

УДК 519.6

JobDigest — подход к исследованию динамических свойств задач на суперкомпьютерных системах

А. В. Адинец¹, П. А. Брызгалов², В. В. Воеводин³,
С. А. Жуматий⁴, Д. А. Никитенко⁵, К. С. Стефанов⁶

¹adinetz@gmail.com, ²petr@parallrel.ru, ³serg@parallel.ru, ⁴dan@parallel.ru, ⁵cstef@parallel.ru

¹⁻⁵НИВЦ МГУ имени М.В.Ломоносова

Поступило в редакцию 22.12.2013

Аннотация. Рассматривается создаваемый комплекс для исследования поведения программ во время выполнения на суперкомпьютерных системах. Отдельно акцентируется внимание на одном из предложенных подходов под названием JobDigest. Данный подход предоставляет пользователю подробный отчет по завершении задачи с описанием ее динамических свойств, что позволяет осуществлять качественный анализ и повышать эффективность выполнения исследуемых приложений.

Ключевые слова. Параллельные вычисления; анализ поведения программ; кластерные системы; эффективность программ

Вместе с ростом масштаба вычислительных систем и решаемых на них задач растет и сложность написания эффективных программ. Причина этого кроется в том, что также возрастает и множество факторов, которые могут влиять на эффективность приложений. Свойства аппаратного и программного обеспечения суперкомпьютера, свойства самой исполняемой программы, взаимное влияние исполняемых программ друг на друга – все это необходимо учитывать, если стремишься добиться высокой эффективности. При этом с развитием области высокопроизводительных вычислений все эти свойства становятся только сложнее, обрастают все новыми факторами. Самые мощные современные кластеры обладают миллионами ядер; как в принципе можно организовать эффективные вычисления такого масштаба? Иерархия памяти с каждым годом становится все сложнее; как добиться ее эффективного использования?

Низкая эффективность выполняемых программ – это не единственная проблема в суперкомпьютерной области. Также важно научиться оценивать и повышать эффективность использования самого суперкомпьютера, т. е. насколько эффективно используется процессорное вре-

мя и другие ресурсы вычислительного комплекса. Наш опыт показывает, что зачастую такие механизмы как алгоритмы планирования выполнения задач и квоты на использование ресурсов суперкомпьютера приводят к снижению эффективности его работы.

Все это приводит к необходимости создания инструмента, который позволит разобраться, где, а главное, почему происходит потеря производительности при выполнении программ и использовании суперкомпьютеров. Понятно, что для получения качественной и многосторонней оценки этот инструмент должен обладать множеством различных данных как о поведении самой задачи, так и о состоянии суперкомпьютера в целом. Подобный инструмент разрабатывается в лаборатории Параллельных информационных технологий Научно-исследовательского вычислительного центра Московского государственного университета имени М.В. Ломоносова и получил название LAPTA – «Laptais a pAckage for Performance moni Toring and Analysis». Данный инструмент создается в рамках выполнения совместного российско-европейского проекта HOPSA (Holistic Performance System Analysis) [1, 2].

В данной статье мы расскажем о разрабатываемом комплексе LAPTA и подробно остановимся на одном из используемых в нем подходов, который предназначен для исследования

поведения задачи во время выполнения. Данный подход изучает динамические свойства задач, исследуемые с помощью средств мониторинга. Его цель состоит в предоставлении как администратору системы, так и пользователю базовых характеристик задачи по ее завершению для получения как качественной, так и детальной оценки каждого отдельно взятого запуска. Данный подход, а также полученный в результате его применения отчет получили название «Jobdigest».

1. АРХИТЕКТУРА КОМПЛЕКСА LARPA

Комплекс LARPA предназначен для всестороннего анализа динамических характеристик параллельных программ и суперкомпьютеров. При проведении анализа учитываются характеристики задачи, начиная от момента постановки в очередь и заканчивая ее завершением. Это позволяет получать полную информацию как о самой задаче, так и о всей совокупности выполняемых на суперкомпьютере задач.

Общая схема комплекса LARPA представлена на рис. 1. На узлах кластера данные с различных аппаратных датчиков собираются при помощи агентов и передаются на уровень агрегации. При помощи соответствующих агентов данные собираются и с системы управления потоком задач (СУПЗ). Модули агрегации сохраняют собранные данные в базы данных при помощи модулей БД. Каждый модуль БД поддерживает работу с одним типом БД. Далее данные из БД могут быть извлечены для анализа при помощи модулей анализа.

Задача модулей анализа – это выполнение обработки данных. При этом выполнение анализа может быть инициализировано из самых разнообразных частей системы. Например, запрос на анализ может прийти из модуля визуализации через центр обработки запросов (ЦОЗ) в результате работы пользователя с системой через веб-браузер. В этом случае основной целью запроса может быть анализ динамики поведения параллельной программы. Анализ данных может быть инициирован системными процессами, например, по таймеру раз в сутки для построения ежедневного отчета о работе суперкомпьютера. Или же данные могут быть запрошены внешними системами интеграции и визуализации. В любом случае инициатор запроса сначала обращается к серверу управления, который при необходимости порождает модули анализа и отслеживает их работу. Через сервер

управления данные не идут – их возвращает модуль анализа напрямую инициатору запроса.

2. ОРГАНИЗАЦИЯ ПОДХОДА JOB DIGEST

В данном разделе мы опишем, каким образом в рамках комплекса LARPA организована работа модуля, реализующего подход Jobdigest.

В момент постановки задачи на счет происходит либо запуск агента, который будет отслеживать изменение статуса задачи в потоке служебной информации от системы управления потоком задач, либо передается идентификатор задачи, если такой агент уже запущен. По окончании работы задачи агент выделяет в потоке вывода СУПЗ признак ее завершения и предоставляет базовую информацию (списки узлов, временной диапазон и т. п.) для дальнейшего формирования отчета Jobdigest. Подробнее о самом формировании отчета будет рассказано отдельно.

Система управления потоком задач ставит задачу на счет и передает на вывод изменения в статусе задачи. В данный момент поддерживаются СУПЗ CLEOи Slurm. Одновременно в постоянном режиме работает система сбора системных данных, данные сохраняются, возможно, с некоторой степенью агрегации для дальнейшего апостериорного анализа. Помимо основных данных самой системы мониторинга могут собираться данные и от других источников, например, от существующих средств сбора трасс-приложений. Для этого предусмотрены соответствующие интерфейсы в системе сбора и сохранения данных статистики.

Данный подход предполагает наличие инфраструктуры для доступа к сохраненным данным системного мониторинга. В частности, такая инфраструктура разрабатывается российской стороной в рамках упомянутого выше совместного проекта HOPSA. Например, разрабатываемый в рамках проекта язык Hoplang [3,4] предназначен для формирования запросов данных мониторинга.

На рис.2 схематично показан процесс формирования отчета на общей схеме комплекса сбора, хранения и обработки данных мониторинга. В качестве СУБД для хранения данных мониторинга используются базы данных Cassandra и MongoDB. Модуль доступа к БД реализован в двух вариантах: на технологиях Pig и Hadoop и на упомянутой выше технологии Hoplang.

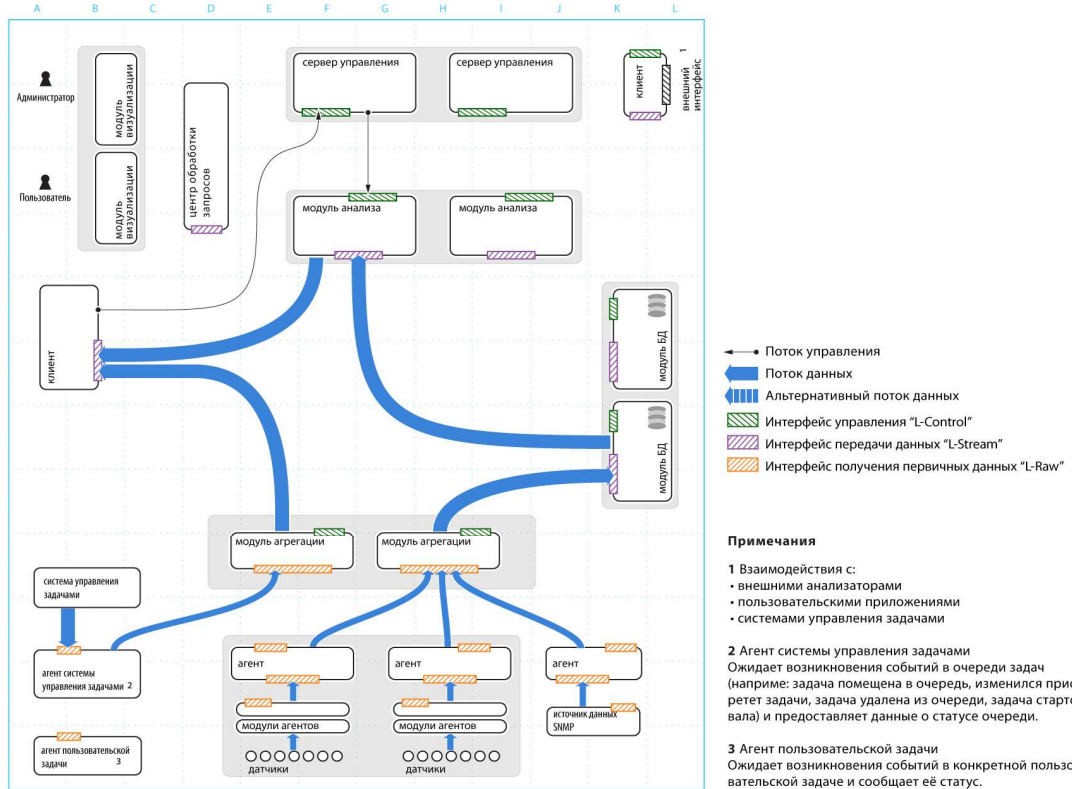


Рис. 1. Архитектура комплекса LAPTA

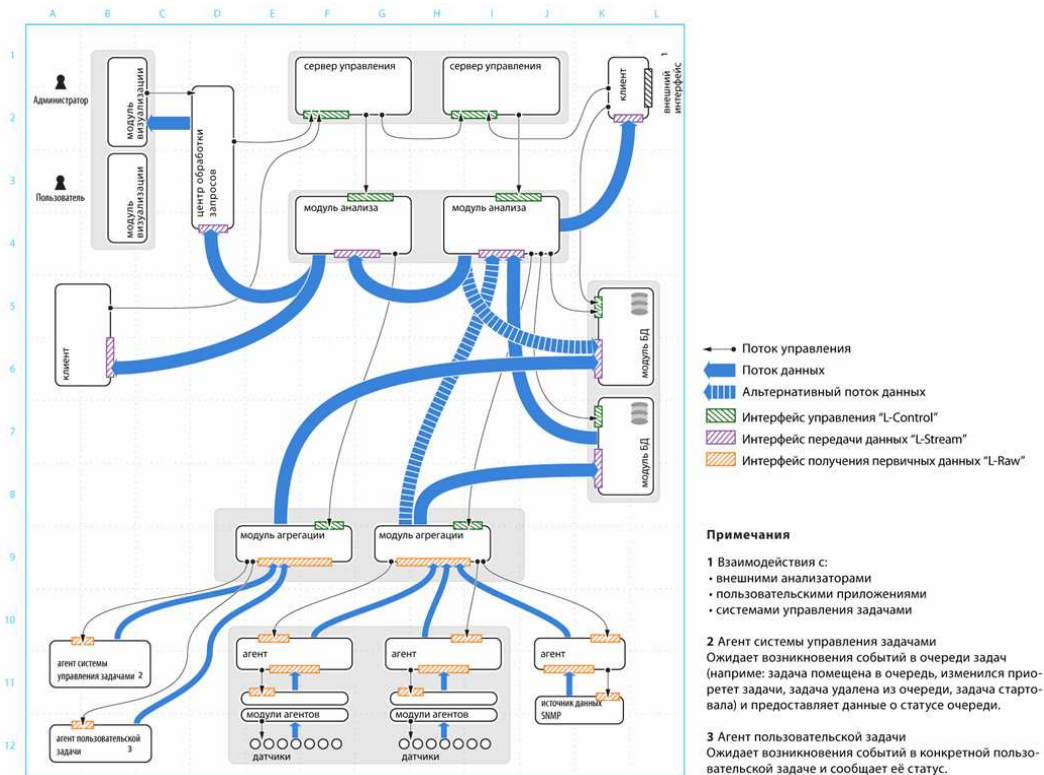


Рис. 2. Механизм создания Jobdigest отчета на общей схеме комплекса LAPTA

И модуль, использующий технологию Pig+Hadoop, и модуль, использующий Hoplang, написаны на языке программирования Ruby. Для описания запросов на получение данных в первом случае используется язык Piglatin, а во втором – Hoplang.

3. ФОРМИРОВАНИЕ ОТЧЕТА JOB DIGEST

В текущей реализации отчет Jobdigest не генерируется автоматически; он создается только в случае необходимости по запросу пользователя. Для составления отчета по определенному событию вызывается скрипт, который формирует запрос к Java-сервлету. По умолчанию запуск скрипта производится в ручном режиме, но он может производиться, например, по окончании работы задачи.

В запросе указываются характеристики завершенной задачи, а также шаблон для формирования отчета. Если Jobdigest уже был ранее сформирован (ответ от Java-сервлета положителен), то информация об отчете записывается в HTML-файл, который потом доступен через Web-браузер. Адрес этого отчета охраняется в отдельном файле вместе с отчетом о задаче, и пользователь может впоследствии открыть его в браузере.

Javaсервлет выполняет набор запросов с параметрами, и полученные данные записываются в CSV-файлы. Набор запросов и параметры к ним указываются в шаблоне, который представляет собой HTML-файл со ссылками на запросы. Сами запросы хранятся в текстовых

файлах, называемых «шаблоны запросов». Запросы выполняются по очереди в том порядке, в каком они расположены в HTML-шаблоне. По окончании обработки запросов генерируется отчет. Он получается из HTML-шаблона путем замены ссылок на шаблоны запросов специальными ссылками на CSV-файлы с результатами. Эти ссылки строятся таким образом, что при открытии отчета в браузере на месте ссылок отображаются диаграммы и таблицы с данными, полученными из соответствующих CSV-файлов.

Общая последовательность обращений внутри модуля визуализации приведена на рис. 3.

Web-сервер в терминах схемы на рис. 4 реализован на языке программирования Java в виде набора сервлетов, работающих под управлением web-сервера Tomcat. Web-сервер общается с модулем БД по протоколу AVRO по технологии удаленного вызова процедур (RPC). Модуль БД принимает запросы от web-сервера и возвращает данные в формате CSV.

Полученные от модуля БД данные web-сервер дополнительно обрабатывает. Для этого используется библиотека GoogleDataSource Library, которая позволяет использовать язык запросов Googlequerylanguage для постобработки полученных от модуля БД данных. Также эта библиотека используется для форматирования данных перед их записью в CSV файлы. Визуализация данных происходит в web-браузере. Для визуализации используется JavaScript и библиотеки jQuery и Highcharts.

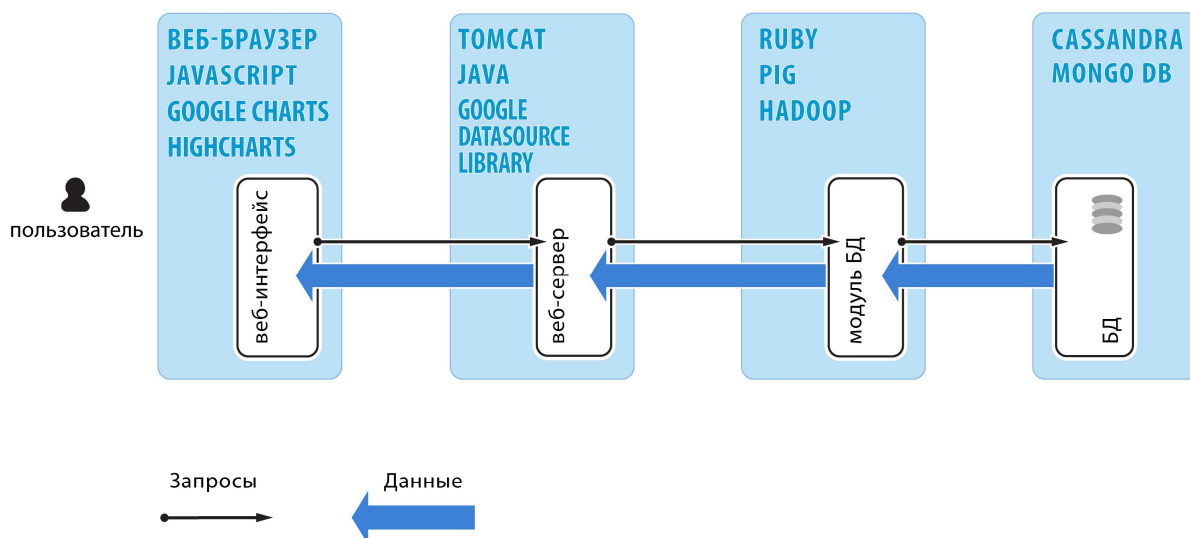


Рис. 3. Последовательность обращений при построении Jobdigest, используемые технологии

Информация о задачах пользователя "username"

214455(namd2) [reports/auxiliary/214455/min_temp.html](#) [reports/auxiliary/214455/job_info.html](#) [передать отчет](#)

Строка запуска: /home/username/exec/namd/namd2_1ei5_196-114H222H227-220T291__D-Ala-pNa.md.CONTINUE.conf
 Число ядер: 8
 Номера узлов: node-39-03
 Дата постановки в очередь: Fri, 18 May 2012 14:19:52
 Дата запуска: Fri, 18 May 2012 14:19:56
 Дата окончания счета: Fri, 18 May 2012 14:20:05
 Время счета: 0 days 0 hours 0 minutes 9 seconds
 Время ожидания: 0 days 0 hours 0 minutes 4 seconds
 Количество процессорочасов(ядра*часы): 0.02

214458(namd2) [reports/auxiliary/214458/min_temp.html](#) [передать отчет](#)

bigmem-1337342936-1942(namd2) [reports/auxiliary/bigmem-1337342936-1942/min_temp.html](#) [передать отчет](#)

bigmem-1337343517-1943(namd2) [reports/auxiliary/bigmem-1337343517-1943/min_temp.html](#) [передать отчет](#)

bigmem-1337344245-1944(namd2) [создать отчет](#)

hdd-1337333568-64986(namd2) [reports/auxiliary/hdd-1337333568-64986/min_temp.html](#) [передать отчет](#)

hdd-1337343218-65005(namd2) [создать отчет](#)

bigmem-1337344708-1947(namd2) [reports/auxiliary/bigmem-1337344708-1947/min_temp.html](#) [передать отчет](#)

regular-1337455929-215378(namd2) [создать отчет](#)

regular-1337456097-215394(namd2) [reports/auxiliary/regular-1337456097-215394/job_info.html](#) [передать отчет](#)

Рис. 4. Страница с информацией по запускам задач пользователем username

Динамические характеристики запусков приложений пользователей доступны администраторам системы в полном объеме, а обычным пользователям доступ предоставляется только к запускам собственных приложений. Пользователям предлагается страница (рис. 4) со списком их приложений, завершивших работу на вычислительной системе. Возле каждого приложения в списке находится ссылка для формирования отчета по работе данного приложения, если он не был сформирован ранее, и ссылка на сам отчет, если он уже был однажды посчитан. По окончании формирования отчета пользователь может просмотреть его в браузере. Основу отчета составляет информация о данных мониторинга выбранного приложения. В отчет можно включить любое количество графиков и диаграмм, отражающих различные параметры работы приложения – загрузка процессора (рис. 5), использование памяти (рис. 6), свопинг, LoadAverage (рис. 6), дисковые операции, загруженность сети Ethernet (рис. 5), InfiniBand и т. д.

Данные графики уже позволяют делать содержательные выводы о поведении программы. Например, на рис. 5 видно, что загрузка CPU в начале выполнения программы практически нулевая. При этом сеть Ethernet не загружена (аналогичная картина и с сетью Infiniband), а это, скорее всего, означает, что в начале программы происходит чтение входных данных. Из этого пользователь может сделать вывод, что начальная загрузка данных занимает достаточно много времени, и может попытаться оптимизировать соответствующий фрагмент программы.

ЗАКЛЮЧЕНИЕ

Разработан и проходит апробацию удобный и эффективный инструмент для качественного анализа каждого отдельного запуска задачи, не требующий специальной подготовки задачи перед постановкой на счет. Это позволит существенно улучшить понимание пользователем того, как задача выполнялась на конкретной системе в реальных условиях. В частности, выделить аномалии, вызванные особенностями параллельной программы, или локализовать по времени или вычислительным узлам влияние других приложений. Во втором случае администратор потенциально сможет сократить это влияние, внося соответствующие изменения в политики распределения задач или обнаружив и устранив какой-то дефект в настройке системного программного обеспечения.

Подход, безусловно, будет развиваться далее, не только совершенствуя текущий функционал, но и расширяясь качественно новыми возможностями. В частности, в подход изначально заложена идея автоматического выделения потенциально узких мест в конкретном запуске задачи, что позволит даже в автоматическом режиме обратить внимание пользователя на потенциальную проблему, а в каких-то случаях и рекомендовать инструментарий для дальнейшего изучения приложения. Автоматическая обработка данных мониторинга и выдача рекомендаций очень важны для пользователей суперкомпьютера, которые могут быть не очень знакомы с тонкостями выполнения задач на суперкомпьютерных системах, а потому самостоятельно могут не заметить какой-то характерный признак неэффективности.



Рис. 5. Профиль загрузки CPU и сети Ethernet

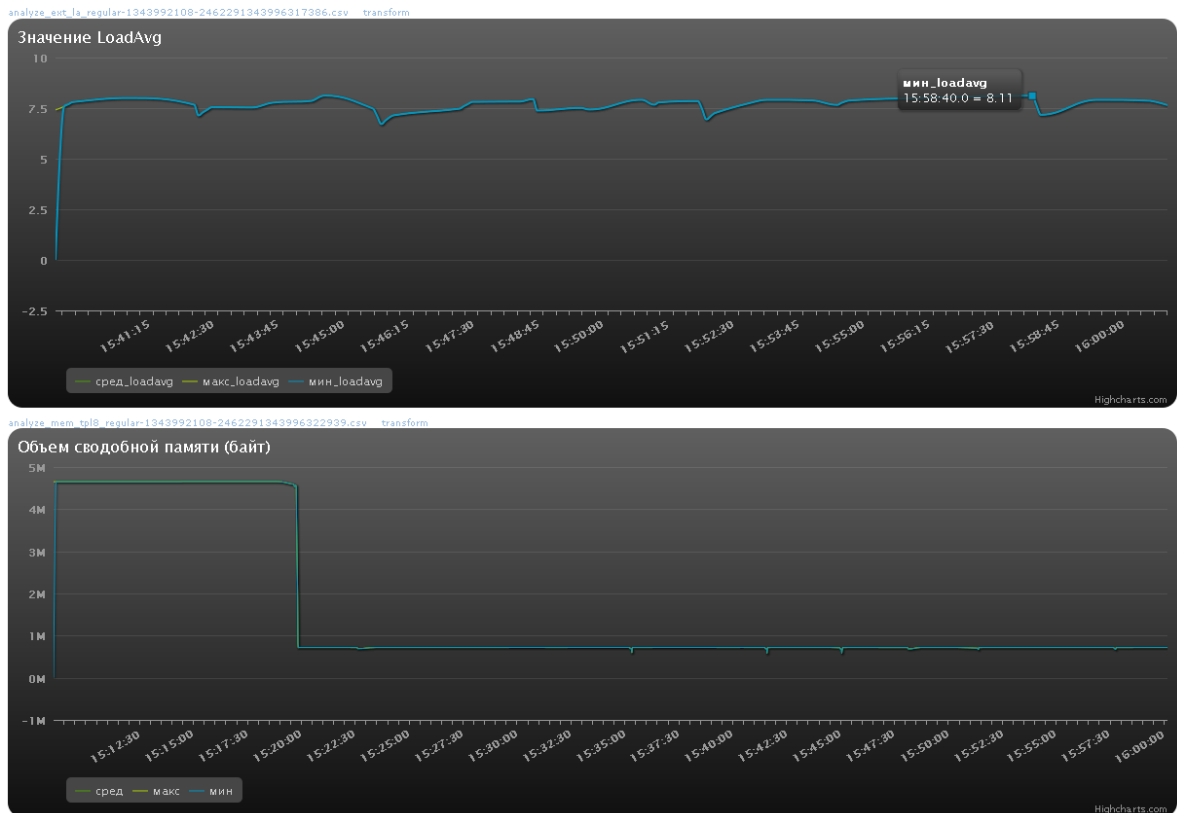


Рис. 6. Профиль использования памяти и LoadAverage

Для дальнейшего комплексного изучения динамических свойств хода выполнения программ в разрабатываемой системе HOPSA предусматриваются интерфейсы взаимодействия (например, через расширенную трассу OTF2) с внешними анализаторами, такими как Vampir, Scalasca и т. п.

Апробация описанного разработанного подхода проводится в Суперкомпьютерном комплексе МГУ имени М.В. Ломоносова.

СПИСОК ЛИТЕРАТУРЫ

1. Разработка методов и инструментальных систем для анализа эффективности работы параллельных программ и суперкомпьютеров (официальный сайт российской части совместного проекта HOPSA) (<http://hopsa.parallel.ru>).

2. HOlistic Performance System Analysis (EU HOPSA website) (<http://vi-hps.org/projects/hopsa/>).

3. Адинец А. В., Жуматий С. А., Никитенко Д. А. Hoplang – язык обработки потоков данных мониторинга // Параллельные вычислительные технологии-2012: сб. тр. междунаrodn. науч. конф. (ПаВТ'12), 2012. С. 351–359.

4. Язык обработки потоков данных для систем кластерного мониторинга Hoplang. URL: <http://github.com/zhum/hoplang>

ОБ АВТОРАХ

Адинец А.В. м.н.с., к.ф.-м.н., НИВЦ МГУ, adinetz@gmail.com

Брызгалов П.А. н.с., к.ф.-м.н., НИВЦ МГУ, pyotr777@guru.ru

Воеводин Вад.В м.н.с., к.ф.-м.н., НИВЦ МГУ, vadim@parallel.ru

Жуматий С.А. с.н.с., к.ф.-м.н., НИВЦ МГУ, serg@parallel.ru

Никитенко Д.А. м.н.с., НИВЦ МГУ, dan@parallel.ru.

METADATA

Title: Job Digest – approach to jobs dynamic properties investigation on supercomputer systems

Authors: A.B. Adinetz¹, P.A. Bryzgalov², S.A. Zhumatiy³, D.A. Nikitenko⁴, K.S. Stefanov⁵

Affiliation: ¹⁻⁵Science Research Computer Center (SRCC) of Moscow State University (MSU), Russia.

Email: ¹adinetz@gmail.com, ²petr@parallel.ru, ³serg@parallel.ru, ⁴dan@parallel.ru, ⁵stef@parallel.ru

Language: Russian.

Source: Vestnik UGATU (Scientific journal of Ufa State Aviation Technical University), 2013, Vol. 17, No. 2 (55), pp. 131-188. ISSN 2225-2789 (Online), ISSN 1992-6502 (Print).

Abstract: Realization of the Web OLAP providing formation of hypercubes "on the fly" from situation-oriented database (SODB) is discussed. The architecture of an OLAP-application on the basis of SODB is considered. The database ER-model as a basis of the conceptual multidimensional model which is setting a set of potential hypercubes is used. Design of hypercubes dimensions and measures are discussed. The approach is illustrated on an example of multidimensional activity model for dissertational councils of scholar institution.

Key words: Web OLAP; multidimensional data model; ER-model; situation-oriented database.

References (English Transliteration):

1. Methods and tools for parallel programs and supercomputers efficiency analysis development (official website of russian part HOPSA project) - <http://hopsa.parallel.ru>

2. HOlistic Performance System Analysis (EU HOPSA website) - <http://vi-hps.org/projects/hopsa/>

3. A.B. Adinetz, S.A. Zhumatiy, D.A. Nikitenko Hoplang — language for monitoring data streams processing // "Parallel Computing Technologies 2012" (PaCT'12) conference proceedings, 2012, s. 351–359.

4. Hoplang, data streams processing language for clusters monitoring. URL: <http://github.com/zhum/hoplang>

About authors:

1. A.B. Adinetz junior researcher, p.h.d., SRCC MSU, adinetz@gmail.com

2. P.A. Bryzgalov researcher, p.h.d., SRCC MSU, pyotr777@guru.ru

3. Vad.V. Voevodin junior searcher, p.h.d., SRCC MSU, vadim@parallel.ru

4. S.A. Zhumatiy leading researcher, p.h.d., SRCC MSU, serg@parallel.ru

5. D.A. Nikitenko junior searcher, м.н.с., SRCC MSU, dan@parallel.ru