

Последние достижения в области высокоскоростной обработки данных ДЗЗ*

Достижения в области компьютерных технологий и совершенствование съемочной аппаратуры позволяют разрабатывать принципиально новые способы сбора, обработки и анализа данных дистанционного зондирования Земли (ДЗЗ). В частности, внедрение сенсоров последнего поколения, используемых для наблюдения Земли и планет, в настоящее время позволяет получать практически непрерывный поток данных высокой размерности. Такой резкий скачок объема получаемой информации вызвал необходимость разработки новых методов обработки данных.

Разработка эффективных методов вычисления, обеспечивающих преобразование больших объемов данных ДЗЗ в доступную информацию, имеет большое значение для развития наук о Земле. Рост объема данных дистанционного зондирования продолжается, в то время как международным организациям и пользователям требуются эффективные системы обмена этими данными и ресурсами. В этих целях за последнее время были проведены исследования по применению методов и систем высокоскоростных вычислений (*high performance computing* — HPC) для решения задач ДЗЗ. HPC предусматривают набор встроенных вычислительных сред и методов программирования, которые могут значительно облегчить решение крупномасштабных задач, в т.ч. многих задач дистанционного зондирования. Например, для многих существующих и перспективных областей применения ДЗЗ в науках о Земле и космосе, а также в различных видах разведки требуется обработка в реальном времени или в режиме, близком к реальному времени. Соответствующие примеры включают экологические исследования, военную разведку, отслеживание и

мониторинг опасностей, таких как пожары в лесах и на целине, нефтяные разливы и прочие типы химического/биологического загрязнения.

Использование систем HPC в приложениях дистанционного зондирования за последние годы получило широкое распространение. Идея об использовании стандартных компьютеров (*commercial off-the-shelf* — COTS), объединенных в кластеры, работающие как «вычислительные группы», привела к созданию многих разработок, основанных на многопроцессорных системах.

Алгоритмы обработки данных ДЗЗ в целом очень хорошо внедряются в многопроцессорные системы, состоящие из кластеров, или сети центральных процессоров, однако эти системы, как правило, являются дорогостоящими и с трудом адаптируются к сценариям бортовой обработки данных, в которых критически важными факторами являются малый вес и малая нагрузка интегрированных компонентов, где требуется снизить вес полезного груза спутника и получать результаты анализа в реальном времени, т.е. во время сбора данных сенсором. Хорошую потенциальную возможность устранения разрыва между бортовым анализом данных ДЗЗ и анализом в реальном времени предоставляют новые специализированные аппаратные средства, такие как программируемые вентильные матрицы (*field programmable gate arrays* — FPGA) и графические процессоры (*graphic processing units* — GPU). Растущую потребность приложений ДЗЗ в скоростных вычислениях могут удовлетворить эти компактные аппаратные средства, преимуществами которых являются небольшой размер и относительно низкая стоимость по сравнению с кла-

* Сокращенный перевод с английского языка статьи «*Recent Developments in High Performance Computing for Remote Sensing: A Review*» (авторы *Craig A. Lee, Samuel D. Gasster, Antonio Plaza, Chein-I Chang, Bormin Huang* — IEEE), опубликованной в *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, Vol. 4, No. 3, September 2011. Перевод подготовлен к публикации Б.А. Дворкиным (Компания «Совзонд»)

стерами или компьютерными сетями. Эти аспекты имеют большое значение при определении задач дистанционного зондирования, для которых важным параметром является вес полезного нагрузки.

В оригинале статьи большое внимание уделяется обзорам специализированной аппаратной архитектуры и кластерным вычислениям, однако мы в нашем переводе ограничимся только двумя разделами: инфраструктуре распределенных вычислений и обсуждению основных задач.

ИНФРАСТРУКТУРА РАСПРЕДЕЛЕННЫХ ВЫЧИСЛЕНИЙ

Область распределенных вычислений сильно развилась за последние сорок лет, начиная с появления компьютерных сетей, в особенности со времени появления сети Интернет. Именно инфраструктура распределенных вычислений (Distributed Computing Infrastructure — DCI) изменила способ мышления общества. Активное использование DCI началось в области образования и науки, но быстро распространилось на другие сферы человеческой деятельности, такие как торговля, развлечения и государственное управление. При таком быстром расширении области применения неудивительно, что возникло множество различных названий и «научных терминов» для обозначения способов реализации DCI, применяемых в различных областях деятельности, различными группами пользователей и отраслями промышленности, что иногда вызывает путаницу в терминологии.

Поэтому перед обсуждением DCI, используемых для дистанционного зондирования, мы рассмотрим фундаментальные понятия и характеристики, необходимые для создания и управления DCI. Это позволит разобраться в изобилии названий, используемых для различных технологий распределенных вычислений, а также классифицировать эти технологии — от очень простых, внедряемых и используемых отдельными исследователями, до очень больших и сложных систем, внедряемых и используемых международными организациями. После этого мы рассмотрим существующие DCI, применяемые для дистанционного зондирования, в рамках теоретической эталонной архитектуры для спутниковых наземных систем. Затем мы определим основные области задач для

DCI будущего.

Терминология распределенных вычислений

В области распределенных вычислений применяется большое количество терминов, что иногда может вызывать путаницу. Перед началом дискуссии о влиянии данной области на ДЗЗ мы хотели бы обсудить терминологию, т.е. то, как применяемые термины связаны друг с другом и что они подразумевают в отношении возможностей и преимуществ систем и инфраструктуры распределенных вычислений.

Системы (а также системы систем) разрабатываются на основе разнообразных архитектур (архитектурных типов); эти архитектуры можно реализовать, используя различные каркасы и опорные технологии. Мы дадим определение и описание этих понятий и их связей. Мы также определим ключевые характеристики каждого из этих понятий, в той мере, в которой оно имеет отношение к общей цели проектирования, построения и использования системы распределенных вычислений, поддерживающей приложения по изучению и дистанционному зондированию Земли.

Система — это набор элементов, объединенных для достижения цели или результатов, которых невозможно достичь, используя эти элементы по отдельности. Системы создаются на основе ряда целей и задач. После создания нескольких систем, основанных на различных наборах целей, можно создать систему систем (system of systems — SOS). Системы разрабатываются на базе понятий архитектур и типов, и именно развертываемая система дает пользователям рабочие возможности для достижения поставленных целей.

Виртуальная организация (virtual organization — VO) — это система систем, состоящая из различных объектов, принадлежащих к различным административным доменам (административный домен определяется как набор ресурсов, управляемых провайдером сервисов, который контролирует доступ к этим ресурсам и сервисам, а также их использование) и физически отделенных друг от друга. VO представляет собой структуру, позволяющую нескольким организациям совместно управлять ресурсами и доступом соответствующих пользователей к этим ресурсам. Каждый член VO имеет роль, определяющую функции, которые он может выполнять, данные, которые он может записывать и читать, и сервисы, которые он

может создавать и использовать. Организаторы VO, т.е. владельцы местных данных и сервисов, могут определять права доступа, связанные с отдельными ролями этих местных ресурсов. В случае предоставления доступа посредством ролей система обладает более высокой способностью к масштабированию, чем в случае предоставления доступа отдельным пользователям.

Архитектура — это набор структур, позволяющий ориентироваться в системе, включающий составные элементы системы, связи между ними, и свойства таких элементов и связей. Архитектура должна обеспечивать описание программных и аппаратных элементов, а также их интерфейсов. Архитектура также должна обеспечивать логическое описание составных элементов системы, контекста системы и взаимодействий между элементами, необходимых для достижения целей или задач системы. Архитектурное описание находится на логическом функциональном уровне («что») и не предоставляет информации о конкретных способах реализации («как»). Инженер-разработчик систем должен привязать архитектуру системы к конкретным способам реализации, используя каркасы и опорные технологии.

Различные архитектурные типы подчеркивают различные цели в системах высокого уровня. Например, облачные вычисления относятся к обеспечению вычислительными ресурсами (предоставление ресурсов по запросу), а распределенные вычисления — к интеграции (распределение физических ресурсов при наличии единственного логического интерфейса в рамках системы). Стиль сервисно-ориентированной архитектуры (service oriented architecture — SOA) относится к предоставлению сервисов и соответствующих процедур, протоколов, интерфейсов и инфраструктур связи, обеспечивающих доступ к этим сервисам и их использование. Каждый сервис имеет четко выраженную функцию, которая является самостоятельной и не зависит от контекста или состояния других сервисов. Примером SOA является сенсорная сеть, в которую входят сервисы исходных данных, предоставляемые провайдером (сервисы провайдеров), и различные типы потребительских данных (сервисы потребительских данных). В сервисах потребительских данных могут использоваться исходные данные (данные уровня 0 по классифика-

ции NASA), предоставленные для различных целей, включая создание различных выходных данных, которые в свою очередь становятся доступными в качестве сервисов данных (например, выходные данные уровней 1–3).

Программные каркасы обеспечивают основную структуру или идею, лежащую в основе проекта системы, который обеспечивает блоки для построения системы или приложения. В данном обзоре мы уделяем основное внимание программным каркасам, но аппаратные каркасы можно рассматривать аналогично. Каркасы предоставляют набор библиотек или классов в качестве фундаментальных блоков, а также набор правил или инструкций, относящихся к составу посредством четко определенных интерфейсов и данных. Каркасы обеспечивают управление выполнением программы, поведение по умолчанию, расширяемость и другие конструкции, необходимые для реализации проекта. Каркасы используются для внедрения связующего программного обеспечения, используемого как «клей» для скрепления физического слоя (конкретные технологии выполнения и базовое аппаратное обеспечение) с логическим слоем. Примеры каркасов: набор инструментов Globus Grid и распределенная файловая система с открытым исходным кодом Grid Data Farm (Gfarm). Каркас Apache Hadoop обеспечивает выполнение приложений на больших компьютерных кластерных системах посредством реализации вычислительной парадигмы Map/Reduce.

Опорные технологии — это основные компоненты (аппаратное и программное обеспечение) и протоколы, позволяющие внедрять каркасы и библиотеки, выражающие данный архитектурный тип. Примеры опорных технологий: веб-сервисы (SOAP), HTTP/HTTPS, сетевые протоколы (TCP, IP), коммерческие аппаратные средства, позволяющие создавать кластеры, высокоскоростные волоконно-оптические сети и т. д.

Термин «инфраструктура распределенных вычислений» (DCI) относится к набору логических, физических и организационных элементов, необходимых для создания и функционирования распределенной системы. Такие системы могут быть распределены логически и физически; цель большинства DCI состоит в четком разъяснении этого различия пользователю через понятие виртуализации.

В следующих разделах мы используем термин

«данные» в нескольких различных контекстах. Этот термин может относиться к данным наблюдения Земли (от исходных данных сенсора до выходных данных с высокой степенью обработки, таких как глобальные карты температуры поверхности моря); он также может относиться к информации, передаваемой в рамках управления и администрирования DCI. Мы не будем определять значение этого термина в каждом конкретном случае, т. к. предполагаем, что это значение ясно из контекста. Существуют также метаданные («данные о данных»), при обсуждении которых мы будем однозначно использовать термин «метаданные».

Возможности, масштабирование и преимущества инфраструктуры распределенных вычислений (DCI)

Возможности. Системы дистанционного зондирования Земли должны обладать конкретными возможностями, обеспечивающими достижение общих целей системы (например, представление калиброванных космических снимков с определением географического местоположения), но они также должны включать набор общих возможностей управления и организации инфраструктуры. Функция, определенная для поддержания возможностей пользователей и администраторов, должна поддерживаться каркасами и базовыми технологиями. Одной из важнейших возможностей, предоставленных в DCI, является разграничение логической и физической организации и функций (виртуализация). Это разграничение освобождает пользователя и приложения от необходимости управления ресурсами и инфраструктурой, позволяя им сосредоточиться на конкретных выполняемых операциях.

Примеры типов возможностей, необходимых для DCI, приведены ниже. В данном обзоре рассматриваются как пользователи-люди, так и клиенты-приложения (интерфейсы прикладного программирования). Для упрощения обсуждения мы будем применять термин «клиент» как к пользователям-людям, так и к клиентам-приложениям. Хотя возможности перечислены отдельно, почти всегда существует потребность в их взаимодействии и взаимной поддержке.

- Открытие ресурсов и каталоги. Требуется, чтобы можно было легко находить ресурсы в системе DCI. Ресурсом обычно считается любой тип дан-

ных или сервисов. В целях облегчения поиска ресурсов для пользователей они должны быть размещены в каталогах в доступных для поиска базах данных с четко определенными интерфейсами и языками запроса. Каталоги, наряду со связанными метаданными и синтаксисом запроса, позволяют клиентам находить и получать доступ к ресурсам на основе логического тождества. В результате запроса может быть получена ссылка или карта от логического тождества до одного или нескольких возможных физических объектов, к которым клиент может получить доступ.

- Функциональная совместимость данных — возможность работы с потенциально разнородными сохраненными данными, подходами и способами реализации в различных административных доменах (типичный пример: данные различных форматов, используемые в различных инфраструктурах облачных вычислений, предоставленных различными провайдерами). Эта возможность является фундаментальным требованием для совместного использования данных и доступа через различные домены. Она должна обеспечивать функциональную совместимость семантики, перевода и преобразования данных, мест происхождения данных и безопасности в различных системах. Для использования этой возможности требуется разработка инструментов и процедур; ее реализация должна быть хорошо понята пользователями.
- Управление сервисами/заданиями/процессами. Для совместного управления различными ресурсами требуется способность управлять запросами на сервисы, заданиями и процессами, определяемыми различными клиентами. Для использования этой возможности требуется разработка механизмов распределения ресурсов, способных обрабатывать запросы на сервисы, создавать экземпляры сервисов, располагать по приоритетам запросы на сервисы и отвечать на соглашения об уровнях сервисов (это может предусматривать определение сроков выполнения заданий).
- Создание экземпляров ресурсов и предоставление (распределение) ресурсов. Для использования этой возможности объединяются различные перспективы распределенных и облачных вычислений. Во многих случаях нежелательно выделять

ресурсы для одного приложения, гораздо чаще требуется обеспечивать распределение и предоставление ресурсов по запросу (по мере поступления запросов). Для этого требуется способность поддержания баланса между запросами и поступлением ресурсов, для чего могут потребоваться модели, прогнозирующие потенциальные всплески спроса, и методы получения доступа к дополнительным ресурсам только в случае, когда они необходимы. Задача распределения ресурсов очень сложна, и активная область исследований, позволяющих найти оптимальные подходы, основывается на различных ограничениях и условиях.

- **Мониторинг.** Эта возможность относится к нескольким уровням операций DCI. Должна быть обеспечена доступность и надежность базовой системы, чтобы предоставлять клиентам сервисы и ресурсы, когда они необходимы. Требуются инструменты не только для мониторинга ресурсов в пределах данного административного домена, но также для их распределения по различным доменам (интеграция данных управления DCI). Ошибки и отказы, связанные с ресурсами, должны проверяться и передаваться во избежание ситуаций с нехваткой ресурсов. Клиентам требуется минимальное качество сервисов для многих задач или составление конкретных соглашений об уровне сервиса (service level agreements — SLA). В этой связи требуется способность осуществлять мониторинг ресурсов и сервисов. Такие административные возможности мониторинга позволяют различным системным администраторам DCI наблюдать за общим состоянием и статусом местных доменов и системы в целом. Наконец, способность осуществлять мониторинг общей безопасности системы имеет ключевое значение с учетом ландшафта угрозы, который существует и продолжает расширяться.
- **Уведомление о событиях.** Эта функция имеет ключевое значение для обеспечения асинхронной связи различных элементов DCI-системы. Уведомления о событиях распространяются в приложениях DCI в различных целях, например, регистрация, мониторинг и проверка и другие события, связанные с изменением состояния ресурса или сервиса. Возможные события включают результаты вычисления, обновления статуса, ошибки и ис-

ключения, а также степень выполнения процесса клиента.

- **Безопасность.** Трудно переоценить важность обеспечения доступности, целостности и безопасности информации как фундаментальной возможности в любой системе DCI. Практически все аспекты работы системы связаны с обеспечением безопасности. Это необходимо для того, чтобы гарантировать клиентам целостность их данных и результатов анализа. Основные компоненты включают механизмы идентификации и авторизации клиентов и процессов. Эти способности должны обеспечивать сложные перекрестные операции различных доменов, например, регистрацию во всей сети путем однократного ввода пароля, при поддержании безопасности системы. Информационная целостность — это способность обеспечить защиту от несанкционированного изменения или уничтожения информации. Учитывая важность данных наблюдения Земли для многих аспектов национальной и международной политики, мы можем сказать, что целостность и происхождение данных имеют критическое значение.
- **Отчетность и проверки.** С учетом того, что многие ресурсы, используемые для построения DCI, поступают из различных источников, включая коммерческие объекты, необходимо отслеживать использование ресурса для платных сервисов. Внутренние средства проверки в любой системе должны отслеживать схемы использования, чтобы определить области, где могут потребоваться дополнительные ресурсы или области, где ресурсы используются недостаточно полно.

Масштабирование. Важным аспектом DCI является способность масштабировать систему в ответ на изменение требований ресурса и требований системы. Для оценки этих изменений необходимо количественно определить характеристики системы в отношении рабочих параметров или параметров масштабирования. Параметры масштабирования касаются работы различных сервисов и ресурсов, предоставляемых системой. Например, время ожидания часто является ключевым рабочим параметром, так как для многих систем требуется почти постоянная поддержка, как в случае борьбы со стихийными бедствиями. В таких системах может быть задано время ожидания

для времени прибытия данных от сенсоров или результатов вычислений прогнозных моделей. Система должна быть способна обеспечивать пропускную способность сети и мощность для поддержания требований времени ожидания, особенно в течение периодов, когда требования ресурса изменяются. В SOA время, требуемое для завершения различных запросов на сервисы, является важной рабочей характеристикой, поскольку оно может зависеть от множества факторов, таких как внутренние сообщения и отбор образцов/ограничения сервисов сенсора, ограничения ресурса для конкретного сервиса, использование и диапазон частот. Системы обычно разрабатываются на основе конкретных требований к рабочим характеристикам, которые могут также включать потенциальный рост спроса со временем. Важным фактором является число клиентов, которое может варьироваться от маленьких групп (≤ 10 пользователей) до крупномасштабных VO (~ 1000 пользователей).

Преимущества. Использование DCI для наблюдения Земли и приложений дистанционного зондирования обеспечивает множество преимуществ благодаря четко определенным концепциям и типам архитектуры, и их реализации с использованием стандартных каркасов и опорных технологий. Концепция виртуализации сенсоров лежит в основе концепции сенсорной сети, упомянутой выше.

Концепция виртуализации является одним из главных преимуществ DCI. Ранее мы упоминали концепцию виртуальной организации, но концепция виртуализации имеет фундаментальное значение для определения способа, которым связующее программное обеспечение, обеспечивает связь между физическим устройством и логическим пользовательским интерфейсом. Цель состоит в том, чтобы освободить пользователей от управления ресурсами, необходимыми для выполнения процесса, и позволить им сосредоточиться на конкретных научных исследованиях. Пользователи могут логически обнаруживать и получать доступ к данным или вычислительным ресурсам и включать их в процесс, не заботясь об их физической реализации. Конечная цель, разумеется, состоит в предоставлении этих возможностей по запросу и удовлетворении требований пользователей, связанных с рабочими характеристиками и сро-

ками. Концепция виртуализации ресурсов может быть применена к любому типу ресурсов — от вычислительной инфраструктуры (CPU, хранение, частота) до источников данных, таких как сенсоры и приборы. Сенсоры могут быть виртуализованы таким образом, чтобы возможности дистанционного зондирования и измерения в месте нахождения были доступны в виде сервисов.

Пользователи могут определять собственные потребности в данных, используя «естественный» синтаксис и семантику, которые система впоследствии транслирует в конкретный процесс, чтобы выполнить этот запрос на данные (например, пользователь задает ограничивающий пространственно-временной прямоугольник с требованиями к пространственно-временному и спектральному разрешению и отбору проб, а система определяет, какой сенсор может выполнить этот запрос наилучшим образом).

Виртуализация также позволяет пользователям искать и находить параметры и данные наблюдений, основанные на характеристиках метаданных, специально предназначенных для такого анализа. Для этого может быть выполнен поиск с использованием пространственного (географического) и временного ограничивающего прямоугольника, характеристик отбора проб (пространственных и временных) и замеров или геофизических параметров. Дополнительное преимущество виртуализации заключается в возможности выполнения модернизации без ущерба или практически без ущерба для доступности.

Еще одно преимущество DCI — функциональная совместимость, достигаемая посредством использования стандартов архитектуры и реализации для протоколов и интерфейсов, таких как протоколы и интерфейсы в SOA. Такая поддержка стандартов обеспечивает комплексность и расширяемость, создавая инфраструктуру, отвечающую функциональным требованиям, указанным ранее. Дополнительные свойства, такие как повторное использование и быстрое развертывание, также являются важными преимуществами подхода DCI. Большое значение имеет способность упорядочивать все необходимые ресурсы по запросу на основе инициирующего события, например, стихийного бедствия, такого как ураган или землетрясения. Такая система может поддерживать готовность с ограниченным использовани-

ем ресурсов до тех пор, пока не потребуется их полное использование. Таким образом, эти ресурсы не «простаивают», а могут использоваться другими приложениями, пока не потребуется доступ, определяемый по схеме приоритетов (события бедствий получают более высокий приоритет, чем плановые научные исследования).

Дополнительное преимущество использования сервисно-ориентированных архитектурных концепций, в сочетании с правильно выбранными каркасами и технологиями реализации, заключается в способности создавать так называемую среду поддержки программирования (Problem Solving Environment — PSE). PSE предназначена для создания каркаса, нацеленного на решение конкретного класса проблем в пределах данной научной области. Каркас представляет инструменты на естественном языке конкретной научной дисциплины таким образом, чтобы пользователь мог упорядочить эти ресурсы при очень слабо выраженной кривой изучения. Каркас может содержать очень мощные средства обработки и анализа данных, объединенные с основными вычислениями и ресурсами данных способом, хорошо понятным для пользователя.

Теоретическая эталонная архитектура наземной системы для ДЗЗ

Теперь, после четкого определения фундаментальных возможностей и терминологии DCI, мы используем системный подход для идентификации возможностей, необходимых для организации сбора

данных с орбитальных и наземных сенсоров, формирования выходных данных, их использования большим распределенным сообществом пользователей и управления предприятием DCI в целом. Для этого мы представляем теоретическую эталонную архитектуру спутниковой наземной системы, показанную на рис. 1, реализованную в виде сервисно-ориентированной архитектуры, где доступ администратора и пользователя осуществляется через инструменты типа браузера. Ключевым аспектом этой эталонной архитектуры является разделение сервисов на предметные сервисы (domain services) и сервисы компании (enterprise services). Предметные сервисы — это сервисы, предназначенные для управления спутниковыми системами (и, возможно, другими сенсорными системами), например, оперативное управление, определение орбиты, определение задач/планов/расписаний (для орбитальных ресурсов), телеметрия и т.д. Специфические особенности таких сервисов не входят в сферу рассмотрения этой работы, поэтому мы не приводим подробного описания этих сервисов.

К предмету обсуждения данной работы непосредственно относятся сервисы компании, обеспечивающие возможности использования и организации всех остальных аспектов инфраструктуры. Четко выделяются сервисы для внесения в каталоги и обнаружения ресурсов, т.е. данных и сервисов. Сервисы исполнения и процесса контролируют срок использования отдельных экземпляров сервисов, а также последовательности выполнения сервисов и все

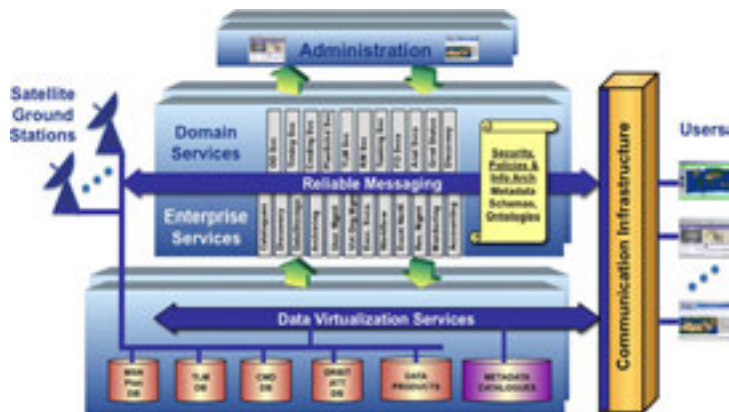


Рис. 1. Теоретическая архитектура спутниковых наземных систем

необходимые передачи данных. Сервисы мониторинга, уведомления о событиях и сервисы отчетности используются для оценки работы системы, мониторинга ошибок и сохранения журнала проверок. Все эти сервисы могут быть динамически распределены из хранилища ресурсов, т.е. облака, в рамках управления ресурсами.

Некоторые другие сервисы охватывают остальные аспекты инфраструктуры. К ним относятся такие сервисы, как надежная передача сообщений, безопасность и процедуры управления, выполняемые во всей пользовательской среде. Следует отметить, что все эти сервисы фактически могут быть распределены по различным участкам. Надежная передача сообщений означает, что в случае, если сообщение не удалось передать, гарантируется создание состояния ошибки, т.е., другими словами, сбой связи не могут остаться незамеченными.

Механизмы безопасности и инфраструктура обеспечивают поддержку целостности и конфиденциальности статических (находящихся на диске) и динамических (находящихся в сети) данных. Для обеспечения целостности используются контрольные суммы и другие методы. К операциям, непосредственно выполняемым большинством пользователей, относятся идентификация и авторизация. В распределенной окружающей среде управление учетными данными пользователей требует совместного управления идентификацией и управления виртуальной организацией. Совместное управление идентификацией предусматривает взаимное доверие пользователей различных организаций. Виртуальные организации (VO) обеспечивают механизм, посредством которого ролевая авторизация может быть осуществлена на основе кода пользователя и его роли в VO, которая может охватывать несколько административных доменов. VO также может быть использована для управления общими данными и общими инструментами со стороны участников VO.

Управление компанией обычно осуществляется посредством процедур. Эти процедуры могут выполняться администраторами или (автоматически) системой. Процедуры использования в основном инициируются посредством ролевой авторизации пользователей. Также существуют процедуры управления системой, определяющие продолжительность обра-

ботки задач, тиражирование данных на участках и т.д.

Отдельный класс составляют сервисы виртуализации данных (data virtualization services). Данные, сформированные орбитальными сенсорами, должны быть собраны, откалиброваны, внесены в каталоги, заархивированы и предоставлены для доступа зарегистрированным пользователям. Базы данных используются для поддержания оперативных данных, таких как планирование задач, телеметрия, контроль, данные о положении на орбите. Большинство пользователей уделяют основное внимание выходным данным и соответствующим метаданным. Эти каталоги могут быть массивными и распределенными. Данные также должны быть заархивированы на неопределенное время. Следовательно, оптимальным способом является виртуализация данных, при которой доступ к данным осуществляется через их атрибуты, т.е. пользователям не требуется знать физическое местоположение данных, формат хранения и т.д. Для этого требуется создание информационной архитектуры, определяющей схемы метаданных и онтологии. Виртуализация данных в информационной архитектуре облегчает определение места происхождения данных, понимание развития данных и их долгосрочное сохранение.

Кроме того, следует отметить, что виртуализация данных также облегчает виртуализацию сенсоров. Доступ к данным по атрибуту может быть применен к тем данным, которые будут сформированы сенсорной сетью, а также к ранее собранному и заархивированному данным. Этот подход обеспечивает чистый логический интерфейс для запросов, сбора и использования данных сенсора, избавляющий пользователя от необходимости изучать технические подробности работы дистанционной сенсорной системы. Многие из указанных проблем обсуждаются более подробно в контексте геокосмических данных в работе.

Обзор ключевых примеров инфраструктуры распределенных вычислений (DCI)

Упомянутая эталонная архитектура предоставляет контекст для обсуждения и оценки ключевых примеров DCI, в которых применяется дистанционное зондирование Земли.

Проект Matsu. Цель проекта Matsu состоит в том, чтобы по запросу обеспечивать возможность оценки бедствий (основанной на облачных вычислениях)

посредством сравнения космических снимков. Этот проект предусматривает сотрудничество исследователей при поддержке Открытого консорциума по облачным системам (Open Cloud Consortium — OCC). OCC управляет распределенной инфраструктурой типа облака, ведущие узлы которой обеспечивают участники OCC и участники рабочей группы по большим объемам данных. Эта инфраструктура представляет собой облако на базе платформы Eucalyptus, содержащее более 300 ядер, 80 Тб памяти, и сетевые соединения с пропускной способностью 10 Гбит/с (с возможностью апгрейда до 80 Гбит/с), сетевое оборудование для которого предоставлено компанией Cisco.

Исходный сценарий обработки данных для проекта Matsu представляет средство прогноза и оценки наводнений в Намибии. На основе относительно простых мэтшапов Web 2.0 в проекте Matsu реализована сенсорная сеть, осуществляющая сбор данных сенсоров из множества источников, включая шесть намибийских речных станций. Matsu также получает данные из онлайн-источников, таких как Глобальная система предупреждения и координации бедствий (Global Disaster Alert and Coordination System), и онлайн-ежедневные маски наводнения, сформированные центром обработки данных MODIS (НАСА).

Что еще более важно, пользователи Matsu могут предложить для сенсоров Hurricane и ALI спутника EO-1 задачу по сбору гиперспектральных изображений для областей интереса. После сбора данных снимки подвергаются радиометрической и геометрической коррекции и сохраняются на облаке OCC. Возможно проведение сравнения изображений в целях оценки наводнения с использованием Hadoop. Окончательные данные предоставляются конечным пользователям, использующим стандартные инструменты сетевого картографирования OGC и обработки охвата сетью.

GENESI-DR и GENESI-DEC. Первоначальная цель проекта GENESI-DR (Ground European Network for Earth Science Interoperations — Digital Repositories; Наземная европейская сеть для взаимодействия исследователей Земли — цифровая база данных) состояла в создании большой распределенной инфраструктуры данных для удовлетворения потребностей международных сообществ. Последующий

проект, GENESI-DEC (Digital Earth Communities — Сообщества цифровых исследований Земли), выполняется до 2012 г. с целью усиления поддержки конкретных сообществ пользователей и прочих существующих архивов данных.

Используя обычный веб-портал и веб-сервисы API, пользователи могут регистрировать собственные наборы данных и предоставлять к ним доступ для других исследователей Земли. Наиболее сложной задачей проекта GENESI-DR было внесение в каталоги неоднородных наборов данных; эта задача решалась путем создания правил метаданных на основе характеристик метаданных для INSPIRE (Европейской инфраструктуры пространственных данных). Последняя была представлена в виде каркасной модели описания ресурсов, использующей общие словари. Благодаря интеграции технологии OpenSearch GENESI-DR поддерживал геокосмические и временные поисковые запросы, основанные на тексте нестандартного формата или на конкретных параметрах метаданных. По завершении проекта GENESI-DR были доступны более двенадцати европейских сайтов и более пятидесяти наборов данных (включая спутниковые наборы данных).

После внедрения базовой инфраструктуры проекта стало очевидно, что необходима модель авторизации, обеспечивающая соблюдение прав интеллектуальной собственности, определенных владельцами данных. Также потребовалось обеспечить возможность регистрации путем однократного ввода пароля во всей сети цифровых баз данных, используемых в разных административных доменах, в целях поддержки перекрестных процессов. Решение этих задач в проекте GENESI-DEC осуществлялось путем использования стандарта OpenID на базе концепции виртуальной организации. Помимо работы с различными сообществами, GENESI-DEC входит в ассоциацию, пропагандирующую концепцию совместного использования данных для общей инфраструктуры GEOSS.

G-POD. Цель проекта распределенной обработки данных по запросу (Grid Processing on Demand — G-POD) заключается в обеспечении обработки данных наблюдения Земли по запросу. Проект G-POD был запущен Европейским космическим агентством в 2002 г. с применением грид-архитектуры, но впослед-

ствии был использован подход на базе облачных вычислений.

G-POD предоставляет собой портал, посредством которого пользователи могут искать данные в каталоге. К требуемым наборам данных можно получить доступ через различные команды. Проект содержит наборы данных, полученных со спутников ERS-1 и ERS-2, а также от сенсоров Envisat ASAR и MERIS. Портал предоставляет сервисы, в которых пользователь может использовать различные инструменты и алгоритмы для обработки наборов данных от уровня 0 (исходные данные сенсора после удаления помех связи) до уровня 3 (геофизические переменные с радиометрической и геометрической калибровкой, привязанные к однородной пространственно-временной системе координат). После запуска заданий обработки можно осуществлять управление этими заданиями и проверять их статус, т.е. какие из них поставлены в очередь, выполняются, завершены и т.д.

G-POD был изначально построен с помощью пакета Globus. Под удобным порталным интерфейсом в проекте G-POD использовались функции GridFTP для передачи наборов данных и GRAM для представления заданий на предварительно сконфигурированных вычислительных ресурсах. Несмотря на традицион-

ную схему процесса с передачей данных в виде пакетов, G-POD обеспечивал возможность обработки по запросу.

Впоследствии ESA использовала Terradue Srl для расширения и коммерциализации G-POD. В результате этой работы в G-POD появилась возможность использовать при необходимости вычислительные узлы Amazon EC2 и блоки хранения S3 без внесения существенных изменений в пользовательский интерфейс. Другими словами, портал перемещает данные наблюдений Земли в блок хранения S3, извлекает из него и управляет сервисами, аналогичными объектам EC2, при этом представляя пользователю тот же интерфейс. Это яркий пример того, что облака используются прежде всего для предоставления ресурсов. На рис. 2 показана страница сервисов G-POD. Кроме портала, доступ к сервисам G-POD обеспечивается также через HTTP и SOAP. Пользователи G-POD получают доступ по сертификатам PKI, выпускаемым администраторами G-POD.

GEO Grid. Цель GEO Grid состоит в обеспечении возможности оценки бедствий; этот проект может считаться прототипом оперативной системы мониторинга стихийных бедствий. GEO Grid объединяет грид-технологии, обеспечивающую надежное управ-



Рис. 2. Страница сервисов G-POD

GEO Grid использует инфраструктуру безопасности сетки (GSI) в сочетании с концепцией VO для реализации масштабируемого механизма авторизации для различных групп пользователей. В настоящее время GEO Grid работает с VO, разработанными для «геологических угроз», а также для «бизнеса, IT и ГИС». Для иллюстрации возможностей GEO Grid рассмотрим виртуальную организацию Сети полевых наблюдений (Field Observation Network — FON), предназначенную для поддержания калибровки и аттестации орбитальных сенсоров, посредством сравнения орбитальных данных с другими источниками данных, например, наземными наблюдениями. Как показано на рис. 3, виртуальная организация FON объединяет данные, полученные из сети наземных обсерваторий (цифровые данные, захватываемые камерой типа «рыбий глаз»), данные полусферического спектрорадиометра и данные солнечного фотометра. FON VO управляет наземными сенсорами на основе стандарта Сервиса наблюдений сенсоров OGC (SOS). Используя порталные сервисы GEO Grid, пользователи могут оценить точность и свойства орбитальных сенсоров.

GEOSS. Цель проекта GEOSS (Global Earth Observation System of Systems — Глобальная система систем наблюдения Земли) состоит в развертывании совместной международной инфраструктуры для совместного использования данных наблюдения Земли во всем мире. Проект поддерживает девять социальных сфер: контроль стихийных бедствий, здравоохранение, энергетика, климат, вода, погода, экосистемы, сельское хозяйство, биологическое разнообразие. Проектом GEOSS руководит Группа наблюдения Земли (Group on Earth Observations — GEO), международное объединение организаций, формирующих и потребляющих данные наблюдения Земли. Текущий рабочий план, определенный для периода 2009–2011 гг., нацелен на построение интегрированной общей инфраструктуры GEOSS (GEOSS Common Infrastructure — GCI). На рис. 4 показана структура GCI в виде сервисно-ориентированной архитектуры. Набор реестров используется для компонентов сервиса, требований пользователей и стандартов функциональной совместимости. Здесь группы пользователей со всего мира могут регистрировать собственные наборы данных и сервисы. В целях

облегчения поиска ресурса Центр обмена информацией GEOSS выполняет глобальный поиск GEOSS на основе зарегистрированных метаданных для всех типов ресурсов, например, систем, сервисов, данных, документов или конкретных типов файлов. Доступ ко всем компонентам системы осуществляется через портал посредством ввода текста нестандартного формата, просмотра социальных сфер или выбора местоположения на интерактивном глобусе.

В качестве члена GEO Комитет по спутникам наблюдения Земли (Committee on Earth Observation Satellites — CEOS) предоставляет космический сегмент для этого проекта и, соответственно, данные, вносимые в каталоги в этих реестрах. Участники CEOS управляют спутниковыми программами, формирующими эти данные в непрерывном режиме.

В целях поддержки GEOSS CEOS разработал концепцию виртуальных группировок спутников, при которой предусмотрено скоординированное управление спутниками и наземными сегментами, используемыми одной или несколькими организациями, что позволит выполнять общие требования наблюдений Земли. Для этого GEO и CEOS проводят ряд совместных мероприятий (CEOS-GEO) в рамках десятилетнего плана GEOSS и текущих рабочих планов на 2009–2011 гг. План включает виртуальные группировки спутников и принципы совместного использования данных в дополнение к поддержке конкретных социальных сфер, таких как глобальный сельскохозяйственный мониторинг.

Обсуждение основных задач

Хотя в приведенном обзоре были продемонстрированы хорошо себя зарекомендовавшие системы, в этой области все еще остается большое количество нерешенных задач. Для определения диапазона и масштаба этих задач мы используем еще один ключевой показательный пример. В 2005 г. ураган «Катрина» унес более 1500 жизней и вызвал материальный ущерб свыше 81 млрд долл. За четыре дня до достижения урагана берега различные системы прогноза урагана давали результаты, показанные на рис. 5.

Очевидно, что прогнозы были ненадежны за четыре дня до урагана и начали приближаться к «истине» только за два дня до бедствия. Что же требуется для построения и развертывания НРС-системы, которая обеспечит смягчение последствий бедствия?



Рис. 5. Прогнозы пути следования урагана «Катрина» за четыре дня до подхода его к берегу в штате Луизиана. Черная линия обозначает фактический маршрут

При рассмотрении такой системы становится понятно, что ее создание представляет очень сложную проблему как с научной, так и с оперативной точки зрения. Для решения базовых научных проблем потребуется существенное повышение знаний о функционировании атмосферных и океанских систем в рамках общей системы Земли, а также разработка соответствующих вычислительных моделей, точно представляющих эти системы. Масштаб этих моделей может потребовать создания более крупномасштабной вычислительной инфраструктуры по сравнению с существующими.

В качестве примера рассмотрим требования DCI для отслеживания урагана с момента зарождения до полной силы. DCI должна обладать способностью усваивать большой объем данных в реальном времени, включая наблюдения со спутников, воздушных судов и наземных систем. Эти данные должны быть переданы в модель прогноза в реальном времени для прогнозирования пути урагана, а затем в различные организации и системы обеспечения решений. Результаты этих моделей отслеживания также должны быть переданы в модели выпадения осадков для оценки скопления воды, которая должна быть передана в модели наводнений, чтобы определить зоны риска для жизни и имущества людей. Для достижения максимальной эффективности это должно быть сделано в объединенных организациях и странах, чтобы обеспечить немедленный доступ к критически важной информации для государственных чиновников, которые будут управлять маршрутами эвакуации, заложением мешков с песком и другими мероприятиями по

смягчению бедствия. Для реализации такой DCI требуется огромная мощность вычислений, экономически нереальная в случае концентрации на единственной указанной цели.

Следовательно, потребуется дополнительно использовать совместные вычислительные ресурсы (включая все типы НРС-платформ, рассмотренных в обзоре). Хотя роль каждого типа архитектуры сильно зависит от рассматриваемого приложения дистанционного зондирования, параллельные кластерные вычисления представляются наиболее подходящими для эффективного извлечения информации из очень больших архивов данных, в т.ч. наборов данных, уже переданных на Землю, в то время как критические по времени ограничения, введенные многими приложениями ДЗЗ (например, приложение, рассмотренное в этом разделе), требуют бортовых средств обработки и зачастую средств обработки в реальном времени, включая специализированные аппаратные архитектуры, такие как GPU и FPGA. Во всех случаях эти вычислительные ресурсы также должны быть доступны по запросу, возможно, из национального облака-ресурса, который может поддерживать объединенные НРС-коды с жесткими сроками обработки.

Очевидно, что такая масштабная сложная система может поддерживать широкий диапазон доменов. С учетом этого мы можем выделить следующие фундаментальные аспекты:

- Своевременность и масштаб по запросу. До недавнего времени выполнение крупномасштабных вычислительных заданий означало передачу задания планировщику и ожидание в очереди заданий. Однако облачные вычисления основаны на принципе получения ресурсов по запросу. Хотя коммерческие облачные вычисления рассчитаны прежде всего на транзакционный стиль вычислений, наблюдается также рост интереса к построению облаков научных приложений, которые могут поддерживать более тесно связанные НРС-коды по запросу. Для обеспечения возможности смягчения бедствий в DCI потребуется распределение ресурсов для поддержания наборов приложений в рабочем процессе, усваивающих данные реального мира в реальном времени и передающих выходные данные в распределенную пользователь-

скую базу. Этот масштаб равнозначен распределению виртуальных центров обработки данных с жесткими ограничениями сроков в реальном времени.

- Обеспечение доступности, целостности и безопасности информации. Такие большие системы могут быть фактически распределены по нескольким центрам обработки данных в различных административных доменах с пересечением не только организационных, но также и национальных границ. Следовательно, требуются совместные системы управления идентификацией типа single sign-on (регистрация во всей сети путем однократного ввода пароля). Ролевая авторизация осуществляется на основе идентификационного кода пользователя и его роли в рамках виртуальной организации. Доверительные отношения, требуемые для управления такими виртуальными организациями, регулируются посредством доверительных федераций, определяющих работу сертификационных организаций, выдающих сертификаты для всех участников такой федерации. Однако, независимо от применяемых механизмов обеспечения безопасности, требуется фундаментальный компромисс между безопасностью, с одной стороны, и работоспособностью и удобством системы, с другой стороны. Выбор правильного уровня безопасности, обеспечивающего оптимальный баланс между работоспособностью и безопасностью системы, всегда является сложной задачей.
- Данные и доступ к данным. В некоторых проектах, рассмотренных в данном обзоре, основное внимание уделяется данным и доступу к данным. При текущем «наводнении данными», т.е. при огромном объеме захватываемых, формируемых и размещаемых в Интернете данных, трудно переоценить важность доступа к данным. Существует ряд стандартов для геокосмических данных, каталогов и сетевого представления, но это всего лишь первый шаг на пути решения данной проблемы. Учитывая широкое разнообразие форматов данных и метаданных, крайне важно разработать улучшенные методы управления информацией и предоставить пользователям простые методы (и инструменты) для получения доступа к этим данным. В случаях, когда наборы данных принадлежат

различным учреждениям, могут использоваться различные методы идентификации, авторизации и доступа к данным, что еще более усложняет проблему доступности данных. В конечном счете, цель состоит в создании цифровых библиотек, в которых обрабатываются и сохраняются наборы текущих данных и данных прошлых периодов с регистрацией места происхождения, доступ к которым осуществляется на основе четко определенного набора стандартов.

- Стандарты и функциональная совместимость. Очевидно, что ни одна из этих систем не может быть реализована без всемирно признанных и принятых стандартов во всех фундаментальных областях, указанных ранее. Число возможных технических стандартов слишком велико, и нет смысла перечислять их в данной работе. В настоящее время разрабатываются новые стандарты для систем облачных вычислений, такие как открытый интерфейс облачных вычислений, открытый формат виртуализации, и интерфейс обработки данных облака, которые в сочетании формируют базу для стандартных облаков IaaS.
- Способы постепенного принятия. С учетом опыта, извлеченного из предыдущих неудачных попыток развертывания крупномасштабных систем в виде монолитного целого, и того факта, что международные стандарты разрабатываются и принимаются не за одну ночь, можно утверждать, что более целесообразным подходом является постепенное развертывание систем. Имея общее представление о будущих больших системах и стандартах для их поддержки, можно поэтапно принимать и развертывать развивающиеся технологии и разрабатывать стандарты для конкретных функций. Такой подход «разработки на ходу» обладает значительными преимуществами. Он позволяет получить ценный опыт и свести к минимуму риск и в то же время пользоваться преимуществами новых вычислительных технологий в плане соотношения цена/производительность. Следовательно, необходимо направить усилия на исследования и испытания экспериментальных систем – это позволит повысить уверенность пользователей и стабильность рынка для всех аспектов НРС-систем, рассмотренных в этом обзоре.