

# textNet: Directed, Multiplex, Multimodal Event Network Extraction from Textual Data

Elise Zufall<sup>1</sup>, Tyler Scott<sup>1</sup>, and <sup>1</sup>

<sup>1</sup> University of California, Davis Department of Environmental Science and Policy

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

## Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: [Open Journals](#) 

## Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

## License

Authors of papers retain copyright  
and release the work under a  
Creative Commons Attribution 4.0  
International License ([CC BY 4.0](#)).

## Introduction

Network measurement in social science typically relies on data collected through surveys and interviews. Document-based measurement is automatable and scalable, providing opportunities for large scale or longitudinal research that are not possible through traditional methods. A number of tools exist to generate networks based on co-occurrence of words within documents (such as the [Nocodefunctions](#) app ([Levallois et al., 2012](#)), the “[textnets](#)” package ([Bail, 2024](#)), [InfraNodus](#) ([Paranyushkin, 2018](#)), and many more). But there is, to our knowledge, no open-source tool that generates network data based on the syntactic relationships between entities within a sentence. *textNet* allows a user to input one or more PDF documents and create arbitrarily complex directed, multiplex, and multimodal network graphs. *textNet* also works on arbitrarily long documents, making it well suited for research applications using long texts such as government planning documents, court proceedings, regulatory impact analyses, and environmental impact assessments.

## Statement of Need

Network extraction from documents has typically required manual coding. Furthermore, existing network extraction methods that use co-occurrence leave a vast amount of data on the table, namely, the rich edge attribute data and directionality of each verb phrase defining the particular relationship between two entities, and the respective roles of the entity nodes involved in that verb phrase. We present an R package, *textNet*, designed to enable directed, multiplex, multimodal network extraction from text documents through syntactic dependency parsing, in a replicable, automated fashion for collections of arbitrarily long documents. The *textNet* package facilitates the automated analysis and comparison of many documents, based on their respective network characteristics. Its flexibility allows for any desired entity categories, such as organizations, geopolitical entities, dates, or custom-defined categories, to be preserved.

## Directed Graph Production

As a syntax-based network extractor, *textNet* identifies source and target nodes. This produces directed graphs that contain information about network flow. Methods based on identifying co-occurring nodes in a document, by contrast, produce undirected graphs. *textNet* also allows the user to code ties based on co-occurrence in a designated piece of text if desired.

## Multiplex Graph Output

Syntax-based measurement encodes edges based on subject-verb-object relationships. *textNet* stores verb information as edge attributes, which allows the user to preserve arbitrarily complex topological layers (of different types of relationships) or customize groupings of edge types to simplify representation.

## 39 Multimodal Graph Output

40 Multimodal networks, or networks where there are multiple categories of nodes, have common  
41 use cases such as social-ecological network analysis of configurations of actors and environmental  
42 features. Existing packages such as the *manynet* package (Hollway, 2024) provide analytical  
43 functions for multimodal network statistics. *textNet* provides a structure for tagging and  
44 organizing arbitrarily complex node labeling schemes that can then be fed into packages for  
45 multi-node network statistical analysis. Node labels can be automated (e.g., the default entity  
46 type tags for an NLP engine such as *spaCy* (Honribal et al., 2021)), customized using a  
47 dictionary, or based on a hybrid scheme of default and custom labels. Any node type is possible  
48 (e.g., species, places, people, concepts, etc.) so this can be adapted to domain specific research  
49 applications by applying dictionaries or using a custom NER model.

## 50 Avoids Saturation

51 Co-occurrence graphs have the tendency to generate saturated subgraphs, since every co-  
52 occurring collection of entities has every possible edge drawn amongst them. By contrast,  
53 *textNet* draws connections not between every entity in the document or even the sentence,  
54 but specifically between pairs of entities that are mediated by an event relationship. This leads  
55 to sparser graphs that preserve the ability for greater structural variance, and correspondingly,  
56 network analysis of structural attributes of the graphs.

## 57 Installation

58 The stable version of this package can be installed from Github, using the *devtools* package  
59 (Wickham et al., 2022):

```
60 devtools::install_github("ucd-cep/textnet")
```

61 The *textNet* package suggests several convenience wrappers of packages such as *spacyr* (Benoit  
62 et al., 2023), *pdftools* (Ooms, 2024), *igraph* (Csárdi et al., 2024), and *network* (Butts et al.,  
63 2023). To use the full functionality of *textNet*, such as pre-processing tools and post-processing  
64 analysis tools, we recommend installing these packages, which for *spacyr* requires integration  
65 with Python. However, the user may wish to preprocess and parse data using their own NLP  
66 engine, and skip directly to the *textnet\_extract()* function, which does not depend on any of  
67 the aforementioned packages. The *textnet\_extract()* function does, however, use functions  
68 from *pbapply* (Solymos et al., 2023), *data\_table* (Barrett et al., 2024), *dplyr* (Wickham et al.,  
69 2023), and *tidyr* (Wickham et al., 2024).

## 70 Overview and Main Functions

71 The package architecture relies on four sets of functions around core tasks:

- 72 ■ [OPTIONAL] Pre-processing: *pdf\_clean()*, a wrapper for the *pdftools::pdf\_text()* func-  
73 tion which includes a custom header/footer text removal feature; and *parse\_text()*,  
74 which is a wrapper for the *spacyr* package and uses the *spaCy* natural language pro-  
75 cessing engine (Honribal et al., 2021) to parse text and perform part of speech tagging,  
76 dependency parsing, and named entity recognition (NER). Alternatively, as described  
77 below, the user can skip this step and load parsed text directly into the package.
- 78 ■ Network extraction: *textnet\_extract()*, which generates a graph database from parsed  
79 text based upon tags and dependency relations
- 80 ■ Disambiguation: tools for cleaning, recoding, and aggregating node and edge attributes,  
81 such as the *find\_acronyms()* function, which can be paired with the *disambiguation()*  
82 function to identify acronyms in the text and replace them with the full entity name.
- 83 ■ Exploration: the *export\_to\_network()* function for exporting the graph database to  
84 *igraph* and *network* objects, *top\_features()* for viewing node and edge attributes, and

85 combine\_networks() for aggregating multiple document-based graphs based on common  
86 nodes.

## 87 Example

88 The following example uses parsed text from the Gravelly Ford Water District Groundwater  
89 Sustainability Plan in the state of California, before and after the plan underwent revisions  
90 required by the California Department of Water Resources. Both versions of the plan were  
91 pre-processed using the optional pdf\_clean() and parse\_text() functions, as shown in the  
92 appendix below and package repository. textNet is designed for modularity with respect to  
93 pdf-to-text conversion and NLP engine. The user can derive plain text by any approach, and  
94 likewise perform event extraction with any NLP engine or large language model (LLM) (more  
95 on LLM extensions below) and bring these data to textNet. The textnet\_extract() function  
96 expects the parsed table to follow specific conventions for column names and speech tagging,  
97 so externally produced data must be converted to standards outlined in the package manual.

## 98 Extract Networks

99 First, we read in the pre-processed data and call textnet\_extract() to produce the network  
100 object:

```
101 library(textNet)
102 old_new_parsed <- textNet::old_new_parsed
103
104 extracts <- vector(mode="list",length=length(old_new_parsed))
105   for(m in 1:length(old_new_parsed)){
106     extracts[[m]] <- textnet_extract(old_new_parsed[[m]],concatenator="_",cl=4,
107                                   keep_entities = c('ORG','GPE','PERSON','WATER'),
108                                   return_to_memory=T, keep_incomplete_edges=T)
109   }
110
111 ## [1] "crawling 802 sentences"
112 ## [1] "crawling 1090 sentences"
```

113 The textnet\_extract() function extracts the entity network. It reads in the result of parse\_text()  
114 as described in the appendix, or another parsing tool with appropriate column names and  
115 tagging conventions. The resulting object consists of a nodelist, an edgelist, a verblist, and  
116 a list of appositives. The nodelist variables are entity\_name, the concatenated name of the  
117 entity; entity\_type, which is a preservation from the entity\_type attribute from the output of  
118 textNet::parse\_text(); and num\_appearances, which is the number of times the entity appears  
119 in the PDF text. (This is not the same as node degree, since there may be multiple edges, or  
120 if keep\_incomplete\_edges is set to false, no edges resulting from a single appearance of the  
121 entity in the document.) The entity\_type attribute represents spaCy's determination of entity  
122 type using its NER recognition, or if a custom parser or NER tool is used, the textnet\_extract()  
123 function will preserve these entity type designations.

124 The file is saved to the provided filename, if provided. It is returned to memory if re-  
125 turn\_to\_memory is set to T. At least one of these return pathways must be established  
126 to avoid an error. In this example, we only keep entity types in the nodelist, edgelist, and  
127 appositivelist that are listed under keep\_entities; namely, "ORG", "GPE", "PERSON", and  
128 "WATER".

129 The resulting object consists of a nodelist, an edgelist, a verblist, and a list of appositives. The  
130 nodelist variables are entity\_name, the concatenated name of the entity; entity\_type, which is  
131 a preservation from the entity\_type attribute from the output of textNet::parse\_text(); and  
132 num\_appearances, which is the number of times the entity appears in the PDF text. The  
133 default entity types are based on spaCy's NER tags, but entity types can be customized as

desired. In this example, we only keep entity types in the nodelist, edgelist, and appositivelist that are listed under keep\_entities; namely, "ORG", "GPE", "PERSON", and "WATER".

### Consolidate Entity Synonyms

In a document, the same real-world entity may be referenced in multiple ways. For instance, the document may introduce an organization using its full name, then use an acronym for the remainder of the document. To have more reliable network results, it is important to consolidate nodes that represent different naming conventions into a single node. The *textNet* package comes with a built-in tool for finding acronyms defined parenthetically within the text. This can be run on the result of pdf\_clean() to generate a table with one column for acronyms and another for the corresponding full names, such that each row is a different instance of a phrase for which an acronym was detected. The use of find\_acronyms() is demonstrated below.

```
old_new_text <- textNet::old_new_text
old_acronyms <- find_acronyms(old_new_text[[1]])
new_acronyms <- find_acronyms(old_new_text[[2]])

print(head(old_acronyms))
```

	name	acronym
	<char>	<char>
## 1:	Central_Valley	CV
## 2:	Total_Dissolved_Solids	TDS
## 3:	California_Code_of_Regulations	CCR
## 4:	Department_of_Water_Resources	DWR
## 5:	Best_Management_Practice	BMP
## 6:	Gravelly_Ford_Water_District	GFWD

The resulting table of acronyms can then be fed into a disambiguation tool, the *textNet* function disambiguate(). This tool is very flexible, allowing a user-defined custom vector or list of strings representing the original entity name to search for in the textnet\_extract object, and another user-defined custom vector or list of strings representing the entity name to which to convert those instances. Additional inputs that may be useful here are names and abbreviations of known federal and state or regional agencies, or other entities that are likely to be discussed in the particular type of document being analyzed. There may also be topic-specific words or phrases that are likely to be discussed in the document. For instance, in Groundwater Sustainability Plans, it is common to discuss entities that involve the term "subbasin," but the spelling of this is not always consistent.

In the example below, we define a "from" vector that includes the acronyms found through the previous step, as well as non-standard spellings of "subbasin." This function is case sensitive, so we have included two alternate cases that are likely to appear in the dataset. The "to" vector includes the full names from the find\_acronyms result, along with the standard spelling of "subbasin".

There are a few rules about defining the "from" and "to" columns. First, the length of "from" and "to" must be identical, since from[[i]] is replaced with to[[i]]. Second, there may not be any duplicated terms in the "from" list, since each string must be matched to a single replacement without ambiguity. It is acceptable to have duplicated terms in the "to" list.

The "from" argument may be formatted as either a vector or a list. However, if it is a list, no element may contain more than one string. The "match\_partial\_entity" argument defaults to F for each element of "from" and "to." However, it can be set to T or F for each individual element. (Replacing an acronym with its full name may only be wise if the entire name of the node is that acronym. Otherwise "EPA" could accidentally match on "NEPAL" and create a nonsense entity called "NEnvironmental\_Protection\_AgencyL". The risk of this for modern,

185 sentence-case documents is decreased, as `disambiguate()` is intentionally a case-sensitive  
 186 function.) For the example below, we set `match_partial_entity` to F for each of the acronyms,  
 187 but to T for the word "Sub\_basin," since "Sub\_basin" may very well be a portion of a longer  
 188 entity, for which we would want to standardize the spelling.

189 Each element in the "from" object must be a single character vector. This is not the case  
 190 for the "to" argument; a user may define elements of "to" to contain multiple character  
 191 vectors in order to convert a single node into multiple nodes. Specifically, there may be some  
 192 cases in which one would want to convert a single node into multiple nodes, each preserving  
 193 the original node's edges to other nodes. For instance, suppose a legal document refers to  
 194 "The\_Defendants" as a shorthand for referring to three individuals involved in the case. In  
 195 the network, it may be desirable for these individuals to be represented as their own separate  
 196 nodes, especially if the network is to be merged with those resulting from other documents,  
 197 where the three defendants may be named separately. To convert this single node into multiple  
 198 nodes that preserve all of their original edges to other entities, `from[[j]]` should be set to  
 199 "The\_Defendants", and `to[[j]]` should be set to a string vector including the individuals' names,  
 200 such as `c("John_Doe", "Jane_Doe", "Emily_Doe")`.

201 The default behavior is to loop through the disambiguation recursively, though by setting  
 202 `recursive` to F, this can be overridden. The difference can be seen in the following example.  
 203 Suppose that the following from list and to list are defined: `from = c("MA", "Mass")`; `to =`  
 204 `c("Mass", "Massachusetts")`. If `recursive = F`, all instances of MA in the original `textnet_extract`  
 205 object would be set to Mass, and all instances of Mass in the original `textnet_extract` object  
 206 would be set to Massachusetts. If `recursive = T`, all instances of MA and Mass in the original  
 207 `textnet_extract` object would be set to Massachusetts. The ability to toggle this behavior can  
 208 be useful when concatenating a large from and to list based on multiple sources.

209 The `disambiguate()` function is designed to be usable even for very large graphs; when  
 210 disambiguating thousands of nodes, a user may choose to use web scraping or another  
 211 automated tool to help generate a long list of "from" and "to" elements by which to merge  
 212 or separate the nodes of the graph. Use of an automated tool to generate "to" and "from"  
 213 columns with hundreds or thousands of elements can lead to uncertainty about the behavior  
 214 of the "to" and "from" columns. Such problems are anticipated and resolved automatically by  
 215 the function. For instance, the function resolves loops such as `from = c("hello", "world")`; `to`  
 216 `= c("world", "hello")` automatically, with a warning summarizing the rows that were removed.  
 217 It also resolves loops resulting from poorly specified partial matching rules on the part of the  
 218 user. This is the only tool we are aware of that can help users troubleshoot user-defined rules  
 219 governing node merging and separation.

220 The `textnet_extract` argument of `disambiguate()` accepts the result of the `textnet_extract()`  
 221 function. The object returned by `disambiguate()` updates the `edgelist$sourcecolumn`, `edgelist$tar-`  
 222 `get column`, and `nodelist$entity_name` column to reflect the new node names.

223 Information about the optional argument `try_drop` can be found in the package documentation.  
 224 When specified, the function merges nodes that differ only by the regex phrase specified by  
 225 `try_drop`, and which become identical upon removal of the regular expression encoded in  
 226 `try_drop`.

```
227 tofrom <- data.table::data.table(  
228   from = c(as.list(old_acronyms$acronym),  
229             list("Sub_basin",  
230                 "Sub_Basin",  
231                 "upper_and_lower_aquifers",  
232                 "Upper_and_lower_aquifers",  
233                 "Lower_and_upper_aquifers",  
234                 "lower_and_upper_aquifers")),  
235   to = c(as.list(old_acronyms$name),  
236           list("Subbasin",
```

```

237         "Subbasin",
238         c("upper_aquifer", "lower_aquifer"),
239         c("upper_aquifer", "lower_aquifer"),
240         c("upper_aquifer", "lower_aquifer"),
241         c("upper_aquifer", "lower_aquifer"))))
242
243     old_extract_clean <- disambiguate(
244       textnet_extract = extracts[[1]],
245       from = tofrom$from,
246       to = tofrom$to,
247       match_partial_entity = c(rep(F, nrow(old_acronyms)), T, T, F, F, F, F))
248
249     tofrom <- data.table::data.table(
250       from = c(as.list(new_acronyms$acronym),
251         list("Sub_basin",
252           "Sub_Basin",
253           "upper_and_lower_aquifers",
254           "Upper_and_lower_aquifers",
255           "Lower_and_upper_aquifers",
256           "lower_and_upper_aquifers")),
257       to = c(as.list(new_acronyms$name),
258         list("Subbasin",
259           "Subbasin",
260           c("upper_aquifer", "lower_aquifer"),
261           c("upper_aquifer", "lower_aquifer"),
262           c("upper_aquifer", "lower_aquifer"),
263           c("upper_aquifer", "lower_aquifer"))))
264
265     new_extract_clean <- disambiguate(
266       textnet_extract = extracts[[2]],
267       from = tofrom$from,
268       to = tofrom$to,
269       match_partial_entity = c(rep(F, nrow(new_acronyms)), T, T, F, F, F, F))

```

## 270 Get Network Attributes

271 A tool that generates an igraph or network object from the textnet\_extract output is included  
 272 in the package as the function export\_to\_network(). It returns a list that contains the igraph  
 273 or network itself as the first element, and an attribute table as the second element. Functions  
 274 from the *sna* (Butts 2024), *igraph* (Csárdi et al. 2024), and *network* packages (Butts et  
 275 al. 2023) are invoked to create a network attribute table of common network-level attributes;  
 276 see package documentation for details.

```

277     old_extract_net <- export_to_network(old_extract_clean, "igraph", keep_isolates = F,
278     new_extract_net <- export_to_network(new_extract_clean, "igraph", keep_isolates = F,
279
280     table <- t(format(rbind(old_extract_net[[2]], new_extract_net[[2]]), digits = 3, scie
281     colnames(table) <- c("old", "new")
282     knitr::kable(table)

```

283 old

284 new

285 num\_nodes

286 88

287 118  
288 num\_edges  
289 163  
290 248  
291 connectedness  
292 0.710  
293 0.677  
294 centralization  
295 0.207  
296 0.325  
297 transitivity  
298 0.109  
299 0.153  
300 pct\_entitytype\_homophily  
301 0.503  
302 0.581  
303 reciprocity  
304 0.245  
305 0.306  
306 mean\_in\_degree  
307 1.85  
308 2.10  
309 mean\_out\_degree  
310 1.85  
311 2.10  
312 median\_in\_degree  
313 1  
314 1  
315 median\_out\_degree  
316 1  
317 1  
318 modularity  
319 0.540  
320 0.525  
321 num\_communities  
322 12  
323 16



324 percent\_vbn

325 0.374

326 0.423

327 percent\_vbg

328 0.0736

329 0.0524

330 percent\_vbp

331 0.1288

332 0.0766

333 percent\_vbd

334 0.0675

335 0.0685

336 percent\_vb

337 0.135

338 0.137

339 percent\_vbz

340 0.221

341 0.242

342 The *ggraph* package (Pedersen and RStudio 2024) has been used to create the two network  
343 visualizations seen here, using a weighted version of the *igraphs* constructed below. We set  
344 `collapse_edges = T` to convert the multiplex graph into its weighted equivalent.

345 `library(ggraph)`

346  
347 `## Warning: package 'ggraph' was built under R version 4.3.2`

348  
349 `## Loading required package: ggplot2`

350  
351 `## Warning: package 'ggplot2' was built under R version 4.3.2`

352  
353 `old_extract_plot <- export_to_network(old_extract_clean, "igraph", keep_isolates = F,`  
354 `new_extract_plot <- export_to_network(new_extract_clean, "igraph", keep_isolates = F,`  
355 `#order of these layers matters`

356 `ggraph(old_extract_plot, layout = 'fr')+`

357  `geom_edge_fan(aes(alpha = weight),`

358  `end_cap = circle(1,"mm"),`

359  `color = "#000000",`

360  `width = 0.3,`

361  `arrow = arrow(angle=15,length=unit(0.07,"inches"),ends = "last",type`

362 `#from Paul Tol's bright color scheme`

363 `scale_color_manual(values = c("#4477AA","#228833","#CCBB44","#66CCFF"))+`

364 `geom_node_point(aes(color = entity_type), size = 1,`

365  `alpha = 0.8))+`

366 `labs(title= "Old Network")+`

367 `theme_void()`



Old Network



Figure 1: Figure 1

```

368 #order of these layers matters
369 ggraph(new_extract_plot, layout = 'fr')+
370   geom_edge_fan(aes(alpha = weight),
371                 end_cap = circle(1,"mm"),
372                 color = "#000000",
373                 width = 0.3,
374                 arrow = arrow(angle=15,length=unit(0.07,"inches"),ends = "last",type
375 #from Paul Tol's bright color scheme
376   scale_color_manual(values = c("#4477AA","#228833","#CCBB44","#66CCEE"))+
377   geom_node_point(aes(color = entity_type), size = 1,
378                  alpha = 0.8)+
379   labs(title= "New Network")+
380   theme_void()

```

New Network



Figure 2: Figure 2

### Explore Edge Attributes

The `top_features()` tool calculates the most common verbs across the entire corpus of documents, as shown below.

```
top_feats <- top_features(list(old_extract_net[[1]], new_extract_net[[1]]))
knitr::kable(head(top_feats[[2]],10))
```

names	avg_fract_of_a_doc
be	0.1043934
include	0.0844053
provide	0.0660870
locate	0.0518875
result	0.0406689
base	0.0274342
receive	

401 0.0254181

402 show

403 0.0223506

404 develop

405 0.0212126

406 make

407 0.0203345

408 Using a syntax-based extraction technique enables the preservation of a rich set of edge  
 409 attributes giving insight into the nature of the relationship between each pair of nodes. The  
 410 edge attributes “head\_verb\_name” and “head\_verb\_lemma,” respectively, indicate the verb  
 411 and infinitive form of the verb mediating the relationship between the source and target  
 412 nodes. The edge attributes “helper\_token” and “helper\_lemma” indicate the presence of  
 413 a helping verb in the verb phrase, while the edge attributes “xcomp\_helper\_lemma” and  
 414 “xcomp\_helper\_token” indicate the presence of an open causal complement in the verb phrase.  
 415 Open causal complements, such as “monitor” in the sentence “The agency is expected to  
 416 monitor the results,” can provide key supplemental information about the relationship between  
 417 the source and target nodes. Additional edge attributes include indicators for verb tense and  
 418 the presence of uncertain “hedging” language in the sentence. Other edge attributes travel  
 419 with the edge to document where in the document, and in which document, the edge occurs.  
 420 For instance, we can summarize the verb tense of edges in the original plan in a table. The  
 421 abbreviations follow [Penn Treebank](#) classifications ([Marcus et al., 1999](#)), such that VB = base  
 422 form, VBD = past tense, VBG = gerund or present participle, VBN = past participle, VBP =  
 423 non-3rd person singular present, and VBZ = 3rd person singular present. The most common  
 424 verb tense used in the plan was VBN, or past participle.

425 `knitr::kable(table(igraph::E(old_extract_net[[1]])$head_verb_tense))`

426 Var1

427 Freq

428 VB

429 22

430 VBD

431 11

432 VBG

433 12

434 VBN

435 61

436 VBP

437 21

438 VBZ

439 36

#### 440 Generate Composite Network

441 The combine\_networks function allows a composite network to be generated from multiple  
442 export\_to\_network() outputs. This function is useful for understanding and analyzing the  
443 overlaps between the network of multiple documents. In this example, a composite network  
444 is not as useful, since these documents are not from two different regions being discussed  
445 but rather are two versions of the same document. However, for illustration purposes, the  
446 composite network is generated below.

447 For best results, composite network generation should not be done without an adequate  
448 disambiguation in Step 4. A function is included that merges the edgelists and nodelists of all  
449 documents. If the same node name is mentioned in multiple documents, the node attributes  
450 associated with the highest total number of edges for that node name are preserved.

```
451 composite_net <- combine_networks(list(old_extract_net[[1]], new_extract_net[[1]]), m
452 ggraph(composite_net, layout = 'fr')+
453   geom_edge_fan(aes(alpha = weight),
454                 end_cap = circle(1,"mm"),
455                 color = "#000000",
456                 width = 0.3,
457                 arrow = arrow(angle=15,length=unit(0.07,"inches"),ends = "last",type
458 #from Paul Tol's bright color scheme
459 scale_color_manual(values = c("#4477AA","#228833","#CCBB44","#66CCFF"))+
460 geom_node_point(aes(color = entity_type), size = 1,
461                 alpha = 0.8)+
462 labs(title= "Composite Network")+
463 theme_void()
```

Composite Network



Figure 3: Figure 3

#### 464 Explore Node Attributes

465 The network objects generated from `export_to_network` can be used to analyze the node  
466 attributes of the graphs. Below we demonstrate several node attribute exploration tools. First,  
467 we use the `top_features()` function to calculate the most common entities across the entire  
468 corpus of documents.

```
469 library(network)
470 library(igraph)
471
472 top_feats <- top_features(list(old_extract_net[[1]], new_extract_net[[1]]))
473 print(head(top_feats[[1]],10))
```

```
474
475 ## # A tibble: 10 × 2
476 ##   names                                avg_fract_of_a_doc
477 ##   <chr>                                <dbl>
478 ## 1 groundwater                        0.180
479 ## 2 gsa                                0.0803
480 ## 3 san_joaquin_river                  0.0692
481 ## 4 gfwd_gsa                           0.0452
482 ## 5 surface_water                      0.0426
483 ## 6 gravelly_ford_water_district        0.0386
484 ## 7 subbasin                           0.0381
485 ## 8 gsp                                0.0293
486 ## 9 madera_subbasin                   0.0259
487 ## 10 north_kings_groundwater_sustainability_agency 0.0254
```

488 Next, we calculate node-level attributes on a weighted version of the networks. First we  
489 prepare the data frames for both the old and new networks. We can include the variable  
490 `num_graphs_in` from our composite network to investigate what kinds of nodes are found in  
491 both plans.

```
492 composite_tbl <- igraph::as_data_frame(composite_net, what = "vertices")
493 composite_tbl <- composite_tbl[,c("name","num_graphs_in")]
494
495 #prepare data frame version of old network, to add composite_tbl variables
496 old_tbl <- igraph::as_data_frame(old_extract_net[[1]], what = "both")
497 #this adds the num_graphs_in variable from composite_tbl
498 old_tbl$vertices <- dplyr::left_join(old_tbl$vertices, composite_tbl)
499
500 ## Joining with `by = join_by(name)`
501
502 #turn back into a network
503 old_net <- network::network(x=old_tbl$edges[,1:2], directed = T,
504                             hyper = F, loops = T, multiple = T,
505                             bipartiate = F, vertices = old_tbl$vertices,
506                             matrix.type = "edgelist")
507
508 #we need a matrix version for some node statistics
509 old_mat <- as.matrix(as.matrix(export_to_network(old_extract_clean, "igraph", keep_i
510
511 #prepare data frame version of new network, to add composite_tbl variables
512 new_tbl <- igraph::as_data_frame(new_extract_net[[1]], what = "both")
513 #this adds the num_graphs_in variable from composite_tbl
514 new_tbl$vertices <- dplyr::left_join(new_tbl$vertices, composite_tbl)
515
516 ## Joining with `by = join_by(name)`
```

```

517 #turn back into a network
518 new_net <- network::network(x=new_tbl$edges[,1:2], directed = T,
519                             hyper = F, loops = T, multiple = T,
520                             bipartiate = F, vertices = new_tbl$vertices,
521                             matrix.type = "edgelist")
522 #we need a matrix version for some node statistics
523 new_mat <- as.matrix(as.matrix(export_to_network(new_extract_clean, "igraph", keep_i
524
525 We can now use these data structures to calculate node statistics, as printed below.
526
527 paths2 <- diag(old_mat %*% old_mat)
528 recip <- 2*paths2 / sna::degree(old_net)
529 totalCC <- as.vector(unname(DirectedClustering::ClustF(old_mat, type = "directed", i
530 closens <- sna::closeness(old_net, gmode = "graph", cmode="suminvundir")
531 between <- sna::betweenness(old_net,gmode = "graph",cmode="undirected")
532 deg <- sna::degree(old_net, gmode = "graph", cmode = "undirected")
533 old_node_df <- dplyr::tibble(name = network::get.vertex.attribute(old_net, "vertex.n
534                                closens,
535                                between,
536                                deg,
537                                recip,
538                                totalCC,
539                                entity_type = network::get.vertex.attribute(old_net,"entity_type"
540                                num_graphs_in = network::get.vertex.attribute(old_net, "num_graph
541
542 paths2 <- diag(new_mat %*% new_mat)
543 recip <- 2*paths2 / sna::degree(new_net)
544 totalCC <- as.vector(unname(DirectedClustering::ClustF(new_mat, type = "directed", i
545 closens <- sna::closeness(new_net, gmode = "graph", cmode="suminvundir")
546 between <- sna::betweenness(new_net,gmode = "graph",cmode="undirected")
547 deg <- sna::degree(new_net, gmode = "graph", cmode = "undirected")
548 new_node_df <- dplyr::tibble(name = network::get.vertex.attribute(new_net, "vertex.n
549                                closens,
550                                between,
551                                deg,
552                                recip,
553                                totalCC,
554                                entity_type = network::get.vertex.attribute(new_net,"entity_type"
555                                num_graphs_in = network::get.vertex.attribute(new_net, "num_graph
556
557 summary(old_node_df)
558
559 ##           name           closens           between           deg
560 ## Length:88           Min.    :0.01149           Min.    :    0.00           Min.    : 0.00
561 ## Class :character     1st Qu.:0.25465           1st Qu.:    0.00           1st Qu.: 0.00
562 ## Mode  :character     Median :0.30134           Median :    0.00           Median : 1.00
563 ##                               Mean  :0.26573           Mean  :   62.41           Mean  : 1.67
564 ##                               3rd Qu.:0.32217           3rd Qu.:   19.66           3rd Qu.: 1.00
565 ##                               Max.   :0.51149           Max.   :1191.82           Max.   :19.00
566 ##           recip           totalCC           entity_type           num_graphs_in
567 ## Min.    :0.00000           Min.    :0.000000           Length:88           Min.    :1.000
568 ## 1st Qu.:0.00000           1st Qu.:0.000000           Class :character     1st Qu.:2.000
569 ## Median :0.00000           Median :0.000000           Mode  :character     Median :2.000
570 ## Mean    :0.0518           Mean    :0.080564                               Mean    :1.864
571 ## 3rd Qu.:0.00000           3rd Qu.:0.003472                               3rd Qu.:2.000

```

```
571 ## Max. :1.0000 Max. :1.000000 Max. :2.000
```

```
572
```

```
573 summary(new_node_df)
```

```
574
```

```
575 ##      name      closens      between      deg
576 ## Length:118      Min.   :0.008547      Min.   : 0.000      Min.   : 0.00
577 ## Class :character      1st Qu.:0.232087      1st Qu.: 0.000      1st Qu.: 0.00
578 ## Mode  :character      Median :0.282051      Median : 0.000      Median : 1.00
579 ##      Mean   :0.246142      Mean   : 82.712      Mean   : 1.78
580 ##      3rd Qu.:0.309829      3rd Qu.: 6.022      3rd Qu.: 1.00
581 ##      Max.   :0.512821      Max.   :2025.067      Max.   :32.00
582 ##      recip      totalCC      entity_type      num_graphs_in
583 ## Min.   :0.000000      Min.   :0.000000      Length:118      Min.   :1.000
584 ## 1st Qu.:0.000000      1st Qu.:0.000000      Class :character      1st Qu.:1.000
585 ## Median :0.000000      Median :0.000000      Mode  :character      Median :2.000
586 ## Mean   :0.04173      Mean   :0.11473      Mean   :1.644
587 ## 3rd Qu.:0.000000      3rd Qu.:0.08808      3rd Qu.:2.000
588 ## Max.   :1.000000      Max.   :1.000000      Max.   :2.000
```

589 The 2x2 table below summarizes the rate at which each entity type is found in both plans.  
 590 Very few nodes in the old version (12 out of 88) are absent from the new version. Conversely,  
 591 a substantial minority of nodes in the new version (42 out of 118) are absent from the old  
 592 version.

```
593 old_node_df$plan_version <- "old"
594 new_node_df$plan_version <- "new"
595 combineddf <- rbind(old_node_df, new_node_df)
596 with(combineddf, table(plan_version, num_graphs_in))
```

```
597
```

```
598 ##      num_graphs_in
599 ## plan_version  1  2
600 ##      new 42 76
601 ##      old 12 76
```

602 We can also investigate differences in network statistics between the two plans. For instance,  
 603 the distribution of degree does not change much between plan versions. The distribution of  
 604 betweenness, likewise, is relatively stable except for person nodes, which are the least common  
 605 nodes in the graph.

```
606 library(gridExtra)
607 library(ggplot2)
608 b1 <- ggplot(old_node_df, aes(x = entity_type, y = deg)) + geom_boxplot() + theme_bw()
609 b2 <- ggplot(new_node_df, aes(x = entity_type, y = deg)) + geom_boxplot() + theme_bw()
610 b3 <- ggplot(old_node_df, aes(x = entity_type, y = log(between+0.01))) + geom_boxplot()
611 b4 <- ggplot(new_node_df, aes(x = entity_type, y = log(between+0.01))) + geom_boxplot()
612
613 grid.arrange(b1, b2, b3, b4, ncol=2)
```



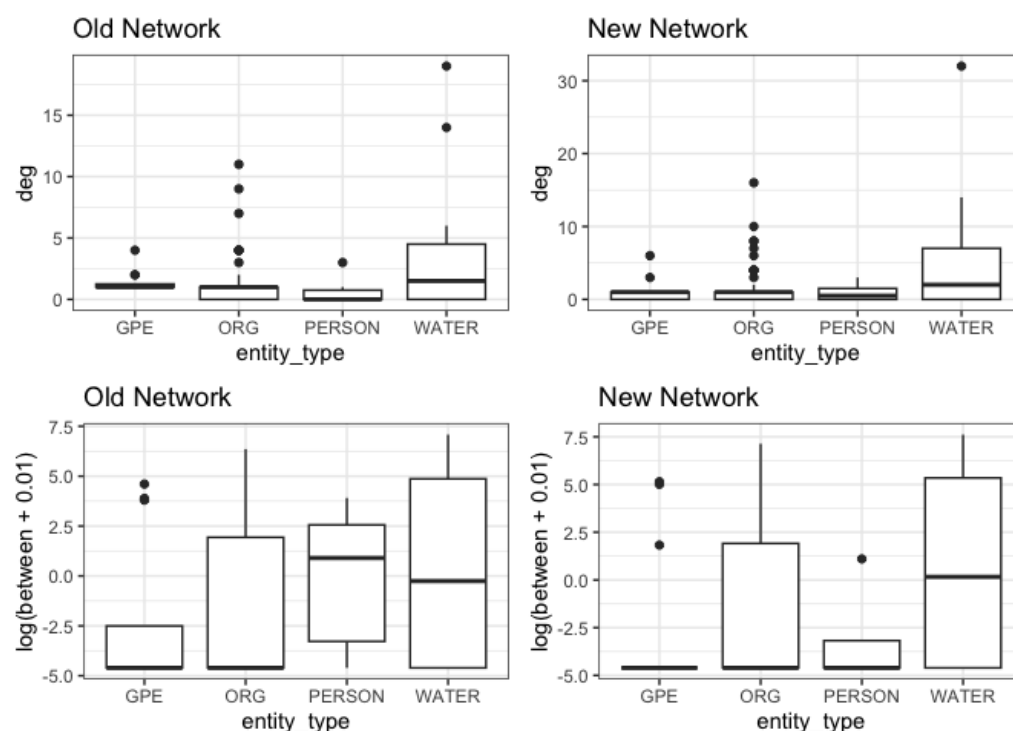


Figure 4: Figure 4

## Potential Further Analyses

The network-level attributes output from `export_to_network` can also be analyzed against exogenous metadata that has been collected separately by the researcher regarding the different documents and their real-world context. The extracted networks, with their collections of verb attributes, node attributes, edge incidences, and edge attributes, can also be analyzed through a variety of tools, such as an Exponential Random Graph Model, to determine the probability of edge formation under certain conditions. A Temporal Exponential Random Graph Model could also shed light on the changes of a document over time, such as the multiple versions of the groundwater sustainability plan in this example.

## Entity Network Extraction Algorithm

The directed network generated by `textNet` represents the collection of all identified entities in the document, joined by edges signifying the verbs that connect them. The user can specify which entity categories should be preserved. The output format is a list containing four data.tables: an edgelist, a nodelist, a verblist, and an appositive list.

The edgelist includes edge attributes such as verb tense, any auxiliary verbs in the verb phrase, whether an open clausal complement (Universal Dependencies code “xcomp”) is associated with the primary verb, whether any hedging words were detected in the sentence, and whether any negations were detected in the sentence.

The returned edgelist by default contains both complete and incomplete edges. A complete edge includes a source, verb, and target. An incomplete edge includes either a source or a target, but not both, along with its associated verb. Incomplete edges convey information about which entities are commonly associated with different verbs, even though they do not reveal information about which other entities they are linked to in the network. These incomplete edges can be filtered out when converting the output into a network object, such as through the `network` package or the `igraph` package. The nodelist returns all entities of the

desired types found in the document, regardless of whether they were found in the edgelist. Thus, the nodelist allows the presence of isolates to be documented, as well as preserving node attributes. The verblist includes all of the verbs found in the document, along with verb attributes imported from *VerbNet* (Kipper-Schuler, 2006). This can be used to conduct analyses of certain verb classifications of interest. Finally, the appositive list is a table of entities that may be synonyms. This list is generated from entities whose universal dependency parsing labels as appositives, and whose head token points to another entity. These pairs are included in the table as potential synonyms. If this feature is used, cleaning and filtering by hand is recommended, as appositives can at times be misidentified by existing NLP tools. An automated alternative we recommend is our `find_acronym` tool, which scans the entire document for acronyms defined parenthetically in-text and compiles them in a table.

This network is directed such that the entities that form the subject of the sentence are denoted as the “source” nodes, and the remaining entities are denoted as the “target” nodes. To identify whether each entity is a “source” or a “target”, we use dependency parsing in the Universal Dependencies format, in which each token in a given sentence has an associated “syntactic head” token from which it is derived. Starting with each entity in the sentence, the chain of syntactic head tokens is traced back until either a subject or a verb is reached. If it reaches a subject first, the entity is considered a “source.” If it reaches a verb first, it is considered a “target.”

To identify the subject, we search for the presence of at least one of the following subject tags: “nsubj” (nominal subject), “nsubjpass” (nominal subject – passive), “csubj” (clausal subject), “csubjpass” (clausal subject – passive), “agent”, and “expl” (expletive). To identify the object, we search for the presence of at least one of the following: “pobj” (object of preposition), “iobj” (indirect object), “dative”, “attr” (attribute), “dobj” (direct object), “oprd” (object predicate), “ccomp” (clausal complement), “xcomp” (open clausal complement), “acomp” (adjectival complement), or “pcomp” (complement of preposition).

If a subject token is reached first (“nsubj,” “nsubjpass,” “csubj,” “csubjpass,” “agent,” or “expl”), this indicates that the original token is doing the verb action. That is, it serves some function related to the subject of the sentence. We designate this by tagging it “source,” since these types of relationships will be used to designate the “from” or “source” nodes in our directed network. If a verb token is reached first (“VERB” or “AUX”), this indicates that the verb action is occurring for or towards the original token, which we denote with the tag “target.” These tokens are potential “to” or “target” nodes in our directed network. Linking the two nodes is an edge representing the verb that connects them in the sentence.

Due to the presence of tables, lists, or other anomalies in the original document, it is possible that a supposed “sentence” has a head token trail that does not lead to a verb as is normatively the case. In these instances, the tokens whose trails terminate with a non-subject, non-verb token are assigned neither “source” nor “target” tags. Finally, an exception is made if an appositive token is reached first, since this indicates that the token in question is merely a synonym or restatement of an entity that is already described elsewhere in the sentence and, accordingly, should not be treated as a separate node. Tokens that lead to appositives are assigned neither “source” nor “target” tags, but are preserved as a separate appositive list.

If a verb phrase in the edgelist does not have any sources, the sources associated with the head token of the verb phrase’s main verb (that is, the verb phrase’s parent verb) are adopted as sources of that verb phrase. As of Version 1.0, *textNet* does not do this recursively, to preserve performance optimization.

The `textNet::textnet_extract()` function returns the full list of open clausal complement lemmas associated with the main verb as an edge attribute: “xcomp\_verb”. The list of auxiliary verbs and their corresponding lemmas associated with the main verb, as well as the list of auxiliary verbs and corresponding lemmas associated with the open clausal complements linked to the main verb, are also included as edge attributes: “helper\_token”, “helper\_lemma”, “xcomp\_helper\_token”, and “xcomp\_helper\_lemma”, respectively.

691 The extraction function also detects hedging words and negations. The function  
692 `textNet::textnet_extract()` produces an edge attribute `"has_hedge"`, which is T if there  
693 is a hedging auxiliary verb (`"may"`, `"might"`, `"can"`, `"could"`) or main verb (`"seem"`, `"ap-`  
694 `pear"`, `"suggest"`, `"tend"`, `"assume"`, `"indicate"`, `"estimate"`, `"doubt"`, `"believe"`) in the verb  
695 phrase.

696 Tense is also detected. The six tenses tagged by *spaCy* in `textNet::parse_text()` are preserved  
697 by `textNet::textnet_extract()` as an edge attribute `"head_verb_tense"`. This attribute can  
698 take on one of six values: `"VB"` (verb, base form), `"VBD"` (verb, past tense), `"VBG"` (verb,  
699 gerund or present participle), `"VBN"` (verb, past participle), `"VBP"` (verb, non-3rd person  
700 singular present), or `"VBZ"` (verb, 3rd person singular present). Additionally, an edge attribute  
701 `"is_future"` is generated by `textNet::textnet_extract()`, which is T if the verb phrase contains an  
702 `xcomp`, has the token `"going"` as a head verb, and a being verb token as an auxiliary verb (i.e. is  
703 of the form `"going to "`) or contains one of the following auxiliary verbs: `"shall"`, `"will"`, `"wo"`,  
704 or `"ll"` (i.e. is of the form `"will "`).

## 705 Acknowledgements

706 The authors gratefully acknowledge the support of the Sustainable Agricultural Systems  
707 program, project award no. 2021-68012-35914, from the U.S. Department of Agriculture's  
708 National Institute of Food and Agriculture and the National Science Foundation's Dynamics  
709 of Integrated Socio-Environmental Systems program, grant no. 2205239.

## 710 Appendix

711 This appendix describes the pre-processing tools available through the *textNet* package, which  
712 enable the user to generate the data frame expected by the `textnet_extract()` function.

### 713 Pre-Processing Step I: Process PDFs

714 This is a wrapper for *pdftools*, which has the option of using `pdf_text` or OCR. We have also  
715 added an optional header/footer removal tool. This optional tool is solely based on carriage  
716 returns in the first or last few lines of the document, so may inadvertently remove portions  
717 of paragraphs. However, not removing headers or footers can lead to improper inclusion of  
718 header and footer material in sentences, artificially inflating the presence of nodes whose entity  
719 names are included in the header and footer. Because of the risk of headers and footers to  
720 preferentially inflate the presence of a few nodes, the header/footer remover is included by  
721 default. It can be turned off if the user has a preferred header/footer removal tool to use  
722 instead, or if the input documents lack headers and footers.

```
723 library(textNet)
724 library(stringr)
725 URL <- "https://sgma.water.ca.gov/portal/service/gspdocument/download/2840"
726 download.file(URL, destfile = "old.pdf", method="curl")
727
728 URL <- "https://sgma.water.ca.gov/portal/service/gspdocument/download/9625"
729 download.file(URL, destfile = "new.pdf", method="curl")
730
731 pdfs <- c("old.pdf",
732           "new.pdf")
733
734 old_new_text <- textNet::pdf_clean(pdfs, keep_pages=T, ocr=F, maxchar=10000,
735                                   export_paths=NULL, return_to_memory=T, suppressWarn = F, auto_headf
736 names(old_new_text) <- c("old", "new")
```

## 737 Pre-Processing Step II: Parse Text

This is a wrapper for the pre-trained multipurpose NLP model *spaCy* (Honnicbal et al., 2021), which we access through the R package *spacyr* (Benoit et al., 2023). It produces a table that can be fed into the `textnet_extract` function in the following step. To initialize the session, the user must define the “RETICULATE\_PYTHON” path, abbreviated as “ret\_path” in *textNet*, as demonstrated in the example below. The page contents processed in the Step 1 must now be specified in vector form in the “pages” argument. To determine which file each page belongs to, the user must specify the `file_ids` of each page. We have demonstrated how to do this below. The package by default does not preserve hyphenated terms, but rather treats them as separate tokens. This can be adjusted.

The user may also specify “phrases\_to\_concatenate”, an argument representing a set of phrases for spaCy to keep together during its parsing. The example below demonstrates how to use this feature to supplement the NER capabilities of spaCy with a custom list of entities. This supplementation could be used to ensure that specific known entities are recognized; for instance, spaCy might not detect that a consulting firm such as “Schmidt and Associates” is one entity rather than two. Conversely, this capability could be leveraged to create a new category of entities to detect, that a pretrained model is not designed to specifically recognize. For instance, to create a public health network, one might include a known list of contaminants and diseases and designate custom entity type tags for them, such as “CONTAM” and “DISEASE”). In this example, we investigate the connections between the organizations, people, and geopolitical entities discussed in the plan with the flow of water in the basin. To assist with this, we have input a custom list of known water bodies in the region governed by our test document and have given it the entity designation “WATER”. This is carried out by setting the variable “phrases\_to\_concatenate” to a character vector, including all of the custom entities. Then, the entity type can be set to the desired category. Note that this function is case-sensitive.

```

763 library(findpython)
764 ret_path <- find_python_cmd(required_modules = c('spacy', 'en_core_web_lg'))
765
766
767 water_bodies <- c("surface water", "Surface water", "groundwater", "Groundwater", "Sa
768 Chowchilla canal", "lower aquifer", "upper aquifer", "upper and lower aquifers", "lower
769
770 old_new_parsed <- textNet::parse_text(ret_path,
771                                     keep_hyph_together = F,
772                                     phrases_to_concatenate = water_bodies,
773                                     concatenator = "_",
774                                     text_list = old_new_text,
775                                     parsed_filenames=c("old_parsed","new_parsed"),
776                                     overwrite = T,
777                                     custom_entities = list(WATER = water_bodies))

```

Another NLP tool may be used instead of the built-in *textNet* function at this phase, as long as the output conforms to spaCy tagging standards: Universal Dependencies tags for the “pos” part-of-speech column (Nivre, 2017), and Penn Treebank tags for the “tags” column (Marcus et al., 1999). The `textnet_extract` function expects the parsed table to follow specific conventions. First, a row must be included for each token. The column names expected by `textnet_extract` are:

- 784   ▪ doc\_id, a unique ID for each page
- 785   ▪ sentence\_id, a unique ID for each sentence
- 786   ▪ token\_id, a unique ID for each token
- 787   ▪ token, the token, generally a word, represented as a string
- 788   ▪ lemma, the canonical or dictionary form of the token

- 789     ▪ pos, a code referring to the token's part of speech, defined according to [Universal](#)
- 790     [Dependencies](#) (Nivre, 2017).
- 791     ▪ tag, a code referring to the token's part of speech, according to [Penn Treebank](#) (Marcus
- 792     [et al., 1999](#)).
- 793     ▪ head\_token\_id, a numeric ID referring to the token\_id of the head token of the current
- 794     row's token
- 795     ▪ dep\_rel, the dependency label according to [ClearNLP Dependency](#) labels (Choi, 2024)
- 796     ▪ entity, the entity type category defined by [OntoNotes 5.0](#) (Weischedel et al., 2012). This
- 797     is represented as a string, ending in “\_B” if it is the first token in the entity or “\_I”
- 798     otherwise).

## 799     References

- 800     Bail, C. (2024). *Cbail/textnets* (Version 0.1.1). <https://github.com/cbail/textnets>
- 801     Barrett, T., Dowle, M., Srinivasan, A., Gorecki, J., Chirico, M., Hocking, T., Schwendinger, B.,
- 802     Stetsenko, P., Short, T., Lianoglou, S., Antonyan, E., Bonsch, M., Parsonage, H., Ritchie,
- 803     S., Ren, K., Tan, X., Saporta, R., Seiskari, O., Dong, X., ... Krylov, I. (2024). *Data.table:*
- 804     *Extension of 'data.frame'* (Version 1.16.2). [https://cran.r-project.org/web/packages/data.](https://cran.r-project.org/web/packages/data.table/index.html)
- 805     [table/index.html](https://cran.r-project.org/web/packages/data.table/index.html)
- 806     Benoit, K., Matsuo, A., Gruber, J., & Council (ERC-2011-StG 283794-QUANTESS), E. R.
- 807     (2023). *Spacyr: Wrapper to the 'spaCy' 'NLP' library* (Version 1.3.0). [https://cran.](https://cran.r-project.org/web/packages/spacyr/index.html)
- 808     [r-project.org/web/packages/spacyr/index.html](https://cran.r-project.org/web/packages/spacyr/index.html)
- 809     Butts, C. T., Hunter, D., Handcock, M., Bender-deMoll, S., Horner, J., Wang, L., Krivitsky, P.
- 810     N., Knapp, B., Bojanowski, M., & Klumb, C. (2023). *Network: Classes for relational data*
- 811     (Version 1.18.2). <https://cran.r-project.org/web/packages/network/index.html>
- 812     Choi, J. (2024, August 7). *ClearNLP dependency labels. ClearNLP guidelines.* [https://github.](https://github.com/clir/clearnlp-guidelines/blob/master/md/specifications/dependency_labels.md)
- 813     [com/clir/clearnlp-guidelines/blob/master/md/specifications/dependency\\_labels.md](https://github.com/clir/clearnlp-guidelines/blob/master/md/specifications/dependency_labels.md)
- 814     Csárdi, G., Nepusz, T., Traag, V., Horvát, S., Zanini, F., Noom, D., Müller, K., Salmon, M.,
- 815     Antonov, M., & details, C. Z. I. igraph author. (2024). *Igraph: Network analysis and*
- 816     *visualization* (Version 2.1.1). <https://cran.r-project.org/web/packages/igraph/index.html>
- 817     Hollway, J. (2024). *Manynet: Many ways to make, modify, map, mark, and measure myriad*
- 818     *networks* (Version 1.2.6). <https://CRAN.R-project.org/package=manynet>
- 819     Honnibal, M., Montani, I., Van Landeghem, S., & Boyd, A. (2021). *spaCy: Industrial-strength*
- 820     *natural language processing in python* (Version 3.1.3). [https://github.com/explosion/](https://github.com/explosion/spaCy/tree/master)
- 821     [spaCy/tree/master](https://github.com/explosion/spaCy/tree/master)
- 822     Kipper-Schuler, K. (2006). *VerbNet* (Version 3.3). University of Colorado Boulder. [https:](https://verbs.colorado.edu/verb-index/vn3.3/index.php)
- 823     [//verbs.colorado.edu/verb-index/vn3.3/index.php](https://verbs.colorado.edu/verb-index/vn3.3/index.php)
- 824     Levallois, C., Clithero, J. A., Wouters, P., Smidts, A., & Huettel, S. A. (2012). Translating
- 825     upwards: Linking the neural and social sciences via neuroeconomics. *Nature Reviews*
- 826     *Neuroscience*, 13(11), 789–797. [https://nocodefunctions.com/cowo/semantic\\_networks\\_](https://nocodefunctions.com/cowo/semantic_networks_tool.html)
- 827     [tool.html](https://nocodefunctions.com/cowo/semantic_networks_tool.html)
- 828     Marcus, M., Santorini, B., Marcinkiewicz, M. A., & Taylor, A. (1999). *Treebank-3 documenta-*
- 829     *tion*. <https://doi.org/https://doi.org/10.35111/gq1x-j780>
- 830     Nivre, J. (2017). *Universal POS tags* (Version 2). Universal Dependencies. [https://](https://universaldependencies.org/u/pos/)
- 831     [universaldependencies.org/u/pos/](https://universaldependencies.org/u/pos/)
- 832     Ooms, J. (2024). *Pdftools: Text extraction, rendering and converting of PDF documents*
- 833     (Version 3.4.1). <https://cran.r-project.org/web/packages/pdftools/index.html>
- 834     Paranyushkin, D. (2018). *InfraNodus*. Nodus Labs. <https://infranodus.com/>

- 835 Solymos, P., Zawadzki, Z., Bengtsson, H., & Team, R. C. (2023). *Pbapply: Adding progress bar*  
836 *to '\*apply' functions* (Version 1.7-2). [https://cran.rstudio.com/web/packages/pbapply/](https://cran.rstudio.com/web/packages/pbapply/index.html)  
837 [index.html](https://cran.rstudio.com/web/packages/pbapply/index.html)
- 838 Weischedel, R., Pradhan, S., Ramshaw, L., Kaufman, J., Franchini, M., El-Bachouti, M.,  
839 Xue, N., Palmer, M., Hwang, J. D., Bonial, C., Choi, J., Mansouri, A., Foster, M.,  
840 Hawwary, A., Marcus, M., Taylor, A., Greenberg, C., Hovy, E., Belvin, R., & Houston, A.  
841 (2012). *OntoNotes release 5.0 with OntoNotes DB tool v0.999 beta*. BBN Technologies.  
842 <https://catalog.ldc.upenn.edu/docs/LDC2013T19/OntoNotes-Release-5.0.pdf>
- 843 Wickham, H., François, R., Henry, L., Müller, K., Vaughan, D., Software, P., & PBC. (2023).  
844 *Dplyr: A grammar of data manipulation* (Version 1.1.4). [https://cran.r-project.org/web/](https://cran.r-project.org/web/packages/dplyr/index.html)  
845 [packages/dplyr/index.html](https://cran.r-project.org/web/packages/dplyr/index.html)
- 846 Wickham, H., Hester, J., Chang, W., Bryan, J., & RStudio. (2022). *Devtools: Tools to make*  
847 *developing r packages easier* (Version 2.4.5). [https://cran.r-project.org/web/packages/](https://cran.r-project.org/web/packages/devtools/index.html)  
848 [devtools/index.html](https://cran.r-project.org/web/packages/devtools/index.html)
- 849 Wickham, H., Vaughan, D., Girlich, M., Ushey, K., Software, P., & PBC. (2024). *Tidyr: Tidy*  
850 *messy data* (Version 1.3.1). <https://cran.r-project.org/web/packages/tidyr/index.html>

DRAFT