

ОТЗЫВ

официального оппонента на диссертацию
Климентова Алексея Анатольевича
«Методы обработки сверхбольших объемов данных в распределенной
гетерогенной компьютерной среде для приложений в ядерной физике и физике
высоких энергий»,
представленную на соискание ученой степени
доктора физико-математических наук по специальности 05.13.11 –
«Математическое и программное обеспечение вычислительных машин,
комплексов и компьютерных сетей»

Компьютерные технологии развиваются очень быстро, проникая во всё новые и новые области. Если раньше наука традиционно держалась на двух составляющих: наука теоретическая и наука экспериментальная, то за последние годы появилась и прочно закрепилась третья составляющая – вычислительная. Наука вычислительная не только содержательно дополняет теорию и эксперимент, во многих случаях продуктивные исследования без нее просто невозможны.

Огромные возможности компьютерных систем, развитие методов математического моделирования, систем и технологий параллельного программирования в совокупности и определяют высокую степень достоверности результатов вычислительного эксперимента, на что и опираются ученые при проведении научных исследований. Вместе с этим, компьютерные системы стали исключительно разнообразными. Мобильные устройства, персональные и офисные компьютеры, серверы, вычислительные кластеры, суперкомпьютеры во всем диапазоне производительности от базового уровня до рекордных петафлопсных систем – всё это сегодня находится в распоряжении исследователей. В дополнение к этому, в настоящее время можно легко подобрать наиболее подходящую процессорную базу:

многоядерные, графические, реконфигурируемые, энергоэффективные и другие типы процессоров, что позволяет скомпоновать компьютерную установку, в наибольшей степени соответствующую решаемой задаче. В дополнение к этому, можно подобрать необходимые по своим характеристикам системы хранения данных, централизованные или распределенные, построенные на разных технологических принципах, определяющих скорость работы, стоимость, совместимость, надежность и другие параметры. В дополнение к этому, можно подобрать необходимое операционное окружение и требуемую инфраструктуру программного обеспечения. И что самое важное, всё сказанное выше сегодня технологически может быть скомпоновано в единую компьютерную среду, объединяя множество отдельных компонент в единый комплекс. Сегодня возможности для объединения есть, и это можно и нужно использовать на практике. Потенциал компьютерных сред огромен, но чтобы им реально воспользоваться нужно решить целое множество сложных вопросов, определяемых колоссальными масштабами и параметрами сред, распределённостью, неоднородностью, динамичностью, различными политиками использования компонент подобных сред. Это самая сложная задача, лежащая на стыке системного программирования, сетевых технологий, компьютерного дизайна, технологий хранения и передачи данных в приложении к масштабным вычислительным экспериментам в ядерной физике и физике высоких энергий. На решение именно этой задачи направлена диссертационная работа А.А.Климентова, что, несомненно, делает **актуальной избранную тему исследований**.

Диссертационная работа организована в виде четырёх глав следующего содержания. В первой главе автор показывает развитие модели computинга, используемой в физических экспериментах в рассматриваемых областях. Важный вывод заключается в том, что по мере развития экспериментальных установок быстрый рост наблюдается не только в объемах обрабатываемых данных. Идет активный рост числа самих участников экспериментов: сегодня

это тысячи людей в сотнях организациях в разных странах, что накладывает жесткие требования на организацию компьютерных сред, обеспечивающих эффективную работу экспериментов и экспериментаторов. Вторая глава посвящена требованиям к вычислительной инфраструктуре, анализу места высокопроизводительных вычислительных систем и суперкомпьютеров в рассматриваемых физических экспериментах. Один из основных выводов – потребности экспериментов в компьютерных ресурсах растут очень быстро, необходимо уметь использовать и суперкомпьютеры, и облачные инфраструктуры, и масштабные учебные сетевые возможности в рамках единой киберинфраструктуры, что, с учетом огромного масштаба по целому ряду параметров, требует нового подхода к ее проектированию. Третья глава диссертации посвящена разработке архитектуры системы управления заданиями в распределенных неоднородных компьютерных средах. Предложены и обоснованы основные модели, в частности, модель данных для общей системы распределенной обработки, описаны методики управления потоками заданий и распределения вычислительных ресурсов, описана структура подсистемы мониторирования для поддержки эксперимента ATLAS. С моей точки зрения это основная глава диссертации, которая читается с большим интересом, и в которой сходятся основные идеи, предложенные в данной работе. Четвертая глава посвящена приложениям и развитию предложенной компьютерной модели как в плане поддержки будущих масштабных физических экспериментов, так и для интеграции с компьютерными ресурсами НИЦ “Курчатовский институт”, и решению альтернативных классов задач, в частности, некоторых классов задач биоинформатики.

Хочу сразу отметить, что по ходу изложения материала автор приводит большое число реальных примеров использования разработок, построенных на основе предложенных подходов, методов и моделей. В этом ряду и использование суперкомпьютера Титан для поддержки эксперимента ATLAS, и поддержка эксперимента COMPASS на ускорителе SPS в ЦЕРН, примеры

использования систем семейства PanDA, примеры реальных систем мониторирования, распределения заданий и многое другое. Высокая эффективность, достигаемая в каждом случае, в сочетании с исключительно высокой нагрузкой, прекрасно подтверждает **высокую степень обоснованности научных положений, выводов и рекомендаций**, сформулированных в диссертации, и говорит о значительном потенциале разработок автора. Более того, обоснованность и достоверность результатов диссертации подтверждены активным использованием результатов научных исследований автора отечественными и зарубежными учеными в практике своей научной деятельности. Представленная работа прошла апробацию и подтвердила правильность сформулированных положений, моделей и методов в ходе больших международных экспериментов, результатами чего пользуются тысячи ученых во всем мире.

Научная новизна диссертационной работы заключается в разработке целого множества методов, моделей и программных систем, позволяющих эффективно обрабатывать большие объемы данных в распределенных и неоднородных компьютерных средах. Принципиально важно, что в данном случае эпитеты “большие”, “распределенные” и “неоднородные” соответствуют максимально сложным вариантам: “большие” объемы – это экзабайты, “распределенные” – это по всему миру, а “неоднородные” – это все доступные классы вычислительных систем от рабочих станций до гигантских облачных инфраструктур и суперкомпьютеров. Чтобы какая-либо программная система работала бы в столь сложных условиях, необходимо закладывать принципиально новые подходы в ее архитектуру и функционирование, что и было сделано автором.

Очень важно, что автору удалось найти хороший баланс между специализированностью и универсальностью созданных средств. Да, автор изначально ориентировался на поддержку масштабных экспериментов в физике высоких энергий и ядерной физике, и эта ориентация позволила ему найти, обосновать и реализовать такие решения, которые эффективно работают в

данной предметной области. Вместе с этим, потенциал распределенной обработки данных востребован в самых разных областях, эти технологии носят универсальный характер, и понимание этого факта позволило автору диссертации, не теряя эффективности работы в исходной области, использовать предложенные подходы для решения иных задач, привнеся универсальность в свои разработки. Это, несомненно, также является сильной стороной данной диссертационной работы, определяющей потенциал её будущих приложений в самых разных областях.

Соответствие содержания и результатов диссертации специальности 05.13.11 подтверждается соответствием по следующим пунктам:

1) Модели, методы, алгоритмы, языки и программные инструменты для организации взаимодействия программ и программных систем.

- Исследованы, разработаны и реализованы базовые принципы подсистемы мониторирования, включая предсказание времени выполнения заданий, и аккаунтинга (учета работы отдельных центров и заданий).

2) Модели и методы создания программ и программных систем для параллельной и распределенной обработки данных, языки и инструментальные средства параллельного программирования.

- Разработана методика управления научными приложениями ФВЭ и ЯФ для суперкомпьютеров с использованием информации о временно свободных ресурсах, повышающая эффективность использования суперкомпьютеров.

3) Модели, методы, алгоритмы и программная инфраструктура для организации глобально распределенной обработки данных.

- На основе методики, методов и архитектуры, разработанных в диссертации, создана глобальная распределенная система для обработки данных на основе динамического управления потоками заданий и динамическим распределением данных с учетом пропускной способности WAN. Реализация такой системы стала ключевым этапом для

дальнейшего развития компьютерной модели и сделала возможным создание гетерогенной киберинфраструктуры, позволив использовать ресурсы суперкомпьютеров и ресурсы «облачных вычислений» наряду с существующей инфраструктурой грид, нивелировав архитектурные различия вычислительных мощностей. Таким образом, разнородные вычислительные ресурсы доступны пользователям в виде единой киберинфраструктуры.

В целом, данная диссертационная работа производит очень хорошее впечатление, представляя масштабное и завершенное исследование, а приведенные мною ниже “замечания” носят, скорее, характер пожеланий, чем замечаний как таковых... И, тем не менее, по работе хотелось бы отметить следующее.

Как я уже писал выше, автор в процессе изложения приводит много примеров реального использования созданных средств в рамках больших физических экспериментов. Один из ключевых моментов, определяющих не только эффективность работы, но и саму возможность функционирования программных систем в столь масштабных условиях – это правильный выбор значений основных внутренних технологических параметров: размеров таблиц, буферов, значений таймаутов и других. В связи с этим хотелось бы в работе увидеть отдельный раздел, посвященный как описанию таких параметров (чтобы понимать, какие свойства действительно важны на практике), так и способу выбора правильных значений (на основе предварительного моделирования, эмпирически, динамический подбор в процессе работы или что-то иное). Неявно в работе это представлено в различных частях текста, но такой материал, собранный воедино в одном месте, безусловно, имеет самостоятельное значение и ценность.

Аналогично, очень хотелось бы видеть критический анализ предложенных подходов и методов, а также анализ функционирования

созданных на их основе программных средств. У автора огромный опыт организации и проведения крупномасштабных вычислительных экспериментов, в ходе которых наверняка становились понятными и “узкие места”, и источники ненадежности, и проблемы с масштабированием и многое другое, что всегда возникает в реальной работе распределенных программных комплексов. В чем основные проблемы? Как они устранились? Какие внутренние параметры программных комплексов наиболее критичны для нормальной работы систем? Как обрабатывались пиковые нагрузки? Какие режимы работы систем особенно важны? На что следует разработчикам подобных систем обращать внимание в будущем? Это не описание “отрицательных” или же “негативных” сторон данной работы – ни в коем случае – это описание степени эластичности предложенных подходов и их поведения в условиях реальной нагрузки.

Обсуждение подобных замечаний не меняет значимости и никак не влияет на оценку качества данной диссертационной работы, выполненной на высоком научном уровне.

На основании выполненных автором серьезных теоретических исследований в диссертации четко изложены научно обоснованные решения, показывающие, что работа является новым крупным достижением в теории и практике организации больших распределенных неоднородных компьютерных сред. Автореферат диссертации верно отражает основное содержание работы. Диссертационная работа А.А. Климентова “Методы обработки сверхбольших объемов данных в распределенной гетерогенной компьютерной среде для приложений в ядерной физике и физике высоких энергий” представляет собой завершенную научно-квалификационную работу по специальности 05.13.11, в которой на основании выполненных автором исследований разработаны теоретические положения, совокупность которых можно квалифицировать как научное достижение, удовлетворяющее критериям, установленным Положением о порядке присуждения ученых степеней, а ее автор, Климентов Алексей

Анатольевич, заслуживает присуждения ученой степени доктора физико-математических наук.

Официальный оппонент,
доктор физико-математических наук, член-корреспондент РАН, профессор,
заместитель директора Научно-исследовательского вычислительного центра
Московского государственного университета имени М. В. Ломоносова (НИВЦ
МГУ)

Воеводин Владимир Валентинович
20 февраля 2018 г.

119234, Москва, Ленинские горы, д. 1, стр. 4
Телефон 8 (495) 939-51-66, 8 (495) 939-54-24
Факс 8 (495) 938-21-36
E-mail: voevodin@parallel.ru

Подпись Вл.В.Воеводина удостоверяю
Директор НИВЦ МГУ, профессор



Тихонравов Александр Владимирович