

Script generated by TTT

Title: Petter: Compilerbau (13.06.2019)

Date: Thu Jun 13 14:18:41 CEST 2019

Duration: 90:31 min

Pages: 35

What if precedences are not enough?

Example (very simplified lambda expressions):

$$\begin{aligned} E &\rightarrow (E)^0 | \text{ident}^1 | L^2 \\ L &\rightarrow \langle \text{args} \rangle \Rightarrow E^0 \\ \langle \text{args} \rangle &\rightarrow (\langle \text{idlist} \rangle)^0 | \text{ident}^1 \\ \langle \text{idlist} \rangle &\rightarrow \langle \text{idlist} \rangle \text{ident}^0 | \text{ident}^1 \end{aligned}$$

E rightmost-derives these forms among others:

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

Naive Idea:

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

149 / 288

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,

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

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,

- 1 there exist several practical $LR(k)$ grammars of $k > 1$, e.g. Java 1.6+ ($LR(2)$), ANSI C, etc.
- 2 often, more lookahead is only exhausted locally
- 3 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:

$$\begin{aligned} S &\rightarrow A \underline{b} b^0 \mid B \underline{b} c^1 \\ A &\rightarrow a A^0 \mid a^1 b \\ B &\rightarrow a B^0 \mid a^1 b \end{aligned}$$

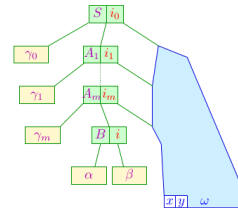
S rightmost-derives one of these forms:

$$a^n \underline{a} b b, a^n \underline{a} b c, a^n \underline{a} A b b, a^n \underline{a} B b c, \underline{A} b b, \underline{B} b c \Rightarrow LR(2)$$

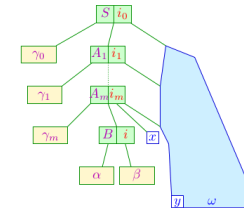
in $LR(1)$, you will have Reduce-/Reduce-Conflicts between the productions $A, 1$ and $B, 1$ under lookahead b

LR(2) to LR(1)

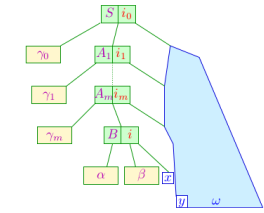
Basic Idea:



Right-context-extraction



Right-context-propagation



in the example:

Right-context is already extracted, so we only perform

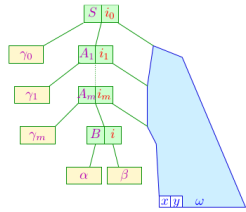
Right-context-propagation:

$$\begin{aligned} S &\rightarrow A b b^0 \mid B \underline{b} c^1 \\ A &\rightarrow a A^0 \mid a^1 b \\ B &\rightarrow a B^0 \mid a^1 b \end{aligned}$$

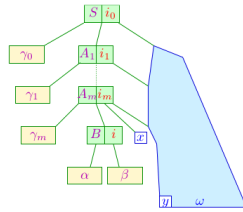
$$\begin{aligned} S &\rightarrow \langle A b \rangle b^0 \mid \langle B b \rangle c^1 \\ \langle A b \rangle &\rightarrow a \langle A b \rangle^0 \mid a b^1 \\ \langle B b \rangle &\rightarrow a \langle B b \rangle^0 \mid a b^1 \end{aligned}$$

LR(2) to LR(1)

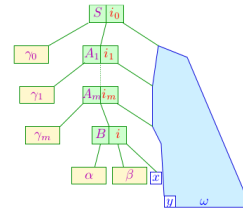
Basic Idea:



Right-context-extraction



Right-context-propagation



in the example:

Right-context is already extracted, so we only perform

Right-context-propagation:

$$\begin{aligned} S &\rightarrow A b b^0 \mid B b c^1 \\ A &\rightarrow a A^0 \mid a^1 b \\ B &\rightarrow a B^0 \mid a^1 b \end{aligned}$$

⇒

$$\begin{aligned} S &\rightarrow \langle A b \rangle b^0 \mid \langle B b \rangle c^1 \\ \langle A b \rangle &\rightarrow a \langle A b \rangle^0 \mid a b^1 \\ \langle B b \rangle &\rightarrow a \langle B b \rangle^0 \mid a b^1 \end{aligned}$$

unreachable

LR(2) to LR(1)

Example 2:

$$\begin{aligned} S &\rightarrow b S S^0 \\ &\quad \mid a^1 \\ &\quad \mid a a c^2 \end{aligned}$$

S rightmost-derives these forms among others:

$$\underline{b} S S, \underline{b} S a, \underline{b} S a a c, \underline{b} a a, \underline{b} a a c a, \underline{b} a a c, \underline{b} a a c a a c, \dots \Rightarrow \text{min. } LR(2)$$

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 ⇒ perform Right-context-extraction

$$\begin{aligned} S &\rightarrow \boxed{b} S S^0 \\ &\quad \mid \boxed{a} \\ &\quad \mid \boxed{a a c^2} \end{aligned}$$

⇒

$$\begin{aligned} S &\rightarrow b S \langle a/S \rangle^0 \mid b S \langle b/S \rangle^0 \\ &\quad \mid a^1 \mid a a c^2 \\ \langle a/S \rangle &\rightarrow \epsilon^0 \mid a c^1 \\ \langle b/S \rangle &\rightarrow S S^0 \end{aligned}$$

LR(2) to LR(1)

Example 2 cont'd:

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

$$\begin{array}{l}
 S \rightarrow bSa \langle a/S \rangle^0 \\
 \quad | bSb \langle b/S \rangle^{0'} \\
 \quad | a^1 | aac^2 \\
 \langle a/S \rangle \rightarrow \epsilon^0 | ac^1 \\
 \langle b/S \rangle \rightarrow Sa \langle a/S \rangle^0 | Sb \langle b/S \rangle^{0'}
 \end{array}
 \Rightarrow
 \begin{array}{l}
 S \rightarrow b \langle Sa \rangle \langle a/S \rangle^0 \\
 \quad | bSb \langle b/S \rangle^{0'} \\
 \quad | a^1 | aac^2 \\
 \langle a/S \rangle \rightarrow \epsilon^0 | ac^1 \\
 \langle Sa \rangle \rightarrow b \langle Sa \rangle \langle \langle a/S \rangle a \rangle^0 \\
 \quad | bSb \langle \langle b/S \rangle a \rangle^{0'} \\
 \quad | a \langle a \rangle^1 | aac \langle a \rangle^2 \\
 \langle \langle a/S \rangle a \rangle \rightarrow a^0 | aca^1 \\
 \langle \langle b/S \rangle a \rangle \rightarrow \langle Sa \rangle \langle \langle a/S \rangle a \rangle^0 | Sb \langle \langle b/S \rangle a \rangle^{0'}
 \end{array}$$

LR(2) to LR(1)

Example 2 finished:

With fresh nonterminals we get the final grammar

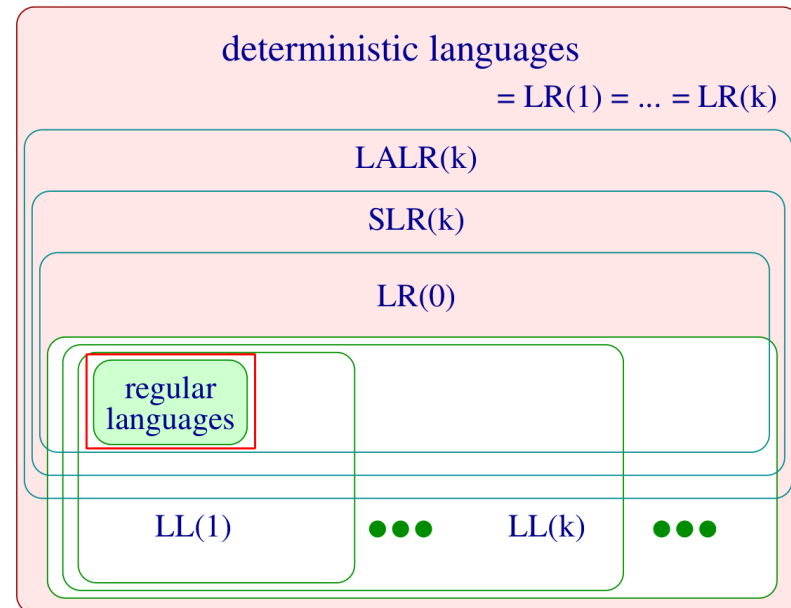
$$\begin{array}{l}
 S \rightarrow bSS^0 \\
 \quad | a^1 \\
 \quad | aac^2
 \end{array}
 \Rightarrow
 \begin{array}{l}
 S \rightarrow bCA^0 | bSbB^1 | a^2 | aac^3 \\
 A \rightarrow \epsilon^0 | ac^1 \\
 B \rightarrow CA^0 | SbB^1 \\
 C \rightarrow bCD^0 | bSbE^1 | aa^2 | aaca^3 \\
 D \rightarrow a^0 | aca^1 \\
 E \rightarrow CD^0 | SbE^1
 \end{array}$$

Special LR(k)-Subclasses

Discussion:

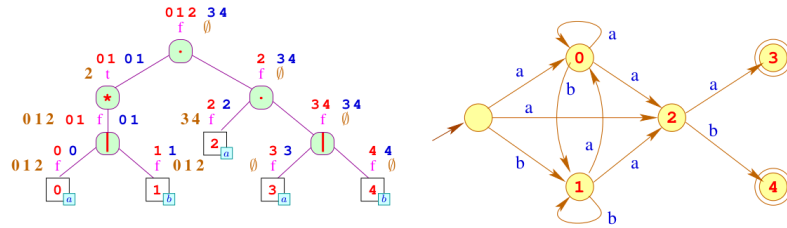
- Our examples mostly were $LR(1)$ – or could be transformed to $LR(1)$
- In general, the canonical $LR(k)$ -automaton has much more states than $LR(G) = LR(G, 0)$
- Therefore in practice, subclasses of $LR(k)$ -grammars are often considered, which only use $LR(G) \dots$
- For resolving conflicts, the items are assigned special lookahead-sets:
 - ① independently on the state itself \Rightarrow Simple $LR(k)$
 - ② dependent on the state itself \Rightarrow $LALR(k)$

Parsing Methods

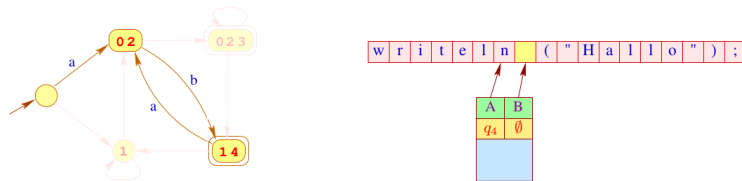


Lexical and Syntactical Analysis:

From Regular Expressions to Finite Automata



From Finite Automata to Scanners

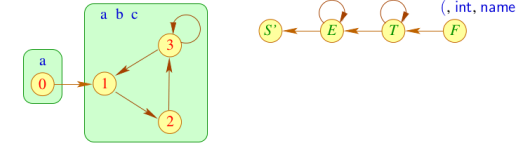


164 / 288

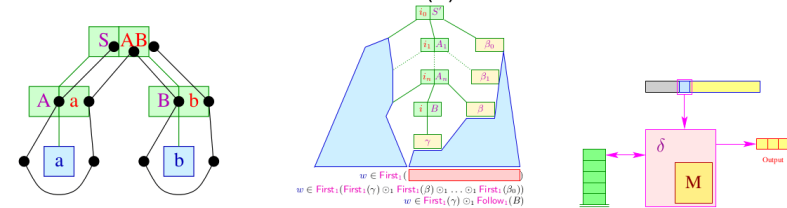
Lexical and Syntactical Analysis:

Computation of lookahead sets:

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



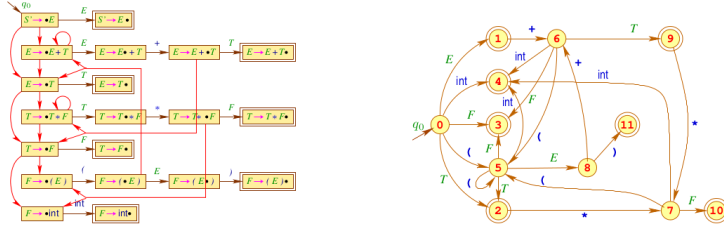
From Item-Pushdown Automata to LL(1)-Parsers:



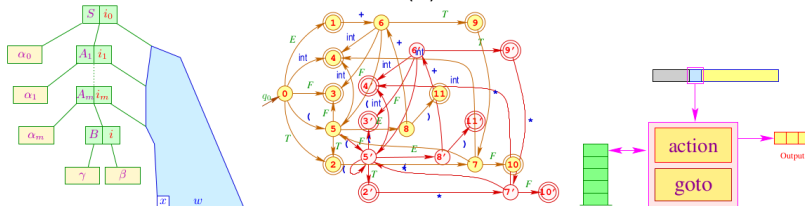
165 / 288

Lexical and Syntactical Analysis:

From characteristic to canonical Automata:



From Shift-Reduce-Parsers to LR(1)-Parsers:



166 / 288

Organizing

- Master or Bachelor in the 6th Semester with 5 ECTS
- Prerequisites
 - Basic Programming: **Java**
 - *Introduction to Theory of Computation*
 - Basic Principles: Operating Systems and System Software
 - Automata Theory
- Delve deeper with
 - Virtual Machines
 - Programm Optimization
 - Programming Languages
 - Labcourse Compiler Construction

Materials:

- TTT-based lecture recordings
- The slides
- Related literature list online (\Rightarrow Wilhelm/Seidl/Hack Compiler Design)
- Tools for visualization of abstract machines (VAM)
- Tools for generating components of Compilers (JFlex/CUP)

2 / 288

Dates:

Lecture: Thursdays 14:15-15:45

Tutorial: Mondays 14:15-15:45, Wednesdays 8:30-10:00

Exam:

- **One Exam** in the summer, *none* in the winter
- planned date: 13.8.19 (no guarantee!)
- Successful mini seminar earns 0.3 bonus

Mini Seminars in week 30:

A 15min talk on a specific compiler related topic, e.g.

1	Introduction to parser combinators	ES
2	Parsing Expression Grammars / Packrat Parsers	
3	LALR(k)	MZ
4	Pager's LR(k)	SK
5	GLR / Tomita parser	MZ, AM
6	Follow-Automata	AM
7	Antimirov-Automata	MT
8	LL(*) and LL-regular grammars	ANTERI
9	Semidivision Procedures for Grammar Uniqueness	CT
10	Island Grammars	CC
11	Learning CFGs with LZW and Sequitur	PK
12	General First _k /Follow _k with Constraint Solvers	AB
13	Reference Attribute Grammars with JastAdd	
14	Extensible Grammars with PPG	PB
15	Language Processing with Kiama/Scala	Laura

Preliminary content

- Regular expressions and finite automata
- Specification and implementation of scanners
- Reduced context free grammars and pushdown automata
- Top-Down/Bottom-Up syntax analysis
- Attribute systems
- Typechecking
- Codegeneration for register machines
- Register assignment
- Optional: Basic optimization

1	Introduction to parser combinators	ES
2	Parsing Expression Grammars / Packrat Parsers	
3	LALR(k)	MZ
4	Pager's LR(k)	SK
5	GLR / Tomita parser	MZ, AM
6	Follow-Automata	AM
7	Antimirov-Automata	MT
8	LL(*) and LL-regular grammars	ANTERI
9	Semidivision Procedures for Grammar Uniqueness	CT
10	Island Grammars	CC
11	Learning CFGs with LZW and Sequitur	PK
12	General First _k /Follow _k with Constraint Solvers	AB
13	Reference Attribute Grammars with JastAdd	
14	Extensible Grammars with PPG	PB
15	Language Processing with Kiama/Scala	Laura

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**

Semantic Analysis

Chapter 1: Attribute Grammars

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

170 / 288

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

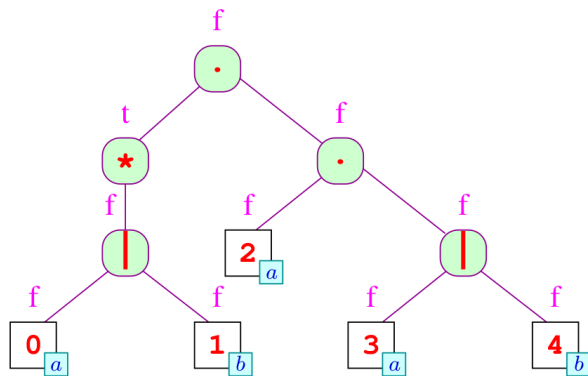
An *attribute grammar* is a CFG extended by

- 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*

170 / 288

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*)

171 / 288

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 *synthetic* attribute
- The *local* dependencies between the attributes are dependent on the *type* of the node

172 / 288

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 **synthetic** attribute
- The **local** dependencies between the attributes are dependent on the **type** of the node

in general:

Definition

An attribute is called

- **synthetic** if its value is always propagated upwards in the tree (in the direction leaf \rightarrow root)
- **inherited** if its value is always propagated downwards in the tree (in the direction root \rightarrow leaf)

172 / 288

Attribute Equations for **empty**

In order to compute an attribute **locally**, we need to specify attribute equations for each node.

These equations depend on the **type** of the node:

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

otherwise:

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

173 / 288

Attribute Equations for **empty**

In order to compute an attribute **locally**, we need to specify attribute equations for each node.

These equations depend on the **type** of the node:

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

otherwise:

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

173 / 288

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**:

\leadsto We use consecutive indices to refer to neighbouring attributes

$\text{attribute}_k[0]$: the attribute of the current root node
 $\text{attribute}_k[i]$: the attribute of the i -th child ($i > 0$)

... in the example:

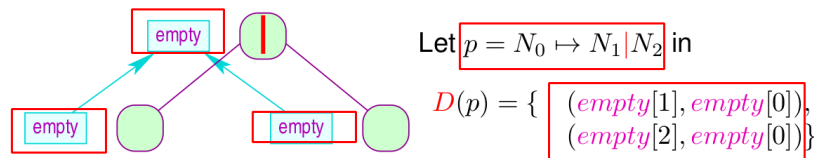
\boxed{x} : $\text{empty}[0] := (x \equiv \epsilon)$
 $\boxed{\mid}$: $\text{empty}[0] := \text{empty}[1] \vee \text{empty}[2]$
 $\boxed{\cdot}$: $\text{empty}[0] := \text{empty}[1] \wedge \text{empty}[2]$
 $\boxed{*}$: $\text{empty}[0] := t$
 $\boxed{?}$: $\text{empty}[0] := t$

174 / 288

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*:



~> arrows point in the direction of information flow

175 / 288

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 thus change with each new 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*

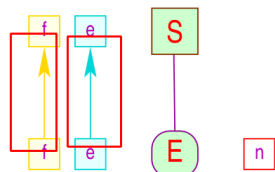
176 / 288

Simultaneous Computation of Multiple Attributes

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

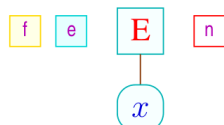
$$\begin{aligned}
 S \rightarrow E : & \quad \begin{aligned} empty[0] &:= empty[1] \\ first[0] &:= first[1] \\ next[1] &:= \emptyset \end{aligned} \\
 E \rightarrow x : & \quad \begin{aligned} empty[0] &:= (x \equiv \epsilon) \\ first[0] &:= \{x \mid x \neq \epsilon\} \\ & // \text{ (no equation for next) } \end{aligned}
 \end{aligned}$$

$D(S \rightarrow E) :$



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

$D(E \rightarrow x) :$



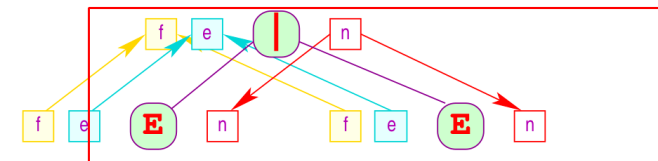
$$D(E \rightarrow x) = \{ \}$$

177 / 288

Regular Expressions: Rules for Alternative

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

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



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

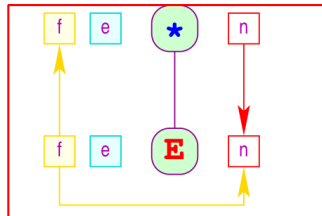
178 / 288

Regular Expressions: Kleene-Star and '?'

$E \rightarrow E^*$: $\text{empty}[0] := t$
 $\text{first}[0] := \text{first}[1]$
 $\text{next}[1] := \text{first}[1] \cup \text{next}[0]$

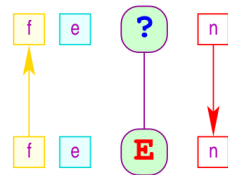
$E \rightarrow E?$: $\text{empty}[0] := t$
 $\text{first}[0] := \text{first}[1]$
 $\text{next}[1] := \text{next}[0]$

$D(E \rightarrow E^*) :$



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

$D(E \rightarrow E?) :$



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

180 / 288

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]

181 / 288

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

181 / 288

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 $X \mapsto X_1 \dots X_k$'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

182 / 288