

**Train AI 2018** 

Marti Hearst
UC Berkeley

With Philippe Laban

# How Many Distinct People Were Mentioned In Accounts of the Cambridge Analytical Scandal?

**Audience Exercise!** 

# How Many Distinct People Were Mentioned In Accounts of the Cambridge Analytical Scandal?

We find > 700

Why Handle The Long Tail?

How to Handle Rare Cases?



# WHY HANDLE RARE CASES?

## **AVOID DAMAGING ERRORS**

Don't erroneously say someone died!

# Who Is John Allen?

In 2014, John Allen died in a plane downed by a Russian-made missile.

John Allen was also a retired U.S. General and Presidential Envoy in 2014.

Confusing these people would have been problematic!

## John Allens are Plentiful!



Results 1 - 20 of 204

#### John L. Allen, Jr. (Q1254163)

American journalist

22 statements, 10 sitelinks - 20:15, 8 May 2018

#### John R. Allen (Q615356)

United States Marine Corps general

25 statements, 14 sitelinks - 20:15, 8 May 2018

#### John S. Allen (Q6012096)

American university president, professor of astronomy 19 statements, 2 sitelinks - 20:15, 8 May 2018

#### John J. Allen, Jr. (Q1700644)

American politician

35 statements, 2 sitelinks - 20:15, 8 May 2018

#### John P. Allen (Q16197280)

Canadian country/rock/bluegrass fiddler 10 statements, 1 sitelink - 20:15, 8 May 2018

#### John Hensleigh Allen (Q15072629)

politician

23 statements, 1 sitelink - 20:15, 8 May 2018

#### John James Allen (Q198304)

American lawyer and judge

26 statements, 4 sitelinks - 20:15, 8 May 2018

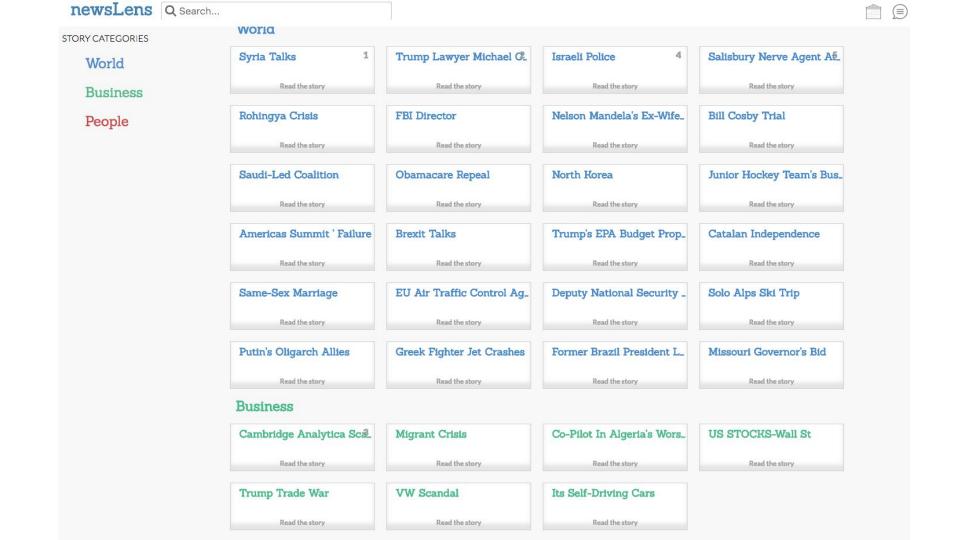
#### John Mills Allen (Q1701104)

American politician

24 statements, 2 sitelinks - 20:15, 8 May 2018

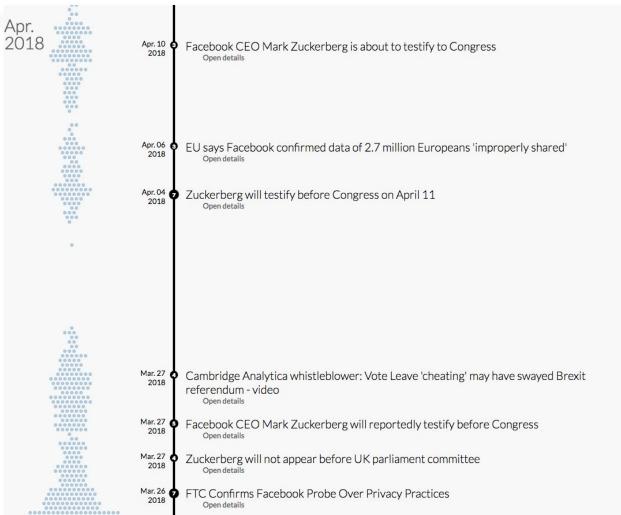
# WHY DO WE CARE?

# newsLens



# Cambridge Analytica Scandal









#### Steve Bannon • 97 mentions

strategist adviser chief strategist white house strategist



#### Brittany Kaiser • 79 mentions

director former director business development director former business development director



#### Darren Grimes • 75 mentions

founder 23-year - old fashion student friend fashion student



#### Nigel Oakes • 73 mentions

founder company parent company executive



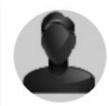
#### Ted Cruz • 71 mentions

sen. senator campaign candidate



#### Arron Banks • 62 mentions

donor founder eu founder campaigner



#### 

Introductions

funder

its

a key financier in the Brexit campaign

who has funded the Leave

British businessman

A prominent campaigner to leave the European Union

EU

the insurance magnate

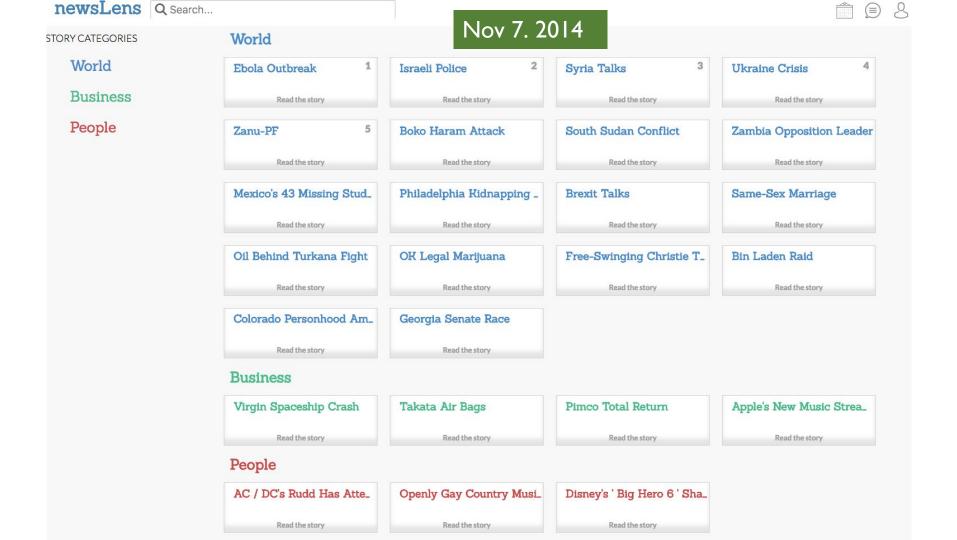
who campaigned for Britain to leave the European Union in a 2016 referendum

EU founder

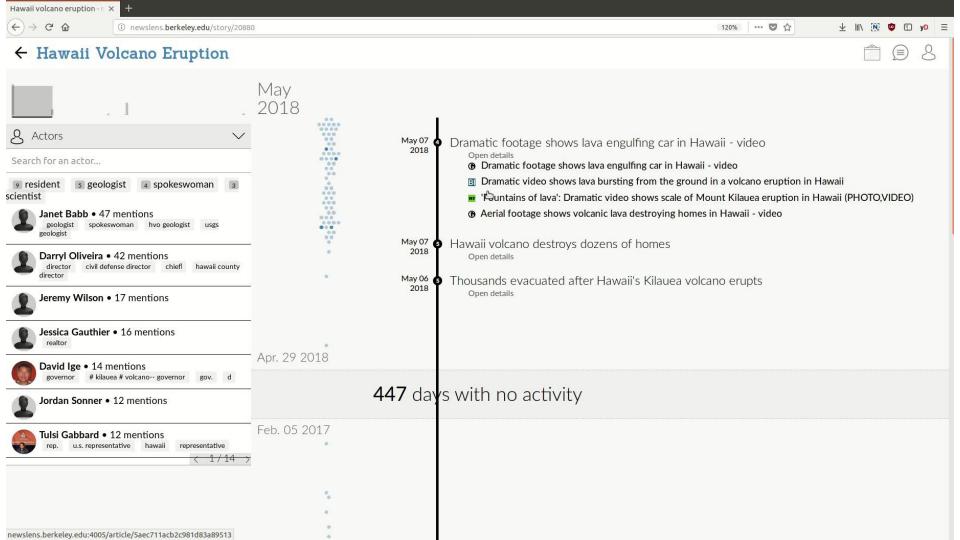
X

## News Lens Shows The Long Tail in News

- Represent disparate articles as stories
- Can go back in time to see past stories
- Shows many points of view at once
- Shows the less famous people in context









Peter Thiel & 1200 mentions in 117 stories

#### Actor

#### Hulk Hogan's Gawker La\_

Read the story

Read the story

His Silicon Valley

Read the story

Read the story

Silicon Valley's Famous S.

Republican Convention

Cambridge Analytica Scal

**FBI** Director

Read the story

Read the story

Read the story

Republican Debate

New Zealand Citizenship

Read the story Donald Trump's Behavior \_

Read the story Silicon Valley's Liberal Bu...

Read the story

Read the story

Republican Debate

Trump Administration

Read the story

Trump Tower Meeting

Google Employee Anti-Di... Read the story

Up To 2 Dozen Parties Ey\_ Read the story

Read the story

Elon Musk's New Company

Read the story

Facebook Shareholder Vote

California Senate Leader \_

Read the story

Read the story

Read the story

Obama Administration

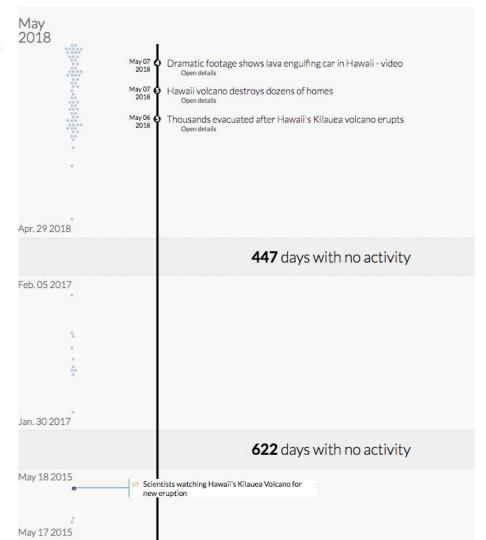
Read the story

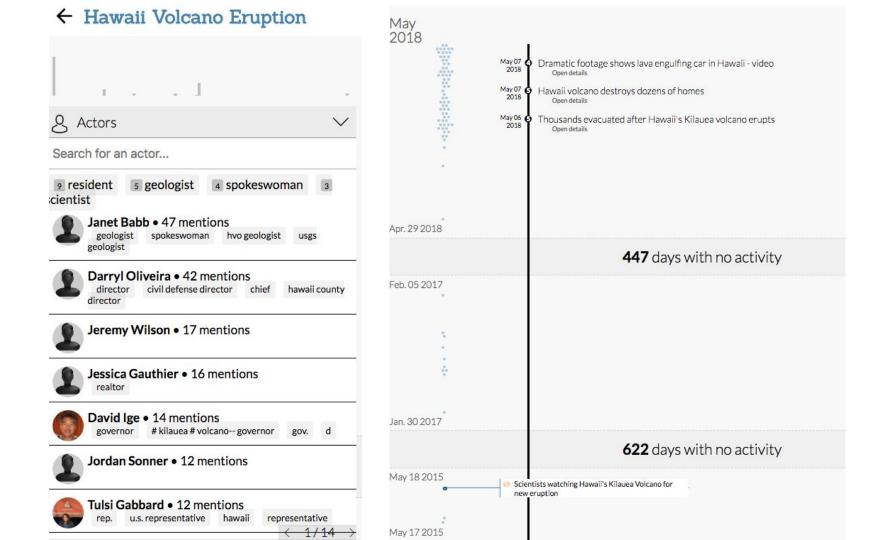
Trump's Immigration Ord...

Clayton's SEC Nomination

Read the story

### Hawaii Volcano Eruption





# Why Handle The Long Tail?

## How to Handle Rare Cases?

Goal: Entity Linking People to Knowledge Base Entries



# Journalists Write Using "Introductions"

A properly sourced news story states some background for anyone who is mentioned.

"Vatican analyst John Allen"

We use an "introduction" as very precise evidence for a person's identity.

# Alternatives / Redundancy Reduce Errors

The first introduction below suggests someone who might not be known on wikidata.

"another victim"

"Oscar winning actress"

"the American actress"

The additional introductions enrich the description.

#### Marine General

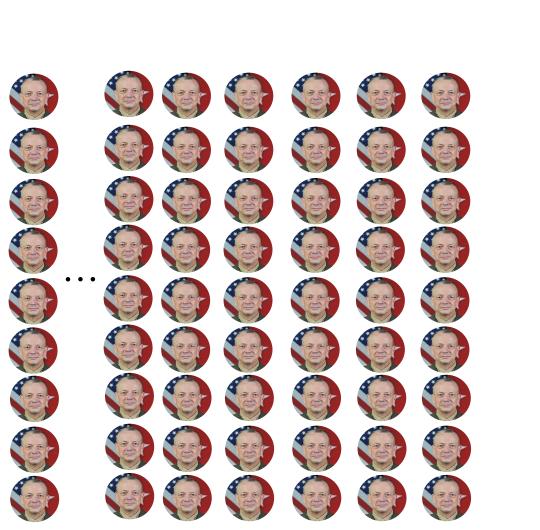


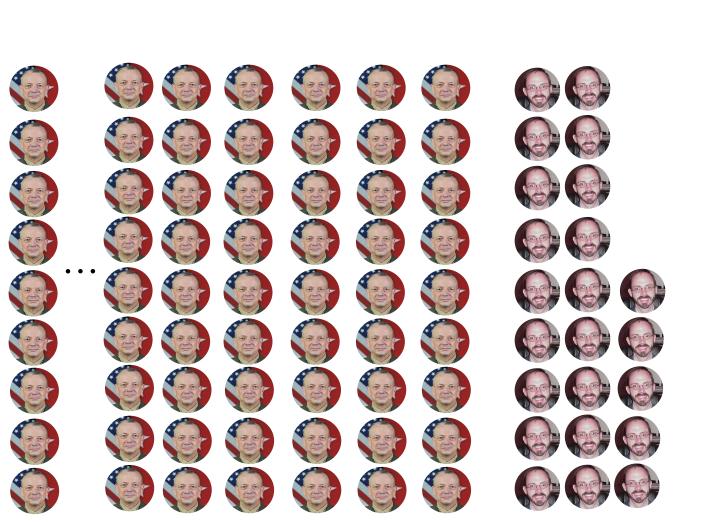
## The top U.S. and NATO commander in Afghanistan Marine Corps Gen.



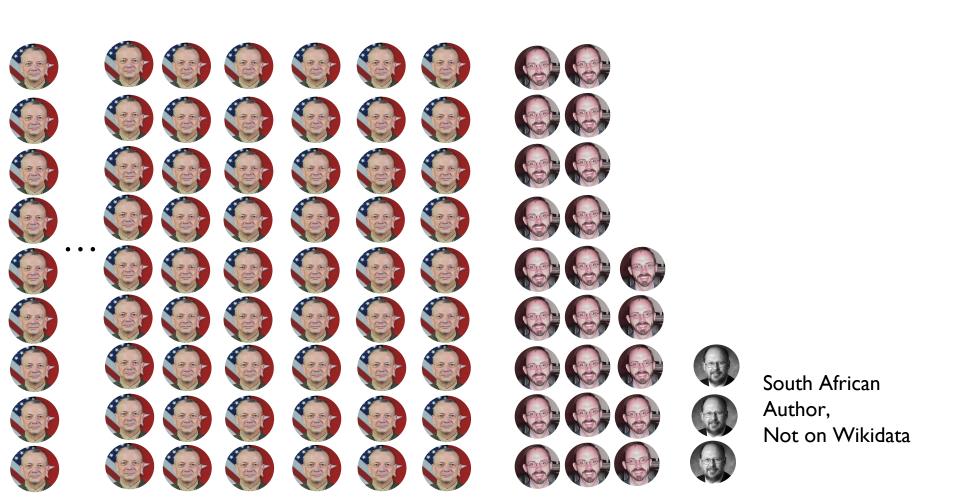
CNN's Vatican analyst

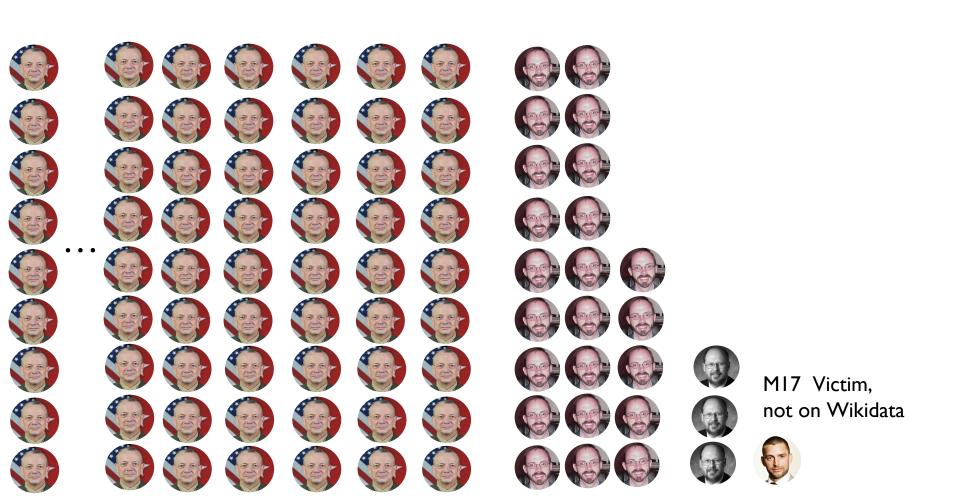


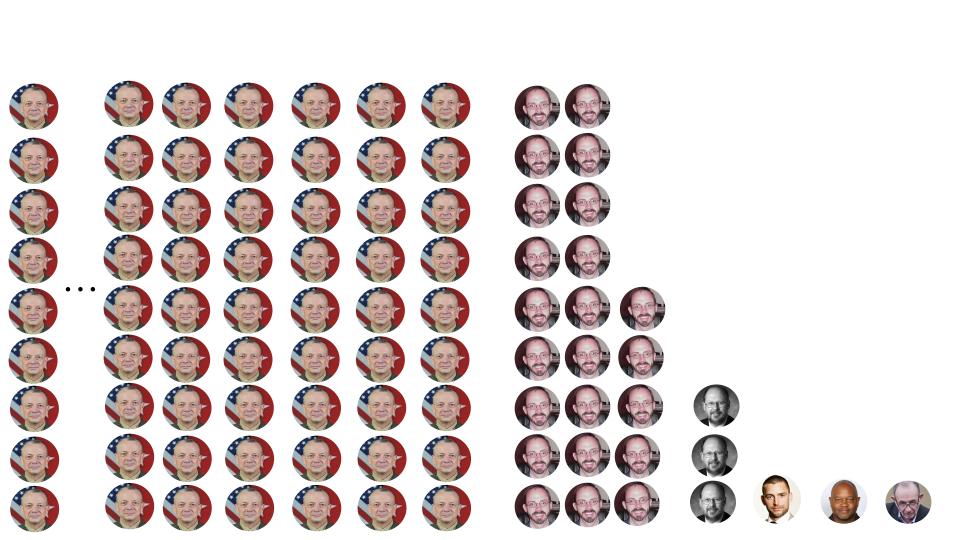


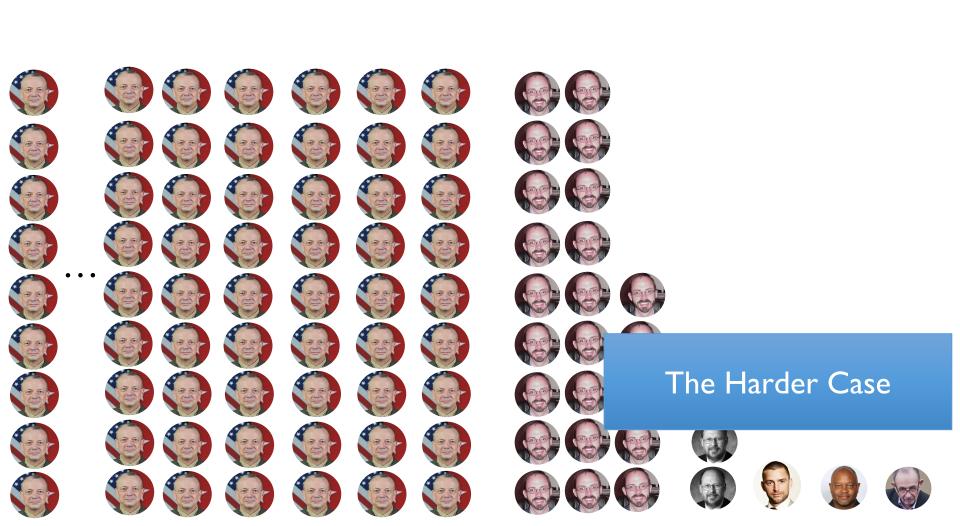


















"our much-loved colleague"

An Amsterdam-based British lawyer

one of the last two of the ten UK victims to be identified

(The Telegraph, July 20, 2017)







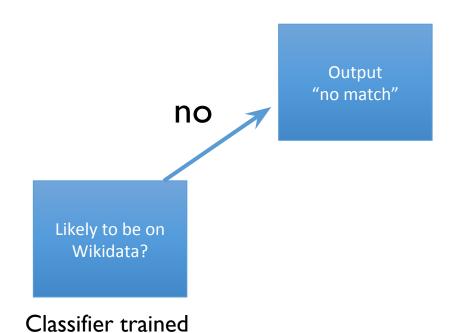
## Not likely to link

"our much-loved colleague"

An Amsterdam-based British lawyer

one of the last two of the ten UK victims to be identified

(The Telegraph, July 20, 2017)



on introductions



"our much-loved colleague"



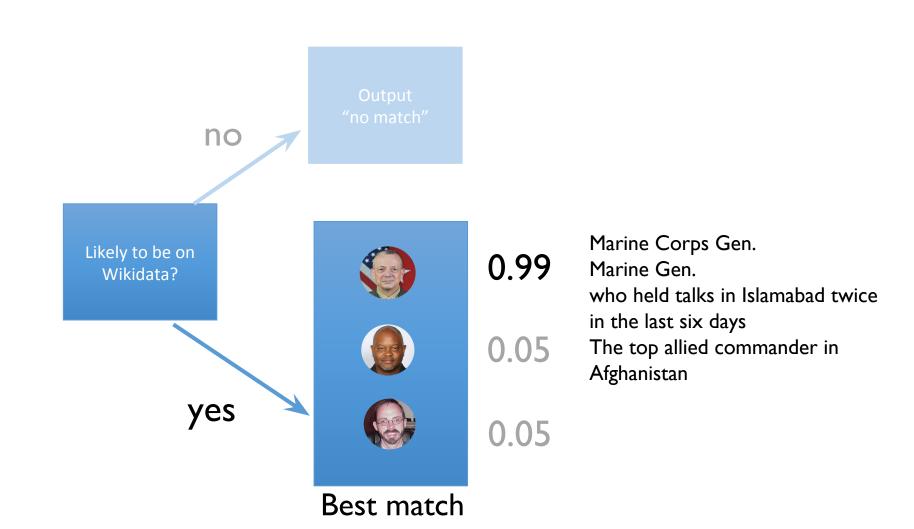


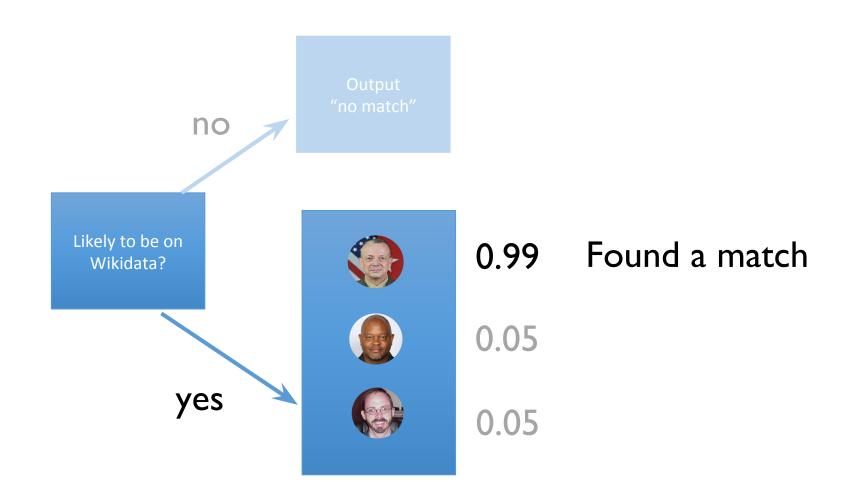
The top U.S. and NATO commander in Afghanistan Marine Corps Gen.

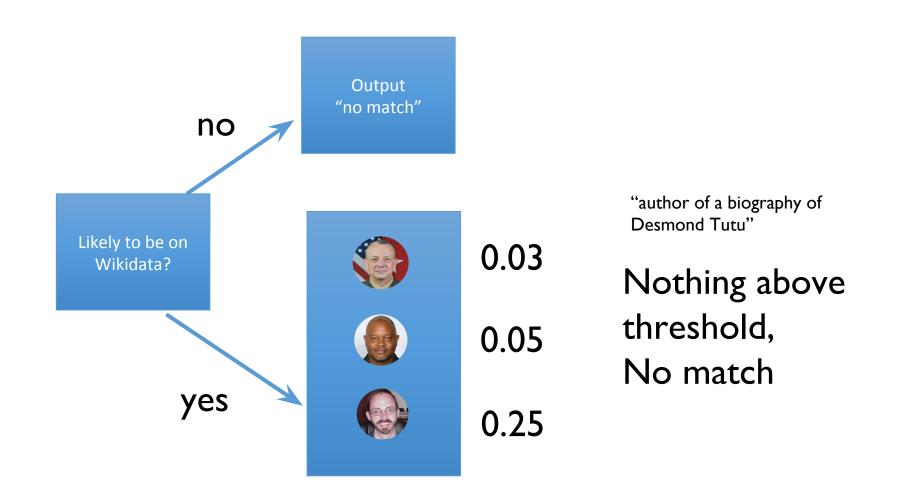
Marine Gen.

who held talks in Islamabad twice in the last six days

The top allied commander in Afghanistan







## A CHALLENGING DATASET

- •Extracted 1000 (article, people) pairs
- •Each person has 2+ potential entries in Wikidata
- •Manually labeled the pairs for ground truth
- •15% (293) are instances that do not match to WD

## **RESULTS**

Easy cases: When there is a match in Wikidata 84.5% correct (91.5 using position in Wikidata) (compared to a popular commercial system of 82.1)

Hard cases: When **no match** in Wikidata **50.8**% correct

(compared to a popular commercial system of **14.5**)

## SUMMARY: How We Handled the Long Tail

### •Identified evidence that is:

- High precision
- Contextually determined
- •Can often label with just one instance, but also can be improved when additional evidence is present.

## •Two-phase classification:

- •First decide yes-or-no for match to KB
- •This greatly improves accuracy on no-match cases.

