# Estimating and Visualizing Language Similarities Using Weighted Alignment and Force-Directed Graph Layout

Gerhard Jäger

April 24, 2012, Avignon

joint work with Armin Buch, David Erschler & Andrei Lupas

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# Force Directed Graph Layout

- method to visualize graphs or similarity matrices in two or three dimensions
- simulation of a physical system:
  - data items $\Leftrightarrow$ physical particles
  - pairwise attractive force between particles proportional to their similarity
  - constant repelling force between any pair of particles
  - *this is just one of many protocols to determine forces*
    - initially, all particles are placed at random
    - in each time step, each particle is move a small amount along the resulting force vector
    - last step is repeated until a stable state is reached
- tends to stabilize in a state where groups of mutually similar items form clusters

# CLANS

- **Cl**uster **An**alysis of **S**equences
- developed by bioinformaticians *Tancred Frickey* and *Andrei Lupas* as exploratory tool to explore evolutionary relationships among protein sequences (Frickey and Lupas 2004)
- similarities of proteins is determined via sequence alignment; resulting matrix is visualized using CLANS
- advantages in comparison to tree-based algorithms:
  - does not *a priori* assume a tree like signal (useful when lateral transfer plays a role)
  - fast (esp. in comparison to character based algorithms)
  - robust (noise in data items does not accumulate)
- general impression so far (Lupas, p.c.):
  - tree algorithms are more precise when evolutionary distances are small; CLANS is more sensitive to weak evolutionary signals

# The Automated Similarity Judgment Program

- Project at MPI EVA in Leipzig around Sören Wichmann
- covers more than 5,000 languages
- basic vocabulary of 40 words for each language, in uniform phonetic transcription
- freely available

**used concepts:** *I, you, we, one, two, person, fish, dog, louse, tree, leaf, skin, blood, bone, horn, ear, eye, nose, tooth, tongue, knee, hand, breast, liver, drink, see, hear, die, come, sun, star, water, stone, fire, path, mountain, night, full, new, name*

# First shot: Levenshtein Distance

- first step: finde minmal edit distance between all translation pairs of the languages to be compared
- e.g. German ↔ Latin

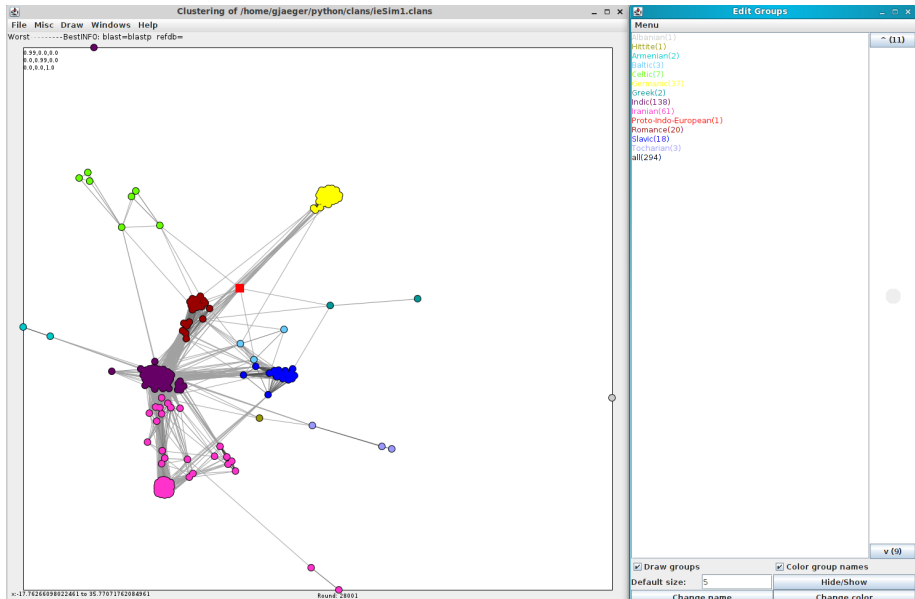$$\begin{array}{cccc} \text{h} & \text{o} & \text{r} & \text{n} \\ | & | & | & | & | \\ \text{k} & \text{o} & \text{r} & \text{n} & \text{u} \end{array}$$

- edit distance $= 2$
- transformation into similarity measure

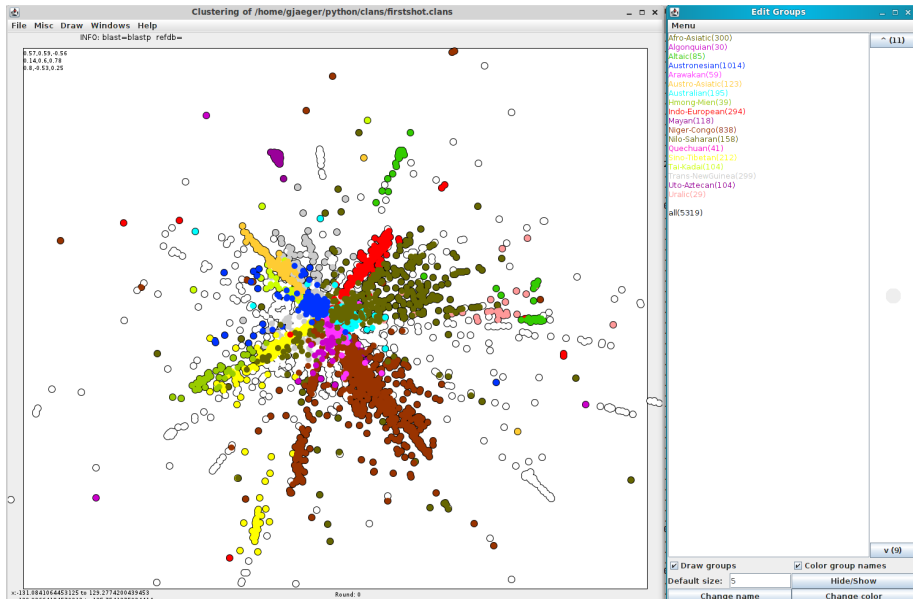$$\mathrm{sim}(x, y) \doteq \frac{2(\max(l(x), l(y)) - d_{Lev}(x, y))}{l(x) + l(y)}$$

- similarity between L1 and L2: average similarity of translation pairs between L1 and L2

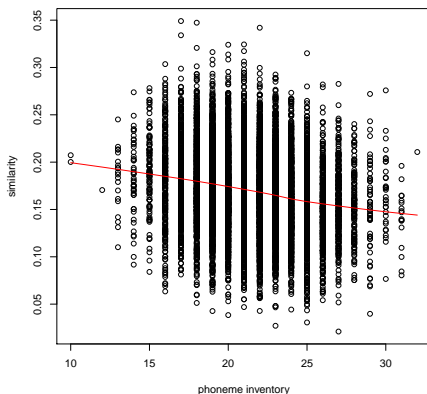# First shot: normalized Levenshtein Distance

# First shot: normalized Levenshtein Distance

# First shot: normalized Levenshtein Distance

- basic problem here: the smaller the sound inventories of the
  languages compared, the higher is the probability of false positives
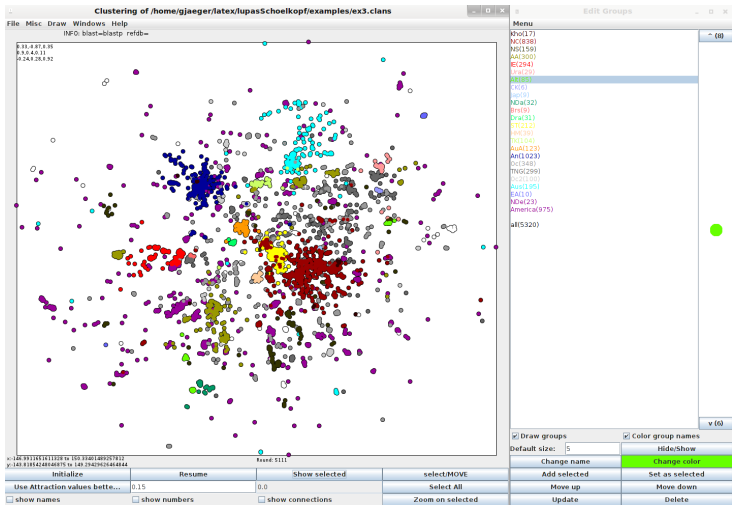
## Benchmark: LDND measure

- Wichmann et al.: doubly normalized Levenshtein distance (**L**evenshtein **D**istance **N**ormalized and **D**ivided)
- normalization for word length

$$\text{nld}(x, y) \doteq \frac{d_{Lev}(x, y)}{\max(l(x), l(y))} \tag{1}$$

- normalization for language specific patterns (including sound inventory size):
  - normalization factor $1/\mu$
  - $\mu_{L_1, L_2}$: mean of $\{nld(x, y) | x \in L_1, y \in L_1, \|x\| \neq \|y\|\}$

$$
\begin{aligned}
\text{ldnd}(x, y, L_1, L_2) &\doteq \frac{nld(x, y)}{\mu_{L_1, L_2}} \\
\text{ldnd}(L_1, L_1) &\doteq \frac{\sum_{x \in L_1, y \in L_2} \{\text{ldnd}(x, y, L_1, L_2) : \|x\| = \|y\|\}}{\#\{x, y : \|x\| = \|y\|\}}
\end{aligned}
$$

# Benchmark: LDND measure

# Needleman-Wunsch-Algorithmus

- Levenshtein distance is somewhat coarse grained

```
h  a  n  t        h  a  n  t
|  |  |  |        |  |  |  |
h  E  n  d        m  a  n  o
```

- simply normalized distance is $0.5$ in both cases
- after second normalization, [hant] even appears somewhat closer to [mano] ($ldnd = 0.54$) than to [hEnd] ($ldnd = 0.55$)
- correspondences $a \sim E$, $t \sim d$ are (according to linguistic criteria like place of articulation) much more natural than $h \sim m$ or $t \sim o$
- German appears equidistant to English and Spanish here, even though the distance to English is clearly smaller

# Weighted Alignment

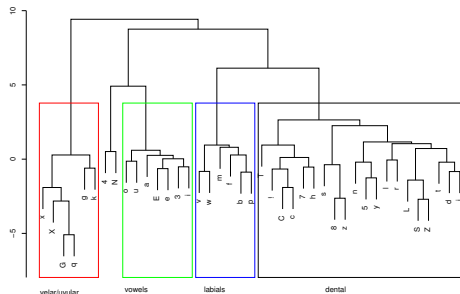- **Needleman Wunsch Algorithm**
    - similar to computation of Levenshtein distance
    - edit operations are *weighted*: algorithm finds optimal alignment, that minimizes total weight
    - $a \sim E$, $d \sim t$ should have lower weight than $t \sim o$
- How to determine these weights?
    - bioinformatics: *log-odds*
    - logarithm of the probability of a replacement, divided by probability of chance co-occurrence of molecula pair in question

# Weighted Alignment

- estimation of correspondence probabilities of two sounds in cognates:
    - large sample of pairwise related languages
    - replacement operation under Levenshtein alignment of translation pairs are counted
    - substantive part of word pairs considered are true cognates: replacement operations thus reflect genuine language change processes
    - replacement of sounds between non-cognates is randomly distributed and boils down to an additive constant in the logarithmic term

# Weighted Alignment

- log odds:
  - $d \sim t$: 0.69
  - $a \sim E$: 0.07
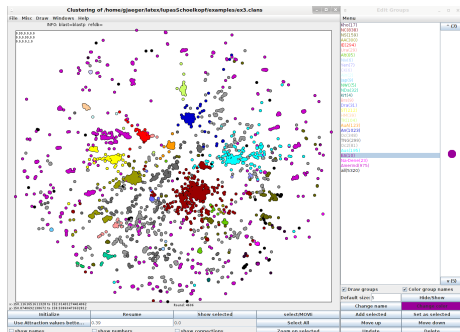  - $h \sim m$: $-0.61$
  - $t \sim o$: $-0.80$

## Weighted Alignment

- total value of optimal alignment is interpreted as similarity between strings
- similarity between languages is computed via $p$-values:

$$\mathrm{nwpv}(x, y, L_1, L_2) \;\doteq\; \frac{\#\{(x', y')|\mathrm{nw}(x', y') \geq \mathrm{nw}(x, y)\}}{\#L_1 \times \#L_2}$$

$$\mathrm{nwpv}(L_1, L_2) \;\doteq\; \frac{\sum_{x \in L_1, y \in L_2} -\log(\mathrm{pv}(x, y))}{\#\{(x, y) : \|x\| = \|y\|\}}$$

# Weighted Alignment

- similarities of English to
    - Dutch 0.60 / 3.38
    - German: 0.68 / 3.16
    - Proto-Indoeuropean: 0.86 / 2.33
    - Latin: 0.88 / 1.85
    - Spanish: 0.93 / 1.59
    - Russian: 0.93 / 1.52
    - Hungarian: 0.95 / 1.30
    - Turkish: 1.03 / 0.83

# Comparison

- two reasonable Gold standards for comparing these two similarity/distance measures:
  - expert judgments on cognacy
  - expert judgments on language classification

# Comparison: cognacy

- **Dyen-Kruskal database:** cognacy judgment for 200-item Swadesh lists from 95 Indo-European languages
- experiment:
  - extract those items from the Dyen-Kruskal database that occur in ASJP
  - define a cognacy estimator based on LDND by finding the optimal cutoff
  - do the same for NWPV
  - compare
- result
  - LDND: optimally achievable *Matthews Correlation Coefficient*: **0.547**
  - NWPV: optimally achievable *Matthews Correlation Coefficient*: **0.574**
    (*+1 means perfect prediction, -1 perfect mis-prediction*)

# Comparison: language classification

- **Ethnologue:** provides taxonomic classification of virtually all living languages
- **Robinson-Foulds metric:**
  - compares two trees over the same set of leafs
  - returns number of partitions that one of the two trees makes and the other doesn't
- **Neighbor Joining Algorithm:** bottom up cluster algorithm to extract an unrooted tree from a distance matrix

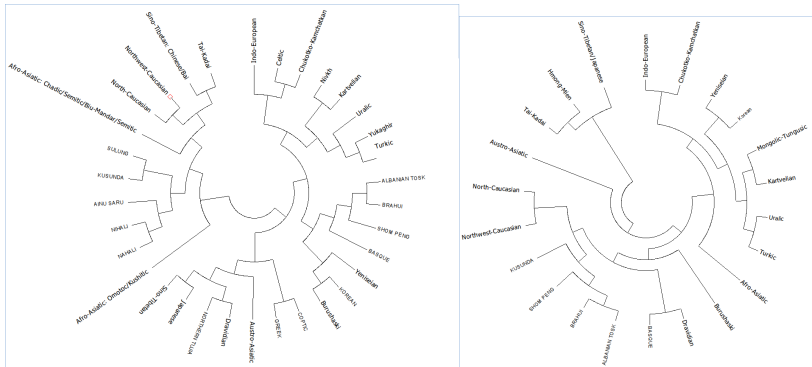# Comparison: language classification

- **Experiment:**
  - compute NJ-tree for all languages in ASJP based on LDND and NWPV distances
  - extract sub-tree of Ethnologue tree for the languages in ASJP
  - compute Robinson-Foulds metric between Ethnologue tree and each of the two NJ trees
- **Outcome:**
  - LDND: **5,522** (4,550 false positives, 972 false negatives)
  - NWPV: **5,476** (4,527 false positives, 949 false negatives)

# Qualitative comparison

- NJ trees for the languages of Eurasia (left: LDND; right: NWPV)

# Qualitative comparison

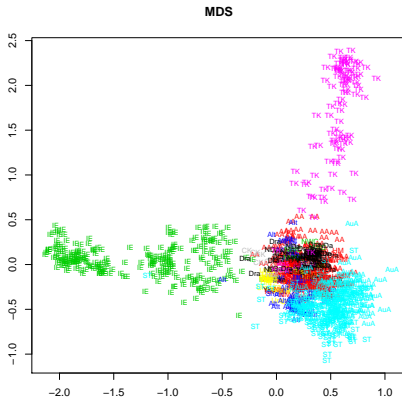- CLANS visualization for the languages of Eurasia (LDND left, NWPV right)

# CLANS and dimensionality reduction

- CLANS performs a kind of (non-deterministic) dimensionality reduction
- How does this relate to more established methods?
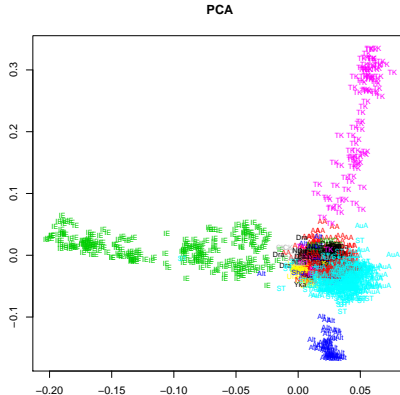
# CLANS vs. Multi-Dimensional Scaling

- MDS applied to NWPV-matrix of the Eurasian languages

# CLANS vs. Principal Component Analysis

- PCA applied to NWPV-matrix of the Eurasian languages



**PCA**

# CLANS and dimensionality reduction

- language families massively vary in size
- MDS and PCA only provide information about the largest families
- CLANS is sensitive to local patterns

# Conclusion

- weighted alignment improves results of lexico-statistical language classification
- more powerful methods from bioinformatics (such as progressive multiple alignment) are likely to lead to further improvement
- CLANS is a useful exploratory tool

Frickey, T. and A. N. Lupas (2004). Clans: a java application for visualizing protein families based on pairwise similarity. *Bioinformatics*, **20**(18):3702–3704.