

RTM using Hadoop and Spark: Is there a case for migration?

L. I. Lumb | ian.lumb@brightcomputing.com | ianlumb@yorku.ca

Motivation

RTM is performance-challenged – Algorithms research remains topical and GPU implementations are delivering key results

Revisit RTM as a Big Data problem – In-memory analytics has the potential to improve performance of data and wavefield manipulations *in concert* with computations and introduce new prospects for imaging conditions

1

Key Performance Challenges

2

• **RTM modeling kernel is compute intensive** – Stable, non-dispersive solution via Finite Difference Modeling requires small time steps, small grid intervals, and higher-order approximations of the spatial derivatives

• **RTM wavefields exceed memory capacity**
Multiple-TB source volumes *must* be stored to disk

e.g., Liu et al., Computers & Geosciences 59 (2013) 17-23

Is there a case for migration?

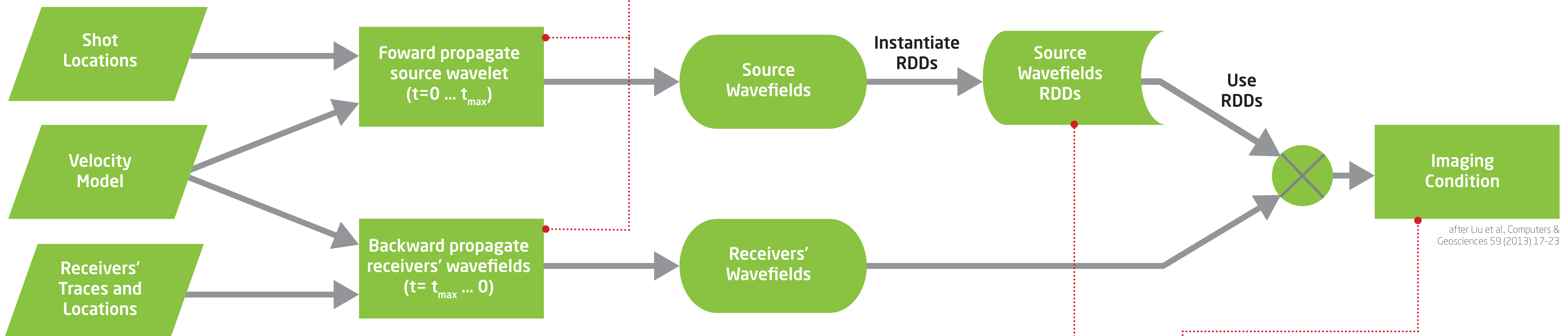
3

In-memory computing via RDDs is promising – Application to gathers and wavefields

Spark provides analytics upside – Imaging conditions other than cross-correlation

Spark may be applicable to modeling kernels

Spark can be easily incorporated into pre-existing IT infrastructures – Compliments existing HPC environments



after Liu et al., Computers & Geosciences 59 (2013) 17-23

RTM via RDDs: Implementation using Spark

Apache Spark is an implementation of RDDs

Make use of HDFS or alternative FS – GPFS, AWS S3, OpenStack Swift or Lustre (expected)

Choose appropriate programming model(s) – Not limited to MapReduce. Iterative and/or interactive (including streaming) programming models are supported. YARN and Mesos support available today, and support for HPC WLMs expected

Deployable on bare metal & clouds – Monitoring/management work-in-progress (e.g., Bright Cluster Manager role for Spark)

Introduces analytics possibilities for RTM – Program in Java (C/C++ via JNA), Scala or Python

Uptake is significant and community is growing

Results are extremely impressive – Exploit CPUs and/or GPUs

after Lumb, insideBIGDATA (in press)

4

RTM via RDDs: Opportunities

5

• Apply RDDs to gathers of seismic data – Partition RDDs optimally for wavefields calculations

• Apply RDDs to source wavefields – Partition RDDs optimally for cross-correlation of forward and reverse time wavefields. Significantly reduce/eliminate disk I/O

• Investigate alternate imaging conditions – Machine-learning and/or graph-analytics algorithms in addition to cross-correlation

Resilient Distributed Datasets (RDDs)

6

Abstraction for *in-memory* computing

Fault-tolerant, parallel data structures that are cluster-ready

Optionally persistent

Can be partitioned for optimal placement

Manipulated via operators

Zaharia et al., NSDI 2012