

University of Wollongong

Research Online

University of Wollongong Thesis Collection
1954-2016

University of Wollongong Thesis Collections

2015

In silico drug discovery targeting Chikungunya virus

Phuong Thuy Viet Nguyen

University of Wollongong

Follow this and additional works at: <https://ro.uow.edu.au/theses>

University of Wollongong

Copyright Warning

You may print or download ONE copy of this document for the purpose of your own research or study. The University does not authorise you to copy, communicate or otherwise make available electronically to any other person any copyright material contained on this site.

You are reminded of the following: This work is copyright. Apart from any use permitted under the Copyright Act 1968, no part of this work may be reproduced by any process, nor may any other exclusive right be exercised, without the permission of the author. Copyright owners are entitled to take legal action against persons who infringe their copyright. A reproduction of material that is protected by copyright may be a copyright infringement. A court may impose penalties and award damages in relation to offences and infringements relating to copyright material.

Higher penalties may apply, and higher damages may be awarded, for offences and infringements involving the conversion of material into digital or electronic form.

Unless otherwise indicated, the views expressed in this thesis are those of the author and do not necessarily represent the views of the University of Wollongong.

Recommended Citation

Nguyen, Phuong Thuy Viet, In silico drug discovery targeting Chikungunya virus, Doctor of Philosophy thesis, School of Chemistry, University of Wollongong, 2015. <https://ro.uow.edu.au/theses/4452>

Research Online is the open access institutional repository for the University of Wollongong. For further information contact the UOW Library: research-pubs@uow.edu.au



School of Chemistry

***IN SILICO* DRUG DISCOVERY TARGETING
CHIKUNGUNYA VIRUS**

Phuong Thuy Viet Nguyen

MSc. in Pharmacy

**This thesis is presented as part of the requirements for the
award of the Degree of Doctor of Philosophy
of the
University of Wollongong**

March 2015

For Mom and Dad

CERTIFICATION

I, Phuong Thuy Viet Nguyen, declare that this thesis, submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the School of Chemistry, Faculty of Science, Medicine and Health, University of Wollongong, is wholly my own work unless otherwise referenced or acknowledged. The document has not been submitted for qualifications at any other academic institution.

Phuong Thuy Viet Nguyen

March 2015

PUBLICATIONS

1. Nguyen, P. T. V.; Yu, H.; Keller, P. A. Discovery of *in silico* hits targetting the nsP3 macrodomain of chikungunya virus. *J. Mol. Model.* **2014**, 20, 1-12.
2. Nguyen, P. T. V.; Yu, H.; Keller, P. A. Identification of chikungunya virus nsP2 protease inhibitors using structure-base approaches. *J. Mol. Graph. Model.* **2015**, 57, 1-8.

ACKNOWLEDGEMENTS

I would like to take this opportunity to acknowledge the people who have helped me to complete my PhD study. I feel fortunate to have completed the PhD thesis in such a friendly environment.

First and foremost, I would like to sincerely thank my supervisors Prof. Paul Keller and Dr Haibo Yu, for their guidance, advice, and support during my PhD journey. It has been a great honour to work under your supervision during the past three years. To Paul, I am always thankful to you for giving me an opportunity to join your group and all of the opportunities you have given me. Thanks so much for being encouraging, friendly, and very patient in every respect. To Haibo, thank you so much for your insights into my project. I have always considered you as one of my main supervisors, and I appreciate all of the comments and suggestions you have given me which led to enlightening discussions which helped me to figure out the problems. I have learnt a lot from you, and am glad to have had the chance to work with you.

I would not have been able to conduct my study without financial contribution from the Vietnamese Government (VIED-MOET), the University of Wollongong (IPTA) scholarships, and then the Matching Postgraduate Research Scholarship. I would also like to acknowledge High Performance Computing (HPC) of UOW for computational and research support. Many thanks to Dr Rui Yang for supporting and caring for the HPC cluster and for fixing things that needed to be fixed.

To the Keller group, I cannot record all of the memories I have of people in the group so must summarize. It is an international group with lovely people from different countries. It was really enjoyable time for me to get to know international friends. Many thanks Paul for the group meetings with research updates and weekly problems. It is very helpful and I have learnt and reviewed the chemistry knowledge from all of you. Thanks to all current and former members of the group for sharing many memorable moments, experiences, and valuable academic discussions; along with the many other aspects of life, as well as for the help and friendship. Special thanks to Mohammed for

my first days with the group; and your support and enthusiasm to help me become familiar with molecular modeling. Thanks to Adel for suggestions and support in the project. Thanks Alex for such a very lovely friendship. Thanks so much to Ari for the creative and promising plans in the future. Many thanks to Matthew for always listening and discussing problems with helped me to find the solutions. Thanks Rudi for lovely chatting in the office and the caring. Huge thanks to Jamie and Nick for helping me to double-check my thesis and giving me encouragement. Thanks also to Andrew, Steve Bailey, Steve Wales, Guy and Josh for the conversations. You all are so lovely and I enjoy the time to talk with you very much.

To the Yu group, I would like to thank Haibo for the group meetings with thanks to all people in the group as well. I learned so much computational knowledge from all of you. I will not forget Tom for always being there to help and the friendship we have shared. I appreciate it very much. Thanks to Tiatian for the help; and Andrew, Bhavani, and Kela for friendly talks.

I greatly appreciate opportunities given to me by Dr Glennys O'Brien and Dr Simon Bedford in teaching undergraduate chemistry students. It gave me a wonderful opportunity to have teaching experience in English which was something I had never done before. I learned a lot from it.

Thanks to all the people in the School of Chemistry for creating such an inspiring and welcoming environment for international students. I am grateful to Ellen Manning for her great support and enthusiasm. Thanks also to Chris Hyland for the concern and lovely chats in the school.

For some of my special friends in the school, I would like to say thanks to Zorik for many discussions about computational chemistry. I could see your passion for science and I learn from it. Thanks to Alex Martyn for a lovely friendship.

To my Vietnamese friends in Wollongong, thanks so much for your friendship over the years. Special thanks to Lien for becoming a good friend during my study here. It was

always relaxing to have chats with you about many things in study and life. Thanks, too to Trang for the friendship.

Outside the school, I am thankful to all Illawara Committee for International Students (ICIS) members for your friendship, your help, and your support. It has been a privilege and a pleasure for me to work with so many ICIS members who are very friendly and enthusiastic which has made my life a happy one here in Wollongong. I always feel relax and refresh from the PhD research after having the social chats with you. Especially, I would like to express my gratitude to Kel Magrath and Bob Colvin who always help, support, and encourage me. My heartfelt thanks to Kel for your love and caring. You are like my dad in Wollongong, thanks so much for that. Many thanks to Bob for helping me with proof-reading and English grammar. My English writing would not be good without your help.

I must not forget to thank my teachers and colleagues in the Faculty of Pharmacy, University of Medicine and Pharmacy in Hochiminh city, Vietnam (my home Univerity). Special thanks to my close colleague, Ms Huong T. Nguyen who is like my older sister, always giving me support. I must also offer thanks to her family for their love and caring. To my close friend, Quyen Duong for the friendship and caring. Also, thanks to Huyen Pham, Loi Huynh and Ha T. Nguyen for their friendships over the years. Thanks for all the sharing emails which encouraged me while I was studying in Wollongong. A big thank to Khac-Minh Thai for his encouragement from proof-reading of the papers work and friendship. It meant so much to me and gave me more motivation to pursue my research.

Finally, I will dedicate my Doctoral thesis to my family for their unconditional love, caring and encouragement. To my parents, thank you for always giving me the freedom to choose my own path, and giving me support to succeed. To my brother and my sister-in-law, thank you for being there whenever I need you. Without your support and belief in me, a higher education in a foreign country would always have remained my dream.

Phuong Thuy Viet Nguyen

March 2015

ABSTRACT

In recent years, there has been an emergence or re-emergence of Chikungunya virus (CHIKV), a member of the alphavirus. The virus is one of the arboviruses, and it is classified as a neglected tropical disease in more than 55 different countries in the world, including many African and Asian countries, Europe, Americas, and Australia. In 2008, it was listed in the US National Institute of Allergy and Infectious Disease (NIAID) category C priority pathogen due to its morbidity and mortality rates. In addition to damaging global health, the virus also imposes a huge economic burden on affected countries. However, there is currently no licensed vaccine or effective drug to combat the disease. Up to now, there have been few studies focusing on finding potential inhibitors of CHIKV. Taking advantage of all available data about CHIKV and a combination of different computational methods, this study aimed to discover and develop an approach leading to identifying inhibitors against this virus. The study targeted the non-structural proteins, nsP3 macrodomain and nsP2 protease, which play crucial roles in the viral replication and transcription (Chapter 2 and Chapter 3), and the envelope glycoprotein complexes responsible for virus entry and attachment (Chapter 4). Initially, this study searched for potential binding pockets of the CHIKV protein structures. A combination of computational tools including molecule docking, virtual screening, molecule dynamics simulations, and binding free energy calculations were used in this approach. A number of lead compounds to fight CHIKV disease were identified. The insights into the interactions between CHIKV inhibitors and their targets were elucidated. Our findings open a way which would be helpful for the further research on antiviral rational drug design, especially design of inhibitors for CHIKV and also contribute to the guidelines for the drug discovery and development.

TABLE OF CONTENTS

CERTIFICATION	ii
PUBLICATIONS.....	iii
ACKNOWLEDGEMENTS	iv
ABSTRACT.....	vii
LIST OF FIGURES	xi
LIST OF TABLES	xvi
ABBREVIATIONS	xix
CHAPTER 1. INTRODUCTION	1
1.1 Overview of Chikungunya virus	1
1.1.1 Global expansion of Chikungunya disease	1
1.1.2 The spread of CHIKV	3
1.1.3 Virus lifecycle	4
1.1.4 Clinical symptoms.....	5
1.1.5 Diagnosis.....	6
1.1.6 CHIKV genome and structure.....	6
1.2 Current treatment for CHIKV infection.....	10
1.2.1 Specific inhibitors of CHIKV	11
1.2.2 Current approach for CHIKV vaccine	15
1.2.3 Alternate approaches	17
1.3 Current research in CHIKV drug discovery	17
1.3.1 Cell-based high throughput screening approaches.....	17
1.3.2 <i>In silico</i> approaches.....	17
1.4 Structure-based drug discovery for CHIKV	18
1.4.1 Overview of computer-aided drug discovery.....	18
1.4.2 Protein-ligand docking assists in drug discovery.....	20
1.4.3 Molecular dynamics simulations in drug discovery	30
1.4.4 Binding free energy calculations.....	36
1.5 Project background and aims	43
CHAPTER 2. DISCOVERY OF INHIBITORS TARGETING CHIKV NON- STRUCTURAL PROTEIN 3.....	45
2.1 Introduction	45

2.1.1	Function and role of the nsP3 of CHIKV	45
2.1.2	Early attempts to discover CHIKV nsP3 inhibitors	46
2.2	Overview of this study	47
2.2.1	Molecular docking and virtual screening	48
2.2.2	MD simulations	49
2.2.3	Binding free energy calculations	50
2.3	Results and discussion	50
2.3.1	Docking results with ADP-ribose	50
2.3.2	Identification of inhibitors for the nsP3 macrodomain	54
2.3.3	MD simulations	61
2.3.4	Binding free energy calculations for the ligands binding to the nsP3 protein	73
2.3.5	Analysis and selecting leads for biological testing	75
2.4	Conclusions	77
CHAPTER 3. DISCOVERY OF INHIBITORS TARGETING CHIKV NON-STRUCTURAL PROTEIN 2		79
3.1	Introduction	79
3.1.1	Functional importance of the nsP2 in drug design	79
3.1.2	Early attempts to discover CHIKV nsP2 inhibitors	80
3.2	Overview of this study	83
3.3	Results and discussion	84
3.3.1	Molecular docking and virtual screening	84
3.3.2	MD simulations	95
3.3.3	Binding free energy calculations for the ligands binding to the nsP2 enzyme	109
3.3.4	Analysis and selecting leads for biological testing	111
3.4	Conclusions	112
CHAPTER 4. DISCOVERY OF INHIBITORS TARGETING CHIKV ENVELOPE GLYCOPROTEINS		114
4.1	Introduction	114
4.1.1	Envelope glycoproteins as potential targets for CHIKV drug discovery ..	114
4.1.2	Early research targeting envelope glycoproteins	115
4.1.3	Comparison of two envelope glycoprotein complexes	116

4.2	Overview of this study	117
4.3	Results and discussion	118
4.3.1	Molecular docking and virtual screenings	118
4.3.2	Identification of potential binding pockets using different methods	125
4.3.3	Analysis of interactions between hit compounds and their binding pockets 126	
4.3.4	Sampling convergence in AutoDock Vina.....	130
4.3.5	Analysis and selection of hit compounds for biological testing	131
4.4	Conclusions	133
CHAPTER 5. DISCUSSION OF FINDINGS		134
5.1	Selection of lead compounds for biological testing for CHIKV	134
5.2	Identification of potential binding sites in a protein target	137
5.3	Evaluation of computational approaches and their combination	138
CHAPTER 6. EXPERIMENTAL PROCEDURES AND METHODS		140
6.1	General procedure for molecular docking and virtual screening	140
6.1.1	Docking protocol.....	141
6.1.2	Prediction of binding sites on protein	144
6.1.3	Applying a docking procedure	145
6.2	Molecular dynamics simulations	147
6.2.1	Simulations for protein.....	149
6.2.2	Simulations for ligands	150
6.3	Binding free energy calculations.....	150
CHAPTER 7. OUTLOOK AND FUTURE DIRECTIONS		152
7.1	NSP3 macrodomain protein	152
7.2	NSP2 protease enzyme.....	154
7.3	Envelope glycoproteins.....	155
REFERENCES.....		158
APPENDICES		181

LIST OF FIGURES

Figure 1.1. Global expansion of CHIKV (taken from the website of Centers for Disease Control and Prevention (http://www.cdc.gov/chikungunya/geo/index.html) with countries and territories where chikungunya cases have been reported with endemic or epidemic, updated 24 th of February, 2015).	2
Figure 1.2. Replication cycles of CHIKV. The picture is reproduced from Rashad A. <i>et al</i> ⁸ with permission).	5
Figure 1.3. Genome organisation of CHIKV (adapted from Singh S. K. and Unni S. K. ²²).	7
Figure 1.4. Structure of the CHIKV, showing the structure of virus particles in a T=4 icosahedral symmetry (A) and the components of a single unit with a nucleocapsid consisting of spikes (made by envelope glycoproteins), virus membrane, and transmembrane (TM) helix; and a capsid covering genome RNA inside (B), (adapted from Sun S. <i>et al.</i> ⁶⁵).	10
Figure 1.5. Less common structures of some potential inhibitors for CHIKV.	14
Figure 1.6. Applications of CADD to the various stages of drug development (adapted from Tang Y. <i>et al.</i> ¹⁷⁰).	19
Figure 1.7. Representation of the strategy used for protein-ligand docking.	21
Figure 1.8. How MD simulations proceed (adapted from Durrant D. <i>et al</i> ²⁴⁶).	31
Figure 1.9. Thermodynamic cycle for calculating relative binding free energies between two ligands bound to the same protein (from Michel J. <i>et al</i> ³⁰⁰).	38
Figure 2.1. X-ray crystal structure of the macrodomain of CHIKV in complex with ADP-ribose. ⁷³	47
Figure 2.2. Schematic diagram of <i>in silico</i> approaches.	47
Figure 2.3. Re-docking ADP-ribose (A) into the active site of the nsP3: (B) The best docking pose of ligand ADP-ribose is represented as a stick model (coloured by atom type) while the protein nsP3 is shown in the solvent surface (coloured by interpolated charge with a probe radius of 1.4 Å). (C) The interactions of this pose and the nsP3 residues show hydrogen bonding interactions at the binding site of the nsP3.	52

Figure 2.4. Superimposition of the ADP-ribose after docking (in red, the top pose) and its structure in the co-crystal structure (in blue) at the active site of nsP3. The heavy-atom RMSD between the two structures is 0.6 Å.	53
Figure 2.5. Representation of three binding pockets identified in the nsP3 with top hit compounds binding in the pockets. Pocket 1 is the ADP-ribose binding site with ligand NCI_25457 (A, in burgundy), NCI_345647_a (B, in red), and NCI_61610 (C, in pink). Pocket 3 shares some residues with Pocket 1 with ligand NCI_670283 (E, in yellow). Pocket 2 is in the other side of Pocket 1 with ligand NCI_127133 (D, in dark green).	56
Figure 2.6. Structures of the top hit compounds, obtained from screenings for the nsP3.	58
Figure 2.7. Binding pose and interactions of hit compounds in the nsP3 macrodomain: (A) NCI_25457 in Pocket 1: HBs with Val113 and π - π interaction with Trp148; (B) NCI_61610 in Pocket 1: HBs with Gly112 and π - π interaction with Trp148; (C) NCI_127133 in Pocket 2: HBs with Asp133; (D) NCI_670283 in Pocket 3: Hydrophobic contacts only. The ligands (in cyan) and the residues surrounding the ligands (in grey) were displayed in sticks and coloured by atoms (carbon in cyan in ligand or in grey in residues, nitrogen in blue, oxygen in red).	61
Figure 2.8. The backbone RMSD profiles for the apo protein nsP3 and its different complexes during MD simulations: (A) Complexes of the nsP3 and top-hit compounds; (B) Complexes of the nsP3 and tenth-hit compounds.	64
Figure 2.9. RMSF values of C $_{\alpha}$ atoms of the apo protein nsP3 and its different complexes during MD simulations: (A) Complexes of the nsP3 and top-hit compounds; (B) Complexes of the nsP3 and tenth-hit compounds.	65
Figure 2.10. Hydrogen bonding interactions between the nsP3 and ligands: (A) Ligand NCI_61610 at Pocket 1 and (B) Ligand NCI_670283 at Pocket 3, with representation of ligands and key residues for interactions surrounding the ligands (in stick).	69
Figure 2.11. Superimposition of the different conformations of ligand and complexed ligand NCI_61610-nsP3 during simulation with the initial structure (red: at 0 ns, grey: at 10 ns, green: at 20 ns, pink: at 30 ns, orange: at 40 ns, and blue: at 50 ns).	71

Figure 3.1. Proposed mechanism of protease catalytic of the nsP2 (adapted from Andrew T. R. <i>et al</i>). ³⁵⁰	80
Figure 3.2. Structures of some potential inhibitors for the nsP2.....	82
Figure 3.3. Binding pose and interactions of compounds 22-25 in the nsP2 protease: (A) Compound 22 in Pocket 1: HBs with Glu1043, Lys1045 and Lys1239; (B) Compound 23 in Pocket 1: HBs with Lys1045 and His1222; (C) Compound 24 in Pocket 1: HBs with Tyr1079 and Asp1246; (D) Compound 25 in Pocket 1: HBs with Glu1204 and His1222. Compound 22, 23 and 25 were in good position in the pocket, except compound 24 was in the rear and nearly out of the pocket. The ligands (in cyan) and the residues surrounding the ligands (in grey) were displayed in sticks and coloured by atoms (carbon in cyan in ligand or in grey in residues, nitrogen in blue, oxygen in red, sulphur in orange).	86
Figure 3.4. Representation of docked structures of top hit compounds in different virtual screenings at five different binding sites of the nsP2 protease with Pocket 4 being the active site of the nsP2 protease. Ligand NCI_61610 (A) and NCI_293778 (B1) in Pocket 1; ligand NCI_293778 (B2) in Pocket 2; ligand NCI_37553 (C) in Pocket 3; ligand NCI_293778 (B3) in Pocket 4; and NCI_293778 (B4) in Pocket 5. Ligand NCI_293778 (B) with different conformations, B1-B4 could bind to different pockets Pocket 1, 2, 4, and 5, respectively.....	89
Figure 3.5. Structures of some top hit compounds for the nsP2.	92
Figure 3.6. Binding poses and interactions of hit compound NCI_293778 at different binding pockets and key residues for interactions at each pocket: (A) At Pocket 1: HBs with Lys1239; (B) At Pocket 2: π - π interactions with Tyr1177; (C) At Pocket 4: π - π interactions with Trp1084; (D) At Pocket 5: Hydrophobic contact only. The ligands (in cyan) and the residues surrounding the ligands (in grey) were displayed in sticks and coloured by atoms (carbon in cyan in ligand or in grey in residues, nitrogen in blue, oxygen in red.	93
Figure 3.7. Superimposition of three crystal structures, namely the CHIKV nsP2 protease (PDB id: 3TRK, in blue), the VEEV nsP2 protease (PDB id: 2HWK, in red), and the structure of SINV (PDB id: 4GUA, in grey). The conserved catalytic residues, cysteine and histidine (in licorice), Cys1013 and His1083 in the CHIKV nsP2 protease (in blue), Cys477 and His546 in the VEEV nsP2 protease (in red); and Cys1021 and His1098 in the SINV structure (in gray) are also shown.	95

Figure 3.8. The backbone RMSD profiles for the apo protein nsP2 and its different complexes during MD simulations from 4 ns to 53 ns: (A) Complexes of the nsP2 and top-hit compounds; (B) Complexes of the nsP2 and tenth-hit compounds.....	98
Figure 3.9. RMSFs values of C _α atoms of the apo protein nsP2 and its different complexes during MD simulations: (A) Complexes of the nsP2 and top-hit compounds; (B) Complexes of the nsP2 and tenth-hit compounds.	99
Figure 3.10. Hydrogen bonding interactions between the nsP2 enzyme and ligands: (A) Ligand NCI_37553 at Pocket 3, (B) Ligand NCI_67436 at Pocket 5; with representation of ligands and key residues for interactions surrounding the ligands (in licorice), showing the interactions maintained between the ligands and residues of protein through strong HBs interactions.....	103
Figure 3.11. The docked structure of ligand NCI_293778vst4 (in green) in Pocket 1 of the nsP2 protease, showing the residues forming Pocket 1 and the ligand-protein interactions (in grey).	104
Figure 3.12. Superimposition of the different conformations of ligand NCI_37553 and complexed with the nsP2 at Pocket 3 of the nsP2 during simulations with respect to the initial structure (red: at 0 ns, grey: at 10 ns, green: at 20 ns, pink: at 30 ns, orange: at 40 ns, and blue: at 50 ns).....	106
Figure 4.1. Structure of the envelope glycoprotein complexes: (A) The immature structure (PDB id: 3N40); (B) The mature structure (PDB id: 3N42). These structures are similar and the only difference is in the furin loop.....	117
Figure 4.2. Representation of docked structures of top hit compounds in different virtual screenings showing the location of the pockets: (A) Pocket 1 with ligand NCI_293778 (L1) and Pocket 2 with ligand NCI_67436 (L2) in the immature structure (PDB id: 3N40); (B) Pockets in the mature structure (PDB id: 3N42): ligand NCI_293778 (conformation L1) in Pocket 1, ligand NCI_67436 (L2) in Pocket 2, ligand NCI_61610 in Pocket 3 (L3), and ligand NCI_293778 (conformation L4) in Pocket 4.	123
Figure 4.3. Hydrogen bonding analysis of compounds with the immature structure in docking: (a) NCI_61610, (b) NCI_84100_a, (c) NCI_116702, (d) NCI_156219_b, (e) NCI_227186_a. The key residues involved in the interactions between glycoproteins and ligands are shown. The ligands (in cyan) and the residues surrounding the ligands (in grey) were displayed in sticks and coloured by atoms	

(carbon in cyan in ligand or in grey in residues, nitrogen in blue, oxygen in red, sulphur in organge).	128
Figure 4.4. Hydrogen bonding analysis of compound with the mature structure: (a) NCI_7524_a, (b) NCI_61610, (c) NCI_156219_b, (d) NCI_227186_b, (e) NCI_84100_b. The key residues involved in the interactions are also shown. The ligands (in cyan) and the residues surrounding the ligands (in grey) were displayed in sticks and coloured by atoms (carbon in cyan in ligand or in grey in residues, nitrogen in blue, oxygen in red, sulphur in organge).	130
Figure 6.1. A diagram of the docking procedure in AutoDock Vina.	141
Figure 6.2. A general scheme for MD simulations.	148

LIST OF TABLES

Table 1.1. Prevalence of CHIKV in terms of year and infected countries.....	2
Table 1.2. Some potential inhibitors for CHIKV.	11
Table 1.3. Potential vaccines for CHIKV.....	16
Table 1.4. Docking programs with corresponding search algorithm and scoring function.	26
Table 1.5. Strengths and weaknesses for docking and MD simulations ²⁸⁰⁻²⁸¹	35
Table 1.6. Values of the β parameter as a function of the chemical nature of the ligand, according to Hansson <i>et al.</i> ³²⁰	42
Table 1.7. Values for the β parameter in Equation 1.17, according to AlmlÖf M. <i>et al.</i> ³²⁶	43
Table 2.1. Poses in the docking of ADP-ribose into the nsP3. RMSD refers to the heavy-atom RMSD from the co-crystal structure for ADP-ribose with the nsP3. .	51
Table 2.2. Analysis of interactions between the best docked of ADP-ribose and the nsP3 macrodomain.	51
Table 2.3. Comparison of the identified hydrogen bonding interactions in the nsP3- ADP-ribose docked complex with the previously published data. In Ref [73], key residues including bonding residues (in bold), were identified by experimental work with the crystal structure of complex nsP3-ADP-ribose (3GPO) while residues in Ref [70] were determined by MD simulations of ADP-ribose in the nsP3 based on the above crystal structure.....	53
Table 2.4. Results of the top ten compounds of different virtual screens for the nsP3. The binding affinities are shown in kcal/mol.....	55
Table 2.5. Pocket residues in the nsP3 macrodomain.	56
Table 2.6. Chemical structures of five top hit compounds for the nsP3 macrodomain and their properties.....	62
Table 2.7. Hydrogen bonding analyses on the trajectories sampled in MD simulations of hit compounds for the nsP3.	66
Table 2.8. Hydrophobic contact analyses on the trajectories sampled in the MD simulations of hit compounds for the nsP3.....	67

Table 2.9. Re-docking results for complex nsP3-NCI_61610 with different conformations of the nsP3 protein taken from the different timepoints in simulations at Pocket 1.	72
Table 2.10. Virtual screening results for blind docking into Pocket 1 with different conformations of the nsP3 taken from the different timepoints in simulations. The binding affinities are shown in kcal/mol.....	72
Table 2.11. Average binding free energies (kcal/mol) of top-hit compounds and tenth-hit compounds for the nsP3 calculated by LIE method using data trajectories from the MD simulations: ΔG^1 (in kcal/mol using $\alpha = 0.18$, $\beta = 0.43$, and $\gamma = 0$) or ΔG^2 (in kcal/mol using $\alpha = 1.043$, $\beta = 0.43$, and $\gamma = 0$). ³²⁸ ΔG is the predicted binding affinity by Vina (in kcal/mol).	74
Table 2.12. Potential lead compounds for the nsP3 proposed for biological testing.	76
Table 3.1. Docking results of the best compounds for the nsP2 taken from previous study [from Singh K. D. <i>et al</i> ³⁵¹] with the binding affinities (kcal/mol).	85
Table 3.2. Results of the top ten hit compounds from the blind dockings for the nsP2. The binding affinities ΔG are in kcal/mol.....	87
Table 3.3. Results of the top ten compounds of focused dockings for the nsP2. The binding affinities are in kcal/mol.	88
Table 3.4. The important residues in each pocket of the nsP2 protease with the key residues in bold involved in forming HBs and hydrophobic contacts between the protein and ligands.	90
Table 3.5. Chemical structures of hit compounds for nsP2 and their properties.....	96
Table 3.6. Hydrogen bonding analyses on the trajectories sampled in the MD simulations for hit compounds in complexed with the nsP2 protease. The HBs with occupancy more than 10% are in highlighted in bold.....	100
Table 3.7. Hydrophobic contact analyses on the trajectories sampled in the MD simulations for hit compounds in complexed with the nsP2 protease.	102
Table 3.8. Re-docking results for complex nsP2-NCI_67436 with different conformations of the nsP2 taken from the different timepoints of simulations at Pocket 5. The binding affinity is in kcal/mol.....	107
Table 3.9. Virtual screening results for blind docking into Pocket 1 with different conformations of the nsP2 taken from the different timepoints of simulations. The binding affinities ΔG are in kcal/mol.....	108

Table 3.10. Average binding free energies (kcal/mol) of top-hit compounds and tenth-hit compounds for the nsP2 calculated by LIE method using data trajectories from the MD simulations: ΔG^1 (in kcal/mol using $\alpha = 0.18$, $\beta = 0.43$, and $\gamma = 0$) or ΔG^2 (in kcal/mol using $\alpha = 1.043$, $\beta = 0.43$, and $\gamma = 0$). ³²⁸ ΔG is the predicted binding affinity by Vina (in kcal/mol).	110
Table 3.11. Potential lead compounds for the nsP2 proposed for biological testing. ..	112
Table 4.1. Results of the top ten compounds of blind dockings with locations of pockets in the immature (PDB id: 3N40) and mature forms (PDB id: 3N42). The binding affinities ΔG are in kcal/mol.	119
Table 4.2. Results of the top ten compounds of blind dockings for the immature structure. The binding affinities are in kcal/mol.	120
Table 4.3. Results of the top ten compounds of blind dockings for the mature structure. The binding affinities are in kcal/mol.	121
Table 4.4. Results of the top ten compounds of focused dockings for the immature and mature structures. The binding affinities are in kcal/mol.	122
Table 4.5. Residues making up the pockets for two structures, the immature (PDB id: 3N40) and mature structure (PDB id: 3N42).	124
Table 4.6. Comparison of locations of identified binding pockets in both structures using different methods, blind dockings and the receptor cavities tool in Accelrys Discovery Studio program, and compared with previous study.	125
Table 4.7. Key residues for interactions in both envelope glycoprotein structures at each of binding site.....	127
Table 4.8. Potential lead compounds for the envelope glycoproteins proposed for biological testing.	132
Table 6.1. Default values of docking parameters.....	143
Table 6.2. Parameters of a grid box in different docking and virtual screenings for the nsP3 protein.....	145
Table 6.3. Parameters of a grid box in different docking and virtual screenings for the nsP2 protease.....	146
Table 6.4. Parameters of a grid box in different docking and virtual screenings for envelope glycoprotein complexes.....	147

ABBREVIATIONS

1D	First Dimension
2D	Two Dimension
3D	Three Dimension
AMBER	Assisted Model Building with Energy Refinement
CADD	Computer-Aided Drug Design
CC ₅₀	50% Cytotoxic Concentration
cDNA	Complementary DNA
CHARMM	Chemistry At Harvard Molecular Mechanics
CHIKV	Chikungunya Virus
cryo-EM	Cryo-Electron Microscopy
DNA	Deoxyribonucleic Acid
dsRNA	Double-stranded RNA
EC ₅₀	50% Effective Concentration
ECSA	East Central South African
EEEV	Eastern Equine Encephalitis Virus
elec	Electrostatic
ELISA	Enzyme-Linked Immunosorbent Assay
FEP	Free Energy Perturbation
GAFF	Generalized AMBER force field
GARD	Generally Applicable Replacement for RMSD
GROMACS	GRONingen MACHine for Chemical Simulations
GTP	Guanosine 5'-Triphosphate
H-A	Hydrogen bond acceptor
HBs	Hydrogen Bonds
H-D	Hydrogen bond donor
HSP	Heat shock protein
IC ₅₀	50% Inhibitory Concentration
ID	Identification
IgG	Immunoglobulin G
IgM	Immunoglobulin M

IMP dehydrogenase	Inosine 5'-Monophosphate Dehydrogenase
IRES	Internal Ribosome Entry Sequence
IU	International Unit
K_d	Dissociation Constant
K_i	Inhibitory Constant
L	Ligand
LIE	Linear Interaction Energy
logP	A calculated octanol-water partition coefficient
MC	Monte Carlo
MD	Molecular dynamics
MM-GBSA	Molecular Mechanics-Generalized Born Surface Area
MM-PBSA	Molecular Mechanics-Poisson Boltzmann Surface Area
MW	Molecular weight
NAMD	Not (just) Another MD Program
NCBI	National Center for Integrative Biomedical Informatics
NCI	National Cancer Institute
NIAID	National Institute of Allergy and Infectious Diseases
NMR	Nuclear Magnetic Resonance
NSAIDs	Non-Steroidal Anti-Inflammatory Drugs
NTP	N: number of particles, T: temperature, P: pressure
OAS3	Oligoadenylate Synthetase 3
OPLS	Optimized Potentials for Liquid Simulations
ORFs	Opening Reading Frames
P	Protein
PDB id	Protein Data Bank Identification Code
pKa	The logarithm of the acid dissociation constant
PL	Protein-Ligand
PME	Particle Mesh Ewald
RC	Replication Complex
RMSD	Root Mean Square Deviation
RMSF	Root Mean Square Fluctuation
RNA	Ribonucleic acid

ROC	Receiver Operating Characteristic
ROC AUCs	Area under ROC curves
RT-LAMP	Real-Time Loop-Mediated Isothermal Amplification
RT-PCR	Reverse Transcription Polymerase Chain Reaction
SASA	Solvent accessible surface area
SFV	Semliki Forest virus
SINV	Sindbis virus
siRNA	Small interfering RNA
SO	Swarm Optimization
TI	Thermodynamic Integration
TIP3P	Transferable intermolecular potential 3 points
TM	Transmembrane
UB	Urey-Bradley
US	United States
vdW	van der Waals
VEEV	Venezuelan Equine Encephalitis Virus
VMD	Visual Molecular Dynamics
VST	Virtual Screening
WNDR	Weighted Desolvation Non-Polar Ratio
ZINC	Zinc Is Not Commercial

CHAPTER 1. INTRODUCTION

1.1 OVERVIEW OF CHIKUNGUNYA VIRUS

1.1.1 GLOBAL EXPANSION OF CHIKUNGUNYA DISEASE

Chikungunya disease is caused by the chikungunya virus (CHIKV), one of the arthropod-borne viruses (arboviruses, or virus spread by mosquitoes) that has emerged¹⁻² or re-emerged³⁻⁴ in recent years. This virus is considered a neglected tropical disease; in 2008, it was listed as a category C priority pathogen by the US National Institute of Allergy and Infectious Diseases (NIAID) due to its morbidity and mortality rates.⁵ The CHIKV produces a dengue-like illness, which may lead to misdiagnosis. The disease is non-fatal,⁶ however its debilitating symptoms such as fever, rash, headache, myalgia (muscle pain), and arthralgia cause enormous health problems,⁷ affecting millions of people in nearly 55 different countries around the world.⁸ Importantly, the polyarthralgia can exist in some infected patients for months.^{7, 9-10} The history of the virus shows the first recorded case was in Tanganyika, Africa, in 1952.¹¹ An explanation for the name CHIKV is that it was derived from a local dialect, meaning “that which bends up”¹¹ in order to describe the stooped posture of patients who suffer from joint pains for weeks to years.^{6, 11} The virus has been largely neglected due to its sporadic re-emergence,^{4, 6} however it has recently attracted interest given a rise in epidemics occurring since 2006 in different countries; from Africa to Asia, Europe, Americas, and Australia, spread by infected travellers.^{6, 10, 12} The expansion of the virus has drawn global attention with the prevalence listed in Table 1.1 and Figure 1.1. Some extensive and expanding epidemics occurred in some large cities, affecting potentially millions of people. However, there is currently no cure for the CHIKV.

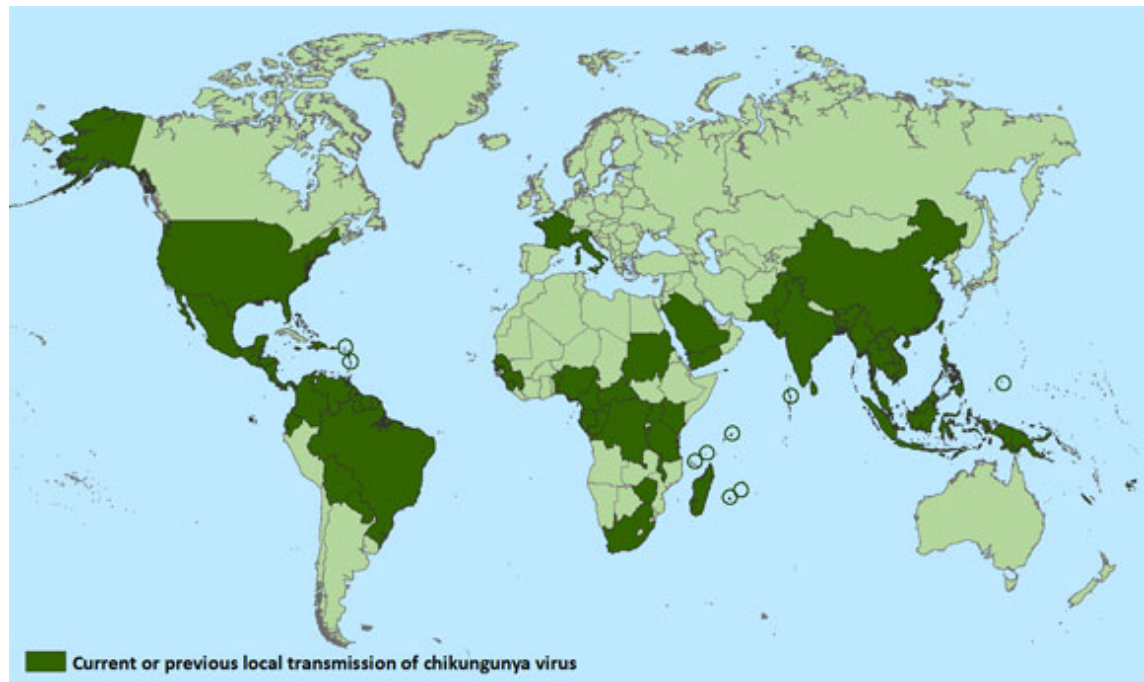


Figure 1.1. Global expansion of CHIKV (taken from the website of Centers for Disease Control and Prevention (<http://www.cdc.gov/chikungunya/geo/index.html>) with countries and territories where chikungunya cases have been reported with endemic or epidemic, updated 24th of February, 2015).

Table 1.1. Prevalence of CHIKV in terms of year and infected countries.

Year	Infected countries or regions	Outbreak
1950s	Tanganyika, ^{11, 13} Uganda, ¹⁴ Thailand, ¹⁵ Philippines ¹⁶	
1960s	India, Sri Lanka, Myanmar and Thailand, ¹⁷ Cambodia, ¹⁸ Vietnam, Laos, ¹⁹ Pakistan, Malaysia, ²⁰ Taiwan, ⁶ and Philippines ¹⁵	≥ 100,000 cases and caused 200 deaths
1970s	India and Southeast Asia ¹⁵	
1980s	Thailand, Philippines, and Indonesia ²¹	
1990s	Malaysia ²²⁻²⁴ (514,271 cases), Congo ²⁵ (50,000 cases)	
2001-2003	Indonesia, Cambodia, Vietnam, Myanmar, Pakistan, Thailand, ¹⁵ the Indian Ocean islands of the Mauritius, Mayotte, Seychelles, Madagascar, and La Reunion ²²	
2004	Comoros ⁵ (5,000 cases)	
2005-2006	The Indian Ocean islands of the Mauritius (15,760 cases), Mayotte (6,346 cases), Seychelles, and La Reunion (244,000 cases/total population of 770,000), ²⁶ Indonesia, ¹⁹ Guinea, ²⁷ India, ²⁸⁻²⁹ US ³⁰	Indian Ocean (300,000 cases with 237 deaths)

Year	Infected countries or regions	Outbreak
2006-2007	Europe, US, Australia, ³¹ Maldives (11,879 cases), Gabon (48,000 cases) ³² , Mayotte (6,346 cases, Mauritius (15,760 cases), Italy (248 cases), Sri Lanka (40,000 cases), Indonesia (15,000 cases), ³³ Cameroon, ³⁴ India ^{29, 35}	India (1.4-6.5 million)
2008-2009	India (95,000 cases in 2008; 68,000 cases in 2009), Singapore (1,033 cases), ³⁶ Thailand, ³⁷⁻³⁸ Malaysia ³⁹ (7,000 cases), Taiwan, ⁴⁰ US, ⁴¹⁻⁴² Indonesia, ³³ Bangladesh ⁴³	
2010	France ⁴⁴ , China, ⁴⁵ Thailand, ¹⁷ Canada, Myanmar ³³	
2011	India, Cambodia ¹⁸	
2012	Australia ⁴⁶ , Bhutan, ⁴⁷ Canada, Cambodia, ¹⁸ Papua New Guinea ⁴⁸	
2013	Caribbean, Canada, Thailand ⁴⁹	
2014	France, ⁵⁰ Europe, US, Caribbean ⁵¹	
2015	US ⁵²	

1.1.2 THE SPREAD OF CHIKV

CHIKV is transmitted from human-mosquito-human (urban cycle) or animal-mosquito-human (sylvatic cycle) by the bite of infected mosquitoes.²² The virus is classified into three genotypes: the Asian, West African, and East Central South African (ECSA).⁶ Genetic analyses revealed that the CHIKV originated in Africa and expanded to Asia. The spread of CHIKV from Asian and African countries to other areas is due to an association of various factors; namely worldwide distribution of the transmission vectors, climatic conditions, and inadequate mosquito control. The global expansion of this disease is solely primarily due to an increase in international travel, with disease transportation from infected travellers.^{1-3, 22} The mosquito species, *Aedes aegypti* was a principal vector in many outbreaks, while *Aedes albopictus*, the Asian tiger mosquito, has been considered a primary re-emergence factor since 2005.^{5-6, 22} The vector switch was found to be a result of reduced populations of *Aedes aegypti*, meaning that viral transmission was primarily caused by the larger *Aedes albopictus* population.³ Furthermore, a number of adaptive mutations have allowed for exploitation of the new epidemic vector.⁵³ In particular, a mutation in the E1 envelope protein, Ala226Val, was responsible for a dramatic increase in CHIKV infectivity for *Aedes albopictus* since

2005 in Africa and Asia.⁵⁴⁻⁵⁶ In addition, a substitution of lysine by glutamic acid at the position 211 of the E1 resulted in adaption of the virus to the *Aedes albopictus*.^{5, 22, 53} Recently, the E2-Ile211Thr substitution was shown to help the virus adapt to *Aedes albopictus*, as this mutation could set up the foundation for the E1-Ala226Val mutation that provides the enhanced infectivity of CHIKV.³³ Also, the E2-Gly60Asp was a determinant factor in CHIKV infectivity for both these species.⁵⁷

1.1.3 VIRUS LIFECYCLE

As with other alphaviruses, the CHIKV life cycle begins with attachment to a host cell via receptor mediated endocytosis in clathrin coated vesicles,^{8, 22, 58} a process described in Figure 1.2. Under the acidic pH of the endosome, there are conformational changes in the structures of the envelope glycoprotein complexes. The complex E1 and E2 heterodimers dissociate to form the E1 trimers. The trimers use the hydrophobic fusion loop to insert into the host membrane, and refold to form a hairpin-like structure. The nucleocapsid and viral genome are then released into the host cell cytoplasm. During translation, the nsP123 (a polyprotein precursor) from the viral genome binds to the nsP4 to form the replication complex (RC). The RC then produces the full length minus strand (negative strand RNA) required for replication. When the nsP123 concentration increases, the nP123 is cleaved to non-structural proteins nsP1, nsP2, nsP3, and nsP4. As these non-structural proteins and the host cell proteins serve as the plus strand (positive strand) in replication, they produce the 26S subgenomic RNAs and genomic (49S) RNAs. The 26S subgenomic positive stranded RNA encodes the polyprotein precursor for structural proteins. The cleavage process takes place in the Golgi complex, and then the products (non-structural proteins and structural proteins) are transported to the plasma membrane. The viral RNA is packaged into the nucleocapsid, and the mature virions bud out of the plasma membrane.

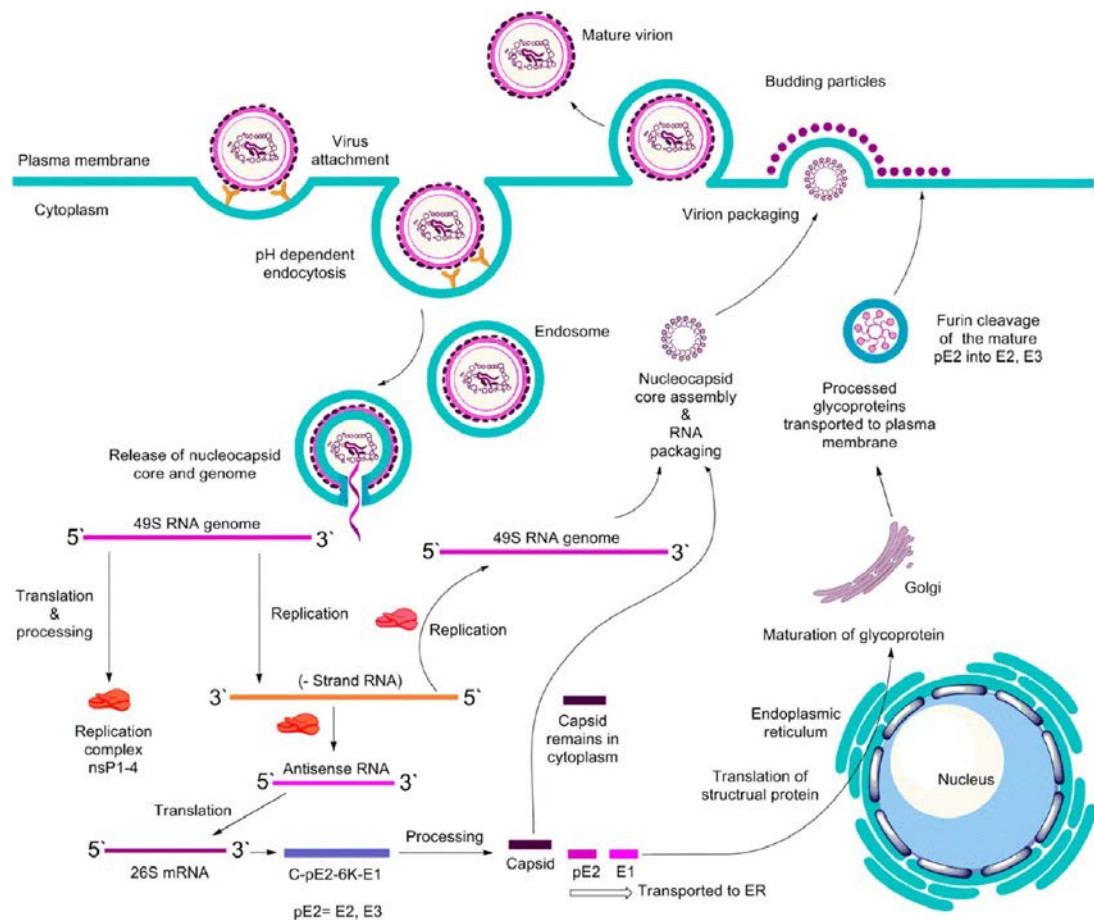


Figure 1.2. Replication cycles of CHIKV. The picture is reproduced from Rashad A. *et al*⁸ with permission).

1.1.4 CLINICAL SYMPTOMS

Following the transmission, CHIKV replicates and expands to liver and joints.⁵ There are similar symptoms with other viruses, such as dengue fever, yellow fever, or Ross River fever. After an incubation period of 2-4 days, acute clinical symptoms start with the onset of high fever, rigors, headache, fatigue, nausea, vomiting, myalgia, and rash.⁵ In particular, intense joint pain (polyarthralgia) is the most characteristic which incapacitate patients.⁵ This acute stage lasts from 1-10 days. The chronic stage depends on the acute stage, and is characterized by polyarthralgia. It can last from weeks to years,⁸ which affects individual patients, and results in social health impacts¹⁰ due to debilitating infection in large working population. Neurological disorders⁵⁹ and eye infection⁶⁰ are also reported in infected patients, while meningoencephalitis and

haemorrhagic manifestations such as haematemesis and malaena leading to death have been reported.⁶

1.1.5 DIAGNOSIS

Since no vaccines or effective therapeutics are available, early detection and proper diagnosis are becoming increasingly important in treating the CHIKV disease. Viral culture is considered a gold standard for CHIKV diagnosis. It is based on inoculation of mosquito cell cultures, mammalian cell cultures, or mice.⁶¹ However, for a rapid diagnosis, techniques such as a detection of reverse transcription polymerase chain reaction (RT-PCR), real-time RT-PCR, real-time loop-mediated isothermal amplification (RT-LAMP) are recommended.^{1, 62} More frequently, serodiagnostic methods,⁶³⁻⁶⁴ such as enzyme-linked immunosorbent assay (ELISA), indirect immunofluorescent method, hemagglutination inhibition, or neutralization techniques were used effectively as reliable techniques in the identification and characterization of CHIKV. The detection of immunoglobulin M (IgM) and immunoglobulin G (IgG) antibodies using capture enzyme-linked immunosorbent assay (MAC-ELISA) is a rapid and reliable technique in serology.

1.1.6 CHIKV GENOME AND STRUCTURE

A full understanding of the structure and genome of CHIKV is crucial for the development of drugs to combat the virus. CHIKV is a positive-sense and single-stranded RNA virus in the alphavirus genus, *Togaviridae* family.⁶⁵ In 1984, Simizu *et al* used African and Asian strains of CHIKV to analyse the structural proteins of CHIKV.⁶⁶ In 2002, Khan *et al* identified the full genomic sequence of CHIKV (S27, African prototype),⁶⁷ opening up further investigations into the elucidation of the structure and genome of this virus. CHIKV's genome includes two opening reading frames (ORFs) that consist of 11,805 nucleotides in total without including the cap at the 5' end, a 1-poly (A) tract, and a poly (A) tail at 3' end.⁶⁷ The first ORF has 2474 amino acids encoding non-structural proteins (nsP1, nsP2, nsP3, and nsP4) at the 5' region, while the second consists of 1244 amino acids encoding structural proteins (the capsid C, envelope glycoproteins E1, E2, E3, and 6K). Between these two ORFs, there is a

junction region. In addition, there is a 7-methylguanosine group capped at the 5' end and the polyadenylated group at the 3' end. The arrangement of the genome can be as follows: 5'-cap-nsP1-nsP2-nsP3-nsP4-(junction region)-C-E3-E2-6K-E1-poly(A)3' (Figure 1.3).

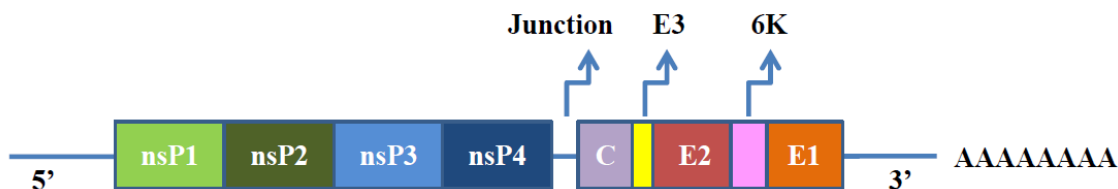


Figure 1.3. Genome organisation of CHIKV (adapted from Singh S. K. and Unni S. K.²²).

The role and function of non-structural and structural proteins of CHIKV have been examined based on information of other alphaviruses, such as Sindbis virus (SINV) and Semliki Forest virus (SFV). According to these studies, the non-structural proteins of CHIKV play an important role in the formation of the transcription/replication complex of the virus and the negative strand synthesis.^{58, 68-70}

Studies of SFV and SINV nsP1 indicate that nsP1 is a multifunctional protein.⁷¹ It is involved in the formation of the virus cap,^{70, 72-73} and directing RNA replication complex to membranes or liposomes,⁷⁴ as well as associating with endosomes and lysosomes at the cytoplasmic surface of membranes.^{72, 75} It has also been implicated with both guanine-7-methyltransferase and guanylyltransferase activities.⁷⁴

Investigations of nsP2 from SINV and SFV show that the nsP2 has multiple enzymatic activities.⁷⁶⁻⁷⁷ Primarily, it belongs to the papain superfamily of cysteine protease.⁷⁸ It has functions of proteolytic enzyme at C-terminal domain whereas at the N-terminal domain, it possesses the activities of ATPase and GTPase,⁷⁶⁻⁷⁷ RNA helicase,⁷⁹ and RNA triphosphatase.⁸⁰ Moreover, it is involved in the regulation of synthesis of the 26S subgenome⁸¹⁻⁸² and activates the switch from early to late stages,⁸² as well as playing a role in the translocation of 50% of the translated nsP2 into the nucleus.⁸³ The nsP2 plays a vital role in the replication of virus; which combined with the exhibition of some degree of sequence specificity,⁷⁶ has resulted in it being an attractive target for drug

design. More recently, the crystal structure of nsP2 protease of CHIKV was solved (PDB id: 3TRK).

The specific functions, roles, and activities of the nsP3 protein remain relatively elusive; however, genetic and functional analysis of the nsP3 from SINV revealed that it is a phosphoprotein participating in the process of synthesis of the minus strand and the subgenome of the virus.^{69, 73, 84} It has been reported that deletion of phosphorylation sites in the SFV nsP3 decreases the level of RNA synthesis.⁸⁵ The nsP3 protein consists of two domains; the N-domain which is highly conserved, and the C-domain which is not.⁷³ There is an “X domain”, or macrodomain (a first 160 amino acid domain with unknown function), present at the N-terminal region of the nsP3 with a determined crystal structure (PDB id: 3GPG).⁷³ The structure on its conserved adenosine binding site was also obtained (PDB id: 3GPO).⁷³ In addition, it was found that the mutation of amino acids at the position Asn10 and Ala24 in ADP-ribose binding of the nsP3 macrodomain in SINV affects the replication and viral RNA synthesis, though it has no effect on the binding region.⁸⁶⁻⁸⁷

Several attempts have been made on SINV to propose the roles of the nsP4 in the RNA-dependent RNA polymerase activity⁸⁸ and in replication, and transcription of the virus.^{88,91} The results showed that the N-terminal Tyr residue of the nsP4 of SINV may be substituted with Phe, Trp, or His without changing the wild-type phenotype in cultured cells. However, other substitutions, except for Met, were lethal or quasilethal.⁹¹ The nsP4 is also stable and remains active during the infection cycle.⁸⁹ Previous research on the mutants reported that the mutations in the nsP4 were Glu191 substituted for Leu and Glu315 to Gly, Val, or Lys, together with one mutation in the nsP1 (Thr349 to Lys); which suppress the minus strand RNA synthesis.⁹¹ Arg183 of SINV nsP4 polymerase was found to have an important role in alphavirus minus strand RNA synthesis.⁹²

Most studies of the structural proteins of CHIKV are based on the biology and pathogenesis of the virus, as infection is mediated by these glycoproteins. The E2 is responsible for receptor binding, while membrane fusion is supported by the E1.⁹³ In 2008, Santhosh *et al*⁵⁶ found the mutation in E1 Ala226Val, leading to the epidemic outbreaks of CHIKV in India. Ongoing insights into the structure of CHIKV have been

revealed by two crystal structures of surface glycoprotein complexes; namely precursor p62-E1 heterodimer (PDB id: 3N40), and the mature E3-E2-E1 (PDB id: 3N42) determined in 2010.⁹³ The E3 protein plays an important role in the proper folding of p62 and the formation of the p62-E1 heterodimer,⁹⁴⁻⁹⁵ but the E3 is not in the component of mature CHIKV.⁶⁶ The 6K associates with the complex p62-E1 and is transported to the plasma membrane before assembly.⁶⁵ The assembly process takes places due to the interactions of the genomic RNA and the nuclear capsid protein. The 6K protein facilitates particle morphogenesis, but it is not stoichiometrically incorporated into virions.⁶⁵ The structural change of envelope glycoproteins in membrane fusion was investigated by Li *et al.*⁹⁶ The resulting roles of the E2 in receptor interactions, and the existence of epitopes, are essential for the design of a vaccine against this virus.

Recently, structural analysis of CHIKV at pseudo-atomic level resolution was reported,⁶⁵ by combining electron cryo-microscopy (cryo-EM) techniques for the whole virus and X-ray crystallography for the component of structural proteins, together with the published crystal structure of the CHKV E1-E2 glycoprotein heterodimer.⁹³ A 5.3 Å resolution cryo-EM map of CHIKV-like particles was interpreted, and the mechanisms of neutralization of antibodies were proposed. The study revealed that like other alphaviruses, CHIKV has an icosahedral spherical structure with T=4 quasi-icosahedral symmetry (diameter of about 60-70 nm, Figure 1.4 A). This structure consists of 80 spikes: including 20 icosahedral “i3” spikes (located on the icosahedral 3-fold axes), and 60 quasi-3-fold “q3” spikes (located in general positions) with a quasi-3-fold axis.⁶⁵ The spikes i3 and q3 are significantly different, possibly indicating different stages of generation of fusogenic E1 trimers. The complete q3 spike combines with one-third of an i3 spike to form a single T=4 icosahedral asymmetric unit (Figure 1.4 A). These spikes are made by the envelope glycoproteins E1, E2, and together with virus membrane, transmembrane (TM) helix, and a capsid covering genome RNA to form a nucleocapsid (Figure 1.4 B). The E1 glycoprotein is composed of 439 amino acids including 404 N-terminal residues, 30 residues which comprise the TM helix, and 5 amino acids which form the cytoplasmic domain, as well as an N-linked glycosylation site at Asn141. The E1 is divided into three domains, namely domain I, II, and III, with domain I located in between domain II and III. The E2 consists of 364 residues, a 26

residue TM helix, and 33 residues cytoplasmic domain. The E2 also has three domains, domain A, B, and C, which are known as distinct immunoglobulin-fold domains. Domain A takes part in receptor-binding process, lying in between domain B and domain C. Domain B covers the fusion loop in domain II of E1 in the mature structure. Under the pH acidic environment, the virus becomes fusogenic by combining three E1 to make a trimer. The fusion loop is then exposed to insert into the host cell membrane. Therefore, hiding the “fusion loop” may interfere with the virus entry and infection of human tissue.

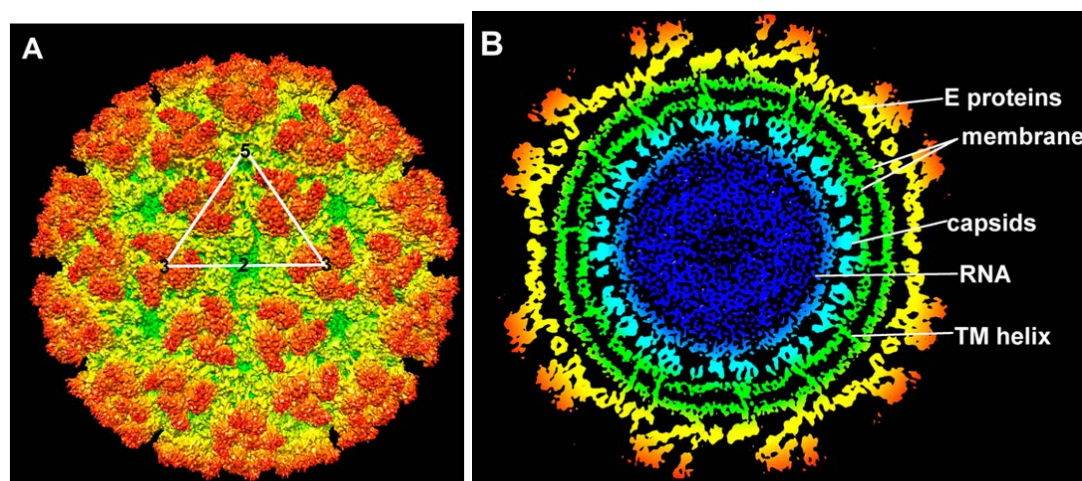


Figure 1.4. Structure of the CHIKV, showing the structure of virus particles in a T=4 icosahedral symmetry (A) and the components of a single unit with a nucleocapsid consisting of spikes (made by envelope glycoproteins), virus membrane, and transmembrane (TM) helix; and a capsid covering genome RNA inside (B), (adapted from Sun S. *et al.*⁶⁵)

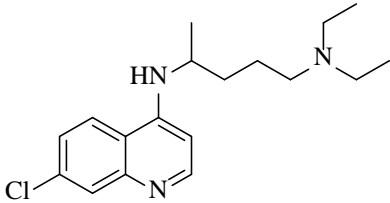
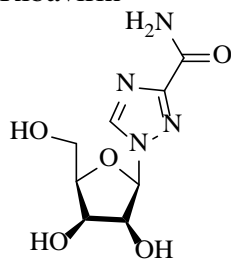
1.2 CURRENT TREATMENT FOR CHIKV INFECTION

There is no effective vaccine or antiviral agent currently available for CHIKV infection.¹⁰ Treatment is mostly based on alleviating symptoms by using analgesics (non-salicylate analgesics), antipyretics,^{1, 6} anti-inflammatory agents (corticosteroids and non-steroidal anti-inflammatory (NSAIDs));^{6, 9, 21} along with taking bed rest, and undertaking an extra fluid diet.¹ Some agents for the treatment of acute CHIKV include Methotrexate; NSAIDs: Rofecoxib, Celecoxib, Parecoxib; corticosteroids: Prednisolon; antirheumatic: Sulfasalazine and Methotrexate; non-salicylate analgesics: Paracetamol, Morphine; traditional herbal medicine: *Fernelia* spp species.

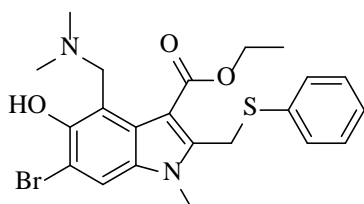
1.2.1 SPECIFIC INHIBITORS OF CHIKV

While some inhibitors have been shown to be effective against CHIKV *in vitro*, currently, there is no antiviral drug available for CHIKV treatment. Further biological testing (*in vivo*) is required. These inhibitors are listed in Table 1.2, along with the current update of inhibitory effects and the clinical trials.

Table 1.2. Some potential inhibitors for CHIKV.

Compound	Inhibitory effects and current update
<p>Chloroquine</p>  <p>(1)</p> <p>IC₅₀ = 1.1 µg/ml</p>	<p>Chloroquine (1) was reported as <i>in vitro</i> antiviral compound more than 35 years ago.⁹⁷ However, recently, a mouse model showed that this compound enhances virus replication, and aggravates the disease.^{9, 98} Chloroquine phosphate has also been used for chronic chikungunya arthritis with anti-inflammatory, rather than antiviral effects.⁹⁹⁻¹⁰¹ Some studies revealed that chloroquine was considered as an entry inhibitor as it could interact with the endosome-mediated internalization process in the infection cycle.^{98, 102} Chloroquine was tested in Phase III clinical trials in France in 2006; however, this compound was terminated in 2007 with no anti-CHIKV effect.⁸</p>
<p>Ribavirin</p>  <p>(2)</p> <p>EC₅₀ = 83.3 µg/ml</p>	<p>Ribavirin (2) is well-known as an antiviral inhibitor <i>in vitro</i>.¹⁰³ The mechanism is varied between different viruses, such as interaction with intracellular viral RNA production, inhibition of inosine 5'-monophosphate dehydrogenase (IMPDH), leading to depletion of cellular GTP pools, and action as a potent mutagen for some RNA viruses (error catastrophe mechanism), even though the precise mechanism is still unclear.⁹ Ribavirin can be used either alone or in combination with α-interferon to get a synergistic effect <i>in vitro</i> against CHIKV (concentration of α-interferon 3.9 IU/ml, and ribavirin 18.75 µg/mL).¹⁰⁴⁻¹⁰⁵ However, up to now, no evidence on clinical efficacy of ribavirin on CHIKV or ribavirin in combination with α-interferon for anti-CHIKV is reported.⁹</p>
<p>Arbidol</p>	<p>Arbidol (3) was developed 20 years ago in Russia for treatment of acute respiratory infections.⁸ This compound elicits a broad effect on RNA, DNA, enveloped, and non-enveloped viruses.¹⁰⁶ The mechanism may be interference with the mediated</p>

Arbidol

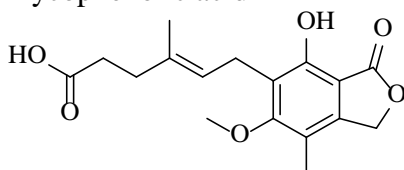


(3)

 $IC_{50} = 12.2 \mu M$

fusion,¹⁰⁷ blocking the viral entry into the target cells through inhibition of glycoprotein conformational changes.¹⁰⁸⁻¹¹⁰ In 2011, arbidol and its derivatives were used for *in vitro* testing for CHIKV. The results reported that arbidol could inhibit CHIKV but very weak effects even though the IC_{50} value is much lower than the toxic concentration ($IC_{50} = 12.2 \mu M$, $CC_{50} \geq 200 \mu g/mL$).¹¹¹ Evidence from *in vivo* studies is required to validate activity of arbidol on CHIKV.⁸

Mycophenolic acid

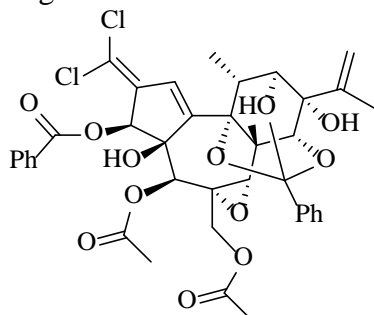


(4)

 $IC_{50} = 0.2 \mu M$

Mycophenolic acid (4) was discovered approximately 100 years ago.¹¹² It is an inhibitor for the enzyme IMPDH involved in *de novo* biosynthesis of guanine nucleotide.¹¹³ It was known as having good activity with antiproliferation, anticancer, as an antiviral agent, and an immunosuppressant.⁸ Recently, this compound showed inhibition of the CHIKV replication using virus-induced cell death.¹¹⁴ However, the compound was reported as suffering from a metabolic drawback associated with rapid conjugation of the C-7 phenolic hydroxyl group with glucuronic acid,¹¹³ thus *in vivo* studies are required.⁸

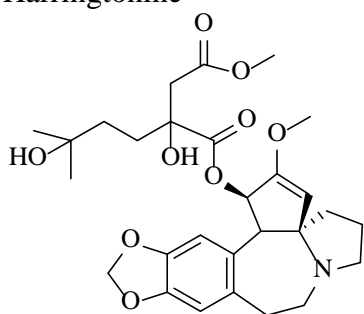
Trigocherrin A



(5)

 $EC_{50} = 1.5 \mu M$

Trigocherrin A (5) is a natural compound, isolated from the bark of *Trigonostemon cherrieri* Veillon (*Euphorbiaceae*),¹¹⁵ a tree in New Caledonia, or the species found in tropical Asia, India, and Sri Lanka to New Guinea.⁸ Recently, in testing CHIKV inhibitory effect, this compound showed inhibition of viral replication function on virus-induced cell death, using a virus-cell-based assay.¹¹⁶

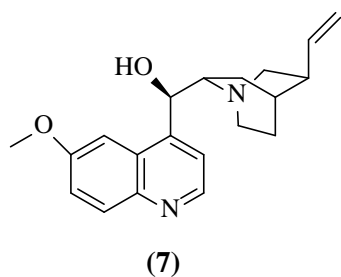
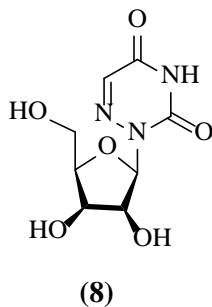
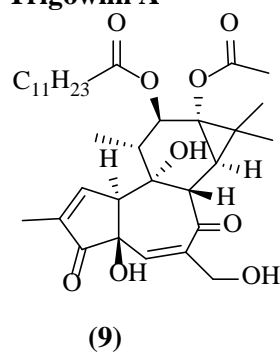
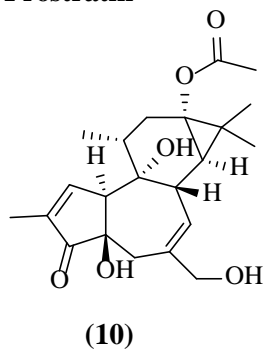
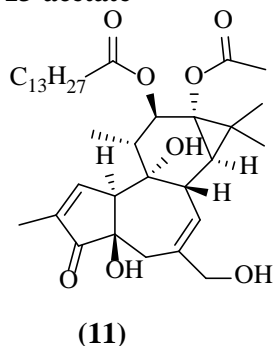
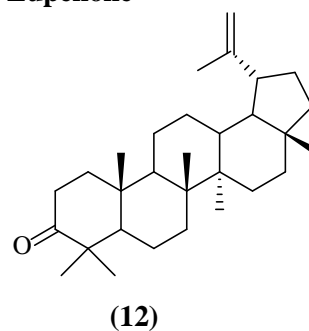
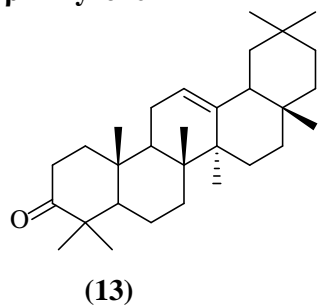
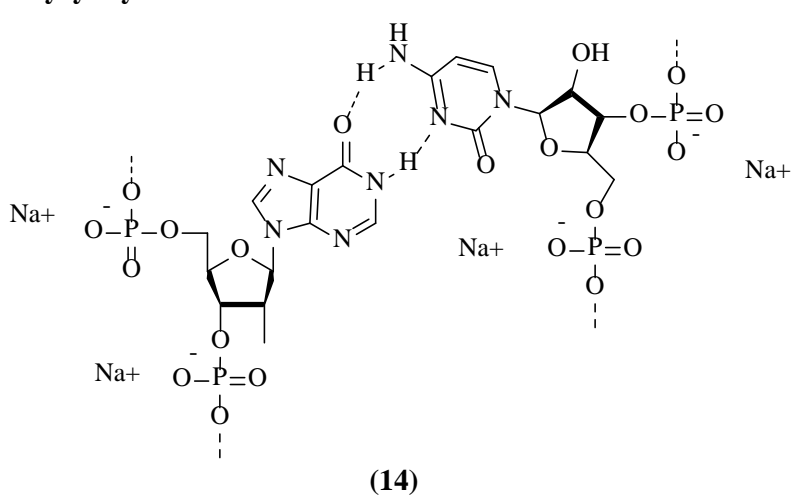
Harringtonine¹¹⁷

(6)

 $EC_{50} = 0.24 \mu M$

Harringtonine (6) is an alkaloid compound, isolated from *Cephalotaxus harringtonia* tree in Japan.⁸ It displayed CHIKV inhibition, by affecting the early stages of infection after cellular endocytosis.¹¹⁷ Also, it was found to affect the CHIKV RNA production inside the infected cell and viral protein expression through the nsP3 and the E2 proteins.¹¹⁷ This compound is still undergoing *in vivo* testing.⁸

Some other compounds are less common, thus requiring further testing to be considered promising lead compounds (Figure 1.5). However, they were not approved for human use. For example, quinine (**7**) inhibited the *in cellulo* growth of CHIKV.⁹ Increasing concentrations of quinine affected the nsP1, as mutations were observed suggesting an impairment of the function of the viral guanylyltransferase activity. Another compound, named 6-azauridine (**8**), showed its inhibition on both DNA and RNA virus replication, and orotidine monophosphate decarboxylase, an enzyme involved in the *de novo* biosynthesis of pyrimidine, cytidine, and thymidine.¹¹⁸ Trigowiin A (**9**) and prostratin (**10**) were isolated from the bark of *Trigonostemon howii* of the *Euphorbiaceae* species in central Vietnam during the biological testing for CHIKV.¹¹⁹ This compound showed a weak anti-CHIKV activity, however, it possessed its selectivity for CHIKV, against the SINV and SFV. Prostratin¹¹⁹ showed better inhibitory effects on CHIKV than (**9**) in a CHIKV inhibition assay. 12-o-tetradecanoylphorbol 13-acetate (**11**) was reported to inhibit CHIKV through the activation of the signal transduction enzyme protein kinase C, also selective for CHIKV inhibition.¹¹⁹ Lupenone (**12**) and β -Amyrone (**13**) are isolated from the leaves of *Anacolosia pervilleana* (Madagascan plant). (**12**) showed moderate anti-CHIKV activity in a virus-cell-based assay,¹²⁰ while (**13**) had a moderate anti-CHIKV activity in the assay. Polycytidylic acid or poly(I:C) (**14**) is a synthetic double-stranded RNA (dsRNA) analogue. It displayed an immunostimulant action as an inducer for the interferon via interaction with the toll-like receptor 3 in CHIKV infection assay.¹²¹

Quinine**6-Azaauridine****Trigowin A****Prostratin****12-O-Tetradecanoylphorbol 13-acetate****Lupenone****β-Amyrone****Polycytidylic acid****Figure 1.5.** Less common structures of some potential inhibitors for CHIKV.

1.2.2 CURRENT APPROACH FOR CHIKV VACCINE

With the widespread distribution of CHIKV, there is a need for a safe and effective vaccine. Currently, there is no effective vaccine for CHIKV although there have been some efforts on this approach. An understanding of antibody-mediated and cell-mediated immune responses is important for vaccine development.⁵ Unfortunately, there is very little information pertaining to the interaction of CHIKV infection and adaptive immune system in re-infection.^{5, 8} Some potential vaccines which need to be developed are listed in Table 1.3.¹²²

Table 1.3. Potential vaccines for CHIKV.

Type of vaccine	Vaccine	Production	Developed year	Status
Inactivated vaccines	Formalin-inactivated vaccine ¹²³	A whole virus grown in monkey cell cultures	1970s	Phase II: discontinued
	Tween 80-ether extraction ¹²⁴	A whole virus grown in monkey cell cultures	1970s	Preclinical
	Vero adapted formalin inactivated vaccine ¹²⁵	A virus on Vero cells	2006	Preclinical
Live-attenuated vaccines	181/clone 25 vaccine strain ¹²⁶	Serial passage CHIKV strain in culture cells	1986	Phase II
	TSI-GSD-218 CHIKV vaccine ¹²⁷	An attenuated strain CHIKV	2000	Underway Phase III
Genetically engineered vaccines	Chimeric virus vaccine ¹²⁸	VEEV (attenuated vaccine strain TC-83), an attenuated strain of EEEV or SINV and the structural protein genes of CHIKV	2008	Preclinical
	Recombinant adenovirus vaccine ¹²⁹	A non-replicating complex adenovirus vector encoding the structural polyprotein cassette of the CHIKV	2010	Preclinical
	CHIKV-IRES vaccine ¹³⁰⁻¹³¹	A cDNA clone generated from the wild-type La Reunion strain	2010	Preclinical
	DNA vaccines ¹³²⁻¹³³	Single or three individual plasmids	2008	Preclinical
	Subunit protein vaccines ¹³⁴⁻¹³⁷	Recombinant CHIKV envelope proteins	2012	Preclinical
	Virus-like particle vaccine ¹³⁸⁻¹³⁹	Non-infectious virus-like particles coated with the same protein	2010	Completed Phase I

1.2.3 ALTERNATE APPROACHES

Some other approaches are targeting virus entry and maturation, viral nucleic acids, and cellular receptors. For example, agents used on each interference can be listed as follows: using furin inhibitors, or decanoyl-RVCR-cholormethyl ketone to impair the maturation of the E2 glycoprotein;^{9, 140} small hairpin RNA molecules to interfere RNAs;¹⁴¹ 2',5'-oligoadenylate synthetase (OAS3),¹⁴²⁻¹⁴³ cellular IMPDH enzyme,¹¹⁴ and visperin¹⁴⁴ in controlling CHIKV replication or human antibodies. Recently, silencing of HSP-90 (a chaperone protein related to heat shock) using small interfering RNA (siRNA) has been shown to disturb CHIKV replication in cultured cells.¹⁴⁵

1.3 CURRENT RESEARCH IN CHIKV DRUG DISCOVERY

1.3.1 CELL-BASED HIGH THROUGHPUT SCREENING APPROACHES

Recently, high throughput screening methods have been developed to identify potential CHIKV inhibitors. For example, a CHIKV replicon and a concomitant screen with SFV surrogate infection model were used to screen 356 natural compounds, and clinically approved drugs.¹⁴⁶ A cell-based high throughput screening assay using resazurin against a kinase inhibitor library of 4,000 compounds, combined with the image-based high content assay approach was applied.¹⁴⁷ A phenotypic assay was also used to identify one natural compound that partially blocks nsP2 activity and inhibits CHIKV replication *in vitro*.¹⁴⁸

1.3.2 IN SILICO APPROACHES

In recent years, *in silico* approaches are promising to save time and the cost of the drug discovery process. With the availability of the crystal structures of several proteins of the CHIKV genome and other related alphaviruses, structure-based approaches using molecular docking and virtual library screening can be applied to identify potential inhibitors (hits) from a large database of compounds. These hit compounds may be experimentally tested, or used to investigate a structure-activity relationship to optimize

the compound's activity. Generally, drug targets are key proteins which is identified through cellular and protein biochemical processes associated with the disease. These biomolecules are known as being involved in signalling or metabolic pathways that are specific to the disease process¹⁴⁹⁻¹⁵⁰ For CHIKV, non-structural proteins and envelope glycoproteins were considered potential targets.⁸ Non-structural proteins of the virus play an integral role in viral replication and transcription, so they are attractive targets for designing potent inhibitors of CHIKV.¹⁵¹ Recently, two crystal structures, namely the nsP2 protease (PDB id: 3TRK) and the nsP3 macrodomain (PDB id: 3GPG)⁷³ have been available for use as a starting point in antiviral research. Homology models for structure nsP4 was also proposed to provide structures for drug design.¹⁵² The envelope glycoprotein complexes were determined by X-ray crystallography, the immature form (PDB id: 3N40) and the mature form (PDB id: 3N42).¹⁵³ These complexes are also of interest as targeting the envelope glycoproteins can affect the virus entry.

1.4 STRUCTURE-BASED DRUG DISCOVERY FOR CHIKV

1.4.1 OVERVIEW OF COMPUTER-AIDED DRUG DISCOVERY

In fighting disease, drug discovery and development are very expensive and time-consuming processes.¹⁵⁴ The drug discovery process involves different stages such as target identification (target ID), hit identification (hit ID), lead discovery and optimization, biological testing (preclinical trials), and clinical trials.^{149, 155} Therefore, utilizing computational techniques not only increases the efficiency and effectiveness of research, but also saves time and reduces costs; in particular during lead discovery and lead optimization.

The field of computer-aided drug design (CADD) was started in the 1960s and in the late 1980s up to now, this has been growing and has become an integral part of drug discovery with the development of computer hardware and software, and the increasing availability of protein structures of biochemical targets of pharmaceutical interest (Figure 1.6).¹⁵⁶ There have been a number of studies reported on describing specific computational methods, clarifying their role and importance, and highlighting recent

advances, as well as their successful applications and challenges.¹⁵⁶⁻¹⁶⁹ These contain many applications of this approach in antiviral drug design, in particular in analysis of viral protein target. The prediction and simulation of conformational, steric, and physicochemical properties helps to elucidate the characteristics of the target. The binding pockets and interactions may be identified using computational approaches. Rationalization of drug action and virtual screening to identify potential inhibitors, together with the structure-activity relationship of these ligands, were applied.¹⁶³

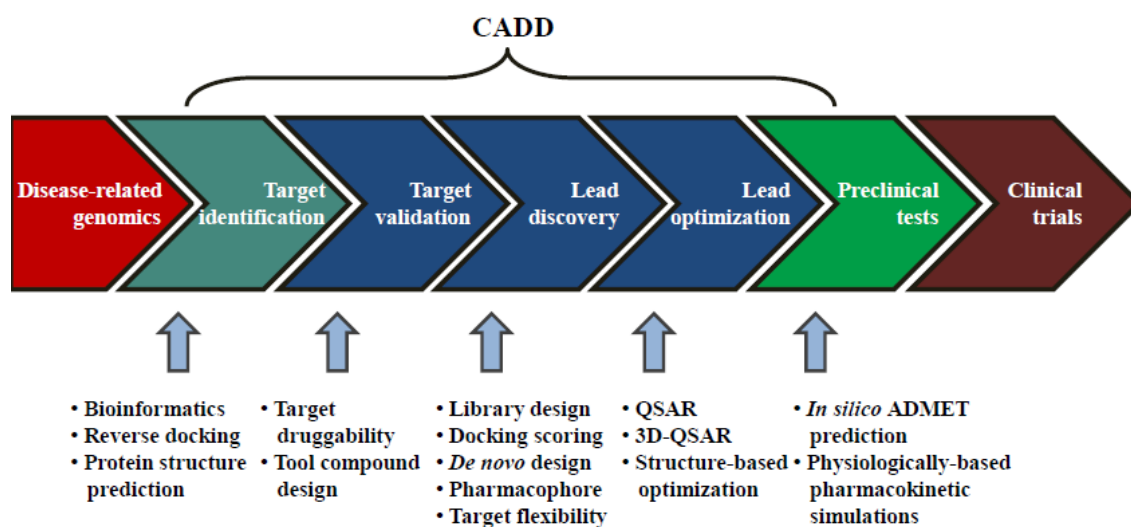


Figure 1.6. Applications of CADD to the various stages of drug development (adapted from Tang Y. *et al.*¹⁷⁰).

Computational approaches in CADD can be classified into two main groups: the ligand-based approaches and the structure-based approaches.

Ligand-based approaches: This is an indirect approach which is useful in the case of a lack of an experimental receptor structure together with the difficulty of designing a reliable model.¹⁷¹ Ligand-based methods are particularly valuable in the early stage drug discovery.¹⁶³ These approaches rely on known active compounds and utilize their similar descriptors taken from their molecular characteristics, properties, and biological activity data to elucidate the structural and physicochemical properties of the ligands.¹⁶³ Molecular characteristics may be physicochemical descriptors in one dimension, 1D (e.g., molecular weight, atom counts, logP, and pKa); two dimensions, 2D (e.g., topological descriptors); and three dimensions, 3D (physicochemical properties such as location, constraints, and shape descriptors). Ligand-based approaches introduce a definition of a pharmacophore as the ensemble of steric and electronic features,

necessary to account for the common molecular interactions with the protein target, and to trigger (or to block) its biological response.^{154, 163} A pharmacophore model describes the three-dimensional chemical features, using pharmacophoric descriptors such as hydrogen bond donors, hydrogen bond acceptors, hydrophobic, aromatic, positive ionizable groups, and negative ionizable groups. There are many different ways to build a pharmacophore model. It can be based on chemical structures of known active compounds from different chemical scaffolds or diverse chemical structures for compounds (with IC₅₀ or K_i values ranging over more than three orders of magnitude). Similarities of molecular properties such as pharmacophore features, shaped-based models or a quantitative structure-activity relationship are also utilized. The underlying assumption here is that ligands having similar physiochemical properties are likely to show comparable activity spectra.

Structure-based approaches: Unlike the ligand-based methods, structural data of a protein target is a prerequisite in a structure-based approach.¹⁶⁶ The full three-dimensional (3D) structure of the protein may be obtained from X-ray crystallography or nuclear magnetic resonance (NMR) techniques. In the case where the structure of protein is unknown, homology modeling can be used based on the similarity of the genomic sequences with other viruses.¹⁵⁰ Protein-ligand docking and structure-based virtual screening are examples of these approaches. They utilize information of the protein structure to identify and optimize drug candidates by examining molecular interactions between ligands and target macromolecules, as detailed below.¹⁵⁴

1.4.2 PROTEIN-LIGAND DOCKING ASSISTS IN DRUG DISCOVERY

Pioneered in the early 1980s,¹⁷² protein-ligand docking has developed as an integral part of drug discovery.¹⁷³ It has become an invaluable tool that assists efficiency in drug discovery to understand interactions of protein-ligand complexes.¹⁷⁴⁻¹⁷⁶ In the 2000s, a series of reviews summarized the methodology, and highlighted the successful applications, recent advances, and challenges of this approach.^{162, 166-167, 173, 177-187}

Given a small molecule and a protein target of a virus, docking attempts to insert the small molecule into a binding site of the protein target. In other words, the aim of

docking is to get the “best match” of a protein-ligand complex,¹⁷⁸ or to accurately predict the orientation and conformation of a small molecule (ligand), i.e. the lowest binding energy (known as binding mode or binding pose), and then estimate binding affinity of the ligand into a known structure (Figure 1.7).

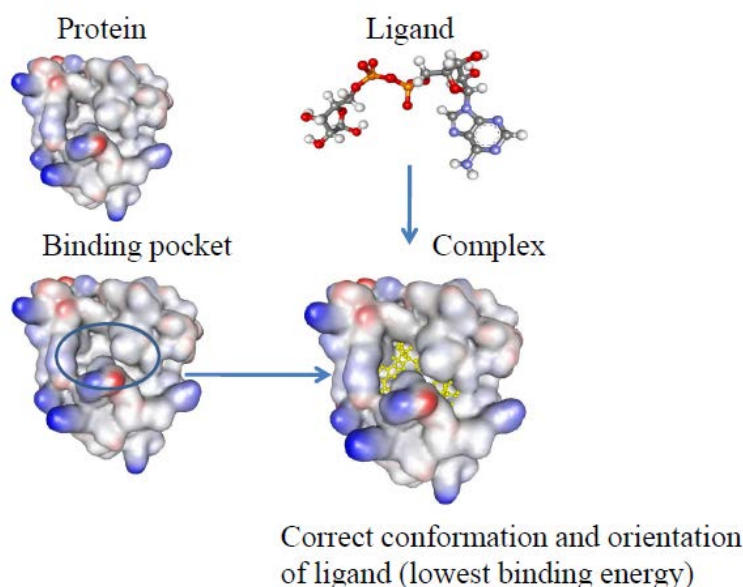


Figure 1.7. Representation of the strategy used for protein-ligand docking.

Therefore, protein-ligand docking includes two processes: docking (geometric sampling of potential ligand/protein binding mode); and scoring (using an equation and specific parameters to estimate a ligand binding affinity). Consequently, a protein-ligand docking program consists of these two essential components: sampling (ligand sampling and protein flexibility) and scoring. There are three main kinds of sampling setups: first, the protein and ligand are kept rigid, only translational and rotational manipulations of the ligand are investigated; second, the protein is rigid and ligand is assumed to be flexible (all degrees of freedom of the ligand are explored); and third, the protein is fully or partly flexible as well as the ligand. Some popular docking programs are DOCK¹⁸⁸, AutoDock,¹⁸⁹ FlexX,¹⁹⁰ GOLD,¹⁹¹ and GLIDE;¹⁹² most of which utilise the flexible ligand and a rigid receptor approach.

To make protein-ligand docking more practical, improvements from docking with rigid structure to partly or full flexible receptor have been undertaken during the last years.^{173,}

^{176-179, 183, 186} However, protein flexibility is still the major hurdle where ligand sampling

is concerned. In addition, there has been no universally applicable scoring function available. Therefore, finding the most appropriate algorithm for the prediction of protein-ligand binding of each specific target is a major focus of research. Many search algorithms and scoring functions have been developed over the years, alongside numerous scoring functions, taking into consideration speed and accuracy.¹⁷³

1.4.2.1 Search algorithms

As previously mentioned, most docking programs account for ligand flexibility. In order to search for precise conformation and configuration, an optimal search algorithm should explore all of degrees of freedom to sample sufficiently so as to include the true binding modes. These algorithms may be classified into three main groups: systematic searches (incremental construction, conformational search, exhaustive or databases), random or stochastic methods (Monte Carlo, genetic algorithms or evolutionary algorithm, Tabu search), and deterministic or simulation searches (molecular dynamics and energy minimization).^{173, 176-177, 179, 186} Some algorithms implement a hybrid approach combining two or all three of the methods.

Systematic search: In this search algorithm, all of the degrees of freedom of ligand are investigated. It can systematically rotate all rotatable bonds of a ligand through 360° using a fixed increment, to obtain all possible combinations for evaluation. This approach is, however, susceptible to a combinatorial explosion in computational cost. Another approach is that the ligand is divided into different fragments or a core fragment and flexible parts (sidechains), and then these components are docked into the active site. The results will be covalently linked together (other incremental search algorithms) or rebuilt the ligand from fragments (Hammerhead algorithm). Libraries of pre-generated conformations may be utilized to tackle the combinatorial explosion problem.

Random (or stochastic) search: The search makes random changes for a single ligand or a population of ligands. A pre-defined probability function is used to evaluate the ligand. The common groups are Monte Carlo (MC) and genetic algorithms. In the MC method, following random changes and energy minimization for the generated conformations, the Boltzmann probability function is used as the criteria of evaluation.

If a change in temperature is also combined to increase the probability, this is called Simulated annealing. In the genetic algorithms or evolutionary algorithms, the theory of the evolution in biological system is applied to search for the correct ligand binding mode.¹⁶² The Tabu search is another method to explore areas of conformational space, and the value of root mean square deviation (RMSD) between the calculated molecular coordinates, and every molecule's previously recorded conformation should be less than a cutoff value to be accepted. Swam optimization tries to search for a search space for the ligand using Swam Intelligent method by taking the information of the best positions of its neighbors.¹⁶²

Deterministic (simulation) method: These algorithms use molecular dynamics (MD) simulation or energy minimization to explore the conformations. The ensembles of populations are generated and are docked rather than a single conformer. MD is the most popular approach, however has weaknesses in its inability to cross high-energy barriers within simulation timescale, leading to arrest of ligands in a local minima of the free energy surface.¹⁷⁷ To overcome this, simulations can be carried out for the different parts of a protein-ligand system at different temperatures, or ligand starting positions. Energy minimization is often used in combination with other approaches.

The search algorithms for protein flexibility in docking are still limited, though there are some reviews and studies on it.¹⁹³⁻¹⁹⁷ MD and MC, are usually applied for a search of protein conformations. Another strategy is the use of rotamer libraries to model protein conformational space based on experimental data, and favourable sidechain conformations. A protein ensemble grid is different approach in which the algorithm uses an ensemble of conformations of protein for docking rather than a single one, and then maps them on a grid representation.

1.4.2.2 Scoring function

For a docking process to be successful, the adopted scoring function is the deterministic factor for obtaining an accurate prediction of conformation of a protein-ligand complex, and a correct ranking of final structures. In other words, scoring is used to predict the binding affinities, and differentiate between correct and incorrect orientations, and

conformations (poses). The ideal docking program should satisfy both computational efficiency and reliability. Therefore, speed and accuracy are two key components of a scoring function. Free-energy simulation can be based on atomistic MD simulation, applied for quantitative modeling of protein-ligand interactions, and binding affinity predictions; however, this approach is expensive. Most docking programs still do not include entropic effects explicitly. There are different scoring functions which are classified in terms of shape, chemical complementation, force fields, empirical results or system knowledge. The three main categories are force-field-based, empirical, and knowledge-based scoring functions.^{162, 166, 173, 176-177, 179, 181-182, 186, 198-199} Beside these scoring functions, some other approaches to improve scoring such as consensus scoring (combining the information from multiple scoring functions), or clustering, and entropy-based scoring methods.¹⁶²

Force-field-based scoring: This type of scoring function tries to model many types of interactions involved in protein-ligand binding by utilizing physics-based functional forms, and parameters derived from experiments or quantum mechanical simulations.¹⁶⁶ There are different types of force field scoring functions. They have similar functional forms but their force field parameters are different, for example G-Score with Tripos force field,²⁰⁰ and AutoDock with AMBER force field.¹⁷⁶ Most of them only consider a single protein conformation to make it possible to omit calculation of the internal protein energy. In molecular mechanics force fields, the binding free energy is calculated as a sum of two energies, receptor-ligand interaction energy (van der Waals (vdW) and electrostatic interactions) and internal ligand energy (steric strain induced by binding). In AMBER force fields, the vdW term is referred to a Lennard-Jones potential function (such as 12-6 Lennard-Jones).¹⁶⁶ The general AMBER force field (GAFF) including some parameters is suitable for simulating small molecules.²⁰¹ The CHARMM force field is similar to AMBER force fields, but has some additional terms. The CHARMM22 force field is usually used for modelling protein. The recent CGenFF force field can be applied as a general force field for small molecules.²⁰² The disadvantage of standard force-field-based scoring function is that it does not include solvation, and entropic term explicitly.¹⁷⁶

Empirical scoring: This type of scoring uses a sum of various empirical energies to estimate the binding energy. A set of weighted empirical energy terms may be composed of vdW, electrostatic, hydrogen bonding energy, desolvation term, entropy, and hydrophobicity term.¹⁶² This scoring is far simpler than force-field-based scoring, however, it depends on the molecular data sets used to perform regression and fitting, leading to reproducibility of the experimental binding affinity data.

Knowledge-based scoring: The binding affinity can be calculated based on the information of experimental protein-ligand complexes data. This is estimated as a sum of free energies of protein-ligand atom-pair interactions (the potential mean force). The frequencies or probability distributions of interatomic distances between two atoms are converted into distance-dependent interaction free energies of protein-ligand atom pairs using the inverse Boltzmann method.¹⁶²

Some popular docking programs together with their searching algorithms and their scoring functions are summarized in Table 1.4.

In conclusion, there are many docking programs with different scoring functions. However, no single one is the best. The ideal scoring function is still not available. The problem is that all of the scoring functions mentioned above, are based on different assumptions and simplifications, and do not fully take into account of entropic and solvation effects. Compromises between the conformational searching algorithms and scoring functions could improve docking algorithms, however, it may not improve binding affinity prediction. Thus, it is difficult to obtain a chemical accuracy, and the results depend on the specific system.

Table 1.4. Docking programs with corresponding search algorithm and scoring function.

Docking program	Fees	Ligand searching algorithm	Scoring function
DOCK ¹⁸⁸	No (academics); Yes (profit)	Incremental build	Force field or contact score
AutoDock ¹⁸⁹	No (academics); Yes (profit)	Lamarckian algorithm	Force field
AutoDock Vina ²⁰³	No (academics); Yes (profit)	Lamarckian algorithm	Combining knowledge-based potentials and empirical scoring functions
FlexX ¹⁹⁰	Yes	Fragmentation and Incremental construction	Empirical score
GOLD ¹⁹¹	Yes	Genetic algorithm	Empirical score
Glide ¹⁹²	Yes	Exhaustive search (Monte Carlo)	Empirical score
FRED ²⁰⁴	Yes	Conformational ensembles (Rigid body docking)	Gaussian score or empirical scores
ICM ²⁰⁵	Yes	Pseudo-Brownian sampling and local minimization (Metropolis Monte Carlo)	Mixed force field and empirical score

1.4.2.3 Setting up a docking protocol

To perform docking, preparations of protein and ligand are indispensable. For the target protein, the structure is usually obtained by X-ray crystallography or by NMR structure determination. The structure may be apo, holo (complexed with another compound) or if the structure is not available, it may be predicted by threading or homology modeling. If the function of protein is unknown, it is crucial to search for possible binding sites in the structure (discussed below). Some crucial factors should be considered to check the structure carefully. Initially, structural integrity is needed to check. Polar hydrogen atoms are also added. The next step is to assign proper protonation and tautomeric states of ionizable residues including aspartate (Asp), glutamic acid (Glu), arginine (Arg), lysine (Lys), and histidine (His). The orientation of asparagine (Asn) and glutamine (Gln) residues also require checking. The active site or binding site needs to be defined before docking, and the treatment of water molecules in docking is considered. Water molecules may affect the formation of the complex as they can form hydrogen bonds with the ligand and the protein. Geometry refinement of the protein/receptor-ligand is required to correct small artefacts in the protein or ligand such as strained bond

lengths/angles or intermolecular steric clashes.¹⁸⁷ Energy minimization of the structure should be used. Preparation of the ligand, like the protein, requires great care. Ligands can be taken from various sources, such as a database or electronic vendor catalogues. Filtering using the drug-like properties is recommended to eliminate molecules, with unfavourable properties: such as poor solubility, pharmacokinetics characteristics relating to adsorption, distribution, metabolism and excretion; oral bioavailability and toxicity. The most popular filter used in drug design is Lipinski “rule of five”,²⁰⁶ that is, molecular weight (MW) lower than 500 daltons, logP less than 5, number of hydrogen bond donors less than 5 and number of hydrogen bond acceptors less than 10.²⁰⁷ Usual, compounds are represented in 2D format and then converted to a 3D representation. Several methods are available for generation of a 3D structure. Special care must take into account of tautomeric and protonation states, stereochemistry of chiral centres. The next step is docking with sampling to generate different binding poses; scoring and ranking. Most docking programs as mentioned previously, focus on docking with rigid docking and flexible ligands. Protein flexibility may be treated for some specific sidechains of protein residues.

1.4.2.4 Prediction of binding sites

In structure-based drug design with molecular docking and virtual screening approaches, the question of where in the structure of the protein the ligand binds is of interest. Therefore, it is a requirement to understand the structure and function of a protein target, in particular, knowledge of locations that small molecules (ligands) could bind in the structure is of key importance to help rationally design of ligands, to fit with high binding affinities and specificity that can modulate the functions of this target.²⁰⁸⁻²¹⁰ A binding process is the sum of many contributions such as environment (pH, ionic strength), and presence of water molecules, in which the shape complementarity of protein and ligand together with their physicochemical properties are balanced.²⁰⁹ Proteins can change conformations upon ligand binding, which may influence the steric accessibility of a binding cleft and can interfere with the ability of an algorithm to identify a potential binding site.²⁰⁸ Small molecules are known to usually bind in the largest pockets on the surface of protein, and their binding free energies is a result of enthalpy-entropy compensation.²⁰⁹

There have been many different approaches developed to identify binding sites of a protein, depending on whether the co-crystallised structure of a protein-ligand complex is known or not.²⁰⁹ Thus, these approaches can use a co-crystallisation of a protein with a ligand if available, use structural or sequence similarity with a known binding site, or use *in silico* prediction methods. For computational methods, some algorithms utilize protein surface analysis to identify ligand binding pockets in the protein. The others use probe clustering and energy contour analysis which is based on analysis of binding energies of probes placed on a grid around the protein.²¹⁰ Other types of characteristics such as surface accessibility, the net charge on the protein residues in a protein as a function of pH, and sequence conservation can be used. MD simulations are also used to generate dynamic ensembles of protein conformations for binding site detection. In general, computational approaches for prediction of binding sites on a protein can be classified into three main groups: geometric-based (using geometrical pocket description or free accessible volume calculation to identify the sterically favoured cavities among all clefts on the protein's surface); energetic-based (finding energetically favoured positions, e.g. by calculation of interactions with different probes and the protein's surface); and knowledge-based (using structure and sequence comparison to generate templates of similar sequence or functionality to identify evolutionarily conserved regions to rank cavities generated with geometric approaches). The geometric-based methods are fast and easy to use while the others are more time-consuming and require user expertise.²⁰⁹

1.4.2.5 Structure-based virtual screening

The most notable application of protein-ligand docking is structure-based virtual screening.²¹¹ It has become increasingly important to improve the speed and efficiency of drug discovery and development.²¹¹⁻²¹⁵ The main purpose is to reduce the large number of compounds from databases and select the most promising compounds for biological testing.¹⁷⁶ It is known as a fast tool for drug design which is valuable to discover lead compounds complementary to experimental methods, for example high-throughput screening. A large volume of studies have shown its importance and successful application in this field of drug design and development.²¹¹⁻²²⁴ Some popular programs available for this purpose are DOCK,¹⁸⁸ FlexX,¹⁹⁰ GOLD,^{191, 225} AutoDock^{189,}

²²⁶ and AutoDock Vina.^{203, 227} However, based on docking, virtual screening has similar challenges in accuracy of scoring and ranking. If many ligands are docked, the time for scoring and ranking in screening a large sample is of utmost important. More reliable scoring functions should include calculations for nonbonded interactions such as cation- π interaction, CH- π interaction, and π - π stacking interactions.²²⁴ The choice of ligand libraries for screening is also of concern. Several databases are available, for instance, the National Cancer Institute (NCI) Diversity Set (dtp.nci.nih.gov/branches/dscb/diversity_explanation.html), NCBI PubChem (pubchem.ncbi.nlm.nih.gov), eMolecules (www.emolecules.com) and ZINC²²⁸ composed of a variety of compounds or libraries of natural products,²²⁹ metabolome,²³⁰⁻²³¹ and nutraceuticals.²³² If in the case of a new target, where no information about binding sites for ligands is known, blind docking for the entire protein is applied to identify sites that ligands bind tightly. Other computational approaches can be used to predict the binding sites as discussed previously. Validation of the approaches is the most critical due to the inaccuracies of the scoring functions, which will affect the results of ranking. The main factors should be taken into account are the quality of the obtained docking poses, and the ability of the methods to discriminate known active and inactive compounds after docking in the same target. False-positive hits or decoy molecules have similar physical structures but are chemically distinct from ‘true’ hits, and can be used as competitive binders to a protein.²³³⁻²³⁵

1.4.2.6 Validation of docking method

Docking is usually validated by the ability to reproduce the experimental data in predicting the binding pose, and binding affinity to distinguish the active and inactive compounds. The commonly accepted criterion for docking success is pose selection by a comparison of the RMSD between the docked structure (top scoring pose) and the experimental structure (the co-crystal structure). With the protein-ligand co-crystal structure, a docking protocol can be evaluated and improved. However, the errors existing in the bound structures due to the poorly refined ligand geometries could lead to misleading interpretation of key binding interactions. They can be errors in the ligand structure (for instance missing atoms, incorrect bond orders or other connectivity issues); or incorrect bond distances, angles or dihedral angles; or conformational errors

(for example, *cis*- or twisted amides, distorted rings, non planar aromatic groups, or groups of planar structures of not being planar); incorrect orientations or bad steric clashes between protein and ligand. Therefore, it has become commonly to apply protein preparation tools to structure.²³⁶ If the RMSD value is no more than 2 Å, it is considered a docking success.^{167, 179, 184, 237-238} Moreover, other criteria to validate docking include Generally applicable replacement for RMSD (GARD),²³⁸ enrichment factors, receiver operating characteristic (ROC) factors,²²⁷ area under ROC curves (ROC AUCs),²³⁹ and binding affinity calculations.²⁴⁰ Using a decoy set of inactive compounds to dock and after ranking by scoring, enrichment is calculated by dividing the ratio of active molecules in the entire dataset by the ratio of active molecules in the top 5 or 10% of the top ranked molecules. The enrichment plot or ROC curves are plotted. The sensitivity of a given docking/scoring combination and specificity are shown in ROC plots.

1.4.3 MOLECULAR DYNAMICS SIMULATIONS IN DRUG DISCOVERY

Given the structure of a macromolecule and its complexes, molecular dynamics (MD) simulation is one of the computational approaches to investigate the motions of a system of particles.²⁴¹ MD simulation of biological macromolecules (proteins) of interest was introduced in 1977.²⁴² Since then, this technique has become a powerful computational method, quickly showing wide application to many fields, and constant improvement quantitatively and qualitatively in structure-based drug discovery to understand drug-receptor interactions.²⁴³⁻²⁴⁶ Especially with the recent advances in algorithms and computer hardware and software, long-timescale from microsecond to millisecond for thousands of atoms has been achieved.²⁴⁷ The interaction for the particles is calculated and included sequentially to the time of simulating. The result of the process (MD trajectory) provides information on the atomic level of positions, velocities and energies. Statistical mechanics related to distribution and motion of atoms are required to connect microscopic simulations with macroscopic observables such as changes of conformations, binding free energy, and mechanism of reaction. MD simulation has been applied to a wide range of biological problems, including protein folding,²⁴⁸⁻²⁴⁹ protein-ligand interactions,²⁵⁰ electron transfer state in photosynthesis,²⁵¹ enzyme reactions,²⁵²⁻²⁵³ determination of protein structures from NMR,²⁵⁴ refinement of

protein X-ray crystal structures,²⁵⁵ and calculation of free energy changes from mutations in proteins.^{245-246, 252, 254, 256}

Classical all-atom MD simulations use a Newtonian equation of motion ($F = ma$) (Equation 1.1), with F is the force on the particle, m is its mass, and a is its acceleration) to simulate the movement of each atom in the system. The MD method starts with the initial set of coordinates obtained from X-ray crystal structures, NMR structures or theoretical models (homology models), or a combination of these. Given the atomic model, the interactions for all atoms must be defined. The Newtonian equation is subsequently applied to calculate the force after the systems are minimized to eliminate high energy interactions, such as steric clash, prior to simulation. The integration of the equations of motion after equilibration generates an ensemble of equilibrated states, including coordinates and velocities of the atoms as a function of time. The simulations require three components: initial coordinates (obtained from experimental structures or from models, or some combination of both), a potential (obtained from a force field and the coordinates), and algorithms for propagation.²⁴⁴ The process of how an MD simulations proceeds is set out in Figure 1.8.

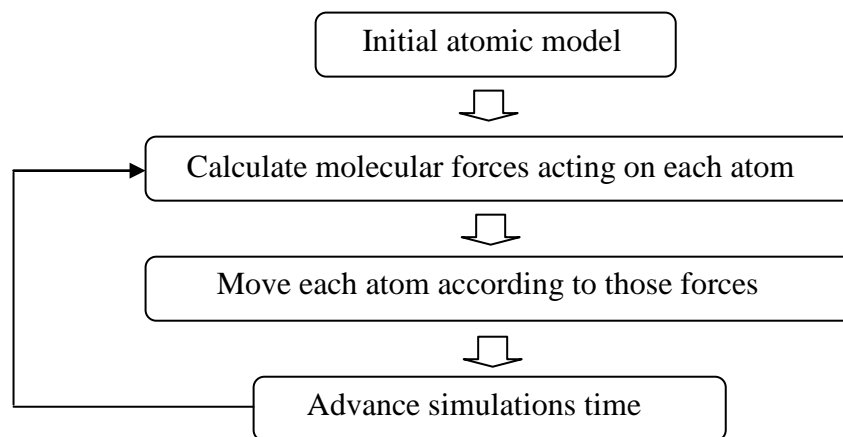


Figure 1.8. How MD simulations proceed (adapted from Durrant D. *et al*²⁴⁶).

1.4.3.1 Force fields

The mathematical functions describing the potential energy of a system and their related parameters are called a “force field”, which is set to describe the interactions between atoms and molecules. A typical molecular CHARMM force field for a molecular system

is estimated as follows in Equation 1.2. In brief, a molecular force field usually include six terms; namely the bond, angles, dihedral, improper dihedral angles, nonbonded, and Urey-Bradley (UB) terms.

$$\begin{aligned}
 U(\{\vec{R}\}) = & \sum_{\text{bonds}} K_b(b - b_0)^2 + \sum_{\text{angles}} K_\theta(\theta - \theta_0)^2 + \sum_{\text{dihedrals}} K_\phi[1 + \cos(n\phi - \delta)] \\
 & + \sum_{\text{impropers}} K_\varphi(\varphi - \varphi_0)^2 + \sum_{\text{UB}} K_{\text{UB}}(r_{1,3} - r_{1,3;0})^2 \\
 & + \sum_{\text{nonbonded}} \epsilon_{ij} \left[\left(\frac{R_{\text{min},ij}}{r_{ij}} \right)^{12} - 2 \left(\frac{R_{\text{min},ij}}{r_{ij}} \right)^6 \right] + \frac{q_i q_j}{4\pi D r_{i,j}}
 \end{aligned}$$

Equation 1.2. Equation used to calculate the atomic forces in MD simulations, from MacKerell A. D. *et al.*,^{202, 257} where K_b , K_θ , K_χ , K_φ and K_{UB} are the bond, angles, dihedrals angles, improper dihedral angles and Urey-Bradley force constants, respectively. b , θ , ϕ , φ , $r_{1,3}$ are the bond length, bond angle, dihedral angle, improper torsion angle, and the Urey-Bradley 1,3-distance, respectively; with the subscript zero means the equilibrium values. The Urey-Bradley is a harmonic term, used for bond-stretching or angle bending in distance between atoms 1 and 3 (the two terminal atoms) in an angle. It is important to more accurately model in vibration spectra.²⁵⁸ n and δ are the values of dihedral multiplicity and phase. The nonbonded includes the van der Waals and electrostatic interactions. The van der Waals interaction is using Lennard-Jones 6-12 potential with $R_{\text{min},ij}$ is the radius in the Lennard-Jones term; q_i and q_j are the partial atomic charge of atom i and j , respectively. ϵ_{ij} is the effective dielectric constant, and r_{ij} is the distance between atoms i and j .

The different force fields have a critical influence on the results of MD simulations.²⁴⁵

The most common currently used force fields are AMBER,²⁵⁹ CHARMM,²⁵⁷ GROMOS²⁶⁰ and OPLS.²⁶¹ The choice of force field usually depends on the preference of the molecular simulation suite.²⁴⁴ Among of the most commonly-used simulation packages are the AMBER,²⁶² CHARMM,²⁶³ GROMACS²⁶⁴ and NAMD.²⁶⁵ These programs share common basic features, but differences lie in their capacities, and underlying philosophies.²⁴⁴

1.4.3.2 Setting up and running MD simulations

To set up and run a MD simulation, some important ingredients should be considered carefully: such as the initial atomic coordinates of the system, the choice of force field, simulation program (integration method), time steps, type of ensemble, boundary conditions, salvation, and time length of the simulation.²⁶⁶⁻²⁶⁷ The length of simulation

time depends on the investigated system and the aim of the study. The integration method will decide parameter of time steps. Any change to these factors will affect the outcomes of the simulation, as well as the requirements of computational time.²⁶⁸ In addition, the type of ensemble must be selected. Traditional MD for biomolecular systems often use NTP (N, number of particles, T, temperature, and P, pressure) or isothermal-isobaric ensemble, which include the constant NTP. Furthermore, in MD simulations, treatment of electrostatic interactions (long range coulombic forces) is very important.²⁶⁷ Nowadays, one can use the Ewald summation method, such as the Particle-Mesh Ewald (PME) for electrostatics.²⁶⁹⁻²⁷⁰

Moreover, boundary conditions and solvent models are required. Due to the influence of interactions at the boundaries of the system on energy calculations, the boundaries must be taken into account. The most common way in simulation of biomolecular systems is using periodic boundary conditions. Periodic boundaries cover the system in a cell (typically a cubic box or a sphere or other geometric shape), and surround it with mirror replica cells of the system. The size of the box must be large enough so that the molecule does not “see” itself.²⁶⁷ The interaction energies can be calculated across the cell boundaries, so the boundary effect is minimized. Two types of solvent models are commonly used, namely implicit (or continuum) and explicit models. The implicit solvation means the solvent is represented with a continuum medium. Two common algorithms used to calculate the solvent electrostatic effects are the Poisson Boltzmann surface area (PBSA) and the Generalized Born surface area (GBSA) model.²⁴⁴ These models are less expensive, but the major problem is that they do not take into account entropic effects, and their difficulty in dealing with heterogeneous environment, which affect many biological processes. On the other hand, the explicit solvation is the model in which the solvent molecules and counterions are treated as explicitly surrounding the biomolecule.²⁴⁴ This method is the most accurate, however it is also time-consuming. Among different water models developed, TIP3P water model is adopted for the CHARMM force field.²⁷¹⁻²⁷²

In general, a MD simulation consists of minimization, equilibration, and production steps. The initial system is minimized to relax potential steric clashes in the structure. The minimized structure is then equilibrated. When equilibration is reached, the values

of average temperature, pressure, and energies are stabilized. The production phase will be the next step and used to calculate the desired properties.

1.4.3.3 Analysis of a MD simulation

The results of MD simulations are analyzed using trajectories to gain insights into the structure; such as protein stability and flexibility. The frequently used analyses are global measures by calculating a RMSD²⁷³ and a root mean square fluctuation (RMSF) or B-factors;²⁷⁴ secondary structure analysis, hydrogen bonds (HBs), and hydrophobic contacts; clustering analysis, quasi-harmonic, and principal component analyses/correlation function, and binding free energy calculation.²⁴⁴

RMSD is usually used to investigate the global stability of the system,²⁷⁵ while RMSF is used to obtain the local structural flexibility and stability. They are calculated as following equations:

$$\text{RMSD} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i^m - x_i^l)^2 + (y_i^m - y_i^l)^2 + (z_i^m - z_i^l)^2}$$

Equation 1.3. The RMSD between atoms of the trajectory frames and the corresponding atoms of the initial structure, where N is the number of atoms, x^m, y^m, z^m are the Cartesian coordinates of the initial structure and x^l, y^l, z^l are the Cartesian coordinates of trajectory at frame t.²⁷⁵

$$\text{RMSF} = \sqrt{\frac{1}{T} \sum_{i=1}^T (x_i - \bar{x})^2}$$

Equation 1.4. The RMSF of an atom, where T is the number of trajectory frames, and \bar{x} is the time-averaged position.²⁷⁵

The resulting trajectories are analyzed by the CHARMM22²⁶³ and VMD (version 1.9.1).²⁷⁶ A simple geometry criterion was used to define a hydrogen bond: the distance between proton donor (D) and acceptor (A) atoms less than 3.5 Å, and the angle D – H ... A greater than 120°. ⁷⁰ If the percentage of HBs occupation is higher than 50%, they are considered as the medium, whereas the strong HBs are determined by HBs occupations of greater than 75%. ⁷⁰ The hydrophobic contacts between the carbon atoms of non-polar parts of residues of proteins were also monitored with a cutoff distance of

4.0 Å.²⁷⁷⁻²⁷⁸ Clustering analysis was conducted using the Clustering Plugin in VMD for all of the snapshots from the trajectories.

Clustering analysis involves the grouping of similar samples of data, joining the ensembles of data into the group to identify the most populated conformations sampled. Structural clustering is a useful method to reduce the sample size for conformational analysis to understand the molecular motion within conformational space.²⁷⁹ Principal-component analysis uses a constructed matrix of atomic fluctuations to find the lowest modes, which represent the most of fluctuations. Quasi-harmonic analysis gives normal modes in the harmonic system, while correlations functions are used to measure the correlation of two fluctuating quantities over the time.²⁴⁴

1.4.3.4 Combining molecular docking and MD simulations

Docking and MD simulations have their own strengths and weaknesses, listed in Table 1.5. Some reviews highlighted the improvement of computational protocols using a combination of MD simulations in docking procedures.²⁸⁰⁻²⁸¹ MD simulations before docking may explore the conformational space of the protein receptor, while using MD simulations after docking helps to optimize the final structures, analyze protein flexibility, and stability of different complexes, account for solvent effects, and obtain accurate energetic properties.¹⁶⁷

Table 1.5. Strengths and weaknesses for docking and MD simulations²⁸⁰⁻²⁸¹

	Docking	MD simulations
Strengths	Fast and inexpensive: docking explore conformation space of ligands in a short time, allowing the scrutiny of large libraries of drug-like compounds at a reasonable cost.	Accurate because it treats both ligand and protein in a flexible way, allowing for an induced fit of the receptor-ligand. The effect of explicit water molecules can be studied directly.
Weaknesses	Lack of or poor flexibility of the protein. Absence of a unique and widely applicable scoring function, necessary to generate a reliable ranking of the final complexes.	Costly and time-consuming. The system can get trapped in local minima.

1.4.4 BINDING FREE ENERGY CALCULATIONS

Given a protein (P) and a ligand (L), the ligand can bind to a protein to form a complex PL. The binding affinity of PL in an equilibrium concentration can be computed simply by using the following equation:

$$K = \frac{[L][P]}{[PL]} \quad (\text{Equation 1.5})$$

This binding affinity can then be related to the free energy of binding (free energy change to changes in enthalpy and entropy) using:

$$\Delta G_{\text{bind}} = -RT \ln K_d = \Delta H - T\Delta S \quad (\text{Equation 1.6})$$

Where ΔG_{bind} is the change in free energy of a binding process, ΔH and ΔS are the corresponding changes in enthalpy and entropy, respectively. R is a gas constant with $R=8.314 \text{ JK}^{-1} \text{ mol}^{-1}$, and T is the temperature of the system in Kelvin degree, and K_d is dissociation constant.

The binding free energies (or binding affinities) are used as a criterion for differentiation of inhibitors from other small molecules (binders), and also selection of strong ligands based on their protein binding strengths. It has become a major interest in structure-based drug design. The computational approaches can be used to estimate the binding free energies together with the experimental assays. However, obtaining accurate values of binding free energies remains a challenge.^{282,283} Upon protein binding, the protein and ligand may be affected by conformational changes influenced by water and ions. Recently, there has been a very large number of approaches developed to solve this problem using different atomistic models.²⁸³⁻²⁹⁸ The approaches can have the simplicity of a scoring function (docking) or the complexity and sophistication of free energy methods.²⁹² Most of them are still under active study, and have different trade-offs between accuracy and computational efficiency.

As mentioned above, docking and scoring use a single bound conformation containing a simplified energy model, such as an empirical force field, with a simple solvent model. Thus, it can provide approximate binding affinities (scores) of ligand and protein. However, the docking results are system-dependent.²⁹² The starting configuration of the protein–ligand complex is also considered.²⁹⁹ Some attempts have been made to

improve docking such as taking account redistribution of ligand charges (potential energy models), solvent models, protein flexibility, and considering changes of configuration entropy.²⁹² In contrast to docking, free energy methods can give more accurate binding energy, but need more computational cost. These approaches generate thermodynamic averages (converged results) using a conformational sampling. They can be classified into two methods: end point and pathway methods. The end point methods require simulation of the bound and free states of the ligands to generate conformations of both states, and compute the binding free energy based on the difference between them. The pathway methods calculate binding free energy using the simulation of many immediate states to sum up all of small changes along a multistep pathway.²⁹² A choice of calculating absolute or relative free energies using implicit or explicit solvent, and the length of the simulation will hugely impact computational time, the accuracy, and efficiency of the calculation. Popular methods include free energy perturbation (FEP), thermodynamic integration (TI), linear interaction energy (LIE), molecular mechanics-PBSA (MM-PBSA), and molecular mechanics-GBSA (MM-GBSA), discussed below.

1.4.4.1 Free energy perturbation (FEP)

The FEP and TI approach are the theoretically rigorous methods, and offer accurate binding free energy; including the absolute binding free energy of a ligand, and relative binding energies between two ligands, X and Y, bound to the same protein (P).³⁰⁰ This approach depends on a thermodynamic cycle (Figure 1.9). The perturbation theory is that different binding free energy between the first state (before binding) and final state (after binding) can be calculated by the formula $\Delta G = -RT\ln\langle e^{-\Delta U/RT} \rangle$ (Equation 1.7) in which ΔU is change of the energy function between two states, and the angle brackets are a Boltzmann average (ensemble average) obtained from MC or MD simulations.³⁰¹

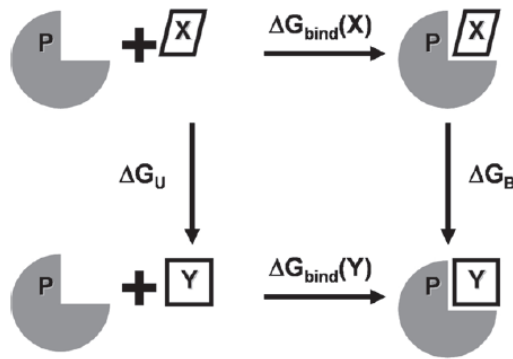


Figure 1.9. Thermodynamic cycle for calculating relative binding free energies between two ligands bound to the same protein (from Michel J. *et al*³⁰⁰).

Based on this theory, the difference of binding free energies of two ligands X and Y bound in the same protein in Figure 1.9 is calculated as follows:

$$\Delta\Delta G_{\text{bind}} = \Delta G_{\text{bind}}(Y) - \Delta G_{\text{bind}}(X) = \Delta G_B - \Delta G_U \text{ (Equation 1.8)}$$

Where $\Delta G_{\text{bind}}(X)$ and $\Delta G_{\text{bind}}(Y)$ are binding free energies for ligand X, Y, respectively; and ΔG_U and ΔG_B are unphysical transmutation free energy from ligand X to ligand Y in the unbound and bound state, respectively.

In the case where ligand X and Y are similar, the values of ΔG_U and ΔG_B are easier to obtain than $\Delta G_{\text{bind}}(X)$ and $\Delta G_{\text{bind}}(Y)$ because the mutation from ligand X to ligand Y is assumed to cause only localized changes. If ligand X and Y are too different, large changes between them may cause sampling problems. So, in the FEP method, the path of transformation is divided gradually into many small steps (intermediate) to allow smooth conversion of ligand X to Y.³⁰⁰ The binding free energy difference between two states, X and Y, is the sum of the contributions from all steps. The formula to calculate the value for each step is:

$$\Delta G = -RT \sum_{i=1}^{N-1} \ln \left\langle \exp \left(-\frac{H(\lambda_{i+1}) - H(\lambda_i)}{RT} \right) \right\rangle_{\lambda_i} \text{ (Equation 1.9)}$$

Where ΔG is the free energy difference between two states, X and Y. λ_i varies from 0 (state X) to 1 (state Y); $H(\lambda_i)$ and $H(\lambda_{i+1})$ represents Hamiltonian of the system at λ_i , λ_{i+1} ; and $\langle \rangle_{\lambda_i}$ indicates an ensemble average. Absolute binding free energies can be obtained from FEP method by setting the interaction potential of the ligand to zero in one of the states. That means transforming the ligand into dummy atoms that do not interact with their surroundings.²⁵⁰

The FEP is currently considered the most powerful and promising approach.²⁹³⁻²⁹⁴ It takes into account entropic contributions to binding affinities arising from solvent effects and protein/ligand flexibility. However, the FEP results could have high precision, but low accuracy.³⁰²⁻³⁰³ Most of the computational time is spent on “perturbations”, meaning uninteresting configurations corresponding from unphysical paths (X to Y), which makes the method difficult for application.²⁹⁵

1.4.4.2 Thermodynamic integration (TI)

Another pathway approach, similar to FEP is TI.³⁰⁰ In the TI method, the difference in free energy between two states, A and B, is calculated based on using multiple intermediate states, defined by a coupling parameter λ .³⁰⁴⁻³⁰⁵ The average of the derivatives of the Hamiltonian at each λ , $H(\lambda)$ is calculated and then TI uses numerical integration over λ to calculate the free energy difference between two states, where λ has the same meaning as in FEP:

$$\Delta G = \int_0^1 \left\langle \frac{\partial H(\lambda)}{\partial \lambda} \right\rangle d\lambda \quad (\text{Equation 1.10})$$

Where ΔG is the Gibbs energy difference between two states, $\langle \rangle$ is an ensemble average obtained at λ . The pathway of intermediates between the states of interest can be parameterized between $\lambda=0$ and $\lambda=1$.²⁸⁴

1.4.4.3 Molecular mechanics-Poisson Boltzmann surface area (MM-PBSA) or Molecular mechanics-Generalised Born surface area (MM-GBSA)

Compared to the pathway methods, the end point methods such as MM-PBSA or MM-GBSA are more computationally efficient, and widely applied for the estimation of the accurate relative binding free energies of related compounds.^{250, 306-309} The methods combine the Poisson Boltzmann (PB) or Generalised Born (GB) electrostatics with molecular mechanics (MM), and solvent accessibility (SA) models, or continuum solvent approaches, to estimate binding energies.^{299, 307-309} An initial MD simulation using continuum solvent approach provides a thermally average ensemble of structures. Several snapshots are then processed, removing all water and counterion molecules, and used to calculate the total binding free energy of the system with the equations:

$$\Delta G_{\text{bind}} = \Delta H - T\Delta S \approx \Delta E_{\text{MM}} + \Delta G_{\text{sol}} - T\Delta S \quad (\text{Equation 1.11})$$

$$\Delta E_{\text{MM}} = \Delta E_{\text{internal}} + \Delta E_{\text{elec}} + \Delta E_{\text{vdw}} \quad (\text{Equation 1.12})$$

$$\Delta G_{\text{solv}} = \Delta G_{\text{PB/GB}} + \Delta G_{\text{SA}} \quad (\text{Equation 1.13})$$

Where ΔE_{MM} is the change of the MM energy in the gas phase, which include $\Delta E_{\text{internal}}$ (corresponding to the bond, angle, and dihedral energies), ΔE_{elec} (electrostatic energy) and ΔE_{vdw} (van der Waals energy); ΔG_{sol} is the solvation free energy which is sum of electrostatic solvation energy, $\Delta G_{\text{PB/GB}}$ (polar contribution, calculated by solvent accessible surface area (SASA)) and the non-electrostatic solvation ΔG_{SA} (nonpolar contribution). The conformational entropy change, $-T\Delta S$, is the most difficult term to evaluate, estimated using quasi-harmonic analysis or normal mode analysis of the trajectory. The entropy change can be assumed to be cancelled if only the relative binding energies of a series of structurally similar compounds is required; however if the absolute energy is important, or if the compounds are notably different, then the contribution to the final free energy cannot be ignored.²⁸⁵

Although the MM-PBSA and MM-GBSA have shown their successful application in biochemical systems,³¹⁰⁻³¹⁴ especially as post-docking methods in virtual screening,³⁰⁸ their performances rely on the system in question. The modification in simulation protocols can affect the approach; such as the sampling strategy of generating snapshots, methods to calculate entropy, and other parameters (charges models, force fields, the solute dielectric constant, and radius parameters in continuum solvent models). In general, the MM-PBSA is more sensitive to the parameters of the systems than MM-GBSA.³⁰⁹

1.4.4.4 Linear interaction energy (LIE)

A new semi-empirical method for calculating binding free energies for ligands from MD simulations has recently been introduced.³¹⁵ This method is based on a linear approximation of polar and non-polar free energy contributions from the MD averages. The idea originated from the problem that with a diversity of compounds with “small perturbations,” it would be very difficult to calculate their binding free energies using FEP. Therefore, the absolute binding free energy of a ligand is calculated as the change

in free energy when the ligand is transferred from aqueous solution (unbound or free state) to its solvated receptor binding site (bound state).³¹⁶ In other words, two simulations are required: one with the ligand free in solution, and one with it bound to solvated receptor. The equation for binding free energy is used as follows:

$$\Delta G_{\text{bind}} = \alpha \Delta \langle V_{l-s}^{\text{vdw}} \rangle + \beta \Delta \langle V_{l-s}^{\text{elec}} \rangle + \gamma \quad (\text{Equation 1.14})$$

Where $\Delta \langle \rangle$ are differences between the averages of the nonbonded van der Waals (vdW) and electrostatic (elec) interactions in the bound or unbound states, collected from MD simulation averages between the ligand and its surrounding environment (l-s). The parameters are the weight coefficients α and β for the non-polar and polar binding energy contributions, respectively; and an additional constant, γ .

Since its initial use, there have been a large number of studies showing LIE as a promising method for computation of binding free energy for protein-ligand.^{285, 295, 315-}

³³⁰ LIE approach is considered to be a good alternative compared to other approaches such as FEP and TI as it estimates the absolute binding free energies; slower than scoring of single conformations, but faster than rigorous FEP approach.³²⁰ When examining the general validity of the electrostatic linear response approximation, the differences in electrostatic response properties between protein and water solvent were investigated by introducing different electrostatic scaling coefficients; β_{wat} and β_{prot} and α_{wat} and α_{prot} . The constant term γ in solvation energy was also discussed and it depends linearly on surface area.^{319, 331-334} So the following LIE equation was:

$$\Delta G_{\text{bind}} = \alpha_{\text{prot}} \langle V_{l-s}^{\text{vdw}} \rangle_{\text{bound}} - \alpha_{\text{wat}} \langle V_{l-s}^{\text{vdw}} \rangle_{\text{unbound}} - \beta_{\text{prot}} \langle V_{l-s}^{\text{elec}} \rangle_{\text{bound}} - \beta_{\text{wat}} \langle V_{l-s}^{\text{elec}} \rangle_{\text{unbound}} + \gamma \quad (\text{Equation 1.15})$$

The general equation with the values $\alpha_{\text{prot}} = \alpha_{\text{wat}}$, $\beta_{\text{prot}} = \beta_{\text{wat}}$ and $\gamma = 0$ gave a significant improvement compared to the original one. Aqvist and Hansson reported the relationships between electrostatic free energies and the solvation energetic for several model compounds in different solvents.³¹⁶ Some deviations from linear response were found, in particular, for neutral dipolar solutes and for uncharged ligands having certain dipolar groups in the case of water solvent through effect of hydrogen bonding network. The optimal value of $\beta = 0.5$ was suggested to be reconsidered, for instance $\beta < 0.5$ to gain more accurate predictions. The α is an empirical constant which can be fitted to experimental binding free energies. In addition, the coefficients α and β converge to the

same values in both bound and unbound states. These findings supported the use of the basic LIE equation 1.14.

A refined LIE model, the FEP-derived model showed that the value of β varies depending on the number of hydroxyl groups. The more hydroxyl groups the compound has, the lower the value of β (Table 1.6).³¹⁶

Table 1.6. Values of the β parameter as a function of the chemical nature of the ligand, according to Hansson *et al.*³²⁰

Chemical nature	β
Charged compounds	0.5
Neutral compounds	0.43
Neutral compounds bearing a single hydroxyl group	0.37
Neutral compounds bearing two or more hydroxyl groups	0.33

So, the final binding energy is calculated as:

$$\Delta G_{\text{bind}} = \alpha \langle V_{\text{bound}}^{\text{vdW}} - V_{\text{unbound}}^{\text{vdW}} \rangle + \beta \langle V_{\text{bound}}^{\text{elec}} - V_{\text{unbound}}^{\text{elec}} \rangle + \gamma \quad (\text{Equation 1.16})$$

Where $\langle V_{\text{bound}}^{\text{elec}} - V_{\text{unbound}}^{\text{elec}} \rangle$ represents the averages change in electrostatic energy between the bound and unbound (free or unbound or just solvent) states, and $\langle V_{\text{bound}}^{\text{vdW}} - V_{\text{unbound}}^{\text{vdW}} \rangle$ the average change in vdW from an aqueous solution to a protein environment. α , β and γ are empirically determined constants. The α , β is for the non-polar, polar contributions, respectively; and are the same values in the bound and unbound state.³¹⁶ Applying a value of $\alpha = 0.18$ has shown to successfully reproduce the experimental binding free energies in a wide variety of ligand-protein systems.³²⁸

β_{FEP} values can be assigned to each chemical group present in the ligand, as shown in Equation 1.17 and the values are provided in Table 1.7. The weighting factors depend on salvation energies of each chemical group. The value of w_i is 1.0 for all neutral groups or 11.0 for the anions and cations.³²⁶ The advantage of this approach is that the β coefficient is flexible and provides higher accuracy, since deviations from the linear response due to chemical groups such as amides, amines, or carboxylic acids is explicitly taken into account.

$$\beta = \beta_0 + \frac{\sum_i w_i \Delta\beta_i}{\sum_i w_i} \quad (\text{Equation 1.17})$$

Table 1.7. Values for the β parameter in Equation 1.17, according to Almlöf M. *et al.*³²⁶

Parameter	Value	Chemical nature
$\Delta\beta_i$	-0.06	Alcohols
$\Delta\beta_i$	-0.04	1°, 2° -Amines
$\Delta\beta_i$	-0.02	1° Amides
$\Delta\beta_i$	-0.03	Carboxylic acid
$\Delta\beta_i$	+0.02	Anions
$\Delta\beta_i$	+0.09	Cations

In addition, a correlation of β and the hydrophobicity of the binding site using a weighted desolvation non-polar ratio (WNDR) has been investigated.³²¹ The result suggested that the β is predictable by calculating the WNDR; in particular for systems in which different ligands bind to different binding sites of the same protein. The parameter γ is influenced by the hydrophobicity of the binding site.³²⁴

1.5 PROJECT BACKGROUND AND AIMS

In recent years, there has been an emergence or re-emergence of some alphaviruses in various countries; in particular the CHIKV. This presents a worldwide threat to human health, and creates an economic burden for the affected countries. However, there are currently no vaccines or effective drugs available for the treatment of CHIKV virus. In addition, there has been little research to find anti-CHIKV compounds. Therefore, a significant need for research into medicines to combat the virus exists. During the past decade, computational approaches have become an increasingly powerful tool in drug discovery and development. It has not only helped scientists succeed in developing many therapeutic compounds for specific diseases, but has also aided the development of time-saving and cost-effective procedures.

With this in mind, we aim to discover and develop an approach leading to the identification of a number of lead compounds to combat CHIKV disease. This study will primarily utilise computational techniques in all stages of the process. We will use all available information of CHIKV, together with a combination of computational tools, to maximize the efficiency of this study. The study has three principal objectives:

1. To identify potential inhibitors for CHIKV using a structure-based approach with molecular docking and virtual screening.
2. To investigate the stability and flexibility of protein-hit compounds complexes with molecular dynamics simulations.
3. To obtain accurate binding free energies from molecular dynamics simulations and provide guidance in rational optimization of hit compounds.

CHAPTER 2. DISCOVERY OF INHIBITORS TARGETING CHIKV NON-STRUCTURAL PROTEIN 3

2.1 INTRODUCTION

2.1.1 FUNCTION AND ROLE OF THE NSP3 OF CHIKV

The nsP3 protein is considered an attractive target for CHIKV drug discovery⁷² due to its participation in the early stages of the transcription processes of viral replication, though the specific functions of the nsP3 protein remain elusive.^{8, 69, 84, 87} The nsP3 is the third non-structural protein in the CHIKV genome. It consists of two domains, the N-domain and the C-domain;^{70, 73} the N-domain is highly conserved, but the C-domain is not.⁷³ The C-domain is phosphorylated on serine and threonine residues on up to 16 positions.^{73, 85, 335-337} The role of this phosphorylation is still unclear, but deletion of the residues involved in the phosphorylated process has been shown to decrease the level of RNA synthesis.^{8, 73, 336} The N-domain contains the X-domain or a macrodomain, the region comprising the first 160 residues; is commonly present in eukaryotic organisms, bacteria, archaea; and also many positive-strand RNA viruses such as hepatitis E, rubella, coronavirus, and alphaviruses.⁷⁰ The alphavirus macrodomain has a highly positively charged patch on the surface, at the crevice of ADP-ribose 1"-phosphate active site and its periphery.⁷³ The other side of the protein, far from the active site, possesses a negative charge. Thus, the nsP3 macrodomain is considered to complex with ADP-ribose derivatives and RNA. It is believed to control the metabolism of ADP-ribose 1"-phosphate and/or other ADP-ribose derivatives with regulatory functions in the cell.^{8, 338}

In addition, studies based on the SINV reported that the nsP3 phosphoprotein is an essential component of the viral replication and transcription process.⁸⁷ Functional analysis of the effects of mutations of nsP3 on RNA synthesis demonstrated that alterations may cause a loss of capacity for minus strand synthesis, or a failure to increase plus strand synthesis. A change of Ala68 to Gly leading to a modification of the His-Ala-Val peptide was predicted to form part of the active site of the conserved

nsP3 macrodomain.⁸⁷ However, no effects on the ADP-ribose binding was found.⁷³ In addition, the mutation of amino acids at the position Asn10 and Ala24 in ADP-ribose binding of nsP3 macrodomain in SINV affected the replication and viral RNA synthesis, without affecting the binding region.⁸⁶⁻⁸⁷

Recent findings revealed that the CHIKV nsP3 was described to have a novel function as a regulator of the cellular stress response.³³⁹ Studies of SINV-infected cells indicated the importance of nsP3 in the interactions with alphavirus-host.³⁴⁰ Functional analysis of SFV at the C-terminal region of nsP3 showed that the mutations, where 10 residues at the C-terminal are lacking, suppresses the establishment of infection; while lacking the 30 C-terminal residues led to reduced synthesis of subgenomic RNA.³⁴¹ The nsP3 macrodomain has been shown to be responsible for SINV and SFV replication in neurons and neurovirulence in mice.⁸⁵⁻⁸⁶

2.1.2 EARLY ATTEMPTS TO DISCOVER CHIKV NSP3 INHIBITORS

The crystal structure of the nsP3 macrodomain of CHIKV was determined in 2010.⁷³ This structure includes four subunits; and the asymmetric unit consists of six-stranded β -sheets and four α -helices (Figure 2.1). The core β -sheet and positions of the α -helices have proven to be highly conserved. Also present in the nsP3 structure is the ligand ADP-ribose.^{70, 73} The active site is in the crevice, between the top of the β -strands 2, 4 and 5 and is surrounded by two loops between β 2- α 1 and β 5- α 3.

Currently, there is only one publication regarding molecular modelling for the CHIKV nsP3 macrodomain based on the crystal structures (PDB id: 3GPG) and its complex with ADP-ribose (PDB id: 3GPO).⁷⁰ The study focused on an understanding of the specific binding of the ADP-ribose to the nsP3 macrodomain of CHIKV, while also comparing with VEEV. The results of MD simulations of the structure with ADP-ribose identified the binding modes and the key residues for interactions between ADP-ribose and the nsP3.⁷⁰ The negatively charged PO_4^{2-} component of ADP-ribose showed the strongest interaction with the protein, and the binding free energies estimated from MD simulations were in good agreement with previous experimental data.

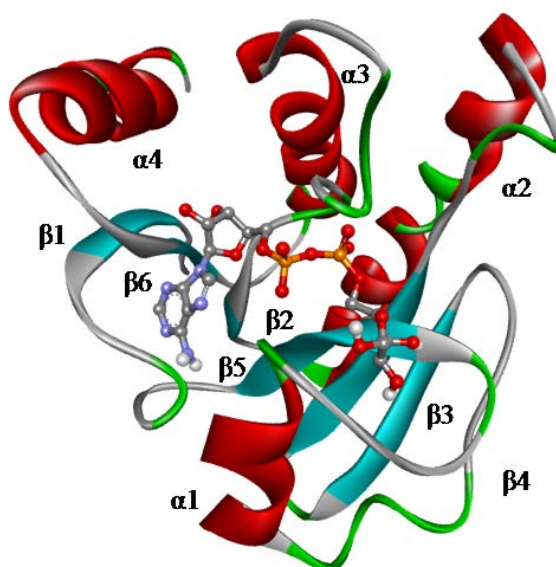


Figure 2.1. X-ray crystal structure of the macrodomain of CHIKV in complex with the ADP-ribose.⁷³

2.2 OVERVIEW OF THIS STUDY

There is very little information available on the nsP3 macrodomain and its inhibitor. Therefore, this study focused on using a combination of computational approaches, including molecular docking, virtual screening, MD simulations, and binding free energy calculations to discover potential lead compounds that inhibit the nsP3 in CHIKV and propose compounds for biological testing (Figure 2.2).

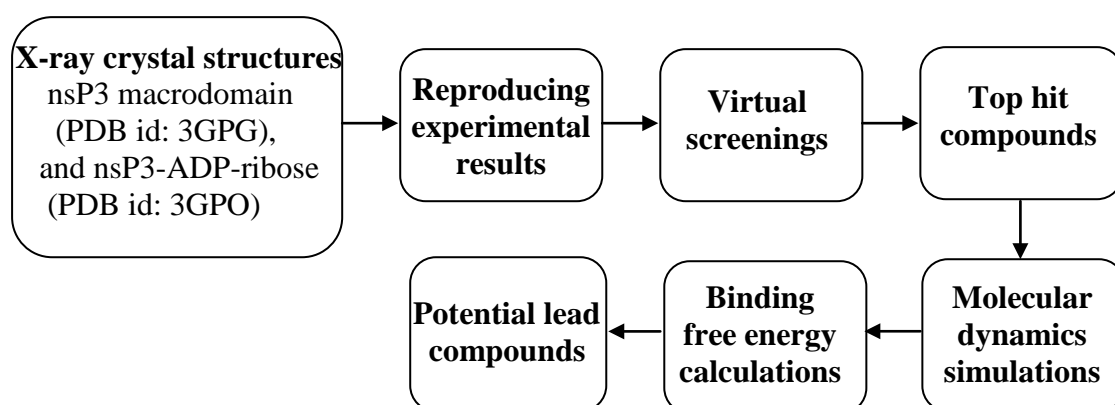


Figure 2.2. Schematic diagram of *in silico* approaches.

2.2.1 MOLECULAR DOCKING AND VIRTUAL SCREENING

A docking protocol was established (details described in the Experimental procedures and methods section), which included the following stages:

The nsP3 macrodomain protein was downloaded from the Protein Data Bank (PDB id: 3GPG) and this was used as the receptor for docking. The protein structure was submitted to WhatIF website to correct manually chirality errors, and the conformation of sidechains of His, Asn, and Gln.

The next step was to carry out energy minimization to relax the structure and remove steric overlaps using the CHARMM22 force field in the Accelrys Discovery Studio (DS) 2.0 software package.³⁴² The steepest descent algorithm with 3,000 steps was applied. With a C_{α} RMSD of 0.59 Å between the minimized structure and the X-ray structure, the minimized structure was utilized for the subsequent docking process. Polar hydrogen atoms were added with AutoDock Tools (version 1.5.4). The ligand ADP-ribose was extracted from the complex crystal structure (PDB id: 3GPO) and prepared by AutoDock Tools for docking in Vina. The other ligands employed for virtual screening were taken from National Cancer Institute (NCI) Diversity Set II.

Parameters of the grid box size and centre of the box were defined, and other parameters including a search space or exhaustiveness (E) and number of binding modes to generate (num_modes), were selected using AutoDock Tools. For docking, the grid box needs to be large enough to accommodate the ligands. Initially, in this case, the grid box size and its location of the binding site were defined based on the place where the ligand ADP-ribose was bound in the X-ray structure. During virtual screening, the location and the size of the grid box were carefully investigated via blind docking (in which the box is sufficiently large to cover the whole protein), and focused docking (in which a smaller box was centred on potential binding sites of interest). Blind docking can reveal potential binding sites in the nsP3. The parameters for docking were chosen such as E and the maximum num_modes set to the default values with $E = 8$ and $\text{num_modes} = 9$.

The docking protocol was validated by re-docking ADP-ribose to compare with the crystal structure of the nsP3-ADP-ribose complex. The re-docking results were evaluated in terms of RMSD value and the binding affinity. After evaluation, the established docking protocol was used for different virtual screenings. In an effort to identify potential inhibitors (hit compounds), 1541 compounds from the NCI Diversity Set II were screened by docking against the nsP3 macrodomain of CHIKV. Based on the starting point of docking with ADP-ribose at ADP-ribose binding site, virtual screening (VST) was carried out with three different setups. The first was a focused docking centred on the ADP-ribose binding site (Pocket 1: VST1 and VST2). The second setup was a blind docking centred either at the middle of the Pocket 1 (VST3) or the protein (VST4) with the box large enough to cover the whole protein. The third setup was a focused docking centred at the predicted binding sites by MetaPocket,³⁴³ (Pocket 2: VST5 and Pocket 3: VST6). The focused dockings at Pocket 2 (VST5) and Pocket 3 (VST6) were also carried out. Interactions of ligands and protein were analyzed from docking results using Accelrys DS 3.5. The Lipinski's rule was also used to give a general drug-likeness information for hits.²⁰⁶

2.2.2 MD SIMULATIONS

The program NAMD²⁶⁵ was used for MD simulations to investigate the stability and flexibility of the hit-target complexes, and study their interactions. The apo protein and the complexes of protein-hit compounds were prepared to run simulations (details of the procedure are given in the Experimental procedures and methods Chapter). The protein atoms were represented with the CHARMM22 force field,²⁵⁷ and the corresponding parameters for the ligands were generated with AmberTools.²⁶² The systems were solvated under periodic boundary condition with explicit solvent model TIP3P and 0.15 M NaCl. The Langevin algorithm was used to maintain the temperature at 298.15 K and pressure at 1 atm. The PME algorithm was used to compute long range electrostatic interactions.²⁶⁹ The cutoff distance for vdW interactions were set at 12 Å and the pair-list distance was 13.5 Å. The minimization process was applied first and followed by equilibrium simulations with weak harmonic restraints on the heavy atoms for 3 ns. The production runs were continued for 50 ns.

The trajectories for analysis were saved every 10 ps. To determine the system stability, the RMSDs of the heavy atoms over 50 ns was calculated with respect to the starting structure versus the simulation time. The RMSF of C α atoms during the simulations was measured to obtain information on local flexibility of the system. The resulting trajectories such as HBs and hydrophobic contact interactions were analyzed by the CHARMM22²⁶³ and VMD (version 1.9.1),²⁷⁶ (details in Experimental procedures and methods Chapter, section 5.2). The Clustering Plugin Tool in VMD is used for clustering analyses.

2.2.3 BINDING FREE ENERGY CALCULATIONS

Having obtained the simulation results of protein and its complexes, the simulations of ligands in solvent (water) were run to apply LIE approach to estimate the absolute binding free energies for ligands in complexes with the protein targets. Preparation of simulations for ligands are described in Chapter Experimental procedures and methods. NAMDEnergy plugin in VMD was utilized to compute the energy components over the frames obtained from the MD simulations.^{265, 276} The α , β , and γ in the LIE equation need to be defined based on the properties such as hydrophobicity of the ligand and binding site to estimate absolute free energies of binding.

2.3 RESULTS AND DISCUSSION

2.3.1 DOCKING RESULTS WITH ADP-RIBOSE

The docking outcomes of ADP-ribose into the nsP3 protein were evaluated and compared with the available co-crystal structure (PDB id: 3GPO). The different ligand conformations were ranked based on their predicted binding affinities with the default scoring function in Vina (Table 2.1). The RMSD value was calculated between the docked structure and the initial structure (Table 2.1). The best docked pose had a binding affinity of -10.2 kcal/mol.

Table 2.1. Poses in the docking of ADP-ribose into the nsP3. RMSD refers to the heavy-atom RMSD from the co-crystal structure for ADP-ribose with the nsP3.

Poses	Binding affinity (kcal/mol)	RMSD (Å)
1	-10.2	0.6
2	-8.9	8.6
3	-8.9	1.8
4	-8.9	5.5
5	-8.6	8.4
6	-8.5	4.4
7	-8.4	5.7
8	-8.3	10.0
9	-7.8	4.5

Details of analysis of interactions between the complex of the best docked of ADP-ribose and the nsP3 are shown in Table 2.2 and Figure 2.3. The residues in the active site of nsP3, namely Ile11, Ala23, Asn24, Asp31, Val33, Leu108, Gly112, Val113, Tyr114, Tyr142, and Arg144, formed HBs with ADP-ribose. Most of the hydrogen bond donors arose from the protein residues, with corresponding acceptors contained in the ADP-ribose. The only exception is that the ribose component of ADP-ribose can be the donor in interactions with Tyr142. In addition, the diphosphate component of ADP-ribose showed the strongest interaction compared to the rest of this ligand with the greatest number of HBs.

Table 2.2. Analysis of interactions between the best docked of ADP-ribose and the nsP3 macrodomain.

Part of ADP-ribose	Number of HBs	Interactions (Å)
Adenine	2	Arg144(HH21)-N1=2.5 Ile11(HN)-N1=2.1
Ribose	2	H9-Tyr142(O)=2.4 Leu108(HN)-O3'=2.4
Diphosphate	5	Val33(HN)-O1A=2.4 Val113(HN)-O2A=2.4 Val113(HN)-O2B=2.4 Gly112(HN)-O2B=2.1 Tyr114(HN)-O2B=2.1
Ter-ribose	3	Asp31(HN)-O1D=1.9 H9-Ala23(O)=2.4 Asn24(HD21)-O3D=2.1

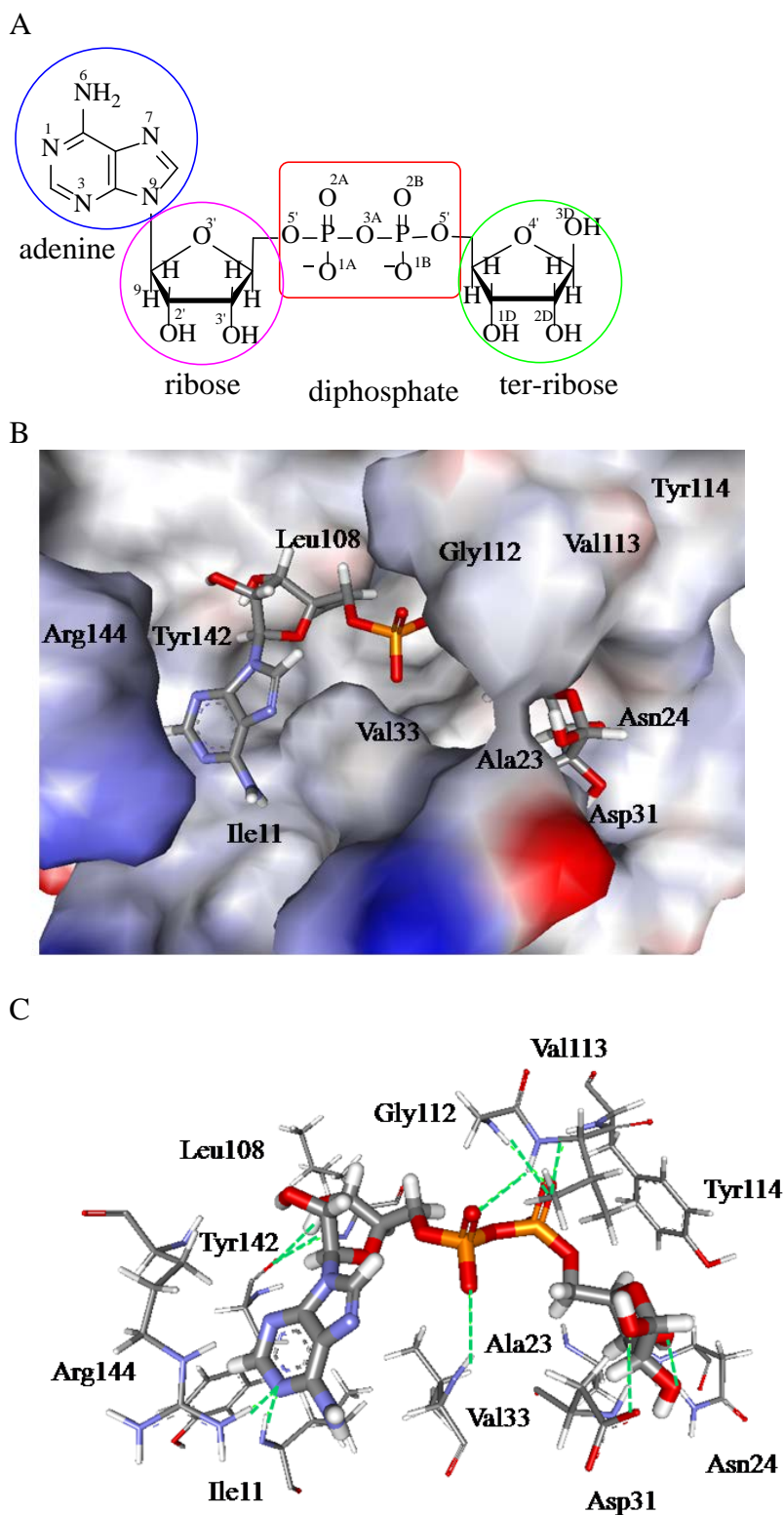


Figure 2.3. Re-docking ADP-ribose (A) into the active site of the nsP3: (B) The best docking pose of ligand ADP-ribose is represented as a stick model (coloured by atom type) while the protein nsP3 is shown in the solvent surface (coloured by interpolated charge with a probe radius of 1.4 Å). (C) The interactions of this pose and the nsP3 residues show hydrogen bonding interactions at the binding site of the nsP3.

Moreover, the accuracy and reliability of docking were evaluated by superimposing both the docked structure and the X-ray structure. The heavy atom RMSD was 0.6 Å, smaller than the 2.0 Å (often used as a criterion for the correct bound structure prediction)³⁴⁴ (Table 2.1 and Figure 2.4), indicating that molecular docking reproduced the binding mode in the co-crystal structure. A comparison of interactions of docking results with published data (Table 2.3) confirms that there was a good agreement with the key interactions, and showed the current docking protocol was able to reproduce the correct pose. The differences in important residues from forming HBs and hydrophobic contacts were acceptable due to different methods used.

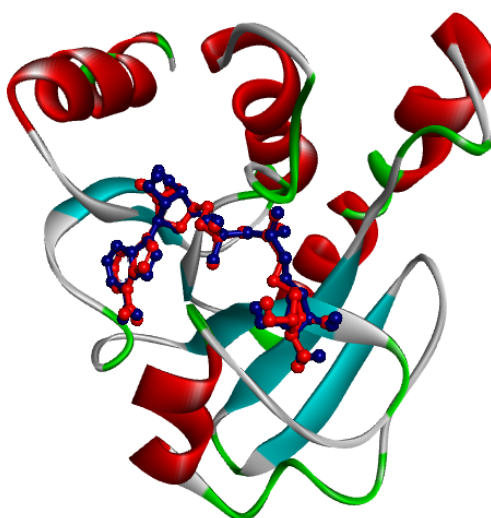


Figure 2.4. Superimposition of the ADP-ribose after docking (in red, the top pose) and its structure in the co-crystal structure (in blue) at the active site of nsP3. The heavy-atom RMSD between the two structures is 0.6 Å.

Table 2.3. Comparison of the identified hydrogen bonding interactions in the nsP3-ADP-ribose docked complex with the previously published data. In Ref [73], key residues including bonding residues (in bold), were identified by experimental work with the crystal structure of complex nsP3-ADP-ribose (3GPO) while residues in Ref [70] were determined by MD simulations of ADP-ribose in the nsP3 based on the above crystal structure.

	Current work	Ref ⁷³	Ref ⁷⁰
HBs	11	11	11
Interacting residues	Ile11 , Ala23, Asn24 , Asp31 , Val33, Leu108, Gly112 , Val113 , Tyr114 , Tyr142, Arg144	Asp10, Ile11 , Asn24 , Asp31 , Thr111, Gly112 , Val113 , Tyr114 , Tyr142, Arg144	Asp10, Ile11 , Asn24 , Asp31 , Val33, Ser110, Thr111 Gly112 , Val113 , Tyr114 , Arg144

2.3.2 IDENTIFICATION OF INHIBITORS FOR THE NSP3 MACRODOMAIN

Results of virtual screenings based on blind docking and focused docking are listed in the Table 2.4. The top ten compounds for each VST and their binding affinities were selected (structures appended in Appendix 1).

In addition to the proposed ADP-ribose binding site (the active site, Pocket 1), two additional binding sites (Pocket 2 and 3) were identified based on blind docking. The residues making up each pocket are listed in Table 2.5.

Pocket 1 and Pocket 3 share a number of interacting residues including Asn24, Asp31, Val33, Gly112, Val113, and Tyr114. Pocket 2 was found on the opposite side and behind Pocket 1. The locations of pockets in the nsP3 and the locations of top hit ligands in the three pockets are illustrated in the Figure 2.5.

Table 2.4. Results of the top ten compounds of different virtual screens for the nsP3. The binding affinities are shown in kcal/mol.

VST1 ^a	VST2 ^b	VST3 ^c	VST4 ^d	VST5 ^e	VST6 ^f
1. NCI_25457 (-10.8)	1. NCI_34567_a (-10.9)	1. NCI_61610 (-11.1)	1. NCI_61610 (-11.1)	1. NCI_127133 (-8.3)	1. NCI_670283 (-10.6)
2. NCI_116702 (-10.7)	2. NCI_37553 (-10.9)	2. NCI_293778 (-11.0)	2. NCI_37553 (-11.0)	2. NCI_293778 (-8.2)	2. NCI_319990 (-10.2)
3. NCI_309892 (-10.3)	3. NCI_25457 (-10.8)	3. NCI_345647_a (-10.9)	3. NCI_345647_a (-10.9)	3. NCI_338042 (-7.6)	3. NCI_80731 (-10.1)
4. NCI_109451 (-10.2)	4. NCI_116702 (-10.7)	4. NCI_25457 (-10.8)	4. NCI_25457 (-10.8)	4. NCI_132232 (-7.5)	4. NCI_84100_b (-10.1)
5. NCI_127133 (-10.2)	5. NCI_58052 (-10.6)	5. NCI_58052 (-10.6)	5. NCI_293778 (-10.8)	5. NCI_310326 (-7.4)	5. NCI_372287_a (-10.0)
6. NCI_328101 (-10.2)	6. NCI_127133 (-10.5)	6. NCI_127133 (-10.5)	6. NCI_127133 (-10.7)	6. NCI_328101 (-7.4)	6. NCI_84100_a (-9.9)
7. NCI_372275_a (-10.2)	7. NCI_293778 (-10.5)	7. NCI_372499_b (-10.3)	7. NCI_116702 (-10.6)	7. NCI_69359_a (-7.3)	7. NCI_97920 (-9.6)
8. NCI_45545 (-10.2)	8. NCI_670283 (-10.5)	8. NCI_37553 (-10.3)	8. NCI_58052 (-10.6)	8. NCI_90737 (-7.3)	8. NCI_58502 (-9.5)
9. NCI_84100_b (-10.2)	9. NCI_328101 (-10.4)	9. NCI_309892 (-10.2)	9. NCI_670283 (-10.5)	9. NCI_122819_a (-7.2)	9. NCI_227186_a (-9.4)
10. NCI_37168 (-10.1)	10. NCI_372499_b (-10.3)	10. NCI_37168 (-10.2)	10. NCI_324623 (-10.3)	10. NCI_400976 (-7.2)	10. NCI_293778 (-9.4)

(a) In VST1, the grid box is fixed at the centre of Pocket 1 (9.7 Å, 43.0 Å, -13.2 Å) with a dimension of 16 Å × 16 Å × 16 Å. **(b)** In VST2, the grid box is fixed at the centre of Pocket 1 (9.7 Å, 43 Å, -13.2 Å) with a dimension of 20 Å × 20 Å × 20 Å. **(c)** In VST3, the grid box is fixed at the centre of Pocket 1 (9.7 Å, 43 Å, -13.2 Å) with a dimension of 50 Å × 50 Å × 56 Å. **(d)** In VST4, the grid box is fixed at the centre of the protein (7.7 Å, 45.3 Å, -5.3 Å) with a dimension of 50 Å × 50 Å × 56 Å. **(e)** In VST5, the grid box is fixed at the centre of Pocket 2 (7.7 Å, 45.4 Å, 11.5 Å) with a dimension of 20 Å × 20 Å × 20 Å. **(f)** In VST6, the grid box is fixed at the centre of Pocket 3 (2.3 Å, 44.6 Å, -18.3 Å) with a dimension of 20 Å × 20 Å × 20 Å.

Table 2.5. Pocket residues in the nsP3 macrodomain.

Pocket	Pocket residues
Pocket 1	Asp10, Ile11, Ala22, Ala23, Asn24, Gly30, Asp31, Gly32, Val33, Cys34, Gly70, Pro107, Leu108, Leu109, Ser110, Thr111, Gly112, Val113, Tyr114, Tyr142, Cys143, Arg144, Asp145, Trp148
Pocket 2	His-1, His0, Ala1, Pro2, Ser3, Tyr4, Phe129, Met132, Asp133, Ser134, Thr135, Asp136, Ala137, Asp138, Val139, Ile156, Gln157, Arg159, Thr160
Pocket 3	Ala22, Ala23, Asn24, Pro25, Arg26, Leu28, Pro29, Gly30, Asp31, Gly32, Val33, Cys34, Pro51, Val52, Gly70, Pro71, Asn72, Tyr76, Leu108, Ser110, Thr111, Gly112, Val113, Tyr114

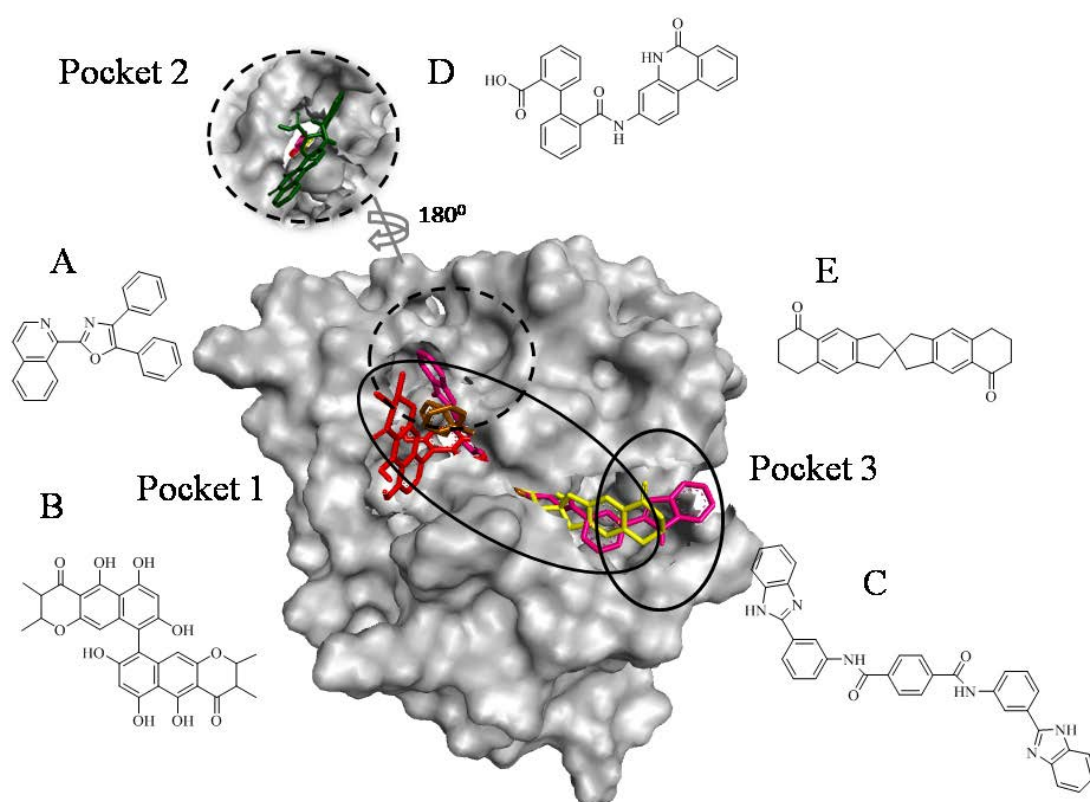


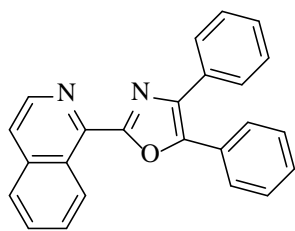
Figure 2.5. Representation of three binding pockets identified in the nsP3 with top hit compounds binding in the pockets. Pocket 1 is the ADP-ribose binding site with ligand NCI_25457 (A, in burgundy), NCI_345647_a (B, in red), and NCI_61610 (C, in pink). Pocket 3 shares some residues with Pocket 1 with ligand NCI_670283 (E, in yellow). Pocket 2 is in the other side of Pocket 1 with ligand NCI_127133 (D, in dark green).

In the focused dockings targeting Pocket 1 (the ADP-ribose binding site, VST1 and VST2), the top hits were NCI_25457 (-10.8 kcal/mol) and NCI_345647_a (-10.9 kcal/mol) (Figure 2.6). Among the top ten hits, four are shared between VST1 and VST2, which differ in the size of the grid box used. It is worth noting that the change in

size of the grid box affected the searching process in Vina, with some new hits being identified. However, the binding affinity values between them are not significantly different from each other. To find more inhibitors, and other potential binding sites in the structure of nsP3, blind docking was used to dock into the entire protein with the grid box centred either at the middle of the ADP-ribose binding site (VST3) or the protein (VST4). For the blind docking (VST3 and VST4), six of the top ten hits are common to VST3 and VST4, and their binding affinities were reproduced within 1.0 kcal/mol. This indicated that the blind dockings are likely to have converged. The results show that most of the top ten ligands fitted well in Pocket 1, and that this pocket can accommodate ligands of different size. However, ligands with bulky structures, such as NCI_293778, NCI_58052 and NCI_61610 (in both VST3 and VST4, Figure 2.6), protruded from the binding site. Therefore, the other pockets surrounding Pocket 1 may serve as alternative binding sites for potential inhibitors.

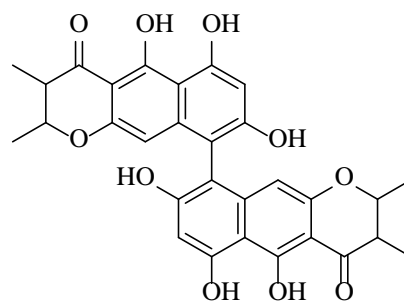
Screenings in VST5 and VST6 produced hits already identified from previous screens, along with some new hits (Table 2.4). For virtual screening, changes to the size of the grid box and its location affected the searching process in Vina. An increase in the dimension of the box is likely to be suitable for larger molecules. For instance, in the blind docking (VST3 and VST4), NCI_61610 (-11.1 kcal/mol) was identified as a top hit, but it did not belong to the top 10 hits in VST1 and VST2. It was also important to note that the majority of hit compounds have a tighter binding in Pocket 1 compared to those in Pocket 2. Most compounds effectively occupied Pocket 2 and Pocket 3 with the significant interactions. Ligands NCI_127133 (-8.3 kcal/mol) and NCI_670283 (-10.6 kcal/mol) (Figure 2.6) bind in Pockets 1, 2, and 3, though in different conformations. Interestingly, the ligand NCI_293778 appeared able to bind in all three pockets, and it may infer that Pocket 1 was more favourable for binding, given the binding affinity of -10.5 kcal/mol compared to -9.4 kcal/mol (Pocket 3) and -8.3 kcal/mol (Pocket 2).

NCI_25457



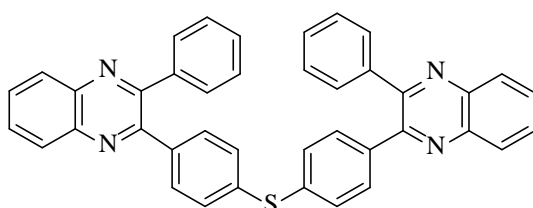
(15)

NCI_345647_a



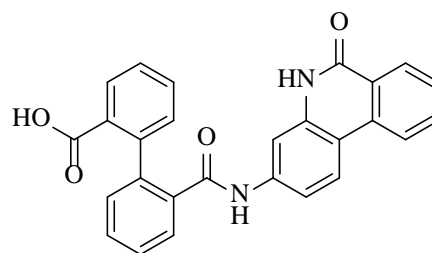
(16)

NCI_293778



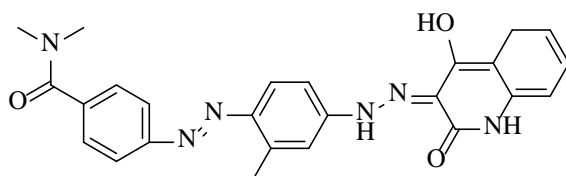
(17)

NCI_127133



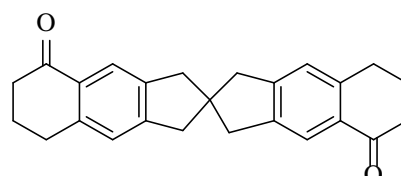
(20)

NCI_58052



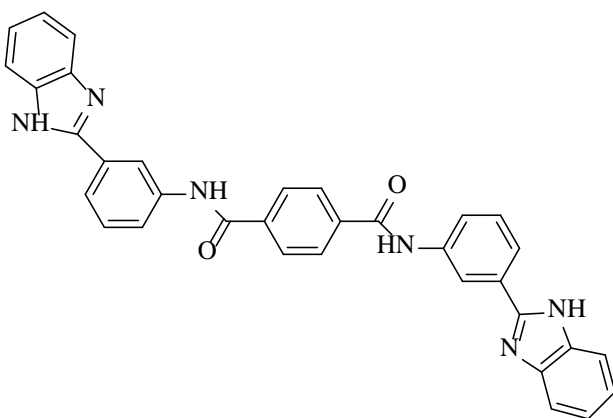
(18)

NCI_670283



(21)

NCI_61610

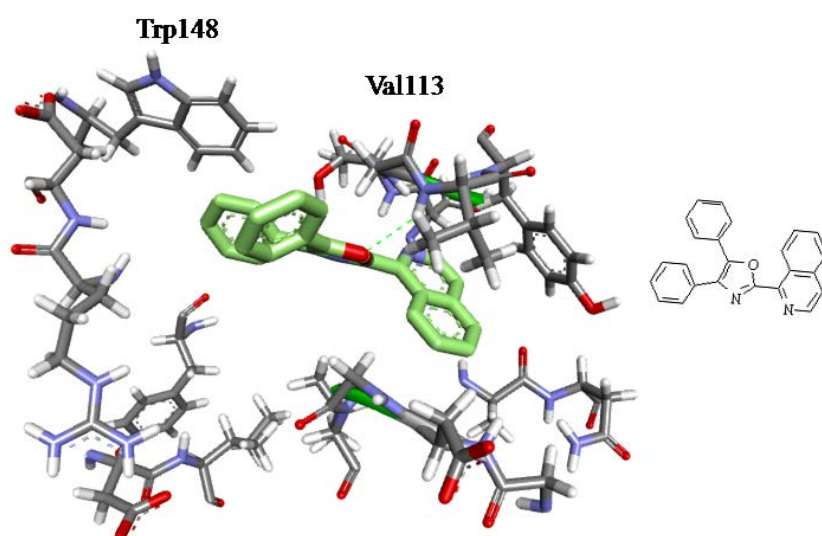


(19)

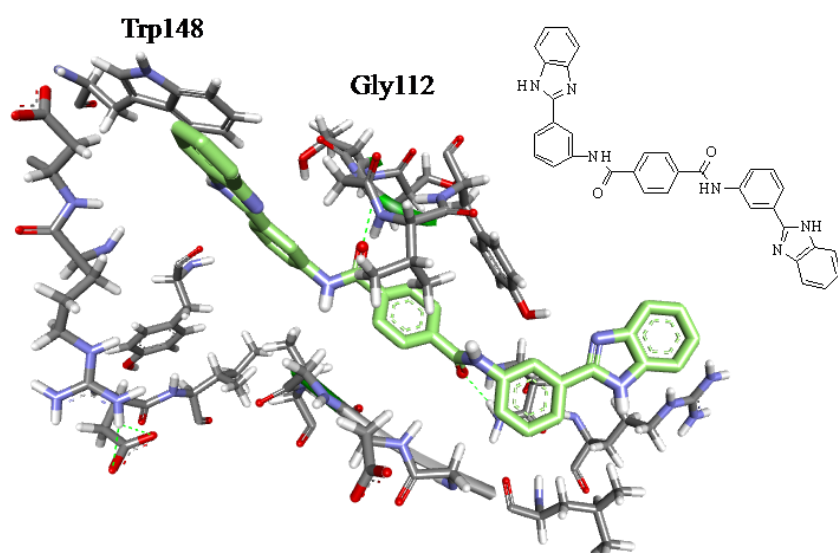
Figure 2.6. Structures of the top hit compounds, obtained from screenings for the nsP3.

Detailed analyses of the interactions between the ligands and protein target were carried out in regard to HBs and hydrophobic contacts (Appendixes 2-7).²²⁴ The results showed that the hydrogen bonding interactions played a more important role in the binding to Pocket 1 and Pocket 3, whereas hydrophobic contacts were responsible for most interactions in the binding to Pocket 2. The contribution of aromatic rings by π -stacking or π -network was emphasized in the enzyme. Most of the ligands can fit very well in the Pocket 1 by forming hydrogen bonds with backbone nitrogen of Val113 or Thr111, and/or interacting through π -stacking or π -network with aromatic ring of Tyr114 or Trp148. For example, top-hit ligands NCI_25457 and NCI_61610 could interact with protein through HBs with residues Val113 and Gly112, respectively; and also π -stacking with Trp148 (Figure 2.7). In addition, the residues located in the region 110-114 play a crucial role in ligand binding to Pocket 1. Mostly, ligands served as hydrogen bonding acceptors while residues of protein were donors. Only interactions between ligands and Thr111, in some cases, this residue can change its role to be an acceptor. In agreement with previous reports,⁷⁰ residues Ser110, Thr111, Gly112, and Tyr114 define Pocket 1 and were key residues in forming interactions with ligands. In addition, we found these residues concurrently define Pocket 3. Among them, Tyr114 formed HBs with ligands, or interacted through the π -network on the aromatic ring with most of ligands in both Pockets 1 and 3. Residues Asp31 and Asn75 contributed in forming hydrogen bonds for ligands in Pocket 3. For ligands bound in the Pocket 2, we found residues Tyr4 or Met132, Asp133, and Thr135 are key residues in forming hydrogen bonding interactions between ligands and protein. Additionally, ligand interactions were observed by π -stacking from the aromatic ring of Tyr4.

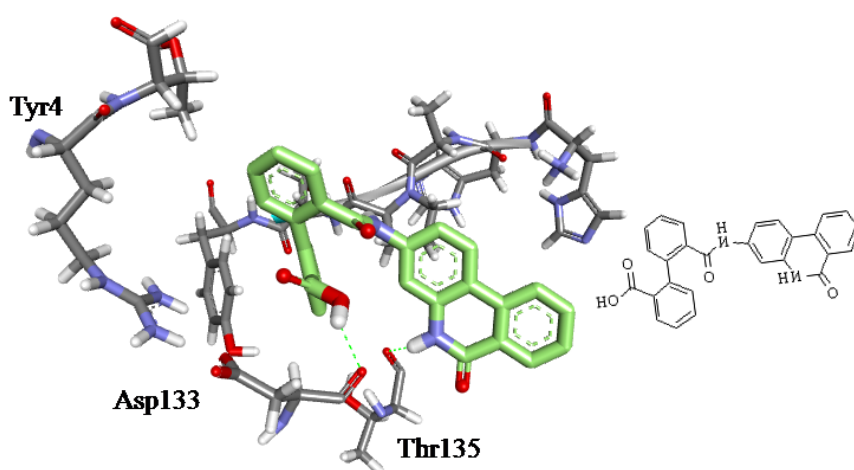
A



B



C



D

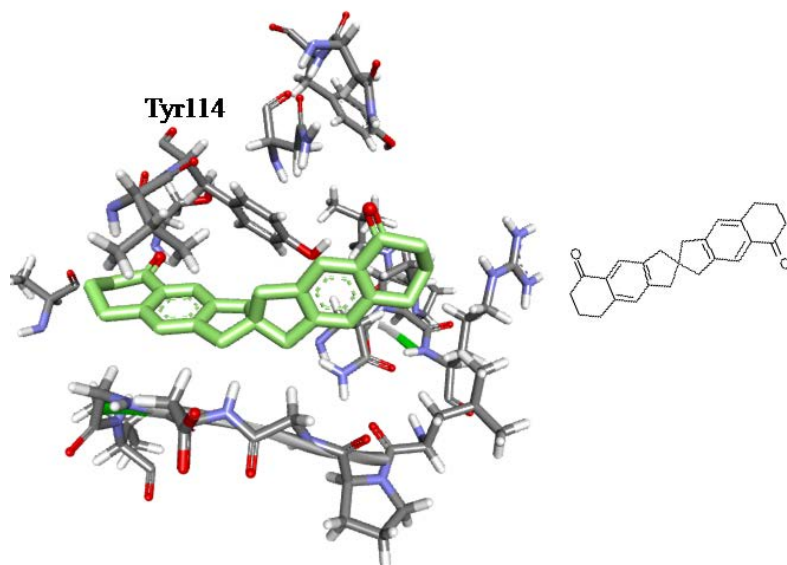
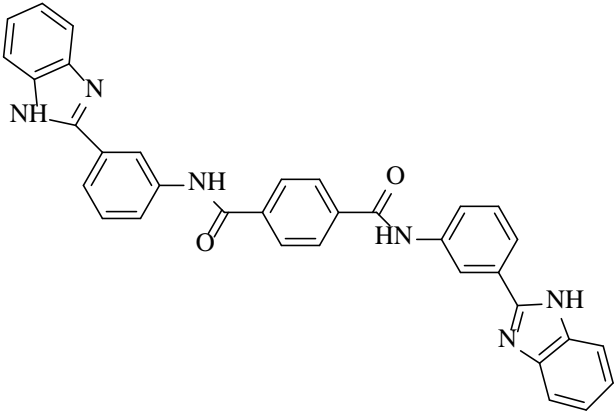
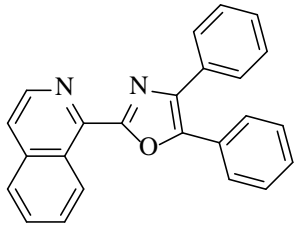
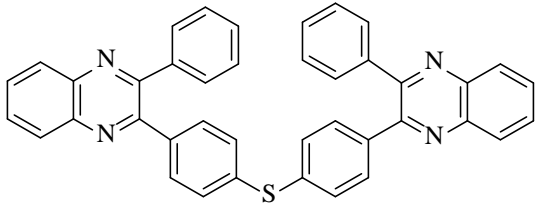
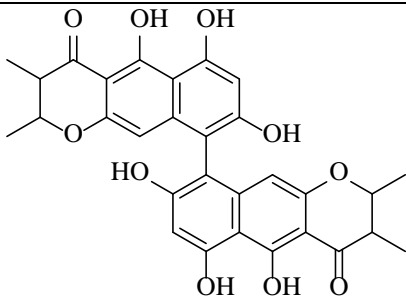
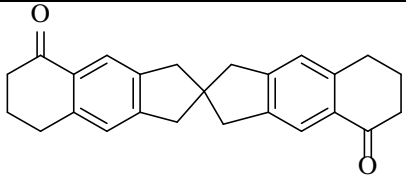


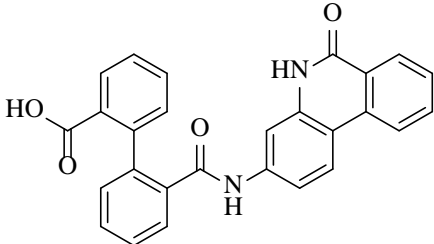
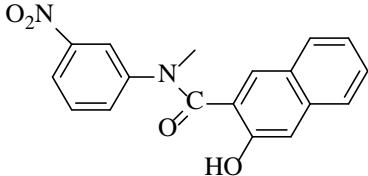
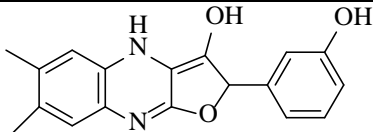
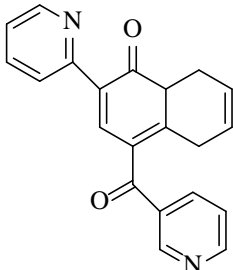
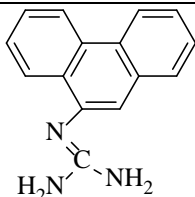
Figure 2.7. Binding pose and interactions of hit compounds in the nsP3 macrodomain: (A) NCI_25457 in Pocket 1: HBs with Val113 and π - π interaction with Trp148; (B) NCI_61610 in Pocket 1: HBs with Gly112 and π - π interaction with Trp148; (C) NCI_127133 in Pocket 2: HBs with Asp133; (D) NCI_670283 in Pocket 3: Hydrophobic contacts only. The ligands (in cyan) and the residues surrounding the ligands (in grey) were displayed in sticks and coloured by atoms (carbon in cyan in ligand or in grey in residues, nitrogen in blue, oxygen in red).

2.3.3 MD SIMULATIONS

MD simulations were undertaken to investigate the stability of the protein and its complexes as well as gain insights into the accurate binding modes of the protein and its inhibitors. The top-hit compounds NCI_61610, NCI_25457, NCI_345647_a, NCI_670283, and NCI_127133; and the tenth-hit compounds NCI_37168, NCI_372499_b, NCI_37168, NCI_324623, NCI_400976, and NCI_293778 from each screening were also subsequently submitted to MD simulations (Table 2.6). MD simulations were carried out with the NAMD package with the CHARMM force field for 50 ns following 3 ns equilibrium simulations.

Table 2.6. Chemical structures of five top hit compounds for the nsP3 macrodomain and their properties.

Compound's name	Structural formula	Binding affinity (kcal/mol)	Lipinski's values
NCI_61610		-11.1	LogP ^a : 5.31 H-D ^b : 4 H-A ^c : 4 MW ^d : 548.60
NCI_25457		-10.8	LogP: 5.37 H-D: 0 H-A: 3 MW: 348.39
NCI_293778		-9.4	LogP: 10.87 H-D: 0 H-A: 4 MW: 594.73
NCI_345647_a		-10.9	LogP: 2.17 H-D: 6 H-A: 10 MW: 546.52
NCI_670283		-10.6	LogP: 4.81 H-D: 0 H-A: 2 MW: 356.45

NCI_127133		-8.3	LogP: 4.52 H-D: 3 H-A: 4 MW: 434.4
NCI_37168		-10.1	LogP: 2.98 H-D: 1 H-A: 5 MW: 322.32
NCI_372499_b		-10.3	LogP: 1.91 H-D: 0 H-A: 3 MW: 308.34
NCI_324623		-10.3	LogP: 1.37 H-D: 0 H-A: 4 MW: 328.37
NCI_400976		-7.2	LogP: 2.50 H-D: 3 H-A: 4 MW: 237.31

(a) A calculated octanol-water partition coefficient; (b) H-D: Hydrogen bond donor; (c) H-A: Hydrogen bond acceptor; (d) MW: Molecular weight.

2.3.3.1 Overall stability of the nsP3 and its complexes

In order to assess overall stability, the values of positional RMSD for backbone from the starting structure are used as a major criterion. Backbone RMSD curves for the nsP3 and its complexes with different ligands with respect to the starting structure after the systems reached equilibrium within 3 ns, are shown in Figure 2.8. The plots showed that most of the systems were relatively stable during the 50 ns simulations within 1-2.5 Å.

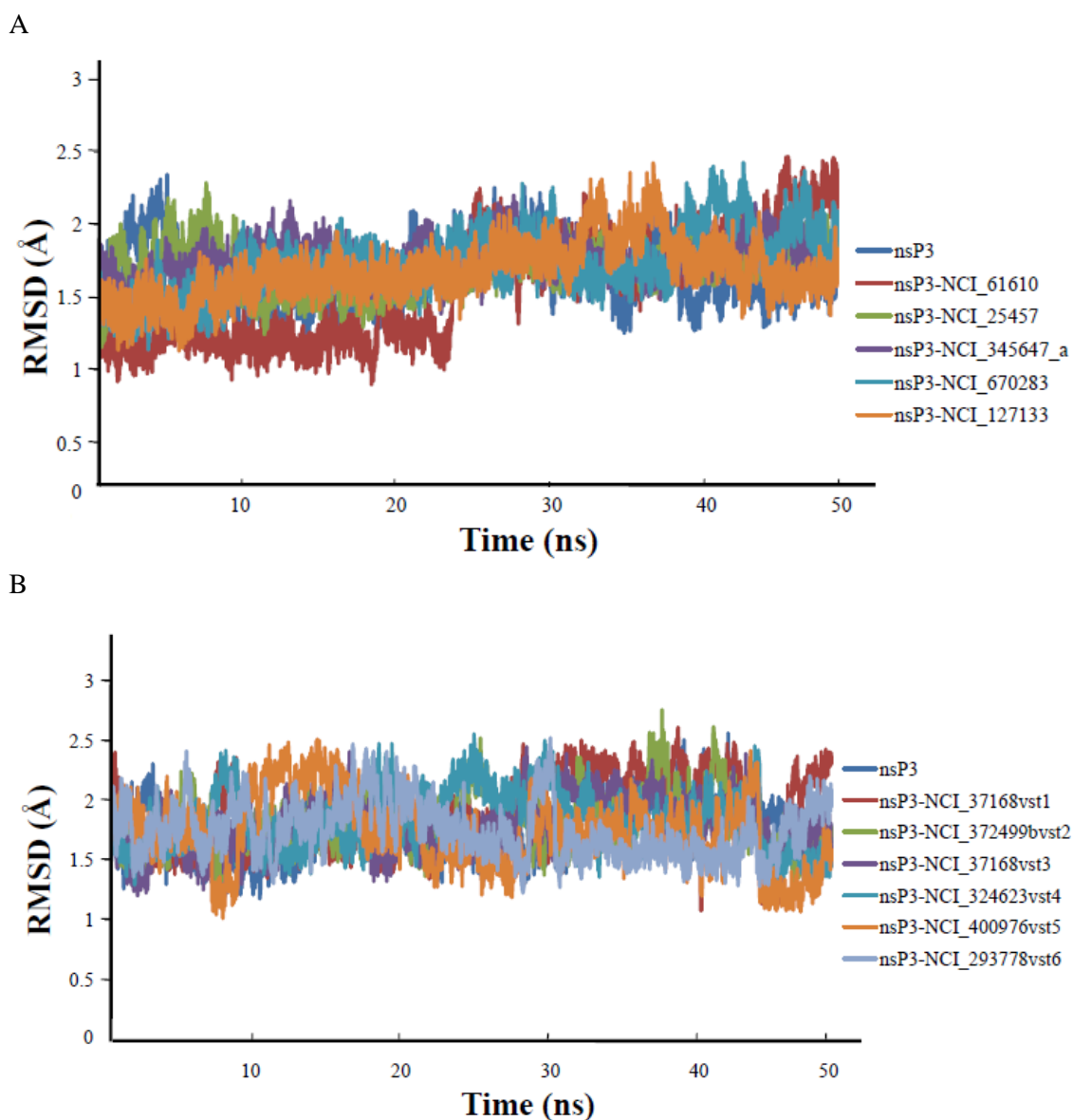


Figure 2.8. The backbone RMSD profiles for the apo protein nsP3 and its different complexes during MD simulations: (A) Complexes of the nsP3 and top-hit compounds; (B) Complexes of the nsP3 and tenth-hit compounds.

2.3.3.2 Investigating the flexibility of the nsP3 and its complexes

To understand the flexibility of the complex, the RMSF of the C α atoms of each residue was calculated from the trajectory data for 50 ns for the protein nsP3 and its complexes. The RMSF profiles presented in Figure 2.9 show they are comparatively similar between the apo protein and the complexes.

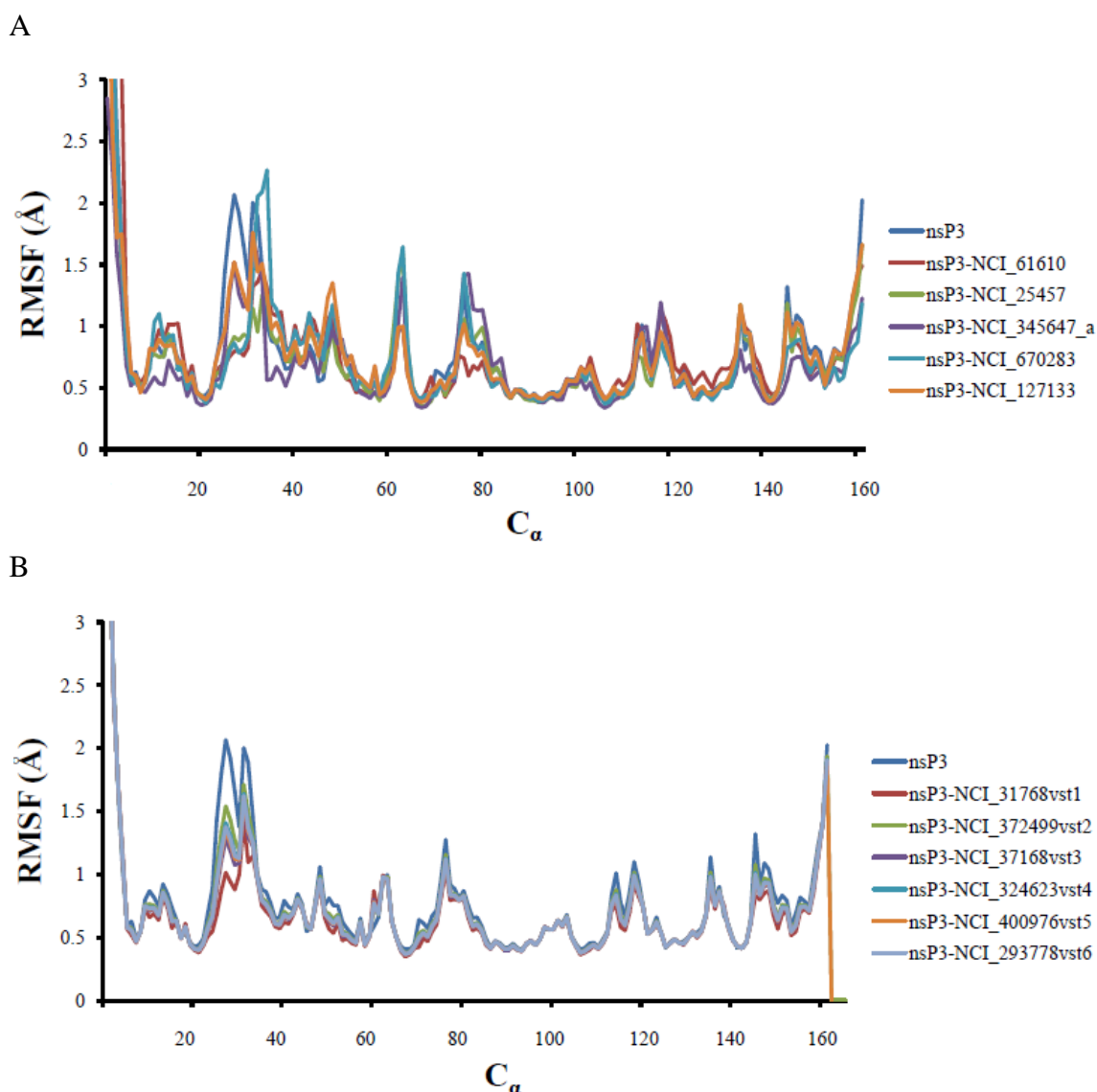


Figure 2.9. RMSF values of C_α atoms of the apo protein nsP3 and its different complexes during MD simulations: (A) Complexes of the nsP3 and top-hit compounds; (B) Complexes of the nsP3 and tenth-hit compounds.

With regards to hit compounds in Figure 2.9 (A) and (B); the residues making up the binding pockets were quite stable during the simulations (the fluctuation within 1.0 Å). Subtle differences were observed for a few regions, including the loop at residue 31-34. It is worth noting that for three ligands NCI_61610, NCI_25457, and NCI_345647_a bound to the protein at Pocket 1, the RMSFs for the binding loop region 31-34 were decreased compared to those in the apo protein. However, this RMSF for this loop was not significantly perturbed for the ligands NCI_670283 and NCI_127133 when bound to Pocket 3 and Pocket 2, respectively.

2.3.3.3 Atomic interaction between the protein nsP3 and ligands

Detailed analysis of the interactions between the ligands and nsP3 were carried out on the HBs interactions and hydrophobic contacts. The outcomes are listed in Table 2.7 and Table 2.8.

Table 2.7. Hydrogen bonding analyses on the trajectories sampled in MD simulations of hit compounds for the nsP3.

	Number	Details of HBs	% occupancy
NCI_61610	5	Asn24 (HD22)-O1	98
		Tyr114 (HN)-O	92
		Gly112 (HN)-O	88
		Thr111 (OG1)-H1	13
		Cys34 (HG1)-O1	10
NCI_25457	3	Val113 (HN)-N	29
		Val33 (HN)-N	21
		Val33 (HN)-O	20
NCI_345647_a	7	Ile11 (HN)-O4	39
		Ile11 (HN)-O6	12
		Gly112 (HN)-O1	18
		Val33 (HN)-O	17
		Gly32 (O)-O5	14
		Thr111 (HN)-O1	13
		Arg144 (HE)-O3	10
NCI_670283	4	Thr111 (HN)-O	77
		Gly112 (HN)-O	65
		Ser110 (HN)-O	25
		Thr111 (HG1)-O	10
NCI_127133	1	Arg159 (HH12)-O2	22
NCI_37168vst1	7	Asn24 (HD22)-N1	18
		Asn24 (HD22)-O2	20
		Asn24 (HD22)-O3	16
		Asp32 (HN)-N1	26
		Asp31 (HN)-O2	18
		Asp31 (HN)-O3	16
		Cys34 (HG1)-O2	10
NCI_372499vst2	4	Asp145 (OD1)-H	13
		Asp145 (OD1)-H1	19
		Asp45 (OD2)-H	19
		Asp145 (OD2)-H1	26
NCI_37168vst3	13	Asn24 (HD22)-N1	25
		Asn24 (HD22)-O2	24
		Asn24 (HD22)-O3	24
		Gly30 (HN)-O3	10
		Cys34 (HN)-N1	12
		Cys34 (HN)-O2	12
		Cys34 (HN)-O3	11
		Leu108 (O)-H1	92

NCI_37168vst3		Thr111 (HN)-O1	92
		Gly112 (HN)-O	90
		Gly112 (HN)-O1	69
		Val113 (HN)-O	49
		Tyr114 (HN)-O	90
NCI_324623vst4	5	Asn24 (HD22)-O	56
		Val33 (HN)-O1	49
		Cys34 (HG1)-N1	11
		Thr111 (HN)-O	24
		Gly112 (HN)-O	34
NCI_400976vst5	1	Tyr4 (HN)-N1	11
NCI_293778vst6	4	Asn24 (HD22)-N	81
		Arg26 (HE)-N2	32
		Arg36 (HH12)-N2	12
		Asp31 (HN)-N1	13

Table 2.8. Hydrophobic contact analyses on the trajectories sampled in the MD simulations of hit compounds for the nsP3.

Ligand	Non-polar part of residues
NCI_61610	Ala22, Pro25, Leu28, Val33, Pro107, Val113, Tyr114, Trp148
NCI_25457	Ala22, Val33, Pro107, Val113, Tyr114, Trp148
NCI_345647_a	Ile11, Val33, Ala36, Val113, Tyr114, Trp148
NCI_670283	Ala22, Leu28, Val33, Pro107, Val113, Tyr114
NCI_127133	Ala1, Pro2, Tyr4
NCI_37168vst1	Ala22, Ala23, Val33, Pro107, Thr111, Val113, Tyr114, Arg144
NCI_372499vst2	Val33, Arg144, Trp148
NCI_37168vst3	Ala22, Ala23, Val33, Pro107, Val113, Tyr114, Trp148
NCI_324623vst4	Ala22, Val33, Phe45, Pro107, Val113, Tyr114
NCI_400976vst5	Ala1, Pro2, Tyr4, Phe129, Arg159
NCI_293778vst6	Ala22, Ala23, Pro25, Arg26, Leu28, Val113, Tyr114

For hit compounds, hydrogen bonding and hydrophobic contact analyses indicated that all investigated ligands are stabilized the protein by a number of HBs (Tables 2.7 and 2.8). Complementary to the docking where the protein was kept rigid, MD simulations revealed that when the ligands bind to the nsP3, the ligand and/or the residues in the binding pockets fluctuate and adapt their structure in order to better accommodate the ligands by optimizing HBs and/or hydrophobic contacts. Most ligands at Pocket 1 and Pocket 3 bound strongly to the nsP3; always displaying strong HBs, particularly with Asn24, Val33, Cys34, Thr111, Gly112, Val113, and Arg114, while Tyr4 is important in Pocket 2.

Moreover, based on the occupancy of HBs for each ligand, the results showed that Pocket 1 could have the highest potential for ligand binding, as Pocket 1 had more HBs with higher occupancy than those in Pocket 3 > Pocket 2. As mentioned above, some residues, such as Val33, Val113, Tyr114, Arg144, and Trp148 (for compounds binding to Pocket 1), displayed noticeable movement upon binding. For the ligands NCI_127133 (Pocket 2), and NCI_670283 (Pocket 3), the results showed the fluctuation of residues Ala1, Pro2, and Tyr4 (at Pocket 2), and residues Val113, and Tyr114 (Pocket 3), were required for a correct fit into the nsP3. In particular, some ligands (NCI_61610 and NCI_37168) established the strong HBs interactions with the nsP3 at Pocket 1, emphasized on the residues at region of 110-114. For instance, occupancy of HBs of NCI_61610 and residues at Pocket 1, namely Asn24, Tyr114, Gly112 was 98%, 92%, and 88%, respectively. Ligand NCI_37168 hydrogen-bonded with residues Thr111 (92%), Gly112 (90%), Tyr114 (90%); and two medium hydrogen bonds between NCI_670283 and Thr111 (77%); and Gly112 (65%) at Pocket 3 were observed. For ligand NCI_61610 at Pocket 1, half of the ligand was quite stable, while the other half was flexible enough to fit well through interacting with Trp148 by π - π interaction and forming HBs with residues in the region 110-114 (illustrated in Figure 2.10). For ligand NCI_670283 at Pocket 3, the part of ligand which interacted with the region of residues from 110-114 was optimized to fit well in the pocket, even though the frequency of HBs interactions were medium with Thr111 (77%), Gly112 (65%), and weak with Ser110 (25%), respectively (Figure 2.10).

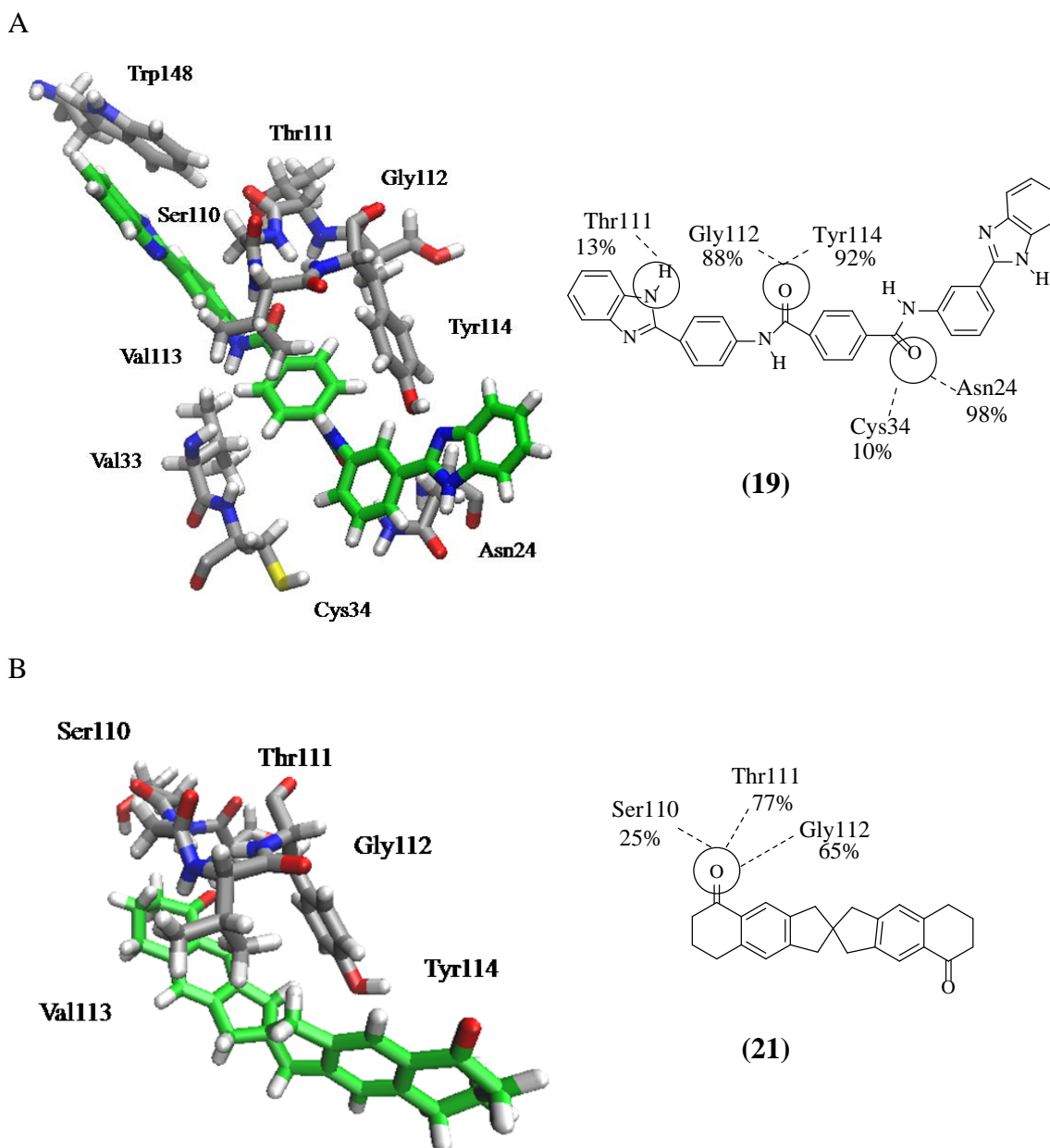


Figure 2.10. Hydrogen bonding interactions between the nsP3 and ligands: **(A)** Ligand NCI_61610 at Pocket 1 and **(B)** Ligand NCI_670283 at Pocket 3, with representation of ligands and key residues for interactions surrounding the ligands (in stick).

Furthermore, the key residues interacting in the region 110-114 of the protein at Pockets 1 and 3, served as hydrogen bond donors, except in the complex nsP3-NCI_61610, where residue Thr111 served as an acceptor. This observation is in close agreement with docking results as well as with earlier simulations and experimental data.⁷⁰ In addition, it emerged that the structure of NCI_61610 and NCI_345647_a are more polar, thus more hydrogens bonds were found in their complexes than others.

The solvent accessible surface areas (SASA) were also calculated to monitor the possible solvation environment change upon ligand binding (Appendix 8). It was expected that the SASA for hydrophobic interacting residues in the complex protein-ligand would be decreased compared to those in apo protein. At Pocket 1, when the ligands bind to protein, it can be seen that the SASA of residue Tyr114 showed a decrease from 70.2 Å² in the apo nsP3 to 63.4 Å² for the nsP3-NCI_61610 complex, and 63.8 Å² for nsP3-NCI_345647_a complex. Also, the SASA of Val33 displayed a reduction from 68.8 Å² in the apo state to 57.4 Å² and 52.9 Å² for the bound state in nsP3-NCI_25457, and NCI_345647_a, respectively. However, changes observed in SASA for Val33, Val113, and Trp148 are not consistent for different ligands. It can be rationalised that these residues were not only able to form hydrophobic contacts, but also form polar hydrogen bonding interactions. Thus, the change in SASA will be compromised by the polar interactions and both of protein and ligands will modulate the SASA values.

2.3.3.4 Clustering analysis

Throughout the simulations, the complex structure of the protein and ligands could vary under the effects of environment, so structural clustering was used to identify the most popular conformation during the simulation, and more important to compare the structures from MD simulation and docking. Clustering analysis was carried out on all of the snapshots from the trajectories, and the clusters were visualized and superimposed with the initial structure. The value of RMSD was used to evaluate the difference between clustering structures and initial structure. The different conformations of protein and its complexes at 0 ns, 10 ns, 20 ns, 30 ns, 40 ns, and 50 ns were superimposed. Slightly fluctuations were found in all complexes, though had very little significance. For example, for the ligand NCI_61610, in Figure 2.11, superimposition of the most popular conformations of protein obtained from simulations and RMSD value for the ligand was 1.2 Å (at 10 ns), 1.1 Å (at 20 ns), 1.8 Å (at 30 ns), 2.2 Å (at 40 ns), and 2.0 Å (at 50 ns), with respect to the initial structure.

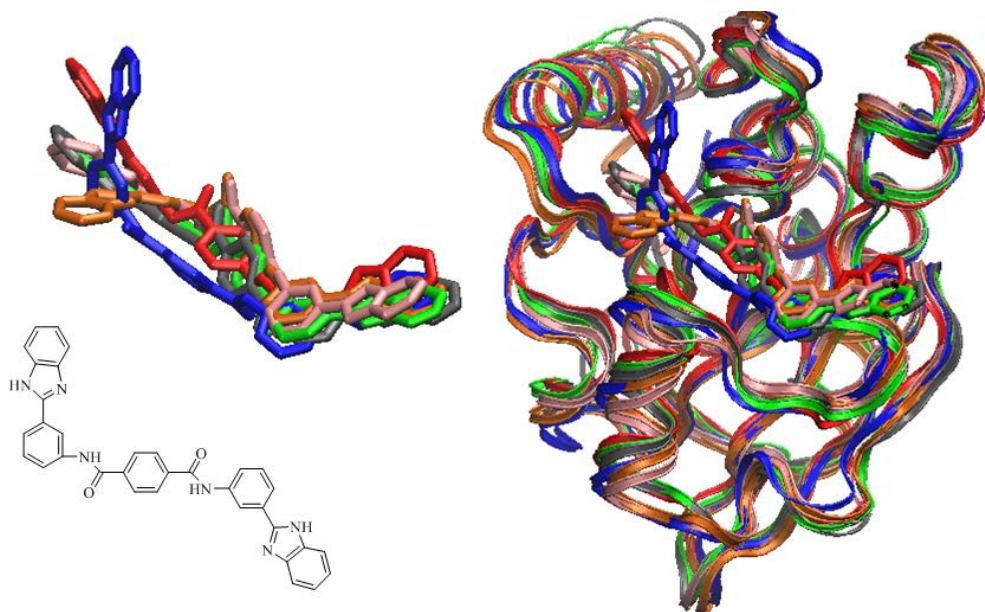


Figure 2.11. Superimposition of the different conformations of ligand and complexed ligand NCI_61610-nsP3 during simulation with the initial structure (red: at 0 ns, grey: at 10 ns, green: at 20 ns, pink: at 30 ns, orange: at 40 ns, and blue: at 50 ns).

2.3.3.5 Combination of MD simulations and docking for the nsP3

In order to probe the effects of the static protein structure used in the docking, multiple docking simulations and virtual screenings were carried out based on the sampled conformations of protein nsP3-NCI_61610 at the different timepoints (5 ns, 10 ns, 15 ns, and 20 ns). This complex was selected as NCI_61610 showed the highest potential for interacting with the protein after analyzing the outcomes of docking and simulations. These results are listed in Table 2.9 and Table 2.10. The binding affinities of these docking runs in Table 2.10 were not significantly different from those obtained from previous docking based on the X-ray structure. Additionally, the binding modes are similar (details of interactions analysis in Appendixes 9-12). For virtual screening, most of the top hits were the same compounds as previous screening indicated although there were some new hits. That indicates that in the case of nsP3, the docking results were not very sensitive for the static structure used.

Table 2.9. Re-docking results for complex nsP3-NCI_61610 with different conformations of the nsP3 protein taken from the different timepoints in simulations at Pocket 1.

nsP3 conformation	Binding affinity (kcal/mol)	Interaction between the inhibitor and residues of protein (with distance in Å)
At 0 ns	-11.1	Tyr114(HH)-N=2.0 Val33(HN)-O=2.3 Asn24(HD21)-O=2.4
At 5 ns	-10.3	Tyr114 (OH)-H1=1.9
At 10 ns	-11.3	Asn24 (HD22)-O=2.1 Tyr114 (OH)-H1=2.3
At 15 ns	-10.6	Asn24 (HD22)-O1=2.5 Ser110 (HN)-O=1.9
At 20 ns	-11.4	Tyr114 (OH)-H1=2.3

Table 2.10. Virtual screening results for blind docking into Pocket 1 with different conformations of the nsP3 taken from the different timepoints in simulations. The binding affinities are shown in kcal/mol.

VST-5ns ^a	VST-10ns ^b	VST-15ns ^c	VST-20ns ^d
1. NCI_293778 (-11.6)	1. NCI_37553 (-11.1)	1. NCI_293778 (-10.5)	1. NCI_37553 (-12.1)
2. NCI_84100_b (-10.7)	2. NCI_293778 (-11.1)	2. NCI_308835 (-10.3)	2. NCI_61610 (-11.9)
3. NCI_84100_a (-10.6)	3. NCI_60785_a (-11.0)	3. NCI_37553 (-10.3)	3. NCI_670283 (-11.7)
4. NCI_80997_b (-10.5)	4. NCI_59620_a (-10.8)	4. NCI_97920 (-10.0)	4. NCI_293778 (-11.3)
5. NCI_37553 (-10.4)	5. NCI_27592_a (-10.7)	5. NCI_84100_b (-9.9)	5. NCI_60785_a (-11.1)
6. NCI_61610 (-10.3)	6. NCI_670283 (-10.7)	6. NCI_59620_a (-9.8)	6. NCI_63680 (-11.1)
7. NCI_670283 (-10.2)	7. NCI_82802_a (-10.7)	7. NCI_84100_a (-9.8)	7. NCI_82802_a (-11.0)
8. NCI_59620_a (-10.2)	8. NCI_328101 (-10.5)	8. NCI_37627 (-9.8)	8. NCI_219894 (-10.8)
9. NCI_308835 (-10.2)	9. NCI_59620_a (-10.5)	9. NCI_60785_b (-9.7)	9. NCI_80997_b (-10.7)
10. NCI_60785_a (-10.2)	10. NCI_308835 (-10.5)	10. NCI_25457 (-9.7)	10. NCI_328101 (-10.7)

(a) VST-5ns, the grid box is fixed at the centre of Pocket 1 (9.7 Å, 43.0 Å, -13.2 Å) with a dimension of 20 Å × 20 Å × 20 Å. (b) In VST-10ns, the grid box is fixed at the centre of Pocket 1 (9.7 Å, 43 Å, -13.2 Å) with a dimension of 20 Å × 20 Å × 20 Å. (c) In VST-15ns, the grid box is fixed at the centre of Pocket 1 (9.7 Å, 43 Å, -13.2 Å) with a dimension of 20 Å × 20 Å × 20 Å. (d) In VST-20ns, the grid box is fixed at the centre of the protein at the centre of Pocket 1 (9.7 Å, 43 Å, -13.2 Å) with a dimension of 20 Å × 20 Å × 20 Å.

2.3.4 BINDING FREE ENERGY CALCULATIONS FOR THE LIGANDS BINDING TO THE NSP3 PROTEIN

The converged trajectories of complexes obtained from simulations were used to calculate the binding free energy for the hit compounds. The average values of vdW and elec interactions for each ligand in the two states, bound state (V_{bound}) and unbound state (V_{unbound}) were calculated (Table 2.11).

In applying the LIE equation, the vdW values of the ligands complexed with the enzyme were more negative than when only in solution, showing that the ligand has more favourable vdW interactions when complexed. The electrostatic interactions of ligands in all the complexes were significantly less favourable than those of the ligand in water. When using the values of empirical parameters α , β and γ , for example $\beta = 0.43$ for neutral compounds, $\alpha = 0.18$ and $\gamma = 0.0$,³²⁸ the value of binding free energies of all of complexes did not agree with the docking results (ΔG^1 in Table 2.11). However, as it has been shown that α and γ need to be recalibrated depending on the different systems; the α value has been suggested to be dependent on the system and the force field used in the LIE calculations while the γ relies on the nature of the binding site. In this case, due to the lack of experimental data of complexes, the chemical nature of the ligands and the binding sites were taken into the consideration. Most binding sites are composed of both polar and non-polar residues, so the magnitude of hydrophobicity in the selection of γ is not easy to define in practice, the exact value of γ does not affect the relative ranking. A larger value of $\alpha = 1.043$ was adopted as it had been shown to provide a better estimate in the study of cytochrome P450-camphor analogue complexes.³²⁴ The results of ΔG^2 , with α set to 1.043, are shown in Table 2.11, and revealed better agreement with the binding affinity obtained from docking. This gave a good explanation when comparing the binding free energies for different ligands at the same pocket. It could be explained that van der Waals and electrostatic results were compromised by hydrogen bonding interactions and hydrophobic contacts between the ligands and the nsP3 that contributes to the binding free energy.

4.4 CONCLUSIONS

Taking advantage of the envelope glycoprotein complexes, virtual screening based on blind docking and focused dockings explored the potential binding pockets and inhibitors for both the immature and mature structures of envelope proteins. Promising hit compounds were identified for two complexes of the envelope glycoproteins. Pocket 2 was a novel binding site in the immature structure. Pocket 2 and Pocket 3 were novel binding sites for the mature complex of the glycoproteins. The key residues involved in stabilizing the complex or participating in the fusion process were confirmed. This study also supported the current docking protocol utilising AutoDock Vina as robust and with a good accuracy, and could be used to identify inhibitors. However, due to larger size of two complexes of glycoproteins, the immature and mature structures, molecular dynamics simulations could not be carried out during the project. Therefore, the results of hit compounds and their binding modes obtained from docking are a good starting point for further studies. Further experiments are required to test the inhibitory effects for the anti-CHIKV envelope glycoproteins compounds.

simulations were used to calculate binding free energy more accurately than the results from docking. However, the challenge in the LIE method lies in the parameterization of the required co-efficients, which rely on the availability of experimental data. The final limitation was the time and computational cost for simulations, such as with large protein, for example envelope glycoproteins, which could be rectified given more time.

In summary, a combination of docking and MD simulations showed a potential approach that balances the computational cost and accuracy. More importantly, the procedure can be applied in discovering therapeutic compounds for other diseases. It can also help to generate the library for CHIKV inhibitors.

In Equation 1.16, α , β , and γ are empirical parameters. α is often set to 0.18 for a wide variety of ligand-protein systems. The β represents a function of the chemical nature of the ligand, so in principle the value of β can be parameterized from explicit solvent free energy calculations of different chemical entities. The γ parameter can be set $\gamma = 0$ or $\gamma \neq 0$ depending on the hydrophobicity of the binding site to estimate absolute free energies of binding.

In future work, experimental studies are required to determine the inhibitory effects of hit compounds on these membrane proteins. It would be useful to investigate the stability and flexibility of the systems. The mechanism of virus entry and virus attachment through the envelope glycoproteins would require further MD simulations experimental data. The binding free energy calculations would be also needed to guide rational antiviral drug design.

