * N + j + 6];`

`C[i * N + j + 7] += A[i * N + k] * B[k * N + j + 7];`

vectorizing matmul

/ goal: vectorize this */*

`C[i * N + j + 0] += A[i * N + k] * B[k * N + j + 0];`

`C[i * N + j + 1] += A[i * N + k] * B[k * N + j + 1];`

...

`C[i * N + j + 6] += A[i * N + k] * B[k * N + j + 6];`

`C[i * N + j + 7] += A[i * N + k] * B[k * N + j + 7];`

// load eight elements from C

`Cij = _mm256_loadu_si256((__m256i*) &C[i * N + j + 0]);`

... // manipulate vector here

// store eight elements into C

`_mm_storeu_si256((__m256i*) &C[i * N + j + 0], Cij);`

vectorizing matmul

/ goal: vectorize this */*

`C[i * N + j + 0] += A[i * N + k] * B[k * N + j + 0];`

`C[i * N + j + 1] += A[i * N + k] * B[k * N + j + 1];`

...

`C[i * N + j + 6] += A[i * N + k] * B[k * N + j + 6];`

`C[i * N + j + 7] += A[i * N + k] * B[k * N + j + 7];`

// load eight elements from B

`Bkj = _mm256_loadu_si256((__m256i*) &B[k * N + j + 0]);`

*... // multiply each by B[i * N + k] here*

vectorizing matmul

/ goal: vectorize this */*

`C[i * N + j + 0] += A[i * N + k] * B[k * N + j + 0];`

`C[i * N + j + 1] += A[i * N + k] * B[k * N + j + 1];`

...

`C[i * N + j + 6] += A[i * N + k] * B[k * N + j + 6];`

`C[i * N + j + 7] += A[i * N + k] * B[k * N + j + 7];`

*// load eight elements starting with B[k * n + j]*

`Bkj = _mm256_loadu_si256((__m256i*) &B[k * N + j + 0]);`

*// load eight copies of A[i * N + k]*

`Aik = _mm256_set1_epi32(A[i * N + k]);`

// multiply each pair

`multiply_results = _mm256_mullo_epi32(Aik, Bkj);`

vectorizing matmul

/ goal: vectorize this */*

`C[i * N + j + 0] += A[i * N + k] * B[k * N + j + 0];`

`C[i * N + j + 1] += A[i * N + k] * B[k * N + j + 1];`

...

`C[i * N + j + 6] += A[i * N + k] * B[k * N + j + 6];`

`C[i * N + j + 7] += A[i * N + k] * B[k * N + j + 7];`

`Cij = _mm256_add_epi32(Cij, multiply_results);`

// store back results

`_mm256_storeu_si256(..., Cij);`

matmul vectorized

```
__m256i Cij, Bkj, Aik, multiply_results;

// Cij = {Ci,j, Ci,j+1, Ci,j+2, ..., Ci,j+7}
Cij = _mm256_loadu_si256((__m256i*) &C[i * N + j]);
// Bkj = {Bk,j, Bk,j+1, Bk,j+2, ..., Bk,j+7}
Bkj = _mm256_loadu_si256((__m256i*) &B[k * N + j]);

// Aik = {Ai,k, Ai,k, ..., Ai,k}
Aik = _mm256_set1_epi32(A[i * N + k]);

// Aik_times_Bkj = {Ai,k × Bk,j, Ai,k × Bk,j+1, Ai,k × Bk,j+2, ..., Ai,k × Bk,j+7}
multiply_results = _mm256_mullo_epi32(Aij, Bkj);

// Cij = {Ci,j + Ai,k × Bk,j, Ci,j+1 + Ai,k × Bk,j+1, ...}
Cij = _mm256_add_epi32(Cij, multiply_results);

// store Cij into C
_mm256_storeu_si256((__m256i*) &C[i * N + j], Cij);
```

vector exercise (2)

```
long A[1024], B[1024];  
...  
for (int i = 0; i < 1024; i += 1)  
    for (int j = 0; j < 1024; j += 1)  
        A[i] += B[i] * B[j];
```

(casts omitted below to reduce clutter:)

```
for (int i = 0; i < 1024; i += 4) {  
    A_part = _mm256_loadu_si256(&A[i]);  
    Bi_part = _mm256_loadu_si256(&B[i]);  
    for (int j = 0; j < 1024; /* BLANK 1 */) {  
        Bj_part = _mm256_/* BLANK 2 */;  
        A_part = _mm256_add_epi64(A_part,  
                                   _mm256_mullo_epi64(Bi_part, Bj_part));  
    }  
    _mm256_storeu_si256(&A[i], A_part);  
}
```

What goes in BLANK 1 and BLANK 2?

- A. `j += 1, loadu_si256(&B[j])` B. `j += 4, loadu_si256(&B[j])`
C. `j += 1, set1_epi64(B[j])` D. `j += 4, set1_epi64(B[j])`

vector exercise 2 explanation

```
for (int i = 0; i < 1024; i += 1)
    for (int j = 0; j < 1024; j += 1)
        A[i] += B[i] * B[j];
/* -- transformed into -- */
for (int i = 0; i < 1024; i += 4)
    for (int j = 0; j < 1024; j += 1) {
        A[i+0] += B[i+0] * B[j];
        A[i+1] += B[i+1] * B[j];
        A[i+2] += B[i+2] * B[j];
        A[i+3] += B[i+3] * B[j];
    }
```

```
/* not the much harder to vectorize: */
for (int i = 0; i < 1024; i += 1)
    for (int j = 0; j < 1024; j += 4) {
        A[i] += B[i] * B[j+0];
        A[i] += B[i] * B[j+1];
        A[i] += B[i] * B[j+2];
        A[i] += B[i] * B[j+3];
    }
```

moving values in vectors?

sometimes values aren't in the right place in vector

example:

have: [1, 2, 3, 4]

want: [3, 4, 1, 2]

there are instructions/intrinsics for doing this
called shuffling/swizzling/permute/...

sometimes might need combination of them

worst-case: could rearrange on stack..., I guess

example shuffling operation (1)

goal: [1, 2, 3, 4] to [3, 4, 1, 2] (64-bit values)

```
/* x = {1, 2, 3, 4} */  
__m256i x = _mm256_setr_epi64x(1, 2, 3, 4);  
__m256i result = _mm256_permute4x64_epi64(  
    x,  
    /* index 2, then 3, then 0, then 1 */  
    2 | (3 << 2) | (0 << 4) | (1 << 6)  
    /* could also write _MM_SHUFFLE(1, 0, 3, 2) */  
);  
/* result = {3, 4, 1, 2} */
```


other vector instructions

multiple extensions to the X86 instruction set for vector instructions

early versions (128-bit vectors): SSE, SSE2, SSE3, SSSE3, SSE4.1, SSE4.2

128-bit vectors

this class (256-bit): AVX, AVX2

not this class (512+-bit): AVX-512

512-bit vectors

also other ISAs have these: e.g. NEON on ARM, MSA on MIPS, AltiVec/VMX on POWER, ...

GPUs are essentially vector-instruction-specialized CPUs

other vector interfaces

intrinsics (our assignments) one way

some alternate programming interfaces

have compiler do more work than intrinsics

e.g. CUDA, OpenCL, GCC's vector instructions

other vector instructions features

more flexible vector instruction features:

- invented in the 1990s

- often present in GPUs and being rediscovered by modern ISAs

reasonable conditional handling

better variable-length vectors

ability to load/store non-contiguous values

some of these features in AVX2/AVX512

alternate vector interfaces

intrinsics functions/assembly aren't the only way to write vector code

e.g. GCC vector extensions: more like normal C code

- types for each kind of vector

- write + instead of `_mm_add_epi32`

e.g. CUDA (GPUs): looks like writing multithreaded code, but each thread is vector “lane”

optimizing real programs

ask your compiler to try first

spend effort where **it matters**

e.g. 90% of program time spent reading files, but optimize computation?

e.g. 90% of program time spent in routine A, but optimize B?

profilers

first step — tool to determine where you spend time

tools exist to do this for programs

example on Linux: `perf`

example

Samples: 37K of event 'cycles', Event count (approx.): 3736755513					
Children	Self	Command	Shared Object	Symbol	
+ 100.00%	0.00%	hclrs-with-debu	hclrs-with-debuginfo	[.]	_start
+ 100.00%	0.00%	hclrs-with-debu	libc-2.31.so	[.]	__libc_start_main
+ 100.00%	0.00%	hclrs-with-debu	hclrs-with-debuginfo	[.]	main
+ 100.00%	0.00%	hclrs-with-debu	hclrs-with-debuginfo	[.]	std::sys_common::backtrace::__rust_begin_short_backt
+ 100.00%	0.00%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hclrs::main
+ 99.99%	9.75%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hclrs::program::RunningProgram::run
+ 60.37%	31.67%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hclrs::ast::SpannedExpr::evaluate
+ 41.34%	23.29%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hashbrown::map::make_hash
+ 18.08%	18.07%	hclrs-with-debu	hclrs-with-debuginfo	[.]	<std::collections::hash::map::DefaultHasher as core:
+ 16.33%	0.68%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hclrs::program::Program::process_register_banks
+ 9.54%	3.15%	hclrs-with-debu	hclrs-with-debuginfo	[.]	std::collections::hash::map::HashMap<K,V,S>::get
+ 9.10%	9.09%	hclrs-with-debu	libc-2.31.so	[.]	__memcpy_avx2_movbe
+ 6.11%	2.10%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hashbrown::map::HashMap<K,V,S>::get_mut
+ 2.32%	0.88%	hclrs-with-debu	hclrs-with-debuginfo	[.]	std::collections::hash::map::HashMap<K,V,S>::get
+ 1.45%	0.52%	hclrs-with-debu	hclrs-with-debuginfo	[.]	hashbrown::map::HashMap<K,V,S>::insert
0.37%	0.11%	hclrs-with-debu	hclrs-with-debuginfo	[.]	<alloc::string::String as core::clone::Clone>::clone
0.19%	0.19%	hclrs-with-debu	libc-2.31.so	[.]	malloc

an infinite loop

```
int main(void) {  
    while (1) {  
        /* waste CPU time */  
    }  
}
```

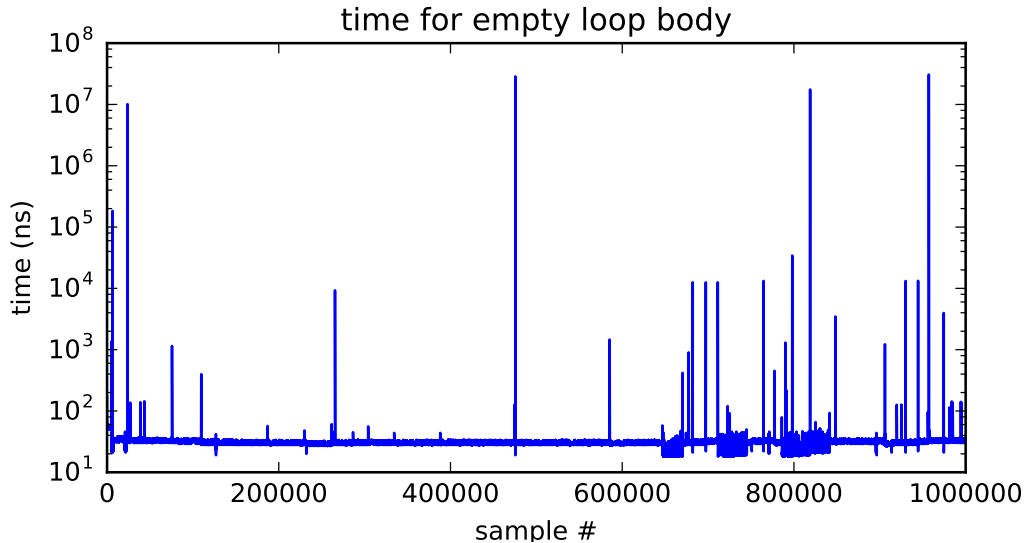
If I run this on a shared department machine, can you still use it?
...if the machine only has one core?

timing nothing

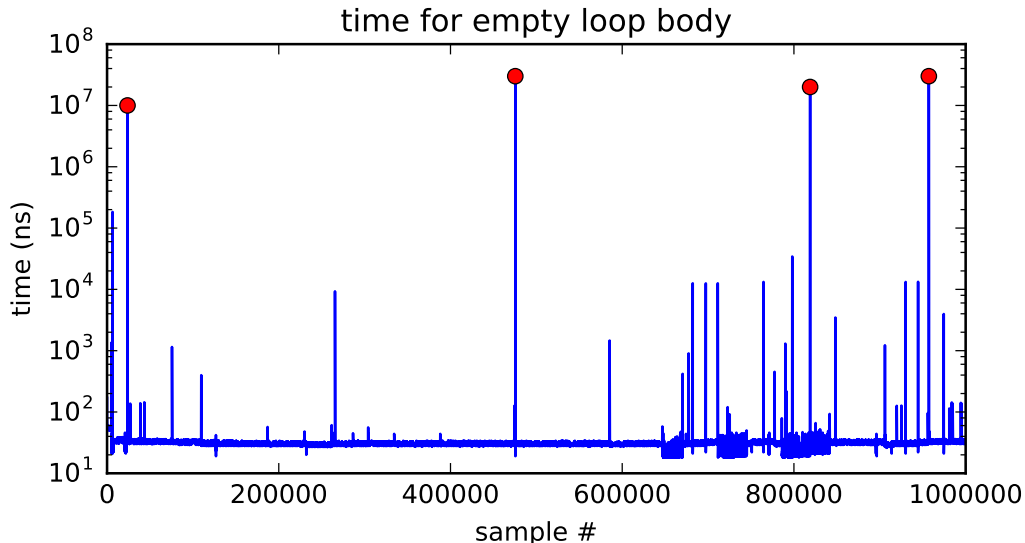
```
long times[NUM_TIMINGS];  
int main(void) {  
    for (int i = 0; i < N; ++i) {  
        long start, end;  
        start = get_time();  
        /* do nothing */  
        end = get_time();  
        times[i] = end - start;  
    }  
    output_timings(times);  
}
```

same instructions — same difference each time?

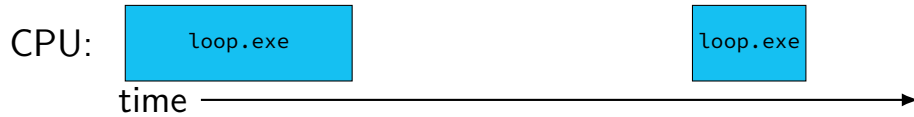
doing nothing on a busy system



doing nothing on a busy system



time multiplexing



time multiplexing



...

```
call get_time
```

// whatever get_time does

```
movq %rax, %rbp
```

———— million cycle delay ————

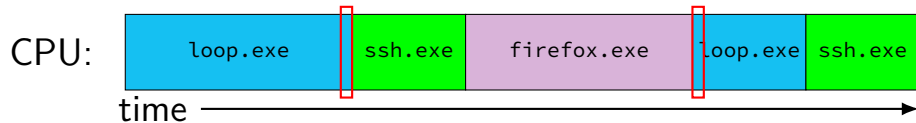
```
call get_time
```

// whatever get_time does

```
subq %rbp, %rax
```

...

time multiplexing



...

```
call get_time
```

```
// whatever get_time does
```

```
movq %rax, %rbp
```

———— million cycle delay ————

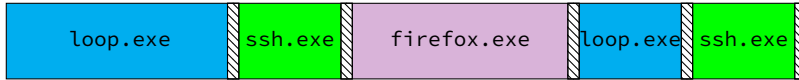
```
call get_time
```


```
// whatever get_time does
```

```
subq %rbp, %rax
```

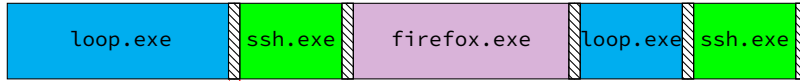
...

time multiplexing really



 = operating system

time multiplexing really



= operating system

exception happens

return from exception

OS and time multiplexing

starts running instead of normal program

mechanism for this: **exceptions** (later)

saves old program counter, registers somewhere

sets new registers, jumps to new program counter

called **context switch**

saved information called **context**

context

all registers values

`%rax %rbx, ..., %rsp, ...`

condition codes

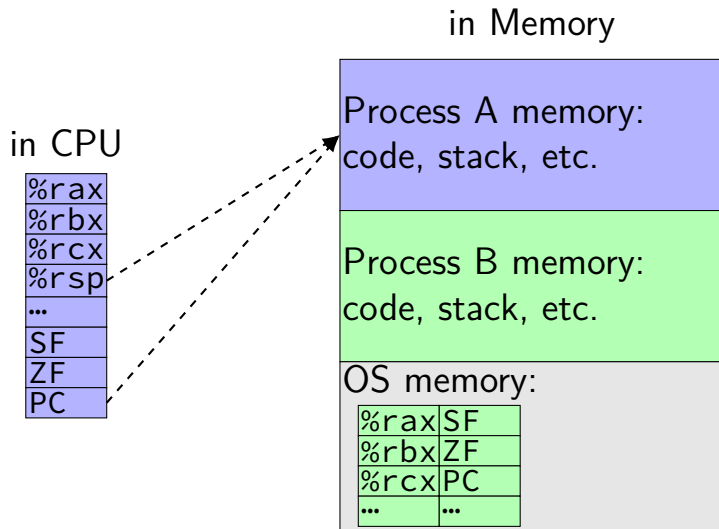
program counter

i.e. all visible state in your CPU except memory

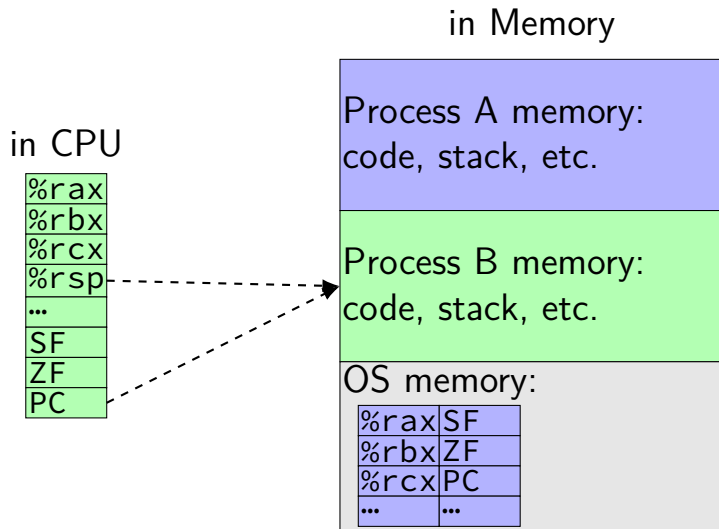
context switch pseudocode

```
context_switch(last, next):  
    copy_preexception_pc last->pc  
    mov rax, last->rax  
    mov rcx, last->rcx  
    mov rdx, last->rdx  
    ...  
    mov next->rdx, rdx  
    mov next->rcx, rcx  
    mov next->rax, rax  
    jmp next->pc
```

contexts (A running)



contexts (B running)



memory protection

reading from another program's memory?

Program A	Program B
<pre>0x10000: .word 42 // ... // do work // ... movq 0x10000, %rax</pre>	<pre><i>// while A is working:</i> movq \$99, %rax movq %rax, 0x10000 ...</pre>

memory protection

reading from another program's memory?

Program A	Program B
<pre>0x10000: .word 42 // ... // do work // ... movq 0x10000, %rax</pre>	<pre><i>// while A is working:</i> movq \$99, %rax movq %rax, 0x10000 ...</pre>

result: %rax is ...

- A. 42 B. 99 C. 0x10000
D. 42 or 99 (depending on timing/program layout/etc)
E. 42 or program might crash (depending on ...)
F. 99 or program might crash (depending on ...)
G. 42 or 99 or program might crash (depending on ...)
H. something else

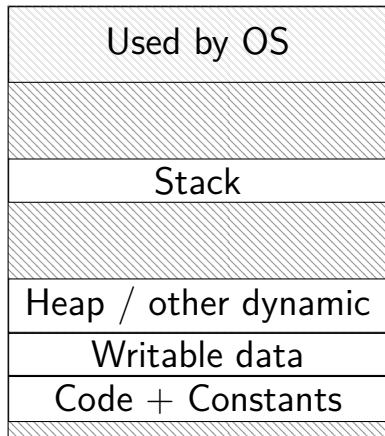
memory protection

reading from another program's memory?

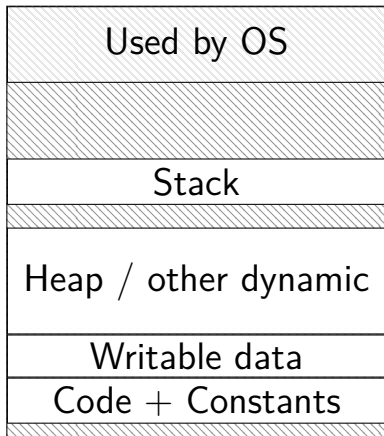
Program A	Program B
<pre>0x10000: .word 42 // ... // do work // ... movq 0x10000, %rax</pre>	<pre><i>// while A is working:</i> movq \$99, %rax movq %rax, 0x10000 ...</pre>
result: %rax is 42 (always)	result: <i>might crash</i>

program memory (two programs)

Program A



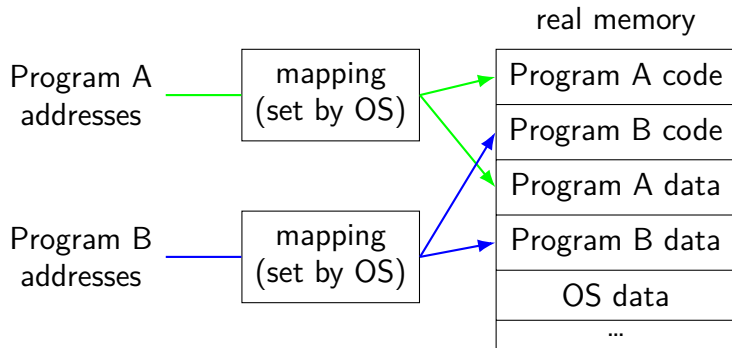
Program B



address space

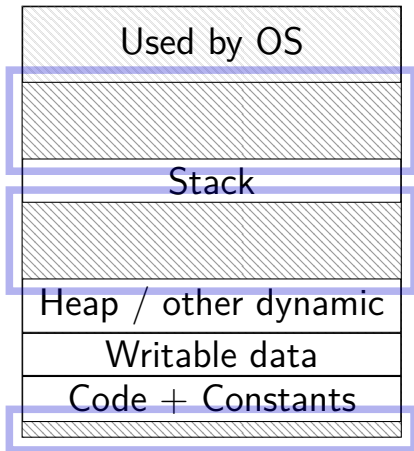
programs have **illusion of own memory**

called a program's **address space**

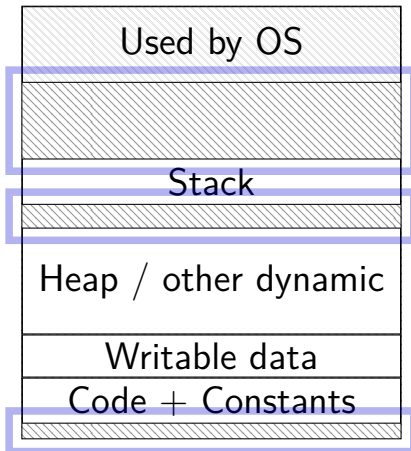


program memory (two programs)

Program A



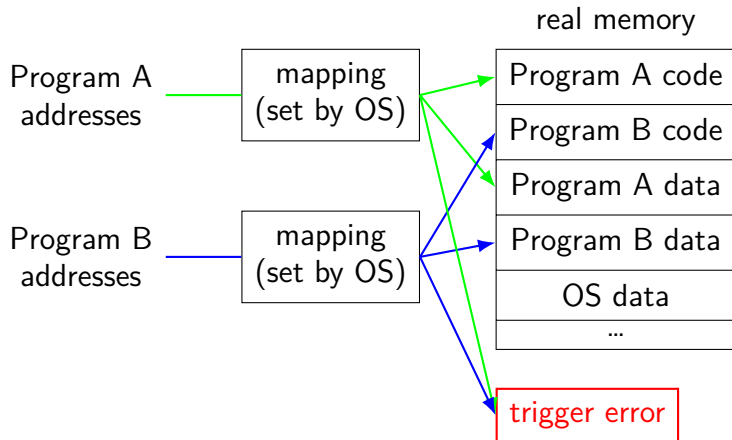
Program B



address space

programs have **illusion of own memory**

called a program's **address space**



address space mechanisms

topic after exceptions

called **virtual memory**

mapping called **page tables**

mapping part of what is changed in context switch

context

all registers values

`%rax %rbx, ..., %rsp, ...`

condition codes

program counter

~~i.e. all visible state in your CPU except memory~~

address space: map from program to real addresses

The Process

process = thread(s) + address space

illusion of **dedicated machine**:

thread = illusion of own CPU

address space = illusion of own memory