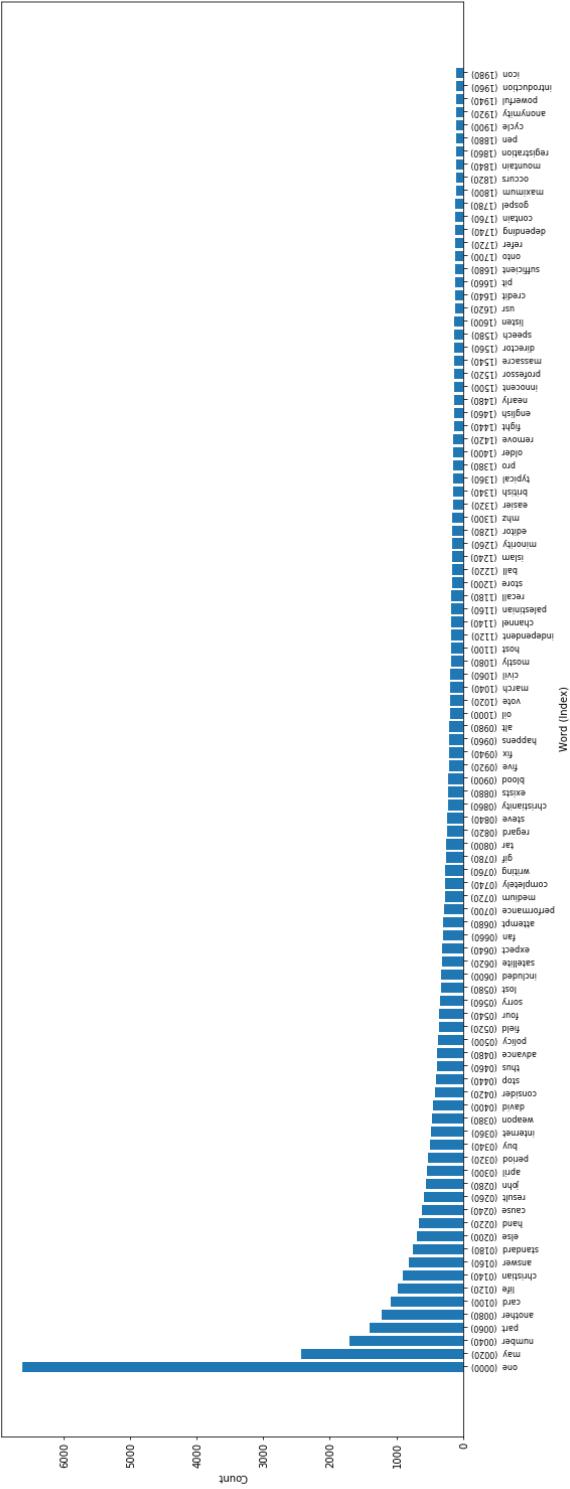
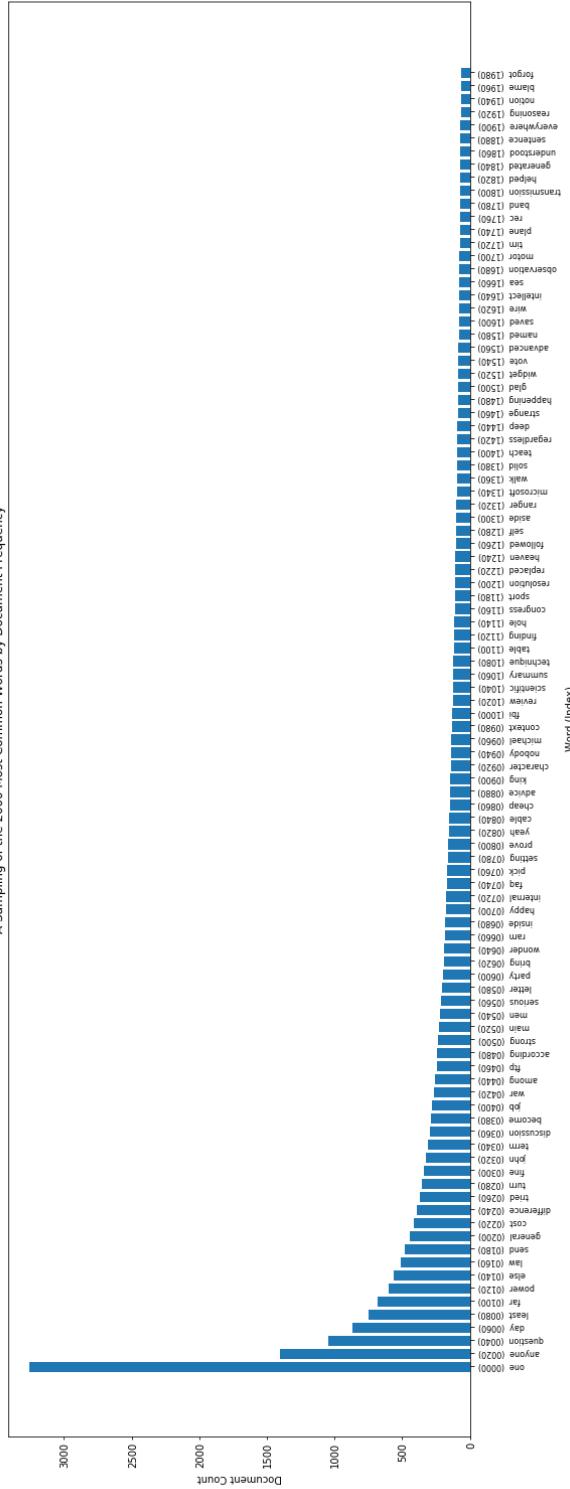


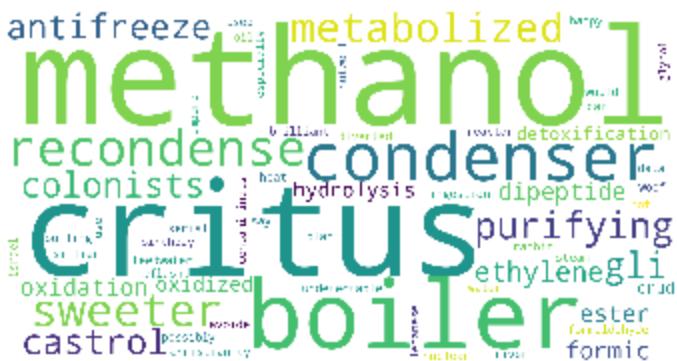
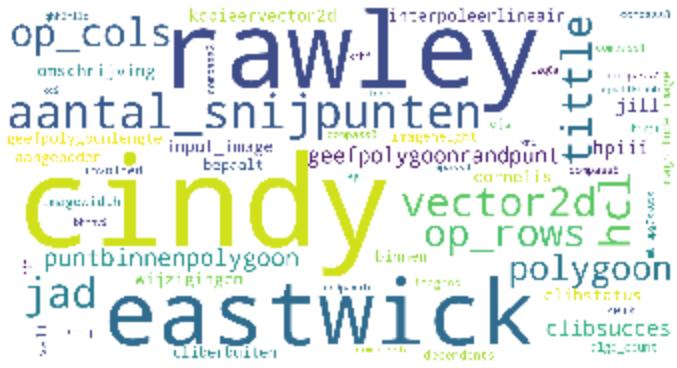
A Sampling of the 2000 Most Common Words in the Dataset



A Sampling of the 2000 Most Common Words by Document Frequency



partition demonstrates charismatic brotherhood salute  
netanyahu promotes toro lust like never seen  
probranham syus truly diversity  
particularism tanta katrina hesh video  
cnn branham tanta diversity  
particularism tanta katrina hesh video  
missoula abilene alerted panache  
binyamin hultman mistress roa murine  
reservist cocker gowan rathole weymouth  
naval scuffling bacon foreclose parkersburg  
cocker gowan rathole neville many things

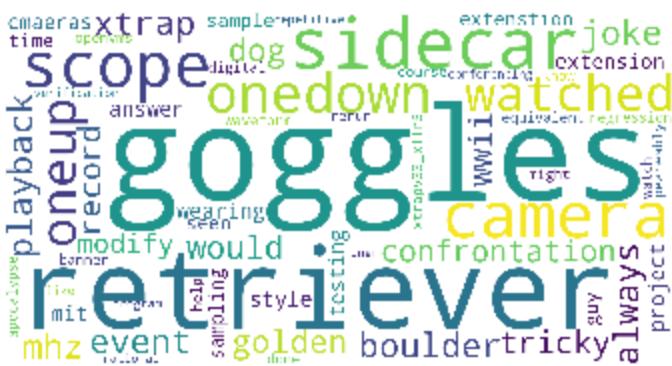
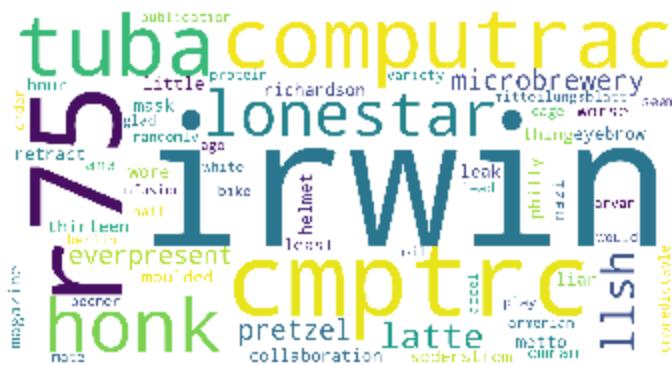


powerkeyinhis m2c excedrin  
backorder sophisicated did1  
sittler imaq  
cartilage vaive  
dalvahn tis rayssd  
linus synergism glenwood

christie recension  
bilinsky kill e esse  
claromontanus bezae reponded

same name sparcclassic  
ctrltest annotation magnifing  
bleah kuiper s  
fairings dialogshell steinn  
bogota bwc karr  
levite shafting yip  
karlamf qb1  
freehand sharen

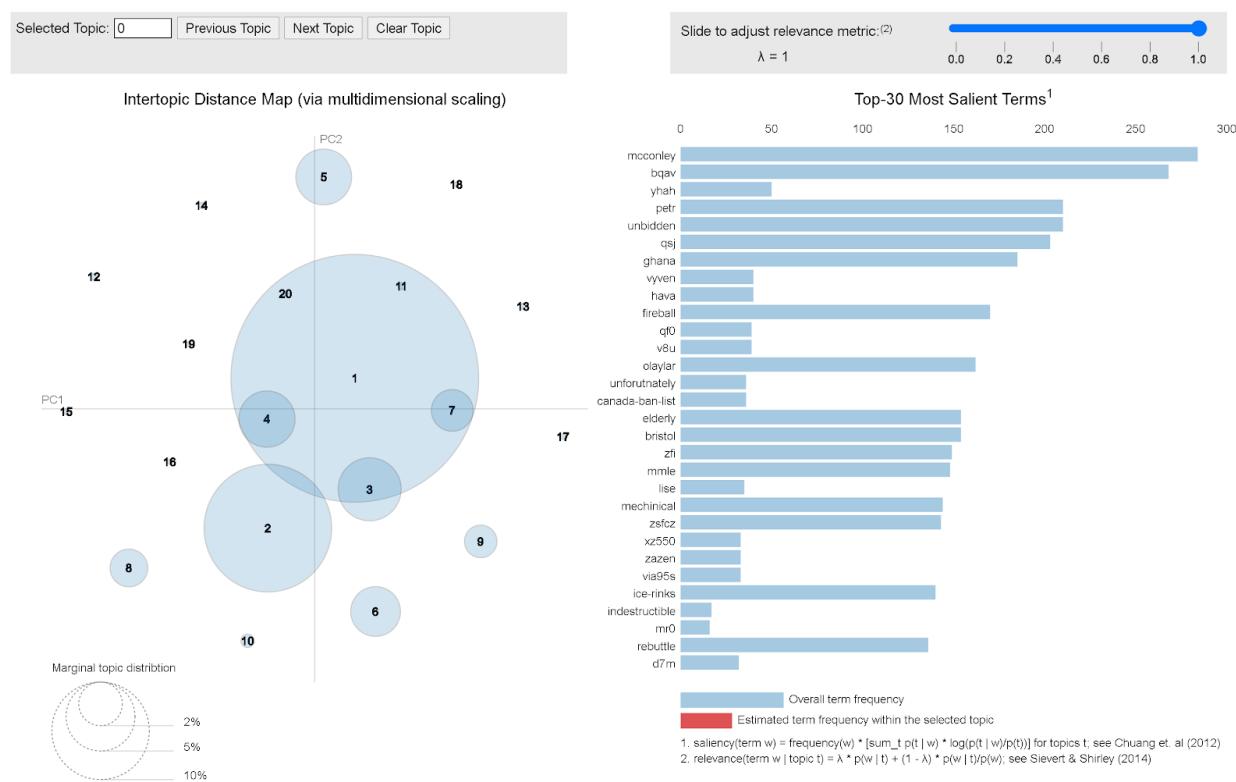
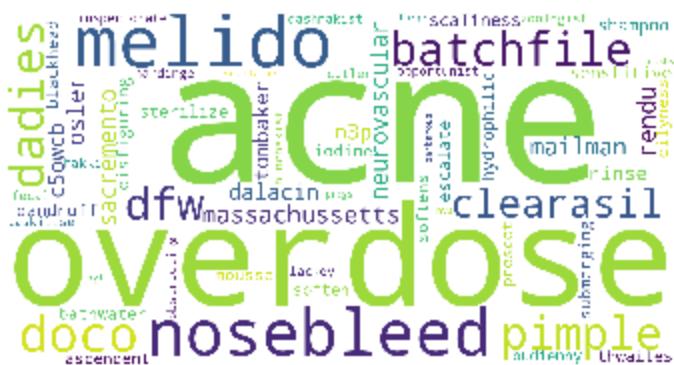
callsign speedisk  
rectum db1spaced  
blunts hite  
nainamo psychoactive malaspina  
rigidly objectbase  
psychoactive



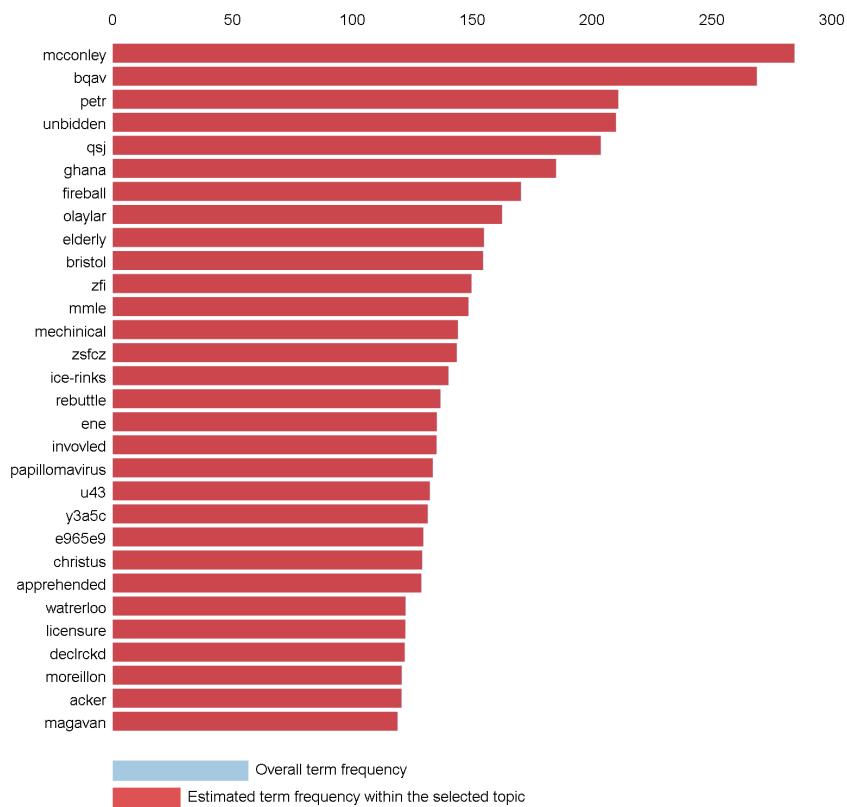
databoss ship timeslip luna cy sleep express modem my system aftwer stick located twin sale maintained spring sport contact type deck cover low legacy automatic upshifts quick sony pedal rev frame stuff econo gear face throttle downshift jurgen whadda recites durant genashor hypoglycemic bloodcount glycemic vit chevrolet indy mazzah nearly whatching胰creatic lois

timeslip armenian kurdish callback timesheet bit software enabler driver multiscope alatlevler agdam use cache fation cyringenocide turek zor lech install boss kedwell superbases stick located twin sale maintained spring sport contact type deck cover low legacy automatic upshifts quick sony pedal rev frame stuff econo gear face throttle downshift jurgen whadda recites durant genashor hypoglycemic bloodcount glycemic vit chevrolet indy mazzah nearly whatching胰creatic lois

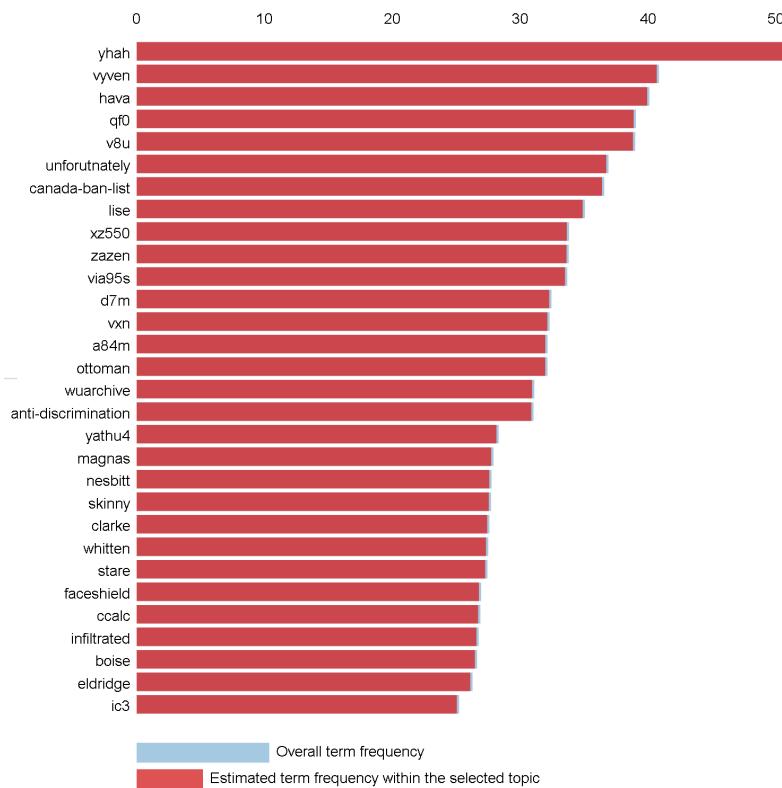
jurgen whadda recites durant genashor hypoglycemic bloodcount glycemic vit chevrolet indy mazzah nearly whatching胰creatic lois



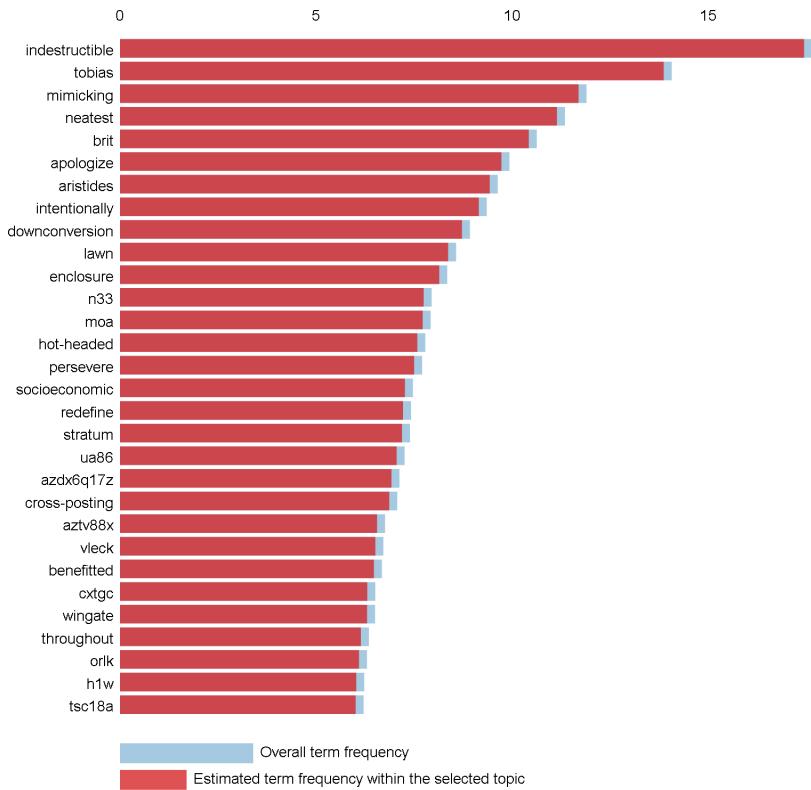
### Top-30 Most Relevant Terms for Topic 1 (64.6% of tokens)



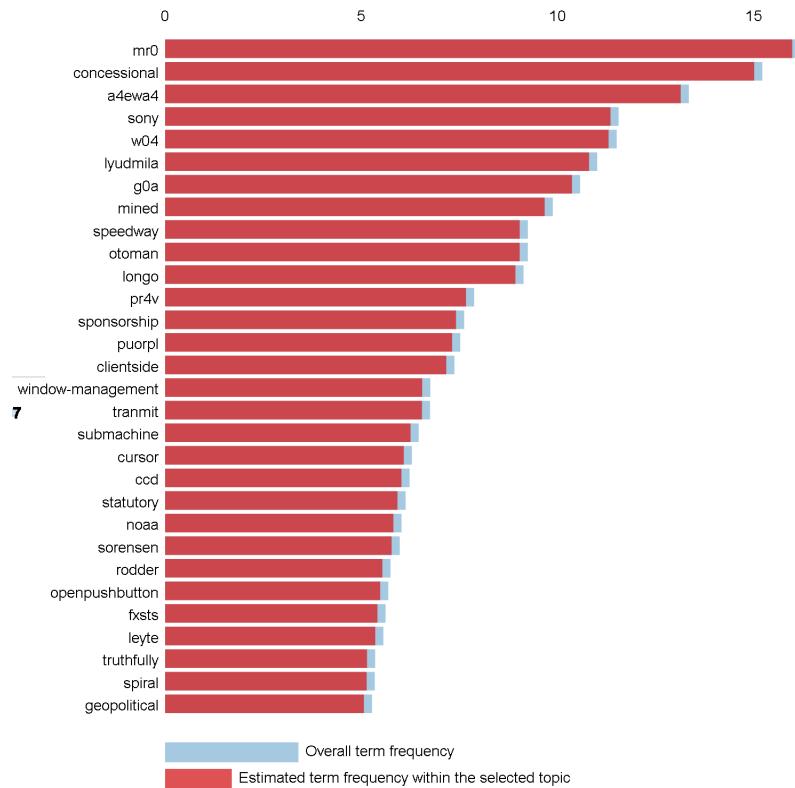
### Top-30 Most Relevant Terms for Topic 2 (17.1% of tokens)



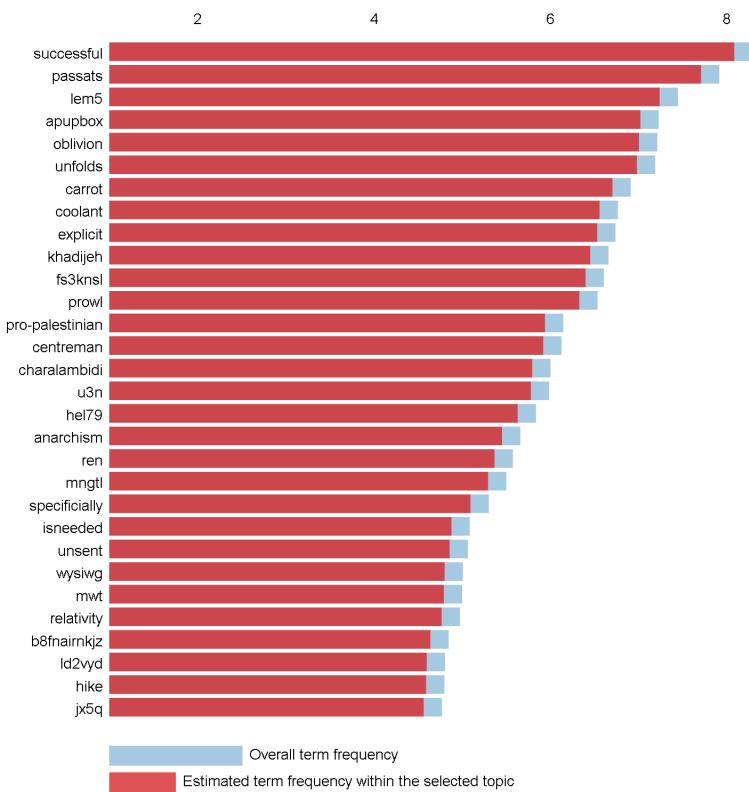
### Top-30 Most Relevant Terms for Topic 3 (4.2% of tokens)



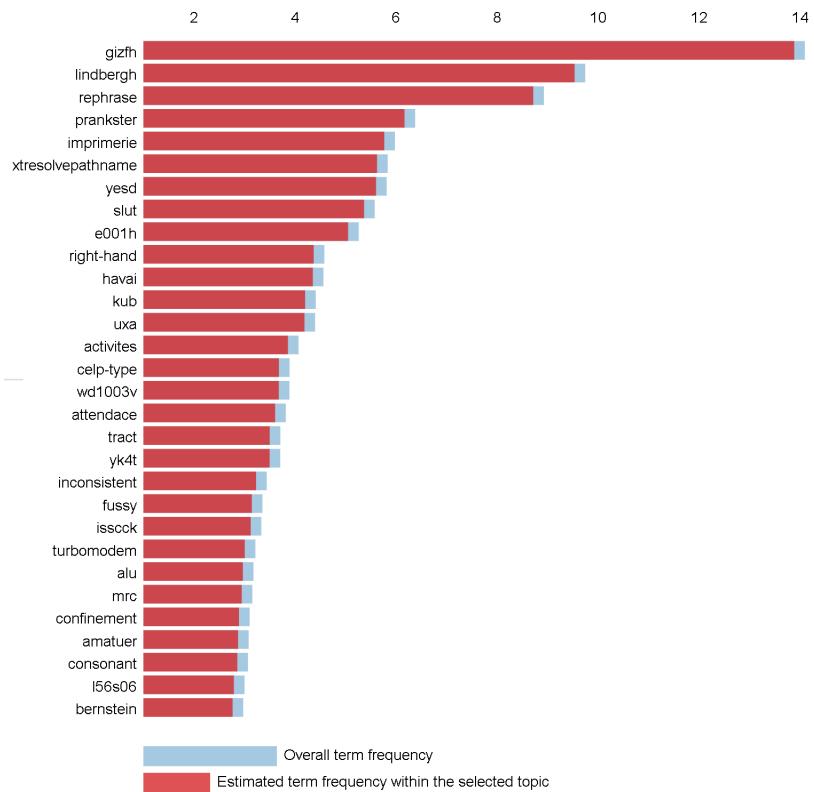
Top-30 Most Relevant Terms for Topic 4 (3.3% of tokens)



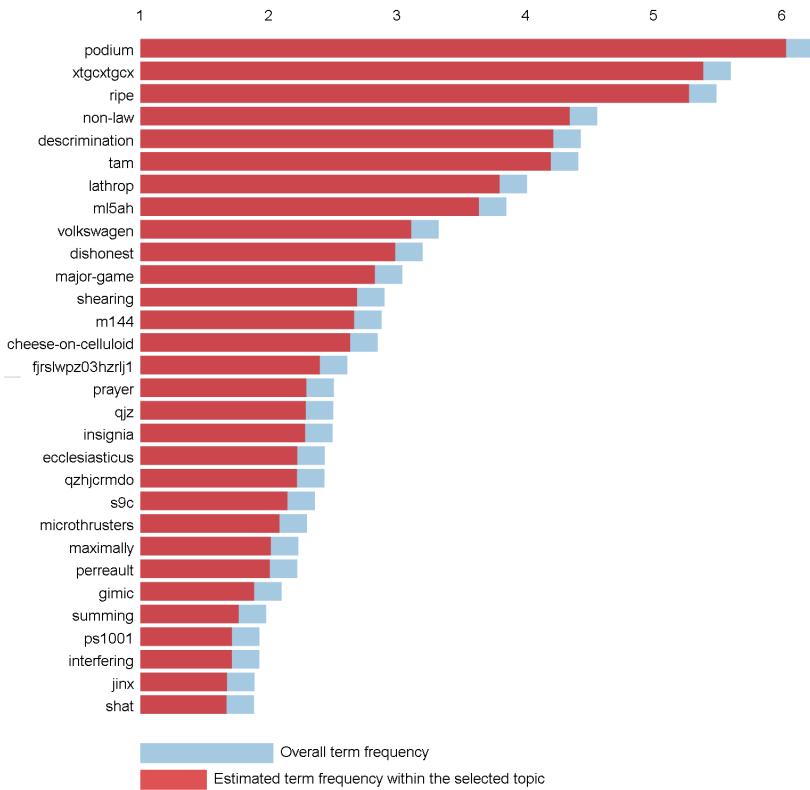
### Top-30 Most Relevant Terms for Topic 5 (3.3% of tokens)



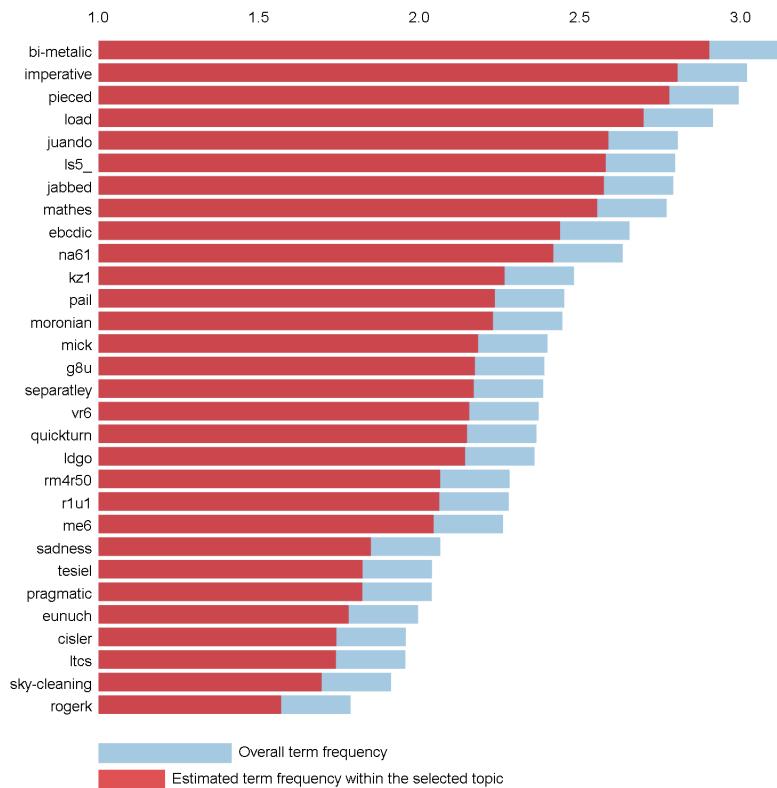
Top-30 Most Relevant Terms for Topic 6 (2.6% of tokens)



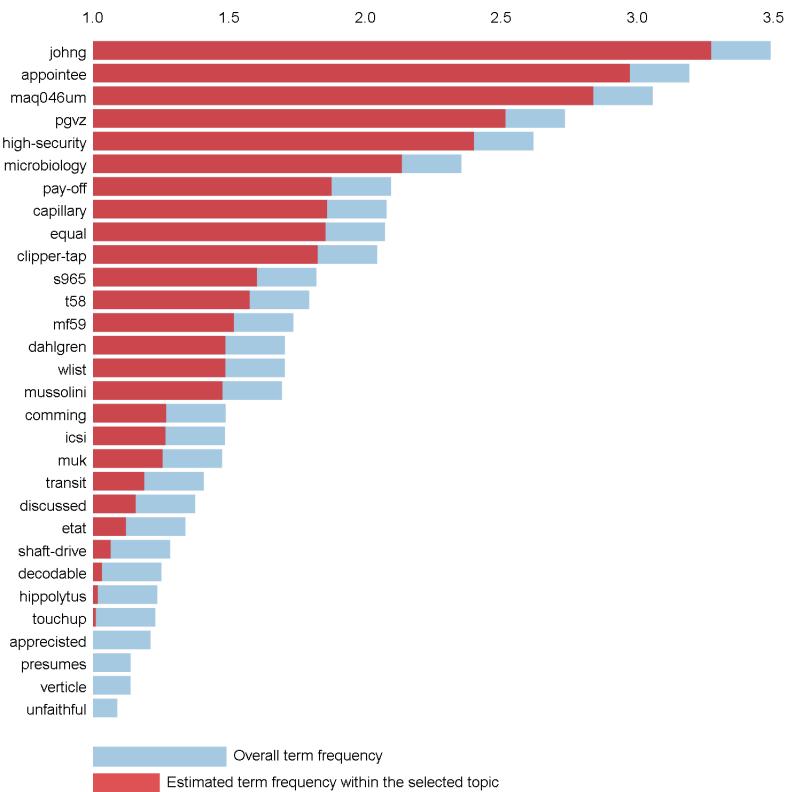
Top-30 Most Relevant Terms for Topic 7 (1.8% of tokens)



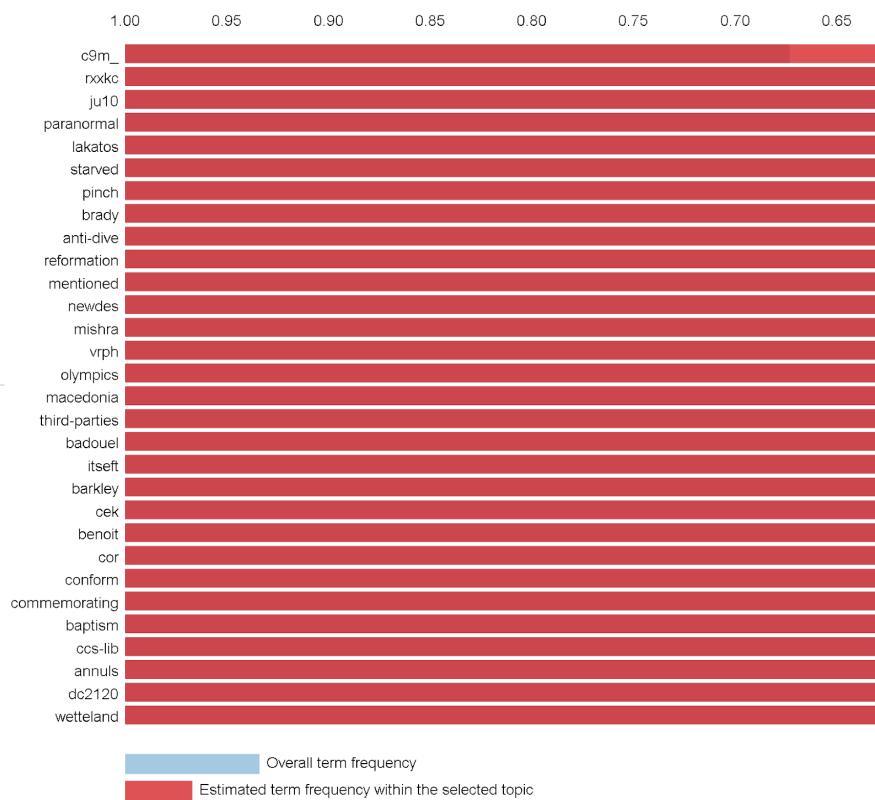
### Top-30 Most Relevant Terms for Topic 8 (1.5% of tokens)



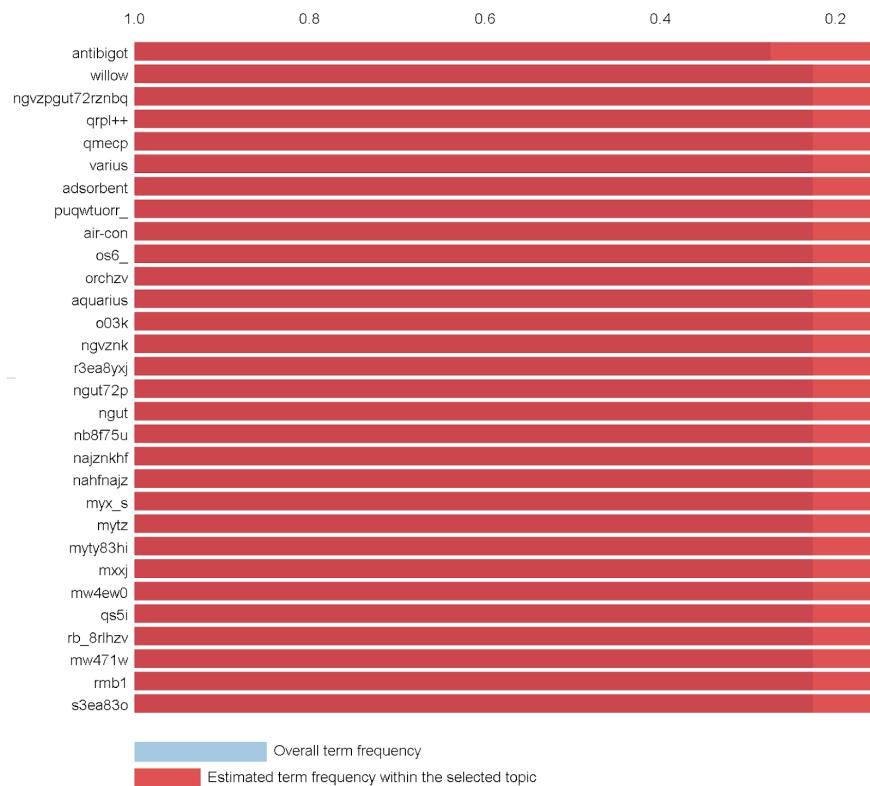
### Top-30 Most Relevant Terms for Topic 9 (1.1% of tokens)



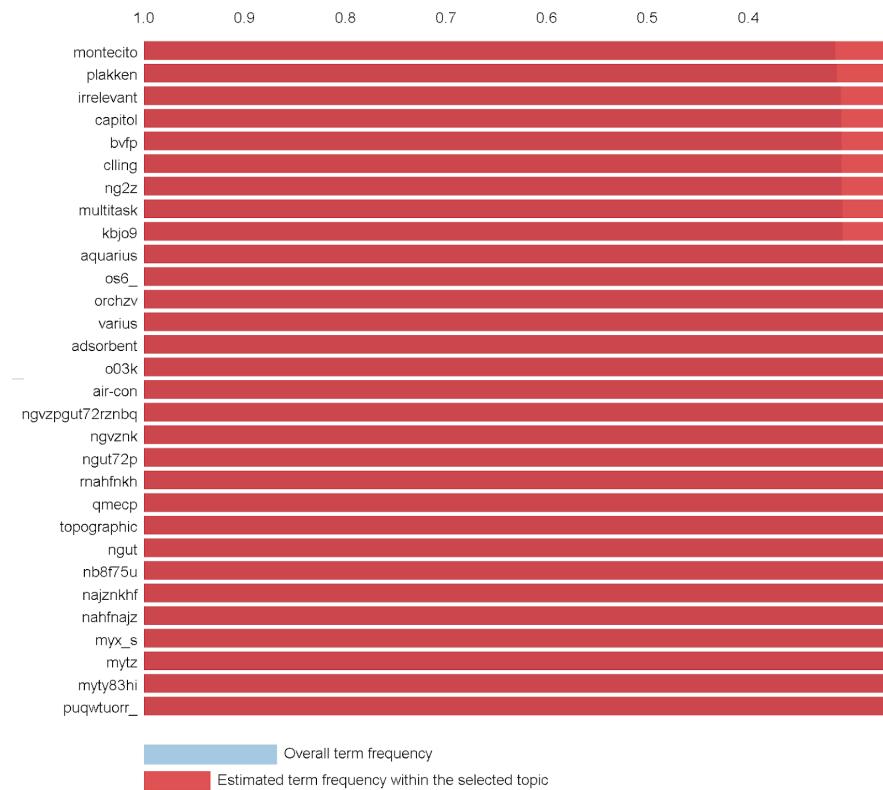
### Top-30 Most Relevant Terms for Topic 10 (0.2% of tokens)



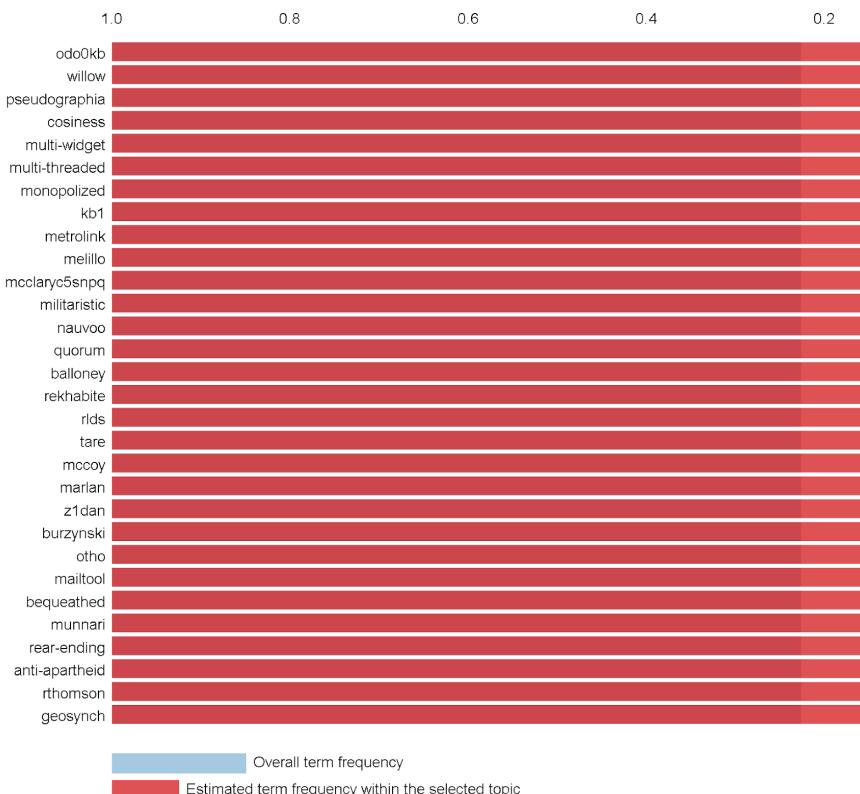
### Top-30 Most Relevant Terms for Topic 11 (0.1% of tokens)



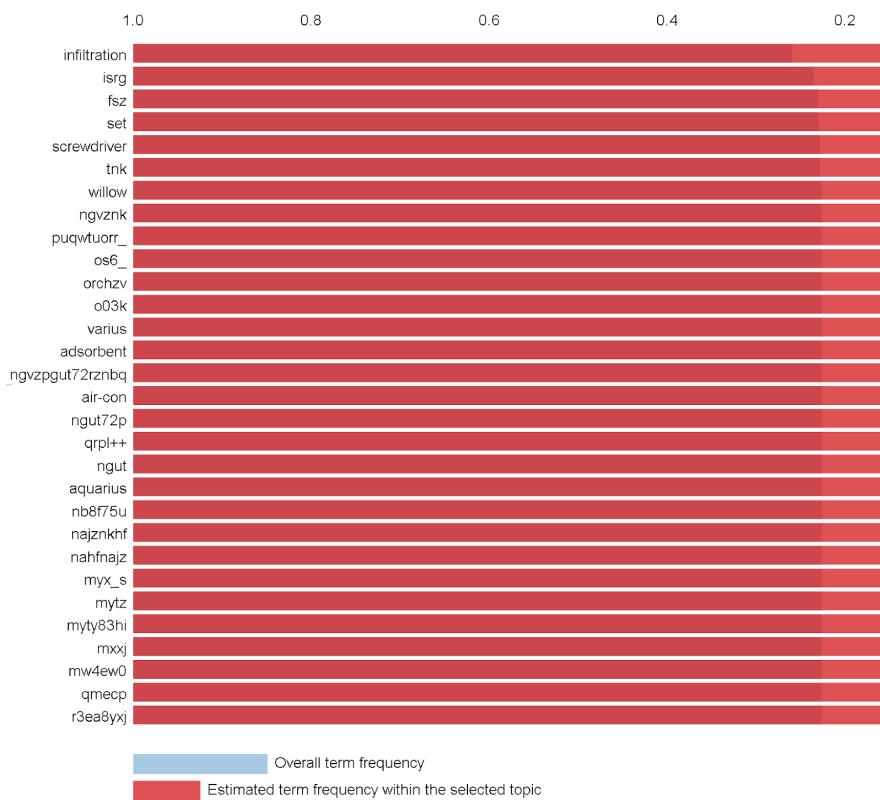
### Top-30 Most Relevant Terms for Topic 12 (0% of tokens)



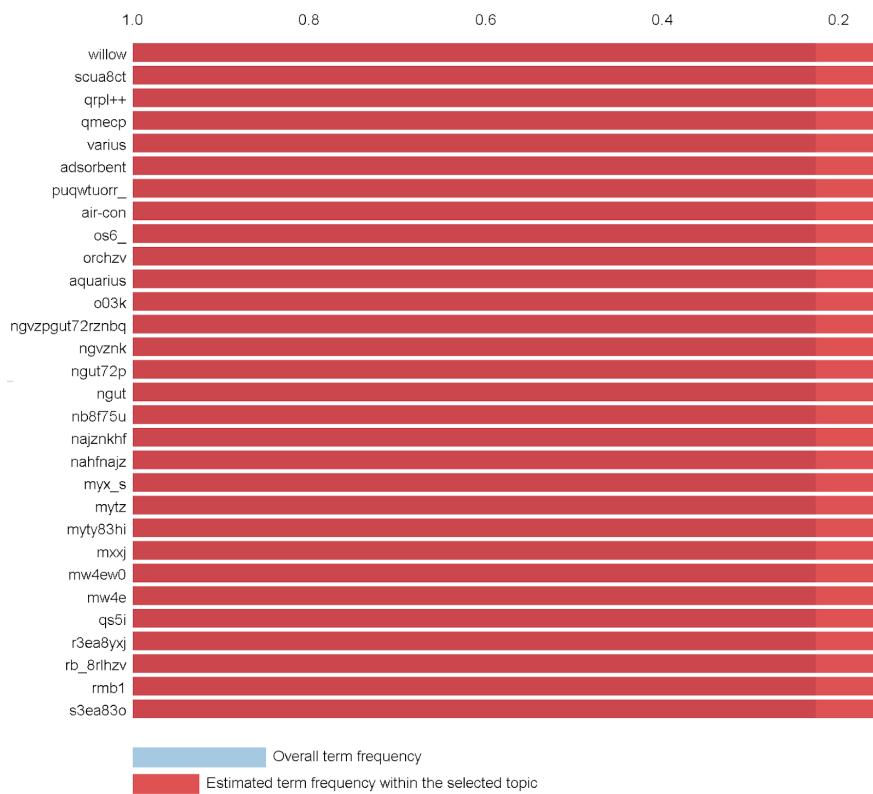
### Top-30 Most Relevant Terms for Topic 13 (0% of tokens)



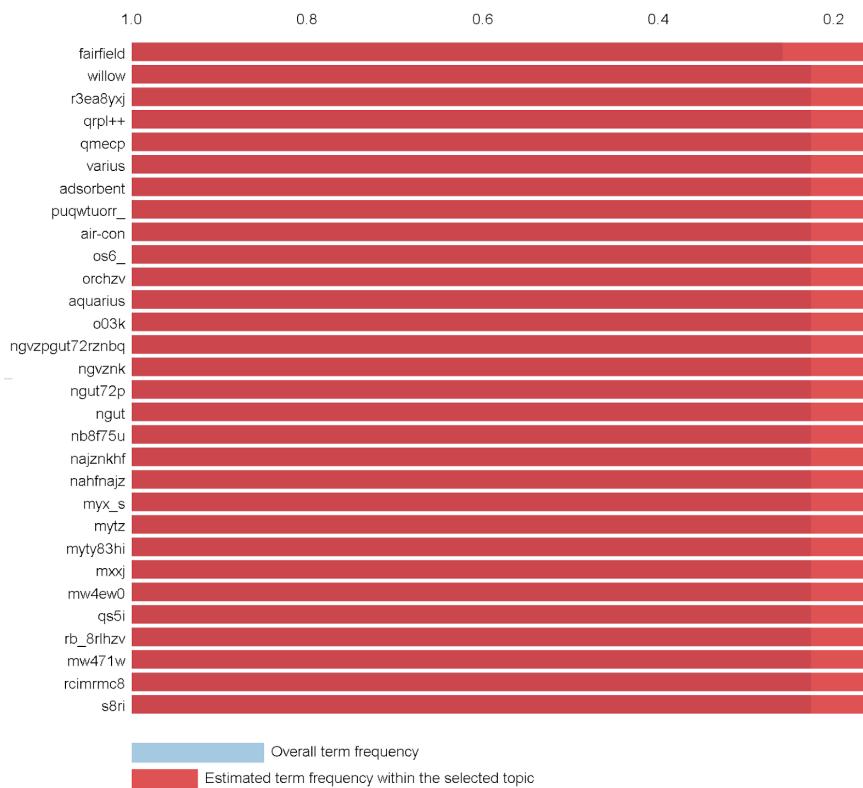
### Top-30 Most Relevant Terms for Topic 14 (0% of tokens)



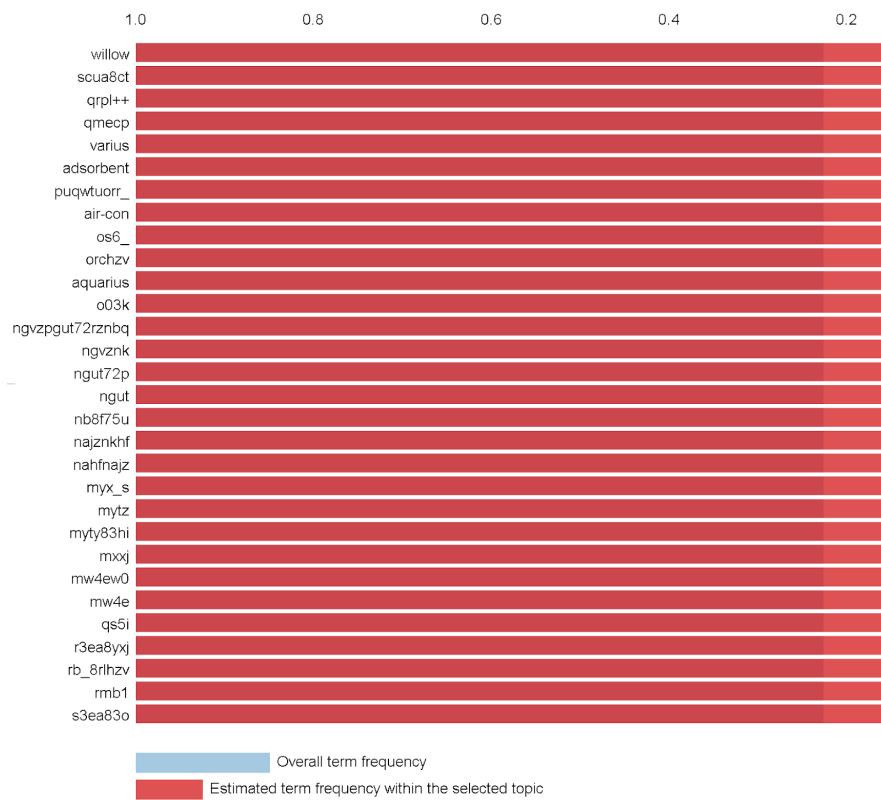
### Top-30 Most Relevant Terms for Topic 15 (0% of tokens)



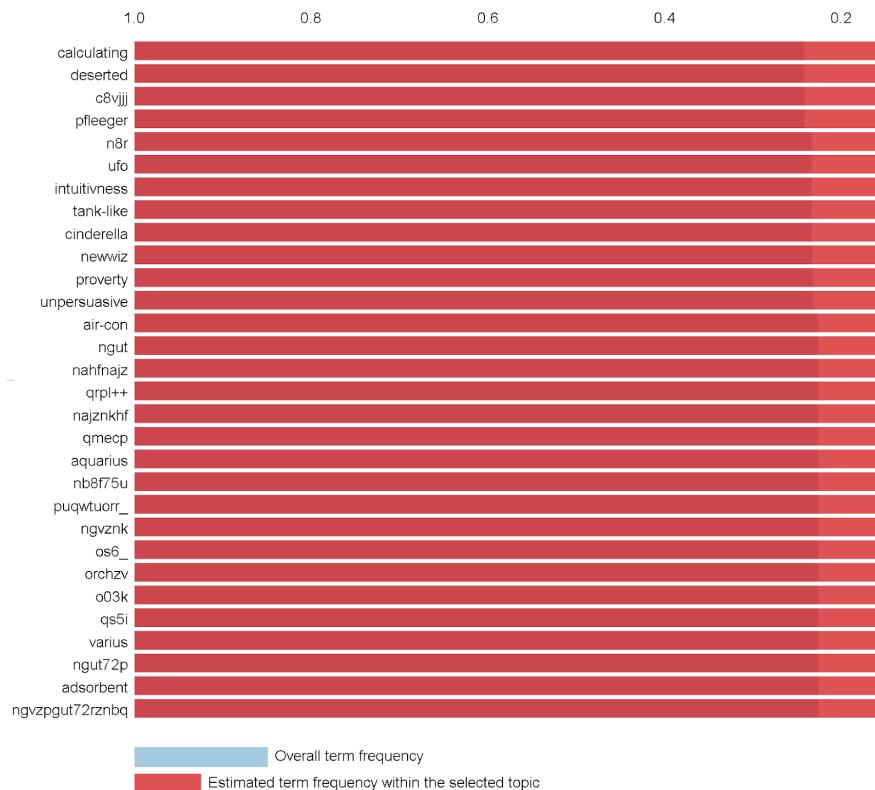
### Top-30 Most Relevant Terms for Topic 16 (0% of tokens)



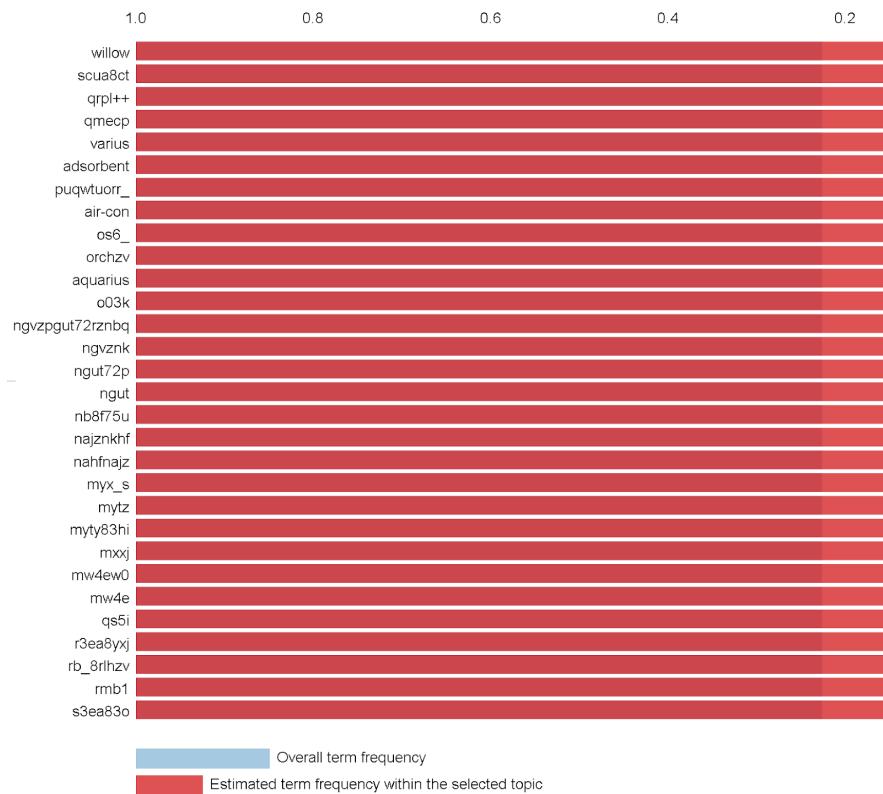
### Top-30 Most Relevant Terms for Topic 17 (0% of tokens)



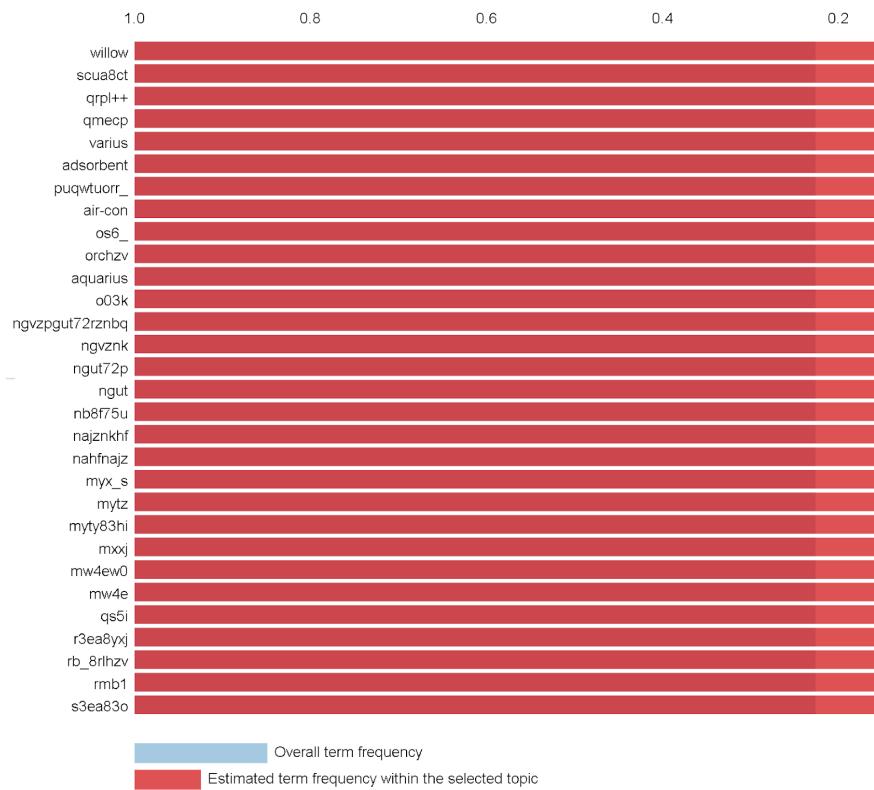
### Top-30 Most Relevant Terms for Topic 18 (0% of tokens)



### Top-30 Most Relevant Terms for Topic 19 (0% of tokens)



### Top-30 Most Relevant Terms for Topic 20 (0% of tokens)



though someone available club excuse anyone  
info power relevant war well  
sent care mail interface direction sure  
bought box know sending sorry  
interior directory trial islamic get  
home thousand round called  
soy died point dala would  
japanese fall amiga  
feature talking remember  
using thousand people obviously  
graphictalking officemuslim

number short force start relationship show used  
moral sexual sex men excuse showing  
know size like whose would ask etc example except  
like size knew expected couple religious  
straight homosexual gay life up get  
homosexual escape done nuclear job much  
marriage male thinks optimism  
outside night alt family experiment break alternative lack

technology escrow number communication conversation  
private algorithm crypto warrant usc device law clinton  
processor privacy data court enforcement agency message search  
block rsa administration telephone random proposal  
rsa public eff security voice scheme  
system bitcipher wiretap cryptography encrypted  
dev cipher clipper serial pgp  
trust public administration encrypted traffic  
bitcipher wiretap secret strong nist, legal  
key encryption government

armenia genocide germany ottoman  
serdar million azerbaijan attack people government  
massacre republic hitler army terrorist source  
turkish azeri argic german  
greek nazi turkey  
plane russian today greece killing town killed buried  
troop island april russia army policy  
nation turkish history  
armenian jew war  
muslim azerbaijan soldier officer village population  
turk village

image even without text  
ftp component electrical  
get nature hank site one  
input neural project weight  
somewhere turn net available  
controller pitt ram bar hello  
terminal software cable external  
spacecraft plan compromise gif  
radio converts sound compatible  
scratches attention like

advance anybody directory graphic time void  
matrix problem run image email  
reclly code.edu pub  
mellin fax version information link application  
send could know looking appreciated info help address display  
appreciated info help address display  
know looking appreciated info help address display  
looking appreciated info help address display  
anyone need ftp list format server using get  
source available color format server using get

seen one people convert wonder application  
version com connect camera error  
like king plug cycle something another problem shareware  
image output choose found domain  
output worth name printer graphic come  
got tell older screen contact luck  
telling current god driver proper come  
laser class port use  
font thanks compare jet print anyone  
window setting net file fit tried thing  
possible resolution

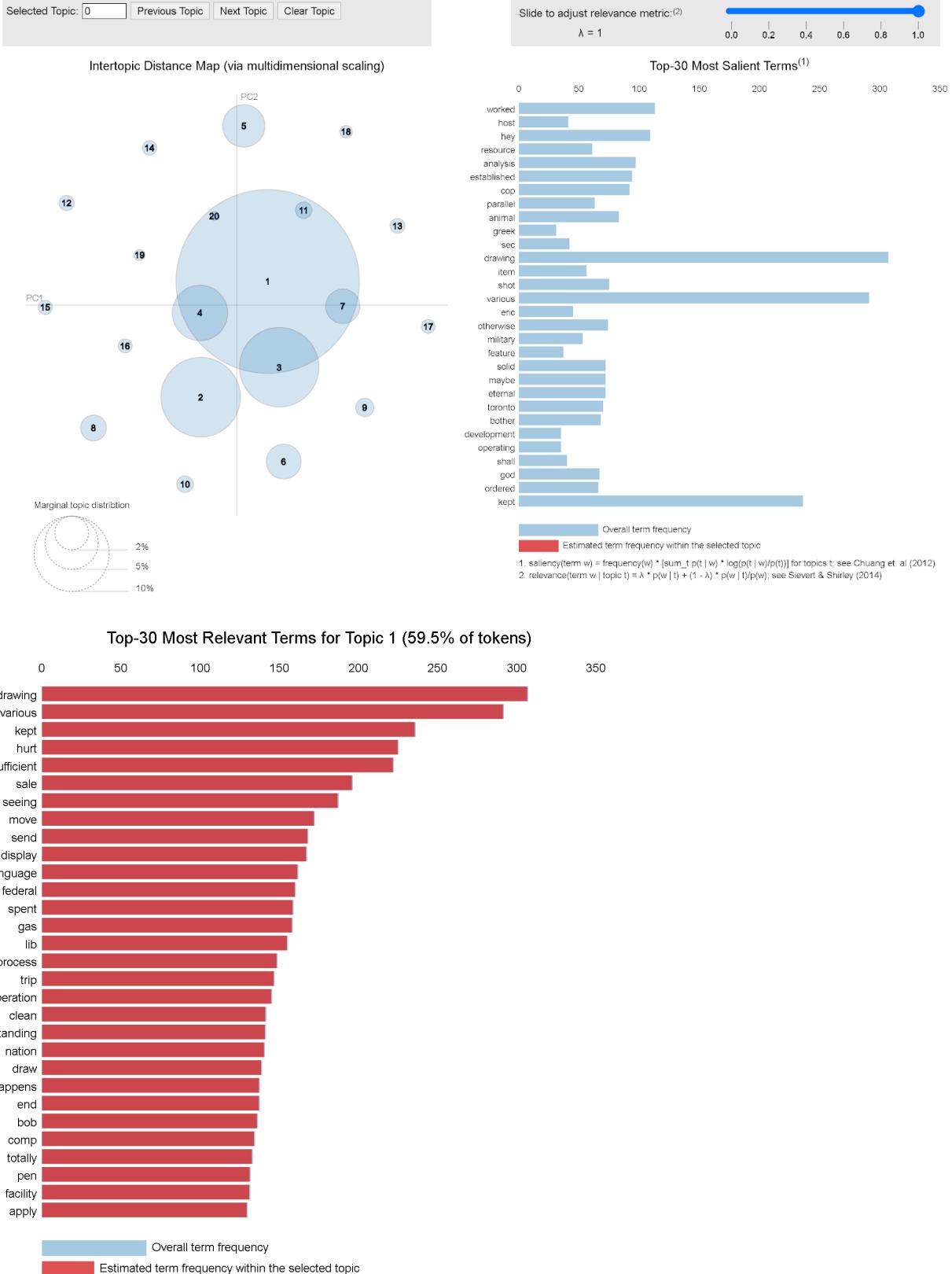
wrote basically monitor  
simms want  
bank tower module job better  
tower memory meg  
power possible mac place shuttle  
bank lead zip killed versatil  
tower cable what ever specific  
shareware window option jumper difference problem  
shareware

new first  
good one people  
thing may like  
think anyone  
point used cone  
use said edu  
work question  
think time  
right two  
right day  
year even  
get know say  
car better many  
post lock case  
way still  
so god also  
problem make  
still since  
script sure  
also well  
since lock  
script never  
make believe  
script probably  
since need  
script could

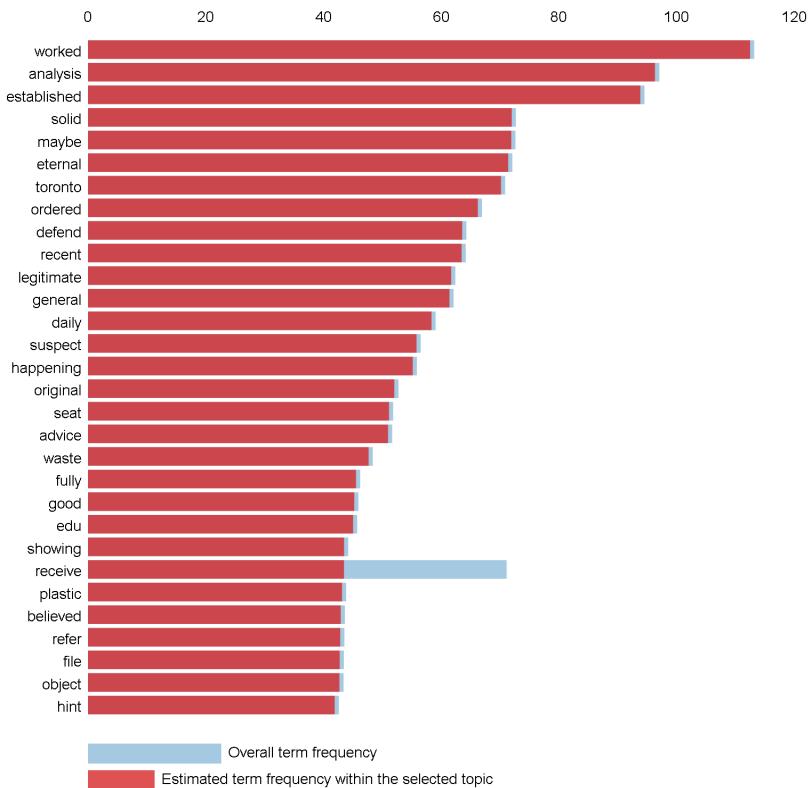
anyone sense buffer car corner window  
regardless answer type vehicle using  
character would could enough gotten  
much control machine export installed host  
wondering variable thanks benefit  
option certainly need either  
left disk write obvious use  
regardless body install problem  
situation xterm require ram box  
larger simms without addition religion  
script door cache

que also library appreciated brother  
greatly knowledge computer help  
street looking seen seen note signal none code  
cram mean speed sorry output  
theory are dad god able need  
assumption voltage father thanks say  
sec right car mile resource  
graphic plus design book student bit  
point information driver mother audio

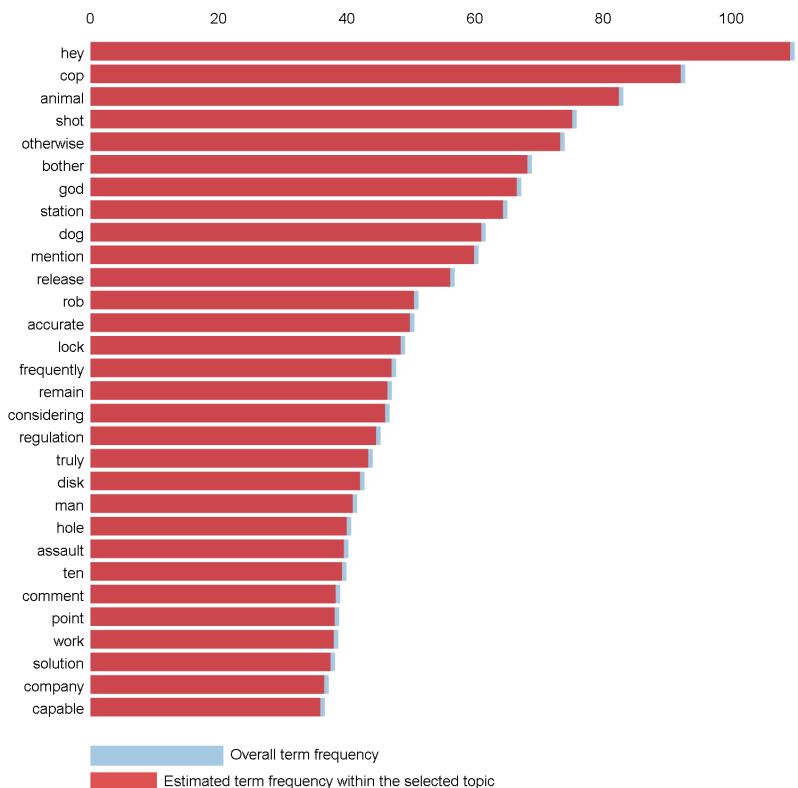
surrender  
effective  
blood pressure  
bank



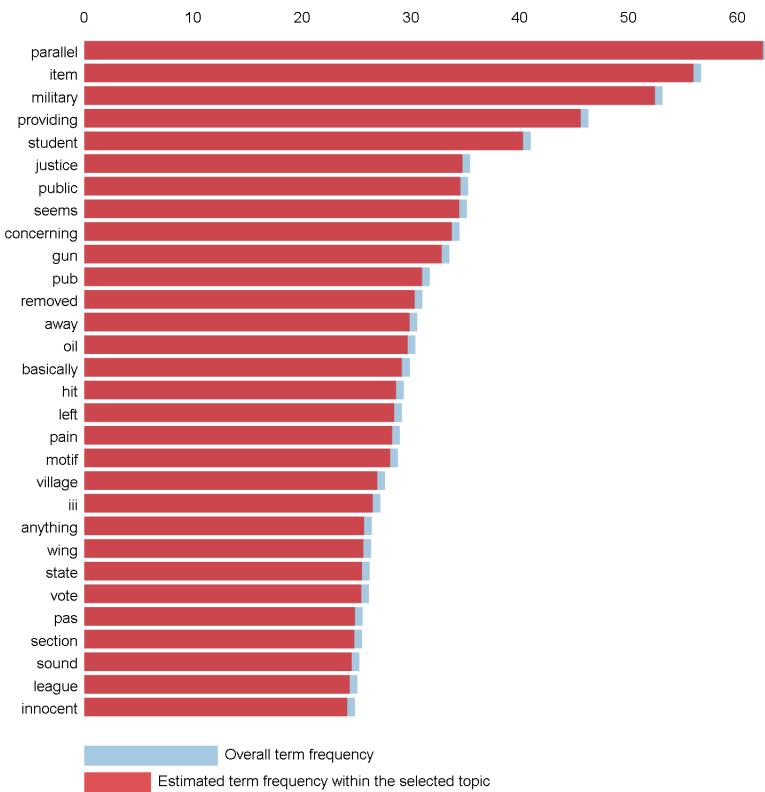
### Top-30 Most Relevant Terms for Topic 2 (11.2% of tokens)



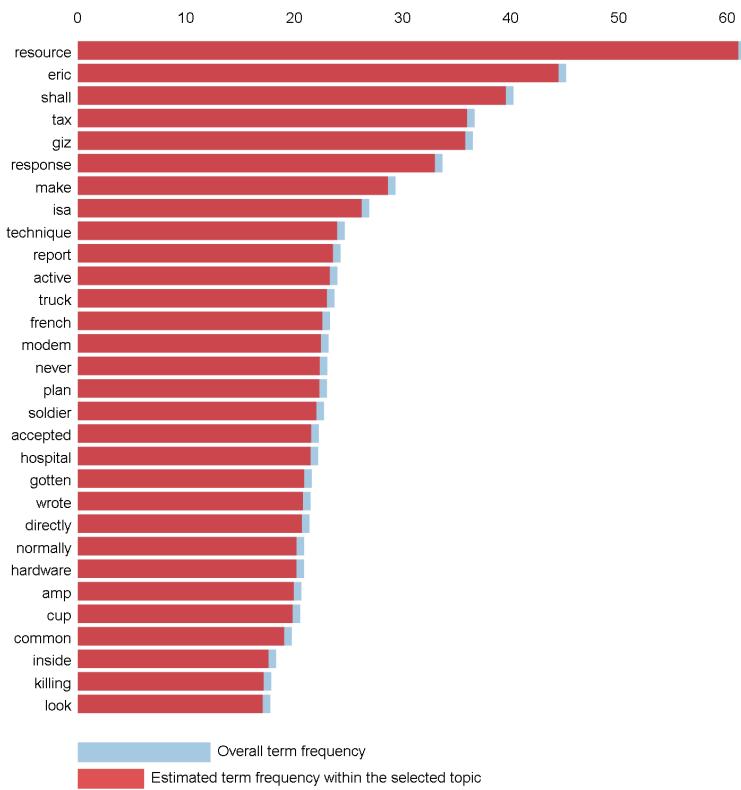
### Top-30 Most Relevant Terms for Topic 3 (11.1% of tokens)



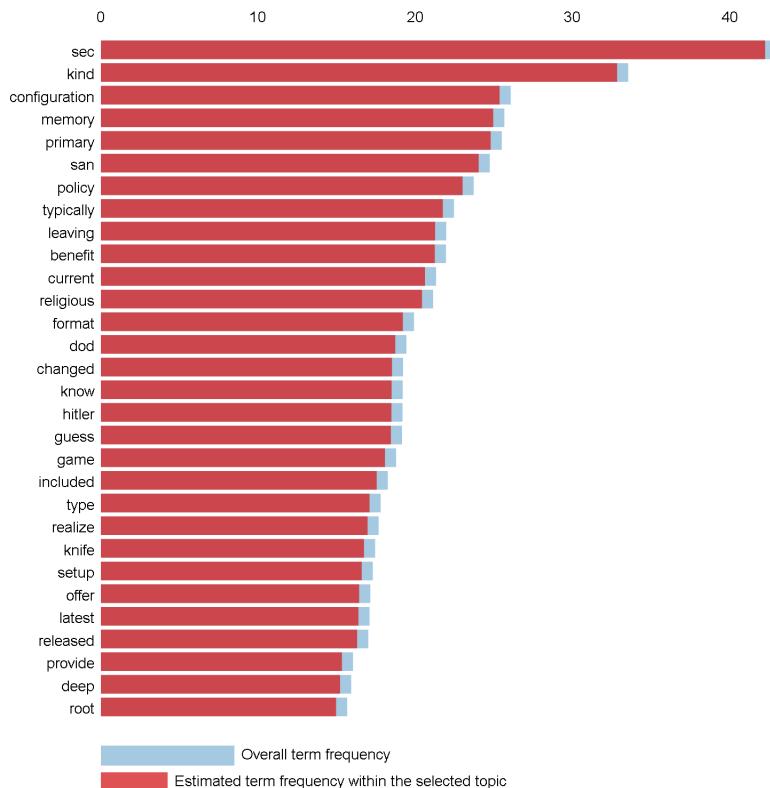
### Top-30 Most Relevant Terms for Topic 4 (5.4% of tokens)



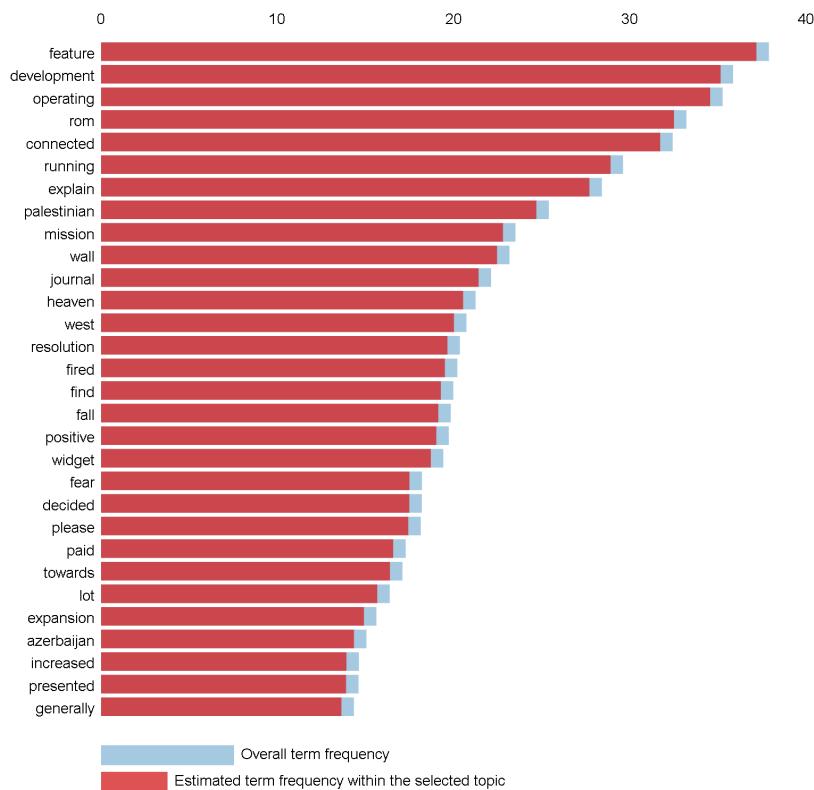
### Top-30 Most Relevant Terms for Topic 5 (3.2% of tokens)



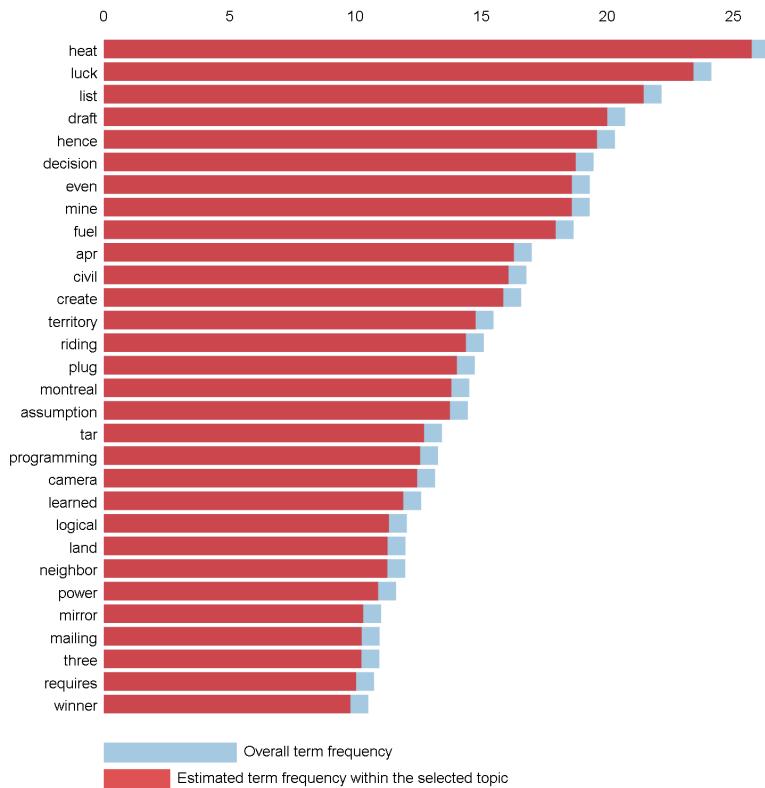
### Top-30 Most Relevant Terms for Topic 6 (2.1% of tokens)



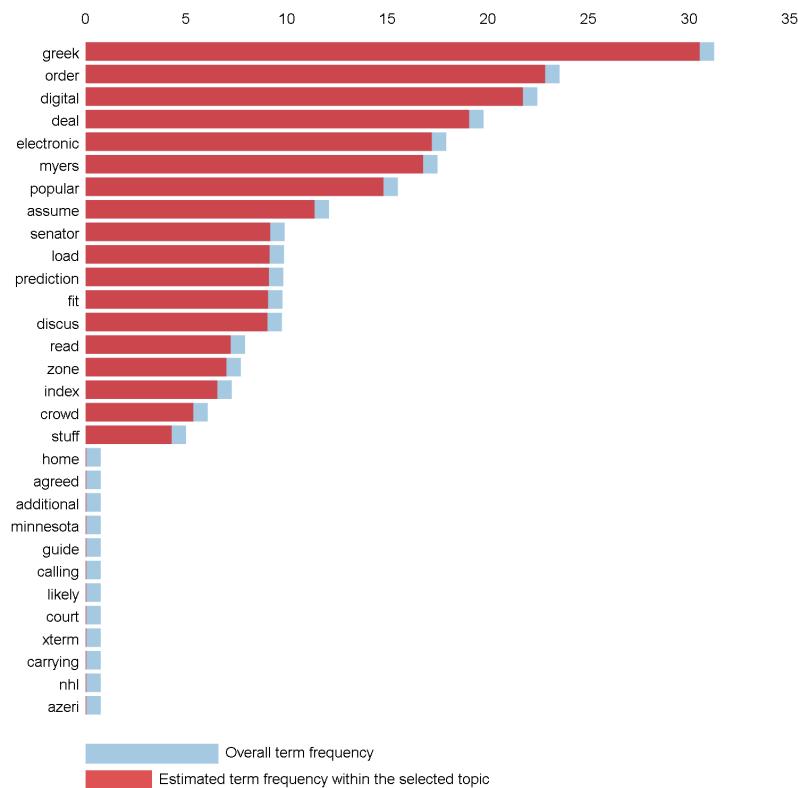
### Top-30 Most Relevant Terms for Topic 7 (2.1% of tokens)



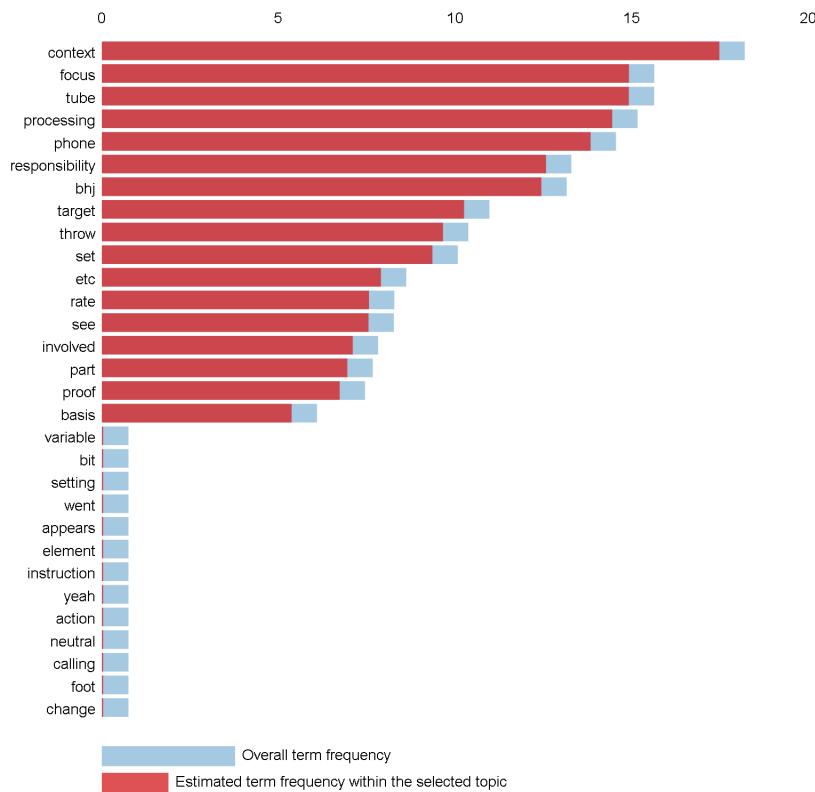
Top-30 Most Relevant Terms for Topic 8 (1.2% of tokens)



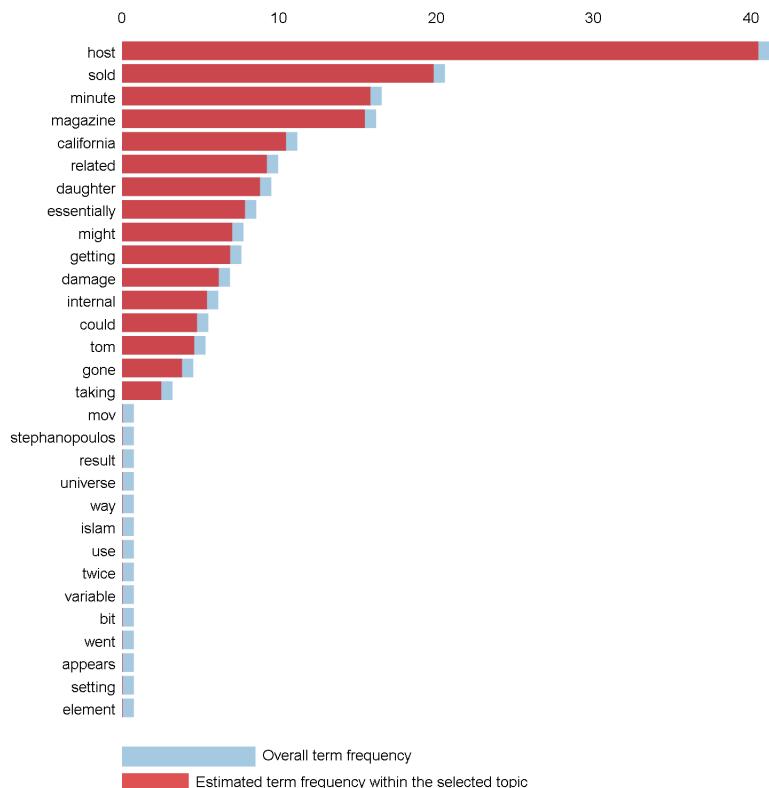
### Top-30 Most Relevant Terms for Topic 9 (0.6% of tokens)



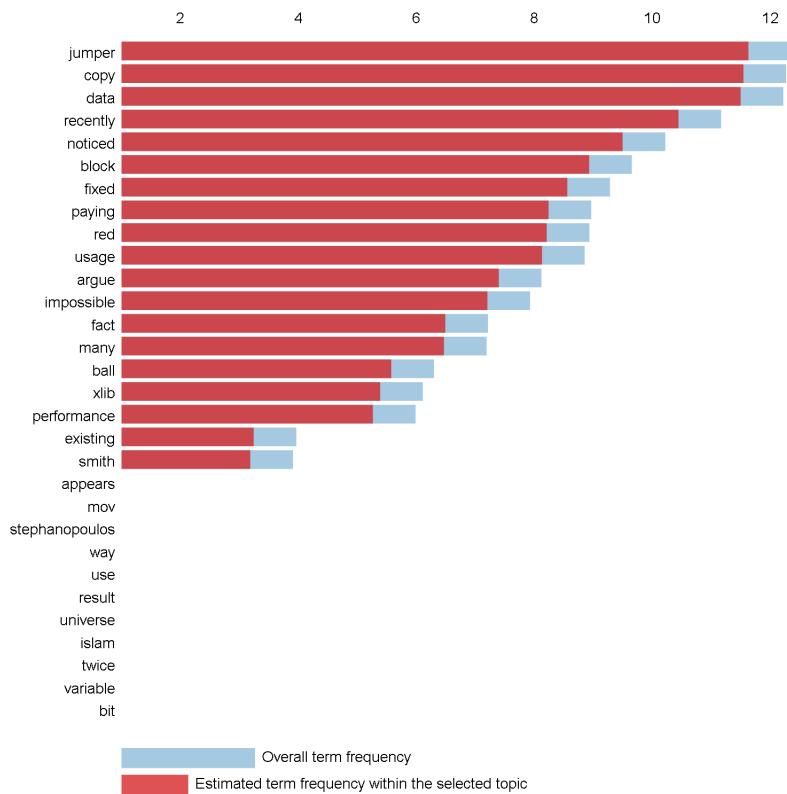
Top-30 Most Relevant Terms for Topic 10 (0.5% of tokens)



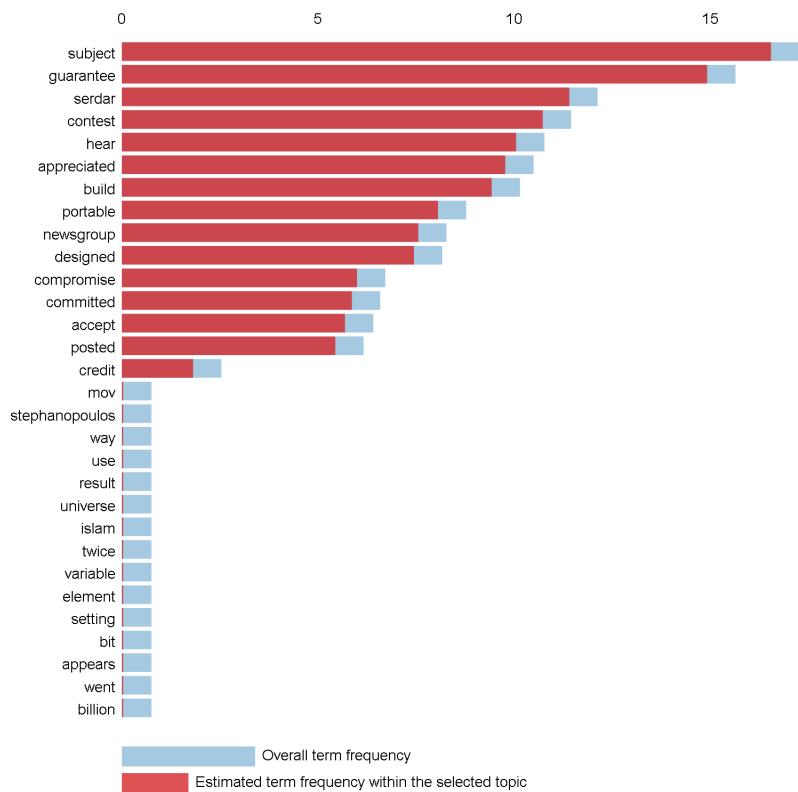
### Top-30 Most Relevant Terms for Topic 11 (0.5% of tokens)



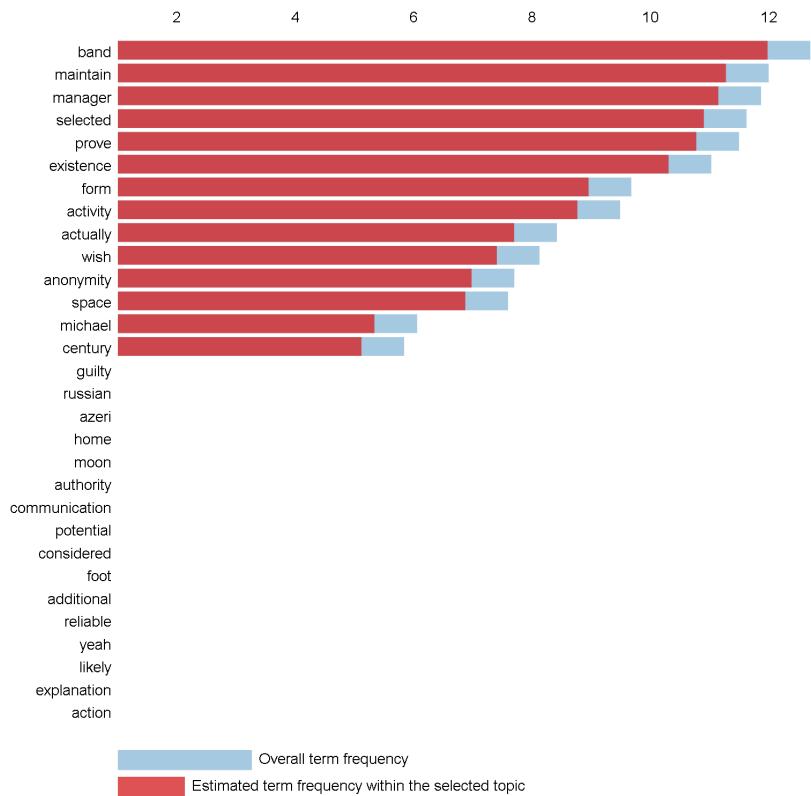
### Top-30 Most Relevant Terms for Topic 12 (0.4% of tokens)



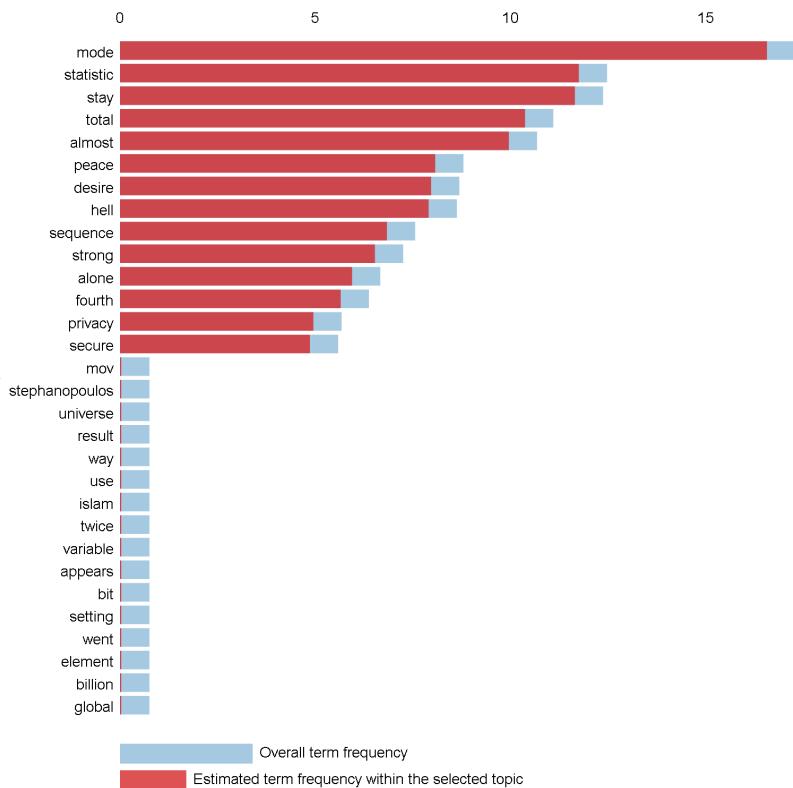
### Top-30 Most Relevant Terms for Topic 13 (0.4% of tokens)



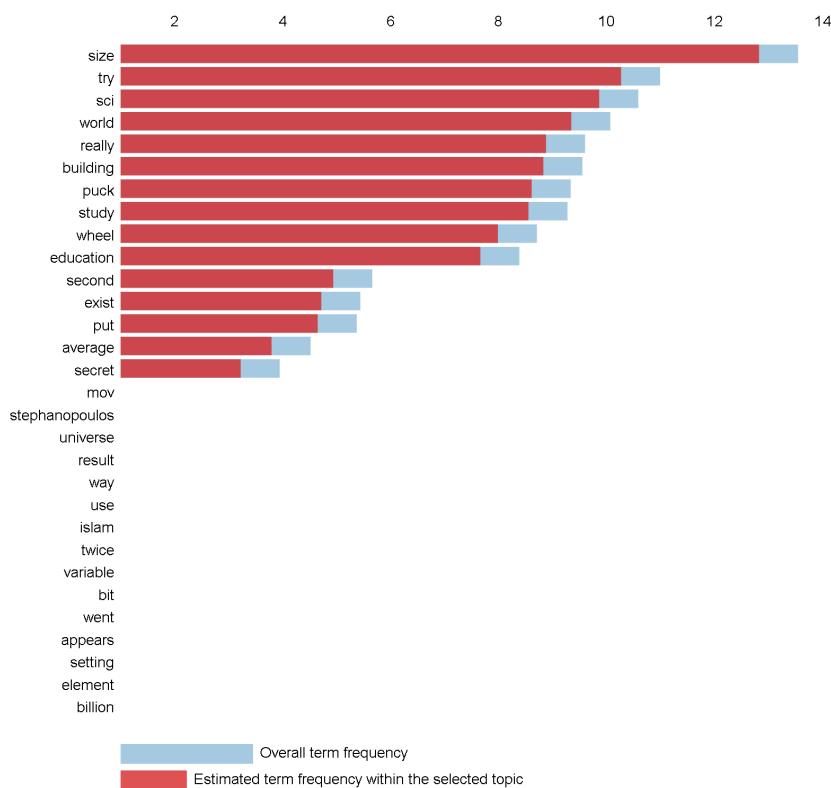
### Top-30 Most Relevant Terms for Topic 14 (0.4% of tokens)



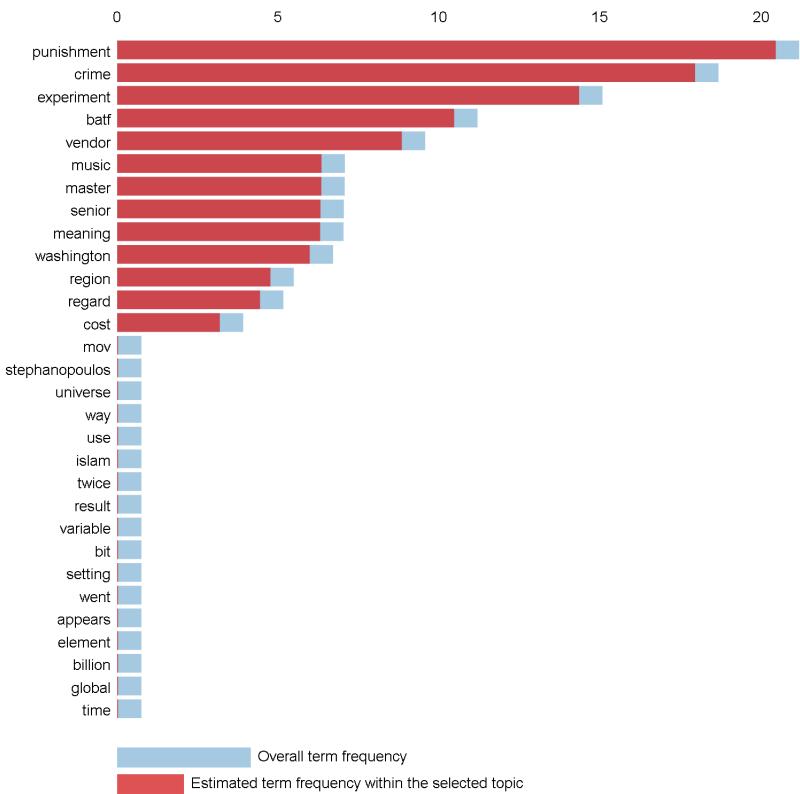
### Top-30 Most Relevant Terms for Topic 15 (0.4% of tokens)



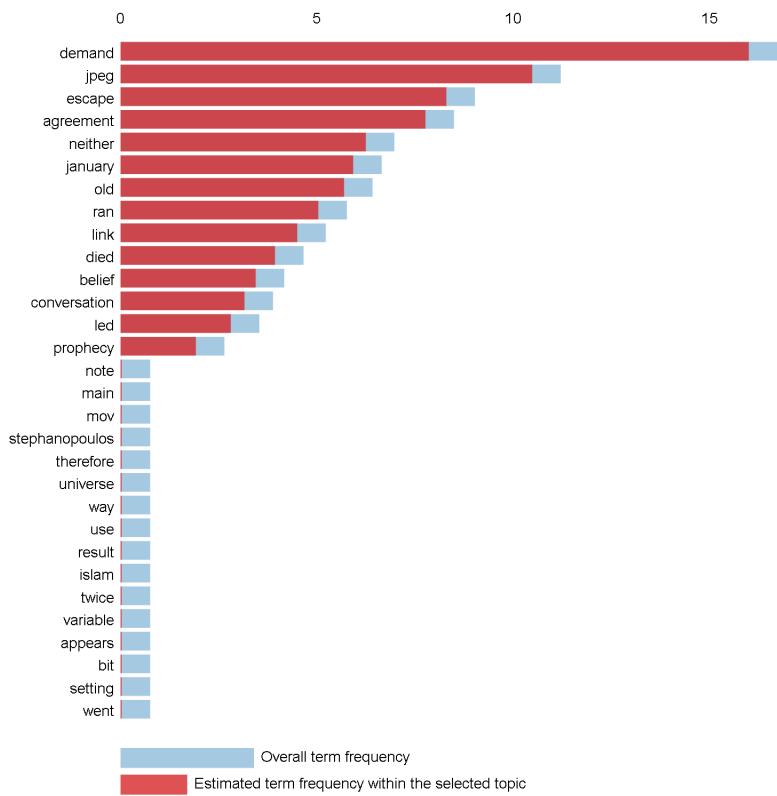
Top-30 Most Relevant Terms for Topic 16 (0.3% of tokens)



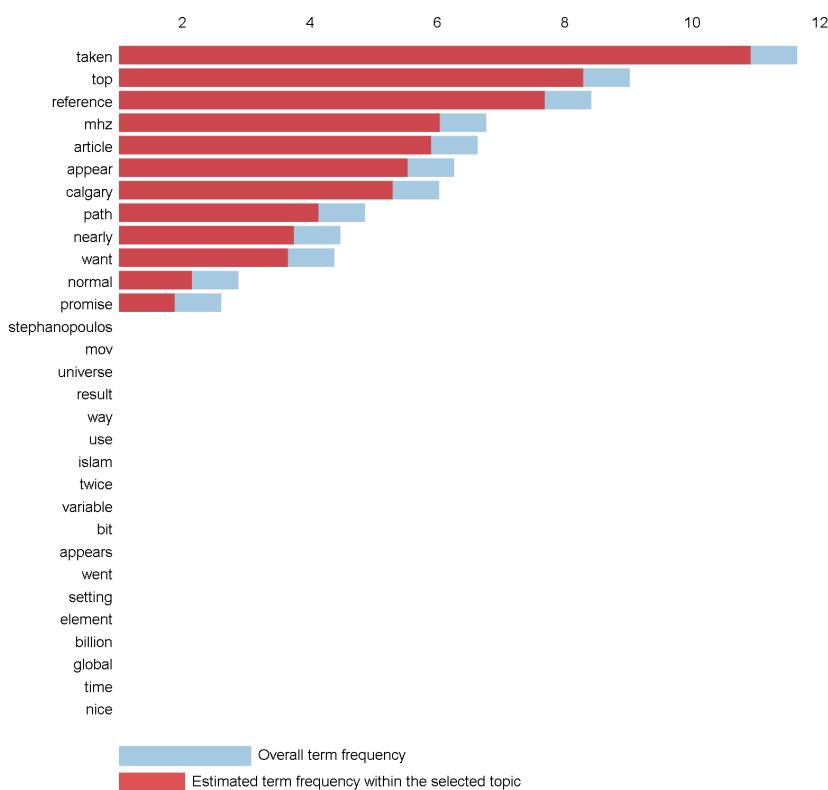
### Top-30 Most Relevant Terms for Topic 17 (0.3% of tokens)



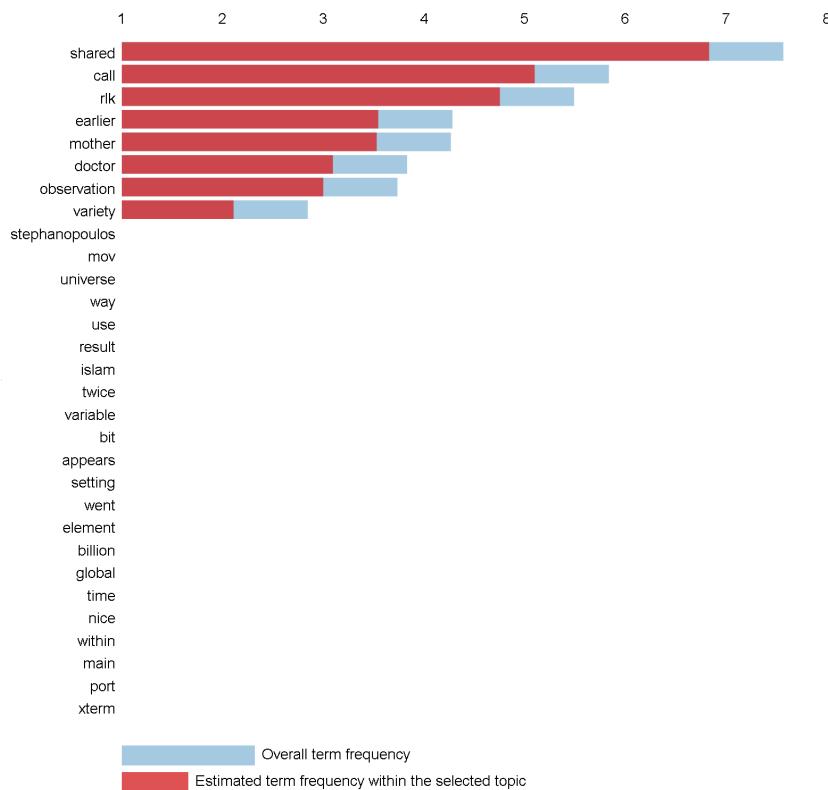
### Top-30 Most Relevant Terms for Topic 18 (0.3% of tokens)



### Top-30 Most Relevant Terms for Topic 19 (0.2% of tokens)



### Top-30 Most Relevant Terms for Topic 20 (0.1% of tokens)



war  
genocide  
army  
soldier  
town  
killed  
turkish  
muslim  
turk  
turkey  
azeri  
republic  
civilian  
group  
million  
iran  
ottoman  
population  
argic  
woman  
territory  
azerbaijan  
village  
soviet  
university  
armenia  
christian  
orthodox  
massacre  
history  
jew  
russian  
jewish  
land

r8f f9d k9jz gizw nuy  
fyn g9p r1k r186  
**bhj** bxn chz qax qtm  
bxn okz b4q q50 nrhj a86  
ghj u31 N nrhj a86  
g9v p4u v9fq45 mg9v  
fij b8e bhjnair v9fq45 mg9v  
b8f c8v  
mb8f yf9 wwiz  
wm4u

A word cloud centered around the word "email". The words are in various sizes and colors (blue, green, red, orange) and are arranged in a roughly circular pattern around the central word.

sale  
format  
cover  
comp  
free  
info  
ftp  
looking  
thanks  
list  
org  
art  
dod  
john  
price  
replica  
jpeg  
address  
university  
image  
manual  
zip  
internet  
gov  
interested  
graphic  
condition  
advance  
please  
send  
copy  
fax  
shipping  
env  
interested

timony  
mil  
overbank  
jite  
gratsys  
liner  
navy  
steфанopoulos  
tomy  
winbench  
override  
frame  
matchhead  
seminar  
lates  
wing  
naval  
lipman  
vessey  
keyand  
critten  
ulf  
disney  
bobby  
louie  
compass  
strap  
lipman  
lade  
keyand  
critten  
ulf  
disney  
bobby  
louie  
rit  
station  
afterline  
tvwan  
tvwan  
schitt  
boot

dtmedin speedo intergraph  
resistor everywhere analog delegate  
polarity plot\_data huntsville uunet  
corpscope gnd B03004  
ingr catbyte b30 nichelodeon  
namaki networking medin ecc ssd  
networking medin ecc ssd  
dr Fey cathode

denizen claus critz ibka trinitarian eternality  
mnemonic caligiuri tul shk  
celp zoroastrian atheisten  
atwood magi taoism prometheus ecclesiastical  
bunny swinburne rsadsi subsistence  
altar aiz aloysius aap unbelief wch  
enviroleague burzynski een mackie prof  
free thought watchmaker zoroaster

disk computer monitor rom controller hard  
speed hardware monitor bus ide  
cpu software ibm mac cable price  
drive system bios serial port thanks  
mode scsi chip data  
motherboard jumper

hitter  
team pit league  
phi wing suck milwaukee cooling  
hit bos ny cooling  
sport cup minor red  
cal murray buf run  
det pitcher lost murray  
season torstl mon nyi chi min  
water cub  
van  
torstl mon nyi chi min pitch

one like  
may used question need many thing  
good first take  
get people year think time  
said problem much work  
something way make know  
two want say  
well new come could right also even  
would  
use key see use point

directory available program  
data display running anonymous color  
set pub sun server using  
image also source  
unix system  
get application support  
user code software  
file window version  
information function screen  
com

nf g3o graeme myhint psiz  
xsize hints trinomials  
gscr t1 acura  
line thick  
xsize hints  
1dl winos y86  
line thick  
detail win  
integra decompress  
zrceasimpladew  
ick  
weatherhead stafford  
sau

lavrencic yoked  
kurtnye abnorm 17y  
uwec vonnegut ijs  
distsq ihr cjackson  
jeezus muskingum loveth  
doctrine\_ pps  
delusion nyeda chsvax  
pavement cheering clubbing  
p-4 pads CNS dreameda  
apocalypse pavement  
pads CNS  
millie knx  
masoretic funeral whl  
kula borut caltrans  
die kula janet  
caltrans 1  
july janet

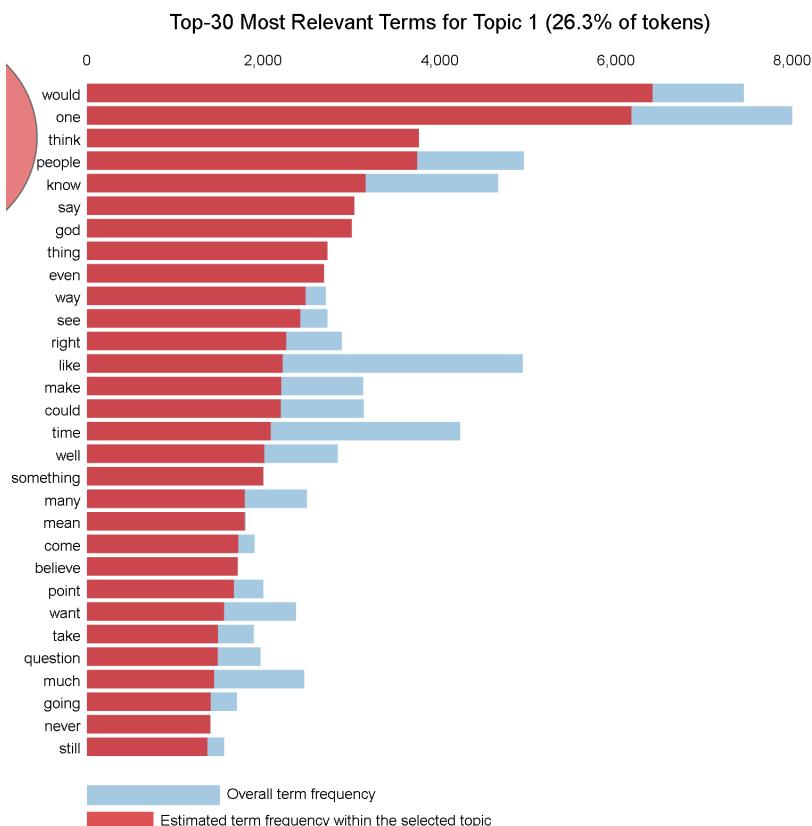
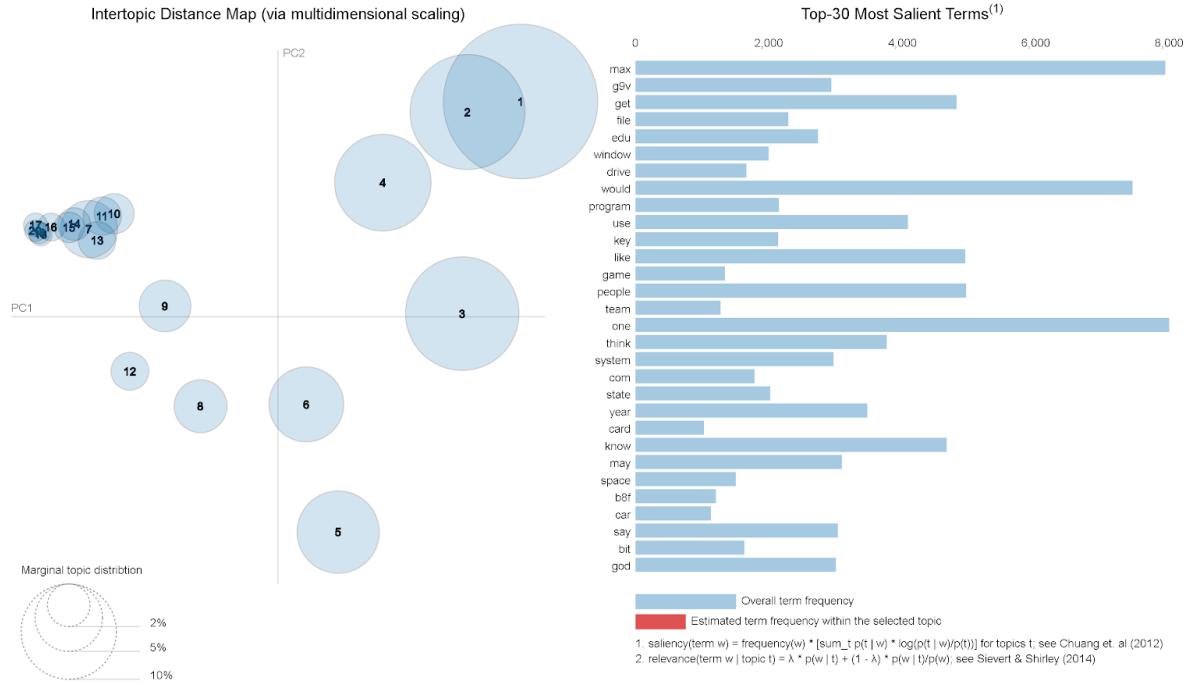
A word cloud centered around the word "space". The words are arranged in a roughly circular pattern around the center word. The size of each word indicates its frequency or importance within the context of space exploration and technology.

semi handgun seattle revolver  
nra death self arm shotgun  
laser batf safety printer control  
rate Limbaugh pistol police  
safety printer control e carry  
assault violent rkba study  
auto weapon armed city  
defense rifle automatic homicide state  
firearm cycle accident  
concealed

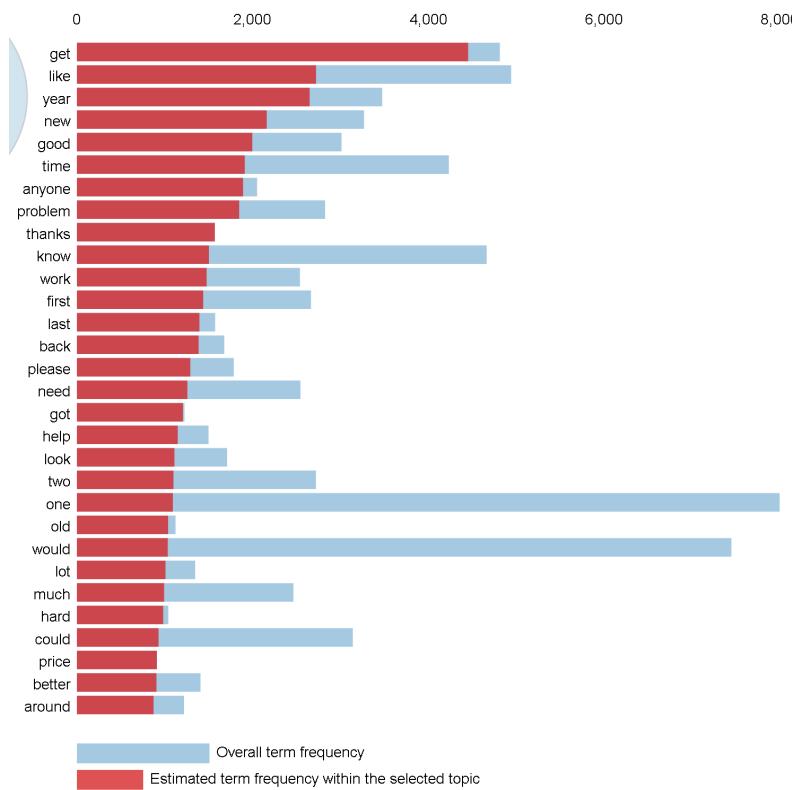
printf file font section name  
info oname nist include filename byte  
line buf title dpy winner open  
size char remark contest printf  
define author check\_io ioccc return  
char int program keyboard stream  
key obfuscated character null uuencode

uvwxyz mcxsqst3p w45gc\_  
w3q shz t3q sy\_x\_scx SCX  
j1 xte s1hz x\_scx neu  
c5 t3s w47ck8 air s  
cahf fdz w3p l45cho  
b9r mc\_ww7 w3sxxi hzv ck  
mc\_ww7 sj1 cx\_scx w44t45cho

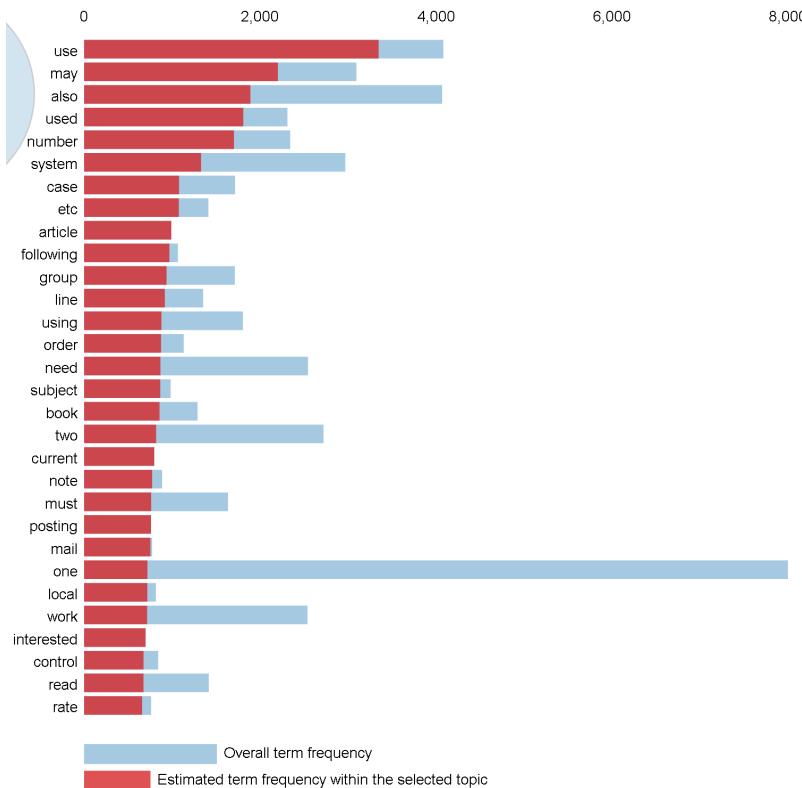
pregnancy laserwriter  
ad1 married doctor  
med married seizure sugar sexual  
homosexual regulated sinus yeast corn  
gay treatment quack rhinac  
regulated gerald vitamin study  
influenza skin ini qur Leishman  
live militia diet diagnosed  
acid infection bullock antimale



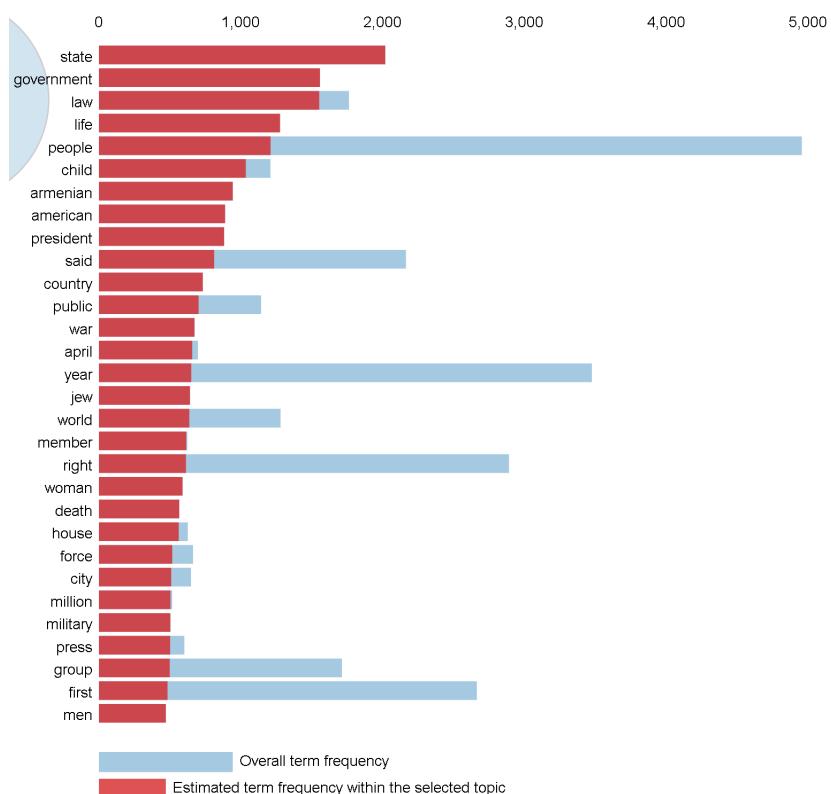
### Top-30 Most Relevant Terms for Topic 2 (14.5% of tokens)



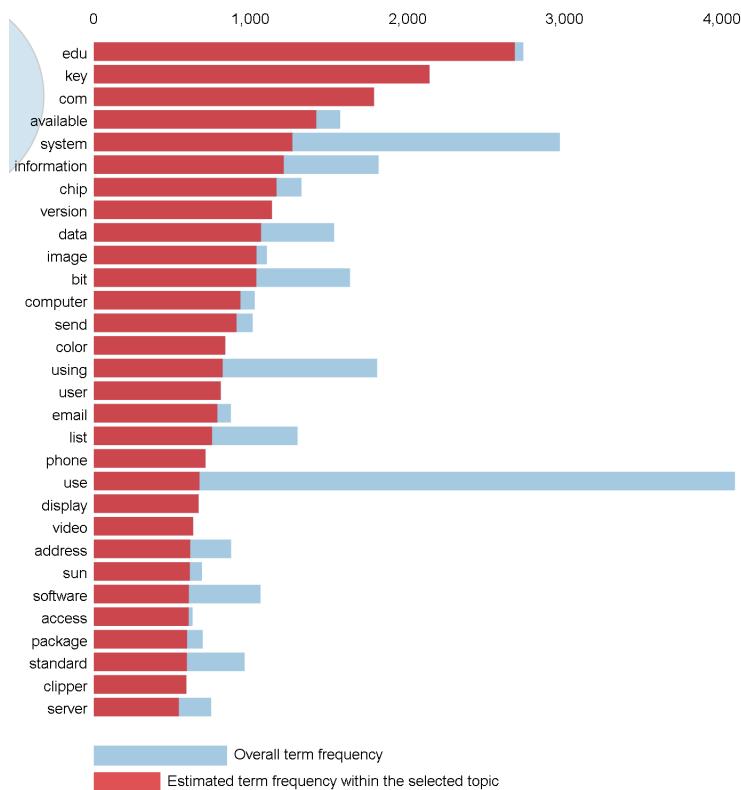
### Top-30 Most Relevant Terms for Topic 3 (14.1% of tokens)



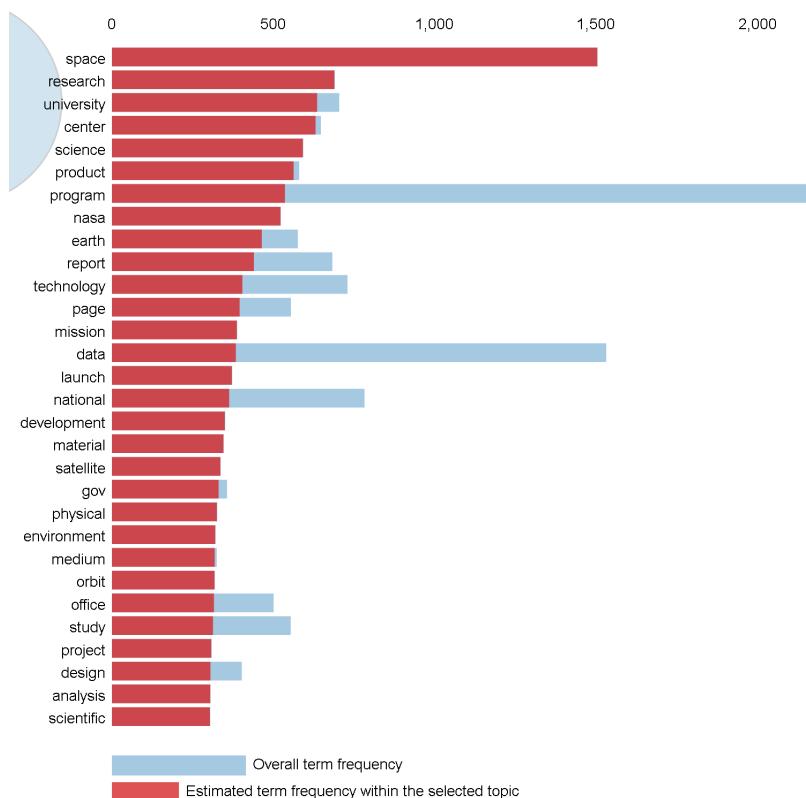
### Top-30 Most Relevant Terms for Topic 4 (10.2% of tokens)



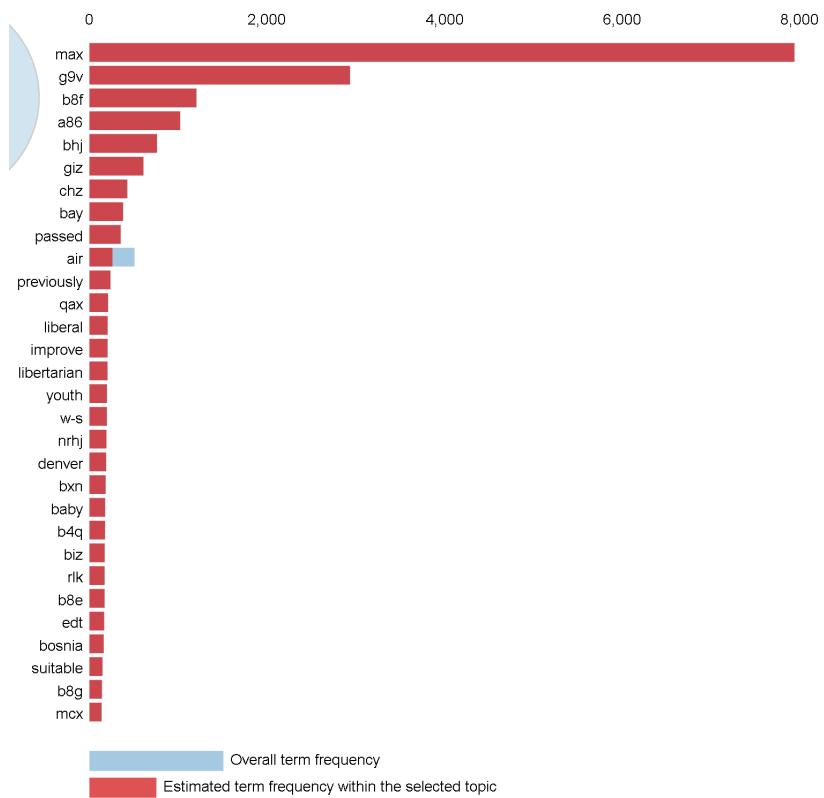
### Top-30 Most Relevant Terms for Topic 5 (7.4% of tokens)



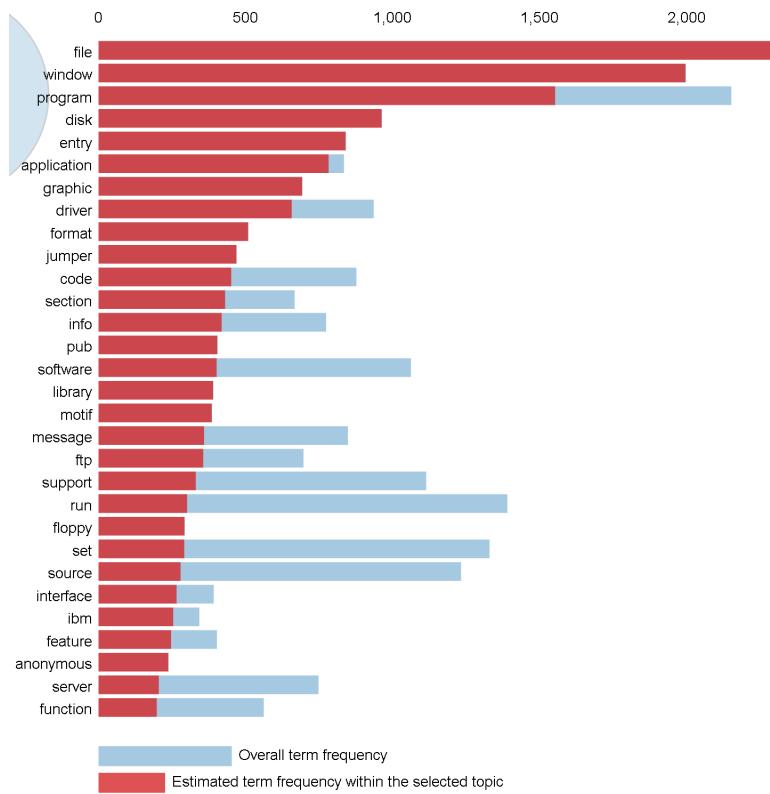
### Top-30 Most Relevant Terms for Topic 6 (6.1% of tokens)



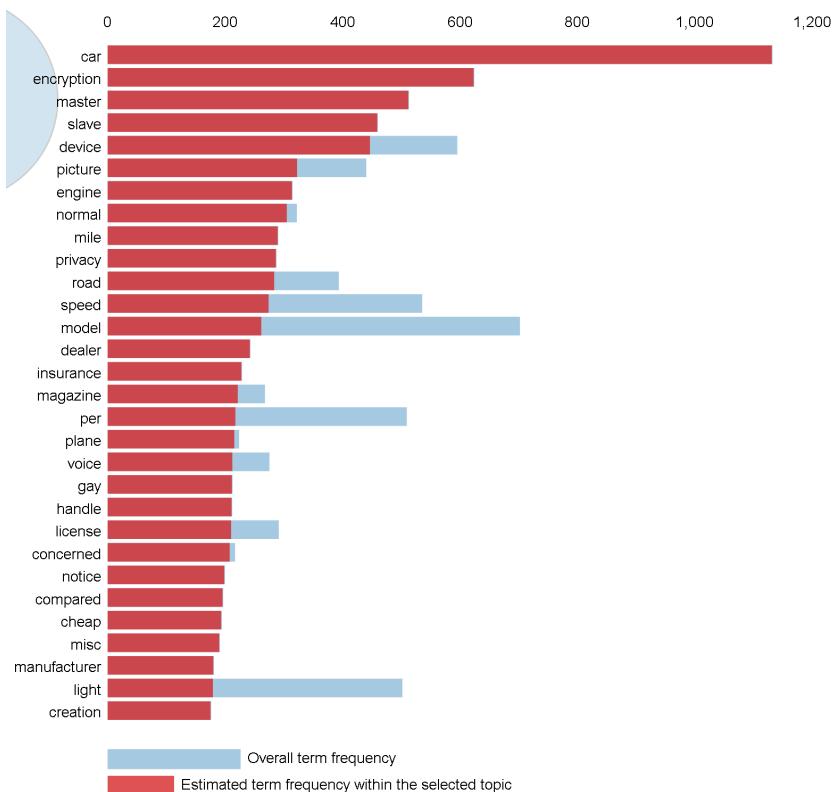
### Top-30 Most Relevant Terms for Topic 7 (3.6% of tokens)



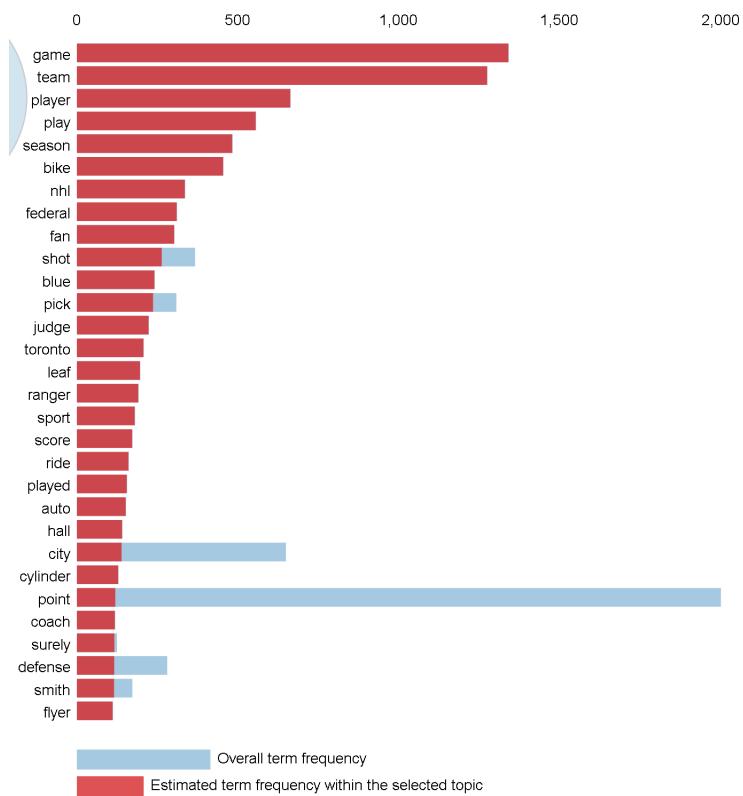
### Top-30 Most Relevant Terms for Topic 8 (3.1% of tokens)



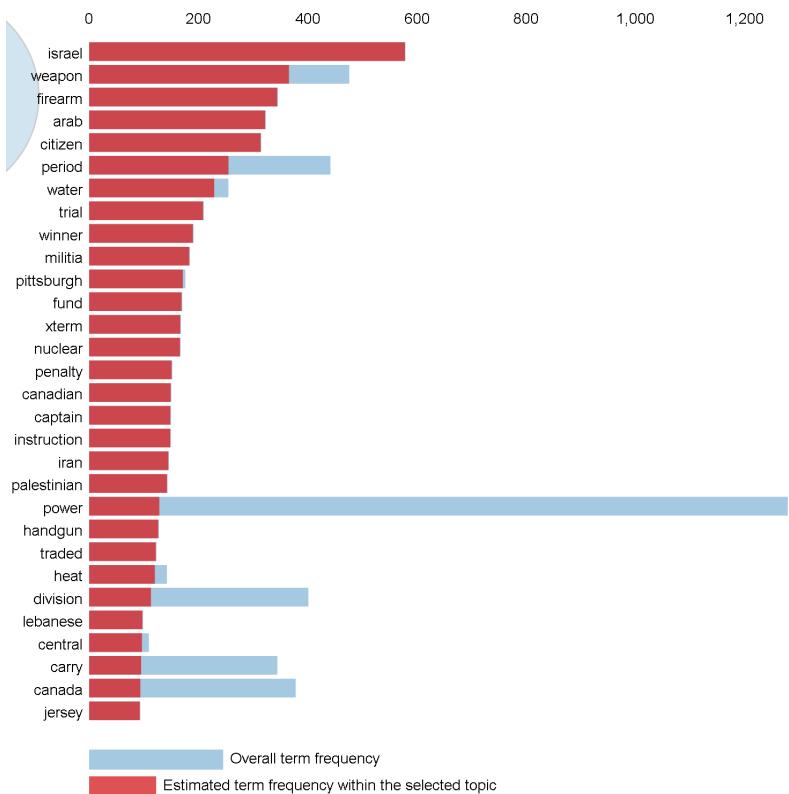
### Top-30 Most Relevant Terms for Topic 9 (2.9% of tokens)



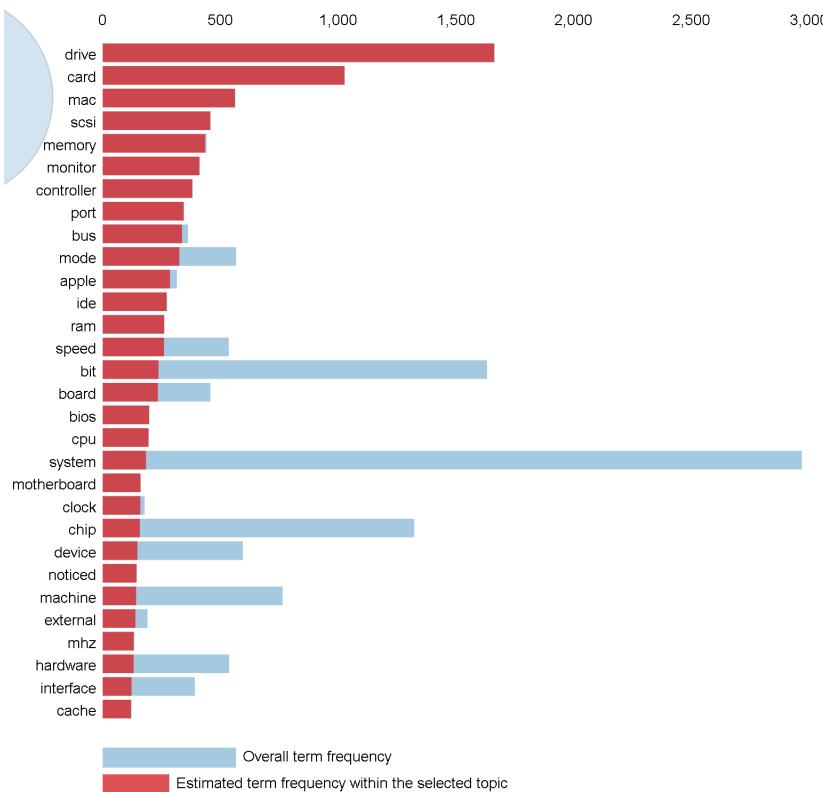
### Top-30 Most Relevant Terms for Topic 10 (1.7% of tokens)



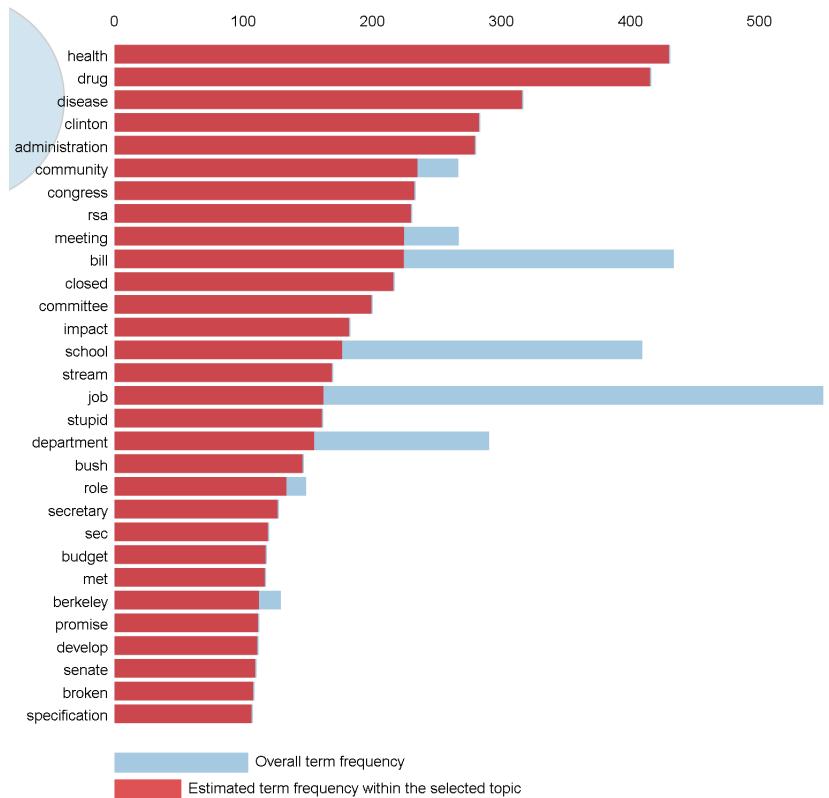
Top-30 Most Relevant Terms for Topic 11 (1.6% of tokens)



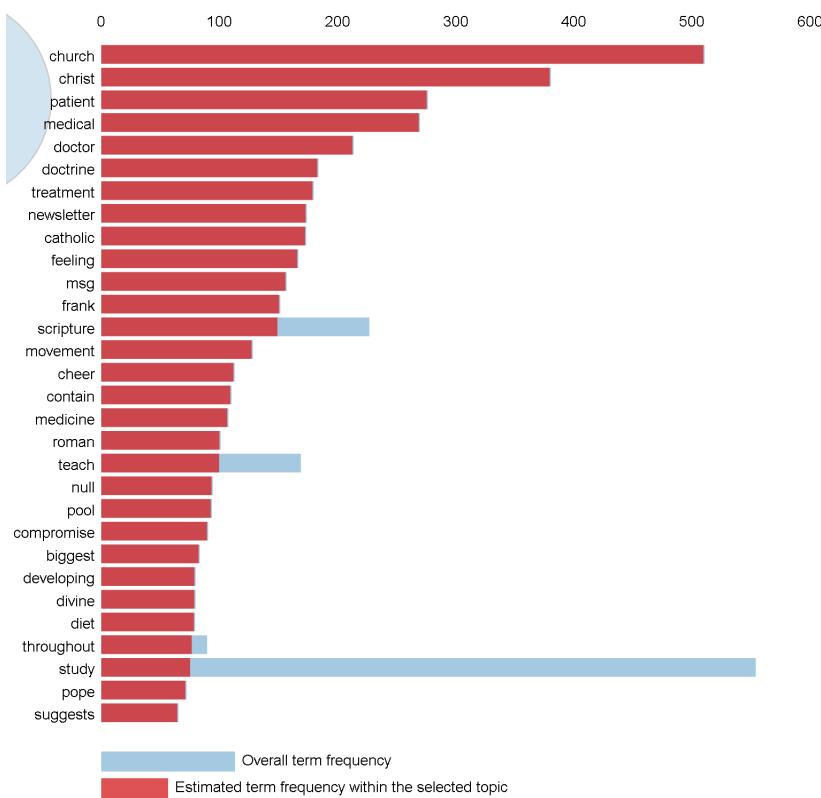
### Top-30 Most Relevant Terms for Topic 12 (1.6% of tokens)



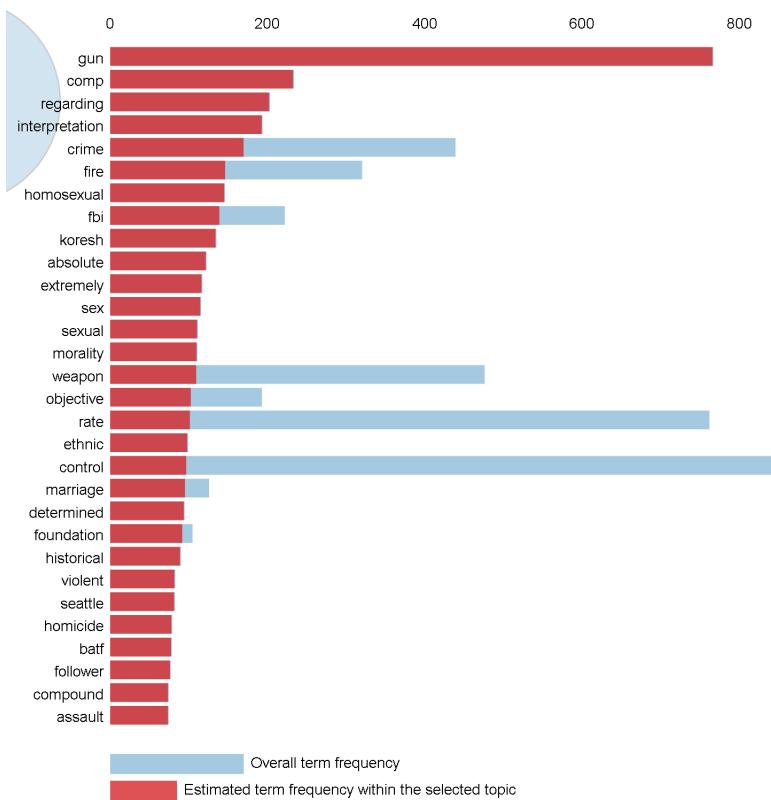
Top-30 Most Relevant Terms for Topic 13 (1.5% of tokens)



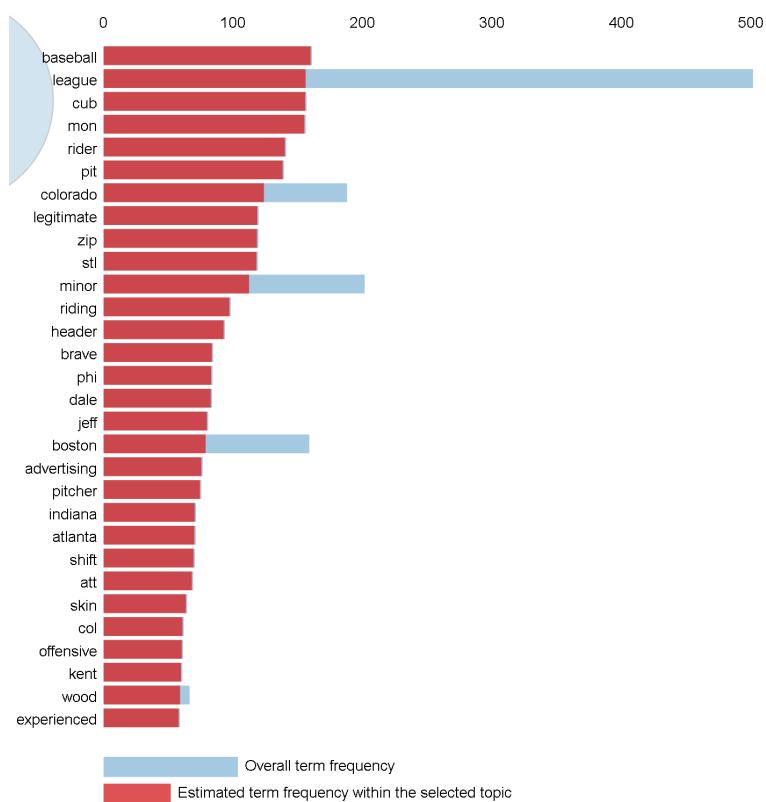
### Top-30 Most Relevant Terms for Topic 14 (1.2% of tokens)



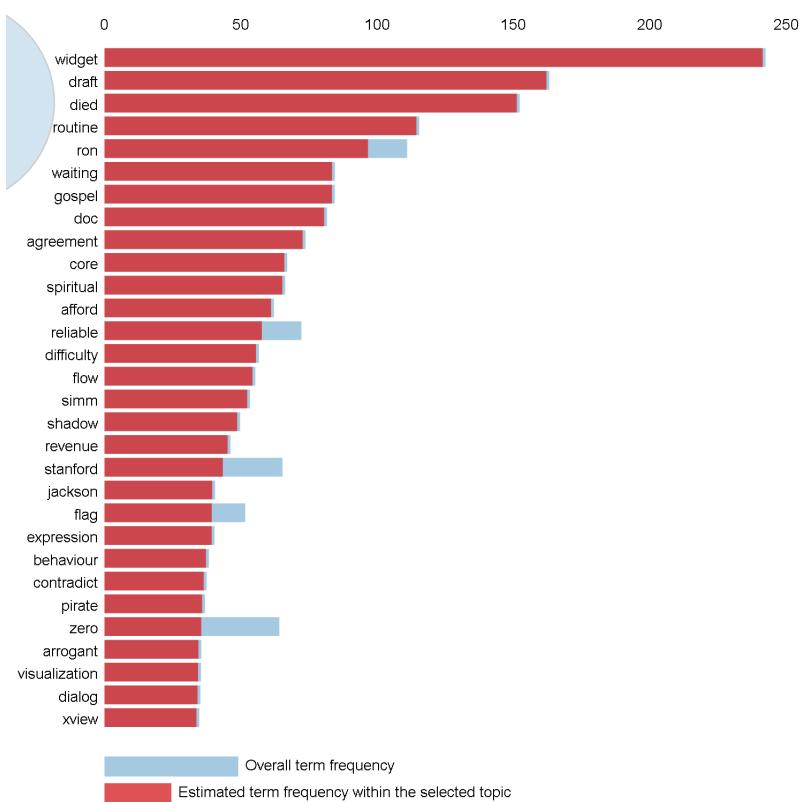
Top-30 Most Relevant Terms for Topic 15 (1% of tokens)



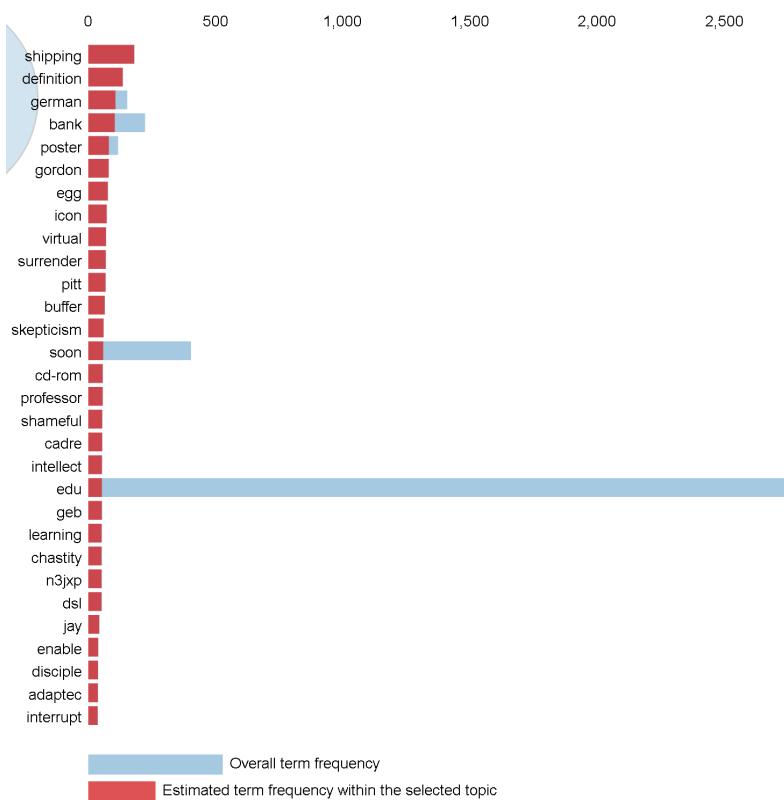
### Top-30 Most Relevant Terms for Topic 16 (0.8% of tokens)



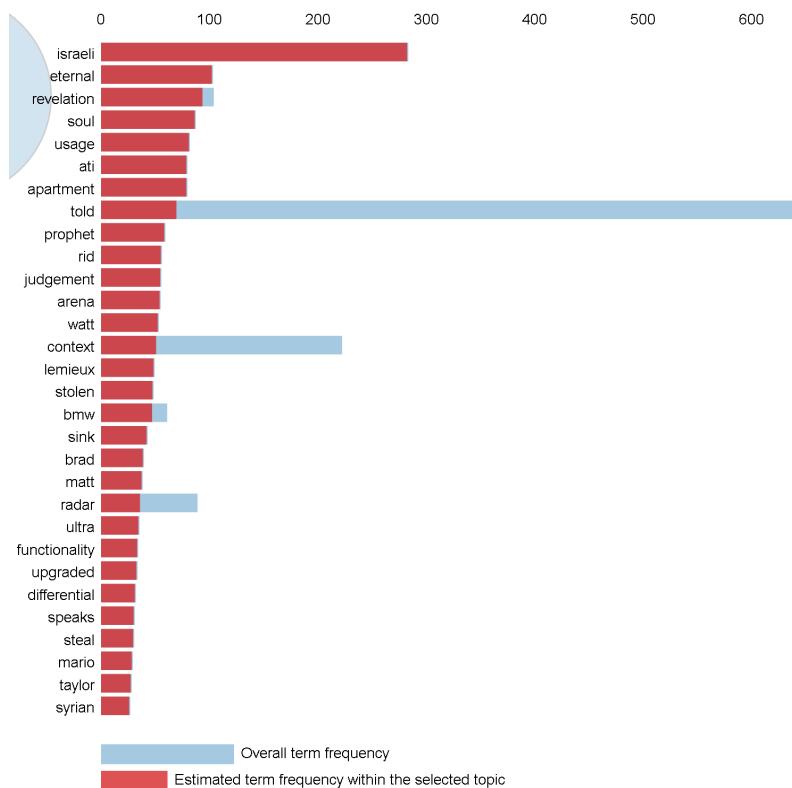
Top-30 Most Relevant Terms for Topic 17 (0.6% of tokens)



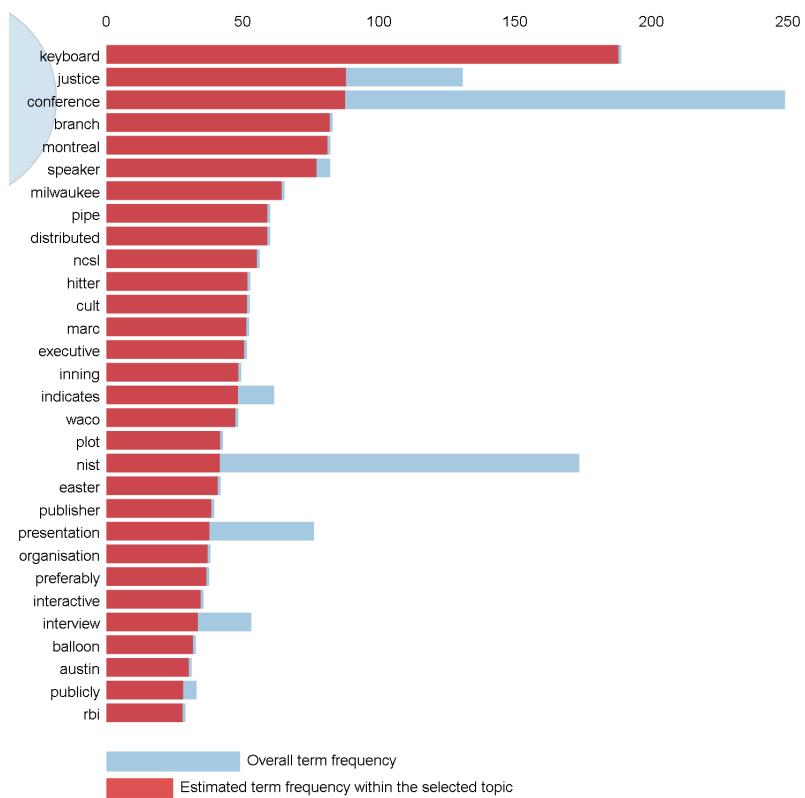
### Top-30 Most Relevant Terms for Topic 18 (0.6% of tokens)



### Top-30 Most Relevant Terms for Topic 19 (0.5% of tokens)



Top-30 Most Relevant Terms for Topic 20 (0.5% of tokens)



also man  
church jesus sin  
one say come book  
religion heaven  
new catholic life  
word lord love hell faith  
bible spirit son john  
christian soul world eternal believe

okz rlk  
chz b4q  
fij end  
max  
g9v bxn b9r a86  
giz biz qax  
bhj nrhj b8e partair

probably anything sure problem  
take even time way maybe  
want thing also say going never  
tell well make something  
day could people  
see anyone tell well  
anyone day  
would could people  
think better little  
know one good  
work someone got really  
right think know one good  
like work someone  
would could people  
think better little  
know one good  
something

palestinian dog state israeli plane  
shift rider south ride  
sea west area territory road policy  
riding country war black arab  
right border motorcycle left picture  
bikes jewish civilian point  
jew attack terrorist adam  
zone

escape recent mirror shell  
object ftp anonymity  
aid comment ray  
number host Nmath abuse  
steve poster contain  
shareware collection  
cause found  
msg good study  
others archive  
family archive  
description  
announced  
sex announced  
msg effect reaction  
**site**  
food anonymous problem

american  
government  
clinton clipper think  
year people well  
know one car  
going right  
tax new system  
would chip  
new system time  
could get  
private system  
time president  
enforcement

turkish went history government  
nazi two told first  
war people turk came  
muslim country turkey home german army  
woman greek said  
men killed world building azerbaijan  
armenian  
year one saw child dead time  
soldier russian jew armenia

send list email advance  
please info pin one  
etc get model computer  
looking also information interested  
used sell port know  
sale help post  
like reply would  
new edu modern mac  
offer need com mail address

many access privacy  
faqkey data public  
mail disease patient  
product also service  
rsa internet may method  
computer network health pub  
sci news general  
network ripem encryption group algorithm part  
information standard medical posting

owner law weapon  
citizen issue arm  
firearm cover handgun  
right city politics  
state auto use  
driver driver safety  
criminal bill defense  
death risk officer likely  
carry homicide dangerous  
brake killed  
rate control amendment  
police amendment

bios driver ide data bus tape master  
adapter card slot mcg  
board floppy rom system transfer  
problem power connector  
motherboard internal upgrade  
isa two switch rom  
motherboard backup boot  
two jumper head  
controller cd  
slave militia  
drive device

use even problem also looked  
sound one time used  
would two way  
get memory  
test fast second  
well anyone know  
chip value since mode  
first idea work  
line serial need read  
could much current actually  
bit speed using

run  
point  
hitter  
**player**  
think  
get  
score  
hockey  
last  
good  
goal  
best  
hit  
well  
time  
first  
would  
got  
made  
stephanopoulos  
light  
wing  
nhl  
fan  
back  
great  
two  
season  
baseball  
win  
home  
one  
league  
**game**  
**team**

interface  
widget  
**window**  
mac  
support  
set  
mouse  
work  
cdw  
machine  
display  
sun  
software  
motif  
also  
font  
run  
**version**  
graphic  
program  
server  
user  
package  
screen  
problem  
running  
monitor  
application  
available  
driver  
code  
card  
running monitor  
**color**  
use  
ram  
file  
video  
system

output  
directory  
size  
contest  
f  
file  
entry  
program  
format  
image  
build  
text  
char  
character  
must  
rule  
name  
information  
number  
title  
source  
byte  
stream  
code  
info  
read  
command  
include  
convert  
line  
input  
use  
section  
may  
line  
use  
section  
directory  
size  
contest  
**file**  
**entry**  
**program**  
**image**  
**format**  
**character**

casereason  
however  
right  
argument  
claim  
believe  
would  
mean  
example  
know  
evidence  
belief  
thing  
person  
said  
statement  
true  
must  
say  
fact  
point  
could  
make  
well  
question  
human  
many  
law  
even  
time  
way  
life  
**people**  
**one**  
**think**

center new space program university

system president national research

science united state launch

study satellite general office shuttle

april press technology nasa

national washington american development available

group information organization report

period

san van mon new

boston francisco

buffalo king pittsburgh det cal philadelphia

toronto but que

andrew frank calgary philadelphia

chicago cub louis

power bay tor league win min stl

shut seconds angeles los

power detroit jose chi

wire engine

water car air

cable hole connected edge

mov one weight tire

back tube electrical

metal neutral heat hot new

front mile circuit

used light ground supply

led voltage motor oil foot

spacecraft audio cou space

orbital sun gas satellite fire

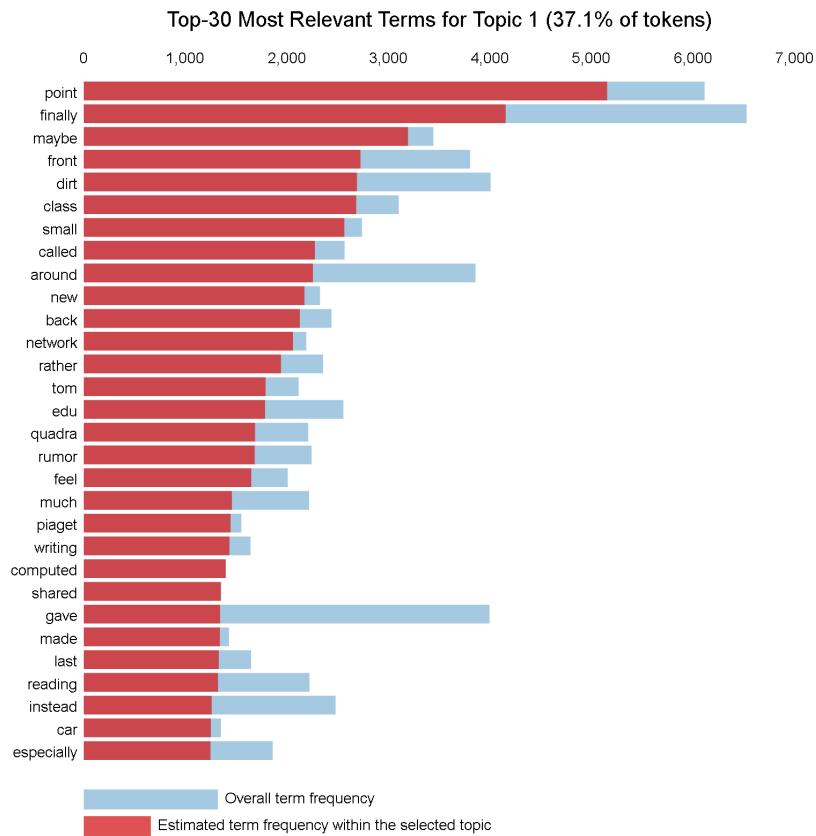
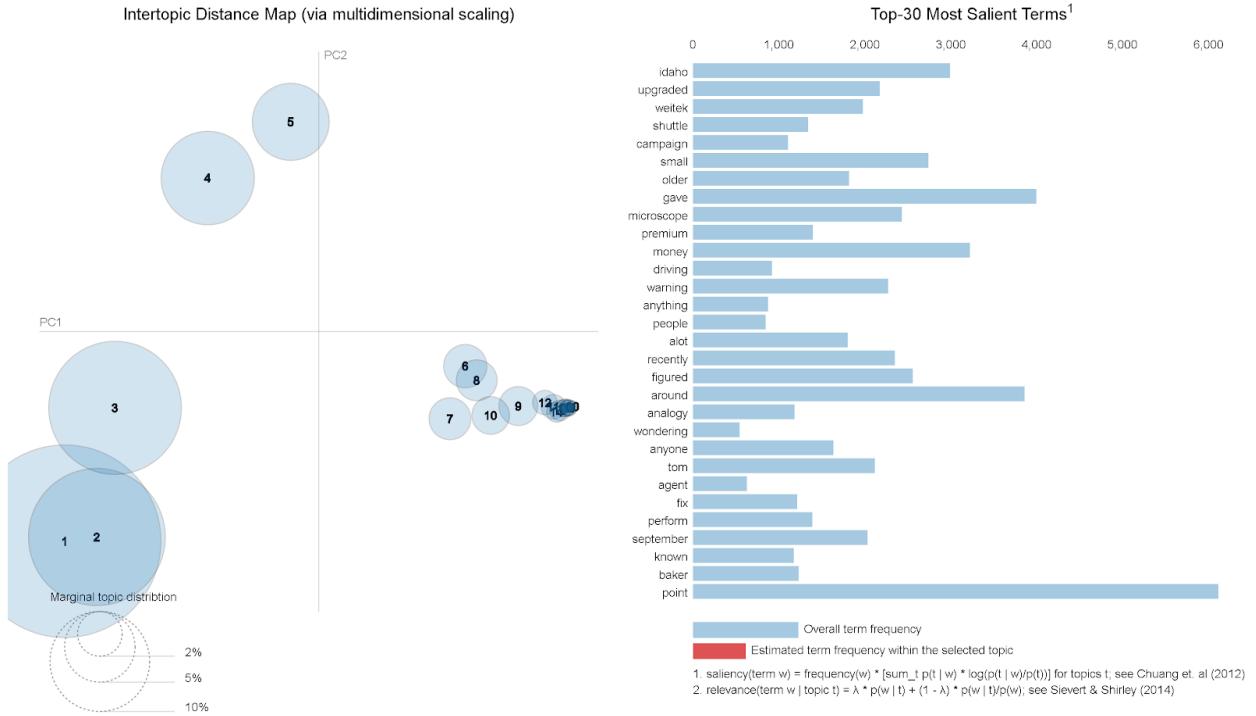
probe detector project radar level image star batf

channel channel project radar level image star batf

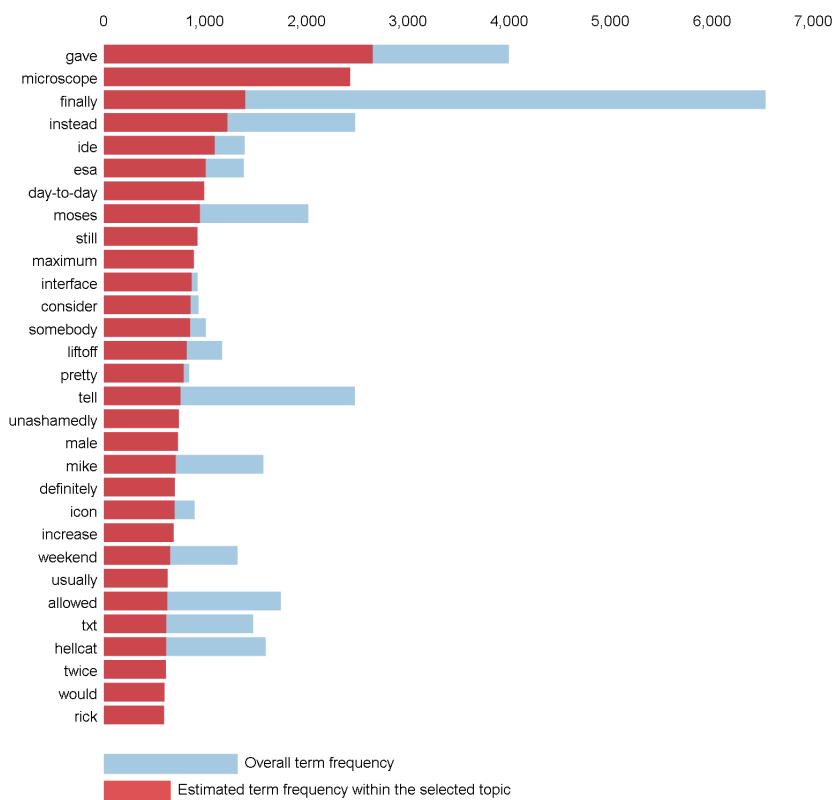
earthy

lunar planet moon would

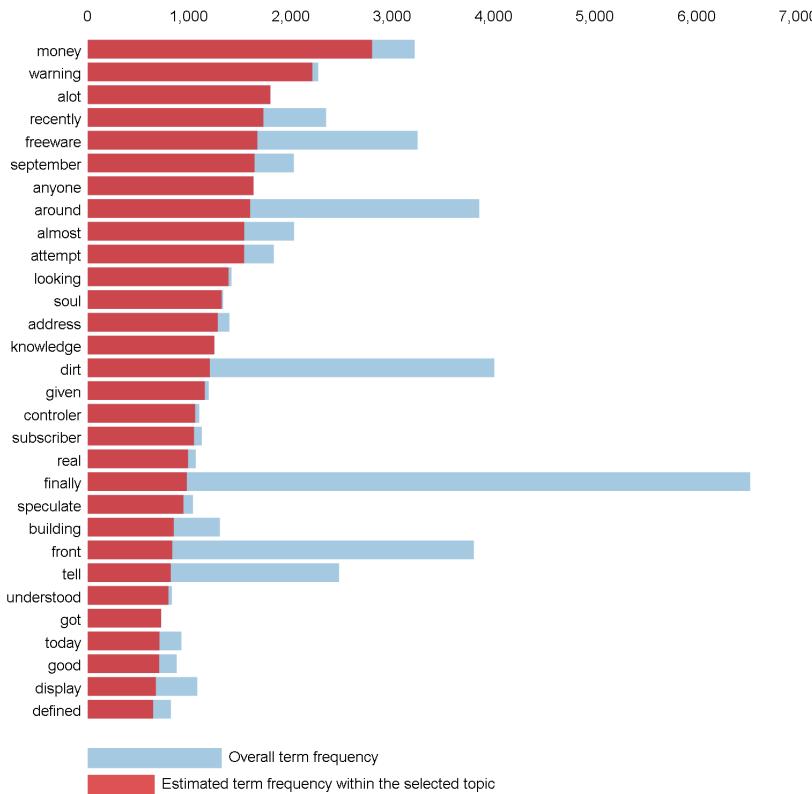
surface koresh radio unit



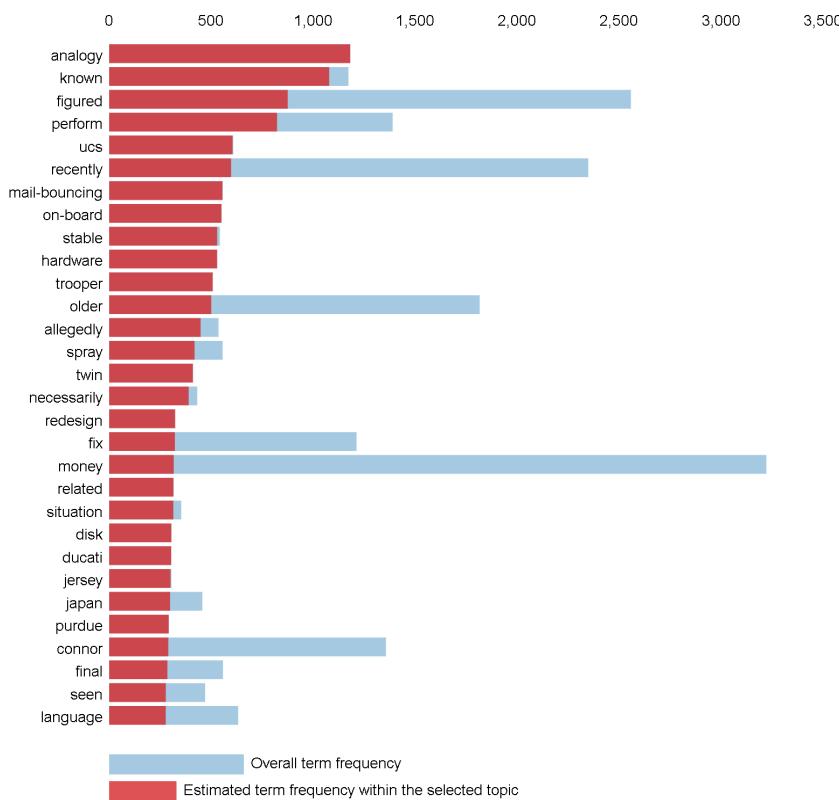
### Top-30 Most Relevant Terms for Topic 2 (18.8% of tokens)



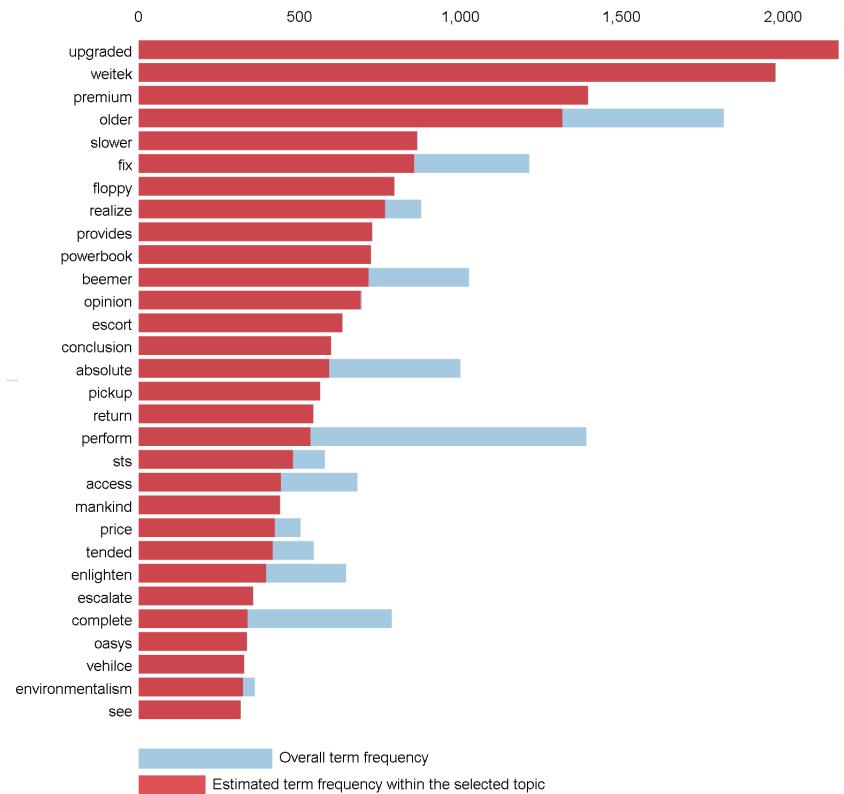
### Top-30 Most Relevant Terms for Topic 3 (17.7% of tokens)



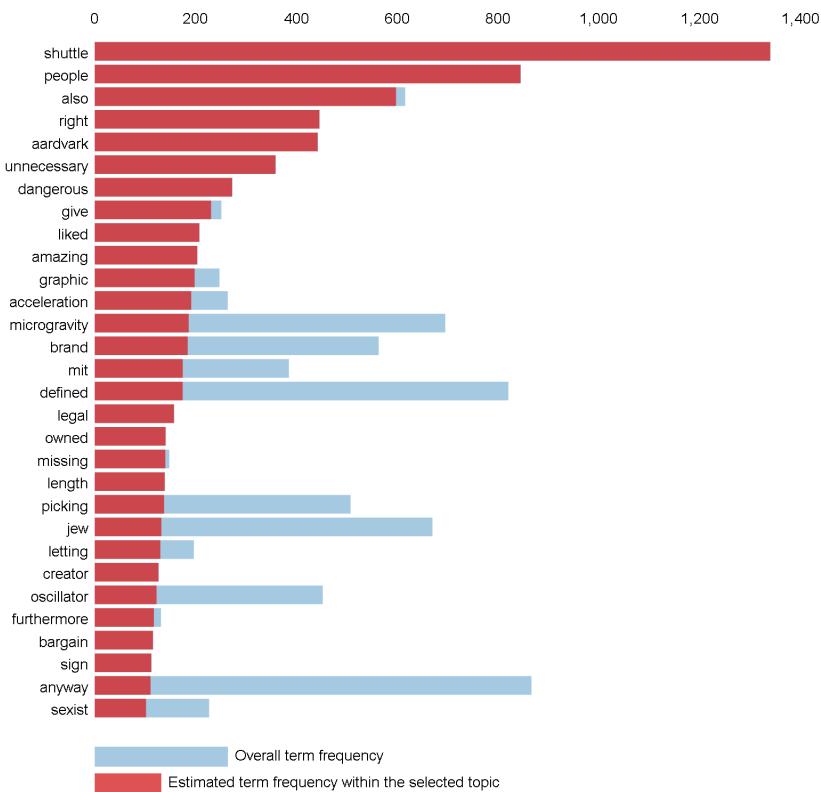
### Top-30 Most Relevant Terms for Topic 4 (8.7% of tokens)



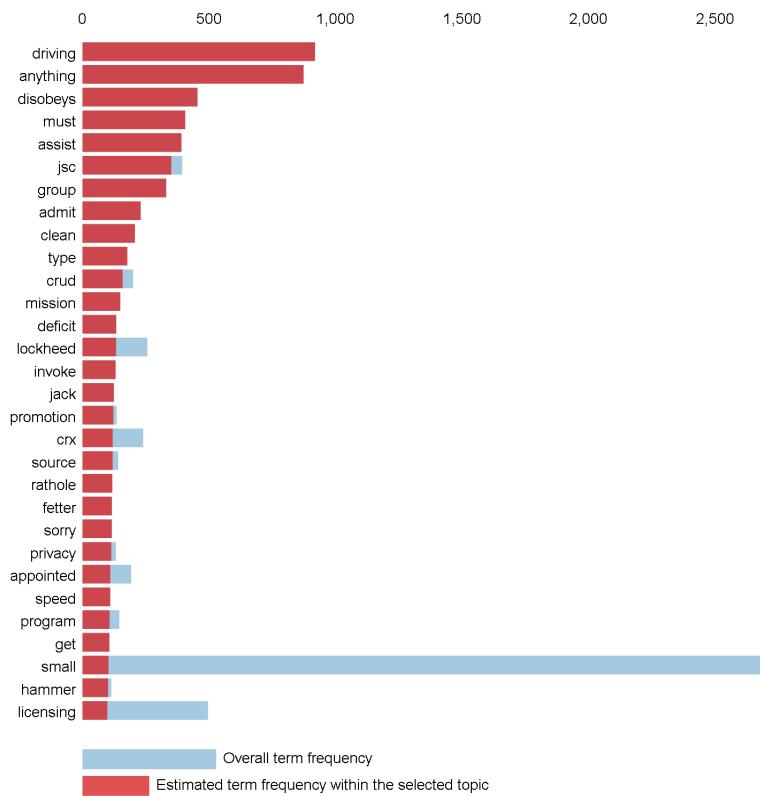
### Top-30 Most Relevant Terms for Topic 5 (5.9% of tokens)



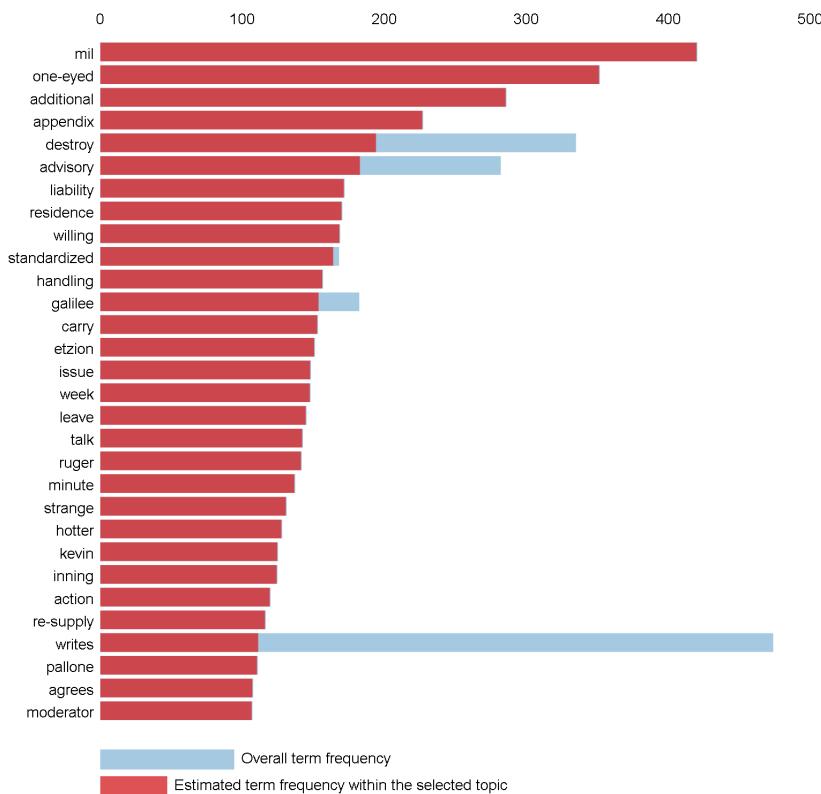
### Top-30 Most Relevant Terms for Topic 6 (1.9% of tokens)



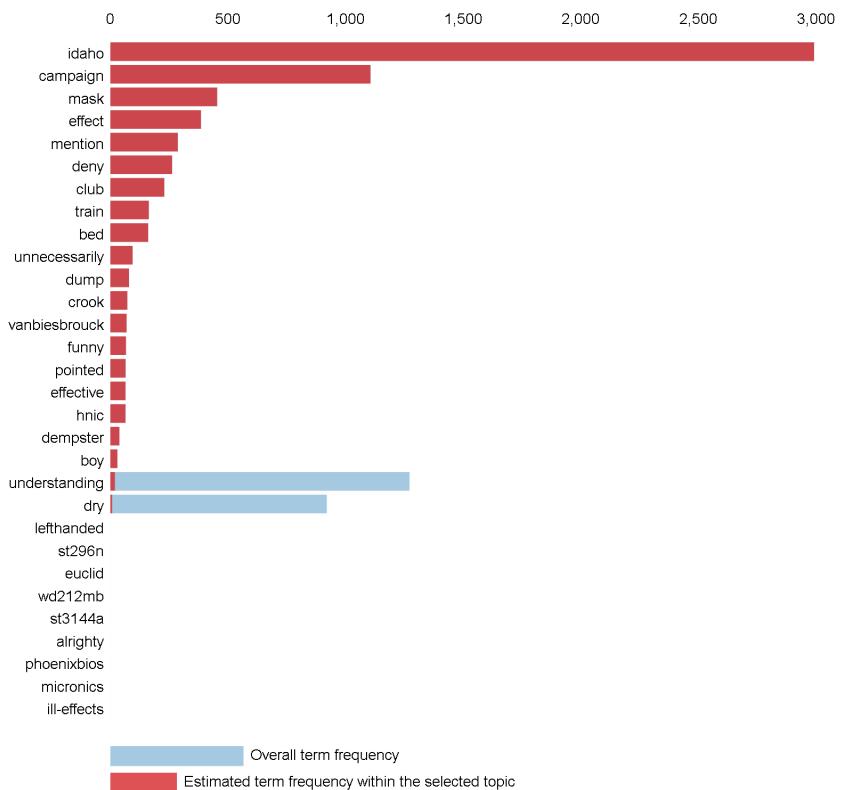
Top-30 Most Relevant Terms for Topic 7 (1.8% of tokens)



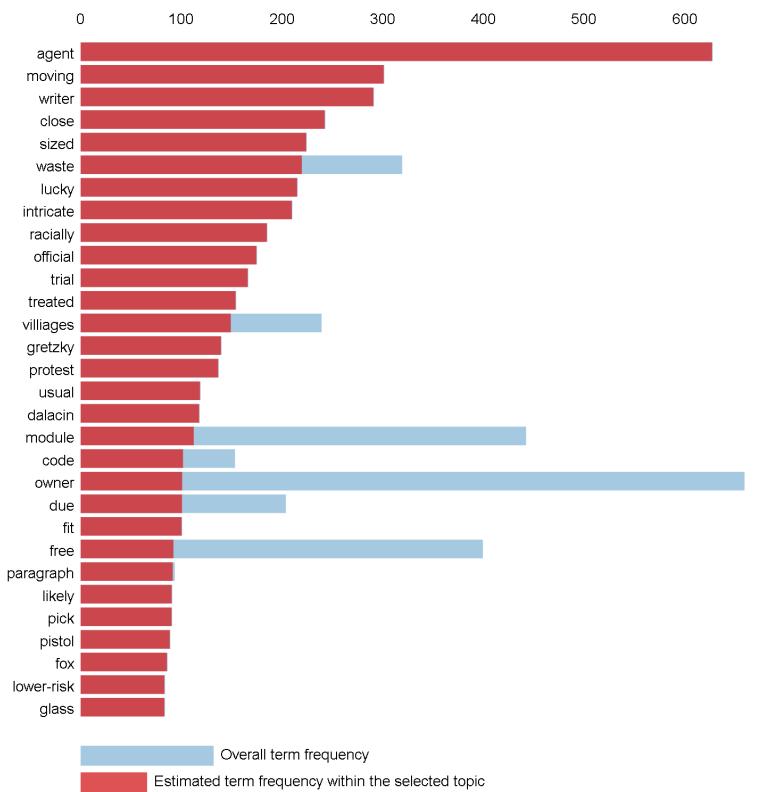
### Top-30 Most Relevant Terms for Topic 8 (1.7% of tokens)



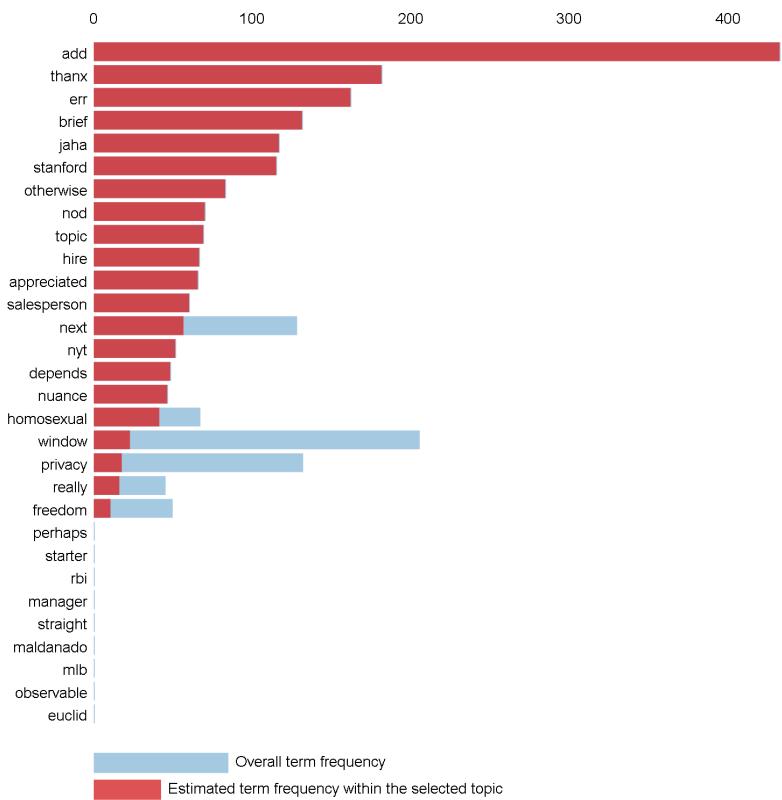
### Top-30 Most Relevant Terms for Topic 9 (1.5% of tokens)



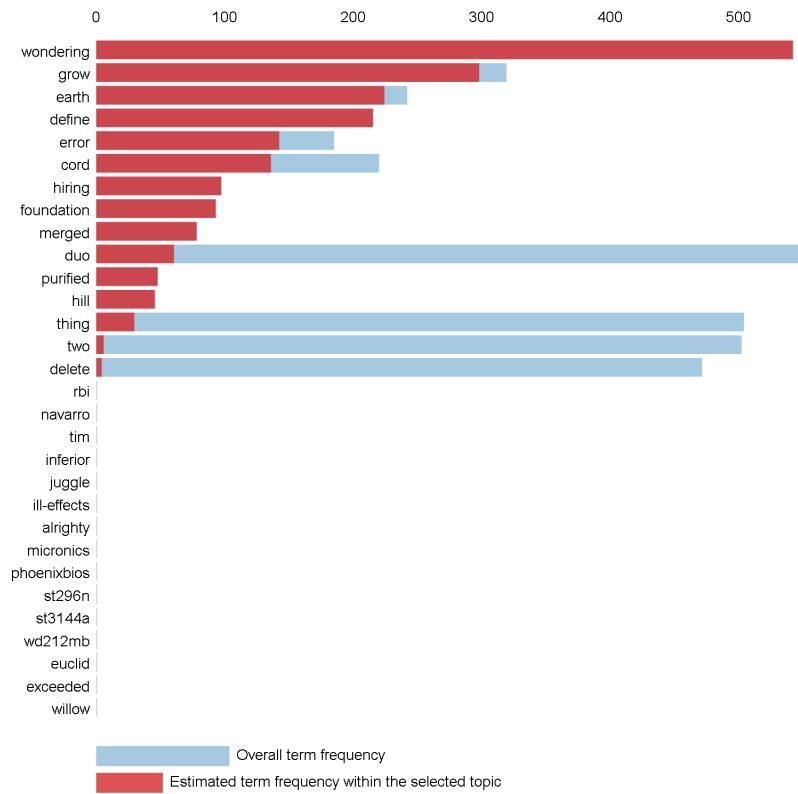
Top-30 Most Relevant Terms for Topic 10 (1.4% of tokens)



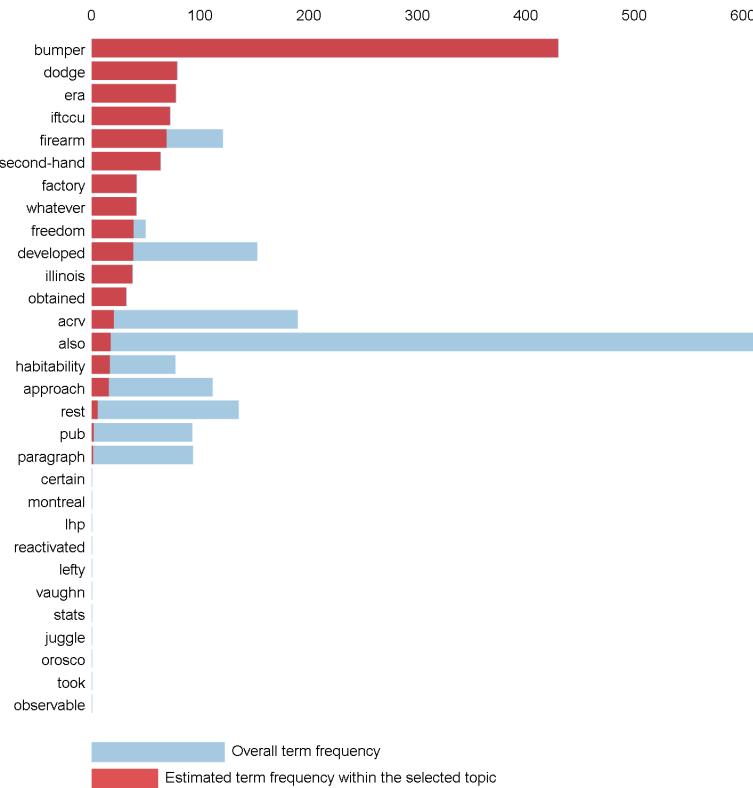
Top-30 Most Relevant Terms for Topic 11 (0.6% of tokens)



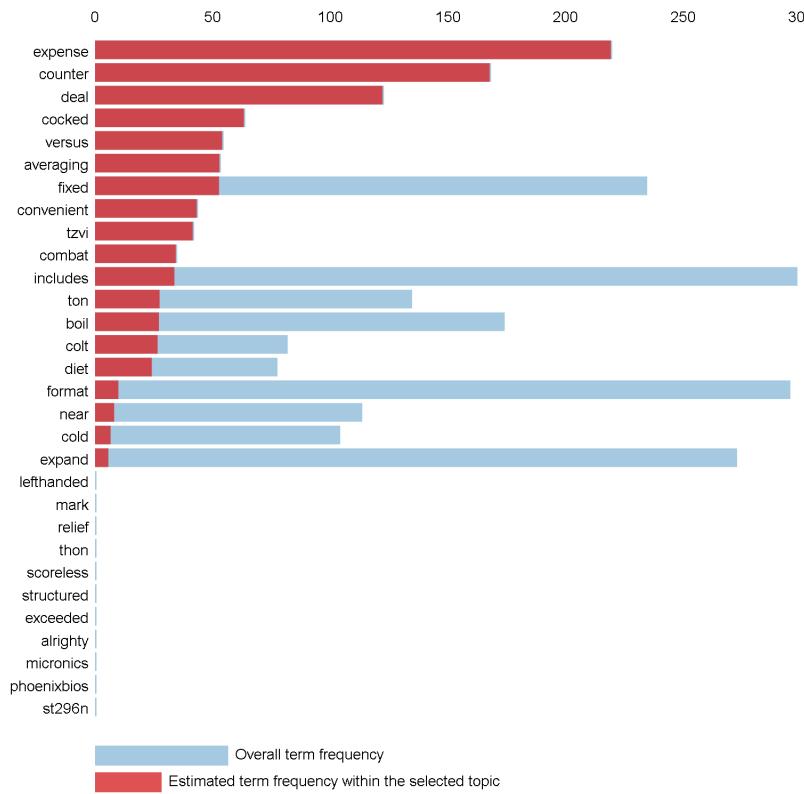
### Top-30 Most Relevant Terms for Topic 12 (0.6% of tokens)



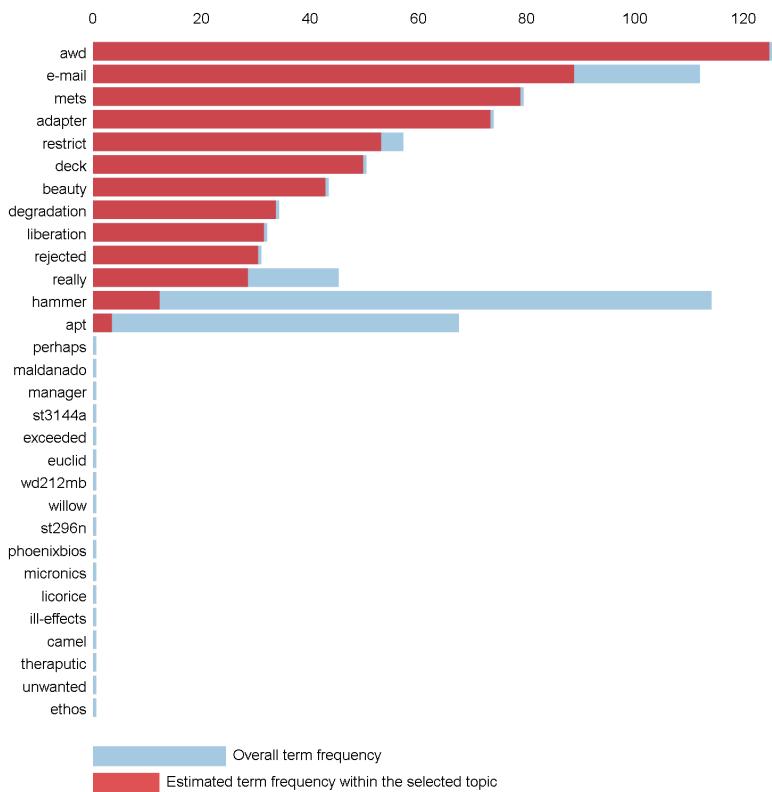
### Top-30 Most Relevant Terms for Topic 13 (0.4% of tokens)



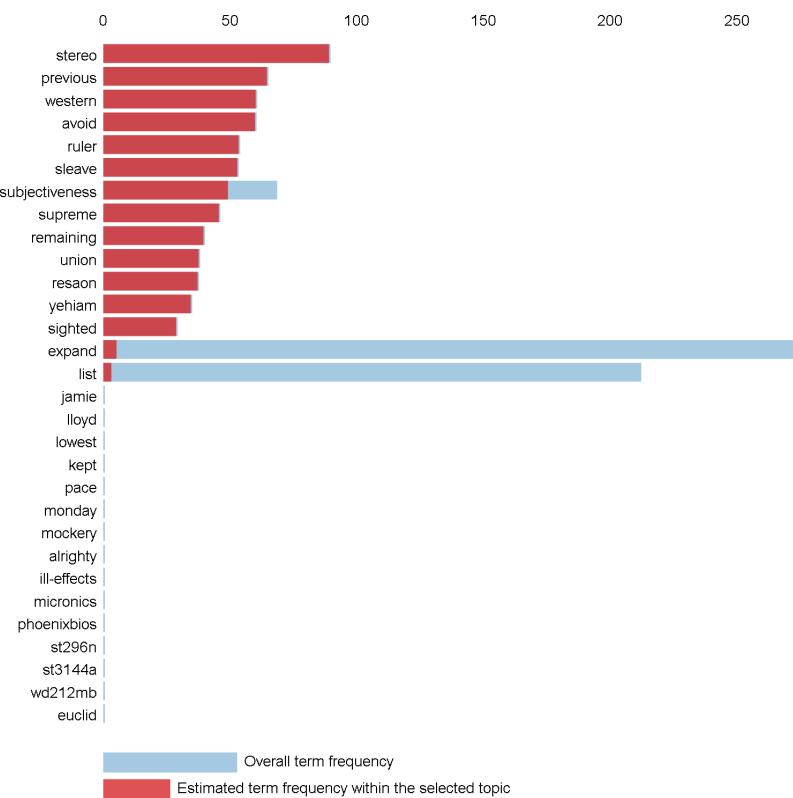
### Top-30 Most Relevant Terms for Topic 14 (0.4% of tokens)



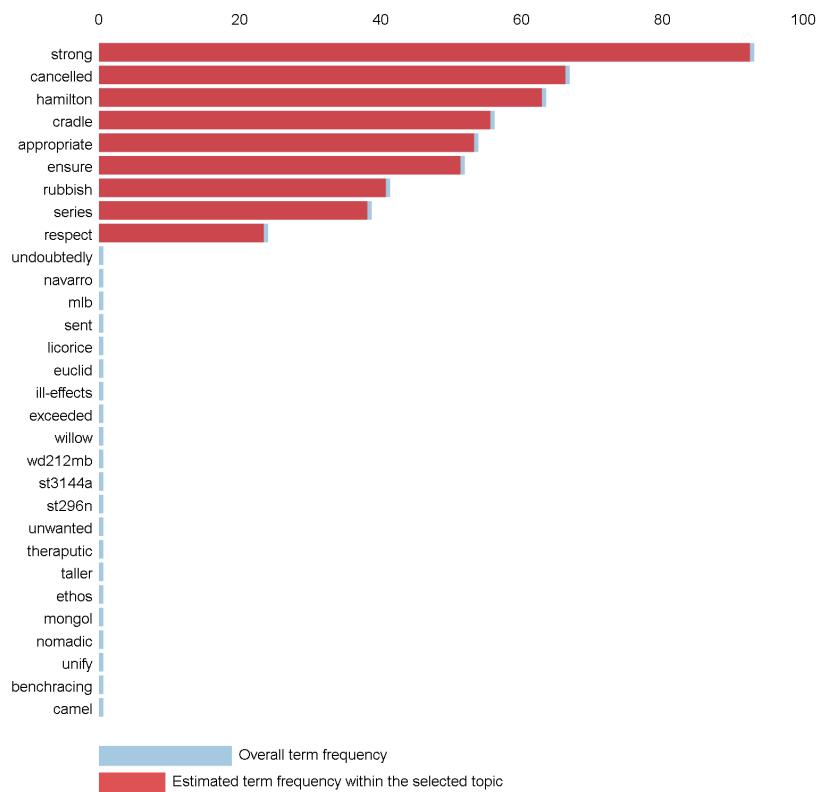
Top-30 Most Relevant Terms for Topic 15 (0.3% of tokens)



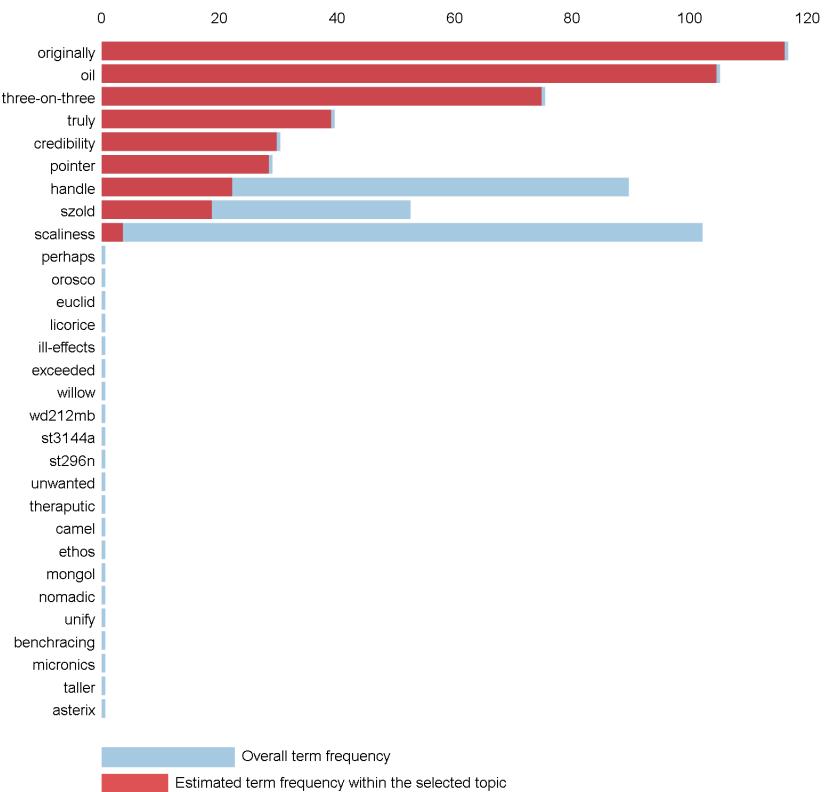
### Top-30 Most Relevant Terms for Topic 16 (0.3% of tokens)



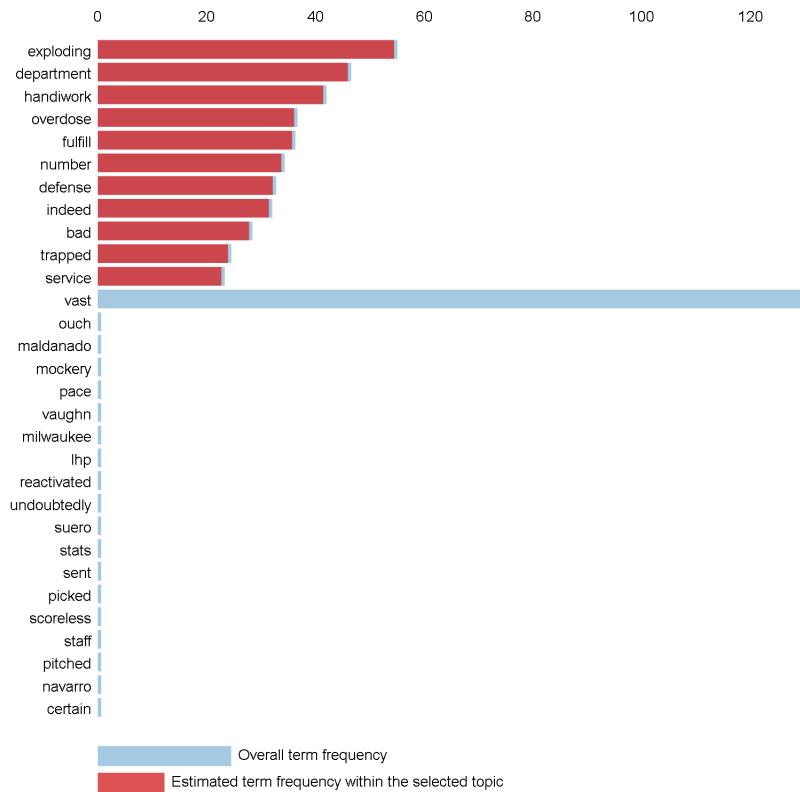
### Top-30 Most Relevant Terms for Topic 17 (0.3% of tokens)



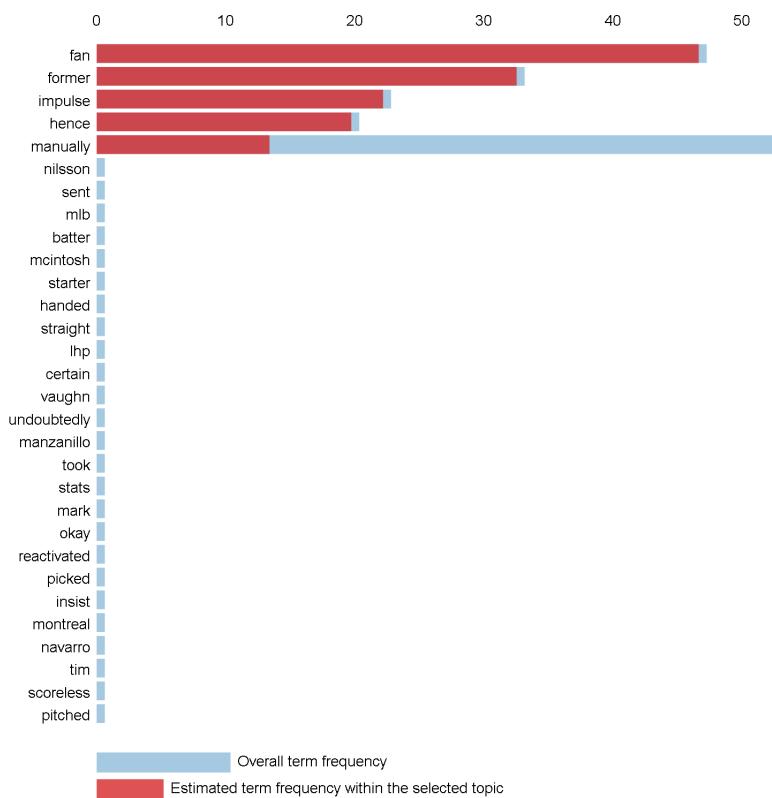
### Top-30 Most Relevant Terms for Topic 18 (0.2% of tokens)



### Top-30 Most Relevant Terms for Topic 19 (0.2% of tokens)

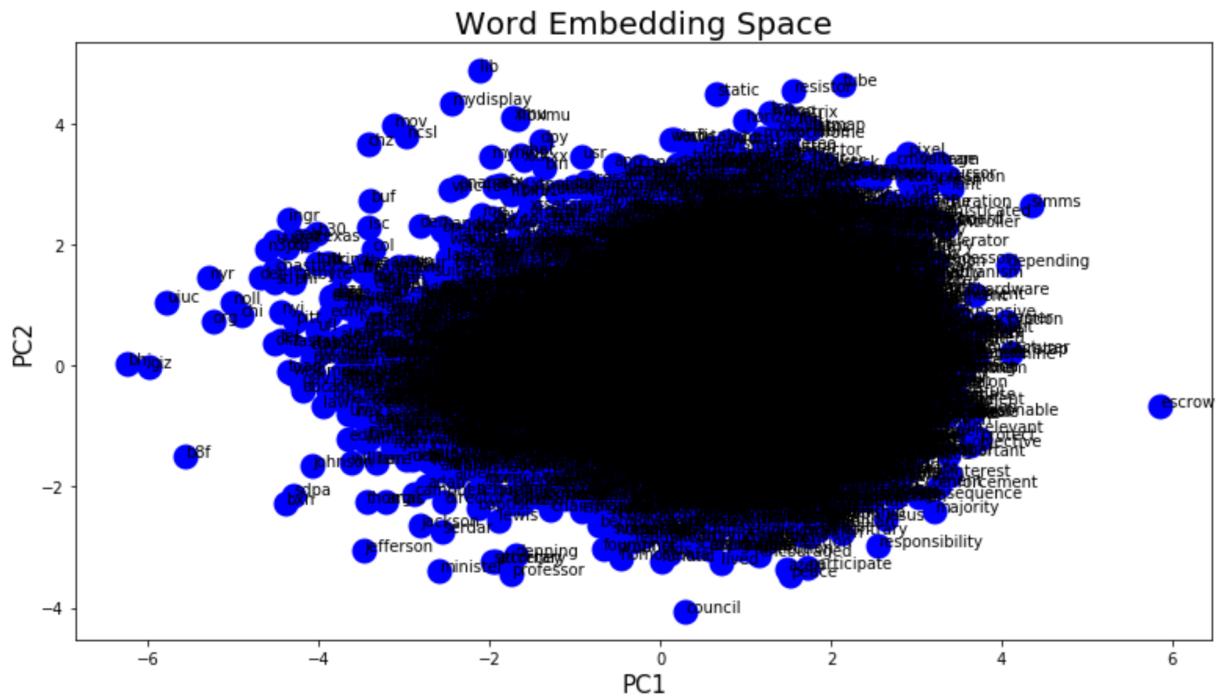


### Top-30 Most Relevant Terms for Topic 20 (0.1% of tokens)

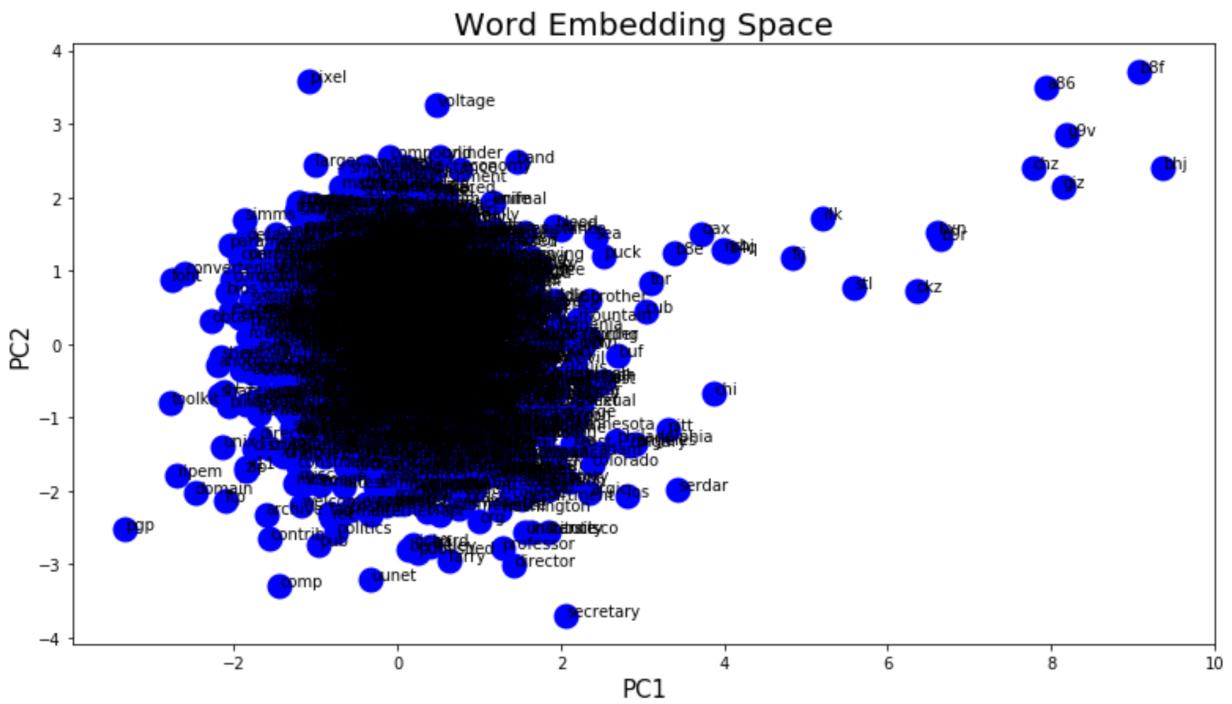


## Appendix F: Word2Vec Visualizations

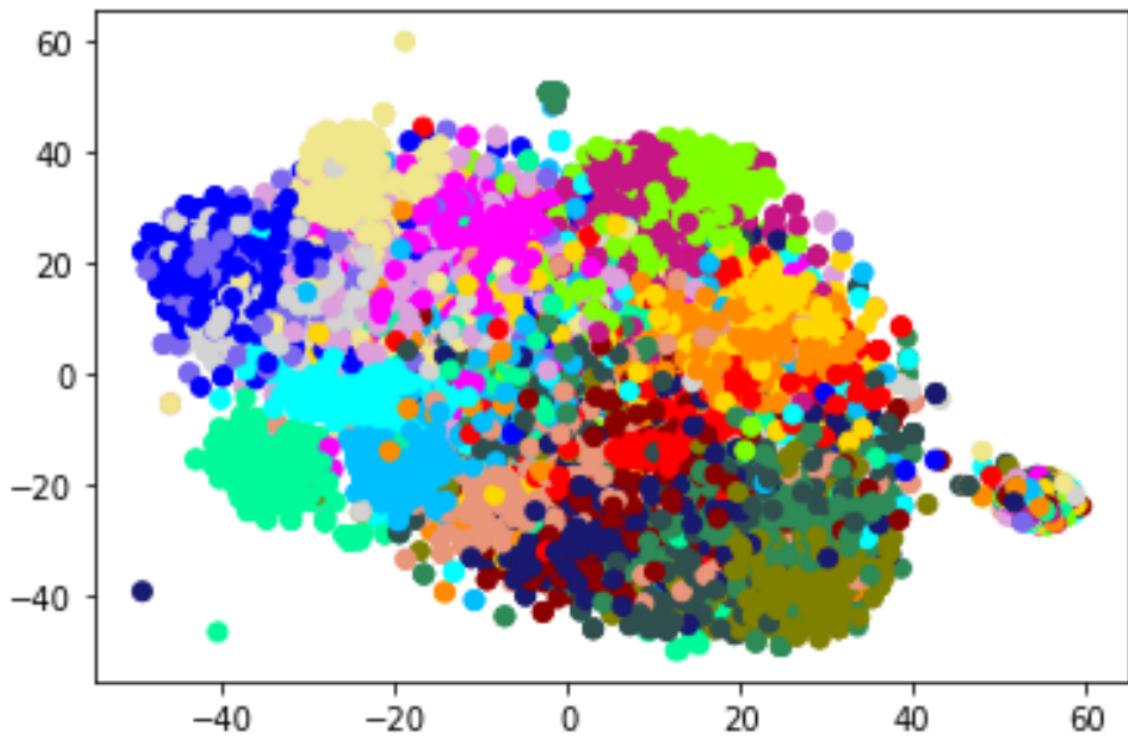
## Complete vocabulary:



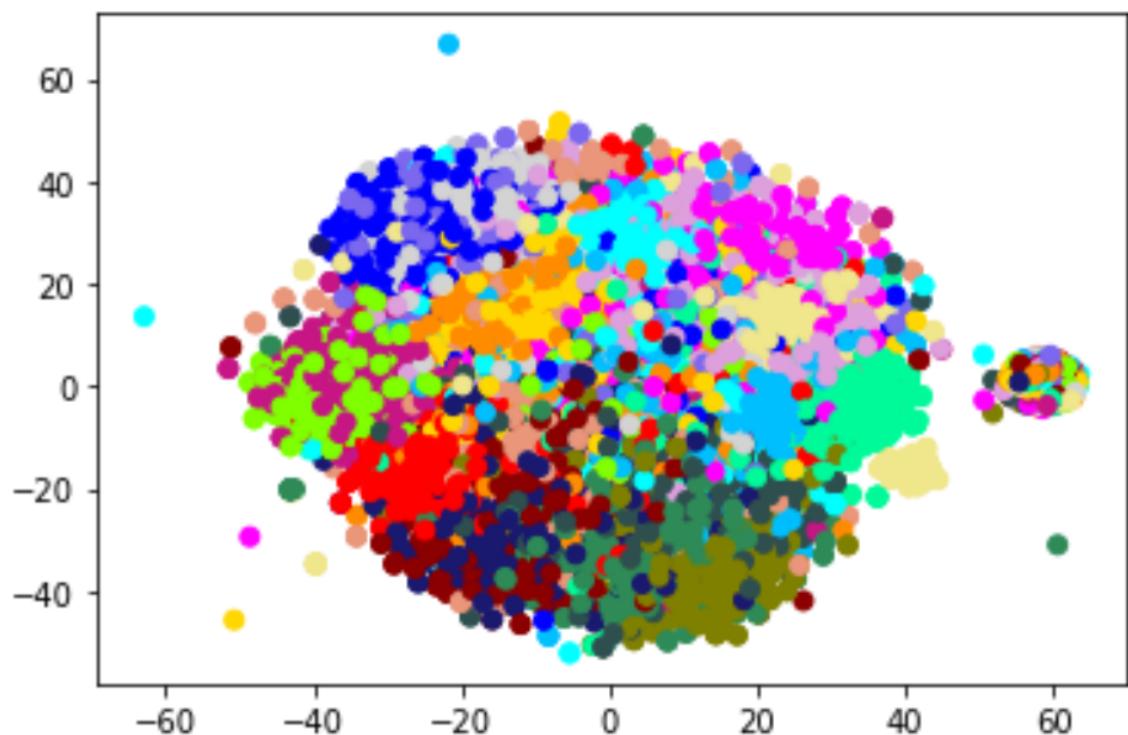
## Top 2k most frequent words:



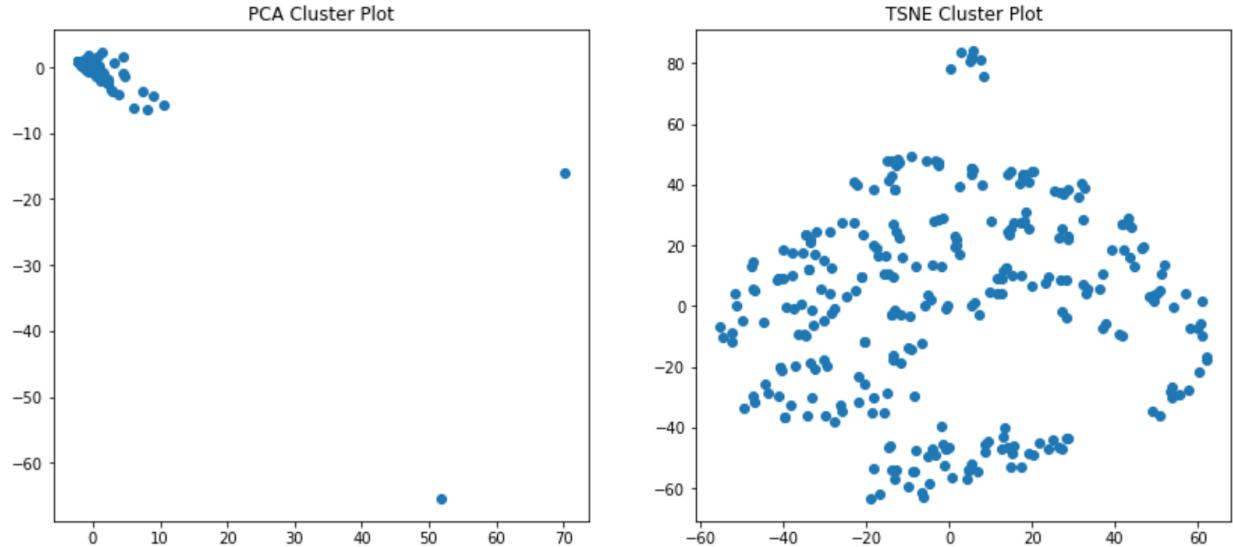
Doc2Vec for the complete dataset



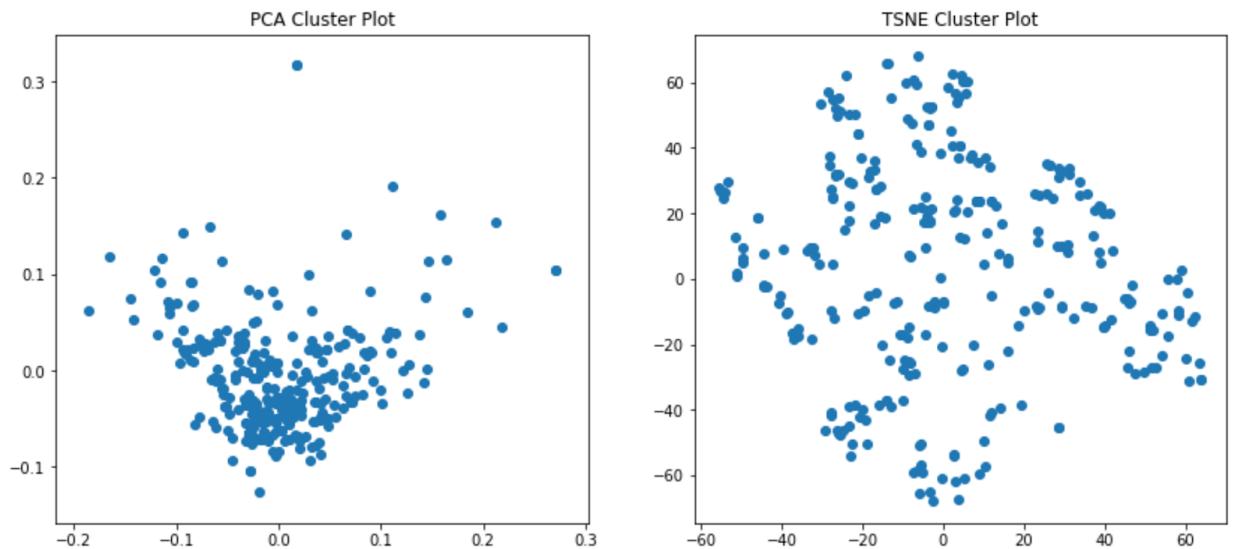
Doc2Vec for the 2k most frequent words



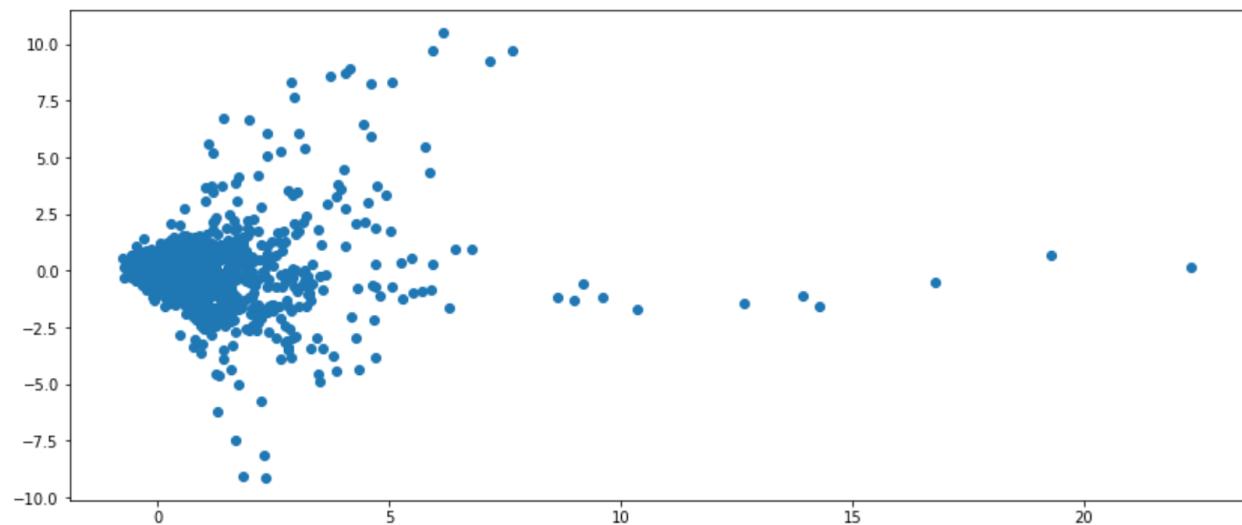
Bag-of-Words:



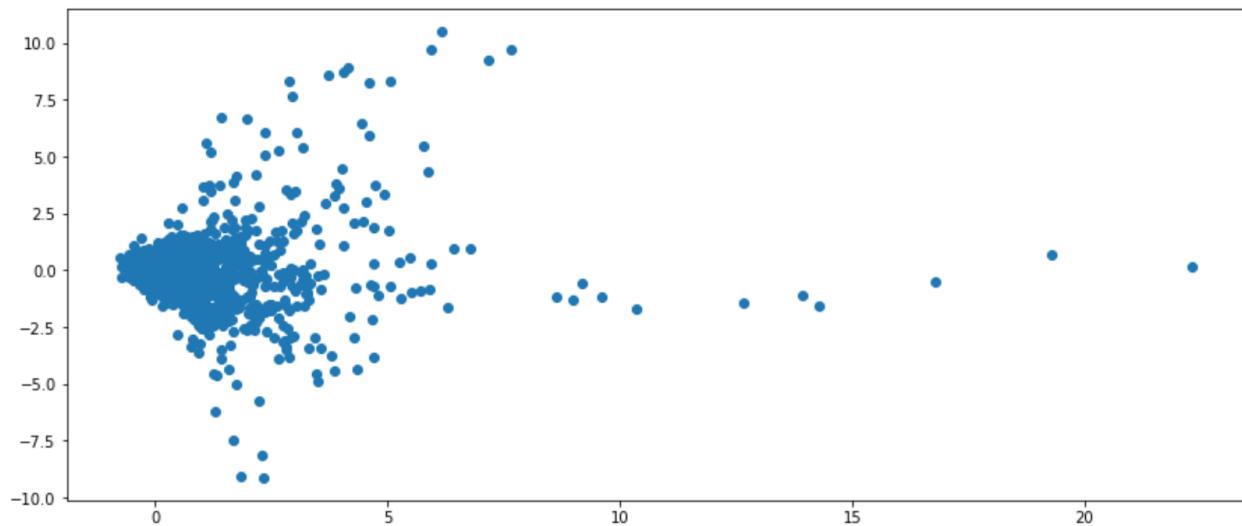
TF-IDF:



Topic Distribution:



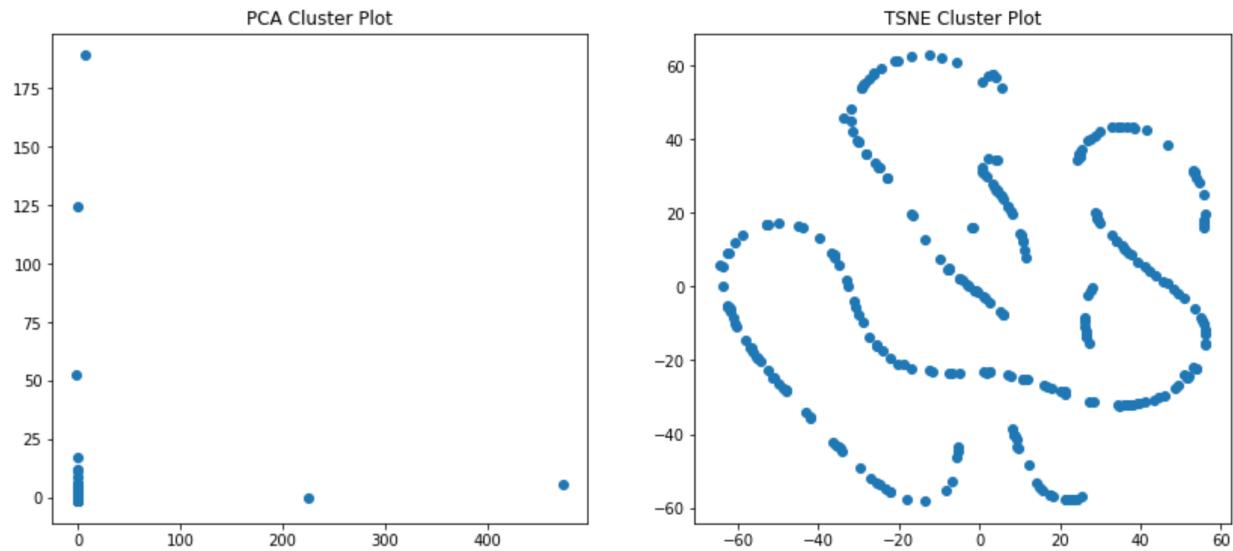
Doc2Vec:



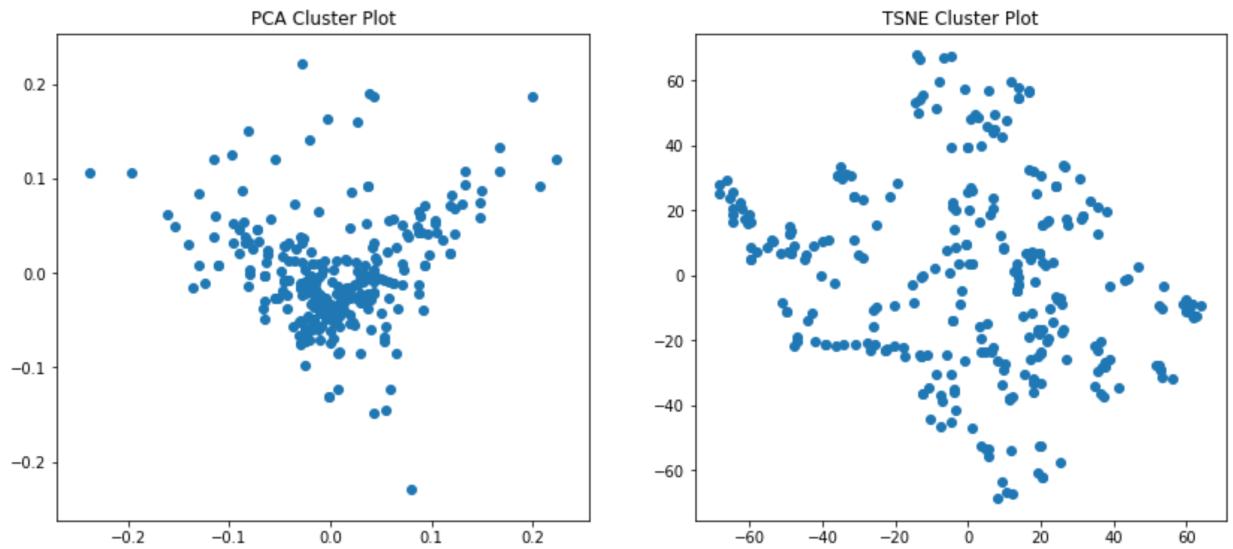
NMI Comparison:

Bag Of Words	TFIDF	Topic Distribution	Doc2Vec
0.00872521360354	0.1624854029375	0.0257184904549024	0.1182999314946

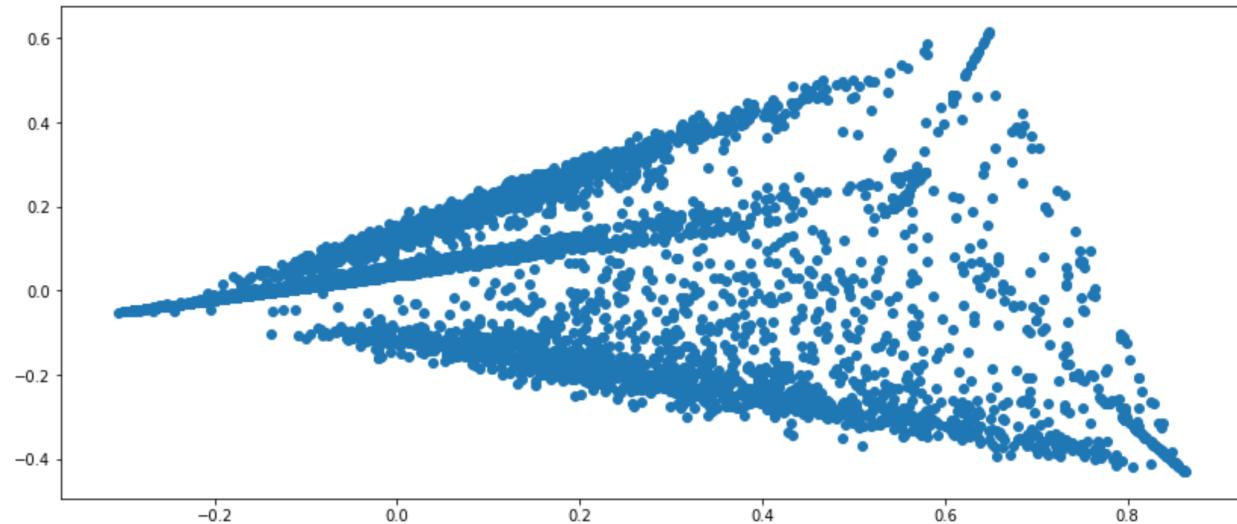
Bag-of-words with 2k vocabulary:



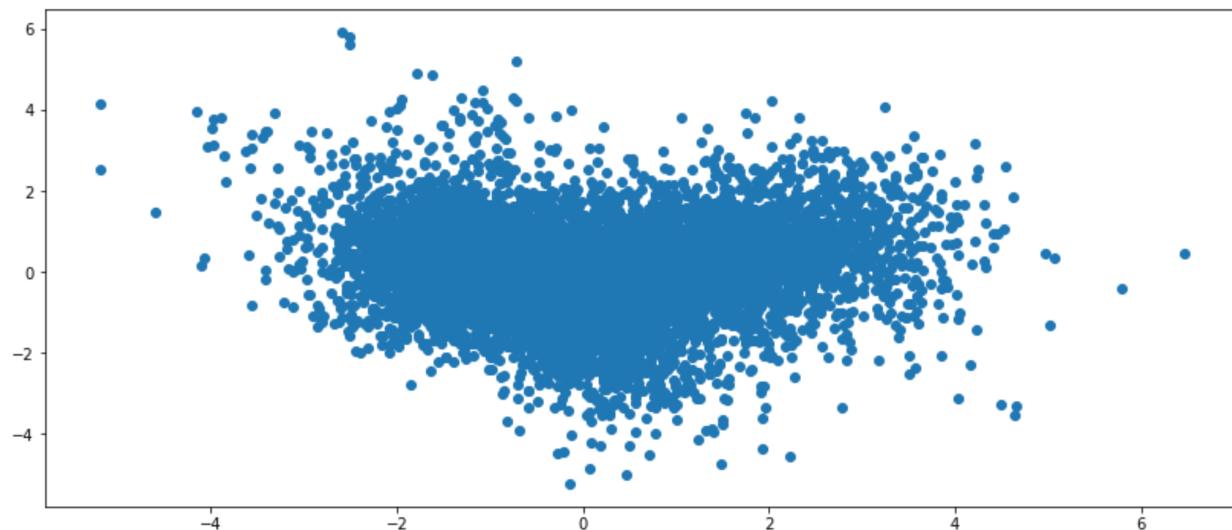
TF-IDF with 2k vocabulary:



Topic Distribution with 2k vocabulary:



Doc2Vec with 2k vocabulary:

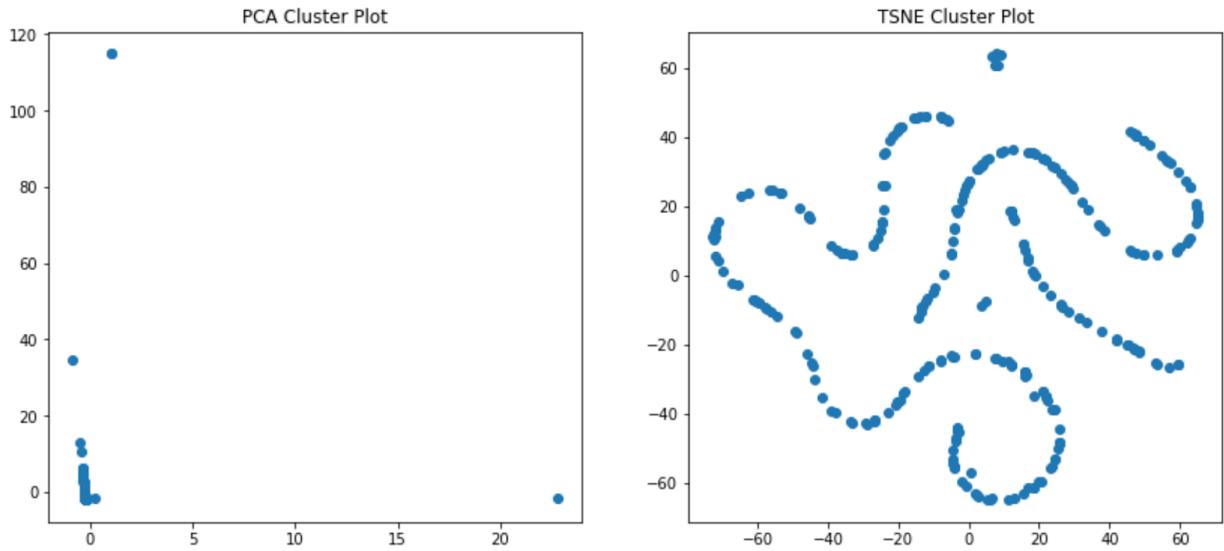


NMI comparison:

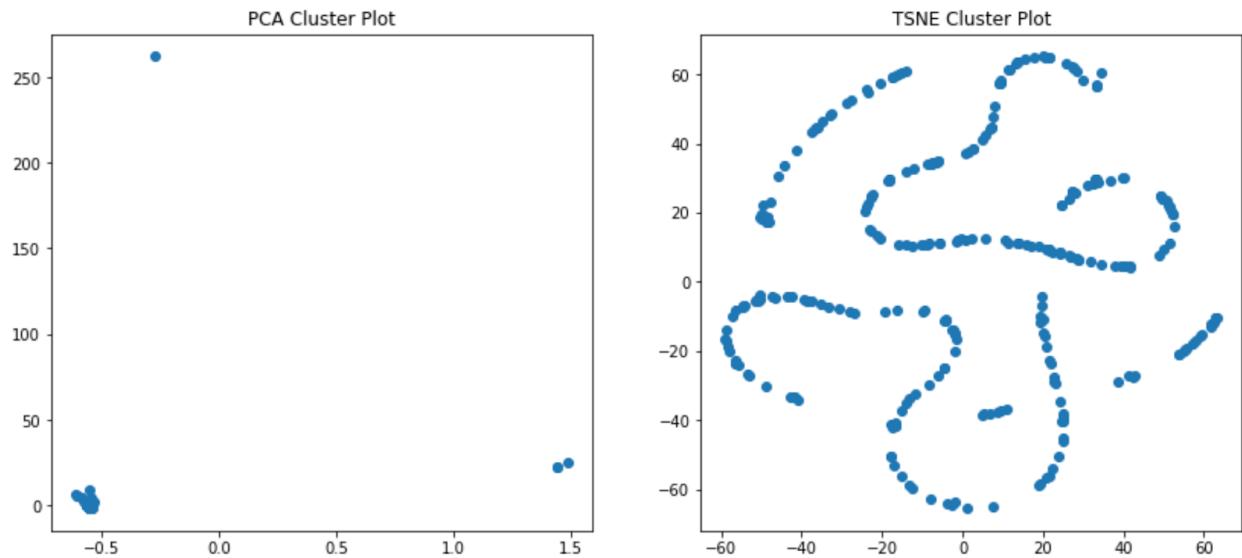
Bag Of Words	TFIDF	Topic Distribution	Doc2Vec
0.00944974065913	0.3171378125559	0.218634846381	0.393188671043

## Appendix I: Results with different number of topics

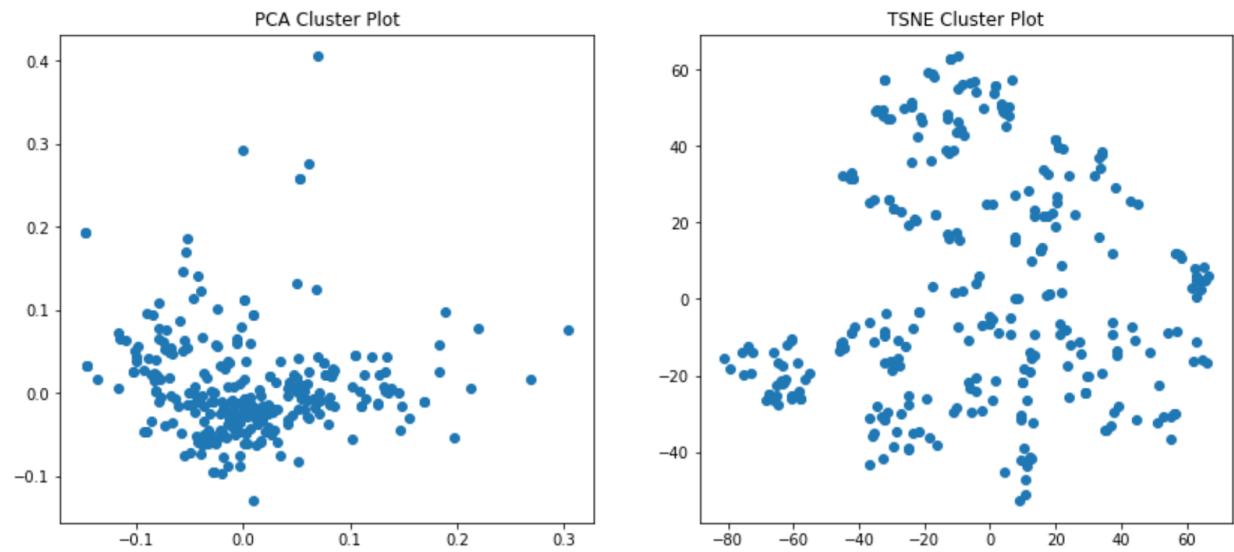
Bag-of-Words with Num\_Topics / 2:



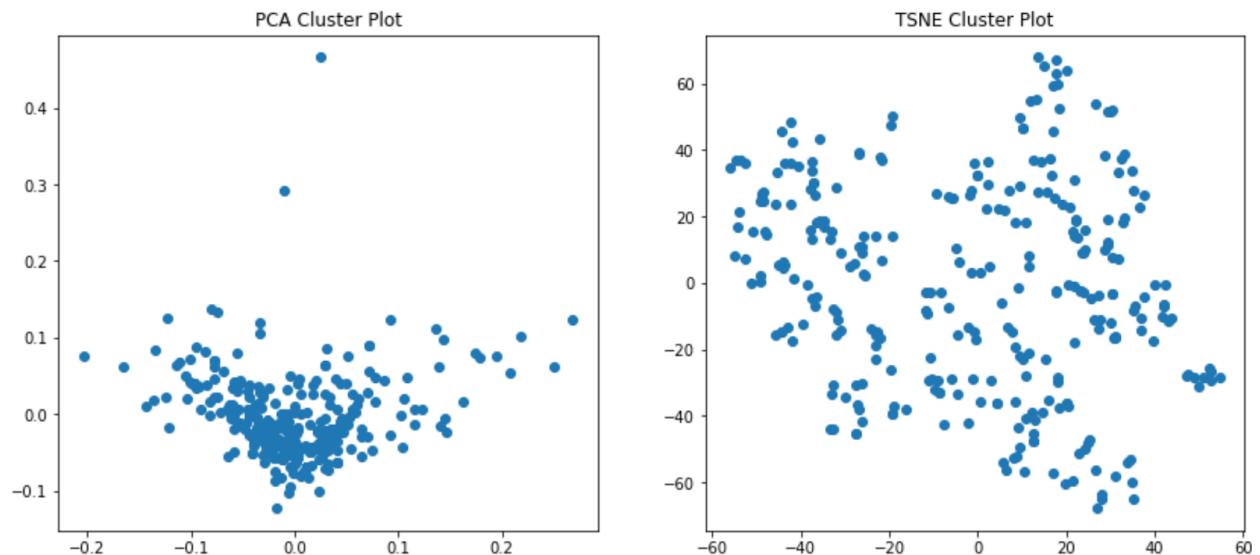
Bag-of-Words with Num\_Topics \* 2:



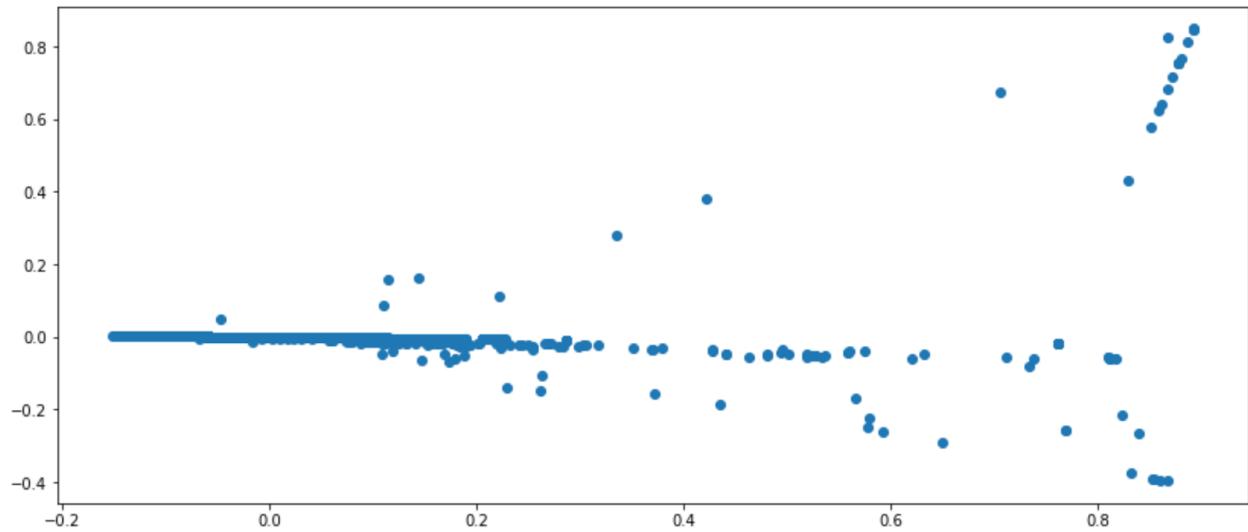
TF-IDF with Num\_Topics / 2:



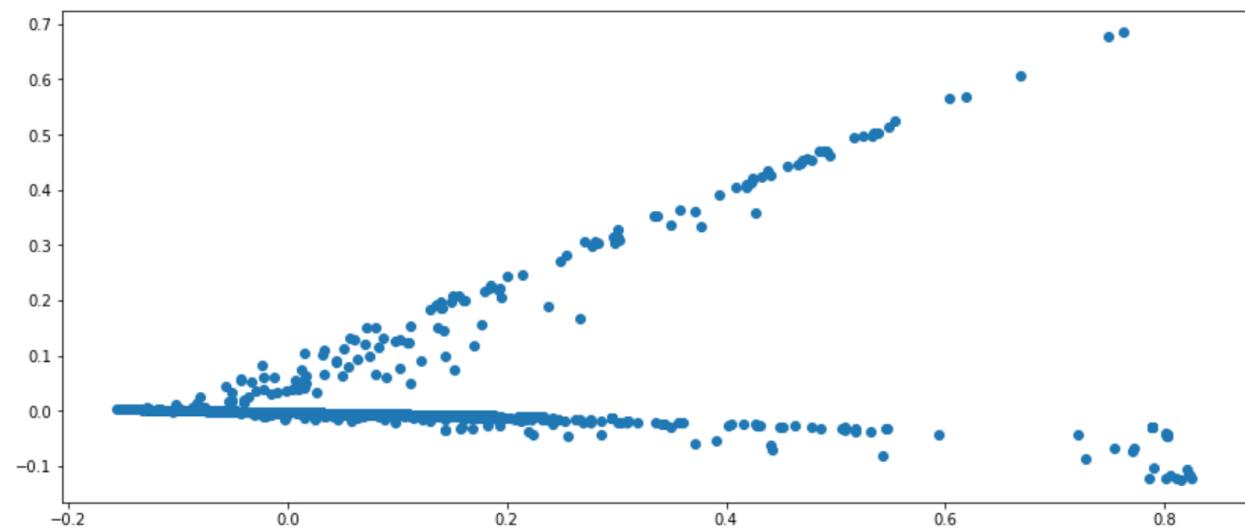
TF-IDF with Num\_Topics \* 2:



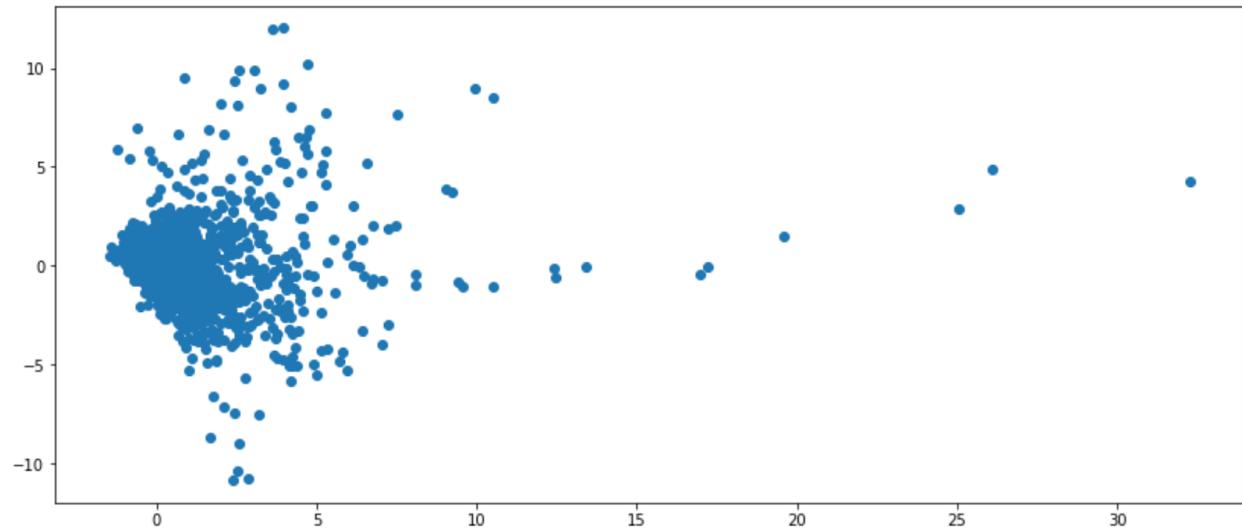
Topic Distribution with Num\_Topics / 2:



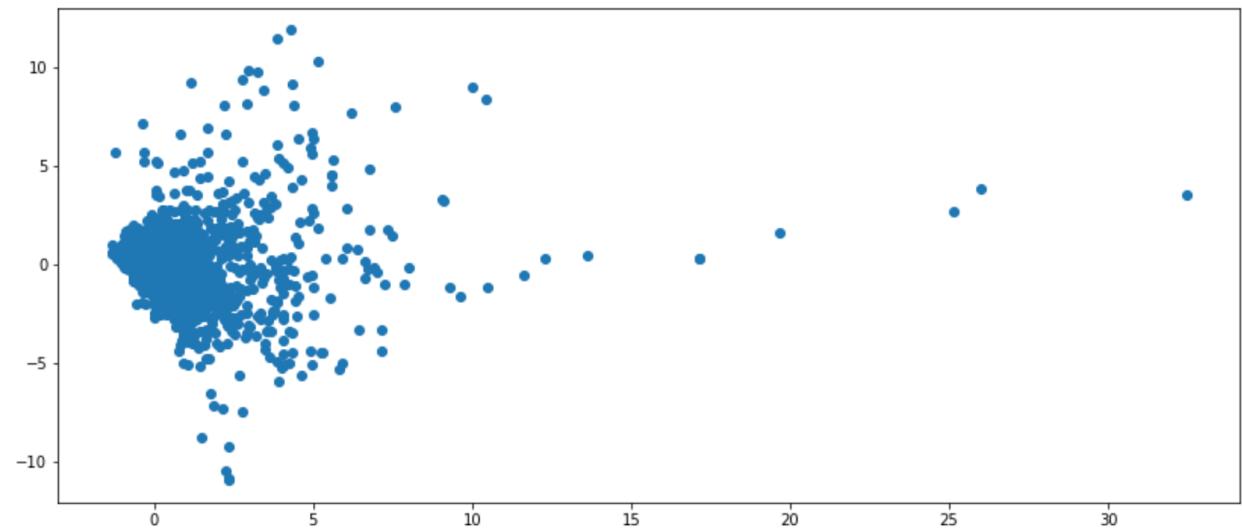
Topic Distribution with Num\_Topics \* 2:



Doc2Vec with Num\_ Topics / 2:



Doc2Vec with Num\_ Topics \* 2:



NMI for Num\_Topics / 2:

Bag Of Words	TFIDF	LdaTopicDistribution	Doc2Vec
0.00424420872286	0.280170443491	0.0213768287910803	0.1020092837936

NMI for Num\_Topics \* 2:

Bag Of Words	TFIDF	LdaTopicDistribution	Doc2Vec
0.02795049095245	0.1161656844488	0.0355767802146235	0.226622236262