

На правах рукописи

Ровнягин Михаил Михайлович

**МЕТОДЫ И СРЕДСТВА РЕШЕНИЯ ЗАДАЧ ПОИСКА  
И ЗАЩИЩЕННОГО ХРАНЕНИЯ ДАННЫХ  
С ПРИМЕНЕНИЕМ ГИБРИДНЫХ ВЫЧИСЛИТЕЛЬНЫХ  
ТЕХНОЛОГИЙ**

05.13.11 – математическое и программное обеспечение  
вычислительных машин, комплексов и компьютерных сетей

05.13.19 – методы и системы защиты информации,  
информационная безопасность

**АВТОРЕФЕРАТ**

диссертации на соискание ученой степени кандидата  
технических наук

Автор: 

Москва – 2015

Работа выполнена в Национальном исследовательском  
ядерном университете «МИФИ»

Научный руководитель: кандидат технических наук, доцент  
ВАСИЛЬЕВ Николай Петрович

Официальные оппоненты: доктор технических наук, профессор  
ГРИБУНИН Вадим Геннадьевич,  
МОУ «Институт инженерной физики»,  
главный научный сотрудник

кандидат технических наук  
ГЕРЦЕНБЕРГЕР Константин Викторович,  
Объединенный институт  
ядерных исследований (ОИЯИ),  
научный сотрудник

Ведущая организация: Учреждение Российской академии наук  
Институт программных систем  
имени А.К. Айламазяна РАН  
(ИПС им. А.К. Айламазяна РАН)

Защита диссертации состоится «16» декабря 2015 г. в 16 часов  
30 минут на заседании диссертационного совета Д 212.130.08 на базе  
Национального исследовательского ядерного университета «МИФИ»  
по адресу: 115409, г. Москва, Каширское шоссе, 31.

С диссертацией можно ознакомиться в библиотеке НИЯУ  
МИФИ и на сайте <http://ods.mephi.ru>.

Отзывы на автореферат в двух экземплярах, заверенные  
печатью, просьба направлять по адресу: 115409, г. Москва, Каширское  
шоссе, 31, диссертационные советы НИЯУ МИФИ (тел. +7 (499) 324-  
84-98).

Автореферат разослан «\_\_\_» \_\_\_\_\_ 2015 г.

Ученый секретарь

диссертационного совета



Горбатов В.С.

## ОБЩАЯ ХАРАКТЕРИСТИКА РАБОТЫ

Актуальность темы исследований. Поиск данных является одной из важнейших проблем в области информатики и вычислительной техники. Как пишет Д. Кнут в своем фундаментальном труде «Искусство программирования»: «Поиск обычно является наиболее «времяемкой» частью многих программ, и замена плохого метода поиска хорошим может значительно увеличить скорость работы программы».

Производительность работы любых систем поиска данных напрямую зависит от объема самих данных, а в настоящее время происходит бурный, лавинообразный рост этих объемов. В силу практически полной компьютеризации всех сфер человеческой деятельности – от научной и до развлекательной – генерация данных, в которых необходимо выполнять поиск, происходит повсеместно и очень высокими темпами. За последние два года человечество накопило больше данных, чем за всю предыдущую историю. Количество пользователей Интернет и других компьютерных сетей к 2016 году втрое превысит население Земли. Сбербанк России за сутки осуществляет более 10 млн транзакций. В социальной сети Twitter каждый месяц выполняется более 30 млрд операций поиска.

В последние годы для обозначения больших и постоянно растущих объемов данных даже возник специальный термин: «Big Data (Большие Данные)». Четкого определения данного термина нет; под этим понимают такие объемы данных, при обработке которых может произойти переход количества в качество.

Классические реляционные СУБД (РСУБД), малоприспособлены в качестве систем хранения Больших Данных. Поэтому большое распространение в настоящее время получила концепция NoSQL (not only SQL) систем, в которых данные хранятся в формате «ключ-значение». Подобные системы хоть и не обладают всей полнотой функциональности РСУБД, позволяют существенно увеличить производительность работы приложений с большими объемами слабоструктурированных данных. Общей чертой подобных систем является высокий показатель пропускной способности, масштабируемый (в большинстве случаев линейно) в зависимости от количества используемых серверов хранения.

Множество компаний по всему миру используют в настоящее время распределенные высокопроизводительные NoSQL-решения. Сервисы таких корпораций как IBM, Amazon, Facebook, Twitter, Google, Oracle и др. напрямую зависят от того, насколько эффективно функционируют системы хранения и поиска данных. Среди важнейших критериев эффективности NoSQL-решений выделяются следующие показатели: пропускная способность, масштабируемость, энергоэффективность, затраты на обслуживание. Как можно заметить, данные показатели перекликаются с требованиями, предъявляемыми к современным суперкомпьютерам.

В настоящее время ведущие позиции в суперкомпьютерных рейтингах занимают гибридные решения, когда в состав вычислительной системы (ВС) входят как обычные CPU, так и специализированные вычислители или дискретные сопроцессоры (GPU, MIC, ПЛИС и др.). Данные устройства позволяют существенно ускорить ряд операций и повысить энергоэффективность ВС. С точки зрения задач хранения и поиска данных, сопроцессоры могут быть полезными для ускорения операций поиска, проверки членства в множестве, выполнения операций дополнительной обработки (шифрование, обработка мультимедиа и др.).

Проблемой поиска информации, в том числе и задачами высокопроизводительного поиска информации в специализированных структурах данных, занимались такие известные ученые, как Воеводин В.В., Т. Кормен, Р. Седжвик и др. Таким образом, представляется целесообразным применить ряд суперкомпьютерных технологий, предназначенных для ускорения вычислений, именно к задаче поиска.

Существующие на данный момент NoSQL-решения не предназначены для использования в гибридных ВС, так как все операции выполняются только центральными процессорами. Кроме того, в них отсутствуют встроенные механизмы защиты хранимых данных от несанкционированного доступа (НСД), в силу чего похищение носителя с данными, например, жесткого диска с сервера хранения, неизбежно приведет к утечке информации. Таким образом, разработка методов и средств решения задач поиска и защищенного хранения данных с применением гибридных вычислительных технологий *является актуальной научной и инженерной задачей.*

**Объектом исследования** являются методы и средства решения задач поиска и защищенного хранения данных, использующиеся в высокопроизводительных распределенных системах хранения данных.

**Предметом исследования** являются методы и программно-аппаратные средства повышения производительности операций поиска, хранения и дополнительной обработки данных (прежде всего, шифрования).

**Целью диссертационной работы** является повышение производительности систем хранения и поиска данных за счет использования гибридных суперкомпьютерных технологий.

**Методы исследования.** В настоящей работе используются методы теории вероятностей и математической статистики, методы теории массового обслуживания, элементы теории алгоритмов, методы проектирования структурных и функциональных моделей распределенных систем (в нотациях IDEF и UML), метод анализа иерархий, методы теории планирования эксперимента.

**Научная новизна.**

1. Предложена классификация современных высокопроизводительных систем хранения и поиска данных, на основе предложенной классификации выделены подсистемы, работу которых можно существенно ускорить, применяя гибридные суперкомпьютерные технологии.
2. Впервые построена математическая модель распределенной системы хранения и поиска данных, использующей гибридные суперкомпьютерные технологии.
3. Впервые предложены структурные и функциональные модели системы хранения и поиска данных, использующей гибридные суперкомпьютерные технологии.
4. Разработаны новые методы повышения производительности поиска данных, основывающиеся на применении новых алгоритмов поиска данных, предназначенных для использования в гибридных вычислительных системах.
5. Предложены новые методы организации дополнительной обработки данных в высокопроизводительных системах хранения и поиска данных с использованием гибридных суперкомпьютерных технологий.

### **Обоснованность и достоверность результатов работы**

подтверждается доказанностью выводов, полученных в результате корректного применения математического аппарата, апробацией основных результатов на отечественных и зарубежных конференциях, публикациями, успешной реализацией предложенных методов и средств в виде прототипа распределенной системы хранения и поиска данных, внедрением результатов в учебный и производственный процесс.

### **Практическая значимость работы:**

1. Для описания структурных и функциональных моделей систем, построенных на основе предлагаемых методов, а также представленных классификаций, используются стандартные нотации IDEF и UML. Подобная унификация упрощает процесс проектирования систем на основе предлагаемых методов.
2. Создан прототип распределенной системы хранения и поиска данных на основе разработанных методов. Экспериментальные исследования прототипа подтверждают эффективность его работы, что делает целесообразным построение промышленных систем хранения данных на базе предложенного подхода.
3. Разработан программный каркас, упрощающий процесс разработки подобных систем, позволяя оснащать систему хранения программными модулями для специализированной обработки информации.
4. Разработанная методика оценки эффективности гибридной распределенной системы хранения и поиска данных может быть использована для анализа других подобных решений.

### **Реализация результатов исследования.**

Результаты исследований использовались при выполнении НИР по анализу защищенности ВК Эльбрус, выполненной по заказу ЗАО МЦСТ и НИР, выполненной в рамках ФЦП «Кадры». Разработанная система и ее программные модули внедрены в учебный процесс кафедры «Компьютерные системы и технологии» НИЯУ МИФИ и производственный процесс ООО «Голден Интернет».

**Апробация.** Основные результаты исследований, проводимых в рамках данной работы, представлены на следующих конференциях: 1-й, 2-й и 3-й Национальный Суперкомпьютерный Форум (г. Переславль-Залесский, 2012-2014), 16-я, 17-я, и 18-я

Международная телекоммуникационная конференция молодых ученых и студентов «Молодежь и наука» (г. Москва, 2013-2015), International Conference «The Radio-Electronic Devices and Systems for The Infocommunication Technologies» RES-2013 (Moscow, 2013), Научная сессия НИЯУ МИФИ 2012, 2014, 2015 (г. Москва, 2014-2015), 2015 IEEE NW Russia Young Researchers in Electrical and Electronic Engineering Conference ElConRusNW-2015, (St. Petersburg, 2015).

**Публикация результатов.** По теме диссертации автор имеет 18 печатных работ в том числе: 4 статьи в журналах и сборниках трудов из списка рекомендованных Высшей аттестационной комиссией, 5 работ, индексируемых базой Scopus, одна из которых индексируется также Web of Science.

**Основными положениями, выносимыми на защиту, являются:**

1. Две классификации: современных высокопроизводительных систем хранения и поиска данных и гибридных вычислительных систем.
2. Математическая модель высокопроизводительной распределенной системы хранения и поиска данных, ориентированной на использование в гибридных вычислительных системах.
3. Структурная и функциональная модели распределенной системы хранения и поиска данных.
4. Методы повышения производительности операций поиска и хранения данных в гибридных вычислительных системах.
5. Способ стохастического преобразования данных, ориентированный на реализацию в гибридных вычислительных системах.
6. Результаты исследований эффективности разработанной системы распределенного хранения и поиска данных.

**Структура и объем диссертационной работы.**

Диссертационная работа состоит из введения, четырех глав, заключения, списка литературы и приложений. Общий объем работы составляет 181 страницу (без учета приложений). Работа содержит 61 иллюстрацию, 24 таблицы. Список литературы состоит из 145 наименований. Диссертация соответствует паспорту специальности

05.13.11 по пунктам 3, 4, 8, 9 областей исследования и специальности  
05.13.19 по пунктам 2 и 13.

### ОСНОВНОЕ СОДЕРЖАНИЕ РАБОТЫ

Во введении обоснована актуальность темы, определены цель, задачи и методы исследования, установлены основные положения диссертационной работы, выносимые на защиту.

В первой главе предложена классификация (рисунок 1) современных высокопроизводительных систем поиска и хранения данных (ВСПиХД).

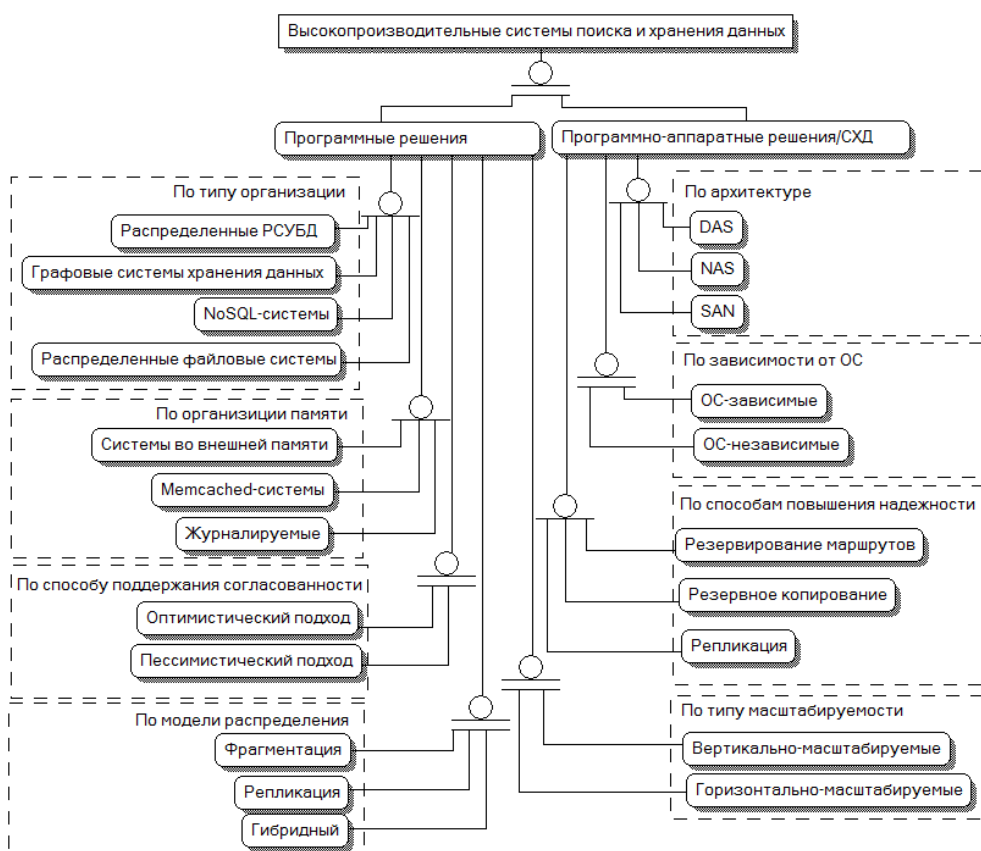


Рисунок. 1. Классификация высокопроизводительных систем поиска и хранения данных



Описаны основные проблемы повышения производительности операций поиска в подобных системах. Сделан вывод о возможности использования гибридных вычислительных технологий в задачах поиска данных и предложена классификация современных гибридных вычислительных систем (рисунок 2).

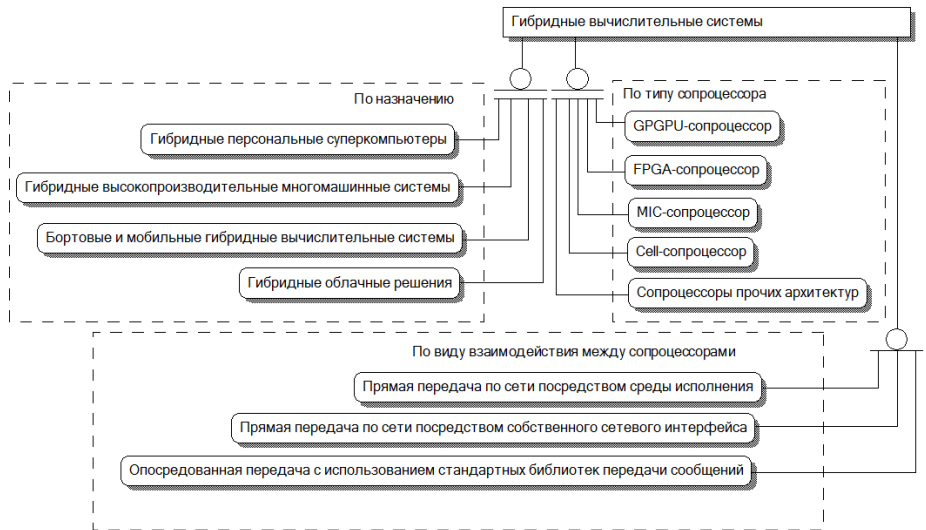


Рисунок 2. Классификация гибридных вычислительных систем

В результате проведенного исследования установлено, что ни одна из существующих систем не поддерживает гибридных вычислительных технологий. Существующие алгоритмы поиска данных с применением гибридных сопроцессоров сильно связаны с архитектурой последних и не являются универсальными.

**Вторая глава** посвящена разработке методов повышения производительности поиска данных в гибридных системах. На основе выполненного в Главе 1 исследования построена обобщенная схема ВСПиХД, состоящей из балансировщика,  $K$  серверов доступа и  $M$  узлов хранения (рисунок 3).

Процесс обработки клиентских поисковых запросов во ВСПиХД состоит из следующих этапов:

- 1) Первоначальный запрос клиента к системе.

- 2) Переадресация запроса на сервер доступа, с балансированием нагрузки.
- 3) Выполнение запросов к данным на узлах хранения. Количество затрагиваемых узлов хранения определяется схемой отображения, установленной на сервере доступа (хеш-кольцо).
- 4) Выполнение запросов к данным во внешнюю память.
- 5) Формирование ответа: передача найденных данных, либо уведомления об их отсутствии.
- 6) Отправка итогового ответа клиенту.

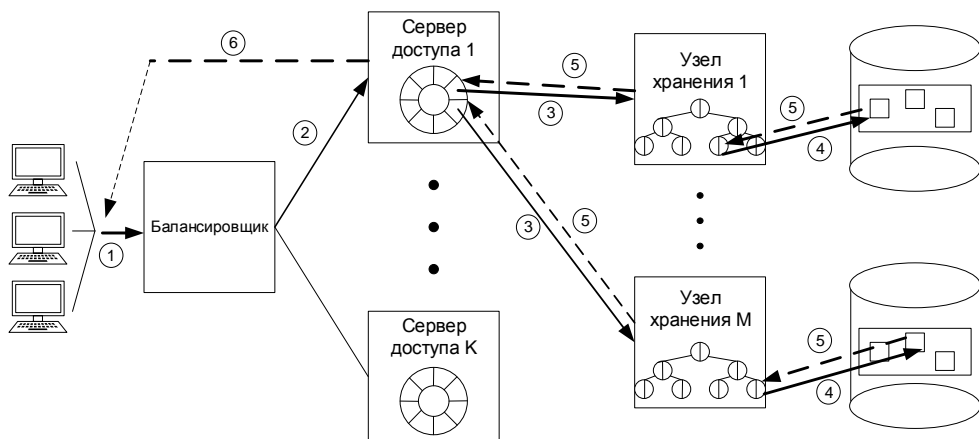


Рисунок 3. Структура современной распределенной ВСПиХД

По аналогичной схеме выполняются операции добавления и удаления данных. В случае добавления нового элемента в систему хранения клиенту возвращается уникальный ключ (идентификатор) сохраненного элемента. При удалении – код возврата операции удаления. Операции обновления реализуются клиентским приложением.

Непосредственно сами данные содержатся на узлах хранения, а соответствующие им метаданные хранятся на серверах доступа. Метаданные могут использоваться для фильтрации запросов, поддержания согласованности и пр.

ВСПиХД как сеть массового обслуживания (СеМО), изображена на рисунке 4.

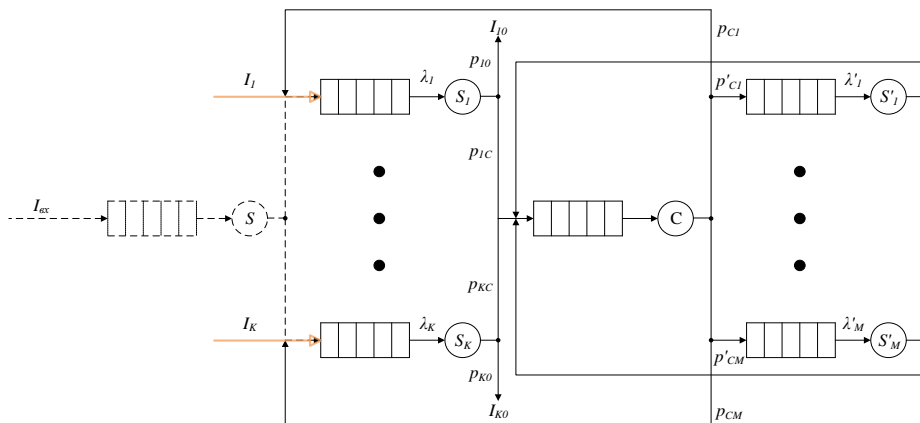


Рисунок 4. ВСПиХД, представленная в виде СеМО

Обозначения:  $S$ ,  $C$ ,  $S_i$  и  $S_i'$  – балансировщик, сеть межсоединений, серверы доступа и узлы хранения, представленные в виде систем массового обслуживания (СМО);  $I_{ek}$ ,  $I_i$  – интенсивности поступления входных потоков заявок в СеМО,  $\lambda_i$ ,  $\lambda_i'$  – интенсивности поступления заявок на СМО,  $p_{ij}$ ,  $p'_{ij}$  – вероятности передачи заявок между СМО.

Внутренняя структура СМО  $S_i$  и  $S_i'$  представлены в виде СеМО на рисунках 5 и 6 соответственно.

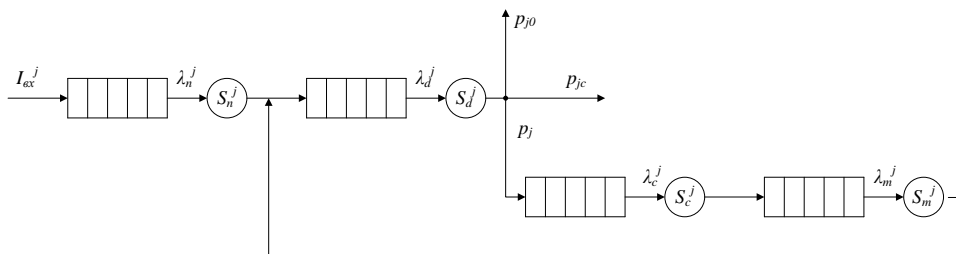


Рисунок 5. Сервер доступа ВСПиХД, представленный в виде СеМО

На рисунках 5-6 введены обозначения подсистем:  $S_n^a$ ,  $S_n^j$  – сетевая подсистема,  $S_d^a$ ,  $S_d^j$  – подсистема управления сетевой активностью,  $S_c^j$

-подсистема обеспечения согласованности,  $S_m^j$  – подсистема разметки пространства ключей,  $S_c^a$  – подсистема классификации запросов,  $S_n^a$  – подсистема предварительной обработки запросов,  $S_d^a$  – подсистема хранения. Вероятности продолжения обработки заявок в СеМО сервера доступа и узла хранения обозначены  $p_j$  и  $p_a$  соответственно.

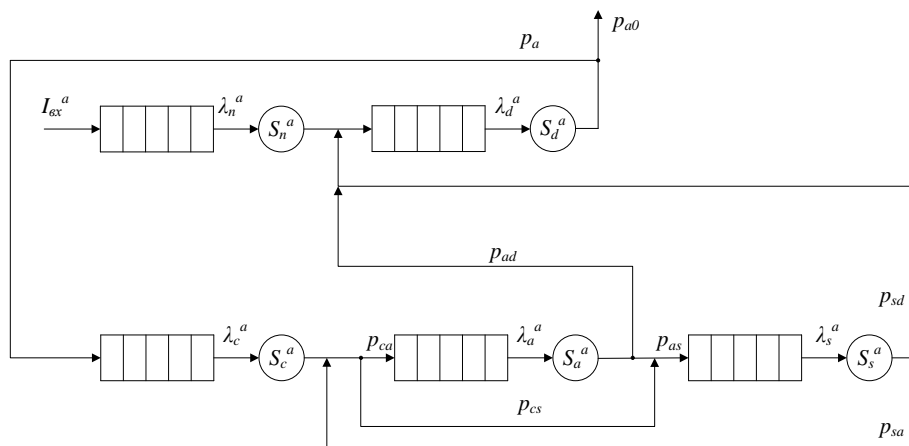


Рисунок 6. Узел хранения ВСПиХД, представленный в виде СеМО

Для прогнозирования возможного увеличения производительности, получаемого за счет использования гибридных сопроцессоров для ускорения операций поиска и дополнительной обработки данных (шифрования), выведено следствие закона Амдала для гибридных систем, представленное формулой (1):

$$A = \frac{1}{(1-\alpha) + c\left(\frac{\alpha(1-\beta)}{n} + q\right) + \frac{\alpha\beta}{p^{COP}}}, \quad (1)$$

где  $A$  – ускорение;  $\alpha$  – коэффициент параллелизма (доля распараллеливаемой части программы ко всей программе);  $n$  – количество параллельных процессов,  $\beta$  – коэффициент использования сопроцессора (доля части кода, распараллеливаемого на сопроцессоре ко всему распараллеливаемому коду);  $p^{COP}$  – количество ядер дискретного

сопроцессора (индекс COP обозначает тип сопроцессора: GPU/FPGA/MIC и др.);  $c$  – коэффициент загрузки ядер CPU,  $q$  – отношение времени смены контекста CPU ко всему времени работы программы.

Таким образом, работу прототипа ВСПиХД, использующего гибридные суперкомпьютерные технологии для ускорения операций поиска и дополнительной обработки данных, можно смоделировать системой уравнений (2):

$$\left\{ \begin{array}{l} \bar{T}_{\text{преб}} = \frac{\sum_{j=1}^K \lambda_j \bar{T}_{\text{пр}j} + \sum_{a=1}^M \lambda_a \bar{T}_{\text{пр}a} + \lambda_c \bar{T}_{\text{пр}c}}{I_{\text{ВХ}}} \\ \bar{T}_{\text{пр}j} = \frac{\lambda_n^j \frac{\bar{T}_{\text{обс}n}}{A_n - \lambda_n^j \bar{T}_{\text{обс}n}} + \lambda_d^j \frac{\bar{T}_{\text{обс}d}}{A_d - \lambda_d^j \bar{T}_{\text{обс}d}} + \lambda_c^j \frac{\bar{T}_{\text{обс}c}}{A_c - \lambda_c^j \bar{T}_{\text{обс}c}} + \lambda_m^j \frac{\bar{T}_{\text{обс}m}}{A_m - \lambda_m^j \bar{T}_{\text{обс}m}}}{I_{\text{ВХ}}^j}, \quad (2) \\ \bar{T}_{\text{пр}a} = \frac{\lambda_n^a \frac{\bar{T}_{\text{обс}n}}{A_n - \lambda_n^a \bar{T}_{\text{обс}n}} + \lambda_d^a \frac{\bar{T}_{\text{обс}d}}{A_d - \lambda_d^a \bar{T}_{\text{обс}d}} + \lambda_c^a \frac{\bar{T}_{\text{обс}c}}{A_c - \lambda_c^a \bar{T}_{\text{обс}c}} + \lambda_a^a \frac{\bar{T}_{\text{обс}a}}{A_a - \lambda_a^a \bar{T}_{\text{обс}a}} + \lambda_s^a \frac{\bar{T}_{\text{обс}s}}{A_s - \lambda_s^a \bar{T}_{\text{обс}s}}}{I_{\text{ВХ}}^a} \end{array} \right.$$

где  $\bar{T}_{\text{преб}}$  выражает время нахождения заявки во ВСПиХД,  $\bar{T}_{\text{пр}j}$  и  $\bar{T}_{\text{пр}a}$  – время пребывания заявки на сервере доступа и узле хранения соответственно, а  $\bar{T}_{\text{обс}}$  – время обслуживания на конкретных СМО, а индексы обозначают соответствующие подсистемы.

В настоящей работе рассматривается способ увеличения производительности операций поиска за счет уменьшения интенсивности поступающих в СеМО заявок. Достигается это благодаря использованию модифицированного фильтра Блума для гибридных систем [14].

Для оценки нижней границы вероятности ложноположительного срабатывания модифицированного фильтра Блума со счетчиками в настоящей работе предложена формула (3):

$$p_{\text{mod}} = \left( 1 - \left( 1 - \frac{L}{m} \right)^{kn} \right)^k, \quad (3)$$

где  $L$  – разрядность счетчика,  $m$  – размер вектора Блума,  $k$  – количество хеш-функций,  $n$  – число элементов хранения. В результате добавления фильтра Блума на сервер доступа изменится вероятность продолжения обработки заявки на сервере доступа  $p_j$ .

Новое значение вероятности определяется следующим соотношением:

$$p_j^{new} = p_j - p_{bn} = p_j - (1 - (p_{mod} + p_{bp} - p_{mod}p_{bp})), \quad (4)$$

где  $p_j^{new}$  – новое значение вероятности продолжения обработки заявки на сервере доступа  $p_{bn}$  – вероятность отрицательного срабатывания фильтра Блума,  $p_{bp}$  – вероятность положительного срабатывания фильтра Блума,  $p_{mod}$  – вероятность ложноположительного срабатывания фильтра Блума. Значения первых двух вероятностей зависят от распределения запросов по ключам в реальной системе, а  $p_{mod}$  определяется соотношением (3).

В результате добавления модифицированного фильтра Блума на серверы доступа работа последних будет описываться уравнениями баланса (5), где  $p_j^{new}$  определяется соотношением (4).

$$\left\{ \begin{array}{l} \lambda_n^j = I_{BX}^j \\ \lambda_d^j = \frac{I_{BX}^j}{p_{jo} + p_{jc}} = \frac{I_{BX}^j}{1 - p_j^{new}} \\ \lambda_c^j = \frac{p_j^{new} I_{BX}^j}{1 - p_j^{new}} \\ \lambda_m^j = \frac{p_j^{new} I_{BX}^j}{1 - p_j^{new}} \end{array} \right. \quad (5)$$

Для ускорения операций поиска данных в работе предложена новая архитектурно-независимая структура данных – стохастическое дерево поиска [8]. В основе принципов работы этой структуры данных лежит тот факт, что математическое

ожидание высоты случайного бинарного дерева поиска с  $n$  ключами равно  $O(\log_2 n)$ . Псевдослучайного характера добавляемых ключей, а значит и сбалансированности двоичного дерева поиска, можно добиться, на основе использования стохастического преобразования поступающих ключей, как показано на рисунке 7.

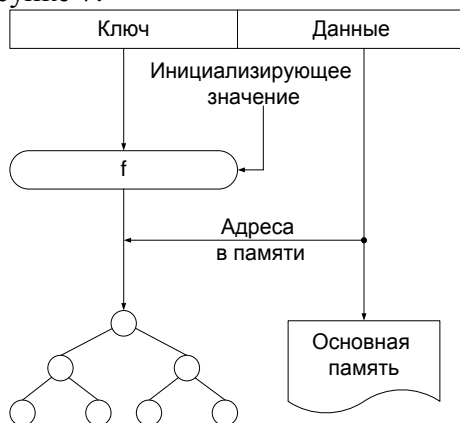


Рисунок 7. Стохастическое бинарное дерево поиска

В качестве функции стохастического преобразования  $f$  могут применяться как классические шифры (в настоящей работе реализована и испытана параллельная версия ГОСТ 28147-89), так и специальные алгоритмы, предназначенные для использования в гибридных суперкомпьютерных системах (в настоящей работе исследовался алгоритм DOZEN [11]). Алгоритм DOZEN – одна из новейших запатентованных российских разработок. Блочный алгоритм DOZEN обладает высокой степенью параллелизма: вся поступающая информация делится на блоки по 512 бит, представляющие собой кубы  $4 \times 4 \times 4$  байт, где каждый полученный куб преобразуется по «слоям» в трех направлениях:  $x$ ,  $y$  и  $z$  (рисунок 8).

В качестве алгоритма стохастического преобразования ключей в работе использовалась параллельная CUDA-версия DOZEN, представленная на рисунке 8.

Параллельная реализация блочного алгоритма шифрования ГОСТ 28147-89 применялась для организации прозрачного шифрования данных (TDE).

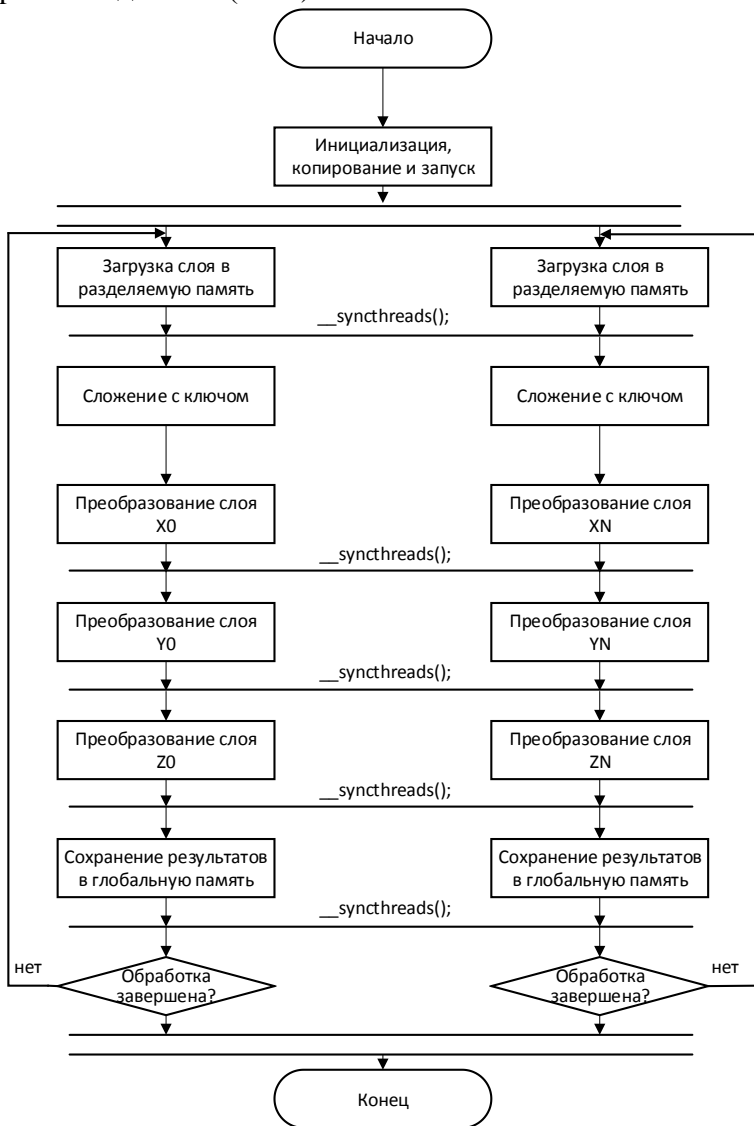


Рисунок 8. Схема параллельной CUDA-версии алгоритма DOZEN



Технология TDE используется в прототипе ВСПиХД для шифрования данных, выгружаемых во внешнюю память. Ключи шифрования разрядностью 256 бит генерируются при помощи ГПСЧ, основанного на алгоритме ГОСТ 28147-89. При этом, для каждой выгружаемой таблицы генерируется уникальный ключ шифрования, сохраняющийся в менеджере ключей (рисунок 9). В момент загрузки таблицы из внешней памяти ключ шифрования извлекается из блока менеджера ключей, после чего данные расшифровываются, а ключ уничтожается.

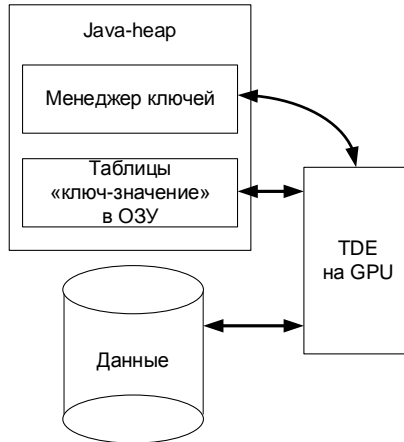


Рисунок 9. TDE по алгоритму ГОСТ 28147-89

**Третья глава** посвящена реализации разработанных во второй главе методов в виде прототипа ВСПиХД для последующего его использования в главе 4 для экспериментального подтверждения эффективности предложенных методов. В состав прототипа ВСПиХД входят компоненты двух видов: узел хранения и сервер доступа. В зависимости от количества узлов в кластерной системе можно инициализировать различное число компонентов прототипа ВСПиХД (рисунок 10).

Прототип ВСПиХД [16] создан с использованием технологий Java 8 и CUDA, взаимодействие которых осуществляется с использованием библиотеки jCuda. Для ускорения операций межузлового взаимодействия использовались асинхронные вызовы из библиотеки NIO.2. В прототипе реализовано ускорение операций хеширования на GPGPU-сопроцессоре для повышения производительности работы модифицированного фильтра Блума и подсистемы разметки типа «хеш-кольцо».

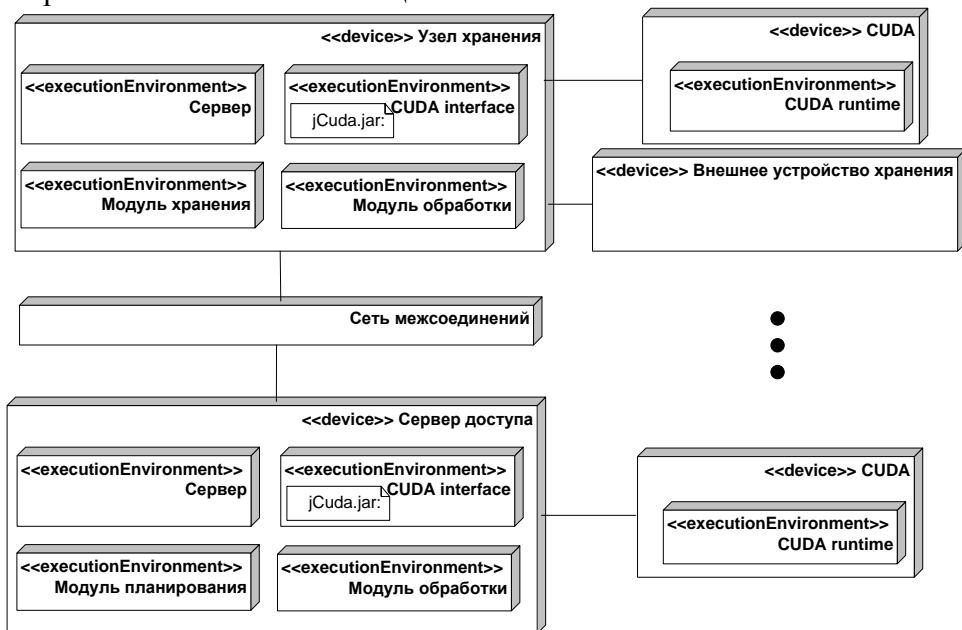


Рисунок 10. UML-диаграмма развертывания прототипа ВСПиХД

GPGPU-сопроцессоры использовались для ускорения алгоритма DOZEN, применяемого в стохастическом бинарном дереве поиска и алгоритма ГОСТ 28147-89, применяемого для зашифрования/расшифрования выгружаемых во внешнюю память данных.

**В четвертой главе** выполнено исследование методов и средств повышения производительности операций поиска и дополнительной обработки данных (TDE), предложенных в главе 2. С использованием метода анализа иерархий Т. Саати выделено два наиболее значимых критерия оценки эффективности: время отклика системы и пропускная способность при обработке потока заявок. Предложена методика проведения испытаний, учитывающая архитектурные особенности прототипа ВСПиХД и принципы проведения испытаний подобных систем, сконфигурирован тест для проведения испытаний. Количество испытаний определено в соответствии с теорией организации эксперимента и законом распределения времени поступления заявок:

$$n \geq \left( \frac{ts_v}{\delta} \right)^2, \quad (6)$$

где  $\delta$  – точность,  $s_v$  – стандартное отклонение по предвыборке,  $t$  – аргумент функции Лапласа, соответствующий заданной доверительно вероятности. Результаты сравнительных испытаний прототипа ВСПиХД и наиболее производительной из нынешних NoSQL-систем Apache Cassandra представлен в тексте диссертации.

Оценка адекватности модели проверялась по критерию Фишера, заключающегося в дисперсионном анализе. Для этого по формулам (7) и (8) рассчитывались дисперсия адекватности и дисперсия воспроизводимости серий экспериментов.

$$S_{\text{ад}}^2 = \frac{\sum_{i=1}^N n_i (\bar{y}_i - \bar{y})^2}{f_1}, \quad (7)$$

$$S_{\{y\}}^2 = \frac{\sum_{i=1}^N \sum_{q=1}^n (y_{iq} - \bar{y}_i)^2}{f_2}, \quad (8)$$

где  $N$  – число разных экспериментов,  $n_i$  – число параллельных опытов в  $i$ -м эксперименте. Например, в одной из серий экспериментов исследовалась зависимость времени отклика от числа ядер CPU. При этом было проведено 3000 параллельных опытов для каждого значения параметра «число ядер CPU» (рисунок 11). Таким образом, для данной серии:

$N = 3$ ,  $n = 3000$ . Прочие обозначения:  $y_{iq}$  – значение, полученное в  $q$ -м опыте  $i$ -го эксперимента (апостериорное значение),  $\bar{y}_i$  – среднее значение, полученное с доверительной вероятностью 0,95 на реальной системе после проведения  $n_i$  параллельных опытов в  $i$ -м эксперименте (среднее апостериорное значение),  $\hat{y}_i$  – значение, полученное при помощи модели в  $i$ -м эксперименте (априорное значение),  $f_1, f_2$  – число степеней свободы.

Получаемое по формуле (9) значение  $F$  сравнивалось с табличным  $F_T$ . В результате для всех серий экспериментов  $F$  оказалось меньше  $F_T$ , что свидетельствует об адекватности модели ВСПиХД, предложенной в главе 2.

$$F = \frac{S_{ад}^2}{S_{\{y\}}^2} \quad (9)$$

Для проведения анализа полученных результатов были построены графики для всех серий проведенных экспериментов. Примеры таких графиков представлены на рисунках 11-13.

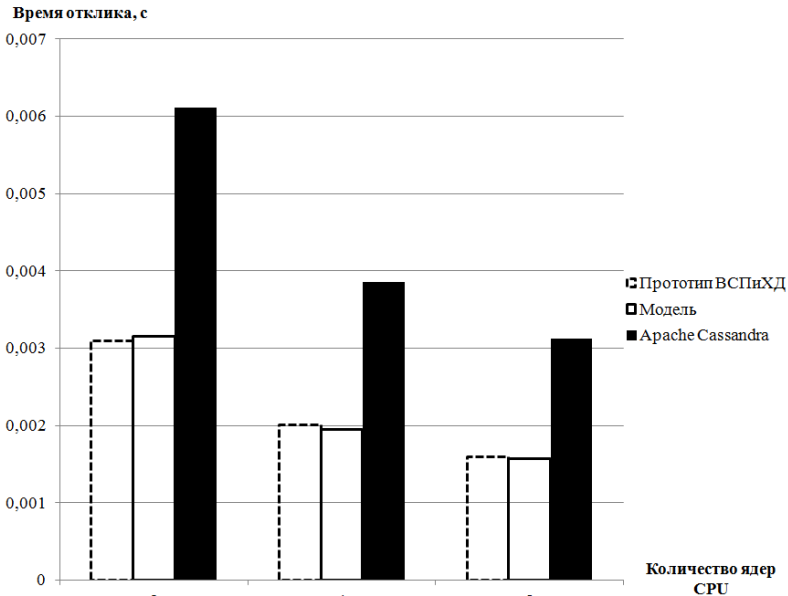
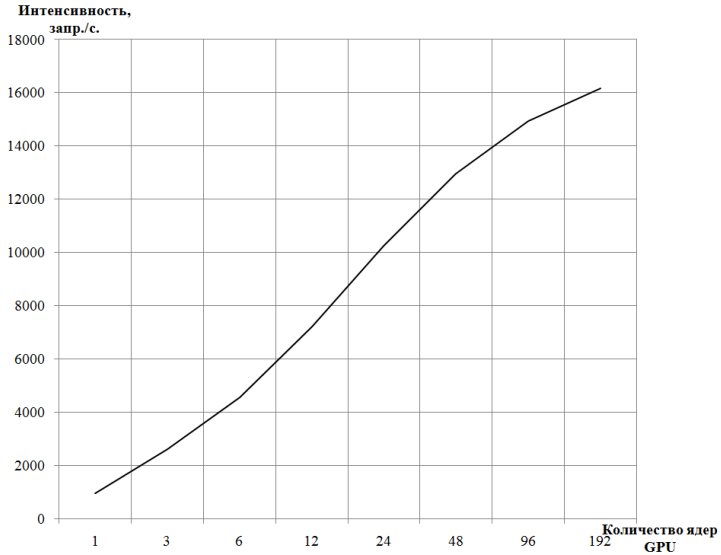
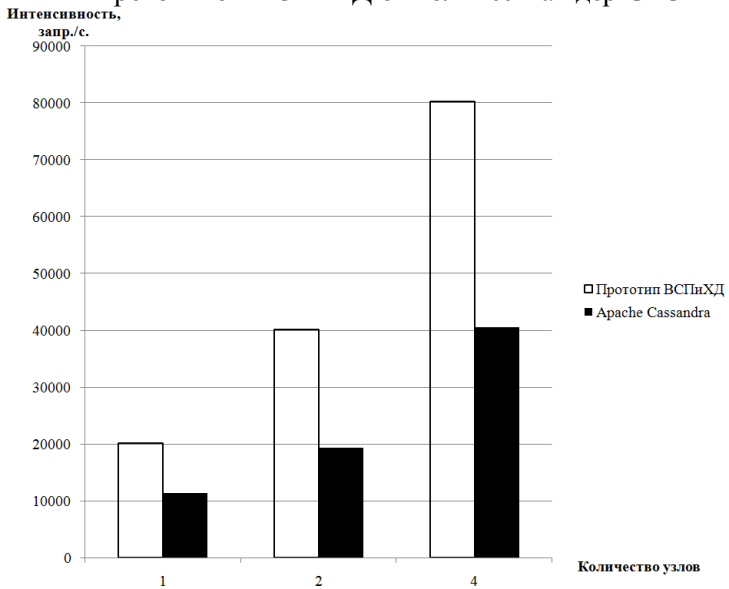


Рисунок 11. Зависимость времени отклика системы от количества ядер CPU



**Рисунок 12. Зависимость интенсивности обработки запросов прототипом ВСПиХД от количества ядер GPU**



**Рисунок 13. Зависимость интенсивности обработки запросов от количества узлов в системе**

## ОСНОВНЫЕ РЕЗУЛЬТАТЫ РАБОТЫ

В диссертации решена важная научная задача разработки методов и средств использования гибридных суперкомпьютерных технологий в распределенных системах поиска и хранения данных с целью повышения их защищенности и производительности. Основные результаты диссертационной работы:

1. Выполнен анализ и предложена классификация современных высокопроизводительных систем хранения и поиска данных. Построена классификация гибридных вычислительных систем. Сделан вывод о том, что ни одна из существующих высокопроизводительных систем хранения и поиска данных не поддерживает технологий гибридных вычислений и встроенных механизмов защиты, а существующие алгоритмы жестко привязаны к конкретным архитектурам.
2. Разработана математическая модель ВСПиХД, предназначенной для реализации в гибридных вычислительных системах. Архитектура современных ВСПиХД формализована в виде сети массового обслуживания, отражающей особенности взаимодействия составных частей ВСПиХД.
3. Разработан метод уменьшения интенсивности входного потока заявок на компоненты ВСПиХД при помощи модифицированного фильтра Блума (ФБ) для гибридных вычислительных систем. Предложена математическая оценка вероятности ложноположительного срабатывания модифицированного ФБ. Разработаны методы уменьшения времени обработки запросов ВСПиХД за счет использования дискретного сопроцессора.
4. Предложена новая высокопроизводительная структура данных – стохастическое дерево поиска. Данная структура эквивалентна сбалансированному двоичному дереву поиска. Независимость от исходного распределения ключей обеспечивается посредством их стохастического преобразования на гибридном сопроцессоре.
5. Разработаны методы организации встроенной системы шифрования данных ВСПиХД с использованием дискретного сопроцессора.
6. Разработаны статическая и динамическая UML-модели ВСПиХД, предназначенной для реализации в гибридных вычислительных системах. Модели выполнены в виде диаграмм прецедентов,

классов и последовательностей, что позволяет использовать их для построения подобных ВСПиХД как промышленного, так и учебного назначения.

7. Разработана архитектура прототипа высокопроизводительной системы поиска и защищенного хранения данных, предназначенная для использования в гибридных вычислительных системах. Архитектура описана при помощи диаграммы развертывания UML. Разработан протокол взаимодействия клиентов с прототипом ВСПиХД на базе механизма TCP-сокетов. Определен формат пакета, набор операций и типы данных, поддерживаемые прототипом ВСПиХД.
8. С использованием технологий Java и CUDA реализован прототип ВСПиХД в составе программных модулей узлов хранения и серверов доступа, использующих разработанную в рамках настоящей работы параллельную версию алгоритма ГОСТ 28147-89 для шифрования данных. Выполнены тестирование и отладка прототипа ВСПиХД, используемого для экспериментального исследования предложенных подходов.
9. Выделены критерии оценки эффективности разработанной системы ВСПиХД. При помощи метода анализа иерархий Т. Саати выделены наиболее значимые критерии: время отклика системы и интенсивность обработки запросов. Выполнена оценка количества испытаний, необходимых для получения экспериментальных данных с доверительной вероятностью 0,95.
10. Разработана методика проведения испытаний, содержащая, в частности, описание этапов проведения нагрузочного тестирования ВСПиХД. Полученные экспериментальные данные были использованы для подтверждения адекватности предложенной математической модели ВСПиХД по Фишеру.
11. Проведен анализ результатов сравнительных экспериментальных испытаний прототипа ВСПиХД и Apache Cassandra (как одной из наиболее производительных современных систем хранения и поиска данных), по итогам которого можно утверждать, что предложенные в настоящей работе методы являются эффективными и позволяют проектировать защищенные системы, превосходящие свои самые производительные аналоги в 1,5-2 раза.

Основываясь на предложенных в работе теоретических оценках, которые получили экспериментальное подтверждение, разработанные методы и средства использования гибридных суперкомпьютерных технологий для повышения производительности систем поиска и защищенного хранения данных позволяют достичь поставленную в работе цель.

**Основные результаты диссертационной работы** изложены в 18 печатных трудах:

1. Ровнягин М.М. Современные средства высокопроизводительной передачи данных в многопроцессорных и многомашинных вычислительных системах // Труды 54-й научной конференции МФТИ. М.: МФТИ. С. 22-23, 2011.
2. Ровнягин М.М., Лебедев М.С., Чудновский А.Л. Современные методы верификации и особенности их применения // Сборник избранных трудов VI Международной научно-практической конференции: учебно-методическое пособие. Под ред. проф. В.А. Сухомлина. М.: ИНТУИТ.РУ, С. 1009-1020, 2011 **(Индексируется РИНЦ)**.
3. Ровнягин М.М. Использование UVM для автономной верификации цифровой аппаратуры // Сборник трудов V Всероссийской научно-технической конференции "Проблемы разработки перспективных микро- и нанoeлектронных систем - 2012" под общ. ред. академика РАН А.Л. Стемпковского. М.: ИППМ РАН, С. 129-132, 2012 **(журнал из перечня ВАК, индексируется РИНЦ)**.
4. Ровнягин М.М. Модифицированный фильтр Блума с применением технологии NVIDIA CUDA для высокопроизводительного поиска данных // Сборник докладов 16-й Международной телекоммуникационной конференции молодых ученых и студентов «Молодежь и наука», М.: МИФИ, С. 62-63, 2013.
5. Васильев Н.П., Ровнягин М.М. Гибридные кластеры как средства организации бюджетных суперкомпьютеров и вычислительных облаков // Автоматизация в промышленности, № 4, С. 51-54, 2013 **(журнал из перечня ВАК, индексируется РИНЦ)**.
6. Ivanov M.A., Rovnyagin M.M. et al. Using Sequential and Parallel Composition for Stochastic Data Processing // International



- Conference «The Radio-Electronic Devices and Systems for The Infocommunication Technologies» (RES-2013), Moscow, P. 152-155, 2013.
7. Vasilyev N.P., Rovnyagin M.M. et al. Modified Bloom filter for high-performance data search in hybrid computing systems // International Conference «The Radio-Electronic Devices and Systems for The Infocommunication Technologies» (RES-2013), Moscow, P. 167-170, 2013.
  8. Ровнягин М.М. Стохастическое бинарное дерево поиска в задачах поиска и хранения данных для высокопроизводительных NoSQL-систем // Сборник докладов 17-й Международной телекоммуникационной конференции молодых ученых и студентов «Молодежь и наука», М.: МИФИ, С. 38-39, 2014.
  9. Кузнецов А.А., Ровнягин М.М. Исследование статистических свойств модифицированного фильтра Блума для гибридных систем // Научная сессия НИЯУ МИФИ-2014. Аннотации докладов. Т.3., М.: НИЯУ МИФИ, 2014. С. 64.
  10. Васильев Н.П., Иванов М.А., Ровнягин М.М. и др. Использование генераторов псевдослучайных чисел при построении алгоритмов стохастического преобразования данных. // Вестник НИЯУ МИФИ, том 3, С. 101-106, 2014 (**Журнал из перечня ВАК, индексируется РИНЦ**).
  11. Ivanov M.A., Rovnyagin M.M. et al. Three-Dimensional Data Stochastic Transformation Algorithms for Hybrid Supercomputer Implementation // 17th IEEE Mediterranean Electrotechnical Conference. 13-16 April 2014 Beirut, Lebanon, pp.451, 457 (**Индексируется Scopus**).
  12. Васильев Н.П., Ровнягин М.М. Бюджетный суперкомпьютер с гибридной CPU/GPU-архитектурой как эффективный инструмент в учебной и научной деятельности. // Вестник НИЯУ МИФИ, Т. 3. № 3. С. 378-384, 2014 (**Журнал из перечня ВАК, индексируется РИНЦ**).
  13. Vasilyev, N.P., Rovnyagin M.M. Hybrid clusters for budget supercomputers and cloud computing // AUTOMATION AND REMOTE CONTROL, Vol. 75, Issue 10, 2014, pp. 1869-1874. (**Индексируется Web of Science, Scopus**).

14. Vasilyev, N.P., Rovnyagin, M.M. et al. Modified Bloom filter for high performance hybrid NoSQL systems // *Life Science Journal*, vol. 11 (SPEC. ISSUE 7), 2014, pp. 457-461 (**Индексируется Scopus**).
15. Васильев Н.П., Ровнягин М.М. VAR - Программная платформа для организации высокопроизводительного поиска данных в гибридных вычислительных системах // Научная сессия НИЯУ МИФИ-2015. Аннотации докладов. В 3 томах. Т.3., М.:НИЯУ МИФИ, 2015. – С.49.
16. Ровнягин М.М. Система VAR. Высокопроизводительный поиск данных с применением GPGPU-технологий // Сборник докладов 18-й Международной телекоммуникационной конференции молодых ученых и студентов «Молодежь и наука», М.: МИФИ, С. 75-76, 2015.
17. Dyumin, A.A.; Kuznetsov, A.A., Rovnyagin, M.M. Evaluation of statistical properties of a modified Bloom filter for heterogeneous GPGPU-systems // *Young Researchers in Electrical and Electronic Engineering Conference (EIconRusNW)*, 2015 IEEE NW Russia , vol., no., pp. 71, 74, 2-4 Feb. 2015, doi: 10.1109/EIconRusNW.2015.7102234 (**Индексируется Scopus**).
18. Dyumin, A.A.; Puzikov, L.A.; Urvanov, G.A.; Chugunkov, I.V., Rovnyagin, M.M. Cloud computing architectures for mobile robotics // *Young Researchers in Electrical and Electronic Engineering Conference (EIconRusNW)*, 2015 IEEE NW Russia , vol., no., pp. 65, 70, 2-4 Feb. 2015 doi: 10.1109/EIconRusNW.2015.7102233 (**Индексируется Scopus**).



