

Article

Covid19-related Scientific Literature Exploration: Short Survey and Comparative Study

Bahaj Adil ^{1,‡}, Safae Lhazmir ^{1,‡}, Mounir Ghogho ^{1,2}, Houda Benbrahim ³

¹ International University of Rabat, TicLAB, Sala el Jadida, Morocco

² University of Leeds, Faculty of Engineering, Leeds, UK

³ Mohamed V University, ENSIAS, Rabat, Morocco

‡ These authors contributed equally to this work.

Abstract: The urgency of the COVID19 pandemic caused a surge in related scientific literature. This surge made the manual exploration of scientific articles time-consuming and inefficient. Therefore, a range of exploratory search applications have been created to facilitate access to the available literature. In this survey, we give a short description of certain efforts in this direction and explore the different approaches that they used.

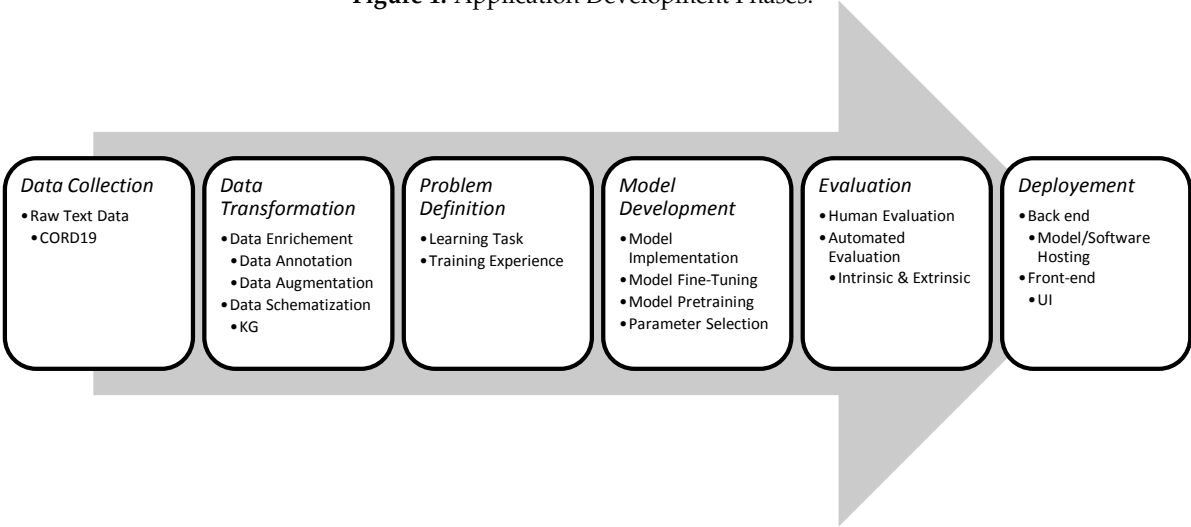
Keywords: COVID-19; Exploratory Search; Machine Learning; Document Retrieval

1. Introduction

Due to the vast expansion of COVID19 literature (more than 20,000 unique papers were published between 1 January and 30 June 2020 [1,2]) there was a need to create information management and retrieval systems for COVID literature. The data science community responded to this urgent need by creating and deploying dozens of applications to provide researchers with easy access to COVID19 literature. These applications mainly focus on text mining [3] and its related tasks (e.g., document retrieval [4], question answering [5], passage retrieval [6], summarization [7]...etc) in order to organize and access relevant knowledge effortlessly. Several public competitions and common tasks, such as the CORD-19 and TREC-COVID initiatives [8,9], further encouraged such efforts.

In this work, we explore COVID literature exploration applications, which we can classify as one of two categories relative to the format of the search results (a) textual search engines, and (b) visual search engines. The first category comprises query-oriented applications, that extract information from COVID19 literature using queries. The second class of applications is used mainly for the bibliometric study of COVID19 literature coupled with visual interactive or static summarization graphs. Each one of these applications goes through the same development phases. Figure 1 shows the most common phases that an application would go through. First, the text data needed by the system must be collected. All the explored applications in this work use the CORD19 [8] dataset (either a version of it or a subset of a version of CORD12). Second, raw data collected may need to be transformed in some cases to meet certain specifications. This can be done by enriching the data in order to make it more representative, or it can be achieved by structuring the available data differently. Third, given the available data and the basic application specifications, a set of learning problems (i.e. question answering, document retrieval, passage retrieval...) need to be defined. Forth, given the defined learning problem, machine learning models are developed and trained to achieve the learning tasks. Fifth, the models are evaluated, either by a human or an automated evaluation process. Sixth, after evaluating the models, they need to be deployed to ensure their accessibility by a larger number of users, and that is by providing an easy-to-use user interface with a reliable model execution backend architecture.

Figure 1. Application Development Phases.



34 Although a previous survey [10] has explored COVID literature search engines, their work has
35 certain limitations that we try to remedy in this work. First, rather than focusing primarily on textual
36 search engines, we explore visual search engines. Second, [10] included a plethora of applications that
37 are not associated with any research papers or technical reports, consequently, we discarded these
38 applications and focused on applications with research papers in order to gain and express a deeper
39 understanding of the methods that they employed. Third, we try to infer some design principles that
40 the authors of the works used to create their system.

41 This work is organized as follows: in the **datasets** section, we describe some datasets that were
42 used in the explored works for various purposes. In the **exploratory search applications** section, we
43 explore the characteristics and design principles of COVID19 exploratory literature search applications.
44 In the **evaluation methods** section, we explore certain methods that were used to evaluate the systems.
45 The **limitations and future improvements** section shows certain limitations of the examined works.
46 The **conclusion** section concludes our work, and the following section gives certain limitations that
47 this work has.

48 **2. Datasets**

49 In this section, we list some of the datasets that were used in the works that we explored. We
50 categorized the datasets relative to their structure into three categories: a) unstructured, b) structured
51 and c) hybrid. Figure 2 represents a taxonomy of the employed datasets.

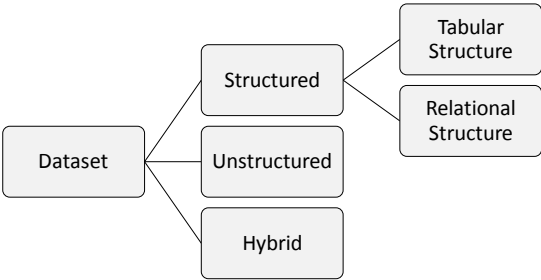


Figure 2. Taxonomy of Data Types.

2.1. Unstructured Datasets

Unstructured data is information that does not have a defined data model. This type of data is mainly textual in nature. The following structured and hybrid datasets have been built using unstructured data. In fact, all the previously mentioned categories were either automatically or manually curated and annotated from different literature databases (e.g., Arxiv, DBLP, Pubmed, bioRxiv, medRxiv), which contain unstructured documents, often in a hard to read format such as PDFs.

2.2. Structured Datasets

We can recognize two kinds of structured data: a) data with tabular structure, where every example shares the same set of variables, and examples are independent of each other and b) data with relational structure, where examples do not necessarily share the same set of variables, examples are inherently typed, that is, each example belong to a predefined group of examples, and examples have a dependency between them, which is implemented practically in the form of links.

The first category contains mainly annotated datasets that are oriented for machine learning purposes, such as training, fine-tuning, or evaluating the created models on specific tasks. The works that we explored use multiple datasets. A later section defines some of the main tasks that the works try to solve. All of these tasks are text-oriented and can fall under the umbrella of information retrieval in general. Annotated datasets such as TREC-COVID [9] and BioASQ [11] were used for document retrieval. These datasets are generally constructed by a set of human curators who were provided with a list of queries (or questions) and a set of supposedly relevant documents, and the goal was to select the most pertinent documents for each query. In addition, multiple datasets have been used to train question answering models such as CovidQA [12], COVID-19 Questions [13], Covid-QA [14], InfoBot Dataset [15], MS-MARCO [16], Med-MARCO [17], Natural Questions [18], SQuAD [19], BioASQ [11], M-CID [20] and QuAC [21]. Other datasets were used to train document summarization models. For example, DUC 2005 [22], 2006 [23], and Debatepedia [24] were used by [25] to train document summarization models. Other datasets, such as GENIA [26], JNLPBA [27], CHEMDNER [28], NCBI Disease Corpus [29], CHEMPROT [30] BC5CDR [31], COV19_729 [32], were used for the named entity recognition (NER) of multiple types of entities, namely, chemicals, genes, proteins, diseases, and other biomedical entities. Relation extraction (RE) was also a task of interest in [32], which was achieved using the CHEMPROT [30] BC5CDR [31] datasets. NER and RE tasks are generally used in knowledge graph construction, where the entities extracted represent nodes and the relations represent edges between nodes. Some of these datasets were curated using data from COVID19 related source documents, e.g., CovidQA [12], COVID-19 Questions [13], Covid-QA [14], InfoBot Dataset [15], TREC-COVID [9].

Concerning data with relational structure, some works used knowledge graphs constructed from COVID19 related literature. In general, the graphs contain four types of entities with multiple properties: 1) a paper entity, which represents a research paper and can be described by a Digital Object Identifier (DOI), title, publication date, and other properties; 2) an author entity, which represents a publication's author, and can be describe by an identifier, a first, middle and last name and other properties of interest; 3) an affiliation entity, which represents the research structure (lab, university, company, etc) to which the author is affiliated, which can be described by an identifier, a name and other properties of interest; 4) a concept entity, which represents a domain knowledge-related notion that exists in a paper. A concept can be represented by one word or a series of words. Concepts can have multiple types of relationships between them depending on the type of concepts. For example, concepts of biomedical types, such as genes, diseases, chemicals, organisms, and proteins can be linked by semantic biomedical relations [32–35], or by syntactic relationships based on their co-occurrence in the same sentence [36]. Figure 3 represents the mentioned entities and how they relate to each other. Tables 1 and 2 offer a more detailed description and these entities and how they are related. It is worth

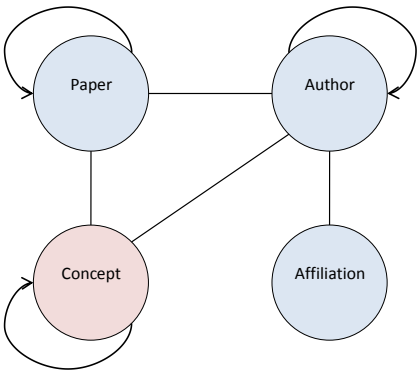


Figure 3. COVID19 Literature Knowledge Graphs Schema.

100 pointing out that not all knowledge graphs respect this schema. Some implement it totally (e.g. CKG
101 [37]), and some implement it partially (e.g. CovEx KG [38]), as shown in tables 1 and 2.

102 Furthermore, it has been observed that the design of certain knowledge graphs are dependent
103 on the tasks they are used for. For instance, for the task of document retrieval, a knowledge graph
104 is generally designed with documents as the central nodes to which other nodes may be linked
105 [37,38]. On the other hand, for the task of question answering, even though the same base data is
106 used, no node holds the document data; instead, documents are ignored and only concept nodes are
107 presented and interlinked [35]. In addition, the granularity of the relationships and the entities is
108 also important, as it was demonstrated in [34,35], where two types of relationships and entities were
109 extracted: (a) coarse-grained and (b) fine-grained. The latter was needed in a question answering
110 task to accommodate the specificity of the entities expressed in user queries, which is not required in
111 other tasks as shown in [32] for the task of link prediction, where the authors discarded fine-grained
112 relationships in favor of more general ones to reduce noise that can hinder the performance of certain
113 models. In the case of network visualization [36] adopted a more flexible approach to KG construction,
114 by extracting a set of entities and saving them, so that they could be later aggregated to create
115 domain-specific networks, which can be visualized. Some tasks such as information extension, which
116 aims at enriching certain information constructs like queries or KGs, do not need directed edges,
117 which is the case for example in Vapur KG [33] and Citation KG [39]. In fact, having undirected edges
118 help explore more complex and unexpected relationships among entities, which was illustrated in a
119 fact-checking application in [40].

Table 1. Examples of Entities Specifications.

Entities	Properties	Description	ID
Paper	title, publication date, journal, Digital Object Identifier (DOI), link	Representation of research paper entities.	E1
Author	identifier, first names, middle names, last names	Representation of the paper authors.	E2
Affiliation	identifier, name, country, city	Representation of a research structure where an author belongs.	E3
Concept	concept identifier, textual value, concept type (gene, disease, topic, chemical ...etc)	Representation of a domain specific concept.	E4

Table 3. Summary of Knowledge Graphs Related to COVID19. .

KG	Usage	Ent.	Rel.
CKG [37]	article recommendations, citation-based navigation, and search result ranking.	E1, E2, E3, E4	R1, R4, R6, R5
CovEx KG [38]	Document Retrieval.	E1, E2, E4	R1, R4, R5, R7
ERLKG [32]	Link prediction.	E4	R3
COVID-KG [35] (aka Blender-KG [41])	QA, Semantic Visualization, Drug Re-purposing	E4	R3
COFIE KG [34]	KG search over relations and entities using a query.	E4	R3
Network Visualization KG [36]	Data Visualization	E4	R3
Vapur KG [33]	Query extension.	E4	R3
Citation KG [39]	Document Ranking.	E1	R1

Table 2. Examples of Relations.

Source Entity	Dest. Entity	Relation	Description	ID
Paper	Paper	cites	This relation connects paper entities with paper references indicating a citation relation.	R1
Author	Author	Co-author	This relation connects an author entity with another author entity indicating a co-authorship relation.	R2
Concept	Concept	relate concepts	This relationship links two concepts with any general relationship that might link them.	R3
Paper	Author	authored by	This relation connects paper entities with author entities and indicates an authorship relation.	R4
Paper	Concept	associated concept	This relation connects paper entities with concept entities.	R5
Author	Affiliation	affiliated with	This relation connects author entities with institution entities.	R6
Author	Concept	research area	This relation connects author entities with concept entities indicating a research area of the author.	R7

2.3. Hybrid

Hybrid datasets have some structure, which can be in the form of tags, but most if not all of the tagged elements have no structure, which generally means that these elements are in a textual format. An example of such datasets is CORD19. The CORD-19 dataset is the centerpiece of COVID19 literature exploration applications. The CORD-19 dataset [8] is a curated set of articles from multiple resources, that were collected to help efforts against the COVID19 pandemic. This dataset was used in a common document retrieval task TREC-COVID, where a set of CORD-19 articles were curated and annotated for their relevance relative to certain user queries. The dataset is ever-expanding, with new articles being added to it intermittently. The dataset is available online at <https://www.kaggle.com/allen-institute-for-ai/CORD-19-research-challenge>.

3. Exploratory Search Applications

3.1. Textual Exploratory Search

Research related to COVID knowledge management and information retrieval (KM&IR) has gained tremendous attention over the past year. We try here to present a concise summary of the

Table 4. Summary of the Datasets. NER refers to Named Entity Recognition, RE refers to Relationship Extraction, SMZ refers to summarization, QA refers to Question Answering, DR refers to Document Retrieval.

Dataset	Application Refs	Tasks	Statistics	URL
TREC-COVID* [9]	[42] [43] [44] [17]	DR	The TREC-COVID dataset has many versions which correspond to TREC-COVID challenges. For example, round three contains a total of 16677 unique journal articles in CORD-19 with their received a relevance annotation.	https://www.kaggle.com/c/trec-covid-information-retrieval/data
CovidQA* [12]	[43] [25] [45]	QA	The dataset contains 147 question-article-answer triples with 27 unique questions and 104 unique articles.	https://github.com/castorini/pygaggle/tree/master/data
COVID-19 Questions* [13]	[13] [45]	QA	The dataset contains 111 question-answer pairs with 53 interrogative and 58 keyword-style queries.	https://drive.google.com/file/d/1z7jW0fovgTfTScCanZvrvtUax1HAMEFV/view?usp=sharing
Covid-QA* [14]	[44] [45]	QA	The dataset consists of 2019 question-article-answer triples.	https://github.com/deepset-ai/COVID-QA
InfoBot Dataset* [15]	[46]	QA, FAQ	2,200 COVID-19 related Frequently asked Question-Answer pairs.	https://covid-19-infobot.org/data/
MS-MARCO [16]	[17]	QA	1,000,000 training instances.	https://microsoft.github.io/msmarco/
Med-MARCO [17]	[17]	QA	79K of the original MS-MARCO questions (9.7%).	https://github.com/Georgetown-IR-Lab/covid-neural-ir/blob/master/med-msmarco-train.txt
Natural Questions [18]	[13]	QA	The public release consists of 307,373 training examples with single annotations; 7,830 examples with 5-way annotations for development data; and a further 7,842 examples with 5-way annotated sequestered as test data.	https://ai.google.com/research/NaturalQuestions/
SQuAD [19]	[13]	QA	The dataset contains 107785 question-answer pairs on 536 articles.	https://rajpurkar.github.io/SQuAD-explorer/
BioASQ [11]	[13]	QA, DR	500 questions with their relevant documents, text span answers and perfect answers.	http://www.bioasq.org/news/golden-datasets-2nd-edition-bioasq-challenge-are-now-available
M-CID [20]	[20]	QA	The dataset is composed out of 6871 natural language utterances across 16 COVID-19 specific intents and 4 languages: English, Spanish, French and German.	https://fb.me/covid_mcid_dataset
QuAC [21]	[44]	QA	14K information-seeking QA dialogs, and 100K questions in total.	http://quac.ai/
GENIA [26]	[38]	NER	2000 abstracts taken from MEDLINE database, and contains more than 400000 words and almost 100000 annotations.	http://www.geniaproject.org/genia-corpus/term-corpus
DUC 2005, 2006 [22,23]	[25]	SMZ	the dataset is composed out of 50 topics.	https://www-nlpir.nist.gov/projects/duc/data.html

Table 4. Summary of the Datasets. NER refers to Named Entity Recognition, RE refers to Relationship Extraction, SMZ refers to summarization, QA refers to Question Answering, DR refers to Document Retrieval (Continued).

Dataset	Application Refs	Tasks	Statistics	URL
Debatepedia [24]	[25]	SMZ	It consists of 10,859 training examples, 1,357 testing and 1,357 validation samples. The average number of words in summary, documents and query is 11.16, 66.4, and 10 respectively.	https://github.com/PrekshaNema25/DiversityBasedAttentionMechanism
JNLPBA [27]	[32]	NER	This dataset contains a subset of the GENIA dataset V3.02. This subset is composed out of 2404 abstracts. The articles were chosen to contain the MeSH terms 'human', 'blood cells' and 'transcription factors', and their publication year ranges over 1990 to 1999.	http://www.geniaproject.org/shared-tasks/bionlp-jnlpba-shared-task-2004
CHEMDNER [28]	[32]	NER	10000 PubMed abstracts that contain a total of 84,355 chemical entities.	https://biocreative.bioinformatics.udel.edu/resources/biocreative-iv/chemdner-corpus/
NCBI Disease Corpus [29]	[32]	NER	793 PubMed abstracts that were annotated. a total of 6892 disease mentions, which are mapped to 790 unique disease concepts were extracted.	https://github.com/spysalo/ncbi-disease
CHEMPROT [30]	[32]	NER, RE	2500 PubMed abstracts, from which 32000 chemical entities and 31000 protein entities were extracted. In addition , 10000 chemical-protein relationships were extracted.	http://www.biocreative.org/accounts/login/?next=/resources/corpora/chemprot-corpus-biocreative-vi/
BC5CDR [31]	[32]	NER, RE	1500 PubMed articles with 4409 annotated chemicals, 5818 diseases and 3116 chemical-disease interactions.	https://github.com/shreyashub/BioFLAIR/tree/master/data/ner
COV19_729* [32]	[32]	NER	The dataset is composed out of dataset 729 example. Each example is a triple comprising an entity, the class that that entity belongs to (i.e. disease, protein, chemical), and a physicians rating of how related those entities are to COVID.	https://github.com/sayantانبасу05/ERKLG

research in this area. The development of search engines goes through certain common steps that are illustrated in figure 4. A search engine's development process begins with the base data or the data that is relevant to the search query. Second, the raw textual data is processed to extract certain elements that are of interest, and transform them. That same raw data can be reorganized in the form of a knowledge graph to satisfy certain specifications such as fast question answering. Afterwards, the tasks that are intended for the search engine should be defined and implemented, followed by an assessment of the efficiency of the system in performing those tasks. Finally, the implemented system needs to be deployed for public access.

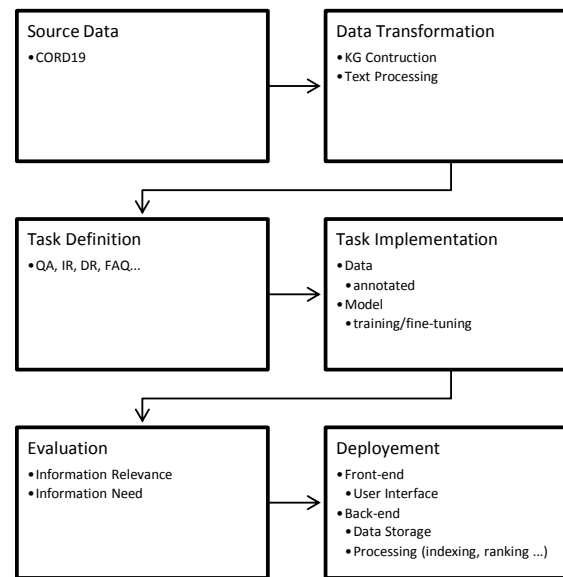


Figure 4. The Development Process of Literature Search Engines.

COVID literature knowledge management and information retrieval systems have multiple axes along which we can study, survey and compare them. We list some of these characteristics in what follows:

- **Tasks:** The tasks are related to textual data, and hence we suppose that we have a text database (or collection or corpus) \mathcal{T} as a string of N symbols drawn from an alphabet (i.e. all possible combinations of letters) Σ . A vocabulary V is the set of unique words used in \mathcal{T} . \mathcal{T} is partitioned into n documents $\{d_1, d_2, \dots, d_n\}$. A document d can be presented as $(w_{d,1}, w_{d,2}, \dots, w_{d,n_d})$ in \mathcal{T} including n_d words from V . Queries are also strings (or sets of strings) composed of symbols drawn from Σ . Symbols in Σ may be letters, bytes, or even words and the documents may be articles, chromosomes or any other texts in which we need to search. In the explored systems, we can identify the following tasks.

– *Document Retrieval (Indexing, Ranking)*: for this task, two sub-tasks can be identified [4].

1. Document Listing: given query $Q = \{q_1, \dots, q_m | q_i \in \Sigma^*, \forall i\}$ and a text $\mathcal{T} \in \Sigma^*$ that is partitioned into n documents, $\{d_1, d_2, \dots, d_n\}$, the aim of this task is to return a list of the documents in which one or multiple tokens of Q appear at least once.
2. Document Ranking: given a query $Q = \{q_1, \dots, q_m | q_i \in \Sigma^*, \forall i\}$, an integer $0 < k \leq N$, and a text $\mathcal{T} \in \Sigma^*$ that is partitioned into n documents $\{d_1, d_2, \dots, d_n\}$, and returns the top- k documents ordered by a similarity measure $S(Q, d_i)$.

– *Passage Retrieval (Indexing, Ranking)*: Given a query Q , and a set of documents D where each document is partitioned into passages, the aim of this task is to find relevant passages for the query [6]. Passage retrieval can also be used for sentence highlighting.

- *Question Answering*: Given a Query $Q = \{q_1, q_2, \dots, q_m\}$ made of m tokens and a passage $P = \{p_1, p_2, \dots, p_k\}$ made of k tokens, the aim of this task is to find an answer span $A = \{a_{\text{start}}, a_{\text{end}}\}$ in P [5].
- *Summarization*: We will opt for the definition presented in [7]. Given a set of documents $D = d_i$ that we will call source documents, summarization aims to generate a text s (called summary) that is coherent and contains a significant amount of relevant information from the source text. Its compression rate $\tau = \frac{c(s)}{c(D)}$ (where $c(x)$ is the word count in x , x can be a sentence or document or any grouping of words) is less than a third of the length of the original document.
- *Topic Modeling*: The aim of topic modeling is to infer a set of K topics capturing a lower-dimensional representation suitable for summarization and prediction tasks [47]. According to [48], Given a text corpus \mathcal{T} with a vocabulary of size V , and the predefined number of topics K , the major tasks of topic modeling can be defined as:
 1. Learning the word representation of topics α : a topic α in a given collection \mathcal{T} is defined as a multinomial distribution over the vocabulary V , i.e., $p(w|\alpha)_{w \in V}$.
 2. Learning the sparse topic representation of documents θ : the topic representation of a document d , θ_d , is defined as a multinomial distribution over K topics, i.e., $p(\alpha_k|\theta_d)_{k=1, \dots, K}$.

In general, the task of topic modeling aims to find K salient topics $\alpha_{k=1, \dots, K}$ from \mathcal{T} and to find the topic representation of each document $\theta_{d=1, \dots, n}$.

- *FAQ Matching*: let F denote the set of question-answer pairs; given F and a user query Q , this task aims to rank question-answer pairs in F . Top k QA pairs with high scores are returned to the user [49].
- *Recommendation*: Given the set of all users \mathcal{C} , and the set of all possible items that can be recommended \mathcal{S} . Let u be a utility function that measures the usefulness of item s to user c , i.e., $u : \mathcal{C} \times \mathcal{S} \rightarrow \mathcal{R}$, where \mathcal{R} is a totally ordered set (e.g., non-negative integers or real numbers within a certain range). The goal of this task is to choose item(s) $s \in \mathcal{S}$ that maximize(s) the utility for each user $c \in \mathcal{C}$ [50].
- **Feedback Loop**: this characteristic is related to the use of user feedback data in any of the mentioned tasks.
- **Representation Level for Text**: in general, text can be represented in two distinct spaces: (a) bag-of-words space, (b) vector space. These representations can be done on one or multiple levels of granularity of textual documents, that is, *Document Level*, *Paragraph Level*, *Sentence Level*, and *Word Level*.
- **Representation Levels for Graphs**: Graphs can also be represented in a frequentist space or low dimensional vectorial space. These representations can be done on one or multiple levels of granularity of graphs, that is, *Full Graph Level*, *Sub-graph Level*, *Node Level*, *Edge Level*. Examples of graph representation in Covid Literature search engines are as follow:
 - Document Sub-graph Embedding: in order to make document level embeddings, [37,43] combined document level textual embeddings with embeddings of documents' related sub-graphs from the bigger KG to recommend similar documents.
- **Novelty**: a research paper is said to have novelty if the authors explored uncharted territories to solve old or new problems.
- **Data Enrichment**: Data enrichment refers, in general, to the process of adding more data to the already existing training data. Data enrichment methods can take two main forms, (a) data augmentation, and (b) data supplementation. The former is characteristic of the set of methods that use the already existing data to generate more data, while the latter encapsulates methods that use external resources in order to supplement the available data. The latter is easy to accomplish as long as there are external resources. There are various data augmentations methods. For example, in CO-Search [42], in order to train a Siamese network, the authors generated negative (*paragraph*, *reference*) pairs based on positive pairs extracted from documents.

• **Search Type:**

- *Keyword*: Keyword search refers to searching using queries composed out of one specific word.
 - *Regular Expression*: In this type of search the query takes the form of regular expressions that annotates textual patterns that we would like to retrieve. For example, [51] used this search strategy to look to drugs with certain properties in a drug re-purposing database.
 - *Open Questions*: this type of search refers to using natural language queries with simple or complex structures.
 - *Keyphrase Search*: this type of search refers to using queries composed of one or multiple keywords, and the order is taken into consideration.
- **KG Traversal**: this refers to the use of knowledge graphs to search for entities or relationships that are relevant to achieving one or multiple tasks.
- **Representation Combination (Rep.Comb.)**: this characteristic exists in one of two cases: (a) the combination of multiple levels of representation to achieve a task, or (b) the combination of KG and textual representation to achieve a task.

Table 5 offers an exhaustive list of search engines and their design specifications. While exploring search engines for COVID19 literature, we noticed multiple characteristics that are elaborated on in what follows:

- **Fast Prototyping and Deployment**: Given the urgent nature of most of the applications, the researcher opted mainly for off-the-shelf technologies that are easy to work with. In addition, except for one application, all the other applications used existing models and algorithms, which can also be attributed to the urgency of the task.
- **Textual Representation Methods**: there are two categories of methods: (a) Bag-of-Words (BOW) models, and (b) Vector Space Models (VSMs). The major difference is that VSMs capture more of the contextual elements of text than the BOW methods, but on the other hand the VSMs are computationally more expensive during training and inference. Some works struck a balance by applying both categories of methods e.g. [38,39,52,53], which is done generally by using a multi-stage ranking scheme that applies the first ranking using BOW models, which is then followed by a re-ranking using a VSM of the output of the previous ranking. Some works compensate for the latency of neural language models [13] by pre-indexing documents offline.
- **Granularity/Levels of Representations**: we also noticed that the works used different levels of granularity, which depends on the intended tasks and the available computational resources. For example, to achieve the task of document retrieval, some works opted for simple document level representations [54], while other works either used more granular representations [13,33,38,44,51,55–57], or a mix of more granular representations with document level representations [17,25,39,42,43,52,53,58].
- **Using KGs**: Knowledge graphs were used in multiple works for different purposes. For example, [43] used a KG (CKG [37]) embedding in tandem with textual representations for document recommendation, while [38] (CovEx KG [38]), [33] (Vapur KG [33]), and [39] (Citation KG [39]) traversed their respective KGs looking for similar entities to retrieve relevant papers. [57] (Blender-KG [35]) used a KG to extend queries and make the search more efficient.
- **Recommendation Modules**: Many search engines [33,38,43] use recommendation modules to offer more user oriented results.
- **Query Transformation/Extension**: Query transformation is also used in many applications to make the queries more expressive, which can help get more relevant results. For example, [54] used an extensive database of medical terms to augment the queries made by novices to search an academic biomedical corpus.
- **Multimedia (e.g., image, video, etc) Grounding**: Multimedia grounding is also used to couple textual data with relevant multimedia content. For example, [55] used a self-supervised method to couple biomedical text with corresponding coordinates in a human atlas. This mapping was

used to conduct two kinds of queries: (a) atlas based document retrieval using textual queries (which contain mentions of body parts), and (b) atlas based document retrieval using 3D atlas coordinates. In addition, [35] associated figures that depict molecular structures in research papers to their chemical entities that exist in a KG by using the captions of the figures. This was done to augment the KG.

3.2. Visual Exploratory Search

While exploring COVID19 literature, researchers can face two kinds of challenges: (a) quantity of the research papers, and (b) the quality of the research papers. Even though the textual exploratory search is a useful literature exploration tool, it is targeted and requires the researcher to know what she/he is looking for in advance, which is not always evident. Consequently, many visual exploratory search tools have been developed to explore the COVID19 literature in a visual, interactive, and general manner, rather than having to go through the tedious process of manually curating the literature. In the context of scientific literature, this can also be used to explore latent structures within the data which may be related to co-authorship networks, citation networks, and other important bibliometric dimensions.

In light of the reviewed literature, we can infer a general process that exploratory visual search applications follow. This process is presented in figure 5. The most important two phases of this process are (a) indicator specification and (b) indicator representation. The former is where one or multiple quantitative (e.g. entity types, topics, affiliation, etc.) or qualitative characteristics (e.g. occurrence/co-occurrence frequency/count) of the data are chosen to be presented and their method of presentation is also specified. The latter phase is where a significant visual representation is chosen for those indicators; for example, qualitative indicators can be presented using colors and quantitative indicators can be presented using distance, surface, or volume variations.

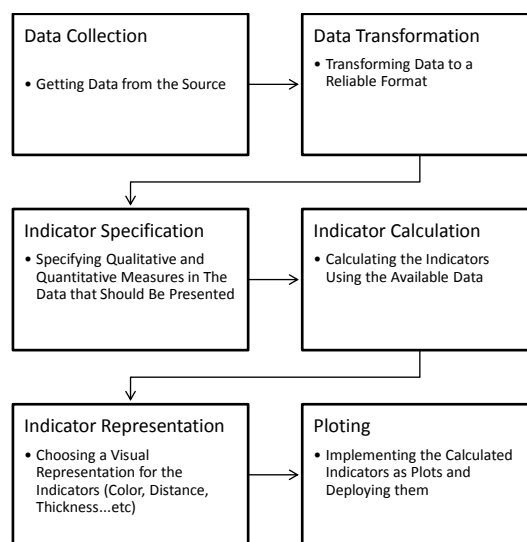


Figure 5. Exploratory Search Application creation Process.

The data used for the exploratory search applications is either CORD-19 [8] or one of the knowledge graphs presented previously. The frequency and count indicators are the most predominantly used, although other indicators are also used. For example, [59] uses topic similarity vectors to cluster similar topics. Multiple plots and visualization tools were used to visualize the indicators; these are summarized in table 6. In addition, some works use certain tasks in the data transformation phase in order to get more relevant data from raw text. The tasks mentioned in the works are information extraction (IE), which is generally attributed to basic textual information extraction, topic modeling, which was used in [59], and NER, which was used in [36,41,60] to extract

Table 5. Search Engine Comparison. "X" signifies "no" and "✓" signifies "yes".

System	CO-Search [42]	AWS Search (ACS) [43]	CORD-19 Drug Repository [51]	COVID-19 CovEx [38]	CovidXapua [52,58]	Yapua [33]	CovidASK [13]	CAiRE-Q [25]	COVID-19 Bridge [17]	EpiBridge [53]	SZ Covid [39]	SLIC [54]	SPIKE [56]	EVIDENCEMINER [57]
Rep.Comb.	Uses Raw Text (Uses KG)	✓(X)	✓(✓)	✓(X)	✓(X)	✓(✓)	✓(X)	✓(X)	✓(X)	✓(X)	✓(✓)	✓(X)	✓(X)	✓(✓)
	Publicly Available	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Feedback Loop	X	X	✓	✓	✓	X	X	X	X	X	X	X	X
	Multistage Ranking	X	X	X	✓	X	X	X	X	✓	✓	X	X	X
Text Representation Levels (KG Level)	KG Traversal	X	✓	X	X	✓	X	X	X	X	✓	X	X	✓
	Document (KG)	✓(X)	✓(X)	X(X)	X(X)	✓(X)	X(X)	✓(X)	X(X)	X(X)	✓(X)	✓(X)	X(X)	X(X)
	Paragraph (Sub-graph)	✓(X)	✓(✓)	X(X)	✓(X)	✓(X)	X(X)	✓(X)	✓(X)	✓(X)	✓(✓)	✓(X)	✓(X)	X(X)
	Sentence (Edge)	✓(X)	X(X)	X(X)	X(X)	X(X)	✓(X)	✓(X)	X(X)	✓(✓)	✓(X)	✓(X)	✓(X)	✓(X)
	Word (Node)	X(X)	X(X)	✓(X)	✓(X)	X(X)	X(X)	X(X)	X(X)	✓(X)	X(X)	X(X)	✓(X)	✓(X)
	n-gram (Node Property)	X(X)	X(X)	X(X)	X(X)	✓(X)	X(X)	X(X)	X(X)	X(X)	X(X)	X(X)	X(X)	X(X)
	Keyphrase (Edge Property)	X(X)	X(X)	✓(X)	X(X)	X(X)	X(X)	X(X)	X(X)	X(X)	X(X)	X(X)	✓(X)	✓(X)
	Inter-Level	✓	X	X	X	✓	X	✓	✓	✓	✓	N	✓	X
	Text & KG	X	✓	X	X	X	X	X	X	X	✓	✓	X	X
	Document Retrieval (Indexing, Ranking)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Tasks	Passage Retrieval (Indexing, Ranking)	✓	✓	X	✓	X	✓	✓	X	✓	✓	X	✓	✓
	Question Answering	✓	✓	X	✓	X	✓	✓	X	X	✓	X	X	X
	Summarization	✓	X	X	X	X	X	✓	X	X	✓	X	X	X
	Topic Modeling	X	✓	✓	X	X	X	X	X	✓	X	X	X	X
Search Type	Recommendation	✓	✓	✓	X	✓	X	X	X	X	X	X	X	X
	FAQ Matching	X	✓	X	X	X	X	X	X	X	X	X	X	X
	Keyword	X	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	Open Questions	✓	✓	X	X	✓	✓	✓	✓	✓	✓	✓	✓	✓
Data Enrichment	Keyphrases	X	✓	✓	X	X	X	✓	X	✓	✓	✓	✓	✓
	Regular Expression	X	X	✓	X	X	X	X	X	X	X	X	✓	X
	Novelty	X	X	X	X	X	X	X	X	X	X	X	X	X
	From External Resources	X	✓	X	✓	✓	✓	✓	✓	X	X	✓	✓	✓
	From Internal Resources	✓	✓	X	X	X	X	✓	✓	X	X	X	X	X

named entities and use their count as an indicator, and network analysis [36] and [60]. In [36], network analysis was used to solve two problems faced during network traversal, namely the problem of network size and the search for deep connections, using a breadth-first-search technique on the network structure. In [36], network analysis was used to detect communities within a co-authorship network, motivated by the need to keep track of what other groups were doing in order to explore new fields and potential collaborations. Reactivity is also an important feature in these tools since it simplifies interactive visual manipulation, which makes the exploration more flexible. Public availability is also looked into and links to the tools are provided if they exist.

4. Evaluation Methods

In general, machine learning models are composed of two main modules, (a) a representation module, and a (b) decision module. The former is responsible for transforming the data from a complex multidimensional space with latent spatial and temporal dependencies to a lower-dimensional, and more abstract space. The second module is used to process the representational modules' output to achieve a task. The training of these modules can be done independently, that is, the representational module can be trained separately in an unsupervised or self-supervised manner, while the combination of the two modules can be trained in a self-supervised, semi-supervised, or fully supervised manner.

The machine learning (ML) models used in the previously explored works, be it search engines related ML models or knowledge graph creation ML models (e.g. named entity recognition models), have to be evaluated to get empirical evidence on their viability. While exploring the literature, we noticed that there are two main evaluation techniques: human evaluation, and automatic evaluation. The former bases its evaluation on relevance judgment of the users, and the latter focuses on information needs in order to evaluate the results. The latter also has two sub-categories of evaluation measures: intrinsic evaluation measures, and extrinsic evaluation measures.

4.1. *Human Evaluation*

Human evaluation is based on quantifying human feedback towards the evaluated application. This type of evaluation is advantageous because of its integral character. Indeed, humans can evaluate more complex applications with multiple interacting modules. For example, in the case of a search engine, a human evaluator can assess the information relevance of the search results in addition to some representational aspects, like highlighting, which are not easy to evaluate automatically [52,58]. However, the downside of the human evaluation method is its irreproducibility, due to the fact that human evaluation is inherently biased and depends on the needs that the evaluators have, their field of expertise, and what they expect from the application. For example, an experienced researcher may find longer spans of text more reliable as answers to a query while a novice would generally prefer direct short answers [44]. This makes performance comparison of multiple applications based on human evaluation generally unreliable.

4.2. *Automatic Evaluation*

Automatic evaluation is the de facto evaluation method in machine learning literature. It is based on using evaluation metrics that quantify the discrepancy that exists between the model output and the wanted output. This is advantageous since it puts multiple applications on an equal footing during evaluation, which is advantageous. On the other hand, automatic evaluation is monolithic, meaning that it only evaluates one aspect of an application at a time (e.g. QA, DR, IR ...) and not the integrality of the application as is the case in human evaluation [52,58]. Furthermore, some aspects, such as ease-of-use and interface interactivity, cannot be evaluated automatically. In addition, the evaluation metrics used can suffer from certain biases that can lessen the validity of the evaluation. For example, [44] have found that automatic metrics such as F1 heavily penalize long answers, as they overlap poorly with the gold annotations, which are mostly short, factual answers.

Table 6. Exploratory Search Applications Summary.

System	Vidar-19 [61]	TopicMaps [59]	Network Visualisations [36]	SciSight [60]	semviz [41]	EvidenceMiner [57]
Available Charts	✗	✗	✗	✗	✗	✓
	✓	✓	✗	✓	✗	✗
	✗	✓	✗	✗	✓	✗
	✗	✗	✗	✗	✓	✗
	✓	✗	✗	✗	✗	✗
	✗	✓	✗	✓	✓	✗
	✓	✗	✗	✗	✗	✗
	✗	✗	✗	✗	✓	✓
	✗	✓	✗	✗	✗	✗
	✗	✗	✗	✓	✗	✗
Indicators	✓	✓	✓	✓	✗	✗
	✓	✓	✓	✓	✓	✓
	✗	✓	✗	✗	✗	✗
	✓	✓	✓	✓	✓	✓
	✗	✓	✗	✓	✗	✗
Related Tasks	✗	✗	✗	✓	✗	✗
	✗	✗	✓	✓	✓	✗
	✗	✗	✓	✓	✗	✗
	✗	✗	✓	✓	✗	✗
Data Source	✓	✓	✗	✓	✓	✓
	✗	✗	✓	✓	✓	✓
Reactivity	✓	✓	✓	✓	✓	✓
Public Availability	✓ ^a	✓ ^b	✓ ^c	✓ ^d	✓ ^e	✓ ^f

^a <https://fran6wol.eu.pythonanywhere.com/>

^b <http://strategicfutures.org/TopicMaps/COVID-19/dimensions.html>, <http://nlp.inspirata.com/NetworkVisualisations/TitleNetwork/>, <https://nlp.inspirata.com/NetworkVisualisations/TopicMaps/COVID-19/cord19.html>

^c <https://nlp.inspirata.com/NetworkVisualisations/TitleNetwork/>, <https://nlp.inspirata.com/NetworkVisualisations/TreatmentNetwork/>, https://nlp.inspirata.com/NetworkVisualisations/

^d <https://scisight.apps.allenai.org/>

^e <https://www.semviz.org/>

^f <https://evidenceminer.firebaseapp.com/analytics?kw=CORONAVIRUS&corpus=covid-19>

As was mentioned before, automatic evaluation measures can be categorized into (a) intrinsic evaluation measures (IEMs), and (b) extrinsic evaluation measures (EEMs). The former measures are generally used to evaluate representation modules separately, and the latter measures are used to evaluate the combined representation and decision downstream model.

4.2.1. *Intrinsic Evaluation*

In the explored works, we only found one example of intrinsic evaluation [32], where KG node embeddings are evaluated by comparing the Pearson and Spearman correlation scores between the ratings and the cosine similarity scores of entities.

4.2.2. *Extrinsic Evaluation*

In contrast to IEM, EEMs are more frequently used. The works that we explored use a plethora of EEMs which depend on the kind of tasks to be evaluated. This type of evaluation is done through multiple evaluation metrics that are task-specific. Multiple evaluation measures and their variants were used. For example, the ROUGE evaluation metric [62] and its variants were used in [25] to evaluate summarization models. The Match [63] [64] method was used in [45] to evaluate QA and IR. Other more standard evaluation metrics such as recall and precision were used for IR tasks [65].

5. Discussion and future research directions

In general, the explored works have certain common limitations. In what follows we summarize a few of them:

- **Evaluation:** Most of the applications (e.g. [44,52,58]) suffer from a monolithic evaluation scheme that focuses on one task in particular and ignores other aspects of the application, especially those related to visual aspects.
- **Feedback Loop:** Some applications (e.g. [52,58]) expressed the importance of including human input in the process of information retrieval, as it tends to balance information need and information relevance.
- **Fact Checking:** Due to the rapid expansion of COVID19 literature, and the existence of many contradictory claims, concerning for example the incubation period of the virus and the optimal social distancing protocol stresses the importance of fact checking applications for COVID19 claims. [66] created a claim verification application for COVID19 literature, which uses a passage and a claim as input and outputs if the claim is true or not given the passage. This type of application needs huge amounts of annotated data, which is particularly cumbersome in the case of COVID19 since it needs skilled specialists to annotate it. Developing semi-supervised or unsupervised techniques would be useful.
- **Extending Data:** Most of the applications (e.g. [54,55]) used limited amounts of data (labeled or not) to perform tasks, either because of the lack of labeled data or because of the lack of computational resources. More data would certainly improve performance.
- **Data Bias:** Some applications (e.g. [55]) can also benefit from reducing data bias, especially gender bias.
- **Smart Querying:** Some applications [57] use query functionalities that tend to be limited to simple word matching. This can be problematic in cases where the intent of the user is not evident in the query. This can be remedied by using embedding based query matching, which uses contextual information for matching the queries to the results.

6. Conclusion

This work represents an exploration of COVID19 literature exploration applications, with emphasis on their design principles and concepts. There are two main types of literature exploration applications, a) exploratory textual search and b) exploratory visual search. The former uses textual

queries made by end-users in order to explore the knowledge base and send the most relevant documents back to the users, while the latter type of application uses visual summaries to offer a structured view of the existing literature.

7. Limitations of this Work

Empirical quantitative evaluation of the systems explored in this work was of interest, but the discrepancies that were found in evaluation results of the same systems in multiple sources (e.g. the results given in [43] are different from those given in [52] for the same system: Covidex), in addition to the unavailable implementation details of some systems, discouraged us from pursuing this objective in this work.

Conflicts of Interest: No conflict of interest.

References

- da Silva, J.A.T.; Tsigaris, P.; Erfanmanesh, M. Publishing volumes in major databases related to Covid-19. *Scientometrics* **2021**, *126*, 831–842.
- Chen, Q.; Allot, A.; Lu, Z.; others. Keep up with the latest coronavirus research. *Nature* **2020**, 579.
- Jo, T. Text mining. *Studies in Big Data*. Cham: Springer International Publishing **2019**.
- Culpepper, J.S.; Navarro, G.; Puglisi, S.J.; Turpin, A. Top-k ranked document search in general text databases. *European Symposium on Algorithms*. Springer, 2010, pp. 194–205.
- Liu, X.; Shen, Y.; Duh, K.; Gao, J. Stochastic answer networks for machine reading comprehension. *arXiv preprint arXiv:1712.03556* **2017**.
- Ganesh, S.; Varma, V. Passage retrieval using answer type profiles in question answering. *Proceedings of the 23rd Pacific Asia Conference on Language, Information and Computation*, Volume 2, 2009, pp. 559–568.
- Torres-Moreno, J.M. *Automatic text summarization*; John Wiley & Sons, 2014.
- Wang, L.L.; Lo, K.; Chandrasekhar, Y.; Reas, R.; Yang, J.; Eide, D.; Funk, K.; Kinney, R.; Liu, Z.; Merrill, W.; others. CORD-19: The Covid-19 Open Research Dataset. *ArXiv* **2020**.
- Voorhees, E.; Alam, T.; Bedrick, S.; Demner-Fushman, D.; Hersh, W.R.; Lo, K.; Roberts, K.; Soboroff, I.; Wang, L.L. TREC-COVID: Constructing a Pandemic Information Retrieval Test Collection. *arXiv preprint arXiv:2005.04474* **2020**.
- Wang, L.L.; Lo, K. Text mining approaches for dealing with the rapidly expanding literature on COVID-19. *Briefings in Bioinformatics* **2020**.
- Tsatsaronis, G.; Schroeder, M.; Paliouras, G.; Almirantis, Y.; Androutsopoulos, I.; Gaussier, E.; Gallinari, P.; Artieres, T.; Alvers, M.R.; Zschunke, M.; others. BioASQ: A Challenge on Large-Scale Biomedical Semantic Indexing and Question Answering. *AAAI fall symposium: Information retrieval and knowledge discovery in biomedical text*. Citeseer, 2012.
- Tang, R.; Nogueira, R.; Zhang, E.; Gupta, N.; Cam, P.; Cho, K.; Lin, J. Rapidly Bootstrapping a Question Answering Dataset for COVID-19. *arXiv preprint arXiv:2004.11339* **2020**.
- Lee, J.; Yi, S.S.; Jeong, M.; Sung, M.; Yoon, W.; Choi, Y.; Ko, M.; Kang, J. Answering questions on covid-19 in real-time. *arXiv preprint arXiv:2006.15830* **2020**.
- Möller, T.; Reina, A.; Jayakumar, R.; Pietsch, M. COVID-QA: A Question Answering Dataset for COVID-19. *Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020*, 2020.
- Poliak, A.; Fleming, M.; Costello, C.; Murray, K.W.; Yarmohammadi, M.; Pandya, S.; Irani, D.; Agarwal, M.; Sharma, U.; Sun, S.; others. Collecting verified covid-19 question answer pairs **2020**.
- Nguyen, T.; Rosenberg, M.; Song, X.; Gao, J.; Tiwary, S.; Majumder, R.; Deng, L. Ms marco: A human-generated machine reading comprehension dataset **2016**.
- MacAvaney, S.; Cohan, A.; Goharian, N. SLEDGE-Z: A Zero-Shot Baseline for COVID-19 Literature Search. *arXiv preprint arXiv:2010.05987* **2020**.
- Kwiatkowski, T.; Palomaki, J.; Redfield, O.; Collins, M.; Parikh, A.; Alberti, C.; Epstein, D.; Polosukhin, I.; Devlin, J.; Lee, K.; others. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics* **2019**, *7*, 453–466.

- 437 19. Rajpurkar, P.; Zhang, J.; Lopyrev, K.; Liang, P. Squad: 100,000+ questions for machine comprehension of
438 text. *arXiv preprint arXiv:1606.05250* **2016**.
- 439 20. Arora, A.; Shrivastava, A.; Mohit, M.; Lecanda, L.S.M.; Aly, A. Cross-lingual Transfer Learning for Intent
440 Detection of Covid-19 Utterances **2020**.
- 441 21. Choi, E.; He, H.; Iyyer, M.; Yatskar, M.; Yih, W.t.; Choi, Y.; Liang, P.; Zettlemoyer, L. Quac: Question
442 answering in context. *arXiv preprint arXiv:1808.07036* **2018**.
- 443 22. Dang, H.T. Overview of DUC 2005. Proceedings of the document understanding conference, 2005, Vol.
444 2005, pp. 1–12.
- 445 23. Hoa, T. Overview of DUC 2006. Document Understanding Conference, 2006.
- 446 24. Nema, P.; Khapra, M.; Laha, A.; Ravindran, B. Diversity driven attention model for query-based abstractive
447 summarization. *arXiv preprint arXiv:1704.08300* **2017**.
- 448 25. Dan, S.; Xu, Y.; Yu, T.; Siddique, F.B.; Barezi, E.; Fung, P. CAiRE-COVID: A Question Answering and
449 Query-focused Multi-Document Summarization System for COVID-19 Scholarly Information Management
450 **2020**.
- 451 26. Kim, J.D.; Ohta, T.; Tateisi, Y.; Tsujii, J. GENIA corpus—a semantically annotated corpus for bio-textmining.
452 *Bioinformatics* **2003**, *19*, i180–i182.
- 453 27. Kim, J.D.; Ohta, T.; Tsuruoka, Y.; Tateisi, Y.; Collier, N. Introduction to the bio-entity recognition task at
454 JNLPBA. Proceedings of the international joint workshop on natural language processing in biomedicine
455 and its applications. Citeseer, 2004, pp. 70–75.
- 456 28. Krallinger, M.; Rabal, O.; Leitner, F.; Vazquez, M.; Salgado, D.; Lu, Z.; Leaman, R.; Lu, Y.; Ji, D.; Lowe,
457 D.M.; others. The ChEMDNER corpus of chemicals and drugs and its annotation principles. *Journal of*
458 *cheminformatics* **2015**, *7*, 1–17.
- 459 29. Doğan, R.I.; Leaman, R.; Lu, Z. NCBI disease corpus: a resource for disease name recognition and concept
460 normalization. *Journal of biomedical informatics* **2014**, *47*, 1–10.
- 461 30. Kringelum, J.; Kjaerulff, S.K.; Brunak, S.; Lund, O.; Oprea, T.I.; Taboureau, O. ChemProt-3.0: a global
462 chemical biology diseases mapping. *Database* **2016**, *2016*.
- 463 31. Li, J.; Sun, Y.; Johnson, R.J.; Sciaky, D.; Wei, C.H.; Leaman, R.; Davis, A.P.; Mattingly, C.J.; Wiegers, T.C.; Lu,
464 Z. BioCreative V CDR task corpus: a resource for chemical disease relation extraction. *Database* **2016**, *2016*.
- 465 32. Basu, S.; Chakraborty, S.; Hassan, A.; Siddique, S.; Anand, A. ERLKG: Entity Representation Learning and
466 Knowledge Graph based association analysis of COVID-19 through mining of unstructured biomedical
467 corpora. Proceedings of the First Workshop on Scholarly Document Processing, 2020, pp. 127–137.
- 468 33. Köksal, A.; Dönmez, H.; Özçelik, R.; Ozkirimli, E.; Özgür, A. Vapur: A Search Engine to Find Related
469 Protein–Compound Pairs in COVID-19 Literature. *arXiv preprint arXiv:2009.02526* **2020**.
- 470 34. Amini, A.; Hope, T.; Wadden, D.; van Zuylen, M.; Horvitz, E.; Schwartz, R.; Hajishirzi, H. Extracting a
471 knowledge base of mechanisms from COVID-19 papers. *arXiv preprint arXiv:2010.03824* **2020**.
- 472 35. Wang, Q.; Li, M.; Wang, X.; Parulian, N.; Han, G.; Ma, J.; Tu, J.; Lin, Y.; Zhang, H.; Liu, W.; others.
473 Covid-19 literature knowledge graph construction and drug repurposing report generation. *arXiv preprint*
474 *arXiv:2007.00576* **2020**.
- 475 36. Cernile, G.; Heritage, T.; Sebire, N.J.; Gordon, B.; Schwering, T.; Kazemlou, S.; Borecki, Y. Network
476 graph representation of COVID-19 scientific publications to aid knowledge discovery. *BMJ Health & Care*
477 *Informatics* **2020**, *28*.
- 478 37. Wise, C.; Ioannidis, V.N.; Calvo, M.R.; Song, X.; Price, G.; Kulkarni, N.; Brand, R.; Bhatia, P.; Karypis,
479 G. COVID-19 knowledge graph: accelerating information retrieval and discovery for scientific literature.
480 *arXiv preprint arXiv:2007.12731* **2020**.
- 481 38. Rahdari, B.; Brusilovsky, P.; Thaker, K.; Chau, H.K. CovEx: An Exploratory Search System for COVID-19
482 Scientific Literature.
- 483 39. Das, D.; Katyal, Y.; Verma, J.; Dubey, S.; Singh, A.; Agarwal, K.; Bhaduri, S.; Ranjan, R. Information retrieval
484 and extraction on covid-19 clinical articles using graph community detection and bio-bert embeddings.
485 Proceedings of the 1st Workshop on NLP for COVID-19 at ACL 2020, 2020.
- 486 40. Ciampaglia, G.L.; Shiralkar, P.; Rocha, L.M.; Bollen, J.; Menczer, F.; Flammini, A. Computational fact
487 checking from knowledge networks. *PloS one* **2015**, *10*, e0128193.
- 488 41. Tu, J.; Verhagen, M.; Cochran, B.; Pustejovsky, J. Exploration and discovery of the COVID-19 literature
489 through semantic visualization. *arXiv preprint arXiv:2007.01800* **2020**.

42. Esteva, A.; Kale, A.; Paulus, R.; Hashimoto, K.; Yin, W.; Radev, D.; Socher, R. Co-search: Covid-19 information retrieval with semantic search, question answering, and abstractive summarization. *arXiv preprint arXiv:2006.09595* **2020**.

43. Bhatia, P.; Arumae, K.; Pourdamghani, N.; Deshpande, S.; Snively, B.; Mona, M.; Wise, C.; Price, G.; Ramaswamy, S.; Kass-Hout, T. AWS CORD19-search: a scientific literature search engine for COVID-19. *arXiv preprint arXiv:2007.09186* **2020**.

44. Otegi, A.; Campos, J.A.; Azkune, G.; Soroa, A.; Agirre, E. Automatic Evaluation vs. User Preference in Neural Textual Question Answering over COVID-19 Scientific Literature. Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020, 2020.

45. Gangi Reddy, R.; Iyer, B.; Arafat Sultan, M.; Zhang, R.; Sil, A.; Castelli, V.; Florian, R.; Roukos, S. End-to-End QA on COVID-19: Domain Adaptation with Synthetic Training. *arXiv e-prints* **2020**, pp. arXiv-2012.

46. Lee, S.; Sedoc, J. Using the Poly-encoder for a COVID-19 Question Answering System. Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020, 2020.

47. Virtanen, S.; Girolami, M. Precision-Recall Balanced Topic Modelling. Advances in Neural Information Processing Systems, 2019, pp. 6750–6759.

48. Qiang, J.; Qian, Z.; Li, Y.; Yuan, Y.; Wu, X. Short text topic modeling techniques, applications, and performance: a survey. *IEEE Transactions on Knowledge and Data Engineering* **2020**.

49. Damani, S.; Narahari, K.N.; Chatterjee, A.; Gupta, M.; Agrawal, P. Optimized Transformer Models for FAQ Answering. Pacific-Asia Conference on Knowledge Discovery and Data Mining. Springer, 2020, pp. 235–248.

50. Adomavicius, G.; Tuzhilin, A. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE transactions on knowledge and data engineering* **2005**, *17*, 734–749.

51. Tworowski, D.; Gorohovski, A.; Mukherjee, S.; Carmi, G.; Levy, E.; Detroja, R.; Mukherjee, S.B.; Frenkel-Morgenstern, M. COVID19 Drug Repository: text-mining the literature in search of putative COVID19 therapeutics. *Nucleic acids research* **2020**, p. 1.

52. Zhang, E.; Gupta, N.; Tang, R.; Han, X.; Pradeep, R.; Lu, K.; Zhang, Y.; Nogueira, R.; Cho, K.; Fang, H.; others. Covidex: Neural ranking models and keyword search infrastructure for the covid-19 open research dataset. *arXiv preprint arXiv:2007.07846* **2020**.

53. Farokhnejad, M.; Pranesh, R.R.; Vargas-Solar, G.; Mehr, D.A. S_Covid: An Engine to Explore COVID-19 Scientific Literature.

54. He, D.; Wang, Z.; Thaker, K.; Zou, N. Translation and expansion: Enabling laypeople access to the COVID-19 academic collection. *Data and Information Management* **2017**, *1*.

55. Grujicic, D.; Radevski, G.; Tuytelaars, T.; Blaschko, M.B. Self-supervised context-aware Covid-19 document exploration through atlas grounding **2020**.

56. Tabib, H.T.; Shlain, M.; Sadde, S.; Lahav, D.; Eyal, M.; Cohen, Y.; Goldberg, Y. Interactive extractive search over biomedical corpora. Proceedings of the 19th SIGBioMed Workshop on Biomedical Language Processing, 2020, pp. 28–37.

57. Wang, X.; Guan, Y.; Liu, W.; Chauhan, A.; Jiang, E.; Li, Q.; Liem, D.; Sigdel, D.; Caufield, J.; Ping, P.; others. Evidenceminer: Textual evidence discovery for life sciences. Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, 2020, pp. 56–62.

58. Zhang, E.; Gupta, N.; Nogueira, R.; Cho, K.; Lin, J. Rapidly deploying a neural search engine for the covid-19 open research dataset: Preliminary thoughts and lessons learned. *arXiv preprint arXiv:2004.05125* **2020**.

59. Le Bras, P.; Gharavi, A.; Robb, D.A.; Vidal, A.F.; Padilla, S.; Chantler, M.J. Visualising COVID-19 Research. *arXiv preprint arXiv:2005.06380* **2020**.

60. Hope, T.; Portenoy, J.; Vasan, K.; Borchardt, J.; Horvitz, E.; Weld, D.S.; Hearst, M.A.; West, J. SciSight: Combining faceted navigation and research group detection for COVID-19 exploratory scientific search. *arXiv preprint arXiv:2005.12668* **2020**.

61. Wolinski, F. Visualization of Diseases at Risk in the COVID-19 Literature. *arXiv preprint arXiv:2005.00848* **2020**.

62. Lin, C.Y. Rouge: A package for automatic evaluation of summaries. Text summarization branches out, 2004, pp. 74–81.

542 63. Chen, D.; Fisch, A.; Weston, J.; Bordes, A. Reading wikipedia to answer open-domain questions. *arXiv*

543 *preprint arXiv:1704.00051* **2017**.

544 64. Karpukhin, V.; Oğuz, B.; Min, S.; Wu, L.; Edunov, S.; Chen, D.; Yih, W.t. Dense Passage Retrieval for

545 Open-Domain Question Answering. *arXiv preprint arXiv:2004.04906* **2020**.

546 65. Zhu, M. Recall, precision and average precision. *Department of Statistics and Actuarial Science, University of*

547 *Waterloo, Waterloo* **2004**, 2, 6.

548 66. Wadden, D.; Lo, K.; Wang, L.L.; Lin, S.; van Zuylen, M.; Cohan, A.; Hajishirzi, H. Fact or Fiction: Verifying

549 Scientific Claims. *arXiv preprint arXiv:2004.14974* **2020**.