



**Script** generated by TTT

Title: Petter: Compilerbau (20.04.2015)

Date: Mon Apr 20 14:13:35 CEST 2015

Duration: 98:35 min

Pages: 51

## Compiler Construction I

Dr. Michael Petter

SoSe 2015

1 / 58

### Organizing

#### Dates:

Lecture: Mo. 14:15-15:45

Tutorial: You can vote on two dates via moodle

#### Exam:

- One Exam in the summer, none in the winter
- Exam managed via TUM-online
- Successful (50% credits) tutorial exercises earns 0.3 bonus

3 / 58

### Organizing

- Master or Bachelor in the 6th Semester with 5 ECTS
- Prerequisites
  - Informatik 1 & 2
  - Theoretische Informatik
  - Technische Informatik
  - Grundlegende Algorithmen
- Delve deeper with
  - Virtual Machines
  - Programoptimization
  - Programming Languages
  - Praktikum Compilerbau
  - Seminars

#### Materials:

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

2 / 58

## Preliminary content

- Basics in regular expressions and automata
- Specification and implementation of scanners
- Reduced context free grammars and pushdown automata
- Bottom-Up Syntaxanalysis
- Attribute systems
- Typechecking *registers*
- Codegeneration for ~~stack~~ machines
- Register assignment
- Basic Optimization

4 / 58

## Topic: Introduction

5 / 58

## Interpreter



### Pro:

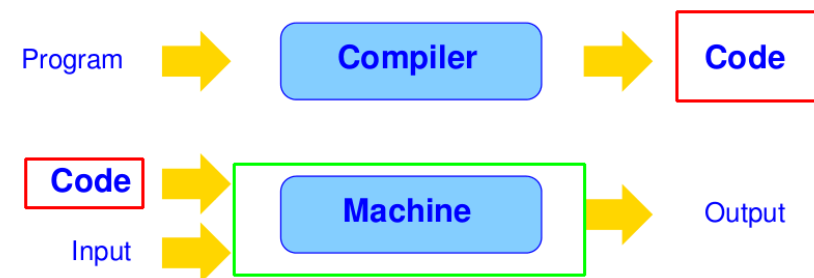
No precomputation on program text necessary  
⇒ no/small Startup-time

### Con:

Program components are analyzed multiple times during the execution  
⇒ longer runtime

6 / 58

## Concept of a Compiler



### Two Phases:

- 1 Translating the program text into a machine code
- 2 Executing the machine code on the input

7 / 58

## Compiler

A precomputation on the program allows

- a more sophisticated variable management
- discovery and implementation of global optimizations

### Disadvantage

The Translation costs time

### Advantage

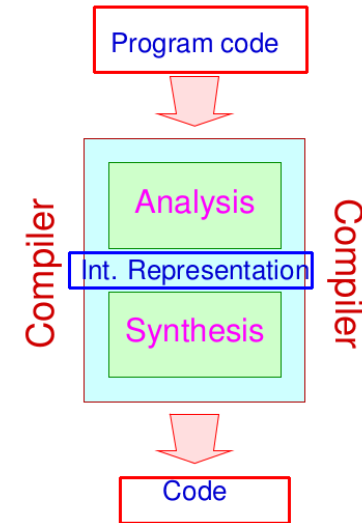
The execution of the program becomes more efficient

⇒ payoff for more sophisticated or multiply running programs.

8 / 58

## Compiler

general Compiler setup:



9 / 58

## Compiler

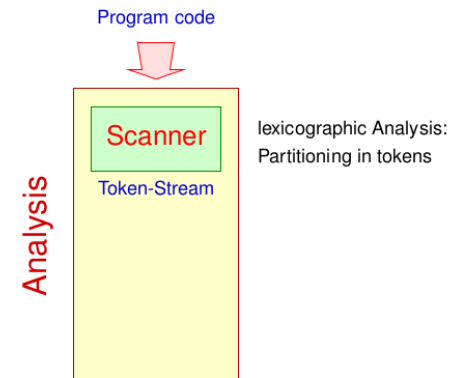
The **Analysis**-Phase is divided in several parts:



9 / 58

## Compiler

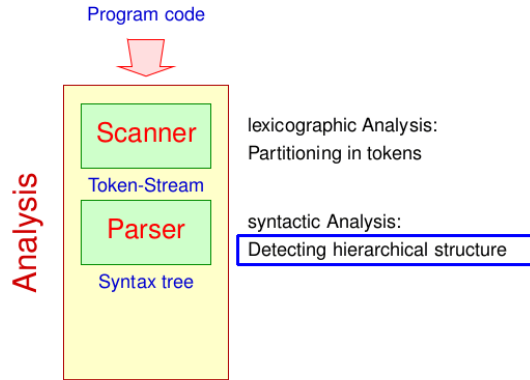
The **Analysis**-Phase is divided in several parts:



9 / 58

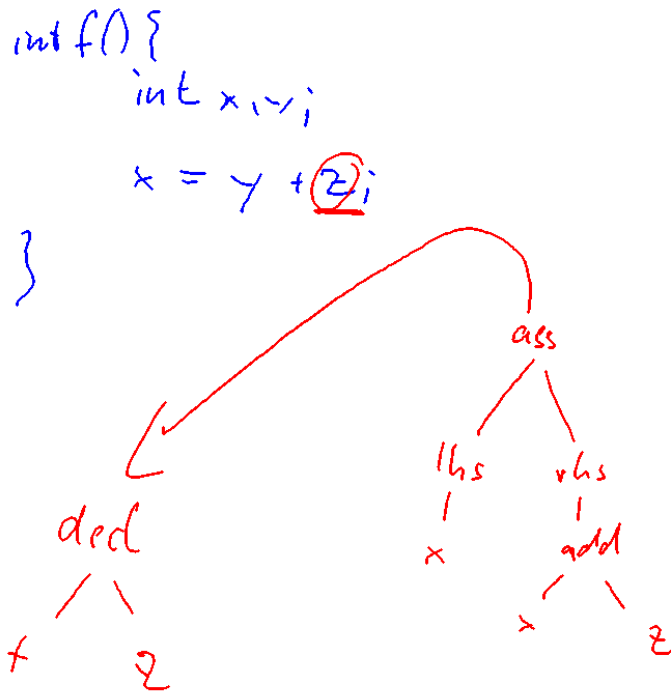
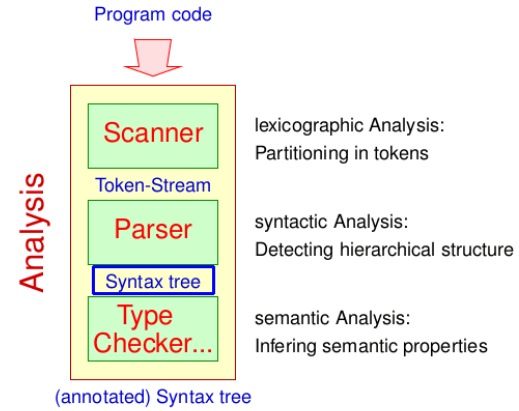
# Compiler

The Analysis-Phase is divided in several parts:



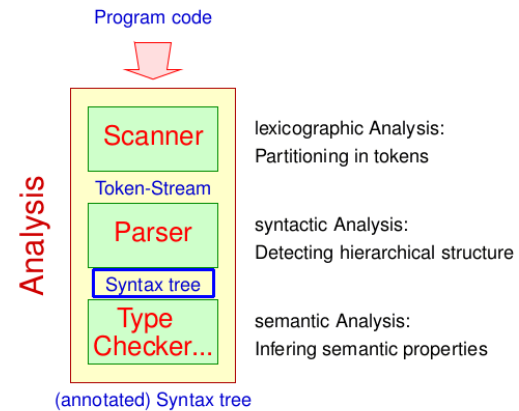
# Compiler

The Analysis-Phase is divided in several parts:

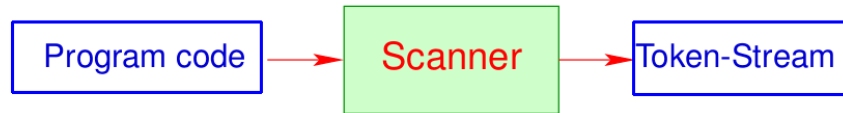


# Compiler

The Analysis-Phase is divided in several parts:

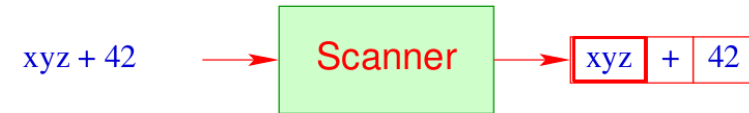


## The lexical Analysis



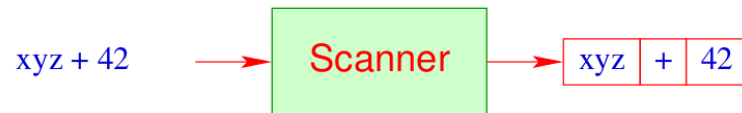
11 / 58

## The lexical Analysis



11 / 58

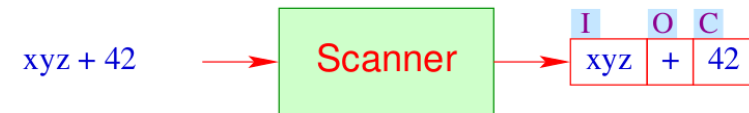
## The lexical Analysis



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

11 / 58

## The lexical Analysis



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

11 / 58

## The Lexical Analysis

Classified tokens allow for further pre-processing:

- Dropping irrelevant fragments e.g. Spacing, Comments,...
- Separating Pragma, i.e. directives vor the compiler, which are not directly part of the program, like include-Statements;
- Replacing of Tokens of particular classes with their meaning / internal representation, e.g.
  - Constants;
  - Names: typically managed centrally in a Symbol-table, evt. compared to reserved terms (if not already done by the scanner) and possibly replaced with an index.

⇒ Siever

12 / 58

## The Lexical Analysis

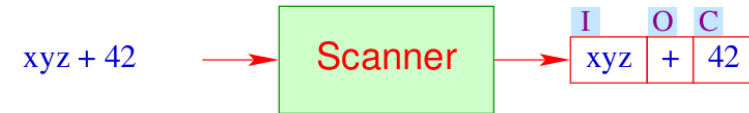
Classified tokens allow for further pre-processing:

- Dropping irrelevant fragments e.g. Spacing, Comments,...
- Separating Pragma, i.e. directives vor the compiler, which are not directly part of the program, like include-Statements;
- Replacing of Tokens of particular classes with their meaning / internal representation, e.g.
  - Constants;
  - Names: typically managed centrally in a Symbol-table, evt. compared to reserved terms (if not already done by the scanner) and possibly replaced with an index.

⇒ Siever

12 / 58

## The lexical Analysis



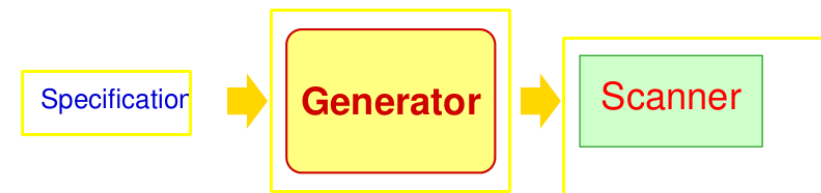
- A Token is a sequence of characters, which together form a unit.
- Tokens are subsumed in classes. For example:
  - Names (Identifiers) e.g. xyz, pi, ...
  - Constants e.g. 42, 3.14, "abc", ...
  - Operators e.g. +, ...
  - reserved terms e.g. if, int, ...

11 / 58

## The Lexical Analysis

### Discussion:

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



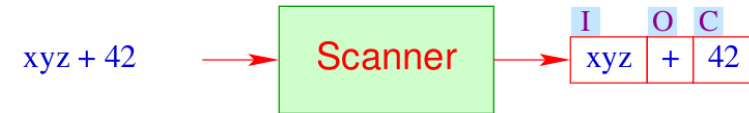
13 / 58

## The Lexical Analysis - Generating:

... in our case:



## The lexical Analysis



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

14 / 58

11 / 58

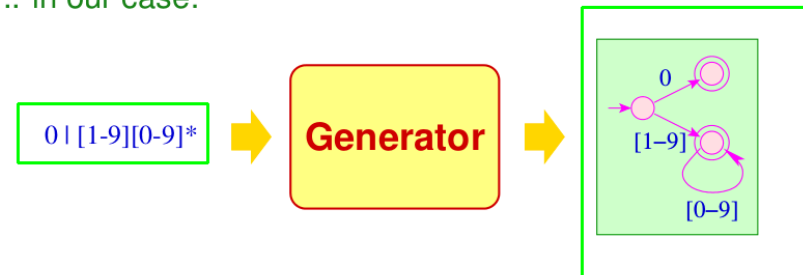
## The Lexical Analysis - Generating:

... in our case:



## The Lexical Analysis - Generating:

... in our case:



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

14 / 58

14 / 58

## Chapter 1: Basics: Regular Expressions

15 / 58

### Regular Expressions

#### Basics

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

#### Definition Regular Expressions

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

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



Stephen Kleene

16 / 58

### Regular Expressions

#### Basics

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

16 / 58

### Regular Expressions

... Example:

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

17 / 58



## Regular Expressions

... Example:

$$\begin{aligned} &((a \cdot b^*) \cdot a) \\ &(a \mid b) \\ &((a \cdot b) \cdot (a \cdot b)) \end{aligned}$$

Attention:

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

$$* > \cdot > |$$

and omit “.”

17 / 58

## Regular Expressions

... Example:

$$\begin{aligned} &((a \cdot b^*) \cdot a) \\ &(a \mid b) \\ &((a \cdot b) \cdot (a \cdot b)) \end{aligned}$$

Attention:

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

$$* > \cdot > |$$

and omit “.”

- Real Specification-languages offer additional constructs:

$$\begin{aligned} e? &\equiv (\epsilon \mid e) \\ e^+ &\equiv (e \cdot e^*) \end{aligned}$$

and omit “ $\epsilon$ ”

17 / 58

## Regular Expressions

Specifications need **Semantics**

...Example:

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

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

$$\begin{aligned} \llbracket \epsilon \rrbracket &= \{\epsilon\} \\ \llbracket a \rrbracket &= \{a\} \\ \llbracket e^* \rrbracket &= (\llbracket e \rrbracket)^* \\ \llbracket e_1 e_2 \rrbracket &= \llbracket e_1 \rrbracket \cup \llbracket e_2 \rrbracket \\ \llbracket e_1 \cdot e_2 \rrbracket &= \llbracket e_1 \rrbracket \cdot \llbracket e_2 \rrbracket \end{aligned}$$

18 / 58

## Keep in Mind:

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

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

19 / 58

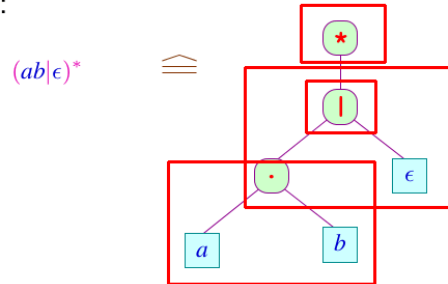
## Keep in Mind:

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

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

$$L_1 \cdot L_2 = \{w_1 w_2 \mid w_1 \in L_1, w_2 \in L_2\}$$

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



**Inner nodes:** Operator-applications;

**Leaves:** particular symbols or  $\epsilon$ .

19/58

## Regular Expressions

**Example:** Identifiers in **Java**:

le = [a-zA-Z\_\\$]

di = [0-9]

Id = {le} ({le} | {di})\*

20/58

## Regular Expressions

**Example:** Identifiers in **Java**:

le = [a-zA-Z\_\\$]

di = [0-9]

Id = {le} ({le} | {di})\*

Float = {di}\* (\. {di} | {di} \.) {di}\* ((e|E) (\+|\-)? {di}+)?

*.9e+80*

20/58

Lexical Analysis

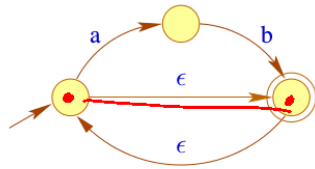
## Chapter 2:

## Basics: Finite Automata

21/58

## Finite Automata

Example:



22 / 58

## Finite Automata



Michael Rabin Dana Scott

### Definition Finite Automata

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

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

23 / 58

## Finite Automata



Michael Rabin Dana Scott

### Definition Finite Automata

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

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

For an NFA, we reckon:

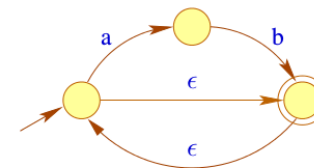
### Definition Deterministic Finite Automata

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

23 / 58

## Finite Automata

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



24 / 58

## Finite Automata

### Definition Finite Automata

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

$Q$	a finite set of states;
$\Sigma$	a finite alphabet of inputs;
$I \subseteq Q$	the set of start states;
$F \subseteq Q$	the set of final states and
$\delta$	the <u>set of transitions</u> (-relation)



For an NFA, we reckon:

### Definition Deterministic Finite Automata

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

23 / 58

## Finite Automata

### Definition Finite Automata

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

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



For an NFA, we reckon:

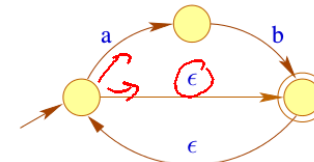
### Definition Deterministic Finite Automata

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

23 / 58

## Finite Automata

Example:



**Nodes:** States;

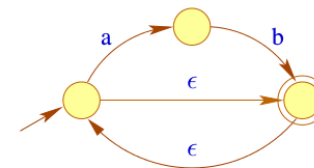
**Edges:** Transitions;

**Labels:** Consumed input;

22 / 58

## Finite Automata

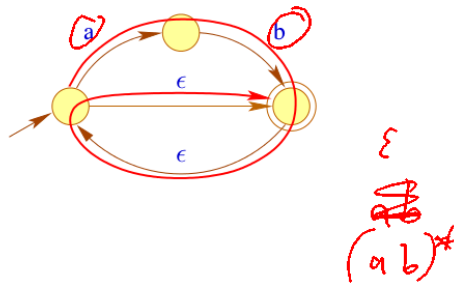
- Computations are paths in the graph.
- Accepting computations lead from  $I$  to  $F$ .
- An **accepted word** is the sequence of labels along an accepting computation ...



24 / 58

# Finite Automata

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



Lexical Analysis

## Chapter 3: Converting Regular Expressions to NFAs

# Finite Automata

Once again, more formally:

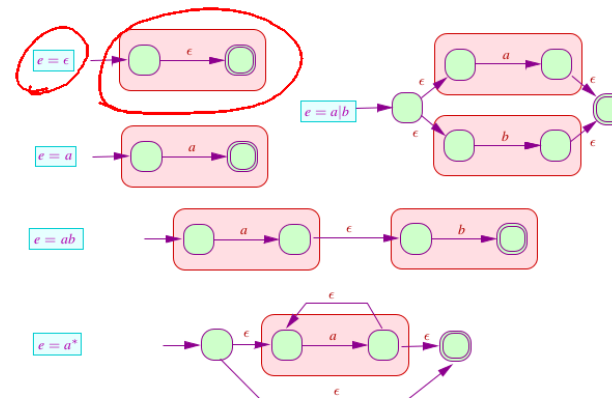
- We define the **transitive closure**  $\delta^*$  of  $\delta$  as the smallest set  $\delta'$  with:
  - $(p, \epsilon, p) \in \delta'$  and
  - $(p, xw, q) \in \delta'$  if  $(p, x, p_1) \in \delta$  and  $(p_1, w, q) \in \delta'$ .

$\delta^*$  characterizes for two states  $p$  and  $q$  the words, along each path between them

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

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

## In linear time from Regular Expressions to NFAs



### Thompson's Algorithm

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



Ken Thompson