# In this session…

- IBM Spark@SETI

- Deep Learning meets SETI Science

- Beyond Classification – Novel Observations and Analytics

- The SETI Project – By The Numbers

- Back to Earth… Spark and Object Store
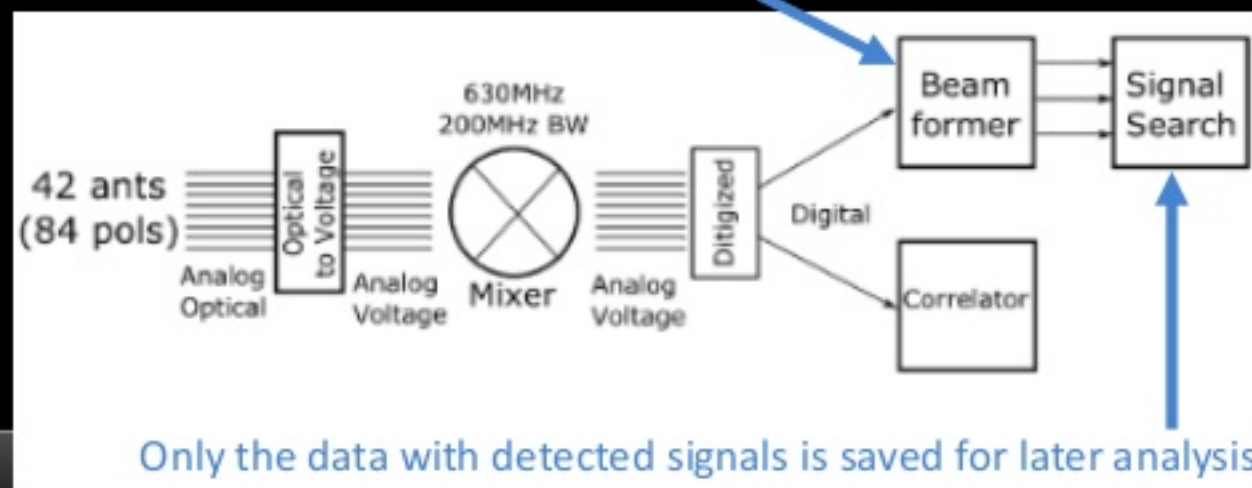
- Conclusions and Take-Aways

# IBM Spark@SETI

## SETI Institute Backgrounder

- Headquartered in Mountain View, CA.  Founded 1984.  150 Scientists, researchers and staff.

- The mission of the SETI Institute is to explore the potential for extra-terrestrial life…. search for narrow band radio signals in the frequency range of 1GHz to 10GHz which could be evidence of intelligence outside our solar system.

- Allen Telescope Array (ATA) – Phased Array Synthetic Dish – 3 Beams

The Allen Telescope Array

42 Receiving Dishes
Each 6m diameter
1GHz to 10GHz

4.5TB data every hour

42 ants
(84 pols)

Analog
Optical

Optical
to Voltage

Analog
Voltage

Mixer

630MHz
200MHz BW

Analog
Voltage

Digitized

Digital

Beam
former

Correlator

Signal
Search

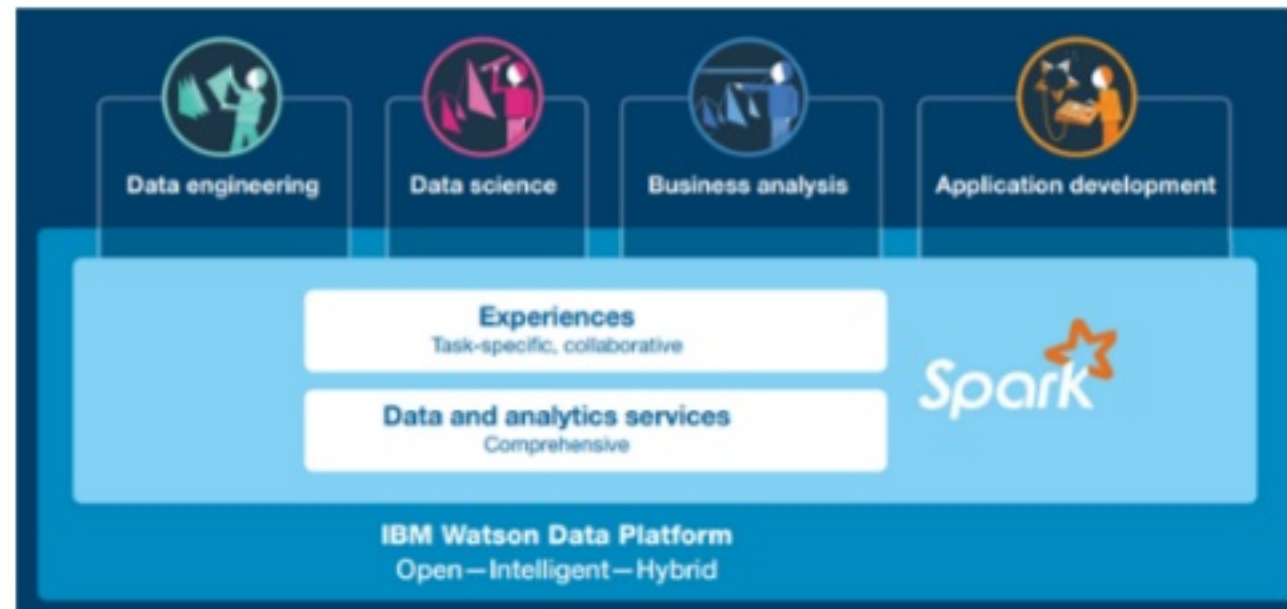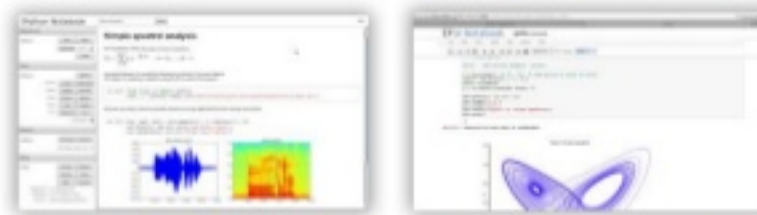Only the data with detected signals is saved for later analysis

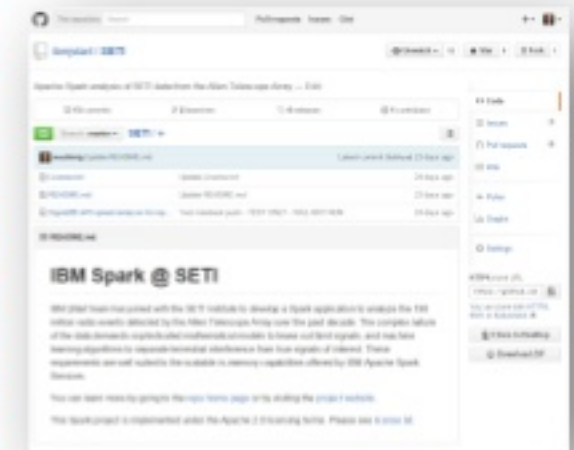# Spark@SETI – Deep analytic capabilities for professional astronomers



**IBM Data Science Experience**

Jupyter

github
SOCIAL CODING

- IBM Spark@SETI GitHub repository

- Jupyter notebooks
- Python libraries for Spark@SETI

Data engineering · Data science · Business analysis · Application development

**Experiences**
Task-specific, collaborative

**Data and analytics services**
Comprehensive

Spark

**IBM Watson Data Platform**
Open – Intelligent – Hybrid

**IBM Spark @ SETI**

STOCATOR

SOFTLAYER

*Import of signal data from SETI radio telescope data archives*

SETI INSTITUTE

**IBM**
**Object Storage**

IBM dashDB

Shared repository of SETI data in Object Store
- 200M rows of signal event data
- 360,000 raw recordings of "signals of interest"
- Large "long duration" observations (~5TB each)
- ~40TB accessible data in storage

# Spark@SETI Deep Learning Signal Classifier

Jupyter notebook showing complex radio signals being classified based on morphology (shape) and other features.

Neural net model was trained on the IBM Cognitive Compute Cluster (750 NVIDIA K80 GPUs) and imported into IBM Data Science Experience

# From raw antennae voltage data streams to classified signals, all with unsupervised machine learning...

# From raw antennae voltage data streams to classified signals, all with unsupervised machine learning…



(NOT REAL - ILLUSTRATIVE ONLY!)

# Beyond Classification...

- Classification is great for intentional "beacon" signals (and for getting rid of RFI - radio frequency interference)

- Searching for beacons assumes we are "important enough" to do that ... we have to look for signals which were not intended for us

- Leakage and Eavesdropping – much harder to detect, data and compute heavy

# TRAPPIST-1 – One stop shopping for ET

- Newly discovered system with 7 rocky planets, three in the habitable zone (mathematically, liquid water could be on the surface)

- Insanely good opportunity for eavesdropping and leakage detection

- Detected by transit (Kepler-K2 )… orbital plane is within 0.3 degrees of perfect alignment with our line of sight = occultations and conjunctions

- Close… only 40 LY

- Recently found… propagated error for orbital predictions is around 2 seconds

- Tight orbits (2-7 days) means LOTS of opportunities to see conjunctions and occultations



Trappist-1 Observation of Planets e and g conjunction @ 15:07UTC and planets f and g at 15:36UTC

- Trappist-1
- Trappist-1 b
- Trappist-1 c
- Trappist-1 d
- Trappist-1 e
- Trappist-1 f
- Trappist-1 g
- Trappist-1 h

Show Now
Show e/f/g Conjunction
Real-Time
Fast-Time

Tue, 18 Apr 2017 21:02:53 UTC
ATA Az: 241.50°, Elev: 22.16°

Earth

# TRAPPIST-1 … so here's the latest news:

- System model implemented in Python and run on IBM DSX / Spark… a "*conjunction and occultation finder*" accurate to < 45 seconds

- Six observations completed and data queued for trickle upload using Stocator

- Occultation ground-to-cloud completed for power curve analysis… don't expect anything definitive until we have data folded several observations together.



Power over time for file t170425-yy-vector.dat

TRAPPIST-1 g OCCULTATION

13:53 UTC
4/25/2017

21:07 UTC

7 hours 14 minutes

# The Spark@SETI Project – By the Numbers

- **200 million** signal events

  - Doppler shift corrected in 22 minutes wall time on 30 executor Spark cluster

- **14 million** complex amplitude files in Object Store

  - Each binary file contains 90 second 'snapshot' of raw antennae voltages

  - 14M files = 1TB of raw signal data... FFT into 14M spectrograms
    followed by feature extraction for clustering **~12 hours**

- Long duration observations = 2 beams @ 2.5TB each

  - Auto-correlation and SWAC analysis looking for wideband... best case is
    order $n \cdot \log(n)$, Symbol-Wise Autocorrelation (SWAC) is order nD, where
    D is the number of delay values calculated (which is generally less than n)

  - Wide-band analysis.... **5TB processed** for wideband detection in approximately **13.5
    hours wall time**.

# How did we make all this work?

- High performance access to large amounts of data in Object Store was the game changer

- Over to Gil…

# Back to the Earth

- What technologies SETI@Spark depends on?
  - Apache Spark
  - Data objects in the variety of formats
  - IBM Cloud Object Storage

# What is an object store?

- Object store is a perfect place to store files (we call them objects)

- Each object consists of rich metadata and the data itself

# Organize data in the object storage

- Data objects are organized inside the buckets or containers

- Each data object may contain a name with delimiters, usually "/"

- This allows to group objects inside buckets (pseudo directories), an analogy to the directories in file systems but without the overhead or scalability limits of lots of directories

data object

bucket

mytalks/year=2016/month=5/day=24/data-palooza.pdf
mytalks/year=2017/month=5/day=24/hadoop-strata.pdf
mytalks/year=2017/month=6/day=07/spark-summit.pdf

# Object storage in more details

- Capable to store huge amounts of unstructured data in any format

- Resilient store – storage data will not be lost

- Object store designed to operate during failures

- Various security models – storage data is safe protected

- Object stores can be easily accessed for write or read flows

- On premise, cloud based, hybrid, etc.

- Analytic job results can be persisted in the object storage

- Object stores allows easily to share data subsets with others

# How Spark access object stores?

Apache Spark utilize various connectors from the Hadoop project

| swift:// | s3a:// | wasb:// |
|---|---|---|
| Hadoop OpenStack connector | Hadoop AWS connector | Hadoop Azure connector |

Object Storage ( S3 API, Swift API, Azure API)

# Example: persist collection as an object

```
val data = Array(1, 2, 3, 4, 5, 6, 7, 8, 9)
val myData = sc.parallelize(data, 9)
myData.saveAsTextFile("s3a://mybucket/mydata.txt")
```

| | API | GET | HEAD | PUT | DELETE |
|---|---|---|---|---|---|
| Hadoop (s3a) | S3 | 158 | 361 | 26 | 16 |

**561 total requests to the object storage**

# Behind the numbers

- We observed that the usage of Hadoop connectors with object store is highly inefficient

- Hadoop connectors adapted for file systems and not object stores

- What is wrong?  - two major reasons

  - The existing algorithms used by Hadoop to achieve fault tolerance for persisting distributed data sets

  - Cost of supporting FS shell operations and treating object store as a file system

# Fault tolerance algorithms in the write flows

- Output committers responsible for persisting data sets generated by MapReduce jobs. Output committers uses temp files and folders for every write operation and then renames them.

- File systems has atomic rename, which perfectly fits into this paradigm.

- Object stores do not support rename natively; use copy and delete instead. This leads to dozens of expensive requests targeted to the object store.

..result/_temporary/0/_temporary/attempt_201702221313_0000_m_000000_0/part-0000

..result/_temporary/0/task_201702221313_0000_m_000000/part-00000

..result/part-00001

# Hadoop FS shell operations

- Hadoop connectors to be 100% compliant with the Hadoop ecosystem must support FS shell operations on files/directories

  ```
  ./bin/hadoop fs –mkdirs hdfs://myhdfs/a/b/c/
  ./bin/hadoop fs –put mydata.txt hdfs://myhdfs/a/b/c/data.txt
  ```

- Object store vendors provide CLI tools that are preferable over Hadoop shell fs commands

- Analytic flows do not need FS shell operations

- The code to enable FS shell indirectly hurts entire analytic flows in the Hadoop connectors by performing operations that are not inherent to the analytic flows

# Why does supporting FS shell affect analytic flows?

Persist distributed data set as an object

/data.txt/_temporary/0/_temporary/attempt_201702221313_0000_m_000000_0/part-0000
/data.txt/_temporary/0/_temporary/attempt_201702221313_0000_m_000001_1/part-0001
......
/data.txt/_temporary/0/_temporary/attempt_201702221313_0000_m_000008_8/part-0008

An opinionated object store connector for Spark can provide significant gains

# Stocator – the next-gen connector

- Advanced object store connector designed for object stores. Doesn't create temp files and folders for write  operations and still provides complete fault tolerance coverage including speculative mode.

- Doesn't use Hadoop modules and interacts with object store directly. This makes Stocator superior faster for write flow and generate many less REST calls

- Released under Apache License 2.0

- Implements Hadoop FileSystem interface.

- No need to modify Spark or Hadoop

https://github.com/SparkTC/stocator

- No need HDFS, comparing to other optimizations to the Hadoop connectors

# Stocator – the next-gen connector

**Stocator adapted for analytic flows**

|  | API | GET | HEAD | PUT | DELETE |
|---|---|:---:|:---:|:---:|:---:|
| **Stocator** | **S3** | **1** | **2** | **11** | **0** |
| Hadoop (s3a) | S3 | 158 | 361 | 26 | 16 |

# Compare performance of Stocator

- Stocator is much faster for write-intensive workloads

- Stocator as good for read-intensive workloads



* 40Gbps in accesser tier

** Comparing Stocator to S3a

# s3a connector is improving*

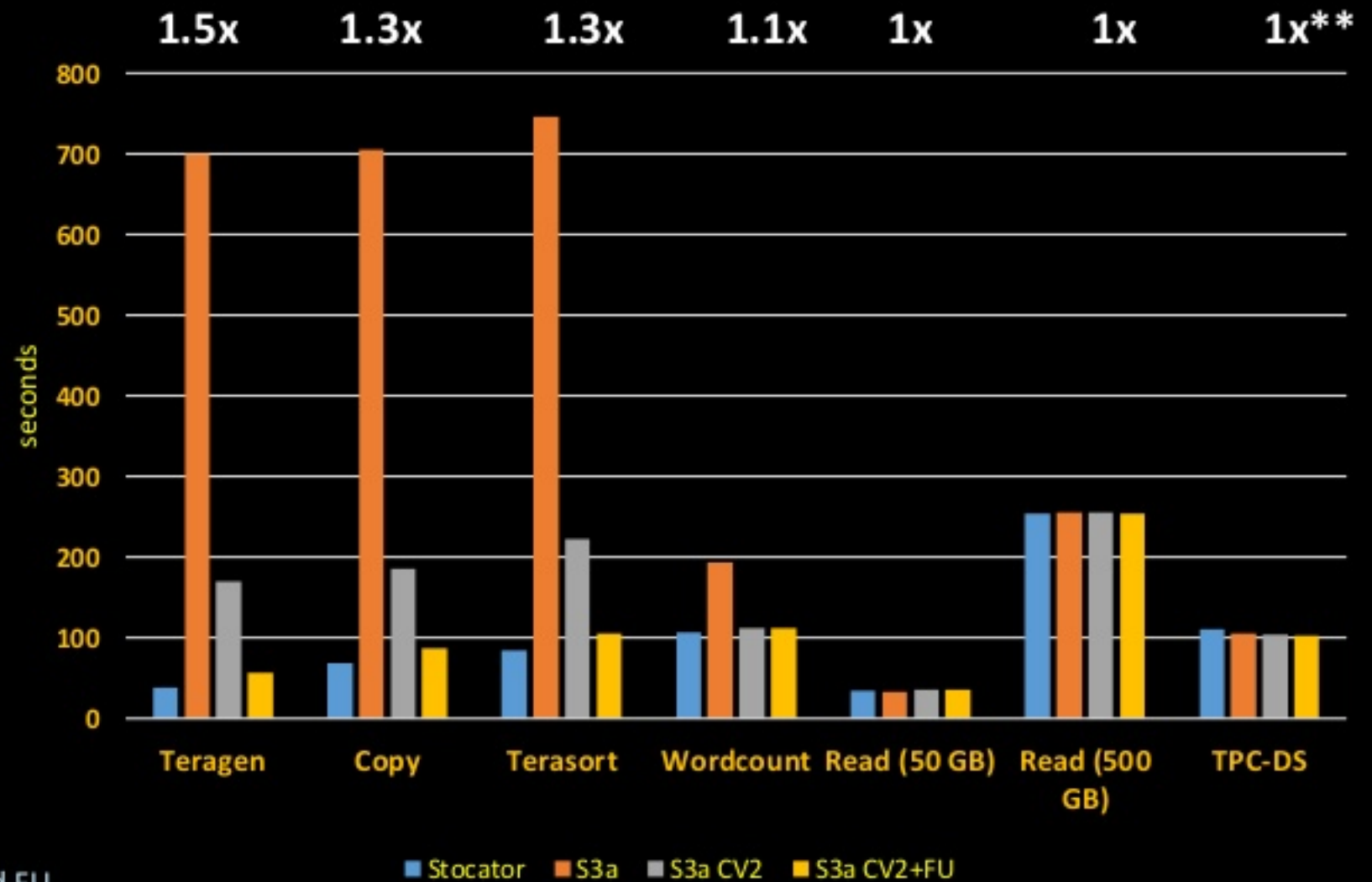- File Output Committer Algorithm 2 halves number of renames (CV2)

- Fast Upload introduces streaming on output (FU)

- Stocator still faster for write-intensive workloads and as good for read-intensive

** Comparing Stocator to S3a with CV2 and FU



**Chart labels:** 1.5x · 1.3x · 1.3x · 1.1x · 1x · 1x · 1x**

Y-axis: seconds (0, 100, 200, 300, 400, 500, 600, 700, 800)

X-axis categories: Teragen · Copy · Terasort · Wordcount · Read (50 GB) · Read (500 GB) · TPC-DS

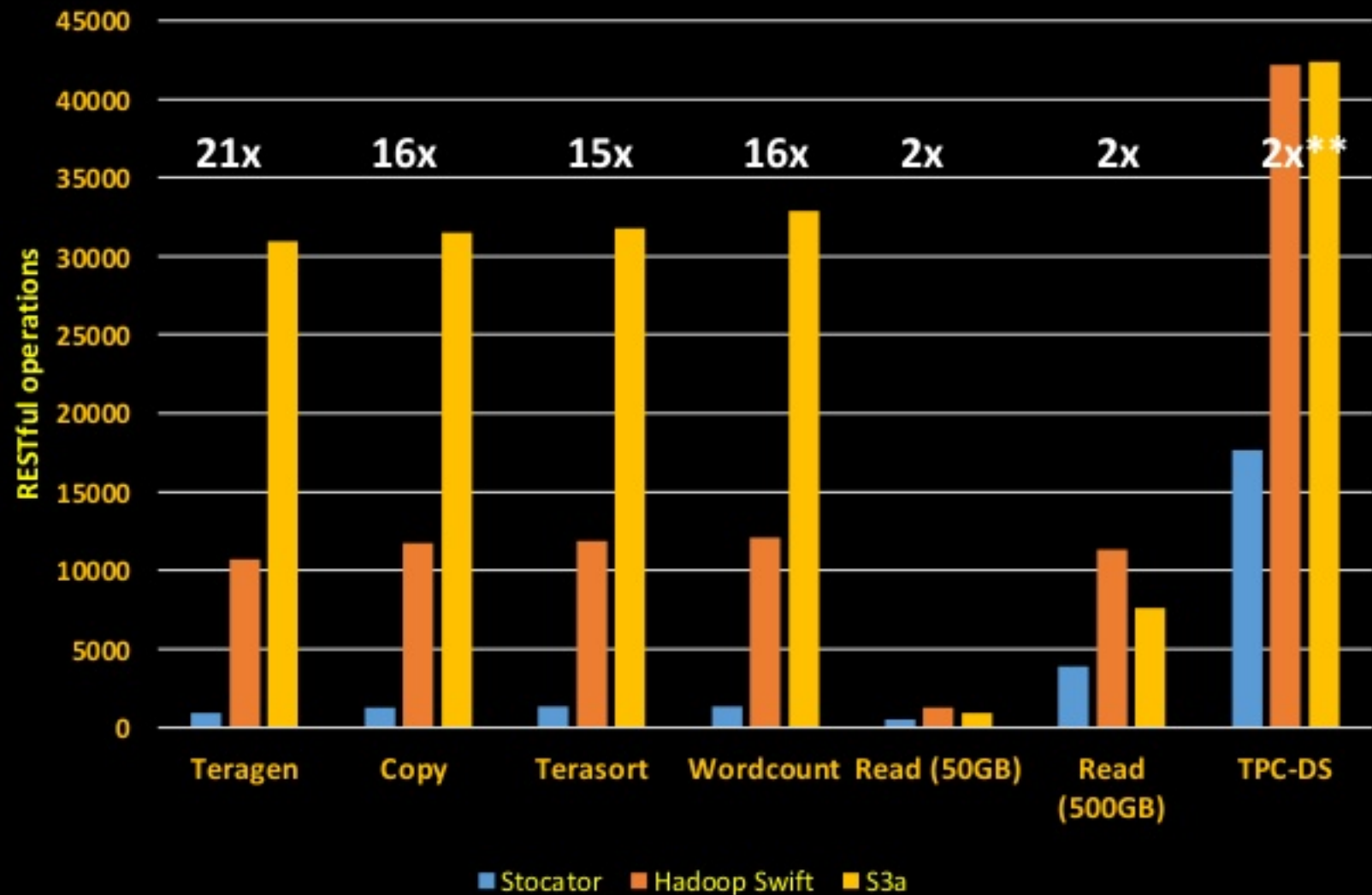Legend: ■ Stocator ■ S3a ■ S3a CV2 ■ S3a CV2+FU

# Compare number of REST operations*

- Stocator does many less REST operations
- Less operations means
  - Lower overhead
  - Lower cost



Chart: RESTful operations (y-axis, 0 to 45000) comparing Stocator, Hadoop Swift, and S3a across benchmarks: Teragen (21x), Copy (16x), Terasort (15x), Wordcount (16x), Read (50GB) (2x), Read (500GB) (2x), TPC-DS (2x**)

Legend: Stocator, Hadoop Swift, S3a

\* 40Gbps in accesser tier

\*\* Comparing Stocator to S3a with CV2 and FU

# Conclusions and Take-aways

- Spark@SETI is a real-world use case with demanding data and compute workload requirements

- Stocator is a critical component of the solution – enabling high performance data access from Object Store up to the Spark Cluster on IBM Cloud

- Stocator has been contributed to open source under Apache License 2.0

# Thank You.

graham.mackintosh@us.ibm.com

gilv@il.ibm.com