# Typology: Lecture III
# Databases and Sampling

### Christian Bentz
*University of Tübingen*

May 3, 2017



WORDS BONES GENES TOOLS
Tracking Linguistic, Cultural, and Biological Trajectories of the Human Past

EVOLAEMP
LANGUAGE EVOLUTION: THE EMPIRICAL TURN

## OVERVIEW

PHYLOGENETIC
  Ethnologue (20th edition)
  Glottolog 3.0

STRUCTURE
  World Atlas of Language Structures (WALS)
  AUTOTYP

PHONETICS/PHONOLOGY
  UPSID
  PHOIBLE

LEXICON
  Automated Similarity Judgement Program (ASJP)
  World Loanword Database

OTHERS
  D-PLACE

SAMPLING BIASES

# ETHNOLOGUE (20TH EDITION)

### What is it?

- "A comprehensive reference work cataloging *all of the worlds known living languages*. Since 1951, the Ethnologue has been an active research project involving hundreds of linguists and other researchers around the world. It is widely regarded to be the most comprehensive source of information of its kind."
- **License**: Pay wall

[Simons, Gary F. and Charles D. Fennig (eds.). (2017). *Ethnologue: Languages of the World, Twentieth edition*. Dallas, Texas: SIL International]

# ETHNOLOGUE

## What's in it?

- Languages: **7099**
- Information: Population sizes, language families and genera, endagerement status, etc.

[Simons, Gary F. and Charles D. Fennig (eds.). (2017). *Ethnologue: Languages of the World, Twentieth edition.* Dallas, Texas: SIL International]

# EXAMPLES



Figure: https://www.ethnologue.com/

# EXAMPLES

**Table 2. Distribution of world languages by number of first-language speakers**

| Population range | Living languages | | | Number of speakers | | |
|---|---|---|---|---|---|---|
| | Count | Percent | Cumulative | Total | Percent | Cumulative |
| 100,000,000 to 999,999,999 | 8 | 0.1 | 0.1% | 2,529,403,578 | 40.20547 | 40.20547% |
| 10,000,000 to 99,999,999 | 82 | 1.2 | 1.3% | 2,480,078,977 | 39.42144 | 79.62691% |
| 1,000,000 to 9,999,999 | 304 | 4.3 | 5.5% | 915,659,448 | 14.55462 | 94.18154% |
| 100,000 to 999,999 | 943 | 13.3 | 18.8% | 296,136,843 | 4.70717 | 98.88870% |
| 10,000 to 99,999 | 1,822 | 25.7 | 44.5% | 61,802,734 | 0.98237 | 99.87107% |
| 1,000 to 9,999 | 1,982 | 27.9 | 72.4% | 7,633,408 | 0.12133 | 99.99241% |
| 100 to 999 | 1,065 | 15.0 | 87.4% | 464,299 | 0.00738 | 99.99979% |
| 10 to 99 | 338 | 4.8 | 92.1% | 12,777 | 0.00020 | 99.99999% |
| 1 to 9 | 140 | 2.0 | 94.1% | 560 | 0.00001 | 100.00000% |
| 0 | 206 | 2.9 | 97.0% | 0 | 0.00000 | 100.00000% |
| Unknown | 212 | 3.0 | 100.0% | | | |
| *Totals* | 7,102 | 100.0 | | 6,291,192,624 | 100.00000 | |

Figure: Ethnologue: Population Sizes

# EXAMPLES

# GLOTTOLOG 3.0

What is it?

- ▸ "Comprehensive **reference information** for the world's languages, especially the lesser known languages."
- ▸ **License**: Open Access

[Hammarström, Harald & Forkel, Robert & Haspelmath, Martin. 2017. Glottolog 3.0. Jena: Max Planck Institute for the Science of Human History.]

# GLOTTOLOG 3.0

### What's in it?

- Languages: **8444**
- Language Trees: **242** families, **188** isolates
- Bibliography: more than **180,000** references

[Hammarström, Harald & Forkel, Robert & Haspelmath, Martin. 2017.
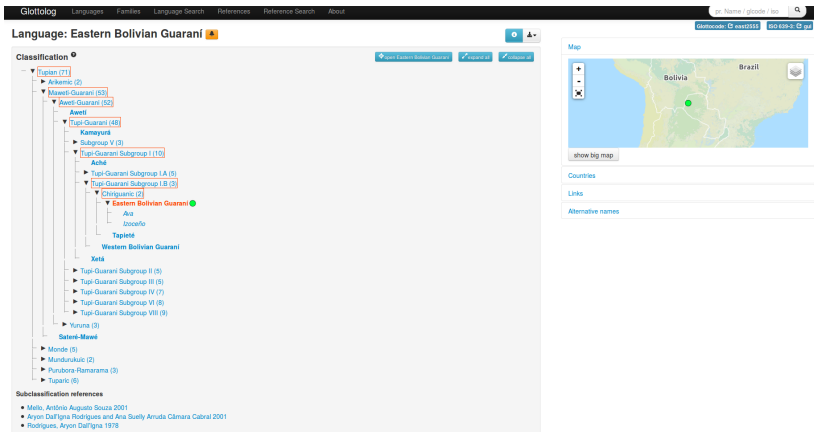Glottolog 3.0. Jena: Max Planck Institute for the Science of Human History.]

# EXAMPLES



Figure: http://glottolog.org/

# EXAMPLES



Figure: Glottolog: Tree for Guarani

# EXAMPLES



Figure: Glottolog: References for Guarani

# THE GLOTTOLOG DATA EXPLORER (HTTPS://CAINESAP.SHINYAPPS.IO/LANGMAP/)



[Caines, Andrew, Christian Bentz, Dimitrios Alikaniotis, Fridah Katushemererwe, and Paula Buttery (2016). The Glottolog Data Explorer: Mapping the worlds languages. Proceedings of the VisLR II Workshop at LREC'16.]

# THE GLOTTOLOG DATA EXPLORER (HTTPS://CAINESAP.SHINYAPPS.IO/LANGMAP/)



[Caines, Andrew, Christian Bentz, Dimitrios Alikaniotis, Fridah Katushemererwe, and Paula Buttery (2016). The Glottolog Data Explorer: Mapping the worlds languages. Proceedings of the VisLR II Workshop at LREC'16.]

# WALS

### What is it?

- "The **World Atlas of Language Structures (WALS)** is a large *database of structural (phonological, grammatical, lexical) properties of languages* gathered from descriptive materials (such as reference grammars) by a team of 55 authors"
- **License**: Open Access

[Dryer & Haspelmath (eds.), 2013]

# WALS

What's in it?

- Languages: **2679** (ca. 37% of the world's languages)
- Features: **192** (Phonology, Morphology, Word Order, etc.)
- Chapters: **151** (Consonant Inventories, Number of Cases, Verbal Inflection, etc.)

[Dryer & Haspelmath (eds.), 2013]

# EXAMPLES



Figure: http://wals.info/

# EXAMPLES



Figure: Chapter 1: Consonant Inventories

PHYLOGENETIC
○○○○○○○○○○○○

STRUCTURE
○○○○●○○○○○○○

PHONETICS/PHONOLOGY
○○○○○○○○○

LEXICON
○○○○○○○

OTHERS
○○○○○○

SAMPLING BIASES

# EXAMPLES



Figure: Feature 1A: Consonant Inventories

# AUTOTYP

### What is it?

- "AUTOTYP is a large-scale research program with goals in both quantitative and qualitative typology. In quantitative typology, we are interested in detecting and explaining geographical distributions of typological features and in producing statistical estimates of universal preferences as well as of genealogical inheritance and areal diffusion potentials."

- **License**: Open Access

[Nichols, Johanna, Alena Witzlack-Makarevich & Balthasar Bickel. 2013. The autotyp genealogy and geography database: 2013 release. http://www.spw.uzh.ch/autotyp/.]

# AUTOTYP

### What's in it?

- ▶ Genealogical and Geographic Information for **2914** languages
- ▶ Phonological, morphological and syntactic information on several hundred languages (though most of this information is published in WALS as well)
- ▶ There is supposed to be a new release this year which should give more extensive grammatical information

[Nichols, Johanna, Alena Witzlack-Makarevich & Balthasar Bickel. 2013. The autotyp genealogy and geography database: 2013 release. http://www.spw.uzh.ch/autotyp/.]

# SSWL

## What is it?

- "The **Syntactic Structures of the World's Languages (SSWL)** is a searchable database that allows users to discover which properties (morphological, syntactic, and semantic) characterize a language, as well as how these properties relate across languages. This system is designed to be free to the public and open-ended. Anyone can use the database to perform queries."
- **License**: Open Access

http://sswl.railsplayground.net/

# SSWL

## What's in it?

**SITE STATISTICS**

Number of Languages: 276
Number of Languages over 90%: 24
Number of Contributors: 414

Number of Properties: 148
Number of Examples: 4555
Number of Property:Value Pairs: 19526

# EXAMPLES

# URIEL

### What is it?

- ▸ "The URIEL knowledge base is a structured compendium of information on language typology and language universals"
- ▸ **License**: Open Access

[Littel, Patrick & Mortensen, David & Levin, Lori (eds.) 2016. URIEL Typological Database Pittsburgh: Carnegie Mellon University (Available online at http://www.cs.cmu.edu/ dmortens/uriel.html, Accessed at 2016-04-20).]

# URIEL

What's in it?

- Geographic and genealogical information
  for **7971** languages
- Binary vectors of grammatical and phonological
  information from PHOIBLE, WALS, and SSWL

[Littel, Patrick & Mortensen, David & Levin, Lori (eds.) 2016. URIEL Typological Database Pittsburgh: Carnegie Mellon University (Available online at http://www.cs.cmu.edu/ dmortens/uriel.html, Accessed at 2016-04-20). ]

# EXAMPLES

## URIEL Typological Database

### Introduction

The URIEL knowledge base is a structured compendium of information on language typology and language universals that is being developed as part of DARPA's LORELEI project.

Download the most recent release here. Read the README in Markdown.

### Releases

- URIEL v0.0
  - Initial release.
- URIEL v0.1
  - Covers more features and languages.
  - Provides more data points.
  - Corrects several crucial bugs.
- URIEL v0.2
  - Covers more features and languages.
  - Provides more data points.
  - Introduces additional data format (.npz) for improved performance.
  - Introduces "mini-grammars" derived from typological features.
- URIEL v0.3.0
  - **Treatment of macrolanguages**. Macrolanguages like Arabic are now included and their feature values assigned in a principled fashion.
  - **Improvements in phylogenetic identifications and geolocations**. These data, located in the geodata directory, are much improved.
  - **Improved feature prediction**. Improvement in the accuracy of feature prediction from 92% to 93%.
  - **Mini-inversion grammars.** Inversion grammars for language pairs in the URIEL knowledge bank now available on request.

    Information on unpacking .tar.xz archives is available here.
- URIEL lang2vec
  - A tool for querying the URIEL database.

### How to Cite

# EXAMPLES

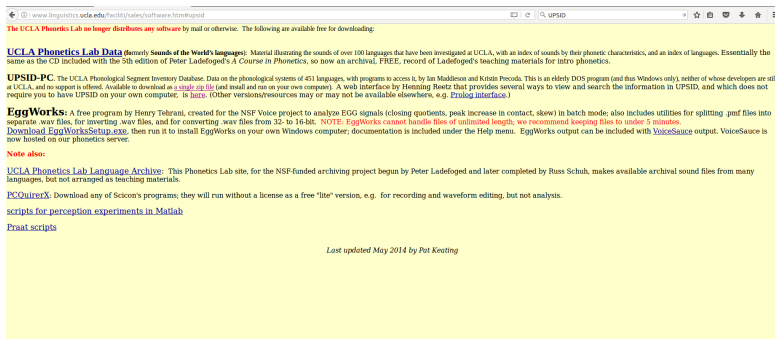| G_CODE | S_SVO | S_SOV | S_VSO | S_VOS | S_OVS | S_OSV | S_SUBJECT_BEFORE_VERB | S_SUBJECT_AFTER_VERB | S_OBJECT_AFTER_VERB | S_OBJECT_BEFORE_VERB | S_SUBJECT_BEFOR |
|---|---|---|---|---|---|---|---|---|---|---|---|
| aaa | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aab | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aac | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aad | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aae | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aaf | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aag | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aah | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aai | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aak | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aal | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aan | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aao | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aap | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aaq | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aar | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | |
| aas | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aat | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aau | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | |
| aaw | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aax | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aaz | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| aba | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abb | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | |
| abc | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abd | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abe | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abf | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abg | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abh | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abi | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abj | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abk | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | |
| abl | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abm | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abn | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abo | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abp | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abq | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abr | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abs | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abt | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | |
| abu | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abv | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abw | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |
| abx | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- | -- |

# UPSID (UCLA PHONOLOGICAL SEGMENT INVENTORY DATABASE)

### What is it?

- "This Database was compiled by Ian Maddieson and Kristin Precoda (cf. Maddieson, 1984) and contains information on the distribution of 919 different segments in 451 languages"
- **License**: Open Access

http://www.linguistics.ucla.edu/faciliti/sales/software.html

# EXAMPLES



http://www.linguistics.ucla.edu/faciliti/sales/software.html

# EXAMPLES

**Segment frequency:**
This is the number of languages that contains a specific segment divided by the number of languages in UPSID expressed in percent. For ex = 0.22 (or, in other words, it only exists in 0.2% of all languages in UPSID). The most frequent segment in UPSID is the bilabial nasal /m/, w different segments in the database and the complete list of all frequencies is rather long. The 20 most frequent consonants and the 10 most

| consonant: | m | k | j | p | w | b | h | g | N | ? | n | s | tS | S | t | f | l | "n | "t | nj |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| in languages: | 425 | 403 | 378 | 375 | 332 | 287 | 279 | 253 | 237 | 216 | 202 | 196 | 188 | 187 | 181 | 180 | 174 | 160 | 152 | 141 |
| frequency: | 94.2 | 89.4 | 83.8 | 83.2 | 73.6 | 63.6 | 61.9 | 56.1 | 52.6 | 47.9 | 44.8 | 43.5 | 41.7 | 41.5 | 40.1 | 39.9 | 38.6 | 35.5 | 33.7 | 31.3 |

| vowel: | i | a | u | E | "o | "e | O | o | e | a~ |
|---|---|---|---|---|---|---|---|---|---|---|
| in languages: | 393 | 392 | 369 | 186 | 181 | 169 | 162 | 131 | 124 | 83 |
| frequency: | 87.1 | 86.9 | 81.8 | 41.2 | 40.1 | 37.5 | 35.9 | 29.0 | 27.5 | 18.4 |

At the other end of the scale there are many segments that occur in one or only few languages:

| Number of segments: | 427 | 117 | 66 | 39 | 27 | 19 | 14 | 14 | 12 | 13 | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| that occur only in | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | languages |
| % of all segments: | 46.46 | 12.73 | 7.18 | 4.24 | 2.94 | 2.07 | 1.52 | 1.52 | 1.31 | 1.41 | |
| cummulative %: | 46.46 | 59.19 | 66.38 | 70.62 | 73.56 | 75.63 | 77.15 | 78.67 | 79.98 | 81.39 | |

That is, the group of sounds that appear in 10 or fewer of the 451 languages make up more than 80% of the 919 sounds in the database.

Interface by Henning Reetz
(http://web.phonetik.uni-frankfurt.de/upsid_info.html)

# EXAMPLES

| | |
|---|---|
| Language name: | PIRAHA |
| UPSID number: | 6802 |
| Alternate name(s): | MURA-PIRAHA~, MURA |
| Classification: | South American, Paezan |
| This language has | 11 segments |
| Its Frequency index is | 0.618020560 (average percentage of segments |
| The language has these sounds: | p b "t k g ? "s h i a "o |
| Comment: | |
| Source(s): | Everett, D.L. 1982. Phonetic rarities in Piraha. |
| | Rodrigues, A.D. 1980. Contribuicoes das lingua |
| | Sheldon, S.N. 1974. Some morphophonemic an |

Report a bug

Interface by Henning Reetz
(http://web.phonetik.uni-frankfurt.de/upsid_info.html)

# PHOIBLE

### What is it?

- "PHOIBLE Online is a repository of cross-linguistic phonological inventory data, which have been extracted from source documents and tertiary databases and compiled into a single searchable convenience sample. The 2014 edition includes 2155 inventories that contain 2160 segment types found in 1672 distinct languages."
- **License**: Open Access

[Moran, Steven & McCloy, Daniel & Wright, Richard (eds.) 2014. PHOIBLE Online. Leipzig: Max Planck Institute for Evolutionary Anthropology. (Available online at http://phoible.org, Accessed on 2017-04-21.)]

# EXAMPLES

# EXAMPLES

# EXAMPLES

['fɔɪ.bɬ]    Home    Inventories    Languages    Segments    Sources

## Segments

Showing 1 to 100 of 2,160 entries

| Name | Representation | Description |
|---|---|---|
| ʦʷ | 1/2155 (0%) | LATIN SMALL LETTER T - COMBINING DOUBLE VERTICAL LINE BELOW - LATIN SMALL LETTER S - MODIFIE |
| x̧ | 1/2155 (0%) | LATIN SMALL LETTER X - COMBINING DOUBLE VERTICAL LINE BELOW |
| k̪ʷ | 1/2155 (0%) | LATIN SMALL LETTER K - COMBINING DOUBLE VERTICAL LINE BELOW - MODIFIER LETTER SMALL W |
| s̪ | 1/2155 (0%) | LATIN SMALL LETTER S - COMBINING DOUBLE VERTICAL LINE BELOW |
| ʃ | 1/2155 (0%) | LATIN SMALL LETTER ESH - COMBINING DOUBLE VERTICAL LINE BELOW |
| χ | 1/2155 (0%) | GREEK SMALL LETTER CHI - COMBINING DOUBLE VERTICAL LINE BELOW |
| ʦ | 1/2155 (0%) | LATIN SMALL LETTER T - COMBINING DOUBLE VERTICAL LINE BELOW - LATIN SMALL LETTER S |
| χʷ | 1/2155 (0%) | LATIN SMALL LETTER X - COMBINING DOUBLE VERTICAL LINE BELOW - MODIFIER LETTER SMALL W |
| ʦʷʼ | 1/2155 (0%) | LATIN SMALL LETTER T - LATIN SMALL LETTER S - MODIFIER LETTER SMALL W - MODIFIER LETTER APOS |
| ʧʷ | 1/2155 (0%) | LATIN SMALL LETTER T - COMBINING MINUS SIGN BELOW - COMBINING DOUBLE VERTICAL LINE BELOW - |
| q | 1/2155 (0%) | LATIN SMALL LETTER Q - COMBINING DOUBLE VERTICAL LINE BELOW |
| qʷ | 1/2155 (0%) | LATIN SMALL LETTER Q - COMBINING DOUBLE VERTICAL LINE BELOW - MODIFIER LETTER SMALL W |
| χʷ | 1/2155 (0%) | GREEK SMALL LETTER CHI - COMBINING DOUBLE VERTICAL LINE BELOW - MODIFIER LETTER SMALL W |

# EXAMPLES

# ASJP

### What is it?

- "The database of the **Automated Similarity Judgment Program (ASJP)** aims to contain 40-item word lists of all the world's languages."
- **License**: Open Access

[Wichmann, Brown, Holman, et al. (eds.), 2013]

# ASJP

What's in it?

- Languages: **4664** (ca. 62% of the world's languages)
- Word lists: **7221** (either 40 or 100 lexical items)

[Wichmann, Sren, Eric W. Holman, and Cecil H. Brown (eds.). (2016). The ASJP Database (version 17).

PHYLOGENETIC
○○○○○○○○○○○○

STRUCTURE
○○○○○○○○○○○○○○○○

PHONETICS/PHONOLOGY
○○○○○○○○○

LEXICON
○○●○○○○○○

OTHERS
○○○○○○

SAMPLING BIASES

# EXAMPLES

PHYLOGENETIC ○○○○○○○○○○○○○
STRUCTURE ○○○○○○○○○○○○○○○○
PHONETICS/PHONOLOGY ○○○○○○○○○
LEXICON ○○○○●○○○○
OTHERS ○○○○○○
SAMPLING BIASES

# EXAMPLES

## Meanings

| ASJP | Home | Wordlists | Meanings | Sources |

### Meanings

Counterparts of the 40 boldfaced meanings can be obtained here for all doculects in the database, to the extent that they are attested. For a few hundred

**I** **you** **we** this that who what not all many **one** **two** big long small woman man **person** **fish** bird **dog** **louse** **tree** seed **tooth** **tongue** claw foot **knee** **hand** belly neck **breast** heart **liver** **drink** eat bite **see** **hear** know sleep **die** kill swim fly w burn **path** **mountain** red green yellow white black **night** hot cold **full** **new** good round dry **name**

# EXAMPLES

## WOLD

### What is it?

- "The **World Loanword Database (WOLD)** provides vocabularies (mini-dictionaries of about 1000-2000 entries) of 41 languages from around the world, with comprehensive information about the loanword status of each word. It allows users to find loanwords, source words and donor languages in each of the 41 languages, but also makes it easy to compare loanwords across languages. "
- **License**: Open Access

[Haspelmath, Martin & Tadmor, Uri (eds.) 2009. World Loanword Database. Leipzig: Max Planck Institute for Evolutionary Anthropology. (Available online at http://wold.clld.org, Accessed on 2017-04-21.) ]

PHYLOGENETIC
○○○○○○○○○○○○○

STRUCTURE
○○○○○○○○○○○○○○○

PHONETICS/PHONOLOGY
○○○○○○○○○

LEXICON
○○○○○○●○○

OTHERS
○○○○○○

SAMPLING BIASES

# EXAMPLES

# EXAMPLES

# D-PLACE
## What is it?

- "D-PLACE, which stands for Database of Places, Language, Culture, and Environment, represents an attempt to bring together this dispersed corpus of information [on cultural and climatic traits associated with societies]. It aims to make it easy for individuals to contrast their own cultural practices with those of other societies, and to consider the factors that may underlie cultural similarities and differences. "
- **License**: Open Access

[Kirby, K.R., Gray, R. D., Greenhill, S. J., Jordan, F. M., Gomes-Ng, S., Bibiko, H-J, et al. (2016). D-PLACE: A Global Database of Cultural, Linguistic and Environmental Diversity. PLoS ONE, 11(7): e0158391. doi:10.1371/journal.pone.0158391.]

# D-PLACE

### What's in it?

- ▶ So far, D-PLACE contains cultural, linguistic, environmental and geographic information for over 1400 human societies. A society in D-PLACE represents a group of people in a particular locality, who often share a language and cultural identity.

# EXAMPLES



**D-PLACE**

Database of Places, Language, Culture and Environment

Max Planck Institute
for the Science of Human History
(MPI SHH)

NESCent
National Evolutionary Synthesis Center

National Evolutionary
Synthesis Center
(NESCent)

# EXAMPLES

**D-PLACE**

SEARCH    RESULTS (1291)

**Societies (1291)** — Return to the search page to refine or add to your results

**RESULT**
Select a viewing mode
for your search results:

Table    Map    Tree    Download

| Name | Dataset | Glottolog Code | Language | Subsistence economy: gathering |
|------|---------|----------------|----------|--------------------------------|
| !Kung | Ethnographic Atlas | juho1239 | Ju'hoan | 76 to 85 percent dependence (Sources) |
| /Xam | Ethnographic Atlas | xamm1241 | Kham | 46 to 55 percent dependence (Sources) |
| Ababda | Ethnographic Atlas | abab1239 | Ababda | Zero to 5 percent dependence (Sources) |
| Abarambo | Ethnographic Atlas | bara1361 | Barambu | Zero to 5 percent dependence (Sources) |
| Abelam | Ethnographic Atlas | ambu1247 | Ambulas | 16 to 25 percent dependence (Source) |
| Abenaki | Ethnographic Atlas | peno1243 | Penobscot | 6 to 15 percent dependence (Sources) |

# EXAMPLES

# EXAMPLES



D-PLACE

Note: trees have been pruned to display only societies present in D-PLACE.



C2: Subsistence economy: agriculture
Dependence on agriculture, relative to other subsistence activities.

- ○ 0-5%
- ● 6-15%
- ● 16-25%
- ● 26-35%
- ● 36-45%
- ● 46-55%
- ● 56-65%
- ● 66-75%
- ● 76-85%
- ● 86-100%
- ○ Missing data

C1 C2
- Latvians
- Lithuanians
- Byelorussians
- Russians
- Ukranians
- Hutsuls
- Bulgarians
- Serbs
- Czechs
- Bihari
- Uttar Pradesh
- Kashmiri
- Punjabi
- Sindhi
- Shina
- Kohistani
- Indo-Iranian
- Bhil
- Baiga
- Gujarati
- Sinhalese
- Vedda
- Bengali
- Chakma
- Hill Bhuiya
- Pahari
- Ossetians
- Ghilzai
- Yusufzai
- Marri
- Kurd
- Basseri
- Iranians
- Hazara
- Bakhtiari
- Nuristani
- Gheg
- Boers
- Dutch
- New Englanders
- Saramaccan
- Ndyuka
- Icelanders
- Armenians
- Ancient Romans
- Spaniards
- Brazilians
- Portuguese
- Walloons
- Haitians
- French Canadians
- Neapolitans
- Romanians
- Moldovans
- Greeks
- Irish

## THE PROBLEM OF SAMPLING (VELUPILLAI, 2012)

At **two levels of analysis** we have to deal with the problem of finding **representative samples**:

- **Within a language**:
  - **corpora**: should be **balanced**, i.e. represent a wide range of registers and styles
  - **experiments**: participants should **represent the population**, e.g. age range, educational background, gender, etc.
- **Across different languages**:
  - it is hard (impossible) to assess a feature across all 7000+ languages, hence we need a **balanced** and **unbiased** sample of them

The **second level** is most relevant for Typology, though note that the first level decides how we represent languages.

## SAMPLING: TYPES OF BIASES

- **(Phylo-)genetic bias**
  languages are non-independent if they share common
  ancestry, i.e. a proto-language. A bias can result from
  over-representation of a specific **family or genus** (e.g.
  Indo-European, Germanic)



Bouckaert et al. (2012)

# SAMPLING: TYPES OF BIASES

- **Areal bias**
  language contact leads to the spread of linguistic features.
  A bias can result when languages from a specific
  **geographic region** (e.g. Balkan Sprachbund) are
  over-represented

## SAMPLING: TYPES OF BIASES

- **Typological bias**
  over-representation of specific **typological features** (e.g.
  languages with/without tone)

- **Cultural bias**
  certain cultures might put more emphasis on encoding
  specific kinds of information in their grammar (e.g. Korean
  and Japanese honorific markers)

- **Bibliographical bias** bias towards **well-described
  languages**. Not all languages are equally well described,
  and for many languages any kind of information is lacking
  all together (e.g. Sentinelese).

## SAMPLING: TYPES OF SAMPLES

- **Variety sample**
  a sample covering **all the parameters of a linguistic variable under investigation**. For example, if we are dealing with word orders we would need to have a sample that covers all the logical possibilities SOV, SVO, OVS, etc. (if they all exist)

- **Convenience sample**
  Based on the **descriptions available** for a typological variable (relation to the biographical bias). Note that still an attempt can be made to balance the sample areally, genetically, etc.

## SAMPLING: TYPES OF SAMPLES

▶ **Probability sample**
A sample that does not have any of the biases named
above, and hence represents only **fully independent
languages**. Strictly speaking, only with such a sample is it
possible to make valid **statistical judgements** about the
probability of occurrence or co-occurrence of certain
typological variables (e.g. how probable is it for a
language to have tone?)

▶ **Random sample**
Drawing languages **randomly without any sensitivity to
biases**. Note that for small samples, a random sample can
be highly biased. However, for bigger samples bias is
increasingly unlikely.

## SAMPLING: PROBLEMS

- **Variety sample:** only works if we know the number of logical possibilities per variable (e.g. for word order SOV, SVO, VSO, OSV, OVS, VOS). Does not work for variables that are open-ended (e.g. number of phonemes, case markers, etc.)

- **Convenience sample:** very likely to be biased in one way or another

- **Probability sample:** very hard (or impossible) to get. For example, Piantadosi & Gibson (2011) argue that to rectify a universal of the type *all languages have feature x*, we would need a sample of **500 independent** languages

- **Random sample:** small samples will be biased, for big samples there is the same problem as for probability samples