

MASARYK UNIVERSITY
FACULTY OF INFORMATICS



Parameter Synthesis from Hypotheses Formulable in CTL Logic

BACHELOR THESIS

Samuel Pastva

Brno, Spring 2015

Declaration

Hereby I declare, that this paper is my original authorial work, which I have worked out by my own. All sources, references and literature used or excerpted during elaboration of this work are properly cited and listed in complete reference to the due source.

Samuel Pastva

Advisor: prof. RNDr. Luboš Brim, CSc.

Acknowledgement

I would like to thank my supervisor...

Abstract

The main focus of this thesis is to provide a distributed algorithm for computation of the parameter synthesis problem for biochemical systems and CTL hypothesis and a usable implementation of said algorithm. The soundness and scalability of the algorithm is successfully demonstrated on biochemical models based on ordinary differential equations.

Keywords

Coloured Model Checking, CTL, Parameter Sythesis Problem, Systems Biology, Distribution

Contents

1	Introduction	1
2	Parameter synthesis problem	3
2.1	<i>Parametrised Kripke Structure</i>	3
2.2	<i>CTL Logic</i>	3
2.3	<i>Parameter Synthesis Problem</i>	4
2.4	<i>CTL Logic and Model Approximation</i>	5
3	State Space Distribution	7
3.1	<i>Kripke Fragments</i>	7
3.2	<i>Assumption Semantics</i>	9
4	Modelling techniques	12
4.1	<i>ODE models</i>	12
4.2	<i>Thomas Networks</i>	12
4.3	<i>State Space Partitioning</i>	12
5	Algorithm	15
5.1	<i>Distributed Environment</i>	15
5.2	<i>Algorithm outline</i>	16
5.3	<i>Temporal Operators</i>	17
5.3.1	<i>Exist Next Operator</i>	17
5.3.2	<i>Exist Until Operator</i>	18
5.3.3	<i>All Until Operator</i>	20
5.4	<i>Merge message buffer</i>	22
6	Implementation	25
6.1	<i>CTL Parser</i>	25
6.2	<i>Model Checker</i>	26
6.3	<i>ODE Abstraction</i>	26
6.3.1	<i>Thomas Network Abstraction</i>	27
6.3.2	<i>Frontend</i>	27
7	Experimental Evaluation	28
7.1	<i>Scaleability</i>	28
7.2	<i>Case study</i>	31
8	Conclusion	33

1 Introduction

With advancements in human knowledge about biochemical processes in living organisms, the complexity of experiments and theories needed to advance this knowledge even further is constantly increasing. This makes the search for easily auditable computational models describing dynamics of biochemical processes a key step towards understanding of the behavioural and physiological phenotypes occurring in biology.

Complex biochemical processes occurring in living organisms are usually dependant on high number of parameters, such as reaction rates or concentration values. While the qualitative aspects of some of these processes are already known, the quantitative aspects of these parameters are usually hard to measure and therefore cannot be easily determined. This leaves values of many parameters in biochemical systems uncertain or completely unknown.

The aim of the *parameter synthesis* is to find a maximum set of valuations of these unknown parameters, such that these valuations meet a stated dynamical constraints. These dynamical constraints can be captured in terms of temporal logic formulae.

A commonly used technique that decides wheter a given model meets specified temporal specification is *model checking* [1]. Several methods for parameter synthesis based on model checking have already been proposed with various modelling techniques and temporal logics in mind [2, 3, 4, 5, 6].

In this work, in order to express biologically relevant hypothesis, we focus on the branching time based logic, CTL. This is because many biological properties, such as multistability, need branching time operators to express them properly. It is very difficult or even impossible to express these properties correctly in linear time based logic, such as LTL.

Another advantage of CTL logic is that the particular procedure of model checking does not rely on initial conditions, effectively identifying all states where given property is satisfied. This leads to a global analysis of the model, as opposed to LTL model checking procedure, which requires a single initial state.

This work proposes a new method of solving the parameter syn-

thesis problem based on the CTL model checking. To tackle the problem of state space and parameter state explosion, we base our algorithm on distributed CTL model checking as defined in [7] and coloured model checking as introduced in [2]. We also provide an implementation on which we demonstrate the soundness and scalability of this new algorithm.

The work is divided into seven chapters. Second and third chapter introduce the notion of Kripke Structures, Kripke Fragments, State Space Partitioning and provide a short overview of employed modelling techniques. Fourth chapter describes the distributed algorithm itself. Fifth chapter discusses the details of the implementation. Finally, sixth chapter is dedicated to scalability and case study discussion.

2 Parameter synthesis problem

2.1 Parametrised Kripke Structure

As an input of our algorithm, we expect a parametrised Kripke Structure as defined in [2]. Parametrised Kripke Structure is a tuple $\mathcal{K} = (\mathcal{P}, S, S_0, \rightarrow, L)$ where

- \mathcal{P} is a finite set of parameters (all possible parameter valuations)
- S is a finite set of states
- $S_0 \subseteq S$ is a set of initial states
- $\rightarrow \subseteq S \times \mathcal{P} \times S$ is a transition relation labeled by parameter valuations
- $L : S \rightarrow 2^{AP}$ is a labeling function from states to sets of atomic propositions which are true in such states

We write $s \xrightarrow{p} s'$ when $(s, p, s') \in \rightarrow$. We also write $s \rightarrow s'$ when $\exists p \in \mathcal{P} : (s, p, s') \in \rightarrow$. Note that fixing a valuation $p \in \mathcal{P}$ reduces the Parametrised Kripke Structure \mathcal{K} to concrete, non-parametrised Kripke Structure $\mathcal{K}(p) = (S, S_0, \xrightarrow{p}, L)$.

2.2 CTL Logic

In order to correctly express various hypotheses in systems biology, the idea of branching time is needed. Example of such hypotheses are given in section 7.2. Therefore, this work uses the Computation Tree Logic (CTL) as means of hypotheses formulation.

CTL syntax is defined inductively upon finite set of atomic propositions AP :

$$\varphi ::= true \mid false \mid Q \mid \neg\varphi_1 \mid \varphi_1 \wedge \varphi_2 \mid EX\varphi_1 \mid E\varphi_1 U \varphi_2 \mid A\varphi_1 U \varphi_2 \quad (2.1)$$

Sometimes, we will use parentheses to make complex formulas easily readable, but they will in no way be used to modify the meaning of formula or priority of operators.

Note that there are also other temporal operators in standard CTL definition. We do not implement those directly in our algorithm. However, we use following equations to transform any general CTL formula prior to computation, so that it only uses operators supported in our algorithm. This way we can achieve a concise algorithm and also support whole CTL logic.

- $AX\varphi = \neg EX\neg\varphi$
- $EF\varphi = E(true)U\varphi$
- $EG\varphi = \neg A(true)U\neg\varphi$
- $AF\varphi = A(true)U\varphi$
- $AG\varphi = \neg E(true)U\neg\varphi$

Note that all of these transformations also preserve number of temporal operators in a formula. Other boolean operators like implication or equivalence can also be derived using similar transformations.

Let φ be a CTL formula. We write $cl(\varphi)$ to denote the set of all sub-formulas of φ and $tcl(\varphi)$ to denote the set of all temporal sub-formulas of φ . By $|\varphi|$ we denote the size of formula φ .

We assume standard CTL semantics over non-parametrised Kripke structures as defined in [1].

2.3 Parameter Synthesis Problem

Parameter synthesis problem is defined in following way. Suppose we are given a parametrised Kripke structure $\mathcal{K} = (\mathcal{P}, S, S_0, \rightarrow, L)$ and a CTL formula φ . For each state $s \in S$ let $P_s = \{p \in \mathcal{P} \mid s \models_{\mathcal{K}(p)} \varphi\}$, where $s \models_{\mathcal{K}(p)} \varphi$ denotes, that φ is satisfied in the state s of $\mathcal{K}(p)$. The parameter synthesis problem requires to compute the function $\mathcal{M}_\varphi^\mathcal{K} : S \rightarrow 2^\mathcal{P}$ such that $\mathcal{M}_\varphi^\mathcal{K}(s) = P_s$. Often we are especially interested in computing the set of all parameters for which the property

holds in some of the initial states: $\cap_{s \in S_0} \mathcal{M}_\varphi^\mathcal{K}(s)$. We can sometimes omit the φ and \mathcal{K} when they are clear from the context.

2.4 CTL Logic and Model Approximation

In model checking, some modeling approaches suffer from over or under approximation. We say that model is over-approximated when all feasible transitions are contained in the model, but it can also contain transitions that are not feasible in the situation the model is describing. Symmetrically, we say that model is under-approximated when all transitions in the model are feasible in the modeled situation, but not all feasible transitions has to be contained in the model.

It is important to discuss this relationship between CTL and approximated models, because it is much more complicated compared to linear-time logic, since CTL allows for universal and existential quantification mixing.

We say that CTL formula is *universal* or that it belongs to ACTL when it only contains universal temporal operators and no negation. Symmetrically, we say that CTL formula is *existential* or that it belongs to ECTL when it only contains existential temporal operators and no negation.

Observe that the truth of ACTL properties is preserved in over-approximated models. In other words, if an ACTL property holds in an over-approximated model, it must also hold in a model without approximation. However, their falsity cannot be guaranteed, because the false transitions may introduce paths that falsify the property in states where it would be normally true. Similarly, the falsity of ECTL properties in over-approximated models is preserved, but the truth is not. In this case, the existence of false transitions can introduce states where ECTL property holds solely due to these false paths.

Symmetrically, for under-approximated models, the falsity of ACTL and the truth of ECTL is preserved. But due to similar arguments, we can't say anything about their counterparts.

If we allow full CTL, in general, we can't make any assumptions about results obtained from either under- or over-approximated systems. This is caused by mixing of existential and universal quantification which leads to results that may be spurious and incomplete at

2. PARAMETER SYNTHESIS PROBLEM

the same time. Therefore, no conclusions can be made without further investigation and validation of such results.

3 State Space Distribution

3.1 Kripke Fragments

Due to the state space explosion, given parametrised Kripke structure can be very large and therefore impossible to fit into memory of one computer. In order to solve parameter synthesis problem for such structures, we have to distribute the state space across several computational nodes. To this end, we adapt the notion of parametrised kripke fragments as described in [7].

A parametrised Kripke structure \mathcal{K} can be divided into several parametrised Kripke fragments $\mathcal{F}_1^{\mathcal{K}}, \mathcal{F}_2^{\mathcal{K}}, \dots, \mathcal{F}_N^{\mathcal{K}}$ using a partition function $f : S \rightarrow \{1, \dots, N\}$. Parametrised Kripke fragment with identifier i over Kripke structure $\mathcal{K} = (\mathcal{P}, S, S_0, \rightarrow, L)$ is then defined as a tuple $\mathcal{F}_i^{\mathcal{K}} = (\mathcal{P}, S_i, I_i, \rightarrow_i, L_i)$ where:

- \mathcal{P} is a finite set of all parameters
- $S_i = \{s \in S \mid f(s) = i \vee \exists s' \in S. ((s \rightarrow s' \vee s' \rightarrow s) \wedge f(s') = i)\}$ is a subset of original state space which belongs to this fragment
- $I_i = \{s \in S_0 \mid s \in S_i\}$ is a set of all initial states that belong to the fragment
- $\rightarrow_i = \{(s, p, s') \in \rightarrow \mid s \in S_i \wedge s' \in S_i \wedge (f(s) = i \vee f(s') = i)\}$ is a subset of original transition relation reduced to only relevant states (not required to be total)
- $L_i = \{(s, l) \in L \mid f(s) = i\}$ is a subset of original labeling function relevant to this fragment

In the following text, we will often omit the superscript \mathcal{K} if it is clear from context.

Intuitively, Kripke fragment represents a subset of the original kripke structure defined by the partition function. It contains all nodes specified by the partition function with all their direct successors and predecessors and all corresponding transitions.

We define a set of border states $border(\mathcal{F}_i^K) = \{s \in S_i \mid \neg \exists (p, s'). s \xrightarrow{p} s'\}$. Intuitively, these states represent the remaining portion of the state space which is stored in memory of other processes and is not directly accessible. We say that state is an *internal state* if it is not a border state. We also define a set of cross edges as $cross(\mathcal{F}_i^K) = \{(s, p, s') \in \rightarrow_i \mid s' \in border(\mathcal{F}_i^K) \vee s \in border(\mathcal{F}_i^K)\}$. Intuitively, these are edges leading from border states to internal states or vice versa.

Note that for a constant partition function $f(s) = 1$ and any given kripke structure, the partitioning results in one fragment with an unchanged state space and an unchanged transition relation (The sets of border states and cross edges are empty for the resulting fragment). Similarly, for every Kripke structure, there exists a partition function for which $N = |S|$ and every resulting fragment contains only one internal state. Under such partitioning, all edges are cross edges.

In worst case (connected graphs), such partitioning results in fragments with $|S| - 1$ border states and one internal state (there is no way to achieve higher state count in fragment than in the original structure). Therefore we can assume that $\sum_{i=1}^N |S_i| \leq |S|^2$, hence the number of introduced border states is at worst quadratic in terms of original state space. The number of internal states remains allways the same.

This increase seems rather high, however, it is important to note that this also requires one process for each original state. In real life scenarios, the number of states per process is usually much higher. Also note that the representation of border state in memory is usually much simpler (and smaller) than representation of internal state, so even a distribution with high border state count can be beneficial in terms of memory consumption.

In terms of edges, in the worst case, each edge has to be contained in two fragments (where either of the end states is an internal state), therefore the total number of edges is at worst doubled.

The number of border states and cross edges is also highly dependent on the partition function. It is usually best to design the partition function with specific model or modelling approach in mind in order to achieve optimal workload distribution. We will discuss different partition functions later in section 4.3.

For each kripke fragment \mathcal{F}_i , we define a function $successors \subseteq S_i \rightarrow S_i \times 2^{\mathcal{P}}$ to denote the set of all successor states and relevant colour sets that enable such transition.

$$successors(s) = \{(s', P) \in S_i \times 2^{\mathcal{P}} \mid s \rightarrow_i s' \wedge P = \{p \in \mathcal{P} \mid (s, p, s') \in \rightarrow_i\}\}$$

Symetrically, we define a function $predecessors \subseteq S_i \rightarrow S_i \times 2^{\mathcal{P}}$ to denote the set of all predecessor states and relevant colour sets that enable such transition.

$$predecessors(s) = \{(s', P) \in S_i \times 2^{\mathcal{P}} \mid s' \rightarrow_i s \wedge P = \{p \in \mathcal{P} \mid (s', p, s) \in \rightarrow_i\}\}$$

3.2 Assumption Semantics

Classic interpretation of CTL formulas is not adequate for Kripke fragments. In order to accommodate for possible non-totality and distributed nature of the Kripke fragments over Kripke structure $\mathcal{K} = (\mathcal{P}, S, S_0, \rightarrow, L)$, we introduce the assumption function $\mathcal{A} : \mathcal{P} \times S \times cl(\varphi) \rightarrow Bool$. The values $\mathcal{A}(s, p, \varphi_1)$ are called assumptions. We use the notation $\mathcal{A}(p, s, \varphi_1) = \perp$ to say that the value of $\mathcal{A}(s, p, \varphi_1)$ is undefined. By \mathcal{A}_{\perp} we denote assumption function which is undefined for all inputs.

Intuitively, $\mathcal{A}(s, p, \varphi_1) = \text{tt}$ when we can assume that φ_1 holds in a state s for parameter valuation p , $\mathcal{A}(s, p, \varphi_1) = \text{ff}$ when we can assume that φ_1 does not hold in state s for parameter valuation p and $\mathcal{A}(s, p, \varphi_1) = \perp$ when we cannot assume anything about validity of φ_1 in state s for parameter valuation p . We write $AS_{\mathcal{F}_i}^{\varphi}$ to denote the set of all assumption functions for a formula φ and parametrised kripke fragment \mathcal{F}_i .

Undefined values are important in CTL semantics over distributed fragments, since such values can be used in places where validity of formula cannot be computed because it depends on information stored in another fragment which has not been received yet. However, a correct model checking algorithm should always provide a definitive answer for all states and parameter valuations in finite time.

For a Kripke fragment $\mathcal{F}_i = (\mathcal{P}, S_i, I_i, \rightarrow_i, L_i)$ and a formula φ , the assumption function is defined inductively on the structure of the formula φ :

$$\mathcal{A}(p, s, Q) = \begin{cases} \text{tt} & Q \in L_i(s) \\ \text{ff} & \text{otherwise} \end{cases}$$

$$\mathcal{A}(p, s, \neg\varphi_1) = \begin{cases} \text{tt} & \mathcal{A}(p, s, \varphi_1) = \text{ff} \\ \text{ff} & \mathcal{A}(p, s, \varphi_1) = \text{tt} \\ \perp & \text{otherwise} \end{cases}$$

$$\mathcal{A}(p, s, \varphi_1 \wedge \varphi_2) = \begin{cases} \text{tt} & \mathcal{A}(p, s, \varphi_1) = \text{tt} \text{ and } \mathcal{A}(p, s, \varphi_2) = \text{tt} \\ \text{ff} & \mathcal{A}(p, s, \varphi_1) = \text{ff} \text{ or } \mathcal{A}(p, s, \varphi_2) = \text{ff} \\ \perp & \text{otherwise} \end{cases}$$

$$\mathcal{A}(p, s, \text{EX}\varphi_1) = \begin{cases} \text{tt} & \exists s' \in S_i : s \xrightarrow{p}_i s' \wedge \mathcal{A}(p, s', \varphi_1) = \text{tt} \\ \text{ff} & \forall s' \in S_i : s \xrightarrow{p}_i s' \Rightarrow \mathcal{A}(p, s', \varphi_1) = \text{ff} \\ \perp & \text{otherwise} \end{cases}$$

$$\mathcal{A}(p, s, \text{E}\varphi_1 \text{U}\varphi_2) = \begin{cases} \text{tt} & \text{exists a p-path } \pi = s_0 s_1 s_2 \dots \text{ with } s = s_0 \text{ such} \\ & \text{that } \exists x < |\pi| \text{ such that (either } \mathcal{A}(p, s_x, \varphi_2) = \\ & \text{tt or } [s_x \in \text{border}(\mathcal{F}_i) \text{ and } \mathcal{A}(p, s_x, \text{E}\varphi_1 \text{U}\varphi_2) = \\ & \text{tt}] \text{), and } \forall 0 \leq y < x : \mathcal{A}(p, s_y, \varphi_1) = \text{tt} \\ \text{ff} & \text{for all p-paths } \pi = s_0 s_1 s_2 \dots \text{ with } s = s_0 \text{ ei-} \\ & \text{ther } \exists x < |\pi| \text{ such that } (\mathcal{A}(p, s_x, \varphi_1) = \text{ff and} \\ & \forall y \leq x : \mathcal{A}(p, s_y, \varphi_2) = \text{ff}) \text{ or } [\forall x < |\pi| : \\ & \mathcal{A}(p, s_x, \varphi_2) = \text{ff and } (|\pi| = \infty \text{ or } (s_{|\pi|-1} \in \\ & \text{border}(\mathcal{F}_i) \text{ and } \mathcal{A}(p, s_{|\pi|-1}, \text{E}\varphi_1 \text{U}\varphi_2) = \text{ff}))] \\ \perp & \text{otherwise} \end{cases}$$

$$\mathcal{A}(p, s, A\varphi_1 U \varphi_2) = \begin{cases} \text{tt} & \text{for all p-path } \pi = s_0 s_1 s_2 \dots \text{ with } s = s_0 \text{ such that} \\ & \exists x < |\pi| \text{ such that [either } \mathcal{A}(p, s_x, \varphi_2) = \text{tt or} \\ & (s_x \in \text{border}(\mathcal{F}_i) \text{ and } \mathcal{A}(p, s_x, A\varphi_1 U \varphi_2) = \text{tt})], \\ & \text{and } \forall 0 \leq y < x : \mathcal{A}(p, s_y, \varphi_1) = \text{tt} \\ \text{ff} & \text{exists a p-path } \pi = s_0 s_1 s_2 \dots \text{ with } s = s_0 \text{ and} \\ & \text{an index } x < |\pi| \text{ such that: } (\mathcal{A}(p, s_x, \varphi_1) = \text{ff} \\ & \text{and } \forall y \leq x : \mathcal{A}(p, s_y, \varphi_2) = \text{ff}) \text{ or } [\forall x < |\pi| : \\ & \mathcal{A}(p, s_x, \varphi_2) = \text{ff and } (|\pi| = \infty \text{ or } (s_{|\pi|-1} \in \\ & \text{border}(\mathcal{F}_i) \text{ and } \mathcal{A}(p, s_{|\pi|-1}, E\varphi_1 U \varphi_2) = \text{ff}))] \\ \perp & \text{otherwise} \end{cases}$$

Here a p-path π from a state s_0 is a sequence $\pi = s_0 s_1 \dots$ such that $\forall j \geq 0 : s_j \in S_i$ and $s_j \xrightarrow{p} s_{j+1}$.

We define a so called semantic function $\mathcal{C}_{\mathcal{F}_i}^\varphi : AS_{\mathcal{F}_i}^\varphi \rightarrow AS_{\mathcal{F}_i}^\varphi$ which takes an input assumption function \mathcal{A}_{in} and computes a new assumption function \mathcal{A} as defined in this section. Note that for a Kripke fragment with total transition relation (without any border states), the value of resulting assumption function does not depend on the input assumptions. Hence for a total Kripke Structure \mathcal{K} , we can solve the parameter synthesis problem by computing the assumption function $\mathcal{C}_{\mathcal{F}}^\varphi(\mathcal{A}_\perp)$ where \mathcal{F} is the single-fragment representation of \mathcal{K} resulting from the trivial partitioning as described in previous section.

We also defined function $initialStates : cl(\phi) \times AS_{\mathcal{F}_i}^\varphi \rightarrow S_i \times 2^{\mathcal{P}}$ which computes a set of states and parameters where truth of given formula is assumed.

$$initialStates(\phi, \mathcal{A}) = \{(s, p) \in S_i \times 2^{\mathcal{P}} \mid \mathcal{A}(s, p, \phi) = \text{tt}\}$$

Similarly, we use a function $validColours : cl(\phi) \times S_i \times AS_{\mathcal{F}_i}^\varphi \rightarrow 2^{\mathcal{P}}$ which returns a set of parameters for which given formula is assumed to be true in given state.

$$validColours(\phi, s, \mathcal{A}) = \{p \in \mathcal{P} \mid \mathcal{A}(s, p, \phi) = \text{tt}\}$$

4 Modelling techniques

To prove soundness and universality of our algorithm, we test it on two different modelling techniques.

4.1 ODE models

First technique is based on discretization of piece-wise multi-affine models defined using Ordinary Differential Equations (ODE). More about this technique can be found in [8]. ODE models have rectangular state space with low number of transitions per state and good locality, because only transitions between adjacent states are allowed. Since parameter space is continuous, it is usually symbolically represented in form of intervals. Unfortunately, ODE models suffer from over-approximation and therefore are not very well suited for verification of properties with mixed existential and universal quantification.

4.2 Thomas Networks

Second technique is based on Thomas Networks [9] which are extension of Boolean Networks. Compared to ODE models, Thomas Networks are inherently discrete. They also do not suffer from over- or under-approximation. The downside of Thomas Networks are parameters, since even small models can contain a very large number of parameters and possible parameter valuations. Usually, it is also hard to find a compact and reasonably fast symbolic representation of said parameters. The influence of parameters on transition system also tends to be much more irregular compared to ODE models. This makes them less suitable for state space distribution.

4.3 State Space Partitioning

One of the most important aspects of efficient distributed model checking algorithms is a suitable state space partitioning. One that minimizes the communication overhead and provides an equalized work-

load distribution. In particular, the partitioning should provide a regular load-balancing, ensuring that each process is responsible for a proportional part of the state space. It should also provide a good locality, minimizing the number of cross transitions and therefore reducing the communication overhead.

The problem of computing the optimal partitioning can introduce a significant overhead into the computation. Therefore, various heuristics are considered to produce partitioning that is easy to compute and at the same time provides reasonable load-balancing and locality. One of such heuristics is a hash based partitioning [10] which is usually used for computer and engineering systems. This partitioning does not require any prior knowledge about the structure of the state space, and therefore can be universally applied to almost any system.

It is based on a hash function which maps each state to a process. This approach usually results in in very good load balancing thanks to the uniformity of the hash function. However, hash based partitioning can't control the locality and therefore may introduce a high number of cross transitions into the system. This can significantly increase communication overhead.

In this work, we exploit the regular structure of the state space of biochemical models in order to produce partitioning that does not suffer from this kind of locality problem. We use structural properties of the rectangular abstraction of the given parametric piece-wise multi-affine ODE model [8]. The approximation is formed by an n -dimensional hyper-rectangular state space defined by n state variables and by a set of thresholds for each variable. The discretization of the state space also ensures that there are only transitions between adjacent states with respect to the hyper-rectangular structure.

Our partitioning decomposes said state space into p hyper-rectangular subspaces such that all subspaces have similar state count (p is the number of processes). This heuristic usually provides good load balancing, since the states are evenly distributed across all processes. However, compared to classic hash based partitioning, it can also provide better locality thanks to the fact, that transitions only occur between adjacent states. This ensures that cross transitions can originate only in states on borders of these hyper-rectangular subspaces. Which in turn provides almost minimal number of cross transitions

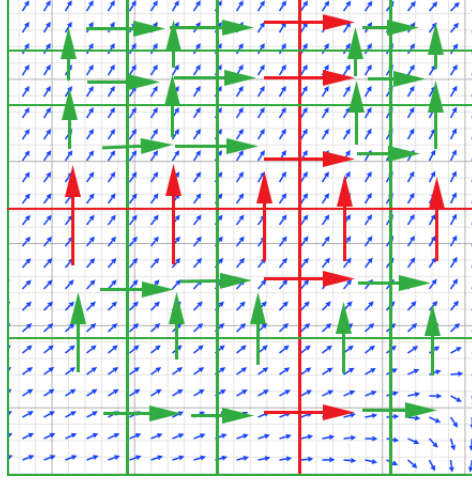


Figure 4.1: Rectangular state space partitioning with $n = 4$. Green arrows are internal transitions, red arrows represent cross transitions.

and therefore significantly reduces the communication overhead.

However, note that this is still only a heuristic and the results can be negatively influenced by the backward connectivity of the state space. Especially in systems with uneven distribution of transitions across states, better partitionings can be constructed that minimize the number of cross transitions while at the same time maintaining good load-balancing. On the other hand, our experiments demonstrate that due to the fact, that we have to consider all possible parametrisations, the connectivity of the state space is significantly increased. Therefore, our heuristic can for majority of models create a partitioning with almost minimal number of cross transitions.

Similar heuristic can be deployed also for the Thomas networks, however, this modeling technique has a less predictable transition system and therefore the minimal number of cross transitions is not that easily achieved.

5 Algorithm

In this chapter, we describe the distributed algorithm that computes the semantic function $\mathcal{C}_{\mathcal{F}_i}^\varphi$.

5.1 Distributed Environment

In this section, we briefly describe the distributed environment assumed by our algorithm, in order to prevent any possible confusion.

We assume a distributed environment with fixed number of reliable processes connected by reliable, order-preserving channels (The order preservation can be relaxed to some extent). We also assume that each process has a fixed identifier and the set of all process identifiers is equal to the result set of the partition function. Each process can communicate directly (using the function `SEND`) with any other process (assuming it knows other process's identifier). We assume that each message can be transmitted in $O(1)$ time and all messages that can't be processed directly are stored in a buffer until they can be processed.

Several parts of the algorithm do not have explicit termination (they terminate by reaching deadlock - no messages are exchanged between processes). In such cases, suitable termination detection algorithm is employed to detect this deadlock and terminate computation properly. Our implementation uses Safra's algorithm [11] for this purpose, but the related code has been skipped for easier readability.

The algorithm is broken into two main parts. First describes the general outline of algorithm and is similar to classic CTL model checking. Second part describes how each of the supported temporal operators is processed. This part contains more detailed description of inter-process communication and operator specific data structures. To better reflect the distributed nature of the algorithm, description of each temporal operator routine is divided into three parts: Process variables, Initialization and Message handler. First section describes data structures stored in process memory available during whole computation. Initialization section is executed exactly once and no messages can be received until it's finished. Message handler defines

what should happen when message is received.

5.2 Algorithm Outline

In order to compute the assumption function \mathcal{A} for the given kripke structure, we compute the value of $\mathcal{C}_{\mathcal{F}_i}^\varphi(\mathcal{A}_\perp)$ by exchanging relevant assumptions between processes. The main idea of the algorithm is described in Algorithm 5.2 and resembles other CTL model checking algorithms.

```

1: procedure CHECKCTL( $\psi, \mathcal{K} = (\mathcal{P}, S, S_0, \rightarrow, L)$ )
2:    $\mathcal{A} \leftarrow \mathcal{A}_\perp$ 
3:   for all  $i < |\psi|$  do
4:     for all  $\varphi$  in  $cl(\psi)$  where  $|\varphi| = i$  do
5:        $\mathcal{A} \leftarrow \text{CHECKOPERATOR}(\varphi, \mathcal{K}, \mathcal{A})$ 
6:       set  $\mathcal{A}(p, s, \varphi) = \text{ff}$  where  $\mathcal{A}(p, s, \varphi) = \perp$ 
7:     end for
8:   end for
9: end procedure

```

The algorithm starts by initializing the assumption function to undefined values. After that, it traverses the structure of formula, starting from smallest formulas and uses previously computed results to process more complex formulas. Function CHECKOPERATOR computes all states and colours where formula φ holds and returns assumption function updated accordingly. This is done using the local information contained in given kripke fragment, assumptions previously computed for smaller formulas and also by communicating with other processes. Note that only assumptions relevant for particular process are computed and returned (each process has information only about it's own state space). After this information has been computed, we can guarantee that the formula holds only in states marked accordingly, so we remove all undefined values by marking them as false.

5.3 Temporal Operators

In this section, we describe how the CHECKOPERATOR is implemented for each of the temporal operators. Note that all of the following algorithms has implicit termination and therefore needs a proper termination detection algorithm to correctly terminate.

We do not provide algorithm for boolean operators and atomic propositions, since these can be trivially deduced from the assumption function semantics.

5.3.1 Exist Next Operator

Algorithm 1 Exists next

```

1: Process variables:
2:  $\mathcal{F}_i = (\mathcal{P}, S_i, I_i, \rightarrow_i, L_i)$  ▷ Kripke fragment
3:  $\phi = \text{EX}\phi_1$  ▷ CTL formula
4:  $\mathcal{A}$  ▷ Initial assumption function
5:  $f$  ▷ Partition function
6: procedure INIT
7:   for all ( $state, colSet$ ) in  $initialStates(\phi_1, \mathcal{A})$  do
8:     for all ( $pred, tranCol$ ) in  $predecessors(state)$  do
9:       SEND( $f(state), (pred, colSet \cap tranCol)$ )
10:    end for
11:  end for
12: end procedure
13: procedure RECEIVE( $dest, colSet$ )
14:   set  $\mathcal{A}(p, dest, \phi) = \text{tt}$  for all  $p \in colSet$ 
15: end procedure

```

If formula is marked as valid in state s for colour p , it means a message containing such state and colour has been received. This can only happen if said state has a successor under colour p where ϕ_1 holds. Therefore no false positive results are produced. Also, for each state where ϕ_1 holds for colour p , a message is sent for all predecessors of such state. Therefore all states where $\text{EX}\phi_1$ holds are labeled accordingly. Since all correct states are labeled and no false positives are possible, the algorithm is correct.

In the worst case, algorithm has to send every colour over every transition, therefore worst case message complexity is $card(\rightarrow_i)$. In practice, message count is usually much lower, since multiple colours can be packed into one message.

Assuming the validity of φ_1 is computed for all states, function $initialStates(\varphi_1, \mathcal{A})$ can be computed in $O(|S_i| \cdot |\mathcal{P}|)$ time. The function $predecessors$ can be pre-computed for all states in $O(card(\rightarrow_i))$ time. The procedure SEND is called at most $card(\rightarrow_i)$ times and the parameter set intersection is at worst linear in the size of parameter space.

This would result in $O(|\mathcal{P}| \cdot card(\rightarrow_i))$ time complexity. However, this can be further reduced to $O(card(\rightarrow_i))$ since we can observe that for every predecessor where $|tranCol| > 1$, only one message is sent instead of $|tranCol|$ messages. The colour set intersection can also be performed in $O(|tranCol|)$ time. This means that the price of set intersection is amortized by the reduced number of transmitted messages.

5.3.2 Exist Until Operator

The EU operator is a little more complex, but again fairly simple. The algorithm starts by computing all states and colours where φ_2 is true and marks them as valid. Starting from these states, a backpropagation of parameter sets along the reversed transitions is performed. During the computation, the propagated parameter set is updated to reflect the validity of φ_1 in examined state and the validity of transitions used along the path. Note that backpropagation is stopped as soon as there is no new information computed ($colSet$ is either empty or equal to already computed assumptions).

Algorithm reaches deadlock, because each message either adds true values to the assumption function or does not trigger sending of new messages. Since we have a finite state and parameter space, we can't add true values to the assumption function forever, and therefore at some point, no new messages will be sent.

Algorithm is correct, because φ is marked as true in some state and parameter only if φ_2 holds there, or φ_1 is true in this state and a message is received from some of the state's successors indicating that φ holds in said successor. Also, if φ is marked as true in some

Algorithm 2 Exists until

```
1: Process variables:
2:  $\mathcal{F}_i = (\mathcal{P}, S_i, I_i, \rightarrow_i, L_i)$  ▷ Kripke fragment
3:  $\varphi = E\varphi_1 U \varphi_2$  ▷ CTL formula
4:  $\mathcal{A}$  ▷ Initial assumption function
5:  $f$  ▷ Partition function
6: procedure INIT
7:   for all  $(state, colSet)$  in  $initialStates(\varphi_2, \mathcal{A})$  do
8:     set  $\mathcal{A}(p, state, \varphi) = \text{tt}$  for all  $p \in colSet$ 
9:     for all  $(pred, tranCol)$  in  $predecessors(state)$  do
10:       $SEND(f(state), (pred, colSet \cap tranCol))$ 
11:    end for
12:  end for
13: end procedure
14: procedure RECEIVE( $s, colSet$ )
15:    $colSet \leftarrow colSet \cap validColours(\varphi_1, s, \mathcal{A})$ 
16:   if  $colSet \neq \emptyset$  and  $colSet \setminus validColours(\varphi, s, \mathcal{A}) \neq \emptyset$  then
17:     set  $\mathcal{A}(p, s, \varphi) = \text{tt}$  for all  $p \in colSet$ 
18:     for all  $(pred, tranCol)$  in  $predecessors(s)$  do
19:       $SEND(f(pred), (pred, colSet \cap tranCol))$ 
20:    end for
21:   end if
22: end procedure
```

state, all of it's predecessors are allways notified about this update.

Worst case message complexity is $card(\rightarrow_i)$ because at worst, we have to send every colour along every transition. Note that we do not use the same parametrised transition twice, because this would mean that φ was also set true twice for the same state and colour (We prevent this using the condition on line 16).

Because all general functions have linear or better complexity and we know that for each transition and parameter, only one message is sent, using similar argument as in analysis of EX, we can show that the time complexity of this algorithm is $O(card(\rightarrow_i))$.

5.3.3 All Until Operator

The AU operator is the most complex one to handle. As opposed to EX, which requires at least one valid successor to be true, AU requires that all successors of the specific node are valid. In order to compute such information, we create a copy of transition relation and call it T .

During the computation, T is modified in such way, so that we can guarantee that if edge is not present in T , this edge does not exist or it leads to a state where either φ_2 or $A\varphi_1 U \varphi_2$ holds. This way, we are keeping track of all edges that can falsify the validity of φ and therefore when a state has no successors in T , we can mark it as valid.

Algorithm reaches deadlock thanks to the same argument as used in analysis of EU.

Algorithm is correct, because a state and colour are marked as true only if they have no successors in T and φ_1 is also true. A state and parametrisation has no successor in T only if all transitions valid for said parametrisation lead to a state where either φ_2 or $A\varphi_1 U \varphi_2$ holds. If a state is marked as true for some parametrisation, or φ_2 is true in it, an appropriate message is sent to all of it's predecessors. This guarantees the removal of this transition from T .

Worst case message complexity of the algorithm is $card(\rightarrow_i)$ thanks to the same argument as used in analysis of EU operator. Worst case time complexity is $O(maxOutDegree(\rightarrow_i) \cdot card(\rightarrow_i))$. This is due to the fact that operation on line 16 needs to check all outcoming edges from state s in order to confirm that they are removed from set T . This is not an optimal time complexity (other methods based on search

Algorithm 3 All until

```
1: Process variables:
2:  $\mathcal{F}_i = (\mathcal{P}, S_i, I_i, \rightarrow_i, L_i)$  ▷ Kripke fragment
3:  $\varphi = A\varphi_1 U \varphi_2$  ▷ CTL formula
4:  $T = \rightarrow_i$  ▷ Uncovered edges
5:  $\mathcal{A}$  ▷ Initial assumption function
6:  $f$  ▷ Partition function
7: procedure INIT
8:   for all  $(state, colSet)$  in  $initialStates(\phi_2, \mathcal{A})$  do
9:     for all  $(pred, tranCol)$  in  $predecessors(state)$  do
10:      SEND( $f(state), (state, pred, colSet \cap tranCol)$ )
11:    end for
12:  end for
13: end procedure
14: procedure RECEIVE( $d, s, colSet$ )
15:    $T \leftarrow T \setminus \{(s, p, d) \mid p \in colSet\}$ 
16:    $colSet \leftarrow \{p \in \mathcal{P} \mid p \in colSet \wedge \forall s' \in S : (s, p, s') \notin T\}$ 
17:    $colSet \leftarrow colSet \cap valid(\varphi_1, s, \mathcal{A})$ 
18:   if  $colSet \neq \emptyset$  and  $colSet \setminus validColours(\varphi, s, \mathcal{A}) \neq \emptyset$  then
19:     set  $\mathcal{A}(p, s, \varphi) = \text{tt}$  for all  $p \in colSet$ 
20:     for all  $(pred, tranCol)$  in  $predecessors(s)$  do
21:       SEND( $f(pred), (to, pred, colSet \cap tranCol)$ )
22:     end for
23:   end if
24: end procedure
```

for strongly connected components can solve this problem in linear time), however, implementing these methods in distributed environment would result in increased communication and synchronization overhead. Also, typical biochemical models have a guaranteed small average maximal out degree, so that it can be effectively considered a constant coefficient.

5.4 Merge Message Buffer

The main argument of coloured model checking efficiency is based on the assumption that operations on parameter sets are in practice less expensive than graph traversal. However, even a simple model structure can easily break the computation into many simultaneous traversals. When two different colour sets are marked as valid in a state due to two distinct transitions, unfortunate timing can easily prevent these two colour sets from merging into one. This leads to two outgoing messages instead of one, which in the end results in two separate graph traversals instead of one.

In some cases, this simply cannot be avoided, since in order to know exactly when to wait for a merge and when to continue the traversal, we would basically have to solve the coloured model checking problem. However, we can take advantage of the fact, that many messages cannot be processed directly at the time of arrival, since message processing can be a costly operation, especially when the transition system is being lazily calculated. Therefore we use a message buffer to store incoming messages that cannot be processed directly.

These buffers can grow quite large during the computation and in extreme cases even outgrow the state space of the transition system itself. Of course, it would be easy to just slow down the creation of messages to prevent the buffers from growing. However, many messages in these buffers contain data that are either duplicate or can be merged into one message while maintaining correct semantics.

In order to reduce the number of unnecessary graph traversals and reduce memory footprint of message buffers while maintaining good performance, we employ a merge message buffer as described

in algorithm 5.4.

In this text, we only provide the implementation for messages that contain destination node and colour set, as used in EX and EU operators. However, the implementation for messages used by AU operator can be obtained trivially by replacing all occurrences of destination node with pair of destination and source nodes.

```

1: MergeBuffer
2:  $Q \leftarrow \text{IterableHashMap}$ 
3: procedure INSERT( $node$ ,  $colours$ )
4:   if  $Q$  hasKey  $node$  then
5:      $current \leftarrow Q.get(node)$ 
6:      $Q.replace(node, colours \cup current)$ 
7:   else
8:      $Q.insert(node, colours)$ 
9:   end if
10: end procedure
11: procedure EMPTY
12:   return  $Q.isEmpty()$ 
13: end procedure
14: procedure TAKE
15:    $val \leftarrow Q.iterator().first()$ 
16:    $Q.remove(val.key)$ 
17:   return  $val$ 
18: end procedure

```

Key property of the buffer is that it is backed by an *IterableHashMap*. Compared to traditional *HashMap*, *IterableHashMap* provides also an iterator on all of its key-value pairs. This gives us the ability to take one (non-deterministic) element out of the map in constant time. Considering a reasonable hash function on the node set, we can guarantee that each operation on the underlying *IterableHashSet* can be done in constant time. Therefore the only interesting operation in terms of time complexity is the union of two colour sets, which can be done in $O(\mathcal{P})$ time.

It is important to note that the resulting buffer does not preserve the FIFO properties of a queue. However, this is not required by the algorithm (actually, performing a random search instead of classic

DFS is considered a valid optimization heuristic).

We do not provide exhaustive benchmark of our implementation, since the main priority of this work is the model checking algorithm itself. However, a comparison with linked queue and circular buffer on models tested in the Section 7.1 showed, that merge queue easily outperforms both of them especially in highly distributed environments. Linked queue was usually two times slower while circular buffer managed to provide approximately 70-90% of the merge buffer performance.

One reason for this, is the fact that more distributed computations have generally more cross transitions. This increases the cost of traversal compared to parameter set operations which in turn makes the merge buffer more effective. Also, the properties of this buffer allow in some cases for super-linear scalability, since greater number of processes can produce higher number of merged traversals.

6 Implementation

Our algorithm is implemented in a proof-of-concept distributed CTL model checker available in the github repository of Sybila organization [12]. In this chapter, we discuss the implementation details of this project.

Although the core model checking module is fully operational and stable, the project is still mainly in experimental phase, since new features, heuristics and modelling approaches are still being considered and reevaluated.

In order to provide greater flexibility and ease of development, majority of the project is implemented in Java. However, several parts still require C++ code reused from similar previous projects. The whole project is built using Gradle build system.

The distributed environment is mainly provided by MPJ, Java implementation of standard MPI interface. However, the model checking algorithm itself does not depend on any concrete communication library and can be easily adapted to any similar distributed communication tool.

The whole project is divided into several modules in order to maximize extensibility and modifiability during future development. This section provides a quick overview of each of these modules.

6.1 CTL Parser

In order to provide user with easy way to input CTL formulas, this module implements a parser for modified formula specification language used in Biodivine [13].

The original parser only supports LTL formulas, therefore the grammar has been modified to allow for CTL operators. However, user familiar with the original LTL syntax should have no problems adapting to the CTL syntax.

The grammar is written in Antlr parser generator and apart from the whole range of CTL and boolean operators supports also boolean and float propositions.

This module also handles the transformation of provided formula into minimal set of operators supported by the main model checker.

6.2 Model Checker

This is the core module implementing the model checking algorithm described in this work. It defines the interfaces and contracts representing the kripke fragment, parameter set, partition function and inter-process communication.

It is completely independent on the implementation of said interfaces, and therefore provides great flexibility and extensibility. This module also provides partial implementations of some of these interfaces to ease the development of new modules.

The termination detection algorithm designed by Safra [11] is also implemented here. The termination detection itself can't be easily overridden, however, the communication is again isolated into one, easily replaceable class, so the support for different communication libraries can be easily provided.

6.3 ODE Abstraction

This module is based on the ODE state space generator from the LTL model checking tool Biodivine [13]. The code responsible for model parsing and abstraction calculation is reused directly and Java Native Interface (JNI) is used to extract resulting model into corresponding Java data structures.

The state space generator responsible for evaluation of atomic propositions and computation of successors/predecessors is also based on code from the Biodivine project. However, this section has been completely rewritten into Java and adapted to the CTL paradigm. This minimizes the number of slow JNI calls to existing C++ code. This new state space generator also features several bug fixes and speed optimizations.

Module also implements the rectangular state space partition function described in section 4.3. The parameter set is implemented using Google Guava Range Set, which provides great flexibility and huge feature set while maintaining good performance.

6.3.1 Thomas Network Abstraction

This module is based on the Thomas Network state space generator from the LTL model checking tool Parsybone [14]. Similarly to the ODE Abstraction module, most of the model parsing and preprocessing is done using the original code and extracted using JNI.

The parameter set is implemented using a EWAH Compressed Bitmap, which provides decent performance even on large parameter sets and is very easy to use.

At the moment of writing, this module does not provide any good partitioning implementation, since much of its functionality is still in development and most of the current implementation is subject to change. However, sequential computation on one processing node is fully supported.

6.3.2 Frontend

Frontend module ties together the functionality of all modules into several runnable utilities. However, the documentation and general outline and output format of these tools have not yet been finalized and is subject to change. Therefore the code in this module should be taken more as an example of usage of different modules.

7 Experimental Evaluation

The aim of this chapter is to demonstrate the soundness of our algorithm on two biologically relevant models. First model is fairly simple and is used for benchmarking purposes, since it can be easily scaled and modified while maintaining similar properties and structure. Second model exhibits non trivial behaviour and is used to demonstrate the ability of our method to correctly detect such behaviour.

7.1 Scaleability

We evaluate the scalability of the algorithm on an ODE model of reversible catalytic reaction with varying number of intermediate enzyme-substrate complexes. Using this model, we can scale the number of intermediate products/variables (N), discretization thresholds (T) and unknown parameters. For simplicity, we assume that each variable is evaluated on the same amount of thresholds, which results in the total state space size of $(T - 1)^N$.

The benchmarks were performed on a cluster of 12 computational nodes, each equipped with 16 GB of RAM memory and a quad-core Intel Xeon 2Ghz processor. In order to focus on the distribution and not multi-core performance, we utilize only one core on each machine. However, note that the algorithm can be easily executed on a multicore machine or a cluster of multicore machines using each core as separate process. No other resource intensive software was running at these machines at the time of benchmarking.

The model itself is described in the figure 7.1. The first line represents a simplified chemical equation of the model. The following lines describe differential equations for every species. The last two lines describe used parameters.

In the subsequent tables, we provide detailed information about the runtime of the distributed algorithm on various modifications of the model and a CTL formula $AG(P \leq 30)$. The value N/A represents a situation when the algorithm ran out of memory. The results have been obtained as an arithmetic mean of several experiments.

They generally illustrate very good scaleability in terms of eval-

$$\begin{array}{c}
\frac{S + E \rightleftharpoons ES_1 \rightleftharpoons \dots \rightleftharpoons ES_k \rightleftharpoons P + E}{\dot{S} = 0.1 \cdot ES_1 - p_1 \cdot E \cdot S} \\
\dot{E} = 0.1 \cdot ES_1 - p_2 \cdot E \cdot S + 0.1 \cdot ES_k - p_2 \cdot E \cdot P \\
E\dot{S}_1 = 0.01 \cdot E \cdot S - p_3 \cdot ES_1 + 0.05 \cdot ES_2 \\
\vdots \\
E\dot{S}_k = 0.1 \cdot ES_{k-1} - p_k \cdot ES_k + 0.01 \cdot E \cdot P \\
\dot{P} = 0.1 \cdot ES_k - p_{k+1} \cdot E \cdot P - 0.1 \cdot P \\
\frac{p_1 = 0.01, p_2 = 0.01, p_3 = 0.2,}{p_k = 0.15, p_{k+1} = 0.01}
\end{array}$$

Figure 7.1: ODE model of reversible catalytic reaction

uated model properties. Note that the number of parameters can change the structure of the transition system, which may result in shorter run times, as demonstrated in table 7.1. Also note that in several occasions, especially when comparing single and dual machine experiments, the algorithm exhibits a super-linear speed-up. There are two main reasons of this behaviour.

First is the merge buffer as described in section 5.4, which can merge several traversals into one and therefore reduce expected overhead. Second is the garbage collector used by the Java Runtime Environment. Single machine experiments generally take up more memory, since the whole state space has to be stored on one computer. This leads to more aggressive garbage collecting and thus results in higher run times. On the other hand, when using high number of computational nodes, the garbage collection has only insignificant impact thanks to the high amount of available memory.

7. EXPERIMENTAL EVALUATION

# of params	1	2	3	4	5	6
# of nodes						
1	6456	N/A	N/A	N/A	N/A	N/A
2	2610	6179	5089	N/A	N/A	N/A
3	1696	4022	3661	3784	N/A	N/A
4	854	1759	2285	2454	N/A	N/A
5	683	1371	1365	1580	1736	N/A
6	499	1095	1019	1254	1609	1350
7	435	861	796	1023	1406	1340
8	292	642	650	853	1134	1118
9	258	439	630	752	983	962
10	232	418	557	637	839	822
11	198	347	516	553	784	679
12	177	329	420	562	759	660

Table 7.1: The runtime in seconds for the model with 6 variables, 13 thresholds and different number of unknown parameters.

# of variables	4	5	6	7
# of nodes				
1	3.3	22	794	N/A
2	3.3	12	489	N/A
3	3.5	9.5	253	N/A
4	3.4	8.2	185	8571
5	2.7	8	112	6608
6	2.4	7.2	101	5291
7	2.7	7.3	77	3024
8	2.3	6.6	64	2630
9	2.3	6.3	55	2366
10	2.5	6.8	52	2081
11	2.7	6.2	47	1999
12	2.2	5.6	41	1828

Table 7.2: The runtime in seconds for the model with 1 unknown parameter, 11 thresholds per the variable and different number of variables.

# of thresholds	10	11	12	13
# of nodes				
1	274	794	1805	6456
2	196	489	1252	2610
3	84	253	570	1696
4	62	185	408	854
5	51	112	280	683
6	40	101	226	499
7	33	77	192	435
8	29	64	161	292
9	26	55	145	258
10	24	52	108	232
11	23	47	94	198
12	22	41	96	177

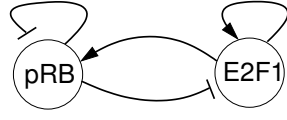
Table 7.3: The runtime in seconds for the model with 1 unknown parameter, 6 variables and different number of thresholds per variable.

7.2 Case study

To demonstrate the applicability of the algorithm, we investigate a well-known ODE model [15] which represents a two-gene regulatory network describing interaction of the tumor suppressor protein *pRB* and the central transcription factor *E2F1* (see Figure 7.2 (left)). This network represents the crucial mechanism governing the transition from G_1 to S phase in the mammalian cell cycle. In the G_1 -phase the cell makes an important decision. In high concentration levels, *E2F1* activates the G_1/S transition mechanism. In low concentration of *E2F1*, committing to S -phase is refused and that way the cell avoids DNA replication.

This mechanism is an example of *bistable switch*. The system makes an irreversible decision to finally reach some of the two stable states. The model is represented by two differential equations as seen in Figure 7.2. In order to discretize this model, a piece-wise multi-affine approximation (PMA) has to be computed [16]. In our experiments, we use 70 thresholds per each variable during this process. However, no significant change in results has been observed for

7. EXPERIMENTAL EVALUATION



$$\begin{aligned}\frac{d[pRB]}{dt} &= k_1 \frac{[E2F1]}{K_{m1} + [E2F1]} \frac{J_{11}}{J_{11} + [pRB]} - \phi_{pRB}[pRB] \\ \frac{d[E2F1]}{dt} &= k_p + k_2 \frac{a^2 + [E2F1]^2}{K_{m2}^2 + [E2F1]^2} \frac{J_{12}}{J_{12} + [pRB]} - \phi_{E2F1}[E2F1]\end{aligned}$$

$$a = 0.04, k_1 = 1, k_2 = 1.6, k_p = 0.05, \phi_{E2F1} = 0.1 \\ J_{11} = 0.5, J_{12} = 5, K_{m1} = 0.5, K_{m2} = 4$$

Figure 7.2: G_1/S transition regulatory network and its ODE model taken from [15].

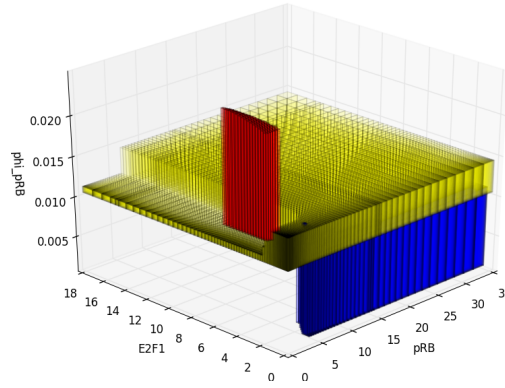


Figure 7.3: Coloured model checking results. Red and blue parts correspond to the high and low stable regions, respectively. Yellow part displays the states where the overall *bistable switch* formula φ holds.

higher threshold counts.

On the resulting model, we can perform a coloured model checking of the formula $\varphi \equiv \text{EFAG}(\text{high}) \wedge \text{EFAG}(\text{low})$ considering the initial parameter space $\phi_{pRB} \in [0.001, 0.025]$. The atomic propositions *low* and *high* characterise the location of expected regions of *E2F1* stability. Based on the results reported in [15] we define the stable regions as *high* $\equiv (E2F1 > 4 \wedge E2F1 < 7.5)$ and *low* $\equiv (E2F1 > 0.5 \wedge E2F1 < 2.5)$.

As seen in figure 7.3, the computation has successfully discovered both stable areas and corresponding parameters, as well as area and parameters from which both stable states can be reached. However, note that due to over-approximation, only the stable areas are guaranteed to be true.

8 Conclusion

The aim of this work was to design and implement a method that can efficiently solve parameter synthesis problem for biochemical models and CTL hypothesis.

We proposed an algorithm based on existing algorithms for distributed CTL model checking and coloured model checking of LTL formulas. The algorithm is given as pseudocode and we also provide a working implementation. The algorithm is not limited to the biochemical models and therefore has general applicability.

We also proposed an efficient state space partitioning heuristic for ODE models and a merge buffer heuristic for reduction of unnecessary state space traversals.

The provided implementation is oriented at biochemical models, provides support for two distinct modelling techniques, full range of CTL operators and features easily extensible modular architecture.

We benchmarked the provided implementation and showed that the algorithm scales well with higher number of computation nodes and provides sufficient performance to investigate models with large state spaces, even though the exponential state space explosion still remains a problem. The soundness of the algorithm and applicability has been demonstrated on a biologically relevant model of bistable switch.

In the future, the implementation can be optimized to provide even better performance and scalability, especially on multicore architectures. Also support for different modelling techniques and partitioning heuristics can be implemented.

Another possible direction of development would be a coloured symbolic model checking, which should provide even better performance compared to our explicit method, especially in terms of memory usage.

Bibliography

- [1] E. M. Clarke, Jr., O. Grumberg, and D. A. Peled, *Model Checking*. Cambridge, MA, USA: MIT Press, 1999.
- [2] J. Barnat, L. Brim, A. Krejci, A. Streck, D. Safranek, M. Vejnar, and T. Vejpustek, "On parameter synthesis by parallel model checking," *IEEE/ACM Trans. Comput. Biol. Bioinformatics*, vol. 9, pp. 693–705, May 2012.
- [3] G. Batt, M. Page, I. Cantone, G. Gössler, P. Monteiro, and H. de Jong, "Efficient parameter search for qualitative models of regulatory networks using symbolic model checking," *Bioinformatics*, vol. 26, no. 18, pp. 603–610, 2010.
- [4] R. Donaldson and D. Gilbert, "A model checking approach to the parameter estimation of biochemical pathways," in *CMSB*, vol. 5307 of *LNCIS*, pp. 269–287, Springer, 2008.
- [5] A. Donzé, G. Clermont, and C. J. Langmead, "Parameter synthesis in nonlinear dynamical systems: Application to systems biology," *Journal of Computational Biology*, vol. 17, no. 3, pp. 325–336, 2010.
- [6] S. K. Jha and C. J. Langmead, "Synthesis and infeasibility analysis for stochastic models of biochemical systems using statistical model checking and abstraction refinement," *Theoretical Computer Science*, vol. 412, no. 21, pp. 2162 – 2187, 2011.
- [7] L. Brim, K. Yorav, and J. Zidkova, "Assumption-based distribution of CTL model checking," *STTT*, vol. 7, no. 1, pp. 61–73, 2005.
- [8] P. Collins, L. C. Habets, J. H. van Schuppen, I. Černá, J. Fabriková, and D. Šafránek, "Abstraction of biochemical reaction systems on polytopes," in *IFAC World Congress*, pp. 14869–14875, IFAC, 2011.
- [9] R. Thomas, "Regulatory networks seen as asynchronous automata: A logical description," *Journal of Theoretical Biology*, vol. 153, no. 1, pp. 1 – 23, 1991.

- [10] F. Lerda and R. Sisto, "Distributed-memory model checking with spin," in *Theoretical and Practical Aspects of SPIN Model Checking*, pp. 22–39, Springer, 1999.
- [11] W. Feijen and A. van Gasteren, "Shmuel safra's termination detection algorithm," in *On a Method of Multiprogramming*, Monographs in Computer Science, pp. 313–332, Springer New York, 1999.
- [12] I. Github, "sybila/distributed-ctl-model-checker," 2015.
- [13] J. Barnat, L. Brim, I. Cerná, S. Drazan, J. Fabriková, J. Láník, D. Šafránek, and H. Ma, "Biodivine: A framework for parallel analysis of biological models," in *Proceedings Second International Workshop on Computational Models for Cell Processes, COMPMOD 2009, Eindhoven, the Netherlands, November 3, 2009.*, pp. 31–45, 2009.
- [14] A. Streck, D. Šafránek, L. Brim, *et al.*, "Parsybone: Parameter synthesizer for boolean networks," 2012.
- [15] M. Swat, A. Kel, and H. Herzog, "Bifurcation analysis of the regulatory modules of the mammalian G1/S transition," *Bioinformatics*, vol. 20, no. 10, pp. 1506–1511, 2004.
- [16] R. Grosu, G. Batt, F. H. Fenton, J. Glimm, C. L. Guernic, S. A. Smolka, and E. Bartocci, "From cardiac cells to genetic regulatory networks," in *CAV*, vol. 6806 of *LNCS*, pp. 396–411, 2011.