

**ОТЗЫВ официального оппонента  
на диссертацию на соискание ученой степени  
кандидата физико-математических наук Шайхисламова Дениса  
Ильгизовича на тему: «Исследование и разработка методов для  
сравнительного анализа суперкомпьютерных приложений на основе  
технологий интеллектуального анализа данных»  
по специальности 2.3.5. Математическое и программное обеспечение  
вычислительных систем, комплексов и компьютерных сетей**

Суперкомпьютерные центры являются важным инструментом для решения фундаментальных научных и прикладных задач в таких областях, как физика, геология, метеорология, химия, биомедицина и многих других. Необходимость их применения обусловлена высокой вычислительной сложностью задач, решение которых на небольших вычислительных системах (и тем более, на отдельных серверах) потребовало бы большое количество времени — от нескольких суток до многих лет.

Однако добиться эффективной эксплуатации суперкомпьютерных центров зачастую очень непросто. Причинами являются, во-первых, чрезвычайно сложная архитектура самих суперкомпьютеров, а во-вторых, сложность разработки высокопроизводительных параллельных алгоритмов и программ. Для обеспечения высокой эффективности работы суперкомпьютера необходимо постоянно следить за его состоянием и рабочими характеристиками, чтобы иметь возможность оперативно отслеживать и устранять негативно влияющие факторы. Одним из важных аспектов такого мониторинга является анализ потока запускаемых пользовательских приложений, поскольку зачастую как сами пользователи, так и системные администраторы не обладают достаточной информацией о наличии проблем с производительностью в выполняемых заданиях.

В диссертационной работе Шайхисламова Д.И. представлен один из способов анализа потока суперкомпьютерных приложений, основанный на выявлении схожих приложений.

Основной идеей применения данных методов является использование доступной информации об уже отработавших приложениях для анализа новых запускаемых приложений. Автором использовалось предположение о том, что схожие (в некоторой метрике) приложения зачастую обладают одними и теми же свойствами и особенностями, корректность которого подтверждается на практике успешным применением методов выделения схожих приложений для решения задач администрирования. Успешное применение предложенных методов на практике для получения новой информации о функционировании суперкомпьютерных систем подтверждает их практическую значимость.

Введение диссертации раскрывает актуальность исследования, связанного с созданием новых подходов для выявления схожих суперкомпьютерных приложений. Сформулированы цели и задачи исследования.

В главе 1 соискатель провел систематический анализ существующих подходов к решению задачи выделения схожих суперкомпьютерных приложений. Особое внимание уделяется доступным данным для анализа суперкомпьютерных приложений, в связи с чем методы разделены на две группы: с использованием статических данных (доступных до запуска приложения) или динамических данных (собираемых во время выполнения приложения).

В главе 2 описан метод выделения схожих приложений на основе статических данных. Предлагается использовать имена функций и переменных, которые извлекаются из исполняемых файлов приложения. Исходные данные преобразуются в вектор имен с помощью модели Doc2Vec, что позволяет их в дальнейшем сравнивать с помощью косинусного сходства и таким образом оценивать степень сходства приложений. Данный метод был протестирован на вручную размеченных данных, на которых скорректированный индекс Рэнда составил 0.79, что показывает хорошее качество работы метода.

В главе 3 описан метод выделения схожих приложений на основе динамических данных. Предлагается использовать данные от системы мониторинга, собираемые во время работы приложения, что позволяет получить временные ряды для различных характеристик работы приложения, как, например, загрузки процессоров или графических ускорителей, нагрузки на сеть Infiniband и т.д. Для оценки схожести приложений предлагается сравнивать полученные многомерные ряды с помощью метода Dynamic Time Warping, который является широко известным методом для сравнения временных рядов с локальными и глобальными временными сдвигами. Тестирование метода проводилось на вручную размеченных данных, точность на которой составила 0.85, что является высоким показателем для задачи такого типа.

В главе 4 описаны различные сценарии применения на практике предложенных в главах 2 и 3 методов. Сценарии применения представляют собой важные для администраторов задачи, которые могут использоваться как для сбора статистики о работе суперкомпьютера, так и для упрощения процесса анализа большого потока запускаемых приложений.

Первый сценарий применения – использование предложенных методов для автоматического определения того, какие программные пакеты задействуются в пользовательских приложениях. Было показано, что эту задачу можно решить с помощью статического метода, используя существующую информацию об использовании программных пакетов ранее изученными заданиями. Динамический метод также может быть применен для этой цели, но качество работы не сильно отличается от статического метода, при этом динамический метод вычислительно затратнее. При апробации метода на суперкомпьютере «Ломоносов-2» было обнаружено на 15% больше заданий, чем с помощью существующего решения.

Во втором сценарии применения был показан способ кластеризации приложений по динамике их работы. Для проведения кластеризации предлагается совместно использовать статический и динамический методы:

сначала выполняется первичная кластеризация с использованием статистического метода, затем каждый полученный кластер разбивается на более мелкие с помощью кластеризации на основе динамического метода. Такой подход существенно ускоряет процесс кластеризации. Подход был протестирован на вручную размеченных данных, где скорректированный индекс Рэнда составил 0.8, что показывает хорошее качество работы. В разделе был показан как оффлайн, так и онлайн алгоритмы кластеризации. Оффлайн алгоритм предполагает кластеризацию всех заданий «с нуля», вне зависимости от предыдущих результатов кластеризации. Онлайн алгоритму же не требуется проводить полную кластеризацию заданий при появлении новых заданий, и прошлые результаты кластеризации только уточняются с учетом этих новых заданий. При апробации оффлайн и онлайн алгоритмов было показано, что результаты их работы практически идентичны – скорректированный индекс Рэнда составил 0.99.

В третьем сценарии применения был продемонстрирован способ предсказания оценок качества использования суперкомпьютерных ресурсов (сами оценки были предложены ранее вне рамок данной работы). Одним из особенностей данных оценок является необходимость сбора дополнительных метрик производительности, которые могут быть не всегда доступны. Отсутствие данных метрик делает расчет оценок качества использования ресурсов суперкомпьютера невозможным, и для восполнения этого пробела был предложен способ предсказания оценок, которому посвящен данный раздел. Было показано, что, используя собранные ранее исторические данные об оценках качества использования ресурсов, можно предсказать данные оценки для новых заданий с помощью методов выделения схожих приложений. В частности, метод позволяет предсказать значения оценок для 89% заданий, со средней абсолютной ошибкой, равной 2.11, что является допустимой погрешностью в данном случае.

**Научная новизна** исследования состоит в том, что автор:

- предложил два новых метода выделения схожих суперкомпьютерных приложений. В первом методе используются данные, получаемые из исполняемых файлов, во втором методе – данные системы мониторинга. Оба метода показали хорошее качество работы;
- предложены различные подходы к анализу потока заданий, выполняемых на суперкомпьютере, на основе вышеуказанных методов выделения схожих приложений: выявление использования программных пакетов, кластеризация потока приложений, а также предсказание оценок качества использования вычислительных ресурсов. Все методы прошли апробацию на реальных данных с суперкомпьютера «Ломоносов-2».

**Достоверность** полученных результатов подтверждается комплексной проверкой: для всех используемых методов проводилась оценка корректности и детальная апробация на реальных данных суперкомпьютерного комплекса как на вручную размеченных данных, так и в сравнении с существующими подходами (например, сравнение с XALT в задаче выявления использования программных пакетов).

По диссертационной работе Д.И. Шайхисламова имеются следующие замечания, которые, однако, не снижают общего высокого научного уровня работы и значимости полученных результатов:

- Задача предсказания оценок качества использования суперкомпьютерных ресурсов относится к разработанному и используемому только в МГУ программному решению; остается открытым вопрос возможности применения этого решения в иных суперкомпьютерных центрах.
- Разработанный метод выделения схожих приложений на основе статических данных использует модель Doc2Vec для оценки сходства приложений. При этом 1) не рассмотрены возможные альтернативы этой

модели 2) обучение модели проводилось только на данных с суперкомпьютера «Ломоносов».

- В работе утверждается, что предлагаемые методы направлены на повышение качества работы суперкомпьютеров, но при этом не приведена количественная оценка улучшения эффективности работы суперкомпьютеров после применения предложенных методов.
- В тексте диссертации не всегда даны точные формулировки используемых понятий. Например, на стр. 37 используется понятие «сплайна» без какой-либо конкретизации его вида (квадратичный, кубический сплайн, и т.д.); на стр. 80 используются неформальные понятия «онлайн» и «оффлайн» метод кластеризации.
- Присутствуют неточности в оформлении: в диссертации – использование прямого шрифта (вместо курсива) и русских слов (вместо кратких обозначений) в формулах; в автореферате – некорректные ссылки на таблицы и рисунки.

### **Заключение**

Диссертация отвечает требованиям, установленным Московским государственным университетом имени М.В.Ломоносова к работам подобного рода. Содержание диссертации соответствует специальности 2.3.5. Математическое и программное обеспечение вычислительных систем, комплексов и компьютерных сетей (по физико-математическим наукам), а также критериям, определенным пп. 2.1-2.5 Положения о присуждении ученых степеней в Московском государственном университете имени М.В.Ломоносова. Диссертационное исследование оформлено согласно требованиям Положения о совете по защите диссертаций на соискание ученой степени кандидата наук, на соискание ученой степени доктора наук Московского государственного университета имени М.В.Ломоносова.

Таким образом, соискатель Шайхисламов Денис Ильгизович заслуживает присуждения ученой степени кандидата физико-математических наук по специальности 2.3.5. Математическое и программное обеспечение вычислительных систем, комплексов и компьютерных сетей.

Официальный оппонент:

доктор технических наук, доцент,  
заведующий кафедрой математического обеспечения и суперкомпьютерных технологий института информационных технологий, математики и механики Федерального государственного автономного образовательного учреждения высшего образования «Национальный исследовательский Нижегородский государственный университет им. Н.И. Лобачевского»

Баркалов Константин Александрович

Дата: 24.02.2026

Контактные данные:

тел.: 7(831)4623356, e-mail: barkalov@vmk.unn.ru  
Специальность, по которой официальным оппонентом защищена диссертация: 05.13.01 – Системный анализ, управление и обработка информации (в науке и промышленности)

Адрес места работы:

603022, г. Н. Новгород, пр. Гагарина, д. 23, корп. 2,  
Федеральное государственное автономное образовательное учреждение высшего образования «Национальный исследовательский Нижегородский государственный университет им. Н.И. Лобачевского», институт информационных технологий, математики и механики, кафедра математического обеспечения и суперкомпьютерных технологий.  
Тел.: 7(831)4623356; e-mail: barkalov@vmk.unn.ru