



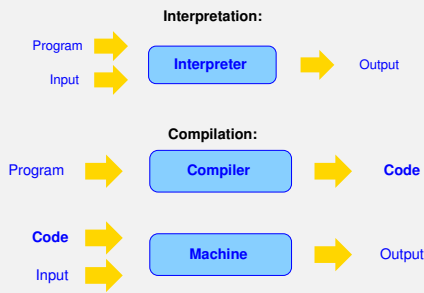
Compiler Construction I

Dr. Michael Petter

SoSe 2020

Topic: Overview

Extremes of Program Execution



Interpretation vs. Compilation

Interpretation

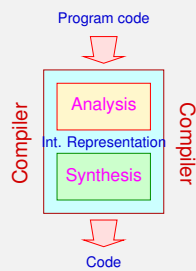
- No precomputation on program text necessary
→ no/small startup-overhead
- More context information allows for specific aggressive optimization

Compilation

- Program components are analyzed once, during preprocessing, instead of multiple times during execution
→ smaller runtime-overhead
- Runtime complexity of optimizations less important than in interpreter

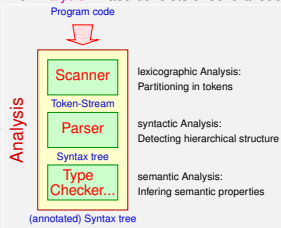
Compiler

General Compiler setup:



Compiler

The Analysis-Phase consists of several subcomponents:

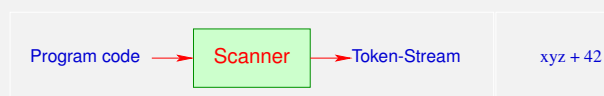


Content on the Way

- Regular expressions and finite automata
- Specification and implementation of scanners
- Context free grammars and pushdown automata
- Top-Down/Bottom-Up syntax analysis
- Attribute systems
- Typechecking
- Codegeneration for register machines

Topic: Lexical Analysis

The Lexical Analysis



- A **Token** is a sequence of characters, which together form a unit.
- Tokens are subsumed in **classes**. For example:
 - **Names (Identifiers)** e.g. `xyz, pi, ...`
 - **Constants** e.g. `42, 3.14, "abc", ...`
 - **Operators** e.g. `+, ...`
 - **Reserved terms** e.g. `if, int, ...`

The Lexical Analysis - Siever

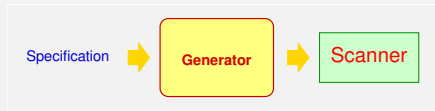
Classified tokens allow for further **pre-processing**:

- **Dropping** irrelevant fragments e.g. **Spacing, Comments,...**
- **Collecting Pragma**s, i.e. directives for the compiler, often implementation dependent, directed at the code generation process, e.g. **OpenMP**-Statements;
- **Replacing** of Tokens of particular classes with their meaning / internal representation, e.g.
 - **Constants**;
 - **Names**: typically managed centrally in a **Symbol**-table, maybe compared to reserved terms (if not already done by the scanner) and possibly replaced with an index or internal format (⇒ **Name Mangling**).

The Lexical Analysis

Discussion:

- Scanner and Siever are often combined into a single component, mostly by providing appropriate callback actions in the event that the scanner detects a token.
- Scanners are mostly not written manually, but **generated** from a specification.



4/49

The Lexical Analysis - Generating:

... in our case:



Specification of Token-classes: Regular expressions;
Generated Implementation: Finite automata + X

5/49

Lexical Analysis

Chapter 1: Basics: Regular Expressions

6/49

Regular Expressions

Basics

- Program code is composed from a finite **alphabet** Σ of input characters, e.g. Unicode
- The sets of textfragments of a token class is in general **regular**.
- Regular languages can be specified by **regular expressions**.

Definition Regular Expressions

The set \mathcal{E}_Σ of (non-empty) **regular expressions** is the smallest set \mathcal{E} with:

- $\epsilon \in \mathcal{E}$ (ϵ a new symbol not from Σ);
- $a \in \mathcal{E}$ for all $a \in \Sigma$;
- $(e_1 | e_2), (e_1 \cdot e_2), e_1^* \in \mathcal{E}$ if $e_1, e_2 \in \mathcal{E}$.



Stephen Kleene

7/49

Regular Expressions

... Example:

$((a \cdot b^*) \cdot a)$
 $(a | b)$
 $((a \cdot b) \cdot (a \cdot b))$

Attention:

- We distinguish between characters $a, 0, \$, \dots$ and **Meta-symbols** $(, |,) \dots$
- To avoid (ugly) parantheses, we make use of **Operator-Precedences**:

$* > \cdot > |$

and omit “.”

- Real Specification-languages offer additional constructs:

$e? \equiv (\epsilon | e)$
 $e^+ \equiv (e \cdot e^*)$

and omit “ ϵ ”

8/49

Regular Expressions

Specification needs **Semantics**

...Example:

Specification	Semantics
$abab$	$\{abab\}$
$a b$	$\{a, b\}$
ab^*a	$\{ab^n a \mid n \geq 0\}$

For $e \in \mathcal{E}_\Sigma$ we define the specified language $[e] \subseteq \Sigma^*$ **inductively** by:

$[e] = \{\epsilon\}$
 $[a] = \{a\}$
 $[e^*] = ([e])^*$
 $[e_1 | e_2] = [e_1] \cup [e_2]$
 $[e_1 e_2] = [e_1] \cdot [e_2]$

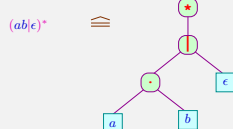
9/49

Keep in Mind:

- The operators $(_)^*, \cup, \cdot$ are interpreted in the context of sets of words:

$(L)^* = \{w_1 \dots w_k \mid k \geq 0, w_i \in L\}$
 $L_1 \cdot L_2 = \{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$

- Regular expressions are internally represented as **annotated ranked trees**:



Inner nodes: Operator-applications;
Leaves: particular symbols or ϵ .

10/49

Regular Expressions

Example: Identifiers in Java:

$le = [a-zA-Z_]\$$
 $di = [0-9]$

$Id = \{le\} \{ \{le\} | \{di\} \}^*$

$Float = \{di\}^* (\cdot \{di\} | \{di\} \cdot) \{di\}^* ((e|E) (\wedge | \wedge) ? \{di\}^+) ?$

Remarks:

- “le” and “di” are **token classes**.
- Defined Names** are enclosed in “{”, “}”.
- Symbols are distinguished from **Meta-symbols** via “\”.

11/49

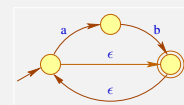
Lexical Analysis

Chapter 2: Basics: Finite Automata

12/49

Finite Automata

Example:



Nodes: States;
Edges: Transitions;
Labels: Consumed input;

13/49

Finite Automata

Definition Finite Automata

A **non-deterministic** finite automaton (NFA) is a tuple $A = (Q, \Sigma, \delta, I, F)$ with:

- Q a finite set of states;
- Σ a finite alphabet of inputs;
- $I \subseteq Q$ the set of start states;
- $F \subseteq Q$ the set of final states and
- δ the set of transitions (-relation)



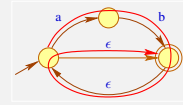
For an NFA, we reckon:

Definition Deterministic Finite Automata

Given $\delta : Q \times \Sigma \rightarrow Q$ a function and $|I| = 1$, then we call the NFA A **deterministic** (DFA).

Finite Automata

- Computations are paths in the graph.
- Accepting computations lead from I to F .
- An **accepted word** is the sequence of labels along an accepting computation ...



Finite Automata

Once again, more formally:

- We define the **transitive closure** δ^* of δ as the smallest set δ' with:

$$(p, \epsilon, p) \in \delta' \text{ and } (p, xw, q) \in \delta' \text{ if } (p, x, p_1) \in \delta \text{ and } (p_1, w, q) \in \delta'$$

δ^* characterizes for a path between the states p and q the words obtained by concatenating the labels along it.

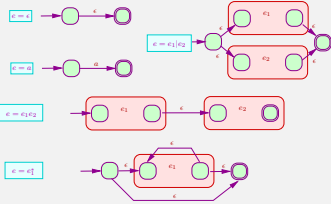
- The set of all accepting words, i.e. A 's **accepted language** can be described compactly as:

$$\mathcal{L}(A) = \{w \in \Sigma^* \mid \exists i \in I, f \in F : (i, w, f) \in \delta^*\}$$

Lexical Analysis

Chapter 3: Converting Regular Expressions to NFAs

In Linear Time from Regular Expressions to NFAs



Thompson's Algorithm

Produces $\mathcal{O}(n)$ states for regular expressions of length n .



A formal approach to Thompson's Algorithm

Berry-Sethi Algorithm / Glushkov Automaton

Produces exactly $n + 1$ states without ϵ -transitions and demonstrates \rightarrow Equality Systems and \rightarrow Attribute Grammars



Idea:

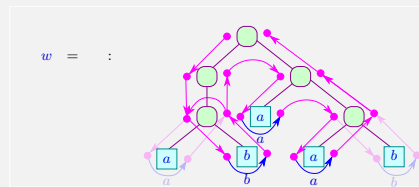
An automaton covering the syntax tree of a regular expression e tracks (conceptually via markers \bullet), which subexpressions e' are reachable consuming the rest of input w .

- markers contribute an entry or exit point into the automaton for this subexpression
- edges for each layer of subexpression are modelled after Thompson's automata



Berry-Sethi Approach

... for example:



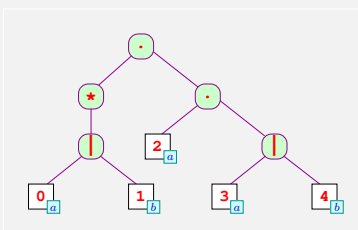
Berry-Sethi Approach

In general:

- Input is only consumed at the leaves.
- Navigating the tree does not consume input \rightarrow ϵ -transitions
- For a formal construction we need **identifiers** for states.
- For a node n 's **identifier** we take the **subexpression**, corresponding to the subtree dominated by n .
- There are possibly **identical subexpressions** in one regular expression. \implies we enumerate the leaves ...

Berry-Sethi Approach

... for example:



Berry-Sethi Approach (naive version)

Construction (naive version):

- States: $\bullet r, r \bullet$ with r nodes of e ;
- Start state: $\bullet e$;
- Final state: $e \bullet$;
- Transitions: for leaves $r \equiv [a|b]$ we require: $(\bullet r, x, r \bullet)$.
- The leftover transitions are:

r	Transitions
$r_1 \mid r_2$	$(\bullet r_1, \epsilon, \bullet r_1)$ $(\bullet r_1, \epsilon, \bullet r_2)$ $(r_1 \bullet, \epsilon, r_1)$ $(r_2 \bullet, \epsilon, r_2)$
$r_1 \cdot r_2$	$(\bullet r_1, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, \bullet r_2)$ $(r_2 \bullet, \epsilon, r_2)$

r	Transitions
r_1	$(\bullet r_1, \epsilon, r_1 \bullet)$ $(\bullet r_1, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, r_1 \bullet)$
$r_1 ?$	$(\bullet r_1, \epsilon, r_1 \bullet)$ $(\bullet r_1, \epsilon, \bullet r_1)$ $(r_1 \bullet, \epsilon, r_1 \bullet)$

Berry-Sethi Approach

Discussion:

- Most transitions navigate through the expression
- The resulting automaton is in general **nondeterministic**

⇒ Strategy for the sophisticated version:
Avoid generating ϵ -transitions

Idea:

Pre-compute helper attributes during **D**(epth)**F**(irst)**S**(earch)!

Necessary node-attributes:

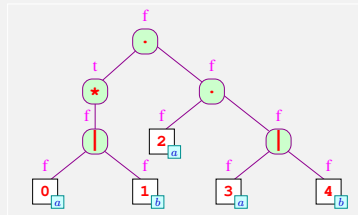
- first** the set of read states below r , which may be reached **first**, when descending into r .
- next** the set of read states, which may be reached **first** in the traversal **after** r .
- last** the set of read states below r , which may be reached **last** when descending into r .
- empty** can the subexpression r consume ϵ ?

24/49

Berry-Sethi Approach: 1st step

$empty[r] = t$ if and only if $\epsilon \in [r]$

... for example:



25/49

Berry-Sethi Approach: 1st step

Implementation:

DFS **post-order** traversal

for leaves $r \equiv [x]$ we find $empty[r] = (x \equiv \epsilon)$.

Otherwise:

$$\begin{aligned} empty[r_1 | r_2] &= empty[r_1] \vee empty[r_2] \\ empty[r_1 \cdot r_2] &= empty[r_1] \wedge empty[r_2] \\ empty[r_1^*] &= t \\ empty[r_1^?] &= t \end{aligned}$$

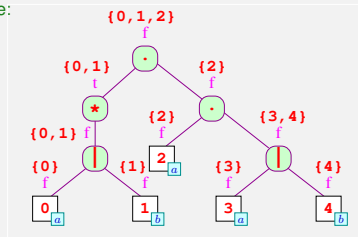
26/49

Berry-Sethi Approach: 2nd step

The **may-set of first reached read states**: The set of read states, that may be reached from r (i.e. while descending into r) via sequences of ϵ -transitions:

$first[r] = \{i \text{ in } r \mid (\epsilon^* r, \epsilon, [x]) \in \delta^*, x \neq \epsilon\}$

... for example:



27/49

Berry-Sethi Approach: 2nd step

Implementation:

DFS **post-order** traversal

for leaves $r \equiv [x]$ we find $first[r] = \{i \mid x \neq \epsilon\}$.

Otherwise:

$$\begin{aligned} first[r_1 | r_2] &= first[r_1] \cup first[r_2] \\ first[r_1 \cdot r_2] &= \begin{cases} first[r_1] \cup first[r_2] & \text{if } empty[r_1] = t \\ first[r_1] & \text{if } empty[r_1] = f \end{cases} \\ first[r_1^*] &= first[r_1] \\ first[r_1^?] &= first[r_1] \end{aligned}$$

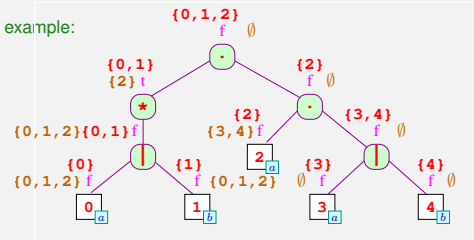
28/49

Berry-Sethi Approach: 3rd step

The **may-set of next read states**: The set of read states reached after reading r , that may be reached next via sequences of ϵ -transitions.

$next[r] = \{i \mid (r, \epsilon, [x]) \in \delta^*, x \neq \epsilon\}$

... for example:



29/49

Berry-Sethi Approach: 3rd step

Implementation:

DFS **pre-order** traversal

For the root, we find: $next[\epsilon] = \emptyset$

Apart from that we distinguish, based on the context:

r	Equalities
$r_1 r_2$	$next[r_1] = next[r]$ $next[r_2] = next[r]$
$r_1 \cdot r_2$	$next[r_1] = \begin{cases} first[r_2] \cup next[r] & \text{if } empty[r_2] = t \\ first[r_2] & \text{if } empty[r_2] = f \end{cases}$ $next[r_2] = next[r]$
r_1^*	$next[r_1] = first[r_1] \cup next[r]$
$r_1^?$	$next[r_1] = next[r]$

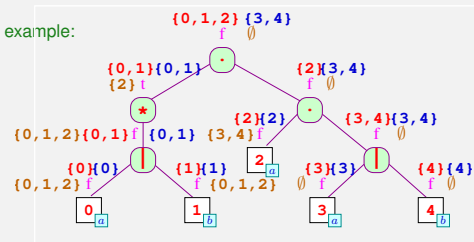
30/49

Berry-Sethi Approach: 4th step

The **may-set of last reached read states**: The set of read states, which may be reached last during the traversal of r connected to the root via ϵ -transitions only:

$last[r] = \{i \text{ in } r \mid ([x], \epsilon, r) \in \delta^*, x \neq \epsilon\}$

... for example:



31/49

Berry-Sethi Approach: 4th step

Implementation:

DFS **post-order** traversal

for leaves $r \equiv [x]$ we find $last[r] = \{i \mid x \neq \epsilon\}$.

Otherwise:

$$\begin{aligned} last[r_1 | r_2] &= last[r_1] \cup last[r_2] \\ last[r_1 \cdot r_2] &= \begin{cases} last[r_1] \cup last[r_2] & \text{if } empty[r_2] = t \\ last[r_2] & \text{if } empty[r_2] = f \end{cases} \\ last[r_1^*] &= last[r_1] \\ last[r_1^?] &= last[r_1] \end{aligned}$$

32/49

Berry-Sethi Approach: (sophisticated version)

Construction (sophisticated version):

Create an automaton based on the syntax tree's new attributes:

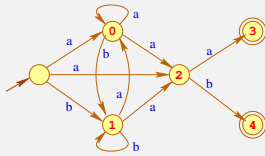
- States: $\{\bullet e\} \cup \{i \bullet \mid i \text{ a leaf not } \epsilon\}$
- Start state: $\bullet e$
- Final states: $last[e]$ if $empty[e] = f$
 $\{\bullet e\} \cup last[e]$ otherwise
- Transitions: $(\bullet e, a, i \bullet)$ if $i \in first[e]$ and i labeled with a .
 $(i \bullet, a, i' \bullet)$ if $i' \in next[i]$ and i' labeled with a .

We call the resulting automaton A_e .

33/49

Berry-Sethi Approach

... for example:



Remarks:

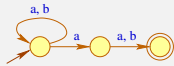
- This construction is known as **Berry-Sethi-** or **Glushkov-construction**.
- It is used for **NFA** to define **Content Models**
- The result may not be, what we had in mind...

34/49

Chapter 4: Turning NFAs deterministic

35/49

The expected outcome:



Remarks:

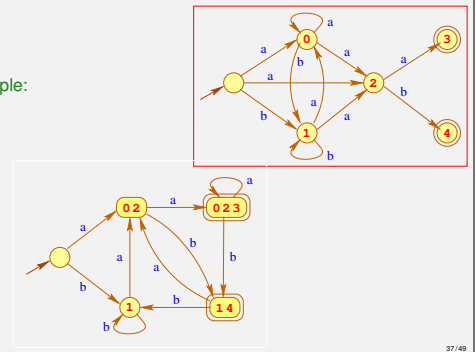
- ideal automaton would be even more compact
(→ *Antimirov automata, Follow Automata*)
- but Berry-Sethi is rather directly constructed
- Anyway, we need a **deterministic** version

⇒ **Powerset-Construction**

36/49

Powerset Construction

... for example:



37/49

Powerset Construction

Theorem:

For every non-deterministic automaton $A = (Q, \Sigma, \delta, I, F)$ we can compute a deterministic automaton $\mathcal{P}(A)$ with

$$\mathcal{L}(A) = \mathcal{L}(\mathcal{P}(A))$$

Construction:

- States:** Powersets of Q ;
- Start state:** I ;
- Final states:** $\{Q' \subseteq Q \mid Q' \cap F \neq \emptyset\}$;
- Transitions:** $\delta_{\mathcal{P}}(Q', a) = \{q \in Q \mid \exists p \in Q' : (p, a, q) \in \delta\}$.

38/49

Powerset Construction

Observation:

There are exponentially many powersets of Q

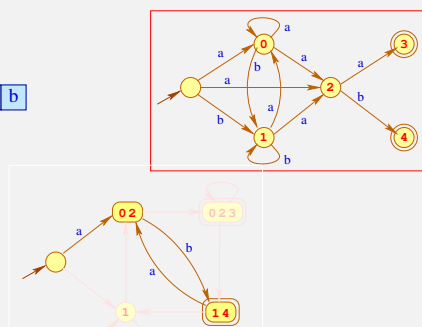
- **Idea:** Consider only **contributing** powersets. Starting with the set $Q_P = \{I\}$ we only add further states **by need** ...
- i.e., whenever we can reach them from a state in Q_P
- However, the resulting automaton can become **enormously huge** ... which is (sort of) not happening in **practice**
- Therefore, in tools like **grep** a regular expression's **DFA** is never created!
- Instead, only the sets, directly necessary for interpreting the input are generated **while processing the input**

39/49

Powerset Construction

... for example:

a b a b



40/49

Remarks:

- For an input sequence of length n , maximally $\mathcal{O}(n)$ sets are generated
- Once a set/edge of the **DFA** is generated, they are stored within a **hash-table**.
- Before generating a new transition, we check this table for already existing edges with the desired label.

Summary:

Theorem:

For each regular expression e we can compute a deterministic automaton $A = \mathcal{P}(A_e)$ with

$$\mathcal{L}(A) = [e]$$

41/49

Chapter 5: Scanner design

42/49

Scanner design

Input (simplified): a set of rules:

$$\begin{aligned} e_1 & \{ \text{action}_1 \} \\ e_2 & \{ \text{action}_2 \} \\ & \dots \\ e_k & \{ \text{action}_k \} \end{aligned}$$

Output: a program,

- ... reading a **maximal prefix** w from the input, that satisfies $e_1 \mid \dots \mid e_k$;
- ... determining the **minimal** i , such that $w \in [e_i]$;
- ... executing **action_i** for w .

43/49

Implementation:

Idea:

- Create the NFA $\mathcal{P}(A_e) = (Q, \Sigma, \delta, q_0, F)$ for the expression $e = (e_1 | \dots | e_k)$;
- Define the sets:

$$F_1 = \{q \in F \mid q \cap \text{last}[e_1] \neq \emptyset\}$$

$$F_2 = \{q \in (F \setminus F_1) \mid q \cap \text{last}[e_2] \neq \emptyset\}$$

$$\dots$$

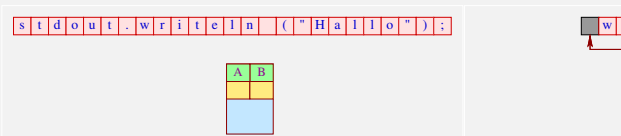
$$F_k = \{q \in (F \setminus (F_1 \cup \dots \cup F_{k-1})) \mid q \cap \text{last}[e_k] \neq \emptyset\}$$
- For input w we find: $\delta^*(q_0, w) \in F_i$ iff the scanner must execute **action_i** for w

44/49

Implementation:

Idea (cont'd):

- The scanner manages two pointers $\langle A, B \rangle$ and the related states $\langle q_A, q_B \rangle \dots$
- Pointer A points to the last position in the input, after which a state $q_A \in F$ was reached;
- Pointer B tracks the current position.



45/49

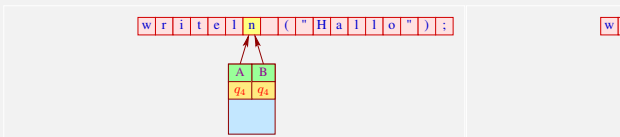
Implementation:

Idea (cont'd):

- The current state being $q_B = \emptyset$, we consume input up to position A and reset:

$$B := A; \quad A := \perp;$$

$$q_B := q_0; \quad q_A := \perp$$



46/49

Extension: States

- Now and then, it is handy to differentiate between particular **scanner states**.
- In different states, we want to recognize different token classes with different precedences.
- Depending on the consumed input, the scanner state can be changed

Example: Comments

Within a comment, identifiers, constants, comments, ... are ignored

47/49

Input (generalized): a set of rules:

```

(state) {
  e1 { action1 yybegin(state1); }
  e2 { action2 yybegin(state2); }
  ...
  ek { actionk yybegin(statek); }
}
  
```

- The statement `yybegin (statei);` resets the current state to `statei`.
- The start state is called (e.g. `flex JFlex`) `YYINITIAL`.

... for example:

```

(YYINITIAL) { "/*" { yybegin(COMMENT); }
(COMMENT) { "*/" { yybegin(YYINITIAL); }
            . | \n { }
            }
  
```

48/49

Remarks:

- "." matches all characters different from "\n".
- For every state we generate the scanner respectively.
- Method `yybegin (STATE);` switches between different scanners.
- Comments might be directly implemented as (admittedly overly complex) token-class.
- Scanner-states are especially handy for implementing **preprocessors**, expanding special fragments in regular programs.


49/49

Topic:

Syntactic Analysis

1/66

Syntactic Analysis

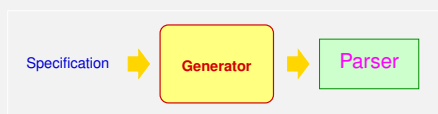


- Syntactic analysis tries to integrate Tokens into larger program units.
- Such units may possibly be:
 - Expressions;
 - Statements;
 - Conditional branches;
 - loops; ...

2/66

Discussion:

In general, parsers are not developed by hand, but **generated** from a specification:

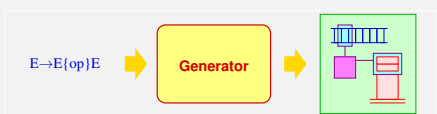


Specification of the hierarchical structure: contextfree grammars
Generated implementation: Pushdown automata + X

3/66

Discussion:

In general, parsers are not developed by hand, but **generated** from a specification:



Specification of the hierarchical structure: contextfree grammars
Generated implementation: Pushdown automata + X

3/66

Chapter 1: Basics of Contextfree Grammars

Basics: Context-free Grammars

- Programs of programming languages can have arbitrary numbers of tokens, but only finitely many **Token-classes**.
- This is why we choose the set of **Token-classes** to be the finite alphabet of terminals T .
- The nested structure of program components can be described elegantly via **context-free grammars**...

Definition: Context-Free Grammar

A context-free grammar (CFG) is a 4-tuple $G = (N, T, P, S)$ with:

- N the set of nonterminals,
- T the set of terminals,
- P the set of productions or rules, and
- $S \in N$ the start symbol



Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \text{ with } A \in N, \alpha \in (N \cup T)^*$$

... for example:

$$\begin{aligned} S &\rightarrow aSb \\ S &\rightarrow \epsilon \end{aligned}$$

Specified language: $\{a^n b^n \mid n \geq 0\}$

Conventions:

In examples, we specify nonterminals and terminals in general implicitly:

- nonterminals are: $A, B, C, \dots, \langle \text{exp} \rangle, \langle \text{stmt} \rangle, \dots;$
- terminals are: $a, b, c, \dots, \text{int}, \text{name}, \dots;$

... a practical example:

$$\begin{aligned} S &\rightarrow \langle \text{stmt} \rangle \\ \langle \text{stmt} \rangle &\rightarrow \langle \text{if} \rangle \mid \langle \text{while} \rangle \mid \langle \text{rexp} \rangle; \\ \langle \text{if} \rangle &\rightarrow \text{if} (\langle \text{rexp} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ \langle \text{while} \rangle &\rightarrow \text{while} (\langle \text{rexp} \rangle) \langle \text{stmt} \rangle \\ \langle \text{rexp} \rangle &\rightarrow \text{int} \mid \langle \text{lexp} \rangle \mid \langle \text{lexp} \rangle = \langle \text{rexp} \rangle \mid \dots \\ \langle \text{lexp} \rangle &\rightarrow \text{name} \mid \dots \end{aligned}$$

More conventions:

- For every nonterminal, we collect the right hand sides of rules and list them together.
- The j -th rule for A can be identified via the pair (A, j) (with $j \geq 0$).

Pair of grammars:

$E \rightarrow E+E$	$E * E$	(E)	name	int
$E \rightarrow E+T$	T			
$T \rightarrow T * F$	F			
$F \rightarrow (E)$	name	int		

$E \rightarrow E+E^0$	$E * E^1$	$(E)^2$	name ³	int ⁴
$E \rightarrow E+T^0$	T^1			
$T \rightarrow T * F^0$	F^1			
$F \rightarrow (E)^0$	name ¹	int ²		

Both grammars describe the same language

Derivation

Grammars are **term rewriting systems**. The rules offer feasible rewriting steps. A sequence of such rewriting steps $\alpha_0 \rightarrow \dots \rightarrow \alpha_m$ is called **derivation**.

... for example:

$$\begin{aligned} E &\rightarrow E + T \\ &\rightarrow T + T \\ &\rightarrow T * E + T \\ &\rightarrow T * \text{int} + T \\ &\rightarrow E * \text{int} + T \\ &\rightarrow \text{name} * \text{int} + T \\ &\rightarrow \text{name} * \text{int} + F \\ &\rightarrow \text{name} * \text{int} + \text{int} \end{aligned}$$

Definition

The rewriting relation \rightarrow is a relation on words over $N \cup T$, with

$$\alpha \rightarrow \alpha' \text{ iff } \alpha = \alpha_1 A \alpha_2 \wedge \alpha' = \alpha_1 \beta \alpha_2 \text{ for an } A \rightarrow \beta \in P$$

The **reflexive** and **transitive** closure of \rightarrow is denoted as: \rightarrow^*

Derivation

Remarks:

- The relation \rightarrow depends on the grammar
- In each step of a derivation, we may choose:
 - * a spot, determining **where** we will rewrite.
 - * a rule, determining **how** we will rewrite.
- The language, specified by G is:

$$\mathcal{L}(G) = \{w \in T^* \mid S \rightarrow^* w\}$$

Attention:

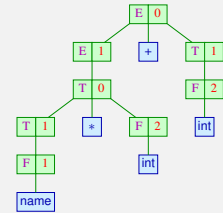
The order, in which disjunct fragments are rewritten is not relevant.

Derivation Tree

Derivations of a symbol are represented as **derivation trees**:

... for example:

$$\begin{aligned} E &\xrightarrow{0} E + T \\ &\xrightarrow{1} T + T \\ &\xrightarrow{0} T * E + T \\ &\xrightarrow{2} T * \text{int} + T \\ &\xrightarrow{1} E * \text{int} + T \\ &\xrightarrow{1} \text{name} * \text{int} + T \\ &\xrightarrow{1} \text{name} * \text{int} + F \\ &\xrightarrow{2} \text{name} * \text{int} + \text{int} \end{aligned}$$



A derivation tree for $A \in N$:
inner nodes: rule applications
root: rule application for A
leaves: terminals or ϵ

The successors of (B, i) correspond to right hand sides of the rule

Special Derivations

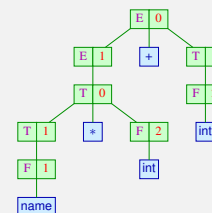
Attention:

In contrast to arbitrary derivations, we find special ones, always rewriting the **leftmost** (or rather **rightmost**) occurrence of a nonterminal.

- These are called **leftmost** (or rather **rightmost**) derivations and are denoted with the index L (or R respectively).
- Leftmost (or rightmost) derivations correspond to a left-to-right (or right-to-left) **preorder**-DFS-traversal of the derivation tree.
- **Reverse** rightmost derivations correspond to a left-to-right **postorder**-DFS-traversal of the derivation tree

Special Derivations

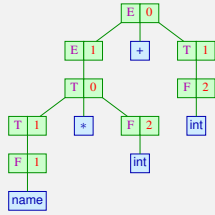
... for example:



Leftmost derivation: $(E, 0) (E, 1) (T, 0) (T, 1) (F, 1) (F, 2) (T, 1) (F, 2)$
Rightmost derivation: $(E, 0) (T, 1) (F, 2) (E, 1) (T, 0) (F, 2) (T, 1) (F, 1)$
Reverse rightmost derivation: $(F, 1) (T, 1) (F, 2) (T, 0) (E, 1) (F, 2) (T, 1) (E, 0)$

Unique Grammars

The concatenation of leaves of a derivation tree t are often called $yield(t)$.
... for example:



gives rise to the concatenation: $name * int + int$.

14/66

Unique Grammars

Definition:

Grammar G is called **unique**, if for every $w \in T^*$ there is maximally one derivation tree t of S with $yield(t) = w$.

... in our example:

$E \rightarrow E+E^0$	$E * E^1$	int^4
$(E)^2$	$name^3$	
$E \rightarrow E+T^0$	T^1	
$T \rightarrow T * F^0$	F^1	
$F \rightarrow (E)^0$	$name^1$	int^2

The first one is ambiguous, the second one is unique

15/66

Conclusion:

- A derivation tree represents a possible hierarchical structure of a word.
- For programming languages, only those grammars with a unique structure are of interest.
- Derivation trees are one-to-one corresponding with leftmost derivations as well as (reverse) rightmost derivations.
- **Leftmost derivations** correspond to a **top-down** reconstruction of the syntax tree.
- **Reverse rightmost derivations** correspond to a **bottom-up** reconstruction of the syntax tree.

16/66

Finger Exercise: Redundant Nonterminals and Rules

Definition:

$A \in N$ is **productive**, if $A \rightarrow^* w$ for a $w \in T^*$

$A \in N$ is **reachable**, if $S \rightarrow^* \alpha A \beta$ for suitable $\alpha, \beta \in (T \cup N)^*$

Example:

$S \rightarrow aBB \mid bD$
 $A \rightarrow Bc$
 $B \rightarrow Sd \mid C$
 $C \rightarrow a$
 $D \rightarrow BD$

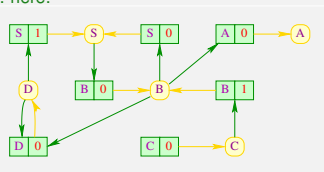
Productive nonterminals: S, A, B, C
Reachable nonterminals: S, B, C, D

17/66

Productive Nonterminals

Idea for Productivity: And-Or-Graph for a Grammar

... here:



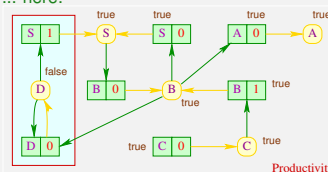
And-nodes: Rules
Or-nodes: Nonterminals
Edges: $((B, i), B)$ for all rules (B, i)
 $(A, (B, i))$ if $(B, i) \equiv B \rightarrow \alpha_1 A \alpha_2$

18/66

Productive Nonterminals

Idea for Productivity: And-Or-Graph for a Grammar

... here:



And-nodes: Rules
Or-nodes: Nonterminals
Edges: $((B, i), B)$ for all rules (B, i)
 $(A, (B, i))$ if $(B, i) \equiv B \rightarrow \alpha_1 A \alpha_2$

18/66

Productive Nonterminals - Algorithm:

```

2^N result = 0; // Result-set
int count[P]; // Rule counter
2^P rhs[N]; // Occurrences in right hand sides

forall (A in N) rhs[A] = 0; // Initialization
forall ((A, i) in P) {
    count[(A, i)] = 0; // Initialization of rhs
    init(A, i);
}

```

Helper function **init** counts the nonterminal-occurrences in right hand sides and protocols them in data structure **rhs**

19/66

Productive Nonterminals - Algorithm (cont.):

```

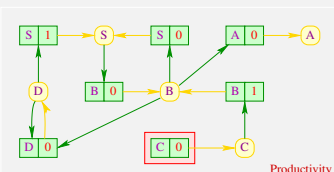
...
2^P W = {r | count[r] = 0}; // Workset
while (W != 0) {
    (A, i) = extract(W);
    if (A not result) {
        result = result union {A};
        forall (r in rhs[A]) {
            count[r]--;
            if (count[r] == 0) W = W union {r};
        }
    }
}

```

Set W contains the rules, whose right hand sides only contain productive nonterminals

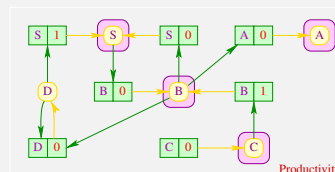
20/66

Productive Nonterminals - in an Example



21/66

Productive Nonterminals - in an Example



21/66

Runtime:

- Initialization of data structures is linear.
- Each rules is added once to W' at most.
- Each A is added once to $result$ at most.
 \implies Runtime is **linear** in the size of the grammar

Correctness:

- If A is added to $result$ in the j -th iteration of the **while**-loop there is a derivation tree for A of height maximally $j - 1$.
- For every derivation tree the root is added once to W'

22/66

Discussion:

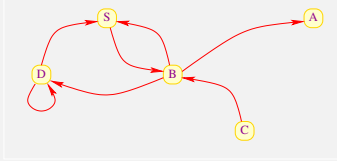
- To simplify the test ($A \in result$), we represent the set $result$ as an **array**.
- W' as well as the sets $rhs[A]$ are represented as **Lists**
- The algorithm also works for finding **smallest** solutions for **Boolean** inequality systems
- $\mathcal{L}(G) \neq \emptyset$ (\rightarrow **Emptiness Problem**) can be reduced to determining productive nonterminals

23/66

Reachable Nonterminals

Idea for Reachability: **Dependency-Graph**

... here:



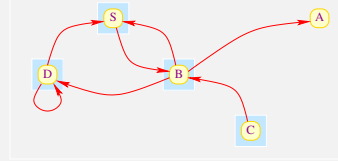
Nodes: Nonterminals
Edges: (A, B) if $B \rightarrow \alpha_1 A \alpha_2 \in P$

24/66

Reachable Nonterminals

Idea for Reachability: **Dependency-Graph**

... here:



Nonterminal A is reachable, if there is a path A to S in the dependency graph

24/66

Reduced Grammars

Conclusion:

- Reachability in directed graphs can be computed via **DFS** in **linear time**.
- This means the set of all reachable and productive nonterminals can be computed in **linear time**.

A Grammar G is called **reduced**, if all of G 's nonterminals are productive and reachable as well...

Theorem:

Each contextfree Grammar $G = (N, T, P, S)$ with $\mathcal{L}(G) \neq \emptyset$ can be converted in **linear time** into a reduced Grammar G' with

$$\mathcal{L}(G) = \mathcal{L}(G')$$

25/66

Reduced Grammars - Construction:

1. Step:

Compute the subset $N_1 \subseteq N$ of all productive nonterminals of G .
 Since $\mathcal{L}(G) \neq \emptyset$ in particular $S \in N_1$.

2. Step:

Construct: $P_1 = \{A \rightarrow \alpha \in P \mid A \in N_1 \wedge \alpha \in (N_1 \cup T)^*\}$

3. Step:

Compute the subset $N_2 \subseteq N_1$ of all productive **and** reachable nonterminals of G .
 Since $\mathcal{L}(G) \neq \emptyset$ in particular $S \in N_2$.

4. Step:

Construct: $P_2 = \{A \rightarrow \alpha \in P \mid A \in N_2 \wedge \alpha \in (N_2 \cup T)^*\}$

Result: $G' = (N_2, T, P_2, S)$

26/66

Reduced Grammars - Example:

$$\begin{aligned} S &\rightarrow aBB \mid bD \\ A &\rightarrow Bc \\ B &\rightarrow Sd \mid C \\ C &\rightarrow a \\ D &\rightarrow BD \end{aligned}$$

27/66

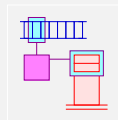
Syntactic Analysis

Chapter 2: Basics of Pushdown Automata

28/66

Basics of Pushdown Automata

Languages, specified by context free grammars are accepted by **Pushdown Automata**:



The pushdown is used e.g. to verify correct nesting of braces.

29/66

Example:

States: 0, 1, 2
Start state: 0
Final states: 0, 2

0	a	11
1	a	11
11	b	2
12	b	2

Conventions:

- We do **not** differentiate between pushdown symbols and states
- The rightmost / upper pushdown symbol represents the state
- Every transition consumes / modifies the upper part of the pushdown

30/66

Definition: Pushdown Automaton

A pushdown automaton (PDA) is a tuple $M = (Q, T, \delta, q_0, F)$ with:

- Q a finite set of states;
- T an input alphabet;
- $q_0 \in Q$ the start state;
- $F \subseteq Q$ the set of final states and
- $\delta \subseteq Q^+ \times (T \cup \{\epsilon\}) \times Q^*$ a finite set of transitions



We define **computations** of pushdown automata with the help of transitions; a particular **computation state** (the current **configuration**) is a pair:

$$(\gamma, w) \in Q^* \times T^*$$

consisting of the **pushdown content** and the **remaining input**.

... for example:

States: 0, 1, 2
Start state: 0
Final states: 0, 2

0	a	11
1	a	11
11	b	2
12	b	2

$(0, aabbb) \vdash (11, aabbb)$
 $\vdash (111, abbb)$
 $\vdash (1111, bbb)$
 $\vdash (112, bb)$
 $\vdash (12, b)$
 $\vdash (2, \epsilon)$

A computation step is characterized by the relation $\vdash \subseteq (Q^* \times T^*)^2$ with

$$(\alpha\gamma, xw) \vdash (\alpha\gamma', w) \text{ for } (\gamma, x, \gamma') \in \delta$$

Remarks:

- The relation \vdash depends on the pushdown automaton M
- The reflexive and transitive closure of \vdash is denoted by \vdash^*
- Then, the language accepted by M is

$$\mathcal{L}(M) = \{w \in T^* \mid \exists f \in F : (q_0, w) \vdash^* (f, \epsilon)\}$$

We accept with a **final state** together with **empty input**.

Definition: Deterministic Pushdown Automaton

The pushdown automaton M is **deterministic**, if every configuration has maximally one successor configuration.

This is exactly the case if for distinct transitions $(\gamma_1, x, \gamma_2), (\gamma'_1, x', \gamma'_2) \in \delta$ we can assume: Is γ_1 a suffix of γ'_1 , then $x \neq x' \wedge x \neq \epsilon \neq x'$ is valid.

... for example:

0	a	11
1	a	11
11	b	2
12	b	2

... this obviously holds

Pushdown Automata

Theorem:

For each context free grammar $G = (N, T, P, S)$ a pushdown automaton M with $\mathcal{L}(G) = \mathcal{L}(M)$ can be built.



The theorem is so important for us, that we take a look at **two** constructions for automata, motivated by both of the special derivations:

- M_G^L to build **Leftmost derivations**
- M_G^R to build **reverse Rightmost derivations**

Syntactic Analysis

**Chapter 3:
Top-down Parsing**

Item Pushdown Automaton

Construction: Item Pushdown Automaton M_G^L

- Reconstruct a **Leftmost derivation**.
- Expand nonterminals using a rule.
- Verify successively, that the chosen rule matches the input.

⇒ The states are now **Items** (= rules with a bullet):

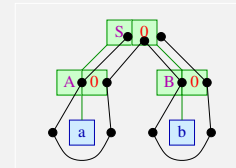
$$[A \rightarrow \alpha \bullet \beta], \quad A \rightarrow \alpha \beta \in P$$

The bullet marks the spot, how far the rule is already processed

Item Pushdown Automaton – Example

Our example:

$$S \rightarrow AB^0 \quad A \rightarrow a^0 \quad B \rightarrow b^0$$



Item Pushdown Automaton – Example

We add another rule $S' \rightarrow S \S$ for initialising the construction:

Start state: $[S' \rightarrow \bullet S \S]$
End state: $[S' \rightarrow S \bullet \S]$
Transition relations:

$[S' \rightarrow \bullet S \S]$	ϵ	$[S' \rightarrow \bullet S \S]$	$[S \rightarrow \bullet AB]$
$[S \rightarrow \bullet AB]$	ϵ	$[S \rightarrow \bullet AB]$	$[A \rightarrow \bullet a]$
$[A \rightarrow \bullet a]$	a	$[A \rightarrow a \bullet]$	
$[S \rightarrow \bullet AB]$	ϵ	$[S \rightarrow A \bullet B]$	
$[S \rightarrow A \bullet B]$	ϵ	$[S \rightarrow A \bullet B]$	$[B \rightarrow \bullet b]$
$[B \rightarrow \bullet b]$	b	$[B \rightarrow b \bullet]$	
$[S \rightarrow A \bullet B]$	ϵ	$[S \rightarrow AB \bullet]$	
$[S' \rightarrow \bullet S \S]$	ϵ	$[S' \rightarrow S \bullet \S]$	

Item Pushdown Automaton

The item pushdown automaton M_G^L has three kinds of transitions:

- Expansions:** $([A \rightarrow \alpha \bullet B \beta], \epsilon, [A \rightarrow \alpha \bullet B \beta] [B \rightarrow \bullet \gamma])$ for $A \rightarrow \alpha B \beta, B \rightarrow \gamma \in P$
- Shifts:** $([A \rightarrow \alpha \bullet a \beta], a, [A \rightarrow \alpha a \bullet \beta])$ for $A \rightarrow \alpha a \beta \in P$
- Reduces:** $([A \rightarrow \alpha \bullet B \beta] [B \rightarrow \gamma \bullet], \epsilon, [A \rightarrow \alpha B \bullet \beta])$ for $A \rightarrow \alpha B \beta, B \rightarrow \gamma \in P$

Items of the form: $[A \rightarrow \alpha \bullet]$ are also called **complete**
 The item pushdown automaton shifts the bullet around the derivation tree ...

Item Pushdown Automaton

Discussion:

- The **expansions** of a computation form a **leftmost derivation**
- Unfortunately, the expansions are chosen **nondeterministically**
- For proving correctness of the construction, we show that for every Item $[A \rightarrow \alpha \bullet B \beta]$ the following holds:

$$([A \rightarrow \alpha \bullet B \beta], w) \vdash^* ([A \rightarrow \alpha B \bullet \beta], \epsilon) \quad \text{iff} \quad B \rightarrow^* w$$
- LL-Parsing** is based on the item pushdown automaton and tries to make the expansions deterministic ...

41/66

Item Pushdown Automaton

Example: $S' \rightarrow S \$ \quad S \rightarrow \epsilon \mid a S b$

The transitions of the according Item Pushdown Automaton:

0	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow \bullet S \$] [S \rightarrow \bullet]$
1	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow \bullet S \$] [S \rightarrow \bullet a S b]$
2	$[S \rightarrow \bullet a S b]$	a	$[S \rightarrow a \bullet S b]$
3	$[S \rightarrow \bullet a S b]$	ϵ	$[S \rightarrow a \bullet S b] [S \rightarrow \bullet]$
4	$[S \rightarrow \bullet a S b]$	ϵ	$[S \rightarrow a \bullet S b] [S \rightarrow \bullet a S b]$
5	$[S \rightarrow \bullet a S b] [S \rightarrow \bullet]$	ϵ	$[S \rightarrow a S \bullet b]$
6	$[S \rightarrow \bullet a S b] [S \rightarrow a S \bullet b]$	ϵ	$[S \rightarrow a S \bullet b]$
7	$[S \rightarrow \bullet a S b] [S \rightarrow a S \bullet b]$	b	$[S \rightarrow a S b \bullet]$
8	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow S \bullet \$]$
9	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow S \bullet \$]$

Conflicts arise between the transitions (0, 1) and (3, 4), resp.

42/66

Topdown Parsing

Problem:

Conflicts between the transitions prohibit an implementation of the item pushdown automaton as deterministic pushdown automaton.

Idea 1: GLL Parsing

For each conflict, we create a virtual copy of the complete configuration and continue computing in parallel.

Idea 2: Recursive Descent & Backtracking

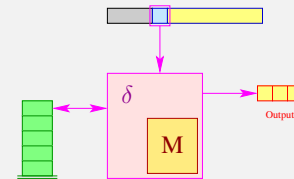
Depth-first search for an appropriate derivation.

Idea 3: Recursive Descent & Lookahead

Conflicts are resolved by considering a lookup of the next input symbols.

43/66

Structure of the LL(1)-Parser:



- The parser accesses a frame of length 1 of the input;
- it corresponds to an item pushdown automaton, essentially;
- table $M[q, w]$ contains the rule of choice.

44/66

Topdown Parsing

Idea:

- Emanate from the item pushdown automaton
- Consider **the next input symbol** to determine the appropriate rule for the next expansion
- A grammar is called **LL(1)** if a unique choice is always possible

Definition:

A reduced grammar is called **LL(1)**, if for each two distinct rules $A \rightarrow \alpha, A \rightarrow \alpha' \in P$ and each derivation $S \rightarrow_i^* u A \beta$ with $u \in T^*$ the following is valid:

$$\text{First}_1(\alpha \beta) \cap \text{First}_1(\alpha' \beta) = \emptyset$$



45/66

Topdown Parsing

Example 1:

$$S \rightarrow \text{if } (E) S \text{ else } S \mid \text{while } (E) S \mid E; \\ E \rightarrow \text{id}$$

is **LL(1)**, since $\text{First}_1(E) = \{\text{id}\}$

Example 2:

$$S \rightarrow \text{if } (E) S \text{ else } S \mid \text{if } (E) S \mid \text{while } (E) S \mid E; \\ E \rightarrow \text{id}$$

... is **not LL(k)** for any $k > 0$.

46/66

Lookahead Sets

Definition: First₁-Sets

For a set $L \subseteq T^*$ we define:

$$\text{First}_1(L) = \{\epsilon \mid \epsilon \in L\} \cup \{u \in T \mid \exists v \in T^* : uv \in L\}$$

Example: $S \rightarrow \epsilon \mid a S b$

$\text{First}_1(\{S\})$
ϵ
$a b$
$a a b b$
$a a a b b b$
...

≡ the yield's prefix of length 1

47/66

Lookahead Sets

Arithmetics:

$\text{First}_1(_)$ is **distributive** with union and concatenation:

$$\begin{aligned} \text{First}_1(\emptyset) &= \emptyset \\ \text{First}_1(L_1 \cup L_2) &= \text{First}_1(L_1) \cup \text{First}_1(L_2) \\ \text{First}_1(L_1 \cdot L_2) &= \text{First}_1(\text{First}_1(L_1) \cdot \text{First}_1(L_2)) \\ &:= \text{First}_1(L_1) \odot_1 \text{First}_1(L_2) \end{aligned}$$

\odot_1 being 1 - concatenation

Definition: 1-concatenation

Let $L_1, L_2 \subseteq T \cup \{\epsilon\}$ with $L_1 \neq \emptyset \neq L_2$. Then:

$$L_1 \odot_1 L_2 = \begin{cases} L_1 & \text{if } \epsilon \notin L_1 \\ (L_1 \setminus \{\epsilon\}) \cup L_2 & \text{otherwise} \end{cases}$$

If all rules of G are productive, then all sets $\text{First}_1(A)$ are non-empty.

48/66

Lookahead Sets

For $\alpha \in (N \cup T)^*$ we are interested in the set:

$$\text{First}_1(\alpha) = \text{First}_1(\{w \in T^* \mid \alpha \rightarrow^* w\})$$

Idea: Treat ϵ separately: $\text{First}_1(A) = F_\epsilon(A) \cup \{\epsilon \mid A \rightarrow^* \epsilon\}$

- Let $\text{empty}(X) = \text{true}$ iff $X \rightarrow^* \epsilon$.
- $F_\epsilon(X_1 \dots X_m) = \bigcup_{i=1}^m F_\epsilon(X_i)$ if $\neg \text{empty}(X_j) \wedge \bigwedge_{i=1}^{j-1} \text{empty}(X_i)$

We characterize the ϵ -free First_1 -sets with an inequality system:

$$\begin{aligned} F_\epsilon(a) &= \{a\} & \text{if } a \in T \\ F_\epsilon(A) &\supseteq F_\epsilon(X_j) & \text{if } A \rightarrow X_1 \dots X_m \in P, \text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_{j-1}) \end{aligned}$$

49/66

Lookahead Sets

for example...

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid \text{name} \mid \text{int} \end{aligned}$$

with $\text{empty}(E) = \text{empty}(T) = \text{empty}(F) = \text{false}$

... we obtain:

$$\begin{aligned} F_\epsilon(S') &\supseteq F_\epsilon(E) & F_\epsilon(E) &\supseteq F_\epsilon(E) \\ F_\epsilon(E) &\supseteq F_\epsilon(T) & F_\epsilon(T) &\supseteq F_\epsilon(T) \\ F_\epsilon(T) &\supseteq F_\epsilon(F) & F_\epsilon(F) &\supseteq \{(\text{, name, int})\} \end{aligned}$$

50/66

Fast Computation of Lookahead Sets

Observation:

- The form of each inequality of these systems is:

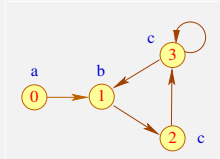
$$x \supseteq y \quad \text{resp.} \quad x \supseteq d$$

for variables x, y und $d \in D$.

- Such systems are called **pure unification problems**
- Such problems can be solved in **linear space/time**.

for example: $D = 2^{\{a,b,c\}}$

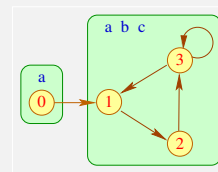
$$\begin{matrix} x_0 \supseteq \{a\} & x_1 \supseteq x_0 & x_1 \supseteq x_3 \\ x_1 \supseteq \{b\} & x_2 \supseteq x_1 & \\ x_2 \supseteq \{c\} & x_3 \supseteq x_2 & \\ x_3 \supseteq \{c\} & x_3 \supseteq x_2 & x_3 \supseteq x_3 \end{matrix}$$



Fast Computation of Lookahead Sets



Frank DeRemer & Tom Pennello



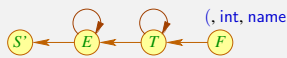
Proceeding:

- Create the **Variable Dependency Graph** for the inequality system.
- Within a **Strongly Connected Component** (\rightarrow Tarjan) all variables have the same value
- Is there no incoming edge for an SCC, its value is computed via the smallest upper bound of all values within the SCC
- In case of incoming edges, their values are also to be considered for the upper bound

Fast Computation of Lookahead Sets

... for our example grammar:

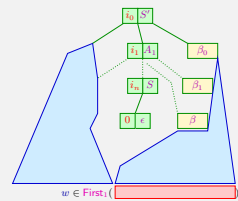
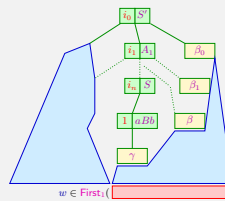
$First_1$:



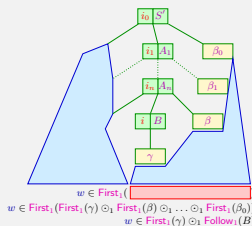
Item Pushdown Automaton as LL(1)-Parser

context is relevant too: $S' \rightarrow S \$ \quad S \rightarrow \epsilon^0 \mid a S b^1$

$First_1(\text{input})$	\$	a	b
S	?	?	?



Item Pushdown Automaton as LL(1)-Parser



Inequality system for $Follow_1(B) = First_1(\beta) \odot_1 \dots \odot_1 First_1(\beta_0)$

- $Follow_1(S) \supseteq \{\$\}$
- $Follow_1(B) \supseteq F_\epsilon(X_j)$ if $A \rightarrow \alpha B X_1 \dots X_m \in P, \text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_{j-1})$
- $Follow_1(B) \supseteq Follow_1(A)$ if $A \rightarrow \alpha B X_1 \dots X_m \in P, \text{empty}(X_1) \wedge \dots \wedge \text{empty}(X_m)$

Item Pushdown Automaton as LL(1)-Parser

Is G an $LL(1)$ -grammar, we can index a lookahead-table with items and nonterminals:

LL(1)-Lookahead Table

We set $M[B, w] = i$ with $B \rightarrow \gamma^i$ if $w \in First_1(\gamma) \odot_1 Follow_1(B)$

... for example: $S' \rightarrow S \$ \quad S \rightarrow \epsilon^0 \mid a S b^1$

$$First_1(S) = \{\epsilon, a\} \quad Follow_1(S) = \{b, \$\}$$

- S-rule 0: $First_1(\epsilon) \odot_1 Follow_1(S) = \{b, \$\}$
- S-rule 1: $First_1(aSb) \odot_1 Follow_1(S) = \{a\}$

	\$	a	b
S	0	1	0

Item Pushdown Automaton as LL(1)-Parser

For example: $S' \rightarrow S \$ \quad S \rightarrow \epsilon^0 \mid a S b^1$

The transitions of the according Item Pushdown Automaton:

0	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow \bullet S \$] [S \rightarrow \bullet]$
1	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow \bullet S \$] [S \rightarrow \bullet a S b]$
2	$[S \rightarrow \bullet a S b]$	a	$[S \rightarrow \bullet a S b]$
3	$[S \rightarrow \bullet a S b]$	ϵ	$[S \rightarrow \bullet a S b] [S \rightarrow \bullet]$
4	$[S \rightarrow \bullet a S b]$	ϵ	$[S \rightarrow \bullet a S b] [S \rightarrow \bullet a S b]$
5	$[S \rightarrow \bullet a S b] [S \rightarrow \bullet]$	ϵ	$[S \rightarrow \bullet a S b]$
6	$[S \rightarrow \bullet a S b] [S \rightarrow \bullet a S b]$	ϵ	$[S \rightarrow \bullet a S b]$
7	$[S \rightarrow \bullet a S b]$	b	$[S \rightarrow \bullet a S b]$
8	$[S' \rightarrow \bullet S \$]$	ϵ	$[S' \rightarrow \bullet S \$]$
9	$[S' \rightarrow \bullet S \$] [S \rightarrow \bullet a S b]$	ϵ	$[S' \rightarrow \bullet S \$]$

Lookahead table:

	\$	a	b
S	0	1	0

Left Recursion

Attention:

Many grammars are not $LL(k)$!

A reason for that is:

Definition

Grammar G is called **left-recursive**, if

$$A \rightarrow^+ A \beta \quad \text{for an } A \in N, \beta \in (T \cup N)^*$$

Example:

$$\begin{matrix} E \rightarrow E + T & | & T \\ T \rightarrow T * F & | & F \\ F \rightarrow (E) & | & \text{name} \mid \text{int} \end{matrix}$$

... is left-recursive

Left Recursion

Theorem:

Let a grammar G be reduced and **left-recursive**, then G is not $LL(k)$ for any k .

Proof:

Let wlog. $A \rightarrow A \beta \mid \alpha \in P$ and A be reachable from S

Assumption: G is $LL(k)$

$$\Rightarrow First_k(\alpha \beta^n \gamma) \cap First_k(\alpha \beta^{n+1} \gamma) = \emptyset$$

- Case 1: $\beta \rightarrow^+ \epsilon$ — Contradiction !!!
- Case 2: $\beta \rightarrow^+ w \neq \epsilon \implies First_k(\alpha w^k \gamma) \cap First_k(\alpha w^{k+1} \gamma) \neq \emptyset$

Right-Regular Context-Free Parsing

Recurring scheme in programming languages: Lists of sth...

$$S \rightarrow b \mid S a b$$

Alternative idea: **Regular Expressions**

$$S \rightarrow (b a)^* b$$

Definition: Right-Regular Context-Free Grammar

A **right-regular context-free grammar** (RR-CFG) is a

4-tuple $G = (N, T, P, S)$ with:

- N the set of nonterminals,
- T the set of terminals,
- P the set of rules with **regular expressions of symbols** as rhs,
- $S \in N$ the start symbol

Example: Arithmetic Expressions

$$\begin{matrix} S \rightarrow E \\ E \rightarrow T (+ T)^* \\ T \rightarrow F (* F)^* \\ F \rightarrow (E) \mid \text{name} \mid \text{int} \end{matrix}$$

Idea 1: Rewrite the rules from G to $\langle G \rangle$:

A	$\rightarrow \langle \alpha \rangle$	if $A \rightarrow \alpha \in P$
$\langle \alpha \rangle$	$\rightarrow \alpha$	if $\alpha \in N \cup T$
$\langle \epsilon \rangle$	$\rightarrow \epsilon$	
$\langle \alpha^* \rangle$	$\rightarrow \epsilon \mid \langle \alpha \rangle \langle \alpha^* \rangle$	if $\alpha \in \text{Regex}_{T,N}$
$\langle \alpha_1 \dots \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \dots \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$
$\langle \alpha_1 \mid \dots \mid \alpha_n \rangle$	$\rightarrow \langle \alpha_1 \rangle \mid \dots \mid \langle \alpha_n \rangle$	if $\alpha_i \in \text{Regex}_{T,N}$

... and generate the according LL(k)-Parser $M_{\langle G \rangle}^L$

Example: Arithmetic Expressions cont'd

S	$\rightarrow E$
E	$\rightarrow T \mid (+T)^* \mid (T + T)^*$
T	$\rightarrow F \mid (*F)^* \mid (F * F)^*$
F	$\rightarrow (E) \mid \text{name} \mid \text{int}$
$\langle T (+T)^* \rangle$	$\rightarrow T \langle (+T)^* \rangle$
$\langle (+T)^* \rangle$	$\rightarrow \epsilon \mid (+T) \langle (+T)^* \rangle$
$\langle +T \rangle$	$\rightarrow +T$
$\langle F (*F)^* \rangle$	$\rightarrow F \langle (*F)^* \rangle$
$\langle (*F)^* \rangle$	$\rightarrow \epsilon \mid (*F) \langle (*F)^* \rangle$
$\langle *F \rangle$	$\rightarrow *F$



Reinhold Heckmann

Definition:

An RR -CFG G is called **RLL(1)**, if the corresponding CFG $\langle G \rangle$ is an LL(1) grammar.

Discussion

- directly yields the table driven parser $M_{\langle G \rangle}^L$ for RLL(1) grammars
- however: mapping regular expressions to recursive productions unnecessarily strains the stack
→ instead directly construct automaton in the style of Berry-Sethi

Idea 2: Recursive Descent RLL Parsers:

Recursive descent RLL(1)-parsers are an alternative to table-driven parsers; apart from the usual function `scan()`, we generate a program frame with the lookahead function `expect()` and the main parsing method `parse()`:

```

int next;
void expect(Set E){
    if ({ε, next} ∩ E = ∅){
        cerr << "Expected" << E << "found" << next;
        exit(0);
    }
    return;
}
void parse(){
    next = scan();
    expect(First1(S));
    S();
    expect({EOF});
}
    
```

Idea 2: Recursive Descent RLL Parsers:

For each $A \rightarrow \alpha \in P$, we introduce:

```

void A(){
    generate(α)
}
    
```

with the meta-program `generate` being defined by structural decomposition of α :

```

generate(r1 ... rk) = generate(r1)
                    expect(First1(r2));
                    generate(r2)
                    :
                    expect(First1(rk));
                    generate(rk)
generate(ε)         = ;
generate(a)         = next = scan();
generate(A)         = A();
    
```

Idea 2: Recursive Descent RLL Parsers:

```

generate(r*)        = while (next ∈ F_r(r)) {
                    generate(r)
                    }
generate(r1 | ... | rk) = switch(next) {
                    labels(First1(r1)) generate(r1) break;
                    :
                    labels(First1(rk)) generate(rk) break;
                    }
labels({α1, ..., αm}) = label(α1) : ... label(αm);
label(α)              = case α
label(ε)              = default
    
```

Topdown-Parsing

Discussion

- A practical implementation of an RLL(1)-parser via recursive descent is a straight-forward idea
- However, **only a subset** of the deterministic contextfree languages can be parsed this way.
- As soon as `First1()` sets are not disjoint any more,
 - Solution 1: For many accessibly written grammars, the alternation between right hand sides happens too early. Keeping the common prefixes of right hand sides joined and introducing a new production for the actual diverging sentence forms often helps.
 - Solution 2: Introduce **ranked** grammars, and decide conflicting lookahead always in favour of the higher ranked alternative
→ relation to LL parsing not so clear any more
→ not so clear for * operator how to decide
 - Solution 3: Going from LL(1) to LL(k)
The size of the occurring sets is rapidly increasing with larger k
Unfortunately, even LL(k) parsers are not sufficient to accept all deterministic contextfree languages. (regular lookahead → LL(*))
- In practical systems, this often motivates the implementation of $k = 1$ only ...

Topic:
Syntactic Analysis - Part II

Syntactic Analysis - Part II

Chapter 1:
Bottom-up Analysis

Shift-Reduce Parser



Donald Knuth

Idea:

We **delay** the decision whether to reduce until we know, whether the input matches the right-hand-side of a rule!

Construction: Shift-Reduce parser M_G^R

- The input is shifted successively to the pushdown.
- Is there a **complete right-hand side** (a **handle**) atop the pushdown, it is replaced (**reduced**) by the corresponding left-hand side

Shift-Reduce Parser

Example:

$S \rightarrow AB$
 $A \rightarrow a$
 $B \rightarrow b$

The pushdown automaton:

States: q_0, f, a, b, A, B, S ;
Start state: q_0
End state: f

q_0	a	$q_0 a$
a	ϵ	A
A	b	Ab
b	ϵ	B
AB	ϵ	S
$q_0 S$	ϵ	f

Shift-Reduce Parser

Construction:

In general, we create an automaton $M_G^R = (Q, T, \delta, q_0, F)$ with:

- $Q = T \cup N \cup \{q_0, f\}$ (q_0, f fresh);
- $F = \{f\}$;
- Transitions:

$$\delta = \{(q, x, qx) \mid q \in Q, x \in T\} \cup \text{// Shift-transitions}$$

$$\{(\alpha, \epsilon, A) \mid A \rightarrow \alpha \in P\} \cup \text{// Reduce-transitions}$$

$$\{(q_0 S, \epsilon, f)\} \text{// finish}$$

Example-computation:

$(q_0, ab) \vdash (q_0 a, b) \vdash (q_0 A, b)$
 $\vdash (q_0 A b, \epsilon) \vdash (q_0 AB, \epsilon)$
 $\vdash (q_0 S, \epsilon) \vdash (f, \epsilon)$

Shift-Reduce Parser

Observation:

- The sequence of reductions corresponds to a **reverse rightmost-derivation** for the input
- To prove correctness, we have to prove:

$$(\epsilon, w) \vdash^* (A, \epsilon) \text{ iff } A \rightarrow^* w$$

- The shift-reduce pushdown automaton M_G^R is in general also **non-deterministic**
- For a deterministic parsing-algorithm, we have to identify computation-states for reduction

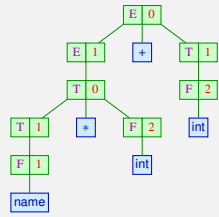
⇒ LR-Parsing

The Pushdown During an RR-Derivation

Idea: Observe a successful run of M_G^R !

Input:
counter + 2 + 40

Pushdown:
(q_0)



$E \rightarrow E+T^0 \mid T^1$
 $T \rightarrow T*F^0 \mid F^1$
 $F \rightarrow (E)^0 \mid \text{name}^1 \mid \text{int}^2$

Result:

Viable Prefixes and Admissible Items

Formalism: use **Items** as representations of **prefixes of righthand sides**

Generic Agreement

In a sequence of configurations of M_G^R

$(q_0 \alpha \gamma, v) \vdash (q_0 \alpha B, v) \vdash^* (q_0 S, \epsilon)$

we call $\alpha \gamma$ a **viable prefix** for the complete item $[B \rightarrow \gamma \bullet]$.

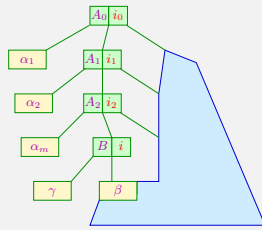
Reformulating the Shift-Reduce-Parsers main problem:

Find the items, for which the content of M_G^R 's stack is the viable prefix....

→ **Admissible Items**

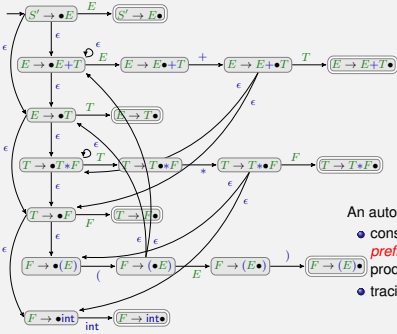
Admissible Items

The item $[B \rightarrow \gamma \bullet \beta]$ is called **admissible** for $\alpha \gamma$ iff $S \rightarrow_R^* \alpha B v$:



... with $\alpha = \alpha_1 \dots \alpha_m$

Characteristic Automaton



An automaton...

- consuming pushdown symbols, i.e. **prefixes of righthand sides** of productions expanding from S
- tracing admissible items in its states

Characteristic Automaton

Observation:

One can now consume the shift-reduce parser's pushdown with the characteristic automaton: If the input $(N \cup T)^*$ for the characteristic automaton corresponds to a viable prefix, its state contains the admissible items.

States: Items

Start state: $[S' \rightarrow \bullet S]$

Final states: $\{[B \rightarrow \gamma \bullet] \mid B \rightarrow \gamma \in P\}$

Transitions:

- (1) $([A \rightarrow \alpha \bullet X \beta], X, [A \rightarrow \alpha X \bullet \beta])$, $X \in (N \cup T), A \rightarrow \alpha X \beta \in P$;
- (2) $([A \rightarrow \alpha \bullet B \beta], \epsilon, [B \rightarrow \bullet \gamma])$, $A \rightarrow \alpha B \beta, B \rightarrow \gamma \in P$;

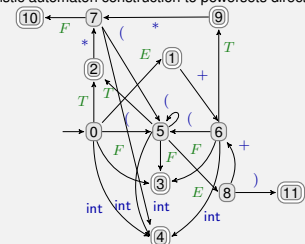
The automaton $c(G)$ is called **characteristic automaton** for G .

Canonical LR(0)-Automaton

The **canonical LR(0)-automaton** $LR(G)$ is created from $c(G)$ by:

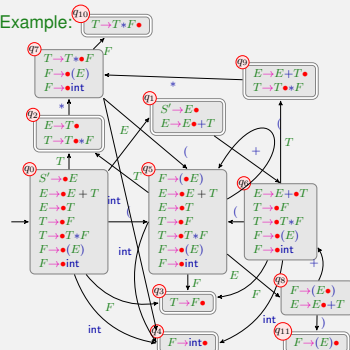
- performing arbitrarily many ϵ -transitions after every consuming transition
- performing the powerset construction
- Idea: or rather apply characteristic automaton construction to powersets directly?

... for example:



Canonical LR(0)-Automaton – Example:

$S' \rightarrow E$
 $E \rightarrow E+T \mid T$
 $T \rightarrow T*F \mid F$
 $F \rightarrow (E) \mid \text{int}$



Canonical LR(0)-Automaton

Observation:

The canonical LR(0)-automaton can be created **directly** from the grammar. For this we need a helper function δ_c^* (ϵ -closure)

$\delta_c^*(q) = q \cup \{[B \rightarrow \bullet \gamma] \mid B \rightarrow \gamma \in P,$
 $[A \rightarrow \alpha \bullet B' \beta'] \in q,$
 $B' \rightarrow \beta \beta\}$

We define:

States: Sets of items;

Start state: $\delta_c^*\{[S' \rightarrow \bullet S]\}$

Final states: $\{q \mid [A \rightarrow \alpha \bullet] \in q\}$

Transitions: $\delta(q, X) = \delta_c^*\{[A \rightarrow \alpha X \bullet \beta] \mid [A \rightarrow \alpha \bullet X \beta] \in q\}$

LR(0)-Parser

Idea for a parser:

- The parser manages a viable prefix $\alpha = X_1 \dots X_m$ on the pushdown and uses $LR(G)$ to identify reduction spots.
- It can reduce with $A \rightarrow \gamma$, if $[A \rightarrow \gamma \bullet]$ is admissible for α

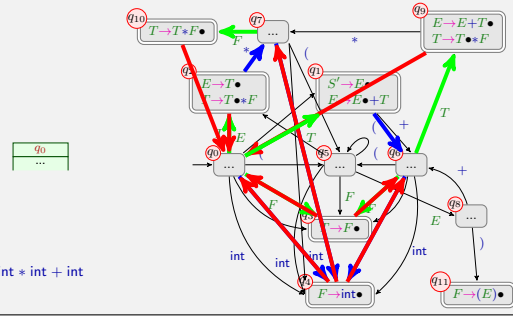
Optimization:

We push the **states** instead of the X_i in order not to process the pushdown's content with the automaton anew all the time.
Reduction with $A \rightarrow \gamma$ leads to popping the uppermost $|\gamma|$ states and continue with the state on top of the stack and input A .

Attention:

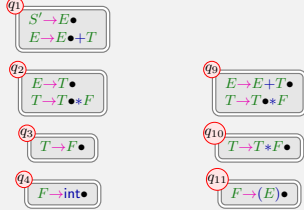
This parser is only **deterministic**, if each final state of the canonical $LR(0)$ -automaton is **conflict free**.

LR(0)-Parser – Example:



LR(0)-Parser

... we observe:



The final states q_1, q_2, q_9 contain more than one admissible item \Rightarrow non-deterministic!

LR(0)-Parser

The construction of the $LR(0)$ -parser:

- States: $Q \cup \{f\}$ (f fresh)
- Start state: q_0
- Final state: f
- Transitions:
 - Shift: $(p, a, p q)$ if $q = \delta(p, a) \neq \emptyset$
 - Reduce: $(p q_1 \dots q_m, \epsilon, p q)$ if $[A \rightarrow X_1 \dots X_m \bullet] \in q_m, q = \delta(p, A)$
 - Finish: $(q_0 p, \epsilon, f)$ if $[S' \rightarrow S \bullet] \in p$

with the canonical automaton $LR(G) = (Q, T, \delta, q_0, F)$.

LR(0)-Parser

Correctness:

we show:

The accepting computations of an $LR(0)$ -parser are one-to-one related to those of a shift-reduce parser M_G^R .

we conclude:

- The accepted language is exactly $\mathcal{L}(G)$
- The sequence of reductions of an accepting computation for a word $w \in T$ yields a **reverse rightmost derivation** of G for w

LR(0)-Parser

Attention: Unfortunately, the $LR(0)$ -parser is in general non-deterministic.

We identify two reasons for a state $q \in Q$:

Reduce-Reduce-Conflict:



Shift-Reduce-Conflict:



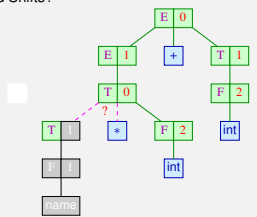
Those states are called $LR(0)$ -unsuited.

Revisiting the Conflicts of the $LR(0)$ -Automaton

What differentiates the particular Reductions and Shifts?

Input:
 $* 2 + 40$

Pushdown:
($q_0 T$)



$E \rightarrow E+T \quad | \quad T$
 $T \rightarrow T * F \quad | \quad F$
 $F \rightarrow (E) \quad | \quad \text{int}$

LR(k)-Grammars

Idea: Consider k -lookahead in conflict situations.

Definition:

The reduced contextfree grammar G is called $LR(k)$ -grammar, if $\alpha \beta w \mid_{|\alpha \beta|+k} = \alpha' \beta' w' \mid_{|\alpha \beta|+k}$ with:

$$\left. \begin{array}{l} S \xrightarrow{R} \alpha A w \rightarrow \alpha \beta w \\ S \xrightarrow{R} \alpha' A' w' \rightarrow \alpha' \beta' w' \end{array} \right\} \text{ follows: } \alpha = \alpha' \wedge \beta = \beta' \wedge A = A'$$

Strategy for testing Grammars for $LR(k)$ -property

- Focus iteratively on all rightmost derivations $S \xrightarrow{R} \alpha X w \rightarrow \alpha \beta w$
- Iterate over $k \geq 0$
 - For each $\gamma = \alpha \beta w \mid_{|\alpha \beta|+k}$ (**handle with k -lookahead**) check if there exists a differently right-derivable $\alpha' \beta' w'$ for which $\gamma = \alpha' \beta' w' \mid_{|\alpha \beta|+k}$
 - if there is none, we have found no objection against k being enough lookahead to disambiguate $\alpha \beta w$ from other rightmost derivations

LR(k)-Grammars

for example:

- (1) $S \rightarrow A \mid B \quad A \rightarrow aAb \mid 0 \quad B \rightarrow aBbb \mid 1$
... is not $LL(k)$ for any k — but $LR(0)$:

Let $S \xrightarrow{R} \alpha X w \rightarrow \alpha \beta w$. Then $\alpha \beta$ is of one of these forms:

$$\underline{A}, \underline{B}, a^n \underline{aAb}, a^n \underline{aBbb}, a^n \underline{0}, a^n \underline{1} \quad (n \geq 0)$$

- (2) $S \rightarrow aAc \quad A \rightarrow Abbb \mid b$

... is also not $LL(k)$ for any k — but again $LR(0)$:

Let $S \xrightarrow{R} \alpha X w \rightarrow \alpha \beta w$. Then $\alpha \beta$ is of one of these forms:

$$a \underline{b}, a \underline{Abbb}, a \underline{Ac}$$

LR(k)-Grammars

for example:

- (3) $S \rightarrow aAc \quad A \rightarrow bba \mid b$... is not $LR(0)$, but $LR(1)$:

Let $S \xrightarrow{R} \alpha X w \rightarrow \alpha \beta w$ with $\{y\} = \text{First}_k(w)$ then $\alpha \beta y$ is of one of these forms:

$$a b^{2n} \underline{bc}, a b^{2n} \underline{bba}c, a \underline{Ac}$$

- (4) $S \rightarrow aAc \quad A \rightarrow bAb \mid b$... is not $LR(k)$ for any $k \geq 0$:

Consider the rightmost derivations:

$$S \xrightarrow{R} a b^n A b^n c \rightarrow a b^n \underline{bb}^n c$$

LR(1)-Parsing

Idea: Let's equip items with 1-lookahead

Definition LR(1)-Item

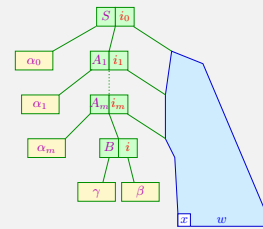
An LR(1)-item is a pair $[B \rightarrow \alpha \bullet \beta, x]$ with

$$x \in \text{Follow}_1(B) = \bigcup \{ \text{First}_1(\nu) \mid S \rightarrow^* \mu B \nu \}$$

Admissible LR(1)-Items

The LR(1)-item $[B \rightarrow \gamma \bullet \beta, x]$ is **admissible** for $\alpha \gamma$ if:

$$S \rightarrow^* \alpha B w \quad \text{with} \quad \{x\} = \text{First}_1(w)$$



... with $\alpha_0 \dots \alpha_m = \alpha$

The Characteristic LR(1)-Automaton

The set of admissible LR(1)-items for viable prefixes is again computed with the help of the finite automaton $c(G, 1)$.

The automaton $c(G, 1)$:

States: LR(1)-items

Start state: $\{S' \rightarrow \bullet S, \$\}$

Final states: $\{[B \rightarrow \gamma \bullet, x] \mid B \rightarrow \gamma \in P, x \in \text{Follow}_1(B)\}$

Transitions: (1) $([A \rightarrow \alpha \bullet X \beta, x], X, [A \rightarrow \alpha X \bullet \beta, x]), X \in (N \cup T)$
 (2) $([A \rightarrow \alpha \bullet B \beta, x], \epsilon, [B \rightarrow \bullet \gamma, x']), A \rightarrow \alpha B \beta, B \rightarrow \gamma \in P, x' \in \text{First}_1(\beta) \circledast_1 \{x\}$

This automaton works like $c(G)$ — but additionally manages a 1-prefix from Follow_1 of the left-hand sides.

The Canonical LR(1)-Automaton

The canonical LR(1)-automaton $LR(G, 1)$ is created from $c(G, 1)$, by performing arbitrarily many ϵ -transitions and then making the resulting automaton **deterministic** ...

But again, it can be constructed **directly** from the grammar; analogously to LR(0), we need the ϵ -closure δ_ϵ^* as a helper function:

$$\delta_\epsilon^*(q) = q \cup \{ [C \rightarrow \bullet \gamma, x] \mid [A \rightarrow \alpha \bullet B \beta', x'] \in q, B \rightarrow^* C \beta, C \rightarrow \gamma \in P, x \in \text{First}_1(\beta \beta') \circledast_1 \{x'\} \}$$

Then, we define:

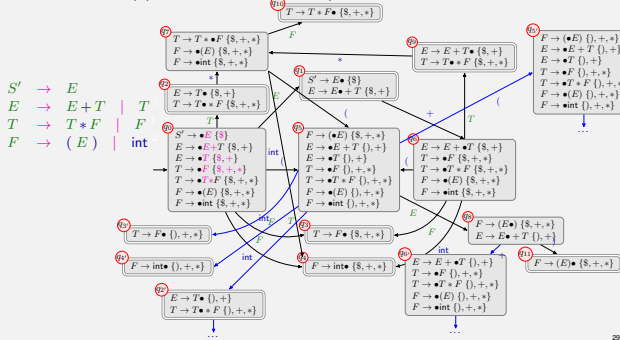
States: Sets of LR(1)-items;

Start state: $\delta_\epsilon^* \{S' \rightarrow \bullet S, \$\}$

Final states: $\{q \mid [A \rightarrow \alpha \bullet, x] \in q\}$

Transitions: $\delta(q, X) = \delta_\epsilon^* \{ [A \rightarrow \alpha \bullet X \beta, x] \mid [A \rightarrow \alpha \bullet X \beta, x] \in q \}$

The Canonical LR(1)-Automaton — for example:

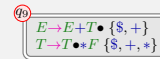


The Canonical LR(1)-Automaton

Discussion:

• In the example, the number of states was almost doubled ... and it can become even worse

• The conflicts in states q_{12}, q_{19} are now resolved !
 e.g. we have:



with:

$$\{\$, +\} \cap (\text{First}_1(*F) \circledast_1 \{\$, +, *\}) = \{\$, +\} \cap \{\$, +, *\} = \emptyset$$

The Action Table:

During practical parsing, we want to represent states just via an integer id. However, when the canonical LR(1)-automaton reaches a final state, we want to know **how to reduce/shift**. Thus we introduce...

The construction of the action table:

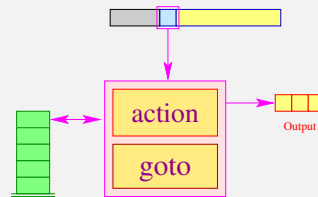
Type: $\text{action} : Q \times T \rightarrow LR(0)\text{-Items} \cup \{s, \text{error}\}$

Reduce: $\text{action}[q, w] = [A \rightarrow \beta \bullet]$ if $[A \rightarrow \beta \bullet, w] \in q$

Shift: $\text{action}[q, w] = s$ if $[A \rightarrow \beta \bullet b \gamma, a] \in q, w \in \text{First}_1(b \gamma) \circledast_1 \{a\}$

Error: $\text{action}[q, w] = \text{error}$ else

The LR(1)-Parser:



• The goto-table encodes the transitions:

$$\text{goto}[q, X] = \delta(q, X) \in Q$$

• The action-table describes for every state q and possible lookahead w the necessary action.

The LR(1)-Parser:

The construction of the LR(1)-parser:

States: $Q \cup \{f\}$ (f fresh)

Start state: q_0

Final state: f

Transitions:

Shift: $(p, a, p q)$ if $a = w,$
 $s = \text{action}[p, a],$
 $q = \text{goto}[p, a]$
Reduce: $(p q_1 \dots q_j \beta_1, \epsilon, p q)$ if $q_j \beta_1 \in F,$
 $[A \rightarrow \beta \bullet] = \text{action}[q_j \beta_1, w],$
 $q = \text{goto}[p, A]$
Finish: $(q_0 p, \epsilon, f)$ if $[S' \rightarrow S \bullet, \$] \in p$

with $LR(G, 1) = (Q, T, \delta, q_0, F)$ and the lookahead w .

The LR(1)-Parser:

Possible actions are:

shift $(A \rightarrow \gamma)$ // Shift-operation
 reduce $(A \rightarrow \gamma)$ // Reduction with callback/output
 error // Error

... for example:

$S' \rightarrow E$
 $E \rightarrow E + T^0 \mid T^1$
 $T \rightarrow T * F^0 \mid F^1$
 $F \rightarrow (E)^0 \mid \text{int}^1$

action	\$	int	()	+	*
q1	S',0				s	
q2	E,1				E,1	s
q3	T,1				T,1	T,1
q4	F,1				F,1	F,1
q5					F,1	F,1
q6	E,0				E,0	s
q7					E,0	s
q8	T,0				T,0	T,0
q9					T,0	T,0
q10	F,0				F,0	F,0
q11					F,0	F,0
q12					F,0	F,0

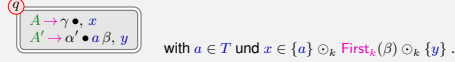
The Canonical LR(1)-Automaton

In general: We identify two conflicts for a state $q \in Q$:

Reduce-Reduce-Conflict:



Shift-Reduce-Conflict:



Such states are now called **LR(1k)-unsuited**

Theorem:

A reduced contextfree grammar G is called **LR(k)** iff the canonical **LR(k)**-automaton $LR(G, k)$ has no **LR(k)**-unsuited states.

Precedences

Many parser generators give the chance to fix Shift-/Reduce-Conflicts by patching the action table either by hand or with **token precedences**.

... for example:

$S' \rightarrow E^0$
 $E \rightarrow E + E^0$
 $\quad \quad \quad E * E^1$
 $\quad \quad \quad (E)^2$
 $\quad \quad \quad \text{int}^3$

Shift-/Reduce Conflict in state 8:

$[E \rightarrow E \bullet + E^0, +]$
 $[E \rightarrow E + E \bullet^0, +]$

$\langle \gamma E + E, + \omega \rangle \Rightarrow$ **Associativity**

+ **left associative**

action	\$	int	()	+	*
q0	S',0			s	s
q1	E,3		E,3	E,3	E,3
q2	s			s	s
q3	s			s	s
q4	s			s	s
q5	E,2		E,2	E,2	E,2
q6	s			s	s
q7	E,1		E,1	?	?
q8	E,0		E,0	E,0	?
q9	s			s	s

Precedences

Many parser generators give the chance to fix Shift-/Reduce-Conflicts by patching the action table either by hand or with **token precedences**.

... for example:

$S' \rightarrow E^0$
 $E \rightarrow E + E^0$
 $\quad \quad \quad E * E^1$
 $\quad \quad \quad (E)^2$
 $\quad \quad \quad \text{int}^3$

Shift-/Reduce Conflict in state 7:

$[E \rightarrow E \bullet * E^1, *]$
 $[E \rightarrow E * E \bullet^1, *]$

$\langle \gamma E * E, * \omega \rangle \Rightarrow$ **Associativity**

* **right associative**

action	\$	int	()	+	*
q0	S',0			s	s
q1	E,3		E,3	E,3	E,3
q2	s			s	s
q3	s			s	s
q4	s			s	s
q5	E,2		E,2	E,2	E,2
q6	s			s	s
q7	E,1		E,1	?	s
q8	E,0		E,0	E,0	?
q9	s			s	s

Precedences

Many parser generators give the chance to fix Shift-/Reduce-Conflicts by patching the action table either by hand or with **token precedences**.

... for example:

$S' \rightarrow E^0$
 $E \rightarrow E + E^0$
 $\quad \quad \quad E * E^1$
 $\quad \quad \quad (E)^2$
 $\quad \quad \quad \text{int}^3$

Shift-/Reduce Conflict in states 8, 7:

$[E \rightarrow E \bullet * E^1, *]$
 $[E \rightarrow E + E \bullet^0, *]$

$\langle \gamma E * E, + \omega \rangle$

$[E \rightarrow E \bullet + E^0, +]$
 $[E \rightarrow E * E \bullet^1, +]$

$\langle \gamma E + E, * \omega \rangle$

* **higher precedence**

+ **lower precedence**

action	\$	int	()	+	*
q0	S',0			s	s
q1	E,3		E,3	E,3	E,3
q2	s			s	s
q3	s			s	s
q4	s			s	s
q5	E,2		E,2	E,2	E,2
q6	s			s	s
q7	E,1		E,1	E,1	s
q8	E,0		E,0	E,0	s
q9	s			s	s

What if precedences are not enough?

Example (very simplified lambda expressions):

$E \rightarrow (E)^0 \text{ident}^1 \mid L^2$
 $L \rightarrow (\text{args}) \Rightarrow E^0$
 $(\text{args}) \rightarrow ((\text{idlist})^0 \mid \text{ident}^1$
 $(\text{idlist}) \rightarrow (\text{idlist})^0 \mid \text{ident}^1$

E rightmost-derives these forms among others:

$(\text{ident}), (\text{ident}) \Rightarrow \text{ident}, \dots \Rightarrow$ at least **LR(2)**

Naive Idea:

poor man's **LR(2)** by combining the tokens $)$ and \Rightarrow during lexical analysis into a single token \Rightarrow .

⚠ in this case obvious solution, but in general not so simple

What if precedences are not enough?

In practice, **LR(k)**-parser generators working with the lookahead sets of sizes larger than $k = 1$ are not common, since computing lookahead sets with $k > 1$ blows up exponentially. However,

- there exist several practical **LR(k)** grammars of $k > 1$, e.g. Java 1.6+ (**LR(2)**)
- often, more lookahead is only exhausted locally
- should we really give up, whenever we are confronted with a Shift-/Reduce-Conflict?



Theorem: LR(k)-to-LR(1)

Any **LR(k)** grammar can be directly transformed into an equivalent **LR(1)** grammar.

LR(2) to LR(1)

... Example:

$S \rightarrow A b b^0 \mid B b c^1$
 $A \rightarrow a A^0 \mid a^1$
 $B \rightarrow a B^0 \mid a^1$

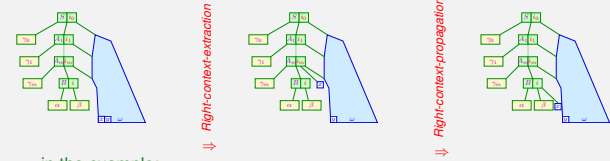
S rightmost-derives one of these forms:

$a^n a b b, a^n a b c, a^n a A b b, a^n a B b c, A b b, B b c \Rightarrow$ **LR(2)**

in **LR(1)**, you will have Reduce-/Reduce-Conflicts between the productions $A, 1$ and $B, 1$ under lookahead b

LR(2) to LR(1)

Basic Idea:



in the example:

Right-context is already extracted, so we only perform **Right-context-propagation**:

$S \rightarrow A b b^0 \mid B b c^1$
 $A \rightarrow a A^0 \mid a^1$
 $B \rightarrow a B^0 \mid a^1$

$S \rightarrow (A b) b^0 \mid (B b) c^1$
 $(A b) \rightarrow a (A b)^0 \mid a b^1$
 $(B b) \rightarrow a (B b)^0 \mid a b^1$

unreachable

LR(2) to LR(1)

Example cont'd:

$S \rightarrow A' b^0 \mid B' c^1$
 $A' \rightarrow a A'^0 \mid a b^1$
 $B' \rightarrow a B'^0 \mid a b^1$

S rightmost-derives one of these forms:

$a^n a b b, a^n a b c, a^n a A' b, a^n a B' c, A' b, B' c \Rightarrow$ **LR(1)**

LR(2) to LR(1)

Example 2:

$S \rightarrow b S S^0$
 $\quad \quad \quad a^1$
 $\quad \quad \quad a a c^2$

S rightmost-derives these forms among others:

$b S S, b S a, b S a a c, b a a, b a a c a, b a a c a c, \dots \Rightarrow$ min. **LR(2)**

in **LR(1)**, you will have (at least) Shift-/Reduce-Conflicts between the items $[S \rightarrow a \bullet, a]$ and $[S \rightarrow a \bullet a c]$

$[S \rightarrow a]$'s right context is a nonterminal \Rightarrow perform **Right-context-extraction**

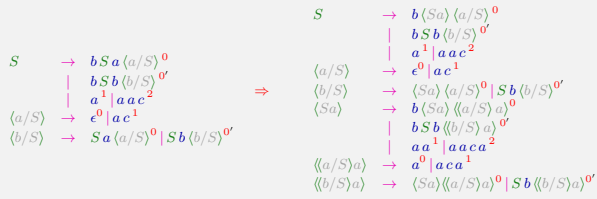
$S \rightarrow b S S^0$
 $\quad \quad \quad a^1$
 $\quad \quad \quad a a c^2$

$S \rightarrow b S a \langle a/S \rangle^0 \mid b S b \langle b/S \rangle^0$
 $\langle a/S \rangle \rightarrow \epsilon^0 \mid a c^1$
 $\langle b/S \rangle \rightarrow S a \langle a/S \rangle^0 \mid S b \langle b/S \rangle^0$

LR(2) to LR(1)

Example 2 cont'd:

[$S \rightarrow a$]'s right context is now terminal $a \Rightarrow$ perform *Right-context-propagation*

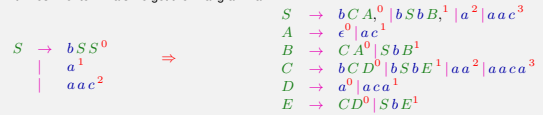


43/49

LR(2) to LR(1)

Example 2 finished:

With fresh nonterminals we get the final grammar



44/49

Chapter 2: LR(k)-Parser Design

45/49

LR(k)-Parser Design

```

S' ::= E:e           { RESULT = e;   ; }
      ;
      E::E:e plus T:t { RESULT = e + t; ; }
      | T:t           { RESULT = t;   ; }
      ;
T ::= T:t times F:f  { RESULT = t * f; ; }
      | F:f           { RESULT = f;   ; }
      ;
F ::= lbrac E:e rbrac { RESULT = e;   ; }
      | intconst:c   { RESULT = c;   ; }
      ;
    
```

Parser Actions

For each rule, specify user code to be executed in case of reduction actions.

- add code sections delimited with { ; } to each variant
- produce results by assigning values to **RESULT**
- add labels to symbols to refer to former results

Implementation Idea: add data stack that

- pushes **RESULT** after each user action
- translates labeled symbols to offset from top of stack based on the position in the rhs

46/49

A Practical Example: Type Definitions in ANSI C

A type definition is a *synonym* for a type expression. In C they are introduced using the **typedef** keyword. Type definitions are useful

- as abbreviation:


```
typedef struct { int x; int y; } point_t;
```
- to construct *recursive* types:

Possible declaration in C: more readable:

```

struct list {
  int info;
  struct list* next;
}
struct list* head;

typedef struct list list_t;
struct list {
  int info;
  list_t* next;
}
list_t* head;
    
```

47/49

A Practical Example: Type Definitions in ANSI C

The C grammar distinguishes **typename** and **identifier**. Consider the following declarations:

```

typedef struct { int x,y } point_t;
point_t origin;
    
```

Idea: in a *parser action* maintain a shared list between parser and scanner to communicate identifiers to report as typenames

Relevant C grammar:

```

declaration → (declarationspecifier)+ declarator ;
declarationspecifier → static | volatile ... typedef
                    | void | char | char ... typename
declarator → identifier | ...
    
```

Problem:

During reduction of the declaration, the scanner eagerly provides a new lookahead token, thus has already interpreted `point_t` in line 2 as **identifier**

48/49

A Practical Example: Type Definitions in ANSI C: Solutions

Relevant C grammar:

```

declaration → (declarationspecifier)+ declarator ;
declarationspecifier → static | volatile ... typedef
                    | void | char | char ... typename
declarator → identifier | ...
    
```

Solution is difficult:

- try to fix the lookahead token class within the scanner-parser-channel Δ a *mess*
- add a rule to the grammar, to make it context-free:

```

typename → identifier ambiguous

Example input: (mytype1) (mytype2);
castexpr → ( typename ) castexpr
postfixexpr → postfixexpr ( expression )

register identifier as typename before lookahead is harmful
declaration → (declarationspecifier)+ declarator { act(); };
    
```

49/49

Topic:

Semantic Analysis

1/67

Semantic Analysis

Scanner and parser accept programs with correct syntax.

- not all programs that are syntactically correct make *sense*
- the compiler may be able to *recognize* some of these
 - these programs are rejected and reported as *erroneous*
 - the language definition defines what *erroneous* means
- **semantic analyses** are necessary that, for instance:
 - check that **identifiers** are known and where they are defined
 - check the **type-correct** use of variables
- **semantic analyses** are also useful to
 - find possibilities to "optimize" the program
 - **warn** about possibly incorrect programs

~ a semantic analysis annotates the syntax tree with **attributes**

2/67

Chapter 1: Attribute Grammars

3/67

Attribute Grammars

- many computations of the semantic analysis as well as the code generation operate on the syntax tree
- what is computed at a given node only depends on the *type* of that node (which is usually a non-terminal)
- we call this a *local* computation:
 - only accesses already computed information from neighbouring nodes
 - computes new information for the current node and other neighbouring nodes

Definition attribute grammar

An **attribute grammar** is a CFG extended by

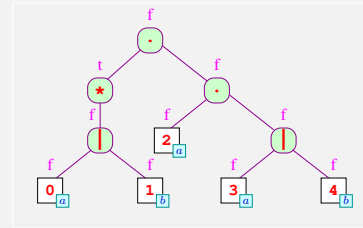
- a set of attributes for each non-terminal and terminal
- local attribute equations

- in order to be able to evaluate the attribute equations, all attributes mentioned in that equation have to be evaluated already
 ~> the nodes of the syntax tree need to be visited in a certain *sequence*

4/67

Example: Computation of the $empty[r]$ Attribute

Consider the syntax tree of the regular expression $(a|b)^*a(a|b)$:



~> equations for $empty[r]$ are computed from bottom to top (aka *bottom-up*)

5/67

Implementation Strategy

- attach an attribute $empty$ to every node of the syntax tree
- compute the attributes in a *depth-first post-order* traversal:
 - at a leaf, we can compute the value of $empty$ without considering other nodes
 - the attribute of an inner node only depends on the attribute of its children
- the $empty$ attribute is a *synthesized* attribute

in general:

Definition

An attribute at N is called

- inherited* if its value is defined in terms of attributes of N 's parent, siblings and/or N itself (root \rightarrow leaves)
- synthesized* if its value is defined in terms of attributes of N 's children and/or N itself (leaves \rightarrow root)

6/67

Example: Attribute Equations for $empty$

In order to compute an attribute *locally*, specify attribute equations for each node depending on the *type* of the node:

In the Example from earlier, we did that intuitively:

for leaves: $r \equiv \boxed{x}$ we define $empty[r] = (x \equiv \epsilon)$
 otherwise:

$$\begin{aligned} empty[r_1 | r_2] &:= empty[r_1] \vee empty[r_2] \\ empty[r_1 \cdot r_2] &:= empty[r_1] \wedge empty[r_2] \\ empty[r_1^+] &:= t \\ empty[r_1^?] &:= t \end{aligned}$$

7/67

Specification of General Attribute Systems

General Attribute Systems

In general, for establishing attribute systems we need a flexible way to *refer to parents and children*:

~> We use consecutive indices to refer to neighbouring attributes

$attributes_x[0]$: the attribute of the current root node
 $attributes_x[i]$: the attribute of the i -th child ($i > 0$)

... the example, now in general formalization:

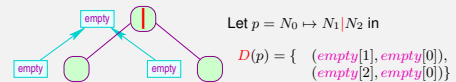
$$\begin{aligned} \boxed{x} &: empty[0] := (x \equiv \epsilon) \\ \boxed{|} &: empty[0] := empty[1] \vee empty[2] \\ \boxed{\cdot} &: empty[0] := empty[1] \wedge empty[2] \\ \boxed{*} &: empty[0] := t \\ \boxed{?} &: empty[0] := t \end{aligned}$$

8/67

Observations

- the *local* attribute equations need to be evaluated using a *global* algorithm that knows about the dependencies of the equations
- in order to construct this algorithm, we need
 - a sequence in which the nodes of the tree are visited
 - a sequence within each node in which the equations are evaluated
- this *evaluation strategy* has to be compatible with the *dependencies* between attributes

We visualize the attribute dependencies $D(p)$ of a production p in a *Local Dependency Graph*:



Let $p = N_0 \rightarrow N_1 | N_2$ in

$$D(p) = \{ (empty[1], empty[0]), (empty[2], empty[0]) \}$$

~> arrows point in the direction of information flow

9/67

Observations

- in order to infer an evaluation strategy, it is not enough to consider the *local* attribute dependencies at each node
- the evaluation strategy must also depend on the *global* dependencies, that is, on the information flow between nodes
- the global dependencies change with each particular syntax tree
- in the example, the parent node is always depending on children only
 ~> a depth-first post-order traversal is possible
- in general, variable dependencies can be much *more complex*

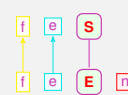
10/67

Simultaneous Computation of Multiple Attributes

Computing $empty$, $first$, $next$ from regular expressions:

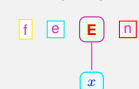
$$\boxed{S \rightarrow E}: \begin{aligned} empty[0] &:= empty[1] \\ first[0] &:= first[1] \\ next[1] &:= \emptyset \end{aligned} \quad \boxed{E \rightarrow x}: \begin{aligned} empty[0] &:= (x \equiv \epsilon) \\ first[0] &:= \{x \mid x \neq \epsilon\} \end{aligned}$$

$D(S \rightarrow E)$:



$$D(S \rightarrow E) = \{ (empty[1], empty[0]), (first[1], first[0]) \}$$

$D(E \rightarrow x)$:



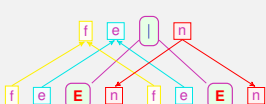
$$D(E \rightarrow x) = \{ \}$$

11/67

Regular Expressions: Rules for Alternative

$$\boxed{E \rightarrow E|E}: \begin{aligned} empty[0] &:= empty[1] \vee empty[2] \\ first[0] &:= first[1] \cup first[2] \\ next[1] &:= next[0] \\ next[2] &:= next[0] \end{aligned}$$

$D(E \rightarrow E|E)$:



$$D(E \rightarrow E|E) = \{ (empty[1], empty[0]), (empty[2], empty[0]), (first[1], first[0]), (first[2], first[0]), (next[0], next[2]), (next[0], next[1]) \}$$

12/67

Regular Expressions: Rules for Concatenation

$$\boxed{E \rightarrow E \cdot E}: \begin{aligned} empty[0] &:= empty[1] \wedge empty[2] \\ first[0] &:= first[1] \cup (empty[1]^? first[2] : \emptyset) \\ next[1] &:= first[2] \cup (empty[2]^? next[0] : \emptyset) \\ next[2] &:= next[0] \end{aligned}$$

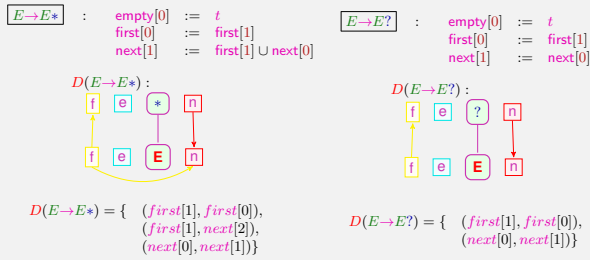
$D(E \rightarrow E \cdot E)$:



$$D(E \rightarrow E \cdot E) = \{ (empty[1], empty[0]), (empty[2], empty[0]), (empty[2], next[1]), (empty[1], first[0]), (first[1], first[0]), (first[2], first[0]), (first[2], next[1]), (next[0], next[2]), (next[0], next[1]) \}$$

13/67

Regular Expressions: Rules for Kleene-Star and Option



14/67

Challenges for General Attribute Systems

Static evaluation

Is there a static evaluation strategy, which is generally applicable?

- an evaluation strategy can only exist, if for *any* derivation tree the dependencies between attributes are *acyclic*
- it is *DEXPTIME*-complete to check for cyclic dependencies [Jazayeri, Odgen, Rounds, 1975]

Ideas

- Let the *User* specify the strategy
- Determine the strategy dynamically
- Automate *subclasses* only

15/67

Subclass: Strongly Acyclic Attribute Dependencies

Idea: For all nonterminals X compute a set $\mathcal{R}(X)$ of relations between its attributes, as an *overapproximation of the global dependencies* between root attributes of every production for X .

Describe $\mathcal{R}(X)$ s as sets of relations, similar to $D(p)$ by

- setting up each production $X \rightarrow X_1 \dots X_k$'s effect on the relations of $\mathcal{R}(X)$
- compute effect on all so far accumulated evaluations of each rhs X_i 's $\mathcal{R}(X_i)$
- iterate until stable

16/67

Subclass: Strongly Acyclic Attribute Dependencies

The 2-ary operator $L[i]$ re-decorates relations from L

$$L[i] = \{ (a[i], b[i]) \mid (a, b) \in L \}$$

π_0 projects only onto relations between root elements only

$$\pi_0(S) = \{ (a, b) \mid (a[0], b[0]) \in S \}$$

$[p]^+$... root-projects the transitive closure of relations from the L_i s and D

$$[p]^+(L_1, \dots, L_k) = \pi_0((D(p) \cup L_1[1] \cup \dots \cup L_k[k])^+)$$

\mathcal{R} maps symbols to relations (global attributes dependencies)

$$\mathcal{R}(X) \supseteq (\bigcup \{ [p]^+(\mathcal{R}(X_1), \dots, \mathcal{R}(X_k)) \mid p : X \rightarrow X_1 \dots X_k \}^+ \mid p \in P)$$

$$\mathcal{R}(X) \supseteq \emptyset \quad | X \in (N \cup T)$$

Strongly Acyclic Grammars

The system of inequalities $\mathcal{R}(X)$

- characterizes the class of strongly acyclic Dependencies
- has a unique least solution $\mathcal{R}^*(X)$ (as $[p]^+$ is monotonic)

17/67

Subclass: Strongly Acyclic Attribute Dependencies

Strongly Acyclic Grammars

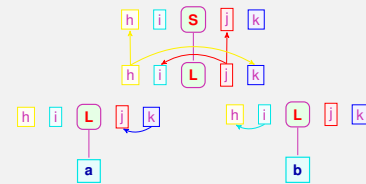
If all $D(p) \cup \mathcal{R}^*(X_1)[1] \cup \dots \cup \mathcal{R}^*(X_k)[k]$ are acyclic for all $p \in G$, G is strongly acyclic.

Idea: we compute the least solution $\mathcal{R}^*(X)$ of $\mathcal{R}(X)$ by a fixpoint computation, starting from $\mathcal{R}(X) = \emptyset$.

18/67

Example: Strong Acyclic Test

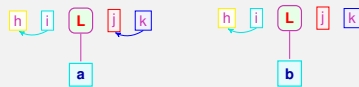
Given grammar $S \rightarrow L, L \rightarrow a \mid b$. Dependency graphs D_p :



19/67

Example: Strong Acyclic Test

Start with computing $\mathcal{R}(L) = [L \rightarrow a]^+(\emptyset) \cup [L \rightarrow b]^+(\emptyset)$:

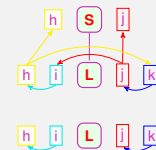


- terminal symbols do not contribute dependencies **check for cycles!**
- transitive closure of all relations in $(D(L \rightarrow a))^+$ and $(D(L \rightarrow b))^+$
- apply π_0
- $\mathcal{R}(L) = \{ (k, j), (i, h) \}$

20/67

Example: Strong Acyclic Test

Continue with $\mathcal{R}(S) = [S \rightarrow L]^+(\mathcal{R}(L))$:

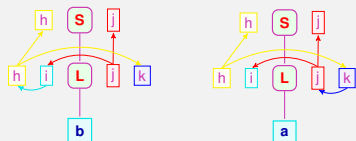


- re-decorate and embed $\mathcal{R}(L)[1]$ **check for cycles!**
- transitive closure of all relations $(D(S \rightarrow L) \cup \{ (k[1], j[1]) \} \cup \{ (i[1], h[1]) \})^+$
- apply π_0
- $\mathcal{R}(S) = \{ \}$

21/67

Strong Acyclic and Acyclic

The grammar $S \rightarrow L, L \rightarrow a \mid b$ has only two derivation trees which are both *acyclic*:



It is *not strongly acyclic* since the over-approximated global dependence graph for the non-terminal L contributes to a cycle when computing $\mathcal{R}(S)$:



22/67

From Dependencies to Evaluation Strategies

Possible strategies:

- let the *user* define the evaluation order
- automatic* strategy based on the dependencies
- consider a *fixed* strategy and only allow an attribute system that can be evaluated using this strategy

23/67

Linear Order from Dependency Partial Order

Possible *automatic* strategies:

- 1 **demand-driven evaluation**
 - start with the evaluation of any required attribute
 - if the equation for this attribute relies on as-of-yet unevaluated attributes, evaluate these recursively
- 2 **evaluation in passes**
 - for each pass, pre-compute a **global strategy** to visit the **nodes** together with a **local strategy** for evaluation **within each node** type
 - ~ **minimize** the number of **visits** to each node

24/67

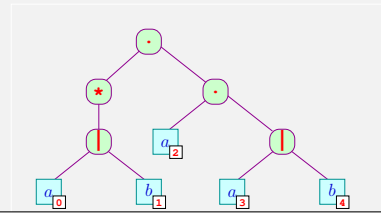
Example: Demand-Driven Evaluation

Compute **next** at leaves a_2, a_3 and b_4 in the expression $(a|b)^*a(a|b)$:

```

[] : next[1] := next[0]
   next[2] := next[0]

[] : next[1] := first[2] ∪ (empty[2]? next[0]: ∅)
   next[2] := next[0]
    
```



25/67

Demand-Driven Evaluation

Observations

- each node must contain a pointer to its parent
- only required** attributes are evaluated
- the evaluation sequence depends – in general – on the actual syntax tree
- the algorithm must track which attributes it has already evaluated
- the algorithm may visit nodes more often than necessary
- ~ the algorithm is **not local**

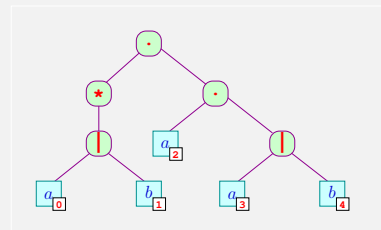
in principle:

- evaluation strategy is dynamic: difficult to debug
- usually all attributes in all nodes are required
- ~ computation of all attributes is often cheaper
- ~ perform evaluation in **passes**

26/67

Implementing State

Problem: In many cases some sort of state is required.
Example: numbering the leaves of a syntax tree



27/67

Example: Implementing Numbering of Leaves

Idea:

- use **helper** attributes **pre** and **post**
- in **pre** we pass the value for the first leaf down (inherited attribute)
- in **post** we pass the value of the last leaf up (synthesized attribute)

```

root: pre[0] := 0
      pre[1] := pre[0]
      post[0] := post[1]
    
```

```

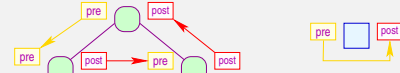
node: pre[1] := pre[0]
      pre[2] := post[1]
      post[0] := post[2]
    
```

```

leaf: post[0] := pre[0] + 1
    
```

28/67

L-Attribution



- the attribute system is apparently strongly acyclic
- each node computes
 - the inherited attributes before descending into a child node (corresponding to a pre-order traversal)
 - the synthesized attributes after returning from a child node (corresponding to post-order traversal)

Definition L-Attributed Grammars

An attribute system is *L-attributed*, if for all productions $S \rightarrow S_1 \dots S_n$ every inherited attribute of S_j where $1 \leq j \leq n$ only depends on

- the attributes of S_1, S_2, \dots, S_{j-1} and
- the inherited attributes of S .

29/67

L-Attribution

Background:

- the attributes of an *L-attributed* grammar can be evaluated during parsing
- important if no syntax tree is required or if error messages should be emitted while parsing
- example: pocket calculator

L-attributed grammars have a fixed evaluation strategy:

a single **depth-first** traversal

- in general: partition all attributes into $A = A_1 \cup \dots \cup A_n$ such that for all attributes in A_i the attribute system is *L-attributed*

- perform a **depth-first** traversal for each attribute set A_i

~ craft attribute system in a way that they can be partitioned into few *L-attributed* sets

30/67

Practical Applications

- symbol tables, type checking/inference, and simple code generation can all be specified using *L-attributed* grammars
- most applications **annotate** syntax trees with additional information
- the nodes in a syntax tree usually have different **types** that depend on the non-terminal that the node represents
- ~ the different types of non-terminals are characterized by the set of attributes with which they are decorated

Example: Def-Use Analysis

- a **statement** may have two attributes containing valid identifiers: one ingoing (inherited) set and one outgoing (synthesised) set
- an **expression** only has an ingoing set

31/67

Implementation of Attribute Systems via a Visitor

- class with a method for every non-terminal in the grammar

```

public abstract class Regex {
    public abstract void accept(Visitor v);
}
    
```

- attribute-evaluation works via **pre-order / post-order callbacks**

```

public interface Visitor {
    default void pre(OrEx re) {}
    default void pre(AndEx re) {}
    ...
    default void post(OrEx re) {}
    default void post(AndEx re) {}
}
    
```

- we pre-define a depth-first traversal of the syntax tree

```

public class OrEx extends Regex {
    Regex l, r;
    public void accept(Visitor v) {
        v.pre(this); l.accept(v); v.inter(this);
        r.accept(v); v.post(this);
    }
}
    
```

32/67

Example: Leaf Numbering

```

public abstract class AbstractVisitor implements Visitor {
    public void pre(OrEx re) { pr(re); }
    public void pre(AndEx re) { pr(re); }
    ... /* redirecting to default handler for bin exprs */
    public void post(OrEx re) { po(re); }
    public void post(AndEx re) { po(re); }
    abstract void po(BinEx re);
    abstract void in(BinEx re);
    abstract void pr(BinEx re);
}

public class LeafNum extends AbstractVisitor {
    public Map<Regex, Integer> pre = new HashMap<>();
    public Map<Regex, Integer> post = new HashMap<>();
    public LeafNum(Regex r) { pre.put(r, 0); r.accept(this); }
    public void pre(Const r) { post.put(r, pre.get(r)+1); }
    public void pr(BinEx r) { pre.put(r.l, pre.get(r)); }
    public void in(BinEx r) { pre.put(r.r, post.get(r.l)); }
    public void po(BinEx r) { post.put(r, post.get(r.r)); }
}
    
```

33/67

Chapter 2: Decl-Use Analysis

Symbol Bindings and Visibility

Consider the following Java code:

```
void foo() {
    int a;
    while (true) {
        double a;
        a = 0.5;
        write(a);
        break;
    }
    a = 2;
    bar();
    write(a);
}
```

- each **declaration** of a variable v causes memory allocation for v
- using v requires knowledge about its memory location
→ determine the declaration v is **bound** to
- a binding is not **visible** when a local declaration of the same name is in scope
in the example the declaration of a is shadowed by the **local declaration** in the loop body

Scope of Identifiers

```
void foo() {
    int a;
    while (true) {
        double a;
        a = 0.5;
        write(a);
        break;
    }
    a = 2;
    bar();
    write(a);
}
```

scope of int a } scope of double a

⚠ administration of identifiers can be quite complicated...

Resolving Identifiers

Observation: each identifier in the AST must be translated into a memory access

Problem: for each identifier, find out what memory needs to be accessed by providing **rapid access** to its **declaration**

Ideas:

- rapid access:** replace every identifier by a **unique** integer
→ integers as keys: comparisons of integers is faster
- link each usage of a variable to the **declaration** of that variable
→ for languages without explicit declarations, create declarations when a variable is first encountered

Rapid Access: Replace Strings with Integers

Idea for Algorithm:

Input: a sequence of strings

Output: sequence of numbers

- table that allows to retrieve the string that corresponds to a number

Apply this algorithm on each identifier during **scanning**.

Implementation approach:

- count the number of new-found identifiers in **int count**
- maintain a **hashtable** $S : \text{String} \rightarrow \text{int}$ to remember numbers for known identifiers

We thus define the function:

```
int indexForIdentifier(String w) {
    if (S(w) ≡ undefined) {
        S = S ⊕ {w ↦ count};
        return count++;
    } else return S(w);
}
```

Implementation: Hashtables for Strings

- allocate an array M of sufficient size m
- choose a **hash function** $H : \text{String} \rightarrow [0, m-1]$ with:
 - $H(w)$ is **cheap** to compute
 - H distributes the occurring words **equally** over $[0, m-1]$

Possible generic choices for sequence types ($\vec{x} = \langle x_0, \dots, x_{r-1} \rangle$):

$$H_0(\vec{x}) = (x_0 + x_{r-1}) \% m$$

$$H_1(\vec{x}) = (\sum_{i=0}^{r-1} x_i \cdot p^i) \% m$$

$$H_2(\vec{x}) = (x_0 + p \cdot (x_1 + p \cdot (\dots + p \cdot x_{r-1} \dots))) \% m$$

for some prime number p (e.g. 31)

✗ The hash value of w **may not be unique!**

- Append (w, i) to a linked list located at $M[H(w)]$
- Finding the index for w , we compare w with all x for which $H(w) = H(x)$

✓ access on average:

- insert: $O(1)$
- lookup: $O(1)$

Example: Replacing Strings with Integers

Input:

```
Peter Piper picked a peck of pickled peppers
If Peter Piper picked a peck of pickled peppers
wheres the peck of pickled peppers Peter Piper picked
```

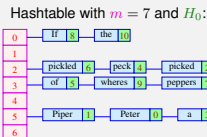
Output:

```
0 1 2 3 4 5 6 7 8 0 1 2 3 4 5 6
7 9 10 4 5 6 7 0 1 2
```

and

0	Peter
1	Piper
2	picked
3	a
4	peck
5	of

6	pickled
7	peppers
8	If
9	wheres
10	the



Refer Uses to Declarations: Symbol Tables

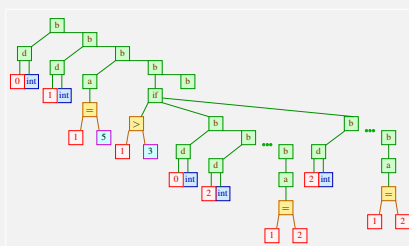
Check for the correct usage of variables:

- Traverse the syntax tree in a suitable sequence, such that
 - each declaration is visited **before** its use
 - the currently visible declaration is the last one visited
→ perfect for an L-attributed grammar
 - equation system for basic block must add and remove identifiers
- for each identifier, we manage a **stack** of declarations
 - if we visit a **declaration**, we push it onto the stack of its identifier
 - upon leaving the **scope**, we remove it from the stack
- if we visit a **usage** of an identifier, we pick the top-most declaration from its stack
- if the stack of the identifier is empty, we have found an undeclared identifier

Example: Decl-Use Analysis via Table of Stacks

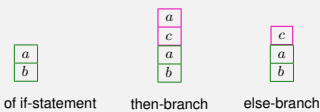
```
1 void f()
2 {
3   int a, b;
4   b = 5;
5   if (b > 3) {
6     int a, c;
7     a = 3;
8     c = a + 1;
9     b = c;
10  } else {
11   int c;
12   c = a + 1;
13   b = c;
14  }
15  b = a + b;
16 }
```

- d declaration
- b basic block
- a assignment



Alternative Implementations for Symbol Tables

- when using a list to store the symbol table, storing a marker indicating the old head of the list is sufficient



- instead of lists of symbols, it is possible to use a list of hash tables → more efficient in large, shallow programs
- an even more elegant solution: **persistent trees** (updates return fresh trees with references to the old tree where possible)
→ a persistent tree t can be passed down into a basic block where new elements may be added, yielding a t' ; after examining the basic block, the analysis proceeds with the unchanged old t

Chapter 3: Type Checking

Goal of Type Checking

In most mainstream (imperative / object oriented / functional) programming languages, variables and functions have a fixed **type**.
for example: `int, void*, struct { int x; int y; }`.

- Types are useful to
- manage **memory**
 - select correct **assembler instructions**
 - to avoid certain **run-time errors**

In imperative and object-oriented programming languages a declaration has to specify a type. The compiler then checks for a type correct use of the declared entity.

Type Expressions

Types are given using type-**expressions**.

The set of type expressions T contains:

- **base types:** `int, char, float, void, ...`
- **type constructors** that can be applied to other types

example for type constructors in C:

- **structures:** `struct { t1 a1; ... tk ak; }`
- **pointers:** `t *`
- **arrays:** `t [n]`
 - the size of an array can be specified
 - the variable to be declared is written between t and $[n]$
- **functions:** `t (t1, ..., tk)`
 - the variable to be declared is written between t and (t_1, \dots, t_k)
 - in ML function types are written as: `t1 * ... * tk → t`

Type Checking

Problem:

Given: A set of type declarations $\Gamma = \{t_1 x_1; \dots t_m x_m\}$
Check: Can an expression e be given the type t ?

Example:

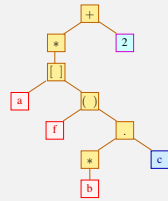
```
struct list { int info; struct list* next; };
int f(struct list* l) { return l; };
struct { struct list* c; } * b;
int* a[11];
```

Consider the expression:

```
*a[f(b->c)]+2;
```

Type Checking using the Syntax Tree

Check the expression `*a[f(b->c)]+2;`:



Idea:

- traverse the syntax tree **bottom-up**
- for each identifier, we lookup its type in Γ
- constants such as 2 or 0.5 have a fixed type
- the types of the inner nodes of the tree are deduced using **typing rules**

Type Systems for C-like Languages

Formally: consider **judgements** of the form:

$$\Gamma \vdash e : t$$

// (in the type environment Γ the expression e has type t)

Axioms:

Const: $\Gamma \vdash c : t_c$ (t_c type of constant c)
Var: $\Gamma \vdash x : \Gamma(x)$ (x Variable)

Rules:

$$\text{Ref: } \frac{\Gamma \vdash e : t}{\Gamma \vdash *e : t*} \quad \text{Deref: } \frac{\Gamma \vdash e : t*}{\Gamma \vdash *e : t}$$

Type Systems for C-like Languages

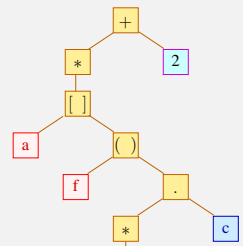
More rules for typing an expression: with subtyping relation \leq

- Array: $\frac{\Gamma \vdash e_1 : t* \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1[e_2] : t}$
- Array: $\frac{\Gamma \vdash e_1 : t[] \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1[e_2] : t}$
- Struct: $\frac{\Gamma \vdash e : \text{struct } \{t_1 a_1; \dots t_m a_m\}}{\Gamma \vdash e.a_i : t_i}$
- App: $\frac{\Gamma \vdash e : t(t_1, \dots, t_m) \quad \Gamma \vdash e_1 : t_1 \dots \Gamma \vdash e_m : t_m}{\Gamma \vdash e(e_1, \dots, e_m) : t}$
- Op \square : $\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash e_1 \square e_2 : t_1 \square t_2}$
- Op $=$: $\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2 \quad t_2 \text{ can be converted to } \leq t_1}{\Gamma \vdash e_1 = e_2 : t_1}$
- Explicit Cast: $\frac{\Gamma \vdash e : t_2 \quad t_2 \text{ can be converted to } \leq t_1}{\Gamma \vdash (t_1) e : t_1}$

Example: Type Checking

Given expression `*a[f(b->c)]+2` and

```
\Gamma = {
struct list { int info; struct list* next; };
int f(struct list* l);
struct { struct list* c; } * b;
int* a[11];
}
```



Example: Type Checking – More formally:

```
\Gamma = {
struct list { int info; struct list* next; };
int f(struct list* l);
struct { struct list* c; } * b;
int* a[11];
}
```

$$\begin{array}{c} \text{VAR } \frac{}{\Gamma \vdash b : \text{struct}\{\text{struct list } *c\}*} \\ \text{DEREF } \frac{}{\Gamma \vdash *b : \text{struct}\{\text{struct list } *c\}} \\ \text{STRUCT } \frac{}{\Gamma \vdash (*b).c : \text{struct list}*} \\ \\ \text{VAR } \frac{}{\Gamma \vdash a : \text{int}*[11]} \quad \text{APP } \frac{\text{VAR } \frac{}{\Gamma \vdash f : \text{int}(\text{struct list})*} \quad \Gamma \vdash (*b).c : \text{struct list}*}{\Gamma \vdash f(b \rightarrow c) : \text{int}*} \\ \text{ARRAY } \frac{}{\Gamma \vdash a[f(b \rightarrow c)] : \text{int}*} \\ \\ \text{DEREF } \frac{}{\Gamma \vdash *a[f(b \rightarrow c)] : \text{int}} \quad \text{CONST } \frac{}{\Gamma \vdash 2 : \text{int}} \\ \text{OP } \frac{}{\Gamma \vdash *a[f(b \rightarrow c)] + 2 : \text{int}} \end{array}$$

but what do we do with \leq ?

Equality of Types =

Summary of Type Checking

- Choosing which rule to apply at an AST node is determined by the type of the child nodes
- determining the rule requires a check for \rightsquigarrow **equality** of types

type equality in C:

- `struct A { }` and `struct B { }` are considered to be different
 - \rightsquigarrow the compiler could re-order the fields of A and B independently (**not** allowed in C)
 - to extend an record A with more fields, it has to be embedded into another record:

```
struct B {
struct A;
int field_of_B;
} extension_of_A;
```
- after issuing `typedef int C;` the types C and `int` are **the same**

Structural Type Equality

Alternative interpretation of type equality (*does not hold in C*):

semantically, two types t_1, t_2 can be considered as *equal* if they accept the same set of access paths.

Example:

```

struct list {
  int info;
  struct list* next;
}
struct list1 {
  int info;
  struct {
    int info;
    struct list1* next;
  } * next;
}
    
```

Consider declarations `struct list* l` and `struct list1* l`. Both allow `l->info` `l->next->info`

but the two declarations of `l` have unequal types in C.

54/67

Algorithm for Testing Structural Equality

Idea:

- track a set of equivalence queries of type expressions
- if two types are *syntactically* equal, we stop and report success
- otherwise, reduce the equivalence query to a several equivalence queries on (hopefully) *simpler* type expressions

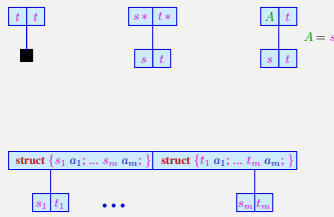
Suppose that recursive types were introduced using type definitions:

```
typedef A t
```

(we omit the Γ). Then define the following rules:

55/67

Rules for Well-Typedness



56/67

Example:

```

typedef struct {int info; A * next;} A
typedef struct {int info; struct {int info; B * next;} * next;} B
    
```

We ask, for instance, if the following equality holds:

```
struct {int info; A * next;} = B
```

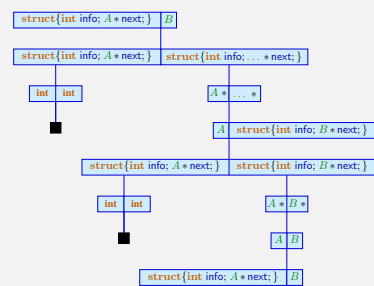
We construct the following deduction tree:

57/67

Proof for the Example:

```

typedef struct {int info; A * next;} A
typedef struct {int info; struct {int info; B * next;} * next;} B
    
```



58/67

Implementation

We implement a function that implements the equivalence query for two types by applying the deduction rules:

- if no deduction rule applies, then the two types are *not equal*
- if the deduction rule for expanding a type definition applies, the function is called recursively with a *potentially larger* type
- in case an equivalence query appears a second time, the types are *equal by definition*

Termination

- the set D of all declared types is finite
- there are no more than $|D|^2$ different equivalence queries
- repeated queries for the same inputs are automatically satisfied

~> termination is ensured

59/67

Subtyping \leq

On the arithmetic basic types `char`, `int`, `long`, etc. there exists a rich *subtype* hierarchy

Subtypes

$t_1 \leq t_2$, means that the values of type t_1

- form a *subset* of the values of type t_2 ;
- can be converted into a value of type t_2 ;
- fulfill the requirements of type t_2 ;
- are assignable to variables of type t_2 .

Example:

assign smaller type (fewer values) to larger type (more values)

```

t1 int x;
t2 double y;
y = x;
t1 <= t2 int <= double
    
```

60/67

Example: Subtyping

Extending the subtype relationship to more complex types, observe:

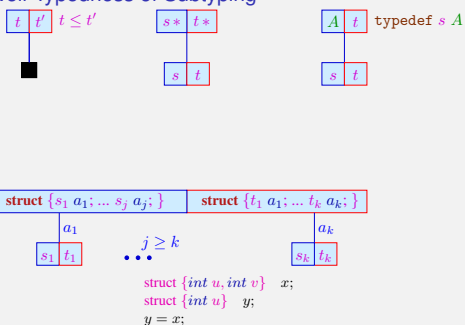
```

string extractInfo( struct { string info; } x ) {
  return x.info;
}
    
```

- we want `extractInfo` to be applicable to all argument structures that return a `string` typed field for accessor `info`
- the idea of subtyping on values is related to subclasses
- we use deduction rules to describe when $t_1 \leq t_2$ should hold...

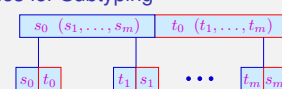
61/67

Rules for Well-Typedness of Subtyping



62/67

Rules and Examples for Subtyping



Examples:

```

struct {int a; int b;} <= struct {float a;}
int (int) <= float (float)
int (float) <= float (int)
    
```

Definition

Given two function types in subtype relation $s_0(s_1, \dots, s_n) \leq t_0(t_1, \dots, t_n)$ then we have

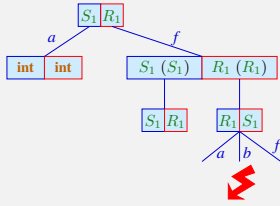
- *co-variance* of the return type $s_0 \leq t_0$ and
- *contra-variance* of the arguments $s_i \geq t_i$ für $1 < i \leq n$

63/67

Subtypes: Application of Rules (I)

Check if $S_1 \leq R_1$:

$R_1 = \text{struct } \{ \text{int } a; R_1(R_1) f; \}$
 $S_1 = \text{struct } \{ \text{int } a; \text{int } b; S_1(S_1) f; \}$
 $R_2 = \text{struct } \{ \text{int } a; R_2(S_2) f; \}$
 $S_2 = \text{struct } \{ \text{int } a; \text{int } b; S_2(R_2) f; \}$

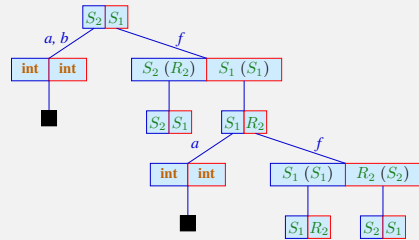


64/67

Subtypes: Application of Rules (II)

Check if $S_2 \leq S_1$:

$R_1 = \text{struct } \{ \text{int } a; R_1(R_1) f; \}$
 $S_1 = \text{struct } \{ \text{int } a; \text{int } b; S_1(S_1) f; \}$
 $R_2 = \text{struct } \{ \text{int } a; R_2(S_2) f; \}$
 $S_2 = \text{struct } \{ \text{int } a; \text{int } b; S_2(R_2) f; \}$

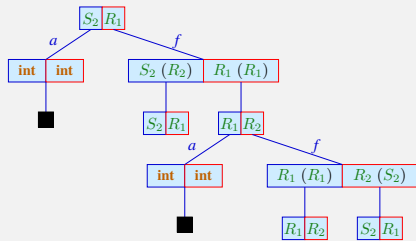


65/67

Subtypes: Application of Rules (III)

Check if $S_2 \leq R_1$:

$R_1 = \text{struct } \{ \text{int } a; R_1(R_1) f; \}$
 $S_1 = \text{struct } \{ \text{int } a; \text{int } b; S_1(S_1) f; \}$
 $R_2 = \text{struct } \{ \text{int } a; R_2(S_2) f; \}$
 $S_2 = \text{struct } \{ \text{int } a; \text{int } b; S_2(R_2) f; \}$



66/67

Discussion

- for presentational purposes, proof trees are often abbreviated by omitting deductions within the tree
- structural sub-types are very powerful and can be quite intricate to understand
- Java generalizes structs to **objects/classes** where a sub-class A inheriting from base class O is a subtype $A \leq O$
- subtype relations between classes must be **explicitly declared**

67/67

Topic: Code Synthesis

1/49

Generating Code: Overview

We inductively generate instructions from the AST:

- there is a rule stating how to generate code for each non-terminal of the grammar
- the code is merely another attribute in the syntax tree
- code generation makes use of the already computed attributes

In order to specify the code generation, we require

- a semantics of the language we are compiling (here: C standard)
- a semantics of the machine instructions

~> we commence by specifying machine instruction semantics

2/49

Code Synthesis

Chapter 1: The Register C-Machine

3/49

The Register C-Machine (R-CMa)

We generate Code for the Register C-Machine.

The **Register C-Machine** is a virtual machine (VM).

- there exists no processor that can execute its instructions
- ... but we can build an interpreter for it
- we provide a visualization environment for the R-CMa
- the R-CMa has no **double, float, char, short** or **long** types
- the R-CMa has no instructions to communicate with the operating system
- the R-CMa has an unlimited supply of registers

The R-CMa is more realistic than it may seem:

- the mentioned restrictions can easily be lifted
- the **Dalvik VM/ART** or the **LLVM** are similar to the R-CMa
- an interpreter of R-CMa can run on any platform

4/49

Virtual Machines

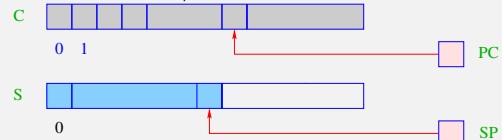
A virtual machine has the following ingredients:

- any virtual machine provides a set of **instructions**
- instructions are executed on virtual hardware
- the virtual hardware is a collection of data structures that is accessed and modified by the VM instructions
- ... and also by other components of the **run-time system**, namely functions that go beyond the instruction semantics
- the **interpreter** is part of the run-time system

5/49

Components of a Virtual Machine

Consider **Java** as an example:



A virtual machine such as the **Dalvik VM** has the following structure:

- **S**: the data store – a memory region in which cells can be stored in LIFO order ~> **stack**.
- **SP**: ($\hat{=}$ **stack pointer**) pointer to the last used cell in **S**
- beyond **S** follows the memory containing the **heap**
- **C** is the memory storing **code**
 - each cell of **C** holds exactly one virtual instruction
 - **C** can only be **read**
- **PC** ($\hat{=}$ **program counter**) address of the instruction that is to be executed next
- **PC** contains 0 initially

6/49

Executing a Program

- the machine loads an instruction from $C[PC]$ into the **instruction register IR** in order to execute it
- before evaluating the instruction, the **PC** is incremented by one


```
while (true) {
  IR = C[PC]; PC++;
  execute (IR);
}
```
- note: the **PC** must be incremented **before** the execution, since an instruction may modify the **PC**
- the loop is exited by evaluating a **halt** instruction that returns directly to the operating system

7/49

Chapter 2: Generating Code for the Register C-Machine

8/49

Simple Expressions and Assignments in R-CMA

- Task:** evaluate the expression $(1 + 7) * 3$
that is, generate an instruction sequence that
- computes the value of the expression and
 - keeps its value accessible in a reproducible way

Idea:

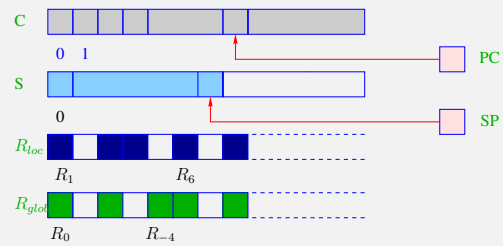
- first compute the value of the sub-expressions
- store the intermediate result in a temporary register
- apply the operator
- loop

9/49

Principles of the R-CMA

The **R-CMA** is composed of a stack, heap and a code segment, just like the **JVM**; it additionally has register sets:

- local** registers are $R_1, R_2, \dots, R_i, \dots$
- global** registers are $R_0, R_{-1}, \dots, R_j, \dots$



10/49

The Register Sets of the R-CMA

The two register sets have the following purpose:

- the **local** registers R_i
 - save temporary results
 - store the contents of local variables of a function
 - can efficiently be stored and restored from the stack
- the **global** registers R_i
 - save the parameters of a function
 - store the result of a function

Note:

for now, we only use registers to store temporary computations

Idea for the translation: use a register counter i :

- registers R_j with $j < i$ are **in use**
- registers R_j with $j \geq i$ are **available**

11/49

Translation of Simple Expressions

Using variables stored in registers; loading constants:

instruction	semantics	intuition
<code>loadc R_i c</code>	$R_i = c$	load constant
<code>move R_i R_j</code>	$R_i = R_j$	copy R_j to R_i

We define the following translation schema (with $\rho x = a$):

$\text{code}_R^i c \rho$	=	<code>loadc R_i c</code>
$\text{code}_R^i x \rho$	=	<code>move R_i R_a</code>
$\text{code}_R^i x = e \rho$	=	<code>code_R^i e \rho</code> <code>move R_a R_i</code>

12/49

Translation of Expressions

Let $\text{op} = \{\text{add}, \text{sub}, \text{div}, \text{mul}, \text{mod}, \text{le}, \text{gr}, \text{eq}, \text{leq}, \text{geq}, \text{and}, \text{or}\}$. The **R-CMA** provides an instruction for each operator op .

$$\text{op } R_i R_j R_k$$

where R_i is the target register, R_j the first and R_k the second argument.

Correspondingly, we generate code as follows:

$$\text{code}_R^i e_1 \text{op } e_2 \rho = \begin{matrix} \text{code}_R^i e_1 \rho \\ \text{code}_R^{i+1} e_2 \rho \\ \text{op } R_i R_i R_{i+1} \end{matrix}$$

Example: Translate $3 * 4$ with $i = 4$:

$$\begin{aligned} \text{code}_R^4 3 * 4 \rho &= \begin{matrix} \text{code}_R^4 3 \rho \\ \text{code}_R^5 4 \rho \\ \text{mul } R_4 R_4 R_5 \end{matrix} \\ \text{code}_R^4 3 * 4 \rho &= \begin{matrix} \text{loadc } R_4 3 \\ \text{loadc } R_5 4 \\ \text{mul } R_4 R_4 R_5 \end{matrix} \end{aligned}$$

13/49

Managing Temporary Registers

Observe that temporary registers are re-used: translate $3 * 4 + 3 * 4$ with $t = 4$:

$$\text{code}_R^4 3 * 4 + 3 * 4 \rho = \begin{matrix} \text{code}_R^4 3 * 4 \rho \\ \text{code}_R^5 3 * 4 \rho \\ \text{add } R_4 R_4 R_5 \end{matrix}$$

where

$$\text{code}_R^i 3 * 4 \rho = \begin{matrix} \text{loadc } R_4 3 \\ \text{loadc } R_{i+1} 4 \\ \text{mul } R_i R_i R_{i+1} \end{matrix}$$

we obtain

$$\text{code}_R^4 3 * 4 + 3 * 4 \rho = \begin{matrix} \text{loadc } R_4 3 \\ \text{loadc } R_5 4 \\ \text{mul } R_4 R_4 R_5 \\ \text{loadc } R_5 3 \\ \text{loadc } R_6 4 \\ \text{mul } R_5 R_5 R_6 \\ \text{add } R_4 R_4 R_5 \end{matrix}$$

14/49

Semantics of Operators

The operators have the following semantics:

<code>add R_i R_j R_k</code>	$R_i = R_j + R_k$
<code>sub R_i R_j R_k</code>	$R_i = R_j - R_k$
<code>div R_i R_j R_k</code>	$R_i = R_j / R_k$
<code>mul R_i R_j R_k</code>	$R_i = R_j * R_k$
<code>mod R_i R_j R_k</code>	$R_i = \text{signum}(R_k) \cdot k$ with $ R_j = n \cdot R_k + k \wedge n \geq 0, 0 \leq k < R_k $
<code>le R_i R_j R_k</code>	$R_i = \text{if } R_j < R_k \text{ then } 1 \text{ else } 0$
<code>gr R_i R_j R_k</code>	$R_i = \text{if } R_j > R_k \text{ then } 1 \text{ else } 0$
<code>eq R_i R_j R_k</code>	$R_i = \text{if } R_j = R_k \text{ then } 1 \text{ else } 0$
<code>leq R_i R_j R_k</code>	$R_i = \text{if } R_j \leq R_k \text{ then } 1 \text{ else } 0$
<code>geq R_i R_j R_k</code>	$R_i = \text{if } R_j \geq R_k \text{ then } 1 \text{ else } 0$
<code>and R_i R_j R_k</code>	$R_i = R_j \& R_k$ // bit-wise and
<code>or R_i R_j R_k</code>	$R_i = R_j R_k$ // bit-wise or

Note: all registers and memory cells contain operands in Z

15/49

Translation of Unary Operators

Unary operators $\text{op} = \{\text{neg}, \text{not}\}$ take only two registers:

$$\text{code}_R^i \text{op } e \rho = \begin{matrix} \text{code}_R^i e \rho \\ \text{op } R_i R_i \end{matrix}$$

Note: We use the same register.

Example: Translate -4 into R_5 :

$$\begin{aligned} \text{code}_R^5 -4 \rho &= \text{code}_R^5 4 \rho \\ \text{code}_R^5 -4 \rho &= \begin{matrix} \text{loadc } R_5 4 \\ \text{neg } R_5 R_5 \end{matrix} \end{aligned}$$

The operators have the following semantics:

<code>not R_i R_j</code>	$R_i \leftarrow \text{if } R_j = 0 \text{ then } 1 \text{ else } 0$
<code>neg R_i R_j</code>	$R_i \leftarrow -R_j$

16/49

Applying Translation Schema for Expressions

Suppose the following function `void f(void) { int x, y, z; x = y+z+3; }` is given:

- Let $\rho = \{x \mapsto 1, y \mapsto 2, z \mapsto 3\}$ be the address environment.
- Let R_1 be the first free register, that is, $i = 4$.

```
code4 x=y+z+3 ρ = code4 y+z+3 ρ
                  move R1 R4
code4 y+z+3 ρ = move R4 R2
                 code5 z+3 ρ
                 add R4 R4 R5
code5 z+3 ρ = move R5 R3
               code6 3 ρ
               mul R5 R5 R6
code6 3 ρ = loadc R6 3
```

→ the assignment `x=y+z+3` is translated as
`move R4 R2; move R5 R3; loadc R6 3; mul R5 R5 R6; add R4 R4 R5; move R1 R4`

Chapter 3: Statements and Control Structures

About Statements and Expressions

General idea for translation: `codei s ρ` : generate code for statement s
`codei e ρ` : generate code for expression e into R_i
 Throughout: $i, i+1, \dots$ are free (unused) registers

For an **expression** $x = e$ with ρ $x = a$ we defined:

```
codei x = e ρ = codei e ρ
               move Ra Ri
```

However, $x = e$; is also an **expression statement**:

- Define:

```
codei e1 = e2; ρ = codei e1 = e2 ρ
```

The temporary register R_i is ignored here. More general:

```
codei e; ρ = codei e ρ
```

- Observation:** the assignment to e_1 is a side effect of the evaluating the expression $e_1 = e_2$.

Translation of Statement Sequences

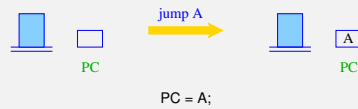
The code for a sequence of statements is the concatenation of the instructions for each statement in that sequence:

```
codei (s ss) ρ = codei s ρ
                  codei ss ρ
codei e ρ = // empty sequence of instructions
```

Note here: s is a statement, ss is a sequence of statements

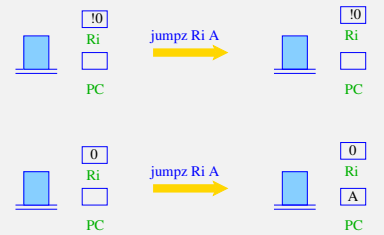
Jumps

In order to diverge from the linear sequence of execution, we need **jumps**:



Conditional Jumps

A conditional jump branches depending on the value in R_i :



if ($R_i == 0$) PC = A;

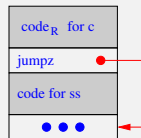
Simple Conditional

We first consider $s \equiv \text{if } (c) \text{ ss}$.
 ...and present a translation without basic blocks.

Idea:

- emit the code of c and ss in sequence
- insert a jump instruction in-between, so that correct control flow is ensured

```
codei s ρ = codei c ρ
            jumpz Ri A
            codei ss ρ
A : ...
```

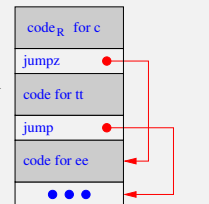


General Conditional



Translation of `if (c) tt else ee`.

```
codei if(c) tt else ee ρ =
codei c ρ
jumpz Ri A
codei tt ρ
jump B
A : codei ee ρ
B :
```



Example for if-statement

Let $\rho = \{x \mapsto 4, y \mapsto 7\}$ and let s be the statement

```
if (x>y) { /* (i) */
x = x - y; /* (ii) */
} else { /* (iii) */
y = y - x; /* (iii) */
}
```

Then `codei s ρ` yields:

```
(i) move Ri R4
     move Ri+1 R7
     gr Ri Ri Ri+1
     jumpz Ri A

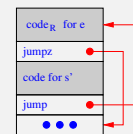
(ii) move Ri R4
     move Ri+1 R7
     sub Ri Ri Ri+1
     move R4 Ri
     jump B

(iii) move Ri R4
       move Ri+1 R4
       sub Ri Ri Ri+1
       move R7 Ri
       jump B
```

Iterating Statements

We only consider the loop $s \equiv \text{while}(e) s'$. For this statement we define:

```
codei while(e) s ρ = A : codei e ρ
                        jumpz Ri B
                        codei s ρ
                        jump A
B :
```



Example: Translation of Loops

Let $\rho = \{a \mapsto 7, b \mapsto 8, c \mapsto 9\}$ and let s be the statement:

```
while (a>0) { /* (i) */
  c = c + 1; /* (ii) */
  a = a - b; /* (iii) */
}
```

Then $code^i s \rho$ evaluates to:

```
(i)      (ii)      (iii)
A:  move Ri R7      move Ri R8      move Ri R7
    loadc Ri+1 0      loadc Ri+1 1      move Ri+1 R8
    gr Ri Ri Ri+1      add Ri Ri Ri+1      sub Ri Ri Ri+1
    jumpz Ri B      move R9 Ri      move R7 Ri
                                     jump A
                                     B:
```

for-Loops

The for-loop $s \equiv \text{for}(e_1; e_2; e_3) s'$ is equivalent to the statement sequence $e_1; \text{while}(e_2) \{s'; e_3\}$ – as long as s' does not contain a **continue** statement.

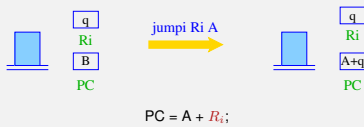
Thus, we translate:

```
code^i for(e1; e2; e3) s ρ = code^iR e1 ρ
                          A: code^iR e2 ρ
                          jumpz Ri B
                          code^i s ρ
                          code^iR e3 ρ
                          jump A
                          B:
```

The switch-Statement

Idea:

- Suppose choosing from multiple options in **constant time** if possible
- use a **jump table** that, at the i th position, holds a jump to the i th alternative
- in order to realize this idea, we need an **indirect jump** instruction



Consecutive Alternatives

Let $\text{switch } s$ be given with k consecutive **case** alternatives:

```
switch (e) {
  case 0: s0; break;
  :
  case k-1: sk-1; break;
  default: sk; break;
}
```

```
code^i s ρ = code^iR e ρ
            check^i 0 k B      B: jump A0
            A0: code^i s0 ρ      :
            jump C              jump Ak
            :                   C:
            Ak: code^i sk ρ
            jump C
```

Define $code^i s \rho$ as follows:

$check^i l u B$ checks if $l \leq R_i < u$ holds and jumps accordingly.

Translation of the $check^i$ Macro

The macro $check^i l u B$ checks if $l \leq R_i < u$. Let $k = u - l$.

- if $l \leq R_i < u$ it jumps to $B + R_i - l$
- if $R_i < l$ or $R_i \geq u$ it jumps to A_k

we define:

```
check^i l u B = loadc Ri+1 l
               geq Ri+2 Ri Ri+1} B: jump A0
               jumpz Ri+2 E
               sub Ri Ri Ri+1} :
               loadc Ri+1 k      jump Ak
               geq Ri+2 Ri Ri+1}
               jumpz Ri+2 D      C:
               E: loadc Ri k
               D: jumpi Ri B
```

Note: a jump $\text{jumpi } R_i B$ with $R_i = u$ winds up at $B + u$, the default case

Improvements for Jump Tables

This translation is only suitable for **certain** **switch**-statement.

- In case the table starts with 0 instead of u we don't need to subtract it from e before we use it as index
- if the value of e is **guaranteed** to be in the interval $[l, u]$, we can omit $check$

General translation of switch-Statements

In general, the values of the various cases may be far apart:

- generate an **if-ladder**, that is, a sequence of **if**-statements
- for n cases, an **if-cascade** (tree of conditionals) can be generated $\sim O(\log n)$ tests
- if the sequence of numbers has small gaps (≤ 3), a jump table may be smaller and faster
- one could generate several jump tables, one for each sets of consecutive cases
- an **if** cascade can be re-arranged by using information from **profiling**, so that paths executed more frequently require fewer tests

Code Synthesis

Chapter 4: Functions

Ingredients of a Function

The definition of a function consists of

- a **name** with which it can be called;
- a specification of its **formal parameters**;
- possibly a **result type**;
- a sequence of **statements**.

In C we have:

```
code^iR f ρ = loadc Ri _f with _f starting address of f
```

Observe:

- function names must have an address assigned to them
- since the size of functions is unknown before they are translated, the addresses of forward-declared functions must be inserted later

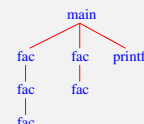
Memory Management in Functions

```
int fac(int x) {
  if (x<=0) return 1;
  else return x*fac(x-1);
}

int main(void) {
  int n;
  n = fac(2) + fac(1);
  printf("%d", n);
}
```

At run-time several **instances** may be active, that is, the function has been called but has not yet returned.

The recursion tree in the example:



Memory Management in Function Variables

The **formal parameters** and the **local variables** of the various instances of a function must be kept separate

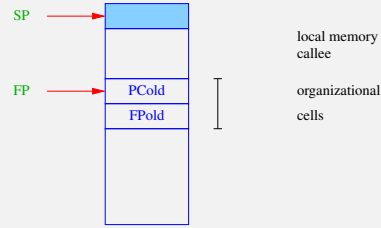
Idea for implementing functions:

- set up a region of memory each time it is called
- in sequential programs this memory region can be allocated on the stack
- thus, each instance of a function has its own region on the stack
- these regions are called **stack frames**

37/49

Organization of a Stack Frame

- stack representation: grows upwards
- SP points to the last used stack cell



- FP $\hat{=}$ frame pointer: points to the last organizational cell
- used to recover the previously active stack frame

38/49

Split of Obligations

Definition

Let f be the current function that calls a function g .

- f is dubbed **caller**
- g is dubbed **callee**

The code for managing function calls has to be split between caller and callee. This split cannot be done arbitrarily since some information is only known in that caller or only in the callee.

Observation:

The space requirement for parameters is only known by the caller:
Example: `printf`

39/49

Principle of Function Call and Return

actions taken on **entering** g :

1. compute the start address of g
2. compute actual parameters in globals
3. backup of caller-save registers
4. backup of FP
5. set the new FP
6. back up of PC and jump to the beginning of g
7. copy actual params to locals

} saveLoc
} mark
} call
} ... } is in g

actions taken on **leaving** g :

1. compute the result into R_0
2. restore FP, SP
3. return to the call site in f , that is, restore PC
4. restore the caller-save registers

} return
} restoreLoc } is in f

40/49

Managing Registers during Function Calls

The two register sets (global and local) are used as follows:

- automatic variables live in **local** registers R_i
- intermediate results also live in **local** registers R_i
- parameters live in **global** registers R_i (with $i \leq 0$)
- global variables: let's suppose there are none

convention:

- the i th argument of a function is passed in register R_{-i}
- the result of a function is stored in R_0
- local registers are saved before calling a function

Definition

Let f be a function that calls g . A register R_i is called

- **caller-saved** if f backs up R_i and g may overwrite it
- **callee-saved** if f does not back up R_i , and g must restore it before returning

41/49

Translation of Function Calls

A function call $g(e_1, \dots, e_n)$ is translated as follows:

```
codeRi g(e1, ... en) ρ = codeRi g ρ
                        codeRi+1 e1 ρ
                        ⋮
                        codeRi+n en ρ
                        move R-1 Ri+1
                        ⋮
                        move R-n Ri+n
                        saveLoc R1 Ri-1
                        mark
                        call Ri
                        restoreLoc R1 Ri-1
```

New instructions:

- **saveLoc** $R_i R_j$ pushes the registers R_i, R_{i+1}, \dots, R_j onto the stack
- **mark** backs up the organizational cells
- **call** R_i calls the function at the address in R_i
- **restoreLoc** $R_i R_j$ pops R_j, R_{j-1}, \dots, R_i off the stack

42/49

Rescuing the FP

The instruction **mark** allocates stack space for the return value and the organizational cells and backs up FP.



S[SP+1] = FP;
SP = SP + 1;

43/49

Calling a Function

The instruction **call** rescues the value of PC+1 onto the stack and sets FP and PC.



SP = SP+1;
S[SP] = PC;
FP = SP;
PC = Ri;

44/49

Result of a Function

The global register set is also used to communicate the result value of a function:

```
codei return e ρ = codeRi e ρ
                  move R0 Ri
                  return
```

alternative without result value:

```
codei return ρ = return
```

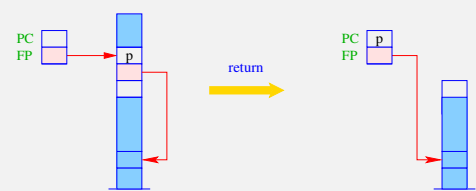
global registers are otherwise not used inside a function body:

- advantage: at any point in the body another function can be called without backing up **global** registers
- disadvantage: on entering a function, all **global** registers must be saved

45/49

Return from a Function

The instruction **return** relinquishes control of the current stack frame, that is, it restores PC and FP.



PC = S[FP];
SP = FP-2;
FP = S[SP+1];

46/49

Translation of Functions

The translation of a function is thus defined as follows:

```
code1 t, f( args ) { decls ss } ρ = move Rt+1 R-1
    ⋮
    move Rt+n R-n
    codet+n+1 ss ρ'
    return
```

Assumptions:

- the function has n parameters
- the local variables are stored in registers R_1, \dots, R_t
- the parameters of the function are in R_{-1}, \dots, R_{-n}
- ρ' is obtained by extending ρ with the bindings in *decls* and the function parameters *args*
- **return** is not always necessary

Are the **move** instructions always necessary?

47/49

Translation of Whole Programs

A program $P = F_1; \dots; F_n$ must have a single `main` function.

```
code1 P ρ = load R1 _main
    mark
    call R1
    halt
    _f1 : code1 F1 ρ ⊕ ρf1
    ⋮
    _fn : code1 Fn ρ ⊕ ρfn
```

Assumptions:

- $\rho = \emptyset$ assuming that we have no global variables
- ρ_{f_i} contain the addresses of the functions up to f_i
- $\rho_1 \oplus \rho_2 = \lambda x. \begin{cases} \rho_2(x) & \text{if } x \in \text{dom}(\rho_2) \\ \rho_1(x) & \text{otherwise} \end{cases}$

48/49

Translation of the `fac`-function

Consider:

```
int fac(int x) {
  if (x<=0)
    return 1;
  else
    return x*fac(x-1);
}

_fac: move R1 R-1 save param.
i=2  move R2 R1   if (x<=0)
    load R3 0
    leq R2 R2 R3
    jumpz R2 _A to else
    load R2 1   return 1
    move R0 R2
    return
    jump _B     code is dead

_A: move R2 R1   x*fac(x-1)
i=3  loadc R3 _fac
i=4  move R4 R1   x-1
i=5  loadc R5 1
i=6  sub R4 R4 R5
i=5  move R-1 R4   fac(x-1)
i=3  save loc R1 R2
    mark
    call R3
    restore loc R1 R2
    move R3 R0
    mul R2 R2 R3
i=4  mul R2 R2 R3
i=3  move R0 R2   return x*...
    return
_B: return
```

49/49