# Apache Mahout

An Extendable Machine Learning Framework
for Flink and Spark (and others)

MAHOUT

# About Me

Trevor Grant

PMC Apache Mahout
& Apache Streams

Open Source Evangelist, IBM

MS Applied Math / MBA

@rawkintrevo

rawkintrevo@apache.org

http://rawkintrevo.org

MAHOUT

# About Mahout

History of Mahout 2008 - 2014

Lucene Subproject

TLP - May 2010

Feb 2014- v 0.9 (Last Map-Reduce)
  - Lots of Hadoop Vendors froze here

Popular as ML on MR.
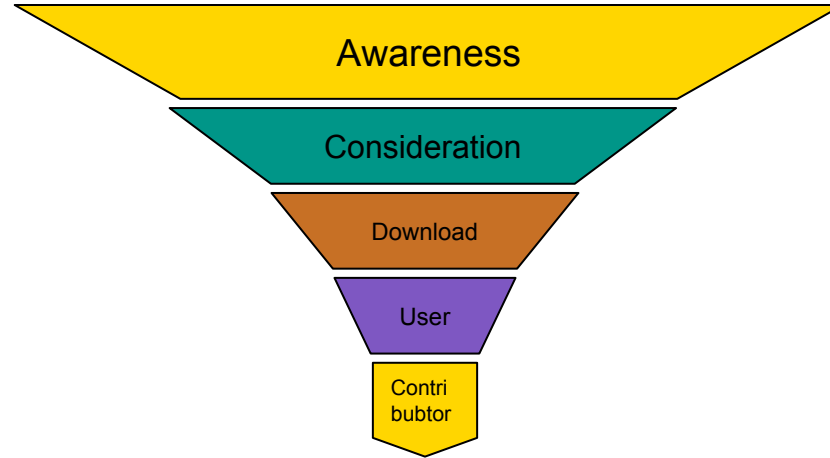
The rise of SparkML / MLLib

MAHOUT

# Rebranding

MAHOUT

# Marketing For Everyone

- You started doing open source so you could code whatever you want
- But you want people to use your open source product
- But your open source product doesn't have a marketing team
- And marketing teams aren't doing open source.
- …
- So you're the marketing team
- …
- Or you're writing a product no one is going to use. /shrug.
- **Exception:**  Your OSS product has a company, and the company pays the marketing team (and you probably have too many contributors).

MAHOUT

# "Marketing Funnel"

# Mahout Evangelism: Rebranding



MAHOUT

# Evangelism: Rebranding

**Problem:** High brand-recognition… for something we don't do anymore

**Opportunity:** Now relevant to many more areas- GPUs, Spark, Flink, etc.

**Solution:** Talks and Blog Posts

**Corollary:** New awareness, and changing perceptions

MAHOUT

# Mahout Evangelism: Talks (Spring '17)

**APACHE MAHOUT'S NEW RECOMMENDER ALGORITHM AND USING GPUS TO SPEED MODEL CREATION** *Pat Ferrel, Andy Palumbo*. GPU Technology Conference. Silicon Valley, CA- May 11, 2017

**EXTENDING MAHOUT-SAMSARA LINEAR ALGEBRA DSL TO SUPPORT GPU CLUSTERS** *Suneel Marthi, Trevor Grant*. GPU Technology Conference. Silicon Valley, CA- May 11, 2017

**Apache Mahout: An Extendable Machine Learning Framework for Spark and Flink** *Trevor Grant*. Apache Big Data. Miami, FL- May 16, 2017

**An Apache Based Intelligent IoT Stack for Transportation** *Trevor Grant, Joe Olsen*. ApacheCon IoT. Miami, FL- May 18, 2017

(+2 at ApacheCon/Apache Big Data but last minute speaker had conflict)

**Apache Mahout: Distributed Matrix Math for Machine Learning** *Andrew Musselman*. MLConf. Seattle, WA- May 19, 2017

**Weekend Project: Real World AirBnB Data Science and Pricing Bot** *Trevor Grant, Andrew Weiner*. Berlin Buzzwords. Berlin, DE- June 13, 2017

**Introduction to Online Machine Learning Algorithms** *Trevor Grant*. Dataworks Summit. San Jose, CA- June 15, 2017

**MAHOUT**

# Evangelism: Talks (Fall '17/ Spring '18)

**Matrix Math at Scale With Apache Mahout and Apache Spark.** *Andrew Musselman*. ODSC. Boston, MA. May 1st, 2018.

**The Magnificent Modular Mahout- An extensible library for distributed math and HPC** *Trevor Grant*. FOSDEM- High Performance Computing. Brussels, BE- February 4th, 2018.

**Do I Know You? Realtime Facial Recognition with an Apache Stack** *Trevor Grant*. Flink Forward, Berlin. Berlin, DE- September 13, 2017

MAHOUT

# Mahout Evangelism: Blog Posts

**Multi-domain predictive AI -** June 26, 2018 *Pat Ferrel PMC*
https://developer.ibm.com/dwblog/2017/mahout-spark-correlated-cross-occurences/


**Introducing Precanned Algorithms in Apache Mahout** - May 2, 2017. *Trevor Grant*
https://rawkintrevo.org/2017/05/02/introducing-pre-canned-algorithms-apache-mahout/

**Getting Started with Apache Mahout** - April 25, 2017. *Trevor Grant*
https://datascience.ibm.com/blog/getting-started-with-apache-mahout-2/

**Correlated Cross-Occurrence (CCO): How to make data behave as we want -** December 1, 2016. *Pat Ferrel*
http://actionml.com/blog/cco

MAHOUT

# Starting to Gain Traction: Non PMC Blog Posts

**Top Skills Data Scientists Need to Learn in 2018** - Feb 3, 2018.

https://insidebigdata.com/2018/02/03/top-skills-data-scientists-need-learn-2018/

**10 Best Machine Learning Software 2018 -** Feb 24, 2018. (Map-Reduce)

https://data-flair.training/blogs/machine-learning-software/

**Encyclopedia of Big Data Technologies -** Jan 29, 2018

https://link.springer.com/referenceworkentry/10.1007/978-3-319-63962-8_144-1

**Machine Learning, predictive AI and IoT… Oh My! -** May 23, 2017 *Lisa Seacat*

https://developer.ibm.com/dwblog/2017/mahout-spark-correlated-cross-occurences/

**Forgot About Mahout- It's Back and Worth Your Attention.** May 18, 2017. *Andrew C. Oliver*

https://www.infoworld.com/article/3197429/machine-learning/forgot-about-mahout-its-back-and-worth-your-attention.html

MAHOUT

# Non-Science/Data Based Graph

# Monitoring Progress: Analytics

# Evangelism

Awareness

Consideration

Download

User

Contribubtor

MAHOUT

# (Old) Website

# New Image

**Problem:** Old Website

**Opportunity:** Websites are easy (technically- design on the other hand...)

**Solution:**  Reboot website with Jekyll Bootstrap (and less emphasis on Map-Reduce)

**Corollary:** Much easier for committers and contributors to update website, add tutorials, etc. Encourages (requires) new features to be submitted with good docs.

MAHOUT

# New Website

# New Logo



MAHOUT

# Relevant Features

MAHOUT

# Killer New Features

GPU Integration

Algo Framework / Precanned Algos

Zeppelin Integration

Mahout-Samsara (Mathematically Expressive Scala DSL)

Etc.

MAHOUT

# New Features

Awareness

Consideration

Download

User

Contri
bubtor

MAHOUT

# Apache Zeppelin integration

MAHOUT

# Huge win

Scala has weak visualization.

R and Python have great visualization!

Zeppelin allows user to hand off variables between interpreters in notebook.

Do work in Mahout (Scala) - Plot in R/Python!

Mahout Gaussian DRM plotted in R



Took 0 sec. Last updated by anonymous at September 28 2016, 1:52:10 PM. (outdated)

MAHOUT

# Huge win- kinda wonky install.

A script exists for assisting with install… needs update for 0.13.0
https://issues.apache.org/jira/browse/ZEPPELIN-2417

Major Issue- Zeppelin out of the box supports Scala 2.11/Spark 2.x

Mahout 0.13.0 will build, but no binaries exist.

User must either build Mahout (for Spark 2.1) or Build Zeppelin (for Spark 1.6)

Fix coming soon in 0.13.1 (profiles and binaries for Spark 2.1/Scala 2.11)
http://mahout.apache.org/docs/0.13.1-SNAPSHOT/tutorials/misc/mahout-in-zeppelin/

MAHOUT

# Or so we thought...

Originally we planned on releasing 0.13.1 weeks after 0.13.0

There were issues with `scopt` being version locked on Scala 2.10 (required for CLI drivers).

Also, this caused issues in `sbt` builds (obvious solution, no one should use sbt?)

Eventually led to...

MAHOUT

# Which Led to

Disagreements on mailing list / Slack.

People blaming people for things.

Other issues one has when fighting with family around the holidays.

MAHOUT

# And Finally...

Massive POM ~~cleaning~~ gutting.

Project structure refactor (MR now in Community, Math components combined).

In progress now.



MAHOUT

# Sometimes Code is like an Appendix

Once useful, but so long ago- no can remember why.

Appendix

**MAHOUT**

# Other Project Integration

# Mathematically Expressive Scala DSL

MAHOUT

# Big Math

**Problem:** May not know Scala/Apache Spark/Apache Flink

**Opportunity:** Do know R

**Solution:** Create an abstracted language- mathematically expressive Scala DSL

MAHOUT

# Big Math

```
implicit val sdc: org.apache.mahout.sparkbindings.SparkDistributedContext = sc2sdc(sc)

val A = drmWrap(rddA)
val B = drmWrap(rddB)

val C = A.t %*% A + A %*% B.t
```

`C` is a `RDD[(Int, org.apache.mahout.Vector)]`

Also have truly distributed matrix decompositions.

MAHOUT

# "Possibly" Coming soon: Mahout for Tensorflow

Problem- Tensorflow is very very ugly.

We can use Tensorflow 2d Matrices with Mahout Scala DSL.

**MAHOUT**

# Big Math



Awareness

Consideration

Download

User

Contri bubtor

MAHOUT

# Engine Neutrality

MAHOUT

# Engine Neutrality

**Problem:** Distributed engines come and go (we learned this first)

**Opportunity:** We learned this first!

**Solution:** Create Engine Neutral Libraries which can bind to new engines

**Corollary:** Implement algorithm once- run it anywhere*

MAHOUT

# We don't actually have Beam bindings (yet), this is just for lulz



HEARD YOU LIKE ABSTRACTING DISTRIBUTED ENGINES

SO I WROTE BEAM BINDINGS FOR MAHOUT NOW YOUR ABSTRACTION IS ABSTRACTED

MAHOUT

# Engine Neutrality



Awareness

Consideration

Download

User

Contri
bubtor

MAHOUT

# Algorithm
# Framework

MAHOUT

# Attracting Contributors

**Problem:** Very small pool of qualified contributors

**Opportunity:** Mathematically expressive scala makes it easy to write and review "math part"

**Solution:** Create templates and tutorials showing users how to add algorithm

MAHOUT

# Drop your algorithms in easily

Algorithms are easy to compose in Mahout (as easy* as R, often can use R implementation for guidance)

Scala classes / package layout may still be overwhelming for our target users

MAHOUT

# Algorithm Template

```scala
3   class Foo[K] extends RegressorFitter[K] {
4
5     def fit(drmX   : DrmLike[K],
6             drmTarget: DrmLike[K],
7             hyperparameters: (Symbol, Any)*): FooModel[K] ={
8       /**
9        * Normally one would have a lot more code here.
10       */
11
12      var model = new FooModel[K]
13      model.summary = "This model has been fit, I would tell you more interesting things- if t
14      model
15    }
16  }
17
18  class FooModel[K] extends RegressorModel[K] {
19
20    def predict(drmPredictors: DrmLike[K]): DrmLike[K] = {
21      drmPredictors.mapBlock(1) {
22        case (keys, block: Matrix) => {
23          var outputBlock = new DenseMatrix(block.nrow, 1)
24          keys -> (outputBlock += 1.0)
25        }
26      }
27    }
28  }
```
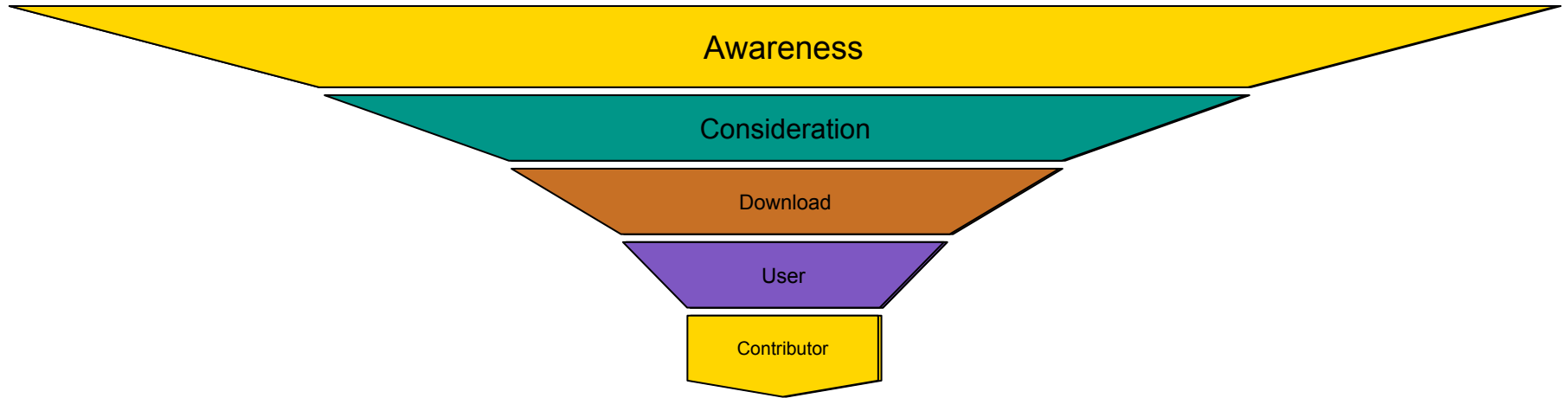
MAHOUT

# Encourage using dev

A lot of hand holding with first time contributors- encourage them to pay it forward.
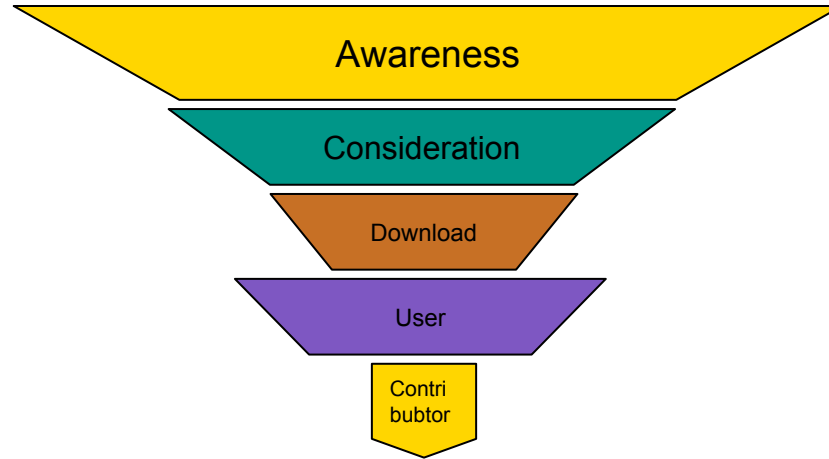
(Mailing lists can be hard.  Guy who did the website, DID the website, took him months to figure out how to subscribe to dev).

**MAHOUT**

# Engine Neutrality



Awareness

Consideration

Download

User

Contributor

MAHOUT

# How this actually progressed...

MAHOUT

# First The Math and Engine Neutrality

Awareness

Consideration

Download

User

Contri
bubtor

2014-2016

MAHOUT

# Then we started our evangelism



2014- Ongoing, big spike 2016/7

Awareness

Consideration

Download

User

Contri
bubtor

MAHOUT

# Then the GPU work and algorithm framework

Q3 2016- ongoing

Awareness

Consideration

Download

User

Contributor

MAHOUT

# New website...

Q2-Q4 2017



Awareness

Consideration

Download

User

Contributor

MAHOUT

# Hooray!



Awareness

Consideration

Download

User

Contributor

MAHOUT

# What's Next

MAHOUT

# Building Algorithms Framework

Framework in place which encourages users to contribute algorithms (already paying off)

Seeking to grow "pre-canned" algorithms collection between now and v 0.14.0

Eventually a "CRAN" like repository of algorithms for Mahout.

MAHOUT

# More work on GPUs

We consider our GPU support a HUGE differentiator among ML packages native to distributed engines (MLLib, FlinkML, etc).

Still opportunities for optimization-

Recent benchmark on unoptimized (still technically PR) CUDA bindings show "significant" speed up on sparse Matrix multiplication

MAHOUT

# 300%

Speed up On Sparse-Sparse Matrix Multiplication on AWS GPU enabled Spark Cluster

Kind of a big deal.

MAHOUT

# More Engine Bindings (Tensorflow)

Create template engine bindings- even if not optimized.

Tutorials on writing new engine bindings.

We feel this is also a huge differentiator.

Possibly "Community" supported engine bindings, not officially supported- but in the trunk, attract "drive-by" contributions from other projects.

MAHOUT

# Getting over version lock...

Some issue had us version locked on Scala 2.10 (Spark 1.6).

Huge refactor of the POMs

"Map-Reduce" is "Community" now.

- We're "soon" (maybe) going to be taking MR PR's again (have not accepted them for 3.5 years)

MAHOUT

# Conclusion

MAHOUT

# Awakening the Giant

Mahout has quietly undergone huge transformation from Map Reduce / Java based Machine Learning to Mathematically Expressive Scala / Engine Neutral / GPU Accelerated
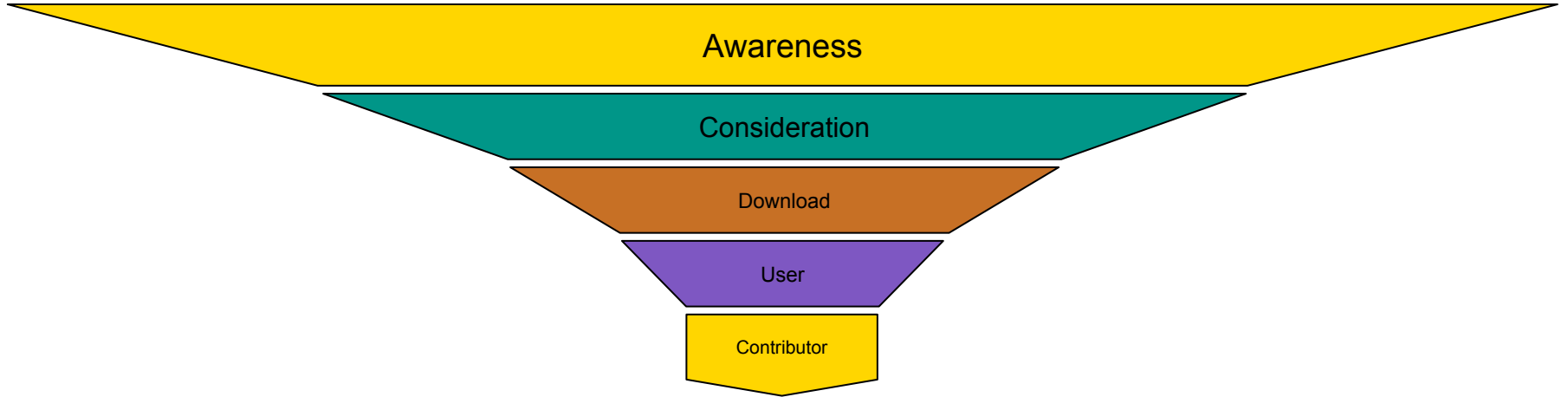
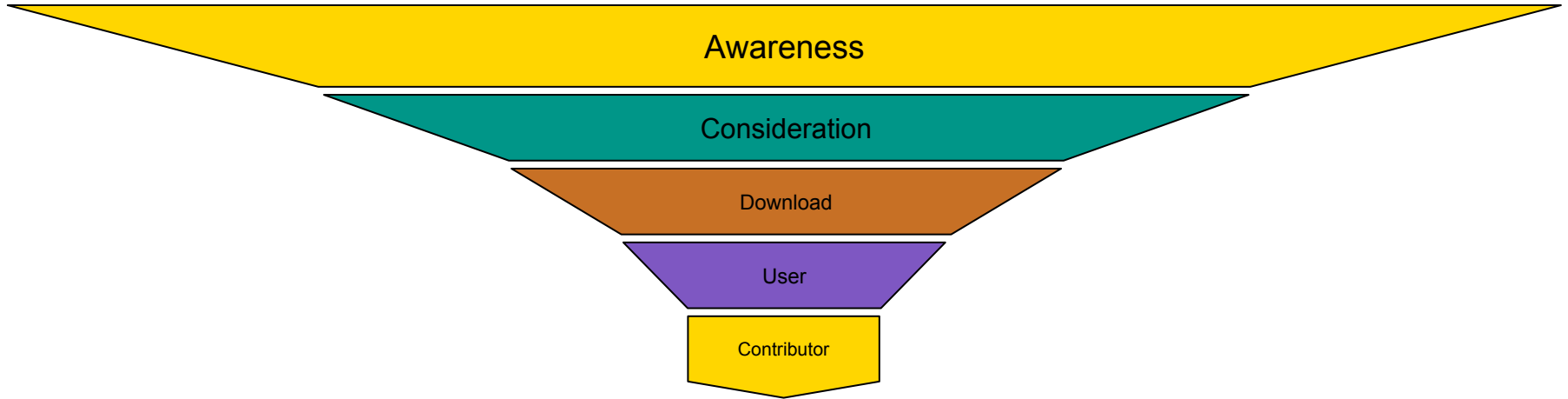Need to let everyone know- I mean you- go tell your friends and tweet and write a blog.

MAHOUT

# "I want you so badly" - The Beatles

# Remember the funnel for your project.

# People can only download if they are aware.

# Questions?

MAHOUT