

Global code motion of congruent instructions

Anonymous Author(s)

Abstract

We present a global-code-motion (GCM) compiler optimization which schedules congruent instructions across the program. Not only GCM saves code size, it exposes redundancies in some cases, it exposes more instruction level parallelism in the basic-block to which instructions are moved, and it enables other passes like loop invariant motion to remove more 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.

We have implemented the pass in LLVM. It is based on Global Value Numbering infrastructure available in LLVM. The experimental results show an average saving of 1% on the total compilation time on SPEC Cpu 2006. GCM enables more inlining and exposes more loop invariant code motion opportunities in majority of the benchmarks. We have also seen execution time improvements in a few of SPEC Cpu 2006 benchmarks viz. mcf and sjeng, moreover, register spills reduced by 2% on the SPEC Cpu 2006 benchmarks suite when compiled for x86_64-linux. GCM is an optimistic algorithm in the sense that it considers all congruent instructions in a function as potential candidates. We make an extra effort to hoist candidates by partitioning the potential candidates in a way to enable partial hoisting in case a common hoisting point for all the candidates cannot be found. 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.

Keywords keyword1, keyword2, keyword3

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. 1999] and brought these techniques in the main stream: nowadays SSA based PRE is available in industrial compilers like GCC and LLVM.

This paper describes 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 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 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 this:

```

111 float fun(float d, float min, float max, float a) {
112     float inv = 1.0f / d;
113     float x = (min - a) * inv;
114     float y = (max - a) * inv;
115     float tmin = x, tmax = y;
116     if (inv < 0) { tmin = y; tmax = x; }
117     return tmax + tmin;
118 }

```

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]. We have seen the performance of an out-of-order processor going up by 15% on a benchmark containing the above motivating example running in a tight loop. In fact, there are several advantages of GCM:

- it helps reduce the code size of the program by replacing multiple instructions with one although the final code size may be higher because of improved inlining heuristics; functions become cheaper to inline when their code size decreases because the inliner heuristics depend on instruction count;
- it exposes more instruction level parallelism 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 and hence reduce the critical path;
- 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;
- 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.

There has been a lot of work both in industry and academia to hoist and sink code out of branches, 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. To the best of our knowledge we have not found any reference which explored global-code-motion of congruent instructions in as much detail as us. A part of our implementation (aggressive code hoisting) is already merged in LLVM trunk, however, a general implementation of global-code-motion of congruent (non redundant) expressions is still missing from GCC and LLVM trunk. The main contributions of this paper are:

- a new optimistic algorithm to move congruent instructions out of branches;
- a cost model to reduce interference and hence, reduce spills;
- a technique to maximize hoisting in an optimistic approach by partitioning the list of potential candidates sorted by their DFS visit number;
- experimental evaluation of our implementation in LLVM which combines SSA based liveness analysis, and ranking of expressions to move very busy expressions in order to reduce register spills and improve performance in some cases.

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.

Dhamdhere [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 AN-TIC_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.

Click [1995] describes aggressive global-code-motion to first schedule 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%.

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

The algorithm for global-code-motion 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. 1989];
- 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;
- 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.

Dominance graphs, Global Value Numbers and Memory SSA are already available in LLVM but there is no facility to infer liveness of virtual registers. For that we implemented DJ-Graph data structure that allows us to calculate MergeSets which is used in the liveness analysis in our implementation of GCM.

3.1 Liveness using DJ-graph

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 and merge-sets to compute liveness of variables as illustrated in [Das et al. 2012]. It is very efficient for computing liveness and does not require any bitvectors to be maintained for each basic-block. The underlying simplicity of liveness computation is due to SSA form, where the values only flow (from def to use) either through dominator edges or the join edges (where we insert a PHI.) The DJ-graph contains both these edges which allows for computation of merge-sets for each basic-block, i.e., a set of all basic-blocks where the values can flow from a particular basic-block. We have implemented the merge-set computation based on a DJ-graph as illustrated in [Das and Ramakrishna 2005]. A simplified version of the merge-set that we implemented is presented here:

```
// Return true if the merge-set of source node of a
// visited J-edge (incoming edge of lnode) is not
// the subset of the merge-set of lnode.
bool
still_inconsistent(Node lnode, JEdges JE,
                  Visited V, MergeSet MergeSet)
for e in (all incoming edges to lnode):
  if e is in JE and e in V:
```

```

331         if MergeSet(Source Node of e)
332             is not a subset of MergeSet(lnode):
333             return true
334     return false
335     Each basic block which has a join-edge (J-Edge) to its suc-
336     cessor will have a merge-set that is a subset of the merge-set
337     of its successor [Das and Ramakrishna 2005]. The function
338     still_inconsistent returns true if this invariant is not satisfied
339     which allows the function constructMergeSet to repeat.
340     // Compute a merge-set breadth first order.
341     bool constructMergeSet_1(BFSList B, JEdges JE,
342                             DomLevel DL)
343     repeat = false
344     V = empty list of visited edges
345     for n in B:
346         for e in (all incoming edges to n):
347             if e is in JE and e not in V:
348                 V(e) = true
349                 let tmp = Source Node of e
350                 let tnode = Target Node of e
351                 let lnode = NULL
352                 while (DL(tmp) >= DL(tnode)):
353                     MergeSet(tmp) = MergeSet(tmp)
354                     U MergeSet(tnode)
355                     U {tnode}
356                 lnode = tmp
357                 tmp = dom-parent(tmp)
358             repeat = still_inconsistent(lnode, JE, MergeSet)
359     return repeat
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385

```

// Construct merge-set of each node in
 // Control-Flow-Graph G of a function.

```

void constructMergeSet(CFG G)
    B = Breadth First Order of G
    JE = JEdges of G
    DL = Path lengths (from root) of each node in G
    do: // Call until a fixed point is reached.
        Repeat = constructMergeSet_1(B, JE, DL)
    while (Repeat)

```

We compute merge sets, of the control flow graph, which remain same throughout the GCM transformation. For GCM we only want to know if a use operand is a kill (to compute changes in register pressure). For that we only need to know whether the use is also required later in the execution path. A variable is live out of a basic-block B if it is used in the merge set of B. We compute the live-out relation on-demand when profitability of hoist/sink candidates is to be evaluated. A simplified version of *isLiveOutUsingMergeSet* is presented here:

```

// Return true if variable A is liveout from N
bool isLiveOutUsingMergeSet(Node N, Variable A)
    if (A is defined in N)
        if (A is used outside any basic-block other than N)
            return true
        return false

```

```

MergeSet(N) = null
// The merge-set of N is the union of
// merge-sets of its successors.
for W in successors(N)
    MergeSet(N) = MergeSet(N) U MergeSet(W)

```

// Iterate over all the uses of A and see
 // if any intersect with the merge-set of N.

```

for T in users(A)
    B = basic_block(T)
    while (B != null and B != def_bb(A))
        if B in MergeSet(N) return true
        B = dom-parent(B)
    return false

```

If a variable is not live-out of a basic-block, it still may be used later in the basic-block. To establish whether a use of a variable is a kill we iterate on all the subsequent instructions in a basic-block checking for uses, the *isKill* function uses *isLiveOutUsingMergeSet* to find out if an operand is killed in an instruction.

```

// Return true if op is the last use of the operand in I
bool isKill(Operand op, Instruction I)
    if (isLiveOutUsingMergeSet(op))
        return false
    for each I1 after I in basic block of I
        if (I1 uses op)
            return false // Not a kill

```


return true

The original algorithm presented in [Das et al. 2012] has mismatched types in terms of uses and nodes because each node can have many instructions and hence many uses. Also, while iterating on the dominator of each user of a variable we may reach to the beginning of a function, in that case the inner while loop needs to terminate. These two cases were missing from the algorithm and we came across them during implementation.

The GCM pass can be broadly divided into the following steps that we will describe in the rest of this section:

- find candidates (congruent instructions) suitable for global-code-motion (Section 3.2),
- 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 3.5) and
- update data structures to continue iterative global-code-motion (Section 3.5).

3.2 Finding candidates to move

The first step is to find a set of congruent instructions [Briggs et al. 1997]. This is performed by a linear scan of all instructions of the program and classifying them by their value numbers. We could compute available and anticipable sets as computed by GCC's code-hoisting but that would be a lot of data structures to maintain at each basic-block level. For GCC it makes sense because their code-hoisting is integrated with GVN-PRE which already has those data structures available.

The current implementation of GVN in LLVM has some limitations when it comes to loads and stores so we compute the GVN of loads and stores separately. Our solution is to value number the address from where the value of a load is to be read from memory. For stores, we value number the address as well as the value to be stored at that address. Another limitation of the current GVN implementation in LLVM is that the instructions dependent on the loads will not get numbered correctly, and so after hoisting all candidates we need to rerun the GVN analysis in order to discover new candidates now available after having hoisted load instructions. This limitation should be addressed in a new implementation of the GVN based on MemorySSA, that would better account for equivalent loads and their dependent instructions.

The process of computing GVN can be on-demand, as we come across an instruction, or precomputed, computing GVN of all the instructions beforehand. Which process to choose is determined by the scope of code-hoisting we want to perform. In a pessimistic approach, we want to hoist a limited set of instructions from the sibling branches as we iterate the DFS tree bottom-up, it is sufficient to compute GVN values on-demand. Whereas, in the optimistic approach, as

described in Section 3.5, we want to move as many instructions as possible, and it would require GVN values to be precomputed.

Once the instructions have been classified into congruence classes, we compute for each group of congruent instructions, a point in the program that is both legal and profitable for the instructions to be moved to.

3.3 Legality check

Since the equality of candidates is purely based on the value numbers, we also need to establish if hoisting them to a common dominator or sinking them to a common post-dominator would be legal. Once a common dominator (post-dominator) is found, we check whether all the use-operands of the set of instructions are available at that position. It is possible to reconstitute or rematerialize the use-operands in some cases when the operands are not available and make it legal to move the instruction.

Subsequently, it is checked that the side-effects of the computations (if any) do not intersect with any side-effects between the instructions to be moved and their hoisting/sinking point. It is also necessary to check if there are indirect branch targets, e.g., landing pad, case statements, goto labels, etc., along the path because it becomes difficult to prove safety checks in those cases. In our current implementation we discard candidates on those paths.

3.3.1 Legality of hoisting scalars

Scalars are the easiest to hoist because we do not have to analyze them for aliasing memory references. As long as all the operands are available, or can be made available by rematerialization,) the scalar computations can be hoisted.

3.3.2 Legality of hoisting loads

The availability of operand to the load (an address) is checked at the hoisting point. If that is not available we try to rematerialize the address if possible. Along the path, from current position of the load instruction backwards on the control flow to the hoisting point, we check whether there are writes to memory that may alias with the load, in which case the candidate is discarded. To iterate on the use-def chains for memory references the MemorySSA infrastructure of LLVM is used.

3.3.3 Legality of hoisting stores

For stores, we check the dependency requirements similar to the hoisting of loads using the MemorySSA of LLVM. We check that the operands of the store instruction are available at the hoisting point, and that there are no aliasing loads or store along the path from the current position to the hoisting point.

3.3.4 Legality of hoisting calls

Call instructions can be divided into three categories: those calls equivalent to purely scalar computations, calls reading from memory, and most of the time, without further information, calls have to be classified as writing to memory, that is the most restrictive form. Each category of call instructions is handled as described for scalar, load, and store instructions.

Hoisting loads/stores across calls also require precise analysis of all the memory addresses accessed by the call. The current implementation being an intraprocedural pass, cannot hoist 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 hoisted across that call.

3.3.5 Legality of sinking expressions

Sometimes, hoisting is not upward-safe [Click 1995], e.g., if the expressions are in a landing pad, sinking of those expressions may reduce code size. Sinking may also reduce the register pressure in some cases, e.g., when the use operands are not kills. 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
B1: i1 = load A
    c1 = i1 + 10
    d1 = i0 + 20
    goto B3
B2: i2 = load A
    c2 = i2 + 10
    d2 = i0 + 20
    goto B3
B3: phi(c1, c2)
    phi(d1, d2)
```

In this example (c1, c2) or (d1, d2) are potential sinkable candidates. Since (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 of '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, just because the expressions are identical and operands are available, it still requires a unique post-dominating PHI to use the exact same values to be legally sinkable.

A general global scheduling algorithm also requires checks for undefined values when introducing a new computation along a path. Since only very busy expressions are moved in the current implementation, there is no need to check for undefinedness resulting due to movement of instructions. This simplifies the implementation.

3.3.6 Barriers

Barriers are based on the concept of pinned instructions [Click 1995] but extended to adapt to LLVM IR. Since a basic-block in LLVM IR is actually an extended basic-block because there might be calls in the middle of a basic-block which might not return. Essentially, any instruction that cannot guarantee progress is marked as a barrier. In the absence of context, as in our current implementation, some instructions which might be safe are still classified as barriers, e.g., volatile loads/stores, calls with missing attributes.

```
// Compute barriers for both hoistable and sinkable
// candidates.
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 there was no barrier in B, any
// instruction satisfying other legality checks
// 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
```

3.3.7 Downward-safety of movable instructions

In order to be able to move an instruction that has side effects, we must guarantee that the side-effects of that instruction before and after its code-motion will produce the same side-effects. For hoisting in particular, we need to prove that on all execution paths from the hoisting point to the end of the program 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 behaviors. The algorithm *compute_downward_safety* computes if hoisting a set of instruction to a common hoisting point is downward-safe [Muchnick 1997]. Since GCM is only dealing with VBEs, computing downward-safety ensures that no undefined behaviors are introduced as a result of code-motion.

```
// Return true if a hoistable instruction is downward
// safe at hoist-point
bool compute_downward_safety(Node HP)
    Worklist = basic-blocks of hoistable instructions
    DT = The dominator tree
    HP = the hoist-point
    for each basic-block B in the DT starting at HP:
```

```

661     if Worklist is empty
662         return false // Path exists! not downward safe
663     if B is in Worklist:
664         remove B from worklist
665         // Available at B => downward safe from B
666         remove subtree with root at B
667     if B is a leaf basic-block
668         // Path exists! not downward safe
669         return false
670     if B dominates hoist-point
671         return false // Back edge
672     return true

```

3.4 Profitability check (Cost models)

After the legality checks have passed, we check if a global-code-motion is profitable. That takes into account the impact global-code-motion 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 across popular benchmarks and at the same time to reduce code size as much as possible. Following are the cost models which are implemented:

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 global-code-motion is less than before it will be moved. Essentially, as long as there is one killed operand, code hoisting will either decrease or preserve the register pressure. Similarly, code-sinking will either decrease or preserve the register pressure as long as there is one operand killed at most. The following example explains how global-code-motion of identical computations can reduce the register pressure.

```

698 B0: b = m
699     c = n
700     if c is true then goto B3 else goto B4

```

```

702 B1: a0 = b<kill> + c<kill>

```

```

704 B2: a1 = b<kill> + c<kill>

```

After hoisting a0 and a1 are removed and a copy of a0 as a01 is placed in B0 just before the branch. Since 'b' and 'c' are killed in 'a0' and 'a1', hoisting the expressions will reduce the register pressure because two registers will be freed.

Ideally, it should be okay to hoist all the instructions and a later live-range-splitting [Cooper and Simpson 1998] pass should make the right decision of rematerializing the instruction should it be beneficial to do so. But the current live-range splitting pass of LLVM is not making the optimal

decision and we have found performance regressions while hoisting aggressively.

Moreover, LLVM has a 'getelementptr' instruction which computes the address where a load or a store would happen. It is a scalar computation and gets hoisted frequently even if the corresponding loads/stores would not be legal to hoist. In order to reduce the register pressure while hoisting loads, we have restricted hoisting of address computations away from their corresponding loads and stores when the loads and stores cannot be moved. This restriction is only to mitigate the limitations of LLVM's register allocator and may be lifted in the future, when the register allocation rematerialization pass has been improved to catch these regressions.

There is a special case for hoisting load instructions when the hoist-point is the predecessor basic-block for all the congruent loads. Even if the register pressure would increase, we prefer to hoist loads in this case. The reason being, that would make the loaded value available by the time it is used. Also, because stores and calls are hoisted the least, see Table 2, the performance does not change much whether they are hoisted or not.

3.4.2 Hoisting an expression across a call

Even hoisting scalars across calls is tricky because it can increase the number of spills. During the frame lowering of calls, 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. So before the call, 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 increases register pressure is not profitable to hoist.

3.4.3 Partitioning the list of hoist candidates to maximize hoisting

In the approach, described in Section 3.5, it is possible that a common hoisting point for all the instructions does not satisfy legality or profitability checks. In these cases, it is still possible to 'partially' hoist a subset of instructions by splitting the set of candidates and finding a closer hoisting point for each subset.

In order to hoist a subset of identical instructions, we partition the list of all candidates in a way to maximize the total number of hoist instructions. By sorting the list of all the candidates in the increasing order of their depth first search discovery time stamp [Cormen et al. 2001] (DFSIn numbers,) we make sure that candidates closer in the list have their common dominator nearby. In Figure-1, if B3, B4, B5, and B6 have identical computations and if for some reason, they cannot be hoisted at B0, then we partition the set of hoistable candidates in their DFSIn order. In this case the DFSIn ordering would be B3, B4, B5, B6 which will allow the instructions in (B3, B4) to be hoisted at B1 and those in (B5, B6) to be hoisted at B2. Even if instructions in (B3, B4)

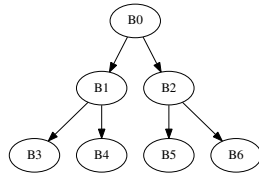


Figure 1. CFG to illustrate partitioning

could not be hoisted, for some reason, to B1, the instructions in (B5, B6) can still be hoisted – if legal and profitable – to B2.

In our current implementation we keep as many candidates in one set as possible (greedy approach.) To start with, first two candidates are tried for hoistability (see Section 3.5,) if that passes then the next instruction is tried (if available) and so on until a non-hoistable instruction is found. The partition starts from that point. We split the list at a point where the legality checks fail to hoist subset of candidates which are legal to hoist and then start finding new hoisting point for the remaining ones.

3.5 Code generation for the optimistic global-code-motion algorithm

Once all the legality and profitability checks are satisfied for a set of congruent instructions, they are suitable candidates for hoisting or sinking. 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 the SSA form is restored by updating the intermediate representation (IR) to reflect the changes. At the same time MemorySSA is also updated to get up-to-date information about memory references.

After one iteration of algorithm runs through the entire function, it creates more opportunities for *higher ranked* computations [Rosen et al. 1988]. Currently, this is a limitation of the GVN analysis pass, and so we rerun the code-hoisting algorithm until there are no more instructions left to be hoisted. Obviously, this is not the most optimal approach and can be improved by ranking the computations [Rosen et al. 1988], or by improving the GVN analysis to correctly populate congruence classes as the program is modified by the code generation.

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.

The amount of hoisting depends on whether we collect GVN of instructions before finding candidates (optimistic) or, on-demand (pessimistic). It also depends on the generality of the GVN algorithm as mentioned earlier in Section 3.2. We have implemented an optimistic global-code-motion of congruent instructions which uses the liveness analysis as

illustrated by Das et al. [2012] and ranking expressions explained by Rosen et al. [1988].

GCM basically consists of two parts, i.e., hoisting and sinking. This implementation only moves congruent instructions. An immediate guarantee of this approach is gain in code size (the final executable may be of larger size because of more inlining.) The algorithm prefers hoisting to sinking. 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.

```
void doGCM()
```

```
  Analyses available:
```

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

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

```
  constructMergeSet(CGF)
```

```
  computeBarriers()
```

```
  do:
```

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

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

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

```
  doCodeSinking()
```

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*).

For hoisting, we find a common dominator dominating all a set of congruent instructions and perform legality checks, as described in Section 3.3. Often times it is not possible to hoist all the instructions to one common dominator, due to legality constraints, e.g., intersecting side-effects, or profitability constraints, e.g., increasing register pressure; in those cases, this algorithm would partition (Section 3.4.3) the list of identical instructions into subsets which can be partially hoisted to their respective common dominators.


```

881 // Find hoistable candidates and do hoisting
882 int doCodeHoisting()
883   For each value number VN with 2 or more instructions
884     sort the instructions according to DFSIn numbers
885     if first two I1, I2 are hoistable
886       && downward-safe at hoist-point
887       && profitable to be hoisted at hoist-point:
888         continue the check for subsequent
889         instructions with same VN
890     else:
891       proceed to check if I2 and I3 (if available)
892       can be hoisted
893
894   For each set of hoistable instructions:
895     move the first one (I1) to the hoist-point
896     update their users others to refer to I1
897     remove all other instructions (I2, ..., Ik)
898     if I1 has memory references:
899       update MemorySSA of I1 and others
900     update statistics
901   return number of hoisted instructions

```

After no more candidates are hoistable, sinking is performed. Sinking is only performed once because currently there are very few sinkable candidates (Table 2) as SimplifyCFG of LLVM (which runs before GCM) has a code-sinking which already sinks several instructions.

```

909 void doCodeSinking()
910   // Find sinkable instructions
911   For each value number VN with 2 or more instructions:
912     if there are two instructions with same immediate
913       successor(S) as post-dominator:
914       if both are only used in the same PHI-node of S
915       && dependencies of both are available at S:
916         mark instructions as sinkable
917
918   // Sink the instructions and update statistics
919   For each pair (I1, I2) of sinkable instructions:
920     Move I1 to the sink point
921     // the sink point is just after all
922     // the PHI-nodes in the post-dominator
923     update all the uses of PHI node (PN) with I1
924     Remove I2 and PN
925     if I1 is a memory reference:
926       update MemorySSA of I1
927       update MemorySSA of removed instructions
928       to point to the MemorySSA reference of I1
929     remove MemorySSA reference of all
930     others which were deleted
931   update statistics
932
933   return number of sinked instructions

```

Code hoisting opens new opportunities for other hoistable candidates which were of higher rank (depended on candidates which got hoisted). Ideally we could iterate on lower ranking expressions first and then proceed to higher ranking expressions in the same iteration but LLVM's GVN infrastructure does not compute equivalence classes in an effective way. We found it simpler to just recompute the value numbers and start finding hoistable candidates again.

3.6 Time complexity of algorithm

The complexity of code hoisting is linear in number of instructions that are candidates for global-code-motion, matching the complexity of PRE on SSA form. The analyses computed for this pass are Global Value Numbering, Computation of DJ-graph and merge-sets, Marking Barriers, all are linear in number of instructions in a function. Liveness analysis is expensive and only performed on-demand for hoistable candidates so it does not affect the overall compilation time by much. 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, we have not seen significant increase in the compilation times (see Section 3). 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.

The analysis is followed by a simple code generation that adds the identified instruction in the destination point and removes all instructions rendered redundant as a result of global-code-motion.

4 Experimental Evaluation

For evaluation of global-code-motion, we built SPEC Cpu 2006 with the patch. All the experiments were conducted on an x86_64-linux machine and at -Ofast optimization level. 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 4, and Table 2. We also measured compile time Table 3 to see the impact of global-code-motion.

Overall there was only a minor change in the performance numbers with a few improvements (429.mcf, 458.sjeng, etc.) and a few regressions (447.dealII, 453.povray, etc.) The global-code-motion pass was run twice in the pass pipeline (at -Ofast), first after EarlyCSE pass, and then after the inliner. Because EarlyCSE removes local redundancies, GCM would have to analyze less redundant instructions, also as inliner creates more redundancies, GCM would have more opportunities to move congruent instructions out of branches.

We see some nice gains in code size at -Os with global-code-motion (403.gcc, 447.dealII, and 483.xalancbmk) as shown in Table 1. The code size listed here is the text size delta (from

Benchmarks	% performance uplift at -Ofast (high better)	size reduction (Bytes) at -Ofast (high better)	size reduction (Bytes) at -Os (high better)
400.perlbench	1.00	-9120	544
401.bzip2	1.00	-1600	128
403.gcc	1.00	-58176	3384
429.mcf	1.03	160	16
433.milc	1.00	-656	224
444.namd	1.00	-3864	-1168
445.gobmk	0.99	-20480	-264
447.dealII	0.95	-27836	50934
450.soplex	0.99	2220	348
453.povray	0.97	164	-472
456.hmmer	0.99	-3464	-352
458.sjeng	1.02	440	1456
462.libquantum	1.00	-96	0
464.h264ref	1.00	-3424	864
470.lbm	1.00	0	0
471.omnetpp	1.00	2688	200
473.astar	1.00	-928	152
482.sphinx3	0.99	3008	-176
483.xalancbmk	1.01	46492	4012
Geomean	0.997		

Table 1. Execution time (ratio) and code size change (Bytes) with and without global-code-motion (GCM) on SPEC Cpu 2006: performance uplift at -Ofast in mcf and sjeng; code size improvement in dealII at -Os.

linux size command) of the final executable for each benchmark. On the other hand, some benchmarks like 473.astar increased in code size by almost 2%. This is because once the code-size of a function decreases (due to GCM,) it becomes cheaper to inline. As we can see in Table 4, four more functions got inlined in 473.astar. 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.

The code-size at -Ofast is very widespread with huge gains in 483.xalancbmk to serious regressions in 403.gcc. This is mostly because the compiler's focus is to improve performance at the cost of code-size at -Ofast, so passes like inlining, loop-unrolling, code-versioning are very aggressive at -Ofast.

Spills are listed in Table 2. 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 the register pressure (except for loads Section 3.4,) overall spills still decreases.

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].

Some useful static metrics are presented in Table 4, which shows how global-code-motion impacted important compiler

Number of spills	base(-Ofast)	GCM(-Ofast)	%change
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
464.h264ref	3379	3454	1.02
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 2. Number of spills before and after global-code-motion(GCM) on SPEC2006 at -Ofast. Lower is better in %age change. Spills reduced by 2%.

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
464.h264ref	151	149	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 3. Change in total number of instructions executed during the compilation with and without global-code-motion (GCM) on SPEC Cpu 2006 at -Ofast. The instructions were counted using valgrind.

optimizations like inlining and 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 400.perlbench, 445.gobmk.

The compile time metrics of GCM are listed in Table 2. The table lists the number of scalars, loads, stores and calls

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.hmmmer	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
464.h264ref	0.5	0.0	13.0	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 4. Change in the static compile time metrics w.r.t. global-code-motion (GCM) on SPEC Cpu 2006 at -Ofast. The % columns show percentage increase in metric w.r.t. baseline. GCM improves LICM and inlining in general.

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 a few sunk instructions even if LLVM has a code sinking and hoisting transforms in the SimplifyCFG pass which runs before GCM: this is because sinking and hoisting in SimplifyCFG have some rather severe limitations to make the analysis and code transform very fast in order to be able to run the pass several times. Another reason is that instructions for which dependencies are not directly available in the successor (Section 3.3.5) is not sunk for now. Sinking those instructions would require a look ahead on the ability to sink the dependencies. We plan to implement this in near future.

The code shown in Section 1 is a reduced example that appears in a hot loop of a proprietary benchmark. When the expressions are hoisted from the conditional clauses, the overall performance of that benchmark improves by 15% on an out-of-order processor due to increased instruction level

parallelism, and better scheduling of the instructions, accommodating for the long latency of the division operation.

5 Conclusions and Future Work

We have presented the global-code-motion of congruent computations in SSA form. We saw that it improves inlining and LICM opportunities, reduces register spills, it also improves performance in some cases. To preserve performance and not hoist/sink too much, we have implemented a register pressure aware cost model as described in Section 3.4. A part of this implementation is already merged in LLVM trunk as GVNHoist.cpp and we posted the GCM implementation for review

There are a few improvements that we would like to make. First would be to integrate this 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. In general it is a good idea to start with lower ranked expressions first such that maximum hoisting can happen in one iteration, however, current implementation does not rank the expressions and iteratively finds a fixed point when no more candidates are available. Even this implementation converges quickly and no compile time regression have been observed. This is not the most optimal approach in terms of compile time complexity and results in multiple data structures to be recomputed. This can be improved by ranking the computations [Rosen et al. 1988]. 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. The first candidate would be the inliner. We have seen difference in inlining decisions in Table 1, when global-code-motion was enabled. Since inliner has several cost-models tuned for the previous pass layout, that would need some improvement.

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
- 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.

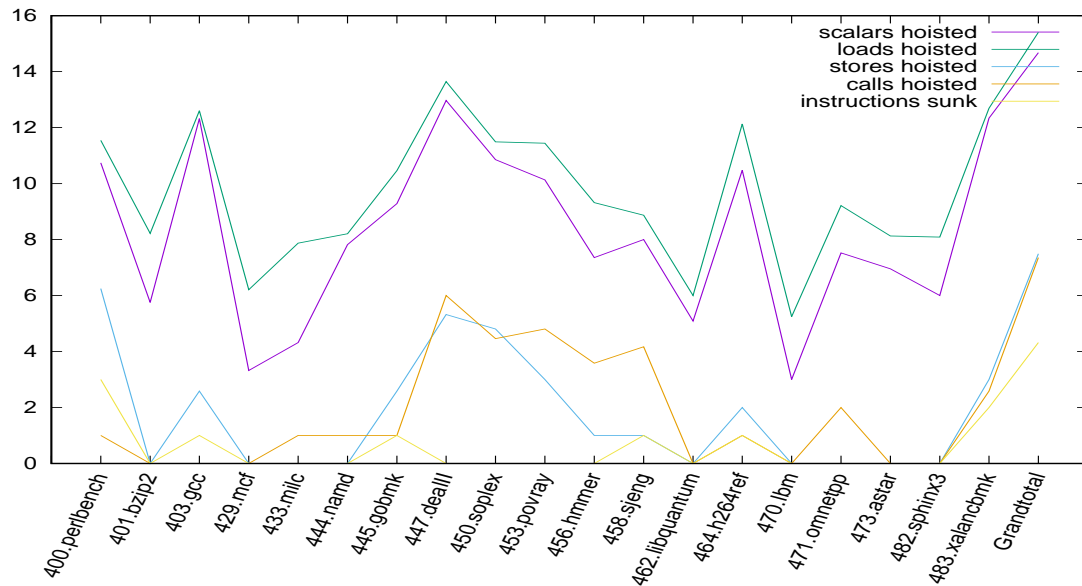


Figure 2. GCM stats on SPEC2006 at -Ofast. Loads are hoisted the most followed by scalars, stores and calls.

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.

Thomas H. Cormen, Charles Eric Leiserson, Ronald L Rivest, and Clifford Stein. 2001. *Introduction to algorithms*. Vol. 6. MIT press Cambridge.

Ron Cytron, Jeanne Ferrante, Barry K Rosen, Mark N Wegman, and F Kenneth Zadeck. 1989. 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.

Dibyendu Das and U Ramakrishna. 2005. A practical and fast iterative algorithm for ϕ -function computation using DJ graphs. *ACM Transactions on Programming Languages and Systems (TOPLAS)* 27, 3 (2005), 426–440.

Dhananjay M Dhamdhare. 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>

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>

Robert Kennedy, Sun Chan, Shin-Ming Liu, Raymond Lo, Peng Tu, and Fred Chow. 1999. 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*.

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.

Ghassan Shobaki, Maxim Shawabkeh, and Najm Eldeen Abu Rmaileh. 2013. Preallocation instruction scheduling with register pressure minimization using a combinatorial optimization approach. *ACM Transactions on Architecture and Code Optimization (TACO)* 10, 3 (2013), 14.

Vugranam C. Sreedhar. 1995. *Efficient Program Analysis Using Dj Graphs*. Ph.D. Dissertation. Montreal, Que., Canada, Canada. UMI Order No. GAXNN-08158.