



# Lessons Learned Deploying a Big Data Geospatial System on AWS

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

## Moving WALDO from a Cluster to the Cloud

© 2016 Applied Research Associates, Inc.



NATIONAL SECURITY



ENERGY & ENVIRONMENT



INFRASTRUCTURE



HEALTH SOLUTIONS

Supported by the Intelligence Advanced Research Projects Activity (IARPA) via Air Force Research Laboratory, contract FA8650-12-C-7214. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes not withstanding any copyright annotation thereon. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, AFRL, or the U.S. Government.

ARA, Inc. (C) 2016 - Approved for FOS40NA2016

# ARA Quick Summary

## Founded 1979, Albuquerque, New Mexico

- Providing highly diversified R&D capabilities
- “Engineering and Science that Matters for Fun and Profit.”

## 1,100 employee owners in the U.S. and Canada

- Strong engineering & scientist focus (67% of company)
- Over 50% hold Advanced Degrees

## Over 30 locations in the U.S. & Canada

- Significant laboratory, manufacturing & testing facilities

**FY15 sales of ~\$210 million**

# ARA Open Positions in North Carolina

- Algorithm Development Systems Engineer
  - *(background in: system integration, inertial navigation, or computer vision - Junior/Mid-level opportunities)*
- Junior to Mid-Level Software Engineer
  - *(background in: Linux, cloud computing, Python & C/C++)*
- Mid-Level Simulation Tools Software Engineer
- Mid-Level Simulation Scientist/Engineer
- Simulation Scientist/Engineer
- Simulation Tools Software Engineer
- iOS or Android Software Developer
- Nuclear Effects Modeling and Simulation Scientist
- Structural Modeling Engineer
- Geospatial Intelligence Analyst
- Senior Imagery Analyst
- Imagery Analyst
- 3D Developer
- RF Communications System Modeler and Analyst

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

# WALDO uses FOSS

QGIS

GDAL/OGR

PostgreSQL/PostGIS

Python (Shapely, libgeographic, boto, fiona ...)

Luigi, Inc. (C) 2016 -- Approved for FOSS4GNA2016

OpenCV

Orfeo Toolbox

GRASS

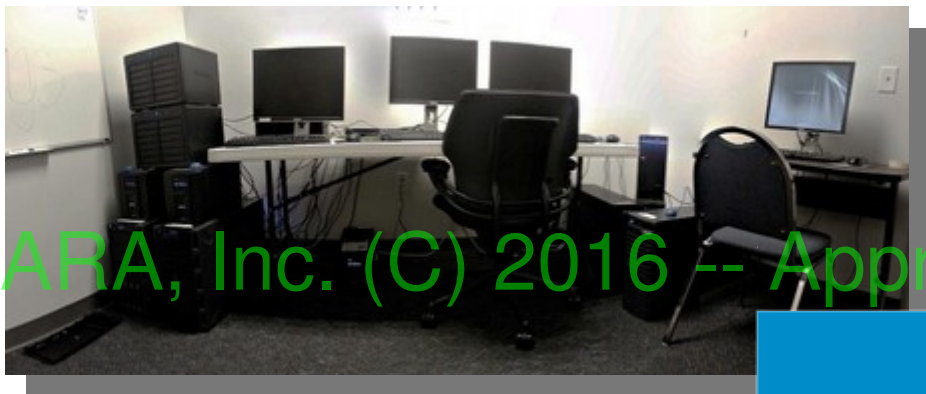
OpenSceneGraph

...and many more

# what did we learn...

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

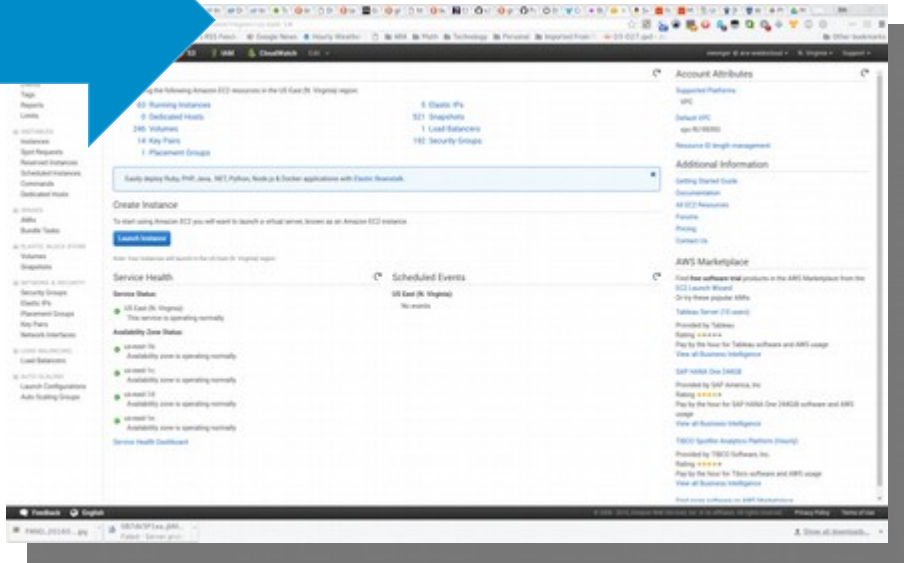
# ...moving WALDO from a cluster to the cloud?



ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016



- 5 nodes
- 84 cores / 168 threads
- 736 GB RAM
- 135 TB storage
- Windows Server 2008 R2
- MS HPC Pack



# aws works

**WALDO knowledgebase production**

ARA, Inc. (C) 2016 – Approved for FOSS4GNA2016  
**scales up by > 3x**

**full system updates delivered  
to customer in three images**

**WALDO is remotely accessible from anywhere**



**aws is easy**

**it really is**

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

**at least it is when you get started**

# aws is expensive

**total spend: \$95,000 (past 12 months)**

...and no one from amazon has ever called me

not even a free prime membership

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

**windows is twice as expensive as linux**

**aws doesn't come with idle resource monitoring**

surprise, surprise

**stopped instances cost too**

33 inst's \* 200GB \* 15 days ~= \$600

# aws is expensive

## use tags to track instances

start early

resource tags can be used to split the check

aws – why can't I enforce this via policy?

## plan to write some custom code

monitor spend, idle instances, etc.

lots of data to do this – just no easy button

## locate (and pay for) a service

# s3 is unreliable...

(and aws cli won't help)

aws s3 sync does not always sync completely

aws s3 sync does not check checksums  
s3cmd does, but it doesn't work on windows

s3fs is even more unpredictable  
what are these 0-byte files?

# s3 is unreliable...

(and aws cli won't help)

occam's solution: always sync thrice

use s3cmd on linux

use boto3 to roll your own sync command  
then put it on Github...please!

better yet—make a key db and roll your own sync

don't use s3fs – even if it tests well

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

# consider performance @ scale

all data (even s3) transits your NAT instance

s3 buckets have rate limits

300 gets/s, 100 puts/s

s3 buckets use key hashes for load-balancing

it's an object-store, not a file-store

ebs has its limits too, but scales easily

watch your iops

provisioned-io is very expensive (\$50/day for a db inst.)

# consider performance @ scale

**distribute data across many buckets**

use mgrs, utm, 1°x1°, etc.

**use a database to store object keys**

eliminate network latency from key lookup

**use ephemeral disks for /tmp volumes**

i'm looking at you, gdal\_translate!

**use pre-provisioned ebs snapshots for clusters**

much better scalability for duplicated data

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016

# security is well done (but complex and time-consuming)

**use a vpc to protect your data**

remember – all data transits the NAT – size it accordingly

vpc data is not encrypted – just segregated

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016  
**iam permissions are incredibly fine-grained...**

...and difficult to manage

**snapshot important instances frequently**

snapshots don't cost much

**use termination protection**

**any two aws employees can access your keys**



# a few remaining thoughts...

**users need to know linux...**

...and be comfortable with the command line

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016  
plan to teach them

**hadoop, spark, etc. is hard-**

-er than webinars make it look  
plan to spend time/money here

**avoid cloud lock-in**

we didn't, but maybe you can  
apache libcloud looks *really* interesting

# where will waldo go next?

**docker / mesos / kubernetes**

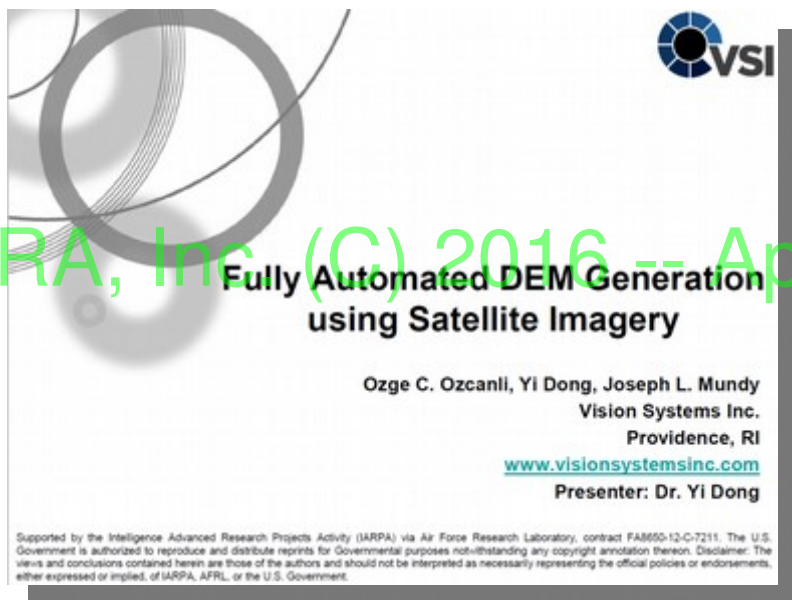
ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016  
**hadoop / spark / geowave**

**deep learning**

**cesium.js**

# Don't miss our other talks!

**Big Data Day – Room 301A  
Today @ 5:20 PM**

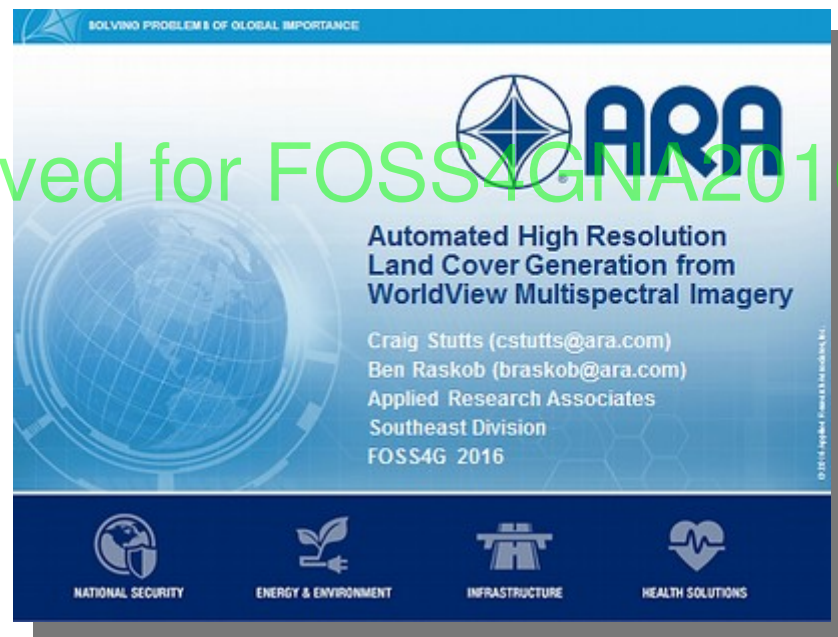


**VSIS**

## Fully Automated DEM Generation using Satellite Imagery

Ozge C. Ozcanli, Yi Dong, Joseph L. Mundy  
Vision Systems Inc.  
Providence, RI  
[www.visionsystemsinc.com](http://www.visionsystemsinc.com)  
Presenter: Dr. Yi Dong

Supported by the Intelligence Advanced Research Projects Activity (IARPA) via Air Force Research Laboratory, contract FAM650-12-C-7211. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, AFRL, or the U.S. Government.



SOLVING PROBLEMS OF GLOBAL IMPORTANCE

# ARA

## Automated High Resolution Land Cover Generation from WorldView Multispectral Imagery

Craig Stutts (cstutts@ara.com)  
Ben Raskob (braskob@ara.com)  
Applied Research Associates  
Southeast Division  
FOSS4G 2016

© 2016 Applied Research Associates, Inc.

NATIONAL SECURITY    ENERGY & ENVIRONMENT    INFRASTRUCTURE    HEALTH SOLUTIONS

**All Things Data – Room 302A  
Thursday @ 11:15**

Sign in and vote at [foss4gna.org](http://foss4gna.org)

# Evaluate the Sessions

ARA, Inc. (C) 2016 -- Approved for FOSS4GNA2016



-1 0 +1