# *TalkBankDB*: A Comprehensive Data Analysis Interface to TalkBank

John Kowalski, Brian MacWhinney

## About TalkBank

**TalkBank**, a **CLARIN B Centre**, is the host for a collection of multilingual multimodal corpora designed to foster fundamental research in the study of human communication.

It contains tens of thousands of audio and video recordings across many languages linked to richly annotated transcriptions, all in the CHAT transcription format.

The origins of TalkBank trace back to 1984 with the creation of the CLAN (Child Language Anlaysis) tools and the associated CHAT transcription format.

The corpus began with annotated media of child language acquisition (CHILDES database) and has expanded to include fourteen annotated media language databases including:
- SLABank for studying second-language acquisition.
- CABank for conversational data.
- ClassBank for study of language in the classroom.
- SamtaleBank for the study of Danish conversations.
- Along with a series of clinical databanks for aphasia, stuttering and other disorders.

The size and scope of TalkBank continues to expand. As of this writing, TalkBank includes over 5TB of richly annotated media.

## TalkBank Usage

| | CHILDES | AphasiaBank | PhonBank | FluencyBank | HomeBank | TalkBank |
|---|---|---|---|---|---|---|
| Age (years) | 30 | 10 | 7 | 1 | 2 | 14 |
| Words (millions) | 59 | 1.8 | 0.8 | 0.5 | audio | 47 |
| Linked Media (TB) | 2.8 | 0.4 | 0.7 | 0.3 | 3.5 | 1.1 |
| Languages | 41 | 6 | 18 | 4 | 2 | 22 |
| Publications | 7000+ | 256 | 480 | 5 | 7 | 320 |
| Users | 2950 | 390 | 182 | 50 | 18 | 930 |
| Web hits (millions) | 5.0 | 0.5 | 0.1 | 0.1 | 0.4 | 1.7 |

## TalkBankDB Goals

The purpose of TalkBankDB is to provide an intuitive on-line interface for researchers to explore TalkBank's media and transcripts, specify data to be extracted, and pass these data on to statistical programs for further analysis. We hope to:

**Increase accessibility:**
Previously, browsing TalkBank required knowing the name of a corpus or area of research, finding its location within the talkbank.org domain (ex: fluency.talkbank.org), then browsing/downloading the media and annotations and installing the CLAN tools. To make these resources more accessible, TalkBankDB provides a single online interface to query across all of TalkBank to find relevant corpora and links to media and transcripts. Users can visually explore data directly in the browser, and if desired, download it for further analysis in a statistical software package.

**Grow community and foster collaboration:**
With the entirety and richness of TalkBank freely accessible from a simple web interface, resources that were previously known only by advanced users will be open to a broader community. Features such as word usage, utterance length, measures of language acquisition speed and ability by demographics can easily be selected, output, plotted, and analyzed through the web interface. By also providing a GitHub account link for users to upload scripts and analyses, the TalkBankDB site provides a single point where users can explore, share their research, and see what others are doing in the TalkBank community.

## Example TalkBankDB WorkFlow

### Build Query to Database



### Click Link in Results to Play Media and View Transcripts



### Explore Data with Visualizations



### Save Data for Import to Statistical Software (R/NumPy/…)



## Technical Details

Creation of the TalkBankDB database relies on the fact that all TalkBank transcripts are pure UTF-8 text files that explicitly implement the CHAT annotation format. These files are then processed by the CHATTER Java program. CHATTER can convert a CHAT file to XML that can be round-tripped back to the file's original CHAT format. The XML format and the associated schema facilitates use of TalkBank corpora by third party programs and systems, eliminating the need to parse complex raw strings.

Since JSON can be used directly by front-end web apps, we eliminate the need of the app to constantly convert XML to JSON and back again by first converting the XML transcripts outputted by CHATTER to JSON using xml-js (Nashwaan, 2018). This tool supports bidirectional XML/JSON conversion. So, combined with CHATTER, we can round-trip from JSON to the original CHAT formatted transcript.

To store our collection of JSON documents, we use MongoDB, a widely-used free and open-source document database. An added benefit of this document database is it makes scaling to increasing data demands easy by allowing the database to scale out across multiple inexpensive machines through "sharding" of the database. This can be difficult to do with relational databases, where often the only option is to "scale up" by purchasing increasingly powerful machines. The scaling-up strategy is not always possible, and can one day be unable to meet the growing size and computational demands of the database.

The front end web interface is written in standard HTML, CSS, and JavaScript to ensure cross-browser support. Visualizations are through the C3.js library, although many others can be used. Care is taken so the JavaScript code is clearly commented and maintainable, following the popular "web component" design pattern common in many large-scale web apps.

## Future Work

*Started only this year, TalkBankDB is still in its early stages. Development is very active, so there are many new features coming soon, including:*

**Query builder:**
- Currently, the query builder ANDs between query categories (corpus, language, media type, gender, age) and ORs within categories. We are working on a UI allowing users to build any logical combination of features (AND, OR, NOT) with expressivity similar to the Corpus Query Language (CQL).

- Additionally, we will expand beyond the current five query categories (ex: morphological properties).

**UI to Add/Edit documents with DB versioning:**
- This will necessitate authenticating users to authorize DB modifications.
  - Additionally, each download will include a version number for reproducibility of analyses.

**Additional Visualizations and Analyses:**
- The code for each visualization is encapsulated in its own JavaScript module (file).
- Many other visualizations can be made following this module's pattern.
  - Users can contribute interesting visualizations.

**Direct access to TalkBankDB from statistical software:**
- As a convenience, we plan to make available R and Python libraries for querying and downloading data from TalkBankDB directly, without having to go through the TalkBankDB web interface, download, then include.

## References

[Chang 2017] Chang, F. (2017) The LuCiD language researcher's toolkit [Computer software]. Retrieved from http://www.lucid.ac.uk/resources/for-researchers/toolkit/

[MacWhinney 2000] MacWhinney, B. (2000). The CHILDES Project: Tools for Analyzing Talk. 3rd Edition. Mahwah, NJ: Lawrence Erlbaum Associates

[MacWhinney 2014] MacWhinney, B. (2014). The childes project: Tools for analyzing talk, volume ii: The database. Psychology Press.

[Nashwaan 2018] Nashwaan, Yousuf, xml-js, (2018) GitHub repository, https://github.com/nashwaan/xml-js

[Sanchez] Sanchez, A., Meylan, S., Braginsky, M., MacDonald, K., Yurovsky, D., & Frank, M. C. (in prep). childes-db: a flexible and reproducible interface to the Child Language Data Exchange System (CHILDES). Manuscript in preparation.

## Acknowledgments