

Applying Lexical Link Analysis to Discover Insights from Public Information on COVID-19

Ying Zhao^{a,1,2} and Charles C. Zhou^{b,1,2}

^aInformation Sciences Department, Naval Postgraduate School, Monterey, CA 93943, USA; ^bQuantum Intelligence, Inc., Monterey, CA 93943, USA

This manuscript was compiled on April 10, 2020

1 **SARS-Cov-2, the deadly and novel virus, which has caused a world-**
2 **wide pandemic and drastic loss of human lives and economic activ-**
3 **ities. An open data set called the COVID-19 Open Research Dataset**
4 **or CORD-19 contains large set full text scientific literature on SARS-**
5 **CoV-2. The Next Strain consists of a database of SARS-CoV-2 viral**
6 **genomes from since 12/3/2019. We applied an unique information**
7 **mining method named lexical link analysis (LLA) to answer the call**
8 **to action and help the science community answer high-priority scien-**
9 **tific questions related to SARS-CoV-2. We first text-mined the CORD-**
10 **19. We also data-mined the next strain database. Finally, we linked**
11 **two databases. The linked databases and information can be used**
12 **to discover the insights and help the research community to address**
13 **high-priority questions related to the SARS-CoV-2's genetics, tests,**
14 **and prevention.**

information mining | text mining | data mining | coronavirus | SARS-CoV-2 | COVID-19 | clade | mutation | genetics | tests | prevention

1 **T**his paper is to answer the call to action to the Nation's
2 data scientists and AI experts to mobilize and develop
3 new text and data mining techniques that can help the science
4 community answer high-priority scientific questions related to
5 SARS-Cov-2, the deadly and novel virus, which has caused
6 a worldwide pandemic and drastic loss of human lives and
7 economic activities.

8 Recently, an open data set, namely the COVID-19 Open Re-
9 search Dataset or CORD-19 (2) containing more than 40,000
10 full text scientific literature on SARS-CoV-2, was released (1, 2)
11 for researchers across technology, academia, and the govern-
12 ment. It is hopeful that machine learning and AI community
13 might employ advanced data sciences to surface unique insights
14 in the body of data and help answer with high-priority ques-
15 tions related to genetics, incubation, treatment, symptoms,
16 and prevention. The corpus is updated weekly as new research
17 is published in peer-reviewed publications and archival services
18 like bioRxiv, medRxiv, and others. In this paper, we show
19 how to apply an unique information mining method lexical
20 link analysis (LLA) to link unstructured and structured data
21 address the research questions below:

- 22 1. How to extract themes and topics that address the high-
23 priority questions of genetics, incubation, treatment,
24 symptoms, and prevention of COVID-19?
- 25 2. How to extract valued information such as information
26 with authority, insights and innovation for combating
27 SARS-Cov-2? What is the authoritative information and
28 insightful information?
- 29 3. What are the timelines of themes and topics across all
30 the research literature?

31 We show LLA that integrates text and data mining into a

single platform so that insights from data can be visualized,
and validated to answer the high-priority questions.

Materials and Methods

Overview of Lexical Link Analysis. Traditionally in social networks, the importance of a network node is a form of high-value information. For example, the leadership role in a social network (7, 8) is measured according to centrality measures (9). Among various centrality measures, a common practice is to sort and rank information based on authority. Current automated methods, such as graph-based ranking used in many search engines (10), require established hyperlinks, citation networks, social networks, or other forms of crowd-sourced collective intelligence. However, these methods are not applicable to situations where there exist no pre-established relationships among network nodes such as the data set SARS-Cov-2. This makes traditional methods difficult to apply. Furthermore, current methods mainly score popular information. High-value information can be totally different depending on specific applications. Popular and authoritative information identified by the current methods can be useful for marketing applications or crowdsourcing applications. Emerging and anomalous information is important for looking for insights and innovation. In paper, we show how to apply game-theoretic framework of lexical link analysis (LLA) to discover and rank high-value information from unstructured and structured data from SARS-Cov-2.

In LLA, a complex system can be expressed in a list of attributes or word features with specific vocabularies or lexicon terms to describe its characteristics. LLA is first a data-driven text analysis method. Fig. 1 shows an example of extracting and learning word pairs, or bi-grams as lexical terms, from text data. Words from a text document are represented as nodes, which form word pairs or bi-grams, via the links between any two nodes. For instance, the node "antiviral" in Fig. 1 is formed with "chain-terminating," "broad-spectrum," etc. as bi-gram word pairs (word features). In contrast to human-annotated word networks, such as WordNet (13), LLA automatically discovers word pairs, divides them into clusters and themes, and displays them as word networks. Fig. 1 shows an example of LLA.

Bi-gram also allows LLA to be extended to structured data (15) including meta-data such as the ones for CORD-19, where a word

Significance Statement

In this paper, we show how to apply an unique information mining method lexical link analysis (LLA) to link unstructured (CORD-19) and structured (Next Strain) data sets to relevant publications, integrate text and data mining into a single platform to discover the insights that can be visualized, and validated to answer the high-priority questions of genetics, incubation, treatment, symptoms, and prevention of COVID-19.

Y.Z. and C.C.Z. designed and performed research, and wrote the paper.

The authors declare no conflict of interest.

¹ Y.Z. and C.C.Z. contributed equally to this work.

² Email correspondence: yzhao@nps.edu or charles.zhou@quantumii.com

72 is an attribute combined with its possible values. LLA is related
 73 to but significantly different from the methods such as bag-of-
 74 words (BOW) methods, Automap (7), Latent Dirichlet Allocation
 75 (LDA) (14), Latent Semantic Analysis (LSA) (16), Probabilistic
 76 Latent Semantic Analysis (PLSA) (17) and can be jointly used with
 77 NEE (20, 21), POS methods (25). LLA is related to unsupervised
 78 learning algorithms such as k-means, principal component analysis
 79 (PCA), and spectral clustering (18). In a social network, the most
 80 connected nodes are typically considered the most important nodes.
 81 However, the uniqueness of LLA is that we consider anomalous
 82 information (word features) might be more interesting. Community
 83 detection algorithms have been illustrated by Newman (11, 12)
 84 in terms of a quality function as the “modularity” measure for a
 85 community (cluster) and optimized using a dendrogram-like greedy
 86 algorithm. The uniqueness of LLA includes new value metrics
 87 to identify high-value information that are not presented in the
 88 other methods. The new value metrics consider a game-theoretic
 89 framework (19, 35).

90 **LLA Applied to CORD-19.** LLA automatically discovers popular (P)
 91 themes, emerging (E) themes, and anomalous (A) themes (as
 92 defined below) as follows:

- 93 • Popular (P) themes: These themes resemble themes generated
 94 from the eigenvalue centrality measure in the network sciences.
 95 The themes represent the main topics in a data set. They can
 96 be insightful information in two folds: 1) These word pairs
 97 are more likely to be shared or cross-validated across multiple
 98 diversified domains, so they are considered authoritative; 2)
 99 These themes could be less interesting because they are already
 100 in public consensus and awareness and can be considered as
 101 popular.
- 102 • Emerging (E) themes: These themes tend to become popular
 103 or authoritative over time. An emerging theme has the
 104 intermediate number of inter-connected word pairs. They are
 105 emerging and important themes and can be high-value for
 106 further investigation.
- 107 • Anomalous (A) themes: These themes may not seem to belong
 108 to the data domain when compared to other themes.
 109 They are interesting outliers and can be high-value for further
 110 investigation.

111 Fig. 2 shows an example of extracted themes from a text data set
 112 where popular themes, e.g. 517(P), emerging theme, e.g. 42(E),
 113 and anomalous theme e.g. 478(A) among others are listed. Fig. 8
 114 shows a drill-down visualization for theme 517(P) labeled “infection
 115 coronavirus, global coronavirus” The labeled nodes in Fig. 8
 116 are the words with the most connections with other words (via bi-gram
 117 measures in LLA).

118 The Next Strain Database.

119 **Data.** Understanding the spread and evolution of a virus such as
 120 SARS-Cov-2 is important for effective public health measures and
 121 surveillance. A global strain tracking tool from nextstrain.org (4, 6),
 122 i.e., which strain evolves from which other strain. Nextstrain consists
 123 of a database of viral genomes, a bioinformatics pipeline for
 124 phylogenetics analysis of many novel viruses such as Zika, Ebola,
 125 and SARS-Cov-2 with an interactive visualization platform. The
 126 visualization integrates sequence data with other data types such as
 127 geographic information, serology, or host species. Nextstrain
 128 compiles the current understanding of phylogenetic analysis into a single
 129 accessible location, open to health professionals, epidemiologists,
 130 virologists and the public alike (5).

131 **Insights.** Fig. 3 shows the phylogenetic tree of the SARS-Cov-2’s
 132 similarity and development around the globe from 12/3/2019 to
 133 3/25/2020 where the 2649 cases of patients’ genomic data uploaded
 134 to the nextstrain.org website. Visually, Fig. 3 shows the strain in
 135 North America and Europe mutated to a different branch of the
 136 tree around 2/25/2020 and is more contagious. Fig. 4 shows the
 137 same data set colored by clades. Tab. 1 shows a split of the 2649
 138 cases to three time periods: before 2/25/2020, from 2/25/2020 to
 139 3/16/2020, after 3/16/2020. We computed the average case per day
 140 for Clade A2, A2a (the new mutation) and Clade A1a and other
 141 types. Clade A2, A2a is the mutated strain because it only has the

142 data starting from 1/28/2020. Clade A1a, etc. and A2, A2a are
 143 statistically different ($p < 0.0001$) measured via average cases per
 144 day for each time period, i.e., 6.29 vs 0.24 (before 2/25/2020), 38.81
 145 vs 46.86 (between 2/25/2020 and 3/16/2020), and 17.89 vs 31.33
 146 (after 3/16/2020). A2a almost started on 2/25/2020 since only 11
 147 cases of A2, A2a happened before 2/25/2020. Clade A2, A2a is more
 148 contagious than Clade A1a, and other since the average cases per day
 149 is much higher, 46.86 and 31.33 between 2/25/2020 and 3/16/2020.
 150 After 3/17/2020, many countries in North America and Europe
 151 implemented “shelter-in-place,” the average cases per day decreased
 152 for both clades, however, the mutated clade decreased statistically
 153 significantly slower (31.33) than A1a, B1, B2, etc. together (17.89).

154 Fig. 5 shows how the first case of Clade A2a linked to earlier
 155 cases of Clade A2. There are two insights as follows:

- 156 • Among the total 1277 cases of Clade A2, A2a, there were only
 157 34 cases in Asia, including four cases in China and 18 cases
 158 in Japan. There were 955 cases in Europe, 216 cases in North

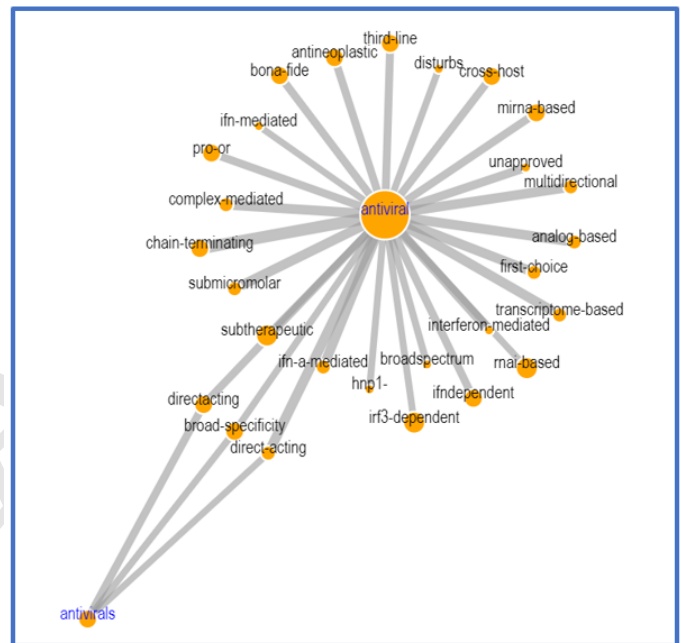


Fig. 1. LLA Example

ID	Theme Keywords
423(P)	in, in, in, in, in
317(P)	infection coronavirus, global, coronavirus
42(E)	antibody limited, antibody, antibodies
374(E)	patients, group, group, canine, patients, canine
233(E)	complete, support, main
423(E)	protein, mutant, gene, mutant, gene, protein, protein, gene
126(E)	data, functional, data, primer
50(E)	assay, rapid, database
478(A)	distinct, membranes, membranes, distinct, structures
114(A)	peptide, inhibitor, peptide, prevent, inhibitor, peptide
451(A)	sequence, ma, ma, final, sequence, final
480(A)	national, increasing, vivo
318(A)	unique, similarly, nm
152(A)	vector, recombinant, plasmid, vector, vector, plasmid, recombinant, vector, plasmid, recombinant
100(A)	usa, university, university, usa, short
464(A)	detect, signaling, receptors
123(A)	chain, medicine, media
204(A)	reduce, weight, staining
515(A)	expressing, strain, kit
471(A)	program, software, numerous

Fig. 2. LLA themes

Table 1. Cases

A1a, B1, B2, etc	Cases	Days	Cases/Per Day
12/24/2019-2/24/2020	396	63	6.29
2/25/2020-3/16/2020	815	21	38.81
3/17/2020-3/25/2020	161	9	17.89
A2, A2a			
1/28/2020-2/24/2020	11	46	0.24
2/25/2020-3/16/2020	984	21	46.86
3/17/2020-3/25/2020	282	9	31.33

word pairs. These matched word pairs (concepts) are also grouped based on these themes. Fig. 8 shows a popular theme of word pair appeared in both data sources. Fig. 9, Fig. 10, and Fig. 11 show examples of emerging themes appeared in both data sources. Fig. 12 shows an example of an anomalous theme appeared in both data sources. Emerging themes are interesting topics for researchers to drill down and discover information in CORD-19 pertinent to high-priority questions of the SARS-Cov-2 of genetics (Group 423(E)), tests (Group 50(E)), and prevention (Group 42(E)).

Conclusion and Challenge to Future. We applied a unique information mining method lexical link analysis to conduct a preliminary study to the call to action and help the science

159 America, and 15 cases in Australia.

- 160 • Among the four cases of Clade A2, A2a in China, three were
- 161 collected from 1/28/2020 to 2/6/2020, they were submitted
- 162 around 3/20/2020. The fourth case was collected on 3/23/2020
- 163 and submitted at the same time.

164 At least the data shows the Clade A2, A2a in Europe and North
165 America are more contagious and virulent.

166 **Linking the Next Strain Data Set to the CORD-19 Data Set.** The next
167 strain data can be linked to the CORD-19 data set using LLA. We
168 first extracted another document data set with documents pertinent
169 to the next strain database. We then fused the two data sources
170 based on themes in Fig. 2.

171 Fig. 6 shows a match matrix of number of matched word pairs
172 from two document sources (the next strain and CORD-19), there
173 are 2012 matched word pairs. Fig. 7 shows the examples of matched

	Match Score	cord_19	nextstrain	Uniqueness Score
1	cord_19	2102.00	2102.00	777591.00
2	nextstrain	2102.00	2102.00	6490.00

Fig. 6. LLA match matrix

- [517]diseases infectious
- [517]tract respiratory
- [517]respiratory acute
- [517]disease infectious
- [517]viruses respiratory
- [517]syndrome respiratory
- [517]virus influenza
- [517]pcr real-time
- [517]health global
- [517]viruses influenza
- [517]control disease
- [517]size sample
- [517]infections respiratory
- [517]symptoms respiratory
- [517]rate mortality
- [517]disease respiratory
- [517]severity disease
- [517]products pcr
- [517]failure respiratory
- [517]pathogens respiratory
- [517]infections tract
- [517]infection tract
- [517]illness respiratory
- [517]influenza seasonal

Fig. 7. LLA matched list

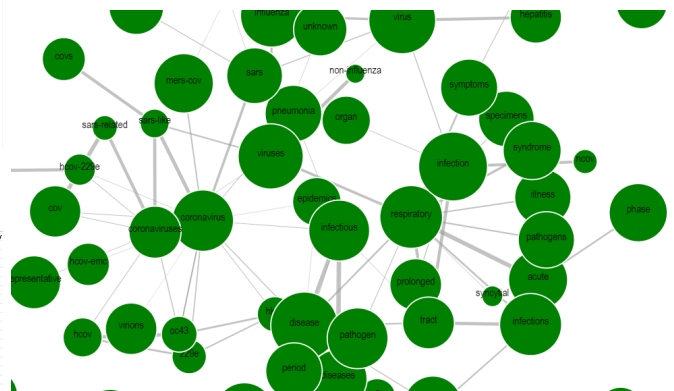


Fig. 8. LLA group 517(P)



Fig. 3. Next strain data by country from 12/3/2019 to 3/25/2020. Source: nextstrain.org

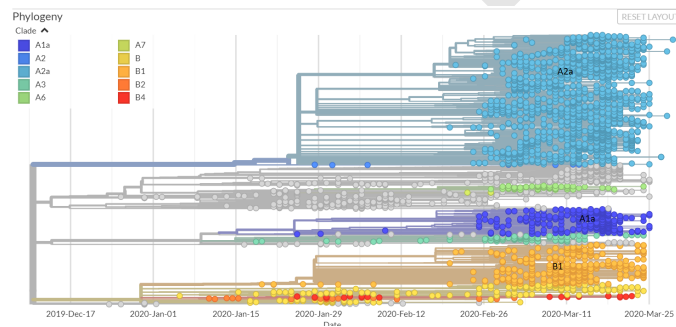


Fig. 4. Next strain data by clade from 12/3/2019 to 3/25/2020. Source: nextstrain.org

Strain	Admin Div	geoid	url	Country	Age	Sex	Submission Date	Originating Region	Clade	Collection Date	Author	genbank_access	Location	Exposure History
Germany/Berlin/1/2020	Bavaria	EPY_IS_405882		Germany	Older	Male	Charter 10h ago	Europe	A2	1/28/2020	Corman et al		Starnberg	
Shanghai/SH08/2020	Shanghai	EPY_IS_415312		China	47 female	One week ago	Shanghai/Asia	Asia	A2	1/28/2020	Wang et al			
Shanghai/SH08/2020	Shanghai	EPY_IS_415388		China	38 Male	One week ago	Shanghai/Asia	Asia	A2	1/31/2020	Wang et al			
Shanghai/SH02/2020	Shanghai	EPY_IS_415318		China	81 Male	One week ago	Shanghai/Asia	Asia	A2	2/6/2020	Wang et al			
Italy/CLC1/2020	Lombardy	EPY_IS_412973		Italy	38 Male	Older	Department	Europe	A2a	2/10/2020	Verfaelliet al			
France/PIE4/2020	Hauts de France	EPY_IS_418238		France	36 Male	3-7 days ago	Centre Hoop	Europe	A2a	2/12/2020	Albert et al		Compiègne	
Italy/UM1/2020	Lombardy	EPY_IS_417445		Italy	76 Male	3-7 days ago	Laboratory	Europe	A2a	2/24/2020	Seidenhofer et al		Milan	
Netherlands/Breda/UM_136364/2020	North Brabant	EPY_IS_413568		Netherlands	71 Male	One month ago	Foundation	Europe	A2a	2/24/2020	Heesterbeek et al		Breda	
Switzerland/17048/2020	Ticino	EPY_IS_413998		Switzerland	70 Male	One month ago	Laboratory	Europe	A2a	2/24/2020	LAURISCHER Florian et al			
Italy/UM1/2020	Lombardy	EPY_IS_417445		Italy	80 Male	3-7 days ago	Laboratory	Europe	A2a	2/24/2020	Seidenhofer et al			
Finland/IN_2/2020	Helsinki	EPY_IS_412975		Finland	24 Female	Older	HUS Diagon	Europe	A2a	2/25/2020	Simonsen et al		Milan	Italy
Germany/Baden-Wuerttemberg_1/2020	Baden-Wuerttemberg	EPY_IS_412952		Germany	Older	State Health	Europe	A2a	2/25/2020	Corman et al			Italy	
Spain/Madrid/15/2020	Madrid	EPY_IS_418235		Spain	24 Male	3-7 days ago	Hospital	South America	A2a	2/25/2020	Alfonso-Caballeros et al			
Brazil/SBR_01/2020	Sao Paulo	EPY_IS_412964		Brazil	61 Male	Older	Hospital	South America	A2a	2/25/2020	Isaqueire Gomes de Jesus et al		Sao Paulo	Italy
Denmark/COV-02/2020	Copenhagen	EPY_IS_415424		Denmark	61 Male	One week ago	Department	Europe	A2a	2/26/2020	Henningsen et al			
Switzerland/GE3895/2020	Geneva	EPY_IS_413997		Switzerland	28 Male	One month ago	Laboratory	Europe	A2a	2/26/2020	LAURISCHER Florian et al			Italy

Fig. 5. The first case of Clade A2a. Source: nextstrain.org

186 community answer high-priority scientific questions related to
 187 SARS-Cov-2. We first text-mined an unstructured database,
 188 i.e. the COVID-19 Open Research Dataset or CORD-19. We
 189 also data-mined a structured database, i.e., the next strain
 190 database, covering the period from 12/3/2019 to 3/25/2020.
 191 Finally, we linked two databases and certain publications,
 192 and discovered the insights and methodologies. The linked docu-
 193 ments from CORD-19 can help address high-priority questions
 194 related to SARS-COV-2's genetics, tests, and prevention.

195 There are some un-answered questions we need to ponder:

- 196 1. If certain data and publications before 12/3/2019 need
 197 to be mined and analysed, can one estimate the outcome of
 198 COVID-19 outbreak after 12/3/2019?
 199 2. Can we use the data and publications up to today to
 200 estimate what will happen one month later?

201 1. COVID-19 Open Research Dataset (CORD-19). 2020. Version 2020-03-13. Retrieved from
 202 <https://pages.semanticscholar.org/coronavirus-research>. doi:10.5281/zenodo.3715506

2. <https://techcrunch.com/2020/03/16/coronavirus-machine-learning-cord-19-chan-zuckerberg-ostp> 203
 3. <https://www.whitehouse.gov/briefings-statements/call-action-tech-community-new-machine-readable-covid-19-dataset/> 204
 4. Nextstrain <https://nextstrain.org/ncov> 205
 5. github.com/nextstrain 206
 6. Hadfield et al. (2018), Nextstrain: real-time tracking of pathogen evolution, *Bioinformatics* (2018) 207
 7. Center for Computational Analysis of Social and Organizational Systems (CASOS). (2009) AutoMap: extract, analyze and represent relational data from texts. Retrieved from <http://www.casos.cs.cmu.edu>. 208
 8. Girvan, M., and Newman, M. E. J. (2002). Community structure in social and biological networks. *Proceedings of the National Academy of Sciences, USA*, 99(12), 7821–7826. 209
 9. Freeman, L.C. (1979). Centrality in social networks I: conceptual clarification. *Social Networks*, 1: 215-239. 210
 10. Brin, S. and Page, L. (1998). The Anatomy of a large-scale hypertextual web search Engine. *Computer Networks and ISDN Systems*, 30:107-117. 211-218

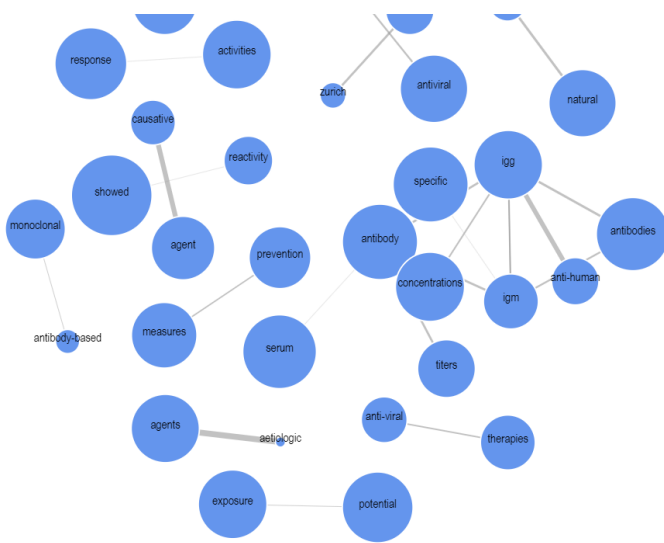


Fig. 9. LLA group 42(E)

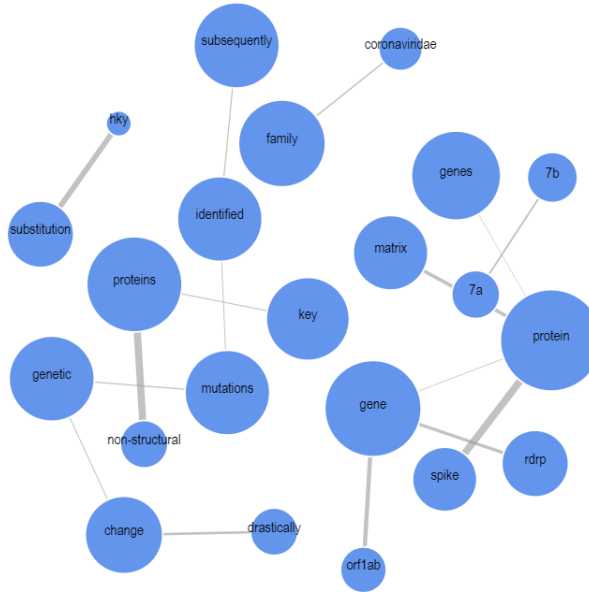


Fig. 11. LLA group 423(E)

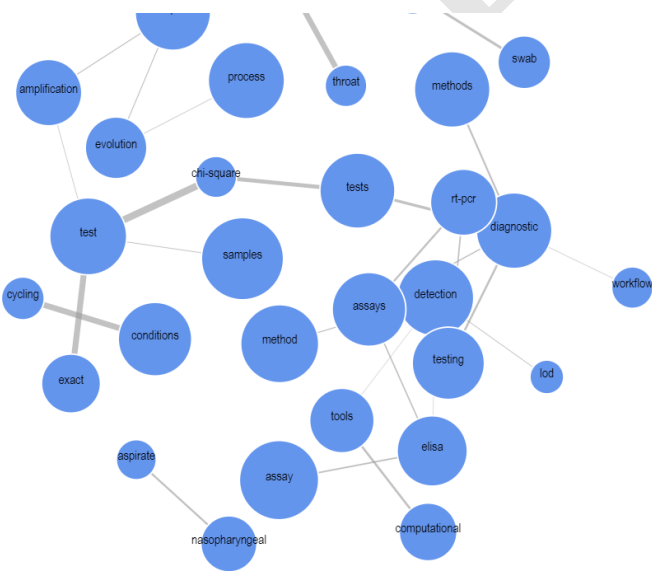


Fig. 10. LLA group 50(E)

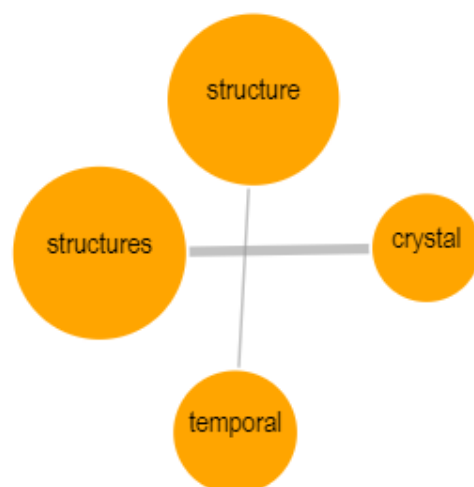


Fig. 12. LLA group 478(A)

- 220 11. Newman, M. E. J. (2003). "Fast algorithm for detecting community structure in networks,"
221 2003. Retrieved from <http://arxiv.org/pdf/cond-mat/0309508.pdf>
- 222 12. Newman, M. E. J. (2006). Finding community structure in networks using the eigenvectors of
223 matrices. *Phys. Rev. E*, vol. 74, no. 3.
- 224 13. Miller, G. A. (1995). WordNet: a lexical database for English. *Communications of the ACM*,
225 38(11).
- 226 14. Blei, D., Ng, A. and Jordan, M. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993-1022,
227 2003. Retrieved from
228 <http://jmlr.csail.mit.edu/papers/volume3/blei03a/blei03a.pdf>.
- 229 15. US patent 8,903,756. (2014). System and method for knowledge pattern search from net-
230 worked agents. 2014. Retrieved from <https://www.google.com/patents/US8903756>
- 231 16. Dumais, S. T., Furnas, G. W., Landauer, T. K. and Deerwester, S. (1988). Using latent seman-
232 tic analysis to improve information retrieval. In *Proceedings of CHI'88: Conference on Human*
233 *Factors in Computing*, 281-285.
- 234 17. Hofmann, T. (1999). Probabilistic latent semantic analysis," In *Proceedings of the Fifteenth*
235 *Conference on Uncertainty in Artificial Intelligence*, Stockholm, Sweden.
- 236 18. Ng, A., Jordan, M., and Weiss, Y. (2002). On spectral clustering: analysis and an al-
237 gorithm. In T. Dietterich, S. Becker, and Z. Ghahramani (Eds.), *Advances in Neural In-*
238 *formation Processing Systems 14* (pp. 849 – 856), (2002). MIT Press. Retrieved from
239 <http://ai.stanford.edu/~ang/papers/nips01-spectral.pdf>
- 240 19. Zhao, Y., Zhou, C., and Huang, S. (2019). Theory and use case of game-theoretic lexical
241 link analysis. In the proceedings of the IEEE/ACM International Conference on Advances in
242 Social Networks Analysis and Mining (ASONAM).
- 243 20. InXight 1997. Retrieved from <http://en.wikipedia.org/wiki/Inxight>
- 244 21. MUC-7: PAPERS: SYSTEMS - Named Entity Tasks Retrieved from [http://www-](http://www-nlpir.nist.gov/related_projects/muc/proceedings/muc_7_toc.html#named)
245 [nlpir.nist.gov/related_projects/muc/proceedings/muc_7_toc.html#named](http://www-nlpir.nist.gov/related_projects/muc/proceedings/muc_7_toc.html#named)
- 246 22. Yatsko, V. A. and Vishnyakov, T. N. 2007. A method for evaluating modern systems of auto-
247 matic text summarization. *Automatic Documentation and Mathematical Linguistics*, 41(3):93-
248 103.
- 249 23. Soergel, D. 1985. *Organizing information: Principles of data base and retrieval systems*. Or-
250 lando, FL: Academic Press.
- 251 24. Turney, P. 2002. Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsuper-
252 vised Classification of Reviews. In *Proceedings of the Association for Computational Linguis-*
253 *tics*, pp. 417–424.
- 254 25. Kristina Toutanova and Christopher Manning. 2000. Enriching the knowledge sources used
255 in a maximum entropy part-of-speech tagger. In *EMNLP/VLC 1999*, pages 63–71.
- 256 26. C.J. van Rijsbergen, S.E. Robertson and M.F. Porter, 1980. *New models in probabilistic in-*
257 *formation retrieval*. London: British Library. (British Library Research and Development Report,
258 no. 5587).
- 259 27. Salton G and Buckley C (1988) Term-weighting approaches in automatic text retrieval. *In-*
260 *formation Processing & Management* 24 (5): 513-523. Self-Organization, SO, wiki, retrieved
261 from <http://en.wikipedia.org/wiki/Self-organization>
- 262 28. Jiampoamarn S and Cercone N (2005) Biological named entity recognition using n-grams
263 and classification methods. In *Proceedings of PACLING*
- 264 29. Bekkerman R and Allan J (2003) Using Bigrams in Text Categorization. Retrieved from
265 <http://people.cs.umass.edu/~ronb/papers/bigrams.pdf>
- 266 30. N-gram, <http://en.wikipedia.org/wiki/N-gram>
- 267 31. sciSpaCy <https://allenai.github.io/scispacy/>
- 268 32. sciBert <https://github.com/allenai/scibert>
- 269 33. Beltagy, I., Lo, K., Cohan, A. (2019). SciBERT: A Pretrained Language Model for Scientific
270 Text. <https://arxiv.org/abs/1903.10676>
- 271 34. Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. (2019). BERT: Pre-
272 training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.
- 273 35. Zhao, Y., Zhou C. (2018). A Game-Theoretic Lexical Link Analysis for Discover-
274 ing High-Value Information from Big Data. In the proceedings the 2018 IEEE/ACM
275 International Conference on Advances in Social Networks Analysis and Mining,
276 Barcelona, Spain, 28-31 Aug. 2018 (ASONAM 2018), page 621 – 625. Retrieved from
277 <https://ieeexplore.ieee.org/document/8508317>.