

A fast algorithm for global code motion of congruent instructions

Anonymous Author(s)

Abstract

We present a fast global-code-motion (GCM) compiler optimization which schedules congruent instructions across the program. We introduce data structures to help exploit flow of global values in sparse representation. Not only GCM improves performance it saves code size as well, it exposes more instruction level parallelism in the basic-block to which instructions are moved, and it enables other compiler optimizations like loop invariant motion to remove redundancies. The cost model to drive the code motion is based on liveness analysis on SSA representation such that the (virtual) register pressure does not increase resulting in 2% fewer spills on the SPEC CPU 2006 benchmark suite. GCM is an optimistic algorithm in the sense that it considers all congruent instructions in a function as potential candidates. We also formalize why register pressure reduces as a result of global-code-motion, and how global-code-motion increases instruction level parallelism thereby enabling more out-of-order execution on modern architectures.

We have implemented the pass in LLVM. The experimental results show performance gains of 8% in 462.libquantum and 5.3% in 471.omnetpp from the SPEC CPU 2006 benchmark suite. It also reduces compilation time by reducing number of instructions to be processed by later compilation stages. GCM enables more inlining and exposes more loop invariant code motion opportunities in the majority of benchmarks.

Keywords Optimizing Compilers, GCC, LLVM, Code Generation, Global Scheduling

1 Introduction

Compiler techniques to remove redundant computations are composed of an analysis phase that detects identical computations in the program and a transformation phase that reduces the number of run-time computations. Classical scalar optimizations like Common Subexpression Elimination (CSE) [Aho et al. 1986] work very well on single basic-blocks (local level) but when it comes to detect redundancies across basic-blocks (global level) these techniques fall short: more complex passes like Global Common Subexpression Elimination (GCSE) and Partial Redundancy Elimination (PRE) have been designed

to handle these cases based on data-flow analysis [Morel and Renvoise 1979]. At first these techniques were described in the classical data-flow analysis framework, but later the use of Static Single Assignment (SSA) representation lowered their cost in terms of compilation time [Briggs and Cooper 1994; Chow et al. 1997; Kennedy et al. 1999a] and brought these techniques in the main stream: SSA based PRE has been implemented in many optimizing compilers like GCC and LLVM.

This paper describes a fast algorithm for global-code-motion (GCM) of congruent instructions [Briggs et al. 1997], a technique that uses the information computed for PRE to detect identical computations but has a transformation phase whose goal differs from PRE: it removes identical computations from different branches of execution. These identical computations in different branches of execution are not redundant computations at run-time and the number of run-time computations is not reduced. It is not a redundancy elimination pass, and thus it has different cost function and heuristics than PRE or CSE. It is similar to global code scheduling [Aho et al. 1986; Click 1995] in the sense that it will only move dynamic computations. Code hoisting, for example, can reduce the critical path length of execution in out-of-order machines. As more instructions are available at the hoisting point, the hardware has more independent instructions to reorder. The following reduced example illustrates how hoisting can improve performance by exposing more ILP.

```
float fun(float d, float min, float max, float a) {  
    float tmin, tmax;  
    float inv = 1.0f / d;  
    if (inv >= 0) {  
        tmin = (min - a) * inv;  
        tmax = (max - a) * inv;  
    } else {  
        tmin = (max - a) * inv;  
        tmax = (min - a) * inv;  
    }  
    return tmax + tmin;  
}
```

In this program the computations of $tmax$ and $tmin$ are identical to the computations of $tmin$ and $tmax$ of sibling branch respectively. Both $tmax$ and $tmin$ depend on inv which depends on a division operation which is

generally more expensive than the addition, subtraction, and multiplication operations. The total latency of computation across each branch is: $C_{div} + 2(C_{sub} + C_{mul})$. For an out-of-order processor with two add units and two multiply units, the total latency of computation is $C_{div} + C_{sub} + C_{mul}$.

Now if the computation of $tmax$ and $tmin$ are hoisted outside the conditionals, the C code version would look like:

```
float fun(float d, float min, float max, float a) {
    float inv = 1.0f / d;
    float x = (min - a) * inv;
    float y = (max - a) * inv;
    float tmin = x, tmax = y;
    if (inv < 0) { tmin = y; tmax = x; }
    return tmax + tmin;
}
```

In this code the division operation can be executed in parallel with the two subtractions because there are no dependencies among them. So the total number of cycles will be $max(C_{div}, C_{sub}) + C_{mul} = C_{div} + C_{mul}$; since C_{div} is usually much greater than C_{sub} [ARM 2014; Intel 2000]. Although it decreases the instruction level parallelism (ILP) of basic blocks from where the instructions were removed, those basic blocks cannot execute for more cycles than the destination basic block. We have seen the performance of a proprietary benchmark (on an out-of-order processor) improve by 15%. In fact, there are several advantages of GCM:

- it helps reduce the code size of the program by replacing multiple instructions with one thereby making the function cheaper to inline because inliner heuristics depend on instruction count;
- it exposes more ILP to the later compiler passes. By hoisting identical computations to be executed earlier, instruction schedulers can move heavy computations earlier in order to avoid pipeline bubbles.
- it reduces branch misprediction penalty on out-of-order processors with speculative execution of branches: by hoisting or sinking expressions out of branches, it can effectively reduce the amount of code to be speculatively executed;
- it reduces interference or register pressure with an appropriate cost-model: we have used a cost model in our implementation which only hoists or sinks when the register pressure is reduced;
- it reduces the total number of instructions to be executed for SIMD architectures which execute all code in branches based on masking or predication [Karrenberg and Hack 2011];
- it may improve loop vectorization by reducing a loop with control flow to a loop with a single basic-block,

should all the instructions in a conditional get hoisted or sunk;

- it enables more loop invariant code motion (LICM): as LICM passes, in general, cannot effectively reason about instructions within conditional branches in the context of loops, code-motion is needed to move instructions out of conditional expressions and expose them to LICM.

The algorithm for GCM uses several common representations of the program that we shortly describe below:

- Dominance (DOM) and Post-Dominance (PDOM) relations [Aho et al. 1986] on a Control Flow Graph (CFG);
- DJ-Graph [Sreedhar 1995] is a data structure that augments the dominator tree with join-edges to keep track of data-flow in a program. We use DJ-Graph to compute liveness of variables as illustrated in [Das et al. 2012];
- Static Single Assignment (SSA) [Cytron et al. 1989a];
- Global Value Numbering (GVN) [Click 1995; Rosen et al. 1988]: to identify similar computations compilers use GVN. Each expression is given a unique number and the expressions that the compiler can prove to be identical are given the same number;
- Rank of an expression [Rosen et al. 1988]: The expressions are ordered according to the amount of dependency they have on other expressions. For example in the expression $a = b + c$, a will have higher rank than rank of b , and c .
- MemorySSA [Novillo 2007]: it is a factored use-def chain of memory operations that the compiler is unable to prove independent. MemorySSA accelerates the access to the alias analysis information.

Dominator tree, Global Value Numbers and Memory SSA are already available in LLVM but there is no facility to infer liveness of virtual registers. So we implemented DJ-Graph data structure that allowed us to calculate MergeSets which was used in the liveness analysis during GCM.

The GCM pass can be broadly divided into the following:

- factor the reachability of expressions (Section 3.1).
- find candidates (congruent instructions) suitable for global-code-motion (Section 4),
- compute a point in the program where it is both legal (Section 3.3) and profitable (Section 3.4) to move the code,
- move the code to hoist point or sink point (Section 4),
- update data structures to continue iterative global-code-motion (Section 4).

There has been a lot of work both in industry and academia related to code hoisting/sinking, and in general global scheduling [Click 1995]. Some relate code-hoisting

to code-size optimization [Rosen et al. 1988] and many [Barany and Krall 2013; Shobaki et al. 2013] use global scheduling to improve performance. Most of the recent work on global scheduling are done using integer linear programming (ILP) which results in prohibitively high compile time which is not suitable for industrial compilers. To the best of our knowledge we have not found any reference which explored global-code-motion with a cost-model to reduce register-pressure. The main contributions of this paper are:

- an optimistic fast algorithm to schedule congruent instructions;
- introducing an analogue of ϕ nodes which allows faster tracking of anticipable instructions;
- a cost model to reduce register spills;
- experimental evaluation of our implementation in LLVM.

2 Related Work

There are a lot of bug reports in GCC and LLVM bugzillas [GCC 2016a; LLVM 2016], showing the interest in having a more powerful code hoist transform. The current LLVM implementation of code hoisting in SimplifyCFG.cpp is very limited to hoisting from identical basic-blocks: the instructions of two sibling basic-blocks are read at the same time, and all the instructions of the blocks are hoisted to the common parent block as long as the compiler is able to prove that the instructions are equivalent. This implementation does not allow for an easy extension: first in terms of compilation time overhead the implementation is quadratic in number of instructions to bisimulate and second, the equivalence of instructions is computed by comparing the operands which is neither general nor scalable.

Click [1995] describes aggressive global-code-motion in sparse representation. It first schedules all the instructions as early as possible. This results in very long live ranges which is mitigated by again scheduling all the instructions as late as possible. They report a speedup of as high as 23% but there is no cost model to limit amount of scheduling and could often introduce redundant computations. It is also not mentioned how safety is ensured while hoisting an instruction when it is not anticipable. The legality checks are only based on barriers (Section 3.2) and availability of used operands. Even introducing a simple integer addition can also result in undefined behavior [Wang et al. 2013]. Except from this paper we found no reference of global scheduling in sparse representation.

Dhamdhare [1988] and Muchnick [1997] mention code hoisting in a data-flow framework. A list of Very Busy Expressions (VBE) computed which are hoisted in a basic-block where the expression is anticipable (all the operands are available.) This algorithm would hoist as

far as possible without regarding the impact on register pressure and as such a cost model will be required. Also the description of VBE is based on the classic data-flow model and an adaptation to a sparse SSA representation is required.

Rosen et al. [1988] also briefly discuss hoisting computations with identical value numbers from immediate successors. Their algorithm iterates on computations of same rank and move the code with identical computations from the sibling branch to a common dominator if they are very busy [Muchnick 1997]. However, the cost-model to mitigate register pressure is missing, also, there is no mention of sinking congruent instructions.

GCC recently got code-hoisting [GCC 2016b] which is implemented as part of GVN-PRE: it uses the set of `ANTIC_IN` and `AVAIL_OUT` value expressions computed for PRE. `ANTIC_IN[B]` contains very busy expressions (VBEs) at basic-block B, i.e., values computed on all paths from B to exit and `AVAIL_OUT[B]` contains values which are already available. The algorithm hoists VBEs to a predecessor. It uses `ANTIC_IN[B]` to know what expressions will be computed on every path from the basic block B to exit, and can be computed in B. It uses `AVAIL_OUT[B]` to subtract out those values already being computed. The cost function is: for each hoist candidate, if all successors of B are dominated by B, then we know insertion into B will eliminate all the remaining computations. It then checks to see if at least one successor of B has the value available. This avoids hoisting it way up the chain to `ANTIC_IN[B]`. It also checks to ensure that B has multiple successors, since hoisting in a straight line is pointless. The algorithm continues top down the dominator tree, working in tandem with PRE until no more hoisting is possible. One advantage of GCC implementation is that it works in sync with the GVN-PRE such that when new hoisting opportunities are created by GVN-PRE, code-hoisting will hoist them.

Barany and Krall [2013] presented a global scheduler with ILP formulation with a goal to minimize register pressure. The results they got were not very promising. It may be because they only used the scheduler for smaller functions (< 1000 instructions); also, they compiled the benchmarks for ARM-Cortex which is more resilient to register pressure because it has more registers compared to X86, for example.

Shobaki et al. [2013] also presented a combinatorial global scheduler with reasonable performance improvements. It is possible that both Barany and Shobaki's implementation will have similar results when compiled for same target architecture. Also, both suffer from the same problem, although Shobaki not so much, of large compile times which would not suit in industrial compilers like GCC and LLVM. With the algorithm described

in the next section, we got a reduction in register spills, improved inliner heuristics, improved compile time, and show performance improvements on SPEC CPU 2006 benchmarks.

3 Global Code Motion of congruent computations

The GCM uses the following data structures (χ and Φ nodes), that are introduced below. They allow us to factor the reachability of values in an efficient way. Together with liveness analysis using the DJ-Graph [Das et al. 2012] and Memory SSA, the entire implementation of global-code-motion only uses sparse representations for efficient query and update.

3.1 Factoring reachability of expressions

Our algorithm factors out the reachability of values such that multiple queries to find reachability of values are fast. This is based on finding the dominance frontiers (DF) and the post-dominance frontiers (PDF) of the CFG. These points, being structural properties of a graph, do not change during the code-motion. A post-dominance frontier is the basic block where anticipability (ANTIC) of an instruction (I) is checked because it is the basic block on which the basic block containing I is control dependent.

So we introduce a data structure (χ nodes) to keep track of values flowing out of a basic block with more than one successor. This approach allows us to hoist instructions to a basic block with multiple successors. The χ node in a basic block B has CHIArgs with three entries: a) the instruction (I) CHIArg refers to, b) the basic block $Dest$ which is post-dominated by the basic block containing I, and c) the GVN of I.

```

struct CHIArg {
    // Edge destination (shows the direction of flow),
    // may not be where the I is.
    BasicBlock *Dest;
    // The instruction CHI-arg refers to.
    Instruction *I;
    // The global value number of I.
    GVNType VN;
};
    
```

A χ keeps information about ‘congruent’ values flowing out of a basic block. It is similar to ϕ [Cytron et al. 1989b] but in the inverse graph, and used for tracking out-flowing values on each edge. An illustration of χ is shown in Figure 1. The GVN for both I1 and I2 is V, the χ node will save the instruction as well as the edge where the value is flowing to. Where $Arg1$, $Arg2$ are of type CHIArg. The χ for both I1 and I2 is inserted in their PDF (A).

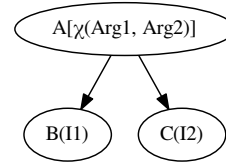


Figure 1. I1, I2 are instructions in B, and C respectively, $Arg1 = \{B, I1, V\}$ and $Arg2 = \{C, I2, V\}$, V is the value number for both I1, and I2

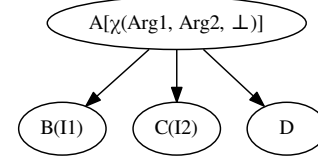


Figure 2. I1, I2 are instructions in B, and C respectively, $Arg1 = \{B, I1, V\}$ and $Arg2 = \{C, I2, V\}$, V is the value number for both I1, and I2. A has missing entry in χ so V is not anticipable

A similar idea has been used in SSAPRE [Kennedy et al. 1999b] for partial redundancy elimination in sparse form. When SSAPRE inserts Φ (expression PHIs), it does so at places where (potentially) multiple expressions may merge. If a Φ has a missing entry, that means the expression is partially available. Subsequent algorithms work on Φ s to find out redundancies. We extend the idea for GCM by propagating the anticipability in both directions.

For code hoisting we insert a χ , to a basic block with multiple successors. Then for each instruction, we register its flow of value to its corresponding χ which can be found in its PDF. Finally, we start walking the CFG and remove the χ if it has a missing entry (see Figure 2), that means the expression is not anticipable in the basic block having incomplete χ . For the remaining χ s anticipability of ‘values’ is ensured by construction. For code-sinking we insert Φ s in basic blocks with multiple predecessors. Then for each instruction we register its flow of value to its corresponding Φ which can be found in its DF. Finally, we start walking the CFG and remove the Φ if it has a missing entry (see Figure 3), that means the expression is not available in the basic block having incomplete Φ . For the remaining Φ s availability of ‘values’ is ensured by construction.

In order to limit the number of spurious χ s and Φ s, [Park and Lee 2008], we insert them only for values with multiple occurrences in the function as they are the interesting candidates for this GCM. We believe this factorization of anticipability/availability of expressions

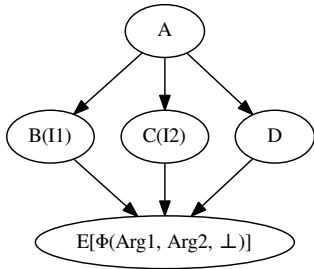


Figure 3. I1, I2 are instructions in B, and C respectively, $\text{Arg1} = \{B, I1, V\}$ and $\text{Arg2} = \{C, I2, V\}$, V is the value number for both I1, and I2. A has missing entry in Φ so V is not available

may be suitable for faster implementation of other global (across basic blocks) analyses/optimizations as well. Before we consider factors affecting the legality of GCM, we introduce the concept of barriers.

3.2 Barriers

Barriers are based on the concept of pinned instructions [Click 1995] but extended to adapt to LLVM IR. A basic-block in LLVM IR is actually an extended basic-block [Hwu et al. 1993] because there might be non-returning calls in the middle of a basic-block. There might be other instructions like stores, where compiler may not guarantee progress of control flow, all such instructions are marked as barriers. A barrier will prevent any instruction occurring lexically below it from getting hoisted. On the other hand, a sink barrier will prevent any instruction above it from getting sunk. If there are no barriers in a basic block, any instruction satisfying other legality checks can be sunk.

// Compute both hoist and sink barriers

void computeBarriers()

for each basic-block B in a function:

barrier_found = false

for each instruction I in B:

if I does not guarantee progress:

mark I as a barrier instruction

barrier_found = true

break;

// Find the last barrier below which instructions

// can be sunk.

if barrier_found:

for each instruction I in B in reverse order:

if I does not guarantee progress:

mark I as a sink barrier

break

Algorithm 1. Compute hoist and sink barriers

In the absence of context, as in our current implementation, some instructions which might be safe are still classified as barriers, e.g., calls with missing attributes. Computing barriers allows for efficient removal of non-movable candidates, this makes the legality check converge faster.

3.3 Legality of movable instructions

Since the equality of movable candidates is purely based on the value numbers, we also need to establish if hoisting them or sinking them would be legal. The completeness of χ and Φ are necessary for GCM but not sufficient. This is because congruence of value numbers only implies that they compute same value if the inputs are same. So for instructions like memory references, they may compute different values should the memory get modified along the path. For that we need to check for intersecting side-effects along the path.

We need to prove that on all execution paths from the initial position of an instruction to its destination, the side effects appear in the same order. In general, we cannot introduce a new computation along any path of execution without checking for undefined behavior [Wang et al. 2013]. It is also necessary to check if there are indirect branch targets, e.g., landing pad, case statements, or other forms of goto labels along the path because it becomes difficult to prove safety checks in those cases.

Scalars are the easiest to hoist because we only need to check for availability of operands. For sinking the scalars we need to make sure that there are no uses of the instruction along the path. For instructions with memory references like loads, stores and calls, Memory SSA allows efficient way of iterating over the use-def chains of memory references on a factored graph. Hoisting loads/stores across calls also requires precise analysis of all the memory addresses accessed by the call. The current implementation being an intraprocedural pass, cannot schedule aggressively across calls. In the presence of pure calls, loads can be hoisted but stores can't. Also, if a call throws exceptions, or if it may not return, nothing can be scheduled across that call.

For sinking, higher ranked expressions would be sunk first. And it would be illegal to sink higher ranked identical expressions if they are not fully anticipable in the common post-dominator. For example:

B0: i0 = load B

if p is true then goto B1 else goto B2

B1: i1 = load A

c1 = i1 + 10

d1 = i0 + 20

```

551      goto B3
552 B2: i2 = load A
553      c2 = i2 + 10
554      d2 = i0 + 20
555      goto B3
556

```

```

557 B3: PHIc1, c2)
558     PHI(d1, d2)
559

```

In this example (c1, c2) or (d1, d2) are potential sinkable candidates. (c1, c2) depend on i1 and i2 respectively which are also in their original basic blocks; c1 and c2 are not fully anticipable in B3. So without knowing the ability to sink 'i1' and 'i2' it would be illegal to sink (c1, c2) to B3. On the other hand (d1, d2) can safely be sunk because their operands are readily available at the sink point, i.e., B3. It should also be noted that, even if the expressions are identical and operands are available, it still requires a both instructions to be post-dominated by a unique Φ to be legally sinkable.

3.4 Profitability check (Cost models)

After the legality checks have passed, we check if a GCM is profitable. That takes into account the impact GCM would have on various parameters that might affect runtime performance, e.g., impact on live-range, gain in the code size. The current implementation makes effort to not regress in performance and still reduce register spills. We implement a cost model to reduce register pressure.

3.4.1 Reduce register pressure

Hoisting upwards will decrease the live-range of its use, if it is a last use (a kill) but increase the live-range of its definition. Conversely, sinking will decrease the live-range of the defined register but increase the live-range for killed operands. If the live-range after GCM will decrease, it will be moved. Essentially, if there is one killed operand, code hoisting will either decrease or preserve the register pressure as long as the hoisting of killed operands does not cross a use. Similarly, code-sinking (of an instructions which defines a register) will either decrease or preserve the register pressure as long as there is one operand killed at most and the instruction is not sunk below the killing use of a non-killed operand. The following example explains how GCM can reduce the register pressure.

```

598 B0:
599     b = m
600     c = n
601     if p is true then goto B1 else goto B2
602 B1:
603     a0 = b<kill> + c<kill>
604

```

```

        goto B3
B2:
    a1 = b<kill> + c<kill>
    goto B3
B3: ...

```

After hoisting, $a0$ and $a1$ are removed and a copy of $a0$ is placed in $B0$. Since b and c are killed in both $a0$ and $a1$, hoisting the expressions will reduce the register pressure because two registers will be freed at the insertion point but only one register will be required to assign the definition of hoisted instruction $a0$. We have reduced performance regressions with this cost model even if LLVM has a live-range splitting [Cooper and Simpson 1998] optimization. Sinking may also reduce the register pressure when the use operands are not kills.

3.4.2 Hoisting an expression across a call

Call expressions are difficult to analyze in an intraprocedural optimization pass like GCM. Even hoisting scalars across calls is tricky because it can increase the number of spills. During the register allocation of call expressions, the argument registers (also known as the caller saved registers) are saved because they might be modified by the callee and after the call they are restored. Due to this frame setup the register pressure is high because the number of available registers are reduced by the number of caller saved registers. In that situation a computation that reduces register pressure may be beneficial to hoist. Similarly, a computation that introduces a register (without killing any other operand) before a call may be beneficial to sink.

4 A fast algorithm for code motion of congruent computations

Once all the legality and profitability checks are satisfied for a set of congruent instructions, they are suitable candidates for code motion. A copy of the computation is inserted at the hoisting/sinking point along with any instructions which needed to be rematerialized. Thereafter, all the computations made redundant by the new copy are removed and affected data structures are updated.

```
void doGCM()
```

```
    Analyses available:
```

```
        Dominator Tree, DFS Numbering, Memory SSA
```

```
    Compute DJ-graph of function
```

```
    constructMergeSet(CFG)
```

```
    computeBarriers()
```

```
    do:
```

```
        Compute GVN of each expression in the function
```

```
        // Repeat hoisting if any candidates were hoisted.
```

```
    while (doCodeHoisting() > 0)
```

```

661 // Code-sinking only once
662 doCodeSinking()

```

Algorithm 2. Global Code Motion

GCM basically consists of two parts, i.e., hoisting and sinking. This implementation only moves congruent instructions which makes it easier to ensure that we do not introduce a redundant computation along any path of execution. If the dependency of a hoistable candidate is in the same basic-block as the candidate, then the dependency must also be hoistable, otherwise hoisting will be illegal or would require a complicated code generation to make it legal. The current algorithm discards cases if the dependency is neither hoistable nor rematerializable. We collect the GVN of all the instructions in the function and iterate on the list of instructions having identical GVNs. The algorithm *doGCM* prefers hoisting (*doCodeHoisting*) to sinking (*doCodeSinking*).

```

679 int doCodeHoisting(VNs)
680 // VNs = map (VN -> list of I with value VN)
681 Sort VNs in ascending order
682
683
684 for each VN in VNs:
685   for each instruction with the same VN:
686     Find its post-dominance-frontier (PDF)
687     Insert the CHI-node at PDF
688
689 for each  $\chi$ -nodes in the CFG:
690   for each  $\chi$ -arg in  $\chi$ -node:
691     Remove  $\chi$ -arg if it fails legality checks
692
693
694 for each  $\chi$ -nodes in the CFG:
695   if  $\chi$  does not have missing entries:
696     Perform code hoisting
697     update data structures and statistics
698
699 return number of hoisted instructions

```

Algorithm 3. Hoist instructions

Code hoisting opens new opportunities for other higher ranked [Rosen et al. 1988] (depended on candidates which got hoisted) which are encountered subsequently as the algorithm visits instructions in the increasing order of their ranks. Some of the congruent higher ranked computations do not get equal value numbers so we rerun the code-hoisting algorithm until there are no more instructions left to be hoisted. The hoisting converges quickly with 2-3 iterations.

```

712 int doCodeSinking()
713 // VNs = map (VN -> list of I with value VN)
714 Sort VNs in ascending order

```

```

716
717 for each VN in VNs:
718   for each instruction with the same VN:
719     Find its dominance-frontier (DF)
720     Insert the  $\Phi$ -node at DF
721
722
723 for each  $\Phi$ -nodes in the CFG:
724   for each  $\Phi$ -arg in PHI-node:
725     Remove  $\Phi$ -arg if it fails legality checks
726
727
728 for each  $\Phi$ -nodes in the CFG:
729   if  $\Phi$  does not have missing entries:
730     Perform code sinking
731     if a  $\phi$  has same entries, reassign
732     update data structures and statistics
733
734 return number of sinked instructions

```

Algorithm 4. Sink instructions

After no more candidates are hoistable, sinking is performed. It is only performed once because it may not open new opportunities for sinkable candidates. For sinking we iterate on the Φ nodes. After each code-motion the SSA form is restored by updating the intermediate representation (IR) to reflect the changes. At the same time MemorySSA is updated as well. Both the data structures can be updated very efficiently. After one full iteration of hoisting, the GVN needs to be recomputed to explore new hoistable candidates. Finally after the transformation is done, we verify a set of post-conditions to establish that program invariants are maintained: e.g., consistency of use-defs, and SSA semantics.

4.1 Time complexity of algorithm

The complexity of code hoisting is linear in the number of instructions that are candidates for global-code-motion, matching the complexity of PRE on SSA form. The analyses computed for this pass like Global Value Numbering, Marking Barriers, are both linear in number of instructions in a function. Computation of DJ-graph and merge-sets is $O(E)$, E being number of edges in the CFG. Liveness analysis is not very expensive [Das et al. 2012] but still only performed on-demand for hoistable candidates. Other analyses like Alias Analysis, MemorySSA and Dominator Tree are already available in the LLVM pass pipeline.

Although we recompute GVN, and Barriers for each iteration of the global-code-motion, there is still gain in the compilation times (see Section 4). We have also provided appropriate compiler flags to expedite global-code-motion by bailing out with fewer iterations, or skip

the liveness based profitability analysis to aggressively move the code as long as they are legal.

5 Experimental Evaluation

For evaluation of global-code-motion, we built SPEC CPU 2006 with and without the GCM optimization. All the experiments were conducted on an x86_64-linux machine and at -Ofast optimization level with inlining and loop unrolling disabled; inliner and loop-unroller vary wildly with GCM optimization so it was difficult to visualize the impact. Each benchmark was run three times and the best result was taken. We collected execution time and code-size which is listed in Table 1, other compiler statistics are listed in Table 2, Table 3, and Figure 4. We also measured compile time as shown in Table 4. The benchmark 464.h264ref caused an internal compiler error with the trunk clang so it could not be reported.

Benchmarks	% performance uplift at -Ofast (high better)	size reduction (Bytes) at -Ofast (high better)
400.perlbench	1.95	984
401.bzip2	-0.84	16
403.gcc	-2.30	2128
429.mcf	-4.36	16
433.milc	-1.86	48
444.namd	-2.40	544
445.gobmk	-0.52	56
447.dealII	0.39	-2560
450.soplex	-0.24	-32
453.povray	-0.05	848
456.hmmer	0.93	296
458.sjeng	-0.58	232
462.libquantum	8.17	-16
470.lbm	2.57	0
471.omnetpp	5.30	152
473.astar	-0.17	48
482.sphinx3	-0.11	16
483.xalancbmk	-2.98	392

Table 1. Execution time (ratio) and code size change (Bytes) with and without GCM on SPEC CPU 2006

Overall, we can see some good improvements in 462.libquantum, 471.omnetpp, etc. and some regressions in 429.mcf, 483.xalancbmk, etc. The global-code-motion pass was run once in the pass pipeline (at -Ofast), after local common subexpression elimination (EarlyCSE) pass. Because EarlyCSE removes local redundancies, GCM would have to analyze less redundant instructions. We believe some performance regressions can be mitigated by running the GCM few more times.

We see some code size reduction in 400.perlbench, 403.gcc, and 453.povray as shown in Table 1. The code size listed here is the text size delta (from size command in linux) of the final executable for each benchmark. On the other hand 447.dealII increased in code size. This is because once the code-size of a function decreases (due to GCM) it becomes cheaper to do some optimizations

like inlining, loop unrolling, copy propagation, rematerialization etc. We tried to minimize the noise by disabling inliner, and loop unroller but still we can see some increase in the code size.

To see the impact of GCM on inliner we ran GCM with inliner enabled and collected the compile-time statistics produced by LLVM (See Table 2). Except for 447.dealII, 450.soplex, 482.sphinx3 and 483.xalancbmk all other benchmarks got more functions inlined. Also, since the pass runs early, it affects many optimizations which rely on the number of instructions, length of the use-def chain, and other metrics.

Benchmarks	Loop Invariant Motion		Functions	
	Hoisted%	Sunk%	Inlined%	Deleted%
400.perlbench	2.1	-2.8	0.2	0.9
401.bzip2	14.4	0.0	4.7	3.1
403.gcc	13.7	27.3	0.9	0.1
429.mcf	-7.3	0.0	0.0	0.0
433.milc	10.2	0.0	8.0	0.0
444.namd	0.1	0.0	0.8	-3.5
445.gobmk	0.6	5.6	5.4	1.7
447.dealII	62.2	-86.7	-0.2	-0.8
450.soplex	3.9	-5.6	-0.2	-0.4
453.povray	-4.8	39.1	0.2	-0.1
456.hmmer	1.0	166.7	0.0	-1.9
458.sjeng	11.5	0.0	6.3	0.0
462.libquantum	6.6	0.0	0.0	0.0
470.lbm	0.0	0.0	0.0	0.0
471.omnetpp	17.0	-6.3	1.1	-0.4
473.astar	9.1	0.0	1.1	0.0
482.sphinx3	5.9	-18.2	-2.0	0.0
483.xalancbmk	14.6	-3.6	-1.0	-1.2

Table 2. Change in the static compile time metrics w.r.t. GCM on SPEC CPU 2006 at -Ofast. The % columns show percentage increase in metric w.r.t. baseline, i.e., without GCM. GCM improves LICM and inlining in general.

Table 2, which shows how global-code-motion impacted important compiler optimizations like inlining and loop invariant code motion (LICM). LICM removes loop invariant code out of the loop. It may hoist the code to the loop's preheader or sink the code to loop's post-dominator [Muchnick 1997]. As we can see number of instructions hoisted increased significantly in several benchmarks, although number of instructions sunk did not change by much. Number of inlined functions also improved for most benchmarks. In general when more functions were inlined, more functions were deleted as in 401.bzip2, 445.gobmk.

Spills are listed in Table 3. Spills were reduced by an average of 2% on the SPEC CPU 2006 benchmark suite. This is a result of cost-model which limits the global-code-motion to those candidates which do not increase

Benchmarks	base(-Ofast)	GCM(-Ofast)	GCM/base
400.perlbench	2539	2493	0.98
401.bzip2	718	707	0.98
403.gcc	5782	5722	0.99
429.mcf	14	17	1.21
433.milc	635	642	1.01
444.namd	3223	3274	1.01
445.gobmk	2166	2173	1.00
447.dealII	11074	10234	0.92
450.soplex	1122	1125	1.00
453.povray	4701	4692	1.09
456.hmmer	1197	1274	1.06
458.sjeng	168	176	1.05
462.libquantum	118	118	1.00
470.lbm	33	33	1.00
471.omnetpp	524	527	1.00
473.astar	180	201	1.12
482.sphinx3	674	670	0.99
483.xalancbmk	5102	4965	0.97
Grand Total	43349	42497	0.98

Table 3. Number of spills with and without GCM on SPEC2006 at -Ofast. Lower is better in GCM/base. Spills reduced by 2%.

the register pressure (except for loads Section 3.4) overall spills still decreases.

Benchmarks	Baseline (in Millions)	GCM (in Millions)	%Decrease (lower is better)
400.perlbench	186	184	0.99
401.bzip2	76	75	0.99
403.gcc	1,079	1,073	0.99
429.mcf	78	77	0.99
433.milc	272	269	0.99
444.namd	78	77	0.99
445.gobmk	261	258	0.99
447.dealII	643	638	0.99
450.soplex	245	242	0.99
453.povray	499	494	0.99
456.hmmer	214	211	0.99
458.sjeng	89	88	0.99
462.libquantum	84	84	0.99
470.lbm	71	71	1.00
471.omnetpp	382	379	0.99
473.astar	78	77	0.99
482.sphinx3	158	156	0.99
483.xalancbmk	22,485	22,455	1.00

Table 4. Change in total number of instructions executed during the compilation with and without GCM on SPEC CPU 2006 at -Ofast. The instructions were counted using valgrind.

The compilation time actually reduced for most of the benchmarks as a result of code motion. This is because, as global-code-motion removes (congruent) instructions, the number of instructions to be processed by later compilation passes also reduces. The compilation time listed here is the total instruction count executed during the compilation of each benchmark as output from ‘valgrind --tool=cachegrind’ [Nethercote and Seward 2007]. We used valgrind because it was a simple and reliable way to show the impact of our optimization.

The compile time metrics of GCM are listed in Figure 4. The table lists the number of scalars, loads, stores and calls hoisted as well as removed. For each category, the number of instructions removed is greater or equal to the number of instructions hoisted because each code-motion is performed only when at least one identical computation is found. Loads are hoisted the most followed by scalars, stores and calls in decreasing order. This was the common trend in all our experiments. One reason why loads are hoisted the most is the early execution of this pass (before mem2reg pass which scalarizes some memory references) in the LLVM pass pipeline. Passes like mem2reg and instcombine might actually remove those loads and the number of hoisted loads may change should the GCM pass be scheduled later. We can see only a few sunk instructions because hoisting is performed before sinking (Section 4) and some of the candidates which can either be hoisted or sunk got hoisted.

6 Conclusions and Future Work

We have presented a fast algorithm to perform global-code-motion of congruent computations in SSA form. We introduced data structures like χ nodes, Φ nodes which helps exploit flow of values in sparse representation. Experimental results show that it is a useful performance and code-size optimization. It also enables inlining and LICM opportunities, reduces register spill. To preserve performance and not hoist/sink too much, we have implemented a register pressure aware cost model as described in Section 3.4. This pass does not hoist fully redundant expressions because they are already removed by a later partial redundancy elimination pass in LLVM.

We would like to explore other places in the optimization pipeline where GCM could improve performance. It may be useful to integrate GCM with LLVM’s new GVN-PRE implementation such that more congruent instructions can be scheduled. This would also expose more redundancies to the GVN-PRE. Also, since GCM runs very early in the pass pipeline, it will be good to evaluate the code size/performance impact when it is run in sync with GVN-PRE just like GCC does. With the implementation of GCM in LLVM, the passes which rely on the code-size or instruction-counts to make optimization decisions need to be revisited for example, the inliner, the loop unroller etc.

References

- Alfred V Aho, Ravi Sethi, and Jeffrey D Ullman. 1986. *Compilers, Principles, Techniques*. Addison Wesley.
- ARM. 2014. ARM Cortex-A57 Software Optimization Guide. (2014). http://infocenter.arm.com/help/topic/com.arm.doc.uan0015b/Cortex_A57_Software_Optimization_Guide_external.pdf

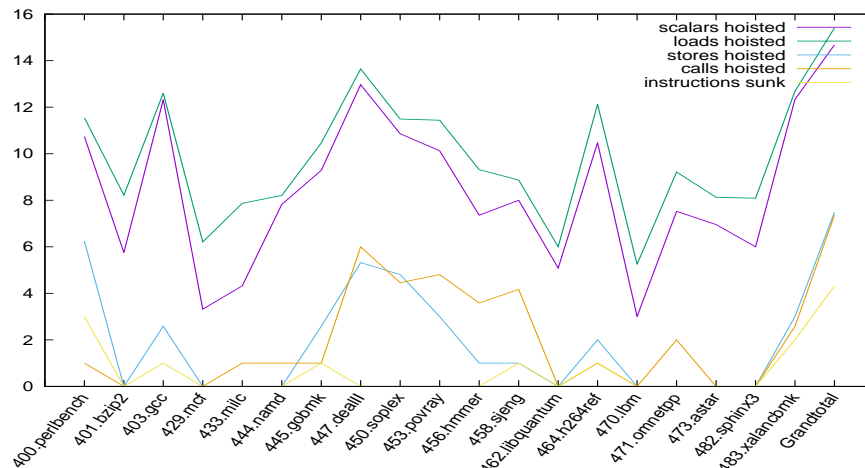


Figure 4. GCM stats on SPEC2006 at -Ofast. Loads are hoisted the most followed by scalars, stores and calls.

- Gergő Barany and Andreas Krall. 2013. Optimal and heuristic global code motion for minimal spilling. In *International Conference on Compiler Construction*. Springer, 21–40.
- Preston Briggs and Keith D Cooper. 1994. Effective partial redundancy elimination. In *ACM SIGPLAN Notices*, Vol. 29. ACM, 159–170.
- Preston Briggs, Keith D Cooper, and L Taylor Simpson. 1997. Value numbering. *Software Practice & Experience* 27, 6 (1997), 701–724.
- Fred Chow, Sun Chan, Robert Kennedy, Shin-Ming Liu, Raymond Lo, and Peng Tu. 1997. A new algorithm for partial redundancy elimination based on SSA form. In *ACM SIGPLAN Notices*, Vol. 32. ACM, 273–286.
- Cliff Click. 1995. Global code motion/global value numbering. In *ACM SIGPLAN Notices*, Vol. 30. ACM, 246–257.
- Keith D Cooper and L Taylor Simpson. 1998. Live range splitting in a graph coloring register allocator. In *International Conference on Compiler Construction*. Springer, 174–187.
- Ron Cytron, Jeanne Ferrante, Barry K Rosen, Mark N Wegman, and F Kenneth Zadeck. 1989a. An efficient method of computing static single assignment form. In *Proceedings of the 16th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. ACM, 25–35.
- Ron Cytron, Jeanne Ferrante, Barry K Rosen, Mark N Wegman, and F Kenneth Zadeck. 1989b. An efficient method of computing static single assignment form. In *Proceedings of the 16th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*. ACM, 25–35.
- Dibyendu Das, B Dupont De Dinechin, and Ramakrishna Upadrasta. 2012. Efficient liveness computation using merge sets and dj-graphs. *ACM Transactions on Architecture and Code Optimization (TACO)* 8, 4 (2012), 27.
- Dhananjay M Dhamdhere. 1988. A fast algorithm for code movement optimisation. *ACM SIGPLAN Notices* 23, 10 (1988), 172–180.
- GCC. 2016a. GCC bugs related to code hoisting. Bug IDs: 5738, 11820, 11832, 21485, 23286, 29144, 32590, 33315, 35303, 38204, 43159, 52256. (2016). <https://gcc.gnu.org/bugzilla>
- GCC. 2016b. GCC code hoisting implementation. (2016). <https://gcc.gnu.org/viewcvs/gcc?view=revision&revision=238242>
- Wen-Mei W Hwu, Scott A Mahlke, William Y Chen, Pohua P Chang, Nancy J Warter, Roger A Bringmann, Roland G Ouellette, Richard E Hank, Tokuzo Kiyohara, Grant E Haab, et al. 1993. The superblock: an effective technique for VLIW and superscalar compilation. In *Instruction-Level Parallelism*. Springer, 229–248.
- Intel. 2000. Intel IA-64 Architecture Software Developer's Manual. (2000). <http://www.intel.com/content/dam/doc/manual/64-ia-32-architectures-optimization-manual.pdf>
- Ralf Karrenberg and Sebastian Hack. 2011. Whole-function Vectorization. In *Proceedings of the 9th Annual IEEE/ACM International Symposium on Code Generation and Optimization (CGO '11)*. IEEE Computer Society, Washington, DC, USA, 141–150. <http://dl.acm.org/citation.cfm?id=2190025.2190061>
- Robert Kennedy, Sun Chan, Shin-Ming Liu, Raymond Lo, Peng Tu, and Fred Chow. 1999a. Partial redundancy elimination in SSA form. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 21, 3 (1999), 627–676.
- Robert Kennedy, Sun Chan, Shin-Ming Liu, Raymond Lo, Peng Tu, and Fred Chow. 1999b. Partial redundancy elimination in SSA form. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 21, 3 (1999), 627–676.
- LLVM. 2016. LLVM bugs related to code hoisting. Bug IDs: 12754, 20242, 22005. (2016). <https://llvm.org/bugs>
- Etienne Morel and Claude Renvoise. 1979. Global optimization by suppression of partial redundancies. *Commun. ACM* 22, 2 (1979), 96–103.
- Steven S. Muchnick. 1997. *Advanced compiler design implementation*. Morgan Kaufmann.
- Nicholas Nethercote and Julian Seward. 2007. Valgrind: a framework for heavyweight dynamic binary instrumentation. In *ACM Sigplan notices*, Vol. 42. ACM, 89–100.
- Diego Novillo. 2007. Memory SSA – a unified approach for sparsely representing memory operations. In *Proc of the GCC Developers' Summit*.
- Jong-Soo Park and Jae-Jin Lee. 2008. A Practical Improvement to the Partial Redundancy Elimination in SSA Form. *Journal of Computing Science and Engineering* 2, 3 (2008), 301–320.
- Barry K Rosen, Mark N Wegman, and F Kenneth Zadeck. 1988. Global value numbers and redundant computations. In *Proceedings of the 15th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 12–27.

1101	Ghassan Shobaki, Maxim Shawabkeh, and Najm Eldeen Abu	UMI Order No. GAXNN-08158.	1156
1102	Rmaileh. 2013. Preallocation instruction scheduling with regis-	Xi Wang, Nikolai Zeldovich, M Frans Kaashoek, and Armando	1157
1103	ter pressure minimization using a combinatorial optimization	Solar-Lezama. 2013. Towards optimization-safe systems: An-	1158
1104	approach. <i>ACM Transactions on Architecture and Code Opti-</i>	alyzing the impact of undefined behavior. In <i>Proceedings of</i>	1159
1105	<i>mization (TACO)</i> 10, 3 (2013), 14.	<i>the Twenty-Fourth ACM Symposium on Operating Systems</i>	1160
1106	Vugranam C. Sreedhar. 1995. <i>Efficient Program Analysis Using DJ</i>	<i>Principles</i> . ACM, 260–275.	1161
1107	<i>Graphs</i> . Ph.D. Dissertation. Montreal, Que., Canada, Canada.		1162
1108			1163
1109			1164
1110			1165
1111			1166
1112			1167
1113			1168
1114			1169
1115			1170
1116			1171
1117			1172
1118			1173
1119			1174
1120			1175
1121			1176
1122			1177
1123			1178
1124			1179
1125			1180
1126			1181
1127			1182
1128			1183
1129			1184
1130			1185
1131			1186
1132			1187
1133			1188
1134			1189
1135			1190
1136			1191
1137			1192
1138			1193
1139			1194
1140			1195
1141			1196
1142			1197
1143			1198
1144			1199
1145			1200
1146			1201
1147			1202
1148			1203
1149			1204
1150			1205
1151			1206
1152			1207
1153			1208
1154			1209
1155			1210