

Compilers

Bert De Saffel

Master in de Industriële Wetenschappen: Informatica Academiejaar 2018–2019

Gecompileerd op 17 februari 2019

Inhoudsopgave

1	Inleiding	2
1.1	Compilers	2
1.2	Basiswerking compilers	2
1.3	Abstract Syntax Tree	3
1.3.1	Contextvrije grammatica's	3
1.3.2	Opbouw AST	4
1.3.3	Interpreter	6
2	Lexicale Analyse	7
2.1	Lexicale tokens	7
2.2	Eindige automaten	7
2.3	Opbouw deterministische eindige automaat	9
2.3.1	Conversie NFA naar DFA	11
3	Parsing	13

Hoofdstuk 1

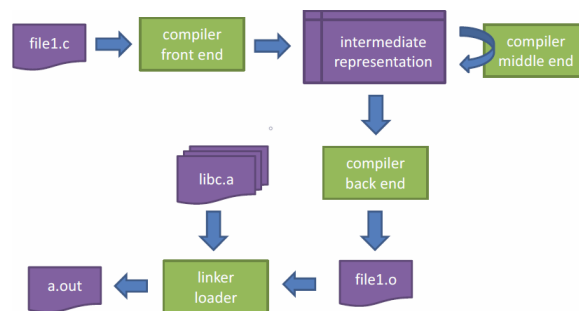
Inleiding

1.1 Compilers

Voorbeelden van functies die een statische compiler moet bevatten:

- Broncode omzetten in uitvoerbare fouten:
 - met dezelfde semantiek
 - zo snel mogelijk
 - en/of zo compact, debugbaar, portable, veilig, ... mogelijk
 - en linkbaar.
- Syntaxfouten moeten herkend worden.

1.2 Basiswerking compilers



Figuur 1.1: De basiswerking van een compiler.

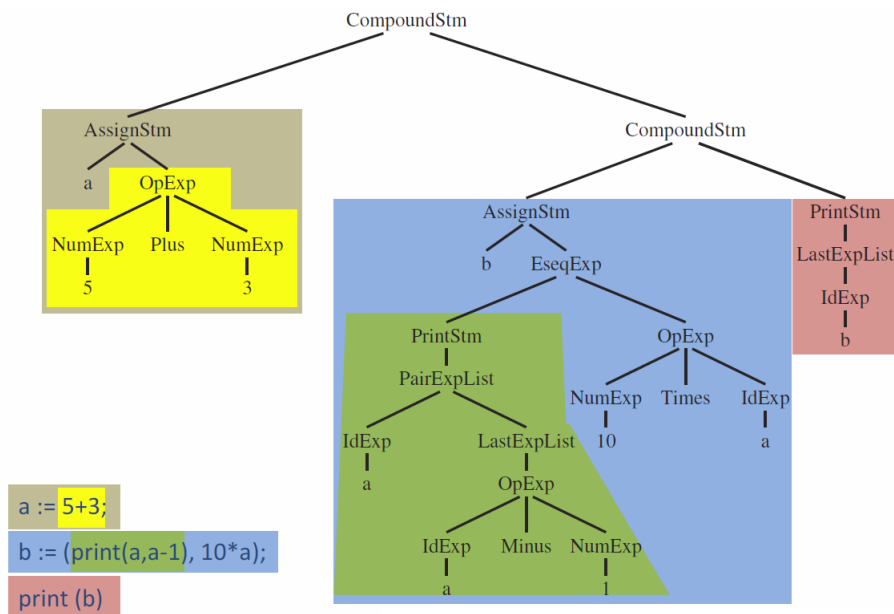
Op figuur 1.1 is de vereenvoudigde basiswerking van een compiler te zien. Een **C** bestand wordt eerst door de compiler front end gestuurd, die het bestand zal omvormen tot een intermediaire representatie. Deze representatie wordt dan door de compiler back end gestuurd om zo assembly of objectcode te genereren. De linker loader zal deze objectcode samenvoegen met eventuele andere libraries om zo een uitvoerbaar programma te hebben.

Q: Waarom wordt de front end en back end opgesplitst?

A: Op die manier is de compiler modulair: Enerzijds moet bij een andere programmeertaal enkel de front end aangepast worden en anderzijds moet bij het wijzigen van de architectuur (de onderliggende processor) enkel de back end aangepast worden.

1.3 Abstract Syntax Tree

De eerste stap van elke compiler is het omvormen van de broncode naar een **Abstract Syntax Tree (AST)**. Veronderstel volgende code, en de daarbijhorende AST die te zien zijn op figuur 1.2. Elke knoop van een AST stelt een bepaalde geldige operatie voor, die onafhankelijk is van de gekozen programmeertaal.



Figuur 1.2: De boomvoorstelling van een eenvoudig, lusloos programma. De gekleurde deelbomen komen overeen met de gekleurde segmenten in de code zelf. Als toekenningsoperator wordt er gekozen voor `:=` dat vanaf nu als één geheel moet beschouwd worden.

1.3.1 Contextvrije grammatica's

Om een AST op te stellen moet de notie van tokens ingevoerd worden. Een token is eenvoudig gezien een bepaald symbool dat een betekenis heeft. De tokens van de code uit figuur 1.2 zijn te zien in tabel 1.1. Uit de theorie van de generatieve grammatica's weten we dat er zowel terminale als niet-terminale tokens bestaan:

- **Terminale tokens** zijn symbolen die een blad voorstellen in de AST. Deze tokens hebben als eigenschap dat ze geen verdere tokens kunnen genereren en vormen dan ook het alfabet van het programma.
- **Niet-terminale tokens**, kortweg niet-terminalen genoemd, zijn de regels die de taal definiëren en zijn de niet-bladeren van de AST. Niet-terminalen hebben als eigenschap dat ze letters van het alfabet kunnen genereren.

symbolen(ascii)	token	waarde
a	id	string a
:=	:=	
5	num	integer 5
+	+	
3	num	integer 3
;	;	
b	id	string b
((
print	print	
-	-	
*	*	
	whitespace	

Tabel 1.1: De tokens die voorkomen uit het programma van figuur 1.2

Op figuur 1.3 zijn een aantal terminalen en niet-terminalen te zien. De niet-terminale token *CompoundStm* bestaat bijvoorbeeld uit twee *Stm* tokens, gescheiden door een punt komma. Deze twee *Stm* tokens kunnen in deze vereenvoudigde programmeertaal enkel een *AssignStm* of *PrintStm* zijn. Bij *AssignStm* wordt er een terminale token verwacht in de vorm van een variabele identifier, gevolgd door de toekenningoperator en een *Exp* token. Enkel deze *Exp* kan nog vier vormen aanneemen: *IdExp*, *NumExp*, enz... Dit wordt uitgewerkt voor de eerste toekenningsoperatie uit figuur 1.2 en

$Stm \rightarrow Stm ; Stm$	(CompoundStm)	$ExpList \rightarrow Exp , ExpList$	(PairExpList)
$Stm \rightarrow id := Exp$	(AssignStm)	$ExpList \rightarrow Exp$	(LastExpList)
$Stm \rightarrow print (ExpList)$	(PrintStm)	$Binop \rightarrow +$	(Plus)
$Exp \rightarrow id$	(IdExp)	$Binop \rightarrow -$	(Minus)
$Exp \rightarrow num$	(NumExp)	$Binop \rightarrow \times$	(Times)
$Exp \rightarrow Exp Binop Exp$	(OpExp)	$Binop \rightarrow /$	(Div)
$Exp \rightarrow (Stm , Exp)$	(EseqExp)		

Figuur 1.3: De rood omkaderde symbolen zijn **terminalen** terwijl de blauw omkaderde **niet-terminalen** zijn.

is te zien op figuur 1.4.

1.3.2 Opbouw AST

Een AST kan nu **bottom-up** opgemaakt worden door volgende procedure uit te voeren:

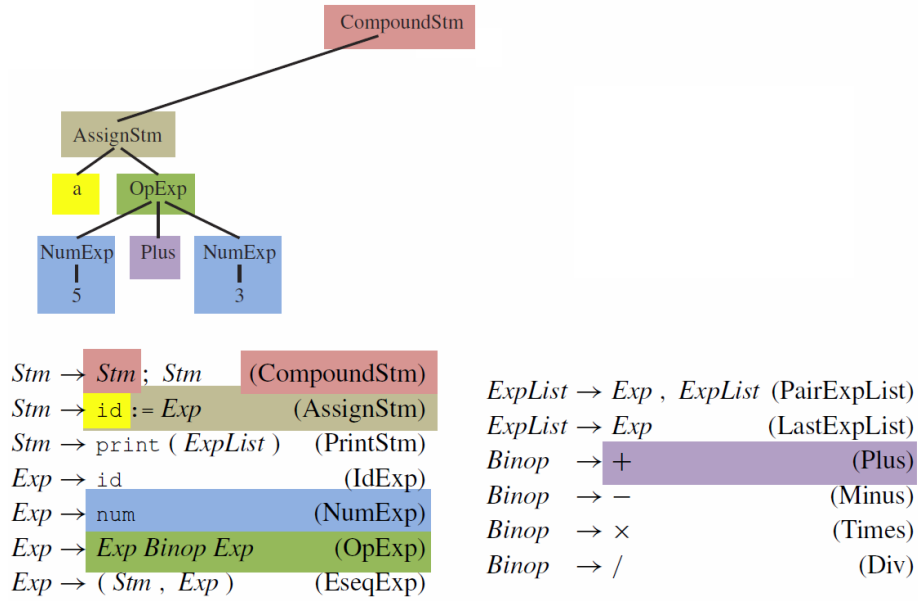
1. Voor elke mogelijke knoop moet er een struct gemaakt worden zoals bijvoorbeeld:

A_stm_ A_exp_ A_expList_

2. Elke struct moet bestaan uit

- een enum voor het precieze token te bepalen,
- een union voor de verschillende combinaties van tokens in het rechter lid en,
- pointers naar kindknoopen.

Dit wordt geïllustreerd in code 1.1.



Figuur 1.4: Illustratie van contextvrije grammatica op de eerste toekenningsoperatie uit figuur 1.2.

```

typedef char * string;
typedef struct A_stm_ * A_stm;
typedef struct A_exp_ * A_exp;
typedef struct A_expList_ * A_expList;

struct A_stm_ {
    enum {A_compoundStm, A_assignStm, A_printStm} kind;
    union {
        struct {A_stm stm1, stm2;} compound;
        struct {string id; A_exp exp;} assign;
        struct {A_expList exps;} print;
    } u;
};

```

Listing 1.1: Voorbeeld van een struct voor een AST.

3. In de constructor worden de knopen aangemaakt, zoals te zien in code 1.2.

```

A_stm A_CompoundStm(A_stm stm1, A_stm stm2){
    A_stm s = malloc(sizeof(*s));
    s->kind = A_compoundStm;
    s->u.compound.stm1 = stm1;
    s->u.compound.stm2 = stm2;
    return s;
}

```

Listing 1.2: Voorbeeld van een constructor voor een AST.

Op deze manier zou de boom uit figuur 1.2 hardgecodeerd kunnen worden, wat natuurlijk geen goede manier is. Het is de taak van een lexer en parser om de constructie van een AST te automatiseren, die respectievelijk in hoofdstuk 2 en 3 behandeld worden.

1.3.3 Interpreter

Uit een AST kan een eenvoudige interpreter geschreven worden. Dit stuk is informatief, en wordt niet gevraagd op het examen.

- Door de boom postorder diepte-eerst te overlopen, wordt de boom in de juiste manier behandeld.
- Het bijhouden van de waarden van variabelen kan via een gelinkte lijst:

```
typedef struct table * Table_;\nstruct table {string id; int value; Table_ tail;};\nTable_ Table(string id; int value; Table_ tail) {\n    Table_ t = malloc(sizeof(*t));\n    t->id = id;\n    t->value = value;\n    t->tail = tail;\n    return t;\n}
```

- Stel nu dat dit de eerste drie regels van een programma zijn:

```
a := 2;\nb := 3;\na := 3;
```

ToDo: onbelangrijk

Hoofdstuk 2

Lexicale Analyse

2.1 Lexicale tokens

- Herkennen van een reeks opeenvolgende karakters die een geheel vormen volgens de syntax van een programmeertaal, zoals o.a:
 - sleutelwoorden: int, float, for, new, ...
 - identifiers: foo, n14, variabelenaam
 - getallen: -37, 0x16L, 10.4, ...
 - operatoren: +, -, *, &, &&, ...
 - andere tokens: { } "; /* */ / () []

- Veronderstel volgende code:

```
float match0(char * s) { /* find a zero */
    if (!strcmp(s, "0.0", 3))
        return 0.;
}
```

, dan worden volgende tokens gegenereerd, waarbij dat sommige tokens een **attribuut** hebben:

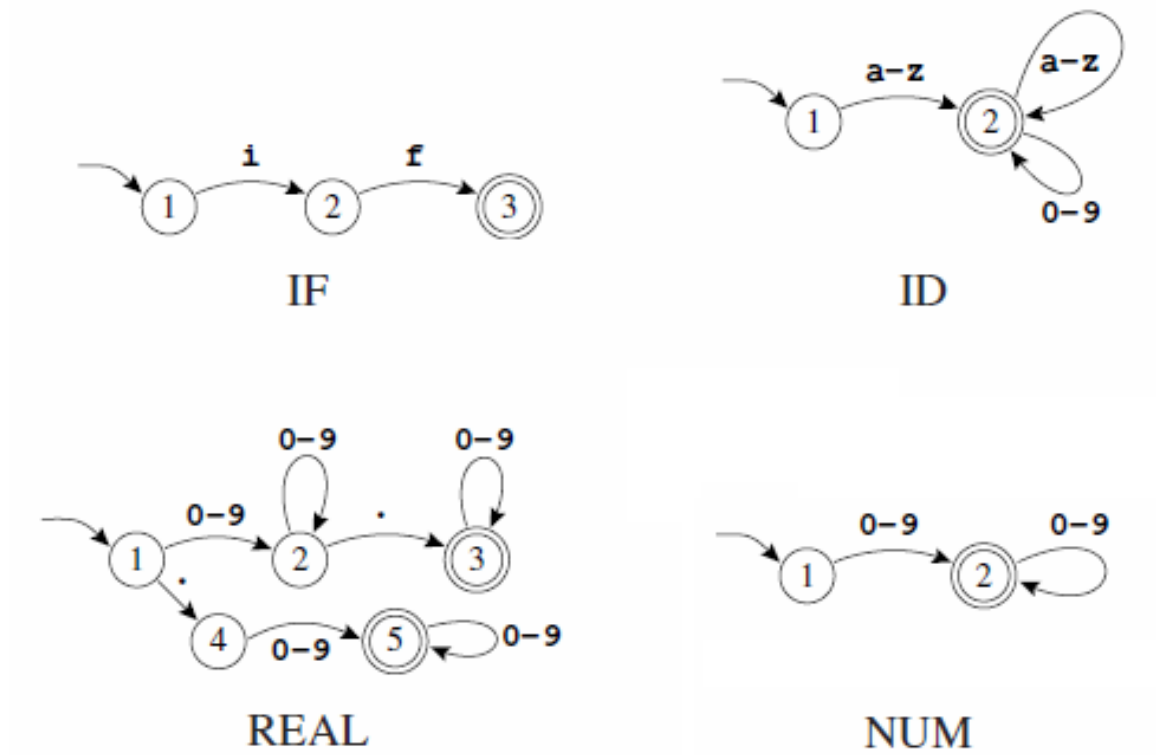
```
FLOAT ID(match0) LPAREN CHAR START ID(s) RPAREN LBRACE IF LPAREN BANG ID(strncmp)
LPAREN ID(s) COMMA STRING(''0.0'') COMMA NUM(3) RPAREN RPAREN RETURN REAL(0.0)
SEMI RBRACE EOF
```

2.2 Eindige automaten

- Er wordt met reguliere expressie gewerkt om te omschrijven welke karaktersequentie met een bepaald token overeenstemmen:

<i>if</i>	{return IF;}
$[a - z][a - z0 - 9]^*$	{return ID;}
$[0 - 9]^+$	{return NUM;}
$([0 - 9]^+ \mid "." [0 - 9]^+)^*$	{return REAL;}

Tabel 2.1: Reguliere expressies voor een aantal tokens.

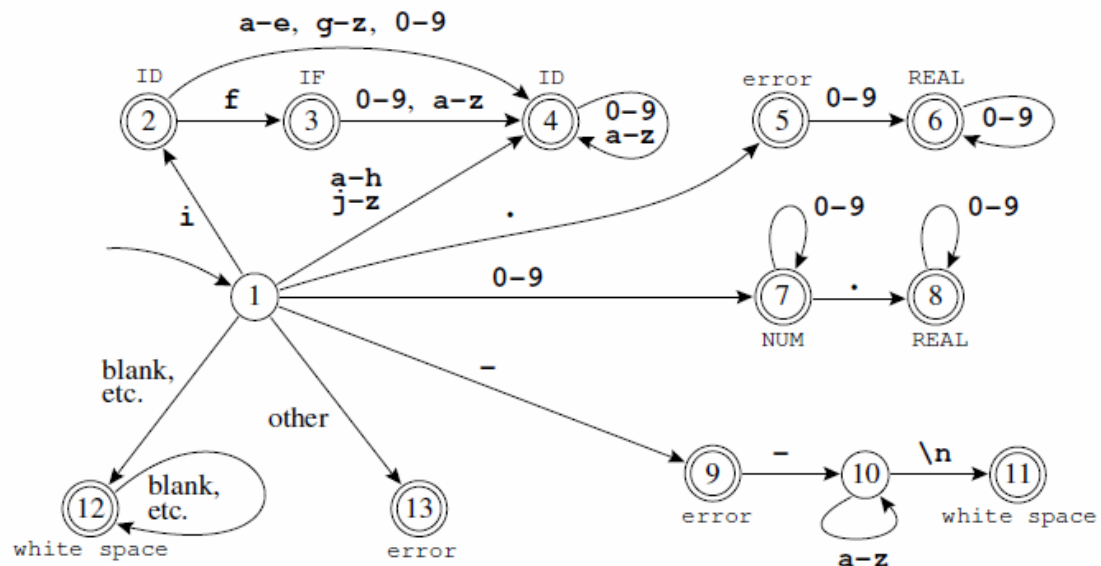


Figuur 2.1: Eindige automaten voor lexicale tokens.

- Met behulp van de constructie van Thompson kan een **niet-deterministische automaat (NFA)** opgebouwd worden uit een reguliere expressie. Op figuur 2.1 zijn de eindige automaten te zien van de reguliere expressies uit tabel 2.1.
- Deze individuele automaten kunnen samengevoegd worden tot een gecombineerde automaat, te zien op figuur 2.2

In dit geval is de gecombineerde automaat al een **deterministische eindige automaat (DFA)** aangezien elke mogelijke staat slechts één transitie heeft voor elke input. Het doel van lexicale analyse is om een DFA op te stellen zodat de tokens op een efficiënte manier kunnen bepaald worden. Een DFA wordt doorgaans geïmplementeerd als een transitietabel:

```
int edges[][256] = { /* ... 0 1 2 ... - ... e f g h i j ... */
/* state 0 */      { ... 0 0 0 ... 0 ... 0 0 0 0 0 0 ... },
/* state 1 */      { ... 7 7 7 ... 9 ... 4 4 4 4 2 4 ... },
/* state 2 */      { ... 4 4 4 ... 0 ... 4 3 4 4 4 4 ... },
/* state 3 */      { ... 4 4 4 ... 0 ... 4 4 4 4 4 4 ... },
/* state 4 */      { ... 4 4 4 ... 0 ... 4 4 4 4 4 4 ... },
/* state 5 */      { ... 6 6 6 ... 0 ... 0 0 0 0 0 0 ... },
/* state 6 */      { ... 6 6 6 ... 0 ... 0 0 0 0 0 0 ... },
/* state 7 */      { ... 7 7 7 ... 0 ... 0 0 0 0 0 0 ... },
/* state 8 */      { ... 8 8 8 ... 0 ... 0 0 0 0 0 0 ... },
/* ... */
};
```



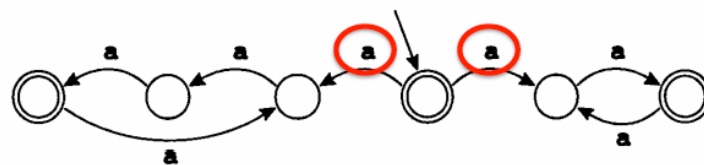
Figuur 2.2: Combinatie van eindige automaten.

2.3 Opbouw deterministische eindige automaat

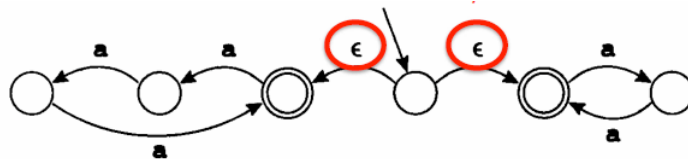
- We starten met een reguliere expressie, die een bepaald token voorstelt:

$$(aaa)^* | (aa)^*$$

- Zoals vermeld zal de constructie van Thompson een niet-deterministische automaat aanmaken van een bepaalde reguliere expressie. Er bestaat de kans dat deze automaat deterministisch is, maar dat is niet altijd zo. In het geval van bovenstaande reguliere expressie ziet de automaat er uit zoals op figuur 2.3 of figuur 2.4.



Figuur 2.3: Een niet-deterministische eindige automaat.



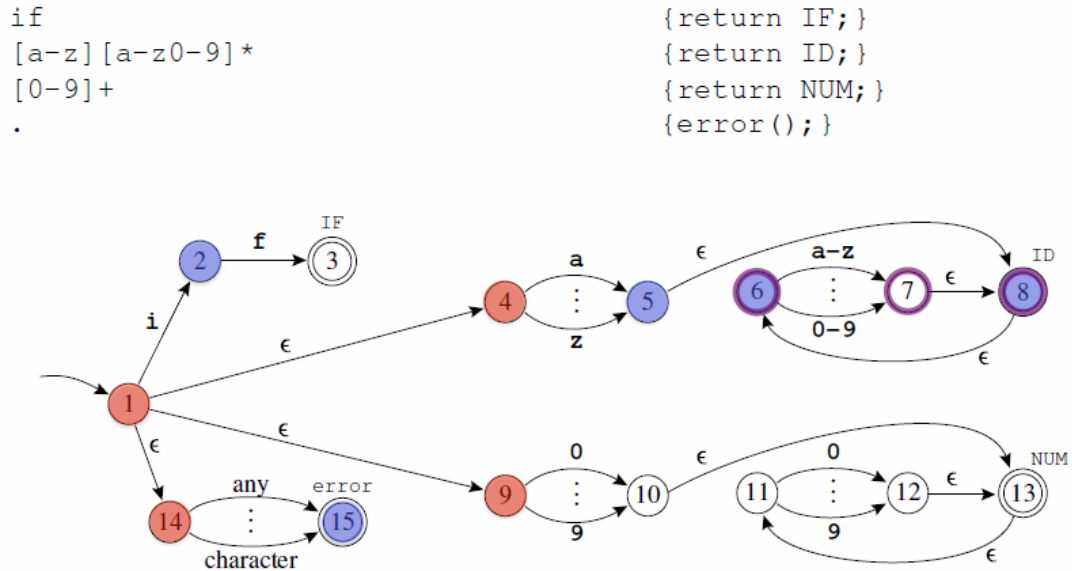
Figuur 2.4: Een niet-deterministische eindige automaat waarbij de eerste transitie kan gebeuren zonder een symbool te verwerken.

Welke richting moeten we nu uit bij `aaaaaaa` voor de eerste `a`? Bij een niet-deterministische automaat moeten we gokken welke de juiste zal zijn.

- Gelukkig kan ook een DFA opgebouwd worden uit een NFA via de deelverzamelingconstructie. Op die manier kan een DFA opgebouwd worden door (i) enkel de reguliere expressies handmatig

te definiëren, (ii) algoritmisch deze reguliere expressies om te vormen tot een NFA, en (iii) algoritmisch deze NFA om te vormen tot een DFA.

- Veronderstel de reguliere expressies en de daarbijhorende NFA in figuur 2.5.

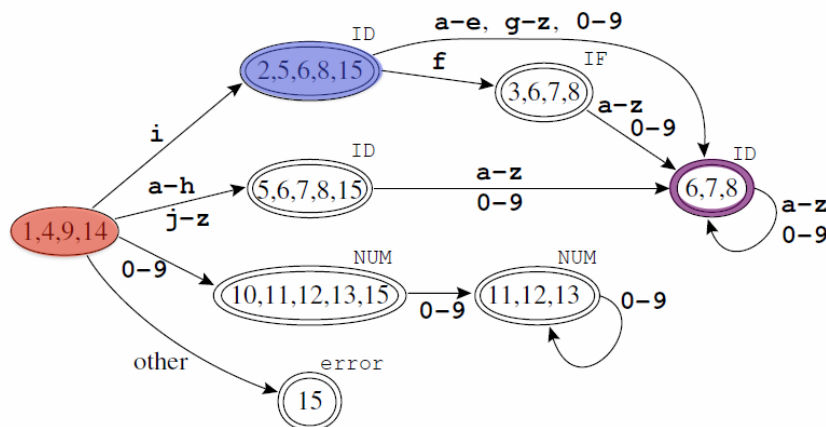


Figuur 2.5: Een aantal reguliere expressies en de daarbijhorende NFA.

Stel dat we nu de string *in* moeten checken:

1. Zonder teken op te eten kunnen we in 1 komen en in zijn ϵ -closure: $\{1, 4, 9, 14\}$.
2. Vanuit $\{1, 4, 9, 14\}$ kunnen we voor *i* naar $\{2, 5, 6, 8, 15\}$.
3. Vanuit $\{2, 5, 6, 8, 15\}$ kunnen we voor *n* naar $\{6, 7, 8\}$.
4. Daarvan is 8 een aanvaardingstoestand voor ID.

Op die manier bekommen we de DFA uit figuur 2.6.



Figuur 2.6: De NFA uit figuur 2.5 geconverteerd naar een DFA.

2.3.1 Conversie NFA naar DFA

- Drie functies:

1. **edge(s, c)** = alle NFA staten bereikbaar uit toestand s over pijlen met transitiesymbool c .
2. **closure(S)** = de kleinste verzameling T voor een subset S waarvoor geldt:

$$T = S \cup \left(\bigcup_{s \in T} \text{edge}(s, \epsilon) \right)$$

Q: Waarom moet dit de kleinste verzameling zijn?

A: De volledige verzameling van toestanden voldoet ook aan deze vergelijking, en dat is een triviaal geval.

Via iteratie kan T berekent worden:

```

T ← S
repeat
  T' ← T
  T ← T' ∪ (⋃s ∈ T' edge(s, ε))
until T = T'
```

- Dit is een voorbeeld van een fixpoint algoritme. Dit wil zeggen dat uiteindelijk $f(x) = x$ geldig is. In het voorbeeld van de functie **closure(S)**, hier genoteerd als **F(x)**, is dit zeker waar:

$$\begin{aligned}
 F(\epsilon) &= S \\
 F(S) &= \dots \\
 F(F(S)) &= \dots \\
 &\dots \\
 F(F(F(\dots))) &= T \\
 F(T) &= T
 \end{aligned}$$

- Aangezien dat uiteindelijk $F(T) = T$ en dat er maar een eindig aantal staten zijn zal het algoritme zeker stoppen.
3. Veronderstel dat we ons bevinden in een set $d = \{s_i, s_k, s_l\}$ van NFA staten s_i, s_k en s_l . Startend vanuit d en het symbool c , bekomen we een nieuwe set van NFA staten:

$$\mathbf{DFAedge(D, c)} = \mathbf{closure} \left(\bigcup_{s \in D} \text{edge}(s, c) \right)$$

Via deze functie, de startstaat s_1 en input string c_1, \dots, c_k kan de NFA simulatie als volgt geschreven worden:

```

d ← closure({s1})
for i ← 1 to k
  d ← DFAedge(d, ci)
```

- De combinatie van deze drie functies leiden tot het algoritme om een NFA om te zetten naar een DFA:

```

states[0] ← {};
states[1] ← closure({s1});
p ← 1;    j ← 0;
```

```
while  $j \leq p$ 
  foreach  $c \in \Sigma$ 
     $e \leftarrow \text{DFAedge}(\text{states}[j], c)$ 
    if  $e == \text{states}[i]$  for some  $i \leq p$ 
      then  $\text{trans}[j, c] \leftarrow i$ 
    else  $p \leftarrow p + 1$ 
       $\text{states}[p] \leftarrow e$ 
       $\text{states}[j, c] \leftarrow p$ 
   $j \leftarrow j + 1$ 
```

De gegenereerde DFA is suboptimaal: vanuit sommige toestanden worden identiek dezelfde strings aanvaard.

Hoofdstuk 3

Parsing