# Further evidence for punctuated language evolution
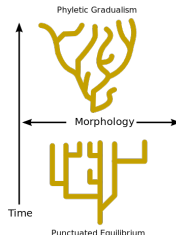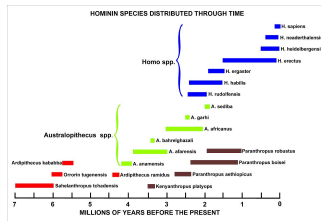
Gerhard Jäger

Tübingen, March 29, 2017

# Punctuated equilibrium



**Gould and Eldredge (1977):**

- surprising lack of **intermediate stages** in fossil record
- possible explanation:
  - evolutionary change occurs primarily during speciation phases
  - when a species is in equilibrium, it neither changes nor speciates

# Punctuated equilibrium

- Possible causal mechanism
  - large population sizes *stabilizes* species
  - mutations, even beneficial ones, rarely reach fixation
  - speciation leads to small populations (bottlenecks) $\rightarrow$ accelerated evolution
- both the existence of the phenomenon and the causal explanation are still contentious in biology



HOMININ SPECIES DISTRIBUTED THROUGH TIME



Phyletic Gradualism

Morphology

Time

Punctuated Equilibrium

# Punctuated language evolution

- Dixon (1997):
  - same logic applies to language change as well
  - rapid tree-like diversification (as in history of IE languages) are the exception in human history
  - Australia prior to European conquest, with an equilibrium between diversification and contact — and hence no tree-like structure — are the rule
- rejected by most historical linguists, especially by experts on Australian languages

# Quantitative approaches 1

- Pagel et al. (2006)
  - amount of evolutionary change is reflected in path lengths of phylogenetic tree (without molecular clock)
  - if punctuational hypothesis is true, there should be positive correlation between length of a path (tip to root) and number of nodes on that path
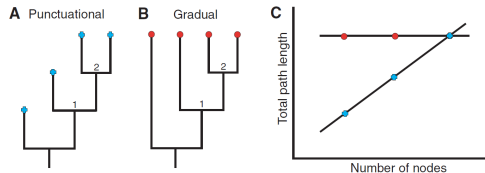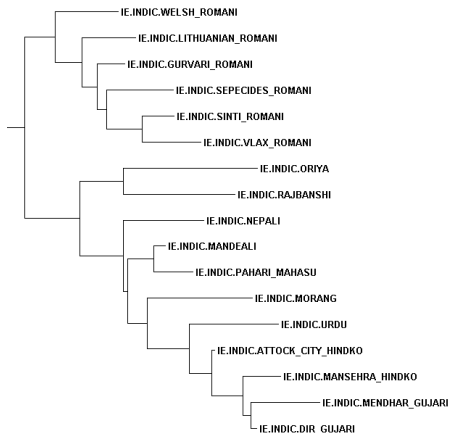  - correlation is tested via phylogenetic regression (PGLS)



**Fig. 1.** Signatures of punctuational and gradual evolution on phylogenetics trees. (**A**) Punctuational evolution presumes a burst of evolution associated with each node of the tree. Path lengths, measured as the sum of branches along a path from the root to the tips of the tree, are proportional to the number of nodes along that path (C). Branches are assumed to be in units of nucleotide substitutions. (**B**) Gradual evolution presumes that change is independent of speciation events. Path lengths do not correlate with the number of nodes along a path (C). (**C**) Punctuational evolution predicts a positive relationship between path length and the number of nodes, whereas gradual evolution does not.

# Quantitative approaches 1

Typical shape of a tree instantiating punctual evolution

# Quantitative approaches 1

- Atkinson et al. (2008):
  - apply this logic to Bayesian trees, based on manual cognacy data, from Austronesian, Bantu, Indo-European and Polynesian
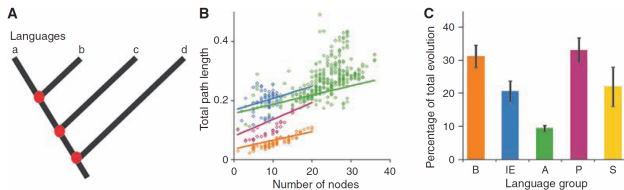


**Fig. 1.** Inferring punctuational language evolution. (**A**) Tree of four languages. If language-splitting events (red nodes) cause bursts of change, the paths from the root to a and b should be longest, followed by c then d (*8*); here, they are all equal. (**B**) Root-to-tip path length plotted against number of nodes along each path for punctuational trees in Bantu (orange), Indo-European (blue), Austronesian (green), and Polynesian (purple). Fitted lines show the relationship between path length and nodes after controlling for phylogeny (*8*). A positive slope is indicative of punctuational evolution. Path lengths for each data set were scaled to account for the number of characters examined. (**C**) Histogram showing the percentage of lexical evolution attributable to punctuational bursts at language-splitting events (mean ± SD) for Bantu (B, orange), Indo-European (IE, blue), Austronesian (A, green), and Polynesian (P, purple) (*8*). For comparison, the percentage of molecular evolution attributable to punctuational effects in biological species is also shown (S, yellow) (*4*).

# Desiderata

- Results only for small number of large and well-studied language families
- Based on manual cognate judgments $\rightarrow$ possible source of implicit bias
- Addressed in Holman and Wichmann (2016) using a different technical approach
- Next part of this talk:
  - Use Atkinson et al.'s method
  - Applied to phylogenies from 6,000+ ASJP doculects,
  - using automatically obtained characters for phylogenetic inference

# The Automated Similarity Judgment Program

- Collaborative data collection project around Cecil Brown, Eric Holman, Søren Wichmann and others
- covers more about 7,000 languages and dialects
- 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*
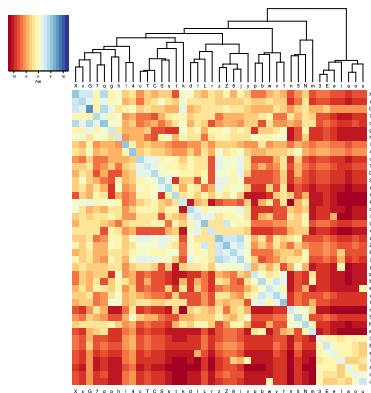
## Automated Similarity Judgment Project

| concept | Latin | English | concept | Latin | English |
|---------|-------|---------|---------|-------|---------|
| *I* | ego | Ei | *nose* | nasus | nos |
| *you* | tu | yu | *tooth* | dens | tu8 |
| *we* | nos | wi | *tongue* | liNgw~E | t3N |
| *one* | unus | w3n | *knee* | genu | ni |
| *two* | duo | tu | *hand* | manus | hEnd |
| *person* | persona, homo | pers3n | *breast* | pektus, mama | brest |
| *fish* | piskis | fiS | *liver* | yekur | liv3r |
| *dog* | kanis | dag | *drink* | bibere | drink |
| *louse* | pedikulus | laus | *see* | widere | si |
| *tree* | arbor | tri | *hear* | audire | hir |
| *leaf* | foly~u* | lif | *die* | mori | dEi |
| *skin* | kutis | skin | *come* | wenire | k3m |
| *blood* | saNgw~is | bl3d | *sun* | sol | s3n |
| *bone* | os | bon | *star* | stela | star |
| *horn* | kornu | horn | *water* | akw~a | wat3r |
| *ear* | auris | ir | *stone* | lapis | ston |
| *eye* | okulus | Ei | *fire* | iNnis | fEir |

# PMI string similarity

- *Pointwise Mutual Information* (PMI) between two sound classes $a$ and $b$:

$$\mathrm{PMI}(a, b) \doteq \log \frac{P(a, b \text{ are homologous})}{P(a)P(b)}$$

- automatically trained from ASJP data (Jäger, 2013)
- PMI similarity between two strings: aggregate PMI score for optimal pairwise alignment of those strings
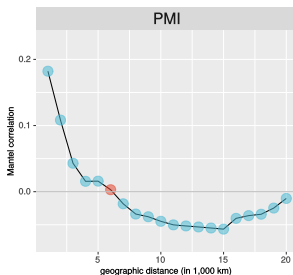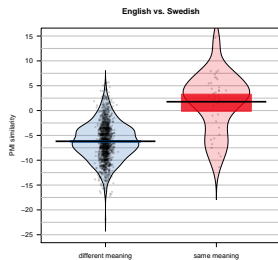
# Calibrated PMI similarity

**English / Swedish**

|      | Ei     | yu     | wi     | w3n    | tu     | fiS     | ... |
|------|--------|--------|--------|--------|--------|---------|-----|
| **yog**  | −**7.77** | 0.75   | −7.68  | −7.90  | −8.57  | −10.50  |     |
| **du**   | −7.62  | **0.33**   | −5.71  | −7.41  | 2.66   | −8.57   |     |
| **vi**   | −2.72  | −2.83  | **4.04**   | −1.34  | −6.45  | 0.70    |     |
| **et**   | −5.47  | −7.87  | −5.47  | −**6.43**  | −1.83  | −4.70   |     |
| **tvo**  | −7.91  | −4.27  | −3.64  | −4.57  | **0.39**   | −6.98   |     |
| **fisk** | −7.45  | −11.2  | −3.07  | −9.97  | −8.66  | **7.58**    |     |

⋮

- values along diagonal give similarity between candidates for cognacy (possibility of meaning change is disregarded)
- values off diagonal provide sample of similarity distribution between non-cognates

# Calibrated PMI similarity



- let $s$ be the PMI-similarity between the English and Swedish word for concept $c$
- **calibrated string similarity**:
  $-\log($probability that random word pairs are more similar than $s$)
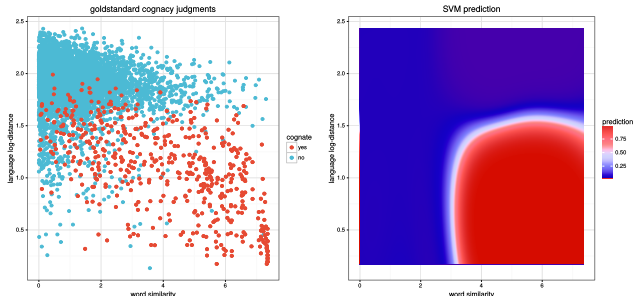- **language similarity:** average word similarity for all concepts

# Cognate clustering

- clustering of ASJP strings into *automatically inferred cognate classes* (Jäger and Sofroniev, 2016; Jäger et al., 2017) (take "cognate" with a grain of salt)
- supervised learning, based on expert cognacy judgments as goldstandard
- sources (only the 40 ASJP concepts were used)

| Dataset | Source | Words | Concepts | Languages | Families | Cognate classes |
|---------|--------|-------|----------|-----------|----------|-----------------|
| ABVD | Greenhill et al. (2008) | 2,306 | 34 | 100 | Austronesian | 409 |
| Afrasian | Militarev (2000) | 770 | 39 | 21 | Afro-Asiatic | 351 |
| Chinese | Běijīng Dàxué (1964) | 422 | 20 | 18 | Sino-Tibetan | 126 |
| Huon | McElhanon (1967) | 441 | 32 | 14 | Trans-New Guinea | 183 |
| IELex | Dunn (2012) | 2,089 | 40 | 52 | Indo-European | 318 |
| Japanese | Hattori (1973) | 387 | 39 | 10 | Japonic | 74 |
| Kadai | Peiros (1998) | 399 | 40 | 12 | Tai-Kadai | 102 |
| Kamasau | Sanders and Sanders (1980) | 270 | 36 | 8 | Torricelli | 59 |
| Mayan | Brown et al. (2008) | 1,113 | 40 | 30 | Mayan | 241 |
| Miao-Yao | Peiros (1998) | 206 | 36 | 6 | Hmong-Mien | 69 |
| Mixe-Zoque | Cysouw et al. (2006) | 355 | 39 | 10 | Mixe-Zoque | 79 |
| Mon-Khmer | Peiros (1998) | 579 | 40 | 16 | Austroasiatic | 232 |
| ObUgrian | Zhivlov (2011) | 769 | 39 | 21 | Uralic | 68 |
| total | | 10,106 | 40 | 318 | 13 | 2,311 |

# Cognate clustering

- calibrated word similarity and language similarity were used as predictors to train a *Support Vector Machine* $\rightarrow$ probability of being cognate for each pair of synonymous ASJP entries
- *Label Propagation* (Raghavan et al., 2007) for clustering
- $0.84$ B-cubed F-score with cross-validation on goldstandard data

# Cognate clustering

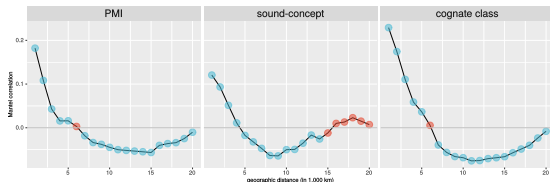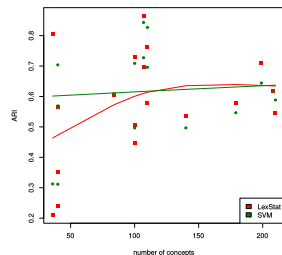| concept | doculect | glot_fam | transcription |
|---------|----------|----------|---------------|
| eye | DORASQUE | Chibchan | oko |
| eye | NORTHERN_LOW_SAXON | Indo-European | ok |
| eye | NORTH_FRISIAN_AMRUM | Indo-European | uk |
| eye | STELLINGWERFS | Indo-European | ok |
| eye | ASSAMESE | Indo-European | soku |
| eye | CHAKMA_UnnamedInSource | Indo-European | sog |
| eye | DALMATIAN | Indo-European | vaklo |
| eye | FRIULIAN | Indo-European | voli |
| eye | ITALIAN | Indo-European | okkyo |
| eye | ITALIAN_GROSSETO_TUSCAN | Indo-European | okyo |
| eye | JUDEO_ESPAGNOL | Indo-European | oxo |
| eye | LATIN | Indo-European | okulus |
| eye | NEAPOLITAN_CALABRESE | Indo-European | woky3 |
| eye | ROMANIAN_2 | Indo-European | oky |
| eye | ROMANIAN_MEGLENO | Indo-European | wokLu |
| eye | SARDINIAN | Indo-European | ogu |
| eye | SARDINIAN_CAMPIDANESE | Indo-European | oxu |
| eye | SARDINIAN_LOGUDARESE | Indo-European | okru |
| eye | SICILIAN_UnnamedInSource | Indo-European | okiu |
| eye | SPANISH | Indo-European | oho |
| eye | TURIA_AROMANIAN | Indo-European | okLu |
| eye | VLACH | Indo-European | okklu |
| eye | BELARUSIAN | Indo-European | voka |
| eye | BOSNIAN | Indo-European | oko |
| eye | BULGARIAN | Indo-European | oko |
| eye | CROATIAN | Indo-European | oko |
| eye | CZECH | Indo-European | oko |
| eye | KASHUBIAN | Indo-European | wokwo |
| eye | LOWER_SORBIAN | Indo-European | voko |
| eye | LOWER_SORBIAN_2 | Indo-European | woko |
| eye | MACEDONIAN | Indo-European | oko |
| eye | OLD_CHURCH_SLAVONIC | Indo-European | oko |
| eye | POLISH | Indo-European | oko |
| eye | SERBOCROATIAN | Indo-European | oko |
| eye | SLOVAK | Indo-European | oko |
| eye | SLOVENIAN | Indo-European | oko |
| eye | UKRAINIAN | Indo-European | oko |
| eye | UPPER_SORBIAN | Indo-European | voCko |
| eye | UPPER_SORBIAN_2 | Indo-European | voko |
| eye | BAINOUK_GUNYAAMOLO | Atlantic-Congo | g3li |
| eye | USINO | Nuclear_Trans_New_Guinea | ogo |

# ASJP word lists → character matrix

**1** **Automatically inferred cognate classes**

- each cluster $cc$ defines one character
- doculect $l$ has value 1 if its word list contains an element of $cc$, undefined if the slot of the corresponding concept is undefined, and 0 else
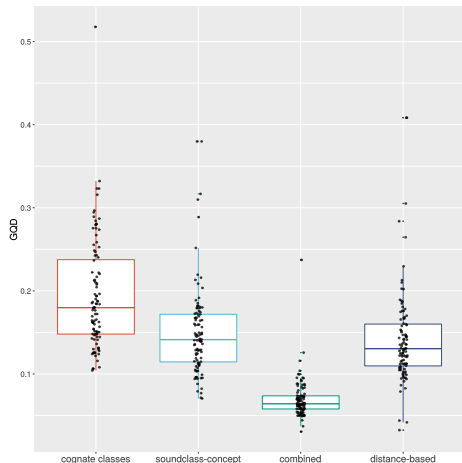
**2** **Soundclass-concept characters**

- each combination $(c, s)$ of an ASJP concept $c$ and an ASJP sound class $s$ is a character
- doculect $l$ has value 1 if one of its entries for $c$ contains $s$, 0 if not, and undefined if there is no entry for $c$

# ASJP word lists → character matrix

- validation
  - correlation with geographic distance
  - phylogenetic inference (Maximum Likelihood) + comparison to Glottolog expert tree on 100 random sample of ASJP doculects, containing between 20 and 400 doculects
- partitioned character-based inference seems to work best

# The node density artifact

- character-based phylogenetic inference tends to under-estimate branch lengths (measured in expected number of mutations) (Webster et al., 2003; Venditti et al., 2006)
- intuitive reason:
  - multiple changes of the same character remain undetected
- effect is stronger on long than on short branches
- leads to spurious correlation between estimated root-to-tip distance and number of intervening nodes
- most pronounced for Maximum Parsimony, but Bayesian and Maximum-Likelihood inference also affected

# Detecting the node density artifact

Webster et al. (2003); Venditti et al. (2006):

- if the relationship is due to the node density artifact, it has a characteristic curved tendency
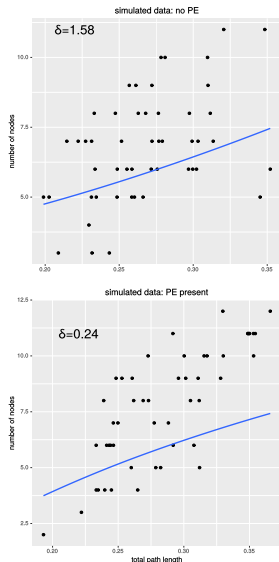
- model:

$$
\begin{aligned}
y &= a + bx^{\delta} + \epsilon \\
x &: \quad \text{total path length} \\
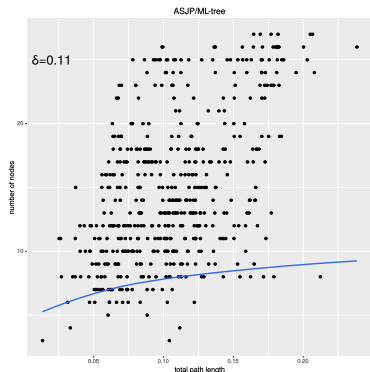y &: \quad \text{number of nodes}
\end{aligned}
$$

- $a, b, \delta$ fitted via phylogenetic generalized least square

- if $b > 0$ and $\delta > 1$, the correlation is due to the node density artifact

(simulated data: courtesy of Søren Wichmann)



simulated data: no PE

$\delta=1.58$



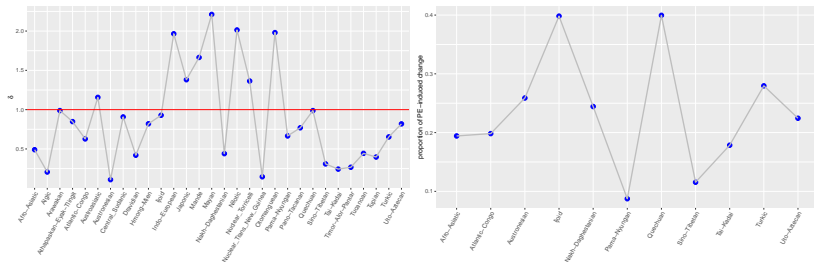simulated data: PE present

$\delta=0.24$

# Application to ASJP data

- 500 randomly selected doculects from ASJP
- Maximum-Likelihood tree/partitioned analysis (cognate-class and sound-concept characters)
- Glottolog expert tree as constraint tree
- results:
  - $\delta = 0.11 \rightarrow$ no node density effect
  - $b > 0$: $p = 1.4 \times 10^{-9}$
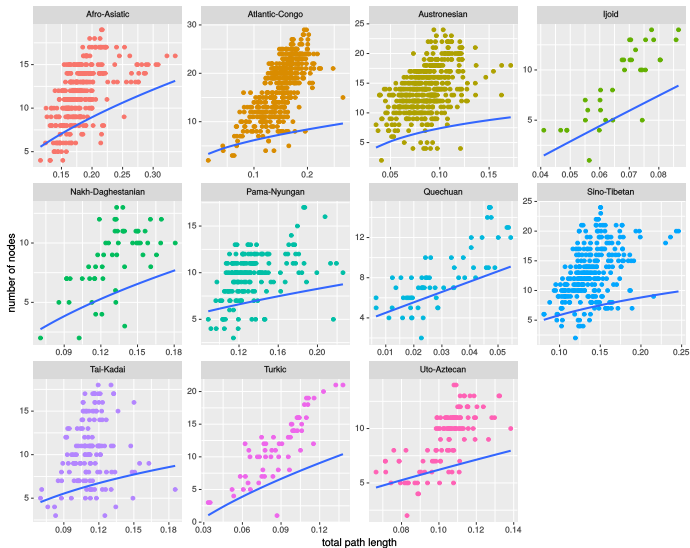  - $20.4\%$ of all change due to punctuation

# Family-wise analysis

- separate model for each Glottolog family with $\geq 30$ doculects in ASJP
- 30 families in total
- if there are more than 500 doculects, 500 doculects randomly selected
- ML tree with Glottolog expert tree constraints
- $\delta < 1$ for 22 families
- significant $b > 1$ ($\alpha = 0.05$, corrected for multiple testing via Holm-Bonferroni method): 11 families
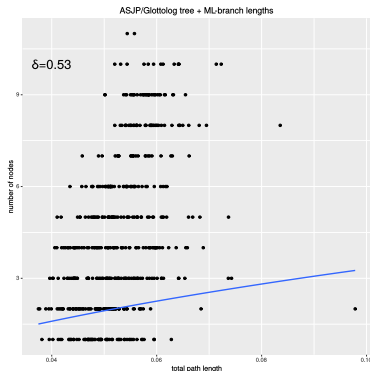
# Family-wise analysis

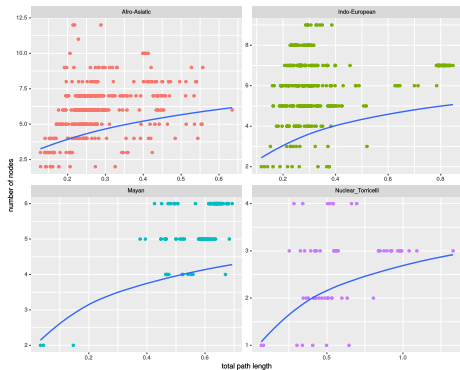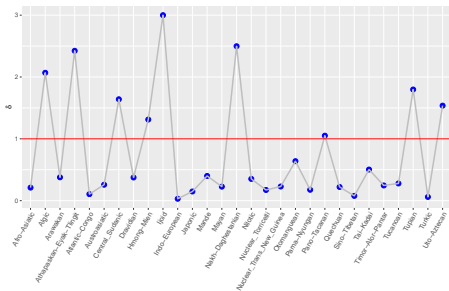| | $\delta$ | $b$ | $p$ | percentage |
|---|---|---|---|---|
| Austronesian | 0.11 | 0.0034 | 0.0 | 25.9 |
| Afro-Asiatic | 0.49 | 0.0042 | 1.1-14 | 19.4 |
| Atlantic-Congo | 0.63 | 0.0035 | 9.1-12 | 19.8 |
| Sino-Tibetan | 0.31 | 0.0024 | 2.0-06 | 11.6 |
| Ijoid | 0.93 | 0.0029 | 1.5-05 | 39.8 |
| Turkic | 0.65 | 0.0031 | 1.7-05 | 28.0 |
| Uto-Aztecan | 0.82 | 0.0027 | 1.3-04 | 22.4 |
| Tai-Kadai | 0.24 | 0.0025 | 2.1-04 | 17.9 |
| Nakh-Daghestanian | 0.44 | 0.0048 | 2.3-04 | 24.4 |
| Quechuan | 0.99 | 0.0018 | 4.9-04 | 40.0 |
| Pama-Nyungan | 0.67 | 0.0023 | 2.9-03 | 8.7 |
| Nuclear_Trans_New_Guinea | 0.15 | 0.0019 | 1.7-02 | 7.0 |
| Central_Sudanic | 0.91 | 0.0041 | 1.7-02 | 22.6 |
| Timor-Alor-Pantar | 0.27 | 0.0030 | 4.9-02 | 17.7 |
| Arawakan | 0.99 | 0.0016 | 7.9-02 | 6.7 |
| Athapaskan-Eyak-Tlingit | 0.85 | 0.0042 | 1.1-01 | 18.9 |
| Dravidian | 0.42 | 0.0025 | 1.1-01 | 15.2 |
| Hmong-Mien | 0.82 | 0.0030 | 1.4-01 | 17.4 |
| Tupian | 0.40 | 0.0010 | 2.5-01 | 4.9 |
| Pano-Tacanan | 0.77 | 0.0011 | 5.2-01 | 5.5 |
| Tucanoan | 0.44 | 0.0003 | 9.1-01 | 2.1 |
| Algic | 0.21 | -0.0001 | 9.7-01 | -0.3 |
| Indo-European | 1.97 | 0.0017 | 3.8-12 | 12.7 |
| Mayan | 2.21 | 0.0024 | 5.2-08 | 28.8 |
| Otomanguean | 1.98 | 0.0042 | 2.7-05 | 20.2 |
| Nilotic | 2.01 | 0.0032 | 3.0-04 | 20.6 |
| Austroasiatic | 1.16 | 0.0018 | 4.5-04 | 10.7 |
| Japonic | 1.38 | 0.0050 | 3.2-03 | 47.6 |
| Mande | 1.66 | 0.0019 | 1.8-02 | 12.0 |
| Nuclear_Torricelli | 1.36 | 0.0020 | 8.7-02 | 7.8 |

# Family-wise analysis

# Polytomies

- phylogenetic inference always produces binary-branching trees
- language diversification possibly involve genuine polytomies $\to$ ML-tree might overestimate number of nodes
- second test:
    - use topology of Glottolog expert tree
    - ML-optimization of branch lengths
- $\delta = 0.53 \to$ no node density effect
- $b > 0$: $p = 3 \times 10^{-5}$



ASJP/Glottolog tree + ML-branch lengths

# Polytomies — family-wise

- $\delta < 1$ for 20 families
- significant $b > 0$ only for four families

# Polytomies — family-wise

| | $\delta$ | $b$ | $p$ |
|---|---|---|---|
| Afro-Asiatic | 0.21 | 0.0108 | 5.7-07 |
| Indo-European | 0.03 | 0.0134 | 8.8-06 |
| Mayan | 0.23 | 0.0560 | 1.4-05 |
| Nuclear_Torricelli | 0.17 | 0.1100 | 1.9-04 |
| Tai-Kadai | 0.50 | 0.0144 | 3.3-03 |
| Atlantic-Congo | 0.11 | 0.0036 | 7.6-03 |
| Pama-Nyungan | 0.18 | 0.0016 | 2.1-02 |
| Nuclear_Trans_New_Guinea | 0.23 | 0.0036 | 2.8-02 |
| Austroasiatic | 0.26 | 0.0097 | 4.4-02 |
| Sino-Tibetan | 0.08 | 0.0047 | 6.2-02 |
| Quechuan | 0.22 | 0.0243 | 1.4-01 |
| Japonic | 0.15 | 0.0339 | 1.7-01 |
| Arawakan | 0.38 | 0.0111 | 2.9-01 |
| Otomanguean | 0.64 | 0.0066 | 3.1-01 |
| Turkic | 0.06 | 0.0098 | 3.4-01 |
| Tucanoan | 0.28 | -0.0235 | 4.1-01 |
| Nilotic | 0.35 | 0.0072 | 5.9-01 |
| Dravidian | 0.37 | 0.0056 | 6.7-01 |
| Mande | 0.40 | 0.0015 | 6.7-01 |
| Timor-Alor-Pantar | 0.25 | 0.0591 | 7.9-01 |
| Nakh-Daghestanian | 2.50 | 0.1239 | 4.4-05 |
| Uto-Aztecan | 1.54 | 0.0149 | 3.8-03 |
| Athapaskan-Eyak-Tlingit | 2.42 | -0.0406 | 1.2-02 |
| Central_Sudanic | 1.64 | 0.0295 | 2.2-02 |
| Hmong-Mien | 1.31 | 0.0602 | 8.6-02 |
| Tupian | 1.80 | -0.0045 | 1.2-01 |
| Pano-Tacanan | 1.05 | 0.0474 | 1.3-01 |
| Algic | 2.07 | -0.0079 | 1.8-01 |
| Ijoid | 3.00 | 0.2448 | 3.3-01 |

# Quantative approaches 2

- Holman and Wichmann (2016):
  - alternative, non-parametric approach
  - various data sources (cognacy data, ASJP/LDND)

## basic idea (cf. graphics to the right)

- B is larger than A

- therefore members of B underwent, on average, more diversification events than members of A

- A and B have same distance (in years) from members of outgroup

- if punctuational hypothesis is true, average distance in amount of change between A and outgroup should be smaller than between B and outgroup
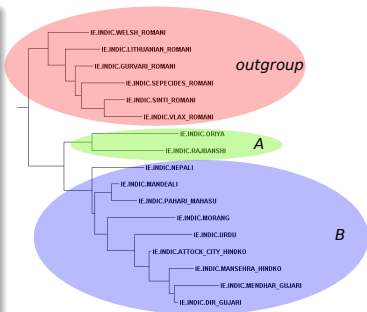
# Quantative approaches 2

- Holman and Wichmann (2016):
  - alternative, non-parametric approach
  - various data sources (cognacy data, ASJP/LDND)

## basic idea (cf. graphics to the right)

- if d(A,outgroup)<d(B,outgroup), this triplet is counted as evidence for puncuation, and vice versa

- Null hypothesis (no punctuational change): probability of a triplet to be evidence for punctuation is $50\%$

- tested on maximal collection of *independent usable triplets* in a given phylogeny

# Independent usable triplets

- a triplet is a local subtree of the shape *((A,B),outgroup)*
- a triplet is *usable* if $A$ and $B$ contain an unequal number of tips
- a group of *independent usable triplets* is a group of usable triplets where no element is contained in another element

# Independent usable triplets

## Algorithm for maximal collectio of independent usable triplets

- the *ratio* of a branching node is the size of its largest daughter, divided by the size of its smallest daughter
- traverse through the tree tip-to-root
- let $N$ the the current node and $C(N)$ be $N$'s *candidate set*
- If $N$ is a tip, $C(N) = \emptyset$, else
- $CC(N) \doteq \bigcup \{C(N') | N' \text{ is a daughter of } N\}$
- If $N$'s ratio $> 1$ and
    - $CC(N) = \emptyset$ or
    - $CC(N) = \{x\}$ and $N$'s ratio $> x$'s ratio:
- $C(N) = \{N\}$, else
- $C(N) = CC(N)$
- $C(\text{root})$ is the maximal collection of independent usable triplets

# Results

- Holman and Wichmann (2016): significant evidence for punctional evolution if amount of change is operationalized as
  - path lengths in an automatically inferred phylogeny
  - difference in cognate class inventory
- no significant effect for LDND

# My results

- "average distance" defined as **median** distance
- PMI-distance and Levenshtein-derived distances
- 5,522 ASJP17-doculects (no ancient, pidgins, creoles, artificial and reconstructed languages)
- Glottolog and ML-inferred topology
- binomial test

| Topology | Distance | positive triplets | total triplets | proportion | $p$-value |
|----------|----------|-------------------|----------------|------------|-----------|
| Glottolog | PMI | 284 | 512 | 0.55 | $0.007 **$ |
| Glottolog | LDPV | 273 | 512 | 0.53 | 0.072 |
| ML | PMI | 513 | 960 | 0.53 | $0.018*$ |
| ML | LDPV | 516 | 960 | 0.54 | $0.011 **$ |

- not a single individual family give significant evidence for PE (if we control for multiple testing)

# Conclusion

- automatically extracted characters give similar results to manually collected ones
- further evidence for punctuated language evolution across the world
- open questions:
  - How much of the effect is due to node denisity artifact?
  - What impact have incomplete lineage sorting, borrowing, hard polytomies on branch length estimation?
  - Is punctuated evolution confined to lexical change?

Quentin D. Atkinson, Andrew Meade, Chris Venditti, Simon J. Greenhill, and Mark Pagel. Languages evolve in punctuational bursts. *Science*, 319(5863):588–588, 2008.

Cecil H. Brown, Eric W. Holman, Søren Wichmann, and Viveka Velupillai. Automated classification of the world's languages: A description of the method and preliminary results. *STUF — Language Typology and Universals*, 4:285–308, 2008.

Běijīng Dàxué. *Hànyǔ fāngyán cíhuì* [Chinese dialect vocabularies]. Wénzì Gǎigé, 1964.

Michael Cysouw, Søren Wichmann, and David Kamholz. A critique of the separation base method for genealogical subgrouping. *Journal of Quantitative Linguistics*, 13(2-3):225–264, 2006.

Robert M. W. Dixon. *The rise and fall of languages*. Cambridge University Press, Cambridge, UK, 1997.

Michael Dunn. Indo-European lexical cognacy database (IELex). URL: http://ielex.mpi.nl/, 2012.

Stephen J. Gould and Niles Eldredge. Punctuated equilibria: the tempo and mode of evolution reconsidered. *Paleobiology*, 3(2):115–151, 1977.

Simon J. Greenhill, Robert Blust, and Russell D. Gray. The Austronesian Basic Vocabulary Database: From bioinformatics to lexomics. *Evolutionary Bioinformatics*, 4:271–283, 2008.

Shirō Hattori. Japanese dialects. In Henry M. Hoenigswald and Robert H. Langacre, editors, *Diachronic, areal and typological linguistics*, pages 368–400. Mouton, The Hague and Paris, 1973.

Eric W. Holman and Søren Wichmann. New evidence from linguistic phylogenetics identifies limits to punctuational change. *Systematic Biology*, 2016. doi: 10.1093/sysbio/syw106.

Gerhard Jäger. Phylogenetic inference from word lists using weighted alignment with empirically determined weights. *Language Dynamics and Change*, 3(2):245–291, 2013.

Gerhard Jäger and Pavel Sofroniev. Automatic cognate classification with a Support Vector Machine. In Stefanie Dipper, Friedrich Neubarth, and Heike Zinsmeister, editors, *Proceedings of the 13th Conference on Natural Language Processing*, volume 16 of *Bochumer Linguistische Arbeitsberichte*, pages 128–134. Ruhr Universität Bochum, 2016.

Gerhard Jäger, Johann-Mattis List, and Pavel Sofroniev. Using support vector machines and state-of-the-art algorithms for phonetic alignment to identify cognates in multi-lingual wordlists. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*. ACL, 2017. to appear.

Kenneth A. McElhanon. Preliminary observations on Huon Peninsula languages. *Oceanic Linguistics*, 6(1):1–45, 1967. ISSN 00298115, 15279421. URL http://www.jstor.org/stable/3622923.

A IU Militarev. *Towards the chronology of Afrasian (Afroasiatic) and its daughter families*. McDonald Institute for Archaelogical Research, Cambridge, 2000.

Mark Pagel, Chris Venditti, and Andrew Meade. Large punctuational contribution of speciation to evolutionary divergence at the molecular level. *Science*, 314(5796):119–121, 2006.

Ilia Peiros. Comparative linguistics in Southeast Asia. *Pacific Linguistics*, 142, 1998.

Usha Nandini Raghavan, Réka Albert, and Soundar Kumara. Near linear

time algorithm to detect community structures in large-scale networks. *Physical Review E*, 76(3):036106, 2007.

Joy Sanders and Arden G Sanders. Dialect survey of the Kamasau language. *Pacific Linguistics. Series A. Occasional Papers*, 56:137, 1980.

Chris Venditti, Andrew Meade, and Mark Pagel. Detecting the node-density artifact in phylogeny reconstruction. *Systematic Biology*, 55(4):637–643, 2006.

Andrea J. Webster, Robert J. H. Payne, and Mark Pagel. Molecular phylogenies link rates of evolution and speciation. *Science*, 301(5632): 478–478, 2003.

Mikhail Zhivlov. Annotated Swadesh wordlists for the Ob-Ugrian group. In George S. Starostin, editor, *The Global Lexicostatistical Database*. RGGU, Moscow, 2011. URL: http://starling.rinet.ru.