

RESEARCH1

2
3
4
5
6
7
8
A tale of lock-free Agents: Towards Software
Transactional Memory in parallel Agent-Based
Simulation

Jonathan Thaler* and Peer-Olaf Siebers9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

Correspondence:

nathan.thaler@nottingham.ac.uk

University of Nottingham, 7301

Vollaton Rd, NG8 1BB

Nottingham, UK

Full list of author information is

available at the end of the article

Abstract

With the decline of Moore's law and the ever increasing availability of cheap massively parallel hardware, it becomes more and more important to embrace parallel programming methods to implement Agent-Based Simulations (ABS). This has been acknowledged in the field a while ago and numerous research on distributed parallel ABS exists, focusing primarily on Parallel Discrete Event Simulation as the underlying mechanism. However, these concepts and tools are inherently difficult to master and apply and often overkill in case implementers simply want to parallelise their own, custom agent-based model implementation. However, with the established programming languages in the field, Python, Java and C++, it is not easy to address the complexities of parallel programming due to unrestricted side effects and the intricacies of low-level locking semantics. Therefore, in this paper we propose the use of a lock-free approach to parallel ABS using Software Transactional Memory (STM) in conjunction with the pure functional programming language Haskell, which in combination, removes some of the problems and complexities of parallel implementations in imperative approaches.

We present two case studies where we compare the performance of lock-based and lock-free STM implementations in two different well known Agent-Based Models, where we investigate both the scaling performance under increasing number of CPUs and the scaling performance under increasing number of agents. We show that the lock-free STM implementations consistently outperform the lock-based ones and scale much better to increasing number of CPUs both on local machines and on Amazon Cloud Services. Further, by utilizing the pure functional language Haskell we gain the benefits of immutable data and lack of unrestricted side effects guaranteed at compile-time, making validation easier and leading to increased confidence in the correctness of an implementation, something of fundamental importance and benefit in parallel programming in general and scientific computing like ABS in particular.

Keywords: Agent-Based Simulation; Software Transactional Memory; Parallel Programming; Haskell

1 Introduction

The future of scientific computing in general and Agent-Based Simulation (ABS) in particular is parallelism: Moore's law is declining as we are reaching the physical

¹limits of CPU clocks. The only option is going massively parallel due to availability¹
²of cheap massive parallel local hardware with many cores, or cloud services like²
³Amazon EC. This trend has been already recognised in the field of ABS as a research³
⁴challenge for *Large-scale ABMS* [1] was called out and as a substantial body of⁴
⁵research for parallel ABS shows [2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13].⁵

⁶ In this body of work it has been established that parallelisation of autonomous⁶
⁷agents, situated in some spacial, metric environment can be particularly challeng-⁷
⁸ing. The reason for this is that the environment constitutes a key medium for the⁸
⁹agents interactions, represented as a *passive* data structure, recording attributes of⁹
¹⁰the environment and the agents [4]. Thus the problem of parallelising ABS boils¹⁰
¹¹down to the problem of how to synchronise access to shared state without violat-¹¹
¹²ing the causality principle and resource constraints [3, 2]. Various researchers have¹²
¹³developed different techniques where most of them are based on the concept of Par-¹³
¹⁴allel Discrete-Event Simulation (PDES). The idea behind PDES is to partition the¹⁴
¹⁵shared space into logical processes which run at their own speed, processing events¹⁵
¹⁶coming from themselves and other logical processes. To deal with inconsistencies¹⁶
¹⁷there exists a conservative approach which does not allow to process events with¹⁷
¹⁸a lower timestamp than the current time of the logical process; and an optimistic¹⁸
¹⁹approach which deals with inconsistencies through rolling back changes to state.¹⁹

²⁰ Adopting PDES to ABS is challenging as agents are autonomous and thus the²⁰
²¹topology can change in every step, making it hard to predict the topology of logical²¹
²²processes in advance [4] and thus posing a difficult problem for parallelisation in²²
²³general [13]. The work [2, 5] discusses this challenge by giving a detailed and in-²³
²⁴depth discussion of the internals and implementation of their powerful and highly²⁴
²⁵complex PDES-MAS system. The rather conceptual work [3] proposes a general,²⁵
²⁶distributed simulation framework for multiagent systems and addresses a number of²⁶
²⁷key problems: decomposition of the environment, load balancing, modelling, com-²⁷
²⁸munication and shared state variables, which the authors mention as the central²⁸
²⁹problem of parallelisation.²⁹

³⁰ In addition, various distributed simulation environments for ABS have been de-³⁰
³¹veloped and their internals published in research papers: the SPADES system [6]³¹
³²manages agents through UNIX pipes using a parallel sense-think-act cycle employ-³²
³³ing a conservative PDES approach; Mace3J [7] a Java based system running on³³

¹single- or multicore workstations implements a message passing approach to paral-¹
²lelism; James II [8] is also a Java based system and focuses on PDEVs simulation²
³with a plugin architecture to facilitate reuse of models; the well known RePast-HPC³
⁴[9, 10] framework is using a PDES engine under the hood.⁴

⁵ The baseline of this body of research is that parallelisation is possible and we⁵
⁶know how to do it. However, the complexity of these parallel and distributed simu-⁶
⁷lation concepts and toolkits is high and the model development effort is hard [12].⁷
⁸Further, this sophisticated and powerful machinery is not always required as ABS⁸
⁹does not always need to be run in a distributed way but the implementers 'sim-⁹
¹⁰ply' want to parallelise their models locally. Although these existing distributed¹⁰
¹¹ABS frameworks could be used for this, they are overkill and more straightforward¹¹
¹²concepts for parallelising ABS would be appropriate. However, for this case there¹²
¹³does not exist much research, and implementers either resort to the distributed¹³
¹⁴ABS frameworks, implement their own low-level concurrency plumbing which can¹⁴
¹⁵be considerably complex - or simply refrain from using parallelism due to the high¹⁵
¹⁶complexity involved and accept a longer execution time. What makes it worse is¹⁶
¹⁷that parallelism always comes with danger of additional, very subtle bugs, which¹⁷
¹⁸might lie dormant, potentially invalidating significant scientific results of the model.¹⁸
¹⁹Therefore something simpler is needed for local parallelism. Unfortunately, the es-¹⁹
²⁰tablished imperative languages in the ABS field, Python, Java, C++, don't make²⁰
²¹adding parallelism easy, due to their inherent use of unrestricted side effects. Fur-²¹
²²ther, they mostly follow a lock-based approach to concurrency which is error prone²²
²³and does not compose.²³

²⁵ This paper proposes Software Transactional Memory (STM) in conjunction with²⁵
²⁶functional programming in Haskell [14] as a new underlying concept for local par-²⁶
²⁷allelisation of ABS. We hypothesise that by using STM in Haskell, implementing²⁷
²⁸local parallel ABS is considerably easier than with lock-based approaches, less error²⁸
²⁹prone and easier to validate. Although STM exists in other languages as well by²⁹
³⁰now, Haskell was one of the first to natively build it into its core. Further, it has³⁰
³¹the unique benefit that it can guarantee the lack of persistent side effects at com-³¹
³²pile time, allowing unproblematic retries of transactions, something of fundamental³²
³³importance in STM. This makes the use of STM in Haskell very compelling. Our³³

¹hypothesis is supported by [15], which gives a good indication how difficult and¹
²complex constructing a correct concurrent program is and shows how much easier,²
³concise and less error-prone an STM implementation is over traditional locking with³
⁴mutexes and semaphores. Further, it shows that STM consistently outperforms the⁴
⁵lock-based implementation.⁵

⁶
⁷ To the best of our knowledge we are the first to *systematically* discuss the use of,⁷
⁸STM in the context of ABS. However, the idea of applying transactional memory to⁸
⁹simulation in general is not new and was already explored in the work [11], where⁹
¹⁰the authors looked into how to apply Intel's *hardware* transactional memory to¹⁰
¹¹simulations in the context of a Time Warp PDES simulation. The results showed¹¹
¹²that their approach generally outperformed traditional locking mechanisms.¹²

¹³ The master thesis [16] investigates Haskell's parallel and concurrency features to¹³
¹⁴implement (amongst others) *HLogo*, a Haskell clone of the NetLogo [17] simulation¹⁴
¹⁵package, focusing on using STM for a limited form of agent interactions. *HLogo* is¹⁵
¹⁶basically a re-implementation of NetLogos API in Haskell where agents run within¹⁶
¹⁷an unrestricted side effect context (known as IO) and thus can also make use of¹⁷
¹⁸STM functionality. The benchmarks show that this approach does indeed result in¹⁸
¹⁹a speedup especially under larger agent populations. Despite the parallelism aspect¹⁹
²⁰our work share, our approach is rather different: we avoid unrestricted side effects²⁰
²¹through IO within the agents under all costs and explore the use of STM more on²¹
²²a conceptual level rather than implementing an ABS library.²²

²³ The aim of this paper is to experimentally investigate the benefits of using STM²³
²⁴over lock-based approaches for concurrent ABS models. Therefore, we follow [15]²⁴
²⁵and compare the performance of lock-based and STM implementations and expect²⁵
²⁶that the reduced complexity and increased performance will be directly applicable²⁶
²⁷to ABS as well. We present two case studies in which we employ an agent-based²⁷
²⁸spatial SIR [18, 19] and the well known SugarScape [20] model to test our hypothesis.²⁸
²⁹The latter model can be seen as one of the most influential exploratory models in²⁹
³⁰ABS which laid the foundations of object-oriented implementation of agent-based³⁰
³¹models. The former one is an easy-to-understand explanatory model which has³¹
³²the advantage that it has an analytical theory behind it which can be used for³²
³³verification and validation.³³

¹ The contribution of this paper is a systematic investigation of the usefulness¹
² of STM over lock-based approaches, therefore giving implementers a new method²
³ of locally parallelising their own implementations without the overhead of a dis-³
⁴ tributed, parallel PDES system or the error-prone low-level locking semantics of a⁴
⁵ custom built parallel implementation. Therefore, our paper directly addresses the⁵
⁶ *Large-scale ABMS* challenge [1], which focuses on efficient modelling and simulat-⁶
⁷ ing large-scale ABS. Further, using STM, which restricts side effects, and makes⁷
⁸ parallelism easier, can help in the validation challenge [1] *H5: Requirement that all*⁸
⁹ *models be completely validated.*⁹

¹⁰ We start with Section 2 where we discuss the concepts of STM and side effects¹⁰
¹¹ in Haskell. In Section 3 we show how to apply STM to ABS in general. Section 4¹¹
¹² contains the first case study using a spatial SIR model whereas Section 5 presents¹²
¹³ the second case study using the SugarScape model. We conclude in Section 6 and¹³
¹⁴ give further research directions in Section 7.¹⁴

¹⁵

¹⁶ 2 Background¹⁶

¹⁷ 2.1 Software Transactional Memory¹⁷

¹⁸ Software Transactional Memory (STM) was introduced by [21] in 1995 as an alter-¹⁸
¹⁹ native to lock-based synchronisation in concurrent programming which, in general,¹⁹
²⁰ is notoriously difficult to get right. This is because reasoning about the interac-²⁰
²¹ tions of multiple concurrently running threads and low level operational details of²¹
²² synchronisation primitives is *very hard*. The main problems are:²²

- ²³ • Race conditions due to forgotten locks;²³
- ²⁴ • Deadlocks resulting from inconsistent lock ordering;²⁴
- ²⁵ • Corruption caused by uncaught exceptions;²⁵
- ²⁶ • Lost wake ups induced by omitted notifications.²⁶

²⁷ Worse, concurrency does not compose. It is very difficult to write two functions²⁷
²⁸ (or methods in an object) acting on concurrent data which can be composed into a²⁸
²⁹ larger concurrent behaviour. The reason for it is that one has to know about internal²⁹
³⁰ details of locking, which breaks encapsulation and makes composition dependent on³⁰
³¹ knowledge about their implementation. Therefore, it is impossible to compose two³¹
³² functions e.g. where one withdraws some amount of money from an account and³²
³³ the other deposits this amount of money into a different account: one ends up³³

¹with a temporary state where the money is in none of either accounts, creating an¹
²inconsistency - a potential source for errors because threads can be rescheduled at²
³any time.³

⁴ STM promises to solve all these problems for a low cost by executing actions⁴
⁵*atomically*, where modifications made in such an action are invisible to other threads⁵
⁶and changes by other threads are invisible as well until actions are committed -⁶
⁷STM actions are atomic and isolated. When an STM action exits, either one of⁷
⁸two outcomes happen: if no other thread has modified the same data as the thread⁸
⁹running the STM action, then the modifications performed by the action will be⁹
¹⁰committed and become visible to the other threads. If other threads have modified¹⁰
¹¹the data then the modifications will be discarded, the action block rolled back and¹¹
¹²automatically restarted.¹²

¹³ STM in Haskell is implemented using optimistic synchronisation, which means¹³
¹⁴that instead of locking access to shared data, each thread keeps a transaction log¹⁴
¹⁵for each read and write to shared data it makes. When the transaction exits, the¹⁵
¹⁶thread checks if it has a consistent view to the shared data by verifying whether¹⁶
¹⁷other threads have written to memory it has read or not.¹⁷

¹⁸ In the paper [22] the authors use a model of STM to simulate optimistic and¹⁸
¹⁹pessimistic STM behaviour under various scenarios using the AnyLogic simulation¹⁹
²⁰package. They conclude that optimistic STM may lead to 25% less retries of trans-²⁰
²¹actions. The authors of [23] analyse several Haskell STM programs with respect to²¹
²²their transactional behaviour. They identified the roll-back rate as one of the key²²
²³metric which determines the scalability of an application. Although STM might²³
²⁴promise better performance, they also warn of the overhead it introduces which²⁴
²⁵could be quite substantial in particular for programs which do not perform much²⁵
²⁶work inside transactions as their commit overhead appears to be high.²⁶

²⁷²⁷

²⁸2.2 Parallelism, Concurrency and STM in Haskell²⁸

²⁹ In our case studies we are using the functional programming language Haskell. The²⁹
³⁰paper of [14] gives a comprehensive overview over the history of the language, how³⁰
³¹it developed and its features and is very interesting to read and get accustomed to³¹
³²the background of the language. Note that Haskell is a *lazy* language which means³²
³³that expressions are only evaluated when they are actually needed.³³

¹2.2.1 Side Effects¹

²One of the fundamental strengths of Haskell is its way of dealing with side effects in²
³functions. A function with side effects has observable interactions with some state³
⁴outside of its explicit scope. This means that the behaviour depends on history and⁴
⁵that it loses its referential transparency character, which makes understanding and⁵
⁶debugging much harder. Examples for side effects are (amongst others): modify a⁶
⁷global variable, await an input from the keyboard, read or write to a file, open a⁷
⁸connection to a server, drawing random numbers, etc.⁸

⁹ The unique feature of Haskell is that it allows to indicate in the *type* of a function⁹
¹⁰that it does have side effects and what kind of effects they are. There are a broad¹⁰
¹¹range of different effect types available, to restrict the possible effects a function can¹¹
¹²have e.g. drawing random numbers, sharing read/write state between functions, etc.¹²
¹³Depending on the type, only specific operations are available, which is then checked¹³
¹⁴by the compiler. This means that a program which tries to read from a file in a¹⁴
¹⁵function which only allows drawing random numbers will fail to compile.¹⁵

¹⁶ Here we are only concerned with two effect types: The IO effect context can be seen¹⁶
¹⁷as completely unrestricted as the main entry point of each Haskell program runs in¹⁷
¹⁸the IO context which means that this is the most general and powerful one. It allows¹⁸
¹⁹all kind of input/output (IO) related side effects: reading/writing a file, creating¹⁹
²⁰threads, write to the standard output, read from the keyboard, opening network²⁰
²¹connections, mutable references, etc. Also the IO context provides functionality²¹
²²for concurrent locks and global shared references. The other effect context we are²²
²³concerned with is STM and indicates the STM context of a function - we discuss it²³
²⁴more in detail below.²⁴

²⁵ A function with a given effect type needs to be executed with a given effect²⁵
²⁶runner which takes all necessary parameters depending on the effect and runs a²⁶
²⁷given effectful function returning its return value and depending on the effect also²⁷
²⁸an effect related result. Note that we cannot call functions of different effect types²⁸
²⁹from a function with another effect type, which would violate the guarantees. A²⁹
³⁰function without any side effect is called *pure*. Calling a *pure* function though is³⁰
³¹always allowed because it has, by definition, no side effects.³¹

³² Although such a type system might seem very restrictive at first, we get a number³²
³³of benefits from making the type of effects we can use explicit. First we can restrict³³

¹the side effects a function can have to a very specific type which is guaranteed at¹
²compile time. This means we can have much stronger guarantees about our program²
³and the absence of potential run time errors. Second, by the use of effect runners,³
⁴we can execute effectful functions in a very controlled way, by making the effect⁴
⁵context explicit in the parameters to the effect runner. 5

6

6

⁷*2.2.2 Parallelism & Concurrency* 7

⁸Haskell makes a very clear distinction between parallelism and concurrency. Paral-⁸
⁹lelism is always deterministic and thus pure without side effects because although⁹
¹⁰parallel code can be run concurrently, it does by definition not interact with data¹⁰
¹¹of other threads. This can be indicated through types: we can run pure functions¹¹
¹²in parallel because for them it doesn't matter in which order they are executed, the¹²
¹³result will always be the same due to the concept of referential transparency. 13

¹⁴Concurrency on the other hand is potentially nondeterministic because of nonde-¹⁴
¹⁵terministic interactions of concurrently running threads through shared data. Al-¹⁵
¹⁶though data in functional programming is immutable, Haskell provides primitives¹⁶
¹⁷which allow to share immutable data between threads. Accessing these primitives is¹⁷
¹⁸only possible from within an `IO` or `STM` context (see below) which means that when¹⁸
¹⁹we are using concurrency in our program, the types of our functions change from¹⁹
²⁰pure to either a `IO` or `STM` effect context. 20

²¹Note that spawning tens of thousands or even millions of threads in Haskell is²¹
²²no problem, because threads in Haskell have a *very* low memory footprint due to²²
²³being lightweight user space threads, also known as green threads, managed by the²³
²⁴Haskell Runtime System, which maps them to physical operating system worker²⁴
²⁵threads [24]. 25

²⁶*2.2.3 STM* 26

²⁷The work of [25, 26] added STM to Haskell which was one of the first programming²⁷
²⁸languages to incorporate STM into its main core and added the ability to compos-²⁸
²⁹able operations. There exist various implementations of STM in other languages as²⁹
³⁰well (Python, Java, C#, C/C++, etc) but we argue, that it is in Haskell with its³⁰
³¹type system and the way how side effects are treated where it truly shines. 31

³²In the Haskell implementation, STM actions run within the `STM` context. This³²
³³restricts the operations to only STM primitives as shown below, which allows to 33

¹enforce that STM actions are always repeatable without persistent side effects because¹
²such persistent side effects (e.g. writing to a file, launching a missile) are not possible²
³in an STM context. This is also the fundamental difference to IO, where all bets are³
⁴off because *everything* is possible as there are basically no restrictions because IO⁴
⁵can run everything. 5

⁶ Thus the ability to *restart* a block of actions without any visible effects is only⁶
⁷possible due to the nature of Haskell's type system: by restricting the effects to STM⁷
⁸only, prevents uncontrolled effects which cannot be rolled back. 8

⁹ STM comes with a number of primitives to share transactional data. Amongst⁹
¹⁰others the most important ones are: 10

- ¹¹ • TVar A transactional variable which can be read and written arbitrarily; 11
- ¹² • TArray A transactional array where each cell is an individual shared data,¹²
¹³ allowing much finer grained transactions instead of e.g. having the whole array¹³
¹⁴ in a TVar; 14
- ¹⁵ • TChan A transactional channel, representing an unbounded FIFO channel; 15
- ¹⁶ • TMVar A transactional *synchronising* variable which is either empty or full.¹⁶
¹⁷ To read from an empty or write to a full TMVar will cause the current thread¹⁷
¹⁸ to block and retry its transaction when the TMVar was updated by another¹⁸
¹⁹ thread. 19

²¹2.2.4 An example 21

²²We provide a short example to demonstrate the use of STM. To show the retry²²
²³semantics more clearly, we use STM within a StateT effect context. A StateT effect²³
²⁴allows to read and write some state, which in this example we simply set to be an²⁴
²⁵Int value. The example code takes a transactional variable TVar which holds an²⁵
²⁶Int and runs within a StateT effect, providing read and write access to an Int,²⁶
²⁷and an STM effect returning an Int: 27

```
28 stmAction :: TVar Int -> StateT Int STM Int 28
28 stmAction v = do
29 -- print a debug output and increment the value in StateT 29
29 Debug.trace "increment!" (modify (+1))
30 -- read from the TVar 30
30 n <- lift (readTVar v) 31
31 -- await a condition: content of the TVar >= 42 31
32 if n < 42 32
33 -- condition not met, therefore retry: block this thread 33
33 -- until the TVar v is written by another thread, then
```

```

1  -- try again
2  then lift retry
3  -- condition met: return content of TVar
4  else return n
5
6  When stmAction is run, it prints an 'increment!' debug message to the console
7  and increments the value in the StateT. Then it awaits a condition for as long as
8  TVar is less than 42 the action will retry whenever it is run. If the condition is met,
9  it will return the content of the TVar. To run stmAction we need to spawn a thread:
10
11 stmThread :: TVar Int -> IO ()
12 stmThread v = do
13   -- the initial state of the StateT effect
14   let s = 0
15   -- run the state with initial value of s (0)
16   let ret = runStateT (stmAction v) s
17   -- atomically run the STM action
18   (a, s') <- atomically ret
19   -- print final result
20   putStrLn("final StateT state      = " ++ show s' ++
21            ", STM computation result = " ++ show a)
22
23 The thread simply runs the StateT effect with the initial value of 0 and then
24 the STM computation through atomically and prints the result to the console.
25 The value of a is the result of stmAction and s' is the final state of the StateT
26 computation. To actually run this example we need the main thread to update the
27 TVar until the condition is met within stmAction:
28
29 main :: IO ()
30 main = do
31   -- create a new TVar with initial value of 0
32   v <- newTVarIO 0
33   -- start the stmThread and pass the TVar
34   forkIO (stmThread v)
35   -- do 42 times...
36   forM_ [1..42] (\i -> do
37     -- use delay to 'make sure' that a retry is happening for every increment
38     threadDelay 10000
39     -- write new value to TVar using atomically, will cause the STM
40     -- thread to wake up and retry
41     atomically (writeTVar v i))
42
43 If we run this program, we will see 'increment!' printed 43 times, followed
44 by 'final StateT state = 1, STM computation result = 42'. This clearly
45 demonstrates the retry semantics where stmAction is retried 42 times and thus
46 prints 'increment!' 43 times to the console. The StateT computation however is
47 always rolled back when a retry is happening. The rollback is easily possible in pure
48 functional programming due to persistent data structures by simply throwing away

```

¹the new value and retrying with the old value. This example also demonstrates¹
²that any IO actions which happen within an STM action are persistent and can²
³obviously not be rolled back. `Debug.trace` is an IO action masked as pure using³
⁴`unsafePerformIO`.⁴

⁶**3 STM in ABS**⁶

⁷In this section we give a short overview of how to apply STM in ABS. We funda-⁷
⁸mentally follow a time-driven approach in both case studies where the simulation⁸
⁹is advanced by some given Δt and in each step all agents are executed. To employ⁹
¹⁰parallelism, each agent runs within its own thread and agents are executed in lock-¹⁰
¹¹step, synchronising between each Δt which is controlled by the main thread. This¹¹
¹²way of stepping the simulation is introduced in [27] on a conceptual level, where¹²
¹³the authors name it *concurrent update-strategy*. See Figure 1 for a visualisation of¹³
¹⁴our concurrent, time-driven lock-step approach.¹⁴

¹⁵ An agent thread will block until the main thread sends the next Δt and runs the¹⁵
¹⁶STM action atomically with the given Δt . When the STM action has been committed,¹⁶
¹⁷the thread will send the output of the agent action to the main thread to signal it¹⁷
¹⁸has finished. The main thread awaits the results of all agents to collect them for¹⁸
¹⁹output of the current step e.g. visualisation or writing to a file.¹⁹

²⁰ As will be described in subsequent sections, central to both case studies is an²⁰
²¹environment which is shared between the agents using a `TVar` or `TArray` primitive²¹
²²through which the agents communicate concurrently with each other. To get the²²
²³environment in each step for visualisation purposes, the main thread can access the²³
²⁴`TVar` and `TArray` as well.²⁴

²⁶**3.1 Adding STM to agents**²⁶

²⁷ A detailed discussion of how to add STM to agents on a technical level is beyond²⁷
²⁸the focus of this paper as it would require to give an in-depth technical explanation²⁸
²⁹of how our agents are actually implemented [19].²⁹

³⁰ However, the concepts are similar to the example in Section 2.2.4. The agent³⁰
³¹behaviour is an STM action and has access to the environment either through a `TVar`³¹
³²or `TArray` and performs read and write operations directly on it. Each agent itself³²
³³is run within its own thread, and synchronises with the main thread. Thus, it takes³³

¹Haskells `MVar` synchronisation primitives to synchronise with the main thread and¹

²simply runs the STM agent behaviour each time it receives the next tick `DTime`:²

```

3agentThread :: RandomGen g                                3
4    => Int -- Number of steps to compute                    4
5    -> SIRAgent g -- Agent behaviour                        5
6    -> g -- Random-number generator of the agent            6
7    -> MVar SIRState -- Synchronisation back to main thread 7
8    -> MVar DTime -- Receiving DTime for next tick          8
9    -> IO ()                                                9
10 agentThread 0 _ _ _ = return () -- all steps computed, terminate thread 10
11 agentThread n agent rng retVar dtVar = do                 11
12   -- wait for dt to compute current step                   12
13   dt <- takeMVar dtVar                                     13
14   -- compute output of current step                         14
15   let agentSTMAction = runAgent agent                      15
16   -- run the agents STM action atomically within IO        16
17   ((ret, agent'), rng') <- atomically agentSTMAction      17
18   -- post result to main thread                             18
19   putMVar retVar ret                                       19
20   -- tail recursion to next step                           20
21   agentThread (n - 1) agent' rng retVar dtVar             21
22
23   Computing a simulation step is quite trivial within the main thread. All agent
24   threads MVars are signalled to unblock, followed by an immediate block on the
25   MVars into which the agent threads post back their result. The state of the current
26   step is then extracted from the environment, which is stored within the TVar which
27   the agent threads have updated:
28
29simulationStep :: TVar SREnv -- environment                20
30    -> [MVar DTime] -- sync dt to threads                 21
31    -> [MVar SIRState] -- sync output from threads         22
32    -> DTime -- time delta                                23
33    -> IO SREnv                                           24
34simulationStep env dtVars retVars dt = do                 25
35   -- tell all threads to compute next tick with the corresponding DTime 26
36   mapM_ ('putMVar' dt) dtVars                             27
37   -- wait for results but ignore them, SREnv contains all states 28
38   mapM_ takeMVar retVars                                  29
39   -- return state of environment when step has finished
40   readTVarIO env

```

²⁹4 Case Study 1: Spatial SIR Model²⁹

³⁰Our first case study is the SIR model which is a very well studied and understood³⁰
³¹compartment model from epidemiology [28], which allows to simulate the dynamics³¹
³²of an infectious disease like influenza, tuberculosis, chicken pox, rubella and measles³²
³³spreading through a population [29].³³

¹ In it, people in a population of size N can be in either one of three states *Suscep-*¹
²*tible*, *Infected* or *Recovered* at a particular time, where it is assumed that initially²
³there is at least one infected person in the population. People interact *on average*³
⁴with a given rate of β other people per time unit and become infected with a given⁴
⁵probability γ when interacting with an infected person. When infected, a person⁵
⁶recovers *on average* after δ time units and is then immune to further infections.⁶
⁷An interaction between infected persons does not lead to re-infection, thus these⁷
⁸interactions are ignored in this model.⁸

⁹ We followed in our agent-based implementation of the SIR model the work [18] but⁹
¹⁰extended it by placing the agents on a discrete 2D grid using a Moore (8 surrounding¹⁰
¹¹cells) neighbourhood [19]. A visualisation can be seen in Figure 2.¹¹
¹²Due to the continuous-time nature of the SIR model, our implementation follows¹²
¹³the time driven [30] approach. This requires us to sample the system with very small¹³
¹⁴ Δt , which means that we have comparatively few writes to the shared environment¹⁴
¹⁵which will become important when discussing the performance results.¹⁵

¹⁷4.1 Experiment Design¹⁷

¹⁸In this case study we compare the performance of three implementations under¹⁸
¹⁹varying numbers of CPU cores and agent numbers. The code of all implementations¹⁹
²⁰can be accessed freely from the [repository](#) [31].²⁰

- ²¹ 1 Sequential - This is the original implementation as discussed in [19] where the²¹
²²discrete 2D grid is shared amongst all agents as read-only data and the agents²²
²³are executed sequentially within the main thread without any concurrency.²³
- ²⁴ 2 STM - This is the same implementation as the *Sequential* but agents run²⁴
²⁵now in the **STM** context and have access to the discrete 2D grid through a²⁵
²⁶transactional variable **TVar**. This means that the agents now communicate²⁶
²⁷indirectly by reads and writes through the **TVar**.²⁷
- ²⁸ 3 Lock-Based - This is exactly the same implementation as the *STM* one, but²⁸
²⁹instead of running in the **STM** context, the agents now run in the **IO** context.²⁹
³⁰They share the discrete 2D grid using a reference and have access to a lock³⁰
³¹to synchronise access to the reference.³¹

³²Each experiment was run until $t = 100$ and stepped using $\Delta t = 0.1$. For each³²
³³experiment we conducted 8 runs on our machine (see Table 1) under no additional³³

OS	Fedora 28 64-bit
RAM	16 GByte
CPU	Intel Quad Core i5-4670K @ 3.4GHz (4th Gen.)
HD	250Gbyte SSD
Haskell	GHC 8.2.2
Java	OpenJDK 1.8.0

Table 1: Machine and Software Specs for all experiments

	Cores	Duration
Sequential	1	72.5
Lock-Based	1	60.6
	2	42.8
	3	38.6
	4	41.6
STM	1	53.2
	2	27.8
	3	21.8
	4	20.8

Table 2: SIR performance comparisons on 51x51 (2,601 agents) grid with varying number of cores. Timings in seconds (lower is better).

workload and report the mean. Further, we checked the visual outputs and the dynamics and they look qualitatively the same as the reference *Sequential* implementation [19]. A rigorous comparison of all three implementations is beyond the focus of this paper but we refer to the use of property-based testing, as shown in [32]. In the experiments we varied the number of agents (grid size) as well as the number of cores when running concurrently - the numbers are always indicated clearly.

4.2 Constant Grid Size, Varying Cores

In this experiment we held the grid size constant to 51 x 51 (2,601 agents) and varied the cores where possible. The results are reported in Table 2.

Comparing the performance and scaling to multiple cores of the *STM* and *Lock-Based* implementations shows that the *STM* implementation significantly outperforms the *Lock-Based* one and scales better to multiple cores. The *Lock-Based* implementation performs best with 3 cores and shows slightly worse performance on 4 cores as can be seen in Figure 3. This is no surprise because the more cores are running at the same time, the more contention for the lock, thus the more likely

Grid Size	STM	Lock-Based (4 cores)	Lock-Based (3 cores)
51 × 51 (2,601)	20.2	41.9	38.6
101 × 101 (1,0201)	74.5	170.5	171.6
151 × 151 (22,801)	168.5	376.9	404.1
201 × 201 (40,401)	302.4	672.0	720.6
251 × 251 (63,001)	495.7	1,027.3	1,117.2

Table 3: Performance comparison of *STM* and *Lock-Based* SIR implementations on varying grid sizes. Timings in seconds (lower is better).

Grid Size	Commits	Retries	Ratio
51 × 51 (2601)	2601000	1306.5	0.0
101 × 101 (10201)	10201000	3712.5	0.0
151 × 151 (22801)	22801000	8189.5	0.0
201 × 201 (40401)	40401000	13285.0	0.0
251 × 251 (63001)	63001000	21217.0	0.0

Table 4: Retry ratios of the SIR *STM* implementation on varying grid sizes on 4 cores.

synchronisation happening, resulting in higher potential for reduced performance. This is not an issue in *STM* because no locks are taken in advance.

4.3 Varying Grid Size, Constat Cores

In this experiment we varied the grid size and used constantly 4 cores. Because in the previous experiment *Lock-Based* performed best on 3 cores, we additionally ran *Lock-Based* on 3 cores as well. The results are reported in Table 3 and plotted in Figure 4.

It is clear that the *STM* implementation outperforms the *Lock-Based* implementation by a substantial factor. Surprisingly, the *Lock-Based* implementation on 4 cores scales just slightly better with increasing agents number than on 3 cores, something we wouldn't have anticipated based on the results seen in Table 2.

4.4 Retries

Of very much interest when using *STM* is the retry ratio, which obviously depends highly on the read-write patterns of the respective model. We used the *stm-stats* [33] library to record statistics of commits, retries and the ratio. The results are reported in Table 4.

	Cores	51x51	251x251
Lock-Based	16	72.5	1830.5
	32	73.1	1882.2
STM	16	8.6	237.0
	32	12.0	248.7

Table 5: SIR *STM* performance on 16 and 32 cores on Amazon EC2. Timings in seconds (lower is better).

Independent of the number of agents we always have a retry ratio of 0.0. This indicates that this model is *very* well suited to STM, which is also directly reflected in the much better performance over the *Lock-Based* implementation. Obviously this ratio stems from the fact, that in our implementation we have *very* few writes, which happen only in case when an agent changes from *Susceptible* to *Infected* or from *Infected* to *Recovered*.

4.5 Going Large-Scale

To test how far we can scale up the number of cores in both the *Lock-Based* and *STM* cases, we ran two experiments, 51x51 and 251x251, on Amazon EC instances with a larger number of cores than our local machinery, starting with 16 and 32 to see if we are running into decreasing returns. The results are reported in Table 5.

As expected, the *Lock-Based* approach doesn't scale up to many cores because each additional core brings more contention to the lock, resulting in an even more decreased performance. This is particularly obvious in the 251x251 experiment because of the much larger number of concurrent agents. The *STM* approach returns better performance on 16 cores but fails to scale further up to 32 where the performance drops below the one with 16 cores. In both STM cases we measured a retry ratio of 0, thus we conclude that with 32 cores we become limited by the overhead of STM transactions [23] because the workload of an STM action in our SIR implementation is quite small.

4.6 Discussion

The timing measurements speak a clear language. Running in STM and sharing state using a transactional variable **TVar** is much more time efficient than both the *Sequential* and *Lock-Based* approach. On 4 cores *STM* achieves a speedup factor of 3.6, nearly reaching the theoretical limit. Obviously both *STM* and *Lock-*

¹*Based* sacrifice determinism, which means that repeated runs might not lead to¹
²same dynamics despite same initial conditions. Still, by sticking to STM, we get²
³the guarantee that the source of this nondeterminism is concurrency within the³
⁴STM context but *nothing else*. This we can not guarantee in the case of the *Lock*-⁴
⁵*Based* approach as all bets are off when running within the IO context. The fact to⁵
⁶have *both* the better performance *and* the stronger static guarantees in the *STM*⁶
⁷approach makes it *very* compelling. 7

8

8

9

9

¹⁰5 Case Study 2: SugarScape 10

¹¹One of the first models in Agent-Based Simulation was the seminal Sugarscape¹¹
¹²model developed by Epstein and Axtell in 1996 [20]. Their aim was to *grow* an¹²
¹³artificial society by simulation and connect observations in their simulation to phe-¹³
¹⁴nomenon observed in real-world societies. In this model a population of agents move¹⁴
¹⁵around in a discrete 2D environment, where sugar grows, and interact with each¹⁵
¹⁶other and the environment in many different ways. The main features of this model¹⁶
¹⁷are (amongst others): searching, harvesting and consuming of resources, wealth and¹⁷
¹⁸age distributions, population dynamics under sexual reproduction, cultural pro-¹⁸
¹⁹cesses and transmission, combat and assimilation, bilateral decentralized trading¹⁹
²⁰(bartering) between agents with endogenous demand and supply, disease processes²⁰
²¹transmission and immunology. 21

²² We implemented the *Carrying Capacity* (p. 30) section of Chapter II of the book²²
²³[20]. There, in each step agents search (move) to the cell with the most sugar they²³
²⁴see within their vision, harvest all of it from the environment and consume sugar²⁴
²⁵based on their metabolism. Sugar regrows in the environment over time. Only one²⁵
²⁶agent can occupy a cell at a time. Agents don't age and cannot die from age. If²⁶
²⁷agents run out of sugar due to their metabolism, they die from starvation and are²⁷
²⁸removed from the simulation. The authors report that the initial number of agents²⁸
²⁹quickly drops and stabilises around a level depending on the model parameters.²⁹
³⁰This is in accordance with our results as we show in Figure 5 and guarantees that³⁰
³¹we don't run out of agents. The model parameters are as follows: 31

- ³² • Sugar Endowment: each agent has an initial sugar endowment randomly uni- 32
- ³³form distributed between 5 and 25 units; 33

- 1 • Sugar Metabolism: each agent has a sugar metabolism randomly uniform dis-¹
tributed between 1 and 5; ²
- 3 • Agent Vision: each agent has a vision randomly uniform distributed between³
1 and 6, same for each of the 4 directions (N, W, S, E); ⁴
- 5 • Sugar Growback: sugar grows back by 1 unit per step until the maximum⁵
capacity of a cell is reached; ⁶
- 7 • Agent Number: initially 500 agents; ⁷
- 8 • Environment Size: 50 x 50 cells with toroid boundaries which wrap around in⁸
both x and y dimension. ⁹

11 5.1 Experiment Design ¹¹

12 We compare four different implementations with the code freely accessible from the ¹²
[repository](#) [34]: ¹³

- 14 1 Sequential - All agents are run after another (including the environment) and¹⁴
the environment is shared amongst the agents using a read and write state¹⁵
context. ¹⁶
- 17 2 Lock-Based - All agents are run concurrently and the environment is shared¹⁷
using a global reference amongst the agents which acquire and release a lock¹⁸
when accessing it. ¹⁹
- 20 3 STM TVar - All agents are run concurrently and the environment is shared²⁰
using a **TVar** amongst the agents. ²¹
- 22 4 STM TArray - All agents are run concurrently and the environment is shared²²
using a **TArray** amongst the agents. ²³

24 *Ordering* The model specification requires to shuffle agents before every step ([20],²⁴
footnote 12 on page 26). In the *Sequential* approach we do this explicitly but in the²⁵
Lock-Based and both *STM* approaches we assume this to happen automatically due²⁶
to race conditions in concurrency, thus we arrive at an effectively shuffled processing²⁷
of agents because we implicitly assume that the order of the agents is *effectively*²⁸
random in every step. The important difference between the two approaches is²⁹
that in the *Sequential* approach we have full control over this randomness but in³⁰
the *STM* and *Lock-Based* not. This has the consequence that repeated runs with³¹
the same initial conditions might lead to slightly different results. This decision³²
leaves the execution order of the agents ultimately to Haskell's Runtime System³³

¹and the underlying operating system. We are aware that by doing this, we make¹
²assumptions that the threads run uniformly distributed (fair) but such assumptions²
³should not be made in concurrent programming. As a result we can expect this fact³
⁴to produces non-uniform distributions of agent runs but we assumed that for this⁴
⁵model this does not has a significance influence. In case of doubt, we could resort to⁵
⁶shuffling the agents before running them in every step. This problem, where also the⁶
⁷influence of nondeterministic ordering on the correctness and results of ABS has to⁷
⁸be analysed, deserves in-depth research on its own. Again we refer to the technique⁸
⁹of property-based testing as shown in [32] as this issue is beyond the focus of this⁹
¹⁰paper. 10

¹¹Note that in the concurrent implementations we have two options for running¹¹
¹²the environment: either asynchronously as a concurrent agent at the same time¹²
¹³with the population agents or synchronously after all agents have run. We must¹³
¹⁴be careful though as running the environment as a concurrent agent can be seen¹⁴
¹⁵as conceptually wrong because the time when the regrowth of the sugar happens¹⁵
¹⁶is now completely random. In this case it could happen that sugar regrows in the¹⁶
¹⁷very first transaction or in the very last, different in each step, which can be seen¹⁷
¹⁸as a violation of the model specifications. Thus we do not run the environment¹⁸
¹⁹concurrently with the agents but synchronously after all agents have run. 19

²⁰We follow [35] and measure the average number of steps per second of the simula-²⁰
²¹tion over 60 seconds. For each experiment we conducted 8 runs on our machine (see²¹
²²Table 1) under no additional workload and report the average. In the experiments²²
²³we varied the number of cores when running concurrently - the numbers are always²³
²⁴indicated clearly. 24

²⁵25 ²⁶5.2 Constant Agent Size 26

²⁷In a first approach we compare the performance of all implementations on varying²⁷
²⁸numbers of cores. The results are reported in Table 6 and plotted in Figure 6. 28

²⁹As expected, the *Sequential* implementation is the slowest, followed by the *Lock*-²⁹
³⁰*Based* and *TVar* approach whereas *TArray* is the best performing one. 30

³¹We clearly see that using a **TVar** to share the environment is a very inefficient³¹
³²choice: *every* write to a cell leads to a retry independent whether the reading agent³²
³³reads that changed cell or not, because the data structure can not distinguish be-³³

	Cores	Steps	Retries
Sequential	1	39.4	N/A
Lock-Based	1	43.0	N/A
	2	51.8	N/A
	3	57.4	N/A
	4	58.1	N/A
STM <i>TVar</i>	1	47.3	0.0
	2	53.5	1.1
	3	57.1	2.2
	4	53.0	3.2
STM <i>TArray</i>	1	45.4	0.0
	2	65.3	0.02
	3	75.7	0.04
	4	84.4	0.05

Table 6: Comparison of steps per seconds (higher is better) and retries (lower is better) of various Sugarscape implementations with constant agent numbers. Using 50x50 grid with 500 initial agents on varying cores.

between individual cells. By using a *TArray* we can avoid the situation where a write to a cell in a far distant location of the environment will lead to a retry of an agent which never even touched that cell. Also the *TArray* seems to scale up by 10 steps per second for every core added. It will be interesting to see how far this will go with the Amazon experiment, as we seem not to hit a limit with 4 cores yet.

The inefficiency of *TVar* is also reflected in the nearly similar performance of the *Lock-Based* implementation which even outperforms it on 4 cores. This is due to very similar approaches because both operate on the whole environment instead of only the cells as *TArray* does. This seems to be a bottleneck in *TVar* reaching the best performance on 3 cores, which then drops on 4 cores. The *Lock-Based* approach seems to reduce its returns on increased number of cores hitting a limit at 4 cores as well.

5.3 Scaling up Agents

So far we kept the initial number of agents at 500, which due to the model specification, quickly drops and stabilises around 200 due to the carrying capacity of the environment as described in the book [20] section *Carrying Capacity* (p. 30).

We now measure the performance of our approaches under increased number of agents. For this we slightly change the implementation: always when an agent dies it

Agents	Sequential	Lock-Based	TVar (3 cores)	TVar (4 cores)	TArray
500	14.4	20.2	20.1	18.5	71.9
1,000	6.8	10.8	10.4	9.5	54.8
1,500	4.7	8.1	7.9	7.3	44.1
2,000	4.4	7.6	7.4	6.7	37.0
2,500	5.3	5.4	9.2	8.9	33.3

Table 7: Steps per second (higher is better) of various Sugarscape implementations with varying agent numbers. Using 50x50 grid with varying number of agents with 4 (and 3) cores except the *Sequential* (1 core) implementation.

spawns a new one which is inspired by the ageing and birthing feature of Chapter III in the book [20]. This ensures that we keep the number of agents roughly constant (still fluctuates but doesn't drop to low levels) over the whole duration. This ensures a constant load of concurrent agents interacting with each other and demonstrates also the ability to terminate and fork threads dynamically during the simulation. Except for the *Sequential* approach we ran all experiments with 4 cores (*TVar* with 3 as well). We looked into the performance of 500, 1,000, 1,500, 2,000 and 2,500 (maximum possible capacity of the 50x50 environment). The results are reported in Table 7 and plotted in Figure 7.

As expected, the *TArray* implementation outperforms all others substantially. Also as expected, the *TVar* implementation on 3 cores is faster than on 4 cores as well when scaling up to more agents. The *Lock-Based* approach performs about the same as the *TVar* on 3 cores because of the very similar approaches: both access the whole environment. Still the *TVar* approach uses one core less to arrive at the same performance, thus strictly speaking outperforming the *Lock-Based* implementation.

What seems to be very surprising is that in the *Sequential* and *TVar* cases the performance with 2,500 agents is *better* than the one with 2,000 agents. The reason for this is that in the case of 2,500 agents, an agent can't move anywhere because all cells are already occupied. In this case the agent won't rank the cells in order of their payoff (max sugar) to move to but just stays where it is. We hypothesize that due to Haskell's laziness the agents actually never look at the content of the cells in this case but only the number which means that the cells themselves are never evaluated which further increases performance. This leads to the better performance in case of *Sequential* and *TVar* because both exploit laziness. In the case of the *Lock-Based*

	Cores	Carrying Capacity	Rebirthing
Lock-Based	16	53.9	4.4
	32	44.2	3.6
STM TArray	16	116.8 (0.23)	39.5 (0.08)
	32	109.8 (0.41)	31.3 (0.18)

Table 8: Sugarscape *STM* performance on 16 and 32 cores on Amazon EC2. Values are steps per second (higher is better). Retry ratio in brackets.

approach we still arrive at a lower performance because the limiting factor are the unconditional locks. In the case of the *TArray* approach we also arrive at a lower performance because it seems that STM perform reads on the neighbouring cells which are not subject to lazy evaluation.

We also measured the average retries both for *TVar* and *TArray* under 2,500 agents where the *TArray* approach shows best scaling performance with 0.01 retries whereas *TVar* averages at 3.28 retries. Again this can be attributed to the better transactional data structure which reduces retry ratio substantially to near-zero levels.

5.4 Going Large-Scale

To test how far we can scale up the number of cores in both the *Lock-Based* and *TArray* cases, we ran the two experiments (carrying capacity and rebirthing) on Amazon EC instances with increasing number of cores starting with 16 and 32 to see if we run into decreasing returns. The results are reported in Table 8.

As expected, the *Lock-Based* approach doesn't scale up to many cores because each additional core brings more contention to the lock, resulting in even more decreased performance. This is particularly obvious in the rebirthing experiment because of the much larger number of concurrent agents. The *TArray* approach returns better performance on 16 cores but fails to scale further up to 32 where the performance drops below the one with 16 cores. We indicated the retry ratio in brackets and see that they roughly double from 16 to 32, which is the reason why performance drops as at this point.

5.5 Comparison with other approaches

The paper [35] reports a performance of 17 steps in RePast, 18 steps in MASON (both non-parallel) and 2,000 steps per second on a GPU on a 128x128 grid. Al-

¹though our *Sequential* implementation, which runs non-parallel as well, outperforms¹
²the RePast and MASON implementations of [35], one must be very well aware that²
³these results were generated in 2008, on current hardware of that time.³

⁴ The very high performance on the GPU does not concern us here as it follows⁴
⁵a very different approach than ours. We focus on speeding up implementations on⁵
⁶the CPU as directly as possible without locking overhead. When following a GPU⁶
⁷approach one needs to map the model to the GPU which is a delicate and non-⁷
⁸trivial matter. With our approach we show that speed up with concurrency is very⁸
⁹possible without the low-level locking details or the need to map to GPU. Also some⁹
¹⁰features like bilateral trading between agents, where a pair of agents needs to come¹⁰
¹¹to a conclusion over multiple synchronous steps, is difficult to implement on a GPU¹¹
¹²whereas this should be not as hard using STM.¹²

¹³ Note that we kept the grid size constant because we implemented the environ-¹³
¹⁴ment as a single agent which works sequentially on the cells to regrow the sugar.¹⁴
¹⁵ Obviously this doesn't really scale up on parallel hardware and experiments which¹⁵
¹⁶we haven't included here due to lack of space, show that the performance goes down¹⁶
¹⁷dramatically when we increase the environment to 128x128 with same number of¹⁷
¹⁸agents. This is the result of Amdahl's law where the environment becomes the limit-¹⁸
¹⁹ing *sequential* factor of the simulation. Depending on the underlying data structure¹⁹
²⁰used for the environment we have two options to solve this problem. In the case of²⁰
²¹the *Sequential* and *TVar* implementation we build on an indexed array, which can²¹
²²be updated in parallel using the existing data-parallel support in Haskell. In the²²
²³case of the *TArray* approach we have no option but to run the update of every cell²³
²⁴within its own thread. We leave both for further research as it is beyond the scope²⁴
²⁵of this paper.²⁵

²⁶5.6 Discussion²⁶

²⁷ This case study showed clearly that besides being substantially faster than the²⁷
²⁸*Sequential* implementation, *STM* implementations are also able to perform consid-²⁸
²⁹erably better than a *Lock-Based* approach even in the case of the Sugarscape model²⁹
³⁰which has a much higher complexity in agent behaviour and dramatically increased³⁰
³¹number of writes to the environment.³¹

Further, this case study demonstrated that the selection of the right transactional¹ data structure is of fundamental importance when using STM. Selecting the right² transactional data structure is highly model-specific and can lead to dramatically³ different performance results. In this case study the *TArray* performed best due to⁴ many writes but in the SIR case study a *TVar* showed good enough results due to the⁵ very low number of writes. When not carefully selecting the right transactional data⁶ structure which supports fine-grained concurrency, a lock-based implementation⁷ might perform as well or even outperform the STM approach as can be seen when⁸ using the *TVar*.⁹

Although the *TArray* is the better transactional data structure overall, it might¹⁰ come with an overhead, performing worse on low number of cores than a *TVar*¹¹ approach but has the benefit of quickly scaling up to multiple cores. Depending on¹² the transactional data structure, scaling up to multiple cores hits a limit at some¹³ point. In the case of the *TVar* the best performance is reached with 3 cores. With¹⁴ the *TArray* we reached this limit around 16 cores.¹⁵

The comparison between the *Lock-Based* approach and the *TArray* implementa-¹⁶ tion is a bit unfair due to a very different locking structure. A more suitable compar-¹⁷ ison would have been to use an indexed Array with a tuple of (*MVar*, *IORef*), hold-¹⁸ ing a synchronisation primitive and reference for each cell to support fine-grained¹⁹ locking on the cell level. This would be a more just comparison to the *TArray* where²⁰ fine-grained transactions happen on the cell level. We hypothesise that STM will²¹ still outperform the lock-based approach but to a lesser degree. We leave the proof²² of this for further research.²³

²⁵6 Conclusion²⁵

In this paper we investigated the potential for using STM for parallel, large scale²⁶ ABS and come to the conclusion that it is indeed a very promising alternative over²⁷ lock-based approaches as our case studies have shown. The STM implementations²⁸ all consistently outperformed the lock-based ones and scaled much better to larger²⁹ number of CPUs. Besides, the concurrency abstractions of STM are very powerful,³⁰ yet simple enough to allow convenient implementation of concurrent agents without³¹ the problems of lock-based implementation. Due to most ABS being primarily pure³² computations, which do not need interactive input from the user, files or network³³

¹during simulation, the fact that no such interactions can occur within an agent¹
²when running within STM is not a problem.²

³ Further, STM primitives map nicely to ABS concepts. When having a shared envi-³
⁴ronment, it is natural either using `TVar` or `TArray`, depending on the environments⁴
⁵nature. Also, there exists the `TChan` primitive, which can be seen as a persistent⁵
⁶message box for agents, underlining the message-oriented approach found in many⁶
⁷agent-based models [36, 37]. Also `TChan` offers a broadcast transactional channel,⁷
⁸which supports broadcasting to listeners which maps nicely to a proactive environ-⁸
⁹ment or a central auctioneer upon which agents need to synchronize. The benefits⁹
¹⁰of these natural mappings are that using STM takes a big portion of burden from¹⁰
¹¹the modeller as one can think in STM primitives instead of low level locks and¹¹
¹²concurrent operational details.¹²

¹³ The strong static type system of Haskell adds another benefit. By running in the¹³
¹⁴STM instead of IO context makes the concurrent nature more explicit and at the¹⁴
¹⁵same time restricts it to purely STM behaviour. So despite obviously losing the¹⁵
¹⁶reproducibility property due to concurrency, we still can guarantee that the agents¹⁶
¹⁷can't do arbitrary IO as they are restricted to STM operations only.¹⁷

¹⁸ Depending on the nature of the transactions, retries could become a bottle neck,¹⁸
¹⁹resulting in a live lock in extreme cases. The central problem of STM is to keep the¹⁹
²⁰retries low, which is directly influenced by the read/writes on the STM primitives.²⁰
²¹By choosing more fine-grained and suitable data structures e.g. using a `TArray`²¹
²²instead of an indexed array within a `TVar`, one can reduce retries and increase²²
²³performance significantly and avoid the problem of live locks as we have shown.²³

²⁴ Despite the indisputable benefits of using STM within a pure functional setting²⁴
²⁵like Haskell, it exists also in other imperative languages (Python, Java and C++,²⁵
²⁶etc) and we hope that our research sparks interest in the use of STM in ABS in²⁶
²⁷general and that other researchers pick up the idea and apply it to the established²⁷
²⁸imperative languages Python, Java, C++ in the ABS community as well.²⁸

³⁰7 Further Research³⁰

³¹ So far we only implemented a tiny bit of the Sugarscape model and left out the later³¹
³²chapters which are more involved as they incorporate direct synchronous commu-³²
³³nication between agents. Such mechanisms are very difficult to approach in GPU³³

¹based approaches [35] but should be quite straightforward in STM using TChan and¹
²retries. However, we have yet to prove how to implement reliable synchronous agent²
³interactions without deadlocks in STM. It might be very well the case that a truly³
⁴concurrent approach is doomed due to the following [38] (Chapter 10. Software⁴
⁵Transactional Memory, *What Can We Not Do with STM?*): “In general, the class⁵
⁶of operations that STM cannot express are those that involve multi-way communi-⁶
⁷cation between threads. The simplest example is a synchronous channel, in which⁷
⁸both the reader and the writer must be present simultaneously for the operation to⁸
⁹go ahead. We cannot implement this in STM, at least compositionally [...]: the op-⁹
¹⁰erations need to block and have a visible effect — advertise that there is a blocked¹⁰
¹¹thread — simultaneously.”. 11

¹²A drawback of STM is that it is not fair because *all* threads, which block on a¹²
¹³transactional primitive, have to be woken up upon a change of the primitive, thus a¹³
¹⁴FIFO guarantee cannot be given. We hypothesise that for most models, where the¹⁴
¹⁵STM approach is applicable, this has no qualitative influence on the dynamics as¹⁵
¹⁶agents are assumed to act conceptually at the same time and no fairness is needed.¹⁶
¹⁷We leave the test of this hypothesis for future research. This is connected to our¹⁷
¹⁸assumption that concurrent execution has no qualitative influence on the dynamics.¹⁸
¹⁹Although repeated runs with same initial conditions might lead to different results¹⁹
²⁰due to nondeterminism, the dynamics follow still the same distribution as the one²⁰
²¹from the sequential implementation. To verify this we can make use the techniques²¹
²²of property-based testing as shown in [32] but we leave it for further research. 22

²³23

²⁴24

²⁴Declarations 24

²⁵Availability of data and materials 25

²⁶The datasets used and/or analysed during the current study are available from the corresponding author on 26
²⁷reasonable request. 27

²⁷27

²⁸Competing interests 28

²⁸The authors declare that they have no competing interests. 28

²⁹29

³⁰Funding 30

³⁰Not applicable. 30

³¹31

³¹Authors' contributions 31

³²JT initiated the idea and the research, did the implementation, experiments, performance measurements, and 32

³³writing. POS supervised the work, gave feedback and supported the writing process. All authors read and approved 33
³³the final manuscript.

¹Acknowledgements 1

²The authors would like to thank J. Hey and M. Handley for constructive feedback, comments and valuable 2
³discussions. 3

⁴Authors' information 4

⁵**JONATHAN THALER** is a Ph.D. student at the University of Nottingham and part of the Intelligent Modelling and 5
⁶Analysis Group (<http://www.cs.nott.ac.uk/~psxjat/>). His main research interest is the benefits and drawbacks 6
⁶of using pure functional programming with Haskell for implementing Agent-Based Simulations. 6

⁷**DR. PEER-OLAF SIEBERS** is an Assistant Professor at the School of Computer Science, University of Nottingham, 7
⁸UK (<http://www.cs.nott.ac.uk/~pszps/>). His main research interest is the application of computer simulation to 8
⁹study human-centric complex adaptive systems. He is a strong advocate of Object Oriented Agent-Based Social 8
⁹Simulation. This is a novel and highly interdisciplinary research field, involving disciplines like Social Science, 9
¹⁰Economics, Psychology, Operations Research, Geography, and Computer Science. His current research focuses on 10
¹¹Urban Sustainability and he is a co-investigator in several related projects and a member of the university's 11
¹¹"Sustainable and Resilient Cities" Research Priority Area management team. 11

¹²References 12

- ¹³1. Macal CM. Everything you need to know about agent-based modelling and simulation. *Journal of Simulation*. 13
¹⁴2016 May;10(2):144–156. Available from: <https://link.springer.com/article/10.1057/jos.2016.7>. 14
- ¹⁵2. Suryanarayanan V, Theodoropoulos G, Lees M. PDES-MAS: Distributed Simulation of Multi-agent Systems. *Procedia Computer Science*. 2013 Jan;18:671–681. Available from: 15
¹⁶<http://www.sciencedirect.com/science/article/pii/S1877050913003748>. 16
- ¹⁷3. Logan B, Theodoropoulos G. The distributed simulation of multiagent systems. *Proceedings of the IEEE*. 2001 17
¹⁸Feb;89(2):174–185. 18
- ¹⁹4. Lees M, Logan B, Theodoropoulos G. Using Access Patterns to Analyze the Performance of Optimistic 19
²⁰Synchronization Algorithms in Simulations of MAS. *Simulation*. 2008 Oct;84(10-11):481–492. Available from: 20
²¹<http://dx.doi.org/10.1177/0037549708096691>. 21
- ²²5. Suryanarayanan V, Theodoropoulos G. Synchronised Range Queries in Distributed Simulations of Multiagent 22
²³Systems. *ACM Trans Model Comput Simul*. 2013 Nov;23(4):25:1–25:25. Available from: 23
²⁴<http://doi.acm.org/10.1145/2517449>. 24
- ²⁵6. Riley PF, Riley GF. Next Generation Modeling III - Agents: Spades — a Distributed Agent Simulation 25
²⁶Environment with Software-in-the-loop Execution. In: *Proceedings of the 35th Conference on Winter* 26
²⁷*Simulation: Driving Innovation*. WSC '03. Winter Simulation Conference; 2003. p. 817–825. Event-place: New 27
²⁸Orleans, Louisiana. Available from: <http://dl.acm.org/citation.cfm?id=1030818.1030926>. 28
- ²⁹7. Gasser L, Kakugawa K. MACE3J: Fast Flexible Distributed Simulation of Large, Large-grain Multi-agent 29
³⁰Systems. In: *Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent* 30
³¹*Systems: Part 2. AAMAS '02*. New York, NY, USA: ACM; 2002. p. 745–752. Event-place: Bologna, Italy. 31
³²Available from: <http://doi.acm.org/10.1145/544862.544918>. 32
- ³³8. Himmelspach J, Uhrmacher AM. Plug'n Simulate. In: *40th Annual Simulation Symposium (ANSS'07)*; 2007. p. 33
³⁴137–143. ISSN: 1080-241X. 34
- ³⁵9. Minson R, Theodoropoulos GK. Distributing RePast agent-based simulations with HLA. *Concurrency and* 35
³⁶*Computation: Practice and Experience*. 2008;20(10):1225–1256. Available from: 36
³⁷<https://onlinelibrary.wiley.com/doi/abs/10.1002/cpe.1280>. 37
- ³⁸10. Gorur BK, Imre K, Oguztuzun H, Yilmaz L. Repast HPC with Optimistic Time Management. In: *Proceedings* 38
³⁹*of the 24th High Performance Computing Symposium. HPC '16*. San Diego, CA, USA: Society for Computer 39
⁴⁰Simulation International; 2016. p. 4:1–4:9. Event-place: Pasadena, California. Available from: 40
⁴¹<https://doi.org/10.22360/SpringSim.2016.HPC.046>. 41
- ⁴²11. Hay J, Wilsey PA. Experiments with Hardware-based Transactional Memory in Parallel Simulation. In: 42
⁴³*Proceedings of the 3rd ACM SIGSIM Conference on Principles of Advanced Discrete Simulation. SIGSIM PADS* 43
⁴⁴*'15*. New York, NY, USA: ACM; 2015. p. 75–86. Event-place: London, United Kingdom. Available from: 44
⁴⁵<http://doi.acm.org/10.1145/2769458.2769462>. 45

12. Abar S, Theodoropoulos GK, Lemarinier P, O'Hare GMP. Agent Based Modelling and Simulation tools: A review of the state-of-art software. *Computer Science Review*. 2017 May;24:13–33. Available from: <https://linkinghub.elsevier.com/retrieve/pii/S1574013716301198>.
13. Cicirelli F, Giordano A, Nigro L. Efficient Environment Management for Distributed Simulation of Large-scale Situated Multi-agent Systems. *Concurr Comput : Pract Exper*. 2015 Mar;27(3):610–632. Available from: <http://dx.doi.org/10.1002/cpe.3254>.
14. Hudak P, Hughes J, Peyton Jones S, Wadler P. A History of Haskell: Being Lazy with Class. In: *Proceedings of the Third ACM SIGPLAN Conference on History of Programming Languages. HOPL III*. New York, NY, USA: ACM; 2007. p. 12–112–55. Available from: <http://doi.acm.org/10.1145/1238844.1238856>.
15. Discolo A, Harris T, Marlow S, Jones SP, Singh S. Lock Free Data Structures Using STM in Haskell. In: *Proceedings of the 8th International Conference on Functional and Logic Programming. FLOPS'06*. Berlin, Heidelberg: Springer-Verlag; 2006. p. 65–80. Available from: http://dx.doi.org/10.1007/11737414_6.
16. Bezirgiannis N. Improving Performance of Simulation Software Using Haskell's Concurrency & Parallelism; 2013.
17. Wilensky U, Rand W. *An Introduction to Agent-Based Modeling: Modeling Natural, Social, and Engineered Complex Systems with NETLogo*. MIT Press; 2015. Available from: <https://www.amazon.co.uk/Introduction-Agent-Based-Modeling-Natural-Engineered/dp/0262731894>.
18. Macal CM. To Agent-based Simulation from System Dynamics. In: *Proceedings of the Winter Simulation Conference. WSC '10*. Baltimore, Maryland: Winter Simulation Conference; 2010. p. 371–382. Available from: <http://dl.acm.org/citation.cfm?id=2433508.2433551>.
19. Thaler J, Altenkirch T, Siebers PO. Pure Functional Epidemics: An Agent-Based Approach. In: *Proceedings of the 30th Symposium on Implementation and Application of Functional Languages. IFL 2018*. New York, NY, USA: ACM; 2018. p. 1–12. Event-place: Lowell, MA, USA. Available from: <http://doi.acm.org/10.1145/3310232.3310372>.
20. Epstein JM, Axtell R. *Growing Artificial Societies: Social Science from the Bottom Up*. Washington, DC, USA: The Brookings Institution; 1996.
21. Shavit N, Touitou D. Software Transactional Memory. In: *Proceedings of the Fourteenth Annual ACM Symposium on Principles of Distributed Computing. PODC '95*. New York, NY, USA: ACM; 1995. p. 204–213. Available from: <http://doi.acm.org/10.1145/224964.224987>.
22. Heindl A, Pokam G. Modeling Software Transactional Memory with AnyLogic. In: *Proceedings of the 2Nd International Conference on Simulation Tools and Techniques. Simutools '09*. ICST, Brussels, Belgium, Belgium: ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering); 2009. p. 10:1–10:10. Available from: <http://dx.doi.org/10.4108/ICST.SIMUTOOLS2009.5581>.
23. Perfumo C, Sönmez N, Stipic S, Unsal O, Cristal A, Harris T, et al. The Limits of Software Transactional Memory (STM): Dissecting Haskell STM Applications on a Many-core Environment. In: *Proceedings of the 5th Conference on Computing Frontiers. CF '08*. New York, NY, USA: ACM; 2008. p. 67–78. Available from: <http://doi.acm.org/10.1145/1366230.1366241>.
24. Marlow S, Peyton Jones S, Singh S. Runtime Support for Multicore Haskell. In: *Proceedings of the 14th ACM SIGPLAN International Conference on Functional Programming. ICFP '09*. New York, NY, USA: ACM; 2009. p. 65–78. Available from: <http://doi.acm.org/10.1145/1596550.1596563>.
25. Harris T, Marlow S, Peyton-Jones S, Herlihy M. Composable Memory Transactions. In: *Proceedings of the Tenth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming. PPOPP '05*. New York, NY, USA: ACM; 2005. p. 48–60. Available from: <http://doi.acm.org/10.1145/1065944.1065952>.
26. Harris T, Peyton Jones S. Transactional memory with data invariants. In: *First ACM SIGPLAN Workshop on Languages, Compilers, and Hardware Support for Transactional Computing (TRANSACT'06)*; 2006. Available from: <https://www.microsoft.com/en-us/research/publication/transactional-memory-data-invariants/>.
27. Thaler J, Siebers PO. *The Art Of Iterating: Update-Strategies in Agent-Based Simulation*. Dublin; 2017. .
28. Kermack WO, McKendrick AG. A Contribution to the Mathematical Theory of Epidemics. *Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*. 1927 Aug;115(772):700–721. Available from: <http://rspa.royalsocietypublishing.org/content/115/772/700>.
29. Enns RH. *It's a Nonlinear World*. 1st ed. Springer Publishing Company, Incorporated; 2010.

30. Meyer R. Event-Driven Multi-agent Simulation. In: Multi-Agent-Based Simulation XV. Lecture Notes in Computer Science. Springer, Cham; 2014. p. 3–16. Available from: https://link.springer.com/chapter/10.1007/978-3-319-14627-0_1.
31. Thaler J. Repository of STM implementations of the agent-based SIR model in Haskell; 2019. Last Access October 11, 2019. <https://github.com/thalerjonathan/haskell-stm-sir>. Available from: <https://github.com/thalerjonathan/haskell-stm-sir>.
32. Thaler J, Siebers PO. Show Me Your Properties! The Potential Of Property-Based Testing In Agent-Based Simulation. Berlin; 2019. .
33. Breitner J. stm-stats library; 2019. Last Access October 11, 2019. <http://hackage.haskell.org/package/stm-stats>. Available from: <http://hackage.haskell.org/package/stm-stats>.
34. Thaler J. Repository of STM implementations of the Sugarscape model in Haskell; 2019. Last Access October 11, 2019. <https://github.com/thalerjonathan/haskell-stm-sugarscape>. Available from: <https://github.com/thalerjonathan/haskell-stm-sugarscape>.
35. Lysenko M, D'Souza RM. A Framework for Megascale Agent Based Model Simulations on Graphics Processing Units. Journal of Artificial Societies and Social Simulation. 2008;11(4):10. Available from: <http://jasss.soc.surrey.ac.uk/11/4/10.html>.
36. Agha G. Actors: A Model of Concurrent Computation in Distributed Systems. Cambridge, MA, USA: MIT Press; 1986.
37. Wooldridge M. An Introduction to MultiAgent Systems. 2nd ed. Wiley Publishing; 2009.
38. Marlow S. Parallel and Concurrent Programming in Haskell. O'Reilly; 2013. Google-Books-ID: k0W6AQAAAJ.

Figures

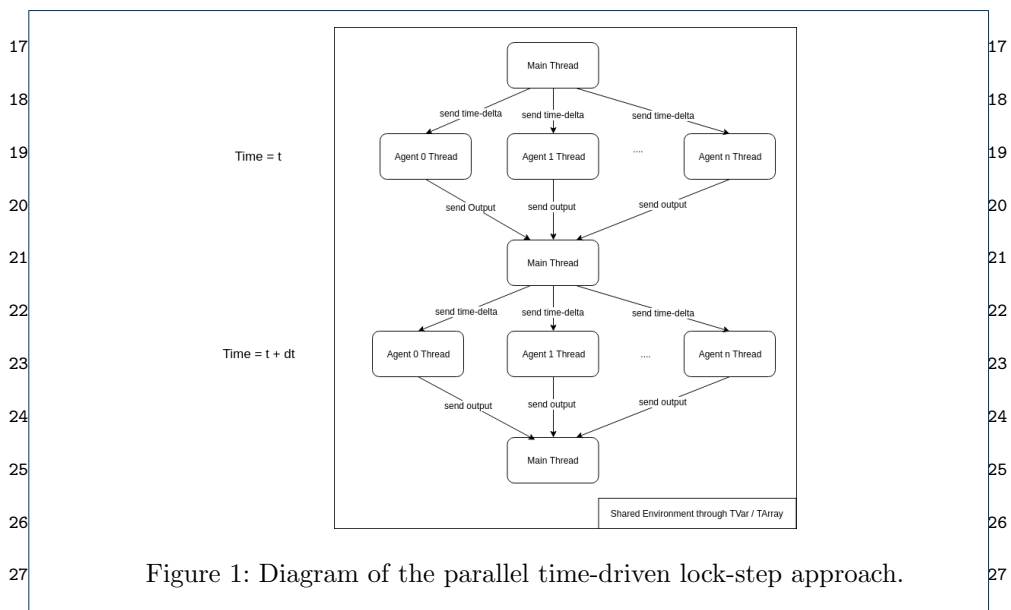


Figure 1: Diagram of the parallel time-driven lock-step approach.

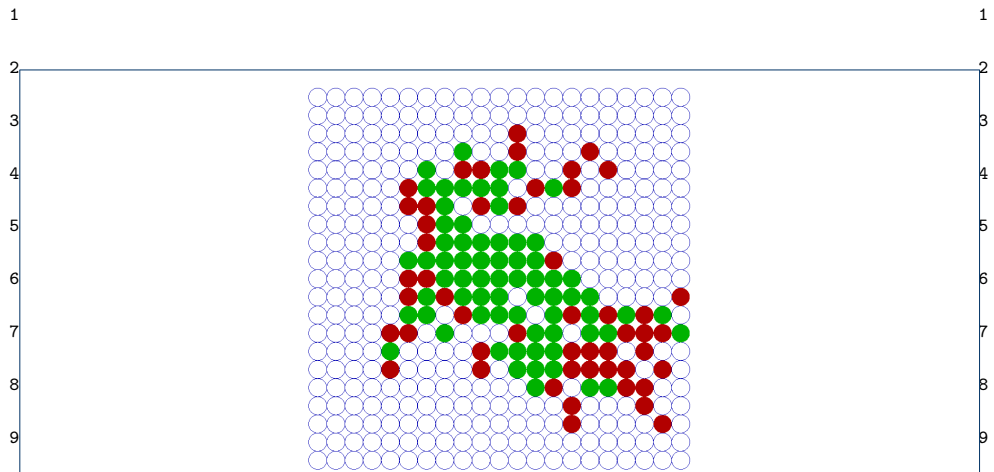


Figure 2: Simulating the spatial SIR model with a Moore neighbourhood, a single infected agent at the center, contact rate $\beta = \frac{1}{5}$, infection probability $\gamma = 0.05$ and illness duration $\delta = 15$. Infected agents are indicated by red circles, recovered agents by green ones. The susceptible agents are rendered as blue hollow circles for better contrast.

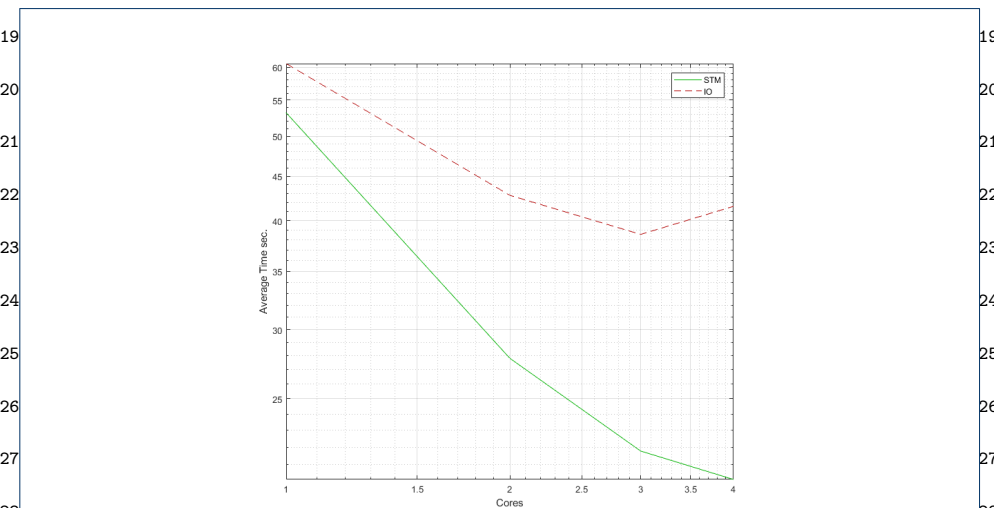


Figure 3: Performance and scaling comparison of *STM* and *Lock-Based* SIR implementations on multiple cores. Note that the Lock-Based implementation performs worse on 4 than on 3 cores due to lock contention.

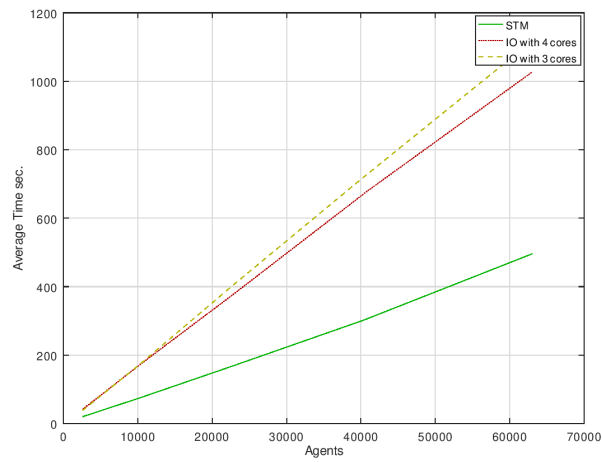
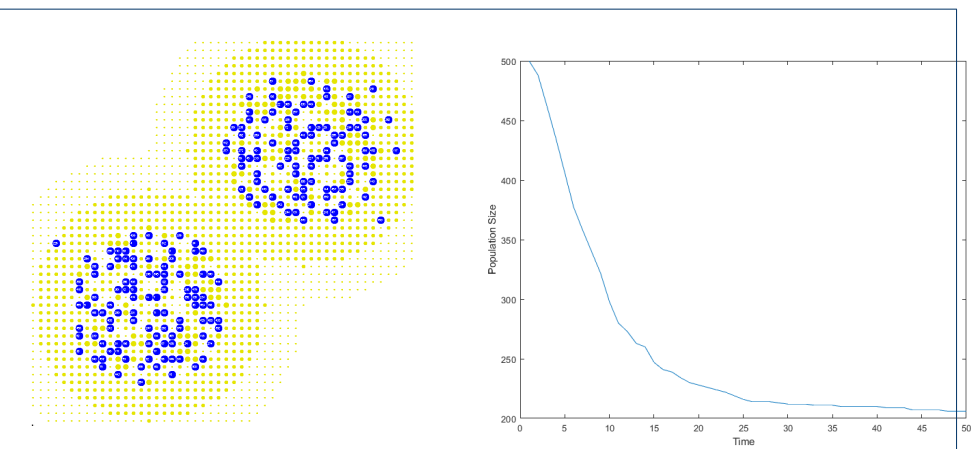


Figure 4: Performance comparison of *STM* and *Lock-Based SIR* implementations on varying grid sizes.



(a) Visualisation of the Sugarscape at $t = 50$

(b) Dynamics population size over 50 steps

Figure 5: Visualisation of our SugarScape implementation and dynamics of the population size over 50 steps. The white numbers in the blue agent circles are the agents unique ids.

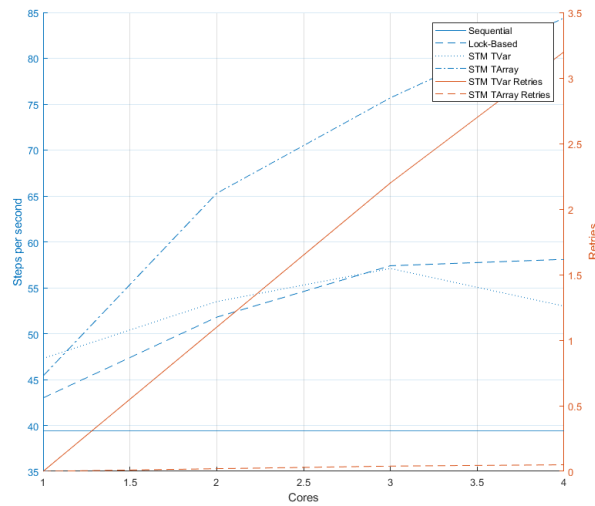


Figure 6: Comparison of steps per seconds (higher is better) and retries of various Sugarscape implementations with constant agent numbers. Using a 50x50 grid and 500 initial agents on varying cores.

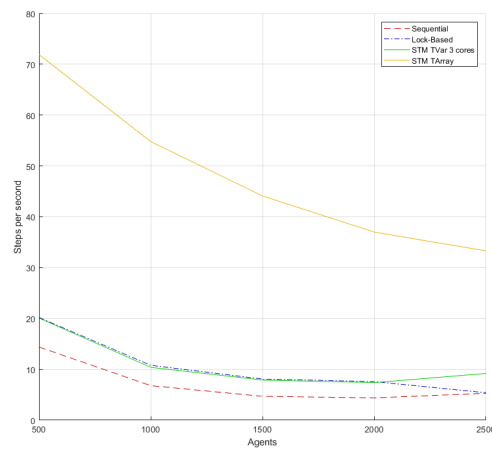


Figure 7: Steps per second (higher is better) of various Sugarscape implementations with varying agent numbers. Using 50x50 grid and varying number of agents with 4 (and 3) cores except the *Sequential* (1 core) implementation.