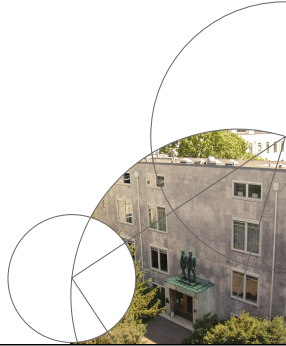
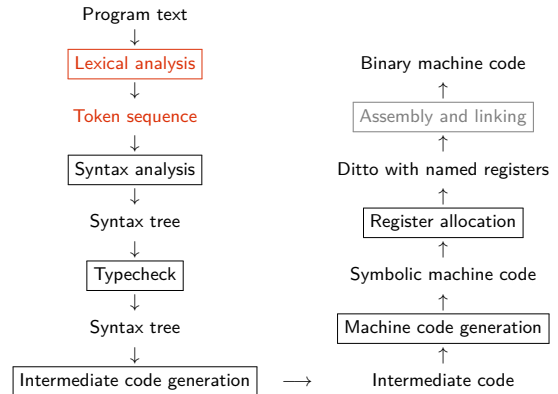


# Compilers (Oversættere): Lexical Analysis

Jost Berthold  
berthold@diku.dk  
Department of Computer Science



## Lexical Analysis (Scanning): First phase of Frontend



## Lexical

of or relating to words or the vocabulary of a language as distinguished from its grammar and construction.

Merriam-Webster Dictionary, m-w.com

- Lexical analyser, also called lexer, scanner or tokeniser.
- splits input (stream of characters) into tokens.
- Token: smallest meaningful unit for a programming language  
Keyword, number, comment, parenthesis, semicolon, ...
- Tokens are classes of concrete input (called lexeme).

## Contents and Goals of this Part

- 1 What is Lexical Analysis?  
Regular Expressions and Languages  
A Tool for Lexical Analysis
- 2 Finite Automata  
Non-deterministic and Deterministic Automata (NFA and DFA)  
Converting and Minimising Automata
- 3 Automata Construction for Lexical Analysis

### Goals:

- Understand regular expressions and the concept of formal languages, and apply them for lexical analysis
- Use scanner generators for lexical analysis
- Construct, convert, and minimise finite automata
- Know automata limitations and use them as arguments

## Lexical analysis – Example

```

An SML program
(* My first
 * SML program.
 *)
val result
= let val x = 10 :: 020 :: 0x30 :: []
  in List.map (fn x => x div 2) x
  end
    
```

Converting textual input into a token sequence

- Input file read as a string
- Contains comments (and preprocessor directives)
- Easy to read with layout/colours.

```

...read by the lexer
(*My first\n.*SML program.\n.*\nval result\n=let val x=10:::020
:::0x30:::[]\nin List.map(fn x=>x div 2)x\nend\n
    
```

- Machine-read. Formatting only helps human readers.
- Lexical analysis: character stream to token sequence

## Lexical analysis – Example

```

An SML program
(* My first
 * SML program.
 *)
val result
= let val x = 10 :: 020 :: 0x30 :: []
  in List.map (fn x => x div 2) x
  end
    
```

Converting textual input into a token sequence

- Input file read as a string
- Contains comments (and preprocessor directives)
- Split into token sequence for machine processing

```

resulting token sequence
[Keyword "val", Id "result", Equal, Keyword "let", Keyword "val", Id "x",
Equal, Int 10, Doublecolon, Int 20, Doublecolon, Int 48, Doublecolon,
LBracket, RBracket, Keyword "in", Id "List", Dot, Id "map", LParen,...
    
```

- Tokens: classes of input lexemes.
- No comments or formatting (only position for error messages)
- Sometimes: special processing instructions (including files).

## Describing Tokens

What are the typical tokens of a programming language?

- Separators and (various) parentheses: ; , { [ ( ) ] } ...
- String and number literals: "Hello", 10, 0x30, 3.14, 1.2e-12, ...
- Operators: \*, +, -, !, :, ||, &&, <=, ==, ...
- Identifiers: x, result, main, map, \_\_a323\_int, ...  
Includes: variable names, user-defined type names, function names, library function and type names
- Keywords: let, val, fun, fn, end, while, for ...  
Keywords are reserved identifiers, constructs for statements and declarations.  
Built-in type names are often considered as keywords.



## Formal Languages

Building a machine to process words and compile programmes requires formal definitions.

### Definition (Formal Language)

Let  $\Sigma$  be an alphabet: a finite set of allowed characters.

- A word over  $\Sigma$  is a string  $w = a_1 a_2 \dots a_n$  of  $n$  characters  $a_i \in \Sigma$ .  
 $n = 0$  is allowed, and results in the empty string  $\varepsilon$ .  
We write  $\Sigma^*$  for the set of all words over  $\Sigma$ .
- A language  $L$  over  $\Sigma$  is a set of words over  $\Sigma$ :  $L \subset \Sigma^*$

Examples (alphabet: small latin letters)

- $\Sigma^*$  and  $\emptyset$ .
- $\{a^n b^n c^n \mid n \in \mathbb{N}\}$
- All C++ keywords: {if, else, return, do, while, int, char...}
- All palindromes (words that are the same backward and forward): {kayak, racecar, mellem, retter, ...}.



## Example Languages: Number Literals in C++

- Integers in decimal format: 123 4 0 8 but not 08 abc
- Integers in octal format: 0123 07 007 but not 08 abc
- Integers in hexadecimal format: 0x123 0xCafE but not 0x 0xG
- Floating point decimals: 0. .123 0123.456
- Scientific notation: 0123E-456 0.E123 .123e+456

A decimal integer is a sequence of digits 0-9 which does not start by 0 or is only a single 0. An octal integer is a sequence of digits starting with 0, followed by any number of digits 0-7.

Floating-point constants have a "mantissa," [...] and an "exponent," [...] The mantissa is specified as a sequence of digits followed by a period, followed by an optional sequence of digits [...]. The exponent, if present, specifies the magnitude [...] using e or E [...] followed by an optional sign (+ or -) and a sequence of digits. If an exponent is present, the trailing decimal point is unnecessary in whole numbers. <http://msdn.microsoft.com/en-us/library/tfh6f0w2.aspx>.

We need a more formal description for automatic processing.



## Regular Expressions

### Definition (Regular Expression)

Let  $\Sigma$  be an alphabet of allowed characters.

The set  $RE(\Sigma)$  of regular expressions over  $\Sigma$  is defined recursively.

- $\varepsilon \in RE(\Sigma)$ : describes the empty word.
- $a \in RE(\Sigma)$  for  $a \in \Sigma$ : describes word  $a$ .

Furthermore, for every  $\alpha, \beta \in RE(\Sigma)$ :

- $\alpha \cdot \beta \in RE(\Sigma)$ : Sequence, one word described by  $\alpha$ , followed by one described by  $\beta$ .
- $\alpha \mid \beta \in RE(\Sigma)$ : Alternative, a word described by  $\alpha$  or by  $\beta$ .
- $\alpha^* \in RE(\Sigma)$ : Repetition, zero or more words described by  $\alpha$ .

- Round parentheses (...) for grouping regular expressions.
- Sequence binds tighter than alternative,  $a(b|c) = ab|ac$



## Example Languages: Number Literals in C

- Integers in decimal format: 123 4 0 8 but not 08 abc
- Integers in octal format: 0123 07 007 but not 08 abc
- Integers in hexadecimal format: 0x123 0xCafE but not 0x 0xG
- Floating point decimals: 0. .123 0123.456
- Scientific notation: 0123E-456 0.E123 .123e+456
- Decimal Numbers: (1|2|...|9)(0|1|2|...|9)\* | 0  
Shorthand – character range: [1-9][0-9]\* | 0
- Octal format: 0 [0-7]\*
- Hexadecimal format: 0 ([xX] [0-9a-fA-F] [0-9a-fA-F] )\*Shorthand – at least once: 0 [xX] [0-9a-fA-F]+
- Floating point numbers: ... (later)



## Common Abbreviations for Regular Expressions

- Character Sets  
 $[a_1 a_2 \dots a_n] := (a_1 \mid a_2 \mid \dots \mid a_n)$   
One of  $a_1, \dots, a_n \in \Sigma$ .
- Character Ranges  
 $[a_1 - a_n] := (a_1 \mid a_2 \mid \dots \mid a_n)$  when  $\{a_i\}$  is ordered.  
One character in the range between  $a_1$  and  $a_n$ .
- Optional Parts  
 $\alpha? := (\alpha \mid \varepsilon)$  for  $\alpha \in RA(\Sigma)$ .  
Optionally a string described by  $\alpha$ .
- Repeated Parts  
 $\alpha^+ := \alpha \alpha^*$  for  $\alpha \in RA(\Sigma)$ .  
At least one string described by  $\alpha$  (maybe more).
- $\alpha\{n, m\} := \overbrace{(\alpha) \dots (\alpha)}^n \overbrace{(\alpha)? \dots (\alpha)?}^{m-n}$  for  $\alpha \in RA(\Sigma)$ ,  $n < m$ .  
A string described by  $\alpha$  between  $n$  and  $m$  times.



### Mosml-lex: Generating Lexical Analysis Programs

```

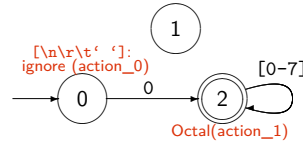
{ (* initial part containing SML code *)
  (* helper functions and data types *)

  data type MyTokens = Decimal of int | Octal of int | ...
  fun decodeInBase (base:int) (s:string) :int = ...
}
let oct = ['0'-'7'] (* a helper definition for reg. expr.s *)
(* Rules section: regular expr., action (returning a token) *)
rule Token = parse
  '0' oct* { Octal (decodeInBase 8 (getLexeme lexbuf)) }
  | ..regex2.. { ...action.2.. (* (SML code) *) }
  ...
;
    
```

### Demo

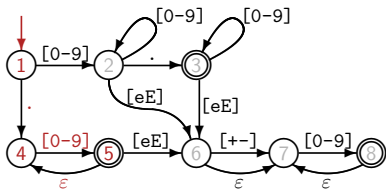
### Under the Hood: States and Actions

Octal numbers, from mosml-lex code:



- Start in state\_0:
- Input whitespace: restart (action\_0)
- Input 0: go to state\_2
- When in state\_2:
- Input 0-7: continue
- otherwise: do action\_1
- state\_1: never reached

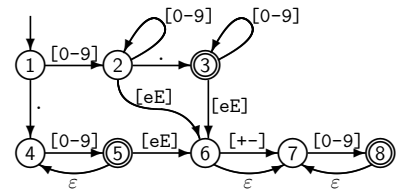
### An Automaton to Analyse Input...



.31 ← **.3e1.** ← .31.4159

- Starting at the pointed state
- At end of input: check if state "accepting"
- Transitions to new states, possibly reading input
- If no transition: stuck, input refused.

### An Automaton to Analyse Input...



.31 ← **.** ← 4159

- Starting at the pointed state
- At end of input: check if state "accepting"
- Transitions to new states, possibly reading input
- If no transition: stuck, input refused.

### Non-deterministic Finite Automaton – Definition

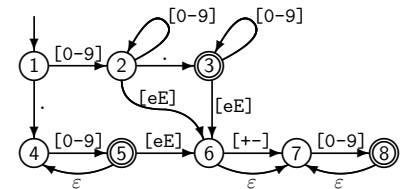
#### Definition (NFA)

Let  $\Sigma$  be an alphabet of (input) characters.  
A Non-deterministic Finite Automaton (NFA) consists of

- A finite set  $S$  of states,
- an alphabet  $\Sigma$  of input characters,
- a start state  $s_0 \in S$ ,
- a set of final states  $F \subset S$
- and a relation  $T \subset S \times (\Sigma \cup \{\epsilon\}) \times S$  describing state transitions (notation:  $s_i \xrightarrow{c} s_j \in T$ )

- Meaning of  $s_1 \xrightarrow{a} s_2 \in T$ : in  $s_1$  with input  $a$ , go to  $s_2$
- Meaning of  $s_1 \xrightarrow{\epsilon} s_2 \in T$ : in  $s_1$ , go to  $s_2$ .
- Several options may exist,  $T$  is a relation.

### Example NFA – Formalised

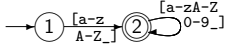


$$\begin{aligned}
 S &= \{1, \dots, 8\}, \Sigma = \{0, 1, \dots, 9, ., e, E\}, s_0 = 1, F = \{3, 5, 8\} \\
 T &= \{1 \xrightarrow{d} 2, 1 \xrightarrow{.} 4, 2 \xrightarrow{d} 2, 2 \xrightarrow{.} 3, 2 \xrightarrow{e} 6, 2 \xrightarrow{E} 6, 3 \xrightarrow{d} 3, 3 \xrightarrow{e} 6, 3 \xrightarrow{E} 6, \\
 &\quad 4 \xrightarrow{d} 5, 5 \xrightarrow{e} 4, 5 \xrightarrow{e} 6, 5 \xrightarrow{E} 6, 6 \xrightarrow{+} 7, 6 \xrightarrow{-} 7, 6 \xrightarrow{\epsilon} 7, 7 \xrightarrow{d} 8, 8 \xrightarrow{\epsilon} 7\} \\
 &\quad (d \in \{0, 1, \dots, 9\})
 \end{aligned}$$

Picture sufficient as definition – Formalisation: machine-processing.

### Other Examples

- Identifiers  
( $\Sigma$ : letters, digits, underscore)  
Starting with a letter or underscore, then any number of letters, underscores, and digits.



$[a-zA-Z_][a-zA-Z0-9_]^*$

- Binary numbers without leading zeros.

$$\Sigma = \{0, 1\}, S = \{0, 1, 2\}$$

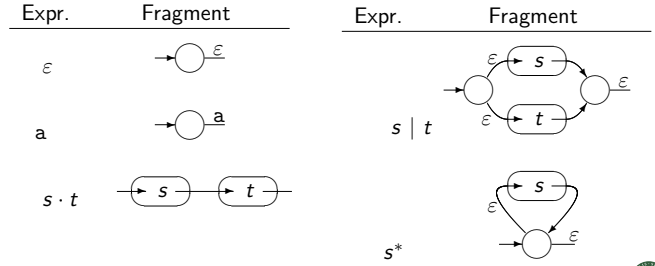
$$s_0 = 0, F = \{1, 2\}$$

$$T = \{0^0 1, 0^1 2, 2^0 2, 2^1 2\}$$

$0 | 1 [01]^*$

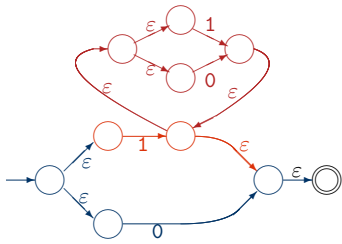
### NFA Construction from Regular Expression

- Define an NFA fragment for every regular expression  
Fragments have exactly one entry (arrow) and exit (line)
- Fragment composition follows expression composition  
A single final state is added at the end of the construction.



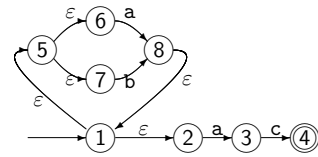
### Construction Examples

$0 | 1 (0 | 1)^*$



### Construction Examples

$(a|b)^*ac$



#### Undesired Non-determinism:

- Many  $\epsilon$  transitions (branch and exit for alternatives)
- Undefined transitions when automaton gets stuck
- Several alternatives for same input

### Deterministic Finite Automaton – Definition

#### Definition (DFA)

Let  $\Sigma$  be an alphabet of (input) characters.  
A Deterministic Finite Automaton (DFA) consists of

- A finite set  $S$  of states,
- an alphabet  $\Sigma$  of input characters,
- a start state  $s_0 \in S$ ,
- a set of final states  $F \subset S$
- and a function  $\delta : S \times \Sigma \rightarrow S$  describing state transitions.

- Meaning of  $\delta(s_1, a) = s_2$ : in  $s_1$  with input  $a$ , go to  $s_2$
- No empty input: always read an input character.
- Only one transition possible,  $\delta$  is a (partial) function.

### Extended Transition Function (whole words)

We define an extended version:  
A walk through the DFA reading a whole word  $w \in \Sigma^*$  at once.

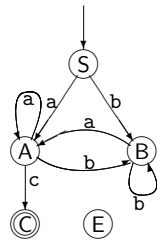
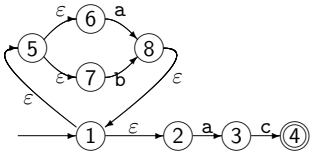
#### Definition (Word Transition Function)

Let  $D = (S, \Sigma, s_0, F, \delta)$  a DFA.  
The word transition function  $\bar{\delta} : S \times \Sigma^* \rightarrow S$  of  $D$  is defined recursively over words:

- For the empty word:  $\bar{\delta}(s, \epsilon) = s$
- For  $a \in \Sigma, w \in \Sigma^*$ :  $\bar{\delta}(s, aw) = \bar{\delta}(\delta(s, a), w)$   
if  $\delta(s, a)$  is defined.

Language accepted by the DFA:  $\{w \in \Sigma^* \mid \bar{\delta}(s_0, w) \in F\} \subset \Sigma^*$ .

### Converting an NFA to a DFA: Idea



- States 1,2,5,6,7 reachable from the start state 1.
- With input a, the NFA can go to state 3 and 8.  
On input b, only state 8 possible.
- States 1,2,5,6,7,8 reachable from 8.
- If in state 3, the NFA can go to state 4 on input c (otherwise nowhere).

S: {1, 2, 5, 6, 7}  
 A: {1, 2, 3, 5, 6, 7, 8}  
 B: {1, 2, 5, 6, 7, 8}  
 C: {4}  
 E: ∅

E is an error state.

### The Subset Construction: Preparation

To formalise the idea, we first define this reachability.

#### Definition ( $\epsilon$ -Closure $\hat{\epsilon}(\cdot)$ )

Let  $N = (S, \Sigma, s_0, F, T)$  a given NFA, and  $M \subset S$  a set of states. The  $\epsilon$ -Closure of  $M$ , written  $\hat{\epsilon}(M)$  contains all states reachable from states in  $M$ . It is recursively defined:

- $M \subset \hat{\epsilon}(M)$
- If  $s \in \hat{\epsilon}(M)$ , then  $\{s' \mid s \xrightarrow{\epsilon} s'\} \subset \hat{\epsilon}(M)$ .

$\hat{\epsilon}(M)$  is the smallest subset of  $S$  that fulfills these conditions.

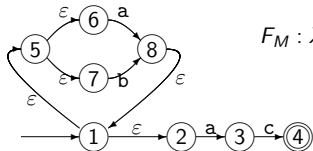
As a set equation:  $X = M \cup \{s' \mid \exists s \in X : s \xrightarrow{\epsilon} s'\}$

Solve this equation by computing a fixed point of  $F_M$ :

$$F_M : X \mapsto M \cup \{s' \mid \exists s \in X : s \xrightarrow{\epsilon} s'\}$$

Starting by  $X_0 = \emptyset$ , compute  $X_i = F(X_{i-1})$  until  $X_n = F(X_n)$   
 (works because  $F$  is monotonic:  $X \subseteq Y \Rightarrow F(X) \subseteq F(Y)$ ).

### $\epsilon$ -Closure $\hat{\epsilon}(\cdot)$ : Example



$$F_M : X \mapsto M \cup \{s' \mid \exists s \in X : s \xrightarrow{\epsilon} s'\}$$

Starting with  $X_0 = \emptyset$ , compute:  
 $X_i = F_M(X_{i-1}) = F_M^i(\emptyset)$   
 ... until  $X_n = F(X_n)$ .

Obviously:  $\hat{\epsilon}(\emptyset) = \emptyset$

$$\hat{\epsilon}(\{1\}) = \{1, 2, 5, 6, 7\}$$

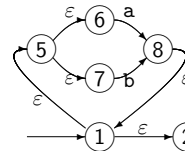
$$X_1 = \{1\}, F(X_1) = X_2 = \{1, 2, 5\},$$

$$F(X_2) = X_3 = \{1, 2, 5, 6, 7\} = F(X_3) = \hat{\epsilon}(\{1\})$$

$$\hat{\epsilon}(\{8\}) = \{1, 2, 5, 6, 7, 8\}$$

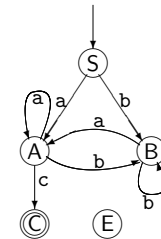
$$\hat{\epsilon}(\{3, 8\}) = \{1, 2, 3, 5, 6, 7, 8\}$$

### $\epsilon$ -Closure $\hat{\epsilon}(\cdot)$ : Example



- $\hat{\epsilon}(1) = S = \{1, 2, 5, 6, 7\}$
- $\hat{\epsilon}(3, 8) = A = \{1, 2, 3, 5, 6, 7, 8\}$
- $\hat{\epsilon}(8) = B = \{1, 2, 5, 6, 7, 8\}$
- $\hat{\epsilon}(\emptyset) = E = \emptyset$
- More transitions: No new states.
- $\hat{\epsilon}(4) = C = \{4\}$

E is an error state.



### The Subset Construction: Definition

#### Theorem (Subset Construction)

Let  $N = (S, \Sigma, s_0, F, T)$  a given NFA.

Define a DFA  $D = (S^d, \Sigma, s_0^d, F_d, \delta)$  as follows:

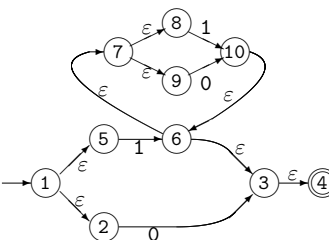
- $S^d = \mathbb{P}(S)$  (all subsets of  $S$ ).
- $s_0^d = \hat{\epsilon}(\{s_0\})$
- $F_d = \{M \subset S \mid M \cap F \neq \emptyset\}$  (subsets with a final NFA state).
- $\delta(s^d, a) = \hat{\epsilon}(\{t \mid s \in s^d, s \xrightarrow{a} t\})$   
 ( $t$  reachable from an  $s \in s^d$  on input  $a$ , and their  $\epsilon$ -Closure).

- This indeed defines a DFA.
- This DFA  $D$  accepts the same language as the NFA  $N$ .

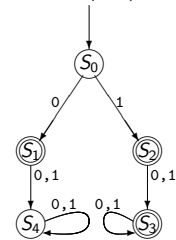
Proof idea: Consider a word  $w = a_1 a_2 \dots a_n$  accepted by the NFA. There is a state sequence  $s_1 \dots s_n$  such that  $(s_{i-1} \xrightarrow{a_i} s_i) \in T$  and  $s_i \in \hat{\epsilon}(\{s'_i\})$  and  $\hat{\epsilon}(\{s_n\}) \cap F \neq \emptyset$ . Therefore:  $s_1 \in s_1^d = \delta(\hat{\epsilon}(\{s_0\}), a_1)$ ,  $s_2 \in s_2^d = \delta(s_1^d, a_2)$ ,  $\dots$ ,  $s_n \in s_n^d = \delta(s_{n-1}^d, a_n)$ .  $\Rightarrow$  There exists an accepting state sequence  $s_1^d s_2^d \dots s_n^d$  in the DFA (since  $s_n^d \in F_d$ ).

### Subset Construction: A Second Example

$0 \mid 1 \mid (0 \mid 1)^*$



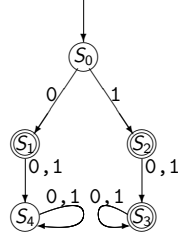
- $S_0 = \{1, 2, 5\}$
- $S_1 = \{3, 4\}$
- $S_2 = \{6, 3, 7, 4, 8, 9\}$
- $S_3 = \{10, 6, 3, 7, 4, 8, 9\}$
- $S_4 = \emptyset$  (error)



## Minimising a DFA

The DFAs we obtain from the subset construction are big!  
 Very often, they contain superfluous states.  
 Want to optimise the DFA (size and runs).  
 But, what exactly does "superfluous" mean?

- States that cannot lead to a final state (dead states).
- States that have identical transitions as others.  
 More generally: States that lead to the same outcome (acceptance, rejection) for any input.

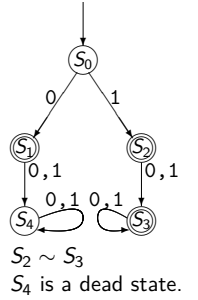


## Minimising a DFA: Preparation

### Definition (DFA State Equivalence)

Let  $D = (S, \Sigma, s_0, F, \delta)$  a DFA.

- A state  $s$  is called **dead** if and only if no final state can be reached from  $s$  with any input. Formally:  
 $s \text{ dead} := \bar{\delta}(s, w) \cap F = \emptyset$  for all  $w \in \Sigma^*$ .
- States  $s$  and  $s' \in S$  are called **equivalent**,  $s \sim s'$ , if and only if both lead to either acceptance or rejection with any input. Formally:  
 $s \sim s' := \bar{\delta}(s, w) \in F \Leftrightarrow \bar{\delta}(s', w) \in F$  for all  $w \in \Sigma^*$ .



## Minimising a DFA: Algorithm

We compute the equivalent states backwards from final states.

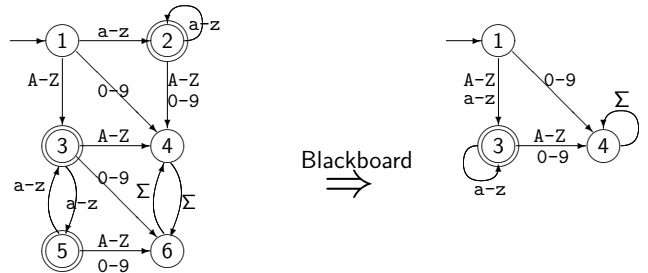
### Algorithm (DFA minimisation)

Let  $D = (S, \Sigma, s_0, F, \delta)$  a DFA. We assume  $\delta$  is total.  
 Determine state equivalence for a minimised DFA as follows:

- 1 Start with two unmarked groups,  $F$  and  $S \setminus F$ .
- 2 While there are unmarked groups:
  - Pick an unmarked group  $G$ .
  - For all  $a \in \Sigma$ , check for all states  $s \in G$  to which group a transition  $\delta(s, a)$  leads.
  - If for any respective input  $a$ , all transitions lead to the same group: mark the group.
  - Otherwise: Split the group into maximal groups that lead to the same group on transitions and unmark all groups.
- 3 Repeat from 2 until all groups are marked.

The resulting groups contain equivalent states  
 (all dead states will be equivalent).

## Minimising a DFA: Example



in the book: no dead states!

## Back to our original question. . .

Result so far: Is an input  $w$  described by the regular expression  $\alpha$ ?  
**Decision problem:** for  $w \in \Sigma^*$ : is  $w$  in the language described by  $\alpha \in RE(\Sigma)$ ?

How can we recognise a whole sequence of tokens?

```
...read by the lexer
(*_My_first\n*_SML_program.\n_*)\nval result\n=let val x_u=10:::020
_u:::0x30_u:::[]\nin List.map(fn_x_u=>_x_u div 2)_x_u\nend\n
```

```
... resulting token sequence
[Keyword "val", Id "result", Equal, Keyword "let", Keyword "val", Id "x",
Equal, Int 10, Doublecolon, Int 20, Doublecolon, Int 48, Doublecolon,
LBracket, RBracket, Keyword "in", Id "List", Dot, Id "map", LParen, ...
```

- Recognise prefixes of input as tokens.
- Restart on remaining input after recognising something.
- Often, several decompositions of the input possible.

## Principles of Longest and First Match

### Principle of Longest Match

A lexical analyser usually outputs the token that consumes the longest part of the input.

This is important when reading in identifiers and numbers (prefixes could otherwise be recognised instead).

### Principle of First Match

Tokens are usually prioritised, so the lexical analyser can decide which token to recognise if two tokens are possible for the same input.

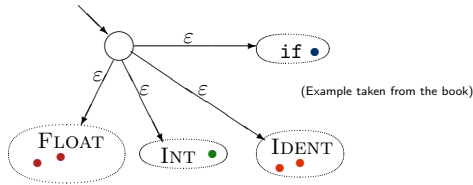
This is especially important when recognising keywords (they could otherwise be recognised as identifiers).

- Define combined NFA with prioritised final states, backtrack.

## Lexical Analysis: Putting it All Together

Construction of the automaton:

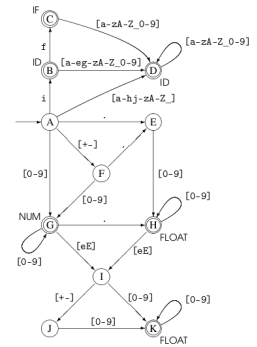
- 1 Define an NFA for each token class.
- 2 Mark final states in each NFA with the respective token name.
- 3 Combine the NFAs using new start state and  $\epsilon$  transitions.
- 4 Construct a small combined DFA, using subset construction and minimisation. Prioritise token classes in final DFA states to decide what to recognise.



## Lexical Analysis: Putting it All Together (2)

Processing input with the automaton:

- 1 Start the DFA in normal mode.
- 2 When reaching a final state: save it, enter read-ahead mode.
- 3 In read-ahead mode:
  - Buffer all input when reading.
  - When reaching a new final state: clear buffer.
  - End of input or DFA stuck: output last final state, restore input from buffer.
- 4 Restart in normal mode until input ends.



Example: Input 12E.3, recognised as NUM, ID, FLOAT



## More About Regular Languages. . .

Regular languages are (by definition) described by regular expressions, but likewise by NFAs, and DFAs. Therefore, we can argue:

- For two regular languages  $L_1, L_2 \subset \Sigma^*$ , their union  $L_1 \cup L_2$  and intersection  $L_1 \cap L_2$  are regular.
- For a regular language  $L \subset \Sigma^*$ , the complement  $\Sigma^* \setminus L$  is regular.
- Regular languages are also closed under common string operations: Prefix, Suffix, Subsequence, Reversal.
- The minimised DFA is uniquely determined. Two regular expressions are thus equivalent if their minimised DFAs are the same (apart from renaming states).

Regular languages are limited. Typically, what requires unbounded memory cannot be expressed as a regular language.

Palindromes (`{kayak, racecar, mellem, retter, ...}`) **not regular**.



## Summary

In this part, you have seen

- **Formal languages:** Sets of words over a finite alphabet.
- **Regular expressions**, describing regular languages (a subset of all formal languages)
- A compiler tool for lexical analysis (mosmlex) . . . and how the tool works internally:
- Deterministic and non-deterministic **finite automata** . . . and how to convert, minimise, and combine them.

