# From linked open data to meta-analysis in historical linguistics (Presentation)

Joonas Kesäniemi, Turo Vartiainen, Tanja Säily, Agata Dominowska, Aatu Liimatta, Terttu Nevalainen

Empirical work on language change is fragmented, and some of the published works in the field are hard to come by. Furthermore, the data on which the research is based is seldom available to the research community. The **Language Change Database** (LCD) is an online resource which draws together earlier corpus-based research on English historical linguistics, with the goal of making it more accessible and cumulative by providing comparative baseline data from earlier studies (Nevalainen et al. 2016).

The information included in the LCD goes far beyond basic bibliographies in that it includes detailed descriptions of the results of the studies along with quantitative data in tabular format. In addition to historical linguistics in general, the data in the LCD is tailored to the needs of researchers interested in statistical modelling, systematic reviews, replication of earlier research and sociolinguistic typologies. The most recent additions to the data model, namely corpus composition data and annotated data tables, also make the LCD a valuable source of structured data for the purposes of meta-analysis.

The LCD data is published with permissive licensing and distributed via open interfaces, which makes it possible to integrate it to existing services or build new applications on top of it. For example, our search application prototype (**Figure 1**) takes advantage of the LCD's classification schemes and hierarchical grammar concepts, allowing the user to focus in on specific kind of content in the database.
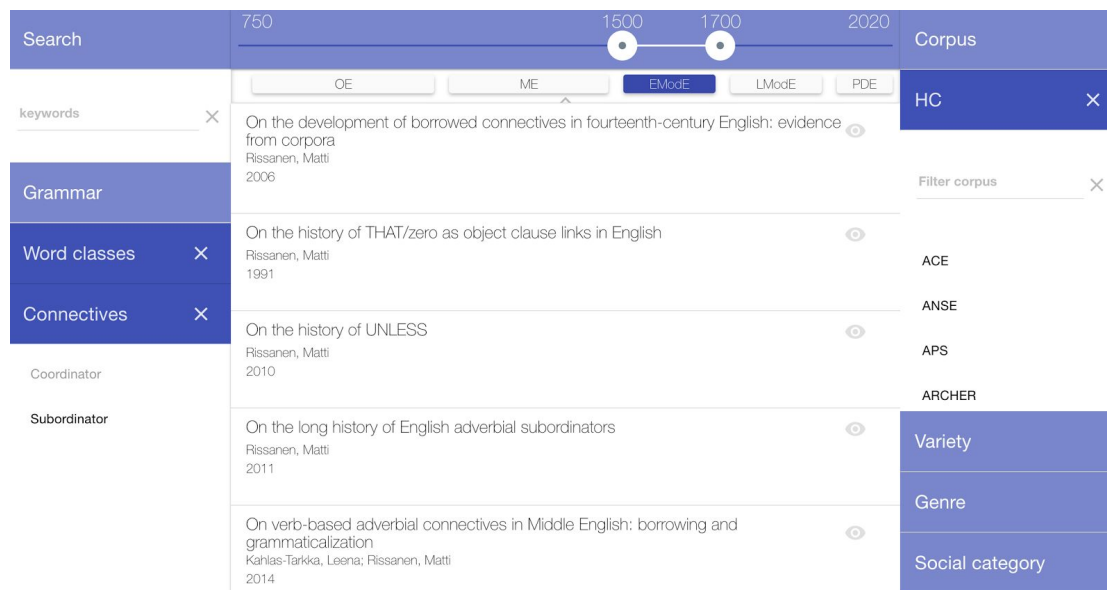


**Figure 1.** Search interface prototype.

The data tables included in the LCD are currently being annotated to make the different kinds of information in them machine-readable. The annotated tables can be used by the **LCD**

**Aggregated Data Analysis** workbench (LADA), an application which provides researchers with a systematic workflow to perform exploratory meta-analyses based on earlier research results. The LADA workflow takes as its input a set of LCD publications collected via the aforementioned search application and deemed relevant to the research question at hand. This initial set of data is then filtered, reviewed (**Figure 2**) and normalized in order to create a new aggregated dataset, which can then be visualised or exported as raw data.
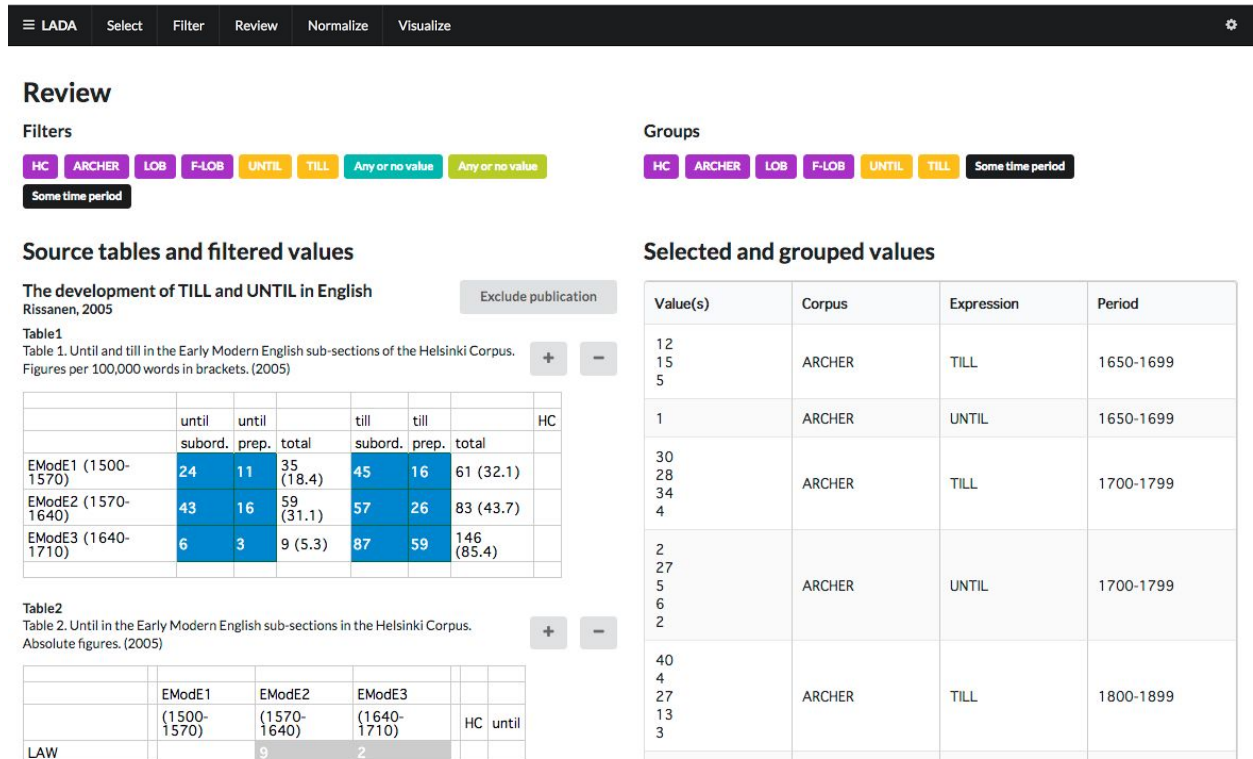


**Figure 2.** LADA review stage. The left-hand column illustrates the filtered tables and individual cells, while the right-hand column shows the new aggregated data table with user-defined groupings and dimensions.

All the original data as well as the datasets that are generated in the different stages of the workflow are stored and interlinked to create comprehensive provenance information about the experiment. The whole experiment can be exported and re-run by other researchers, for example, to validate, extend or comment on the results.

The LCD aims to be a collaborative, trusted and accumulative source of linked research data, and the LADA meta-analysis tool is an example of how this innovative approach to data-driven aggregation of empirical findings can be used in the context of historical linguistics. The LADA tool is currently at an early prototype stage, but we hope to expand it with more analytical functionality in the future.

**Reference**
Nevalainen, Terttu, Turo Vartiainen, Tanja Säily, Joonas Kesäniemi, Agata Dominowska and Emily Öhman. 2016. Language Change Database: A new online resource. *ICAME Journal* 40: 77–94. doi:10.1515/icame-2016-0006