## Lecture 17

## Dynamic Code Optimization

I. Motivation & Background
II. Overview
III. Partial Method Compilation
IV. Partial Dead Code Elimination
V. Partial Escape Analysis
VI. Results

"Partial Method Compilation Using Dynamic Profile Information",
John Whaley, OOPSLA 01

---

### I. Beyond Static Compilation

1) Profile-based Compiler: high-level → binary, static
   – Uses (dynamic=runtime) information collected in profiling passes

2) Interpreter: high-level, emulate, dynamic

3) Dynamic compilation / code optimization: high-level → binary, dynamic
   – interpreter/compiler hybrid
   – supports cross-module optimization
   – can specialize program using runtime information
     • without separate profiling passes

---

### 1) Dynamic Profiling Can Improve Compile-time Optimizations

• Understanding common dynamic behaviors may help guide optimizations
   – e.g., control flow, data dependences, input values

```
void foo(int A, int B) {
    …
    while (…) {
        if (A > B)              ← How often is this condition true?
            *p = 0;             ← How often does *p == val[i]?
        C = val[i] + D;         ← Is this loop invariant?
        E += C – B;
        …
    }
}
```
What are typical values of A, B?

• Profile-based compile-time optimizations
   – e.g., speculative scheduling, cache optimizations, code specialization

---

### Profile-Based Compile-time Optimization

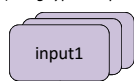1. Compile statically    2. Collect profile *(using typical inputs)*    3. Re-compile, using profile



• Collecting control-flow profiles is relatively inexpensive
   – profiling data dependences, data values, etc., is more costly
• Limitations of this approach?

1

## Instrumenting Executable Binaries

1. Compile statically     2. Collect profile
                          *(using typical inputs)*

**How to perform the instrumentation?**

- prog1.c → compiler → runme.exe
- input1 → runme.exe *(instrumented)* → execution profile → binary instrumentation tool

1. The compiler could insert it directly
2. A **binary instrumentation tool** could modify the executable directly
   - that way, we don't need to modify the compiler
   - compilers that target the same architecture (e.g., x86) can use the same tool

Carnegie Mellon

---

## Binary Instrumentation/Optimization Tools

- Unlike typical compilation, the input is a binary (not source code)
- One option: **static** binary-to-binary rewriting

  runme.exe → tool → runme_modified.exe

- <u>Challenges</u> (with the static approach):
  - what about dynamically-linked shared libraries?
  - if our goal is optimization, are we likely to make the code faster?
    - a compiler already tried its best, and it had source code (we don't)
  - if we are adding instrumentation code, what about time/space overheads?
    - instrumented code might be slow & bloated if we aren't careful
    - optimization may be needed just to keep these overheads under control

- <u>Bottom line</u>: the purely static approach to binary rewriting is rarely used

Carnegie Mellon

---

## 2) (Pure) Interpreter

- One approach to dynamic code execution/analysis is an **interpreter**
  - <u>basic idea</u>: a software loop that grabs, decodes, and emulates each instruction

```
while (stillExecuting) {
    inst = readInst(PC);
    instInfo = decodeInst(inst);
    switch (instInfo.opType) {
        case binaryArithmetic: …
        case memoryLoad: …
        …
    }
    PC = nextPC(PC,instInfo);
}
```

- <u>Advantages</u>:
  - also works for dynamic programming languages (e.g., Java)
  - easy to change the way we execute code on-the-fly (SW controls everything)
- <u>Disadvantages</u>:
  - runtime overhead!
    - *each dynamic instruction is emulated individually by software*

Carnegie Mellon

---

## A Sweet Spot?

- Is there a way that we can combine:
  - the flexibility of an interpreter (analyzing and changing code dynamically); and
  - the performance of direct hardware execution?

- <u>Key insights</u>:
  - increase the granularity of interpretation
    - ~~instructions~~ → chunks of code (e.g., procedures, basic blocks)
  - dynamically *compile* these chunks into directly-executed optimized code
    - store these compiled chunks in a software code cache
    - jump in and out of these cached chunks when appropriate
    - these cached code chunks can be updated!
  - invest more time optimizing code chunks that are clearly hot/important
    - easy to instrument the code, since already rewriting it
    - must balance (dynamic) compilation time with likely benefits
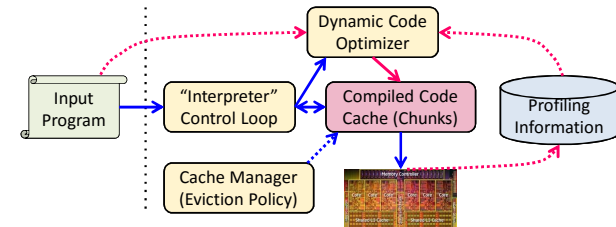
Carnegie Mellon

2

## 3) Dynamic Compiler

```
while (stillExecuting) {
    if (!codeCompiledAlready(PC)) {
        compileChunkAndInsertInCache(PC);
    }
    jumpIntoCodeCache(PC);
    // compiled chunk returns here when finished
    PC = getNextPC(…);
}
```

- This general approach is widely used:
  - Java virtual machines
  - dynamic binary instrumentation tools (Valgrind, Pin, Dynamo Rio)
  - hardware virtualization

---

## Components in a Typical Just-In-Time (JIT) Compiler



- Cached chunks of compiled code run at hardware speed
  - returns control to "interpreter" loop when chunk is finished
- Dynamic optimizer uses profiling information to guide code optimization
  - as code becomes hotter, more aggressive optimization is justified
    → replace the old compiled code chunk with a faster version

---

## II. Overview of Dynamic Compilation / Code Optimization

- Interpretation/Compilation/Optimization policy decisions
  - Choosing what and how to compile, and how much to optimize

- Collecting runtime information
  - Instrumentation
  - Sampling

- Optimizations exploiting runtime information
  - Focus on frequently-executed code paths
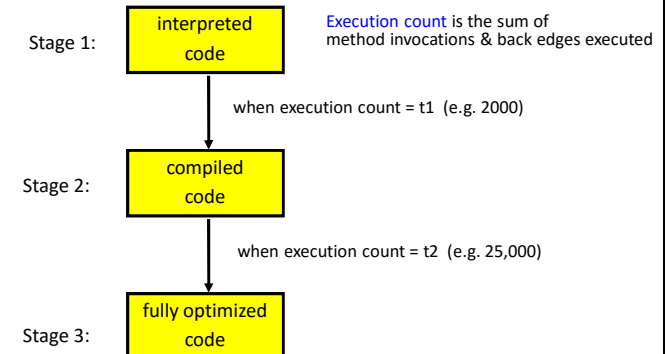
---

## Dynamic Compilation Policy

- $\Delta T_{total} = T_{compile} - (n_{executions} * T_{improvement})$
  - If $\Delta T_{total}$ is negative, our compilation policy decision was effective.

- We can try to:
  - Reduce $T_{compile}$ (faster compile times)
  - Increase $T_{improvement}$ (generate better code: but at cost of increasing $T_{compile}$)
  - Focus on large $n_{executions}$ (compile/optimize hot spots)

- 80/20 rule: Pareto Principle
  - 20% of the work for 80% of the advantage

3

## Latency vs. Throughput

- <u>Tradeoff</u>: startup speed vs. execution performance

| | Startup speed | Execution performance |
|---|---|---|
| Interpreter | Best | Poor |
| 'Quick' compiler | Fair | Fair |
| Optimizing compiler | Poor | Best |

**Carnegie Mellon**

---

## Multi-Stage Dynamic Compilation System

Stage 1:  [ interpreted code ]

*Execution count* is the sum of method invocations & back edges executed

when execution count = t1  (e.g. 2000)

Stage 2:  [ compiled code ]

when execution count = t2  (e.g. 25,000)

Stage 3:  [ fully optimized code ]

**Carnegie Mellon**

---

## Granularity of Compilation: Per Method?

- Methods can be large, especially after inlining
  - Cutting/avoiding inlining too much hurts performance considerably

- Compilation time is proportional to the amount of code being compiled
  - Moreover, many optimizations are not linear

- Even "hot" methods typically contain some code that is rarely/never executed

**Carnegie Mellon**

---

## Example: SpecJVM98 db

```
void read_db(String fn) {
  int n = 0, act = 0; int b; byte buffer[] = null;
  try {
    FileInputStream sif = new FileInputStream(fn);
    n = sif.getContentLength();
    buffer = new byte[n];
    while ((b = sif.read(buffer, act, n-act))>0) {
      act = act + b;
    }
    sif.close();
    if (act != n) {
      /* lots of error handling code, rare */
    }
  } catch (IOException ioe) {
    /* lots of error handling code, rare */
  }
}
```

Hot loop →

**Carnegie Mellon**

4

## Example: SpecJVM98 db

```
void read_db(String fn) {
  int n = 0, act = 0; int b; byte buffer[] = null;
  try {
    FileInputStream sif = new FileInputStream(fn);
    n = sif.getContentLength();
    buffer = new byte[n];
    while ((b = sif.read(buffer, act, n-act))>0) {
      act = act + b;
    }
    sif.close();
    if (act != n) {
      /* lots of error handling code, rare */
    }
  } catch (IOException ioe) {
    /* lots of error handling code, rare */
  }
}
```

Lots of rare code!

---

## Optimize hot "regions", not methods

- Optimize only the most frequently executed segments within a method
  - Simple technique:
    - Track execution counts of basic blocks in Stages 1 & 2
    - Any basic block executing in Stage 2 is considered to be not rare

- Beneficial secondary effect of improving optimization opportunities on the common paths

- No need to profile any basic block executing in Stage 3
  - Already fully optimized

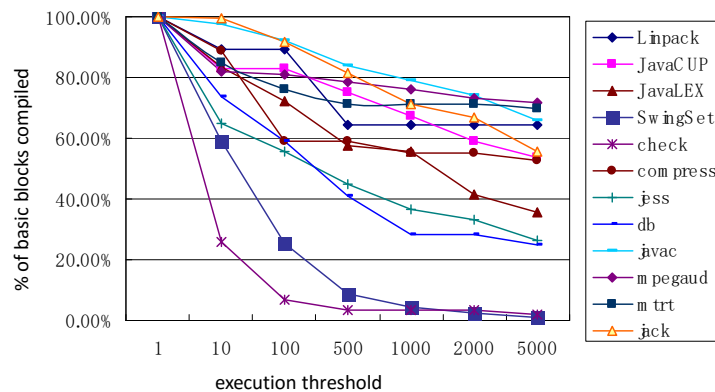Stage 1: interpreted code

Stage 2: compiled code

Stage 3: fully optimized code

---

## % of Basic Blocks in Methods that are Executed > Threshold Times
(hence would get compiled under per-method strategy)



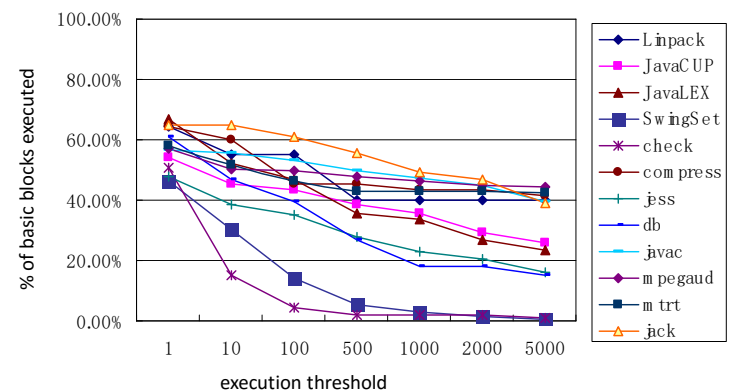Legend: Linpack, JavaCUP, JavaLEX, SwingSet, check, compress, jess, db, javac, mpegaud, mtrt, jack

y-axis: % of basic blocks compiled (0.00% – 100.00%)
x-axis: execution threshold (1, 10, 100, 500, 1000, 2000, 5000)

---

## % of Basic Blocks that are Executed > Threshold Times
(hence get compiled under per-basic-block strategy)



Legend: Linpack, JavaCUP, JavaLEX, SwingSet, check, compress, jess, db, javac, mpegaud, mtrt, jack

y-axis: % of basic blocks executed (0.00% – 100.00%)
x-axis: execution threshold (1, 10, 100, 500, 1000, 2000, 5000)

## Dynamic Code Transformations

- Compiling partial methods
- Partial dead code elimination
- Partial escape analysis

**Carnegie Mellon**

---

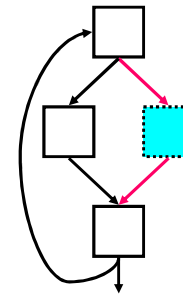## III. Partial Method Compilation

1. Based on profile data, determine the set of rare blocks
   - Use code coverage information from the first compiled version

Goal: Program runs correctly with white blocks compiled and blue blocks interpreted

What are the challenges?
- How to transition from white to blue
- How to transition from blue to white
- How to compile/optimize ignoring blue

**Carnegie Mellon**

---

## Partial Method Compilation

2. Perform live variable analysis
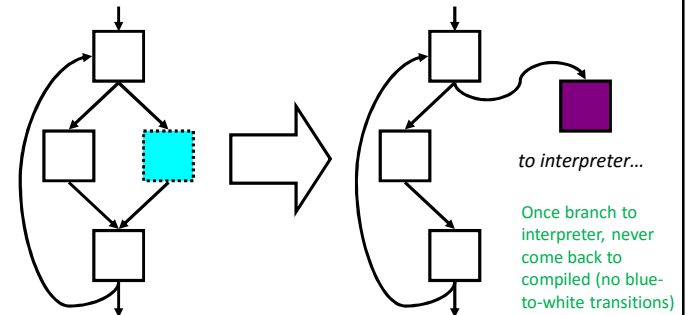   - Determine the set of live variables at rare block entry points

*live: x,y,z*

**Carnegie Mellon**

---

## Partial Method Compilation

3. Redirect the control flow edges that targeted rare blocks, and remove the rare blocks

*to interpreter...*

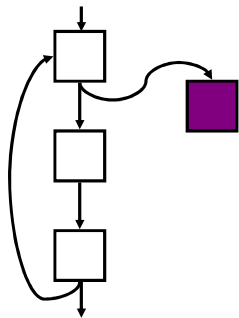Once branch to interpreter, never come back to compiled (no blue-to-white transitions)

**Carnegie Mellon**

6

## Partial Method Compilation

4. Perform compilation normally
   - Analyses treat the interpreter transfer point as an unanalyzable method call
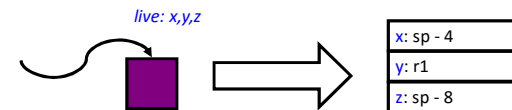
**Carnegie Mellon**

---

## Partial Method Compilation

5. Record a map for each interpreter transfer point
   - In code generation, generate a map that specifies the location, in registers or memory, of each of the live variables
   - Maps are typically < 100 bytes

*live: x,y,z*

| x: sp - 4 |
|-----------|
| y: r1 |
| z: sp - 8 |

**Carnegie Mellon**

---

## IV. Partial Dead Code Elimination

- Move computation that is only live on a rare path into the rare block, saving computation in the common case

**Carnegie Mellon**

---

## Partial Dead Code Example

```
x = 0;
if (rare branch 1){
    ...
    z = x + y;
    ...
}
if (rare branch 2){
    ...
    a = x + z;
    ...
}
```

```
if (rare branch 1) {
    x = 0;
    ...
    z = x + y;
    ...
}
if (rare branch 2) {
    x = 0;
    ...
    a = x + z;
    ...
}
```

**Carnegie Mellon**

7

## V. Escape Analysis

- Escape analysis finds objects that do not escape a method or a thread
  - "Captured" by method:
    - can be allocated on the stack or in registers
  - "Captured" by thread:
    - can avoid synchronization operations
- All Java objects are normally heap allocated, so this is a big win
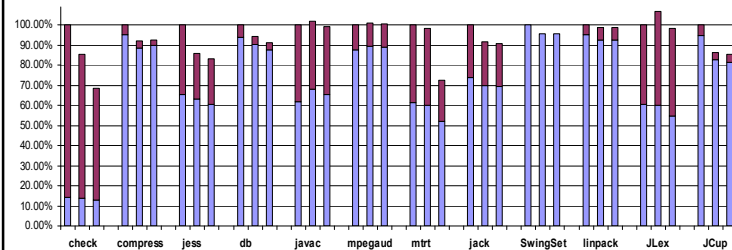
Carnegie Mellon

---

## Partial Escape Analysis

- Stack allocate objects that don't escape in the common blocks

- Eliminate synchronization on objects that don't escape the common blocks

- If a branch to a rare block is taken:
  - Copy stack-allocated objects to the heap and update pointers
  - Reapply eliminated synchronizations

Carnegie Mellon

---

## VI. Results: Run Time Improvement

First bar: original (Whole method opt)
Second bar: Partial Method Comp (PMC)
Third bar: PMC + opts
   Bottom bar: Execution time if code was compiled/opt. from the beginning

Carnegie Mellon

---

## Summary: Beyond Static Compilation

1) Profile-based Compiler: high-level → binary, static
   - Uses (dynamic=runtime) information collected in profiling passes

2) Interpreter: high-level, emulate, dynamic

3) Dynamic compilation / code optimization: high-level → binary, dynamic
   - interpreter/compiler hybrid
   - supports cross-module optimization
   - can specialize program using runtime information
     - without separate profiling passes
     - for what's hot on this particular run

Carnegie Mellon

## Looking Ahead

- Friday: No class

- Monday & Wednesday: "Recent Research on Optimization"
  - Student-led discussions, in groups of 2, with 20 minutes/group
  - Read 3 papers on a topic, and lead a discussion in class
  - See "Discussion Leads" tab of course web page for topics, sign-up sheet, instructions

- Spring Break

- Monday March 14
  - Homework #3 due
  - Meetings to discuss project proposal ideas