

Data Cleaning With The Google

Bob Lannon
Senior NLP Analyst, Verilogue



Data is Open! Long Live Data!



NOAA NATIONAL OCEANIC AND
ATMOSPHERIC ADMINISTRATION
UNITED STATES DEPARTMENT OF COMMERCE

Google

Fusion Tables



infochimps

SPARC®



ANONYMOUS



open
data
PHILLY



How?

KIDNAPPING THREAT REP BY *JAYSH AL-MAHDI* BAGHDAD (ZONE 15)
(ROUTE UNKNOWN): 0 INJ/*DAN*
REF: BAGSTAT

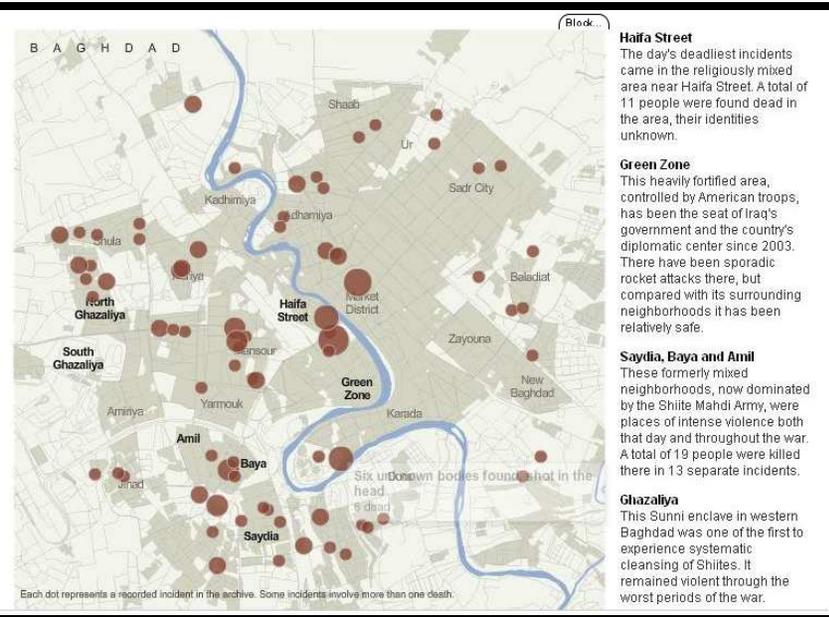
DOI: 22 DEC 06

TITLE: ALLEGED JAYSH AL-MAHDI PLANS TO KIDNAP U.S. SOLDIERS
IN BAGHDAD, IRAQ

AS OF EARLY DECEMBER 2006, JAYSH AL-MAHDI ALLEGELY PLANNED TO ATTACK U. S. HUMVEES TRAVELING IN TWO TO THREE CAR CONVOYS WITH THE INTENT TO KIDNAP U.S. SOLDIERS IN BAGHDAD, IRAQ. JAYSH AL-MAHDI ALLEGEDLY PLANNED TO CONDUCT THE KIDNAPPINGS SOMETIME AROUND THE NEW YEAR, 30 DECEMBER 2006 TO EARLY JANUARY 2007. A SENIOR JAYSH AL-MAHDI COMMANDER, HASAN ((SALIM)), ORDERED A SUBORDINATE, SHAYKH AZHAR AL-((DULAYMI)), TO PLAN AND EXECUTE THE ATTACK. DULAYMI PLANNED TO TARGET U.S. CONVOYS CONSISTING OF TWO TO THREE HUMVEES AS THEY TRAVELED INTO UNDERGROUND TUNNELS IN AL-QADISIYA DISTRICT, BAGHDAD. DULAYMI SPECIFICALLY PLANNED TO USE THE AL-QADISIYA AREAS TO STAGE HIS ATTACKS INCLUDING AL-BALADIYAH DISTRICT STREETS, SADR CITY AND PALESTINE STREET, THE AREA OF AL-SHA'AB WHICH WAS LOCATED NORTH EAST OF SADR CITY, AND AL-SULAYKH. DULAYMI PLANNED TO USE FAKE ROAD BLOCKS AND TUNNELS TO STOP THE HUMVEES AND THEN ATTACK THEM WITH IMPROVISED EXPLOSIVE DEVICES (IEDS) AND SMALL ARMS TO RENDER THE VEHICLES IMMOBILE. ONCE THE VEHICLES WERE STOPPED AND INCAPACITATED, DULAYMI REPORTEDLY PLANNED TO TARGET AND KIDNAP A FEW U.S. SOLDIERS AS POSSIBLE IN AN ATTEMPT TO INCREASE THE NUMBER OF HOSTAGES. DULAYMI WAS ORDERED TO TAKE THE CONVOYS TO SADR CITY. DULAYMI WAS PREVIOUSLY A SUNNI BUDDHIST AND TOOK TO SHI'A WHILE STUDYING IN AN NAJAF UNIVERSITY IN IRAQ IN 1998.

A Deadly Day In Baghdad

Violence peaked in December 2006, just two months before American troops arrived as part of what was later called "the surge." At right are the details of one of the city's deadliest days, Dec. 20. There were 114 separate episodes of violence that day, resulting in the deaths of about 160 Iraqi citizens and police officers.



Case Study: Physician Reporting

- ▶ Please indicate any additional conditions that the patient is suffering from:

<input type="checkbox"/>	Acute infection
<input checked="" type="checkbox"/>	Allergic rhinitis
<input type="checkbox"/>	Allergies
<input type="checkbox"/>	Alzheimer's disease
<input checked="" type="checkbox"/>	Anemia
<input type="checkbox"/>	Angina
<input type="checkbox"/>	Anxiety
<input type="checkbox"/>	Arthritis
<input checked="" type="checkbox"/>	Asthma
<input type="checkbox"/>	Atrial fibrillation
<input checked="" type="checkbox"/>	Attention deficit hyperactivity disorder (ADHD)
<input type="checkbox"/>	Benign prostatic hyperplasia (BPH)
<input type="checkbox"/>	Bipolar disorder
<input type="checkbox"/>	Cancer
<input type="checkbox"/>	Chronic fatigue syndrome
<input type="checkbox"/>	Chronic obstructive pulmonary disease (COPD)
<input type="checkbox"/>	Chronic otitis media

...

Other (Please Specify): _____

Zone of Pain



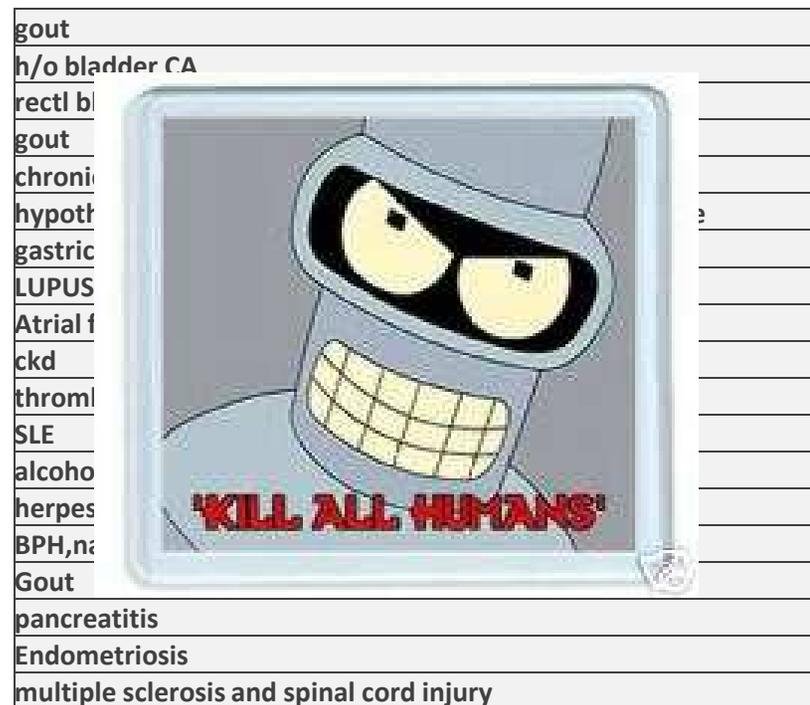
Case Study: Physician Reporting

- ▶ **Could the survey have been implemented better? (Yes)**
 - ▶ Auto-complete on entry
 - ▶ Post-processing inputs and asking entrant for confirmation
- ▶ **Each solution to this problem has problems**
 - ▶ Auto-complete fails
 - ▶ User spends more time filling out form (I hate you Amtrak.com)
- ▶ **When things don't work seamlessly:**
 - ▶ Same messy data
 - ▶ ...OR incomplete entries (most tragic)



End product

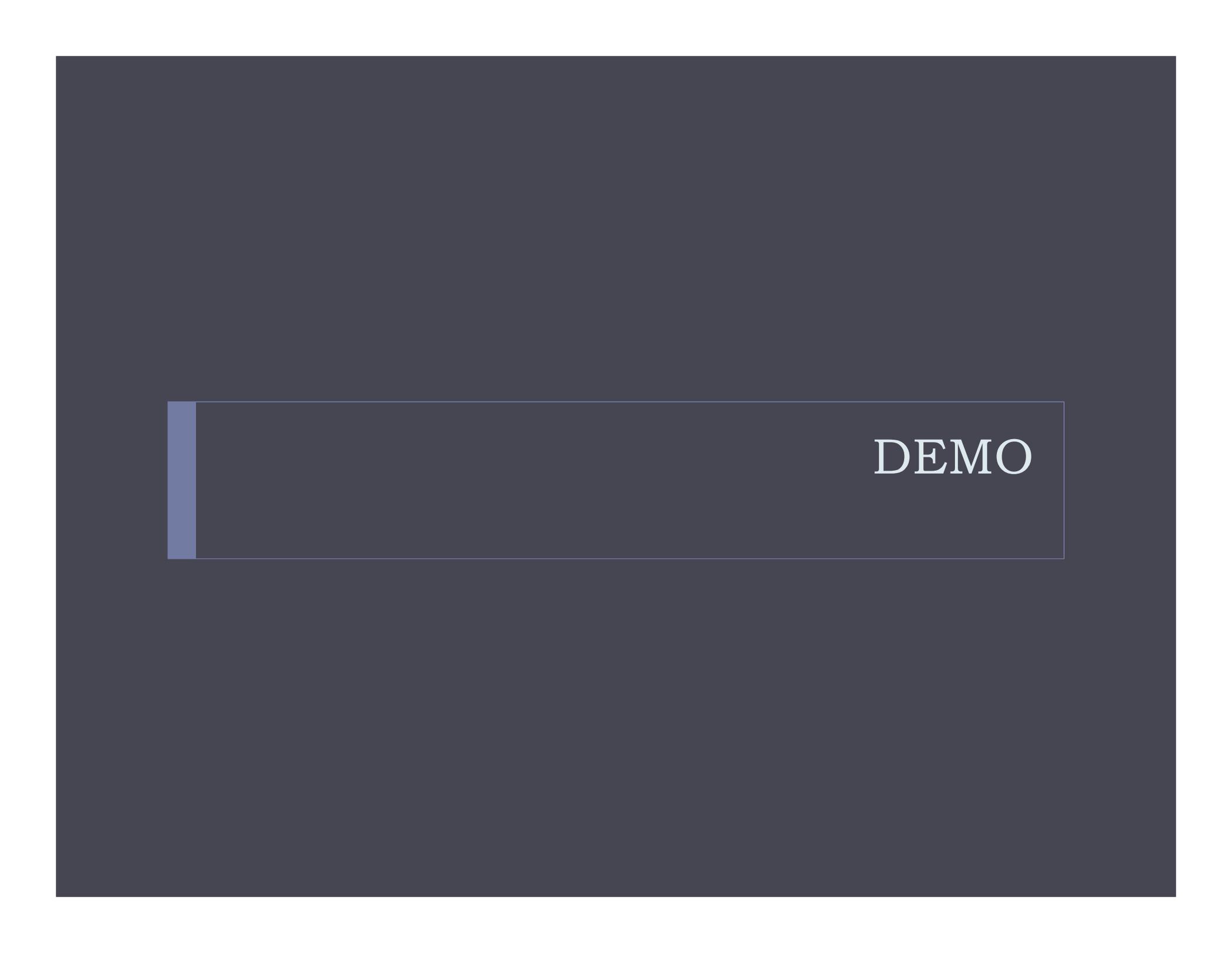
- ▶ What you will receive will probably be either
 - ▶ riddled with incomplete info
 - ▶ full of typos, errors and inconsistencies



But Never Fear!

- ▶ Google Refine is Here!





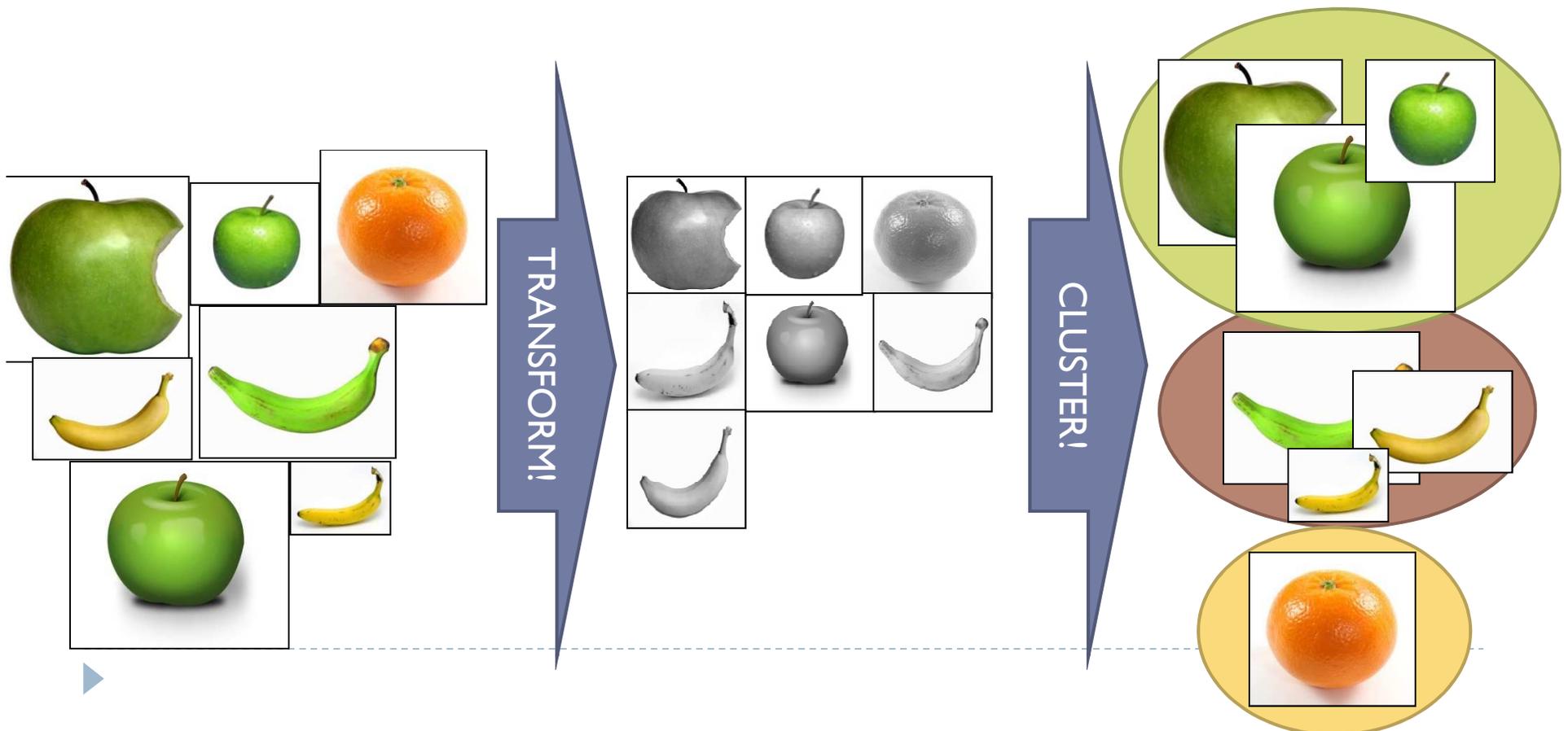
DEMO



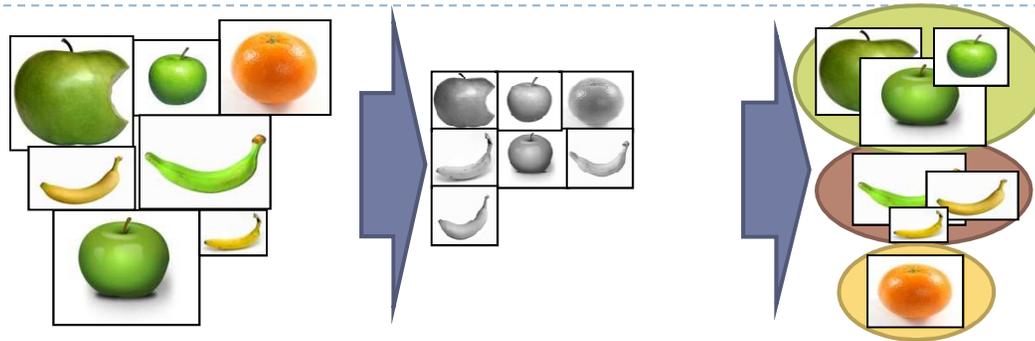
Clustering

Key Collision

- ▶ For each string of letters
 - ▶ Transform each object to eliminate uninformative information
 - ▶ Decide whether the transformed strings group together



Key Collision: Fingerprint



"Godel"
 "Gödel"
 "Johnny 5 "
 "Johnny 5"
 "Johnny-5"
 "It's a mad mad mad mad world"
 "It's a mad, mad, mad, world"
 "Bob Lannon"
 "Lannon, Bob"

TRANSFORM

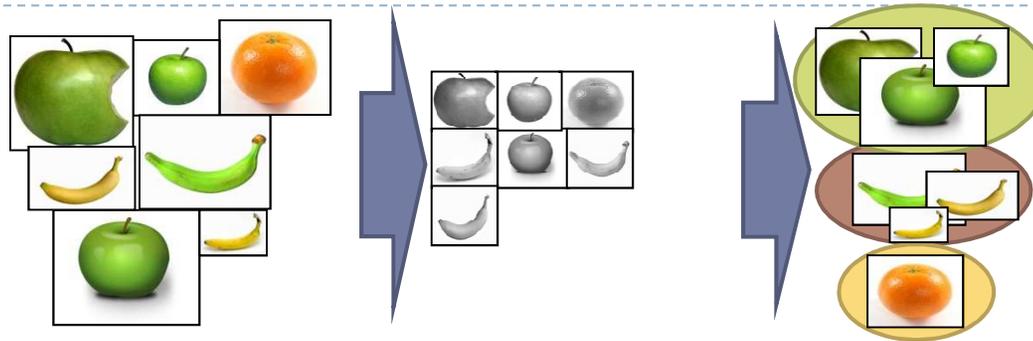
"godel"
 "godel"
 "johnny 5 "
 "johnny 5"
 "johnny 5"
 "johnny 5"
 "a it's mad world"
 "a it's mad world"
 "bob lannon"
 "bob lannon"

CLUSTER

"Godel"
 "Gödel"
 "Johnny 5 "
 "Johnny 5"
 "Johnny-5"
 "It's a mad mad mad
 mad world"
 "It's a mad, mad, mad,
 world"
 "Bob Lannon"
 "Lannon, Bob"



Key Collision: Metaphone



Hypothyroidism

Hupothiroidism

Hyperthyroidism

Huperthoroidism

Throat Cancer

Thyroid cancer



{HP0R, HPTR}

{HP0R, HPTR}

{HPR0, HPRT}

{HPR0, HPRT}

{0RTK, TRTK}

{0RTK,TRTK}



Hypothyroidism

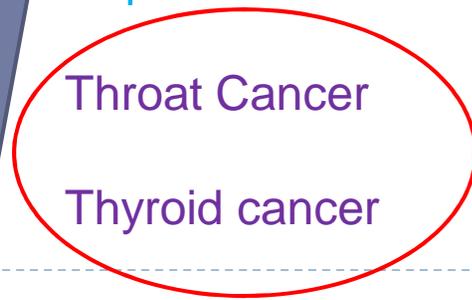
Hupothiroidism

Hyperthyroidism

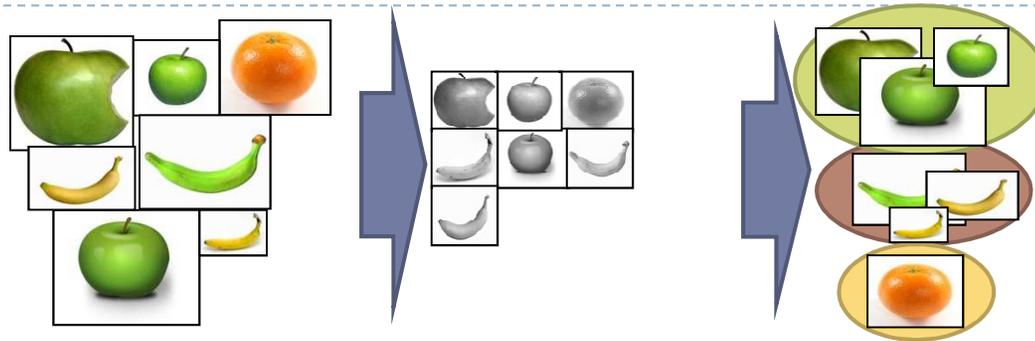
Huperthoroidism

Throat Cancer

Thyroid cancer



Key Collision: N-Gram (size = 1)



Psoriasis

Proriasis

Psoriaasis

Psoriais

psoriiasis

TRANSFORM

aioprs

aioprs

aioprs

aioprs

aioprs

CLUSTER

Psoriasis

Proriasis

Psoriaasis

Psoriais

psoriiasis



Nearest Neighbors (earmuffs)

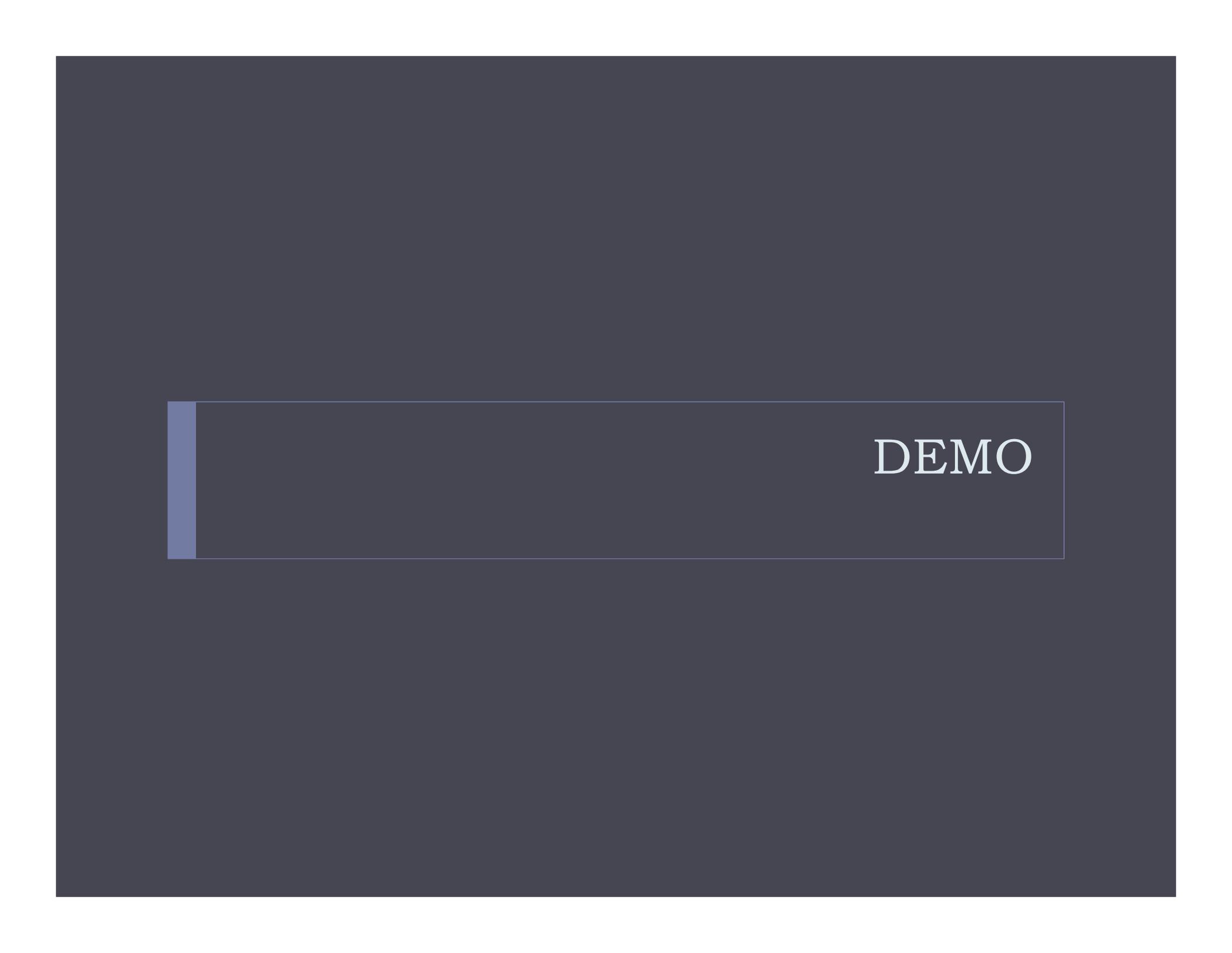
- ▶ **Plotting by attributes**

- ▶ Assign each string to a point in some high-dimensional space
- ▶ look for distances between points

- ▶ **Methods**

- ▶ Levenshtein Distance (like some spell checks)
- ▶ Prediction by Partial Matching (like some better spell checks)

The image shows a software interface for a Nearest Neighbors search. It features four input fields: 'Method' set to 'nearest neighbor', 'Distance Function' set to 'levenshtein', 'Radius' set to '1.0', and 'Block Chars' set to '6'. The 'Radius' field is circled in red, and the 'Block Chars' field is circled in blue. A red box with a line pointing to the 'Radius' field contains the text: 'Increasing gets you more matches, but potentially of lower quality'. A blue box with a line pointing to the 'Block Chars' field contains the text: 'Decreasing gets more matches without sacrificing quality, but will take more time (lower than 3 is a waste of time)'. A dashed horizontal line is present below the interface, with a blue triangle pointing to the right on the left side.



DEMO

Link Party!

- ▶ **Google Refine:** <http://code.google.com/p/google-refine>
 - ▶ How to think: <http://code.google.com/p/google-refine/wiki/UserGuide>
 - ▶ Documentation: <http://code.google.com/p/google-refine/wiki/DocumentationForUsers>
- ▶ **Tutorials:**
 - ▶ Google: http://www.youtube.com/watch?v=B70J_H_zAWM
 - ▶ Open Corporates: <http://vimeo.com/17924204>
 - ▶ ProPublica: <http://www.propublica.org/nerds/item/using-google-refine-for-data-cleaning>
 - ▶ A Bunch More: <http://googlerefine.blogspot.com/p/tutorials-section.html>
- ▶ **More on clustering:** <http://code.google.com/p/google-refine/wiki/ClusteringInDepth>

