





#### ArrayUDF: User-Defined Scientific Data Analysis on Arrays

<u>Bin Dong</u><sup>1</sup>, Kesheng Wu<sup>1</sup>, Surendra Byna<sup>1</sup>, Jialin Liu<sup>1</sup> Weijie Zhao<sup>2</sup>, Florin Rusu<sup>2</sup> <sup>1</sup>LBNL, Berkeley, CA <sup>2</sup>UC Merced, Merced, CA

HPDC 2017, Washington D.C. June 28, 2017.

#### Scientific Activities Evolved into Big Data Analysis

Example: scientific projects for dark matter/energy, supernovae, etc.



#### **Data Size for Sky Survey Porjects**





## Q1: How many data analysis operations that have been and will be developed ?





## Q1: How many data analysis operations that have been and will be developed ?



Large population

#### Implication from popular data analysis languages

	Total amount	Avg Growth	
Python modules	110,486	69/day	
R packages	10,877	8/day	

\*Data from <a href="http://www.modulecounts.com/">http://www.modulecounts.com/</a> on June 21 2017





## Q1: How many data analysis operations that have been and will be developed ?

Large population

Q2: What are the functions of these data analysis operations ?





#### Q1: How many data analysis operations that can be used to extract scientific meaning? Large population Q2: What are the functions of these data analysis operations? Variety Smoothin IFFERENCE Optimization Learning Scatterpl ression Filtering UNION Sug ence Estimato Prediction **Back-Propagati Deep Belief Networ** Evaluating integra RESTRI Differential equati Locally Classificat inear Regression Local Naive I Naive I Bayesian N Stepwise Regression Mul-Averaged One-Depend

**FRKELEY** 

# Two Methods to Embrace Large Population & High variety in Data Analysis Operations





#### **UDF** is at Heart of Modern Big Data system

Examples: MapReduce in Apache Hadoop and Spark







#### MapReduce Doesn't Fit Scientific Data Analysis

Reason 1: Most Scientific Data are Multi-dimensional Arrays



Converting array to (key, value) is expensive



Pictures Credit: Kyle Hemes, Peter Nugent, Suren Byna, etc.



#### MapReduce Doesn't Fit Scientific Data Analysis

Reason 1: Most Scientific Data are Multi-dimensional Arrays Reason 2: Most Scientific Data Analysis Operations Own Structure Locality Property

Structure Locality:

The analysis operation on a single cell accesses its neighborhood cells

Map deals with a single element at a time

- Reduce requires to duplicate an each cell for all neighborhood cells
- Reduce only happens after expensive shuffle





2D Poisson Equation Solver (Discrete)







#### ArrayUDF: User-Defined Scientific Data Analysis on Arrays

- Stencil based User-defined Function
  Structural locality aware array operations
- Native Multidimensional Array Data Model
  - In-situ data processing in scientific data formats, e.g., HDF5
- Optimal Chunking and Ghost Zone Method
  - Efficiently parallel array processing on HPC system





### **Stencil Definition**

• Stencil S Definition for *d*-dimensional array

$$S = \{ c_{i_1 + \delta_1, \bullet \bullet \bullet, i_d + \delta_d, | \delta \in [L, R], L \in [-i, 0], R \in [0, N - i] \}$$

Materialized structure locality

Flexible neighborhood expression

- Centre of Stencil:  $S_{0,\bullet\bullet\bullet,0} = C_{i_1,\bullet\bullet\bullet,i_d}$
- Element at offsets  $\delta_{1}, \bullet \bullet \bullet, \delta_{d}$ :

$$S_{\delta_1,\bullet\bullet\bullet,\delta_d} = C_{i_1+\delta_1,\bullet\bullet\bullet,i_d+\delta_d}$$

2D Example:





#### **Stencil based Computing Model**

$$c'_{i_{1}, \dots, i_{d}} = f\left(\{s_{\delta_{1}, \dots, \delta_{d}} \mid \delta \in [L, R]\}\right)$$
  
The cell in array A'  
at coordinate  $i_{1} \dots i_{d}$   
The stencil set in array A  
at offset:  $\delta_{1} \dots \delta_{d}$ 

- A = A' when in-situ updates
- Dim(A) = Dim(A') in most cases





### **Stencil based Computing Model vs Others**

	Input	Output	UDF
SQL	Tuple <i>t</i>	Tuple <i>t</i> '	$t'=\boldsymbol{f}(t)$
SciDB	Cell c	Cell c'	c'= <b>f</b> (c)
MapReduce	KeyValue kv	KeyValue kv"	kv'=Map(kv) kv''=Reduce(kv' <sub>1</sub> , kv' <sub>2,</sub> )
ArrayUDF	Stencil s	Cell c'	c'= <b>f</b> (s)

*vs.* MapReduce:

ArrayUDF generalizes its two steps as a single step on array *vs.* SciDB:

ArrayUDF supports structure-locality based computing on array





#### **Examples of ArrayUDF**



**Example 2: Vorticity computation** 



## **UDF Support System: Chunking**

- Chunking enables parallel and out-of-cores processing
  - general chunking (layout unknown)
    - minimize ghost cells







#### **UDF Support System: Ghost Zone handling**

- ArrayUDF processes chunks in parallel and/or in out-of-core manner
- Ghost zone avoids communications between chunks
- •Ghost zone size might be unknown in advance
  - UDF source code might be unavailable
- •Trail-run is used to find ghost zone size:
  - Run UDF on a sample Stencil instance, that collects the offsets applied within UDF





#### **Evaluations**

- Hardware:
  - Edison, a Cray XC30 supercomputer at NERSC
  - 5576 computing nodes, 24 cores/node, 64GB DDR3 Memory
- Software
  - ArrayUDF 0.0.1
  - Spark 1.5.0
  - SciDB 16.9

- RasDaMan 9.5 (sequential version)
- EXTASCID, hand-optimized version
- Hand-optimized C/C++ code

- Workloads
  - Two synthetic data sets (i.e., 2D and 3D array) for micro benchmarks
    - Chunking strategy, trail-run, etc.
  - Four real scientific data sets (i.e., S3D, MSI , VPIC , CoRTAD)
    - Overall performance tests /w generic UDF interface



### **Chunking Strategy Evaluation**

- general chunking (for average cases)
  - minimize ghost cells # to reduce I/O cost
- layout-aware chunking (for layout special case)
  - maximize contiguous disk read
  - ignore the impact of ghost zone





# Comparison with peer systems with standard "window" operators

- "window" comes from SCiDB and RasDaMan
- "window" supports certain structure locality but lack the link to
  UDF function
  BasDaMan(Sequential)
  Spark

Poisson equation solver /w Stencil S of ArrayUDF 2D : 4S(0,0)-S(-1,0)-S(0,1)-S(1,0)-S(-1,0)3D : 6S(0, 0,0)-S(-1,0,0)-S(0, 1, 0)-S(1,0, 0)-S(-1,0,0)-S(0,0,-1)-S(0, 0,1)



- ArrayUDF has close performance to hand-optimized code
- ArrayUDF is at least 13X faster than peer systems





# Comparison with Spark in supporting real applications operations



# of local cells used by an UDF /w generic interface



#### Conclusions

- ArrayUDF: User-Defined Scientific Data Analysis on Arrays
  - Stencil based UDF for structural locality-aware operations
  - •Native array model & In-situ array processing in HDF5, etc.
  - Optimal chunking and ghost zone methods for large array
- ArrayUDF provides close performance to hand-optimized code
- ArrayUDF is as much as 2070X faster than Spark
- ArrayUDF source code: <u>https://bitbucket.org/arrayudf/</u>
- Future work
  - Python and other language interface
  - Communication optimizations



### Acknowledgments

- Nicholas Chaimov from University of Oregon for suggestions to set up Spark on Editon at NERSC
- Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, support for the SDS project (Program manager: Dr. Lucy Nowell) under contract number DE-AC02-05CH11231



National Energy Research Scientific Computing Center









Thanks

#### Bin Dong dbin@lbl.gov http://crd.lbl.gov//dongbin





#### **Trail-run overhead**

- Detect ghost zone size automatically
  - Run the UDF on a single Stencil but the UDF might access more neighborhood cells

	The number of cells used by UDF							
Data sets	4	8	16	32	64	128	256	
S1	0.37	0.38	0.46	0.48	0.54	0.59	0.80	
S2	0.48	0.52	0.65	0.75	0.79	0.84	1.04	

Unit: microsecond

≈ 1 ms when 256 cells are used in the UDF



