

Micropublication¹

Michael Cysouw
MPI-EVA Leipzig

¹ Footnotes for the 21st Century

Keep it simple, keep it local

- Instead of doing everything in one system, **split up large databases** into independent modules
- Cooperation can take place by sharing some parts, while other aspects remain proprietary
- The challenge is to define small tasks that can be shared independently of other parts

Proposals

- Languoids & Doculects
- Micropublication
- How can this become part of your database ?!

What is a Linguoid?

Linguoids are *language-like entities*, including

- All kinds of lects:
 - ▶ language, dialect, sociolect, idiolect, stylistic register, etc.
- Genealogical groupings:
 - ▶ all levels, from dialect cluster to stock
- Geographic groupings:
 - ▶ sprachbund, spread zone, macro area, climate zone, continent, etc.

Why languoid?

- To allow us to move forward to investigate “languages” while **avoiding the insoluble problem** of deciding what a “language” is
- A languoid can be catalogued separately from the specification as to what kind of languoid it is
 - ▶ Dialect or language?
 - ▶ Language or small family?
 - ▶ Genealogical or areal group?
 - ▶ Different register or different language?

Defining Linguoids

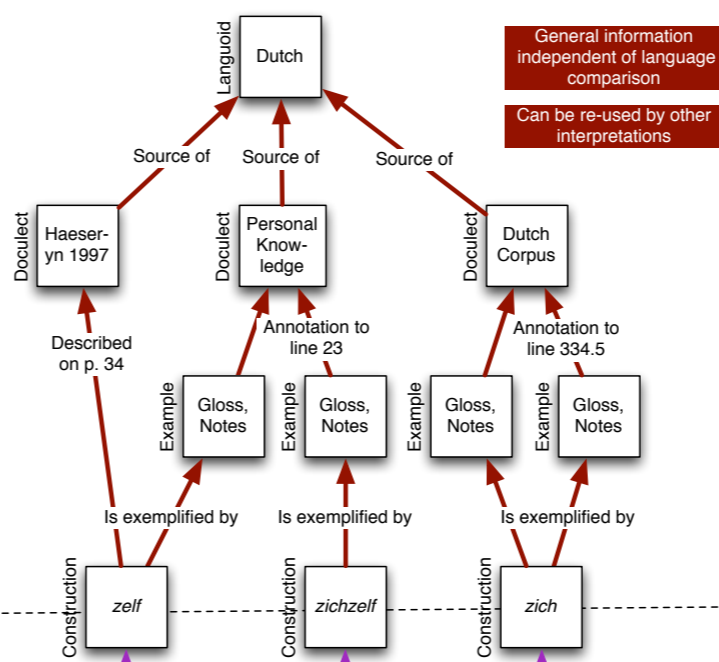
- Linguoids are defined as a set of linguoids
 - ▶ Recursion ends at doculects
- *Doculect*: variety as instantiated by any available documentation
 - ▶ Grammatical description (grammar, article)
 - ▶ Dictionary, wordlist
 - ▶ Inscription, transcription, recording, questionnaire
 - ▶ Description of personal knowledge
 - ▶ Language notes in a traveller's diary
 - ▶ Name given in an ancient text, or in a census

Micropublications

(in the context of linguistic typology)

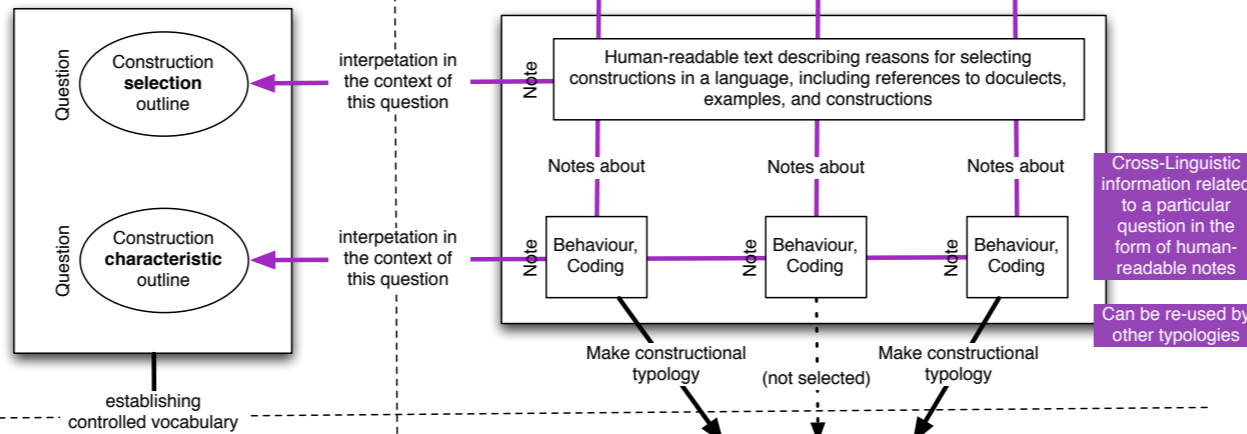
- **Annotations of doculects**
- Interpretation of individual doculect in the context of a research question
- Small notes normally not important enough to be published
- Traditionally, such content would at most appear in a footnote or an appendix

1. Language-particular Information



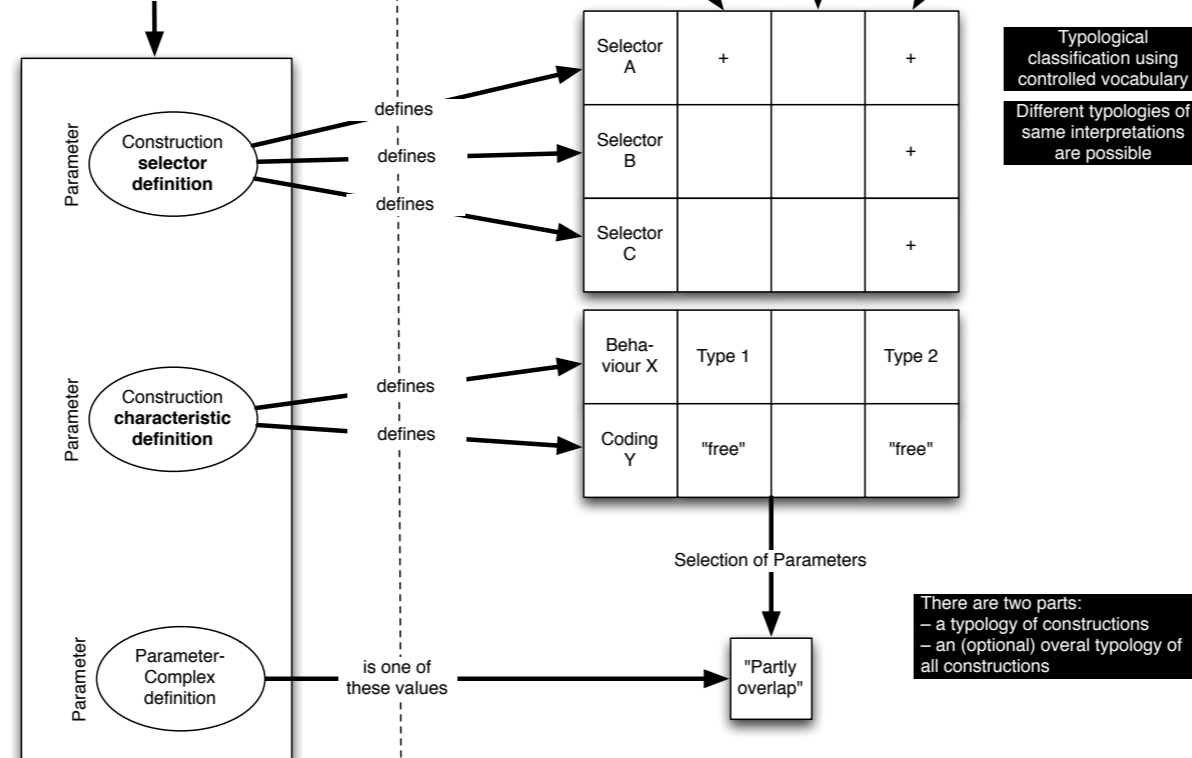
General information independent of language comparison
Can be re-used by other interpretations

2. Cross-linguistic Interpretation



Cross-Linguistic information related to a particular question in the form of human-readable notes
Can be re-used by other typologies

3. Typology

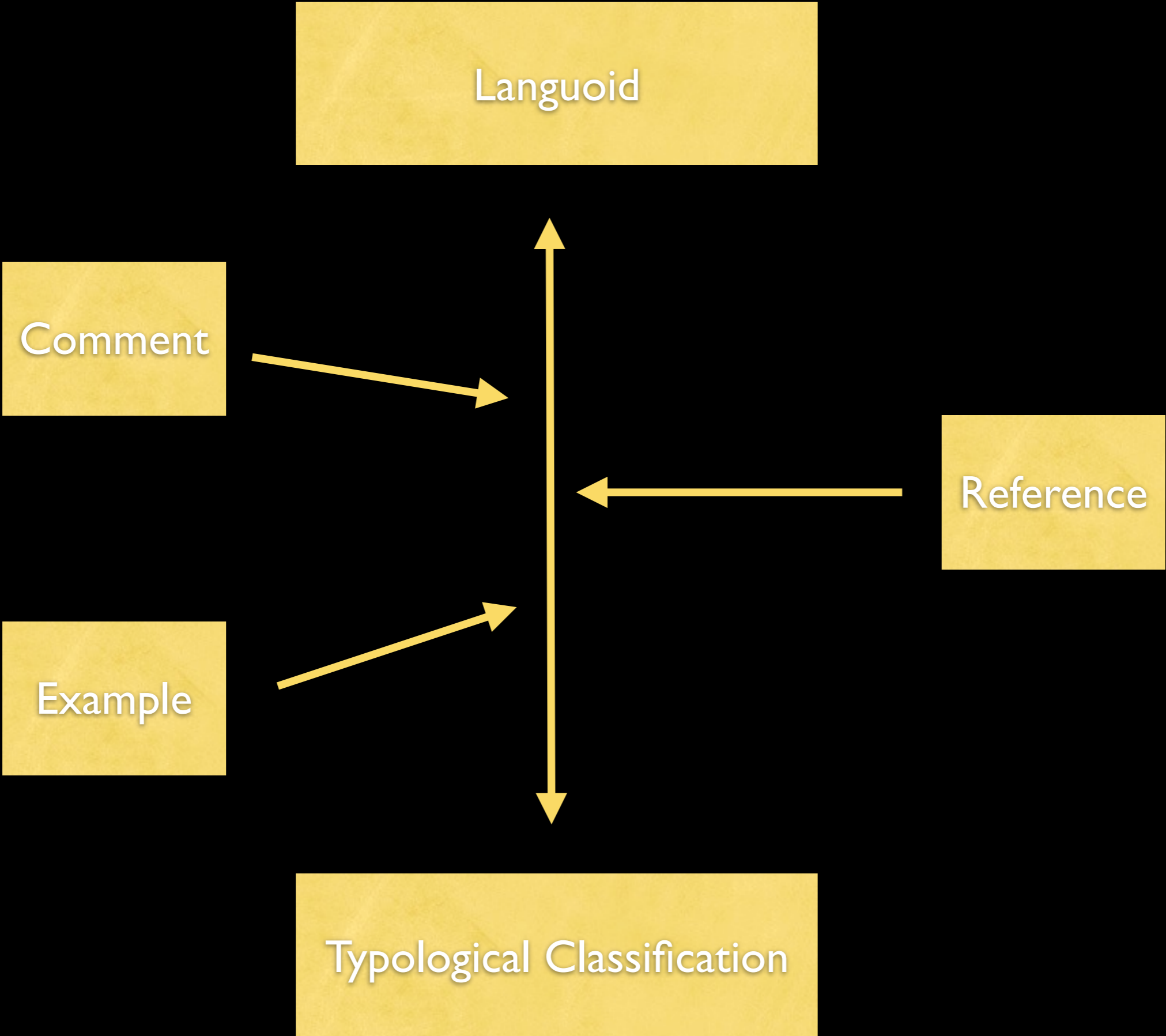


Typological classification using controlled vocabulary
Different typologies of same interpretations are possible

There are two parts:
- a typology of constructions
- an (optional) overall typology of all constructions

Description of comparative concept

Information for each Languoid



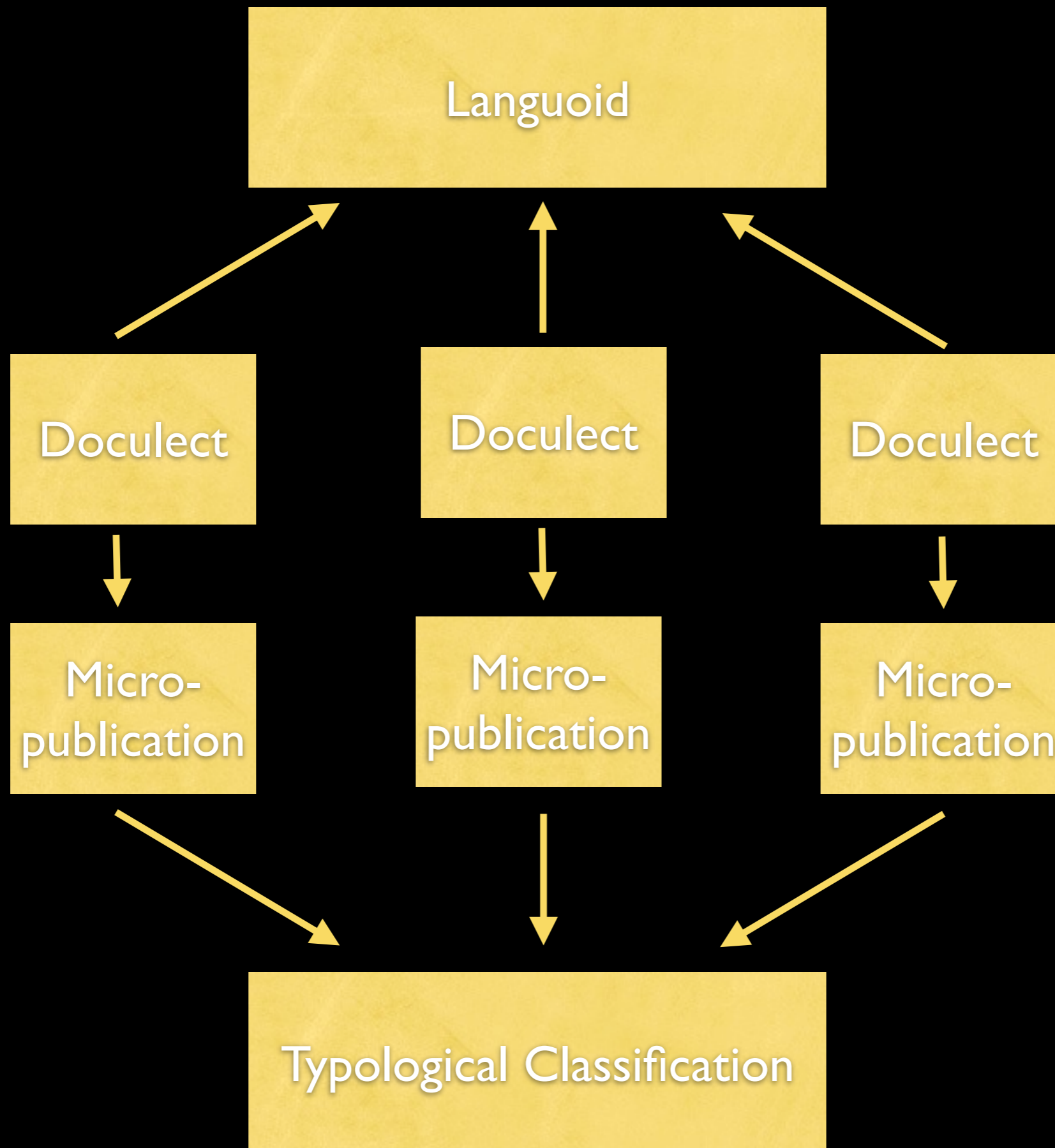
Languoid

Comment

Reference

Example

Typological Classification



Typological databases

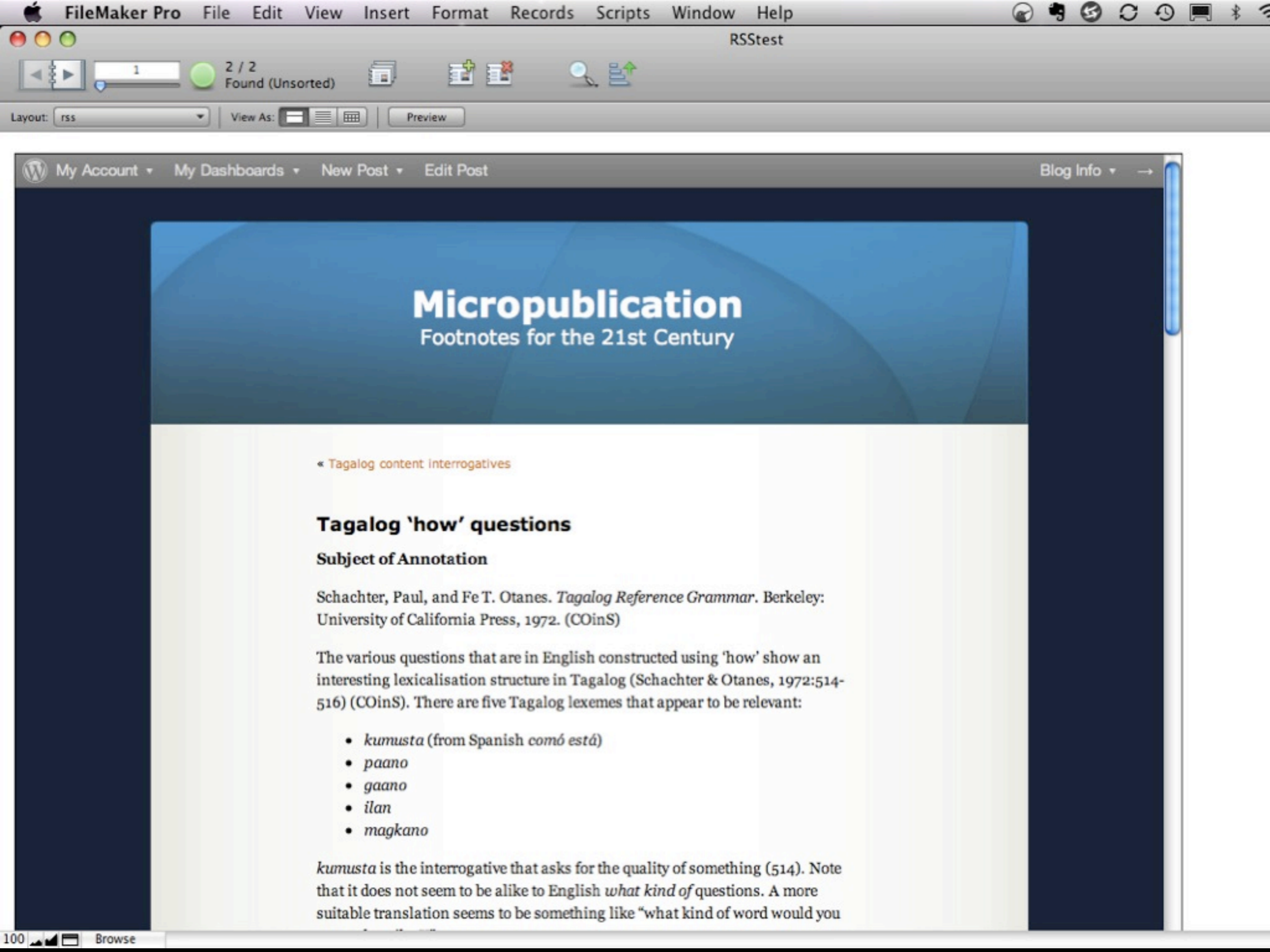
- The core of typological databases are:
 - ▶ Doculects
 - ▶ Micropublications (i.e. annotations of doculects)
- This information is reusable !
- Other levels are more debatable:
 - ▶ Languoids
 - ▶ Typological classifications

Format of Micropublication

- Standard-compliant
- Authored
- Fixed
- Stable
- Referenced

Structure of Micropublication

- Written in natural language
- Annotation to a restricted number of doculects (preferably only one)
- Include snippets !



Navigation icons: back, forward, search, and document management symbols.

Layout: rss View As: [List View Icon] [Table View Icon] Preview

Micropublication

Footnotes for the 21st Century

« [Tagalog content interrogatives](#)

Tagalog 'how' questions

Subject of Annotation

Schachter, Paul, and Fe T. Otanes. *Tagalog Reference Grammar*. Berkeley: University of California Press, 1972. (COinS)

The various questions that are in English constructed using 'how' show an interesting lexicalisation structure in Tagalog (Schachter & Otanes, 1972:514-516) (COinS). There are five Tagalog lexemes that appear to be relevant:

- *kumusta* (from Spanish *cómo está*)
- *paano*
- *gaano*
- *ilan*
- *magkano*

kumusta is the interrogative that asks for the quality of something (514). Note that it does not seem to be alike to English *what kind of* questions. A more suitable translation seems to be something like "what kind of word would you