

Relevance Feedback Search Based on Automatic Annotation and Classification of Texts

Rafael Leal ✉ 

HELDIG – Helsinki Centre for Digital Humanities, University of Helsinki, Finland

Joonas Kesäniemi ✉ 

Semantic Computing Research Group (SeCo), Aalto University, Finland

Mikko Koho ✉ 

HELDIG – Helsinki Centre for Digital Humanities, University of Helsinki, Finland

Semantic Computing Research Group (SeCo), Aalto University, Finland

Eero Hyvönen ✉ 

Semantic Computing Research Group (SeCo), Aalto University, Finland

HELDIG – Helsinki Centre for Digital Humanities, University of Helsinki, Finland

Abstract

The idea behind Relevance Feedback Search (RFBS) is to build the search query as an iterative process where it is gradually refined or redirected interactively, based on the search results of the previous round. This can be helpful in situations where the end user cannot easily formulate their information needs at the outset as a well-focused query, or more generally as a way to filter and focus search results. This paper concerns the application of this search paradigm using textual documents in the legal domain. The aim is to integrate keyword based on word extraction and unsupervised classification into a RFBS framework and apply it to a use case of exploring legal documents. We focus on the Natural Language Processing (NLP) methods underlying the framework and application, where an automatic annotation tool is used for extracting document keywords as ontology concepts, which are then transformed into word embeddings to form vectorial representations of the texts. An unsupervised classification system that uses similar techniques to classify the documents into broad thematic classes is also presented. This classification functionality is evaluated using two different datasets. As a use case for this framework, an application perspective in the semantic portal *LawSampo – Finnish Legislation and Case Law on the Semantic Web* is described. This online demonstrator uses a dataset of 82 145 sections in 3725 statutes of Finnish legislation and a dataset of 13 470 court decisions.

2012 ACM Subject Classification [Computing methodologies Information extraction](#); [Applied computing Document searching](#); [Information systems Clustering and classification](#)


Keywords and phrases relevance feedback, keyword extraction, zero-shot text classification, word embeddings, LawSampo

Digital Object Identifier [10.4230/OASICS...](#)

1 Introduction

In many search situations, the information need of the user cannot be formulated precisely. The search query must then be gradually refined and the results re-evaluated in a series of successive rounds in order to achieve a satisfactory outcome. This paper describes language technology methods that can be used in the application of this kind of iterative and interactive search, the Relevance Feedback Search (RFBS) paradigm [1, Ch. 5], to searching and exploring textual documents. We outline a search system, based on pre-trained models and algorithms, that integrates keyword extraction and unsupervised categorization of documents into RFBS. We also present a case study with an implementation of this framework to the legal domain as part of the *LawSampo – Finnish Legislation and Case Law*

 © Rafael Leal, Joonas Kesäniemi, Mikko Koho and Eero Hyvönen;
licensed under Creative Commons License CC-BY 4.0

 OpenAccess Series in Informatics
Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

45 on the *Semantic Web*¹ [5] system.

■ **Algorithm 1** Proposed relevance feedback search

Result: Result set RS in documents D satisfying the user's information need

```

/* Initialization */
1 AQ := ε ; // Active free-form text query (initially the empty string)
2 AK := ∅ ; // Active Keywords set (initially empty)
3 AC := ∅ ; // Active Categories set (initially empty)
4 RS := {d0, ..., dn} ; // Result Set (initially contains the whole document
  domain)
46 /* RFBS loop */
5 while RS is not a satisfying result do
6   SK := KeywordsOf(RS) ; // Suggested keywords based on RS
7   SC := CategoriesOf(RS) ; // Suggested categories based on RS
8   AQ := ModifyQ(AQ) ; // User optionally modifies AQ
9   AK := ModifyK(AK, KW) ; // User optionally modifies AK based on SK
10  AC := ModifyC(AC, SC) ; // User optionally modifies AC based on SC
11  RS := Search(AQ, AK, AC);
12 end

```

47 The proposed RFBS process is outlined in Algorithm 1. The function *Search*(*AQ*, *AK*, *AC*)
48 in line 11 is used for searching the documents based on the search query *AQ*, active keywords
49 *AK*, and active categories *AC*. The idea is that the *categories* are used for setting larger
50 thematic contexts for search that can then be refined using keywords. For example, categories
51 may refer to different phases or situations during the life of the end users, such as childhood
52 or golden age [17], or societal contexts, such as health or environmental issues. In the While
53 loop, the system computes in each iteration a set of new categories and keywords based on
54 the search results RS as suggestions for the user to consider in the next round of searching.
55 The functions *ModifyQ*, *ModifyK* and *ModifyC* in lines 8-10 allow the user to make optional
56 modifications for the next search. In this way, the process is expected to converge gradually
57 towards more and more satisfying results in RS.

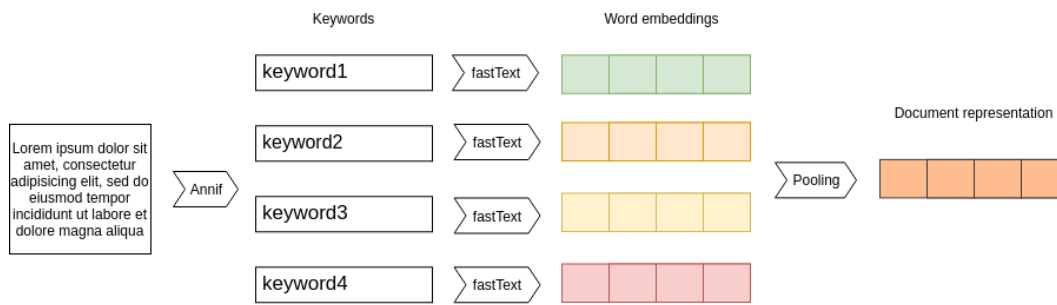
58 A novel idea in the proposed approach is to combine implicit and explicit feedback
59 methods [13] by using topical classification of the documents and keyword concepts extracted
60 from the search results. User feedback on the topics and keywords is used to generate new
61 search queries, thus guiding the iterative search process.

62 In order to avoid the challenges of traditional text-based search, for example morphological
63 variation of words in highly inflectional languages such as Finnish, and semantic difficulties
64 such as synonymy and polysemy, the presented system works on a semantic level. This is
65 done based on ontological concepts and themes extracted automatically from the texts, which
66 are processed via a word embedding algorithm. A zero-shot classification sub-system based
67 on the same approach is responsible for categorizing the documents.

68 The implementation of the system uses Finnish legal documents, in particular sections of
69 the law and court decisions, with the aim to help users to find jurisprudence related to their
70 problematic situations in life. For the data, we use documents from the Semantic Finlex
71 data service² [10], refined further for the LawSampo system and data service [5]. This search

¹Project homepage: <https://seco.cs.aalto.fi/projects/lakisampo/en/>

²<https://data.finlex.fi>



■ **Figure 1** Simplified visual overview of the document representation building algorithm

72 prototype constitutes one application perspective³ of the LawSampo system. It is designed
 73 to work with Finnish language content, but the methods presented here can be applied to
 74 documents in any language, if similar ontologies and language processing tools are available.

75 In the next sections, the methods and software needed for knowledge extraction of
 76 thematic categories and keywords from the texts are described, alongside some evaluation
 77 results. Subsequently, the data underlying our case study is explained, as well as the RFBS
 78 application in LawSampo. Finally, contributions of our work and related works are discussed,
 79 and directions for further research suggested.

80 2 Framework for Knowledge Extraction and Search

81 All documents in the target dataset are semantically annotated in two ways: 1) by extracting
 82 keyword concepts based on a keyword ontology and 2) by classifying the documents into a set
 83 of larger thematic categories. The search framework relies on three main components: how
 84 the system represents the documents internally, how they are classified in different categories,
 85 and how the search is performed. These will be explained in more detail in the following
 86 subsections.

87 2.1 Document representation

88 The core of this system lies in its representations of the documents, which are built using
 89 a three-step approach: 1) representative keywords are extracted, then 2) transformed into
 90 their respective word embeddings, and finally 3) combined via mean pooling to form the
 91 document representation vector. This process is visualized in Figure 1.

92 Keyword extraction is performed via Annif⁴ [20], a subject indexing tool developed by the
 93 National Library of Finland. It is capable of using different algorithms to return suggested
 94 keywords and their respective weights for a given input text. Annif developers also provide a
 95 REST API containing various projects that combine different algorithms. They are all trained
 96 on bibliographical metadata from the works found on the Finna portal⁵ which publishes
 97 information about objects in Finnish archives, museums and libraries, in a vein similar to
 98 Europeana. Since the training data is labelled with terms from the General Finnish Ontology
 99 (YSO)⁶ [16], the API returns keywords identified by unique YSO URIs. This also means that

³In Sampo portals, the different ways for accessing the data are called “perspectives”.

⁴<https://annif.org/>

⁵<https://www.finna.fi>

⁶<https://finto.fi/ys0/en/>

100 the keywords are already lemmatized.

101 Our present system uses the *yso-fi* pre-trained project on the Annif REST API. It
 102 combines, with equal weights, the results of three algorithms: the venerable TF-IDF, which
 103 uses word counts to determine the most representative words in each document; Maui, a
 104 topic indexing tool; and Parabel, an extreme multi-label classification algorithm capable of
 105 handling millions of categories. The first two directly match salient terms to a vocabulary,
 106 while the last is also able to find indirect correlations between words [20]. This mix provides
 107 results that are not only grounded on the text of the documents but also able to extrapolate
 108 their specific wording. The settings we use are a limit of 100 and a threshold of 0.01, which
 109 offer a good balance between rigour and broadness.

110 The keyword weights obtained from Annif are saved into the *keyword matrix*, a document-
 111 term structure. It can be further transformed via the TF-IDF algorithm in order to reduce
 112 the weight of keywords repeated throughout the dataset, since terms that occur in fewer
 113 documents are in general semantically more representative [14].

114 Once the texts are distilled into a set of representative keywords, our system uses fastText
 115 [3] to obtain word embeddings for each of them. This approach improves the task of
 116 unsupervised categorization by avoiding creating representations at the sentence level [7].
 117 FastText improves on the word2vec embedding algorithm by breaking up the words into a bag
 118 of n-grams, each with their own vectorial representation. When building the final embedding,
 119 these components are also taken into account. This mechanism provides the means to build
 120 representations for words that are not in the training data by breaking them up into smaller,
 121 already-seen chunks. This is especially important for morphologically rich languages such as
 122 Finnish, which present unseen word forms more often than analytic languages.

123 Representations are calculated via the `get_sentence_vector` function from the fastText
 124 Python package, which averages the l2 norm for each word representation in the sentence.
 125 The system uses a pre-trained language model offered by the fastText developers [8]. Since it
 126 is extensively trained, this model is able to represent more accurately words both present
 127 in and absent from the dataset—the latter is essential considering that Annif is capable of
 128 suggesting latent keywords.

129 The final representation of each document is obtained via mean pooling of their keyword
 130 embeddings. All document representations are then collected into a $D^{d \times 300}$ *document matrix*,
 131 where d is the total number of documents and 300 is the number of dimensions in the
 132 pre-trained fastText model.

133 2.2 Document Classification

134 A similar process is used to build a $C^{c \times 300}$ category matrix, where c is the total number
 135 of category labels. Although the keyword extraction algorithms are trained on categorical
 136 metadata, as explained in Section 2.1, the classification mechanism that employs it can be
 137 considered unsupervised, since it is built without any training.

138 Category labels must be provided by the end user, since at this point the system does not
 139 recognise broad topics automatically. The labels are pre-processed by stripping commas and
 140 conjunctions and used as input for the Annif→fastText document representation pipeline of
 141 Figure 1.

142 However, Annif’s results can be somewhat erratic when it comes to very short texts,
 143 and the category labels often contain a single word. Filtering out unrelated terms from
 144 the set of keyword suggestions can thus produce better results. This is done via cosine
 145 similarity: a matrix containing word embeddings for the entire category label text as well as
 146 its individual words is built, and then compared to another matrix containing *semantically*

147 *reinforced* embeddings for all the keywords suggested by Annif (this reinforcement technique
 148 is described below). If the maximum cosine similarity between a suggestion and any of
 149 the query vectors is under a certain threshold (0.25 is the default), the keyword is rejected.
 150 The word embeddings for the original label and the resulting keywords suggestions, pooled
 151 together, form the category representation vector. A diagram representing this process can
 152 be seen in Figure 2.

153 Since Annif returns YSO concepts, their vectorial representation can be semantically
 154 reinforced. This simple technique is based on ontology relations: main and alternative labels
 155 (*prefLabel* and *altLabel*), exactly matching (*exactMatch*) and closely matching (*closeMatch*)
 156 entities linked to the target concept are fetched; the final representation is calculated as the
 157 average embedding of this expanded set of entities.

158 As an example: the concept *shares*, whose URI is <http://www.yso.fi/onto/koko/p12994>,
 159 returns a wealth of associated concepts: broader, narrower, related, exactly matching and
 160 closely matching concepts, as well as sections named "in other languages" and "entry terms",
 161 which are alternative labels for the concept. Out of those, the framework chooses only the
 162 following:

- 163 ■ ***prefLabel*** (the main label): *shares*
- 164 ■ ***altLabel*** (entry terms): *share*, *stocks*
- 165 ■ ***closeMatch*** (closely matching categories): *Stocks* (linked to the Library of Congress
 166 concept)
- 167 ■ ***exactMatch*** (exactly matching categories): *share*, *shares* (linked to other Finnish
 168 ontologies)

169 The associated terms are then collected in a set $S = \{shares, share, stock\}$ and the final
 170 representation will be an average of the word embeddings for the terms in S . In this specific
 171 case, the procedure helps disambiguate this concept from another *shares* entry⁷, which is a
 172 technical objects also known as *drill coulters* or *ploughshares*. However, this technique can
 173 help build better representations even without disambiguation, by reinforcing the category
 174 labels with the labels of semantically similar concepts.

175 The dot product between the category matrix and the document matrix builds a $V^{c \times d}$
 176 *category-document matrix*, which stores the strength of each category-document pair.

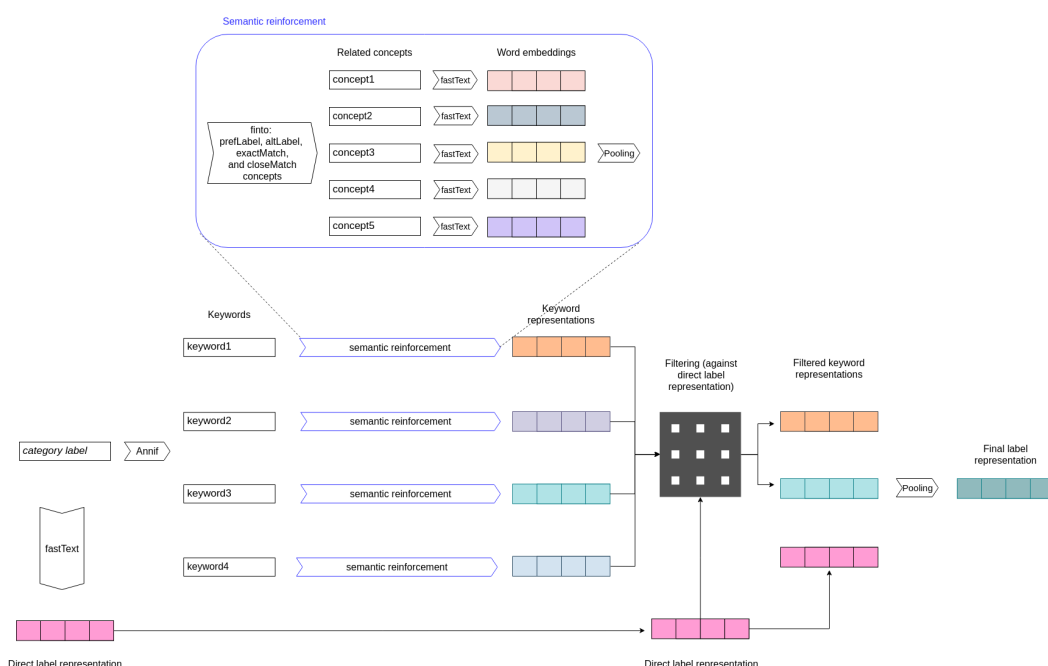
177 2.3 Search

178 The results obtained by the search system of our RFBS framework depends on four elements:
 179 text query, category, positive keywords, and negative keywords. They may work separately
 180 or in tandem. This algorithm corresponds to the *Search()* method in line 11 of Algorithm 1,
 181 where both positive and negative keyword sets are represented by AK .

182 The free-form search query, in case it is given, is the first element to be resolved. It is
 183 transformed into embeddings in the standard way, via `Annif→fastText`. The results can be
 184 filtered as explained in section 2.2, however without semantically reinforcement in order to
 185 save computational time. The most similar documents are calculated via cosine similarity
 186 between the query representation and the document matrix, and the results are collected in
 187 a *result list*.

188 Positive keywords and categories are executed next, in case they are set. First, scores
 189 for each document are calculated as the sum of their respective weights (from the keyword

⁷<http://www.yso.fi/onto/koko/p48662>



■ **Figure 2** Simplified visual overview of the building of category label representations using semantic reinforcement and filtering

190 matrix) for all the positive keywords in question. These scores are used either as weights
 191 for an already-existing result list or to create a *result list* from scratch. Category scores are
 192 simply the respective row of the category-document matrix. Similarly to the previous step,
 193 these scores are used either as weights or to create a result list. Finally, negative keywords
 194 are used to exclude all documents from the result list containing any of the negative keywords
 195 chosen.

196 Next, the system calculates keywords candidates and category candidates. These are
 197 equivalent to methods *KeywordsOf()* and *CategoriesOf()* in lines 6 and 7 in Algorithm 1.
 198 Keywords are calculated by averaging the respective *keyword matrix* scores for all documents
 199 in the result list. Additionally, category candidates are calculated by averaging the category
 200 scores for the relevant documents in the category-document matrix.

201 **3 Evaluation of the Automatic Annotation and Search**

202 In this section, we present an evaluation of the framework components, focusing on the
 203 classification system.

204 **3.1 Classification**

205 The classification component of this system has been evaluated with two different datasets:

- 206 ■ The Yle dataset, containing 5096 news articles from Yle (Yleisradio Oy), the Finnish
 207 Broadcasting Company⁸, classified in 11 categories according to their main tag: *politiikka*,

⁸<https://yle.fi>

208 ‘politics’, *talous* ‘economy’, *kulttuuri* ‘culture’, *luonto* ‘nature’, *tiede* ‘science’, *terveys*
 209 ‘health’, *liikenne* ‘transportation’, *urheilu* ‘sports’, *sää* ‘weather’, *parisuhde* ‘intimate
 210 relationships’ and *rikokset* ‘crimes’). This dataset is equivalent to *Yle 1* and *Yle 2* in [7]
 211 combined⁹. It has a mean length, in characters, of 2110.8 ± 1534.5 .

212 ■ The Minilex dataset: a collection of 2567 legislation-related questions classified in 13
 213 categories (*asunto*, *kiinteistö* ‘housing, real state’, *immateriaalioikeus* ‘intellectual prop-
 214 erty’, *irtain omaisuus* ‘chattel’, *lainat*, *velat* ‘loans, debts’, *liikenne* ‘transportation’,
 215 *oikeudenkäynti* ‘legal proceedings’, *perheoikeus*, *perintöoikeus* ‘family law, inheritance
 216 law’, *rikokset* ‘crimes’, *sopimus*, *vahingonkorvaus* ‘contracts, compensation’, *työsuhde*,
 217 *virkasuhde* ‘employment, public service’, *ulosotto*, *konkurssi* ‘debt recovery, bankruptcy’,
 218 *vuokra* ‘rent’, *yrietykset*, *yhteisöt* ‘companies, organisations’) from the legal services website
 219 Minilex¹⁰. As can be observed, there is semantic overlap among some of the categories
 220 in this dataset, such as “debt recovery, bankruptcy” / “loans, debts”; or “housing, real
 221 estate” / “rent” / “contracts, compensation”; or even “transportation” / “crimes”. The
 222 mean length of his dataset is 447.3 ± 236.1 characters.

223 Two different measures were taken: the F_1 score and the rank of the gold label among the
 224 predictions made by the classifier. The latter springs from the fact that this system is used
 225 as a multi-label classifier, so not having the gold tag as the best prediction is not necessarily
 226 a consequential result.

227 The system fares better with the Yle dataset, with an F_1 score of 0.737 and a mean rank
 228 of 0.7 (‘0’ being the gold label, ‘1’ the second prediction, etc; lower is thus better), while the
 229 Minilex dataset obtains an F_1 score of 0.56 and a mean rank of 1.557. Using semantically
 230 reinforced vectors gives a small boost to these numbers: Yle gets an F_1 of 0.742 and a
 231 mean rank of 0.67, and Minilex respectively 0.572 and 1.496. Table 1 details the counts and
 232 cumulative distributions of the results.

233 These numbers show that the present classification method is capable of categorizing the
 234 gold labels as the top prediction in 57–74% of the cases and among the top 3 predictions
 235 in 80–90% of the cases, depending on the dataset used. This variation can possibly be
 236 attributed to a combination of the length of the documents and the quality of the category
 237 labels. These results are not in line with state-of-the art supervised algorithms, which reach
 238 results above 95% and nearing 100% for their top predictions (cf. survey in [9]). However,
 239 training data in Finnish is hard to obtain, and without it a supervised system cannot be
 240 built. The main advantages of the classification algorithm presented in this paper are its
 241 flexible nature – since it only requires a list of category labels in order to work – and its
 242 straightforward integration into the search framework, since it uses the same underlying
 243 technologies.

244 3.2 Document representation

245 These results also show strong consistency on the part of Annif: it is capable of coherently
 246 assigning keywords to both documents and category labels, so that most documents can
 247 be correctly classified when transposed to a vector space of embeddings. FastText also
 248 demonstrate reliability when transporting semantic meanings to vectorial representations.

249 A more decisive evaluation of this component is planned as part of our future work.

⁹Excluding one article, whose main tag has changed. This data is not redistributable.

¹⁰<https://www.minilex.fi>. Their terms of use allows for non-commercial use of the data, but not for its redistribution.

rank	Yle				Minilex			
	Standard		Reinforced		Standard		Reinforced	
	COUNT	CUMD	COUNT	CUMD	COUNT	CUMD	COUNT	CUMD
0	3757	0.737	3779	0.742	1438	0.56	1469	0.572
1	560	0.847	568	0.853	406	0.718	388	0.723
2	269	0.9	264	0.905	179	0.788	187	0.796
3	147	0.929	151	0.934	102	0.828	104	0.837
4	115	0.951	96	0.953	90	0.863	78	0.867
5	91	0.969	104	0.974	59	0.886	64	0.892
6	68	0.983	57	0.985	86	0.919	86	0.926
7	38	0.99	30	0.991	61	0.943	62	0.95
8	25	0.995	22	0.995	56	0.965	44	0.967
9	19	0.999	18	0.999	39	0.98	38	0.982
10	7	1.0	7	1.0	32	0.993	28	0.993
11	-	-	-	-	18	1.0	18	1.0
12	-	-	-	-	1	1.0	1	1.0

■ **Table 1** Count and cumulative distribution (CUMD) of results in the classification task for the Yle and Minilex datasets, in both their standard and reinforced versions

250 3.3 Search

251 No formal evaluation of the search system has been carried out so far. However, Section 4
252 contains tests and insights about the reliability and usability of this component.

253 4 Use Case: Search Engine for Finnish Legislation and Case Law

254 This section describes, via a concrete use case, how the RFBS has been adapted to the legal
255 domain as part of the semantic portal LawSampo [5]. LawSampo is the first Sampo portal to
256 add RFBS-based functionality to complement the search features of previous Sampos, which
257 are mostly facet-based.

258 4.1 The Data

259 LawSampo contains data about Finnish legislation and case law as a harmonized Linked Data
260 knowledge graph. This knowledge graph is based on Semantic Finlex data [10], filtered and
261 transformed into a simpler data model. This was done with the aim of hiding the inherent
262 complexity of legal documentation while keeping the data relevant to anyone interested in
263 the topic. These data transformations were implemented as SPARQL CONSTRUCT queries,
264 and the data is enriched with information about referenced EU legislation, as well as the
265 generated annotations of subject keywords category labels.

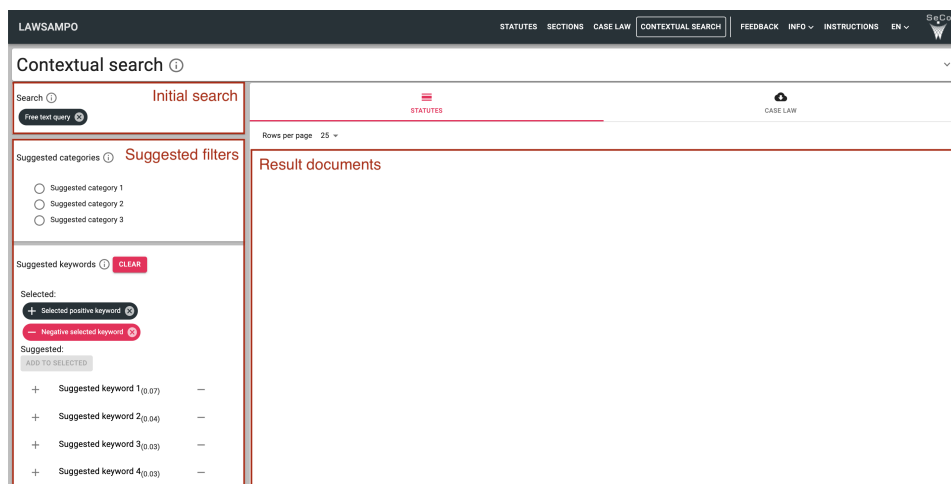
266 The set of category labels used in LawSampo’s RFBS system is based on the Minilex
267 dataset discussed in Section 3. However, in order to avoid the semantic overlap present
268 among some of the categories and to expand the field of possible categories, they were refined
269 with some input from experts at the Ministry of Justice of Finland. The resulting set contains
270 12 categories: *asuminen*, *kiinteistö* ‘housing, real state’, *ihmisoikeudet*, *perusoikeudet* ‘human
271 rights, basic rights’, *omaisuus*, *kaupankäynti*, *kuluttajansuoja* ‘property, commerce, consumer
272 protection’, *julkishallinto*, *valtiohallinto* ‘public administration’, *rahoitus* ‘finance’, *verotus*
273 ‘taxation’, *yrietykset*, *yhteisöt*, *työelämä* ‘business, organizations, working life’, *liikenne*, *kulje-*

274 *tus* ‘traffic’, *perheoikeus*, *perintöoikeus* ‘family law, inheritance’, *rikosasiat*, *oikeudenkäynti*
 275 ‘crime, legal proceedings’, *koulutus* ‘education’, *ympäristö* ‘environment’.

276 The consolidated legislation consists of 3725 statutes and their 82 145 sections, with each
 277 section consisting of the most current version of the full-text contents in Finnish. The case
 278 law dataset consists of 13 470 court decisions and their full-text contents.

279 4.2 LawSampo User Interface for RFBS

280 The LawSampo portal implements the RFBS Algorithm 1 in its “Contextual Search” applica-
 281 tion perspective¹¹. LawSampo allows the user to initialize the RFBS system either by setting
 282 a free-form text query or by selecting one of the provided document categories from the list
 283 presented in Section 4.1. Figure 3 illustrates the user interface after the user has started
 284 a search via a text query. After the initial search, the application switches to an iterative
 285 mode corresponding to the while-do loop (lines 5-12) of Algorithm 1, where the search is
 286 fine-tuned by managing a set of active filters through categories and keywords suggested
 287 by the system. Since no category was selected in the initial phase in Figure 3, the user is
 288 presented with a top-three list of suggested categories (obtained in line 7 in the algorithm)
 289 in addition to 20 suggested keywords (line 6). In this implementation, the values used for
 290 the initial search cannot be changed during the iterative loop: if the initial filter (query or
 291 category) is removed, the iteration stops and the search returns to its empty initial state.



■ **Figure 3** LawSampo’s contextual search UI. The user has made an initial text query and selected one positive and one negative keyword. Next, the user can continue the search by using the suggested categories and keywords in order to modify the active filters. The number shown in subscript next to the keywords is a normalized relevance score

292 Table 2 shows a simple example of how, using the same query, “right of redemption”,
 293 different categories affect the resulting documents. With the given query, it is not surprising
 294 that *Act on the Redemption of Immovable Property and Special Rights* is the most relevant
 295 hit. More interestingly, *Water Act* becomes the second-most referenced document when
 296 the “human rights, basic rights” category is active, whereas the “crime, legal proceedings”

¹¹In the Sampo model, the user is provided with several independent but interlinked applications that use a shared underlying knowledge graph

	crime, legal proceedings	human rights, basic rights	property, commerce, consumer protection
Water act	1	2	1
Building act	1		
Real Estate Formation Act	1		1
Act on the Residential and Commercial Property Information System			1
Act on the Redemption of Immoveable Property and Special Rights	2	3	2

■ **Table 2** How does the selected category affect the resulting documents? This table shows statute results for the query *lunastusoikeus* ‘right of redemption’ with three different categories and a result size of five. The values represent the number of sections returned from each statute

297 category surfaces the *Building Act*. Finally, the “property, commerce and consumer protection”
298 category adds *Act on the Residential and Commercial Property Information System* to the
299 returned documents.

300 Suggested keywords are shown with a relevance score calculated on the basis of the current
301 list of resulting documents. Keywords can be added to the active filters (line 9 in Algorithm
302 1) in either positive or negative mode using the plus (+) or minus (-) buttons respectively. A
303 negative selection excludes any documents containing the given keyword from the result list,
304 as explained in Section 2.3.

305 The result list view can be toggled between statutes and case law documents as two
306 parallel tabs. The results are updated whenever the user changes any of the active filters, i.e.,
307 textual query, category or selected positive and negative keywords. The documents returned
308 are statute sections or case law abstracts. The user is given the possibility to skim through
309 them and to follow links to the full documents. Since the statutes search works at the section
310 level, the results can contain multiple documents from the same statute.

311 4.3 Example Search Scenario

312 This section presents as an example of how, using LawSampo’s “Contextual Search”, the
313 following information need can be satisfied:

314 *I’ve been thinking about making a huge renovation to our house. However, I’m worried*
315 *that because our house is in a such good area, the city might expropriate the lot. If*
316 *that happens, what kind of payout could I be looking at?*

317 Let us begin the search with a simple query based on the information need described
318 above: *pakkolunastus tontti* ‘expropriation lot’. The query is executed against the statutes by
319 default and, for this example, the maximum number of results is set to 10. A summary of
320 the RFBS process for each search iteration can be found in Table 3.

321 The document list resulting from the first iteration (RS_1) does not look very promising:
322 the only vaguely relevant document is *Erämaalaki* ‘Wilderness Act’, which describes how
323 the state can expropriate land in wilderness areas in order to build roads. The result list
324 also contains multiple non-relevant documents related to forced auction and water systems.
325 For the next iteration, we add from the suggested keywords *lunastus* ‘redemption’ and

	Iteration 0	Iteration 1	Iteration 2	Iteration 3
Query	"expropriation lot"			
Selected keywords		+ redemption + right for re- redemption - alluvial land - compulsory auction	+ redemption + right for re- redemption - alluvial land - compulsory auction	+ compensation for redemption - alluvial land - compulsory auction
Document type	Statutes	Statutes	Case law	Statutes
Suggested keywords	real property, <i>redemption</i> , <i>right for redemption</i> cadastral procedures, land surveying, <i>alluvial land</i> , <i>compulsory auction</i>	savings banks, railways, limited companies, shares, town planning, land use policy, company law	roads, construction, municipalities, land use planning, land acquisition, land use, <i>compensation for redemption</i>	indemnities, prices, value (properties), interest (economics), owners

■ **Table 3** Relevance feedback development during the first search iterations in the example scenario. *Iteration 0* is the initial search, which is done without any selected keyword filters. Positive and negative keywords are denoted with '+' and '-' characters respectively. The *Suggested keywords* row only show a subset of the keywords suggested by LawSampo's contextual search. The italic typeface is used to mark keywords selected for the next iteration.

326 *lunastusoikeus* 'right of redemption' as positive keywords and *vesijättömaa* 'alluvial land'
327 and *pakkohuutokauppa* 'compulsory auction' as negative ones to our set of active filters (*AK*
328 in Algorithm 1).

329 The second set of results (RS_2) is already more useful. There are two more references
330 to expropriation related to roads and railroads that can be considered somewhat relevant,
331 but also a match to the highly useful *Land use and Building Act*. The results also indicate
332 that there the positive keyword filters might not work as intended, since the results include
333 documents related to *limited companies*, which deal with the wrong kind of "redemption"
334 with respect to our information need. The results even contain a section from the *Saving*
335 *Back Act*, most likely due to the use of the term *redemption* in the document.

336 We can test our intuition about the keyword "redemption" by switching over to the case
337 law view to verify if the results are similar: the case law results (RS_3) with the same set
338 of active filters are indeed consistent with the results in the statute view, with a couple of
339 only indirectly relevant documents. However, the suggested keywords (AK_3) now include
340 *lunastuskorvaukset* 'compensation for redemption', which matches our information need
341 perfectly.

342 Finally, let us replace "redemption" and "right of redemption" with "compensation from
343 redemption" and swap back to the statutes view to retrieve one more result list (RS_4).
344 This time, all returned documents can be considered useful, including three references to a
345 document titled *Act on the Redemption of Immovable Property and Special Rights*.

346 **5 Discussion**

347 This section overviews earlier related research, summarizes the contributions made by this
348 paper, and outlines paths for future research.

349 **5.1 Related Work**

350 Various methods exist for relevance feedback search [1, 13]. Teevan et al. [22] enrich web
351 search with relevance feedback based on a constructed user profile. Peltonen et al. [11]
352 combine visual intent modeling with exploratory relevance feedback search. Tang et al. [21]
353 have used topic modeling in academic literature search. Song et al. [19] employed topic
354 modeling with relevance search, based on implicit feedback from the topics of the user web
355 search history. In [6], RFBS is combined with topic modelling [2]. However, the method
356 of combining automatic document classification and keyword extractions with relevance
357 feedback search, as described in this paper, is novel.

358 Regarding zero-shot classification methods, most of them work by training a classifier
359 and then adapting it to a new set of categories [4, 12, 23], while [25] also integrates a
360 knowledge graph into their algorithm. Among unsupervised classification models, [15] use
361 skip-gram word embeddings to calculate the semantic similarity between a label and the
362 given documents, which is also the basis of our work. In contrast, [24] treats zero-shot as an
363 entailment problem.

364 **5.2 Contributions**

365 This paper argued for using a combination of topical classification and ontological keywords
366 as a semantic basis for RFBS when exploring textual documents from complex domains,
367 such as legislation and case law. A method for accomplishing this was presented, as well as
368 an implementation of it for testing and evaluating purposes.

369 The content annotation results shown in Section 3 indicate that the proposed classification
370 system, despite its unsupervised nature, is capable of classifying documents correctly 74%
371 of the time (or 90% within the first three predictions) when the classes are semantically
372 non-overlapping and the texts are long enough.

373 Section 2.3 illustrated how the RFBS algorithm suggested in this paper can be used
374 in practice, illustrating how LawSampo's Contextual Searcher perspective can be used to
375 navigate documents from a semantically complex domain successfully in an iterative fashion.
376 Automatically suggested keywords mitigate the burden on the user of coming up with
377 suitable queries and can provide valuable feedback even when the user selects non-optimal
378 keywords. We have not yet performed a more general evaluation of the search functionality
379 besides testing the system in selected individual problems. Nevertheless, the experiments
380 presented in this paper suggest that a combination of ontological keyword annotations and
381 topical classifications with word embeddings can create a useful semantic basis for the RFBS
382 paradigm when searching and exploring textual legal documents.

383 **5.3 Future Work**

384 More research will be done in order to improve the vectorial representation of the documents.
385 As it stands, these representations are entirely based on each document's set of keywords,
386 which add depth (by emphasizing these keywords) but subtract breadth (other details of the
387 text) to them. We plan to pursue two main lines of research in the future: one line aims
388 at improving the set of representative keywords by both filtering out unrelated suggestions

389 and adding new ones via ontology relations, named entity recognition and other keyword
 390 extraction algorithms, provided they can be integrated into the system; the second line aims
 391 at investigating alternative representation vectors partly based on whole-text embeddings. A
 392 third promising line of research consists in the automatic identification of major topics in the
 393 data as an alternative to the user-provided category list. This can possibly be accomplished by
 394 capitalizing on existing clustering methods and on ontological relations among the keywords.

395 LawSampo is the first portal in the the Sampo series of systems¹² to take advantage of
 396 this search functionality based on relevance feedback. Similar methods are planned as part of
 397 the upcoming new Sampos as well. These include especially the ParliamentSampo system¹³,
 398 which incorporates over 900 000 parliamentary debate speeches [18] from the Parliament of
 399 Finland (1907–2021), documents that are related to the legislative texts found in LawSampo.

400 ——— References ———

- 401 1 R. Baeza-Yates and B. Ribeiro-Neto. *Modern Information Retrieval (2nd Ed.)*. Addison-Wesley
 402 Longman Publishing Co., Inc., 2011.
- 403 2 David M. Blei. Probabilistic topic models. *Commun. ACM*, 55(4):77–84, April 2012. URL:
 404 <http://doi.acm.org/10.1145/2133806.2133826>, doi:10.1145/2133806.2133826.
- 405 3 Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching Word Vec-
 406 tors with Subword Information. *Transactions of the Association for Computational Linguistics*,
 407 5:135–146, 12 2017. doi:10.1162/tac1_a_00051.
- 408 4 Yann N. Dauphin, Gokhan Tur, Dilek Hakkani-Tur, and Larry Heck. Zero-Shot Learning for
 409 Semantic Utterance Classification. *ICLR 2014*, 2014. URL: <http://arxiv.org/abs/1401.0509>,
 410 arXiv:1401.0509.
- 411 5 Eero Hyvönen, Minna Tamper, Arttu Oksanen, Esko Ikkala, Sami Sarsa, Jouni Tuominen,
 412 and Aki Hietanen. LawSampo: A semantic portal on a linked open data service for finnish
 413 legislation and case law. In *The Semantic Web: ESWC 2020 Satellite Events. Revised Selected*
 414 *Papers*, pages 110–114. Springer–Verlag, 2019.
- 415 6 Mikko Koho, Erkki Heino, Arttu Oksanen, and Eero Hyvönen. Toffee - semantic media search
 416 using topic modeling and relevance feedback. In *Proceedings of the ISWC 2018 Posters &*
 417 *Demonstrations, Industry and Blue Sky Ideas Tracks*. CEUR Workshop Proceedings, October
 418 2018. Vol 2180. URL: <http://ceur-ws.org/Vol-2180/>.
- 419 7 Rafael Leal. Unsupervised zero-shot classification of Finnish documents using pre-trained
 420 language models. Master’s thesis, University of Helsinki, Department of Digital Humanities,
 421 2020. URL: <http://urn.fi/URN:NBN:fi:hulib-202012155147>.
- 422 8 Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhersch, and Armand Joulin.
 423 Advances in pre-training distributed word representations. In *Proceedings of the International*
 424 *Conference on Language Resources and Evaluation (LREC 2018)*, 2018.
- 425 9 Shervin Minaee, Nal Kalchbrenner, E. Cambria, Narjes Nikzad, M. Chenaghlu, and Jianfeng
 426 Gao. Deep Learning Based Text Classification: A Comprehensive Review. *ArXiv*, 2020.
- 427 10 Arttu Oksanen, Jouni Tuominen, Eetu Mäkelä, Minna Tamper, Aki Hietanen, and Eero
 428 Hyvönen. Semantic Finlex: Transforming, publishing, and using finnish legislation and case
 429 law as linked open data on the web. In G. Peruginelli and S. Faro, editors, *Knowledge of the*
 430 *Law in the Big Data Age*, volume 317 of *Frontiers in Artificial Intelligence and Applications*,
 431 pages 212–228. IOS Press, 2019. ISBN 978-1-61499-984-3 (print); ISBN 978-1-61499-985-0
 432 (online). URL: <http://doi.org/10.3233/FAIA190023>.

¹²<https://seco.cs.aalto.fi/applications/sampo/>

¹³<https://seco.cs.aalto.fi/projects/semparl/en/>

- 433 11 Jaakko Peltonen, Jonathan Strahl, and Patrik Floréen. Negative relevance feedback for
434 exploratory search with visual interactive intent modeling. In *Proceedings of the 22nd International
435 Conference on Intelligent User Interfaces*, pages 149–159. ACM, 2017.
- 436 12 Anthony Rios and Ramakanth Kavuluru. Few-Shot and Zero-Shot Multi-Label Learning for
437 Structured Label Spaces. *EMNLP*, 2018. doi:10.18653/v1/D18-1352.
- 438 13 Gerard Salton and Chris Buckley. Improving retrieval performance by relevance feedback.
439 *Journal of the American Society for Information Science*, 41(4):288, 1990.
- 440 14 Gerard Salton and Michael J. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill
441 Computer Science Series. McGraw-Hill, 1983.
- 442 15 Prateek Veeranna Sappadla, Jinseok Nam, Eneldo Loza Mencía, and Johannes Fürnkranz.
443 Using semantic similarity for multi-label zero-shot classification of text documents. In *ESANN*,
444 2016.
- 445 16 Katri Seppälä and Eero Hyvönen. Asiasanaston muuttaminen ontologiaksi. Yleinen suomalainen
446 ontologia esimerkkinä FinnONTO-hankkeen mallista (Changing a keyword thesaurus into
447 an ontology. General Finnish Ontology as an example of the FinnONTO model). Technical
448 report, March 2014. URL: <https://www.doria.fi/handle/10024/96825>.
- 449 17 Teemu Sidoroff and Eero Hyvönen. Semantic e-government portals - a case study. In
450 *Proceedings of the ISWC-2005 Workshop Semantic Web Case Studies and Best Prac-
451 tices for eBusiness SWCASE05*, 2005. URL: [https://seco.cs.aalto.fi/publications/2005/
452 sidoroff-hyvonen-semantic-e-government-2005.pdf](https://seco.cs.aalto.fi/publications/2005/sidoroff-hyvonen-semantic-e-government-2005.pdf).
- 453 18 Laura Sinikallio, Senka Drobac, Minna Tamper, Rafael Leal, Mikko Koho, Jouni Tuominen,
454 Matti La Mela, and Eero Hyvönen. Plenary debates of the Parliament of Finland as linked
455 open data and in Parla-CLARIN markup, March 2021. Paper submitted for evaluation, LDK
456 2021.
- 457 19 Wei Song, Yu Zhang, Ting Liu, and Sheng Li. Bridging topic modeling and personalized
458 search. In *Proceedings of the 23rd International Conference on Computational Linguistics:
459 Posters*, pages 1167–1175. Association for Computational Linguistics, 2010.
- 460 20 Osma Suominen. Annif: DIY automated subject indexing using multiple algorithms. *LIBER
461 Quarterly*, 29(1):1–25, 07 2019. doi:10.18352/lq.10285.
- 462 21 Jie Tang, Ruoming Jin, and Jing Zhang. A topic modeling approach and its integration into
463 the random walk framework for academic search. In *Data Mining, 2008. ICDM'08. Eighth
464 IEEE International Conference on*, pages 1055–1060. IEEE, 2008.
- 465 22 Jaime Teevan, Susan T. Dumais, and Eric Horvitz. Personalizing search via automated analysis
466 of interests and activities. In *Proc. of the 28th Annual International ACM SIGIR Conference*,
467 SIGIR '05, pages 449–456. ACM, 2005.
- 468 23 Zhiquan Ye, Yuxia Geng, Jiaoyan Chen, Jingmin Chen, Xiaoxiao Xu, SuHang Zheng, Feng
469 Wang, Jun Zhang, and Huajun Chen. Zero-shot Text Classification via Reinforced Self-training.
470 In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*,
471 pages 3014–3024. Association for Computational Linguistics, 2020. doi:10.18653/v1/2020.
472 acl-main.272.
- 473 24 Wenpeng Yin, Jamaal Hay, and Dan Roth. Benchmarking Zero-shot Text Classification:
474 Datasets, Evaluation and Entailment Approach. In *Proceedings of the 2019 Conference on
475 Empirical Methods in Natural Language Processing and the 9th International Joint Confer-
476 ence on Natural Language Processing (EMNLP-IJCNLP)*, pages 3914–3923. Association for
477 Computational Linguistics, 2019. doi:10.18653/v1/D19-1404.
- 478 25 Jingqing Zhang, Piyawat Lertvittayakumjorn, and Yike Guo. Integrating Semantic Knowledge
479 to Tackle Zero-shot Text Classification. In *Proceedings of the 2019 Conference of the North
480 American Chapter of the Association for Computational Linguistics: Human Language Tech-
481 nologies, Volume 1 (Long and Short Papers)*, pages 1031–1040. Association for Computational
482 Linguistics, 2019. doi:10.18653/v1/N19-1108.