# A New Web Interface to Facilitate Access to Corpora: Development of the ASLLRP Data Access Interface*

**lrec 2012 istanbul** May 27, 2012

Workshop on the Representation and Processing of Sign Languages: Interactions between Corpus & Lexicon

**Carol Neidle**
carol@bu.edu
**Boston University**

**Christian Vogler**
christian.vogler@gallaudet.edu
**Gallaudet University**

## Abstract

A significant obstacle to broad utilization of corpora is the difficulty in gaining access to the specific subsets of data and annotations that may be relevant for particular types of research. With that in mind, we have developed a Web-based Data Access Interface (DAI), to provide access to the expanding datasets of the American Sign Language Linguistic Research Project (ASLLRP). The DAI facilitates browsing the corpora, viewing videos and annotations, searching for phenomena of interest, and downloading selected materials from the website. The Web interface, compared to providing videos and annotation files off-line, also greatly increases access by people who have no prior experience in working with linguistic annotation tools, and it opens the door to integrating the data with third-party applications on the desktop and in the mobile space.

## Currently available data sets

**National Center for Sign Language and Gesture Resources (NCSLGR) Corpus:** ASL videos collected and linguistically annotated at Boston U.

* **Synchronized video files** (compressed, uncompressed formats) show signing from front (stereo views) and side, including face close-up.
* **Linguistic annotations of manual and non-manual components** of the signing, carried out using SignStream®, are available in XML format.
  * Start and end points of each sign, identified by unique gloss label.
  * Parts of speech.
  * Start and end points of non-manual behaviors, also labeled with respect to the linguistic information they convey (e.g., negation, question marking, etc.).
* **Language materials: Total of 1,887 linguistically annotated utterances**, including 1,920 distinct canonical signs (grouping close variants) and 11,861 sign tokens. Contained in 19 short narratives (1,002 utterances) plus 885 additional elicited utterances from Deaf native signers of ASL (most from 4 signers).

## New download capability

User can browse/search the dataset, and select materials of interest, to be added to the "download cart" for subsequent download.

## Best practices and lessons learned

* **Presentation of data**
  * Presenting signs as still images (ideally of both start and end frames for a given sign) saves users and annotators time and effort.
    ◇ The user can detect at a glance whether item is of interest.
    ◇ This facilitates consistency checks.
* **Resource management**
  * Keep metadata separate from file names and assets.
    It is much easier to keep metadata consistent and up-to-date if encoded in a centralized spreadsheet or database rather than encoded in filenames, as metadata can change.
  * Designate only one asset as the authoritative source on metadata, and auto-generate other assets from there.
  * Separate file location and names.
    Files can move, as systems are upgraded, or redundancy is built in. If the location is encoded separately, only this part needs to be updated, rather than every link to a file. The DAI uses a 2-part schema: <url prefix> <path to file> where URL prefix points to a location on the server. Moving the collection to a different location entails updating only a single row in the table that contains the affected URL prefix.
  * Be mindful of cross-platform issues (e.g., restrictions on filenames).
* **Development processes**
  * Plan for continuity.
  * Use version control on all source files *and* third-party libraries.
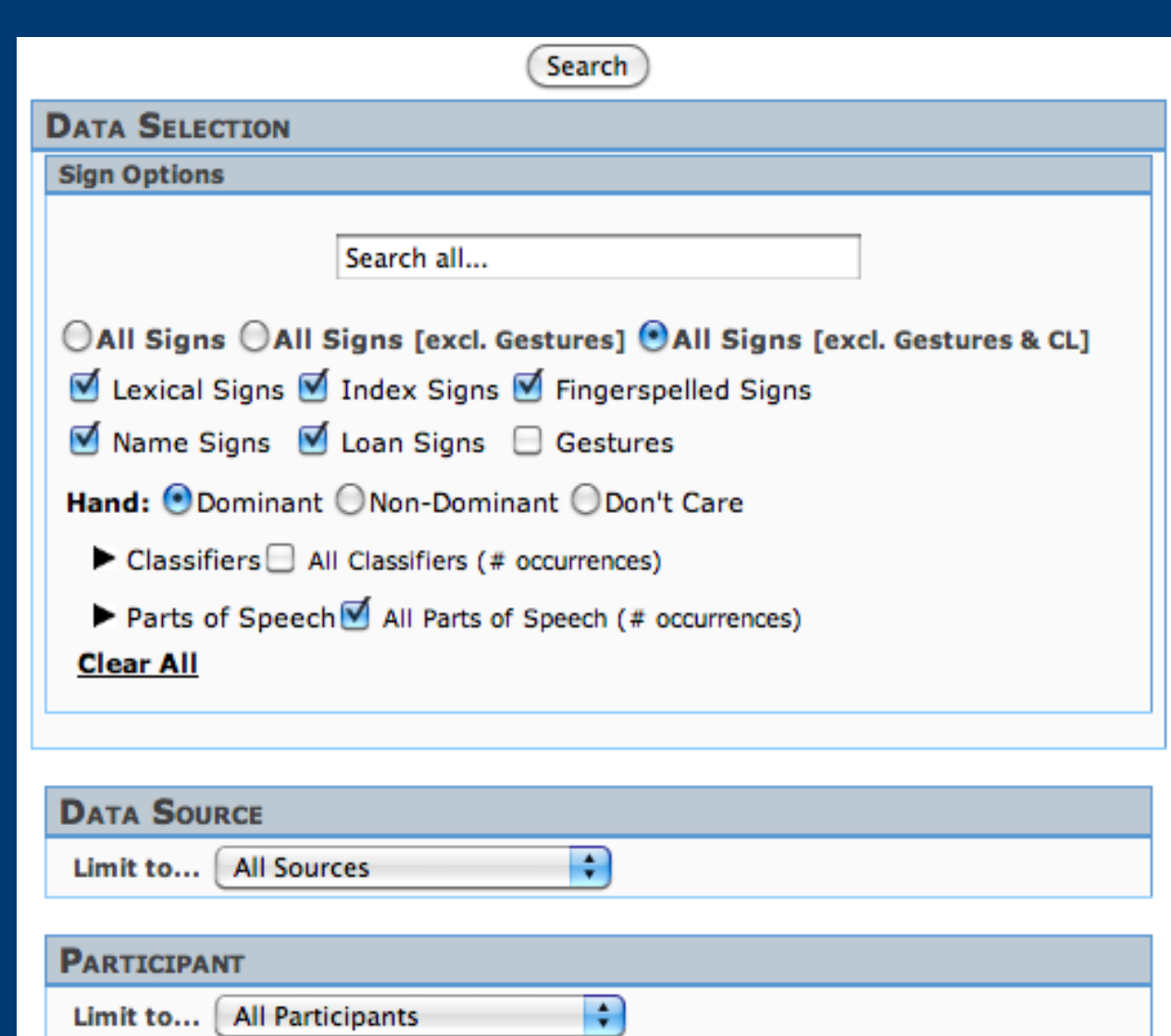* **Database design**
  * Think queries, not data format.
    The types of queries that need to be supported drive database design and tables. They inform every decision that pertains to the tables, the relationships among tables, database views, and choice of indices, and can result in a representation of the annotations that is markedly different from the one chosen for the annotation file format.
  * Use a collection of tags.
    Standardized tagging of annotations provides a powerful and efficient way to search for specific linguistic phenomena.

## Resources

* **Database Access Interface (DAI):** http://secrets.rutgers.edu/dai/queryPages/
* **XML file format and DTD:**
  http://www.bu.edu/asllrp/ncslgr-for-download/download-info.html
* For **documentation of annotation conventions**:
  see reports 11 and 13 from http://www.bu.edu/asllrp/reports.html.

## Functionalities of the Interface



(1) Sample search query
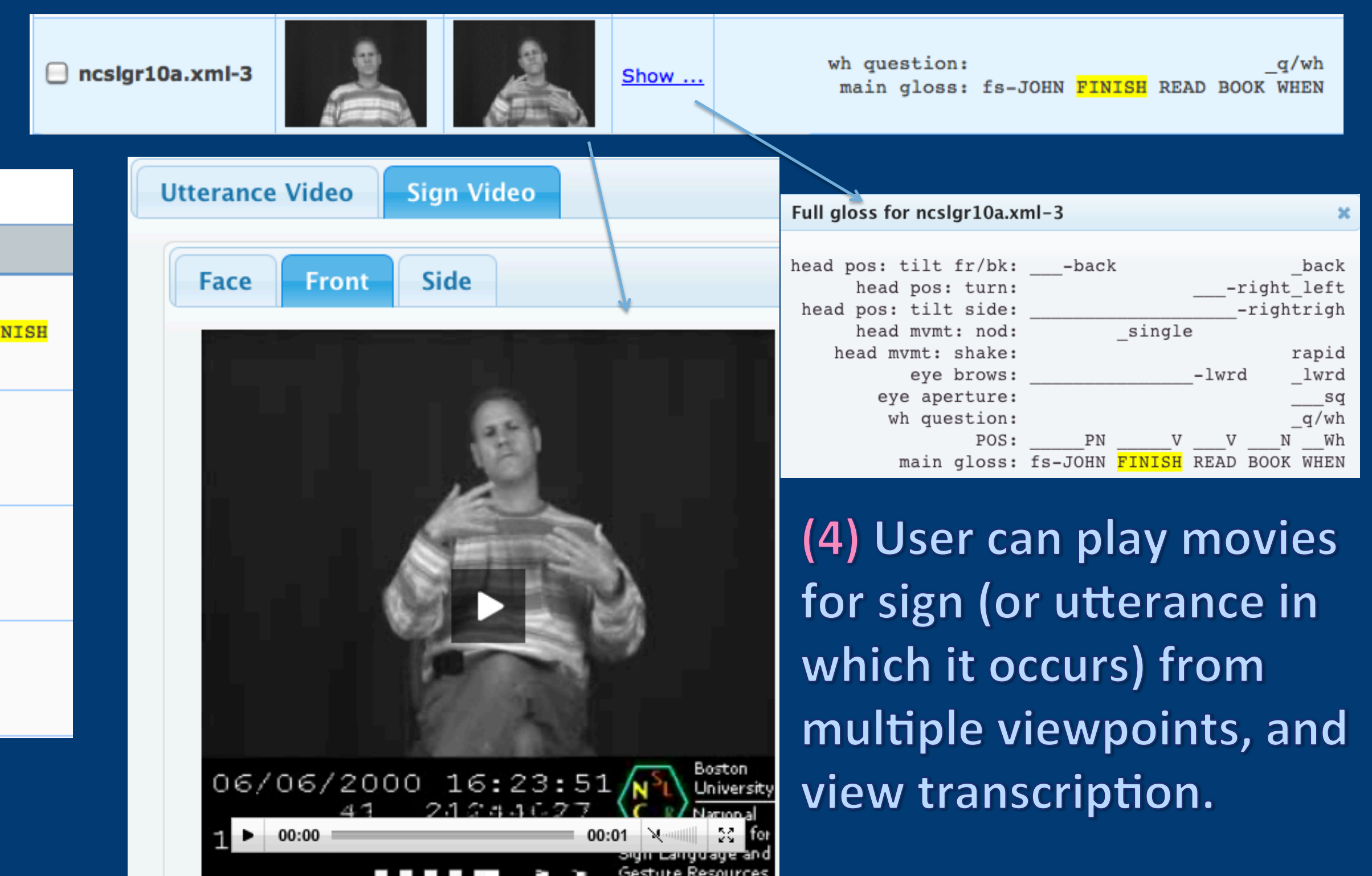


(2) Subset of results



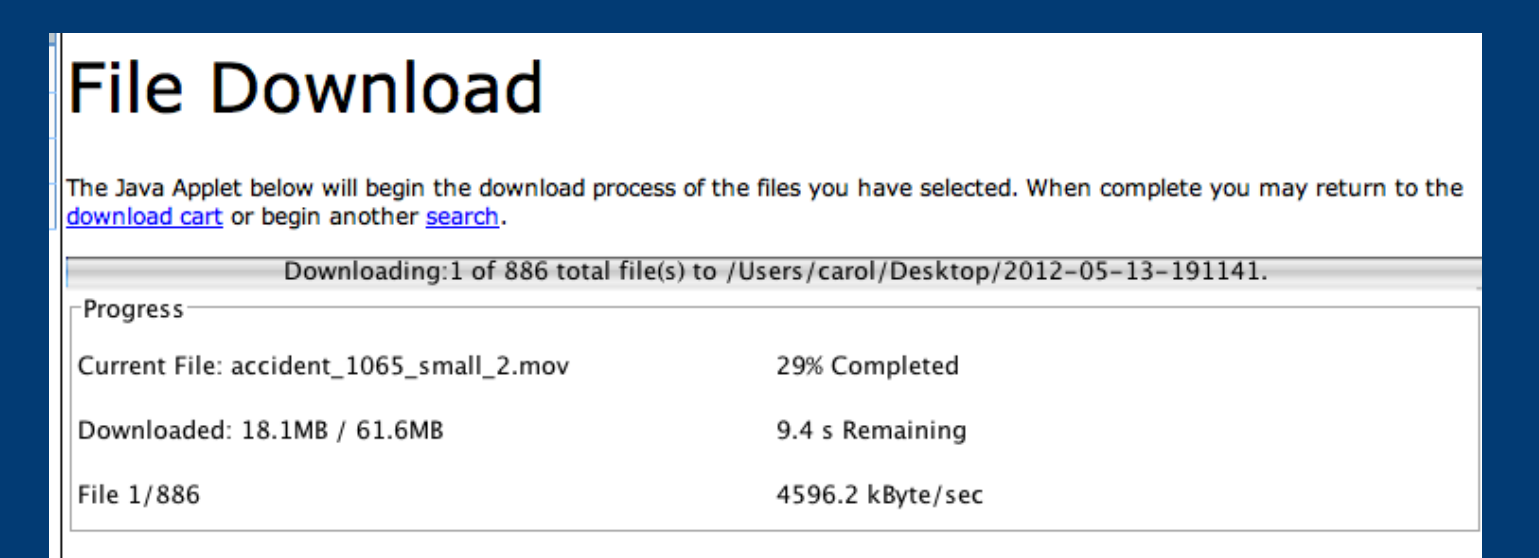(3) Sentences with FINISH: selections for download



(4) User can play movies for sign (or utterance in which it occurs) from multiple viewpoints, and view transcription.



(5) Download cart



(6) Download in progress

## Plans for future development

* Integration of other types of corpora (e.g., ASLLVD corpus — cf. Neidle, Thangali & Sclaroff, "Challenges in Development of the American Sign Language Lexicon Video Dataset (ASLLVD) Corpus," LREC 2012).
* Functionalities to enable additional types of searches – based on linguistic information, statistical frequencies, and video properties.
* Providing annotations in additional formats, e.g., ELAN EAF format.
* Display of various kinds of statistical information.
* Integration of new technologies as they become available (e.g., to provide accessibility for mobile devices).

## *Acknowledgments