# What makes a corpus easy to (re)use?

The case of language acquisition corpora

Robert Schikowski & Steven Moran

erc

# Introduction

University of
Zurich UZH

**Department of Comparative Linguistics**

# The ACQDIV project

- ACQDIV = **Acq**uisition processes in maximally **div**erse languages
- ERC project, 2014 - 2018/19
- Central question: What's universal in language acquisition?
- Method: Compare acquisition processes in 10 languages that we know to be very different wrt some macrotypological parameters („maximum diversity sampling")
- Data:
  - 11 corpora containing data from children and surrounding adults
  - media recorded in a natural environment (family, home etc.)
  - longitudinal organisation (periodical recordings across ~ 1 year)
  - transcribed, glossed, and (mostly) translated

University of
Zurich UZH

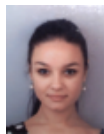**Department of Comparative Linguistics**

# The ACQDIV languages

## The ACQDIV corpora

| Language | Format | Session MD | Speaker MD |
|----------|--------|------------|------------|
| Turkish | Quasi-CHAT | Quasi-CHAT | Quasi-CHAT |
| Japanese (1&2) | Talkbank XML | Talkbank XML | Talkbank XML |
| Indonesian | Toolbox | CHAT | XLS |
| Yucatec | Quasi-CHAT | Quasi-CHAT | Quasi-CHAT |
| Inuktitut | Quasi-CHAT | CHAT | CHAT |
| Chintang | Toolbox | IMDI | IMDI |
| Sesotho | Talkbank XML | Talkbank XML | Talkbank XML |
| Russian | Toolbox | IMDI | IMDI |
| Dene | Toolbox | CSV | CSV |

## The ACQDIV team



Andreas Gerster

Katia Mažara

Danica Pajović

Cazim Hysi

Dagmar Jung

Steven Moran

Sabine Stoll

Robert Schikowski

Laura Canedo

Carolin Remensberger

## Talk map

- Transforming corpus data: general steps
- Reusing corpus text: criteria for easy transformability
  - consistency
  - separation of independent contents
  - documentation
  - explicit coding
- Reusing other kinds of corpus data
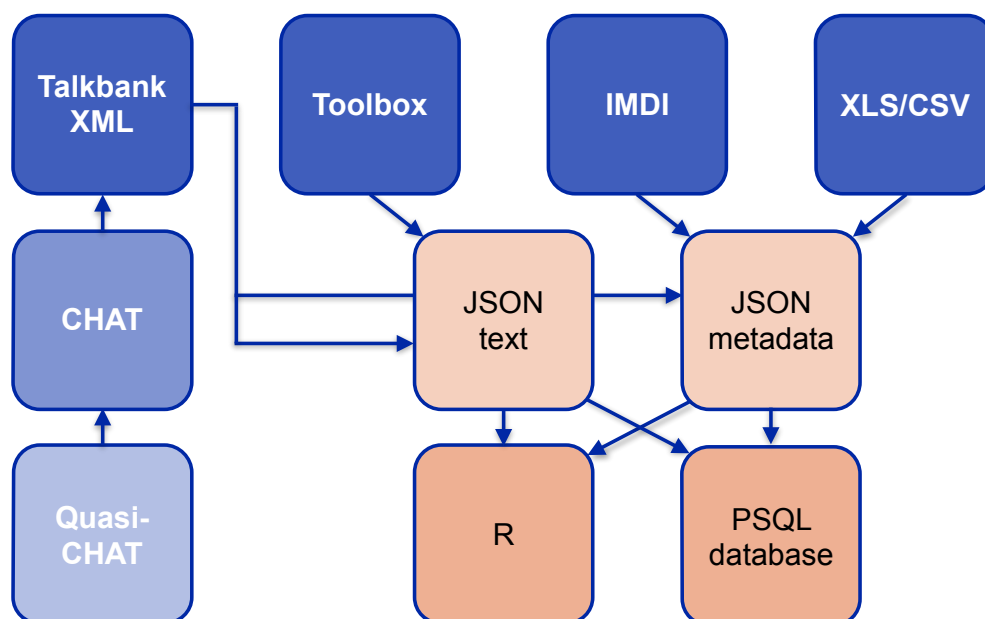- Conclusion

# Transforming corpus data

# Standards are nice, but…

- Corpus standards are a big help for (re)using data, but:
  - all standards only provide a general framework inside of which more or less variation is possible
  - not all standards are equally easy to process
  - data can always be broken if they are not regularly tested for wellformedness
  - linguistic analysis (esp. quantitative analysis) often requires rearrangement of data
  - older data might not be standardized
- ➔ Data transformation is a normal component of any research project involving corpus data!
- ➔ Usability depends on how easy data are to transform, and this question is partially independent of standardization

# Transformations in ACQDIV

# Knowledge required for transforming a corpus

**University of Zurich** UZH

**Department of Comparative Linguistics**

# Reusing corpus text

**Department of Comparative Linguistics**

# Syntactic vs. semantic structure

- Corpora have a complex semantic structure with the following basic elements:
    - **tiers:** mutually independent types of information coded in the corpus, e.g. transcription or translation.
    - **tier content:** the content of individual tiers, e.g. the words constituting a transcription or translation
    - **ordering relations:** temporal or textual precedence between elements on one tier, e.g. word 1 precedes word 2
    - **hierarchical relations:** associations between elements on different tiers, e.g. morphemes are contained in words
- The semantic structure of a corpus is formally coded by the **corpus syntax** (for instance, XML is one kind of corpus syntax).

**Department of Comparative Linguistics**

# Structures and corpus design

- For making a corpus easy to transform, the following points are important:
    - There should be a **1:1 relationship** between elements of the semantic and elements of the syntactic structure.
    - Semantically **independent tiers** should also be kept independent on the syntactic level, i.e. their contents should not be intermingled.
    - The correspondences between the semantic and the syntactic structure should be **richly documented**.
    - The syntactic coding should be **explicit** rather than implicit.
- But why are these points important?
- Some examples from our project experience…

# 1:1 relations between syntax/semantics

- When the rule of 1:1 relations between syntactic and semantic elements is violated, the coding becomes **inconsistent**:
  - Several syntactic elements correspond to 1 semantic element
    ➔ semantic elements with deviant coding may be overlooked; generating syntax may be problematic
  - 1 syntactic element corresponds to several semantic elements
    ➔ syntactic element becomes ambiguous and can't be dealt with automatically at all!

# Example 1: Tier names

# Example 2: Tier names

# Example 3: Line breaks

# Example 4: Coding words

# Keeping tiers independent

- When tiers that are semantically at least partially independent are not kept separately in the corpus syntax, bad things may happen:
  - It becomes difficult or impossible to pull the tiers apart in the semantics.
  - Content belonging to one tier may mistakenly be interpreted as belonging to a different tier.
  - In the worst case all involved tiers may end up being unanalyzable.

## Example 1: Transcriptions and psycholinguistics

## Example 2: Translations and metacomments

## Example 3: Metadata
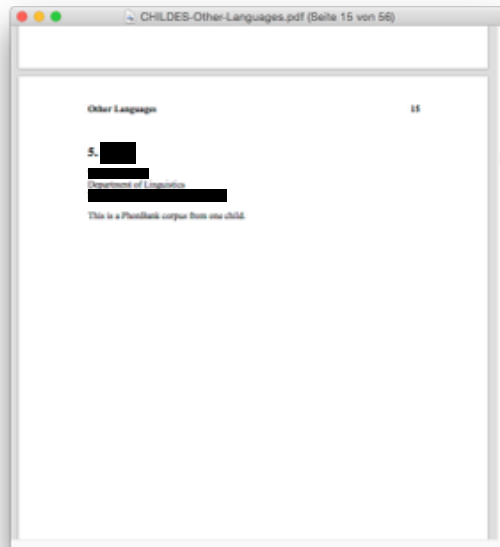
## Documentation

- Documentation matters:
  - Even a very messy corpus may become somehow processable when good documentation is available (although documentation does not help against inconsistency).
  - When documentation is missing, even an exemplary corpus may take a long time to understand and (re)use.
- Documentation is best published together with corpus data and regularly updated.

# Example: The CHAT manual

# Explicit coding

- The syntax of a corpus can be consistent, tidy, and well documented but still be awkward to use when the coding is not explicit.
- Explicit coding includes:
  - visible coding of semantic relations (no use of invisible characters such as spaces or line breaks)
  - as much information as possible is marked by standardised labels
- Implicit coding brings some dangers with it:
  - confusion on the side of the user who hasn't internalised the coding
  - invisible characters are especially prone to inconsistency and other errors

# Example 1: Hierarchical relations

# Example 2: Metadata

# Reusing other kinds of corpus data

University of
Zurich UZH

**Department of Comparative Linguistics**

## Directory structure

- A corpus directory structure that is too complicated can create various problems:
    - outsiders won't understand what's located where
    - the corpus creators themselves might forget, leading to confusion and often more serious problems such as file duplicates
- ➔ **A good directory structure is simple:**
    - it is not too deep (generally one additional level below the top level should be enough)
    - it has transparent folder names
    - folders are disjunct so that the possibility for duplicates is excluded. For instance, sorting files into folders by file types is good because normally no file can belong to two types; sorting files by speakers is bad because one file can be associated with multiple speakers

# File naming

- A bad file naming system can create various problems:
  - data can't be located or eventually even get lost
  - files may get duplicated without anybody noticing
  - it may become opaque which files are instances of the same more abstract corpus unit („text" or „session")
- ➜ **A good file naming system is transparent:**
  - it indicate a few important points for each file (e.g. in the case of language acquisition: recording date and target child of the session)
  - it names all files identically or at least similarly that belong to the same abstract unit
  - it is practical for sorting. For instance, starting files with a recording date in the format YYYY-MM-DD makes it easy to sort files chronologically.

# Encodings

- Encodings are important for linguists because:
  - most languages apart from English use „special characters" (à, щ, 字…)
  - if a file with special characters is opened in an environment which doesn't recognize or even breaks the intended encoding, some or all of them will be 文字化けd.
- ➜ **A good encoding is standardised:**
  - Ideally use UTF-8 (even more ideally with some additions ignored here).
  - If this is impossible for some reason, use an encoding that can code all the graphs required by your language and that is widespread

# Conclusion

- When it comes to reusing corpora, normally some kind of data transformation will be necessary.
- Standards facilitate transformation, especially when they are kept.
- But: some underlying properties of corpora are even more important than that. In this talk we have focussed on:
  - 1:1 relations between syntactic and semantic structure (aka consistency)
  - keeping semantically independent tiers apart in the syntactic representation
  - providing documentation
  - preferring explicit over implicit coding