

# On Extending a Full-Sharing Design with Batched Scheduling

Miguel Areias and Ricardo Rocha

CRACS & INESC TEC and Faculty of Sciences, University of Porto  
Rua do Campo Alegre, 1021, 4169-007 Porto, Portugal  
{miguel-areias,ricroc}@dcc.fc.up.pt

**Abstract. Keywords:** Multithreading, Tabling, Concurrency, Batched Scheduling

## 1 Introduction

Tabling [4] is an implementation technique that overcomes some limitations of traditional Prolog systems in dealing with redundant sub-computations and recursion. Tabling consists of storing intermediate answers for subgoals so that they can be reused when a repeated subgoal appears during the resolution process. Tabling has become a popular and successful technique thanks to the groundbreaking work in the XSB Prolog system and in particular in the SLG-WAM engine [7], the most successful engine of XSB. The success of SLG-WAM led to several alternative implementations that differ in the execution rule, in the data-structures used to implement tabling, and in the changes to the underlying Prolog engine. Implementations of tabling are now widely available in systems like Yap Prolog, B-Prolog, ALS-Prolog, Mercury, Ciao Prolog and more recently Picat.

Multithreading in Prolog is the ability to concurrently perform computations, in which each computation runs independently but shares the program clauses. When multithreading is combined with tabling, we have the best of both worlds, since we can exploit the combination of higher procedural control with higher declarative semantics. To the best of our knowledge, XSB [6] and Yap [2] were the only Prolog systems that were able to combine both multithreading with tabling. Yap implements a SWI-Prolog compatible multi-threading library [8]. Yap's threads have their own execution stacks and only share the code area where predicates, records, flags and other global non-backtrackable data are stored. For tabled evaluation, a thread views its tables as private but, at the engine level, Yap has three designs that vary from a *No-Sharing* (NS) design, where each thread allocates fully private tables for each new tabled subgoal called during its computation, to a *Full-Sharing* (FS), where threads share the complete table space.

In this paper we discuss a novel approach for combining the FS design with batched scheduling and we present a performance analysis comparison between local scheduling with batched scheduling. Experimental results show that despite

the extra data structures required to combine FS design with batched scheduling, the execution time of combination is still quite competitive when comparing against the FS with local scheduling.

The remainder of the paper is organized as follows. First, we briefly introduce some background and related work. Then, we describe our approach for combining batched scheduling with the FS design. Next, we show the most important implementation details and finally, we discuss experimental results and we end by outlining some conclusions.

## 2 Background

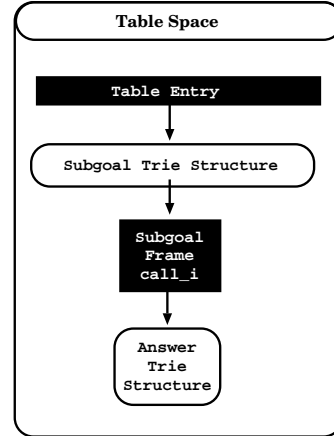
This section introduces some background needed for the following sections.

### 2.1 Table Space Organization

The basic idea behind tabling is straightforward: programs are evaluated by storing answers for tabled subgoals in an appropriate data space, called the *table space*. Similar calls to tabled subgoals are not re-evaluated against the program clauses, instead they are resolved by consuming the answers already stored in their table entries. During this process, as further new answers are found, they are stored in their tables and later returned to all similar calls.

Figure 1 shows the table space organization of the Yap system. At the entry point we have the *table entry* data structure. This structure is allocated when a tabled predicate is being compiled, so that a pointer to the table entry can be included in its compiled code. This guarantees that further calls to the predicate will access the table space starting from the same point. Below the table entry, we have the *subgoal trie structure*. Each different tabled subgoal call to the predicate at hand corresponds to a unique path through the subgoal trie structure, always starting from the table entry, passing by several subgoal trie data units, the *subgoal trie nodes*, and reaching a leaf data structure, the *subgoal frame*.

The subgoal frame stores additional information about the subgoal and acts like an entry point to the *answer trie structure*. Each unique path through the answer trie data units, the *answer trie nodes*, corresponds to a different answer to the entry subgoal.



**Fig. 1.** Table space organization

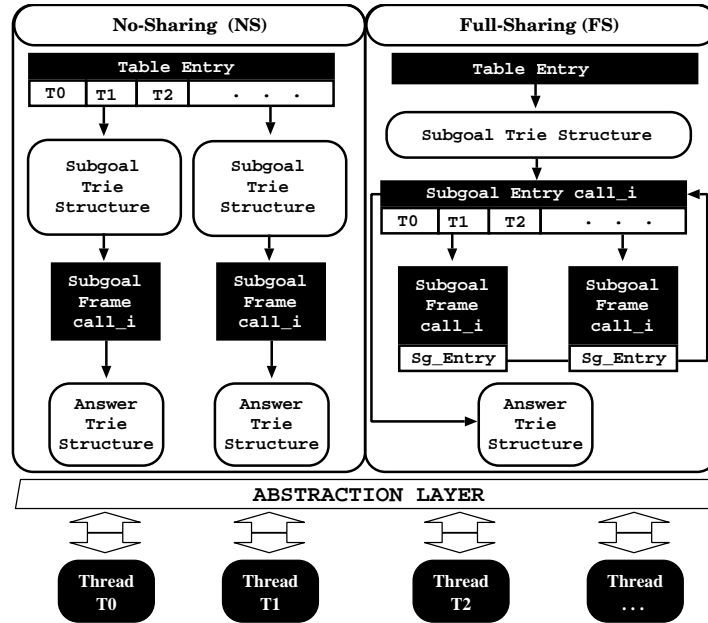
### 2.2 Yap's Multithreaded Tabling Support

Yap implements a SWI-Prolog compatible multi-threading library [8]. Like in SWI-Prolog, Yap's threads have their own execution stacks and only share the

code area where predicates, records, flags and other global non-backtrackable data are stored.

In a multithreaded tabling evaluation, a thread views its tables as private but, at the engine level, Yap has three designs that vary from a *No-Sharing* (NS) design, where each thread allocates fully private tables for each new tabled subgoal called during its computation, to a *Full-Sharing* (FS), where threads share the complete table space. Figure 1 shows Yap's multithreaded tabling support with an abstraction layer to which threads interact (the threads represented in the bottom of the figure), and at the engine level the configuration of the NS and FS designs in a situation where several different threads evaluated the same tabled subgoal call  $i$ .

In the NS design, each thread allocates fully private tables for each new subgoal called during its computation. One can observe that the table entry data structure still stores the common information for the predicate (such as the arity or the evaluation strategy), and then each thread  $t$  has its own cell  $T_t$  inside a *bucket array* which points to the private data structures.



**Fig. 2.** Yap's multithreaded tabled space using the NS and FS designs

In the FS design, part of the subgoal frame information (the *subgoal entry* data structure) and the answer trie structure are now shared among all threads. The previous subgoal frame data structure was split into two: the *subgoal entry* stores common information for the subgoal call (such as the pointer to the

shared answer trie structure); the remaining information (the *subgoal frame* data structure) remains private to each thread.

### 2.3 Batched Scheduling

The decision about the evaluation flow is determined by the *scheduling strategy*. Different strategies may have a significant impact on performance, and may lead to a different ordering of solutions to the query goal. Arguably, the two most successful tabling scheduling strategies are *local scheduling* and *batched scheduling* [5].

Local scheduling strategy schedules the evaluation of a program in a breath-first manner. It favors the backtracking first with completion instead of the forward execution, leaving the consumption of answers for last. Thus, it only allows a Cluster of Dependent Subgoals (CDS) to return answers only after the completion point has been reached [5]. In other words, the local scheduling tries to keep a CDS as minimal as possible. When new answers are found, they are added to the table space and the computation fails as consequence, tabled subgoals inside a CDS propagate their answers to outside the CDS only after its completion point is found. Local scheduling causes a sooner completion of subgoals, which creates less complex dependencies between them.

On the other hand, batched scheduling schedules the evaluation of a program in a depth-first manner. It favors the forward execution first instead of backtracking, leaving the consumption of answers and completion for last. It thus tries to delay the need to move around the search tree by batching the return of answers. When new answers are found for a particular tabled subgoal, they are added to the table space and the execution continues. For some situations, this results in creating dependencies to older subgoals, therefore enlarging the current CDS [7] and delaying the completion point to an older generator node.

## 3 Full-Sharing with Batched Scheduling

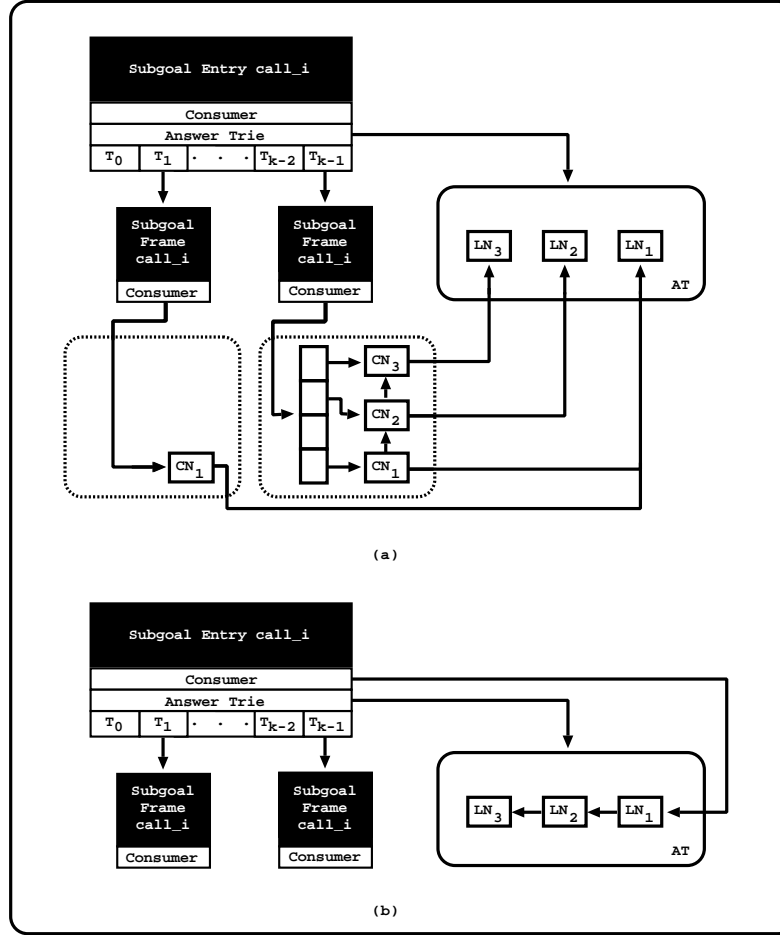
In this section, we show how we have extended the FS design to support batched scheduling.

### 3.1 Our Approach

The key idea of our approach, which we named Privately-consumed Answer Chaining (PAC), is then to extend the *FS* design with batched scheduling, by chaining privately for each subgoal call, the answers that were already consumed by a thread. Since the procedure is private, it will only affect the thread that is doing it. At the end, when the evaluation is complete, i.e, when a subgoal call is marked as complete, we put one of the private chain as public, so that from that point on all threads can use that chain in complete (only reading) mode.

Figure 3 shows the key data structures for supporting the implementation of the PAC procedure during the evaluation of a tabled subgoal call  $i$  using the

FS design. The FS design, uses a subgoal entry data structure to store common information for a subgoal call and a subgoal frame (*SF*) data structure to store private information about the execution of each thread. The PAC procedure works at the subgoal frame level, which is private to each thread.



**Fig. 3.** The *FS* design with the PAC procedure - (a) private chaining and (b) public chaining

Figure 3(a) shows then a situation where two threads,  $T_1$  and  $T_{k-1}$ , are sharing the same subgoal entry call  $i$  when the subgoal is still under evaluation, i.e., the subgoal is not yet complete. The current state of the evaluation shows an answer trie with 3 answers found for the subgoal call  $i$ . For the sake of simplicity, we are omitting the internal answer trie nodes and we are only showing the leaf nodes  $LN_1$ ,  $LN_2$  and  $LN_3$ , in the figure. With the PAC procedure, the

leaf nodes are not chained in the answer trie data structure. Now, the chaining process is done privately, and for that, we use the subgoal frame structure of each thread. On the subgoal frame structure we added a new field, called *consumer*, to store the answers found within the execution of the thread. In order to minimize the impact of the *PAC* optimization, each node within the new consumer answer structure has two fields: (i) an entry pointer, which points to the corresponding leaf node in the answer trie data structure; (ii) a next pointer to chain the answers within the consumer structure. To maintain a good performance, when the number of nodes exceeds a certain threshold, we use a hash trie mechanism design similar to the one presented in the work [3]. However, since this mechanism is private to each thread, it does not require any of the tools that were necessary to support concurrency. In particular, on each hash trie level, we have removed the tools necessary to support concurrency, such as useless pointers and compare-and-swap operations. We have chosen this hashing mechanism, because it showed a good balance between lookup and insert operations [3], but the major reason was mostly because of the integration in the TabMalloc memory allocator [1].

Going back to Figure 3(a), the consumer answer structures represent then two different situations where threads can be evaluating a subgoal call. Thread  $T_1$  has only found one answer and it is using a direct consumer answer chaining to access the node  $LN_1$ . Thread  $T_{k-1}$  was already found three answers for the subgoal call and it is already using the hash trie mechanism within its consumer answer structure. The consumer nodes are chained between themselves, thus that consumer nodes belonging to thread  $T_{k-1}$  can consume the answers as in the original mechanism.

Figure 3(b) shows the state of the subgoal call after completion (recall that after completion of a subgoal call, the threads use loader nodes to consume the answers). When a thread  $T$  completes a subgoal call, it frees its private consumer structures, but before doing that, it checks whether another thread as already marked the subgoal as completed. If no other thread has done that, then thread  $T$  not only follows its private chaining mechanism as it would for freeing its private nodes, but also, follows the pointers to the answer trie leaf nodes in order to reproduce the chain inside the answer trie. Since this procedure is done inside a critical region, no more than one thread can be doing this chaining process. Thus, in Figure 3(b), we are showing a situation where the subgoal call is completed and both threads  $T_1$  and  $T_{k-1}$  have already removed their consumer answer structures and chained the leaf nodes inside the answer trie.

### 3.2 Implementations Details

At the implementation level, the major difference between local and batched scheduling is in the tabling operation *tabled new answer*, where we decide what to do when an answer is found during the evaluation. This operation checks whether a newly found answer is already in the corresponding answer trie structure and, if not, inserts it. The usage of batched scheduling with the FS design requires further support since with batched scheduling, answers are immediately propagated and we have to ensure that the propagation of an answer occurs on

all subgoal calls one and only once. To do so, we take advantage of the private chaining procedure, presented in the previous subsection, as a way to keep, for every subgoal call of every thread, track of all the answers that were already propagated. This requires minor changes to the *tabled new answer* tabling operation. Algorithm 1 shows how we have extended the tabled new answer operation to support the FS design with batched scheduling.

---

**Algorithm 1** `tabled_new_answer(answer ANS, subgoal frame SF)`

---

```

1: leaf  $\leftarrow$  check_insert_answer_trie(ANS, SF)
2: chain  $\leftarrow$  check_insert_consumer_chain(leaf, SF)
3: if is_answer_marked_as_found(chain) = True then
4:     return failure
5: else {the answer is new}
6:     mark_answer_as_found(chain)
7:     if local_scheduling_mode(SF) then
8:         return failure
9:     else {batched scheduling mode}
10:    return proceed

```

---

The algorithm receives two arguments: the new answer found during the evaluation (*ANS*) and the subgoal frame which corresponds to the call at hand (*SF*). The algorithm begins by checking/inserting the given *ANS* into the answer trie structure, which will return the leaf node for the path representing *ANS* (line 1).

Then, it checks/inserts the given *leaf* node into the private consumer chain for the current thread, which will return the corresponding chain node (line 2). Next in line 3, it tests whether the chain node already existed in the consumer chain, i.e., if it was inserted or not by the current check/insert operation in order to return failure (line 4), or it proceeds with marking the answer *ANS* has found (line 6). At the end (lines 7 to 10), it returns failure if local scheduling is active (line 8), otherwise, the batched scheduling is active, thus it propagates the answer *ANS* (line 10).

## 4 Performance Analysis

We now present experimental results about the usage of the batched scheduling on the *NS*, *SS* and *FS* designs. The environment for our experiments was a machine with 32-Core AMD Opteron (TM) Processor 6274 (2 sockets with 16 cores each) with 32G of main memory, each processor with caches *L1*, *L2* and *L3* respectively with the sizes of 64K, 2048K and 6144K, running the Linux kernel is the 3.16.7-200.fc20.x86\_64, with Yap 6.3 compiled with gcc 4.8.

Table 1 shows the overhead ratios, when compared with the NS design with 1 thread (running with local scheduling), for the NS and FS designs (all running TabMalloc), when running 1, 8, 16, 24 and 32 threads with local and batched

scheduling on the five sets of benchmarks. For each design, the table has then two columns, a column with **Local** that shows results for the local scheduling and a column with **Batched** with results for batched scheduling. The overhead results presented in both **Local** and **Batched** columns use as base time the execution times presented in the **NS** column of the Table ??.

**Table 1.** Overhead ratios, when compared with the NS design with 1 thread (running with local scheduling without TabMalloc) for the NS and FS designs (with TabMalloc), when running 1, 8, 16, 24 and 32 threads with local and batched scheduling on the five sets of benchmarks (best ratios by row and by design for the Minimum, Average and Maximum are in bold)

Threads		NS		FS	
		Local	Batched	Local	Batched
1	Min	<b>0.53</b>	0.55	1.01	<b>0.95</b>
	Avg	<b>0.78</b>	0.82	<b>1.30</b>	1.46
	Max	1.06	<b>1.05</b>	<b>1.76</b>	2.33
	StD	0.15	0.14	0.22	0.44
8	Min	0.66	<b>0.63</b>	1.16	<b>0.99</b>
	Avg	<b>0.85</b>	0.88	<b>1.88</b>	1.95
	Max	<b>1.12</b>	1.14	<b>2.82</b>	3.49
	StD	0.13	0.14	0.60	0.79
16	Min	0.85	<b>0.75</b>	1.17	<b>1.06</b>
	Avg	<b>0.98</b>	1.00	<b>1.97</b>	2.08
	Max	<b>1.16</b>	1.31	<b>3.14</b>	3.69
	StD	0.09	0.17	0.65	0.83
24	Min	<b>0.91</b>	0.93	1.16	<b>1.09</b>
	Avg	<b>1.15</b>	1.16	<b>2.06</b>	2.19
	Max	1.72	<b>1.60</b>	<b>3.49</b>	4.08
	StD	0.20	0.21	0.70	0.91
32	Min	1.05	<b>1.04</b>	1.33	<b>1.26</b>
	Avg	1.51	<b>1.49</b>	<b>2.24</b>	2.41
	Max	<b>2.52</b>	2.63	<b>3.71</b>	4.51
	StD	0.45	0.45	0.74	1.02

By observing Table 1, we can see that, for one thread, on average, local scheduling is slightly better than batched on the two designs. For the NS design we have 0.78 and 0.82 and for the FS design we have 1.30 and 1.46 average overhead ratios, for the local and batched scheduling strategies, respectively. As we scale the number of threads, one can observe that, for the NS design both scheduling strategies have similar minimum, average and maximum overhead ratios. For the FS design, the best minimum overhead ratio on 1, 8, 16, 24 and 32 threads is always the batched scheduling. For the average and maximum overhead ratio, local scheduling is always better than batched scheduling. However, for the average ratio, the results between both strategies is quite close. We consider this to be quite significant, once the FS design with batched scheduling



requires the support of the PAC procedure. For the maximum overhead ratio, the difference between both strategies is higher than for the average overhead ratio, however we hope to reduce the difference with further improvements in the PAC procedure.

## 5 Conclusions and Further Work

### References

1. Areias, M., Rocha, R.: An Efficient and Scalable Memory Allocator for Multi-threaded Tabled Evaluation of Logic Programs. In: International Conference on Parallel and Distributed Systems. pp. 636–643. IEEE Computer Society (2012)
2. Areias, M., Rocha, R.: Towards Multi-Threaded Local Tabling Using a Common Table Space. *Journal of Theory and Practice of Logic Programming*, International Conference on Logic Programming, Special Issue 12(4 & 5), 427–443 (2012)
3. Areias, M., Rocha, R.: A lock-free hash trie design for concurrent tabled logic programs. *International Journal of Parallel Programming* pp. 1–21 (2015)
4. Chen, W., Warren, D.S.: Tabled Evaluation with Delaying for General Logic Programs. *Journal of the ACM* 43(1), 20–74 (1996)
5. Freire, J., Swift, T., Warren, D.S.: Beyond Depth-First: Improving Tabled Logic Programs through Alternative Scheduling Strategies. In: International Symposium on Programming Language Implementation and Logic Programming. pp. 243–258. No. 1140 in LNCS, Springer (1996)
6. Marques, R., Swift, T.: Concurrent and Local Evaluation of Normal Programs. In: International Conference on Logic Programming. pp. 206–222. No. 5366 in LNCS, Springer (2008)
7. Sagonas, K., Swift, T.: An Abstract Machine for Tabled Execution of Fixed-Order Stratified Logic Programs. *ACM Transactions on Programming Languages and Systems* 20(3), 586–634 (1998)
8. Wielemaker, J.: Native Preemptive Threads in SWI-Prolog. In: International Conference on Logic Programming. pp. 331–345. No. 2916 in LNCS, Springer (2003)