















| Loop Interchange (cont) | |
|---|--|
| Example | |
| do i = 1, n | do j = 1,n |
| do $j = 1, n$ | do $i = 1, n$ |
| $\mathbf{x} = \mathbf{A}(\mathbf{i}, \mathbf{j})$ | $\mathbf{x} = \mathbf{A}(\mathbf{i},\mathbf{j})$ |
| enddo This array has stride | enddo This array now has stride 1 |
| enddo ⁿ access | enddo access |
| | |
| | |
| | |
| | |
| May 4, 2015 Loop Trar | sformations 9 |

Legality of Loop Interchange



(=,=)

The dependence is loop independent, so it is unaffected by interchange

(=,<)

The dependence is carried by the j loop.

After interchange the dependence will be (<,=), so the dependence will still be carried by the j loop, so the dependence relations do not change.

(<,=)

The dependence is carried by the i loop.

After interchange the dependence will be (=,<), so the dependence will still be carried by the i loop, so the dependence relations do not change.











| Motivation | | |
|--------------------|--|--|
| - Reduces loop or | verhead | |
| - Improves effecti | veness of other transformations | |
| - Code sched | uling | |
| – Code sched | uning | |
| The Transformatic | | |
| | | |
| – Make n copies of | of the loop: n is the unrolling factor | |
| – Adjust loop bou | nds accordingly | |
| | | |
| | | |
| | | |
| | | |



Loop Balance Problem - We'd like to produce loops with the right balance of memory operations and floating point operations - The ideal balance is machine-dependent -e.g. How many load-store units are connected to the L1 cache? -e.g. How many functional units are provided? **Example** do j = 1,2*n- The inner loop has 1 memory do i = 1, moperation per iteration and 1 floating A(j) = A(j) + B(i) point operation per iteration -If our target machine can only enddo support 1 memory operation for enddo every two floating point operations, this loop will be memory bound What can we do? May 4, 2015 Loop Transformations



| Unroll the Outer Loop | |
|---|------------------------------|
| do j = 1,2*n by 2 | |
| do i = 1,m | |
| A(j) = A(j) + | B(1) |
| enddo | |
| do i = $1,m$ | |
| A(j+1) = A(j+1) + | B(i) |
| enddo 🔶 | |
| enddo | |
| | Jam the inner loops |
| – The inner loop has 1 load per | do $j = 1,2*n$ by 2 |
| iteration and 2 floating point | do i = 1,m |
| operations per iteration | A(j) = A(j) + B(i) |
| - We reuse the loaded value of B | (i) $A(j+1) = A(j+1) + B(i)$ |
| - The Loop Balance matches the | enddo |
| | onddo |















Motivation

Limitations of static analysis

- Programs can have values and invariants that are known at runtime but unknown at compile time. Static compilers cannot exploit such values or invariants
- Many of the motivations for profile-guided optimizations apply here

Basic idea

- Perform translation at runtime when more information is known
- Traditionally, two types of translations are done
 - Runtime code generation
 - Partial evaluation

May 4, 2015



| Interpreters: | The program being interpreted is runtime cons | tant |
|--|--|------|
| Simulators: | The subject of simulation (circuit, cache, network is runtime constant | ork) |
| Graphics renderers: | The scene to render is runtime constant | |
| Scientific simulations: | Matrices can be runtime constants | |
| Extensible OS kernels: | Extensions to the kernel can be runtime consta | nt |
| Examples – A cache simulator migi – A partially evaluated si special case where the | ht take the line size as a parameter imulator might produce a faster simulator for the line size is 16 | |
| May 4, 2015 | Loop Transformations | 30 |



Dynamic Compilation with DyC

DyC [Auslander, *et al* 1996]

- Apply ideas of Partial Evaluation
- Perform some of the Partial Evaluation at runtime
 - Can handle more runtime constants than Partial Evaluation
- Reminiscent of link-time register allocation in the sense that the compilation is performed in stages

Tradeoffs

- Must overcome the run-time cost of the dynamic compiler
 - Fast dynamic compilation: low overhead
 - High quality dynamically generated code: high benefit
- Ideal: dynamically translate code once, execute this code many times
- Implication: don't dynamically translate everything
 - Only perform dynamic translation where it will be profitable

May 4, 2015







The Dynamic Compiler

The Stitcher

- Follows directives, which are produced by the static compiler, to copy code templates and to fill in holes with appropriate constants
- The resulting code becomes part of the executable code and is hopefully executed many times











The Need for Annotations

Automatic dynamic compilation is difficult

- Which variables are runtime constant over which pieces of code?
- Complicated by aliases, side effects, pointers that can modify memory
- Which loops are profitable to unroll?
- Estimating **profitability** is the difficult part

Annotation errors

- Lead to incorrect dynamic compilation
 - -e.g., Incorrect code if a value is not really a runtime constant



Detecting Runtime Constants

Simple data-flow analysis

Propagates initial runtime constants through the dynamic region using the following transfer functions

| $-\mathbf{x} = \mathbf{y}$ | x is a constant iff y is a constant |
|---|--|
| -x = y op z | x is a const iff y and z are consts and op is an idempotent, side-effect free, non-trapping op |
| $-\mathbf{x} = \mathbf{f}(\mathbf{y}_1, \dots, \mathbf{y}_n)$ | x is a const iff the y _i are consts and f is an idempotent, side-effect free, non-trapping function |
| -x = *p | x is a constant iff p is constant |
| – x = dynamic *p | x is not constant |
| May 4, 2015 | Loop Transformations 44 |



Optimizations

Integrated optimizations

- For best quality code, optimizations should be performed across dynamic region boundaries, *e.g.*, global CSE, global register allocation
- Optimizations can be performed both before and after the dynamic region has been split into setup and template codes

Restrictions on optimizing split code

- Instructions with holes cannot be moved outside of their dynamic region
- Holes cannot be treated as legal values outside of the dynamic region.
 (*e.g.*, Copy propagation cannot propagate values of holes outside of dynamic regions)
- Holes are typically viewed as constants throughout the dynamic region, but induction variables become constant for only a given iteration of an unrolled loop





Performance Results

Two measures of performance

- Asymptotic improvement: speedup if overhead were 0
- Break even point: the fewest number of iterations at which the dynamic compilation system is profitable

| calculator | | 916 interpretations |
|---------------------|-----|----------------------|
| matrix multiply | 1.6 | 31,392 scalar ×'s |
| sparse mat multiply | 1.8 | 2645 matrix ×'s |
| event dispatcher | 1.4 | 722 event dispatches |
| quicksort | 1.2 | 3050 records |

Evaluation

Today's discussion

- Simple caching scheme

- Setup once, reuse thereafter
- More sophisticated schemes are possible
 - Can cache multiple versions of code
 - Can provide eager, or speculative, specialization
 - Can allow different dynamic regions for different variables

Subsequent progress on DyC

- More sophisticated language and compiler [Grant, et al 1999]
 - More complexity is needed
 - Extremely difficult to annotate the applications
- Automated insertion of annotations [Mock, et al 2000]
 - Use profiling to obtain value and frequency information

May 4, 2015

Loop Transformations