

# Code Generation for Data Processing

## Lecture 5: Analyses and Transformations


Alexis Engelke

Chair of Data Science and Engineering (I25)  
School of Computation, Information, and Technology  
Technical University of Munich

Winter 2022/23

# Program Transformation: Motivation

- ▶ “User code” is often not very efficient
- ▶ Also: no need to, compiler can (often?) optimize better
  - ▶ More knowledge: e.g., data layout, constants after inlining, etc.
- ▶ Allows for more pragmatic/simple code
  
- ▶ Generating “better” IR code on first attempt is expensive
  - ▶ What parts are actually used? How to find out?
- ▶ Transformation to “better” code must be done *somewhere*
  
- ▶ Optimization is a misnomer: we don’t know whether it improves code!
  - ▶ Many transformations are driven by heuristics
- ▶ Many types of optimizations are well-known<sup>9</sup>

<sup>9</sup>FE Allen and J Cocke. *A catalogue of optimizing transformations*. 1971. .

# Dead Block Elimination

- ▶ CFG not necessarily connected
- ▶ E.g., consequence of optimization
  - ▶ Conditional branch → unconditional branch
- ▶ Removing dead blocks is trivial
  1. DFS traversal of CFG from entry, mark visited blocks
  2. Remove unmarked blocks

# Optimization Example 1

```
define i32 @fac(i32 %0) {
  br label %for.header
for.header: ; preds = %for.body, %1
  %a = phi i32 [ 1, %1 ], [ %a.new, %for.body ]
  %b = phi i32 [ 0, %1 ], [ %b.new, %for.body ]
  %i = phi i32 [ 0, %1 ], [ %i.new, %for.body ]
  %cond = icmp sle i32 %i, %0
  br i1 %cond, label %for.body, label %exit
for.body: ; preds = %for.header
  %a.new = mul i32 %a, %i
  %b.new = add i32 %b, %i
  %i.new = add i32 %i, 1
  br label %for.header
exit: ; preds = %for.header
  %absum = add i32 %a, %b
  ret i32 %a
}
```

# Simple Dead Code Elimination (DCE)

- ▶ Look for trivially dead instructions
    - ▶ No users or side-effects
    - ▶ Calls *might* be removed
1. Add all instructions to work queue
  2. While work queue not empty:
    - 2.1 Check for deadness
    - 2.2 If dead, remove and add all operands to work queue

**Warning:** Don't implement it this naively, this is inefficient

# Applying Simple DCE

```
define i32 @fac(i32 %0) {
eff.: cf   br label %for.header
for.header: ; preds = %for.body, %1
users: 1   %a = phi i32 [ 1, %1 ], [ %a.new, %for.body ]
users: 1   %b = phi i32 [ 0, %1 ], [ %b.new, %for.body ]
users: 4   %i = phi i32 [ 0, %1 ], [ %i.new, %for.body ]
users: 1   %cond = icmp sle i32 %i, %0
eff.: cf   br i1 %cond, label %for.body, label %exit
for.body: ; preds = %for.header
users: 1   %a.new = mul i32 %a, %i
users: 1   %b.new = add i32 %b, %i
users: 1   %i.new = add i32 %i, 1
eff.: cf   br label %for.header
exit: ; preds = %for.header

eff.: cf   ret i32 %a
}
```

# Dead Code Elimination

- ▶ Problem: unused value cycles
  - ▶ Idea: find “value sinks” and mark all needed values as live
    - ▶ Sink: instruction with side effects (e.g., store, control flow)
1. Only mark instrs. with side effects as live
  2. Populate work list with newly added live instrs.
  3. While work list not empty:
    - 3.1 Mark dead operand instructions as live and add to work list
  4. Remove instructions not marked as live

# Applying Liveness-based DCE

```
define i32 @fac(i32 %0) {  
  live   br1 label %for.header  
for.header: ; preds = %for.body, %1  
  live   %a = phi i32 [ 1, %1 ], [ %a.new, %for.body ]  
  
  live   %i = phi i32 [ 0, %1 ], [ %i.new, %for.body ]  
  live   %cond = icmp sle i32 %i, %0  
  live   br2 i1 %cond, label %for.body, label %exit  
for.body: ; preds = %for.header  
  live   %a.new = mul i32 %a, %i  
  
  live   %i.new = add i32 %i, 1  
  live   br3 label %for.header  
exit: ; preds = %for.header  
  
  live   ret i32 %a  
}
```

Work list (stack)



## Optimization Example 2

```
define i32 @foo(i32 %0, ptr %1, ptr %2) {  
    %4 = zext i32 %0 to i64  
    %5 = getelementptr inbounds i32, ptr %1, i64 %4  
    %6 = load i32, ptr %5, align 4  
    %7 = zext i32 %0 to i64  
    %8 = getelementptr inbounds i32, ptr %2, i64 %7  
    %9 = load i32, ptr %8, align 4  
    %10 = add nsw i32 %6, %9  
    ret i32 %10  
}
```

# Common Subexpression Elimination (CSE) – Attempt 1

- ▶ Idea: find/eliminate redundant computation of same value
- ▶ Keep track of previously seen values in hash map
- ▶ Iterate over all instructions
  - ▶ If found in map, remove and replace references
  - ▶ Otherwise add to map
- ▶ Easy, right?

# CSE Attempt 1 – Example 1

```
define i32 @foo(i32 %0, ptr %1, ptr %2) {  
→ ht      %4 = zext i32 %0 to i64  
→ ht      %5 = getelementptr inbounds i32, ptr %1, i64 %4  
→ ht      %6 = load i32, ptr %5, align 4  
dup %4    %7 = zext i32 %0 to i64  
→ ht      %8 = getelementptr inbounds i32, ptr %2, i64 %7%4  
→ ht      %9 = load i32, ptr %8, align 4  
→ ht      %10 = add nsw i32 %6, %9  
→ ht      ret i32 %10  
}
```

- ▶ Obsolete instr. can be killed immediately, or in a later DCE

## CSE Attempt 1 – Example 2

```
define i32 @square(i32 %a, i32 %b) {  
  entry:  
→ ht    %cmp = icmp slt i32 %a, %b  
→ ht    br i1 %cmp, label %if.then, label %if.end  
  if.then: ; preds = %entry  
→ ht    %add1 = add i32 %a, %b  
→ ht    br label %if.end  
  if.end: ; preds = %if.then, %entry  
→ ht    %condvar = phi i32 [ %add1, %if.then ], [ %a, %entry ]  
dup %add1 %add2 = add i32 %a, %b  
→ ht    %res = add i32 %condvar, %add2%add1  
→ ht    ret i32 %res  
}
```

---

Instruction does not dominate all uses!

error: input module is broken!

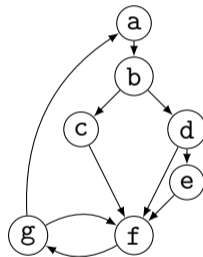
# Domination

- ▶ Remember: CFG  $G = (N, E, s)$  with digraph  $(N, E)$  and entry  $s \in N$
  - ▶ Dominate:  $d \text{ dom } n$  iff every path from  $s$  to  $n$  contains  $d$ 
    - ▶ Dominators of  $n$ :  $DOM(n) = \{d \mid d \text{ dom } n\}$
  - ▶ Strictly dominate:  $d \text{ sdom } n \Leftrightarrow d \text{ dom } n \wedge d \neq n$
  - ▶ Immediate dominator:  
 $\text{idom}(n) = d : d \text{ sdom } n \wedge \nexists d'. d \text{ sdom } d' \wedge d' \text{ sdom } n$
- $\Rightarrow$  All strict dominators are always executed before the block
- $\Rightarrow$  All values from dominators available/usable
- $\Rightarrow$  All values not from dominators **not** usable

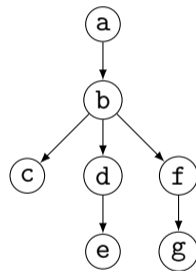
# Dominator Tree

- ▶ Tree of immediate dominators
- ▶ Allows to iterate over blocks in pre-order/post-order
- ▶ Answer  $a \text{ sdom } b$  quickly

Control Flow Graph





Dominator Tree



# Dominator Tree: Construction

- ▶ Naive: inefficient (but reasonably simple)<sup>10</sup>
  - ▶ For each block: find a path from the root – superset of dominators
  - ▶ Remove last block on path and check for alternative path
  - ▶ If no alternative path exists, last block is idom
  
- ▶ Lengauer–Tarjan: more efficient methods<sup>11</sup>
  - ▶ Simple method in  $\mathcal{O}(m \log n)$ ; sophisticated method in  $\mathcal{O}(m \cdot \alpha(m, n))$   
( $\alpha(m, n)$  is the inverse Ackermann function, grows *extremely* slowly)
  - ▶ Used frequently in compilers<sup>12</sup>

<sup>10</sup>ES Lowry and CW Medlock. “Object code optimization”. In: *CACM* 12.1 (1969), pp. 13–22. 

<sup>11</sup>T Lengauer and RE Tarjan. “A fast algorithm for finding dominators in a flowgraph”. In: *TOPLAS* 1.1 (1979), pp. 121–141. 

<sup>12</sup>Example: <https://github.com/WebKit/WebKit/blob/aabfacb/Source/WTF/wtf/Dominators.h>

# Dominator Tree: Implementation

- ▶ Per node store: *idom*, idom-children, DFS pre-order/post-order number
- ▶ Get immediate dominator: ...lookup *idom*
- ▶ Iterate over all dominators/-dominated by: ...trivial
- ▶ Check whether  $a$  sdom  $b$ <sup>13</sup>
  - ▶  $a.preNum < b.preNum \wedge a.postNum > b.postNum$
  - ▶ After updates, numbers might be invalid: recompute or walk tree
- ▶ Problem: dominance of unreachable blocks ill-defined  $\rightsquigarrow$  special handling



## CSE Attempt 2

- ▶ Option 1:
  - ▶ For identical instructions, store all
  - ▶ Add dominance check before replacing
  - ▶ Visit nodes in reverse post-order (i.e., topological order)
  
- ▶ Option 2:<sup>14</sup>
  - ▶ Do a DFS over dominator tree
  - ▶ Use scoped hashmap to track available values

Does this work? Yes.

## CSE: Hashing an Instruction (and Beyond)

- ▶ Needs hash function *and* “relaxed” equality
- ▶ Idea: combine opcode and operands/constants into hash value
  - ▶ Use pointer or index for instruction result operands
- ▶ Canonicalize commutative operations
  - ▶ Order operands deterministically, e.g., by address
- ▶ Identities:  $a+(b+c)$  vs.  $(a+b)+c$

## Global Value Numbering – or: advanced CSE

- ▶ Hash-based approach only catches trivially removable duplicates
- ▶ Alternative: partition values into *congruence classes*
  - ▶ Congruent values are guaranteed to always have the same value
- ▶ Optimistic approach: values are congruent unless proven otherwise
- ▶ Pessimistic approach: values are not congruent unless proven
- ▶ Combinable with: reassociation, DCE, constant folding
- ▶ Rather complex, but can be highly beneficial<sup>15</sup>

<sup>15</sup>K Gargi. "A sparse algorithm for predicated global value numbering". In: *PLDI. 2002*, pp. 45–56.

# Simple Transformations: Inlining

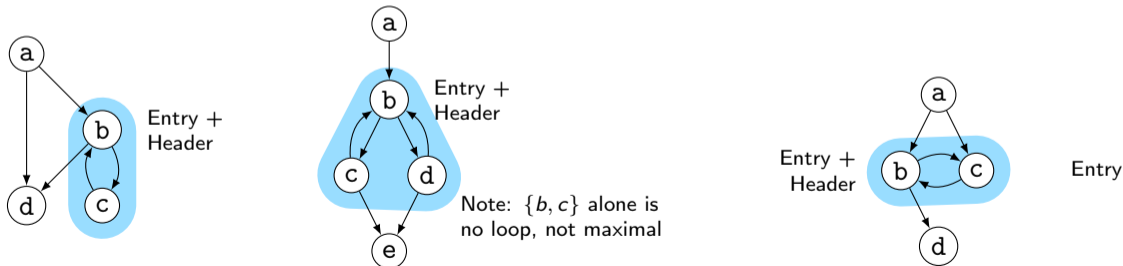
- ▶ Estimate whether inlining is beneficial
  - ▶ Savings of avoided call/computations/branches; cost of increased size
- ▶ Copy original function in place of the call
  - ▶ Split basic block containing function call
- ▶ Replace returns with branches and  $\phi$ -node to/at continuation point
- ▶ Move `alloca` to beginning or save stack pointer
  - ▶ Prevent unbounded stack growth in loops
  - ▶ LLVM provides `stacksave/stackrestore` intrinsics
- ▶ Exceptions may need special treatment

# Simple Transformations: Mem2Reg and SROA

- ▶ Mem2reg: promote `alloca` to SSA values/phis
  - ▶ Condition: only `load/store`, no address taken
  - ▶ Essentially just SSA construction
- ▶ SROA: scalar replacement of aggregate
  - ▶ Separate structure fields into separate variables
  - ▶ Also promote them to SSA

# Loops

- ▶ Loop: maximal SCC  $L$  with at least one internal edge<sup>16</sup>  
(strongly connected component (SCC): all blocks reachable from each other)
  - ▶ Entry: block with an edge from outside of  $L$
  - ▶ Header  $h$ : first entry found (might be ambiguous)
- ▶ Loop nested in  $L$ : loop in subgraph  $L \setminus \{h\}$



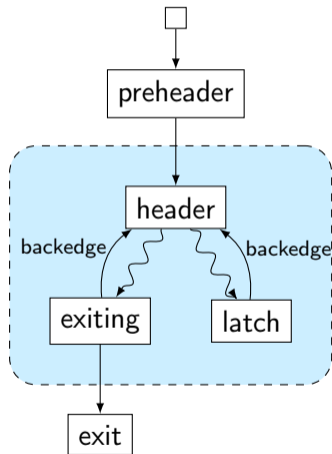
# Natural Loops

- ▶ Natural Loop: loop with single entry
  - ⇒ Header is unique
  - ⇒ Header dominates all block
  - ⇒ Loop is reducible
- ▶ Backedge: edge from block to header
- ▶ Predecessor: block with edge into loop
- ▶ Preheader: unique predecessor

## Formal Definition


Loop  $L$  is reducible iff  $\exists h \in L . \forall n \in L . h \text{ dom } n$


CFG is reducible iff all loops are reducible



# Finding Natural Loops

- ▶ Modified version<sup>17</sup> of Tarjan's algorithm<sup>18</sup>
- ▶ Iterate over dominator tree in post order
- ▶ Each block: find predecessors dominated by the block
  - ▶ None  $\rightsquigarrow$  no loop header, continue
  - ▶ Any  $\rightsquigarrow$  loop header, these edges *must* be backedges
- ▶ Walk through predecessors until reaching header again
  - ▶ All blocks on the way must be part of the loop body
  - ▶ Might encounter nested loops, update loop parent

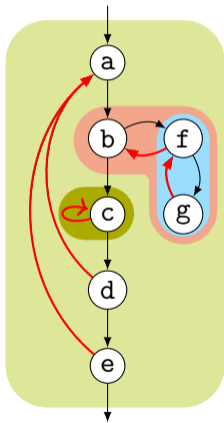
<sup>17</sup>G Ramalingam. "Identifying loops in almost linear time". In: *TOPLAS* 21.2 (1999), pp. 175–188. .

<sup>18</sup>R Tarjan. "Testing flow graph reducibility". In: *STOC*. 1973, pp. 96–107. .

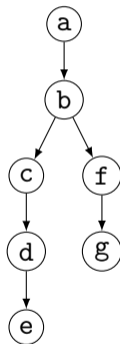


# Finding Natural Loops: Example

Control Flow Graph



Dominator Tree



Loop Info

Loop **A**:  $\{c\}$   
header:  $c$ ; parent:  $D$

Loop **B**:  $\{f, g\}$   
header:  $f$ ; parent:  $C$

Loop **C**:  $\{b, f, g\}$   
header:  $b$ ; parent:  $D$

Loop **D**:  $\{a, b, c, d, e, f, g\}$   
header:  $a$ ; parent:  $NULL$

# Loop Invariant Code Motion (LICM)

- ▶ Analyze loops, iterate over loop tree in post-order
  - ▶ I.e., visit inner loops first
- ↑ Hoist:<sup>19</sup> iterate over blocks of loop in reverse post-order
  - ▶ For each movable inst., check for loop-defined operands
  - ▶ If not, move to preheader (create one, if not existent)
  - ▶ Otherwise, add inst. to set of values defined inside loop
- ↓ Sink: Iterate over blocks of loop in post-order
  - ▶ For each movable inst., check for users inside loop
  - ▶ If none, move to unique exit (if existent)

<sup>19</sup><https://github.com/bytecodealliance/wasmtime/blob/bd6fe11/craneflift/codegen/src/licm.rs>

# Transformations and Analyses in LLVM: Passes

- ▶ Transformations and analyses organized in *passes*
- ▶ Pass can operate on Module/(CGSCC)/Function/Loop
- ▶ Analysis pass: takes input IR and returns analysis result
  - ▶ May also use results of other analyses; results are cached
- ▶ Transformation pass: takes input IR and returns preserved analyses
  - ▶ Can use analyses, which are re-run when outdated
- ▶ Pass manager executes passes on same granularity
  - ▶ Otherwise, use adaptor: `createFunctionToLoopPassAdaptor`  
(and preferably combine multiple smaller passes into a separate pass manager)

## Using LLVM (New) Pass Manager

```
void optimize(llvm::Function* fn) {
    llvm::PassBuilder pb;
    llvm::LoopAnalysisManager lam{};
    llvm::FunctionAnalysisManager fam{};
    llvm::CGSCCAnalysisManager cgam{};
    llvm::ModuleAnalysisManager mam{};
    pb.registerModuleAnalyses(mam);
    pb.registerCGSCCAnalyses(cgam);
    pb.registerFunctionAnalyses(fam);
    pb.registerLoopAnalyses(lam);
    pb.crossRegisterProxies(lam, fam, cgam, mam);

    llvm::FunctionPassManager fpm{};
    fpm.addPass(llvm::DCEPass());
    fpm.addPass(llvm::createFunctionToLoopPassAdaptor(llvm::LoopRotatePass()));
    fpm.run(*fn, fam);
}
```

# Writing a Pass for LLVM's New PM – Part 1

```
#include "llvm/IR/PassManager.h"
#include "llvm/Passes/PassBuilder.h"
#include "llvm/Passes/PassPlugin.h"

class TestPass : public llvm::PassInfoMixin<TestPass> {
public:
    llvm::PreservedAnalyses run(llvm::Function &F,
                               llvm::FunctionAnalysisManager &AM) {
        // Do some magic
        llvm::DominatorTree *DT = &AM.getResult<llvm::DominatorTreeAnalysis>(F);
        // ...
        llvm::errs() << F.getName() << "\n";
        return llvm::PreservedAnalyses::all();
    }
};
// ...
```

## Writing a Pass for LLVM's New PM – Part 2

```
extern "C" ::llvm::PassPluginLibraryInfo LLVM_ATTRIBUTE_WEAK
llvmGetPassPluginInfo() {
    return { LLVM_PLUGIN_API_VERSION, "TestPass", "v1",
            [] (llvm::PassBuilder &PB) {
                PB.registerPipelineParsingCallback(
                    [] (llvm::StringRef Name, llvm::FunctionPassManager &FPM,
                        llvm::ArrayRef<llvm::PassBuilder::PipelineElement>) {
                        if (Name == "testpass") {
                            FPM.addPass(TestPass());
                            return true;
                        }
                        return false;
                    });
            } });
}
```

```
c++ -shared -o testpass.so testpass.cc -lLLVM -fPIC
```

```
opt -load-pass-plugin=$PWD/testpass.so -passes=testpass input.ll | llvm-dis
```

# Analyses and Transformations – Summary

- ▶ Program Transformation critical for performance improvement
- ▶ Code not necessarily better
- ▶ Analyses are important to drive transformations
  - ▶ Dominator tree, loop detection, value liveness
- ▶ Important optimizations
  - ▶ Dead code elimination, common sub-expression elimination, loop-invariant code motion
- ▶ Compilers often implement transformations as passes
- ▶ Analyses may be invalidated by transformations, needs tracking

# Analyses and Transformations – Questions

- ▶ Why is “optimization” a misleading name for a transformation?
- ▶ How to find unused code sections in a function’s CFG?
- ▶ Why is a liveness-based DCE better than a simple, user-based DCE?
- ▶ What is a dominator tree useful for?
- ▶ What is the difference between an irreducible and a natural loop?
- ▶ How to find natural loops in a CFG?
- ▶ How does the algorithm handle irreducible loops?
- ▶ Why is sinking a loop-invariant inst. harder than hoisting?