Опыт использования GRID-технологий в системе обработки данных Спутникового Центра ДВО РАН

П.В. Бабяк, Г.В. Тарасов

Институт автоматики и процессов управления ДВО РАН 690041 Владивосток, Радио 5
E-mails: vanger.paul@gmail.com, george@dvo.ru

В работе рассмотрены принципы интеграции больших вычислительных ресурсов в систему обработки спутниковых данных с использованием GRID-технологий. На примере двух Центров Коллективного Пользования (ЦКП) ДВО РАН построена распределенная GRID-сеть для проведения расчетов полей доминантных ориентацией термических контрастов (ДОТК). Экспериментальные расчеты показали многократное увеличение скорости обработки, что особенно актуально в условиях круглосуточного мониторинга. Полученный успешный опыт позволяет использовать построенную GRID-сеть для различных задач обработки спутниковых данных, а не только для задачи расчетов полей ДОТК.

Работа выполнена по программам Президиума РАН №1 и №2 и поддержана грантом РФФИ №08-07-227.

Ключевые слова: ГРИД-вычисления, параллельные вычисления, распределённые системы, ДОТК.

Введение

Общеизвестно, что задачи обработки данных дистанционного зондирования (ДДЗ) Земли относятся к классу ресурсоемких задач. Многими авторами [1, 2] неоднократно отмечалось, что с каждым годом увеличивается объем принимаемой и обрабатываемой информации, полученной со спутников. Также существенно увеличиваются требования, предъявляемые к результатам обработки спутниковых данных, как следствие увеличивается сложность алгоритмов, повышается их вычислительная емкость. Перспективным направлением развития систем обработки ДДЗ является использование высокопроизводительной вычислительной техники и соответствующего параллельного программного обеспечения. Однако это направление не всегда удается эффективно реализовать на практике, так как центрам обработки спутниковых данных зачастую сложно иметь собственные мощные вычислительные ресурсы. Логичным выходом в этой ситуации является использование сторонних ресурсов. Одной из таких технологий, позволяющей организовать распределенную вычислительную систему большой мощности, является GRID.

Примером вычислительно ёмкого алгоритма обработки спутниковых данных служит задача построения полей доминантных ориентацией термических контрастов (ДОТК). Суть алгоритма состоит в определении для некоторой окрестности точки величины ДОТК и его статистической значимости [3]. В зависимости от параметров настройки алгоритма его усредненное время счета для одного стандартного региона (акватория Охотского моря) составляет более 12 часов. В условиях круглосуточного мониторинга окружающей среды для этой задачи не представляется возможным организовать оперативную обработку данных без привлечения мощных параллельных вычислительных комплексов. Так как значение ДОТК в точке зависит только от температур точек окрестности, а размер окрестности (точный или максимальный, в зависимости от способа вычисления ДОТК) известен до начала вычислений, то можно использовать простой метод распараллеливания по данным. При этом были созданы программные средства, разделяющие исходные изображения на примерно равные части с наложением частей друг на

друга, и объединяющие результатов работы алгоритма. Входом и выходом алгоритма являются файлы с исходным изображением и результатом расчетом соответственно.

В данной работе на примере задачи расчета полей ДОТК рассматриваются принципы интеграции больших вычислительных ресурсов в систему обработки спутниковых данных с использованием GRID-технологий. Результатом работы является распределенная вычислительная GRID-сеть, построенная на базе двух центров коллективного пользования ДВО РАН: ЦКП регионального спутникового мониторинга окружающей среды (Центр РСМ) и ЦКП «Дальневосточный Вычислительный Ресурс» (ДВВР). Центр РСМ является научно-исследовательским учреждением, осуществляющим исследовательскую и информационную деятельность в целях изучения природных и техногенных процессов в Дальневосточном регионе. Центр РСМ базируется на распределенной системе обработки данных. Система обеспечивает доступ к ключевым источникам информации (в различных форматах), надежно функционирует в режиме реального времени и позволяет пользователю управлять процессом обработки. ЦКП ДВВР организован на базе суперкомпьютерного центра ИАПУ ДВО РАН и является одним из научно-технических подразделений ДВО для развития и поддержки параллельных вычислений в сфере науки и образования. В настоящее время суммарная производительность вычислительных ресурсов ДВВР оценивается в 1 терафлоп.

Ресурсы GRID-сети ЦКП ДВВР

В настоящее время основные ресурсы ДВВР объединены в опорный GRID-узел. Структура узла предусматривает его всестороннее расширение, как по объему интегрируемых ресурсов, так и по количеству подключаемых пользователей. Описание ресурсов ДВВР и краткая схема их взаимосвязи показана на рисунке 1. В рамках центра функционирует три вычислительных кластера различной архитектуры, предназначенные для решения различного класса задач. Кластер МВС1000/16 ориентирован на выполнение учебных и тестовых задач с небольшим объемом данных и временем счета и чаще всего используется для отладочных целей. Кластеры МВС15К и МВС1000/17 образуют основную вычислительную мощность центра и ориентированы на проведение длительных, массивных расчетов. Помимо вычислительной техники в составе ДВВР функционирует несколько вспомогательных серверов. Сервер доступа обеспечивает централизованный доступ пользователей к ресурсам центра. Файл-сервер предоставляет дисковое пространство для хранения пользовательской информации, исходных данных и обеспечивает результатов вычислений. Интернет-сервер информационную пользователей и осуществляет мониторинг производительности и доступности всех ресурсов. ДВВР имеет прямое подключение к опорной научно-образовательной сети Приморского научного центра ДВО, которая объединяет все ведущие научные и образовательные учреждения. В частности, данная сеть обеспечивает гигабитный канал связи между ДВВР и Центром РСМ. Краткая характеристика вычислительных кластеров представлена в таблице 1.

Таблица 1.	Краткая х	карактеристи	ка вычислительны	х ресурсов
------------	-----------	--------------	------------------	------------

Название кластера	Аппаратные характеристики	Производительнос ть	
MBC1000/16	16 узлов (16 процессоров)	~9 гигафлоп	
	Pentium III 800MHz, память 512Мб		
MBC1000/17	7 узлов (7 процессоров, 28 ядер)	~140 гигафлоп	
WIDC1000/17	Intel Core 2 Quad 2.33GHz, память 4Гб	- 140 гм афлоп	
MBC15K	84 узла (168 процессоров)	~920 гигафлоп	
MIDCISK	PowerPC 970FX 2.2GHz, память 4Гб	~920 гигафлоп	

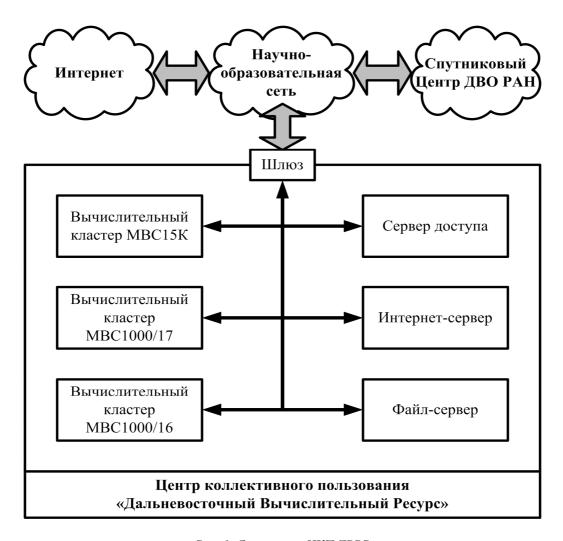


Рис. 1. Структура ЦКП ДВВР

Выполнение вычислительных задач пользователей на кластерах контролируется соответствующими локальными системами пакетной обработки заданий (СПОЗ), которые определяют политику предоставления ресурсов пользователям в рамках вычислительного кластера. В настоящее время СПОЗ реализованы на базе свободного пакета Torque (OpenPBS). Основные действия пользователя по запуску своих задач заключаются в выполнении стандартных шагов: подготовить и закачать входные данные, подготовить описание задачи (паспорт) и запустить задачу на счет, дождаться завершения и забрать результаты вычислений.

GRID-инфраструктура ЦКП ДВВР построена на базе открытого пакета «Globus Toolkit 4» (далее просто GT). Идеологической основой создания GRID является технология Web-сервисов [4, 5]. Пакет GT устанавливается и настраивается на каждом узле GRID индивидуально. Конечная настройка зависит от назначения данного узла в GRID и требуемых выполняемых функций. GT имеет модульную структуру и содержит набор связанных компонент, взаимодействующих с помощью общей платформы Web-сервисов. В структуре пакета можно выделить несколько разделов, реализующих основные исполняемые компоненты: «Common», «Security», «Data», «Execution», «Monitor». Раздел «Common» является общей платформой создания и запуска Web-сервисов и обеспечивает их взаимодействие. Раздел «Security» отвечает за безопасность передачи информации, идентификацию и авторизацию пользователей, ресурсов и сервисов. Раздел «Data» отвечает за управление и передачу файлов, организацию распределенных файловых каталогов. «Execution» отвечает за выделение вычислительных ресурсов и управление этими ресурсами. «Мопіtог» отвечает за сбор информации о ресурсах и мониторинг их работы. Типичный узел GRID всегда использует некоторый набор исполняемых компонент из перечисленных разделов. На рисунке 2 показана схема распределения Web-сервисов по ресурсам ДВВР.

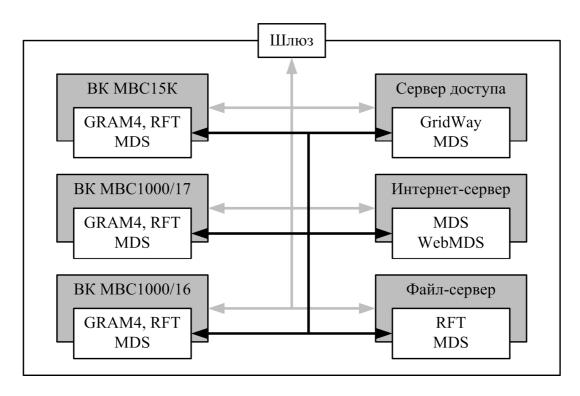


Рис. 2. Схема распределения GRID-сервисов по ресурсам ДВВР

Кратко опишем используемые Web-сервисы из пакета GT. MDS (Monitoring and Discovery System) – система обнаружения и мониторинга ресурсов. Задача данного сервиса – хранить описание состояния всех ресурсах GRID. WebMDS – это компонента доступа к описаниям ресурсов через Интернет-браузер. GRAM (Globus Resource Allocation and Management) – система распределения и управления ресурсами. Сервис GRAM работает на вычислительном ресурсе. GRAM принимает на вход паспорт описания задачи пользователя, которую нужно выполнить, связывается через адаптер с локальной СОПЗ и передает ей необходимые инструкции для запуска задачи. Дополнительно поддерживается проверка состояния выполнения задачи, ее удаление и сбор информации о состоянии СПОЗ в целом. RFT (Reliable File Transfer) - систем надежной передачи файлов. Сервер RFT реализует передачу файлов между различными узлами GRID. Клиент RFT формирует заявку серверу на передачу файла. Заявка состоит в указании двух URL на файлы: откуда и куда копировать. GridWay – основная задача этого сервиса – организовать планирование и учет использования ресурсов в рамках GRID-сети. В обязанности GridWay входит обработка заявок от пользователей на выполнение задач, выделение соответствующих ресурсов под эти задачи из числа свободных, контроль использования ресурсов и другие. В процессе своего функционирования GridWay обращается к: GRAM для осуществления запуска задачи на некотором узле GRID, RFT для передачи исходных данных и результатов, если это прописано в паспорте задачи, и MDS для получения централизованной информации о состоянии локальных систем очередей ресурсов. Четвертой компонентой GridWay является диспетчер расписания, который получает от GridWay необходимую информацию о ресурсах и заявках на использование ресурсов, и выдает названия узлов, которые в данный момент можно использовать для выполнения имеющихся задач. Предусмотрено два способа взаимодействия с GridWay: через интерфейс командной строки и через программный интерфейс DRMAA (Distributed Resource Management Application API).

Распределенная система обработки данных спутникового мониторинга

Распределенная система состоит из двух частей: подсистемы запуска сценариев обработки и подсистемы мониторинга. Принципиальная схема распределенной системы обработки представлена на рисунке 3. Подсистема запуска предназначена для выполнения сценариев

обработки (последовательностей процедур, связанных командами управления) на удаленных обрабатывающих компьютерах. Сервер управления обработкой состоит из трех основных компонентов: обработчика сообщений, пула задач и диспетчера. Его основная задача заключается в получении сценариев от инициатора и запуске их на соответствующих удаленных компьютерах. Инициатором может быть как пользователь, так и другой сценарий. Это позволяет создавать схемы для обработки на нескольких компьютерах. Для взаимодействия с сервером используются два типа сообщений: командные (сигнал о необходимости запуска сценария, сигнал аварийной остановки процедуры и т.п.) и информационные (например, сигнал о процессе или его завершении). Для выполнения сценария обработки инициатор ставит в очередь запрос на его запуск (см. рис. 3). Обработчик сообщений создает и помещает в пул задачу в соответствии с хранимом в базе данных описанием сценария и переданными параметрами, и возвращает ее идентификатор инициатору. Созданная задача из пула попадает в планировщик, который запускает ее на первом свободном компьютере, подходящем для данного сценария. Все информационные сообщения от сервера и от выполняемых задач поступают в топик, из которого могут быть прочитаны инициатором. Интерфейс системы допускает как асинхронный, так и синхронный, то есть с ожиданием завершения работы, режим запуска сценария. Протокол передачи данных, в рамках подсистемы обработки, не формализован и определяется самим сценарием. В описываемой системе для передачи данных по сети используется протоколы SMB/CIFS и FTP. Для удаленного запуска и контроля процесса обработки на стороне обрабатывающих компьютеров используется открытая реализация протокола SSH (пакет OpenSSH).

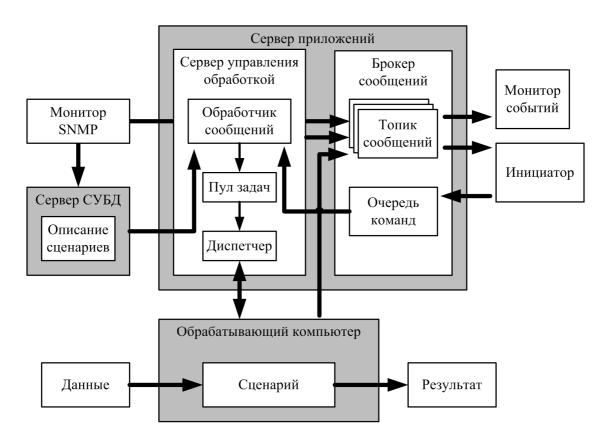


Рис. 3. Обмен сообщениями в распределенной системе обработки

Подсистема мониторинга предназначена для контроля состояния узлов (доступности, уровня загрузки процессора, наличия свободного места на дисках и т.д.) и процессов, происходящих в системе (работы диспетчера и сценариев обработки). Подсистема состоит из монитора событий, отображающего информацию посредством графического интерфейса и отдельных модулей-

мониторов, предназначенных для отслеживания различных событий. Взаимодействие между мониторами и другими подсистемами осуществляется через очередь сообщений типа издатель/подписчик. Реакция монитора на появление нового сообщения (графический индикатор, звуковой сигнал и т.п.) настраивается пользователем при помощи встроенного скриптового языка.

Интеграция ресурсов ЦКП ДВВР в систему обработки спутниковых данных

Суть интеграции сторонних для системы обработки ресурсов состоит в формировании для сервера управления обработкой абстракции «обрабатывающего компьютера», реализация которого опирается на GRID-сервисы для выполнения функций обработки данных в распределенной вычислительной сети. Таким образом, наряду с уже имеющимися компьютерами обработки, в системе обработки появляется дополнительный «компьютер», обладающий большим вычислительным потенциалом. Полезная мощность такого «GRID-компьютера» определяется, вопервых, объемом ресурсов, сосредоточенных в GRID-сети, во-вторых, текущей загрузкой этих ресурсов. Таким образом, необходимо было решить две основные задачи. С одной стороны обеспечить взаимодействие «GRID-компьютера» с сервером управления обработкой. С другой стороны обеспечить доступность распределенных вычислительных ресурсов имеющейся GRID-сети. Процесс выполнение задач был разбит на три основных этапа: подготовительный, основной и этап тестирования.

На первом этапе была подготовлена аппаратная платформа «GRID-компьютера». Для этого использовался двухпроцессорный сервер на базе процессоров Intel Xeon. Отметим, что в общем случае нет необходимости выделять отдельный сервер только под задачи интеграции с GRID. В нашем случае был взят рабочий компьютер обработки, дополнительной функциональностью которого стал интерфейс доступа к GRID-ресурсам. Основное требование к компьютеру — наличие достаточного объема оперативной памяти (4 Гб и более). В случае, если бы для задач интеграции с GRID использовался новый компьютер, то в нашем случае его предварительно необходимо было бы оснастить системным программным обеспечением: операционная система Linux, библиотеки работы с протоколами FTP и SMB/CIFS, пакет OpenSSH, Java Development Kit версии 1.4 и выше.

На второй, основном, этапе была произведена установка и настройка необходимых пакетов программ для взаимодействия с сервером управления обработкой и удаленными GRID-ресурсами. Настройка взаимодействия с сервером обработки не составляет труда и заключается в настройке SSH в соответствии с определенными требованиями сервера обработки (необходимо активировать PubkeyAuthentication режим авторизации на основе открытого и закрытого ключей). Настройка взаимодействия с GRID существенно более трудоемкий и сложный процесс. Был выполнен следующий набор действий:

- 1. На заданном компьютере обработки установлен пакет GT и получен сертификат авторизованного GRID-узла. Технически, именно получение компьютером сертификата делает его полноценным участником GRID-сети. Здесь мы не будем останавливаться на деталях этого процесса. Определение понятия сертификата, его значение и все операции с ним детально описаны в документации [6].
- 2. Активированы три основных сервиса: MDS для мониторинга состояния этого узла со стороны Интернет-сервера ДВВР, RFT для обеспечения обмена файлами между узлами GRID-сети и клиент GridWay для формирования заданий на обработку.
- 3. Определен пользователя, от имени которого будет осуществляться работа в GRID-сети. Получить сертификат на этого пользователя.
- 4. Разработан набор сценариев ориентированных на использование клиента GridWay для выполнения заданий на удаленных ресурсах GRID.

Весь технологический цикл обработки с использованием GRID показан на рисунке 4. Пунктирной линией отделены соответствующие ресурсы Спутникового Центра и ДВВР. Серым фоном подсвечены те ресурсы, которые являются членами объединенной GRID-сети. Запрос инициатора обработки (см.

- рис. 3) проходит следующие стадии. Первым делом выполняются обязательные стадии системы обработки, связанные с брокером сообщений и определением сценария обработки. Далее выбранный сценарий обработки диспетчером направляется на обрабатывающий компьютер, интегрированный в GRID-сеть. Сценарий обработки для GridWay представляет собой скрипт на языке командного интерпретатора BASH. Результатом работы сценария является команда (или набор команд) клиента сервиса GridWay, которая запускается локально на обрабатывающем компьютере и обеспечивает передачу необходимых запросов в сервер GridWay. Командами могут быть:
- 1. Инициализация новой задачи и передача в сервер GridWay для последующего исполнения. Результатом этой команды является уникальный идентификатор задачи.
- 2. Запрос состояния выполнения задачи по ее идентификатору. В настоящее время сценарий обработки распознает четыре основных состояния задачи в GridWay: ожидает, выполняется, успешно завершена и завершена с ошибкой. Эти результаты передаются в соответствующий топик брокера, откуда могут быть прочитаны инициатором или монитором для принятия дальнейших решений.
- 3. Ожидание завершения задачи. В этом случае сценарий закончит работу только тогда, когда задача перейдет либо в состояние успешно завершена, либо ошибка.
- 4. Принудительная остановка задачи и ее удаление из GridWay. В этом случае в систему обработки будет передано ошибочное состояние задачи и ее дальнейшая обработка зависит от действий инициатора.

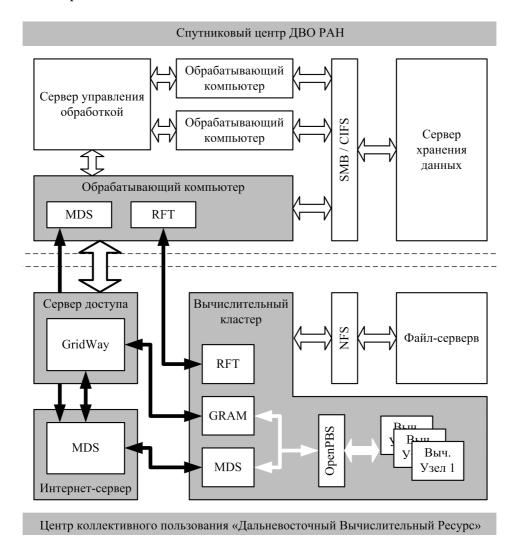


Рис. 4. Структура взаимодействия системы обработки спутниковых данных с GRID-сетью

Каждая команда обладает некоторым набором параметром. Последние три команды имеют всего один параметра – уникальный идентификатор задачи, полученный в результате выполнения первой команды. Параметром первой команды является паспорт задачи, предварительно сгенерированный сценарием. Паспорт задачи формируется в соответствии с требованиями GridWay и описывает основные характеристики задачи. А именно, путь и название исполняемого файла алгоритма ДОТК, имя исходного файла и имя файла с результатами обработки, дополнительные параметры настройки алгоритма. Пример паспорта задачи показан на рис. 5.

Рис. 5. Паспорт задачи

После получения очередной заявки от клиента GridWay сервер разбирает паспорт задачи и определяет объем требуемых ресурсов. В настоящее время для одной подзадачи расчета полей ДОТК объем ресурсов составляет один процессор. Рассмотрим пример на рисунке 5. Здесь переменная EXECUTABLE хранит имя исполняемого файла для запуска алгоритма. Для всех значений переменных допускается макроподстановка значений системных и служебных переменных. В данном случае \$GW_USER хранит имя пользователя, от имени которого будет выполняться задача. ARGUMENTS описывает входные параметры для исполняемого файла. Значение ТҮРЕ говорит о том, что это обычная однопоточная задача. NP определяет количество необходимых процессов для задачи. INPUT FILES и OUTPUT FILES определяет имена файлов которые необходимо скопировать до и после выполнения алгоритма. До выполнения копируется исходный файл из архива, доступного по SMB. После выполнения в архив копируется результат обработки. Значения параметров INPUT и OUTPUT используются службой RFT. Параметр REQUIREMENTS определяет набор ресурсов, необходимых для запуска. В данном случае задано использовать 900-гигафлопный кластер с архитектурой PowerPC64. На основе этих значений GridWay при наличии свободных ресурсов формирует соответствующий запрос в GRAM-сервис заданного GRID-узла (mvs15k.cc.dvo.ru), который и производит выполнение задачи, а также необходимые пред- и пост- копирования файлов данных. В результате выполнения рассмотренного примера взят файл с именем N87426_3.prc. Этот файл будет пропущен через набор GRID-ресурсов. В результате чего будет получен файл N87426_3.st0 с результатами обработки.

Тестирование распределенной GRID-сети

Последний этап интеграции GRID-ресурсов в систему обработки спутниковых данных заключался в проведении серии тестовых запусков обработки. Для тестирования использовался массив изображений, полученный за сутки мониторинга региона акватории Охотского моря. Входной объем данных в среднем составляет около 15-20 файлов и объемом порядка 50 мегабайт. Каждый файл представляет собой изображение в меркаторской проекции для которой необходимо рассчитать поля ДОТК. Для заданного региона разрешение изображения проекции составляет порядка 1500х2000 пикселей. Параллельная обработка данных изображений осуществляется нарезкой изображения на 10 частей, что за сутки составляет порядка 160 обрабатываемых изображений (полоса размером 150х2000

пикселей). Время вычисления полей для одной полосы составляет порядка 1 часа для компьютера с реальной производительностью 3 гигафлопа. Соответственно, расчет всего набора изображений составляет ориентировочно 160 часов или 7 суток.

Применение GRID-технологий для проведения расчетов ДОТК позволило сократить время расчетов в десятки раз. На 10 процессорах расчет 160 изображений занял примерно 8 часов. Увеличение числа процессоров до 160 штук позволил сократить время счета до 4 часов. Нелинейный характер зависимости времени счета от числа изображений связаны с ограничениями планировщика GridWay. Текущая настройка GridWay позволяет проводить не более 25 одновременных заданий от одного пользователя GRID. Если увеличить пропускную способность планировщика GridWay по количеству одновременно выполняемых заданий, то можно получить фактически линейное ускорение времени обработки. При этом суммарное время обработки всех изображений складывается из максимальное время обработки полосы, плюс коммуникационные издержки (менее 1%). Соответственно, в лучшем случае работа алгоритма обработки для 160 полос составит около 30 минут.

Заключение

В результате выполненных работ решена задача интеграции внешних вычислительных ресурсов в систему обработки спутниковых данных с использованием GRID-технологий. Экспериментальные расчеты для задач определения полей ДОТК показали высокую эффективность применения подобных вычислительных ресурсов. Полученный положительный опыт применения GRID-технологий в системе обработки позволяет в будущем перенести в GRID расчеты и других более сложных задач.

Литература

- 1. Левин В.А., А.И. Алексанин, М.Г. Алексанина, П.В. Бабяк. Состояние дел и перспективы развития ЦКП регионального спутникового мониторинга окружающей среды ДВО РАН в области современных информационных и телекоммуникационных технологий // Открытое образование, 2008. №4. С. 23-29.
- 2. *Кудашев Е.Б.*, *Филонов А.Н*. Технологии и стандарты интеграции сервисов, каталогов и баз данных дистанционного исследования Земли из космоса // Тр. 9-й Всерос. научн. конф. «Электронные библиотеки: перспективные методы и технологии, электронные библиотеки» (RCDL'2007). Переславль-Залесский, 2007. С. 273-279.
- 3. *Алексанина М.Г.* Автоматическое выделение поверхностных структур океана по инфракрасным данным спутников NOAA // Исслед. Земли из космоса, 1997. №3. С.44–51.
- 4. **The Open Grid Services Architecture, Version 1.0**. I. Foster, H. Kishimoto, A. Savva, D. Berry, A. Djaoui, A. Grimshaw, B. Horn, F. Maciel, F. Siebenlist, R. Subramaniam, J. Treadwell, J. Von Reich. Informational Document, Global Grid Forum (GGF), January 29, 2005.
- 5. **The WS-Resource Framework**. K. Czajkowski, D. F. Ferguson, I. Foster, J. Frey, S. Graham, I. Sedukhin, D. Snelling, S. Tuecke, W. Vambenepe. March 5, 2004.
- 6. Интернет-ресурс. Globus Toolkit 4.2.1 Release Manuals. http://www.globus.org/toolkit/docs/latest-stable/. Документация к пакету Globus Toolkit.

Use experience in applying GRID-technologies to data processing system of Satellite Center FEBRAS

P.V. Babyak, G.V. Tarasov

Institute of Automation and Control Processes feb RAS 690041 Vladivostok, Radio 5
E-mails: yanger.paul@gmail.com, george@dvo.ru

In the article principles of integration massive parallel computing resources into satellite data processing system based on GRID-technologies are considered. On example of different resources of two Joint Centers (JC) FEB RAS the distributed GRID-network is constructed for carrying out calcaulations of dominant oriantation of thermal contrasts (DOTC). Experimental calculations have shown multuiple speedup in data processing that is especially actual in the round-the-clock monitoring conditions. The received successfule experience allows to use constructed GRID-network for various massive processing tasks, not just DOTC fields calculations.

Keywords: GRID-computing, parallel computing, distributed systems, DOTC.