# SciSpark 201

Searching for MCCs

# Agenda for 201:

- Access your SciSpark & Notebook VM (personal sandbox)
- Quick recap. of SciSpark Project
  - What is Spark?
- SciSpark Extensions
  - sciTensor:  N-dimensional arrays in Spark
  - sRDD:  scientific RDD's
- Implementation of Mesoscale Convective Complexes (MCC) Analytics
  - Find cloud elements, analyze evolution, connect the graph
- Tutorial Notebooks:
  - 201-1:  Using sciTensor and sRDD
  - 201-2:  More SciSpark API & Lazy Evaluation
  - 201-3:  End-to-end analysis of Mesoscale Convective Systems
  - 201-4:  MCS Cloud Element Analytics

# Claim your SciSpark VM

- Go to:  http://is.gd/scispark
  - You'll have to click through!
- Enter your name next to an IP address to claim it
- Enter the URL in your browser to access the SciSpark Notebook
  - Chrome or Firefox preferred

# Agenda for 301:

- Access your SciSpark & Notebook VM (personal sandbox)
- Quick recap. of SciSpark Project
- PDF Clustering Use Case
  - K-means Clustering of atmospheric behavior using a full probability density function (PDF) or just moments: histogram versus (std dev, skewness)
- SciSpark Challenges & Lessons Learned
- Tutorial Notebooks:
  - 301-1, 2, 3: PDF Clustering in PySpark, in Scala, Side by Side
  - 301-4: Regional Climate Model Evaluation (RCMES)     -  (optional)
  - 301-5: CMDA Workflow Example
  - 301-6: Lessons Learned Notebook
- Wrap-up Discussion: How do you do X in SciSpark?
  - Bring your Use Cases and questions

# Discussion in 301:
## How do you do X in SciSpark?

- What are your Use Cases?
  - Bring them to SciSpark 301
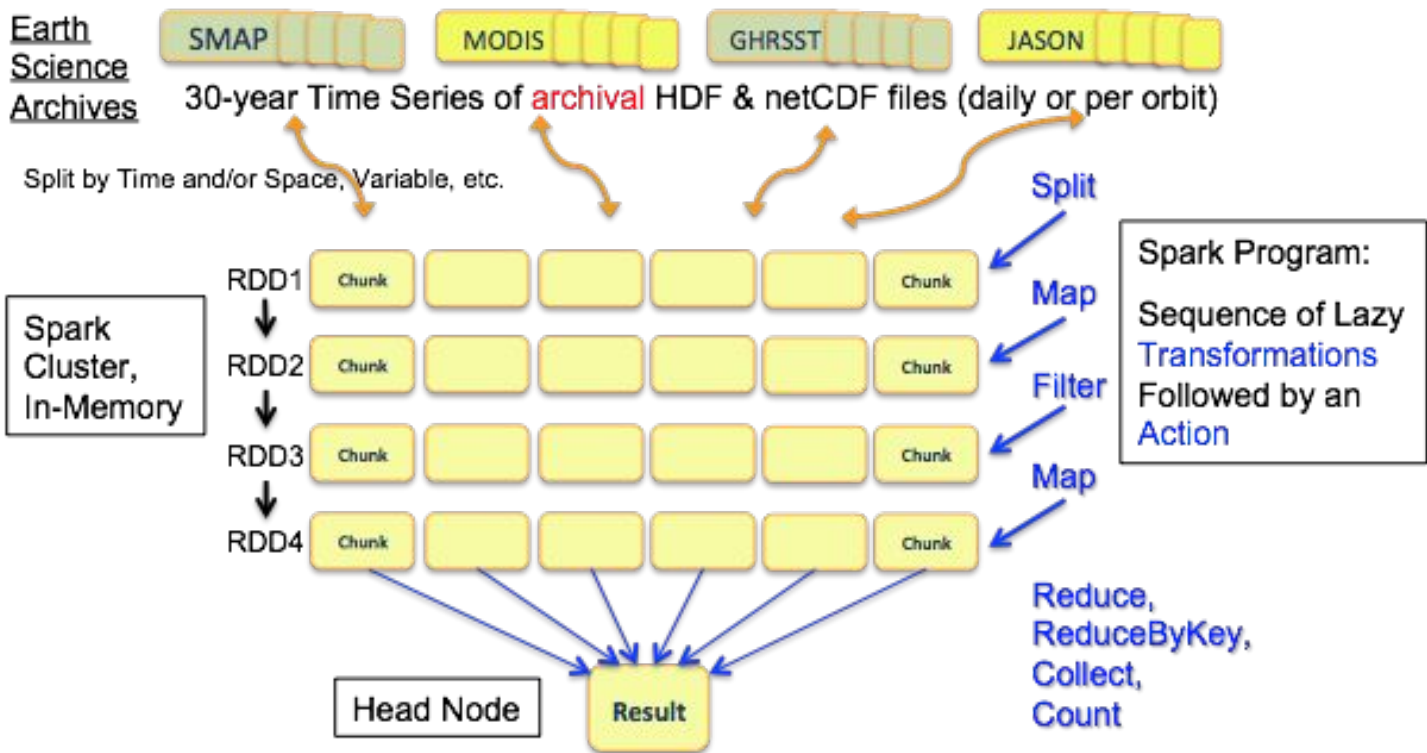- Who is using Spark or Hadoop?
- Questions?

# SciSpark:  A Two-Pronged Approach to Spark

- Extend Native Spark on the JVM (scala, java)
  - Handle Earth Science geolocated arrays (variables as a function of time, lat, lon)
  - Provide netCDF & OPeNDAP data ingest
  - Provide array operators like numpy
  - Implement two complex Use Cases:
    - Mesoscale Convective Systems (MCS) - discover cloud elements, connect graph
    - PDF Clustering of atmospheric state over N. America

- Exploit PySpark (Spark jobs using python gateway)
  - Exploit numpy & scipy ecosystem
  - Port and parallelize existing python analysis codes in PySpark
    - Statistical Rollups over Long Time-Series
    - Regional Climate Model Evaluation System (RCMES)
    - Many others

# Sources of Parallel Work (Independent Data)

- Parallelize Over Time:
  - Variables (grids) over a Long 10-30 Year Time-Series (keep grid together)
- Parallelize Over Space (and/or Altitude):
  - By Pixel: Independent Time-Series, Massive Parallelism
  - By Spatial Tiles: Subsetting, Area Averaging, Stencils (NEXUS)
  - By Blocks: Blocked Array operations (harder algorithms)
- Parallelize Over Variable, Model, Metrics, Parameters, etc.
  - Ensemble Model evaluation with multiple metrics
  - Search Parameter Space with many analysis anruns

# What is Apache Spark? - In-memory Map-Reduce

# Apache Spark Transformations & Actions

| | | | |
|---|---|---|---|
| **Transformations** | $map(f : T \Rightarrow U)$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $filter(f : T \Rightarrow Bool)$ | : | $RDD[T] \Rightarrow RDD[T]$ |
| | $flatMap(f : T \Rightarrow Seq[U])$ | : | $RDD[T] \Rightarrow RDD[U]$ |
| | $sample(fraction : Float)$ | : | $RDD[T] \Rightarrow RDD[T]$ (Deterministic sampling) |
| | $groupByKey()$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$ |
| | $reduceByKey(f : (V, V) \Rightarrow V)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $union()$ | : | $(RDD[T], RDD[T]) \Rightarrow RDD[T]$ |
| | $join()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (V, W))]$ |
| | $cogroup()$ | : | $(RDD[(K, V)], RDD[(K, W)]) \Rightarrow RDD[(K, (Seq[V], Seq[W]))]$ |
| | $crossProduct()$ | : | $(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$ |
| | $mapValues(f : V \Rightarrow W)$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, W)]$ (Preserves partitioning) |
| | $sort(c : Comparator[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| | $partitionBy(p : Partitioner[K])$ | : | $RDD[(K, V)] \Rightarrow RDD[(K, V)]$ |
| **Actions** | $count()$ | : | $RDD[T] \Rightarrow Long$ |
| | $collect()$ | : | $RDD[T] \Rightarrow Seq[T]$ |
| | $reduce(f : (T, T) \Rightarrow T)$ | : | $RDD[T] \Rightarrow T$ |
| | $lookup(k : K)$ | : | $RDD[(K, V)] \Rightarrow Seq[V]$ (On hash/range partitioned RDDs) |
| | $save(path : String)$ | : | Outputs RDD to a storage system, $e.g.$, HDFS |

Table 2: Transformations and actions available on RDDs in Spark. Seq[T] denotes a sequence of elements of type T.

# SciSpark Project Contributions

- **Parallel Ingest** of Science Data from netCDF and HDF
  - Using OPeNDAP and Webification URL's to slice arrays
- Scientific RDD's for large arrays (sRDD's)
  - Bundles of 2,3,4-dimensional arrays keyed by name
  - Partitioned by time and/or space
- More Operators
  - ArraySplit by time and space, custom statistics, etc.
- Sophisticated Statistics and Machine Learning
  - Higher-Order Statistics (skewness, kurtosis)
  - Multivariate PDF's and histograms
  - Clustering, Graph algorithms
- Partitioned Variable Cache (in development)
  - Store named arrays in distributed Cassandra db or HDFS
- Interactive Statistics and Plots
  - "Live" code window submits jobs to SciSpark Cluster
  - Incremental statistics and plots "stream" into the browser UI

# Spark Split Methods

- Sequence files:  text files, log files, sequences of records
  - If in HDFS, already split across the cluster (original use for Hadoop & Spark)
- sc.parallelize(list, numPartitions=32)
  - Split a list or array into chunks across the cluster
- sc.textFile(path, numPartitions=32),    sc.binaryFile()
  - Split a text (or binary) file into chunks of records
  - Each chunk is set of lines or records
- sc.wholeTextFiles(), sc.wholeBinaryFiles()
  - Split list of files across the cluster, whole files are kept intact
- Custom Partition methods
  - Write your own

# SciSpark Extensions for netCDF

- sciSparkContext.NetcdfFile([path or DAP URL], numPartitions=32)
    - Read multiple variables (arrays) & attributes into a sciTensor
- sciSparkContext.NetcdfDFSFile([list of netCDF files in local dir. or HDFS])
    - Split a list or array into chunks across the cluster
- sciSparkContext.openPath([all netCDF files nested under top directory])
- sciSparkContext.readMERGFile()
    - Read a custom binary file

# What is the goal of scientific RDD (sRDD)?

**The goal of the sRDD was to:**
    Create a scientific partition aware dataset that could be:
1. ETL'ed from HDF, NetCDF
2. Split by Time/Region
3. Mapped e.g., in Regrid, or Metric
4. Save/Cache'd

**Example Scenarios**
Two scenarios demonstrating intelligent data caching and access in SciSpark.
A) a multi-stage operation to generate a time split, B) a multi-stage operation to select data from Shark (now SparkSQL), and cluster by deviation from mean.
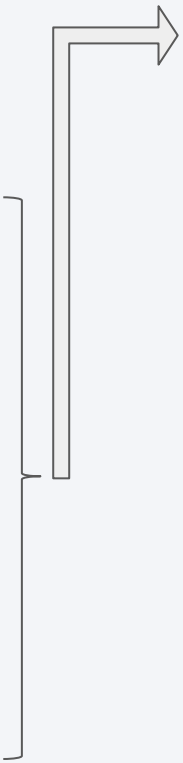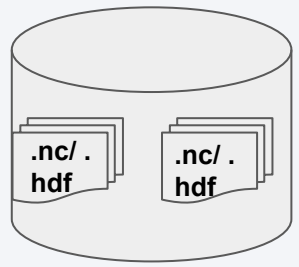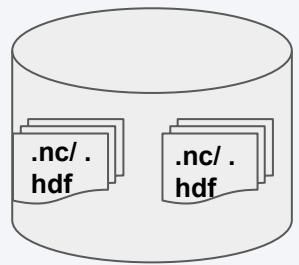Automatically capture: **provenance**

# Why SciSpark?

To understand the climate variability and change, and extreme weather events we need to be able perform complex analysis on voluminous and heterogeneous highly-resolved datasets. Analysis involves:

- Acquiring and ingesting data
- Resolving profiles and point sources within spatial and/or temporal coverage
- Gridding the data (for analysis and/or comparison)
- Quantifying changes in resolutions according to regridding techniques Calculating a range of statistics
- Visualizing data at all steps
- Ensuring reproducibility and data provenance



Work funded under 14-AIST-14-0034 - PI: Christian Mattmann

**Typical atmospheric scientist's analysis workflow**

Get data from remote source(s)

.nc/ .hdf

.nc/ .hdf

.nc/ .hdf

.nc/ .hdf

Store data locally

.nc/ .hdf

.nc/ .hdf

Load & regrid data

.nc/ .hdf

.nc/ .hdf

Regridded data

.nc/ .hdf

.nc/ .hdf

Calculated data

.nc

Visualize and analyze

while studying feature do
    for *each time*
        do calc. (simple/complex)
        **write to disk**
    end for
end while

# Challenges of handling Earth Science data in Hadoop environment

- Storage on the Apache Hadoop Distributed FileSystem (HDFS)
  - How to chunk up the file?
  - How to re-assemble it?
- How to access the data?
  - Loading **self describing, multidimensional & multivariable** data
  - How to keep that metadata?
  - How to extract variable from file?
  - Handling multidimensional data as Hadoop env is Java based
- How to compute on the data?
  - Matrix and scalar operations
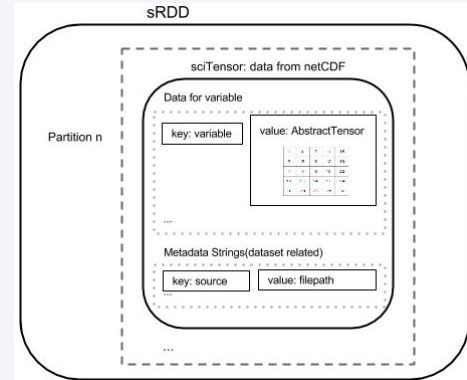  - Basic, complex statistics and user-defined algorithms

Not **SciSpark**'s focus
See: SciHadoop, sciDB

# The sRDD

- The sRDD is an immutable, partitioned collection of **sciTensor** elements, that can be operated on in parallel within the cluster environment.
- The sRDD **extends** the RDD i.e. it has all the non-private variable & methods available to the RDD e.g. getPartitions, persist, transformations: map, filter, etc.
- All transformations in the sRDD are **lazy**, which helps to optimize the disk and memory usage of the cluster.
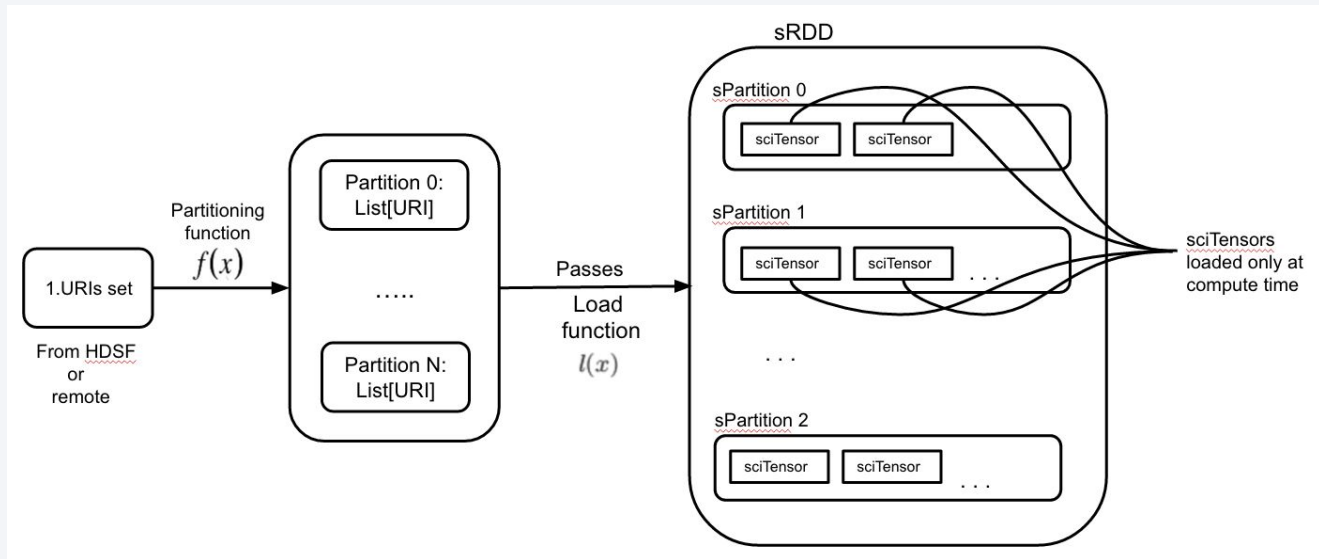
# The sciTensor

- The sciTensor is the datatype used in SciSpark to read in data from hierarchical data formats such as netCDF. The sciTensor:
    - is a self-documented array
    - keeps a list of variable arrays
    - keeps metadata in a hashmap
    - users can query the hashmap and perform matrix operations



- The sciTensor is required in the SciSpark environment because:
    - SciSpark/ Spark is built on the Hadoop environment that is Java based, which inherently does not handle multidimensional arrays well (as compared to C, Fortran, Python)
- The sciTensor is read into the sRDD (i.e. the memory of the cluster) and the data within is operated on via arithmetic and relational operations.

# Loading data into the sciTensor

- Data can be loaded into the sciTensor from:
    - Apache Hadoop Distributed FileSystem (HDFS)
    - OpenDap server
    - Local Filesystem

# SciSpark usecase: Mesoscale convective complexes

| PHYSICAL CHARACTERISTICS | |
|---|---|
| **SPATIAL RESOLUTION** | A – Cloud shield with continuously low IR temperature $\leq$ -32°C must have an area $\geq$ 100 000km$^2$<br>B- Interior cold cloud region with temperature $\leq$ -52°C must have an area $\geq$ 50 000km$^2$ |
| **MAXIMUM EXTENT** | Contiguous cold cloud shield (IR temperature $\leq$ -32°C) reaches maximum size |
| **SHAPE** | Eccentricity $\geq$ 0.7 at time of maximum extent |
| **DURATION** | A and B must be met for a period $\geq$ 6hrs |



2. Locations of MCCs based upon 1-3 year regional samples of satellite imagery (from Laing and Fritsch, 1997). Locations are shown for the time of maximum extent of the cold cloud shield. Additional data concerning the samples is given in Appendix A.
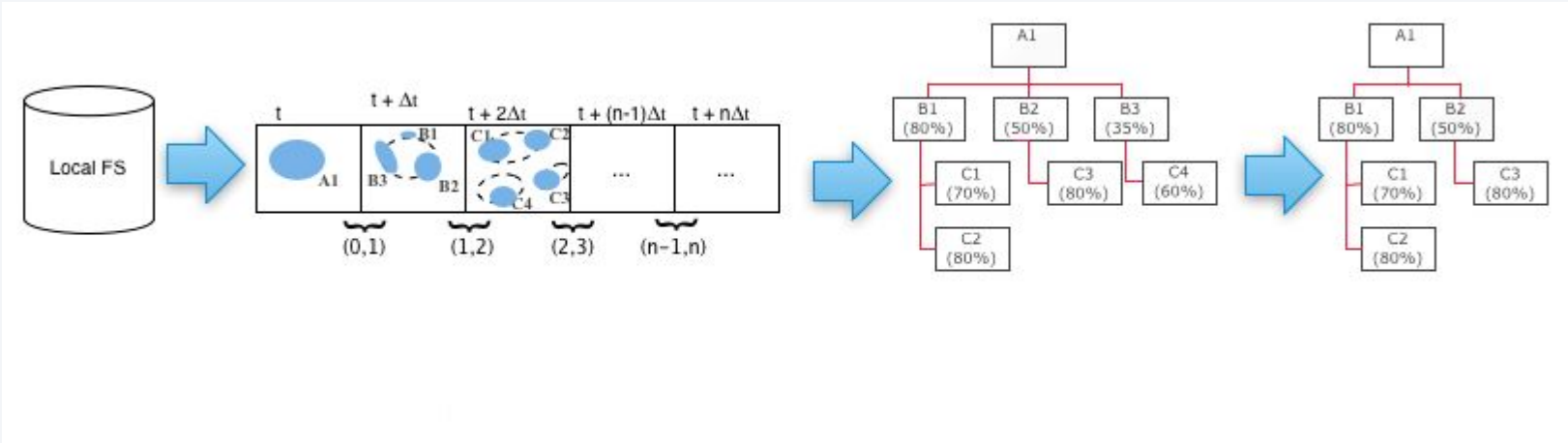
- A subset of mesoscale convective systems (MCSs)
- More than 700 MCCs occur globally, with 66% in the NH
- Maddox 1980 criteria for identification

# Finding MCSs



- Uses brightness temperature data.
- Nodes are areas with a given brightness temperature value, and of a given size.
- Edges are determined by area overlaps between nodes within consecutive time periods.

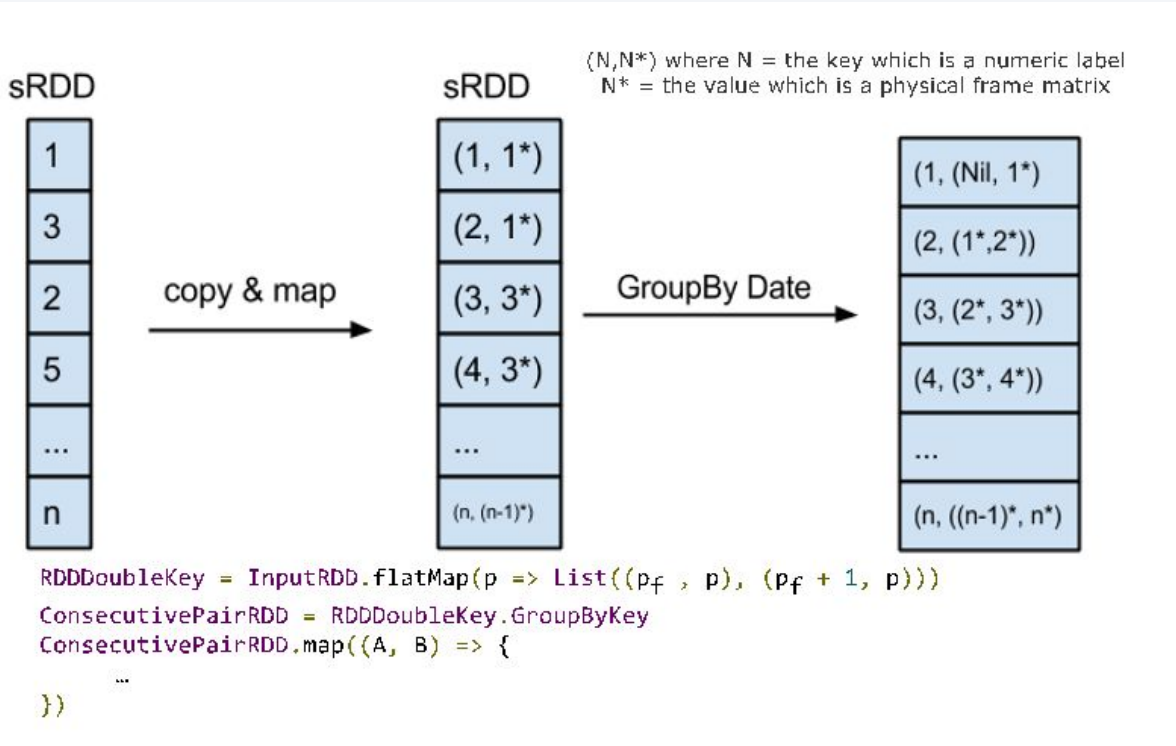# Finding cloud elements and overlaps

- Identify areas with a given brightness temperature value, and of a given size - nodes (in a parallelized manner)
- Edges are determined by area overlaps between nodes within consecutive time periods (this can also be parallelized in time).



**Figure 3. Graph construction from Whitehall et al. [26] GTG algorithm for MCC identification over Niamey, Niger.**

# Finding graph edges between two time epochs

- Bullet
- 



$(N, N*)$ where $N$ = the key which is a numeric label
$N*$ = the value which is a physical frame matrix

sRDD

| 1 |
| 3 |
| 2 |
| 5 |
| ... |
| n |

copy & map →

sRDD

| $(1, 1*)$ |
| $(2, 1*)$ |
| $(3, 3*)$ |
| $(4, 3*)$ |
| ... |
| $(n, (n-1)*)$ |

GroupBy Date →

| $(1, (Nil, 1*))$ |
| $(2, (1*, 2*))$ |
| $(3, (2*, 3*))$ |
| $(4, (3*, 4*))$ |
| ... |
| $(n, ((n-1)*, n*))$ |

```
RDDDoubleKey = InputRDD.flatMap(p => List((p_f , p), (p_f + 1, p)))
ConsecutivePairRDD = RDDDoubleKey.GroupByKey
ConsecutivePairRDD.map({A, B) => {
    ...
}})
```

# Traversing graph by recursive merge

- Smartly partition: nodes with greatest connectivity
- For each partition: eliminate only "full" graphs
- Merge consec. partitions and repeat above step
- Iterate until one partition or nodes = 0

# Complexities of this use case:

- Keeping time sequence in the sRDD: Our graph has an implicit time dependence, thus we need to consider the how the problem will be parallelized
- Keeping data & metadata flowing through the algorithm: As nodes are determined, there is some "irregular" subsetting of the initial dataset.
- Traversing the graph in the SciSpark engine. Our graph is:
  - Sparse, disconnected and directed.
  - There are no cycles

# Notebook 201-3:

## End-to-end MCS part 1

- ○ Loading data from HDFS
- ○ Maintaining the time sequence
- ○ Creating the graph object: Identifying cloud elements and edges
- ○ Visualizing the cloud elements from the run leveraging SparkSQL

# Notebook 201-4:

## End-to-end MCS part 2

- Traversing the graph object to identify MCSs (i.e. subgraphs)
- Visualizing the graph and results (D3)

# Notebook Exercises

- Try some of the exercises
- Just Play
- Ask Questions
- Notebook exercises:
  - 201-1: Using sciTensor and sRDD
  - 201-2: More SciSpark API & Lazy Evaluation
  - 201-3: End-to-end analysis of Mesoscale Convective Systems
  - 201-4: MCS Cloud Element Analytics