

Compiler construction

Martin Steffen

January 21, 2016

Contents

1 Scanning	1
1.1 Intro	1
1.2 Regular expressions	5
1.3 DFA	12
1.4 Implementation of DFA	17
1.5 NFA	18
1.6 From regular expressions to DFAs	19
1.7 Thompson’s construction	20
1.8 Determinization	22
1.9 Scanner generation tools	25

Abstract

This is the *handout* version of the slides. It contains basically the same content, only in a way which allows more compact printing. Sometimes, the overlays, which make sense in a presentation, are not fully rendered here.

Besides the material of the slides, the *handout* versions also contain additional remarks and background information which may or may not be helpful in getting the bigger picture.

1 Scanning

19. 01. 2016

1.1 Intro

1. Scanner section overview

(a) what’s a **scanner**?

- Input: source code.¹
- Output: sequential stream of **tokens**
- (b) • *regular expressions* to describe various token classes
- (deterministic/nondeterministic) finite-state automata (FSA, DFA, NFA)
- implementation of FSA
- regular expressions → NFA
- NFA ↔ DFA

2. What’s a scanner?

- other names: lexical scanner, **lexer**, tokenizer

(a) A scanner’s functionality Part of a compiler that takes the source code as input and translates this stream of characters into a stream of **tokens**.

(b) More info

- char’s typically language independent.²
- *tokens* already language-specific.³
- works always “left-to-right”, producing one *single token* after the other, as it scans the input⁴

¹The argument of a scanner is often a *file name* or an *input stream* or similar.

²Characters are language-independent, but perhaps the encoding may vary, like ASCII, UTF-8, also Windows-vs.-Unix-vs.-Mac newlines etc.

³There are large commonalities across many languages, though.

⁴No theoretical necessity, but that’s how also humans consume or “scan” a source-code text. At least those humans trained in e.g. Western languages.

- it “segments” char stream into “chunks” while at the same time “classifying” those pieces \Rightarrow **tokens**

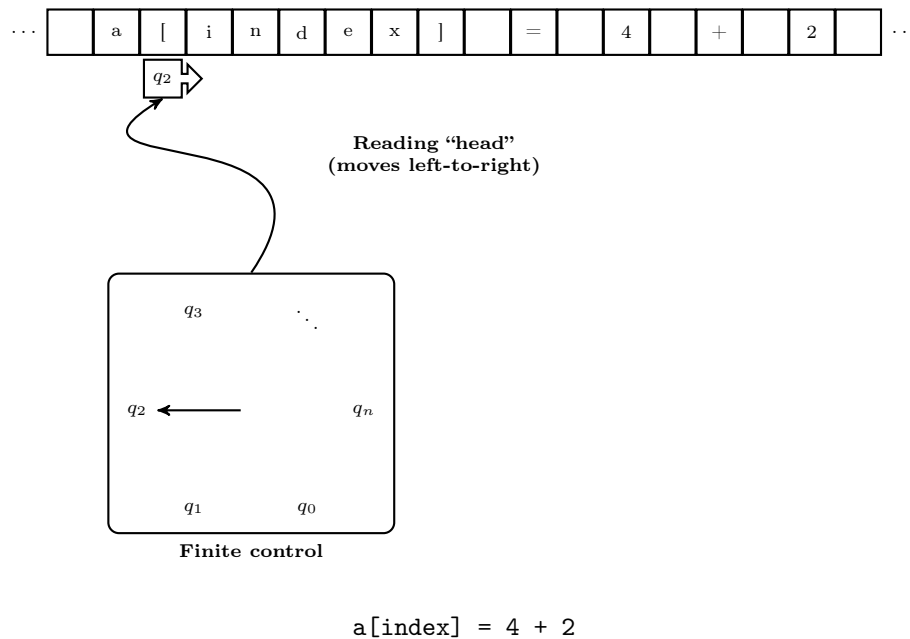
3. Typical responsibilities of a scanner

- segment & classify char stream into tokens
- typically described by “rules” (and **regular expressions**)
- typical language aspects covered by the scanner
 - describing *reserved words* or *key words*
 - describing format of *identifiers* (= “strings” representing variables, classes ...)
 - comments (for instance, between // and NEWLINE)
 - *white space*
 - * to segment into tokens, a scanner typically “jumps over” white spaces and afterwards starts to determine a new token
 - * not only “blank” character, also TAB, NEWLINE, etc.
- lexical rules: often (explicit or implicit) *priorities*
 - *identifier* or *keyword*? \Rightarrow keyword
 - take the *longest* possible scan that yields a valid token.

4. “Scanner = Regular expressions (+ priorities)”

- (a) Rule of thumb Everything about the source code which is so simple that it can be captured by **reg. expressions** belongs into the scanner.

5. How does scanning roughly work?



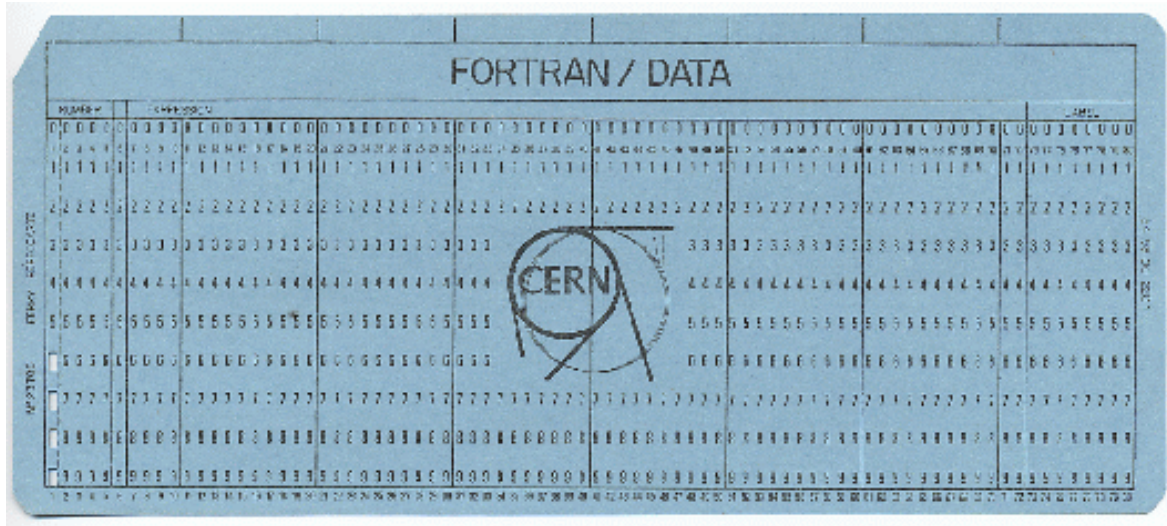
6. How does scanning roughly work?

- usual invariant in such pictures (by convention): arrow or head points to the *first* character to be *read next* (and thus *after* the last character having been scanned/read last)
- in the scanner *program* or procedure:
 - analogous invariant, the arrow corresponds to a *specific variable*
 - contains/points to the next character to be read
 - name of the variable depends on the scanner/scanner tool
- the *head* in the pic: for illustration, the scanner does not really have a “reading head”
 - remembrance of Turing machines, or
 - the old times when perhaps the program data was stored on a tape.⁵

⁵Very deep down, if one still has a magnetic disk (as opposed to SSD) the secondary storage still has “magnetic heads”, only that one typically does not parse *directly* char by char from disk...

7. The bad old times: Fortran

- in the days of the pioneers
 - main memory was *smaaaaaaaaaall*
 - compiler technology was not well-developed (or not at all)
 - programming was for *very* few “experts”.⁶
 - Fortran was considered a very high-level language (wow, a language so complex that you had to compile it ...)



8. (Slightly weird) lexical aspects of Fortran Lexical aspects = those dealt with a scanner

- **whitespace** *without* “meaning”:

I F (X 2. EQ. 0) TH E N vs. IF (X2. EQ.0) THEN

- no *reserved* words!

IF (IF.EQ.0) THEN THEN=1.0

- general obscurity tolerated:

D099I=1,10 vs. D099I=1.10

```
DO 99 I=1,10
  —
  —
99 CONTINUE
```

9. Fortran scanning: remarks

- Fortran (of course) has evolved from the pioneer days ...
- no keywords: nowadays mostly seen as *bad* idea⁷
- treatment of white-space as in Fortran: not done anymore: THEN and TH EN *are* different things in all languages
- however:⁸ both considered “the same”:

(a) :BMCOL:B_{block}:

```
if_then...
```

⁶There was no computer science as profession or university curriculum.

⁷It’s mostly a question of language *pragmatics*. The lexers/parsers would have no problems using `while` as variable, but humans tend to have.

⁸Sometimes, the part of a lexer / parser which removes whitespace (and comments) is considered as separate and then called *screener*. Not very common though.

(b) :BMCOL:B_{block}:

```
if_____then_...
```

- (c)
- since concepts/tools (and much memory) were missing, Fortran scanner and parser (and compiler) were
 - quite simplistic
 - syntax: designed to “help” the lexer (and other phases)

10. A scanner classifies

- “good” classification: depends also on later phases, may not be clear till later
- (a) Rule of thumb Things being treated equal in the syntactic analysis (= parser, i.e., subsequent phase) should be put into the same category.
- (b) End
- terminology not 100% uniform, but most would agree:
- (c) Lexemes and tokens **Lexemes** are the “chunks” (pieces) the scanner produces from segmenting the input source code (and typically dropping whitespace). **Tokens** are the result of /classifying those lexemes.
- (d) End
- token = token name × token value

11. A scanner classifies & does a bit more

- token data structure in *OO* settings
 - token themselves defined by classes (i.e., as instance of a class representing a specific token)
 - token values: as attribute (instance variable) in its values
- often: scanner does slightly *more* than just classification
 - store names in some *table* and store a corresponding index as attribute
 - store text constants in some *table*, and store corresponding index as attribute
 - even: *calculate* numeric constants and store value as attribute

12. One possible classification

name/identifier	abc123
integer constant	42
real number constant	3.14E3
text constant, string literal	"this is a text constant"
arithmetic op's	+ - * /
boolean/logical op's	and or not (alternatively /\ \/)
relational symbols	<= < >= > = == !=
all other tokens:	{ } () [] , ; := . etc.
every one it its own group	

- this classification: not the only possible (and not necessarily complete)
- note: *overlap*:
 - "." is here a token, but also part of real number constant
 - "<" is part of "<="

13. One way to represent tokens in C

```
typedef struct {
    TokenType tokenval;
    char * stringval;
    int numval;
} TokenRecord;
```

If one only wants to store one attribute:

```

typedef struct {
    Tokentype tokenval;
    union
    { char * stringval;
      int numval
    } attribute;
} TokenRecord;

```

14. How to define lexical analysis and implement a scanner?

- even for complex languages: lexical analysis (in principle) not hard to do
 - “manual” implementation straightforwardly possible
 - *specification* (e.g., of different token classes) may be given in “prosa”
 - however: there are straightforward formalisms and efficient, rock-solid tools available:
 - easier to specify unambiguously
 - easier to communicate the lexical definitions to others
 - easier to change and maintain
 - often called **parser generators** typically not just generate a scanner, but code for the next phase (parser), as well.
- (a) Prosa specification A precise *prosa* specification is not so easy to achieve than one might think. For ASCII source code or input, things are basically under control. But what if dealing with unicode? Checking “legality” of user input to avoid SQL injections or format string attacks is largely done by lexical analysis/scanning. If you “specify” in English: “ *Backlash is a control character and forbidden as user input* ”, which characters (besides `char 92` in ASCII) in Chinese Unicode represents actually other versions of *backslash*?
- (b) Parser generator The most famous pair of lexer+parser is called “compiler compiler” (lex/yacc = “yet another compiler compiler”) since it generates (or “compiles”) an important part of the front end of a compiler, the lexer+parser. Those kind of tools are seldomly called compiler compilers any longer.

1.2 Regular expressions

1. General concept: How to generate a scanner?

- (a) **regular expressions** to describe language’s *lexical* aspects
- like whitespaces, comments, keywords, format of identifiers etc.
 - often: more “user friendly” variants of reg-expr are supported to specify that phase
- (b) *classify* the lexemes to tokens
- (c) translate the reg-expressions \Rightarrow NFA.
- (d) turn the NFA into a *deterministic* FSA (= DFA)
- (e) the DFA can straightforwardly be implemented
- Above steps are done automatically by a “lexer generator”
 - lexer generators help also in other user-friendly ways of specifying the lexer: defining *priorities*, assuring that the longest possible token is given back, repeat the process to generate a sequence of tokens⁹
 - Step 2 is actually *not* covered by the classical Reg-expr = DFA = NFA results, it’s something extra.

2. Use of regular expressions

- **regular languages**: fundamental class of “languages”
- **regular expressions**: standard way to describe regular languages
- origin of regular expressions: one starting point is Kleene [Kleene, 1956] but there had been earlier works outside “computer science”
- Not just used in compilers
- often used for flexible “ *searching* ”: simple form of [pattern matching](#)
- e.g. input to search engine interfaces

⁹Maybe even prepare useful error messages if scanning (not scanner generation) fails.

- also supported by many editors and text processing or scripting languages (starting from classical ones like `awk` or `sed`)
- but also tools like `grep` or `find`

```
find . -name "*.tex"
```

- often *extended* regular expressions, for user-friendliness, not theoretical expressiveness.

(a) Remarks

Kleene was a famous mathematician and influence on theoretical computer science. Funnily enough, regular languages came up in the context of neuro/brain science. See the following for the origin of the terminology. Perhaps in the early years, people liked to draw connections between between biology and machines and used metaphors like “electronic brain” etc.

3. Alphabets and languages

Definition 1 (Alphabet Σ). Finite set of elements called “letters” or “symbols” or “characters”

Definition 2 (Words and languages over Σ). Given alphabet Σ , a **word** over Σ is a finite sequence of letters from Σ . A **language** over alphabet Σ is a *set* of finite *words* over Σ .

- in this lecture: we avoid terminology “symbols” for now, as later we deal with e.g. symbol tables, where symbols means something slightly different (at least: at a different level).
- Sometimes Σ left “implicit” (as assumed to be understood from the context)
- practical examples of alphabets: ASCII, Norwegian letters (capital and non-capitals) etc.

(a) Remark: Symbols in a symbol table

In a certain way, symbols in a symbol table can be seen similar than symbols in the way we are handled by automata or regular expressions now. They are simply “atomic” (not further dividable) members of what one calls an alphabet. On the other hand, in practical terms inside a compiler, the symbols here live on a different level. Typically, they are *characters*, i.e., the alphabet is a so-called character set like for instance ASCII. The lexer, as stated, segments and classifies the sequence of characters and hands over the result of that process to the parser. The results is a squence of *tokens*, which is what the parser has to deal with later. It’s on that parser-level, that the pieces (notably the identifiers) can be treated as atomic pieces of some language, and what is known as the symbol table typically operates on symbols at that level, not at the level of individual characters.

4. Languages

- note: Σ is finite, and words are of *finite* length
- languages: in general *infinite* sets of words
- Simple examples: Assume $\Sigma = \{a, b\}$
- *words* as finite “sequences” of letters
 - ϵ : the empty word (= empty sequence)
 - ab means “ first a then b ”
- sample languages over Σ are
 - $\{\}$ (also written as \emptyset) the empty set
 - $\{a, b, ab\}$: language with 3 finite words
 - $\{\epsilon\}$ ($\neq \emptyset$)
 - $\{\epsilon, a, aa, aaa, \dots\}$: infinite languages, all words using only a ’s.
 - $\{\epsilon, a, ab, aba, abab, \dots\}$: alternating a ’s and b ’s
 - $\{ab, bbab, aaaaa, bbabbabab, aabb, \dots\}$: ?????

(a) Remarks

Remark 1 (Words and strings). *In terms of a real implementation: often, the letters are of type *character* (like type *char* or *char32* ...) words then are “sequences” (say arrays) of characters, which may or may not be identical to elements of type *string*, depending on the language for implementing the compiler. In a more conceptual part like here we do not write words in “string notation” (like “ab”), since we are dealing abstractly with sequences of letters, which, as said, may not actually be strings in the implementation. Also in the more conceptual parts, it’s often good enough when handling alphabets with 2 letters, only, like $\Sigma = \{a, b\}$ (with one letter, it gets unrealistically trivial and results may not carry over to the many-letter alphabets). After all, computers are using 2 bits only, as well ...*

(b) Finite and infinite words

There are important applications dealing with infinite words, as well, or also infinite alphabets. For traditional scanners, one mostly is happy with finite Σ 's and especially sees no use in scanning infinite "words".

5. How to describe languages

- language mostly here in the abstract sense just defined.
- the "dot-dot-dot" (...) is not a good way to describe to a computer (and many humans) what is meant (what was meant in the last example?)
- enumerating explicitly all allowed words for an infinite language does not work either

(a) Needed A **finite** way of describing infinite languages (which is hopefully efficiently implementable & easily readable)

Beware

Is it a priori clear to expect that *all* infinite languages can even be captured in a finite manner?

- small metaphor

$$2.727272727 \dots \quad 3.1415926 \dots \quad (1)$$

(b) Remarks

Remark 2 (Programming languages as "languages"). *Well, Java etc., seen syntactically as all possible strings that can be compiled to well-formed byte-code, also is a language in the sense we are currently discussing, namely a set of words over unicode. But when speaking of the "Java-language" or other programming languages, one typically has also other aspects in mind (like what a program does when it is executed), which is not covered by thinking of Java as an infinite set of strings.*

Remark 3 (Rational and irrational numbers). *The illustration here with the two numbers is partly meant as that: an illustration drawn from a field you may know. The first number from equation (1) is a rational number. It corresponds to the fraction*

$$\frac{30}{11} . \quad (2)$$

That fraction is actually an acceptable finite representation for the "endless" notation 2.72727272... using "...". As one may remember, it may pass as a decent definition of rational numbers that they are exactly those which be represented finitely as fractions of two integers, like the one from equation (2). We may also remember that it is characteristic for the "endless" notation as the one from equation (1), that for rational numbers, it's periodic. Some may have learnt the notation

$$2.\overline{72} \quad (3)$$

for finitely representing numbers with a periodic digit expansion (which are exactly the rationals). The second number, of course, is π , one of the most famous numbers which do not belong to the rationals, but to the "rest" of the reals which are not rational (and hence called irrational). Thus it's one example of a "number" which cannot be represented by a fraction, resp. in the periodic way as in (3).

Well, fractions may not work out for π (and other irrationals), but still, one may ask, whether π can otherwise be represented finitely. That, however, depends on what actually one accepts as a "finite representation". If one accepts a finite description that describes how to construct ever closer approximations to π , then there is a finite representation of π . That construction basically is very old (Archimedes), it corresponds to the limits one learns in analysis, and there are computer algorithms, that spit out digits of π as long as you want (of course they can spit them out all only if you had infinite time). But the code of the algo who does that is finite.

The bottom line is: it's possible to describe infinite "constructions" in a finite manner, but what exactly can be captured depends on what precisely is allowed in the description formalism. If only fractions of natural numbers are allowed, one can describe the rationals but not more.

A final word on the analogy to regular languages. The set of rationals (in, let's say, decimal notation) can be seen as language over the alphabet $\{0, 1, \dots, 9, .\}$, i.e., the decimals and the "decimal point". It's however, a language containing infinite words, such as $2.727272727 \dots$. The syntax $2.\overline{72}$ is a finite expression but denotes the mentioned infinite word (which is a decimal representation of a rational number). Thus, coming back to the regular languages resp. regular expressions, $2.\overline{72}$ is similar to the Kleene-star, but not the same. If we write $2.(72)^$, we mean the language of finite words*

$$\{2, 2.72, 2.727272, \dots\} .$$

In the same way as one may conveniently define rational number (when represented in the alphabet of the decimals) as those which can be written using periodic expressions (using for instance overline), regular

languages over an alphabet are simply those sets of finite words that can be written by regular expressions (see later). Actually, there are deeper connections between regular languages and rational numbers, but it's not the topic of compiler constructions. Suffice to say that it's not a coincidence that regular languages are also called rational languages (but not in this course).

6. Regular expressions

Definition 3 (Regular expressions). A *regular expression* is one of the following

- (a) a *basic* regular expression of the form \mathbf{a} (with $a \in \Sigma$), or ϵ , or \emptyset
- (b) an expression of the form $r \mid s$, where r and s are regular expressions.
- (c) an expression of the form rs , where r and s are regular expressions.
- (d) an expression of the form r^* , where r is a regular expression.
- (e) an expression of the form (r) , where r is a regular expression.

Precedence (from high to low): $*$, concatenation, $|$

(a) Regular expressions

In other textbooks, also the notation $+$ instead of $|$ for “alternative” or “choice” is a known convention. The $|$ seems more popular in texts concentrating on *grammars*. Later, we will encounter *context-free* grammars (which can be understood as a generalization of regular expressions) and the $|$ -symbol is consistent with the notation of alternatives in the definition of rules or productions in such grammars. One motivation for using $+$ elsewhere is that one might wish to express “parallel” composition of languages, and a conventional symbol for parallel is $|$. We will not encounter parallel composition of languages in this course. Also, regular expressions using lot of parentheses and $|$ seems slightly less readable for humans than using $+$.

Regular expressions are a language in itself, so they have a syntax and a semantics. One could write a lexer (and parser) to parse a regular language. Obviously, tools like parser generators *do* have such a lexer/parser, because their input language are regular expression (and context free grammars, besides syntax to describe further things).

7. A concise definition later introduced as (notation for) context-free grammars:

$$\begin{aligned}
 r &\rightarrow \mathbf{a} \\
 r &\rightarrow \epsilon \\
 r &\rightarrow \emptyset \\
 r &\rightarrow r \mid r \\
 r &\rightarrow rr \\
 r &\rightarrow r^* \\
 r &\rightarrow (r)
 \end{aligned}
 \tag{4}$$

8. Same again

(a) Notational conventions

Later, for CF grammars, we use capital letters to denote “variables” of the grammars (then called *non-terminals*). If we like to be consistent with that convention, the definition looks as follows:

(b) Grammar

$$\begin{aligned}
 R &\rightarrow \mathbf{a} \\
 R &\rightarrow \epsilon \\
 R &\rightarrow \emptyset \\
 R &\rightarrow R \mid R \\
 R &\rightarrow RR \\
 R &\rightarrow R^* \\
 R &\rightarrow (R)
 \end{aligned}
 \tag{5}$$

9. Symbols, meta-symbols, meta-meta-symbols ...

- regexps: notation or “language” to describe “languages” over a given alphabet Σ (i.e. subsets of Σ^*)
 - language being described \Leftrightarrow language used to describe the language
- \Rightarrow language \Leftrightarrow meta-language
- here:
 - regular expressions: notation to describe regular languages

- English resp. context-free notation:¹⁰ notation to describe regular expression
- for now: carefully use *notational convention* for precision

10. Notational conventions

- notational conventions by *typographic* means (i.e., different fonts etc.)
- not easy discernible, but: difference between
 - \mathbf{a} and a
 - ϵ and ϵ
 - \emptyset and \emptyset
 - $|$ and $|$ (especially hard to see :-))
 - ...
- later (when gotten used to it) we may take a more “relaxed” attitude toward it, assuming things are clear, as do many textbooks
- Note: in compiler *implementations*, the distinction between language and meta-language etc. is very real (even if not done by typographic means ...)

(a) Remarks

Remark 4 (Regular expression syntax). *Later there will be a number of examples using regular expressions. There is a slight “ambiguity” about the way regular expressions are described (in this slides, and elsewhere). It may remain unnoticed (so it’s unclear if one should point it out). On the other had, the lecture is, among other things, about scanning and parsing of syntax, therefore it may be a good idea to reflect on the syntax of regular expressions themselves.*

In the examples shown later, we will use regular expressions using parentheses, like for instance in $b(ab)^$. One question is: are the parentheses (and) part of the definition of regular expressions or not? That depends a bit. In more mathematical texts, typically one would not care, one tells the readers that parantheses will be used for disambiguation, and leaves it at that (in the same way one would not tell the reader it’s fine to use “space” between different expressions (like $a | b$ is the same expression as $a | b$). Another way of saying that is that textbooks, intended for human readers, give the definition of regular expressions as abstract syntax as opposed to concrete syntax. It’s thereby assumed that the reader can interpret parentheses as grouping mechanism, as is common elsewhere as well and they are left out from the definition not to clutter it. Of course, computers (i.e., compilers), are not as good as humans to be educated in “commonly understood” conventions (such as “parentheses are not really part of the regular expressions but can be added for diambiguation”. Abstract syntax corresponds to describing the output of a parser (which are abstract syntax trees). In that view, regular expressions (as all notation represented by abstract syntax) denote trees. Since trees in texts are more difficult (and space-consuming) two write, one simply use the usual linear notation like the $b(\mathbf{ab})^*$ from above, with parentheses and “conventions” like precedences, to disambiguate the expression. Note that a tree representation represents the grouping of sub-expressions in its structure, so for grouping purposes, parentheses are not needed in abstract syntax.*

Of course, if one wants to implement a lexer or use one of the available ones, one has to deal with the particular concrete syntax of the particular scanner. There, of course, characters like '(' and ')' (or tokens like LPAREN or RPAREN) might occur.

Using concepts which will be discussed in more depth later, one may say: whether parentheses are considered as part of the syntax of regular expressions or not depends on the fact whether the definition is wished to be understood as describing concrete syntax trees or /abstract syntax trees!

See also Remark 5 later, which discusses further “ambiguities” in this context.

11. Same again once more

$$\begin{array}{ll}
 R \rightarrow \mathbf{a} \mid \epsilon \mid \emptyset & \text{basic reg. expr.} \\
 \mid R \mid R \mid RR \mid R^* \mid (R) & \text{compound reg. expr.}
 \end{array} \tag{6}$$

Note:

- symbol $|$: as symbol of regular expressions
- symbol $\mathbf{|}$: meta-symbol of the CF grammar notation
- The meta-notation use here for regular expressions will be the subject of later chapters

¹⁰To be careful, we will (later) distinguish between context-free languages on the one hand and notations to denote context-free languages on the other, in the same manner that we *now* don’t want to confuse regular languages as concept from particular notations (specifically, regular expressions) to write them down.

12. Semantics (meaning) of regular expressions

Definition 4 (Regular expression). Given an alphabet Σ . The meaning of a regexp r (written $\mathcal{L}(r)$) over Σ is given by equation (7).

$$\begin{array}{lll}
 \mathcal{L}(\emptyset) & = & \{\} & \text{empty language} \\
 \mathcal{L}(\epsilon) & = & \epsilon & \text{empty word} \\
 \mathcal{L}(a) & = & \{a\} & \text{single "letter" from } \Sigma \\
 \mathcal{L}(r \mid s) & = & \mathcal{L}(r) \cup \mathcal{L}(s) & \text{alternative} \\
 \mathcal{L}(r^*) & = & \mathcal{L}(r)^* & \text{iteration}
 \end{array} \tag{7}$$

- conventional *precedences*: $*$, concatenation, $|$.
- Note: left of “=”: reg-expr *syntax*, right of “=”: semantics/meaning/math ¹¹

13. Examples

In the following:

- $\Sigma = \{a, b, c\}$.
- we don't bother to “boldface” the syntax

words with exactly one b	$(a \mid c)^* b (a \mid c)^*$
words with max. one b	$((a \mid c)^* \mid ((a \mid c)^* b (a \mid c)^*))^*$
words of the form $a^n b a^n$, i.e., equal number of a 's before and after 1 b	$(a \mid c)^* (b \mid \epsilon) (a \mid c)^*$

14. Another regexp example

- (a) words that do *not* contain two b 's in a row.

$$\begin{array}{ll}
 (b (a \mid c))^* & \text{not quite there yet} \\
 ((a \mid c)^* \mid (b (a \mid c))^*)^* & \text{better, but still not there} \\
 & = \text{(simplify)} \\
 ((a \mid c) \mid (b (a \mid c)))^* & = \text{(simplify even more)} \\
 (a \mid c \mid ba \mid bc)^* & \\
 (a \mid c \mid ba \mid bc)^* (b \mid \epsilon) & \text{potential } b \text{ at the end} \\
 (notb \mid notb b)^* (b \mid \epsilon) & \text{where } notb \triangleq a \mid c
 \end{array}$$

- (b) Remarks

Remark 5 (Regular expressions, disambiguation, and associativity). *Note that in the equations in the example, silently allowed ourselves some “sloppiness” (at least for the nitpicking mind). The slight ambiguity depends on how we exactly interpret definition of regular expressions. Remember also Remark 4 on the preceding page, discussing the (non-)status of parentheses in regular expressions. If we think of Definition ?? on page ?? as describing abstract syntax and a concrete regular expression as representing an abstract syntax tree, then the constructor $|$ for alternatives is a binary constructor. Thus, the regular expression $a \mid c \mid ba \mid bc$ which occurs in the previous example is ambiguous. What is meant would be one of the following*

$$a \mid (c \mid (ba \mid bc)) \tag{8}$$

$$(a \mid c) \mid (ba \mid bc) \tag{9}$$

$$((a \mid c) \mid ba) \mid bc, \tag{10}$$

corresponding to 3 different trees, where occurrences of $|$ are inner nodes with two children each, i.e., subtrees representing subexpressions. In textbooks, one generally does not want to be bothered by writing all the parentheses. There are typically two ways to disambiguate the situation. One is to state (in the text) that the operator, in this case $|$, associates to the left (alternatively it associates to the right). That would mean that the “sloppy” expression without parentheses is meant to represent either (8) or (10), but not (9). If one really wants (9), one needs to indicate that using parentheses. Another way of finding an excuse for the sloppiness is to realize that it (in the context of regular expressions) does not matter, which of the three trees (8) – (10) is actually meant. This is specific for the setting here, where the symbol $|$ is semantically

¹¹Sometimes confusingly “the same” notation.

represented by set union \cup (cf. Definition 4 on the preceding page) which is an associative operation on sets. Note that, in principle, one may choose the first option —disambiguation via fixing an associativity— also in situations, where the operator is not semantically associative. As illustration, use the ‘-’ symbol with the usual intended meaning of “subtraction” or “one number minus another”. Obviously, the expression

$$5 - 3 - 1 \tag{11}$$

now can be interpreted in two semantically different ways, one representing the result 1, and the other 3. As said, one could introduce the convention (for instance) that the binary minus-operator associates to the left. In this case, (11) represents $(5 - 3) - 1$.

Whether or not in such a situation one wants symbols to be associative or not is a judgement call (a matter of language pragmatics). On the one hand, disambiguating may make expressions more readable by allowing to omit parenthesis or other syntactic markers which may make the expression or program look cumbersome. On the other, the “light-weight” and “programmer-friendly” syntax may trick the unexpecting programmer into misconceptions about what the program means, if unaware of the rules of associativity (and other priorities). Disambiguation via associativity rules and priorities is therefore a double-edged sword and should be used carefully. A situation where most would agree associativity is useful and completely unproblematic is the one illustrated for $|$ in regular expression: it does not matter anyhow semantically. Decisions concerning using a-priori ambiguous syntax plus rules how to disambiguate them (or forbid them, or warn the user) occur in many situations in the scanning and parsing phases of a compiler.

Now, the discussion concerning the “ambiguity” of the expression $(a | c | ba | bc)$ from equation (??) concentrated on the $|$ -construct. A similar discussion could obviously be made concerning concatenation (which actually here is not represented by a readable concatenation operator, but just by juxtaposition (=writing expressions side by side)). In the concrete example from (??), no ambiguity wrt. concatenation actually occurs, since expressions like ba are not ambiguous, but for longer sequences of concatenation like abc , the question of whether it means $a(bc)$ or $a(bc)$ arises (and again, it’s not critical, since concatenation is semantically associative).

Note also that one might think that the expression suffering from an ambiguity concerning combinations of operators, for instance, combinations of $|$ and concatenation. For instance, one may wonder if $\mathbf{ba | bc}$ could be interpreted as $(ba) | (bc)$ and $b(a | (bc))$ and $b(a | b)c$. However, on page 10, we stated precedences or priorities 4 on the previous page, stating that concatenation has a higher precedence over $|$, meaning that the correct interpretation is $(ba) | (bc)$. In a text-book the interpretation is “suggested” to the reader by the typesetting $ba | bc$ (and the notation it would be slightly less “helpful” if one would write $ba|bc\dots$ and what about the programmer’s version $\mathbf{a_b|a_c}$?). The situation with precedence is one where difference precedences lead to semantically different interpretations. Even if there’s a danger therefore that programmers/readers mis-interpret the real meaning (being unaware of precedences or mixing them up in their head), using precedences in the case of regular expressions certainly is helpful, The alternative of being forced to write, for instance

$$((a(b(cd))) | (b(a(ad)))) \text{ for } abcd | baad$$

is not even appealing to hard-core Lisp-programmers.

A final note: all this discussion about the status of parentheses or left or right associativity in the interpretation of (for instance mathematical) notation is mostly is over-the-top for most mathematics or other fields where some kind of formal notations or languages are used. There, notation is introduced, perhaps accompanied by sentences like “parentheses or similar will be used when helpful” or “we will allow ourselves to omit parentheses if no confusion may arise”, which means, the educated reader is expected to figure it out. Typically, thus, one glosses over too detailed syntactic conventions to proceed to the more interesting and challenging aspects of the subject matter. In such fields one is furthermore sometimes so used to notational traditions (“multiplication binds stronger than addition”), perhaps established since decades or even centuries, that one does not even think about them consciously. For scanner and parser designers, the situation is different; they are requested to come up with the notational (lexical and syntactical) conventions of perhaps a new language, specify them precisely and implement them efficiently. Not only that: at the same time, one aims at a good balance between explicitness (“Let’s just force the programmer to write all the parentheses and grouping explicitly, then he will get less misconceptions of what the program means (and the lexer/parser will be easy to write for me. . .)”) and economy in syntax, leaving many conventions, priorities, etc. implicit without confusing the target programmer.

15. Additional “user-friendly” notations

$$\begin{aligned} r^+ &= rr^* \\ r^? &= r | \epsilon \end{aligned}$$

Special notations for *sets* of letters:

- $[0 - 9]$ range (for ordered alphabets)
- $\sim a$ not a (everything except a)
- \cdot all of Σ

naming regular expressions (“regular definitions”)

- $digit = [0 - 9]$
- $nat = digit^+$
- $signedNat = (+|-)nat$
- $number = signedNat(“.” nat)?(E signedNat)?$

1.3 DFA

1. Finite-state automata

- simple “computational” machine
- (variations of) FSA’s exist in many flavors and under different names
- other rather well-known names include finite-state machines, finite labelled transition systems,
- “state-and-transition” representations of programs or behaviors (finite state or else) are wide-spread as well
 - state diagrams
 - Kripke-structures
 - I/O automata
 - Moore & Mealy machines
- the logical behavior of certain classes of electronic circuitry with internal memory (“flip-flops”) is described by finite-state automata.¹²

(a) Remarks

Remark 6 (Finite states). *The distinguishing feature of FSA (as opposed to more powerful automata models such as push-down automata, or Turing-machines), is that they have “finitely many states”. That sounds clear enough at first sight. But one has to be a bit more careful. First of all, the set of states of the automaton, here called Q , is finite and fixed for a given automaton, all right. But actually, the same is true for pushdown automata and Turing machines! The trick is: if we look at the illustration of the finite-state automaton earlier, where the automaton had a head. The picture corresponds to an accepting use of an automaton, namely one that is fed by letters on the tape, moving internally from one state to another, as controlled by the different letters (and the automaton’s internal “logic”, i.e., transitions). Compared to the full power of Turing machines, there are two restrictions, things that a finite state automaton cannot do*

- *it moves on one direction only (left-to-right)*
- *it is read-only.*

All non-finite state machines have some additional memory they can use (besides $q_0, \dots, q_n \in Q$). Push-down automata for example have additionally a stack, a Turing machine is allowed to write on the tape, thus using it as external memory.

2. FSA

Definition 5 (FSA). A FSA \mathcal{A} over an alphabet Σ is a tuple $(\Sigma, Q, I, F, \delta)$

- Q : finite set of states
- $I \subseteq Q, F \subseteq Q$: initial and final states.
- $\delta \subseteq Q \times \Sigma \times Q$ transition relation
- final states: also called *accepting* states
- transition relation: can *equivalently* be seen as function $\delta : Q \times \Sigma \rightarrow 2^Q$: for each state and for each letter, give back the *set* of successor states (which may be empty)
- more suggestive notation: $q_1 \xrightarrow{a} q_2$ for $(q_1, a, q_2) \in \delta$
- We also use freely —self-evident, we hope— things like

$$q_1 \xrightarrow{a} q_2 \xrightarrow{b} q_3$$

¹²Historically, design of electronic circuitry (not yet chip-based, though) was one of the early very important applications of finite-state machines.

3. FSA as scanning machine?

- FSA have slightly unpleasant properties when considering them as describing an actual program (i.e., a scanner procedure/lexer)
- given the “theoretical definition” of acceptance:

Mental picture of a scanning automaton

The automaton eats one character after the other, and, when reading a letter, it moves to a successor state, if any, of the current state, depending on the character at hand.

- 2 problematic aspects of FSA
 - **non-determinism**: what if there is more than one possible successor state?
 - **undefinedness**: what happens if there’s no next state for a given input
- the second one is *easily* repaired, the first one requires more thought

(a) Non-determinism

Sure, one could try **backtracking**, but, trust us, you don’t want that in a scanner. And even if you think it’s worth a shot: how do you scan a program directly from magnetic tape, as done in the bad old days? Magnetic tapes can be rewound, of course, but winding them back and forth all the time destroys hardware quickly. How should one scan network traffic, packets etc. on the fly? The network definitely cannot be rewound. Of course, buffering the traffic would be an option and doing then backtracking using the buffered traffic, but maybe the packet-scanning-and-filtering should be done in hardware/firmware, to keep up with today’s enormous traffic bandwidth. Hardware-only solutions have no dynamic memory, and therefore actually *are* ultimately finite-state machine with no extra memory.

4. DFA: deterministic automata

Definition 6 (DFA). A *deterministic, finite automaton* \mathcal{A} (DFA for short) over an alphabet Σ is a tuple $(\Sigma, Q, I, F, \delta)$

- Q : finite set of states
- $I = \{i\} \subseteq Q, F \subseteq Q$: initial and final states.
- $\delta : Q \times \Sigma \rightarrow Q$ transition function
- transition function: special case of transition relation:
 - deterministic
 - left-total¹³

5. Meaning of an FSA

Semantics

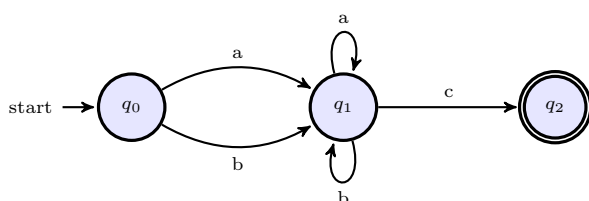
The intended **meaning** of an FSA over an alphabet Σ is the set consisting of all the finite words, the automaton **accepts**.

Definition 7 (Accepting words and language of an automaton). A word $c_1c_2\dots c_n$ with $c_i \in \Sigma$ is *accepted* by automaton \mathcal{A} over Σ , if there exists states q_0, q_2, \dots, q_n all from Q such that

$$q_0 \xrightarrow{c_1} q_1 \xrightarrow{c_2} q_2 \xrightarrow{c_3} \dots q_{n-1} \xrightarrow{c_n} q_n ,$$

and were $q_0 \in I$ and $q_n \in F$. The *language* of an FSA \mathcal{A} , written $\mathcal{L}(\mathcal{A})$, is the set of all words \mathcal{A} accepts

6. FSA example

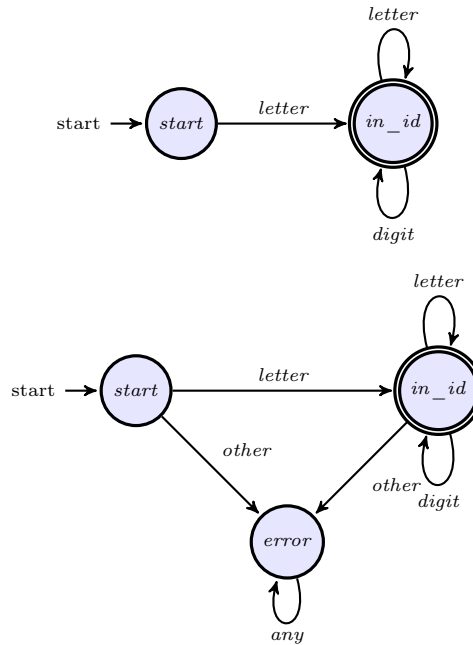


¹³That means, for each pair q, a from $Q \times \Sigma$, $\delta(q, a)$ is defined. Some people call an automaton where δ is not a left-total but a deterministic relation (or, equivalently, the function δ is not total, but partial) still a deterministic automaton. In that terminology, the DFA as defined here would be deterministic *and* total.

7. Example: identifiers

Regular expression

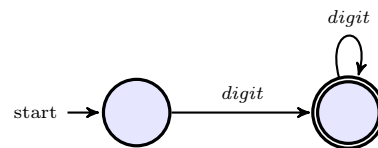
$$\text{identifier} = \text{letter}(\text{letter} \mid \text{digit})^* \tag{12}$$



- transition *function*/relation δ *not* completely defined (= *partial* function)

8. Automata for numbers: natural numbers

$$\begin{aligned} \text{digit} &= [0 - 9] \\ \text{nat} &= \text{digit}^+ \end{aligned} \tag{13}$$

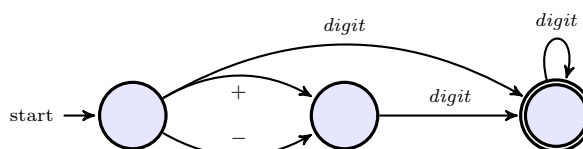


(a) Remarks

One might say, it's not really the natural numbers, it's about a decimal *notation* of natural numbers (as opposed to other notations, for example Roman numeral notation). Note also that initial zeroes are allowed here. It would be easy to disallow that.

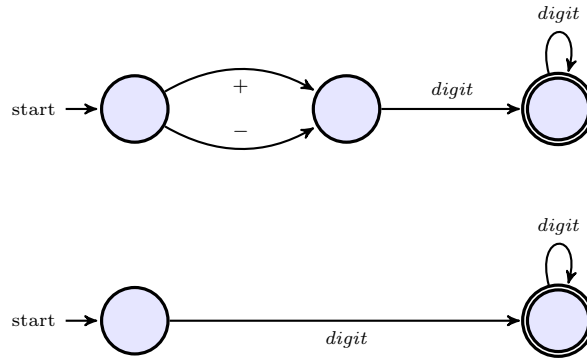
9. Signed natural numbers

$$\text{signednat} = (+ \mid -)\text{nat} \mid \text{nat} \tag{14}$$



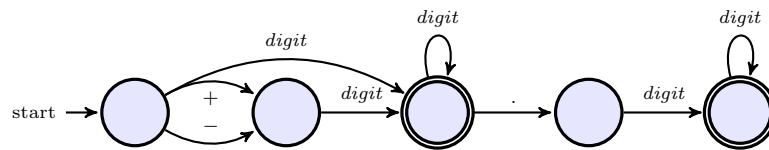
- (a) Remarks Again, the automaton is *deterministic*. It's easy enough to come up with this automaton, but the *non-deterministic* one is probably more straightforward to come by with. Basically, one informally does two “constructions”, the “alternative” which is simply writing “two automata”, i.e., one automaton which consists of the union of the two automata, basically. In this example, it therefore has two initial states (which is disallowed obviously for deterministic automata). Another implicit construction is the “*sequential composition*”.

10. Signed natural numbers: non-deterministic



11. Fractional numbers

$$frac = signednat("." nat)? \tag{15}$$



12. Floats

$$\begin{aligned} digit &= [0 - 9] \\ nat &= digit^+ \\ signednat &= (+ | -)nat | nat \\ frac &= signednat("." nat)? \\ float &= frac(E signednat)? \end{aligned} \tag{16}$$

- Note: no (explicit) recursion in the definitions
- note also the treatment of *digit* in the automata.

(a) Recursion

Remark 7 (Recursion). [Louden, 1997] points out that regular expressions do not contain recursion. This is for instance stated at the beginning of Chapter 3 in [Louden, 1997] where the absence of recursion in regular expressions is pointed out as the main distinguishing feature of regular expressions compared to context-free grammars (or even more expressive grammars of formalisms, see later).

While considering regular expressions as being “without recursion”, not everyone would agree. Looking at the defining equations in (16), the series of equations “culminates” in the one for floats, the last one listed. Furthermore, each equation makes use on its right-hand side only of definitions defined strictly before that equation (“strict” means, that a category defined in the left-hand side equation may also not depend directly on itself by mentioning the category being defined on the defining right-hand side). In that sense, the definition clearly is without recursion, and for context-free grammars, that restriction will not longer apply. This absence of at least explicit recursion when defining a regular expression for floats allows that one can consider the definitions as given as simply useful “abbreviations” to assist the reader’s or designer’s

understanding of the definition. They therefore play the role of macro definitions as for instance supported by C-style preprocessors: the “real” regular expression can easily be obtained by literally replacing the “macro names” as they appear in some right-hand sides by their definition, until all of them have disappeared and one has reached the “real” regular expression, using only the syntax supported by the original, concise definition of regular expressions.¹⁴ Textual replacement, by the way, is also the way, pre-processors deal with macro definitions. Clearly this easy way of replacing mentioning of left-hand sides by their corresponding right-hand sides works only in absence of recursive definitions.

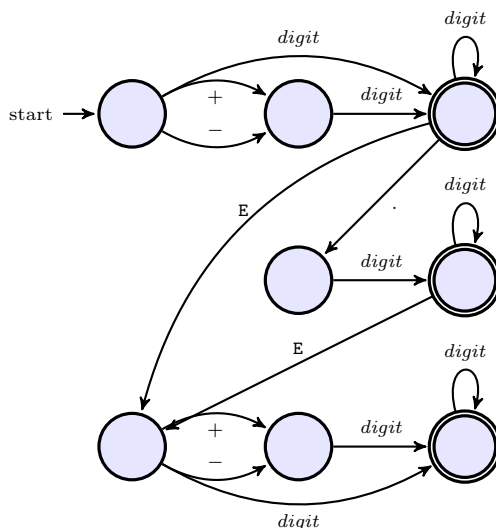
That was the argument supporting the case that regular expressions don’t contain recursion. There is, however, a different angle to approach the issue. Recursion, very generally, is a way to describe infinite structures, behavior etc. Regular languages are, in general, infinite, i.e., infinite sets of words, and the way to capture those infinite sets in a finite way is via the Kleene star. In the automata, infinitely different words are represented by “loops”. Thus, the Kleene star does allow to express (implicitly) a form of recursion, even if it’s a more restricted form than that allowed by context-free grammars. The point may become more clear if we replace the definition for natural numbers from equation (13), using + for “one or more iterations” by the following recursive one:

$$\text{nat} = \text{digit nat} \mid \text{digit} . \quad (17)$$

Compared to the more general definitions for context-free grammars later, the recursive mentioning of nat on the right-hand side of the definition for regular language is restricted. The restriction will become clearer once we have covered context-free grammars in the context of parsing. Suffices for now, that the restrictions are called left-linear grammars (alternatively right-linear grammars, both are equally expressive), where linear refers to the fact that at most one of the “meta-variables” in a grammar (such as nat from above) allowed to occur on the right-hand side of a rule, and right-linear would mean, it’s allowed to occur only at the end of a right-hand side.¹⁵ Another, basically equivalent, word is that the definitions may use tail-recursion (but not general recursion).

To summarize: the dividing line between regular languages vs. context-free languages may well be described as allowing tail-recursion vs. general recursion (not as without or with recursion as in [Louden, 1997]). Finally, sometimes, when programming, one distinguishes between iteration (when using a loop construct) on the one hand and recursion on the other, where iteration corresponds to a restricted form of recursion, namely tail-recursion: tail-recursion is a form of recursion, where no stack is needed (and which therefore here can be handled by finite-state automata), in contrast to context-free grammars, which cannot be handled by FSA’s, one needs equip them with a stack, after which they are called push-down automata, most oftenly.

13. DFA for floats

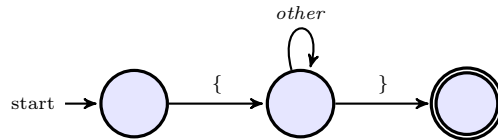


14. DFAs for comments

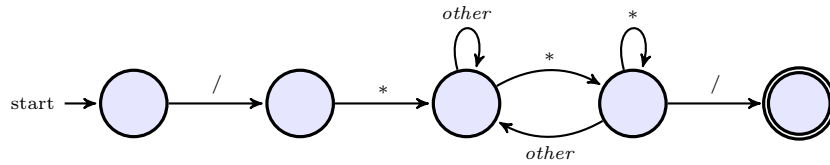
Pascal-style

¹⁴Additional syntactic material that is added for the convenience of the programmer *without* adding expressivity of the language and which can *easily* be “expanded way” is also known as syntactic sugar.

¹⁵As a fine point, to avoid confusion later: The definition from equation (17) would count as *two* rules in a grammar, not one, corresponding to the two alternatives. The restrictions for linearity etc. apply *per rule/branch* individually.

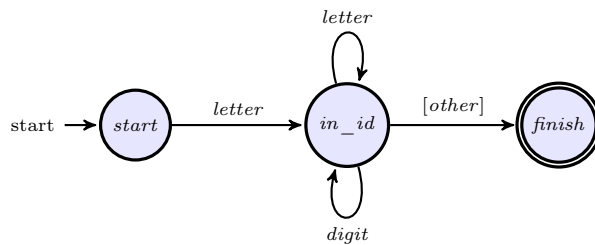


C, C++, Java



1.4 Implementation of DFA

1. Implementation of DFA (1)



2. Implementation of DFA (1): "code"

```

{ starting state }

if the next character is a letter
then
  advance the input;
  { now in state 2 }
  while the next character is a letter or digit
  do
    advance the input;
    { stay in state 2 }
  end while;
  { go to state 3, without advancing input }
  accept;
else
  { error or other cases }
end
  
```

3. Explicit state representation

```

state := 1 { start }
while state = 1 or 2
do
  case state of
  1: case input character of
      letter: advance the input;
          state := 2
      else state := .... { error or other };
      end case;
  2: case input character of
  
```

```

letter , digit: advance the input;
                    state := 2; { actually unnessessary }
else
    state := 3;
end case;
end case;
end while;
if state = 3 then accept else error;

```

4. Table representation of a DFA

state \ input char	letter	digit	other
1	2		
2	2	2	3
3			

state \ input char	letter	digit	other	accepting
1	2			no
2	2	2	[3]	no
3				yes

5. Table-based implementation

```

state := 1 { start }
ch := next input character;
while not Accept[state] and not error(state)
do
while state = 1 or 2
do
    newstate := T[state, ch];
    {if Advance[state, ch]
    then ch:=next input character};
    state := newstate
end while;
if Accept [state] then accept;

```

1.5 NFA

1. Non-deterministic FSA

Definition 8 (NFA (with ϵ transitions)). A *non-deterministic* finite-state automaton (NFA for short) \mathcal{A} over an alphabet Σ is a tuple $(\Sigma, Q, I, F, \delta)$, where

- Q : finite set of states
- $I \subseteq Q, F \subseteq Q$: initial and final states.
- $\delta : Q \times \Sigma \rightarrow 2^Q$ transition function

In case, one uses the alphabet $\Sigma + \{\epsilon\}$, one speaks about an NFA with ϵ -transitions.

- in the following: NFA mostly means, allowing ϵ transitions¹⁶
- ϵ : treated *differently* than the “normal” letters from Σ .
- δ can *equivalently* be interpreted as *relation*: $\delta \subseteq Q \times \Sigma \times Q$ (transition relation labelled by elements from Σ).

(a) Finite state machines

¹⁶It does not matter much anyhow, as we will see.

Remark 8 (Terminology (finite state automata)). *There are slight variations in the definition of (deterministic resp. non-deterministic) finite-state automata. For instance, some definitions for non-deterministic automata might not use ϵ transitions, i.e., defined over Σ , not over $\Sigma + \{\epsilon\}$. Another word for FSAs are finite-state machines. Chapter 2 in [Louden, 1997] builds in ϵ -transitions into the definition of NFA, whereas in Definition 8, we mention that the NFA is not just non-deterministic, but “also” allows those specific transitions. Of course, ϵ -transitions lead to non-determinism as well, in that they correspond to “spontaneous” transitions, not triggered and determined by input. Thus, in the presence of ϵ -transition, and starting at a given state, a fixed input may not determine in which state the automaton ends up in.*

Deterministic or non-deterministic FSA (and many, many variations and extensions thereof) are widely used, not only for scanning. When discussing scanning, ϵ transitions come in handy, when translating regular expressions to FSA, that’s why [Louden, 1997] directly builds them in.

2. Language of an NFA

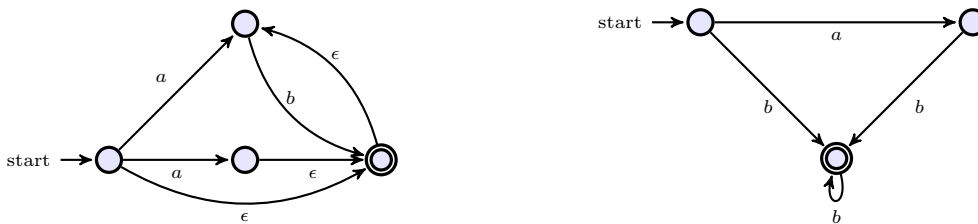
- Remember $\mathcal{L}(\mathcal{A})$ (Definition 7 on page 13)
- applying definition directly to $\Sigma + \{\epsilon\}$: accepting words “containing” letters ϵ
- as said: special treatment for ϵ -transitions/ ϵ -“letters”. ϵ rather represents **absence** of input character/letter.

Definition 9 (Acceptance with ϵ -transitions). A word w over alphabet Σ is *accepted* by an NFA with ϵ -transitions, if there exists a word w' which is accepted by the NFA with alphabet $\Sigma + \{\epsilon\}$ according to Definition 7 and where w is w' with all occurrences of ϵ **removed**.

- (a) Alternative (but equivalent) intuition \mathcal{A} reads one character after the other (following its transition relation). If in a state with an outgoing ϵ -transition, \mathcal{A} can move to a corresponding successor state *without* reading an input symbol.

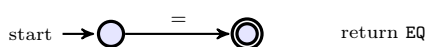
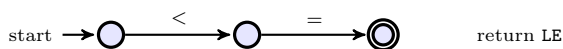
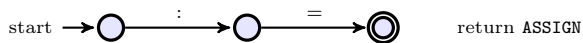
3. NFA vs. DFA

- *NFA*: often easier (and smaller) to write down, esp. starting from a reg expression.
- Non-determinism: not *immediately* transferable to an *algo*

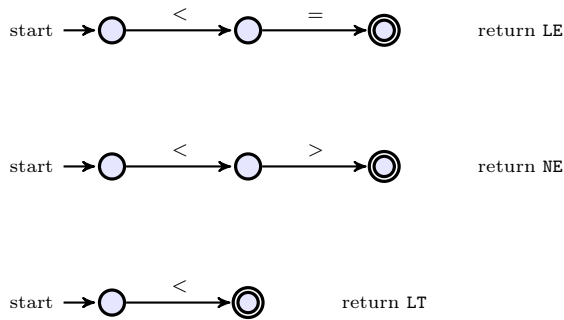


1.6 From regular expressions to DFAs

1. Why non-deterministic FSA? Task: recognize $:=$, $<=$, and $=$ as three different tokens:



2. What about the following 3 tokens?



1.7 Thompson's construction

1. Regular expressions \rightarrow NFA

- needed: a *systematic* translation
 - conceptually easiest: translate to DFA
 - postpone determinization for a second step
 - (postpone minimization for later as well)
- (a) Compositional construction [Thompson, 1968] Design goal: The NFA of a compound regular expression is given by taking the NFA of the immediate subexpressions and connecting them appropriately.
- construction slightly¹⁷ simpler, if one uses automata with **one** start and one accepting state
- \Rightarrow ample use of ϵ -transitions

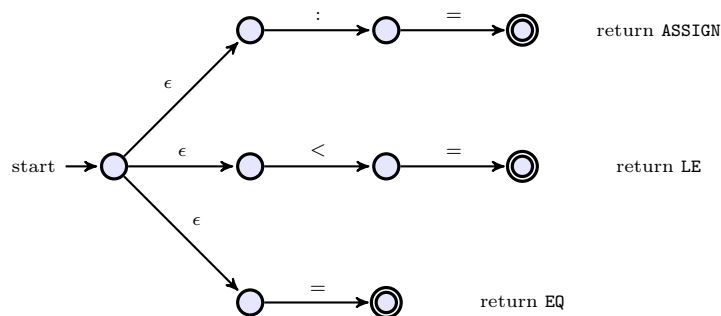
(b) Compositionality

Remark 9 (Compositionality). *Compositional concepts (definitions, constructions, analyses, translations ...) are immensely important and pervasive in compiler techniques (and beyond). One example already encountered was the definition of the language of a regular expression (see Definition 4 on page 10). The design goal of a compositional translation here is the underlying reason why to base the construction on non-deterministic machines.*

Compositionality is also of practical importance (“component-based software”). In connection with compilers, separate compilation and (static / dynamic) linking (i.e. “composing”) of separately compiled “units” of code is a crucial feature of modern programming languages/compilers. Separately compilable units may vary, sometimes they are called modules or similarly. Part of the success of C was its support for separate compilation (and tools like make that helps organizing the (re-)compilation process). For fairness sake, C was by far not the first major language supporting separate compilation, for instance FORTRAN II allowed that, as well, back in 1958.

Btw., Ken Thompson, the guy who first described the regex-to-NFA construction discussed here, is one of the key figures behind the UNIX operating system and thus also the C language (both went hand in hand). Not suprisingly, considering the material of this section, he is also the author of the grep -tool (“globally search a regular expression and print”). He got the Turing-award (and many other honors) for his contributions.

2. Illustration for ϵ -transitions

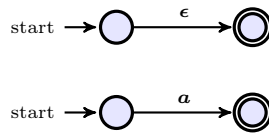


¹⁷does not matter much, though.

3. Thompson's construction: basic expressions

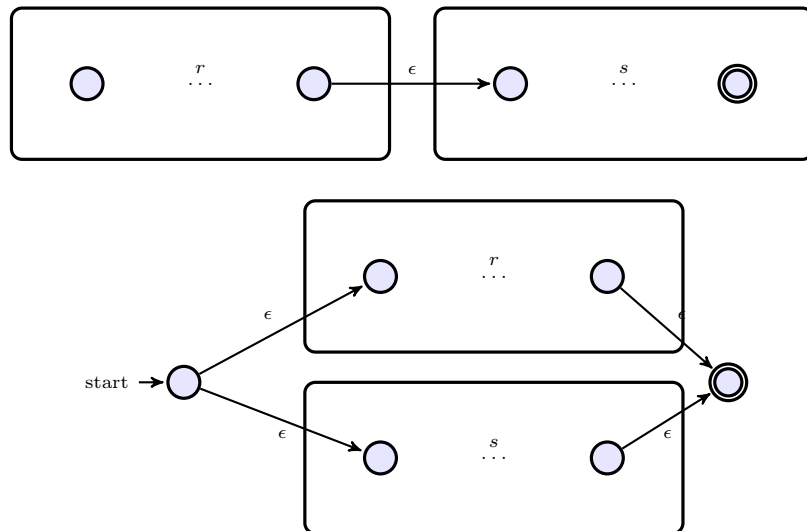
basic regular expressions

basic (= non-composed) regular expressions: ϵ , \emptyset , a (for all $a \in \Sigma$)

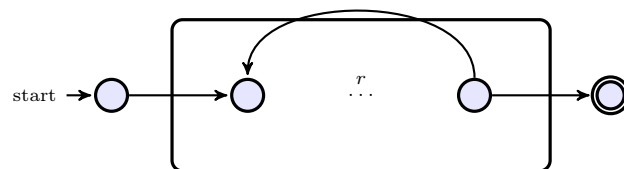


(a) Remarks The \emptyset is slightly odd: it's sometimes not part of regular expressions. If it's lacking, then one cannot express the empty language, obviously. That's not nice, because then the regular languages are not closed under complement. Also: obviously there exists an automaton with an empty language. Therefore, \emptyset should be part of the regular expressions, even it practically it does not play much of a role.

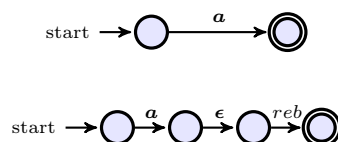
4. Thompson's construction: compound expressions

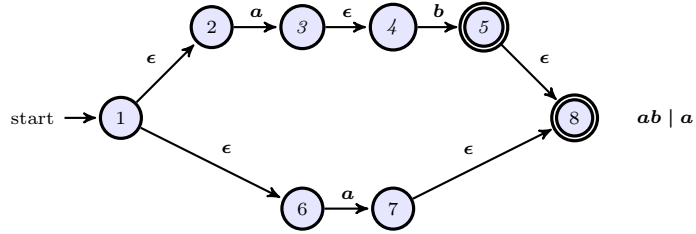


5. Thompson's construction: compound expressions: iteration



6. Example





1.8 Determinization

1. Determinization: the subset construction

(a) Main idea

- Given a non-det. automaton \mathcal{A} . To construct a DFA \mathcal{A}' : instead of *backtracking*: explore all successors “at the same time” \Rightarrow
- each state q' in $\overline{\mathcal{A}}$: represents a *subset* of states from \mathcal{A}
- Given a word w : “feeding” that to $\overline{\mathcal{A}}$ leads to *the* state representing *all* states of \mathcal{A} /reachable/via w .

(b) Items

- side remark: this construction, known also as *powerset* construction, seems straightforward enough, but: analogous constructions works for some other kinds of automata, as well, but for others, the approach does *not* work.¹⁸
- Origin [Rabin and Scott, 1959]

2. Some notation/definitions

Definition 10 (ϵ -closure, a -successors). Given a state q , the ϵ -closure of q , written $close_\epsilon(a)$, is the set of states reachable via zero, one, or more ϵ -transitions. We write q_a for the set of states, reachable from q with one a -transition. Both definitions are used analogously for sets of states.

(a) *emptyword*-closure

Remark 10 (ϵ -closure). [Louden, 1997] does not sketch an algorithm but it should be clear that the ϵ -closure is easily compilable for a given state, resp. a given finite set of states. Some textbooks also denote λ instead of ϵ , and consequently speak of λ -closure. And in still other contexts (mainly not in language theory and recognizers), silent transitions are marked with τ .

It may be obvious but: the set of states in the ϵ -closure of a given state are not “language-equivalent”. However, the union of languages for all states from the ϵ -closure corresponds to the language accepted with the given state as initial one. However, the language being accepted is not the property which is relevant here in the determinization. The ϵ -closure is needed to capture the set of all states reachable by a given word. But again, the exact characterization of the set need to be done carefully. The states in the set are also not equivalent wrt. their reachability information: Obviously, states in the ϵ -closure of a given state may be reached by more words. The set of reaching words for a given state, however, is not in general the intersection of the sets of corresponding words of the states in the closure.

3. Transformation process: sketch of the algo **Input:** NFA \mathcal{A} over a given Σ

Output: DFA $\overline{\mathcal{A}}$

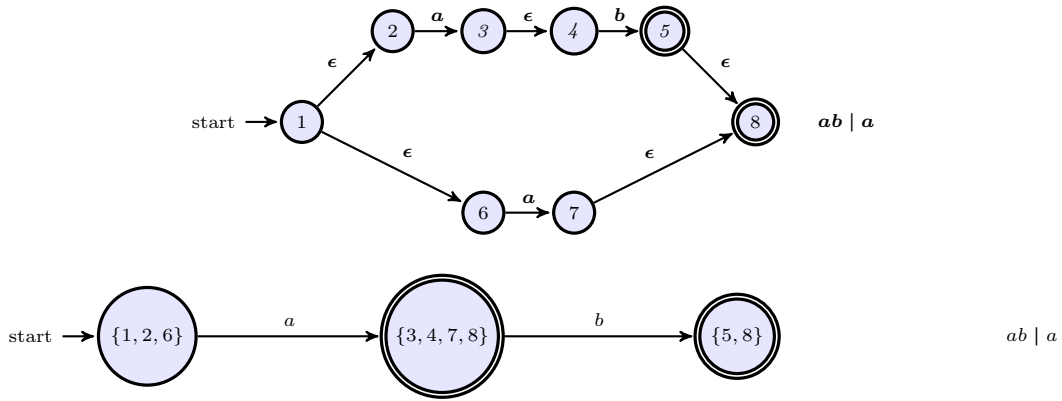
- the *initial* state: $close_\epsilon(I)$, where I are the initial states of $\overline{\mathcal{A}}$
- for a state Q' in $\overline{\mathcal{A}}$: the a -successor of Q is given by $close_\epsilon(Q_a)$, i.e.,

$$Q \xrightarrow{a} close_\epsilon(Q_a) \tag{18}$$

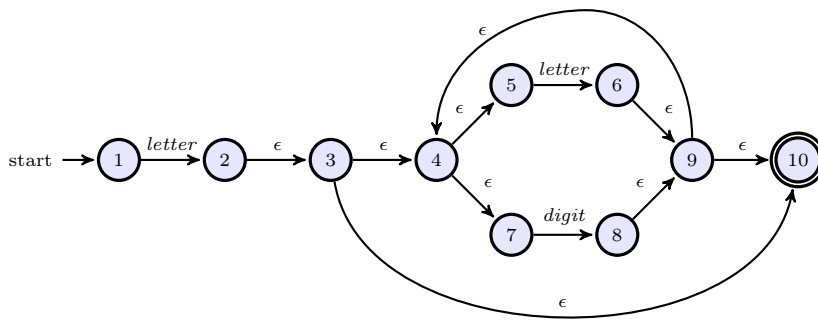
- repeat step 2 for all states in $\overline{\mathcal{A}}$ and all $a \in \Sigma$, until no more states are being added
- the *accepting* states in $\overline{\mathcal{A}}$: those containing *at least one* accepting states of \mathcal{A} .

¹⁸For some forms of automata, non-deterministic versions are strictly more expressive than the deterministic one.

4. Example $ab \mid a$



5. Example: identifiers Remember regexpr from equation (12)



6. Minimization

- automatic construction of DFA (via e.g. Thompson): often many superfluous states
- goal: “combine” states of a DFA without changing the accepted language

(a) Properties of the minimization algo

Canonicity: all DFA for the same language are transformed to the *same* DFA

Minimality: resulting DFA has *minimal* number of states

(b) Itemize

- “side effects”: answers to *equivalence* problems
 - given 2 DFA: do they accept the same language?
 - given 2 regular expressions, do they describe the same language?
- modern version: [Hopcroft, 1971].

7. Hopcroft’s partition refinement algo for minimization

- starting point: *complete* DFA (i.e., *error-state* possibly needed)
- first idea: *equivalent* states in the given DFA may be *identified*
- **equivalent:** when used as starting point, accepting the same language
- **partition refinement:**
 - works “the other way around”
 - instead of collapsing equivalent states:
 - * start by “collapsing as much as possible” and then,
 - * iteratively, detect *non-equivalent* states, and then *split* a “collapsed” state

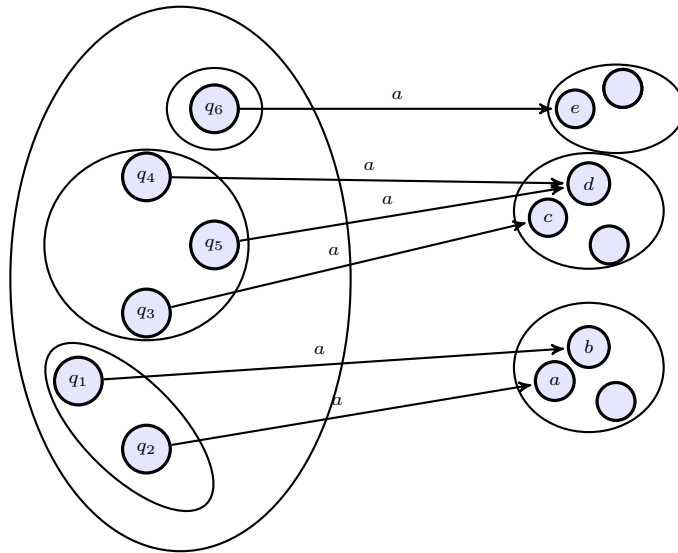
* stop when no violations of “equivalence” are detected

- *partitioning* of a set (of states):
- *worklist*: data structure of to keep non-treated classes, termination if worklist is empty

8. Partition refinement: a bit more concrete

- **Initial** partitioning: 2 partitions: set containing all *accepting* states F , set containing all *non-accepting* states $Q \setminus F$
- **Loop** do the following: pick a current equivalence class Q_i and a symbol a
 - if for all $q \in Q_i$, $\delta(q, a)$ is member of the *same* class $Q_j \Rightarrow$ consider Q_i as done (for now)
 - else:
 - * **split** Q_i into Q_i^1, \dots, Q_i^k s.t. the above situation is repaired for each Q_i^l (but don't split more than necessary).
 - * be aware: a split may have a “cascading effect”: other classes being fine before the split of Q_i need to be reconsidered \Rightarrow *worklist* algo
- **stop** if the situation stabilizes, i.e., no more split happens (= worklist empty, at latest if back to the original DFA)

9. Split in partition refinement: basic step

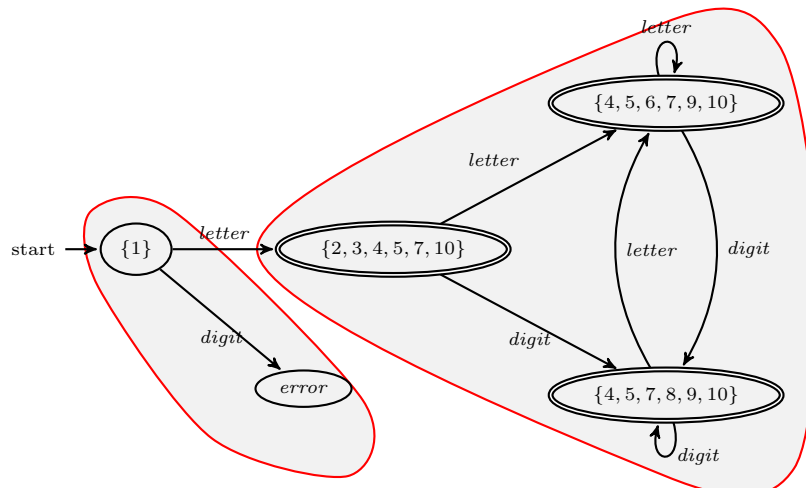


- before the split $\{q_1, q_2, \dots, q_6\}$
- after the split $\{q_1, q_2\}, \{q_3, q_4, q_5\}, \{q_6\}$

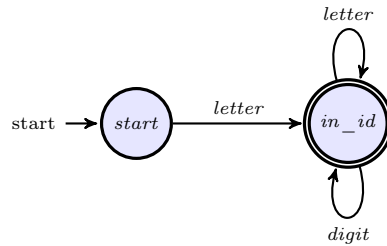
(a) Note

The pic shows only one letter a , in general one has to do the same construction for all letters of the alphabet.

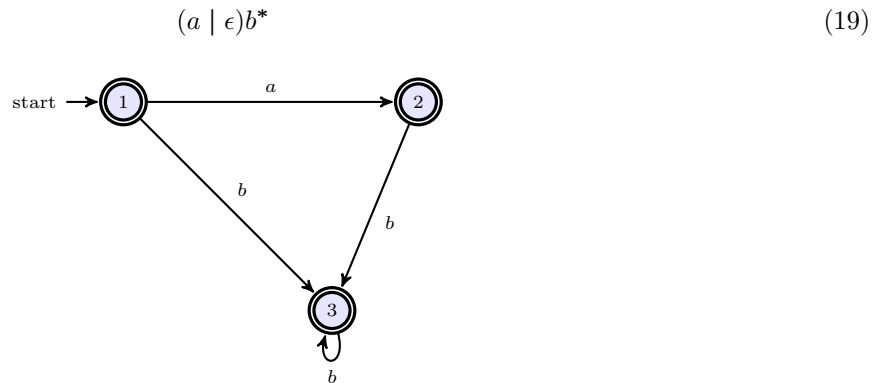
10. Completed automaton



11. Minimized automaton (error state omitted)



12. Another example: partition refinement & error state



1.9 Scanner generation tools

1. Tools for generating scanners

- scanners: simple and well-understood part of compiler
- hand-coding possible
- mostly better off with: generated scanner
- standard tools **lex** / **flex** (also in combination with *parser* generators, like **yacc** / **bison**)
- variants exist for many implementing languages
- based on the results of this section

2. Main idea of (f)lex and similar

- output of lexer/scanner = input for parser
- programmer specifies regular expressions for each **token**-class and corresponding actions¹⁹ (and whitespace, comments etc.)
- the spec.\ language offers some conveniences (extended regexpr with priorities, associativities etc) to ease the task
- automatically translated to NFA (e.g. Thompson)
- then made into a deterministic DFA (“subset construction”)
- minimized (with a little care to keep the token classes separate)
- implement the DFA (usually with the help of a *table* representation)

¹⁹Tokens and actions of a parser will be covered later. For example, identifiers and digits as described but the reg. expressions, would end up in two different token classes, where the actual string of characters (also known as *lexeme*) being the value of the token attribute.

References

- [Hopcroft, 1971] Hopcroft, J. E. (1971). An $n \log n$ algorithm for minimizing the states in a finite automaton. In Kohavi, Z., editor, *The Theory of Machines and Computations*, pages 189–196. Academic Press, New York.
- [Kleene, 1956] Kleene, S. C. (1956). Representation of events in nerve nets and finite automata. In *Automata Studies*, pages 3–42. Princeton University Press.
- [Louden, 1997] Loudon, K. (1997). *Compiler Construction, Principles and Practice*. PWS Publishing.
- [Rabin and Scott, 1959] Rabin, M. and Scott, D. (1959). Finite automata and their decision problems. *IBM Journal of Research Developments*, 3:114–125.
- [Thompson, 1968] Thompson, K. (1968). Programming techniques: Regular expression search algorithm. *Communications of the ACM*, 11(6):419.

Index

- Σ , 6
- $\mathcal{L}(r)$ (language of r), 10
- \emptyset -closure, 22
- ϵ , 19
- ϵ (empty word), 18
- ϵ transition, 18
- ϵ -closure, 22
- ϵ -transition, 20

- accepting state, 13
- alphabet, 6
 - ordered, 11
- automaton
 - accepting, 13
 - language, 13
 - push down, 16
 - semantics, 13

- bison, 25
- blank character, 2

- character, 2
- classification, 4
- comment, 16
- compiler compiler, 5
- compositionality, 20
- context-free grammar, 8, 9

- determinization, 22
- DFA, 1
 - definition, 13
- digit, 14
- disk head, 2

- encoding, 2

- final state, 13
- finite state machine, 19
- flex, 25
- floating point numbers, 16
- Fortran, 3
- Fortran, 3
- FSA, 1, 12
 - definition, 12
 - scanner, 13
 - semantics, 13

- grammar
 - left-linear, 16
 - right-linear, 16
- grep, 20

- Hopcroft's partition refinement algorithm, 24

- I/O automaton, 12
- identifier, 2, 4
- inite-state automaton, 12
- initial state, 13
- irrational number, 7

- keyword, 2, 4
- Kripke structure, 12

- labelled transition system, 12
- language, 6
 - of an automaton, 13
- left-linear grammar, 16
- letter, 6
- lex, 25
- lexem
 - and token, 4
- lexeme, 25
- lexer, 1
 - classification, 4
- lexical scanner, 1

- macro definition, 16
- Mealy machine, 12
- meaning, 13
- minimization, 23
- Moore machine, 12

- NFA, 1, 18
 - language, 19
- non-determinism, 13
- non-deterministic FSA, 18
- number
 - floating point, 15
 - fractional, 15
- numeric costants, 4

- parser generator, 5, 25
- partition refinement, 24
- partitioning, 24
- powerset construction, 22
- pragmatics, 4, 11
- pre-processor, 16
- priority, 5
- push-down automaton, 16

- rational language, 8
- rational number, 7
- recursion, 16
 - tail, 16
- regular definition, 11
- regular expression, 1, 5
 - language, 10
 - meaning, 10
 - named, 11
 - precedence, 10
 - recursion, 15
 - semanticsx, 10
 - syntax, 10
- regular expressions, 8
- reserved word, 2, 4
- right-linear grammar, 16

- scanner, 1
- scanner generator, 25
- screener, 4
- semantics, 13
- separate compilation, 20
- stack, 16
- state diagram, 12

- string literal, 4
- subset construction, 22
- successor state, 13
- symbol, 6
- symbol table, 6
- symbols, 6
- syntactic sugar, 16

- tail recursion, 16
- Thompon's construction, 20
- token, 4, 25
- tokenizer, 1
- transition function, 13
- transition relation, 13
- Turing machine, 2

- undefinedness, 13

- whitespace, 2, 4
- word, 6
 - vs. string, 6
- worklist, 24

- yacc, 25