# Program Optimization

15-213: Introduction to Computer Systems
10th Lecture, Oct. 1, 2015

**Instructors:**

Randal E. Bryant and David R. O'Hallaron

Maltrakteret af:

Finn Schiermer Andersen

# Vigtigt

- **Bemærk skifte fra beskrivelse af programudførelse som sekventiel til parallel**

- **Udførelse er nu en 2-trins proces**

  - Konstruktion af en afhængighedsgraf på basis af sekventiel semantik af koden

  - Afvikling af operationer I rækkefølge angivet af afhængighedsgrafen, men iøvrigt I parallel

- **Selvfølgeligt underlagt en masse begrænsinger**

# Today

- **Overview**
- **Generally Useful Optimizations**
  - Code motion/precomputation
  - Strength reduction
  - Removing unnecessary procedure calls
- **Optimization Blockers**
  - Procedure calls
  - Memory aliasing (next lecture)
- **Exploiting Instruction-Level Parallelism**
- **Dealing with Conditionals**

# Performance Realities

- ***There's more to performance than asymptotic complexity***

- **Constant factors matter too!**
  - Easily see 10:1 performance range depending on how code is written
  - Must optimize at multiple levels:
    - algorithm, data representations, procedures, and loops
- **Must understand system to optimize performance**
  - How programs are compiled and executed
  - How modern processors + memory systems operate
  - How to measure program performance and identify bottlenecks
  - How to improve performance without destroying code modularity and generality

# Optimizing Compilers

- **Provide efficient mapping of program to machine**
  - register allocation
  - code selection and ordering (scheduling)
  - dead code elimination
  - eliminating minor inefficiencies
- **Don't (usually) improve asymptotic efficiency**
  - up to programmer to select best overall algorithm
  - big-O savings are (often) more important than constant factors
    - but constant factors also matter
- **Have difficulty overcoming "optimization blockers"**
  - potential memory aliasing
  - potential procedure side-effects

# Limitations of Optimizing Compilers

- **Operate under fundamental constraint**
  - Must not cause any change in program behavior
    - Except, possibly when program making use of nonstandard language features
  - Often prevents it from making optimizations that would only affect behavior under pathological conditions.
- **Behavior that may be obvious to the programmer can  be obfuscated by languages and coding styles**
  - e.g., Data ranges may be more limited than variable types suggest
- **Most analysis is performed only within procedures**
  - Whole-program analysis is too expensive in most cases
  - Newer versions of GCC do interprocedural analysis within individual files
    - But, not between code in different files
- **Most analysis is based only on *static* information**
  - Compiler has difficulty anticipating run-time inputs
- **When in doubt, the compiler must be conservative**

# Generally Useful Optimizations

- **Optimizations that you or the compiler should do regardless of processor / compiler**

- **Code Motion**
  - Reduce frequency with which computation performed
    - If it will always produce same result
    - Especially moving code out of loop

```
void set_row(double *a, double *b,
   long i, long n)
{
   long j;
   for (j = 0; j < n; j++)
      a[n*i+j] = b[j];
}
```

→

```
   long j;
   int ni = n*i;
   for (j = 0; j < n; j++)
      a[ni+j] = b[j];
```

# Reduction in Strength

- Replace costly operation with simpler one
- Shift, add instead of multiply or divide

  `16*x -->   x << 4`

  - Utility machine dependent
  - Depends on cost of multiply or divide instruction
    - On Intel Nehalem, integer multiply requires 3 CPU cycles
- Recognize sequence of products

```
for (i = 0; i < n; i++) {
  int ni = n*i;
  for (j = 0; j < n; j++)
    a[ni + j] = b[j];
}
```

⟶

```
int ni = 0;
for (i = 0; i < n; i++) {
  for (j = 0; j < n; j++)
    a[ni + j] = b[j];
  ni += n;
}
```

# Optimization Blocker #1: Procedure Calls
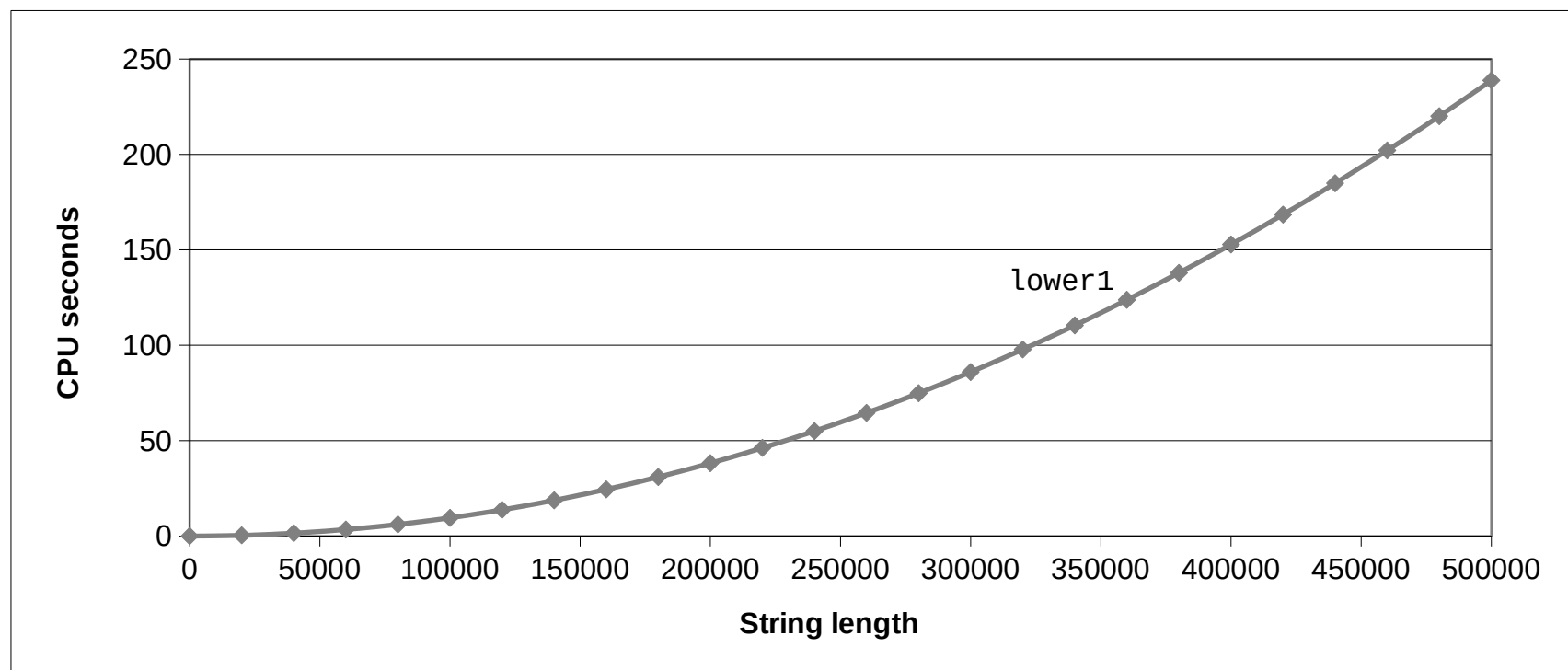
- **Procedure to Convert String to Lower Case**

```
void lower(char *s)
{
  size_t i;
  for (i = 0; i < strlen(s); i++)
    if (s[i] >= 'A' && s[i] <= 'Z')
      s[i] -= ('A' - 'a');
}
```

- Extracted from 213 lab submissions, Fall, 1998

# Lower Case Conversion Performance

- Time quadruples when double string length
- Quadratic performance

# Improving Performance

```
void lower(char *s)
{
  size_t i;
  size_t len = strlen(s);
  for (i = 0; i < len; i++)
    if (s[i] >= 'A' && s[i] <= 'Z')
      s[i] -= ('A' - 'a');
}
```
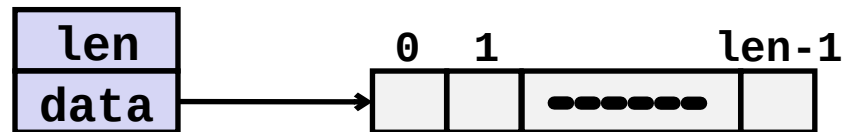
- Move call to `strlen` outside of loop
- Since result does not change from one iteration to another
- Form of code motion
- Compiler can't do it, unless code for strlen is inlined
- And perhaps not even then

# Exploiting Instruction-Level Parallelism

- **Need general understanding of modern processor design**
  - Hardware can execute multiple instructions in parallel
- **Performance limited by data dependencies**
- **Simple transformations can yield dramatic performance improvement**
  - Compilers often cannot make these transformations
  - Lack of associativity and distributivity in floating-point arithmetic

# Benchmark Example: Data Type for Vectors

```
/* data structure for vectors */
typedef struct{
    size_t len;
    data_t *data;
} vec;
```



- **Data Types**
  - Use different declarations for data_t
  - int
  - long
  - float
  - double

```
/* retrieve vector element
   and store at val */
int get_vec_element
  (*vec v, size_t idx, data_t *val)
{
    if (idx >= v->len)
        return 0;
    *val = v->data[idx];
    return 1;
}
```

# Benchmark Computation

```
void combine1(vec_ptr v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

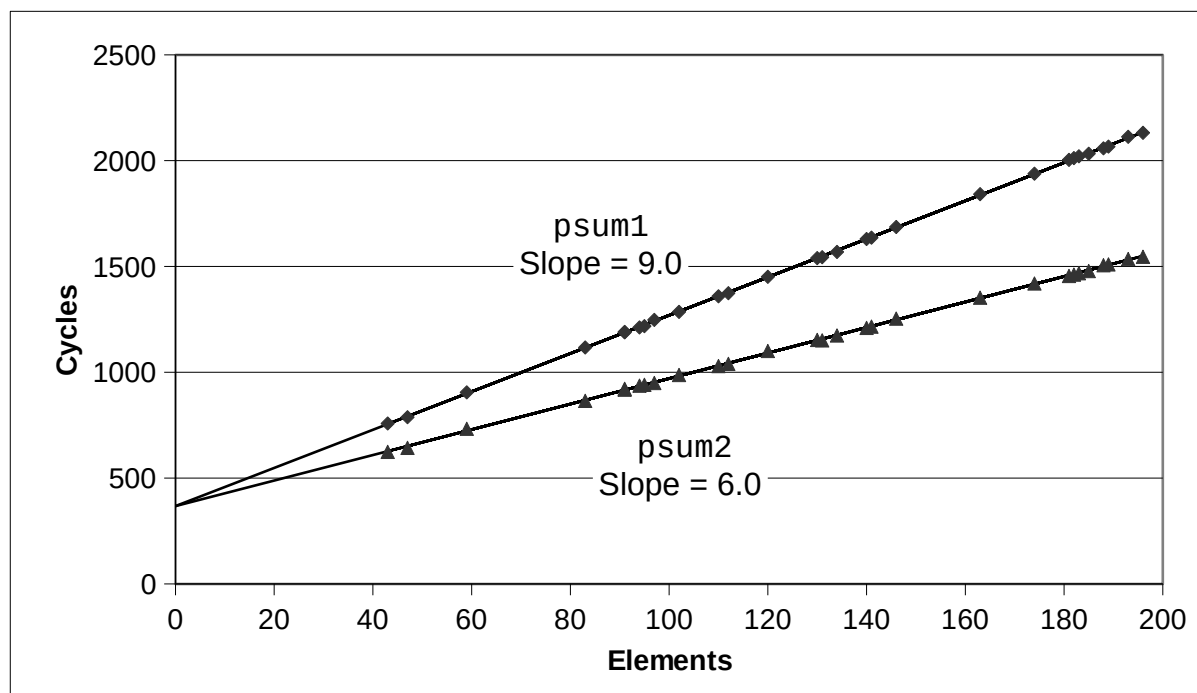**Compute sum or product of vector elements**

## ■Data Types
- Use different declarations for `data_t`
- `int`
- `long`
- `float`
- `double`

## ■Operations
- Use different definitions of `OP` and `IDENT`
- `+ / 0`
- `* / 1`

# Cycles Per Element (CPE)

- **Convenient way to express performance of program that operates on vectors or lists**
- **Length = n**
- **In our case: CPE = cycles per OP**
- **T = CPE*n + Overhead**
  - CPE is slope of line

# Benchmark Performance

```
void combine1(vec_ptr v, data_t *dest)
{
    long int i;
    *dest = IDENT;
    for (i = 0; i < vec_length(v); i++) {
        data_t val;
        get_vec_element(v, i, &val);
        *dest = *dest OP val;
    }
}
```

**Compute sum or product of vector elements**

| Method | Integer | | Double FP | |
|---|---|---|---|---|
| Operation | Add | Mult | Add | Mult |
| Combine1 unoptimized | 22.68 | 20.02 | 19.98 | 20.18 |
| Combine1 –O1 | 10.12 | 10.12 | 10.17 | 11.14 |

# Basic Optimizations

```
void combine4(vec_ptr v, data_t *dest)
{
  long i;
  long length = vec_length(v);
  data_t *d = get_vec_start(v);
  data_t t = IDENT;
  for (i = 0; i < length; i++)
    t = t OP d[i];
  *dest = t;
}
```

- **Move vec_length out of loop**
- **Avoid bounds check on each cycle**
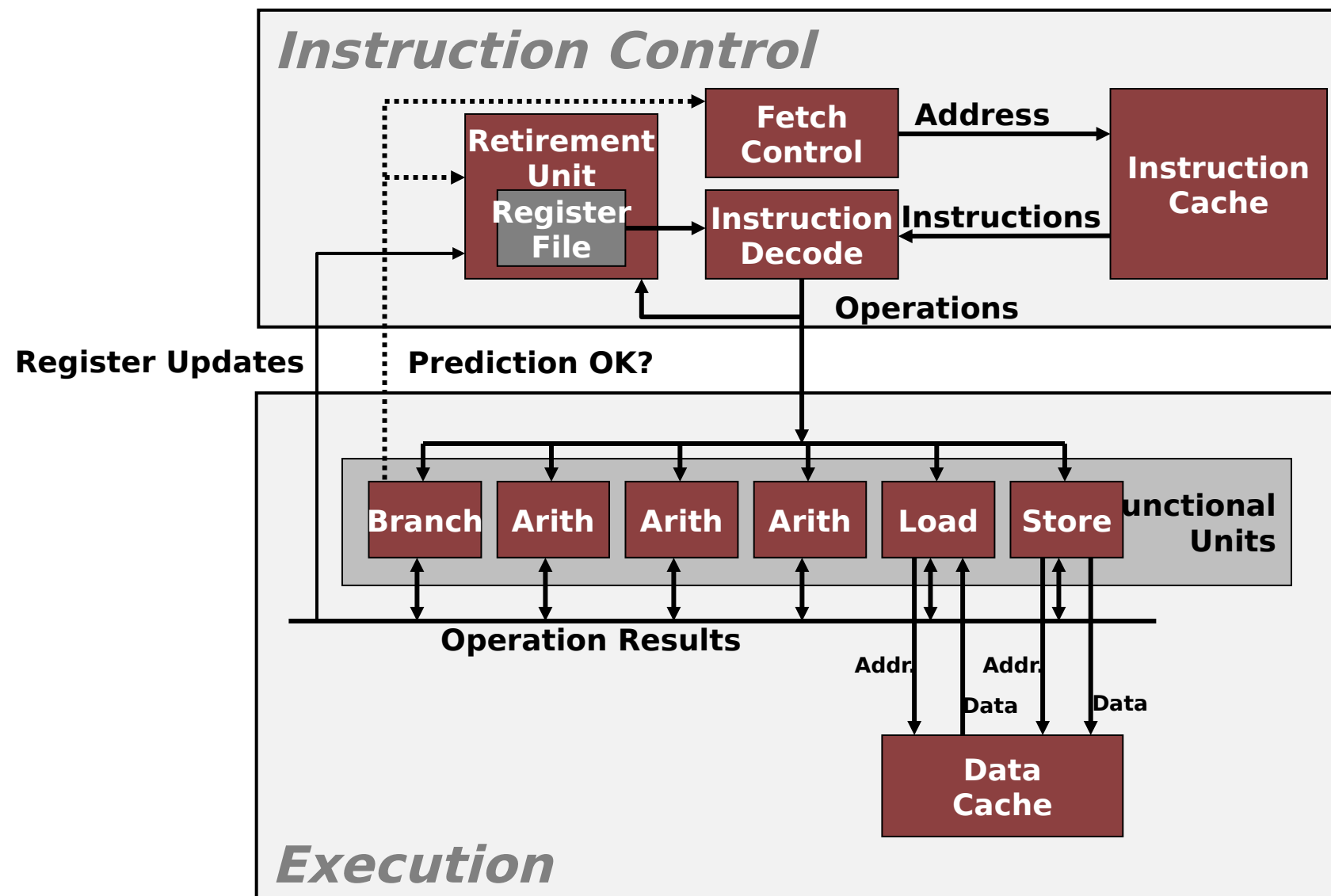- **Accumulate in temporary**

# Effect of Basic Optimizations

```
void combine4(vec_ptr v, data_t *dest)
{
  long i;
  long length = vec_length(v);
  data_t *d = get_vec_start(v);
  data_t t = IDENT;
  for (i = 0; i < length; i++)
    t = t OP d[i];
  *dest = t;
}
```

| Method | Integer | | Double FP | |
|---|---|---|---|---|
| Operation | Add | Mult | Add | Mult |
| Combine1 –O1 | 10.12 | 10.12 | 10.17 | 11.14 |
| Combine4 | 1.27 | 3.01 | 3.01 | 5.01 |

■ **Eliminates sources of overhead in loop**

# Modern CPU Design



**Instruction Control**

- Retirement Unit / Register File
- Fetch Control
- Instruction Decode
- Instruction Cache
- Address
- Instructions
- Operations

**Register Updates** | **Prediction OK?**

**Execution**

- Branch | Arith | Arith | Arith | Load | Store — Functional Units
- Operation Results
- Addr. | Addr.
- Data | Data
- Data Cache

# Superscalar Procesor

- **Definition: A superscalar processor can issue and execute *multiple instructions in one cycle*. The instructions are retrieved from a sequential instruction stream and are usually scheduled dynamically.**

- **Benefit: without programming effort, superscalar processor can take advantage of the *instruction level parallelism* that most programs have**

- **Most modern CPUs are superscalar.**
- **Intel: since Pentium (1993)**

# Haswell CPU

- ▪ 8 Total Functional Units
- **Multiple instructions can execute in parallel**
  2 load, with address computation
  1 store, with address computation
  4 integer
  2 FP multiply
  1 FP add
  1 FP divide
- **Some instructions take > 1 cycle, but can be pipelined**

| *Instruction* | *Latency* | *Cycles/Issue* |
|---|---|---|
| Load / Store | 4 | 1 |
| Integer Multiply | 3 | 1 |
| **Integer/Long Divide** | **3-30** | **3-30** |
| Single/Double FP Multiply | 5 | 1 |
| Single/Double FP Add | 3 | 1 |
| **Single/Double FP Divide** | **3-15** | **3-15** |

# x86-64 Compilation of Combine4

- ## Inner Loop (Case: Integer Multiply)

```
.L519:                             # Loop:
  movq   (%rax,%rdx,4), %r15  # t = t * d[i] (I)
  addq   $1, %rdx                 # i++
  imul   %r15, %ecx              # t = t * d[i] (II)
  cbg    %rbp,%rdx,.L519         # If lim>i, goto Loop
```

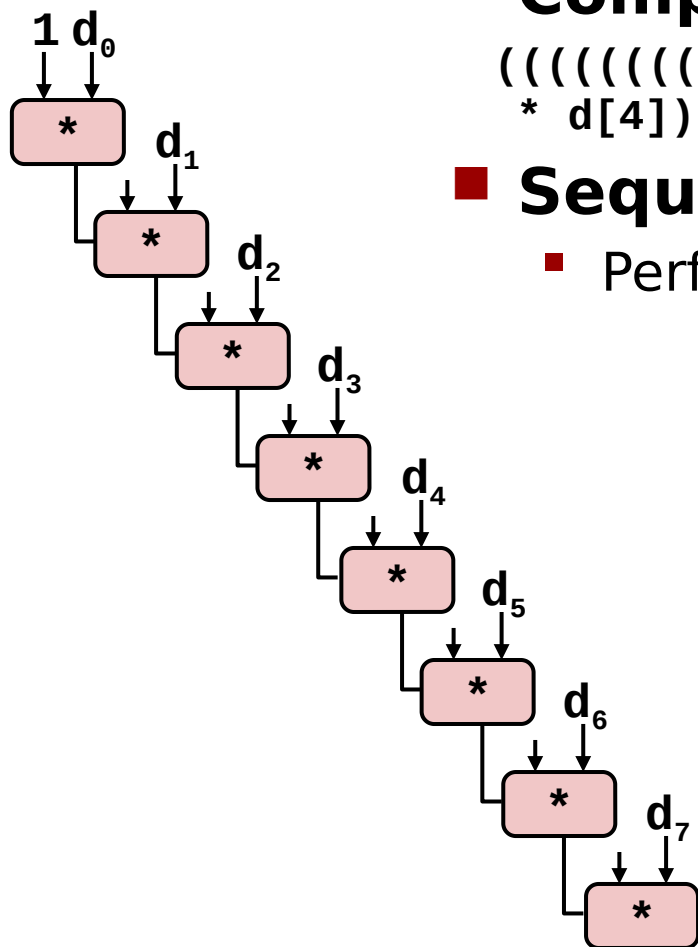| Method | Integer | | Double FP | |
|---|---|---|---|---|
| Operation | Add | Mult | Add | Mult |
| Combine4 | 1.27 | 3.01 | 3.01 | 5.01 |
| Latency Bound | 1.00 | 3.00 | 3.00 | 5.00 |

# Combine4 = Serial Computation (OP = *)



- **Computation (length=8)**
  `(((((((((1 * d[0]) * d[1]) * d[2]) * d[3])`
  `* d[4]) * d[5]) * d[6]) * d[7])`
- **Sequential dependence**
  - Performance: determined by latency of OP

# Loop Unrolling (2x1)

```
void unroll2a_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x = (x OP d[i]) OP d[i+1];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x = x OP d[i];
    }
    *dest = x;
}
```

- **Perform 2x more useful work per iteration**

# Effect of Loop Unrolling

| Method | Integer | | Double FP | |
|---|---|---|---|---|
| Operation | Add | Mult | Add | Mult |
| Combine4 | 1.27 | 3.01 | 3.01 | 5.01 |
| Unroll 2x1 | 1.01 | 3.01 | 3.01 | 5.01 |
| Latency Bound | 1.00 | 3.00 | 3.00 | 5.00 |

- **Helps integer add**
  - Achieves latency bound

```
x = (x OP d[i]) OP d[i+1];
```

- **Others don't improve.** *Why?*
  - Still sequential dependency

# Loop Unrolling with Reassociation (2x1a)

```
void unroll2aa_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x = x OP (d[i] OP d[i+1]);
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x = x OP d[i];
    }
    *dest = x;
}
```
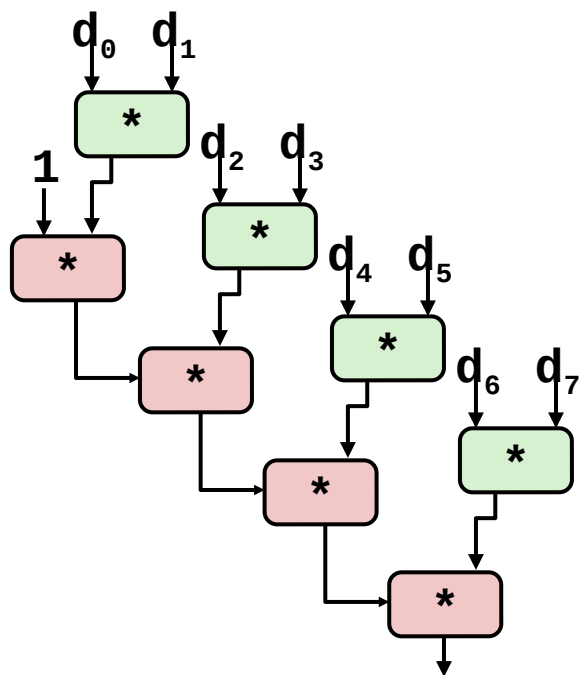
**Compare to before**

```
x = (x OP d[i]) OP d[i+1];
```

- **Can this change the result of the computation?**
- **Yes, for FP.** *Why?*

# Reassociated Computation

```
x = x OP (d[i] OP d[i+1]);
```



- **What changed:**
  - Ops in the next iteration can be started early (no dependency)

- **Overall Performance**
  - N elements, D cycles latency/op
  - (N/2+1)*D cycles: **CPE = D/2**

# Loop Unrolling with Separate Accumulators (2x2)

```
void unroll2a_combine(vec_ptr v, data_t *dest)
{
    long length = vec_length(v);
    long limit = length-1;
    data_t *d = get_vec_start(v);
    data_t x0 = IDENT;
    data_t x1 = IDENT;
    long i;
    /* Combine 2 elements at a time */
    for (i = 0; i < limit; i+=2) {
        x0 = x0 OP d[i];
        x1 = x1 OP d[i+1];
    }
    /* Finish any remaining elements */
    for (; i < length; i++) {
        x0 = x0 OP d[i];
    }
    *dest = x0 OP x1;
}
```

## ■ Different form of reassociation

# Effect of Separate Accumulators

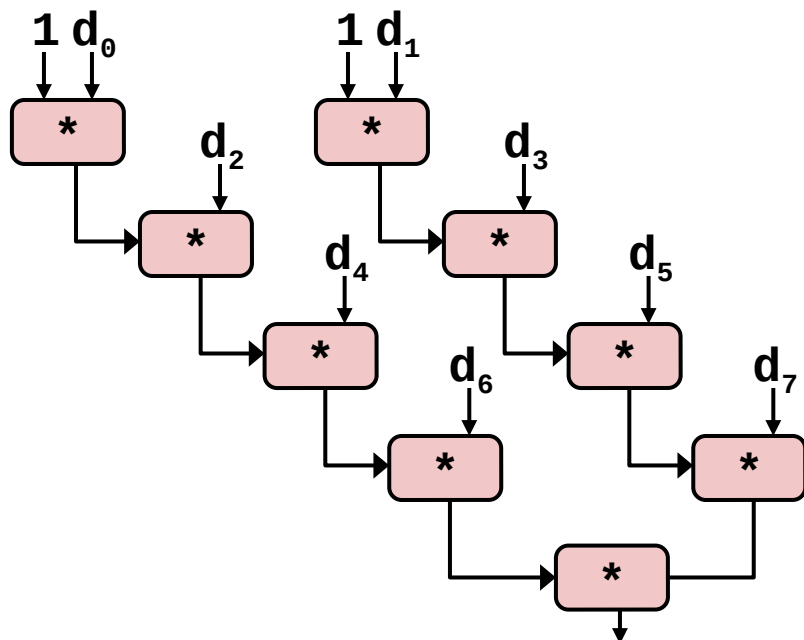| Method | Integer | | Double FP | |
|---|---|---|---|---|
| Operation | Add | Mult | Add | Mult |
| Combine4 | 1.27 | 3.01 | 3.01 | 5.01 |
| Unroll 2x1 | 1.01 | 3.01 | 3.01 | 5.01 |
| Unroll 2x1a | 1.01 | 1.51 | 1.51 | 2.51 |
| Unroll 2x2 | 0.81 | 1.51 | 1.51 | 2.51 |
| Latency Bound | 1.00 | 3.00 | 3.00 | 5.00 |
| Throughput Bound | 0.50 | 1.00 | 1.00 | 0.50 |

- **Int + makes use of two load units**

```
x0 = x0 OP d[i];
x1 = x1 OP d[i+1];
```

- **2x speedup (over unroll2) for Int *, FP +, FP ***

# Separate Accumulators

```
x0 = x0 OP d[i];
x1 = x1 OP d[i+1];
```

- **What changed:**
  - Two independent "streams" of operations

- **Overall Performance**
  - N elements, D cycles latency/op
  - Should be (N/2+1)*D cycles:
    **CPE = D/2**
  - CPE matches prediction!

  *What Now?*

# Unrolling & Accumulating

- **Idea**
  - Can unroll to any degree L
  - Can accumulate K results in parallel
  - L must be multiple of K

- **Limitations**
  - Diminishing returns
    - Cannot go beyond throughput limitations of execution units
  - Large overhead for short lengths
    - Finish off iterations sequentially

# Unrolling & Accumulating: Double *

- **Case**
  - Intel Haswell
  - Double FP Multiplication
  - Latency bound: 5.00.  Throughput bound: 0.50

| FP * | Unrolling Factor L | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **K** | 1 | 2 | 3 | 4 | 6 | 8 | 10 | 12 |
| 1 | 5.01 | 5.01 | 5.01 | 5.01 | 5.01 | 5.01 | 5.01 | |
| 2 | | 2.51 | | 2.51 | | 2.51 | | |
| 3 | | | 1.67 | | | | | |
| 4 | | | | 1.25 | | 1.26 | | |
| 6 | | | | | 0.84 | | | 0.88 |
| 8 | | | | | | 0.63 | | |
| 10 | | | | | | | 0.51 | |
| 12 | | | | | | | | 0.52 |

*Accumulators*

# Unrolling & Accumulating: Int +

- **Case**
  - Intel Haswell
  - Integer addition
  - Latency bound: 1.00.  Throughput bound: 0.50

| FP * | Unrolling Factor L | | | | | | | |
|------|------|------|------|------|------|------|------|------|
| **K** | 1 | 2 | 3 | 4 | 6 | 8 | 10 | 12 |
| 1 | 1.27 | 1.01 | 1.01 | 1.01 | 1.01 | 1.01 | 1.01 | |
| 2 | | 0.81 | | 0.69 | | 0.54 | | |
| 3 | | | 0.74 | | | | | |
| 4 | | | | 0.69 | | 1.24 | | |
| 6 | | | | | 0.56 | | | 0.56 |
| 8 | | | | | | 0.54 | | |
| 10 | | | | | | | 0.54 | |
| 12 | | | | | | | | 0.56 |

*Accumulators*

# Achievable Performance

| Method | Integer | | Double FP | |
|---|---|---|---|---|
| Operation | Add | Mult | Add | Mult |
| Best | 0.54 | 1.01 | 1.01 | 0.52 |
| Latency Bound | 1.00 | 3.00 | 3.00 | 5.00 |
| Throughput Bound | 0.50 | 1.00 | 1.00 | 0.50 |

- **Limited only by throughput of functional units**
- **Up to 42X improvement over original, unoptimized code**

# What About Branches?

- **Challenge**
  - Instruction Control Unit must work well ahead of Execution Unit

    to generate enough operations to keep EU busy

```
404663:   mov     $0x0,%rax          ⎫
404668:   mov     (%rdi),%r15        ⎬  Executing
40466b:   cbge    %r15,%rsi,404685   ⎭
40466d:   mov     0x8(%rdi),%rax        How to continue?


  . . .


404685:   something :-)
```
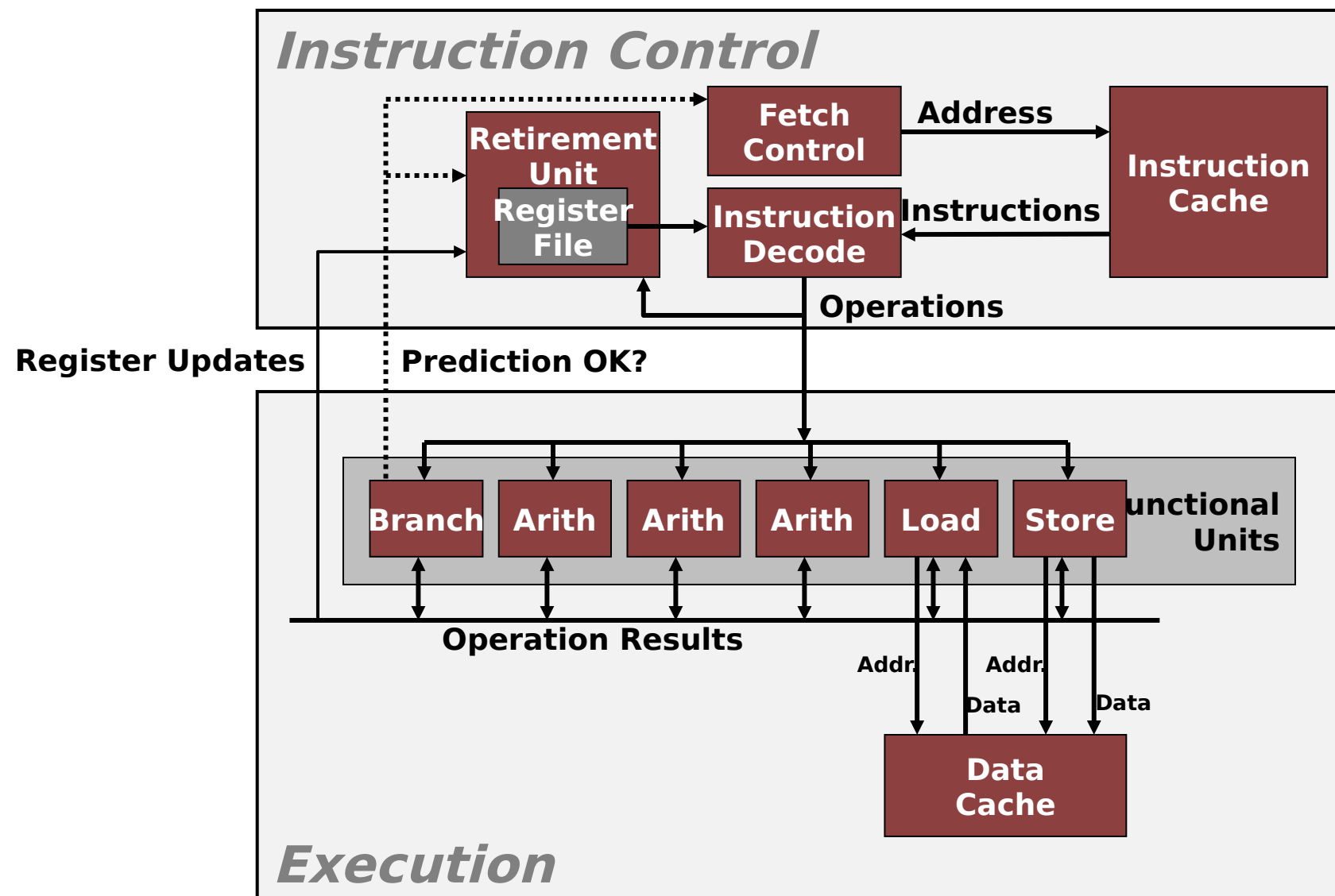
  - When encounters conditional branch, cannot reliably determine where to continue fetching

# Modern CPU Design

# Branch Outcomes

▪ **When encounter conditional branch, cannot determine where to continue fetching**

  ▪ Branch Taken: Transfer control to branch target
  ▪ Branch Not-Taken: Continue with next instruction in sequence

▪ **Cannot resolve until outcome determined by branch/integer unit**

```
404663:   mov      $0x0,%eax
404668:   mov      (%rdi),%r15
40466b:   cbge     %r15,%rsi,404685
40466d:   mov      0x8(%rdi),%rax

 . . .

404685:   something else
```

**Branch Not-Taken**

**Branch Taken**

# Branch Prediction

- **Idea**
  - Guess which way branch will go
  - Begin executing instructions at predicted position
    - But don't actually modify register or memory data

```
404663:   mov     $0x0,%eax
404668:   mov     (%rdi),%r15
40466b:   cbge    %r15,%rsi,404685
40466d:   mov     0x8(%rdi),%rax

 . . .

404685:   something else
```

**Predict Taken**

**Begin Execution**

# Branch Prediction Through Loop

*sorry – x86 not prime :-(*

```
401029:   vmulsd   (%rdx),%xmm0,%xmm0
40102d:   add      $0x8,%rdx
401031:   cmp      %rax,%rdx
401034:   jne      401029
```
*i = 98*

*Assume
vector length = 100*

**Predict Taken (OK)**

```
401029:   vmulsd   (%rdx),%xmm0,%xmm0
40102d:   add      $0x8,%rdx
401031:   cmp      %rax,%rdx
401034:   jne      401029
```
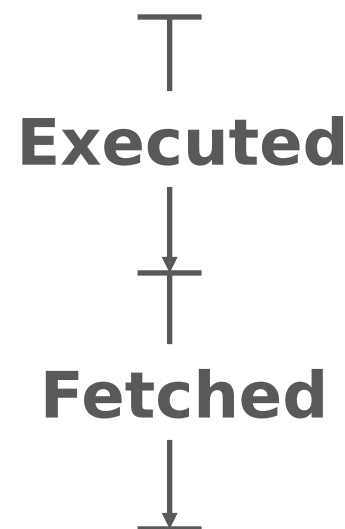*i = 99*

**Predict Taken
(Oops)**

```
401029:   vmulsd   (%rdx),%xmm0,%xmm0
40102d:   add      $0x8,%rdx
401031:   cmp      %rax,%rdx
401034:   jne      401029
```
*i = 100*

**Read
invalid
location**

**Executed**

```
401029:   vmulsd   (%rdx),%xmm0,%xmm0
40102d:   add      $0x8,%rdx
401031:   cmp      %rax,%rdx
401034:   jne      401029
```
*i = 101*

**Fetched**

# Branch Misprediction Invalidation

```
401029:   vmulsd (%rdx),%xmm0,%xmm0
40102d:   add    $0x8,%rdx
401031:   cmp    %rax,%rdx
401034:   jne    401029
```
*i = 98*

*Assume vector length = 100*

**Predict Taken (OK)**

```
401029:   vmulsd (%rdx),%xmm0,%xmm0
40102d:   add    $0x8,%rdx
401031:   cmp    %rax,%rdx
401034:   jne    401029
```
*i = 99*

**Predict Taken (Oops)**

```
401029:   vmulsd (%rdx),%xmm0,%xmm0
40102d:   add    $0x8,%rdx
401031:   cmp    %rax,%rdx
401034:   jne    401029
```
*i = 100*

**Invalidate**

```
401029:   vmulsd (%rdx),%xmm0,%xmm0
40102d:   add    $0x8,%rdx
401031:   cmp    %rax,%rdx
401034:   jne    401029
```
*i = 101*

# Branch Misprediction Recovery

```
401029:   vmulsd (%rdx),%xmm0,%xmm0
40102d:   add     $0x8,%rdx
401031:   cmp     %rax,%rdx
401034:   jne     401029
401036:   jmp     401040
 . . .
401040:   vmovsd %xmm0,(%r12)
```

*i = 99*

**Definitely not taken**

**Reload Pipeline**

## Performance Cost
- Multiple clock cycles on modern processor
- Can be a major performance limiter

# Getting High Performance

- **Good compiler and flags**
- **Don't do anything stupid**
  - Watch out for hidden algorithmic inefficiencies
  - Write compiler-friendly code
    - Watch out for optimization blockers: procedure calls & memory references
  - Look carefully at innermost loops (where most work is done)

- **Tune code for machine**
  - Exploit instruction-level parallelism
  - Avoid unpredictable branches
  - Make code cache friendly (Covered earlier in course)