

Greibach Normal Form

Alexander Haberl and Tobias Nipkow and Akihisa Yamada

September 12, 2025

Abstract

This theory formalizes Hopcroft and Ullman’s algorithm [3] to transform a set of productions into Greibach Normal Form (GNF) [2]. We concentrate on the essential property of the GNF: every production starts with a terminal; the tail of a rhs may contain further terminals. The complexity of the algorithm can be exponential.

```
theory Greibach_Normal_Form
imports
  Context_Free_Grammar.Context_Free_Grammar
  Fresh_Identifiers.Fresh_Nat
begin
```

```
declare relpowp.simps(2)[simp del]
```

1 Aux Lemmas

```
lemma Nts_mono:  $G \subseteq H \implies \text{Nts } G \subseteq \text{Nts } H$ 
<proof>
```

```
lemma derivern_prepend:  $R \vdash u \Rightarrow r(n) v \implies R \vdash p @ u \Rightarrow r(n) p @ v$ 
<proof>
```

```
lemma Lang_subset_if_Ders_subset:  $\text{Ders } R \ A \subseteq \text{Ders } R' \ A \implies \text{Lang } R \ A \subseteq \text{Lang } R' \ A$ 
<proof>
```

```
lemma Eps_free_deriven_Nil:
   $\llbracket \text{Eps\_free } R; R \vdash l \Rightarrow (n) \ \square \rrbracket \implies l = \square$ 
<proof>
```

```
lemma nts_syms_empty_iff:  $\text{nts\_syms } w = \{\} \longleftrightarrow (\exists u. w = \text{map } Tm \ u)$ 
<proof>
```

lemma *non_word_has_last_Nt*: $nts_syms\ w \neq \{\} \implies \exists u\ A\ v. w = u @ [Nt\ A]$
 $@\ map\ Tm\ v$
 $\langle proof \rangle$

lemma *nts_syms_rev*: $nts_syms\ (rev\ w) = nts_syms\ w$
 $\langle proof \rangle$

Sentential form that is not a word has a first *Nt*.

lemma *non_word_has_first_Nt*: $nts_syms\ w \neq \{\} \implies \exists u\ A\ v. w = map\ Tm\ u$
 $@\ Nt\ A\ \# v$
 $\langle proof \rangle$

If there exists a derivation from *u* to *v* then there exists one which does not use productions of the form $A \rightarrow A$.

lemma *no_self_loops_derives*: $P \vdash u \Rightarrow l(n)\ v \implies \{p \in P. \neg(\exists A. p = (A, [Nt\ A]))\} \vdash u \Rightarrow l^* v$
 $\langle proof \rangle$

A decomposition of a derivation from a sentential form to a word into multiple derivations that derive words.

lemma *derivern_snoc_Nt_Tms_decomp1*:
 $R \vdash p @ [Nt\ A] \Rightarrow r(n)\ map\ Tm\ q$
 $\implies \exists pt\ At\ w\ k\ m. R \vdash p \Rightarrow (k)\ map\ Tm\ pt \wedge R \vdash w \Rightarrow (m)\ map\ Tm\ At \wedge (A,$
 $w) \in R$
 $\wedge q = pt @ At \wedge n = Suc(k + m)$
 $\langle proof \rangle$

A decomposition of a derivation from a sentential form to a word into multiple derivations that derive words.

lemma *word_decomp1*:
 $R \vdash p @ [Nt\ A] @\ map\ Tm\ ts \Rightarrow (n)\ map\ Tm\ q$
 $\implies \exists pt\ At\ w\ k\ m. R \vdash p \Rightarrow (k)\ map\ Tm\ pt \wedge R \vdash w \Rightarrow (m)\ map\ Tm\ At \wedge (A,$
 $w) \in R$
 $\wedge q = pt @ At @ ts \wedge n = Suc(k + m)$
 $\langle proof \rangle$

Sentential form that derives to terminals and has a *Nt* in it has a derivation that starts with some rule acting on that *Nt*.

lemma *derived_start_sent*:
 $R \vdash u @ Nt\ V\ \# w \Rightarrow (Suc\ n)\ map\ Tm\ x \implies \exists v. (V, v) \in R \wedge R \vdash u @ v @ w$
 $\Rightarrow (n)\ map\ Tm\ x$
 $\langle proof \rangle$

definition *nts_syms_list* :: $('n, 't)syms \Rightarrow 'n\ list \Rightarrow 'n\ list$ **where**
 $nts_syms_list\ sys = foldr\ (\lambda sy\ ns. case\ sy\ of\ Nt\ A \Rightarrow List.insert\ A\ ns\ |\ Tm\ _ \Rightarrow ns)\ sys$

definition *nts_prods_list* :: ('n,'t)prods \Rightarrow 'n list **where**
nts_prods_list ps = foldr ($\lambda(A,sys) ns. List.insert A (nts_syms_list sys ns)$) ps []

lemma *set_nts_syms_list*: $set(nts_syms_list sys ns) = nts_syms sys \cup set ns$
 <proof>

lemma *set_nts_prods_list*: $set(nts_prods_list ps) = nts ps$
 <proof>

lemma *distinct_nts_syms_list*: $distinct(nts_syms_list sys ns) = distinct ns$
 <proof>

lemma *distinct_nts_prods_list*: $distinct(nts_prods_list ps)$
 <proof>

fun *freshs* :: ('a::fresh) set \Rightarrow 'a list \Rightarrow 'a list **where**
freshs X [] = [] |
freshs X (a#as) = (let a' = fresh X a in a' # *freshs* (insert a' X) as)

lemma *length_freshs*: $finite X \implies length(freshs X as) = length as$
 <proof>

lemma *freshs_disj*: $finite X \implies X \cap set(freshs X as) = \{\}$
 <proof>

lemma *freshs_distinct*: $finite X \implies distinct (freshs X as)$
 <proof>

This theory formalizes a method to transform a set of productions into Greibach Normal Form (GNF) [2]. We concentrate on the essential property of the GNF: every production starts with a *Tm*; the tail of a rhs can contain further terminals. This is formalized as *GNF_hd* below. This more liberal definition of GNF is also found elsewhere [1].

The algorithm consists of two phases:

- *solve_tri* converts the productions into a *triangular* form, where Nt *Ai* does not depend on Nts *Ai*, ..., *An*. This involves the elimination of left-recursion and is the heart of the algorithm.
- *expand_tri* expands the triangular form by substituting in: Due to triangular form, *A0* productions satisfy *GNF_hd* and we can substitute them into the heads of the remaining productions. Now all *A1* productions satisfy *GNF_hd*, and we continue until all productions satisfy *GNF_hd*.

This is essentially the algorithm given by Hopcroft and Ullman [3], except that we can drop the conversion to Chomsky Normal Form because of our more liberal *GNF_{hd}*.

2 Function Definitions

Depend on: A depends on B if there is a rule $A \rightarrow B w$:

definition *dep_on* :: $('n, 't) Prods \Rightarrow 'n \Rightarrow 'n \text{ set}$ **where**
dep_on $R A = \{B. \exists w. (A, Nt B \# w) \in R\}$

GNF property: All productions start with a terminal.

definition *GNF_hd* :: $('n, 't) Prods \Rightarrow \text{bool}$ **where**
GNF_hd $R = (\forall (A, w) \in R. \exists t. hd w = Tm t)$

GNF property expressed via *dep_on*:

definition *GNF_hd_dep_on* :: $('n, 't) Prods \Rightarrow \text{bool}$ **where**
GNF_hd_dep_on $R = (\forall A \in Nts R. dep_on R A = \{\})$

abbreviation *lrec_Prods* :: $('n, 't) Prods \Rightarrow 'n \Rightarrow 'n \text{ set} \Rightarrow ('n, 't) Prods$ **where**
lrec_Prods $R A S \equiv \{(A', Bw) \in R. A' = A \wedge (\exists w B. Bw = Nt B \# w \wedge B \in S)\}$

abbreviation *subst_hd* :: $('n, 't) Prods \Rightarrow ('n, 't) Prods \Rightarrow 'n \Rightarrow ('n, 't) Prods$ **where**
subst_hd $R X A \equiv \{(A, v @ w) \mid v w. \exists B. (A, Nt B \# w) \in X \wedge (B, v) \in R\}$

Expand head: Replace all rules $A \rightarrow B w$ where $B \in Ss$ ($Ss =$ solved Nts in *triangular* form) by $A \rightarrow v w$ where $B \rightarrow v$. Starting from the end of Ss .

fun *expand_hd* :: $'n \Rightarrow 'n \text{ list} \Rightarrow ('n, 't) Prods \Rightarrow ('n, 't) Prods$ **where**
expand_hd $A [] R = R$ |
expand_hd $A (S \# Ss) R =$
 (let $R' = \text{expand_hd } A Ss R$;
 $X = \text{lrec_Prods } R' A \{S\}$;
 $Y = \text{subst_hd } R' X A$
 in $R' - X \cup Y$)

lemma *Rhss_code*[*code*]: $Rhss P A = snd \text{ ' } \{Aw \in P. fst Aw = A\}$
 <proof>

declare *expand_hd.simps(1)*[*code*]

lemma *expand_hd_Cons_code*[*code*]: $\text{expand_hd } A (S \# Ss) R =$
 (let $R' = \text{expand_hd } A Ss R$;
 $X = \{w \in Rhss R' A. w \neq [] \wedge hd w = Nt S\}$;
 $Y = (\bigcup (B, v) \in R'. \bigcup w \in X. \text{if } hd w \neq Nt B \text{ then } \{\} \text{ else } \{(A, v @ tl w)\})$
 in $R' - (\{A\} \times X) \cup Y$)
 <proof>

Remove left-recursions: Remove left-recursive rules $A \rightarrow A w$:

definition $rm_lrec :: 'n \Rightarrow ('n, 't) Prods \Rightarrow ('n, 't) Prods$ **where**
 $rm_lrec\ A\ R = R - \{(A, Nt\ A\ \# \ v) \mid v. \text{True}\}$

lemma $rm_lrec_code[code]$:

$rm_lrec\ A\ R = \{Aw \in R. \text{let } (A', w) = Aw \text{ in } A' \neq A \vee w = [] \vee hd\ w \neq Nt\ A\}$
 $\langle proof \rangle$

Make right-recursion of left-recursion: Conversion from left-recursion to right-recursion: Split A -rules into $A \rightarrow u$ and $A \rightarrow A\ v$. Keep $A \rightarrow u$ but replace $A \rightarrow A\ v$ by $A \rightarrow u\ A', A' \rightarrow v, A' \rightarrow v\ A'$.

The then part of the if statement is only an optimisation, so that we do not introduce the $A \rightarrow u\ A'$ rules if we do not introduce any A' rules, but the function also works, if we always enter the else part.

definition $rrec_of_lrec :: 'n \Rightarrow 'n \Rightarrow ('n, 't) Prods \Rightarrow ('n, 't) Prods$ **where**
 $rrec_of_lrec\ A\ A'\ R =$

$(\text{let } V = \{v. (A, Nt\ A\ \# \ v) \in R \wedge v \neq []\};$
 $U = \{u. (A, u) \in R \wedge \neg(\exists v. u = Nt\ A\ \# \ v)\}$
 $\text{in if } V = \{\} \text{ then } R - \{(A, [Nt\ A])\} \text{ else } (\{A\} \times U) \cup (\bigcup_{u \in U}. \{(A, u@[Nt\ A'])\}) \cup (\{A'\} \times V) \cup (\bigcup_{v \in V}. \{(A', v@[Nt\ A'])\}))$

lemma $rrec_of_lrec_code[code]$: $rrec_of_lrec\ A\ A'\ R =$

$(\text{let } RA = Rhss\ R\ A;$
 $V = tl\ ' \{w \in RA. w \neq [] \wedge hd\ w = Nt\ A \wedge tl\ w \neq []\};$
 $U = \{u \in RA. u = [] \vee hd\ u \neq Nt\ A\}$
 $\text{in if } V = \{\} \text{ then } R - \{(A, [Nt\ A])\} \text{ else } (\{A\} \times U) \cup (\bigcup_{u \in U}. \{(A, u@[Nt\ A'])\}) \cup (\{A'\} \times V) \cup (\bigcup_{v \in V}. \{(A', v@[Nt\ A'])\}))$
 $\langle proof \rangle$

Solve left-recursions: Solves the left-recursion of $Nt\ A$ by replacing it with a right-recursion on a fresh $Nt\ A'$. The fresh $Nt\ A'$ is also given as a parameter.

definition $solve_lrec :: 'n \Rightarrow 'n \Rightarrow ('n, 't) Prods \Rightarrow ('n, 't) Prods$ **where**
 $solve_lrec\ A\ A'\ R = rm_lrec\ A\ R \cup rrec_of_lrec\ A\ A'\ R$

lemmas $solve_lrec_defs = solve_lrec_def\ rm_lrec_def\ rrec_of_lrec_def\ Let_def\ Nts_def$

Solve triangular: Put R into triangular form wrt As (using the new $Nts\ As'$). In each step $A\ \# \ As$, first the remaining Nts in As are solved, then A is solved. This should mean that in the result of the outermost $expand_hd\ A\ As$, A only depends on A . Then the A rules in the result of $solve_lrec\ A\ A'$ are already in GNF. More precisely: the result should be in *triangular* form.

fun $solve_tri :: 'a\ list \Rightarrow 'a\ list \Rightarrow ('a, 'b)\ Prods \Rightarrow ('a, 'b)\ Prods$ **where**

$solve_tri\ []\ _ \ R = R$
 $solve_tri\ (A\ \# \ As)\ (A'\ \# \ As')\ R = solve_lrec\ A\ A'\ (expand_hd\ A\ As\ (solve_tri\ As\ As'\ R))$

Triangular form wrt $[A1, \dots, An]$ means that Ai must not depend on Ai , \dots , An . In particular: $A0$ does not depend on any Ai , its rules are already

in GNF. Therefore one can convert a *triangular* form into GNF by backwards substitution: The rules for A_i are used to expand the heads of all $A(i+1), \dots, A_n$ rules, starting with A_0 .

```
fun triangular :: 'n list  $\Rightarrow$  ('n  $\times$  ('n, 't) sym list) set  $\Rightarrow$  bool where
  triangular [] R = True |
  triangular (A#As) R = (dep_on R A  $\cap$  ({A}  $\cup$  set As) = {})  $\wedge$  triangular As R
```

Remove self loops: Removes all productions of the form $A \rightarrow A$.

```
definition rm_self_loops :: ('n,'t) Prods  $\Rightarrow$  ('n,'t) Prods where
  rm_self_loops P = P - {x  $\in$  P.  $\exists$  A. x = (A, [Nt A])}
```

Expand triangular: Expands all head-Nts of productions with a Lhs in As (*triangular* (*rev* As)). In each step $A\#As$ first all Nts in As are expanded, then every rule $A \rightarrow B w$ is expanded if $B \in$ set As . If the productions were in *triangular* form wrt *rev* As then A_i only depends on $A(i+1), \dots, A_n$ which have already been expanded in the first part of the step and are in GNF. Then the all A -productions are also in GNF after expansion.

```
fun expand_tri :: 'n list  $\Rightarrow$  ('n,'t)Prods  $\Rightarrow$  ('n,'t)Prods where
  expand_tri [] R = R |
  expand_tri (A#As) R =
    (let R' = expand_tri As R;
      X = lrec_Prods R' A (set As);
      Y = subst_hd R' X A
    in R' - X  $\cup$  Y)
```

```
declare expand_tri.simps(1)[code]
```

```
lemma expand_tri_Cons_code[code]: expand_tri (S#Ss) R =
  (let R' = expand_tri Ss R;
    X = {w  $\in$  Rhss R' S. w  $\neq$  []  $\wedge$  hd w  $\in$  Nt ' (set Ss)};
    Y = ( $\bigcup$  (B,v)  $\in$  R'.  $\bigcup$  w  $\in$  X. if hd w  $\neq$  Nt B then {} else {(S,v @ tl w)})
  in R' - ({S}  $\times$  X)  $\cup$  Y)
<proof>
```

The main function *gnf_hd* converts into *GNF_hd*:

```
definition gnf_hd :: ('n::fresh,'t)prods  $\Rightarrow$  ('n,'t)Prods where
  gnf_hd ps =
    (let As = nts_prods_list ps;
      As' = freshs (set As) As
    in expand_tri (As' @ rev As) (solve_tri As As' (set ps)))
```

3 Some Basic Lemmas

3.1 Eps_free preservation

```
lemma Eps_free_expand_hd: Eps_free R  $\Longrightarrow$  Eps_free (expand_hd A Ss R)
  <proof>
```

lemma *Eps_free_solve_lrec*: $Eps_free\ R \implies Eps_free\ (solve_lrec\ A\ A'\ R)$
 ⟨proof⟩

lemma *Eps_free_solve_tri*: $Eps_free\ R \implies length\ As \leq length\ As' \implies Eps_free\ (solve_tri\ As\ As'\ R)$
 ⟨proof⟩

lemma *Eps_free_expand_tri*: $Eps_free\ R \implies Eps_free\ (expand_tri\ As\ R)$
 ⟨proof⟩

3.2 Lemmas about *Nts* and *dep_on*

lemma *dep_on_Un[simp]*: $dep_on\ (R \cup S)\ A = dep_on\ R\ A \cup dep_on\ S\ A$
 ⟨proof⟩

lemma *expand_hd_preserves_neq*: $B \neq A \implies (B, w) \in expand_hd\ A\ Ss\ R \longleftrightarrow (B, w) \in R$
 ⟨proof⟩

Let R be epsilon-free and in *triangular* form wrt Bs . After *expand_hd* $A\ Bs\ R$, A depends only on what A depended on before or what one of the $B \in Bs$ depends on, but A does not depend on the Bs :

lemma *dep_on_expand_hd*:
 $\llbracket Eps_free\ R; triangular\ Bs\ R; distinct\ Bs; A \notin set\ Bs \rrbracket$
 $\implies dep_on\ (expand_hd\ A\ Bs\ R)\ A \subseteq (dep_on\ R\ A \cup (\bigcup_{B \in set\ Bs} dep_on\ R\ B)) - set\ Bs$
 ⟨proof⟩

lemma *dep_on_subs_Nts*: $dep_on\ R\ A \subseteq Nts\ R$
 ⟨proof⟩

lemma *Nts_expand_hd_sub*: $Nts\ (expand_hd\ A\ As\ R) \subseteq Nts\ R$
 ⟨proof⟩

lemma *Nts_solve_lrec_sub*: $Nts\ (solve_lrec\ A\ A'\ R) \subseteq Nts\ R \cup \{A'\}$
 ⟨proof⟩

lemma *Nts_solve_tri_sub*: $length\ As \leq length\ As' \implies Nts\ (solve_tri\ As\ As'\ R) \subseteq Nts\ R \cup set\ As'$
 ⟨proof⟩

3.3 Lemmas about *triangular*

lemma *tri_Snoc_impl_tri*: $triangular\ (As\ @\ [A])\ R \implies triangular\ As\ R$
 ⟨proof⟩

If two parts of the productions are *triangular* and no *Nts* from the first part depend on ones of the second they are also *triangular* when put together.

lemma *triangular_append*:

$$\llbracket \text{triangular } As \ R; \text{triangular } Bs \ R; \forall A \in \text{set } As. \text{dep_on } R \ A \cap \text{set } Bs = \{\} \rrbracket$$

$$\implies \text{triangular } (As @ Bs) \ R$$

$$\langle \text{proof} \rangle$$

4 Function *solve_tri*: Remove Left-Recursion and Convert into Triangular Form

4.1 Basic Lemmas

Mostly about rule inclusions in *solve_lrec*.

lemma *solve_lrec_rule_simp1*: $A \neq B \implies A \neq B' \implies (A, w) \in \text{solve_lrec } B \ B'$
 $R \longleftrightarrow (A, w) \in R$
 $\langle \text{proof} \rangle$

lemma *solve_lrec_rule_simp3*: $A \neq A' \implies A' \notin \text{Nts } R \implies \text{Eps_free } R$
 $\implies (A, [\text{Nt } A']) \notin \text{solve_lrec } A \ A' \ R$
 $\langle \text{proof} \rangle$

lemma *solve_lrec_rule_simp7*: $A' \neq A \implies A' \notin \text{Nts } R \implies (A', \text{Nt } A' \ \# \ w) \notin \text{solve_lrec } A \ A' \ R$
 $\langle \text{proof} \rangle$

lemma *solve_lrec_rule_simp8*: $A' \notin \text{Nts } R \implies B \neq A' \implies B \neq A$
 $\implies (B, \text{Nt } A' \ \# \ w) \notin \text{solve_lrec } A \ A' \ R$
 $\langle \text{proof} \rangle$

lemma *dep_on_expand_hd_simp2*: $B \neq A \implies \text{dep_on } (\text{expand_hd } A \ As \ R) \ B = \text{dep_on } R \ B$
 $\langle \text{proof} \rangle$

lemma *dep_on_solve_lrec_simp2*: $A \neq B \implies A' \neq B \implies \text{dep_on } (\text{solve_lrec } A \ A' \ R) \ B = \text{dep_on } R \ B$
 $\langle \text{proof} \rangle$

4.2 Triangular Form

expand_hd preserves *triangular*, if it does not expand a *Nt* considered in *triangular*.

lemma *triangular_expand_hd*: $\llbracket A \notin \text{set } As; \text{triangular } As \ R \rrbracket \implies \text{triangular } As \ (\text{expand_hd } A \ Bs \ R)$
 $\langle \text{proof} \rangle$

Solving a *Nt* not considered by *triangular* preserves the *triangular* property.

lemma *triangular_solve_lrec*: $\llbracket A \notin \text{set } As; A' \notin \text{set } As; \text{triangular } As \ R \rrbracket$
 $\implies \text{triangular } As \ (\text{solve_lrec } A \ A' \ R)$
 $\langle \text{proof} \rangle$

Solving more Nts does not remove the *triangular* property of previously solved Nts.

lemma *part_triangular_induct_step*:

$\llbracket \text{Eps_free } R; \text{ distinct } ((A \# As) @ (A' \# As')); \text{ triangular } As \text{ (solve_tri } As \text{ As' } R) \rrbracket$
 $\implies \text{ triangular } As \text{ (solve_tri } (A \# As) \text{ (A' \# As')} \text{ } R)$
 $\langle \text{proof} \rangle$

Couple of small lemmas about *dep_on* and the solving of left-recursion.

lemma *rm_lrec_rem_own_dep*: $A \notin \text{dep_on } (\text{rm_lrec } A \text{ } R) \text{ } A$

$\langle \text{proof} \rangle$

lemma *rrec_of_lrec_has_no_own_dep*: $A \neq A' \implies A \notin \text{dep_on } (\text{rrec_of_lrec } A \text{ } A' \text{ } R) \text{ } A$

$\langle \text{proof} \rangle$

lemma *solve_lrec_no_own_dep*: $A \neq A' \implies A \notin \text{dep_on } (\text{solve_lrec } A \text{ } A' \text{ } R) \text{ } A$

$\langle \text{proof} \rangle$

lemma *solve_lrec_no_new_own_dep*: $A \neq A' \implies A' \notin \text{Nts } R \implies A' \notin \text{dep_on } (\text{solve_lrec } A \text{ } A' \text{ } R) \text{ } A'$

$\langle \text{proof} \rangle$

lemma *dep_on_rm_lrec_simp*: $\text{dep_on } (\text{rm_lrec } A \text{ } R) \text{ } A = \text{dep_on } R \text{ } A - \{A\}$

$\langle \text{proof} \rangle$

lemma *dep_on_rrec_of_lrec_simp*:

$\text{Eps_free } R \implies A \neq A' \implies \text{dep_on } (\text{rrec_of_lrec } A \text{ } A' \text{ } R) \text{ } A = \text{dep_on } R \text{ } A - \{A\}$

$\langle \text{proof} \rangle$

lemma *dep_on_solve_lrec_simp*:

$\llbracket \text{Eps_free } R; A \neq A' \rrbracket \implies \text{dep_on } (\text{solve_lrec } A \text{ } A' \text{ } R) \text{ } A = \text{dep_on } R \text{ } A - \{A\}$

$\langle \text{proof} \rangle$

lemma *dep_on_solve_tri_simp*: $B \notin \text{set } As \implies B \notin \text{set } As' \implies \text{length } As \leq \text{length } As'$

$\implies \text{dep_on } (\text{solve_tri } As \text{ As' } R) \text{ } B = \text{dep_on } R \text{ } B$

$\langle \text{proof} \rangle$

Induction step for showing that *solve_tri* removes dependencies of previously solved Nts.

lemma *triangular_dep_on_induct_step*:

assumes $\text{Eps_free } R \text{ length } As \leq \text{length } As' \text{ distinct } ((A \# As) @ (A' \# As')) \text{ triangular } As \text{ (solve_tri } As \text{ As' } R)$

shows $\text{dep_on } (\text{solve_tri } (A \# As) \text{ (A' \# As')} \text{ } R) \text{ } A \cap (\{A\} \cup \text{set } As) = \{\}$

$\langle \text{proof} \rangle$

theorem *triangular_solve_tri*: $\llbracket \text{Eps_free } R; \text{ length } As \leq \text{length } As'; \text{ distinct } (As @ As') \rrbracket$

$\implies \text{triangular } As \text{ (solve_tri } As \text{ } As' \text{ } R)$
 $\langle \text{proof} \rangle$

lemma *dep_on_solve_tri_Nts_R*:
 $\llbracket \text{Eps_free } R; B \in \text{set } As; \text{distinct } (As @ As'); \text{set } As' \cap \text{Nts } R = \{\}; \text{length } As \leq \text{length } As' \rrbracket$
 $\implies \text{dep_on } (\text{solve_tri } As \text{ } As' \text{ } R) \text{ } B \subseteq \text{Nts } R$
 $\langle \text{proof} \rangle$

lemma *triangular_unused_Nts*: $\text{set } As \cap \text{Nts } R = \{\} \implies \text{triangular } As \text{ } R$
 $\langle \text{proof} \rangle$

The newly added Nts in *solve_lrec* are in *triangular* form wrt *rev As'*.

lemma *triangular_rev_As'_solve_tri*:
 $\llbracket \text{set } As' \cap \text{Nts } R = \{\}; \text{distinct } (As @ As'); \text{length } As \leq \text{length } As' \rrbracket$
 $\implies \text{triangular } (\text{rev } As') \text{ (solve_tri } As \text{ } As' \text{ } R)$
 $\langle \text{proof} \rangle$

The entire set of productions is in *triangular* form after *solve_tri* wrt $As @ (\text{rev } As')$.

theorem *triangular_As_As'_solve_tri*:
assumes *Eps_free* *R* $\text{length } As \leq \text{length } As'$ *distinct* $(As @ As')$ $\text{Nts } R \subseteq \text{set } As$
shows *triangular* $(As @ (\text{rev } As')) \text{ (solve_tri } As \text{ } As' \text{ } R)$
 $\langle \text{proof} \rangle$

4.3 *solve_lrec* Preserves Language

4.3.1 $\text{Lang } R \text{ } A \subseteq \text{Lang } (\text{solve_lrec } B \text{ } B' \text{ } R) \text{ } A$

If there exists a derivation from *u* to *v* then there exists one which does not use productions of the form $A \rightarrow A$.

lemma *rm_self_loops_derivels*: **assumes** $P \vdash u \Rightarrow l(n) \text{ } v$ **shows** *rm_self_loops* $P \vdash u \Rightarrow l^* \text{ } v$
 $\langle \text{proof} \rangle$

Restricted to productions with one lhs (*A*), and no $A \rightarrow A$ productions if there is a derivation from *u* to $A \# v$ then *u* must start with Nt *A*.

lemma *lrec_lemma1*:
assumes $S = \{x. (\exists v. x = (A, v) \wedge x \in R)\}$ *rm_self_loops* $S \vdash u \Rightarrow l(n) \text{ } Nt \text{ } A \# v$
shows $\exists u'. u = Nt \text{ } A \# u'$
 $\langle \text{proof} \rangle$

Restricted to productions with one lhs (*A*), and no $A \rightarrow A$ productions if there is a derivation from *u* to $A \# v$ then *u* must start with Nt *A* and there exists a prefix of $A \# v$ s.t. a left-derivation from $[Nt \text{ } A]$ to that prefix exists.

lemma *lrec_lemma2*:

assumes $S = \{x. (\exists v. x = (A, v) \wedge x \in R)\} \text{ Eps_free } R$
shows $\text{rm_self_loops } S \vdash u \Rightarrow l(n) \text{ Nt } A \# v \Longrightarrow$
 $\exists u' v'. u = \text{Nt } A \# u' \wedge v = v' @ u' \wedge (\text{rm_self_loops } S) \vdash [\text{Nt } A] \Rightarrow l(n) \text{ Nt } A \# v'$
 $\langle \text{proof} \rangle$

Restricted to productions with one lhs (A), and no $A \rightarrow A$ productions if there is a left-derivation from $[\text{Nt } A]$ to $A \# u$ then there exists a derivation from $[\text{Nt } A]$ to $u @ [\text{Nt } A]$ and if $u \neq []$ also to u in $\text{solve_lrec } A \text{ A}' R$.

lemma *lrec_lemma3*:

assumes $S = \{x. (\exists v. x = (A, v) \wedge x \in R)\} \text{ Eps_free } R$
shows $\text{rm_self_loops } S \vdash [\text{Nt } A] \Rightarrow l(n) \text{ Nt } A \# u$
 $\Longrightarrow \text{solve_lrec } A \text{ A}' S \vdash [\text{Nt } A] \Rightarrow (n) u @ [\text{Nt } A] \wedge$
 $(u \neq [] \longrightarrow \text{solve_lrec } A \text{ A}' S \vdash [\text{Nt } A] \Rightarrow (n) u)$
 $\langle \text{proof} \rangle$

A left derivation from p ($\text{hd } p = \text{Nt } A$) to q ($\text{hd } q \neq \text{Nt } A$) can be split into a left-recursive part, only using left-recursive productions $A \rightarrow A \# w$, one left derivation step consuming $\text{Nt } A$ using some rule $A \rightarrow B \# v$ where $B \neq \text{Nt } A$ and a left-derivation comprising the rest of the derivation.

lemma *lrec_decomp*:

assumes $S = \{x. (\exists v. x = (A, v) \wedge x \in R)\} \text{ Eps_free } R$
shows $\llbracket \text{hd } p = \text{Nt } A; \text{hd } q \neq \text{Nt } A; R \vdash p \Rightarrow l(n) q \rrbracket$
 $\Longrightarrow \exists u w m k. S \vdash p \Rightarrow l(m) \text{ Nt } A \# u \wedge S \vdash \text{Nt } A \# u \Rightarrow l w \wedge \text{hd } w \neq \text{Nt } A \wedge$
 $R \vdash w \Rightarrow l(k) q \wedge n = m + k + 1$
 $\langle \text{proof} \rangle$

Every derivation resulting in a word has a derivation in $\text{solve_lrec } B \text{ B}' R$.

lemma *tm_derive_impl_solve_lrec_derive*:

assumes $\text{Eps_free } R \text{ B} \neq \text{B}' \text{ B}' \notin \text{Nts } R$
shows $\llbracket p \neq []; R \vdash p \Rightarrow (n) \text{ map } \text{Tm } q \rrbracket \Longrightarrow \text{solve_lrec } B \text{ B}' R \vdash p \Rightarrow * \text{ map } \text{Tm } q$
 $\langle \text{proof} \rangle$

corollary *Lang_incl_Lang_solve_lrec*:

$\llbracket \text{Eps_free } R; B \neq B'; B' \notin \text{Nts } R \rrbracket \Longrightarrow \text{Lang } R \text{ A} \subseteq \text{Lang } (\text{solve_lrec } B \text{ B}' R)$
 A
 $\langle \text{proof} \rangle$

4.3.2 $\text{Lang } (\text{solve_lrec } B \text{ B}' R) \text{ A} \subseteq \text{Lang } R \text{ A}$

Restricted to right-recursive productions of one Nt ($A' \rightarrow w @ [\text{Nt } A']$) if there is a right-derivation from u to $v @ [\text{Nt } A']$ then u ends in $\text{Nt } A'$.

lemma *rrec_lemma1*:

assumes $S = \{x. \exists v. x = (A', v @ [\text{Nt } A']) \wedge x \in \text{solve_lrec } A \text{ A}' R\} S \vdash u$
 $\Rightarrow r(n) v @ [\text{Nt } A']$

shows $\exists u'. u = u' @ [Nt A']$
 $\langle proof \rangle$

solve_lrec does not add productions of the form $A' \rightarrow Nt A'$.

lemma *solve_lrec_no_self_loop*: $Eps_free R \implies A' \notin Nts R \implies (A', [Nt A']) \notin solve_lrec A A' R$
 $\langle proof \rangle$

Restricted to right-recursive productions of one Nt ($A' \rightarrow w @ [Nt A']$) if there is a right-derivation from u to $v @ [Nt A']$ then u ends in Nt A' and there exists a suffix of $v @ [Nt A']$ s.t. there is a right-derivation from $[Nt A']$ to that suffix.

lemma *rrec_lemma2*:
assumes $S = \{x. (\exists v. x = (A', v @ [Nt A']) \wedge x \in solve_lrec A A' R)\}$ $Eps_free R$ $A' \notin Nts R$
shows $S \vdash u \Rightarrow r(n) v @ [Nt A']$
 $\implies \exists u' v'. u = u' @ [Nt A'] \wedge v = u' @ v' \wedge S \vdash [Nt A'] \Rightarrow r(n) v' @ [Nt A']$
 $\langle proof \rangle$

Restricted to right-recursive productions of one Nt ($A' \rightarrow w @ [Nt A']$) if there is a restricted right-derivation in *solve_lrec* from $[Nt A']$ to $u @ [Nt A']$ then there exists a derivation in R from $[Nt A]$ to $A \# u$.

lemma *rrec_lemma3*:
assumes $S = \{x. (\exists v. x = (A', v @ [Nt A']) \wedge x \in solve_lrec A A' R)\}$ $Eps_free R$
 $A' \notin Nts R$ $A \neq A'$
shows $S \vdash [Nt A'] \Rightarrow r(n) u @ [Nt A'] \implies R \vdash [Nt A] \Rightarrow (n) Nt A \# u$
 $\langle proof \rangle$

A right derivation from $p@[Nt A']$ to q (*last* $q \neq Nt A'$) can be split into a right-recursive part, only using right-recursive productions with Nt A' , one right derivation step consuming Nt A' using some rule $A' \rightarrow as@[Nt B]$ where $Nt B \neq Nt A'$ and a right-derivation comprising the rest of the derivation.

lemma *rrec_decomp*:
assumes $S = \{x. (\exists v. x = (A', v @ [Nt A']) \wedge x \in solve_lrec A A' R)\}$ $Eps_free R$
 $A \neq A'$ $A' A' \notin Nts R$
shows $\llbracket A' \notin nts_syms p; last\ q \neq Nt A'; solve_lrec A A' R \vdash p @ [Nt A'] \Rightarrow r(n) q \rrbracket$
 $\implies \exists u\ w\ m\ k. S \vdash p @ [Nt A'] \Rightarrow r(m) u @ [Nt A']$
 $\wedge solve_lrec A A' R \vdash u @ [Nt A'] \Rightarrow r\ w \wedge A' \notin nts_syms w$
 $\wedge solve_lrec A A' R \vdash w \Rightarrow r(k) q \wedge n = m + k + 1$
 $\langle proof \rangle$

Every word derived by *solve_lrec* $B B' R$ can be derived by R .

lemma *tm_solve_lrec_derive_impl_derive*:

assumes $Eps_free\ R\ B \neq B'\ B' \notin Nts\ R$
shows $\llbracket p \neq []; B' \notin nts_syms\ p; (solve_lrec\ B\ B'\ R) \vdash p \Rightarrow (n)\ map\ Tm\ q \rrbracket \implies$
 $R \vdash p \Rightarrow^* map\ Tm\ q$
 $\langle proof \rangle$

corollary $Lang_solve_lrec_incl_Lang$:
assumes $Eps_free\ R\ B \neq B'\ B' \notin Nts\ R\ A \neq B'$
shows $Lang\ (solve_lrec\ B\ B'\ R)\ A \subseteq Lang\ R\ A$
 $\langle proof \rangle$

corollary $solve_lrec_Lang$:
 $\llbracket Eps_free\ R; B \neq B'; B' \notin Nts\ R; A \neq B' \rrbracket \implies Lang\ (solve_lrec\ B\ B'\ R)\ A =$
 $Lang\ R\ A$
 $\langle proof \rangle$

4.4 $expand_hd$ Preserves Language

Every rhs of an $expand_hd\ R$ production is derivable by R .

lemma $expand_hd_is_deriveable$: $(A, w) \in expand_hd\ B\ As\ R \implies R \vdash [Nt\ A]$
 $\Rightarrow^* w$
 $\langle proof \rangle$

lemma $expand_hd_incl1$: $Lang\ (expand_hd\ B\ As\ R)\ A \subseteq Lang\ R\ A$
 $\langle proof \rangle$

This lemma expects a set of quadruples $(A, a1, B, a2)$. Each quadruple encodes a specific Nt in a specific rule $A \rightarrow a1 @ Nt\ B \# a2$ (this encodes Nt B) which should be expanded, by replacing the Nt with every rule for that Nt and then removing the original rule. This expansion contains the original productions Language.

lemma $exp_includes_Lang$:
assumes S_props : $\forall x \in S. \exists A\ a1\ B\ a2. x = (A, a1, B, a2) \wedge (A, a1 @ Nt\ B \# a2) \in R$
shows $Lang\ R\ A$
 $\subseteq Lang\ (R - \{x. \exists A\ a1\ B\ a2. x = (A, a1 @ Nt\ B \# a2) \wedge (A, a1, B, a2) \in S\})$
 $\cup \{x. \exists A\ v\ a1\ a2\ B. x = (A, a1 @ v @ a2) \wedge (A, a1, B, a2) \in S \wedge (B, v) \in R\})\ A$
 $\langle proof \rangle$

lemma $expand_hd_incl2$: $Lang\ (expand_hd\ B\ As\ R)\ A \supseteq Lang\ R\ A$
 $\langle proof \rangle$

theorem $expand_hd_Lang$: $Lang\ (expand_hd\ B\ As\ R)\ A = Lang\ R\ A$
 $\langle proof \rangle$

4.5 *solve_tri* Preserves Language

lemma *solve_tri_Lang*:

$\llbracket \text{Eps_free } R; \text{length } As \leq \text{length } As'; \text{distinct}(As @ As'); Nts\ R \cap \text{set } As' = \{\}; A \notin \text{set } As' \rrbracket$
 $\implies \text{Lang } (\text{solve_tri } As\ As'\ R) A = \text{Lang } R A$
 $\langle \text{proof} \rangle$

5 Function *expand_hd*: Convert Triangular Form into GNF

5.1 *expand_hd*: Result is in *GNF_hd*

lemma *dep_on_helper*: $\text{dep_on } R\ A = \{\} \implies (A, w) \in R \implies w = [] \vee (\exists T\ wt. w = Tm\ T\ \# wt)$
 $\langle \text{proof} \rangle$

lemma *GNF_hd_iff_dep_on*:

assumes *Eps_free* *R*
shows $\text{GNF_hd } R \longleftrightarrow (\forall A \in Nts\ R. \text{dep_on } R\ A = \{\})$ (**is** $?L=?R$)
 $\langle \text{proof} \rangle$

lemma *helper_expand_tri1*: $A \notin \text{set } As \implies (A, w) \in \text{expand_tri } As\ R \implies (A, w) \in R$
 $\langle \text{proof} \rangle$

If none of the expanded Nts depend on *A* then any rule depending on *A* in *expand_tri* *As* *R* must already have been in *R*.

lemma *helper_expand_tri2*:

$\llbracket \text{Eps_free } R; A \notin \text{set } As; \forall C \in \text{set } As. A \notin (\text{dep_on } R\ C); B \neq A; (B, Nt\ A\ \# w) \in \text{expand_tri } As\ R \rrbracket$
 $\implies (B, Nt\ A\ \# w) \in R$
 $\langle \text{proof} \rangle$

In a triangular form no Nts depend on the last Nt in the list.

lemma *triangular_snoc_dep_on*: $\text{triangular } (As@[A])\ R \implies \forall C \in \text{set } As. A \notin (\text{dep_on } R\ C)$
 $\langle \text{proof} \rangle$

lemma *triangular_helper1*: $\text{triangular } As\ R \implies A \in \text{set } As \implies A \notin \text{dep_on } R\ A$
 $\langle \text{proof} \rangle$

lemma *dep_on_expand_tri*:

$\llbracket \text{Eps_free } R; \text{triangular } (\text{rev } As)\ R; \text{distinct } As; A \in \text{set } As \rrbracket$
 $\implies \text{dep_on } (\text{expand_tri } As\ R) A \cap \text{set } As = \{\}$
 $\langle \text{proof} \rangle$

Interlude: *Nts* of *expand_tri*:

lemma *Lhss_expand_tri*: $Lhss (expand_tri As R) \subseteq Lhss R$
 $\langle proof \rangle$

lemma *Rhs_Nts_expand_tri*: $Rhs_Nts (expand_tri As R) \subseteq Rhs_Nts R$
 $\langle proof \rangle$

lemma *Nts_expand_tri*: $Nts (expand_tri As R) \subseteq Nts R$
 $\langle proof \rangle$

If the entire *triangular* form is expanded, the result is in GNF:

theorem *GNF_hd_expand_tri*:
assumes *Eps_free* *R triangular* (*rev As*) *R distinct As Nts R* $\subseteq set As$
shows *GNF_hd* (*expand_tri As R*)
 $\langle proof \rangle$

Any set of productions can be transformed into GNF via *expand_tri* (*solve_tri*).

theorem *GNF_of_R*:
assumes *assms*: *Eps_free R distinct (As @ As') Nts R* $\subseteq set As$ *length As* \leq
length As'
shows *GNF_hd* (*expand_tri (As' @ rev As) (solve_tri As As' R)*)
 $\langle proof \rangle$

5.2 *expand_tri* Preserves Language

Similar to the proof of Language equivalence of *expand_hd*.

All productions in *expand_tri As R* are derivable by *R*.

lemma *expand_tri_prods_deirvable*: $(B, bs) \in expand_tri As R \implies R \vdash [Nt B] \Rightarrow^* bs$
 $\langle proof \rangle$

Language Preservation:

lemma *expand_tri_Lang*: $Lang (expand_tri As R) A = Lang R A$
 $\langle proof \rangle$

6 Function *gnf_hd*: Conversion to *GNF_hd*

All epsilon-free grammars can be put into GNF while preserving their language.

Putting the productions into GNF via *expand_tri* (*solve_tri*) preserves the language.

lemma *GNF_of_R_Lang*:
assumes *Eps_free R length As* \leq *length As'* *distinct (As @ As') Nts R* $\cap set As' = \{\}$ *A* $\notin set As'$
shows $Lang (expand_tri (As' @ rev As) (solve_tri As As' R)) A = Lang R A$
 $\langle proof \rangle$

Any epsilon-free Grammar can be brought into GNF.

theorem *GNF_hd_gnf_hd*: $\text{eps_free } ps \implies \text{GNF_hd } (\text{gnf_hd } ps)$
 $\langle \text{proof} \rangle$

lemma *distinct_app_freshs*: $\llbracket As = \text{nts_prods_list } ps; As' = \text{freshs } (\text{set } As) \ As \rrbracket \implies$
 $\text{distinct } (As @ As')$
 $\langle \text{proof} \rangle$

gnf_hd preserves the language:

theorem *Lang_gnf_hd*: $\llbracket \text{eps_free } ps; A \in \text{nts } ps \rrbracket \implies \text{Lang } (\text{gnf_hd } ps) \ A =$
 $\text{lang } ps \ A$
 $\langle \text{proof} \rangle$

Two simple examples:

lemma *gnf_hd* $[(0, [\text{Nt}(0::\text{nat}), \text{Tm } (1::\text{int})])] = \{(1, [\text{Tm } 1]), (1, [\text{Tm } 1, \text{Nt } 1])\}$
 $\langle \text{proof} \rangle$

lemma *gnf_hd* $[(0, [\text{Nt}(0::\text{nat}), \text{Tm } (1::\text{int})]), (0, [\text{Tm } 2])] =$
 $\{(0, [\text{Tm } 2, \text{Nt } 1]), (0, [\text{Tm } 2]), (1, [\text{Tm } 1, \text{Nt } 1]), (1, [\text{Tm } 1])\}$
 $\langle \text{proof} \rangle$

Example 4.10 [3]: *P0* is the input; *P1* is the result after Step 1; *P3* is the result after Step 2 and 3.

lemma

let
 $P0 =$
 $[(1::\text{int}, [\text{Nt } 2, \text{Nt } 3]), (2, [\text{Nt } 3, \text{Nt } 1]), (2, [\text{Tm } (1::\text{int})]), (3, [\text{Nt } 1, \text{Nt } 2]),$
 $(3, [\text{Tm } 0])];$
 $P1 =$
 $[(1, [\text{Nt } 2, \text{Nt } 3]), (2, [\text{Nt } 3, \text{Nt } 1]), (2, [\text{Tm } 1]),$
 $(3, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4]), (3, [\text{Tm } 0, \text{Nt } 4]), (3, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2]),$
 $(3, [\text{Tm } 0]),$
 $(4, [\text{Nt } 1, \text{Nt } 3, \text{Nt } 2]), (4, [\text{Nt } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4])];$
 $P2 =$
 $[(1, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4, \text{Nt } 1, \text{Nt } 3]), (1, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 1, \text{Nt } 3]),$
 $(1, [\text{Tm } 0, \text{Nt } 4, \text{Nt } 1, \text{Nt } 3]), (1, [\text{Tm } 0, \text{Nt } 1, \text{Nt } 3]), (1, [\text{Tm } 1, \text{Nt } 3]),$
 $(2, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4, \text{Nt } 1]), (2, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 1]),$
 $(2, [\text{Tm } 0, \text{Nt } 4, \text{Nt } 1]), (2, [\text{Tm } 0, \text{Nt } 1]), (2, [\text{Tm } 1]),$
 $(3, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4]), (3, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2]),$
 $(3, [\text{Tm } 0, \text{Nt } 4]), (3, [\text{Tm } 0]),$
 $(4, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4, \text{Nt } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4]), (4, [\text{Tm } 1, \text{Nt } 3,$
 $\text{Nt } 2, \text{Nt } 4, \text{Nt } 1, \text{Nt } 3, \text{Nt } 3, \text{Nt } 2]),$
 $(4, [\text{Tm } 0, \text{Nt } 4, \text{Nt } 1, \text{Nt } 3, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4]), (4, [\text{Tm } 0, \text{Nt } 4, \text{Nt } 1, \text{Nt } 3,$
 $\text{Nt } 3, \text{Nt } 3, \text{Nt } 2]),$
 $(4, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4]), (4, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 3, \text{Nt } 2]),$
 $(4, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2, \text{Nt } 1, \text{Nt } 3, \text{Nt } 3, \text{Nt } 2, \text{Nt } 4]), (4, [\text{Tm } 1, \text{Nt } 3, \text{Nt } 2,$
 $\text{Nt } 1, \text{Nt } 3, \text{Nt } 3, \text{Nt } 2]),$

$(4, [Tm\ 0, Nt\ 1, Nt\ 3, Nt\ 3, Nt\ 2, Nt\ 4]), (4, [Tm\ 0, Nt\ 1, Nt\ 3, Nt\ 3, Nt\ 2])]$
in
 $solve_tri\ [3,2,1]\ [4,5,6]\ (set\ P0) = set\ P1 \wedge expand_tri\ [4,1,2,3]\ (set\ P1)$
 $= set\ P2$
 $\langle proof \rangle$

7 Complexity

Our method has exponential complexity, which we demonstrate below. Alternative polynomial methods are described in the literature [1].

We start with an informal proof that the blowup of the whole method can be as bad as 2^{n^2} , where n is the number of non terminals, and the starting grammar has $4n$ productions.

Consider this grammar, where a and b are terminals and we use nested alternatives in the obvious way:

$$A0 \rightarrow A1\ (a \mid b) \mid A2\ (a \mid b) \mid \dots \mid An\ (a \mid b) \mid a \mid b$$

$$A(i+1) \rightarrow Ai\ (a \mid b)$$

Expanding all alternatives makes this a grammar of size $4n$.

When converting this grammar into triangular form, starting with $A0$, we find that $A0$ remains the same after *expand_hd*, and *solve_lrec* introduces a new additional production for every $A0$ production, which we will ignore to simplify things:

Then every *expand_hd* step yields for Ai these number of productions:

- (1) 2^{i+1} productions with rhs $Ak\ (a \mid b)^{i+1}$ for every $k \in [i+1, n]$,
- (2) 2^{i+1} productions with rhs $(a \mid b)^{i+1}$,
- (3) 2^{i+1} productions with rhs $Ai\ (a \mid b)^{i+1}$.

Note that $(a \mid b)^{i+1}$ represents all words of length $i+1$ over $\{a, b\}$. Solving the left recursion again introduces a new additional production for every production of form (1) and (2), which we will again ignore for simplicity. The productions of (3) get removed by *solve_lrec*. We will not consider the productions of the newly introduced nonterminals.

In the triangular form, every Ai has at least 2^{i+1} productions starting with terminals (2) and 2^{i+1} productions with rhs starting with Ak for every $k \in [i+1, n]$.

When expanding the triangular form starting from An , which has at least the 2^{i+1} productions from (2), we observe that the number of productions of Ai (denoted by $\#Ai$) is $\#Ai \geq 2^{i+1} * \#A(i+1)$ (Only considering the productions of the form $A(i+1)\ (a \mid b)^{i+1}$). This yields that $\#Ai \geq 2^{i+1} * 2^{((i+2) + \dots + (n+1))} = 2^{((i+1) + (i+2) + \dots + (n+1))}$. Thus $\#A0 \geq 2^{(1 + 2 + \dots + n + (n+1))} = 2^{((n+1)*(n+2)/2)}$.

Below we prove formally that *expand_tri* can cause exponential blowup.

Bad grammar: Constructs a grammar which leads to a exponential blowup when expanded by *expand_tri*:

```
fun bad_grammar :: 'n list  $\Rightarrow$  ('n, nat)Prods where
  bad_grammar [] = {}
| bad_grammar [A] = {(A, [Tm 0]), (A, [Tm 1])}
| bad_grammar (A#B#As) = {(A, Nt B # [Tm 0]), (A, Nt B # [Tm 1])}  $\cup$ 
  (bad_grammar (B#As))
```

lemma bad_gram_simp1: $A \notin \text{set } As \implies (A, Bs) \notin (\text{bad_grammar } As)$
 <proof>

lemma expand_tri_simp1: $A \notin \text{set } As \implies (A, Bs) \in R \implies (A, Bs) \in \text{expand_tri } As \ R$
 <proof>

lemma expand_tri_iff1: $A \notin \text{set } As \implies (A, Bs) \in \text{expand_tri } As \ R \longleftrightarrow (A, Bs) \in R$
 <proof>

lemma expand_tri_insert_simp:
 $B \notin \text{set } As \implies \text{expand_tri } As \ (\text{insert } (B, Bs) \ R) = \text{insert } (B, Bs) \ (\text{expand_tri } As \ R)$
 <proof>

lemma expand_tri_bad_grammar_simp1:
 $\text{distinct } (A\#As) \implies \text{length } As \geq 1$
 $\implies \text{expand_tri } As \ (\text{bad_grammar } (A\#As))$
 $= \{(A, Nt (\text{hd } As) \# [Tm 0]), (A, Nt (\text{hd } As) \# [Tm 1])\} \cup (\text{expand_tri } As \ (\text{bad_grammar } As))$
 <proof>

lemma finite_bad_grammar: $\text{finite } (\text{bad_grammar } As)$
 <proof>

lemma finite_expand_tri: $\text{finite } R \implies \text{finite } (\text{expand_tri } As \ R)$
 <proof>

The last Nt expanded by *expand_tri* has an exponential number of productions.

lemma bad_gram_last_expanded_card:
 $\llbracket \text{distinct } As; \text{length } As = n; n \geq 1 \rrbracket$
 $\implies \text{card } (\{v. (\text{hd } As, v) \in \text{expand_tri } As \ (\text{bad_grammar } As)\}) = 2^n$
 <proof>

The productions resulting from *expand_tri* (*bad_grammar*) have at least exponential size.

theorem expand_tri_blowup: **assumes** $n \geq 1$

shows *card (expand_tri [0..<n] (bad_grammar [0..<n]))* $\geq 2^n$
<proof>

end

References

- [1] N. Blum and R. Koch. Greibach normal form transformation revisited. *Inf. Comput.*, 150(1):112–118, 1999.
- [2] S. A. Greibach. A new normal-form theorem for context-free phrase structure grammars. *J. ACM*, 12(1):42–52, 1965.
- [3] J. E. Hopcroft and J. D. Ullman. *Introduction to Automata Theory, Languages and Computation*. Addison-Wesley, 1979.