



**HAL**  
open science

## Searching for Discriminative Metadata of Heterogenous Corpora

Gaël Guibon, Isabelle Tellier, Sophie Prévost, Mathieu Constant, Kim Gerdes

► **To cite this version:**

Gaël Guibon, Isabelle Tellier, Sophie Prévost, Mathieu Constant, Kim Gerdes. Searching for Discriminative Metadata of Heterogenous Corpora. Fourteenth International Workshop on Treebanks and Linguistic Theories (TLT14), Dec 2015, Varsovie, Poland. pp.72-82. hal-01250981

**HAL Id: hal-01250981**

**<https://hal.science/hal-01250981>**

Submitted on 5 Jan 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Searching for Discriminative Metadata of Heterogenous Corpora

Gael Guibon (1), Isabelle Tellier (1,2)  
Sophie Prevost (1), Matthieu Constant (3) and Kim Gerdes (2,4)

(1) Lattice CNRS

(2) universit  Paris 3 - Sorbonne Nouvelle

(3) universit  Paris-Est, LIGM

(4) LPP CNRS

E-mails: gael.guibon@gmail.com, isabelle.tellier@univ-paris3.fr,  
Matthieu.Constant@u-pem.fr, sophie.prevost@ens.fr, kim@gerdes.fr

November 25, 2015

## Abstract

In this paper, we use machine learning techniques for part-of-speech tagging and parsing to explore the specificities of a highly heterogeneous corpus. The corpus used is a treebank of Old French made of texts which differ with respect to several types of metadata: production date, form (verse/prose), domain, and dialect. We conduct experiments in order to determine which of these metadata are the most discriminative and to induce a general methodology.

## 1 Introduction

Labeled data used to train POS taggers or syntactic parsers by supervised machine learning technics are usually rather homogenous: The texts they come from share a common origin and most of their metadata. Yet, most actual text collections that NLP tools have to handle today are heterogenous in many ways. The consequence of this inadequacy is well known: Programs trained on homogenous texts by supervised machine learning approaches do not perform well when applied to new texts that differ from them in any important aspects, such as subject, genre or sub-language. To address this problem, domain adaptation has become an important issue in machine learning NLP.

In this paper, we explore a very heterogenous corpus of Old French but the problem we tackle is not exactly domain adaptation. We want to use machine learning for *corpus exploration*, i.e. as a way to search for the most discriminative metadata of our texts. As a matter of fact, they belong to a *highly heterogenous*

*treebank* and vary in dialect, domain, production date, and form (verse and prose). In this context, new questions arise: Which of these properties (metadata, in the following) influence most the language in which they are written? How can we best train a POS tagger and a parser on this treebank, in order to annotate a new yet unlabeled text of Old French? Is it better to use, as training data, a small homogenous corpus similar to the new text or a large heterogeneous dissimilar one? These questions can also be relevant for other heterogeneous corpora and cross domain applications, so our purpose will also be to provide a general methodology.

## 2 Syntactic Reference Corpus of Medieval French

The SRCMF <sup>1</sup>[11] is a treebank of Old French texts enriched with POS tags (chosen among 60 distinct possible ones) and fine-grained dependency structures (labeled with 31 distinct syntactic functions) which were built manually during an ANR-DFG joint research project (2009-2012). The corpus consists of 15 texts (245 000 words) whose electronic versions are stemming from the "Base de Francais Médiéval" (BFM)<sup>2</sup> [4] and the "Nouveau Corpus d'Amsterdam" (NCA)<sup>3</sup> [5]. The selection of the included texts is based on criteria such as date, form (verse/prose), domain (historical, literary...) and dialect. From the SRCMF we choose 10 texts, whose metadata are shown in Table 1.

Text	Date	Words	Form	Dialect	Domain
<i>Vie Saint Légier</i>	late 10c.	1388	verse	n/a	religious
<i>Vie de Saint Alexis</i>	1050	4804	verse	normand	religious
<i>Chanson de Roland</i>	1100	28 766	verse	normand	literary
<i>Lapidaire en prose</i>	Mid. 12c.	4708	prose	anglo-norm.	didactical
<i>Yvain</i> , Chr. de Troyes	1177-1181	41 305	verse	champenois	literary
<i>La Conquête de Constantinople</i> , R. de Clari	>1205	33 534	prose	picard	historical
<i>Queste del Saint Graal</i>	1220	40 417	prose	n/a	literary
<i>Aucassin et Nicolette</i>	late 12c.-early 13c.	9844	verse & prose	picard	literary
<i>Miracles</i> from Gautier de Coinci	1218-1227	17 360	verse	picard	religious
<i>Roman de la Rose</i> from Jean de Meun	1269-1278	19 339	verse	n/a	didactical

Table 1: Texts of the Corpus and their Metadata

The syntactic annotation is based on a dependency model [12, 9], which distinguishes between, on the one hand, syntactic units and, on the other hand, different functions (such as "subject", "object", "adverbial", "auxiliary", "modifier"...),

<sup>1</sup><http://srcmf.org/>

<sup>2</sup><http://bfm.ens-lyon.fr/>

<sup>3</sup><http://www.uni-stuttgart.de/lingrom/stein/corpus/>

which specify the relation between a head and the structures depending on it.

### 3 Most Discriminative Metadata in SRCMF

In this section we propose a new experimental strategy to explore SRCMF based on metadata-related experiments, in order to find the metadata that are the most discriminative to tag and parse a new text.

#### 3.1 Protocol of our Experiments

Our experiments are not designed to search for the best parser and configurations, this has already been studied before [8, 3]. But, to evaluate the impact of each metadata, it is important to obtain comparable results. So, we define a general protocol for all the experiments whose results are reported in this paper. For each metadata value, we process as follows:

- The data are split into a training set (of equal size for each possible value) and a test set made of the remaining content;
- On all training and test sets, the lemmas are predicted by *TreeTagger* trained on the *Nouveau Corpus d'Amsterdam* (NCA) parameters<sup>4</sup> [10];
- A CRF part-of speech tagger [6] is trained on the training set using *Wapiti* 1.4.0 [7] with the same templates taking into account the contexts of words and lemmas, word endings, etc. as in [3];
- A Dependency parser is trained on the training set (with Gold POS labels) with *Mate-tools* (anna-3.61<sup>5</sup>) [1];
- The lemmatizer, the tagger and the parser are successively applied to the test set. The tagger is evaluated by its accuracy, the parser by the classical UAS and LAS values.

#### 3.2 Summary of the Experiments

**Dates.** We first conduct these experiments using the production dates of the texts, which are usually considered as discriminative metadata. This approach seems all the more obvious for the SRCMF whose texts date from the late 11th to the 13th century, which allows us to make a comparison between centuries. Nonetheless, as the 11th century sub-corpus is too small (less than 7 000 words), we restrict ourselves to the other two. For this time slicing, we also have to remove *Aucassin*

---

<sup>4</sup>[http://bfm.ens-lyon.fr/article.php3?id\\_article=324](http://bfm.ens-lyon.fr/article.php3?id_article=324)

<sup>5</sup>Available at <https://code.google.com/p/mate-tools/downloads/detail?name=anna-3.61.jar&can=2&q=>

from the data, as it is not clear whether it has been written in the 12th or the 13th century.

The results show that applying a model trained on the same century gives higher scores, especially in LAS. In fact, even if high UAS, LAS, and accuracy match with the correlation of century between train and test sets, the similarity between sub-corpora shows the contrary, with far more shared words when centuries differ, which is quite unexpected. This induces that the Out-Of-Vocabulary rate is not enough to predict the results. To confirm this intuition, we conducted few additional iterations of this experiment following the same method, in which texts from a century have been mixed, and a training set of 50 000 words<sup>6</sup> taken at random among them has been built. The performances reached by the models learned from these training sets lead to the same conclusions as previously. It is important to point out that even if both century based sub-corpora are composed of more than one domain and both verse and prose texts, the data from the 13th century are all written in Picard or in an undefined dialect. So, as one century does not vary in dialects, the date value is not independent from other metadata values. Moreover, as stated before, we could only use two centuries, which limits the scope of the results.

	Words	Units	Sentence	Sentence length
<b>12th century [train]</b>	50002	7337	5430	9
<b>13th century [train]</b>	50009	6934	3767	13
<b>12th century [test]</b>	24777	4844	2685	9
<b>13th century [test]</b>	60638	7597	4538	13

Table 2: Characteristics of the Time Sliced Sub-Corpora

**Forms.** The impact of the production date has been clearly shown but other experiments have to be conducted to compare it with the effect of the other metadata. We now consider text forms, i.e. if they are written in verse or prose. In the SR-CMF, most of the texts are in verse but some of them are in prose and one text (*Aucassin et Nicolette*) is written in both forms. It is well acknowledged that verses are syntactically more constrained than prose, and moreover they have a richer lexicon (based on a higher amount of different units): We thus wonder if this could have an impact, in particular on parsing. This effect combines with the fact that sentences in prose are usually longer than those in verse. The results in Table 5 contradict the idea that a different text form induces a significantly different model of dependency parsing. Indeed, the prose test corpus obtains higher UAS and LAS even from the verse model. Compared to using centuries, this time shared lexicon and known words show a greater similarity between corpora of the same text form.

According to our results, text form seems to be useful to discriminate texts from a corpus but only when it comes to part-of-speech tagging, as dependency parsing

<sup>6</sup>Because of the necessity to keep complete sentences, the size of randomly built training sets may vary by few words

Train \ Test		12th century [test]	13th century [test]
<b>12th century</b>	UAS	<b>88.81</b>	83.14
	LAS	<b>79.91</b>	71.93
	ACC	<b>94.69</b>	89.62
	Unknown   known words	<b>91.39</b>   08.61	78.72   <b>21.28</b>
	Different   shared lexicon	<b>61.20</b>   38.80	28.59   <b>71.41</b>
	Unknown   known words UAS	<b>81.05</b>   <b>90.00</b>	71.60   85.13
	Unknown   known words LAS	<b>66.42</b>   <b>81.47</b>	54.18   75.72
	Unknown   known words ACC	<b>87.29</b>   <b>95.39</b>	78.14   92.73
<b>13th century</b>	UAS	82.24	<b>89.07</b>
	LAS	69.24	<b>80.75</b>
	ACC	88.67	<b>94.62</b>
	Unkonwn   known words	73.83   <b>26.17</b>	<b>92.25</b>   07.75
	Different   shared lexicon	33.96   <b>66.04</b>	<b>50.12</b>   49.88
	Unkonwn   known words UAS	<b>76.94</b>   86.61	74.35   <b>88.77</b>
	Unkonwn   known words LAS	56.75   75.84	<b>57.96</b>   <b>80.46</b>
	Unkonwn   known words ACC	80.13   91.69	<b>85.31</b>   <b>95.41</b>

Table 3: Results Using Time Slicing

results seem irrelevant to our purpose. This makes the form metadata useful, but less reliable than the time slicing, when it comes to corpus exploration.

	Words	Units	Sentence	Sentence length
<b>prose [train]</b>	41910	4320	4320	14
<b>verse [train]</b>	41907	6840	6840	7
<b>prose [test]</b>	36749	4370	4370	12
<b>verse [test]</b>	34478	4417	4417	10

Table 4: Characteristics of the Text Form Based Sub-Corpora

**Domains.** Moving on to domain-related corpora, we now have more than two corpora to compare.

The notion of "domain" usually corresponds to what the texts are about. In the context of SRCMF, it is more related to the *literary genre* of these texts. Both notions do not exactly coincide but, in both cases, texts from the same domain should share some specific content words. As content words are less frequent and more ambiguous than grammatical words, we expect their presence or absence in both the training and test sets to affect the parsing results.

Domain adaptation is a prolific research field in machine learning and previous works have shown that it is possible to obtain better results by focusing on a specific domain rather than using a global approach [2]. This is why we could expect the domain value to have a great impact on the results.

A training set of about 16 000 words is first extracted from each of our four domain-specific sub-corpora, to ensure balanced training data. Test sets are made of the remaining content (Table 6). As expected, when the training and the test sets come from the same domain, the results, given in Table 7, are (in average) better

Train \ Test		Prose [test]	Verse [test]
<b>Prose</b>	UAS	<b>85.47</b>	76.33
	LAS	<b>74.96</b>	62.96
	ACC	<b>91.36</b>	83.61
	Unknown   known words	16.49   <b>83.51</b>	<b>21.26</b>   78.74
	Different   shared lexicon	57.02   <b>42.98</b>	<b>77.05</b>   22.95
	Unknown   known words UAS	<b>73.76</b>   <b>87.78</b>	65.87   79.15
	Unknown   known words LAS	<b>55.48</b>   <b>78.81</b>	46.37   67.44
	Unknown   known words ACC	<b>77.33</b>   <b>94.14</b>	76.78   85.46
<b>Verse</b>	UAS	<b>83.12</b>	82.79
	LAS	<b>71.52</b>	71.40
	ACC	90.06	<b>90.78</b>
	Unknown   known words	<b>18.81</b>   81.19	14.03   <b>85.97</b>
	Different   shared lexicon	<b>66.47</b>   33.53	42.52   <b>57.48</b>
	Unknown   known words UAS	<b>73.43</b>   <b>85.37</b>	72.39   84.49
	Unknown   known words LAS	55.45   <b>75.24</b>	<b>55.62</b>   73.98
	Unknown   known words ACC	81.02   <b>92.15</b>	<b>84.13</b>   91.86

Table 5: Verse [train] Shows Better Results on Prose [test] than on Verse [test], while Results on Known Words and Shared Lexicon Suggest the Contrary

and higher scores go along with the proportion of shared lexicon. This suggests that the domain is indeed a discriminative metadata. Moreover, results are quite stable except for the historical trainset, which shows an even greater gap when both sets come from the same domain, with an increase of about 30% in LAS. But, it is proper to remind that the historical corpora is made of only one text (*La Conquete de Constantinople*).

	Words	Units	Sentence	Sentence length
<b>Didactical [train]</b>	16003	3820	1238	12
<b>Historical [train]</b>	16007	2298	1108	14
<b>Literary [train]</b>	16009	3529	1526	10
<b>Religious [train]</b>	16011	3645	1470	10
<b>Didactical [test]</b>	8013	2374	680	11
<b>Historical [test]</b>	17528	2414	1249	14
<b>Literary [test]</b>	104323	10828	10209	10
<b>Religious [test]</b>	7541	2182	708	10

Table 6: Domain-Specific Corpora’s Charact.

**Dialects.** We then finally evaluate dialects as a discriminative metadata. A dialect speaker understands, at least partly, another dialect of the same language, as dialects share large parts of lexicon and grammar. In Table 9, we aim to determine whether or not the same holds for the dialects of our corpus. We use three sub-corpora based on the three distinct dialects, each training set being approximately made up of 20 000 words.

We observe a huge increase in performance (of about 10 points in LAS, UAS, and accuracy) while applying a model on a same dialect. When the training and

Train \ Test	Didactical[test]	Historical[test]	Literary[test]	Religious[test]
<b>Didactical [train]</b>				
UAS	<b>81.78</b>	78.88	80.11	70.05
LAS	<b>71.23</b>	67.28	66.67	55.04
ACC	<b>90.75</b>	87.58	87.08	80.80
Unknown   known words	16.53   <b>83.47</b>	<b>31.15</b>   68.85	26.08   73.92	30.58   69.42
Different   shared lexicon	50.19   <b>49.81</b>	78.05   21.95	<b>83.85</b>   16.15	69.67   30.33
Unknown   known w. UAS	<b>71.68</b>   <b>83.78</b>	69.12   83.29	70.13   83.63	59.80   74.57
Unknown   known w. LAS	<b>53.93</b>   <b>74.66</b>	52.89   73.79	50.69   72.29	38.46   62.34
Unknown   known w. ACC	80.89   <b>92.70</b>	<b>81.30</b>   90.43	77.53   90.44	66.96   86.89
<b>Historical [train]</b>				
UAS	67.49	<b>90.07</b>	73.03	32.29
LAS	51.12	<b>82.20</b>	57.30	45.08
ACC	72.74	<b>95.66</b>	76.67	69.93
Unknown   known w.	41.09   58.91	08.08   <b>91.92</b>	38.66   61.34	<b>42.57</b>   57.43
Different   shared lexicon	81.94   18.06	46.67   <b>53.33</b>	<b>90.46</b>   09.54	79.84   20.16
Unknown   known w. UAS	58.08   74.05	<b>80.16</b>   <b>90.94</b>	65.06   78.05	52.80   69.33
Unknown   known w. LAS	38.24   60.11	<b>63.70</b>   <b>83.92</b>	45.20   64.93	31.56   55.10
Unknown   known w. ACC	62.67   79.77	<b>87.50</b>   <b>96.38</b>	66.95   82.80	57.20   79.38
<b>Literary [train]</b>				
UAS	77.22	82.02	<b>84.79</b>	73.09
LAS	64.07	70.79	<b>73.63</b>	59.01
ACC	85.10	88.95	<b>91.93</b>	83.25
Unknown   known w.	27.01   72.99	<b>27.35</b>   72.65	14.42   <b>85.58</b>	27.16   72.84
Different   shared lexicon	68.17   31.83	73.58   26.42	<b>75.36</b>   24.64	65.96   <b>34.04</b>
Unknown   known w. UAS	66.17   81.31	74.07   85.02	<b>74.28</b>   <b>86.56</b>	61.18   77.53
Unknown   known w. LAS	46.03   70.74	<b>57.21</b>   75.90	56.25   <b>76.55</b>	40.67   65.84
Unknown   known w. ACC	73.61   89.35	80.72   92.04	<b>82.50</b>   <b>93.51</b>	69.04   88.55
<b>Religious [train]</b>				
UAS	74.99	79.76	79.52	<b>80.72</b>
LAS	61.61	67.94	65.94	<b>69.35</b>
ACC	83.31	87.62	85.91	<b>90.16</b>
Unknown   known w.	29.01   70.99	<b>29.50</b>   70.50	26.58   73.42	14.07   <b>85.93</b>
Different   shared lexicon	71.66   28.34	76.56   23.44	<b>85.05</b>   14.95	43.47   <b>56.53</b>
Unknown   known w. UAS	63.98   79.48	<b>70.17</b>   <b>83.77</b>	69.16   83.28	68.61   82.70
Unknown   known w. LAS	44.10   68.76	<b>53.31</b>   <b>74.06</b>	49.00   72.07	49.58   72.59
Unknown   known w. ACC	71.73   88.05	<b>81.80</b>   90.06	75.87   89.55	75.87   <b>92.50</b>

Table 7: Experiments Using Domain Based corpora

test sets do not stem from the same dialect, the shared lexicon is small. This could be due to the size of the SRCMF compared to contemporary language corpora, but more probably it is due to the heterogeneity of the texts in SRCMF, in particular concerning morpho-syntax and spelling, as shown in [3]. As a confirmation of the importance of shared lexicon: With an average sentence length of 7 words only, the *Normand* corpus should be easier to parse than the other dialects. Results do not clearly exhibit such differences (except for UAS), probably because the rate of known words is lower when evaluated on the *Normand* test set.

In any case, dialect turns out to be the most discriminative metadata among those evaluated, when it comes to predicting parsing results on Old French.



	Words	Units	Sentence	Sentence length
<b>Champenois [train]</b>	20005	3283	1821	10
<b>Normand [train]</b>	20028	3799	2639	7
<b>Picard [train]</b>	20011	3565	1538	13
<b>Champenois [test]</b>	21301	3440	1970	10
<b>Normand [test]</b>	13542	5503	1784	7
<b>Picard [test]</b>	40727	3042	3205	12

Table 8: Dialectal Corpora Characteristics

Train \ Test	Champenois[test]	Normand[test]	Picard[test]
<b>Champenois [train]</b>			
UAS	<b>86.07</b>	78.61	76.66
LAS	<b>76.30</b>	61.93	63.63
ACC	<b>93.41</b>	81.17	84.02
Unknown/known words	10.23   <b>89.77</b>	<b>51.05</b>   48.95	31.20   68.80
Different/shared lexicon	51.09   <b>48.91</b>	<b>82.09</b>   17.91	79.56   20.44
Unknown/known words UAS	<b>73.83</b>   <b>87.46</b>	72.83   84.63	66.38   81.32
Unknown/known words LAS	<b>59.14</b>   <b>78.25</b>	51.34   72.98	46.29   71.49
Unknown/known words ACC	<b>84.57</b>   <b>94.41</b>	72.59   90.12	67.99   91.30
<b>Normand [train]</b>			
UAS	74.54	<b>88</b>	73.77
LAS	59.31	<b>77.96</b>	60.48
ACC	81.12	<b>93.31</b>	82.55
Unknown/known words	34.14   65.86	11.25   <b>88.75</b>	<b>38.77</b>   61.23
Different/shared lexicon	82.24   17.76	43.90   <b>56.10</b>	<b>87.05</b>   12.95
Unknown/known words UAS	64.37   79.81	<b>78.53</b>   <b>89.20</b>	64.19   79.84
Unknown/known words LAS	45.30   66.58	<b>60.21</b>   <b>80.21</b>	46.86   69.11
Unknown/known words ACC	72.54   85.57	<b>82.01</b>   <b>94.74</b>	72.50   88.92
<b>Picard [train]</b>			
UAS	77.35	79.41	<b>85.14</b>
LAS	63.46	63.20	<b>75.90</b>
ACC	84.40	82.11	<b>93.25</b>
Unknown/known words	24.58   75.42	<b>45.57</b>   54.23	11.16   <b>88.84</b>
Different/shared lexicon	74.51   25.49	<b>82.42</b>   17.58	60.03   <b>39.97</b>
Unknown/known words UAS	66.15   81.00	<b>72.60</b>   85.24	71.49   <b>86.86</b>
Unknown/known words LAS	47.03   68.81	51.47   72.98	<b>55.29</b>   <b>78.49</b>
Unknown/known words ACC	75.34   87.34	72.93   89.78	<b>80.56</b>   <b>94.85</b>

Table 9: Experiments Using the Dialectal Corpus Segmentation

To go even further with this metadata, we make a complementary “leave one out” experiment using the dialect-based corpus segmentation, in order to determine which of its three distinct values is the most different from a machine learning point of view. The results in Table 10 show that the *normand* dialect seems to be the most remote one with the lowest amount of shared lexicon and known words. It leads to the lowest accuracy and LAS, however for both unknown and known words its UAS and LAS are higher than for the other two corpora. This can only be possible because it contains a higher rate of unknown words. This means that low proximity

between corpora does not necessary brings lower unknown words recognition.

The other 2 tested on	Champenois	Normand	Picard
UAS	84.18	<b>84.41</b>	82.06
LAS	<b>72.64</b>	70.41	70.98
ACC	<b>89.86</b>	86.56	88.54
unknown/known w. different/shared lex.	17.04   <b>82.96</b> 63.08   <b>36.92</b>	<b>37.45</b>   62.55 <b>75.97</b>   24.03	25.05   74.95 71.55   28.45
unknown/known w. uas	75.25   86.01	<b>79.16</b>   <b>87.56</b>	73.56   84.91
unknown/known w. las	58.82   75.49	<b>59.84</b>   <b>76.74</b>	56.39   75.85
unknown/known w. acc	<b>84.29</b>   91.01	79.66   90.69	76.97   <b>92.41</b>

Table 10: Leave One Out experiments Using the Dialectal Corpus Segmentation

## 4 Towards a General Methodology

This exploration based on metadata opens interesting perspectives. Given a new text with its associated metadata as an input, one can expect to develop a general methodology to find the best tagging or parsing model. Let the associated metadata be a set of attribute-value pairs. For instance, assume that the new input text has the following set of attribute-value pairs: *century=13th*, *domain=literary*, *form=verse*, *dialect=picard*. From the previous experimental results, one can find the best tagging/parsing model among those available. The selected parsing model for the input text would be the one leading to the best LAS among those associated with the text metadata. In our example, we have four possible candidate models, each one associated with an attribute-value pair:

- *century=13th* (LAS=80.75 trained on *13th century [train]*)
- *domain=literary* (LAS=73.63, trained on *literary [train]*)
- *form=verse* (LAS=71.40 trained on *verse [train]*)
- *dialect=picard* (LAS=75.90, trained on *picard [train]*)

In this example, the best parsing model for *century=13th* seems to be the one trained on *13th century [train]*, reaching a LAS of 80.75 (cf. Table 7).

This proposed methodology is still rudimentary. We are aware that, in our corpus, the metadata are correlated and interfere with each other. For instance, in SRCMF, texts written in *Picard* are all from the 13th century. This cannot be avoided due to the lack of available texts, let alone tagged corpora, of Old French. Furthermore, there is an unavoidable part of an arbitrary in the way metadata values are defined (the various distinguished domains, the time slicing, knowing that changes do not occurs specifically at the turn of two centuries...). The different sizes of the training sets are also another issue for such a general method.

## 5 Conclusion

In this paper, we have shown that machine learning could serve as a very effective corpus exploration strategy. Each experiment helps us to better understand the specificities of our highly heterogenous corpus. The originality of the approach we have followed is that it is focused on *metadata discrimination*. Machine learning engineering is usually more concerned with feature selection or parameter optimization, applied to stable training and test sets to gain better overall results. Here, the machine learning devices used are stable, but we vary the way training and test corpora are built, in order to evaluate their influence on the final result.

This work can be extended in various ways. With SRCMF, the metadata selection we have started out with should be investigated further, with the goal to build a complete decision tree for a given new text whose metadata are known. The ideal decision tree would provide the best possible labeled sub-corpus to use as a training set, to build the best possible model for this given text. To achieve this goal, the correlations between distinct metadata should also be investigated further. We are prudent concerning the generalizability of some of our conclusions, because the impacts of metadata are mixed with other factors: size of the available sub-corpora, lexical variation, effect of the combination of metadata... More experiments are necessary to clarify each of them.

Nevertheless, we believe that this global approach could be relevant in many contexts, as heterogenous corpora are increasingly becoming an important subject of parsing technologies.

## References

- [1] Bernd Bohnet. Very high accuracy and fast dependency parsing is not a contradiction. In *The 23rd International Conference on Computational Linguistics (COLING 2010)*, Beijing, China, 2010.
- [2] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. Domain adaptation for large-scale sentiment classification: A deep learning approach. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 513–520, 2011.
- [3] Gaël Guibon, Isabelle Tellier, Matthieu Constant, Sophie Prévost, and Kim Gerdes. Parsing poorly standardized language dependency on old french. In *13th Treebank and Language Theory (TLT)*, 2014.
- [4] Céline Guillot, Alexei Lavrentiev, and Christiane Marchello-Nizia. Document 2. les corpus de français médiéval : Etat des lieux et perspectives. *Revue française de linguistique appliquée*, XII:125–128, 2007.
- [5] Pierre Kunstmann and Achim Stein. Le nouveau corpus d’amsterdam. In *"Actes de l’atelier de Lauterbad"*, pages 9–27, 2007.

- [6] John Lafferty, Andrew McCallum, and Fernando C.N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning 2001 (ICML 2001)*, pages 282–289, Seattle, Washington, 2001.
- [7] Thomas Lavergne, Olivier Cappé, and François Yvon. Practical very large scale CRFs. In *Proceedings the 48th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 504–513, July 2010.
- [8] Francesco Mambrini and Marco Passarotti. Will a parser overtake achilles? first experiments on parsing the ancient greek dependency treebank. In *Eleventh International Workshop on Treebanks and Linguistic Theories*, pages 133–144. Edições Colibri, 2012.
- [9] Alain Polguère and al. *Dependency in linguistic description*, volume 111. John Benjamins Publishing, 2009.
- [10] Achim Stein. Parsing heterogeneous corpora with a rich dependency grammar. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Hrafn Loftsson, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland, may 2014. European Language Resources Association.
- [11] Achim Stein and Sophie Prévost. Syntactic annotation of medieval texts: the syntactic reference corpus of medieval french (srcmf). In Tübingen: Narr, editor, *New Methods in Historical Corpus Linguistics*. 2013. Corpus Linguistics and International Perspectives on Language, CLIP Vol. 3, P. Bennett, M. Durrell, S. Scheible and R. Whitt (eds).
- [12] Lucien Tesnière. *Eléments de syntaxe structurale*. Librairie C. Klincksieck, 1959.