

FACULTY OF INFORMATICS, MASARYK UNIVERSITY

Multi-Threaded Nested DFS

BACHELOR'S THESIS

Petr Ročkai

Brno, spring 2007

Declaration

Thereby I declare that this thesis is my original work, which I have created on my own. All sources and literature used in writing the thesis, as well as any quoted material, are properly cited, including full reference to its source.

Advisor: RNDr. Jiří Barnat, PhD.

Acknowledgements

I would like to thank my advisor RNDr. Jiří Barnat, PhD. for the help and motivation he provided throughout my work on this thesis and prof. RNDr. Luboš Brim, CSc. for valuable input. My thanks go also to the fine people in the Parallel and Distributed Systems Laboratory and especially those that have contributed to the DiViNE project. Finally, I would like to thank my family and Lucy – for the support and all the little things.

Abstract

Recent development in computer hardware has brought more wide-spread emergence of shared-memory, multi-core systems. These architectures offer opportunities to speed up various tasks – among others LTL model checking. We have implemented several parallel algorithms for shared-memory LTL model checking, based on Nested DFS and distributed-memory algorithms. To achieve good scalability, we have devised and experimentally evaluated several implementation techniques, which we present in the thesis.

Keywords

Model Checking, Parallel, Multi-Threading, Multi-Core, Shared Memory, Symmetric Multiprocessing

Contents

1	Introduction	6
2	Parallel LTL Model-Checking Algorithms	8
2.1	Maximal Accepting Predecessor	8
2.2	One Way Catch Them Young	9
2.3	Negative Cycle Detection	9
2.4	Back-Level Edges	10
2.5	Cycle Locality	10
3	Implementation Techniques	12
3.1	Shared Memory Platform	12
3.2	Implementing Algorithms in Shared Memory	13
3.3	Communication	14
3.4	Memory Allocation	16
3.5	Efficient Termination Detection	17
4	Implemented Algorithms	20
4.1	Nested DFS	20
4.2	Nested DFS with VCL	20
4.3	OWCTY	22
4.4	Nested DFS and OWCTY	23
5	Experiments	25
5.1	Methodology	25
5.2	Results	25
5.3	Nested DFS and VCL	27
5.4	Comparison with SPIN	27
6	Conclusions	31

Chapter 1

Introduction

With the arrival of 64-bit technology the traditional space limitations in explicit-state model checking¹, may become less severe. Instead, time could quickly become an important bottleneck. This naturally raises interest in using parallelism to fight the “time-explosion” problem.

Much of the extensive research on the parallelisation of model checking algorithms followed the distributed-memory programming model and the algorithms were parallelised for networks of workstations, largely due to easy access to networks of workstations. Recent shift in architecture design toward multicores has intensified research pertaining to shared-memory paradigm as well.

In [11] G. Holzmann proposed an extension of the SPIN [12] model-checker for dual-core machines. The algorithms keep their linear time complexity and the liveness checking algorithm supports full LTL. The algorithm for checking safety properties scales well to N-core systems. The algorithm for liveness checking, which is based on the original SPIN’s nested DFS algorithm, is however unable to scale to N-core systems. It is still an open problem to do verification of general liveness properties on N-cores with linear time complexity.

A different approach to shared-memory model checking is presented in [14], based on CTL* translation to Hesitant Alternating Automata. The proposed algorithm uses so-called non-emptiness game for deciding validity of the original formula and is therefore largely unrelated to the algorithms based on fair-cycle detection.

In this work, we show several general techniques useful in writing parallel shared-memory model checking algorithms, as well as an implementation of basic versions of several such algorithms, with special focus on those based on the venerable Nested DFS. This lends those algorithms linear time execution and on-the-fly characteristic. Furthermore, a shared-memory version of the

¹A comprehensive introduction to model checking and the approach based on fair-cycle detection, as used in this thesis, may be found in [8].

One Way Catch Them Young algorithm, originally intended for distributed-memory systems, is implemented and compared to Nested DFS. Finally, the implemented algorithms are experimentally evaluated on a few example models and properties and the results are compared to an MPI-based implementation and to the SPIN modelchecker.

Many of the techniques and results from this thesis are also presented in [2]. Source code for the implemented algorithms and accompanying infrastructure is available on the attached CD.

Chapter 2

Parallel LTL Model-Checking Algorithms

Efficient parallel solution of many problems often requires approaches radically different from those used to solve the same problems sequentially. Classical examples are list rankings, connected components, depth-first search in planar graphs etc. In the area of LTL model-checking the best known enumerative *sequential* algorithms based on fair-cycle detection are the *Nested DFS* algorithm [9, 13] (implemented, e.g., in the model checker SPIN [12]) and *SCC-based algorithms* originating in Tarjan's algorithm for the decomposition of the graph into strongly connected components (SCCs) [19]. However, both algorithms rely on inherently sequential depth-first search postorder, hence it is difficult to adapt them to parallel architectures. Consequently, different techniques and algorithms are needed. Unlike LTL model-checking, the reachability analysis is a verification problem with efficient parallel solution. The reason is that the exploration of the state space can be implemented e.g. using breadth-first search. In the following, we sketch four parallel algorithms for enumerative LTL model checking that are, more or less, based on performing repeated parallel reachability to detect reachable accepting cycles. The reader is kindly asked to consult the original sources for the details.

2.1 Maximal Accepting Predecessor

The main idea of the algorithm [4, 6] is based on the fact that every accepting vertex lying on an accepting cycle is its own predecessor. An algorithm that is directly derived from the idea, would require expensive computation as well as space to store all proper accepting predecessors of all (accepting) vertices. To remedy this obstacle, the MAP algorithm stores only a single representative of all proper accepting predecessor for every vertex. The representative is chosen as the *maximal accepting predecessor* accordingly

to a presupposed linear ordering \prec of vertices (given e.g. by their memory representation). Clearly, if an accepting vertex is its own maximal accepting predecessor, it lies on an accepting cycle. Unfortunately, it can happen that all the maximal accepting predecessors lie out of accepting cycles. In that case, the algorithm removes all accepting vertices that are maximal accepting predecessors of some vertex, and recomputes the maximal accepting predecessors. This is repeated until an accepting cycle is found or there are no more accepting vertices in the graph.

The time complexity of the algorithm is $\mathcal{O}(a^2 \cdot m)$, where a is the number of accepting vertices. One of the key aspects influencing the overall performance of the algorithm is the underlying ordering of vertices used by the algorithm. It is not possible to compute the optimal ordering in parallel, hence heuristics for computing a suitable vertex ordering are used.

2.2 One Way Catch Them Young

The next algorithm [7] is an extended enumerative version of the **One Way Catch Them Young Algorithm** [10]. The idea of the algorithm is to repeatedly remove vertices from the graph that cannot lie on an accepting cycle. The two removal rules are as follows. First, a vertex is removed from the graph if it has no successors in the graph (the vertex cannot lie on a cycle), second, a vertex is removed if it cannot reach an accepting vertex (a potential cycle the vertex lies on is non-accepting). The algorithm performs removal steps as far as there are vertices to be removed. In the end, either there are some vertices remaining in the graph meaning that the original graph contained an accepting cycle, or all vertices have been removed meaning that the original graph had no accepting cycles.

The time complexity of the algorithm is $\mathcal{O}(h \cdot m)$ where $h = h(G)$. Here the factor m comes from the computation of elimination rules while the factor h relates to the number of global iterations the removal rules must be applied. Also note, that an alternative algorithm is obtained if the rules are replaced with their backward search counterparts.

2.3 Negative Cycle Detection

The idea behind the **Negative Cycle Algorithm** [5] is a transformation of the LTL model checking problem to the problem of negative cycle detection. Every edge of the graph outgoing from a non-accepting vertex is labeled with 0 while every edge outgoing from an accepting vertex is labeled with -1 . Clearly, the graph contains a negative cycle if and only if it has an accepting cycle.

The algorithm exploits the *walk to root* strategy to detect the presence of a negative cycle. The strategy involves construction of the so called *parent*

graph that keeps the shortest path to the initial vertex for every vertex of the graph. The parent graph is repeatedly checked for the existence of the path. If the shortest path does not exist for a given vertex, then the vertex is a part of negative, thus accepting, cycle. The worst case time complexity of the algorithm is $\mathcal{O}(n \cdot m)$.

2.4 Back-Level Edges

An edge (u, v) is called a *back-level edge* if it does not increase the distance of the target vertex v from the initial vertex of the graph. The key observation connecting the cycle detection problem with the back-level edge concept, as used in the **Back-Level Edges Algorithm** [1], is that every cycle contains at least one back-level edge. Back-level edges are, therefore, used as triggers to start a procedure that checks whether the edge is a part of an accepting cycle. However, this is too expensive to be done completely for every back-level edge. Therefore, several improvements and heuristics are suggested and integrated within the algorithm to decrease the number of tested edges and speed-up the cycle test.

The BFS procedure which detects back-level edges runs in time $\mathcal{O}(m+n)$. In the worst case, each back-level edge has to be checked to be a part of a cycle, which requires linear time $\mathcal{O}(m+n)$ as well. Since there is at most m back-level edges, the overall time complexity of the algorithm is $\mathcal{O}(m \cdot (m+n))$.

2.5 Cycle Locality

All the algorithms allow for an efficient implementation on a parallel architecture. The implementation is based on partitioning the graph (its vertices) into disjoint parts. Suitable partitioning is important to benefit from parallelisation.

One particular technique, that is specific to automata based LTL model checking, is *cycle locality preserving* problem decomposition. The graph (product automaton) originates from synchronous product of the property and system automata. Hence, vertices of product automaton graph are ordered pairs. An interesting observation is that every cycle in a product automaton graph emerges from cycles in system and property automaton graphs. Let A, B be Büchi automata and $A \otimes B$ their synchronous product. If C is a strongly connected component in the automaton graph of $A \otimes B$, then A -projection of C and B -projection of C are (not necessarily maximal) strongly connected components in automaton graphs of A and B , respectively.

As the property automaton originates from the LTL formula to be verified, it is typically quite small and can be pre-analysed. In particular, it is possi-

ble to identify all strongly connected components of the property automaton graph. A partition function may then be devised, that respects strongly connected components of the property automaton and therefore preserves cycle locality. The partitioning strategy is to assign all vertices that project to the same strongly connected component of the property automaton graph to the same sub-problem. Since no cycle is split among different sub-problems it is possible to employ localised Nested DFS algorithm to perform local accepting cycle detection simultaneously.

Yet another interesting information can be drawn from the property automaton graph decomposition. Maximal strongly connected components can be classified into three categories:

Type F: (*Fully Accepting*) Any cycle within the component contains at least one accepting vertex. (There is no non-accepting cycle within the component.)

Type P: (*Partially Accepting*) There is at least one accepting cycle and one non-accepting cycle within the component.

Type N: (*Non-Accepting*) There is no accepting cycle within the component.

Realizing that vertex of a product automaton graph is accepting only if the corresponding vertex in the property automaton graph is accepting it is possible to characterise types of strongly connected components of product automaton graph according to types of components in the property automaton graph. This classification of components into types N , F , and P can be used to gain additional improvements that may be incorporated into the above algorithms.

Chapter 3

Implementation Techniques

It is a well known fact, that a distributed-memory, parallel algorithm is straightforwardly transformed into a shared-memory one. However, there are several inefficiencies involved in this direct translation. Several traits of shared memory architecture may be leveraged to improve real-world performance of such implementations. In this section, we present our approaches at the challenges of shared-memory architecture and its specific characteristics.

3.1 Shared Memory Platform

We work with a model based on threads that share all memory, although they have separate stacks in their shared address space and a special thread-local storage to store thread-private data. Our working environment is POSIX, with its implementation of threads as lightweight processes. Switching contexts among different threads is cheaper than switching contexts among full-featured processes with separate address spaces, so using more threads than there are CPUs in the system incurs only a minor penalty.

Critical Sections, Locking and Lock Contention. In a shared-memory setting, access to memory, that may be used for writing by more than a single thread, has to be controlled through use of mutual exclusion, otherwise, race conditions will occur. This is generally achieved through use of a “mutual exclusion device”, so-called mutex. A thread wishing to enter a critical section has to lock the associated mutex, which may block the calling thread if the mutex is locked already by some other thread. An effect called resource or lock contention is associated with this behaviour. This occurs, when two or more threads happen to need to enter the same critical section (and therefore lock the same mutex), at the same time. If critical sections are long or they are entered very often, contention starts to cause observable performance degradation, as more and more time is spent waiting for mutexes.

Role of Processor Cache. There are two fairly orthogonal issues associated with processor cache. First, cache coherence which is implemented by hardware, but its efficiency is affected by programmer, and cache efficiency, which mostly depends on data structures and algorithms employed.

Cache coherence poses an efficiency penalty when there are many processors writing to same area of memory. This is largely avoided by the distributed algorithm, however, locking and access to shared data structures have no other choice. Cache coherence on modern architectures works at a level of cache lines, roughly 64 byte chunks of memory that are fetched from main memory into cache at once.

Modern mutex implementations ensure that the mutex is the only thing present on a given cache line, so it does not affect other data, and, more importantly, it ensures that two mutexes never share a cache line, which would pose a performance penalty.

Recent development in multi-core platforms deals with cache coherence problem in a different, more efficient manner, namely, by sharing the level two cache among two or more cores, therefore reducing the cache coherence overhead significantly. Yet, with the current state of technology, this still does not mitigate the overhead completely.

3.2 Implementing Algorithms in Shared Memory

The above considerations bring us to the actual algorithm implementation and the associated techniques we came up with. They are all designed to reduce communication overhead, exploiting traits of shared memory systems that are not available in distributed memory environments. Consequently, the main goal is to improve scalability of the implementation, which is inversely proportional to communication overhead and its growth with increasing number of threads. That said, keeping in mind the possibility to scale beyond shared memory systems, we try to keep the implementation in a shape that would make a combined tool to work efficiently on clusters of multi-CPU machines achievable.

When we venture into a strictly shared-memory implementation, one may pose a question, whether a different approach of using a standard serial algorithm modified to allow parallelisation at lower levels of abstraction would give a scalable, efficient program for multi-CPU and/or multi-core systems. Our efforts at extracting such a micro-parallelism in our code-base have been largely fruitless, due high synchronisation cost relative to amount of work we were able to perform in parallel. Although we intend to do more research on this topic, we do not expect significant results. This is partially supported by the observation, that modern CPUs and compilers already do a fair amount of work at extracting micro-parallelism from the code to execute efficiently on the internally parallel microarchitecture of a single core.

```

type FIFO of T:
  type Node:
    buffer: array of T
    next: pointer to Node
    read, write: integer
    nodeSize: integer (size of buffer)
    head, tail: pointer to Node
    writeLock: mutex

```

Figure 3.1: FIFO representation

In the following sections, we explore the possibilities to build on existing distributed memory approaches, in the vein of statically-partitioned graphs, reducing the overhead using idioms only possible due to locality of memory.

3.3 Communication

Generally, in a distributed computation, all communication is accomplished by passing messages – eg. using a library like MPI for cluster message passing. However, in communication-intensive programs, or those sensitive to communication delay, using general-purpose message passing interface is fairly inefficient.

In shared memory, most of the communication overhead can be eliminated by using more appropriate communication primitives, like high-performance, contention- and lock- free FIFOs (First In, First Out queues). We have adopted a variant of the two-lock algorithm – a decent compromise between performance on one hand and simplicity and portability on the other – presented in [17]. Our modifications involve improved cache-efficiency (by using a linked list of memory-continuous blocks, instead of linked list of single items) and only using a single write-lock, instead of a pair of locks, one for reading and one for writing, since there is ever only one thread reading, while there may be several trying to write.

Every thread involved in the computation owns a single instance of the FIFO and all messages for this thread are pushed onto this single queue. This may introduce a source of resource-contention, where many processes are trying to append messages to a single queue, but considering the message distribution in our system, this turns out to be a negligible problem in practice. With different patterns of communication, a complete lock-free design may be more appropriate (one is given in [17]).

Representation and pseudocode for enqueue and dequeue algorithms are found in Figures 3.1, 3.2 and 3.3, respectively. The correctness, linearizability and liveness proofs as given in [17] are straightforwardly adapted to our implementation and thus left out.

Require: f is a FIFO of T instance, x of type T is an element to enqueue
Ensure: f contains x as its last element

```

lock(  $f.writeMutex$  )
if  $f.tail.write = f.nodeSize$  then
     $t \leftarrow$  newly allocated Node, all fields 0
else
     $t \leftarrow f.tail$ 
end if
 $t.buffer[t.write] \leftarrow x$ 
 $t.write \leftarrow t.write + 1$ 
if  $f.tail \neq t$  then
     $f.tail.next = t$ 
     $f.tail = t$ 
end if
unlock(  $f.writeMutex$  )

```

Figure 3.2: FIFO enqueue

Require: f is a non-empty FIFO instance
Ensure: front element of f is dequeued and then returned

```

if  $f.head.read = f.nodeSize$  then
     $f.head \leftarrow f.head.next$ 
end if
 $f.head.read \leftarrow f.head.read + 1$ 
return  $f.head.buffer[f.head.read - 1]$ 

```

Figure 3.3: FIFO dequeue

Alternatives to our implementation, which may be more appropriate in different settings, include a ring-buffer FIFO implementation (if there is a bound on the amount of in-flight data known beforehand, the ring-buffer implementation may be more efficient) and possibly an algorithm based on swapping incoming and outgoing queues (which could be easily implemented as a pointer swap). The latter gives results comparable to the described FIFO method, although the code and locking behaviour is much more complex and error-prone, which made us opt for the simpler FIFO implementation.

3.4 Memory Allocation

In a distributed computation, every process has simply its own memory which it fully manages. In a shared memory, however, we prefer to manage the memory as a single shared area, since an equal partitioning of available memory and separate management may fall short of efficient resource usage. However, this poses some challenges, especially in allocation-intensive environment like ours.

First, a naïve approach of protecting the allocation routines with a simple mutual exclusion is highly prone to resource contention. Fortunately, modern general-purpose allocator implementations refrain from this idea and have a generally non-contending behaviour on allocation. However, releasing memory back for reuse is more complex to achieve without introducing contention, in a setting where it is often the case that thread other than the one allocating the chunk tries to release it.

There are known general-purpose solutions to this problem, eg. [16], however they are currently not in widespread use, therefore we have to refrain from the above-mentioned pattern of releasing memory from different than allocating thread, in order to avoid contention and the accompanying slowdown.

The message-passing implementation we employ is pointer-based, in other words, the message sent is only a pointer and the payload (actual interesting message content) is allocated on the shared heap and it may be either reused or released by the receiving thread. Observe however, that releasing the associated memory in the receiving thread will introduce the situation which we are trying to avoid.

We side-step the issue by adding a new communication FIFO to each thread (recall that our communication induces only low overhead and virtually no contention). When a receiving thread decides that the message content needs to be disposed of, instead of doing it itself, sends the message back to the originating thread using the second FIFO. The originating thread then, at convenient intervals, releases the memory in a single batch, having an interesting side-effect of slightly improving cache-efficiency.

3.5 Efficient Termination Detection

Since our algorithms rely on work distribution among several largely independent threads, similarly to a distributed algorithm, we need a specific algorithm for shared-memory termination detection, that would pose minimal overhead and minimal serialisation.

One possible solution is presented in [15], which does not use locking and relies on the system to provide an enqueue-with-wakeup primitive. However, in our system, we have primitives available that support a somewhat different approach: implementation of sleeping/wakeup primitives already relies on locking and we leverage this inherent locking in our termination detection algorithm.

The POSIX threading library offers a mechanism called “condition signalling”, which we use to implement thread sleeping and wakeup. A “condition” is a device that allows to be waited-for by its owning thread and “signalling a condition” from another thread will cause the waiting thread to wake up and continue execution. However, this device in itself is race-prone, since the condition may be signalled just before the owner goes to sleep, leading to a deadlock – another signal may never come. Therefore, the condition is always protected by a mutex, which is always locked through the execution of the owner thread and is only atomically unlocked when the thread enters sleep state and atomically reclaimed before waking up.

Since the available mutex implementation allows a lock-or-fail behaviour, as opposed to lock-or-wait which is usually employed for protecting critical sections, we can use the condition device to implement an efficient termination detection algorithm.

Observe, that at any time when a thread is idle, its condition-protecting mutex is unlocked and conversely, whenever the thread is busy, this mutex is locked. So the termination detection algorithm first tries to lock condition mutexes of all worker threads, one by one, using the lock-or-fail behaviour. Then, it proceeds to check the queues. If it succeeded locking all threads and all queues are empty, termination has occurred. Pseudocode for the algorithm is shown in Figure 3.4.

We run the termination detection in a dedicated scheduler thread, which also wakes up threads that have pending work – i.e., if it has successfully grabbed any locks, queues belonging to those locked threads are checked, and if any is found to be non-empty, the thread is awakened. After every run, all grabbed locks are released again.

Moreover, although this algorithm works correctly as-is, it is rather inefficient if left running in a loop. Therefore, the scheduler thread goes to sleep after every iteration, and is woken up by any worker thread that goes idle. This requires a slight modification to the algorithm above, since it adds a race-condition, where the last thread going to sleep wakes up the scheduler, which then runs the algorithm before the calling thread manages to go to

Require: *threads*: array of Thread, Thread contains *idleMutex* and *idleCondition*, *fifo*

Ensure: termination has occurred iff true is returned

mutex: Mutex, *cond*: Condition, *held*: array of Boolean

busy \leftarrow false

for *t* in *threads* **do**

if trylock(*t.idleMutex*) **then**

held[*t*] \leftarrow true

else

held[*t*] \leftarrow false

busy \leftarrow true

end if

end for

for *t* in *threads* **do**

if not empty(*t.fifo*) **then**

busy \leftarrow true

if *held*[*t*] **then**

 signal(*t.idleCondition*)

end if

end if

end for

for *t* in *threads* **do**

 unlock(*t.idleMutex*)

end for

return not *busy*

Figure 3.4: Termination Detection in Shared Memory

sleep, assuming termination did not happen and going to sleep, at which point the system deadlocks, as everyone is idle.

An alternative approach would be to synchronously execute the termination detection algorithm in the thread that has become idle, but due to the nature of the system, the above is more practical code-wise and only incurs very insignificant overhead.

Chapter 4

Implemented Algorithms

4.1 Nested DFS

The dual-threaded Nested DFS algorithm we implemented is based on the design presented in [11], i.e. the outer DFS runs in one of the threads, while another thread executes the nested DFS. There is a FIFO from the outer thread to inner and the outer DFS pushes accepting states to the FIFO in post-order. The pseudocode for this approach can be found in Figures 4.1 and 4.2.

Since there is a clear master-slave relationship between the two threads, the termination conditions are very simple. When the master (outer DFS) finishes its work, it waits for the slave thread (nested DFS) to become idle and then terminates.

4.2 Nested DFS with VCL

Further parallelisation of Nested DFS can be obtained from SCC decomposition of the property automaton. This yields, as described in Section 2.5, components of the product automaton. Out of those components, only those of type P (partially accepting) need to be verified using Nested DFS. The N (non-accepting) components may be explored using BFS reachability and those of type F (fully accepting) may be checked with single cycle detection, which, since we prefer BFS order, may be accomplished with single pass of OWCTY elimination.

This makes both N and F component checking easily parallelisable, as those may be checked using BFS-based algorithms. It also adds parallelism for the case of multiple P components, since Nested DFS can be run in parallel on each of these independent components. There is however no benefit for properties, which are composed of single P component. The P component checking may be done using either parallel or serial Nested DFS.

The current implementation only differentiates between N and non- N

```

procedure NESTEDDFS( $S, I$ )
  if  $S \in \text{NestedSeen}$  then return
  else
     $\text{NestedSeen} \leftarrow \text{NestedSeen} \cup S$ 
  end if
  if  $I = \text{Invalid}$  then
    for  $\text{succ} \leftarrow \text{successor of } S$  do
      NESTEDDFS( $\text{succ}, S$ )
    end for
  else
    for  $\text{succ} \leftarrow \text{successor of } S$  do
      if  $S = I$  then
        echo “accepting cycle found”
      else
        NESTEDDFS( $\text{succ}, I$ )
      end if
    end for
  end if
end procedure

procedure OUTERDFS( $\text{state}$  instance of State,  $\text{fifo}$  instance of FIFO)
  if  $\text{state} \in \text{Seen}$  then return
  end if
   $\text{Seen} \leftarrow \text{Seen} \cup \text{state}$ 
  for  $\text{succ} \leftarrow \text{successor of } \text{state}$  do
    OUTERDFS( $\text{state}$ )
    if  $\text{state}$  is accepting then
       $\text{fifo.push}(\text{state})$ 
      wake up other end of  $\text{fifo}$ 
    end if
  end for
end procedure

```

Figure 4.1: Algorithmic procedures of Nested DFS

```

procedure NESTEDDFS THREAD(fifo instance of FIFO)
  if fifo not empty then
     $S \leftarrow (fifo.pop)$ 
    NESTEDDFS( $S, Invalid$ )
  else
    sleep
  end if
end procedure
procedure MAIN(initial initial state of the system)
  fifo  $\leftarrow$  new instance of FIFO
  asynchronously execute NESTEDDFS THREAD(fifo)
  OUTERDFS(initial, fifo)
  terminate nested DFS thread
end procedure

```

Figure 4.2: Nested DFS

components, using BFS reachability for N components and Nested DFS on the rest. All of the N components in the system are merged into one pool of states, statically partitioned among available reachability worker threads. Each P and F component gets its own instance of dual-thread Nested DFS.

This setting requires somewhat more complex termination detection than the dual-core Nested DFS variant. All the worker threads take part in the termination detection algorithm given in Section 3.5. Obviously, the termination detection cannot be done separately over the identified components, since there are edges crossing their boundaries.

Furthermore, it is a nontrivial problem to determine optimal number of threads to use with VCL. In the experimental section, we have allocated one CPU per non- N component and rest was allocated for the shared pool of N components. This approach however depends upon having enough CPUs to cover all the components, and in case there are almost as many non- N components as CPUs, this approach becomes suboptimal.

4.3 OWCTY

As can be seen from the pseudocode (refer to Figure 4.3), the main OWCTY loop consists of few steps, namely, reachability, elimination and reset. All of them can be parallelised, but only on their own, which requires a barrier after each of them. Only reachability and elimination run in parallel in the current code, reset is to be implemented.

The algorithm uses a BFS state space visitor to implement both reachability and elimination. The underlying BFS is currently implemented using a partition function, i.e., every state is unambiguously assigned to one of

```

Require: initial is initial state
 $S \leftarrow \text{REACHABILITY}(\textit{initial})$ 
 $\textit{old} \leftarrow \emptyset$ 
while  $S \neq \textit{old}$  do
     $\textit{old} \leftarrow S$ 
     $S \leftarrow \text{RESET}(S)$ 
     $S \leftarrow \text{REACHABILITY}(S)$ 
     $S \leftarrow \text{ELIMINATION}(S)$ 
end while
return  $S \neq \emptyset$ 

```

Figure 4.3: OWCTY Pseudocode

the threads. The framework in which the algorithm is implemented offers a multi-threaded BFS implementation based on this kind of state-space partitioning. The algorithm itself is only presented with resulting transition and node-expansion events, unconcerned with the partitioning or communication details.

The barriers are implemented using the termination detection algorithm presented – the computation is initiated by the main thread and the termination detection is then executed in this same thread, which also doubles as a scheduler. When the step terminates, the main thread prepares the next step, spawns the worker threads and initiates the computation again. Since the hash table is always thread-private, i.e. owned exclusively by a single thread, the main thread has to transfer the hash table among different threads in the serial portion of computation. This is nonetheless done cheaply (few pointer operations only) so is probably not worth parallelising.

4.4 Nested DFS and OWCTY

Since OWCTY over weak automata behaves linearly, the VCL decomposition can be used to partition the state space into a set of weak components (N and F) and the remaining P components. The weak components can be all merged together into a shared pool (like with N components in the Nested DFS with VCL case) and a single-pass OWCTY can be run on this part of the graph. Nested DFS can be used for the P components.

Algorithm designed this way is asymptotically optimal for all inputs, while using the most scalable method available for any given portion of the product graph.

In comparison with VCL, it should achieve better load-balancing over the weak part of the graph, while reducing transition locality among threads (i.e., it requires more communication than VCL, but achieves more even distribution of states to processors). It also, compared to Nested DFS with

VCL, makes the problem of thread allocation slightly easier to tackle, since all of the weak part of the graph can be allocated to an arbitrary number of cores (whereas in VCL, the N and F subgraphs require separate core allocation).

Note that this is only a theoretical result, as the proposed scheme has not been implemented so far.

Chapter 5

Experiments

5.1 Methodology

The main testing machine we have used is a 16-way AMD Opteron 885 (8 CPU units with 2 cores each). All timed programs were compiled using gcc 4.1.2 20060525 (Red Hat 4.1.1-1) in 32-bit mode, using -O3. This limits addressable memory to 3GB, which was enough for our testing. The machine has 64GB of memory installed, which means that none of the runs were affected by swapping.

For this paper, our main concern is speed and scalability, therefore we focus on these two parameters. Measurement was done using standard UNIX `time` command, which measures real and cpu times used by program. Note that the cpu time given in tables equals to a sum of times spent by individual processors, thus for parallel computations the value of cpu time should exceed the value of real time.

For the experimental evaluation we implemented algorithms upon the state generator from DiVINE [3]. All the models we have used are listed in Table 5.1 including the verified properties. The models come from the **BEEM** database [18] that contains the models in DiVINE-native modeling language as well as in ProMeLa. We used ProMeLa models for comparison with the **SPIN** model checker. The models are not extremely large, although, their size is sufficient for the time spent on parsing and initialization to be negligible.

5.2 Results

First, we have measured run-times of algorithms presented in Section 2 that were implemented using DiVINE framework and `mpich2` library compiled for shared-memory architecture. As shown in Figure 5.1 these implementations do not exhibit desired scalability on shared-memory architecture, even though they all scale well in a distributed memory environment. Some algorithms have scaled up to 4 cores, but using more cores did not bring

Acronym	Description	Property (LTL formula)
<i>elevator₁</i>	Motivated by elevator promela model from distribution of SPIN. The cab controller chooses the next floor to be served as the next requested floor in the direction of the last cab movement. If there is no such floor then the controller consider the oposite direction. (3 floors)	If level 0 is requested, the cab passes the level without serving it at most once. $G(r0 \implies (\neg l_0 U (l_0 U (\neg l_0 U (l_0 U (l_0 \wedge open))))))$
<i>elevator₂</i>	Same model as elevator, with slightly adjusted parameters to increase state space size. No formula was used with this model.	<i>N/A</i>
<i>leader</i>	Leader election algorithm based on filters. A filter is a piece of code that satisfy the two following conditions: a) if <i>m</i> processes enter the filter, then at most <i>m</i> /2 processes exit; b) if some process enter the filter, then at least one of them exits. (5 processes)	Eventually a leader will be elected. $F(leader)$
<i>rether</i>	Software-based, real-time Ethernet protocol whose purpose is to provide guaranteed bandwidth and deterministic, periodic network access to multimedia applications over commodity Ethernet hardware. It is a contention-free token bus protocol for the datalink layer of the ISO protocol stack. (5 Nodes)	Infinitely many NRT actions of Node 0. $G(F(nact0))$
<i>peterson</i>	Peterson's mutual exclusion protocol for N processes. (N=4)	Someone is in critical section infinitely many times. $G(F(SomeoneInCS))$
<i>anderson</i>	Anderson's mutual exclusion protocol for N processes. (N=6)	<i>N/A</i>

Table 5.1: Models and verified properties.

any speedup. Under this setting, every MPI node is executed in a separate process.

We have performed more experiments to evaluate the efficiency of techniques introduced in Section 3. We have implemented parallel breadth-first search based reachability and the OWCTY algorithm. Run-times of the thread-optimised BFS reachability and of the thread-optimised implementation of OWCTY algorithm are reported in Figure 5.2.

The thread-optimised implementations display better scalability behaviour, since adding cores reduces computation time at least up to 12 cores, for some models even to 16 cores. Between 12 and 16 cores, the communication overhead reaches a limiting threshold, so adding more does not bring any further speedup and may even impede a slight performance setback.

The actual threshold and curve steepness is generally affected by the partition function used, as well as relative cost of cross transitions. The latter issue partially explains why the MPI versions of algorithms have scalability problems, since the cross transition cost is in this case much higher than in the multi-threaded version.

5.3 Nested DFS and VCL

We have measured runtimes of single-threaded, dual-threaded and VCL-based Nested DFS, using 1, 2 and 16 threads, respectively. The results are to be seen in Table 5.2. Also, a 15-thread (+ 1 management thread) OWCTY runtime is given in the table. These times however come from a different version of the code-base than those in the above comparisons, so the OWCTY timings were re-done with the same version as Nested DFS ones, although the timings did not exhibit any remarkable deviations.

We can see, that performance of VCL varies with model used. In practice, OWCTY seems to outperform both dual-core Nested DFS as well as VCL-based one. This still does not invalidate the advantage of NDFS where cycle (and therefore a counterexample) may be found without generating the full state space. Also, as of this writing, we do not have a dual-core version of OWCTY with satisfactory performance, so for dual-core machines, NDFS is probably the first choice.

5.4 Comparison with SPIN

Since the multi-core version of SPIN was not publicly available, in order to make a direct comparison, we run a single reachability on the product automaton graph with **SPIN**. As SPIN was running only the first procedure of the Nested DFS algorithm we get a good lower bound on runtime of the multi-core **SPIN** implementation. **SPIN** was used with parameters `-m10000000 -w27` to get the best performance. We have not observed any

Model	single-core NDFS	dual-core NDFS	VCL NDFS	OWCTY
elevator ₁	1:12.4	0:51.4	0:35.1	0:27.8
leader	0:16.8	0:9.6	0:11.5	0:7.2
peterson	0:30.8	0:22.8	0:19.9	0:9.4
rether	1:4.5	0:48.0	0:26.8	0:21.3

Table 5.2: Parallel Nested DFS in DiVINE.

Model	SPIN reachability	BFS reachability	OWCTY
elevator ₁	0:14.4	0:12.1	0:26.8
peterson	0:7.4	0:4.2	0:9.2

Table 5.3: Comparison with SPIN

performance penalty from using bigger stack or hash table than strictly necessary.

Table 5.3 gives runtimes for **SPIN**, multi-threaded BFS reachability, and OWCTY cycle detection algorithm, both performed on 16 cores.

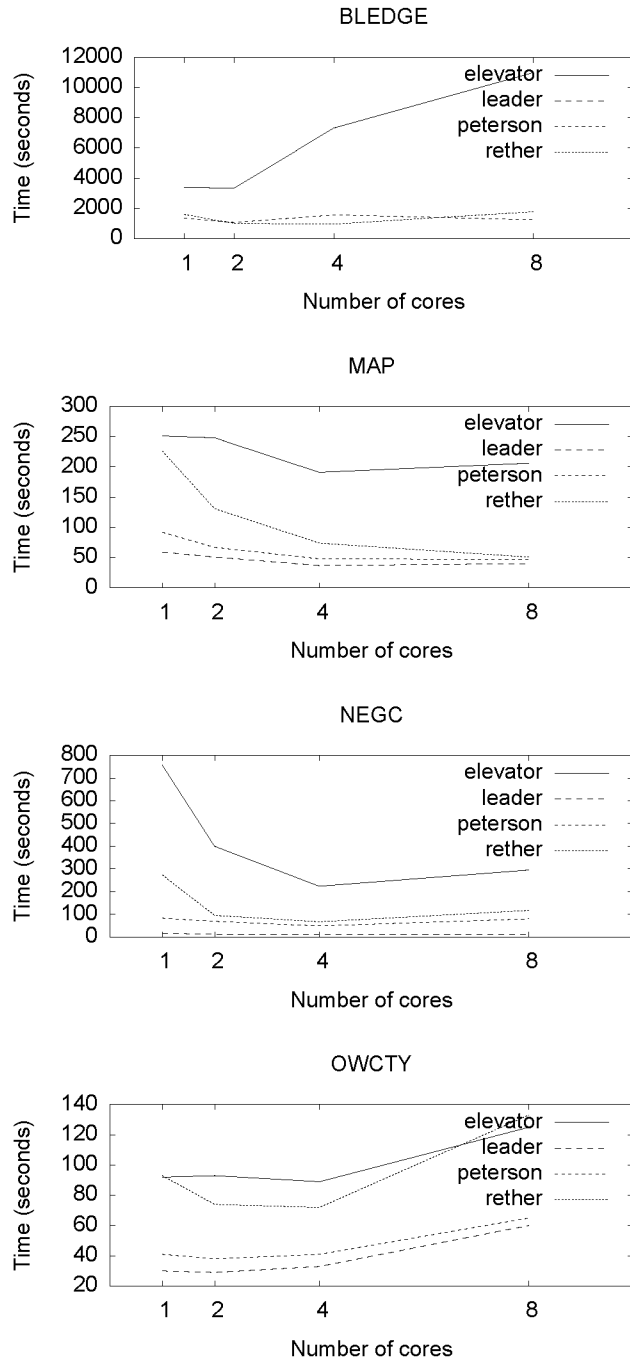


Figure 5.1: Scalability of BLEDGE, MAP, NEGC, and OWCTY algorithms implemented using DiVINE and MPI compiled for shared-memory architecture.

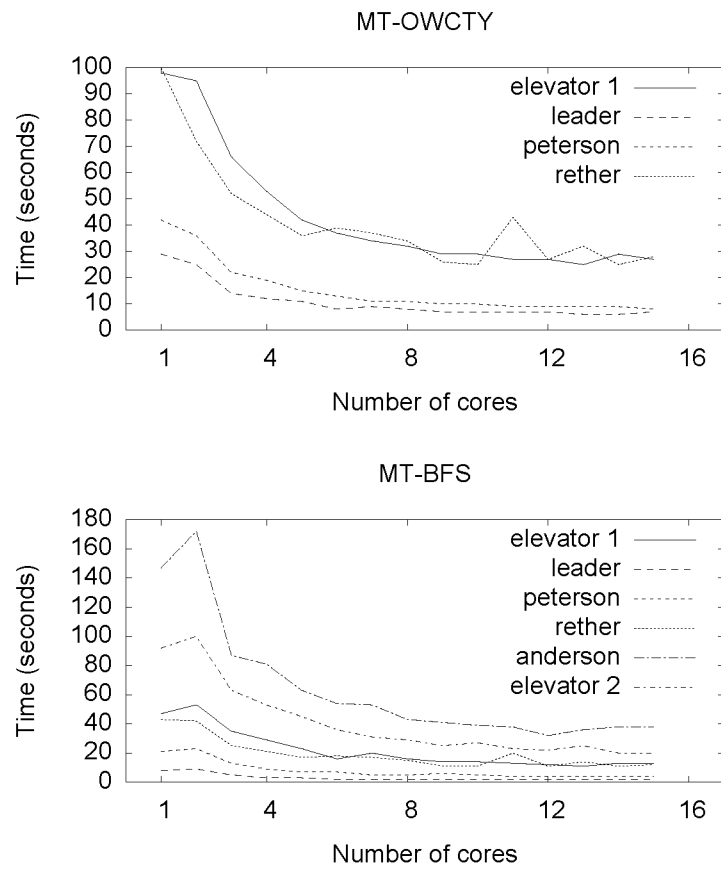


Figure 5.2: Scalability of multi-threaded OWCTY and BFS reachability.

Chapter 6

Conclusions

The current solutions available based on Nested DFS perform really well for most cases. In general, at least a dual-core version is always available, for all kinds of LTL verification problems. For many properties, more CPUs and/or cores can be employed using VCL partitioning of the state space.

On the contrary, the OWCTY algorithm has super-linear complexity for the general case, although for majority of LTL properties verified in practice, it is linear. Scalability of OWCTY is superior to that of either of the Nested DFS implementations. It does not depend on property decomposition to be nontrivial to scale beyond 2 cores. However, for cases where both original and dual-thread Nested DFS shine – automata composed of a single partially accepting component, OWCTY may be outperformed by the dual-core version of Nested DFS.

While Nested DFS seems like the ultimate choice for dual-core systems right now, OWCTY appears to be the superior candidate when multitude of cores is available, or when ability to execute efficiently in distributed memory becomes a factor.

For reachability analysis, parallel BFS seems like the best pick, since it both performs and scales well, while offering the possibility of efficient distributed memory implementation.

From the profiling work we have done, it is clear that the main execution time bottleneck of DiVINE is its state generator. Improvements in this area should reduce the absolute running times, but will likely negatively affect relative scalability. Therefore, we will continue to work on reducing parallel execution overhead, to maintain or even improve current scalability.

In the pursue of scalability, we also intend to explore alternative approaches to state-space partitioning, non-partitioning approaches and usefulness of load-balancing in this context.

Bibliography

- [1] J. Barnat, L. Brim, and J. Chaloupka. Parallel Breadth-First Search LTL Model-Checking. In *Proc. 18th IEEE International Conference on Automated Software Engineering*, pages 106–115. IEEE Computer Society, 2003.
- [2] J. Barnat, L. Brim, and P. Ročkal. Scalable Multi-Core LTL Model-Checking. In *SPIN 2007, to appear*, LNCS, 2007.
- [3] J. Barnat, L. Brim, I. Černá, P. Moravec, P. Ročkal, and P. Šimeček. DiVinE – A Tool for Distributed Verification (Tool Paper). In *Computer Aided Verification*, volume 4144/2006 of *LNCS*, pages 278–281. Springer Berlin / Heidelberg, 2006.
- [4] L. Brim, I. Černá, P. Moravec, and J. Šimša. Accepting predecessors are better than back edges in distributed ltl model-checking. In *5th International Conference on Formal Methods in Computer-Aided Design (FMCAD'04)*, volume 3312 of *LNCS*, pages 352–366. Springer-Verlag, 2004.
- [5] L. Brim, I. Černá, P. Krčál, and R. Pelánek. Distributed LTL model checking based on negative cycle detection. In *Proc. of Foundations of Software Technology and Theoretical Computer Science (FST TCS 2001)*, volume 2245 of *LNCS*, pages 96–107. Springer, 2001.
- [6] L. Brim, I. Černá, P. Moravec, and J. Šimša. How to order vertices for distributed ltl model-checking based on accepting predecessors. In *Proceedings of the 4th International Workshop on Parallel and Distributed Methods in verification (PDMC 2005)*, volume 1M0545, pages 1–12, 2005.
- [7] I. Černá and R. Pelánek. Distributed explicit fair cycle detection (set based approach). In T. Ball and S.K. Rajamani, editors, *Model Checking Software. 10th International SPIN Workshop*, volume 2648 of *Lecture Notes in Computer Science*, pages 49 – 73. Springer Verlag, 2003.
- [8] E. Clarke, O. Grumberg, and D. Peled. *Model Checking*. MIT press, 1999.

- [9] C. Courcoubetis, M.Y. Vardi, P. Wolper, and M. Yannakakis. Memory-Efficient Algorithms for the Verification of Temporal Properties. *Formal Methods in System Design*, 1:275–288, 1992.
- [10] K. Fisler, R. Fraer, G. Kamhi, M. Y. Vardi, and Z. Yang. Is there a best symbolic cycle-detection algorithm? In *Proc. Tools and Algorithms for the Construction and Analysis of Systems*, volume 2031 of *LNCS*, pages 420–434. Springer-Verlag, 2001.
- [11] G. Holzmann. The Design of a Distributed Model Checking Algorithm for SPIN. In *FMCAD, Invited Talk*, 2006.
- [12] G. J. Holzmann. *The Spin Model Checker: Primer and Reference Manual*. Addison-Wesley, 2003.
- [13] G.J. Holzmann, D. Peled, and M. Yannakakis. On Nested Depth First Search. In *The SPIN Verification System*, pages 23–32. American Mathematical Society, 1996. Proc. of the 2nd SPIN Workshop.
- [14] C. Inggs and H. Barringer. Ctl* model checking on a shared memory architecture. *Formal Methods in System Design*, 29(2):135–155, 2006.
- [15] Ho-Fung Leung and Hing-Fung Ting. An optimal algorithm for global termination detection in shared-memory asynchronous multiprocessor systems. *IEEE Transactions on Parallel and Distributed Systems*, 8(5):538–543, 1997.
- [16] Maged M. Michael. Scalable lock-free dynamic memory allocation. *SIGPLAN Not.*, 39(6):35–46, 2004.
- [17] Maged M. Michael and Michael L. Scott. Simple, fast, and practical non-blocking and blocking concurrent queue algorithms. In *Symposium on Principles of Distributed Computing*, pages 267–275, 1996.
- [18] R. Pelánek. BEEM: BEncmarks for Explicit Model checkers. <http://anna.fi.muni.cz/models/index.html>, February 2007.
- [19] R. Tarjan. Depth First Search and Linear Graph Algorithms. *SIAM Journal on Computing*, pages 146–160, Januar 1972.