

## **Отзыв официального оппонента**

на докторскую диссертацию Климентова Алексея Анатольевича “Методы обработки сверхбольших объемов данных в распределенной гетерогенной компьютерной среде для приложений в ядерной физике и физике высоких энергий”, представленной на соискание ученой степени доктора физико-математических наук по специальности 05.13.11 – математическое и программное обеспечение вычислительных машин, комплексов и компьютерных сетей

**Об авторе.** Алексей Анатольевич Климентов является высококвалифицированным специалистом в области сетевых технологий обработки данных физических экспериментов. Его работы получили большую известность в Европе, России и США. За последние почти двадцать лет им опубликовано несколько десятков научных работ, посвященных решению “открытых задач” обработки сверхбольших объемов данных на Большом адронном коллайдере (БАК) в Европейском центре ядерных исследований (ЦЕРН). Решение указанных выше задач имеет принципиальное значение с точки зрения доступности экспериментальных данных для физиков всего мира (в работе ЦЕРН принимают участие более трехсот крупнейших университетов и национальных лабораторий по всему миру).

А.А. Климентов является одним из немногих российских специалистов, разрабатывающих и успешно внедряющих в практику ЦЕРН системные решения, обеспечивающие возможности регистрации и обработки экспериментальных данных по мере увеличения мощности и светимости БАК. Наиболее значимые работы А.А. Климентова относятся к системе управления данными в экспериментах ATLAS и ALICE.

Совокупность системных исследований и разработок Климентова А.А. позволяет считать его одним из соавторов уникальной компьютерной сети грид, в которой реализованы базовые принципы обработки данных в научных мегапроектах не только по физике элементарных частиц, но и по другим тематическим направлениям. Например, в проекте по биоинформатике, выполняющемуся в Национальном исследовательском центре “Курчатовский институт”. Работы по этим проектам ведутся под

научным руководством А.А. Климентова (по гранту в соответствии с постановлением Правительства РФ № 220).

Докторская диссертация Климентова А.А. является не только обобщением ранее выполненных им работ, на практике доказавших их эффективность, но и принципиальные предложения автора по дальнейшему развитию мировой сети обработки данных физических экспериментов с использованием системы грид, вычислительных центров, оснащенных суперкомпьютерами, и ресурсов облачных систем. Подобные системы и объемы данных автор называет “системами эксабайтного масштаба” ( $10^{18}$ ).

Помимо собственных исследований и научного руководства международными проектами Алексей Анатольевич ведет большую научно-организационную работу, участвуя в программных и организационных комитетах международных конференций и школ для студентов, аспирантов и молодых ученых в России и за рубежом.

### **Общая характеристика диссертации.**

**Актуальность темы.** Развитие информационных технологий автор рассматривает не как самоцель, а как средства обеспечения фундаментальных исследований структуры материи. Поэтому актуальность темы диссертации автор обосновывает объективными требованиями современных исследовательских мегапроектов и, в частности, требованиями экспериментальных комплексов для исследований в области ядерной физики и физики высоких энергий, например, таких как БАК (ЦЕРН, Европа), NICA (ОИЯИ, Россия), электронный синхротрон (DESY, Германия).

На стр. 13 авторефера приведена таблица роста характеристик ускорителей за последние 60 лет и соответствующий ему рост требований к компьютерным системам обработки данных. Из таблицы видно, что современные требования по обрабатываемым объемам данных находятся в промежутке между сотнями петабайт и эксабайтами.

“Астрономические масштабы” экспериментальных данных в современных системах фундаментальных исследований в области физики высоких энергий являются важнейшим, но не единственным аргументом актуальности развития “инженерно-физического компьютеринга”. Вторым важным требованием, на которое указывает автор диссертации, является глобализация научных мегапроектов. Это означает, что интерфейсная часть сетевой технологии должна обеспечивать равные и быстрые возможности

доступа к экспериментальным данным всем участникам мегапроектов из любого региона в мире.

Этими объективными требованиями автор обосновывает актуальность выбранной им тематики исследований и разработок.

**Цель и задачи работы.** Целеполагание диссертации основано на реализации идеи использования для обработки данных не только существующую и успешно работающую систему грид (WLCG – World LHC Computing Grid), но и ресурсы суперкомпьютерных вычислительных центров, ресурсы облачных систем и университетских кластеров. Т.е. речь идет о создании глобальной сети обработки научных данных. Автор показывает, что для достижения этой цели необходимо учитывать не только производительность вычислительных узлов сети, но и пропускную способность каналов связи между ними. Для достижения поставленной цели автор провел полную классификацию типов задач, решаемых при планировании и реализации крупномасштабных экспериментов и задач обработки экспериментальных данных и доступа к ним профессиональных пользователей. Весь последующий текст диссертации по существу является ответом на вопрос: “Как эти задачи можно решать в архитектуре глобальной сети, оптимально используя все ресурсы по производительности вычислений, объемам памяти и пропускной способности каналов связи?” Далее в диссертации автор показывает, что создание глобальной сети для обработки данных по научным мегапроектам, является генеральным направлением современного инженерно-физического компьютеринга и приводит в списке литературы многочисленные ссылки на работы, проводимые в этом направлении в разных странах. При этом в диссертации центральное место занимают ссылки на работы ЦЕРН, ОИЯИ, НИЦ КИ, МГУ и НИЯУ МИФИ. Судя по публикациям, это организации, с которыми автор имел длительное и плодотворное научное-техническое сотрудничество.

Таким образом, объектом исследований автора являются глобальные гетерогенные компьютерные инфраструктуры для обработки данных научных мегапроектов.

Предметом исследований и разработок являются инструментальные программно-технические средства, позволяющие решать все комплексы задач математического моделирования и обработки данных научных мегапроектов, при оптимальном использовании ресурсов глобальной сети по производительности вычислений, объемам памяти и пропускной способности каналов связи узлов глобальной гетерогенной сети.

**Структура диссертации и логика изложения.** Диссертация состоит из введения, 4 глав, заключения, списка литературы из 115 наименований, полный объем работы составляет 238 страниц.

Введение содержит все положения, отвечающие требованиям ВАК к общей характеристику диссертации. Этот раздел работы является поучительным, т.к. автор обладает широкой эрудицией в области сетевых технологий и сумел кратко и точно описать развитие ситуации в этой научно-технической области за последние 60 лет ее бурного развития.

Логика изложения научного материала диссертации состоит в следующем. В общей характеристике работы представлены 5 новых крупных научных положений, выносимых на защиту. Затем в последующих главах приводится подробное описание этих результатов и доказательства их новизны, научной значимости и практической полезности. Таким образом, несмотря на большое количество сугубо специальных вопросов, которые рассматривает автор, работа сохраняет цельность и доступность для понимания. Этому также способствуют визуализация архитектурных и функциональных моделей и строгое определение всех понятий, используемых в работе.

### **Содержание диссертации по главам.**

**Глава 1.** Глава имеет аналитический характер. В ней анализируются достоинства и недостатки первой системы обработки данных на начальном этапе работы Большого адронного коллайдера. Система была построена по иерархическому принципу, имела 3 уровня обработки данных и получила название MONARC. При всех ограничениях системы MONARC, на которые указывает автор, она стала значительным шагом в развитии компьютеринга в области физики частиц. Более 200 центров в 60 странах мира вошли в консорциум WLCG, был получен первый опыт по распределенной обработке данных. В диссертации проанализированы принципиальные недостатки системы MONARC, связанные с иерархической архитектурой и отсутствием классификации данных, в частности, с выделением класса “популярных” данных (чаще используемых пользователями). Проведенный автором анализ позволил сформулировать новые принципы построения системы распределенной обработки данных в гетерогенной компьютерной сети, используя тип связи узлов “каждый с каждым” через ресурс WAN (глобальная компьютерная сеть). Таким образом была реализована “смешанная модель”, получившая название grid-mesh. В настоящее время грид является основной рабочей системой обработки данных в глобальных

научных мегапроектах. В России наиболее мощным и высоко загруженным узлом грид является узел, созданный в ОИЯИ, город Дубна, под руководством профессора Коренькова В.В. Полноценный узел грид создан также в НИЦ КИ. Оба центра имеют классификацию «центр первого уровня» (Tier-1) по классификации, принятой в консорциуме WLCG.

Вклад автора в создание мировой компьютерной сети грид состоит в следующем. Разработаны методы расширения понятия “вычислительный ресурс” за счет включения в его состав ресурса WAN. Это дало возможность отказаться от иерархической компьютерной модели и реализовать “смешанную модель” для распределенной обработки и управления данными в грид-среде. Разработаны методы определения популярности классов данных, а также разработана модель динамического управления данными в распределенной среде для сверхбольших объемов данных.

Несомненно, что выше указанный вклад автора в создание мировой компьютерной системы грид, является существенным.

**Глава 2.** Во второй главе описывается нормативный прогноз характеристик информационного обслуживания экспериментов на втором и последующих этапах работы БАК, этапы работы при увеличении энергии и светимости коллайдера. Анализируя очереди на выполнение заданий обработки и анализа данных при пиковых режимах работы коллайдера, автор приходит к выводу о недостатке мощности, предоставляемой в рамках WLCG. Разумный выход из этой ситуации автор видит в привлечении к обработке данных временно свободных ресурсов суперкомпьютерных вычислительных центров, облачных структур и университетских кластеров. Формулируются проблемы реализации, указанной выше идеи и предлагаются принципиальные методы их решения.

Существенный вклад автора в решение этих проблем состоит в следующем. Обоснованы принципы распределенной системы для обработки данных, которая могла бы работать с динамически изменяющимися вычислительными ресурсами и использовать мощности, доступные в течение относительно коротких временных интервалов. Показана необходимость расширения модели компьютинга и введения понятия ВЦ без “дискового пространства”, поскольку временно привлеченные к вычислениям ресурсы не предоставляют дисковое пространство для постоянного хранения данных (речь идет об объемах памяти порядка петабайт).

**Глава 3.** Третья глава посвящена подробному решению задач создания глобальной системы для обработки данных с временным привлечением

свободных вычислительных ресурсов. Большой объем представленного здесь материала не позволяет подробно рассмотреть и оценить все результаты проведенной автором теоретической работы в рамках отзыва на диссертацию. Поэтому упомянем только принципиально новые научно-технические предложения автора. Реализация идеи создания упомянутой глобальной системы с переменными вычислительными ресурсами, зависящими от времени, неизбежно привела автора к необходимости разработки динамической модели управления вычислительным процессом, в отличие от статической модели в системе грид. Это обстоятельство, в свою очередь, привело к необходимости иной структуризации трех основных типов потоков заданий: потоков заданий, выполняемых системой для распределенной обработки данных эксперимента; потоков заданий, выполняемых системой по требованию физических групп эксперимента; потоков заданий, выполняемых отдельными пользователями. При этом изменение ресурсов, выполняющих задания, должно быть “скрыто” от пользователя. Виртуальная машина, выделенная для выполнения заданий пользователя, должна оставаться для него неизменной. Решение указанных общесистемных задач не должно влиять на многочисленные прикладные инструменты обработки экспериментальных данных, уже накопленных пользователями в системе грид.

Обобщая материалы третьей главы, вклад автора в разработку глобальной системы обработки научных данных можно сформулировать следующим образом. Предложены новые принципы построения и методы реализации архитектуры глобальной системы для обработки данных в гетерогенной компьютерной среде, которые позволяют эффективно использовать вычислительные ресурсы и снимают противоречие по доступу к ресурсу между экспериментами, группами пользователей и отдельными учеными. Разработанная автором архитектура глобальной системы и методы управления вычислительным процессом в ней обеспечивают обработку данных в экспериментальном диапазоне, что является крупнейшим научно-техническим достижением.

**Глава 4.** Четвертая глава посвящена дальнейшему развитию компьютерной модели, интеграции суперкомпьютеров и ресурсов облачных вычислений с распределенными вычислительными ресурсами грид. В этой главе автор обобщает и анализирует все результаты своей работы по созданию глобальной системы обработки научных данных. При этом ставится вопрос о способности глобальной системы обеспечить обработку экспериментальных данных нового поколения в ближайшие 10 лет. Результаты анализа приводят автора к положительному ответу на данный вопрос.

Материалы заключительной главы диссертации показывают, что создание глобальной системы обработки научных данных относится не к отдаленному будущему, а к современному генеральному направлению инженерно-физического компьютеринга. Автором приводятся убедительные примеры эффективного решения задач обработки больших данных (Big Data) с использованием суперкомпьютеров и облачных структур. Создаваемые автором новые информационные технологии успешно используются не только в области физики высоких энергий, но и в других тематических направлениях, по которым работал автор, например, в задачах биоинформатики.

**Заключение.** Совокупным результатом диссертации является создание глобальной распределенной системы 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 GRID, нивелировав архитектурные различия вычислительных мощностей.

Подобный результат, безусловно, соответствует квалификации доктора физико-математических наук по специальности 05.13.11 – математическое и программное обеспечение вычислительных машин, комплексов и компьютерных сетей.

Теоретическое значение работы состоит в решении принципиально новых открытых задач в области информационных технологий, а высокая практическая значимость, полученных автором результатов, состоит в создании принципиально новой глобальной компьютерной сети для обработки данных фундаментальных исследований в области ядерной физики и физики высоких энергий.

В диссертации обобщены результаты 20-летней работы автора, отраженные в опубликованных свыше 150 печатных работах, в том числе в изданиях, индексируемых в базах данных Web of Science и Scopus.

Алексей Анатольевич Климентов является известным в мире специалистом в области компьютерных технологий обработки данных крупномасштабных физических экспериментов.

У оппонента нет существенных замечаний по работе.

Автореферат правильно и полно отражает содержание диссертации.

Диссертация соответствует критериям, установленным п. 9 Положения о присуждении учёных степеней (утверждено постановлением Правительства Российской Федерации от 24 сентября 2013 г. № 842) для учёной степени доктора наук.

Считаю, что Алексей Анатольевич Климентов заслуживает присуждения ученой степени доктора физико-математических наук по специальности 05.13.11 – математическое и программное обеспечение вычислительных машин, комплексов и компьютерных сетей за создание глобальной распределенной системы для обработки данных на основе динамического управления потоками заданий и динамическим распределением данных с учетом пропускной способности WAN.

Официальный оппонент, заведующий кафедрой анализ конкурентных систем, № 65, Федерального государственного автономного образовательного учреждения высшего образования Национальный исследовательский ядерный университет «МИФИ», профессор, доктор технических наук, лауреат Премии Правительства РФ по науке и технике

Оныкий Борис Николаевич

E-mail bonykij@mephi.ru  
Номер моб. телефона: 8 (903) 9730885  
Адрес: 115409, г. Москва, Каширское шоссе, д.31, НИЯУ МИФИ

Подпись удостоверяет

Заместитель начальника Управления по научной и технической документации и архивному делу  
документационного центра Национального исследовательского ядерного университета высшего образования  
МИФИ



*Марина Григорьевна Стасова*