

# Mästarprov Syntax

Emilie Ruixin Cao, ercao@kth.se

24 mars 2025

## 1 Grammatik

Grammatiken  $G$  definieras:

$$\begin{array}{ll} S \rightarrow q & P \rightarrow S \\ S \rightarrow \mathbf{ay} P & P \rightarrow P * P \\ S \rightarrow \mathbf{ez} P & P \rightarrow ?P \\ S \rightarrow S * S & P \rightarrow \mathbf{xt} P \\ S \rightarrow ?S & P \rightarrow P \mathbf{uv} P \end{array}$$

a) Visa att grammatiken  $G$  är tvetydig genom att göra två härledningar för en sträng  $u$  i  $L$  som resulterar i olika syntaxträd. Du behöver inte visa syntaxträden.

Ett exempel på en sträng som har två olika härledningar som också leder till två olika syntaxträd är strängen  $?q*?q*?q$ .

**Första härledningen (högerassociativt):**

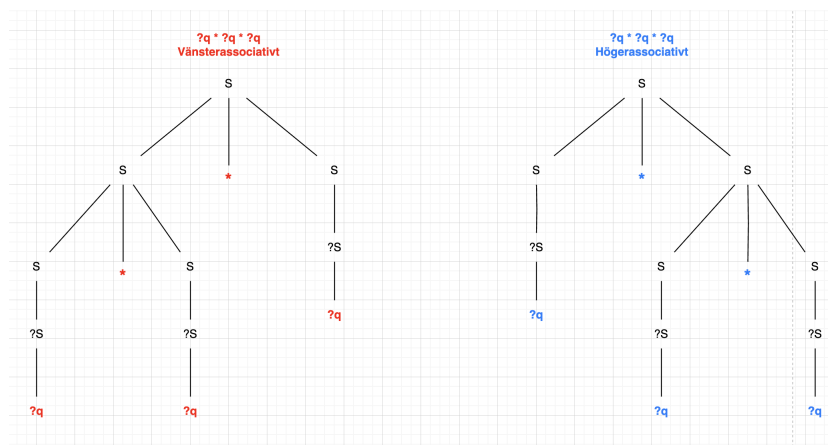
$$\begin{aligned} S &\rightarrow S * S \rightarrow S * S * S \rightarrow S * S * ?S \rightarrow S * S * ?q \\ &\rightarrow S * ?S * ?q \rightarrow S * ?q * ?q \rightarrow ?S * ?q * ?q \rightarrow ?q * ?q * ?q \end{aligned}$$

**Andra härledningen (vänsterassociativt):**

$$\begin{aligned} S &\rightarrow S * S \rightarrow S * S * S \rightarrow ?S * S * S \rightarrow ?q * S * S \\ &\rightarrow ?q * ?S * S \rightarrow ?q * ?q * S \rightarrow ?q * ?q * ?S \rightarrow ?q * ?q * ?q \end{aligned}$$

Dessa två härledningar resulterar i olika syntaxträd eftersom uttrycket  $S * S$  kan antingen byta ut det vänstra  $S$ :et mot ännu ett  $S * S$  så att det blir  $(S * S) * S$  eller så byts det högra  $S$ :et mot ännu ett  $S * S$  så att det blir  $S * (S * S)$ . Detta kommer att leda till två olika syntaxträd där den ena är en spegelbild av den andra. Anledningen till denna skillnad i syntaxträd beror på att regeln  $S * S$  kan tillämpas på två sätt. Antingen så expanderar man med det vänstra  $S$ :et först vilket leder till att uttrycket grupperas som  $(S * S) * S$  eller så expanderar man det högra  $S$ :et först så att uttrycket grupperas som  $S * (S * S)$ . När det vänstra  $S$ :et expanderas först

blir syntaxträdet vänsterassociativt, medan expansion av högra S:et först gör trädet högerassociativt. Därmed blir det två olika syntaxträd och grammatiken är tvetydig enligt boken sida 128. De två olika syntaxträden visas i bilden nedan.



b) Du ska nu göra en grammatik för ett språk som är en utökning av L med slutsymboler för vänsterparentes, “(”, och högerparentes, “)”.

Startsymbolen i denna grammatik är ' $\langle s \rangle$ '.

$$\begin{aligned}
 \langle s \rangle &::= '(< q > \langle s' \rangle) ' \\
 \langle s' \rangle &::= ' ' \mid '* ' \langle q > \langle s' \rangle \\
 \langle q \rangle &::= CAPS \mid 'ay ' \langle p \rangle \mid 'ez ' \langle p \rangle \mid '? ' \langle s \rangle \\
 \langle p \rangle &::= '(< t > \langle p' \rangle) ' \\
 \langle p' \rangle &::= ' ' \mid '* ' \langle t > \langle p' \rangle \mid 'uv ' \langle t > \langle p' \rangle \\
 \langle t \rangle &::= \langle i \rangle \mid '(< s > ' ) ' \\
 \langle i \rangle &::= '? ' \langle p \rangle \mid 'xt ' \langle p \rangle
 \end{aligned}$$

**CAPS:** [A-Z]+ i den definierade grammatiken. Alltså en eller flera kapitaliserade bokstäver.

Ovanstående grammatik är icke-tvetydig och kan parsas med en rekursiv medåkare med endast en lookahead eftersom de olika utfallen har en unik terminal som 'förutspår' vilken icke-terminal som ska köras närmast. Dessutom så är alla produktioner ej vänsterrekursiva då vänsterrekursion kan leda till oändliga loopar vid rekursion. Ex. kan  $S \rightarrow S * S$  leda till oändliga loopar då den första produktionen S pekar på är sig själv och därmed finns risken att den aldrig terminerar. Enligt drakboken s.68 så kan man eliminera vänsterrekursion genom följande teknik: Givet regeln  $A \rightarrow A\alpha$  så kan det skrivas om till två

regler  $A \rightarrow \beta R$  och  $R \rightarrow \alpha R \mid \epsilon$ .

**c) Motivera kort varför en hypotetisk implementation av en parser med rekursiv medåkning utgående från din grammatik i b) endast behöver inspektera en slutsymbol (token) framåt i en sträng för att avgöra vilken regel som ska användas för en given ickeslutsymbol.**

Alla utfall i min grammatik i b) kan parsas genom att inspektera en slutsymbol fram eftersom alla olika utfall har en unik slutsymbol som första token. Exempelvis om man börjar på  $\langle s \rangle$  och parsern sedan läser av  $\langle q \rangle$  först så kommer den kunna välja mellan fyra alternativ varav varje alternativ börjar på ett unikt ASCII-tecken. Alltså går det att avgöra nästa operation bara baserat på en slutsymbol framåt då ex. slutsymbolen  $?$  i produktionen för  $\langle q \rangle$  innebär att nästa operation är  $\langle s \rangle$  eller om nästa slutsymbol är  $a$  och då är nästa produktion  $\langle p \rangle$ . Sedan kommer den också evaluera  $\langle s' \rangle$ .  $\langle s' \rangle$  kan också välja mellan två olika alternativ med unika slutsymboler: den tomma strängen eller stjärnan  $(*)$ . Läser den av stjärna så kommer den evaluera  $\langle q \rangle$  igen och sedan  $\langle s' \rangle$  med rekursiva anrop tills den når den tomma strängen och terminerar. Produktionen i  $\langle p \rangle$  fungerar på samma sätt förutom att  $\langle p' \rangle$  kan evaluera tre olika fall: antingen börjar den en produktion med **uv**, eller en produktion med stjärnan  $*$ , eller så är det tomma strängen. Om  $\langle p' \rangle$  inte leder till tomma strängen så kommer det göras ett anrop på  $\langle t \rangle$  som antingen läser av en parentes och startar produktionen  $\langle s \rangle$  inom parentesen, annars kallar den på  $\langle i \rangle$  som då läser av om det är **xt** eller  $?$  som nästa symbol. Då kallar den på  $\langle p \rangle$  igen. I övrigt så är kallas  $\langle p' \rangle$  rekursivt i regeln  $\langle p' \rangle \rightarrow \text{EFTER anropet på } \langle t \rangle$ .

För att detta ska funka och inte vara tvetydigt så följer det att den inte kan vara vänsterrekursiv eftersom detta innebär att den kan fastna i en oändlig rekursion. Tekniken som används för att eliminera vänsterrekursion härleds på sådant sätt att om man har  $A \rightarrow A\alpha \mid \beta$  som är vänsterrekursiv så kan man skriva om den på formen  $A \rightarrow \beta R$ ,  $R \rightarrow \alpha R \mid \epsilon$ . Denna teknik förklaras i drakboken på sida 68. Enligt drakboken (s.223) så behöver tre krav uppfyllas för att det ska vara LL(1) – alltså en rekursiv medåkare med en lookahead.

Första kravet är att för ett uttryck  $A \rightarrow \alpha \mid \beta$  ska varken  $\alpha$  eller  $\beta$  kunna härleda en sträng som börjar på samma slutsymbol  $a$ . Detta stämmer i den konstruerade grammatiken i b) eftersom alla möjliga utfall leder till strängar som har en unik slutsymbol först. Exempelvis så leder  $\langle s' \rangle ::= '' \mid '*' < q > < s' \rangle$  till att den ena produktionen ger  $\epsilon$  som första slutsymbol och den andra produktionen ger att  $'*$  är nästa slutsymbol. Detsamma syns i produktionen för  $\langle p' \rangle$  där  $\epsilon$ ,  $*$ , och **uv** är de tre olika första slutsymbolerna.

Det andra kravet för uttrycket är att max en av produktionerna  $\alpha$  eller  $\beta$  får producera en tom sträng vilket också stämmer då ex.  $\langle s' \rangle$  kan härleda en tom sträng endast från ena ledet som den består av. Detsamma gäller för  $\langle p' \rangle$  också där det bara är en av produktionerna som ger tomma strängen. För  $\langle t \rangle$  så är det ingen produktion som ger tomma strängen vilket också uppfyller villkoret då det är MAX en produktion som kan leda till tomma strängen. Alltså

uppfylls även krav två i enlighet med drakbokens tre krav på en LL1 grammatik.

Den tredje regeln säger att om en grammatikregel har mer än en produktion och en av produktionerna ger den tomma strängen så måste de resterande produktionernas första slutsymbol vara disjunkt från huvudets följande symbol. Alltså för uttrycket  $A \rightarrow \alpha \mid \beta$  där  $\beta \rightarrow \epsilon$  så måste  $\text{FIRST}(\alpha)$  vara disjunkt från  $\text{FOLLOW}(A)$ . Om vi kollar på exempelvis  $\langle s' \rangle$  där den första produktionen ger  $\epsilon$  så är  $\text{FIRST}(*' \langle q \rangle \langle s' \rangle) = \{*\}$ .  $\text{FOLLOW}(\langle s' \rangle)$  blir symbolen som kommer efter produktionen  $\langle s' \rangle$ . Produktionen  $\langle s' \rangle$  förekommer endast i regeln  $\langle s \rangle ::= '(' \langle q \rangle \langle s' \rangle ')'$ . Detta innebär att  $\text{FOLLOW}(\langle s' \rangle)$  alltid kommer att vara en högerparentes  $)$  då  $\langle s' \rangle$  är sluten i stängda paranteser. Då  $\langle s \rangle$  också är startsymbolen lär  $\text{FOLLOW}(\langle s' \rangle)$  också innehålla  $\$$  som markerar slutet på input. Vi kan se att  $\text{FOLLOW}(\langle s' \rangle)$  inte innehåller  $*$  och är skilt från  $\text{FIRST}(*' \langle q \rangle \langle s' \rangle)$ . Alltså är de disjunkta och kommer vara kompatibla med en LL(1) parser. Samma argument kan dras för  $\langle p' \rangle$  (där epsilon kan härledas i en av produktionerna) som har  $\text{FIRST}(*' \langle t \rangle \langle p' \rangle) = *$  och  $\text{FIRST}('uv' \langle t \rangle \langle p' \rangle) = 'u'$ .  $\text{FOLLOW}(\langle p' \rangle)$  förekommer endast i regeln  $\langle p \rangle ::= '(' \langle t \rangle \langle p' \rangle ')'$  och där blir  $\text{FOLLOW}(\langle p' \rangle) = )$  då denna också är sluten inom parentes. Vi kan se att dessa två mängder är disjunkta. Alltså är grammatiken tvetydig och kan parsas med en rekursiv medåkare som har en lookahead. Enligt drakboken så kan endast grammatiker som är **entydiga** parsas med en LL(1) parser och därmed följer även att grammatiken är entydig.

**d) Ge ett exempel på en sträng  $v$  i ditt språk från b), och en sträng  $w$  från  $L$ , sådana att om alla vänsterparenteser och högerparenteser tas bort från  $v$  så erhålls  $w$ .**

Ett exempel på en sträng  $v$  i mitt språk från b) är strängen  $'(ay((q)))'$ , om man tar bort alla vänster- och högerparenteser från  $v$  så erhålls strängen  $w$   $'ayq'$  från språket  $L$ . Härledning för båda språken ges nedan:

I språket  $L$  blir strängen  $w = 'ayq'$ :

$$S \rightarrow ayP \rightarrow ayS \rightarrow ayq$$

Vi kan alltså se att strängen  $w = 'ayq'$  går att härleda för språket  $L$ . Detta motsvarar strängen  $v = '(ay((q)))'$  utan paranteser i min grammatik från b). För att härleda att strängen  $v$  accepteras i min grammatik görs följande härledning:

$$\begin{aligned} \langle s \rangle &\rightarrow (\langle q \rangle \langle s' \rangle) \rightarrow (ay \langle p \rangle \langle s' \rangle) \rightarrow (ay(\langle t \rangle \langle p' \rangle) \langle s' \rangle) \\ &\rightarrow (ay((\langle s \rangle) \langle p' \rangle) \langle s' \rangle) \rightarrow (ay((\langle q \rangle \langle s' \rangle) \langle p' \rangle) \langle s' \rangle) \\ &\rightarrow (ay((q \langle s' \rangle) \langle p' \rangle) \langle s' \rangle) \rightarrow (ay((q\epsilon) \langle p' \rangle) \langle s' \rangle) \\ &\rightarrow (ay((q)\epsilon) \langle s' \rangle) \rightarrow (ay((q))\epsilon) \rightarrow (ay((q))) \end{aligned}$$

Här ser vi att man kan härleda strängen  $v = 'ay((q))'$  med hjälp av grammatiken från uppgift b). Tar vi bort alla vänster- och högerparanteser så återstår endast strängen 'ayq' vilket motsvarar strängen  $w$  från språket  $L$ .

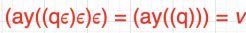
**e) Ge en härledning i din grammatik från b) som visar att strängen  $v$  från d) verkligen tillhör språket som din grammatik definierar. För att undvika långa upprepningar i härledningen kan du använda tydligt definierade förkortningar.**

En härledning av strängen  $v$  visas i motiveringen för uppgift d), men nedan lyder denna härledning igen:

$$\begin{aligned} < s > \rightarrow (< q > < s' >) \rightarrow (ay < p > < s' >) \rightarrow (ay(< t > < p' >) < s' >) \\ &\rightarrow (ay((< s >) < p' >) < s' >) \rightarrow (ay((< q > < s' >) < p' >) < s' >) \\ &\rightarrow (ay((q < s' >) < p' >) < s' >) \rightarrow (ay((q\epsilon) < p' >) < s' >) \\ &\rightarrow (ay((q)\epsilon) < s' >) \rightarrow (ay((q)\epsilon) \rightarrow (ay((q))) \end{aligned}$$

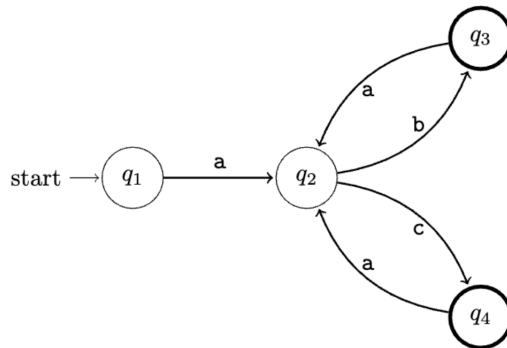
**f) Ge ett syntaxträd för strängen  $v$  i d) och motivera kort varför det bara finns ett enda syntaxträd för  $v$  i din grammatik. Det räcker att du använder namnen på dina tokens i trädet, dvs du behöver inte ange deras representation som sekvenser i ASCII.**

Ett syntaxträd för strängen  $v = 'ay((q))'$  ges av följande:



## 2 Automater och Reguljära Uttryck

$$\begin{aligned} R(x, L_1) &= \{w \mid xw \in L_1\} \text{ f\"or } x \in \Sigma \\ S(L_1) &= L_1^* \\ K(L_1, L_2) &= \{uv \mid u \in L_1 \text{ och } v \in L_2\} \\ P(L_1, L_2) &= \{u \mid u \in L_1 \text{ eller } u \in L_2\} \\ Z(L_1, L_2) &= \{uv \mid \mathbf{cbu} \in L_1 \text{ och } v \in L_2\} \\ B &= \{\mathbf{ab}\} \\ C &= \{\mathbf{ac}\} \end{aligned}$$



Låt  $M$  vara språket som känns igen av DFA:n över  $\Sigma$  med följande grafiska representation.

**a) Uttryck språket  $Z(L, M)$  endast i termer av operationerna  $R$ ,  $S$ ,  $K$ ,  $P$ , och språken  $B$ ,  $C$  och  $L$ . Med andra ord, ge en definition av  $Z(L, M)$  med hjälp av  $B$ ,  $C$  och  $L$  som inte använder mängdnotation eller andra språkoperationer än  $R$ ,  $S$ ,  $K$  och  $P$ .**

Språket  $M$  kan beskrivas på följande vis med hjälp av de definierade operatorerna  $R$ ,  $S$ ,  $K$ , och  $P$  samt språken  $B$ ,  $C$ , och  $L$ :

$$M = K(P(B, C), S(P(B, C)))$$

Den grafiska representationen av DFA:n visar att språket  $M$  accepterar strängar som börjar på antingen 'ab' eller 'ac'. Utöver detta så kan strängarna fortsätta med noll eller flera iterationer av antingen 'ab' eller 'ac'. Regeln  $K$  tar in två språk som parametrar och returnerar ett språk som består av strängarna  $uv$  där  $u$  tillhör det första språket som input och  $v$  tillhör det andra språket som input. I detta fall är det  $P(B, C)$  som är det första språket som är input vilket innebär att den första delsträngen  $u$  ska tillhöra språket som alla strängar som både finns med i språket  $B$  och språket  $C$ . Detta innebär att  $u$  kommer att vara antingen 'ab' (språket  $B$ ) eller 'ac' (språket  $C$ ). Det andra språket som är input är  $S(P(B, C))$  vilket innebär att delsträngen  $v$  kommer att bestå av noll eller flera iterationer av antingen 'ab' eller 'ac' då detta är språket som returneras från  $P(B, C)$ . Kleenestjärnan i operatoren  $S$  innebär att alla strängar i inputspråket kan konkateneras noll eller flera gånger. Om man då sedan applicerar operatoren  $K$  på dessa på inputspråk så får man alltså en konkatenering mellan 'ab' eller 'ac' och  $ab^*$  eller  $ac^*$ .

Med detta i åtanke kan man alltså uttrycka  $Z(L, M)$  på följande vis:

$$Z(L, M) = K(R(b, R(c, L)), K(P(B, C), S(P(B, C))))$$

Det vi gör är ännu en konkatenering med  $K$ . Denna gång är det första inputspråket  $R(b, R(c, L))$  och det andra inputspråket är  $M$ , som definierats

tidigare. Eftersom  $K(L_1, L_2)$  definieras som mängden av alla konkateneringar där den första delsträngen  $u$  tillhör  $L_1$  och den andra delsträngen  $v$  tillhör  $L_2$ , betyder detta att den andra delsträngen  $v$  i konkateneringen bör tillhöra språket  $M$ . Det första inputspråket i  $K$ , alltså  $R(b, R(c, L))$ , betecknar alla strängar i språket  $L$  som börjar på bokstäverna 'cb'. Detta inputspråk kan delas upp i en inre och yttre operator:

**1. Inre operatör:**  $R(c, L)$

$R(c, L)$  innehåller alla strängar  $u$  sådana att OM vi lägger till ett 'c' framför  $u$ , så tillhör  $cu$  språket  $L$ . Detta betyder att  $R(c, L)$  innehåller exakt de strängar i  $L$  där vi kan anta att ett 'c' redan har lästs bort.

**2. Yttre operator:**  $R(b, R(c, L))$

$R(b, R(c, L))$  innehåller alla strängar  $u$  sådana att OM vi lägger till ett 'b' framför  $u$ , så tillhör  $bu$  mängden  $R(c, L)$ . Eftersom vi redan vet att  $R(c, L)$  består av alla  $u$  där  $cu \in L$ , så betyder detta att  $bu$  måste vara en sådan sträng att  $c bu \in L$ .

När man då använder operatörn  $K$  för konkatenering så kan man alltså applicera detta som det första inputspråket för att definiera att den första delsträngen måste tillhöra språket  $L$  samt börja med bokstäverna 'cb'. Det andra inputspråket är våran definition av  $M$  som vi kan se ovan.

**b) Antag att språket  $L$  känns igen av DFA:n  $(Q, \Sigma, \delta, q, F)$ . Ge en matematisk definition av en DFA som känner igen språket  $\{u \mid cbu \in L\}$ . Med matematisk definition av en DFA menas en definition av en tupel  $(Q', \Sigma, \delta', q', F')$  och inte en grafisk representation.**

En matematisk definition av en DFA som känner igen språket är:

$$L' = (Q, \Sigma, \delta, s, F) \\ s = \delta(\delta(q, c), b)$$

Denna DFA känner igen språket  $\{u \mid cbu \in L\}$  eftersom starttillståndet (den fjärde parametern i tupeln) ska börja i ett tillstånd där den redan har 'gått igenom' bokstäverna 'cb'. För att kunna säkerställa detta så kan man applicera det boken kallar för *transition function*  $\delta$  (s.48) två gånger. Funktionen  $\delta$  har två inparametrar: ett tillstånd och en symbol. Alltså är tillståndet  $\delta(a, q)$  det tillståndet  $p$  man når när man övergår från tillstånd  $q$  till  $p$  med en kant markerad med 'a' emellan dem. Detta är relevant till denna matematiska definition av  $L'$  som känner igen språket  $\{u \mid cbu \in L\}$  eftersom vi måste säkerställa att starttillståndet innefattar att strängen har 'cb' som startsymboler. Först applicerar vi  $\delta$  på symbolen  $c$  och tillståndet  $q$  som är starttillståndet för språket  $L$ . Genom denna applikation av tillståndsfunktionen får vi ett nytt tillstånd  $q_1$  i vilket vi har 'ackumulerat' ett 'c' i strängen redan. Från detta



tillstånd  $q_1 = \delta(q, c)$  så applicerar vi  $\delta$  igen fast med bokstaven 'b' och på det nya tillståndet  $q_1$ . Alltså  $\delta(q_1, b)$  vilket ger oss det nästa tillståndet vi når om vi ackumulerar ett 'b' i vår sträng. Då  $\delta(\delta(q, c), b)$  returnerar det tillstånd vi är på då vi redan ackumulerat 'cb' i starten av strängen så motsvarar det alltså starttillståndet till  $L'$ . Sedan är det samma mängd tillstånd  $Q$  som i  $L$ :s definition, samma alfabet  $\Sigma = \{a, b, c\}$  som från  $L$ , samma funktion  $\delta$  som från  $L$ , och samma accepterande tillstånd  $F$  som från  $L$ . Detta då det nya språket är uppbyggt på samma sätt som språket  $L$ , det viktigaste att skilja på är starttillståndet. I detta fall är starttillståndet det tillstånd där vi redan har checkat att de första två symbolerna ska vara 'cb'. Vi kan kalla detta nya starttillstånd för  $s = \delta(\delta(q, c), b)$ .

**c) Antag att din DFA från b) är sådan att alla tillstånd i  $F'$  har en övergång för a till ett loopande tillstånd  $q_0 \in Q' \setminus F'$  (dvs  $q_0$  är ett icke-accepterande tillstånd som har övergångar till sig själv för alla  $x \in \Sigma$ ). Förklara kortfattat hur du med hjälp av denna DFA kan bygga en DFA som känner igen  $Z(L, M)$ .**

Eftersom  $Z(L, M)$  består av en konkatenerad sträng med två delsträngar från olika språk är det lämpligt att koppla ihop DFA:n för  $L'$  och DFA:n för  $M$ . Detta kan göras med hjälp av något som kallas *product* automaton enligt boken s.43. Det som händer är att varje tillstånd representeras av en tupel  $(p, q)$  där  $p$  är ett tillstånd från  $L'$ :s DFA och  $q$  är ett tillstånd från  $M$ :s DFA. Starttillståndet kan då exempelvis skrivas som en tupel  $(q_{l_{start}}, q_{m_{start}})$  istället för bara  $q_{l_{start}}$ . För samtliga tillstånd  $q$  som tillhör DFA:n till  $L'$  så kan man beteckna tillstånden med tupeln  $(q, q_{m_{start}})$ . Alltså att medan vi behandlar strängen  $cbu$  som tillhör språket  $L'$  så kommer samtidigt  $M$ :s DFA alltid konstant förbli i ett starttillstånd  $q_{m_{start}}$ . När vi sedan når ett accepterande tillstånd i DFA:n för  $L'$  så kommer det finnas en övergång markerad 'a' till tillståndet  $q_0$  i DFA:n för  $L'$ , men samtidigt så når den även tillståndet  $q_2$  i DFA:n för  $M$  (se figuren på  $M$  för tydliggörande). Då får vi alltså tupeln  $(q_0, q_2)$  när vi initialt når tillståndet  $q_0$  i DFA för  $L'$ . För alla resterande tillstånd i DFA:n för  $M$  så kommer tillståndet i DFA:n för  $L'$  vara  $q_0$ . Så exempelvis när DFA:n för  $M$  vandrar från  $q_2$  till  $q_3$  så kommer tupeln gå från  $(q_0, q_2)$  till  $(q_0, q_3)$ . På samma sätt kan man resonera för hela DFA:n för  $M$ . Samma övergångar i DFA:n för  $M$  som finns på figuren i uppgiftsbeskrivningen tillhandahålls i den nya DFA:n för  $Z(L, M)$  men varje tillstånd definieras som en tupel där  $q_0$  är konstant för DFA:n för  $L'$  medan tillståndet i  $M$ :s DFA skiftar. På så sätt kan man bygga upp en DFA som känner igen  $Z(L, M)$ , men då måste man som sagt använda sig av det så kallade *product automaton*. Den informationen tillhandahålls i boken s.43-44.

### 3 Språkklasser

Låt  $\Sigma_1$  och  $\Sigma_2$  vara alfabet, och låt  $f$  vara en funktion från  $\Sigma_1^*$  till  $\Sigma_2^*$  sådan att  $f(uv) = f(u)f(v)$  för alla  $u \in \Sigma_1^*$  och  $v \in \Sigma_1^*$ . Antag

att  $L$  är ett reguljärt språk över  $\Sigma_2$ . Antag vidare att  $L'$  är ett språk över  $\Sigma_1$  som definieras av en (kontextfri) grammatik. Vilken är den minsta språkklass som tagits upp i kursen som språket  $\{w \in \Sigma_1^* \mid f(w) \in L \text{ eller } w \in L'\}$  då med säkerhet tillhör? Motivera ditt svar noggrant.

Språket  $S = \{w \in \Sigma_1^* \mid f(w) \in L \text{ eller } w \in L'\}$  tillhör med säkerhet språkklassen 'kontextfria språk'. Detta beror på att funktionen  $f$  är en så kallad homomorfism. Definitionen ges i kursboken s.141.

Definitionen av en homomorfism ges av boken som: Funktionen  $h$  är en homomorfism på alfabetet  $\Sigma$  där  $w = a_1a_2a_3\dots a_n$  är en sträng med symboler från  $\Sigma$  om  $h(w) = h(a_1)h(a_2)\dots h(a_n)$ . Vi kan därmed se att funktionen  $f$  är en homomorfism då vi kan se att för  $w = uv$  så  $f(w) = f(uv) = f(u)f(v)$ . Enligt bokens definition så är reguljära språk slutna under en homomorfism ("... the regular languages are closed under homomorphism", s.141). Mer formellt står det i boken: "If  $L$  is a regular language over alphabet  $\Sigma$ , and  $h$  is a homomorphism on  $\Sigma$ , then  $h(L)$  is also regular."

Detta innebär att i vårt fall, där  $L$  är givet som ett reguljärt språk i uppgiften, så innebär det att  $f(L)$  också är ett reguljärt språk. Om man sedan kollar vidare i boken så står det att inversen av en homomorfism också bevarar ett reguljärt språk. Enligt boken definieras inversen till en homomorfism  $h$  på följande vis: ".../  $h^{-1}(L)$  /.../ is the set of strings  $w$  in  $\Sigma^*$  such that  $h(w)$  is in  $L$ ." I vårt språk  $S$  så säger vi att antingen så  $f(w) \in L$  eller  $w \in L'$ . Uttrycket  $f(w) \in L$  kan kopplas bokens definition av invers homomorfism så att  $f^{-1}(L)$  är mängden strängar  $w$  i  $\Sigma_2^*$  så att  $f(w)$  tillhör språket  $L$ . Alltså definierar det första uttrycket  $f(w) \in L$  ett reguljärt språk i enlighet med boken.

Om man sedan evaluerar det andra påståendet i eller-uttrycket,  $w \in L'$ , så kan man genast konstatera att detta uttryck definierar ett kontextfritt språk eftersom det är givet i uppgiften att  $L'$  är definierat av en kontextfri grammatik. Enligt föreläsningen 'syntax 4' på slide 12 så står det att kontextfria språk är en språkklass som består av alla språk som kan beskrivas med en kontextfri grammatik vilket precis är vad  $L'$  är. Därmed kan man konstatera att alla strängar så att  $w \in L'$  betecknar ett språk av klassen kontextfritt språk.

Eftersom språket  $S$  definieras av ett 'eller'-uttryck så är språket  $S$  en union mellan  $f(w) \in L$  och  $w \in L'$ , då kan man se det som en union mellan ett reguljärt språk och ett kontextfritt språk. Enligt Wenn-diagrammet från föreläsningen 'syntax 4' slide 12 så är reguljära språk en delmängd av kontextfria språk. Detta innebär alltså att unionen mellan språken med säkerhet kommer att vara ett kontextfritt språk.

