# Workshop Proceedings

# 7[th] Workshop on the Representation and Processing of Sign Languages:
# Corpus Mining

# Language Resources and Evaluation Conference (LREC)
# Portorož, Slovenia, 28 May 2016

# Editors and Workshop Organizers

Eleni Efthimiou      Institute for Language and Speech Processing, Athens GR

Stavroula-Evita Fotinea      Institute for Language and Speech Processing, Athens GR

Thomas Hanke      Institute of German Sign Language, University of Hamburg, Hamburg DE

Julie Hochgesang      Gallaudet University, Washington US

Jette Kristoffersen      Centre for Sign Language, University College Capital, Copenhagen DK

Johanna Mesch      Stockholm University, Stockholm SE

# Workshop Programme Committee

Penny Boyes Braem      Center for Sign Language Research, Basel CH

Annelies Braffort      LIMSI/CNRS, Orsay FR

Onno Crasborn      Radboud University, Nijmegen NL

Athanasia-Lida Dimou      Institute for Language and Speech Processing, Athens GR

Sarah Ebling      Institute of Computational Linguistics, University of Zurich, Zurich CH

Eleni Efthimiou      Institute for Language and Speech Processing, Athens GR

Michael Filhol      CNRS–LIMSI, Université Paris-Saclay, Orsay FR

Stavroula-Evita Fotinea      Institute for Language and Speech Processing, Athens GR

Thomas Hanke      Institute of German Sign Language, University of Hamburg, Hamburg DE

Julie Hochgesang      Gallaudet University, Washington D.C. US

Matt Huenerfauth      Rochester Institute of Technology, Rochester, NY US

Hernisa Kacorri      City University New York (CUNY), New York NY US

Athanasios Katsamanis      Computer Vision, Speech Communication and Signal Processing Group, National Technical University of Athens, Athens GR

Jette Kristoffersen      Centre for Sign Language, University College Capital, Copenhagen DK

John McDonald      DePaul University, Chicago IL US

Johanna Mesch      Stockholm University, Stockholm SE

Carol Neidle      Boston University, Boston MA US

Rosalee Wolfe      DePaul University, Chicago IL US

# Community Input on Re-consenting for Data Sharing

**Deborah Chen Pichler[1,3], Julie Hochgesang[1,3], Doreen Simons[2,3], Diane Lillo-Martin[2,3]**

[1]Gallaudet University; [2]The University of Connecticut, [3]Haskins Laboratories
[1]Washington, DC 20002; [2]Storrs, CT 06269-1145
E-mail: {deborah.pichler, julie.hochgesang}@gallaudet.edu; {diane.lillo-martin, doreen.simons}@uconn.edu

## Abstract

Development of large sign language corpora is on the rise, and online sharing of such corpora promises unprecedented access to high quality sign language data, with significant time-saving benefits for sign language acquisition research. Yet data sharing also brings complex logistical challenges for which few standardized practices exist, particularly with regard to the protection of participant rights. Although some ethical guidelines have been established for large-scale archiving of spoken or transcribed language data, not all of these are feasible for sign language video data, especially given the relatively small and historically vulnerable communities from which sign language data are typically collected. Our primary focus is the process of re-consenting participants whose original informed consent did not address the possibility of sharing their video data. We describe efforts to develop ethically sound, community-supported practices for data sharing and archiving, summarizing feedback collected from two focus groups including a cross-section of community stakeholders. Finally, we discuss general themes that emerged from the focus groups, placing them in the wider context of similar discussions previously published by other researchers grappling with these same issues, with the goal of contributing to best-practices guidelines for data archiving and sharing in the sign language research community.

Keywords: language documentation and long-term accessibility for sign language data; experiences in building sign language corpora, ASL, child acquisition

## 1.   Introduction

Development of large sign language corpora is on the rise, and online sharing of such corpora promises unprecedented access to high quality sign language data. For researchers studying early language development, having ready access to longitudinal video data means that many research questions can be tested immediately, on data from multiple children, without the time-consuming prerequisite of subject recruitment, filming and video annotation over the relevant age range. Considering the time and effort required to collect and process longitudinal data from just a single child, the time-saving benefits of shared online corpora clearly has potential to revolutionize the way sign language acquisition research is conducted ("economization of resources" as described in Himmelmann, 2006).

Yet the same long-term data infrastructure that promises such accessibility also brings with it complex logistical challenges for which few standardized practices currently exist. Some of the greatest challenges revolve around the protection of participant rights. Although some ethical guidelines have been established for large-scale archiving of spoken or transcribed language data (e.g. the CHILDES database; MacWhinney, 2000), not all of these are feasible for sign language video data, especially given the relatively small and historically vulnerable communities from which sign language data are typically collected.

Our primary focus in this paper is the process of re-consenting participants whose original informed consent did not address the possibility of sharing their video data. We describe our efforts to develop ethically sound, community-supported practices for data sharing and archiving. Our discussion is focused on video data collected two decades ago from a longitudinal spontaneous production study of the acquisition of American Sign Language (ASL), but the issues and recommendations outlined here are equally relevant to any situation in which video data are shared with a wider audience than initially intended. Below, we introduce the set of longitudinal video data that we plan to share, and outline the anticipated steps for obtaining re-consent from filming participants. We then summarize outcomes of two focus group events in which we sought feedback from a cross-section of community stakeholders. Finally, we discuss general themes that emerged from the focus groups, placing them in the wider context of similar discussions previously published by other researchers grappling with these same issues, with the goal of contributing to best-practices guidelines for data archiving and sharing in the sign language research community.

## 2.   Background

Our immediate context for addressing the issues of this discussion is a body of naturalistic video footage collected longitudinally from four deaf children and their deaf families, between ages 1;05-4;02 (years; months) (Lillo-Martin and Chen Pichler, 2008). The children were filmed in their homes or other familiar locations at intervals ranging from one week to two months. Because all four children were under the age of 5 at the time of filming, their parents provided signed consent for the children's participation. The video data have been painstakingly annotated in different ways over the past twenty years, and "basic transcription" will soon be available for a large portion of the sessions, including ID glosses for individual signs and free translations for all utterances by the target children and their various interlocutors. A screenshot of an example transcript for our project along with text balloons exemplifying our annotation conventions is shown in Figure 1. These basic transcriptions, along with their accompanying video files, are slated for digital archiving in the future at a databank that will be monitored and restricted to academic use, from where they can be shared with researchers pursuing a wide variety of topics related to sign language development.
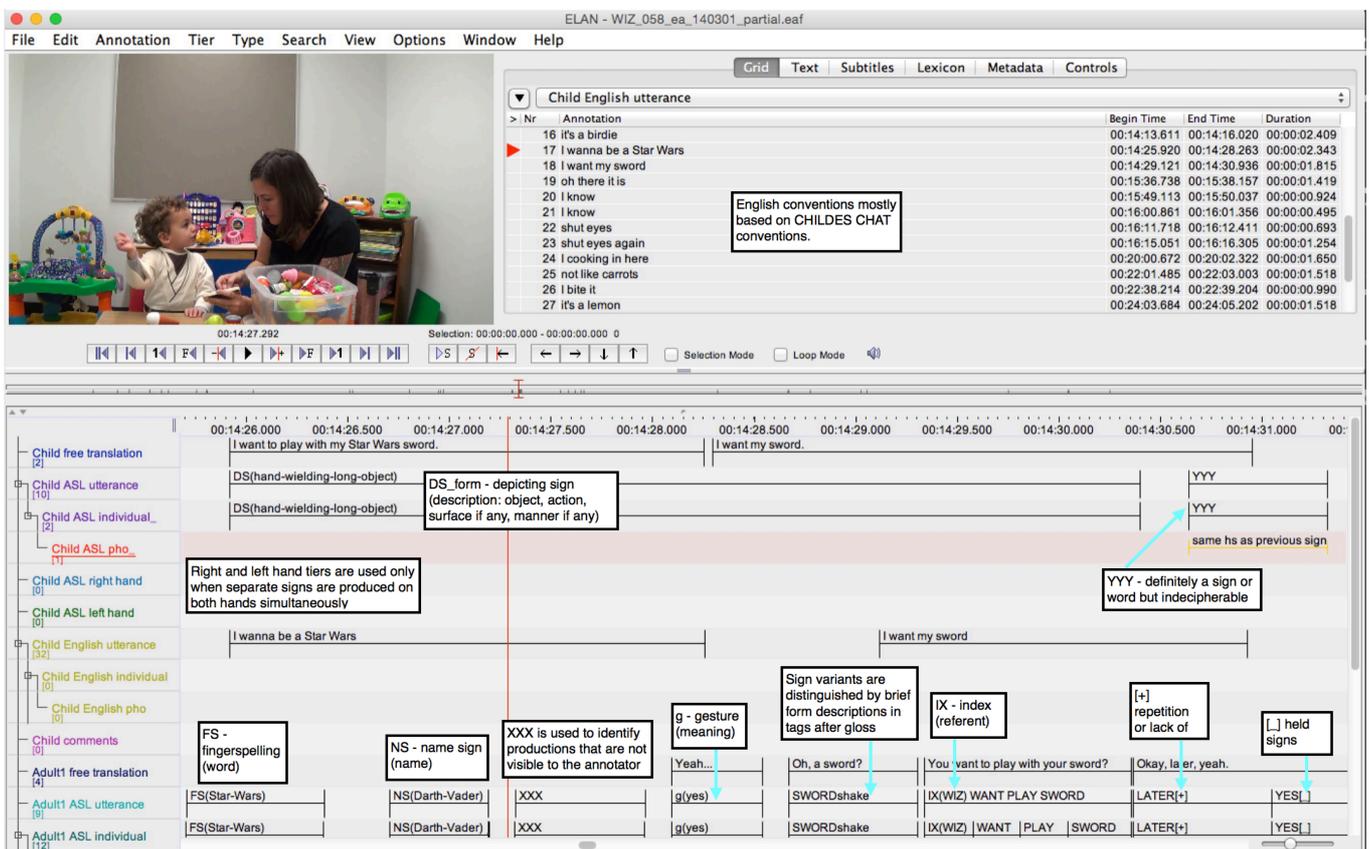
Figure 1. Example of our project's "basic transcript" with ID glosses and free translations

Subsequently, other similar video data of sign language acquisition may also be shared in this way.

However, before sharing the video data and basic transcripts, we must first locate and obtain consent from individuals appearing on video. This re-consenting is no trivial task, given that the data were collected between 1991 and 1999, since which time the target children have grown to adulthood and moved away. The task is further complicated by the many individuals who interacted with the target children on our video footage, ranging from research assistants and the children's immediate family members, to occasional friends and neighbors who appear only sporadically on camera – and in some cases only a portion of their bodies may be visible because they are largely out of the camera's range. Informed consent procedures at the time did not require signed consent from anyone beyond the target children (or parents granting consent in place of target children) so we do not have contact information for most of these "incidental appearances." Thus we must also establish guidelines for determining who requires (re-)consent and what must be done if individuals can not be located or do not grant consent for their video footage to be archived and shared. And finally, we need to determine what measures are deemed necessary by the stakeholder community before they will be comfortable with data archiving and sharing. Individual preferences vary widely, and it is clearly not possible to accommodate the wishes of everyone. Nevertheless, stakeholders in the Deaf community have traditionally had little input on issues of how their video data are used and shared in the long term, so their inclusion in this discussion is critically important.

## 3. Focus Groups

In view of the questions raised here, we convened two focus groups to collect community feedback on issues of data sharing and re-consent. The focus groups took place at Gallaudet University in Washington, D.C., and the American School for the Deaf in West Hartford, Connecticut. Participants were selected from the following groups, identified as stakeholders because of their participation (actual or potential) as subjects or parents of longitudinal filming, and/or their interest as researchers in collecting or analyzing sign language longitudinal data:

1. Deaf of Deaf adults who participated or could have participated in longitudinal video collection for research purposes when they were children
2. Signing family members of Deaf or Koda children
3. Researchers interested in sign language video data
4. Current and former research assistants on projects related to sign language

Each focus group began with a summary of the important role longitudinal data have played in acquisition research and the value of sharing data more widely. It was emphasized that the videos would be shared through online archives maintained by academic institutions, in stark contrast to unmonitored online sharing on YouTube or other forms of social media. Then participants were presented with question prompts targeting selected issues surrounding video data sharing:

- How comfortable are people in the Deaf community with the idea of their videos

30

appearing online? Has the rise of social media made people become more accepting of their videos online or more cautious?

- What types of footage from longitudinal, spontaneous filming might potentially be embarrassing for subjects? How should such cases be addressed?
- If longitudinal videos are shared beyond the original research team, who should have access to them, and what are some ways in which researchers could potentially use them?
- When an outside researcher requests access to shared longitudinal video data, what information should be collected from them? How much of this information do you think should be available to subjects appearing in the shared videos?
- Of the many individuals aside from the target children who appear in longitudinal videos, who should be contacted for re-consent? Does any appearance on video warrant re-consent, or only those that exceed a specific level of frequency and prominence? Should members of the research team (including assistants hired to film and interact with children on video, but not necessarily to analyze the resulting footage) be contacted for re-consent, or is their consent tacit in their role as filmers and experimenters?
- At what age should children be expected to give (re-)consent for longitudinal filming and/or sharing of previously filmed video data?
- What should be done with "orphan works" (O'Meara and Good, 2010), or data from an individual that has died, or otherwise cannot be reached, or does not give consent for video to be used? Should procedures differ for primary subjects (e.g. target children and those with whom they most often interacted) and incidental subjects (e.g. classmates in the background, occasional visitors)?
- What measures and safeguards should researchers establish for Deaf community members to feel comfortable with longitudinal video collection and sharing?

## 4. General Findings (Themes)

The two focus group discussions yielded a wide range of opinions concerning some of the target topics listed in (2), and almost no opinions concerning others, with high variability across viewpoints. This was all expected, given the disparate composition of the focus groups. Nevertheless, several important points emerged for which a consensus or dominant opinion could be identified across one or both focus groups or subsets of focus group participants. In this section, we summarize discussion on three of those points.

### 4.1 General Acceptance of Video Data Sharing

Video data sharing was generally viewed positively, with participants regarding research as valuable to the community. Participants acknowledged the utility of sharing data and were in principle supportive. At the same time, many participants emphasized that Deaf communities are small and participants are never anonymous on film, so sharing of sign language data

requires a higher level of precautions than is typical in the majority spoken language community. Not surprisingly, younger (i.e. 18-25 years old) participants reported less anxiety about the idea of their videos being online, and parents reported much more anxiety about their young children's videos being online than their own. Having personal contact with researchers and periodic updates emerged as a crucial mitigating factor; several parents reported that they trusted the researchers who filmed their children because these researchers had met with them in person to explain what their research goals were and periodically presented updates on their findings. Periodic updates were regarded as more than simple professional courtesy, as expressed by one parent who recalled, "On the consent forms, I checked 'Yes, yes, yes...' straight down the whole page, everything was fine with me. But I expect to be contacted every now and then with updates. Don't come and film us then just disappear." Periodic updates from researchers not only inform participants on what has been done with their data and keeps contact information updated, they also provide a tangible illustration of how the video data benefit their community.

### 4.2 Deciding Who Should be Re-consented and When

Opinions were split on when children should be allowed or expected to give consent. On the one hand, parents were in favor of respecting the wishes of their children, even those who are still minors, if they chose to withdraw their video data from research analysis. At the same time, participants recognized that children and teenagers may not yet fully understand or appreciate the importance of research, so they suggested not destroying any data in these cases, but simply suspending further analysis of them until the subject reached 18 years and had the opportunity to revise their preferences on the consent forms. Opinions were also mixed on whether research assistants on longitudinal filming projects should be re-consented in the same way as target children and their families, or whether individuals in the former group a priori give consent for their video footage to be analyzed and shared when they accept their positions as filmers and experimenters. Clearly, this topic warrants much further discussion; a conservative approach would include them in the re-consenting process.

### 4.3 Measures That Increase Comfort Level with Video Data Sharing

Some participants were willing to allow data to be used by authorized researchers in any scientifically appropriate way. However, others expressed the opinion that videos in which they or their children appear should only be used for research towards certain goals (e.g., promotion of the use of a sign language). Those participants supported the idea of data archives collecting information from any researchers requesting access to video data and making it available to subject families. Suggested information included the researcher's name and institution, research history, and involvement in the Deaf community. This information might be posted on a user list associated with the data archive, to which participants could have access.

Preferences expressed by families on their original video consent forms regarding such questions as whether clips can be used in scientific presentations should naturally extend to any researcher obtaining the video data

from an archive, and families should have the option to change or update those preferences at any time (Harris, Holmes and Mertens, 2009). One parent declared, "It's an exchange. We as research participants give up our privacy in allowing ourselves to be filmed, but in return, the researchers must respect our preferences and wishes." Focus group participants also emphasized the responsibility that researchers have to train their students and assistants in responsible conduct of research, again citing the small size of the Deaf community as reason to be extra sensitive when sharing data from signing families.

As for embarrassing moments on video (e.g. children's temper tantrums (Figure 2); parents losing their temper; parents or children caught on tape in various states of undress; occasional voicing from individuals who normally refrain from using their voice in public, etc.), participants generally agreed that families should have the right to request that certain segments of the video be excised, particularly if those clips have little research significance.



Figure 2. Temper tantrum visible on camera

Since researchers must perform a general review of all video footage as a prerequisite for archiving in a shared database, focus group participants agreed that identification and deletion of embarrassing segments could be undertaken at the same time, and would not necessarily need to involve participant families. Families could indicate the types of activities they would like to have excised (e.g., breastfeeding), and they would then trust the judgment of the researchers in finding and excising appropriate segments.

Another practice that was heavily favored by focus group participants was to include three options in the video release form regarding permission for video clips to be shared or shown in public: a) broad permission; b) only after the subject has the opportunity to view them and give consent for each one; or c) no public viewing (see example item in (1)).

(1) May we show short clips of video footage including you as part of scientific publications resulting from this research?
___ Yes, you may do so without further approval from me.
___ Yes, but only with my prior approval of each clip that you plan to share.
___ No.

This option was included in the video release form signed by all focus group participants, and many commented that it made them more comfortable about giving consent for

their data to be shared. One participant stated, "If the only choices I have on a video release form are "Yes" and "No," that's tough to make a decision...That third option gives me the opportunity to say, "Oh yeah, that clip is fine, you can do whatever you want with it, " or "No, no, that clip is embarrassing, I'd rather keep it private." Having that option eliminates a lot of deliberation, it's really nice."

## 5. Preliminary Recommendations, With Consideration of Previous Proposals

The goal of the two focus group discussions was to sample the varied opinions surrounding archiving and sharing of potentially sensitive video data involving Deaf children and their families, in the hope that they would direct us in the development of guidelines for best practices in this area. As mentioned earlier, other researchers have previously raised similar issues, providing us with a broader context in which to consider our focus group findings. In the field of spoken language acquisition, large databases of archived longitudinal data already exist, perhaps the most notable being the CHILDES database (MacWhinney, 2000). Digital archives of sign language video data are also increasingly common, including at least two that include longitudinal data from child signers (VALID Data Archive (Klatter et al., 2014) and the IPROSLA data sets (Crasborn et al., 2015)). Outside the domain of language acquisition, many researchers working with endangered spoken languages have established archives of digital language resources as part of language documentation and maintenance efforts (Himmelmann, 2006). All of these groups have wrestled with issues of re-consent and data sharing, and in this section, we will comment on how our findings fit with the discussions that have already emerged from those other groups.

Some of the issues listed in (2) are relatively circumscribed, and for those, it is fairly easy to identify existing practices that are directly relevant for our purposes. For example, with respect to (re-)consenting minors, researchers archiving NGT (Sign Language of the Netherlands) data follow the practice supported by our focus groups of not collecting consent from children until they are 18 years old or older; until then, children's parents give consent for them. Baker (2012) notes that once children reach 18, they must have the right to withdraw consent for their data to be used, if they so choose, also in line with the sentiments of our focus groups.

Our focus group participants' views also generally aligned with previous proposals on the topic of data anonymization. In fact, there appears to be virtually unanimous agreement that total anonymization, long taken as a standard practice for medical data, is not feasible for language data that include audio and/or video components. Distortion of faces or voices compromises the usefulness of language data too dramatically to be a viable option (Crasborn, 2010; O'Meara and Good, 2010; Baker, 2012). Crasborn (2010) suggests that development of life-like sign avatars may offer a solution in the future, but for now, a better solution is to accept the fact that sign participants' identities will not be anonymous and establish guidelines to ensure that participants are aware of this fact, and have options to deal with various related

eventualities. Crasborn (2010) details a series of steps taken by researchers on the NGT corpus project to ensure that participants were fully aware of the implications of their video being freely available online for perpetuity before consenting. Additionally, participants were given a DVD copy of their video data after filming, with instructions to review it carefully and inform researchers if there were any segments that they wished to be excluded from sharing. In theory this practice could be applied to the cases of embarrassing footage discussed by our focus groups, although asking participants to review all their footage, which might involve over a hundred videos from a single family, would be less feasible for longitudinal data. An attractive alternative proposed during our focus groups was discussed earlier in the section on "Measures that increase comfort level with video data sharing".

Establishing graded access levels for shared data is another current best practice that would address potential concerns about data sharing, since participants may accept their footage being viewed by the researcher(s) who originally collected the data, but object to the same footage being shared with a wider audience. Among our focus groups, the idea of graded access levels was suggested as a measure that would increase stakeholders' acceptance of large-scale data sharing. Figure 3 shows an example of graded access, taken from the Endangered Languages Archives at SOAS (http://elar.soas.ac.uk, last accessed March 2016).
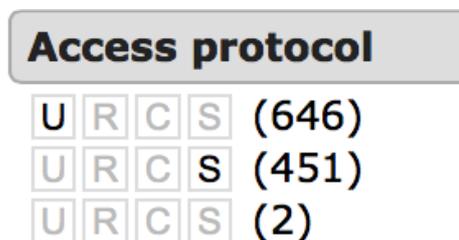


Figure 3. Screenshot of "Access Protocol" from Endangered Languages Archives at SOAS

In the access protocol featured in Figure 2, "U" is for "ordinary user", "R" for "researcher", "C" for "community member", and "S" for "subscriber" (see http://www.elar-archive.org/using-elar/access-protocol.php for more, last accessed March 2016).

One approach to applying graded access is to focus on the qualifications of the researcher requesting access to the data. This is the type of system that was discussed at our focus groups, and it appears to be a common option for other language archives, too (as demonstrated in Figure 2). However, this kind of graded access typically prioritizes access for individuals with ties to the community from which the data were recorded, raising the thorny issue of defining "community" (Leopold, 2013; O'Meara and Good, 2010; Harris, Holmes and Mertens 2009). As a case in point, some focus group participants suggested that only researchers with verifiable ties to the Deaf community (e.g. issuance from a Deaf family, ability to sign, history of working on projects with Deaf community members) should receive full access to archived sign data, presumably because such individuals

are most highly aware of the potential harm that could come to the community if the data are misused. But like any social construct, the boundaries of "the Deaf community" are fluid, depending on who defines them, so determining which researchers possess the requisite community ties will not be straightforward.

The notion of "community" is still relevant, but possibly less problematic, under the second approach to applying graded access, which involves categorizing the data themselves into different levels of accessibility. This is the approach proposed by the CHILDES database (http://childes.psy.cmu.edu, last accessed March 2016), in which each corpus is assigned to one of nine levels, with data at the least restrictive level being fully public and viewable/downloadable without prior registration. At the more restricted levels, researchers may be required to register, submit nondisclosure forms, obtain explicit approval from the original data collectors, or view the data only under direct supervision of someone from the original data collection team (MacWhinney, 2000). For data at any given level, researchers all have the same access, regardless of qualifications. Requirements for data at more restricted levels are made explicit in the database, so prospective researchers can consider them when deciding which corpora to request access to. However graded access is implemented, it would be helpful to establish an Advisory Board to work together with researchers to develop these guidelines. Also, regardless of which type of graded access is instituted in the end, video release forms should also offer the various options for public sharing of specific segments of the data, mentioned earlier.

For the remaining topics in list (2), there was still very little consensus after the focus group discussions: how to deal with "orphaned works;" how to define incidental appearances and whether or not the same re-consent procedures extend to them; whether or not to extend the same re-consent procedures to former research assistants. Opinions on these topics varied widely, some of it probably reflecting age and location. Continued dialogue on these topics is an important step towards developing clear, diversified and actionable protocols, especially since many focus group participants felt that the community has traditionally had very little input on the collection or use of sign language data by researchers. Indeed, the importance of sustained and transparent communication between researchers and research participants can not be overstated, as it is lays the basis for joint efforts across these groups to develop guidelines for video archiving and sharing that are culturally sensitive and balance the benefits of increased access for sign language research with the need to protect individual participant rights.

## 6. Acknowledgments

have participated in our research in the past, as well as the research assistants who have worked with us.

# 7. References

Crasborn, O. (2010). What Does" Informed Consent" Mean in the Internet Age?: Publishing Sign Language Corpora as Open Content. *Sign Language Studies*, *10* 2), pp. 276--290.

Crasborn, O. van Zuilen, M., Fikkert, P., van den Bogaerde, B. & Baker, A. (2015). IPROSLA: Data sets on the acquisition of NGT as a first language. Pposter presented at 2nd ICSLA (International Congress on Sign Language Acquisition), Amsterdam.

Harris, R., Holmes, H. M., & Mertens, D. M. (2009). Research ethics in sign language communities. *Sign Language Studies*, *9*(2), pp. 104--131.

Himmelmann, N. (2006). Language documentation: What is it and what is it good for? In J. Gippert, NP. Himmelmann, & U. Mosel (Eds.), *Essentials of Language Documentation*. New York: Mouton de Gruyter, p 1--30.

Klatter-Folmer, J., Van Hout, R., Van den Heuvel, H., Fikkert, P., Baker, A., De Jong, J., Wijnen, F., Sanders, E. and Trilsbeek, P. (2014). Vulnerability in acquisition, language impairments in Dutch: Creating a VALID data archive. In *Proceedings of LREC 2014: 9th International Conference on Language Resources and Evaluation (http://www.lrec-conf.org/proceedings/lrec2014/index.html)*. pp. 357--364.

Lillo-Martin, D., & Chen Pichler, D. (2008). Development of sign language acquisition corpora. In *Proceedings of LREC 2008: 3rd Workshop on the Representation and Processing of Sign Languages, 6th Language Resources and Evaluation Conference (http://www.lrec-conf.org/proceedings/lrec2008/)*. pp. 129--133.

MacWhinney, B. (2000). *The CHILDES Project: Tools for analyzing talk. 3rd Edition. Vol. 2: The Database*. Mahwah, NJ: Lawrence Erlbaum Associates.

O'Meara, C., & Good, J. (2010). Ethical issues in legacy language resources. *Language & Communication*, *30* (3), pp. 162--170.

# 8. Language Resource References

Baker, A. (2012) Ethics issues in the archiving of sign language acquisition data. Project deliverable from IPROSLA (Integrating and Publishing Resources of Sign Language Acquisition) Archive (http://www.ru.nl/publish/pages/637809/iprosla_deliverable_d2_v3.pdf).