



**University of  
Zurich**<sup>UZH</sup>

**Zurich Open Repository and  
Archive**

University of Zurich  
Main Library  
Strickhofstrasse 39  
CH-8057 Zurich  
[www.zora.uzh.ch](http://www.zora.uzh.ch)

---

Year: 2016

---

## **Building a Corpus of Multi-lingual and Multi-format International Investment Agreements**

Sugisaki, Kyoko; Volk, Martin; Polanco, Rodrigo; Alschner, Wolfgang; Skougarevskiy, Dmitriy

**Abstract:** In this paper, we present an on-going research project whose aim is to develop a new database of international investment agreements that complements existing endeavors. In particular, this paper describes our efforts to build a standardized corpus of multi-lingual and multi-format agreement texts in order to enable researchers in the fields of international law and economics systematically investigate investment treaties.

Posted at the Zurich Open Repository and Archive, University of Zurich  
ZORA URL: <https://doi.org/10.5167/uzh-127362>

Originally published at:

Sugisaki, Kyoko; Volk, Martin; Polanco, Rodrigo; Alschner, Wolfgang; Skougarevskiy, Dmitriy (2016). Building a Corpus of Multi-lingual and Multi-format International Investment Agreements. In: 29th International Conference on Legal Knowledge and Information Systems (Jurix), Nice, December 2016 - December 2016.

# Building a Corpus of Multi-lingual and Multi-format International Investment Agreements

Kyoko SUGISAKI <sup>a,1</sup>, Martin VOLK <sup>a</sup>, Rodrigo POLANCO <sup>b</sup>,  
Wolfgang ALSCHNER <sup>b,c</sup> and Dmitriy SKOUGAREVSKIY <sup>c,d</sup>

<sup>a</sup>*Institute of Computational Linguistics, University of Zurich, Switzerland*

<sup>b</sup>*World Trade Institute, University of Bern, Switzerland*

<sup>c</sup>*Graduate Institute of International and Development Studies, Geneva, Switzerland*

<sup>d</sup>*Institute for the Rule of Law, European University, St. Petersburg, Russia*

**Abstract.** In this paper, we present an on-going research project whose aim is to develop a new database of international investment agreements that complements existing endeavors. In particular, this paper describes our efforts to build a standardized corpus of multi-lingual and multi-format agreement texts in order to enable researchers in the fields of international law and economics systematically investigate investment treaties.

**Keywords.** Natural language processing, Annotation of International investment agreements, OCR, Language identification

## 1. Introduction

International investment agreements (henceforth: IIAs) are treaties between two or more countries that are mainly designed for the “protection and liberalization of foreign investment.”[1] Since 1990, the number of IIAs has been growing [1] and worldwide, more than 3300 treaties are in force. These treaties are “a key instrument in the strategies of most countries, in particular developing countries, to attract foreign investment.”[2] However, they come in a broad variety of formats and languages. The majority of treaties are written in English. Yet, many IIAs are still exclusively in local languages of the contracting countries. Existing databases contain either only treaties in English,<sup>2</sup> or all in original language, and are incomplete.<sup>3</sup> Retrieving such treaties is often difficult for practitioners and scholars.

---

<sup>1</sup>All authors have been supported by Swiss Network for International Studies (SNIS) grant 37-740.

<sup>2</sup>The examples are Kluwer Arbitration “BITs” [<http://www.kluwerarbitration.com/CommonUI/BITs-countries.aspx>] and Oxford’s “Investment Claims” [<http://oxia.oup.com/>]

<sup>3</sup>The database of UNCTAD [<http://investmentpolicyhub.unctad.org/IIA>] is the most advanced and includes 3523 treaties. However, the database misses 772 treaties (22% of the total)

In our project,<sup>4</sup> we collect IIAs, and make them centrally available in one single language (English) and easily accessible by creating a machine-readable structured database. We have so far processed 1026 international investment agreements in HTML, Microsoft Word, and PDF format in more than 30 languages. A few thousand more are in the pipeline.<sup>5</sup> To bring these heterogeneous agreement texts into one format and one language, we have processed PDF documents through optical character recognition (OCR), translated multi-lingual texts into English, and converted unstructured text data into structured data.<sup>6</sup>

This paper describes the project phases from digitization to XML standardization. We focus on the language technology challenges in using OCR and language identification (Section 2), and introduce domain-specific XML mark-up for international investment treaties and automatic text structure recognition for the XML conversion (Section 3).

## 2. OCR and Automatic Language Identification

Out of the available source documents, 908 were PDFs. About half of these were not digital born but image scans. Therefore, we processed our PDF documents through a state-of-the-art OCR software.<sup>7</sup>

The conversion of treaties from more than 20 languages posed challenges because the software does not automatically identify the language. We thus first processed all documents without specific language settings. This led to unsatisfactory outputs, because the OCR software uses language-specific dictionaries to guess words.<sup>8</sup> Entering the languages of such a large number of documents manually was not a viable option. We therefore used a  $n$ -gram-based language identification system [3].<sup>9</sup>

Our PDF data collection, in particular, contains multi-lingual documents that include more than one language at the level of sentences, columns, and pages in a document (cf. Figure 1). To handle these three types of code-switching,<sup>10</sup> we segmented the texts into 15 words and changed the letters to lower case. Then we used these units as input to the language identification system. In this way, multi-lingual texts were processed as a unit, short enough to handle sentence-level code-switching but also long enough to handle  $n$ -gram in language identification. To obtain a high level of accuracy of language identification, we used the confidence score (i.e. probability estimate for the predicted language)

---

<sup>4</sup>The Swiss Network for International Studies project *Diffusion of International Law: A Textual Analysis of International Investment Agreements*: [http://www.snis.ch/project\\_diffusion-international-law-textual-analysis-international-investment-agreements](http://www.snis.ch/project_diffusion-international-law-textual-analysis-international-investment-agreements)

<sup>5</sup>Work on collecting IIAs is still in progress in the project, therefore we can not provide a total number of treaties. Yet we have so far collected 3329 treaties.

<sup>6</sup>By unstructured data, we refer to the fact that the texts are not marked with structure in the source data. The treaties themselves are very structured but without markups not accessible for machines.

<sup>7</sup>We used Abbyy Recognition Server: <https://www.abbyy.com/recognition-server/>

<sup>8</sup>The problem was also caused, because the OCR system is not restricted to the character set of a specific language.

<sup>9</sup>We used the off-the-shelf `langid.py` system [<https://pypi.python.org/pypi/langid/1.1.5>]. This tool is suitable to our source legal texts, as the  $n$ -gram is based on government-related documents among other documents (e.g. Wikipedia and Reuters).

<sup>10</sup>While code-switching usually refers to switching language within a sentence, we use the term in a broader sense.

returned by the language identification system and set a threshold of 0.75. We observed that the language of single-language documents was identified with higher confidence if more than 70% of the units in a document were recognized as the same language. In multi-lingual documents, on the other hand, we found that 25% of the units needed to be identified as one language to provide a good result.<sup>11</sup> Otherwise, we manually examined and corrected the languages of the documents.

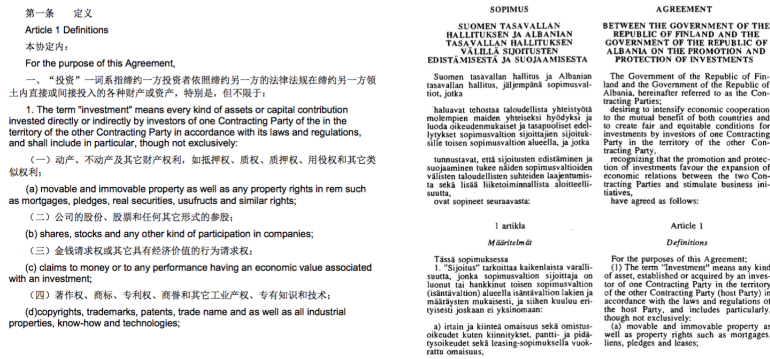


Figure 1. Examples of sentence-level (left) and column-level (right) multi-lingual documents

### 3. XML and Automatic Text Structure Recognition

Once we had all the sources in machine-readable format, we augmented them with mark-ups for (1) layout, (2) text structures, and (3) linguistic information.

Even though our source data comes in a variety of formats, the layout of various texts possesses a strong commonality; they contain text blocks and paragraphs. Therefore, we first converted all our multi-format documents into an XML mark-up for their layout structures. This XML layout is language-independent.

Secondly, we enhanced these layout mark-ups to also reflect document structure. International investment agreements comprise preface (such as title page and table of contents), preamble, text body (i.e. article and paragraph), conclusion (i.e. signatures), and some times attachments (e.g. annex). Hence, we divided each document into these five text zones and annotated the zones with XML mark-up.

#### Listing 1: Example of XML mark-up for document structures

```
<?xml version="1.0" encoding="utf-8"?>
<treaty>
  <main language="en">
    <preface/>
    <preamble><p>The Government of the Republic of Turkey and ...</p></preamble>
    <body>
      <div type="article" num="1" title="Definitions"><intro><p>For the purpose ...</p></intro></div>
    </body>
    <conclusion><p>IN WITNESS WHEREOF, ...</p></conclusion>
```

<sup>11</sup>Most of our documents contain less than 3 languages

```
<attachments/>
</main>
</treaty>
```

For the segmentation of these five text zones, we compiled a set of corresponding, typical linguistic feature patterns. For instance, a preamble typically begins with “The Government of the [county name]...” and ends with “... have agreed as follows;” or “... to conclude the following agreement;” Signatures begin with “In witness thereof ...” or “Done in duplicate at ...” Our linguistic feature method is similar to those for legal definition extraction [4], norm classification [5], and content zone identification of German court decisions [6].

In the text body and attachments, we automatically structured texts at the level of articles and paragraphs. This segmentation was again based on the surface patterns of enumeration structures, such as [alpha-numeric character(s)] + period, to recognize the beginning of items of enumerations such as *I.* or *Ia.*

Thirdly, we further enriched this XML mark-up with linguistic information for English. For this purpose, we used a tool of natural language processing [7]. Texts were automatically segmented into sentences, then the sentences were tokenized and augmented with part-of-speech tags, named entities, and dependency grammar structures. This linguistic information allows us to extract further information, for example, meta information such as the date of the signatures or legal definitions (cf. [4]).

#### 4. Conclusion

In this paper, we described our approach to constructing a corpus of international investment agreements. We illustrated that language technology can contribute significantly to reducing the manual effort when building a corpus based on multi-lingual and multi-format source data. In our future work, we will develop an information retrieval system that enables investment law scholars, arbitrators, negotiators, and other legal practitioners to easily retrieve investment treaties.

#### References

- [1] Kenneth J. Vandeveld. A brief history of international investment agreements. *U.C.-Davis Journal of International Law & Policy*, 12(1), 2005.
- [2] The role of international investment agreements in attracting foreign direct investment to developing countries. UNCTAD series on international investment policies for development, United Nations, New York and Geneva, 2009.
- [3] Kenneth Heafield, Rohan Kshirsagar, and Santiago Barona. Language identification and modeling in specialized hardware. In *Proceedings of the ACL-IJCNLP 2015*, 2015.
- [4] Stephan Walter. Linguistic description and automatic extraction of definitions from German court decisions. In *Proceedings of the LREC*, 2008.
- [5] Emile de Maat and Radboud Winkels. A next step towards automated modelling of sources of law. In *Proceedings of the ICAIL-09*, 2009.
- [6] Florian Kuhn. A description language for content zones of German court decisions. In *The proceedings of Workshop SPLeT-2010*, 2008.
- [7] Christopher D. Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven J. Bethard, and David McClosky. The Stanford CoreNLP natural language processing toolkit. In *Proceedings of the ACL: System Demonstrations*, 2014.