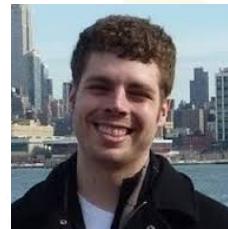


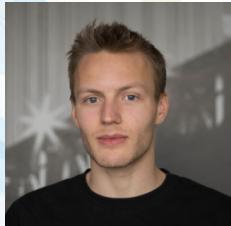
# 3 Mini Talks of 10 minutes each



**The Machine  
Intelligence Landscape:  
A Venture Capital  
Perspective**  
David Beyer



**Algorithms for Anti-Money  
Laundering**  
Richard Minerich



**The future of global,  
trustless transactions on  
the largest graph:  
blockchain**  
Olaf Carlson-Wee



# The Machine Intelligence Landscape: A Venture Capital Perspective

# An Explosion of Data



2002

**5 Exabytes**

Growing by a  
factor of  
**7,040x**



2020

**35,200 Exabytes**

Source: EMC, Forbes estimates

# What is Enabling This Transition?



- Open Source: It's now cheap enough to store ALL the data vs. being selective
- AWS: infrastructure/scale for the masses
- Big Data laying groundwork for Machine Intelligence – the “so what?”

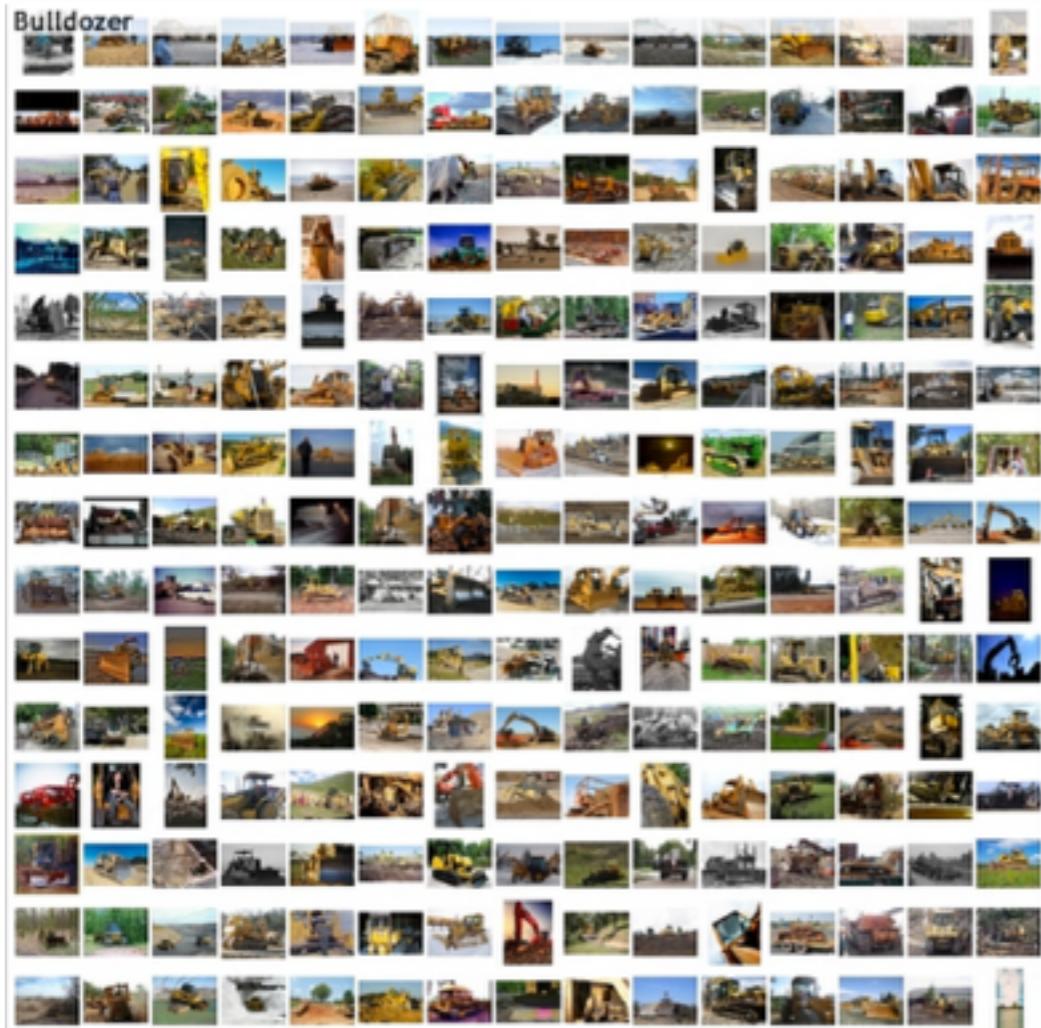
# What is Machine Learning?

“A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E.”

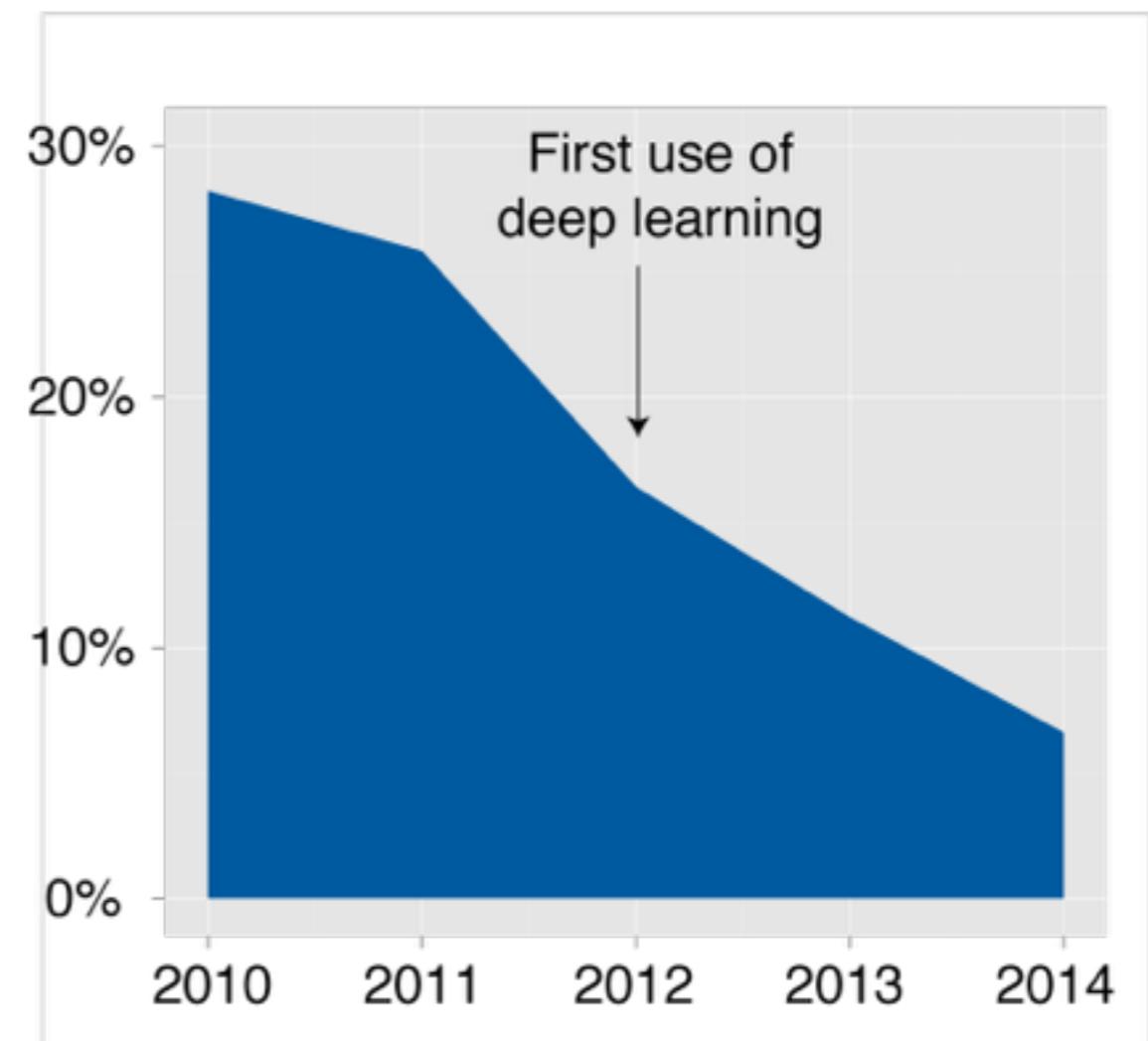
- Tom Mitchell

# Deep Learning Becomes Very, Very Good

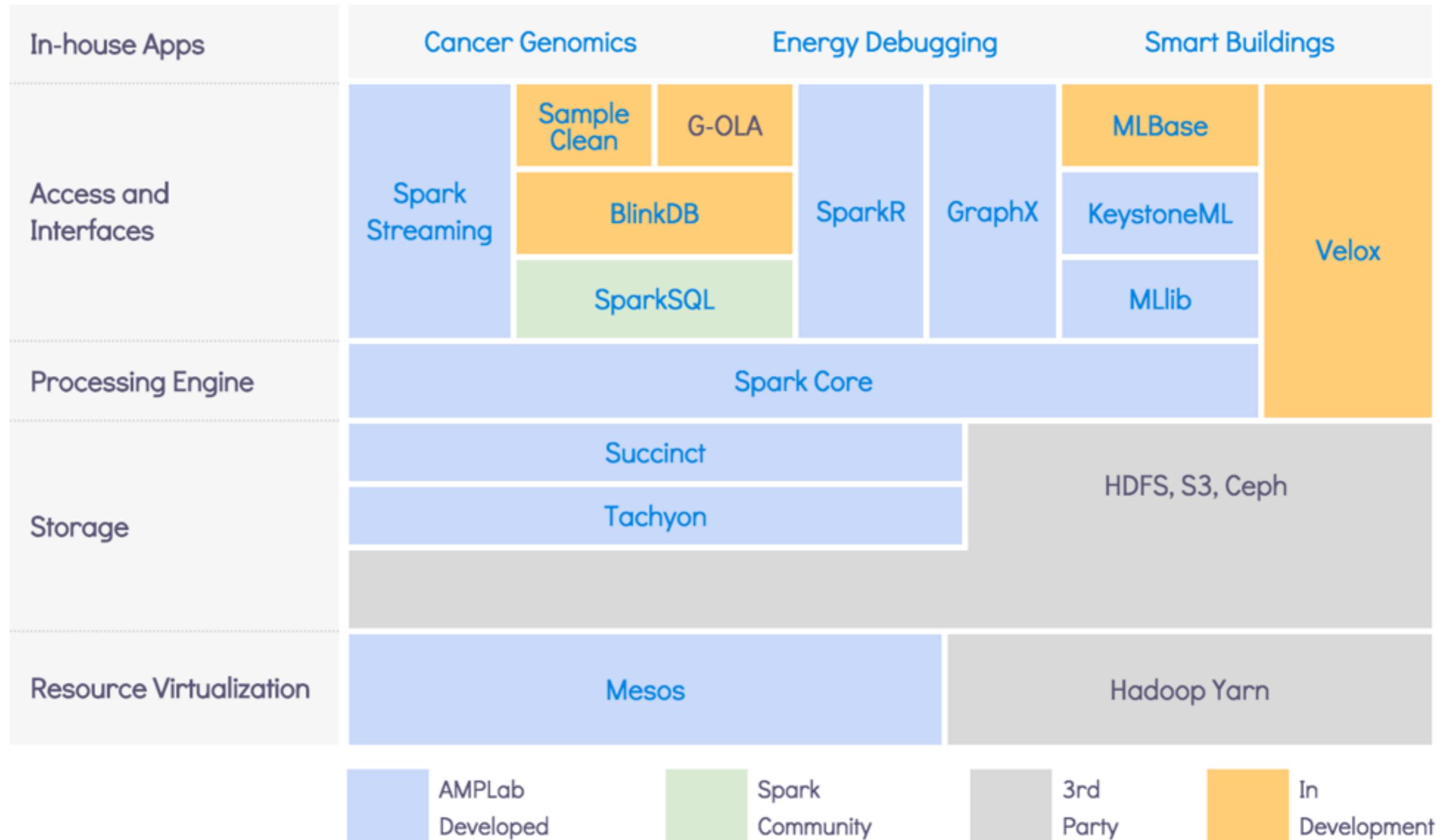
ImageNet examples



Objection classification error rate

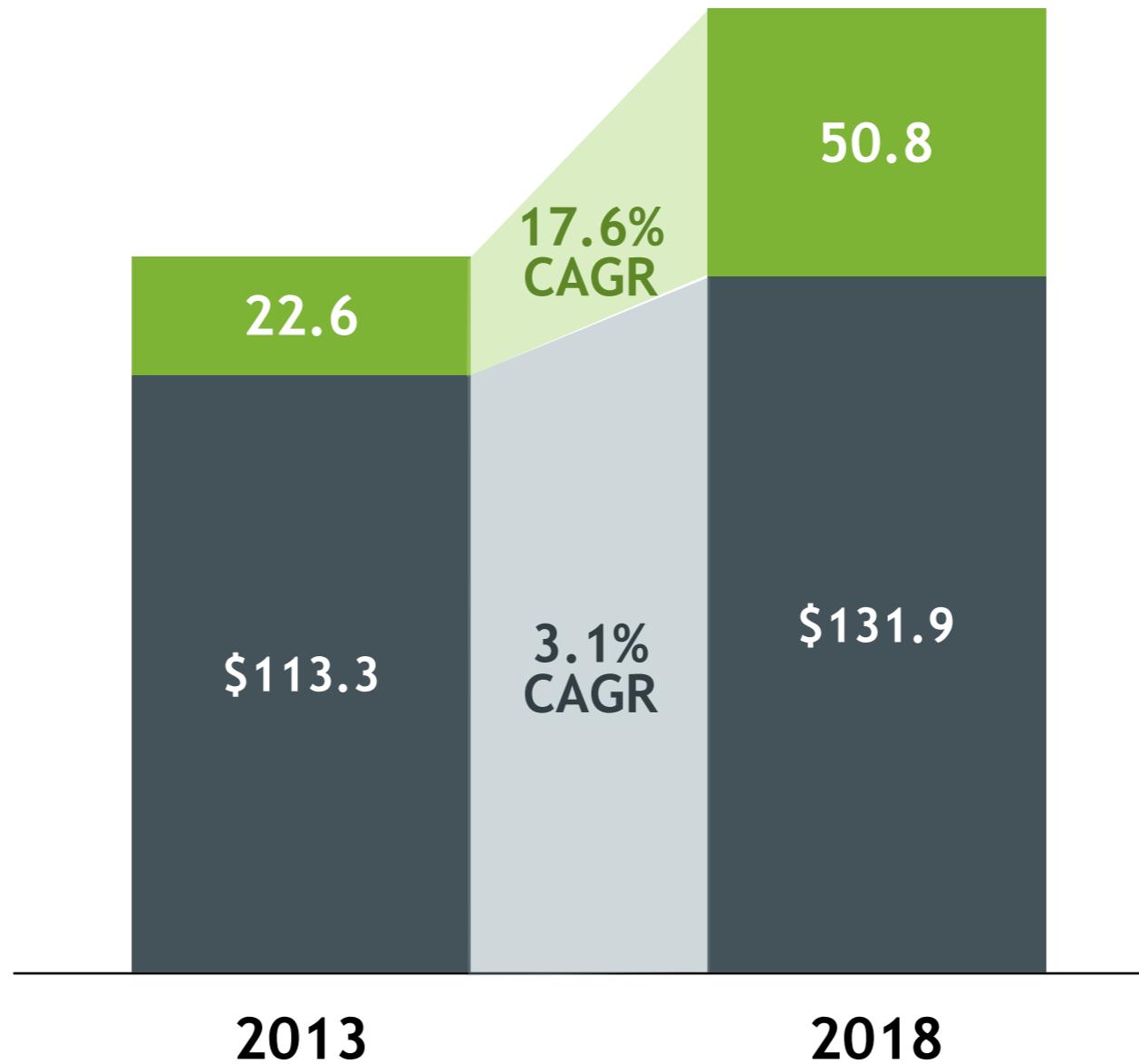


# Recent Innovations in Distributed Computing



# SaaS is Crushing Traditional Enterprise SW

## World Wide Enterprise Applications Market\* (\$B)



\*according to IDC estimates

# Not Much has Changed In 15 Years

File Edit View Navigate Query Tools Help

ORACLE

Saved Queries: All Accounts

Account:

Home Accounts Contacts Opportunities Quotes Sales Orders Service Administration - Business Process

Accounts Home | Accounts List | Global Accounts Hierarchy List | Charts | Account Explorer | Account D&B Explorer | Service Explorer | Accounts Administration

My Accounts Menu New Edit Delete Query Collaborate Create Team Space 1 - 10 of 18+ +

Account Name	Site	Main Phone #	Status	URL	DUNS #	Team Space	Industries
Marriott International HQ	HQ	(800) 234-5000	Gold	www.marriott.com			hotels & motels
3Com	Headquarters	(773) 326-5000	Gold	www.3com.com	842079576		manufacturing industries
3Com Distribution	UK	+0283456857	Active		842079576		management consulting
3Com Research	US	(415) 329-6500	Active	www.3com.com	099956906		manufactured hardware
9 Telecom	France	+33155206242	Active				steel pipe & tubes
AMCO Communications	Chicago, IL	(847) 491-2300	Active	www.amco.net			
Acer America, Inc.	San Jose, Ca	(408) 922-2957	Current Customer	www.acer.com			
Acer Stores	Clayton	(925) 745-2000	Active	www.acerstore			
Aegis	Warehouse		Active				
Air France	France	+33141567800	Active				

Marriott International HQ

Menu New Delete Query

Account Name: Marriott International HQ Site: HQ Account Team

Address: 10400 Fernwood Rd Address Line 2: Main Phone #

City: Bethesda State: MD Main Fax #

Zip Code: 20817 Country: USA URL

Opportunity

Edge Emergency Generator

Opportunity Detail

Opportunity Owner	Jason Venable (Change)	Amount	\$275,000.00
Account Name	Edge Communications	Close Date	11/6/2007
Opportunity Name	Edge Emergency Generator	Stage	Id: Decision Makers
Type	Existing Customer - Replacement	Probability (%)	60%

Other Information

Description	Budget Low	\$262
Next Step	Budget High	\$780

Additional Information

Order Number	Main Competitor(s)	John Deere, Mitsubishi, Hawkpowers
Tracking Number	Delivery Installation Status	
Custom Links	Delivery Status	

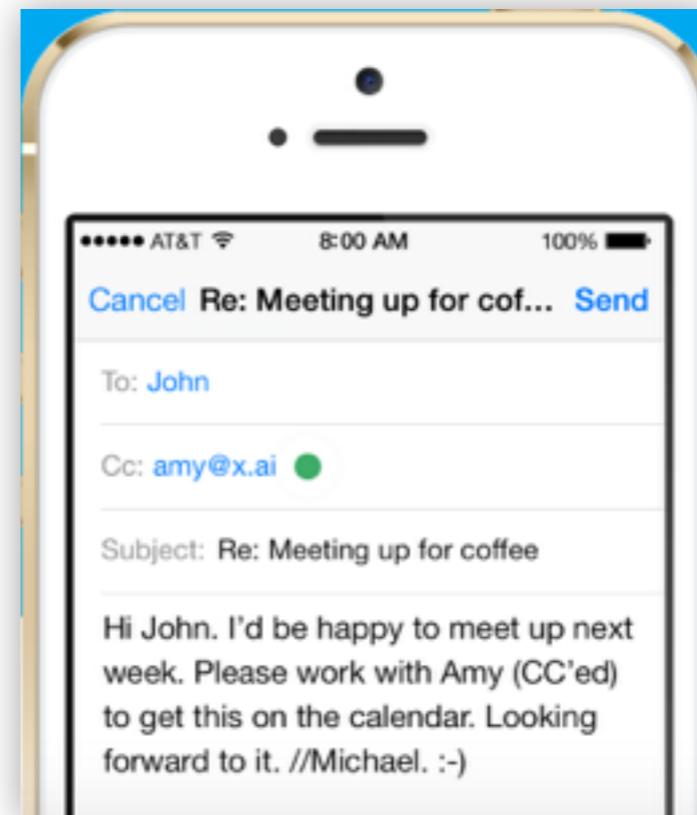
System Information

Created By	Jason Venable, 11/25/2008 3:08 PM	Last Modified By	Jason Venable, 3/9/2010 10:51 PM
------------	-----------------------------------	------------------	----------------------------------

Products (Standard)

Action	Product	Quantity	Sales Price	Date	Line Description	Start Workflow
Edit Del	GenWatt Diesel 1000kW	1.00	\$100,000.00			
Edit Del	GenWatt Diesel 200kW	1.00	\$25,000.00			
Edit Del	GenWatt Gasoline 2000kW	1.00	\$150,000.00			

# Samantha (Her) vs. Amy (x.ai)



# Machine Intelligence LANDSCAPE

## CORE TECHNOLOGIES

### ARTIFICIAL INTELLIGENCE



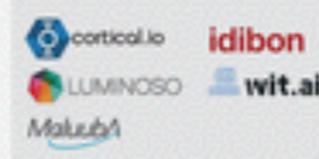
### DEEP LEARNING



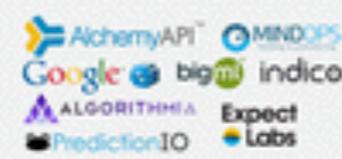
### MACHINE LEARNING



### NLP PLATFORMS



### PREDICTIVE APIs



### IMAGE RECOGNITION



### SPEECH RECOGNITION



## RETHINKING ENTERPRISE

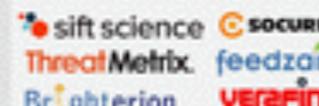
### SALES



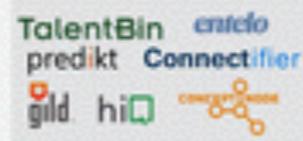
### SECURITY / AUTHENTICATION



### FRAUD DETECTION



### HR / RECRUITING



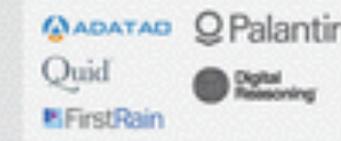
### MARKETING



### PERSONAL ASSISTANT



### INTELLIGENCE TOOLS



## RETHINKING INDUSTRIES

### ADTECH



### AGRICULTURE



### EDUCATION



### FINANCE



### LEGAL



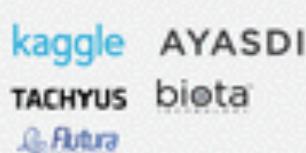
### MANUFACTURING



### MEDICAL



### OIL AND GAS



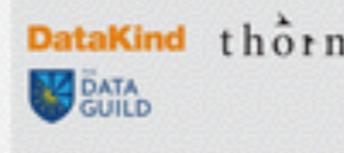
### MEDIA / CONTENT



### CONSUMER FINANCE



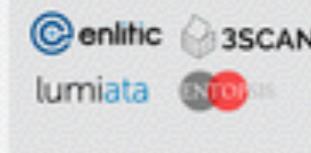
### PHILANTHROPIES



### AUTOMOTIVE



### DIAGNOSTICS



### RETAIL

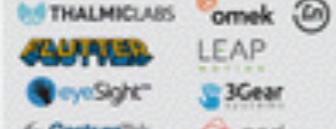


## RETHINKING HUMANS / HCI

### AUGMENTED REALITY



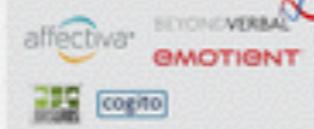
### GESTURAL COMPUTING



### ROBOTICS



### EMOTIONAL RECOGNITION



### HARDWARE



### DATA PREP



### DATA COLLECTION



# THE FUTURE OF BLOCKCHAIN COMPUTING



OLAF CARLSON-WEE  
HEAD OF RISK, COINBASE

# Algorithms For Finding The Bad Guys

---

Richard Minerich, Director of R&D at  
[@Rickasaurus](https://twitter.com/Rickasaurus)



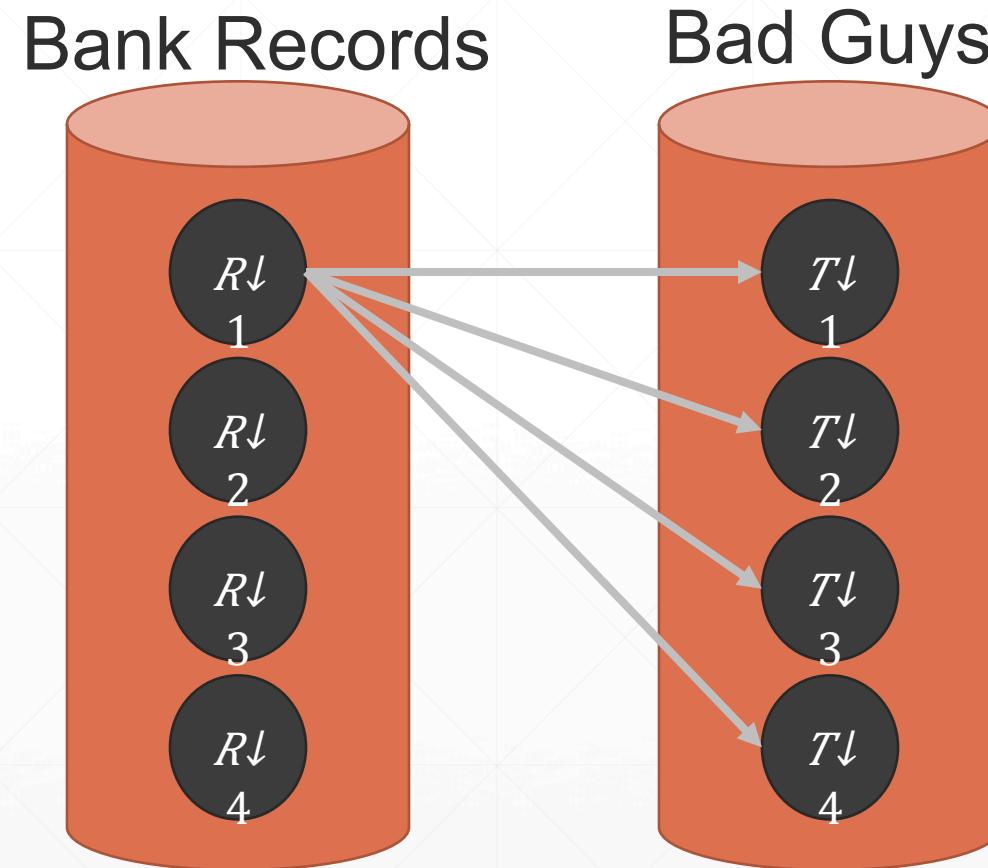
# Entity Resolution in Theory

Example Variations:

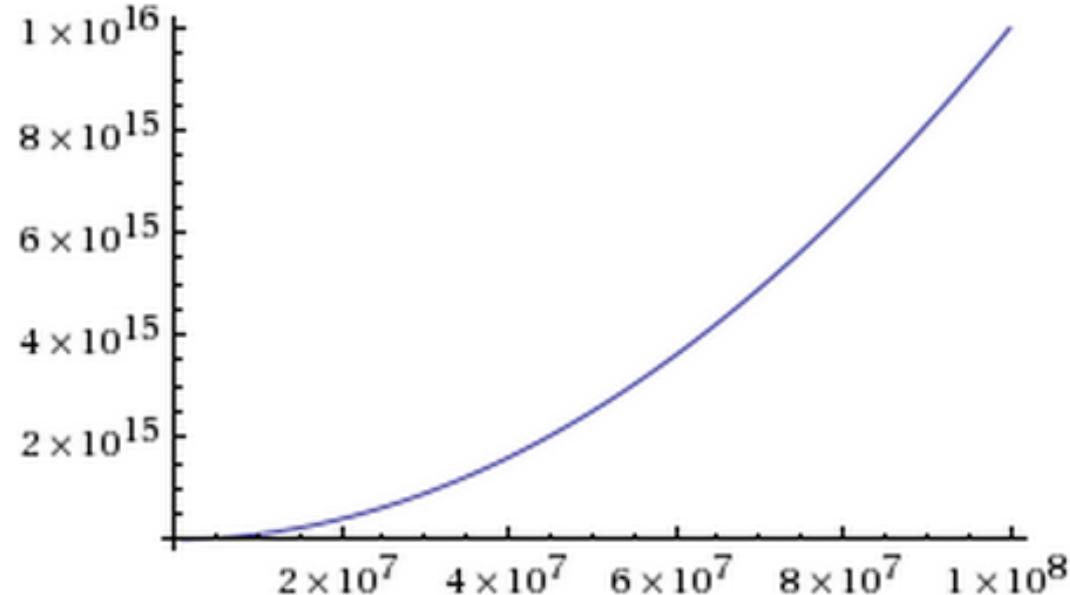
- Aggregating Products
- Finding Medical Records
- Resolving Paper Authors
- Census
- Finding Bad Guys
- Database Deduping

Different tradeoffs per domain

$P(R \downarrow n \text{ represents the same entity as } T \downarrow m)$



# Too Many Comparisons!



- 100 Million x 100 Million = 10 quadrillion pairs
- 86,400,000 milliseconds per day
- One pair per ms: ~116,000,000 days to compute (~317K years)

---

# How do we beat N\*M? Blocking Algorithms.

Input:

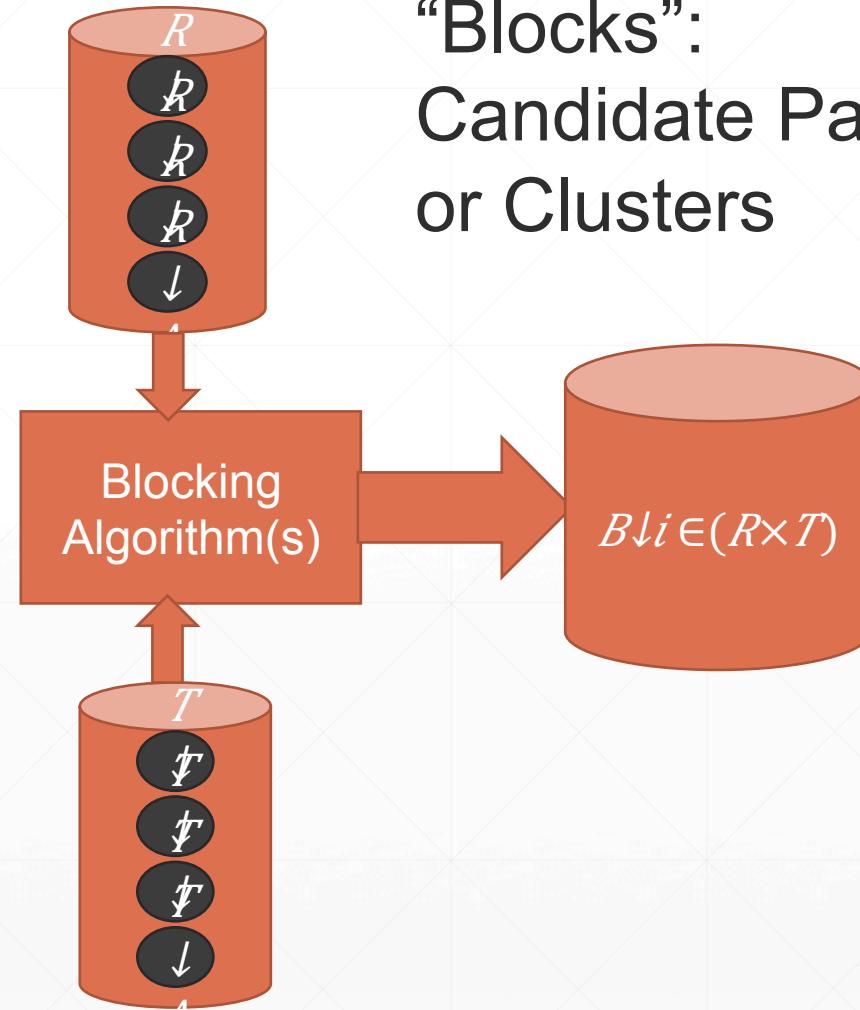
- Source Records R
- Target Records T

Output:

- Blocks of Similar Records

$$B \downarrow i \in (R \times T)$$

“Blocks”:  
Candidate Pairs  
or Clusters



# Simplest: Key-based Blocking

RecID	GivenName	Surname	Postcode	Suburb
r1	peter	christen	2010	north sydney
r2	paul	smith	2600	canberra
r3	pedro	kristen	2000	sydney
r4	pablo	smyth	2700	canberra sth

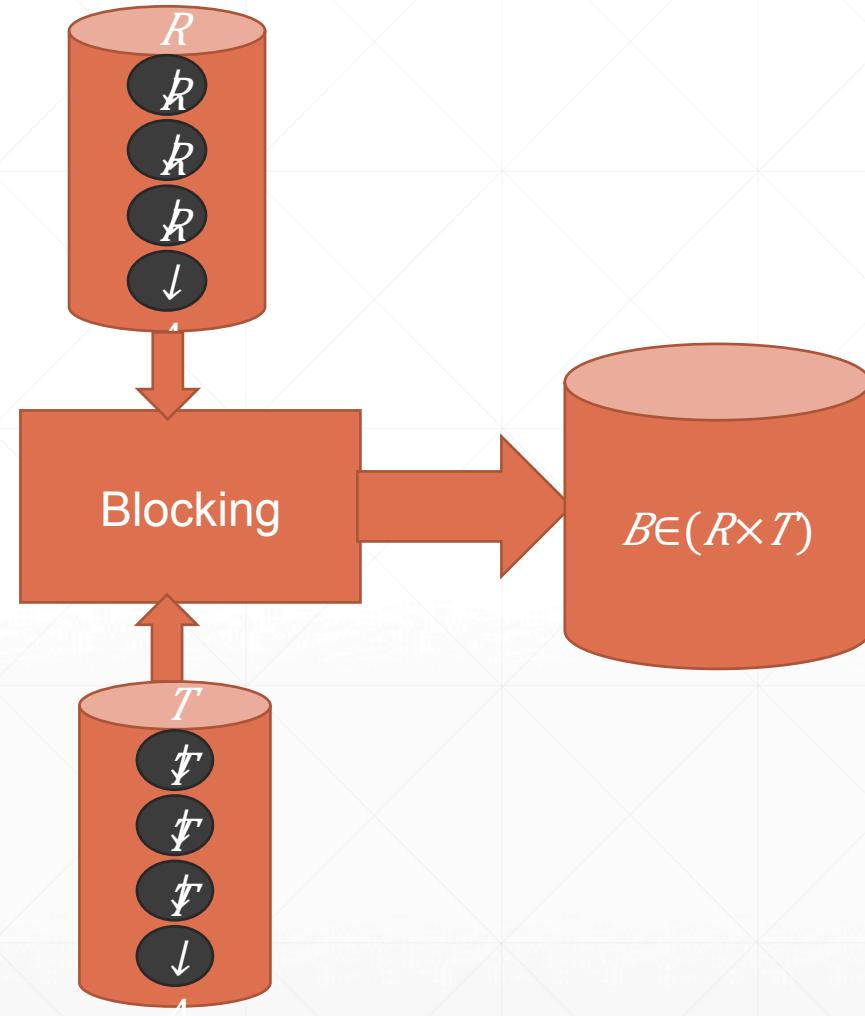
- Nothing's easier than a table lookup!
- Many ways to key, choosing is hard
- Small errors can cause misses
- What about missing data?

RecID	PC+Sndx(GiN)	Fi2D(PC)+DMe(SurN)	La2D(PC)+Sndx(SubN)
r1	2010-p360	<b>20-krst</b>	10-n632
r2	2600-p400	26-sm0	<b>00-c516</b>
r3	2000-p360	<b>20-krst</b>	00-s530
r4	2700-p140	27-sm0	<b>00-c516</b>

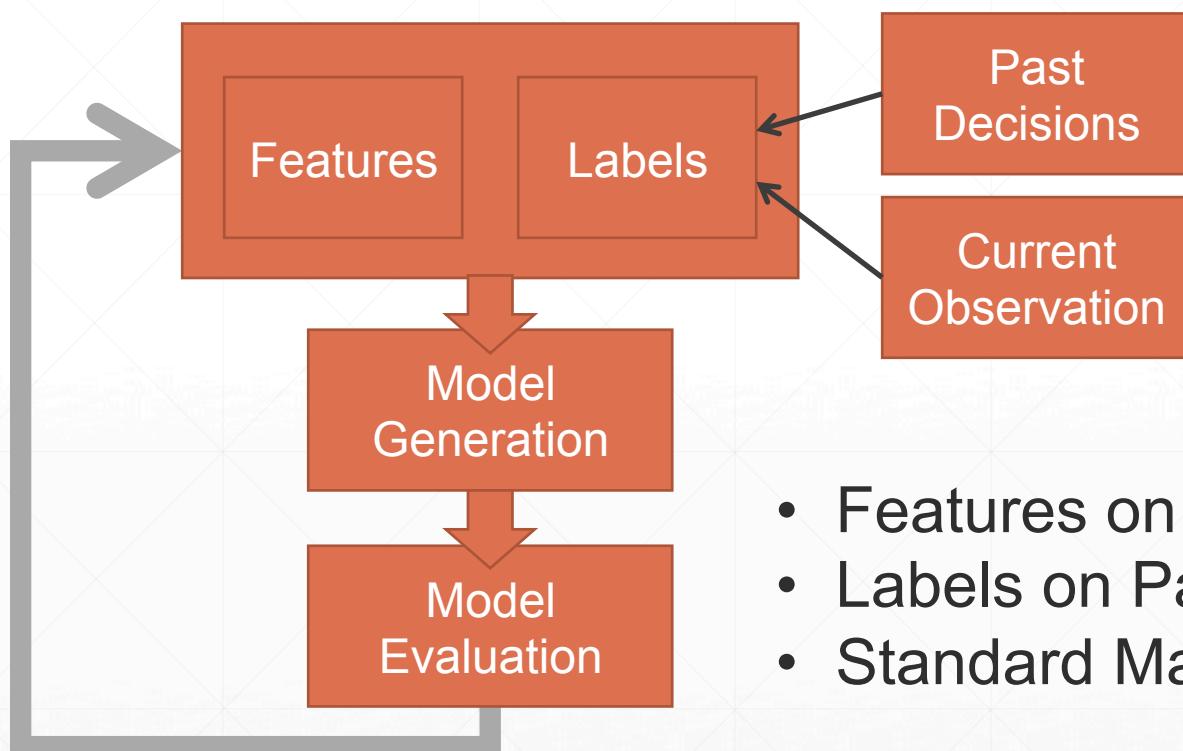
# Many ways to Block

- Sorted Neighborhood
- Suffix Arrays/Trees
- Various kinds of Q-gram Indices
- Metric Space Embedding
- Semantic Hashing
- Cluster-based approaches

Best to use a mix.



# The Basics of Pairwise Scoring



NAME	LARRY O' BRIAN
STATE	CANADA
CITY	Montreal
STATE	Quebec
ADDRESS	121 Buffalo Drive
ZIP	H3G 1Z2
DOB	10/24/80

- Features on the Similarity of Fields
- Labels on Pairs are True/False
- Standard Machine Learning Techniques Apply

# Simplest: Empirical Summed Similarity

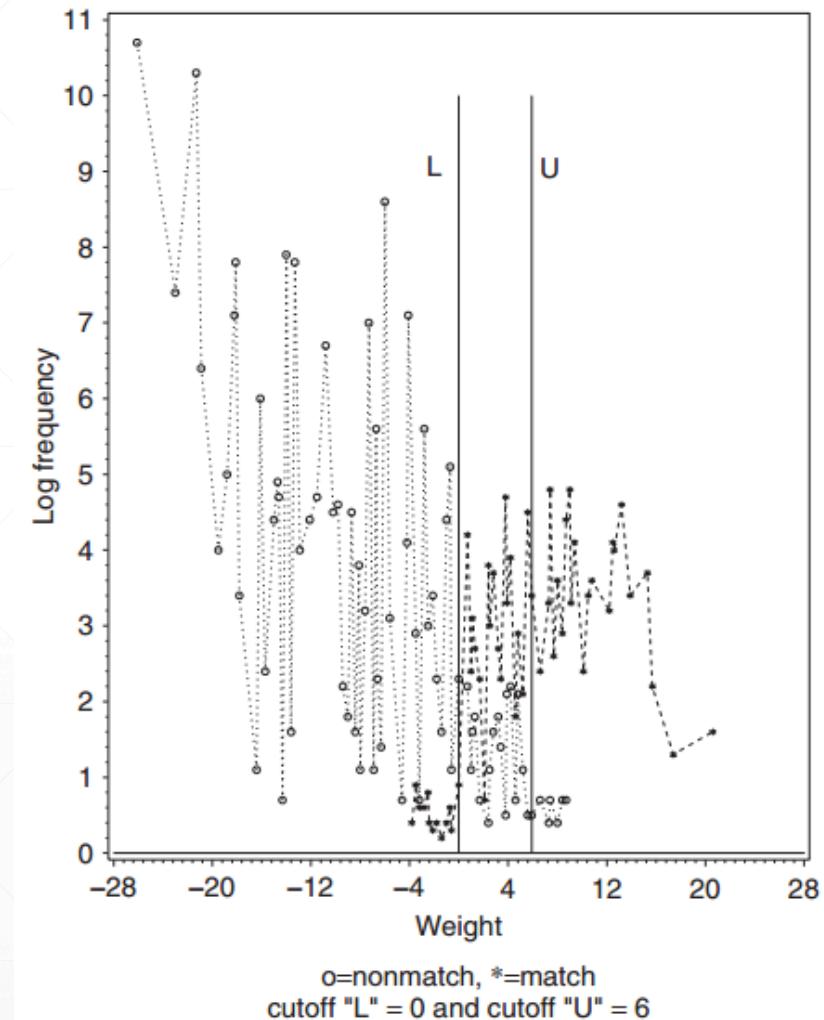
- F: feature functions (0 .. m) : (a,b)  $\rightarrow$  [0, 1]
- W: feature weights (0 .. m) : {0+}
- $SimSum(a,b) = \sum_{i=0}^m f_i w_i$

Thresholds such that:

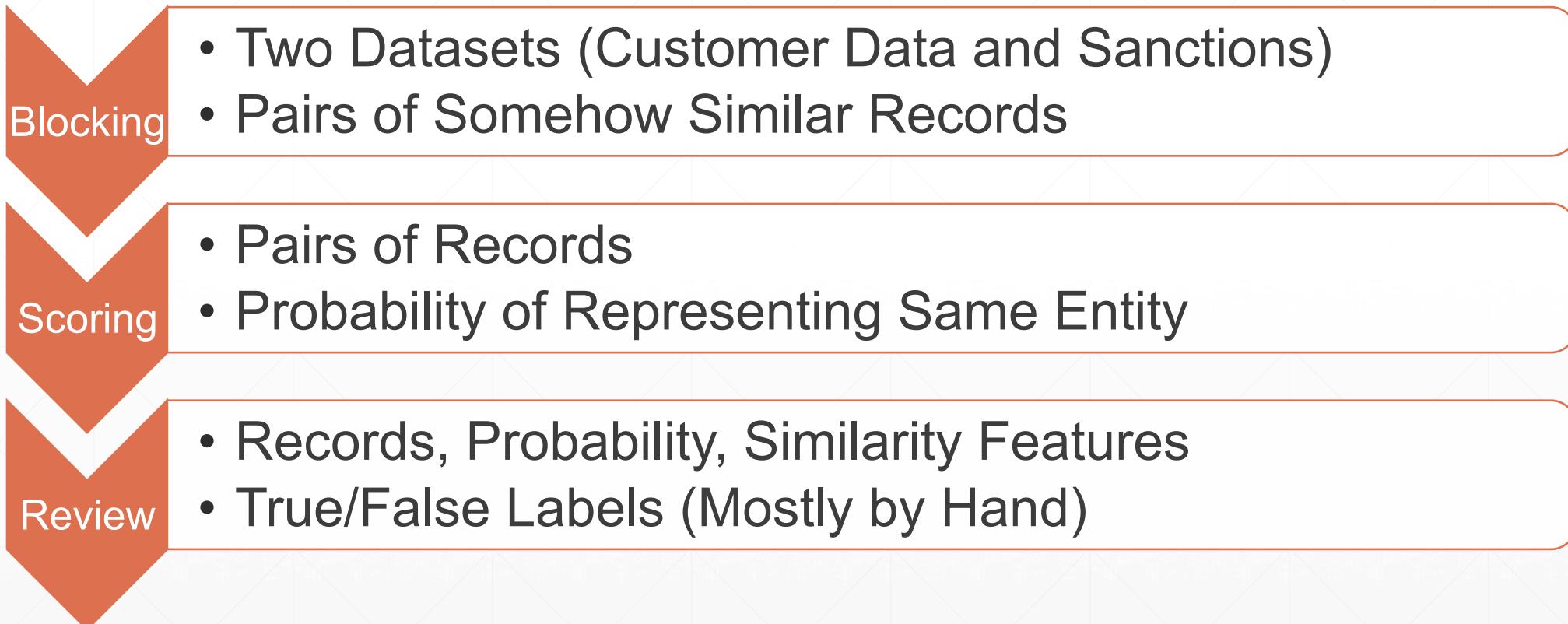
Match:  $SimSum(a,b) \geq \text{Upper}$

Review:  $\text{Lower} \leq SimSum(a,b) \leq \text{Upper}$

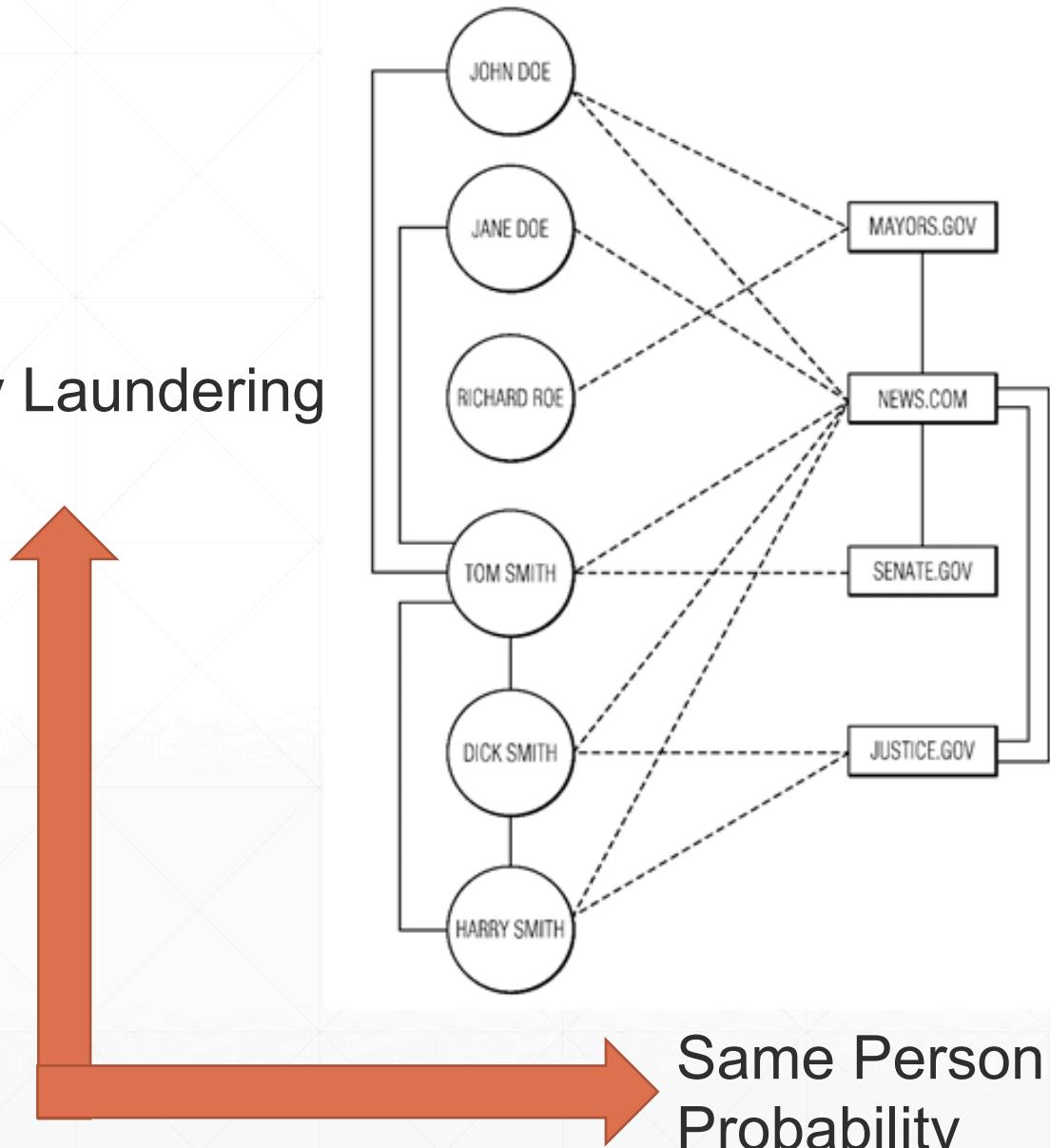
Discard:  $SimSum(a,b) \leq \text{Lower}$



# The Pairwise Entity Resolution Process



# Money Laundering Risk



# Another Dimension to aid review: **Risk vs Probability**

## Useful in many domains:

- How Likely To Launder Money?
  - How much money do I lose if I get the product wrong?
  - Risk of Incorrect Medical Diagnosis

# Simplest Ranking? You Can “Learn to Rank” with Regression.

- The features are the difference in would-be regression features
- The value to predict is the difference in label rank

Select 2 labeled samples randomly =>  $(x_1, y_1)$   $(x_2, y_2)$

$$x = x_1 - x_2$$
$$y = y_1 - y_2$$

	Sample 1	Sample 2	Result
Names?	1	1	0
Addresses?	1	0	1
DOB?	0	1	-1
Same Person?	0	0	0

# Page Rank: The Easy Parts

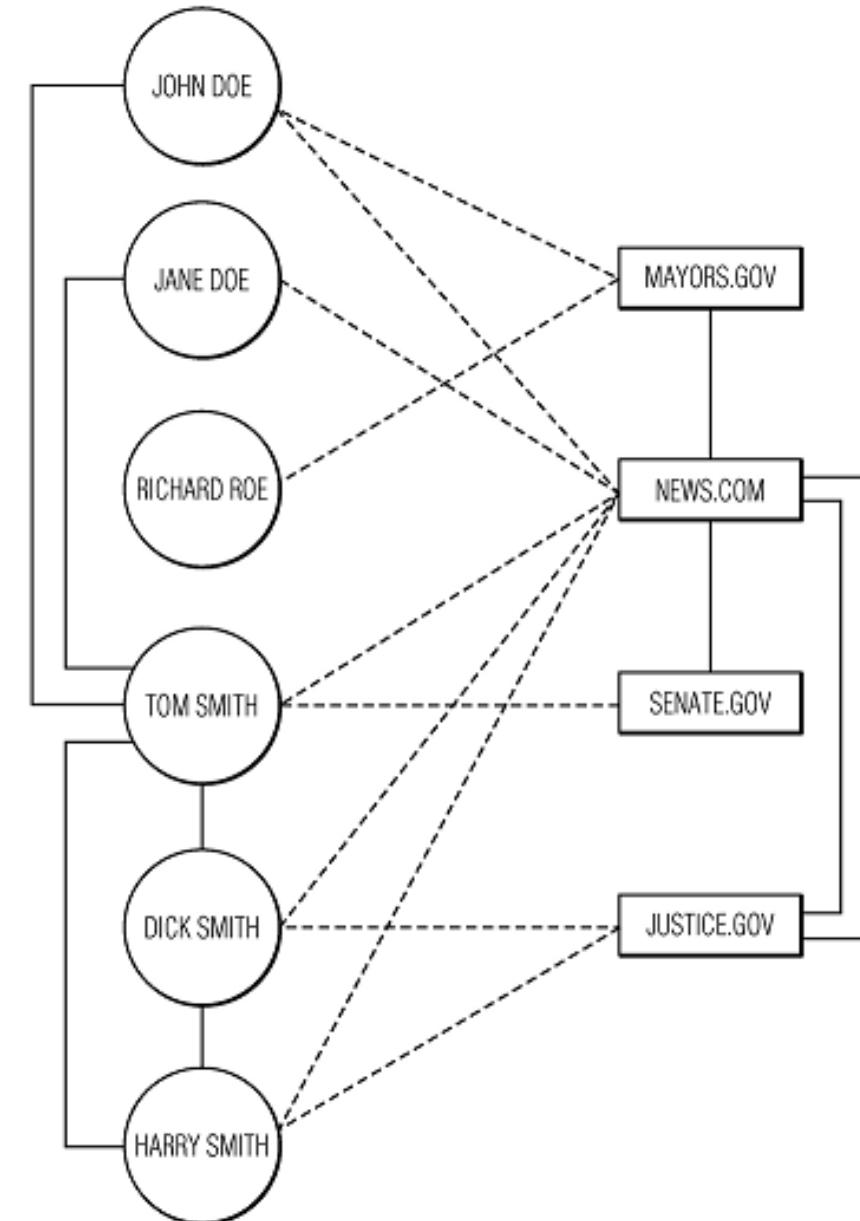
$$\mathbf{R} = \begin{bmatrix} (1-d)/N \\ (1-d)/N \\ \vdots \\ (1-d)/N \end{bmatrix} + d \begin{bmatrix} \ell(p_1, p_1) & \ell(p_1, p_2) & \cdots & \ell(p_1, p_N) \\ \ell(p_2, p_1) & \ddots & & \vdots \\ \vdots & & \ell(p_i, p_j) & \\ \ell(p_N, p_1) & \cdots & & \ell(p_N, p_N) \end{bmatrix} \mathbf{R}$$

```
//Calculate the ranking given a matrix and initial vector
let private calcRanking (A:matrix) (E:Vector<_>) (x:float) (y:float) (iterations:int) =
    let rec calc R n =
        //Calculate new matrix
        let R' = x*(A*R) + y*E

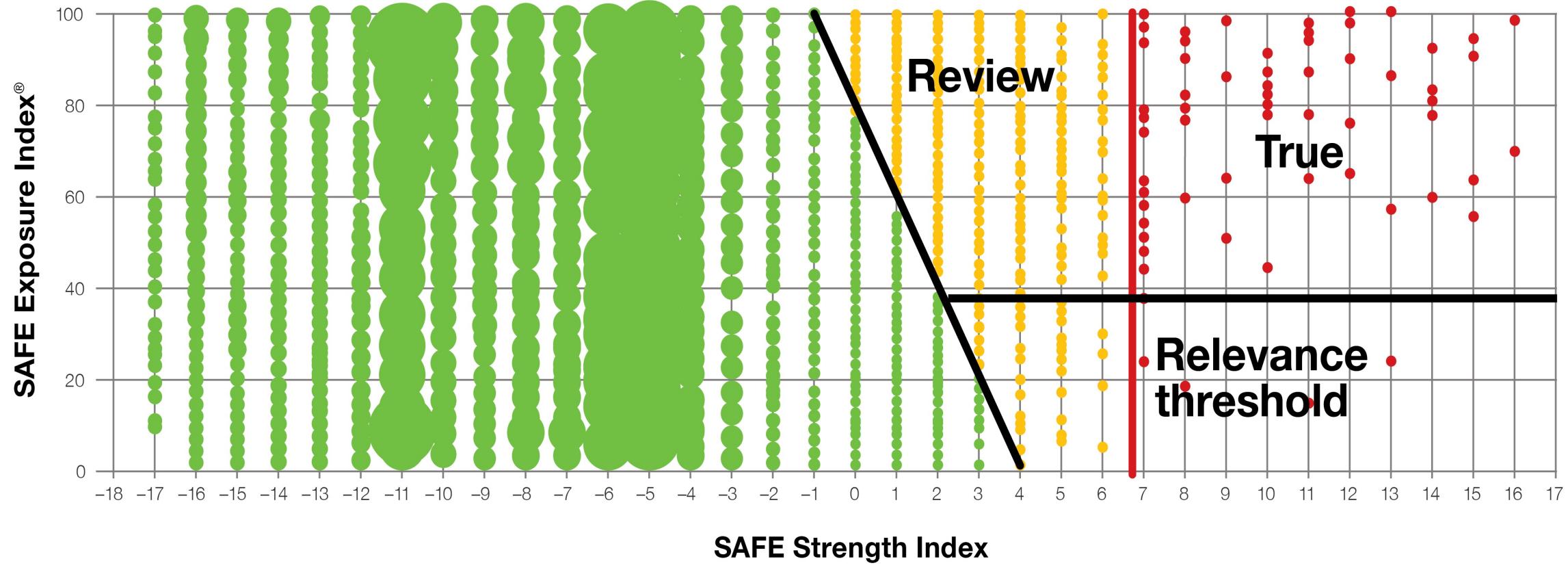
        //Decide when to return results
        if n = iterations then ((Vector.norm(R' - R)), R)
        else calc R' (n + 1)
    calc E 1
```

## Page Rank: The Hard Parts

- Domains, Websites, Pages in Context
- Determining Initial Risk for Sources
- 27 Pages of Data Transformation Code
- Fluctuation with no changes
- Prediction and Explainability



# Combining Ranking and Probability: Big Picture



## Thank You! Questions?

You can read more on my blog at:  
<http://richardminerich.com>

Contact me on twitter:  
[@Rickasaurus](https://twitter.com/Rickasaurus)



Email me with questions:

[rick@bayardrock.com](mailto:rick@bayardrock.com)

Check out the NYC F# User Group:  
<http://www.meetup.com/nyc-fsharp>



---

# Questions