Compilers - Chapter 8: Loop scheduling optimisations Part 1: Why mess with the order of loop execution?

- Lecturer:
  - Paul Kelly (p.kelly@imperial.ac.uk)

PAGE 3			
DEPARTMENT	COURSE	DESCRIPTION	PREREQS
COMPUTER SCIENCE	CPSC 432	INTERMEDIATE COMPILER DESIGN, WITH A FOCUS ON DEPENDENCY RESOLUTION.	CPSC 432
0	Ocean Lumm	interior colucito projeti	OUT IN

(It turns out that this cartoon is good for almost any compilers topic, definitely this one) https://xkcd.com/754/

# "Restructuring" compilers

- The optimisations we have studied so far reduce the number of instructions that need to be executed at runtime
  - This is fundamentally a good idea!
- But sometimes we can get a performance improvement by thinking about the order in which loops are executed
- Why might that be?

# "Restructuring" compilers

- The optimisations we have studied so far reduce the number of instructions that need to be executed at runtime
  - This is fundamentally a good idea!
- But sometimes we can get a performance improvement by thinking about the order in which loops are executed
- Why might that be?
  - We might be able to use vector instructions
    - So different iterations of a loop are being executed at the same time

### "Restructuring" and "parallelizing" compilers

- The optimisations we have studied so far reduce the number of instructions that need to be executed at runtime
  - This is fundamentally a good idea!
- But sometimes we can get a performance improvement by thinking about the order in which loops are executed
- Why might that be?
  - We might be able to use vector instructions
    - So different iterations of a loop are being executed at the same time
  - We might be able to use multiple cores
    - So different iterations of a loop might be assigned to different threads running on different CPUs

### "Restructuring" and "parallelizing" compilers

- The optimisations we have studied so far reduce the number of instructions that need to be executed at runtime
  - This is fundamentally a good idea!
- But sometimes we can get a performance improvement by thinking about the order in which loops are executed
- Why might that be?
  - We might be able to use vector instructions
    - So different iterations of a loop are being executed at the same time
  - We might be able to use multiple cores
    - So different iterations of a loop might be assigned to different threads running on different CPUs
  - We might be able to improve how the cache is used
    - We will come to this later!

### Vector instruction set extensions

- Example: Intel's AVX512
- Extended registers ZMM0-ZMM31, 512 bits wide
  - Can be used to store 8 doubles, 16 floats, 32 shorts,
    64 bytes
  - So instructions are executed in parallel in 64,32,16 or 8 "lanes"

### Vector instruction set extensions

- Example: Intel's AVX512
- Extended registers ZMM0-ZMM31, 512 bits wide
  - Can be used to store 8 doubles, 16 floats, 32 shorts,
     64 bytes
- Example: vaddps zmm0 zmm1 zmm2

- "Add Packed Single Precision Floating-Point Values"



In *one instruction* we add 16 32-bit floating point values from zmm1 and 16 32-bit values from zmm2







### Can we get the compiler to vectorise?

C++ source #1 🖉	x86-64 gcc 14.2 (Editor #1) 🖉 🗙
A• 🖬 +• 1⁄ ,≅ 📌	x86-64 gcc 14.2 🔹 🗹 📀 -march=znver4 -Ofast -fopt-infc
Get + ▲	A ▼ ✿ Output ▼ Filter ▼ 🛢 Libraries 🖌 Overrides + Add new ▼ 🖌
<pre>1 float c[1024]; 2 float a[1024]; 3 float b[1024]; 4 5 void addabc() { 6 for (int i=0; i&lt;1024; i++) 7 c[i] = a[i] + b[i]; 8 } 9</pre>	1       addabc():         2       xor       eax, eax         3       .L2:         4       vmovaps zmm0, ZMMWORD PTR <u>a[rax]</u> 5       vaddps zmm0, zmm0, ZMMWORD PTR <u>b[rax]</u> 6       add         7       vmovaps ZMMWORD PTR <u>c[rax-64]</u> , zmm0         8       cmp rax, 4096         9       jne <u>.L2</u> 10       vzeroupper         11       ret         12       b:         13       .zero 4096         14       a:         15       .zero 4096
	17 .zero 4096

In sufficiently simple cases, no problem: Gcc reports: addcba.c:6:20: optimized: loop vectorized using 64 byte vectors

#### Can we get the compiler to vectorise?

Tell the compiler to generate code for AMD Zen 4 which has AVX512



In sufficiently simple cases, no problem: Gcc reports: addcba.c:6:20: optimized: loop vectorized using 64 byte vectors



Iteration #0

Iteration #1

Iteration #2

### How do we know a loop is parallel?

- To use vector instructions, we need to verify that different iterations of the loop are truly parallel
- In this case we can easily see that the dependence arrows do not cross iteration boundaries



This case was easy: for (int i=0; i<4; i++)

c[i] = a[i] + b[i];

• Source code:

#### How much does it help? First: *without* vectorisation

- Processor: AMD Ryzen 9 7940HS ("maple10")
- Compiler command line:

c[i] = a[i] + b[i];

for (int i=0; i<size; i++)</pre>

gcc -01 addcba-perf.c

• Generated code – not vectorised:

.L3:		
	movss	(%rsi,%rax), %xmmO
	addss	(%rcx,%rax), %xmm0
	movss	<pre>%xmm0, (%rdi,%rax)</pre>
	addq	\$4, %rax
	cmpq	%rdx, %rax
	jne	.13

- Performance: 4.8 GFLOPS (4.8\*10<sup>9</sup> single precision floating-point operations/second)
- Time per loop iteration: 0.21ns (one clock cycle at 4.8GHz, 1 result per iteration)

• Source code:

### How much does it help?

- This time with vectorisation
- Processor: AMD Ryzen 9 7940HS ("maple10")
- Compiler command line:

gcc -Ofast -march=znver4 addcba-perf.c

Generated code:

.L4:		
	vmovaps	(%r8,%rax), %zmm1
	vaddps	(%rdi,%rax), %zmm1, %zmm0
	vmovaps	<pre>%zmm0, (%rsi,%rax)</pre>
	addq	\$64, %rax
	cmpq	<pre>%rax, %rdx</pre>
	jne	. L4

- Performance: 34.8 GFLOPS (single precision)
- Time per loop iteration: 0.45ns (two clock cycles, 16 results per iteration)

• Source code:

## How much does it help? This time *with* vectorisation

- Processor: AMD Ryzen 9 7940HS ("maple10")
- Compiler command line:

gcc -Ofast -march=znver4 addcba-perf.c

Generated code:

.L4:		
	vmovaps	(%r8,%rax), %zmm1
	vaddps	(%rdi,%rax), %zmm1, %zmm0
	vmovaps	%zmm0, (%rsi,%rax)
	addq	\$64, %rax
	cmpq	%rax, %rdx
	jne	. L4

#### Speed of light is 30cm/ns.

So this machine completes about three iterations in the time it takes the light to get from your computer screen to your eyes

- Performance: 34.8 GFLOPS (single precision)
- Time per loop iteration: 0.45ns (two clock cycles, 16 results per iteration)

Compilers - Chapter 8: Loop scheduling optimisations Part 2: Determining whether a loop can be executed in parallel

- Lecturer:
  - Paul Kelly (p.kelly@imperial.ac.uk)

PAGE 3			
DEPARTMENT	COURSE	DESCRIPTION	PREREQS
COMPUTER SCIENCE	CPSC 432	INTERMEDIATE COMPILER DESIGN, WITH A FOCUS ON DEPENDENCY RESOLUTION.	CPSC 432
6	00001	COMPANY COMPANY DECK	01.171.101

(It turns out that this cartoon is good for almost any compilers topic, definitely this one) <u>https://xkcd.com/754/</u> 24

#### June 25

### But that example was obviously parallel?

Ρ

- To use vector instructions, we need to verify that different iterations of the loop are truly parallel
- This case was easy:

for (int i=0; i<1024; i++) c[i] = a[i] + b[i];

• How about this one?

for (int i=0; i<1024; i++) Q c[i] = c[i-1] + b[i];



# But that example was obviously parallel?

- To use vector instructions, we need to verify that different iterations of the loop are truly parallel
- This case was easy:

for (int i=0; i<1024; i++) **P** c[i] = a[i] + b[i];

• How about this one?

for (int i=0; i<1024; i++) Q c[i] = c[i-1] + b[i];

• And this?

for (int i=0; i<1024; i+=2) **R** c[i] = c[i-1] + b[i];

• Consider this example:

for (int i=1; i<8; i++) Q c[i] = c[i-1] + b[i];

• What does it *do*?

• Consider this example:

for (int i=1; i<8; i++) Q c[i] = c[i-1] + b[i];

• When executed we get:

$$c[1] = c[0] + b[1];$$
  

$$c[2] = c[1] + b[2];$$
  

$$c[3] = c[2] + b[3];$$
  

$$c[4] = c[3] + b[4];$$
  

$$c[5] = c[4] + b[5];$$
  

$$c[6] = c[5] + b[6];$$
  

$$c[7] = c[6] + b[7];$$

Consider this example:
 for (int i=1; i<8; i++)</li>

Each iteration produces a value that is used in the next iteration

• When executed we get:

c[i] = c[i-1] + b[i];

$$c[1] = c[0] + b[1];$$
  

$$c[2] = c[1] + b[2];$$
  

$$c[3] = c[2] + b[3];$$
  

$$c[4] = c[3] + b[4];$$
  

$$c[5] = c[4] + b[5];$$
  

$$c[6] = c[5] + b[6];$$
  

$$c[7] = c[6] + b[7];$$

The dependence arrows go from one iteration to the next

The dependence is *carried* by the loop

- Consider this example:
   for (int)i=1; i<8; i++)</li>
   c[i] = c[i-1] + b[i];
- When executed we get:

$$c[1] = c[0] + b[1];$$
  

$$c[2] = c[1] + b[2];$$
  

$$c[3] = c[2] + b[3];$$
  

$$c[4] = c[3] + b[4];$$
  

$$c[5] = c[4] + b[5];$$
  

$$c[6] = c[5] + b[6];$$
  

$$c[7] = c[6] + b[7];$$



#### June 25

# Iteration #3

#### 31

## "Loop-carried dependence"

- Consider this example: for (int i=1; i<8; i++) c[i] = c[i-1] + b[i];
- When executed we get:

$$c[1] = c[0] + b[1];$$
  

$$c[2] = c[1] + b[2];$$
  

$$c[3] = c[2] + b[3];$$
  

$$c[4] = c[3] + b[4];$$
  

$$c[5] = c[4] + b[5];$$
  

$$c[6] = c[5] + b[6];$$
  

$$c[7] = c[6] + b[7];$$



- Consider this example:
   for (int)i=1; i<8; i++)</li>
   c[i] = c[i-1] + b[i];
- When executed we get:

$$c[1] = c[0] + b[1];$$
  

$$c[2] = c[1] + b[2];$$
  

$$c[3] = c[2] + b[3];$$
  

$$c[4] = c[3] + b[4];$$
  

$$c[5] = c[4] + b[5];$$
  

$$c[6] = c[5] + b[6];$$
  

$$c[7] = c[6] + b[7];$$



# So we need a compiler algorithm

- To determine whether there is a loop-carried dependence
- To distinguish, for example, P, Q and R:

- No loop-carried dependence
- So iterations can be executed in parallel
- So vectorisable
- loop-carried dependence
- So iterations cannot be executed in parallel
- So not vectorisable
  - No loop-carried dependence
- So iterations can be executed in parallel
- So vectorisable

# So we need a compiler algorithm

- To determine whether there is a loop-carried dependence
- To distinguish, for example, P, Q and R:

- No loop-carried dependence
- So iterations can be executed in parallel
- So vectorisable
- loop-carried dependence
- So iterations cannot be executed in parallel
- So not vectorisable
- No loop-carried dependence
- So iterations can be executed in parallel
- So vectorisable
- (though actually generating efficient vector code for this might be a bit tricky?)



Compiler Explorer Editor Diff View More-C++ source #1 X A- Add new...-1 void add(float \*\_\_restrict\_\_ c, float \* restrict a, 2 3 float \* restrict b, 4 int N) 5 for (int i=0; i <= N; i++)</pre> 6 c[i]=a[i]+b[i]; 7 8

#### If the alignment of the operand pointers is not known:

gcc reports: test.c:6:3: note: loop vectorized test.c:6:3: note: loop peeled for vectorization to enhance alignment test.c:6:3: note: loop turned into non-loop; it never loops. test.c:6:3: note: loop with 3 iterations completely unrolled test.c:1:6: note: loop turned into non-loop; it never loops. test.c:1:6: note: loop with 4 iterations completely unrolled

x86-64 gcc 5.4 (Editor #1, Compiler #1) C++ X x86-64 acc 5.4 -O3 -fopt-info // \s+ Intel Demangle 11010 test1 Necx, Nec pushq %r13 pushq %r12 pushq Krbp pushq Krbp movo %rsi, %rax leal andl 1(%rcx), %r90 \$15, %eax shrq \$2, %rax negq andl cmpl \$3, %ear %r9d, %ea %r9d, %ear cmova cmpl \$4, %r9d %r9d, %ea (%rsi), %om Three copies of the noncmpl \$1, %eax mov1 \$1, %r8d addss (%rdx), %xmm0 movss %xmm0, (%rdi) je .L5 movss 4(%rsi), %ommi cmpl \$2, %eax movl \$2, %r8d vectorised loop body to addss 4(%rdx), %onni novss %xmm0, 4(%rdi) je .L5 align the start address of movss 8(%rsi), %xmmi cmpl \$3, %eax \$3, %r8d 8(%rdx), %xmm the vectorised code on a movss %xmm0, 8(%rdi) je .L5 movss 12(%rsi), %xmm0 movl \$4, %r8d addss 12(%rdx), %xmm0 32-byte boundary movss %xmm0, 12(%rdi movl %ecx, %ebx movl %eax, %r11d -4(%r9), %r18d leal subl shrl \$2, %rl0d addl \$1, %rl0d cmpl leal \$2, %ebx 0(,%r10,4), %ebp jbe .L7 leaq 0(,%r11,4), %rax xorl %ebx, %ebx leaq (%rsi,%rax), %r1: (%rsi,%rax), %r13 leag (%rdx,%rax), %r1 leag (%rdi,%rax), %r11 movups (%r12,%rax), %xmm addl Basically the same \$1, %ebx 0(%r13,%rax), %xmm0 %xmm0, (%r11,%rax) \$16, %rai %ebx, %r10d vectorised code as before .L9 %ebp, %r8d %ebp, %r9d .L1 movss (%rsi,%rax,4), %xmm8 addss (%rdx,%rax,4), %xmm8 Three copies of the nonmovss %xmm0, (%rdi,%rax,4 leal cmpl 1(%r8), %eas Neax, Necx .L1 vectorised loop body to cltq Thhe \$2, %r8d movss (%rsi,%rax,4), %xmm cmol. ScRd. Secu addss (%rdx,%rax,4), %xmm8 mop up the additional movss %xmm0, (%rdi,%rax,4) jl .L1 movslq %r8d, %r8 movss (%rsi,%r8,4), %onm8
addss (%rdx,%r8,4), %onm8 iterations in case N is not movss %xmm0, (%rdi,%r8,4) popq %rbx popq %rbp popq %r12 popq %r13 divisible by 4 net. testl %eax, %eax jne .L3 %r8d, %r8d xorl 1mp

101 102

A-

C++

Share - Other -

		0.	L	
	ia 🗸	сл		٦

Compiler Explorer Editor Diff View More-

CTT SU	urce #1	^			
Δ_	B Sav	heo l/a	<b>- - - -</b>	Add now	

Save/Load	+ Add new▼
1	<pre>void add(float *c,</pre>
2	float *a,
3	<pre>float *b,</pre>
4	int N)
5	{
6	<pre>for (int i=0; i &lt;= N; i++)</pre>
7	c[i]=a[i]+b[i];
8	}

#### If the pointers might be aliases:

#### gcc reports:

test.c:6:3: note: loop vectorized

test.c:6:3: note: loop versioned for vectorization because of possible aliasing

test.c:6:3: note: loop peeled for vectorization to enhance alignment test.c:6:3: note: loop turned into non-loop; it never loops. test.c:6:3: note: loop with 3 iterations completely unrolled test.c:1:6: note: loop turned into non-loop; it never loops. test.c:1:6: note: loop with 3 iterations completely unrolled

Check whether the memory regions pointed to by c, b and a might overlap

Three copies of the nonvectorised loop body to align the start address of the vectorised code on a 32-byte boundary

Basically the same vectorised code as before

Three copies of the nonvectorised loop body to mop up the additional iterations in case N is not divisible by 4

Non-vector version of the loop for the case when c might overlap with a or b

C++ • x86-64 gcc 5 4 (Editor #1, Compiler #1) C++ X

test1 leck, leck js .127 krew is(licit), True is(licit), True (lest i(licit), True i(licit), Srid cmpq Kras, Krdi setho Nai orl Mean, Krdi setho Nai orl Kris, Krdi setho Nai orl Kris, Krdi setho Nai

Xriid, Xea testb %al, %r8b

> \$8, %r9d Xrsi, Xra Xr13 \$15, %ea

 point
 2-10

 point
 2-10

je .L4 novss 8(Nrsi), Xxmmi novl \$3, Xr8d addss 8(Nrdx), Xxmmi

movss %xmm0, 8(%rdi)

Reen, Xr8d S2, Xrax Xr18d, Xr18d -4(Dr9), Hr11d (Nr51,Xrax), Nr1 Nebx, Nebx Nrd1, Nrax S2, Xr11d S1, Xr11d 0(,Nr11,4), Nebp leaq leaq xorl addq shrl

movups (%r12,%r10), %xm addl \$1, %ebx addps 0(%r13,%r10), %xm

movups (xr12,xr18), xxmme addps qd(xr13,Xr18), Xxmme movups Xxmm8, (Xr13,Xr18), Xxmme addg \$16, Xr18 cmp1 Xr11d, Xebx jb .L7

cmpl kebp, kred
je .11
movslq Nrdd, Nrax
movsds (krsi,Nrax,4), Nomed
movss (krsi,Nrax,4), Nomed
movss Nomed, (krdi,Nrax,4)
leal 1(krs), Neax
cmpl Keax, Kecx
jl .11

\$7. 102 (%rsi,%rax,4), %xmm8 %r8d, %ecx addss (%rdx,%rax,4),%xmm novss %xmm0, (%rdi,%rax,4 1 .L1 Novslq Xrad, Xrax

(Mrsi,Mrax,4), Norme (Mrdx,Mrax,4), Norme www. Symp. (Sedi Seav

.L7 addl Xebp, Xr8d cmpl Xebp, Xr8d je .L3

cltq

popq %rbx popq %rbp popq %r12 popq %r13

xorl Xeax, Xe

moves (Xrsi,Xrax,4), Xxmm0 addss (Krot,Xrax,4), Xxmm0 moves Xxmm0, (Krdi,Xrax,4) addg S1, Xrax cmp1 Reax, Xecx S12 -L12

rep ret

-O3 -fopt-info

\s+ Intel Demangle ■ Libraries + Add new...-

x86-64 gcc 5.4

A-11010

# What do we see?

 Actually exploiting vectorisation is a bit tricky even when the dependence analysis is easy

- In the following slides we start with an easily-vectorizable example
- And look at some of the things that make it complicated



Compiler returned: 0



C https://godbolt.org G 🔍 Startpage Search En... 📭 Home - BBC News 🛛 I Staff travel and exp... 🔽 Shareable Whiteboa... 🔘 The Add... ▼ More ▼ Templates EXPLORER C++ source #1 Ø € Ŕ 🕝 C++ 12 // icx -Ofast -march=znver4 -qopt-report #define SIZE 10240 #define ALIGN \_\_attribute\_\_ ((aligned(64))) 3 4 int ALIGN A[SIZE]; 5 int ALIGN ind[SIZE]; 6 int ALIGN C[SIZE]; 7 int ALIGN D[SIZE]; 8 9 10 #define IB 32 #define JB 32 11 12 void P() 13 14 int i, j; 15 16 17 for (i=0; i<SIZE; i++) {</pre> C[ind[i]] += A[ind[i]] - D[i]; 18 19 20

#### **Incrementing through indirection: ind[i]**

- 1. Load a vector ind[i:i+16]
- 2. Gather a vector A[ind[i:i+16]
- 3. Subtract the D[i] values:
- 4. RHS[0:16]=A[ind[i:i+16]] D[i:i+16]
- 5. Gather the LHS[0:16] = C[ind[i:i+16]]
- 6. Add (+=): LHS[0:16] += RHS[0:16]
- 7. Scatter: C[ind[i:i+16]] = LHS[0:16]





#### **Incrementing through indirection: ind[i]**

- 1. Load a vector ind[i:i+16]
- 2. Gather a vector A[ind[i:i+16]
- 3. Subtract the D[i] values:
- 4. RHS[0:16]=A[ind[i:i+16]] D[i:i+16]
- 5. Gather the LHS[0:16] = C[ind[i:i+16]]
- 6. Add (+=): LHS[0:16] += RHS[0:16]
- 7. Scatter: C[ind[i:i+16]] = LHS[0:16]

What would happen if there were duplicate indices in ind?





#### **Incrementing through indirection: ind[i]**

- 1. Load a vector ind[i:i+16]
- 2. Gather a vector A[ind[i:i+16]
- 3. Subtract the D[i] values:
- 4. RHS[0:16]=A[ind[i:i+16]] D[i:i+16]
- 5. Gather the LHS[0:16] = C[ind[i:i+16]]
- 6. Add (+=): LHS[0:16] += RHS[0:16]
- 7. Scatter: C[ind[i:i+16]] = LHS[0:16]

What would happen if there were duplicate indices in ind?



It's not parallel! We have to sum two (or more) different values into the same C element


#### Incrementing through indirection: ind[i]

- 1. Load a vector ind[i:i+16]
- 2. Gather a vector A[ind[i:i+16]
- 3. Subtract the D[i] values:
- 4. RHS[0:16]=A[ind[i:i+16]] D[i:i+16]
- 5. Gather the LHS[0:16] = C[ind[i:i+16]]
- 6. Add (+=): LHS[0:16] += RHS[0:16]
- 7. Scatter: C[ind[i:i+16]] = LHS[0:16]



71

72

73

74

75

· · · · · vpcmpnead ·

vzeroupper ret

LBB0\_8:

•••••vptest••ymm6,•ymm2

jmp ..... LBB0 7

•••••• **k1**, •ymm6, •ymm2

vpconflictq instruction checks for duplicate values in ind[i:i+16]

If found, we branch to a loop over each distinct value

Roughly...

This is addressed by AVX512 "conflict detect" instructions which enable us to catch duplicates and serialise where needed

# **Health warning**

- Automatic discovery of parallelism has a bad reputation
  - Deservedly! It looks great on simple examples
  - But real code has complexity that means it often just doesn't happen
- But in some application domains it can really work
- And some programming languages make it easier, maybe!
  - Functional languages lack anti- and output-dependences (but tend to add higher-order functions and lazy evaluation)
  - Some languages control pointer ownership and aliasing
  - Some programming models discourage explicit loops and explicit elementwise subscripting

# So: we need a compiler algorithm to determine whether a loop is parallel...

47

How?

Dependence

M Define:

IN(S): set of memory locns which might be read by some execn of statement S

OUT(S): set of memory locns which might be written by some execn of statement S

Reordering is constrained by dependences; ("S1 must write something before S2 can read it") Mare are four types: Data ("true") dependence: S1 δ S2 ("S1 must read something • OUT(S1)  $\cap$  IN(S2) before S2 overwrites it")  $\bullet$ Anti dependence: S1<sup> $\overline{\delta}$ </sup> S2 ("If S1 and S2 might both •  $IN(S1) \cap OUT(S2)$ write to a location, S2 must Output dependence: S1 δ° S2 write after S1") •  $OUT(S1) \cap OUT(S2)$ **Control dependence:** S1  $\delta^{c}$  S2 ("S1 determines whether S2 should execute")

These are static analogues of the dynamic RAW, WAR, WAW and control hazards which have to be considered in processor architecture

#### Recall:

### **Loop-carried dependences**

S1 : A[0] := 0 for I = 1 to 8 S2 : A[I] := A[I-1] + B[I]

#### What does this loop do?



#### Recall:

## **Loop-carried dependences**

S1 : A[0] := 0 for I = 1 to 8 S2 : A[I] := A[I-1] + B[I]

What does this loop do?





In this case, there is a data dependence

- This is a loop-carried dependence the dependence spans a loop iteration
- This loop is inherently sequential

#### **Loop-carried dependences**

S1 : A[0] := 0 for I = 1 to 8 S2 : A[I] := A[I-1] + B[I]

Loop carried:

**Recall**:

- S2<sup>1</sup> : A[1] := A[0] + B[1]
- S2<sup>2</sup> : A[2] := A[1] + B[2]
- S2<sup>3</sup>: A[3] := A[2] + B[3]
- S2<sup>4</sup> : A[4] := A[3] + B[4]
- S2<sup>5</sup> : A[5] := A[4] + B[5]
- S2<sup>6</sup> : A[6] := A[5] + B[6]
- S2<sup>7</sup> : A[7] := A[6] + B[7]
- S2<sup>8</sup> : A[8] := A[7] + B[8]

Dependences cross, from one iteration to next



# What is a loop-carried dependence?

- Consider two iterations I<sup>1</sup> and I<sup>2</sup>
- A dependence occurs between two statements  $S_p$  and  $S_q$  (not necessarily distinct), when an assignment in  $S_p^{11}$  refers to the same location as a use in  $S_q^{12}$
- In the example,

 $S_1$ : A[0] := 0 for I = 1 to 5  $S_2$ : A[I] := A[I-1] + B[I]

- The assignment is "A[I<sup>1</sup>] := ..."
- The use is "... := A[I<sup>2</sup>-1] ..."
- These refer to the same location when I<sup>1</sup> = I<sup>2</sup>-1
- Thus  $I^1 < I^2$ , ie the assignment is in an earlier iteration

# Notation: $S_2 \delta_{c} S_2$

## **Definition**: The dependence equation

#### A dependence occurs

- between two statements  $S_p$  and  $S_q$  (not necessarily distinct),
- when there exists a pair of loop iterations I<sup>1</sup> and I<sup>2</sup>,
- such that a memory reference in  $S_p$  in  $I^1$  may refer to the same location as a memory reference in  $S_q$  in  $I^2$ .
- This might occur if  $S_{\rm p}$  and  $S_{\rm q}$  refer to some common array A
- Suppose  $S_p$  refers to  $A[\phi_p(I)]$   $(\phi_p(I) \text{ is some subscript} expression involving I)$
- Suppose  $S_q$  refers to  $A[\phi_q(I)]$

solution to the equation

• A dependence of some kind occurs between  $S_p$  and  $S_q$  if there exists a

$$\phi_p(I^1) = \phi_q(I^2) \cdot$$
for integer values of I<sup>1</sup> and I<sup>2</sup> lying within the loop bounds

# Types of dependence

If a solution to the dependence equation exists, a dependence of some kind occurs

The dependence type depends on what solutions exist

- The solutions consist of a set of pairs (l<sup>1</sup>,l<sup>2</sup>)
- · We would appear to have a data dependence if

```
A[\phi_p(I)] \in OUT(S_p)
```

and

 $A[\phi_q(\mathbf{I})] \in IN(S_q)$ 

- But we only really have a data dependence if the assignments *precede* the uses, ie
  - $S_p \delta_{c} S_q$
  - if, for each solution pair  $(I^1, I^2)$ ,  $I^1 < I^2$

#### Dependence versus anti-dependence

• If the *uses* precede the *assignments*, we actually have an *anti*-dependence, ie



if, for each solution pair  $(\mathbf{I}^1, \mathbf{I}^2), \mathbf{I}^1 > \mathbf{I}^2$ 

- In this case we do have a constraint on execution order
- Because we (may) have to read a value before it (may) be overwritten
- And this anti-dependence is loop-carried
- Anti-dependences prevent re-ordering, and multi-thread parallelism

### Dependence versus anti-dependence

 If there are some solution pairs (I<sup>1</sup>,I<sup>2</sup>) with I<sup>1</sup> < I<sup>2</sup> and some with I<sup>1</sup> > I<sup>2</sup>, we write

$$S_p \ \delta_* \ S_q$$

This represents that we know we must respect execution ordering, even though the compiler is unable to classify the dependence fully

 If, for all solution pairs (I<sup>1</sup>, I<sup>2</sup>), I<sup>1</sup> = I<sup>2</sup>, there are dependences *within* an iteration of the loop, but there are no loop-carried dependences:

$$S_p \delta_{=} S_q$$

#### **Dependence distance**

In many common examples, the set of solution pairs is characterised easily:

- **Definition:** dependence distance
  - If, for all solution pairs (I<sup>1</sup>, I<sup>2</sup>),

 $I^1 = I^2 - k$ 

then the dependence distance is **k** 

• For example in the loop we considered earlier,

$$S_1$$
: A[0] := 0  
for I = 1 to 5  
 $S_2$ : A[I] := A[I-1] + B[I]

We find that  $S_2 \quad \delta_1 \quad S_2$  with dependence distance 1.

• ((of course there are many cases where the difference is not constant and so the dependence cannot be summarised this way)).

#### **Reuse distance**

When optimising for cache performance, it is sometimes useful to consider the re-use relationship,

• IN(S<sub>1</sub>)  $\cap$  IN(S<sub>2</sub>)

- Here there is no dependence it doesn't matter which read occurs first
- Nonetheless, cache performance can be improved by minimising the reuse distance
- The reuse distance is calculated essentially the same way

🕨 Eg

for I = 5 to 100

```
S1: B[I] := A[I] * 2
```

S2: C[I] := A[I-5] \* 10

Here we have a loop-carried reuse with distance 5

Compilers - Chapter 8: Loop scheduling optimisations Part 3: Dependence analysis in nested loops

- Lecturer:
  - Paul Kelly (p.kelly@imperial.ac.uk)

PAGE 3				
DEPARTMENT	COURSE	DESCRIPTION	PREREQS	
COMPUTER SCIENCE	CPSC 432	INTERMEDIATE COMPILER DESIGN, WITH A FOCUS ON DEPENDENCY RESOLUTION.	CPSC 432	
	Conce Lund	COMPANY COMPANY DECK	DUTE IN	

### **Nested loops**

- Up to now we have looked at single loops
- Now let's generalise to loop "nests"
- We begin by considering a very common dependence pattern, called the "wavefront":

for 
$$I_1 = 0$$
 to 3 do  
for  $I_2 = 0$  to 3 do  
 $S : A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$ 

Dependence structure?

### **Nested loops**

I is  $I_1$ J is  $I_2$ 

- Up to now we have looked at single loops
- Now let's generalise to loop "nests"
- We begin by considering a very common dependence pattern, called the "wavefront":

Dependence structure?

#### System of dependence equations

Consider the dependence equations for this loop nest:

for  $I_1 = 0$  to 3 do for  $I_2 = 0$  to 3 do  $S: A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$ 

There are two potential dependences arising from the three references to A, so two systems of dependence equations to solve:

1. Between A[
$$I_1^1, I_2^1$$
] and A[ $I_1^2 - 1, I_2^2$ ]:  

$$\begin{cases}
I_1^1 = I_1^2 - 1 \\
I_2^1 = I_2^2
\end{cases}$$

2. Between A[
$$I_1^1$$
,  $I_2^1$ ] and A[ $I_1^2$ ,  $I_2^2 - 1$ ]:  

$$\begin{cases} I_1^1 = I_1^2 \\ I_2^1 = I_2^2 - 1 \end{cases}$$

• The same loop:

Iteration space graph

for 
$$I_1 = 0$$
 to 3 do  
for  $I_2 = 0$  to 3 do  
 $S: A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$ 

 For humans the easy way to understand this loop nest is to draw the *iteration space graph* showing the iteration-toiteration dependences:



 The diagram shows an arrow for each solution of each dependence equation. • The same loop:

#### Iteration space graph

for  $I_1 = 0$  to 3 do for  $I_2 = 0$  to 3 do  $S: A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$ 

 For humans the easy way to understand this loop nest is to draw the *iteration space graph* showing the iteration-toiteration dependences:



 The diagram shows an arrow for each solution of each dependence equation. Is there any parallelism?



- The inner loop is not vectorisable since there is a dependence chain linking successive iterations.
  - (to use a vector instruction, need to be able to operate on each element of the vector in parallel)
- Similarly, the outer loop is not parallel



- The inner loop is not vectorisable since there is a dependence chain linking successive iterations.
  - (to use a vector instruction, need to be able to operate on each element of the vector in parallel)
- Similarly, the outer loop is not parallel



- The inner loop is not vectorisable since there is a dependence chain linking successive iterations.
  - (to use a vector instruction, need to be able to operate on each element of the vector in parallel)
- Similarly, the outer loop is not parallel
- This loop nest has two dependence distance vectors:
  - (1,0) carried by the outer loop
     Direction vector: (<,=)</li>
  - (0,1) carried by the inner loop Direction vector: (=,<)</li>



- The inner loop is not vectorisable since there is a dependence chain linking successive iterations.
  - (to use a vector instruction, need to be able to operate on each element of the vector in parallel)
- Similarly, the outer loop is not parallel
- This loop is *interchangeable*: the top-to-bottom, left-to-right execution order is also valid since all dependence constraints (as shown by the arrows) are still satisfied.
- Interchanging the loop does not improve vectorisability or parallelisability



# Can you think of a loop like this that *cannot* safely be interchanged?



# Can you think of a loop like this that *cannot* safely be interchanged?









### Interchange: condition

- A loop is *interchangeable* if all dependence constraints (as shown by the arrows) are still satisfied by the top-tobottom, left-to-right execution order
- How can you tell whether a loop can be interchanged?
- Look at its dependence direction vectors:
  - Is there a dependence direction vector with the form (<,>) ?
  - ie there is a dependence distance vector (k<sub>1</sub>,k<sub>2</sub>) with k<sub>1</sub>>0 and k<sub>2</sub><0 ?</li>
  - If so, interchange would be invalid

Because the arrows would be traversed backwards

All other dependence directions are OK.

Consider this variation on the wavefront loop:	Skewing	
for k <sub>1</sub> := 0 to 3 do		
for $k_2 := k_1$ to $k_1$ +3 do		
$S : A[k_1,k_2-k_1] := A[k_1-1,k_2-k_1]+A[k_1$	,k <sub>2</sub> -k <sub>1</sub> -1]	

- The inner loop's control variable runs from  $k_1$  to  $k_1+3$ .
- The iteration space of this loop has 4<sup>2</sup> iterations just like the original loop.
- If we draw the iteration space with each iteration  $S^{K_1,K_2}$  at coordinate position  $(K_1,K_2)$ , it is skewed to form a lozenge shape:

#### racarvas camantics **ZOWIDA**

					y pro		S SCHIAHUUS		
for k	1 := (	) to 3	do						
for	k <sub>2</sub> :=	= k <sub>1</sub> t	o k <sub>1</sub> +:	3 do					
S: $A[k_1,k_2-k_1]$ := $A[k_1-1,k_2-k_1]+A[k_1,k_2-k_1-1]$									
$\varsigma 00$	$\varsigma 01$	$\varsigma$ 02	<b>ç</b> 03				To see that this		
$\mathcal{O}$							loop performs the		
$A_{OO}$	$A_{01}$	$A_{02}$	$A_{03}$				same		
	$S^{11}$	$S^{12}$	$S^{13}$	$S^{14}$			computation, lets		
	$A_{10}$	$A_{11}$	$A_{12}$	A <sub>13</sub>			work out its		
		S <sup>22</sup>	S <sup>23</sup>	S <sup>24</sup>	S <sup>25</sup>		dependence		
		$A_{20}$	$A_{21}$	A <sub>22</sub>	A <sub>23</sub>		First label each		
			S <sup>33</sup>	S <sup>34</sup>	S <sup>35</sup>	S <sup>36</sup>	iteration with the		
			$A_{30}$	A <sub>31</sub>	A <sub>32</sub>	$A_{33}$	element of A to		
The loop body is which it assi							which it assigns		

The loop body is

 $A[k_1,k_2-k_1] := A[k_1-1,k_2-k_1] + A[k_1,k_2-k_1-1]$ 

• E.g. iteration S<sub>23</sub> does: A[2,1] := A[1,1] + A[2,0] Skewing doesn't actually change the order in which the loop body is executed

75

Thus the dependence structure of the skewed loop is shown by marking the iteration space with all the dependences:



Can this loop nest be interchanged?



#### Interchange after skewing

Thus the dependence structure of the skewed loop is shown by marking the iteration space with all the dependences:



#### Transposed execution order

- You can think of loop interchange as changing the way the iteration space is traversed
- Alternatively, you can think of it as a change to the way the runtime code instances are mapped onto the iteration space
- Traversal is always lexicographic – ie left-toright, top-down


- The inner loop is now vectorisable, since it has no loop-carried dependence
- The skewed iteration space has N rows and 2N-1 columns, but still only N<sup>2</sup> actual statement instances.



# Skewing and interchange: summary $\overline{S_{S^{02}}}$

- - for  $I_2 = 0$  to 3 do  $S: A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$
- Original loop interchangeable but not vectorisable.
- We skewed inner loop by outer loop by factor 1.
- Still not vectorisable, but interchangeable.
- Interchanged, skewed loop is vectorisable.
- Bounds of new loop not simple!

for  $k_2 := 0$  to  $2N_2 - 2$  do for  $k_1 := \max(0, K_2 - N_2 + 2)$  to  $\min(K_2, N_1)$  do  $S : A[k_1, k_2 - k_1] := A[k_1 - 1, k_2 - k_1] + A[k_1, k_2 - k_1 - 1]$ 

- Is skewing ever invalid?
- Does skewing affect interchangeability?
- Does skewing affect dependence distances?
- Can you predict value of skewing?

# Summary: dependence

- Dependence equation for single loop:
  - Suppose  $S_p$  refers to  $A[\phi_p(I)]$
  - Suppose S<sub>q</sub> refers to  $A[\phi_q(I)]$
  - A dependence of some kind occurs between  ${\rm S}_{\rm p}$  and  ${\rm S}_{\rm q}$  if there exists a solution to the equation

 $\boldsymbol{\varphi}_{p}(\mathbf{I}^{1}) = \boldsymbol{\varphi}_{q}(\mathbf{I}^{2})$ 

- for integer values of  $I^1$  and  $I^2$  lying within the loop bounds
- For multidimensional arrays, and nested for-loops, we generalise this to a system of simultaneous dependence equations for two iterations,  $(I_1^1, I_2^1)$  and  $(I_1^2, I_2^2)$
- Iteration space graph, lexicographic schedule of execution
- Arrows in graph show solutions to dependence equation
- Dependence distance vectors characterise families of congruent arrows

# Summary: transformations

Not explained yet

- A loop can be executed in parallel if it has no loop-carried dependence
- A loop nest can be interchanged if the transposed dependence distance vectors are lexicographically forward
- Strip-mining is always valid
- Tiling = strip-mining + interchange
- 🕨 Skewing is always valid
- Skewing can expose parallelism by aligning parallel iterations with one of the loops
- Skewing can make interchange (and therefore tiling) valid

Student question: "why is antidependence a dependence?"



Loop-carried true dependence:

for i

A[i] = A[i-1] + B[i]

Loop-carried anti-dependence:

for i

```
A[i] = A[i+1] + B[i]
```

# "Loop-carried anti-dependence"

- Consider this example:
   for (int)i=0; i<7; i++)</li>
   c[i] = c[i+1] + b[i];
- When executed we get:

$$c[0] = c[1] + b[1];$$
  

$$c[1] = c[2] + b[2];$$
  

$$c[2] = c[3] + b[3];$$
  

$$c[3] = c[4] + b[3];$$
  

$$c[4] = c[5] + b[4];$$
  

$$c[4] = c[5] + b[5];$$
  

$$c[5] = c[6] + b[6];$$
  

$$c[6] = c[7] + b[7];$$

Each iteration uses a value which is overwritten in the next iteration

We need the use to happen before the overwrite

So we have a *precedence* requirement due to an anti-dependence

The anti-dependence arrows go from one iteration to the next

The anti-dependence is *carried* by the loop

### Implementing shared-memory parallel loop

Barrier(): block until

all threads reach this

point



FetchAndAdd....

```
for (i=0; i<N; i++) {
A[i] = A[i] + B[i];
```

#### Thread #0

```
if (myThreadId() == 0)
i = 0;
barrier();
while (true) {
    local_i = FetchAndAdd(&i);
    if (local_i >= N) break;
    A[local_i] = A[local_i] + B[local_i];
}
barrier();
```

### Thread #1

```
if (myThreadId() == 0)
i = 0;
barrier();
while (true) {
    local_i = FetchAndAdd(&i);
    if (local_i >= N) break;
    A[local_i] = A[local_i] + B[local_i];
}
barrier();
```

# Thread #0 gets some sequence of iterations to do

Thread #1 gets some sequence of iterations to do

June 25

```
for (i=0; i<N; i++) {
A[i] = A[i+1] + B[i];
```

# What could possibly go wrong?

Thread	<b>#0</b>
--------	-----------

```
if (myThreadId() == 0)
    i = 0;
    barrier();
    while (true) {
        local_i = FetchAndAdd(&i);
        if (local_i >= N) break;
        A[local_i] = A[local_i+1] + B[local_i];
    }
    barrier();
```

### Thread #1

```
if (myThreadId() == 0)
i = 0;
barrier();
while (true) {
    local_i = FetchAndAdd(&i);
    if (local_i >= N) break;
    A[local_i] = A[local_i+1] + B[local_i];
}
```

```
barrier();
```

Thread #0 gets some sequence of iterations to do

Thread #1 gets some sequence of iterations to do

```
for (i=0; i<N; i++) {
A[i] = A[i+1] + B[i];
```

### What could possibly go wrong?

This example has a loop-carried anti-dependence. We must read from A before overwriting A

Thread #0	Thread #1
if (myThreadId() == 0)	if (myThreadId() == 0)
i = 0;	i = 0;
barrier();	barrier();
while (true) {	while (true) {
local_i = FetchAndAdd(&i);	local_i = FetchAndAdd(&i);
if (local_i >= N) break;	if (local_i >= N) break;
A[local_i] = A[local_i+1] + B[local_i];	A[local_i] = A[local_i+1] + B[local_i];
}	}
barrier();	barrier();

Thread #0 gets some sequence of iterations to do, eg: 0, 2, 4, 6...

Thread #1 gets some sequence of iterations to do, eg: 1, 3, 5, 7...

### Feeding curiosity: solving the dependence equation (not examinable)

from z3 import \*

N=100

i1 = Int("i1")

i2 = Int("i2")

# consider a loop like this:

# for i = 1 to N

# a[phi1(i)] = a[phi2(i)] + b[i]

# So the dependence equation is

```
# exists i1, i2: 1<i<n s.t. phi1(i1) == phi2(i2)</pre>
```

def DependenceTest(bounds, dependence\_equation):

s = Solver()

s.add( bounds, dependence\_equation )

if s.check() == unsat:

print ("No dependence is present")

#### else:

```
print("Dependence is found, for example when:")
```

m = s.model()

```
print ("i1 = %s (LHS)" % m[i1])
```

```
print ("i2 = %s (RHS)" % m[i2])
```

#### Example 1:

print("for i = 1 to N")
print(" a[i] = a[i-1] + b[i]")

DependenceTest( And(i1>=1, i1<N, i2>=1, i2<N),

i1 == i2-1)



for i = 1 to N a[i] = a[i-1] + b[i] Dependence is found, for example when: i1 = 1 (LHS) i2 = 2 (RHS)

#### Just add the constraints and call the solver

### Feeding curiosity: solving the dependence equation

Example 1: def DependenceTest(bounds, dependence equation): s = Solver()print("for i = 1 to N") s.add(bounds, dependence equation) if s.check() == unsat: print(" a[i] = a[i-1] + b[i]")print ("No dependence is present") DependenceTest( And(i1>=1, i1<N, i2>=1, i2<N), else: print("Dependence is found, for example when:") i1 == i2 - 1) m = s.model()print ("i1 = %s (LHS)" % m[i1]) print ("i2 = %s (RHS)" % m[i2]) # Is there a loop-carried true dependence?  $s_2 = Solver()$ for i = 1 to N s2.add( bounds, dependence\_equation, i1<i2 ) if s2.check() == unsat: a[i] = a[i-1] + b[i]print ("No loop-carried true dependence is present") Dependence is found, for example when: else: print("Loop-carried true dependence found, for example when:") i1 = 1 (LHS) m = s2.model()i2 = 2 (RHS) print ("i1 = %s" % m[i1]) print ("i2 = %s" % m[i2]) Loop-carried true dependence found, for example # Is there a loop-carried anti-dependence? when: s3 = Solver()s3.add( bounds, dependence equation, i1>i2 ) i1 = 1if s3.check() == unsat: i2 = 2print ("No loop-carried anti-dependence is present") else: No loop-carried anti-dependence is present print("Loop-carried anti-dependence found, for example when:") m = s3.model()print ("i1 = %s" % m[i1]) Extend to distinguish loop-carried true and anti-dependencies print ("i2 = %s" % m[i2])

94

### Feeding curiosity: solving the dependence equation

def DependenceTest(bounds, dependence equation): s = Solver()s.add(bounds, dependence equation) if s.check() == unsat: print ("No dependence is present") else: print("Dependence is found, for example when:") m = s.model()print ("i1 = %s (LHS)" % m[i1]) print ("i2 = %s (RHS)" % m[i2]) # Is there a loop-carried true dependence?  $s_2 = Solver()$ s2.add( bounds, dependence equation, i1<i2 ) if s2.check() == unsat: print ("No loop-carried true dependence is present") else: print("Loop-carried true dependence found, for example when:") m = s2.model()print ("i1 = %s" % m[i1]) print ("i2 = %s" % m[i2]) # Is there a loop-carried anti-dependence? s3 = Solver()s3.add( bounds, dependence equation, i1>i2 ) if s3.check() == unsat: print ("No loop-carried anti-dependence is present") else: print("Loop-carried anti-dependence found, for example when:") m = s3.model()print ("i1 = %s" % m[i1])

print ("i2 = %s" % m[i2])

#### Example 2:

print("for i = 1 to N") print(" a[i] = a[i] + b[i]")

DependenceTest( And(i1>=1, i1<N, i2>=1, i2<N),

i1 == i2)



for i = 1 to N a[i] = a[i] + b[i] Dependence is found, for example when: i1 = 1 (LHS) i2 = 1 (RHS) No loop-carried true dependence is present No loop-carried anti-dependence is present

95

#### In this case the dependence is present but not loop-carried

June 25

### Feeding curiosity: solving the dependence equation

def DependenceTest(bounds, dependence equation): s = Solver()s.add(bounds, dependence equation) if s.check() == unsat: print ("No dependence is present") else: print("Dependence is found, for example when:") m = s.model()print ("i1 = %s (LHS)" % m[i1]) print ("i2 = %s (RHS)" % m[i2]) # Is there a loop-carried true dependence?  $s_2 = Solver()$ s2.add( bounds, dependence equation, i1<i2 ) if s2.check() == unsat: print ("No loop-carried true dependence is present") else: print("Loop-carried true dependence found, for example when:") m = s2.model()print ("i1 = %s" % m[i1]) print ("i2 = %s" % m[i2]) # Is there a loop-carried anti-dependence? s3 = Solver()s3.add( bounds, dependence equation, i1>i2 ) if s3.check() == unsat: print ("No loop-carried anti-dependence is present") else: print("Loop-carried anti-dependence found, for example when:") m = s3.model()print ("i1 = %s" % m[i1]) print ("i2 = %s" % m[i2])

### Example 3:

print("for i = 1 to N") print(" a[2\*i] = a[2\*i-1] + b[i]") DependenceTest( And(i1>=1, i1<N, i2>=1, i2<N), 2\*i1 == 2\*i2-1 )



for i = 1 to N a[2\*i] = a[2\*i-1] + b[2\*i] No dependence is present

### **Feeding curiosity: solving the dependence equation** Example 4:

def DependenceTest(bounds, dependence equation): s = Solver()s.add(bounds, dependence equation) if s.check() == unsat: print ("No dependence is present") else: print("Dependence is found, for example when:") m = s.model()print ("i1 = %s (LHS)" % m[i1]) print ("i2 = %s (RHS)" % m[i2]) # Is there a loop-carried true dependence?  $s_2 = Solver()$ s2.add( bounds, dependence equation, i1<i2 ) if s2.check() == unsat: print ("No loop-carried true dependence is present") else: print("Loop-carried true dependence found, for example when:") m = s2.model()print ("i1 = %s" % m[i1]) print ("i2 = %s" % m[i2]) # Is there a loop-carried anti-dependence? s3 = Solver()s3.add( bounds, dependence equation, i1>i2 ) if s3.check() == unsat: print ("No loop-carried anti-dependence is present") else: print("Loop-carried anti-dependence found, for example when:") m = s3.model()print ("i1 = %s" % m[i1]) print ("i2 = %s" % m[i2])

#### print("for i = 1 to N") print(" a[3\*i] = a[5\*i-10] + b[i]")DependenceTest( And(i1>=1, i1<N, i2>=1, i2<N), 3\*i1 == 5\*i2-20) for i = 1 to N a[3\*i] = a[5\*1-20] + b[i]Dependence is found, for example when: i1 = 5 (LHS) i2 = 7 (RHS) Loop-carried true dependence found, for example when: i1 = 5 $i^2 = 7$ Loop-carried anti-dependence found, for example when:

i1 = 15

i2 = 13

#### In this case we have both true and anti-dependences: weird!

S2 : A[i] := A[i-1] + B[i] sometimes still be parallelised

Appears to be inherently sequential



S2 : A[i] := A[i-1] + B[i] sometimes still be parallelised

Appears to be inherently sequential

But parallel is possible:





"Parallel scan" or "parallel prefix sum"

S2 : A[i] := A[i-1] + B[i] sometimes still be parallelised

Appears to be inherently sequential

**But parallel implementation is possible** 



"Parallel scan" or "parallel prefix sum"

S2 : A[i] := A[i-1] + B[i] sometimes still be parallelised

Appears to be inherently sequential

But parallel implementation is possible



We can see that the last element is computed with a reduction tree

S2 : A[i] := A[i-1] + B[i] sometimes still be parallelised

Appears to be inherently sequential

But parallel implementation is possible



All the elements are computed by reduction trees of depth log(N) – for example element 7

S1 : A[0] := 0 for i = 1 to 8

S2 : A[i] := A[i-1] + B[i]

Appears to be inherently sequential

But parallel implementation is possible



"Parallel scan" or "parallel prefix sum"

# **Feeding curiosity**

- This is the "naïve" parallel scan
- It does more work than the sequential scan – but it does use parallelism
- There are "workefficient" parallel scans
- Eg see Mark Harris, GPU Gems Ch39

https://developer.nvidia.com/gp ugems/gpugems3/part-vi-gpucomputing/chapter-39-parallelprefix-sum-scan-cuda

### **Compilers** - That wraps it up!

We have seen....

- How to build a simple non-optimising compiler for a simple imperative language
- With functions
- With local variables, static variables, heap data, inheritance
- We have seen how an optimising compiler might work
  - intermediate representations, lowering,
  - dataflow analysis, register allocation, code motion optimisations,
  - instruction selection
  - SSA
- Dependence analysis and parallelisation
  - Loop-carried dependence
  - Dependence distance
  - Vectorisation and parallelisation
- You have been introduced to a world of topics fundamental to how your code actually gets executed, and what can be done to make it efficient

Compilers - Chapter 8: Loop scheduling optimisations Part 4: Representing loop transformations as matrix multiplications

- Lecturer:
  - Paul Kelly (p.kelly@imperjal.ac.uk)

This section is not examinable



### Matrix representation of loop transformations

To skew the inner loop by the outer loop by factor 1 we adjust the loop bounds, and replace I<sub>1</sub> by K<sub>1</sub>, and I<sub>2</sub> by K<sub>2</sub>-K<sub>1</sub>. That is,

 $(K_1, K_2) = (I_1, I_2) \cdot U$ 

where U is a 2 x 2 matrix

$$\mathbf{U} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$$

• That is,

lite-

$$(K_1, K_2) = (I_1, I_2) \cdot U = (I_1, I_2 + I_1)$$

### Matrix representation of loop transformations

To skew the inner loop by the outer loop by factor 1 we adjust the loop bounds, and replace I<sub>1</sub> by K<sub>1</sub>, and I<sub>2</sub> by K<sub>2</sub>-K<sub>1</sub>. That is,

$$(K_1, K_2) = (I_1, I_2) \cdot U$$

where U is a 2 x 2 matrix

$$\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \qquad \mathbf{U}^{-1} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix}$$

• That is,

$$(K_1, K_2) = (I_1, I_2) \cdot U = (I_1, I_2 + I_1)$$

The inverse gets us back again:

$$(I_1, I_2) = (K_1, K_2) \cdot U^{-1} = (K_1, K_2 - K_1)$$

 Matrix U maps each statement instance S<sup>I1I2</sup> to its position in the new iteration space, S<sup>K1K2</sup>:

### Original iteration space:

for 
$$I_1 = 0$$
 to 3 do  
for  $I_2 = 0$  to 3 do  
 $S: A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$   
for  $k_1 := 0$  to 3 do  
for  $k_2 := k_1$  to  $k_1+3$  do  
 $S: A[k_1, k_2-k_1] := A[k_1-1, k_2-k_1]+A[k_1, k_2-k_1-1]$ 

The subscripts are mapped back using U<sup>-1</sup>

 $(K_1, K_2) = (I_1, I_2) \cdot U = (I_1, I_2 + I_1)$ 

 $(I_1, I_2) = (K_1, K_2) \cdot U^{-1} = (K_1, K_2 - K_1)$ 

The matrix representation is not examinable

# Using matrices to reason about dependence

Recall that:

• There is a dependence between two iterations  $(I_1^1, I_2^1)$  and  $(I_1^2, I_2^2)$  if there is a memory location which is assigned to in iteration  $(I_1^1, I_2^1)$ , and read in iteration  $(I_1^2, I_2^2)$ .

((unless there is an intervening assignment))

- If  $(I_1^1, I_2^1)$  precedes  $(I_1^2, I_2^2)$  it is a *data*-dependence.
- If  $(I_1^2, I_2^2)$  precedes  $(I_1^1, I_2^1)$  it is a *anti*-dependence.
- If the location is assigned to in both iterations, it is an *output*-dependence.
- The dependence distance vector  $(D_1, D_2)$  is  $(I_1^1 I_1^2, I_2^1 I_2^2)$ .

# **Transforming dependence vectors**

- If there is a dependence between two iterations  $(I_1^1, I_2^1)$  and  $(I_1^2, I_2^2)$
- Then iterations  $(I_1^1, I_2^1)$ . U and  $(I_1^2, I_2^2)$ . U will also read and write the same location
- The transformation U is valid iff

 $(I_1^1, I_2^1)$ . U precedes  $(I_1^2, I_2^2)$ . U whenever there is a dependence between  $(I_1^1, I_2^1)$  and  $(I_1^2, I_2^2)$ .

 In the transformed loop the dependence distance vector is also transformed, to

 $(D_1, D_2) . U$ 

• U is a valid transformation if all the program's dependence distance vectors are still "forward" when transformed by U

# **Transforming dependence vectors**

• What do we mean by "precedes"?

```
Definition: Lexicographic ordering:
(I^1,J^1) precedes (I^2,J^2)
if I^1 < I^2, or I^1 = I^2 and J^1 < J^2
```

- "Lexicographic" is dictionary order both "baz" and "can" precede "cat"
- So (1,2) precedes (1,3)
- But (0,3) precedes (1,4)
- A dependence distance vector (D<sub>1</sub>,D<sub>2</sub>) is lexicographically "forward" if it precedes (0,0)

# Example: loop given earlier

- Before transformation we had two dependences:
- 1. Distance: (1,0), direction: (<,.)
- 2. Distance: (0,1), direction: (.,<)
- After transformation by matrix

 $\mathbf{U} = \left[ \begin{array}{cc} 1 & 1 \\ 0 & 1 \end{array} \right]$ 

- (i.e. skewing of inner loop by outer) we get:
- 1. Distance: (1,0).U = (1,1), direction: (<,<)
- 2. Distance: (0,1).U = (0,1), direction: (.,<)

for  $I_1 = 0$  to 3 do for  $I_2 = 0$  to 3 do  $S: A[I_1, I_2] := A[I_1 - 1, I_2] + A[I_1, I_2 - 1]$ for  $k_1 := 0$  to 3 do for  $k_2 := k_1$  to  $k_1+3$  do  $S: A[k_1,k_2-k_1] := A[k_1-1,k_2-k_1]+A[k_1,k_2-k_1-1]$  $I_2$  :  $\frac{I_1}{0}$ S<sup>03</sup> 1 2  $S^{30}$  $S^{31}$  $S^{32}$  $S^{33}$ 

We can also represent loop interchange by a matrix transformation.

After transforming the skewed loop by matrix

$$\mathbf{V} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

(i.e. loop interchange) we get:

- 1. Distance: (1,0).U.V = (1,1).V = (1,1), direction: (<,<)
- 2. Distance: (0,1).U.V = (0,1).V = (1,0), direction: (<,.)
- The transformed iteration space is the transpose of the skewed iteration space:

# Summary

- $(I_1, I_2)$ . U maps each statement instance  $(I_1, I_2)$  to its new position  $(K_1, K_2)$  in the transformed loop's execution sequence
- $(D_1, D_2)$ . U gives new dependence distance vector, giving test for validity
- Captures skewing, interchange and reversal
- Compose transformations by matrix multiplication

 $U_1 \cdot U_2$ 

- Resulting loop's bounds may be a little tricky
  - Efficient algorithms exist [Banerjee90] to maximise parallelism by skewing and loop interchanging
  - Efficient algorithms exist to optimise cache performance by finding the combination of blocking, block size, interchange and skewing which leads to the best reuse [Wolf91]

Restructuring compilers - conclusions:

- Restructuring compilers can find parallelism
- And enhance locality

### For a very restricted class of programs

For-loops over arrays with array subscripts that are simple ("affine") expressions involving loop control variables

But for this restricted class there is a rather elegant theory (the "polyhedral" or "polytope" model, <u>http://en.wikipedia.org/wiki/Polytope\_model</u>)

Extending beyond this is a big research problem

Current compilers (GCC, Clang/LLVM, Intel, Microsoft etc) can do some of this, in theory – but are often defeated by program complexity

#### Textbooks covering restructuring compilers

### References

- Michael Wolfe. High Performance Compilers for Parallel Computing. Addison Wesley, 1996.
- Steven Muchnick, Advanced Compiler Design and Implementation. Morgan Kaufmann, 1997.
- Ken Kennedy and Randy Allen, Optimizing Compilers for Modern Architectures. Morgan Kaufmann, 2001.

#### Research papers:

 D. F. Bacon and S. L. Graham and O. J. Sharp, "Compiler Transformations for High-Performance Computing". ACM Computing Surveys V26 N4 Dec 1994 <u>http://doi.acm.org/10.1145/197405.197406</u>

U. Banerjee. Unimodular transformations of double loops. In Proceedings of the Third Workshop on Programming Languages and Compilers for Parallel Computing, Irvine, CA. Pitman/MIT Press, 1990.

M.E. Wolf and M.S. Lam. A data locality optimizing algorithm. In Proceedings of the ACM SIGPLAN '91 Conference on Programming Language Design and Implementation, volume 26, pages 30-44, Toronto, Ontario, Canada, June 1991.

### Polyhedral:

The Polyhedral Model (aka Polytope method) takes the ideas in this lecture much further, notably

- Automatic code generation using the matrix model shown here
- Checking validity of such transformations
- Calculating parallelism and locality metrics for alternative versions
  - https://en.wikipedia.org/wiki/Polytope\_model
  - https://normrubin.github.io/lectures/poly\_final.html
  - https://polyhedral.info/
  - https://www.impact-workshop.org/

# **Matrix transpose**

# **Feeding curiosity**

### Try this link to the the Compiler Explorer:

### **Collision detect**

### Try this link to the the Compiler Explorer:

64+icx+2025.0.0+(Compiler+%231)',t:'0')),k:33.3333333333333333,l:'4',n:'0',o:",s:0,t:'0')),l:'2',n:'0',o:",t:'0')),version:4

# **Feeding curiosity**

# Ask me about....

# Loop interchange for locality

- For i, j, k matrix multiply vs
- For i, k, j matrix multiply
- See /homes/phjk/ToyPrograms/ACA24-25/MM compare speed of versions MM1.c, MM2.c

# • Tiling for locality

- For the transpose example shown in the last chapter
- For matrix multiply (see version MM3.c)

# Stencils and convolutions

• skewed, split, diamond

# **Graphs and unstructured meshes**