



Departamento de Inteligencia Artificial  
Escuela Técnica Superior de Ingenieros Informáticos

PhD Thesis

# Semantically-enabled Browsing of Large Multilingual Document Collections

Author: Carlos Badenes-Olmedo  
Supervisors: Prof. Dr. Oscar Corcho

xxx, 2020



Tribunal nombrado por el Sr. Rector Magfco. de la Universidad Politécnica de Madrid,  
el día XX de xxx de 2020.

Presidente: Dr. Xxx Xxx

Vocal: Dra. Xxx Xxx

Vocal: Dr. Xxx Xxx

Vocal: Dra. Xxx Xxx

Secretario: Dr. Xxx Xxx

Suplente: Dra. Xxx Xxx

Suplente: Dr. Xxx Xxx

Realizado el acto de defensa y lectura de la Tesis el día X de xxx de 2020 en la Escuela  
Técnica Superior de Ingenieros Informáticos

Calificación: \_\_\_\_\_

EL PRESIDENTE

VOCAL 1

VOCAL 2

VOCAL 3

EL SECRETARIO



A mis padres.  
A Beatriz.  
A Martín y Alonso.



## Agradecimientos

xxxxxx





## Abstract

XXXXXX



## Resumen

xxxxxxx



# Contents

<b>List of Figures</b>	<b>xvii</b>
<b>List of Tables</b>	<b>xix</b>
<b>Acronyms</b>	<b>xxi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Contributions . . . . .	4
1.2 Thesis Structure . . . . .	5
1.3 Publications . . . . .	5
<b>2 Related Work</b>	<b>9</b>
2.1 Text Annotations . . . . .	9
2.2 Topic Distributions . . . . .	10
2.2.1 Distance Measures . . . . .	13
2.3 Multilingual Topic Alignment . . . . .	17
<b>3 Research Objectives</b>	<b>21</b>
3.1 Research Hypotheses . . . . .	22
3.2 Research Challenges . . . . .	24
3.2.1 Topic-model Programming Interface . . . . .	24
3.2.2 Explainable Topic-based Associations . . . . .	25
3.2.3 Large-scale Comparisons of Topic Distributions . . . . .	25
3.2.4 Automatic Cross-lingual Topic Alignment . . . . .	26
3.3 Research Methodology . . . . .	26
3.3.1 Scalable Creation and Inference of Topics . . . . .	28
3.3.2 Explainable Topic-based Associations . . . . .	28

3.3.3	Large-scale Comparisons of Topic Distributions . . . . .	30
3.3.4	Automatic Cross-lingual Topic Alignment . . . . .	30
<b>4</b>	<b>Scalable Creation and Inference of Topics</b>	<b>33</b>
4.1	Document Workflow . . . . .	33
4.2	librAIry . . . . .	34
4.2.1	Functional Features . . . . .	34
4.2.1.1	Resources . . . . .	35
4.2.1.2	Event-based Paradigm . . . . .	36
4.2.1.3	Linked Data Principles . . . . .	36
4.2.2	Framework Architecture . . . . .	36
4.2.2.1	Event-Bus . . . . .	36
4.2.2.2	API . . . . .	37
4.2.2.3	Storage . . . . .	37
4.2.2.4	Modules . . . . .	38
4.3	Model-as-a-Service . . . . .	39
4.4	Summary . . . . .	39
<b>5</b>	<b>Explainable Topic-based Associations</b>	<b>41</b>
5.1	Topic Relevance . . . . .	41
5.2	Topic-based Clustering . . . . .	41
5.3	Summary . . . . .	41
<b>6</b>	<b>Large-scale Comparisons of Topic Distributions</b>	<b>43</b>
6.1	Document Similarity . . . . .	43
6.2	Hashing Topic Distributions . . . . .	43
6.3	Summary . . . . .	43
<b>7</b>	<b>Cross-lingual Document Similarity</b>	<b>45</b>
7.1	Synset-based Representational Space . . . . .	45
7.2	Cross-lingual Models . . . . .	45
7.3	Summary . . . . .	45

<b>8</b>	<b>Evaluation</b>	<b>47</b>
8.1	Evaluation Metrics . . . . .	47
8.2	Text Representativeness . . . . .	47
8.3	Large-scale Text Processing . . . . .	47
8.4	Topic-based Clustering . . . . .	49
8.5	Cross-lingual Similarity . . . . .	49
8.6	Conclusions . . . . .	49
<b>9</b>	<b>Experiments</b>	<b>51</b>
9.1	Polypharmacy and Drug-drug Interactions . . . . .	51
9.2	Corpus Viewer . . . . .	51
9.3	ODS Classifier . . . . .	51
9.4	Drugs4Covid . . . . .	51
<b>10</b>	<b>Conclusions</b>	<b>53</b>
10.1	Assumptions and Restrictions . . . . .	53
10.2	Contributions . . . . .	53
10.3	Impact . . . . .	53
10.4	Limitations . . . . .	53
10.5	Future Work . . . . .	53
	<b>Bibliography</b>	<b>55</b>





# List of Figures

2.1	Distance values between 10 pair of documents from topic models with 100-to-2000 dimensions. The Kullback-Liebler(a), Jensen-Shannon Divergence(b), Hellinger(c) and S2JSD(c) metrics are considered. . . . .	15
3.1	Research dimensions of the thesis. The first ones must be overcome before reaching higher dimensions. . . . .	29
4.1	Domain deleted flow. . . . .	35
4.2	Resource states. . . . .	37
4.3	Modules. . . . .	38



# List of Tables

3.1	Hypotheses and research dimensions. . . . .	23
3.2	Open Research Challenges and Hypotheses. . . . .	27
3.3	Research and technical objectives and their related challenges. . . . .	32



# Acronyms

**API:** Application Programming Interface

**CQ:** Competency Question

**GUI:** Graphical User Interface

**IDE:** Integrated Development Environment

**LD:** Linked Data

**LOD:** Linked Open Data

**UML:** Unified Modeling Language

**URI:** Uniform Resource Identifier

**URL:** Uniform Resource Locator

**WUI:** Web User Interface



# Chapter 1

## Introduction

Huge amounts of data, in the form of textual documents, are daily produced in digital format. Every second more than two thousand blog entries are published, nine thousand tweets are written, and more than two million emails are sent in the Internet<sup>1</sup>. The number of scientific publications per year has increased by 8-9% in the last decade (Rob Johnson and Mabe, 2018). More than one million papers, about two per minute, were submitted to the PubMed database, the leading database of references and abstracts on life sciences and biomedical topics, in the last year. Statistics on judicial activity is similar. More than 168,000 procedural documents and 3,000 judicial notices were published in the Official Journal of the European Union in 2019<sup>2</sup>. Unlike the academic domain where articles are mostly published in English, legal documents are usually available in multiple languages. The Court of Justice of the European Union had to translated over 1 million texts into its 24 official languages, with 552 possible language combinations, in just one year. These numbers make it virtually impossible for an expert in academic or legal domain to stay abreast by only reading a few articles nowadays. Navigate the growing torrent of textual data and explore its content is not only necessary, but has become a second job that experts must reconcile with their daily tasks.

Some initiatives based on document retrieval techniques are underway in order to facilitate text review in big collections. Major digital publishers specialized in scientific<sup>3</sup>,

---

<sup>1</sup><https://www.internetlivestats.com/one-second>

<sup>2</sup><https://curia.europa.eu>

<sup>3</sup><https://www.nature.com>

technical<sup>4</sup>, and medical<sup>5</sup> content provide search engines to make it easier to browse their collections of scientific articles. Given a few keywords, a list of relevant papers is retrieved and offered for reading. Legal documents are also exploited with similar solutions. The spanish<sup>6</sup>, american<sup>7</sup> and european<sup>8</sup> patent and intellectual property registration offices, for example, allow exploring their patent collections by search engines guided by keywords and/or categories. These categories are available because documents are manually categorized by their authors according to the International Patent Classification (IPC) system. It contains approximately 70,000 different codes for different technical areas. This label-based browsing<sup>9</sup> has been also adopted by several academic search engines<sup>10</sup> to organize papers by research areas, or even by evaluation tasks to browse state-of-the-art methods<sup>11</sup>. In natural language processing domain, for example, research papers are organized into 256 tasks such as 'knowledge representation', 'question-answer', 'machine translation' and so on. However, while there are initiatives to normalize research areas, the use of keywords by authors in the form of tags to categorize their scientific papers is still insufficient and some text processing tasks are necessary to set labels to articles following an uniform criteria. One of the main reasons that limit its widespread use is the difficulty that authors have in picking labels that describe their research work in sufficient detail.

Along with searches based on keywords and categories, the third key aspect aimed at facilitating the exploration of documentary corpus is the provision of related texts. A documentary exploration does not stop when a relevant article is found, but starts from its content shaping the area of interest. Most academic<sup>12</sup> and legal<sup>13</sup> search engines provide a list of related documents for each text and offer navigating through them. The relationship can be of reference, when documents are cited by others, or of content, when documents share a thematic area. The chains of articles derived from that related content can lead to more complex structures when cross-relations are

---

<sup>4</sup><https://www.elsevier.com>

<sup>5</sup><https://pubmed.ncbi.nlm.nih.gov>

<sup>6</sup><https://www.oepm.es>

<sup>7</sup><https://www.uspto.gov/>

<sup>8</sup><https://www.epo.org>

<sup>9</sup><https://patents.google.com>

<sup>10</sup><https://academic.microsoft.com>

<sup>11</sup><https://paperswithcode.com>

<sup>12</sup><https://www.semanticscholar.org>

<sup>13</sup><https://patents.google.com>



assumed. A document can be related to another that, in turn, is related to a third one that can be also related to the first article. This content-guided exploration helps to browse document collections by areas of interest not necessarily aligned with a list of predefined categories. A visual overview of an academic field, for example, can be provided by showing graphs of articles with similar content<sup>14</sup>.

While these initiatives are valuable efforts to address access to huge amounts of documents, they are still insufficient to examine the content offered by their texts. On an individual level, the knowledge derived from a text comes from the concepts evoked by its words Griffiths et al. (2007). On a collective level, the knowledge derived from a document collection emerges from the relationships among its texts Kenter and de Rijke (2015). Exploring a textual corpora and therefore acquiring some knowledge of its texts requires understanding how its documents are organized through the concepts evoked by its words. It is necessary to focus on why some texts are related to others, and what concepts are key to those relationships. But analyzing and comparing texts on a large scale requires addressing some challenges imposed by external conditions that have appeared in recent years:

- Complexity: The huge number of documents has forced a reconsideration of the way to compare them in order to be able to deal with big collections. The time required to compute each comparison should be reduced as much as possible.
- Efficiency: The algorithms, besides being accurate enough, must be also efficient in order to be applied on a large scale. Brute-force techniques cannot be applied to compare all items in a huge corpus.
- Transparency: The associations between documents must be explained in such a way that the relationship itself provides knowledge about the content of the texts. It is not enough that one text is related to another, it is necessary to explain why it is so.
- Multilinguality: In addition, the increasing availability of texts written in different languages also makes it necessary to address comparison in multilingual

---

<sup>14</sup><https://www.connectedpapers.com>

collections. In these collections, external translation systems cannot be considered, since they increase processing costs and potentially introduce a bias in the relationships obtained.

This work is aimed at facilitating the exploration of huge collections of documents that contain texts written in multiple languages. We address the problem of comparing them on a large scale enabling a semantic-aware exploration through their content. Our proposal automatically discovers thematic associations between texts using probabilistic topic models and organizes document collections so that it can be efficiently and transparently browsed through the related content regardless of their language.

## 1.1 Contributions

The work presented in this thesis makes the following contributions:

- **Large-scale Topic-based Text-processing Pipeline:** We define a scalable text processing pipeline following web standards and software best practices for the creation and exploitation of probabilistic topic models.
- **Topic Model-as-a-Service:** We propose a format to distribute and reuse probabilistic topic models.
- **Hierarchical Thematic Annotations:** We present a method to annotate texts by topic hierarchies automatically inferred from their content.
- **Massive Document Comparisons:** We leverage multi-level topic annotations to efficiently index and retrieve related documents while allowing the exploration of the collection by the themes inferred from its texts.
- **Cross-lingual Document Relations:** We introduce a technique to transform probabilistic topics from different languages into a single representation space based on shared concepts where texts can be thematically related regardless of the language used.

## 1.2 Thesis Structure

The thesis is structured as follows:

*Chapter 2* analyses the state of the art and describes the main concepts handled throughout the thesis. *Chapter 3* presents the research problems and hypotheses that guide our work, and details the methodology that has been followed. *Chapter 4* describes the software architecture proposed to analyze huge document collections and the format suggested to distribute and reuse topic models on which the work presented in this thesis is built. *Chapter 5* details the text annotation algorithm from probabilistic topics. *Chapter 6* shows how to efficiently store and search documents from large collections when they are annotated with topic hierarchies. *Chapter 7* explains the method to relate texts written in different languages from their main topics without the need for translation. This approach is evaluated in *Chapter 8*, where the results are explained in detail. *Chapter 9* provides information on real-world projects where contributions from this thesis have been used. Finally, *Chapter 10* describes conclusions and future lines of work.

## 1.3 Publications

The following publications support the research work presented in this thesis:

- *Chapter 4*:
  - **Carlos Badenes-Olmedo**, José Luis Redondo-Garcia, and Oscar Corcho. Distributing Text Mining tasks with libRAIry. Proceedings of the 17th ACM Symposium on Document Engineering (DocEng). Association for Computing Machinery, Valletta, Malta. 2017.
  - Victoria Kosa, Alyona Chugunenko, Eugene Yuschenko, **Carlos Badenes-Olmedo**, Vadim Ermolayev, and Aliaksandr Birukou. Semantic saturation in retrospective text document collections. Information and Communication Technologies in Education, Research, and Industrial Applications (ICTERI) PhD Symposium, vol. 1851, pages 1-8. CEUR-WS. 2017
  - Victoria Kosa, David Chaves-Fraga, Dmitriy Naumenko, Eugene Yuschenko, **Carlos Badenes-Olmedo**, Vadim Ermolayev, Aliaksandr Birukou, Nick

Bassiliades, Hans-Georg Fill, Vitaliy Yakovyna, Heinrich C. Mayr, Mykola Nikitchenko, Grygoriy Zholtkevych, and Aleksander Spivakovsky. Cross-Evaluation of Automated Term Extraction Tools by Measuring Terminological Saturation. *Information and Communication Technologies in Education, Research, and Industrial Applications*, pages 135-163. Springer International Publishing. 2018

- *Chapter 5:*

- **Carlos Badenes-Olmedo**, José Luis Redondo-García, and Oscar Corcho. Efficient Clustering from Distributions over Topics. *Proceedings of the 9th International Conference on Knowledge Capture (K-CAP)*, Article 17, 1–8. Association for Computing Machinery, Austin, TX, USA. 2017.
- José Manuel Gómez-Pérez, Ronald Denaux, Daniel Vila and **Carlos Badenes-Olmedo**. Hybrid Techniques for Knowledge-based NLP - Knowledge Graphs meet Machine Learning and all their Friends. *Proceedings of Workshops and Tutorials of the 9th International Conference on Knowledge Capture (K-CAP)*, 69–70. CEUR-WS, Austin, TX, USA. 2017
- **Carlos Badenes-Olmedo**, Jose Luis Redondo-Garcia, and Oscar Corcho. An initial Analysis of Topic-based Similarity among Scientific Documents based on their Rhetorical Discourse Parts. *Proceedings of the 1st Workshop on Enabling Open Semantic Science (SemSci) co-located with 16th International Semantic Web Conference (ISWC 2017)*, 15-22. Vienna, Austria. 2017.

- *Chapter 6:*

- **Carlos Badenes-Olmedo**, José Luis Redondo-García, and Oscar Corcho. Large-scale Semantic Exploration of Scientific Literature Using Topic-based Hashing Algorithms. *Semantic Web*, vol. Pre-press, no. Pre-press, pp. 1-16. 2020

- *Chapter 7:*

- **Carlos Badenes-Olmedo**, José Luis Redondo-García, and Oscar Corcho. Scalable Cross-lingual Document Similarity through Language-specific

Concept Hierarchies. Proceedings of the 10th International Conference on Knowledge Capture (K-CAP). Association for Computing Machinery, 147–153. Marina Del Rey, CA, USA. 2019

- **Carlos Badenes-Olmedo**, José Luis Redondo-García, and Oscar Corcho. Legal document retrieval across languages: topic hierarchies based on synsets. arXiv e-prints, arXiv:1911.12637. 2019
- Ahmet Soylu, Oscar Corcho, Brian Elvesaeter, **Carlos Badenes-Olmedo**, Francisco Yedro, Matej Kovacic, Matej Posinkovic, Ian Makgill, Chris Taggart, Elena Simperl, Till C. Lech, and Dumitru Roman. Enhancing Public Procurement in the European Union through Constructing and Exploiting an Integrated Knowledge Graph. Proceedings of the 19th International Semantic Web Conference (ISWC), (under revision). 2020



## Chapter 2

# Related Work

Texts usually contain noisy, non-relevant information and keeping only what can bring value for the involved agents (general consumers, experts, companies, investors...) becomes a challenge. In order to facilitate the exploration of large and multilingual document collections we need to process texts in a way that is computationally affordable and enables a semantic-aware exploration of the knowledge inside it. Let's look at the existing techniques and methods involved in tasks related to this objective.

### 2.1 Text Annotations

A necessary first step before using documents for knowledge-intensive tasks is to process them following different techniques to leverage their content. Recent studies Westergaard et al. (2017) Sciences and life (2016) have shown that mining full-text articles give consistently better results than only using sections or summaries. Given the size limitations and concise nature of summaries, they often omit descriptions or results that are considered to be less relevant but still are important for some IR tasks Divoli et al. (2012). Since this behavior is present in many other domains, our interest is focused on processing full texts, not summaries or parts of texts, so we have to take it into account during the whole process.

The annotation of human-readable documents is a well-known problem in the Artificial Intelligence (AI) domain in general and Information Retrieval (IR) and Natural Language Processing (NLP) fields in particular. There already exist a broad set of tools and frameworks able to analyze text for automatically producing such annotations, at

very different levels of granularity: from minimal units such as terms and entities, to descriptors at the level of the entire collection such as summaries or topics. spaCy<sup>15</sup>, NLTK, Stanford CoreNLP Manning et al. (2014) and IXA pipes Agerri et al. (2014) are well-known frameworks for text annotation using Part-of-Speech (PoS) tagging or Named Entity Recognition (NER). Mallet<sup>16</sup>, SparkLDA<sup>17</sup> or Gensim are also widely used libraries to perform more advanced tasks such as topic modeling or clustering over document collections.

Although they are all widely used resources, their design has not paid special attention to facilitating their interoperability. The main reason that limits the reuse and expansion possibilities of NLP-based tools is their strong technological dependence. Mallet has been designed as a Java library, so that spaCy cannot be integrated and used for text processing because its technology is based on Python. This example can even be extended to the models that Mallet generates, since they cannot be used from other tools either as they are distributed in a proprietary format. To the best of our knowledge, the efforts that have been made go in the direction of creating ecosystems that integrate resources<sup>18</sup>, rather than creating open environments that follow standards to be (re)used freely. We are focused on the transversal problem of making those standalone tools coexisting under a same solution in an open environment. In this thesis we propose *reusable text annotation models and scalable document processing pipelines to integrate them*.

## 2.2 Topic Distributions

Once texts have been processed, vector space models (VSM) are created from their content. The objective is twofold, on the one hand to make a huge collection manageable since we move from having lots of words for each text to only one vector per document, and on the other hand to have representations based on metric spaces where calculations can be made, for example comparisons by measuring vector distances. The definition and number of dimensions for each vector are key aspects in a VSM. Traditional retrieval tasks over large collections of textual documents highly rely on individual features like

---

<sup>15</sup><https://spacy.io>

<sup>16</sup><http://mallet.cs.umass.edu>

<sup>17</sup><https://spark.apache.org/mllib/>

<sup>18</sup><https://onnx.ai/>



term frequencies (TF) Hearst and Hall (1999). A representational space is created where each word in the vocabulary is projected by a separate and orthogonal dimension. Term Frequency-Inversed Document Frequency (TF-IDF) relativizes the relevance of each term with respect to the entire corpus. The loss of semantic information and the high-number of dimensions are the main drawbacks of these approaches that lead to the emergence of other techniques. New ways of characterizing documents based on the automatic generation of models surfacing the main subjects covered in the corpus are developed during recent years. Among them, text embedding proposes transforming texts into low-dimensional vectors by prediction methods based on (i) word sequences or (ii) bag-of-words. The first approach assumes words with similar meanings tend to occur in similar contexts. It considers word order relevant and is based on Neural Models (NM) that learn word vectors from pairs of target and context words, where context words are taken as words observed to surround a target word. Document vectors are usually created by averaging the word vectors they contain or by considering them as target and context items. The second approach does not consider the order of the words to be relevant, but their frequency is. It assumes words with similar meanings will occur in similar documents. This second approach is used in our work since we are not only interested in representing words and documents, but we also seek structures that can provide knowledge about the collection as a whole.

Probabilistic Topic Models (PTM) Blei et al. (2010) are statistical methods based on bag-of-words that analyze the words of the original texts to discover the themes that run through them, how those themes are connected to each other, or how they change over time. PTM do not require any prior annotations or labeling of the documents. The topics emerge, as hidden structures, from the analysis of the original texts. These structures are topics distributions, per-resource topic distributions or per-resource per-word topic assignments. In turn, a topic is a distribution over terms that is biased around those words associated to a single theme. This interpretable hidden structure annotates each resource in the collection and these annotations can be used to perform deeper analysis about relationships between resources. Topic-based representations bring a lot of potential when applied over different IR tasks, as evidenced by recent works in different domains such as scholarly (Gatti et al., 2015), health (Lu et al., 2016; Tapi Nzali et al., 2017), legal (Greene and Cross, 2016; O’Neill et al., 2017), news (He et al., 2017) and social networks (Cheng et al., 2014). Topic modeling provides us an

algorithmic solution to organize and annotate large collections of textual documents according to their topics.

The simplest generative topic model is *Latent Dirichlet Allocation* (LDA) Blei et al. (2003). This and other models such as *Latent Semantic Analysis* (LSA) Deerwester et al. (1990) or *Probabilistic Latent Semantic Analysis* (pLSA) Hofmann (2001) are part of the field known as topic modeling. They are well-known latent variable models for high dimensional data, such as the bag-of-words representation for textual data or any other count-based data representation. They try to capture the intuition that documents can exhibit multiple themes. Each document exhibits each topic in different proportion, and each word in each document is drawn from one of the topics, where the selected topic is chosen from the per-document distribution over topics. All the documents in a collection share the same set of topics, but each document exhibits these topics in a different proportion.

Documents are represented as a vector of counts with  $W$  components, where  $W$  is the number of words in the vocabulary. Each document in the corpus is modeled as a mixture over  $K$  topics, and each topic  $k$  is a distribution over the vocabulary of  $W$  words. Formally, a *topic* is a multinomial distribution over words of a fixed vocabulary representing some concept. Depending on the function used to describe that distribution there are different algorithms to create topic models. While LSA and pLSA propose a singular value decomposition, LDA, influenced by the generative Bayesian framework to avoid some of the over-fitting issues that were observed with pLSA, suggests the use of a Dirichlet function. It is a continuous multivariate probability distribution parameterized by a vector of positive reals whose elements sum to 1. It is *continuous* because the relative likelihood for a random variable to take on a given value is described by a probability density function, and is *multivariate* because it has a list of variables with unknown values. In fact, the Dirichlet distribution is the conjugate prior of the categorical distribution and multinomial distribution and is responsible for, unlike LSA and pLSA, LDA can infer topic distributions in texts that have not been used during training.

It is not a restrictive clustering model, where each document is assigned to one cluster, but allows documents to exhibit multiple topics. Since LDA is unsupervised, the topics covered in a set of documents are discovered from the own corpus. The mixed-membership assumptions lead to sharper estimates of word co-occurrence patterns, key

to this thesis that proposes a *thematic and low-dimensional feature space suitable for big real-world data sets, where documents are only described by their most relevant topics.*

### 2.2.1 Distance Measures

In a *Topic Model* the feature vector is a topic distribution expressed as vector of probabilities. Taking into account this premise, the similarity between two topic-based resources will be based on the distance between their topic distributions, which can be also seen as two probability mass functions. A commonly used metric is the *Kullback-Liebler* (KL) divergence. However, it presents two major problems: (1) when a topic distribution is zero, KL divergence is not defined and (2) it is not symmetric, which does not fit well with semantic similarity measures that are usually symmetric Rus et al. (2013).

*Jensen-Shannon* (JS) divergence Rao (1982) Lin (1991) solves these problems considering the average of the distributions as below Celikyilmaz et al. (2010):

$$JS(p, q) = \sum_{i=1}^K p_i * \log \frac{2 * p_i}{p_i + q_i} + \sum_{i=1}^K q_i * \log \frac{2 * q_i}{q_i + p_i} \quad (2.1)$$

where  $K$  is the number of topics and  $p, q$  are the topics distributions

It can be transformed into a similarity measure as follows Dagan et al. (1999) :

$$sim_{JS}(D_i, D_j) = 10^{-JS(p, q)} \quad (2.2)$$

where  $D_i, D_j$  are the documents and  $p, q$  the topic distributions of each of them.

*Hellinger* (He) distance is also symmetric and is used along with JS divergence in various fields where a comparison between two probability distributions is required Blei and Lafferty (2007) Hall et al. (2008) Boyd-Graber and Resnik (2010):

$$He(p, q) = \frac{1}{\sqrt{2}} \cdot \sqrt{\sum_{i=1}^K (\sqrt{p_i} - \sqrt{q_i})^2} \quad (2.3)$$

It can be transformed into a similarity measure by subtracting it from 1 Rus et al. (2013) such that a zero distance means max. similarity score and vice versa:

$$sim_{He}(D_i, D_j) = 1 - He(p, q) \quad (2.4)$$

However, all these metrics are not well-defined distance metrics, that is, they do not satisfy triangle inequality Charikar (2002). This inequality considers  $d(x, z) \leq d(x, y) + d(y, z)$  for a metric  $d$  Griffiths et al. (2007). It places strong constraints on distance measures and on the locations of points in a space given a set of distances. As a metric axiom the triangle inequality must be satisfied in order to take advantage of the inferences that can be deduced from it. Thus, if similarity is assumed to be a monotonically decreasing function of distance, this inequality avoids the calculation of all pairs of similarities by considering that if  $x$  is similar to  $y$  and  $y$  is similar to  $z$ , then  $x$  must be similar to  $z$ .

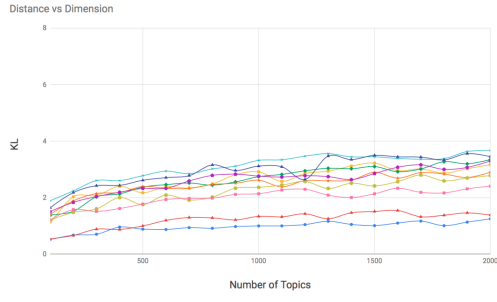
S2JSD was introduced by Endres and Schindelin (2003) to satisfy the triangle inequality. It is the square root of two times the  $JS$  divergence:

$$S2JSD(P, Q) = \sqrt{2 * JS(P, Q)} \quad (2.5)$$

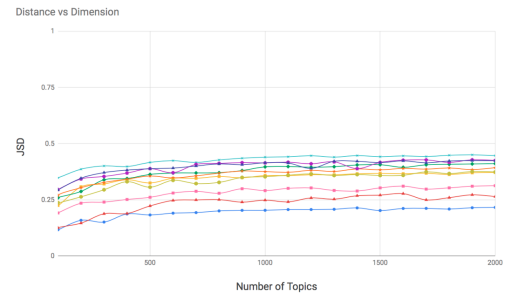
Making sense out of the similarity score is not easy. As shown in figure 2.1, given a set of pairs of documents, their similarity scores vary according to the number of topics. So the distances between those pairs fluctuate from being more to less distant when changing the number of topics.

Distances based on topic distributions between documents generally increase as the number of dimensions of the space increases. This is due to the fact that as the number of topics describing the model increases, the more specific the topics will be. Topics shared by a pair of documents can be broken down into more specific topics that are not shared by those documents. Similarity between pairs of documents is then dependent on the model used to represent them when considering this type of metrics.

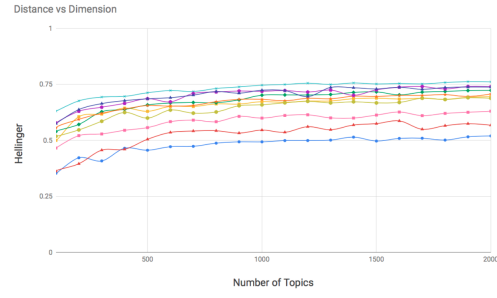
Each topic is drawn from a Dirichlet distribution with parameter  $\beta$ , while each document's mixture is sampled from a Dirichlet distribution with parameter  $\alpha$ . These two priors,  $\alpha$  and  $\beta$ , are also known as hyper-parameters of a topic model and they set the probability that a document or a word, respectively, contains more than one topic. We know that absolute distances between documents vary when we tune those hyper-parameters differently, but we also see that "relative distances" also change. Imagine that we have two documents, A and B, and one topic model, M1. The distance from the topic distribution of A to B is less than from A to C. However, in a second topic model, M2, trained with the same documents as M1 but with different hyper-parameters, the distance from the topic distribution of A to C is less than to B (cross-lines in fig 2.1).



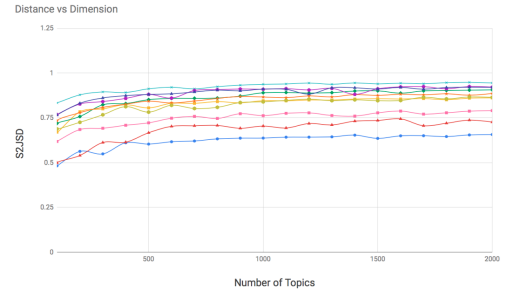
(a)



(b)



(c)



(d)

**Figure 2.1:** Distance values between 10 pair of documents from topic models with 100-to-2000 dimensions. The Kullback-Liebler(a), Jensen-Shannon Divergence(b), Hellinger(c) and S2JSD(c) metrics are considered.

This behaviour highlights the difficulty of establishing absolute similarity thresholds and the complexity to measure distances taking into account all dimensions. Distance thresholds should be model-dependent rather than general and metrics flexible enough to handle dimensional changes.

In addition, document similarity comparisons are too costly to be performed in huge collections of data and require more efficient approaches than having to calculate all pairwise similarities. Using a naive approach creating a similarity matrix with all document comparisons takes  $O(n^2)$  time (where  $n$  is the number of documents), so obtaining all possible pairs of similarities in a large collection of documents (e.g. a corpus of 32 million patents) can be unfeasible because of the exponential cost of comparing every pair of elements. Many different approaches have been proposed to reduce this complexity. For instance, computation can be approximated by a nearest neighbors (ANN) search problem. ANN search is an optimization problem that finds nearest neighbors of a given query in a metric space of  $n$  points. Due to the low storage cost and fast retrieval speed, hashing is one of the most popular solutions for ANN search (Zhen et al., 2016).

This technique transforms data points from the original feature space into a binary-code space, so that similar data points have larger probability of collision (i.e. having the same hash code). This type of formulation for the document similarity comparison problem has proven to yield good results in the metric space (Krstovski and Smith, 2011) due to the fact that ANN search has been designed to handle distance metrics (e.g. cosine, Euclidean, Manhattan). But distance metrics between topic distributions should be information-theoretically motivated metrics (e.g. Hellinger, Kullback-Leibler divergence, Jensen-Shannon divergence) since they compare density functions.

These challenges can be tackled by hashing methods based on clusters of topics to measure similarity, instead of directly using their weights. Hashing methods transform the data points from the original feature space into a binary-code Hamming space, where the similarities in the original space are preserved. They can learn hash functions (data-dependent) or use projections (data-independent) from the training data Wang et al. (2016). Data-independent methods unlike data-dependent ones do not need to be re-calculated when data changes, i.e. adding or removing documents to the collection. Taking large-scale scenarios into account (e.g. Document clustering, Content-based Recommendation, Duplicate Detection), this is a key feature along with the ability to

infer hash codes individually (for each document) rather than on a set of documents. Data-independent hashing methods depend on two key elements: (1) data type and (2) distance metric. For vector-type data, as introduced in section ??, based on  $l_p$  distance with  $p \in [0, 2)$  lots of hashing methods have been proposed, such as p-stable Locality-Sensitive Hashing (LSH) Datar et al. (2004), Leech lattice LSH Andoni and Indyk (2006), Spherical LSH Terasawa and Tanaka (2007), and Beyond LSH Andoni et al. (2014). Based on the  $\theta$  distance many methods have been developed such as Kernel LSH Kulis and Grauman (2012) and Hyperplane hashing Vijayanarasimhan et al. (2014). But only few methods handle density metrics in a simplex space, where topic distributions are projected. A first approach transformed the  $H_e$  divergence into an Euclidean distance so that existing ANN techniques, such as LSH and k-d tree, could be applied Krstovski et al. (2013). But this solution does not consider the special attributions of probability distributions, such as Non-negative and Sum-equal-one. Recently, a hashing schema Mao et al. (2017) has been proposed taking into account the symmetry, non-negativity and triangle inequality features of the S2JSD metric for probability distributions. For set-type data, Jaccard Coefficient is the main metric used. Some examples are K-min Sketch Li et al. (2012), Min-max hash Ji et al. (2013), B-bit minwise hashing Li and König (2010) and Sim-min-hash Zhao et al. (2013).

All of them have demonstrated efficiency in the search for similar documents, but none of them allows the search for documents (1) by thematic areas or (2) by similarity levels, nor they offer (3) an explanation about the similarity obtained beyond the vectors used to calculate it. Binary-hash codes drop a very precious information: the topic relevance. This thesis proposes a *hash function-based approach that allows efficiently searching for related documents while maintaining topic-based annotation, giving the reasons why two documents are related.*

## 2.3 Multilingual Topic Alignment

When the IR task is also cross-language, document retrieval must be independent of the language of the user’s query. At execution time, the query in the source language is typically translated into a target language with the help of a dictionary or a machine-translation system. But for many languages we may not have access to translation

dictionaries or a full translation system, or they can be expensive to apply in an online search system. In such situations it is useful to rely on smaller annotation units derived from the text so the full content does not need to be translated, for instance by finding correspondences with regard to the topics discussed.

They are mainly based on Latent Dirichlet Allocation (LDA) Blei et al. (2003), adding supervised associations between languages by using *parallel* corpus, with sentence-aligned documents (e.g. Europarl<sup>19</sup> corpora), or *comparable* corpus, with theme-aligned documents (e.g. Wikipedia<sup>20</sup> articles), in multiple languages. These requirements restrict the kind of corpora that can be used for training since large parallel corpora are rare in most of the use cases, especially for languages with fewer resources. Wikipedia, for example, contains texts in 304 languages but 255 of them have less than 3% of articles<sup>21</sup>. Therefore, the requirement of parallel/comparable corpora for multilingual topic models limits their usage in many situations. In addition, these models rely on associations between documents prior to training. So in order to incorporate new languages or update the existing associations, the model must be re-trained with documents from all languages, making it difficult to scale to large corpora Hao et al. (2018) Moritz and Èchler (2017).

Multilingual probabilistic topic models (MuPTM) Vulić et al. (2015) have recently emerged as a group of language-independent generative machine learning models that can be used on large-volume theme-aligned multilingual text. They are based on LDA, adding supervised associations between languages by using *parallel* corpus, with sentence-aligned documents (e.g. Europarl<sup>22</sup> corpora), or *comparable* corpus, with theme-aligned documents (e.g. Wikipedia<sup>23</sup> articles), in multiple languages. Once a MuPTM has been generated, documents can be represented by data points in a single feature space based on topics to detect similarities among them exploiting inference results and using distance metrics. Due to its generic language-independent nature and the power of inference on unseen documents, MuPTM's have found many interesting applications in many different cross-lingual tasks. They have been used on cross-lingual event clustering De Smet and Moens (2009), document classification

---

<sup>19</sup><https://ec.europa.eu/jrc/en/language-technologies/dcep>

<sup>20</sup><https://www.wikipedia.org/>

<sup>21</sup>[https://meta.wikimedia.org/wiki/List\\_of\\_Wikipedias](https://meta.wikimedia.org/wiki/List_of_Wikipedias)

<sup>22</sup><https://ec.europa.eu/jrc/en/language-technologies/dcep>

<sup>23</sup><https://www.wikipedia.org/>



De Smet et al. (2011) Ni et al. (2011), semantic similarity of words Mimno et al. (2009) Vulić and Moens (2012), information retrieval Vulić and Moens (2013) Ganguly et al. (2012), document matching Platt et al. (2010) Zhu et al. (2013), and others.

However, the requirement of parallel/comparable corpora limits their usage in many situations. There are not many document collections that can be used for training since large parallel corpora are rare in most of the use cases, especially for languages with fewer resources. Moreover, in order to incorporate new languages or update the existing associations, these models must be re-trained with documents from all languages at the same time, making it difficult to scale to large corpora Hao et al. (2018) Moritz and Îchler (2017). We take MuPTM a step further, to make them cross-lingual through representations based on topic hierarchies. Documents from multi-language corpora are described by expressions of multi-lingual concepts and can then be efficiently browsed and related without the need for translation or parallel or comparable corpora. In this thesis we propose to *automatically learn cross-lingual topics to browse multi-lingual document collections*.



## Chapter 3

# Research Objectives

The work presented in this thesis aims to facilitate the exploration of huge collections of multilingual documents through thematic associations inferred from their content. Each of the challenges arising from this objective defines a working dimension and guides the research carried out in this thesis.

The first dimension focuses on **scalability**, in order to create the text processing flows that are required to create or apply learning models. The workload required to process a corpus varies according to the number of documents, the length of texts and the kind of knowledge (annotations) that need to be inferred from the text. If the design of the workflow is scalable, there is no need to modify the processing logic when working with larger collections of documents, since adding a reasonable amount of computational resources is enough to perform it. These resources can be machines (i.e horizontal scaling) or processing units (e.g CPU, RAM) in an existing machine (i.e vertical scaling).

The second dimension covers the **representativeness** of the text annotations when projected into spaces where they are manipulated. The idea behind these spaces is to represent documents as points (or vectors in a vector space) that are close together when the texts are semantically similar, and far apart when they are semantically distant. The ability of these spaces to create meaningful representations is also studied in this work.

In the third dimension, data structures that efficiently **sort** texts from their representations based on probabilistic topics are studied. Divisions of space into semantically-related regions are convenient to allow browsing large document collections. The *rep-*

*representativeness* covered in the previous dimension enables the interpretation of the relations and regions obtained.

And finally, the fourth dimension handles the **multilingualism** of collections that contain documents in several languages. On a multilingual space, documents are described and related across languages.

This chapter introduces our main hypothesis (section 3.1), and the associated research challenges (section 3.2), and presents the research methodology (section 3.3).

### 3.1 Research Hypotheses

We define our main hypothesis as follows:

**Hypothesis 1** *Large multilingual document collections can be automatically analyzed to discover appropriate thematic relations that facilitate a semantically-enabled text browsing.*

Our hypothesis can be divided into four different sub-parts, which are related to the aforementioned scalability, representativeness, sorting, and multilinguality dimensions respectively. First, by *distributing both natural language processing tasks and representational models we can efficiently process huge collections of documents (H1.1).*

Second, we can *semantically relate documents by comparing their most relevant topics (H1.2).* Furthermore, for this purpose we hypothesize that the use of *topic hierarchies (H1.2.1)* and *similarity metrics based on relevance levels (H1.2.2)* can help quantifying the semantic distance between texts. Third, by *dividing the representational space into regions based on topics and relevance levels we can search for related documents without having to calculate all pairwise comparisons and without losing the ability to rely on topics for further processing (H1.3).*

And finally, *by abstracting the topic representations into concept-based descriptions across languages we can relate documents in various languages without having to translate them (H1.5).*

A summary of the hypotheses and how they tackle our research dimensions can be found in Table 3.1.

Hypothesis	Research Dimension
H1: Large multilingual document collections can be automatically analyzed to be semantically-browsed through thematic relations	D1: Scalability, D2: Representativeness, D3: Sorting, D4: Multilingualism
H1.1: it is possible to efficiently annotate documents on a large scale by distributing natural language processing tasks and representation models	D1: Scalability
H1.2: it is possible to semantically relate texts from their most relevant topics	D2: Representativeness
H1.3: it is possible to find documents with similar topic distributions without calculating all pairwise comparisons and without losing the ability to explore them through their topics	D3: Sorting
H1.4: it is possible to relate documents in different languages without having to translate them using language agnostic concepts from their main topics	D4: Multilingualism

**Table 3.1:** Hypotheses and research dimensions.

## 3.2 Research Challenges

Several research challenges emerge from these hypotheses. First, in order to facilitate reusing existing topic models by processing systems with different architectures and technological stacks, we need to define *topic-model programming interfaces*. Second, in order to describe and thematically relate documents, we must address how to produce *explainable topic-based associations*. Third, by working with huge collections of documents described by topics, we need to handle *large-scale comparisons of topic distributions*. Finally, in order to explore multilingual document collections from shared topic-based representational spaces, we have to provide *automatic cross-lingual topic alignment*. Each of these research challenges are described below and covered throughout this thesis.

### 3.2.1 Topic-model Programming Interface

Although some initiatives to standardize the format of machine-learning models and to provide tools that facilitate their transformation among the most widespread proprietary formats already exist in the literature, there are still some software restrictions that can limit their reuse. These models may hold certain software dependencies that e.g. force using a specific version of a programming language (python2 vs python3<sup>24</sup>) or an operating system (e.g., linux kernel vs on-cloud environments<sup>25</sup>) to load them or to launch the service that deploys them (e.g., ONNX<sup>26</sup>). This limits their ability to be reused in domains that are not familiar with these technological stacks. *Integrating pre-trained topic models into general-purpose systems is not easy (RCInterface1).*

Topic models, as many other machine learning models, may be distributed in a proprietary or standard format with software dependencies or by directly providing the data. However, *there is no standard way to specify the topics and the operations that can be performed on them (RCInterface2).* Sometimes topics are described by the top ten or five most relevant words, and occasionally these word lists are not accompanied by weights, making a density-based analysis impossible. These differences in presenting the models can sometimes limit their reusability if they cannot infer new topic distributions even when the learning algorithm allows for it.

---

<sup>24</sup><https://www.python.org>

<sup>25</sup><https://vespa.ai>

<sup>26</sup><https://onnx.ai>

### 3.2.2 Explainable Topic-based Associations

In order to facilitate the exploration of document collections, vector space models are often used to semantically relate texts based on their word distributions. These models first create a dictionary with the words used in the collection, and then represent documents by vectors whose dimensions correspond to each word in the dictionary. In large collections, these models need to be adapted to make operations on vectors more manageable. As a result, a new abstraction method based on topics emerged that reduces the dimensions of vectors. Topics are described by word distributions over the entire vocabulary and documents by vectors containing topic distributions. Despite the extensive use of these representation models, *there is no common criteria for identifying the most representative topics in a document (RCExplainable1)*.

In addition, since similarity metrics over this representation space are based on accumulating the difference in topic densities, *it is difficult to explain the distance between topic distributions (RCExplainable2)*. And, unless a minimum distance threshold is defined or a n-top topics agreed, *there is no common criterion for determining whether two documents are related (RCExplainable3)*.

### 3.2.3 Large-scale Comparisons of Topic Distributions

There are many scenarios where finding related documents in a huge corpus is desirable (e.g. a researcher doing literature review, or an R&D manager analyzing project proposals). Experts can benefit from discovering those connections to achieve these goals, but brute-force pairwise comparisons are not computationally adequate when the size of the corpus is too large. Some algorithms in the literature divide the search space into regions containing potentially similar documents, which are later processed separately from the rest in order to reduce the number of pairs compared. However, *there are no mechanisms that efficiently partition the topic-based search space without compromising the ability for thematic exploration (RCComparison1)*.

In addition, documents from the same region should be compared and *there are no similarity metrics that compare partial distributions of topics (RCComparison2)*.

### 3.2.4 Automatic Cross-lingual Topic Alignment

With the ongoing growth in the number of texts in different languages, we need annotation methods that enable browsing multi-lingual corpora. As discussed in section 2, multilingual probabilistic topic models have recently emerged as a group of semi-supervised machine learning models that can be used to perform thematic explorations on collections of texts in multiple languages. However, *there are no approaches that abstract the representation of probabilistic topics in language-independent spaces without translating texts or aligning documents (RCCrossLingual1)*. Existing approaches require parallel or comparable training data to create a language-independent space.

A summary of the challenges covered in this work and how they map to the hypotheses is presented in table 3.2.

## 3.3 Research Methodology

The research presented in this thesis is based on four dimensions or research areas as discussed in section 3.2. Each one is motivated by different research problems that we need to solve in order to achieve our ultimate goal of making it easier to explore large multilingual document collections through their topics. Once a dimension is tackled, the next one is considered, and so on. This iterative and incremental methodology allows refining the research results by evaluating them with more experiments and addressing increasingly complex research problems.

Figure 3.1 shows the dimensions on which the research of this thesis has been built. The top of the pyramid is only reached once the lower dimensions are dealt with successfully. They are presented as a chain of four steps. The first step describes the motivation to perform a given task coming from real-world problems that we had to deal with, and is represented by a brown arrow. In the context of this task, the research problem arises and is framed by a pink arrow. For each of them a solution is proposed and evaluated according to a specific criterion. The proposed solution is represented by a green arrow and the evaluation with a blue arrow. Once a proposal has been validated, the next dimension of the pyramid is achievable and all the previous research problems are added to the new research problem as conditions to be taken into account.



Research Challenge	Hypotheses
RCInterface1: integrating pre-trained topic models into general-purpose systems is not easy	H1.1: documents can be efficiently annotated on a large scale by distributing natural language processing tasks and representation models
RCInterface2: there is no standard presentation of topics that facilitates their reuse	H1.1: documents can be efficiently annotated on a large scale by distributing natural language processing tasks and representation models
RCExplainable1: there is no common criteria for identifying the most representative topics in a document	H1.2: texts can be semantically related from their most relevant topics, H1.3: documents with similar topic distributions can be found without calculate all pairwise comparisons and without losing the ability to explore them through their topics
RCExplainable2: it is difficult to understand the distance between topic distributions	H1.2: texts can be semantically related from their most relevant topics
RCExplainable3: there is no common criterion for determining whether documents are related	H1.2: texts can be semantically related from their most relevant topics
RCComparison1: there are no mechanisms that efficiently partition the topic-based search space without compromising the ability for thematic exploration	H1.3: documents with similar topic distributions can be found without calculate all pairwise comparisons
RCComparison2: there are no similarity metrics that compare partial distributions of topics	H1.3: documents with similar topic distributions can be found without calculate all pairwise comparisons
RCCrossLingual1: there are no approaches to abstract probabilistic topics in language-independent spaces without translating texts or aligning documents	H1.4: documents in different languages can be related without having to translate them using language agnostic concepts from their main topics

**Table 3.2:** Open Research Challenges and Hypotheses.

Technical objectives (i.e., develop a new resource) or research objectives (i.e., discover the solution to a problem) guide the solution proposal before moving on to the next dimension. They are presented below, organized by the research problem associated with each dimension.

### 3.3.1 Scalable Creation and Inference of Topics

This first dimension arose when we had to analyze a huge collection of documents describing research and innovation projects to discover which research areas are being addressed, measure their presence in the collection, and characterize them so their presence can be inferred in unseen documents. Such a high volume of data made difficult to process it manually, so we needed to automate the required processing to draw insights from it. Probabilistic topics allow describing research areas, so we defined a *distributed text-processing model for creating large probabilistic topic models (RO1)* and a *web service template to distribute them (RO2)*. In this way, the models themselves could be easily integrated into scalable text processing pipelines. As a result, we created a *platform for large-scale text analysis (TO1)*, and produced a *model-as-a-service repository with pre-trained topic models (TO2)*. The efficiency of this solution was validated by processing a corpus of 100,000 documents collected from the CORDIS dataset<sup>27</sup>, which contains descriptions of projects funded by the European Union under a framework programme since 1990 (Badenes-Olmedo et al., 2017b).

The main contributions under this dimension are described in Chapter 4 as follows:

- a software architecture to process big volumes of textual documents in a distributed and decoupled manner;
- the definition of a model-as-a-service template for probabilistic topic models;
- an implementation of the architecture, libRAIry, following those design principles.

### 3.3.2 Explainable Topic-based Associations

In the second dimension we needed to browse scientific papers through their content-based relations. The problem of massively annotating documents with topic distributions came up. We had to *create annotations based on topic models in a way that*

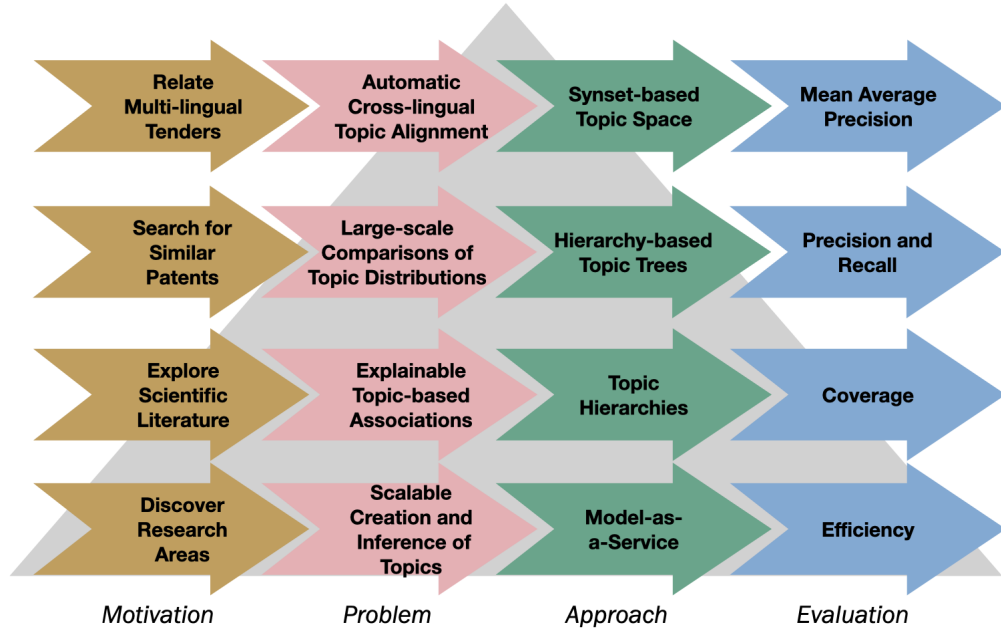
---

<sup>27</sup><https://data.europa.eu/euodp/es/data/dataset/cordisH2020projects>

was computationally affordable and enabled a semantic-aware exploration of the knowledge inside them (**RO3**). Once documents were annotated, a metric that compares documents and facilitates their interpretation from topic annotations (**RO4**) was required. As a result, we integrated the annotation method into the topic model service (**TO3**) and implemented a text comparison metric based on partial representations of topics. These proposals were validated by classifying 500,000 scientific articles from the Open Research Corpus<sup>28</sup> in domains such as Computer Science, Neuroscience and Biomedicine (Badenes-Olmedo et al., 2017a, 2019a, 2017c).

The main contributions under this dimension are described in Chapter 5 as follows:

- a clustering algorithm based on probabilistic topic distributions;
- a hash function to transform topic distributions into topic hierarchies;
- a similarity metric based on topic sets.



**Figure 3.1:** Research dimensions of the thesis. The first ones must be overcome before reaching higher dimensions.

<sup>28</sup><https://allenai.org/data/open-research-corpus>

### 3.3.3 Large-scale Comparisons of Topic Distributions

This dimension covered the search for similar documents based on their most relevant topics. Thanks to having dealt with the above two dimensions, large collections of documents could be annotated with topic hierarchies and text distances could be measured from their annotations. Now, the aim was to find similar documents without losing the exploratory capacity offered by topics. Similarity comparisons were too costly to be performed in such huge collections of data and required more efficient approaches than having to calculate all pairwise similarities. We applied *techniques based on approximate nearest-neighbors to organize documents in regions with similar topic hierarchies (RO5)*. As a result, we developed *a system to automatically find similar documents (TO4)*. It was validated on a collection of one million texts retrieved from the United States patents corpus<sup>29</sup>. The relations between patents derived from their manual categorization were compared with those automatically obtained from their topic distributions (Badenes-Olmedo et al., 2019a, 2020).

The main contributions under this dimension are described in Chapter 6 as follows:

- a data structure to partition the search space and organize documents described by topic hierarchies;
- a corpus browser that leverages these representations to automatically relate documents.

### 3.3.4 Automatic Cross-lingual Topic Alignment

Finally, a new dimension on top of the previous ones emerged to relate texts coming from different languages. In particular, since document relations were based on their topics, this dimension was focused on aligning topics without supervision from models trained with texts in different languages. Since each language defined its own vocabulary, the topics were model-specific and could not be directly compared. We abstracted the *topic representations to create a single space out of the particularities of the language (RO6)*. This approach was validated on the English, Spanish, French, Italian and Portuguese editions of the JCR-Acquis<sup>30</sup> corpora and revealed promising results on

---

<sup>29</sup><https://www.uspto.gov/ip-policy/economic-research/research-datasets>

<sup>30</sup><https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis>

classifying and sorting documents by similar content across languages (Badenes-Olmedo et al., 2019a,b).

The main contributions under this dimension are described in Chapter 7, as follows:

- an algorithm to represent probabilistic topics using concept sets;
- a repository of aligned topic models from the English, Spanish, French, Italian and Portuguese editions of the JRC-Acquis corpus.

Table 3.3 summarizes the research objectives (ROs), technical objectives (TOs) and connects them with the research challenges (RCs) from Table 3.2.

Research Objective	Research Challenge
RO1: Define a distributed text-processing model for creating large probabilistic topic models	RCInterface1
RO2: Define a template to package probabilistic topic models as web services	RCInterface2
RO3: Define annotations based on topics that enable a semantic-aware exploration of the knowledge inside a corpus	RCExplainable1
RO4: Define a metric based on topic annotations that compares documents and facilitates their interpretation	RCExplainable2, RCExplainable3
RO5: Define nearest-neighbor techniques to organize documents in regions with similar topic hierarchies	RCComparison1, RCComparison2
RO6: Define a transformation of the topic-based annotations to create a unique representational space out of the particularities from each language	RCCrossLingual1
TO1: Create a platform for large scale text processing	RCInterface1, RCInterface2
T02: Create a repository of Topic-based web services	RCInterface2
T03: Integrate the annotation method based on topic hierarchies into the topic model service	RCExplainable2, RCComparison2
T04: Create a system capable of finding similar document automatically	RCExplainable2, RCExplainable3, RCComparison1, RCCrossLingual1

**Table 3.3:** Research and technical objectives and their related challenges.

## Chapter 4

# Scalable Creation and Inference of Topics

This chapter presents *librAIry*<sup>31</sup>, our text management platform that combines natural language processing techniques with automatic learning algorithms to analyze large collections of documents. It serves as a technological framework where we can implement the advances of our research and measure its performance.

### 4.1 Document Workflow

Given the huge amount of textual data about any domain that is daily being produced or captured in any imaginable domain, it becomes crucial to provide mechanisms for programmatically processing this raw data so we can make sense out of it: discarding all the noisy, non-relevant information and keeping only the data that can bring value for the involved agents (general consumers, experts, companies, investors...).

While some specific tools already allow for advanced sense-making operations, others opt for composing a solution where different analysis techniques are integrated under a uniform data schema. However, this integration involves significant efforts on reconciling data sources, coordinating processing operations, and efficiently exploiting results from the execution of those techniques. There is the need for a more flexible paradigm where tools and algorithms for textual document analysis, from different programming languages and technologies, can operate independently and in a collaborative manner

---

<sup>31</sup><http://librairy.linkeddata.es>

creating a common document oriented workflow through their actions. In the context of the scientific publications, the personalized recommendation of research papers based on their content is a key novel feature for performing a smart selection of relevant resources over very big collections of scientific content. From the set of values and different attributes extracted from the papers and by generating advanced knowledge models about the information they contain we can bridge across the different relevant pieces of information and allow users to navigate them in a more efficient and powerful way. This knowledge about a specific document is frequently acquired by different techniques focused on revealing certain aspects of it, that are later combined to achieve one particular task.

The architecture presented in this thesis aims to ease the way different software modules work together and lays the foundation for efficiently process big volumes of textual documents in a distributed, decoupled manner.

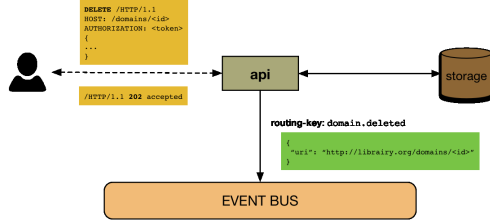
## 4.2 librAIry

*librAIry* is a framework where different text mining tools, available in various languages and technologies, can operate in a distributed, high-performance and isolated manner creating a common workflow through their actions. Instead to work towards a pre-defined sequence of actions, synchronization across modules is achieved through the aggregation of the operations executed by them in response to an emergent chain of events. This raises both technical and functional challenges to coordinate multiple executions. From the technical point of view, isolated environments and communication mechanisms are provided so initially dissimilar tools can be executed with maximum guarantees. From the functional point of view, all executions are coordinated to reach a final result as aggregation of partial results derived from each execution.

### 4.2.1 Functional Features

The architecture is articulated around three main concepts: (1) the **resource** such as *document*, a *part-of-a-document*, or a *domain*. (2) the **actions** performed over them: *create*, *update* or *delete* a resource. And (3) the new **state** that is reached by the resource after an action is performed, such as *created*, *updated* or *deleted*. An **event** is a message containing details about those three aspects, published on a shared





**Figure 4.1:** Domain deleted flow.

event-bus available for all the modules deployed in the framework. This will, in turn, allow that any module can perform actions on one or more resources in response to a new state reached by a given resource. Actions executed in parallel from distributed environments.

#### 4.2.1.1 Resources

Two main kinds of resources are considered: those derived from external sources such as (1) *documents* from textual files (e.g. a research paper), (2) *parts* from logical divisions of a *document* (e.g. rhetorical classes or sections), and (3) *domains* from sets of *documents* (e.g. a conference or journal), and those derived from processing the previous ones such as *annotations*.

To better illustrate this model, consider to explore the research papers published at the SIGGRAPH conference in 2016<sup>32</sup>. First, every paper will be materialized as a new *document* containing the full-text. Immediately after, the *document* will be automatically associated to several *parts*, each of them grouping sentences by rhetorical class (e.g. approach, background, challenge, future work and outcome) and by section (e.g. abstract, introduction). Finally, a new *domain* will be created grouping all these *documents*. Different analysis will be performed extending the initial set of resources with more annotations at several representational levels: at *document level*, full-text based annotations are provided such as named-entities, compounds and descriptive tags. At *relational level*, connection between resources are found (e.g. semantic similarity-based relationships). And finally, at *domain level* annotations such as tags and summaries are composed describing the corpus of *documents*.

<sup>32</sup><http://s2016.siggraph.org>

#### 4.2.1.2 Event-based Paradigm

An event illustrates a performed action, i.e. a resource and its new state. It follows the Representational State Transfer (REST) Fielding and Taylor (2002) paradigm, but taking into account the state reached after an action, i.e. *created*, *deleted* or *updated*. Thus, an event contains the resource type and the new state reached by a specific resource.

#### 4.2.1.3 Linked Data Principles

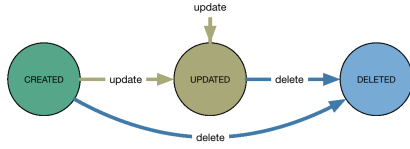
Data in *librAIry* is individually addressable and linkable Turchi et al. (2012) following the Linked Data principles defined by T. Berners-Lee Bizer et al. (2009). Thus, resources (i.e. a *domain*, a *document*, a *part* or an *annotations*) have: (1) a URI as name, (2) a retrievable (or dereferenceable) HTTP URI so that it can be looked up, (3) a useful information provided by using standard notation (e.g. JavaScript Object Notation (JSON)) when it is looked up by URI, and (4) links to other URIs so that other resources can be discovered from it.

### 4.2.2 Framework Architecture

Following a publisher/subscriber approach, all the modules in the framework can publish and read events to notify and to be notified about the state of a resource. Therefore, the system flow is not unique and is not explicitly implemented, instead distributed and emergent flows can appear according to particular actions on resources.

#### 4.2.2.1 Event-Bus

We use the Advanced Message Queuing Protocol (AMQP) as the messaging standard in *librAIry* to avoid any cross-platform problem and any dependency to the selected message broker. This protocol defines: *exchanges*, *queues*, *routing-keys* and *binding-keys* to communicate publishers and consumers. A message sent by a publisher to an exchange is tagged with a routing-key. Consumers matching that routing-key with the binding-key used to link the queue to that exchange will receive the message. In *librAIry* this key follows the structure: *resource.status*. Since a wildcard-based definition can be used to set the key, this paradigm allow modules both listening to individual type events



**Figure 4.2:** Resource states.

(e.g. `domains.created` for new domains), or multiple type events (e.g. `#.created` for all new resources).

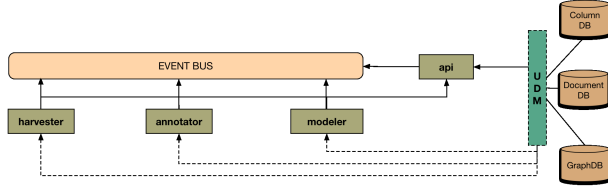
#### 4.2.2.2 API

A HTTP-Rest Application Program Interface (API) was designed for interaction with end-users. Any external operation motivated by a user will be handled here. Some of them, usually those related to reading operations, will be completely managed by this module getting all the data from the internal storage. However, those operations implying a modification of the status of some resource (e.g. creation of a *document*), may be also performed by other modules listening for that type of event asynchronously. This module publishes to the following routing-keys: *domain.(created;updated;deleted)*, *document.(created;updated;deleted)*, *part.(created;updated;deleted)*, and *annotation.(created;updated;deleted)*.

#### 4.2.2.3 Storage

Multiple types of data can be handled in this ecosystem. Inspired in the Data Access Object (DAO) pattern, we have created a Unified Data Manager (UDM) providing access to any type of data used in the system. Three types of databases have been considered:

- **column-oriented database:** Focused on unique identified and/or *structured data*. This storage allow us searching key elements across resources.
- **document-oriented database:** Focused on indexing raw text. This storage allow us to execute advanced search operations over all the information gathered about a textual resource.
- **graph database:** Focused on relations. This storage allow us exploring resources through the relationships between them.



**Figure 4.3:** Modules.

#### 4.2.2.4 Modules

The modules composing *libRAIry* have been designed following the microservices architectural style. A module is a cohesive (i.e. it implements only functionalities strongly related to the concern that it is meant to model Dragoni et al. (2016)) and independent process working on the framework with a specific purpose. This purpose is defined by both the routing-key and the binding-key associated to the events handled by the module.

These are the main types of modules identified in *libRAIry*:

- **Harvester:** creates system resources such as *documents*, *parts* and *domains*, from local or remote located textual files.
  - Listening for: nothing
  - Publishing to: *document.(created)*,  
*part.(created)*, *domain.(created;updated)*
- **Annotator:** retrieves named-entities, compounds, lemmas and other annotations resulting of Natural Language Processing (NLP) task execution from *documents* and *parts*.
  - Listening for: *document.(created;updated)*,  
*part.(created;updated)*
  - Publishing to: *annotation.(created;deleted)*
- **Modeler:** builds representational models from a given *domain*.
  - Listening for: *domain.(created;updated)*
  - Publishing to: *annotation.(created;deleted)*

### 4.3 Model-as-a-Service

...

### 4.4 Summary

In *librAIry*, existing algorithms and tools coming from different technologies can work collaboratively to process and analyze large collections of textual resources which has been successful applied to some real scenarios <sup>33</sup>.

A new model definition based on the previously mentioned principle of maximizing information re-usability and minimize irrelevant data is being studied to create a more fine-grained resource design. New domains, in the sense of particular vocabularies or specific textual formats, are also being analyzed to be included into the system via specific harvesters or more precise annotators. Moreover, a template-based mechanism oriented to facilitate the integration of new tools and techniques into the system is being built to make easier to develop new modules as well as increasing the available modules at Docker-Hub.

---

<sup>33</sup><http://drinventor.dia.fi.upm.es>



## Chapter 5

# Explainable Topic-based Associations

### 5.1 Topic Relevance

### 5.2 Topic-based Clustering

### 5.3 Summary





## Chapter 6

# Large-scale Comparisons of Topic Distributions

### 6.1 Document Similarity

### 6.2 Hashing Topic Distributions

### 6.3 Summary



## Chapter 7

# Cross-lingual Document Similarity

### 7.1 Synset-based Representational Space

### 7.2 Cross-lingual Models

### 7.3 Summary



## Chapter 8

# Evaluation

### 8.1 Evaluation Metrics

### 8.2 Text Representativeness

### 8.3 Large-scale Text Processing

*librAIry* has been used in some real scenarios such as a research-paper repository for the European project DrInventor<sup>34</sup>, a support to decision makers for analyzing patents and public aids for the ICT sector, and also as a book recommender for an online content platform. This has allowed us to identify some weak and strong points of the framework and iterate over the architecture to come with the described solution.

The following modules have been developed<sup>35</sup>: (1) a ***general-purpose harvester*** which retrieves text and meta-information from PDF files in local or remote file-system; (2) a ***research paper-oriented harvester*** focused on collecting and processing more specific textual files (e.g. scientific papers) creating both *documents* and *parts* inferred from the rhetorical classes of the paper; (3) a ***Stanford CoreNLP-based Annotator*** which discovers named-entities, compounds and lemmas from *documents* and *parts*; (4) a ***Topic Modeler*** based on Latent Dirichlet Allocation (LDA) which creates probabilistic topic models for each *domain* in the framework. They are annotated with the set of topics (i.e. ranked list of words) discovered from the corpus, and both *documents* and *parts* of that domain are also annotated by the vector of probabilities to belong

---

<sup>34</sup><http://drinventor.eu>

<sup>35</sup><https://github.com/librairy>

to these topics. It uses the Spark implementation of the algorithm; and (5) a **Word Embedding Modeler** which creates a *word2vec* model from the *documents* contained in a *domain*.

Due to linear scalability and high performance features, Cassandra has been used to support the column-oriented storage functionality, Elasticsearch as document-oriented storage and Neo4j as graph-oriented storage.

All modules in *librAIry* have been packaged as Docker <sup>36</sup> containers and uploaded to Docker-Hub <sup>37</sup> to facilitate the installation of the system.

Maximizing information re-usability and minimize irrelevant data, becomes specially important when the system handles large collections of data (around million of documents). Fine-grained resource definitions have been key to achieve this, so modules execute actions only when really necessary. When a new *domain* is created, for instance, a new Topic Model is trained for that *domain* and is used to calculate the semantic similarity between the *documents* (and the *parts*) in that domain. If a new *document* (or *part*) is added to that *domain*, the model is trained again and the semantic similarities are re-calculated. However, this becomes unfeasible when the domain is frequently updated and it is composed by a large number of documents. One solution has been to define a new type of resource between domains and documents, models, that describes the representational state (e.g. topic model) of a collection of documents. Thus the model is only re-trained when a significant amount of *documents* are added to the sampling data set and not to the entire *domain*. This less transient model is used to calculate semantic similarities between the *document* collection (and *parts*) inside a *domain* in a more efficient way. Following this more precise execution of tasks, the routing-keys should include the URI of the implied resource into the definition, not only in the content of the message. It would allow modules listening to both the type of a resource or to a specific resource (or subsets, via regular expressions).

While the storage modules are always used to save/update/delete a resource, they are not always required from the end-user. The graph storage, for instance, makes sense when a path between two *documents* or *parts* is requested for a given *domain*. However, some *domains* are not intended to be explored by their linked resources. A

---

<sup>36</sup><https://www.docker.com>

<sup>37</sup><https://hub.docker.com/u/librairy/>

more fine/grained definition of resources will allow graph-storage being only used when necessary.

On the other hand, distributed execution of NLP tasks (not only in threads, but also in machines) has proved to be especially useful to handle large collection of *documents*. It requires less processing time than a monolithic solution (e.g. CoreNLP application) and it also provides a dynamic load balancing between modules.

## **8.4 Topic-based Clustering**

## **8.5 Cross-lingual Similarity**

## **8.6 Conclusions**





## Chapter 9

# Experiments

### 9.1 Polypharmacy and Drug-drug Interactions

...

### 9.2 Corpus Viewer

...

### 9.3 ODS Classifier

...

### 9.4 Drugs4Covid

...



## Chapter 10

# Conclusions

### 10.1 Assumptions and Restrictions

...

### 10.2 Contributions

...

### 10.3 Impact

...

### 10.4 Limitations

...

### 10.5 Future Work

...



# Bibliography

- Agerri, R., Bermudez, J., and Rigau, G. (2014). IXA pipeline: Efficient and Ready to Use Multilingual NLP tools. In *Proceedings of the 9th Language Resources and Evaluation Conference (LREC)*, pages 26–31, Reykjavik, Iceland. 10
- Andoni, A. and Indyk, P. (2006). Near-Optimal Hashing Algorithms for Approximate Nearest Neighbor in High Dimensions. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, pages 459–468. IEEE. 17
- Andoni, A., Indyk, P., Nguyá»...n, H. L., and Razenshteyn, I. (2014). Beyond Locality-Sensitive Hashing. In *Proceedings of the Twenty-Fifth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1018–1028. 17
- Badenes-Olmedo, C., Redondo-García, J., and Corcho, O. (2017a). An Initial Analysis of Topic-based Similarity among Scientific Documents Based on their Rhetorical Discourse Parts. In Garijo, D., van Hage, W., Kauppinen, T. and Kuhn, T., and Zhao, J., editors, *Proceedings of the First Workshop on Enabling Open Semantic Science co-located with 16th International Semantic Web Conference (ISWC)*, volume 1931 of *CEUR Workshop Proceedings*, pages 15–22. CEUR-WS.org. 29
- Badenes-Olmedo, C., Redondo-Garcia, J., and Corcho, O. (2017b). Distributing Text Mining tasks with librAIry. In *17th ACM Symposium on Document Engineering (DocEng)*. 28
- Badenes-Olmedo, C., Redondo-García, J., and Corcho, O. (2019a). Legal document retrieval across languages: topic hierarchies based on synsets. *arXiv e-prints*, page arXiv:1911.12637. 29, 30, 31

- Badenes-Olmedo, C., Redondo-García, J., and Corcho, O. (2019b). Scalable cross-lingual document similarity through language-specific concept hierarchies. In *Proceedings of the 10th International Conference on Knowledge Capture*, pages 147–153. 31
- Badenes-Olmedo, C., Redondo-García, J., and Corcho, O. (2020). Large-Scale Semantic Exploration of Scientific Literature using Topic-based Hashing Algorithms. *Semantic Web*. 30
- Badenes-Olmedo, C., Redondo-Garcia, J. L., and Corcho, O. (2017c). Efficient Clustering from Distributions over Topics. In *9th International Conference on Knowledge Capture (K-CAP)*, page 8. 29
- Bizer, C., Heath, T., and Berners-Lee, T. (2009). Linked data-the story so far. *International journal on Semantic Web and Information Systems*, 5(3):1–22. 36
- Blei, D., Carin, L., and Dunson, D. (2010). Probabilistic topic models. *IEEE Signal Processing Magazine*, 27(6):55–65. 11
- Blei, D. M. and Lafferty, J. D. (2007). A correlated topic model of Science. *The Annals of Applied Statistics*, 1(1):17–35. 13
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3(4-5):993–1022. 12, 18
- Boyd-Graber, J. and Resnik, P. (2010). Holistic Sentiment Analysis Across Languages: Multilingual Supervised Latent Dirichlet Allocation. *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, (October):45–55. 13
- Celikyilmaz, a., Hakkani-Tur, D., and Tur, G. (2010). LDA Based Similarity Modeling for Question Answering. In *Proceedings of the NAACL HLT 2010 Workshop on Semantic Search*, pages 1–9. 13
- Charikar, M. S. (2002). Similarity estimation techniques from rounding algorithms. In *Proceedings of the thirty-fourth annual ACM symposium on Theory of computing - STOC '02*, page 380. ACM Press. 14

- Cheng, X., Yan, X., Lan, Y., and Guo, J. (2014). BTM : Topic Modeling over Short Texts. *IEEE Transactions on Knowledge and Data Engineering*, 26(12):2928–2941. 11
- Dagan, I., Lee, L., and Pereira, F. C. N. (1999). Similarity-Based Models of Word Cooccurrence Probabilities. *Machine Learning*, 34(1-3):43–69. 13
- Datar, M., Immorlica, N., Indyk, P., and Mirrokni, V. S. (2004). Locality-sensitive hashing scheme based on p-stable distributions. In *Proceedings of the twentieth annual symposium on Computational geometry - SCG '04*, page 253. ACM Press. 17
- De Smet, W. and Moens, M.-F. (2009). Cross-language linking of news stories on the web using interlingual topic modelling. In *Proceedings of the 2nd ACM workshop on Social web search and mining*, page 57. 18
- De Smet, W., Tang, J., and Moens, M.-F. (2011). Knowledge Transfer across Multilingual Corpora via Latent Topics. In *Advances in Knowledge Discovery and Data Mining*, pages 549–560. 18
- Deerwester, S., Dumais, S. T., Furnas, G. W., Landauer, T. K., and Harshman, R. (1990). Indexing by Latent Semantic Analysis. *Journal of the American Society for Information Science*, 41:391–407. 12
- Divoli, A., Nakov, P., and Hearst, M. A. (2012). Do peers see more in a paper than its authors? *Advances in Bioinformatics*. 9
- Dragoni, N., Giallorenzo, S., Lafuente, A. L., Mazzara, M., Montesi, F., Mustafin, R., and Safina, L. (2016). Microservices: yesterday, today, and tomorrow. *CoRR*, abs/1606.0:1–17. 38
- Endres, D. and Schindelin, J. (2003). A new metric for probability distributions. *IEEE Transactions on Information Theory*, 49(7):1858–1860. 14
- Fielding, R. T. and Taylor, R. N. (2002). Principled Design of the Modern Web Architecture. *ACM Transactions on Internet Technology*, 2(2):407–416. 36
- Ganguly, D., Leveling, J., and Jones, G. (2012). Cross-Lingual Topical Relevance Models. In *Proceedings of COLING 2012*, pages 927–942. 19

- Gatti, C. J., Brooks, J. D., and Nurre, S. G. (2015). A Historical Analysis of the Field of OR/MS using Topic Models. *CoRR*, abs/1510.0. 11
- Greene, D. and Cross, J. P. (2016). Exploring the political agenda of the european parliament using a dynamic topic modeling approach. *Political Analysis*, 25(1):77–94. 11
- Griffiths, T. L., Steyvers, M., and Tenenbaum, J. B. (2007). Topics in semantic representation. *Psychological Review*, 114(2):211–244. 3, 14
- Hall, D., Jurafsky, D., and Manning, C. D. (2008). *Studying the History of Ideas Using Topic Models*. 13
- Hao, S., Boyd-Graber, J. L., and Paul, M. J. (2018). Lessons from the Bible on Modern Topics: Adapting Topic Model Evaluation to Multilingual and Low-Resource Settings. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, pages 1090–1100.
- He, J., Li, L., and Wu, X. (2017). A self-adaptive sliding window based topic model for non-uniform texts. In *Proceedings - IEEE International Conference on Data Mining, ICDM*, volume 2017-Novem, pages 147–156. 11
- Hearst, M. a. and Hall, S. (1999). Untangling Text Data Mining. In *the 37th Annual Meeting of the Association for Computational Linguistics*, pages 1–13. 11
- Hofmann, T. (2001). Unsupervised Learning by Probabilistic Latent Semantic Analysis. *Machine Learning*, 42(1-2):177–196. 12
- Ji, J., Li, J., Yan, S., Tian, Q., and Zhang, B. (2013). Min-Max Hash for Jaccard Similarity. In *2013 IEEE 13th International Conference on Data Mining*, pages 301–309. IEEE. 17
- Kenter, T. and de Rijke, M. (2015). Short Text Similarity with Word Embeddings Categories and Subject Descriptors. *Proc. 24th ACM Int. Conf. Inf. Knowl. Manag. (CIKM 2015)*, pages 1411–1420. 3



- Krstovski, K. and Smith, D. A. (2011). A Minimally Supervised Approach for Detecting and Ranking Document Translation Pairs. In *Workshop on Statistical MT*. 16
- Krstovski, K., Smith, D. A., Wallach, H. M., and McGregor, A. (2013). Efficient Nearest-Neighbor Search in the Probability Simplex. In *Proceedings of the 2013 Conference on the Theory of Information Retrieval - ICTIR '13*, pages 101–108, New York, New York, USA. ACM Press. 17
- Kulis, B. and Grauman, K. (2012). Kernelized Locality-Sensitive Hashing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(6):1092–1104. 17
- Li, P. and König, C. (2010). b-Bit minwise hashing. In *Proceedings of the 19th international conference on World wide web - WWW '10*, page 671. ACM Press. 17
- Li, P., Owen, A. B., and Zhang, C.-H. (2012). One Permutation Hashing. *Advances in Neural Information Processing*. 17
- Lin, J. (1991). Divergence Measures Based on the Shannon Entropy. *IEEE Transactions on Information Theory*, 37(1):145–151. 13
- Lu, H. M., Wei, C. P., and Hsiao, F. Y. (2016). Modeling healthcare data using multiple-channel latent Dirichlet allocation. *Journal of Biomedical Informatics*, 60:210–223. 11
- Manning, C., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S., and McClosky, D. (2014). The Stanford CoreNLP Natural Language Processing Toolkit. In *Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. 10
- Mao, X., Feng, B.-S., Hao, Y.-J., Nie, L., Huang, H., and Wen, G. (2017). S2JSD-LSH: A Locality-Sensitive Hashing Schema for Probability Distributions. In *AAAI*. 17
- Mimno, D., Wallach, H. M., Naradowsky, J., Smith, D. A., and McCallum, A. (2009). Polylingual Topic Models. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing: Volume 2 - Volume 2*, EMNLP '09, pages 880–889, Stroudsburg, PA, USA. Association for Computational Linguistics. 18

- Moritz, M. and Hchler, M. B. (2017). Ambiguity in Semantically Related Word Substitutions: an Investigation in Historical Bible Translations. In *Proceedings of the NoDaLiDa 2017 Workshop on Processing Historical Language*, pages 18–23.
- Ni, X., Sun, J.-T., Hu, J., and Chen, Z. (2011). Cross Lingual Text Classification by Mining Multilingual Topics from Wikipedia. In *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining*, pages 375–384. 18
- O’Neill, J., Robin, C., O’Brien, L., and Buitelaar, P. (2017). An analysis of topic modelling for legislative texts. *CEUR Workshop Proceedings*, 2143. 11
- Platt, J. C., Toutanova, K., and Yih, W.-t. (2010). Translingual Document Representations from Discriminative Projections. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, EMNLP ’10*, pages 251–261, Stroudsburg, PA, USA. Association for Computational Linguistics. 19
- Rao, C. R. (1982). Diversity: Its Measurement, Decomposition, Apportionment and Analysis. *Sankhyā: The Indian Journal of Statistics, Series A*, 44(1):1–22. 13
- Rob Johnson, A. W. and Mabe, M. (2018). The stm report: An overview of scientific and scholarly journal publishing fifth edition. 1
- Rus, V., Niraula, N., and Banjade, R. (2013). Similarity Measures Based on Latent Dirichlet Allocation. In *Computational Linguistics and Intelligent Text Processing*, pages 459–470. 13
- Sciences, E. R. S. f. P. and life (2016). Harnessing the power of content - Extracting value from scientific literature: the power of mining full-text articles for pathway analysis Harnessing the Power of content. 9
- Tapi Nzali, M. D., Bringay, S., Lavergne, C., Mollevi, C., and Opitz, T. (2017). What Patients Can Tell Us: Topic Analysis for Social Media on Breast Cancer. *JMIR medical informatics*, 5(3):e23. 11
- Terasawa, K. and Tanaka, Y. (2007). Spherical LSH for Approximate Nearest Neighbor Search on Unit Hypersphere. In *Algorithms and Data Structures*, pages 27–38. 17

- Turchi, S., Ciofi, L., Paganelli, F., Pirri, F., and Giuli, D. (2012). Designing EPCIS through Linked Data and REST principles. *Software, Telecommunications and Computer Networks ({SoftCOM})*, 2012 20th International Conference on, pages 1–6. 36
- Vijayanarasimhan, S., Jain, P., and Grauman, K. (2014). Hashing Hyperplane Queries to Near Points with Applications to Large-Scale Active Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(2):276–288. 17
- Vulić, I., De Smet, W., Tang, J., and Moens, M. F. (2015). Probabilistic topic modeling in multilingual settings: An overview of its methodology and applications. *Information Processing and Management*, 51(1):111–147. 18
- Vulić, I. and Moens, M.-F. (2012). Detecting Highly Confident Word Translations from Comparable Corpora Without Any Prior Knowledge. In *Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics*, pages 449–459. 19
- Vulić, I. and Moens, M.-F. (2013). A Unified Framework for Monolingual and Cross-Lingual Relevance Modeling Based on Probabilistic Topic Models. In *Advances in Information Retrieval*, pages 98–109. 19
- Wang, J., Liu, W., Kumar, S., and Chang, S.-F. (2016). Learning to Hash for Indexing Big Data-A Survey. *Proceedings of the IEEE*, 104(1):34–57. 16
- Westergaard, D., Stærfeldt, H.-h., Tønsberg, C., Jensen, L. J., and Brunak, S. (2017). Text mining of 15 million full-text scientific articles. *bioRxiv*. 9
- Zhao, W.-L., Jégou, H., and Gravier, G. (2013). Sim-min-hash. In *Proceedings of the 21st ACM international conference on Multimedia - MM '13*, pages 577–580. 17
- Zhen, Y., Gao, Y., Yeung, D.-Y., Zha, H., and Li, X. (2016). Spectral Multimodal Hashing and Its Application to Multimedia Retrieval. *IEEE Transactions on Cybernetics*, 46(1):27–38. 16
- Zhu, Z., Li, M., Chen, L., and Yang, Z. (2013). Building Comparable Corpora Based on Bilingual {LDA} Model. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 278–282. 19