



ELSEVIER

Available online at www.sciencedirect.com

ScienceDirect

Electronic Notes in
Theoretical Computer
Science

Electronic Notes in Theoretical Computer Science 177 (2007) 153–168

www.elsevier.com/locate/entcs

Implementing Dynamic-Cut in TOY ¹

R. Caballero² Y. García-Ruiz³

*Departamento de Sistemas Informáticos y Programación
Universidad Complutense de Madrid
Madrid, Spain*

Abstract

This paper presents the integration of the optimization known as *dynamic cut* within the functional-logic system TOY . The implementation automatically detects deterministic functions at compile time, and includes in the generated code the test for detecting at run-time the computations that can actually be pruned. The outcome is a much better performance when executing deterministic functions including either or-branches in their definitional trees or extra variables in their conditions, with no serious overhead in the rest of the computations. The paper also proves the correctness of the criterion used for detecting deterministic functions w.r.t. the semantic calculus CRWL.

Keywords: determinism, functional-logic Programming, program analysis, programming language implementation.

1 Introduction

Nondeterminism is one of the characteristic features of Logic Programming shared by Functional-Logic Programming. It allows elegant algorithm definitions, increasing the expressiveness of programs. However, this benefit has an associated drawback, namely the lack of efficiency of the computations. There are two main reasons for this:

- The complexity of the search engine required by nondeterministic programs, which slows down the execution mechanism.
- The possible occurrence of redundant subcomputations during a computation.

In the Logic Programming language Prolog, the second point is partially solved by introducing a non-declarative mechanism, the so-called *cut*. Programs using cuts are much more efficient, but at the price of becoming less declarative.

¹ This work has been funded by the projects TIN2005-09207-C03-03 and S-0505/TIC/0407.

² Email: rafa@sip.ucm.es

³ Email: ygruiz@gmail.com

In the case of Functional-Logic Programming the situation is somehow alleviated by the *demand driven strategy* [2,8], which is based on the use of *definitional trees* [1,8]. Given any particular program function, the strategy uses the structure of the left-hand sides of the program rules in order to reduce the number of redundant subcomputations. The implementation of modern Functional-Logic languages such as *TOY* [9] or Curry [6] is based on this strategy. Our proposal also relies on the demand driven strategy, but introduces a safe and declarative optimization to further improve the efficiency of deterministic computations. This optimization is the *dynamic cut*, first proposed by Rita Loogen and Stephan Winkler in [10]. In [4,3] the same ideas were adapted to a setting including non-deterministic functions and a demand driven strategy, showing by means of examples the efficiency of the optimization.

However, in spite of being well-known and accepted as an interesting optimization, the dynamic cut had not been implemented in any real system up to now. In this paper we present this implementation in the functional-logic system *TOY* (available at <http://toy.sourceforge.net>).

The dynamic cut considers two special fragments of code:

- (i) Rules with existential variables in the conditions.
- (ii) Sets of overlapping rules occurring in deterministic functions.

As we will explain in section 3, computations involving these fragments of code can be safely pruned if certain dynamic conditions are fulfilled.

A key point of the optimization is detecting deterministic functions. The information about deterministic functions is required not only at compile time but also at run-time, when it is used for checking dynamically if the cut must take place in a particular computation. As previous works [10,4,3] have shown, this dynamic test is necessary for ensuring the correctness of the cut, i.e. that the optimization does not affect the set of solutions of any goal.

The determinism analysis performed by the system follows the well-known criterion of non-ambiguity already introduced in [10]. From the theoretical point of view, the novelty of this paper w.r.t. previous work is that we have proved formally the correctness of such a criterion w.r.t. the semantic calculus CRWL, proposed as suitable logic foundation for Functional-Logic Programming in [5]. Of course, completeness cannot be established because determinism is an undecidable property [13]. For that reason we also allow the user to annotate explicitly some functions as deterministic.

The paper is organized as follows. The next section introduces the non-ambiguity criterion for detecting deterministic functions and the correctness theorem. Section 3 shows by means of examples the cases where the optimization will be applied. Section 4 presents the steps followed during the implementation of the dynamic cut in *TOY*, and Section 5 finalizes presenting some conclusions.

2 Detecting Deterministic Functions in Functional-Logic Programs

This section proves the correctness of the non-ambiguity condition used for detecting deterministic functions w.r.t. the semantic calculus CRWL [5].

2.1 The CRWL calculus

CRWL is an inference system consisting of six inference rules:

$$\begin{array}{ll}
 \mathbf{BT} \text{ Bottom: } \frac{}{e \rightarrow \perp} & \mathbf{RF} \text{ Reflexivity: } \frac{}{X \rightarrow X} \\
 \mathbf{DC} \text{ Decomposition } \frac{e_1 \rightarrow t_1 \dots e_m \rightarrow t_m}{c e_1 \dots e_m \rightarrow c t_1 \dots t_m} & c \in CD^n \cup FS^{n+1}, m \leq n, t_i \in CTerm_{\perp} \\
 \mathbf{FA} \text{ Function Application: } \frac{e_1 \rightarrow t_1 \dots, e_n \rightarrow t_n \quad C \quad r \rightarrow a \quad a a_1 \dots a_k \rightarrow t}{f e_1 \dots e_n a_1 \dots a_k \rightarrow t} (k \geq 0) \\
 & \text{if } t \neq \perp, (f t_1 \dots t_n \rightarrow r \Leftarrow C) \in [R]_{\perp} \\
 \mathbf{JN} \text{ Join: } \frac{e_1 \rightarrow t \quad e_2 \rightarrow t}{e_1 == e_2} & t \in CTerm
 \end{array}$$

The notation $[R]_{\perp}$ in rule **FA** represents the set of all the possible instances of program rules, where each particular instance is obtained from some function defining rule in \mathcal{R} , by some substitution of (possibly partial) terms in place of variables. See [5] for a detailed description of this and related calculi.

2.2 Deterministic Functional-Logic Functions

Before defining and characterizing deterministic functions we need to establish briefly some basic notions and terminology. We refer to [5] for more detailed definitions. We assume a *signature* $\Sigma = \langle DC, FS \rangle$, where DC and FS are ranked sets of *constructor symbols* resp. *function symbols*. Given a countably infinite set \mathcal{V} of variables, we build $CTerms$ (using only variables and constructors) and $Terms$ (using variables, constructors and function symbols). We extend Σ with a special nullary constructor \perp (0-arity constructor) obtaining a new signature Σ_{\perp} and we will write $Term_{\perp}$ and $CTerm_{\perp}$ (*partial terms*) for the corresponding sets of terms in this extended signature.

A *TOY* program \mathcal{P} is composed of data type declarations, type alias, infix operators, function type declarations and a set of defining rules for functions symbols. Each defining rule for a function $f \in FS$ has a *left-hand side*, a *right-hand side* and an optional *condition*:

$$\underbrace{f t_1 \dots t_n}_{\text{left-hand side}} \rightarrow \underbrace{r}_{\text{right-hand side}} \Leftarrow \underbrace{C}_{\text{condition}}$$

where $t_1 \dots t_n$ must be linear $CTerms$ and C must consist of finitely many (possibly zero) joinability statements $e_1 == e_2$ with $e_1, e_2 \in Term$. A natural *approximation ordering* \sqsubseteq for partial terms can be defined as the least partial ordering over $Term_{\perp}$ satisfying the following properties:

- $\perp \sqsubseteq t$, for all $t \in Term_{\perp}$

- $X \sqsubseteq X$, for all variable X
- if $t_1 \sqsubseteq s_1, \dots, t_n \sqsubseteq s_n$, then $c \ t_1 \dots t_n \sqsubseteq c \ s_1 \dots s_n$, for all $c \in DC^n$ and $t_i, s_i \in CTerm_{\perp}$.

A *partially ordered set* (poset in short) with bottom is a set S equipped with a partial order \sqsubseteq and a least element \perp (w.r.t. \sqsubseteq). $D \subseteq S$ is a *directed set* iff for all $x, y \in D$ there exists $z \in D$ such that $x \sqsubseteq z, y \sqsubseteq z$. A subset $A \subseteq S$ is a *cone*, iff $\perp \in A$ and for all $x \in A, y \in S \ y \sqsubseteq x \Rightarrow y \in A$. An *ideal* $I \subseteq S$ is a directed cone. The program semantics is defined by the semantic calculus CRWL presented in [5]. CRWL (Constructor Based ReWriting Logic) is a theoretical framework for the lazy functional logic programming paradigm. Given any program P , CRWL proves statements of the form $e \rightarrow t$ with $e \in Term_{\perp}$ and $t \in CTerm_{\perp}$. We denote by $\mathcal{P} \vdash_{CRWL} e \rightarrow t$ that the statement $e \rightarrow t$ can be proved in CRWL w.r.t. P . The intuitive idea is that t is a valid approximation of e in \mathcal{P} . The *denotation* of any $e \in Term_{\perp}$, written $\llbracket e \rrbracket$, is defined as: $\llbracket e \rrbracket = \{t \in CTerm_{\perp} \mid P \vdash_{CRWL} e \rightarrow t\}$.

Now we are ready for presenting the formal definition of deterministic function in our setting.

Definition 2.1 (Deterministic Functions)

Let f be a function defined in a program \mathcal{P} . We say that f is a deterministic function iff $\llbracket f \bar{t}_n \rrbracket$ is an ideal for every \bar{t}_n s.t. t_i is a $CTerm_{\perp}$ for all $i = 1 \dots n$.

We call a function *non-deterministic*, if it does not fulfill the previous definition. The intuitive idea behind a deterministic function is that it returns at most one result for any arbitrary ground parameters [7]. In addition, in a lazy setting whenever a function returns some value t , it is expected to return all the less defined terms $s \sqsubseteq t$ as well. The previous definition of deterministic function takes this idea into account. Consider for instance the following small program:

```
data pair = pair int int      f 1 = pair 1 2      g 1 = 1      g 1 = 2
```

Using CRWL it can be proved that $\llbracket f \ 1 \rrbracket = \{\perp, \text{pair } \perp \ \perp, \text{pair } 1 \ \perp, \text{pair } \perp \ 2, \text{pair } 1 \ 2\}$, $\llbracket f \ t \rrbracket = \{\perp\}$ if $t \neq 1$, $\llbracket g \ 1 \rrbracket = \{\perp, 1, 2\}$, $\llbracket g \ t \rrbracket = \{\perp\}$ if $t \neq 1$. Then g is a non-deterministic function because for the parameter 1 the set $\{\perp, 1, 2\}$ is not an ideal, in particular because it is not directed: taking $x = 1, y = 2$ it is not possible to find $z \in \{\perp, 1, 2\}$ s.t. $x \sqsubseteq z, z \sqsubseteq 2$. On the other hand, it is easy to check that f is a deterministic function.

2.3 Non-ambiguous functions

The definition 2.1 is only a formal definition and cannot be used in practice. In [4] an adaptation of the non-ambiguity condition of [11] is presented, which we will use as an easy mechanism for the effective recognition of deterministic functions. Although not all the deterministic functions are non-ambiguous, the non-ambiguity criterion will be enough for detecting several interesting deterministic functions.

Definition 2.2 (Non-ambiguous functions)

Let \mathcal{P} be a program defining a set of functions G . We say that $F \subseteq G$ is a set of

non-ambiguous functions if every $f \in F$ verifies:

- (i) If $f \bar{t}_n = e \Leftarrow C$ is a defining rule for f , then $\text{var}(e) \subseteq \text{var}(\bar{t})$ and all function symbols in e belong to F .
- (ii) For any pair of variants of defining rules for f , $f \bar{t}_n = e \Leftarrow C$, $f \bar{t}'_n = e' \Leftarrow C'$, one of the following two possibilities holds:
 - (a) Left-hand sides do not overlap, that is, the terms $(f \bar{t}_n)$ and $(f \bar{t}'_n)$ are not unifiable.
 - (b) If θ is the m.g.u. of $f \bar{t}_n$ and $f \bar{t}'_n$, then $e\theta \equiv e'\theta$.

In [3,4] the inclusion of the set on non-ambiguous functions in the set of deterministic functions was claimed. Here, and thank to the previous formal definition, we will be able to prove the result.

Before that we need some auxiliar lemmata. The proofs of these results are tedious but straightforward using induction on the structure of the CRWL-proofs and are not included for the sake of the space. The first two lemmata establish substitution properties that will play an important role in the proof. The lemmata use the symbol $CSubst$ for the set of all the *c-substitutions*, which are mappings $\theta : \mathcal{V} \rightarrow CTerm$, and the notation $CSubst_{\perp}$ for the set of all the *partial c-substitutions* $\theta : \mathcal{V} \rightarrow CTerm_{\perp}$ defined analogously. We note as $t\theta$ the result of applying the substitution θ to the term t .

Lemma 2.3 *Let $t \in CTerm$, $s \in CTerm_{\perp}$ be such that $t \sqsubseteq s$. There there exists a substitution $\theta \in CSubst_{\perp}$ verifying $t\theta = s$.*

Lemma 2.4 *Let $t, t' \in CTerm$ be such that: 1) t, t' are linear, 2) $\text{var}(t) \cap \text{var}(t') = \emptyset$ and 3) There exists $\gamma = \text{m.g.u.}(t, t')$. Let $s \in CTerm_{\perp}$ be a term and $\theta, \theta' \in CSubst_{\perp}$ such that $t\theta \sqsubseteq s$, $t'\theta' \sqsubseteq s$. Then there exists a substitution θ'' s.t. $t\gamma\theta'' = t'\gamma\theta'' = s$.*

Lemma 2.5 *Let \mathcal{P} be aprogram and $e \in Term_{\perp}$. Then:*

- i) *Let $t, t' \in CTerm_{\perp}$ be such that $\mathcal{P} \vdash_{CRWL} e \rightarrow t$ and $t' \sqsubseteq t$. Then $\mathcal{P} \vdash_{CRWL} e \rightarrow t'$.*
- ii) *Let \mathcal{P} be a program and $e \in Term_{\perp}$ and $\theta \in CSubst_{\perp}$ be s.t. $\mathcal{P} \vdash_{CRWL} e\theta \rightarrow t$. Then $\mathcal{P} \vdash_{CRWL} e\theta' \rightarrow t$ for all θ' s.t. $\theta \sqsubseteq \theta'$.*
- iii) *Let \bar{e}_n s.t. $e_i \in Term_{\perp}$ for all $i = 1 \dots n$, and s.t. $\mathcal{P} \vdash_{CRWL} e \bar{e}_n \rightarrow t$, and $a \in Term_{\perp}$ such that $e \sqsubseteq a$. Then $\mathcal{P} \vdash_{CRWL} a \bar{e}_n \rightarrow t$.*
- iv) *$\llbracket e \rrbracket$ is a cone.*

Now we are ready to prove that non-ambiguous functions are deterministic.

Theorem 2.6 . *Let \mathcal{P} be a program and f be a non-ambiguous function defined in \mathcal{P} . Then f is deterministic.*

Proof. In order to check that f is a deterministic function, we must prove that $\llbracket f \bar{t}_n \rrbracket$ is an ideal, i.e.:

- $\llbracket f \bar{t}_n \rrbracket$ is a cone by lemma 2.5 item iv).

- $\llbracket f \bar{t}_n \rrbracket$ is a directed set. We prove a more general result: *Consider $e \in Term_\perp$ and suppose that all the function symbols occurring in e are correspond to non-ambiguous functions. Then, $\llbracket e \rrbracket$ is a directed set.*

Let be. $t, t' \in Cterm_\perp$ verifying $(R1) : \mathcal{P} \vdash_{CRWL} e \rightarrow t$ and $(R2) : \mathcal{P} \vdash_{CRWL} e \rightarrow t'$. We prove that exists $s \in Cterm_\perp$ s.t.: a) $t \sqsubseteq s$, b) $t' \sqsubseteq s$ and c) $\mathcal{P} \vdash_{CRWL} e \rightarrow s$ by induction on the depth l of a CRWL -proof for $e \rightarrow t$:

$l = 0$. Three possible CWRL-inference rules:

- **BT.** Then $t = \perp$ and defining $s = t'$ we have: a) $\perp \sqsubseteq s$, b) $t' \sqsubseteq s$ and c) $\mathcal{P} \vdash_{CRWL} e \rightarrow s$ (by $(R2)$).
- **RF.** Then the proof for $(R1)$ must be of the form $X \rightarrow X$, and hence $e = X$ and $t = X$. Then t' only can be X or \perp (otherwise no CRWL inference could be applied and $(R2)$ would not hold). We define s as X and then: a) $t \sqsubseteq X$ b) $t' \sqsubseteq X$ c) $\mathcal{P} \vdash_{CRWL} e \rightarrow s$ by $(R1)$.
- **DC.** Then $e = c$, $t = c$, with $c \in DC^0$. Then t' must be either c or \perp . In any case defining s as c the result holds.

$l > 0$ There are three possible inference rules applied at the first step of the proof:

- **DC.** Then $e = c e_1 \dots e_m$, $t = c t_1 \dots t_m$ with $c \in DC^n \cup FS^{n+1}$, $m \leq n$. Analogously $t' = c t'_1 \dots t'_m$ and the first inference rules of any proof for $(R1)$ y $(R2)$ must be of the form:

$$(R1) : \frac{e_1 \rightarrow t_1 \dots e_m \rightarrow t_m}{c e_1 \dots e_m \rightarrow c t_1 \dots t_m} \quad (R2) : \frac{e_1 \rightarrow t'_1 \dots e_m \rightarrow t'_m}{c e_1 \dots e_m \rightarrow c t'_1 \dots t'_m}$$

The proofs for $\mathcal{P} \vdash_{CRWL} e_i \rightarrow t_i$ and $\mathcal{P} \vdash_{CRWL} e_i \rightarrow t'_i$ have a maximum depth of $l - 1$. Therefore by induction hypotheses exists $s_i \in Cterm_\perp$ satisfying $t_i, t'_i \sqsubseteq s_i$, and $\mathcal{P} \vdash_{CRWL} e_i \rightarrow s_i$ for all $1 \leq i \leq m$. Then defining $s = c s_1 \dots s_m$, $t \sqsubseteq s$, $t' \sqsubseteq s$ hold and $\mathcal{P} \vdash_{CRWL} e \rightarrow s$ with a proof starting with a DC inference.

- **JN.** Very similar to the previous case.
- **AF.** Then e is of the form $f \bar{e}_n$ with $e_i \in CTerm_\perp$ for $i = 1 \dots n$. Moreover n is greater or equal to the program arity of f . Hence an AF inference must have been applied at the first step of any proof of $(R2)$. In each case a suitable instance $(I1)$ y $(I2)$ must have been used. We call θ and θ' to the substitutions associated to the first and to the second instance respectively, $\theta, \theta' \in CSubst_\perp$.

The first inference step of each proof will be of the following form:

$$(1) : \frac{e_1 \rightarrow t_1 \theta, \dots, e_k \rightarrow t_k \theta, C\theta, r\theta \rightarrow a, a e_{k+1} \dots e_n \rightarrow t}{f e_1 \dots e_k e_{k+1} \dots e_n \rightarrow t}$$

$$(2) : \frac{e_1 \rightarrow t'_1 \theta', \dots, e_k \rightarrow t'_k \theta', C'\theta', r'\theta' \rightarrow a', a' e_{k+1} \dots e_n \rightarrow t'}{f e_1 \dots e_k e_{k+1} \dots e_n \rightarrow t'}$$

with $(k > 0)$, $t, t' \neq \perp$ and the rule instances:

$$I_1 : (f t_1 \dots t_k \rightarrow r \Leftarrow C)\theta \in [R]_\perp \quad I_2 : (f t'_1 \dots t'_k \rightarrow r' \Leftarrow C')\theta' \in [R]_\perp$$

Now we consider separately two cases: a) I_1 e I_2 correspond to the same program rule, and b) each instance correspond to a different program rule. The first case is easy to check and does not rely on the non-ambiguity criterion. For the sake

of the space we only include the proof of the case b).

Assume that I_1, I_2 are instances of two different program rules. By the non-ambiguity criterion there exists $\gamma = \text{m.g.u.}$ $(f \bar{t}_k, f \bar{t}'_k)$, i.e. $t_i \gamma = t'_i \gamma$ for $i = 1 \dots k$ and $r \gamma = r' \gamma$. Calling u_i to $t_i \gamma = t'_i \gamma$, the rule instances can be seen as: $(f u_1 \dots u_k \rightarrow r'' \Leftarrow C \gamma)$ and $(f u_1 \dots u_k \rightarrow r'' \Leftarrow C' \gamma)$. Now we must look for some $s \in CTerm_{\perp}$ such that: a) $t \sqsubseteq s$, b) $t' \sqsubseteq s$ and c) $\mathcal{P} \vdash_{\text{CRWL}} f \bar{e}_n \rightarrow s$ for some substitution θ'' . The proof of c) can be of one of these two forms

$$(4) : \frac{e_1 \rightarrow u_1 \theta'', \dots, e_k \rightarrow u_k \theta'', C \gamma \theta'', r'' \theta'' \rightarrow a'', a'' e_{k+1} \dots e_n \rightarrow s}{f e_1 \dots e_k e_{k+1} \dots e_n \rightarrow s}$$

$$(5) : \frac{e_1 \rightarrow u_1 \theta'', \dots, e_k \rightarrow u_k \theta'', C' \gamma \theta'', r'' \theta'' \rightarrow a'', a'' e_{k+1} \dots e_n \rightarrow s}{f e_1 \dots e_k e_{k+1} \dots e_n \rightarrow s}$$

We observe that γ unifies the heads and fusions the right-hand sides, but it doesn't relation C y C' . We consider the form (4) (the (5) is analogous). From the premises of (1) y (2) we know that $\mathcal{P} \vdash_{\text{CRWL}} e_i \rightarrow t_i \theta$ and $\mathcal{P} \vdash_{\text{CRWL}} e_i \rightarrow t'_i \theta'$ for $i = 1 \dots k$. By induction hypotheses exists $s_i \in CTerm_{\perp}$ s.t.: a) $t_i \theta \sqsubseteq s_i$, b) $t'_i \theta' \sqsubseteq s_i$, and c) $\mathcal{P} \vdash_{\text{CRWL}} e_i \rightarrow s_i$. Since t_i, t'_i are unified by γ , we can apply the Lemma 2.4. Then there exist substitutions θ_i which we can restrict to the variables in u_i s.t. $u_i \theta_i = s_i$. (u_1, \dots, u_k) is a linear tuple because (t_1, \dots, t_k) and (t'_1, \dots, t'_k) are both linear. Then we can define a substitution θ'' as:

$$\theta''(X) = \begin{cases} \theta_i(X) & \text{if } X \in \text{var}(t_i, t'_i) \text{ for some } i, 1 \leq i \leq k \\ \theta(X) & \text{otherwise} \end{cases}$$

ensuring that there exist CRWL -proofs of $e_i \rightarrow u_i \theta''$ for all $i = \{1, \dots, k\}$ in (4) (this is because $u_i \theta_i = u_i \theta''$).

Checking that rest of the premises of (4) also have CRWL -proof requires similar arguments.

□

The non-ambiguity condition characterizes a set of functions F as deterministic. This is because the value of a function may depend on other functions, and in general this dependence can be mutual. In practice the implementation starts with an empty set F of non-ambiguous functions, adding at each step to F those functions that satisfy the definition and that only depend on functions already in F . This is done until a fix-point for F is reached.

Although most of the deterministic functions that occur in a program are non-ambiguous as well, there are some functions which are not detected. This happens for instance in the function f of following example: $f \ 1 = 1 \quad f \ 1 = g \ 1 \quad g \ 1 = 1$. It would be useful to use additional determinism criteria, such as those based on abstract interpretation proposed in [12], but the detection of deterministic function will be still incomplete. For that reason the system allows the programmer to distinguish deterministic functions annotating them by using $-->$ instead of $=$, as in the following example: $f \ 1 --> 1 \quad f \ 1 --> g \ 1 \quad g \ 1 = 1$, which indicates that f is deterministic. The non-annotated functions like g will be analyzed following the non-ambiguity criterion.

% P_1 : 'Parallel' multiplication		% P_2 : 'Classical' multiplication	
data nat	= zero s nat	data nat	= zero s nat
add zero Y	= Y	add zero Y	= Y
add (s X) Y	= s (add X Y)	add (s X) Y	= s (add X Y)
multi zero _	= zero	multi zero _	= zero
multi _ zero	= zero	multi (s X) Y	= add Y (add X Y)
multi (s X) (s Y)	= s (add X (add Y (multi X Y)))		
power N zero	= s zero	power N zero	= s zero
power N (s M)	= multi N (power N M)	power N (s M)	= multi N (power N M)
odd zero	= false	odd zero	= false
odd (s zero)	= true	odd (s zero)	= true
odd (s (s N))	= odd N	odd (s (s N))	= odd N
toNat N	= if (N==0) then zero else s (toNat (N-1))	toNat N	= if (N==0) then zero else s (toNat (N-1))

Fig. 1. Two methods for multiplying

X	Y	P_1	P_2
0	100000	0	0
0	50000	0	0
100	1000	2.7	2.7
400	400	4.1	4.1
1000	100	4.9	-
50000	0	0	3.5
100000	0	0	-
multi (toNat X) (toNat Y)			

N	P_1	P_2
10^4	0.7	0
10^5	6.1	0
10^6	60.0	0
10^7	-	0
odd (power zero (toNat N)) <u>without</u> dynamic cut		

N	P_1	P_2
10^4	0	0
10^5	0	0
10^6	0	0
10^7	0	0
odd (power zero (toNat N)) <u>with</u> dynamic cut		

Fig. 2. Comparative tables

3 Pruning Deterministic Computations

In this section we present briefly the two different situations where the dynamic cut can be introduced.

3.1 Deterministic Functions Defined through Overlapping Program Rules

Sometimes deterministic functions can be defined in a natural way by using overlapping rules. Consider for instance the two programs of Figure 1. Both programs contain functions for computing arithmetic using Peano's representation. The function `toNat` is used for easily converting positive numbers of type `int` to their Peano representation. The only difference between P_1 and P_2 is the method for multiplying numbers. The function `multi` at P_2 , which we have called 'classical' reduces the first argument before each recursive call until it becomes zero. The method `multi` of P_1 , which we have called 'parallel', reduces both arguments before the recursive call. Observe that the first two rules of `multi` in P_1 are overlapping. However it is easy to check that it is a non-ambiguous and hence a deterministic function.

The first table at Figure 2 shows the time⁴ required for computing the first answer for goals of the form `multi (toNat X) (toNat Y) == R` in both programs.

The symbol `_` means that the system has run out of memory for the goal. From this data it is clear that the parallel `multi` of P_1 behaves better than its classical counterpart of P_2 . The reason is that in P_1 the computation of `multi` reduces the two arguments simultaneously saving both time and space. However this kind of 'parallel' definition is not used very often in Functional-Logic Programming because programmers know that overlapping rules can produce unexpected behaviors due to the backtracking mechanism. Indeed using P_1 a goal like `multi zero zero == R` has two solutions, both giving R the value `zero`, instead of only one as expected (and as the program P_2 does). Such redundant computations can affect the efficiency of other computations. The central table of Figure 2 contains the time required by both programs for checking if the N -th power of zero is odd without the dynamic cut optimization. The goal returns `no` in both cases as expected, but we observe that now P_1 behaves rather worse than P_2 , even running out of memory for large enough numbers. This is because the subgoal `power zero (toInt N)` needs to compute N multiplications, and in P_1 this means N redundant computations. Thus using P_1 without dynamic cut the goal `odd (power zero (toInt N))` will check N times if `zero` is odd, while in P_2 this is done only once. The dynamic cut solves this situation, detecting that `multi` in P_1 is a deterministic function and cutting the possibility of using the second rule of `multi` if the first one has succeeded producing a result (and satisfying some conditions explained below). The third table, at the right of Figure 2 has been obtained after activating the dynamic cut. The problem of the redundant computations has been solved. It is worth pointing out that the data of the first table do not change after activating the optimization, because all the goals considered produce only one answer, and the dynamic cut optimization only has effect on the second and posteriors answers.

3.2 Existential variables in conditions

Consider now the program of Figure 3. It includes a simple representation of DNA molecules, which are build by two chains of nucleotides. The nucleotides of the two strands are connected in compatible pairs, defined in the program through function `compatible`. The function `dna` detects if its two input parameters represent two strands that can be combined in a DNA molecule. Function `dnaPart` checks if the two input sequences `S1` and `S2` contain some subsequences `P1` and `P2` of length `L` that can occur associated in a DNA molecule. This function relies in function `part` which checks if the parameter `X` is a sublist of length `L` of the list `Y`. The functions `++` and `length`, represent respectively the concatenation of lists and the number of elements in a list. Consider the following session in the system *TOY*:

```
Toy> dnaPart (repeat 1000 adenine) (repeat 1000 thymine) 5
yes.    Elapsed time: 844 ms.
more solutions? y
yes.    Elapsed time: 40390 ms.
```

⁴ All the results displayed in seconds, obtained on a computer at 2.13 GHz with 1 Gb of RAM

```

data nucleotides = adenine | guanine | cytosine | thymine

compatible adenine thymine = true
compatible thymine adenine = true
compatible guanine cytosine = true
compatible cytosine guanine = true

dna [] [] = true
dna [N1|R1] [N2|R2] = true <== compatible N1 N2, (dna R1 R2)

dnaPart S1 S2 L = true <== part P1 S1 L , part P2 S2 L, dna P1 P2

part X Y L = true <== (U ++ X) ++ V == Y, length X == L

```

Fig. 3. Detecting DNA strands

The goal `dnaPart (repeat 1000 adenine) (repeat 1000 thymine) 5` asks if in two strands of 1000 nucleotides of adenine and thymine respectively it is possible to find two subsequences of 5 nucleotides, one from each strand, which can occur associated in a DNA molecule. The answer given by the system after 0.8 seconds is **yes** (actually *all* the subsequences of n elements of the first strand are compatible with all the subsequences of n elements of the second strand). If the user asks for a second answer, the same redundant answer **yes** is obtained after more than 40 seconds. The second answer is useless because it doesn't provide new information, and greatly affects the efficiency. It can be argued that there is no point in asking for a second answer after the first, but this situation can occur as subcomputations of a bigger computation and cannot be avoided in general.

Examining the code we find out easily the source of the redundant computation: the condition of function `part` includes two existential variables `U` and `V`. When the user asks for more solutions the backtracking mechanism looks for new values of the variables satisfying the conditions. But this is unnecessary because the rule already has returned **true** and cannot return any new value. The dynamic cut will avoid this redundant computation. Here is the same goal running after activating the dynamic cut optimization in *TOY*:

```

Toy>dnaPart (repeat 1000 adenine) (repeat 1000 thymine) 5
yes.      Elapsed time: 844 ms.
more solutions ? y
no.       Elapsed time: 0 ms.

```

Now the system detects automatically that there are no more possible solutions after the first one, reducing the 40 seconds to 0. The interested reader can find in [4] more experimental results. The experiments in that paper were tested introducing manually the code for the dynamic cut before the optimization was part of the system. However the results have been confirmed by the current implementation.

3.3 Dynamic conditions for the cut

From the previous examples one could consider that the cut can be introduced safely in the code of functions `multi` and `part` without taking into account any run-time test. But the cut also depends on dynamic conditions. There are two situations

that must be taken into account before applying the cut:

i) Variable bindings.

Consider the goal: $\text{multi } X \text{ zero} == R$, with X a logical variable. Using the program P_1 of Figure 1 this goal produces two answers: $\{ X \mapsto \text{zero}, R \mapsto \text{zero} \}$ and $\{ R \mapsto \text{zero} \}$. The first answer is obtained using the first rule for **multi** and the second answer through the second rule. Introducing a cut after the first answer would be unsafe; the second answer is not redundant, but gives new information w.r.t. the first one. As it includes no binding for X it can be interpreted as '*for every X , the equality $\text{multi } X \text{ zero} == \text{zero}$ holds*', and therefore subsumes the first answer.

ii) Non deterministic functions computed.

Suppose we include a new function **zeroOrOne** in the program P_1 of Figure 1 defined as:

zeroOrOne = zero **zeroOrOne** = s zero

Then a goal like $\text{multi zeroOrOne (s zero)} == R$ will return two answers: $\{ R \mapsto \text{zero} \}$ and $\{ R \mapsto \text{s zero} \}$. Introducing the cut after the first answer would be again unsafe. But in this case it is not because it prevents the use of the second rule, but because it would avoid the backtracking of the non-deterministic function **zeroOrOne** that leads to the application of the third rule of **multi**, yielding the second answer.

Therefore the cut must not take place if after obtaining the first result of the deterministic function any of the variables in the input arguments has been bound or a non-deterministic function has been computed. As we will see in the following paragraph the implementation generates a dynamic test for checking these conditions before introducing the cut.

4 Implementing the Dynamic Cut

4.1 Compiling programs into Prolog

The *TOY* compiler transforms *TOY* programs into Prolog programs following ideas described in [8]. A main component of the operational mechanism is the computation of *head normal forms* (hnf) for expressions. The translation scheme can be divided into three phases:

- 1) Higher order *TOY* programs are translated into programs in first order syntax.
- 2) Function calls $f(e_1, \dots, e_n)$ occurring in the first order *TOY* program rules are replaced by Prolog terms of the form $\text{susp}(f(e_1, \dots, e_n), R, S)$ called *suspensions*. The logical variable S is a flag which is bound to a concrete value, say hnf, once the suspension is evaluated. R contains the result of evaluating the function call. Its value is meaningful only if $S == \text{hnf}$ holds.
- 3) Finally the Prolog clauses are generated, adding code for *strict equality* and *hnf* (to compute head normal forms). Each n -ary function f is translated into a Prolog predicate $f(X_1, \dots, X_n, H)$. When computing a hnf for an unevaluated suspension $\text{susp}(f(X_1, \dots, X_n), R, S)$, a call $f(X_1, \dots, X_n, H)$ will occur in order to obtain in H the desired head normal form.

We are particularly interested in the third phase (code generation), since it will be affected by the introduction of dynamic cuts. Before looking more closely at this phase we need to introduce briefly our notation for definitional trees.

4.2 Definitional Trees in *TOY*

Before generating the code for any function the compiler builds its associated definitional tree. In our setting the definitional tree dt of a function f , can be of one of the following three forms:

- $dt(f) = f(\bar{t}_n) \rightarrow \text{case } X \text{ of } \langle c_1(\bar{X}_{m_1}) : dt_1; \dots; c_k(\bar{X}_{m_k}) : dt_k \rangle$, where X is the variable at position u in $f(\bar{t}_n)$ and $c_1 \dots c_k$ are constructor symbols, with dt_i a definitional tree for $i = 1 \dots k$.
- $dt(f) = f(\bar{t}_n) \rightarrow \text{or } \langle dt_1 \mid \dots \mid dt_k \rangle$, with dt_i a definitional tree for $i = 1 \dots k$.
- $dt(f) = f(\bar{t}_n) \rightarrow \text{try } (r \Leftarrow C)$, with $f \bar{t}_n = r \Leftarrow C$ corresponding to an instance of a program rule for f .

In each case we say that the tree has a case/or/try node at the root, respectively. A more precise definition together with the algorithm that produces a definitional tree from a function definition can be found in [8]. The only difference is that we do not allow 'multiple tries', i.e. *try* nodes including several program rules, replacing them by *or* nodes with multiple *try* child nodes, one for each rule included in the initial multiple *try*. The tree obtained by this modification is obviously equivalent and will be more suitable for our purposes. As an example of a definitional tree, consider again the definition of function *multi* in the program P_1 of Figure 1.

Its definitional tree, denoted as $dt(\text{multi})$, is defined in *TOY* as:

```
dt(multi) =  multi(A,B) → or (
               multi(A,B) → case A of
                   { zero      : multi (zero, B) → try (zero) % 1st rule
                     ; s(X)    : multi (s(X),B) → case B of
                                   { s(Y) : multi (s(X),s(Y)) → try (s (add X (add Y (multi(X,Y)))) ) % 3rd rule
                                   }
                   }
               | multi(A,B) → case B of { zero: multi (A,zero) → try (zero) } % 2nd rule
```

4.3 Definitional trees with cut

From the definitional tree dt of each function the *TOY* system generates a definitional tree with cut, dtc . Definitional trees with cut have the same structure as usual definitional trees. The only difference is that they rename some *or* and *try* nodes as *orCut* and *tryCut*, respectively. We define a function Γ transforming a definitional tree dt into its corresponding definitional tree with cut straightforwardly by distinguishing cases depending on the root node of dt :

- $\Gamma(f(\bar{t}_n) \rightarrow \text{case } X \text{ of } \langle c_1(X_{m_1}) : dt_1; \dots; c_k(X_{m_k}) : dt_k \rangle) =$
 $f(\bar{t}_n) \rightarrow \text{case } X \text{ of } \langle c_1(X_{m_1}) : \Gamma(dt_1); \dots; c_k(X_{m_k}) : \Gamma(dt_k) \rangle$
- $\Gamma(f(\bar{t}_n) \rightarrow \text{or } \langle dt_1 \mid \dots \mid dt_k \rangle) =$

$f(\bar{t}_n) \rightarrow \text{orCut } \langle \Gamma(dt_1) \mid \dots \mid \Gamma(dt_k) \rangle$, if f is deterministic.

- $\Gamma(f(\bar{t}_n) \rightarrow \text{or} \langle dt_1 \mid \dots \mid dt_k \rangle) = f(\bar{t}_n) \rightarrow \text{or} \langle \Gamma(dt_1) \mid \dots \mid \Gamma(dt_k) \rangle$, if f is non-deterministic.
- $\Gamma(f(\bar{t}_n) \rightarrow \text{try } (r \Leftarrow C) = f(\bar{t}_n) \rightarrow \text{tryCut } (r \Leftarrow C)$ if some existential variable occurs in C (i.e. some variable occurs in C but not in the rest of program rule).
- $\Gamma(f(\bar{t}_n) \rightarrow \text{try } (r \Leftarrow C) = f(\bar{t}_n) \rightarrow \text{try } (r \Leftarrow C)$ if no existential variable occurs in C .

For instance the dt of function `multi` displayed above is transformed into the following definitional tree with cut dct (denoted $dct(\text{multi})$):

```
dct(multi) =  multi(A,B) → orCut (
    multi(A,B) → case A of
      ( zero      : multi (zero, B) → try (zero) % 1st rule
      ; s(X)      : multi (s(X), B) → case B of
          ( s(Y) : multi (s(X), s(Y)) → try (s (add C (add D (multi(C,D)))))) ) % 3rd rule
      | multi(A,B) → case B of ( zero: multi (A, zero) → try (zero) ) % 2nd rule
```

Notice that the only difference corresponds to the root, which has been transformed into a `orCut` node because `multi` is a deterministic function.

4.4 Generating the code

Now we can describe the function $\text{prolog}(f, dct)$ which generates the code for a function f from its definitional tree with cut dct . The function definition depends on the node found at the root of dct . There are five possibilities:

Case 1. $dct = f(\bar{s}) \rightarrow \text{case } X \text{ of } \langle c_1(X_{m_1}) : dct_1; \dots; c_m(X_{m_k}) : dct_m \rangle$. Then:

$$\text{prolog}(g, dct) = \{g(\bar{s}, H) : - \text{hnf}(X, HX), g'(\bar{s}\sigma, H).\} \cup \\ \text{prolog}(g', dct_1) \cup \dots \cup \text{prolog}(g', dct_m)$$

where $\sigma = X/HX$ and g' is a new function symbol. The first call to `hnf` ensures that the position indicated by X is already in head normal form, and therefore can be used in order to distinguish the different alternatives.

Case 2. $dct = f(\bar{s}) \rightarrow \text{or} \langle dct_1 \mid \dots \mid dct_m \rangle$. Then:

$$\text{prolog}(g, dct) = \{g(\bar{s}, H) : - g_1(\bar{s}, H).\} \cup \dots \cup \{g(\bar{s}, H) : - g_m(\bar{s}, H).\} \cup \\ \text{prolog}(g_1, dct_1) \cup \dots \cup \text{prolog}(g_m, dct_m)$$

where g_1, \dots, g_m are new function symbols. In this case each new function symbol represents one of the non-deterministic choices.

Case 3. $dct = f(\bar{s}) \rightarrow \text{orCut} \langle dct_1 \mid \dots \mid dct_m \rangle$. Then

$$\begin{aligned}
prolog(g, dtc) = & \{g(\bar{s}, H) : -varlist(\bar{s}, V_s), g'(\bar{s}, H), \\
& (checkvarlist(V_s), !; \text{ true}).\} \cup \\
& \{g'(\bar{s}, H) : -\{g_1(\bar{s}, H).\} \cup \dots \cup \{g'(\bar{s}, H) : -g_m(\bar{s}, H).\} \cup \\
& prolog(g_1, dtc_1) \cup \dots \cup prolog(g_m, dtc_m)
\end{aligned}$$

where g' , g_1, \dots, g_m are new function symbols. Observe the differences with the case 2:

- A new function g' is used as an intermediate auxiliary function between g and the non-deterministic choices.
- g starts calling a predicate `varlist`. This predicate, whose definition is tedious but straightforward, returns in its second parameter V_s a list containing all the logical variables in the input parameters, including those used as flags for detecting the evaluation of suspensions of non-deterministic functions.
- After g' succeeds, i.e. after an or-branch has produced a result, the test for the dynamic cut is performed. This test, represented by predicate `checkvarlist`, checks if any of the variables in the list produced by `varlist` has been bound. This will mean that either an input logical variable has been bound or a non-deterministic function has been evaluated. In any of these cases the cut is avoided. Otherwise the dynamic cut, which is implemented as an ordinary Prolog cut, is safely performed. The definition of `checkvarlist` is simple:

`checkVarList([]).`

`checkVarList([X|Xs]):- var(X), \+varInList(X,Xs), checkVarList(Xs).`

The literal `\+varInList(X,Xs)`, checks if the variable X occurs twice in the list, detecting bindings among variables of the list.

Case 4. $dtc = \mathbf{try} \ (e \Leftarrow l_1 == r_1, \dots, l_n == r_n)$. Then

$$prolog(g, dtc) = \{g(\bar{s}, H) : -equal(l_1, r_1), \dots, equal(l_n, r_n), hnf(e, H).\}$$

If all equalities in the conditions are satisfied the program rule returns the head normal form of its right-hand side e .

Case 5. $dtc = \mathbf{tryCut} \ (e \Leftarrow l_1 == r_1, \dots, l_n == r_n)$. Then

$$\begin{aligned}
prolog(g, dtc) = & \{g(\bar{s}, H) : -varlist((\bar{s}, e), V_s), \\
& equal(l_1, r_1), \dots, equal(l_n, r_n), \\
& (checkvarlist(V_s), !; \text{ true}), \\
& hnf(e, H).\}
\end{aligned}$$

This case is similar to the case of the `orCut`. The main difference is that in this case we also collect the possible new variables of the right-hand side, because if the condition binds any of them the cut must be discarded.

4.5 Examples

Now we show the Prolog code generated by *TOY* for some of the function examples presented through the paper:

- Prolog code for function part of Figure 3:

```
part(A, B, C, true):- varList( [A, B, C ], Vs ),
                      equal(susp( ++, [ susp(++ , [D,A]),J]),B),
                      equal(susp(length, [A]), C),
                      (checkVarList(Vs), !; true).
```

This corresponds to the implementation of a *tryCut* node. In this example *varList* only looks for variables and non-deterministic functions in the parameters A, B and C, because the right-hand side of this rule is the ground term *true*.

- Prolog code for function multi of Figure 1

```
multi(A, B, H):- varList([A,B], Vs),
                 multi'(A, B, H),
                 (checkVarList(Vs), ! ; true ).

multi'(A, B, H):- hnf(A, F),
                  multi'_1(F, B, H).
multi'(A, B, zero):- hnf(B, zero).

multi'_1(zero, B, zero).
multi'_1(s(X), B,
s(susp(add, [X,susp(add, [Y,susp(multi, [X,Y])])]))):- hnf(B, s(Y)).
```

The code of this example corresponds to the implementation of an *orCut* node. The two branches are represented here by the two clauses for *multi'* (corresponding to function g' in the case 3 of the previous subsection). The cut is introduced if the first alternative, which corresponds to a *case* node with two possibilities, succeeds.

5 Conclusions

In this paper we have presented the implementation of the dynamic cut optimization in the Functional-Logic system *TOY*. The optimization improves dramatically the efficiency of the computations in the situations explained in the paper. Moreover, we claim that in practice it allows the use of some elegant and expressive function definitions that were disregarded due to their inefficiency up to now.

The cut is introduced automatically by the system following the next steps:

- (i) The deterministic functions of the program are detected using the non-ambiguity criterion. The correctness of the criterion is ensured by theorem 2.6. Also the user can indicate explicitly that any function is deterministic.
- (ii) The definitional tree associated to each program function is examined. The *or* nodes occurring in deterministic functions are labeled during this process as *or-cut* nodes. Also the *try* nodes corresponding to program rules including existential variables in the conditions are labeled as *try-cut* nodes.
- (iii) During the code generation the system will generate the dynamic cut code for *or-cut* and *try-cut* nodes. However the cut only will be performed if the dynamic conditions explained in subsection 3.3 are fulfilled.

We think that a similar scheme might also be used for incorporating the dynamic cut to the Prolog-based implementations of the Curry language [6].

Currently the dynamic cut must be turned on in *TOY* by typing the command `/cut` at the prompt. However, we have checked that the optimization produces almost no overhead in the cases where it cannot be applied, and we plan to provide it activated by default in the future versions of the system.

References

- [1] Antoy, S., *Definitional trees*, in: *Int. Conf. on Algebraic Logic Programming (ALP'92)*, number 632 in LNCS (1992), pp. 143–157.
- [2] Antoy, S., R.Echahed and M. Hanus, *A needed narrowing strategy*, *Journal of the ACM* **47** (2000), pp. 776–822.
- [3] Caballero, R. and F. López-Fraguas, *Dynamic-cut with definitional trees*, in: *Proceedings of the 6th International Symposium on Functional and Logic Programming, FLOPS 2002*, number 2441 in LNCS (2002), pp. 245–258.
- [4] Caballero, R. and F. López-Fraguas, *Improving deterministic computations in lazy functional logic languages*, *Journal of Functional and Logic Programming* **2003** (2003).
- [5] González-Moreno, J., M. Hortalá-González, F. López-Fraguas and M. Rodríguez-Artalejo, *An approach to declarative programming based on a rewriting logic*, *The Journal of Logic Programming* **40** (1999), pp. 47–87.
- [6] Hanus, M., *Curry: An Integrated Functional Logic Language (version 0.8.2. march 28, 2006)*, Available at: <http://www.informatik.uni-kiel.de/~curry/papers/report.pdf> (2006).
- [7] Henderson, F., Z. Somogyi and T. Conway, *Determinism analysis in the mercury compiler* (1996). URL citeseer.ist.psu.edu/henderson96determinism.html
- [8] Loogen, R., F. López-Fraguas and M. Rodríguez-Artalejo, *A demand driven computation strategy for lazy narrowing*, in: *Int. Symp. on Programming Language Implementation and Logic Programming (PLILP'93)*, number 714 in LNCS (1993), pp. 184–200.
- [9] Loogen, R., F. López-Fraguas and M. Rodríguez-Artalejo, *Toy: a multiparadigm declarative system*, in: *Int. Symp. RTA'99*, number 1631 in LNCS (1999), pp. 244–247.
- [10] Loogen, R. and S. Winkler, *Dynamic detection of determinism in functional-logic languages*, in: *Int. Symp. on Programming Language Implementation and Logic Programming (PLILP'91)*, number 528 in LNCS (1991), pp. 335–346.
- [11] Loogen, R. and S. Winkler, *Dynamic detection of determinism in functional logic languages*, in: J. Maluszynski and M. Wirsing, editors, *Programming Language Implementation and Logic Programming: Proc. of the 3rd International Symposium PLILP'91, Passau*, Springer, Berlin, Heidelberg, 1991 pp. 335–346.
- [12] Peña, R. and C. Segura, *Non-determinism analyses in a parallel-functional language*, *Journal of Logic Programming* **2004** (2005), pp. 67–100.
- [13] Sawamura, H. and T. Takeshima, *Recursive Unsolvability of Determinacy, Solvable Cases of Determinacy and Their Applications to Prolog Optimization*, in: *Proceedings of the Symposium on Logic Programming*, 1985, pp. 200–207.