

DD1361

Programmeringsparadigm

Formella Språk & Syntaxanalys
Föreläsning 4

Per Austrin

2016-11-16

Idag

Stackautomater

Olika klasser av språk och grammatiker

Parsegeneratorer

Sammanfattning / Inför KS

Idag

Stackautomater

Olika klasser av språk och grammatiker

Parsegeneratorer

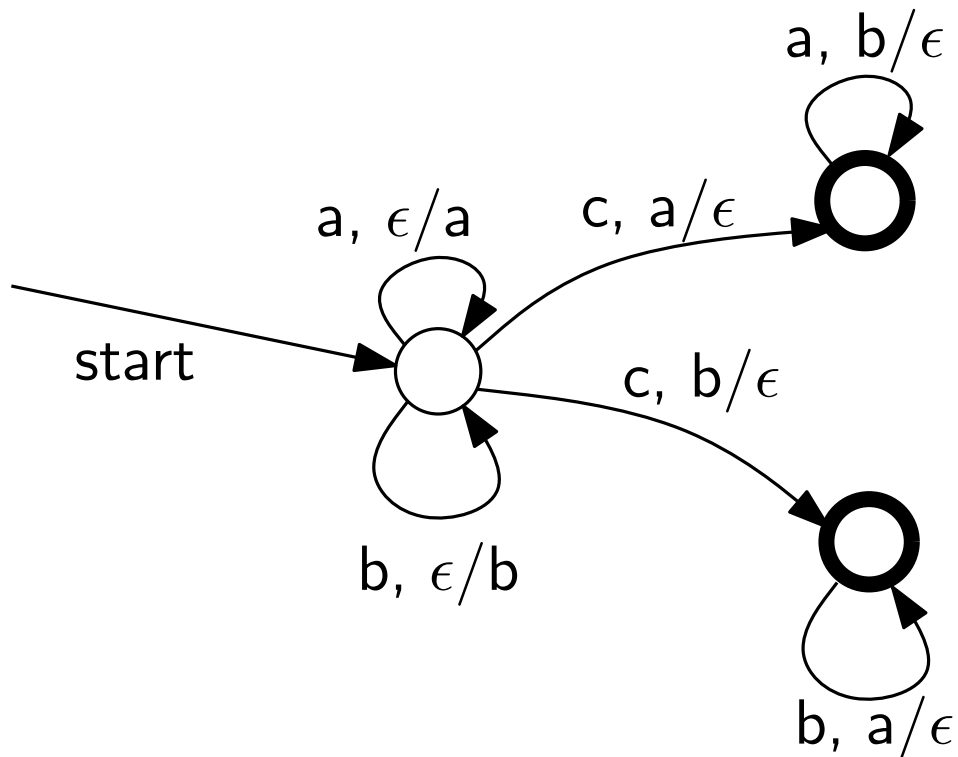
Sammanfattning / Inför KS

Stackautomater

Stackautomat (PDA, Push-Down Automaton):
som DFA, men har ett **obegränsat** minne i form av en stack

Stackautomater

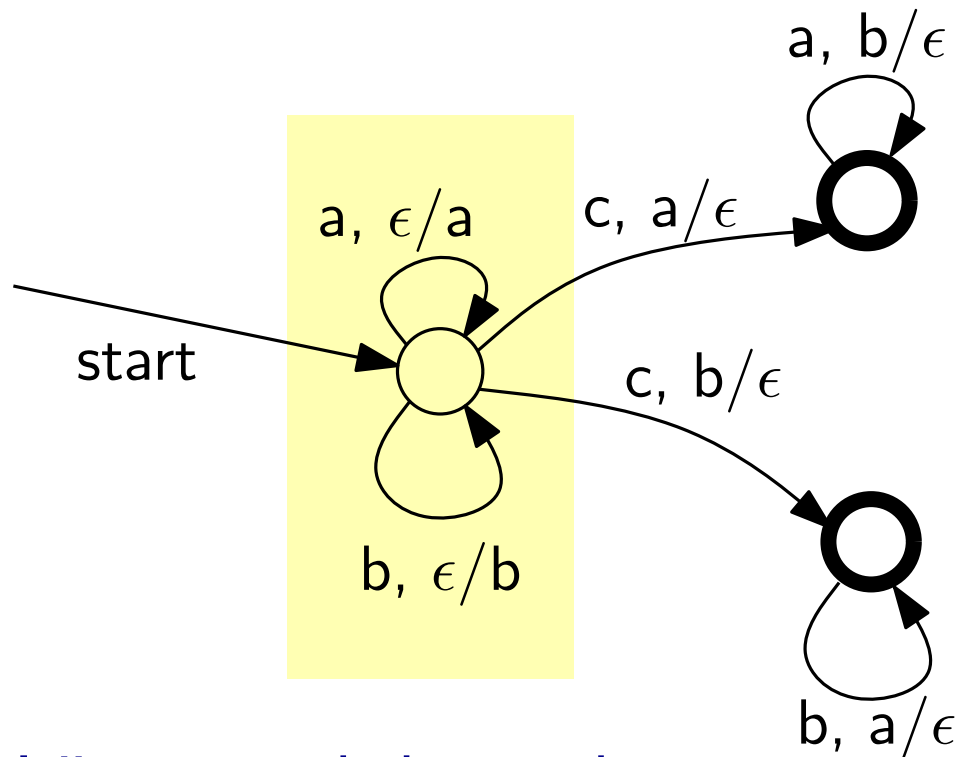
Stackautomat (PDA, Push-Down Automaton):
som DFA, men har ett **obegränsat** minne i form av en stack



$x, y/z$: läs x , poppa y
från stacken, pusha z

Stackautomater

Stackautomat (PDA, Push-Down Automaton):
som DFA, men har ett **obegränsat** minne i form av en stack

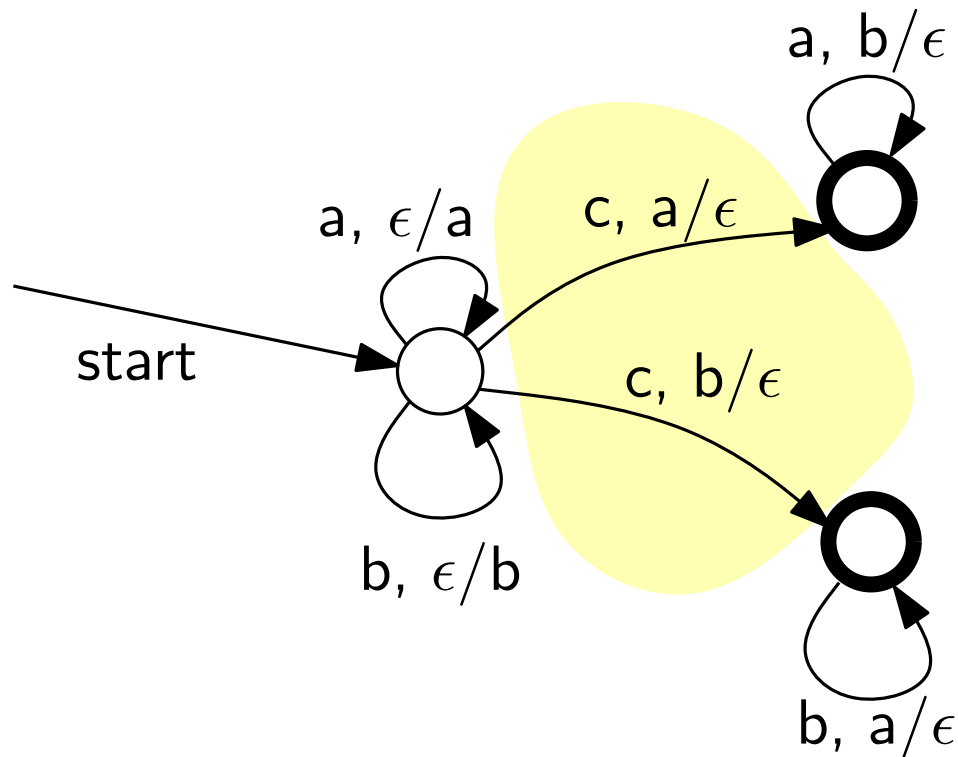


$x, y/z$: läs x , poppa y
från stacken, pusha z

Läs a:n och b:n och
lägg på stacken

Stackautomater

Stackautomat (PDA, Push-Down Automaton):
som DFA, men har ett **obegränsat** minne i form av en stack

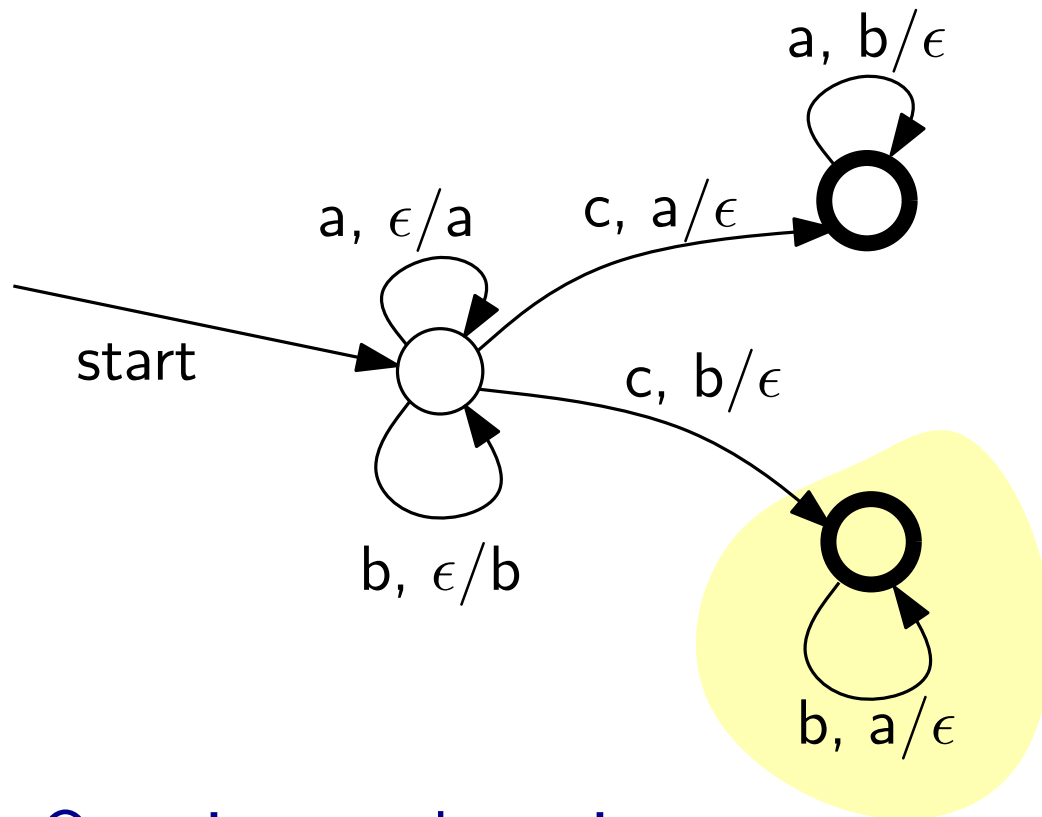


$x, y/z$: läs x , poppa y
från stacken, pusha z

När det kommer ett c , hoppa
till något av de accepterande
tillstånden beroende på om det
översta tecknet på stacken är a
eller b

Stackautomater

Stackautomat (PDA, Push-Down Automaton):
som DFA, men har ett **obegränsat** minne i form av en stack

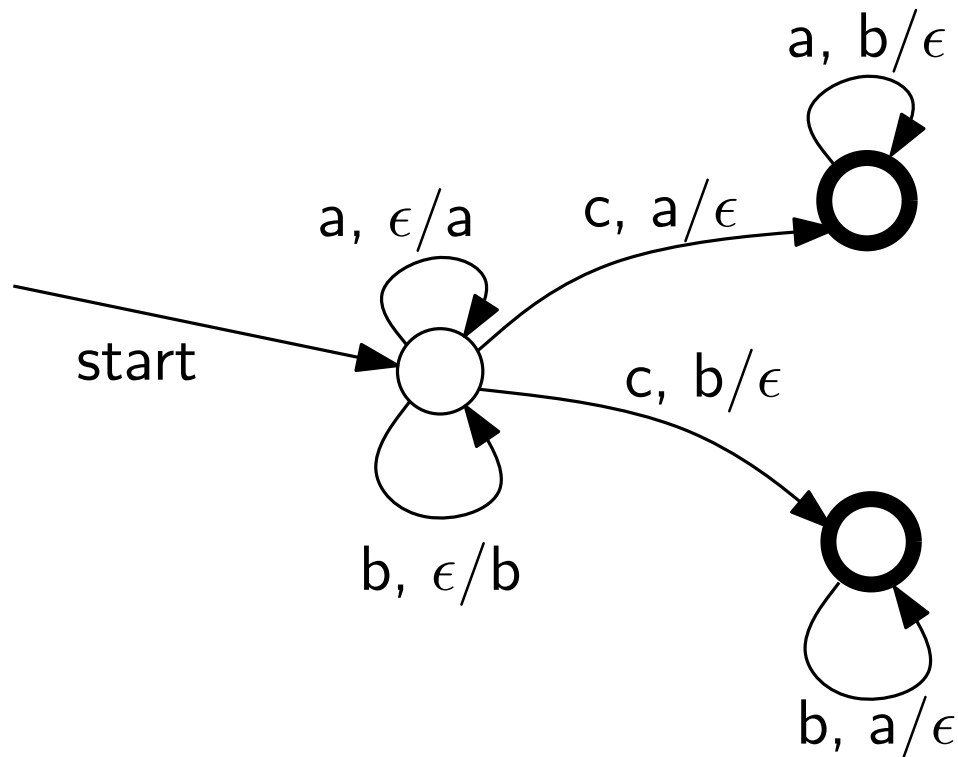


$x, y/z$: läs x , poppa y
från stacken, pusha z

Om sista tecknet innan c var ett b : fortsätt läsa b :n från
indata och kolla att det ligger minst lika många a :n på stacken

Stackautomater

Stackautomat (PDA, Push-Down Automaton):
som DFA, men har ett **obegränsat** minne i form av en stack

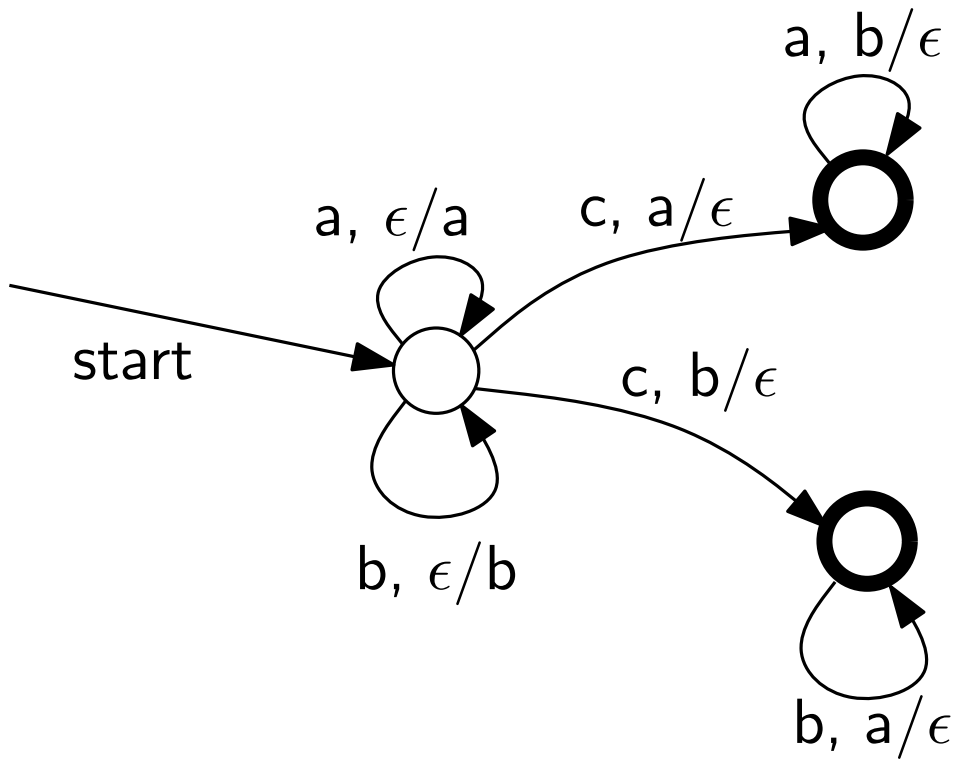


$x, y/z$: läs x , poppa y
från stacken, pusha z

Automaten accepterar om den är i ett accepterande tillstånd
OCH stacken är tom när indata är slut

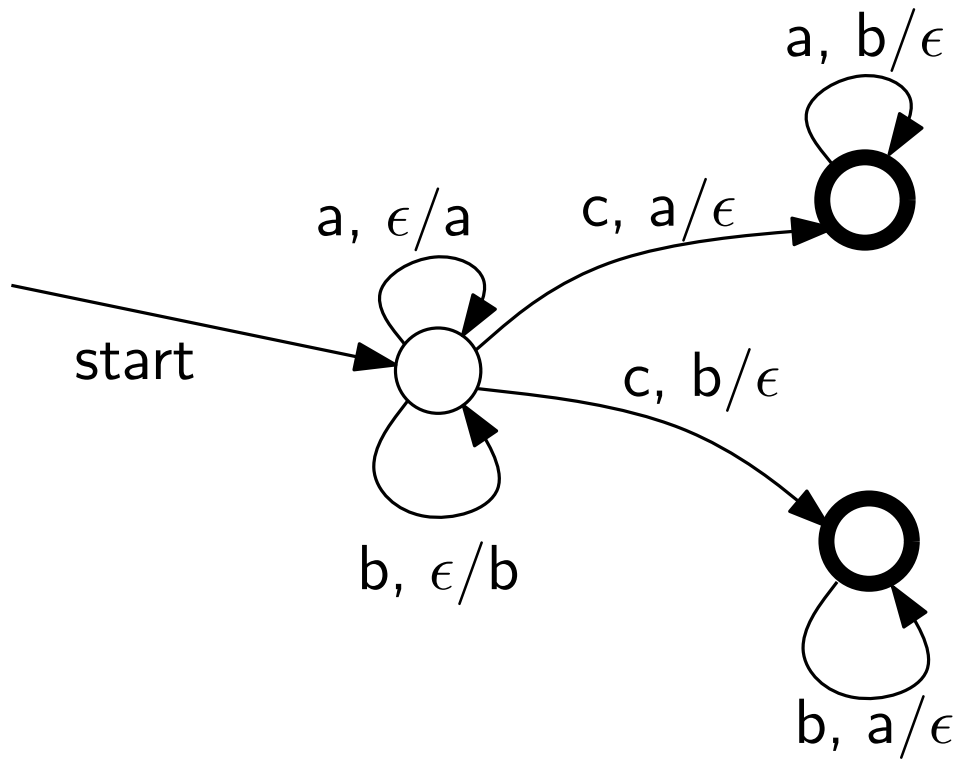
Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



Exempelkörning, stackautomat

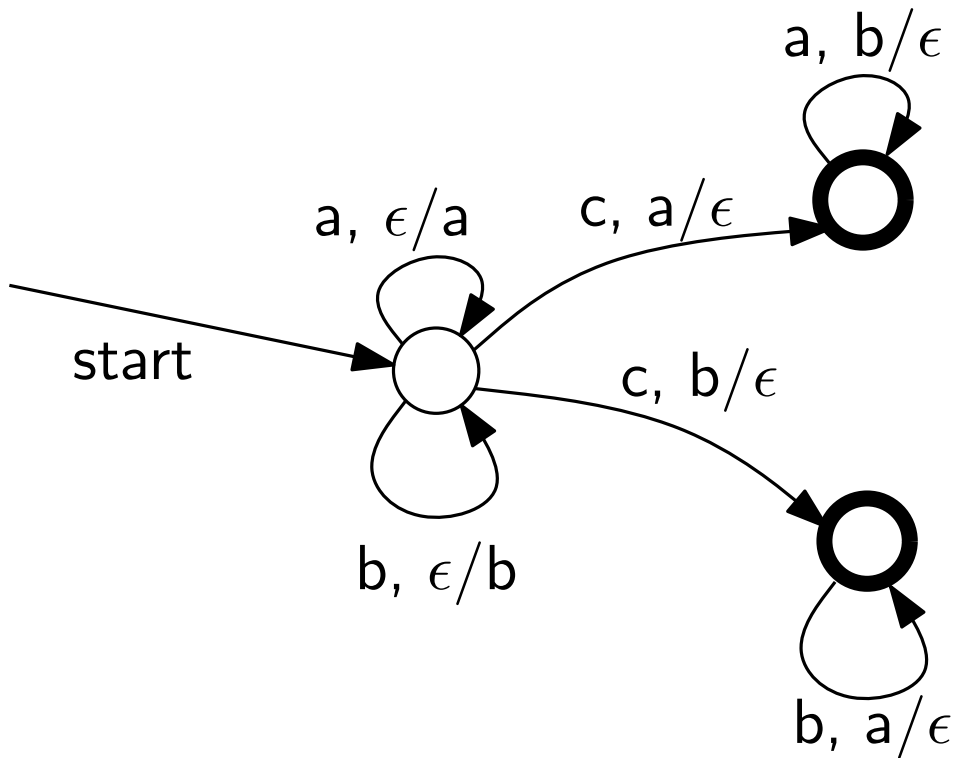
$x, y/z$: läs x , poppa y
från stacken, pusha z



Indata:
baabcbbb

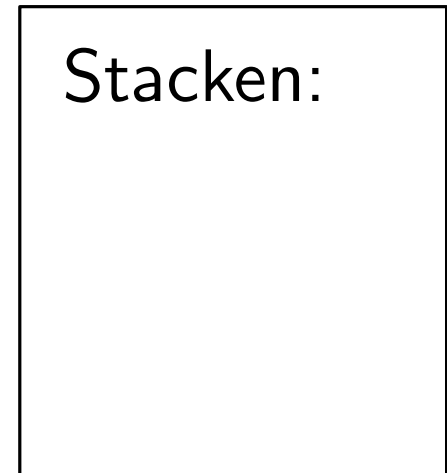
Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



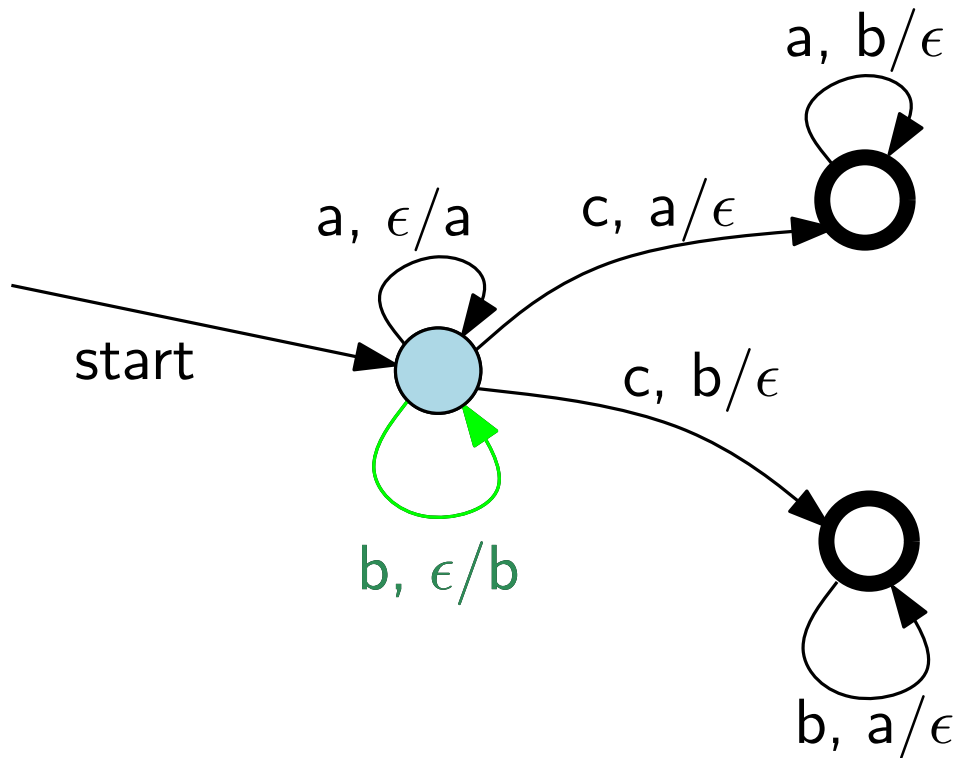
Indata:
baabcbbb

Stacken:



Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



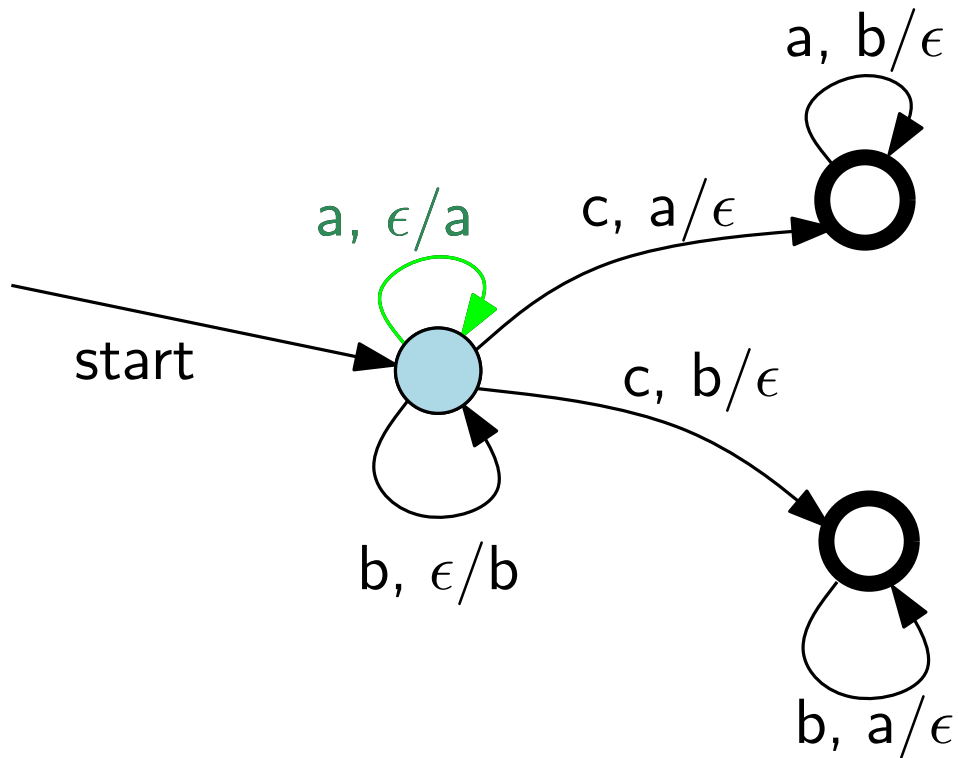
Indata:
baabcbbb

Stacken:

b pushas

Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



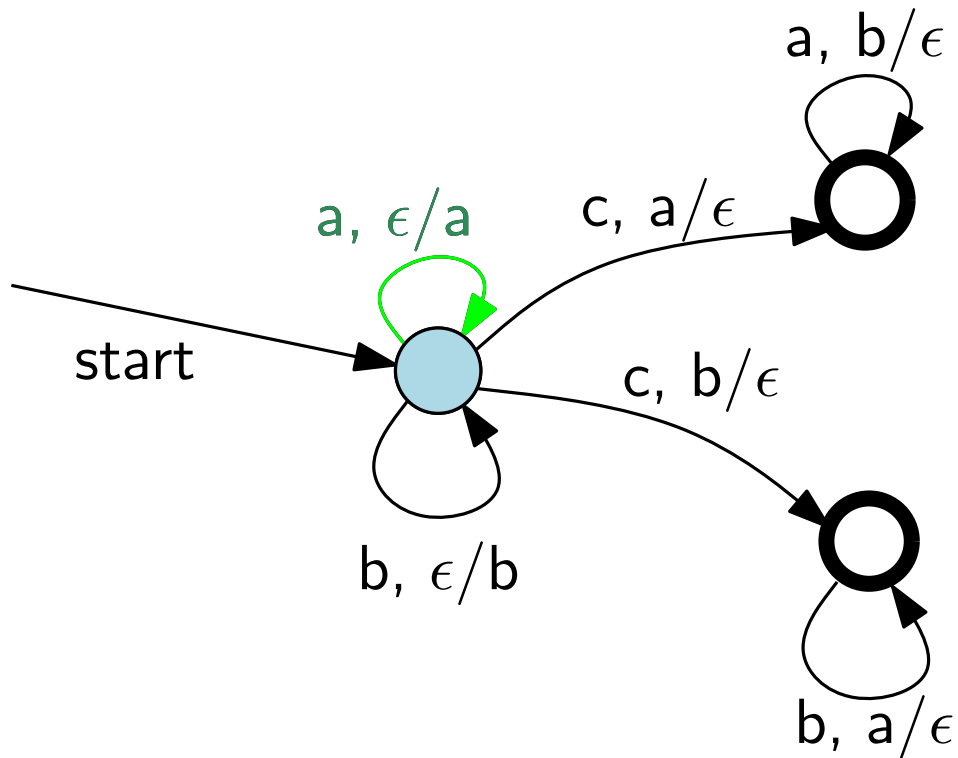
Indata:
baabcbbb

Stacken:

a pushas
b

Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



Indata:
ba**a**bcbbb

Stacken:

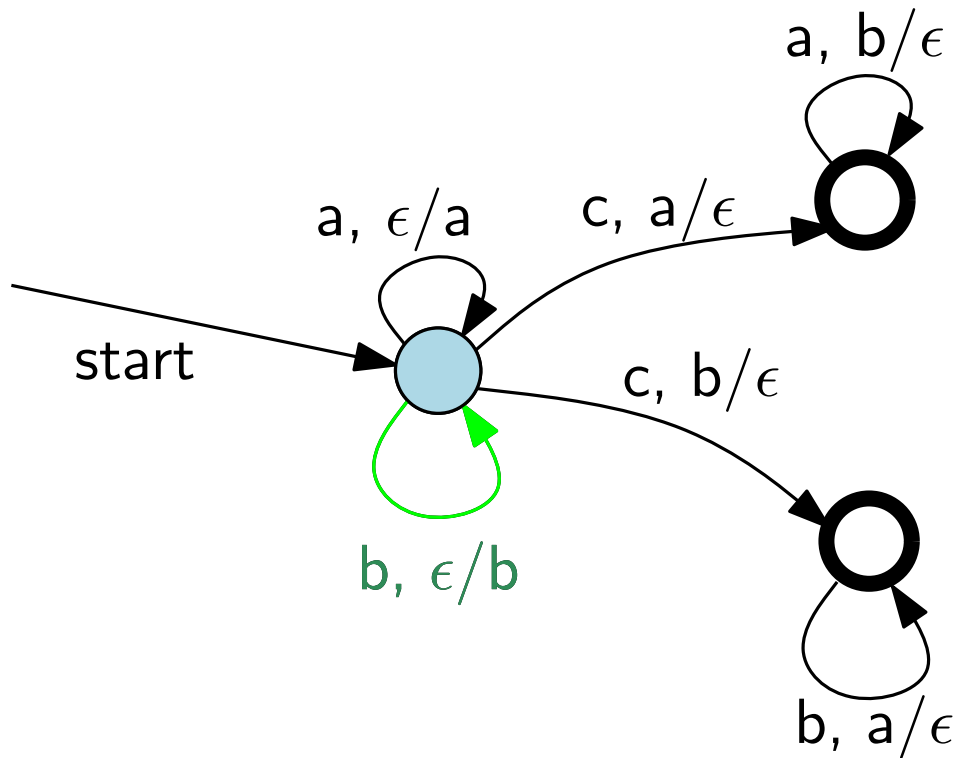
a pushas

a

b

Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z

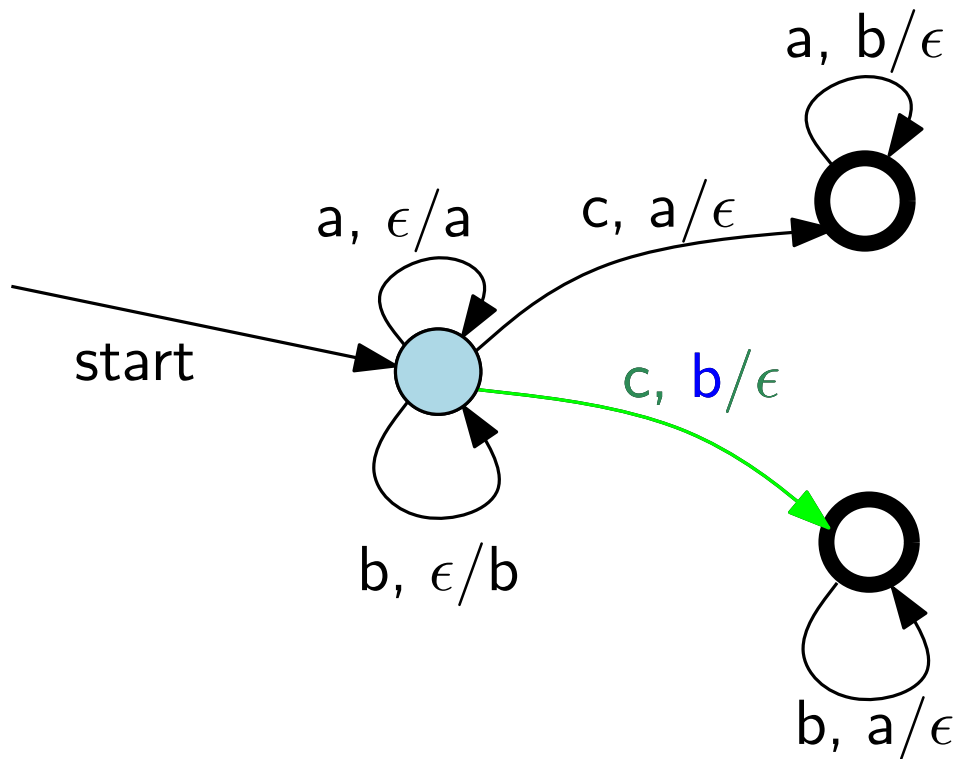


Indata:
baab**c**bbb

Stacken:
b pushas
a
a
b

Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z

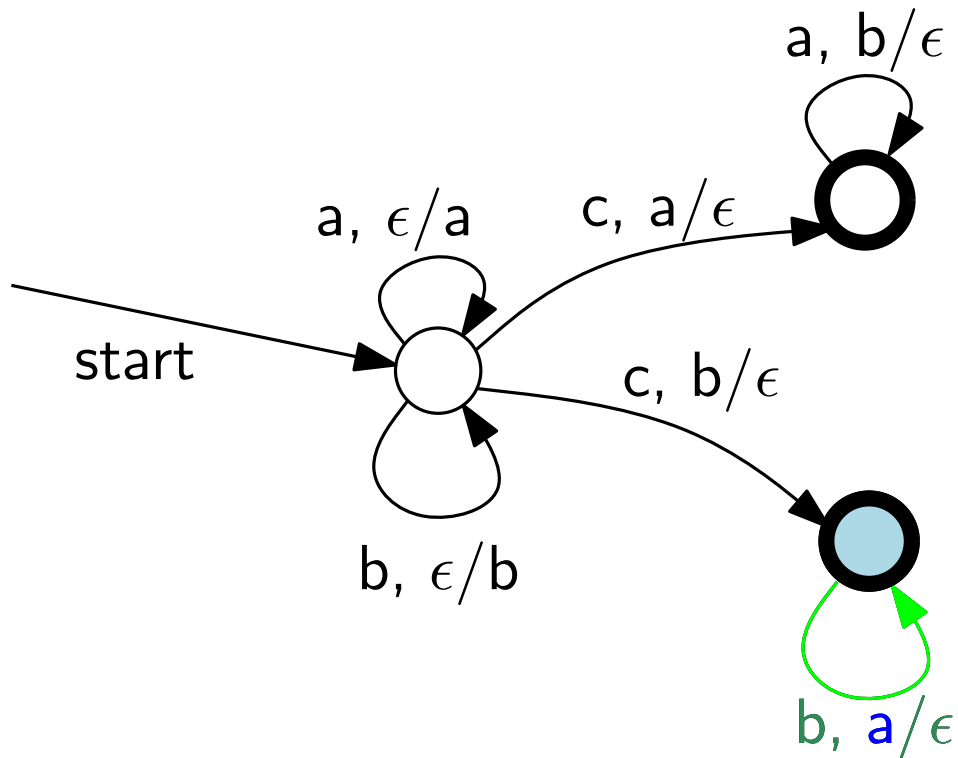


Indata:
baab**c**bbb

Stacken:
b poppas
a
a
b

Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



Indata:
baabc**b**bb

Stacken:

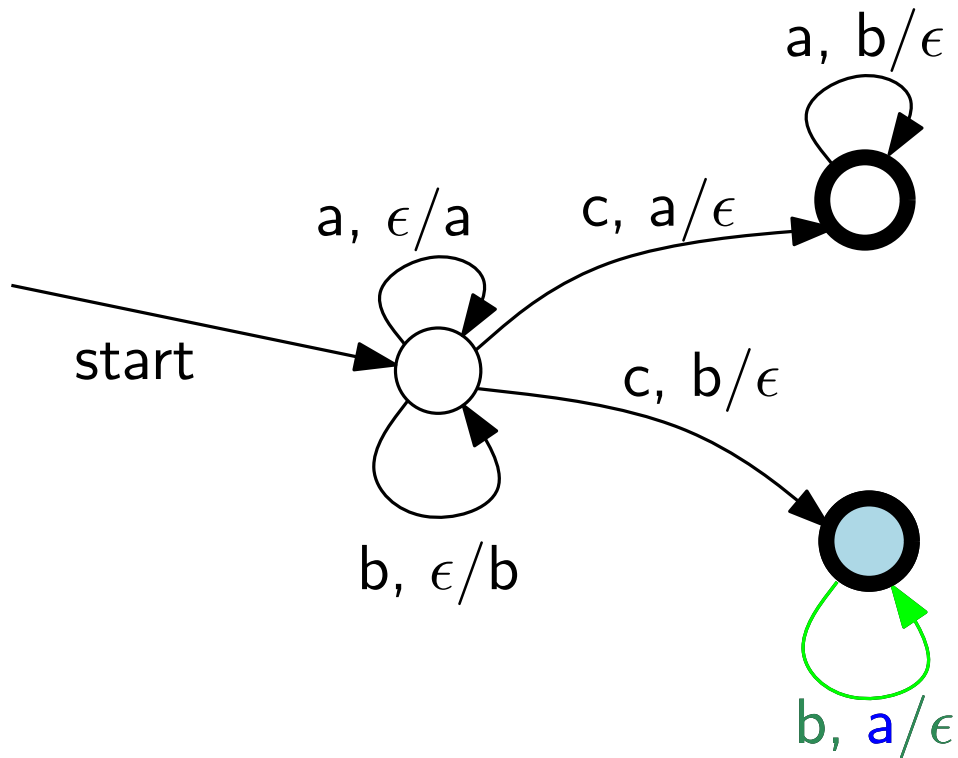
a poppas

a

b

Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



Indata:
baabcb**b**

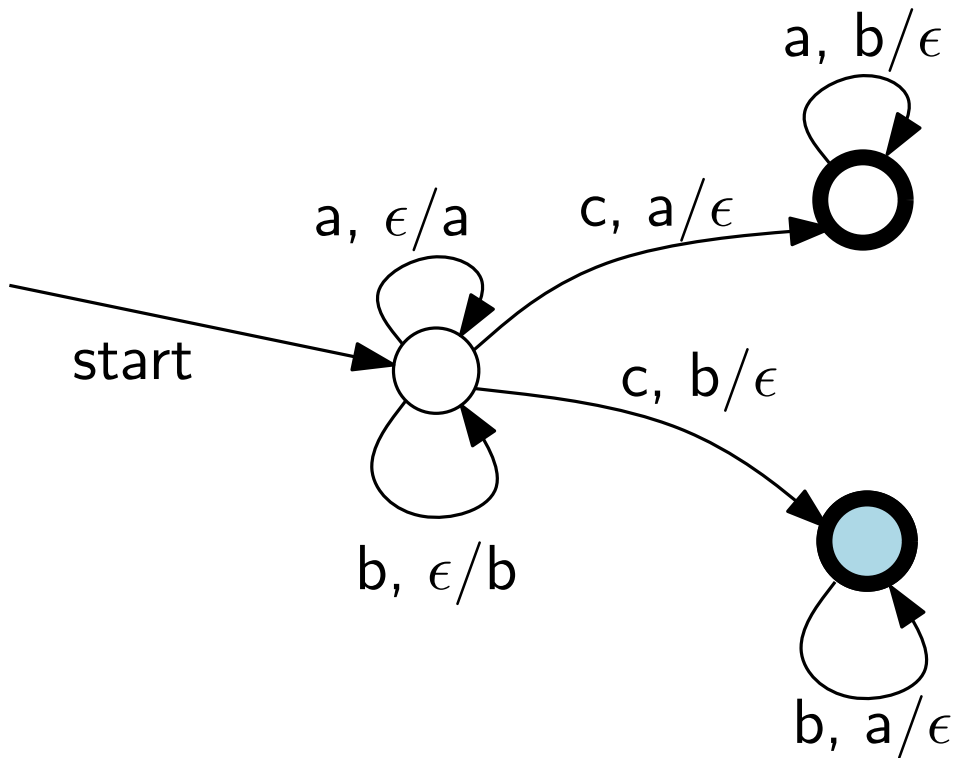
Stacken:

a poppas

b

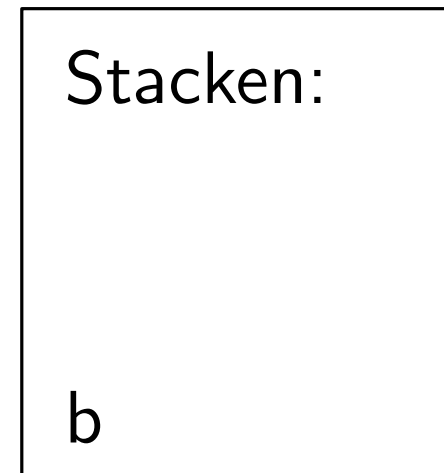
Exempelkörning, stackautomat

$x, y/z$: läs x , poppa y
från stacken, pusha z



Indata:
baabcbb**b**

Övergång saknas, för tecknet b måste
vi ha ett a högst upp på stacken
→ automaten accepterar inte



Stackautomat för balanserade parentesuttryck

En grammatik för balanserade parentesuttryck:

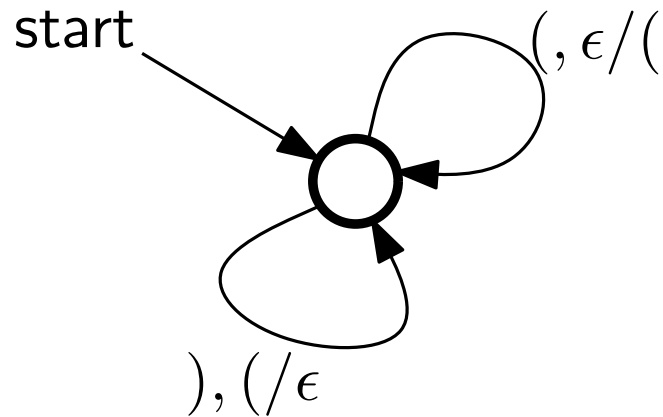
$$\text{Expr} \rightarrow \epsilon \mid (\text{Expr})\text{Expr}$$

Stackautomat för balanserade parentesuttryck

En grammatik för balanserade parentesuttryck:

$$\text{Expr} \rightarrow \epsilon | (\text{Expr})\text{Expr}$$

Vi kan också lätt konstruera en PDA för språket, vi behöver faktiskt bara ett enda tillstånd! (Plus det implicita fail-tillståndet.)

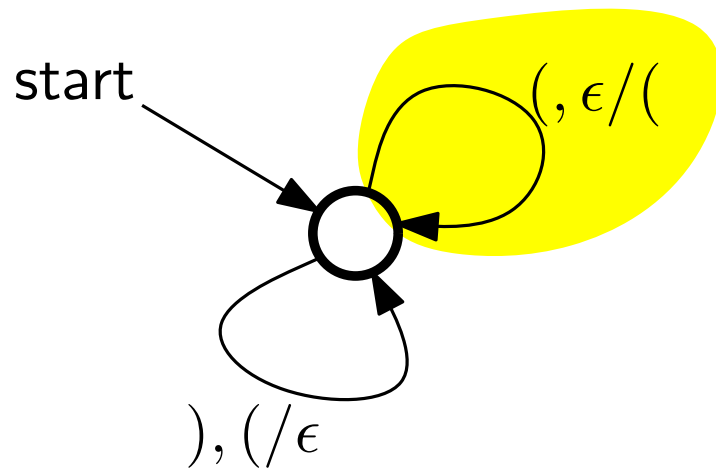


Stackautomat för balanserade parentesuttryck

En grammatik för balanserade parentesuttryck:

$$\text{Expr} \rightarrow \epsilon | (\text{Expr})\text{Expr}$$

Vi kan också lätt konstruera en PDA för språket, vi behöver faktiskt bara ett enda tillstånd! (Plus det implicita fail-tillståndet.)



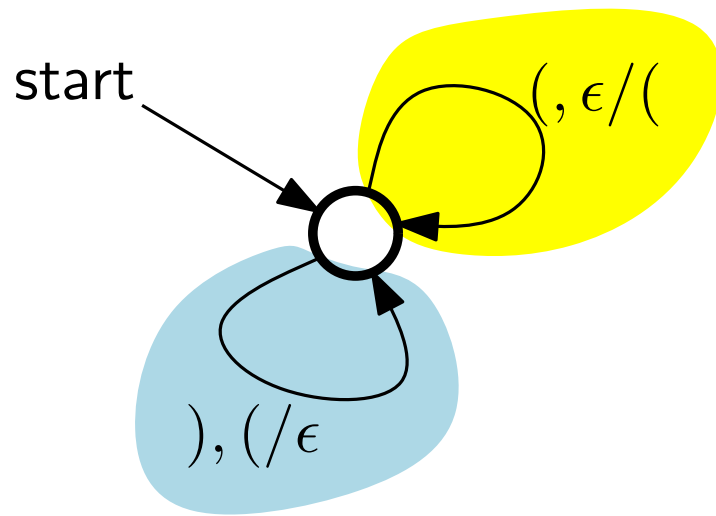
När vi ser en vänsterparentes, pusha på vänsterparentes på stacken

Stackautomat för balanserade parentesuttryck

En grammatik för balanserade parentesuttryck:

$$\text{Expr} \rightarrow \epsilon | (\text{Expr})\text{Expr}$$

Vi kan också lätt konstruera en PDA för språket, vi behöver faktiskt bara ett enda tillstånd! (Plus det implicita fail-tillståndet.)



När vi ser en vänsterparentes, pusha på vänsterparentes på stacken

När vi ser högerparentes, poppa en vänsterparentes från stacken

Från grammatiker till stackautomater

Parentesuttryck var enkelt att bygga stackautomat för.

Från grammatiker till stackautomater

Parentesuttryck var enkelt att bygga stackautomat för.

Kan man alltid konvertera en grammatik till en stackautomat?

Från grammatiker till stackautomater

Parentesuttryck var enkelt att bygga stackautomat för.

Kan man alltid konvertera en grammatik till en stackautomat?

Nej, det kan man inte!

Från grammatiker till stackautomater

Parentesuttryck var enkelt att bygga stackautomat för.

Kan man alltid konvertera en grammatik till en stackautomat?

Nej, det kan man inte!

Exempel: Palindrom över $\{a, b\}$

Palin $\rightarrow \epsilon \mid a \mid b \mid a$ Palin $a \mid b$ Palin b

Det existerar inte någon PDA som känner igen detta språk.

Från grammatiker till stackautomater

Parentesuttryck var enkelt att bygga stackautomat för.

Kan man alltid konvertera en grammatik till en stackautomat?

Nej, det kan man inte!

Exempel: Palindrom över $\{a, b\}$

Palin $\rightarrow \epsilon \mid a \mid b \mid a$ Palin $a \mid b$ Palin b

Det existerar inte någon PDA som känner igen detta språk.

Intuition: PDA:n borde, när den kommit exakt halvvägs in i strängen, börja matcha av tecken mot de som den redan sett, men det finns inget sätt för den att veta när den är halvvägs.

Från grammatiker till stackautomater

Parentesuttryck var enkelt att bygga stackautomat för.

Kan man alltid konvertera en grammatik till en stackautomat?

Nej, det kan man inte!

Exempel: Palindrom över $\{a, b\}$

$\text{Palin} \rightarrow \epsilon \mid a \mid b \mid a \text{ Palin} \mid b \text{ Palin}$

Det existerar inte någon PDA som känner igen detta språk.

Intuition: PDA:n borde, när den kommit exakt halvvägs in i strängen, börja matcha av tecken mot de som den redan sett, men det finns inget sätt för den att veta när den är halvvägs.

Man kan bevisa detta med ett pumping-lemma för PDA:er, liknande beviset vi såg för DFA:er (beviset ingår ej i kursen).

Grammatiker vs. stackautomater

Om stackautomater inte är tillräckligt kraftfulla för att kunna användas till alla grammatiker, vad är poängen?

Grammatiker vs. stackautomater

Om stackautomater inte är tillräckligt kraftfulla för att kunna användas till alla grammatiker, vad är poängen?

Poängen är att de kan användas för “de flesta” grammatiker, t.ex. (nästan?) alltid de man får när man konstruerar programspråk.

Idag

Stackautomater

Olika klasser av språk och grammatiker

Parsergeneratorer

Sammanfattning / Inför KS

Olika kraftfulla grammatiker

De grammatiker vi använt oss av är vad som kallas för *kontextfria grammatiker*

(Det finns mer generella grammatiker som kallas för *kontextkänsliga*)

Olika kraftfulla grammatiker

De grammatiker vi använt oss av är vad som kallas för *kontextfria grammatiker*

(Det finns mer generella grammatiker som kallas för *kontextkänsliga*)

Vi har sett två verktyg för att parsas kontextfria grammatiker:

- Rekursiv medåkning (förra föreläsningen)
- Stackautomater (nyss)

Olika kraftfulla grammatiker

De grammatiker vi använt oss av är vad som kallas för *kontextfria grammatiker*

(Det finns mer generella grammatiker som kallas för *kontextkänsliga*)

Vi har sett två verktyg för att parsas kontextfria grammatiker:

- Rekursiv medåkning (förra föreläsningen)
- Stackautomater (nyss)

Inget av dem är tillräckligt kraftfullt för att kunna hantera alla kontextfria grammatiker, men de är tillräckliga för att hantera de flesta språk man faktiskt vill skriva en parser för.

Olika kraftfulla grammatiker

De grammatiker vi använt oss av är vad som kallas för *kontextfria grammatiker*

(Det finns mer generella grammatiker som kallas för *kontextkänsliga*)

Vi har sett två verktyg för att parse kontextfria grammatiker:

- Rekursiv medåkning (förra föreläsningen)
LL-grammatiker
LL-språk
- Stackautomater (nyss)
Deterministiska kontextfria grammatiker
Deterministiska kontextfria språk

Inget av dem är tillräckligt kraftfullt för att kunna hantera alla kontextfria grammatiker, men de är tillräckliga för att hantera de flesta språk man faktiskt vill skriva en parser för.

Perspektiv

Reguljära
språk

Perspektiv

T.ex. giltiga e-post-adresser



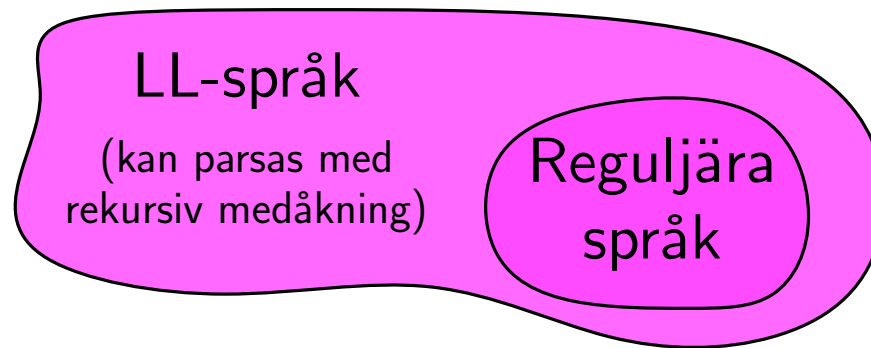
Perspektiv

OBS: varje prick här är *ett språk*, dvs *en mängd strängar*

T.ex. giltiga e-post-adresser

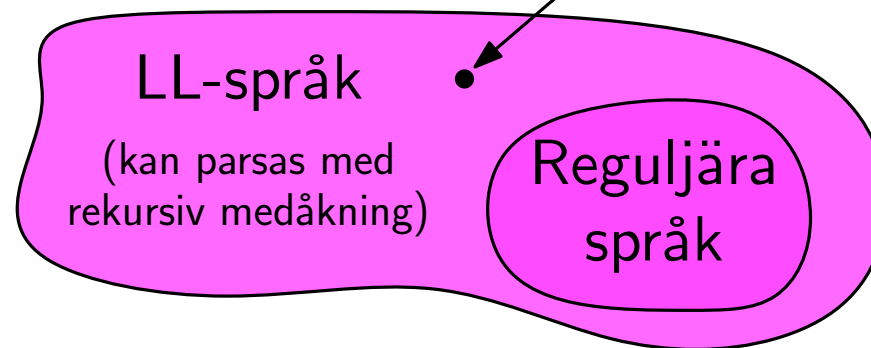


Perspektiv

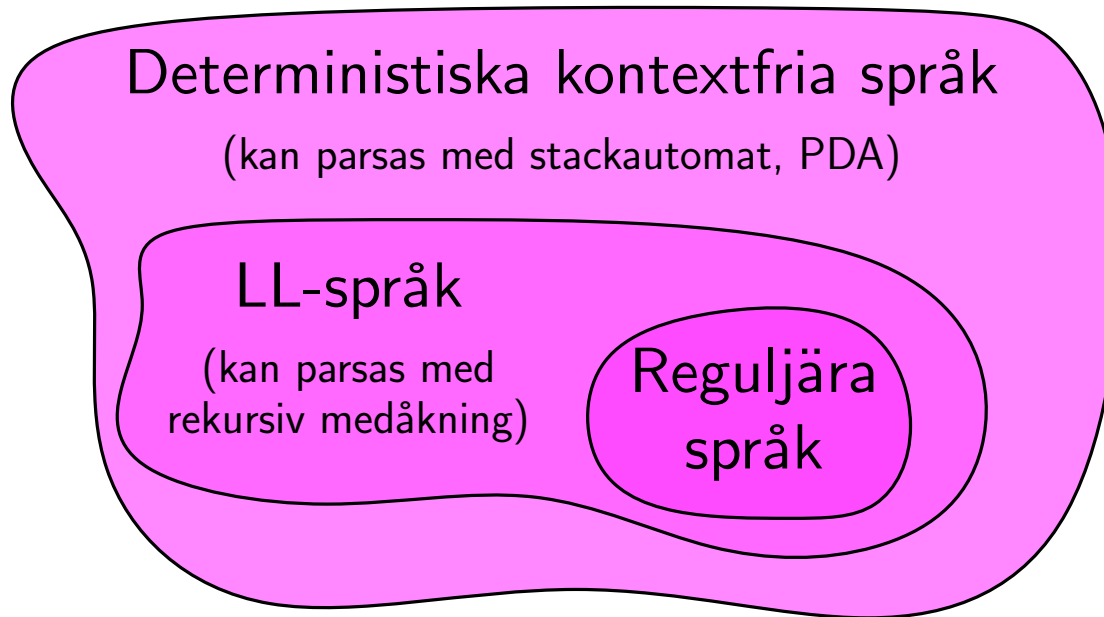


Perspektiv

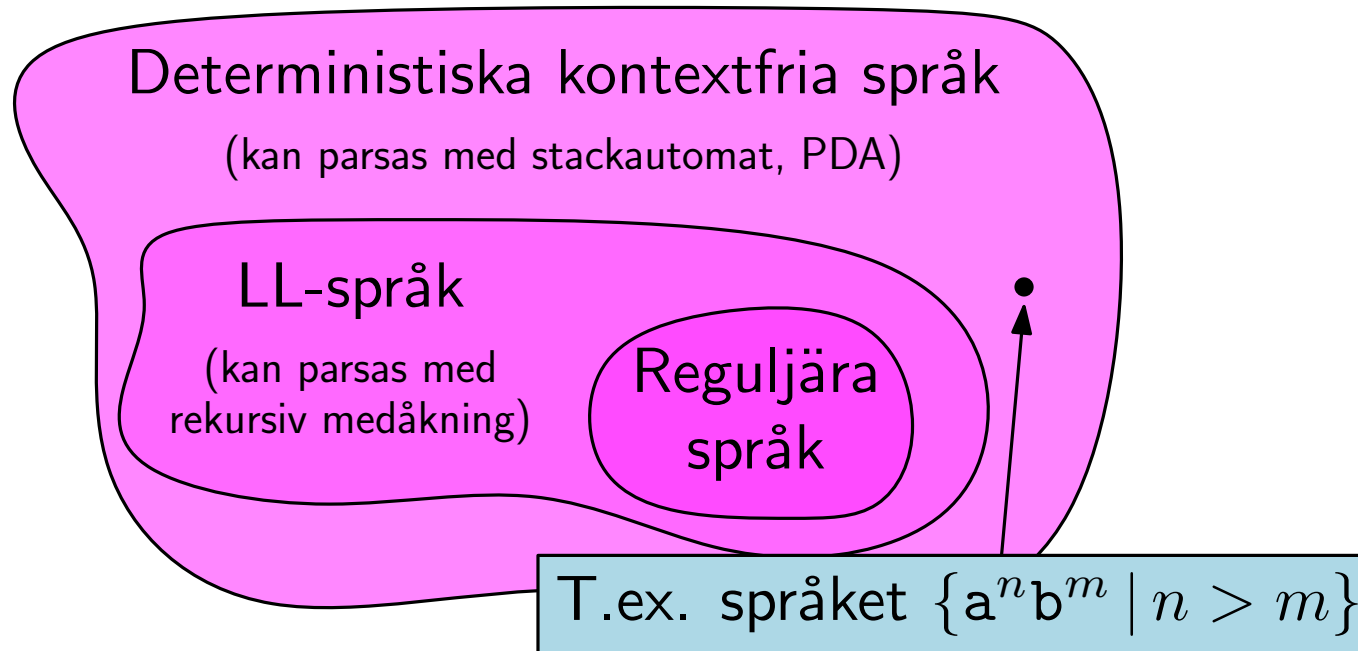
T.ex. balanserade parentesuttryck



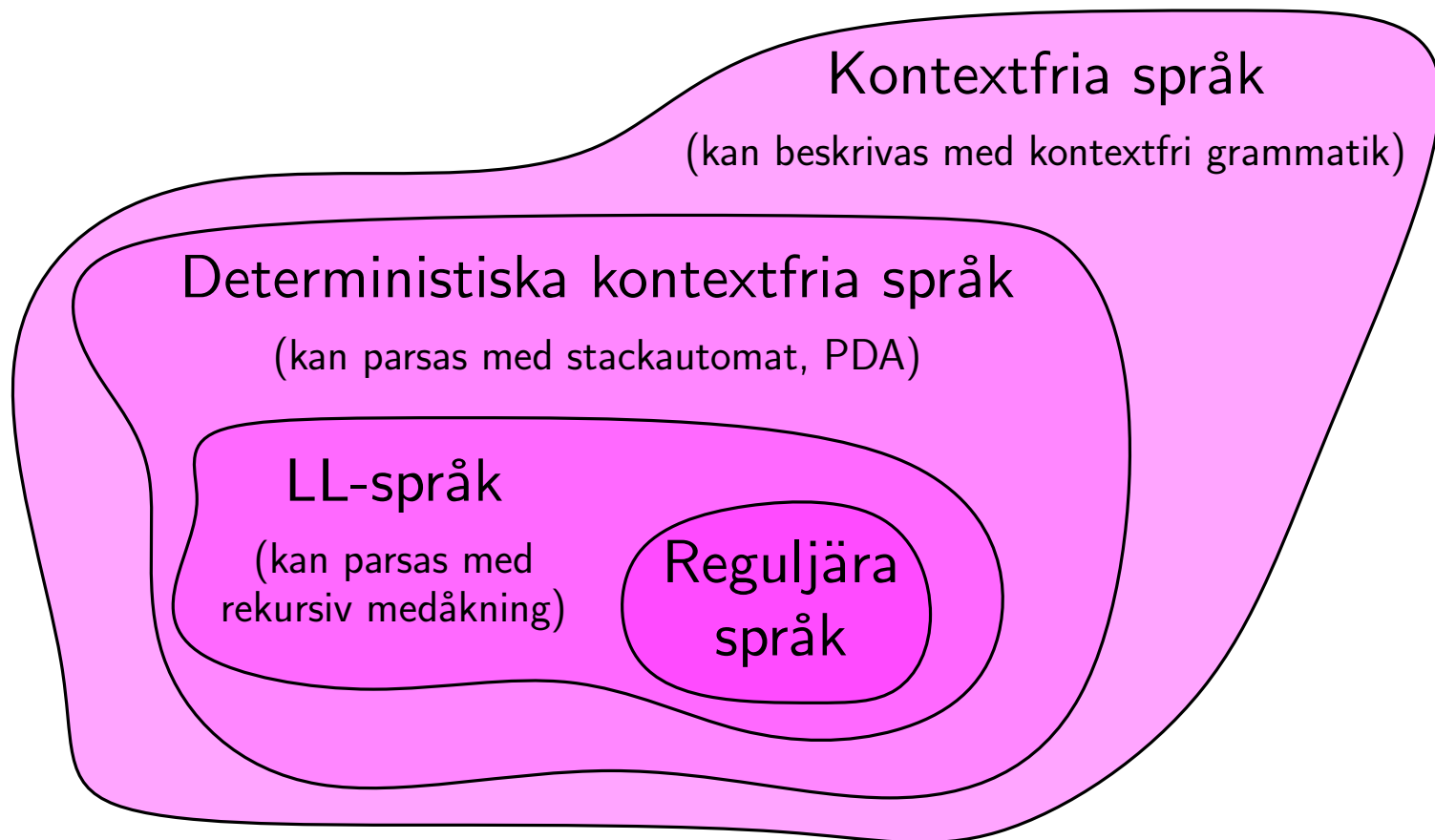
Perspektiv



Perspektiv

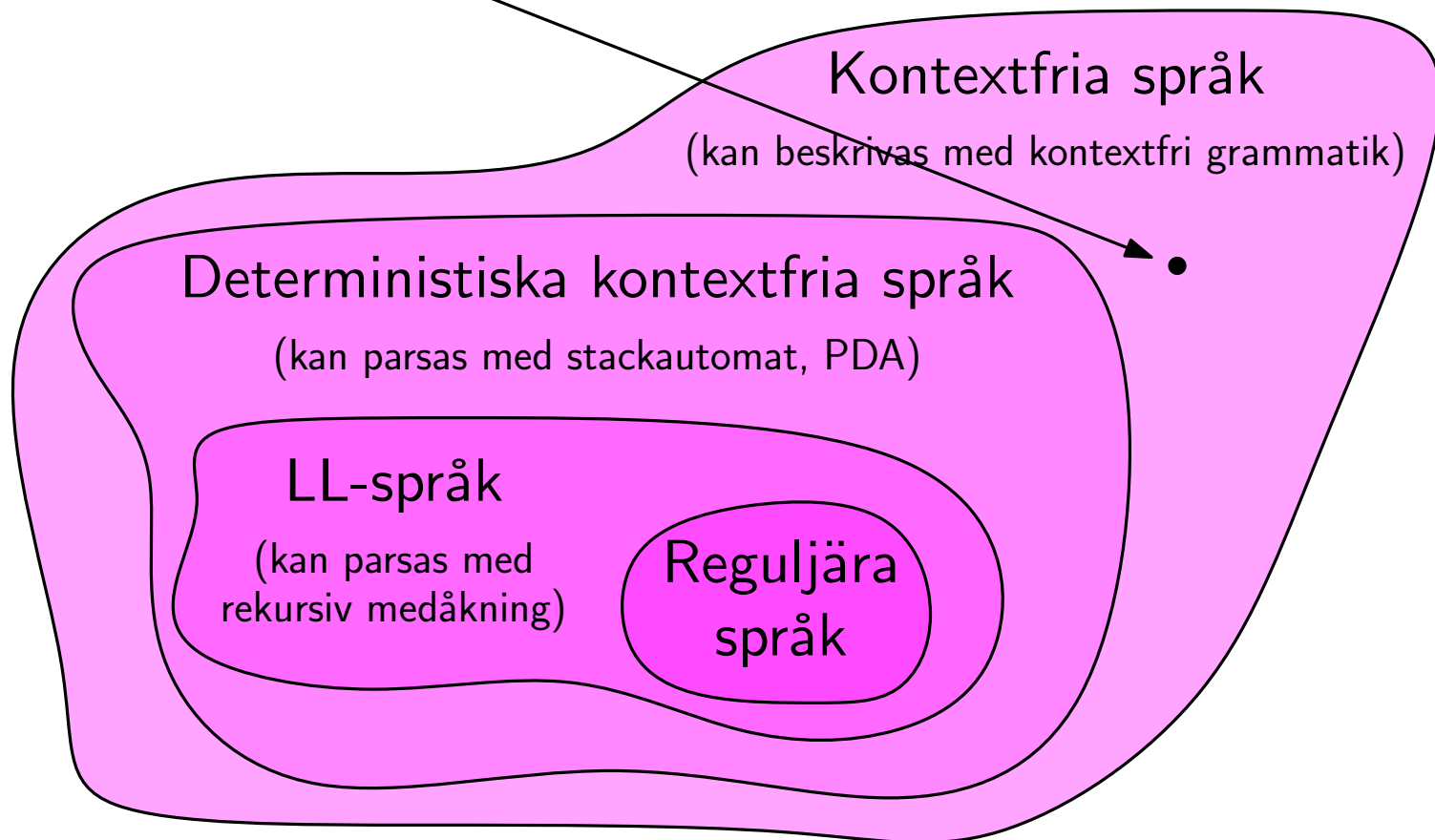


Perspektiv

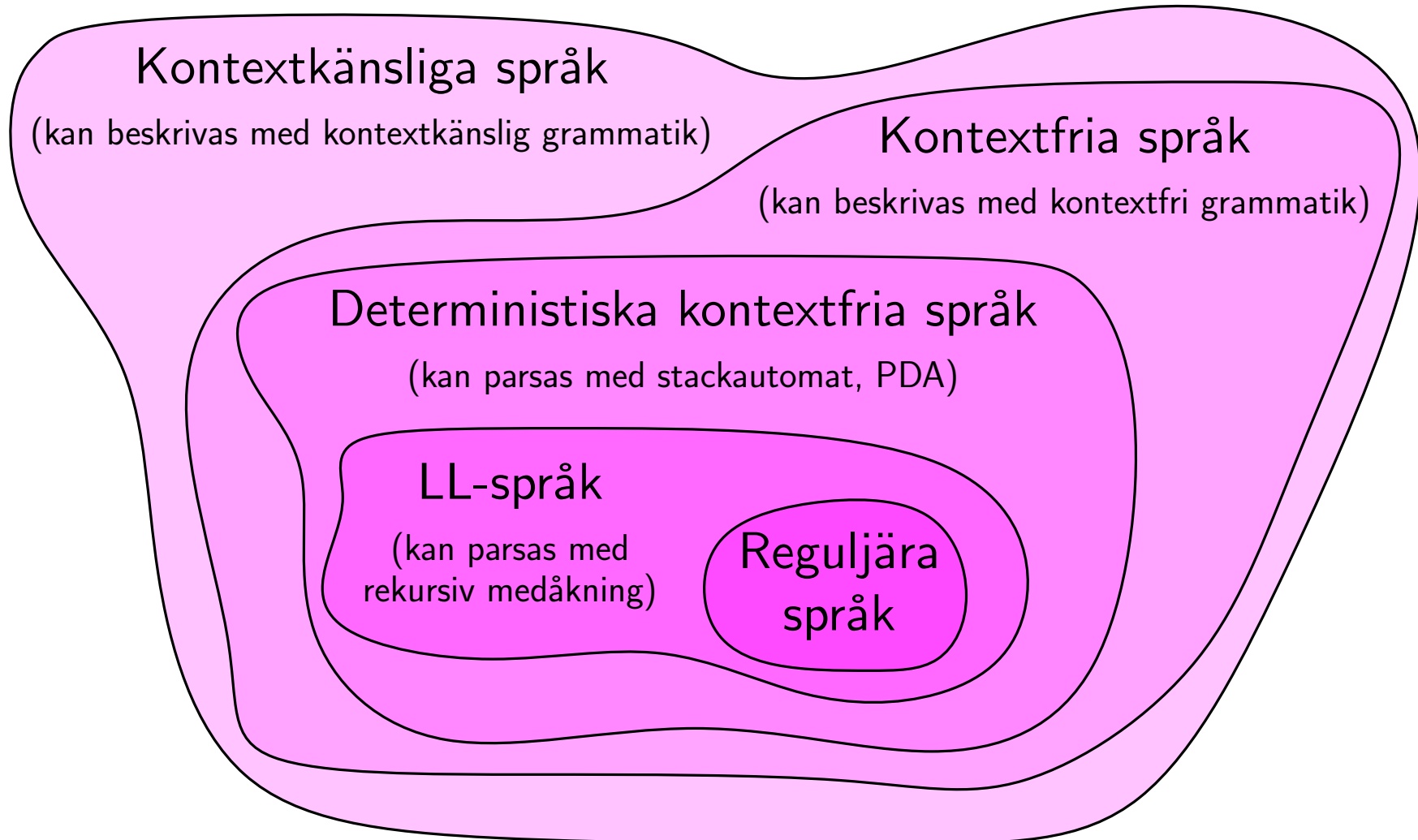


Perspektiv

T.ex. palindrom över $\{a, b\}$

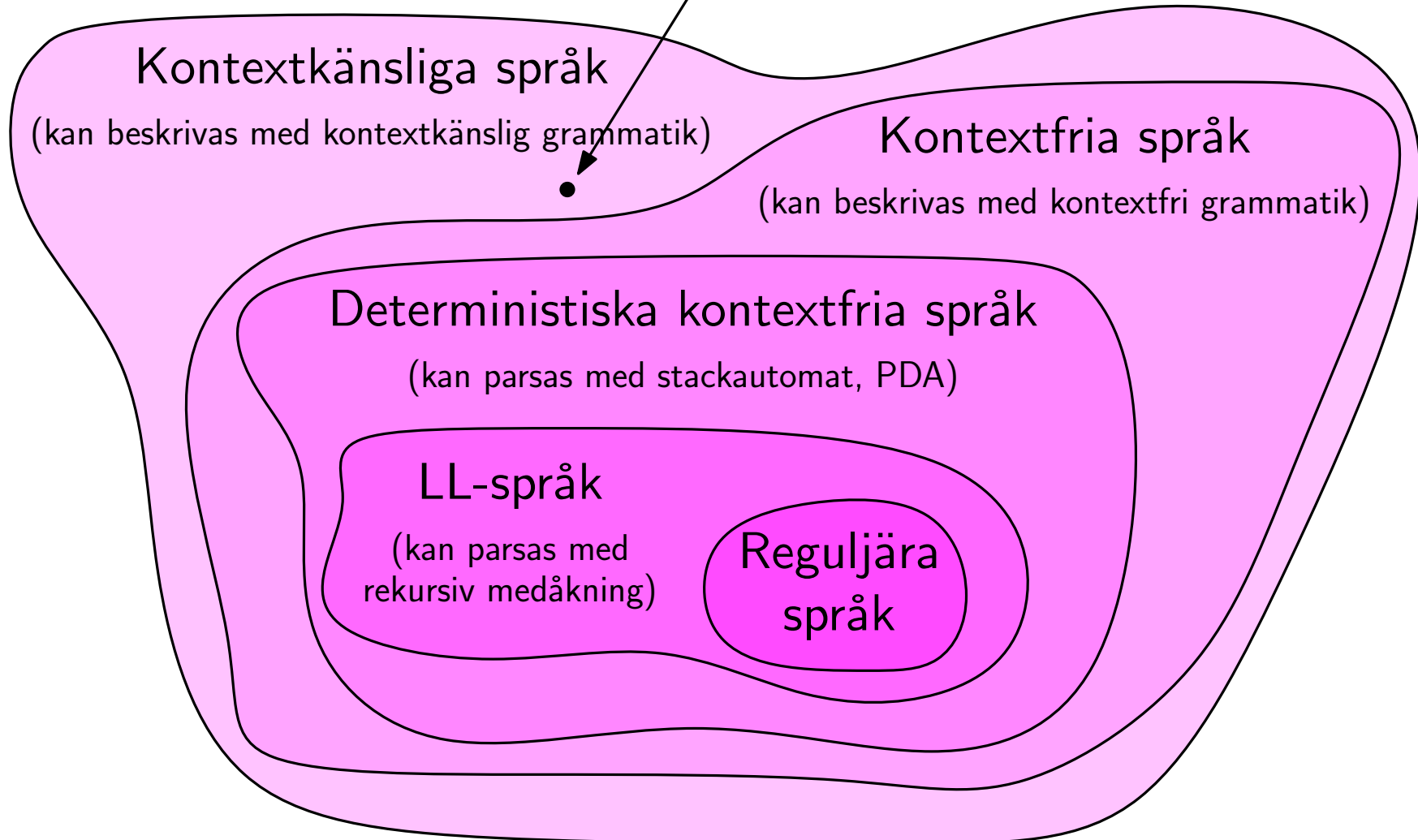


Perspektiv

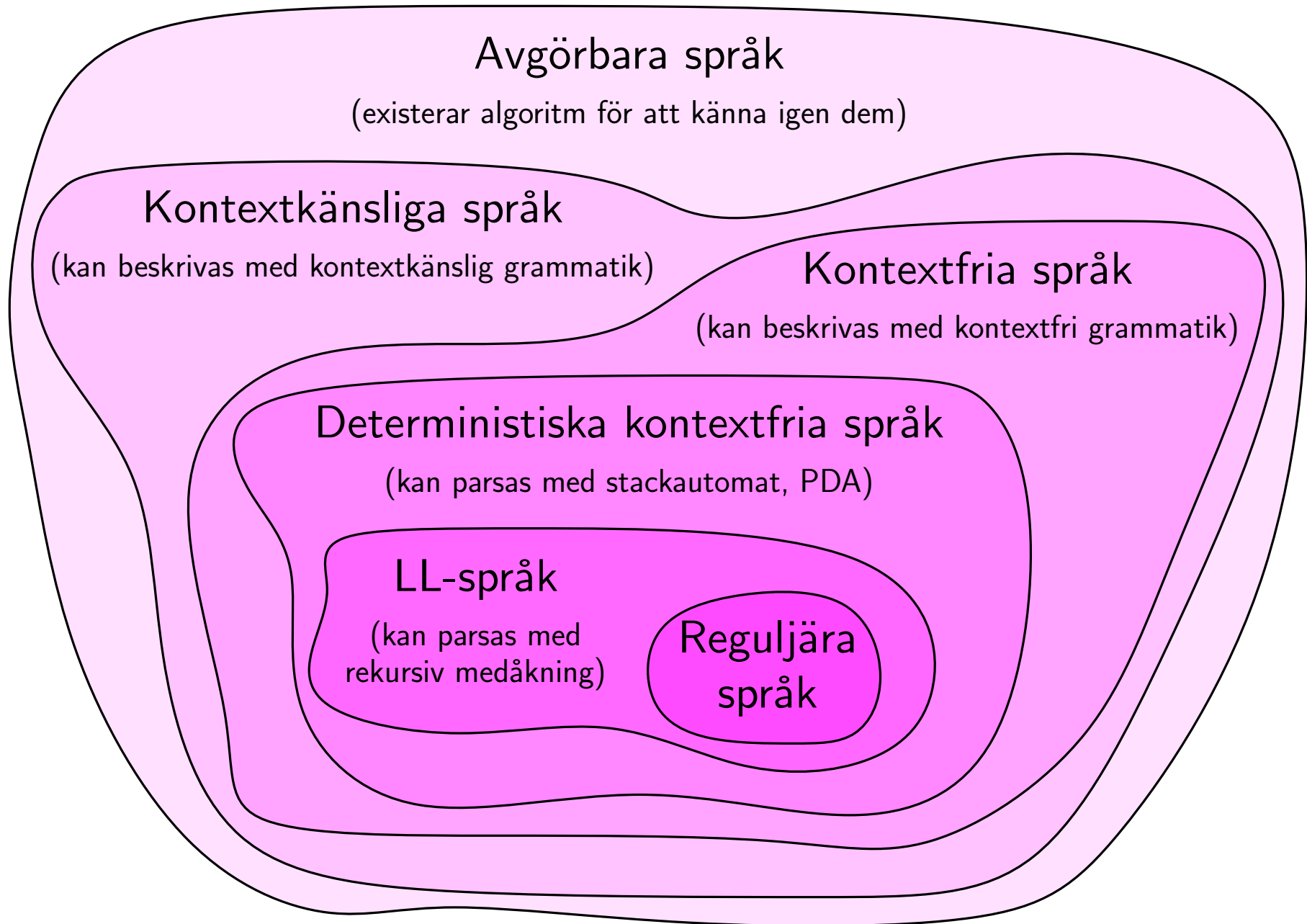


Perspektiv

T.ex. $\{a^n b^n c^n \mid n \geq 1\}$



Perspektiv



Perspektiv

Alla språk

Avgörbara språk

(existerar algoritm för att känna igen dem)

Kontextkänsliga språk

(kan beskrivas med kontextkänslig grammatik)

Kontextfria språk

(kan beskrivas med kontextfri grammatik)

Deterministiska kontextfria språk

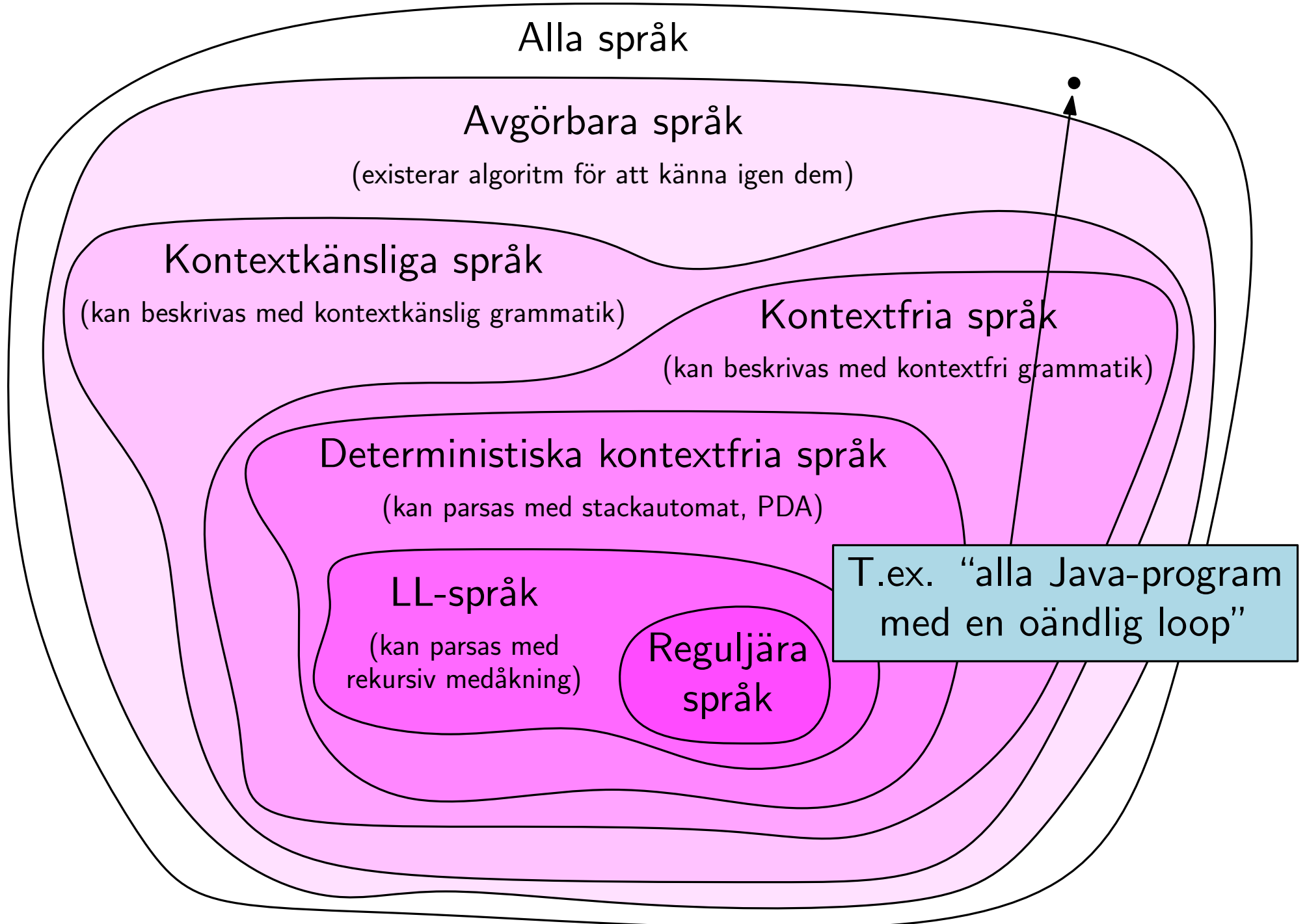
(kan parsas med stackautomat, PDA)

LL-språk

(kan parsas med
rekursiv medåkning)

Reguljära
språk

Perspektiv



Idag

Stackautomater

Olika klasser av språk och grammatiker

Parsegeneratorer

Sammanfattning / Inför KS

Parsegenereratorer

De flesta vill nog inte sätta sig ner och skriva en parser för ett helt programmeringsspråk för hand.

Parsegenereratorer

De flesta vill nog inte sätta sig ner och skriva en parser för ett helt programmeringsspråk för hand.

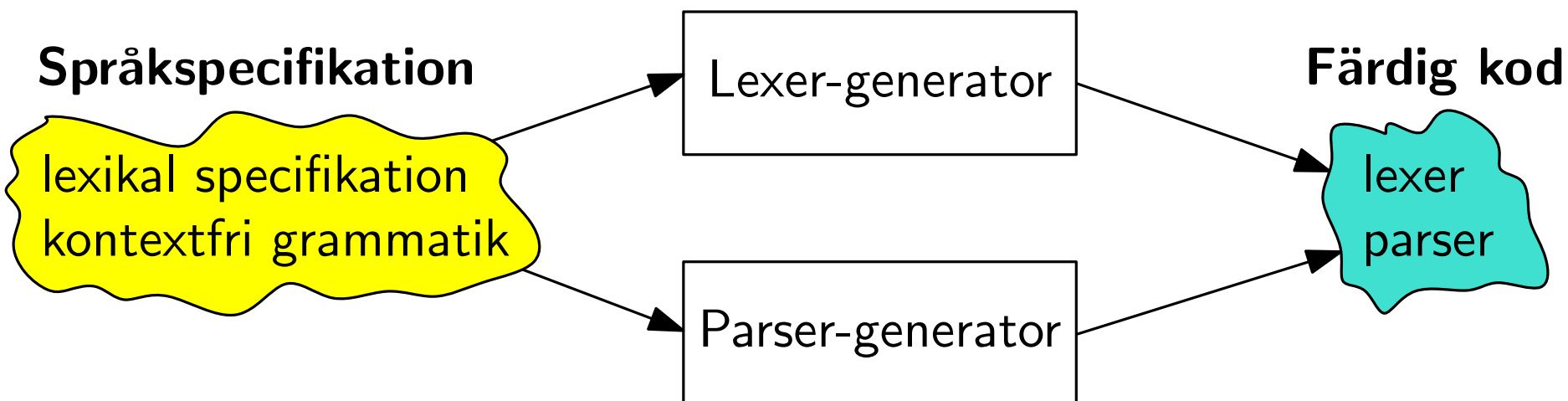
Ofta ganska mekaniskt men många smådetaljer man lätt kan göra fel på → perfekt lämpat för en dator att göra istället!

Parsergeneratorer

De flesta vill nog inte sätta sig ner och skriva en parser för ett helt programmeringsspråk för hand.

Ofta ganska mekaniskt men många smådetaljer man lätt kan göra fel på → perfekt lämpat för en dator att göra istället!

Det finns många verktyg som automatiskt genererar en lexer och en parser från en språkspecifikation.



Exempel på parsergeneratorer

- Lex/Flex och Yacc/Bison, klassiska unix-verktyg, producerar C/C++-kod
(hanterar de flesta deterministiska kontextfria grammatiker – specifikt, en delmängd som kallas LALR-grammatiker)
- JFlex och Cup, Java-versioner av Flex och Yacc
- DCG i Prolog, inbyggt i språket
(hanterar alla kontextfria grammatiker, men använder backtracking i Prolog så kan vara långsamt)
- ANTLR, kan generera LL-parser i många olika språk (C/C++, C#, Java, JavaScript, Python...)
- Parsec i Haskell

Exempel på parsergeneratorer

- Lex/Flex och Yacc/Bison, klassiska unix-verktyg, producerar C/C++-kod
(hanterar de flesta deterministiska kontextfria grammatiker – specifikt, en delmängd som kallas LALR-grammatiker)
- JFlex och Cup, Java-versioner av Flex och Yacc
- DCG i Prolog, inbyggt i språket
(hanterar alla kontextfria grammatiker, men använder backtracking i Prolog så kan vara långsamt)
- ANTLR, kan generera LL-parser i många olika språk (C/C++, C#, Java, JavaScript, Python...)
- Parsec i Haskell

Binära träd igen

Förra veckan skrev vi en rekursiv medåknings-parser för binära träd, definierade enligt följande grammatik.

BinTree \rightarrow Leaf LParen Number RParen |

Branch LParen BinTree Comma BinTree RParen

Slutsymboler:

Leaf: strängen "leaf"

Branch: strängen "branch"

Number: [0-9]+

LParen, RParen, Comma: parenteser och kommatecken

Binära träd i Cup

Fil Parser.cup:

```
import java_cup.runtime.*;

// Deklarera slut-symboler
terminal BRANCH;
terminal LEAF;
terminal LPAREN;
terminal RPAREN;
terminal COMMA;
terminal Integer NUM; // Ett tal har data av typen Integer

// Deklarera icke-slutsymboler, bara en i det här fallet
non terminal ParseTree BinTree; // Har data av typen ParseTree

BinTree ::= LEAF LPAREN NUM:t RPAREN
           { : RESULT = new LeafNode(t); : }
           | BRANCH LPAREN BinTree:left COMMA BinTree:right RPAREN
           { : RESULT = new BranchNode(left, right); : }
           ;
```

Binära träd i Cup

Fil Parser.cup:

```
import java_cup.runtime.*;
```

```
// Deklarera slut-symboler
terminal BRANCH;
terminal LEAF;
terminal LPAREN;
terminal RPAREN;
terminal COMMA;
terminal Integer NUM; // Ett tal har data av typen Integer
```

```
// Deklarera icke-slutsymboler, bara en i det här fallet
non terminal ParseTree BinTree; // Har data av typen ParseTree
```

```
BinTree ::= LEAF LPAREN NUM:t RPAREN
           { : RESULT = new LeafNode(t); : }
         | BRANCH LPAREN BinTree:left COMMA BinTree:right RPAREN
           { : RESULT = new BranchNode(left, right); : }
         ;
```

Binära träd i Cup

Fil Parser.cup:

```
import java_cup.runtime.*;
```

```
// Deklarera slut-symboler  
terminal BRANCH;  
terminal LEAF;  
terminal LPAREN;  
terminal RPAREN;  
terminal COMMA;  
terminal Integer NUM; // Ett tal har data av typen Integer
```

```
// Deklarera icke-slutsymboler, bara en i det här fallet  
non terminal ParseTree BinTree; // Har data av typen ParseTree
```

```
BinTree ::= LEAF LPAREN NUM:t RPAREN  
          { : RESULT = new LeafNode(t); : }  
          | BRANCH LPAREN BinTree:left COMMA BinTree:right RPAREN  
          { : RESULT = new BranchNode(left, right); : }  
          ;
```

Binära träd i Cup

Fil Parser.cup:

```
import java_cup.runtime.*;
```

```
// Deklarera slut-symboler
terminal BRANCH;
terminal LEAF;
terminal LPAREN;
terminal RPAREN;
terminal COMMA;
terminal Integer NUM; // Ett tal har data av typen Integer
```

```
// Deklarera icke-slutsymboler, bara en i det här fallet
non terminal ParseTree BinTree; // Har data av typen ParseTree
```

```
BinTree ::= LEAF LPAREN NUM:t RPAREN
           { : RESULT = new LeafNode(t); : }
         | BRANCH LPAREN BinTree:left COMMA BinTree:right RPAREN
           { : RESULT = new BranchNode(left, right); : }
         ;
```

Binära träd i Cup

Fil Parser.cup:

```
import java_cup.runtime.*;
```

```
// Deklarera slut-symboler
terminal BRANCH;
terminal LEAF;
terminal LPAREN;
terminal RPAREN;
terminal COMMA;
terminal Integer NUM; // Ett tal har data av typen Integer
```

```
// Deklarera icke-slutsymboler, bara en i det här fallet
non terminal ParseTree BinTree; // Har data av typen ParseTree
```

```
BinTree ::= LEAF LPAREN NUM:t RPAREN
           { : RESULT = new LeafNode(t); : }
           | BRANCH LPAREN BinTree:left COMMA BinTree:right RPAREN
           { : RESULT = new BranchNode(left, right); : }
           ;
```

Kör “cup -parser Parser Parser.cup”, genererar
Parser.java och sym.java

Lexikal analys för binära träd i JFlex

Fil Lexer.lex:

```
import java.lang.System;
import java_cup.runtime.Symbol;

%%
%cup
%class Lexer

%%

branch { return new Symbol(sym.BRANCH); }
leaf { return new Symbol(sym.LEAF); }
"(" { return new Symbol(sym.LPAREN); }
")" { return new Symbol(sym.RPAREN); }
, { return new Symbol(sym.COMMA); }
[0-9]+ { return new Symbol(sym.NUM, Integer.parseInt(ytext())); }
[ \t\n] { }
```


Lexikal analys för binära träd i JFlex

Fil Lexer.lex:

```
import java.lang.System;  
import java_cup.runtime.Symbol;
```

```
%%  
%cup  
%class Lexer
```

```
%%
```

```
branch { return new Symbol(sym.BRANCH); }  
leaf { return new Symbol(sym.LEAF); }  
"(" { return new Symbol(sym.LPAREN); }  
")" { return new Symbol(sym.RPAREN); }  
, { return new Symbol(sym.COMMA); }  
[0-9]+ { return new Symbol(sym.NUM, Integer.parseInt(yytext())); }  
[ \t\n] { }
```

Lexikal analys för binära träd i JFlex

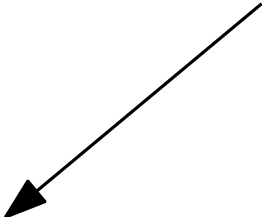
Fil Lexer.lex:

```
import java.lang.System;  
import java_cup.runtime.Symbol;
```

```
%%  
%cup  
%class Lexer
```

Slutsymbolsnamnen från Parser.cup

```
%%
```



```
branch { return new Symbol(sym.BRANCH); }  
leaf { return new Symbol(sym.LEAF); }  
"(" { return new Symbol(sym.LPAREN); }  
")" { return new Symbol(sym.RPAREN); }  
, { return new Symbol(sym.COMMA); }  
[0-9]+ { return new Symbol(sym.NUM, Integer.parseInt(ytext())); }  
[ \t\n] { }
```

Lexikal analys för binära träd i JFlex

Fil Lexer.lex:

```
import java.lang.System;  
import java_cup.runtime.Symbol;
```

```
%%  
%cup  
%class Lexer
```

Slutsymbolsnamnen från Parser.cup

```
%%  
  
branch { return new Symbol(sym.BRANCH); }  
leaf { return new Symbol(sym.LEAF); }  
"(" { return new Symbol(sym.LPAREN); }  
")" { return new Symbol(sym.RPAREN); }  
, { return new Symbol(sym.COMMA); }  
[0-9]+ { return new Symbol(sym.NUM, Integer.parseInt(ytext())); }  
[ \t\n] { }
```

Kör "jflex Lexer.lex", genererar Lexer.java

Lexikal analys för binära träd i JFlex

Fil Lexer.lex:

```
import java.lang.System;
import java_cup.runtime.Symbol;
```

```
%%
%cup
%class Lexer
```

Slutsymbolsnamnen från Parser.cup

```
%%
```

```
branch { return new Symbol(sym.BRANCH); }
leaf { return new Symbol(sym.LEAF); }
"(" { return new Symbol(sym.LPAREN); }
")" { return new Symbol(sym.RPAREN); }
, { return new Symbol(sym.COMMA); }
[0-9]+ { return new Symbol(sym.NUM, Integer.parseInt(ytext())); }
[ \t\n] { }
```

Kör “jflex Lexer.lex”, genererar Lexer.java

Fullständigt exempel med Main-klass upplagt på kurshemsidan

Idag

Stackautomater

Olika klasser av språk och grammatiker

Parsergeneratorer

Sammanfattning / Inför KS

Sammanfattning

Formella språk

- Språkklasser (ex. reguljära språk)
- Informella språkbeskrivningar (ex. “alla e-post-adresser”)
- Formella språkbeskrivningar/algorithm (ex. automater, reguljära uttryck)

Reguljära språk

- Reguljära uttryck och DFA
- Skillnad mellan reguljära språk och regex
- Begränsningar (språk som inte är reguljära)

Kontextfria språk

- Kontextfria grammatiker
- Stackautomater, PDA
- Härledningar/syntaxträd
- Tvetydighet

Lexikal analys

Rekursiv medåkning

Förkortningar

DFA *Deterministic Finite Automaton*, ändlig automat.

Den enklaste typen av automat, lika uttrycksfull som reguljära uttryck.

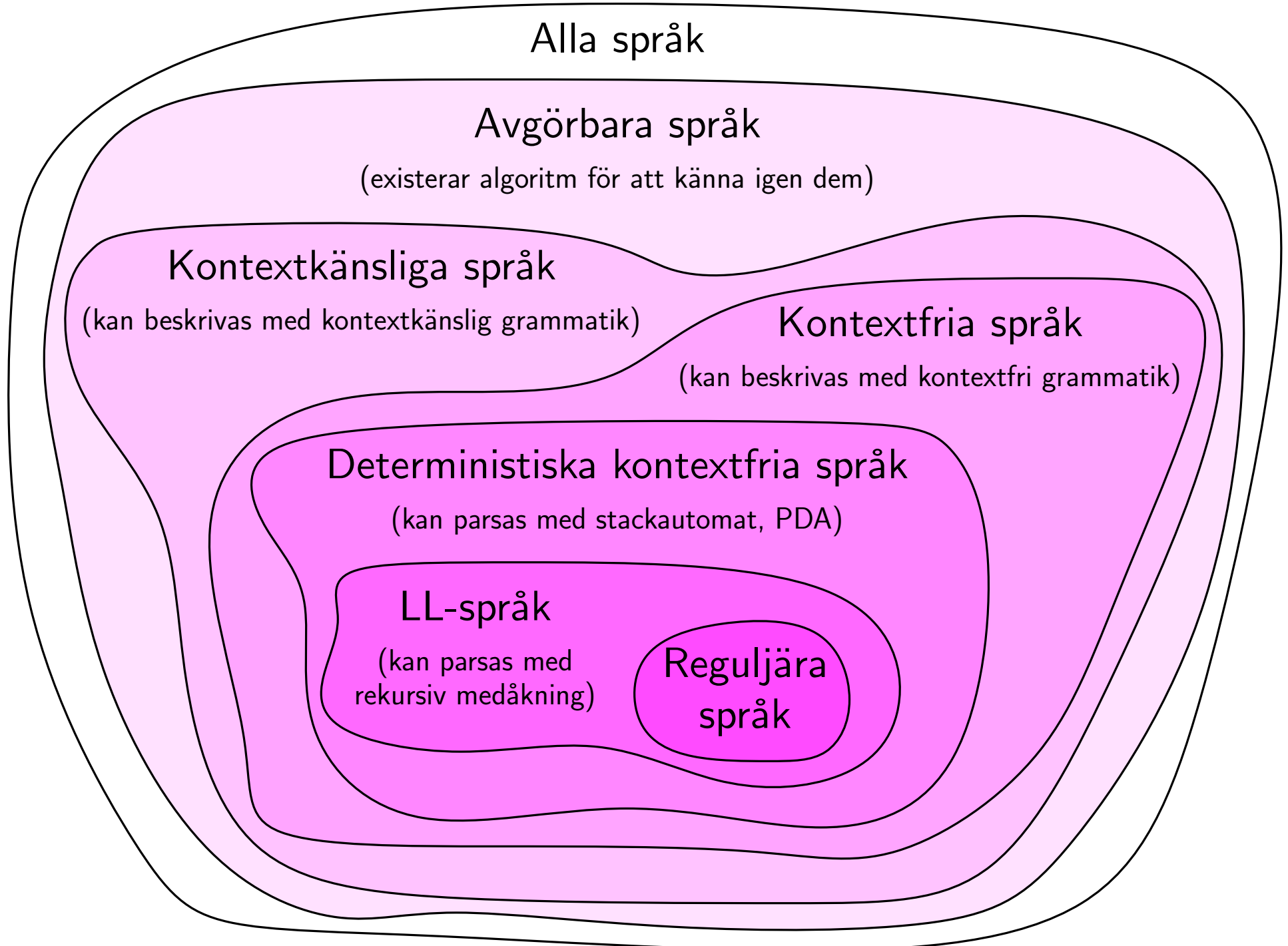
PDA *Push-Down Automaton*, stackautomat.

Som DFA fast med oändligt stort minne i form av en stack.

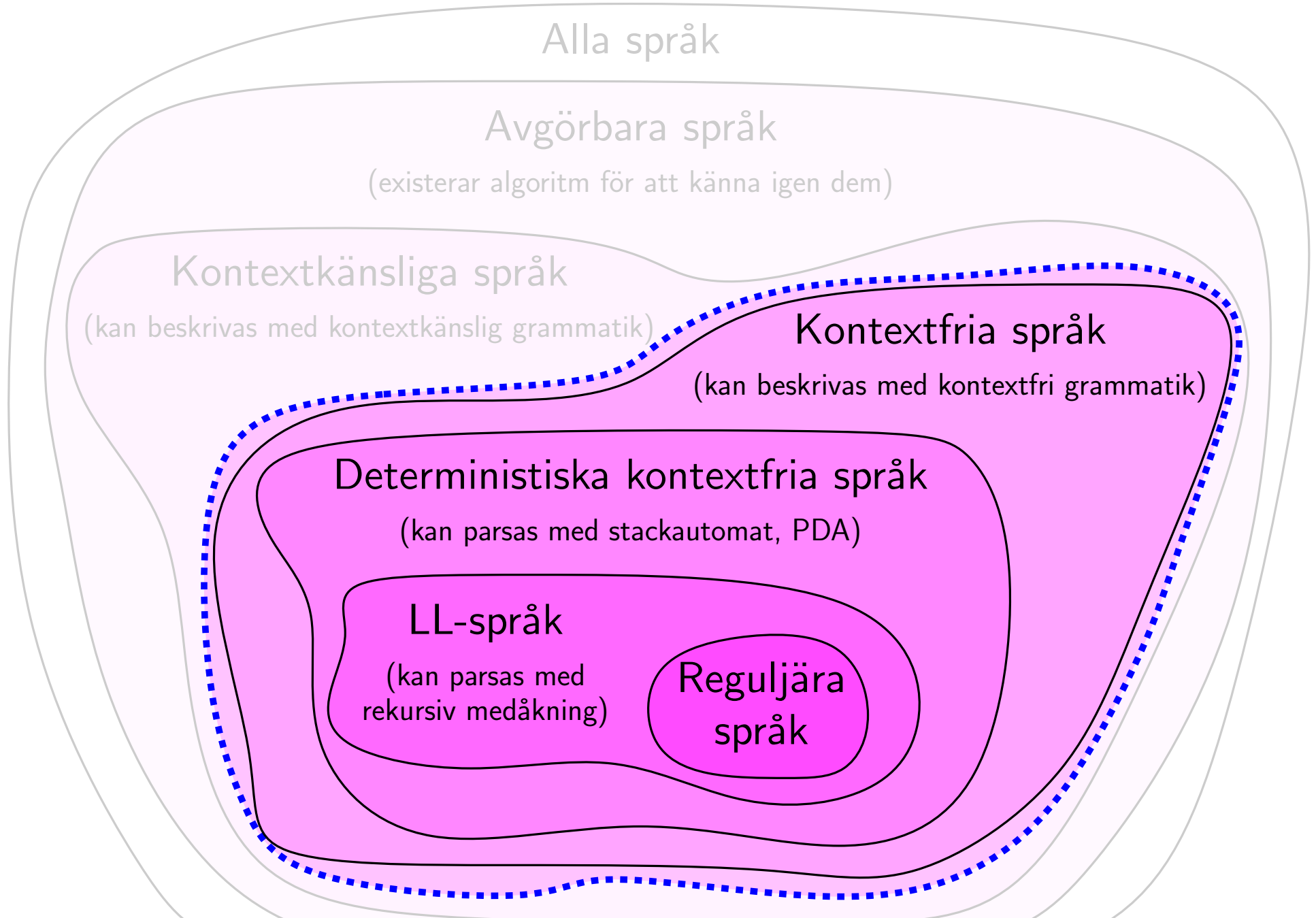
LL *Left-to-right Leftmost-derivation*.

- LL-parser: läser strängen vänster till höger och expanderar icke-slutsymboler från vänster till höger.
- LL-grammatik: grammatik som kan parsas av LL-parser

Perspektiv igen



Perspektiv igen



Vad vi har pratat om / Vad som ingår i kursavsnittet

Vad plugga på?

KS/Tenta kommer inte ta upp något som jag inte pratat om på någon av föreläsningarna.

KS samt tentorna från förra året är någorlunda representativa.

Under "Kursmaterial" på hemsidan finns länkar till föreläsningsvideos från en del liknande kurser, samt lite annat material, som man kan titta på.

Typiska KS-uppgifter

Grundläggande relationer mellan olika saker, vilka sätt att beskriva språk som är mer eller mindre kraftfulla, etc

Typiska KS-uppgifter

Grundläggande relationer mellan olika saker, vilka sätt att beskriva språk som är mer eller mindre kraftfulla, etc

Givet en DFA/PDA och en sträng förklara vad automaten gör med strängen som indata

Typiska KS-uppgifter

Grundläggande relationer mellan olika saker, vilka sätt att beskriva språk som är mer eller mindre kraftfulla, etc

Givet en DFA/PDA och en sträng förklara vad automaten gör med strängen som indata

Givet en grammatik och en sträng, konstruera ett syntaxträd / avgör om det finns flera syntaxträd / etc

Typiska KS-uppgifter

Grundläggande relationer mellan olika saker, vilka sätt att beskriva språk som är mer eller mindre kraftfulla, etc

Givet en DFA/PDA och en sträng förklara vad automaten gör med strängen som indata

Givet en grammatik och en sträng, konstruera ett syntaxträd / avgör om det finns flera syntaxträd / etc

Skriv ett reguljärt uttryck/kontextfri grammatik eller konstruera en DFA för ett givet språk

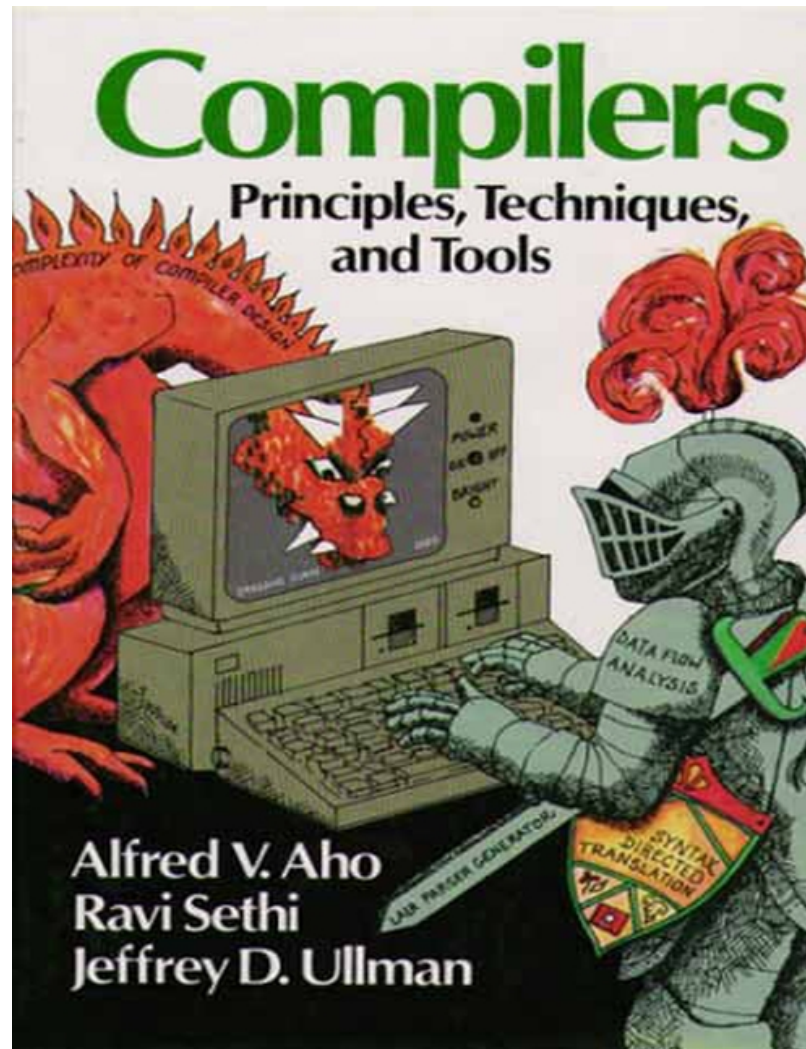
Fortsättningskurser

- DD1352 Algoritmer, datastrukturer och komplexitet
Mer om de högre nivåerna i språkhierarkin
- DD2372 Automater och språk
Mer på djupet om automater och formella språk.
- DD2488 Kompilatorkonstruktion
Skriv en hel kompilator!
- DD2418 Språkteknologi
Naturlig språkbehandling, inriktad mot text.
- ID2202 Kompilatorer och exekveringsmiljöer
Tekniker för implementation av programspråk.
(Ges i Kista)

För den som vill läsa mer

Aho, Sethi, Ullman:

Compilers: Principles, Techniques and Tools



Lycka till på kontrollskrivningen!