

# ВОПРОСЫ ОБРАБОТКИ СЕЙСМИЧЕСКИХ ДАННЫХ В «ОБЛАЧНЫХ» ВЫЧИСЛИТЕЛЬНЫХ СРЕДАХ

Е.А. Курин

## ВВЕДЕНИЕ

Наиболее достоверным и информативным современным геофизическим методом является сейсморазведка, основанная на принципе эхо-локации. В некоторой точке пространства (чаще всего — вблизи поверхности земли) возбуждаются упругие колебания, вызывающие волновые процессы в земной коре. Волна, отражённая от глубинной залежи полезных ископаемых, несёт в себе информацию о свойствах этого объекта. Задачей обработки сейсмических данных является извлечение этой информации и построение модели и изображения изучаемой среды.

Практически все современные алгоритмы для обработки сейсмических данных, особенно для построения скоростных моделей и глубинных сейсмических изображений, являются чрезвычайно требовательными к вычислительным ресурсам. Некоторые эффективные процедуры построения моделей и изображений требуют применения самых мощных на сегодняшний день вычислительных устройств – суперкомпьютеров. Существующие программные пакеты для обработки данных сейсморазведки рассчитаны на работу на специализированных вычислительных центрах. Реализованные алгоритмы, как правило, используют технологию MPI для обмена данными между вычислительными процессами по быстрой коммуникационной сети, например, Infiniband, а при чтении-записи данных полагаются на наличие кластерной системы хранения данных, например, производства Panasas или NetApp. Подобные аппаратно-программные решения имеют высокую цену, и применение наиболее эффективных, но требовательных к вычислительным ресурсам, методов обработки остаётся прерогативой крупных сервисных компаний. В связи с этим представляется разумным поиск альтернативных способов организации вычислений при обработке больших объёмов сейсмических данных. В настоящей работе будут рассмотрены общие вопросы реализации некоторых алгоритмов обработки для одной из альтернатив – вычислений в «облачных» вычислительных средах.

## «ОБЛАЧНЫЕ» ВЫЧИСЛИТЕЛЬНЫЕ СРЕДЫ ДЛЯ ТЕХНИЧЕСКИХ ВЫЧИСЛЕНИЙ

В последние годы на рынке услуг по аренде вычислительных ресурсов («облачные» инфраструктуры, IaaS, PaaS) появилось достаточно много предложений. Среди них есть как те, что базируются на специально основанных для этого вычислительных ресурсах (например, Amazon Elastic Compute Cloud [1]), так и те, что основаны на аппаратно-программных средствах, имеющих какое-либо иное основное предназначение (например, поисковые системы).

В первом случае пользователю обеспечены достаточно удобные возможности по организации собственной вычислительной среды, и *теоретически*, можно использовать существующее ПО. Однако экономические расчёты [2] показывают, что при некоторых условиях стоимость использования таких ресурсов в задачах обработки научных данных сопоставима со стоимостью владения собственным оборудованием. Кроме того, даже специализированные провайдеры имеют весьма ограниченные ресурсы (например, Cluster Compute Instances в Amazon EC2), которые могут быть эффективно использованы для работы программ, разработанных для использования «традиционных» высокопроизводительных систем.

Во втором случае, владельцы вычислительных ресурсов стремятся обеспечить лучшую загрузку имеющихся вычислительных мощностей, решающих основную задачу. Поэтому стоимость аренды машинного времени у них существенно, на порядок, ниже цен специализированных провайдеров. Широкому практическому использованию подобных вычислительных ресурсов препятствуют как организационные причины, так и технические. К первым относится, в первую очередь, осторожное отношение компаний-заказчиков обработки к размещению их данных в вычислительных системах, доступных возможным злоумышленникам через сеть Интернет. Отсутствие эффективных программных реализаций алгоритмов с использованием технологий облачных вычислений является главным техническим препятствием. Однако положительным моментом является то, что большинство подобных провайдеров использует сходные между собой подходы к организации процесса вычислений, основанные на технологии MapReduce [3], а также на применении той или иной распределённой файловой системы. Наибольшее распространение получила вычислительная среда Hadoop [4], реализующая технологию MapReduce и распределённую файловую систему HDFS. Следует отметить, что и основная часть ресурсов, предоставляемых специализированными провайдерами, также оптимизирована для работы MapReduce-приложений.

## РЕАЛИЗАЦИЯ ПОТОКОВ И АЛГОРИТМОВ ОБРАБОТКИ

Рассмотрим реализацию некоторых обрабатывающих процедур для эффективного выполнения в среде Hadoop. В отличие от одно- и многоканальных процедур стандартной обработки сейсмических данных, некоторые алгоритмы для подавления волн-помех и построения глубинных изображений имеют сложную

коммуникационную структуру, причём некоторые из них состоят из нескольких шагов, каждый из которых имеет собственную коммуникационную структуру. Наши исследования показывают, что такие процедуры следует реализовать как последовательный вызов Map, Reduce, Map+Reduce блоков.

В начале рассмотрим организацию вычислений при стандартной обработке сейсмических данных. На рисунке 1 показана возможная реализация графа стандартной обработки (получение суммарного временного разреза ОСТ) данных при помощи одного из под-проектов Hadoop – Cascading [5]. Все коммуникации осуществляются при помощи средств Cascading, а обрабатываемые процедуры вызываются пользовательскими функциями (UDF). Так как программные средства Cascading реализованы на языке Java, а процедуры обработки - на языке C, то необходимо организовать взаимодействие этих частей. Оно осуществляется посредством вызова вычислительных программ как независимых процессов и обмена данными через стандартные потоки ввода/вывода.

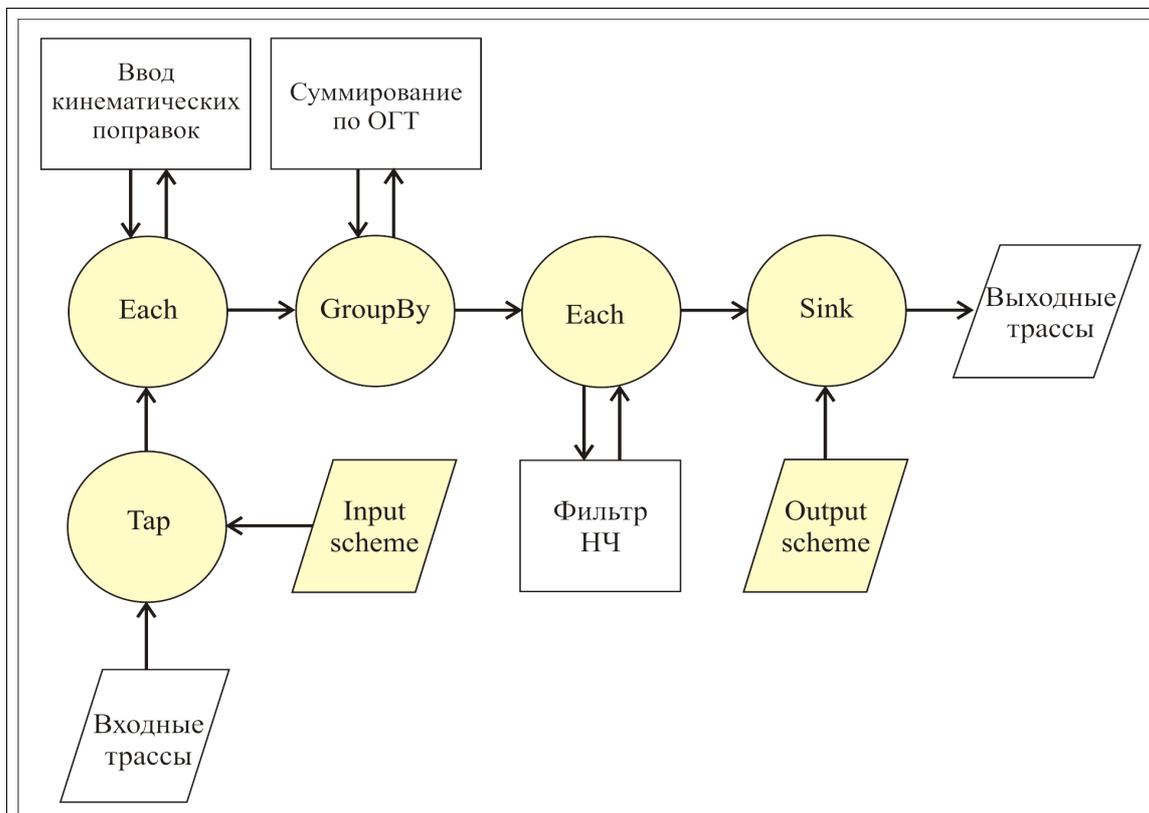


Рис.1. Применение Hadoop/Cascading для организации потока стандартной обработки данных сейсморазведки

На первом этапе в каждый элемент данных (сейсмическую трассу) независимо от других элементов вводятся кинематические поправки, представляющие собой переменное по времени масштабирование записи. После этого осуществляется суммирование поднаборов трасс, имеющих общее значение положения так называемой общей средней точки (ОСТ), то есть выполняется операция редукции по ключу ОСТ. После этого каждая результирующая трасса независимо проходит через операцию полосовой фильтрации и записывается на диск. Отметим, что схожий подход применяется в экспериментальном проекте Seismic Hadoop [6], инициированном компанией Cloudera, где в качестве базовой надстройки над Hadoop используются программные средства Crunch.

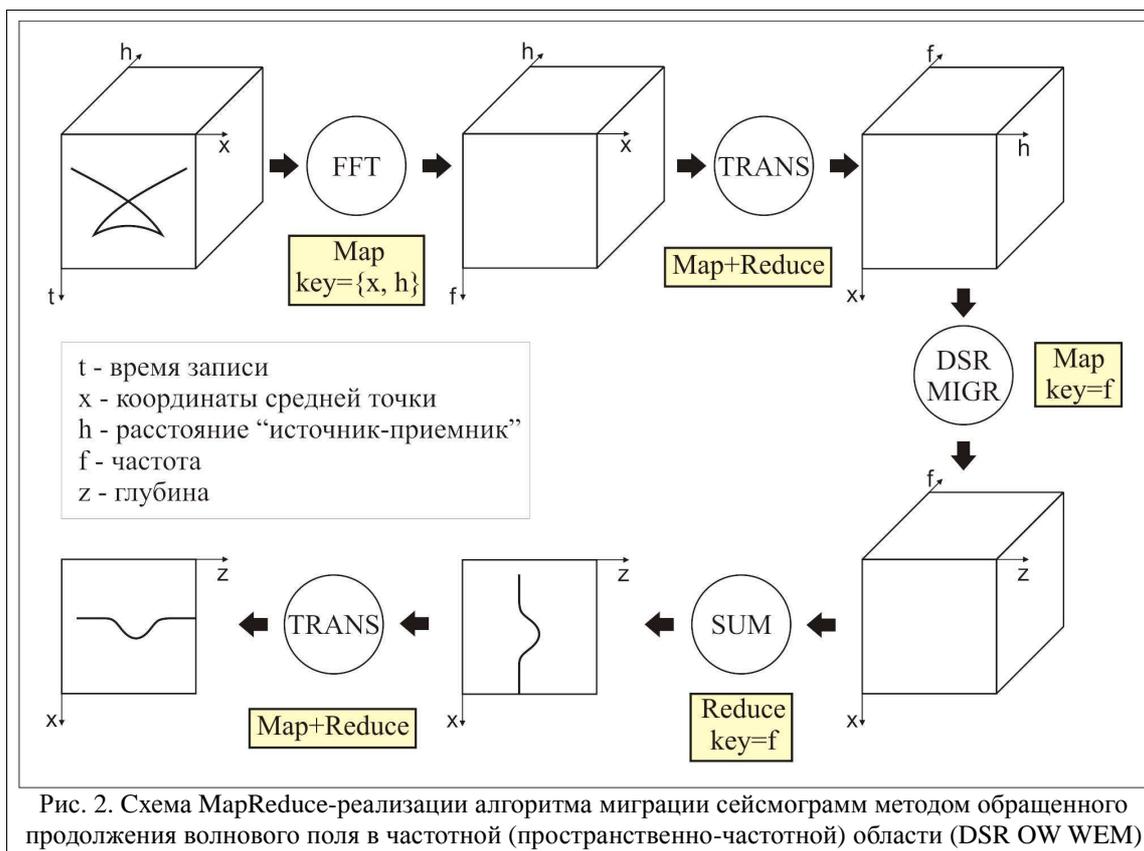


Рис. 2. Схема MapReduce-реализации алгоритма миграции сейсмограмм методом обращенного продолжения волнового поля в частотной (пространственно-частотной) области (DSR OW WEM)

Реализация сложных процедур, например, глубинной сейсмической миграции до суммирования, где для получения одного элемента выходных данных требуется доступ ко всему набору входных данных, требует нетривиальных подходов к организации потоков данных и вычислений. В качестве примера на рисунке 2 показана упрощенная схема MapReduce-реализации алгоритма миграции сейсмограмм методом обращенного продолжения волнового поля в частотной (пространственно-частотной) области (double-square-root one-way wavefield-extrapolation migration, [7]). Здесь производится обращенное продолжение волнового поля с заданным шагом по глубине и получение изображения как значения волнового поля на заданной глубине в момент времени, равном 0. Как видно из рисунка, алгоритм требует транспонирования многомерных массивов, что для старших индексов выполняется средствами сортировки MapReduce-среды. К сожалению, количество источников в литературе, посвященных решению подобных проблем, невелико. Отметим лишь работу [8], описывающую MapReduce-реализацию миграции Кирхгофа в модели средних скоростей, которая представляет собой один из простейших вариантов сейсмической миграции.

## ЗАКЛЮЧЕНИЕ

Использование распределенных («облачных») вычислительных сред для обработки сейсмических данных может дать существенный экономический эффект. Широкому практическому внедрению использования подобных вычислительных ресурсов препятствуют, в том числе, и технические причины, в первую очередь, отсутствие эффективных алгоритмов решения ресурсоемких задач. В настоящей работе предложен подход к реализации подобных алгоритмов.

Такие свойства, как отказоустойчивость работы программ в вычислительных средах, аналогичных Hadoop, и обеспечение высокой агрегированной скорости доступа к файлам, делают использование подобных программ привлекательным и для работы на специализированных центрах обработки сейсмических данных при решении задач, требующих длительного времени выполнения.

Рассмотренные программные реализации процедур и организации потоков обработки данных требуют проведения тщательных экспериментальных исследований на вычислительных системах различной конфигурации, в том числе, с эмуляцией отказов отдельных элементов систем. Проведение подобных экспериментов наряду с разработкой алгоритмов составляет предмет дальнейших исследований по теме настоящей работы.

Автор признателен М.С.Денисову, Е.Л.Музыченко, Е.Е.Куриной («ГЕОЛАБ») за содействие в процессе исследований и подготовки работы.

Работа проводилась при финансовой поддержке Министерства образования и науки Российской Федерации, Государственный контракт № 07.514.12.4007.

ЛИТЕРАТУРА:

1. <http://aws.amazon.com/ec2>
2. <http://cloud-computing-economics.com>
3. <http://en.wikipedia.org/wiki/MapReduce>
4. <http://hadoop.apache.org>
5. <http://www.cascading.org>
6. <http://www.cloudera.com/blog/2012/01/seismic-data-science-hadoop-use-case/>
7. Biondi B.L., 2006, 3D seismic imaging: SEG, Investigations in geophysics, 14.
8. Rizvandi N. B., Bolori A.J., Kamyabpour N., Zomaya A.Y., 2011, MapReduce Implementation of Prestack Kirchhoff Time Migration (PKTM) on Seismic Data: 12th International Conference on Parallel and Distributed Computing, Applications and Technologies, 200-206.