

Information Storage and Retrieval

CSCE 670
Texas A&M University
Department of Computer Science & Engineering
Instructor: Prof. James Caverlee

Text Analytics: Topic Models
II April 2017

Text Analytics

Data Mining View: Explore patterns in textual data

- Find latent topics

- Find topical trends

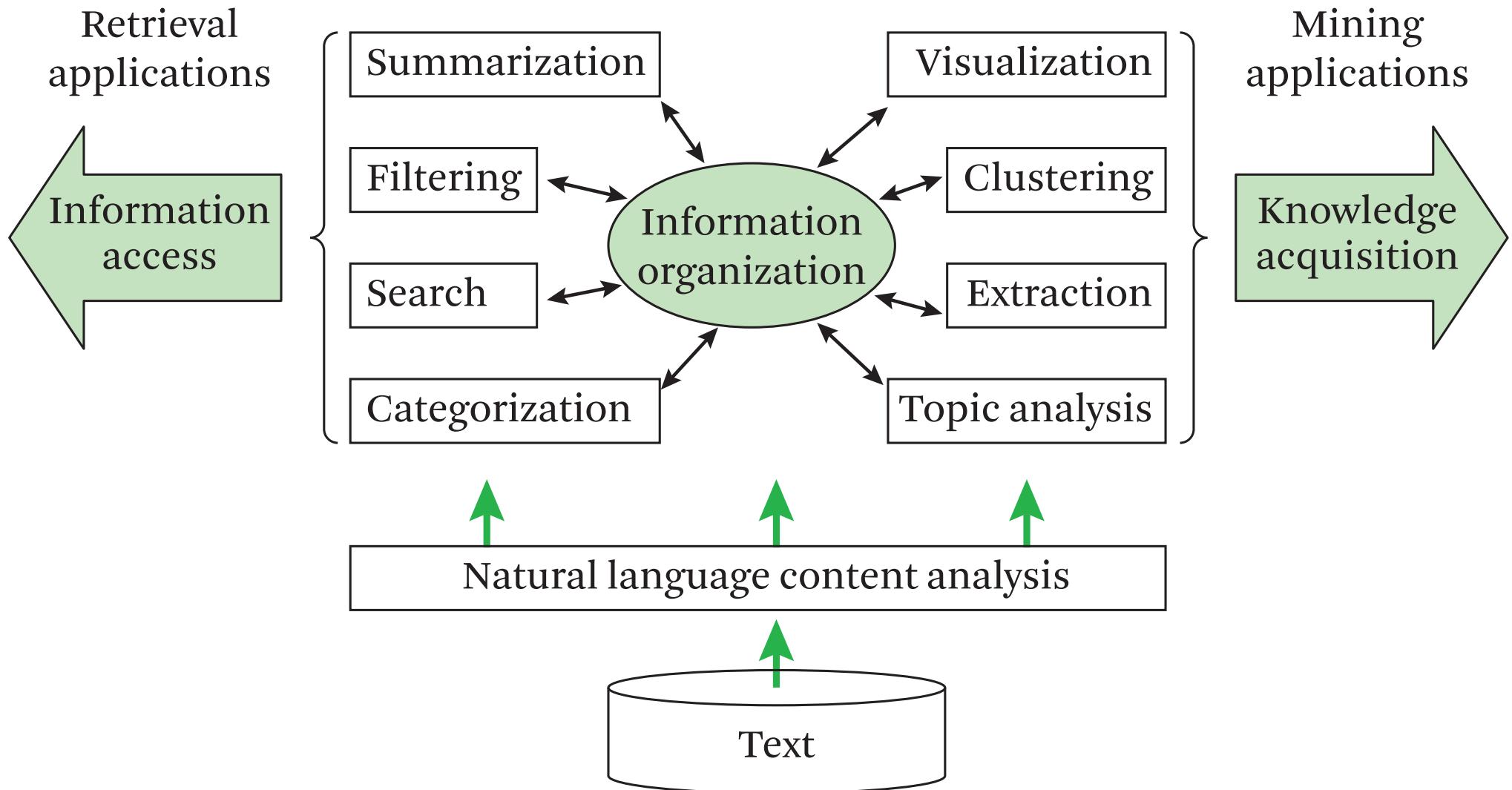
- Find outliers and other hidden patterns

Natural Language Processing View: Make inferences based on partial understanding of natural language text

- Information extraction

- Question answering

Text Analytics



Summarization

ⓘ textsummarization.net/text-summarizer



TextSummarization

Text Summarizer Online

Text Summarization API ▾

Text Summarization Result

Original URL/Text

Four score and seven years ago our fathers brought forth on this continent, a new nation, conceived in Liberty, and dedicated to the proposition that all men are created equal. Now we are engaged in a great civil war, testing whether that nation, or any nation so conceived and so dedicated, can long endure. We are met on a great battle-field of that war. We have come to dedicate a portion of that field, as a final resting place for those who here gave their lives that that nation might live. It is altogether fitting and proper that we should do this. But, in a larger sense, we can not dedicate -- we can not consecrate -- we can not hallow -- this ground. The brave men, living and dead, who struggled here, have consecrated it, far above our poor power to add or detract. The world will little note, nor long remember what we say here, but it can never forget what they did here. It is for us the living, rather, to be dedicated here to the unfinished work which they who fought here have thus far so nobly advanced. It is rather for us to be here dedicated to the great task remaining before us -- that from these honored dead we take increased devotion to that cause for which they gave the last full measure of devotion -- that we here highly resolve that these dead shall not have died in vain -- that this nation, under God, shall have a new birth of freedom -- and that government of the people, by the people, for the people, shall not perish from the earth.

Summarized Text

Four score and seven years ago our fathers brought forth on this continent, a new nation, conceived in Liberty, and dedicated to the proposition that all men are created equal.

We are met on a great battle-field of that war.

Text Categorization

Google Corporation®
123 Buckingham Palace Road
London SW1W 9SH
United Kingdom

Good day Sir/Madam.

You have successfully been picked as one of our 12 Lucky Winners in this months Lottery Draw, Please see attached file for more details.

Best regards
Sundar Pichai
CEO Google Inc.
sundarpichai@gpromo-team.com

Text Categorization

The 2nd edition of iOS Apps for Masterminds is now available. This is the first book to teach how to program applications for iOS 10 with Swift 3 and Xcode 8.

iOS Apps for Masterminds is great for a semester course, a complement to an undergraduate program, or to get you familiar with the technology. The book gradually introduces the concepts to guide students step by step on how to create an application from scratch. It goes from basic programming concepts like variables and memory management to more complex subjects like Core Data and iCloud (more than 800 pages of content).

Read the preview at Amazon
www.amazon.com/dp/1537517880/

The For Masterminds books are also the less expensive and more comprehensive books on the subjects. Perfect for a low-budget course.

To see the Table of Contents and examples included in the book, visit our website
www.formasterminds.com

For additional information or promo codes to download the eBook for free, please reply this email or contact me at info@jdgauchat.com.

Text Categorization

U.S. Presidential Debate



Two tweets are shown. The first tweet is from "Bloomberg Politics" (@bpolitics) with a purple logo, stating: "We could uncover a dog fighting camp in @realDonaldTrump's backyard & people would find a way to justify it. #NeverTrump #debate2016". The second tweet is from "CPD" (@debates) with a yellow logo, stating: "Live stream of the Presidential Debate: Hillary Clinton vs Donald Trump - thedmvdaily.com/live-stream-of... via @TheDMV р Daily".

A video thumbnail showing a stage setup for a debate, featuring a podium and a backdrop with stars.

Text Categorization

How would you represent each “document”?

How would you compare them?

What approach (“classifier”) would you use?

How would you evaluate your approach?

Text Categorization

Other examples?

Text Categorization

Other examples?

Authorship attribution

Positive or negative movie review

Age / gender identification

...

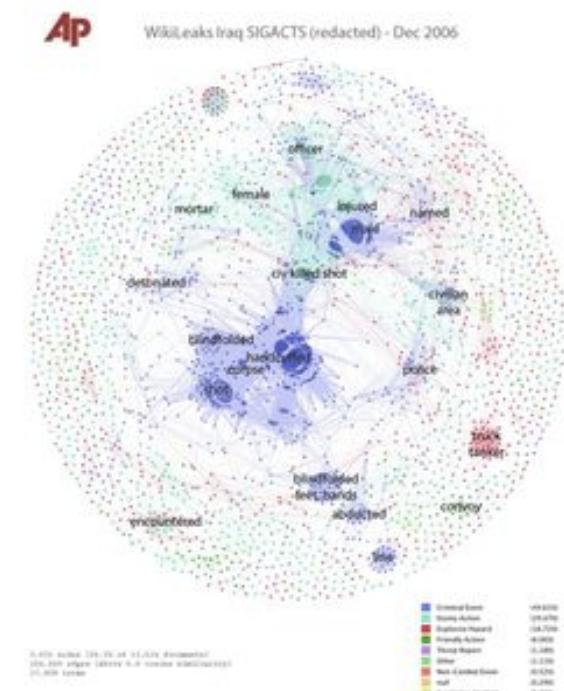
Document Clustering

Suppose you're an analyst for:

GoldmanSachs looking for **major shifts in a company** (e.g., Theranos)

the **NSA** looking for **evidence of a terror cell operating in a particular country**

... ?



Document Clustering

How would you represent each “document”?

How would you compare them?

What approach (“clustering algorithm”) would you use?

How would you evaluate your approach?

Information Extraction

Find and understand limited relevant parts of texts

Gather information from many pieces of text

Produce a structured representation of relevant information:

relations (in the database sense), a.k.a., a knowledge base

Goals:

Organize information so that it is useful to people

Put information in a semantically precise form that allows further inferences to be made by computer algorithms



Sherry Escalante <sherry.escalante@tamu.edu>

Oct 7 (2 days ago)



to AM-ENGINEER-FA.

Engineering Faculty,

As you know, service on a faculty search committee requires you to take the training session offered by the Dean of Faculties. They will be offering a

session for engineering faculty and staff on **Wednesday, October 26th**
from 2:00-4:00 pm in room 217 of the Civil Engineering Office Building
(CEOI)

[Add to Calendar](#)

, please register here by 10/21 :



what year was texas founded

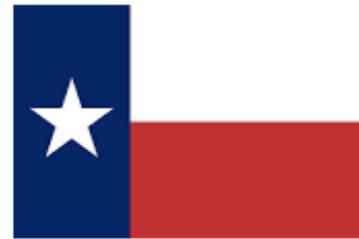


All News Images Shopping Maps More ▾ Search tools

About 7,490,000 results (0.69 seconds)

Texas / Founded

December 29, 1845



Feedback

People also ask

What is the history of Texas? ▾

When was Texas first settled? ▾

What is the state of Texas? ▾

Why is the capital of Texas Austin? ▾

[Texas - Wikipedia, the free encyclopedia](#)

<https://en.wikipedia.org/wiki/Texas> ▾ [Wikipedia](#) ▾

Information Extraction

How would you approach this problem?

Text Visualization

The 2007 State of the Union Address

Over the years, President Bush's State of the Union address has averaged almost 5,000 words each, meaning the the President has delivered over 34,000 words. Some words appear frequently while others appear only sporadically. Use the tools below to analyze what Mr. Bush has said.

 or choose a word here.

Use of the phrase "Economy" in past State of the Union Addresses

2001*	2002	2003	2004	2005	2006	2007
3	4	10	14	11	15	7



The word in context

Government has a role, and an important role. Yet, too much government crowds out initiative and hard work, private charity and the private ECONOMY. Our new governing vision says government should be active, but limited; engaged, but not overbearing. And my budget is based on that philosophy.

— 2001 (Paragraph 6 of 73)

Next Instance of 'Economy'

Compared with other words

2001*	2002	2003	2004	2005	2006	2007
-------	------	------	------	------	------	------

Economy

3	4	10	14	11	16	7
---	---	----	----	----	----	---

Afghanistan

-	13	3	5	3	2	4
---	----	---	---	---	---	---

Economy(ic)

6	7	13	17	14	23	8
---	---	----	----	----	----	---

Insurance

2	-	5	6	5	3	14
---	---	---	---	---	---	----

Iraq/Iraqi(s)

-	2	21	24	27	16	34
---	---	----	----	----	----	----

Iran

-	2	3	5	3	6	5
---	---	---	---	---	---	---

Oil

-	5	-	-	-	3	9
---	---	---	---	---	---	---

Social Security

15	2	2	2	18	3	2
----	---	---	---	----	---	---

CHAPTERS



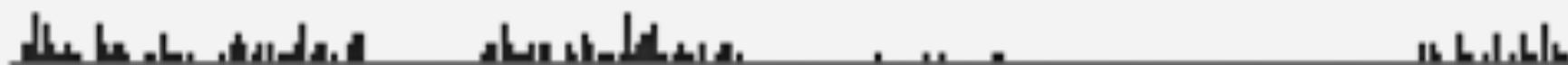
SENTIMENT



BILBO
554 TIMES



GANDALF
187 TIMES



DWALIN
21 TIMES



BALIN
67 TIMES



FILI
49 TIMES



KILI
37 TIMES



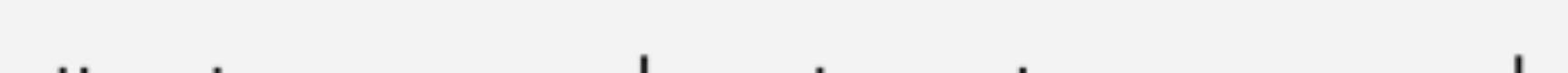
DORI
34 TIMES



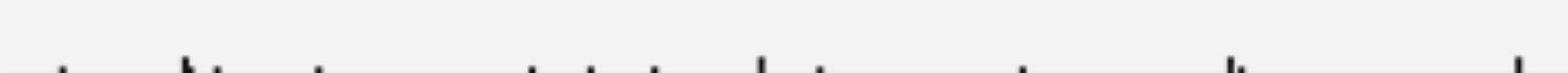
NORI
17 TIMES



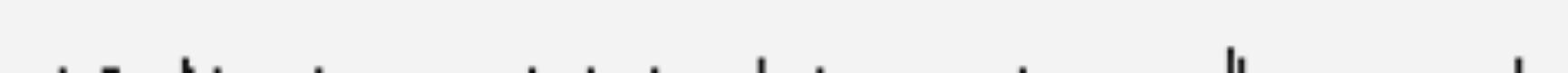
ORI
11 TIMES



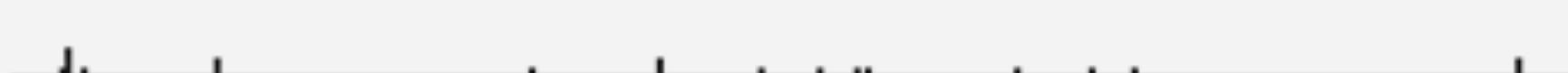
ÓIN
18 TIMES

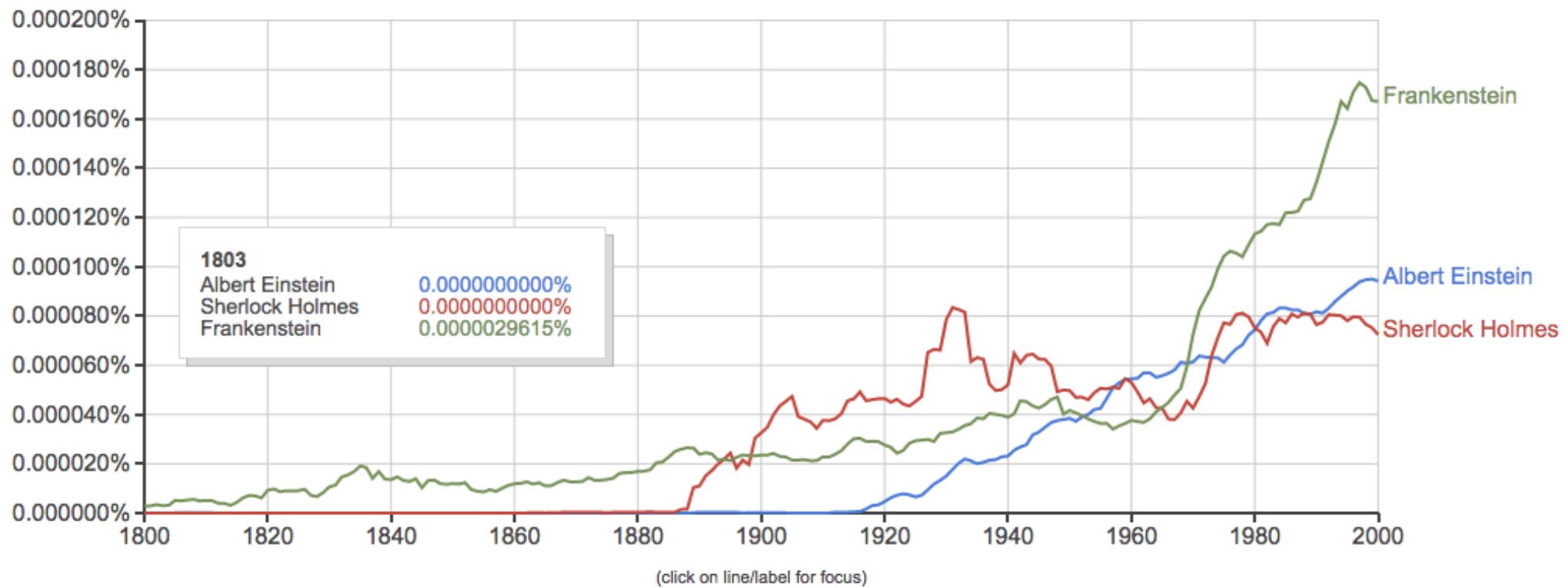


GLÓIN
23 TIMES



BIFUR
19 TIMES





<https://books.google.com/ngrams>

English

Version 20120701

total counts

1-grams 0 1 2 3 4 5 6 7 8 9 a b c d e f g h i j k l m n o other p pos punctuation q r s t u v w x y z

2-grams 0 1 2 3 4 5 6 7 8 9 _ADJ _ADP _ADV _CONJ _DET _NOUN _NUM _PRON _PRT _VERB _a _aa _ab _ac _ad _ae _af _ag _ah _ai _aj _ak _al _am _an _ao _ap _aq _ar _as _at _au _av
aw ax ay az b_ ba bb bc bd be bf bg bh bi bj bk bl bm bn bo bp bq br bs bt bu bv bw bx by bz c_ ca cb cc cd ce cf cg ch ci cj ck cl cm cn co cp cq cr cs ct cu cv cw cx cy cz d_ da db dc
dd de df dg dh di dj dk dl dm dn do dp dq dr ds dt du dv dw dx dy dz e_ ea eb ec ed ee ef eg eh ei ej ek el em en eo ep eq er es et eu ev ew ex ey ez f_ fa fb fc fd fe ff fg fh fi fj fk fl fm
fn fo fp fq fr fs ft fu fv fw fx fy fz g_ ga qb qc gd ge gf gg gh qj gj qk gl gm gn go gp qq qr gs qt qu gv gw qx qy gz h_ ha hb hc hd he hf hg hh hi hj hk hl hm hn ho hp hq hr hs ht hu hv
hw hx hy hz l_ ia jb jc id ie if ig ih ii jj ik ll jm in io ip iq ir is it iu lv lx ly lz l_ ja jb jc jd ie if jg jh jj ll jm jn jo jp qj jr is jt ju jv jw jx ly zk_ ka kb kc kd ke kf kg kh ki kj kk kl km kn ko kp
kq kr ks kt ku kv kw kx ky kz l_ la lb lc ld le lf lg lh ii jj lk ll lm in lo lp lq lr ls lt lu lv lx ly lz m_ ma mb mc md me mf mg mh mi mj mk ml mm mn mo mp mq mr ms mt mu mv mw mx
my mz n_ na nb nc nd ne nf ng nh ni nj nk nl nm nn no np nq nr ns nt nu nv nw nx ny nz o_ oa ob oc od oe of og oh oi oj ok ol om on oo op oq or os ot other ou ov ow ox oy oz p_ pa
pb pc pd pe pf pg ph pi pj pk pl pm pn po pp pq pr ps pt pu punctuation pv pw px py pz q_ qa qb qc qd ge gf gg qh qj qk gl qm an qo qp qq qr qs qt qu gv gw qx qy qz r_ ra rb rc rd
re rf rg rh ri rj rk rl rm rn ro rp rq rr rs rt ru rv rw rx ry rz s_ sa sb sc sd se sf sg sh si sj sk sl sm sn so sp sq sr ss st su sv sw sx sy sz t_ ta tb tc td te tf tg th ti tj tk tl tm tn to tp tq tr ts tt u
tv tw tx ty tz u_ ua ub uc ud ue uf ug uh ui uj uk ul um un uo up uq ur us ut uu uv uw ux uy uz v_ va vb vc vd ve vf vg vh vi vj vk vl vm vn vo vp vq vr vs vt vu vv vw vx vy vz w_ wa wb
wc wd we wf wg wh wi wj wk wl wn wo wp wq wr ws wt wu ww wx wy wz x_ xa xb xc xd xe xf xq xh xi xj xk xl xm xn xo xp xq xr xs xt xu xv xw xx xy xz y_ ya yb yc yd ye yf yg
yh yi yj yk yl ym yn yo yp yq yr ys yt yu yv yw yx yv yz z_ za zb zc zd ze zf zg zh zi zj zk zl zm zn zo zp zq zr zs zt zu zv zw zx zy zz

3-grams 0 1 2 3 4 5 6 7 8 9 _ADJ _ADP _ADV _CONJ _DET _NOUN _NUM _PRON _PRT _VERB _a _aa _ab _ac _ad _ae _af _ag _ah _ai _aj _ak _al _am _an _ao _ap _aq _ar _as _at _au _av
aw ax ay az b_ ba bb bc bd be bf bg bh bi bj bk bl bm bn bo bp bq br bs bt bu bv bw bx by bz c_ ca cb cc cd ce cf cg ch ci cj ck cl cm cn co cp cq cr cs ct cu cv cw cx cy cz d_ da db dc
dd de df dg dh di dj dk dl dm dn do dp dq dr ds dt du dv dw dx dy dz e_ ea eb ec ed ee ef eg eh ei ej ek el em en eo ep eq er es et eu ev ew ex ey ez f_ fa fb fc fd fe ff fg fh fi fj fk fl fm
fn fo fp fq fr fs ft fu fv fw fx fy fz g_ ga gb qc gd ge gf gg gh gi gj qk gl gm gn go gp qq qr gs gt qu gv gw qx qy gz h_ ha hb hc hd he hf hg hh hi hj hk hl hm hn ho hp hq hr hs ht hu hv
hw hx hy hz l_ ia ib jc id ie if lg lh ii jk ll jm ln lo jp iq lr ls it lu lv lw ix ly lz l_ ja jb jc jd je if jq jh ii jk ll jm ln jo jp iq jr ls jt ju lv lw ix ly lz k_ ka kb kc kd ke kf kg kh ki kj kk kl km kn ko kp
kq kr ks kt ku kv kw kx ky kz l_ la lb lc ld le lf lg lh ii jk ll lm ln lo lp lq lr ls lt lu lv lw lx ly lz m_ ma mb mc md me mf mg mh mi mj mk ml mm mn mo mp mq mr ms mt mu mv mw mx
my mz n_ na nb nc nd ne nf ng nh ni nj nk nl nm nn no np nq nr ns nt nu nv nw nx ny nz o_ oa ob oc od oe of og oh oi oj ok ol om on oo op oq or os ot other ou ov ow ox oy oz p_ pa
pb pc pd pe pf pg ph pi pj pk pl pm pn po pp pq pr ps pt pu punctuation pv pw px py pz q_ qa qb qc qd qe qf qq qh qj qk ql qm qn qo qp qq qr qs qt qu qv qw qx qy qz r_ ra rb rc rd
re rf rg rh ri rj rk ll rm rn ro rp rq rr rs rt ru rv rw rx ry rz s_ sa sb sc sd se sf sg sh si sj sk sl sm sn so sp sq sr ss st su sv sw sx sy sz t_ ta tb tc td te tf tg th ti tj tk tl tm tn to tp tq tr ts tt u
tv tw tx ty tz u_ ua ub uc ud ue uf ug uh ui uj uk ul um un uo up uq ur us ut uu uv uw ux uy uz v_ va vb vc vd ve vf vg vh vi vj vk vl vm vo vp vq vr vs vt vu vv vw vx vy vz w_ wa wb
wc wd we wf wg wh wi wj wk wl wm wn wo wp wq wr ws wt wu ww wx wy wz x_ xa xb xc xd xe xf xg xh xi xj xk xl xm xn xo xp xq xr xs xt xu xv xw xx xy xz y_ ya yb yc yd ye yf yg
yh yi yj yk yl ym yn yo yp yq yr ys yt yu yv yw yx yv yz z_ za zb zc zd ze zf zg zh zi zj zk zl zm zn zo zp zq zr zs zt zu zv zw zx zy zz

4-grams 0 1 2 3 4 5 6 7 8 9 _ADJ _ADP _ADV _CONJ _DET _NOUN _NUM _PRON _PRT _VERB _a _aa _ab _ac _ad _ae _af _ag _ah _ai _aj _ak _al _am _an _ao _ap _aq _ar _as _at _au _av _aw _ax _ay _az _b _ba _bb _bc _bd _be _bf _bg _bh _bi _bj _bk _bl _bm _bn _bo _bp _bq _br _bs _bt _bu _bv _bw _bx _by _bz _c _ca _cb _cc _cd _ce _cf _cg _ch _ci _cj _ck _cl _cm _cn _co _cp _cq _cr _cs _ct _cu _cv _cw _cx _cy _cz _d _da _db _dc _dd _de _df _dg _dh _di _dj _dk _dl _dm _dn _do _dp _dq _dr _ds _dt _du _dv _dw _dx _dy _dz _e _ea _eb _ec _ed _ee _ef _eg _eh _ei _ej _ek _el _em _en _eo _ep _eq _er _es _et _eu _ev _ew _ex _ey _ez _f _fa _fb _fc _fd _fe _ff _fg _fh _fi _fj _fk _fl _fm _fn _fo _fp _fq _fr _fs _ft _fu _fv _fw _fx _fy _fz _g _ga _gb _gc _gd _ge _gf _gg _gh _gi _gk _gl _gm _gn _go _gp _gq _gr _gs _gt _gu _gv _gw _gx _gy _gz _h _ha _hb _hc _hd _he _hf _hg _hh _hi _hj _hk _hl _hm _hn _ho _hp _hq _hr _hs _ht _hu _hv _hw _hx _hy _hz _i _ia _ib _ic _id _ie _if _ig _ih _ii _ij _ik _il _im _in _io _ip _iq _ir _is _it _iu _iv _iw _ix _iy _iz _j _ja _jb _jc _jd _je _if _ig _jh _ji _jk _il _jm _in _jo _jp _iq _jr _is _jt _ju _jv _jw _jx _jy _k _ka _kb _kc _kd _ke _kf _kg _kh _ki _kj _kk _kl _km _kn _ko _kp _kq _kr _ks _kt _ku _kv _kw _kx _ky _kz _l _la _lb _lc _ld _le _lf _lg _lh _lj _lk _ll _lm _ln _lo _lp _lq _lr _ls _lt _lu _lv _lx _ly _lz _m _ma _mb _mc _md _me _mf _mg _mh _mi _mj _mk _ml _mm _mn _mo _mp _mq _mr _ms _mt _mu _mv _mw _mx _my _mz _n _na _nb _nc _nd _ne _nf _ng _nh _ni _nj _nk _nl _nm _nn _no _np _nq _nr _ns _nt _nu _nv _nw _nx _ny _nz _o _oa _ob _oc _od _oe _of _og _oh _oi _oj _ok _ol _om _on _oo _op _oq _or _os _ot _other _ou _ov _ow _ox _oy _oz _p _pa _pb _pc _pd _pe _pf _pg _ph _pi _pk _pl _pm _pn _po _pp _pq _pr _ps _pt _pu _punctuation _pv _pw _px _py _pz _q _qa _qb _qc _qd _qe _gf _gg _qh _qi _qj _ak _al _qm _qn _qo _qp _qq _qr _qs _gt _qu _qv _gw _qx _qy _qz _r _ra _rb _rc _rd _re _rf _rg _rh _ri _rj _rk _rl _rm _rn _ro _rp _rq _rr _rs _rt _ru _rv _rw _rx _ry _rz _s _sa _sb _sc _sd _se _sf _sg _sh _si _sj _sk _sl _sm _sn _so _sp _sq _sr _ss _st _su _sv _sw _sx _sy _sz _t _ta _tb _tc _td _te _tf _tg _th _ti _tj _tk _tl _tn _to _tp _qr _tr _ts _tt _tu _tv _tw _tx _ty _tz _u _ua _ub _uc _ud _ue _uf _ug _uh _ui _ui _uk _ul _um _un _uo _up _ug _ur _us _ut _uu _uv _uw _ux _uy _uz _v _va _vb _vc _vd _ve _vf _vg _vh _vi _vk _vl _vn _vo _vp _va _vr _vs _vt _vu _vv _vw _vx _vv _yz _w _wa _wb

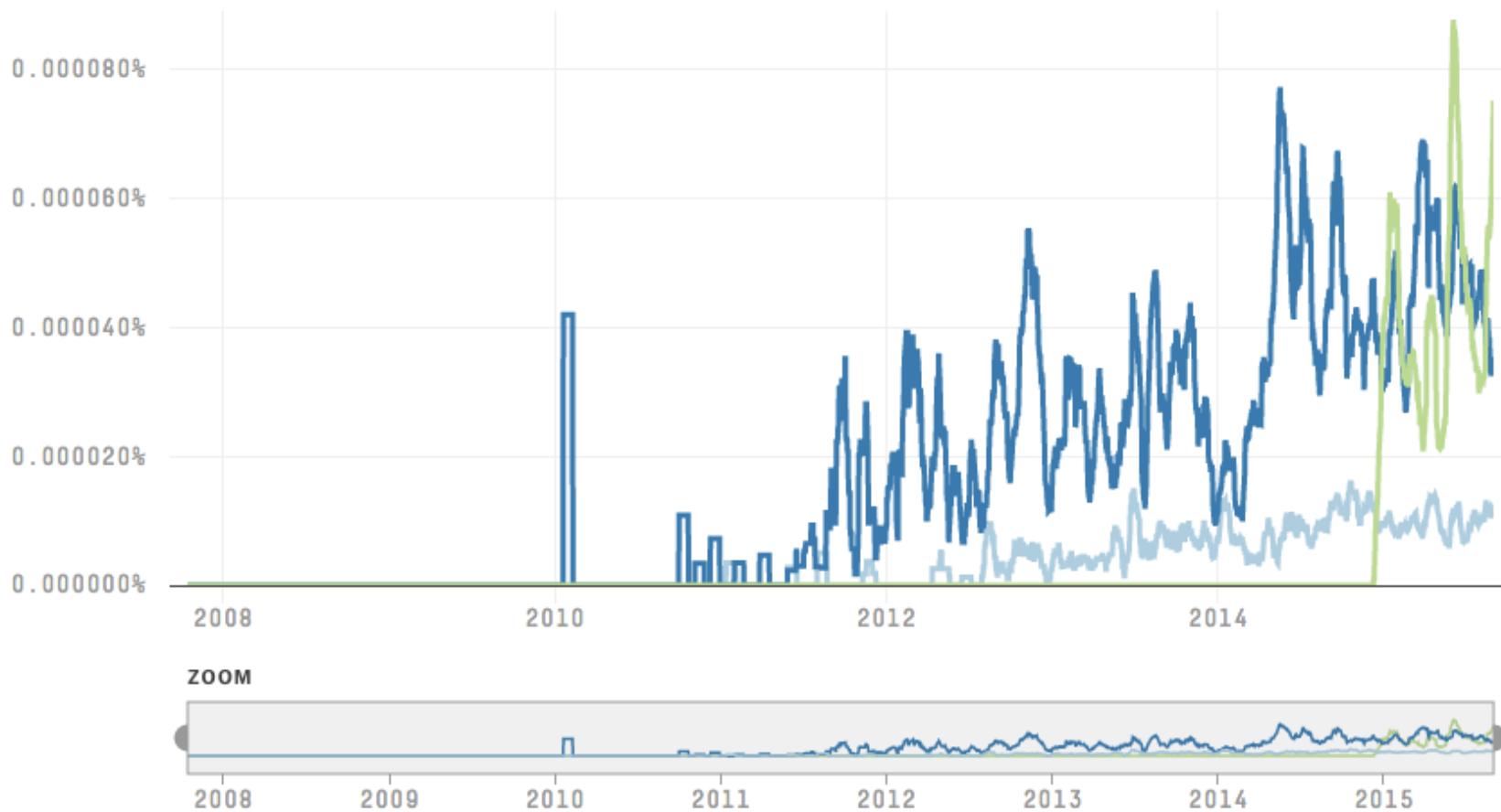
MANHANDLING X

MANSPLAINING X

MANSPREADING X

Enter terms

SEARCH



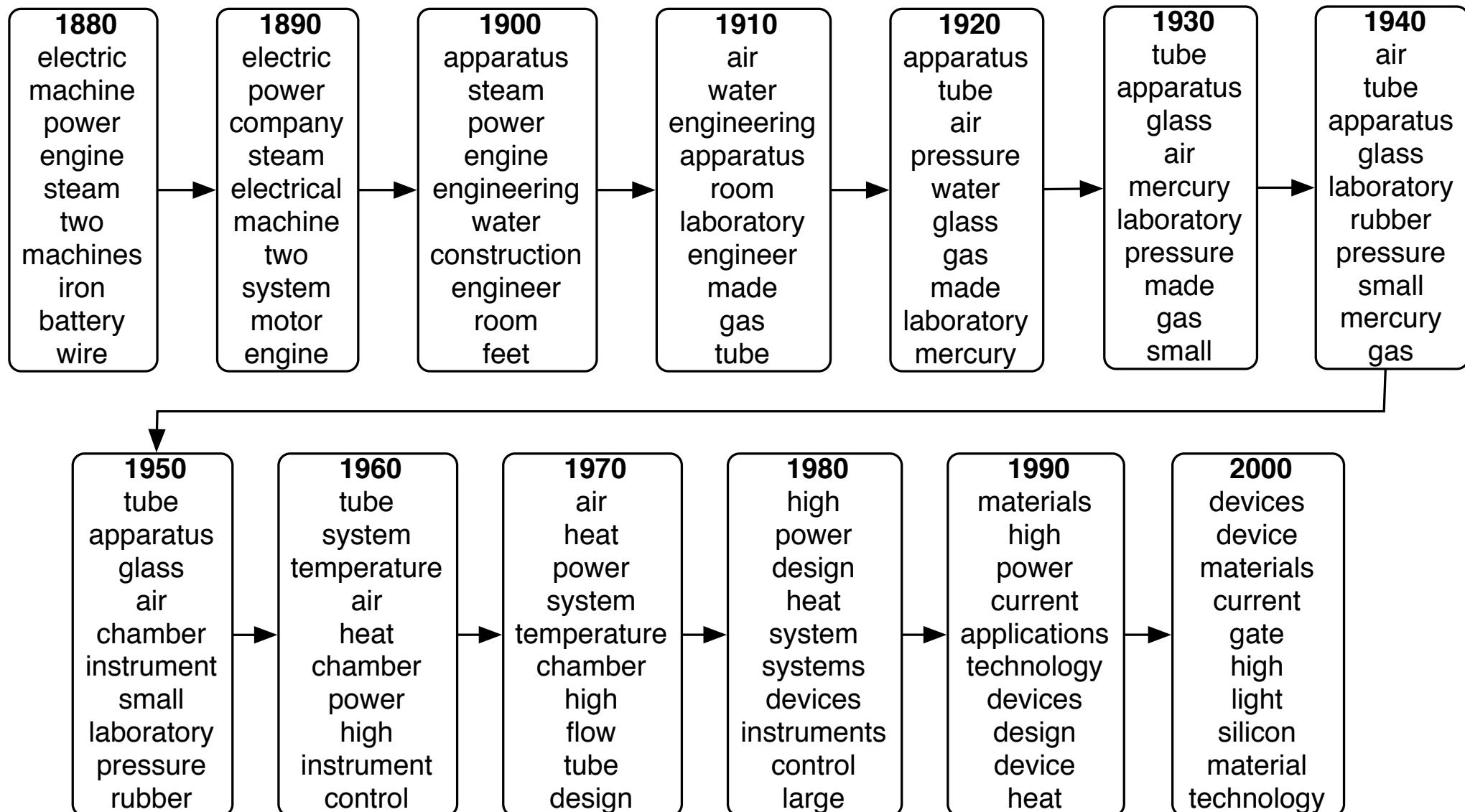
Topic Models

<http://www.cs.columbia.edu/~blei/papers/Blei2012.pdf>



Input: An unorganized collection of documents

Output: An organized collection, and a description of how



LDA = Latent Dirichlet Allocation

by Blei, Ng, and Jordan at Berkeley

introduced in 2003

hugely influential

tons of implementations



Dave Blei: Professor at Columbia, ACM Fellow



Andrew Ng: Professor at Stanford, founded Google Brain, founded Coursera, formerly Chief Scientist at Baidu



Michael Jordan: Professor at Berkeley, hugely influential, (dad is an Aggie!)

Topics

gene 0.04
dna 0.02
genetic 0.01
...,

life 0.02
evolve 0.01
organism 0.01
...,

brain 0.04
neuron 0.02
nerve 0.01
...,

data 0.02
number 0.02
computer 0.01
...,

Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

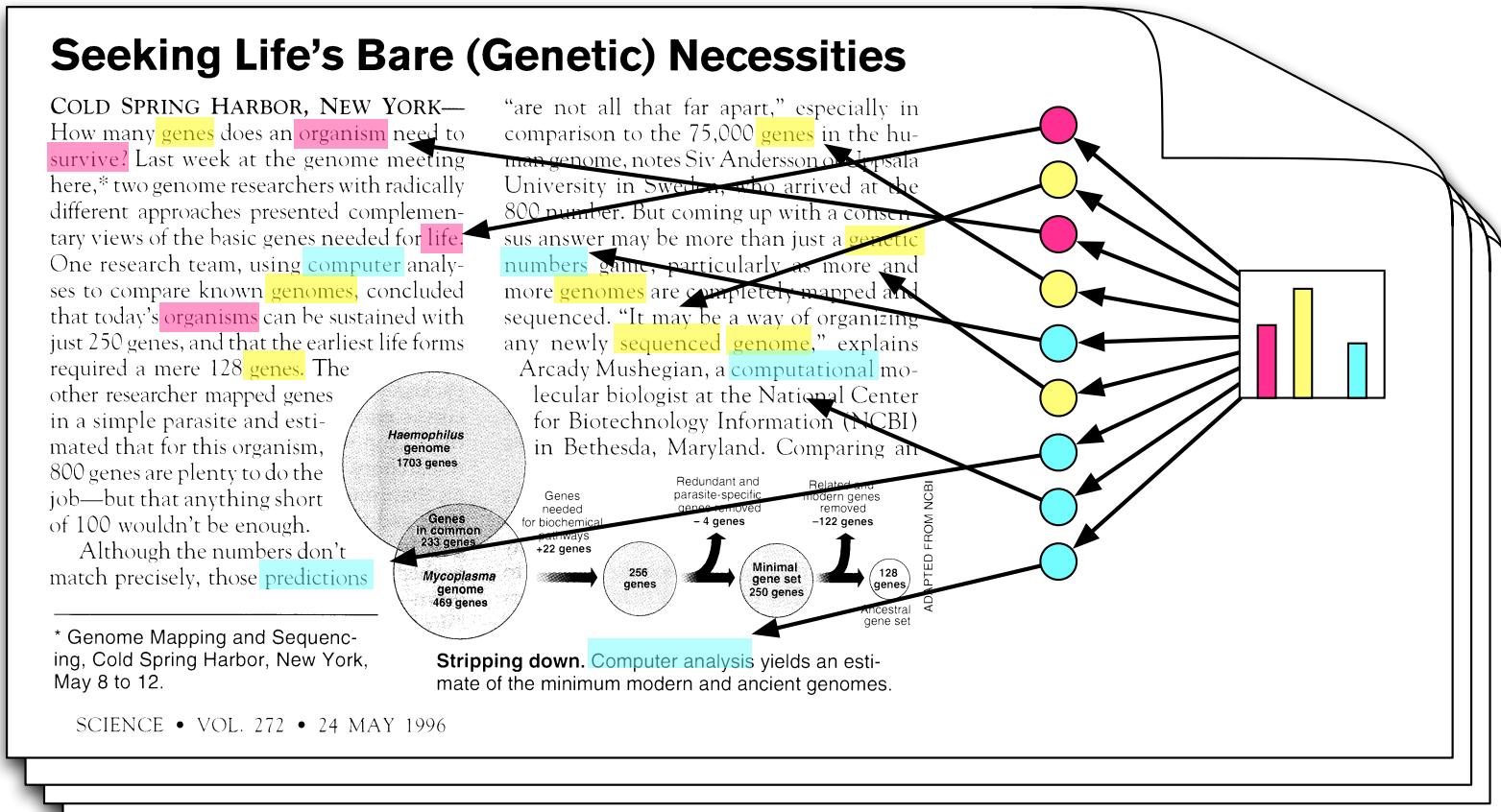
Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

Topic proportions and assignments



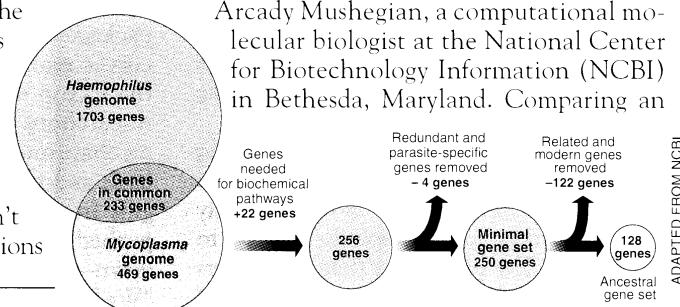
Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

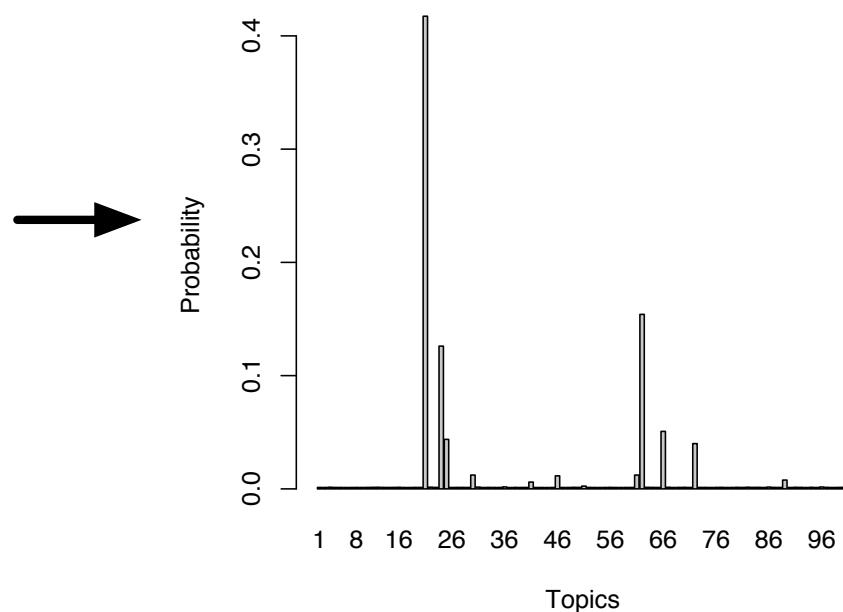
Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

ADAPTED FROM NCBI



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.



* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

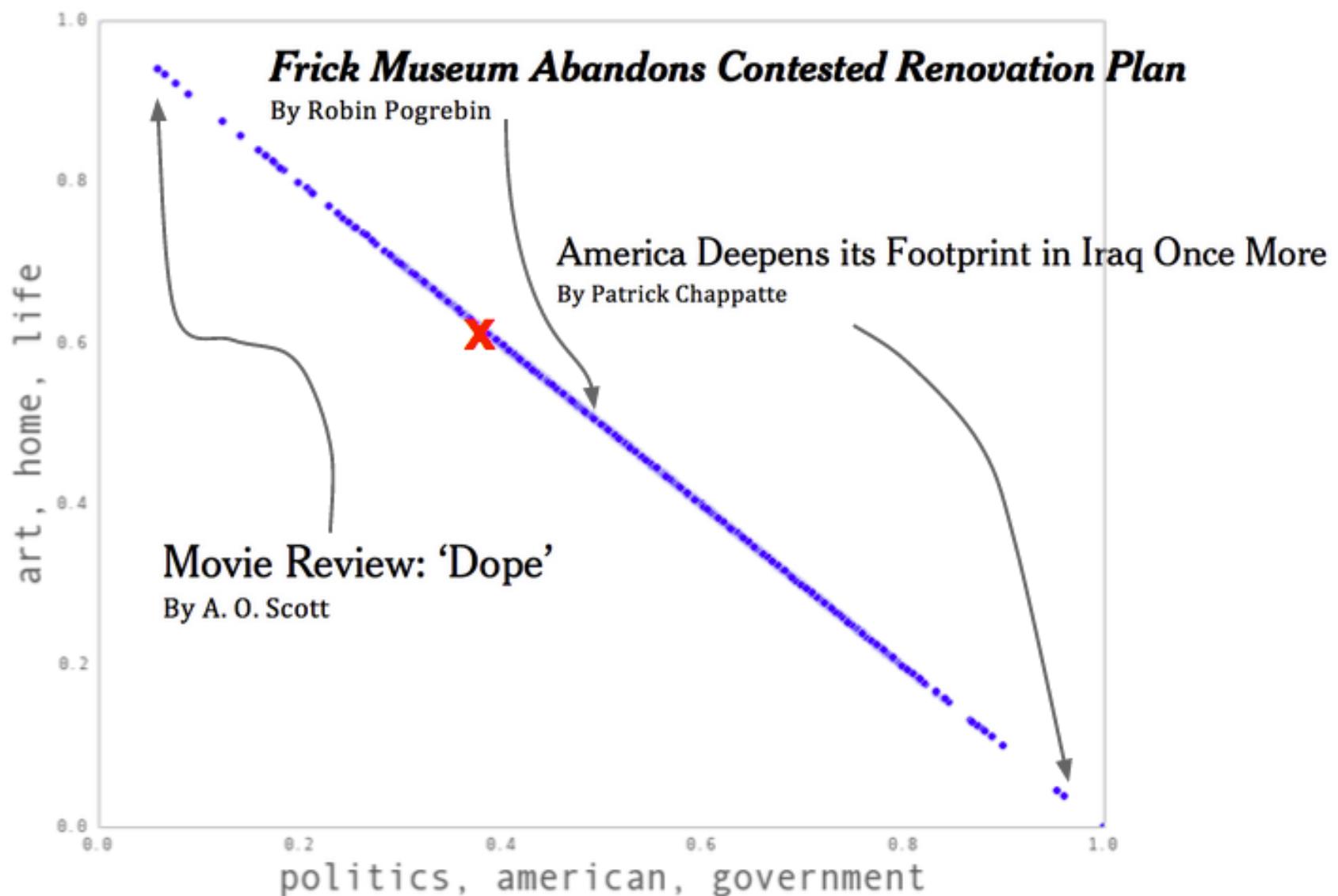
human
genome
dna
genetic
genes
sequence
gene
molecular
sequencing
map
information
genetics
mapping
project
sequences

evolution
evolutionary
species
organisms
life
origin
biology
groups
phylogenetic
living
diversity
group
new
two
common

disease
host
bacteria
diseases
resistance
bacterial
new
strains
control
infectious
malaria
parasite
parasites
united
tuberculosis

computer
models
information
data
computers
system
network
systems
model
parallel
methods
networks
software
new
simulations

LDA in action ...



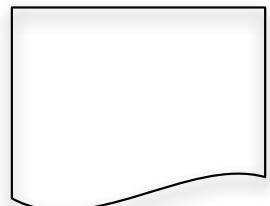
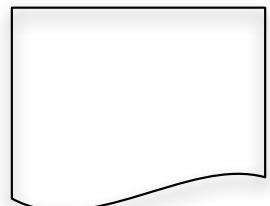
Topics



Documents



Topic proportions and assignments



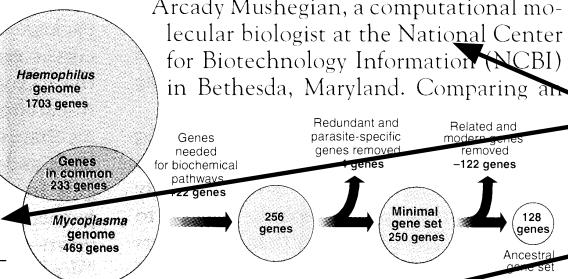
Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all



Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

Consider the following “documents”

I like to eat broccoli and bananas.

I ate a banana and spinach smoothie for breakfast.

Chinchillas and kittens are cute.

My sister adopted a kitten yesterday.

Look at this cute hamster munching on a piece of broccoli.



= 30% broccoli, 15% bananas, 10% breakfast, 10% munching, ... (at which point, you could interpret topic A to be about food)



= 20% chinchillas, 20% kittens, 20% cute, 15% hamster, ... (at which point, you could interpret topic B to be about cute animals)

LDA is a generative model

LDA represents documents as mixtures of topics that spit out words with certain probabilities. Here's how:

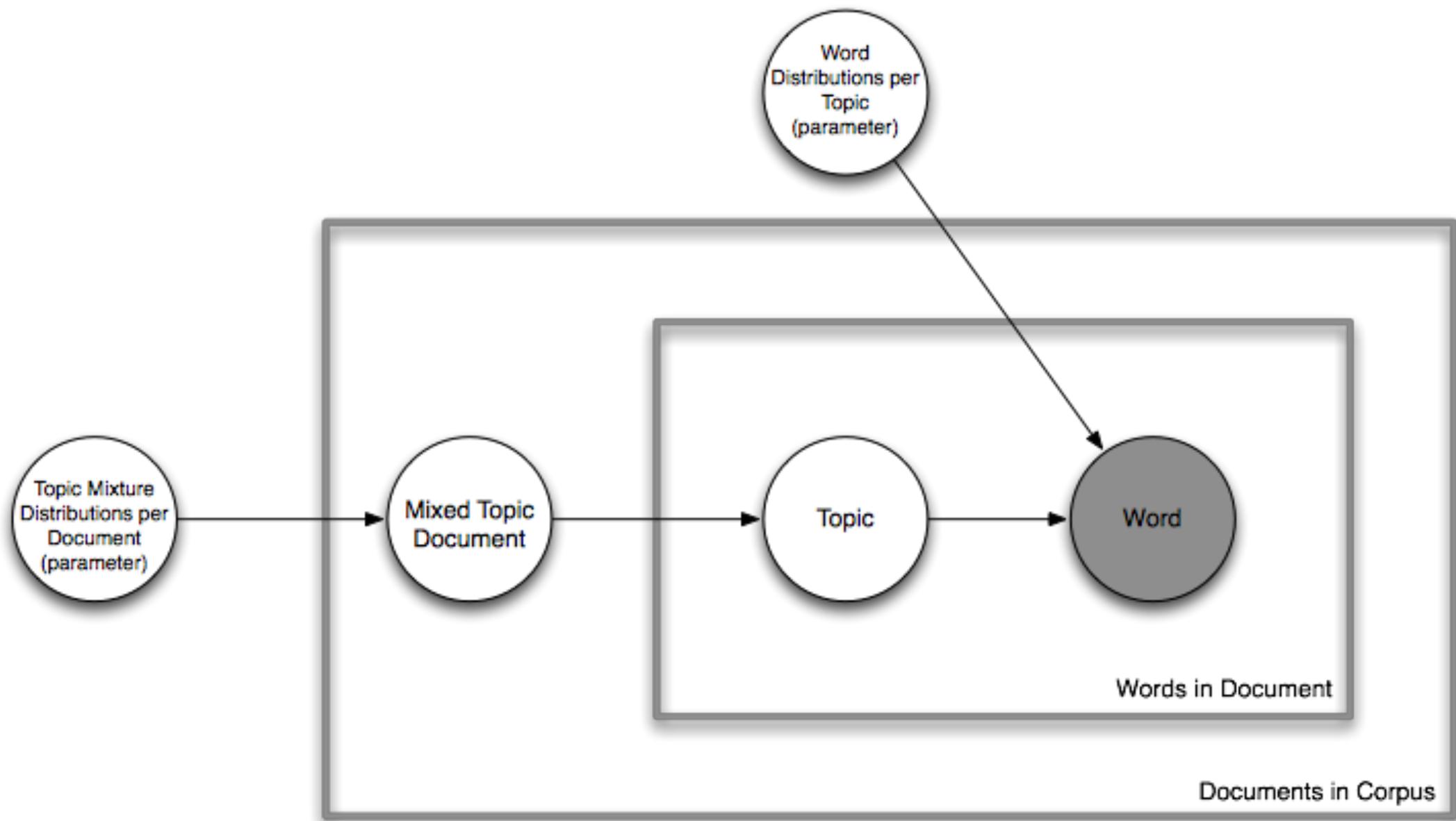
Decide on the number of words N the document will have (say, according to a Poisson distribution).

Choose a topic mixture for the document (according to a Dirichlet distribution over a fixed set of K topics). For example, assuming that we have the two food and cute animal topics above, you might choose the document to consist of 1/3 food and 2/3 cute animals.

Generate each word w_i in the document by:

First picking a topic: for example, you might pick the food topic with 1/3 probability and the cute animals topic with 2/3 probability).

Using the topic to generate the word itself. For example, if we selected the food topic, we might generate the word “broccoli” with 30% probability, “bananas” with 15% probability, and so on.



LDA is a generative model

Pick 5 to be the number of words in D.

Decide that D will be 1/2 about food and 1/2 about cute animals.

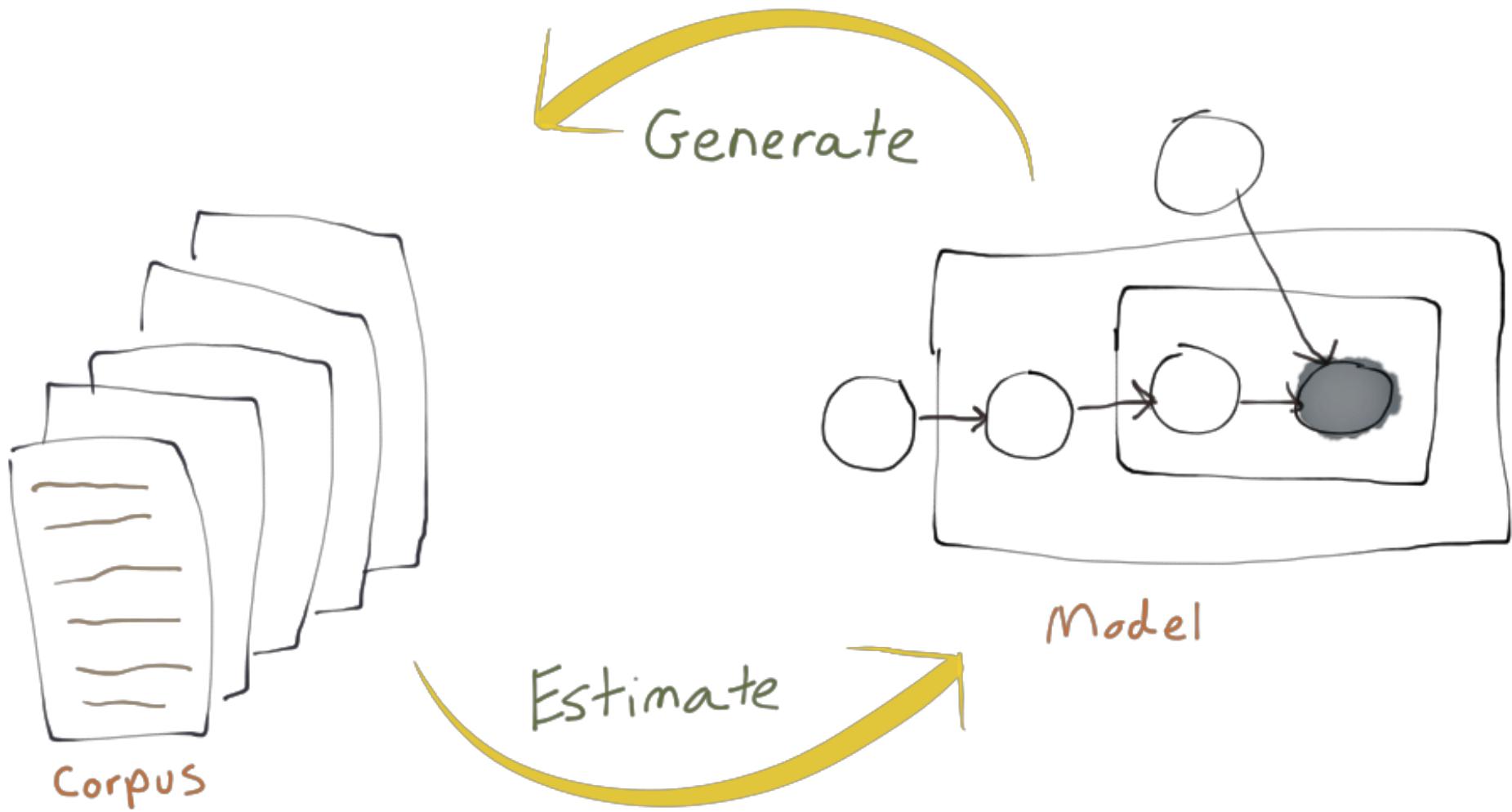
Pick the first word to come from the food topic, which then gives you the word “broccoli”.

Pick the second word to come from the cute animals topic, which gives you “panda”.

Pick the third word to come from the cute animals topic, giving you “adorable”.

Pick the fourth word to come from the food topic, giving you “cherries”.

Pick the fifth word to come from the food topic, giving you “eating”.



Learning

Go through each document, and randomly assign each word in the document to one of the K topics.

Notice that this random assignment already gives you both topic representations of all the documents and word distributions of all the topics (albeit not very good ones).

So to improve on them, for each document d:

Go through each word w in d:

And for each topic t, compute two things: 1) $p(\text{topic } t \mid \text{document } d)$ = the proportion of words in document d that are currently assigned to topic t, and 2) $p(\text{word } w \mid \text{topic } t)$ = the proportion of assignments to topic t over all documents that come from this word w. Reassign w a new topic, where we choose topic t with probability $p(\text{topic } t \mid \text{document } d) * p(\text{word } w \mid \text{topic } t)$

In other words, in this step, we're assuming that all topic assignments except for the current word in question are correct, and then updating the assignment of the current word using our model of how documents are generated.

After repeating the previous step a large number of times, you'll eventually reach a roughly steady state where your assignments are pretty good. So use these assignments to estimate the topic mixtures of each document (by counting the proportion of words assigned to each topic within that document) and the words associated to each topic (by counting the proportion of words assigned to each topic overall).