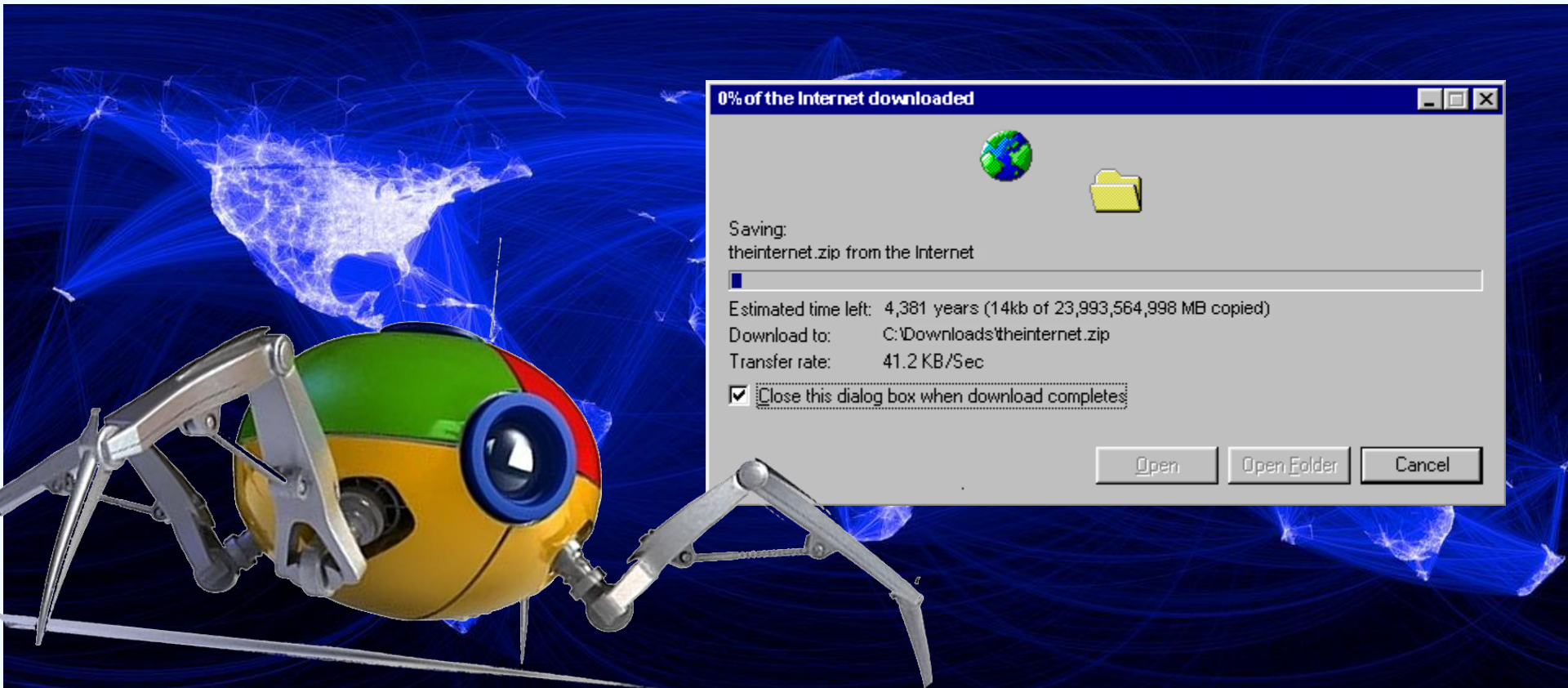# CC5212-1
## PROCESAMIENTO MASIVO DE DATOS
## OTOÑO 2017

## Lecture 7: Information Retrieval II

Aidan Hogan
aidhog@gmail.com

# How does Google know about the Web?

# Inverted Index: Example

en.wikipedia.org/wiki/Fruitvale_Station

## Fruitvale Station

From Wikipedia, the free encyclopedia

***Fruitvale Station*** is a 2013 American [drama film](#) written and directed by [Ryan Coogler](#).

Inverted index:

| Term List | Posting List |
|-----------|--------------|
| a | (1,2,…) |
| american | (1,5,…) |
| and | (1,2,…) |
| by | (1,2,…) |
| directed | (1,2,…) |
| drama | (1,16,…) |
| … | … |

# Apache Lucene

- Inverted Index
  - They built one so you don't have to!
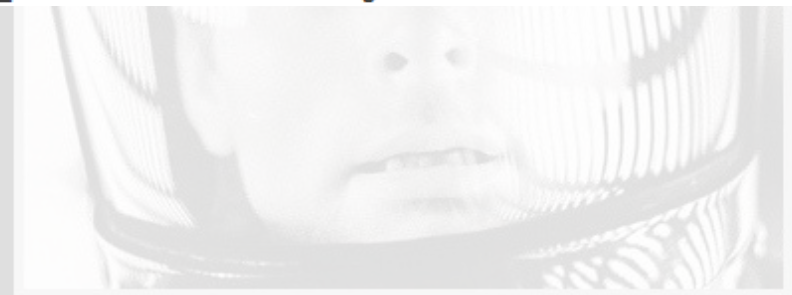  - Open Source in Java



My God. It's full of win.
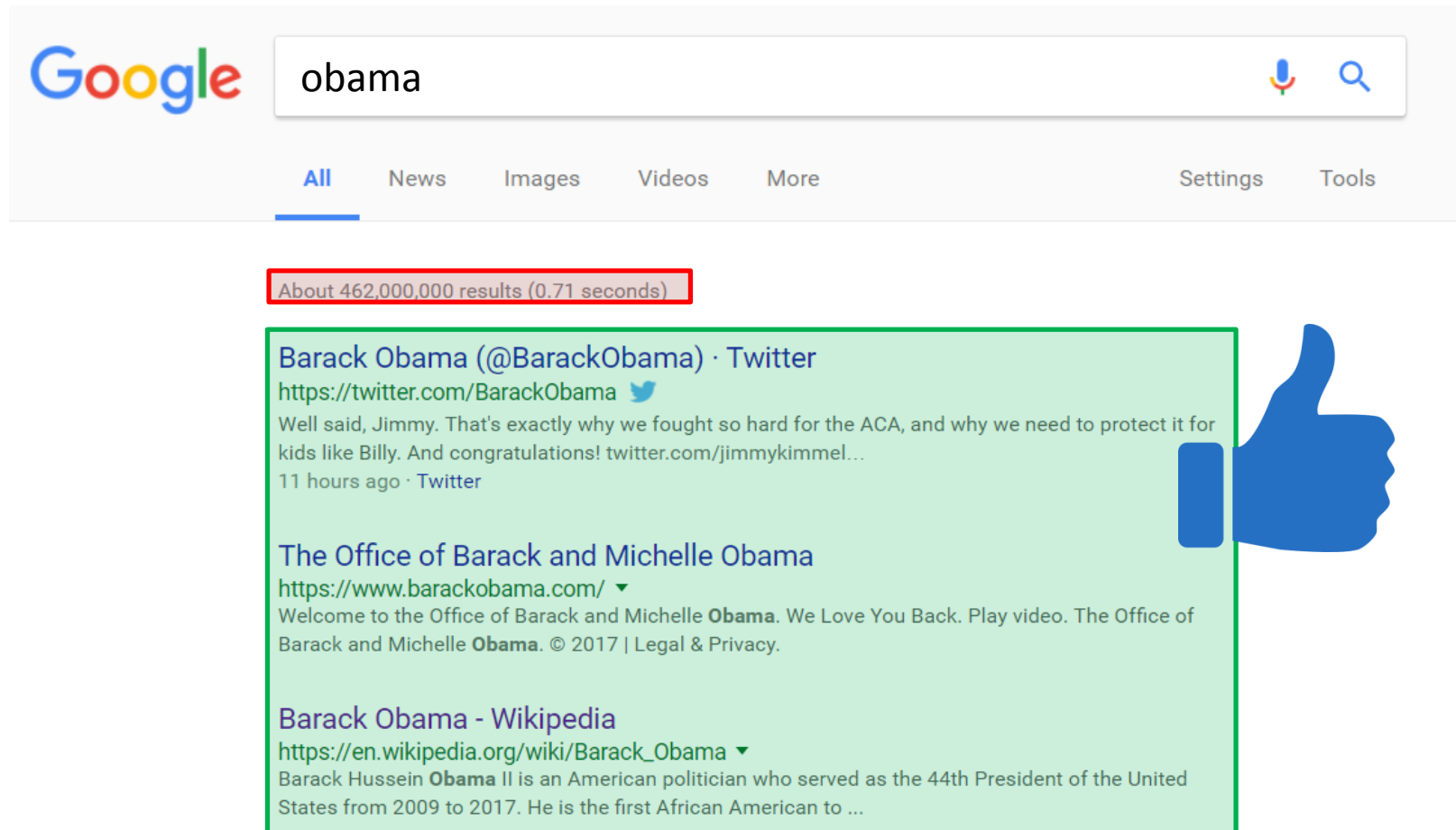
# Apache Lucene



```
Tasks    Console ✕
SearchWikiIndex [Java Application] C:\Program Files\Java\jre1.8.0_65\bin\javaw.exe (03-05-2017 12:45:22 a. m.)
Opening directory at  lucene
Enter a keyword search phrase:
obama
Running query: obama
Parsed query: TITLE:obam^5.0 ABSTRACT:obam
Matching documents: 255
Showing top 10 results
1       http://es.wikipedia.org/wiki/Obama_Republican    Obama Republican
2       http://es.wikipedia.org/wiki/Obama_(Fukui)       Obama (Fukui)
3       http://es.wikipedia.org/wiki/Republicanos_por_Obama      Republicanos por Obama
4       http://es.wikipedia.org/wiki/Engonga_Obame       Engonga Obame
5       http://es.wikipedia.org/wiki/Barack_Obama        Barack Obama
6       http://es.wikipedia.org/wiki/Michelle_Obama      Michelle Obama
7       http://es.wikipedia.org/wiki/Cartel_%22Hope%22_de_Obama Cartel "Hope" de Obama
8       http://es.wikipedia.org/wiki/Transición_presidencial_de_Barack_Obama     Transición presidencial de Barack Obama
9       http://es.wikipedia.org/wiki/Por_qué_Obama_ganará_en_2008_y_en_2012      Por qué Obama ganará en 2008 y en 2012
10      http://es.wikipedia.org/wiki/Ricardo_Mangue_Obama_Nfubea        Ricardo Mangue Obama Nfubea
```





My God. It's full of win.

# INFORMATION RETRIEVAL: RANKING

# How Does Google Get Such Good Results?

Google

**obama**

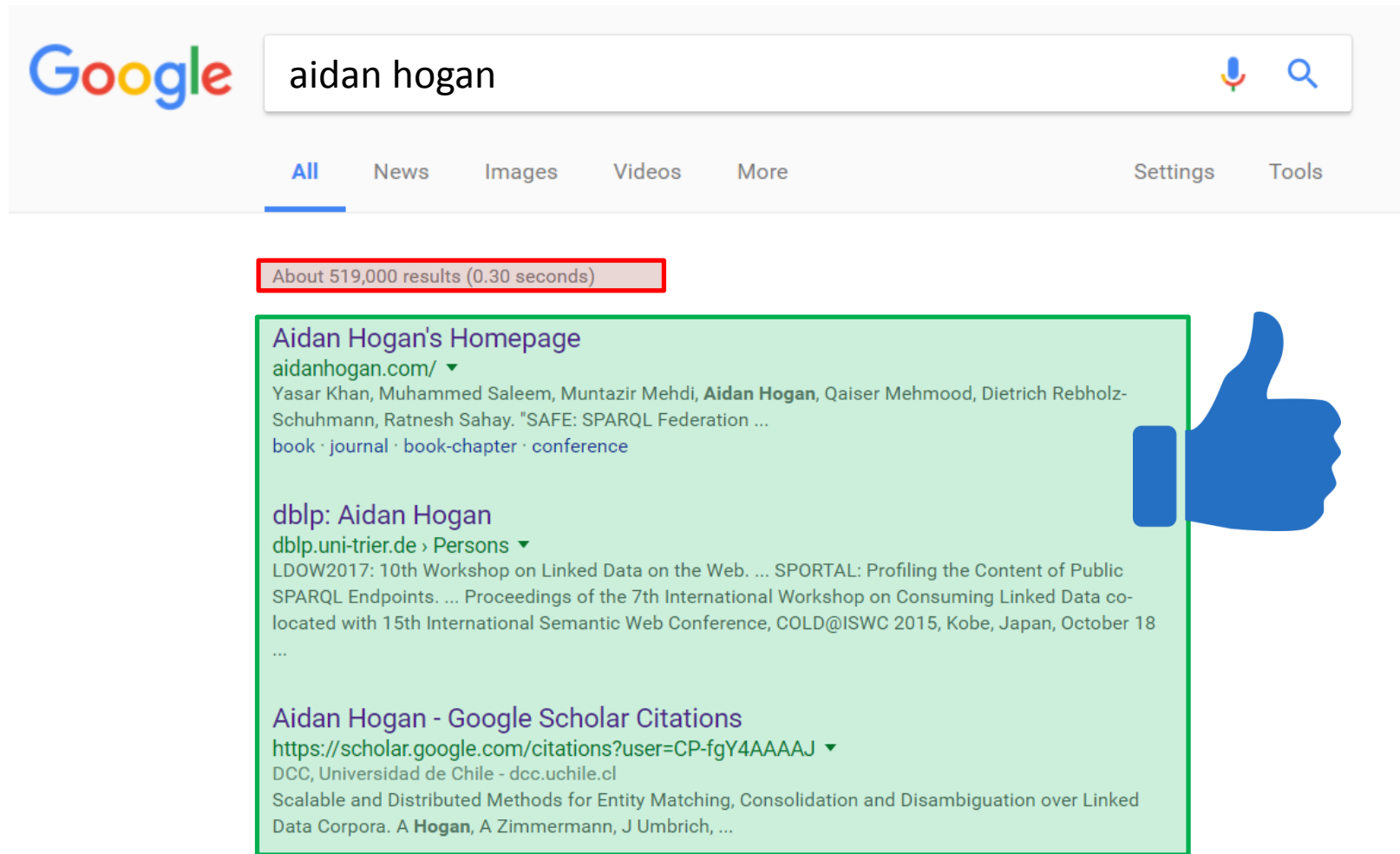All    News    Images    Videos    More        Settings    Tools

About 462,000,000 results (0.71 seconds)

### Barack Obama (@BarackObama) · Twitter
https://twitter.com/BarackObama
Well said, Jimmy. That's exactly why we fought so hard for the ACA, and why we need to protect it for kids like Billy. And congratulations! twitter.com/jimmykimmel...
11 hours ago · Twitter

### The Office of Barack and Michelle Obama
https://www.barackobama.com/ ▾
Welcome to the Office of Barack and Michelle **Obama**. We Love You Back. Play video. The Office of Barack and Michelle **Obama**. © 2017 | Legal & Privacy.

### Barack Obama - Wikipedia
https://en.wikipedia.org/wiki/Barack_Obama ▾
Barack Hussein **Obama** II is an American politician who served as the 44th President of the United States from 2009 to 2017. He is the first African American to ...

# How Does Google Get Such Good Results?

# How does Google Get Such Good Results?

# Two Sides to Ranking: Relevance

# Two Sides to Ranking: Importance

# RANKING: RELEVANCE

# Example Query

# Matches in a Document

# Matches in a Document



Browser window showing the Wikipedia article "Braveheart" with a find bar reading "movie  3 de 16".

**freedom**
- 7 occurrences

**movie**
- 16 occurrences

# Matches in a Document

# Usefulness of Words



**movie**
- occurs very frequently

**freedom**
- occurs frequently

**wallace**
- occurs occassionally

# Estimating Relevance

- Rare words more important than common words
  - wallace (49M) more important than freedom (198M) more important than movie (835M)


- Words occurring more frequently in a document indicate higher relevance
  - wallace (88) more matches than movie (16) more matches than freedom (7)

# Relevance Measure: TF–IDF

- ## TF: Term Frequency
  - Measures occurrences of a term in a document
  - $\text{tf}(t, d)$ ... various options
    - Raw count of occurrences
      $$\text{tf}(t, d) = \text{count}(t, d)$$
    - Logarithmically scaled
      $$\text{tf}(t, d) = \log(\text{count}(t, d) + 1)$$
    - Normalised by document length
      $$\text{tf}(t, d) = \frac{\text{count}(t,d)}{\sum_{t' \in d} \text{count}(t',d)}$$
      $$\text{tf}(t, d) = \frac{\text{count}(t,d)}{\max_{t' \in d} \text{count}(t',d)}$$
    - A combination / something else ☺

# Relevance Measure: TF–IDF

- **IDF**: **I**nverse **D**ocument **F**requency
  - How common a term is across **all** documents
  - $\mathrm{idf}(t, D)$ ...
    - Logarithmically scaled document occurrences

$$\mathrm{idf}(t, D) = \log\left(\frac{|D|}{|\{d \in D : t \in d\}| + 1}\right)$$

    - Note: The more rare, the larger the value

# Relevance Measure: TF–IDF

- TF–IDF: Combine Term Frequency and Inverse Document Frequency:

$$\text{tf-idf}(t, d) = \text{tf}(t, d) \times \text{idf}(t, D)$$

- Score for a query
  - Let query $q = (t_1, \ldots, t_n)$
  - Score for a query: $score(q, d) = \sum_{t \in q} \text{tf-idf}(t, d)$

  (There are other possibilities)

# Relevance Measure: TF–IDF

Google    movie freedom wallace

WIKIPEDIA
The Free Encyclopedia

Article  Talk          Read  Edit  View history

*Braveheart*

From Wikipedia, the free encyclopedia

## Term Frequency

$$\text{tf}(t,d) = \text{count}(t,d)$$

## Inverse Document Frequency

$$\text{idf}(t,D) = \log_2\left(\frac{|D|}{|\{d \in D : t \in d\}| + 1}\right)$$

$$\text{tf-idf}(t,d) = \boxed{\text{tf}(t,d)} \times \boxed{\text{idf}(t,D)}$$

| $t$ | $\text{tf}(t,d)$ |
|---|---|
| movie | 16 |
| freedom | 7 |
| wallace | 43 |

# Relevance Measure: TF–IDF



## Term Frequency

$$\text{tf}(t,d) = \text{count}(t,d)$$

## Inverse Document Frequency

$$\text{idf}(t,D) = \log_2\left(\frac{|D|}{|\{d \in D : t \in d\}| + 1}\right)$$

$$\text{tf-idf}(t,d) = \boxed{\text{tf}(t,d)} \times \boxed{\text{idf}(t,D)}$$

| $t$ | $\text{tf}(t,d)$ | $|\{d \in D : t \in d\}|$ |
|---|---|---|
| movie | 16 | 835,000,000 |
| freedom | 7 | 198,000,000 |
| wallace | 43 | 49,200,000 |

# Relevance Measure: TF–IDF

Google | movie freedom wallace

WIKIPEDIA
The Free Encyclopedia

Article  Talk      Read  Edit  View history

*Braveheart*

From Wikipedia, the free encyclopedia

## Term Frequency

$$\mathrm{tf}(t, d) = \mathrm{count}(t, d)$$

## Inverse Document Frequency

$$\mathrm{idf}(t, D) = \log_2\left(\frac{|D|}{|\{d \in D : t \in d\}| + 1}\right)$$

$$\text{tf-idf}(t, d) = \boxed{\mathrm{tf}(t, d)} \times \boxed{\mathrm{idf}(t, D)}$$

| $t$ | $\mathrm{tf}(t, d)$ | $\lvert\{d \in D : t \in d\}\rvert$ | $\frac{\lvert D\rvert}{\lvert\{d\in D\,:\,t\in d\}\rvert + 1}$ |
|---|---|---|---|
| movie | 16 | 835,000,000 | |
| freedom | 7 | 198,000,000 | |
| wallace | 43 | 49,200,000 | |

Google | the

Web   Images   News   Books   More ▾   Search tools

About 11,410,000,000 results (0.27 seconds)

$$|D| = 11,410,000,000$$

# Relevance Measure: TF–IDF

Google    movie freedom wallace

WIKIPEDIA
The Free Encyclopedia

Article  Talk          Read  Edit  View history

*Braveheart*

From Wikipedia, the free encyclopedia

## Term Frequency

$$\mathrm{tf}(t, d) = \mathrm{count}(t, d)$$

## Inverse Document Frequency

$$\mathrm{idf}(t, D) = \log_2\left(\frac{|D|}{|\{d \in D : t \in d\}| + 1}\right)$$

$$\mathrm{tf\text{-}idf}(t, d) = \boxed{\mathrm{tf}(t, d)} \times \boxed{\mathrm{idf}(t, D)}$$

| $t$ | $\mathrm{tf}(t, d)$ | $\lvert\{d \in D : t \in d\}\rvert$ | $\frac{\lvert D \rvert}{\lvert\{d \in D : t \in d\}\rvert + 1}$ |
|---|---|---|---|
| movie | 16 | 835,000,000 | 13.66 |
| freedom | 7 | 198,000,000 | 57.62 |
| wallace | 43 | 49,200,000 | 231.91 |

Google    the                                                    🎤    🔍

Web    Images    News    Books    More ▾    Search tools

About 11,410,000,000 results (0.27 seconds)

$$|D| = 11,410,000,000$$

# Relevance Measure: TF–IDF

## Term Frequency

$$\mathrm{tf}(t, d) = \mathrm{count}(t, d)$$

## Inverse Document Frequency

$$\mathrm{idf}(t, D) = \log_2 \left( \frac{|D|}{|\{d \in D : t \in d\}| + 1} \right)$$

$$\mathrm{tf\text{-}idf}(t, d) = \boxed{\mathrm{tf}(t, d)} \times \boxed{\mathrm{idf}(t, D)}$$

| $t$ | $\mathrm{tf}(t, d)$ | $|\{d \in D : t \in d\}|$ | $\frac{|D|}{|\{d \in D : t \in d\}| + 1}$ | $\mathrm{idf}(t, D)$ |
|---|---|---|---|---|
| movie | 16 | 835,000,000 | 13.66 | 3.77 |
| freedom | 7 | 198,000,000 | 57.62 | 5.84 |
| wallace | 43 | 49,200,000 | 231.91 | 7.85 |

# Relevance Measure: TF–IDF



## Term Frequency

$$\mathrm{tf}(t,d) = \mathrm{count}(t,d)$$

## Inverse Document Frequency

$$\mathrm{idf}(t,D) = \log_2\left(\frac{|D|}{|\{d \in D : t \in d\}| + 1}\right)$$

$$\text{tf-idf}(t,d) = \boxed{\mathrm{tf}(t,d)} \times \boxed{\mathrm{idf}(t,D)}$$

| $t$ | $\mathrm{tf}(t,d)$ | $\|\{d \in D : t \in d\}\|$ | $\dfrac{|D|}{|\{d\in D\,:\,t\in d\}|+1}$ | $\mathrm{idf}(t,D)$ | $\text{tf-idf}(t,d)$ |
|---|---|---|---|---|---|
| movie | 16 | 835,000,000 | 13.66 | 3.77 | 60.36 |
| freedom | 7 | 198,000,000 | 57.62 | 5.84 | 40.94 |
| wallace | 43 | 49,200,000 | 231.91 | 7.85 | 337.87 |

# Relevance Measure: TF–IDF

## Term Frequency

$$\text{tf}(t, d) = \text{count}(t, d)$$

## Inverse Document Frequency

$$\text{idf}(t, D) = \log_2 \left( \frac{|D|}{|\{d \in D : t \in d\}| + 1} \right)$$

$$\text{tf-idf}(t, d) = \boxed{\text{tf}(t, d)} \times \boxed{\text{idf}(t, D)}$$

| $t$ | $\text{tf}(t, d)$ | $|\{d \in D : t \in d\}|$ | $\frac{|D|}{|\{d \in D : t \in d\}| + 1}$ | $\text{idf}(t, D)$ | $\text{tf-idf}(t, d)$ |
|---|---|---|---|---|---|
| movie | 16 | 835,000,000 | 13.66 | 3.77 | 60.36 |
| freedom | 7 | 198,000,000 | 57.62 | 5.84 | 40.94 |
| wallace | 43 | 49,200,000 | 231.91 | 7.85 | 337.87 |

$$\text{score}(q, d) = \sum_{t \in q} \text{tf-idf}(t, d)$$

$$\text{score}\big((\text{movie, freedom, wallace}), \text{http://en.wikipedia.org/Braveheart}\big) \approx 439.17$$

# Vector Space Model (a mention)

| $t$ | $\mathrm{tf}(t, d)$ |
|---|---|
| movie | 16 |
| freedom | 7 |
| wallace | 43 |

$$l = \sqrt{\sum_{t \in q} \mathrm{tf}(t, d)^2}$$

# Vector Space Model (a mention)

| $t$ | $\mathrm{tf}(t, d)$ | $\mathrm{tf}(t, d)^2$ |
|---|---|---|
| movie | 16 | 256 |
| freedom | 7 | 49 |
| wallace | 43 | 1,894 |

$$l = \sqrt{\sum_{t \in q} \mathrm{tf}(t, d)^2}$$

# Vector Space Model (a mention)

| $t$ | $\mathrm{tf}(t,d)$ | $\mathrm{tf}(t,d)^2$ | $\frac{\mathrm{tf}(t,d)}{l}$ |
|---|---|---|---|
| movie | 16 | 256 | 0.34 |
| freedom | 7 | 49 | 0.15 |
| wallace | 43 | 1,894 | 0.92 |

$$l = \sqrt{\sum_{t \in q} \mathrm{tf}(t,d)^2}$$

$$\vec{v}(d) = (0.34, 0.15, 0.92)$$



Dividing by $l$ normalises length of vector to 1

# Vector Space Model (a mention)

- ## Cosine Similarity

$$\text{sim}(d, d') = \vec{v}(d) \cdot \vec{v}(d')$$

| $t$ | $\vec{v}(d)$ | $\vec{v}(d')$ | $\times$ |
|---------|------|------|------|
| movie | 0.34 | 0.49 | 0.17 |
| freedom | 0.15 | 0.82 | 0.12 |
| wallace | 0.93 | 0.30 | 0.28 |

$\Sigma$

$$\text{sim}(d, d') \approx \boxed{0.57}$$

- ## Note:

$$\mathbf{a} \cdot \mathbf{b} = |\mathbf{a}|\,|\mathbf{b}|\cos(\angle(\mathbf{a}, \mathbf{b}))$$

$$|\vec{v}(d)| = |\vec{v}(d')| = 1$$

$\vec{v}(d) = (0.34, 0.15, 0.92)$

$\vec{v}(d') = (0.49, 0.82, 0.30)$

Hence the similarity is the cosine of the **angle** between the vectors

# Relevance Measure: TF–IDF

- TF–IDF: Combine Term Frequency and Inverse Document Frequency:

$$\text{tf-idf}(t, d) = \text{tf}(t, d) \times \text{idf}(t, D)$$

- Score for a query
  - Let query $q = (t_1, \ldots, t_n)$
  - Score for a query: $score(q, d) = \sum_{t \in q} \text{tf-idf}(t, d)$

  (There are other possibilities)

  ... we could also use cosine similarity between query and document using TF–IDF weights

# Two Sides to Ranking: Relevance

# Field-Based Boosting

- Not all text is equal: titles, headers, etc.

# Anchor Text

- See how the Web views/tags a page

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN"
 "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">
<html>
<head>
  <title>What I watched last night ...</title>
</head>
<body>
<p>Last night I was pretty bored so I made some popcorn and watched
<a href="http://en.wikipedia.org/Braveheart">a movie about William Wallace called Braveheart</a>.
Set in Scotland it has lots of blood and gore.
</p>
</body>
</html>
```

# Anchor Text

- See how the Web views/tags a page

# Lucene uses relevance scoring



```
Tasks  Console ✕

SearchWikiIndex [Java Application] C:\Program Files\Java\jre1.8.0_65\bin\javaw.exe (03-05-2017 12:45:22 a. m.)
Opening directory at  lucene
Enter a keyword search phrase:
obama
Running query: obama
Parsed query: TITLE:obam^5.0 ABSTRACT:obam
Matching documents: 255
Showing top 10 results
1       http://es.wikipedia.org/wiki/Obama_Republican    Obama Republican
2       http://es.wikipedia.org/wiki/Obama_(Fukui)       Obama (Fukui)
3       http://es.wikipedia.org/wiki/Republicanos_por_Obama      Republicanos por Obama
4       http://es.wikipedia.org/wiki/Engonga_Obame       Engonga Obame
5       http://es.wikipedia.org/wiki/Barack_Obama        Barack Obama
6       http://es.wikipedia.org/wiki/Michelle_Obama      Michelle Obama
7       http://es.wikipedia.org/wiki/Cartel_%22Hope%22_de_Obama Cartel "Hope" de Obama
8       http://es.wikipedia.org/wiki/Transición_presidencial_de_Barack_Obama    Transición presidencial de Barack Obama
9       http://es.wikipedia.org/wiki/Por_qué_Obama_ganará_en_2008_y_en_2012     Por qué Obama ganará en 2008 y en 2012
10      http://es.wikipedia.org/wiki/Ricardo_Mangue_Obama_Nfubea        Ricardo Mangue Obama Nfubea
```



My God. It's full of win.

# RANKING: IMPORTANCE

# Two Sides to Ranking: Importance

Google  obama

Web    Images    News    Videos    More ▾    Search tools

About 48,100,000 results (0.26 seconds)

How could we determine that Barack Obama is more important ⑦
than Mount Obama as a search result on the Web?

Images for **mount obama**                                    Report images

re images for **mount ob**    >

**Mount Obama** Natio... ...rk | Antigua a...
antigua**mountobama** com/
Jun 16, 2011 – As the **Mount Obama** Committee contin...
...Area, the committee organized a site visit to the ...

# Link Analysis

Which will have more links from other pages?
The Wikipedia article for Mount Obama?
The Wikipedia article for Barack Obama?

# Link Analysis

- Consider links as votes of confidence in a page
- A hyperlink is the open Web's version of …



(… even if the page is linked in a negative way.)

# Link Analysis

So if we just count links to a page we can determine its importance and we are done? ⊘

# Link Spamming

semanticweb.com™

The Voice of Semantic Technology Busine
Big Data, Linked Data, Smart Data

Home | Events | Media | Industry Verticals | Answers | Jo

Questions   Tags   Users   Badges

○ Qu

## [deleted] Kala Jadu Specialist +91961

-1

black magic specialist baba ji call now +919610897260

http://www.blackmagicspecialist.net.in

java

edit | close | undelete | more ▼

Claritin Clomid Combivent Confido Copegus Cordarone Coreg Coumadin Cozaar Crestor Cyklokapron Cymbalta Cystone Cytotec Danazol Deltasone Depakote Desyrel Detrol Diabecon Diakof Diarex Didronel Differin Dilantin Diovan Dostinex Elavil Elimite Emsam Endep Eurax Evecare Evista Exelon Famvir Feldene Femara Femcare Flomax Flonase Flovent Fosamax Gasex Geodon Geriforte Herbolax High Love Himcocid Himcolin Himcospaz Himplasia Hoodia Hytrin Hyzaar Imdur Imitrex Inderal Ismo Isoptin Isordil Kamagra Karela Keftab Koflet Kytril Lamictal Lamisil Lanoxin Lariam Lasix Lasuna Leukeran Levaquin Levlen Levothroid Lincocin Lioresal Lisinopril Liv.52 Lopid Lopressor Loprox Lotensin Lotrisone Loxitane Lozol Lukol Lynoral Maxaquin Menosan Mentat Mentax Mevacor Mexitil Miacalcin Micardis Mobic Monoket Motrin Myambutol Mycelex-G Mysoline Naprosyn Neurontin Nicotinell Nimotop Nirdosh Nizoral Nolvadex Nonoxinol Noroxin Omnicef Ophthacare Oxytrol Pamelor Parlodel Paxil Penisole Phentrimine Pilex Plan B Plavix Plendil Pletal Prandin Pravachol Prednisone Premarin Prevacid Prilosec Prinivil Procardia Prograf Prometrium Propecia Proscar Protonix Proventil Prozac Purim Purinethol Quibron-T Relafen Renalka Reosto Requip Retin-A Revia Rhinocort Rimonabant Risperdal Rocaltrol Rogaine Rumalaya Sarafem Septilin Serevent Serophene Seroquel Shallaki Shoot Sinequan Singulair Snoroff Sorbitrate Speman Starlix StretchNil Stromectol Styplon Sumycin Superman Sustiva Synthroid Tenormin Topamax Trandate Tricor Trimox Triphala Tulasi Urispas V-Gel Vantin Vasodilan Vasotec Ventolin Viramune Vytorin Xeloda Xenacore Zanaflex Zantac Zebeta Zelnorm Zerit Yerba Diet Wellbutrin SR Women Attracting Pheromones Women's Intimacy Enhancer Women's Intimacy Enhancer Cream Virility Gum Vitamin A & D Viagra + Cialis Viagra + Cialis + Levitra Viagra Jelly Viagra Soft + Cialis Soft Viagra Soft Tabs Ultimate Male Enhancer Toprol XL Touch-Up Kit Tentex Royal Tentex Forte Tiberius Erectus Zero Nicotine 2 Complete Professional Whitening Kits 2 Sets Of Moldable Mouth Trays 36 Beauty Acne-n-Pimple Cream ActoPlus Met Superloss Multi SleepWell (Herbal XANAX) Shuddha Guggulu Rythmol SR Rumalaya Forte Pulmicort Inhaler Professional Plasma Tooth Whitening Kit Premium Diet Patch Penis Growth Oil Penis Growth Pack Penis Growth Patch Penis Growth Pills Orgasm Enhancer Norpace CR Mental Booster Men Attracting Pheromones Menopause Gum Male Enhancement Oil Male Enhancement Patch Male Enhancement Pills Male Sexual Tonic InnoPran XL Hoodia Weght Loss Gum Hoodia Weight Loss Patch Human Growth Hormone Agent Glucotrol XL Green Tea Grifulvin V Gyne-Lotrimin Hair Loss Cream Herbal Maxx Herbal Phentermine Flagyl ER Female Sexual Tonic Female Viagra Epivir-HBV Diet Maxx Deluxe Handheld Plasma Whitening Tool Deluxe Whitening System With Plasma Lamp Coral Calcium Cialis Jelly Cialis Soft Tabs Calcium Carbonate Bust Enhancer Beconase AQ Anatrim Diet Pills Advair Diskus Advanced Gain Pro Breast Augmentation Breast Enhancement Breast Enhancement Gel Breast Enhancement Gum Breast Intense Buy Trazodone Buy Celebrex Buy Alprazolam Buy Tramadol Buy Fioricet Buy Soma Buy Cialis Buy Carisoprodol Buy Levitra Buy Ultram Buy Ambien Buy Viagra Buy Xanax Buy Phentermine Buy Valium Buy Diazepam Generic Celebrex Generic Alprazolam Generic Tramadol Generic Fioricet Generic Soma Generic Cialis Generic Carisoprodol Generic Levitra Generic Ultram Generic Ambien

# Link Importance

So which should count for more?
A link from http://en.wikipedia.org/wiki/Barack_Obama?
Or a link from http://freev1agra.com/shop.html?

# Link Importance

Maybe we could consider links from some
domains as having more "vote"?

# PageRank

# PageRank

- Not just a count of inlinks

  – A link from a more important page is more important

  – A link from a page with fewer links is more important

  ∴ A page with lots of inlinks from important pages (which have few outlinks) is more important

# PageRank is Recursive

- ## Not just a count of inlinks

  - A link from a <u>more important</u> page is <u>more important</u>

  - A link from a pag

  ∴ A page with lots
      have few out

# PageRank Model

- The Web: a directed graph

$G = (V, E)$

Vertices (*pages*)

Edges (*links*)



0.225

**0.265**

0.138

0.127

0.172

0.074

Which vertex is most important?

$V = \{a, b, c, d, e, f\}$

$E = \{(a, e), (a, f), (b, d), (c, b), (d, a), (d, c), (d, f), (e, b), (e, d), (e, f), (f, a)\}$

# PageRank Model

- The Web: a directed graph

$$G = (V, E)$$

Vertices (*pages*)

Edges (*links*)



$$\text{out}(v) \coloneqq \{v' \in V \mid (v, v') \in E\}$$

$$\text{in}(v) \coloneqq \{v' \in V \mid (v', v) \in E\}$$

$$\text{rank}_0(v) \coloneqq \frac{1}{|V|}$$

$$\text{rank}_i(v) \coloneqq \sum_{v' \in \text{in}(v)} \frac{\text{rank}_{i-1}(v')}{|\text{out}(v')|}$$

HUH?

# PageRank Model

$\text{rank}_1(f) = \frac{1}{6} \times \frac{1}{3}$

f

$\text{rank}_0(e) = \frac{1}{6}$

$|\text{out}(e)| = 3$

$\text{rank}_1(b) = \frac{1}{6} \times \frac{1}{3}$

e

b

d

$\text{rank}_1(d) = \frac{1}{6} \times \frac{1}{3}$

$G = (V, E)$

Vertices
(*pages*)

Edges
(*links*)

$$\text{out}(v) \coloneqq \{v' \in V \mid (v, v') \in E\}$$

$$\text{in}(v) \coloneqq \{v' \in V \mid (v', v) \in E\}$$

$$\text{rank}_0(v) \coloneqq \frac{1}{|V|}$$

$$\text{rank}_i(v) \coloneqq \sum_{v' \in \text{in}(v)} \frac{\text{rank}_{i-1}(v')}{|\text{out}(v')|}$$

# PageRank Model

$$\text{rank}_2(a) = \frac{1}{6} \times \frac{1}{3} \times 1 + \frac{1}{6} \times \frac{1}{3} \times \frac{1}{3}$$

$$G = (V, E)$$

Vertices (*pages*)

Edges (*links*)

$$\text{rank}_1(f) = \frac{1}{6} \times \frac{1}{3} \qquad \text{rank}_1(a) = \frac{1}{6} \times 1 + \frac{1}{6} \times \frac{1}{3}$$

$$\text{rank}_0(e) = \frac{1}{6}$$

$$|\text{out}(e)| = 3$$

$$\text{rank}_1(e) = 0$$

$$\text{rank}_1(b) = \frac{1}{6} \times \frac{1}{3} + 1 \times \frac{1}{6}$$

$$\text{out}(v) \coloneqq \{v' \in V \mid (v, v') \in E\}$$

$$\text{in}(v) \coloneqq \{v' \in V \mid (v', v) \in E\}$$

$$\text{rank}_0(v) \coloneqq \frac{1}{|V|}$$

$$\text{rank}_i(v) \coloneqq \sum_{v' \in \text{in}(v)} \frac{\text{rank}_{i-1}(v')}{|\text{out}(v')|}$$

$$\text{rank}_1(d) = \frac{1}{6} \times \frac{1}{3} \qquad \text{rank}_0(c) = \frac{1}{6}$$

$$|\text{out}(c)| = 1$$

$$\text{rank}_1(c) = \frac{1}{6} \times \frac{1}{3}$$

$$\text{rank}_2(c) = \frac{1}{6} \times \frac{1}{3} \times \frac{1}{3}$$

# PageRank: Random Surfer Model



= someone surfing the web, clicking links randomly

- What is the probability of being at page $x$ after $n$ hops?

# PageRank: Random Surfer Model



= someone surfing the web, clicking links randomly

- What is the probability of being at page *x* after *n* hops?
- *Initial state:* surfer equally likely to start at any node

# PageRank: Random Surfer Model

= someone surfing the web, clicking links randomly

- What is the probability of being at page *x* after *n* hops?
- *Initial state:* surfer equally likely to start at any node
- PageRank applied iteratively for each hop: score indicates probability of being at that page after that many hops

What would happen with g over time?

# PageRank: Random Surfer Model



= someone surfing the web, clicking links randomly

- What is the probability of being at page *x* after *n* hops?
- *Initial state:* surfer equally likely to start at any node
- PageRank applied iteratively for each hop: score indicates probability of being at that page after than many hops
- If the surfer reaches a page without links, the surfer randomly jumps to another page

# PageRank: Random Surfer Model



= someone surfing the web, clicking links randomly

- What is the probability of being at page *x* after *n* hops?

- *Initial state:* surfer equally likely to start at any node

- PageRank applied iteratively for each hop: score indicates probability of being at that page after than many hops

- If the surfer reaches a page without links, the surfer randomly jumps to another page

What would happen with g and i over time?

# PageRank: Random Surfer Model



= someone surfing the web, clicking links randomly

- What is the probability of being at page *x* after *n* hops?
- *Initial state:* surfer equally likely to start at any node
- PageRank applied iteratively for each hop: score indicates probability of being at that page after than many hops
- If the surfer reaches a page without links, the surfer randomly jumps to another page

What would happen with g and i over time?

# PageRank: Random Surfer Model

= someone surfing the web, clicking links randomly

- What is the probability of being at page *x* after *n* hops?
- *Initial state:* surfer equally likely to start at any node
- PageRank applied iteratively for each hop: score indicates probability of being at that page after than many hops
- If the surfer reaches a page without links, the surfer randomly jumps to another page
- The surfer will jump to a random page at any time with a probability 1 − *d ... this avoids traps and ensures convergence!*

# PageRank Model: Final Version

- The Web: a directed graph

$$G = (V, E)$$

Vertices (*pages*)

Edges (*links*)

$$\text{out}(v) \coloneqq \{v' \in V \mid (v, v') \in E\}$$

$$\text{in}(v) \coloneqq \{v' \in V \mid (v', v) \in E\}$$
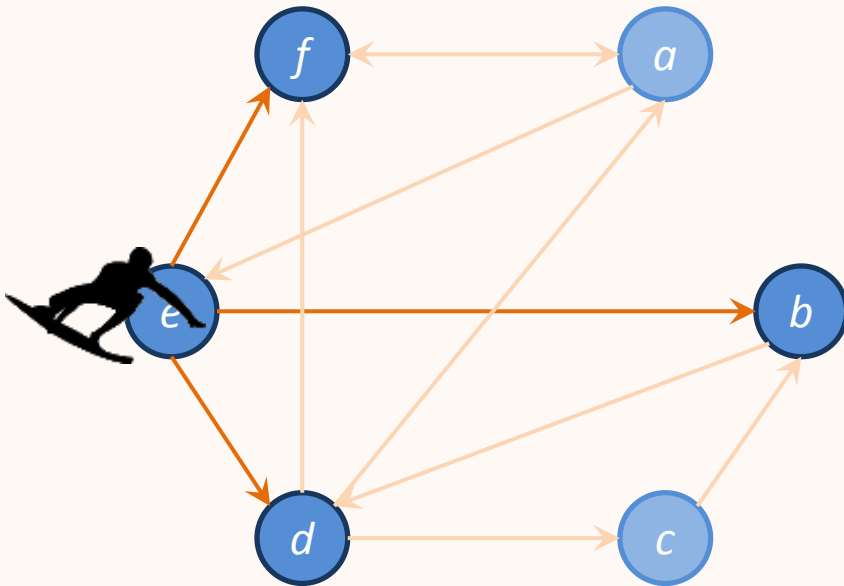
$$\text{rank}_0(v) \coloneqq \frac{1}{|V|}$$

$$V' \coloneqq \{v \in V : |\text{out}(v)| = 0\}$$

$$V'' \coloneqq \{v \in V : |\text{out}(v)| \neq 0\}$$

$d$ is the follow-a-link probability
typically $(d = 0.85)$

$$\text{rank}_i(v) \coloneqq d \times \sum_{u \in \text{in}(v)} \frac{\text{rank}_{i-1}(u)}{|\text{out}(u)|} + \sum_{v' \in V'} \frac{\text{rank}_{i-1}(v')}{|V|} + (1-d) \times \sum_{v'' \in V''} \frac{\text{rank}_{i-1}(v'')}{|V|}$$

# PageRank: Benefits



- ✓ More robust than a simple link count
- ✓ Fewer ties than link counting
- ✓ Scalable to approximate (for sparse graphs)
- ✓ Convergence guaranteed

# Two Sides to Ranking: Importance

# GOOGLE: A GUESS

# How Modern Google ranks results (maybe)



**Weighting of Thematic Clusters of Ranking Factors in Google**

(based on survey responses by 128 SEO professionals in June 2013)

**Domain-Level, Keyword-Agnostic Features**
(e.g. domain name length, TLD extension, domain HTTP response time, etc.) — 5.21%

**Domain Level Keyword Usage**
(e.g. exact-match keyword domains, partial-keyword matches, etc.) — 6.98%
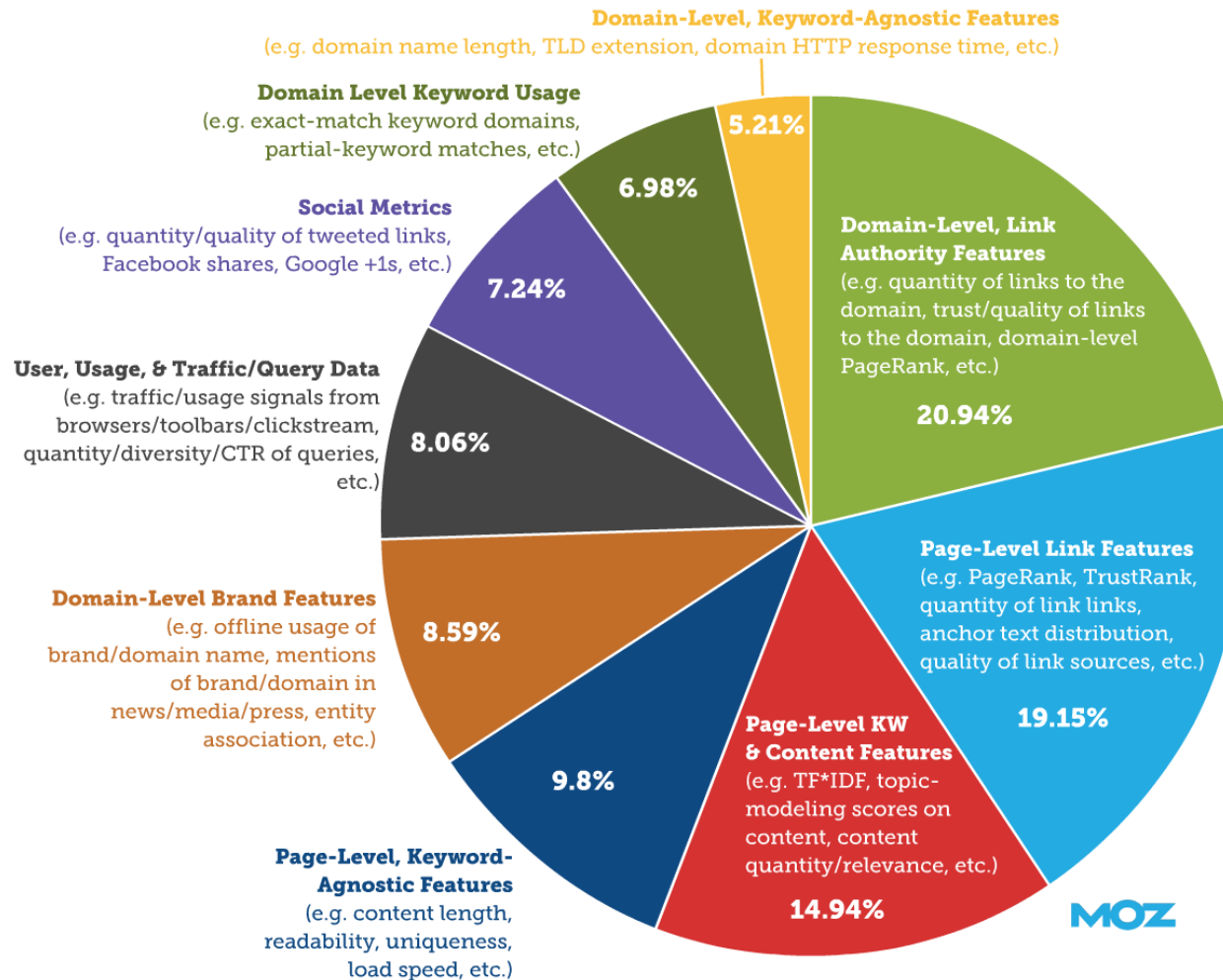
**Social Metrics**
(e.g. quantity/quality of tweeted links, Facebook shares, Google +1s, etc.) — 7.24%

**User, Usage, & Traffic/Query Data**
(e.g. traffic/usage signals from browsers/toolbars/clickstream, quantity/diversity/CTR of queries, etc.) — 8.06%

**Domain-Level Brand Features**
(e.g. offline usage of brand/domain name, mentions of brand/domain in news/media/press, entity association, etc.) — 8.59%

**Page-Level, Keyword-Agnostic Features**
(e.g. content length, readability, uniqueness, load speed, etc.) — 9.8%

**Domain-Level, Link Authority Features**
(e.g. quantity of links to the domain, trust/quality of links to the domain, domain-level PageRank, etc.) — 20.94%

**Page-Level Link Features**
(e.g. PageRank, TrustRank, quantity of link links, anchor text distribution, quality of link sources, etc.) — 19.15%

**Page-Level KW & Content Features**
(e.g. TF*IDF, topic-modeling scores on content, content quantity/relevance, etc.) — 14.94%

MOZ

*According to survey of SEO experts, not people in Google*

# How Modern Google ranks results (maybe)

Weighting of Thematic Clusters of Ranking Factors in Google
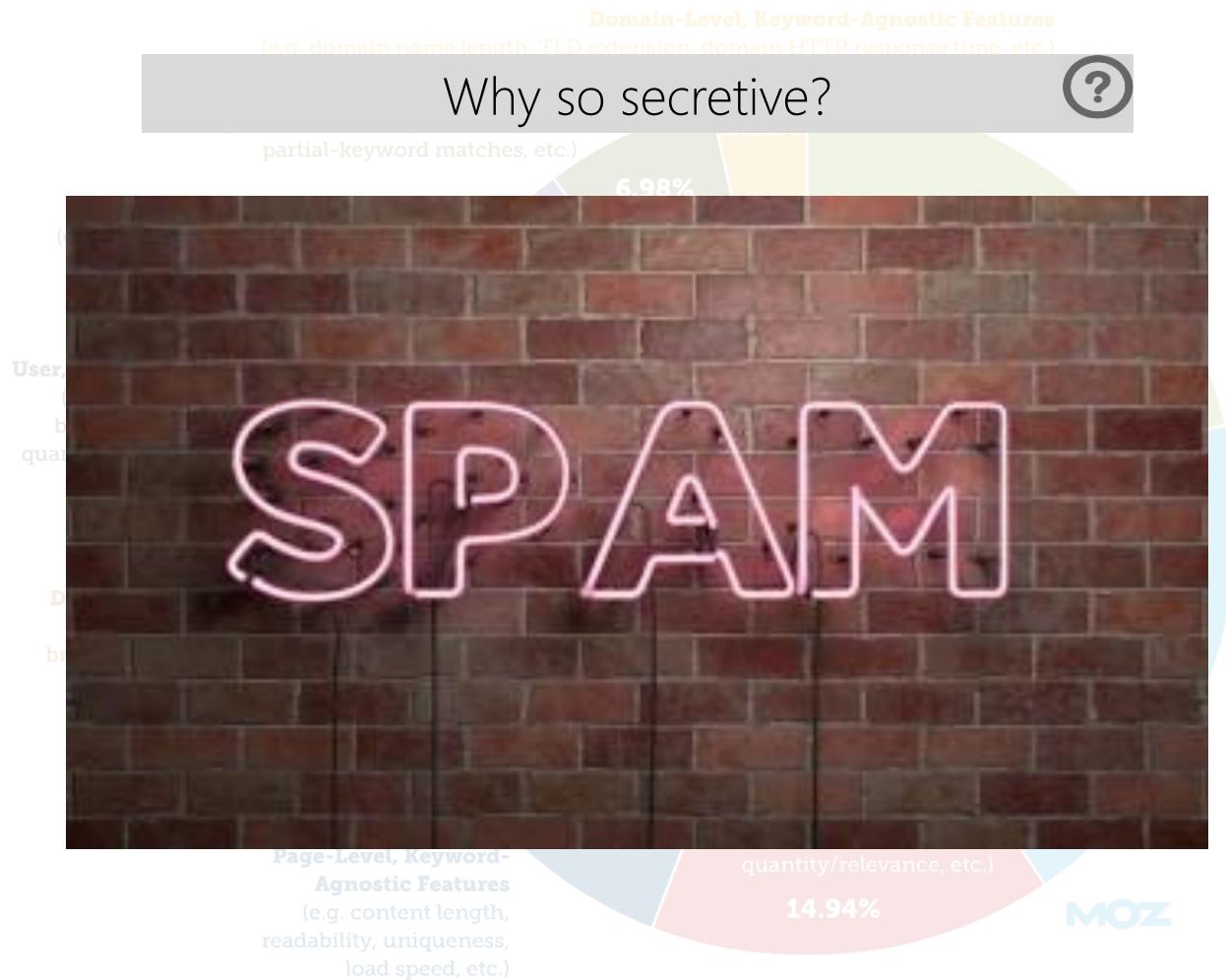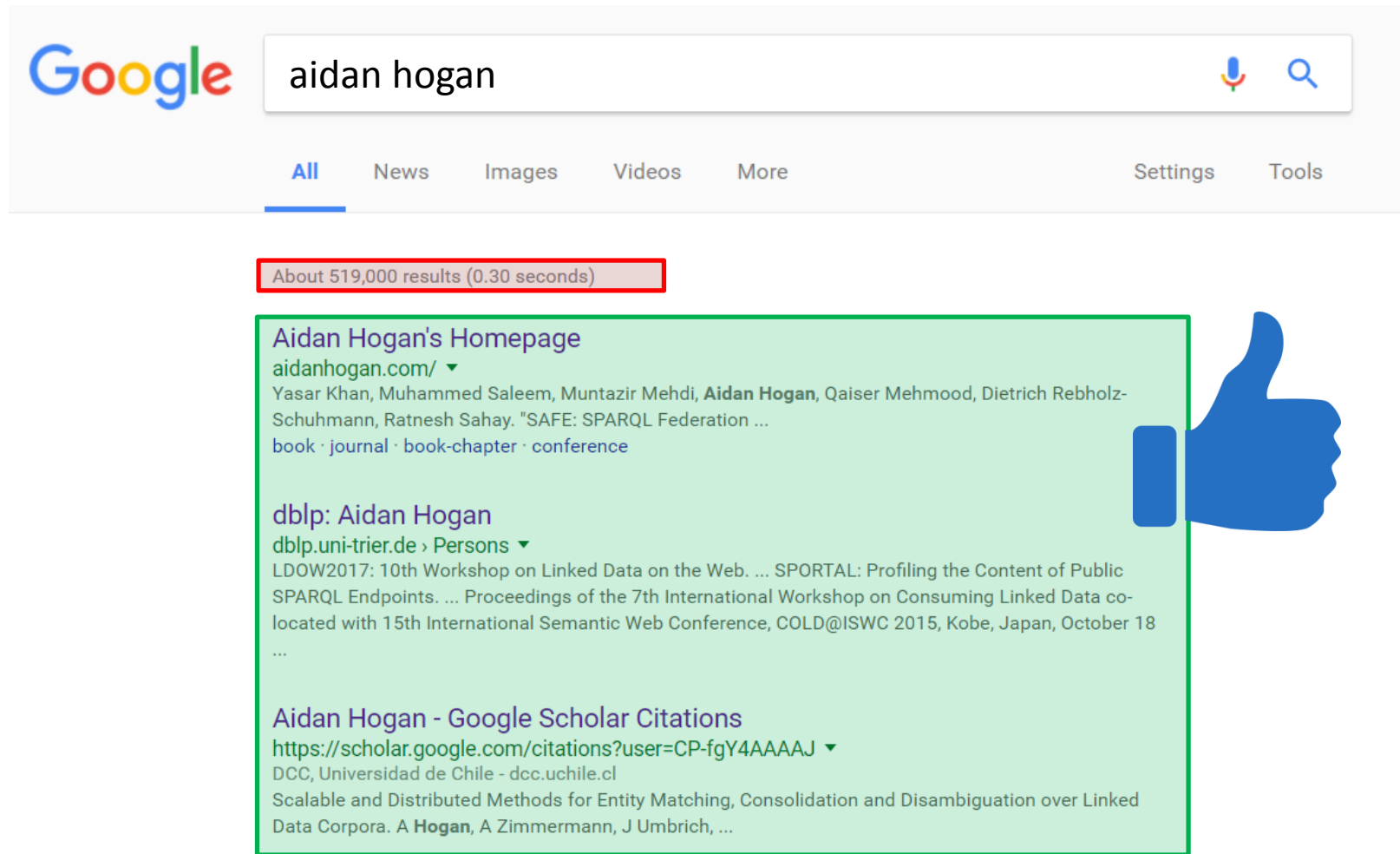(based on survey responses by 128 SEO professionals in June 2013)

Domain-Level, Keyword-Agnostic Features
(e.g. domain name length, TLD extension, domain HTTP response time, etc.)

partial-keyword matches, etc.)

6.98%

Why so secretive? ⓘ

User,

quar

Page-Level, Keyword-
Agnostic Features
(e.g. content length,
readability, uniqueness,
load speed, etc.)

quantity/relevance, etc.)

14.94%

MOZ

*According to survey of SEO experts, not people in Google*

# INFORMATION RETRIEVAL: RECAP

# How Does Google Get Such Good Results?

# Ranking in Information Retrieval

- Relevance: Is the document relevant for the query?
  - Term Frequency * Inverse Document Frequency
  - Cosine similarity

- Importance: Is the document a popular one?
  - Links analysis
  - PageRank

# Ranking: Science or Art?

Questions?