

Bereza Andrew Nicolaevich

The Volgodonsk Institute of Service (branch) of the South-Russian State University of Economy and Service.

E-mail: anbirch@mail.ru.

6, Volgodonsk, Chernikova Street, Volgodonsk, 347360, Russia.

Phone: +79281574449.

The Department of Information; Cand. of Eng. Sc.; Associate Professor; Head of Department.

Liyshov Maksim Vasilievich

E-mail: max185@mail.ru.

Phone: +79604591974.

Master of Engineering.

УДК 004.382.2

А.И. Дордопуло, И.И. Левин, Д.А. Сорокин

**РЕАЛИЗАЦИЯ ДОКИНГА ДЛЯ МОЛЕКУЛЯРНОГО МОДЕЛИРОВАНИЯ
НА РЕКОНФИГУРИРУЕМЫХ ВЫЧИСЛИТЕЛЬНЫХ СИСТЕМАХ**

Статья посвящена описанию методов оптимизации фрагментов задачи и адаптации архитектуры реконфигурируемой вычислительной системы под структуру решаемой задачи при аппаратной реализации докинга – метода молекулярного моделирования. Отличительной особенностью описываемого решения по сравнению с известными реализациями является функционально завершенное решение полной задачи докинга на реконфигурируемой вычислительной системе, обеспечивающее согласованность функционирования всех фрагментов задачи в едином вычислительном контуре.

Аппаратная реализация; докинг; суперкомпьютерное молекулярное моделирование; реконфигурируемые вычислительные системы; лиганд.

A.I. Dordopulo, I.I. Levin, D.A. Sorokin

**DOCKING REALIZATION FOR MOLECULAR MODELING
ON RECONFIGURABLE COMPUTER SYSTEMS**

The paper is devoted to description of methods of task fragments optimization and adaptation of architecture of reconfigurable computer system to the structure of the solving task of docking (method of molecular modeling) hardware realization. In comparison with existing realizations, the distinctive feature of the viewed solution is all-in-one solution of complete problem of docking on reconfigurable computer system, providing coordinated functioning of all fragments of the task in a single computer system.

Hardware realization; docking; supercomputer molecular modeling; reconfigurable computer systems; ligand.

Введение. Создание новых лекарственных средств и повышение эффективности существующих препаратов является актуальной сферой применения высокопроизводительной вычислительной техники. Правильный выбор перспективных соединений-кандидатов для исследуемого препарата в значительной степени определяет материальные затраты, а также продолжительность и эффективность последующих этапов исследований нового лекарства, занимающих в среднем около 5 лет.

Для выбора соединений-кандидатов широко используются методы молекулярного моделирования, одним из которых является докинг молекулы-ингибитора (лиганда) в активный центр молекулы-мишени (белка). Процедура докинга [1,2] представляет собой перебор пространственных конфигураций молекулы-лиганда, для каждой из которых выполняется оценка энергии связывания

молекулы-лиганда и молекулы-белка. Докинг обеспечивает достаточную точность получаемых решений, но характеризуется высокой вычислительной сложностью, которая обусловлена большим числом входных параметров для модели взаимодействующих молекул и принадлежностью задачи к классу NP [2]. Поэтому на практике при решении задачи докинга широко используют суперкомпьютерные технологии моделирования, эмпирические методы сокращения пространства перебора [3] и, в частности, генетический алгоритм [2,4]. Использование кластерных многопроцессорных вычислительных систем позволяет сократить время решения задачи докинга, но не приводит к росту производительности, пропорциональному числу задействованных вычислительных узлов.

Для более эффективного решения задачи докинга необходимы средства адаптации архитектуры вычислительной системы к структуре решаемой задачи, которыми обладают, в частности, реконфигурируемые вычислительные системы (РВС), построенные на основе программируемых логических интегральных схем (ПЛИС) [5], успешно применяемые для решения вычислительно-трудоемких задач. Известны [6–10] реализации отдельных фрагментов докинга на РВС [6–10], но качественно новое решение задачи докинга, удовлетворяющее требованиям по быстродействию, возможно при согласованной аппаратной реализации всех фрагментов задачи в едином вычислительном контуре.

Описание математической модели и структуры решаемой задачи. Настоящая статья описывает аппаратную реализацию метода докинга лигандов в активный центр белка-мишени, обоснование и подробное математическое описание которого представлено в работе [2]. Физическим смыслом докинга является поиск оптимального положения лиганда в поле белка, которое характеризуется минимумом энергии связывания комплекса «лиганд-белок», что в вычислительном плане представляет собой расчет значения энергии такого комплекса, зависящего от пространственной конфигурации и ориентации лиганда в многомерном пространстве.

Поиск оптимального положения лиганда в поле белка является итерационным процессом, в котором можно выделить следующие этапы: создание очередной пространственной конфигурации лиганда (с учетом его ориентации относительно центра белка), вычисление энергетических характеристик созданной конфигурации и оценку конфигурации лиганда по критериям задачи.

Наибольшая вычислительная трудоемкость характерна для первых двух этапов. Так, на этапе создания очередной пространственной конфигурации выполняются расчет конфигурации лиганда с учетом вращения фрагментов молекулы (гибкий докинг) и позиционирование лиганда в многомерном пространстве с учетом вращения молекулы как целого, что требует неоднократного пересчета координат каждого атома. На этапе вычисления энергетических характеристик созданной конфигурации выполняется расчет общей энергии связывания, состоящей из трех основных слагаемых, вычисление каждого из которых является вычислительно трудоемким и зависит от рассчитанных на предыдущем этапе координат и типа каждого атома лиганда.

На этапе оценки конфигурации лиганда по критериям задачи выполняется сравнение рассчитанных энергетических характеристик текущей конфигурации с допустимым по условию пороговым значением и с лучшим из рассчитанных значений для определения целесообразности хранения и дальнейшего использования полученной конфигурации.

Для решения задачи используются следующие параметры:

- ◆ число атомов лиганда не превышает 200, каждый атом характеризуется своими пространственными координатами x , y и z (32-разрядные вещественные числа);

- ◆ число точек внутреннего вращения (торсионные степени свободы) не превышает 20, каждая точка внутреннего вращения характеризуется углом в интервале $[0, \pi]$ и положением в молекуле (номера двух атомов, характеризующие неподвижную и подвижную части молекулы);
- ◆ угол вращения молекулы как единого целого в интервале $[0, \pi]$ задается с помощью системы кватернионов.

Таким образом, максимальное число параметров, характеризующих конфигурацию лиганда, составляет 27, что позволяет с большим запасом учитывать как существующие, так и большинство планируемых к созданию лекарственных препаратов и белков.

Структура вычислений в задаче докинга представлена на рис.1. Для позиционирования лиганда в поле белка осуществляются расчет координат атомов лиганда с учетом поворота фрагментов молекулы в точках внутреннего вращения (блок R) и расчет координат всех атомов лиганда относительно центра молекулы белка с учетом поступательного и вращательного движения лиганда как целого (блок RT). Затем производится расчет суммарной энергии связывания E_{all} (блок E) для текущей конфигурации, значение которой поступает на блок оценки конфигурации лиганда (блок MPS). В зависимости от результатов оценки параметры текущей конфигурации либо сохраняются в списке лучших конфигураций и передаются на блок создания новых конфигураций (блок GEN), либо в случае неудовлетворительного значения суммарной энергии связывания не сохраняются и не учитываются.

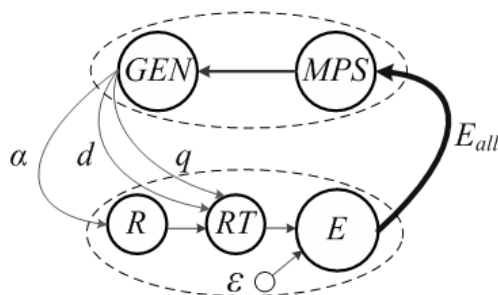


Рис. 1. Основные фрагменты задачи

Вычислительно-трудоемкими фрагментами в структуре задачи являются фрагменты R, RT, E и Gen, оперирующие 32-разрядными значениями с плавающей запятой.

Аппаратной платформой для решения задачи докинга является плата вычислительного модуля (ПВМ) 16V5-75, содержащая 16 вычислительных ПЛИС XC5VLX110 по 11 млн. эквивалентных вентилях в каждой, 2 Гбайта распределенной памяти. Тактовая частота составляет 250 МГц, а производительность для 32-разрядной арифметики с плавающей запятой составляет 140 Гфлопс. Внешний вид ПВМ 16V5-75 представлен на рис. 2.

Для эффективного решения задачи докинга на PBC необходимо реализовать этапы создания новых конфигураций лиганда, вычисления энергетических характеристик конфигурации и оценки конфигурации лиганда, обеспечив при этом как их согласованную работу, так и сбалансированную загрузку фиксированного аппаратного ресурса – одной ПВМ 16V5-75.

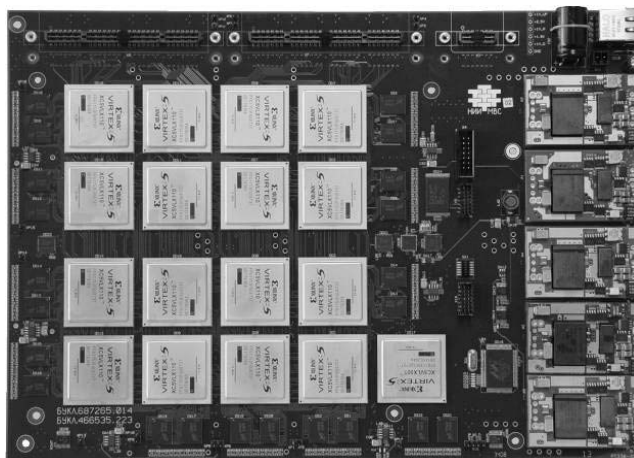


Рис. 2. ПБМ 16V5-75

Оптимизация фрагментов задачи и адаптация архитектуры РВС под структуру решаемой задачи. Для решения задачи докинга на ПБМ 16V5-75 необходимо в едином вычислительном контуре построить эффективный вычислительный конвейер, удовлетворяющий условиям задачи по критерию «производительность/объем аппаратного ресурса», для чего, как правило, используется структурный [5] метод организации вычислений на РВС.

Предварительный анализ структурной реализации задачи докинга показал, что требуемый вычислительный ресурс значительно превышает имеющийся на ПБМ 16V5-75, поэтому для успешного решения задачи необходимо выполнить оптимизацию организации вычислений для вычислительно-трудоемких фрагментов и адаптировать ресурсы РВС под структуру решаемой задачи.

Для удовлетворения заданным параметрам по занимаемому ресурсу и быстродействию к фрагментам задачи докинга были применены следующие разработанные методы оптимизации:

1) редукция вычислительной структуры графа фрагмента задачи путем эквивалентных математических преобразований исходных формул, позволяющая сократить занимаемый аппаратный ресурс реконфигурируемой вычислительной системы;

2) использование предвычисленных массивов для хранения результатов повторяющихся операций, значения которых не меняются для данного итерационного процесса, позволяющее рационально использовать ресурс внутренней памяти ПЛИС и сократить объем используемого вычислительного оборудования;

3) согласованное распараллеливание подграфов графа задачи, обеспечивающее сбалансированную загрузку вычислительных устройств для существенно переменных по интенсивности информационно-значимых потоков данных, возрастающих на два десятичных порядка в процессе решения по сравнению с входным потоком данных;

4) использование специальной структуры хранения данных с распределением по банкам динамической памяти, позволяющее в 9-10 раз ускорить доступ к памяти.

Рассмотрим применение разработанных методов на примере структурной реализации вычислительно-трудоемких фрагментов задачи.

Вычислительно-трудоемкий фрагмент R выполняет расчет внутренней геометрии лиганда для текущей конфигурации путем преобразования координат каждого атома лиганда для каждой точки вращения. В вычислительной структуре задачи это соответствует циклическим вызовам этого фрагмента в последовательной програм-

ме, число которых определяется числом торсионных связей вращения, максимальное значение которых равно 20. Согласно принципам структурной реализации вычислений [5] циклические вызовы процедуры последовательной программы соответствуют параллельной по числу точек вращения реализации фрагментов R. Время вычисления преобразований R в этом случае можно вычислить по формуле

$$t_R = N_{atom} \times t_{atom}, \quad (1)$$

где N_{atom} – число атомов в лиганде, t_{atom} – время вычисления преобразований для одного атома лиганда. При 20-кратном распараллеливании фрагмента R t_{atom} составит один такт ПВМ (один такт работы ПВМ 16V5-75 на частоте 250 МГц составляет 4 нс). Следовательно, для максимального количества атомов $N_{atom}=200$ время вычислений t_R составит 200 тактов.

Преобразование декартовых координат атомов при прямой структурной реализации вычислений потребует 76 устройств, реализующих 32-разрядные математические операции в стандарте IEEE754. Структурная реализация с 20-кратным распараллеливанием преобразования R задействует 1520 устройств, что превышает доступный ресурс ПВМ 16V5-75 и приводит к необходимости сокращения занимаемого ресурса с помощью методов оптимизации организации вычислений во фрагменте R.

Применение редукции вычислительного графа (метод 1) для фрагмента R и предвычисленные значения, хранящиеся во внутренней памяти ПЛИС (метод 2), позволяют сократить аппаратный ресурс для одной реализации фрагмента R с 77 до 22 устройств. При требуемом 20-кратном распараллеливании фрагмента R теперь понадобится 440 устройств, реализующих 32-разрядные математические операции в стандарте IEEE754. Редукция вычислительного графа фрагмента R приводит к построению конвейера с последовательной обработкой координат атома лиганда, что увеличивает время обработки в 3 раза ($t_{atom}=3$, а t_R для максимального числа атомов составит 600 тактов), но позволяет сократить необходимый ресурс на построение вычислительной структуры R в 3,5 раза при сохранении точности вычислений.

Структурная реализация вычислительно-трудоемкого фрагмента RT, выполняющего пересчет декартовых координат атомов лиганда с учетом вращательных степеней свободы лиганда как целого, требует 51 устройство, реализующее 32-разрядные математические операции в стандарте IEEE754. Для согласования скорости обработки в фрагментах R и RT в последнем целесообразно использовать конвейерные вычисления с последовательной обработкой координат атома лиганда. Это позволит сократить оборудование в 3 раза (17 устройств) и обеспечить согласованную работу этих фрагментов задачи, т.е. $t_{RT}=t_R=600$ тактов.

Одним из наиболее вычислительно-трудоемких фрагментов является блок E, реализующий расчет энергии связывания комплекса «лиганд-белок» и содержащий расчет нескольких компонентов суммарной энергии по формуле:

$$E_{total} = E_{lig-prot} + E_{inner}, \quad (2)$$

где $E_{lig-prot}$ – энергия лиганда в поле протеина,

E_{inner} – внутренняя энергия лиганда.

Энергия лиганда в поле протеина представляет собой сумму трех составляющих E_0, E_1, E_2 , рассчитанных на трехмерных сетках потенциалов по MMFF94 для всех атомов лиганда [2]. Вычисление каждой составляющей ведётся методом трилинейной интерполяции. Структурная реализация вычислений $E_{lig-prot}$ требует 55 устройств, реализующих 32-разрядные математические операции в стандарте IEEE754. Время обработки одного лиганда определяется по формуле:

$$t_{E_{lig-prot}} = N \times t_{atom},$$

где t_{atom} – время чтения коэффициентов для одного атома лиганда из массивов трехмерных сеток потенциалов.

Для хранения массивов трехмерных сеток потенциалов необходимо около 200 Мб памяти, поэтому требуется задействовать распределенную память ПВМ 16V5-75, организованную на микросхемах типа SDRAM DDR2.

При вычислении энергии лиганда в поле протеина $E_{lig-prot}$ для каждого атома выполняется чтение коэффициентов из массивов трехмерных сеток потенциалов, что приводит к чтению из памяти по произвольному адресу, время которого определяется технологическими задержками для памяти типа DDR2. Память типа DDR2 имеет максимальную пропускную способность при «линейном чтении» в пределах строки (содержащей 512 32-разрядных слов), а поскольку переключение между строками микросхемы памяти при произвольной адресации будет происходить постоянно, то время обработки запросов по чтению будет максимальным.

Время обработки запроса к одному из массивов составит, в среднем, 15 тактов, а обработка одного атома лиганда $t_{atom} \approx 45$ тактов соответственно, что составит примерно 72000 тактов работы фрагмента расчета энергии лиганда в поле протеина для лиганда с максимальным числом атомов $N_{atom}=200$.

Для решения задачи в едином вычислительном контуре возникает необходимость согласования темпа обработки в фрагментах R , RT и $E_{lig-prot}$. Для этого к фрагменту $E_{lig-prot}$ применим методы оптимизации 3 и 4.

Первое преобразование состоит в сокращении объема используемых для данного лиганда данных в массивах трехмерных сеток потенциалов: в распределенную память ПВМ 16V5-75 загружаются только те коэффициенты, которые соответствуют типам атомов обрабатываемого лиганда, число сочетаний которых, как правило, для одного лиганда не превышает 8. Это позволяет уменьшить требования к объему памяти для массивов трехмерных сеток потенциалов до значения 32 Мб.

Второе преобразование состоит в объединении массивов трехмерных сеток потенциалов в единый массив, каждый элемент которого представляет собой кортеж из трех подряд идущих энергетических коэффициентов массивов трехмерных сеток потенциалов. Контроллер распределенной памяти, установленный на ПВМ 16V5-75, в случае нелинейного чтения позволяет организовать режим, при котором происходит чередование обращений к «строкам» между «банками» памяти, что позволяет сократить время обработки одного запроса при равномерном (или близком к нему) обращении к каждому банку памяти. При организации такого режима теоретическое время обработки одного атома лиганда t_{atom} составит 4 такта. На практике, среднее реальное время обработки одного атома не превышает 4,75 такта, что связано с наличием в атомах лиганда повторяющихся групп атомов водорода, углерода, кислорода и др.

Третье преобразование заключается в согласованном распараллеливании вычислений $E_{lig-prot}$ по восьми интерполяционным точкам для каждого атома лиганда. Это позволяет эффективно задействовать имеющуюся распределенную память ПВМ 16V5-75. Упорядоченный массив трехмерных сеток потенциалов дублируется для каждой из 8-ми точек и к нему организуется независимый канал доступа. Требуемый для такой организации вычислений объем памяти возрастает до 256 Мб, но в то же время увеличивает скорость обработки $E_{lig-prot}$ также в 8 раз.

Как показали практические исследования, при такой оптимизации время обработки лиганда $N_{atom}=200$ не превышает $t_{E_{lig-prot}} \approx 950$ тактов. Это значение времени обработки будем считать опорным для согласования скорости обработки фрагментов задачи докинга при объединении в единый вычислительный контур, реализация которого будет представлена во второй части статьи.

БИБЛИОГРАФИЧЕСКИЙ СПИСОК

1. Молекулярная стыковка: <http://ru.wikipedia.org/wiki/Докинг>.
2. Романов А.Н., Кондакова О.А., Григорьев Ф.В. и др. Компьютерный дизайн лекарственных средств: программа докинга SOL // Вычислительные методы и программирование. – 2008. – Т. 9. – С. 213-233.
3. Генетический алгоритм http://ru.wikipedia.org/wiki/Генетические_алгоритмы.
4. AutoDock. <http://autodock.scripps.edu/>.
5. Каляев И.А., Левин И.И., Семерников Е.А., Шмойлов В.И. Реконфигурируемые мультиконвейерные вычислительные структуры // Изд. 2-е, перераб., доп. / Под общ. ред. И.А. Каляева. – Ростов-на-Дону: Изд-во ЮНЦ РАН, 2009. – 344 с.
6. Van Court T. FPGA acceleration of rigid molecule interactions / T. Van Court, Y. Gu, M. Herbordt // Int. Conf. Field Programmable Logic and Applications (FPL 2004). – Antwerpen, Belgium, 2004. – P. 862-867.
7. Van Court T., Gu Y., Mundada M.C., Herbordt M.C. Rigid molecular docking: FPGA reconfiguration for alternative force laws // J Appl. Signal Processing v. 2006, 2006. – P. 1-10.
8. Herbordt M.C., Gu Y., Van Court T., Model J., Sukhwani B., Chiu M. Computing models for FPGA-based accelerations with case studies in molecular modeling // Porceed. of the Reconfigurable systems summer institute (RSSI 2008), 2008.
9. Sukhwani B. Acceleration of a production rigid molecule docking code / B. Sukhwani, M. Herbordt // Int. Conf. Field Programmable Logic and Applications (FPL 2008). – Heidelberg, Germany, 2008. – P. 341-346.
10. Sukhwani B., Herbordt M.C. FPGA acceleration of rigid-molecule docking codes // IET Computers & digital techniques (ACM-TRETS), 2009 (accepted for publication).

Статью рекомендовал к опубликованию д.т.н., профессор Я.Е. Ромм.

Левин Илья Израилевич

Научно-исследовательский институт многопроцессорных вычислительных систем имени академика А.В. Каляева федерального государственного автономного образовательного учреждения высшего профессионального образования «Южный федеральный университет».

E-mail: levin@mvs.tsure.ru.

347922, г. Таганрог, ул. Ленина, д. 224/1, кв. 65.

Тел.: 88634623226.

Заместитель директора по науке; д.т.н.

Сорокин Дмитрий Анатольевич

E-mail: jotun@inbox.ru.

347922, г. Таганрог, переулок Украинский, д. 21, кв. 30.

Тел.: 88634393820.

Научный сотрудник.

Дордопуло Алексей Игоревич

Учреждение Российской академии наук «Южный научный центр РАН».

E-mail: scorpio@mvs.tsure.ru.

347900, г. Таганрог, 10-й переулок, 114/1, кв. 6.

Тел.: 88634368651.

Старший научный сотрудник; к.т.н.

Levin Ilya Israilevich

Kalyaev Scientific Research Institute of Multiprocessor Computer Systems at Southern Federal University.

E-mail: levin@mvs.tsure.ru.

224/1, Lenin Street, Ap. 65, Taganrog, 347922, Russia.

Phone: +78634623226.

Deputy Director of Science; Dr. of Eng. Sc.

Sorokin Dmitry Anatolievich

E-mail: jotun@inbox.ru.

21, Ukrainskiy Lane, Ap. 30, Taganrog, 347922, Russia.

Phone: +78634393820.

Scientific Associate.

Dordopulo Alexey Igorevich

Southern Scientific Centre of the Russian Academy of Sciences.

E-mail: scorpio@mvs.tsure.ru.

114/1, 10th Lane, Ap. 6, Taganrog, 347900, Russia.

Phone: +78634368651.

Senior Staff Scientist; Cand. of Eng. Sc.

УДК 681.3

А.Э. Саак

**АНАЛИЗ ВЗАИМОДЕЙСТВИЯ ПОЛЬЗОВАТЕЛЕЙ
И ОБСЛУЖИВАЮЩЕЙ КОМПЬЮТЕРНОЙ СИСТЕМЫ**

Анализируется взаимодействие пользователей и обслуживающей компьютерной системы (МВС, Grid- системы) в форме комбинаторных многомерных моделей экспериментов спроса и предложения для некоторых основополагающих дисциплин обслуживания. Предлагается вариантный признак равновесия целочисленных сред в форме совпадения вариантных мощностей входа-выхода в системе компьютерного обслуживания множественного типа. На основе данных моделей исследуются явления переполнения спроса относительно общего ресурса предложений формализуемые усечением комбинаторных экспериментов.

Пропускная способность многопроцессорных и Grid- систем; однородно-ресурсное диспетчирование; комбинаторные эксперименты спроса-предложения.

A.E. Saak

**THE ANALYSIS OF AN INTERACTION OF USERS AND THE COMPUTER
SERVICE SYSTEM**

An interaction of users and the computer service system (MPS, Grid- system) in the form of combinatorial multidimensional models of demand and supply experiments for some fundamental service procedures is analyzed. It is suggested the variant sign of the equilibrium of integer-valued surroundings in the form of coincidence of input-output variant capacities in the multiplex type computer service system. Phenomena of demand overflow relative to shared supply resource that are formalizable by the truncation of combinatorial experiments are explored on basis of these models.

The capacity of multiprocessor systems and Grid- systems, uniformly resource dispatching control, the demand and supply combinatorial experiments.