# SUPPLEMENTAL MATERIAL FOR „VISUALIZING EXPANDED QUERY RESULTS"

**Michael Mazurek and Manuela Waldner**
**TU Wien**

Here, we present detailed supplemental material to the EuroVis paper "Visualizing Expanded Query Results". Each section number refers to the associated section in the paper.

## 8.1 TASKS

We selected 10 ambiguous topics from the TREC web tracks 2009 – 2014. For each topic, we performed the given query and presented two different sub-topics to the users. We always chose the main description, as well as one dissimilar sub-topic, which could be resolved by ConceptNet.

```
<topic number="99" type="ambiguous">
  <query>satellite</query>
  <description>
    Find background information about man-made satellites.
  </description>
  <subtopic number="1" type="inf">
    Find background information about man-made satellites.
  </subtopic>
  <subtopic number="2" type="nav">
    Find satellite maps and geographic images.
  </subtopic>
  <subtopic number="3" type="nav">
    Find providers of satellite television.
  </subtopic>
  <subtopic number="4" type="inf">
    Find information about satellite telephones.
  </subtopic>
  <subtopic number="5" type="nav">
    Find providers of satellite internet service.
  </subtopic>
  <subtopic number="6" type="nav">
    Find providers of satellite radio systems.
  </subtopic>
</topic>
```

Below is the complete list of topics and sub-topic descriptions selected for the study:
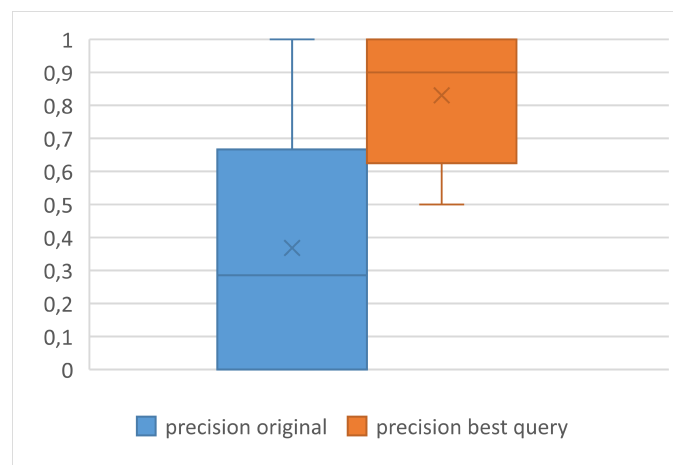
| Task set | Query | Description | Precision in original query | Query with highest precision | Total number of relevant hits |
|---|---|---|---|---|---|
| 1 | grilling | Find recipies for grilling.* | 0.67 | Q8: 1.0 | 27 |

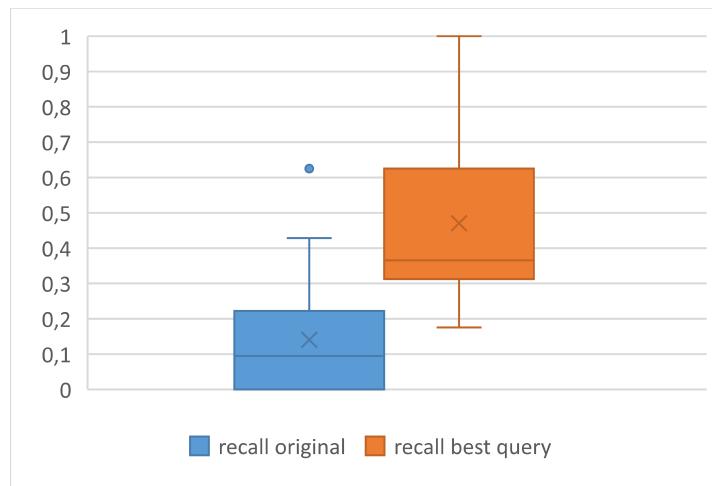| 1 | rock art | Find information on cave paintings all around the world.* | 0.75 | Q0: 0.75 | 20 |
|---|---|---|---|---|---|
| 1 | iron | Find information about iron as an essential nutrient. | 0.56 | Q0: 0.56 | 8 |
| 1 | Worm | Find information about worms in nature. | 0.56 | Q4: 0.9 | 40 |
| 1 | Kiwi | Find information on kiwi fruit. | 0.33 | Q4,Q5: 1.0 | 32 |
| 2 | Worm | Find information about computer worms, viruses, and spyware. | 0.00 | Q1: 1.0 | 17 |
| 2 | Pvc | Find information about PVC pipes and fittings. | 0.29 | Q4: 1.0 | 16 |
| 2 | Kiwi | Find information on kiwi birds. | 0.11 | Q2: 1.0 | 31 |
| 2 | Keyboard review | Find reviews of computer keyboards. | 0.88 | Q1, Q5: 1.0 | 57 |
| 2 | Joints | Find an explanation of the different types of joints used in woodworking. | 0.00 | Q9: 0.8 | 8 |
| 3 | grilling | Find information on different type of barbecue grills.* | 0.00 | Q4: 0.5 | 7 |
| 3 | Pvc | How are premature ventricular contractions treated? | 0.14 | Q6: 1.0 | 16 |
| 3 | Joints | Find information about joints in the human body. | 1.00 | Q0, Q6: 1.0 | 32 |
| 3 | Satellite | Find providers of satellite television hardware.* | 0.00 | Q6: 0.6 | 10 |
| 3 | Dog heat | What is the effect of excessive heat on dogs? | 0.00 | Q8, Q9: 0.8 | 24 |
| 4 | rock art | Where can I learn about rock painting or buy a rock-painting kit? | 0.25 | Q1,Q5: 7 | 21 |
| 4 | Iron | Find information about the element iron (Fe). | 0.44 | Q5: 0.7 | 33 |
| 4 | Keyboard review | Find reviews of electronic keyboards and digital pianos. | 0.00 | Q7: 0.5 | 6 |
| 4 | Satellite | Find background information about man-made satellites. | 0.67 | Q1, Q2: 0.8 | 27 |
| 4 | Dog heat | Find information on dogs' reproductive cycle. What does it mean when a dog is "in heat"? | 1.00 | Q0: 1.0 | 21 |

* Slightly modified from original TREC description.

Mind, that the number of document surrogates for the original query is lower than 10 in our implementation, because the original query results are parsed directly from the Google results list, which is usually shorter than the list of 10 document surrogates delivered by the API (average number of document surrogates for the original query is 8.7 in our experiment).

The average precision of the *best* queries (see second to last column in table above) is 0.83 (median: 0.9). The average precision of the original query is 0.37 (median: 0.29):



The average recall of the best queries is 0.47 (median: 0.37). The average recall of the original query is 0.14 (median: 0.09):

## 8.2 APPARATUS AND PROCEDURE

The complete procedure for the study was as follows:

- consent form
- demographic questionnaire
- task description
- for each interface:
  - warm-up task (queries: spider [program] and jaguar [cat])
  - task set with five sub-topics shown in random order
- questionnaire

Each task was preceded by a presentation of the query, together with its description:



After clicking the Query-button, the Google result page with the respective visualization (here Parallel Tag Clouds) was shown:

Task assignments to interfaces and presentation order of interfaces was balanced using a Graeco-Latin Square:

T = text, Q=Euler Diagram, P = Parallel Tag Clouds, L = Parallel Lists

| User | Con.1 | Task1 | Con.2 | Task2 | Con.3 | Task3 | Con.4 | Task4 |
|------|-------|-------|-------|-------|-------|-------|-------|-------|
| 1 | T | 1 | Q | 3 | P | 4 | L | 2 |
| 2 | Q | 2 | T | 4 | L | 3 | P | 1 |
| 3 | P | 3 | L | 1 | T | 2 | Q | 4 |
| 4 | L | 4 | P | 2 | Q | 1 | T | 3 |
| 5 | T | 1 | Q | 3 | P | 4 | L | 2 |
| 6 | Q | 2 | T | 4 | L | 3 | P | 1 |
| 7 | P | 3 | L | 1 | T | 2 | Q | 4 |
| 8 | L | 4 | P | 2 | Q | 1 | T | 3 |
| 9 | T | 1 | Q | 3 | P | 4 | L | 2 |
| 10 | Q | 2 | T | 4 | L | 3 | P | 1 |
| 11 | P | 3 | L | 1 | T | 2 | Q | 4 |
| 12 | L | 4 | P | 2 | Q | 1 | T | 3 |
| 13 | T | 1 | Q | 3 | P | 4 | L | 2 |
| 14 | Q | 2 | T | 4 | L | 3 | P | 1 |
| 15 | P | 3 | L | 1 | T | 2 | Q | 4 |
| 16 | L | 4 | P | 2 | Q | 1 | T | 3 |

## 8.4 PARTICIPANTS

Responses from the demographic questionnaire:



Age
16 responses



Sex
16 responses

## Eye Sight

16 responses



- normal — 50%
- corrected to normal — 50%
- not corrected

## Domain (e.g., computer science, biology)

16 responses



| Computer Science | IT | chemistry | computer science | med vis |
|---|---|---|---|---|
| 1 (6.3%) | 2 (12.5%) | 1 (6.3%) | 11 (68.8%) | 1 (6.3%) |

## How often do you have to use online search engines (e.g., Google or Bing) for your work?

16 responses



- daily — 100%
- very often, but not every day
- sometimes
- never

How often do you auto-complete a query suggested by the search engine? (see example below)

16 responses



How often do you select an expanded query suggested by the search engine? (see example below)

16 responses



How familiar are you with visualization techniques (e.g., bar charts or pie charts)?

16 responses



## 8.5 ANALYSIS AND RESULTS

16 users conducted 5 tasks with 4 different interfaces, each, resulting in 320 samples for task completion time (TCT) and number of hits.

We first analyzed TCT and the number of hits for outliers. There were no outliers for number of hits, but we removed 18 samples (from 6 different users), because they were outliers in terms of TCT. Below, there is a box plot before outlier removal:

The same box plot after removing outliers:



For the remaining 302, we computed precision (number of relevant hits in selected query divided by number of document surrogates in the selected query) and recall (number of relevant hits in the selected query divided by the overall number of relevant hits in all queries). We aggregated the 302 per user and condition, resulting in 64 average TCT, precision, and recall values.

For precision, recall, and TCT, we conducted a repeated measures ANOVA with condition as within-subjects factor. In case of significance, we performed pairwise Bonferroni-corrected post-hoc comparisons between the 4 conditions. For the user ratings, we performed a Friedman test with Bonferroni-corrected Wilcoxon Signed-Rank post-hoc comparisons.

## Task Completion Time

There is a significant difference of TCT between the conditions:

**Tests of Within-Subjects Effects**

Measure:   TCT

| Source | | Type III Sum of Squares | df | Mean Square | F | Sig. | Partial Eta Squared |
|---|---|---|---|---|---|---|---|
| condition | Sphericity Assumed | 4098,352 | 3 | 1366,117 | 34,923 | ,000 | ,700 |
| | Greenhouse-Geisser | 4098,352 | 2,323 | 1764,284 | 34,923 | ,000 | ,700 |
| | Huynh-Feldt | 4098,352 | 2,774 | 1477,368 | 34,923 | ,000 | ,700 |
| | Lower-bound | 4098,352 | 1,000 | 4098,352 | 34,923 | ,000 | ,700 |
| Error(condition) | Sphericity Assumed | 1760,327 | 45 | 39,118 | | | |
| | Greenhouse-Geisser | 1760,327 | 34,844 | 50,520 | | | |
| | Huynh-Feldt | 1760,327 | 41,611 | 42,304 | | | |
| | Lower-bound | 1760,327 | 15,000 | 117,355 | | | |

Condition 4 (text) had a significantly lower TCT than all other conditions:

**Pairwise Comparisons**

Measure:   TCT

| (I)condition | (J)condition | Mean Difference (I-J) | Std. Error | Sig.[b] | 95% Confidence Interval for Difference[b] | |
|---|---|---|---|---|---|---|
| | | | | | Lower Bound | Upper Bound |
| 1 | 2 | -5,618 | 2,604 | ,286 | -13,526 | 2,290 |
| | 3 | -,926 | 1,781 | 1,000 | -6,334 | 4,483 |
| | 4 | 15,633[*] | 2,639 | ,000 | 7,621 | 23,644 |
| 2 | 1 | 5,618 | 2,604 | ,286 | -2,290 | 13,526 |
| | 3 | 4,692 | 2,280 | ,344 | -2,230 | 11,614 |
| | 4 | 21,251[*] | 2,128 | ,000 | 14,791 | 27,711 |
| 3 | 1 | ,926 | 1,781 | 1,000 | -4,483 | 6,334 |
| | 2 | -4,692 | 2,280 | ,344 | -11,614 | 2,230 |
| | 4 | 16,558[*] | 1,642 | ,000 | 11,573 | 21,544 |
| **4** | **1** | **-15,633[*]** | 2,639 | ,000 | -23,644 | -7,621 |
| | **2** | **-21,251[*]** | 2,128 | ,000 | -27,711 | -14,791 |
| | **3** | **-16,558[*]** | 1,642 | ,000 | -21,544 | -11,573 |

Based on estimated marginal means

*. The mean difference is significant at ,05 level.

b. Adjustment for multiple comparisons: Bonferroni.

## Precision

There is no significant difference in precision between the conditions:

**Tests of Within-Subjects Effects**

Measure: precision

| Source | | Type III Sum of Squares | df | Mean Square | F | Sig. | Partial Eta-Squared |
|---|---|---|---|---|---|---|---|
| condition | Sphericity Assumed | ,109 | 3 | ,036 | **1,992** | **,129** | **,117** |
| | Greenhouse-Geisser | ,109 | 2,297 | ,048 | 1,992 | ,146 | ,117 |
| | Huynh-Feldt | ,109 | 2,735 | ,040 | 1,992 | ,135 | ,117 |
| | Lower Bound | ,109 | 1,000 | ,109 | 1,992 | ,179 | ,117 |
| Error(condition) | Sphericity Assumed | ,824 | 45 | ,018 | | | |
| | Greenhouse-Geisser | ,824 | 34,451 | ,024 | | | |
| | Huynh-Feldt | ,824 | 41,031 | ,020 | | | |
| | Lower Bound | ,824 | 15,000 | ,055 | | | |

In addition, we also compared the precision between the *original* query (i.e., the precision without performing any query expansion), the *best* query (i.e., the optimal solution of all sub-topics), and the aggregated precision per task set and visualization of the users' *selected* queries from all conditions in the visualization.

We compared the precision of selected queries to the original query and the best query using a Kruskal-Wallis H test. The difference is significant:

**Test Statistics[a,b]**

|  | precision |
|---|---|
| Chi-square | **25,404** |
| df | 2 |
| Asymp. Sig. | **,000** |

a. Kruskal-Wallis-Test

b. Grouping Variable: query

We therefore compared the precision of the selected query to the original and the best query using Mann-Whitney U tests.

There is a significant difference between the selected query and the original query:

**Test Statistics[a]**

|  | precision |
|---|---|
| Mann-Whitney-U | 344,000 |
| Wilcoxon-W | 554,000 |
| Z | **-3,111** |
| Asymp. Sig. (2-sided) | **,002** |

a. Grouping Variable: query

There is also a significant difference between the selected query and the best query:

**Test Statistics[a]**

|  | precision |
|---|---|

| | |
|---|---|
| Mann-Whitney-U | 263,000 |
| Wilcoxon-W | 2343,000 |
| Z | **-3,964** |
| Asymp. Sig. (2-sided) | **,000** |

a. Grouping Variable: query

This means that users could improve the precision of the retrieved documents by expanding the query. However, the selection of query expansions was not optimal.

## Recall

There is also no difference between the conditions for recall:

**Tests of Within-Subjects Effects**

Measure: recall

| Source | | Type III Sum of Squares | df | Mean Square | F | Sig. | Partial Eta-Squared |
|---|---|---|---|---|---|---|---|
| condition | Sphericity Assumed | ,025 | 3 | ,008 | **,919** | **,439** | **,058** |
| | Greenhouse-Geisser | ,025 | 1,941 | ,013 | ,919 | ,439 | ,058 |
| | Huynh-Feldt | ,025 | 2,225 | ,011 | ,919 | ,439 | ,058 |
| | Lower Bound | ,025 | 1,000 | ,025 | ,919 | ,439 | ,058 |
| Error(condition) | Sphericity Assumed | ,407 | 45 | ,009 | | | |
| | Greenhouse-Geisser | ,407 | 29,113 | ,014 | | | |
| | Huynh-Feldt | ,407 | 33,371 | ,012 | | | |
| | Lower Bound | ,407 | 15,000 | ,027 | | | |

## User Ratings

User ratings were not normally distributed. We therefore compared the ratings using a Friedman test. The Friedman test was significant:

**Statistics for Test<sup>a</sup>**

| N | 16 |
|---|---|
| **Chi-Square** | **15,396** |
| **df** | **3** |
| **Asymp. Sig.** | **,002** |

a. Friedman-Test

We therefore conducted pairwise Wilcoxon Signed-Rank post-hoc comparisons with a Bonferroni-adjusted α of 0.0083. The results show that PTC was rated significantly lower than text and lists.

**Statistics for Test[a]**

|  | euler - text | **PTC - text** | lists - text | PTC - euler | lists - euler | **lists - PTC** |
|---|---|---|---|---|---|---|
| Z | -,537[b] | **-2,829[b]** | ,000[c] | -2,300[b] | -,476[d] | **-3,119[d]** |
| Asymp. Sig. (2-sided) | ,591 | **,005** | 1,000 | ,021 | ,634 | **,002** |

a. Wilcoxon-Test

b. Based on positive ranks.

c. The sum of negative ranks equals the sum of positive ranks.

d. Based on negative ranks.

## User Feedback

Here, we list all user comments given for the final questionnaire:

**Text list (text)**

- easy to scan quickly, look for specifying keywords related to original query
- Fast overview, but not very detailed
- clean and linear alignment. quick overview based on the words. the sense is made by own experience and association of words
- simple but with less content
- i m quite used to the current way of queries in google, but i disliked that i couldnt actually see what each quey meant (i was still not sure what i should select)
- + Easy to understand - not much additional information
- simple for fast seeing, bad for choices
- fast to read
- fast to find query, very used to it
- it's fast
- very easy to understand, not cluttered, visualization not appealing

**Compact Euler Diagram (euler)**

- hard to read
- Much too complex

- spatial alignment of blocks allows me to mentally associate the blocks with a position in space. I liked the clustering/bundelling of keywords together into blocks. this makes it easy to mentally divide between the meaning of words inside the blocks.
- straight forward but need to study how to use it
- i really liked that the most important terms in the queries were very easy to see at a glance, and with the links between the different queries, i could choose very easily what was the best result for me
- + good representation of query expansion - slow / optically not that great
- good: interaction within the keywords, con:need a bit of time to understand how the representation works
- blocks provide a good overview, takes more time to understand
- a little bit crowded
- also fast finding
- best overview of connected terms
- visualization is nice, but confusing

### Parallel Tag Clouds (PTC)
- hard to read
- hard to make connections
- initially, I was not aware that there are pillars associated with the expanded queries. the association is not quite easily visible, however highlighting helps. i did not know what the size of words is encoding. also the gradient confused me.
- all keywords are listed but hard to find its relationship
- i liked the fact that i could immediately see the most important keywords and based on the connections between the different queries i could choose the 'final' result
- - not very intuitive, requires a lot of space on screen
- nice idea - too chaotic
- lot of redundant information
- a little bit confusing
- too many words to read, size of text distracts if the searched term is not very large.
- didn't get a good overview
- good overview, easy to select and understand, connections are interesting

### Lists View (lists)
- nice to see what keywords the different suggested queries had in common
- appealing view, easier to make connections
- i liked the left view, since it allows me to see an aggregated query and the optional expansion with the right view where more keywords can be seen. the connection are a bit hard to see without highlighting, so it needs interaction to expand. also, i am not sure how scalable this is, when there is a lot of different suggestions (however, this is true for most other visualizations)
- more text about the keywords but difficult to read when more items are associated
- i think that this is a good way of having an overview of the most important queries and what they mean, but i had to go through all of them to decide which one was the best match
- + very good representation of query expansion, easy to understand and handle - takes a lot of space on screen
- easy to see and follow
- I like it

- sometimes hard to identify which words on the right side belong together
- right column is rather confusing
- nice visualization, easy to understand and select, interesting info about overlaps

## Categorized Utterances

Below, we list the categories revealed during open coding, and the number of positive (blue) and negative (orange) utterances associated with these categories for each condition:

# 8 STIMULI

Here, all stimuli for each topic and interface condition from the user study are listed:

## dog heat

Parallel List



Euler Diagram

## Text

List of alternate Queries for **dog heat**

**dog companionship heat**

**dog heat warm**

**dog heat energy**

**dog heat warmth**

**dog a loyal friend heat**

**dog heat fire**

**dog mammal heat**

**dog pet heat**

**dog animal heat**

## Parallel Tag Clouds



| dog heat | dog companionship heat | dog heat warm | dog heat energy | dog heat warmth | dog a loyal friend heat | dog heat fire | dog mammal heat | dog pet heat | dog animal heat |

# Grilling

Parallel List

| | |
|---|---|
| grilling | barbecue grills patio charcoal outdoor |
| grill framework | grillwork grillworks material decorative wood |
| grill cook | steak steaks grilled recipe lab |
| grill grillwork | framework project liferay river based |
| grill barbecue | recipes charcoal grills outdoor grilling |
| grill examine | grille germantown words barrier noun |
| grill grille | jobs cook turkey resume salary |
| grill cooking | cruel cruelty gifs live unusual |
| grill grilling steak | examine cross interrogate put synonyms |
| grill cruel | grille capital animals root bar |

Euler Diagram

Text

List of alternate Queries for **grilling**

**grill framework**

**grill cook**

**grill grillwork**

**grill barbecue**

**grill examine**

**grill grille**

**grill cooking**

**grill grilling steak**

**grill cruel**

Parallel Tag Clouds

## Iron

Parallel List

| | |
|---|---|
| iron | element facts fe earth basic |
| iron a metal | journal steel research jax international |
| iron ironing your clothes | clothes ironing robot board start |
| iron clothes | ni appliances steam home irons |
| iron steel | elements earth crust abundant fraction |
| iron an element | clothes fist fashion clothing ironing |
| iron an abundant element upon the earth | golf ball hit hitting shots |
| iron home appliance | scrap liberty steel earth metal |
| iron hitting a golf ball | cloth inch ironing pressing buy |
| iron pressing cloth | pressing cloth ironing press minky |

Euler Diagram

Text

**iron a metal**

**iron ironing your clothes**

**iron clothes**

**iron steel**

**iron an element**

**iron an abundant element upon the earth**

**iron home appliance**

**iron hitting a golf ball**

**iron pressing cloth**

Parallel Tag Clouds



| abundant | abundant | board | clothes | element | abundant | abundant | appliances | ball | board |
| basic | crust | cloth | clothing | international | basic | crust | buy | golf | buy |
| earth | earth | clothes | fashion | jax | crust | earth | clothes | hit | cloth |
| element | element | clothing | fist | journal | earth | element | home | hitting | clothes |
| elements | fe | ironing | home | metal | element | elements | ironing | irons | home |
| facts | home | metal | ironing | research | elements | facts | irons | shots | inch |
| fe | international | robot | irons | steel | facts | fe | metal | | ironing |
| metal | liberty | start | metal | | fe | fraction | ni | | minky |
| steel | metal | | press | metal | metal | metal | press | | press |
| | scrap | | | | fe | steel | steam | | pressing |
| | steel | | | | | | | | start |

iron    iron a metal    iron ironing your clothes    iron clothes    iron steel    iron an element   iron an abundant element upon the earth   iron home appliance    iron hitting a golf ball    iron pressing cloth

# Joints

Parallel List

| | |
|---|---|
| joints | integrated defense probabilistic organization association |
| joint anatomy | collective bakersfield menu reviews bargaining |
| joint cigarette | carpentry woodworking wood carpenters apprenticeship |
| joint spot | bones joints sacroiliac skeleton synovial |
| joint common | articulatory articulation systems bus system |
| joint articulatory system | cigarette smoking cigarettes marijuana lungs |
| joint a skeleton | tenants property tenancy difference prayer |
| joint integrated | spot chiropractic steel high ultrasonic |
| joint collective | tenants common tenancy individual jtic |
| joint carpentry | spot welding interface bonding welded |

Euler Diagram

Text

Parallel Tag Clouds



| joints | joint anatomy | joint cigarette | joint spot | joint common | joint articulatory system | joint a skeleton | joint integrated | joint collective | joint carpentry |

# keyboard review

Parallel List

| | |
|---|---|
| keyboard review | music debussy rameau play home |
| keyboard device review | multi device bluetooth devices logitech |
| keyboard review accounting | legal legalboard law lawyers adaptxt |
| keyboard review criticism | gaming anandtech trends mechanical professional |
| keyboard a computer system review | system output virtual computer solar |
| keyboard a computer review | sellers top list computer devices |
| keyboard review exercise | mechanical typing oct ultra cover |
| keyboard playing music review | accounting accounts university journal research |
| keyboard review law | macbook apple blackberry musicverse step |
| keyboard typing review | exercise identification piano desk gloves |

Euler Diagram

Text

Parallel Tag Clouds



keyboard review    keyboard device review  keyboard review accounting  keyboard review criticism  keyboard a computer system  keyboard a computer review  keyboard review exercise  keyboard playing music review  keyboard review law    keyboard typing review

## Kiwi

Parallel List



| kiwi | bus login cookies member experience |
| kiwi new zealand | zealander pete confirms citizen barnaby |
| kiwi ratite | mantelli apteryx brown conservation haastii |
| kiwi edible fruit | ratite ratites birds metabolic rates |
| kiwi Chinese gooseberry | inhabitant definition dictionary dollar synonyms |
| kiwi a fruit | kiwifruit benefits health fruit vitamin |
| kiwi inhabitant | baskets skin mango female red |
| kiwi apteryx | kindergarten wien plural harry kinderinwien.at |
| kiwi Kiwi | van camper bus campers smaller |
| kiwi New Zealander | chinese gooseberry kiwifruit fruit china |

Euler Diagram

Text

List of alternate Queries for **kiwi**

**kiwi new zealand**

**kiwi ratite**

**kiwi edible fruit**

**kiwi Chinese gooseberry**

**kiwi a fruit**

**kiwi inhabitant**

**kiwi apteryx**

**kiwi Kiwi**

**kiwi New Zealander**

Parallel Tag Clouds



| kiwi | kiwi new zealand | kiwi ratite | kiwi edible fruit | kiwi Chinese gooseberry | kiwi a fruit | kiwi inhabitant | kiwi apteryx | kiwi Kiwi | kiwi New Zealander |

# pvc

Parallel List

| pvc | premature ventricular contractions beat heart |
|---|---|
| pvc polyvinyl chloride | chloride polyvinyl vinyl made polymerization |
| pvc iv catheter | ventricular premature extrasystoles een contractie |
| pvc extrasystole | virtual circuit permanent medicine chloride |
| pvc abs | abs difference fittings pipe systems |
| pvc artificial substance | catheters peripheral iv catheter intravenous |
| pvc premature ventricular contraction | polymers artificial substances synthetic vinyl |

Euler Diagram

## Text

List of alternate Queries for **pvc**

**pvc polyvinyl chloride**

**pvc iv catheter**

**pvc extrasystole**

**pvc abs**

**pvc artificial substance**

**pvc premature ventricular contraction**

## Parallel Tag Clouds



| pvc | pvc polyvinyl chloride | pvc iv catheter | pvc extrasystole | pvc abs | pvc artificial substance | pvc premature ventricular contraction |
|---|---|---|---|---|---|---|
| catheter | chloride | catheter | beat | abs | artificial | beat |
| chloride | made | catheters | catheter | chloride | chloride | contractions |
| circuit | pipe | intravenous | contractie | difference | made | heart |
| heart | polymerization | iv | contractions | fittings | pipe | premature |
| intravenous | polyvinyl | medicine | een | pipe | polymerization | ventricular |
| made | synthetic | peripheral | extrasystoles | polyvinyl | polymers | |
| medicine | vinyl | | heart | systems | polyvinyl | |
| peripheral | | | premature | | substances | |
| permanent | | | ventricular | | synthetic | |
| pipe | | | | | vinyl | |
| polyvinyl | | | | | | |
| premature | | | | | | |
| synthetic | | | | | | |
| ventricular | | | | | | |
| virtual | | | | | | |

# rock art

Parallel List

| rock art | painted painting crafting forms humble |
|---|---|
| rock stone art | grab round balancing sculpture subs |
| rock roll art | canada canadian reconciliation inuksuk list |
| rock art paintings | museum alta valley area blowing |
| rock and roll art | music festival cellars castoro whale |
| rock art painting | trust prehistoric tara follow african |
| rock art sculpture | roll punk target bang show |
| rock canada art | balancing beautiful stone rocks ideas |
| rock art museum | blowing museum preserve van witte |
| rock music art | roll punk bang show artists |

Euler Diagram

**rock music art**
castoro whale festival cellars

**rock art sculpture**
subs round
grab

**rock roll art**
music
punk
show
bang target
artists
roll
museum

**rock canada art**
list canadian
inuksuk
reconciliation
canada

**rock stone art**
beautiful
painted area
balancing sculpture

**rock art**
painting
ideas
rocks
stone
prehistoric
valley
african
follow trust
tara

**rock art paintings**
crafting humble
forms

**rock art museum**
preserve witte van
alta blowing

Text

List of alternate Queries for **rock art**

**rock stone art**

**rock roll art**

**rock art paintings**

**rock and roll art**

**rock art painting**

**rock art sculpture**

**rock canada art**

**rock art museum**

**rock music art**

Parallel Tag Clouds

# Satellite

Parallel List

| | |
|---|---|
| satellite | equipment internet review operating iowa |
| satellite artificial satellite | air pollution sentinel quality monitoring |
| satellite follower | orbiter mars reconnaissance mro education |
| satellite air | spy secret military spacex sep |
| satellite orbiter | planet company mars feb named |
| satellite spying | artificial fireworks distant geodesy phantom |
| satellite equipment | celestial body bodies planets navigation |
| satellite apparatus | follower formation leader flying motion |
| satellite planet | orbiter solar mission esa lro |
| satellite celestial body | equipment tv dishes dish receivers |

Euler Diagram

Text

List of alternate Queries for **satellite**

**satellite artificial satellite**

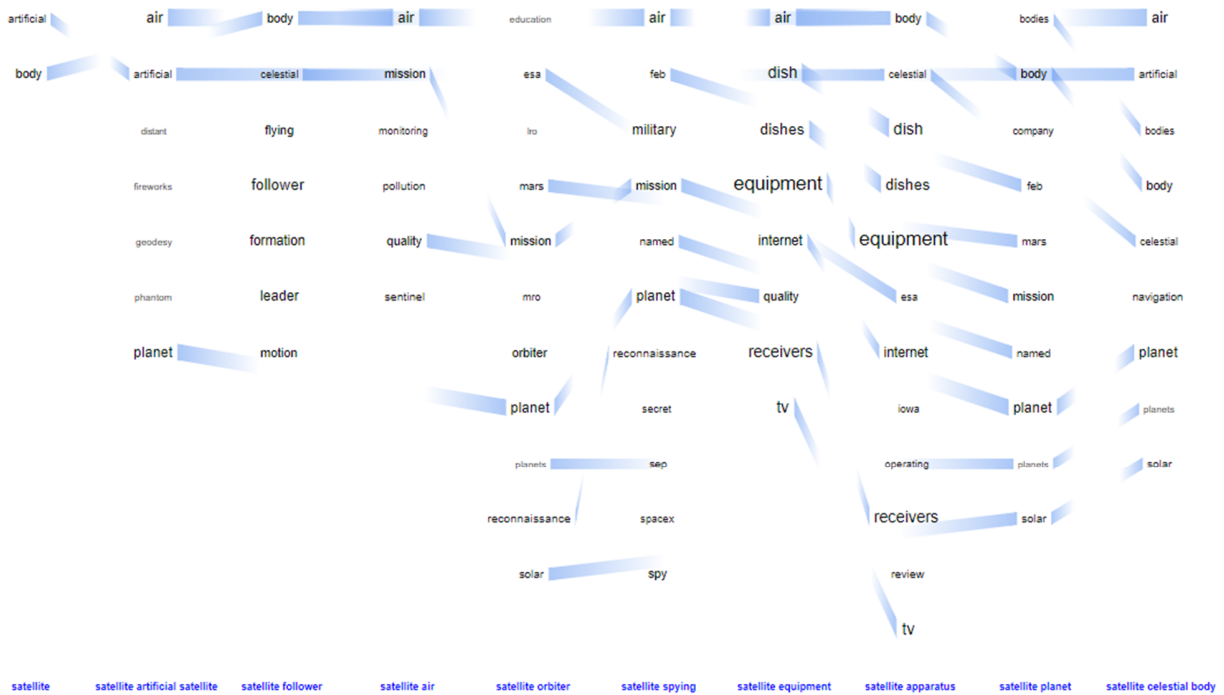**satellite follower**

**satellite air**

**satellite orbiter**

**satellite spying**

**satellite equipment**

**satellite apparatus**

**satellite planet**

**satellite celestial body**

Parallel Tag Clouds



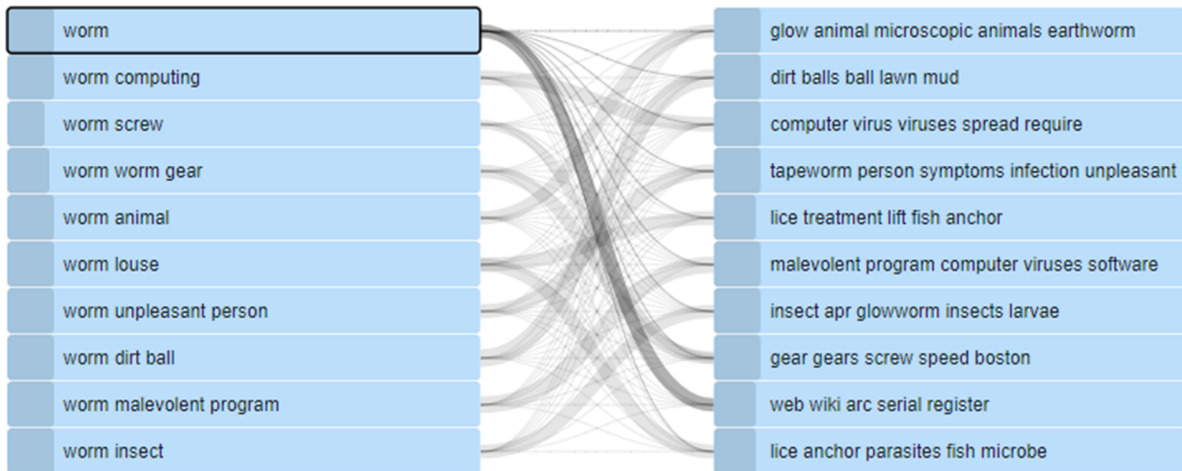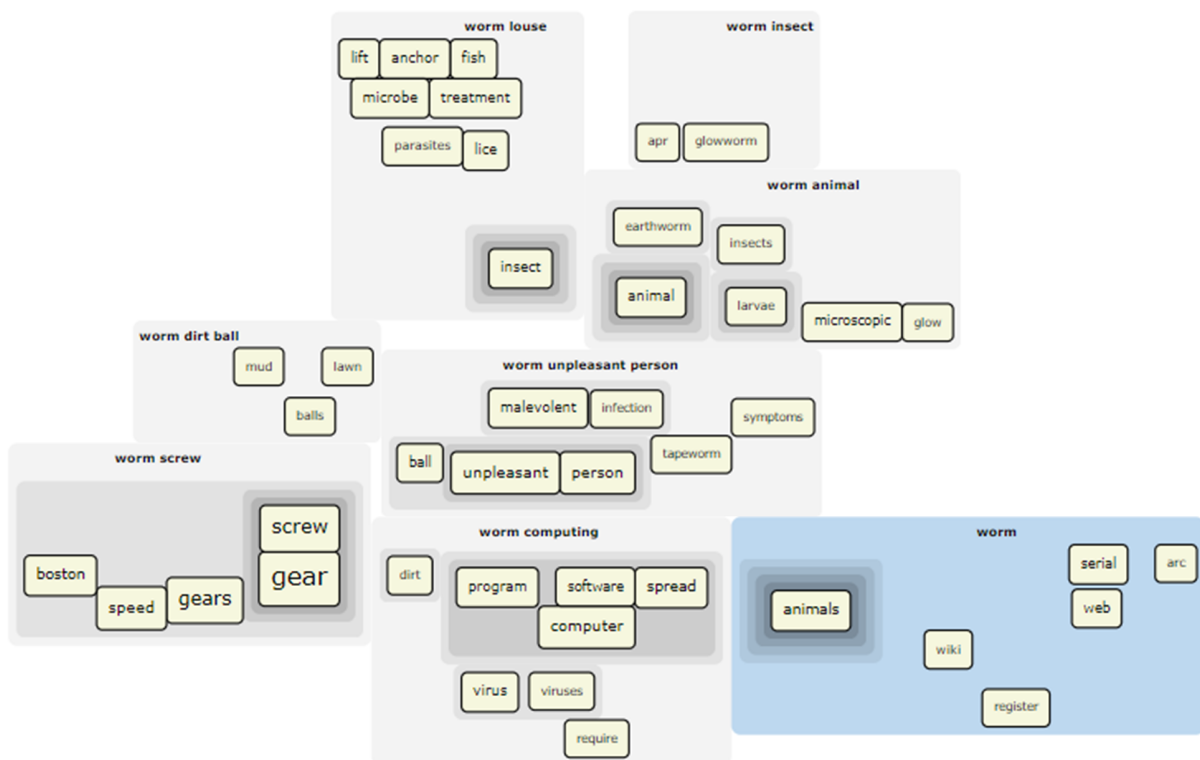| satellite | satellite artificial satellite | satellite follower | satellite air | satellite orbiter | satellite spying | satellite equipment | satellite apparatus | satellite planet | satellite celestial body |

# Worm

Parallel List

| worm | glow animal microscopic animals earthworm |
|---|---|
| worm computing | dirt balls ball lawn mud |
| worm screw | computer virus viruses spread require |
| worm worm gear | tapeworm person symptoms infection unpleasant |
| worm animal | lice treatment lift fish anchor |
| worm louse | malevolent program computer viruses software |
| worm unpleasant person | insect apr glowworm insects larvae |
| worm dirt ball | gear gears screw speed boston |
| worm malevolent program | web wiki arc serial register |
| worm insect | lice anchor parasites fish microbe |

Euler Diagram

Text

**worm computing**

**worm screw**

**worm worm gear**

**worm animal**

**worm louse**

**worm unpleasant person**

**worm dirt ball**

**worm malevolent program**

**worm insect**

Parallel Tag Clouds

# 7 PERFORMANCE TESTS

We performed a benchmark test of the three visualizations with five different queries. Currently, the natural language processing step is the biggest bottleneck in our pipeline. Also, the layout of ComED consumes considerable processing time. We therefore pre-computed all visualizations for the user study, so that the rendering step was the only limiting step.

The table below shows the average computation times (in milliseconds) for n=9 expansions of five queries (kiwi, worm, iron, grilling, and rock art) using k=10 topics and m=5 key terms per topic.

| Pipeline step | Time (ms) |
|---|---|
| Expansion term retrieval from ConceptNet | 500 |
| Retrieval of document surrogates from Google Custom Search | 1120 |
| Part-of-Speech tagging and stop word removal | 3280 |
| Topic modeling (20 - 1000 NMF iterations, convergence value 0.001) | 80-3530 |
| Computation of visualization data structures | 20 |
| ComED layout computation (background script) | 3040 |
| Rendering of PTC and List View | 20 |
| Rendering of ComED | 1700 |