

Navigating

Unstructured Data

Matt Brender Developer Advocate

This presentation includes information that is confidential and proprietary to Basho Technologies and should not be forwarded or distributed without Basho's prior written consent. © 2014. Basho Technologies, Inc. All Rights Reserved.

```
"text": "Woot! #GOTOChgo",
"entities": {
  "hashtags": ["#GOTOChgo"],
  "symbols": [],
  "urls": [],
  "user mentions": [{
    "screen name": "mjbrender",
    "name": "Matt Brender",
    "id": 4948123,
    "id str": "42424242",
    "indices": [81, 92]
 }, {
    "screen name": "mjbrender",
    "name": "Matt Brender",
    "id": 376825877,
    "id str": "376825877",
    "indices": [121, 132]
 }]
```

screen_name	text	hashtags	id	id_str
mjbrender		#GOTOChgo	376825877	4242342





Big Data



Big Data

Big Data

Big Data is





And it's a distributed systems problem





Ergo, NoSQL



NoSQL is







For Good Reason





Consistency Level Conflict Resolution Partitioning Strategy



Brewer's Conjecture

Cap theorem states that a distributed system can at most support 2 out of these 3 properties





Consistency Level









Conflict Resolution

Last Write Wins vs. Causal Context



Conflict Resolution











Partition Strategy





A few tactics







What kind of questions do you need to ask?









No Updates. Only Truth.



Making Sense of stream ->-> ->processing Martin Kleppmann @martinkl



Immutable commands

Add To Cart (cust=123, prod = 888, quantity = 1) Update CartQuantity (cust = 123, prod = 888, quantity = 3) Update Cart Quantity (cust=123, prod = 888, quantity = 2) Checkout Cart (cust = 123)

Architectures



Report on this













Error Analysis? NoSQL + Solr



Patterns? Multi-client writes to NoSQL & HDFS







NoSQL + ETL process => Message Queue => Spark and/or Hadoop M/R



















LAMDA



Source: Marz, N. & Warren, J. (2013) Big Data. Manning.



In Review



You can't analyze what you don't have.

And you don't want an analysis system to be unreliable.





everything works at small scale



Summary

• NoSQL

A collection of highly scalable, highly available systems that fall within CAP theorem

- Unstructured Data
 Data that cannot or should
 not be put into a relational
 database
- Immutable Data Modeling data as raw states as oppose to relative updates

- Distributed System
 A tough problem to solve
- **Big Data** Any computation that requires multiple computers



Summary

Lambda Architecture
 A strategy to provide data
 processing in batch and real time simaltaneously

• Solr

Apache project for indexing text for search

 Kafka
 Distributed scalable pub/sub messaging queue
 Hadoop

A framework that allows for the distributed processing of large data sets across clusters

 Spark
 A fast, general engine for large-scale data processing

Storm
 A distributed real-time computation system



Many thanks for all the concepts I stole:

- Martin Kleppman's post on Streaming
- <u>Confluent</u>'s Blog
- Highly Scalable Blog



Thank You!





*