

# 6.035 Phase 4 Documentation

*Decaf Optimizing Compiler*

An Bo Chen, Alex Dang, Youran (Yoland) Gao  
{anbochen, alexdang, youran}@mit.edu

## Design & Implementation

### Common Subexpression Elimination (CSE)

**Enable this optimization with flag:** `--opt cse`

Common subexpression elimination is performed as a CFG pass by first running an available expressions (AE) fix point algorithm similar to the one seen in lecture and then visiting all basic blocks in the program. Note that expressions that contain global variables do not participate in common subexpression elimination due to the possibility of side effects and the expressions in the bit vector are only bin op expressions. We use AE to find the set of the available expressions at the beginning of each basic block, which would be `IN[n]` for each block `n`. In the `GEN[n]` procedure, we make sure to not generate available expressions that we also kill in the same block.

At the beginning of the CSE analysis, we have `exprToTemp`, a mapping of `CfgExpression` to their corresponding unique temp variable, `varToExprs`, a mapping of `CfgVariable` to a set of `CfgExpressions` that have the variable as one of its operands, and `IN[n]` for each block `n`. For each assign statement in block `n`, if the right hand expression is available via `IN[n]` or a previously seen expression in the block and has not been killed yet, we replace right hand expression with the corresponding temp variable using `exprToTemp`. If not, then we append a new assign statement afterwards with the corresponding temp variable = left hand variable of the current assign statement. Regardless of what the right hand expression is, we kill all available expressions and seen expressions corresponding to the left hand variable using `varToExprs`.

Below is an example of a simple Decaf program with common subexpression enabled and the associated assembly instruction counts by type. Note that the number of assembly instructions increases slightly because we add some new assign statements when we save variables to temps. This enables Copy Propagation to propagate the use in the copy statements and Dead Code Elimination to remove extraneous assign statements. You can run `./run.sh --target cfg --opt cse tests/codegen/input/_cse-<TEST NUMBER>.dcf` to check out the new CFG applied with CSE. Note that variables without initializers are initialized with a default value.

Decaf	Unoptimized	Optimized
<pre>void main(){     int a, b, c, d;      c = a + b;     d = a + b; }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d, x)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; _t = a + b     &gt; c = _t     &gt; _t_2 = a + b     &gt; d = _t_2 }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; _t = a + b     &gt; _tcse = _t     &gt; c = _t     &gt; _t_2 = _tcse     &gt; d = _t_2 }</pre>
Count of Assembly Instructions:	<pre>movq: 13 jmp: 2 add: 1 enter: 1 leave: 1 ret: 1</pre>	<pre>movq: 15 jmp: 2 add: 1 enter: 1 leave: 1 ret: 1</pre>

Here is an example Decaf program with conditional branching:

Decaf	Unoptimized	Optimized
<pre>void main(){     int a, b, c, d, e, f;     bool x;     a = 1;     b = 2;     x = true;     c = a + b;     if ( x ) {         a = 1;         d = a + b;     }     else {         e = a + b;     } }</pre>	<pre>main() =&gt; void {     _main_1 @entry     (a, b, c, d, e, f, x)     &gt; a = 0 (default)     &gt; b = 0 (default)     &gt; c = 0     &gt; d = 0     &gt; e = 0     &gt; f = 0     &gt; x = 0     &gt; a = 1     &gt; b = 2     &gt; x = 1     &gt; _t = a + b</pre>	<pre>main() =&gt; void {     _main_1 @entry     (a, b, c, d, e, f, x)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; e = 0     &gt; f = 0     &gt; x = 0     &gt; a = 1     &gt; b = 2     &gt; x = 1     &gt; _t = a + b</pre>

	<pre> &gt; c = _.t x ? T-&gt; _.main_2 F-&gt; _.main_3  _.main_2 &gt; a = 1 &gt; _.t_2 = a + b  &gt; d = _.t_2 --&gt; _.main_4  _.main_3 &gt; _.t_3 = a + b &gt; e = _.t_3 --&gt; _.main_4  _.main_4 @exit } </pre>	<pre> &gt; _.tcse = _.t &gt; c = _.t x ? T-&gt; _.main_2 F-&gt; _.main_3  _.main_2 &gt; a = 1(killed a + b) &gt; _.t_2 = a + b &gt; _.tcse = _.t_2 &gt; d = _.t_2 --&gt; _.main_4  _.main_3 &gt; _.t_3 = _.tcse &gt; e = _.t_3 --&gt; _.main_4  _.main_4 @exit } </pre>
Count of Assembly Instructions:	<pre> movq: 25 jmp: 4 add: 3 cmpq: 1 enter: 1 je: 1 jne: 1 leave: 1 ret: 1 </pre>	<pre> movq: 29 jmp: 4 add: 2 cmpq: 1 enter: 1 je: 1 jne: 1 leave: 1 ret: 1 </pre>

## Copy Propagation (CP)

**Enable this optimization with flag:** `--opt cp`

Copy propagation is performed as a CFG pass by visiting all basic blocks in the program. We then visit all the statements in each basic block. Firstly, on all assign statements, the left hand variable is killed, regardless of what the right hand side is. When there is an assign statement such that both the left and right sides are of type `CfgVariable`, it is possible to add this mapping to `tmpToVars`. Global variables do not participate in copy propagation due to the possibility of side effects.

Then, for most uses of variables, we check the `tmpToVars` map to see if any copies can be propagated. Function parameters do not participate in copy propagation due to the possibility of side effects.

Below is an example of a simple Decaf program with copy propagation enabled and the associated assembly instruction counts by type. Note that the number of assembly instructions

do not change because we are only changing the variable that is being copied. This enables Dead Code Elimination to remove extraneous assign statements, but CP itself does not affect the generated code size directly. You can run `./run.sh --target cfg --opt cp tests/codegen/input/_cp-<TEST NUMBER>.dcf` to check out the new CFG applied with CP.

Decaf	Unoptimized	Optimized
<pre>void main(){     int a, b, c, d, x;      a = x;     b = a;     c = b;     d = b + c; }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d, x)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; x = 0     &gt; a = x     &gt; b = a     &gt; c = b     &gt; _t = b + c     &gt; d = _t }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d, x)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; x = 0     &gt; a = x     &gt; b = x     &gt; c = x     &gt; _t = x + x     &gt; d = _t }</pre>
<b>Count of Assembly Instructions:</b>	<pre>movq: 16 jmp: 2 add: 1 enter: 1 leave: 1 ret: 1</pre>	<pre>movq: 16 jmp: 2 add: 1 enter: 1 leave: 1 ret: 1</pre>

## Constant Propagation (CSP)

**Enable this optimization with flag:** `--opt csp`

Constant propagation is performed in the same file as copy propagation using almost exactly the same procedure, except with constants. If variable is assigned to a constant, it is added to the `tmpToVars` map, and when the left-hand variable is reassigned, it is removed from the map.

Then, for most `CfgExpression` types, we check the `tmpToVars` map to see if any constants can be propagated. Function parameters and global variables do not participate in constant propagation due to the possibility of side effects.

Below is an example of a simple Decaf program with constant propagation enabled and the associated assembly instruction counts by type. You can run `./run.sh --target cfg --opt csp tests/codegen/input/_cp-<TEST NUMBER>.dcf` to check out the new CFG applied with CSP.

Decaf	Unoptimized	Optimized
<pre>void main(){     int a, b, c, d;      a = 5;     b = a;     c = b;     d = c + b; }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; a = 5     &gt; b = a     &gt; c = b     &gt; _t = c + b     &gt; d = _t }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; a = 5     &gt; b = 5     &gt; c = 5     &gt; _t = 5 + 5     &gt; d = _t }</pre>
<b>Count of Assembly Instructions:</b>	<b>movq: 14</b> jmp: 2 add: 1 enter: 1 leave: 1 ret: 1	<b>movq: 12</b> jmp: 2 add: 1 enter: 1 leave: 1 ret: 1

## Algebraic Simplification

**Enable this optimization with flag:** `--opt asimple`

Algebraic simplification is implemented as a pass over the CFG. On each pass, it scans for `CfgExpression` that have `CfgValue` that are able to be simplified using simple algebraic rules. If the two sides are of type `CfgImmediate`, the value is computed and returned as another immediate. Otherwise, the following simplification rules are checked:

### *Algebraic Simplification rules:*

```
0 + x = x, a + 0 = x
0 - x = -x, a - 0 = x
x * 1 = x, 1 * x = x, x * 0 = 0, 0 * x = 0
x / 1 = x, 0 / x = 0, x / 0 = 0 (undef)
x % 1 = 0
```

### *Boolean Simplification rules:*

```
x && true = x, true && x = x
x && false = false, false && x = false
x || true = true, true || x = true
x || false = x, false || x = x
```

Below is an example of the algebraic simplification being done on a simple Decaf program and the associated assembly instruction counts by type.

Decaf	Unoptimized	Optimized
<pre> void main(){     int a, b;     bool c;      a = 1 + 2;     b = 2 * 5;     c = 10 == b; } </pre>	<pre> main() =&gt; void {     _main_1 @entry @exit     (a, b, c)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; <b>_t = 1 + 2</b>     &gt; a = <b>_t</b>     &gt; <b>_t_2 = 2 * 5</b>     &gt; b = <b>_t_2</b>     &gt; <b>_t_3 = 10 == b</b>     &gt; c = <b>_t_3</b> } </pre>	<pre> main() =&gt; void {     _main_1 @entry @exit     (a, b, c)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; <b>_t = 3</b>     &gt; a = <b>_t</b>     &gt; <b>_t_2 = 10</b>     &gt; b = <b>_t_2</b>     &gt; <b>_t_3 = 10 == b</b>     &gt; c = <b>_t_3</b> } </pre>
Count of Assembly Instructions:	<pre> movq: 18 jmp: 2 <del>add: 1</del> cmov: 1 cmpq: 1 enter: 1 <del>imul: 1</del> leave: 1 ret: 1 </pre>	<pre> movq: 12 jmp: 2  cmov: 1 cmpq: 1 enter: 1  leave: 1 ret: 1 </pre>

## Dead Code Elimination (DCE)

**Enable this optimization with flag:** `--opt dce`

Dead Code elimination is implemented on the function scope level. Within this scope, it first passes over all blocks and computes the variables that are assigned to, and (upwards) used within each block. Then, it traverses the control flow graph to compute liveness as in class (though it differs slightly from the slides in that it does account for global variables). Then, it passes over all blocks and eliminates the dead code using the liveness information. There is a special caveat that dead code where variables are assigned to function calls are not eliminated, but just converted to the function call without the assignment since the function may have side effects.

Finally, it repeats the above steps once again until the code stops changing, since it is possible that eliminating some dead code exposes more dead code. There is an optimization possible here in terms of performance- I had an idea where the liveness computation could be modified so as to focus only on the “essential” statements i.e. returns and function calls- but this was not done due to time constraints.

Perhaps the next step to further improve the DCE step would be to implement a heuristic which determines whether a Decaf function is pure or not- currently DCE assumes all function calls

may have side effects and thus does not eliminate some dead code containing pure function calls.

Below is an example of DCE working on a simple Decaf program:

Decaf	Unoptimized	Optimized
<pre>import printf; int get_int ( int x ) {     return x; } void main ( ) {     int a, b, c, d, e;     a = get_int ( 2 );     b = get_int ( 3 );     c = a + b;     d = c + b;     e = a * b;     printf ( "%d\n", e ); }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d, e)     &gt; a = 0     &gt; b = 0     &gt; c = 0     &gt; d = 0     &gt; e = 0     &gt; _t = get_int(2)     &gt; a = _t     &gt; _t_2 = get_int(3)     &gt; b = _t_2     &gt; _t_3 = a + b     &gt; c = _t_3     &gt; _t_4 = c + b     &gt; d = _t_4     &gt; _t_5 = a * b     &gt; e = _t_5     &gt; printf(_str_1, e) }</pre>	<pre>main() =&gt; void {     _main_1 @entry @exit     (a, b, c, d, e)     &gt; _t = get_int(2)     &gt; a = _t     &gt; _t_2 = get_int(3)     &gt; b = _t_2     &gt; _t_5 = a * b     &gt; e = _t_5     &gt; printf(_str_1, e) }</pre>
Count of Assembly Instructions:	<pre>movq: 30 jmp: 4 call: 3 <del>add: 2</del> enter: 2 leave: 2 ret: 2 imul: 1 leaq: 1</pre>	<pre>movq: 17 jmp: 4 call: 3 enter: 2 leave: 2 ret: 2 imul: 1 leaq: 1</pre>

## Difficulties with Phase 4

### Known Problems

- Global copy and constant propagation do not work properly at merge points. The temporary fix is to disable copy/constant propagation between blocks. We ran out of time to fix it as it involves implementing splitting.

## Phase 3 Fixes

### Phase 3 Private Tests Fixes

1. We failed to consider that some declarations are re-assigned on every iteration. For example, declarations inside of `while` and `for` loops.

*We fixed this by separating declarations from assignments, so that they are reassigned on each iteration in a `CfgAssignStatement`*

2. There was an assembly bug with division and modulo with the instructions generated.

*We needed to use the `cqto` assembly instruction because we needed to sign extend `%rax` to the octoword `%rdx:%rax`*

## Highlight for Feedback

- *Are there any good resources for helping us implement global copy propagation?*

## Contributions

In Phase 4, An Bo wrote most of the Common Subexpression Elimination (CSE) optimization and Yoland wrote most of Copy Propagation (CP) and Algebraic Simplification optimizations. An Bo and Yoland both worked together to debug both CP and CSE. Alex wrote most of the Dead Code Elimination code. We collaborated in writing the documentation.