

Linear Tabling Strategies and Optimizations

Neng-Fa Zhou

CUNY Brooklyn College & Graduate Center
zhou@sci.brooklyn.cuny.edu

Taisuke Sato

Tokyo Institute of Technology
sato@cs.titech.ac.jp

Yi-Dong Shen

*State Key Laboratory of Computer Science, Institute of Software,
Chinese Academy of Sciences*
ydshen@ios.ac.cn

submitted 11 April 2006; revised 4 October 2006, 17 May 2007; accepted 23 May 2007

Abstract

Recently there has been a growing interest of research in tabling in the logic programming community because of its usefulness in a variety of application domains including program analysis, parsing, deductive databases, theorem proving, model checking, and logic-based probabilistic learning. The main idea of tabling is to memorize the answers to some subgoals and use the answers to resolve subsequent variant subgoals. Early resolution mechanisms proposed for tabling such as OLDT and SLG rely on suspension and resumption of subgoals to compute fixpoints. Recently, the iterative approach named linear tabling has received considerable attention because of its simplicity, ease of implementation, and good space efficiency. Linear tabling is a framework from which different methods can be derived based on the strategies used in handling looping subgoals. One decision concerns when answers are consumed and returned. This paper describes two strategies, namely, *lazy* and *eager* strategies, and compares them both qualitatively and quantitatively. The results indicate that, while the lazy strategy has good locality and is well suited for finding all solutions, the eager strategy is comparable in speed with the lazy strategy and is well suited for programs with cuts. Linear tabling relies on depth-first iterative deepening rather than suspension to compute fixpoints. Each cluster of inter-dependent subgoals as represented by a top-most looping subgoal is iteratively evaluated until no subgoal in it can produce any new answers. Naive re-evaluation of all looping subgoals, albeit simple, may be computationally unacceptable. In this paper, we also introduce semi-naive optimization, an effective technique employed in bottom-up evaluation of logic programs to avoid redundant joins of answers, into linear tabling. We give the conditions for the technique to be safe (i.e. sound and complete) and propose an optimization technique called *early answer promotion* to enhance its effectiveness. Benchmarking in B-Prolog demonstrates that with this optimization linear tabling compares favorably well in speed with the state-of-the-art implementation of SLG.

KEYWORDS: Prolog, Semi-naive evaluation, Recursion, Tabling, Memoization, Linear tabling.

1 Introduction

The SLD resolution used in Prolog may not be complete or efficient for programs in the presence of recursion. For example, for a recursive definition of the transitive closure of a relation, a query may never terminate under SLD resolution if the program contains left-recursion or the graph represented by the relation contains cycles even if no rule is left-recursive. For a natural definition of the Fibonacci function, the evaluation of a subgoal under SLD resolution spawns an exponential number of subgoals, many of which are variants. The lack of completeness and efficiency in evaluating recursive programs is problematic: novice programmers may lose confidence in writing declarative programs that terminate and real programmers have to reformulate a natural and declarative formulation to avoid these problems, resulting in cluttered programs.

Tabling (Tamaki and Sato 1986; Warren 1992) is a technique that can get rid of infinite loops for bounded-term-size programs and redundant computations in the execution of recursive programs. The main idea of tabling is to memorize the answers to subgoals and use the answers to resolve their variant descendants. Tabling helps narrow the gap between declarative and procedural readings of logic programs. It not only is useful in the problem domains that motivated its birth, such as program analysis (Dawson et al. 1996), parsing (Eisner et al. 2004; Johnson 1995; Warren 1999), deductive databases (Liu 1999; Ramakrishnan and Ullman 1995; Sagonas et al. 1994), and theorem proving (Nielsen et al. 2004; Pientka 2003), but also has been found essential in several other problem domains such as model checking (Ramakrishnan 2002) and logic-based probabilistic learning (Sato and Kameya 2001; Zhou et al. 2003). This idea of caching previously calculated solutions, called *memoization*, was first used to speed up the evaluation of functions (Michie 1968). OLDT (Tamaki and Sato 1986) is the first resolution mechanism that accommodates the idea of tabling in logic programming and XSB is the first Prolog system that successfully supports tabling (Sagonas and Swift 1998). Tabling has become a practical technique thanks to the availability of large amounts of memory in computers. It has become an embedded feature in a number of other logic programming systems such as B-Prolog (Zhou et al. 2000; Zhou et al. 2004), Mercury (Somogyi and Sagonas 2006), TALS (Guo and Gupta 2001), and YAP (Rocha et al. 2005b).

OLDT, and SLG (Chen and Warren 1996) alike, is non-linear in the sense that the state of a consumer must be preserved before execution backtracks to its producer. This non-linearity requires freezing stack segments (Sagonas and Swift 1998) or copying stack segments into a different area (Demoen and Sagonas 1999) before backtracking takes place. Linear tabling is an alternative tabling scheme (Shen et al. 2001; Zhou et al. 2000; Zhou and Sato 2003; Zhou et al. 2004). The main idea of linear tabling is to use iterative computation of looping subgoals rather than suspension and resumption of them as is done in OLDT to compute fixpoints. This basic idea dates back to the ET* algorithm (Dietrich 1987). The DRA method proposed in (Guo and Gupta 2001) is based on the same idea but employs different strategies for handling looping subgoals and clauses. In linear tabling, a cluster of inter-dependent subgoals as represented by a *top-most looping subgoal* is iteratively

evaluated until no subgoal in it can produce any new answers. Linear tabling is relatively easy to implement on top of a stack machine thanks to its linearity, and is more space efficient than OLDT since the states of subgoals need not be preserved.

Linear tabling is a framework from which different methods can be derived based on the strategies used in handling looping subgoals. One decision concerns when answers are consumed and returned. The *lazy* strategy postpones the consumption of answers until no answers can be produced. It is in general space efficient because of its locality and is well suited for all-solution search programs. The *eager* strategy, in contrast, prefers answer consumption and return over production. It is well suited for programs with cuts. These two strategies have been compared in SLG-WAM as two scheduling strategies called *local* and *single-stack* (Freire et al. 1998). This paper gives a comprehensive analysis of these two strategies and compares their performance experimentally.

Linear tabling relies on iterative evaluation of top-most looping subgoals to compute fixpoints. Naive re-evaluation of all looping subgoals may be computationally expensive. *Semi-naive optimization* is an effective technique used in bottom-up evaluation of Datalog programs (Bancilhon and Ramakrishnan 1986; Ullman 1988). It avoids redundant joins by ensuring that the join of the subgoals in the body of each rule must involve at least one new answer produced in the previous round. The impact of semi-naive optimization on top-down evaluation had been unknown before (Zhou et al. 2004). In this paper, we also propose to introduce semi-naive optimization into linear tabling. We have made efforts to properly tailor semi-naive optimization to linear tabling. In our semi-naive optimization, answers for each tabled subgoal are divided into three regions as in bottom-up evaluation, but answers are consumed sequentially until exhaustion not incrementally as in bottom-up evaluation so that answers produced in a round are consumed in the same round. We have found that incremental consumption of answers does not fit linear tabling since it may require more iterations to reach fixpoints. Moreover, consuming answers incrementally may cause redundant consumption of answers. We further propose a technique called *early promotion* of answers to reduce redundant consumption of answers. Our benchmarking shows that this technique gives significant speed-ups to some programs.

An efficient tabling system has been implemented in B-Prolog,¹ in which the lazy strategy is employed by default but the eager strategy can be used through declarations for subgoals that are in the scopes of cuts or are not required to return all the answers. Our tabling system not only consumes considerably less stack space than XSB for some programs but also compares favorably well in speed with XSB.

The theoretical framework of linear tabling is given in (Shen et al. 2001). The main objective of this paper is to propose evaluation strategies and their optimizations for linear tabling. The remainder of the paper is structured as follows: In the next section we define the terms used in this paper. In Section 3 we give the linear

¹ www.bprolog.com

tabling framework and the two answer consumption strategies. In Section 4 we introduce semi-naive optimization into linear tabling and prove its completeness. In Section 5 we describe the implementation of our tabling system and also show how to implement semi-naive optimization. In Section 6 we compare the tabling strategies experimentally, evaluate the effectiveness of semi-naive optimization, and also compare the performance of B-Prolog with XSB. In Section 7 we survey the related work and in Section 8 we conclude the paper.

2 Preliminaries

In this section we give the definitions of the terms to make this paper as much self-contained as possible. The reader is referred to (Lloyd 1988) for a description of SLD resolution. In this paper, we always assume the top-down strategy for selecting clauses and the left-to-right computation rule.

Let P be a program. Tabled predicates in P are explicitly declared and all the other predicates are assumed to be non-tabled. A subgoal of a tabled predicate is called a *tabled subgoal*. Tabled predicates are transformed into a form that facilitates execution: each rule ends with a dummy subgoal named $memo(H)$ where H is the head, and each tabled predicate contains a dummy ending rule whose body contains only one subgoal named $check_completion(H)$. For example, given the definition of the transitive closure of a relation,

```
:-table p/2.
p(X,Y):-p(X,Z),e(Z,Y).
p(X,Y):-e(X,Y).
```

The transformed predicate is as follows:

```
p(X,Y):-p(X,Z),e(Z,Y),memo(p(X,Y)).
p(X,Y):-e(X,Y),memo(p(X,Y)).
p(X,Y):-check_completion(p(X,Y)).
```

A table is used to record subgoals and their answers. For each subgoal and its variants, there is an entry in the table that stores the state of the subgoal (e.g., complete or not) and an answer table for holding the answers generated for the subgoal. Initially, the answer table is empty.

Definition 1

Let t_1 and t_2 be two terms with no shared variables. The term t_1 *subsumes* t_2 if there exists a substitution θ such that $t_1\theta=t_2$. The two terms t_1 and t_2 are called *variants* if they subsume each other.

Definition 2

Let $G = (A_1, A_2, \dots, A_k)$ be a goal. The first subgoal A_1 is called the *selected subgoal* of the goal. G' is *derived* from G by using a tabled *answer* F if there exists a unifier θ such that $A_1\theta = F$ and $G' = (A_2, \dots, A_k)\theta$. G' is *derived* from G by using a rule " $H : -B_1, \dots, B_m$ " if $A_1\theta = H\theta$ and $G' = (B_1, \dots, B_m, A_2, \dots, A_k)\theta$. A_1 is said to be the *parent* of B_1, \dots, B_m . The relation *ancestor* is defined recursively from the parent relation.

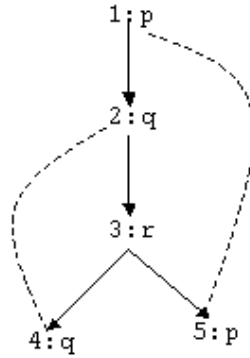


Fig. 1. A top-most looping subgoal.

Definition 3

A tabled subgoal that occurs first in the construction of an SLD tree is called a *pioneer*, and all subsequent variants are called *followers* of the pioneer. Let G_0 be a given goal, and $G_0 \Rightarrow G_1 \Rightarrow \dots \Rightarrow G_n$ be a *derivation* where each goal is derived from the goal immediately preceding it. Let $G_i \Rightarrow \dots \Rightarrow G_j$ be a sub-sequence of the derivation where $G_i = (A\dots)$ and $G_j = (A'\dots)$. The sub-sequence forms a *loop* if A and A' are variants. The subgoals A and A' are called *looping subgoals*. In particular, A is called the *pioneer looping subgoal* and A' is called the *follower looping subgoal* of the loop.

Notice that the pioneer and follower looping subgoals are not required to have the ancestor-descendent relationship, and thus a derivation that contains two variant subgoals may not be a *real* loop. Consider, for example, the goal “ $p(X), p(Y)$ ” where p is defined by facts. The derivation “ $p(X), p(Y) \Rightarrow p(Y)$ ” is treated as a loop although the selected subgoal $p(Y)$ in the second goal is not a descendant of $p(X)$.

Definition 4

A subgoal A is said to be *dependent* on another subgoal A' if A' occurs in a derived goal from A , i.e., $A \Rightarrow \dots \Rightarrow (A'\dots)$. Two subgoals are said to be *inter-dependent* if they are dependent on each other. Inter-dependent subgoals constitute a *cluster*, which is called a *strongly connected component* elsewhere (Sagonas and Swift 1998). A subgoal in a cluster is called the *top-most* subgoal of the cluster if none of its ancestors is included in the cluster.

Unless a cluster contains only a single subgoal, its top-most subgoal must also be a looping subgoal. For example, the subgoals at the nodes in the SLD tree in Figure 1 constitute a cluster and the subgoal p at node 1 is the top-most looping subgoal of the cluster.

3 Linear Tabling and Answer Consumption Strategies

Linear tabling takes a transformed program and a goal, and tries to find a path in the SLD tree that leads to an empty goal. The primitive $table_start(A)$ is executed when a tabled subgoal A is encountered. Just as in SLD resolution, linear tabling explores the SLD tree in a depth-first fashion, taking special actions when $table_start(A)$, $memo(A)$, and $check_completion(A)$ are encountered. Backtracking is done in exactly the same way as in SLD resolution. When the current path reaches a dead end, meaning that no action can be taken on the selected subgoal, execution backtracks to the latest previous goal in the path and continues with an alternative branch. When execution backtracks to the top-most looping subgoal of a cluster, however, we cannot fail the subgoal even after all the alternative clauses have been tried. In general, the evaluation of a top-most looping subgoal must be iterated until its fixpoint is reached. We call each iteration of a top-most looping subgoal a *round*.

Various linear tabling methods can be devised based on the framework. A linear tabling method comprises strategies used in the three primitives: $table_start(A)$, $memo(A)$, and $check_completion(A)$. In linear tabling, a pioneer subgoal has two roles: one is to produce answers into the table and the other is to return answers to its parent through its variables. Different strategies can be used to produce and return answers. The *lazy strategy* gives priority to answer production and the *eager strategy* prefers answer consumption over production. In the following we define the three primitives in both strategies.

3.1 The lazy strategy

The lazy strategy postpones the consumption of answers until no answers can be produced. In concrete, for top-most looping subgoals no answer is returned until they are complete, and for other pioneer subgoals answers are consumed only after all the rules have been tried.

3.1.1 $table_start(A)$

This primitive is executed when a tabled subgoal A is encountered. The subgoal A is registered into the table if it is not registered yet. If A 's state is *complete* meaning that A has been completely evaluated before, then A is resolved by using the answers in the table.

If A is a pioneer, meaning that it is encountered for the first time in the current path, then different actions are taken depending on A 's state. If A 's state is *evaluated* meaning that A has occurred before in a different path during the current round, then it is resolved by using answers. Otherwise, if A has never occurred before during the current round, it is resolved by using rules. In this way, a pioneer subgoal needs to be evaluated only once in each round.

If A is a follower of some ancestor A_0 , meaning that a loop has been encountered,²

² As to be discussed later, A_0 must be an ancestor of A under the lazy strategy.

then it is resolved by using the answers in the table. After the answers are exhausted, A fails. Failing A is unsafe in general since it may have not returned all of its possible answers. For this reason, the top-most looping subgoal of the cluster of A needs be iterated until no new answer can be produced.

3.1.2 *memo(A)*

This primitive is executed when an answer is found for the tabled subgoal A . If the answer A is already in the table, then just fail; otherwise fail after the answer is added into the table. The failure of *memo* postpones the return of answers until all rules have been tried.

3.1.3 *check_completion(A)*

This primitive is executed when the subgoal A is being resolved by using rules and the dummy ending rule is being tried. If A has never occurred in a loop, then A 's state is set to *complete* and A is failed after all the answers are consumed.

If A is a top-most looping subgoal, we check if any new answers are produced during the last iteration of the cluster under A . If so, A is re-evaluated by calling *table_start(A)* after all the dependent subgoals's states are initialized. Otherwise, if no new answer is produced, A is resolved by using answers after its state and all its dependent subgoals' states are set to *complete*. Notice that a top-most looping subgoal does not return any answers until it is complete.

If A is a looping subgoal but not a top-most one, A will be resolved by using answers after its state is set to *evaluated*. Notice that A 's state cannot be set to *complete* since A is contained in a loop whose top-most subgoal has not been completely evaluated. For example, in Figure 1, q reaches its fixpoint only after the top-most looping subgoal p reaches its fixpoint.

As described in the definition of *table_start(A)*, an *evaluated* subgoal is never evaluated using rules again in the same round. This optimization is called *subgoal optimization* in (Zhou and Sato 2003). If evaluating a subgoal produces some new answers then the top-most looping subgoal will be re-evaluated and so will the subgoal; and if evaluating a subgoal does not produce any new answer, then evaluating it again in the same round would not produce any new answers either. Therefore, the subgoal optimization is safe.

3.1.4 *Example*

Consider the following program, where $p/2$ is tabled, and the query $p(a, Y0)$.

$p(X, Y) :- p(X, Z), e(Z, Y), memo(p(X, Y)).$ (p1)

$p(X, Y) :- e(X, Y), memo(p(X, Y)).$ (p2)

$p(X, Y) :- check_completion(p(X, Y)).$ (p3)

$e(a, b).$

$e(b, c).$

The following shows the steps that lead to the production of the first answer:

```

1: p(a, Y0)
   ↓ apply p1
2: p(a, Z1), e(Z1, Y0), memo(p(a, Y0))
   loop found, backtrack to goal 1
1: p(a, Y0)
   ↓ apply p2
3: e(a, Y0), memo(p(a, Y0))
   ↓ apply e(a,b)
4: memo(p(a, b))
   ↓ add answer p(a,b)

```

After the answer $p(a, b)$ is added into the table, $\text{memo}(p(a, b))$ fails. The failure forces execution to backtrack to $p(a, Y0)$.

```

1: p(a, Y0)
   ↓ apply p3
5: check_completion(p(a, Y0))

```

Since $p(a, Y0)$ is a top-most looping subgoal which has not been completely evaluated yet, $\text{check_completion}(p(a, Y0))$ does not consume the answer in the table but instead starts re-evaluation of the subgoal.

```

1: p(a, Y0)
   ↓ apply p1
6: p(a, Z1), e(Z1, Y0), memo(p(a, Y0))
   ↓ use answer p(a,b)
7: e(b, Y0), memo(p(a, Y0))
   ↓ apply e(b,c)
8: memo(p(a, c))

```

When the follower $p(a, Z1)$ is encountered this time, it consumes the answer $p(a, b)$. The current path leads to the second answer $p(a, c)$. On backtracking, the goal numbered 6 becomes the current goal.

```

6: p(a, Z1), e(Z1, Y0), memo(p(a, Y0))
   ↓ use answer p(a,c)
9: e(c, Y0), memo(p(a, Y0))

```

Goal 9 fails. Execution backtracks to the top goal and tries the clause p3 on it.

```

1: p(a, Y0)
   ↓ apply p3
10: check_completion(p(a, Y0))

```

Since the new answer $p(a, c)$ is produced in the last round, the top-most looping subgoal $p(a, Y0)$ needs to be evaluated again. The next round produces no new answer and thus the subgoal's state is set to *complete*. After that the top-most subgoal returns the answers $p(a, b)$ and $p(a, c)$.

3.1.5 Properties of the lazy strategy

Under the lazy strategy, answers are not returned immediately after they are produced but are returned via the table after all clauses are tried. No answer is returned for a top-most looping subgoal until the subgoal is complete.

All loops are guaranteed to be real: for any loop $G_i = (A \dots) \Rightarrow \dots \Rightarrow G_j = (A' \dots)$ where A and A' are variants, A must be an ancestor of A' . Because each cluster of inter-dependent subgoals is completely evaluated before any answers are returned to outside of the cluster, the lazy strategy has good locality and is thus suited for finding all solutions. For example, when the subgoal $p(Y)$ is encountered in the goal “ $p(X), p(Y)$ ”, the subtree for $p(X)$ must have been explored completely and thus needs not be saved for evaluating $p(Y)$.

The cut operator cannot be handled efficiently under the lazy strategy. The goal “ $p(X), !, q(X)$ ” produces all the answers for $p(X)$ even though only one is needed.

3.2 The eager strategy

The eager strategy prefers answer consumption and return over production. For a pioneer, answers are used first and rules are used only after all available answers are exhausted, and moreover a new answer is returned to its parent immediately after it is added into the table. The following describes how the three primitives behave under the eager strategy.

3.2.1 *table_start(A)*

Just as in the lazy strategy, A is registered if it is not registered yet. A is resolved by using the tabled answers if A is complete or A is a follower of some former variant subgoal. If A is a pioneer, being encountered for the first time in the current round, it is resolved by using answers first, and then rules after all existing answers are exhausted.

3.2.2 *memo(A)*

If the answer A is already in the table, then this primitive fails; otherwise, this primitive succeeds after adding the answer A into the table. Notice that A is returned immediately after it is added into the table. If A is not new, then it must have been returned before.

3.2.3 *check_completion(A)*

If A is a top-most looping subgoal, just as in the lazy strategy, we check whether any new answers are produced during the last iteration of A . If so, A is evaluated again by calling *table_start(A)*. Otherwise, if no new answer is produced, this primitive fails after A 's and all its dependent subgoals' states are set to *complete*. If A is a looping subgoal but not a top-most one, this primitive fails after A 's state is set to *evaluated*. An *evaluated* subgoal is never evaluated using rules

again in the same round. Notice that unlike under the lazy strategy, the primitive $check_completion(A)$ never returns any answers under the eager strategy. As described above, all the available answers must have been returned by $table_start(A)$ and $memo(A)$ by the time $check_completion(A)$ is executed.

3.2.4 Example

Because of the need to re-evaluate a top-most looping subgoal, redundant solutions may be observed for a query. Consider, for example, the following program and the query “ $p(X), p(Y)$ ”.

```
p(1) :- memo(p(1)).           (r1)
p(2) :- memo(p(2)).           (r2)
p(X) :- check_completion(p(X)). (r3)
```

The following derivation steps lead to the return of the first solution $(1,1)$ for (X,Y) .

```
1: p(X), p(Y)
   ↓ use r1
2: memo(p(1)), p(Y)
   ↓ add answer p(1)
3: p(Y)
   ↓ loop found, use answer p(1)
```

When the subgoal $p(Y)$ is encountered, it is treated as a follower and is resolved using the tabled answer $p(1)$. After that the first solution $(1,1)$ is returned to the top query. When execution backtracks to $p(Y)$, it fails since it is a follower and no more answer is available in the table. Execution backtracks to $p(X)$, which produces and adds the second answer $p(2)$ into the table.

```
1: p(X), p(Y)
   ↓ use r2
4: memo(p(2)), p(Y)
   ↓ add answer p(2)
5: p(Y)
   ↓ use answer p(1)
```

When $p(Y)$ is encountered this time, there are two answers $p(1)$ and $p(2)$ in the table. So the next two solutions returned are $(2,1)$ and $(2,2)$. When execution backtracks to goal 1, the dummy ending rule is applied.

```
1: p(X), p(Y)
   ↓ use r3
6: check_completion(p(X)), p(Y)
```

Since new answers are added into the table during this round, the subgoal $p(X)$ needs to be evaluated again, first using answers and then using rules. The second round produces no answer but returns the four solutions $(1,1)$, $(1,2)$, $(2,1)$ and $(2,2)$ among which only $(1,2)$ has not been observed before.

3.2.5 Properties of the eager strategy

Since answers are returned eagerly, a pioneer and a follower may not have an ancestor-descendant relationship. Because of the existence of this kind of *fake* loops and the necessity of iterating the evaluation of top-most looping subgoals, redundant solutions may be observed. In the previous example, the solutions $(1,1)$, $(2,1)$ and $(2,2)$ are each observed twice. Provided that the top-most looping subgoal $p(X)$ did not return the answer $p(1)$ again in the second round, the solution $(1,2)$ would have been lost.

The eager strategy is more suited than the lazy strategy for single-solution search. For certain applications such as planning it is unreasonable to find all answers either because the set is infinite or because only one answer is needed. For these applications the eager strategy is more effective than the lazy one. Cuts are handled more efficiently under the eager strategy.

4 Semi-naive Optimization

The basic linear tabling framework described in the previous section does not distinguish between new and old answers. The problem with this naive method is that it redundantly joins answers of subgoals that have been joined in early rounds. Semi-naive optimization (Ullman 1988) reduces the redundancy by ensuring that at least one new answer is involved in the join of the answers for each rule. In this section, we introduce semi-naive optimization into linear tabling and identify sufficient conditions for it to be complete. We also propose a technique called *early answer promotion* to further avoid redundant consumption of answers. This optimization works with both the lazy and eager strategies.

4.1 Preparation

To make semi-naive optimization possible, we divide the answer table for each tabled subgoal into three regions:

<i>old</i>	<i>previous</i>	<i>current</i>
------------	-----------------	----------------

The names of the regions indicate the rounds during which the answers in the regions are produced: *old* means that the answers were produced before the previous round, *previous* the answers produced during the previous round, and *current* the answers produced in the current round. The answers stored in *previous* and *current* are said to be *new*. Before each round is started, answers are promoted accordingly: *previous* answers become *old* and *current* answers become *previous*.

In our optimization, answers are consumed *sequentially*. For a subgoal, either all the available answers or only new answers are consumed. This is unlike in bottom-up evaluation where answers are consumed *incrementally*, i.e., answers produced in a round are not consumed until the next round. As will be discussed later, incremental consumption of answers as is done in bottom-up evaluation does avoid certain redundant joins but does not fit linear tabling since it may require more rounds to reach fixpoints.

A predicate p *calls* a predicate q if: (1) if q occurs in the body of at least one rule in the definition of p (p calls q *directly*); or (2) q does not occur in the body of any rule in the definition of p but there exists a predicate in the body of a rule in the definition of p that calls q (p calls q *indirectly*). The calling relationship constitutes a graph called a *call graph*.

For a given program, we find a level mapping from the predicate symbols in the program to the set of integers to represent the *call graph* of the program. Let m be a level mapping. We extend the notation to assume that $m(p(\dots)) = m(p/n)$ for any subgoal $p(\dots)$ of arity n .

Definition 5

For a given program, a level mapping m represents the *call graph* if: for each rule “ $H:-A_1, \dots, A_n$ ” in the program, $m(H) > m(A_i)$ iff the predicate of A_i does not call (either directly or indirectly) the predicate of H , and $m(H) = m(A_i)$ iff the predicates of H and A_i call each other.

The level mapping as defined divides predicates in a program into several strata. The predicate at each stratum depends only on those on the lower strata. The level mapping is an abstract representation of the dependence relationship of the subgoals that may occur in execution. If two subgoals A and A' occur in a loop, then it is guaranteed that $m(A) = m(A')$.

Definition 6

Let “ $H:-A_1, \dots, A_k, \dots, A_n$ ” be a rule in a program and m be the level mapping that represents the call graph of the program. A_k is called the *last depending subgoal* of the rule if $m(A_k) = m(H)$ and $m(H) > m(A_i)$ for $i > k$.

The last depending subgoal A_k is the last subgoal in the body that may depend on the head to become complete. Thus, when the rule is re-executed on a subgoal, all the subgoals to the right of A_k that have occurred before must already be complete.

Definition 7

Let “ $H:-A_1, \dots, A_n$ ” be a rule in a program and m be a level mapping that represents the call graph of the program. If there is no depending subgoal in the body, i.e., $m(H) > m(A_i)$ for $i = 1, \dots, n$, then the rule is called a *base rule*.

4.2 Semi-naive optimization

Theorem 1

Let “ $H:-A_1, \dots, A_k, \dots, A_n$ ” be a rule where A_k is the last depending tabled subgoal, and C be a subgoal that is being resolved by using the rule in an iteration of a top-most looping subgoal T . For a combination of answers of A_1, \dots , and A_{k-1} , if C has occurred in an early round and the combination does not contain any new answers, then it is safe to let A_k consume new answers only.

Proof

Because A_k is the last depending subgoal, the subgoals A_{k+1}, \dots , and A_n must have been completely evaluated when C is re-evaluated. Let $A_{k_{old}}$ and $A_{k_{new}}$ be the *old* and *new* answers of the subgoal A_k , respectively. For a combination of answers of A_1, \dots , and A_{k-1} , if the combination does not contain new answers then the join of the combination and $A_{k_{old}}$ must have been done and all possible answers for C that can result from the join must have been produced during the previous round because the subgoal C has been encountered before. Therefore only new answers in $A_{k_{new}}$ should be used. \square

Corollary 1

Base rules need not be considered in the re-evaluation of any subgoals.

Semi-naive optimization would be unsafe if it were applied to new subgoals that have never been encountered before. The following example illustrates this possibility:

```
?- p(X,Y).

:-table p/2.
p(X,Y) :- p(X,Z),q(Z,Y). (C1)
p(b,c) :- p(X,Y). (C2)
p(a,b). (C3)

:-table q/2.
q(c,d) :- p(X,Y),t(X,Y). (C4)

t(a,b). (C5)
```

In the first round of $p(X,Y)$ the answer $p(a,b)$ is added to the table by C3, and in the second round the rule C2 produces the answer $p(b,c)$ by using the answer produced in the first round. In the third round, the rule C1 generates a new subgoal $q(c,Y)$ after $p(X,Z)$ consumes $p(b,c)$. If semi-naive optimization were applied to $q(c,Y)$, then the subgoal $p(X,Y)$ in C4 could consume only the new answer $p(b,c)$ and the third answer $p(b,d)$ would be lost.

4.3 Analysis

Semi-naive optimization can lower the complexity of evaluation for some programs. Consider the following example created by David S. Warren:³

```
:-table p/2.
p(X,Y) :- p(X,Z),c(Z,a,Y).
p(X,Y) :- p(X,Z),c(Z,b,Y).
p(X,X).
```

which detects if a given string represented as facts $c(I, S, J)$ ($J = I + 1, S = a$ or $S = b$) is a sentence of the regular expression $(a|b)^*$. For a string $(ab)^{n/2}$, the query $p(0, n)$ needs $n/2$ rounds to reach the fixpoint. With semi-naive optimization, the variants of $p(X, Z)$ in the bodies consume only new answers, and therefore the program takes linear time. Without semi-naive optimization, however, the program would take $O(n^2)$ time since the variants of $p(X, Z)$ would consume all existing answers.

In our semi-naive optimization, answers produced in the current round are consumed immediately rather than postponed to the next round as in the bottom-up version, and answers are promoted each time a new round is started. This way of consuming and promoting answers may cause certain redundancy.

Consider the conjunction (P, Q) . Assume Q_o , Q_p , and Q_c are the sets of answers in the three regions (respectively, *old*, *previous*, and *current*) of the subgoal Q when Q is encountered in round i . Assume also that P had been complete before round i and P_a is the set of answers. The join $P_a \bowtie (Q_p \cup Q_c)$ is computed for the conjunction in round i . Assume Q'_o , Q'_p , and Q'_c are the sets of answers in the three regions when Q is encountered in round $i+1$. Since answers are promoted before round $i+1$ is started, we have:

$$\begin{aligned} Q'_o &= Q_o \cup Q_p \\ Q'_p &= Q_c \cup \alpha \end{aligned}$$

where α denotes the new answers produced for Q after the conjunction (P, Q) in round i . When the conjunction (P, Q) is encountered in round $i+1$, the following join is computed.

$$P_a \bowtie (Q'_p \cup Q'_c) = P_a \bowtie (Q_c \cup \alpha \cup Q_c')$$

Notice that the join $P_a \bowtie Q_c$ is computed in both round i and $i+1$.

We could allow last depending subgoals to consume answers incrementally as is done in bottom-up evaluation, but doing so may require more rounds to reach fixpoints. Consider the following example, which is the same as the one shown above but has a different ordering of clauses:

```
?- p(X,Y).

:-table p/2.
```

³ Personal communications.

```

p(a,b).                (C1)
p(b,c) :- p(X,Y).      (C2)
p(X,Y) :- p(X,Z),q(Z,Y). (C3)

:-table q/2.
q(c,d) :- p(X,Y),t(X,Y). (C4)

t(a,b).                (C5)

```

In the first round, C1 produces the answer $p(a,b)$. When C2 is executed, the subgoal in the body cannot consume $p(a,b)$ since it is produced in the current round. Similarly, C3 produces no answer either. In the second round, $p(a,b)$ is moved to the *previous* region, and thus can be consumed. C2 produces a new answer $p(b,c)$. When C3 is executed, no answer is produced since $p(b,c)$ cannot be consumed. In the third round, $p(a,b)$ is moved to the *old* region, and $p(b,c)$ is moved to the *previous* region. C3 produces the third answer $p(b,d)$. The fourth round produces no new answer and confirms the completion of the computation. So in total four rounds are needed to compute the fixpoint. If answers produced in the current round are consumed in the same round, then only two rounds are needed to reach the fixpoint.

4.4 Early promotion of answers

As discussed above, sequential consumption of answers may cause redundant joins. In this subsection, we propose a technique called *early promotion* of answers to reduce the redundancy.

Definition 8

Let Q be the first follower that exhausts its answers in the current round. Then all the answers of Q in the *current* region are promoted to the *previous* region once being consumed by Q .

Consider again the conjunction (P,Q) where Q is the first follower that exhausts its answers. The answers in the current region Q_c are promoted to the *previous* region after Q has consumed all its answers in round i . By doing so, the join $P_a \bowtie Q_c$ will not be recomputed in round $i+1$ since Q_c must have been promoted to the *old* region in round $i+1$.

Consider, for example, the following program:

```

?- p(X,Y).

:-table p/2.
p(a,b).                (C1)
p(b,c) :- p(X,Y).      (C2)

```

Before C2 is executed in the first round, $p(a,b)$ is in the *current* region. Executing C2 produces the second answer $p(b,c)$. Since the subgoal $p(X,Y)$ in C2 is the first

follower that exhausts its answers in the current round, it is qualified to promote its answers. So the answers $p(\mathbf{a}, \mathbf{b})$ and $p(\mathbf{b}, \mathbf{c})$ are moved from the *current* region to the *previous* region immediately after being consumed by $p(\mathbf{X}, \mathbf{Y})$. As a result, the potential redundant consumption of these answers by $p(\mathbf{X}, \mathbf{Y})$ is avoided in the second round since they will all be transferred to the *old* region before the second round starts.

Theorem 2

Early promotion does not lose any answers.

Proof

First note that although answers are tabled in three disjoint regions, all tabled answers will be consumed except for some last depending subgoals that would skip the answers in their *old* regions (see Theorem 1). Assume, on the contrary, that applying early promotion loses answers. Then there must be a last depending subgoal A_k in a rule “ $H:-A_1, \dots, A_k, \dots, A_n$ ” and a tabled answer A for A_k such that A has been moved to the *old* region before being consumed by A_k so that A will never be consumed by A_k . Assume A is produced in round i by a variant of A_k . We distinguish between the following two cases:

1. The last depending subgoal A_k is not selected in round i . In round $j(j > i)$, A_k is selected either because H is new or some $A_s(s < k)$ consumes a new answer. By Theorem 1, A_k will consume all answers in the three regions, including the answer A .
2. Otherwise, A must be produced by A_k itself or a variant subgoal of A_k that is selected either *before* or *after* A_k in round i . If A is produced by A_k itself or *before* A_k is selected, then the answer will be consumed by A_k since promoted answers will remain new by the end of the round. If A is produced by a variant *after* A_k is selected, then the answer cannot be promoted because A_k exhausts its answers *before* the variant. In this case, the answer A will remain new in the next round and will thus be consumed by A_k .

Both of the above two cases contradict our assumption. The proof then concludes.

□

5 Implementation

Changes to the Prolog machine ATOAM (Zhou 1996) are needed to implement linear tabling. In this section we describe the changes to the data structures and the instruction set. To make the paper self-contained, we first give an overview of the ATOAM architecture.

5.1 An overview of ATOAM

The ATOAM uses all the data areas used by the WAM. The *heap* stores terms created during execution. The register **H** points to the top of the heap. The *trail* stack stores updates that must be undone upon backtracking. The register **T** points to

the top of the trail stack. The *control* stack stores frames associated with predicate calls.

Unlike in the WAM where arguments are passed through argument registers, arguments in the ATOAM are passed through stack frames and only one frame is used for each predicate call. Each time a predicate is invoked by a call, a frame is placed on top of the local stack unless the frame currently at the top can be reused. Frames for different types of predicates have different structures. For standard Prolog, a frame is either *determinate* or *nondeterminate*. A nondeterminate frame is also called a *choice point*. The register **AR** points to the current frame and the register **B** points to the latest *choice point*.

A determinate frame has the following structure:

A1 . . An	Arguments
AR	Pointer to the parent frame
CP	Continuation program pointer
BTM	Bottom of the frame
TOP	Top of the frame
Y1 . . Ym	Local variables

Where **BTM** points to the bottom of the frame, i.e., the slot for the first argument, and **TOP** points to the top of the frame, i.e., the slot just next to that for the last local variable⁴. The **TOP** register points to the next available slot on the stack. The **BTM** slot is not in the original version (Zhou 1996). This slot is introduced for supporting garbage collection and co-routining. The **AR** register points to the **AR** slot of the current frame. Arguments and local variables are accessed through offsets with respect to the **AR** slot. An argument or a local variable is denoted as $y(I)$ where I is the offset. Arguments have positive offsets and local variables have negative offsets. It is the caller's job to place the arguments and fill in the **AR**, and **CP** slots. The callee fills in the **BTM** and **TOP** slots and initializes the local variables.

A choice point frame contains, in addition to the slots in a determinate frame, four slots located between the **TOP** slot and local variables:

CPF	Backtracking program pointer
H	Top of the heap
T	Top of the trail
B	Parent choice point

The **CPF** slot stores the program pointer to continue with when the current branch fails. The slot **H** points to the top of the heap when the frame is allocated. As in the WAM, a new register, called **HB**, is used as an alias for $B \rightarrow H$. When a variable is bound, it must be trailed if it is older than **B** or **HB**.

5.2 The extension of ATOAM for tabling

A new data area, called *table area*, is introduced for memorizing tabled subgoals and their answers. The *subgoal table* is a hash table that stores all the tabled subgoals

⁴ It is a convention in the literature that the stack is assumed to grow downwards

encountered in execution. For each tabled subgoal and its variants, there is an entry in the table, which is a record containing the following information:

SubgoalCopy
PioneerAR
State
TopMostLoopingSubgoal
DependentSubgoals
AnswerTable

The field `SubgoalCopy` points to the copy of the subgoal in the table area. In the copy, all variables are numbered. Therefore all variants of the subgoal are identical.

The field `PioneerAR` points to the frame of the pioneer, which is needed for implementing cuts. When the choice point of a tabled subgoal is cut off before the subgoal reaches completion, the field `PioneerAR` will be set to `NULL`. When a variant of the subgoal is encountered again after, the subgoal will be treated as a pioneer.

The field `State` indicates whether the subgoal is a looping subgoal, whether the answer table has been revised, and whether the subgoal is *complete* or *evaluated*. When execution backtracks to a top-most looping subgoal, if the *revised* bit is set, then another round will be started for the subgoal. A top-most looping subgoal becomes complete if this *revised* bit is unset after a round. At that time, the subgoal and all of its dependent subgoals will be set to *complete*. As described in 3.1.3, an *evaluated* subgoal is never evaluated again using rules in each round.

The `TopMostLoopingSubgoal` field points to the entry for the top-most looping subgoal, and the field `DependentSubgoals` stores the list of subgoals on which this subgoal depends. When a top-most looping subgoal becomes complete, all of its dependent subgoals turn to complete too.

The field `AnswerTable` points to the answer table for this subgoal, which is also a hash table. Hash tables expand dynamically. Let `g` be the pointer to the record for a subgoal in the table. The first answer in the answer table is referenced as `g->AnswerTable->FirstAnswer` and the last answer is referenced as `g->AnswerTable->LastAnswer`. In the beginning, the answer table is empty and both `FirstAnswer` and `LastAnswer` reference a dummy answer.

The frame for a tabled predicate contains the following two slots in addition to those slots stored in a choice point frame:

SubgoalTable
CurrentAnswer

The `SubgoalTable` points to the subgoal table entry, and the `CurrentAnswer` points to the last answer that has been consumed. The next answer can be reached from this reference on backtracking. When a frame is created, the slot `CurrentAnswer` is initialized to be `g->AnswerTable->FirstAnswer` where `g` is the pointer to the record for the tabled subgoal.

Three new instructions, namely `table_start`, `memo`, and `check_completion`, are introduced into the ATOAM for encoding the three table primitives. Figure 2 shows the compiled code of an example program.

```

% :-tabled p/2.
% p(X,Y):-p(X,Z),e(Z,Y).
% p(X,Y):-e(X,Y).
p/2: table_start 2,1
      fork r2
      para_value y(2)
      para_var y(-13)
      call p/2      % p(X,Z)
      para_value y(-13)
      para_value y(1)
      call e/2      % e(Z,Y)
      memo
r2:   fork r3
      para_value y(2)
      para_value y(1)
      call e/2      % e(X,Y)
      memo
r3:   check_completion p/2

```

Fig. 2. Compiled code of an example program.

The `table_start` instruction takes two operands: the arity (2) and the number of local variables (1). The `fork` instruction sets the CPF slot to hold the address to backtrack to on failure. The parameter passing instructions (`para_value` and `para_var` in this example) pass arguments to the callee's frame. The `memo` instruction is executed after an answer has been found. The `check_completion` instruction takes the entrance (`p/2`) as an operand so that the predicate can be re-entered when it needs re-evaluation.

5.3 Implementing semi-naive optimization

To implement semi-naive optimization, we add the following two pointers into the record for each tabled subgoal:

LastOldAnswer
LastPrevAnswer

where the pointer `LastOldAnswer` points to the last answer in the old region and the pointer `LastPrevAnswer` points to the last answer in the previous region. The `check_completion` instruction resets the pointers for all the tabled subgoals in the current cluster before it starts the next round:

```

for each subgoal g in the current cluster {
  g->LastOldAnswer = g->LastPrevAnswer;
  g->LastPrevAnswer = g->AnswerTable->LastAnswer;
}

```

The `memo` instruction is changed so that early promotion of answers is performed if the condition for promotion is met. Let `g` be the pointer to the tabled subgoal. If the subgoal has exhausted all its answers in the table and early promotion has

never be done before on the subgoal in the same round, then answers in the current region are promoted to the previous region:

```
g->LastPrevAnswer = g->AnswerTable->LastAnswer
```

The promoted answers will be moved to the old region before the start of the next round.

A bit vector is added into the frame for each tabled predicate to indicate if any new answer has been consumed by any tabled subgoal. Semi-naive optimization can be applied only if no tabled subgoal in the predicate has consumed any new answer.

A new instruction, called `last_depending_tabled_call`, is introduced to encode last depending tabled subgoals. In the example shown in Figure 2, the “`call p/2`” instruction is changed to “`last_depending_tabled_call p/2`” to enable semi-naive optimization. The `last_depending_tabled_call` instruction has the same behavior as the `call` instruction, but the callee can check the type of the instruction to see if it is invoked by a last depending tabled subgoal.

Let `g` be the pointer to the current tabled subgoal. The `table_start` instruction sets the `CurrentAnswer` slot of the frame to `g->LastOldAnswer` so that the subgoal consumes only new answers if: (1) the parent frame is a tabled frame; (2) no bit in the bit vector in the parent frame is set, which means that no tabled subgoal has consumed any new answer; and (3) the predicate is invoked by a `last_depending_tabled_call` instruction. If any of these condition is not satisfied, the `CurrentAnswer` slot is set to `g->AnswerTable->FirstAnswer` and all the answers will be consumed by the subgoal.

6 Performance Evaluation

We empirically compared the two answer consumption strategies and evaluated the effectiveness of semi-naive optimization. We also compared the performance of B-Prolog (version 6.9) with XSB (version 3.0). A Linux machine with 750MHz Intel process and 512GB RAM was used in the experiment. Benchmarks from three different sources were used:⁵ Datalog programs shown in Figure 3 with randomly generated graphs; the CHAT benchmark suite (Demoen and Sagonas 1999); and a parser, called *atr*, for the Japanese language defined by a grammar of over 860 rules (Uratani et al. 1994). This section presents the experimental results and reports the statistics to support the results. This section also gives experimental results on the Warren’s example for which SLG as implemented in XSB has lower time complexity than linear tabling when semi-naive optimization ceases to be effective.

6.1 Comparison of the two answer-consumption strategies

Table 1 compares the two answer-consumption strategies in terms of speed and stack space⁶ efficiencies. The difference of these two strategies in terms of CPU

⁵ The benchmarks are available from probp.com/bench.tar.gz.

⁶ The total usage of the local, global and trail stacks.

```

tcl:  tcl(X,Y):-edge(X,Y).
      tcl(X,Y):-tcl(X,Z),edge(Z,Y).

tcr:  tcr(X,Y):-edge(X,Y).
      tcr(X,Y):-edge(X,Z),tcr(Z,Y).

tcn:  tcn(X,Y):-edge(X,Y).
      tcn(X,Y):-tcn(X,Z),tcn(Z,Y).

sg:   sg(X,X).
      sg(X,Y):-edge(X,XX),sg(XX,YY),edge(Y,YY).

```

Fig. 3. Datalog programs.

Table 1. Comparison of the lazy and eager strategies.

program	CPU time		Stack space	
	Lazy	Eager	Lazy	Eager
tcl	1	1.02	1	1.00
tcr	1	0.96	1	1.00
tcn	1	0.90	1	1.00
sg	1	0.89	1	1.02
cs_o	1	1.17	1	1.36
cs_r	1	1.09	1	1.36
disj	1	1.06	1	1.41
gabriel	1	1.08	1	1.18
kalah	1	1.17	1	2.03
pg	1	2.28	1	3.59
peep	1	0.99	1	2.88
read	1	0.85	1	2.22
atr	1	1.03	1	1.06
<i>average</i>	1	1.12	1	1.62

time is small on average. This result implies that for programs with cuts declaring the use of the eager strategy would not cause significant slow-down. The difference in the usage of stack space is more significant than in CPU time. This is because, as discussed before, the lazy strategy has better locality than the eager strategy.

6.2 Effectiveness of semi-naive optimization

Table 2 shows the effectiveness of semi-naive optimization in gaining speed-ups under both strategies. Without this optimization, the system would consume over 30% more CPU time on average under either strategy. Our experiment also indicates that on average over 95% of the gains in speed are attributed to the *early promotion* technique.

Table 2. *Effectiveness of semi-naive optimization.*

program	CPU time ($\frac{nosemi}{semi}$)	
	Lazy	Eager
tcl	2.00	1.89
tcr	1.22	1.19
tcn	1.68	1.74
sg	1.22	1.51
cs_o	1.10	1.10
cs_r	1.09	1.10
disj	1.52	1.46
gabriel	1.32	1.15
kalah	1.52	1.41
pg	1.21	1.05
peep	1.09	1.11
read	1.98	1.27
atr	1.00	1.00
<i>average</i>	1.38	1.31

6.3 Comparison with XSB

Table 3 compares BP with XSB on time and stack space efficiencies. For XSB, the stack space is the total of the maximum amounts of global, local, trail, choice point, and SLG completion stack spaces. The default setting, namely, the SLG-WAM and the local scheduling strategy, is used. BP is faster than XSB on the Datalog programs and the parser but slower than XSB on the CHAT benchmark suite; and BP consumes considerably less stack space than XSB on some of the programs (*tcr*, *tcn*, *sg*, and *atr*).

The results must be interpreted with two differences of the two compared systems taken into account: On the one hand, BP is on average more than twice as fast as XSB for standard Prolog programs, and on the other hand the trie data structure used in XSB (Ramakrishnan et al. 1998) is far more advanced than hash tables used in BP for managing the table area. It is unclear to what extent each difference contributes to the overall efficiency.

The YAP implementation of SLG-WAM is up to twice as fast as XSB (Somogyi and Sagonas 2006) on the transitive closure and same-generation benchmarks with both chain and cyclic graphs. This entails that the BP implementation of linear tabling is comparable in speed with the most sophisticated implementation of SLG-WAM for the Datalog benchmarks.

The empirical data on the usage of table space are not reported. BP constantly consumes less table space than XSB for the benchmarks. In BP, both subgoal and answer tables are maintained as dynamic hashtables. In XSB, in contrast, tables are maintained as tries (Ramakrishnan et al. 1998). The usage of table space is independent of the strategies and optimizations. Both BP and XSB would consume the same amount of table space if the same data structure were employed.

Table 3. Comparison of B-Prolog and XSB.

program	BP (Lazy)	XSB	
		CPU time	Stack space
tcl	1	1.85	0.81
tcr	1	1.46	33.41
tcn	1	1.31	32.84
sg	1	1.47	109.12
cs_o	1	0.37	0.57
cs_r	1	0.35	0.73
disj	1	0.68	0.82
gabriel	1	0.61	2.05
kalah	1	1.00	0.58
pg	1	0.76	1.85
peep	1	0.37	2.97
read	1	0.69	11.12
atr	1	2.26	21.24

6.4 Statistics on iterations

Table 4 reports the statistics on the maximum (max its.) and average (ave. its.) numbers of iterations for tabled subgoals to reach their fixpoints.⁷ The column #subgoals shows the number of tabled subgoals. While for some programs, the maximum number of iterations performed is high (e.g., the maximum number for *atr* is 6), the average numbers are quite low.

The necessity of re-evaluating looping subgoals has been blamed for the low speed of iteration-based tabling systems (Zhou et al. 2000; Guo and Gupta 2001). Our new findings indicate that re-evaluation is not a dominant factor for the benchmarks. This statistics well explain why an implementation of linear tabling could achieve comparable speed performance with SLG-WAM for the benchmarks.

6.5 The complexity issue

The following is a slightly changed version of the Warren's example which disables semi-naive optimization:

```
:-table p/2.
p(X,Y) :- q(X,Z),c(Z,a,Y).
p(X,Y) :- q(X,Z),c(Z,b,Y).
p(X,X).

q(X,Y) :- p(X,Y).
```

Since the last depending subgoals $q(X,Z)$ in $p/2$ are not tabled, semi-naive optimization cannot be applied to $p/2$. For a string $(ab)^{n/2}$, the query $p(0,n)$ needs

⁷ Each subgoal has a counter which is initialized when the subgoal is tabled and is incremented each time the subgoal is resolved using rules. Note that semi-naive optimization may reduce the work of each iteration but has no effect on the number of iterations needed to reach the fixpoint.

Table 4. *Statistics on iterations.*

program	#subgoals	max its.	ave. its.
tcl	1	2	2.00
tcr	51	2	1.96
tcn	51	2	1.98
sg	153	2	1.32
cs_o	76	2	1.14
cs_r	76	2	1.16
disj	74	2	1.20
gabriel	59	2	1.20
kalah	102	3	1.24
pg	48	2	1.13
peep	49	3	1.29
read	131	5	1.34
atr	7139	6	1.81

$n/2$ iterations to reach the fixpoint. Since in each iteration the subgoal $q(X,Z)$ is rewritten into $p(X,Z)$ which returns all existing answers, the total time taken is $O(n^2)$. In contrast, the program takes only $O(n)$ time under SLG. For the size $n=5000$, it took BP 3.5 seconds to run the program while XSB only 15 milliseconds. For the original version of the program to which semi-naive optimization is applicable, it took BP only 7 milliseconds.

7 Related Work

There are three different tabling schemes, namely OLDT and SLG (Tamaki and Sato 1986; Sagonas and Swift 1998), CAT (Demoen and Sagonas 1998; Somogyi and Sagonas 2006), and iteration-based tabling including linear tabling (Shen et al. 1999; Shen et al. 2001; Zhou et al. 2000; Zhou and Sato 2003; Zhou et al. 2004) and DRA (Guo and Gupta 2001). SLG (Chen and Warren 1996) is a formalization based on OLDT for computing well-founded semantics for general programs with negation. The basic idea of using iterative deepening to compute fixpoints dates back to the ET* algorithm (Dietrich 1987).

In SLG-WAM, a consumer fails after it exhausts all the existing answers and its state is preserved by freezing the stack so that it can be reactivated after new answers are generated. The CAT approach does not freeze the stack but instead copies the stack segments between the consumer and its producer into a separate area so that backtracking can be done normally. The saved state is reinstalled after a new answer is generated. CHAT (Demoen and Sagonas 1999) is a hybrid approach that combines SLG-WAM and CAT.

Linear tabling relies on iterative computation of looping subgoals to compute fixpoints. Linear tabling is probably the easiest scheme to implement since no effort is needed to preserve states of consumers and the garbage collector can be kept untouched for tabling. Linear tabling is also the most space-efficient scheme since no extra space is needed to save states of consumers. Nevertheless, linear tabling

without optimization could be computationally more expensive than the other two schemes.

The DRA method (Guo and Gupta 2001) is also iteration based, but it identifies looping clauses dynamically and iterates the execution of looping clauses to compute fixpoints. While in linear tabling iteration is performed on only top-most looping subgoals, in DRA iteration is performed on every looping subgoal. In ET* (Dietrich 1987), every tabled subgoal is iterated even if it does not occur in a loop. Besides the difference in answer consumption strategies and optimizations, the linear tabling scheme described in this paper differs from the original version (Zhou et al. 2000; Shen et al. 2001) in that followers fail after they exhaust their answers rather than steal their pioneers' choice points. This strategy is originally adopted in the DRA method.

The two consumption strategies have been compared in XSB (Freire et al. 1998) as two scheduling strategies. The lazy strategy is called *local scheduling* and the eager strategy is called *single-stack scheduling*. Another strategy, called *batched scheduling*, is similar to local scheduling but top-most looping subgoals do not have to wait until their clusters become complete to return answers. Their experimental results indicate that local scheduling constantly outperforms the other two strategies on stack space and can perform asymptotically better than the other two strategies on speed. The superior performance of local scheduling is attributed to the saving of freezing stack segments. Although our experiment confirms the good space performance of the lazy strategy, it gives a counterintuitive result that the eager strategy is as fast as the lazy strategy. This result implies that the cost of iterative evaluation is considerably smaller than that of freezing stack segments, and for predicates with cuts the eager strategy can be used without significant slow-down. In our tabling system, different answer consumption strategies can be used for different predicates. The tabling system described in (Rocha et al. 2005a) also supports mixed strategies.

Semi-naive optimization is a fundamental idea for reducing redundancy in bottom-up evaluation of logic database queries (Bancilhon and Ramakrishnan 1986; Ullman 1988). As far as we know, its impact on top-down evaluation had been unknown before (Zhou et al. 2004). OLDT (Tamaki and Sato 1986) and SLG (Sagonas and Swift 1998) do not need this technique since it is not iterative and the underlying delaying mechanism successfully avoids the repetition of any derivation step. An attempt has been made by Guo and Gupta (Guo and Gupta 2001) to make incremental consumption of tabled answers possible in DRA. In their scheme, answers are also divided into three regions but answers are consumed incrementally as in bottom-up evaluation. Since no condition is given for the completeness and no experimental result is reported on the impact of the technique, we are unable to give a detailed comparison.

Our semi-naive optimization differs from the bottom-up version in two major aspects: Firstly, no differentiated rules are used. In the bottom-up version differentiated rules are used to ensure that at least one new answer is involved in the join of answers for each rule. Consider, for example, the clause:

$$H : -P, Q.$$

The following two differentiated rules are used in the evaluation instead of the original one:

$$H : -\Delta P, Q.$$

$$H : -P, \Delta Q.$$

Where ΔP denotes the new answers produced in the previous round for P. Using differentiated rules in top-down evaluation can cause considerable redundancy, especially when the body of a clause contains non-tabled subgoals.

The second major difference between our semi-naive optimization and the bottom-up version is that answers in our method are consumed sequentially until exhaustion, not incrementally as in bottom-up evaluation. A tabled subgoal consumes either all available answers or only new answers including answers produced in the current round. Neither incremental consumption nor sequential consumption seems satisfactory. Incremental consumption avoids redundant joins but may require more rounds to reach fixpoints. In contrast, sequential consumption never need more rounds to reach fixpoints but may cause redundant joins of answers. The early promotion technique alleviates the problem of sequential consumption. By promoting answers early from the *current* region to the *previous* region, we can considerably reduce the redundancy in joins.

Semi-naive optimization may lower time complexities in bottom-up evaluation (Bancilhon and Ramakrishnan 1986). The same result holds to the top-down version as demonstrated by Warren’s example. Our experimental results show that semi-naive optimization gives an average speed-up of over 30% to linear tabling if answers are promoted early, and almost no speed gain if no answer is promoted early. In linear tabling, only looping subgoals need to be iteratively evaluated. For non-looping subgoals, no re-evaluation is necessary and thus semi-naive optimization has no effect at all on the performance. Most of the looping subgoals in our chosen benchmarks reach their fixpoints after 2-3 iterations. In general, more iterations are needed to reach fixpoints in bottom-up evaluation. In addition, in bottom-up evaluation, the order of the joins can be optimized and no further joins are necessary once a participating set is known to be empty. In contrast, in linear tabling joins are done in strictly chronological order. For a conjunction (P, Q, R) , the join $P \bowtie Q$ is computed even if no answer is available for R . Because of all these factors, semi-naive optimization is not as effective in linear tabling as in bottom-up evaluation.

Our semi-naive optimization requires the identification of last depending subgoals. For this purpose, a level mapping is used to represent the call graph of a given program. The use of a level mapping to identify optimizable subgoals is analogous to the idea used in the stratification-based methods for evaluating logic programs (Apt et al. 1988; Chen and Warren 1996; Przymusiński 1989). In our level mapping, only predicate symbols are considered. It is expected that more accurate approximations can be achieved if arguments are considered as well.

Semi-naive optimization does not solve all the problems of recomputation in linear tabling. Recall the Warren’s example:

```
:-table p/2.  
p(X,Y) :- p(X,Z),c(Z,a,Y).  
p(X,Y) :- p(X,Z),c(Z,b,Y).  
p(X,X).
```

Assume there is a very costly non-tabled subgoal preceding $p(X,Z)$, then the subgoal has to be executed in each iteration even with semi-naive optimization. This example demonstrates the acuteness of the problem of recomputation because the number of iterations needed to reach the fixpoint is not constant. One treatment would be to table the subgoal to avoid recomputation, as suggested in (Guo and Gupta 2001), but tabling extra predicates can cause other problems such as over consumption of table space.

8 Conclusion

In this paper we have described two answer consumption strategies (namely, *lazy* and *eager* strategies) and semi-naive optimization for linear tabling. We have compared the two strategies both qualitatively and quantitatively. Our results indicate that, while the lazy strategy has better space efficiency than the eager strategy, the eager strategy is comparable in speed with the lazy strategy. This result implies that for all-solution search programs the lazy strategy should be adopted and for partial-solution search programs including programs with cuts the eager strategy should be used.

We have tailored semi-naive optimization to linear tabling and have given sufficient conditions for it to be complete. Moreover, we have proposed a technique called *early answer promotion* to reduce redundant consumption of answers. Our experimental result indicates that semi-naive optimization gives significant speed-ups to some programs.

Linear tabling has several attractive advantages including its simplicity, ease of implementation, and good space efficiency. Early implementations of linear tabling were several times slower than XSB. This paper has demonstrated for the first time that linear tabling with optimization is as competitive as SLG on time efficiency as well for the benchmarks.

Semi-naive optimization does not solve all the problems of recomputation in linear tabling. There are programs for which recomputation can be costly, even leading to higher complexities. The future work is to identify the patterns of such programs and find methods to deal with them.

9 Acknowledgement

The preliminary results of this article appear in ACM PPDP'03 and PPDP'04. Taisuke Sato is supported in part by CREST, and Yi-Dong Shen is supported in part by the National Natural Science Foundation of China grants numbered 60673103 and 60421001.

References

- APT, K., BLAIR, H. A., AND WALKER, A. 1988. Towards a theory of declarative knowledge. In *Foundations of deductive databases and logic programming*, J. Minker, Ed. Morgan Kaufmann, 89–142.
- BANCILHON, F. AND RAMAKRISHNAN, R. 1986. An amateur’s introduction to recursive query processing strategies. *Proc. of ACM SIGMOD ’86*, 16–52.
- CHEN, W. AND WARREN, D. S. 1996. Tabled evaluation with delaying for general logic programs. *Journal of the ACM* 43, 1, 20–74.
- DAWSON, S., RAMAKRISHNAN, C. R., AND WARREN, D. S. 1996. Practical program analysis using general purpose logic programming systems — A case study. *ACM SIGPLAN Notices* 31, 5, 117–126.
- DEMOEN, B. AND SAGONAS, K. 1998. CAT: The copying approach to tabling. In *Proceedings of Programming Language Implementation and Logic Programming (PLILP)*. LNCS 1490, 21–35.
- DEMOEN, B. AND SAGONAS, K. 1999. CHAT: The copy-hybrid approach to tabling. In *Proceedings of Practical Aspects of Declarative Programming (PADL)*. LNCS 1551, 106–121.
- DIETRICH, S. W. 1987. Extension tables: Memo relations in logic programming. In *IEEE Fourth Symposium on Logic Programming*. 264–272.
- EISNER, J., GOLDLUST, E., AND SMITH, N. A. 2004. Dyna: A declarative language for implementing dynamic programs. In *Proc. of the 42nd Annual Meeting of ACL*.
- FREIRE, J., SWIFT, T., AND WARREN, D. S. 1998. Beyond depth-first: Improving tabled logic programs through alternative scheduling strategies. *Journal of Functional and Logic Programming*.
- GUO, H.-F. AND GUPTA, G. 2001. A simple scheme for implementing tabled logic programming systems based on dynamic reordering of alternatives. In *Proceedings International Conference on Logic Programming (ICLP)*. LNCS 2237, 181–195.
- JOHNSON, M. 1995. Memoization of top down parsing. *Computational Linguistics* 21, 3.
- LIU, M. 1999. Deductive database languages: Problems and solutions. *ACM Computing Surveys* 31, 1, 27–62.
- LLOYD, J. W. 1988. *Foundation of Logic Programming*, 2 ed. Springer-Verlag.
- MICHIE, D. 1968. “memo” functions and machine learning. *Nature*, 19–22.
- NIELSON, F., NIELSON, H. R., SUN, H., BUCHHOLTZ, M., HANSEN, R. R., PILEGAARD, H., AND SEIDL, H. 2004. The succinct solver suite. In *Proc. Tools and Algorithms for the Construction and Analysis of Systems: 10th International Conference (TACAS)*, LNCS 2988. 251–265.
- PIENTKA, B. December 2003. Tabled higher-order logic programming. Ph.D. thesis, Technical Report CMU-CS-03-185.
- PRZYMUSINSKI, T. C. 1989. Every logic program has a natural stratification and an iterated least fixed point model. In *PODS ’89. Proceedings of the Eighth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems*. ACM Press, 11–21.
- RAMAKRISHNAN, C. 2002. Model checking with tabled logic programming. In *ALP News Letter*. ALP.
- RAMAKRISHNAN, I., RAO, P., SAGONAS, K., SWIFT, T., AND WARREN, D. 1998. Efficient access mechanisms for tabled logic programs. *J. Logic Programming* 38, 31–54.
- RAMAKRISHNAN, R. AND ULLMAN, J. D. 1995. A survey of deductive database systems. *Journal of Logic Programming* 23, 2, 125–149.
- ROCHA, R., SILVA, F., AND COSTA, V. S. 2005a. Dynamic mixed-strategy evaluation of tabled logic programs. In *Proceedings International Conference on Logic Programming (ICLP)*. 250–264.

- ROCHA, R., SILVA, F., AND COSTA, V. S. 2005b. On applying or-parallelism and tabling to logic programs. *Theory and Practice of Logic Programming (TPLP)* 5, 1 & 2, 161–205.
- SAGONAS, K. AND SWIFT, T. 1998. An abstract machine for tabled execution of fixed-order stratified logic programs. *ACM Transactions on Programming Languages and Systems* 20, 3, 586–634.
- SAGONAS, K., SWIFT, T., AND WARREN, D. S. 1994. XSB as a deductive database. *SIGMOD Record (ACM Special Interest Group on Management of Data)* 23, 2, 512–512.
- SATO, T. AND KAMEYA, Y. 2001. Parameter learning of logic programs for symbolic-statistical modeling. *Journal of Artificial Intelligence Research*, 391–454.
- SHEN, Y.-D., YUAN, L., YOU, J., AND ZHOU, N.-F. 2001. Linear tabulated resolution based on Prolog control strategy. *Theory and Practice of Logic Programming (TPLP)* 1, 1, 71–103.
- SHEN, Y.-D., YUAN, L.-Y., YOU, J.-H., AND ZHOU, N.-F. 1999. Linear tabulated resolutions for the well-founded semantics. In *Proceedings of Logic Programming and Non-monotonic Reasoning*. 192–205.
- SOMOGYI, Z. AND SAGONAS, K. 2006. Tabling in mercury: Design and implementation. In *Proceedings of Practical Aspects of Declarative Programming (PADL)*. Springer-Verlag, to appear.
- TAMAKI, H. AND SATO, T. 1986. OLD resolution with tabulation. In *Proceedings of the Third International Conference on Logic Programming*, E. Shapiro, Ed. 84–98.
- ULLMAN, J. D. 1988. *Database and Knowledge-Base Systems*. Vol. 1 & 2. Computer Science Press.
- URATANI, N., TAKEZAWA, T., MATSUO, H., AND MORITA, C. 1994. ATR integrated speech and language database. Technical Report TR-IT-0056, ATR Interpreting Telecommunications Research Laboratories. In Japanese.
- WARREN, D. S. 1992. Memoing for logic programs. *Comm. of the ACM, Special Section on Logic Programming* 35, 93–111.
- WARREN, D. S. 1999. *Programming in Tabled Prolog*. DRAFT 1 (<http://www.cs.sunysb.edu/~warren/xsbbook/book.html>).
- ZHOU, N.-F. 1996. Parameter passing and control stack management in Prolog implementation revisited. *ACM Transactions on Programming Languages and Systems* 18, 6, 752–779.
- ZHOU, N.-F. AND SATO, T. 2003. Efficient fixpoint computation in linear tabling. In *Fifth ACM-SIGPLAN International Conference on Principles and Practice of Declarative Programming*. 275–283.
- ZHOU, N.-F., SATO, T., AND HASIDA, K. 2003. Toward a high-performance system for symbolic and statistical modeling. In *IJCAI Workshop on Learning Statistical Models from Relational Data*. 153–159.
- ZHOU, N.-F., SHEN, Y.-D., AND SATO, T. 2004. Semi-naive evaluation in linear tabling. In *Fifth ACM-SIGPLAN International Conference on Principles and Practice of Declarative Programming*. 90–97.
- ZHOU, N.-F., SHEN, Y.-D., YUAN, L.-Y., AND YOU, J.-H. 2000. Implementation of a linear tabling mechanism. In *Proceedings of Practical Aspects of Declarative Programming (PADL)*. LNCS 1753, 109–123.