

**Tackling The Challenges of Big Data**  
**Big Data Collection**

**Michael Stonebraker**  
 Professor  
 Massachusetts Institute of Technology

---

---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
 Introduction

**Michael Stonebraker**  
 Professor  
 Massachusetts Institute of Technology

---

---

---

---

---



---

---

---

**Data Curation**

- Ingest
- Validate
- Transform
- Correct
- Consolidate (dedup)
- And visualize information to be integrated

 Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology 

---

---

---

---

---

---

---

---

## Data Warehouse Roots

- Retail sector started integrating sales data into a data warehouse in the early 1990's
- Average system was 2X budget and 2X late
- Because of data integration headaches
- However, warehouse paid for itself within 6 months with smarter buying decisions!



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Issues

- sold \$100K of widgets to IBM, Inc.
- sold 800K Euros of m-widgets to IBM, SA
- Translate currencies
- Is IBM, SA the same as IBM, Inc?
- Are m-widgets the same as widgets?



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## The Pile-On

- Essentially all enterprises followed suit and built warehouses of customer facing data
- Serviced by so-called Extract-Transform-and Load (ETL) tools



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Architecture

- [http://en.wikipedia.org/wiki/Data\\_integration](http://en.wikipedia.org/wiki/Data_integration)



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Traditional Wisdom - ETL

- **Human defines a global schema**
- **Assign a programmer to each data source to:**
  - Understand it
  - Write local to global mapping (in a scripting language)
  - Write cleaning routine
  - Run the ETL
- **Scales to (maybe) 25 data sources**



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Traditional ETL Methodology – Schema Mapping

[http://www.google.com/imgres?imgurl=http://www.xmlschema.info/images/shots/map\\_xml\\_thumb.gif&imgrefurl=http://www.xmlschema.info/xml\\_schema\\_mapping.html&h=469&w=600&sz=63&tbid=70oECAvg0TMwkM:&tbnh=102&tbnw=131&zoom=1&usg=\\_\\_AsYa-CEcZeR6MX9IV8JqvdpX9RM=&docid=MbKKOGXt3LED1M&sa=X&ei=XDGJUpi9OYevsQT5IYGAAg&ved=0CDQQ9QEwAg](http://www.google.com/imgres?imgurl=http://www.xmlschema.info/images/shots/map_xml_thumb.gif&imgrefurl=http://www.xmlschema.info/xml_schema_mapping.html&h=469&w=600&sz=63&tbid=70oECAvg0TMwkM:&tbnh=102&tbnw=131&zoom=1&usg=__AsYa-CEcZeR6MX9IV8JqvdpX9RM=&docid=MbKKOGXt3LED1M&sa=X&ei=XDGJUpi9OYevsQT5IYGAAg&ved=0CDQQ9QEwAg)



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Traditional ETL Methodology – Data Transformation

<http://www.informatica.com/us/products/enterprise-data-integration/powercenter/>



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Current Situation

- **Enterprises want to integrate more and more data sources**
  - Miller beer example
  - Novartis example
  - Goby example (we will see this again)
- **Traditional ETL won't scale!!!!**
- **Point-projects in departments -- Staffed by a data scientist**
  - Brand manager deciding marketing spend
  - Augmenting demographic customer data for department use
- **Traditional ETL way too heavy-weight**



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## The Rest of This Module

- Curation example
- Low end (individual data scientist) support
- Enterprise support



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
Introduction

**THANK YOU**



---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**

**Michael Stonebraker**  
Professor  
Massachusetts Institute of Technology



---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
Issues

**Michael Stonebraker**  
Professor  
Massachusetts Institute of Technology



---

---

---

---

---

---

---

## The Problem

Demo of Goby.com



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## In Summary

- Data is dirty!!!!
- Sometimes not clear how to clean it
  - 2 restaurants at the same address: food court or one went out of business??
- Transformations may be a big problem



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
 Issues

**THANK YOU**




---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**

**Michael Stonebraker**  
 Professor  
 Massachusetts Institute of Technology

---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
 New Ideas - 1

**Michael Stonebraker**  
 Professor  
 Massachusetts Institute of Technology

---

---

---

---



---

---

---

**Startups in This Space (probably a bunch more)**

- Paxata
- Trifacta (commercial Data Wrangler)
- Cambridge Semantics
- Data Tamer
- ClearStory
- Attivio

 Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology 

---

---

---

---

---

---

---

## At Least Two Foci

- Support for the individual data scientist
- Enterprise data integration



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Support for the Data Scientist

- [Wrangler video](#)
- <http://vis.stanford.edu/wrangler/>



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Summary

- Expect more systems in this space
- At low prices
- Market will be gated by the availability of data scientists
  - Insurance example



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---



**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
New Ideas - 1

**THANK YOU**



---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**

**Michael Stonebraker**  
Professor  
Massachusetts Institute of Technology



---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
New Ideas - 2

**Michael Stonebraker**  
Professor  
Massachusetts Institute of Technology



---

---

---

---

---

---

---

## Data Tamer Goals

- Do the "long tail"
  - Better/cheaper/faster than the ad-hoc techniques being used currently
- By inverting the normal ETL architecture
  - Machine learning and statistics
  - Ask for human help only when automatic algorithms are unsure

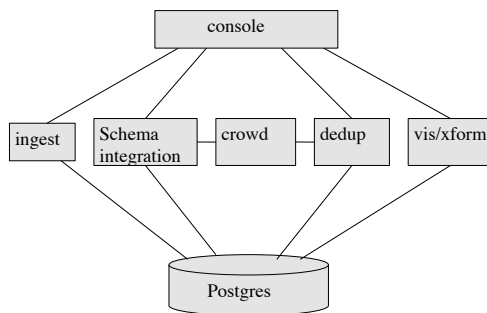


Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Data Tamer Architecture



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Data Tamer -- Ingest

- Assumes (for now) a data source is a collection of records, each a collection of (attribute-name, value) pairs.
- Loaded into Postgres



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Data Tamer – Schema Integration

- Must be told whether there is a predefined partial or complete global schema or nothing
- Starts integrating data sources
  - Using synonyms, templates, and authoritative tables for help
  - 1st couple of sources require asking the crowd for answers
  - System gets better and better over time



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Data Tamer – Schema Integration

- Inner loop is a collection of experts
  - T-test on the data
  - Cosine similarity on attribute names
  - Cosine similarity on the data
- Scores combined heuristically
- After modest training, get 90% of the matching attributes on Goby and Novartis automatically
  - Cuts human cost dramatically



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## Data Tamer – Crowd Sourcing

- Hierarchy of experts
- With specializations
- With algorithms to adjust the “expertness” of experts
- And a marketplace to perform load balancing
- Currently doing a large scale evaluation at Novartis
  - Late flash: it works!!!!



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

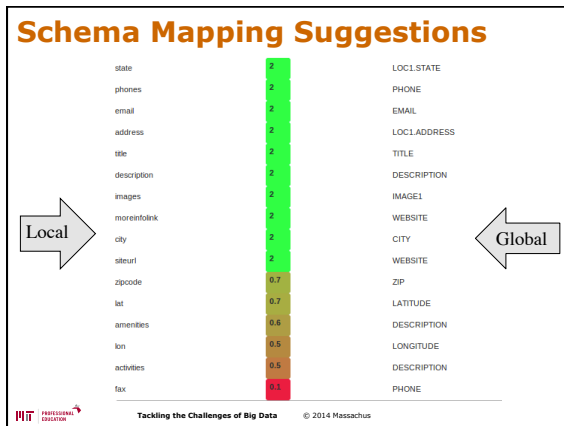
---

---

---

---

---




---

---

---

---

---

---

---

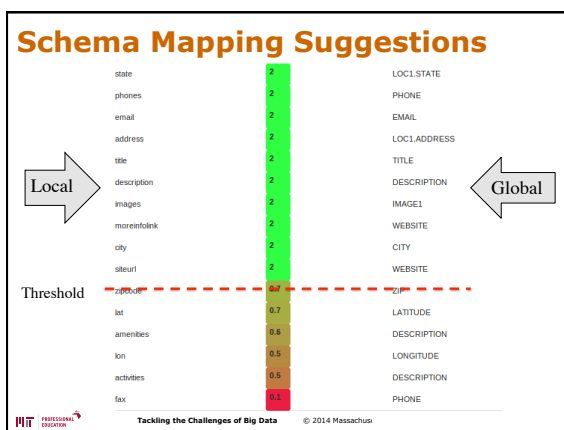
---

---

---

---

---




---

---

---

---

---

---

---

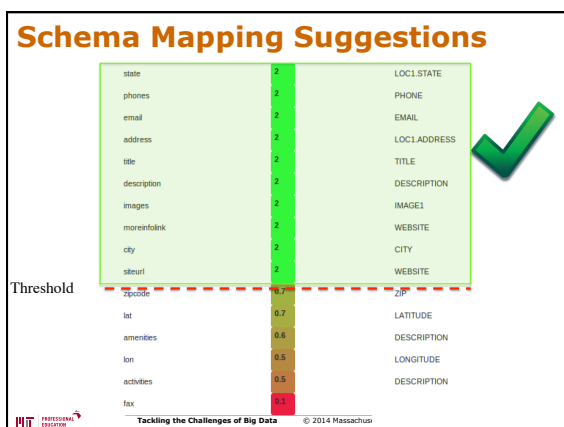
---

---

---

---

---




---

---

---

---

---

---

---

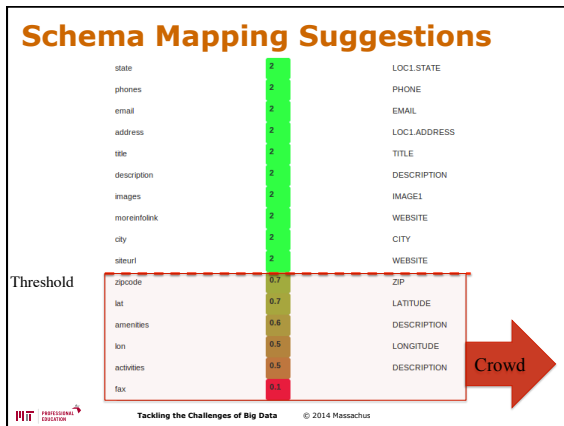
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

### Data Tamer – Entity Consolidation

- On tables defined by schema integration module
- Entity matching on all attributes, weighted by value presence and distribution
- Basically a data clustering problem
- With a first pass to try to identify “blocks” of records
  - Otherwise  $N \times 2$  in the number of records
- Wildly better than Goby; a bit better than domain-specific Verisk module

MIT PROFESSIONAL EDUCATION Tackling the Challenges of Big Data © 2014 Massachusetts Institute of Technology

---

---

---

---

---

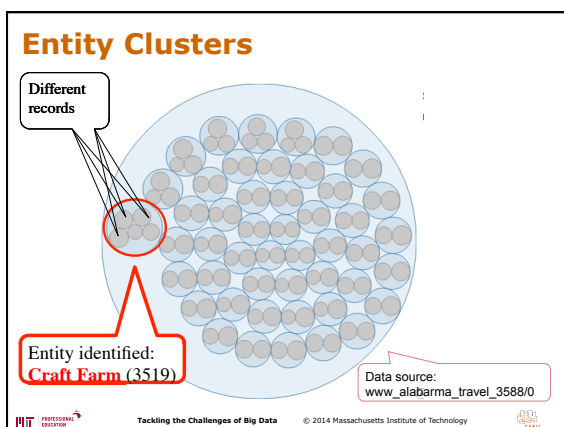
---

---

---

---

---




---

---

---

---

---

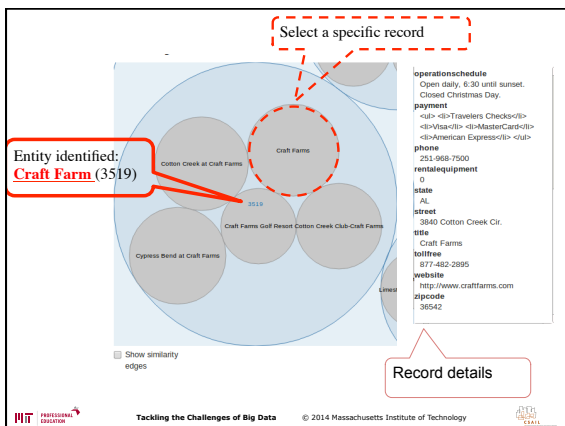
---

---

---

---

---




---

---

---

---

---

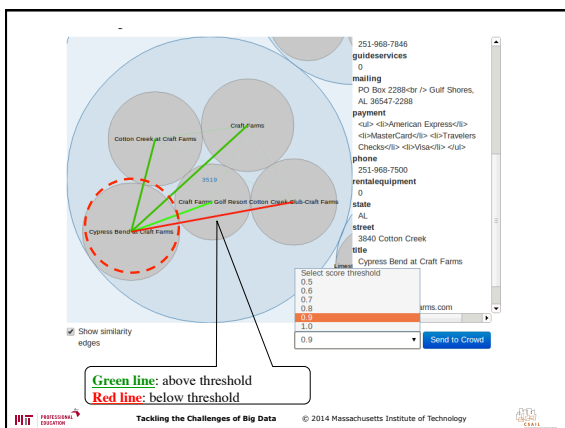
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

## Data Tamer Future

- Text
- Relationships
- Hierarchical data (maybe)
- Adaptors
- Better algorithms
- User-defined operations

---

---

---

---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
New Ideas - 2

**THANK YOU**



---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**

**Michael Stonebraker**  
Professor  
Massachusetts Institute of Technology



---

---

---

---


---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
Summary

**Michael Stonebraker**  
Professor  
Massachusetts Institute of Technology



---

---

---

---

---

---

---

## The Way Forward

- Enterprises will want to integrate more and more data sources
  - This is the number one headache of most of them
  - Too expensive to do manually with a programmer
- Remains to be seen what fraction of the market can be aided by Data Tamer-style tools
  - Initial results are encouraging



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

## The Way Forward

- Cleaning will be a big issue forever
  - How clean does your data need to be?
- I imagine a big database of transformations
  - Pick the one that you need
- Data scientists will have to familiar with this stuff
  - Not just stat and data management



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology




---

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Collection**  
**Data Cleaning & Integration**  
 Summary

**THANK YOU**




---

---

---

---

---

---

---



**Tackling The Challenges of Big Data**

**Big Data Collection**  
**Data Cleaning & Integration**

**Michael Stonebraker**

Professor

Massachusetts Institute of Technology



---

---

---

---

---

---

---