# *Some Computational Tools for Language Typology*

## Dik Bakker
## Lancaster University

# Tools for Typology

**<u>Period 1990 – 2009:</u>**

**Computer programs <-> Typological projects**

# Tools for Typology

**Period 1990 – 2009:**

**Computer programs <-> Typological projects**

**a. Language sampling**

# Tools for Typology

**Period 1990 – 2009:**

**Computer programs <-> Typological projects**

**a. Language sampling**
**b. Inference of universal implications**

# Tools for Typology

**<u>Period 1990 – 2009:</u>**

**Computer programs <-> Typological projects**

**a. Language sampling**
**b. Inference of universals**
**c. <span style="color:red">Lexical classification of languages</span>**

# Tools for Typology

**Period 1990 – 2009:**

**Computer programs <-> Typological projects**

**a. Language sampling**
**b. Inference of universals**
**c. Lexical classification of languages**
**d. <span style="color:red">Language contact and borrowing</span>**

# Tools for Typology

**<u>Period 1990 – 2009:</u>**

**Computer programs <-> Typological projects**

**a. Language sampling**
**b. Inference of universals**
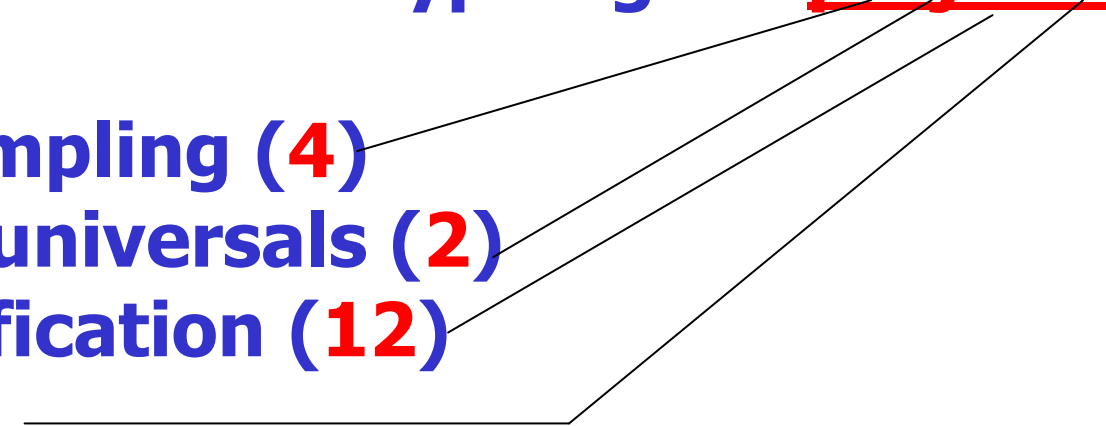**c. Lexical classification of languages**
**d. Borrowing**

**...**

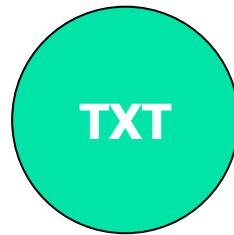# Tools for Typology

**Period 1990 – 2009:**

**Computer programs <-> Typological <span style="color:red">projects</span>**

**a. Language sampling (4)**
**b. Inference of universals (2)**
**c. Lexical classification (12)**
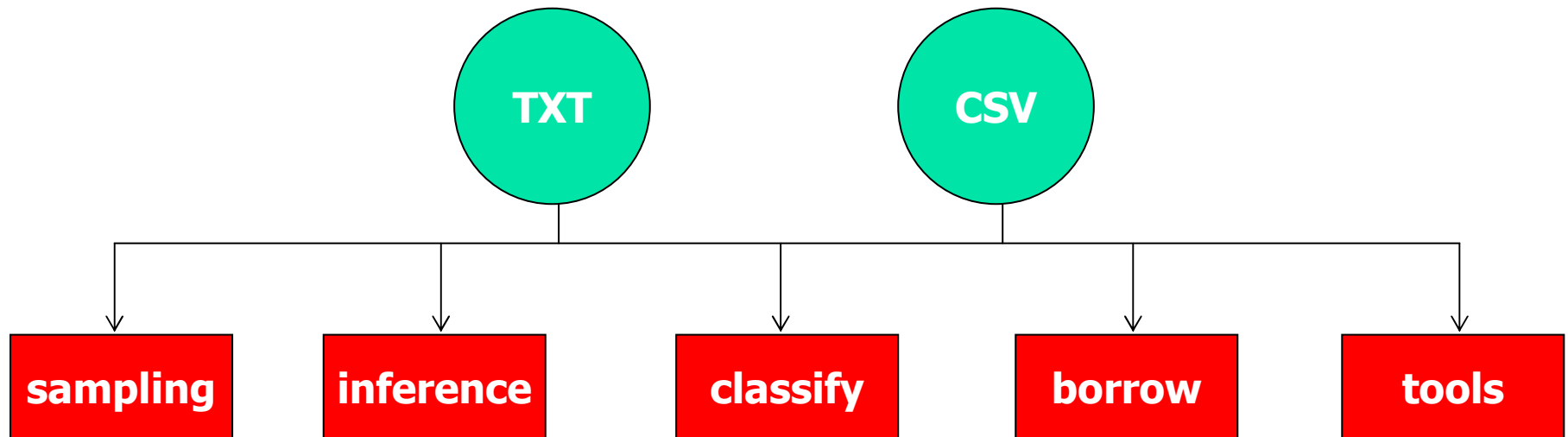**d. Borowing (3)**

# Tools for Typology

**TXT**

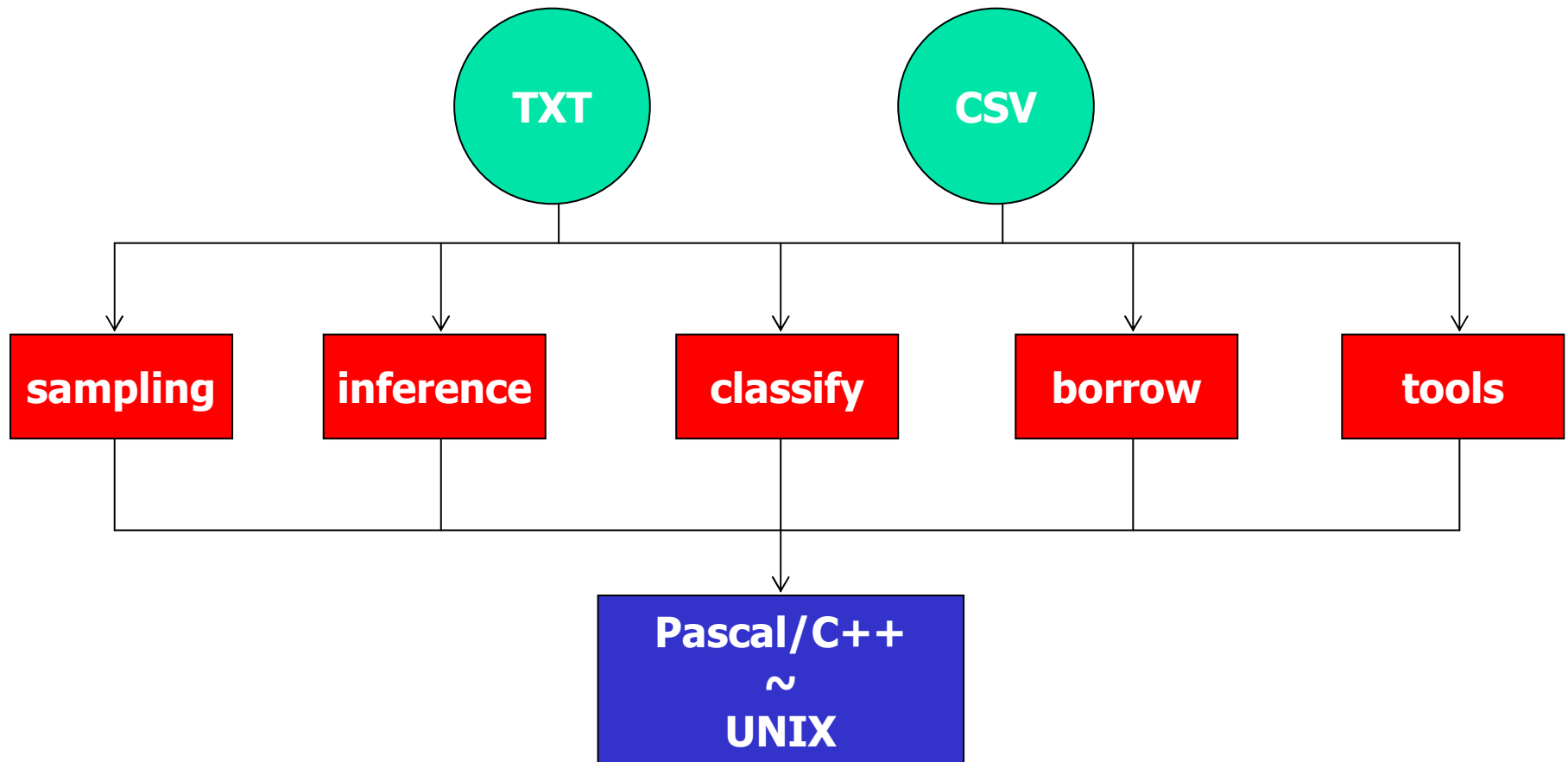# Tools for Typology

**TXT**

**CSV**

# Tools for Typology

# Tools for Typology

# Tools for Typology

**<u>Points today:</u>**

# Tools for Typology

**Points today:**

**1. Give an impression of local software ( ∞ )**

# Tools for Typology

**Points today:**

**1. Give an impression of local software ( ∞ )**

**2. How to make it accessible?**

# Tools for Typology

**Overview:**

# Tools for Typology

**Overview:**
**1. Sampling**

# Tools for Typology

**Overview:**
**1. Sampling**
**2. Inference of universals**

# Tools for Typology

**Overview:**
**1. Sampling**
**2. Inference of universals**
**3. Lexical classification**

# Tools for Typology

**<u>Overview:</u>**
**1. Sampling**
**2. Inference of universals**
**3. Lexical classification**
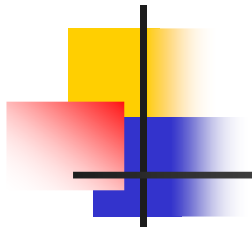**4. Nothing about Borrowing:**

Bakker, D., J. Gómez-Rendón & E. Hekking (2008).
'Spanish meets Guaraní, Otomí and Quichua: a multilingual
  confrontation'.
In Th. Stolz, D. Bakker & R. Palomo (eds)
*Aspects of Language Contact*. Mouton de Gruyter, 165-238.

# Tools for Typology

**Overview:**
**1. Sampling**
**2. Inference of universals**
**3. Lexical classification**
**4. Accessibility**

# 1. Language Sampling

**Together with:**
**Kees Hengeveld (Amsterdam)**
**Peter Kahrel (Amsterdam)**
**Jan Rijkhoff (Aarhus)**

**Reference:**
**Rijkhoff J. & D. Bakker (1998).**
**'Language sampling'.**
*Linguistic Typology* **2-3, 263-314.**

# Language sampling

**Typological project: typically 50 – 500 languages**

**Question: how to select?**

# Language sampling

**General issues:**

# Language sampling

**General issues:**

- **Many features more or less tight to genetic relationships**

# Language sampling

**General issues:**

**- Many features tight to genetic relationships**

**- Areal and contact phenomena**

# Language sampling

**General issues:**

- **Many features tight to genetic relationships**

- **Areal and contact phenomena**

- **Distribution of some linguistic features and relations between them are well-known, of (most) others not at all**

# Language sampling

**General issues:**

- **Many features tight to genetic relationships**

- **Areal and contact phenomena**

- **Only some distributions and relations well-known**

- **Bibliographic gaps**

# Language sampling

**Three types of samples:**

# Language sampling

**Three types of samples:**

**1. Random sample**

**→ Only when each language same chance**

# Language sampling

**Three types of samples:**

**1. Random sample**

**2. Probability sample**

→ **Measures chance on occurrence of certain feature value, or of language type**

# Language sampling

**Three types of samples:**

**1. Random sample**

**2. Probability sample**

**→ Measure chance certain feature value/type**

**Genetic and areal bias:**
**independency ~ (in)stability**

# Language sampling

**Three types of samples:**

**1. Random sample**

**2. Probability sample**

**3. Variety sample**

→ **Exploration of unknown feature/type:**
               **maximum variation**

# Language sampling

**Three types of samples:**

1. **Random sample** → **large**

2. **Probability sample** → **small**

3. **Variety sample** → **intermediate - large**

# Language sampling

**Three types of samples:**

**1. Random sample** → **large**

**2. Probability sample** → **small**

**3. Variety sample** → **intermediate - large**

# Language sampling

**Variety sample:**

**Maximize variety ~ maximize diversity factor:**

# Language sampling

**Variety sample:**

**Maximize variety ~ maximize diversity factor:**

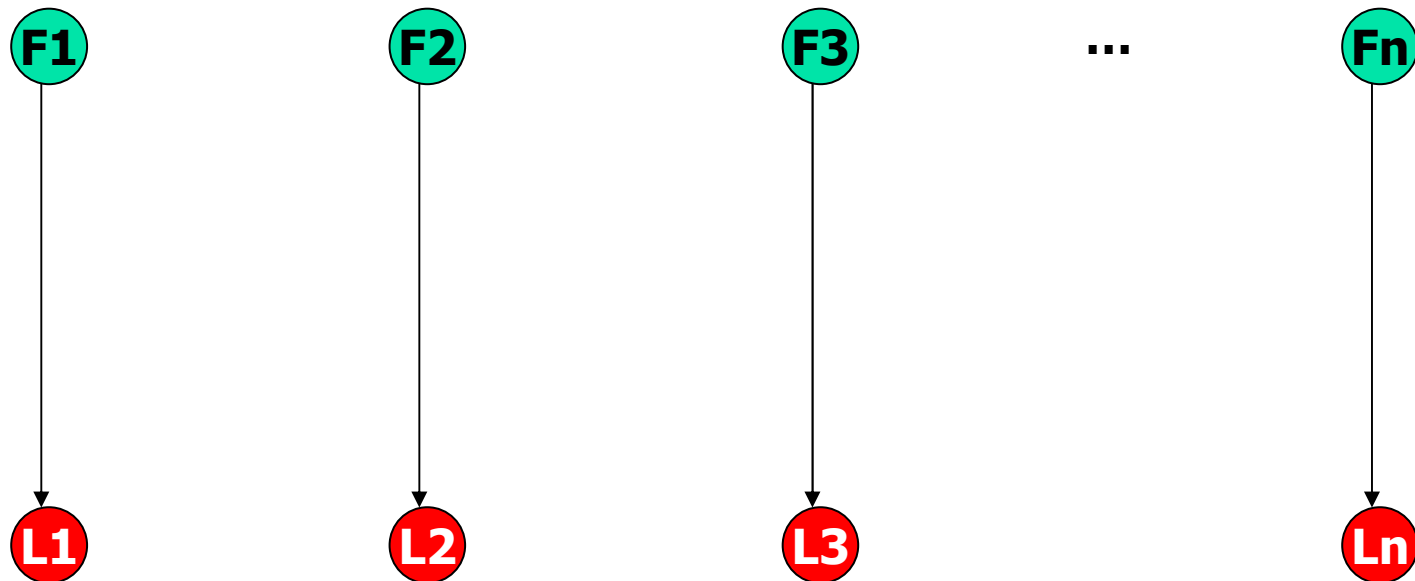- **language(s) from all families**

# Language sampling

**Variety sample:**

**Maximize variety ~ maximize diversity factor:**

- **language(s) from all families**

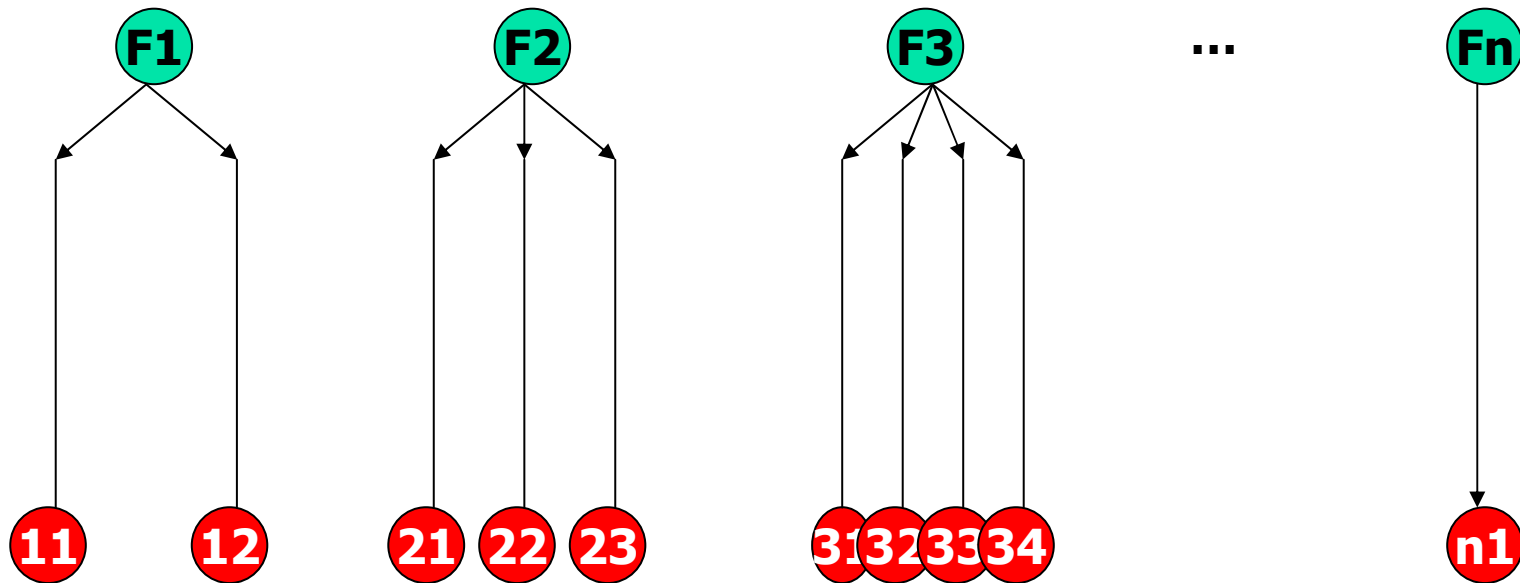- **from as many subgroupings as fit in sample size**

# DV method

**Sample size = n (minimum)**

F1          F2          F3          **...**          Fn

L1          L2          L3                           Ln

**(any language from family for which documentation available)**

# DV method

**Sample size > n**

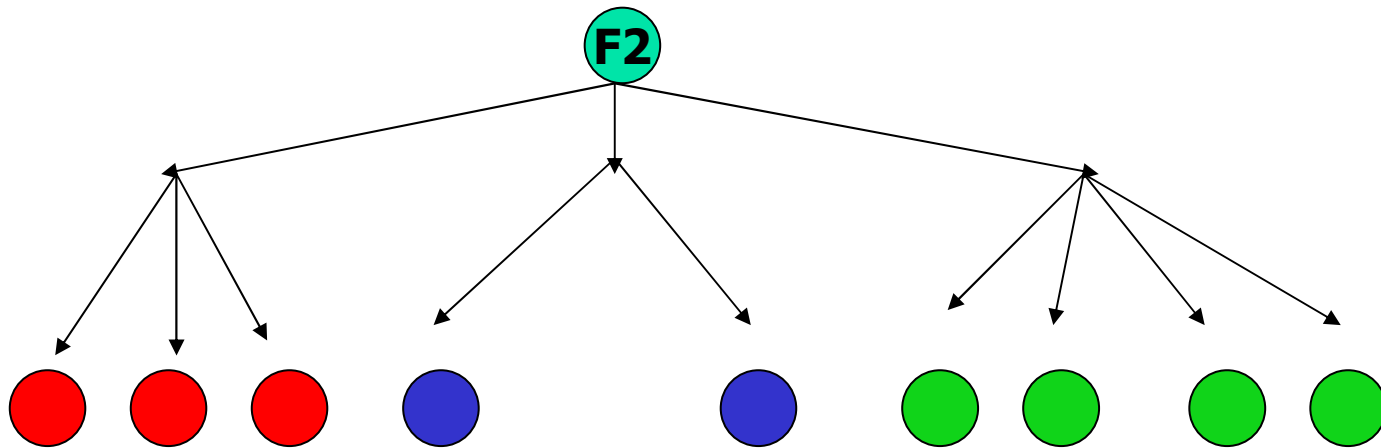F1 → 11, 12

F2 → 21, 22, 23

F3 → 31, 32, 33, 34

... Fn → n1

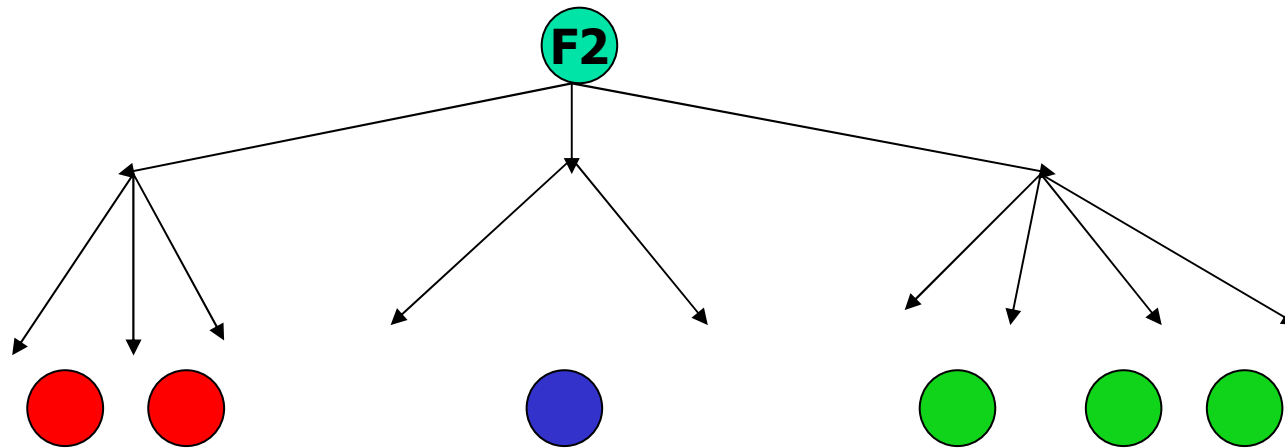**(any language from group for which documentation available)**

# DV method

**Sample size >>> n**



**(any language for which documentation available)**

# DV method

**Sample size >> n**



**(any language for which documentation available)**
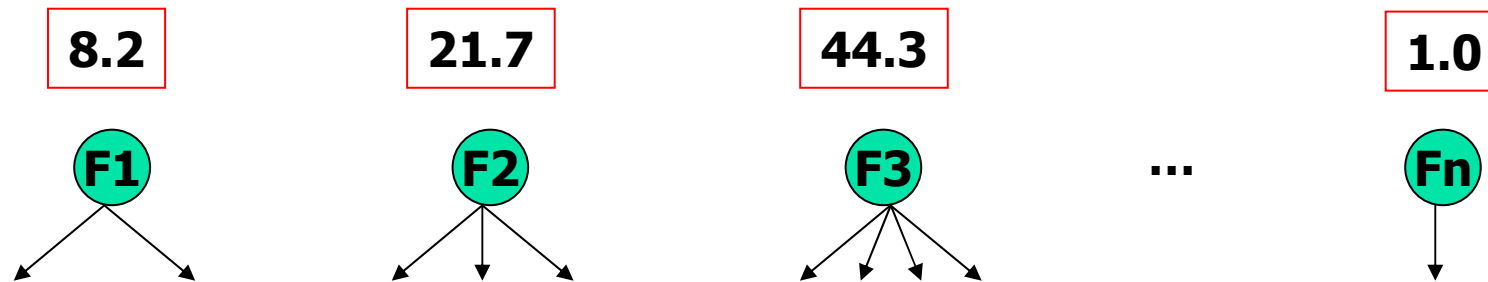
# DV method

F1   F2   F3   ...   Fn

**Diversity Value** per node, based on:
- **NOT number of daughter nodes**
- **NOT on number of daughter languages, but**
- **Internal Complexity (breadth per level, diminishing)**

# DV method

| 8.2 | 21.7 | 44.3 | | 1.0 |

F1      F2      F3     ...     Fn

**Diversity Value** **per node, based on internal complexity**
**a. Per family**

# DV method



| 8.2 | 21.7 | 44.3 | 1.0 |

F1      F2      F3      ...      Fn

| 22.1 | 12.5 | 84.2 | 41.7 | 3.2 | | 1.0 |

**Diversity Value** per node, based on internal complexity
a. **Per family**
b. **Recursively per lower node**

# DV method

## Procedure:

**1. Choose language classification (Ethn/Ruh/Voeg)**

# DV method

**Procedure:**

**1. Choose language classification (Ethn/Ruh/Voeg)**

**2. Calculate DV value per node (all tree-like)**

# DV method

**Procedure:**

1. **Choose language classification**

2. **Calculate DV value per node**

3. **Establish sample size (minimum = n of families)**

# DV method

**Procedure:**

1. **Choose language classification**

2. **Calculate DV value per node**

3. **Establish sample size**

4. **Assign languages to families weighted by DV (> 0)**

# DV method

**Procedure:**

**1. Choose language classification**

**2. Calculate DV value per node**

**3. Establish sample size**

**4. Assign languages to families weighted by DV**

**5. Recursively assign languages to lower groups**

# DV method

**Procedure:**

**1. Choose language classification**

**2. Calculate DV value per node**

**3. Establish sample size**

**4. Assign languages to families weighted by DV**

**5. Recursively assign languages to lower groups**

**6. Stop when no languages left to assign**

# DV method

**Procedure:**

**1. Choose language classification**

**2. Calculate DV value per node**

**3. Establish sample size**

**4. Assign languages to families weighted by DV**

**5. Recursively assign languages to lower groups**

**6. Stop when no languages left to assign**

**7. Optional: select language names (random / criteria)**

# DV method: results

```
Classification: Ruhlen91
Criterion 1: Diversity Value
Sample size: 50 ( 0.95 % of 5273, min=30)
```

# DV method: results

```
Classification: Ruhlen91
Criterion 1: Diversity Value
Sample size: 50 ( 0.95 % of 5273, min=30)

Afro-Asiatic (55.53/6/258)        2
Altaic (15.07/2/62)               1
Korean-Japanese (3.00/3/4)        1
Australian (67.58/30/262)         3
Austric (137.41/3/1186)           5
    Austro-Tai (106.03/2/1027)         3
        Austronesian (118.17/4/970)         2
        Daic (4.67/2/57)                    1
    Austroasiatic (28.08/2/155)        1
    Miao-Yao (2.00/2/4)                1
…
```

# DV method: results

```
Classification: Ruhlen91
Criterion 1: Diversity Value
Sample size: 50 ( 0.95 % of 5273, min=30)

Afro-Asiatic (55.53/6/258)        2
Altaic (15.07/2/62)               1
Korean-Japanese (3.00/3/4)        1
Australian (67.58/30/262)         3
Austric (137.41/3/1186)           5
    Austro-Tai (106.03/2/1027)        3
        Austronesian (118.17/4/970)           2
        Daic (4.67/2/57)                      1
    Austroasiatic (28.08/2/155)       1
    Miao-Yao (2.00/2/4)               1
…
```

# DV method: results

```
Classification: Ruhlen91
Criterion 1: Diversity Value
Sample size: 50 ( 0.95 % of 5273, min=30)

Afro-Asiatic (55.53/6/258)          2
Altaic (15.07/2/62)                 1
Korean-Japanese (3.00/3/4)          1
Australian (67.58/30/262)           3
Austric (137.41/3/1186)             5
    Austro-Tai (106.03/2/1027)          3
        Austronesian (118.17/4/970)          2
        Daic (4.67/2/57)                     1
    Austroasiatic (28.08/2/155)         1
    Miao-Yao (2.00/2/4)                 1
…
```

# DV method

**Options:**

# DV method

**Options:**

## 1. Random selection of languages under nodes

# DV method

```
Classification: Ethnologue15
Criterion 1: Diversity Value
   Sample size: 150 ( 2.06 % of 7299)

   Austronesian (192.99/12/1268)  5
       Unclassified (1.00/0/1)       1
           1. Ketangalan (G)
       East Formosan (3.00/3/5)      1
           Central (1.00/0/2)          1
           2. Amis        (G)
       Bunun (1.00/0/1)              1
           3. Bunun       (G)
       Western Plains (2.00/2/2)     1
           Thao (1.00/0/1)             1
           4. Thao        (G)
…
```

# DV method

**Options:**

**1. Random selection of languages under nodes**

**2. Stratification on basis of feature values**

# DV method

**Options:**

**1. Random selection of languages under nodes**

**2. Stratification on basis of feature values**

**→ Problem: bibliographic bias**

# DV method

**Options:**

**1. Random selection of languages under nodes**

**2. Stratification on basis of feature values**
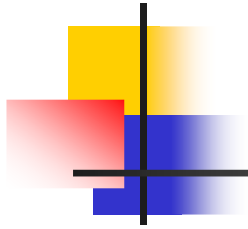
**3. Evaluate existing samples**

# DV method

**Options:**

**1. Random selection of languages under nodes**

**2. Stratification on basis of feature values**

**3. Evaluate existing samples**

**Program has been used for a large number of studies (MA, PhD, articles, books)**

# 2. Inference of Universals

**Together with:**
**Anna Siewierska (Lancaster)**

**Reference:**
**Bakker, D. (2008).**
**'LINFER: inferring implications from the WALS database'.**
*STUF* **61-3, 186-198.**

# UNIVERSALS

Greenberg (1963):

# UNIVERSALS

Greenberg (1963):

Absolute: Universal 3

Languages with dominant VSO order are *always* prepositional.

# UNIVERSALS

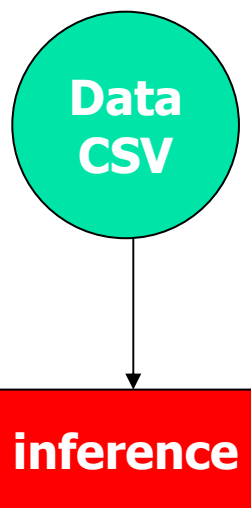Greenberg (1963):

Absolute: Universal 3

Languages with dominant VSO order are *always* prepositional.

Statistical: Universal 4

*With overwhelmingly greater than chance frequency*, languages with normal SOV order are postpositional.
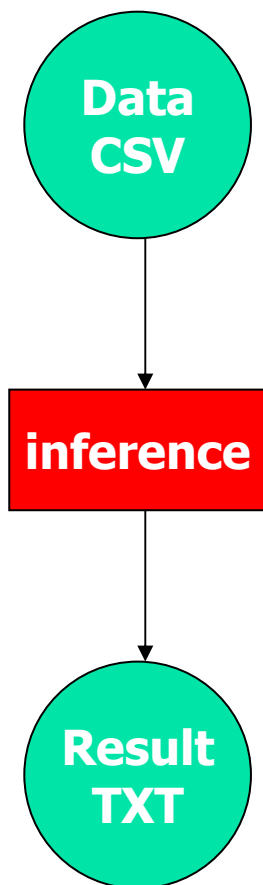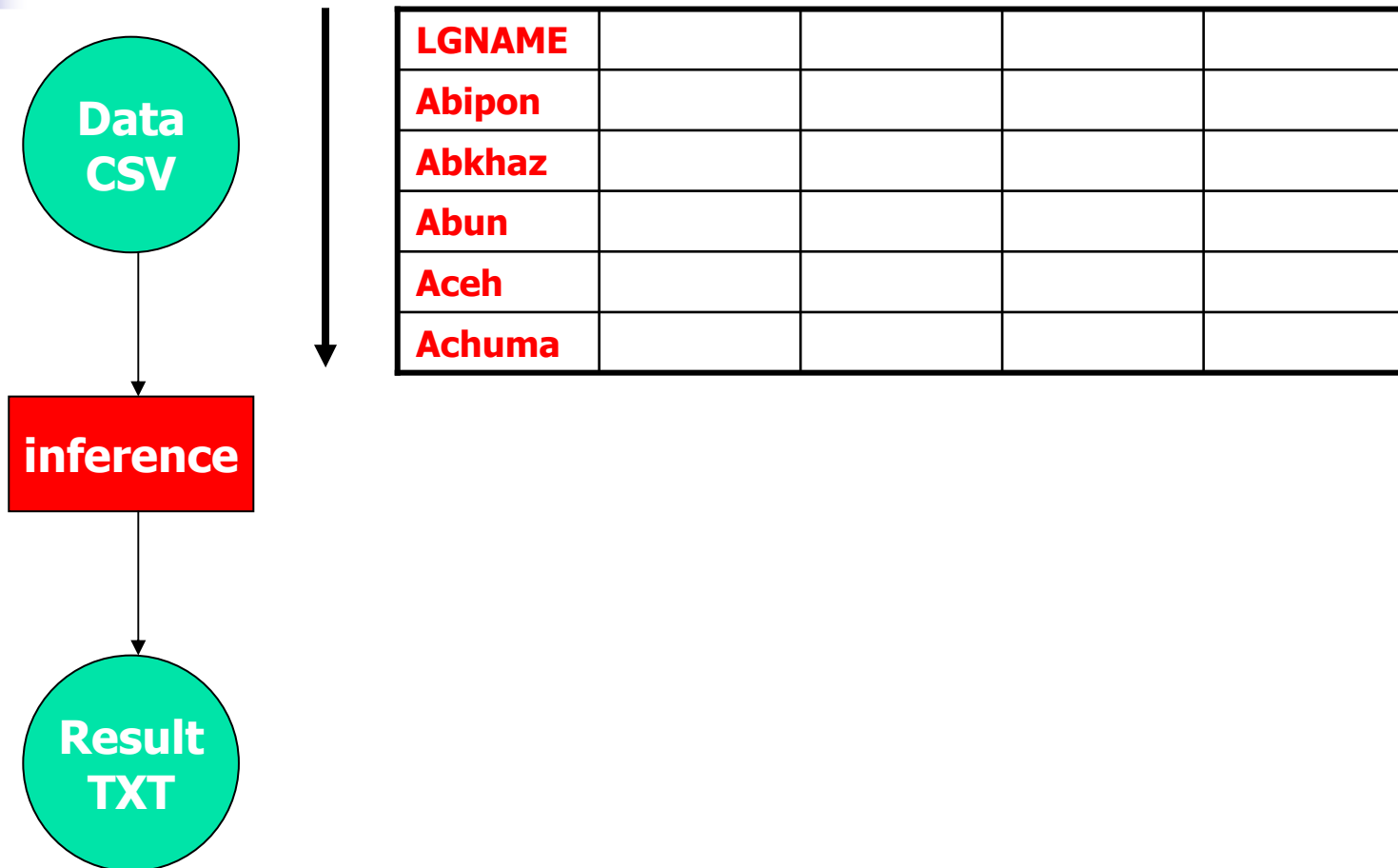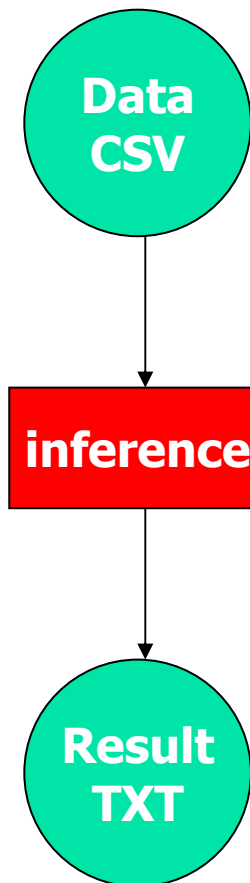
# LINFER

**Data CSV**

**inference**

# LINFER

**Data CSV**

↓

**inference**

↓

**Result TXT**

# LINFER

**Data CSV**

**inference**

**Result TXT**

| LGNAME | | | | |
|--------|---|---|---|---|
| Abipon | | | | |
| Abkhaz | | | | |
| Abun | | | | |
| Aceh | | | | |
| Achuma | | | | |

# LINFER

**Data CSV** → **inference** → **Result TXT**

| LGNAME | SmrkP | SmrkV | SmrkN | SmrkH |
|--------|-------|-------|-------|-------|
| Abipon | 123 | No | Sgpl | No |
| Abkhaz | 123 | No | Sgpl | No |
| Abun | 12 | No | Sg | No |
| Aceh | 12 | Yes | Nonum | Irr |
| Achuma | 123 | No | Sgdupl | Yes |

# LINFER

**Data CSV**

**inference**

**Result TXT**

| LGNAME | SmrkP ➤ | SmrkV | SmrkN | SmrkH |
|--------|---------|-------|-------|-------|
| Abipon | **123** | **No** | Sgpl | No |
| Abkhaz | **123** | **No** | Sgpl | No |
| Abun | 12 | No | Sg | No |
| Aceh | 12 | Yes | Nonum | Irr |
| Achuma | **123** | **No** | Sgdupl | Yes |

**SmrkP = 123  →  SmrkV = No (ABS)**

Tools for Typology

# LINFER

**Data CSV**

**inference**

**Result TXT**

| LGNAME | SmrkP ← | SmrkV | SmrkN | SmrkH |
|--------|---------|-------|-------|-------|
| Abipon | **123** | **No** | Sgpl | No |
| Abkhaz | **123** | **No** | Sgpl | No |
| Abun | **12** | **No** | Sg | No |
| Aceh | 12 | Yes | Nonum | Irr |
| Achuma | **123** | **No** | Sgdupl | Yes |

SmrkP = 123 → SmrkV = No (ABS)

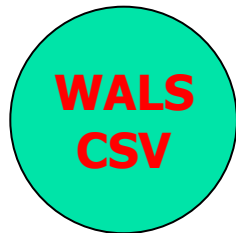**SmrkV = No → SmrkP = 123 (0.75)**

# LINFER

- Automatic inference of implications (A & S):

→ **generate + test**

# LINFER: on WALS

**WALS CSV**

↓

**inference**

Haspelmath, M., M. Dryer, D. Gil & B. Comrie (eds) (2005).
*The World Atlas Of Language Structures*.
Oxford: Oxford University Press
WALS Online:  http://wals.info/

Number of languages:        2558
Number of variables:         143

# LINFER

A first run:

All languages                     2558
All variables                     139 (minus SignLgs)

# LINFER

First run:

All languages                      2558
All variables                      139

Results:

Potential implications:            413,886
Accepted implications:             1,385 ( = 0.33%)

# LINFER

```
1. CORSEX=1 <=> CORNUM=1      n=144
   0.56 [3] - 0.56 [5]
```

```
[Fr=0.563, Fa=1.000, Fc=1.000, Fn=1.000, chi2<0.5%] EQUIV
```

```
Sex-based and Non-sex-based Gender System: No gender system
        <=>
Number of Genders: None
```

# LINFER

1. **COR**SEX=1   **OR**NUM=1     **n=144**
   0.56 [3...]

[Fr=0.563, Fa=1.00...   Fn=1.000, chi2<0.5%] **EQUIV**

Sex-based and Non-sex-based   System: **No gender system**
      **<=>**
Number of Genders: **None**

**TRIVIAL!!!**

# LINFER

15. **KAY**BCC=4 **<=> VES**TAM=2   n=   2
    0.13 [7] - 0.13 [4]

[Fr=0.133, Fa=1.000, Fc=1.000, Fn=1.000, chi2<0.5%] **EQUIV**

Number of Basic Colour Categories:
  **7 or between 7 and 8 categories**
      <=>
Suppletion According to Tense and Aspect:
  **Suppletion according to aspect**

# LINFER

15. **KAY**BCC=~~YES~~TAM=2          ( n=  2 )
    0.13 [~~~~]

[Fr=0.133, Fa=1.00~~~~n=1.000, chi2<0.5%] **EQUIV**


Number of Basic Colour Categ~~~~
  **7 or between 7 and 8 categori~~~~**
      **<=>**
Suppletion According to Tense and Aspec~~~~
  **Suppletion according to aspect**

**INSIGNIFICANT!!!**

# LINFER

57. **DRY**RPO=4 => **DRY**REL=1    **n=291**
      0.49 [5] – 0.73 [7]

[Fr=0.486, Fa=0.983, Fc=0.655, Fn=0.511, chi2<0.5%] **STAT**

Relationship between the Order of Object:
  **Verb-object and prepositional (VO&Prep)**
        =>
Order of Relative Clause and Noun:
  **Relative clause follows noun (NRel)**

# LINFER

57. **DRY**RPO=4 => **DRY**REL=1    ( **n=291** )
      0.49 [5] – 0.73 [7]

[Fr=0.486, Fa=0.983, Fc=0.655, Fn=0.511, chi2<0.5%] **STAT**


Relationship between the Order of Object:
  **Verb-object and prepositional (VO&Prep)**
      =>
Order of Relative Clause and Noun:
  **Relative clause follows noun (NRel)**

**VO & Prep → NRel**

# LINFER

57. **DRY**RPO=4 => **DRY**REL=1     **n=291**
      0.49 [5] – 0.73 [7]

[Fr=0.486, **Fa=0.983, Fc=0.655,** Fn=0.511, chi2<0.5%] **STAT**

**Relationship between the Order of Object:**
  **Verb-object and prepositional (VO&Prep)**
      **=>**
**Order of Relative Clause and Noun:**
  **Relative clause follows noun (NRel)**

**VO & Prep → NRel**

**EXC: cnt hak mnd squ tuk**

# LINFER

- Automatic inference of implications (Abs & Stat)

# LINFER

- Automatic inference of implications (Abs & Stat)
- Ordered from 'strongest' to 'weakest'

# LINFER

- Automatic inference of implications (A & S)
- Ordered from 'strongest' to 'weakest'
- Filtering thresholds

# LINFER

- Automatic inference of implications (A & S)
- Ordered from 'strongest' to 'weakest'
- Filtering thresholds
- Selection on subsamples of languages

# LINFER

- Automatic inference of implications (A & S)
- Ordered from 'strongest' to 'weakest'
- Filtering thresholds
- Selection on subsamples of languages
- Grouping of variables and values

# LINFER

- Automatic inference of implications (A & S)
- Ordered from 'strongest' to 'weakest'
- Filtering thresholds
- Selection on subsamples of languages
- Grouping of variables and values
- Analysis of exceptions

# LINFER

**EXPLANATION COUNTEREXAMPLES:**

    9. BICEXP=5 => NICMTP=1     n= 60
        0.48 [5] - 0.86 [3]


**[Fr=0.481, Fa=0.952, Fc=0.531, Fn=0.221, chi2<0.5%] STAT**


**Exponence of Selected Inflectional Form:**
**No case**
        **=>**
**M-T Pronouns**
**No M-T pronouns**

**EXC: fre grb lkt**

# LINFER

```
     9. BICEXP=5 => NICMTP=1      n= 60
```

EXC: fre grb lkt

** Possible explaining factors: **
fre:
NICMTP=2 (M-T pronouns, paradigmatic)
HAAEVC=5 (Separate particle)
MADUVU=3 (Uvular continuants only)
grb:
NICMTP=2 (M-T pronouns, paradigmatic)
lkt:
NICMTP=2 (M-T pronouns, paradigmatic)

# LINFER

- Automatic inference of implications (A & S)
- Ordered from 'strongest' to 'weakest'
- Filtering thresholds
- Selection on subsamples of languages
- Grouping of variables and values
- Analysis of exceptions
- Chaining of implications (AND/OR)

VO & Prep → NRel

# LINFER

Two major questions:

# LINFER

Two major questions:

1. When is an implication statistically reliable?

# LINFER

Two major questions:

1. When is an implication statistically reliable?

2. When is an implication linguistically interesting?

# LINFER

**57. DRYRPO=4 => DRYREL=1   n=291**

**[Fr=0.486, Fa=0.983, Fc=0.655, Fn=0.511, chi2<0.5%] STAT**

Relevance: proportion of values for premisse ($p / \Sigma p_i$)
Applicability: proportion of counterexamples ($p \rightarrow \neg q$)
Coverage: proportion of non-premisse languages with conclusion ($\neg p \rightarrow q$)
Dominance: proportion of languages with relevant value for variables ($p / q$)
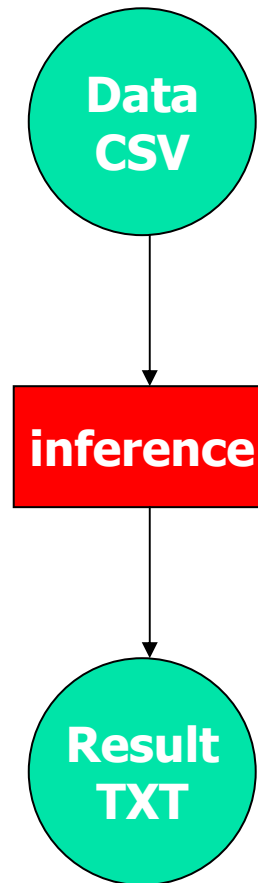Negation: proportion of languages with reverse implication ($\neg p \rightarrow \neg q$)
Chi2: for n x m tables (not tetrachoric)
Fisher Exact: when tetrachoric and 1 empty cell
Other statistics: < export data >

# LINFER

**57. DRYRPO=4 => DRYREL=1    n=291**

**[Fr=0.486, Fa=0.983, Fc=0.655, Fn=0.511, <span style="color:red">chi2<0.5%</span>] STAT**

Relevance: proportion of values for premisse ($p$ / $\Sigma p_i$)
Applicability: proportion of counterexamples ($p \rightarrow \neg q$)
Coverage: proportion of non-premisse languages with conclusion ($\neg p \rightarrow q$)
Dominance: proportion of languages with relevant value for variables ($p$ / $q$)
Negation: proportion of languages with reverse implication ($\neg p \rightarrow \neg q$)
**Chi2: for n x m tables (not tetrachoric)**
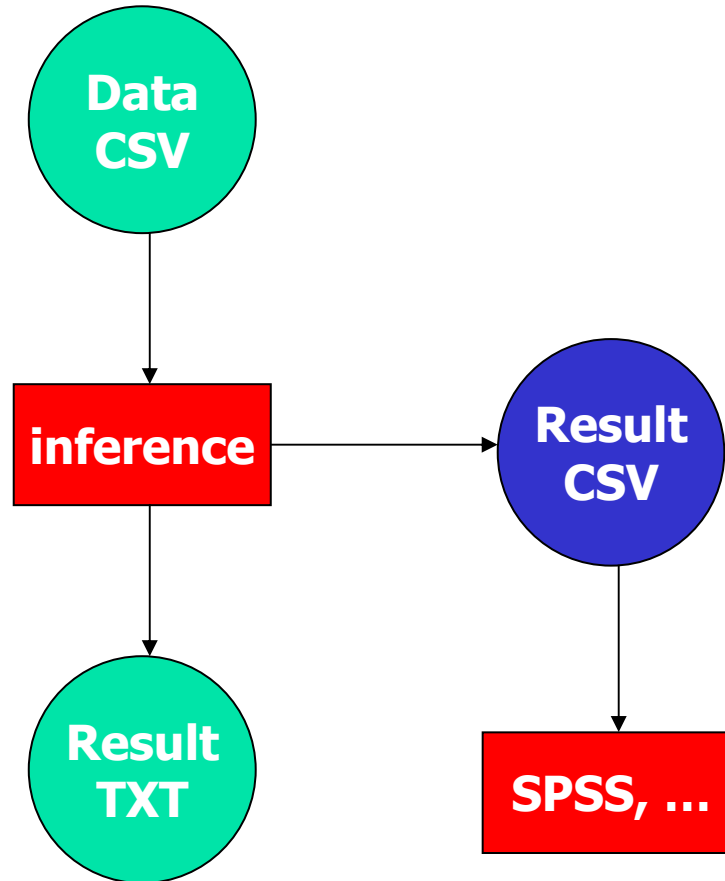**Fisher Exact: when tetrachoric and 1 empty cell**
Other statistics: <export data >

# LINFER

**57. DRYRPO=4 => DRYREL=1   n=291**

**[Fr=0.486, Fa=0.983, Fc=0.655, Fn=0.511, chi2<0.5%] STAT**

Relevance: proportion of values for premisse ($p / \Sigma p_i$)
Applicability: proportion of counterexamples ($p \rightarrow \neg q$)
Coverage: proportion of non-premisse languages with conclusion ($\neg p \rightarrow q$)
Dominance: proportion of languages with relevant value for variables ($p / q$)
Negation: proportion of languages with reverse implication ($\neg p \rightarrow \neg q$)
Chi2: for n x m tables (not tetrachoric)
Fisher Exact: when tetrachoric and 1 empty cell
**Other statistics: <export data >**

# LINFER

**Data CSV**

↓

**inference**

↓

**Result TXT**

# LINFER

Data
CSV

inference

Result
CSV

Result
TXT

SPSS, …

# LINFER: DIACHRONY?

**DIACHRONY?:**

99. **CYS**VRB=3 => **CYS**IND=3    n= 75
    0.40 [5] – 0.60 [5]

[Fr=0.395, Fa=0.949, Fc=0.625, Fn=0.628, chi2<0.5%] **STAT**

Inclusive/Exclusive Distinction in **Verba**:
  No inclusive/exclusive opposition
     =>
Inclusive/Exclusive Distinction in **Pronoun**:
  No inclusive/exclusive opposition

**EXC: abk cle map mrd**

# LINFER: RELATED?

**RELATED?:**

187. **HAS**NPL=6 => **BRO**FIN=2      n= **49**
         0.49 [6] – 0.93 [2]


[Fr=0.490, Fa=1.000, Fc=0.527, Fn=0.137, chi2<1.0%] ABS

Occurrence of Nominal Plurality:
  **Plural in all nouns, always obligatory**
         =>
Finger and Hand:
  **Different words denote 'hand' and 'finger'**

# LINFER

WALS(-like) database, observations:

# LINFER

WALS(-like) database:

- Less than 1:1000 logically possible implications are of potential interest

# LINFER

WALS(-like) database:

- Less than 1:1000 logically possible implications are of potential interest
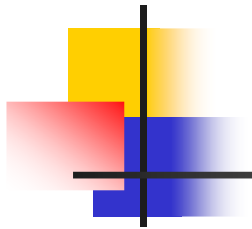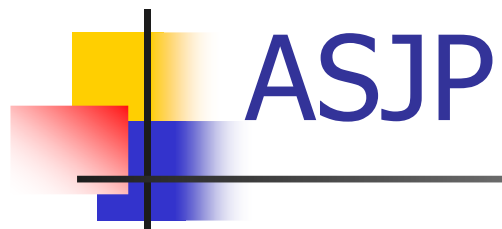- Most equivalences are trivial

# LINFER

WALS(-like) database:

- Less than 1:1000 logically possible implications are of potential interest
- Most equivalences are trivial
- Many statistically valid implications are hard to interpret linguistically

# LINFER

WALS(-like) database:

- Less than 1:1000 logically possible implications are of potential interest
- Most equivalences are trivial
- Many statistically valid implications are hard to interpret linguistically
- Need for definition: **interesting universal**

# 3. Lexical Language Classification

# ASJP

**Project ASJP (= Automated Similarity Judgment Program)**

# ASJP

**ASJP are:**     Sören Wichmann (BRD; Netherlands)
Viveka Velupillai (BRD)
André Müller (BRD)
Robert Mailhammer (BRD)
Hagen Jung (BRD)
Eric Holman (USA)
Anthony Grant (UK)
Dmitry Egorov (Russia)
Pamela Brown (USA)
Cecil Brown (USA)
Dik Bakker (UK; Netherlands)

# ASJP

**ASJP are:**    Sören Wichmann (BRD; Netherlands)
Viveka Velupillai (BRD)
André Müller (BRD)
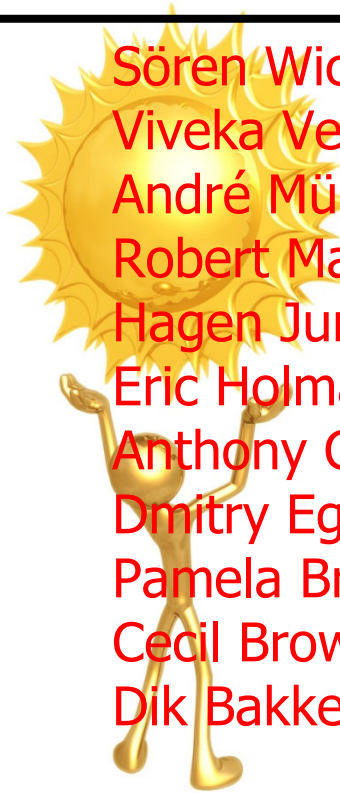Robert Mailhammer (BRD)
Hagen Jung (BRD)
Eric Holman (USA)
Anthony Grant (UK)
Dmitry Egorov (Russia)
Pamela Brown (USA)
Cecil Brown (USA)
Dik Bakker (UK; Netherlands)

# ASJP

**Reference:**
**Bakker, D., A. Müller, V. Velupillai, S. Wichmann, C. H. Brown, P. Brown, D. Egorov, R. Mailhammer, A. Grant, E. W. Holman (2009). 'Adding typology to lexicostatistics: a combined approach to language classification'.**
***Linguistic Typology* 13, 167-179.**

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**Overall goal:**
**Automatic reconstruction of language relationships
(lexical, grammatical → genetic, areal, typological, …)**

# ASJP

**Project:**
ASJP (**A**utomated **S**imilarity **J**udgment **P**rogram)

**Overall goal:**
Automatic reconstruction of language relationships

**Basis:**
Distance matrix between individual languages based
    on lexical features

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**Overall goal:**
**Automatic reconstruction of language relationships**

**Basis:**
**Distance matrix between individual languages based on lexical features**

**Method:**
**Lexicostatistics: mass comparison of *basic* lexical items,**

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**Overall goal:**
**Automatic reconstruction of language relationships**

**Basis:**
**Distance matrix between individual languages based
on lexical features**

**Method:**
**Lexicostatistics: mass comparison of basic lexical items,
extended by all relevant data available**

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**


**As in traditional lexicostatistics, but:**

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**As in traditional lexicostatistics, but:**

**1. use of computational algorithms and tools**

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**As in traditional lexicostatistics, but:**

**1. use of computational algorithms and tools**

**2. methodology from classification in biology**

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**WWW**

# ASJP

**Project:**
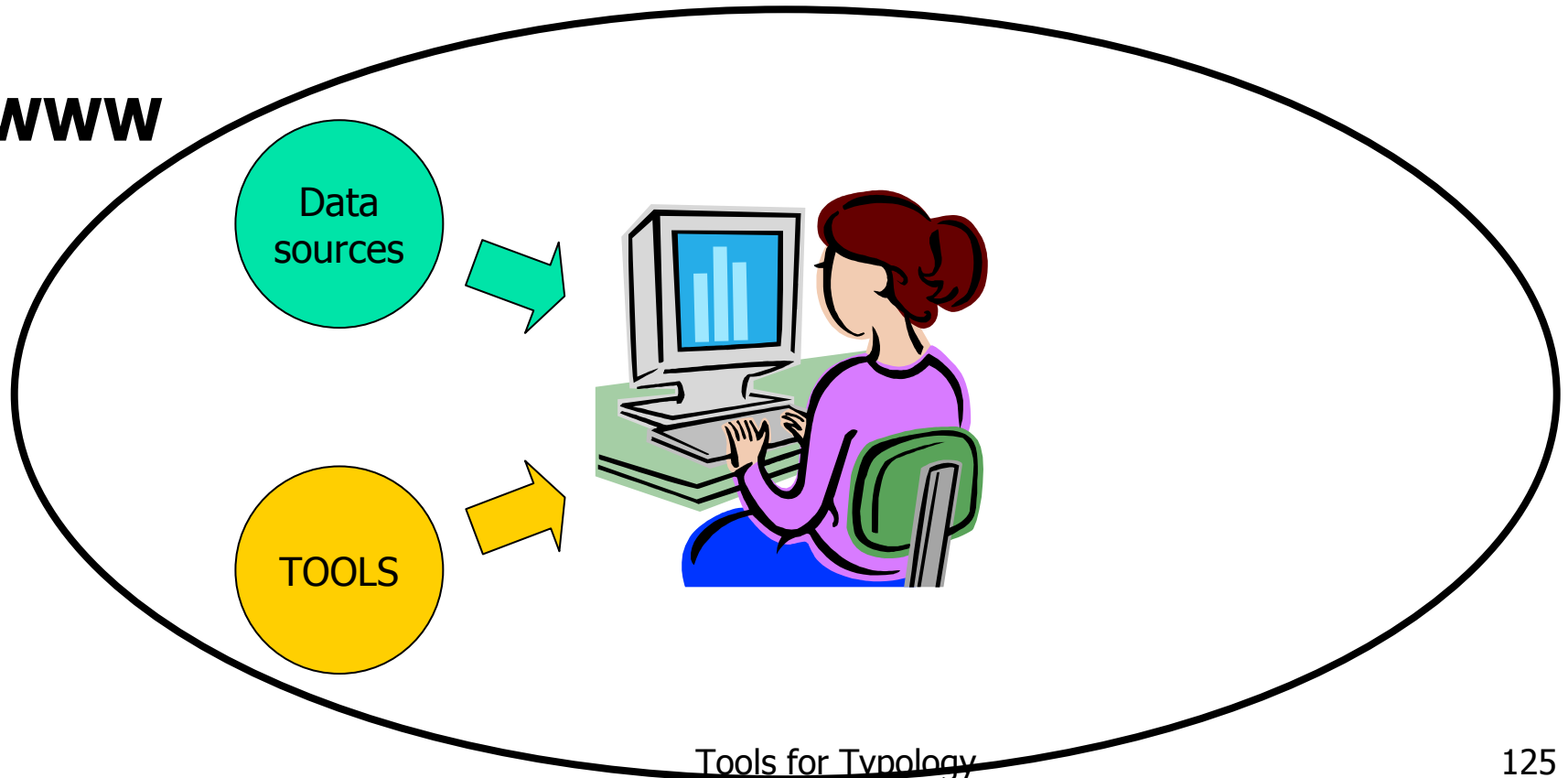**ASJP (Automated Similarity Judgment Program)**

**WWW**



Data
sources

# ASJP

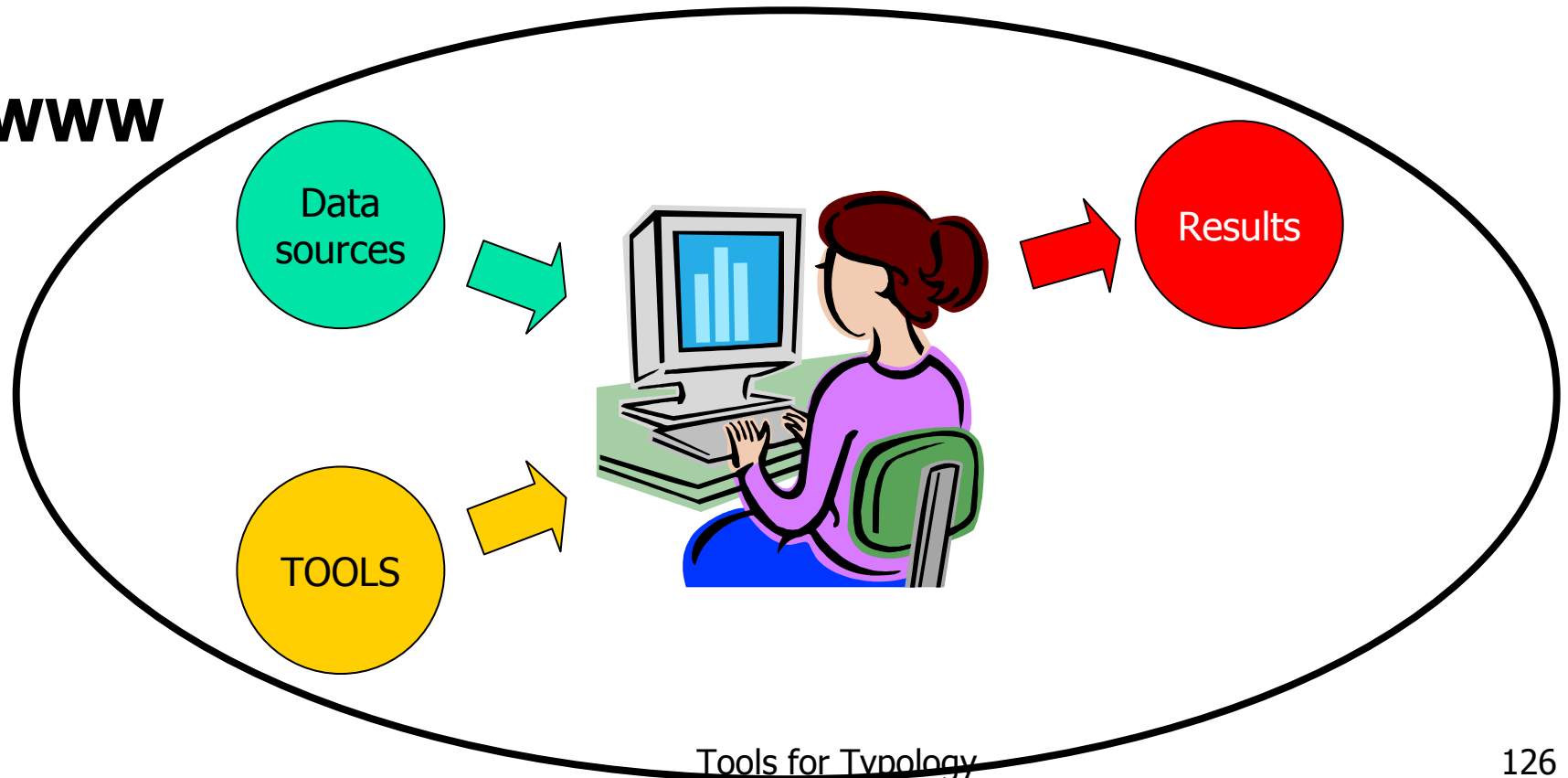**ASJP (Automated Similarity Judgment Program)**

**WWW**

Data sources

TOOLS

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**



**www**

Data sources

TOOLS

Results

Tools for Typology

126

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**WWW**



Data sources

TOOLS

Results

Data Base

Tools for Typology

# ASJP

**Project:**
**ASJP (Automated Similarity Judgment Program)**

**WWW**

Data sources

TOOLS

Results

Data Base

# ASJP

**ASJP (Automated Similarity Judgment Program)**

**WWW**

Tools for Typology

129

# Overview ASJP system

# Overview ASJP system

LEX

# Overview ASJP system

LEX

ASJP software

# Overview ASJP system



LEX → ASJP software → distance matrix

# Overview ASJP system

LEX

↓

ASJP software

↓

distance matrix

| DUTCH | ENGLISH | 53.3 |
|-------|---------|------|
| DUTCH | FRENCH | 72.7 |
| DUTCH | MANDARIN | 93.8 |
| … | | |

# Overview ASJP system

LEX
↓
ASJP software
↓
distance matrix
↓
CLASSIF software

# Overview ASJP system



LEX

ASJP software

distance matrix

CLASSIF software

**Existing Expert Classifications:**

ETHN
WALS
EXPRT

LEX

*EVALUATION*

ASJP software

distance matrix

STAT software

CLASSIF software

Tools for Typology

## Existing Expert Classifications:

ETHN
WALS
EXPRT

LEX

ASJP software

EVALUATION

distance
matrix

STAT
software

CLASSIF
software

Tools for Typology

Tools for Typology

HIST FACTS

GEO GRAPH

ETHN WALS EXPRT

LEX

ASJP software

distance matrix

MAP software

STAT software

CLASSIF software

Tools for Typology

140

Tools for Typology

# Lexical items

**Data: Word list Morris Swadesh (1955):**

**100 basic meanings**

| | | | | |
|---|---|---|---|---|
| 1. I | 21. dog | 41. nose | 61. die | 81. smoke |
| 2. you | 22. louse | 42. mouth | 62. kill | 82. fire |
| 3. we | 23. tree | 43. tooth | 63. swim | 83. ash |
| 4. this | 24. seed | 44. tongue | 64. fly | 84. burn |
| 5. that | 25. leaf | 45. claw | 65. walk | 85. path |
| 6. who | 26. root | 46. foot | 66. come | 86. mountain |
| 7. what | 27. bark | 47. knee | 67. lie | 87. red |
| 8. not | 28. skin | 48. hand | 68. sit | 88. green |
| 9. all | 29. flesh | 49. belly | 69. stand | 89. yellow |
| 10. many | 30. blood | 50. neck | 70. give | 90. white |
| 11. one | 31. bone | 51. breasts | 71. say | 91. black |
| 12. two | 32. grease | 52. heart | 72. sun | 92. night |
| 13. big | 33. egg | 53. liver | 73. moon | 93. hot |
| 14. long | 34. horn | 54. drink | 74. star | 94. cold |
| 15. small | 35. tail | 55. eat | 75. water | 95. full |
| 16. woman | 36. feather | 56. bite | 76. rain | 96. new |
| 17. man | 37. hair | 57. see | 77. stone | 97. good |
| 18. person | 38. head | 58. hear | 78. sand | 98. round |
| 19. fish | 39. ear | 59. know | 79. earth | 99. dry |
| 20. bird | 40. eye | 60. sleep | 80. cloud | 100. name |

# Lexical items: further reduction

**Early ASJP analyses have shown:**

→**It is not necessary to take all 100 words,**

   **but rather: the MOST STABLE subset**

# Lexical items: further reduction

**Early ASJP analyses have shown:**

→**It is not necessary to take all 100 words,**

    **but rather: the MOST STABLE subset**

**Least formal variation in accepted classifications**

**(e.g. Dryer's Genera; specialized classifications)**

| GERMANIC | FISH |
|---|---|
| AFRIKAANS | fis |
| BERNESE_GERMAN | fiS |
| BRABANTIC | fis |
| CIMBRIAN | fiS |
| DANISH | fesk |
| DUTCH | vis |
| ENGLISH | fiS |
| FAROESE | fiskur |
| FRANS_VLAAMS | fiS |
| FRISIAN_WESTERN | fisk |
| GOTHIC | fisks |
| ICELANDIC | fiskir |
| JAMTLANDIC | fisk |
| LIMBURGISH | vES |
| LUXEMBOURGISH | feS |
| NORTH_FRISIAN_AMRUM | fask |

| | | |
|---|---|---|
| **GERMANIC** | **FISH** | |
| **AFRIKAANS** | **fis** | |
| **BERNESE_GERMAN** | **fiS** | |
| **BRABANTIC** | **fis** | |
| **CIMBRIAN** | **fiS** | |
| **DANISH** | **fesk** | |
| **DUTCH** | **vis** | |
| **ENGLISH** | **fiS** | |
| **FAROESE** | **fiskur** | **1 proto form** |
| **FRANS_VLAAMS** | **fiS** | |
| **FRISIAN_WESTERN** | **fisk** | |
| **GOTHIC** | **fisks** | |
| **ICELANDIC** | **fiskir** | |
| **JAMTLANDIC** | **fisk** | |
| **LIMBURGISH** | **vES** | |
| **LUXEMBOURGISH** | **feS** | |
| **NORTH_FRISIAN_AMRUM** | **fask** | |

| GERMANIC | TREE |
|---|---|
| AFRIKAANS | bom |
| BERNESE_GERMAN | boum |
| BRABANTIC | bu3m |
| DANISH | trE7 |
| DUTCH | bom |
| ENGLISH | tri |
| FAROESE | trEa |
| FRANS_VLAAMS | bom |
| FRISIAN_WESTERN | bi3m\|by~Em |
| GOTHIC | bagms\|triu |
| ICELANDIC | th~ry~E |
| JAMTLANDIC | tre |
| LIMBURGISH | boum |
| LUXEMBOURGISH | bam |
| NORTH_FRISIAN_AMRUM | bum |
| NORTHERN_LOW_SAXON | bom |
| NORWEGIAN_BOKMAAL | tre |

| GERMANIC | TREE |
|---|---|
| AFRIKAANS | bom |
| BERNESE_GERMAN | boum |
| BRABANTIC | bu3m |
| DANISH | trE7 |
| DUTCH | bom |
| ENGLISH | tri |
| FAROESE | trEa |
| FRANS_VLAAMS | bom |
| FRISIAN_WESTERN | bi3m\|by~Em |
| GOTHIC | bagms\|triu |
| ICELANDIC | th~ry~E |
| JAMTLANDIC | tre |
| LIMBURGISH | boum |
| LUXEMBOURGISH | bam |
| NORTH_FRISIAN_AMRUM | bum |
| NORTHERN_LOW_SAXON | bom |
| NORWEGIAN_BOKMAAL | tre |

**2 forms**

| FIN-UGRIC | FISH |
|---|---|
| FINNISH | kala |
| ESTONIAN | kala |
| KARELIAN | kolo |
| KILDIN_SAAMI | kuly |
| KOMI_PERMYAK | Ceri |
| KOMI_ZYRIAN | cyeri |
| LULE_SAAMI | kuole |
| MEADOW_MARI | kol |
| MORDVIN(MOKSHA) | kEl |
| NORTH_SAAMI | guoli |
| SKOLT_SAAMI | kuel |
| SOUTH_SAAMI | gueli3 |
| UDMURT | cyorig |
| VEPS | kala |
| NENETS | xaly |
| SELKUP | q3l3 |
| CSANGO | hol |
| HUNGARIAN | hal |

| FIN-UGRIC | FISH |
|---|---|
| **FINNISH** | **kala** |
| **ESTONIAN** | **kala** |
| **KARELIAN** | **kolo** |
| **KILDIN_SAAMI** | **kuly** |
| **KOMI_PERMYAK** | **Ceri** |
| **KOMI_ZYRIAN** | **cyeri** |
| **LULE_SAAMI** | **kuole** |
| **MEADOW_MARI** | **kol** |
| **MORDVIN(MOKSHA)** | **kEl** |
| **NORTH_SAAMI** | **guoli** |
| **SKOLT_SAAMI** | **kuel** |
| **SOUTH_SAAMI** | **gueli3** |
| **UDMURT** | **cyorig** |
| **VEPS** | **kala** |
| **NENETS** | **xaly** |
| **SELKUP** | **q3l3** |
| **CSANGO** | **hol** |
| **HUNGARIAN** | **hal** |

**1 proto form**

Tools for Typology

| FIN-UGRIC | TREE |
|---|---|
| FINNISH | puu |
| INARI_SAAMI | muoro |
| KARELIAN | pu |
| KILDIN_SAAMI | mur |
| KOMI_PERMYAK | pu |
| KOMI_ZYRIAN | pu |
| LULE_SAAMI | muora |
| MEADOW_MARI | puSeNxe |
| MORDVIN(MOKSHA) | SuftE |
| NORTH_SAAMI | muoro |
| SKOLT_SAAMI | mu3r\|mw3r |
| SOUTH_SAAMI | moer3 |
| UDMURT | pispu |
| VEPS | pu |
| NENETS | pya |
| SELKUP | po |
| CSANGO | fo |
| HUNGARIAN | fa |

| FIN-UGRIC | TREE | |
|---|---|---|
| FINNISH | puu | |
| INARI_SAAMI | muoro | |
| KARELIAN | pu | |
| KILDIN_SAAMI | mur | |
| KOMI_PERMYAK | pu | |
| KOMI_ZYRIAN | pu | |
| LULE_SAAMI | muora | |
| MEADOW_MARI | puSeNxe | |
| MORDVIN(MOKSHA) | SuftE | **4 forms** |
| NORTH_SAAMI | muoro | |
| SKOLT_SAAMI | mu3r\|mw3r | |
| SOUTH_SAAMI | moer3 | |
| UDMURT | pispu | |
| VEPS | pu | |
| NENETS | pya | |
| SELKUP | po | |
| CSANGO | fo | |
| HUNGARIAN | fa | |

Tools for Typology

# Lexical items: further reduction

**Early analyses have shown:**

**Most stable 40/100 item subset gives:**

# Lexical items: further reduction

**<u>Early analyses have shown:</u>**

**Most stable 40/100 item subset gives:**

- **at least the same results as > 40**

# Lexical items: further reduction

**<u>Early analyses have shown:</u>**

**Most stable 40/100 item subset gives:**

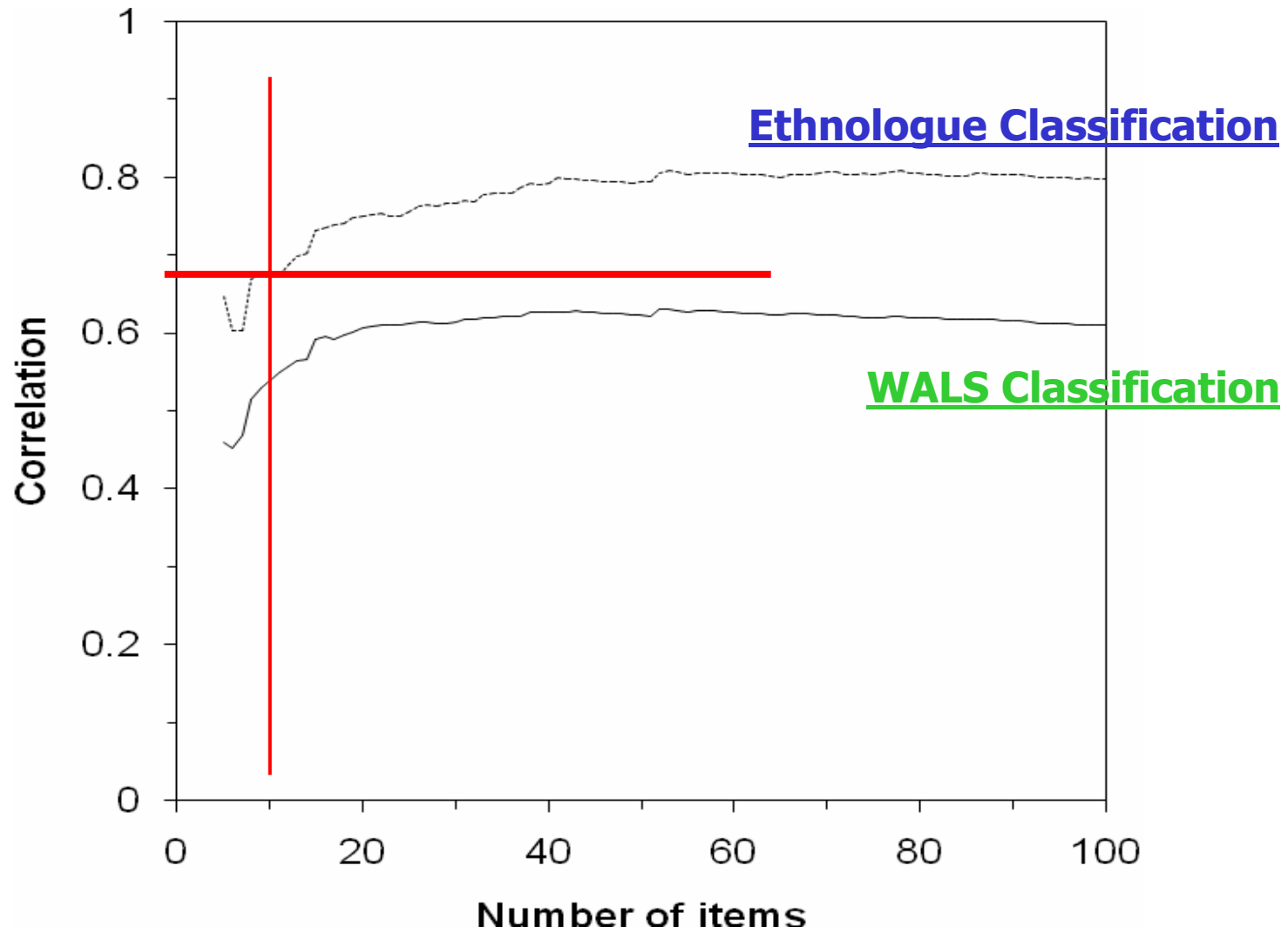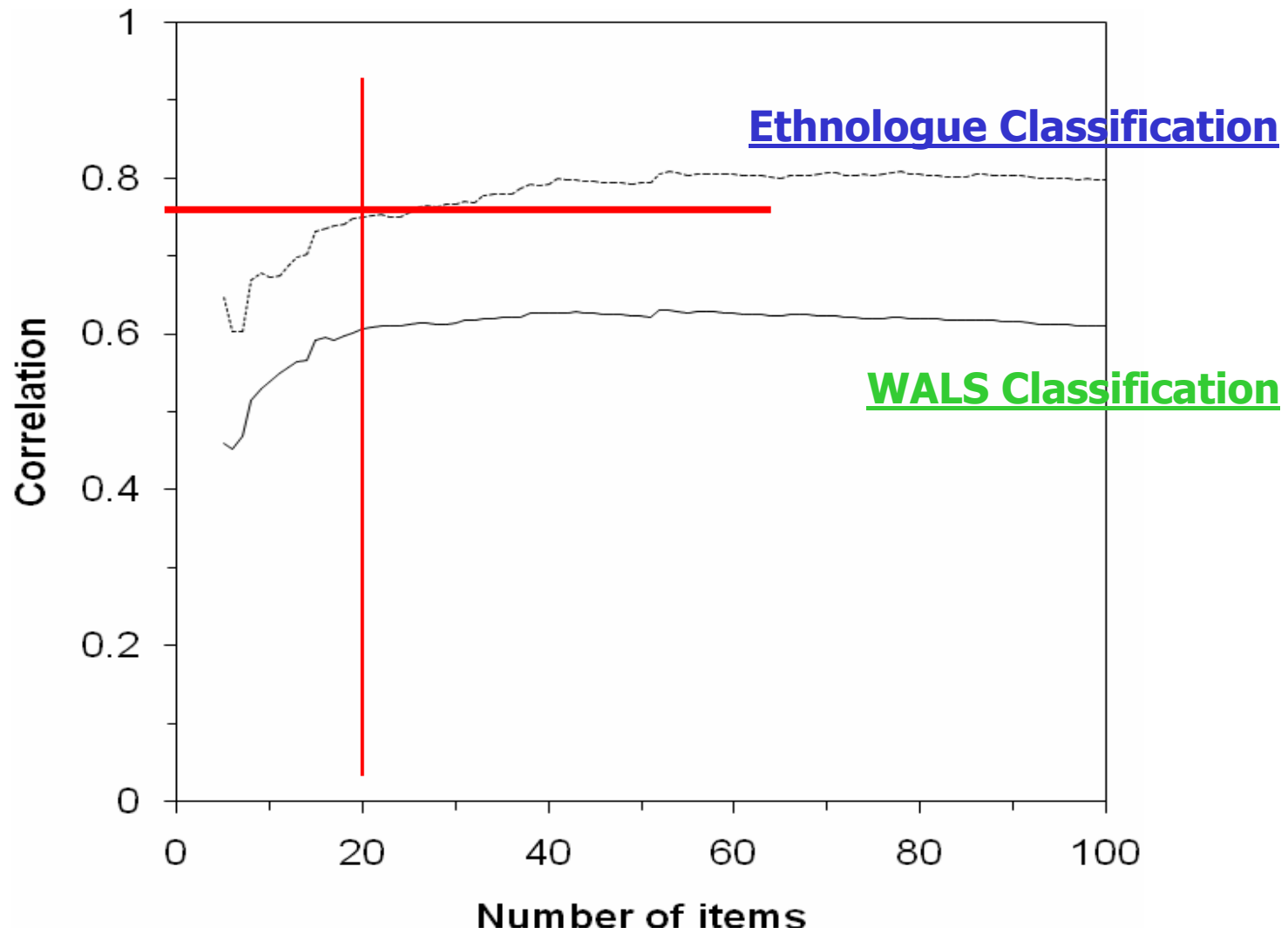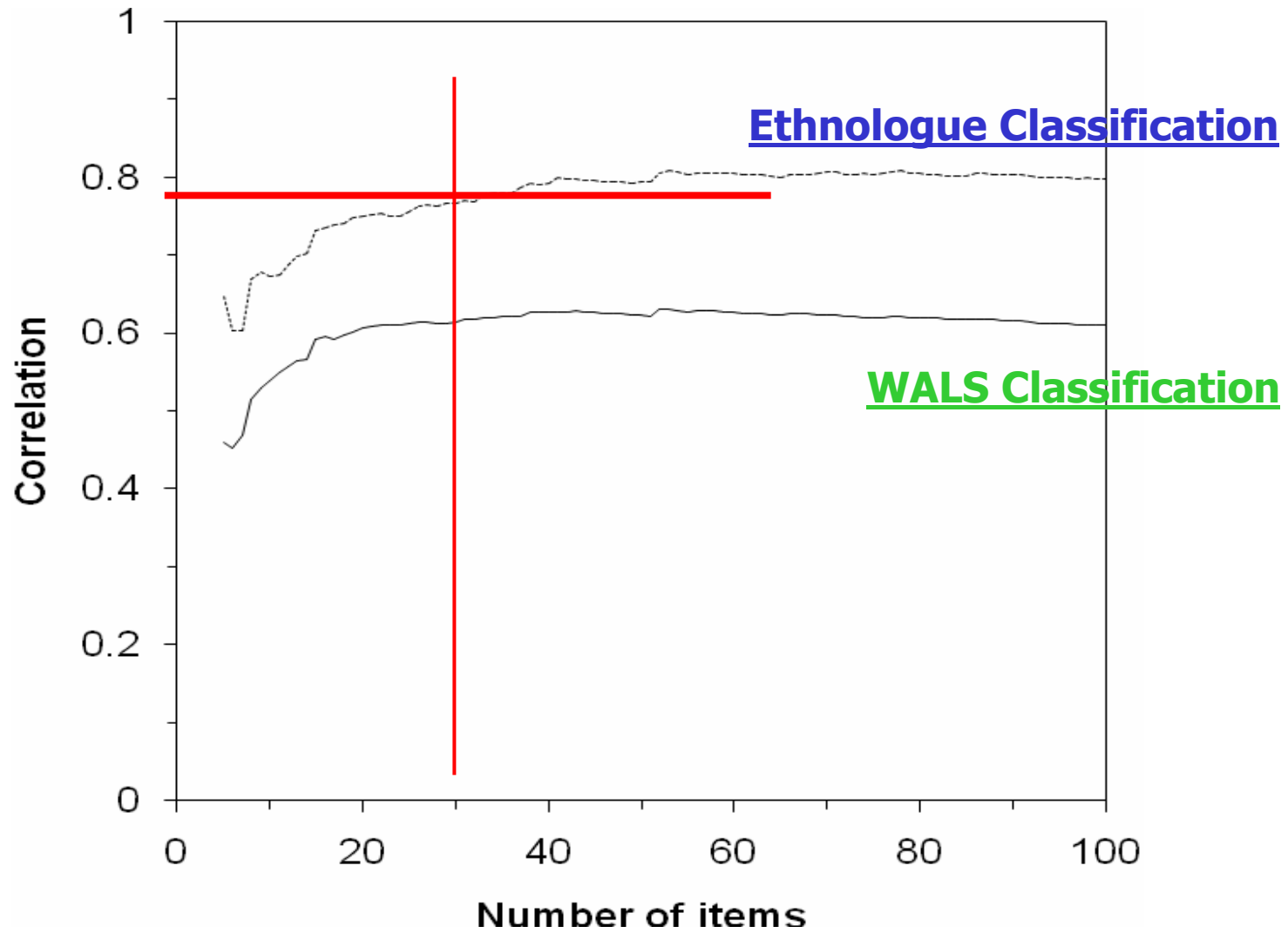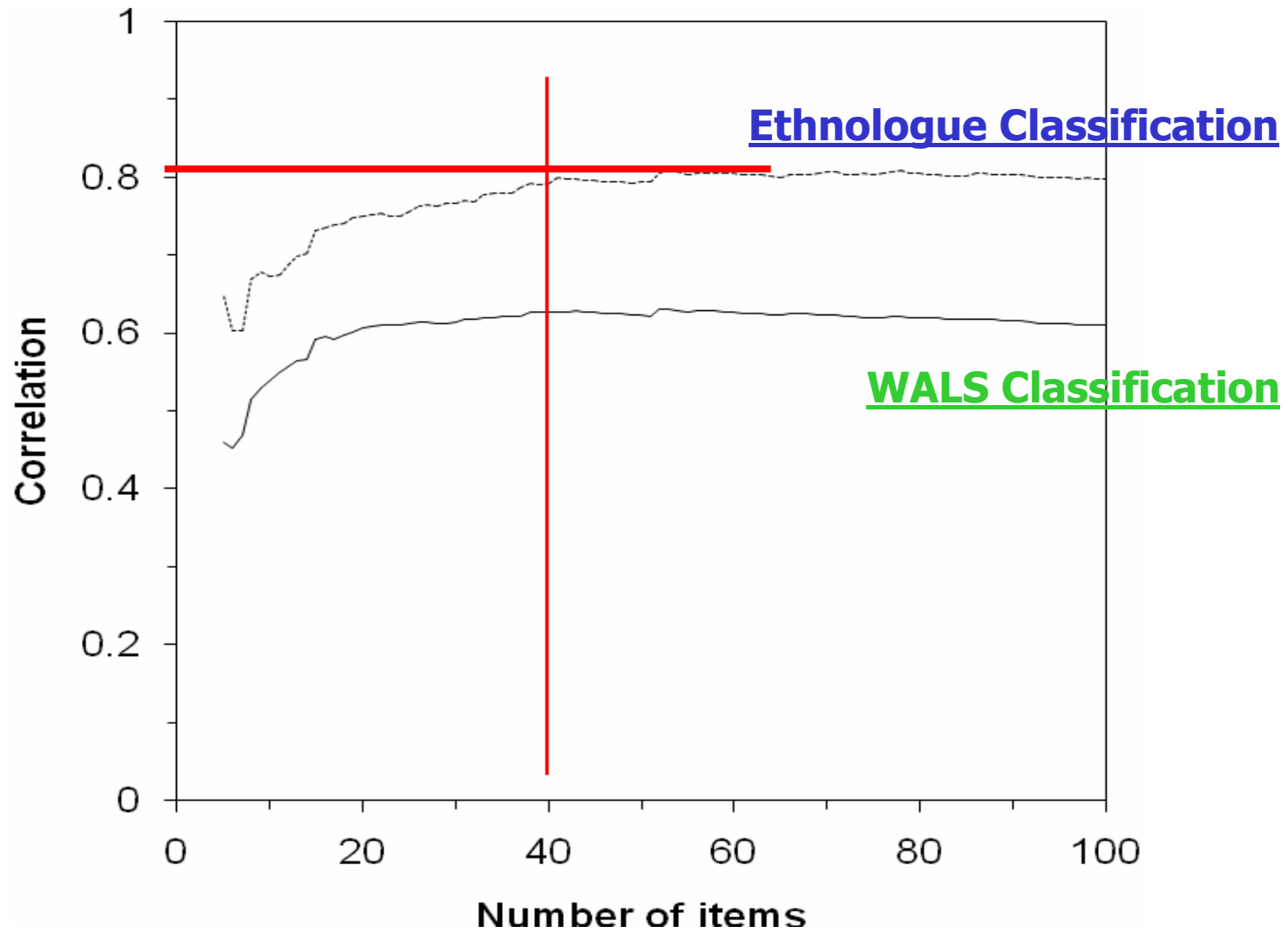- **at least the same results as > 40**

- **better results than < 40**

**Ethnologue Classification**

**WALS Classification**

Tools for Typology

| | | | | |
|---|---|---|---|---|
| I | dog | nose | die | smoke |
| you | louse | mouth | kill | fire |
| we | tree | tooth | swim | ash |
| this | seed | tongue | fly | burn |
| that | leaf | claw | walk | path |
| who | root | foot | come | mountain |
| what | bark | knee | lie | red |
| not | skin | hand | sit | green |
| all | flesh | belly | stand | yellow |
| many | blood | neck | give | white |
| one | bone | breast | say | black |
| two | grease | heart | sun | night |
| big | egg | liver | moon | hot |
| long | horn | drink | star | cold |
| small | tail | eat | water | full |
| woman | feather | bite | rain | new |
| man | hair | see | stone | good |
| person | head | hear | sand | round |
| fish | ear | know | earth | dry |
| bird | eye | sleep | cloud | name |

| | | | | |
|---|---|---|---|---|
| **I** | **dog** | **nose** | **die** | smoke |
| **you** | **louse** | mouth | kill | **fire** |
| **we** | **tree** | **tooth** | swim | ash |
| this | seed | **tongue** | fly | burn |
| that | **leaf** | claw | walk | **path** |
| who | root | foot | **come** | **mountain** |
| what | bark | **knee** | lie | red |
| not | **skin** | **hand** | sit | green |
| all | flesh | belly | stand | yellow |
| many | **blood** | neck | give | white |
| **one** | **bone** | **breast** | say | black |
| **two** | grease | heart | **sun** | **night** |
| big | egg | **liver** | moon | hot |
| long | **horn** | **drink** | **star** | cold |
| small | tail | eat | **water** | **full** |
| woman | feather | bite | rain | **new** |
| man | hair | **see** | **stone** | good |
| **person** | head | **hear** | sand | round |
| **fish** | **ear** | know | earth | dry |
| bird | **eye** | sleep | cloud | **name** |

**40 Most Stable**

# Lexical items: further reduction

**Early analyses have shown:**

- **Most stable 40/100 item subset gives optimal results**

→ **Less work**

# Lexical items: further reduction

**Early analyses have shown:**

- **Most stable 40/100 item subset gives optimal results**

→ **Less work**

→ **Less missing data**

# Lexical items: further reduction

**Early analyses have shown:**

- **Most stable 40/100 item subset gives optimal results**

→ **Less work**
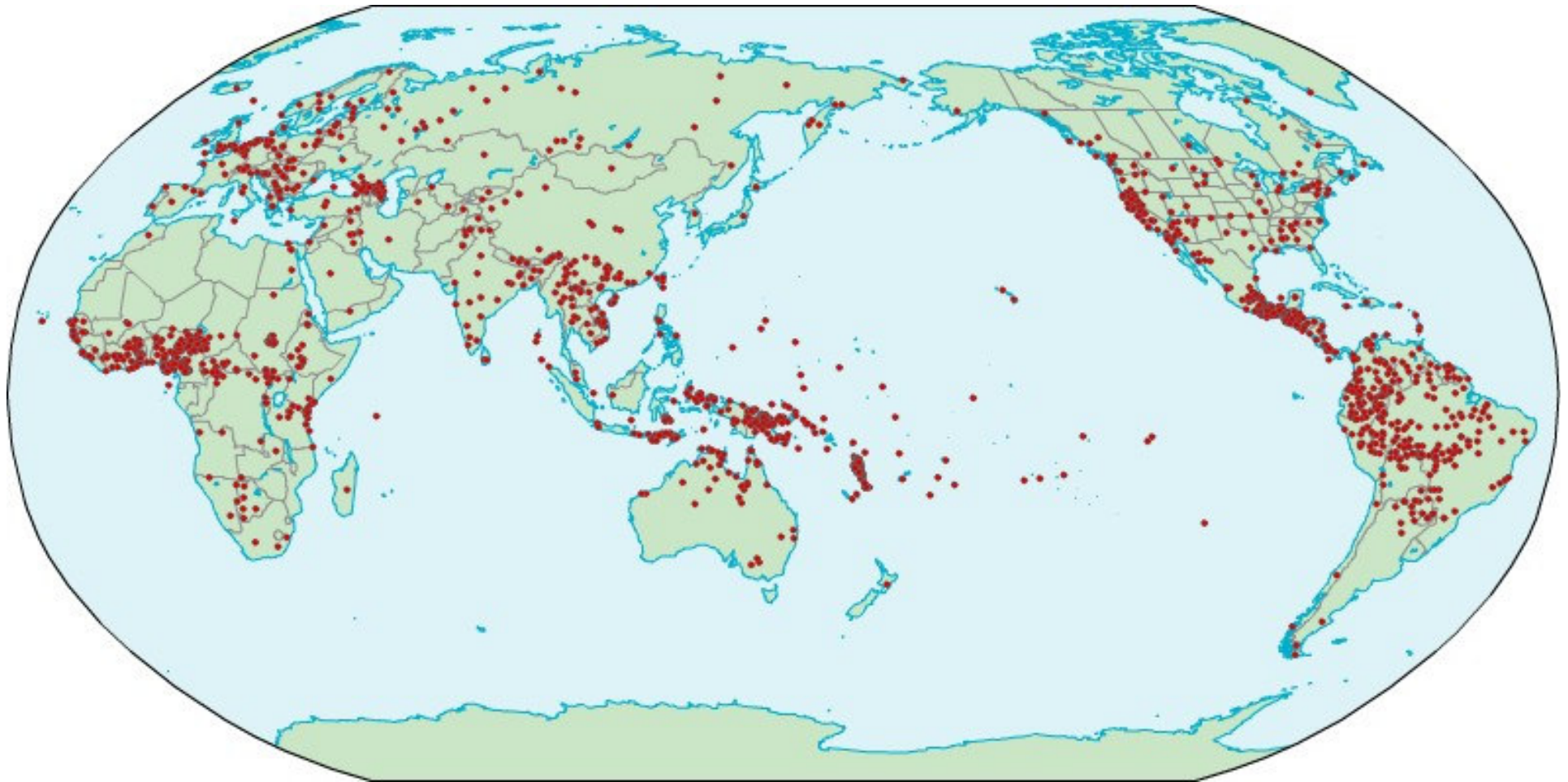
→ **Less missing data**

→ **Faster processing; combinatorial explosion:**

**40 : 100 ~ 2.5 * 2.5 = 6.3**

# Current sample

**3500 languages \* 40 lexical items**

**Languages currently sampled**

# Processing problems ...

**3500 languages \* 40 lexical items:**

**~ 10.000.000.000 comparisons .....**

# Processing problems ...

**3500 languages * 40 lexical items:**

**~ 10.000.000.000 comparisons ….. (10G)**

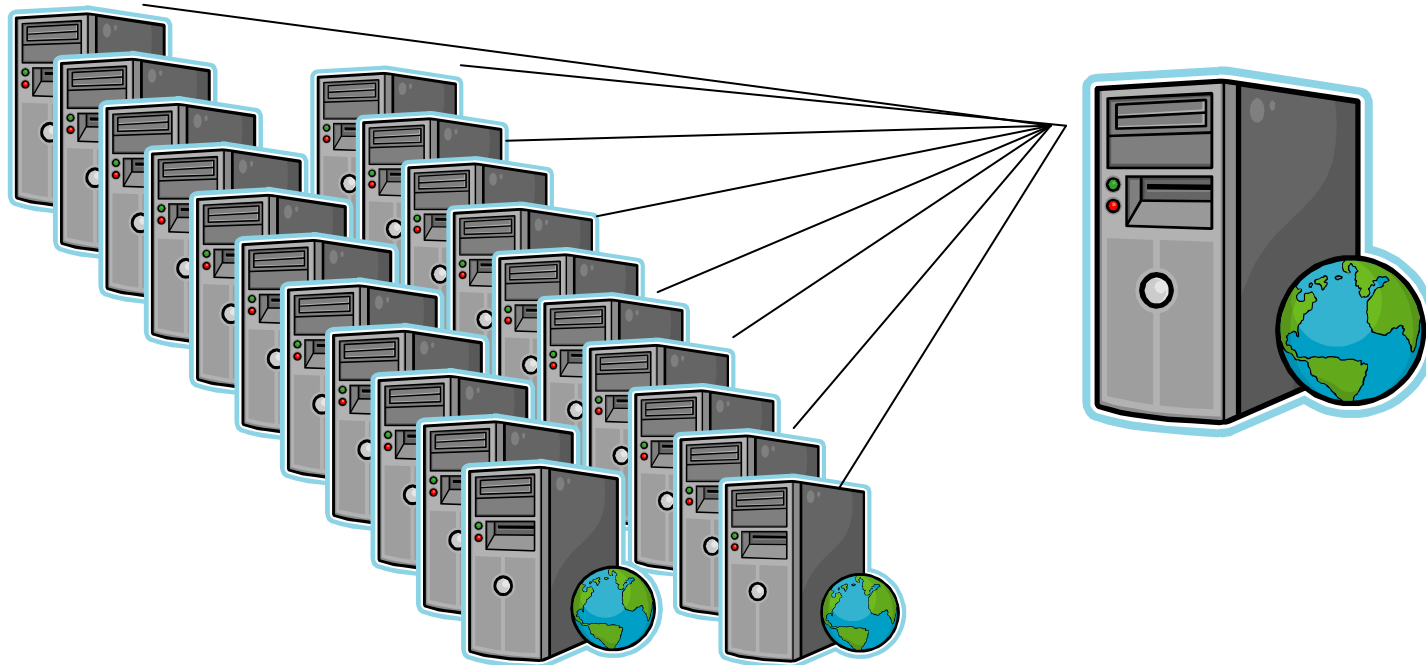**→ comparison at the phoneme level**

**for feature level: ~ 250.000.000.000 (0.25T)**

# Processing problems …

# Solution: parallel processing



**100 times faster**

# Lexical items: transcription

**First phase of project (2007):**

**Problems with full phonological (IPA) representation of words:**

# Lexical items: transcription

**First phase of project (2007):**

**Problems with full IPA representation of words:**

**- data entry via keyboard**

# Lexical items: transcription

**First phase of project (2007):**

**Problems with full IPA representation of words:**

**- data entry via keyboard**

**- simple programming languages (Fortran; Pascal)**

# Lexical items: transcription

**First phase of project (2007):**

**Problems with full IPA representation of words:**

**- data entry via keyboard**

**- simple programming languages (Fortran; Pascal)**

**→ Recoding to simplified ASJPcode (keyboard)**

# Lexical items: transcription

**ASJPcode:**

**7 Vowels**

**34 Consonants**

All other phonemes to 'closest sound'

# Lexical items: transcription

**ASJPcode:**

**7 Vowels**

**34 Consonants**

All other
phonemes
to
'closest sound'

**Symbols for:**  **Nasalization**
**Labialization**
**Palatalization**
**Aspiration**
**Glottalization**

## Abaza (Caucasian):

## Meaning

PERSON

LEAF

SKIN

HORN

NOSE

TOOTH

# Abaza (Caucasian):

| Meaning | IPA |
|---------|-----|
| PERSON | ʕʷɨtʃ'ʲʷʕʷɨs |
| LEAF | bɣʲɨ |
| SKIN | tʃʷazʲ |
| HORN | tʃ'ʷɨʕʷa |
| NOSE | pɨntsʼa |
| TOOTH | pɨts |

# Abaza (Caucasian):

| Meaning | IPA | | ASJPcode |
|---------|-----|---|----------|
| PERSON | ʕʷɨtʃʼʲʷʕʷɨs | ⟶ | Xw3Cw"yXw3s |
| LEAF | bɣʲɨ | ⟶ | bxy3 |
| SKIN | tʃʷazʲ | ⟶ | Cwazy |
| HORN | tʃʼʷɨʕʷa | ⟶ | Cw"3Xwa |
| NOSE | pɨntsʼa | ⟶ | p3nc"a |
| TOOTH | pɨts | ⟶ | p3c |

# Loss of information?

**<u>Experiment with Caucasian (39 lgs):</u>**

# Loss of information?

**<u>Experiment with Caucasian (39 lgs):</u>**

**- Full IPA does not score better for separating**
**language families**

# Loss of information?

**Experiment with Caucasian (39 lgs):**

**- Full IPA does not score better for separating language families**

**- For *precise genetic classification* IPA is even less accurate than ASJP code (too specific?)**

# Comparing words

**Most important measure: Levenshtein Distance**

# Comparing words

**Levenshtein Distance (LD)**

**a. between 2 words:**

# Comparing words

**Levenshtein Distance (LD)**

**a. between 2 words:**

**number of transformations (=changes & additions)
to get from the shorter form to the longer one**

# Comparing words

**Levenshtein Distance (LD)**

**a. between 2 words:**

**number of transformations (=changes & additions) to get from the shorter form to the longer one**

**b. between 2 languages:**

**mean LD for all common pairs**

# Comparing words

**Two problems with <span style="color:red">simple LD</span>:**

# Comparing words

**Two problems:**

**1. Value depends on length of longest word**

# Comparing words

**1. Value depends on length of longest word**

**C A T**

**D O G**

**x x x = 3**

# Comparing words

**1. Value depends on length of longest word**

**C A T**          **E L E P H A N T**

**D O G**          **D O G**

x x x = ③          x x x x x x x x = ⑧

# Comparing words

**1. Value depends on length of longest word**

**→ Normalize: LDN = ( LD / L$_{max}$ )**

# Comparing words

**1. Value depends on length of longest word**

C A T                E L E P H A N T

D O G                D O G

x x x = 3/3 = (1.0)       x x x x x x x x = 8 / 8 = (1.0)

# Comparing words

**<u>Two problems:</u>**

1. **Value depends on length of longest word**

→ **Normalize: LDN = ( LD / $L_{max}$ )**

2. **Differences between lgs in phonological overlap**

# Comparing words

**2. Differences between lgs in phonological overlap**

**DUTCH ~ ENGLISH:**      **mean LDN: 0.55**

# Comparing words

**2. Differences between lgs in phonological overlap**

**DUTCH ~ ENGLISH:**      **mean LDN: 0.55**
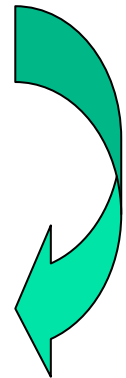
**DUTCH ~ MANDARIN:**    **mean LDN: 0.91**

# Comparing words

**2. Differences between lgs in phonological overlap**

**DUT ~ ENG:**     mean LDN: **0.55**

          mean LDN *other* words: **0.89**

**DUT ~ MAN:**     mean LDN: **0.91**

          mean LDN *other* words: **0.93**

# Comparing words

**2. Differences between lgs in phonological overlap**

**DUT ~ ENG:**       **mean LDN: 0.55 / 0.89**

                     **mean LDN *other* words: 0.89**

**DUT ~ MAN:**       **mean LDN: 0.91 / 0.93**

                     **mean LDN *other* words: 0.93**

# Comparing words

**2. Differences between lgs in phonological overlap**

**DUT ~ ENG:**       mean LDN: **0.55 / 0.89 = 0.62**

**DUT ~ MAN:**       mean LDN: **0.91 / 0.93 = 0.99**

# Comparing words

**Two problems:**

**1. Value depends on length of longest word**
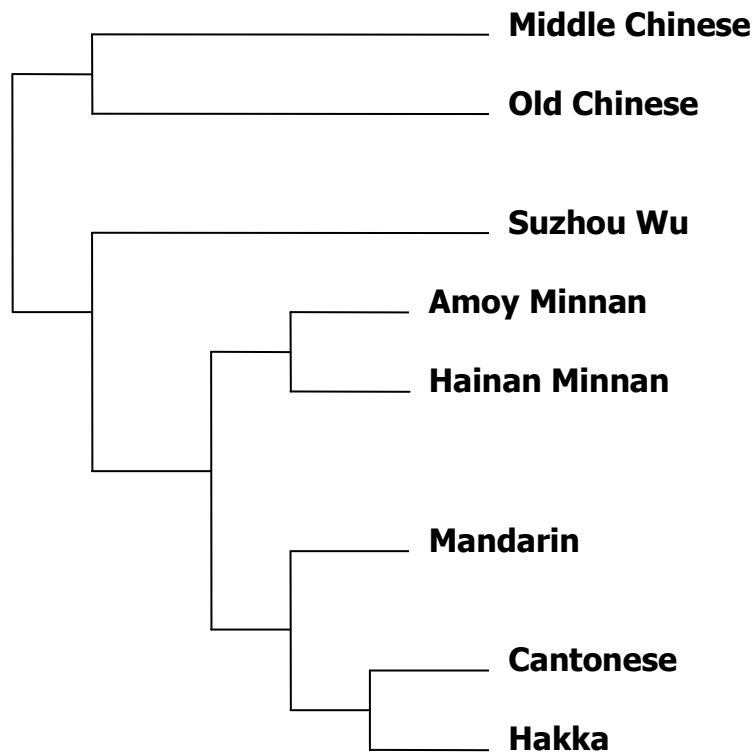
**→ Normalize: LDN = ( LD / $L_{max}$ )**

**2. Differences between lgs in phonological overlap**

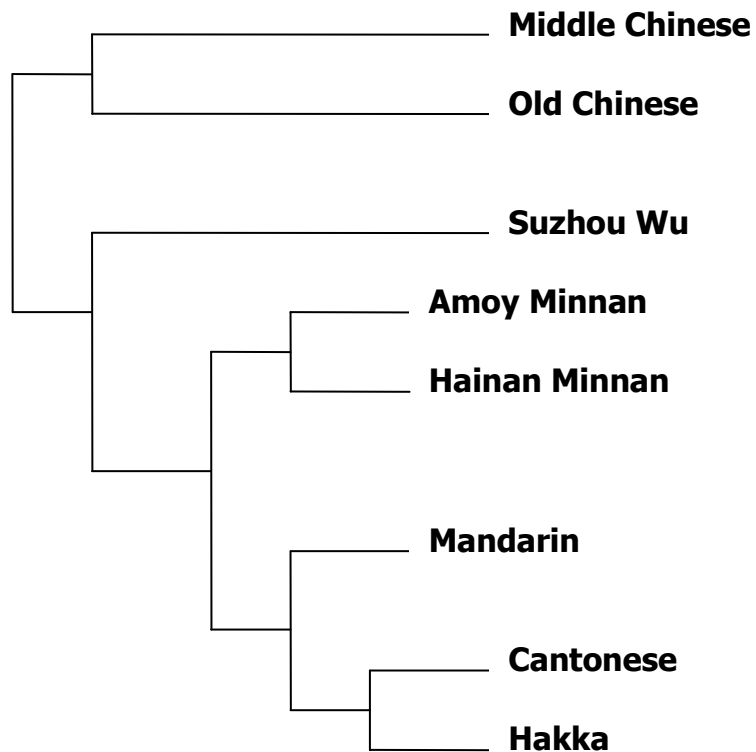**→ Eliminate ' background noise':**

$$LDND = ( LDN / LDN_{different\ pairs} )$$

# Sino-Tibetan: Chinese

```
┌──────────────────── Middle Chinese
│
├──────────────────── Old Chinese
│
│   ┌──────────────── Suzhou Wu
│   │
│   │       ┌──────── Amoy Minnan
│   │   ┌───┤
│   │   │   └──────── Hainan Minnan
└───┤   │
    │   │   ┌──────── Mandarin
    └───┤   │
        │   │   ┌──── Cantonese
        └───┤   │
            └───┤
                └──── Hakka
```

**ASJP tree based on lexical relations**

# Sino-Tibetan: Chinese



ALL & ONLY

- Middle Chinese
- Old Chinese
- Suzhou Wu
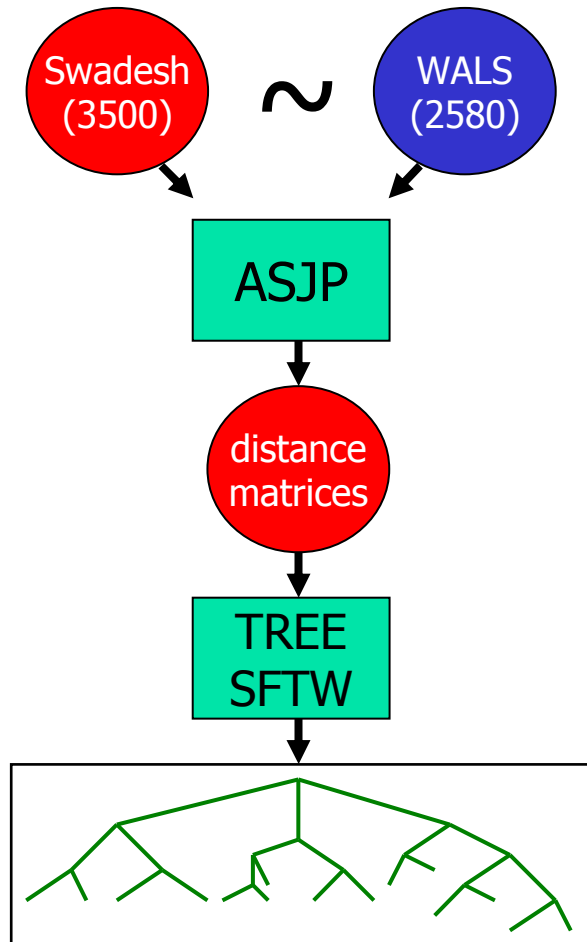- Amoy Minnan
- Hainan Minnan
- Mandarin
- Cantonese
- Hakka

**ASJP tree based on lexical relations**

# Sino-Tibetan: Chinese



Genetic classification in Thurgood & LaPolla (eds)

Tools for Typology

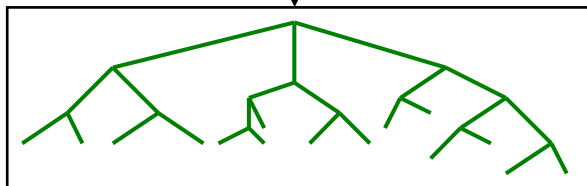# Lexical plus typological data

Swadesh (3500) ~ WALS (2580)

ASJP

distance matrices

TREE SFTW

Tools for Typology

'SWALSH'

ASJP

distance matrices

TREE SFTW

Tools for Typology

# Improving the fit



**Etnologue 2005**

**Only WALS**  **Percentage LDND**  **Only ASJP**

Tools for Typology

# Improving the fit



**Etnologue 2005**

**ONLY 550 LANGUAGES !!!**

Correlation

0.25000

0.20000

0.17500

0    25    50    75    100

**Only WALS**          **Percentage LDND**          **Only ASJP**

Tools for Typology                                        208

# Lexical items: transcription

**Second phase of project (2009-10):**

**Replace ASJP code by full IPA representations**

# Lexical items: transcription

**Second phase of project (2009-10):**

**Problems with full IPA representation solved:**

# Lexical items: transcription

**<u>Second phase of project (2009-10):</u>**

**Problems with full IPA representation solved:**

**1. scan/download/... full IPA representations**

# Lexical items: transcription

**Second phase of project (2009-10):**

**Problems with full IPA representation solved:**

**1. scan/download/… full IPA representations**

**2. automatic conversion IPA to integer (Python)**

# Lexical items: transcription

**Second phase of project (2009-10):**

**Problems with full IPA representation solved:**

**1. scan/download/... full IPA representations**

**2. automatic conversion IPA to integer (Python)**

**3. (semi-)automatic recoding to ASJPcode: transduction on the basis of a formal grammar**

# Lexical items: transcription

**Abaza (Caucasian):**
**Meaning**:     PERSON

# Lexical items: transcription

Abaza (Caucasian):
Meaning:       PERSON

IPA:              ʕʷɨtʃʼjʷʕʷɨs

# Lexical items: transcription

Abaza (Caucasian):
Meaning:        PERSON

IPA:            ʕʷɨtʃˈʲʷʕʷɨs

Decimal:        **661,695,616,679,700,690,695,661,695,616,115**

# Lexical items: transcription

Abaza (Caucasian):

Meaning: PERSON

IPA: ʕʷɨtʃʼʲʷʕʷɨs

Decimal: **661,695,616,679,700,690,695,661,695,616,115**

**'a' <- 661, 895, 416, …**    *formal grammar*

# Lexical items: transcription

Abaza (Caucasian):
Meaning:      PERSON

IPA:          ʕʷɨtʃ'ʲʷʕʷɨs

Decimal:      **661,695,616,679,700,690,695,661,695,616,115**

**'a' <- 661, 895, 416, ...**      *formal grammar*

**ASJP++code**

# Lexical items: transcription

IPA:  ʕʷɨtʃʼjʷʕʷɨs

Decimal: **661,695,616,679,700,690,695,661,695,616,115**

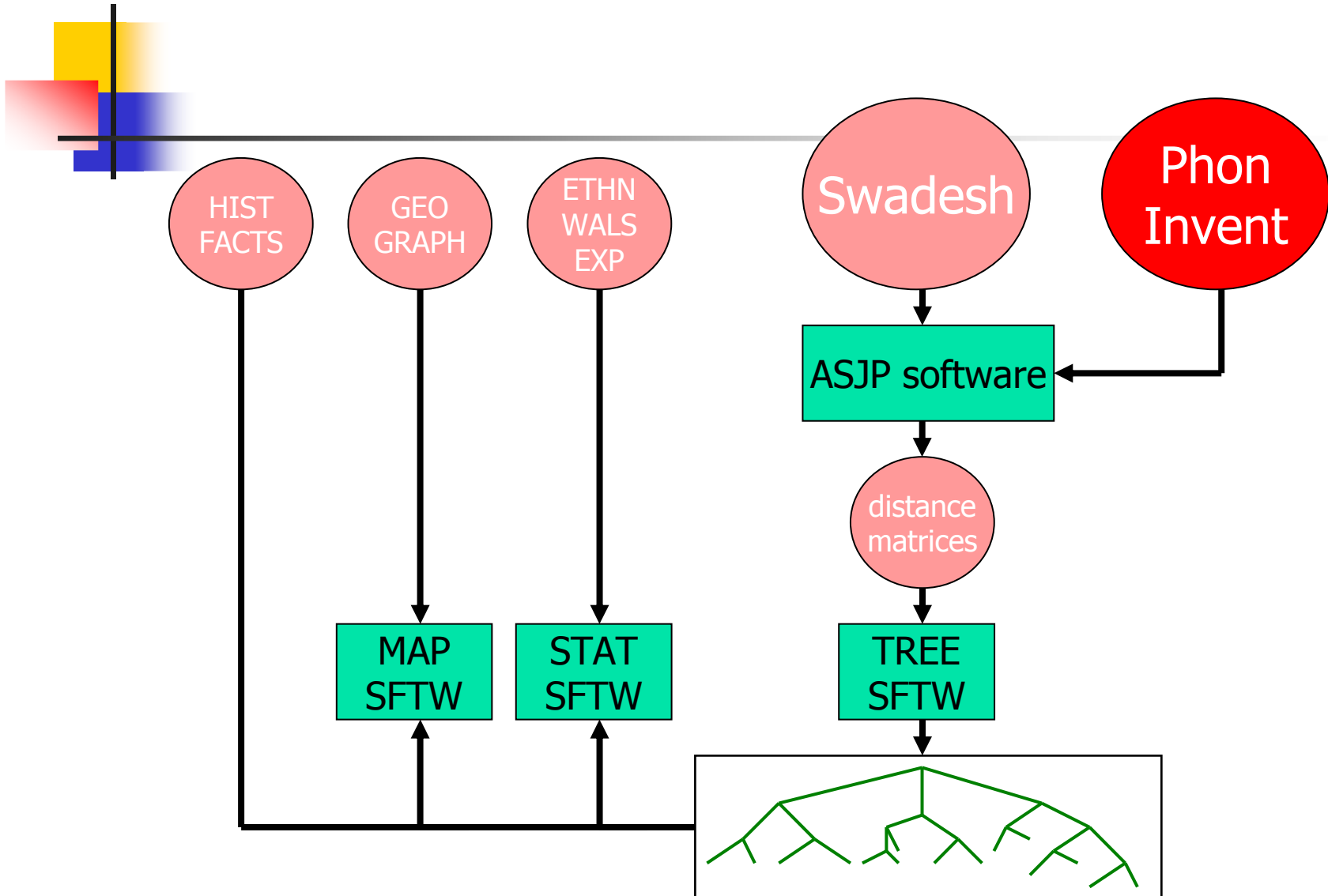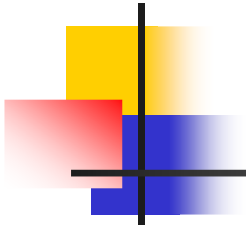**'a' <- 661, 895, 416, …**
**…**
**'a' [+Vow, +Low, +Middle]**
**'b' [+Cons, +Labial, +Plosive, +Voice ]**

*formal grammar*
*+*
*phonological*
*features*
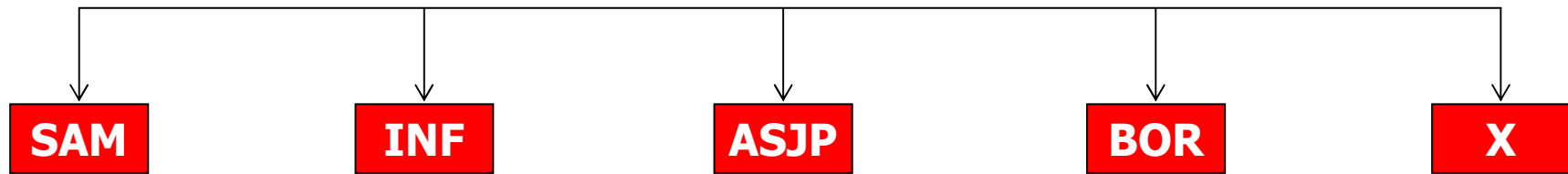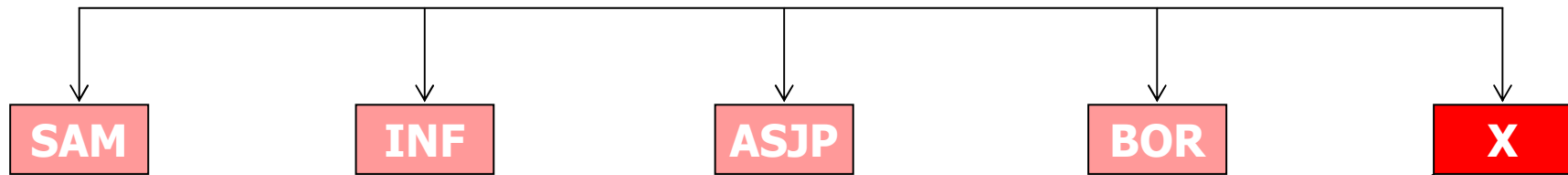
ASJP⁺⁺code: ( comparison of phonological features )

Tools for Typology

# 5. Accessibility

# Accessibility

```
┌──────────┬──────────┬──────────┬──────────┬──────────┐
│          │          │          │          │          │
▼          ▼          ▼          ▼          ▼
```

**SAM**    **INF**    **ASJP**    **BOR**    **X**

# Accessibility

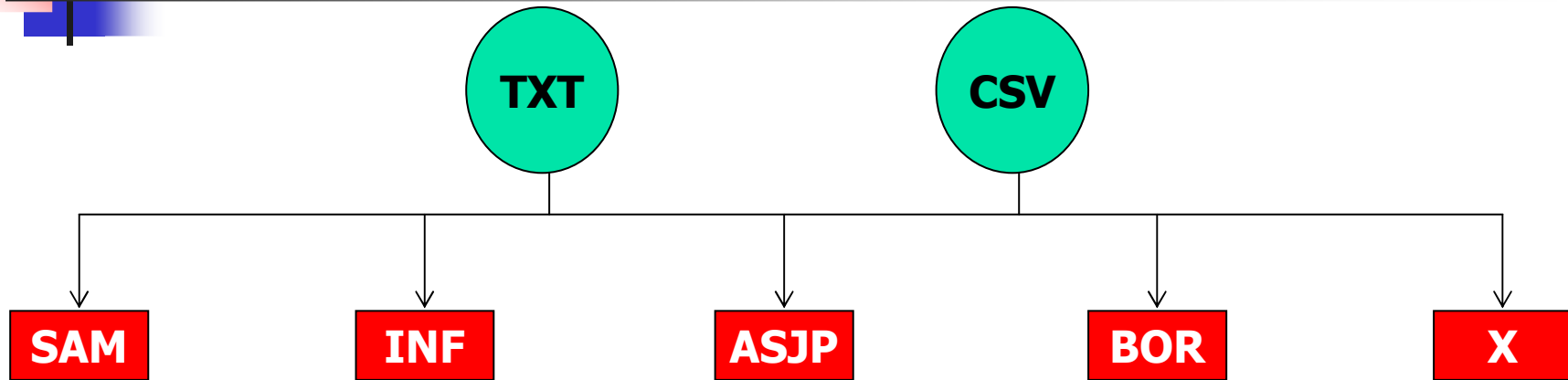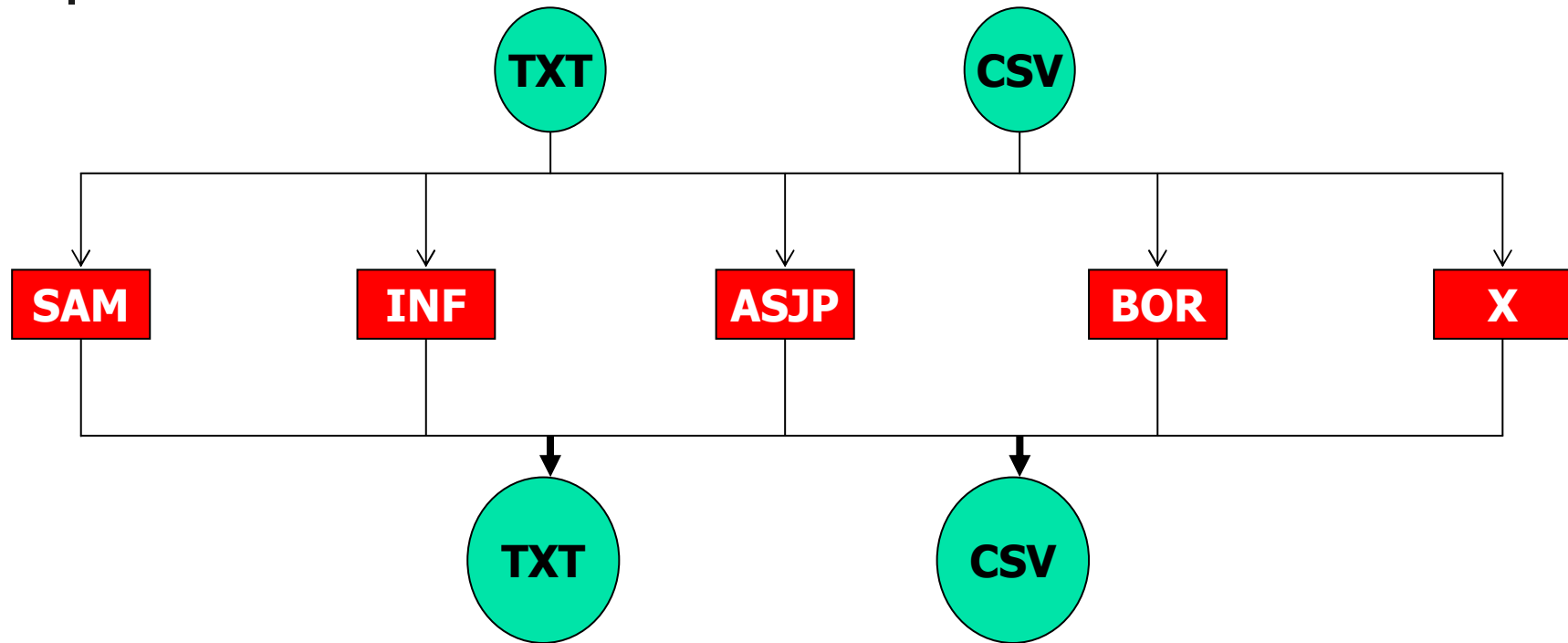| SAM | INF | ASJP | BOR | X |

**Small Tools:**
**Lookup Ethnologue code**
**Affiliation**
**Linguistic variables**
**...**

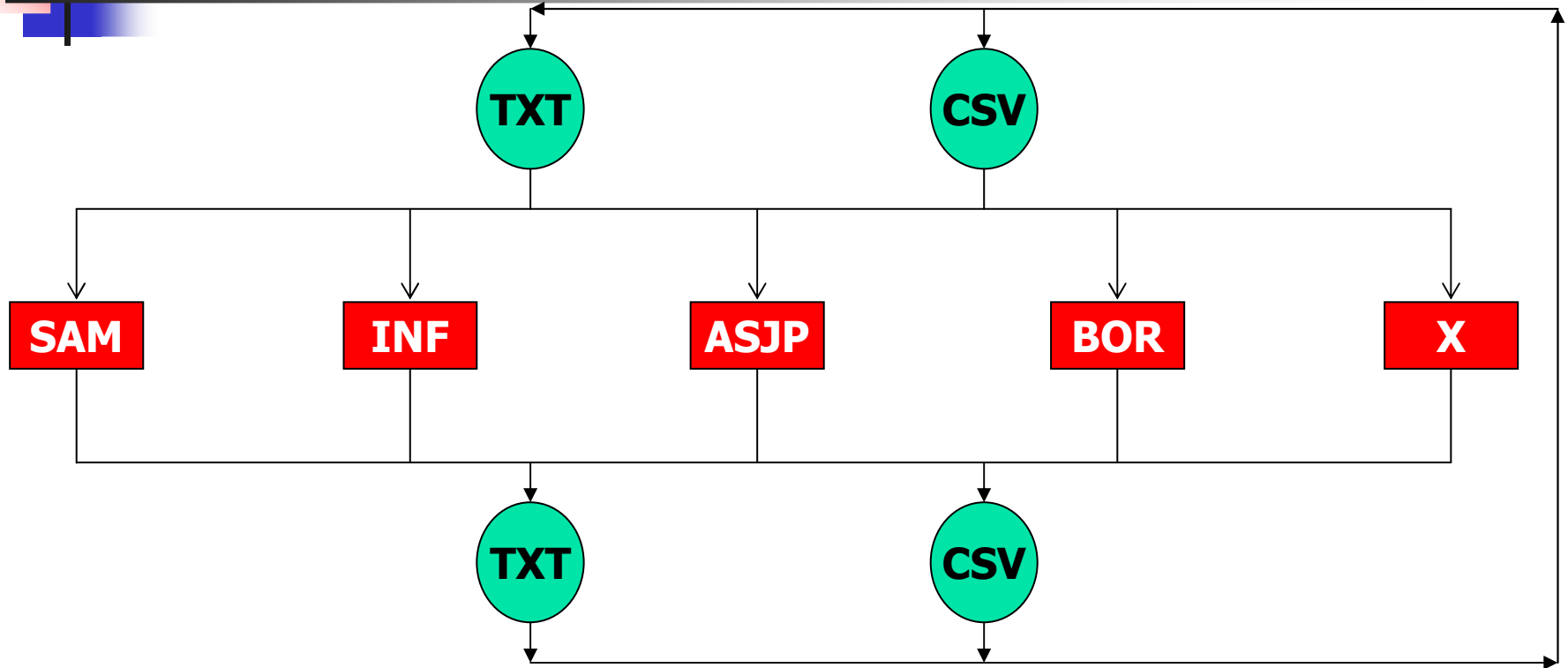# Access: data, internal



**Generally accepted data structures (Unicode; UTF8)**
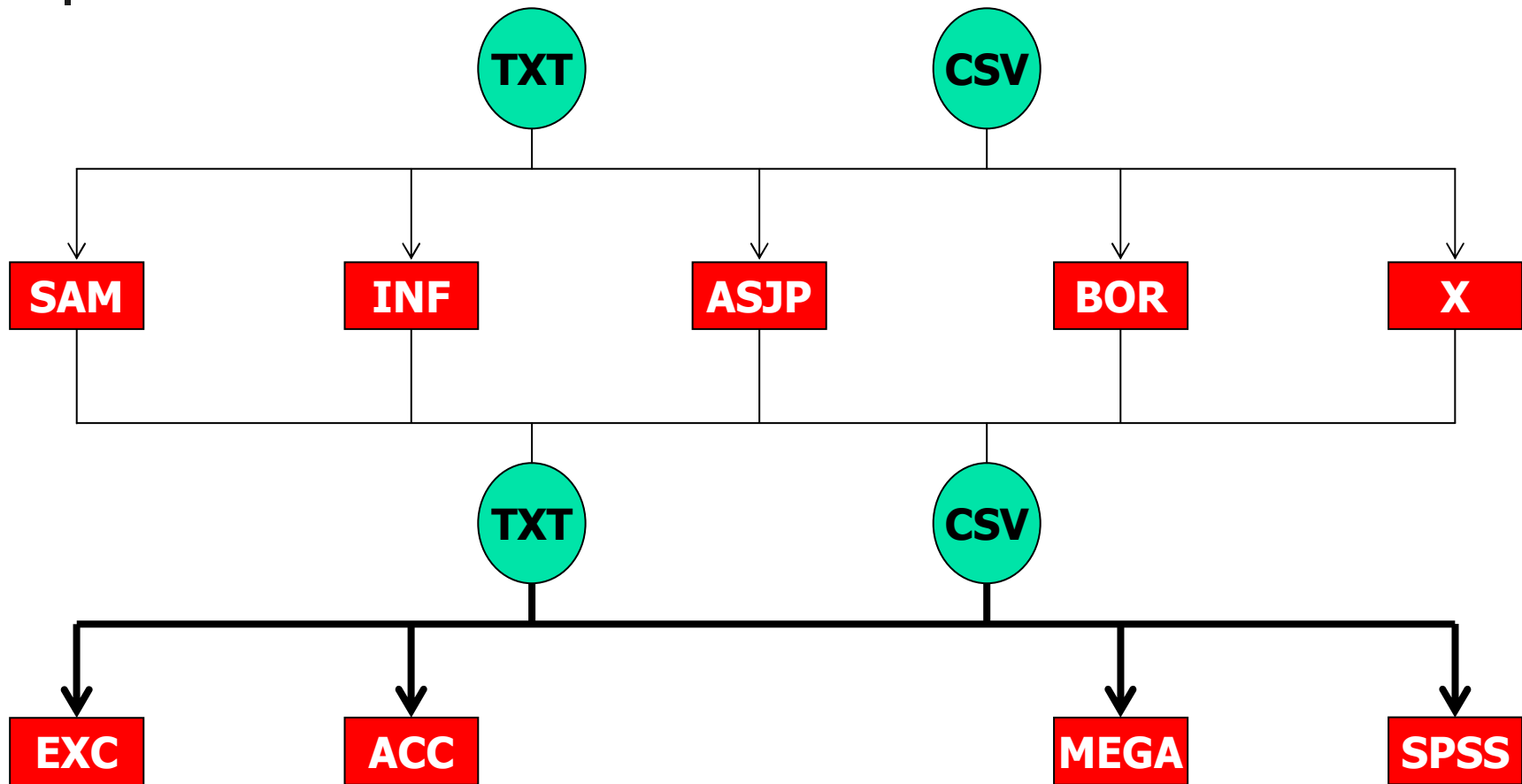
# Access: data, internal



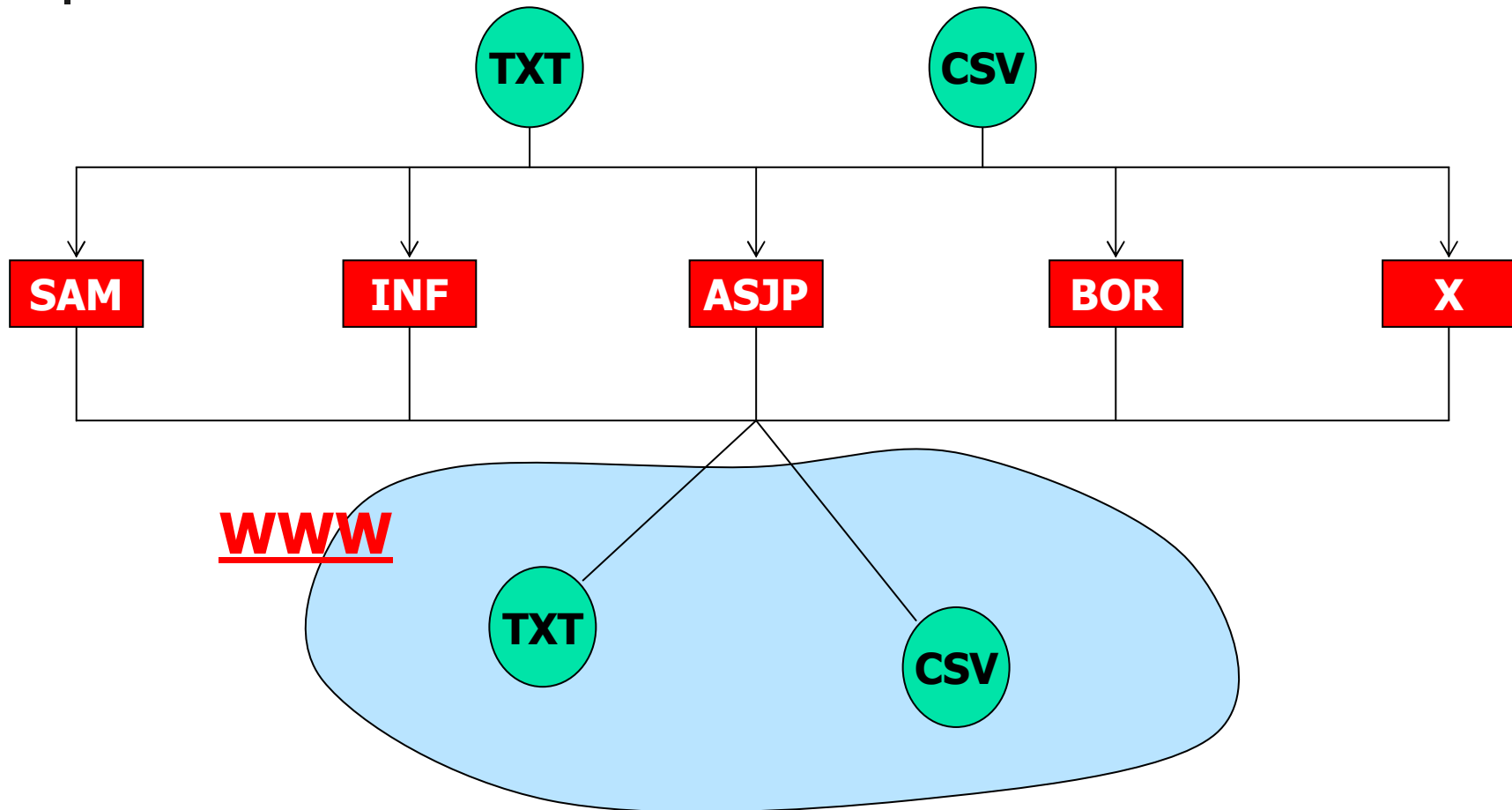**Generally accepted data structures (Unicode; UTF8)**

# Access: data, internal

# Access: data, external

# Access: data, universal?

# Access: software: in/external?

TXT

CSV

SAM INF ASJP BOR X

Pascal/C++
~
UNIX

## 1. Too big / slow for Windows (?)

# Access: software: in/external?

```
        ( TXT )                    ( CSV )
          │                          │
   ┌──────┼──────────┬──────────────┼────────┐
   ▼      ▼          ▼              ▼         ▼
[ SAM ] [ INF ]   [ ASJP ]       [ BOR ]   [ X ]
    \      \         │             /        /
     \      \   ┌─────────────┐   /        /
      \      \  │  Pascal/C++ │  /        /
       \      \ │      ~      │ /        /
        \      \│    UNIX     │/        /
                └─────────────┘
```

**1. Too big / slow for Windows (?)**
**2. No user interface**

# Access: software: in/external?

```
        TXT                    CSV

SAM      INF      ASJP      BOR       X

              Pascal/C++
                  ~
                UNIX
```

**1. Too big / slow for Windows (?)**
**2. No user interface**

# Access: software: in/external?

```
           ( TXT )                    ( CSV )
             │                          │
    ┌────────┼────────┬─────────────────┼────────┐
    ▼        ▼        ▼                 ▼        ▼
 ┌─────┐  ┌─────┐  ┌──────┐         ┌──────┐  ┌───┐
 │ SAM │  │ INF │  │ ASJP │         │ BOR  │  │ X │
 └─────┘  └─────┘  └──────┘         └──────┘  └───┘
     \        \       │                /        /
      \        \   ┌──────────────┐   /        /
       \        \  │  Pascal/C++  │  /        /
        _____\ │      ~       │ /_____/
                   │     UNIX     │
                   └──────────────┘
```

**1. Too big / slow for Windows (?)**
**2. No user interface**   ⟹   . . . . . . . . . .

# Access: software: in/external?

```
        ( TXT )                    ( CSV )
          |                          |
   ┌──────┼──────────┬──────────┬────┴─────┐
   ↓      ↓          ↓          ↓          ↓
[ SAM ] [ INF ]   [ ASJP ]   [ BOR ]    [ X ]
```

**Pascal/C++ ~ UNIX**

**1. Too big / slow for Windows (?)**
**2. No user interface**

**There must be more of such out there, some**

*useful* **for the linguistic community, but:**

# Accessibility requirements

**a. platform**
  **- accessible from WWW**
  **- programming language**

# Accessibility requirements

**a. platform**

      **- accessible from WWW**

      **- programming language**


**b. 'human' interface**

      **- interactive interface < - > actual application**

      **- user documentation**

# Accessibility requirements

**a. platform**

       **- accessible from WWW**
       **- programming language**

**b. 'human' interface**

       **- interactive interface < - > actual application**
       **- user documentation**

**c. data structure**

       **- TXT, CSV → HTML, Java Script, … (?)**

# Accessibility requirements

**a. platform**
- **accessible from WWW**
- **programming language**

**b. 'human' interface**
- **interactive interface < - > actual application**
- **user documentation**

**c. data structure**
- **TXT, CSV → HTML, Java Script, … (?)**

**d. maintenance**
- **programmer documentation**

**?**