# Compiling Tools and Resources for Studying of Luxemburgish Language and beyond.

*Joshgun Sirajzade*
University of Luxembourg
joshgun.sirajzade@uni.lu

After collaborating and working in different projects, which directly applied computational methods, at the Institute of Luxembourgish Language and Literature under Prof. Dr. Peter Gilles, University of Luxembourg[1] I would like not only to summarize the tools and resources we built in the last years but also to give a theoretical analysis of how computational methods can actually affect or even change our research and vice versa, enabling us to see new correlations.[2] For a long time the responsibility of the available technologies for the Digital Humanities remained in the information technology sector, and they were merely considered as routines or tools (Ramsay and Rockwell 2012). As the technology is crucial for our research today, I would like to highlight the necessity to question and develop the technologies we use in the Digital Humanities ourselves. However, looking back I see that some decades ago every project used – just to name one example – some different sort of data encoding and was doing a tedious work in order to query and exchange the data, whereas today we have interoperable strategies like XML and lots of ready to use tools available, which make our work faster and easier. XML has become pretty famous in the Digitals Humanities for some time now, and the same applies to databases, programming languages, web technologies etc., too. Digital Humanities needs its own tools and I think there should also be a theoretical part about it inside of Digital Humanities, while it is also a good thing to leave the low-level part of it like memory control etc. to computer science. What Digital Humanities is interested in is high-level programming for its own tasks. Thus, the question has two tails; a) what tools are needed in Digital Humanities today and b) how can the hermeneutics of Digital Humanities be used to improve these tools? I think if we combine these two strategies, we might one day receive better answers to our questions, which conventional humanities were not able to give. But before providing an insight into theoretical issues, let me introduce some empirical work, on which my theoretical thoughts are based.

The Idea behind the platform Phraseolux[3] was about gathering Luxembourgish phraseologies and storing them in a relational database. The gathered phraseologies were annotated by their lexical constitutes and semantical fields of their meaning. By storing them in a database it was not only possible to get precise frequencies about the material very quickly (and that in many possible ways, e.g. number of phraseologies as a token and type or number of word tokens and types etc.) and do some advanced statistics but also to achieve a better semantical understanding of the phraseologies

---

through their connection by the semantical field. By doing so, it was possible to see which semantical fields have more phraseological units, thus about what people make the most metaphorical thinking. Another interesting point was visualizing these data and seeing how recipients interpret them (See figure 1). Having the data organized digitally in a database opened so many different opportunities that not all possibilities were able to be drawn upon and the institute is still making use of the data again and again, even after the project was officially finished in 2014. Technically the project made use of scripting web languages like PHP, database software MySQL (later one of its successors MariaDB) and JavaScript libraries like d3.js for visualizing the data.
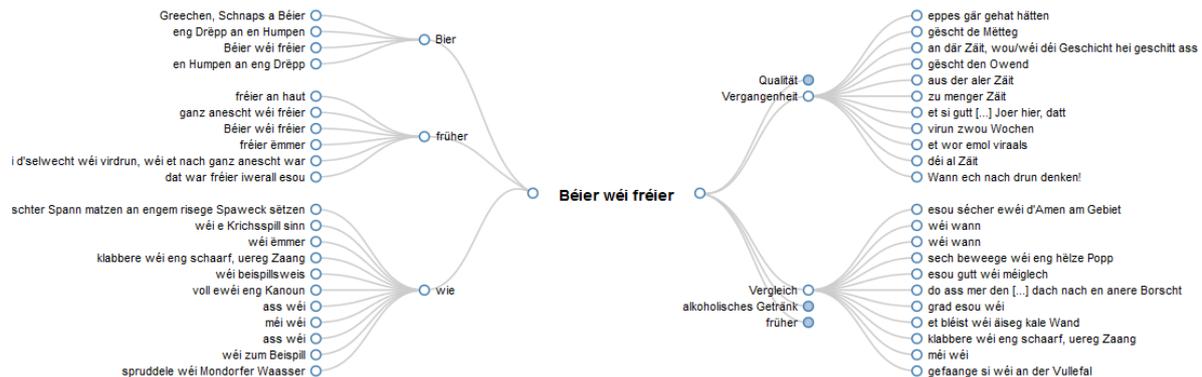


Figure 1 Lexical constituents and semantical fields of 'Béier wéi fréier'(beer as before). By clicking on them (Vergleich (comparing) or Vergangeheit (the past)) one can see other phraseologies, which have same annotations.

Currently I am working on the project WBLUX2 — Luxembourgish word formation. Historical preconditions and contrastive analysis.[4] The project aims to study word formation in the Luxembourgish language. This time we decided to use a different database, namely the so called XML-Database eXist, because our empirical text sources were already available in XML. As a programming language we took Java this time. Script languages are often simple to use and pretty fast to write. But the produced code is often not easily readable, and errors are difficult to find. In my opinion if you work in a bigger project, you need to use a language, in which code can be maintained easily.[5] Java is one of the most common programming languages, has lots of tools and so called libraries with high level abstraction – which I am going to discuss later – that makes it a good choice for long living software. For the study of word formation, we used some tools and strategies, which come both from digital humanities and corpus linguistics. First of all our texts were coded in TEI, so everyone, who is familiar with TEI can take look at our data and knows immediately what it is about. Secondly, we used strategies like Tokenization, Standardization and POS-Tagging, which derive from corpus linguistics (see figure 2 and 3). And again it is about linking the data, organizing it smartly, in order to get best results. For example, if you are searching for an ending -er like in writer, you also get better, because a language sign is multifunctional. But if you use information from POS-Tagging and search only in nouns for this ending, you will get better results. By doing so, it is possible to build annotations on one another, which can be used both in GUI for users as a search tools and for automatic tasks as well.

---

[4] Project Leader: Prof. Dr. Peter Gilles, Prof. Dr. Heinz Sieburg, Funding: University of Luxembourg.
[5] Sustainability is a big issue in building tools, which can also be 'thorny' for any project (Terras 2012).

In order to make theoretical assumptions the situation in computer science has to be reviewed first. In the history of software development it has always been the similar case to using existing information in order to produce new ones. There it is called abstraction[6] and has two perspectives: control abstraction and data abstraction. Although the term 'abstraction' can also mean something different in the humanities, here it refers to the combination of commands e.g. subroutines, modules or software components, in order to generate new commands for prompt and frequent usage. As example one can imagine, if I want to open and see a picture, a computer needs to perform many steps starting with reading the data from the hard disc to the memory, and then sending it to the monitor to turn on certain pixels in certain colors etc. Basically, if you need to do it frequently, you abstract these steps to one command, instead of describing these small steps every time again. In fact, you can break down every action on a computer until it goes to ones and zeros. However, as a digital humanist you are often not even really interested in these tiny steps, but rather in using existing abstract tools (like XSLT or certain Java and Python libraries) for certain tasks to do the job.

There are some well-known advantages of humanities going digital, like accessibility over internet, or the fact that everything can easily be copied and multiplied. However, one particular advantage has not yet been discussed extensively. This is the ability of abstraction, thus using software parts and information in order to build new ones. And this is not only control and data abstraction, but also concerns other fields like linking or visualizing the data. A good example here would be hyperlinks, which are almost predated from today's perspective. But in the history of the internet, these were developed as a (new) concept and technical realization. So where are we going now – popups, widget switches, drag and drops to name some other newer methods of acting on internet. Which technologies are best for certain parts of digital humanities? Knowing the history of how digital media evolved so far, we can take a look to the future, and not only predict things which are going to come, but actively take part in their development.

---

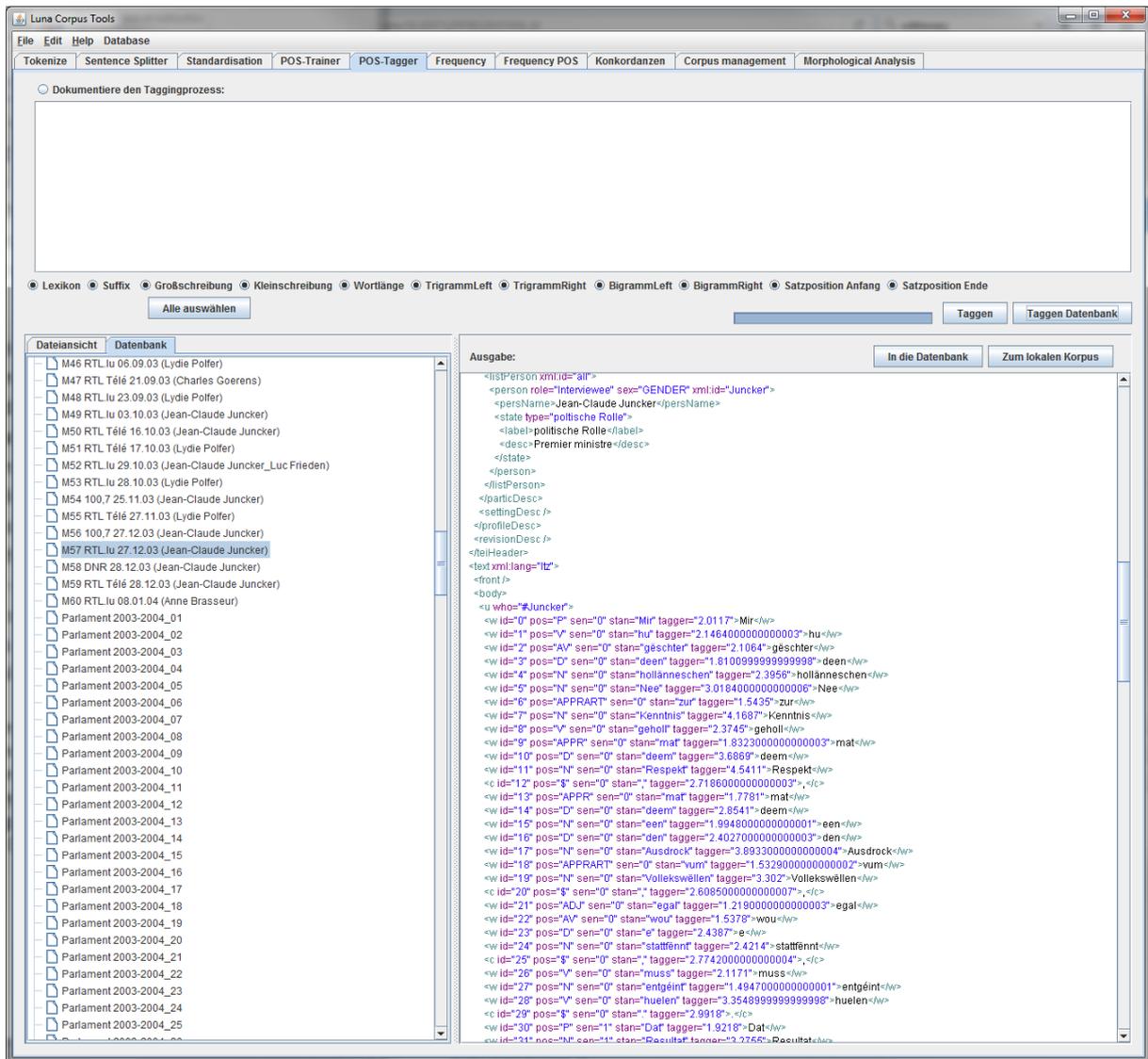[6] See about „Abstract. Simplify. Generalize" In: (Mcpherson 2012).

**Figure 2** The Java Swing GUI to the bundle of tools, like Tokeniser, POS-Tagger etc. The Tool has eXist-Database at the backend and supports XML
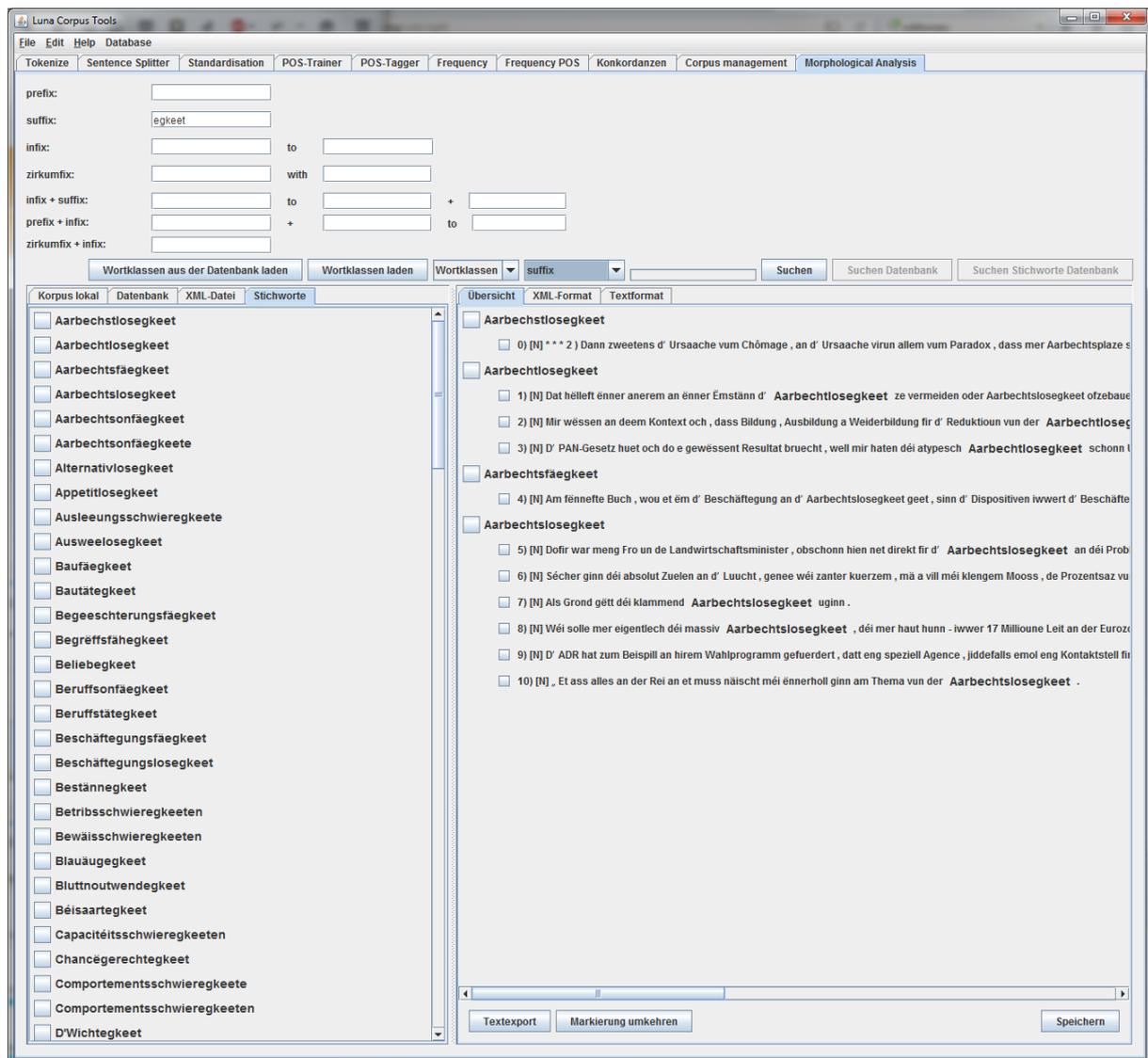
**Figure 3 The Part of the GUI where users can search for affixes with combination of part of speech**

## References

Fitzpatrick, K. (2012). The Humanities, Done Digitally. In: Matthew K. Gold (Ed.) Debates in the Digital Humanities. University of Minnesota Press. Minneapolis. p 12-13.

InfoLux. (2016) Fuerschungsportal iwwert d'Lëtzebuergescht. Institut fir lëtzebuergesch Sprooch- a Literaturwëssenschaft. Available: infolux.uni.lu. Last accessed 18th April 2016.

Mcpherson, T (2012). Why Are the Digital Humanities So White? or Thinking the Histories of Race and Computation. In: Matthew K. Gold (Ed.) Debates in the Digital Humanities. University of Minnesota Press. Minneapolis. p 146.

PhraseoLux. (2016). Datenbank zur Phraseologie des Luxemburgischen . Available: http://phraseolux.uni.lu. Last accessed 18th April 2016.

Ramsay, S. and Rockwell, G. (2012). Developing Things: Notes toward an Epistemology of Building in the Digital Humanities. In: Gold, Matthew K. (Ed.) Debates in the Digital Humanities. University of Minnesota Press. Minneapolis. p 75.

Terras M. (2012). Present, Not Voting: Digital Humanities in the Panopticon. In David M. Berry (ed.) Understanding Digital Humanities. Palgrave Macmillan, London. p 172.