

1 Introduction and Interpretation

1.1 Organization

[Slide 2] Module “Code Generation for Data Processing”

Learning Goals

- Getting from an intermediate code representation to machine code
- Designing and implementing IRs and machine code generators
- Apply for: JIT compilation, query compilation, ISA emulation

Prerequisites

- Computer Architecture, Assembly ERA, GRA/ASP
- Databases, Relational Algebra GDB
- Beneficial: Compiler Construction, Modern DBs

[Slide 3] Topic Overview

Introduction

- Introduction and Interpretation
- Compiler Front-end

Intermediate Representations

- IR Concepts and Design
- LLVM-IR
- Analyses and Optimizations

Compiler Back-end

- Instruction Selection
- Register Allocation
- Linker, Loader, Debuginfo

Applications

- JIT-compilation + Sandboxing
- Query Compilation
- Binary Translation

[Slide 4] Lecture Organization

- Lecturer: Dr. Alexis Engelke engelke@in.tum.de
- Time slot: Thu 10-14, 02.11.018
- Material: <https://db.in.tum.de/teaching/ws2425/codegen/>

Exam

- Written exam, 90 minutes, **no retake**, date TBD
- (Might change to oral on very low registration count)

[Slide 5] Exercises

- Regular homework, often with programming exercise
- Submission via POST request (see assignments)
 - Grading with $\{*, +, \sim, -\}$, feedback on best effort
- Exercise session modes:
 - Present and discuss homework solutions
 - Hands-on programming or analysis of systems (needs laptop)

Grade Bonus

- Requirement: $N - 2$ “sufficiently working” homework submissions **and** one presentations of homework in class (depends on submission count)
- Bonus: grades in $[1.3; 4.0]$ improved by 0.3/0.4

[Slide 6] Why study compilers?

- Critical component of every system, functionality and performance
 - Compiler mostly *alone* responsible for using hardware well
- Brings together many aspects of CS:
 - Theory, algorithms, systems, architecture, software engineering, (ML)
- New developments/requirements pose new challenges
 - New architectures, environments, language concepts, . . .
- High complexity!

[Slide 7] Compiler Lectures @ TUM

Compiler Construction IN2227, SS, THEO	Program Optimization IN2053, WS, THEO	Virtual Machines IN2040, SS, THEO
Front-end, parsing, semantic analyses, types	Analyses, transformations, abstract interpretation	Mapping programming paradigms to IR/bytecode
Programming Languages CIT3230000, WS	Code Generation CIT3230001, WS	
Implementation of advanced language features	Back-end, machine code generation, JIT comp.	

[Slide 8] Why study code generation?

- Frameworks (LLVM, ...) exist and are comparably good, but often not good enough (performance, features)
 - Many systems with code gen. have their own back-end
 - E.g.: V8, WebKit FTL, .NET RyuJIT, GHC, Zig, QEMU, Umbra, ...
- Machine code is not the only target: bytecode
 - Often used for code execution
 - E.g.: V8, Java, .NET MSIL, BEAM (Erlang), Python, MonetDB, eBPF, ...
 - Allows for flexible design
 - But: efficient execution needs machine code generation

[Slide 9] Proebsting's Law

“Compiler advances double computing power every 18 years.”

– Todd Proebsting, 1998^a

^a<http://proebsting.cs.arizona.edu/law.html>

- Still optimistic; depends on number of abstractions

The performance increases compilers can make on existing code are typically low. However, optimizing compilers gain more abilities in simplifying needlessly complex code, enabling the use of more abstractions and therefore higher level code. These abstractions are removed/optimized during compilation, enabling languages to promote these as *zero-cost abstractions*. They do, however, have a cost: compile times.

Also note that some of these “zero-cost” abstractions actually *do* have some runtime cost. For example, the mere possibility of C++ exceptions can cause less efficient machine code and might prevent optimizations due to the more complex control flow possibilities.

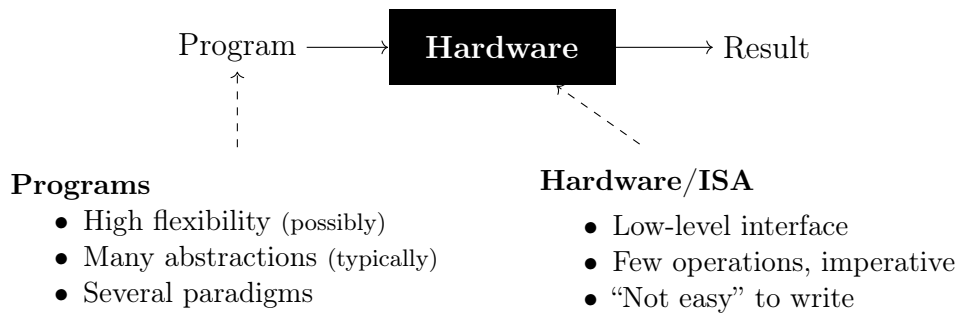
1.2 Overview**[Slide 10] Motivational Example: Brainfuck**

- Turing-complete esoteric programming language, 8 operations
 - Input/output: . ,
 - Moving pointer over infinite array: < >
 - Increment/decrement: + -
 - Jump to matching bracket if (not) zero: []

+++++[->+++++<]>.

- Execution with pen/paper? ☹

[Slide 11] Program Execution



[Slide 12] Motivational Example: Brainfuck – Interpretation

- Write an interpreter!

```

unsigned char state[10000];
unsigned ptr = 0, pc = 0;
while (prog[pc])
  switch (prog[pc++]) {
  case '.': putchar(state[ptr]); break;
  case ',': state[ptr] = getchar(); break;
  case '>': ptr++; break;
  case '<': ptr--; break;
  case '+': state[ptr]++; break;
  case '-': state[ptr]--; break;
  case '[': state[ptr] || (pc = matchParen(pc, prog)); break;
  case ']': state[ptr] && (pc = matchParen(pc, prog)); break;
  }
  
```

[Slide 13] Program Execution

Compiler



- Translate program to other lang.
- Might optimize/improve program
- C, C++, Rust → machine code
- Python, Java → bytecode

Multiple compilation steps can precede the “final interpretation”

Interpreter



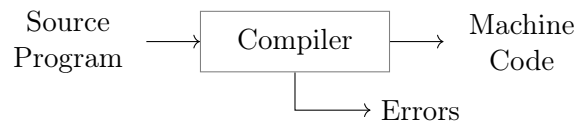
- Directly execute program
- Computes program result
- Shell scripts, Python bytecode, machine code (conceptually)

1.3 High-Level Structure of Compilers

[Slide 14] Compilers

- Targets: machine code, bytecode, or other source language
- Typical goals: better language usability and performance
 - Make lang. usable at all, faster, use less resources, etc.
- Constraints: specs, resources (comp.-time, etc.), requirements (perf., etc.)
- Examples:
 - “Classic” compilers source \rightarrow machine code
 - JIT compilation of JavaScript, WebAssembly, Java bytecode, ...
 - Database query compilation
 - ISA emulation/binary translation

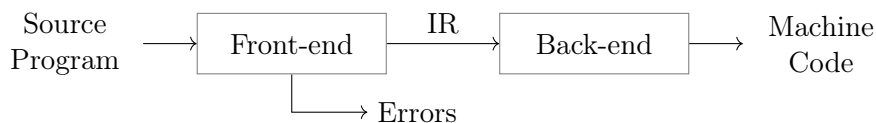
[Slide 15] Compiler Structure: Monolithic



- Inflexible architecture, hard to retarget

Some languages like C are designed to be compilable in a single pass without building any intermediate representation of the code between source and assembly. Single-pass compilers exist, but often have very limited possibilities to transform the code. They might not even know basic code properties, e.g., the size of the stack frame, during compilation of a function.

[Slide 16] Compiler Structure: Two-phase architecture



Front-end

- Parses source code
- Detect syntax/semantical errors
- Emit *intermediate representation* encode semantics/knowledge
- Typically: $\mathcal{O}(n)$ or $\mathcal{O}(n \log n)$

Back-end

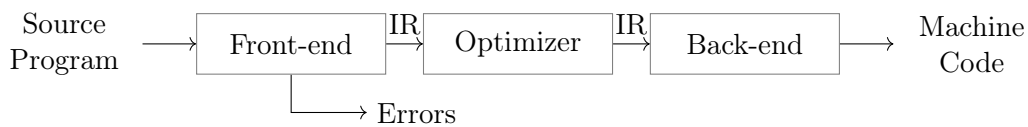
- Translate IR to target architecture
- Can assume valid IR (\rightsquigarrow no errors)
- Possibly one back-end per arch.

- Contains \mathcal{NP} -complete problems

After parsing, all information is encoded in the IR, including references to source code constructs for debugging support. The input source code is (at least conceptually) no longer needed.

In practice, there are very rare cases where the back-end can also raise errors. This can happen, for example, when some very architecture-specific constraints might be hard to verify during parsing (e.g., inline assembly constraints in combination with available registers).

[Slide 17] Compiler Structure: Three-phase architecture



- Optimizer: analyze/transform/rewrite program inside IR
-
- Conceptual architecture: real compilers typically much more complex
 - Several IRs in front-end and back-end, optimizations on different IRs
 - Multiple front-ends for different languages
 - Multiple back-ends for different architectures

Example Clang/LLVM (will be covered in more detail later): Clang parses the input into an abstract syntax tree (IR 1), uses this for semantic analyses; then Clang transforms the code into LLVM-IR (IR 2), which is primarily used for optimization; then the LLVM back-end transforms the code further into LLVM's Machine IR (IR 3), executes some low-level optimizations and register allocation there; the assembly printer of the back-end then lowers the code further to LLVM's machine code representation (IR 4), before finally emitting machine code. Some optimizations inside this pipeline, e.g. vectorization, might even build further representation of the code.

Why are compilers using so many different code representations? Different transformations work best at different abstraction levels. Diagnosing unused variables, for example, requires information about the source code. Optimization of arithmetic computations is easier in a data-flow-focused representation, where no explicit variables exist. Low-level modifications, like folding operations into complex addressing modes of the ISA, need a code representation where ISA instructions are already present.

[Slide 18] Compiler Front-end

1. Tokenizer: recognize words, numbers, operators, etc.

Re

- Example: $a+b*c \rightarrow \text{ID}(a) \text{ PLUS ID}(b) \text{ TIMES ID}(c)$

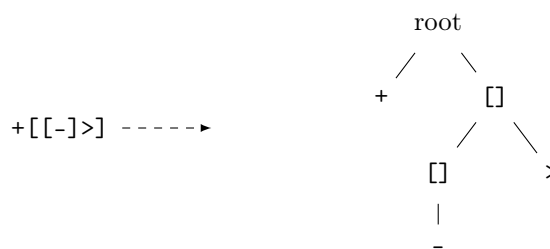
2. Parser: build (abstract) syntax tree, check for syntax errors *CFG*
 - Syntax Tree: describe grammatical structure of complete program Example: `expr("a", op("+"), expr("b", op("*"), expr("c")))`
 - Abstract Syntax Tree: only relevant information, more concise Example: `plus("a", times("b", "c"))`
3. Semantic Analysis: check types, variable existence, etc.
4. IR Generator: produce IR for next stage
 - This might be the AST itself

[Slide 19] Compiler Back-end

1. Instruction Selection: map IR operations to target instructions
 - Use target features: special insts., addressing modes, ...
 - Still using virtual/unlimited registers
2. Instruction Scheduling: optimize order for target arch.
 - Start memory/high-latency earlier, etc.
 - Requires knowledge about micro-architecture
3. Register Allocation: map values to fixed register set/stack
 - Use available registers effectively, minimize stack usage

1.4 Interpretation**[Slide 20] Motivational Example: Brainfuck – Front-end**

- Need to skip comments
- Bracket searching is expensive/redundant
- Idea: “parse” program!
- Tokenizer: yield next operation, skipping comments
- Parser: find matching brackets, construct AST



[Slide 21] Motivational Example: Brainfuck – AST Interpretation

- AST can be interpreted recursively

```
struct node { char kind; unsigned cldCnt; struct node* cld; };
struct state { unsigned char* arr; size_t ptr; };
void donode(struct node* n, struct state* s) {
    switch (n->kind) {
        case '+': s->arr[s->ptr]++; break;
        // ...
        case '[': while (s->arr[s->ptr]) children(n, s); break;
        case 0: children(n, s); break; // root
    }
}
void children(struct node* n, struct state* s) {
    for (unsigned i = 0; i < n->cldCnt; i++) donode(n->cld + i, s);
}
```

[Slide 22] Motivational Example: Brainfuck – Optimization

- Inefficient sequences of +/ - / < / > can be combined
 - Trivially done when generating IR
- Fold patterns into more high-level operations

In-Class Exercise:

Look at some Brainfuck programs. Which patterns are beneficial to fold?

[Slide 23] Motivational Example: Brainfuck – Optimization

- Fold offset into operation
 - `right(2) add(1) = addoff(2, 1) right(2)`
 - Also possible with loops
- Analysis: does loop move pointer?
 - Loops that keep position intact allow more optimizations
 - Maybe distinguish “regular loops” from arbitrary loops?
- Get rid of all “effect-less” pointer movements
- Combine arithmetic operations, disambiguate addresses, etc.

[Slide 24] Motivational Example: Brainfuck – Bytecode

- Tree is nice, but rather inefficient \rightsquigarrow flat and compact bytecode
- Avoid pointer dereferences/indirections; keep code size small
- Maybe dispatch two instructions at once?
 - `switch (ops[pc] | ops[pc+1] << 8)`
- Superinstructions: combine common sequences to one instruction

Dispatching multiple instructions at once can be problematic due to the explosion of cases that need to be implemented (often results in large jump tables and lots of code with resulting cache misses and branch mispredictions). Often, it is advisable to not always switch over multiple neighbored instructions, but instead combine common sequences into superinstructions.

[Slide 25] Motivational Example: Brainfuck – Threaded Interpretation

- Simple switch-case dispatch has lots of branch misses
- Threaded interpretation: at end of a handler, jump to next op

```
struct op { char op; char data; };
struct state { unsigned char* arr; size_t ptr; };
void threadedInterp(struct op* ops, struct state* s) {
    static const void* table[] = { &&CASE_ADD, &&CASE_RIGHT, };
#define DISPATCH do { goto *table[(++pc)->op]; } while (0)

    struct op* pc = ops;
    DISPATCH;

CASE_ADD: s->arr[s->ptr] += pc->data; DISPATCH;
CASE_RIGHT: s->arr += pc->data; DISPATCH;
}
```

With threaded interpretation there is not a single indirect jump instruction inside the dispatcher, but one indirect jump instruction per operation. Each of these indirect jumps then occupies a different branch prediction slot in the CPU. If an operation of type X is typically followed by an operation of type Y, with threaded interpretation the CPU has a much better chance of correctly predicting the dispatch branch to the next operation, because the indirect jump at the end of operation X typically jumps to operation Y. Without threaded interpretation, there would be only a single indirect branch, which is much harder to predict.

Threaded interpretation is especially useful on older and less powerful CPUs. Recent CPUs (e.g., Intel since Skylake, AMD since Zen 3, Apple Silicon) store the history of indirect branches and use this for better prediction. On such processors, threaded interpretation might not improve performance (or gains might be lower).

[Slide 26] Fast Interpretation

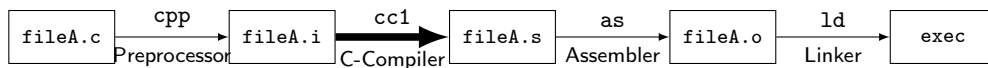
- Key technique to “avoid” compilation to machine code
- Preprocess program into efficiently executable bytecode
 - Easily identifiable opcode, homogeneous structure
 - Can be linear (fast to execute), but trees also work
 - Match bytecode ops with needed operations \rightsquigarrow fewer instructions
- Perhaps optimize – if it’s worth the benefit
 - Fold constants, combine instructions, ...
 - Consider superinstructions for common sequences

- For very cold code: avoid transformations at all

1.5 Context of Compilation

[Slide 27] Compiler: Surrounding – Compile-time

- Typical environment for a C/C++ compiler:



- Calling Convention: interface with other objects/libraries
- Build systems, dependencies, debuggers, etc.
- Compilation target machine (hardware, VM, etc.)

[Slide 28] Compiler: Surrounding – Run-time

- OS interface (I/O, ...)
- Memory management (allocation, GC, ...)
- Parallelization, threads, ...
- VM for execution of virtual assembly (JVM, ...)
- Run-time type checking
- Error handling: exception unwinding, assertions, ...
- Reflection, RTTI

[Slide 29] Motivational Example: Brainfuck – Runtime Environment

- Needs I/O for . and ,
- Error handling: unmatched brackets
- Memory management: infinitely sized array

In-Class Exercise:

How to efficiently emulate an infinitely sized array?

[Slide 30] Compilation point: AoT vs. JIT

Ahead-of-Time (AoT)

- All code has to be compiled
- No dynamic optimizations
- Compilation-time secondary concern

Just-in-Time (JIT)

- Compilation-time is critical
- Code can be compiled on-demand
 - Incremental optimization, too
- Handle cold code fast

- Dynamic specializations possible
- Allows for `eval()`

Various hybrid combinations possible

[Slide 31] Introduction and Interpretation – Summary

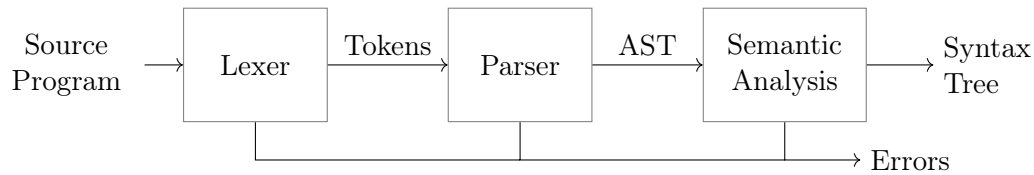
- Compilation vs. interpretation and combinations
- Compilers are key to usable/performant languages
- Target language typically machine code or bytecode
- Three-phase architecture widely used
- Interpretation techniques: bytecode, threaded interpretation, ...
- JIT compilation imposes different constraints

[Slide 32] Introduction and Interpretation – Questions

- What is typically compiled and what is interpreted? Why?
 - PostScript, C, JavaScript, HTML, SQL
- What are typical types of output languages of compilers?
- How does a compiler IR differ from the source input?
- What is the impact of the language paradigm on optimizations?
- What are important factors for an efficient interpreter?
- What are key differences between AoT and JIT compilation?

2 Compiler Front-end

[Slide 34] Compiler Front-end



- Typical architecture: separate lexer, parser, and context analysis
 - Allows for more efficient lexical analysis
 - Smaller components, easier to understand, etc.
- Some languages: preprocessor and macro expansion

2.1 Lexing

[Slide 35] Lexer

- Convert stream of chars to stream of words (*tokens*)
- Detect/classify identifiers, numbers, operators, ...
- Strip whitespace, comments, etc.

`a+b*c` → ID(a) PLUS ID(b) TIMES ID(c)

- Typically representable as regular expressions

[Slide 36] Typical Token Kinds

- Punctuators `() [] { } ; = + += | ||`
- Identifiers `abc123 main`
- Keywords `void int __asm__`
- Numeric constants `123 0xab1 5.7e3 0x1.8p1 09.1f`
- Char constants `'a' u'œ'`
- String literals `"abc\x12\n"`
- Internal `EOF COMMENT UNKNOWN INDENT DEDENT`
 - Comments might be useful for annotations, e.g. `// fallthrough`

Indentation-based languages like Python need separate tokens for indent/dedent, the indentation level is tracked in the lexer. Parsing numbers may need special care to correctly handle all possible cases of integer and floating-point numbers.

[Slide 37] Lexer Implementation

```
struct Token { enum Kind { IDENT, EOF, PLUS, PLUSEQ, /*...*/ };
  std::string_view v; Kind kind; };
Token next(std::string_view v) {
  if (v.empty()) return Token{v, Token::EOF};
  if (v.starts_with("+=")) return Token{"+"sv, Token::PLUSEQ};
  if (v.starts_with("+")) return Token{"+"sv, Token::PLUS};
  switch (v[0]) {
  case ' ', '\n', '\t': return next(v.substr(1)); // skip whitespace
  case 'a' ... 'z', 'A' ... 'Z', '_': {
    Token t = // ... parse identifier, e.g. using regex
    if (auto kind = isKeyword(t.v)) return Token{*kind, t.v};
    return t;
  }
  case '0' ... '9': // ... parse number
  default: return Token{v.substr(0, 1), Token::ERROR};
  }
}
```

This is just a minimal and non-optimized implementation to illustrate the concept. Performance-focused implementations do not use explicit regular expressions but write the state machine into code.

The struct `Token` has room for improvement. First, a `string_view` is unnecessarily large with 16 bytes, most tokens are smaller than 2^{16} bytes. Some tracking of the source locations is advisable for attaching diagnostics to their origin inside the code, for example by storing a file ID and the byte offset into the file. By tracking the byte offsets of line breaks, the line number can be reconstructed in $\mathcal{O}(\log n)$ from the byte offset.

Another optimization strategy is string interning, where identifiers are converted into unique integers (or pointers) during parsing. During later phases, comparing interned strings is much more efficient, as it is just an integer/pointer comparison. Another benefit is that the entire input file does not need to be kept in memory during parsing.

[Slide 38] Lexing C??= main() <%

```
  // yay, this is C99??/
  puts("hi_world!");
  puts("what's_up??!");
%>
```

Output: what's up|

- Trigraphs for systems with more limited encodings/char sets
- Digraphs to provide a more readable alternative...

Besides digraphs, trigraphs, and the preprocessor, C has another weird property: identifier names can be split by `\`, which concatenates two lines. It is necessary to construct the “real” identifier first. To simplify memory management in such cases, a bump pointer allocator (allocate large chunks of memory from the OS, then simply bump the end pointer for every allocation) can be useful to store such constructed names.

[Slide 39] Lexer Implementation

- Essentially a DFA (for most languages)
 - Set of regexes \rightarrow NFA \rightarrow DFA
- Respect whitespace/separators for operators, e.g. `+` and `+=`
- Automatic tools (e.g., flex) exist; most compilers do their own
- Keywords typically parsed as identifiers first
 - Check identifier if it is a keyword; can use perfect hashing
- Other practical problems
 - UTF-8 homoglyphs; trigraphs; pre-processing directives

A tool to generate perfect hash tables from a set of keywords is `gperf`. Example, compile with `gperf -L C++ -C -E -t <input>`:

```
struct keyword {char* name; int val; }
%%
int, 1
char, 2
void, 3
if, 4
else, 5
while, 6
return, 7
```

2.2 Parsing

[Slide 40] Parsing

- Convert stream of tokens into (abstract) syntax tree
- Most programming languages are context-sensitive
 - Variable declarations, argument count, type match, etc. \rightsquigarrow separated into semantic analysis
- Syntactically valid: `void foo = doesntExist / "abc";`
- Grammar usually specified as CFG

[Slide 41] Context-Free Grammar (CFG)

- Terminals: basic symbols/tokens
- Non-terminals: syntactic variables

- Start symbol: non-terminal defining language
- Productions: non-terminal \rightarrow series of (non-)terminals

```
stmt  $\rightarrow$  whileStmt | breakStmt | exprStmt
whileStmt  $\rightarrow$  while ( expr ) stmt
breakStmt  $\rightarrow$  break ;
exprStmt  $\rightarrow$  expr ;
expr  $\rightarrow$  expr + expr | expr * expr | expr = expr | ( expr ) | number
```

[Slide 42] Hand-written Parsing – First Try

- One function per non-terminal
- Check expected structure
- Return AST node
- Need look-ahead!

```
NodePtr parseBreakStmt() {
    consume(Token::BREAK);
    consume(Token::SEMICOLON);
    return newNode(Node::BreakStmt);
}
NodePtr parseWhileStmt() {
    consume(Token::WHILE);
    consume(Token::LPAREN);
    NodePtr expr = parseExpr();
    consume(Token::RPAREN);
    NodePtr body = parseStmt();
    return newNode(Node::WhileStmt,
        {expr, body});
}
NodePtr parseStmt() {
    // whoops!
}
```

[Slide 43] Hand-written Parsing – Second Try

- Need look-ahead to distinguish production rules
- Consequences for grammar:
 - No left-recursion
 - First n terminals must allow distinguishing rules
 - $LL(n)$ grammar; n typically 1 \Rightarrow Not all CFGs (easily) parseable (but most programming langs. are)
- Now... expressions

```
NodePtr parseBreakStmt() { /*...*/ }
NodePtr parseWhileStmt() { /*...*/ }

NodePtr parseStmt() {
    Token t = peekToken();
    if (t.kind == Token::BREAK)
        return parseBreakStmt();
}
```

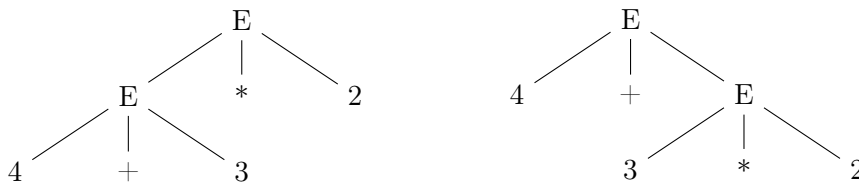


```

if (t.kind == Token::WHILE)
    return parseWhileStmt();
// ...
NodePtr expr = parseExpr();
consume(Token::SEMICOLON);
return newNode(Node::ExprStmt,
    {expr});
}

```

[Slide 44] Ambiguity

$$expr \rightarrow expr + expr \mid expr * expr \mid expr = expr \mid (expr) \mid \text{number}$$
Input: $4 + 3 * 2$ 

The grammar, as specified, is ambiguous, there are two possible ways to parse the input.

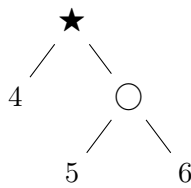
[Slide 45] Ambiguity – Rewrite Grammar?

$$primary \rightarrow (expr) \mid \text{number}$$

$$expr \rightarrow primary + expr \mid primary * expr \mid primary = expr \mid primary$$
Input: $4 + 3 * 2$ Input: $4 * 3 + 2$ 

The grammar is no longer ambiguous, but the result might not be expected, conventionally, multiplication has a stronger binding than addition.

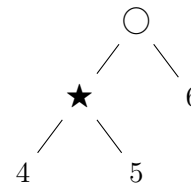
[Slide 46] Ambiguity – PrecedenceInput: $4 \star 5 \circ 6$



If $prec(○) > prec(★)$ or equal prec. and $★$ is right-assoc.

Examples:

- $4 + 5 \cdot 6$ ($prec(\cdot) > prec(+)$)
- $a = b = c$ ($=$ is right-assoc.)
 $b = c$ should be executed first



If $prec(○) < prec(★)$ or equal prec. and $★$ is left-assoc.

Examples:

- $4 + 5 < 6$ ($prec(<) < prec(+)$)
- $a + b - c$ ($+$ is left-assoc.)
 $a + b$ should be executed first

[Slide 47] Hand-written Parsing – Expression Parsing

- Start with basic expr.:
- Number, variable, etc.
- Parenthesized expr.
 - Parse full expression
 - Next token must be)
- Unary expr: followed by expr. with higher prec.
 - $- < \text{unary} - < [] / ->$

```

NodePtr parseExpr(unsigned minPrec=0);
NodePtr parsePrimaryExpr() {
  switch (Token t = next(); t.kind) {
  case Token::IDENT:
    return makeNode(Node::IDENT, t.v);
  case Token::NUMBER: // ...
  case Token::MINUS:
    // Only exprs with high precedence
    return makeNode(Node::UMINUS,
      {parseExpr(UNARY_PREC)});
  case Token::LPAREN: // ...
    // ...
  }
}

```

[Slide 48] Hand-written Parsing – Expression Parsing

- Only allow ops. with higher prec. on the right child
 - Right-assoc.: allow same
- Lower prec.: return + insert higher up in the tree

```

OpDesc OPS[] = { // {prec, rassoc}
  [Token::MUL] = {12, false},
  [Token::ADD] = {11, false},

```

```
[Token::EQ] = {2, true},
[Token::QUEST] = {3, true}, // ?:
}
NodePtr parseExpr(unsigned minPrec=1) {
  auto lhs = parsePrimaryExpr();
  while (auto op = OPS[next().kind];
         op.prec >= minPrec) {
    // ... handle (, [, ?: ...
    auto newPrec = op.rassoc ?
      op.prec : op.prec + 1;
    auto rhs = parseExpr(newPrec);
    lhs = makeNode(op.nodeKind,
      {lhs, rhs});
  }
  return lhs;
}
```

In-Class Exercise:

$a = 3 * 2 + 1;$ $a = b + c + d = 1;$ $a ? 1 : b ? 2 : 3;$

Example for input: $a = 3 * 2 + 1;$

	Rec. Depth 1	Rec. Depth 2	Rec. Depth 3
minPrec	1		
lhs	a		
op (prec/assoc)	= (2/r)		
minPrec	1	2	
lhs	a	3	
op (prec/assoc)	= (2/r)	* (12/l)	
minPrec	1	2	13
lhs	a	3	2
op (prec/assoc)	= (2/r)	* (12/l)	+ (11/l)
minPrec	1	2	
lhs	a	3*2	
op (prec/assoc)	= (2/r)	+ (11/l)	
minPrec	1	2	12
lhs	a	3*2	1
op (prec/assoc)	= (2/r)	+ (11/l)	; (0/-)
minPrec	1	2	
lhs	a	(3*2)+1	
op (prec/assoc)	= (2/r)	; (0/-)	
minPrec	1		
lhs	a=((3*2)+1)		
op (prec/assoc)	; (0/-)		

[Slide 49] Top-down vs. Bottom-up Parsing

Top-down Parsing

- Start with top rule
- Every step: choose expansion
- LL(1) parser
 - Left-to-right, Leftmost Derivation
- “Easily” writable by hand
- Error handling rather simple
- Covers many prog. languages

Bottom-up Parsing

- Start with text
- Reduce to non-terminal
- LR(1) parser
 - Left-to-right, Rightmost Derivation
 - Strict super-set of LL(1)
- Often: uses parser generator

- Error handling more complex
- Covers nearly all prog. languages

[Slide 50] Parser Generators

- Writing parsers by hand can be large effort
- Parser generators can simplify parser writing a lot
 - Yacc/Bison, PLY, ANTLR, ...
- Automatic generation of parser/parsing tables from CFG
 - Finds ambiguities in the grammar
 - Lexer often written by hand
- Used heavily in practice, unless error handling is important

[Slide 51] Bison Example – part 1

```
%define api.pure full
#define api.value.type {ASTNode*}
%param { Lexer* lexer }
%code{
static int yylex(ASTNode** lvalp, Lexer* lexer);
}
%token NUMBER
%token WHILE "while"
%token BREAK "break"

// precedence and associativity
%right '='
%left '+'
%left '*'
```

[Slide 52] Bison Example – part 2

```
%%
stmt : WHILE '(' expr ')' stmt { $$ = mkNode(WHILE, $1, $2); }
      | BREAK ';'              { $$ = mkNode(BREAK, NULL, NULL); }
      | expr ';'                { $$ = $1; }
      ;
expr  : expr '+' expr           { $$ = mkNode('+', $1, $2); }
      | expr '*' expr          { $$ = mkNode('*', $1, $2); }
      | expr '=' expr          { $$ = mkNode('=', $1, $2); }
      | '(' expr ')'           { $$ = $1; }
      | NUMBER
      ;
%%
static int yylex(ASTNode** lvalp, Lexer* lexer) {
    /* return next token, or YYEOF/... */ }
```

Compile with `bison -dg input.ypp`, it will emit a C++ header, the implementation file, and also a graph showing the state machine of the parser.

[Slide 53] Parsing in Practice

- Some use parser generators, e.g. Python some use hand-written parsers, e.g. GCC, Clang, Swift, Go
- Optimization of grammar for performance
 - Rewrite rules to reduce states, etc.
- Useful error-handling: complex!
 - Try skipping to next separator, e.g. `;` or `,`
- Programming languages are not always context-free
 - C: `foo* bar;`
 - May need to break separation between lexer and parser

In fact, many compilers^a use hand-written parsers, because they allow for better error messages a more graceful handling of syntax errors, leading to more reported errors during a single (failing) compilation.

^a<https://notes.eatonphil.com/parser-generators-vs-handwritten-parsers-survey-2021.html>

[Slide 54] Parsing C++

- C++ is not context-free (inherited from C): `T * a;`
- C++ is ambiguous: `Type (a), b;`
 - Can be a declaration or a comma expression
- C++ templates are Turing-complete¹
- C++ *parsing* is hence *undecidable*²
 - Template instantiation combined with `C T * a` ambiguity

2.3 Semantic Analysis

[Slide 55] Semantic Analysis

- Syntactical correctness $\not\Rightarrow$ correct program `void foo = doesntExist / ++"abc";`
- Needs context-sensitive analysis:
 - Variable existence, storage, accessibility, ...
 - Function existence, arguments, ...
 - Operator type compatibility
 - Attribute allowance
- Additional type complexity: inference, polymorphism, ...

¹TL Veldhuizen. *C++ templates are Turing complete*. 2003. URL: <http://port70.net/~nsz/c/c%2B%2B/turing.pdf>.

²J Haberman. *Parsing C++ is literally undecidable*. 2013. URL: <https://blog.reverberate.org/2013/08/parsing-c-is-literally-undecidable.html>.

[Slide 56] Semantic Analysis: Scope Checking with AST Walking

- Idea: walk through AST (in DFS-order) and validate on the way
- Keep track of scope with declared variables
 - Might need to keep track of defined types separately

In-Class Exercise:

How to implement the scope data structure?

- For identifiers: check existence and get type
- For expressions: check types and derive result type
- For assignment: check lvalue-ness of left side
- *Might* be possible during AST creation
- Needs care with built-ins and other special constructs

There are two ways of implementing a scoped hash table:

- Chain of hash maps: $Scope = (Map[Name \rightarrow Type] \text{ names}, Scope \text{ parent})$. This is, however, very slow for deeply nested scopes, as all hash maps of the parent scopes must be queried. Hash map lookups are fairly expensive.
- Hash map of lists: $Map[Name \rightarrow List[Tuple[Depth, Type]]]$. For every identifier, the type at a given scope nesting depth is stored. Invalidation can be implemented with an epoch counter for every depth. The downside is that this hash map can grow very large, as entries are never removed.

[Slide 57] Semantic Analysis and Post-Parsing Transformations

- Check for error-prone code patterns
 - Completeness of `switch`, out-of-range constants, unused variables, ...
- Check method calls, parameter types
- Duplicate code for templates
- Make implicit value conversions explicit
- Handle attributes: visibility, warnings, etc.
- Mangle names, split functions (OpenMP), ABI-specific setup, ...
- Last step: generate IR code

2.4 Miscellaneous**[Slide 58] Parsing Performance**

Is parsing/front-end performance important?

- Not necessarily: normal compilers
 - Some languages (e.g., Rust) need unbounded time *for parsing*

- Somewhat: JIT compilers
 - Start-up time is generally noticeable
- Somewhat more: Developer tools
 - Imagine: waiting for seconds just for updated syntax highlighting
 - Often uses tricks like incremental updates to parse tree

[Slide 59] Data Types

- Important part of programming languages
- Might have large variety and compatibility
 - Numbers, Strings, Arrays, Compound Types (struct/union), Enum, Templates, Functions, Pointers, ...
 - Class hierarchy, Interfaces, Abstract Classes, ...
 - Integer/float compatibility, promotion, ...
- Might have implicit conversions

[Slide 60] Data Types: Implementing Classes

- Simple `class/struct`: trivial, just bunch of fields
 - Methods take (pointer to) `this` as implicit parameter
- Single inheritance: also trivial – extend struct at end
- Virtual methods: store vtable in object representation
 - vtable = table of function pointers for virtual methods
 - Each sub-class has their own vtable
- Multiple inheritance is much more involved
- Dynamic casts: needs run-time type information (RTTI)

[Slide 61] Recommended Lectures

AD IN2227 “Compiler Constructions” covers parsing/analysis in depth

AD CIT3230000 “Programming Languages” covers dispatching/mixins/...

[Slide 62] Compiler Front-end – Summary

- Lexer splits input into tokens
 - Essentially Regex-Matching + Keywords; rather simple
- Parser constructs (abstract) syntax tree from tokens
 - Top-down vs. bottom-up parsing
 - Typical: top-down for control flow; bottom-up for expressions
 - Respect precedence and associativity for operators
- Semantic analysis ensures meaningful program

- Some data structures are complex to implement
- Some programming languages are more difficult to parse

[Slide 63] Compiler Front-end – Questions

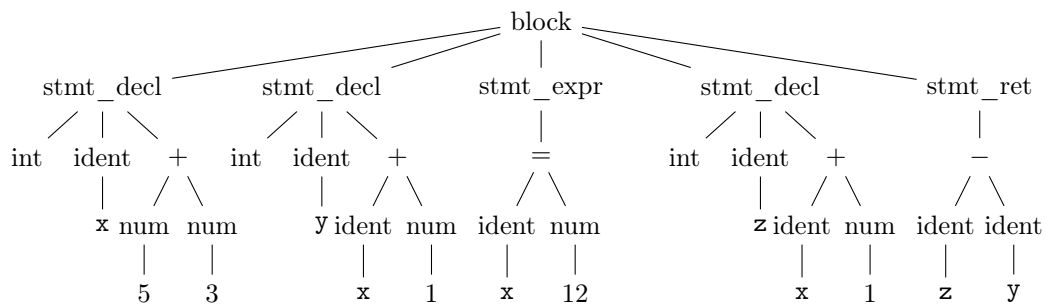
- What are typical components of a compiler front-end?
- What output does the lexer produce?
- How does a parser disambiguate rules?
- What is the typical way to handle operator precedence?
- Why are not all programming languages describable using CFGs?
- How to implement classes with virtual functions?

3 Intermediate Representations

[Slide 65] Intermediate Representations: Motivation

- So far: program parsed into AST
- + Great for language-related checks
- + Easy to correlate with original source code (e.g., errors)
- Hard for analyses/optimizations due to high complexity
 - variable names, control flow constructs, etc.
 - Data and control flow implicit
- Highly language-specific

[Slide 66] Intermediate Representations: Motivation



Question: how to optimize? Is $x+1$ redundant? \rightsquigarrow hard to tell ☹️

In this representation, it is very easy to see that the two $+1$ operations have different operands on the left side and are therefore not trivially redundant.

[Slide 67] Intermediate Representations: Motivation

```

x1 ← 5 + 3
y1 ← x1 + 1
x2 ← 12
z1 ← x2 + 1
tmp1 ← z1 - y1
return tmp1
  
```

Question: how to optimize? Is $x+1$ redundant? \rightsquigarrow No! 😊

[Slide 68] Intermediate Representations

- Definitive program representation inside compiler
 - During compilation, only the (current) IR is considered

In practice, there are, of course, exceptions to the general rule; sometimes an IR contains references to a previous/higher-level IR. An example is LLVM's low-level Machine IR, which only represents single functions and therefore references to global variables use the higher-level LLVM IR.

- Goal: simplify analyses/transformations
 - *Technically*, single-step compilation is possible for, e.g., C ... but optimizations are hard without proper IRs
- Compilers *design* IRs to support frequent operations
 - IR design can vary strongly between compilers
- Typically based on **graphs** or **linear instructions** (or both)

[Slide 69] Compiler Design: Effect of Languages – Imperative

- Step-by-step execution of program modification of state
- Close to hardware execution model
- Direct influence of result

- Tracking of state is complex
- Dynamic typing: more complexity
- Limits optimization possibilities

```
void addvec(int* a, const int* b) {
    for (unsigned i = 0; i < 4; i++)
        a[i] += b[i]; // vectorizable?
}
func:
    mov [rdi], rsi
    mov [rdi+8], rdx
    mov [rdi], 0 // redundant?
    ret
```

Tracking state, especially when memory is involved, is one of the main challenges during optimization. In the first example, the loop is not easily vectorizable, because `a` and `b` could point to the same underlying array (e.g., with `addvec(buf + 1, buf)`).

[Slide 70] Compiler Design: Effect of Languages – Declarative

- Describes execution target
- Compiler has to derive good mapping to imperative hardware

- Allows for more optimizations
- Mapping to hardware non-trivial

- Might need more stages
- Preserve semantic info for opt!

- Programmer has less “control”

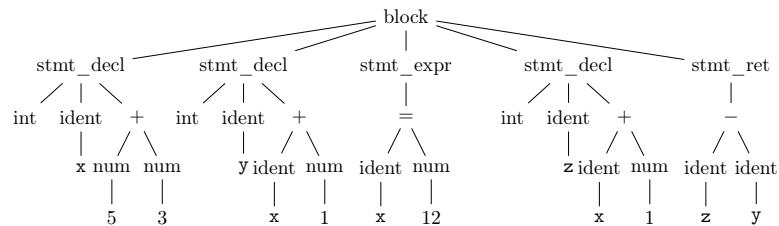
```

select s.name
from studenten s
where exists (select 1
              from hoeren h
              where h.matrno=s.matrno)
let rec fac = function
| 0 | 1 -> 1
| n -> n * fac (n - 1)

```

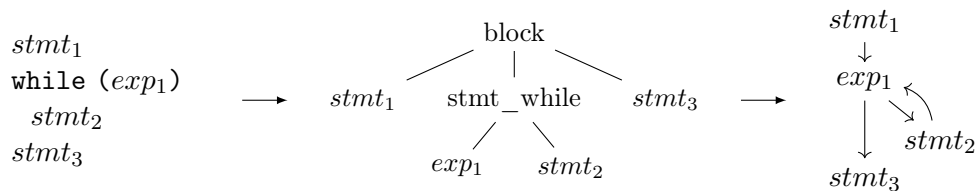
[Slide 71] Graph IRs: Abstract Syntax Tree (AST)

- Code representation close to the source
- Representation of types, constants, etc. might differ
- Storage might be problematic for large inputs



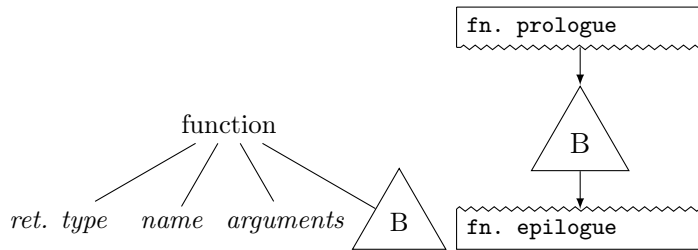
[Slide 72] Graph IRs: Control Flow Graph (CFG)

- Motivation: model control flow between different code sections
- Graph nodes represent **basic blocks**
 - Basic block: sequence of branch-free code (modulo exceptions)
 - Typically represented using a linear IR

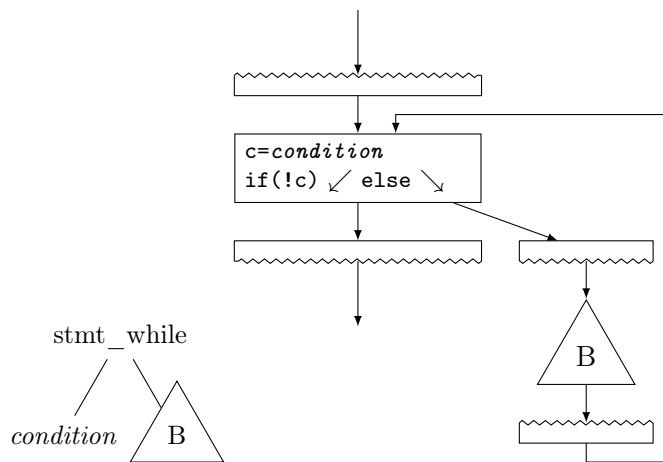


[Slide 73] Build CFG from AST – Function

- Idea: Keep track of current insert block while walking through AST



[Slide 74] Build CFG from AST – While Loop



Written in pseudo-code:

```

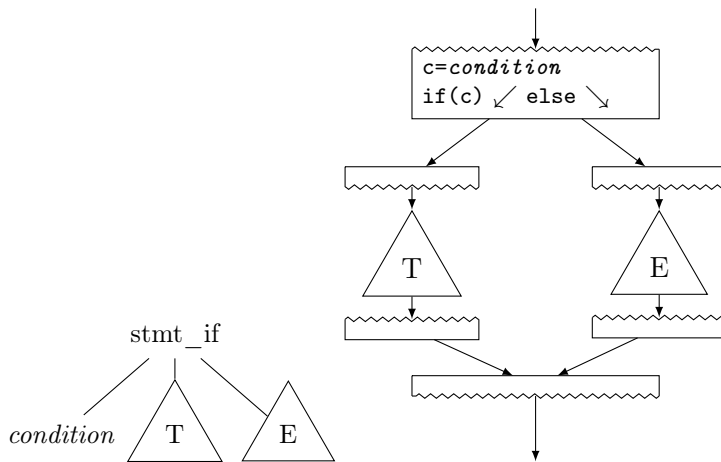
IRValue generateCFG(ASTNode* node, BasicBlock*& insPos) {
  switch (node->kind()) {
  case ASTNode::Function:
    insPos = generatePrologue(node);
    generateCFG(node->child(0), insPos);
    generateEpilogue(insPos);
    return nullptr;
  case ASTNode::Block:
    for (ASTNode* child : node->children())
      generateCFG(child, insPos);
    return nullptr;
  case ASTNode::While: {
    BasicBlock* cond = newBlock();
    BasicBlock* body = newBlock();
    BasicBlock* end = newBlock();
    branchTo(insPos, cond);
    insPos = cond;
    IRValue brcond = generateCFG(node->child(0), insPos);
    // NB: generateCFG can modify insPos
    branchToCond(insPos, brcond, body, end);
    insPos = body;
    generateCFG(node->child(1), insPos);
    branchTo(insPos, cond);
    insPos = end;
    return nullptr;
  }
  }
}
  
```

```

}
// ...
}
}

```

[Slide 75] Build CFG from AST – If Condition



[Slide 76] Build CFG from AST: Switch

Linear search

```

t ← exp
if t == 3: goto B3
if t == 4: goto B4
if t == 7: goto B7
if t == 9: goto B9
goto BD

```

- + Trivial
- Slow, lot of code

Binary search

```

t ← exp
if t == 7: goto B7
elif t > 7:
    if t == 9: goto B9
else:
    if t == 3: goto B3
    if t == 4: goto B4
goto BD

```

- + Good: sparse values
- Even more code

Jump table

```

t ← exp
if 0 ≤ t < 10:
    goto table[t]
goto BD

```

```

table = {
    BD, BD, BD, B3,
    B4, BD, ... }

```

- + Fastest
- Table can be large, needs ind. jump

[Slide 77] Build CFG from AST: Break, Continue, Goto

- break/continue: trivial
 - Keep track of target block, insert branch
- goto: also trivial
 - Split block at target label, if needed
 - But: may lead to irreducible control flow graph (see later)

[Slide 78] CFG: Formal Definition

- Flow graph: $G = (N, E, s)$ with a digraph (N, E) and entry $s \in N$

- Each node is a basic block, s is the entry block
- $(n_1, n_2) \in E$ iff n_2 might be executed immediately after n_1
- All $n \in N$ shall be reachable from s (unreachable nodes can be discarded)
- Nodes without successors are end points

[Slide 79] CFG from C – Example

In-Class Exercise:

Derive the CFG for the these functions. Assume a `switch` instruction exists.

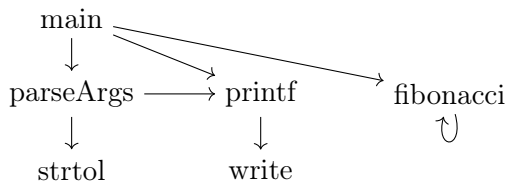
```

int fn1() {
    if (a()) {
        while (b()) {
            c();
            if (d())
                continue;
            e();
        }
    } else {
        f();
    }
}

int fn2() {
    a();
    do switch (c()) {
    case 1:
        while (d()) {
            e();
        }
    case 2:
        f();
    }
    default:
        g();
    } while (h());
    return b();
}
    
```

[Slide 80] Graph IRs: Call Graph

- Graph showing (possible) call relations between functions
- Useful for interprocedural optimizations
 - Function ordering
 - Stack depth estimation
 - ...



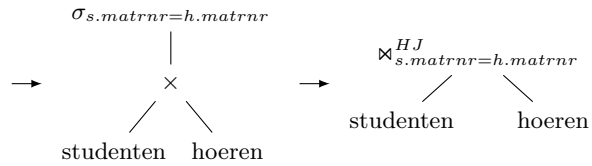
[Slide 81] Graph IRs: Relational Algebra

- Higher-level representation of query plans
 - Explicit data flow
- Allow for optimization and selection actual implementations
 - Elimination of common sub-trees
 - Joins: ordering, implementation, etc.


```

SELECT s.name, h.vorlnr
FROM studenten s, hoeren h
WHERE s.matrnr = h.matrnr

```



[Slide 82] Linear IRs: Stack Machines

- Operands stored on a stack
- Operations pop arguments from top and push result
- Typically accompanied with variable storage
- Generating IR from AST: trivial
- Often used for bytecode, e.g. Java, Python

- + Compact code, easy to generate and implement
- Performance, hard to analyze

```

push 5
push 3
add
pop x
push x
push 1
add
pop y
push 12
pop x
push x
push 1
add
pop z

```

[Slide 83] Linear IRs: Register Machines

- Operands stored in registers
- Operations read and write registers
- Typically: infinite number of registers
- Typically: three-address form

– $dst = src1 \ op \ src2$

- Generating IR from AST: trivial
- E.g., GIMPLE, eBPF, Assembly

```

x ← 5 + 3
y ← x + 1
x ← 12
z ← x + 1
tmp1 ← z - y
return tmp1

```

[Slide 84] Example: High GIMPLE

```
int foo(int n) {
  int res = 1;
  while (n) {
    res *= n * n;
    n -= 1;
  }
  return res;
}
int fac (int n)
gimple_bind < /* <-- still has lexical scopes
  int D.1950;
  int res;

  gimple_assign <integer_cst, res, 1, NULL, NULL>
  gimple_goto <<D.1947>>
  gimple_label <<D.1948>>
  gimple_assign <mult_expr, _1, n, n, NULL>
  gimple_assign <mult_expr, res, res, _1, NULL>
  gimple_assign <plus_expr, n, n, -1, NULL>
  gimple_label <<D.1947>>
  gimple_cond <ne_expr, n, 0, <D.1948>, <D.1946>>
  gimple_label <<D.1946>>
  gimple_assign <var_decl, D.1950, res, NULL, NULL>
  gimple_return <D.1950>
>
$ gcc -fdump-tree-gimple-raw -c foo.c
```

[Slide 85] Example: Low GIMPLE

```
int foo(int n) {
  int res = 1;
  while (n) {
    res *= n * n;
    n -= 1;
  }
  return res;
}
int fac (int n)
{
  int res;
  int D.1950;

  gimple_assign <integer_cst, res, 1, NULL, NULL>
  gimple_goto <<D.1947>>
  gimple_label <<D.1948>>
  gimple_assign <mult_expr, _1, n, n, NULL>
  gimple_assign <mult_expr, res, res, _1, NULL>
  gimple_assign <plus_expr, n, n, -1, NULL>
  gimple_label <<D.1947>>
  gimple_cond <ne_expr, n, 0, <D.1948>, <D.1946>>
  gimple_label <<D.1946>>
  gimple_assign <var_decl, D.1950, res, NULL, NULL>
  gimple_goto <<D.1951>>
  gimple_label <<D.1951>>
  gimple_return <D.1950>
}
$ gcc -fdump-tree-lower-raw -c foo.c
```

[Slide 86] Example: Low GIMPLE with CFG

```
int foo(int n) {
  int res = 1;
  while (n) {
    res *= n * n;
    n -= 1;
  }
  return res;
}
int fac (int n) {
  int res;
  int D.1950;
  <bb 2> :
  gimple_assign <integer_cst, res, 1, NULL, NULL>
  goto <bb 4>; [INV]
  <bb 3> :
  gimple_assign <mult_expr, _1, n, n, NULL>
  gimple_assign <mult_expr, res, res, _1, NULL>
  gimple_assign <plus_expr, n, n, -1, NULL>
  <bb 4> :
  gimple_cond <ne_expr, n, 0, NULL, NULL>
  goto <bb 3>; [INV]
  else
  goto <bb 5>; [INV]
  <bb 5> :
  gimple_assign <var_decl, D.1950, res, NULL, NULL>
  <bb 6> :
gimple_label <<L3>>
  gimple_return <D.1950>
}
$ gcc -fdump-tree-cfg-raw -c foo.c
```

[Slide 87] Linear IRs: Register Machines

- Problem: no clear def–use information
 - Is $x + 1$ the same?
 - Hard to track actual values!
- How to optimize?

⇒ Disallow mutations of variables

```
x    ←  5  +  3
y    ←  x  +  1
x    ←  12
z    ←  x  +  1
tmp1 ←  z  -  y
return tmp1
```

[Slide 88] Single Static Assignment: Introduction

- Idea: disallow mutations of variables, value set in declaration
- Instead: create new variable for updated value

- SSA form: every computed value has a unique definition
 - Equivalent formulation: each name describes result of one operation

<pre> x ← 5 + 3 y ← x + 1 x ← 12 z ← x + 1 tmp1 ← z - y return tmp1 </pre>	→	<pre> v1 ← 5 + 3 v2 ← v1 + 1 v3 ← 12 v4 ← v3 + 1 v5 ← v4 - v2 return v5 </pre>
---	---	--

[Slide 89] Single Static Assignment: Control Flow

- How to handle diverging values in control flow?
- Solution: Φ -nodes to merge values depending on predecessor
 - Value depends on edge used to enter the block
 - All Φ -nodes of a block execute concurrently (ordering irrelevant)

```

entry : x ← ...
        if (x > 2) goto cont
then  : x ← x * 2
cont  : return x
                
```

→

```

entry : v1 ← ...
        if (v1 > 2) goto cont
then  : v2 ← v1 * 2
cont  : v3 ←  $\Phi$ (entry : v1, then : v2)
        return v3
                
```

[Slide 90] Example: GIMPLE in SSA form

```

int foo(int n) {
  int res = 1;
  while (n) {
    res *= n * n;
    n -= 1;
  }
  return res;
}

int fac (int n) { int res, D.1950, _1, _6;
<bb 2> :
gimple_assign <integer_cst, res_4, 1, NULL, NULL>
goto <bb 4>; [INV]
<bb 3> :
gimple_assign <mult_expr, _1, n_2, n_2, NULL>
                
```

```

gimple_assign <mult_expr, res_8, res_3, _1, NULL>
gimple_assign <plus_expr, n_9, n_2, -1, NULL>
<bb 4> :
# gimple_phi <n_2, n_5(D)(2), n_9(3)>
# gimple_phi <res_3, res_4(2), res_8(3)>
gimple_cond <ne_expr, n_2, 0, NULL, NULL>
  goto <bb 3>; [INV]
else
  goto <bb 5>; [INV]
<bb 5> :
gimple_assign <ssa_name, _6, res_3, NULL, NULL>
<bb 6> :
gimple_label <<L3>>
  gimple_return <_6>
}
$ gcc -fdump-tree-ssa-raw -c foo.c

```

[Slide 91] SSA Construction – Local Value Numbering

- Simple case: inside block – keep mapping of variable to value

Code

```

x    ← 5 + 3
y    ← x + 1
x    ← 12
z    ← x + 1
tmp1 ← z - y
return tmp1

```

SSA IR

```

v1 ← add 5, 3
v2 ← add v1, 1
v3 ← const 12
v4 ← add v3, 1
v5 ← sub v4, v2
ret v5

```

Variable Mapping

```

x → v3
y → v2
z → v4
tmp1 → v5

```

[Slide 92] SSA Construction – Across Blocks

- SSA construction with control flow is non-trivial
- Key problem: find value for variable in predecessor
- Naive approach: Φ -nodes for all variables everywhere
 - Create empty Φ -nodes for variables, populate variable mapping
 - Fill blocks (as on last slide)
 - Fill Φ -nodes with last value of variable in predecessor

- Why is this a bad idea? \Rightarrow don't do this!
 - Extremely inefficient, code size explosion, many dead Φ

[Slide 93] SSA Construction – Across Blocks (“simple”¹)

- Key problem: find value in predecessor
- Idea: *seal* block once all direct predecessors are known
 - For acyclic constructs: trivial
 - For loops: seal header once loop block is generated
- Current block not sealed: add Φ -node, fill on sealing
- Single predecessor: recursively query that
- Multiple preds.: add Φ -node, fill now

Confer the (very readable) paper for a more formal specification of the algorithm.
The removal of trivial and redundant Φ -nodes is not strictly required.

[Slide 94] SSA Construction – Example

```
int foo(int n) {
  int res = 1;
  while (n) {
    res *= n * n;
    n -= 1;
  }
  return res;
}
```

```
func foo(v1)
entry: sealed; varmap: n → v1, res → v2
      v2 ← 1
header: sealed; varmap: n →  $\phi_1$ , res →  $\phi_2$ 
       $\phi_1$  ←  $\phi(\text{entry: } v_1, \text{body: } v_6)$ 
       $\phi_2$  ←  $\phi(\text{entry: } v_2, \text{body: } v_5)$ 
      v3 ← equal  $\phi_1$ , 0
      br v3, cont, body
body:  sealed; varmap: n → v6, res → v5
      v4 ← mul  $\phi_1$ ,  $\phi_1$ 
      v5 ← mul  $\phi_2$ , v4
      v6 ← sub  $\phi_1$ , 1
      br header
cont:  sealed; varmap: res →  $\phi_2$ 
      ret  $\phi_2$ 
```

[Slide 95] SSA Construction – Example

¹M Braun et al. “Simple and efficient construction of static single assignment form”. In: *CC*. 2013, pp. 102–122. URL: https://link.springer.com/content/pdf/10.1007/978-3-642-37051-9_6.pdf.

In-Class Exercise:

Construct an IR in SSA form for the following C code.

```
int phis(int a, in b){
  a = a * b;
  if (a > b * b) {
    int c = 1;
    while (a > 0)
      a = a - c;
  } else {
    a = b * b;
  }
  return a;
}
```

[Slide 96] SSA Construction – Pruned/Minimal Form

- Resulting SSA is *pruned* – all ϕ are used
- But not *minimal* – ϕ nodes might have single, unique value
- When filling ϕ , check that multiple real values exist
 - Otherwise: replace ϕ with the single value
 - On replacement, update all ϕ using this value, they might be trivial now, too
- Sufficient? Not for irreducible CFG
 - Needs more complex algorithms² or different construction method³

AD IN2053 “Program Optimization” covers this more formally

[Slide 97] SSA: Implementation

- Value is often just a pointer to instruction
- ϕ nodes placed at beginning of block
 - They execute “concurrently” and on the edges, after all
- Variable number of operands required for ϕ nodes
- Storage format for instructions and basic blocks
 - Consecutive in memory: hard to modify/traverse
 - Array of pointers: $\mathcal{O}(n)$ for a single insertion...
 - Linked List: easy to insert, but pointer overhead

Is SSA a graph IR?

Only if instructions have no side effects, consider `load`, `store`, `call`, ...

These *can* be solved using explicit dependencies as SSA values, e.g. for memory

²M Braun et al. “Simple and efficient construction of static single assignment form”. In: *CC*. 2013, pp. 102–122. URL: https://link.springer.com/content/pdf/10.1007/978-3-642-37051-9_6.pdf.

³R Cytron et al. “Efficiently computing static single assignment form and the control dependence graph”. In: *TOPLAS* 13.4 (1991), pp. 451–490. URL: <https://dl.acm.org/doi/pdf/10.1145/115372.115320>.

[Slide 99] Intermediate Representations – Summary

- An IR is an internal representation of a program
- Main goal: simplify analyses and transformations
- IRs typically based on graphs or linear instructions
- Graph IRs: AST, Control Flow Graph, Relational Algebra
- Linear IRs: stack machines, register machines, SSA
- Single Static Assignment makes data flow explicit
- SSA is extremely popular, although non-trivial to construct

[Slide 100] Intermediate Representations – Questions

- Who designs an IR? What are design criteria?
- Why is an AST not suited for program optimization?
- How to convert an AST to another IR?
- What are the benefits/drawbacks of stack/register machines?
- What benefits does SSA offer over a normal register machine?
- How do ϕ -instructions differ from normal instructions?

4 LLVM-IR

4.1 Overview

[Slide 102] LLVM¹

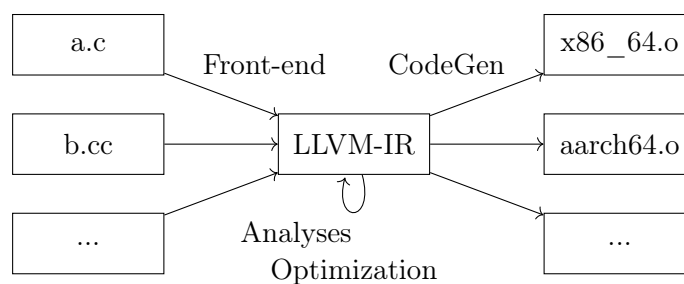
LLVM “Core” Library

- Optimizer and compiler back-end
- “Set of compiler components”
 - IRs: LLVM-IR, SelDag, MIR
 - Analyses and Optimizations
 - Code generation back-ends
- Started from Chris Lattner’s master’s thesis
- Used for C, C++, Swift, D, Julia, Rust, Haskell, ...

LLVM Project

- Umbrella for several projects related to compilers/toolchain
 - LLVM Core
 - Clang: C/C++ front-end for LLVM
 - libc++, compiler-rt: runtime support
 - LLDB: debugger
 - LLD: linker
 - MLIR: experimental IR framework

[Slide 103] LLVM: Overview



- Independent front-end derives LLVM-IR, LLVM does opt. and code gen.
- LTO: dump LLVM-IR into object file, optimize at link-time

¹C Lattner and V Adve. “LLVM: A compilation framework for lifelong program analysis & transformation”. In: *CGO*. 2004, pp. 75–86. URL: <http://www.llvm.org/pubs/2004-01-30-CGO-LLVM.pdf>.

The single IR allows multiple front-ends to reuse the same back-end infrastructure. Thus, generating LLVM-IR provides an easy way to target a wide range of architectures.

For link-time optimization, the LLVM-IR is stored in the object files instead of the machine code. At link-time, a linker plugin detects these files, merges the LLVM-IR from all object files, and then runs the actual compilation as part of the linking step. We will look at LTO again later when discussing object file generation and linking.

4.2 LLVM-IR

[Slide 104] LLVM-IR: Overview

- SSA-based IR, representations textual, bitcode, in-memory
- Hierarchical structure
 - Module
 - Functions, global variables
 - Basic blocks
 - Instructions
- Strongly/strictly typed

```
define dso_local i32 @foo(i32 %0) {
  %2 = icmp eq i32 %0, 0
  br i1 %2, label %10, label %3

3: ; preds = %1, %3
  %4 = phi i32 [ %7, %3 ], [ 1, %1 ]
  %5 = phi i32 [ %8, %3 ], [ %0, %1 ]
  %6 = mul nsw i32 %5, %5
  %7 = mul nsw i32 %6, %4
  %8 = add nsw i32 %5, -1
  %9 = icmp eq i32 %8, 0
  br i1 %9, label %10, label %3

10: ; preds = %3, %1
  %11 = phi i32 [ 1, %1 ], [ %7, %3 ]
  ret i32 %11
}
```

[Slide 105] LLVM-IR: Data types

- First class types:
 - `i<N>` – arbitrary bit width integer, e.g. `i1`, `i25`, `i1942652`
 - `ptr/ptr addrspace(1)` – pointer with optional address space
 - `float/double/half/bfloat/fp128/...`
 - `<N x ty>` – vector type, e.g. `<4 x i32>`
- Aggregate types:
 - `[N x ty]` – constant-size array type, e.g. `[32 x float]`
 - `{ ty, ... }` – struct (can be packed/opaque), e.g. `{i32, float}`

- Other types:
 - `ty (ty, ...)` – function type, e.g. `{i32, i32} (ptr, ...)`
 - `void`
 - `label/token/metadata`

Although structure types can be used in various places in the IR, e.g., a single instruction to load a large structure from memory, this is strongly discouraged: LLVM is not optimized for this and both code quality and compile times get considerably worse. Only use struct types for globals and to implement multiple return values.

[Slide 106] LLVM-IR: Modules

- Top-level entity, one compilation unit – akin to C/C++
- Contains global values, specified with linkage type
- Global variable declarations/definitions


```
@externInt = external global i32, align 4
@globVar = global i32 4, align 4
@staticPtr = internal global ptr null, align 8
```
- Function declarations/definitions


```
declare i32 @readPtr(ptr)
define i32 @return1() {
    ret i32 1
}
```
- Global named metadata (discarded during compilation)

[Slide 107] LLVM-IR: Functions

- Functions definitions contain all code, not nestable
- Single return type (or `void`), multiple parameters, list of basic blocks
 - No basic blocks \Rightarrow function declaration
- Specifiers for `callconv`, section name, other attributes
 - E.g.: `noinline/alwaysinline, noreturn, readonly`
- Parameter and return can also have attributes
 - E.g.: `noalias, nonnull, sret(<ty>)`

[Slide 108] LLVM-IR: Basic Block

- Sequence of instructions
 - ϕ nodes come first
 - Regular instructions come next
 - Must end with a terminator
- First block in function is entry block Entry block cannot be branch target

[Slide 109] LLVM-IR: Instructions – Control Flow and Terminators

- Terminators end a block/modify control flow
 - `ret <ty> <val>/ret void`
 - `br label <dest>/br i1 <cond>, label <then>, label <else>`
 - `switch/indirectbr`
 - `unreachable`
 - Few others for exception handling
- Not a terminator: `call`

Although `call` does modify control flow in some sense, the assumption is that every function call returns ordinarily. When special control flow for exceptions is needed, the `invoke` instruction is used, which specifies one basic block as successor for the ordinary case and one basic block for the exceptional case.

[Slide 110] LLVM-IR: Instructions – Arithmetic-Logical

- `add/sub/mul/udiv/sdiv/urem/srem`
 - Arithmetic uses two's complement
 - Division corner cases are *undefined behavior*
- `fneg/fadd/fsub/fmul/fdiv/frem`
- `shl/lshr/ashr/and/or/xor`
 - Out-of-range shifts have an undefined result
- `icmp <pred>/fcmp <pred>/select <cond>, <then>, <else>`
- `trunc/zext/sext/fptrunc/fpext/fptoui/fptosi/uitofp/sitofp`
- `bitcast`
 - Cast between equi-sized datatypes by reinterpreting bits

Technically, out-of-range shifts return `poison`, see below.

[Slide 111] LLVM-IR: Instructions – Memory and Pointer

- `alloca <ty>` – allocate addressable stack slot
- `load <ty>, ptr <ptr>/store <ty> <val>, ptr <ptr>`
 - May be `volatile` (e.g., MMIO) and/or `atomic`
- `cmpxchg/atomicrmw` – similar to hardware operations
- `ptrtoint/inttoptr`
- `getelementptr` – address computation on `ptr/structs/arrays`

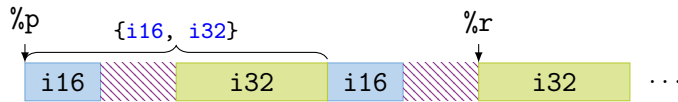
[Slide 112] LLVM-IR: `getelementptr` Examples

- `%r = getelementptr i32, ptr %p, i64 3`



Equivalent in C: `&((int*) p)[3]`

- `%r = getelementptr {i16, i32}, ptr %p, i64 1, i32 1`



Equivalent in C: `&((struct {short _0; int _1;}*) p)[1]._1`

- Also works with nested structs and arrays

[Slide 113] LLVM-IR: undef and poison

- `undef` – unspecified value, compiler may choose any value
 - `%b = add i32 %a, i32 undef → i32 undef`
 - `%c = and i32 %a, i32 undef → i32 %a`
 - `%d = xor i32 %b, i32 %b → i32 undef`
 - `br i1 undef, label %p, label %q → undefined behavior`
- `poison` – result of erroneous operations
 - Delay *undefined behavior* on illegal operation until actually relevant
 - Allows to speculatively “execute” instructions in IR
 - `%d = shl i32 %b, i32 34 → i32 poison`

[Slide 114] LLVM-IR: Intrinsics

- Not all operations provided as instructions
- Intrinsic functions: special functions with defined semantics
 - Replaced during compilation, e.g., with instruction or lib call
- Benefit: no changes needed for parser/bitcode/... on addition
- Examples:
 - `declare iN @llvm.ctpop.iN(iN <src>)`
 - `declare {iN, i1} @llvm.sadd.with.overflow.iN(iN %a, iN %b)`
 - `memcpy, memset, sqrt, returnaddress, ...`

[Slide 115] LLVM-IR: Tools

- clang can emit LLVM-IR bitcode `clang -O -emit-llvm -c test.c -o test.bc`
- llvm-dis disassembles bitcode to textual LLVM-IR `clang -O -emit-llvm -c test.c -o - | llvm-dis`
- llc compiles LLVM-IR (textual or bitcode) to assembly `clang -O -emit-llvm -c test.c -o - | llc clang -O -emit-llvm -c test.c -o - | llvm-dis | llc`

Example Listings omitted – they would span several slides

[Slide 116] LLVM-IR: Example

```
define dso_local <4 x float> @foo2(<4 x float> %0, <4 x float> %1) {
  %3 = alloca <4 x float>, align 16
  %4 = alloca <4 x float>, align 16
  store <4 x float> %0, ptr %3, align 16
  store <4 x float> %1, ptr %4, align 16
  %5 = load <4 x float>, ptr %3, align 16
  %6 = load <4 x float>, ptr %4, align 16
  %7 = fadd <4 x float> %5, %6
  ret <4 x float> %7
}
```

[Slide 117] LLVM-IR: Example

```
define dso_local i32 @foo3(i32 %0, i32 %1) {
  %3 = tail call { i32, i1 } @llvm.smul.with.overflow.i32(i32 %0, i32 %1)
  %4 = extractvalue { i32, i1 } %3, 1
  %5 = extractvalue { i32, i1 } %3, 0
  %6 = select i1 %4, i32 -2147483648, i32 %5
  ret i32 %6
}
```

[Slide 118] LLVM-IR: Example

```
define dso_local i32 @sw(i32 %0) {
  switch i32 %0, label %4 [
    i32 4, label %5
    i32 5, label %2
    i32 8, label %3
    i32 100, label %5
  ]
2: ; preds = %1
  br label %5
3: ; preds = %1
  br label %5
4: ; preds = %1
  br label %5
5: ; preds = %1, %1, %4, %3, %2
  %6 = phi i32 [ %0, %4 ], [ 9, %3 ], [ 32, %2 ], [ 12, %1 ], [ 12, %1 ]
  ret i32 %6
}
```

[Slide 119] LLVM-IR: Example**In-Class Exercise:**

```
@a = private unnamed_addr constant [7 x i32] [i32 12, i32 32, i32 12,
                                             i32 12, i32 9, i32 12, i32 12], align 4

define dso_local i32 @f(i32 %0) {
  %2 = add i32 %0, -4
  %3 = icmp ult i32 %2, 7
  br i1 %3, label %4, label %13
4: ; preds = %1
```

```

%5 = trunc i32 %2 to i8
%6 = lshr i8 83, %5
%7 = and i8 %6, 1
%8 = icmp eq i8 %7, 0
br i1 %8, label %13, label %9
9: ; preds = %4
%10 = sext i32 %2 to i64
%11 = getelementptr @inbounds [7 x i32], ptr @a, i64 0, i64 %10
%12 = load i32, ptr %11, align 4
br label %13
13: ; preds = %1, %4, %9
%14 = phi i32 [ %12, %9 ], [ %0, %4 ], [ %0, %1 ]
ret i32 %14
}

```

4.3 API

[Slide 120] LLVM-IR API

- LLVM offers two APIs: C++ and C
 - C++ is the full API, exposing nearly all internals
 - C API is more limited, but more stable
- Nearly all major versions have breaking changes
- Some support for multi-threading:
 - All modules/types/... associated with an LLVMContext
 - Different contexts may be used in different threads

[Slide 121] LLVM-IR C++ API: Basic Example

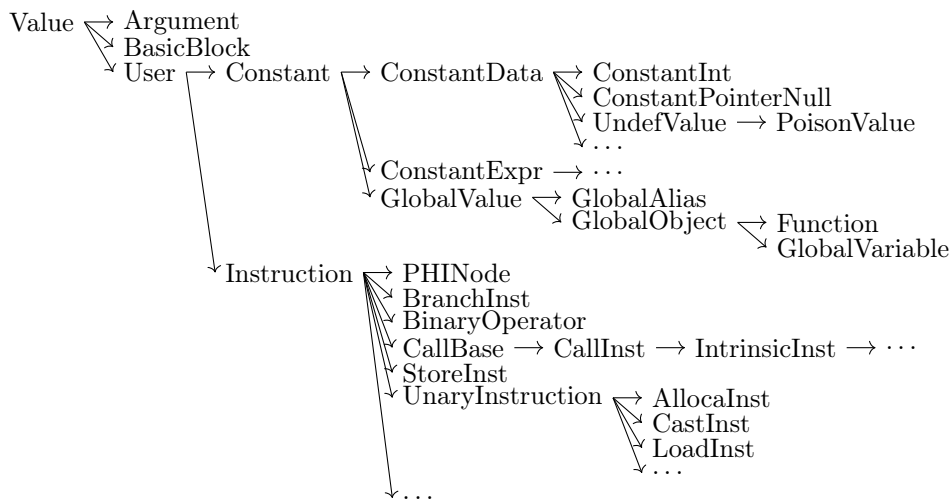
```

#include <llvm/IR/IRBuilder.h>
int main(void) {
    llvm::LLVMContext ctx;
    auto modUP = std::make_unique<llvm::Module>("mod", ctx);

    llvm::Type* i64 = llvm::Type::getInt64Ty(ctx);
    llvm::FunctionType* fnTy = llvm::FunctionType::get(i64, {i64}, false);
    llvm::Function* fn = llvm::Function::Create(fnTy,
        llvm::GlobalValue::ExternalLinkage, "addOne", modUP.get());
    llvm::BasicBlock* entryBB = llvm::BasicBlock::Create(ctx, "entry", fn);

    llvm::IRBuilder<> irb(entryBB);
    llvm::Value* add = irb.CreateAdd(fn->getArg(0), irb.getInt64(1));
    irb.CreateRet(add);
    modUP->print(llvm::outs(), nullptr);
    return 0;
}

```

[Slide 122] LLVM-IR API: Almost Everything is a Value... (excerpt)

See LLVM Doxygen^a for a full graph.

^ahttps://llvm.org/doxygen/classllvm_1_1Value.html

[Slide 123] LLVM-IR API: Programming Environment

- LLVM implements custom RTTI
 - `isa<>`, `cast<>`, `dyn_cast<>`
- LLVM implements a multitude of specialized data structures
 - E.g.: `SmallVector<T, N>` to keep N elements stack-allocated
 - Custom vectors, sets, maps; see manual²
- Preferably uses `ArrayRef`, `StringRef`, `Twine` for references
- LLVM implements custom streams instead of std streams
 - `outs()`, `errs()`, `dbgs()`

Many of these data types are used for efficiency. Standard C++ RTTI is inefficient while LLVM's implementation is very flexible, fast, and has a low memory usage in data structures.

`SmallVector` is preferred over `std::vector` not just because of the inline storage, but also because (for non-`char` types) it only uses 32-bit integers for length/capacity (lower memory usage, often sufficient) and grows more efficiently for trivially movable data structures.

`Twine` is a lazily evaluated string. For example, when specifying `Twine("foo") + 5`, on-stack data structures are constructed to represent this sequence, but the resulting string is constructed only when and if it is actually used. This also allows constructing strings directly into target buffers.

²<https://www.llvm.org/docs/ProgrammersManual.html>

Standard C++ streams are not just inefficient, implementations also tend to inject global constructors in all files. Therefore, LLVM has its own stream implementation. With `raw_svector_ostream` and `raw_string_ostream`, a `raw_ostream` can be used to write into a `SmallVector` or `std::string`.

[Slide 124] LLVM-IR API: Use Tracking

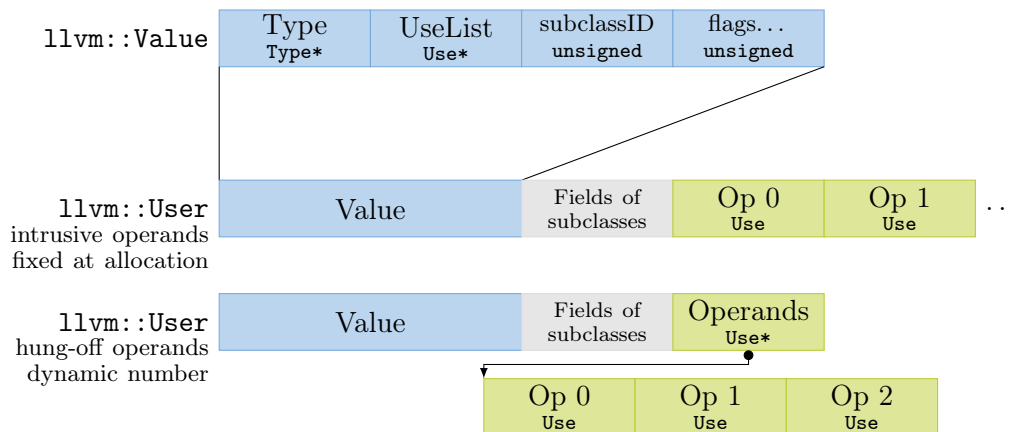
- Values track their users


```
llvm::Value* v = /* ... */;
for (llvm::User* u : v->users())
    if (auto i = llvm::dyn_cast<llvm::Instruction>(u))
        // ...
```
- Simplifies implementation of analyses
- Allows for easy replacement:


```
- inst->replaceAllUsesWith(replVal);
```

4.4 IR Implementation

[Slide 125] LLVM IR Implementation: Value/User

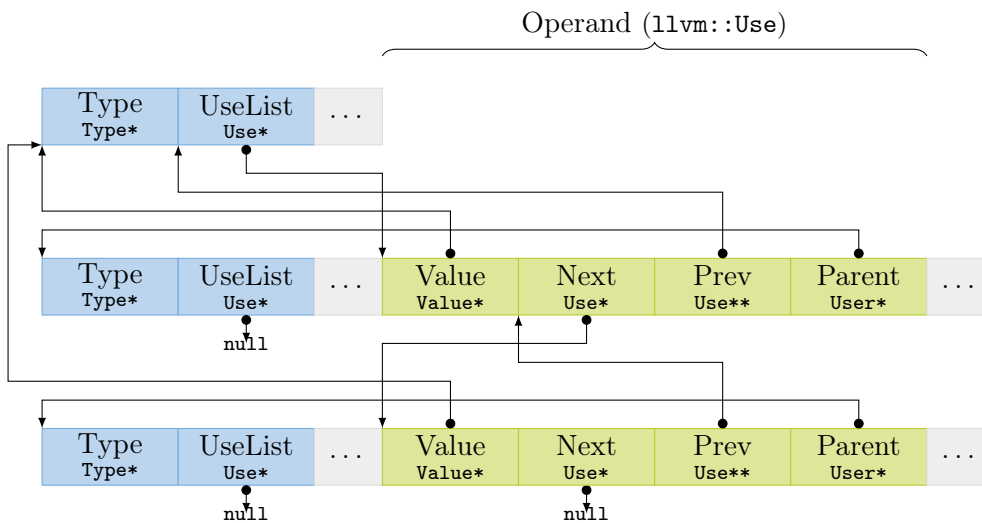


`PHINode` additionally stores n `BasicBlock*` after the operands, but aren't users of blocks.

Every LLVM `Instruction` is a separate heap allocation. As the number of operands is typically known when constructing the instruction, they are allocated after the instruction data structure (this is implemented by `User`).

It can happen that the number of operands increases beyond the allocated storage, for example, when a `PHINode` gets more operands than initially expected. In such cases, the operand list gets hung off into a separate allocation.

As a special case, `PHINode` needs to store the associated `BasicBlocks` in addition to the merged values. The blocks are stored after the operands, but are not operands themselves.

[Slide 126] LLVM IR Implementation: Use

The use list is a doubly-linked list. Starting from `Value::UseList`, one can find all used by following the `Use::Next` pointer. The `Use::Prev` pointer does not point to the previous `Use`, but the previous `Use::Next` pointer or the `Value::UseList` — this way, unlinking does not need to distinguish the special case of the beginning of the use list.

A `Use` also has a pointer to the actual `Value`, so that when inspecting an operand one can actually find the operand itself.

There is also a `Parent` pointer, which points to the `User` which owns the operand: when iterating over the use list, this is the only way to find out which instruction (`User`) uses the value.

In sum, an LLVM-IR operand is quite large, using 32 bytes on a 64-bit system. In addition to every instruction being a separate heap allocation and every operand update requires updating the use list (less data locality), the IR data structures are (in absolute terms) not very efficient — despite being fairly optimized for the use cases they serve.

[Slide 127] LLVM IR Implementation: Instructions/Blocks

- `Instruction` and `BasicBlock` have pointers to parent and next/prev
 - Linked list updated on changes and used for iteration
 - Instructions have cached *order* (integer) for fast “comes before”
- `BasicBlock` successors: blocks used by terminator
- `BasicBlock` predecessors:
 - Iterate over users of block – these are terminators (and `blockaddress`)
 - Ignore non-terminators, parent of using terminator is predecessor
 - Same predecessor might be duplicated (\rightsquigarrow `getUniquePredecessor()`)
- Finding first non- ϕ requires iterating over ϕ -nodes

4.5 IR Design

[Slide 128] LLVM and IR Design

- LLVM provides a decent general-purpose IR for compilers
- But: not ideal for all purposes
 - High-level optimizations difficult, e.g. due to lost semantics
 - Several low-level operations only exposed as intrinsics
 - IR rather complex, high code complexity
 - High compilation times, not very efficient data structures
- Thus: heavy trend towards custom IRs

[Slide 129] IR Design: High-level Considerations

- Define purpose!
- Structure: SSA vs. something else; control flow
 - Control flow: basic blocks/CFG vs. structured control flow
 - Remember: SSA can be considered as a DAG, too
 - SSA is easy to analyse, but non-trivial to construct/leave
- Broader integration: keep multiple stages in single IR?
 - Example: create IR with high-level operations, then incrementally lower
 - Model machine instructions in same IR?
 - Can avoid costly transformations, but adds complexity

[Slide 130] IR Design: Operations

- Data types
 - Simple type structure vs. complex/aggregate types?
 - Keep relation to high-level types vs. low-level only?
 - Virtual data types, e.g. for flags/memory?
- Instruction format
 - Single vs. multiple results?
 - Strongly typed vs. more generic result/operand types?
 - Operand number – fixed vs. dynamic?

[Slide 131] IR Design: Operations

- Allow instruction side effects?
 - E.g.: memory, floating-point arithmetic, implicit control flow
- Operation complexity and abstraction
 - E.g.: `CheckBounds`, `GetStackPtr`, `HashInt128`
 - E.g.: `load` vs. `MOVQconstidx4`
- Extensibility for new operations (e.g., new targets, high-level ops)

[Slide 132] IR Design: Implementation

- Maintain user lists?
 - Simplifies optimizations, but adds considerable overhead
 - Replacement can use `copy` and lazy canonicalization
 - User *count* might be sufficient alternative
- Storage layout: operation size and locations
 - For performance: reduce heap allocations, small data structures
- Special handling for arguments vs. all-instructions?
- Metadata for source location, register allocation, etc.
- SSA: ϕ nodes vs. block arguments?

[Slide 133] IR Example: Go SSA

- Strongly typed
 - Structured types decomposed
- Explicit memory side-effects
- Also High-level operations
 - `IsInBounds`, `VarDef`
- Only one type of value/instruction
 - `Const64`, `Arg`, `Phi`
- No user list, but user count
- Also used for arch-specific repr.

```
env GOSSAFUNC=fac go build test.go
b1:
  v1 (?) = InitMem <mem>
  v2 (?) = SP <uintptr>
  v5 (?) = LocalAddr <*int> {~r1} v2 v1
  v6 (7) = Arg <int> {n} (n[int])
  v8 (?) = Const64 <int> [1] (res[int])
  v9 (?) = Const64 <int> [2] (i[int])
Plain -> b2 (+9)
b2: <- b1 b4
  v10 (9) = Phi <int> v9 v17 (i[int])
  v23 (12) = Phi <int> v8 v15 (res[int])
  v12 (+9) = Less64 <bool> v10 v6
If v12 -> b4 b5 (likely) (9)
b4: <- b2
  v15 (+10) = Mul64 <int> v23 v10 (res[int])
  v17 (+9) = Add64 <int> v10 v8 (i[int])
Plain -> b2 (9)
b5: <- b2
  v20 (12) = VarDef <mem> {~r1} v1
  v21 (+12) = Store <mem> {int} v5 v23 v20
Ret v21 (+12)
```

[Slide 134] LLVM-IR – Summary

- LLVM is a modular compiler framework
- Extremely popular and high-quality compiler back-end

- Primarily provides optimizations and a code generator
- Main interface is the SSA-based LLVM-IR
 - Easy to generate, friendly for writing front-ends/optimizations
- IR design depends on purpose and integration constraints

[Slide 135] LLVM-IR – Questions

- What is the structure of an LLVM-IR module/function?
- Which LLVM-IR data types exist? How do they relate to the target architecture?
- How do semantically invalid operations in LLVM-IR behave?
- What is special about intrinsic functions?
- How to derive LLVM-IR from C code using Clang?
- How does LLVM's `replaceAllUsesWith` work? How could this work without building/maintaining user lists?
- How can an SSA-based IR make side effects explicit?
- How would you design an IR for optimizing Brainfuck?

5 Analyses and Transformations

5.1 Motivation

[Slide 137] Program Transformation: Motivation

- “User code” is often not very efficient
- Also: no need to, compiler can (often?) optimize better
 - More knowledge: e.g., data layout, constants after inlining, etc.
- Allows for more pragmatic/simple code
- Generating “better” IR code on first attempt is expensive
 - What parts are actually used? How to find out?
- Transformation to “better” code must be done *somewhere*
- Optimization is a misnomer: we don’t know whether it improves code!
 - Many transformations are driven by heuristics
- Many types of optimizations are well-known¹

5.2 Dead Code Elimination

[Slide 138] Dead Block Elimination

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

[Slide 139] Optimization Example 1

```
define i32 @fac(i32 %0) {  
  br label %for.header  
for.header: ; preds = %for.body, %1  
  %a = phi i32 [ 1, %1 ], [ %a.new, %for.body ]  
  %b = phi i32 [ 0, %1 ], [ %b.new, %for.body ]  
  %i = phi i32 [ 0, %1 ], [ %i.new, %for.body ]
```

¹FE Allen and J Cocke. *A catalogue of optimizing transformations*. 1971. URL: <https://www.clear.rice.edu/comp512/Lectures/Papers/1971-allen-catalog.pdf>.

```
%cond = icmp sle i32 %i, %0
br i1 %cond, label %for.body, label %exit
for.body: ; preds = %for.header
  %a.new = mul i32 %a, %i
  %b.new = add i32 %b, %i
  %i.new = add i32 %i, 1
  br label %for.header
exit: ; preds = %for.header
  %absum = add i32 %a, %b
  ret i32 %a
}
```

[Slide 140] Simple Dead Code Elimination (DCE)

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

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

[Slide 141] Applying Simple DCE

```
define i32 @fac(i32 %0) {
eff.: cf   br label %for.header
for.header: ; preds = %for.body, %1
users: 2   %a = phi i32 [ 1, %1 ], [ %a.new, %for.body ]
users: 2   %b = phi i32 [ 0, %1 ], [ %b.new, %for.body ]
users: 4   %i = phi i32 [ 0, %1 ], [ %i.new, %for.body ]
users: 1   %cond = icmp sle i32 %i, %0
eff.: cf   br i1 %cond, label %for.body, label %exit
for.body: ; preds = %for.header
users: 1   %a.new = mul i32 %a, %i
users: 1   %b.new = add i32 %b, %i
users: 1   %i.new = add i32 %i, 1
eff.: cf   br label %for.header
exit: ; preds = %for.header
users: 0   %absum = add i32 %a, %b
eff.: cf   ret i32 %a
}
```

In this example, the instruction `%absum` can be removed. This reduces the number of users of `%a` and `%b` by 1. As no other instructions have a user count of 0 after this change, the algorithm terminates.

[Slide 142] Dead Code Elimination

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

[Slide 143] Applying Liveness-based DCE

```

define i32 @fac(i32 %0) {
live  br1 label %for.header
for.header: ; preds = %for.body, %1
live  %a = phi i32 [ 1, %1 ], [ %a.new, %for.body ]

live  %i = phi i32 [ 0, %1 ], [ %i.new, %for.body ]
live  %cond = icmp sle i32 %i, %0
live  br2 i1 %cond, label %for.body, label %exit
for.body: ; preds = %for.header
live  %a.new = mul i32 %a, %i

live  %i.new = add i32 %i, 1
live  br3 label %for.header
exit: ; preds = %for.header

live  ret i32 %a
}

```

Work list (stack)

This algorithm finds the dead value cycle of %b from the previous example. (Refer to the slide deck for the animated version.)

[Slide 144] Liveness-based DCE: Work List Implementation**In-Class Exercise:**

- What operations are performed on a work list?
 - Insert instruction
 - Remove any instruction
 - Test whether instruction is contained
 - Get and remove next instruction to handle

- How to implement an efficient work list?

[Slide 145] Optimization Example 2

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

[Slide 146] Common Subexpression Elimination (CSE) – Attempt 1

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

[Slide 147] CSE Attempt 1 – Example 1

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

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

[Slide 148] CSE Attempt 1 – Example 2

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

```
→ ht      ret i32 %res
          }
```

Instruction does not dominate all uses! error: input module is broken!

5.3 Dominator Tree

[Slide 149] Domination

- Remember: CFG $G = (N, E, s)$ with digraph (N, E) and entry $s \in N$
- Dominate: $d \text{ dom } n$ iff every path from s to n contains d
 - Dominators of n : $DOM(n) = \{d \mid d \text{ dom } n\}$
- Strictly dominate: $d \text{ sdom } n \Leftrightarrow d \text{ dom } n \wedge d \neq n$
- Immediate dominator: $\text{idom}(n) = d : d \text{ sdom } n \wedge \nexists d'. d \text{ sdom } d' \wedge d' \text{ sdom } n$

⇒ All strict dominators are always executed before the block

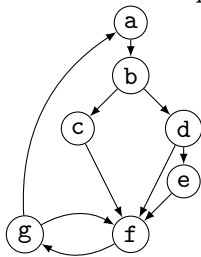
⇒ All values from dominators available/usable

⇒ All values not from dominators **not** usable

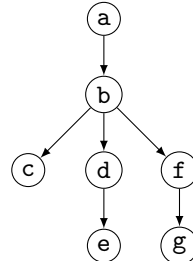
[Slide 150] Dominator Tree

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

Control Flow Graph



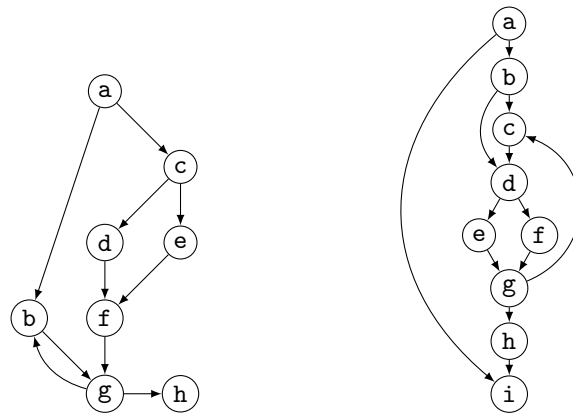
Dominator Tree



[Slide 151] Dominator Tree – Example

In-Class Exercise:

Construct the dominator tree for the following CFGs (entry at a):

**[Slide 152] Dominator Tree: Construction**

- Naive: inefficient (but reasonably simple)²
 - For each block: find a path from the root – superset of dominators
 - Remove last block on path and check for alternative path
 - If no alternative path exists, last block is idom
- Lengauer–Tarjan: more efficient methods³
 - Simple method in $\mathcal{O}(m \log n)$; sophisticated method in $\mathcal{O}(m \cdot \alpha(m, n))$ ($\alpha(m, n)$ is the inverse Ackermann function, grows *extremely* slowly)
 - Used in some compilers⁴
- Semi-NCA: $\mathcal{O}(n^2)$, but lower constant factors⁵

Most notable, LLVM doesn't use the Lengauer–Tarjan algorithm. Instead, they use the Semi-NCA algorithm, which has $\mathcal{O}(n^2)$ runtime, but lower constant factors and is therefore substantially faster for certain (typical) inputs^a.

^aJ Kuderski. “Dominator Trees and incremental updates that transcend times”. In: *LLVM Dev Meeting*. Oct. 2017. URL: https://llvm.org/devmtg/2017-10/slides/Kuderski-Dominator_Trees.pdf.

[Slide 153] Dominator Tree: Implementation

- Per node store: *idom*, idom-children, DFS pre-order/post-order number
- Get immediate dominator: ...lookup *idom*
- Iterate over all dominators/dominated by: ...trivial

²ES Lowry and CW Medlock. “Object code optimization”. In: *CACM* 12.1 (1969), pp. 13–22. URL: <https://dl.acm.org/doi/pdf/10.1145/362835.362838>.

³T Lengauer and RE Tarjan. “A fast algorithm for finding dominators in a flowgraph”. In: *TOPLAS* 1.1 (1979), pp. 121–141. URL: <https://dl.acm.org/doi/pdf/10.1145/357062.357071>

⁴Example: <https://github.com/WebKit/WebKit/blob/aabfacb/Source/WTF/wtf/Dominators.h>

⁵L Georgiadis. “Linear-Time Algorithms for Dominators and Related Problems”. PhD thesis. Princeton University, Nov. 2005

- Check whether a sdom b ⁶
 - $a.preNum < b.preNum \wedge a.postNum > b.postNum$
 - After updates, numbers might be invalid: recompute or walk tree
- Problem: dominance of unreachable blocks ill-defined \rightsquigarrow special handling

5.4 Common Subexpression Elimination

[Slide 154] CSE Attempt 2

- Option 1:
 - For identical instructions, store all
 - Add dominance check before replacing
 - Visit nodes in reverse post-order (i.e., topological order)
- Option 2:⁷
 - Do a DFS over dominator tree
 - Use scoped hashmap to track available values

Does this work? Yes.

[Slide 155] CSE: Hashing an Instruction (and Beyond)

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

[Slide 156] Global Value Numbering – or: advanced CSE

- Hash-based approach only catches trivially removable duplicates
- Alternative: partition values into *congruence classes*
 - Congruent values are guaranteed to always have the same value
- Optimistic approach: values are congruent unless proven otherwise
- Pessimistic approach: values are not congruent unless proven
- Combinable with: reassociation, DCE, constant folding
- Rather complex, but can be highly beneficial⁸

⁶PF Dietz. “Maintaining order in a linked list”. In: *STOC*. 1982, pp. 122–127. URL: <https://dl.acm.org/doi/pdf/10.1145/800070.802184>.

⁷P Briggs, KD Cooper, and LT Simpson. *Value numbering*. Tech. rep. CRPC-TR94517-S. Rice University, 1997. URL: <https://www.cs.rice.edu/~keith/Promo/CRPC-TR94517.pdf.gz>.

⁸K Gargi. “A sparse algorithm for predicated global value numbering”. In: *PLDI*. 2002, pp. 45–56.

5.5 Simple Transformations

[Slide 157] Simple Transformations: Inlining

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

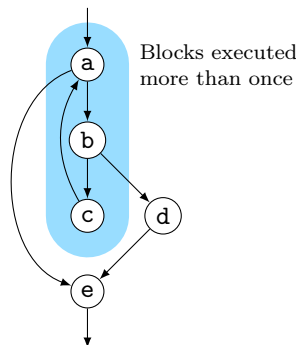
[Slide 158] Simple Transformations: Mem2Reg and SROA

- Mem2reg: promote `alloca` to SSA values/`phis`
 - Condition: only `load/store`, no address taken
 - Essentially just SSA construction
 - Not run in default pipeline, subsumed by SROA
- SROA: scalar replacement of aggregate
 - Separate structure fields into separate variables
 - Also promote them to SSA

5.6 Loop Analysis

[Slide 159] What is a Loop?

```
void func() {
  while (a()) {
    if (b()) {
      d();
      break;
    }
    c();
  }
  e();
}
```



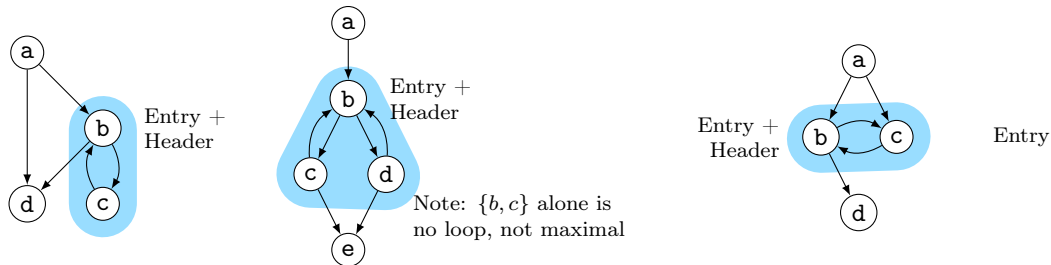
- Loops in source code \neq loops in CFG
 - `d` is *not* part of loop: executed at most once
- \rightsquigarrow Need algorithm to find loops in CFG

[Slide 160] Loops

- Loop: maximal SCC L with at least one internal edge⁹ (strongly connected component (SCC): all blocks reachable from each other)

⁹P Havlak. "Nesting of reducible and irreducible loops". In: *TOPLAS* 19.4 (1997), pp. 557–567. URL: <https://dl.acm.org/doi/pdf/10.1145/262004.262005>.

- Entry: block with an edge from outside of L
- Header h : first entry found (might be ambiguous)
- Loop nested in L : loop in subgraph $L \setminus \{h\}$

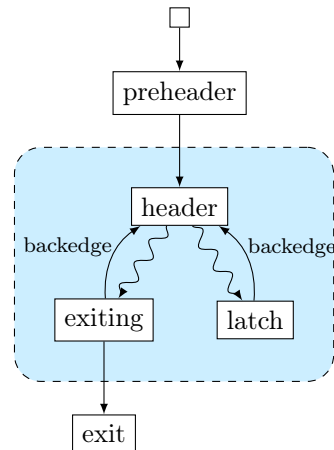


[Slide 161] Natural Loops

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

Formal Definition

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



[Slide 162] Finding Natural Loops

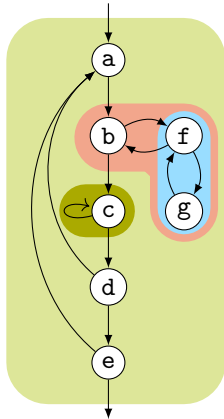
- Modified version¹⁰ of Tarjan's algorithm¹¹
- Iterate over dominator tree in post order
- Each block: find predecessors dominated by the block
 - None \rightsquigarrow no loop header, continue
 - Any \rightsquigarrow loop header, these edges *must* be backedges
- Walk through predecessors until reaching header again
 - All blocks on the way must be part of the loop body
 - Might encounter nested loops, update loop parent

¹⁰G Ramalingam. "Identifying loops in almost linear time". In: *TOPLAS* 21.2 (1999), pp. 175–188. URL: <https://dl.acm.org/doi/pdf/10.1145/316686.316687>.

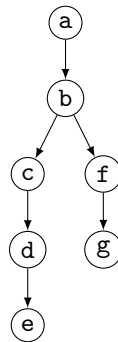
¹¹R Tarjan. "Testing flow graph reducibility". In: *STOC*. 1973, pp. 96–107. URL: <https://dl.acm.org/doi/pdf/10.1145/800125.804040>.

[Slide 163] Finding Natural Loops: Example

Control Flow Graph



Dominator Tree



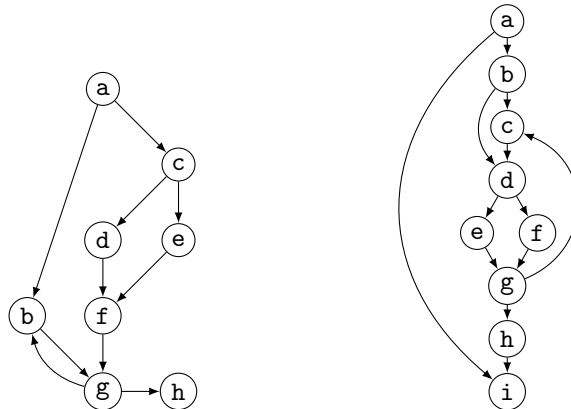
Loop Info

- Loop **A**: {c}
header: c; parent: D
- Loop **B**: {f,g}
header: f; parent: C
- Loop **C**: {b,f,g}
header: b; parent: D
- Loop **D**: {a,b,c,d,e,f,g}
header: a; parent: NULL

[Slide 164] Loop Analysis – Example

In-Class Exercise:

Apply the previous algorithm to find loops in the following CFGs (entry at a):



[Slide 165] Loop Invariant Code Motion (LICM)

- Analyze loops, iterate over loop tree in post-order
 - I.e., visit inner loops first
- ↑ Hoist:¹² iterate over blocks of loop in reverse post-order
 - For each movable inst., check for loop-defined operands
 - If not, move to preheader (create one, if not existent)
 - Otherwise, add inst. to set of values defined inside loop
- ↓ Sink: Iterate over blocks of loop in post-order

¹²<https://github.com/bytecodealliance/wasmtime/blob/bd6fe11/cranelift/codegen/src/licm.rs>

- For each movable inst., check for users inside loop
- If none, move to unique exit (if existent)

5.7 LLVM Passes

[Slide 166] Transformations and Analyses in LLVM: Passes

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

[Slide 167] Using LLVM (New) Pass Manager

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

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

[Slide 168] Writing a Pass for LLVM's New PM – Part 1

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

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

```
    llvm::errs() << F.getName() << "\n";
    return llvm::PreservedAnalyses::all();
}
};
// ...
```

[Slide 169] Writing a Pass for LLVM's New PM – Part 2

```
extern "C" ::llvm::PassPluginLibraryInfo LLVM_ATTRIBUTE_WEAK
llvmGetPassPluginInfo() {
    return { LLVM_PLUGIN_API_VERSION, "TestPass", "v1",
            [] (llvm::PassBuilder &PB) {
                PB.registerPipelineParsingCallback(
                    [] (llvm::StringRef Name, llvm::FunctionPassManager &FPM,
                        llvm::ArrayRef<llvm::PassBuilder::PipelineElement>) {
                        if (Name == "testpass") {
                            FPM.addPass(TestPass());
                            return true;
                        }
                        return false;
                    });
            } };
}
++ -shared -o testpass.so testpass.cc -LLVM -fPIC
opt -S -load-pass-plugin=$PWD/testpass.so -passes=testpass input.ll
```

[Slide 170] Analyses and Transformations – Summary

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

[Slide 171] Analyses and Transformations – Questions

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