# FIELD NOTES

## FROM EXPEDITIONS IN THE CLOUD

---

### DR. MATT WOOD
### GM, PRODUCT STRATEGY

# GOOD MORNING

---

MATTHEW@AMAZON.COM

@MZA

# THANK YOU

# TECHNOLOGY ADOPTION

# "WHAT IS CLOUD COMPUTING?"

# "WHAT IS BIG DATA?"

what is big data

Upload

**Content**    Users

Uploaded Anytime ▾    All File Types ▾    English ▾

Page 1 of 431,710 results for **what is big data**

### BIG Data
### What is it?
BROUGHT TO YOU BY THE BESTSELLING AUTHOR OF...

Bernard Marr
**What is** Big Data?

25 slides, 280 likes

### What is Big Data?

David Wellman
**What is** big data?

54 slides, 54 likes

edureka!
What is *Big Data* **and** Why learn *Hadoop*

View Hadoop Courses at:

Slide 1    Twitter @edurekaIN, Facebook /edurekaIN, use #AskEdureka for Questions

Edureka!
**What is** Big Data and Why Learn Hadoop

20 slides, 3 likes

### WHAT IS
### BIG DATA?

What is Big Data?

# BOXED IN

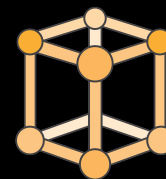**Canonical data stores**

**Streaming data**

**Task clusters**

# SPARK WITH AMAZON EMR

Production workloads on AWS

**yelp**

Machine learning &
ad targeting

**HEARST**

Web analytics

**The Washington Post**

Ad targeting &
recommendations

**QUIXEY**

App search

**CROWDSTRIKE**

Security event
streaming

**gumgum**

Revenue forecasting

**RADIUS**

Predictive marketing

**krux**

Personalization

S3

Clicks, videos,
interactives

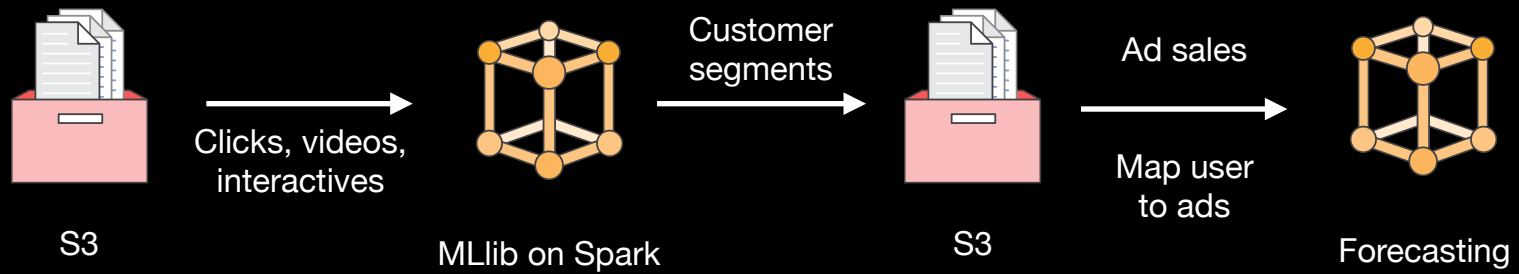MLlib on Spark

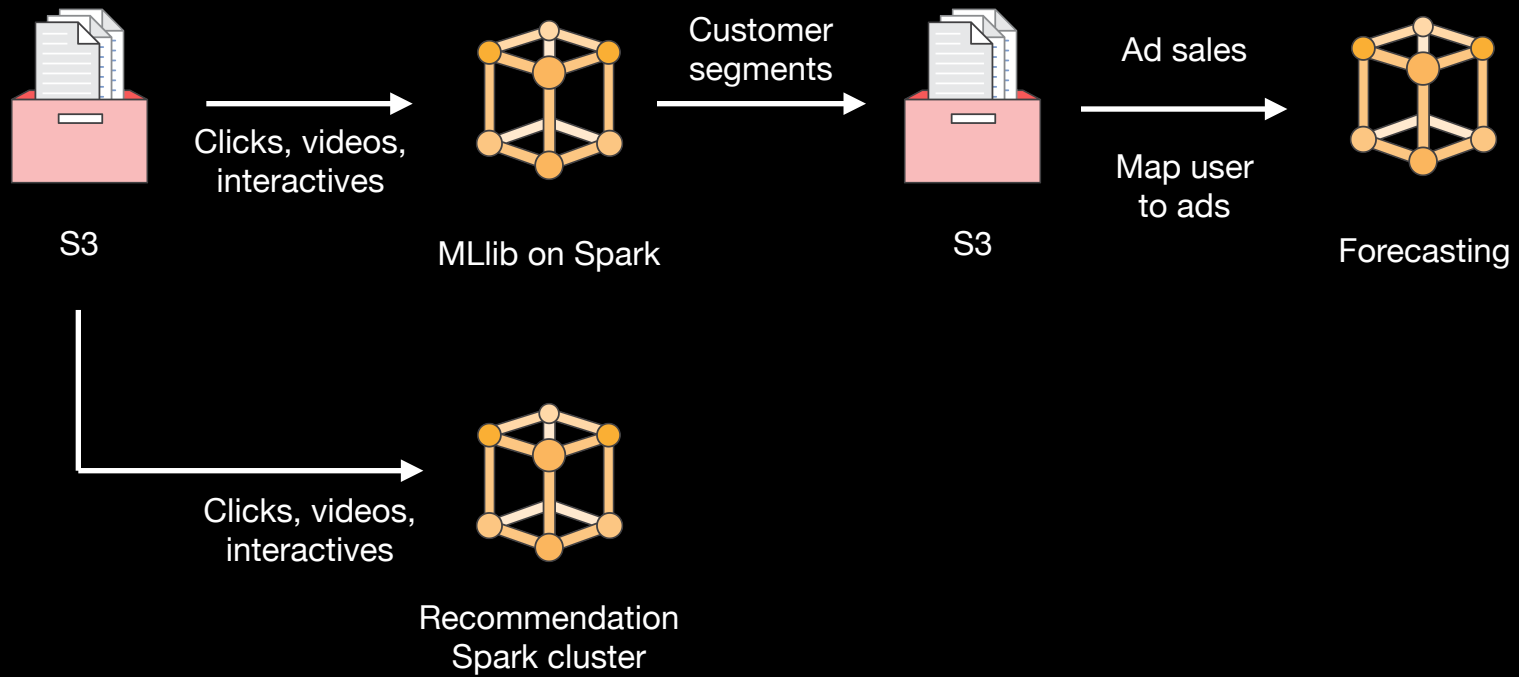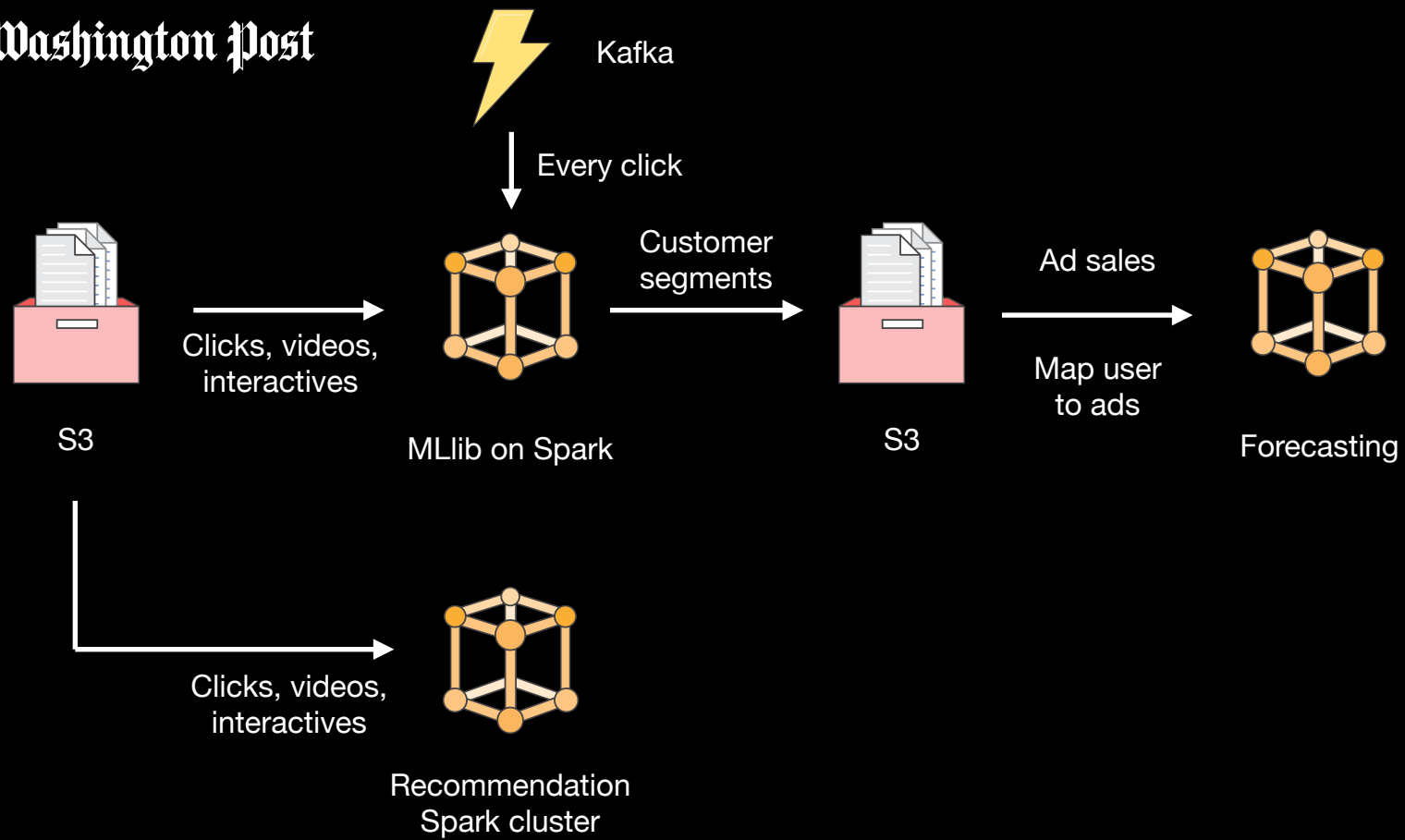The Washington Post

S3 → Clicks, videos, interactives → MLlib on Spark → Customer segments → S3 → Ad sales / Map user to ads → Forecasting

S3 → Clicks, videos, interactives → Recommendation Spark cluster

# HEARST

Node.js app

Elastic Beanstalk

**HEARST**

Node.js app
Elastic Beanstalk

— Click stream → Kinesis

— 5 minute windows / 100MB → Spark Streaming / EMR cluster

— JSON and CSV → S3

ElasticSearch ← Redshift ← EMR

HEARST

Node.js app
Elastic Beanstalk

Click stream

Kinesis

5 minute windows

100MB

Spark Streaming
EMR cluster

JSON and CSV

S3

ElasticSearch

Redshift

EMR

# SPARK ON EMR

Provision and scale managed Spark clusters, as first class citizens

# RAPID PROVISIONING OF ELASTIC CLUSTERS

Provision new clusters in minutes

High memory, high CPU, high IO instances

Add or remove capacity on running clusters

Access cluster instances directly

Clusters run within a VPC

# DIRECT ACCESS TO DATA ON S3

Access objects directly on Amazon S3

Multiple clusters can access canonical data in S3

No need to copy or manage the data on the cluster

Server-side and client-side encryption with customer controlled keys

Mix S3 and HDFS on a cluster

# INTEGRATION WITH THE SPOT MARKET

Bid on under utilized capacity on EC2

"Name your price" clusters

Very low cost at high scale

Lowest cost for time insensitive workloads

Also on-demand and reserve capacity pricing

# SPARK ON EMR ☀New

No additional cost. Available today.

aws.amazon.com/emr/spark

# THANK YOU

MATTHEW@AMAZON.COM

@MZA