



Introduction

We present a diachronic parallel corpus on the Credit Suisse Bulletin, the world's oldest banking magazine, and three of its applications: vocabulary changes, the translation of the partitive particle, and trends and changes in society over time.

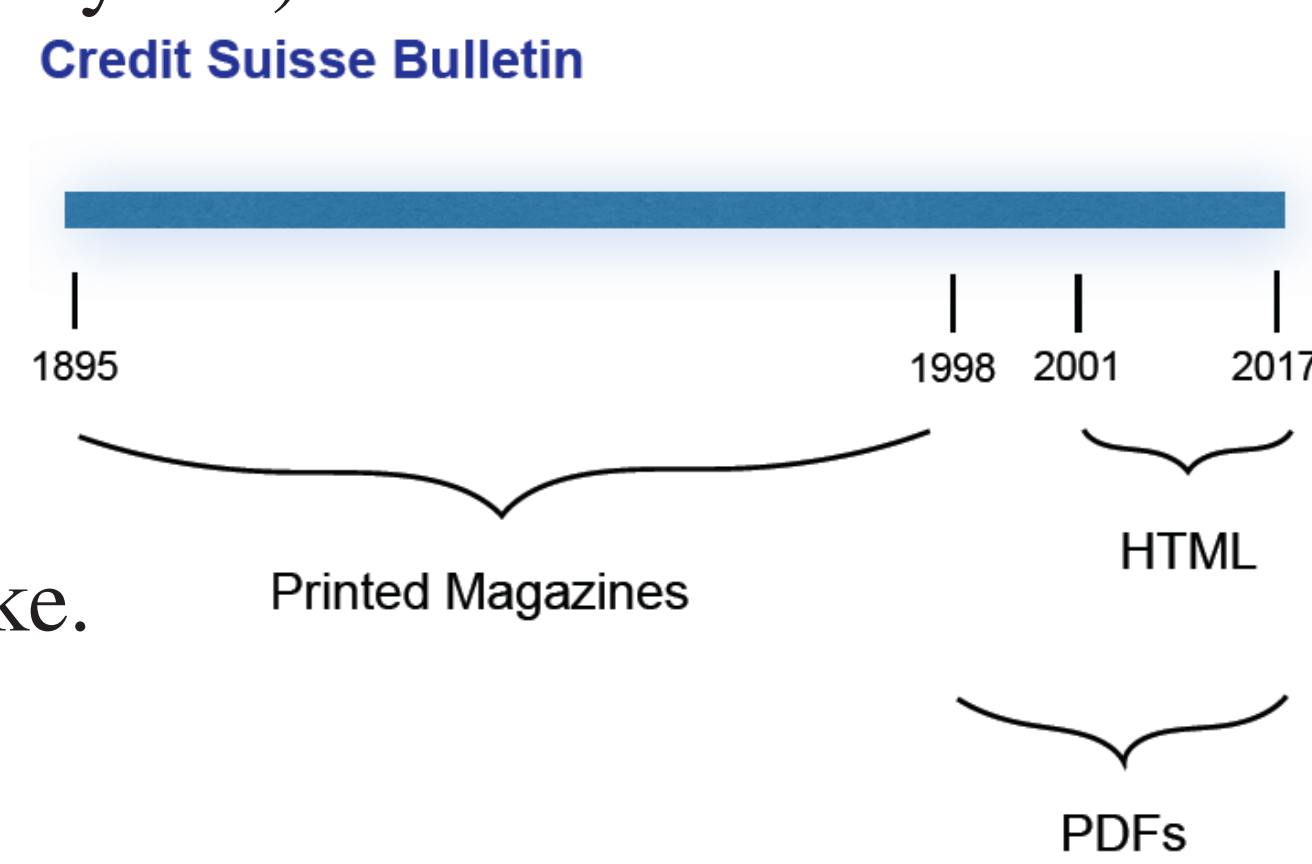
CS Bulletin Corpus

THE CREDIT SUISSE BULLETIN CORPUS has been published since 1895 in German, with translations in French and partly in English and Italian. The German and French part contain over 20 million words each, the English and Italian part about 10 million words each. Our parallel, multilingual corpus fills a gap in parallel corpora with respect to genre (magazine articles), domain (banking and economy articles), and its time span (120 years):

It can be freely accessed and downloaded.

Its characteristics allow researchers to conduct a multitude of orthogonal types of research, turning it into a treasure for historians, sociologists, and linguists alike.

We showcase three applications:



Application: Vocabulary

DOCUMENT CLASSIFICATION is applied to the corpus split into an early (until 1969) and late (from 1970 on) period. We use logistic regression on a monogram bag-of-words model. The classification accuracy is 98.5%. Many of the features show developments in banking vocabulary, language, orientation of the magazine. We show selected features ranked by decreasing weight.

Table with 2 columns: Late Period and Early Period. Each column lists Rank, word, Freq, and weight. Words include zeitfrage, foto, bulletin, anleger, kunde, suisse, partner, ordern, weltweit, perspektive, spreitenbach, warenhaus, mio, zinstrend, grafik, veränderung, welt, krieg, aktienindex, finanziell, reingewinn, indessen, aktiengesellschaft, bern, wechselstube, industrie, zinsfuss, herr, bundesbahn, million, valuta, westdeutschland, bankier, pfandbrief.

Selected features (see Rank) of early and late period.

Application: Translation

THE ITALIAN PARTITIVE PARTICLE ne is translated to French en in about half of the instances, and direct translations to German are rarer (dafür, davon, dazu, da, daraus), as one can see in our online browser multilingwis:

Help About "Credit Suisse Bulletin" Institute of Computational Linguistics / University of Zurich - SPARCLING project - Code Repository

More often it is (1) dropped or (2) partly translated or (3,4) a fixed expression

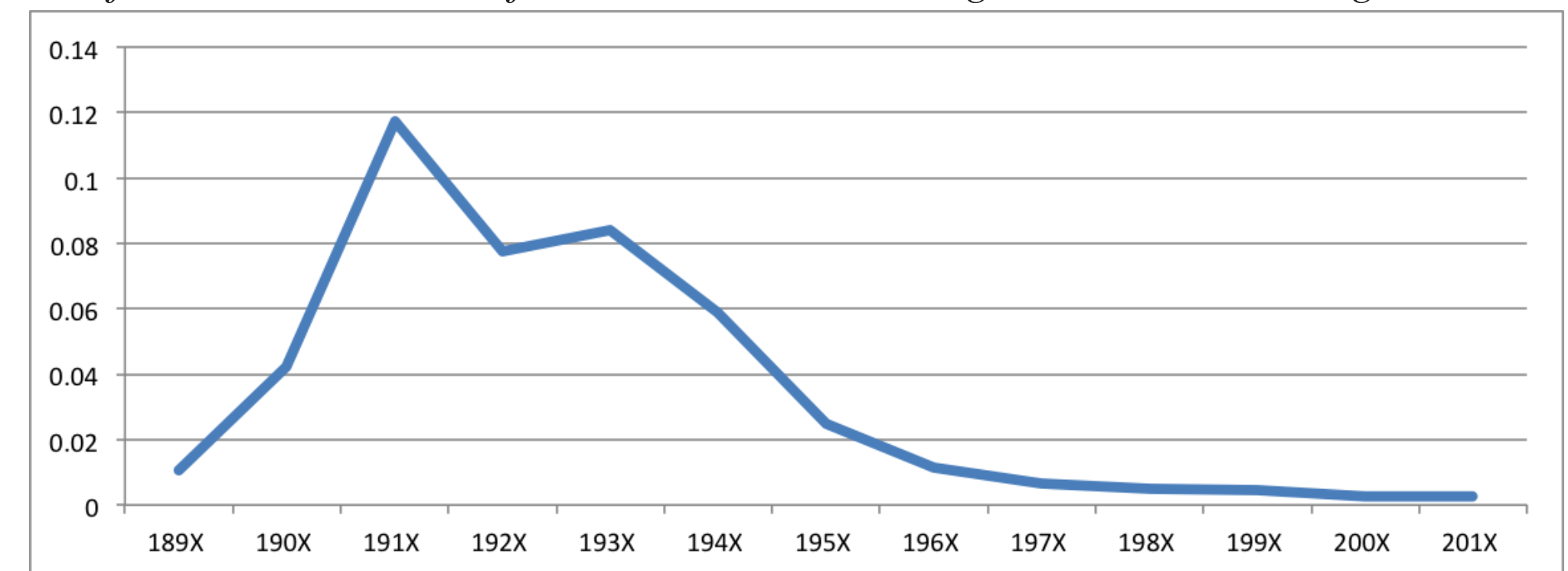
Table with 4 columns showing examples of translation variations for the particle 'ne' in German, French, and Italian.

Application: Changes in Society

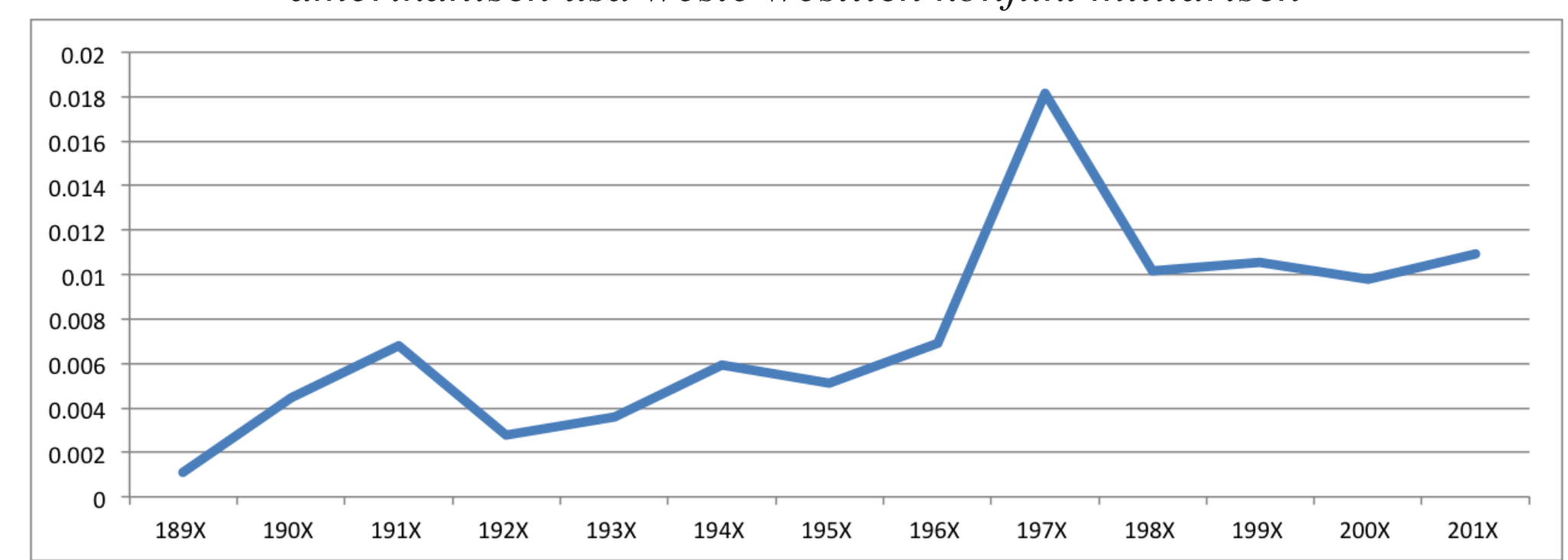
TOPICS IN SOCIETY are changing, and also what a modern banking magazine wants to be seen as reporting on.

We use topic modeling to detect hidden topics and investigate how they change over time. With 50 topics we observe e.g. the following changes in selected topics (we give keywords and proportion of decade):

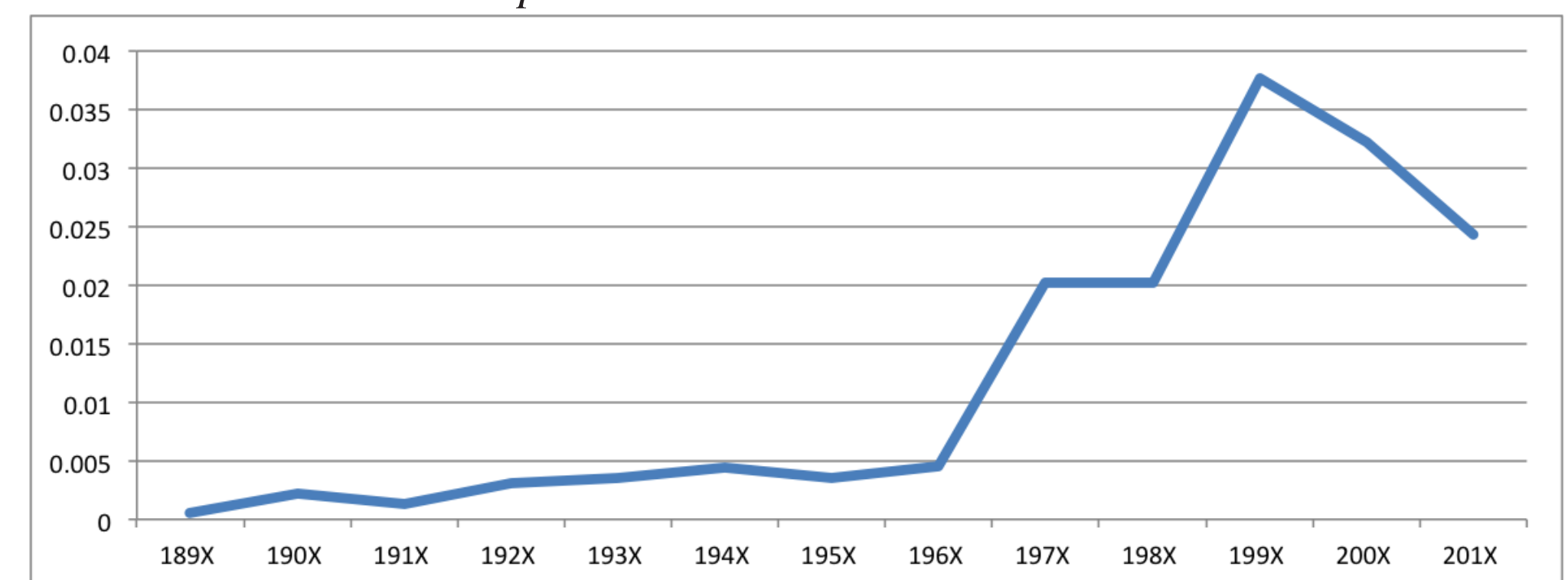
land staat milliarde million krieg frankreich england deutschland vereinigt ausland regierung amerikanisch französisch wirtschaftlich teil international gold deutsch zeit englisch



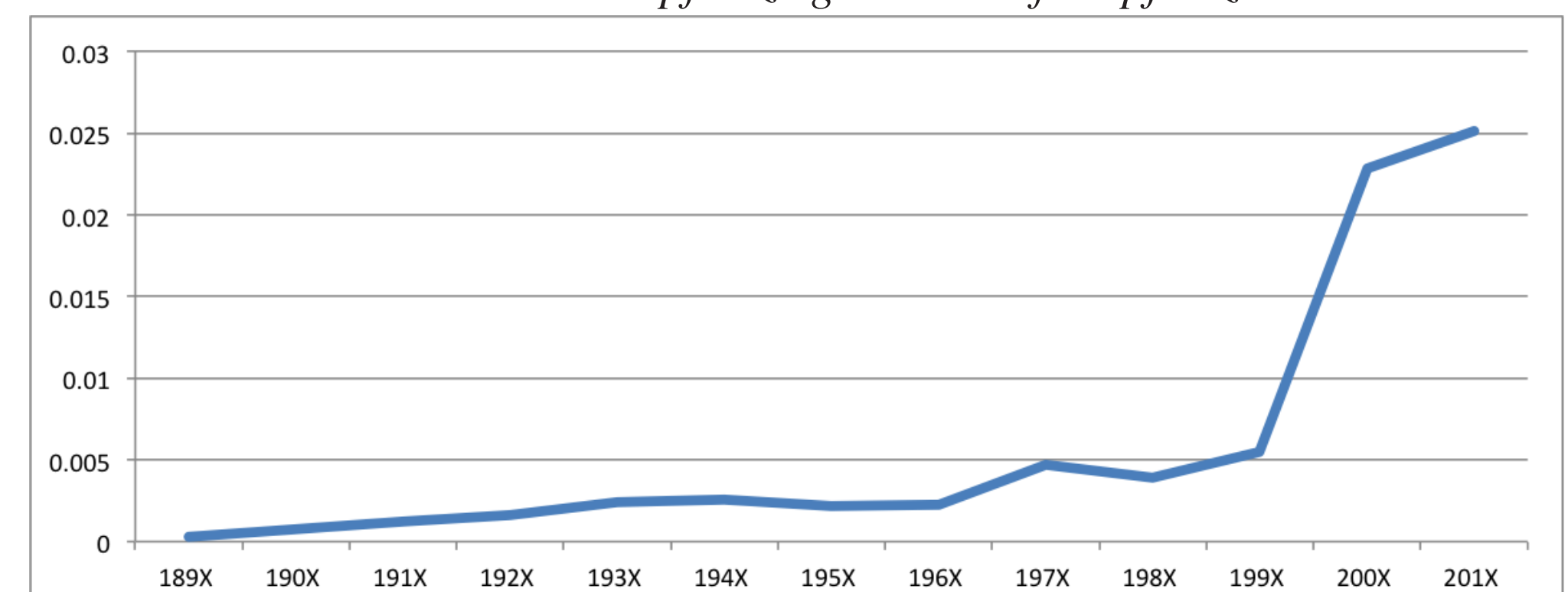
politisch prääsident politik partei regierung wahl staat land osten westen europa krieg frankreich sowjetunion amerikanisch usa weste westlich konflikt militärisch



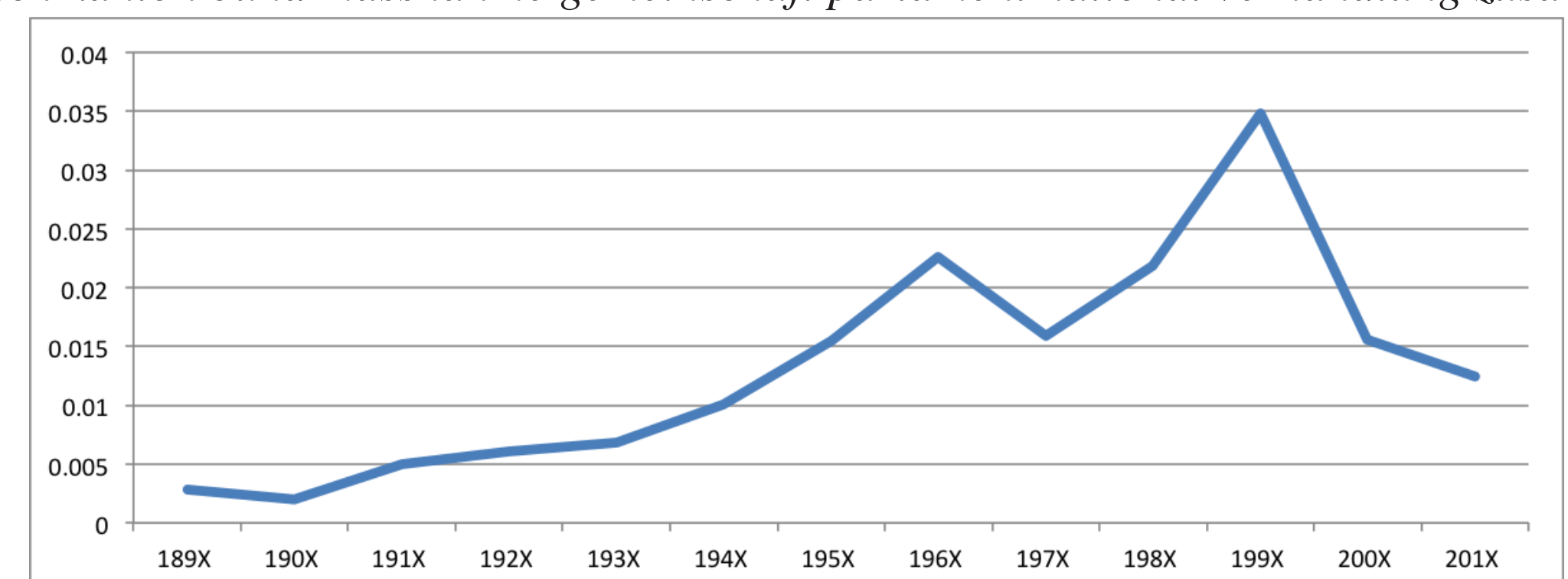
jahrhundert kunst werk bild künstler ausstellung geschichte alt museum sammlung zeit münze magazin zürich buch spät bulletin modern art stadt



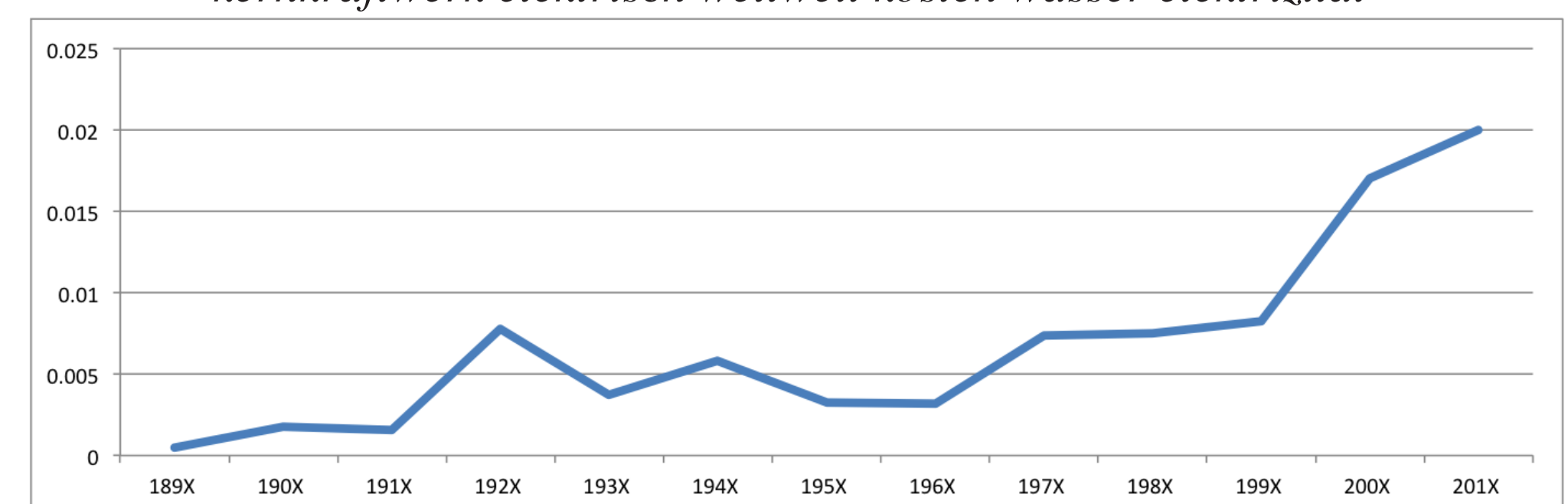
mensch tier wein arzt wasser patient essen krankheit körper menschlich forschung forschner gesund medizinisch natur medikament pflanze gesundheit foto pflanzen



europäisch schweiz land staat bundesrat international politisch schweizerisch frage gemeinsam wichtig eidgenössisch kanton bund massnahme gemeinschaft parlament national verhandlung zusammenarbeit



energie bau bauen projekt wald strom gebäude umwelt anlage erneuerbar nachhaltig stadt holz ökologisch kernkraftwerk elektrisch weltweit kosten wasser elektrizität



Conclusions

- The CS Bulletin corpus is a multi-purpose corpus. As a diachronic multilingual parallel corpus it allows statistical analyses of changes in language and society over the span of 120 years. Aligned sentences serve as data set for translation tools. We have shown how text classification and topic modeling tools can be used on our corpus to visualize societal trends.

References

Our corpus can be downloaded at http://pub.cl.uzh.ch/projects/b4c/en/corpora.php Volk/Amrhein/Aeppli/Müller/Ströbel. 2016. "Building a Parallel Corpus on the World's Oldest Banking Magazine". Proceedings of KONVENS 2016, Bochum.