

LECTURE 08

Text Visualization

A collection of numerous index cards, each featuring a single word printed in black ink. The words are scattered across a plain white background. Some words are oriented vertically, while others are horizontal. The words include: live, un, they, am, men, puppy, sad, screaming, is, her, frantic, ing, do, pes, blood, leg, you, from, at, eggs, sag, y, need, live, has, soar, delirious, then, think, still, whisper, delicate, to, not, drunk, these, spra, have, finger, as 'e, and, an, have, me, pole, through, could, breast, cook, this, honey, language, let, your, sleep, there, take, boy, must, &, produce, some, head, mean, size, tongue, out, trust, my, said, above, moon, the, smooth, like, read, behind, rock, pan, ride, the, sympt, egg, go, stare, am, skin, blow, all, juice, day, bed, petal, drool, friend, ask, just, drive.

Text And Document

- ▶ We have huge resources of text information:
 - ▶ Libraries
 - ▶ Email archives
 - ▶ WWW documents
 - ▶ Twitter feeds

Can we visualize these resources?

Levels Of Text Representations

- ▶ **Lexical Levels**
 - ▶ Transforming a string of characters into a sequence of atomic entities called *tokens* (e.g., characters, words, phrases)
- ▶ **Syntactic Level**
 - ▶ Identifying and tagging each token's function.
- ▶ **Semantic level**
 - ▶ Extraction of meaning and relationship between pieces of knowledge derived from the structure identified in the syntactic level.



[Google Search](#)

[I'm Feeling Lucky](#)

Google offered in: [日本語](#)

Vector Space Model

The Vector Space Model

There is a great deal of controversy about the safety of genetically engineered foods. Advocates of biotechnology often say that the risks are overblown. “There have been 25,000 trials of genetically modified crops in the world, now, and not a single incident, or anything dangerous in these releases,” said a spokesman for Adventa Holdings, a UK biotech firm. During the 2000 presidential campaign, then-candidate George W. Bush said that “study after study has shown no evidence of danger.” And Clinton Administration Agriculture Secretary Dan Glickman said that “test after rigorous scientific test” had proven the safety of genetically engineered products.

The Vector Space Model

- ▶ The paragraph contains 98 string tokens, 74 terms, and 48 terms when stop words are removed. Here is a sample of the term vector that would be generated by the pseudocode:

genetically	said	safety	engineered	study	test	great	deal	controversy	foods
3	3	2	2	2	2	1	1	1	1

Tasks Of Text Analysis

- ▶ Which documents are similar to a specific one?
- ▶ Which documents are relevant to a given collection?
- ▶ Which documents are relevant to a given search query?
- ▶ Finding patterns or structures
- ▶ Finding a document's main theme

Single Document

Tag Clouds

The tag cloud illustrates the frequency of specific terms across three main categories:

- Food and Agriculture:** author, biotechnology, contained crops, danger, diet, dr, earthsave, eating, engineered, foods, ge genetically, incident, labeled, life, monitoring, monsanto, products, proven, releases, researcher, risks, rissler, safety, sequence, soybeans, study, test, trials, unfortunately, validity, vegetarian, wall, wild, world.
- Music and Culture:** 00s, 60s, 70s, 80s, 90s, acoustic, albums i own, alternative, alternative metal, american, anime, atmospheric, avant-garde, awesome, beautiful, black metal, blues, blues rock, britpop, brutal death metal, canadian, celtic, chill, chillout, christian, classic, classic rock, classical, comedy, country, cover, dance, dark ambient, darkwave, death metal, disco, doom metal, downtempo, drum and bass, dub, east, electronic, electronica, emo, experimental, favorite, favorite songs,芬兰语, electro, female vocalists, finnish, folk, folk metal, folk rock, french, favorites, favourite, favourites, female, female vocalist, funk, german, goth, gothic, gothic metal, gothic rock, grindcore, grunge, guitar, hard rock, hardcore, heavy, hip hop, house, idm, indie, indie pop, indie rock, industrial, industrial metal, metal, stening, ebay, electro.
- Cultural and Social:** indie, indie pop, indie rock, industrial, industrial metal.

Wordle



<http://www.wordle.net/>

Word Tree

17
hits

the

most part you wouldn't know [if you were eating them] but the point being that you wouldn't need to know .
point being that you wouldn't need to know .

safety of genetically engineered foods .
products .

risks are overblown .

world , now , and not a single incident , or anything dangerous in these releases , " said a spokesman for adventa holdings , a uk
[conari press , 2550 ninth st .

2000 presidential campaign , then - candidate george w .

case ?

union of concerned scientists , dr .

epa , she is one of the nation's leading authorities on the environmental risks of genetically engineered foods .

nation's leading authorities on the environmental risks of genetically engineered foods .

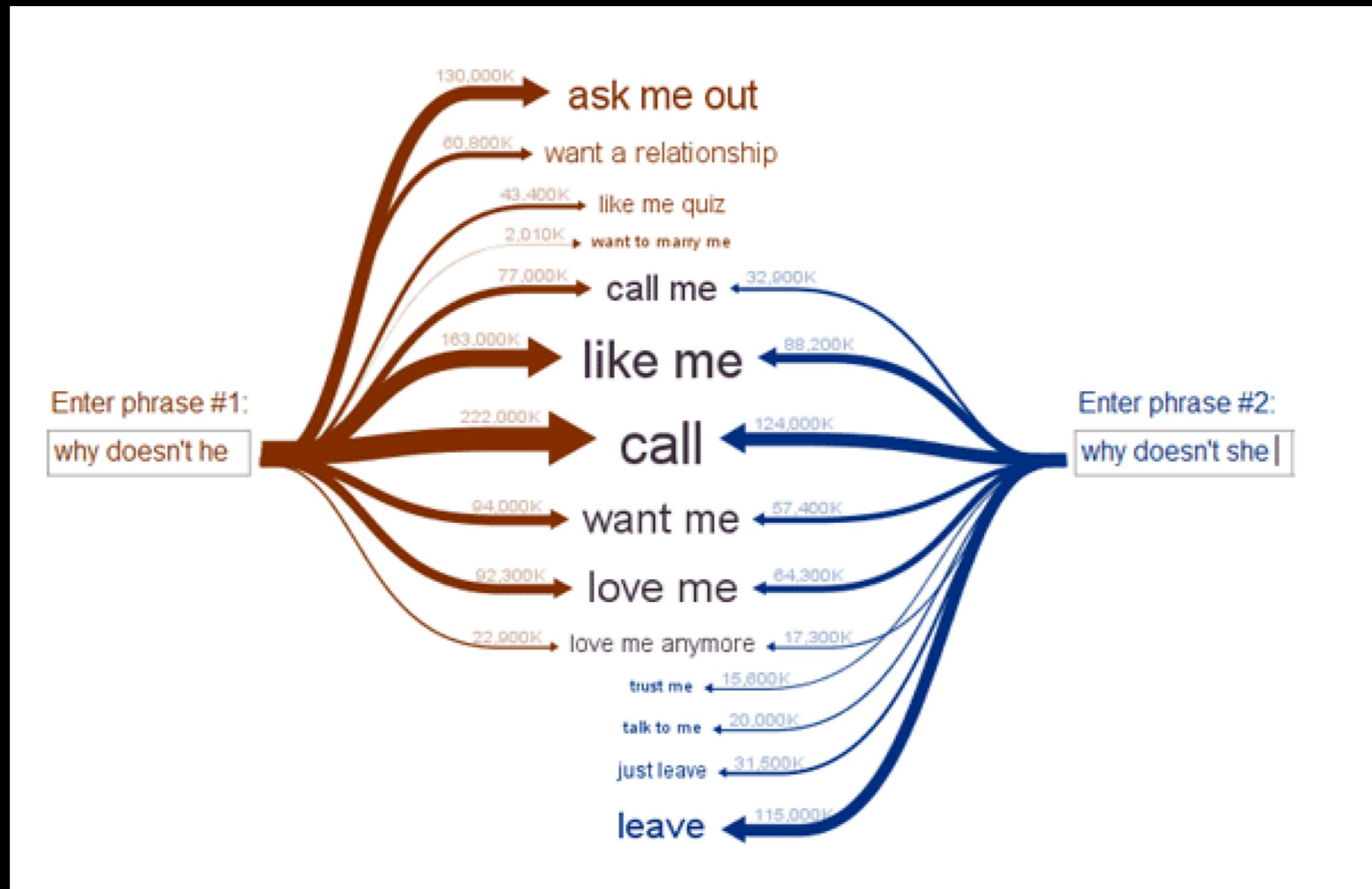
environmental risks of genetically engineered foods .

trials and studies .

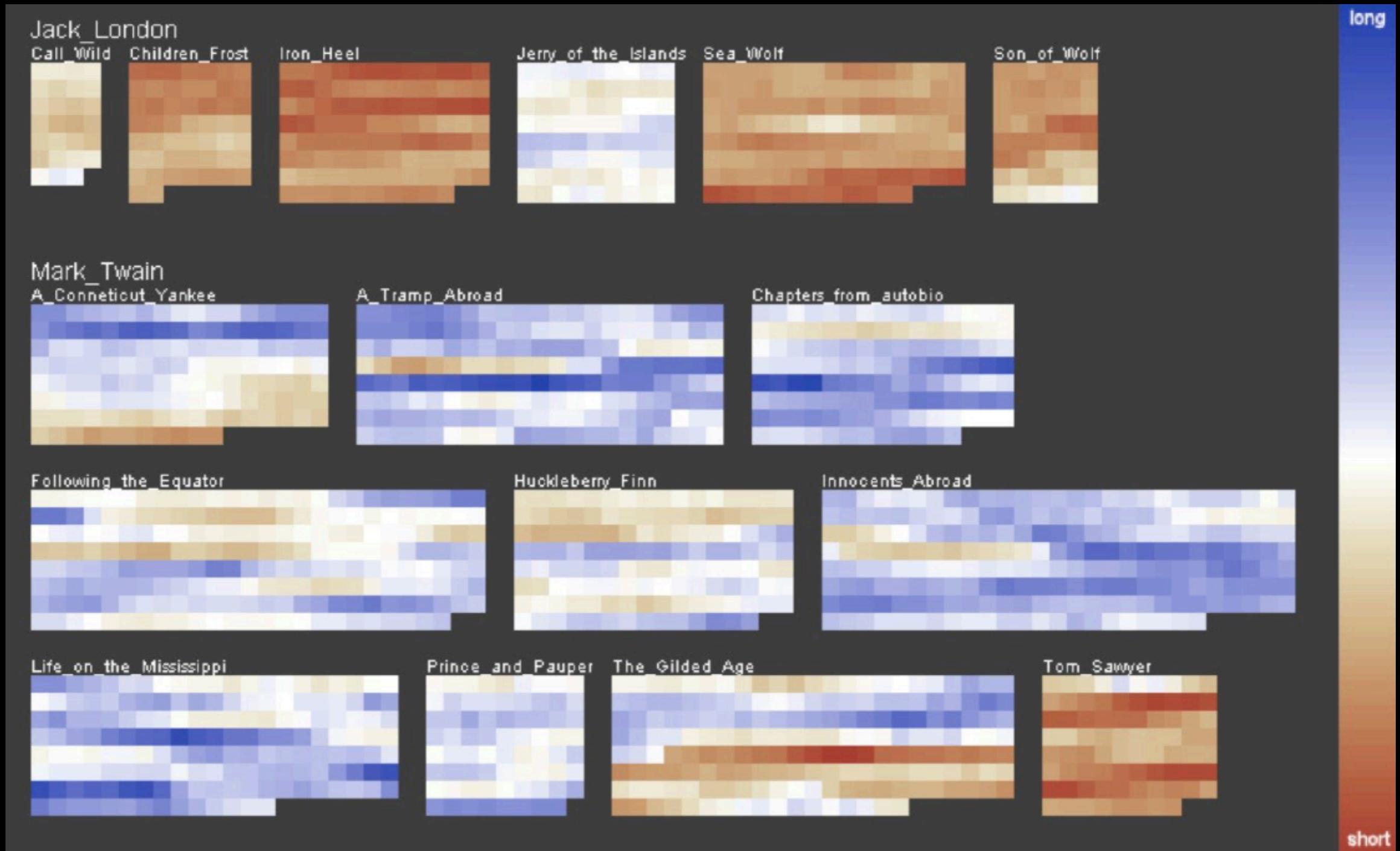
wall street journal found that 16 of 20 vegetarian foods labeled as being " free " of genetically engineered products actually contained ge soybeans .

validity of all their safety tests ?

author of diet for a new america and founder of earthsave international .

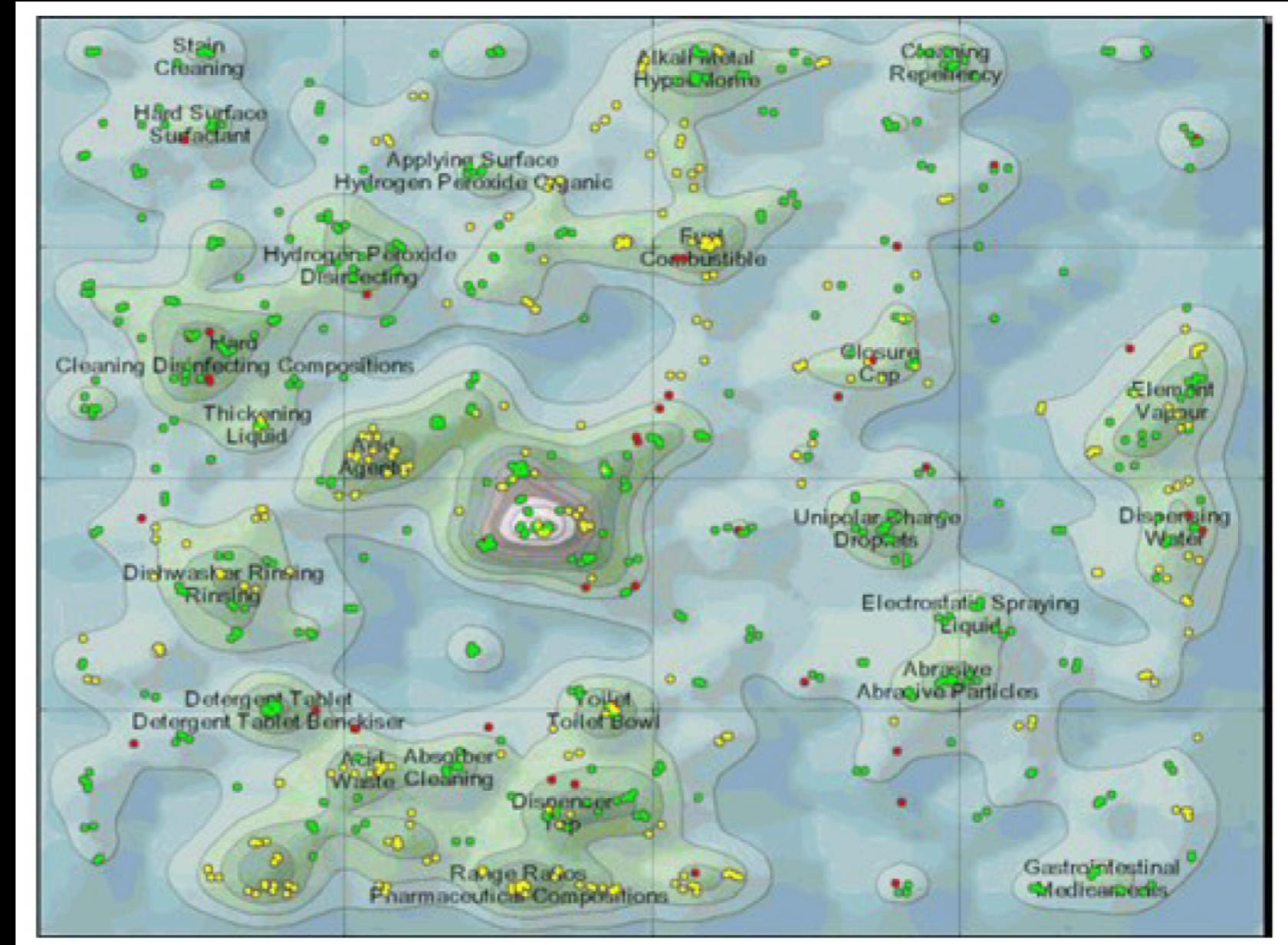
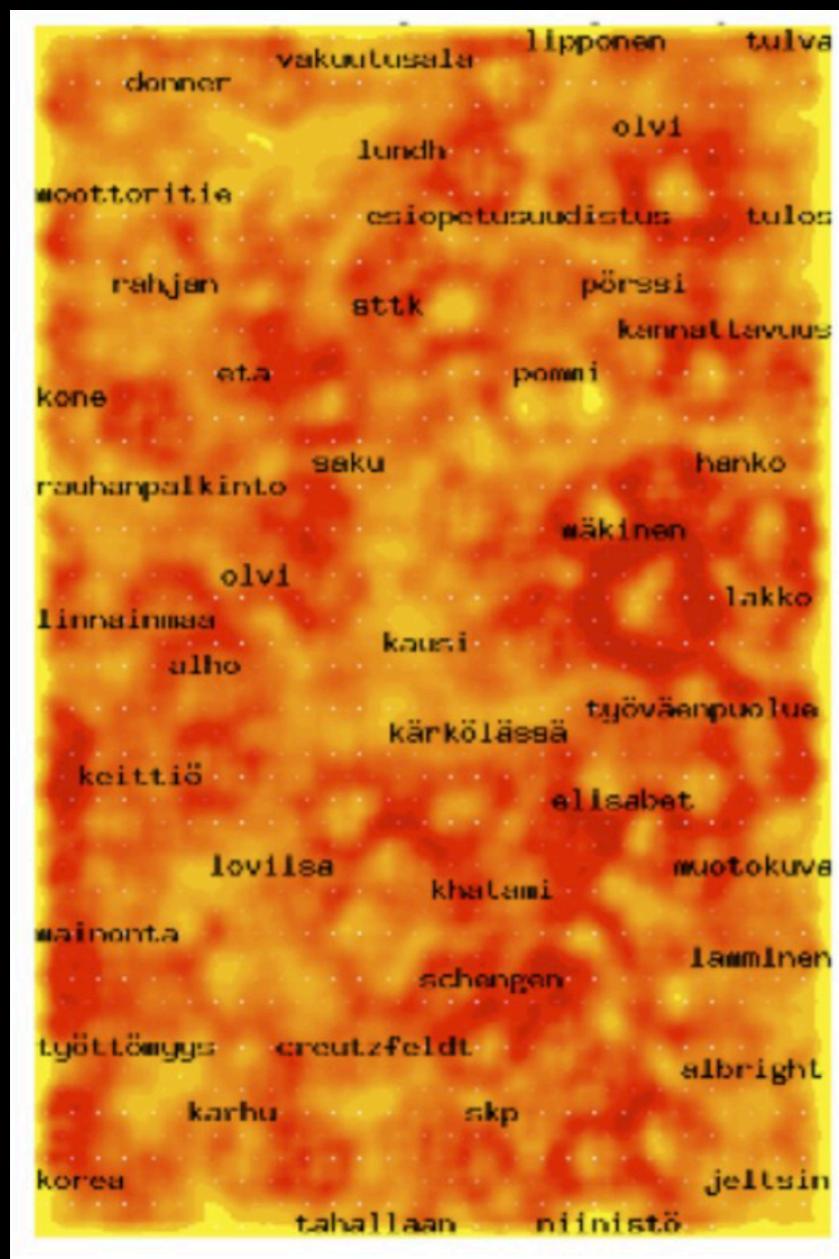


Literature Fingerprinting



Sentence Length

Document Collection



The labels show the topical areas, and color represents the number of documents, with light areas containing more.

Joined Analysis

Google NGram

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

