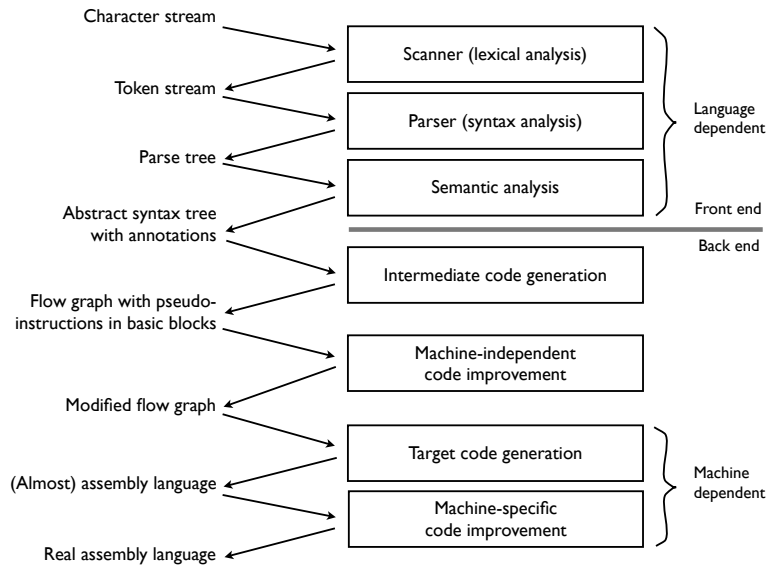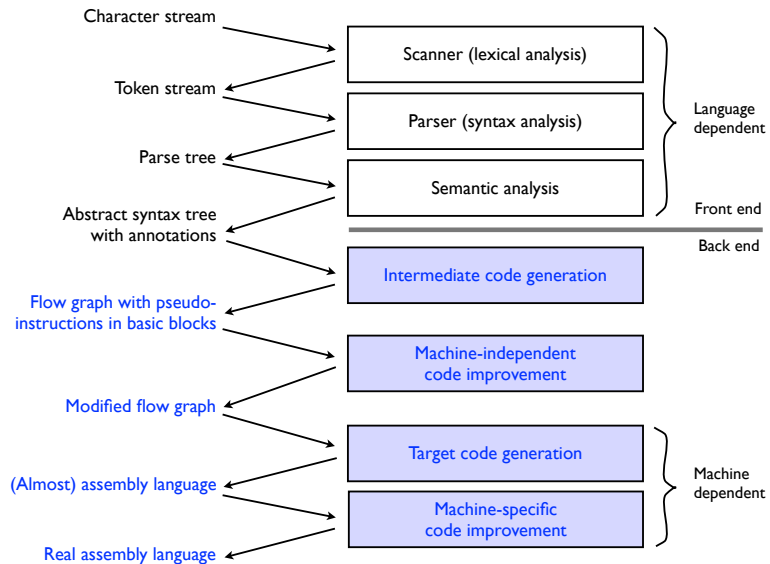# Building a runnable program
## Chapter 14

# Phases of compilation

- Front end
  - analyses source code
  - is language dependent
  - fairly uniform
- Back end
  - produces target code
  - machine dependent
  - very non-uniform
- Intermediate
  - between the front and back end
  - independent of language and machine

# Structure of compiler

Character stream

Scanner (lexical analysis)

Token stream

Parser (syntax analysis)

Parse tree

Semantic analysis

Abstract syntax tree
with annotations

Language
dependent

Front end

Back end

Intermediate code generation

Flow graph with pseudo-
instructions in basic blocks

Machine-independent
code improvement

Modified flow graph

Target code generation

(Almost) assembly language

Machine-specific
code improvement

Real assembly language

Machine
dependent

---

# Structure of compiler

Character stream

Scanner (lexical analysis)

Token stream

Parser (syntax analysis)

Parse tree

Semantic analysis

Abstract syntax tree
with annotations

Language
dependent

Front end

Back end

Intermediate code generation

Flow graph with pseudo-
instructions in basic blocks

Machine-independent
code improvement

Modified flow graph

Target code generation

(Almost) assembly language

Machine-specific
code improvement

Real assembly language

Machine
dependent

# Control flow graph

- Nodes
  - correspond to "basic blocks" of the syntax tree
  - a basic block is a maximal-length sequence of operations with no branches in or out
- Edges
  - represent interblock control flow
- Operations in blocks
  - instructions of an idealised machine, with an unlimited number of registers ("virtual" registers).

# Machine-independent code improvement

- Within nodes
  - eliminate redundant loads, stores and arithmetic computations
- Between nodes
  - remove unnecessary repeated computations
    - e.g. an expression in a loop, the value of which will not change
  - these improvements can cause restructuring of the graph

# Target code generation

- Translates nodes into instructions of target machine

- Strings these together with appropriate branches generated from graph edges

- Still relies on virtual registers

- Target code generators can be generated automatically from a formal description of a target machine

# Machine-specific code improvement

- Register allocation

  - maps unlimited virtual registers to physical registers

  - if there are insufficient physical registers, generates extra loads and stores to maintain virtual registers in real memory

- Instruction scheduling

  - reorder instructions to keep pipeline as full as possible

# Intermediate Forms

- High-level
  - based on trees (e.g. syntax trees), directed acyclic graphs, or stacks
- Medium-level
  - control flow graphs containing pseudo-assembly language instructions, quadruples consisting of
    - two operands
    - an operator
    - a destination
- Low-level
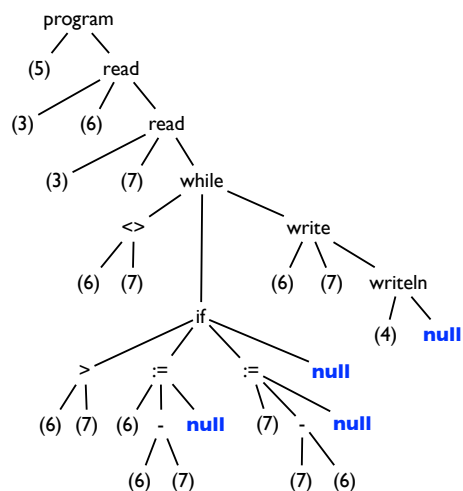  - assembly language of target machine.

# Choice of intermediate form

- Single language/single target
  - distinction between front and back end less clear
  - code improvement more at the lower levels
- Single language/several targets
  - do as much code improvement as possible on a high- or medium-level intermediate form
- Multilanguage compiler family
  - Intermediate form is independent of both language and target
  - For $n$ languages, $m$ machines - $n$ front ends, $m$ back ends instead of $n \times m$ different compilers
  - e.g. GCC compilers on Unix

## Intermediate form representation

- Intermediate form needs linear representation for storing in a file
    - tree-based IFs
        - linearised by ordered traversal
    - control flow graphs
        - replace pointers with offsets from the beginning of the file

## Syntax Tree and Symbol Table



```
program gcd (input, output);
var i, j : integer;
begin
    read (i, j);
    while i <> j do
        if i > j then i := i - j
        else j := j - i;
    writeln (i)
end
```

| Index | Symbol | Type |
|-------|---------|---------|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

# Control flow graph
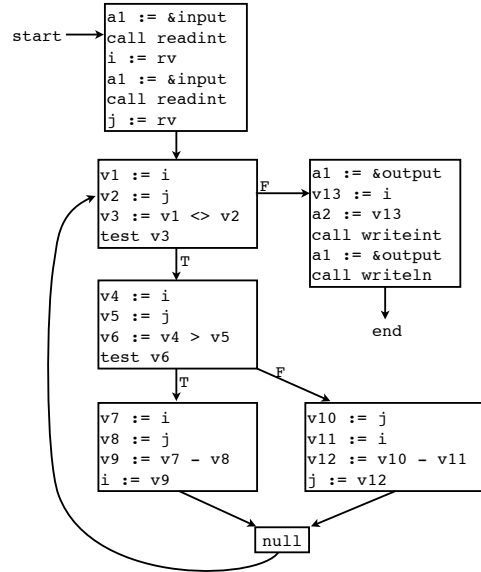
```
program gcd (input, output);
var i, j : integer;
begin
    read (i, j);
    while i <> j do
        if i > j then i := i - j
        else j := j - i;
    writeln (i)
end
```

start →
```
a1 := &input
call readint
i := rv
a1 := &input
call readint
j := rv
```

```
v1 := i
v2 := j
v3 := v1 <> v2
test v3
```
F →
```
a1 := &output
v13 := i
a2 := v13
call writeint
a1 := &output
call writeln
```
→ end

T
```
v4 := i
v5 := j
v6 := v4 > v5
test v6
```

T
```
v7 := i
v8 := j
v9 := v7 - v8
i := v9
```

F
```
v10 := j
v11 := i
v12 := v10 - v11
j := v12
```
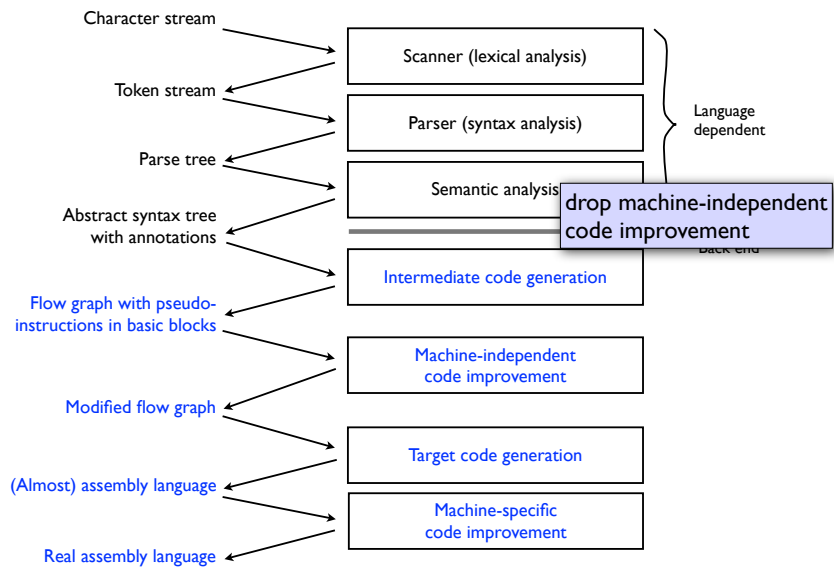
```
null
```

- v1, v2,... are general-purpose virtual registers
- a1, a2 are argument registers
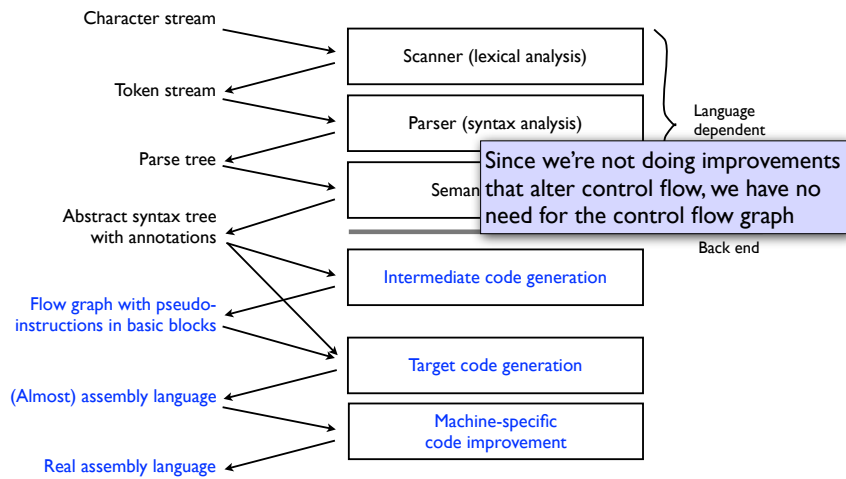- rv is a return-value register

# Code Generation

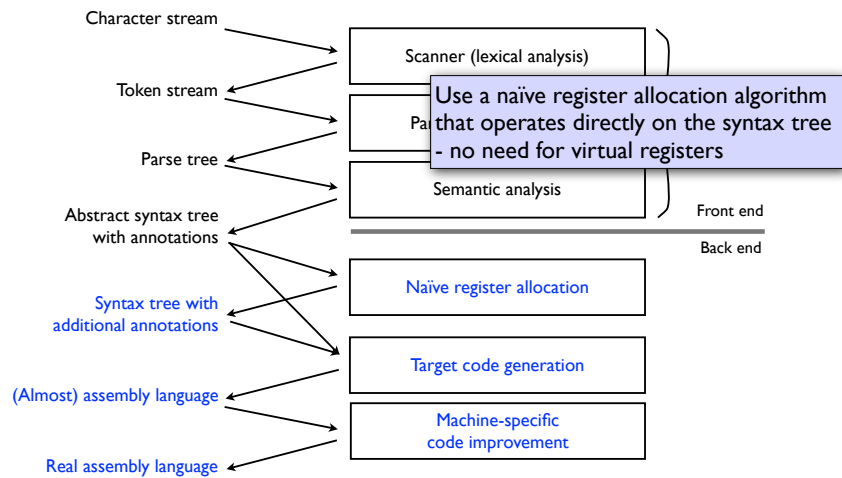- Back end in our earlier compiler structure diagram is complex, so we focus on a simpler model, as follows…

# Structure of compiler

Character stream

Token stream

Parse tree

Abstract syntax tree
with annotations

Flow graph with pseudo-
instructions in basic blocks

Modified flow graph

(Almost) assembly language

Real assembly language

Scanner (lexical analysis)

Parser (syntax analysis)

Semantic analysis

Language
dependent

drop machine-independent
code improvement

Back end

Intermediate code generation

Machine-independent
code improvement

Target code generation

Machine-specific
code improvement

# Structure of compiler

Character stream

Token stream

Parse tree

Abstract syntax tree
with annotations

Flow graph with pseudo-
instructions in basic blocks

(Almost) assembly language

Real assembly language

Scanner (lexical analysis)

Parser (syntax analysis)

Seman

Language
dependent

Since we're not doing improvements
that alter control flow, we have no
need for the control flow graph

Back end

Intermediate code generation

Target code generation

Machine-specific
code improvement

# Structure of compiler

Character stream

Token stream

Parse tree

Abstract syntax tree
with annotations

Syntax tree with
additional annotations

(Almost) assembly language

Real assembly language

| Scanner (lexical analysis) |
| Parser (syntax analysis) |
| Semantic analysis |

Front end

Back end

| Naïve register allocation |
| Target code generation |
| Machine-specific code improvement |

> Use a naïve register allocation algorithm that operates directly on the syntax tree - no need for virtual registers

---

# Structure of compiler

Character stream

Token stream

Parse tree

Abstract syntax tree
with annotations

Syntax tree with
additional annotations

(Almost) **Assembly language**

Real assembly language

| Scanner (lexical analysis) |
| Parser (syntax analysis) |
| Semantic analysis |

Front end

Back end

| Naïve register allocation |
| Target code generation |
| Machine-specific code improvement |

> No virtual registers, so no need to map to real registers.
> Drop instruction scheduling.

# Structure of compiler

Character stream

Scanner (lexical analysis)

Token stream

Parser (syntax analysis)

Parse tree

Semantic analysis

Abstract syntax tree
with annotations

Language
dependent

Front end
—————————————————
Back end

Naïve register allocation

Syntax tree with
additional annotations

Target code generation

Assembly language

# Code Generation

- From symbol table

  - Code for storage management
    - stack frame offsets for local variables and parameters
    - space for global variables

| Index | Symbol   | Type    |
|-------|----------|---------|
| 1     | integer  | type    |
| 2     | textfile | type    |
| 3     | input    | 2       |
| 4     | output   | 2       |
| 5     | gcd      | program |
| 6     | i        | 1       |
| 7     | j        | 1       |

```
        .data
i:      .word
j:      .word
        .text
```
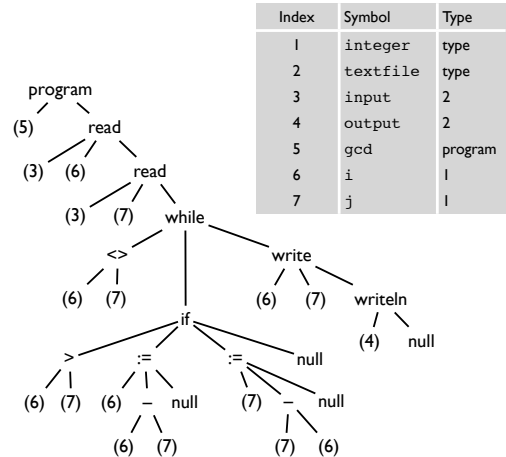
# Syntax tree to code

- Uses a *tree grammar*
  - similar to a CFG, but describes the structure of syntax trees
  - used, with addition of attributes and rules, to annotate a syntax tree
    - see earlier discussion about uses for attribute grammars
- "A : B" on the LHS of a production means that A is a kind of B, and can appear anywhere a B is expected in a RHS
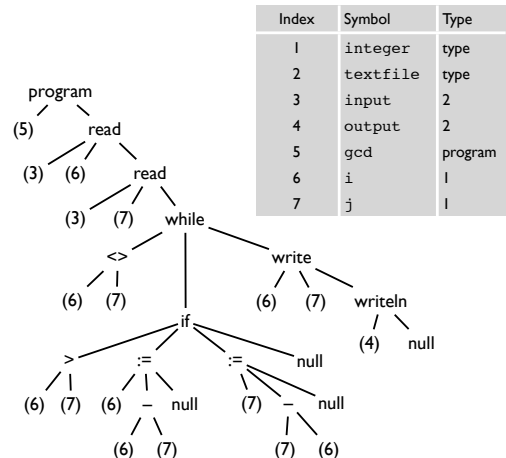
# Syntax tree to code

- Tree grammar for the GCD example
  - Synthesized attributes
    - code - in program, expr, stmt
    - stp - in id, points to symbol table entry for identifier
    - reg - in expr, points to register
  - Inherited attribute
    - nfree_reg - in expr, stmt, is the next register available
  - Registers
    - assume a fixed and limited set
    - special purpose, for arguments (a1,a2…), return value (rv), stack pointer (sp) etc.
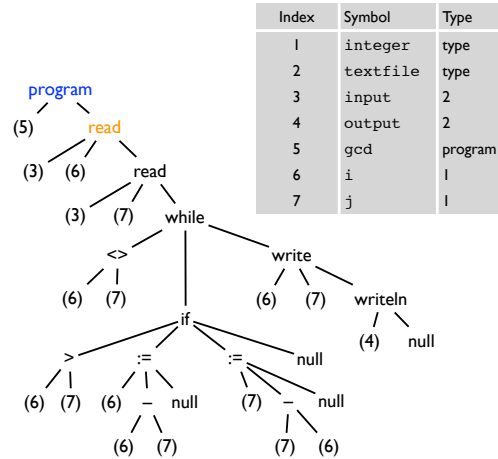    - general purpose: r1,…,r$k$

| Index | Symbol   | Type    |
|-------|----------|---------|
| 1     | integer  | type    |
| 2     | textfile | type    |
| 3     | input    | 2       |
| 4     | output   | 2       |
| 5     | gcd      | program |
| 6     | i        | 1       |
| 7     | j        | 1       |

---
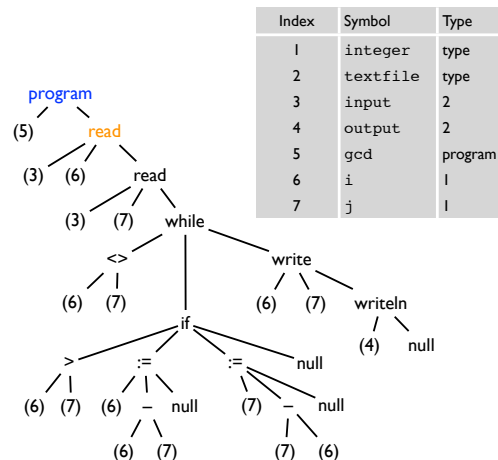
program → id stmt
- ➤ stmt.nfree_reg := 0
- ➤ program.code := ["main:"] + stmt.code + ["goto exit"]
- ➤ program.name := id.stp → name

| Index | Symbol   | Type    |
|-------|----------|---------|
| 1     | integer  | type    |
| 2     | textfile | type    |
| 3     | input    | 2       |
| 4     | output   | 2       |
| 5     | gcd      | program |
| 6     | i        | 1       |
| 7     | j        | 1       |

## Slide 25

program → id stmt
- stmt.nfree_reg := 0
- program.code := ["main:"] + stmt.code + ["goto exit"]
- program.name := id.stp → name

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |



**stmt**

| | nfree_reg | code |
|---|---|---|
| | 0 | |

---

## Slide 26

program → id stmt
- stmt.nfree_reg := 0
- program.code := ["main:"] + stmt.code + ["goto exit"]
- program.name := id.stp → name

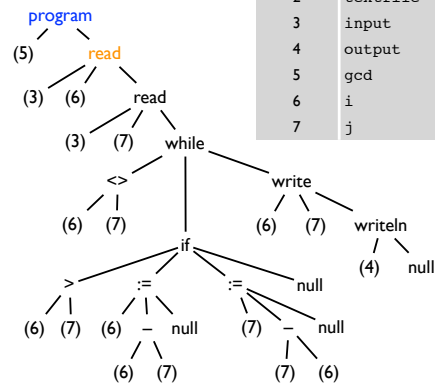| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |



**stmt**

| | nfree_reg | code |
|---|---|---|
| | 0 | |

**program**

| | name | code |
|---|---|---|
| | | main:<br>stmt.code<br>goto exit |

## Slide (page 27)

$$\text{program} \rightarrow \text{id stmt}$$

➤ stmt.nfree_reg := 0
➤ program.code := ["main:"] + stmt.code + ["goto exit"]
➤ program.name := id.stp → name

| Index | Symbol | Type |
|-------|----------|---------|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

program
(5)  read
(3) (6)  read
(3) (7)  while
<>        write
(6) (7)   (6) (7)  writeln
if                  (4)  null
>    :=    :=   null
(6) (7) (6) – null  (7) – null
(6) (7)             (7) (6)

**stmt**

| nfree_reg | code |
|-----------|------|
| 0 | |

**program**

| name | code |
|------|------|
| gcd | main:<br>    stmt.code<br>    goto exit |

## Slide (page 28)

$$\text{read:stmt}_1 \rightarrow \text{id}_1 \ \text{id}_2 \ \text{stmt}_2$$

➤ $\text{stmt}_2.\text{nfree\_reg} := \text{stmt}_1.\text{nfree\_reg}$
➤ $\text{stmt}_1.\text{code} := [\text{"a1} := \&\text{"} \ \text{id}_1.\text{stp} \rightarrow \text{name}]$
   $+ [\text{"call" if } \text{id}_2.\text{stp} \rightarrow \text{type} = \text{int then "readint" else} \ldots]$
   $+ [\text{id}_2.\text{stp} \rightarrow \text{name "} := \text{rv"}] + \text{stmt}_2.\text{code}$

| Index | Symbol | Type |
|-------|----------|---------|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

program
(5)  read
(3) (6)  read
(3) (7)  while
<>        write
(6) (7)   (6) (7)  writeln
if                  (4)  null
>    :=    :=   null
(6) (7) (6) – null  (7) – null
(6) (7)             (7) (6)

**stmt**

| nfree_reg | code |
|-----------|------|
| 0 | |

**program**

| name | code |
|------|------|
| gcd | main:<br>    stmt.code<br>    goto exit |

# Slide 29

$read{:}stmt_1 \rightarrow id_1\ id_2\ stmt_2$

➤ $stmt_2.nfree\_reg := stmt_1.nfree\_reg$

➤ $stmt_1.code := [\text{"a1} := \&\text{"}\ id_1.stp{\rightarrow}name]$
  $+ [\text{"call"}\ \text{if}\ id_2.stp{\rightarrow}type = int\ \text{then "readint" else} \ldots]$
  $+ [id_2.stp{\rightarrow}name\ \text{":= rv"}] + stmt_2.code$

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |



stmt

| | nfree_reg | code |
|---|---|---|
| | 0 | |

stmt

| | nfree_reg | code |
|---|---|---|
| | 0 | |

program

| | name | code |
|---|---|---|
| | gcd | main:<br>    stmt.code<br>    goto exit |

---

# Slide 30

$read{:}stmt_1 \rightarrow id_1\ id_2\ stmt_2$

➤ $stmt_2.nfree\_reg := stmt_1.nfree\_reg$

➤ $stmt_1.code := [\text{"a1} := \&\text{"}\ id_1.stp{\rightarrow}name]$
  $+ [\text{"call"}\ \text{if}\ id_2.stp{\rightarrow}type = int\ \text{then "readint" else} \ldots]$
  $+ [id_2.stp{\rightarrow}name\ \text{":= rv"}] + stmt_2.code$

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |



stmt

| | nfree_reg | code |
|---|---|---|
| | 0 | |

stmt

| | nfree_reg | code |
|---|---|---|
| | 0 | |

program

| | name | code |
|---|---|---|
| | gcd | main:<br>    stmt.code<br>    goto exit |

## Slide 31

$read{:}stmt_1 \rightarrow id_1\ id_2\ stmt_2$

- ➤ $stmt_2.nfree\_reg := stmt_1.nfree\_reg$
- ➤ $stmt_1.code := [\text{"a1 := \&"}\ id_1.stp{\rightarrow}name]$
  $+ [\text{"call" if } id_2.stp{\rightarrow}type = int \text{ then "readint" else } \ldots]$
  $+ [id_2.stp{\rightarrow}name\ \text{":= rv"}] + stmt_2.code$

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | a1 := &input<br>call readint<br>i := rv<br>stmt.code |

**program**

| name | code |
|---|---|
| gcd | main:<br>stmt.code<br>goto exit |

## Slide 32

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | a1 := &input<br>call readint<br>i := rv<br>stmt.code |

**program**

| name | code |
|---|---|
| gcd | main:<br>stmt.code<br>goto exit |

## Slide 33

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | a1 := &input<br>call readint<br>i := rv<br>stmt.code |

**program**

| name | code |
|---|---|
| gcd | main:<br><br><br>stmt.code<br>goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

## Slide 34

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**program**

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>stmt.code<br>goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

## Slide 35

read:stmt$_1$ → id$_1$ id$_2$ stmt$_2$

➤ stmt$_2$.nfree_reg := stmt$_1$.nfree_reg
➤ stmt$_1$.code := ["a1 := &" id$_1$.stp→name]
  + ["call" if id$_2$.stp→type = int then "readint" else …]
  + [id$_2$.stp→name ":= rv"] + stmt$_2$.code

| Index | Symbol | Type |
|-------|---------|---------|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**stmt**

| nfree_reg | code |
|-----------|------|
| 0 | |

**stmt**

| nfree_reg | code |
|-----------|------|
| 0 | |

**program**

| name | code |
|------|------|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>stmt.code<br>goto exit |

program
(5) read
(3) (6) read
(3) (7) while
<> write
(6) (7) (6) (7) writeln
if (4) null
> := := null
(6) (7) (6) − null (7) − null
(6) (7) (7) (6)

---

## Slide 36

read:stmt$_1$ → id$_1$ id$_2$ stmt$_2$

➤ stmt$_2$.nfree_reg := stmt$_1$.nfree_reg
➤ stmt$_1$.code := ["a1 := &" id$_1$.stp→name]
  + ["call" if id$_2$.stp→type = int then "readint" else …]
  + [id$_2$.stp→name ":= rv"] + stmt$_2$.code

| Index | Symbol | Type |
|-------|---------|---------|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**stmt**

| nfree_reg | code |
|-----------|------|
| 0 | |

**stmt**

| nfree_reg | code |
|-----------|------|
| 0 | |

**program**

| name | code |
|------|------|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>stmt.code<br>goto exit |

program
(5) read
(3) (6) read
(3) (7) while
<> write
(6) (7) (6) (7) writeln
if (4) null
> := := null
(6) (7) (6) − null (7) − null
(6) (7) (7) (6)

## Slide 37

read:stmt$_1$ → id$_1$ id$_2$ stmt$_2$

➤ stmt$_2$.nfree_reg := stmt$_1$.nfree_reg
➤ stmt$_1$.code := ["a1 := &" id$_1$.stp→name]
  + ["call" if id$_2$.stp→type = int then "readint" else …]
  + [id$_2$.stp→name ":= rv"] + stmt$_2$.code

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

stmt

| nfree_reg | code |
|---|---|
| 0 | |

stmt

| nfree_reg | code |
|---|---|
| 0 | a1 := &input<br>call readint<br>j := rv<br>stmt.code |

program

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>stmt.code<br>goto exit |

---

## Slide 38

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

stmt

| nfree_reg | code |
|---|---|
| 0 | |

stmt

| nfree_reg | code |
|---|---|
| 0 | a1 := &input<br>call readint<br>j := rv<br>stmt.code |

program

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>stmt.code<br>goto exit |

## Slide 39

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | a1 := &input<br>call readint<br>j := rv<br>stmt.code |

**program**

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br><br>stmt.code<br>goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

Tree:

program
(5) read
(3) (6) read
(3) (7) while
<> write
(6) (7) (6) (7) writeln
if (4) null
> := := null
(6) (7) (6) − null (7) − null
(6) (7) (7) (6)

## Slide 40

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**program**

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>a1 := &input<br>call readint<br>j := rv<br>stmt.code<br>goto exit |

Tree:

program
(5) read
(3) (6) read
(3) (7) while
<> write
(6) (7) (6) (7) writeln
if (4) null
> := := null
(6) (7) (6) − null (7) − null
(6) (7) (7) (6)

## Slide 41

while:stmt$_1$ → expr stmt$_2$ stmt$_3$
- expr.nfree_reg := stmt$_1$.nfree_reg
- stmt$_2$.nfree_reg := stmt$_1$.nfree_reg
- stmt$_3$.nfree_reg := stmt$_1$.nfree_reg
- L1 := new_label() ; L2 := new_label();
  stmt$_1$.code := ["goto" L1] + [L2 ":"] + stmt$_2$.code + [L1 ":"] +
  expr.code + ["if" expr.reg "goto" L2] + stmt$_3$.code

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**expr**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**program**

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>a1 := &input<br>call readint<br>j := rv<br>stmt.code<br>goto exit |

---

## Slide 42

while:stmt$_1$ → expr stmt$_2$ stmt$_3$
- expr.nfree_reg := stmt$_1$.nfree_reg
- stmt$_2$.nfree_reg := stmt$_1$.nfree_reg
- stmt$_3$.nfree_reg := stmt$_1$.nfree_reg
- L1 := new_label() ; L2 := new_label();
  stmt$_1$.code := ["goto" L1] + [L2 ":"] + stmt$_2$.code + [L1 ":"] +
  expr.code + ["if" expr.reg "goto" L2] + stmt$_3$.code

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**expr**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**program**

| name | code |
|---|---|
| gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>a1 := &input<br>call readint<br>j := rv<br>stmt.code<br>goto exit |

## Slide 43

$while:stmt_1 \rightarrow expr\ stmt_2\ stmt_3$
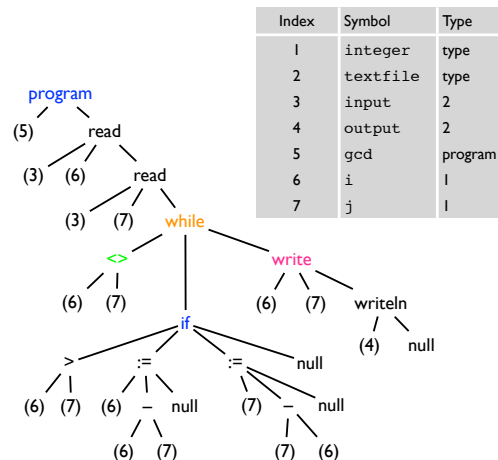
➤ expr.nfree_reg := stmt₁.nfree_reg
➤ stmt₂.nfree_reg := stmt₁.nfree_reg
➤ stmt₃.nfree_reg := stmt₁.nfree_reg
➤ L1 := new_label() ; L2 := new_label();
  stmt₁.code := ["goto" L1] + [L2 ":"] + stmt₂.code + [L1 ":"] +
  expr.code + ["if" expr.reg "goto" L2] + stmt₃.code

| | nfree_reg | code |
|---|---|---|
| **stmt** | | |
| | 0 | |

| | nfree_reg | code |
|---|---|---|
| **stmt** | | |
| | 0 | |

| | nfree_reg | code |
|---|---|---|
| **expr** | | |
| | 0 | |

| | nfree_reg | code |
|---|---|---|
| **stmt** | | |
| | 0 | |

| | name | code |
|---|---|---|
| **program** | | |
| | gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>a1 := &input<br>call readint<br>j := rv<br>stmt.code<br>goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

program
(5) read
(3) (6) read
(3) (7) while
<> write
(6) (7) (6) (7) writeln
if (4) null
> := := null
(6) (7) (6) — null (7) — null
(6) (7) (7) (6)

## Slide 44

$while:stmt_1 \rightarrow expr\ stmt_2\ stmt_3$

➤ expr.nfree_reg := stmt₁.nfree_reg
➤ stmt₂.nfree_reg := stmt₁.nfree_reg
➤ stmt₃.nfree_reg := stmt₁.nfree_reg
➤ L1 := new_label() ; L2 := new_label();
  stmt₁.code := ["goto" L1] + [L2 ":"] + stmt₂.code + [L1 ":"] +
  expr.code + ["if" expr.reg "goto" L2] + stmt₃.code

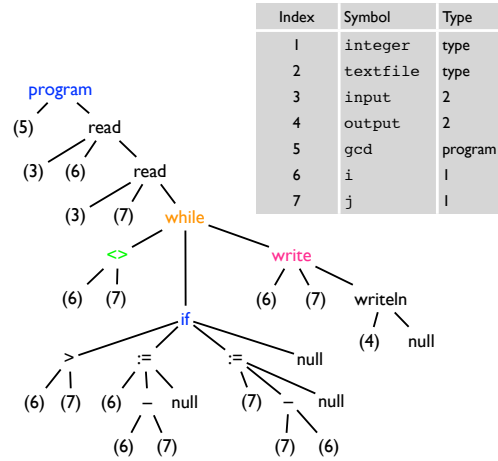| | nfree_reg | code |
|---|---|---|
| **stmt** | | |
| | 0 | |

| | nfree_reg | code |
|---|---|---|
| **stmt** | | |
| | 0 | |

| | nfree_reg | code |
|---|---|---|
| **expr** | | |
| | 0 | |

| | nfree_reg | code |
|---|---|---|
| **stmt** | | |
| | 0 | |

| | name | code |
|---|---|---|
| **program** | | |
| | gcd | main:<br>a1 := &input<br>call readint<br>i := rv<br>a1 := &input<br>call readint<br>j := rv<br>stmt.code<br>goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

program
(5) read
(3) (6) read
(3) (7) while
<> write
(6) (7) (6) (7) writeln
if (4) null
> := := null
(6) (7) (6) — null (7) — null
(6) (7) (7) (6)

## Slide 45

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**expr**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | goto x1 |
| | x2: |
| |   stmt.code |
| | x1: |
| |   expr.code |
| |   if expr.reg goto x2 |
| |   stmt.code |

**program**

| name | code |
|---|---|
| gcd | main: |
| |   a1 := &input |
| |   call readint |
| |   i := rv |
| |   a1 := &input |
| |   call readint |
| |   j := rv |
| |   stmt.code |
| |   goto exit |

while:stmt$_1$ → expr stmt$_2$ stmt$_3$

➤ expr.nfree_reg := stmt$_1$.nfree_reg
➤ stmt$_2$.nfree_reg := stmt$_1$.nfree_reg
➤ stmt$_3$.nfree_reg := stmt$_1$.nfree_reg
➤ L1 := new_label() ; L2 := new_label();
stmt$_1$.code := ["goto" L1] + [L2 ":"] + stmt$_2$.code + [L1 ":"] +
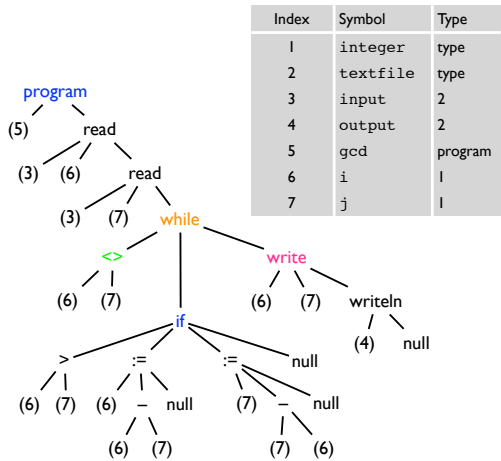expr.code + ["if" expr.reg "goto" L2] + stmt$_3$.code

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

## Slide 46

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**expr**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | goto x1 |
| | x2: |
| |   stmt.code |
| | x1: |
| |   expr.code |
| |   if expr.reg goto x2 |
| |   stmt.code |

**program**

| name | code |
|---|---|
| gcd | main: |
| |   a1 := &input |
| |   call readint |
| |   i := rv |
| |   a1 := &input |
| |   call readint |
| |   j := rv |
| |   stmt.code |
| |   goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

# Page 47

**expr**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | goto x1 |
| | x2: |
| |   stmt.code |
| | x1: |
| |   expr.code |
| |   if expr.reg goto x2 |
| |   stmt.code |

**program**

| name | code |
|---|---|
| gcd | main: |
| |   a1 := &input |
| |   call readint |
| |   i := rv |
| |   a1 := &input |
| |   call readint |
| |   j := rv |
| | |
| |   stmt.code |
| |   goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

# Page 48

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**stmt**

| nfree_reg | code |
|---|---|
| 0 | |

**expr**

| nfree_reg | code |
|---|---|
| 0 | |

**program**

| name | code |
|---|---|
| gcd | main: |
| |   a1 := &input |
| |   call readint |
| |   i := rv |
| |   a1 := &input |
| |   call readint |
| |   j := rv |
| |   goto x1 |
| | x2: |
| |   stmt.code |
| | x1: |
| |   expr.code |
| |   if expr.reg goto x2 |
| |   stmt.code |
| |   goto exit |

| Index | Symbol | Type |
|---|---|---|
| 1 | integer | type |
| 2 | textfile | type |
| 3 | input | 2 |
| 4 | output | 2 |
| 5 | gcd | program |
| 6 | i | 1 |
| 7 | j | 1 |

# Register allocation

- Rule evaluation, for each subtree
  - determine registers to be used
  - generate code
- In naïve register allocation, use a register stack
- If you run out of physical registers, spill registers into memory

# Address space organization

- Two kinds of object code
  - relocatable
    - input to linker
    - several files linked together into an executable program
  - executable
    - input to loader
    - can be loaded into memory and run

# A relocatable object file

- Includes the following
  - import table
    - identifies instructions that refer to named locations presumed to lie in other files.
  - relocation table
    - identifies instructions that refer to locations in current file, which must be modifed at link time to reflect the position of current file in the executable program
  - export table
    - lists names and addresses of locations in current file that can be referred to in other files.

# Contents of a running program

- code

- constants

- initialised data

- uninitialised data

- stack
  - small initial size, grown by the OS
- heap
  - small initial size, grown on demand
- files - mapped into memory

# Assembly

- Translates assembly code into executable code
  - replaces opcodes and operands with machine language encodings
  - replaces symbolic names with actual addresses
- Modern assemblers may perform some machine-specific code improvement
  - instruction scheduling
  - register allocation
  - peephole optimisation
    - fixes suboptimal patterns of instructions within a small window

# Constructing instructions

- Many assemblers
  - extend the instruction set in minor ways to make assembly language easier for humans to read
    - e.g. pseudo-instructions in MIPS
  - have directives for
    - segment switching, e.g. .text in MIPS
    - data generation, e.g. .byte, .half, .word, .float
    - symbol identification, e.g. .globl
    - alignment, e.g. .align

# Assigning addresses to names

- Pass 1
  - Identify internal and external (imported) symbols
    - symbol in `.globl` directive is put in export table
    - any symbol mentioned in an instruction but not defined is put in import table, with instruction offset
    - any instruction or data with a value depending on placement of file into executable program is put in the relocation table
  - assign locations to the internal symbols
- Pass 2
  - Produce object code

# Linker

- Joins together compilation units
  - static linker - prior to program execution
  - dynamic linker - during execution
- Two subtasks
  - Relocation
  - Resolution of external references

# Relocation

- Phase 1
  - gather compilation units together
  - choose an order for them in memory
  - note addresses of each unit in memory
- Phase 2
  - replace unresolved external references with addresses
  - modify instructions that need to be relocated to reflect addresses of their units