

**CONTROL-ORIENTED MODELING OF DISCRETE
CONFIGURATION MOLECULAR SCALE PROCESSES:
APPLICATIONS IN POLYMER SYNTHESIS AND THIN
FILM GROWTH**

A Thesis
Presented to
The Academic Faculty

by

Cihan Oguz

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
School of Chemical & Biomolecular Engineering

Georgia Institute of Technology
December 2007

Copyright © 2007 by Cihan Oguz

**CONTROL-ORIENTED MODELING OF DISCRETE
CONFIGURATION MOLECULAR SCALE PROCESSES:
APPLICATIONS IN POLYMER SYNTHESIS AND THIN
FILM GROWTH**

Approved by:

Professor D. W. Hess,
Committee Chair
School of Chemical and Biomolecular
Engineering
Georgia Institute of Technology

Professor Martha A. Gallivan, Advisor
School of Chemical & Biomolecular
Engineering
Georgia Institute of Technology

Professor J. H. Lee
School of Chemical and Biomolecular
Engineering
Georgia Institute of Technology

Professor M. Li
Materials Science and Engineering
Georgia Institute of Technology

Professor P. J. Ludovice
School of Chemical and Biomolecular
Engineering
Georgia Institute of Technology

Date Approved: 26 October 2007

To my grandfather

Tahir Ozgoren (1913-2006)

ACKNOWLEDGEMENTS

First and foremost, I would like to thank my advisor Professor Martha Grover Gallivan for her continuous support, and guidance in the last five years. She has set a great example for me to follow, in my academic and personal life. I would also like to acknowledge the members of my thesis committee for their help and advice: Professor Ludovice, Professor Lee, and Professor Hess from my department, and Professor Li from Materials Science and Engineering.

I thank my coworkers in the group: Andres Felipe Hernandez Moreno, Jonathan Rawlston, Paul Wissmann, and Rentian Xiong for their support, kindness, and enthusiasm for research. Their warm personalities, and professionalism has created a great working environment for me.

My friends in Georgia Tech have been the greatest gifts through this long journey. I would especially like to thank Gozde and Cerag for their incredible support when I was going through the most difficult times of my life. I can never forget Gozde's taking care of me after my recent surgery. I also thank Ugur, Mehmet, Ozgul, Burcak, and Ulas for the great times I had outside the school. Special thanks to my office mates Yeny and Chris for being such good listeners and to my former roommate Tibor for stopping by Atlanta every year.

I can never thank my parents enough for providing me the best education possible, and teaching me values I need through my whole life. I would like to thank my whole family for their encouragement and support. It was an honor to have a research collaboration with my aunt Emel Yilgor, and my uncle Iskender Yilgor. And finally, I would like to dedicate this work to my grandfather Tahir Ozgoren, who recently passed away.

TABLE OF CONTENTS

| | |
|---|------|
| DEDICATION | iii |
| ACKNOWLEDGEMENTS | iv |
| LIST OF TABLES | vii |
| LIST OF FIGURES | viii |
| SUMMARY | xii |
| I INTRODUCTION | 1 |
| 1.1 Motivation | 2 |
| 1.2 Outline of the Thesis | 4 |
| II PRELIMINARIES | 8 |
| 2.1 KMC simulations and their challenges | 8 |
| 2.1.1 The general KMC algorithm | 9 |
| 2.1.2 Challenges in KMC simulations | 10 |
| 2.2 Hyperbranched polymerization | 12 |
| 2.2.1 Past work in the molecular modeling of hyperbranched polymers | 13 |
| 2.3 Thin film deposition | 16 |
| 2.3.1 Past work in the molecular modeling of thin film deposition | 16 |
| III MODELING AND EXPERIMENTAL VALIDATION OF STRUCTURE DEVELOPMENT IN HYPERBRANCHED POLYMERS | 22 |
| 3.1 Synthesis of hyperbranched polyurethaneureas in solution | 24 |
| 3.1.1 Experimental Procedure | 24 |
| 3.1.2 Simulation Procedure | 24 |
| 3.1.3 Results | 28 |
| 3.2 Synthesis of hyperbranched polyetheresters in melt | 33 |
| 3.2.1 Experimental Procedure | 34 |
| 3.2.2 Simulation Procedure | 34 |

| | | |
|-------|--|-----|
| 3.2.3 | Results | 37 |
| 3.3 | State space modeling of hyperbranched polymers | 48 |
| 3.4 | Conclusions | 53 |
| IV | MODEL REDUCTION OF STOCHASTIC MOLECULAR SIMULATIONS | 56 |
| 4.1 | Characterization of the state space | 58 |
| 4.2 | Reduction of the state dimension | 62 |
| 4.3 | Discretization of state space | 64 |
| 4.3.1 | Results | 66 |
| 4.4 | Model Identification | 74 |
| 4.4.1 | Simple Cell Mapping | 74 |
| 4.4.2 | K-nearest neighbor algorithm | 78 |
| 4.5 | Conclusions | 79 |
| V | PERFORMANCE EVALUATION OF THE DYNAMIC MODEL | 80 |
| 5.1 | Local Error Quantification | 80 |
| 5.2 | Global Error Quantification | 83 |
| 5.3 | Optimization of the final film structure and the deposition time | 92 |
| 5.4 | Conclusions | 96 |
| VI | EVALUATION OF THE MODEL REDUCTION PARAMETERS | 98 |
| 6.1 | Effects of the state dimension | 98 |
| 6.2 | Effects of map size on the model predictions | 100 |
| 6.3 | Effects of state space exploration | 109 |
| 6.4 | Conclusions | 118 |
| VII | CONTRIBUTIONS AND FUTURE WORK | 119 |
| 7.1 | Contributions | 119 |
| 7.2 | Future work | 123 |
| | REFERENCES | 126 |
| | VITA | 137 |

LIST OF TABLES

| | | |
|---|---|-----|
| 1 | Statistics of three SOMs trained using different data sets. | 70 |
| 2 | Mean and standard deviation values of E_{SSC} for three test simulation sets. | 89 |
| 3 | Mean values of E_q , E_x , $E_{SSC'}$ and $E_{SSC''}$. E_q , which is computed for each data vector and the prototype vector of the data vector's best matching unit on SOM, is the difference between data vector and prototype vector during SOM training. The other errors are the normalized versions of E_{SSC} | 92 |
| 4 | Statistics of SOMs trained with differently sized data vectors. | 100 |
| 5 | Statistics of differently sized SOMs. | 101 |

LIST OF FIGURES

| | | |
|----|--|----|
| 1 | Two important modeling aspects related to molecular simulations. | 2 |
| 2 | Model inversion problem. | 6 |
| 3 | Different synthesis methods using $A_2 + B_3$ polymerization. | 23 |
| 4 | Chemical structures of monomeric and oligomeric A_2 and B_3 | 25 |
| 5 | Cyclization versus non-cyclization reactions. | 26 |
| 6 | Dendritic, linear and terminal B_3 units. | 28 |
| 7 | Comparison of experimental and simulation results on the development of average polymer molecular weight as a function of A_2 addition and cyclization ratio for a $1000 \times 1000 A_2 \times B_3$ system. (a) Number average molecular weight and (b) weight average molecular weight. Experimental data: (o) experiment with 25% solids by weight, and (*) 10% solids by weight. Simulation data: $\gamma = 0$ (solid line), $\gamma = 0.01$ (dashed line), $\gamma = 0.1$ (dotted line), $\gamma = 1$ (dash-dotted line). | 30 |
| 8 | Kinetic Monte Carlo simulations on the polydispersity index (PDI) (a) and the degree of branching (DB) (b) as a function of A_2 addition. $\gamma = 0$ (solid line), $\gamma = 0.01$ (dashed line), $\gamma = 0.1$ (dotted line), $\gamma = 1$ (dash-dotted line). | 31 |
| 9 | Number of cyclization events per molecule, as predicted by the kinetic Monte-Carlo simulations as a function of A_2 addition, at different levels of cyclization. $\gamma = 0.01$ (dashed line), $\gamma = 0.1$ (dotted line), $\gamma = 1$ (dash-dotted line). | 33 |
| 10 | Comparison of simulation and experiment [117] (*). In the simulations, the cyclization ratio γ is varied. $\gamma = 0$ (solid line), $\gamma = 10^{-3}$ (dashed line), $\gamma = 10^{-2}$ (dotted line), $\gamma = 10^{-1}$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D . Agreement between experiments and simulations is not achieved for f_D at any value of γ | 39 |
| 11 | Simulation predictions of extent of cyclization. In the simulations, the cyclization ratio γ is varied: $\gamma = 10^{-3}$ (dashed line), $\gamma = 10^{-2}$ (dotted line), $\gamma = 10^{-1}$ (dash-dotted line). For the $\gamma = 0$ case (solid line), $EOC=0$ | 40 |

| | | |
|----|--|----|
| 12 | Comparison of simulation and experiment [117] (*). In the simulations, $\rho = k_1/k_2 = k_2/k_3$. $\rho = 1$ (solid line), $\rho = 1.5$ (dashed line), $\rho = 2$ (dotted line), $\rho = 10$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D | 42 |
| 13 | Comparison of simulation and experiment [117] (*). In the simulations, $\rho_{12} = k_1/k_2$ and $k_2 = k_3$. $\rho_{12} = 1$ (solid line), $\rho_{12} = 1.5$ (dashed line), $\rho_{12} = 2$ (dotted line), $\rho_{12} = 10$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D | 43 |
| 14 | Simulated evolution with end-capping reagents added at the beginning of the process, with $A_2 : B_3 : E = 1 : 1 : 1$. The molecular weight of E (PPG-M-1000) is 1200 g/mol. As in Figure 13, $\rho_{12} = k_1/k_2$ and $k_2 = k_3$. $\rho_{12} = 1$ (solid line), $\rho_{12} = 1.5$ (dashed line), $\rho_{12} = 2$ (dotted line), $\rho_{12} = 10$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D . (Note: dendritic units are calculated here based on the number of A - B reactions. E - B reactions are not considered in the calculation since they do not lead to further branching. | 45 |
| 15 | Simulated evolution with end-capping reagents added at the beginning of the process, with variation in stoichiometry: $A_2 : B_3 : E = 1 : 1 : 1$ (solid line) and $A_2 : B_3 : E = 1 : 0.9 : 1$ (dotted line with markers). The molecular weight of E (PPG-M-1000) is 1200 g/mol. $\rho_{12} = k_1/k_2 = 1.5$ and $k_2 = k_3$. (a) Weight-average molecular weight M_w (b) Polydispersity index PDI (c) Fraction of dendritic units f_D | 46 |
| 16 | Two isomeric polymers with the same molecular weight and different branching. | 49 |
| 17 | Types of events that take place during the MBE deposition of GaAs. | 57 |
| 18 | Schematic of the modeling approach. | 57 |
| 19 | Relationship between the crystallographic directions and two dimensional coordinates used for constructing the lattice. Dashed region is the $\beta 2(2 \times 4)$ reconstruction of GaAs. Dark and white atoms are As and Ga atoms, respectively. The atomic radius decreases with increased depth. | 60 |
| 20 | Normalized eigenvalues versus principal components. | 67 |
| 21 | Comparison of the reconstructions of $SSC_{up,down,i}$ with the original data using 2 and 5 modes. The region enclosed with the gray dashed line represents the surface structures with a size of less than 20 lattice units. | 68 |

| | | |
|----|---|----|
| 22 | U-matrix for a 3×1 SOM. | 71 |
| 23 | Values of film surface coverage (in monolayers) for each node of SOM3. | 72 |
| 24 | Trajectory of a training simulation on SOM3. Ga flux is kept constant at the minimum flux (0.06 ML/s) during this simulation. | 72 |
| 25 | Trajectory of a training simulation on SOM3. Ga flux is kept constant at the maximum flux (0.20 ML/s) during this simulation. | 73 |
| 26 | SOM1 (generated using the Training Simulation Set 1) and its U-matrix. | 74 |
| 27 | SOM3 (generated using the Training Simulation Sets 1, 2 and 3) and its U-matrix. | 75 |
| 28 | Schematic of simple cell mapping. | 76 |
| 29 | Computation of the cell mapping error. | 81 |
| 30 | Computation of the distribution of the cell mapping error. | 82 |
| 31 | Cumulative distribution function of the cell mapping error. | 82 |
| 32 | Flux profile of a test simulation. | 83 |
| 33 | Trajectories of the KMC test simulation (red line) and the prediction (black line) on SOM3. | 84 |
| 34 | Reconstruction of the SSC function with the prototype vector of the SOM Node 186 and the original KMC simulation data at final film coverage. | 85 |
| 35 | Trajectories of a training simulation (red line) and the prediction (black line) on SOM3. Ga flux is kept constant at 0.06 ML/s during this simulation. | 86 |
| 36 | Trajectories of a training simulation (red line) and the prediction (black line) on SOM3. Ga flux is kept constant at 0.20 ML/s during this simulation. | 86 |
| 37 | The evolution of E_{SSC} for the test simulation with $k=1$ and $k=2$ | 88 |
| 38 | The mean value of E_{SSC} at different film coverage levels for three test simulation sets and the training data. | 89 |
| 39 | Cumulative distribution function of the E_{SSC} for three test simulation sets. | 90 |
| 40 | Cumulative distribution function of different types of error. | 92 |
| 41 | A portion of the optimal surface structure. The initial surface structure has regular trenches (dark areas), and as the deposition is performed, clusters form in and on top of the trenches. | 94 |

| | | |
|----|--|-----|
| 42 | Optimal flux profile computed by the dynamic model. | 95 |
| 43 | Trajectories of the KMC simulation (red line) with the optimal flux profile and the prediction (black line) on SOM3. | 96 |
| 44 | Reconstructions of the SSC function with the prototype vectors of the SOM nodes 182, 184 and the original simulation data. | 97 |
| 45 | Periodic flux profile in a test simulation. | 104 |
| 46 | Trajectories of a test KMC simulation (red line) and the prediction (black line) on SOM3. | 104 |
| 47 | Trajectories of a test KMC simulation (red line) and the prediction (black line) on SOM4. | 105 |
| 48 | The value of the relative quantization error E_{rq} at different film coverage levels for the test simulation with the flux profile given in Figure 45. | 106 |
| 49 | A portion of $SSC_{up,down,i}$ from the KMC simulation data at 0.10 ML film coverage, and its reconstructed form using the prototype vectors of its best matching units on SOM3 and SOM4. | 106 |
| 50 | The value of the state prediction error E_x at different film coverage levels for the test simulation with the flux profile given in 45. | 107 |
| 51 | A portion of $SSC_{up,down,i}$ coming from the KMC simulation data at 0.14 ML film coverage and the prediction of the low and high dimensional models. | 108 |
| 52 | Cumulative distribution functions of the prediction error for 1210 test simulations with the low dimensional model, and the 1024 exploration simulations with the low and high dimensional models. | 110 |
| 53 | Cumulative distribution functions of the prediction error for 1210 test simulations, and the 568 exploration simulations with the small model at a film coverage of 0.08 ML. | 116 |
| 54 | Cumulative distributions function of the prediction error for 1210 test simulations, and the 240 exploration simulations with the small model at a film coverage of 0.10 ML. | 117 |

SUMMARY

Integration of accurate, and detailed molecular simulation models with system tools is a growing research area, because molecular simulations enable design and optimization of processes at a molecular scale. This opens up possibilities to control quality of products used in many applications, such as microelectronic devices. Such a model-based approach is already common for macroscopic processes, but not for microscopic or molecular models and structures. That is because there are two major issues that need to be resolved for the integration of molecular simulations and systems tools (e.g. dynamic analysis and optimization techniques) to take place.

The first issue is the development of accurate molecular simulation models for processing, using experimental data. In this thesis, we focus on stochastic molecular simulations since continuum models can not successfully capture the dynamics associated with interactions between molecules, and molecular dynamics simulations can not predict over long processing time scales. We bring forward a stochastic model of a process used in applications such as printer inks, and automotive coatings: synthesis of hyperbranched polymers using difunctional A_2 oligomers, and trifunctional B_3 monomers. Currently, the kinetics of this process, as well as the process conditions affecting the polymer structure and properties, are not clearly understood. There are many possible different reaction mechanisms that can influence the polymer structure development in this system. Even though some of these mechanisms, such as cyclization reactions, have been taken into account in the past, a comprehensive simulation study for hyperbranched polymers does not exist. Also, previous approaches in modeling hyperbranched polymers lack close integration with experimental data. Because of these reasons, we developed a kinetic Monte Carlo (KMC) model, which

takes into account a wide range of reaction mechanisms, using experimental data. By implementing our simulation model, we compare the effects of different synthesis routes on the properties of polymers, such as average molecular weight, and the degree of branching. These synthesis routes involve melt polymerization (with no solvent), and solution polymerization. In melt phase, we consider the influence of cyclization and endcapping reactions, and unequal reactivity of different monomer units on the polymer structure. On the other hand, in solution polymerization, we consider the effects of monomer concentration. These models are an important step toward rational process design to achieve desired polymer structures.

The second major issue in the integration of molecular simulations with dynamic analysis and optimization tools is the development of reduced order models from high dimensional stochastic simulation models. These molecular simulations possess a fine level of description of the process physics compared to the macroscopic (continuum) models. However, due to their high computational cost, it is not feasible to employ them for optimization and control tasks. Hence, reduced order models are needed for design purposes. Current model reduction algorithms in the literature contain too many assumptions about system dynamics to alleviate the computational burden associated with the reduction. These assumptions include large separation between the time scales of low and high order moments describing the system state and existence of steady state operating points for linearization purposes. In this thesis, we present a novel model reduction algorithm that can potentially be applied to any nonlinear dynamic system even when these assumptions are not valid. Epitaxial deposition of gallium arsenide (GaAs), used in manufacturing high performance microelectronic devices, is used to illustrate our model reduction algorithm. First, the state space of the simulation model is characterized by running a limited set of simulations with various material flux profiles, and characterizing surface snapshots during these simulations. The characterization of the dynamic state is performed by using a correlation function

describing the correlation between the steps on film surfaces. The rationale is that only a small fraction of all possible film surface configurations are accessible in this process, and a limited set of simulations can access these configurations. The model reduction algorithm consists of applying principal component analysis (PCA) to high dimensional simulation data for reduction of the state dimension, self organizing map (SOM) for grouping similar surface snapshots, and simple cell mapping (SCM) for identifying the transitions between different surface configuration groups. After identifying the model, its predictive ability is evaluated by a large set of simulations with highly dynamic material flux profiles. The computed distribution of the prediction error, which lacks in most model reduction studies for nonlinear systems described by molecular simulations, demonstrates the accuracy the model. As another part of the evaluation, we critically evaluate the effects of model reduction parameters, such as dimension of the state, number of surface configuration groups, and the quality of the data used to identify the model on its predictive ability under a wide range of process conditions. Dynamic process optimization, which is an important systems task, is also performed by using the reduced order model to compute the material flux profile (with minimum deposition time) that leads to the optimized film structure. For this optimization, the reduced order model provides 11 orders of magnitude reduction in the computational time, compared to the high dimensional molecular simulations. As a part of the future work, we also discuss how our model reduction scheme could be applied to the hyperbranched polymerization process, through a dynamic state description with monomer-monomer correlation functions. Our reduction approach is potentially applicable to any molecular simulation, as long as an appropriate state description is used. Challenges include the difficulty of fully validating the reduced order model. In order to achieve rationale material design using molecular simulations, close integration of experiments and modeling is critical.

CHAPTER I

INTRODUCTION

The objective of this thesis is to propose modeling techniques to enable the design and optimization of material systems which require descriptions via molecular simulations. These kinds of systems are quite common in materials and engineering research [123]. The first step in performing design and optimization tasks on such systems is the development of accurate simulation models from experimental data. In the first part of this thesis, we present a novel simulation model for the hyperbranched polymerization process of difunctional A_2 oligomers, and B_3 monomers. Unlike the previous models developed by other groups [23, 31, 107, 106], our model is able to simulate the evolution of the polymer structure development under a wide range of synthesis routes, and in the presence of cyclization and endcapping reactions. Furthermore, our results are in agreement with the experimental data, and add insight into the underlying kinetic mechanisms of this polymerization process.

The second major step in our work is the development of reduced order process models that are suitable for design and optimization tasks, using simulation data. We illustrate our approach on a stochastic simulation model of epitaxial thin film deposition process [54]. Compared to the widely used approach called equation-free modeling [36], our method requires fewer assumptions about the dynamic system. The assumptions required in equation-free modeling include a wide separation between the time scales of low and high order moments describing the system state, and the accuracy of the time derivatives of system properties computed from molecular simulation data, despite the potentially large amount of fluctuations in stochastic simulations [8]. Unlike the recent similar studies [8, 120], our study also includes the

analysis of prediction error which is important to evaluate the predictions of the reduced order model, compared to the high dimensional molecular simulations. Hence, we address two major issues in this thesis: development of simulation models from molecular experimental data, and derivation of reduced order models from molecular simulation data. These two aspects of modeling, which are illustrated in Figure 1 are both necessary to design and optimize processing conditions of materials for which continuum level descriptions are not available or accurate enough.

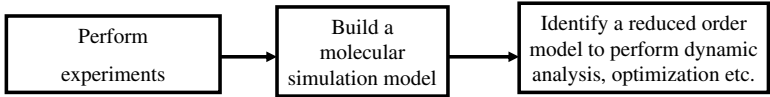


Figure 1: Two important modeling aspects related to molecular simulations.

1.1 *Motivation*

We are currently in an era where many materials design problems require models with explicit atomic scale interactions. An example is the design of integrated circuits. The material properties are functions of processing conditions such as temperature, and pressure within a reactor. Hence, process models that describe the relationship between inputs (processing conditions) and outputs (material properties) are needed for analyzing and designing material systems. When continuum assumptions are not valid at very small scales, developing accurate and realistic continuum models from material and energy balances closed through constitutive relations is not possible. Therefore, descriptions of processes at finer levels than macroscopic ones are required. As a result of this, and the growth in computational power over the last two decades, dynamics of material systems are commonly simulated by molecular simulations (e.g. molecular dynamics (MD) and kinetic Monte Carlo (KMC) simulations). For material systems which evolve in seconds or larger time scales, KMC simulations are widely used due to the inaccessibility of these time scales by MD simulations. Even though the general KMC algorithm is well established over the years, some major challenges

include the determination of microscopic events to take into account in the simulations (e.g. reaction, diffusion, and adsorption events), and the computation of the kinetic parameters of these events [123].

Hyperbranched polymers are materials with a large potential to be used in many applications due to their favorable properties (e.g. low viscosity and high solubility). However, a comprehensive KMC model for hyperbranched polymerization does not exist in the literature. In this thesis, we present such a simulation model by identifying the rates of several proposed reaction mechanisms in the polymerization of difunctional A_2 oligomers, and B_3 monomers. Compared to the existing models [23, 31, 107, 106], our model enables the evaluation of different synthesis routes by giving insights about the polymer structure development. Such an approach is necessary to design polymerization processes that are targeting specific polymer structure and properties in the future.

For material systems with dynamics that are appropriately described only by molecular simulations, design and optimization is usually not straightforward due to the high computational load of these simulations. This computational load originates from the disparity between the time and length scales of microscopic processes in the simulations and the evolutions of macroscopic material properties. Since these design and optimization problems are multiscale in nature, they can not be solved directly by employing molecular simulations. The common practice in the literature for solving these kinds of problems is the derivation of reduced order models from simulation data. These reduced order models provide a mapping between the process inputs and outputs (material properties). However, current approaches for model reduction and engineering design are only limited to highly simplified KMC simulation models (with 1 or 2 dimensional lattices, small system sizes, and just a few types of reaction/adsorption/diffusion events) [4, 32, 33, 86]. In past model reduction studies, a very low order representation of the system state is assumed (e.g. surface roughness

of a thin film [32, 33]) ignoring the effects of higher order moments (e.g. coordinates of each atom on the thin film surface that would give a certain surface roughness) on the material properties. In these studies, it is also assumed that the higher order moments become functionals of the low order ones to perform computations of the time derivatives of the system properties easily by using short simulations [8]. This assumption leads to a huge reduction in the number of possible states that the system can be in (dimension of the state space), which can be overly simplistic for non-trivial, realistic systems. The time derivatives obtained from short simulations are generally used to propagate the system dynamics for longer times by extrapolation (without running extra simulations) [36], hence reducing the computational load for dynamic optimization. However, due to the extensive amount of noise in stochastic molecular simulations, the derivative information coming from the short simulations may not be accurate, and can cause large errors when predicting the evolution of the system state [8]. Another deficiency of these studies is the lack of evaluation of the model reduction error, which is necessary to validate any model reduction approach. In order to address all these issues, we propose a novel model reduction algorithm for design and optimization that offers high dimensional state descriptions, time derivative free predictions, and a comprehensive approximation error analysis. This algorithm is illustrated on epitaxial deposition of gallium arsenide (GaAs) thin films, using a pair correlation function between steps on the film surface to describe the dynamic state of the system during the simulations.

1.2 Outline of the Thesis

The remaining chapters of this thesis are organized as follows: In Chapter 2, we give an overview of the KMC simulations, their general algorithm, and common challenges in their implementation. These challenges include the development of KMC models, their integration with control and optimization tools, and possible inefficiencies due

to the wide separation between the time scales of microscopic processes taking place during these simulations. Then, we describe two material systems that are used as case studies in this thesis: Hyperbranched polymers and epitaxially deposited thin films. These descriptions are followed by the past modeling work done to understand the dynamics associated with the processing in these two material systems.

In Chapter 3, we present a KMC simulation model for the hyperbranched polymerization of difunctional A_2 oligomers, and B_3 monomers. This model, whose kinetic parameters are derived from experimental data, is employed to study the effects of using different synthesis routes on the polymer structure development. These synthesis routes include melt polymerization of A_2 and B_3 , and solution polymerization with dropwise addition of A_2 into B_3 . We show that factors such as the extent of cyclization reactions, varying reactivity of the B groups, and the presence of monofunctional endcapping reactions can influence the polymer properties (e.g. molecular weight, and degree of branching) significantly.

Even though high order molecular simulations provide a description of the state evolution at a finer level than macroscopic (continuum) models, they do not provide a reduced order process model that can easily be inverted, which is necessary for process design and optimization. This is illustrated in Figure 2. We identify the main challenge associated with this approach as the limited measurements, which are necessary to validate the reduced order model. In addition, to build a reduced order model, one must describe the dynamic state of the system in the simulations. In the last part of Chapter 3, we discuss the ways dynamic state can be described in the hyperbranched polymerization simulations. Then, in Chapter 4, we present a novel model reduction algorithm that enables the use of molecular simulations for design and optimization. This algorithm generates a reduced order process model from high dimensional KMC simulation data. In order to illustrate our algorithm, we use an existing KMC model for the epitaxial deposition of gallium arsenide (GaAs). Unlike

the other model reduction studies in the literature, our study includes a comprehensive evaluation of the reduction error. Also, it does not rely on numerous assumptions, such as a very low order representation of the system state.

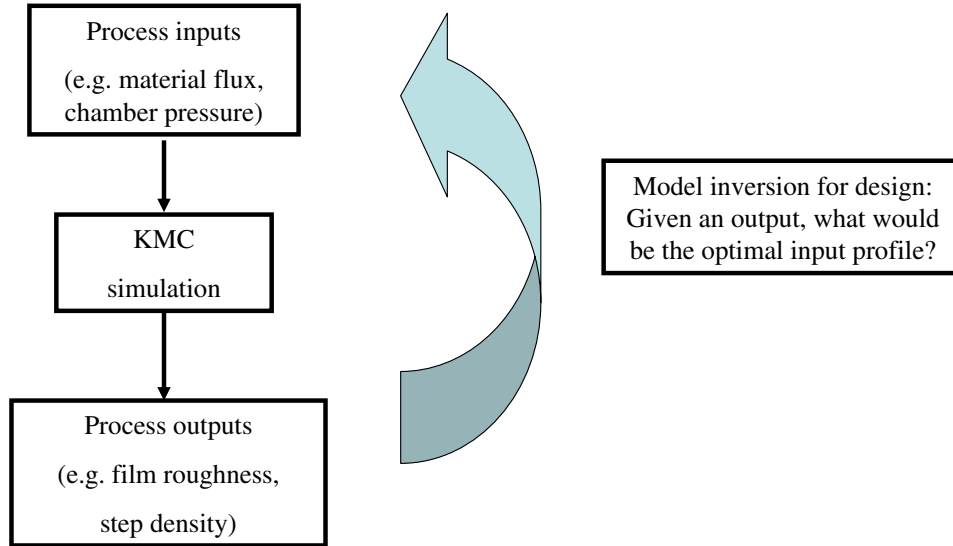


Figure 2: Model inversion problem.

In model reduction studies, it is important to derive the distribution of the error originating from the approximations associated with the order reduction (from high order KMC simulations to reduced order dynamic models). In Chapter 5, we address this issue by evaluating the predictive ability of our reduced order model. First, we run a large set of simulations with highly dynamic and random input profiles. Then, we compute the prediction error for these simulations, and the cumulative distribution function. Results show that the overall prediction error is quite low. In the last part of this chapter, we perform dynamic optimization using the reduced order model. The optimal flux profile that reaches the optimized film structure is computed by an approach motivated by dynamic programming. Here, the reduced order model provides a reduction of 11 orders of magnitude in the computational time, compared to the full KMC simulations.

In Chapter 6, we evaluate the effects of model reduction parameters such as the

dimension of the state, the number of configurations in the state space, and the quality of the training data that is used to build the dynamic model on the accuracy of the model predictions. Here, we derive two differently sized models with respect to the number of configurations in their state space and compare their performance. We identify that the dominant source of prediction error in the model is the quality of the training data. Finally, we present a method for improving the quality of the training data by exploring the state space of the simulations systematically.

CHAPTER II

PRELIMINARIES

In this chapter, we present a general kinetic Monte Carlo (KMC) simulation algorithm, and discuss some major challenges of KMC simulations. Then, we give an overview of the two material systems that are simulated using KMC methods: hyperbranched polymers and thin films. That is followed by a brief review of the work done in modeling these two classes of materials.

2.1 KMC simulations and their challenges

Monte Carlo (MC) simulations are most commonly used to compute equilibrium properties. However, Kinetic Monte Carlo (KMC) simulations are also able to describe evolution in time. The KMC algorithm was simultaneously proposed by Bortz et al. [10] and Gillespie [38, 37]. It is based on the assumption that the evolution of a system is the result of a series of discrete events. An event can be any microscopic process, such as adsorption (atom attaching to a surface), surface diffusion (atom hopping on a surface), desorption (atom leaving a surface), or a reaction between two molecules.

KMC simulations are an alternative to solving the master equation, which describes changes in a system with respect space and time. The master equation is expressed as:

$$\frac{dP_H(t)}{dt} = \sum_{H'} k^{H' \rightarrow H} P_{H'}(t) - \sum_{H'} k^{H \rightarrow H'} P_H(t) \quad (1)$$

where t and $P_H(t)$ are the time and the probability of the system being in configuration H . Here, the probability of a transition from configuration H' to H is $k^{H' \rightarrow H}$.

Due to the large number of possible system configurations, and the large dimension of the state vector (i.e. set of variables describing the state of a configuration) in many physical systems, it is not possible to solve the master equation analytically. For example, in a thin film deposition simulation with a substrate having 50×50 sites; and a maximum surface height of 10 for each site, the number of possible surface configurations is $10^{50 \times 50}$. In this case, the analytical solution of the master equation would require solving $10^{50 \times 50}$ coupled differential equations (one for each configuration). A KMC simulation simply provides a single stochastic realization of the master equation, and by averaging a large set of simulation results, one can obtain the expected system behavior under a specific set of conditions.

2.1.1 The general KMC algorithm

Let's consider a thin film deposition process with the events of adsorption, desorption, and diffusion. In this process, the rate of an adsorption event is the product of species flux and surface area, whereas the rates of a hopping or desorption event α is defined by $r_\alpha = v \times \exp(-B_\alpha/(k_B/T))$ where v is the attempt frequency, k_B is Boltzmann's constant, T is the temperature of the surface, and B_α is the energy barrier associated with event α . Following are the three steps of the KMC algorithm:

- The initial arrangement of atoms or molecules is specified and the possible transitions from this configuration are evaluated. Then, based on generated random numbers, one transition and its location are selected. The selection criterion is expressed as:

$$\frac{\sum_{j=1}^{r-1} k_j}{\sum_{j=1}^R k_j} \leq \zeta_1 \leq \frac{\sum_{j=1}^r k_j}{\sum_{j=1}^R k_j}$$

where k_i is the rate proportional to both the event rate r_i associated with the energy barrier β_i , and the number of sites at which event i can take place. In the above expression; r , which is the selected event number, can take any integer value from 1 to R where R is the total number of the possible events, and it depends on the current configuration of the system. The

random number ζ_1 , which is used for event selection, comes from a uniform distribution from 0 to 1.

- After this selection, based on a second random number ζ_2 , a site at which the event r takes place is selected. The selection criterion is similar to the previous one: $\frac{s-1}{S} \leq \zeta_2 \leq \frac{s}{S}$ where S is the total number of sites where event r can take place and s is the site number which is picked. After that, the event is executed at the selected site. It should be noted that, for a given event number, all sites have an equal probability for selection and no searching is needed. This results in saving computational time.
- Incrementing the simulation clock is the final step of the algorithm. The following formula with a third uniformly distributed random number ζ_3 between 0 and 1 is used: $\Delta t = \frac{-\ln \zeta_3}{\sum_{j=0}^R k_j}$.

These steps are repeated until a certain time (or surface coverage in thin film deposition) is reached during a simulation run.

2.1.2 Challenges in KMC simulations

Following are the major challenges in implementing KMC simulations:

- Determining the parameters in the simulations: In order to carry out KMC simulations, one needs a full list of the events that can take place in a simulation step, together with the rates of these events. This information can be extracted by first principles methods (e.g. molecular dynamics simulations or density functional theory calculations [70]), or alternatively from experimental data [54].

In order to simulate the epitaxial deposition of GaAs, we use an existing KMC model in the literature [54]. The kinetic parameters in this model (e.g. activation energies of diffusion, and desorption events) are functions of interactions

energies between Ga-Ga and As-As pairs on the surface. Derivation of these kinetic parameters is based on the inequalities of the interaction parameters, and their consistency with the STM images. In another study, Kratzer et al. used ab initio energy (or first principles) minimization calculations to derive the kinetic parameters of this process [70, 69]. We selected the first KMC model [54] for our simulations since it is applicable for a wider range of surface configurations. For the KMC model of our second process, hyperbranched polymerization, we propose several reaction mechanisms and determine the rate constants of reactions (e.g. cyclization and noncyclization reactions) from experimental data [117, 118], such as the polymer molecular weight.

- Large separation between time scales of events: In KMC simulations, there is usually a large separation between time scales of individual events. For example, in thin film deposition while a diffusion event happens very frequently (due to its low activation energy), adsorption is quite rare (e.g. millions of diffusion events for an adsorption event). Since we always aim for a certain surface coverage in these simulations, and adsorption is the only event that increases surface coverage, millions of diffusion events need to be executed in order to observe significant surface coverage change on the film surface. Also, after the execution of each event, events lists, which store the number of sites at which events can occur, need to be updated. All these factors lead to a high computational cost that needs to be dealt with efficient search (for event type and event site) and update algorithms. The high computational cost also prohibits the direct use of KMC simulations for design and optimization purposes. Therefore, reduced order process models that describe the mapping between process inputs and outputs need to be derived from a limited set of simulations.

In engineering, direct inversion of an algebraic model (Figure 2) is usually not

possible, which is why numerical methods are developed. They involve making predictions to efficiently search for the best solution (e.g. optimal operating conditions). But it is difficult to use them with KMC simulation models, because (1) each prediction (simulation) can take a long time, and (2) derivatives of system variables are often needed. However, derivatives are not directly available in KMC simulations, and are difficult to approximate numerically due to noise in these simulations. In order to address this problem, optimization methods such as Nelder-Mead and Hooke-Jeeves algorithms, which do not require derivatives, are used in past studies [4, 8]. These methods only use the evaluation of the objective functions (using the KMC simulation) for optimization instead of the derivatives. Another alternative to these deterministic optimization methods is simulated annealing [64], which is a stochastic method that is widely used in process engineering [123]. One example is a study by Raimondeau et al., where the authors used simulated annealing to determine the parameters of a reduced order model obtained from the KMC simulations of catalytic oxidation reaction of CO on platinum [102].

2.2 Hyperbranched polymerization

The shape and topology of organic molecules has a profound effect on their properties [2]. In the last two decades, synthetic polymer chemists have introduced a new class of highly branched macromolecules which are composed of multifunctional monomers, and are classified as either dendrimers or hyperbranched polymers. These versatile materials display several unique properties when compared with their linear analogs. Their favorable properties include low solution and melt viscosities, high solubilities and the presence of very large number of functional end groups that offers the possibility for further modification for various specialty applications. Dendrimers are typically synthesized using multi-step reactions and they offer superior control of

molecular size, shape, and functionality. On the other hand, hyperbranched polymers are less ordered but are easier to synthesize.

A number of excellent reviews, which describe the synthetic methodologies for the preparation of a wide variety of hyperbranched and dendritic polymeric systems through condensation, addition or ring-opening reactions, are available [35, 53, 57, 60, 124, 131]. These include polyesters [47, 75, 113, 114], polyamides [44, 59, 92], poly(ester-amides) [119], polyphenylenes [63, 62], poly(ether sulfones) [17, 83, 87], poly(etherketones) [46, 91], polyethers [109, 112], polyurethanes [13, 27, 34, 76, 77, 108], polyureas [34, 76], polycarbonates [9] and others [35, 53, 57, 60, 124, 131].

2.2.1 Past work in the molecular modeling of hyperbranched polymers

Modeling studies of polymer networks have long been used to explain experimental observations in an adhoc fashion, and modeling is often used to steer the discovery of synthetic methods and the formation of novel architectures. Early work of Flory [28] and Stockmayer [110, 111], on the step growth of multifunctional monomers was based on the assumption that there was no cycle formation in these polymerization processes, which enabled the calculation of molecular weight using an infinite series solution. The models of Flory and Stockmayer are useful because they enable quantitative predictions of polymer properties. While it is intuitive that a branched polymer that is composed of A_2 and B_3 monomers will ultimately gel, the model generates the quantitative prediction of gelation at 87% A conversion for the stoichiometry of $A_2 : B_3 = 1 : 1$ [111]. More recently, this method has been extended to include the effect of cycle forming reactions on the gel point [14].

In order to predict the time-evolution of polymerization, kinetic models based on mass-action kinetics are often employed. Ordinary differential equations are used to describe the concentration of each type of branching unit [101, 51]. For example, a trifunctional monomer may exist in four states: no reactions (free), one reaction

(terminal), two reactions (linear), and fully reacted (dendritic). These models are typically of low dimension and may be solved analytically in some cases, or numerically in others. A disadvantage of this modeling approach is that no information is provided on the molecular weight distribution of the resulting polymer.

When prediction of the molecular weight and its distribution are desired, population balance models are commonly used, in which the concentration of polymers of each possible size is computed [22, 23, 43, 101, 132]. Therefore, these models have high dimension. For systems of linear polymers composed of a single monomer type, the kinetics is fully described by the number of monomers in the polymer. Thus, if one models the concentration of polymers up to a maximum size of 1000 monomeric units, then the dimension of the model is also 1000. However, in branched polymers, more information is needed to describe the kinetics, such as the number of reactive end groups in the polymer. When several descriptors are needed to describe each polymer, the dimension of the population balance model grows rapidly [22]. For the model considered in this thesis, six descriptors are used to describe each polymer: the number of A_2 monomers, the number of B_3 monomers, the number of unreacted A groups in the polymer, the number and type of unreacted B groups (linear or terminal), and the number of end-capping agents. Solving a population balance model with six descriptors would be very intensive computationally.

Generating functions and the method of moments are often used to reduce the dimension of the population balance models. However, this approach becomes much more difficult as the number of descriptors grows, since moments must be included for each descriptor. The difficulty of this approach is illustrated by the recent paper of Dusek, Duskova-Srmckova, and Voit [22], in which unequal reactivities and monofunctional reagents were considered separately, but not simultaneously.

Cycle formation was neglected by Flory and Stockmayer, but it is a major factor in the structural development of dendritic and hyperbranched polymers [23, 42, 43,

72, 73]. The method of moments has been extended to describe cycle formation in hyperbranched polymers composed of AB_2 monomers [23] and the $AB_2 + B_3$ system [31]. This extension is enabled because in these systems, a maximum of one cycle is possible, at which point the polymer cannot grow further. In contrast, in the $A_2 + B_3$ system, there is no limit on the maximum number of cycles. Infinite series solutions have been developed to predict the gel point in the $A_2 + B_3$ system, but do not include additional effects such as unequal reactivities of the groups in B_3 , and the effect of monofunctional reagents [14].

As an alternative to mass action kinetic models that use concentration variables, Monte Carlo simulations have been performed for hyperbranched polymers, so that more realistic kinetics can be included and the structural information can be obtained. For example, the Wiener index can be computed for each polymer in the simulation, which is related to viscosity [127]. In the Monte Carlo simulations, individual monomers are reacted with each other to build up the polymer, using random numbers to select each event. Each monomer is tracked throughout the simulation, and information regarding its connection to other monomers is stored. Monte Carlo simulations may be used to describe only the connectivity of the polymers, without describing their spatial positions [23, 31, 107, 118] or lattice Monte Carlo simulations can be performed in which each monomer is associated with a spatial position in the lattice [15]. Direct comparison between experimental data and Monte Carlo simulations has been limited to date, but these simulations can be valuable in interpreting experimental results [23, 118]. As computational resources grow, Monte Carlo simulations become an increasingly attractive alternative to population balance modeling, since their major drawback has been the amount of computation required.

Another modeling approach is to describe the structural development of hyperbranched polymers using atomistically detailed algorithms [2, 127]. These studies

are limited to the growth of single polymers at a time because of the high computational cost, so they are not as useful for predicting the molecular weight distribution. However, if the conformation of the polymer changes throughout the reaction and this strongly influences the kinetics, then it may be necessary to include this level of detail.

2.3 Thin film deposition

Thin film deposition, which involves the deposition of a material (precursor) onto a substrate, is a critical step in manufacturing integrated circuits and MEMS devices. The thickness of a thin film can vary from tens of microns to single atomic layers depending on the deposition process, and the application in which the thin film is used.

In a thin film deposition process, the precursor can be in the gas-phase (e.g. chemical vapor deposition), or liquid phase (e.g. chemical deposition with plating). The production of the precursor can be achieved by the ionization of the atoms of a metal (e.g. in sputtering), or by gas-phase reactions (e.g. in chemical vapor deposition) [105].

As device size becomes smaller, film thickness and tolerance approach the atomic scale. Therefore, uniform and smooth thin film surfaces are desired for high device performance. Surface processing is commonly used in developing integrated circuits in order to build features with dimensions of 100 nm and below [50, 126]. Some other applications of surface processing are mechanical coatings [30], thermal coatings [88] and MEMS devices [12].

2.3.1 Past work in the molecular modeling of thin film deposition

The structure of a thin film is usually a strong function of process inputs such as pressure and temperature. As a result, process models that describe the relationship between inputs and outputs are needed to control the microstructure of thin films.

If the morphological features of the thin film surface are much larger than the mean free path of the particles, continuum assumptions can be used for modeling. Such an approach implemented by Kardar et al. [58] is based on solving a partial differential equation that describes the evolution of surface height. When the continuum assumptions are not valid at such small scales (i.e. the mean free path of particles are comparable to the feature size of a material), developing accurate continuum models is not possible. An alternative approach is simulating the dynamics of the system by means of molecular simulations, such as molecular dynamics (MD) and Monte Carlo (KMC) simulations. In MD simulations, Newton's equations of motion are solved for the position of each atom in a system. Given the initial location of each atom, a potential energy function is used to compute the interaction between atom pairs and this information allows the computation of the trajectory of each atom over a time interval. Due to the computational demands of MD, it is not possible to simulate thin film growth which usually takes minutes or hours using this approach. On the other hand, KMC simulation is a stochastic method that simulating the evolution based on randomly generated numbers. Some examples of KMC models used for thin film growth include Gilmer's work on Al thin films [55] and the work of Srolovitz on diamond deposition [6].

In some cases, a wide separation of length scales associated with different kinds of reaction or diffusion phenomena can make the multiscale modeling of thin film growth necessary [123]. For example, for the epitaxial film growth, stochastic KMC can be used to capture the film surface evolution as a result of microscopic processes taking place on the surface (e.g. diffusion of individual atoms). In addition, a continuum model can be necessary to capture the dynamics in the bulk phase. In a past study, a PDE model for the fluid phase was coupled with a KMC model to investigate the transitions between different growth modes by Lam and Vlachos [78, 122]. Similar studies by Pricer, Drews et al. also involved using hybrid models that are composed

of continuum and KMC models for the electrodeposition of Cu [20, 21, 78, 99, 100, 122]. In this case, the continuum model provided the flux values as an input to the KMC simulations, whereas KMC simulations input the concentration values into the continuum model.

Even though the simulation models give insight about process dynamics, their high computational load prohibits the direct use of them for dynamic analysis and optimization tasks [123]. Several approaches have been developed in the past few years to integrate molecular simulations with dynamic analysis and optimization through reduced order modeling. One approach is the derivation of stochastic time evolution equations from the probabilistic master equation [39], under assumptions that are mostly applicable for well-mixed reacting systems. An alternative approach for constructing stochastic differential equations from molecular simulations of film growth is the generation of stochastic spatially distributed PDEs for the height profile of a surface [24]. These models describe continuum behavior and dynamics, and thus are not appropriate for modeling surface structure at atomic scales. This modeling approach, with stochastic PDEs, has been used recently to control the surface roughness in a thin film deposition process on a one-dimensional lattice [85]. On the other hand, controller design based on molecular simulations is an active area of research [93]. In a recent study [104], a reduced order stochastic model obtained from KMC and finite difference simulation data has been used to design a feedback-feedforward controller in order to maintain the current density during the copper electrodeposition process at a constant level.

In order to explicitly model and predict discrete atomic scale structure, a model is required that does not average over the small length scales. One approach to address this issue is equation-free computing, which was first used for stability and bifurcation analysis [36]. A low-dimensional system state was assumed (based on macroscopic

arguments), and short simulations were run from specific initial conditions to approximate time-derivatives. In this methodology, also called the timestepper, the evolution of a process is captured using short simulations. The input of a timestepper is a certain initial condition at time t , and the output is the state at time $t + dt$, where dt is the duration of the simulation clock. If the dynamics of a system is described by deterministic differential equations, the timestepper would simply be the integration of these equations for a dt amount of time, and it would return the final system state. When such equations do not exist in closed form, a timestepper can use small scale realizations of system dynamics (molecular simulations). In that case, the following steps need to be taken by the timestepper:

1. Lifting: In this step, a certain macroscopic initial condition is used to generate a set of microscopic system configurations with that initial condition. For example, the lifting step can involve generating multiple surface configurations that have the same surface coverage in a thin film deposition process. In other words, these configurations would be microscopically different in terms of the location of atoms, but they would be at the same surface coverage level.
2. Evolution: The system is simulated using microscopic scale realizations (molecular simulations). Starting from each microscopic system configuration generated in the lifting step, molecular simulations are run under certain input conditions.
3. Restriction: Using the simulation results obtained in the evolution step, a macroscopic system state is computed. This usually requires averaging the results of molecular simulations run in the evolution step.

In order for timestepper approach to be feasible for a physical system, an important condition should exist. Generating a distribution of microscopic system configurations with a common macroscopic state should not be computationally too demanding. Even though generating that kind of configuration sets is straightforward

when the state description is compact, lifting may involve extremely high computational demand for systems described by high dimensional state vectors (e.g. pair correlation functions describing the relative orientation of atoms on a film surface). Also, for the timestepper approach to be successful, high order moments of microscopic realizations, which can rapidly change, need to converge onto a slow manifold even if these realizations are initialized inaccurately. The inaccurate initialization can possibly occur for many physical systems since the number of possible microscopic configurations consistent with a certain macroscopic property can be infinitely large. In that case, the system state predicted after the restriction step might be inaccurate.

More recently, timestepper method has also been applied for optimization [8]. However, the reduction in computational time achieved by this method, compared to running full molecular simulations, and may not be sufficient to make the approach practical when many predictions over long time intervals are required.

In an alternative approach presented by Gallivan and Murray [33], molecular simulations were used to construct Markov models, with discrete states describing groups of similar configurations. This grouping strategy was based on the similarity of the roughness and coverage values, and enabled computation of the optimal temperature profile by penalizing the surface roughness and temperature changes during the thin film deposition. The construction of this explicit low-order model reduced the computational load by four orders of magnitude relative to the KMC simulations. This significant reduction was critical for making the dynamic simulation and optimization feasible over macroscopic processing times.

The construction of a dynamic model relies on the existence and knowledge of the system state. While a low-order state was selected in previous studies using physical [33] or mathematical arguments [36], the detailed structure in a molecular simulation may not even have a low-order representation. For example, high order statistical moments may not always be slaved to the low order statistical moments (coarse variables)

that are describing the system state, as was assumed in equation-free computing [36]. In order to address these issues, the study by Oguz and Gallivan [97] proposed the use of high-dimensional step-step correlation functions, which provide a more detailed state description of the film surface, compared to the surface roughness alone. This study demonstrated the implementation of principal component analysis (PCA) for reducing the state dimension and self organizing map (SOM) for automated grouping of the similar states in the state space. In a more recent study, Varshney and Armaou [121] also used spatial correlation functions to characterize the state of their film growth simulation, and used equation-free computing to simulate the dynamics [36].

CHAPTER III

MODELING AND EXPERIMENTAL VALIDATION OF STRUCTURE DEVELOPMENT IN HYPERBRANCHED POLYMERS

Highly branched polymers, which include dendritic, hyperbranched or multibranched polymers, are interesting and versatile materials and display several unique properties when compared with their linear analogs. These include low solution and melt viscosities, high solubilities and the presence of very large number of functional end groups that offers the possibility for further modification for various specialty applications.

Discussion of the synthetic methodologies for preparation of a wide range of hyperbranched and dendritic polymers can be found in several extensive review articles [57, 60, 125]. Many of the past experimental studies have focused on the synthesis and characterization of AB_n type monomers ($n \geq 1$), particularly with AB_2 monomers [35, 57, 60, 125]. However, very few of the AB_n type monomers are commercially available due to their lack of symmetrical functionality and tendency to react prematurely [125]. As a result, $A_2 + B_3$ polymerization has recently been the subject of extensive research [25, 26, 56, 68, 73, 74, 84, 130] since it provides an alternative and more convenient way to synthesize highly branched polymers. In contrast to polymerization of AB_n type monomers, these systems offer a wider range of molecular structures depending on the monomeric types and processing conditions. For example, $A_2 + B_3$ polymerization has been performed by heating a mixture of $A_2 + B_3$ (mixed together or one pot) [117] and also by dropwise addition of A_2 into B_3 [118]. The molar ratio of A_2 to B_3 can also be varied [22, 28, 125]. The third main synthesis option (dropwise addition of B_3 into A_2) is known to lead very high degree

of branching that leads to crosslinking at very low conversions [95]. Therefore, it is not a preferred synthesis method and we focus on the other two methods. All of the mentioned synthesis methods, in which only the A and B functional groups react, are illustrated in Figure 3.

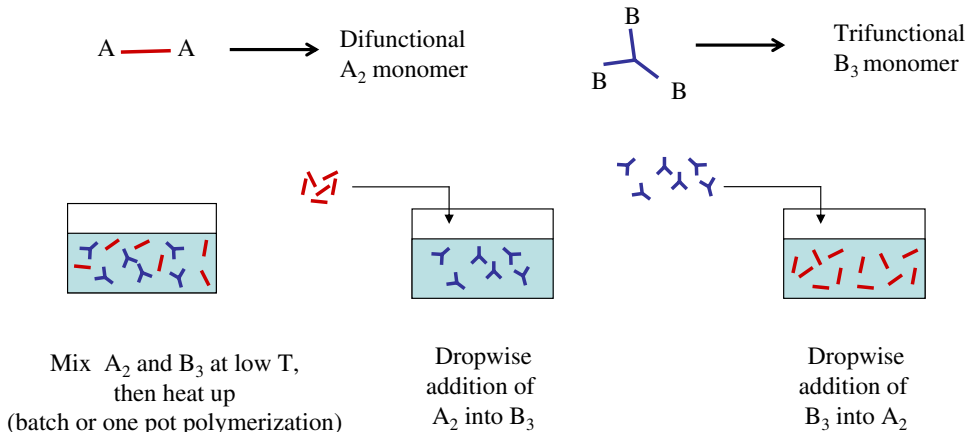


Figure 3: Different synthesis methods using $A_2 + B_3$ polymerization.

In this chapter, we present a Kinetic Monte Carlo (KMC) model, and use it to interpret the experimental findings and elucidate the underlying kinetics of hyperbranched polymerization of $A_2 + B_3$ systems. In the first section, we study a solution polymerization process by dropwise addition of A_2 into B_3 in which the solution concentration influences the extent of cyclization, which affects the gel point and polymer properties significantly. In the second section, we study a melt polymerization process of one pot $A_2 + B_3$ polymerization with no solvent, where the extent of cyclization is negligible. But in this case, due to the aromatic nature of the B_3 monomers used, it is possible to observe reduced B group reactivity after the reaction of a B group in a free B_3 (with 3 unreacted B groups). Therefore, we mainly consider the effects of unequal reactivities, and also endcapping reactions on the structure development of hyperbranched polymers for the melt polymerization case. In the final section, we discuss how the high dimensional and computationally expensive KMC simulation model and experimental data can be used to derive a reduced order model for design

and optimization of the hyperbranched polymerization process.

3.1 Synthesis of hyperbranched polyurethaneureas in solution

In this section, we present a KMC simulation model for the hyperbranched polymerization of difunctional A_2 oligomers, and B_3 monomers with dropwise A_2 addition into a B_3 solution. Using this model and the experimental data, we study the effects of solution concentration on the polymer structure development. By changing the monomer concentration and the extent of cyclization, we obtain polymers with different molecular weight and degree of branching.

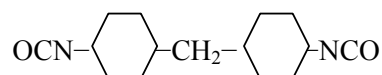
3.1.1 Experimental Procedure

Experiments were performed by Serkan Unal in Dr. Timothy Long's group in Virginia Tech. Isocyanate end-capped Poly(tetramethylene oxide)-glycol (A_2 in this system) was prepared in bulk at 80°C under the catalytic action of 100 ppm of dibutyltin dilaurate (T-12). Polymerization reactions for the preparation of hyperbranched polymers based on PTMO were carried out in THF/IPA (25/75; v/v) solutions, at room temperature, under strong agitation. During the reactions oligomeric A_2 (isocyanate end-capped PTMO) solution was always added into B_3 (polyoxyalkylenetriamine, also known as TRI) solution drop-wise. In order to monitor the growth in the molecular weight of the products, samples were withdrawn from the reactor at different amounts of A_2 addition and endcapped with cyclohexyl isocyanate prior to analysis by size exclusion chromatography (SEC) using a multiple angle laser light scattering (MALLS) detector. Structures of A_2 and B_3 are shown in Figure 4.

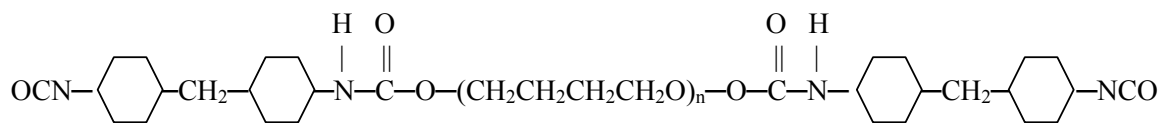
3.1.2 Simulation Procedure

Similar to the experimental procedure followed, initially, N molecules of B_3 are present in the system, and molecules of A_2 are then added sequentially during each

Bis(4-isocyanatohexyl)methane (HMDI)



HMDI and capped PTMO-2000



Polyoxyalkylenetriamine (TRI)

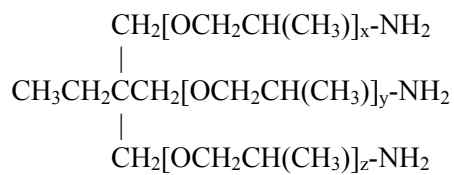


Figure 4: Chemical structures of monomeric and oligomeric A_2 and B_3 .

simulation run.

The simulations consist of three steps. First, an A_2 monomer is added to the system. An unreacted B group is then selected, and is reacted with one of the two A groups. Each unreacted B group in the system has an equal probability of being selected, independent of molecular structure. In the third step, the remaining A group is reacted with another B group. When no cyclization is allowed, then the A group and the B group must be selected from different molecules, but each eligible B group has the same probability of selection.

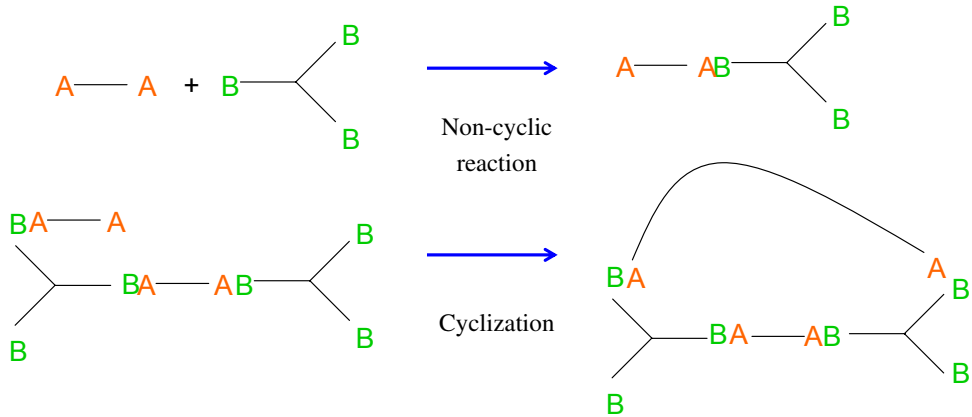


Figure 5: Cyclization versus non-cyclization reactions.

Cyclization is an important factor in step-growth reactions leading to the formation of dendritic and hyperbranched macromolecules [23, 42, 43, 72]. In the simulation studies cyclization, which is illustrated versus non-cyclization reactions in Figure 5, was taken into account in the following manner: an A group and a B group in the same molecule may react, but the selection probability for each B group is not equal. Instead, there is one selection probability for each B group in the same molecule as the A group, and a different probability for each B group not in that molecule. The selection probabilities are calculated from rates, using the kinetic Monte-Carlo (KMC) simulation algorithm of Bortz and co-workers [10], in which the selection probability of each event is proportional to its rate. In the simulations, a variable cyclization parameter γ is defined, such that $\gamma = (k_c/k_{nc})/N$, in which k_c is the per

end group rate of cyclization reaction, and k_{nc} is the rate when a non-cyclization reaction occurs. In the simulations, rather than the individual values of k_c and k_{nc} , their ratio γ is critical for structure evolution. A high γ value would indicate dilute conditions in which the functional groups in the same molecule would be more likely to react with each other than the functional groups in other molecules. Other than dilution, faster mixing can also change the γ value. With fast mixing, the γ value would be low since faster mixing would increase the chance of intermolecular (non-cyclization) reactions. During calculations, the cyclization parameter γ was varied and the development of molecular characteristics such as number average molecular weight (M_n), weight average molecular weight (M_w), polydispersity (PDI), degree of branching (DB), and cycles per molecule were determined. These characteristics are known to have significant effects on the rheological and thermal properties of polymers [89]. DB was calculated using: $DB = (D + T)/(D + L + T)$, where D , L and T indicate dendritic, linear and terminal units in the polymer that are illustrated in Figure 6. The dependence of cyclization probability on conversion is not explicitly built in the KMC model used in simulations, because rate constants do not depend on conversion or molecular structure. However, as the monomer conversion increases, cyclization events become more likely due to the smaller number of molecules and the higher number of unreacted groups per molecule. Consequently, an increase in cyclization probability with conversion is implicitly built into the simulation model.

In the simulations presented here, the simulation size $N = 1000$. Smaller and larger simulation sizes of $N = 100$, $N = 700$ and $N = 1300$ were also studied. The simulations with $N = 100$ differ significantly from the larger simulations, but the simulations with $N = 700$, $N = 1000$ and $N = 1300$ agree quantitatively, suggesting that the results reported in this study ($N = 1000$) are not dependent on the system size. The only exception occurs when there is no cyclization. In this case the simulations with $N = 700$ and $N = 1000$ differ near full conversion, mainly because the

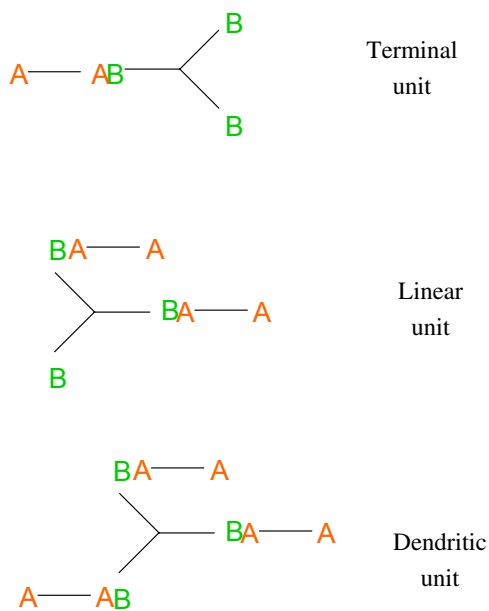


Figure 6: Dendritic, linear and terminal B_3 units.

molecular weight is equal to the total weight in the system. However, this regime is not relevant to the experimental data and is not reported.

To mimic various levels of cyclization in the polymers formed, the simulations were carried out at different cyclization ratios, such as $\gamma = 0$ (no cyclization), $\gamma = 0.01$ and $\gamma = 0.1$ (low cyclization) and $\gamma = 1$ (very high cyclization). Molecular weights of A_2 and B_3 are taken as 2500 and 440 g/mol, to mimic the experimental system based on isocyanate capped PTMO-2000 and TRI.

Similar to the experiments, during simulations oligomeric A_2 is added into B_3 slowly up to a stoichiometric ratio of $A_2 : B_3 = 1.15$ or 115% A_2 .

3.1.3 Results

Figures 7(a) and (b) show the Monte Carlo simulation results on the development of the number average M_n and weight average M_w molecular weights as a function of A_2 addition, for a 1000×1000 ($A_2 \times B_3$) system. As depicted in Figure 7(a), regardless of the cyclization ratio, a slow increase in M_n values were observed until a fairly large amount of A_2 (about 75%) is added into the system. This is followed

by a sharp increase for polymers in which the cyclization is not taken into account, where M_n value is theoretically expected to reach infinity, as predicted by Flory's theory [29, 28]. As expected, the growth in the number average molecular weight M_n is severely limited for polymers that show moderate to high level of cyclization. As depicted in Figure 7(a), even at the fairly low cyclization ratio of $\gamma = 0.01$, at 100% A_2 addition, a very dramatic reduction in M_n is observed, where it only reaches to about 60,000 g/mol. At a cyclization ratio of 0.1, M_n is further reduced and reaches to only about 20,000 g/mol at 100% A_2 addition. Simulations performed assuming the highest cyclization ratio of $\gamma = 1$ clearly show the formation of very low molecular weight products, which is expected. As clearly demonstrated in Figure 7(b), cyclization has less of an effect on the development of M_w . Even at a cyclization ratio of $\gamma = 0.1$, M_w reaches to very high values. Only at the highest cyclization ratio of $\gamma = 1$, similar to M_n , there is a dramatic reduction in M_w . Simulations clearly indicate that cyclization delays the onset of gel formation well beyond the theoretical A_2 conversion of 75%. In order to make a direct comparison, Figures 7(a) and (b) give the experimental results on M_n and M_w in addition to the results of Monte Carlo simulations. It is interesting to note that experimental M_n and M_w values obtained for polymerizations at 25% solids agree fairly well with simulations, where cyclization ratio is low ($\gamma = 0 - 0.01$). Even more interestingly, experimental M_n and M_w values obtained for polymerizations at 10% solids agree very well with simulations where degree of cyclization is higher ($\gamma = 0.1$).

Figures 8(a) and (b) show simulation results on the polydispersity index (PDI) and the degree of branching (DB) as a function of A_2 addition. As depicted in Figure 8(a), PDI shows a gradual increase as A_2 is added into the system and reacted with B_3 . As expected, after about 60% A_2 addition there is a dramatic increase in PDI for all systems, except for the case of very high cyclization, or when $\gamma = 1$. Simulation results on the degree of branching (DB) as a function of cyclization parameter γ

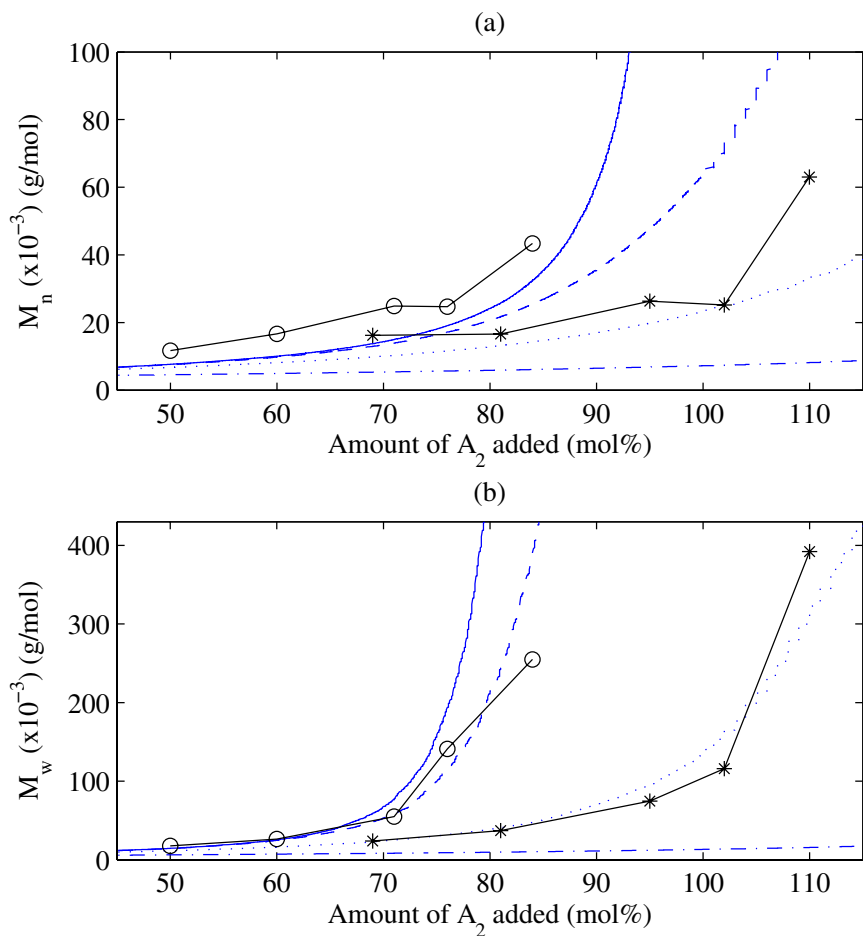


Figure 7: Comparison of experimental and simulation results on the development of average polymer molecular weight as a function of A_2 addition and cyclization ratio for a $1000 \times 1000 A_2 \times B_3$ system. (a) Number average molecular weight and (b) weight average molecular weight. Experimental data: (o) experiment with 25% solids by weight, and (*) 10% solids by weight. Simulation data: $\gamma = 0$ (solid line), $\gamma = 0.01$ (dashed line), $\gamma = 0.1$ (dotted line), $\gamma = 1$ (dash-dotted line).

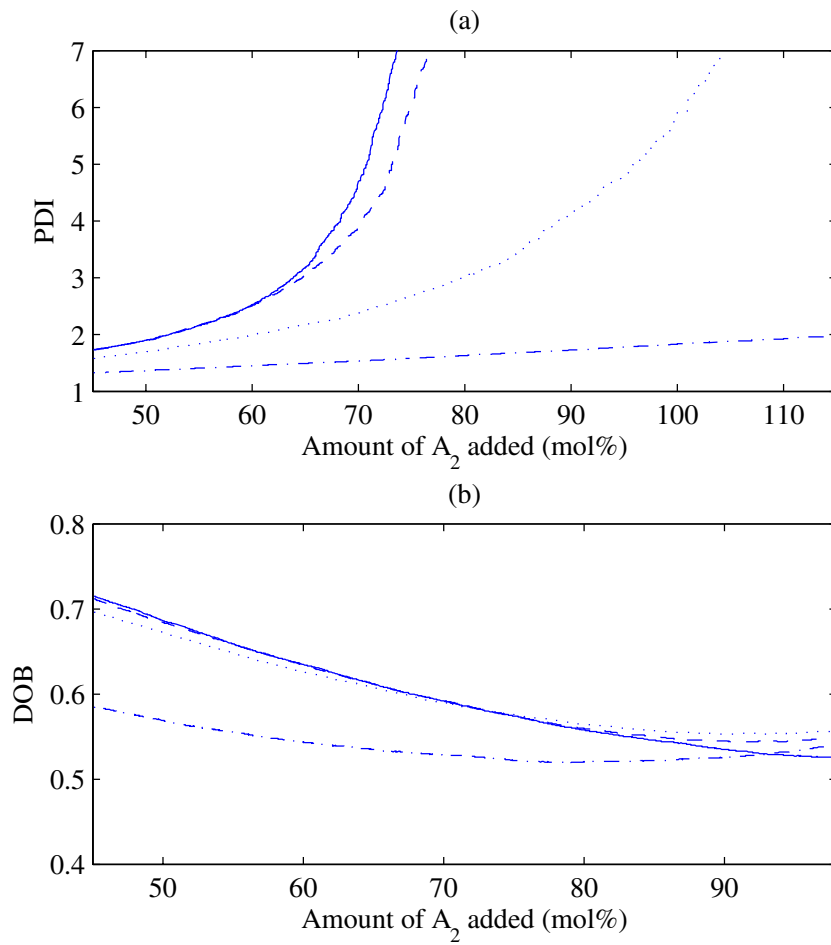


Figure 8: Kinetic Monte Carlo simulations on the polydispersity index (PDI) (a) and the degree of branching (DB) (b) as a function of A_2 addition. $\gamma = 0$ (solid line), $\gamma = 0.01$ (dashed line), $\gamma = 0.1$ (dotted line), $\gamma = 1$ (dash-dotted line).

and amount of A_2 added are depicted in Figure 8(b). It is important to note that at the beginning of the simulations there is only B_3 species in the system, so $DB = 1.0$. As A_2 is added into the system and reacts with B_3 , DB starts going down slowly and levels off around 0.5 – 0.6 after about 80% A_2 addition for simulations where cyclization parameter is low ($\gamma = 0 - 0.1$). Interestingly, for the system with highest cyclization parameter ($\gamma = 1$), the behavior is quite different. Monte Carlo simulations indicate that in this system with very high cyclization probability, even at low A_2 additions, DB starts around 0.75 and gradually moves to about 0.5-0.6 similar to the other systems. Frey has pointed out that DB statistically approaches 0.5 for the case of polymerization of AB_2 monomers [51], calculated using $DB = (D+T)/(D+L+T)$ [47]. Our simulations resulted in a DB value of 0.53 at complete A_2 addition without cyclization. This is in excellent agreement with simulations of Frey [51] and also with the DB values observed experimentally in $A_2 + B_3$ systems [17]. With the inclusion of cyclization, DB goes up slightly from this value of 0.53, as shown in Figure 8(b).

Another important characteristic of such highly branched polymers, the number of cyclization events per molecule, as a function of A_2 conversion is shown in Figure 9. For small cyclization ratios ($\gamma = 0.01$ and $\gamma = 0.1$) and low conversions, the number of cyclic species per molecule are negligible until about 80% A_2 addition. As the amount of A_2 exceeds 80% and high molecular weight polymers are obtained, cyclization increases and reaches to about 2 per molecule. Interestingly, simulation results given in Figure 9 on the amount of cyclization per molecule, for high cyclization ratio (e.g. $\gamma = 1$) and high conversions seems to be somewhat contradictory to expectations, since they are much smaller. However, since the molecular weight of the polymer formed is strongly suppressed due to extensive cyclization, in these simulations polymers formed have very low molecular weights (Figures 7(a) and (b)). In other words, in these cases only very small molecules which also have a smaller

total number of cycles are formed.

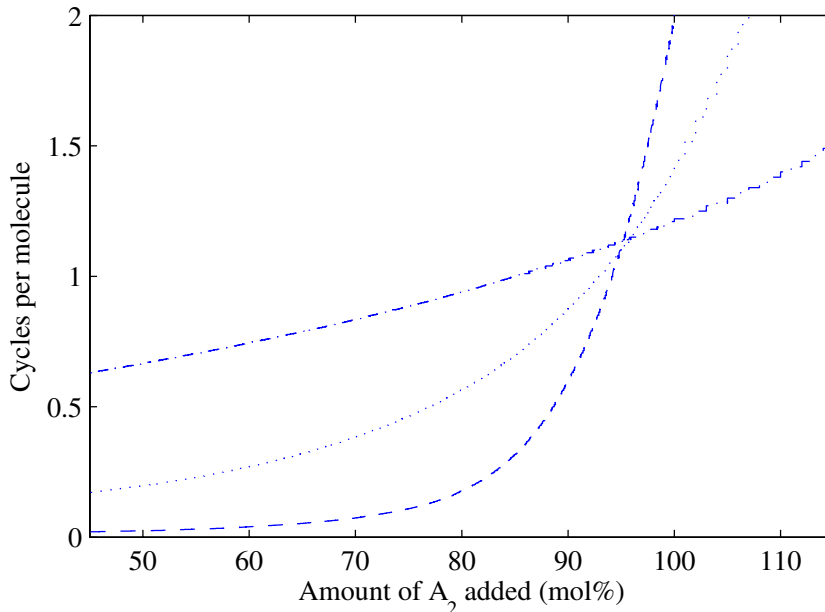


Figure 9: Number of cyclization events per molecule, as predicted by the kinetic Monte-Carlo simulations as a function of A_2 addition, at different levels of cyclization. $\gamma = 0.01$ (dashed line), $\gamma = 0.1$ (dotted line), $\gamma = 1$ (dash-dotted line).

3.2 *Synthesis of hyperbranched polyetheresters in melt*

In the previous section, we used a simple KMC simulation model to explore the role of cyclization in the dropwise A_2 addition into a B_3 in a solvent. Herein we consider a similar statistical framework to explore a wider range of phenomena, for a batch $A_2 + B_3$ reaction in a melt. The Monte Carlo simulations are used to interpret the experimentally measured number-averaged molecular weight, weight-averaged molecular weight, and the density of branched units, by considering the effects of the cyclization reactions, unequal reactivities, and endcapping on the structure development of hyperbranched polymers.

3.2.1 Experimental Procedure

Highly branched poly(ether ester)s were synthesized in the melt phase using an oligomeric $A_2 + B_3$ polymerization strategy. Condensation of poly(propylene glycol) (A_2 oligomer) and trimethyl 1,3,5-benzenetricarboxylate (B_3 monomer) generated highly branched structures. The conversion and degree of branching were measured with ^1H NMR spectroscopy, and the molecular weight (number-average, weight-average, and polydispersity) was characterized by size exclusion chromatography (SEC) using a multiple angle light scattering (MALLS) detector, at six points throughout the polymerization process. Additional experiments were also performed in which two different monofunctional endcapping reagents were added for the purpose of delaying the gel point. This experimental work is by Unal and Long [117], and can be consulted for further detail on the experimental procedures. In this thesis, we interpret these experiments through a comparison with Monte Carlo simulations, utilizing alternative assumptions and kinetic models. Because the error in the SEC measurements is approximately 10%, our goal is not to achieve exact agreement between the experiments and modeling, but to compare the trends and magnitudes.

3.2.2 Simulation Procedure

This study focuses on the use of a model to interpret the experimental data. Each simulation starts with N monomers of A_2 and N monomers of B_3 in the system, since the monomers were polymerized with a 1 : 1 molar ratio during the experiments. At each step in the simulation, all of the available (unreacted) functional groups are listed. Then, using a random number, an A group and a B group are selected from the list and the reaction is executed. This is followed by updating the list of available A and B groups, molecular weights of the molecules, and the number of dendritic, linear and terminal units in the system.

The probability of selecting a particular pair of A and B groups is proportional

to the reaction rate for that pair, which yields the correct time-evolution of the system [10]. Therefore, a model is needed for the reaction rates of various events. The first effect that is considered in this effort is the formation of cycles through intramolecular reactions (5). Polymerization was performed in the melt, which was hypothesized to minimize cycle formation due to high concentration of reactants, so cyclization reactions are taken into account in the simulations to address this hypothesis. While selecting an A group (from A_2 oligomers) and a B group (from B_3 monomers) for reaction at each simulation step, some pairs are favored more than the others. If the reaction of an $A - B$ pair leads to cycle formation, the selection probability of that particular pair is promoted by the cyclization parameter γ , such that $\gamma = (k_c/k_{nc})/N$, where k_c is the rate of cyclization reactions for each $A - B$ pair and k_{nc} is the rate of non-cyclization reactions for each $A - B$ pair. N appears in the parameter γ since the number of intermolecular reactions in the simulations grows as N^2 , while intramolecular reactions are initially proportional to N . Although γ is constant throughout each simulation, the overall number of cycle-forming reactions increases with conversion, because the number of possible cycle-forming reactions increases with molecular weight. We use different values of γ in the simulations to explore the effect of cycle formation on the resulting polymer.

In addition to cyclization and non-cyclization reactions, we also consider end-capping reactions between E groups (from monofunctional end-capping reagents) and B groups, with a rate constant of k_e . The ratio $\epsilon = k_e/k_{nc}$ is then the second parameter in our kinetic model. In a second set of simulations, γ and ϵ have been varied to observe their effects on the development of molecular characteristics such as number average molecular weight (M_n), weight average molecular weight (M_w), polydispersity index ($PDI = M_w/M_n$) and the fraction of dendritic units (f_D). f_D is calculated using: $f_D = D/(D + T + L)$, where D, L and T indicate the number of dendritic, linear and terminal units in the system (illustrated in Figure 6). We plot

f_D in the simulation results (as opposed to f_L , f_T), or DB since f_D is the quantity extracted from the ^1H NMR measurements [117].

Depending on the properties of the monomers, such as their molar mass or electrostatic interactions, unequal reactivities of their groups can also affect the structural development of hyperbranched polymers. End groups usually have higher reactivities than the groups along the length of the chains because of the lower kinetic excluded volume effect [89]. This causes the linear units to have lower reactivities than the terminal units. In order to simulate unequal reactivities of the B groups, we define a third parameter ρ . For each unreacted B group in a B_3 monomer, we check the reaction state of the other two B groups in the monomer. We then consider three possible cases of unequal reactivity. For this purpose, we define k_1 to be the rate of reaction of a B group in a free B_3 monomer, k_2 to be the rate of a B group in a terminal unit, and k_3 for a B group in a linear unit. We assign reaction rates in the ratio of $\rho = k_1/k_2 = k_2/k_3$. k_1 is expected to be enhanced relative to k_2 due to the greater mobility of the free B_3 monomer and its ability to diffuse through the polymer, while k_2 may be different from k_3 due to blocking, free volume, and electrostatic considerations. In order to isolate these two effects, we have also performed simulations with $\rho_{12} = k_1/k_2$ and $k_2 = k_3$, and then also with $\rho_{23} = k_2/k_3$ and $k_1 = k_2$.

For the simulation results presented in this study, the system size is $N = 10,000$. Smaller and larger simulation sizes of $N = 1000, 3000, 5000$ and 7000 have also been used. The simulation results with $N = 1000$ differ significantly from the simulations with larger N . On the other hand, simulations with larger N agreed quantitatively. This suggests that the trends reported in this study are not dependent on the system size. Additionally, for the case of no cycle formation and equal B_3 reactivity, we compared our weight-averaged molecular weight with the analytical theory of Stockmayer [111], and the error is approximately 1%. Clearly, if one is interested in describing the approach to gelation with no bound on the molecular weight, then the system

size would also need to approach a macroscopic number of monomers (10^{23}), and studies have been performed to quantify this tradeoff [107]. However, in these simulations, only molecular weights up to 500,000 g/mol are presented, since this is the range of the experimental data. Furthermore, the error in the SEC measurements is approximately 10%, so an error of 1% in the model predictions is not significant.

3.2.3 Results

Monte Carlo simulations have been carried out in order to assess the effects of cyclization, unequal reactivities, and end-capping reactions on the polymer structure development. Molecular weights of A_2 (PPG-1000) and B_3 (TMT) are 1060 and 252 g/mol, whereas the molecular weights of two types of end-capping reagents, PPG-M-1000 and dodecanol, are 1200 g/mol and 187 g/mol, respectively. The simulation system containing $10,000 \times 10,000$ A_2 and B_3 monomers yielded molecular weights in the same range as those observed in the experiments. Similar to the experiments, simulations have been initialized with an equal number of A_2 and B_3 monomers in the system. The simulation results plotted are averages over 50 independent realizations.

3.2.3.1 Effects of cyclization reactions

The melt polymerization was previously considered to be sufficiently concentrated so that the cycle formation would be negligible [117]. We first investigate the extent of cyclization via the simulations. Extent of cyclization (EOC) is defined as the fraction of reactions between A and B that is intramolecular, and is a quantity that has been measured previously using MALDI-ToF for linear polymers but is less reliably measured in hyperbranched polymers due to the many possible isomers [71, 74].

Figures 10(a) and (b) show the experimental [117] and simulated evolution of the weight-average molecular weight M_w and the polydispersity PDI as a function of A_2 conversion, for the $10,000 \times 10,000$ $A_2 + B_3$ system with different γ values. The reactivity ratio ρ is 1 and no monofunctional reagents are present. For all γ

values, a slow increase in M_w and PDI values was observed until about 80% A_2 conversion. Above 80% A_2 conversion, a sharp increase in M_w and PDI takes place in all the systems, except for the one with the highest level of cyclization ($\gamma = 1$). The experimental data for PDI and M_w are most consistent with a value of γ around 10^{-2} . The ideal limit of no cycle formation, modeled by Flory, is the solid curve with $\gamma = 0$.

M_w and PDI development in systems with cyclization ratios in the range of $\gamma = 0$ to 10^{-2} all agree reasonably well with the experimental data. Due to the variability of the experimental measurements, the goal of the modeling is not to match the experiments exactly, but to assess the magnitude of the various effects. The simulation with $\gamma = 1$ has a lower M_w and PDI at high conversions, compared to the experimental data. This suggested that the extent of cyclization, which is the fraction of reactions that are cycle forming, was quite low during the experiment. Figure 11 shows that the system with $\gamma = 1$ reaches an extent of cyclization of 0.09 at 90% A_2 conversion. Interestingly, even such a low extent of cyclization dramatically suppressed the M_w and PDI as illustrated in Figures 10(a) and (b). Figure 11 also shows that with $\gamma = 10^{-3}$ or $\gamma = 10^{-2}$, the extent of cyclization is less than 3%. This result supports the original hypothesis, that melt polymerization would suppress the effect of cycle formation on molecular weight and gelation.

Another important characteristic of hyperbranched polymers is the fraction of dendritic units f_D , which is directly proportional to the extent of branching in the system. The development of f_D at different γ levels is shown in Figure 10(c) as a function of A_2 conversion. As γ is increased from $\gamma = 10^{-1}$ to $\gamma = 1$, an increase in f_D is observed at A_2 conversions of 60% and above. This trend was expected, since cyclization reactions enhance the number of dendritic groups via the formation of small, fully reacted polymers with no free groups. However, the evolution of f_D is not consistent with the experimental data for any value of γ . For all γ , the simulations

predict a higher fraction of dendritic units.

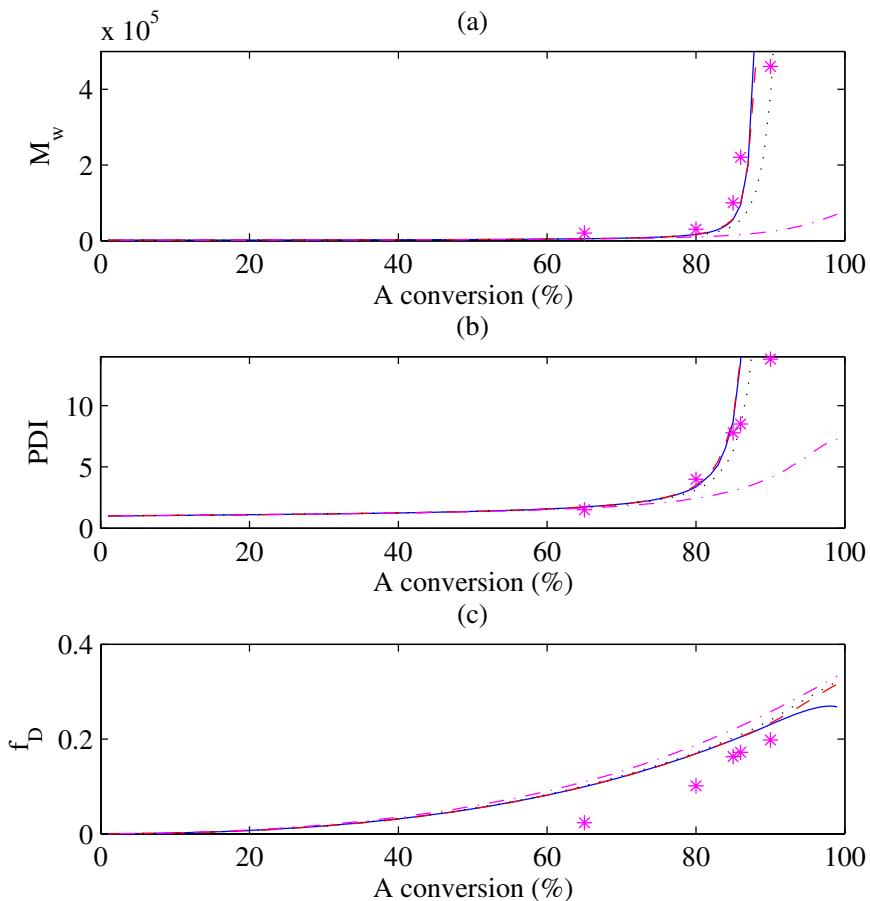


Figure 10: Comparison of simulation and experiment [117] (*). In the simulations, the cyclization ratio γ is varied. $\gamma = 0$ (solid line), $\gamma = 10^{-3}$ (dashed line), $\gamma = 10^{-2}$ (dotted line), $\gamma = 10^{-1}$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D . Agreement between experiments and simulations is not achieved for f_D at any value of γ .

3.2.3.2 Effects of unequal reactivities

The disagreement of f_D between the experiments and simulations suggests that there is an additional effect that suppresses the amount of branching during the experiments, other than the effect of cyclization reactions. It could be attributed to the lower reactivity of free B groups in linear units relative to the free B groups in the

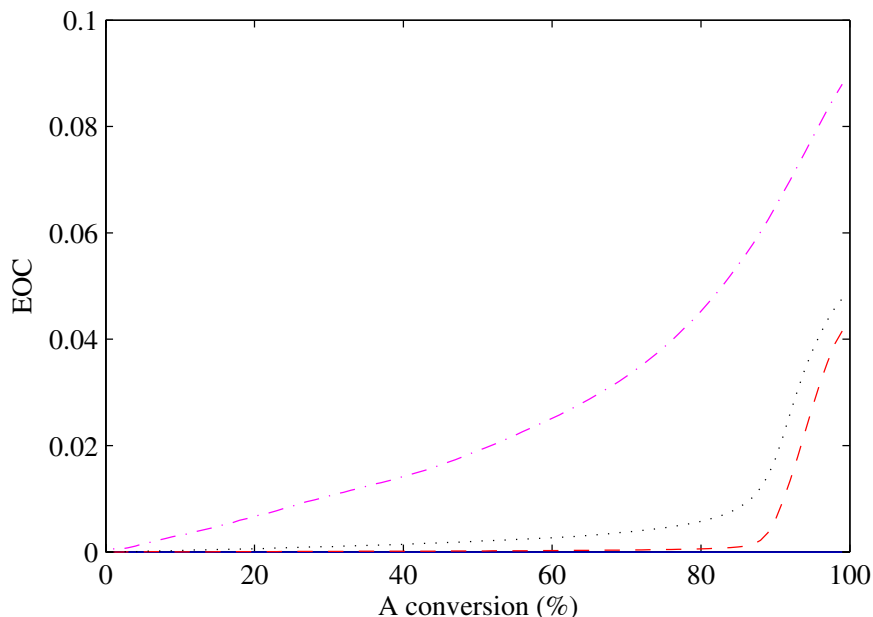


Figure 11: Simulation predictions of extent of cyclization. In the simulations, the cyclization ratio γ is varied: $\gamma = 10^{-3}$ (dashed line), $\gamma = 10^{-2}$ (dotted line), $\gamma = 10^{-1}$ (dash-dotted line). For the $\gamma = 0$ case (solid line), $\text{EOC}=0$.

terminal units or completely unreacted B_3 monomers. However, this would also reduce the molecular weight. In order to assess this trade-off quantitatively, additional simulations have been performed during which the reactivity of free B groups is modified, based on the overall state of the B_3 monomer. In these simulations presented in this section, $\gamma = 0$ since the previous section demonstrated that cyclization did not play a major role in this system.

Figure 12 shows the evolution of the simulations with different levels of reactivity ratio ρ . With $\rho = 1$, the simulations are identical to those in Figure 10, while larger ρ reduces the amount of branching and also the molecular weight and distribution. The reduction in molecular weight was expected since the polymers are becoming more linear, but the reduction of molecular weight and polydispersity is not as dramatic as it is with cyclization. Even for the extreme value of $\rho = 10$, gelation is delayed but not completely suppressed. In Figure 12(c), the increase in ρ causes a decrease in f_D , as expected, but the shape of the curves does not match the experimental data.

In the data, the fraction of dendritic units rises quite high at late conversion, but is very low at earlier conversions. The 1 : 1 molar ratio of A_2 and B_3 is important in understanding this behavior. Due to the 1 : 1 molar ratio, there is an excess of B groups, so at full A conversion, only 2/3 of the B groups have reacted. If the reactivity of the third B group is strongly reduced, then there will be few dendritic units that form, but this is not observed in the experiments. The fact that the experiments eventually reach a large value of f_D near that predicted with $\rho = 1$ instead suggests that a B group in a linear unit has a similar reactivity to a B group in a terminal unit. Therefore, $\rho_{23} = k_2/k_3 = 1$ where k_2 is the rate of a B group in a terminal unit, and k_3 is the rate of a B group in a linear unit. This was also supported by our simulations with ρ_{23} greater than one. In these simulations, compared to the experimental data, very low f_D values were observed.

The low values of f_D around 60% conversion are more consistent with a suppression of the reaction of terminal units relative to free units, as shown in Figure 13 for various levels of ρ_{12} . Recall that $\rho_{12} = k_1/k_2$, with $k_2 = k_3$. The high reactivity of the free B_3 could be due to its mobility, as well as the fact that B groups in polymers may be partially blocked by other portions of the polymer. This trend in f_D is much more consistent with the experimental data. A high value of $\rho_{12} = 10$ best matches the f_D measurements, while a somewhat lower value of ρ_{12} near 1.5 – 2 agrees best in the molecular weight distribution. Our conclusion based on the simulations in Figures 12 and 13 is that a suppressed reactivity of the third B group in the linear unit is not consistent with the observed data. The more consistent explanation is that the free B_3 monomers are more mobile and therefore react faster than B_3 in polymer.

3.2.3.3 *Effects of end-capping reagents*

The third and final simulation study considers the addition of monofunctional reagents to the $A_2 + B_3$ system. Stockmayer’s theoretical studies of highly branched polymers

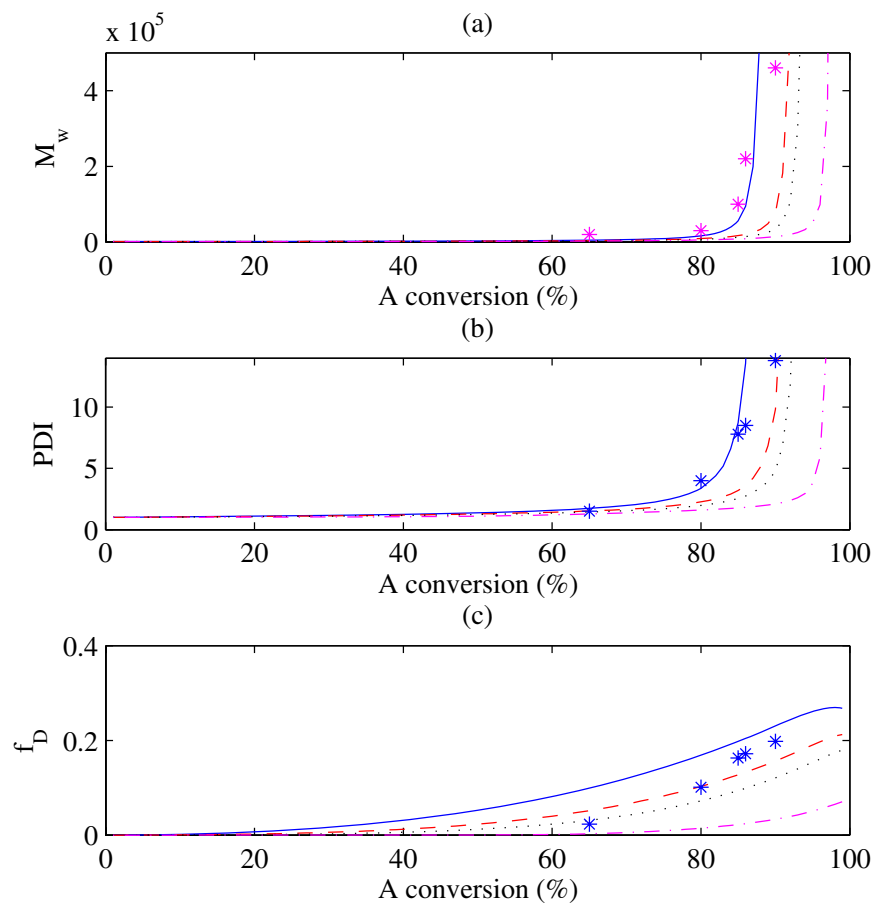


Figure 12: Comparison of simulation and experiment [117] (*). In the simulations, $\rho = k_1/k_2 = k_2/k_3$. $\rho = 1$ (solid line), $\rho = 1.5$ (dashed line), $\rho = 2$ (dotted line), $\rho = 10$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D .

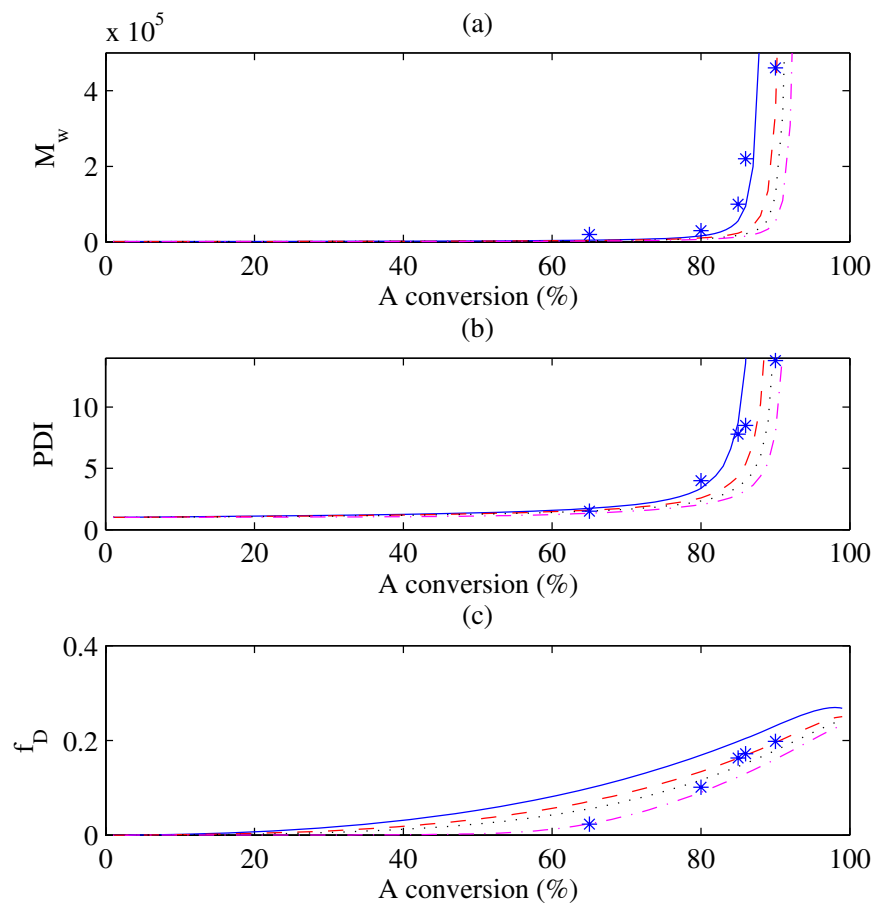


Figure 13: Comparison of simulation and experiment [117] (*). In the simulations, $\rho_{12} = k_1/k_2$ and $k_2 = k_3$. $\rho_{12} = 1$ (solid line), $\rho_{12} = 1.5$ (dashed line), $\rho_{12} = 2$ (dotted line), $\rho_{12} = 10$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D .

indicated that addition of a monofunctional end-capping reagent should shift the gel point to higher monomer conversion values [28, 111, 110]. Thus, delaying the gel point by terminating some of the B functionalities is the main motivation behind using end-capping reagents in this system. A molar ratio of monomers as $A_2 : B_3 : E = 1 : 1 : 1$ was used in the experiments to ensure that residual B end groups, which would be expected in an $A_2 : B_3 = 1 : 1$ system at full A_2 conversion, do not remain in the system at full conversion.

Figure 14 provides a comparison of the change in the evolution of M_w for the addition of PPG-M-1000 (1200 g/mol) at the beginning of the reaction. The A conversion that is plotted also includes the conversion of the end-capping reagent, since that is measured by NMR. Curves are shown for the values of ρ_{12} considered previously in Figure 13. A primary observation is that the addition of the end-capping reagents has the larger effect when the reactivity ratio is also large. Furthermore, this effect is only observed when the end-capping reagents also have a higher reactivity than the A groups in A_2 ($\epsilon \gg 1$). This might be the case for the dodecanol reagent, if the end-capping reagent has a higher diffusivity than the A_2 due to its lower molar mass.

In the simulations presented in Figures 14 and 15, we set $\epsilon = 1000$, although the results are similar for $\epsilon = 10$. When $\epsilon \approx 1$, the simulations predict that the end-capping reagents have a negligible effect on the polymer structure. By reacting with the excess B groups, they only add their extra mass to the polymer. In the opposite limit, when $\epsilon \gg 1$ and $\rho_{12} \gg 1$, the E groups react quickly with the free B_3 monomers, after which the EB_3 units begin reacting with A_2 . Each B_3 is thus bonded to only two A_2 monomers, so the polymers have a linear structure.

In the experiments, it was observed that the gel point was completely suppressed up to 98% conversion of each monomer, and the measured g contraction factor from

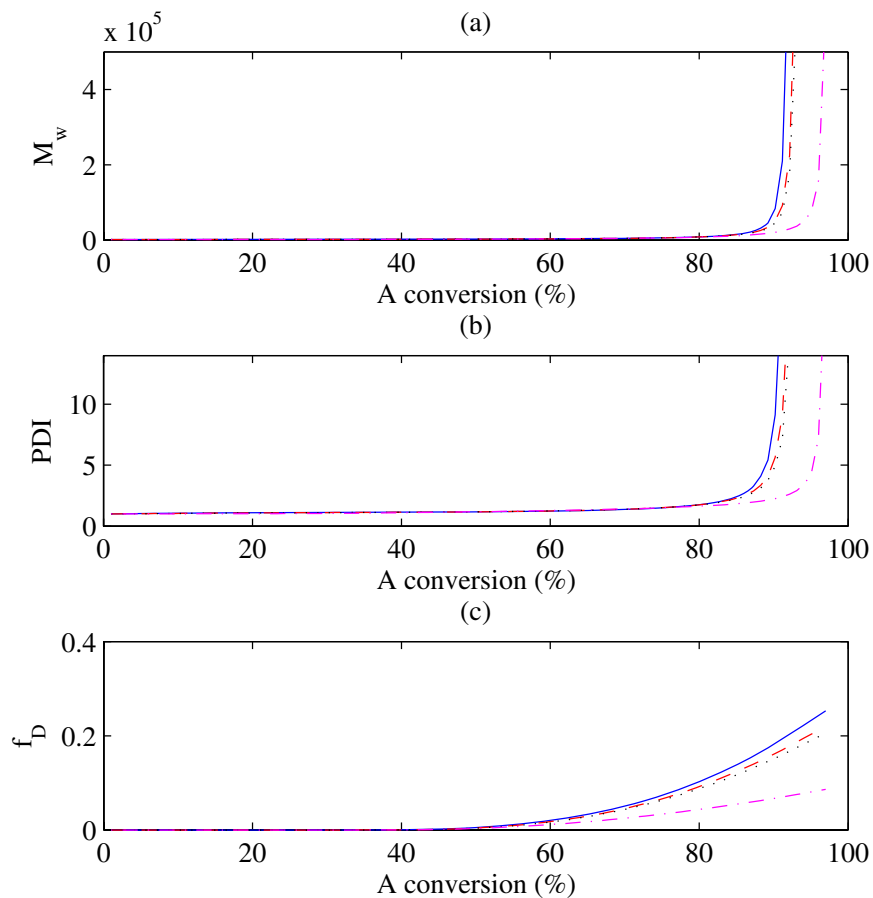


Figure 14: Simulated evolution with end-capping reagents added at the beginning of the process, with $A_2 : B_3 : E = 1 : 1 : 1$. The molecular weight of E (PPG-M-1000) is 1200 g/mol. As in Figure 13, $\rho_{12} = k_1/k_2$ and $k_2 = k_3$. $\rho_{12} = 1$ (solid line), $\rho_{12} = 1.5$ (dashed line), $\rho_{12} = 2$ (dotted line), $\rho_{12} = 10$ (dash-dotted line). (a) Weight-average molecular weight M_w (b) polydispersity index PDI (c) fraction of dendritic units f_D . (Note: dendritic units are calculated here based on the number of A - B reactions. E - B reactions are not considered in the calculation since they do not lead to further branching.)

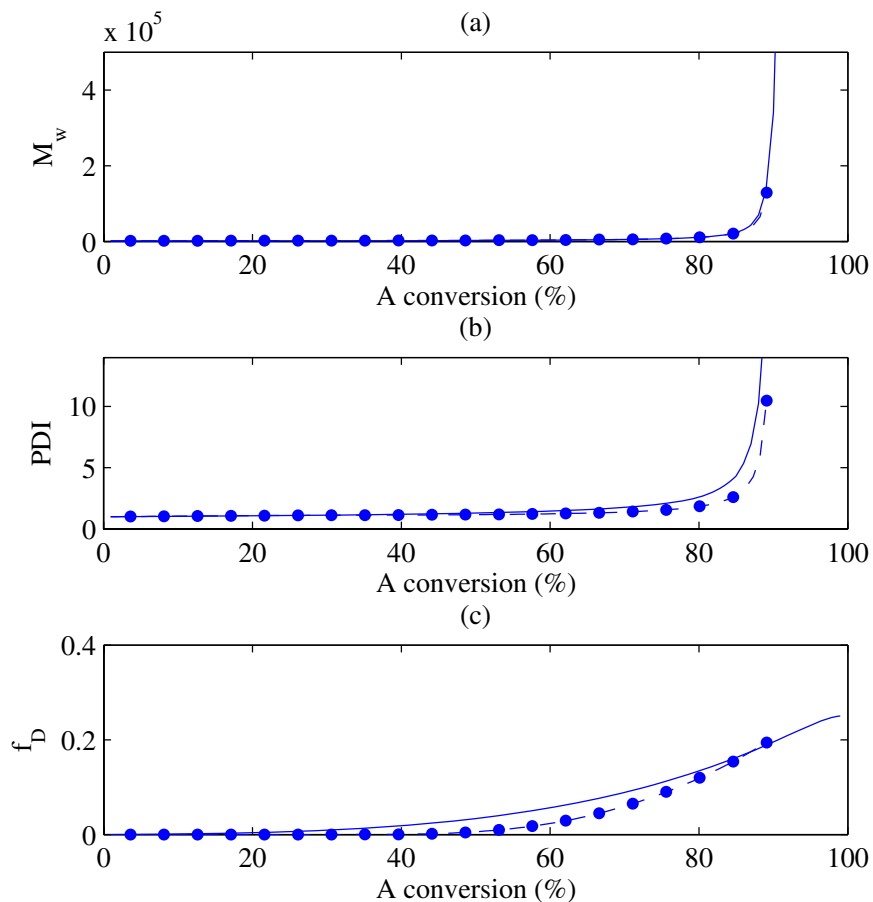


Figure 15: Simulated evolution with end-capping reagents added at the beginning of the process, with variation in stoichiometry: $A_2 : B_3 : E = 1 : 1 : 1$ (solid line) and $A_2 : B_3 : E = 1 : 0.9 : 1$ (dotted line with markers). The molecular weight of E (PPG-M-1000) is 1200 g/mol. $\rho_{12} = k_1/k_2 = 1.5$ and $k_2 = k_3$. (a) Weight-average molecular weight M_w (b) Polydispersity index PDI (c) Fraction of dendritic units f_D .

GPC was more consistent with a highly branched polymer (large f_D) [117]. Our simulations do predict a suppression of molecular weight with the end-capping reagents, but gelation is only delayed and not completely suppressed. At the extreme value $\rho_{12} = 10$, a significant reduction in f_D is also implied, but at lower values of ρ_{12} high levels of branching are still predicted.

In the presentation of the experimental results [117], it was suggested that ester interchange of the A_2 with the monofunctional reagents might account for the observed reduction in molecular weight. While this interchange would cause a randomization of the polymer, at the high conversion of 98% it does not provide a consistent explanation for the extreme reduction in molecular weight. At 98% conversion, most of the A_2 monomers freed up by ester interchange would have reacted with another B_3 monomer. Other effects not included in the model could be causing the suppression of gelation observed in the experiments, such as the spatial distribution of the monomers in the polymer. These effects could be exacerbated when the end-capping reagents are added, since all B groups must eventually react, even those buried or blocked in the center of the spherical polymer. Possibly, at high conversion, such groups are more likely to undergo cyclization reactions, which would suppress the molecular weight.

An alternative explanation was also suggested by our simulations. We observed that the simulation results are extremely sensitive to the stoichiometry near $A_2 : B_3 : E = 1 : 1 : 1$. In particular, if there is a reduction in the amount of B_3 , then not all of the A_2 groups will be able to react with B groups, and the molecular weight will be reduced. This may be a particular issue in the experiments, since B_3 loss may be facilitated by the nitrogen purge at the final polymerization temperature of 180°C. Due to the uncertainty in the final stoichiometry measurements, it is not possible to eliminate this effect. The simulations suggest that that our chosen stoichiometry of 1 : 1 : 1 is not a robust operating point, due to the extreme sensitivity of the molecular

weight on the stoichiometry. Figure 15 shows a comparison of the simulations with $A_2 : B_3 : E = 1 : 1 : 1$ and with $1 : 0.9 : 1$. At this stoichiometry, 90% A_2 conversion is the maximum that can be achieved, and gelation is completely suppressed at full B_3 conversion.

3.3 State space modeling of hyperbranched polymers

So far in this chapter, we have investigated the ways that synthesis route affects basic polymer structure properties such as molecular weight and degree of branching. Even though these results can help process design to a certain extent by giving insight about differences in polymer structure development under a range of polymerization conditions, that is not enough to obtain a process model that can be used for optimization of the synthesis route, which would be necessary to target and reach precisely defined polymeric structures. Even with a very accurate KMC simulation model, one can not perform such an optimization due to the high computational load of these simulations. This model inversion problem (i.e. for a given output, what is the optimal input profile?) was previously illustrated in Figure 2. Reduced order models which can adequately describe the relationship between process inputs and polymer properties and also can be easily inverted, are necessary for accomplishing design and optimization.

During hyperbranched polymerization, the system status can be described by a vector including the number averaged molecular weight, degree of branching, and the polydispersity index. However, hyperbranched polymers can have a high number of geometrical isomers due to the stochastic reaction of each added monomer into the system. These geometrical isomers can have the same molecular weight and branching properties, but a wide range of aspect ratios. The disperse nature of the aspect ratio of geometrical isomers can potentially lead to differences in solubility, and also a wide range in solid-state packing and physical properties [61]. Because

of these factors, variables other than number averaged molecular weight and degree of branching are needed to differentiate between isomers while describing the status of a polymerization system. Topological indices are such descriptors that are widely implemented using the graph theory [45]. In chemical graphs, atoms are represented by vertices, and bonds are drawn as lines that connect the vertices. In graph theory terminology, these lines are referred to as the edges of a graph.

Topological indices can detect very small differences in branching, and are useful for modeling quantitative structure-property relationships. For example, the Wiener index [129, 128] of organic substances is well correlated with properties such as heat of formation, density, boiling point, and viscosity [18]. This index gives the total number of bonds between all atom pairs in a molecule along the shortest path between the atoms. In general, for a set of isomers, Wiener index decreases when number of branches and branch length increases. In other words, the most compact molecule in a set of isomers has the smallest Wiener index. Therefore, Wiener index is maximum for the linear chain in a set of isomers.

The following is an example that shows the calculation of this index for the molecule A given in Figure 16.

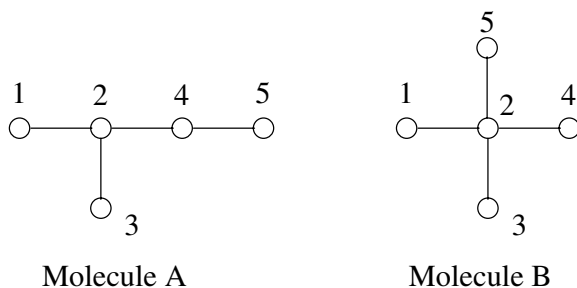


Figure 16: Two isomeric polymers with the same molecular weight and different branching.

- Atom pairs with interdistance of 1: 1-2, 2-4, 4-5, 2-3 (4 pairs)
- Atom pairs with interdistance of 2: 1-4, 2-5, 3-4, 1-3 (4 pairs)

- Atom pairs with interdistance of 3: 1-5 ,3-5 (2 pairs)

Wiener index $W = (4 \times 1) + (4 \times 2) + (2 \times 3) = 18$

Similarly, W can be computed for Molecule B:

- Atom pairs with interdistance of 1: 1-2, 2-3, 2-4, 2-5 (4 pairs)
- Atom pairs with interdistance of 2: 1-4, 1-5, 1-3, 3-4, 3-5, 4-5 (6 pairs)
- Atom pairs with interdistance of 3: None

Wiener index $W = (4 \times 1) + (6 \times 2) = 16$

Since Molecule B is more branched (less linear) it has a lower Wiener index than Molecule A.

Another topological index that is used to quantify polymer branching is the complexity index K [7]. K is defined as the total number of subgraphs in a molecular graph. These subgraphs include vertices, edges, two-edge subgraphs, three-edge subgraphs, etc., and finally the complete molecular graph itself. As a molecule becomes more complex (more branched), the number of subgraphs that it contains increases, resulting in a higher K value. The calculation of K for Molecule A in Figure 16 is given below:

- List of vertices: 1, 2, 3, 4, 5 (Total=5)
- List of edges: 1-2, 2-4, 2-3, 3-5 (Total=4)
- List of 2-edge subgraphs: 1-2-4, 1-2-3, 2-3-5, 4-2-3 (Total=4)
- List of 3-edge subgraphs: 1-2-3-4, 1-2-3-5, 4-2-3-5 (Total=3)
- List of 4-edge subgraphs: whole graph (Total=1)

$K = 5 + 4 + 4 + 3 + 1 = 17$

Similarly, K can be computed for Molecule B:

- List of vertices: 1, 2, 3, 4, 5 (Total=5)
- List of edges: 1-2, 2-4, 2-3, 3-5 (Total=4)
- List of 2-edge subgraphs: 1-2-4, 1-2-3, 1-2-5, 3-2-4, 3-2-5, 4-2-5 (Total=6)
- List of 3-edge subgraphs: 1-2-3-4, 1-2-3-5, 2-3-4-5, 1-2-4-5 (Total=4)
- List of 4-edge subgraphs: whole graph (Total=1)

$$K=5+4+6+4+1=20$$

In this case, molecule B has a higher index since it is more branched.

Degeneracy of a topological index (TI) occurs when two chemical graphs have the identical TI. In that case, the TI will not be adequate to differentiate between these graphs. Since this problem is widely seen with basic TI's such as W and K that are described above, alternative TI's that can capture subtle differences between very similar graphs are developed. Some of these TI's are the Hosoya index, Randić's connectivity index, Kier Hall topological indices and Balaban's index. A detailed discussion of these indices can be found elsewhere [5].

An alternative way to characterize the dynamic system in the KMC simulations is using pair correlation functions [3, 90]. In our case, the pairs would be formed by monomer units. First, we classify the B_3 monomers in the system with at least one reacted B as terminal B_3 units (with 2 free or unreacted B groups), linear B_3 units (with 1 free B group) and dendritic units (with no free B groups). This classification was previously illustrated in Figure 6. Let m_i be the status of monomer i which can take one of the 3 values (D for a dendritic unit, L for a linear unit, and T for a terminal unit). Since m_i can take any of these values, a monomer pair $m_i m_j$ can take 6 values (DD, DT, DL, TT, TL, LL). As a result, monomer pair correlation function (MPC) can be defined as $MPC_{i,j}(r) = \sum_{j=1}^n m_i m_j$ where r is the shortest topological distance between monomers i and j in terms of the number of bonds,

and n is the total number of monomer units in a molecule. The distance r can take any value from 0 to $n - 1$. For a linear molecule with n monomer units, the distance between the end points of the graph would be $n - 1$. One can get a detailed characterization of the polymer system by computing the *MPC* function of each molecule. However, one drawback of this approach would be the high dimensional state description, especially for a system with large number of molecules since each molecule's state vector would have $6 \times (n - 1) \times n^2$ variables. This drawback can be dealt by reducing the dimensions of the data while preserving a high extent of information in it. For example, the whole molecular weight range of the polymers in the system can be divided into bins, and the average *MPC* function for each MW bin can be computed. Also, this characterization should be carried out under a wide range of process conditions to characterize the state space of the system as much as possible.

Once the typical system states are found, a reduced order model can be developed by identifying the transitions between these typical states. Such an approach has recently been used by Gallivan [32] and the author [96, 97] for KMC simulations of a thin film deposition process in which the state was described using macroscopic variables such as film roughness and island density [32], and step-step correlation functions [96, 97]. A challenge in this approach is the possibility of having a very large number of accessible configurations in the state space (e.g. polymer systems with specific *MPC* function distributions). In those studies, authors suggest grouping similar states and representing each group with a single configuration and state vector before identifying the transitions between different states. Since system identification requires running simulations starting from each configuration, this grouping strategy can be very helpful for making the identification of the model computationally tractable.

After the system identification is completed, and the reduced order model that

describes the relationship between process conditions and the state is obtained, another challenge would be the validation of the model using experimental data. As described in Section 3.2.1, experimental tools for characterizing the molecular weight distribution (size exclusion chromatography) and the degree of branching (^1H NMR spectroscopy) are available. The reduced order model can be used to design process inputs that target polymeric structures with specified properties, and these inputs can also be used in experiments. Comparison of the experimental data and the reduced order model's predictions would give insight about the weaknesses of the model (e.g. regions of the state space of the polymerization process with low data density where model predictions are inaccurate). However, this procedure is not straightforward, and it would definitely involve an iterative process in which the model would be continually improved using new experimental data. Once the model is validated, it can be used for optimization of the process. One example would be the identification of the optimal monomer feed rate profile that involves minimum number of switches of the feed rate to synthesize a specified polymer product at a given monomer conversion.

In this section, we outlined a model reduction scheme that could be used to derive a reduced order process model from high dimensional KMC simulation data. In the later chapters of this thesis, we formalize this methodology and demonstrate it for a thin film deposition process.

3.4 Conclusions

In the first section of this chapter, formation of highly branched, segmented polyurethaneureas based on oligomeric $A_2 + B_3$ approach, where A_2 is slowly added onto B_3 , were investigated by experimental studies and kinetic Monte-Carlo simulations. SEC results clearly demonstrated the formation of high molecular weight segmented copolymers with very high polydispersity values, typical of highly branched polymers. When polymerization reactions are conducted in dilute solutions no gelation was observed

even at stoichiometric ratios of A_2/B_3 well beyond the theoretical gel point of 0.75. This is attributed to high degree of cyclization in dilute solutions causing the suppression of molecular weight. Results obtained from KMC simulations supported this hypothesis.

In the second section of this chapter, the formation of highly branched poly (ether ester)s by the melt condensation of an A_2 oligomer with a B_3 monomer has been studied. The simulations demonstrated that unequal reactivities can play an important role in the structure development of hyperbranched polymers, even when it has a little impact on the molecular weight. The results also indicate that the presence of end-capping reagents delays the gel point. However, the effect of end-capping agents also depends strongly on the ratios of the various monomers and their reactivity ratios. These results are motivating our further study of the role of end-capping reagents in the $A_2 + B_3$ system.

In the last section, we suggested a model reduction approach to convert the high dimensional KMC simulation model into a reduced order process model that can be used for process design and optimization to target specific polymeric structures. The suggested approach involves exploration of the state space using KMC simulations, identification of typical states, and the transitions between them. Challenges of this approach, such as describing the dynamic state in the simulations, and validation of the reduced order model, are also discussed. We conclude that pair correlation functions are very suitable for the state description as long as their high dimension can be reduced while preserving information. On the other hand, model development is very likely to involve a recursive process using new experimental data which is necessary to explore a wide range of polymeric structures.

As the final conclusion of this chapter, KMC simulations provide a tool for quantitatively assessing the effects of simple reaction mechanisms on molecular structure evolution, enabling the consideration of a broader range of mechanisms than with

analytical models. They relate process inputs to molecular structure and thus, can also enable the design of molecular structure via design of the process.

CHAPTER IV

MODEL REDUCTION OF STOCHASTIC MOLECULAR SIMULATIONS

The second application considered in this thesis is the epitaxial growth of gallium arsenide (GaAs). GaAs is deposited by various methods including ultra-high vacuum molecular beam epitaxy (MBE) and chemical vapor deposition. There are many advantages to using GaAs over silicon in transistors. GaAs has a higher saturated electron velocity and higher electron mobility, allowing devices to function at frequencies excess of 250 GHz. Also, GaAs devices generate less fundamental noise than silicon devices at ultrahigh radio frequencies. Having a direct band gap, GaAs can be used to emit light unlike silicon, which has an indirect band gap and is a poor light emitter.

MBE deposition of GaAs occurs with a rate of approximately one layer per second. Hence, the growth morphology in this process develops on the order of seconds. Even with a further increase in computer power, such a time scale will not be accessible to simulations using molecular dynamics (MD). Even though MD can capture atomic vibrations on the order of picoseconds, slower events like an atom overcoming an energy barrier and moving to a new site in the crystal lattice will be infrequent when this technique is used. However, the slower events will be dominant in terms of capturing the dynamics of the process. Therefore, we use Kinetic Monte Carlo (KMC) simulations to reach the time scales of these slower events. A KMC simulation model on a zincblende crystal structure has been developed to describe epitaxial growth of GaAs [54]. This model includes the rates (derived from experiments) of over one thousand possible events taking place on the surface. These events include adsorption,

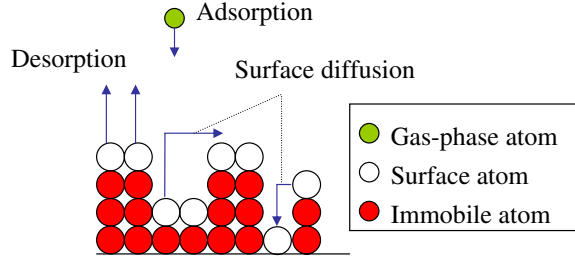


Figure 17: Types of events that take place during the MBE deposition of GaAs.

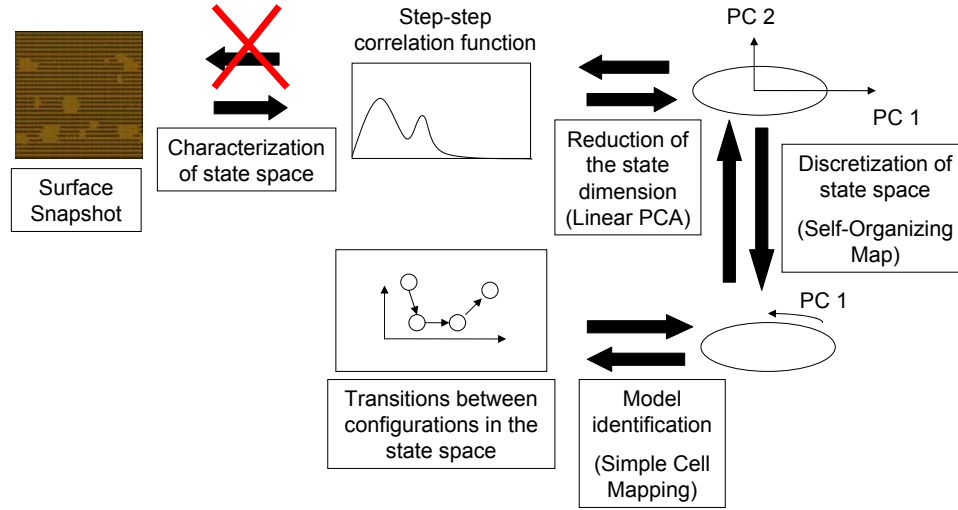


Figure 18: Schematic of the modeling approach.

desorption, and surface diffusion of gallium and arsenic species as shown in Figure 17.

In this chapter, we propose an algorithm for developing reduced order models from high dimensional stochastic simulation models. These molecular simulations possess a fine level of description of the process physics compared to the macroscopic (continuum) models. However, due to their high computational cost, it is not feasible to employ molecular simulations for optimization and control tasks. Hence, reduced order models are needed. Our modeling approach, which consists of four steps, is illustrated in Figure 18. Characterization of the state space, reduction of state dimension, discretization of the state space, and model identification gives us a reduced order dynamic model.

4.1 *Characterization of the state space*

The first critical step in our approach is to define the state space of the system, which is a non-trivial task for molecular simulations. When the state of the stochastic simulation is defined as the position of each atom in the system, each realization of the stochastic simulation gives a completely different state. However, control and optimization objectives are generally based on overall system properties that define the material structure. In order for a reduced order model to meet these objectives, overall system properties should not change significantly between realizations that are performed under the same nominal process inputs. Therefore, we consider symmetries in the stochastic simulations instead of the exact position of each atom in the system and use relative positions of each atom type (e.g. gallium or arsenide) to characterize the system state. This concept is commonly used in analyzing atomic scale simulations through the pair correlation function [3, 90]. Even though it is still a high dimensional description, this function captures the frequency of atom pairs that are certain distances apart from each other. When simulating the thin film deposition processes, a height-height correlation function can also be used to characterize the microscopic state. This function is a special type of pair correlation function based on the height of the surface atoms. It can be computed by taking the Fourier transform of the surface height for the purpose of removing symmetries. However, this function is extremely sensitive to noise in the KMC simulations. In our GaAs demonstration, step-step correlation (SSC) function is used to describe the microscopic state of the system during the simulations since it captures the discrete changes due to the location of atomic steps with a small number of modes. This approach is equivalent to taking the Fourier transform of the derivative of the surface height instead of the absolute value of it. The SSC function gives the distance and orientation between pairs of steps on the surface, where a step is defined as a change in height from one atomic surface site to the next.

As shown in Figure 19, we define the 2D coordinates (i, j) of each atom using two crystallographic directions [54], where the third coordinate k is the height of the atom. As illustrated in Figure 19, i , j and k increase by one when we move by one lattice unit in $[1\bar{1}0]$, $[\bar{1}10]$ and $[001]$ directions, respectively. Defining the surface height to be $\mathbf{h}(i, j)$ for the atomic site (i, j) , the presence of an up step in direction i can be computed as:

$$s_{u,i}(i, j) = \begin{cases} 1 & \text{if } h(i+1, j) > h(i, j); \\ 0 & \text{if } otherwise. \end{cases}$$

On the other hand, an up step in direction j is computed using the following relation:

$$s_{u,j}(i, j) = \begin{cases} 1 & \text{if } h(i, j+1) > h(i, j); \\ 0 & \text{if } otherwise. \end{cases}$$

Down steps are computed similarly by reversing the inequalities. Then, the SSC function in direction i is defined to be $SSC_{k,l,i}(r) = \sum_{j=1}^{n_i} s_{k,i}(i, j) s_{l,i}(i, j)$ where r is the distance between the steps, and k and l may each take values of u or d to denote the up and down steps, respectively. Since the occurrence of every type of step pair (up-up, up-down, down-up and down-down) is counted on the surface in four different directions i , $-i$, j and $-j$, there are sixteen types of functions in the whole SSC function of a surface snapshot. This spatial correlation function is high-dimensional and may contain redundant information. The data is also noisy. Noise is reduced by performing multiple realizations under identical conditions and averaging the results. After averaging, PCA is used to determine the number of independent variables needed to fully determine the SSC function. This technique is widely used to eliminate linear correlations among the variables in data sets. It is a crucial step in our study. The reduced state dimension makes it possible to construct a compact dynamic model.

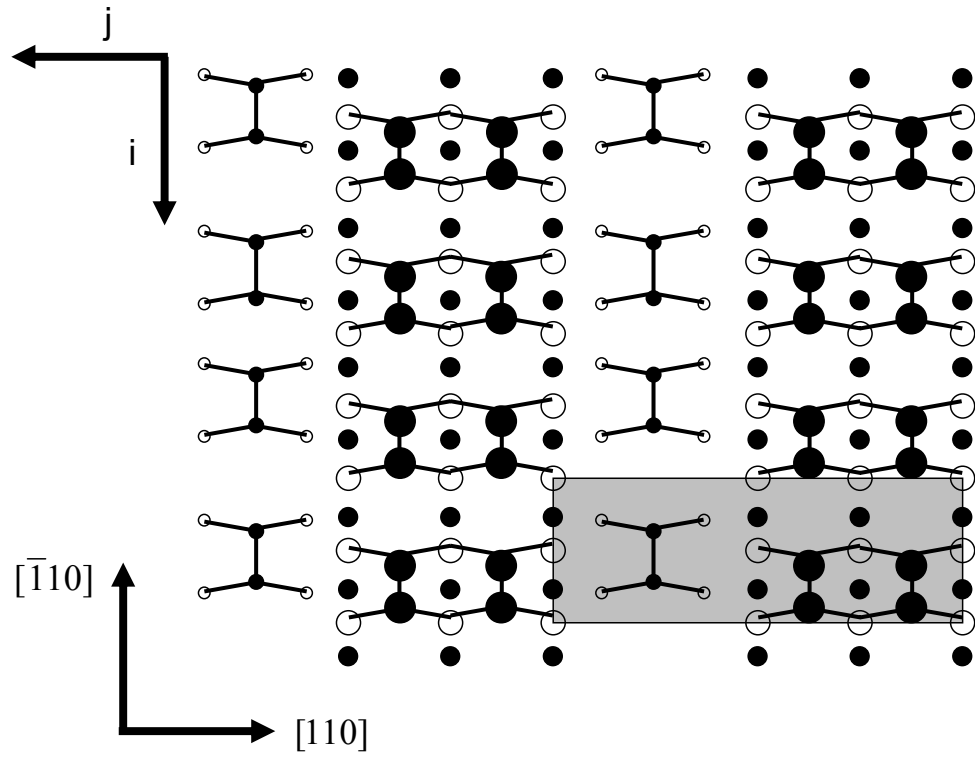


Figure 19: Relationship between the crystallographic directions and two dimensional coordinates used for constructing the lattice. Dashed region is the $\beta_2(2 \times 4)$ reconstruction of GaAs. Dark and white atoms are As and Ga atoms, respectively. The atomic radius decreases with increased depth.

The KMC simulations have been carried out using the kinetic barriers calculated by Itoh [54] with the following parameters:

- Growth temperature: 580°C
- Film deposition interval: 0.20 monolayers (ML)
- Incident arsenic dimer (As_2) flux: 0.4 ML/s
- Incident gallium (Ga) flux: Varied between 0.06-0.20 ML/s (0.06, 0.08, 0.10, ..., 0.20) flux range where the model is valid.
- Lattice size: 300x300 (90000 surface atoms) with periodic boundary conditions

The starting surface in our simulations is the thermodynamic ground state of GaAs(001), the $\beta 2(2 \times 4)$ reconstruction shown in Figure 19, which prevails in a wide range of growth conditions.

Snapshots of surfaces have been recorded at surface increments of 0.01 ML, starting from an initial surface in the $\beta 2(2 \times 4)$ configuration, up to 0.20 ML coverage. Eight constant input simulations were performed at Ga fluxes of 0.06 ML/s, 0.08 ML/s, 0.10 ML/s, ..., 0.20 ML/s and we will refer this as Training Simulation Set 1. In Training Simulation Set 2, we again have 8 simulations, but this time the flux is kept constant up to 0.10 ML coverage (middle of the deposition), and the flux is switched to a different value at that coverage point. To explore the KMC state space even further, we perform 60 additional simulations (Third Training Simulation Set), where two flux switches are made at 0.07 and 0.14 ML coverage points.

The SSC function has been computed as $\mathbf{S} \in \mathbb{R}^{d_S}$ for each snapshot. \mathbf{S} is obtained by combining the 16 different portions of the SSC function, which has been evaluated for four different types of step pairs in four different directions on the surface ($4 \times 4 = 16$ combinations). Each portion has 300 variables (one for each interatomic distance or lattice unit), which is equal to the lattice size. Therefore, d_S has a value of 4800

($300 \times 16 = 4800$). We have a total number of snapshots $n_s=1521$ (76 simulations, 20 snapshots for each simulation, plus the snapshot of the initial state of the system). \mathbf{S} of all snapshots in the training data are collected in $D \in \mathbb{R}^{4800 \times 1521}$. Before performing PCA, D is first transformed into $D' \in \mathbb{R}^{4800 \times 1521}$ with the following elements:

$$D'_{i,j} = D_{i,j} - D_{mean,i}. \quad (2)$$

Here, $D_{mean} \in \mathbb{R}^{4800}$ is defined as:

$$D_{mean,i} = \frac{\left(\sum_{j=1}^{n_s} D_{i,j}\right)}{n_s} \quad (i = 1, 2, \dots, 4800). \quad (3)$$

In order to complete the pre-processing of D , D' is transformed into $D'' \in \mathbb{R}^{4800 \times 1521}$.

$$D''_{i,j} = D'_{i,j}/D_{std,i} \quad (i = 1, 2, \dots, 4800, j = 1, 2, \dots, 1521) \quad (4)$$

$D_{std} \in \mathbb{R}^{4800}$ is defined as:

$$D_{std,i} = \sqrt{\frac{\sum_{j=1}^{n_s} (D_{i,j} - D_{mean,i})^2}{n_s}} \quad (5)$$

Hence, by pre-processing D , the variance of each variable in the SSC function is made equally important before PCA is carried out. In other words, small features with low variance will not be neglected.

4.2 *Reduction of the state dimension*

The dimensionality of a data set is defined as the number of variables which are used to describe each object (e.g. surface snapshot). However, there might be significant correlations between these variables. Principal component analysis (PCA) is a widely used method for reducing the dimensionality of a data set by eliminating linear correlations [79].

A principal component is a linear combination of the variables:

$$PC_i = \sum_{j=1}^v a_{i,j} b_j \quad (6)$$

where PC_i is the i^{th} principal component, $a_{i,j}$ is the coefficient of the variable b_j and v is the number of variables. The first principal component of a data set is the linear combination of variables which gives the best fit straight line when the data is plotted in the v -dimensional space.

In general, to represent all the variance in the data, all the principal components are needed. On the other hand, in many cases, a few principal components may be enough to explain a significant portion of the data variation. If the number of the principal components which can capture most of the variance of data is one or two, then a graphical representation is possible.

To compute the principal components, singular value decomposition can be performed on covariance of a data matrix B . When there are s number of observations and v number of variables, B matrix will have s columns and v rows. The covariance matrix of B is $Z = B^T B$ (if B is a zero-mean and unit-variance matrix). The first principal component is the eigenvector of Z with the largest eigenvalue, the second principal component corresponds to the second largest eigenvalue, and so on. The amount of variance captured by mode i is proportional to the i^{th} eigenvalue λ_i . In other words, principal component i accounts for $\lambda_i / \sum_{j=1}^v \lambda_j$ of the variance in the original data. Therefore, the first m principal components account for $\sum_{j=1}^m \lambda_j / \sum_{j=1}^v \lambda_j$ of the total variation in the data. Principal component analysis is straightforward to implement, but it is only able to find linear correlations. Nonlinear correlations may also exist within the data set.

In this study, PCA has been performed through the singular value decomposition of matrix D'' by computing the singular values of this matrix. The squares of the singular values correspond to the eigenvalues of the covariance matrix of D'' , and the

ratio of each eigenvalue to the sum of all eigenvalues (normalized eigenvalue) is plotted against the principal components. The point on the plot, where a sudden decay of the normalized eigenvalue is seen, gives the minimum number of principal components (the minimum dimension n) that can reconstruct the data effectively. At this point, we project each snapshot's SSC function onto the first n principal components D'' . We define $\mathbf{x} \in \mathbb{R}^n$ as the coefficient set with these new coordinates. Each \mathbf{x} characterizes a particular snapshot. As a result of PCA, small features in the SSC function, which do not contribute to the variance of the D'' significantly, are eliminated. We note that this could create a problem later while grouping similar surface structures according to first n principal components, because these small features may possess valuable information about the differences between the surface structures and their evolution in time.

4.3 Discretization of state space

We use self-organizing map (SOM) for the discretization of the state space. SOM is a neural-network model and algorithm that is widely employed for visualization of high dimensional data. Identification and monitoring of complex process states, which are sometimes very hard to interpret and analyze, are among the most common engineering applications of the SOM [66].

Let's assume that we have some data where each data sample or input vector is described by an n -dimensional $\mathbf{y}(t)$ where t is the sample index. In order to map the data onto a two dimensional array, which can represent the data in two dimensions, SOM algorithm can be used. Each node in SOM contains a model vector \mathbf{m}_i which is the same size as each $\mathbf{y}(t)$. Initial values of the model vector components can be assigned randomly or along a two-dimensional subspace spanned by the two principal eigenvectors of the input data [67]. Each input vector $\mathbf{y}(t)$ is mapped onto a particular node represented by \mathbf{m}_i , which matches best with $\mathbf{y}(t)$. The matching of an input

vector with a map unit is based on some metric (e.g. Euclidean distance between $\mathbf{y}(t)$ and \mathbf{m}_i). The SOM algorithm is made up of repeating the following steps:

- An input vector $\mathbf{y}(t)$ is compared with all the model vectors \mathbf{m}_i , and the model vector that matches $\mathbf{y}(t)$ best (minimum Euclidean distance) is selected as the best matching unit. This unit is also called the winner. The input vector for which \mathbf{m}_i is the winner, is called a “hit” for \mathbf{m}_i .
- The model vector of the winner node and a number of its neighbors are changed towards the input vector $\mathbf{y}(t)$. The following relation is used to make that update:

$$\mathbf{m}_i(k+1) = \begin{cases} \mathbf{m}_i(k) + \alpha(k)[\mathbf{y}(t) - \mathbf{m}_i(k)] & \text{if } i \in N_c(k); \\ \mathbf{m}_i(k) & \text{otherwise.} \end{cases}$$

where k is the discrete time index of the model vectors, $\alpha \in [0, 1]$ is a scalar that defines the learning rate and $N_c(k)$ specifies the radius of the neighborhood of the winner node on the map, which will be updated once $\mathbf{y}(t)$'s best matching unit \mathbf{m}_i , is found. The goal of the SOM learning process is that, for each sample input vector $\mathbf{y}(t)$, the winner and the nodes in its neighborhood are changed closer to $\mathbf{y}(t)$. At the beginning of this learning process, neighborhood radius is generally taken as a large value, and it shrinks during the process. Using this approach, global order is obtained at the beginning. Towards the end, radius gets smaller and local corrections of the model vectors are made. Also, the rate $\alpha(k)$ decreases as the map evolves.

After performing PCA and obtaining a coefficient set to characterize each surface snapshot, we use SOM to eliminate some of the nonlinear correlations within the data set by employing MATLAB's SOM Toolbox. Snapshots with very similar microstructures are grouped by SOM, and snapshots in the same group are viewed as equivalent when identifying the dynamic model in the next step. The PCA data is discretized

as a part of this step, since snapshots are grouped among the nodes of SOM. As a result, the number of discrete states is finite. The computational load for the model identification is also reduced with the grouping achieved by SOM. Because, during model identification, rather than the transitions between each surface structure, only the transitions between the groups are computed.

Various training procedures for SOM are described elsewhere [65]. In this study, prototype vectors are initialized along the first two principal components of the training data. Also, SOM training is accomplished by sequential, rather than the batch-wise comparison of snapshots to the overall map, and then updating the prototype vectors $\mathbf{m} \in \mathbb{R}^n$ of the map nodes to match and organize the snapshots. Once the map is trained, each snapshot is associated with the node that has the closest prototype vector, as measured by the Euclidean distance. The default number of nodes in the map is computed by the SOM toolbox using a heuristic formula which is a function of the number of data points (snapshots) in the data matrix. The map size can also be increased to provide a finer discretization of the state space. Another important SOM parameter, the ratio of the side lengths of the map, is set equal to the ratio of the two largest eigenvalues of the data matrix.

4.3.1 Results

In this study, PCA is used to find the minimum dimension that can represent the microscopic state of the surface snapshots recorded during the simulations. This method consists of computing the variance captured by each principal component of the entire data set and selecting the most important principal components that can reconstruct the data well. Figure 20 shows the normalized eigenvalues of each principal component. A knee shape is observed after the second principal component and the first two principal components capture nearly all of the variance within the data set.

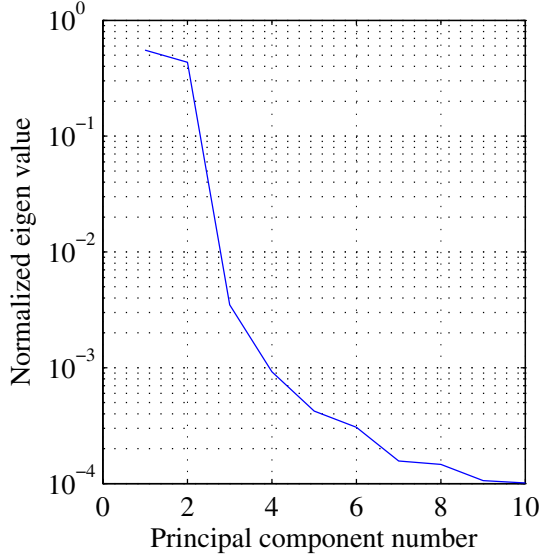
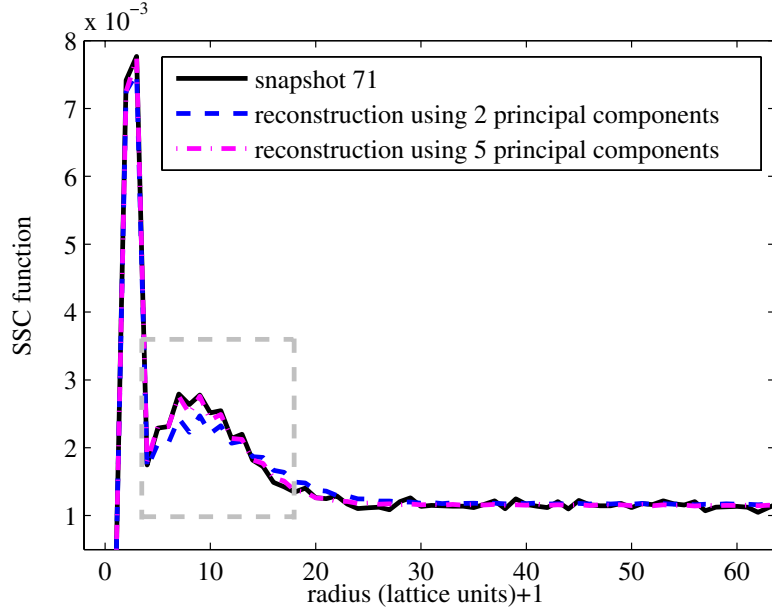


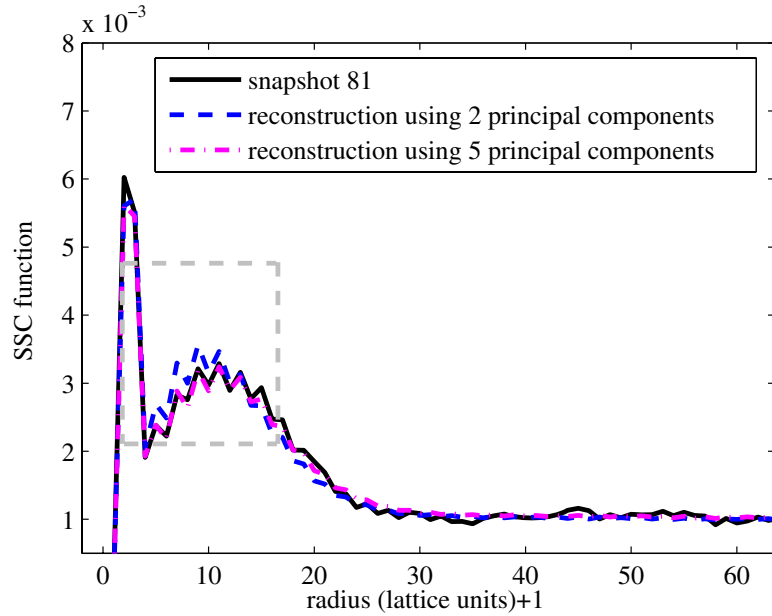
Figure 20: Normalized eigenvalues versus principal components.

However, data reconstructions with 2 and even more principal components showed us that at least 5 components are needed in order to effectively represent our data. Specifically, 5 principal components reconstructed the small clusters of atoms on the surface (with a size of less than 20 lattice units) much better than 2 principal components. Figure 21 illustrates this comparison for two different snapshots. Here, $SSC_{up,down,i}$ at a radius of 1 represents the adatom density, where as $SSC_{up,down,i}$ at a radius of 2 is the dimer density. The oscillations in the size range of 5 to 15 lattice units represent the density of islands with that size range. In the longer range, correlations between the atoms are much lower because of the clustering due to diffusion. The rise in the $SSC_{up,down,i}$ near a radius of 300 is due to the relatively high density of voids compared to the other longer range correlations (step pairs) that have smaller distances between them.

As already defined in Section 2.2, the coefficient set \mathbf{x} for each surface snapshot is obtained by taking the inner product of the snapshot's SSC function with the $n = 5$ principal components of the data matrix D'' . This reduced representation, which includes the coefficient sets of every snapshot, is used to train an SOM. For training,



(a) This surface is obtained under 0.12 ML/s Ga flux and has a 0.10 ML surface coverage.



(b) This surface is obtained under 0.12 ML/s Ga flux and has a 0.20 ML surface coverage.

Figure 21: Comparison of the reconstructions of $SSC_{up,down,i}$ with the original data using 2 and 5 modes. The region enclosed with the gray dashed line represents the surface structures with a size of less than 20 lattice units.

we used the default map size, which is a function of the number of surface snapshots in the training data. Each SOM node has its own prototype vector having the same dimension as each snapshot in the data set ($n = 5$). As a result, each map node also has its own SSC function \mathbf{S} , which is the PCA reconstruction of its prototype vector \mathbf{v} . Thus, each map node is associated with a type of surface.

The quality of the resulting map can be evaluated by calculating the average quantization error over the input data, given by $\|\mathbf{x} - \mathbf{v}_{BMU}\|_2$, where \mathbf{v}_{BMU} is the prototype vector for the best matching unit, for the snapshot with the projection vector \mathbf{x} [65]. In words, the Euclidean distance between each surface snapshot and its best matching SOM node is computed, and the average of these distances gives the average quantization error.

In order to see the effects of training data on the map quality, we generated multiple maps while keeping the map size at 192 nodes. For each map, the training phase has taken approximately 3 minutes using a Pentium 4 processor with a speed of 3 GHz. The first map SOM1 is trained with 161 surface snapshots (from Training Simulation Set 1), the second map SOM2 is trained with 321 surface snapshots (from Training Simulation Sets 1 and 2) and SOM3 is trained using all of the 1521 surface snapshots (from Training Simulation Sets 1, 2 and 3). Table 1 shows the statistics of three maps with 192 nodes, which were trained using these three different data sets. SOM1 and SOM2 have similar statistics, except that SOM2, which has been trained with a larger data set, has a larger number of surface snapshot groups. This suggests that Training Simulation Set 2, where the flux was changed in the middle of the simulations at 0.10 ML surface coverage, enabled us to see surface structures that had not been accessible through constant input simulations (Training Simulation Set 1). On Table 1, as we move from SOM2 to SOM3, we again observe an increase in the number of surface snapshot groups. Therefore, we conclude that Training Simulation Set 3 (during which the flux has been changed twice during the simulations) explored

the KMC state space even further, producing additional surface structures compared with Training Simulation Sets 1 and 2. SOM3 has an average quantization error of 5.1, whereas the average prototype vector size is 53.0. If needed, the average quantization error can be reduced by adding more nodes to the map, but we were able to obtain a useful model with SOM3.

Table 1: Statistics of three SOMs trained using different data sets.

| Map | Training sim. sets | Quant. error | Proto. vec. size | Topog. error | Number of groups |
|------|--------------------|--------------|------------------|--------------|------------------|
| SOM1 | 1 | 4.1 | 54.1 | 0.01 | 119 |
| SOM2 | 1, 2 | 4.8 | 54.4 | 0.02 | 142 |
| SOM3 | 1, 2, 3 | 5.1 | 53.0 | 0.01 | 175 |

Another SOM metric is the topographic error, which is the proportion of all the data vectors for which first and second best matching units are not adjacent on the map [65]. SOM3 also has a topographic error of 0.01, which means that for only 15 of the 1521 snapshots, the first and second best matching units were not adjacent on the SOM. Topology preservation is not required for a good state representation, but it aids in visualization of the dynamics on the map.

To illustrate the clustering of nodes, which are represented by prototype vectors, a graphic display called unified distance matrix (U-matrix) is used [116]. Let’s say we have a 3×1 sized map represented by the prototype vectors $\mathbf{m}(1)$, $\mathbf{m}(2)$ and $\mathbf{m}(3)$. The U-matrix will be a 5×1 vector with the following elements: $\mathbf{u}(1)$, $\mathbf{u}(1, 2)$, $\mathbf{u}(2)$, $\mathbf{u}(2, 3)$, $\mathbf{u}(3)$ as shown in Figure (22).

In Figure 22, $\mathbf{u}(i, j)$ is the distance between prototype vectors $\mathbf{m}(i)$ and $\mathbf{m}(j)$ and $\mathbf{u}(k)$ is the mean of the surrounding values, e.g. $\mathbf{u}(2) = (\mathbf{u}(1, 2) + \mathbf{u}(2, 3))/2$. In a real U-matrix with a color index, each U-matrix unit $u(i)$ or $\mathbf{u}(i, j)$ would have a color according to its magnitude, here all map units have the same color which means $\mathbf{u}(1) = \mathbf{u}(1, 2) = \mathbf{u}(2) = \mathbf{u}(2, 3) = \mathbf{u}(3)$. The U-matrix shows graphically

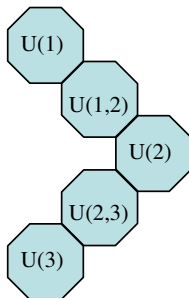


Figure 22: U-matrix for a 3×1 SOM.

which neighboring nodes are similar or different with respect to the prototype vectors representing them. This helps in the visualization of data clustering.

Figure 23 shows the film coverage levels associated with the SOM nodes. As a reminder, film coverage levels of surface snapshots are recorded at increments of 0.01 monolayers (ML), starting from the initial surface configuration at 0 ML, up to 0.20 ML. In Figure 23, 17 map nodes (out of 192 nodes) are empty since they were not the best matching units for any surface snapshot. Also, the film coverage increases as we move to the right portion of the map. As an approximation, it can be said that each column of the map corresponds to a specific film coverage level. On the other hand, Figures 24 and 25, which show the trajectories of simulations with constant flux at the minimum and maximum settings, indicate that each row of the map is associated with a constant flux simulation trajectory.

Figures 26 and 27 show SOM1 and SOM3. SOM1, which is trained with only the constant flux simulation data, has a 20 % narrower color index range in its U-matrix ($[2.21, 45.3]$ for SOM1 versus $[3.24, 56.9]$ for SOM3). This is expected since varying flux simulations generated a wider range of surface configurations leading to a wider color index for the U-matrix of SOM3. U-matrices for both maps indicate that both maps are quite uniform in the sense that neighboring map nodes have quite similar prototype vectors since the blue color, which corresponds to low Euclidean distances between the neighboring prototype vectors, dominates in both figures. Low

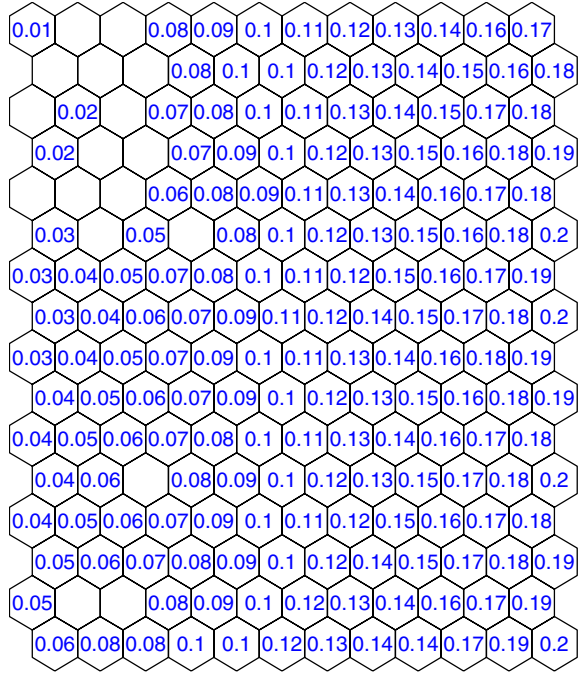


Figure 23: Values of film surface coverage (in monolayers) for each node of SOM3.

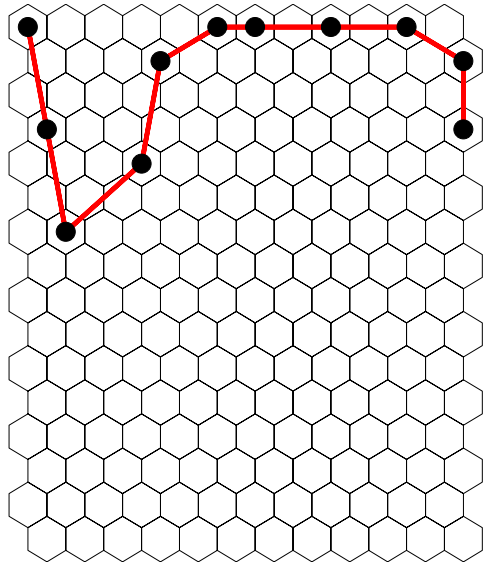


Figure 24: Trajectory of a training simulation on SOM3. Ga flux is kept constant at the minimum flux (0.06 ML/s) during this simulation.

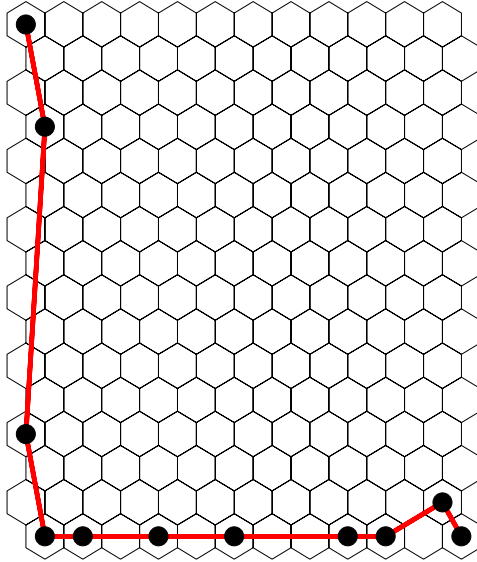


Figure 25: Trajectory of a training simulation on SOM3. Ga flux is kept constant at the maximum flux (0.20 ML/s) during this simulation.

topographic error for these maps (Table 1) previously supported this conclusion.

SOM1 and SOM3 also include 5 component planes each of which represent an element of prototype vector $\mathbf{m} \in \mathbb{R}^5$. The component planes shown under the titles of “Variable1” to “Variable5” are the values of the first to the fifth elements of the prototype vectors. As seen on the maps (Figures 26 and 27) and according to the Figures 24 and 25, the first element of the prototype vectors roughly change with the flux variable as we move from the top of the maps to their bottom portions. On the other hand, second element of these vectors are analogous to film coverage value that changes from left portion to the right portion of the map, but do not change as we move vertically on the map. Figure 23 previously supported this conclusion. The remaining three elements of the prototype vectors (Variable3, Variable4 and Variable5) have much lower values than the first two variables since the first two principal components were able to capture more than 99% of variance in the simulation data. Hence, the grouping of snapshots is mainly determined by the first two variables. As a result, using SOM, we are able to visualize the system dynamics in a 2-dimensional

fashion, where two major axes of the map correspond to gallium flux and surface coverage variables (or the portions of the training data vectors on the first two principal components).

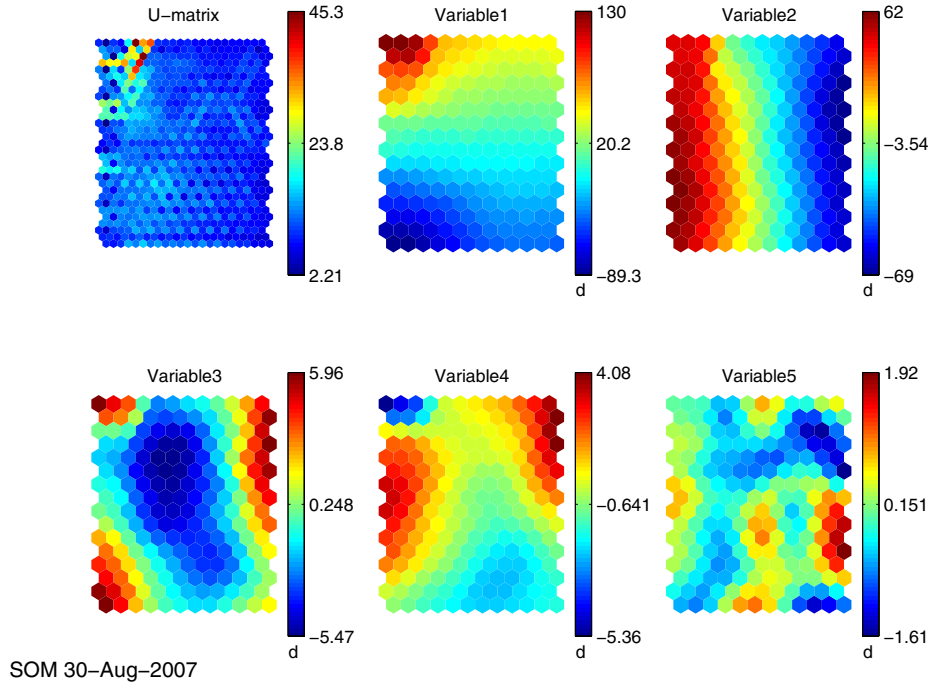


Figure 26: SOM1 (generated using the Training Simulation Set 1) and its U-matrix.

At this point, three SOMs have been generated. Only one of them, SOM3, will be used to identify a dynamic model, since it has been trained with the most extensive amount of training data.

4.4 Model Identification

4.4.1 Simple Cell Mapping

Characterization of the surface snapshots from all of the training simulations gives the system’s state space. As explained in the previous section, the surface snapshots are grouped into discrete states using SOM, according to their similarities in terms of surface morphology. These groups represent the cells, and simple cell mapping (SCM) or generalized cell mapping (GCM) can be used to obtain a global view of

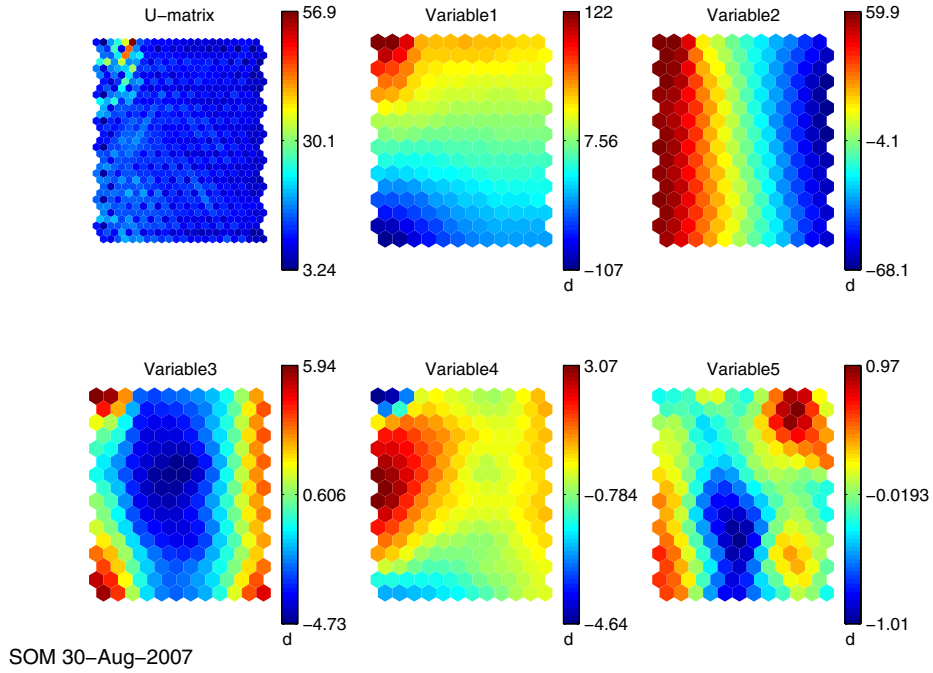


Figure 27: SOM3 (generated using the Training Simulation Sets 1, 2 and 3) and its U-matrix.

the behavior of the system [52]. A cell map is formed by dividing the state space into a finite number of cells (using SOM) and approximating the behavior of the system by means of transitions between the cells. SCM is a deterministic approach in which each cell is mapped into exactly one other cell for a particular input. In GCM, each cell can have several image cells. In other words, a cell can be mapped to several other cells. Using this stochastic approach, we can assign a probability of the system being in a cell at a specific time and extract the dynamic properties. In the current study, we implement the SCM approach to identify the dynamic model for the process, because SCM provides a deterministic model to describe the evolution of the surface structure under different material flux profiles. The following are the steps used for model identification:

- One surface configuration is selected from a particular map node.
- A new simulation is started from that surface configuration and run for an

incremental coverage interval with one of the flux settings.

- The final surface structure is recorded and the configuration group, which has the closest microstructure to that final structure, is obtained. In other words, the flux-dependent transition (from one configuration group to another) is found.

When the above procedure, also illustrated in Figure 28, is repeated for every map node and each flux setting, we identify a reduced order dynamic model. The identification procedure requires running $146 \times 8 = 1168$ KMC simulations (for 146 map nodes and 8 flux settings), which takes approximately 7.3 days using a computer cluster made up of 16 computers, each having an Intel Xeon processor with a speed of 2.66 GHz.

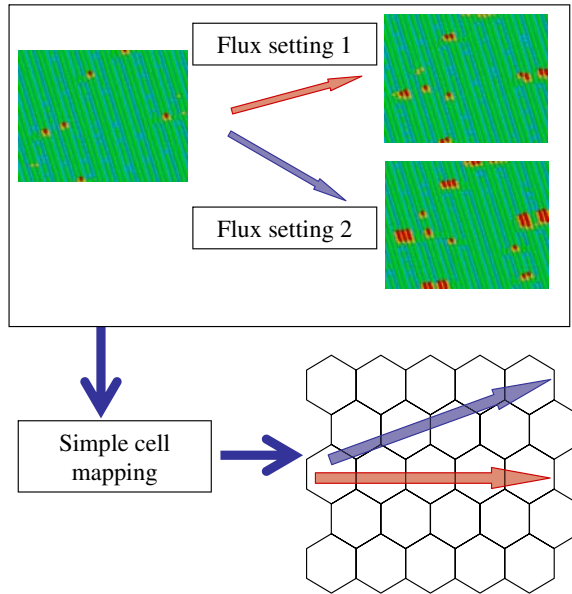


Figure 28: Schematic of simple cell mapping.

The reduced order model is described as:

$$\mathbf{p}[j + 1] = A(u[j])\mathbf{p}[j] \quad (7)$$

where j is the time step (or coverage level) number, $u[j] \in \mathbb{R}$ is the input (or flux) value, $\mathbf{p} \in \mathbb{R}^m$ is the probability vector describing the configuration group that the system is currently in and $A \in \mathbb{R}^{m \times m}$ is the transition matrix. The value of m (dimension of \mathbf{p}) is equal to the number of configuration groups in the state space (the number of cells in the cell map). In our model, there are 149 surface configurations. 146 of the number of initial configurations (with coverage values of 0.02 to 0.18 ML coverages) used for the model identification, and the remaining 3 configurations are associated only with 0.20 ML surface coverage. These 3 states are assumed to transition into themselves since 0.20 ML is the highest possible surface coverage value. We define \mathbf{p} as:

$$\mathbf{p}_i = \begin{cases} 1 & \text{if } i = l; \\ 0 & \text{if } i \neq l. \end{cases}$$

Here, l is the number of the current system configuration. The transition matrix A is a function of the input. Therefore, in the reduced order model, we have eight transition matrices for eight flux settings. A is a highly sparse matrix. The sum of the elements in each column of A is equal to one, because for each column, only one element of A is non-zero, whereas the other elements are equal to zero. In other words, there is only one possible transition from each configuration group under a particular input. This construction of A is due to the deterministic nature of the simple cell map. If the system makes a transition from configuration s to t at time step j with the flux setting $u[j]$, this transition can be represented by a transition matrix with the following properties:

$$A_{k,s} = \begin{cases} 1 & \text{if } k = t; \\ 0 & \text{if } k \neq t. \end{cases}$$

At time step j , the surface properties of the system are given by $\mathbf{x}[j]$ (set of projection coefficients defined in Section 4.2):

$$\mathbf{x}[j] = X\mathbf{p}[j] \quad (8)$$

where $X \in \mathbb{R}^{n \times m}$ has the $\mathbf{v} \in \mathbb{R}^n$ (prototype vector) of each configuration group (SOM node). In other words, \mathbf{v} of the configuration group i is $\mathbf{v}_i = \{X_{1,i}, X_{2,i}, X_{3,i}, X_{4,i}, X_{5,i}\}$.

4.4.2 K-nearest neighbor algorithm

Cell mapping can be generalized to enable interpolation between cells. In order to carry out the interpolation, we use the k-nearest neighbor algorithm [80] and predict the evolution of \mathbf{x} for a given flux profile. When $k=1$, this algorithm produces results identical to simple cell mapping and no interpolation is made. The following steps represent the k-nearest neighbor algorithm when $k=2$:

1. Let \mathbf{x}_{old} be the vector representing the surface state at an initial film coverage. This vector is compared with all the prototype vectors on the SOM and its best matching unit (BMU) and second BMU are found based on the Euclidean distances between these vectors. Let these two SOM nodes have prototype vectors \mathbf{v}_1 and \mathbf{v}_2 , where \mathbf{v}_1 and \mathbf{v}_2 represent the first BMU and second BMU, respectively.
2. Compute the distances $d_1 = \|\mathbf{x}_{old} - \mathbf{v}_1\|_2$ and $d_2 = \|\mathbf{x}_{old} - \mathbf{v}_2\|_2$.
3. Compute the weights associated with these two distance values $w_1 = (1/d_1)/(1/d_1 + 1/d_2)$ and $w_2 = (1/d_2)/(1/d_1 + 1/d_2)$.
4. Predict the system state at the next step from $\mathbf{x}_{new} = w_1 \cdot \mathbf{v}_3 + w_2 \cdot \mathbf{v}_4$, where \mathbf{x}_{new} is the state vector prediction for the next step, \mathbf{v}_3 and \mathbf{v}_4 are the prototype vectors of the nodes determined from cell mapping results (\mathbf{v}_1 transitions into \mathbf{v}_3 , and \mathbf{v}_2 transitions into \mathbf{v}_4 under the particular flux setting).
5. Set $\mathbf{x}_{old} = \mathbf{x}_{new}$ and repeat steps 2-4 to make predictions at higher film coverage values. At each coverage value, the predicted \mathbf{x} can be converted to the high

dimensional SSC function by multiplying the elements of \mathbf{x} with the eigenvectors obtained from the training data D'' .

The algorithm described above is the k-nearest neighbor algorithm with $k=2$. Here, the interpolation between prototype vectors provides the opportunity to improve the accuracy of the prediction. On the other hand, when $k=1$, k-nearest neighbor algorithm is identical to SCM. In that case, the second BMU is not identified and the weights are not computed.

4.5 Conclusions

In this chapter, we proposed a novel model reduction algorithm that makes use of the high dimensional KMC simulation data to derive reduced order process models. The proposed algorithm is used on high dimensional KMC simulations of epitaxial GaAs thin film deposition. This algorithm consists of three steps: applying principal component analysis (PCA), self organizing map (SOM) and simple cell mapping (SCM) to identify a dynamic process model. First, a spatial correlation function is used to describe the state space of the system by the characterization of surface snapshots generated under different material flux profiles. Then, a minimum state dimension that can represent the system state is found using PCA. After this reduction step, SOM is employed to group similar surface structures. SOM results led to the identification of a dynamic model through the computation of flux-dependent transitions between surface configuration groups with simple cell mapping (SCM). This dynamic model is evaluated in the next chapter.

CHAPTER V

PERFORMANCE EVALUATION OF THE DYNAMIC MODEL

In model reduction studies, it is important to derive error bounds originating from the approximations associated with the order reduction (from high order KMC simulations to reduced order dynamic models). In this chapter, we address this issue by evaluating the predictive ability of our reduced order model. We also use our model for dynamic optimization of the thin film deposition process.

5.1 Local Error Quantification

The core hypothesis of our model reduction approach is that the surface structures in the same configuration group (SOM node) should show similar dynamic behavior under identical process input (material flux). SCM, which is used to extract a dynamic process model in this study, is also useful for testing this hypothesis. We define the cell mapping error (*CME*), which is illustrated in Figure 29, to quantify the different dynamic behaviors of surface structures which belong to the same configuration group. The following is the procedure we used to find *CME*:

1. Randomly, select one surface structure from a map node and identify where this structure is mapped on the SOM under a particular flux setting (first cell mapping).
2. Randomly, select another surface structure from the same map node and run another simulation starting with this new surface structure and perform SCM again (second cell mapping).

3. Compute $CME = \|S_1 - S_2\|_2 / (\|S_1 + S_2\|_2 / 2)$, where S_1 and S_2 are the SSC functions of the SOM nodes coming from the first and second cell mapping, respectively. It should be noted that these functions are reconstructed from the prototype vectors of the map nodes.
4. Repeat steps 1-3 for all SOM nodes under all flux settings.

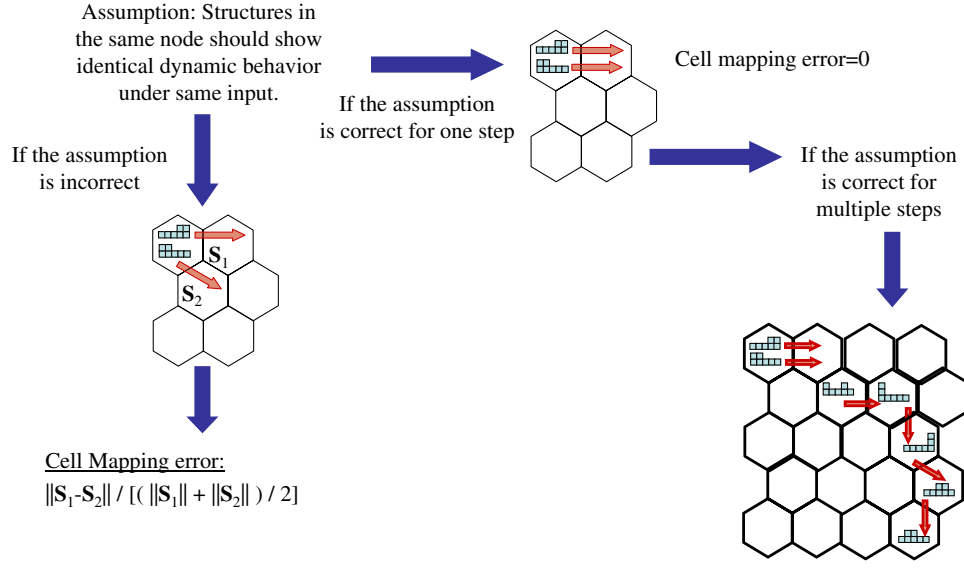


Figure 29: Computation of the cell mapping error.

CME is computed by picking two different surface structures from each node and performing SCM under each of the eight flux settings. In order to characterize the distribution of CME , its cumulative distribution function is computed as illustrated in Figure 30. This function is defined as $CDF(y) = P(CME \leq y)$ and gives the value of the probability that $CME \leq y$ for a given y . Figure 31 shows the CDF of the CME . 52% of the mappings are identical, with both surfaces evolving to the same map node. In those cases, $CME = 0$. Also, with a 0.90 probability $CME \leq 0.0075$, supporting the fact that there is a very strong chance for surface structures in the same group to show similar dynamic behavior. In order to reduce CME , a larger SOM can be used. The trade-off is that this would increase the dimension of the cell map and the computational load associated with the system identification step.

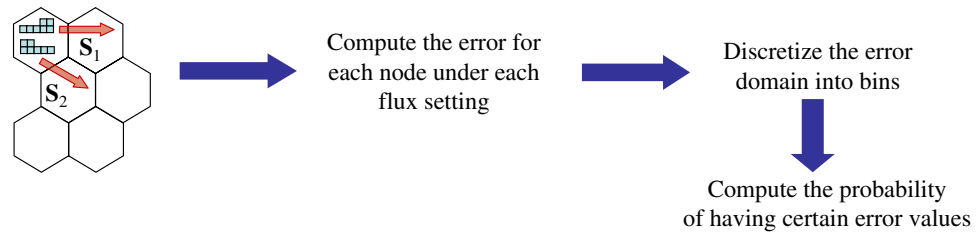


Figure 30: Computation of the distribution of the cell mapping error.

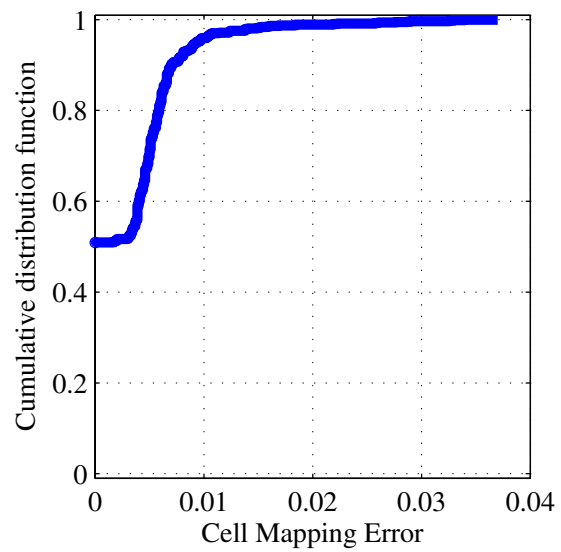


Figure 31: Cumulative distribution function of the cell mapping error.

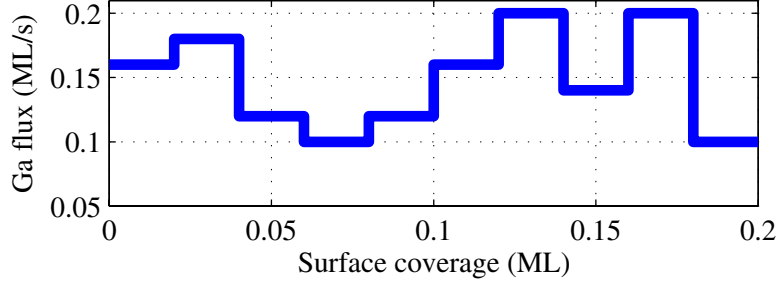


Figure 32: Flux profile of a test simulation.

5.2 Global Error Quantification

In Chapter 4, a dynamic model is constructed by finding the flux dependent transitions between the nodes of SOM3 under all flux settings. In order to test the predictive ability of our model, 1210 test simulations are performed. In these simulations, we split the film coverage domain (from 0 to 0.20 ML) into 10 equal intervals. The gallium flux value is kept constant in each interval. However, its value randomly changes when moving from one coverage interval to the next. This strategy provides very different flux profiles than the ones used to generate the training data, which involved a maximum of two flux switches. Figure 32 shows the flux profile for one of these simulations during which the flux is randomly changed at each step.

As described in Chapter 3, surface snapshots from the test simulations at different coverage levels (10 surface snapshots from 0.02 ML to 0.20 ML film coverage) are characterized using SSC functions, and the state vector (or coefficient set) $\mathbf{x} \in \mathbb{R}^n$ of each surface snapshot is computed. Then, each \mathbf{x} is matched with an SOM node based on the criteria of minimum Euclidean distance between \mathbf{x} and the prototype vectors of the SOM nodes. In other words, the best matching unit for each \mathbf{x} is sought.

For the test simulation with the flux profile shown in Figure 32, the predicted trajectory and the KMC simulation trajectory are given in Figure 33. Here, the prediction comes from cell mapping or k-nearest neighbor algorithm (with $k=1$), and each hexagon represents an SOM node corresponding to a surface structure group.

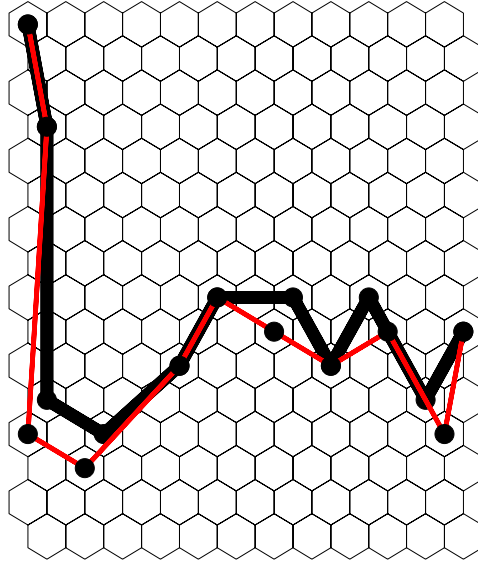


Figure 33: Trajectories of the KMC test simulation (red line) and the prediction (black line) on SOM3.

The simulation follows a path starting from the SOM Node 1 (initial surface structure) and moves to the right hand side of the SOM as the film coverage increases. Since each simulation starts from SOM Node 1 and the transitions from each node under all flux settings are known, the prediction of the simulation trajectory (using the dynamic model) is straightforward. Figure 33 indicates that there is an agreement between the predicted trajectory and the trajectory of the KMC simulation. Similar results were obtained for all test simulations. In some instances, these trajectories passed through neighboring SOM nodes. However, as Table 1 indicates, SOM3 has a very low topographic error. Hence, neighboring SOM nodes are similar, so slight differences between these trajectories do not jeopardize the accuracy of the state predictions.

SOM Node 186 represents the predicted film structure at the final coverage value (0.20 ML). Figure 37 shows the SSC function reconstructed from the prototype vector of this SOM node. The reconstruction agrees with the KMC simulation's SSC function. Hence, our dynamic model is capable of predicting the final film structure quite well. This figure also illustrates that the noise in the KMC simulation data is considerably reduced when the SSC function is reconstructed from the prototype

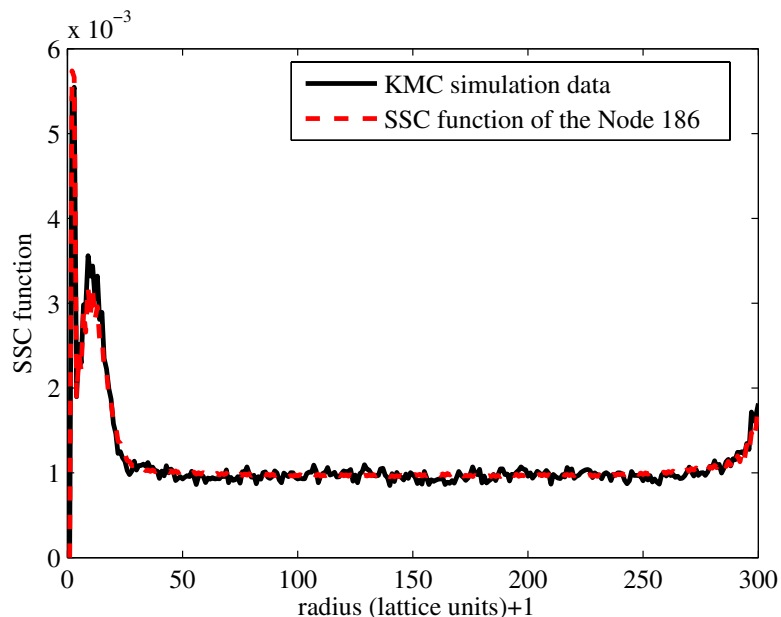


Figure 34: Reconstruction of the SSC function with the prototype vector of the SOM Node 186 and the original KMC simulation data at final film coverage.

vector.

We also tested the consistency of the dynamic model with the training data by reproducing the trajectories of simulations in this data set. From the results, we concluded that the dynamic model was able to reproduce trajectories of the simulations in Training Simulation Sets 1, 2 and 3 well. Figures 35 and 36 shows the estimated and real trajectories of two constant flux simulations (with the minimum and maximum flux settings 0.06 ML/s and 0.20 ML/s) supporting our conclusion. These figures also show that rougher surfaces (produced with higher material flux) are collected at the bottom of the SOM. On the contrary, lower flux simulations follow trajectories through upper regions of the SOM. With the lower material flux, there are more diffusion events per an adsorption event. Therefore, more diffusion events in the lower flux conditions lead to smoother surfaces compared to higher flux conditions with less number of diffusion events per an adsorption event.

In order to evaluate the performance of the reduced order dynamic model, we quantify the prediction error for the test simulations, which were not included in

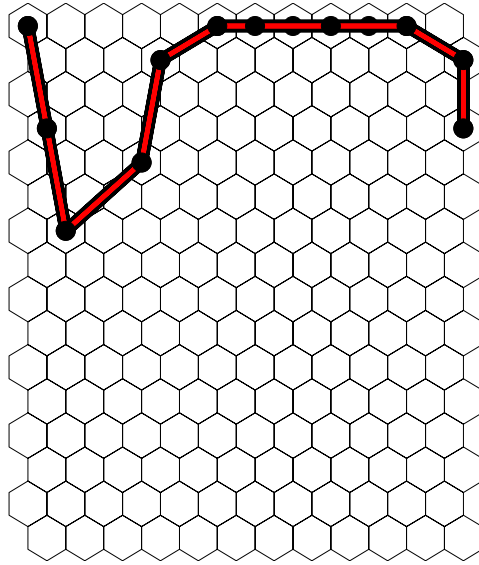


Figure 35: Trajectories of a training simulation (red line) and the prediction (black line) on SOM3. Ga flux is kept constant at 0.06 ML/s during this simulation.

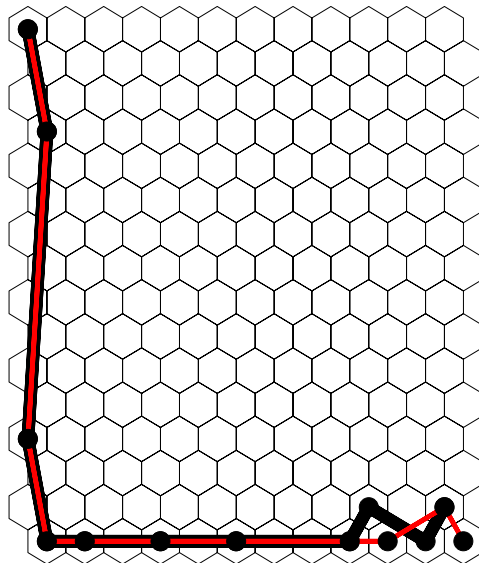


Figure 36: Trajectories of a training simulation (red line) and the prediction (black line) on SOM3. Ga flux is kept constant at 0.20 ML/s during this simulation.

the training data. The global (multi-step) error associated with the predicted SSC function is defined as:

$$E_{SSC} = \|\mathbf{S}_s - \mathbf{S}_p\|_2 / \|\mathbf{S}_s\|_2, \quad (9)$$

where \mathbf{S}_s is computed from the KMC simulation data and \mathbf{S}_p is reconstructed from the predicted state vector \mathbf{v} (prototype vector of the SOM node). The values of E_{SSC} at different film coverage levels (for the test simulation with the flux profile in Figure 32) are less than 0.018 according to Figure 37. This figure compares the error in the predictions made using the k-nearest neighbor algorithm with different k values. For this particular test simulation, predictions are less accurate with k=2. Because, for that simulation, the best matching unit of the predicted state shows a different dynamic behavior than the second best matching unit. This suggests a high cell mapping error for those map units. However, for more than 90% of the test simulations, predictions did not change when k value was increased from 1 to 2. This is because, at low film coverage levels, different surface configuration groups (corresponding to different SOM nodes) do not possess significantly different surface features (or different \mathbf{x}) and they are mapped to the same SOM node at subsequent coverage levels. Therefore, throughout the rest of this section, we only report the results obtained with k=1, which is identical to simple cell mapping.

As a part of the global error quantification, we also computed E_{SSC} at nine evenly distributed film coverage values (0.04 ML, 0.06 ML, ..., 0.20 ML) for three large test simulation sets (with random input profiles), none of which had been in the SOM training data. It should be noted that the SOM node 1, which represents the initial film structure at 0 ML, maps to the same SOM node at 0.02 ML with all flux settings. Hence, the prediction error at 0.02 ML is not computed. The three test simulation sets include 300, 600 and 1210 simulations, respectively. Figure 38 compares the mean values of E_{SSC} at different film coverage levels for these three test simulation sets and

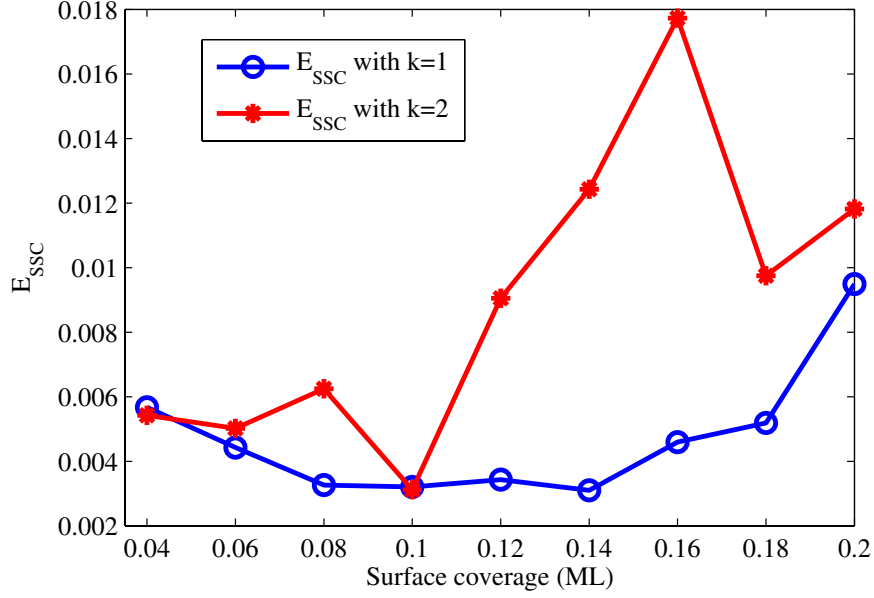


Figure 37: The evolution of E_{SSC} for the test simulation with $k=1$ and $k=2$.

also the Training Simulation Set 3. The mean E_{SSC} is below 0.012 at all coverage levels. It is interesting to note that the dynamic model was able to predict the final film structure in the training simulations with less error than the test simulations. This is because of the fact that the data from the training simulations were used as an input in the system identification, whereas our dynamic model was not as familiar with some of the surface structures produced during the test simulations, especially at high film coverage.

As shown in Table 2, the mean values of E_{SSC} for the predictions associated with the test simulation sets are around 0.006, but the standard deviation values are comparable to the mean values indicating a wide distribution of the prediction error. In order to get a more clear idea about the distribution of E_{SSC} , its cumulative distribution function (CDF) is computed for the three test simulation sets and also the Training Simulation Set 3. Figure 39 shows that the CDF curves of the test simulation sets are very similar and enlarging the size of the simulation set does not change the distribution of the prediction error significantly. Also, the probability of having an E_{SSC} less than 0.01 is around 90%, which again supports the fact that the

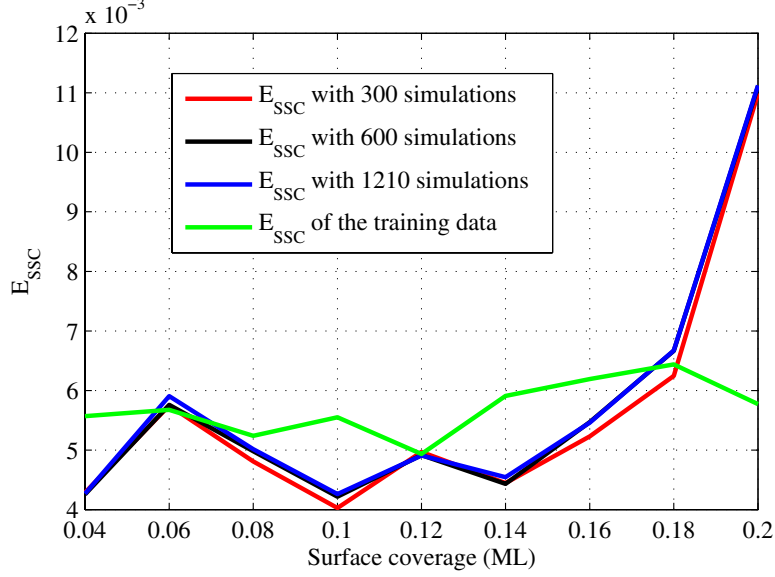


Figure 38: The mean value of E_{SSC} at different film coverage levels for three test simulation sets and the training data.

reduced order dynamic model has a good prediction capability.

Table 2: Mean and standard deviation values of E_{SSC} for three test simulation sets.

| Test simulation set | Mean | Standard deviation |
|---------------------|--------|--------------------|
| 1 | 0.0057 | 0.0035 |
| 2 | 0.0058 | 0.0037 |
| 3 | 0.0058 | 0.0037 |

In the last part of this section, we compare different types of prediction error to understand the major factors contributing to the inaccuracies in the model predictions. We define two alternatives to E_{SSC} in order to investigate the effect of normalization (using D_{mean} and D_{std} defined in Section 2.1.) on the prediction error. The first alternative to E_{SSC} is

$$E_{SSC'} = \|\mathbf{S}'_s - \mathbf{S}'_p\|_2 / \|\mathbf{S}'_s\|_2 \quad (10)$$

where \mathbf{S}'_s and \mathbf{S}'_p are the normalized versions of \mathbf{S}_s and \mathbf{S}_p :

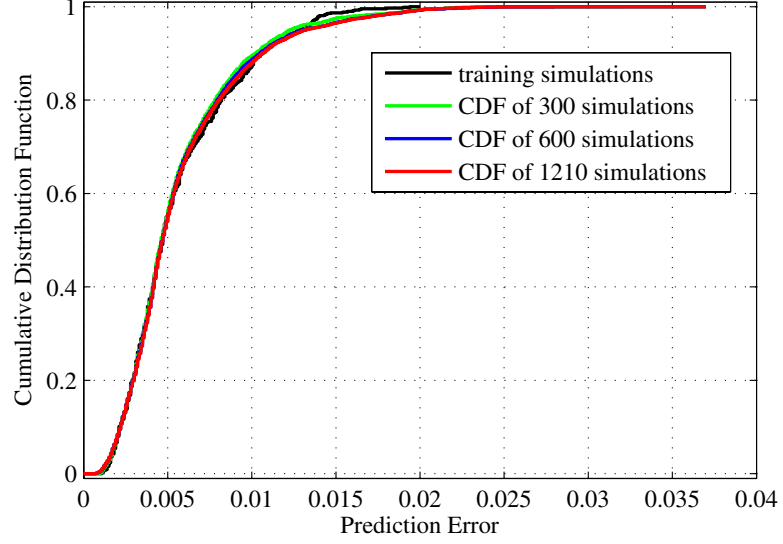


Figure 39: Cumulative distribution function of the E_{SSC} for three test simulation sets.

$$\mathbf{S}'_{s,i} = \mathbf{S}_{s,i} - D_{mean,i} \quad (i = 1, 2, \dots, 4800) \quad (11)$$

$$\mathbf{S}'_{p,i} = \mathbf{S}_{p,i} - D_{mean,i} \quad (i = 1, 2, \dots, 4800) \quad (12)$$

$D_{mean,i}$ was already defined in Section 2.1. The second alternative to E_{SSC} is:

$$E_{SSC''} = \|\mathbf{S}''_s - \mathbf{S}''_p\|_2 / \|\mathbf{S}''_s\|_2 \quad (13)$$

where \mathbf{S}''_s and \mathbf{S}''_p are the normalized versions of \mathbf{S}'_s and \mathbf{S}'_p , respectively:

$$\mathbf{S}''_{s,i} = \mathbf{S}'_{s,i} / D_{std,i} \quad (i = 1, 2, \dots, 4800) \quad (14)$$

$$\mathbf{S}''_{p,i} = \mathbf{S}'_{p,i} / D_{std,i} \quad (i = 1, 2, \dots, 4800) \quad (15)$$

$D_{std,i}$ was also defined in Section 2.1. Table 3 shows that mean value of $E_{SSC''}$ (0.2675) is much higher than that of $E_{SSC'}$ (0.1758). Comparison of the CDF of $E_{SSC'}$ and $E_{SSC''}$ given in Figure 40 also shows a higher $E_{SSC''}$ for any value of the CDF.

This indicates that features with large variance are predicted more accurately than the ones with small variance. It follows then, that the features with small variance are more noise-corrupted in the KMC simulations.

The quantization error, which is defined in Section 2.3.1, provides a minimum bound on the prediction error. The average quantization error (E_q) for SOM3 is 0.1094. This error is due to the discretization of the state space, and does not include any additional error propagated from one step to the next one through the dynamics of the thin film deposition process. Therefore, we define another kind of prediction error as

$$E_x = \|\mathbf{x}_s - \mathbf{x}_p\|_2 / \|\mathbf{x}_s\|_2. \quad (16)$$

Here, $\mathbf{x}_p \in \mathbb{R}^5$ represents the predicted state and $\mathbf{x}_s \in \mathbb{R}^5$ is computed from the KMC simulation data:

$$\mathbf{x}_{p,i} = \mathbf{S}''_p \cdot \mathbf{U}_i^T \quad (i = 1, 2, \dots, 5) \quad (17)$$

$$\mathbf{x}_{s,i} = \mathbf{S}''_s \cdot \mathbf{U}_i^T \quad (i = 1, 2, \dots, 5) \quad (18)$$

where $\mathbf{U}_i \in \mathbb{R}^{4800}$ is the i^{th} principal component of D'' .

As shown in Table 3, the mean value of E_x , which is computed from the prediction of states in Test Simulation Set 3, is twice as high as E_q incurred during the training of SOM3. From this result, it can be concluded that the propagation error due to dynamics is not negligible and accounts for approximately half of E_x . Hence, the other half of the prediction error comes from the discretization of the state space. Figure 40, which compares the CDF of E_x and E_q , also supports this conclusion. According to this figure, with a 0.50 probability, $E_q \leq 0.10$ and $E_x \leq 0.18$.

Table 3: Mean values of E_q , E_x , $E_{SSC'}$ and $E_{SSC''}$. E_q , which is computed for each data vector and the prototype vector of the data vector's best matching unit on SOM, is the difference between data vector and prototype vector during SOM training. The other errors are the normalized versions of E_{SSC} .

| Error type | Mean |
|-------------|--------|
| E_q | 0.1094 |
| E_x | 0.2277 |
| E_{SSC} | 0.0057 |
| $E_{SSC'}$ | 0.1758 |
| $E_{SSC''}$ | 0.2675 |

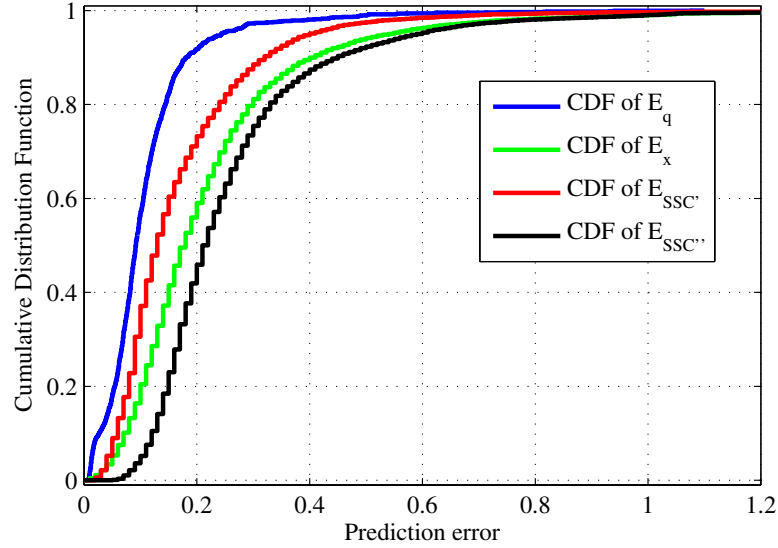


Figure 40: Cumulative distribution function of different types of error.

5.3 Optimization of the final film structure and the deposition time

In this section, the dynamic model is used to minimize the deposition time to reach a particular surface configuration. One desired structure is a very regular surface with many large islands and a very low Ga adatom (isolated Ga atom) density. This surface structure can be identified by minimizing:

$$F = a_i - b_i c_i \quad (19)$$

where a_i , b_i , and c_i are the values of Ga adatom density, typical island size and the number of islands with the typical island size for surface configuration i , respectively. From the Training Simulation Sets 1, 2 and 3, we extract a_i , b_i , and c_i values of the surfaces at 0.20 ML surface coverage.

According to equation (19), the optimal surface structure is Snapshot 861 (from Simulation 43). A portion of that surface structure is shown in Figure 41. Snapshot 861 is matched with the Node 182 of SOM3 during the training. This node is accessible through a constant input simulation with a flux of 0.08 ML/s. Identification of the flux profile that would reach Node 182 in the minimum amount of time can be posed as a dynamic programming problem [49].

Let the decision variables $d_q (q = 1, 2, 3, \dots, 11)$ be the immediate destinations on different stages. In this problem, we have 11 stages for 11 film coverage levels (0 ML, 0.02 ML, 0.04 ML, ..., 0.20 ML) and d_q corresponds to the map node number at stage q . Thus the route (trajectory) of the deposition is $d_1, d_2, d_3, \dots, d_{11}$, where $d_1 = 1$ and $d_{11} = 182$ since the initial surface structure is represented by SOM Node 1 and the final (optimal) surface structure is in Node 182.

Let $f_q(s, d_q)$ be the total cost of the best overall policy for the remaining stages, given that we are in state s (number of the map node we are currently in), ready to start stage q , and d_q (the number of the map node we are moving to) is selected as our immediate destination. Here, the total cost is the deposition time and each deposition interval is 0.02 ML long. Given s and q , let d_q^* denote any value of d_q that minimizes $f_q(s, d_q)$ and let $f_q^*(s)$ be the corresponding minimum value of $f_q(s, d_q)$. Thus,

$$f_q^*(s) = \min_{d_q} f_q(s, d_q) = f_q(s, d_q^*) \quad (20)$$

where

$$f_q(s, d_q) = c_{s, d_q} + f_{q+1}(s, d_{q+1}) \quad (21)$$

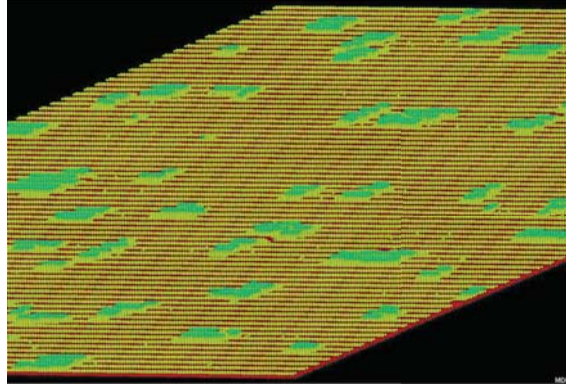


Figure 41: A portion of the optimal surface structure. The initial surface structure has regular trenches (dark areas), and as the deposition is performed, clusters form in and on top of the trenches.

Here, the cost c_{s,d_q} , is the time incurred while moving from s to d_q , given by:

$$c_{s,d_q} = L/F_{s,d_q}, \quad (22)$$

where L is the length of a single coverage interval (0.02 ML) and F_{s,d_q} is the value of the gallium flux that provides a transition from state s to d_q . The ultimate destination reached at the end of stage 11, $f_{11}^*(182) = 0$. The objective is to find $f_1^*(1)$ and the corresponding route. Dynamic programming can solve this problem by successively finding $f_{10}^*(s)$, $f_9^*(s)$, $f_8^*(s)$.. $f_2^*(s)$ and using $f_2^*(s)$ to solve for $f_1^*(s)$. This is achieved by eliminating some of the suboptimal trajectories as we move from $f_{10}^*(s)$ to $f_1^*(s)$. Because of the limited state space obtained by grouping similar surface configuration groups, it is possible to solve this dynamic optimization problem (finding the optimal flux profile) using exhaustive enumeration (without eliminating the suboptimal paths at each step) in a short amount of time.

We used eight flux settings (0.06,0.08,...0.20 ML/s) and found the optimal input profile that would give us minimum deposition time to get to SOM Node 182. Each input profile is a sequence of 10 flux values, and there is a flux value for each coverage interval. Having eight flux settings and 10 coverage intervals, there are 8^{10} possible input profiles. According to the dynamic model, only 20% of these profiles are able to

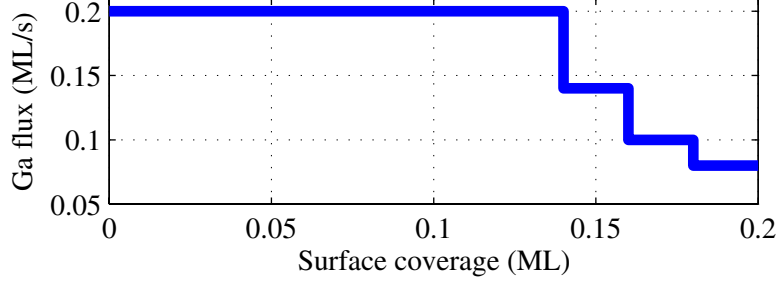


Figure 42: Optimal flux profile computed by the dynamic model.

reach SOM Node 182. Running each simulation takes about 24 hrs utilizing an Intel Xeon processor with a speed of 2.66 GHz, so it would have taken 2.9 million years to run all of the 8^{10} simulations with a single processor. However, using our dynamic model, it took only 5 minutes to predict the evolution of the film structure during these simulations. The minimum cost was obtained with the input profile shown in Figure 42.

Figure 43 shows the real and the estimated trajectories of this KMC simulation on SOM3. Again, we have a good agreement between these two trajectories. This particular input profile provided a 48% reduction in the deposition time to reach optimal structure when compared with the constant input KMC simulation under 0.08 ML/s Ga flux. The values of E_{SSC} , $E_{SSC'}$ and $E_{SSC''}$ for the prediction of the final film structure are 0.0159, 0.2207 and 0.2312, respectively. These values are within the range of cumulative distribution functions shown in Figures 39 and 40.

In order to visualize the accuracy of this prediction, we also plotted a portion of the SSC functions, which belonged to the actual KMC simulation data, its best matching unit on the map (Node 184) and Node 182 (predicted film structure). According to Figure 44, the simulation data is much noisier than the reconstructions. Also, the lower peak value of the SSC function (around radius=12) of the SOM Node 182 is slightly off compared to the one coming from the KMC simulation. However, the plateau corresponding to the number of step pairs which are distanced with more than 20 lattice units is captured well with the prediction. Hence, these two SSC

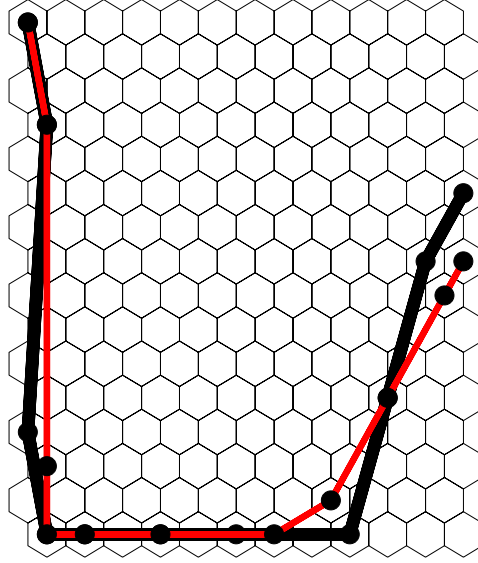


Figure 43: Trajectories of the KMC simulation (red line) with the optimal flux profile and the prediction (black line) on SOM3.

functions are very similar. This similarity indicates that the dynamic model, once again, does a good job in terms of predicting the final film structure. Furthermore, this prediction can be improved by increasing the size of the training data set and number of SOM nodes in the cell map, or possibly by placing a greater weight on the important, but small sized features, during the training of the SOM.

5.4 *Conclusions*

In this chapter, we evaluated the performance of the dynamic model derived from KMC simulation data of epitaxial GaAs deposition process. This evaluation involved quantification of different types of error associated with the model reduction approach. Analysis of the cell mapping error (CME), which is the one-step prediction error, showed that the structures within the same configuration groups show very similar dynamic behavior under same input conditions. The global error associated with this model reduction approach has also been characterized using 1210 test simulations with highly dynamic input profiles and turned out to be fairly low (less than 0.006 on average for 10890 predictions). Furthermore, the minimization of the deposition time

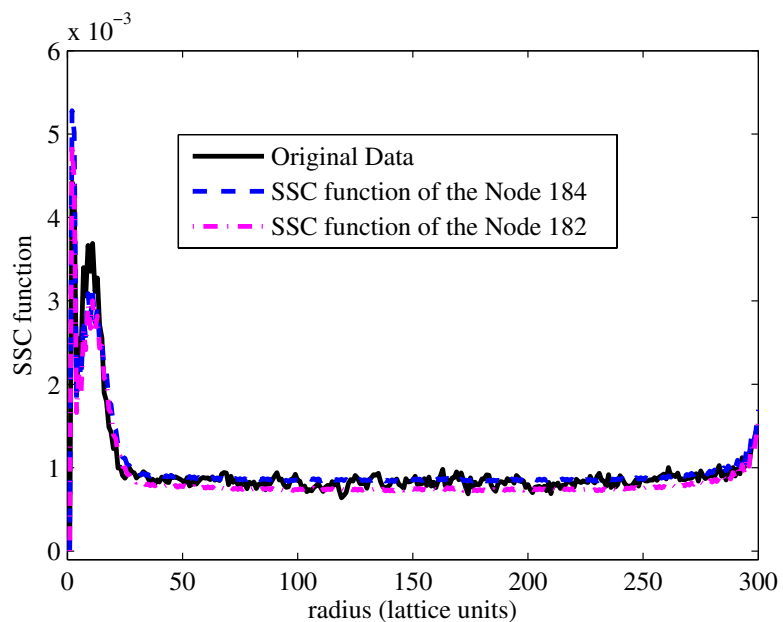


Figure 44: Reconstructions of the SSC function with the prototype vectors of the SOM nodes 182, 184 and the original simulation data.

to reach a desired film structure has also been achieved using the compact dynamic model. The reduced order model provides an 11 orders of magnitude reduction in the computational time compared to the high dimensional molecular simulations. Unlike the existing approaches used for dynamic optimization using stochastic simulation models [36], our approach does not involve frequent calls to the molecular simulation code which is computationally expensive.

CHAPTER VI

EVALUATION OF THE MODEL REDUCTION PARAMETERS

In this chapter, we present a critical evaluation of the model reduction parameters which can affect the predictive ability of the model. The parameters we consider here are the dimension of the state, number of configuration groups in the model, and the training data that is used for system identification. In the final section of this chapter, we focus on the effects of training data on model quality, and propose a novel state space exploration method that could improve the model's predictive ability under a wide range of process conditions.

6.1 Effects of the state dimension

In Chapter 4, we used a training data set to characterize the state space of the epitaxial deposition of GaAs. This characterization involved collecting a series of surface snapshots during kinetic Monte Carlo (KMC) simulations of this thin film deposition process under various input profiles. Since the dimension of the step-step correlation (SSC) function used for characterization was too high for system identification purposes, we have used principal component analysis (PCA) for reducing the state dimension. This reduction was done by eliminating the principal components which do not contribute to the variance of the data matrix significantly. After analyzing the eigenvalues of the covariance matrix (Figure 20) and data reconstructions with different number of principal components (Figure 21), we have retained 5 principal components, since that was enough to reconstruct the data effectively, capturing more than 99% of the variance. The state of each snapshot is computed by projecting its

SSC function onto the first 5 modes, and keeping the projection coefficients as $\mathbf{x} \in \mathbb{R}^5$, rather than the whole projection.

The number of principal components retained while reducing the state dimension can affect the results of the next step in the model reduction algorithm: discretization of the state space using self organizing map (SOM). After PCA is performed and the state vector \mathbf{x} is computed for each surface snapshot, SOM groups similar snapshots. Grouping is based on the Euclidean distances between state vectors (i.e. snapshots in similar states are grouped together), and snapshots in the same group are represented by only one surface structure and a vector \mathbf{v} (prototype vector of the SOM node) while identifying the dynamic model. This reduces the computational load for system identification significantly since only the transitions between surface configuration groups (not the individual surface configurations) are identified. In order to study the effects of state dimension on the SOM results, we trained SOMs with differently sized data vectors \mathbf{x} . Table 4 compares the SOM results when 2, 5, 10 and 20 principal components are used to compute the state of the surface snapshots in the training data. SOM7, which is trained with two dimensional state vectors, has the lowest prototype vector size since principal components 3, 4, and 5 are ignored in that case. However, rest of the map statistics is quite similar. We also look at the way surface snapshots are grouped by these maps. Since there are too many snapshots to consider, we only focus on the ones at the highest coverage of 0.20 ML. According to our analysis, these 76 snapshots are divided into 10 groups, and the snapshots within a group are in the same SOM node in all of the maps. For example, snapshots 1, 13, 20, 24, 32, 36, 39, and 47 are grouped together in these 4 maps. Only 2 of the 76 snapshots move between different groups in different maps. The similarity of the statistics and grouping of SOM3, SOM5, SOM6 and SOM7 indicates that the number of principal components retained does not have a significant effect on the SOM results, and therefore the system identification step. In fact, even 2

principal components are enough for consistent grouping of SOM. However, at least 5 principal components are needed to reconstruct SSC functions of the snapshots effectively (discussed in Section 4.3.1), which is necessary for identifying an accurate dynamic model. For other problems, very high order principal components might be necessary to differentiate between configurations that are slightly different than each other in terms of their state. In those cases, during SOM training, elements within the state vectors that are associated with the high order principal components could be weighted higher. Therefore, they would have a significant effect on the SOM grouping despite the very low percent variance captured by the high order principal components associated with them.

Table 4: Statistics of SOMs trained with differently sized data vectors.

| Map | State dimension | Quant. error | Proto. vec. size | Topog. error | Number of groups |
|------|-----------------|--------------|------------------|--------------|------------------|
| SOM3 | 5 | 5.1 | 53.0 | 0.013 | 175 |
| SOM5 | 10 | 5.4 | 53.3 | 0.018 | 177 |
| SOM6 | 20 | 5.7 | 53.0 | 0.012 | 175 |
| SOM7 | 2 | 4.2 | 52.8 | 0.013 | 175 |

6.2 *Effects of map size on the model predictions*

In Chapter 4, we used an SOM with 192 map nodes to group 1521 snapshots in order to build a dynamic model. When the map size is increased, a finer discretization of the state space can be performed. For such a larger map with 775 nodes (SOM4), statistics are compared with the smaller map (SOM3) in Table 5. SOM4 has a much lower quantization error since it has more nodes and thus fewer snapshots per node. On the other hand, its topographic error is 3% which is slightly higher than SOM3's topographic error (1%) since there is a higher chance of first and second best matching units of a snapshot not being adjacent when the number of map nodes increases. Here

a topographic error of 3% means that, for 45 of the total 1521 surface snapshots in the training data, the first and second best matching units on the map were not adjacent. Even though the larger map (SOM4) has a slightly higher topographic error than the smaller one (SOM3), both can be regarded as high quality maps with low quantization and topographic errors. Low quantization error shows that SOM nodes are able to represent the data vectors that are matched with them. In this case, approximation of the data vectors in a map node by a single prototype vector is justified. On the other hand, low topographic error is desirable for preservation of the topology (i.e. similarity of the neighboring SOM nodes on the map), and visualization of the system dynamics. For example, an SOM with low topographic error can provide a comparison of the KMC simulation trajectories and the trajectories predicted by the dynamic model based on their paths on the SOM.

Table 5: Statistics of differently sized SOMs.

| Map | Map size | Quant. error | Proto. vec. size | Topog. error | Number of groups |
|------|----------|--------------|------------------|--------------|------------------|
| SOM3 | 192 | 5.1 | 53.0 | 0.01 | 175 |
| SOM4 | 775 | 2.8 | 53.0 | 0.03 | 554 |

One important model reduction parameter is the computational load for the system identification step. As explained in Section 4.4.1, system identification required simulations starting from 146 initial conditions for SOM3 which took approximately 7.3 days using 16 Intel Xeon processors with 2.66 GHz speed that were run in parallel. The number of initial conditions is 286 for SOM4, doubling the computational time necessary to identify the transitions between different partitions of the state space. As the initial conditions to be used, we consider SOM nodes associated only with even film coverage values (0 ML, 0.02 ML, 0.04 ML,..., 0.18 ML). This leads to the elimination of the SOM nodes associated only with odd coverage values (0.01 ML, 0.03

ML, 0.05 ML, ..., 0.19 ML). Also, SOM nodes which exclusively contain snapshots at the largest film coverage in the training data (0.20 ML) are not considered since the maximum film coverage level that can be simulated by the KMC model is 0.20 ML. Therefore, running simulations with the initial condition of 0.20 ML film coverage would lead to 0.22 ML, which would violate the maximum coverage level constraint. Furthermore, states described by SOM nodes which only contain snapshots at 0.20 ML (3 nodes of SOM3, and 19 nodes of SOM4) are assumed to be absorbing states [103] that are impossible to leave. As a result, a total of 149 surface configurations are considered in the low dimensional model derived from SOM3, with 146 configurations up to 0.18 ML coverage, and 3 configurations at 0.20 ML. On the other hand, the high dimensional model obtained from SOM4 has 305 surface configurations where 286 of them belong to film structures up to 0.18 ML coverage and the remaining 19 are at 0.20 ML. Therefore, the number of initial conditions used for system identification is less than the number of snapshot groups for both maps (Table 5).

While generating surface snapshots for the training data, our sampling interval was 0.01 ML. In other words, during the training simulations, a surface snapshot was recorded at every 0.01 ML film coverage starting from 0.01 ML until 0.20 ML. After using SOM to discretize the state space by these simulations, we identified the transitions between different cells or partitions of the state space (SOM nodes). This procedure, which is called cell mapping, is described in Section 4.4.1. Simple cell mapping involves running short simulations (short bursts of the KMC simulation model) starting from each surface configuration described by an SOM node, and identifying the final structure in these simulations, and their best matching units on the SOM. Using this technique, all possible transitions within the state space are represented as transitions on the SOM. One important thing to note is that the interval used in these short simulations was 0.02 ML (each short simulation was run for 0.02 ML starting from the initial condition). The length of this interval is twice

the length of the coverage interval used while collecting data (0.01 ML). We made this selection to avoid the occurrence of absorbing states in the model before the maximum film coverage level. An absorbing state is defined as a state which is impossible to leave [103] since it transitions into itself repeatedly. Because of the structure of the reduced order model (described in Section 4.4.1), if the structure of a surface does not change significantly during a short burst used for identifying the transition out of that state, and stays in its starting map node, then the transition matrix (described by equation (6) in Section 4.4.1) would dictate the repeated occurrence of this process. In this case, the film state would get stuck in the same SOM node. Since the system trajectory is required to move from left to right on SOM3 and SOM4 as the film coverage increases, such a case would be unrealistic. Therefore, instead of using the sampling interval of 0.01 for the length of the short simulations for system identification, we used twice that value to ensure continuously moving trajectories on the map as the film coverage increases.

In this section, we compare the predictive abilities of the low and high dimensional dynamic models (obtained from SOM3 and SOM4, respectively), and quantify how the extra computational load necessary to identify the latter improved the model performance. Dynamic input profiles used in Training Simulation Sets 1, 2 and 3 (described in Section 4.1) mostly involved small changes between intermediate values of gallium flux using 8 discrete flux settings (0.06 ML/s, 0.08 ML/s, 0.10 ML/s... 0.20 ML/s). Since the dynamic model was trained using the data coming from these simulations, a good way to test the predictive accuracy of the model under different conditions is using flux switches between the maximum and minimum settings. Because, low flux settings produce smooth surfaces compared to the high flux conditions that result in less number of diffusion events per an adsorption, and therefore rough surfaces. Because of this, flux switches between two extreme values can potentially produce new behavior which is not observed in the training data generated with small

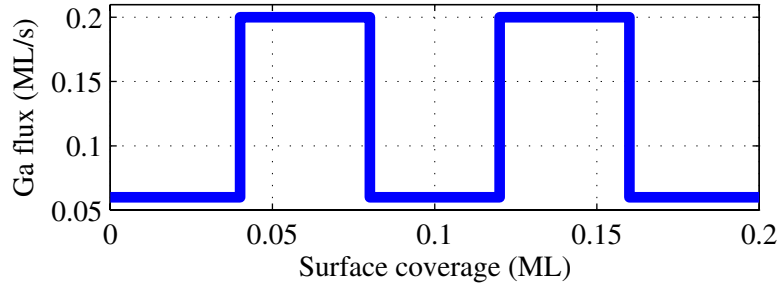


Figure 45: Periodic flux profile in a test simulation.

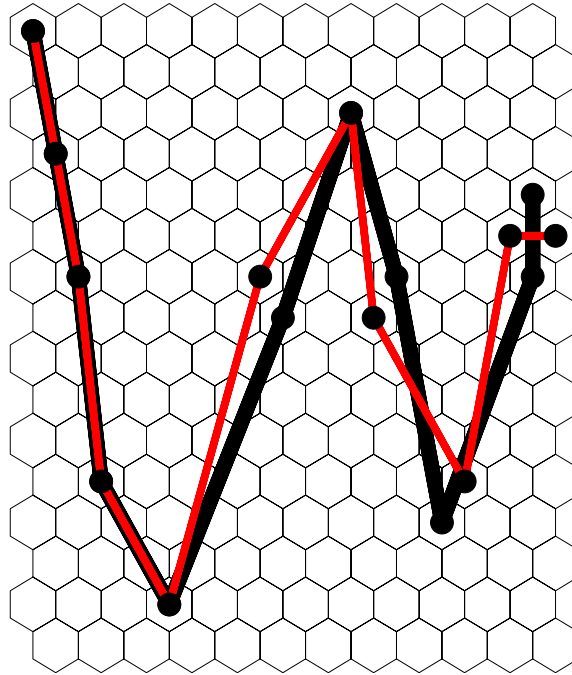


Figure 46: Trajectories of a test KMC simulation (red line) and the prediction (black line) on SOM3.

changes in the flux. Figure 45 shows a flux profile of a simulation in which switches between the extreme flux settings are made. The comparison of the trajectories of this KMC simulation and the predicted trajectory on both maps (Figures 46 and 47) show that there is an agreement between these trajectories using both models. The low value of the quantization error and topographic error for these maps, together with the similarity of the simulation trajectories, and the predicted ones suggest that both models perform well for this specific test simulation even though the high dimensional one required twice the computational load for its identification.

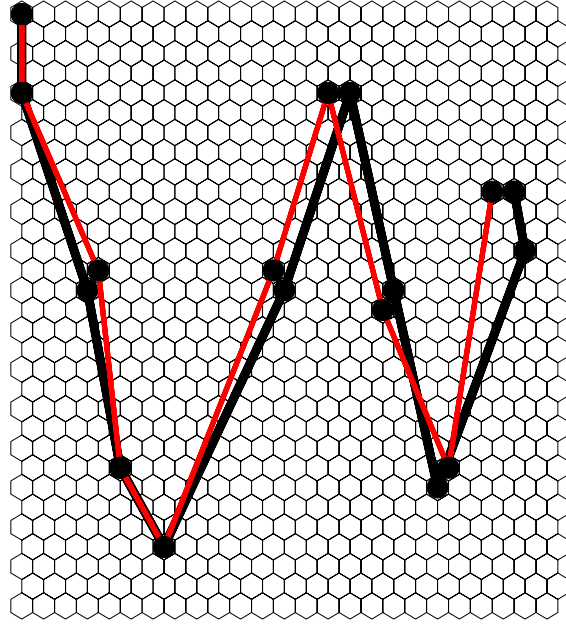


Figure 47: Trajectories of a test KMC simulation (red line) and the prediction (black line) on SOM4.

A better way to compare the performance of the models identified from SOM3 and SOM4 is the quantification of the model error. First, we define the relative quantization error over the test simulation data as $E_{rq} = \|\mathbf{x}_s - \mathbf{v}_{BMU}\|_2 / \|\mathbf{x}_s\|_2$. In words, it is the Euclidean distance between the state vector of a surface snapshot and the prototype vector of its best matching unit with respect to the norm of the snapshot's state vector. Figure 48 shows that E_{rq} is lower with the larger SOM (SOM4) at almost all coverage levels. This was expected since the larger map has more nodes that can match with a wider range of possible system states. As a result, the larger map has a lower approximation error associated with the discretization of the state space. Figure 49 indicates that at 0.10 ML film coverage, SOM4 does a slightly better job than SOM3 in terms of reconstructing the KMC simulation data. In this figure, a portion of $SSC_{up,down,i}$ coming from the KMC simulation data is compared with its reconstructed form using the prototype vectors of its best matching units on SOM3 and SOM4.

A similar trend is observed when we compare the state prediction error $E_x =$

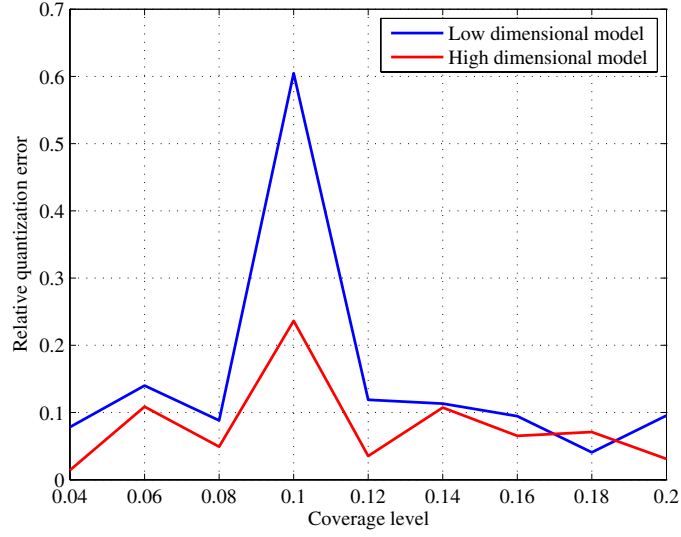


Figure 48: The value of the relative quantization error E_{rq} at different film coverage levels for the test simulation with the flux profile given in Figure 45.

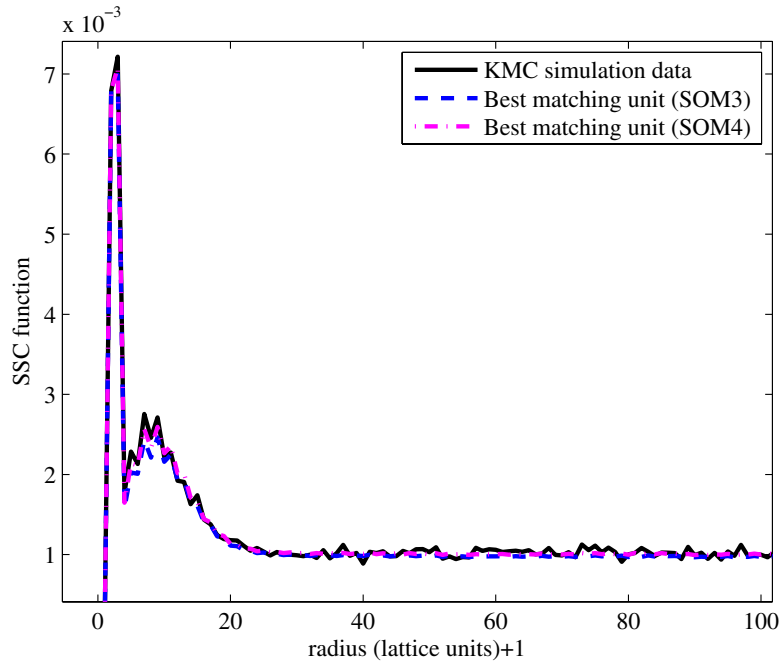


Figure 49: A portion of $SSC_{up,down,i}$ from the KMC simulation data at 0.10 ML film coverage, and its reconstructed form using the prototype vectors of its best matching units on SOM3 and SOM4.

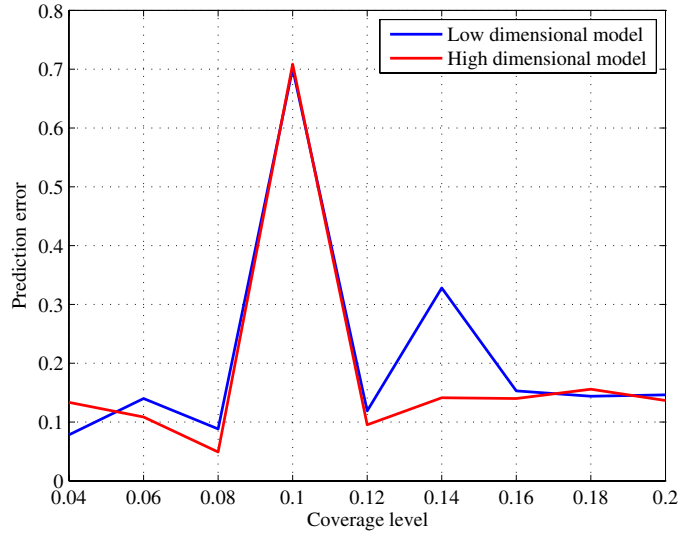


Figure 50: The value of the state prediction error E_x at different film coverage levels for the test simulation with the flux profile given in 45.

$\|\mathbf{x}_s - \mathbf{x}_p\|_2 / \|\mathbf{x}_s\|_2$ of the high and low dimensional models identified from SOM4 and SOM3, respectively (Figure 50). State prediction error represents the normalized difference between the predicted state (\mathbf{x}_p) and the state in the KMC simulation (\mathbf{x}_p). According to Figure 50, high dimensional model gives more accurate predictions at almost all coverage levels. Again for illustration purposes, we compare $SSC_{up,down,i}$ from the KMC simulation data and the predictions of both models, and observe that the high dimensional model does a slightly better job in terms of capturing the SSC features (Figure 51).

Results presented in Figures 48, 49, 50 and 51 are only for a single test simulation that has drastic changes in its flux profile. In order to compare the distributions of the prediction error with the large and the small models, we vary the flux using the minimum and maximum flux settings and 10 coverage intervals, and explore the state space much more aggressively compared to the 1210 test simulations we used in Chapter 4, which included the intermediate flux settings as well. In these simulations, the flux value is kept constant at 0.06 ML/s (minimum setting) and

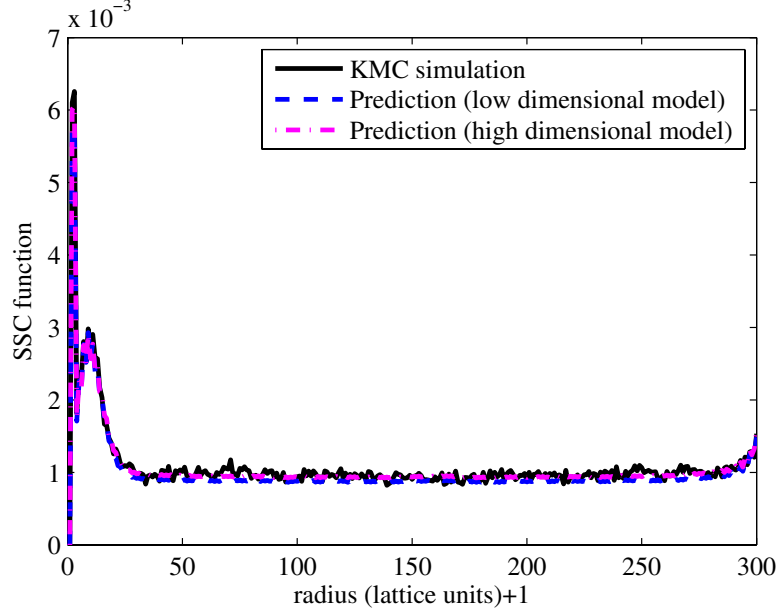


Figure 51: A portion of $SSC_{up,down,i}$ coming from the KMC simulation data at 0.14 ML film coverage and the prediction of the low and high dimensional models.

0.20 ML/s (maximum setting) in each coverage interval of 0.02 ML, and we consider all possible flux sequences (2^{10} simulations). We use E_{SSC} (previously defined in Section 6.2.) to quantify the error. As a reminder this error function is defined as $E_{SSC} = \|\mathbf{S}_s - \mathbf{S}_p\|_2 / \|\mathbf{S}_s\|_2$ where \mathbf{S}_s is the computed SSC function from the KMC simulation data, and \mathbf{S}_p is reconstructed from the predicted state vector using the first 5 principal components of the training data. As shown in Figure 52, the cumulative distribution function of E_{SSC} of the high dimensional model is slightly above the one for the low dimensional model indicating that the high dimensional model has a better prediction capability under the conditions that involve switches between minimum and maximum flux settings. However, despite the much higher computational load necessary to derive the high dimensional model, the improvement in the predictive ability is minimal. According to Figure 52, with a 50% probability, $E_{SSC} \leq 0.025$ for the low dimensional model. On the other hand, with the same probability, $E_{SSC} \leq 0.024$ for the high dimensional model. Therefore, the low dimensional model is a better alternative in terms of the trade-off between computational efficiency and predictive

ability.

Figure 52 also shows another interesting comparison between the CDFs of the prediction error of the low dimensional model with two different sets of flux profiles. According to this figure, when flux profiles that involve switches between minimum and maximum flux settings (0.06 ML/s and 0.20 ML/s) are used instead of all 8 flux settings that include intermediate flux values (0.06 ML/s, 0.08 ML/s, 0.10 ML/s... 0.20 ML/s), prediction error increases significantly. Because the training data did not include many simulations in which the flux value was switched between extreme flux settings, and was unfamiliar with some of the surface configurations that are generated under such extreme conditions. Since we have used a data based modeling approach for model identification, the quality of the training data in terms of its ability to include the accessible states as much as possible, is very important for model's prediction accuracy under a wide range of process conditions. Therefore, for a better model with higher prediction accuracy, it would be desirable to include training simulations with cleverly designed input profiles in order to excite system dynamics as much as we can for system identification purposes. This is an issue about state space exploration which we talk about in the next section.

6.3 Effects of state space exploration

It is well known in process control that while using models derived from input-output data (i.e. empirical models), the prediction accuracy of the model can be highly dependent on the training data used for system identification [98, 115]. The identification signals that are used to generate training data determine the response of the system, and therefore the information relevant to system dynamics that will be captured. One commonly used approach to generate input profiles that can excite a wide range of dynamic behavior is the pseudo-random binary sequence (PRBS) technique with two input levels (e.g. high and low input settings). However, it is

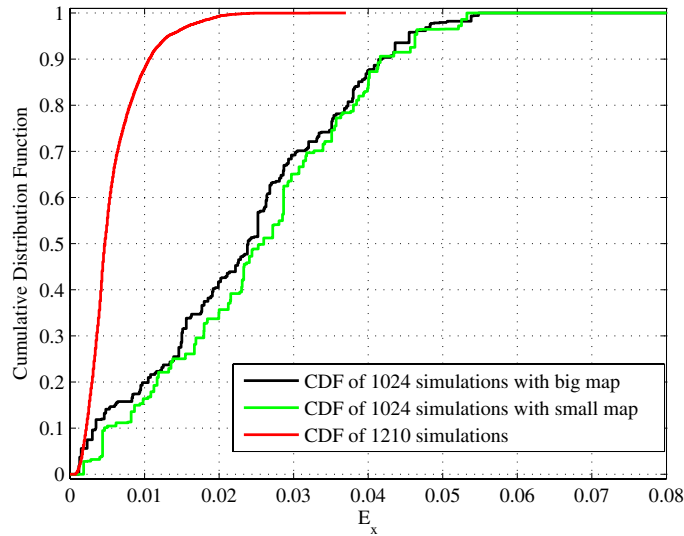


Figure 52: Cumulative distribution functions of the prediction error for 1210 test simulations with the low dimensional model, and the 1024 exploration simulations with the low and high dimensional models.

known that this technique is not capable of identifying nonlinear systems which are very common in process control [40]. This is due to the PRBS technique’s inability to excite nonlinearities in system dynamics [19, 94]. In general, for linear model identification, input signals should have amplitudes that are small enough to prevent nonlinear effects [41]. On the other hand, for nonlinear model identification, amplitudes should be large enough to excite the nonlinearities of a process. Another important factor in nonlinear system identification is the frequency of the changes in the input level, which should be as much as possible while staying within feasible levels of frequency. Therefore, multilevel PRS inputs, where the number of input settings is much larger than two, are commonly used for nonlinear system identification. One example is a study where 31 levels of inputs are used [1]. In most cases, such a high number of input levels is not necessary to identify reliable models [11]. In fact, “plant friendliness” is defined as a technique’s ability to generate the input-output data that can minimize the uncertainty in model parameters in the shortest time while having the input and output variance within desired ranges. Such an approach

is Model On Demand used by Braun et al. [11] to perform system identification for a rapid thermal processing reactor, and a pilot scale brine-water mixing tank. Unlike the traditional nonlinear system identification techniques, which make use of very large data sets representing the whole state space (global modeling), Model on Demand is developed as a local modeling technique. With this data, models are built 'on demand' when needed, around the neighborhood of the available data. The identification of the model is posed as an optimization problem where the goal is to minimize the estimation error. Since this optimization is done locally, and without requiring a non-convex optimization scheme, error is smaller compared to global approaches. Also, Model on Demand allows incorporation of the *a priori* information into the design of multi-level PRS signals and minimizes the number of simulations needed for system identification. However, this local modeling approach is computationally more demanding than global approaches because of the neighborhood search around each data point to identify the model by regression. Even though such local approaches can have higher computational demands, their superior performance can outweigh this disadvantage. In a recent thesis by Lee [81], a dynamic programming approach for nonlinear control using input-output data, and local approximators was advocated (dynamic programming was described earlier in Section 5.3). In that study, the author approximates the cost-to-go function around the data points in the state space with adequate data density to control an adiabatic CSTR. Similar to Braun [11] and Lee [81], our approach is a local modeling approach since we divide the state space into cells (SOM nodes), and each cell is approximated by a prototype vector computed by the states grouped within that particular cell.

In the previous section, we demonstrated that the quality of the training data can be improved by using input profiles that involve switches between the extreme values of the gallium flux. When such simulations were not included extensively in the training data, and the dynamic model was used to predict the behavior of the system

in these simulations, prediction accuracy decreased significantly as shown in Figure 52. The sensitivity of model predictions to the quality of the training data can be a problem when the dynamic model is used to solve an optimal control problem such as minimization of the deposition time to reach a certain film surface configuration. This problem was described earlier in Section 5.3. When a model identified from input-output data is reliable in only a certain fraction of the state space, optimal inputs need to be searched cautiously. Because, in certain regions of the input space with little or no training data, extrapolations for model prediction can lead to poor predictions. Furthermore, if the local fitting error is high, interpolation can also cause the same problem. In such cases, searching for the optimal input values can be restricted to the parts of the state space with adequate amount of training data. One such approach in the context of neural networks was proposed by Leonard et al. [82] where a validity index that quantifies the reliability of model predictions is used for both “local goodness of fit” in the case of interpolations and extrapolations.

One approach to address the model inaccuracies, while performing optimization, is penalizing predictions of the model in certain regions of the state space. Lee recently proposed [81] imposing a penalty term in the objective function to limit the search for the optimal inputs within the explored regions of the state space. In our study, we address this problem by solving the optimal control of deposition (minimization of the deposition time to reach a specific thin film structure) by restricting our search within the state space discretized by SOM. This can prevent the extrapolation problem. On the other hand, in order to prevent poor interpolation, we can include a penalty term in the cost function (deposition time for a given sequence of inputs) for the transitions in the state space with high cell mapping error. This type of error is described earlier in Section 5.1. It quantifies the local error which originates from representing multiple surface structures in an SOM node using a single prototype vector. If the cell mapping error is very low, then all surface structures grouped

together in the SOM node are expected to show very similar or identical dynamic behavior under same input conditions. In that case, the reduction in the number of possible states by grouping similar surface configurations is justified. For the SOM nodes associated with high cell mapping error, by including a penalty term in the cost function, we can penalize the poor local approximation and obtain a more reliable solution to the problem of minimizing the thin film deposition time.

Another alternative approach to address the issue of model inaccuracies in certain regions of the state space is continually updating the dynamic model by systematic explorations of the state space. Similar approaches in the past contain studies of Hernandez et al. [48], and Chikkula et al. [16] who updated the model parameters with new identification data in a recursive fashion. More recently, using a model free approach, Lee [81] proposed a technique which involved adding the cost-to-go functions in unexplored regions of the state space to prevent inaccurate extrapolations around the existing data points. In this section, we address the same problem by proposing a novel state space exploration method that involves updating the principal components that span the explored state space at each film coverage level. In this approach, after characterizing surface snapshots and updating the principal components, SOM is used to identify surfaces with state vectors that are significantly different than the prototype vectors of their best matching units. If the quantization error for a snapshot (i.e. the Euclidean distance between the state vector of the snapshot and the prototype vector of its best matching unit) is much higher than the average quantization error of the whole SOM, we use this system state as an initial condition for state space exploration by running short simulations (for 0.02 ML) until the next film coverage level. Because such a difference between the state and prototype vectors would indicate the SOM's inability to effectively represent the surface configuration with the prototype vector of its node, and this might lead to prediction inaccuracies since all surface structures matched with and SOM node are represented by a single

prototype vector.

Following is the strategy we propose for systematic state space exploration of epitaxial GaAs deposition with KMC simulations:

- Starting at the initial surface configuration ($\beta 2(2 \times 4)$ reconstruction of GaAs shown in Figure 19), we run 8 simulations under all flux settings (0.06 ML/s , 0.08 ML/s, 0.10 ML/s,..., 0.20 ML/s) for a film coverage interval of 0.02 ML, and record surface snapshots generated at the end of these simulations.
- Starting from each of the 8 surface snapshots at 0.02 ML described above, we run 8 eight more simulations under all flux settings until the next film coverage level (0.04 ML). When these simulations are finalized, a total of 64 more surface snapshots are recorded (8 initial conditions \times 8 simulations per initial condition = 64 surface snapshots.)
- Once again, starting from each of the 64 surface snapshots generated in the previous step, we run 8 simulations under all flux settings to reach the next film coverage level (0.06 ML). This gives us 512 (8 \times 64) more surface snapshots. At this point, we collect the SSC functions of all snapshots generated until 0.06 ML (8 + 64 + 512 = 584 surface snapshots) and the initial surface configuration (a total of 585 snapshots) in a data matrix. Then, the principal components of this data matrix are found, and the state vector of each snapshot is computed by using the first 5 principal components. Details of this procedure are given in Section 4.2.
- Following the state space characterization, SOM is trained with the state vectors of 585 surface snapshots to identify surface configuration groups. The average quantization error of this SOM (average Euclidean distance between state vectors and the prototype vectors) is 7.23, whereas average prototype vector norm is 54.1 (13.4% average quantization error relative to the average prototype vector

norm). Using the SOM results, we find that for 71 of the 512 surface snapshots at 0.06 ML coverage, quantization error is higher than the selected minimum threshold (30% average quantization error relative to the average prototype vector norm). This threshold must be greater than the average value for the whole SOM, but it should not be too high. Because, in that case, very few surface structures (or none) would be picked as starting points for further state space exploration.

- Starting from each of the 71 initial conditions picked in the previous step, we run 8 simulations under all flux settings until 0.08 ML film coverage and characterize the final surface snapshots ($71 \times 8 = 568$ snapshots). Since these structures are obtained from a very thorough exploration of the state space (unlike the training simulations), we characterize the distribution of the prediction error E_{SSC} (defined in section 5.2) to see if the model performs poorly while predicting the evolution of the system under input conditions that were possibly not included in the training data. Figure 53 compares the CDF of the E_{SSC} for these state space exploration simulations and the 1210 test simulations which did not involve systematic exploration of the state space. According to this figure, the model performs equally well in both cases. Therefore, we conclude that at this coverage level, accessible surface configurations in the KMC simulation model are explored well enough with the training simulations and further exploration might not be necessary. One reason for this behavior is the similarity of the surface snapshots that are generated under a wide range of flux conditions at early coverage levels. This was previously demonstrated in Figure 23 where the number of SOM nodes that contained surface snapshots at early coverages (0 ML-0.08 ML) were fewer than the rest of the nodes. Because, these fewer nodes were able to group more snapshots per node compared to the other nodes associated with 0.10 ML coverage level and higher due to the similarity of the

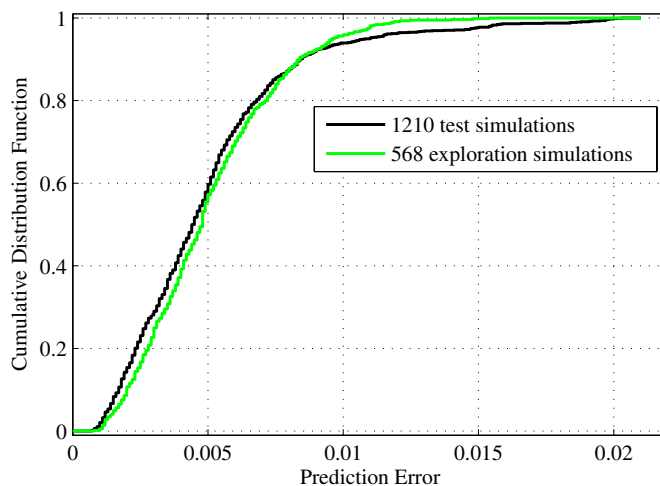


Figure 53: Cumulative distribution functions of the prediction error for 1210 test simulations, and the 568 exploration simulations with the small model at a film coverage of 0.08 ML.

state vectors of surface snapshots generated earlier in the training simulations.

- In order to carry out the state space exploration further, we collect all the surface snapshots generated until 0.08 ML, characterize them using SSC functions, compute their states, and use SOM once again for grouping similar snapshots. The resulting SOM has an average quantization error 6.60 with an average prototype vector size of 57.4 (11% average quantization error relative to the average prototype vector size). The threshold we choose for identifying new surface configurations (20% error) is about twice the average. We identify 30 surface snapshots (at 0.08 ML coverage level) for which the quantization error is above that threshold as new surface configurations which possibly have not been explored previously, and not included in the training data. These 30 structures are used as initial conditions in the next step for further state space exploration.
- Starting from each of the 30 initial conditions selected in the previous step, we run 8 simulations under all flux settings until 0.10 ML film coverage and

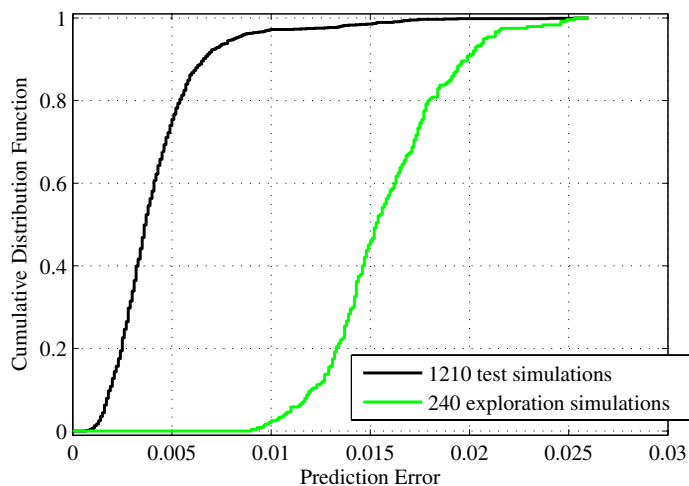


Figure 54: Cumulative distributions function of the prediction error for 1210 test simulations, and the 240 exploration simulations with the small model at a film coverage of 0.10 ML.

characterize the final surface snapshots ($30 \times 8 = 240$ snapshots). The gallium flux profiles that were used to generate these structures are used as inputs to the dynamic model and the prediction error E_{SSC} is computed once again as the difference between the predictions and the actual KMC simulation results. Figure 54 shows that the CDF of the E_{SSC} for these predictions is below the one for the 1210 test simulations which had randomly generated input profiles. We conclude that the state space exploration until 0.10 ML generated some surface configurations that the dynamic model was unfamiliar with, that led to poor predictive performance. This indicates that a systematic state space exploration method based on updating the principal components that span the state space and expanding the SOM (or cell map) with new states can be useful to improve model’s predictive ability under a wide range of process conditions.

6.4 *Conclusions*

In this chapter, we presented a critical evaluation of the model reduction parameters which can affect the predictive ability of the model. These parameters are the dimension of the state, number of configuration groups in the model, and the training data that is used for system identification. First, we show that projecting the data onto more than five principal components do not change the data vectors, and therefore the grouping of similar data vectors by the SOM. Then, we compare the predictive ability of two dynamic models that include different number of surface configuration groups (SOM nodes) and conclude that the high dimensional model generates very similar predictions despite the fact that it is computationally twice as demanding to identify, compared to the low dimensional one. In the final section of this chapter, we focus on the effects of training data on model quality, and propose a novel state space exploration method that can improve the model gradually by updating the principal components that span the state space, and expanding the SOM at each film coverage level.

CHAPTER VII

CONTRIBUTIONS AND FUTURE WORK

7.1 *Contributions*

This thesis proposed modeling techniques for the design and optimization of material systems that require atomic scale models for simulation. The necessity to capture detailed interactions in these systems leads to high dimensional state description, and a high number of possible events in the system that take place at a wide range of length and time scales. These factors result in prohibitive computational costs when one tries to use molecular simulation models directly for design and optimization. Therefore, in addition to development of accurate simulation models, reduced order models with lower computational cost are also necessary. The major contributions in this thesis are:

1. A KMC simulation model for hyperbranched polymerization process, and the investigation of the effects of different synthesis routes on the polymer structure using simulations and experiments.
2. A model reduction algorithm that enables the use of high dimensional molecular simulation data for the dynamic optimization of epitaxial deposition of gallium arsenide.

In Chapter 3, we presented a KMC simulation model for the hyperbranched polymerization of difunctional A_2 oligomers, and B_3 monomers. This model, whose kinetic parameters are derived from experimental data, is employed to study the effects of the synthesis routes on the polymer structure development. Comparing the experimental and simulation results, we showed that in melt polymerization (with no solvent),

the extent of cyclization reactions is negligible. We also observed that the polymer structure is strongly affected by the high reactivity of free B_3 units (compared to the partially reacted B_3 units), and the presence of endcapping reactions in the medium. In the rest of Chapter 3, we focus on the solution polymerization of an $A_2 + B_3$ system with dropwise addition of difunctional A_2 into B_3 . Our modeling and experimental studies showed that dilution of the B_3 solution promotes cyclization reactions, delays the gel point, and suppresses the polymer molecular weight. These results indicated that KMC simulations provide a tool for understanding the effects of simple reactions mechanisms on the polymer structure by enabling the consideration of a broader range of mechanisms compared to analytical models. Furthermore, the results also suggested that simulations and experiments can be potentially used for designing and controlling the process with process inputs such as solution concentration, or monomer feed rates.

Molecular simulations generally have high dimensional state description (e.g. coordinates of thousands of atoms), and high dimensional state space (number of possible system configurations). Therefore, they are not compatible with the existing design, and optimization tools. In order to address this issue, in the last section of Chapter 3, we suggested a model reduction approach to derive a reduced order process model for hyperbranched polymerization from KMC simulations and experimental data. Description of the dynamic state in the simulations, and experimental validation of the reduced order model have been identified as the major challenges of this approach. In Chapter 4, we formalized this novel model reduction algorithm that makes use of the high dimensional simulation data (representing the evolution of the system properties in a set of simulations) to derive a reduced order process model. To demonstrate our algorithm, we employed an existing KMC simulation model for the epitaxial deposition of gallium arsenide (GaAs). By implementing principal component analysis (PCA) to reduce the dimensions of the state vector of each surface configuration,

self-organizing map (SOM) for grouping similar surface configurations, and simple cell mapping and k-nearest neighbor method (for identifying the transitions between the surface configuration groups), we obtained a reduced order model of the process. This algorithm does not rely on a wide separation between the time scales of low and high order system statistics, existence of steady state points, or the accuracy of the time derivatives of system properties, which are generally computed from short simulations to be used for optimization in other studies.

In the first part of Chapter 5, to validate our model reduction algorithm, we computed the distribution of the error resulting from the reduction. First, we generated simulation data by recording and characterizing surface snapshots in simulations which have input profiles different than the ones used for identifying the model. These new simulations were called test simulations, whereas the simulations already used for model identification were called training simulations. After predicting the evolution of the system state in the test simulations, we computed the distribution of the prediction error. Results showed that the overall prediction error is quite low even though the training simulation set was small (76 simulations with constant and varying gallium flux profiles) compared to the test simulation set (1210 simulations with randomly varying gallium flux profiles).

In the second part of Chapter 5, we performed dynamic optimization of the thin film deposition to determine the optimal flux profile that can reach a given final film structure in the minimum amount of time. Unlike the other approaches for dynamic optimization of molecular simulations, our approach does not need the computation of the derivatives of system variables which would require running extra simulations at each step of the optimization in the presence of noise originating from the stochastic nature of the simulations. Since we identify the state space, and the transitions between different partitions of the state space, it was also possible to perform dynamic optimization. In our model, each partition (cell) of the state space is represented

by a typical surface configuration rather than all the surface configurations within that partition. Therefore, while performing dynamic optimization, it was enough to just consider the transitions between a limited number of configurations. This simplification led to an 11 orders of magnitude reduction in computational time for the optimization compared to the high dimensional KMC model.

In Chapter 6, we evaluated the model reduction parameters such as the dimension of the state, cell map (SOM) size (number of configuration groups in the state space), and the quality of the training data. Despite the high computational load associated with its identification, the high dimensional model did not significantly improve our prediction capability. We also showed that the major source of prediction error was the quality of the training data. Finally, in the last section of this chapter, we presented a novel state space exploration method that can be used to identify the unexplored regions of the state space by updating the principal components, and training the SOM at each coverage level.

The following is the list of publications resulted from this work:

- Journal publications
 - Oguz, C., Cakir, S., Yilgor, E., Gallivan, M. A., and Yilgor, I., “Investigation of the influence of polymerization procedure on polymer topology and properties in branched polymers obtained by $A_2 + B_3$ approach.” in preparation to be submitted to Polymer.
 - Oguz, C. and Gallivan, M. A., “Optimization of a thin film deposition process using a dynamic model extracted from molecular simulations.” accepted for publication in Automatica.
 - Oguz, C., Unal, S., Long, T. E., and Gallivan M. A., “Interpretation of molecular structure and kinetics in melt condensation of A_2 oligomers, B_3 monomers, and monofunctional reagents,” *Macromolecules*, vol. 40, pp.

6529-6534, Sep. 2007.

- Oguz, C. and Gallivan, M. A., “A data-driven approach for reduction of molecular simulations,” *International Journal Of Robust And Nonlinear Control*, vol. 15, pp. 727-743, Oct. 2005.
- Unal, S., Oguz, C., Yilgor, E., Gallivan, M., Long, T. E., and Yilgor, I., “Understanding the structure development in hyperbranched polymers prepared by oligomeric $A_2 + B_3$ approach: comparison of experimental results and simulations,” *Polymer*, vol. 46, pp. 4533-4543, June 2005.
- Peer reviewed conference proceedings
 - Oguz, C. and Gallivan, M. A., “Identification and Evaluation of a Dynamic Model for a Thin Film Deposition Process,” *Proceedings of the 2007 American Control Conference*, (2007) 4124-4129.
 - Oguz, C. and Gallivan, M. A., “Identification of a Dynamic Model for a Thin Film Deposition Process,” *Proceedings of the 2006 IEEE International Joint Conference on Neural Networks*, (2006) 937-980.
 - Oguz, C. and Gallivan, M. A., “Dynamics of materials processing at the molecular scale,” *Proceedings of the 2004 International Symposium on Nonlinear Theory and its Applications* (2004) 135-138.

7.2 Future work

Following are the main directions for future work based on the modeling framework presented in this thesis:

1. Extension of the model reduction approach to dynamic optimization of hyperbranched polymerization: The first step for this extension would be the identification of the state for the polymer system. Other than molecular weight

and degree of branching, more descriptive variables, such as topological indices with low degeneracy, and monomer-monomer correlation functions would be necessary. Then, the state space of this polymerization process can be explored using different process conditions (e.g. solution concentration, and monomer feed rate). Exploration can be carried out by using experiments and KMC simulations. After dividing the state space into different partitions (cells), and a low order input-output model can be obtained by identifying the transitions between them. The next challenging step in this approach would be the validation of the reduced order model with more experiments. Similar to the way we have used a reduced order model for optimizing the thin film deposition, hyper-branched polymerization can also be optimized to target polymeric structures with certain molecular weight, degree of branching and other properties.

2. Recursively updating the model by systematic exploration of the state space: In our approach, the reduced order model is obtained from the training simulations with random input profiles. This may cause the model to be inaccurate in some regions of the state space, which are poorly explored. However, a systematic exploration of the state space can overcome this problem. For example, for the epitaxial deposition process, the principal components that span the state space can be updated after each surface coverage interval. This can be followed by the discretization of the state space by self organizing map (SOM). At that point, states that are significantly different from the prototype vector of their map node can be identified to continue simulations in these directions which require further exploration as described in Section 6.3. The new data, coming from the further exploration, can be used to update the dynamic model by expanding the cell map (or the dynamic model) in a recursive fashion at every coverage level until the state space is thoroughly explored, and the transitions between the explored states are found using cell mapping.

3. Dynamic optimization of molecular systems with continuous and multivariable input space: In this study, we used the dynamic model developed for the epitaxial gallium arsenide deposition for dynamic optimization of the process with only eight discrete flux settings. However, using interpolation techniques, such as k-nearest neighbor algorithm, the evolution of the state with under intermediate flux values can be predicted. Furthermore, dynamic optimization can be carried out for multivariable and continuous input space by implementing gradient optimization techniques [49].

REFERENCES

- [1] A., B. H. and M., Z., “Design of pseudorandom perturbation signals for frequency-domain identification of nonlinear systems,” in *11th IFAC Symposium on System Identification*, vol. 3, pp. 1635–1640, 1997.
- [2] AERTS, J., “Prediction of intrinsic viscosities of dendritic, hyperbranched and branched polymers,” *Computational And Theoretical Polymer Science*, vol. 8, no. 1-2, pp. 49–54, 1998.
- [3] ALLEN, M. P. and TILDESLEY, D. J., *Computer Simulation of Liquids*. Oxford University Press, 1996.
- [4] ARMAOU, A., KEVREKIDIS, I. G., and THEODOROPOULOS, C., “Equation-free gaptooth-based controller design for distributed complex/multiscale processes,” *Computers & Chemical Engineering*, vol. 29, pp. 731–740, Mar. 2005.
- [5] BALABAN, A. T., “Can topological indices transmit information on properties but not on structures?,” *Journal Of Computer-Aided Molecular Design*, vol. 19, pp. 651–660, Sept. 2005.
- [6] BATTAILE, C. C., SROLOVITZ, D. J., and BUTLER, J. E., “A kinetic monte carlo method for the atomic-scale simulation of chemical vapor deposition: Application to diamond,” *Journal Of Applied Physics*, vol. 82, pp. 6293–6300, Dec. 1997.
- [7] BERTZ, S. H. and SOMMER, T. J., “Rigorous mathematical approaches to strategic bonds and synthetic analysis based on conceptually simple new complexity indices,” *Chemical Communications*, pp. 2409–2410, Dec. 1997.
- [8] BINDAL, A., IERAPETRITOU, M. G., BALAKRISHNAN, S., ARMAOU, A., MAKEEV, A. G., and KEVREKIDIS, I. G., “Equation-free, coarse-grained computational optimization using timesteppers,” *Chemical Engineering Science*, vol. 61, pp. 779–793, Jan. 2006.
- [9] BOLTON, D. H. and WOOLEY, K. L., “Synthesis and characterization of hyperbranched polycarbonates,” *Macromolecules*, vol. 30, pp. 1890–1896, Apr. 1997.
- [10] BORTZ, A. B., KALOS, M. H., and LEBOWITZ, J. L., “New algorithm for monte-carlo simulation of ising spin systems,” *Journal Of Computational Physics*, vol. 17, no. 1, pp. 10–18, 1975.

- [11] BRAUN, M. W., RIVERA, D. E., and STENMAN, A., "A 'model-on-demand' identification methodology for non-linear process systems," *International Journal Of Control*, vol. 74, pp. 1708–1717, Dec. 2001.
- [12] BROWN, I. G., ANDERS, A., DICKINSON, M. R., MACGILL, R. A., and MONTEIRO, O. R., "Recent advances in surface processing with metal plasma and ion beams," *Surface & Coatings Technology*, vol. 112, pp. 271–277, Feb. 1999.
- [13] BRUCHMANN, B. and SCHREPP, W., "The aa* plus b*b-2 approach - a simple and convenient synthetic strategy towards hyperbranched polyurea-urethanes," *E-Polymers*, p. 014, Apr. 2003.
- [14] CAIL, J. I. and STEPTO, R. F. T., "The gel point and network formation - theory and experiment," *Polymer Bulletin*, vol. 58, pp. 15–25, Jan. 2007.
- [15] CAMERON, C., FAWCETT, A. H., HETHERINGTON, C. R., MEE, R. A. W., and MCBRIDE, F. V., "Step growth of an ab(2) monomer, with cycle formation," *Journal Of Chemical Physics*, vol. 108, pp. 8235–8251, May 1998.
- [16] CHIKKULA, Y. and LEE, J. H., "Robust adaptive predictive control of non-linear processes using nonlinear moving average system models," *Industrial & Engineering Chemistry Research*, vol. 39, pp. 2010–2023, June 2000.
- [17] CZUPIK, M. and FOSSUM, E., "Manipulation of the molecular weight and branching structure of hyperbranched poly(arylene ether phosphine oxide)s prepared via an a(2)+b-3 approach," *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 41, pp. 3871–3881, Dec. 2003.
- [18] DIUDEA, M. V. and GUTMAN, I., "Wiener-type topological indices," *Croatica Chemica Acta*, vol. 71, pp. 21–51, Mar. 1998.
- [19] DOYLE, F., I. P. R. O. B., *Identification and control using Volterra models*. New York, NY: Springer-Verlag, 2002.
- [20] DREWS, T. O., GANLEY, J. C., and ALKIRE, R. C., "Evolution of surface roughness during copper electrodeposition in the presence of additives - comparison of experiments and monte carlo simulations," *Journal Of The Electrochemical Society*, vol. 150, pp. C325–C334, May 2003.
- [21] DREWS, T. O., WEBB, E. G., MA, D. L., ALAMEDA, J., BRAATZ, R. D., and ALKIRE, R. C., "Coupled mesoscale - continuum simulations of copper electrodeposition in a trench," *Aiche Journal*, vol. 50, pp. 226–240, Jan. 2004.
- [22] DUSEK, K., DUSKOVA-SMRCKOVA, M., and VOIT, B., "Highly-branched off-stoichiometric functional polymers as polymer networks precursors," *Polymer*, vol. 46, pp. 4265–4282, May 2005.

- [23] DUSEK, K., SOMVARSKY, J., SMRCKOVA, M., SIMONSICK, W. J., and WILCZEK, L., "Role of cyclization in the degree-of-polymerization distribution of hyperbranched polymers - modelling and experiments," *Polymer Bulletin*, vol. 42, pp. 489–496, Apr. 1999.
- [24] EDWARDS, S. F. and WILKINSON, D. R., "The surface statistics of a granular aggregate," *Proceedings Of The Royal Society Of London Series A-Mathematical Physical And Engineering Sciences*, vol. 381, no. 1780, pp. 17–31, 1982.
- [25] EMRICK, T., CHANG, H. T., and FRECHET, J. M. J., "An $a(2)+b-3$ approach to hyperbranched aliphatic polyethers containing chain end epoxy substituents," *Macromolecules*, vol. 32, pp. 6380–6382, Sept. 1999.
- [26] FANG, J. H., KITA, H., and OKAMOTO, K., "Hyperbranched polyimides for gas separation applications. 1. synthesis and characterization," *Macromolecules*, vol. 33, pp. 4639–4646, June 2000.
- [27] FEAST, W. J., RANNARD, S. P., and STODDART, A., "Selective convergent synthesis of aliphatic polyurethane dendrimers," *Macromolecules*, vol. 36, pp. 9704–9706, Dec. 2003.
- [28] FLORY, P. J., "Fundamental principles of condensation polymerization," *Chemical Reviews*, vol. 39, no. 1, pp. 137–197, 1946.
- [29] FLORY, P. J., "Molecular size distribution in three dimensional polymers .6. branched polymers containing a-r-bf-1 type units," *Journal Of The American Chemical Society*, vol. 74, no. 11, pp. 2718–2723, 1952.
- [30] FRAZIER, A. B., WARRINGTON, R. O., and FRIEDRICH, C., "The miniaturization technologies - past, present, and future," *Ieee Transactions On Industrial Electronics*, vol. 42, pp. 423–430, Oct. 1995.
- [31] GALINA, H. and LECHOWICZ, J. B., "Kinetic and monte-carlo modelling of hyperbranched polymerisation," *E-Polymers*, p. 012, Mar. 2002.
- [32] GALLIVAN, M. A., *Modeling and Control of Epitaxial Thin Film Growth*. PhD thesis, California Institute of Technology, 2003.
- [33] GALLIVAN, M. A. and MURRAY, R. M., "Reduction and identification methods for markovian control systems, with application to thin film deposition," *International Journal Of Robust And Nonlinear Control*, vol. 14, pp. 113–132, Jan. 2004.
- [34] GAO, C. and YAN, D., "“ $a(2)+cbn$ ” approach to hyperbranched polymers with alternating ureido and urethano units," *Macromolecules*, vol. 36, pp. 613–620, Feb. 2003.
- [35] GAO, C. and YAN, D., "Hyperbranched polymers: from synthesis to applications," *Progress In Polymer Science*, vol. 29, pp. 183–275, Mar. 2004.

- [36] GEAR, C. W., KEVREKIDIS, I. G., and THEODOROPOULOS, C., “‘coarse’ integration/bifurcation analysis via microscopic simulators: micro-galerkin methods,” *Computers & Chemical Engineering*, vol. 26, pp. 941–963, Aug. 2002.
- [37] GILLESPIE, D. T., “General method for numerically simulating stochastic time evolution of coupled chemical-reactions,” *Journal Of Computational Physics*, vol. 22, no. 4, pp. 403–434, 1976.
- [38] GILLESPIE, D. T., “Exact stochastic simulation of coupled chemical-reactions,” *Journal Of Physical Chemistry*, vol. 81, no. 25, pp. 2340–2361, 1977.
- [39] GILLESPIE, D. T., “The chemical langevin equation,” *Journal Of Chemical Physics*, vol. 113, pp. 297–306, July 2000.
- [40] GODFREY, K., *Perturbation signals for system identification*. Hertfordshire, UK: Prentice Hall International, 1993.
- [41] GODFREY, K. R., BARKER, H. A., and TUCKER, A. J., “Comparison of perturbation signals for linear system identification in the frequency domain,” *Iee Proceedings-Control Theory And Applications*, vol. 146, pp. 535–548, Nov. 1999.
- [42] GONG, C. G., MIRAVET, J., and FRECHET, J. M. J., “Intramolecular cyclization in the polymerization of ab(x) monomers: Approaches to the control of molecular weight and polydispersity in hyperbranched poly(siloxysilane),” *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 37, pp. 3193–3201, Aug. 1999.
- [43] HANSELMANN, R., HOLTER, D., and FREY, H., “Hyperbranched polymers prepared via the core-dilution slow addition technique: Computer simulation of molecular weight distribution and degree of branching,” *Macromolecules*, vol. 31, pp. 3790–3801, June 1998.
- [44] HAO, J. J., JIKEI, M., and KAKIMOTO, M. A., “Synthesis and comparison of hyperbranched aromatic polyimides having the same repeating unit by ab(2) self-polymerization and a(2)+b-3 polymerization,” *Macromolecules*, vol. 36, pp. 3519–3528, May 2003.
- [45] HARARY, F., *Graph Theory*. Reading, MA: Addison-Wesley, 1969.
- [46] HAWKER, C. J. and CHU, F. K., “Hyperbranched poly(ether ketones): Manipulation of structure and physical properties,” *Macromolecules*, vol. 29, pp. 4370–4380, June 1996.
- [47] HAWKER, C. J., LEE, R., and FRECHET, J. M. J., “One-step synthesis of hyperbranched dendritic polyesters,” *Journal Of The American Chemical Society*, vol. 113, pp. 4583–4588, June 1991.

- [48] HERNANDEZ, E. and ARKUN, Y., “Control of nonlinear-systems using polynomial arma models,” *Aiche Journal*, vol. 39, pp. 446–460, Mar. 1993.
- [49] HILLIER, F. S. and LIEBERMAN, G. J., *Introduction to Operations Research*. New York: McGraw-Hill, 2005.
- [50] HOLLOWAY, P. H., , and MCGUIRE, G. E., *Handbook of compound semi-conductors: growth, processing, characterization, and devices*. Park Ridge, New Jersey: Noyes Publications, 1995.
- [51] HOLTER, D., BURGATH, A., and FREY, H., “Degree of branching in hyper-branched polymers,” *Acta Polymerica*, vol. 48, pp. 30–35, Jan. 1997.
- [52] HSU, C. S., *Cell-to-cell mapping : a method of global analysis for nonlinear systems*. New York: Springer-Verlag, 1987.
- [53] HULT, A., JOHANSSON, M., and MALMSTROM, E., “Hyperbranched polymers,” *Branched Polymers Ii*, vol. 143, pp. 1–34, 1999.
- [54] ITOH, M., “Atomic-scale homoepitaxial growth simulations of reconstructed iii-v surfaces,” *Progress In Surface Science*, vol. 66, pp. 53–153, Feb. 2001.
- [55] JARAIZ, M., RUBIO, E., CASTRILLO, P., PELAZ, L., BAILON, L., BARBOLLA, J., GILMER, G. H., and RAFFERTY, C. S., “Kinetic monte carlo simulations: an accurate bridge between ab initio calculations and standard process experimental data,” *Materials Science In Semiconductor Processing*, vol. 3, pp. 59–63, Mar. 2000.
- [56] JIKEI, M., CHON, S. H., KAKIMOTO, M., KAWAUCHI, S., IMASE, T., and WATANEBE, J., “Synthesis of hyperbranched aromatic polyamide from aromatic diamines and trimesic acid,” *Macromolecules*, vol. 32, pp. 2061–2064, Mar. 1999.
- [57] JIKEI, M. and KAKIMOTO, M., “Hyperbranched polymers: a promising new class of materials,” *Progress In Polymer Science*, vol. 26, pp. 1233–1285, Oct. 2001.
- [58] KARDAR, M., PARISI, G., and ZHANG, Y. C., “Dynamic scaling of growing interfaces,” *Physical Review Letters*, vol. 56, pp. 889–892, Mar. 1986.
- [59] KIM, Y. H., “Lyotropic liquid-crystalline hyperbranched aromatic polyamides,” *Journal Of The American Chemical Society*, vol. 114, pp. 4947–4948, June 1992.
- [60] KIM, Y. H., “Hyperbranched polymers 10 years after,” *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 36, pp. 1685–1698, Aug. 1998.
- [61] KIM, Y. H. and WEBSTER, O., “Hyperbranched polymers (reprinted from star and hyperbranched polymers, pg 201-238, 1999),” *Journal Of Macromolecular Science-Polymer Reviews*, vol. C42, no. 1, pp. 55–89, 2002.

- [62] KIM, Y. H. and WEBSTER, O. W., "Water-soluble hyperbranched polyphenylene - a unimolecular micelle," *Journal Of The American Chemical Society*, vol. 112, pp. 4592–4593, May 1990.
- [63] KIM, Y. H. and WEBSTER, O. W., "Hyperbranched polyphenylenes," *Macromolecules*, vol. 25, pp. 5561–5572, Oct. 1992.
- [64] KIRKPATRICK, S., GELATT, C. D., and VECCHI, M. P., "Optimization by simulated annealing," *Science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [65] KOHONEN, T., "The self-organizing map," *Proceedings Of The Ieee*, vol. 78, pp. 1464–1480, Sept. 1990.
- [66] KOHONEN, T., OJA, E., SIMULA, O., VISA, A., and KANGAS, J., "Engineering applications of the self-organizing map," *Proceedings Of The Ieee*, vol. 84, pp. 1358–1384, Oct. 1996.
- [67] KOHONEN, T., *Self-organizing maps*. New York: Springer, 1995.
- [68] KOMBER, H., VOIT, B., MONTICELLI, O., and RUSSO, S., "H-1 and c-13 nmr spectra of a hyperbranched aromatic polyamide from p-phenylenediamine and trimesic acid," *Macromolecules*, vol. 34, pp. 5487–5493, July 2001.
- [69] KRATZER, P., MORGAN, C. G., and SCHEFFLER, M., "Density-functional theory studies on microscopic processes of gas growth," *Progress In Surface Science*, vol. 59, pp. 135–147, Sept. 1998.
- [70] KRATZER, P., MORGAN, C. G., and SCHEFFLER, M., "Model for nucleation in gas homoepitaxy derived from first principles," *Physical Review B*, vol. 59, pp. 15246–15252, June 1999.
- [71] KRICHELDORF, H. R., LOMADZE, N., POLEFKA, C., and SCHWARZ, G., "Multicyclic poly(ether ester)s by polycondensation of oligo(ethylene glycol)s and trimesoyl chloride," *Macromolecules*, vol. 39, pp. 2107–2112, Mar. 2006.
- [72] KRICHELDORF, H. R. and SCHWARZ, G., "Cyclic polymers by kinetically controlled step-growth polymerization," *Macromolecular Rapid Communications*, vol. 24, pp. 359–381, Apr. 2003.
- [73] KRICHELDORF, H. R., VAKHTANGISHVILI, L., and FRITSCH, D., "Synthesis and functionalization of poly(ether sulfone)s based on 1,1,1-tris(4-hydroxyphenyl) ethane," *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 40, pp. 2967–2978, Sept. 2002.
- [74] KRICHELDORF, H. R., VAKHTANGISHVILI, L., SCHWARZ, G., and KRUGER, R. P., "Cyclic hyperbranched poly(ether ketone)s derived from 3,5-bis(4-fluorobenzoyl)phenol," *Macromolecules*, vol. 36, pp. 5551–5558, July 2003.

- [75] KRICHELDORF, H. R., ZANG, Q. Z., and SCHWARZ, G., “New polymer syntheses .6. linear and branched poly(3-hydroxy-benzoates),” *Polymer*, vol. 23, no. 12, pp. 1821–1829, 1982.
- [76] KUMAR, A. and MEIJER, E. W., “Novel hyperbranched polymer based on urea linkages,” *Chemical Communications*, pp. 1629–1630, Aug. 1998.
- [77] KUMAR, A. and RAMAKRISHNAN, S., “Hyperbranched polyurethanes with varying spacer segments between the branching points,” *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 34, pp. 839–848, Apr. 1996.
- [78] LAM, R. and VLACHOS, D. G., “Multiscale model for epitaxial growth of films: Growth mode transition,” *Physical Review B*, vol. 6403, p. 035401, July 2001.
- [79] LAROSE, D. T., *Data Mining Methods and Models*. Wiley-IEEE Press, 2006.
- [80] LAROSE, D. T., *Discovering Knowledge in Data*. Hoboken, New Jersey: John Wiley and Sons, Inc., 2005.
- [81] LEE, J. M., *A study on architecture, algorithms, and applications of approximate dynamic programming-based approach to optimal control*. PhD thesis, Georgia Institute of Technology, 2004.
- [82] LEONARD, J. A., KRAMER, M. A., and UNGAR, L. H., “A neural network architecture that computes its own reliability,” *Computers & Chemical Engineering*, vol. 16, pp. 819–835, Sept. 1992.
- [83] LIN, Q. and LONG, T. E., “Synthesis and characterization of a novel ab(2) monomer and corresponding hyperbranched poly(arylene ether phosphine oxide)s,” *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 38, pp. 3736–3741, Oct. 2000.
- [84] LIN, Q. and LONG, T. E., “Polymerization of a(2) with b-3 monomers: A facile approach to hyperbranched poly(aryl ester)s,” *Macromolecules*, vol. 36, pp. 9809–9816, Dec. 2003.
- [85] LOU, Y. M. and CHRISTOFIDES, P. D., “Feedback control of surface roughness using stochastic pdes,” *Aiche Journal*, vol. 51, pp. 345–352, Jan. 2005.
- [86] MAKEEV, A. G., MAROUDAS, D., and KEVREKIDIS, I. G., ““coarse” stability and bifurcation analysis using stochastic simulators: Kinetic monte carlo examples,” *Journal Of Chemical Physics*, vol. 116, pp. 10083–10091, June 2002.
- [87] MARTINEZ, C. A. and HAY, A. S., “Synthesis of poly(aryl ether) dendrimers using an aryl carbonate and mixtures of metal carbonates and metal hydroxides,” *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 35, pp. 1781–1798, July 1997.

- [88] MATTHES, B., BROSZEIT, E., ZUCKER, O., and GAUER, P., "Investigation of thin aln films for piezolayer-field effect transistor applications," *Thin Solid Films*, vol. 226, pp. 178–184, Apr. 1993.
- [89] MCKEE, M. G., UNAL, S., WILKES, G. L., and LONG, T. E., "Branched polyesters: recent advances in synthesis and performance," *Progress In Polymer Science*, vol. 30, pp. 507–539, May 2005.
- [90] MCQUARRIE, D. A., *Statistical Mechanics*. University Science Books, 2000.
- [91] MILLER, T. M., NEENAN, T. X., KWOCK, E. W., and STEIN, S. M., "Dendritic analogs of engineering plastics - a general one-step synthesis of dendritic polyaryl ethers," *Journal Of The American Chemical Society*, vol. 115, pp. 356–357, Jan. 1993.
- [92] MONTICELLI, O., MARIANI, A., VOIT, B., KOMBER, H., MENDICHI, R., PITTO, V., TABUANI, D., and RUSSO, S., "Hyperbranched aramids by the a(2)+b-3 versus ab(2) approach: influence of the reaction conditions on structural development," *High Performance Polymers*, vol. 13, pp. S45–S59, June 2001.
- [93] MURRAY, R. M., "Future directions in control, dynamics, and systems: Overview, grand challenges, and new courses," *European Journal Of Control*, vol. 9, no. 2-3, pp. 144–158, 2003.
- [94] NOWAK, R. D. and VANVEEN, B. D., "Random and pseudorandom inputs for volterra filter identification," *Ieee Transactions On Signal Processing*, vol. 42, pp. 2124–2135, Aug. 1994.
- [95] OGUZ, C., CAKIR, S., YILGOR, E., A., G. M., and I., Y., "Influence of polymerization procedure on the topology of highly branched polymers in a2+b3 systems: A modeling study." Submitted to Division of Polymer Chemistry for the 235th ACS National Meeting in 2008, New Orleans, LA.
- [96] OGUZ, C. and GALLIVAN, M. A., "Optimization of a thin film deposition process using a dynamic model extracted from molecular simulations." accepted for publication in *Automatica*.
- [97] OGUZ, C. and GALLIVAN, M. A., "A data-driven approach for reduction of molecular simulations," *International Journal Of Robust And Nonlinear Control*, vol. 15, pp. 727–743, Oct. 2005.
- [98] PARKER, R. S., HEEMSTRA, D., DOYLE, F. J., PEARSON, R. K., and OGUNNAIKE, B. A., "The identification of nonlinear models for process control using tailored "plant-friendly" input sequences," *Journal Of Process Control*, vol. 11, pp. 237–250, Apr. 2001.

- [99] PRICER, T. J., KUSHNER, M. J., and ALKIRE, R. C., "Monte carlo simulation of the electrodeposition of copper - i. additive-free acidic sulfate solution," *Journal Of The Electrochemical Society*, vol. 149, pp. C396–C405, Aug. 2002.
- [100] PRICER, T. J., KUSHNER, M. J., and ALKIRE, R. C., "Monte carlo simulation of the electrodeposition of copper - ii. acid sulfate solution with blocking additive," *Journal Of The Electrochemical Society*, vol. 149, pp. C406–C412, Aug. 2002.
- [101] RADKE, W., LITVINENKO, G., and MULLER, A. H. E., "Effect of core-forming molecules on molecular weight distribution and degree of branching in the synthesis of hyperbranched polymers," *Macromolecules*, vol. 31, pp. 239–248, Jan. 1998.
- [102] RAIMONDEAU, S., AGHALAYAM, P., MHADESHWAR, A. B., and VLACHOS, D. G., "Parameter optimization of molecular models: Application to surface kinetics," *Industrial & Engineering Chemistry Research*, vol. 42, pp. 1174–1183, Mar. 2003.
- [103] ROSS, S. M., *Introduction to Probability Models*. Academic Press, 2006.
- [104] RUSLI, E., DREWS, T. O., MA, D. L., ALKIRE, R. C., and BRAATZ, R. D., "Robust nonlinear feedback-feedforward control of a coupled kinetic monte carlo-finite difference simulation," *Journal Of Process Control*, vol. 16, pp. 409–417, Apr. 2006.
- [105] SMITH, D. L., *Thin-Film Deposition: Principles and Practice*. McGraw-Hill Professional, 1995.
- [106] SOMVARSKY, J. and DUSEK, K., "Kinetic monte-carlo simulation of network formation .1. simulation method," *Polymer Bulletin*, vol. 33, pp. 369–376, Aug. 1994.
- [107] SOMVARSKY, J. and DUSEK, K., "Kinetic monte-carlo simulation of network formation .2. effect of system size," *Polymer Bulletin*, vol. 33, pp. 377–384, Aug. 1994.
- [108] SPINDLER, R. and FRECHET, J. M. J., "Synthesis and characterization of hyperbranched polyurethanes prepared from blocked isocyanate monomers by step-growth polymerization," *Macromolecules*, vol. 26, pp. 4809–4813, Aug. 1993.
- [109] SRINIVASAN, S., TWIEG, R., HEDRICK, J. L., and HAWKER, C. J., "Heterocycle-activated aromatic nucleophilic substitution of ab(2) poly(aryl ether phenylquinoxaline) monomers .3.," *Macromolecules*, vol. 29, pp. 8543–8545, Dec. 1996.

- [110] STOCKMAYER, W. H., "Theory of molecular size distribution and gel formation in branched-chain polymers," *Journal Of Chemical Physics*, vol. 11, pp. 45–55, Feb. 1943.
- [111] STOCKMAYER, W. H., "Molecular distribution in condensation polymers," *Journal Of Polymer Science*, vol. 9, no. 1, pp. 69–71, 1952.
- [112] SUNDER, A., HANSELMANN, R., FREY, H., and MULHAUPT, R., "Controlled synthesis of hyperbranched polyglycerols by ring-opening multibranching polymerization," *Macromolecules*, vol. 32, pp. 4240–4246, June 1999.
- [113] TROLLSAS, M., ATTHOFF, B., CLAESSON, H., and HEDRICK, J. L., "Hyperbranched poly(epsilon-caprolactone)s," *Macromolecules*, vol. 31, pp. 3439–3445, June 1998.
- [114] TROLLSAS, M. and HEDRICK, J. L., "Hyperbranched poly(epsilon-caprolactone) derived from intrinsically branched ab(2) macromonomers," *Macromolecules*, vol. 31, pp. 4390–4395, June 1998.
- [115] TSAI, P. F., CHU, J. Z., JANG, S. S., and SHIEH, S. S., "Developing a robust model predictive control architecture through regional knowledge analysis of artificial neural networks," *Journal Of Process Control*, vol. 13, pp. 423–435, Aug. 2003.
- [116] ULTSCH, A. and SIEMON, H. P., "Kohonen's self organizing feature maps for exploratory data analysis," in *Proceedings of ICNN'90, International Neural Network Conference*, pp. 305–308, 1990.
- [117] UNAL, S. and LONG, T. E., "Highly branched poly(ether ester)s via cyclization-free melt condensation of a(2) oligomers and b-3 monomers," *Macromolecules*, vol. 39, pp. 2788–2793, Apr. 2006.
- [118] UNAL, S., OGUZ, C., YILGOR, E., GALLIVAN, M., LONG, T. E., and YILGOR, I., "Understanding the structure development in hyperbranched polymers prepared by oligomeric a(2)+b-3 approach: comparison of experimental results and simulations," *Polymer*, vol. 46, pp. 4533–4543, June 2005.
- [119] VAN BENTHEM, R. A. T. M., MEIJERINK, N., GELADE, E., DE KOSTER, C. G., MUSCAT, D., FROEHLING, P. E., HENDRIKS, P. H. M., VERMEULEN, C. J. A. A., and ZWARTKRUIS, T. J. G., "Synthesis and characterization of bis(2-hydroxypropyl)amide-based hyperbranched polyesteramides," *Macromolecules*, vol. 34, pp. 3559–3566, May 2001.
- [120] VARSHNEY, A. and ARMAOU, A., "Multiscale optimization using hybrid pde/kmc process systems with application to thin film growth," *Chemical Engineering Science*, vol. 60, pp. 6780–6794, Dec. 2005.

- [121] VARSHNEY, A. and ARMAOU, A., "Identification of macroscopic variables for low-order modeling of thin-film growth," *Industrial & Engineering Chemistry Research*, vol. 45, pp. 8290–8298, Dec. 2006.
- [122] VLACHOS, D. G., "The role of macroscopic transport phenomena in film microstructure during epitaxial growth," *Applied Physics Letters*, vol. 74, pp. 2797–2799, May 1999.
- [123] VLACHOS, D. G., "A review of multiscale analysis: Examples from systems biology, materials engineering, and other fluid-surface interacting systems," *ADVANCES IN CHEMICAL ENGINEERING*, vol. 30, pp. 2–63, 2005.
- [124] VOIT, B., "New developments in hyperbranched polymers," *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 38, pp. 2505–2525, July 2000.
- [125] VOIT, B., "Hyperbranched polymers - all problems solved after 15 years of research?," *Journal Of Polymer Science Part A-Polymer Chemistry*, vol. 43, pp. 2679–2699, July 2005.
- [126] WELDON, M. K., QUEENEY, K. T., ENG, J., RAGHAVACHARI, K., and CHABAL, Y. J., "The surface science of semiconductor processing: gate oxides in the ever-shrinking transistor," *Surface Science*, vol. 500, pp. 859–878, Mar. 2002.
- [127] WIDMANN, A. H. and DAVIES, G. R., "Simulation of the intrinsic viscosity of hyperbranched polymers with varying topology. 1. dendritic polymers built by sequential addition," *Computational And Theoretical Polymer Science*, vol. 8, no. 1-2, pp. 191–199, 1998.
- [128] WIENER, H., "Structural determination of paraffin boiling points," *Journal Of The American Chemical Society*, vol. 69, no. 1, pp. 17–20, 1947.
- [129] WIENER, H., "Relation of the physical properties of the isomeric alkanes to molecular structure - surface tension, specific dispersion, and critical solution temperature in aniline," *Journal Of Physical And Colloid Chemistry*, vol. 52, no. 6, pp. 1082–1089, 1948.
- [130] YAN, D. Y. and GAO, C., "Hyperbranched polymers made from a(2) and bb'(2) type monomers. 1. polyaddition of 1-(2-aminoethyl)piperazine to divinyl sulfone," *Macromolecules*, vol. 33, pp. 7693–7699, Oct. 2000.
- [131] YATES, C. R. and HAYES, W., "Synthesis and applications of hyperbranched polymers," *European Polymer Journal*, vol. 40, pp. 1257–1281, July 2004.
- [132] ZHOU, Z. P. and YAN, D. Y., "Distribution function of hyperbranched polymers formed by ab(2) type polycondensation with substitution effect," *Polymer*, vol. 47, pp. 1473–1479, Feb. 2006.

VITA

Cihan Oguz was born in Ankara, Turkey in 1981. He finished high school in METU College in 1998, and got his B.S. degree in Chemical Engineering from Middle East Technical University in June 2002 as an Honor Student. In August 2002, he started the PhD program in the School of Chemical and Biomolecular in Georgia Institute of Technology. At Georgia Tech, his research was focused on modeling of thin film deposition and hyperbranched polymerization processes. His dissertation title was "Control-oriented modeling of discrete configuration molecular scale processes: Applications in polymer synthesis and thin film growth". He defended his thesis on October 26, 2007. After graduating with a PhD degree in Chemical and Biomolecular Engineering, he will work as a post-doctoral researcher in Dr. Hana El Samad's group at UCSF's Biochemistry and Biophysics Department.