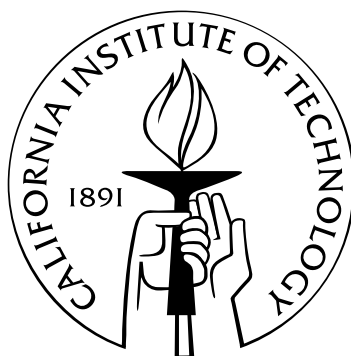


The Multistrand Simulator: Stochastic Simulation of the Kinetics of Multiple Interacting DNA Strands

Thesis by
Joseph Malcolm Schaeffer

In Partial Fulfillment of the Requirements
for the Degree of
Master of Science



California Institute of Technology
Pasadena, California

2012
(Submitted February 4, 2012)

Acknowledgements

Thanks to my advisor Erik Winfree, for his enthusiasm, expertise and encouragement. The models presented here are due in a large part to helpful discussions with Niles Pierce, Robert Dirks, Justin Bois and Victor Beck. Two undergraduates, Chris Berlind and Jashua Loving, did summer research projects based on Multistrand and in the process helped build and shape the simulator. There are many people who have used Multistrand and provided very helpful feedback for improving the simulator, especially Josh Bishop, Nadine Dabby, Jonathan Othmer, and Niranjana Srinivas. Thanks also to the many past and current members of the DNA and Natural Algorithms group for providing a stimulating environment in which to work.

There are many medical professionals to which I owe my good health while writing this thesis, especially Dr. Jeanette Butler, Dr. Mariel Tourani, Cathy Evaristo, and the staff of the Caltech Health Center, especially Alice, Divina and Jeannie.

I want to acknowledge all my family and friends for their support. A journey is made all the richer for having good company, and I would not have made it nearly as far without all the encouragement.

Finally, I must thank my wife Lorian, who has been with me every step of this journey and has shared all the high points and low points with her endless love and support.

Abstract

DNA nanotechnology is an emerging field which utilizes the unique structural properties of nucleic acids in order to build nanoscale devices, such as logic gates, motors, walkers, and algorithmic structures. These devices are built out of DNA strands whose sequences have been carefully designed in order to control their secondary structure - the hydrogen bonding state of the bases within the strand (called “base-pairing”). This base-pairing is used to not only control the physical structure of the device, but also to enable specific interactions between different components of the system, such as allowing a DNA walker to take steps along a prefabricated track. Predicting the structure and interactions of a DNA device requires good modeling of both the thermodynamics and the kinetics of the DNA strands within the system. Thermodynamic models can be used to make equilibrium predictions for these systems, allowing us to look at questions like “Is the walker-track interaction a well-formed and stable molecular structure?”, while kinetics models allow us to predict the non-equilibrium dynamics, such as “How quickly will the walker take a step?”. While the thermodynamics of multiple interacting DNA strands is a well-studied model which allows for both analysis and design of DNA devices, previous work on the kinetics models only explored the kinetics of how a single strand folds on itself.

The kinetics of a set of DNA strands can be modeled as a continuous time Markov process through the state space of all secondary structures. Due to the exponential size of this state space it is computationally intractable to obtain an analytic solution for most problem sizes of interest. Thus the primary means of exploring the kinetics of a DNA system is by simulating trajectories through the state space and aggregating data over many such trajectories. We developed the the **Multistrand** kinetics simulator, which extends the previous work by including the multiple strand version of the thermodynamics model (a core component for calculating parameters of the kinetics model), extending the thermodynamics model to include terms accounting for the fixed volume simulations, and

by adding new kinetic moves that allow interactions between distinct strands. Furthermore, we prove that our modified thermodynamic and kinetic models are exactly equivalent to the canonical thermodynamics model when the simulation is run for sufficiently long time to reach equilibrium.

The kinetic simulator was implemented in C++ for time critical components and in Python for input and output routines as well as post-processing of trajectory data. A key contribution of this work was the development of data structures and algorithms that take advantage of local properties of secondary structures. These algorithms enable the efficient reuse of the basic objects that form the system, such that only a very small part of the state's neighborhood information needs to be recalculated with every step. Another key addition was the implementation of algorithms to handle the new kinetic steps that occur between different DNA strands, without increasing the time complexity of the overall simulation. These improvements led to a reduction in worst case time complexity of a single step being just quadratic in the input size (the number of bases in the simulated system), rather than cubic, and also led to additional improvements in the average case time complexity.

What data does the simulation produce? At the very simplest, the simulation produces a full kinetic trajectory through the state space - the exact states it passed through, and the time at which it reached them. A small system might produce trajectories that pass through hundreds of thousands of states, and that number increases rapidly as the system gets larger! Going back to our original question, the type of information a researcher hopes to get out of the data could be very simple: "How quickly does the walker take a step?", with the implied question of whether it's worth it to actually order the particular DNA strands composing the walker, or go back to the drawing board and redesign the device. One way to acquire that type of information is to look at the first time in the trajectory where we reached the "walker took a step" state, and record that information for a large number of simulated trajectories in order to obtain a useful answer. We designed and implemented new simulation modes that allow the full trajectory data to be condensed as it's generated into only the pieces the user cares about for their particular question. This analysis tool also required the development of flexible ways to talk about states that occur in trajectory data; if someone wants data on when the walker took a step, we have to be able to express that in terms of the Markov process states which meet that condition.

Contents

Acknowledgements	iii
Abstract	iv
1 Introduction	1
2 System	3
2.1 Strands	3
2.2 Complex Microstate	4
2.3 System Microstate	4
3 Energy	6
3.1 Energy of a System Microstate	6
3.2 Energy of a Complex Microstate	8
3.3 Computational Considerations	9
4 Kinetics	11
4.1 Basics	11
4.2 Unimolecular Transitions	13
4.3 Bimolecular Transitions	14
4.4 Transition Rates	15
4.5 Unimolecular Rate Models	16
4.6 Bimolecular Rate Model	17
5 The Simulator : Multistrand	19
5.1 Data Structures	19
5.1.1 Energy Model	19

5.1.2	The Current State: Loop Structure	20
5.1.3	Reachable States: Moves	22
5.2	Algorithms	23
5.2.1	Move Selection	24
5.2.2	Move Update	26
5.2.3	Move Generation	27
5.2.4	Energy Computation	28
5.3	Analysis	29
6	Multistrand : Output and Analysis	32
6.1	Trajectory Mode	32
6.1.1	Testing: Energy Model	36
6.1.2	Testing: Kinetics Model	36
6.2	Macrostates	37
6.2.1	Common Macrostates	40
6.3	Transition Mode	42
6.4	First Passage Time Mode	45
6.4.1	Comparing Sequence Designs	48
6.4.2	Systems with Multiple Stop Conditions	50
6.5	Fitting Chemical Reaction Equations	51
6.5.1	Fitting Full Simulation Data to the k_{eff} model	53
6.6	First Step Mode	54
6.6.1	Fitting the First Step Model	54
6.6.2	Analysis of First Step Model Parameters	55
A	Data Structures	58
A.1	Overview	58
A.2	SimulationSystem	59
A.3	StrandComplex	60
A.4	SComplexList	61
A.5	Loop	63
A.6	Move	66
A.7	MoveTree	66

A.8	EnergyModel	68
A.9	StrandOrdering	69
A.10	Options	71
B	Algorithms	72
B.1	Main Simulation Loop	72
B.2	Initial Loop Generation	73
B.3	Energy Computation	76
B.4	Move Generation	77
B.5	Move Update	80
B.6	Move Choice	83
B.7	Efficiency	86
C	Equivalence between Multistrand’s thermodynamics model and the NU- PACK thermodynamics model.	89
C.1	Population Vectors	90
C.2	Indistinguishability Macrostates	90
C.3	Macrostate Energy	93
C.4	Partition Function	94
C.5	Proof of equivalence between Multistrand’s partition function and the NU- PACK partition function	96
C.5.1	Proof Outline	96
C.5.2	System with a single complex	97
C.5.3	System with multiple copies of a single complex type	99
C.5.4	Full System	102
D	Strand Orderings for Pseudoknot-Free Representations	106
D.1	Representation Theorem	110
	Bibliography	111

List of Figures

3.1	Secondary structure divided into loops.	8
4.1	Adjacent Microstate Diagram	12
5.1	Representation of Secondary Structures	21
5.2	Move Data Structure	23
5.3	Move Update (example)	26
5.4	Move Generation (example)	28
5.5	Full comparison vs Kinfold 1.0	30
6.1	Trajectory Data	33
6.2	Three way branch migration system	34
6.3	Trajectory Output after 0.01 s simulated time.	35
6.4	Trajectory Output after 0.05 s simulated time.	36
6.5	Example Macrostate	39
6.6	Hairpin Folding Pathways	42
6.7	First Passage Time Data, Design B	47
6.8	First Passage Time Data, Design A	49
6.9	First Passage Time Data, Sequence Design Comparison	49
6.10	First Passage Time Data, 6 Base Toeholds	50
6.11	Starting Complexes and Strand labels	52
6.12	Final Complexes and Strand labels	52
A.1	Relationships between data structure components	58
D.1	Polymer Graph Representation	107
D.2	Polymer Graph Changes (Break Move)	108

D.3 Polymer Graph Changes (Join Move) 109

Chapter 1

Introduction

DNA nanotechnology is an emerging field which utilizes the unique structural properties of nucleic acids in order to build nanoscale devices, such as logic gates [18], motors [4, 1], walkers [19, 1, 20], and algorithmic structures [14, 23]. These devices are built out of DNA strands whose sequences have been carefully designed in order to control their secondary structure - the hydrogen bonding state of the bases within the strand (called “base-pairing”). This base-pairing is used to not only control the physical structure of the device, but also to enable specific interactions between different components of the system, such as allowing a DNA walker to take steps along a prefabricated track. Predicting the structure and interactions of a DNA device requires good modeling of both the thermodynamics and the kinetics of the DNA strands within the system. Thermodynamic models can be used to make equilibrium predictions for these systems, allowing us to look at questions like “Is the walker-track interaction a well-formed and stable molecular structure?”, while kinetics models allow us to predict the non-equilibrium dynamics, such as “How quickly will the walker take a step?”. While the thermodynamics of multiple interacting DNA strands is a well-studied model [5] which allows for both analysis and design of DNA devices [24, 6], previous work on secondary structure kinetics models only explored the kinetics of how a single strand folds on itself [7].

The kinetics of a set of DNA strands can be modeled as a continuous time Markov process through the state space of all secondary structures. Due to the exponential size of this state space it is computationally intractable to obtain an analytic solution for most problem sizes of interest. Thus the primary means of exploring the kinetics of a DNA system is by simulating trajectories through the state space and aggregating data over many such trajectories. We present here the **Multistrand** kinetics simulator, which extends the previous

work [7] by using the multiple strand thermodynamics model [5] (a core component for calculating transition rates in the kinetics model), adding new terms to the thermodynamics model to account for stochastic modeling considerations, and by adding new kinetic moves that allow bimolecular interactions between strands. Furthermore, we prove that this new kinetics and thermodynamics model is consistent with the prior work on multiple strand thermodynamics models [5].

The **Multistrand** simulator is based on the Gillespie algorithm [8] for generating statistically correct trajectories of a stochastic Markov process. We developed data structures and algorithms that take advantage of local properties of secondary structures. These algorithms enable the efficient reuse of the basic objects that form the system, such that only a very small part of the state’s neighborhood information needs to be recalculated with every step. A key addition was the implementation of algorithms to handle the new kinetic steps that occur between different DNA strands, without increasing the time complexity of the overall simulation. These improvements lead to a reduction in worst case time complexity of a single step and also led to additional improvements in the average case time complexity.

What data does the simulation produce? At the very simplest, the simulation produces a full kinetic trajectory through the state space - the exact states it passed through, and the time at which it reached them. A small system might produce trajectories that pass through hundreds of thousands of states, and that number increases rapidly as the system gets larger. Going back to our original question, the type of information a researcher hopes to get out of the data could be very simple: “How quickly does the walker take a step?”, with the implied question of whether it’s worth it to actually purchase the particular DNA strands composing the walker to perform an experiment, or go back to the drawing board and redesign the device. One way to acquire that type of information is to look at the first time in the trajectory where we reached the “walker took a step” state, and record that information for a large number of simulated trajectories in order to obtain a useful answer. We designed and implemented new simulation modes that allow the full trajectory data to be condensed as it’s generated into only the pieces the user cares about for their particular question. This analysis tool also required the development of flexible ways to talk about states that occur in trajectory data; if someone wants data on when the walker took a step, we have to be able to express that in terms of the Markov process states which meet that condition.