

# Chomsky-Schützenberger Representation Theorem

Moritz Roos and Tobias Nipkow

October 13, 2025

## Abstract

The Chomsky-Schützenberger Representation Theorem says that any context-free language is the homomorphic image of the intersection of a regular language and a Dyck language.

## Contents

<b>1</b>	<b>Overview of the Proof</b>	<b>2</b>
<b>2</b>	<b>Production Transformation and Homomorphisms</b>	<b>4</b>
2.1	Brackets . . . . .	4
2.2	Transformation . . . . .	5
2.3	Homomorphisms . . . . .	6
<b>3</b>	<b>The Regular Language</b>	<b>7</b>
3.1	$P1$ . . . . .	8
3.2	$P2$ . . . . .	9
3.3	$P3$ . . . . .	9
3.4	$P4$ . . . . .	10
3.5	$P5$ . . . . .	11
3.6	$P7$ and $P8$ . . . . .	11
3.7	$Reg$ and $Reg\_sym$ . . . . .	12
<b>4</b>	<b>Showing Regularity</b>	<b>13</b>
4.1	An automaton for $\{xs. \textit{ successively } Q \ xs \wedge xs \in \textit{brackets } P\}$ .	14
4.2	Regularity of $P2$ , $P3$ and $P4$ . . . . .	16
4.3	An automaton for $P1$ . . . . .	16
4.4	An automaton for $P5$ . . . . .	18
<b>5</b>	<b>Definitions of <math>L</math>, <math>\Gamma</math>, <math>P'</math>, <math>L'</math></b>	<b>20</b>
<b>6</b>	<b>Lemmas for <math>P' \vdash A \Rightarrow^* x \longleftrightarrow x \in R_A \cap Dyck\_lang \ \Gamma</math></b>	<b>20</b>
<b>7</b>	<b>Showing <math>h(L') = L</math></b>	<b>21</b>

```

theory Chomsky_Schuetzenberger
imports
  Context_Free_Grammar.Parse_Tree
  Context_Free_Grammar.Chomsky_Normal_Form
  Finite_Automata_HF.Finite_Automata_HF
  Dyck_Language_Syms
begin

```

This theory proves the Chomsky-Schützenberger representation theorem [1]. We closely follow Kozen [2] for the proof. The theorem states that every context-free language  $L$  can be written as  $h(R \cap \text{Dyck\_lang } \Gamma)$ , for a suitable alphabet  $\Gamma$ , a regular language  $R$  and a word-homomorphism  $h$ .

The Dyck language over a set  $\Gamma$  (also called it's bracket language) is defined as follows: The symbols of  $\Gamma$  are paired with  $[$  and  $]$ , as in  $[_g$  and  $]_g$  for  $g \in \Gamma$ . The Dyck language over  $\Gamma$  is the language of correctly bracketed words. The construction of the Dyck language is found in theory *Chomsky\_Schuetzenberger.Dyck\_Language\_Syms*.

## 1 Overview of the Proof

A rough proof of Chomsky-Schützenberger is as follows: Take some context-free grammar for  $L$  with productions  $P$ . Wlog assume it is in Chomsky Normal Form. Now define a new language  $L'$  with productions  $P'$  in the following way from  $P$ :

If  $\pi = A \rightarrow BC$  let  $\pi' = A \rightarrow [^1_\pi B ]^1_p [^2_\pi C ]^2_p$ , if  $\pi = A \rightarrow a$  let  $\pi' = A \rightarrow [^1_\pi ]^1_p [^2_\pi ]^2_p$ , where the brackets are viewed as terminals and the old variables  $A, B, C$  are again viewed as nonterminals. This transformation is implemented by the function *transform\_prod* below. Note brackets are now adorned with superscripts 1 and 2 to distinguish the first and second occurrences easily. That is, we work with symbols that are triples of type  $\{[,]\} \times \text{old\_prod\_type} \times \{1,2\}$ .

This bracketing encodes the parse tree of any old word. The old word is easily recovered by the homomorphism which sends  $[^1_\pi$  to  $a$  if  $\pi = A \rightarrow a$ , and sends every other bracket to  $\varepsilon$ . Thus we have  $h(L') = L$  by essentially exchanging  $\pi$  for  $\pi'$  and the other way round in the derivation. The direction  $\supseteq$  is done in *transfer\_parse\_tree*, the direction  $\subseteq$  is done directly in the proof of the main theorem.

Then all that remains to show is, that  $L'$  is of the form  $R \cap \text{Dyck\_lang } \Gamma$  (for  $\Gamma := P \times \{1, 2\}$ ) and the regularity of  $R$ .

For this,  $R := R_S$  is defined via an intersection of 5 following regular languages. Each of these is defined via a property on words  $x$ :

*P1*  $x$ : after a  $]^1_p$  there always immediately follows a  $[^2_p$  in  $x$ . This especially means, that  $]^1_p$  cannot be the end of the string.

*successively P2*  $x$ : a  $]^2_\pi$  is never directly followed by some  $[$  in  $x$ .

*successively P3*  $x$ : each  $[^1_{A \rightarrow BC}$  is directly followed by  $[^1_{B \rightarrow \_}$  in  $x$  (last letter isn't checked).

*successively P4*  $x$ : each  $[^1_{A \rightarrow a}$  is directly followed by  $]^1_{A \rightarrow a}$  in  $x$  and each  $[^2_{A \rightarrow a}$  is directly followed by  $]^2_{A \rightarrow a}$  in  $x$  (last letter isn't checked).

*P5*  $A$   $x$ : there exists some  $y$  such that the word begins with  $[^1_{A \rightarrow y}$ .

One then shows the key theorem  $P' \vdash A \rightarrow^* w \iff w \in R_A \cap \text{Dyck\_lang } \Gamma$ :

The  $\rightarrow$ -direction (see lemma  $P' \text{\_imp\_Reg}$ ) is easily checked, by checking that every condition holds during all derivation steps already. For this one needs a version of  $R$  (and all the conditions) which ignores any Terminals that might still exist in such a derivation step. Since this version operates on symbols (a different type) it needs a fully new definition. Since these new versions allow more flexibility on the words, it turns out that the original 5 conditions aren't enough anymore to fully constrain to the target language. Thus we add two additional constraints *successively P7* and *successively P8* on the symbol-version of  $R_A$  that vanish when we ultimately restricts back to words consisting only of terminal symbols. With these the induction goes through:

(*successively P7\_sym*)  $x$ : each  $Nt$   $Y$  is directly preceded by some  $Tm$   $[^1_{A \rightarrow YC}$  or some  $Tm$   $[^2_{A \rightarrow BY}$  in  $x$ ;

(*successively P8\_sym*)  $x$ : each  $Nt$   $Y$  is directly followed by some  $]^1_{A \rightarrow YC}$  or some  $]^2_{A \rightarrow BY}$  in  $x$ .

The  $\leftarrow$ -direction (see lemma  $\text{Reg\_and\_dyck\_imp\_P'}$ ) is more work. This time we stick with fully terminal words, so we work with the standard version of  $R_A$ : Proceed by induction on the length of  $w$  generalized over  $A$ . For this, let  $x \in R_A \cap \text{Dyck\_lang } \Gamma$ , thus we have the properties *P1*  $x$ , *successively Pi*  $x$  for  $i \in \{2,3,4,7,8\}$  and *P5*  $A$   $x$  available. From *P5*  $A$   $x$  we have that there exists  $\pi \in P$  s.t.  $\text{fst } \pi = A$  and  $x$  begins with  $[^1_\pi$ . Since  $x \in \text{Dyck\_lang } \Gamma$  it is balanced, so it must be of the form  $x = [^1_\pi y ]^1_\pi r1$  for some balanced  $y$ . From *P1*  $x$  it must then be of the form  $x = [^1_\pi y ]^1_\pi [^2_\pi r1'$ . Since  $x$  is balanced it must then be of the form  $x = [^1_\pi y ]^1_\pi [^2_\pi z ]^2_\pi r2$  for some balanced  $z$ . Then  $r2$  must also be balanced. If  $r2$  was not empty it would begin with an opening bracket, but *P2*  $x$  makes this impossible - so  $r2 = []$  and as such  $x = [^1_\pi y ]^1_\pi [^2_\pi z ]^2_\pi$ . Since our grammar is in CNF, we can consider the following case distinction on  $\pi$ :

Case 1:  $\pi = A \rightarrow BC$ . Since  $y, z$  are balanced substrings of  $x$  one easily checks  $Pi\ y$  and  $Pi\ z$  for  $i \in \{1, 2, 3, 4\}$ . From  $P3\ x$  (and  $\pi = A \rightarrow BC$ ) we further obtain  $P5\ B\ y$  and  $P5\ C\ z$ . So  $y \in R_B \cap Dyck\_lang\ \Gamma$  and  $z \in R_C \cap Dyck\_lang\ \Gamma$ . From the induction hypothesis we thus obtain  $P' \vdash B \rightarrow^* y$  and  $P' \vdash C \rightarrow^* z$ . Since  $\pi = A \rightarrow BC$  we then have  $A \rightarrow^1_{\pi'} [^1_{\pi} B ]^1_{\pi} [^2_{\pi} C ]^2_{\pi} \rightarrow^* [^1_{\pi} y ]^1_{\pi} [^2_{\pi} z ]^2_{\pi} = x$  as required.

Case 2:  $\pi = A \rightarrow a$ . Suppose we didn't have  $y = []$ . Then from  $P4\ x$  (and  $\pi = A \rightarrow a$ ) we would have  $y = ]^1_{\pi}$ . But since  $y$  is balanced it needs to begin with an opening bracket, contradiction. So it must be that  $y = []$ . By the same argument we also have that  $z = []$ . So really  $x = [^1_{\pi} ]^1_{\pi} [^2_{\pi} ]^2_{\pi}$  and of course from  $\pi = A \rightarrow a$  it holds  $A \rightarrow^1_{\pi'} [^1_{\pi} ]^1_{\pi} [^2_{\pi} ]^2_{\pi} = x$  as required.

From the key theorem we obtain (by setting  $A := S$ ) that  $L' = R_S \cap Dyck\_lang\ \Gamma$  as wanted.

Only regularity remains to be shown. For this we use that  $R_S \cap Dyck\_lang\ \Gamma = (R_S \cap brackets\ \Gamma) \cap Dyck\_lang\ \Gamma$ , where  $brackets\ \Gamma (\supseteq Dyck\_lang\ \Gamma)$  is the set of words which only consist of brackets over  $\Gamma$ . Actually, what we defined as  $R_S$ , isn't regular, only  $(R_S \cap brackets\ \Gamma)$  is. The intersection restricts to a finite amount of possible brackets, that are used in states for finite automaton for the 5 languages that  $R_S$  is the intersection of.

Throughout most of the proof below, we implicitly or explicitly assume that the grammar is in CNF. This is lifted only at the very end.

## 2 Production Transformation and Homomorphisms

A fixed finite set of productions  $P$ , used later on:

```

locale locale_P =
fixes P :: ('n,'t) Prods
assumes finiteP: ‹finite P›

```

### 2.1 Brackets

A type with 2 elements, for creating 2 copies as needed in the proof:

```

datatype version = One | Two

```

```

type_synonym ('n,'t) bracket3 = (('n, 't) prod × version) bracket

```

```

abbreviation open_bracket1 :: ('n, 't) prod ⇒ ('n,'t) bracket3 ([^1_ [1000])
where

```

```

    [^1_p ≡ (Open (p, One))

```

```

abbreviation close_bracket1 :: ('n,'t) prod ⇒ ('n,'t) bracket3 ([^1_ [1000]) where

```

$]^1_p \equiv (\text{Close } (p, \text{One}))$

**abbreviation**  $\text{open\_bracket2} :: ('n, 't) \text{ prod} \Rightarrow ('n, 't) \text{ bracket3 } ([^2\_ [1000])$  **where**  
 $[^2_p \equiv (\text{Open } (p, \text{Two}))$

**abbreviation**  $\text{close\_bracket2} :: ('n, 't) \text{ prod} \Rightarrow ('n, 't) \text{ bracket3 } (]^2\_ [1000])$  **where**  
 $]^2_p \equiv (\text{Close } (p, \text{Two}))$

Version for p = (A, w) (multiple letters) with bsub and esub:

**abbreviation**  $\text{open\_bracket1}' :: ('n, 't) \text{ prod} \Rightarrow ('n, 't) \text{ bracket3 } ([^1\_)$  **where**  
 $[^1_p \equiv (\text{Open } (p, \text{One}))$

**abbreviation**  $\text{close\_bracket1}' :: ('n, 't) \text{ prod} \Rightarrow ('n, 't) \text{ bracket3 } (]^1\_)$  **where**  
 $]^1_p \equiv (\text{Close } (p, \text{One}))$

**abbreviation**  $\text{open\_bracket2}' :: ('n, 't) \text{ prod} \Rightarrow ('n, 't) \text{ bracket3 } ([^2\_)$  **where**  
 $[^2_p \equiv (\text{Open } (p, \text{Two}))$

**abbreviation**  $\text{close\_bracket2}' :: ('n, 't) \text{ prod} \Rightarrow ('n, 't) \text{ bracket3 } (]^2\_)$  **where**  
 $]^2_p \equiv (\text{Close } (p, \text{Two}))$

Nice LaTeX rendering:

**notation** (*latex output*)  $\text{open\_bracket1 } ([^1\_)$   
**notation** (*latex output*)  $\text{open\_bracket1}' ([^1\_)$   
**notation** (*latex output*)  $\text{open\_bracket2 } ([^2\_)$   
**notation** (*latex output*)  $\text{open\_bracket2}' ([^2\_)$   
**notation** (*latex output*)  $\text{close\_bracket1 } (]^1\_)$   
**notation** (*latex output*)  $\text{close\_bracket1}' (]^1\_)$   
**notation** (*latex output*)  $\text{close\_bracket2 } (]^2\_)$   
**notation** (*latex output*)  $\text{close\_bracket2}' (]^2\_)$

## 2.2 Transformation

**abbreviation**  $\text{wrap1} :: \langle 'n \Rightarrow 't \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ syms} \rangle$  **where**

$\langle \text{wrap1 } A \ a \equiv$   
 $[ \ Tm \ [^1(A, [Tm \ a]),$   
 $\quad Tm \ ]^1(A, [Tm \ a]),$   
 $\quad Tm \ [^2(A, [Tm \ a]),$   
 $\quad Tm \ ]^2(A, [Tm \ a]) \ ] \rangle$

**abbreviation**  $\text{wrap2} :: \langle 'n \Rightarrow 'n \Rightarrow 'n \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ syms} \rangle$  **where**

$\langle \text{wrap2 } A \ B \ C \equiv$   
 $[ \ Tm \ [^1(A, [Nt \ B, Nt \ C]),$   
 $\quad Nt \ B,$   
 $\quad Tm \ ]^1(A, [Nt \ B, Nt \ C]),$   
 $\quad Tm \ [^2(A, [Nt \ B, Nt \ C]),$   
 $\quad Nt \ C,$

$$Tm ]^2 (A, [Nt B, Nt C]) ] \rangle$$

The transformation of old productions to new productions used in the proof:

**fun** *transform\_rhs* ::  $\langle 'n \Rightarrow ('n, 't) \text{ syms} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ syms} \rangle$  **where**  
 $\langle \text{transform\_rhs } A \text{ } [Tm \ a] = \text{wrap1 } A \ a \rangle \mid$   
 $\langle \text{transform\_rhs } A \text{ } [Nt \ B, Nt \ C] = \text{wrap2 } A \ B \ C \rangle$

The last equation is only added to permit us to state lemmas about

**fun** *transform\_prod* ::  $\langle 'n, 't \rangle \text{ prod} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ prod}$  **where**  
 $\langle \text{transform\_prod } (A, \alpha) = (A, \text{transform\_rhs } A \ \alpha) \rangle$

## 2.3 Homomorphisms

Definition of a monoid-homomorphism where multiplication is (@):

**definition** *hom\_list* ::  $\langle ('a \text{ list} \Rightarrow 'b \text{ list}) \Rightarrow \text{bool} \rangle$  **where**  
 $\langle \text{hom\_list } h = (\forall a \ b. \ h \ (a \ @ \ b) = h \ a \ @ \ h \ b) \rangle$

**lemma** *hom\_list\_Nil*:  $\text{hom\_list } h \Longrightarrow h \ [] = []$   
 $\langle \text{proof} \rangle$

The homomorphism on single brackets:

**fun** *the\_hom1* ::  $\langle ('n, 't) \text{ bracket3} \Rightarrow 't \text{ list} \rangle$  **where**  
 $\langle \text{the\_hom1 } [^1 (A, [Tm \ a])] = [a] \rangle \mid$   
 $\langle \text{the\_hom1 } \_ = [] \rangle$

The homomorphism on single bracket symbols:

**fun** *the\_hom\_sym* ::  $\langle ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow ('n, 't) \text{ sym list} \rangle$  **where**  
 $\langle \text{the\_hom\_sym } (Tm \ [^1 (A, [Tm \ a])]) = [Tm \ a] \rangle \mid$   
 $\langle \text{the\_hom\_sym } (Nt \ A) = [Nt \ A] \rangle \mid$   
 $\langle \text{the\_hom\_sym } \_ = [] \rangle$

The homomorphism on bracket words:

**fun** *the\_hom* ::  $\langle ('n, 't) \text{ bracket3 list} \Rightarrow 't \text{ list} \rangle$  (h) **where**  
 $\langle \text{the\_hom } l = \text{concat } (\text{map } \text{the\_hom1 } l) \rangle$

The homomorphism extended to symbols:

**fun** *the\_hom\_syms* ::  $\langle ('n, ('n, 't) \text{ bracket3}) \text{ syms} \Rightarrow ('n, 't) \text{ syms} \rangle$  **where**  
 $\langle \text{the\_hom\_syms } l = \text{concat } (\text{map } \text{the\_hom\_sym } l) \rangle$

**notation** *the\_hom* (h)

**notation** *the\_hom\_syms* (hs)

**lemma** *the\_hom\_syms\_hom*:  $\langle \text{hs } (l1 \ @ \ l2) = \text{hs } l1 \ @ \ \text{hs } l2 \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *the\_hom\_syms\_keep\_var*:  $\langle \text{hs } [(Nt \ A)] = [Nt \ A] \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *the\_hom\_syms\_tms\_inj*:  $\langle \text{hs } w = \text{map } Tm \ m \implies \exists w'. w = \text{map } Tm \ w' \rangle$

*<proof>*

Helper for showing the upcoming lemma:

**lemma** *helper*:  $\langle \text{the\_hom\_sym } (Tm \ x) = \text{map } Tm \ (\text{the\_hom1 } x) \rangle$

*<proof>*

Show that the extension really is an extension in some sense:

**lemma** *h\_eq\_h\_ext*:  $\langle \text{hs } (\text{map } Tm \ x) = \text{map } Tm \ (h \ x) \rangle$

*<proof>*

**lemma** *the\_hom1\_strip*:  $\langle (\text{the\_hom\_sym } x') = \text{map } Tm \ w \implies \text{the\_hom1 } (\text{destTm } x') = w \rangle$

*<proof>*

**lemma** *the\_hom1\_strip2*:  $\langle \text{concat } (\text{map } \text{the\_hom\_sym } w') = \text{map } Tm \ w \implies \text{concat } (\text{map } (\text{the\_hom1 } \circ \text{destTm}) \ w') = w \rangle$

*<proof>*

**lemma** *h\_eq\_h\_ext2*:

**assumes**  $\langle \text{hs } w' = (\text{map } Tm \ w) \rangle$

**shows**  $\langle h \ (\text{map } \text{destTm } w') = w \rangle$

*<proof>*

### 3 The Regular Language

The regular Language *Reg* will be an intersection of 5 Languages. The languages 2, 3, 4 are defined each via a relation *P2*, *P3*, *P4* on neighbouring letters and lifted to a language via *successively*. Language 1 is an intersection of another such lifted relation *P1'* and a condition on the last letter (if existent). Language 5 is a condition on the first letter (and requires it to exist). It takes a term of type '*n*' (the original variable type) as parameter.

Additionally a version of each language (taking symbols as input) is defined which allows arbitrary interspersions of nonterminals.

As this interspersions weakens the description, the symbol version of the regular language (*Reg\_sym*) is defined using two additional languages lifted from *P7* and *P8*. These vanish when restricted to words only containing terminals.

As stated in the introductory text, these languages will only be regular, when constrained to a finite bracket set. The theorems about this, are in the later section *Showing Regularity*.

### 3.1 $P1$

$P1$  will define a predicate on string elements. It will be true iff each  $]_p^1$  is directly followed by  $]_p^2$ . That also means  $]_p^1$  cannot be the end of the string.

But first we define a helper function, that only captures the neighbouring condition for two strings:

```
fun  $P1'$  ::  $\langle ('n, 't) \text{ bracket3} \Rightarrow ('n, 't) \text{ bracket3} \Rightarrow \text{bool} \rangle$  where
   $\langle P1' ]_p^1 [^2_p' = (p = p') \rangle$  |
   $\langle P1' ]_p^1 y = \text{False} \rangle$  |
   $\langle P1' x y = \text{True} \rangle$ 
```

A version of  $P1'$  for symbols, i.e. strings that may still contain Nt's:

```
fun  $P1\_sym$  ::  $\langle ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow \text{bool} \rangle$ 
where
   $\langle P1\_sym (Tm ]_p^1) (Tm [^2_p') = (p = p') \rangle$  |
   $\langle P1\_sym (Tm ]_p^1) y = \text{False} \rangle$  |
   $\langle P1\_sym x y = \text{True} \rangle$ 
```

```
lemma  $P1'D[simp]$ :
   $\langle P1' ]_p^1 r \longleftrightarrow r = [^2_p \rangle$ 
 $\langle \text{proof} \rangle$ 
```

Asserts that  $P1'$  holds for every pair in  $xs$ , and that  $xs$  doesn't end in  $]_p^1$ :

```
fun  $P1$  ::  $\langle ('n, 't) \text{ bracket3 list} \Rightarrow \text{bool} \rangle$  where
   $\langle P1 xs = ((\text{successively } P1' xs) \wedge (\text{if } xs \neq [] \text{ then } (\nexists p. \text{last } xs = ]_p^1) \text{ else True})) \rangle$ 
```

Asserts that  $P1'$  holds for every pair in  $xs$ , and that  $xs$  doesn't end in  $Tm ]_p^1$ :

```
fun  $P1\_sym$  where
   $\langle P1\_sym xs = ((\text{successively } P1\_sym xs) \wedge (\text{if } xs \neq [] \text{ then } (\nexists p. \text{last } xs = Tm ]_p^1) \text{ else True})) \rangle$ 
```

```
lemma  $P1\_for\_tm\_if\_P1\_sym[dest!]$ :  $\langle P1\_sym (\text{map } Tm x) \implies P1 x \rangle$ 
 $\langle \text{proof} \rangle$ 
```

```
lemma  $P1I[intro]$ :
  assumes  $\langle \text{successively } P1' xs \rangle$ 
  and  $\langle \nexists p. \text{last } xs = ]_p^1 \rangle$ 
  shows  $\langle P1 xs \rangle$ 
 $\langle \text{proof} \rangle$ 
```

```
lemma  $P1\_symI[intro]$ :
  assumes  $\langle \text{successively } P1\_sym xs \rangle$ 
  and  $\langle \nexists p. \text{last } xs = Tm ]_p^1 \rangle$ 
  shows  $\langle P1\_sym xs \rangle$ 
 $\langle \text{proof} \rangle$ 
```



**lemma**  $P1\_symD[dest]$ :  $\langle P1\_sym\ xs \implies successively\ P1'\_sym\ xs \rangle \langle proof \rangle$

**lemma**  $P1D\_not\_empty[intro]$ :  
**assumes**  $\langle xs \neq [] \rangle$   
**and**  $\langle P1\ xs \rangle$   
**shows**  $\langle last\ xs \neq ]^1_p \rangle$   
 $\langle proof \rangle$

**lemma**  $P1\_symD\_not\_empty'[intro]$ :  
**assumes**  $\langle xs \neq [] \rangle$   
**and**  $\langle P1\_sym\ xs \rangle$   
**shows**  $\langle last\ xs \neq Tm\ ]^1_p \rangle$   
 $\langle proof \rangle$

**lemma**  $P1\_symD\_not\_empty$ :  
**assumes**  $\langle xs \neq [] \rangle$   
**and**  $\langle P1\_sym\ xs \rangle$   
**shows**  $\langle \nexists p. last\ xs = Tm\ ]^1_p \rangle$   
 $\langle proof \rangle$

### 3.2 $P2$

$A\ ]^2_\pi$  is never directly followed by some  $[\cdot$ :

**fun**  $P2 :: \langle ('n, 't)\ bracket3 \Rightarrow ('n, 't)\ bracket3 \Rightarrow bool \rangle$  **where**  
 $\langle P2\ (Close\ (p, Two))\ (Open\ (p', v)) = False \rangle \mid$   
 $\langle P2\ (Close\ (p, Two))\ y = True \rangle \mid$   
 $\langle P2\ x\ y = True \rangle$

**fun**  $P2\_sym :: \langle ('n, ('n, 't)\ bracket3)\ sym \Rightarrow ('n, ('n, 't)\ bracket3)\ sym \Rightarrow bool \rangle$   
**where**  
 $\langle P2\_sym\ (Tm\ (Close\ (p, Two)))\ (Tm\ (Open\ (p', v))) = False \rangle \mid$   
 $\langle P2\_sym\ (Tm\ (Close\ (p, Two)))\ y = True \rangle \mid$   
 $\langle P2\_sym\ x\ y = True \rangle$

**lemma**  $P2\_for\_tm\_if\_P2\_sym[dest]$ :  $\langle successively\ P2\_sym\ (map\ Tm\ x) \implies successively\ P2\ x \rangle$   
 $\langle proof \rangle$

### 3.3 $P3$

Each  $[^1_{A \rightarrow BC}$  is directly followed by  $[^1_{B \rightarrow \_}$ , and each  $[^2_{A \rightarrow BC}$  is directly followed by  $[^1_{C \rightarrow \_}$ :

**fun**  $P3 :: \langle ('n, 't)\ bracket3 \Rightarrow ('n, 't)\ bracket3 \Rightarrow bool \rangle$  **where**  
 $\langle P3\ [^1_{(A, [Nt\ B, Nt\ C])}\ (p, ((X, y), t)) = (p = True \wedge t = One \wedge X = B) \rangle \mid$   
 $\langle P3\ [^2_{(A, [Nt\ B, Nt\ C])}\ (p, ((X, y), t)) = (p = True \wedge t = One \wedge X = C) \rangle \mid$   
 $\langle P3\ x\ y = True \rangle$

Each  $[^1_{A \rightarrow BC}$  is directly followed  $[^1_{B \rightarrow \_}$  or  $Nt\ B$ , and each  $[^2_{A \rightarrow BC}$  is directly followed by  $[^1_{C \rightarrow \_}$  or  $Nt\ C$ :

**fun**  $P3\_sym :: \langle ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow \text{bool} \rangle$   
**where**  
 $\langle P3\_sym (Tm \llbracket^1 (A, [Nt B, Nt C]) \rrbracket) (Tm (p, ((X, y), t))) = (p = \text{True} \wedge t = \text{One} \wedge X = B) \rangle \mid$   
 — Not obvious: the case  $(Tm \llbracket^1 (A, [Nt B, Nt C]) \rrbracket) Nt X$  is set to  $\text{True}$  with the catch all  
 $\langle P3\_sym (Tm \llbracket^1 (A, [Nt B, Nt C]) \rrbracket) (Nt X) = (X = B) \rangle \mid$   
 $\langle P3\_sym (Tm \llbracket^2 (A, [Nt B, Nt C]) \rrbracket) (Tm (p, ((X, y), t))) = (p = \text{True} \wedge t = \text{One} \wedge X = C) \rangle \mid$   
 $\langle P3\_sym (Tm \llbracket^2 (A, [Nt B, Nt C]) \rrbracket) (Nt X) = (X = C) \rangle \mid$   
 $\langle P3\_sym x y = \text{True} \rangle$

**lemma**  $P3D1[dest]$ :  
**fixes**  $r :: \langle ('n, 't) \text{ bracket3} \rangle$   
**assumes**  $\langle P3 \llbracket^1 (A, [Nt B, Nt C]) \rrbracket r \rangle$   
**shows**  $\langle \exists l. r = \llbracket^1 (B, l) \rrbracket \rangle$   
 $\langle \text{proof} \rangle$

**lemma**  $P3D2[dest]$ :  
**fixes**  $r :: \langle ('n, 't) \text{ bracket3} \rangle$   
**assumes**  $\langle P3 \llbracket^2 (A, [Nt B, Nt C]) \rrbracket r \rangle$   
**shows**  $\langle \exists l. r = \llbracket^1 (C, l) \rrbracket \rangle$   
 $\langle \text{proof} \rangle$

**lemma**  $P3\_for\_tm\_if\_P3\_sym[dest]$ :  $\langle \text{successively } P3\_sym (\text{map } Tm x) \Longrightarrow \text{successively } P3 x \rangle$   
 $\langle \text{proof} \rangle$

### 3.4 $P4$

Each  $\llbracket^1_{A \rightarrow a}$  is directly followed by  $\rrbracket^1_{A \rightarrow a}$  and each  $\llbracket^2_{A \rightarrow a}$  is directly followed by  $\rrbracket^2_{A \rightarrow a}$ :

**fun**  $P4 :: \langle ('n, 't) \text{ bracket3} \Rightarrow ('n, 't) \text{ bracket3} \Rightarrow \text{bool} \rangle$  **where**  
 $\langle P4 (\text{Open} ((A, [Tm a]), s)) (p, ((X, y), t)) = (p = \text{False} \wedge X = A \wedge y = [Tm a] \wedge s = t) \rangle \mid$   
 $\langle P4 x y = \text{True} \rangle$

Each  $\llbracket^1_{A \rightarrow a}$  is directly followed by  $\rrbracket^1_{A \rightarrow a}$  and each  $\llbracket^2_{A \rightarrow a}$  is directly followed by  $\rrbracket^2_{A \rightarrow a}$ :

**fun**  $P4\_sym :: \langle ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow \text{bool} \rangle$   
**where**  
 $\langle P4\_sym (Tm (\text{Open} ((A, [Tm a]), s))) (Tm (p, ((X, y), t))) = (p = \text{False} \wedge X = A \wedge y = [Tm a] \wedge s = t) \rangle \mid$   
 $\langle P4\_sym (Tm (\text{Open} ((A, [Tm a]), s))) (Nt X) = \text{False} \rangle \mid$   
 $\langle P4\_sym x y = \text{True} \rangle$

**lemma**  $P4D[dest]$ :  
**fixes**  $r :: \langle 'n, 't \rangle \text{ bracket3} \rangle$   
**assumes**  $\langle P4 \text{ (Open ((A, [Tm a]), v)) } r \rangle$   
**shows**  $\langle r = \text{Close ((A, [Tm a]), v)} \rangle$   
 $\langle \text{proof} \rangle$

**lemma**  $P4\_for\_tm\_if\_P4\_sym[dest]$ :  $\langle \text{successively } P4\_sym \text{ (map Tm } x) \implies \text{successively } P4 \text{ } x \rangle$   
 $\langle \text{proof} \rangle$

### 3.5 $P5$

$P5 \ A \ x$  holds, iff there exists some  $y$  such that  $x$  begins with  $[^1_{A \rightarrow y}$ :

**fun**  $P5 :: \langle 'n \Rightarrow ('n, 't) \text{ bracket3 list} \Rightarrow \text{bool} \rangle$  **where**  
 $\langle P5 \ A \ [] = \text{False} \rangle \mid$   
 $\langle P5 \ A \ ([^1_{(X,x)} \# xs) = (X = A) \rangle \mid$   
 $\langle P5 \ A \ (x \# xs) = \text{False} \rangle$

$P5\_sym \ A \ x$  holds, iff either there exists some  $y$  such that  $x$  begins with  $[^1_{A \rightarrow y}$ , or if it begins with  $Nt \ A$ :

**fun**  $P5\_sym :: \langle 'n \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ syms} \Rightarrow \text{bool} \rangle$  **where**  
 $\langle P5\_sym \ A \ [] = \text{False} \rangle \mid$   
 $\langle P5\_sym \ A \ (\text{Tm } [^1_{(X,x)} \# xs) = (X = A) \rangle \mid$   
 $\langle P5\_sym \ A \ ((Nt \ X) \# xs) = (X = A) \rangle \mid$   
 $\langle P5\_sym \ A \ (x \# xs) = \text{False} \rangle$

**lemma**  $P5D[dest]$ :  
**assumes**  $\langle P5 \ A \ x \rangle$   
**shows**  $\langle \exists y. \text{hd } x = [^1_{(A,y)} \rangle$   
 $\langle \text{proof} \rangle$

**lemma**  $P5\_symD[dest]$ :  
**assumes**  $\langle P5\_sym \ A \ x \rangle$   
**shows**  $\langle (\exists y. \text{hd } x = \text{Tm } [^1_{(A,y)}) \vee \text{hd } x = Nt \ A \rangle$   
 $\langle \text{proof} \rangle$

**lemma**  $P5\_for\_tm\_if\_P5\_sym[dest]$ :  $\langle P5\_sym \ A \ (\text{map Tm } x) \implies P5 \ A \ x \rangle$   
 $\langle \text{proof} \rangle$

### 3.6 $P7$ and $P8$

$(\text{successively } P7\_sym) \ w$  iff  $Nt \ Y$  is directly preceded by some  $\text{Tm } [^1_{A \rightarrow YC}$  or  $\text{Tm } [^2_{A \rightarrow BY}$  in  $w$ :

**fun**  $P7\_sym :: \langle ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ sym} \Rightarrow \text{bool} \rangle$   
**where**  
 $\langle P7\_sym \ (\text{Tm } (b, (A, [Nt \ B, Nt \ C]), v)) \ (Nt \ Y) = (b = \text{True} \wedge ((Y = B \wedge v = \text{One}) \vee (Y = C \wedge v = \text{Two}))) \rangle \mid$

$\langle P\gamma\_sym\ x\ (Nt\ Y) = False \rangle \mid$   
 $\langle P\gamma\_sym\ x\ y = True \rangle$

**lemma**  $P\gamma\_symD[dest]$ :  
**fixes**  $x:: \langle ('n, ('n, 't)\ bracket3)\ sym \rangle$   
**assumes**  $\langle P\gamma\_sym\ x\ (Nt\ Y) \rangle$   
**shows**  $\langle (\exists A\ C. x = Tm\ ]^1(A, [Nt\ Y, Nt\ C]) \rangle \vee (\exists A\ B. x = Tm\ ]^2(A, [Nt\ B, Nt\ Y]) \rangle$   
 $\langle proof \rangle$

$(successively\ P8\_sym)\ w$  iff  $Nt\ Y$  is directly followed by some  $]^1_{A \rightarrow YC}$   
or  $]^2_{A \rightarrow BY}$  in  $w$ :

**fun**  $P8\_sym :: \langle ('n, ('n, 't)\ bracket3)\ sym \Rightarrow ('n, ('n, 't)\ bracket3)\ sym \Rightarrow bool \rangle$   
**where**  
 $\langle P8\_sym\ (Nt\ Y)\ (Tm\ (b, (A, [Nt\ B, Nt\ C]), v)) = (b = False \wedge (Y = B \wedge v = One) \vee (Y = C \wedge v = Two)) \rangle \mid$   
 $\langle P8\_sym\ (Nt\ Y)\ x = False \rangle \mid$   
 $\langle P8\_sym\ x\ y = True \rangle$

**lemma**  $P8\_symD[dest]$ :  
**fixes**  $x:: \langle ('n, ('n, 't)\ bracket3)\ sym \rangle$   
**assumes**  $\langle P8\_sym\ (Nt\ Y)\ x \rangle$   
**shows**  $\langle (\exists A\ C. x = Tm\ ]^1(A, [Nt\ Y, Nt\ C]) \rangle \vee (\exists A\ B. x = Tm\ ]^2(A, [Nt\ B, Nt\ Y]) \rangle$   
 $\langle proof \rangle$

### 3.7 $Reg$ and $Reg\_sym$

This is the regular language, where one takes the Start symbol as a parameter, and then has the searched for  $R := R_A$ :

**definition**  $Reg :: \langle 'n \Rightarrow ('n, 't)\ bracket3\ list\ set \rangle$  **where**  
 $\langle Reg\ A = \{x. (P1\ x) \wedge$   
 $(successively\ P2\ x) \wedge$   
 $(successively\ P3\ x) \wedge$   
 $(successively\ P4\ x) \wedge$   
 $(P5\ A\ x)\} \rangle$

**lemma**  $RegI[intro]$ :  
**assumes**  $\langle (P1\ x) \rangle$   
**and**  $\langle (successively\ P2\ x) \rangle$   
**and**  $\langle (successively\ P3\ x) \rangle$   
**and**  $\langle (successively\ P4\ x) \rangle$   
**and**  $\langle (P5\ A\ x) \rangle$   
**shows**  $\langle x \in Reg\ A \rangle$   
 $\langle proof \rangle$

**lemma**  $RegD[dest]$ :  
**assumes**  $\langle x \in Reg\ A \rangle$   
**shows**  $\langle (P1\ x) \rangle$   
**and**  $\langle (successively\ P2\ x) \rangle$

```

and  $\langle \text{successively } P3 \ x \rangle$ 
and  $\langle \text{successively } P4 \ x \rangle$ 
and  $\langle P5 \ A \ x \rangle$ 
 $\langle \text{proof} \rangle$ 

```

A version of *Reg* for symbols, i.e. strings that may still contain Nt's. It has 2 more Properties *P7* and *P8* that vanish for pure terminal strings:

```

definition Reg_sym ::  $\langle 'n \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ syms set} \rangle$  where
 $\langle \text{Reg\_sym } A = \{x. (P1\_sym \ x) \wedge$ 
   $(\text{successively } P2\_sym \ x) \wedge$ 
   $(\text{successively } P3\_sym \ x) \wedge$ 
   $(\text{successively } P4\_sym \ x) \wedge$ 
   $(P5\_sym \ A \ x) \wedge$ 
   $(\text{successively } P7\_sym \ x) \wedge$ 
   $(\text{successively } P8\_sym \ x)\} \rangle$ 

```

```

lemma Reg_symI[intro]:
assumes  $\langle P1\_sym \ x \rangle$ 
and  $\langle \text{successively } P2\_sym \ x \rangle$ 
and  $\langle \text{successively } P3\_sym \ x \rangle$ 
and  $\langle \text{successively } P4\_sym \ x \rangle$ 
and  $\langle P5\_sym \ A \ x \rangle$ 
and  $\langle \text{successively } P7\_sym \ x \rangle$ 
and  $\langle \text{successively } P8\_sym \ x \rangle$ 
shows  $\langle x \in \text{Reg\_sym } A \rangle$ 
 $\langle \text{proof} \rangle$ 

```

```

lemma Reg_symD[dest]:
assumes  $\langle x \in \text{Reg\_sym } A \rangle$ 
shows  $\langle P1\_sym \ x \rangle$ 
and  $\langle \text{successively } P2\_sym \ x \rangle$ 
and  $\langle \text{successively } P3\_sym \ x \rangle$ 
and  $\langle \text{successively } P4\_sym \ x \rangle$ 
and  $\langle P5\_sym \ A \ x \rangle$ 
and  $\langle \text{successively } P7\_sym \ x \rangle$ 
and  $\langle \text{successively } P8\_sym \ x \rangle$ 
 $\langle \text{proof} \rangle$ 

```

```

lemma Reg_for_tm_if_Reg_sym[dest]:  $\langle (\text{map } Tm \ x) \in \text{Reg\_sym } A \implies x \in \text{Reg } A \rangle$ 
 $\langle \text{proof} \rangle$ 

```

## 4 Showing Regularity

```

context locale_P
begin

```

```

abbreviation brackets:: $\langle ('n, 't) \text{ bracket3 list set} \rangle$  where
 $\langle \text{brackets} \equiv \{bs. \forall (\_, p, \_) \in \text{set } bs. p \in P\} \rangle$ 

```

This is needed for the construction that shows P2,P3,P4 regular.

**datatype** 'a state = start | garbage | letter 'a

**definition** allStates :: ⟨('n,'t) bracket3 state set ⟩ **where**  
 ⟨allStates = { letter (br,(p,v)) | br p v. p ∈ P } ∪ {start, garbage}⟩

**lemma** allStatesI: ⟨p ∈ P ⟹ letter (br,(p,v)) ∈ allStates⟩  
 ⟨proof⟩

**lemma** start\_in\_allStates[simp]: ⟨start ∈ allStates⟩  
 ⟨proof⟩

**lemma** garbage\_in\_allStates[simp]: ⟨garbage ∈ allStates⟩  
 ⟨proof⟩

**lemma** finite\_allStates\_if:  
 shows ⟨finite( allStates)⟩  
 ⟨proof⟩

**end**

#### 4.1 An automaton for {xs. successively Q xs ∧ xs ∈ brackets P}

**locale** successivelyConstruction = locale\_P **where** P = P **for** P :: ('n,'t) Prods  
 +  
**fixes** Q :: ('n,'t) bracket3 ⇒ ('n,'t) bracket3 ⇒ bool — e.g. P2  
**begin**

**fun** succNext :: ⟨('n,'t) bracket3 state ⇒ ('n,'t) bracket3 ⇒ ('n,'t) bracket3 state⟩  
**where**  
 ⟨succNext garbage \_ = garbage⟩ |  
 ⟨succNext start (br', p', v') = (if p' ∈ P then letter (br', p',v') else garbage)⟩ |  
 ⟨succNext (letter (br, p, v)) (br', p', v') = (if Q (br,p,v) (br',p',v') ∧ p ∈ P ∧  
 p' ∈ P then letter (br',p',v') else garbage)⟩

**theorem** succNext\_induct[case\_names garbage startp startnp letterQ letternQ]:  
**fixes** R :: ('n,'t) bracket3 state ⇒ ('n,'t) bracket3 ⇒ bool  
**fixes** a0 :: ('n,'t) bracket3 state  
**and** a1 :: ('n,'t) bracket3  
**assumes** ∧u. R garbage u  
**and** ∧br' p' v'. p' ∈ P ⟹ R state.start (br', p', v')  
**and** ∧br' p' v'. p' ∉ P ⟹ R state.start (br', p', v')  
**and** ∧br p v br' p' v'. Q (br,p,v) (br',p',v') ∧ p ∈ P ∧ p' ∈ P ⟹ R (letter  
 (br, p, v)) (br', p', v')  
**and** ∧br p v br' p' v'. ¬(Q (br,p,v) (br',p',v') ∧ p ∈ P ∧ p' ∈ P) ⟹ R (letter  
 (br, p, v)) (br', p', v')  
**shows** R a0 a1  
 ⟨proof⟩

**abbreviation** *aut* **where**  $\langle aut \equiv \langle dfa.states = allStates,$   
 $init = start,$   
 $final = (allStates - \{garbage\}),$   
 $next = succNext \rangle \rangle$

**interpretation** *aut* : *dfa aut*  
 $\langle proof \rangle$

**lemma** *nextl\_in\_allStates*[*intro,simp*]:  $\langle q \in allStates \implies aut.nextl\ q\ ys \in all-$   
 $States \rangle$   
 $\langle proof \rangle$

**lemma** *nextl\_garbage*[*simp*]:  $\langle aut.nextl\ garbage\ xs = garbage \rangle$   
 $\langle proof \rangle$

**lemma** *drop\_right*:  $\langle xs@ys \in aut.language \implies xs \in aut.language \rangle$   
 $\langle proof \rangle$

**lemma** *state\_after1*[*iff*]:  $\langle (succNext\ q\ a \neq garbage) = (succNext\ q\ a = letter\ a) \rangle$   
 $\langle proof \rangle$

**lemma** *state\_after\_in\_P*[*intro*]:  $\langle succNext\ q\ (br,\ p,\ v) \neq garbage \implies p \in P \rangle$   
 $\langle proof \rangle$

**lemma** *drop\_left\_general*:  $\langle aut.nextl\ start\ ys = garbage \implies aut.nextl\ q\ ys =$   
 $garbage \rangle$   
 $\langle proof \rangle$

**lemma** *drop\_left*:  $\langle xs@ys \in aut.language \implies ys \in aut.language \rangle$   
 $\langle proof \rangle$

**lemma** *empty\_in\_aut*:  $\langle [] \in aut.language \rangle$   
 $\langle proof \rangle$

**lemma** *singleton\_in\_aut\_iff*:  $\langle [(br,\ p,\ v)] \in aut.language \longleftrightarrow p \in P \rangle$   
 $\langle proof \rangle$

**lemma** *duo\_in\_aut\_iff*:  $\langle [(br,\ p,\ v), (br',\ p',\ v')] \in aut.language \longleftrightarrow Q\ (br,p,v)$   
 $(br',p',v') \wedge p \in P \wedge p' \in P \rangle$   
 $\langle proof \rangle$

**lemma** *trio\_in\_aut\_iff*:  $\langle (br,\ p,\ v) \# (br',\ p',\ v') \# zs \in aut.language \longleftrightarrow Q$   
 $(br,p,v)\ (br',p',v') \wedge p \in P \wedge p' \in P \wedge (br',p',v') \# zs \in aut.language \rangle$   
 $\langle proof \rangle$

**lemma** *aut\_lang\_iff\_succ\_Q*:  $\langle (successively\ Q\ xs \wedge xs \in brackets) \longleftrightarrow (xs \in$   
 $aut.language) \rangle$   
 $\langle proof \rangle$

**corollary** *regular\_successively\_inter\_brackets*:  $\langle \text{regular } \{xs. \text{ successively } Q \text{ } xs \wedge xs \in \text{brackets}\} \rangle$   
 $\langle \text{proof} \rangle$

**end**

## 4.2 Regularity of $P2$ , $P3$ and $P4$

**context** *locale\_P*  
**begin**

**lemma** *P2\_regular*:  
**shows**  $\langle \text{regular } \{xs. \text{ successively } P2 \text{ } xs \wedge xs \in \text{brackets}\} \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *P3\_regular*:  
 $\langle \text{regular } \{xs. \text{ successively } P3 \text{ } xs \wedge xs \in \text{brackets}\} \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *P4\_regular*:  
 $\langle \text{regular } \{xs. \text{ successively } P4 \text{ } xs \wedge xs \in \text{brackets}\} \rangle$   
 $\langle \text{proof} \rangle$

## 4.3 An automaton for $P1$

More Precisely, for the *if not empty, then doesnt end in (Close\_,1)* part.  
Then intersect with the other construction for  $P1'$  to get  $P1$  regular.

**datatype**  $P1\_State = \text{last\_ok} \mid \text{last\_bad} \mid \text{garbage}$

The good ending letters, are those that are not of the form  $(\text{Close } \_ , 1)$ .

**fun** *good where*  
 $\langle \text{good } ]^1_p = \text{False} \rangle \mid$   
 $\langle \text{good } (br, p, v) = \text{True} \rangle$

**fun** *next1* ::  $\langle P1\_State \Rightarrow ('n, 't) \text{ bracket3} \Rightarrow P1\_State \rangle$  **where**  
 $\langle \text{next1 } \text{garbage } \_ = \text{garbage} \rangle \mid$   
 $\langle \text{next1 } \text{last\_ok } (br, p, v) = (\text{if } p \notin P \text{ then garbage else (if good } (br, p, v) \text{ then last\_ok else last\_bad)}) \rangle \mid$   
 $\langle \text{next1 } \text{last\_bad } (br, p, v) = (\text{if } p \notin P \text{ then garbage else (if good } (br, p, v) \text{ then last\_ok else last\_bad)}) \rangle$

**theorem** *next1\_induct*[*case\_names garbage startp startnp letterQ letternQ*]:  
**fixes**  $R :: P1\_State \Rightarrow ('n, 't) \text{ bracket3} \Rightarrow \text{bool}$   
**fixes**  $a0 :: P1\_State$   
**and**  $a1 :: ('n, 't) \text{ bracket3}$   
**assumes**  $\bigwedge u. R \text{ garbage } u$   
**and**  $\bigwedge br \ p \ v. p \notin P \implies R \text{ last\_ok } (br, p, v)$   
**and**  $\bigwedge br \ p \ v. p \in P \wedge \text{good } (br, p, v) \implies R \text{ last\_ok } (br, p, v)$



**and**  $\bigwedge br\ p\ v. p \in P \wedge \neg(good\ (br, p, v)) \implies R\ last\_ok\ (br, p, v)$   
**and**  $\bigwedge br\ p\ v. p \notin P \implies R\ last\_bad\ (br, p, v)$   
**and**  $\bigwedge br\ p\ v. p \in P \wedge good\ (br, p, v) \implies R\ last\_bad\ (br, p, v)$   
**and**  $\bigwedge br\ p\ v. p \in P \wedge \neg(good\ (br, p, v)) \implies R\ last\_bad\ (br, p, v)$   
**shows**  $R\ a0\ a1$   
 $\langle proof \rangle$

**abbreviation**  $p1\_aut$  **where**  $\langle p1\_aut \equiv (dfa.states = \{last\_ok, last\_bad, garbage\},$   
 $init = last\_ok,$   
 $final = \{last\_ok\},$   
 $next = next1) \rangle$

**interpretation**  $p1\_aut : dfa\ p1\_aut$   
 $\langle proof \rangle$

**lemma**  $next1\_garbage\_iff[simp]: \langle p1\_aut.next1\ last\_ok\ xs = garbage \longleftrightarrow xs \notin$   
 $brackets \rangle$   
 $\langle proof \rangle$

**lemma**  $lang\_descr\_full:$   
 $\langle (p1\_aut.next1\ last\_ok\ xs = last\_ok \longleftrightarrow (xs = [] \vee (xs \neq [] \wedge good\ (last\ xs) \wedge$   
 $xs \in brackets))) \wedge$   
 $(p1\_aut.next1\ last\_ok\ xs = last\_bad \longleftrightarrow ((xs \neq [] \wedge \neg good\ (last\ xs) \wedge xs \in$   
 $brackets))) \rangle$   
 $\langle proof \rangle$

**lemma**  $lang\_descr: \langle xs \in p1\_aut.language \longleftrightarrow (xs = [] \vee (xs \neq [] \wedge good\ (last$   
 $xs) \wedge xs \in brackets)) \rangle$   
 $\langle proof \rangle$

**lemma**  $good\_iff[simp]: \langle (\forall a\ b. last\ xs \neq ]^1_{(a, b)}) = good\ (last\ xs) \rangle$   
 $\langle proof \rangle$

**lemma**  $in\_P1\_iff: \langle (P1\ xs \wedge xs \in brackets) \longleftrightarrow (xs = [] \vee (xs \neq [] \wedge good\ (last$   
 $xs) \wedge xs \in brackets)) \wedge successively\ P1'\ xs \wedge xs \in brackets \rangle$   
 $\langle proof \rangle$

**corollary**  $P1\_eq: \langle \{xs. P1\ xs \wedge xs \in brackets\} =$   
 $\{xs. successively\ P1'\ xs \wedge xs \in brackets\} \cap \{xs. xs = [] \vee (xs \neq [] \wedge good$   
 $(last\ xs) \wedge xs \in brackets)\} \rangle$   
 $\langle proof \rangle$

**lemma**  $P1'\_regular:$   
**shows**  $\langle regular\ \{xs. successively\ P1'\ xs \wedge xs \in brackets\} \rangle$   
 $\langle proof \rangle$

**corollary**  $aux\_regular: \langle regular\ \{xs. xs = [] \vee (xs \neq [] \wedge good\ (last\ xs) \wedge xs \in$   
 $brackets)\} \rangle$   
 $\langle proof \rangle$

**corollary** *regular\_P1*:  $\langle \text{regular } \{xs. P1\ xs \wedge xs \in \text{brackets}\} \rangle$   
 $\langle \text{proof} \rangle$

**end**

#### 4.4 An automaton for *P5*

**locale** *P5Construction* = *locale\_P* **where**  $P=P$  **for**  $P :: ('n, 't) \text{Prods} +$   
**fixes**  $A :: 'n$   
**begin**

**datatype** *P5\_State* = *start* | *first\_ok* | *garbage*

The good/ok ending letters, are those that are not of the form (*Close* \_\_ , 1).

**fun** *ok* **where**  
 $\langle \text{ok } (\text{Open } ((X, \_), \text{One})) = (X = A) \rangle$  |  
 $\langle \text{ok } \_ = \text{False} \rangle$

**fun** *next2* ::  $\langle P5\_State \Rightarrow ('n, 't) \text{ bracket3} \Rightarrow P5\_State \rangle$  **where**  
 $\langle \text{next2 } \text{garbage } \_ = \text{garbage} \rangle$  |  
 $\langle \text{next2 } \text{start } (br, (X, r), v) = (\text{if } (X, r) \notin P \text{ then garbage else } (\text{if } \text{ok } (br, (X, r), v) \text{ then first\_ok else garbage})) \rangle$  |  
 $\langle \text{next2 } \text{first\_ok } (br, p, v) = (\text{if } p \notin P \text{ then garbage else first\_ok}) \rangle$

**theorem** *next2\_induct*[*case\_names* *garbage startnp start\_p\_ok start\_p\_nok first\_ok\_np first\_ok\_p*]:

**fixes**  $R :: P5\_State \Rightarrow ('n, 't) \text{ bracket3} \Rightarrow \text{bool}$   
**fixes**  $a0 :: P5\_State$   
**and**  $a1 :: ('n, 't) \text{ bracket3}$   
**assumes**  $\bigwedge u. R \text{ garbage } u$   
**and**  $\bigwedge br\ p\ v. p \notin P \implies R \text{ start } (br, p, v)$   
**and**  $\bigwedge br\ X\ r\ v. (X, r) \in P \wedge \text{ok } (br, (X, r), v) \implies R \text{ start } (br, (X, r), v)$   
**and**  $\bigwedge br\ X\ r\ v. (X, r) \in P \wedge \neg \text{ok } (br, (X, r), v) \implies R \text{ start } (br, (X, r), v)$   
**and**  $\bigwedge br\ X\ r\ v. (X, r) \notin P \implies R \text{ first\_ok } (br, (X, r), v)$   
**and**  $\bigwedge br\ X\ r\ v. (X, r) \in P \implies R \text{ first\_ok } (br, (X, r), v)$   
**shows**  $R\ a0\ a1$   
 $\langle \text{proof} \rangle$

**abbreviation** *p5\_aut* **where**  $\langle p5\_aut \equiv (\text{dfa.states} = \{\text{start}, \text{first\_ok}, \text{garbage}\},$   
 $\text{init} = \text{start},$   
 $\text{final} = \{\text{first\_ok}\},$   
 $\text{next} = \text{next2}) \rangle$

**interpretation** *p5\_aut* : *dfa* *p5\_aut*  
 $\langle \text{proof} \rangle$

**corollary** *next2\_start\_ok\_iff*:  $\langle \text{ok } x \wedge \text{fst}(\text{snd } x) \in P \longleftrightarrow \text{next2 } \text{start } x = \text{first\_ok} \rangle$

$\langle \text{proof} \rangle$

**lemma** *empty\_not\_in\_lang[simp]*:  $\langle [] \notin p5\_aut.language \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *singleton\_in\_lang\_iff*:  $\langle [x] \in p5\_aut.language \longleftrightarrow ok\ (hd\ [x]) \wedge [x] \in brackets \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *singleton\_first\_ok\_iff*:  $\langle p5\_aut.nextl\ start\ ([x]) = first\_ok \vee p5\_aut.nextl\ start\ ([x]) = garbage \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *first\_ok\_iff*:  $\langle xs \neq [] \implies p5\_aut.nextl\ start\ xs = first\_ok \vee p5\_aut.nextl\ start\ xs = garbage \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *lang\_descr*:  $\langle xs \in p5\_aut.language \longleftrightarrow (xs \neq [] \wedge ok\ (hd\ xs) \wedge xs \in brackets) \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *in\_P5\_iff*:  $\langle P5\ A\ xs \wedge xs \in brackets \longleftrightarrow (xs \neq [] \wedge ok\ (hd\ xs) \wedge xs \in brackets) \rangle$   
 $\langle \text{proof} \rangle$

**corollary** *aux\_regular*:  $\langle regular\ \{xs.\ xs \neq [] \wedge ok\ (hd\ xs) \wedge xs \in brackets\} \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *regular\_P5*:  $\langle regular\ \{xs.\ P5\ A\ xs \wedge xs \in brackets\} \rangle$   
 $\langle \text{proof} \rangle$

**end**

**context** *locale\_P*  
**begin**

**corollary** *regular\_Reg\_inter*:  $\langle regular\ (brackets \cap Reg\ A) \rangle$   
 $\langle \text{proof} \rangle$

A lemma saying that all *Dyck\_lang* words really only consist of brackets (trivial definition wrangling):

**lemma** *Dyck\_lang\_subset\_brackets*:  $\langle Dyck\_lang\ (P \times \{One, Two\}) \subseteq brackets \rangle$   
 $\langle \text{proof} \rangle$

**end**

## 5 Definitions of $L, \Gamma, P', L'$

**locale** *Chomsky\_Schuetzenberger\_locale* = *locale\_P* **where**  $P = P$  **for**  $P :: ('n, 't) Prods$   
 $+$   
**fixes**  $S :: 'n$   
**assumes**  $CNF\_P$ :  $\langle CNF\ P \rangle$

**begin**

**lemma**  $P\_CNFE[dest]$ :  
**assumes**  $\langle \pi \in P \rangle$   
**shows**  $\langle \exists A\ a\ B\ C. \pi = (A, [Nt\ B, Nt\ C]) \vee \pi = (A, [Tm\ a]) \rangle$   
 $\langle proof \rangle$

**definition**  $L$  **where**

$\langle L = Lang\ P\ S \rangle$

**definition**  $\Gamma$  **where**

$\langle \Gamma = P \times \{One, Two\} \rangle$

**definition**  $P'$  **where**

$\langle P' = transform\_prod\ 'P \rangle$

**definition**  $L'$  **where**

$\langle L' = Lang\ P'\ S \rangle$

## 6 Lemmas for $P' \vdash A \Rightarrow^* x \longleftrightarrow x \in R_A \cap Dyck\_lang\ \Gamma$

**lemma**  $prod1\_snds\_in\_tm$  [*intro*, *simp*]:  $\langle (A, [Nt\ B, Nt\ C]) \in P \implies snds\_in\_tm\ \Gamma\ (wrap2\ A\ B\ C) \rangle$   
 $\langle proof \rangle$

**lemma**  $prod2\_snds\_in\_tm$  [*intro*, *simp*]:  $\langle (A, [Tm\ a]) \in P \implies snds\_in\_tm\ \Gamma\ (wrap1\ A\ a) \rangle$   
 $\langle proof \rangle$

**lemma**  $bal\_tm\_wrap1$  [*iff*]:  $\langle bal\_tm\ (wrap1\ A\ a) \rangle$   
 $\langle proof \rangle$

**lemma**  $bal\_tm\_wrap2$  [*iff*]:  $\langle bal\_tm\ (wrap2\ A\ B\ C) \rangle$   
 $\langle proof \rangle$

This essentially says, that the right sides of productions are in the Dyck language of  $\Gamma$ , if one ignores any occuring nonterminals. This will be needed for  $\rightarrow$ .

**lemma**  $bal\_tm\_transform\_rhs$  [*intro!*]:  
 $\langle (A, \alpha) \in P \implies bal\_tm\ (transform\_rhs\ A\ \alpha) \rangle$

$\langle \text{proof} \rangle$

**lemma** *snds\_in\_tm\_transform\_rhs*[intro!]:  
 $\langle (A, \alpha) \in P \implies \text{snds\_in\_tm } \Gamma \text{ (transform\_rhs } A \text{ } \alpha) \rangle$   
 $\langle \text{proof} \rangle$

The lemma for  $\rightarrow$

**lemma** *P'\_imp\_bal*:  
**assumes**  $\langle P' \vdash [Nt \ A] \Rightarrow^* x \rangle$   
**shows**  $\langle \text{bal\_tm } x \wedge \text{snds\_in\_tm } \Gamma \ x \rangle$   
 $\langle \text{proof} \rangle$

Another lemma for  $\rightarrow$

**lemma** *P'\_imp\_Reg*:  
**assumes**  $\langle P' \vdash [Nt \ T] \Rightarrow^* x \rangle$   
**shows**  $\langle x \in \text{Reg\_sym } T \rangle$   
 $\langle \text{proof} \rangle$

This will be needed for the direction  $\leftarrow$ .

**lemma** *transform\_prod\_one\_step*:  
**assumes**  $\langle \pi \in P \rangle$   
**shows**  $\langle P' \vdash [Nt \ (\text{fst } \pi)] \Rightarrow \text{snd } (\text{transform\_prod } \pi) \rangle$   
 $\langle \text{proof} \rangle$

The lemma for  $\leftarrow$

**lemma** *Reg\_and\_dyck\_imp\_P'*:  
**assumes**  $\langle x \in (\text{Reg } A \cap \text{Dyck\_lang } \Gamma) \rangle$   
**shows**  $\langle P' \vdash [Nt \ A] \Rightarrow^* \text{map } Tm \ x \rangle \langle \text{proof} \rangle$

## 7 Showing $h(L') = L$

Particularly  $\supseteq$  is formally hard. To create the witness in  $L'$  we need to use the corresponding production in  $P'$  in each step. We do this by defining the transformation on the parse tree, instead of only the word. Simple induction on the derivation wouldn't (in the induction step) get us enough information on where the corresponding production needs to be applied in the transformed version.

**abbreviation**  $\langle \text{roots } ts \equiv \text{map } \text{root } ts \rangle$

**fun** *wrap1\_Sym* ::  $\langle 'n \Rightarrow ('n, 't) \text{ sym} \Rightarrow \text{version} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ tree list} \rangle$   
**where**  
 $\text{wrap1\_Sym } A \text{ (Tm } a) \text{ } v = [ \text{Sym } (Tm \text{ (Open ((A, [Tm } a]), v))), \text{Sym } (Tm \text{ (Close ((A, [Tm } a]), v)))] \mid$   
 $\langle \text{wrap1\_Sym } \_ \_ \_ = [] \rangle$

**fun** *wrap2\_Sym* ::  $\langle 'n \Rightarrow ('n, 't) \text{ sym} \Rightarrow ('n, 't) \text{ sym} \Rightarrow \text{version} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ tree} \Rightarrow ('n, ('n, 't) \text{ bracket3}) \text{ tree list} \rangle$  **where**

$\text{wrap2\_Sym } A \text{ (Nt } B) \text{ (Nt } C) \text{ } v \text{ } t = [\text{Sym } (\text{Tm } (\text{Open } ((A, [\text{Nt } B, \text{Nt } C]), v))), t$   
 $, \text{Sym } (\text{Tm } (\text{Close } ((A, [\text{Nt } B, \text{Nt } C]), v)))] \mid$   
 $\langle \text{wrap2\_Sym } \_ \_ \_ \_ \_ = [] \rangle$

**fun** *transform\_tree* :: ('n,'t) tree  $\Rightarrow$  ('n,('n,'t) bracket3) tree **where**  
 $\langle \text{transform\_tree } (\text{Sym } (\text{Nt } A)) = (\text{Sym } (\text{Nt } A)) \rangle \mid$   
 $\langle \text{transform\_tree } (\text{Sym } (\text{Tm } a)) = (\text{Sym } (\text{Tm } [\text{SOME } A. \text{True}, [\text{Tm } a]))) \rangle \mid$   
 $\langle \text{transform\_tree } (\text{Rule } A [\text{Sym } (\text{Tm } a)]) = \text{Rule } A ((\text{wrap1\_Sym } A (\text{Tm } a)$   
 $\text{One}) @ (\text{wrap1\_Sym } A (\text{Tm } a) \text{ Two})) \rangle \mid$   
 $\langle \text{transform\_tree } (\text{Rule } A [t1, t2]) = \text{Rule } A ((\text{wrap2\_Sym } A (\text{root } t1) (\text{root } t2)$   
 $\text{One } (\text{transform\_tree } t1)) @ (\text{wrap2\_Sym } A (\text{root } t1) (\text{root } t2) \text{ Two } (\text{transform\_tree}$   
 $t2))) \rangle \mid$   
 $\langle \text{transform\_tree } (\text{Rule } A y) = (\text{Rule } A []) \rangle$

**lemma** *root\_of\_transform\_tree*[intro, simp]:  $\langle \text{root } t = \text{Nt } X \implies \text{root } (\text{transform\_tree } t) = \text{Nt } X \rangle$   
 $\langle \text{proof} \rangle$

**lemma** *transform\_tree\_correct*:  
**assumes**  $\langle \text{parse\_tree } P \text{ } t \wedge \text{fringe } t = w \rangle$   
**shows**  $\langle \text{parse\_tree } P' (\text{transform\_tree } t) \wedge \text{hs } (\text{fringe } (\text{transform\_tree } t)) = w \rangle$   
 $\langle \text{proof} \rangle$

**lemma**  
*transfer\_parse\_tree*:  
**assumes**  $\langle w \in \text{Ders } P \text{ } S \rangle$   
**shows**  $\langle \exists w' \in \text{Ders } P' \text{ } S. w = \text{hs } w' \rangle$   
 $\langle \text{proof} \rangle$

This is essentially  $h(L') \supseteq L$ :

**lemma** *P\_imp\_h\_L'*:  
**assumes**  $\langle w \in \text{Lang } P \text{ } S \rangle$   
**shows**  $\langle \exists w' \in L'. w = h \text{ } w' \rangle$   
 $\langle \text{proof} \rangle$

This lemma is used in the proof of the other direction ( $h(L') \subseteq L$ ):

**lemma** *hom\_ext\_inv*[simp]:  
**assumes**  $\langle \pi \in P \rangle$   
**shows**  $\langle \text{hs } (\text{snd } (\text{transform\_prod } \pi)) = \text{snd } \pi \rangle$   
 $\langle \text{proof} \rangle$

This lemma is essentially the other direction ( $h(L') \subseteq L$ ):

**lemma** *L'\_imp\_h\_P*:  
**assumes**  $\langle w' \in L' \rangle$   
**shows**  $\langle h \text{ } w' \in \text{Lang } P \text{ } S \rangle$   
 $\langle \text{proof} \rangle$

## 8 The Theorem

The constructive version of the Theorem, for a grammar already in CNF:

**lemma** *Chomsky\_Schuetzenberger\_CNF*:

⟨regular (brackets  $\cap$  Reg  $S$ )  
 $\wedge L = h \text{ ‘ } ((brackets \cap Reg S) \cap Dyck\_lang \Gamma)$   
 $\wedge hom\_list (h :: ('n, 't) bracket3 list \Rightarrow 't list)$ ⟩  
 ⟨proof⟩

**end**

Now we want to prove the theorem without assuming that  $P$  is in CNF. Of course any grammar can be converted into CNF, but this requires an infinite type of nonterminals (because the conversion to CNF may need to invent new nonterminals). Therefore we cannot just re-enter *locale\_P*. Now we make all the assumption explicit.

The theorem for any grammar, but only for languages not containing  $\varepsilon$ :

**lemma** *Chomsky\_Schuetzenberger\_not\_empty*:

**fixes**  $P :: \langle ('n :: infinite, 't) Prods \rangle$  **and**  $S :: 'n$   
**defines**  $\langle L \equiv Lang P S - \{\emptyset\} \rangle$   
**assumes** *finiteP*:  $\langle finite P \rangle$   
**shows**  $\langle \exists (R :: ('n, 't) bracket3 list set) h \Gamma. regular R \wedge L = h \text{ ‘ } (R \cap Dyck\_lang \Gamma) \wedge hom\_list h \rangle$   
 ⟨proof⟩

The Chomsky-Schützenberger theorem that we really want to prove:

**theorem** *Chomsky\_Schuetzenberger*:

**fixes**  $P :: \langle ('n :: infinite, 't) Prods \rangle$  **and**  $S :: 'n$   
**defines**  $\langle L \equiv Lang P S \rangle$   
**assumes** *finite*:  $\langle finite P \rangle$   
**shows**  $\langle \exists (R :: ('n, 't) bracket3 list set) h \Gamma. regular R \wedge L = h \text{ ‘ } (R \cap Dyck\_lang \Gamma) \wedge hom\_list h \rangle$   
 ⟨proof⟩

**no\_notation** *the\_hom* (h)

**no\_notation** *the\_hom\_syms* (hs)

**end**

## References

- [1] N. Chomsky and M. Schützenberger. The algebraic theory of context-free languages. In P. Braffort and D. Hirschberg, editors, *Computer Programming and Formal Systems*, volume 26 of *Studies in Logic and the Foundations of Mathematics*, pages 118–161. Elsevier, 1959.

- [2] D. Kozen. *Automata and computability*. Undergraduate texts in computer science. Springer, 1997.