

# Verifying String Replacing Procedures by Supercompilation

Antonina Nepeivoda

Program Systems Institute of RAS

Verification and Program Transformation, April 2, 2019

# Verification of String-Replace Algorithms

- Strings are the main data type for the majority of scripting languages (Python, JavaScript, PHP). String-replace algorithms are used either in sanitization and implicit browser transduction ([Bjorner 2009], [Yu 2014], [Trinh 2016], [Abdulla 2017]).
- Replace-all algorithms for strings are Turing complete [Markov 1960].

```
var b = new B();
var doc = document;
var w = "wr"+ "ite";
doc[w](b.f);
```

```
function json_decode(str) {
  str = str.replace(/ip/g, "ip address");
  str = str.replace(/dom/g, "domain");
  return str;}
function json_show(str) {
  var arr = JSON.parse(str);
  var c = arr[0].content.split("&");
  var s = c[0]++c[1];
  document.getElementById("info").innerHTML = s;}

res = json_decode(input);
json_show(res);
```

# Main Contribution

The existing methods of supercompilation are developed and implemented to verify safety properties of programs over strings of an arbitrary length.

- The program state is a pair

$$\{[\text{Stack}], [\text{Constraints}]\},$$

where the set of constraints contains word equations and inequalities of a restricted form.

- The methods to unfold, fold, and generalize the program states with constraints are developed.
- The constraints are generated either from the program syntax or from the program behaviour on computation paths.

# Verification of Safety Properties by Transformation

Consider a program  $\mathcal{P}(u)$ , where  $u = \langle u_1, \dots, u_n \rangle$ ,  $u \in \mathcal{D}$ , and a parameterized input point  $\mathcal{P}(u')$ ,  $u' \in \mathcal{D}'$ ,  $\mathcal{D}' \subset \mathcal{D}$ . Given a safety property  $q$ ,

- specialize  $\mathcal{P}(u')$  to residual program  $\mathcal{P}'(u')$ ;
- check a simple explicit syntactic property of  $\mathcal{P}'(u')$  showing that  $q$  holds.

Then the safety property of  $\mathcal{P}$  is verified on the data set  $\mathcal{D}'$ .

## Example

The input points are  $q(\mathbf{G}(u))$  and  $\mathbf{H}(u)$  respectively.

Input Program	Residual Program
$q(\alpha_1) = \mathbf{T};$	$\mathbf{H}(\delta_0) = \mathbf{T};$
$q(\alpha_2) = \mathbf{F};$	$\mathbf{H}(\delta_1) = \mathbf{H}(\xi_1);$
$\mathbf{G}(\beta_1) = \mathbf{G}(\gamma_1);$	$\mathbf{H}(\delta_2) = \mathbf{H}(\xi_2);$
$\dots$	$\dots$
$\mathbf{G}(\beta_k) = \mathbf{G}(\gamma_k);$	$\mathbf{H}(\delta_m) = \mathbf{H}(\xi_m);$

The residual program does not contain rules returning  $\mathbf{F}$ . Hence, the input program is safe.

# Markov Normal Algorithms

Let  $\Sigma$  be an alphabet,  $\Phi, \Psi \in \Sigma^+$ ,  $x$  and  $y$  be variables taking values from  $\Sigma^*$ .  
 A *Markov normal algorithm* is a list of rewriting rules [Markov 1960].

$$\begin{array}{c|c}
 x++\Phi++y \rightarrow x++\Psi++y & x++\Phi++y \rightarrow x++\Psi++y \\
 \textit{normal rule} & \textit{terminating rule}
 \end{array}$$

The order of the rewriting process is controlled by:

- *Screening*. A string is matched against the left-hand sides of the rules from the beginning to the end of the rule list.
- *Markov's rule*. If a string has several occurrences of  $\Phi$ , we choose the first one to apply the rule.

## Example

$$\begin{array}{l}
 x++\mathbf{A}++\mathbf{A}++y \rightarrow x++\mathbf{A}++\mathbf{B}++y \\
 x++\mathbf{B}++\mathbf{A}++y \rightarrow x++\mathbf{B}++\mathbf{B}++y \\
 x++\mathbf{B}++y \rightarrow x++\mathbf{B}++y
 \end{array}$$

A run (the associative concatenation  $++$  is omitted):

$$\mathbf{AAAAB} \rightarrow \mathbf{ABAAB} \rightarrow \mathbf{ABABB} \rightarrow \mathbf{ABBBB} \rightarrow \mathbf{ABBBB}$$

# Presentation Language: Sample Example

$$\begin{array}{ll}
 \text{Go}(x) & = \text{Test}(\text{F}(\varepsilon, x)); \\
 \text{F}(x, \mathbf{AB} y) & = \text{F}(x, y); & \text{Test}(x \mathbf{AB} y) & = \mathbf{F}; \\
 \text{F}(x, y c \mathbf{AB} z) & = \text{F}(x y, c z); & \text{Test}(x) & = \mathbf{T}; \\
 \text{F}(x, y) & = x y;
 \end{array}$$


---

$\varepsilon$  — the empty string,  $x, y, z$  — string type variables,  $c$  — character type variable.

The entry point:  $\text{Go}(v)$ .

A run example:

$$\text{Go}(\mathbf{BAABB}) \rightarrow \text{Test}(\text{F}(\varepsilon, \mathbf{BAABB})) \rightarrow \text{Test}(\text{F}(\mathbf{B}, \mathbf{AB})) \rightarrow \text{Test}(\text{F}(\mathbf{B}, \varepsilon)) \rightarrow \text{Test}(\mathbf{B}) \rightarrow \mathbf{T}$$

The rule application order is the same as for Markov normal algorithms.

# Presentation Language: Sample Example

$$\begin{array}{lcl}
 \text{Go}(x) & = & \text{Test}(F(\varepsilon, x)); \\
 F(x, \mathbf{AB} y) & = & F(x, y); \quad \text{Test}(x \mathbf{AB} y) = \mathbf{F}; \\
 F(x, y \mathbf{c} \mathbf{AB} z) & = & F(x y, \mathbf{c} z); \quad \text{Test}(x) = \mathbf{T}; \\
 F(x, y) & = & x y;
 \end{array}$$


---

Does F a correct exhaustive deletion? Does not!  
 The program is not safe.

$$\begin{aligned}
 \text{Go}(\mathbf{AAABBB}) &\rightarrow \text{Test}(F(\varepsilon, \mathbf{AAABBB})) \rightarrow \text{Test}(F(\mathbf{A}, \mathbf{ABB})) \\
 &\rightarrow \text{Test}(F(\mathbf{A}, \mathbf{B})) \rightarrow \text{Test}(\mathbf{AB}) \rightarrow \mathbf{F}
 \end{aligned}$$

# More Examples of String-Replace Algorithms

Which of the following programs does a correct exhaustive deletion?

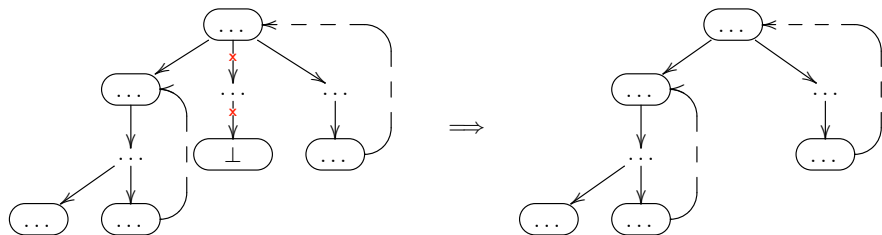
$G_0(x)$	$=$	$\text{Test}(G(\varepsilon, x));$		
$G(\varepsilon, \mathbf{AB}x)$	$=$	$G(\varepsilon, x);$	$\text{Test}(x \mathbf{AB}y)$	$=$ <b>F</b> ;
$G(xc, \mathbf{AB}y)$	$=$	$G(x, cy);$	$\text{Test}(x)$	$=$ <b>T</b> ;
$G(x, yc \mathbf{AB}z)$	$=$	$G(xy, cz);$		
$G(x, y)$	$=$	$xy;$		

---

$G_0(x)$	$=$	$\text{Test}(H(\varepsilon, x));$		
$H(x, y \mathbf{AAB}z)$	$=$	$H(x, y \mathbf{A}z);$	$\text{Test}(x \mathbf{AB}y)$	$=$ <b>F</b> ;
$H(x, y \mathbf{AB}z)$	$=$	$H(xy, z);$	$\text{Test}(x)$	$=$ <b>T</b> ;
$H(x, y)$	$=$	$xy;$		



# Brief Introduction to Supercompilation

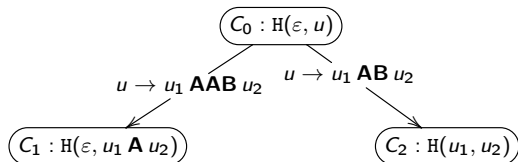


Supercompilation is a program specialization method based on unfold/fold operations. Introduced by V. F. Turchin using the string operating language Refal as the input language [Turchin 1986].

One of the ways of verification by supercompilation is to prune unreachable branches in the unfolding tree.

# Unfolding for String Data Type

The arrow  $\rightarrow$  stands for the narrowing relation.



**Problem:** The narrowing for  $C_1$  is an instance of the narrowing for  $C_2$ .

We use the negative constraints to preserve information in the program states.

# Configurations

$\Sigma$  — an alphabet,  $\mathcal{V}_c$  — a set of the character type variables.

## Definition

A *linear pattern*  $P$  is an expression  $\Phi_0 ++ z_1 ++ \Phi_1 ++ \dots ++ z_n ++ \Phi_n$ , where  $\Phi_i \in (\Sigma \cup \mathcal{V}_c)^*$ , and all the string type variables  $z_i$  are *distinct*.

Let  $\text{Match}(\Delta, P)$  return **T** if  $\Delta$  matches against  $P$  and **F** otherwise.

A *linear word inequality* is a predicate  $\mathcal{I} = \neg \text{Match}(u, P)$ , where  $u$  is a parameter. Briefly:  $u \neq P$ .

## Definition

A *configuration* is a pair  $\{\gamma, \mathcal{Q}\}$ , where  $\mathcal{Q}$  is a constraint of the form  $\bigwedge_{i=1}^n \bigvee_{j=1}^{k_i} \mathcal{I}_j^i$ , and  $\gamma$  is a parameterized expression.

## Example

$c ++ z_0 ++ c$ ,  $c \in \mathcal{V}_c$  is a linear pattern; while  $z_0 ++ \mathbf{A} ++ z_0$  is not.

# Simplifying Linear Patterns

Given a pattern  $P$  we say that a variable  $c \in \mathcal{V}_c$  is *floating* in  $P$ , if

- $c$  occurs in  $P$  only once;
- $c$  is either preceded or succeeded by a string type variable.

We say that a linear pattern  $P$  is *in a normal form*, if

- it does not contain two or more string variables in a row;
- all the floating variables from  $\mathcal{V}_c$  are both preceded and succeeded by string type variables in  $P$ .

## Example

In the pattern  $z_0 c_1 \mathbf{A} c_2 \mathbf{B} c_3 z_1 z_2 c_3$ , only  $c_1$  is a floating variable.

The normal form is  $z_0 c_1 z_1 \mathbf{A} c_2 \mathbf{B} c_3 z_2 c_3$ .

# Simplifying Linear Patterns

The pattern language  $\mathcal{L}(P)$  is a set of all strings  $\Delta$  s.t.  $\text{Match}(\Delta, P) = \mathbf{T}$ . We consider the case of an infinite  $\Sigma$ .

Inclusion of pattern languages for linear patterns with no character type variables occurs iff a substitution of pattern variables exists [Angluin, 1980].

Does not hold for all linear patterns:  $\mathcal{L}(c x) \neq \mathcal{L}(x c)$ .

## Theorem (Inclusion for normalized patterns)

*Given two linear patterns  $P_1$  and  $P_2$  in the normal form,  $\mathcal{L}(P_1) \subseteq \mathcal{L}(P_2)$  iff a substitution  $\sigma$  of variables exists s.t.  $\sigma(P_2) = P_1$ .*

# Folding

$L(C_i)$  is a set of computation paths starting from  $C_i$ .  $u, w, v$  are the string type parameters.

Given two configurations  $C_i = \{\gamma_i(w_1, \dots, w_{k_i}), Q_i\}$ ,  $i = 1, 2$ ,  $C_2$  is folded with  $C_1$  if there is a substitution  $\sigma$  of parameters of  $C_1$  s.t.

- $\gamma_1(\sigma(w_1), \dots, \sigma(w_{k_1})) \equiv \gamma_2(w_1, \dots, w_{k_2})$ ;
- $Q_2$  implies  $\sigma(Q_1)$ .

Hence,  $L(C_2) \subseteq L(C_1)$ .

## Example

The configuration  $\{F(\mathbf{B}, v_1 \mathbf{A} \mathbf{B} v_2), v_1 \neq z_1 \mathbf{A} z_2\}$  is folded with  $\{F(u_0, u_1 \mathbf{A} u_2), u_0 \neq z_0 \mathbf{A} z_1\}$ .

The substitution is  $\sigma(u_0) = \mathbf{B}$ ,  $\sigma(u_1) = v_1$ ,  $\sigma(u_2) = \mathbf{B} v_2$ .

Obviously,

$$v_1 \neq z_1 \mathbf{A} z_2 \Rightarrow (\mathbf{B} \neq z_0 \mathbf{A} z_1)$$

# Generalization

A *generalization* of configurations  $\{\gamma_1, Q_1\}$ ,  $\{\gamma_2, Q_2\}$  is a configuration  $\{\gamma', Q'\}$  and substitutions  $\sigma_1, \sigma_2$  s.t.

- $\gamma'(\sigma_1(u_1), \dots, \sigma_1(u_n)) \equiv \gamma_1(w_1, \dots, w_{k_1})$ ,  
 $\gamma'(\sigma_2(u_1), \dots, \sigma_2(u_n)) \equiv \gamma_2(w_1, \dots, w_{k_2})$ ;
- $Q_1 \Rightarrow \sigma_1(Q')$ ,  $Q_2 \Rightarrow \sigma_2(Q')$ .

Hence  $L(C_1) \subseteq L(C')$  &  $L(C_2) \subseteq L(C')$ .

## Example

$\{F(\varepsilon, v), \mathbf{T}\}$  and  $\{F(s, u), s \neq \mathbf{A}\}$ , where  $s$  is a character type parameter, can be generalized to  $\{F(w_1, w_2), w_1 \neq z_0 \mathbf{A} z_1\}$ .

$$\begin{aligned} \sigma_1(w_1) &= \varepsilon; & \sigma_2(w_1) &= s; \\ \sigma_1(w_2) &= v; & \sigma_2(w_2) &= u; \\ \mathbf{T} &\Rightarrow \varepsilon \neq z_0 \mathbf{A} z_1; & s \neq \mathbf{A} &\Rightarrow s \neq z_0 \mathbf{A} z_1. \end{aligned}$$

# Generalization Problems for Strings

Many possible generalizations, not comparable with each other.

		$\Phi_1 = \mathbf{B}$	$\Phi_2 = \mathbf{B B B}$
	$\Psi = \mathbf{B}u$	$\sigma_1(u) = \varepsilon$	$\sigma_2(u) = \mathbf{B B}$
OR	$\Psi = u\mathbf{B}$	$\sigma_1(u) = \varepsilon$	$\sigma_2(u) = \mathbf{B B}$
OR	$\Psi = u\mathbf{B}u$	$\sigma_1(u) = \varepsilon$	$\sigma_2(u) = \mathbf{B}$
OR	$\Psi = \mathbf{B}uu$	$\sigma_1(u) = \varepsilon$	$\sigma_2(u) = \mathbf{B}$
OR	$\Psi = uu\mathbf{B}$	$\sigma_1(u) = \varepsilon$	$\sigma_2(u) = \mathbf{B}$

...and an infinite list of the other generalizations.

The notion of a *minimal* generalization instead of *the* most specific generalization:

- Two or more string parameters in a row are forbidden.
- For all  $u$ ,  $\sigma_2(u) \trianglelefteq \sigma_1(u)$ .<sup>1</sup>
- For no  $u$ ,  $\sigma_2(u) = \varepsilon$ .

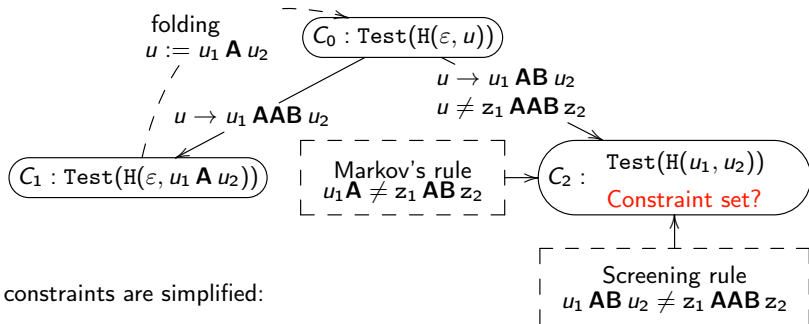
<sup>1</sup> $\trianglelefteq$  is the homeomorphic embedding relation.



## Supercompilation Example

$$\begin{array}{lcl}
 \text{Go}(x) & = & \text{Test}(\text{H}(\varepsilon, x)); \\
 \text{H}(x, y \mathbf{AAB} z) & = & \text{H}(x, y \mathbf{A} z); \quad \text{Test}(x \mathbf{AB} y) = \mathbf{F}; \\
 \text{H}(x, y \mathbf{AB} z) & = & \text{H}(x y, z); \quad \text{Test}(x) = \mathbf{T}; \\
 \text{H}(x, y) & = & x y;
 \end{array}$$


---



The constraints are simplified:

$$u_1 \mathbf{A} \neq z_1 \mathbf{AB} z_2 \Rightarrow u_1 \neq z_1 \mathbf{AB} z_2$$

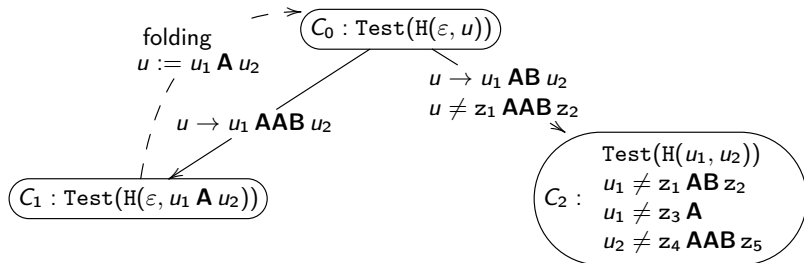
$$u_1 \mathbf{AB} u_2 \neq z_1 \mathbf{AAB} z_2 \Rightarrow (u_1 \neq z_3 \mathbf{A}) \ \& \ (u_2 \neq z_4 \mathbf{AAB} z_5) \ \& \ (u_1 \neq z_6 \mathbf{AAB} z_7).$$

Since  $u_1 \neq z_1 \mathbf{AB} z_2 \Rightarrow u_1 \neq z_6 \mathbf{AAB} z_7$ , the last constraint is removed.

## Supercompilation Example

$$\begin{array}{ll}
 \text{Go}(x) & = \text{Test}(\text{H}(\varepsilon, x)); \\
 \text{H}(x, y \mathbf{AAB} z) & = \text{H}(x, y \mathbf{A} z); & \text{Test}(x \mathbf{AB} y) & = \mathbf{F}; \\
 \text{H}(x, y \mathbf{AB} z) & = \text{H}(x y, z); & \text{Test}(x) & = \mathbf{T}; \\
 \text{H}(x, y) & = x y;
 \end{array}$$


---



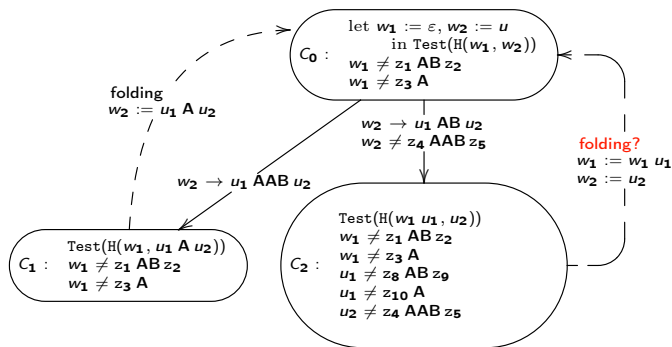
Generalization of  $C_0$  and  $C_2$  is constructed:  $\{\text{Test}(\text{H}(w_1, w_2)), \text{??}\}$ .

$\sigma_1(w_1) = \varepsilon$ ,  $\sigma_2(w_1) = u_1$ ;  $\sigma_1(w_2) = u$ ,  $\sigma_2(w_2) = u_2$ .

No constraints on  $u$ , hence no constraints on  $w_2$  in the generalization.

$\varepsilon \neq z_1 \mathbf{AB} z_2$  and  $\varepsilon \neq z_3 \mathbf{A}$ , hence both constraints on  $u_1$  in  $C_2$  hold for  $w_1$ .

## Supercompilation Example

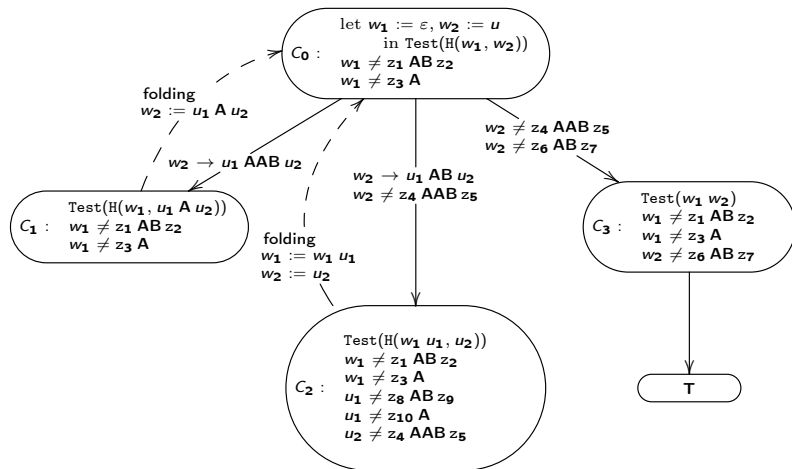


The implication is checked.

$$(w_1 \neq z_1 \mathbf{AB} z_2) \ \& \ (w_1 \neq z_3 \mathbf{A}) \ \& \ (u_1 \neq z_8 \mathbf{AB} z_9) \ \& \ (u_1 \neq z_{10} \mathbf{A}) \\ \Rightarrow (w_1 u_1 \neq z_1 \mathbf{AB} z_2 \ \& \ w_1 u_1 \neq z_3 \mathbf{A}).$$

It holds, hence the folding occurs.

## Entire Residual Graph of Supercompilation



The residual graph does not contain configurations returning **F**.  
The exhaustive deletion algorithm given in the function  $H$  is safe.

## Two Residual Programs

$\text{Go}(x)$	$=$	$\text{Test}(\text{H}(\varepsilon, x));$	$\text{Go}(x)$	$=$	$\text{Test}_0(x);$
$\text{H}(x, y \mathbf{AAB} z)$	$=$	$\text{H}(x, y \mathbf{A} z);$	$\text{Test}_0(x \mathbf{AAB} y)$	$=$	$\text{Test}_0(x \mathbf{A} y);$
$\text{H}(x, y \mathbf{AB} z)$	$=$	$\text{H}(x y, z);$	$\text{Test}_0(x \mathbf{AB} y)$	$=$	$\text{Test}_1(x, y);$
$\text{H}(x, y)$	$=$	$x y;$	$\text{Test}_0(x)$	$=$	$\mathbf{T};$
$\text{Test}(x \mathbf{AB} y)$	$=$	$\mathbf{F};$	$\text{Test}_1(z, x \mathbf{AB} y)$	$=$	$\text{Test}_1(z x, y);$
$\text{Test}(x)$	$=$	$\mathbf{T};$	$\text{Test}_1(z, x)$	$=$	$\mathbf{T};$

---

$\text{Go}(x)$	$=$	$\text{Test}(\text{G}(\varepsilon, x));$	$\text{Go}(x)$	$=$	$\text{Test}_0(x);$
$\text{G}(\varepsilon, \mathbf{AB} x)$	$=$	$\text{G}(\varepsilon, x);$	$\text{Test}_0(\mathbf{AB} x)$	$=$	$\text{Test}_0(x);$
$\text{G}(x c, \mathbf{AB} y)$	$=$	$\text{G}(x, c y);$	$\text{Test}_0(x c \mathbf{AB} y)$	$=$	$\text{Test}_1(x c, y);$
$\text{G}(x, y c \mathbf{AB} z)$	$=$	$\text{G}(x y, c z);$	$\text{Test}_0(x)$	$=$	$\mathbf{T};$
$\text{G}(x, y)$	$=$	$x y;$	$\text{Test}_1(\mathbf{A}, \mathbf{B} x)$	$=$	$\text{Test}_0(x);$
$\text{Test}(x \mathbf{AB} y)$	$=$	$\mathbf{F};$	$\text{Test}_1(x c \mathbf{A}, \mathbf{B} y)$	$=$	$\text{Test}_1(x c, y);$
$\text{Test}(x)$	$=$	$\mathbf{T};$	$\text{Test}_1(x, \mathbf{AB} y)$	$=$	$\text{Test}_1(x, y);$
			$\text{Test}_1(x, y c \mathbf{AB} z)$	$=$	$\text{Test}_2(x, y c, z);$
			$\text{Test}_1(x, y)$	$=$	$\mathbf{T};$
			$\text{Test}_2(x, \mathbf{A}, \mathbf{B} z)$	$=$	$\text{Test}_1(x, z);$
			$\text{Test}_2(x, y c \mathbf{A}, \mathbf{B} z)$	$=$	$\text{Test}_2(x, y c, z);$
			$\text{Test}_2(x, y, \mathbf{AB} z)$	$=$	$\text{Test}_2(x, y, z);$
			$\text{Test}_2(x, y, z_1 c \mathbf{AB} z_2)$	$=$	$\text{Test}_2(x, y z_1 c, z_2);$
			$\text{Test}_2(x, y, z)$	$=$	$\mathbf{T};$

# Termination Issue

## Theorem

*For every term  $\gamma$ , the set of generalizations of  $\gamma$  is finite (due to the choice of the minimal generalization).*

## Theorem

*The length of every linear word inequality  $u \neq P$  generated by the algorithm is uniformly bounded by the size of the input program.*

# Conclusion

The supercompilation methods are developed and implemented for programs modelling Markov normal algorithms.

- multiple occurrences of the character type variables;
- repeated string type variables in the right-hand sides;
- the constraints may be generated both from the syntax or from a program behaviour.

## Future work

An entire supercompilation method for a language allowing multiple occurrences of string type variables in the patterns.

# The model supercompiler

The page of the model supercompiler MSCP-A:

[http://refal.botik.ru/mscp/mscp-a\\_eng.html](http://refal.botik.ru/mscp/mscp-a_eng.html)