



Towards Use And Reuse Driven Big Data Management

Zhiwu Xie¹, Yinlin Chen¹, Julie Speer¹, Tyler Walters¹, Pablo A Tarazaga², and Mary Kasarda²

¹University Libraries and ²Department of Mechanical Engineering

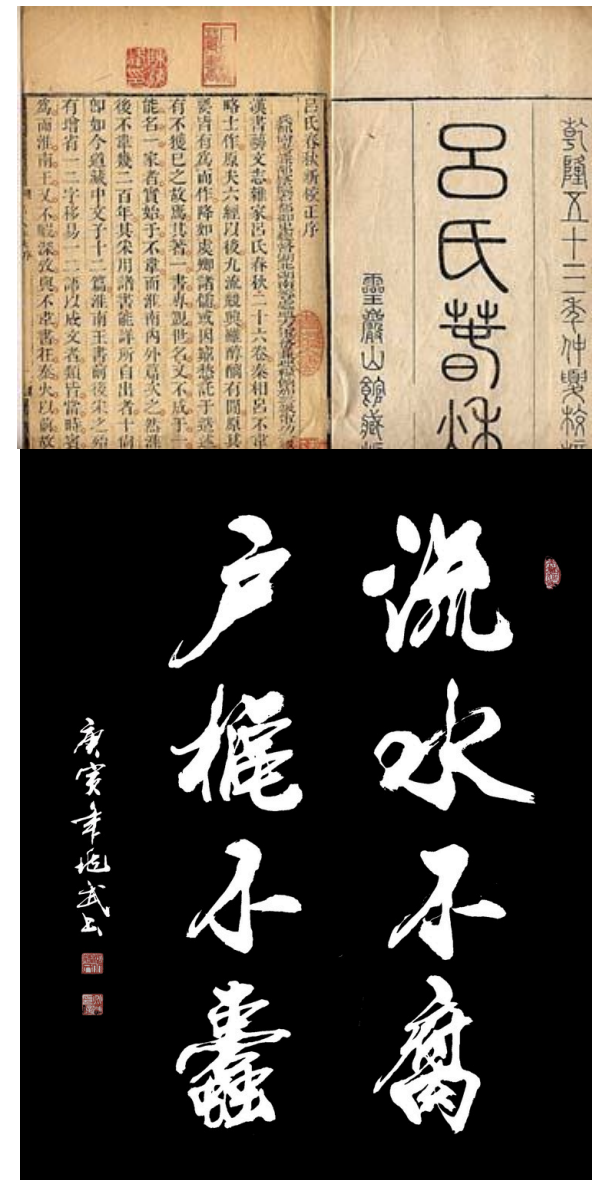
Virginia Polytechnic Institute and State University

Blacksburg, USA

June 23, 2015, JCDL 2015, Knoxville, TN

“...running water is never stale and a door-hinges never get worm-eaten...”

-- Lü's Annals, c. 239 BCE



Research Data Management

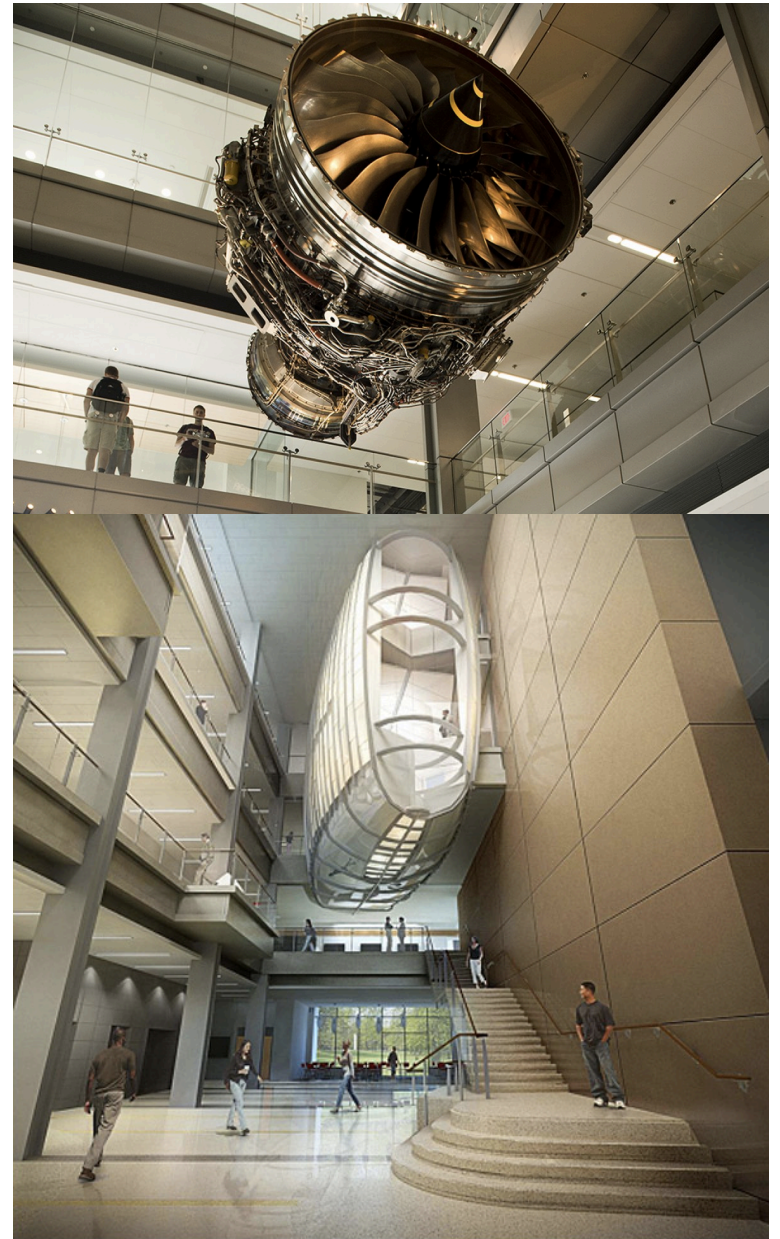
- What are the roles of the academic and research library?

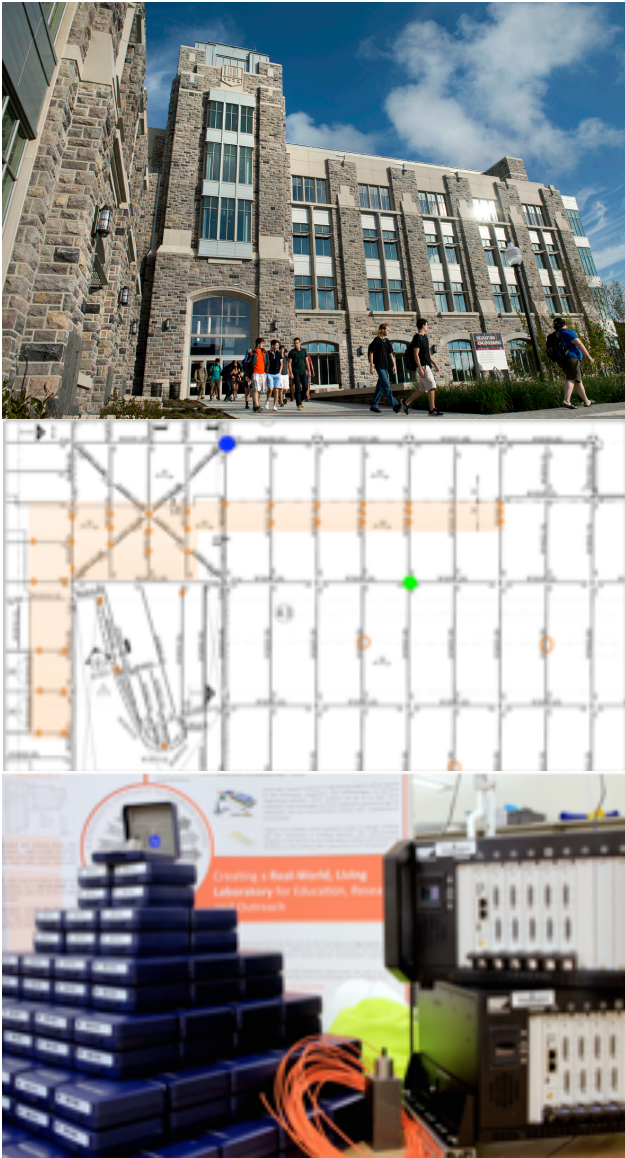


Chris 73 / Wikimedia Commons

Outline

- Motivation
- IT Infrastructure Choices
- System Architecture
- Prototype
- Validation
- Lessons Learned



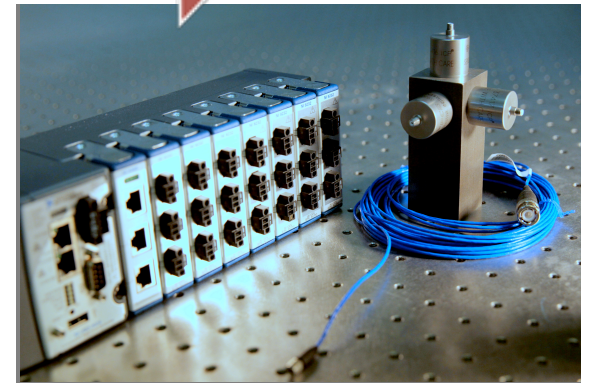
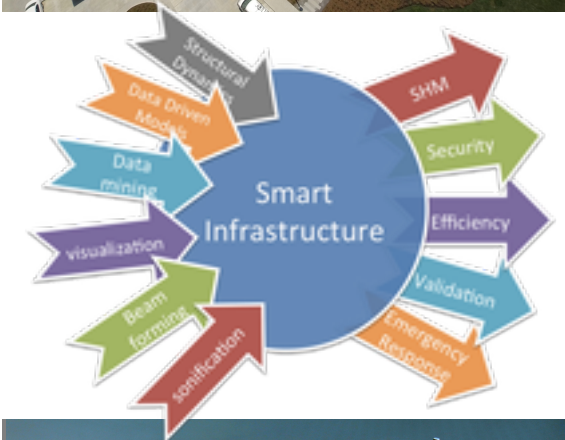


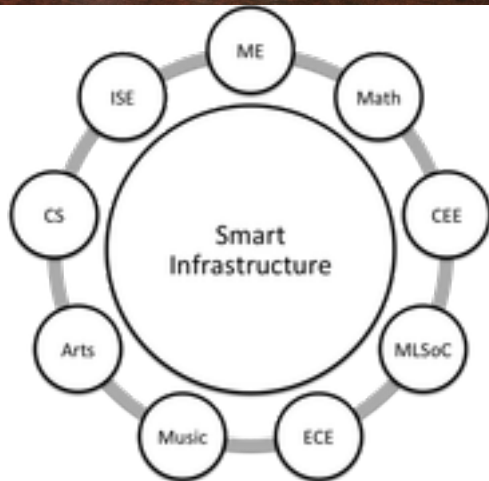
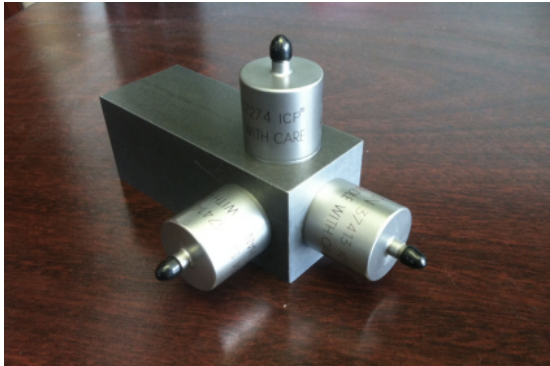
Motivation

- Goodwin Hall Living Lab
- A 160,000-sf new building wired with >240 different sensors
- Sensor mounts were directly welded to the structural steel during the building construction
- Sensors are strategically positioned and sufficiently sensitive to detect human movements
- Will be the most instrumented building for vibration

Motivation

- Designed as a multi-purpose living laboratory
- Opportunities for multi- and cross-disciplinary exploration and discovery
- > 40 researchers and educators in various disciplines and institutes expressed interests in using the data
- VT libraries is tasked with building the digital libraries to manage the data and support these activities.





Challenges

- A small team facing big data challenges
 - Cut across small and big science dichotomy
 - Cut across institutional boundaries
 - How to cultivate self-organizing, ad hoc user communities and support groups?

Challenges

- Multi- and cross- disciplinary exploration and discovery
 - Not a disciplinary repo
 - Must not dictate usage patterns
 - Difficult to index and search
 - Support on-demand analytics
- Infrastructure mismatch
 - Moving data is less favorable
 - Data repositories typically lack big data processing capabilities

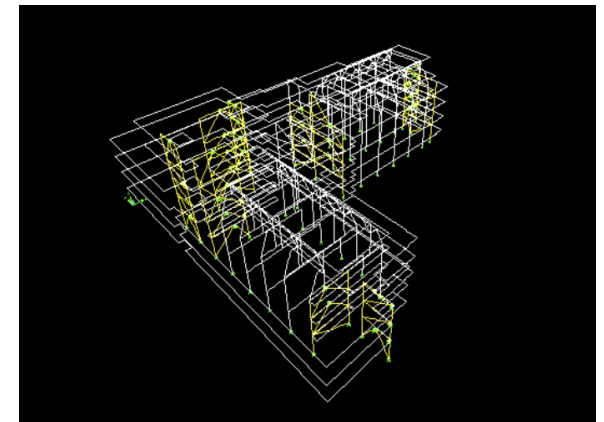
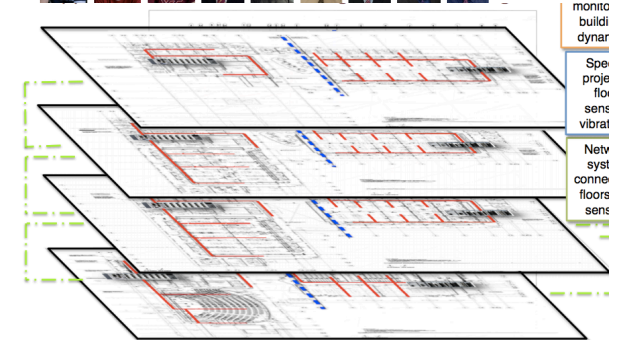
VT Smart Infrastructure Laboratory (VT SIL) established to create a synergistic environment for smart infrastructure education and research



Central organization for handling logistics, data and project management as well as access to instrumentation and raw data

VirginiaTech

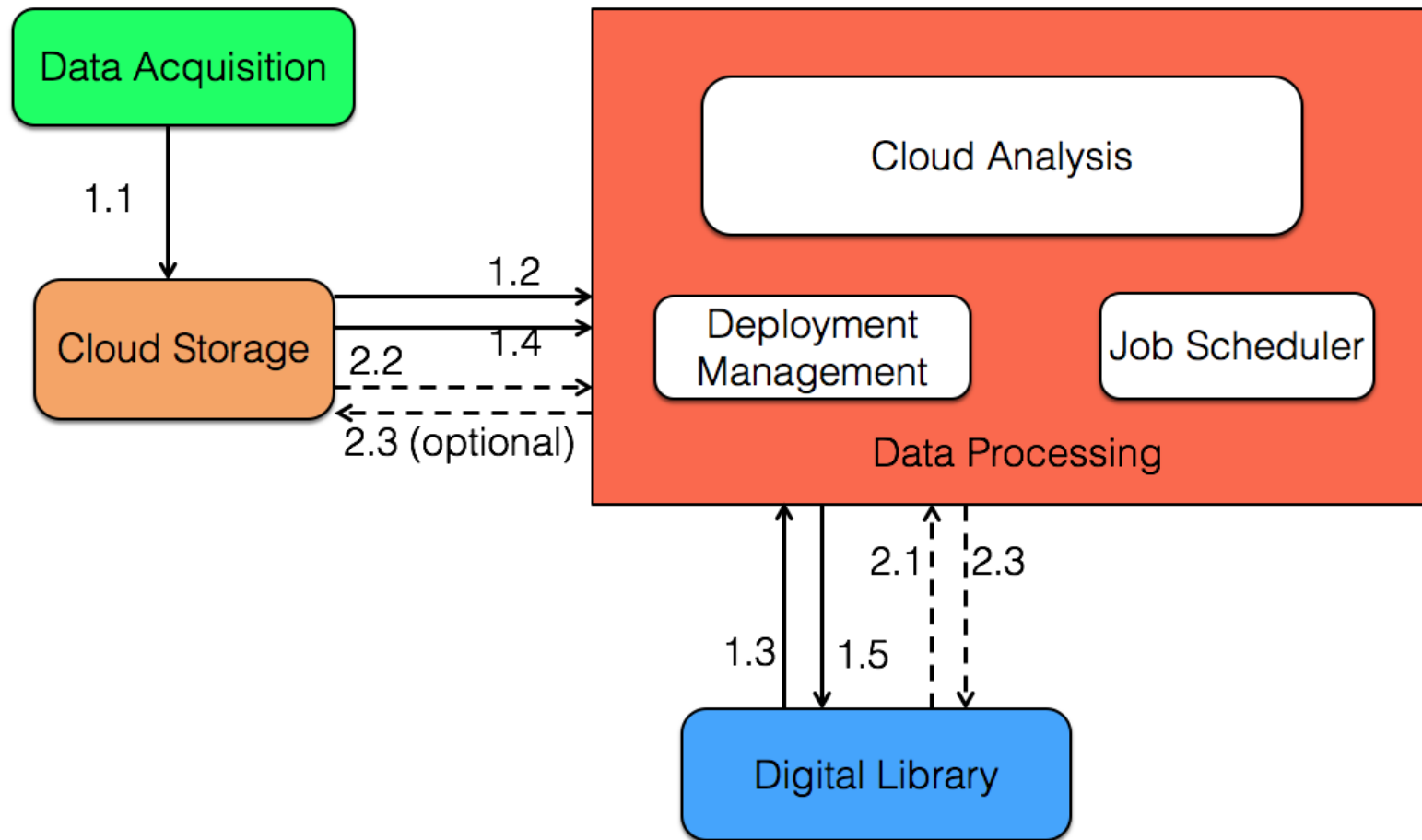
Smart Infrastructure Laboratory



Key IT Infrastructure Requirements

- Sufficient storage
- Highly scalable and elastic processing capabilities
- Co-located storage and compute nodes
 - A reference point: took 28 hours to move 4.3TB to a supercomputer center 5 miles away
- Cloud as a viable option

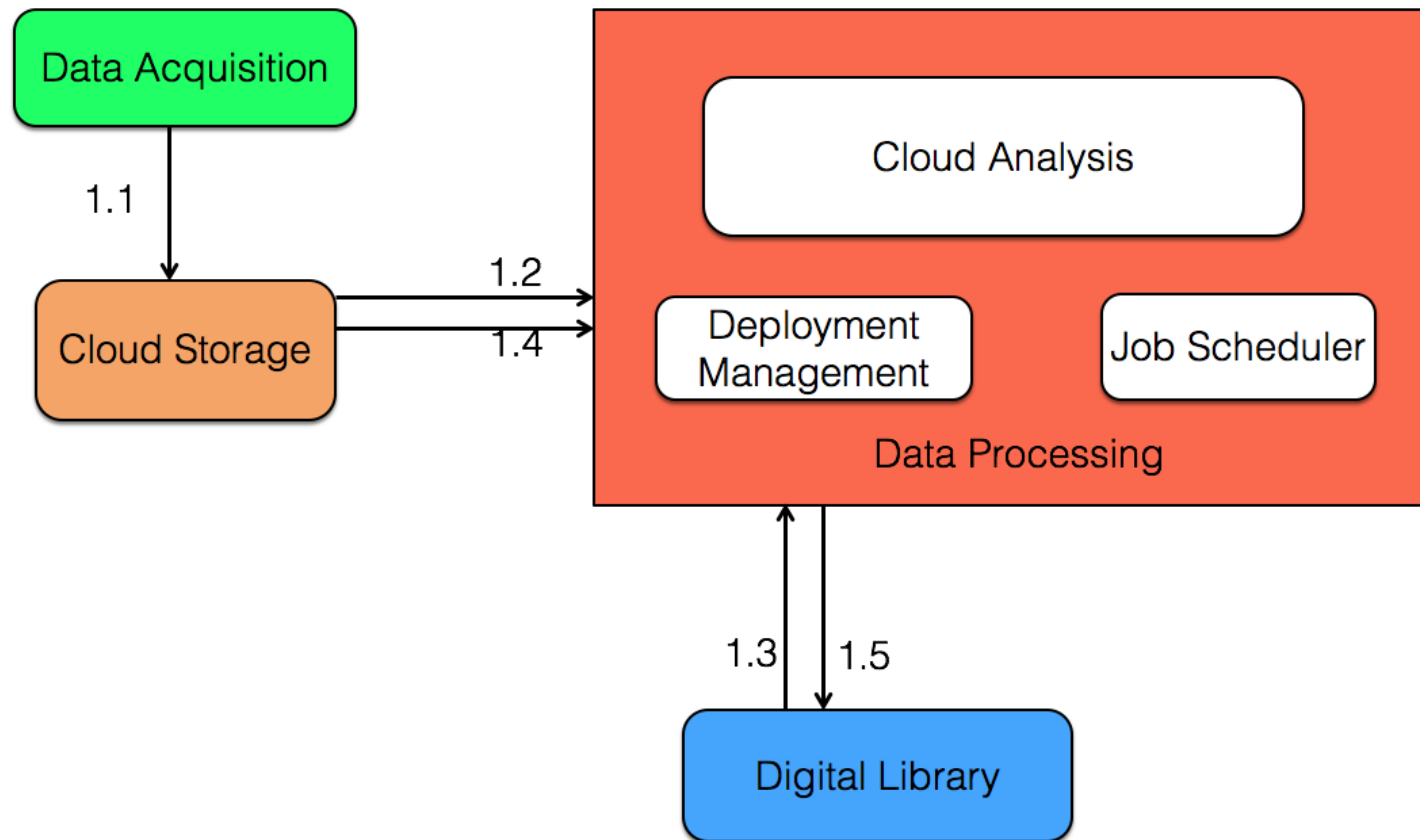
System Architecture



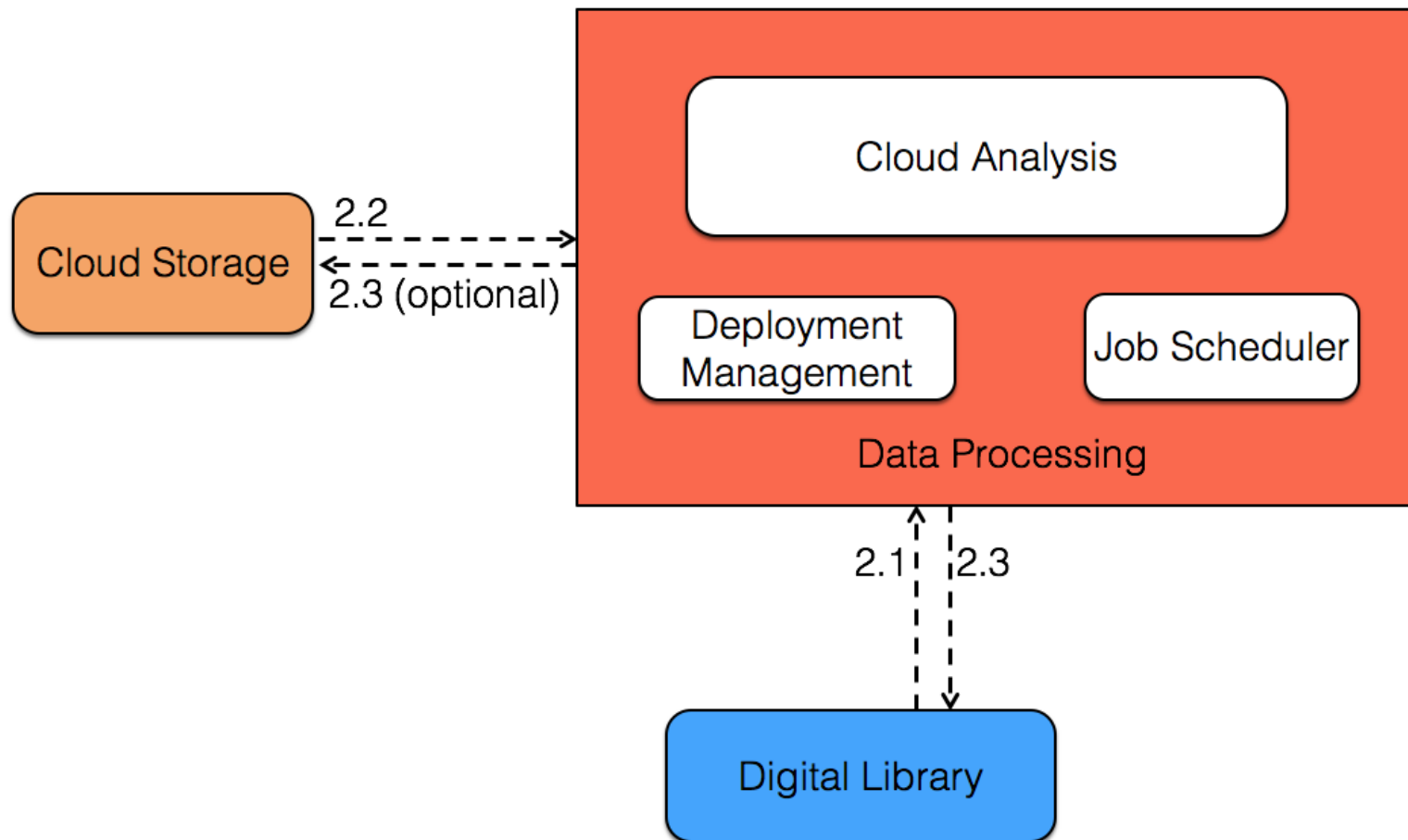
System Modules

- Data acquisition: Collect and process sensor data, store in a local storage server and upload to cloud storage module.
- Cloud storage: Raw sensor data and derivate data
- Data processing: Cloud analysis, deployment management, and job scheduler.
- Digital library: Bitstreams metadata, software codes, configuration files and user's info.

Data Ingestion Workflow



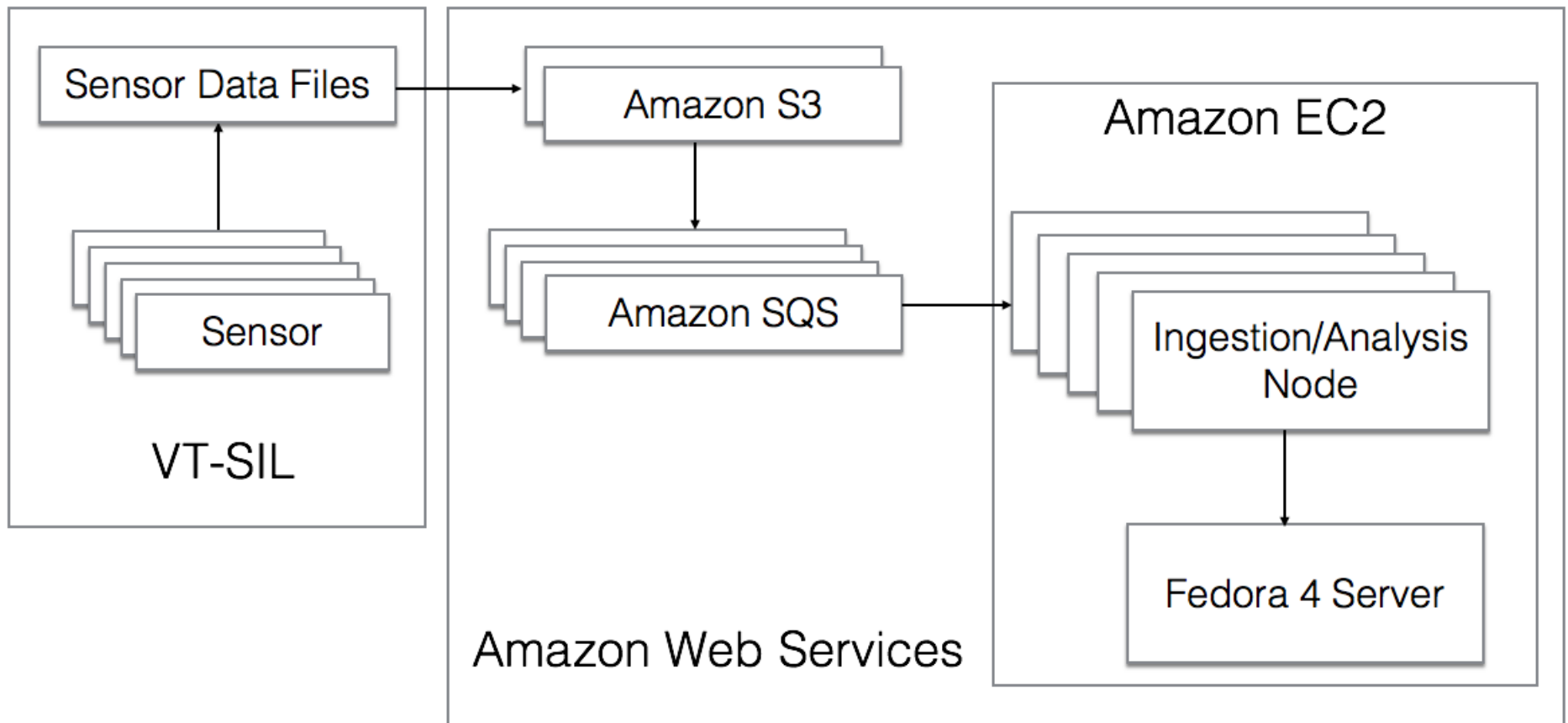
Data Analysis Workflow



Prototype

- Data acquisition and Cloud storage: Amazon S3
- Data processing:
 - Cloud analysis: Amazon EC2. (Workers - m3.large, Fedora 4 - r3.8xlarge)
 - Deployment management: Manually deploy using command line
- Job scheduler: Amazon SQS
- Digital Library: Fedora Repository

Prototype



Utilize AWS

- Sensor data upload to Amazon S3.
- Create queues in Amazon SQS

View/Delete Messages in seb2

Cancel 

View up to: messages Poll queue for: seconds *Polling for new messages once every 2 seconds.*

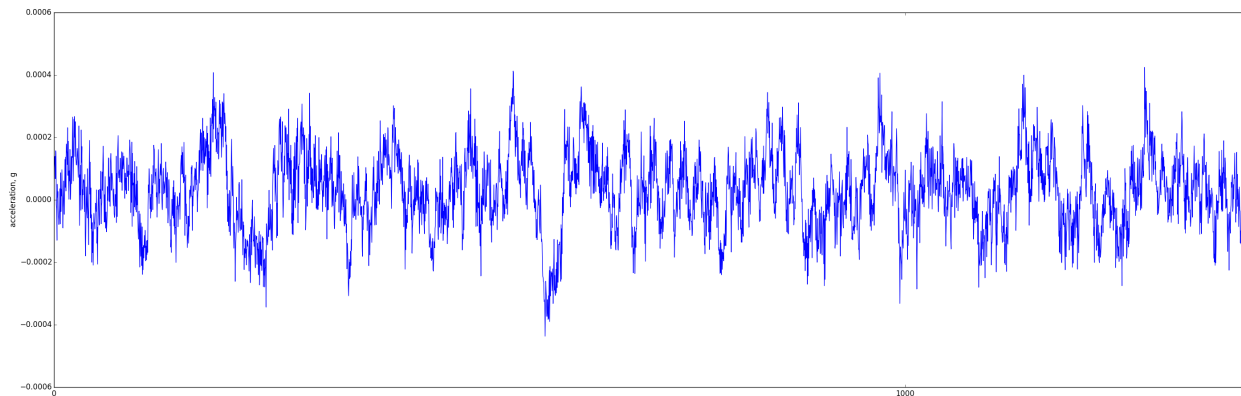
Start Polling for Messages

Stop Now

Delete	Body	Size	Sent	Receive Count
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "e84d3922-6817-545c-b870-77225321b79f", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:26 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "5b070303-99ed-5fff-a4eb-4d03a2db37d1", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:26 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "ed57339e-bc0c-5bf8-87ea-7d9e873b4bc8", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:27 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "f030059a-c451-5308-bd27-4b9aff66fcf0", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:27 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "5de32dd9-2f53-5b3f-9bac-63d9146ba2db", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:28 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "c6971288-bb87-5e31-a13e-5c8880bc00d2", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:31 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "ed8bd9f4-a901-5b3d-bf87-97f4923c7897", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:33 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "20849468-4981-5491-babe-da6872bf643b", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:34 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "b5f35133-fe41-5b4c-be0e-1f33c57e6185", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:35 GMT-04:00	1
<input type="checkbox"/>	{ "Type" : "Notification", "MessageId" : "d7e61912-5858-5e3d-9033-5712631885d1", "TopicArn" : More Details	1.7 KB	2015-06-21 00:09:36 GMT-04:00	1

Worker Operations

- Python codes that extract metadata from sensor data in HDF5 format.
- Mathematical operations:
 - Calculate max, min, mean, and median
 - Split sensor data
 - Draw charts from sensor data



Experiment Setup

- 24 hours of accelerometer measurements. (130GB)
- Three test cases:
 1. Simple data ingestion without data copying into the library
 2. Simple data analysis with additional metadata write to the library
 3. Complicated operations: Split data file into 6-second segments, create image, store image into library, metadata creation and etc.
- Setup 1, 2, 4, 8, 16 nodes to run these three test cases.

Test Cases

- Test case 1: Move data from S3 to workers, extract metadata, then create a container in Fedora 4, and write metadata info in each container.
- Test case 2: All operations in test case 1 + calculate the min, max, mean, and median value for each file and write these computing result into Fedora 4
- Test case 3: All operations in test case 2 + split data every 6-second, draw 172800 image files, and deposit to the Fedora 4

Fedora Record Example

data_channel_10_03-Sep-2014-1706.h5

[Home](#) / [ef](#) / [84](#) / [44](#) / [6a](#) / ef84446a-fcc6-4984-9e02-af9068a55f29

UUID

2e220349-062a-4aa0-8e91-1db9341aaa9a

Created at

2015-06-21T04:43:22.994Z by bypassAdmin

Last Modified at

2015-06-21T04:45:42.303Z by bypassAdmin

Children 150

1. <http://54.88.250.130:8080/fcrepo4/rest/ef/84/44/6a/ef84446a-fcc6-4984-9e02-af9068a55f29/41/04/a0/19/4104a019-0165-4a48-ab03-fefbef942256>

Properties

ns001: **channel**

12

ns001: **samplecount**

2304000

ns001: **samplerate**

2560

dc: **source**

https://s3.amazonaws.com/sebexdata/data_channel_10_03-Sep-2014-1706.h5

dc: **title**

data_channel_10_03-Sep-2014-1706.h5

fedora: **created**

2015-06-21T04:43:22.994Z

Experiment Results

Average time in second spent to copy a file from
Amazon S3 to Amazon EC2, both in the US East Region

	Number of Workers				
	1	2	4	8	16
1	0.2150	0.2525	0.2665	0.2045	0.2077
2	0.2186	0.2134	0.2675	0.2192	0.2073
3	0.2257	0.2412	0.2233	0.2086	0.2063

Experiment Results

Time in second spent to complete the test case

	Number of Workers				
	1	2	4	8	16
1	213.77	105.54	52.74	25.80	12.93
2	625.92	320.71	157.03	77.13	38.27
3	81564.42	40619.44	20284.53	10113.05	5059.34

Lessons Learned

- Opportunistic cost and cost sharing
- Organizational Changes
- Cultural changes

Questions

zhiwuxie@vt.edu @zxie