

# Анализ основных тенденций в области хранения данных

Е.О. Тютляева, А.А. Московский

**Аннотация.** В статье анализируется тенденция возникновения и развития все большего числа научно-исследовательских проектов, накапливающих и анализирующих архивы данных большого объема (0,1–100 ПБ). Проводится анализ возможностей современных суперЭВМ в части подсистемы хранения данных, на основе примерных характеристик вычислительных машин из верхней части рейтинга Top 500. Проведен краткий анализ современных проектов развития средств хранения с перспективой применения на машинах экзафлопсного класса.

**Ключевые слова:** системы хранения данных; приложения, обрабатывающие большие объемы данных; тенденции развития.

## Введение

Благодаря ускоренному развитию микроэлектроники, непрерывно возрастающее количество наблюдательных приборов с все более высоким разрешением позволяет получать большие объемы данных (достигающие сотен терабайт и петабайт) в самых различных сферах человеческой деятельности, включая естественные науки, социологию и экономику. Современные мощности хранения позволяют сохранять и архивировать данные, которые могут представлять интерес для последующего исследования.

Под данными подразумеваются самые различные необработанные информационные материалы, включая не только снимки дистанционного зондирования, персональные медицинские данные, полные сырые данные различных наблюдений и экспериментов, но и базы данных социальных сетей, различных магазинов и прочую статистическую информацию. Кроме того, с ростом технологий возрастают требования к разрешению данных, и новые вычислительные эксперименты предполагают более широкие временные и пространственные диапазоны обрабатываемых сырых данных.

Как утверждается в отчете «Большие данные: следующий рубеж для инноваций, соревнования и производительности» [1] – данные сегодня становятся важнейшим фактором продукции наравне с материальными активами и человеческим капиталом. В обозримом будущем экспоненциальный рост объема данных должен продолжиться в связи с возрастающей интенсивностью представления данных и сбора информации. Параллельно будет развиваться комплексное представление информации, объем социальных коммуникаций, количество информации, представленной в Интернете. Большие объемы данных имеют значительный потенциал для того, чтобы стать значительной ценностью для бизнеса и пользователей.

По утверждениям ИВМ, каждый день создаются около 15 ПБ новых данных [2]. Это и научные данные, и сведения о проведенных операциях-транзакциях, и новые фотографии и отчеты в социальных сетях. Согласно информации из газет, для создания фильма «Аватар», потребовалась система хранения данных (далее в тексте СХД) более чем на 1 ПБ [3]. В социальную сеть Facebook каждый день добавляется около 12 ТБ данных (после сжатия) [4].

В книге «The Fourth Paradigm: Data-Intensive Scientific Discovery» [5] исследование огромных массивов данных называют четвертой парадигмой науки, после экспериментальной, теоретической и вычислительной парадигм, сменявших друг друга на разных стадиях ее развития. Автором высказано предположение, что выявление закономерностей в больших массивах данных становится основным инструментом для исследования и получения новых знаний в передовых областях науки в наше время. К примеру, в [6] упоминается, что с 2001 . до 2009 г. количество баз данных, зарегистрированных в Nucleic Acids Research, увеличилось с 218 до 1170.

В то же время растущие объемы социальных данных способствуют расширению числа и масштабов исследовательских задач в области менеджмента, исследования рынка и социальной активности в виде аналогичных задач анализа взаимосвязей и закономерностей.

С увеличением интенсивности работы с данными во всех областях человеческой деятельности характеристики подсистемы ввода-вывода и управления данными становятся одними из наиболее проблемных компонент современных информационных и вычислительных систем.

В данной статье мы постараемся провести краткий обзор ряда приложений, требующих интенсивной работы с данными, существующих решений в области высокопроизводительного хранения данных; оценить тенденции развития СХД, характеристик подсистем хранения данных наиболее мощных суперЭВМ и их соответствие реальным требованиям актуальных приложений.

### **Приложения, требующие интенсивной обработки больших объемов данных**

В настоящее время проблема приложений, связанных с интенсивной работой с большими объемами данных, находится на переднем крае науки. В англоязычной литературе такие задачи получили название «data intensive», что можно перевести как «оперирующий большими объемами данных» (далее в тексте ОБОД). ОБОД

называют те вычислительные задачи, в которых хранение, обработка и анализ значительных объемов данных становится первостепенной проблемой [7].

Сложность при обработке больших объемов данных порождает технологические проблемы как на уровне подсистемы хранения (скорость чтения/записи, надежность, доступный объем), так и на уровне обработки (доступные полосы пропускания оперативной памяти, возможный темп запросов в ОЗУ).

Для оценки готовности системы к обработке значительных объемов данных в 2010 году был анонсирован новый рейтинг, Graph500 [8], который является первой серьезной попыткой дополнить список ТОП-500 оценкой возможностей системы оперировать большими объемами данными. Текущие тесты производительности, которые используются для построения рейтинга ТОП-500, не позволяют оценить пригодность высокопроизводительной установки для ОБОД приложений. Несмотря на то, что в ранжировании в новом списке пока приняли участие только 29 установок, уже очевидно, что результаты данного исследования значительно отличаются от рейтинга TOP-500. К примеру, лидирующий в Graph-500 суперкомпьютер Intrepid занимает всего 15 место в ТОП-500 и обгоняет по скорости работы с данными Jaguar (19 место в Graph-500 против 3 места TOP-500), Nopper (4 Graph-500 против 8 TOP-500), Jugene (2 Graph-500 против 12 TOP-500) и Lomonosov (3 Graph-500 против 13 TOP-500). Данный рейтинг позволяет оценить, прежде всего, готовность оперативной памяти системы к обработке приложений, интенсивно работающих с большими объемами данных. Тем не менее, даже этот рейтинг не дает возможности оценить системы хранения, характеристики которых играют ключевую роль для современных ОБОД приложений. Интересным курьезом является вхождение в текущую редакцию рейтинга машины из 1 узла суперЭВМ Kraken, использовавшего для хранения данных задачи, обычно размещаемых в ОЗУ, высокоскоростную СХД на твердотельных накопителях Fusion IO [55].

На уровне развития СХД перед исследователями стоят принципиально новые задачи. Наиболее значительным проектом, ставящим

высокую планку для современных высокопроизводительных СХД, стоит назвать Большой Адронный Коллайдер в CERN. Миллионы сенсоров БАК генерируют около петабайта данных в секунду. В связи с тем, что современные мощности хранения не способны поддержать такие объемы хранения данных, большая часть измерений отфильтровывается на основе простейших правил. «Цель — не потерять ничего интересного». Тем не менее, даже после фильтрации и предобработки данных, коллайдер производит до 25 ПБ данных в год. Всего в центре обработки и хранения данных CERN 34 ПБ магнитных носителей и 45.3 ПБ на дисковых носителях [9].

Астрономия - наука, непосредственно связанная с обработкой больших объемов данных, - также демонстрирует наличие актуальных проектов этой области. Приборы для астрономических наблюдений позволяют получать данные со все более высоким разрешением, исследователи заинтересованы в долгосрочном хранении полученных архивов для возможности последующих исследований. Кроме того, астрономические данные в большинстве своем не имеют ограничений приватности или коммерческой тайны, научное сообщество заинтересовано в общедоступности полученных данных, новых задачах, исследованиях и экспериментах, что накладывает дополнительные требования к распределенности и доступности данных.

Одним из примеров может быть система телескопов панорамного обзора и быстрого реагирования Pan-STARRS, нацеленная на обнаружение и изучение приближающихся к Земле объектов, включая астероиды и кометы, которые могут оказаться опасными для нашей планеты. Одной из особенностей проекта является новаторская цифровая камера, позволяющая получать изображения 38000×38000 пикселей [10]. Каждое изображение, сделанное одной Pan-STARRS камерой, содержит около 2 ГБ данных. В режиме полного обзора объем необработанных данных телескопа за ночь достигает нескольких терабайт.

В докладе «Вычислительные вызовы в астрономии волн тяготения» [11] упоминается, что в проекте участвует 4 детектора, располо-

женные на двух континентах, которые собирают необходимые данные для дальнейшего анализа и моделирования. Каждый детектор обладает скоростью передачи данных 10 МБ/с. В докладе сказано, что годовой объем данных для 3-х детекторов составляет 947 ТБ. Основные проблемы, которые стоят перед исследователями, заключаются в управлении научными данными (в частности, полученными в результате анализа), управлении научными потоками, а также в недостатке квалифицированных кадров для разработки требуемой инфраструктуры.

Другим примером проекта, интенсивно работающего с данными, является NEEShub: киберинфраструктура данных для моделирования землетрясений [12]. Цель данного проекта - создание национальной многопользовательской исследовательской инфраструктуры для поддержки исследований и инноваций по минимизации ущерба от землетрясений и цунами. Под управлением проекта находятся значительные объемы разнообразных научных данных, включая изображения, видео, текст и т.п. По состоянию на март 2011 года в реализованной киберинфраструктуре находилось 417 проектов и почти 1 миллион файлов, объем данных менее 1 ПБ.

Проект, посвященный изучению нейронных связей, в настоящее время позволил получить 10 ТБ данных — результатов моделирования, представляющих нейронные связи для около 1/80000 мозга мыши. Данный проект имеет значительный потенциал для масштабирования. Следующей целью является моделирование кубического миллиметра, что составит 1/1000 мозга мыши и займет предположительно более 1 ПБ. Моделирование целой мыши по предположениям исследователей потребует системы хранения объемом в эксабайт [13].

Большие объемы данных накапливаются и в климатологии. Например, немецкий центр по изучению климата (DKRZ, г. Гамбург) оснащен не только мощными суперЭВМ (более 150 Тфлопс), но и средствами визуализации данных, такими как специализированные комнаты, а также многоуровневой системой хранения данных общим объемом около 60 ПБ [54].

Приведенные примеры свидетельствуют о том, что сегодня задачи, связанные с передовыми направлениями в науке, работают с

большими объемами данных и уже сейчас имеют очень высокие и обоснованные требования к объему и производительности СХД. Построение соответствующих систем хранения и преодоление возникающих барьеров должно являться одной из ключевых задач современной суперкомпьютерной отрасли.

## Существующие решения в области хранения данных

### Стандартные

С точки зрения «обычных пользователей», бизнеса, наиболее удобными представляются «коробочные» версии высокопроизводительных систем хранения данных, представляющие собой настроенный и готовый к работе программно-аппаратный комплекс. Все ведущие поставщики ИТ-решений (IBM, HP, Oracle и другие) имеют в своих продуктовых линейках оригинальные либо заимствованные комплексы хранения данных. Существует и ряд специализированных компаний, которые успешно поставляют подобные хранилища «под ключ», такие как EMC. Приведем несколько примеров.

Компания Dell Terascale [14] предоставляет высокопроизводительные решения в области хранения, стандартная конфигурация которых может предоставлять емкость 768 ТБ для пользовательских данных под управлением высокопроизводительной файловой системы Lustre и с поддержкой на 3 года. Другой лидирующей компанией на рынке США в области предоставления и поддержки высокопроизводительных систем хранения является компания Xugatex [15], которая также предоставляет высокопроизводительные хранилища с параллельным доступом. На российском рынке существует достаточно широкий спектр предложений.

### НРС

Для описанных выше задач, используются более сложные, единичные разработки, которые представляют собой сложную совокупность инженерных решений и программных продуктов.

Одним из способов оценить состояние суперкомпьютерного рынка является изучение статистики, предоставляемой рейтингом ТОП-500

[16]. ТОП-500 — это поддерживаемый в актуальном состоянии с 1993 года список пятисот самых мощных высокопроизводительных компьютеров в мире. Применительно к рассматриваемой теме, мы можем при помощи данного списка получить адекватный список СХД, обладающих достаточной надежностью, масштабируемостью и производительностью для использования на ведущих суперкомпьютерах мира. Изучение технических характеристик СХД, представленных на ведущих суперкомпьютерах мира, может позволить оценить тенденции развития систем хранения, отметить основные проблемы и намеченные пути их решения.

Согласно последнему рейтингу ТОП-500, который вышел в июне 2011 года, самой высокопроизводительной машиной мира является японский компьютер K computer. Несмотря на то, что этот компьютер уже лидирует в списке ТОП-500, согласно планам он будет окончательно сдан в эксплуатацию только в 2012 году.

Для данного суперкомпьютера разрабатывается система сверхвысокой масштабируемости FEFS [17]. Оперативная память одного узла K компьютера более 1 ПБ, всего узлов планируется более 80000. Предполагаемая файловая система должна обладать экстремально большой емкостью (от 100 ПБ до 1 ЭБ), значительным количеством клиентов (100 тыс. до 1 млн.) и серверов (1 до 10 тыс.).

Предполагаемые характеристики СХД:

- пропускная способность одиночного потока (~1 Гб/с) параллельного ввода-вывода (~Тб/с);
- сокращенное время ожидания открытия файла (~10 тыс. оп/с);
- всегда доступный файловый сервис, даже если какая-то часть системы сломана/недоступна.

Предполагаемая файловая система отражает планы и перспективы в направлении построения высокопроизводительных систем хранения для суперкомпьютеров в эру ОБОД приложений. Характеристики уже реализованных систем хранения на остальных установках первой десятки отражают реальное состояние СХД на сегодняшний день.

Рассмотрим таблицы, отражающие состояние СХД на ведущих суперкомпьютерах мира в 2011, 2006 и 2001 годах и проанализируем полученные сведения.

Табл. 1. Состояние СХД на ведущих суперкомпьютерах мира в июне 2011 года

	Имя	Объем	Пропускная способность	Файловая система	Дополнительная информация
1	K computer	(от 100 ПБ до 1 ЕБ) ождается	(~ГБ/с) Параллельного ввода-вывода (~ТБ/с). ождается	FEFS	Япония
2	Tianhe-1A [18]	1PB (2 PB по некоторым данным)		Lustre	Китай
3	Jaguar [19]	10 ПБ	240 ГБ/с	Spider (Lustre extension)	США
4	Nebulae	-	-	-	Китай
5	TSUBAME2.0 [20]	15 ПБ, иерархическое		7.13PB (Lustre + NFS Home)	Япония Дополнительно доступно 8 ПБ СХД на магнитных лентах
6	Cielo - Cray XE6 [21]	10 ПБ (в разработке)	160 ГБ/с (в разработке)	PANASAS (в разработке)	США
7	Pleiades [22]	Всего доступно 6.9 ПБ total		7 файловых систем Lustre	США
8	Hopper [23]	2 ПБ рабочей памяти	35 ГБ/с	Lustre	США Дополнительно доступны все глобальные файловые системы NERSC, к примеру, HPSS на 59 PB.
9	Tera-100 [24]	20 ПБ	500 ГБ/с	Lustre	Франция
10	Roadrunner [25]	2 ПБ	~60 ГБ/с	PANASAS	США
11	Kraken XT5 [26]	3.3 ПБ		NFS, HPSS	США
12	JUGENE [27]	5.3 ПБ	66 ГБ/с	GPFS	Германия
13	Lomonosov [28]	500 ТБ + 300 ТБ + 1 ПБ			Россия Трехуровневая СХД, включающая 500 ТБ T-Platforms ReadyStorage SAN, 300TB NAS storage и 1 ПБ на магнитных лентах
14	BlueGene/L [29]	1,89 ПБ			США, содержит 1,024 ГБ/с соединений с глобальной файловой системой
15	Intrepid [30]	~8 ПБ	35 ГБ/с	GPFS	США

Табл. 2. Состояние СХД на ведущих суперкомпьютерах мира в июне 2006 года

	Имя	Объем	Пропускная способность	Файловая система	Дополнительная информация
1	BlueGene/L - eServer Blue Gene Solution				США
2	BGW - eServer Blue Gene Solution [31]	60 ТБ		GPFS	США Дополнительно используется 500 ТБ IBM 3494 на магнитных лентах

	Имя	Объем	Пропускная способность	Файловая система	Дополнительная информация
3	ASC Purple [32]	1.6 ПБ (2 ПБ на 2007 г.)	102 ГБ/с		США, эта система показывала высокую пропускную способность, и позволила преодолеть так называемый «гигабайтовый барьер», выражающийся в неспособности интерконнекта большого суперкомпьютера «насытить» процессор данными
4	Columbia [33]	650 ТБ RAID storage			США; Дополнительно 10 ПБ на магнитных лентах
5	Tera-10 [34]	1 PB	100 ГБ/с	Lustre	Франция
6	Thunderbird [35]	120 ТБ 50 ТБ	6.0 ГБ/с 4.0 ГБ/с	Lustre PANASAS	США, две файловые системы показаны в двух строках
7	TSUBAME Grid Cluster [36]	1 ПБ (2007)	8 ГБ/с	Lustre	Япония, первая промышленная система объединившая программный RAID Linux и Lustre.
8	JUBL	-	-	-	Германия
9	Red Storm [37]	340 ТБ, 1753 ТБ к 2008	Цель — 50.0 ГБ/с для каждого цвета	Lustre	США
10	Earth-Simulator	240 ТБ HDD RAID			Япония Иерархическое хранилище, 1.5 ПБ кассетных накопителей на магнитных лентах
11	MareNostrum [38]	280 ТБ			Испания, «самый красивый суперкомпьютер мира»
12	Stella				Нидерланды
13	Jaguar - Cray XT3 [39]	600 ТБ		Lustre	США
14	Thunder [40]	200 ТБ	6.4 ГБ/с	Lustre	США
15	Blue Protein				Япония

Табл. 3. Состояние СХД на ведущих суперкомпьютерах мира в июне 2001 года

	Name	Volume	Bandwidth	FS	Additional Info
1	ASCI White [41]	160 ТБ	-	GPFS	США
2	SP Power3	-	-	-	США
3	ASCI Red [42]	12.5 ТБ RAID			США, Дополнительно было хранилище на магнитных лентах
4	ASCI Blue-Pacific SST [43]	62.5 ТБ - RAID5 0 глобальная файловая система; 17 ТБ – локальные диски	6.6 ГБ/с – глобальная; 11 ГБ/с – локальная файловая система.	GPFS	США, иерархическое хранилище, HDD на узлах.
5	SR8000/MPP				Япония, для вычислений с высокой точностью
6	ASCI Blue Mountain [44]	76 ТБ			США

Проанализируем полученные таблицы. Как известно, в базовые сведения, которые сообщаются в TOP-500 о каждом суперкомпьютере, информация о конфигурации системы хранения не входит, что согласуется с природой теста LINPACK, результаты которого не зависят от СХД. В связи с этим в заполнении таблиц есть пробелы, т.к. производители некоторых установок не публикуют данную информацию.

Тем не менее, в первую очередь следует отметить, что масштабы систем хранения претерпели не столь колоссальные изменения за 5 лет с 2006–2011. В 2006 году в первой десятке суперкомпьютеров лидирующей была СХД суперкомпьютера ASC Purple, которая обладала объемом в 1.6 ПБ и пропускной способностью в 102 ГБ/с (Табл. 2). Между тем, в 2011 году в первой десятке ведущих суперкомпьютеров мира уверенно держит место китайский суперкомпьютер Tianhe-1A с СХД, достигающей по различным данным размера от 1 до 2 ПБ, т.е. сравнимой с системой хранения суперкомпьютера Purple (Табл. 1). В первой десятке также можно наблюдать суперкомпьютеры с пропускной способностью ввода/вывода не достигающей 100 ГБ/с (Hopper, Roadrunner — из тех для которых эти данные доступны), хотя этот барьер был также преодолен в 2006 году. Самыми лучшими характеристиками из первой десятки TOP-500 обладает система хранения суперкомпьютера Tera-100 (20 ПБ — объем, 500 ГБ/с — пропускная способность, т.е. объем в 12.5 раз, а пропускная способность в 4.9 раз больше чем у лучшего хранилища в 2001 году).

Для сравнения, теоретическая пиковая производительность с июня 2006 года изменилась с (18.20–280.60 TFLOPS) до (557.06–8773.63 TFLOPS) (Табл. 2). (Первое значение — минимальная теоретическая пиковая производительность системы из 15, второе — максимальная, из TOP-500), т.е. лучшая пиковая производительность увеличилась в 31 раз.

Кроме того, нельзя не отметить принципиальный разброс в объемах систем хранения за 2006 год (от 0,060 ПБ до 1.6 ПБ). В рейтинге за 2011 год разброс менее принципиален (от 1 ПБ до 20 ПБ), все системы (для которых доступна информация), вошедшие в первую десятку, обладают системой хранения с объемом, превы-

шающим 1 ПБ. Намечившаяся тенденция к выравниванию характеристик систем хранения показывает, что наличие адекватной системы хранения становится все более важным для современного суперкомпьютера. Следует также отметить, что суперкомпьютерные центры более развитых стран — США, стран Европы — обладают значительно превосходящими емкостями хранения по сравнению с машинами из Китая, хотя последние и могут занимать более высокое положение в рейтинге Linpack.

Тем не менее, заметный рост масштабов систем хранения суперЭВМ значительно ниже роста вычислительных мощностей. Сложно однозначно назвать причину, можно лишь сформулировать ряд предположений.

1) При построении рейтинга TOP-500 не учитываются характеристики системы хранения. Тем не менее, именно рейтинг TOP-500 имеет ключевое значение в мире высокопроизводительных вычислений и представляет наибольший экономический и даже политический интерес для производителей, реальных пользователей.

2) Особенно актуальными задачи с интенсивной работой с данными стали именно сейчас. Это связано как с улучшением характеристик приборов наблюдения, получением данных с более высоким временным и пространственным разрешением, так и с накоплением архивов цифровых данных наблюдений, архивов данных от социальных сетей и экономических баз данных в беспрецедентных ранее масштабах.

3) Влияние на рост масштабов систем хранения могут оказывать технологические проблемы. Это и проблемы в области отказоустойчивости, надежности, обеспечение надлежащей пропускной способности, особенно в области системного ПО, включая файловые системы. В частности, лидирующую позицию на сегодняшний день среди СХД для суперкомпьютеров занимает ФС Lustre. Между тем увеличение масштабируемости Lustre обходится в миллионы долларов и годы разработки.

## Exascale

Исследование возможностей увеличения масштабируемости, пропускной способности и

надежности данных ведется также в рамках инициатив по созданию вычислительного кластера экса-класса.

Roger Haskin из исследовательской группы IBM General Parallel File System предполагает, что увеличение масштабов суперкомпьютера до Exascale завершит извлечение файловой системы и хранилища из суперкомпьютера, т.е. система хранения будет существовать отдельно, аналогично файловому серверу, соединенная с вычислительным суперкомпьютером при помощи высокопроизводительных коммутационных решений. Он считает, что встроенные узлы ввода-вывода не предоставляют удовлетворительного объема памяти для оперирования данными и обладают рядом других недостатков [45].

Известна концепция, согласно которой на высокомасштабируемых суперкомпьютерах будет применяться иерархическая система хранения. Одной из наиболее интересных разработок в этом направлении является файловая система Colibri (руководитель разработки Peter Braam в компании Xyratex [46]). В системе предполагается промежуточный уровень - прокси - для быстрого сохранения большого объема данных, который будет предоставлять начальную пропускную способность и нижний уровень для традиционного хранения необходимого объема. Предполагается, что уровень прокси будет состоять из высокоскоростных твердотельных накопителей передового уровня технологии, а нижний уровень из более традиционных дисков. Идея буферизации данных (предоставления промежуточного аппаратного слоя между оперативной памятью и хранилищем) также разрабатывалась исследователями из национальной лаборатории США Аргон [47].

Другую позицию представляют исследователи из университета Токио, которые предполагают, что независимые системы ввода-вывода не демонстрируют надлежащую масштабируемость. Они предполагают, что использование развивающихся устройств хранения, таких как solid-state disks (SSDs) или Storage Class Memories (SCM), перспективны для организации ввода-вывода, увеличения производительности и оптимизации энергопотребления. Базируясь на данных технологических изменений, японские ученые предлагают исследовать воз-

можности активного хранения данных, уменьшения нагрузки на сохранение метаданных, анализ, организацию и перераспределение данных [48].

Наибольшие технологические проблемы связаны не с разработкой аппаратной базы, а с изменением концепций программного обеспечения для поддержки беспрецедентного уровня масштабируемости. Вышеупомянутая файловая система Colibri обещает переопределить стандартные парадигмы хранения данных. В системе будет использоваться концептуально отличная модель данных — Модель хранения объектов. Основную единицу будет представлять из себя объект-«контейнер», аналогичный Логическому Тому, но обладающий дополнительной операцией «вложение». Предполагается, что контейнеры можно будет вкладывать один в другой (вырожденный случай — переименовать) без разбора содержания. Такой подход позволяет значительно снизить накладные расходы на работу с метаданными и делает концепцию иерархического хранилища очень эффективной. База данных размещения контейнеров предполагается более абстрактной, чем в текущих системах, базирующейся не на таблицах, а на формулах-зависимостях.

В системе предполагается вести журнал ошибок транзакций, который будет хранить информацию обо всех совершенных действиях для предоставления возможности отладки и мониторинга. Исследователями предлагаются также идеи об использовании опыта торрент-систем для повышения эффективности чтения и методы интеллектуального кэширования (использование опыта предыдущего запуска вычисляемой задачи и соответствующее перемещение данных, которые могут потребоваться в «быструю» память — твердотельный прокси слой).

## Облачные вычисления

ОБОД приложения могут быть реализованы и в модели облачных вычислений. К традиционным преимуществам относятся перенос больших начальных затрат на покупку и поддержание дорогостоящего оборудования и организации центра данных, на «плоскую» систему оплаты услуг облачных инфраструктур. Когда речь идет о долгосрочном хранении дан-



ных и проведении различных исследований над сырыми данными в какой-то конкретной области науки, может идти речь о создании специализированного облака, к которому будут иметь доступ профильные специалисты.

Соответственно, с возросшими требованиями к разрешению, объемам и обработке данных, повышаются требования и к предоставляемой облачной системой инфраструктуре системы хранения. В рамках этих требований модифицируются и создаются новые облачные платформы. Одним из примеров создающихся облачных платформ может служить платформа VISION Cloud [49], программа по созданию которой составлена с октября 2010 по сентябрь 2013. Цель данного проекта — создание мощной инфраструктуры для предоставления надежных и эффективных ОБОД сервисов хранения, упрощение сближения информационных и коммуникационных технологий, СМТ и телекоммуникаций. В рамках данного проекта развивается более абстрагированная модель системы хранения, чем традиционные файловые системы, схожая с описанной моделью файловой системы Colibri.

Полный обзор облачных проектов выходит за рамки данной статьи, но это активно развивающееся направление, включающее в себя объединение суперкомпьютеров, создание вычислительных сетей и масштабных центров данных.

## Перспективы

Наиболее перспективным представляется развитие иерархических систем хранения данных с использованием твердотельных (SSD) дисков в качестве одного из уровней хранения. По ряду факторов, включая стоимость, производительность и надежность, SSD диски еще не могут полностью заменить HDD диски в высокопроизводительных СХД. Тем не менее, SSD диски могут и должны занять соответствующее место в иерархии хранения.

Эффективному использованию SSD дисков в высокопроизводительных СХД посвящен ряд проектов. К примеру, в статье [50] представлена комбинированная система хранения с SSD и HDD дисками с улучшенной производительностью, в проекте [51] показаны перспективы

использования SSD дисков для хранения контрольных точек.

Colibri — уже названный выше проект, который подразумевает использование SSD в качестве промежуточного слоя в иерархии хранения. Как показывает исследование рейтинга Top500, иерархические СХД, использующие магнитные носители в качестве одного из уровней хранения, популярны и сейчас. Добавление нового уровня SSD может повысить энергоэффективность и производительность СХД.

В ряде исследований можно заметить наметившиеся тенденции к приближению части вычислений к местам хранения данных за счет реализации технологий активного хранения. Такие предложения в рамках экза-исследований предлагают исследователи из университета Токио. Более детально эта мысль была рассмотрена в презентации [52] под названием «Киберкирпичи», по аналогии с «активными дисками» Jim'a Gray. Для построения «кирпича» предлагается использовать материнскую плату Zotac Atom/ION, двухъядерный процессор Atom и 7.7 TB на SSD накопителях.

Достоинствами данной системы являются низкое энергопотребление и способность выполнить часть операций по обработке данных непосредственно в пределах данного блока.

В ИПС РАН также проводились работы по исследованию возможностей активного хранения с использованием ФС Lustre, был получен прирост производительности [53].

## Заключение

Сделанная выборка ОБОД приложений и определяемых ими требований к высокопроизводительным СХД позволяет предположить, что современных мощностей хранения недостаточно для удовлетворения запросов развивающейся науки. Количество данных возрастает быстрее, чем современные мощности хранения могут позволить сохранить и поддержать, в связи с чем приходится применять различные техники сжатия, фильтрации или просто удаления уже исследованных данных о проведенных экспериментах, что порождает риск потери ценной информации, которая могла бы пригодиться для дальнейших исследований.

Многие проекты, работающие со значительным объемом данных, называют цифры в сотни петабайт данных, в некоторых запросы доходят до эксабайта.

Увеличение объемов СХД приводит к проблемам надежности, производительности, изменению концепций работы с данными и метаданными.

Высокопроизводительные системы хранения данных, способные обеспечить пропускную способность выше 100 ГБ/с и объем более 1 ПБ, должны в ближайшем будущем войти в нашу жизнь как «стандартные решения». Подобные мощности и объемы могут потребоваться для проведения маркетинговых исследований, создания фильмов, отслеживания социальных движений и т.п.

Для задач, решаемых передовыми отраслями науки, подобных мощностей уже недостаточно. Для К компьютера исследователи ставят себя цели в сотни петабайт, разработчики эксафлопсного проекта также называют цифры от 500–1000 ПБ с пропускной способностью 30–60 ТБ/с.

Тем не менее, ведущие системы из актуального рейтинга Top500 сегодня обладают объемом систем хранения до 20 ПБ, при этом в первой десятке можно увидеть системы с хранилищем, не превышающим 2 ПБ.

Основные проблемы остаются в области системного ПО, которое нуждается в увеличении пределов масштабирования, изменении ряда концепций, поддержки новых архитектурных решений. Ряд перспективных проектов в этом направлении позволяет предположить, что в ближайшем будущем большую популярность получат параллельные файловые системы с объектной моделью хранения данных. Многие исследователи видят перспективным развитие технологий активного хранения, выполнения хотя бы части операций по предобработке сырых данных непосредственно на узлах хранения, используя доступные вычислительные мощности.

Таким образом, можно выделить три наиболее перспективные тенденции в области высокопроизводительного хранения данных.

1. Изменение модели хранения данных (Объектная модель, приближение метаданных к данным, абстрагирование таблиц размещения);

2. Иерархическая система хранения данных (с уровнем SSD-накопителей);

3. Использование новых концепций (активное хранение, приближение части вычислений к местам хранения данных).

Вероятно также, что развитию высокопроизводительных СХД могло бы способствовать создание рейтинга, в чем-то аналогичного Graph 500, позволяющего сравнивать между собой высокопроизводительные системы хранения и анализа данных.

## Литература

1. McKinsey & Company: Big data: The next frontier for innovation, competition, and productivity, URL: [http://www.mckinsey.com/mgi/publications/big\\_data/pdfs/MGI\\_big\\_data\\_full\\_report.pdf](http://www.mckinsey.com/mgi/publications/big_data/pdfs/MGI_big_data_full_report.pdf)
2. IBM's Top Storage Predictions for 2011, Январь 2011, StorageNewsletter.com, URL: <http://www.storagenewsletter.com/news/miscellaneous/ibm-top-storage-predictions-for-2011>
3. Avatar takes 1 petabyte of storage space, Январь, 2010, <http://www.devilsduke.com/avatar-takes-1-petabyte-of-storage-space/608/>
4. Facebook has the world's largest Hadoop cluster!, май 2010, URL: <http://hadoopblog.blogspot.com/2010/05/facebook-has-worlds-largest-hadoop.html>
5. The Fourth Paradigm: Data-Intensive Scientific Discovery, 2009, URL: <http://research.microsoft.com/en-us/collaboration/fourthparadigm>
6. Goble, C. and De Roure, D.: The impact of workflow tools on data-centric research. In: Data Intensive Computing: The Fourth Paradigm of Scientific Discovery, 2009.
7. Data Intensive Computing, <http://dicomputing.pnnl.gov/>
8. The Graph 500 list, URL: <http://www.graph500.org/index.html>
9. Loek Essers: CERN pushes storage limits as it probes secrets of universe, URL: <http://news.idg.no/cw/art.cfm?id=FF726AD5-1A64-6A71-CE987454D9028BDF>
10. University of Hawaii: World's Largest Digital Camera Installed on Maui Telescope, август, 2007, URL: [http://www.ifa.hawaii.edu/info/press-releases/GPC/gigapixel\\_camera-8-07.html](http://www.ifa.hawaii.edu/info/press-releases/GPC/gigapixel_camera-8-07.html)
11. Duncan Brown, Syracuse University: Computational Challenges in Gravitational Wave Astronomy, URL: <http://www.psc.edu/data-analytics/proceedings/BrownSlides.pdf>

12. Hacker T. J., Eigenmann, R., Irfanoglu, A., Pujol, S., Rathje, E., Catlin, A., Bahchi, S.: Developing an Effective Cyberinfrastructure for Earthquake Engineering: The NEEShub, In IEEE Computing in Science & Engineering, 2011 (Invited Paper.)
13. Arthur W. Wetzel, Greg Hood: Connectomics: Challenges in Reconstructing Neural Circuitry from Massive Serial Section Electron Microscopy Datasets; Data-Intensive Analysis, Analytics and Informatics TeraGrid/Blue Waters Symposium, Апрель 2011, URL: <http://www.psc.edu/data-analytics/proceedings/WetzelSlides.pdf>
14. Dell | Terascale HPC Storage Solution, URL: <http://www.terascale.com/dell-terascale-hss.html>
15. Xyratex — Advancing Digital Storage Innovation, URL: <http://www.xyratex.com/>
16. TOP-500 supercomputer sites, URL: <http://top500.org/>
17. Shinji Sumimoto: An Overview of Fujitsu's Lustre Based File System, Apr.12 2011, URL: [http://www.olcf.ornl.gov/wp-content/events/lug2011/4-12-2011/230-300 Shinji\\_Sumimoto\\_LUG2011-FJ-20110407-pub.pdf](http://www.olcf.ornl.gov/wp-content/events/lug2011/4-12-2011/230-300 Shinji_Sumimoto_LUG2011-FJ-20110407-pub.pdf)
18. Tianhe-1 Pflor Supercomputer, URL: <http://nsc-tj.gov.cn/en/show.asp?id=191>
19. Arthur S. Bland, Ricky A. Kendall, Douglas B. Kothe, James H. Rogers, Galen M. Shipman, Oak Ridge National Laboratory: Jaguar: The World's Most Powerful Computer, CUG 2009 Proceedings, URL: <http://www.nccs.gov/wp-content/uploads/2010/01/Bland-Jaguar-Paper.pdf>
20. Satoshi Matsuoka: TSUBAME2.0: A Tiny and Greenest Petaflops Supercomputer, Nov 2010, URL: [http://www.nvidia.com/content/PDF/sc\\_2010/theater/Matsuoka\\_SC10.pdf](http://www.nvidia.com/content/PDF/sc_2010/theater/Matsuoka_SC10.pdf)
21. Garth Gibson: Data Systems @ Scale, Carnegie Mellon University, 9 февраля 2011, URL: <http://www.cs.cmu.edu/~pl/CNOSSG/Gibson-CNOSSG-Feb9.pdf>
22. Pleiades Supercomputer, NAS Division Website, URL: <http://www.nas.nasa.gov/hecc/resources/pleiades.html>
23. Hopper, National Energy Research Scientific Computing Center (NERSC), URL: <http://www.nersc.gov/users/computational-systems/hopper/>
24. Peter Sayer: Bull Bills Tera 100 as Europe's First Petaflop Computer, IDG News, Май 2010, URL: [http://www.pcworld.com/businesscenter/article/197454/bull\\_bills\\_tera\\_100\\_as\\_europes\\_first\\_petaflop\\_computer.html](http://www.pcworld.com/businesscenter/article/197454/bull_bills_tera_100_as_europes_first_petaflop_computer.html)
25. Brent Welch: Exascale Distributed File Systems, MSST, Май, 2010, URL: <http://storageconference.org/2010/Presentations/MSST/8.Welch.pdf>
26. Kraken, National Institute for Computational Sciences (NICS), URL: <http://www.nics.tennessee.edu/computing-resources/kraken>
27. N. Attig, F. Berberich, U. Detert, N. Eicker, T. Eickermann, P. Gibbon, W. Gurich, W. Homberg, A. Illich, S. Rinke, M. Stephan, K. Wolkersdorfer, and T. Lippert: Entering the petaflop-era - new developments in supercomputing. In G. Munster, D. Wolf, and M. Kremer, editors, NIC Symposium 2010, volume 3, pages 1-12. IAS Series, 2010
28. MSU SUPERCOMPUTERS: «LOMONOSOV», URL: <http://hpc.msu.ru/?q=node/59>
29. BlueGene/L Configuration, Lawrence Livermore National Laboratory, URL: [https://asc.llnl.gov/computing\\_resources/bluegenel/configuration.html](https://asc.llnl.gov/computing_resources/bluegenel/configuration.html)
30. Jing Fu, Ning Liu: Scalable Parallel I/O Alternatives for Massively Parallel Partitioned Solver Systems, URL: <http://cmes.colorado.edu/courses/hpc/ipdps-lssp-parallel-io-04-23-2010-1.ppt>
31. BGW, TOP-500 supercomputer sites, URL: <http://top500.org/system/7466>
32. Clint Boulton: IBM: The Power of Purple, Mapr, 2006, URL: <http://www.internetnews.com/ent-news/article.php/590236/IBM-The-Power-of-Purple.htm>
33. Columbia, TOP-500 supercomputer sites, URL: <http://top500.org/system/7288>
34. Peter Bojanic: LUSTRE ROADMAP and FUTURE PLANS, Sun HPC Consortium, Июнь, 2008, URL: [http://www.hpcuserforum.com/presentations/Tucson/SUN%20%20Lustre\\_Update-080615.pdf](http://www.hpcuserforum.com/presentations/Tucson/SUN%20%20Lustre_Update-080615.pdf)
35. Jerry D. Smith II: Thunderbird Capacity Computing System, Sandia National Laboratories, May 3, 2006, URL: <http://www.linuxclustersinstitute.org/conferences/archive/2006/PDF/ThunderbirdUpdate.pdf>
36. Syuuichi Ihara: TOKYO TECH TSUBAME GRID STORAGE IMPLEMENTATION, Sun BluePrints™ On-Line, May 2007, Part No 820-2187-10, Revision 1.0, 5/22/07, URL: <http://www.filibeto.org/sun/lib/blueprints/820-2187.pdf>
37. Red Storm upgrade lifts Sandia supercomputer to 2nd in world, but 1st in scalability, say researchers, ноябрь, 2006, URL: <http://share.sandia.gov/news/resources/releases/2006/red-storm.html>
38. MareNostrum, TOP-500 supercomputer sites, URL: <http://top500.org/system/8242>
39. Jaguar, TOP-500 supercomputer sites, URL: <http://top500.org/system/7938>
40. Robin Goldstone: The Roar of Thunder: LLNL Goes Itanium in a Big Way, Lawrence Livermore National Laboratory, Presented to Gelato.org, Май, 2004, UCRL-PRES-204277, URL: [http://www.gelato.org/pdf/Illinois/gelato\\_IL2004\\_goldstone\\_llnl.pdf](http://www.gelato.org/pdf/Illinois/gelato_IL2004_goldstone_llnl.pdf)
41. ASCI White, URL: [https://computation.llnl.gov/casc/sc2001\\_fliers/ASCI\\_White/ASCI\\_White01.html](https://computation.llnl.gov/casc/sc2001_fliers/ASCI_White/ASCI_White01.html)
42. ASCI Red, TOP-500 supercomputer sites, URL: <http://www.top500.org/system/4428>
43. Mark Seager: An ASCI Terascale Simulation Environment Implementation, UCRL-JC-134806 PREPRINT, Mannheim Supercomputer '99 Conference, June 11, 1999, URL: <https://e-reports-ext.llnl.gov/pdf/235862.pdf>
44. Overview of the Advanced Simulation and Computing Program (ASCI), UKHEC, URL:

- <http://www.ukhec.ac.uk/publications/reports/asci.pdf>
45. Roger Haskin: Exascale Storage Challenges, 2010, IBM Corp, URL: <http://institute.lanl.gov/hec-fsio/conferences/2010/resentations/day3/Haskin-HECFsIO-2010-ExascaleChallenges.pdf>
  46. Peter Braam: Exascale File Systems, Scalability in ClusterStor's Colibri System, 2010, URL: [http://www.teratec.eu/forum\\_2010/Presentations/A5\\_Braam\\_ClusterStor\\_Forum\\_Teratec\\_2010.pdf](http://www.teratec.eu/forum_2010/Presentations/A5_Braam_ClusterStor_Forum_Teratec_2010.pdf)
  47. Rob Ross: Storage in an Exascale World, Argonne National Laboratory, URL: <http://storageconference.org/2010/Presentations/SNAPI/1.Ross.pdf>
  48. Yutaka Ishikawa: Towards Exascale File I/O, University of Tokyo, Japan, 2009/05/21, <http://www.exascale.orgmediawiki/images/6/65/ExascaleFile-io-ishikawa071309.pdf>
  49. Mirko Lorenz: Vision Cloud: The Fact Sheet, 20.12.2010, URL: <http://www.visioncloud.eu/content.php?s=30,47>
  50. Youngjae Kim, Aayush Gupta, Bhuvan Uргаonkar, Piotr Berman, and Anand Sivasubramaniam: HybridStore: A Cost-Efficient, High-Performance Storage System Combining SSDs and HDDs, Proceedings of the the IEEE International Symposium on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS), Singapore, July 2011.
  51. N. Kämmer, S. Gerhold, A. Weggerle, C. Himpel, P. Schulthess: Pageserver: High-Performance SSD-based Checkpointing of Transactional Distributed Memory, Proceedings of the 2nd International Conference on Computing Engineering and Applications (ICCEA 2010), Bali, Indonesia, 2010.
  52. Alex Szalay: Extreme Data-Intensive Computing, The Johns Hopkins University, 19 May 2011, URL: <http://salsahpc.indiana.edu/tutorial/slides/0726/szalay-bigdata-2010.pdf>
  53. Шевчук Е. В., Тютляева Е. О., Московский А. А. 2009. Система активного хранения данных на базе библиотеки динамического распараллеливания TSim. // Научный сервис в сети Интернет: масштабируемость, параллельность, эффективность. Труды Всероссийской научной конференции, 21–26 сентября 2009 г. Новороссийск, — М.: Изд-во МГУ имени М.В. Ломоносова, 2009 с. 226–230 (CD) ISBN 978-5-211-05697-8.
  54. DKRZ brochure (2009) «The power to understand: Supercomputing for Climate System Science».
  55. <http://www.graph500.org/junc2011.html>, позиция 7.

**Тютляева Екатерина Олеговна.** Инженер Института программных систем имени А.К. Айламазяна РАН. Окончила НОУ ИПС Университет города Переславля имени А.К. Айламазяна в 2009 году. Автор более 12 научных работ. Область научных интересов: системы хранения данных, высокопроизводительный ввод-вывод, отказоустойчивость. E-mail: [ordi@xgl.pereslavl.ru](mailto:ordi@xgl.pereslavl.ru)

**Московский Александр Александрович.** Директор по науке ЗАО «РСК СКИФ». Окончил химический факультет МГУ имени М.В. Ломоносова в 1997 году. Кандидат химических наук. Автор более 20 научных работ. Область научных интересов: молекулярное моделирование, высокопроизводительные вычисления. E-mail: [moskov@rsc-skif.ru](mailto:moskov@rsc-skif.ru)