# Bivalent verb classes in the languages of Northern Eurasia: genetic and areal factors

Sergey Say
St. Petersburg State University &
Institute for Linguistic Studies, RAS
serjozhka@yahoo.com

# Structure of the talk

- Background and aims
- Data collection
- Distance metrics
- Results
- Conclusions

# Structure of the talk

- Background and aims
- Data collection
- Distance metrics
- Results
- Conclusions

# Background and aims

- Inspiration: wordlist-based typological studies into valency patterns
    - experiencer encoding [Bossong, 1998; Haspelmath 2001]
    - Split-S: A-like vs. P-like vs. G-like [Nichols 2008]
    - Causative~Inchoative alternation and valence orientation [Nedjalkov 1969; Haspelmath 1993; Nichols et al. 2004; WATP]

    and especially:
    - Transitivity hierarchies, cf. Wichmann's [2015] and Haspelmath's [2015] wordlist-based reassessment of Tsunoda's [1981, 1985] hierarchy, & other studies within the ValPaL project [Malchukov & Comrie (eds.) 2015]

# Background and aims

- Typical problems
  - short wordlists (4-70 verbs) ≈ only major patterns
  - *tertium comparationis*, especially if sets of values are pre-established (e.g. agent-like vs. dative-like vs. patient-like experiencer)

- Consequences
  - we know which verbs are most likely to be transitive, but:
  - we don't know much about internal organization of **minor** (non-canonical) valency classes
  - and the ways in which **genetic and areal factors** affect valency class systems

# Background and aims

- Research questions
  - To what extent are valency class systems similar in areally and genetically related languages?
  - How can we identify and **measure** these similarities?
  - What is the depth of genetic effects = **how stable** are valency class systems?
  - What is the **granularity** of areal effects? Cf.:

  The scale of geographical patterning is the size of the areal unit – local, subcontinental, larger than continental, global – within which the geagraphical distribution of a feature displays some clear and describable pattern. For example, ... nominal classes tend to cluster areally and form hotbeds which are generally smaller than continental in size (subcontinental) [Nichols 1992: 185]

# Background and aims

- **Bivalent** verbs
  - Because bivalent verbs are especially prone to show deviant valency behaviour [Bickel et al. 2014] and here, language-internal lexical distributions can be especially complex

- **130** verb meanings
  - Because we need many meanings in order to discern finer signals in the data

- Just one macro-area: **Northern Eurasia**
  - Because this it is possible to rely exclusively on **primary data** (it is not feasible to extract reliable data on as many as 130 verbs from published sources)
  - and still have a relatively dense grid of languages covered

# Background and aims

- It comes at a price
  - convenience sample: I depend on availability of experts and speakers
  - the wordlist can be biased in many ways
  - cross-validation is problematic
  - some meanings can be marginal or non-attested in some languages

# Structure of the talk

- Background and aims
- Data collection
- Distance metrics
- Results
- Conclusions

9

# Data collection: questionnaire

- 130 predicates
- Predicates are provided with contexts in order to make cross-linguistic comparison more accurate

#21     (Peter was crossing the river in a boat)
       'Peter       **reached**     the bank'
       A                              P

#22.     (The wall was covered with fresh paint)
       'Peter       **touched**     the wall' (and got dirty)
       A                              P

# Data collection: questionnaire

- Predicates
  - only predicates that can be expected to be bivalent
  - many predicates that are known to tend to deviate from the transitive prototype

- Translations
  - elicited from native speakers (some exceptions, e.g. Latin)
  - annotated for argument coding devices (flagging and indexing) by language experts
  - variation in argument realization, synonyms etc. are disregarded: one pattern annotated for each predicate

# Data collection: questionnaire

- Valency classes: two verbs belong to the same valency class iff their two arguments are coded by identical devices respectively

Armenian (Eastern)

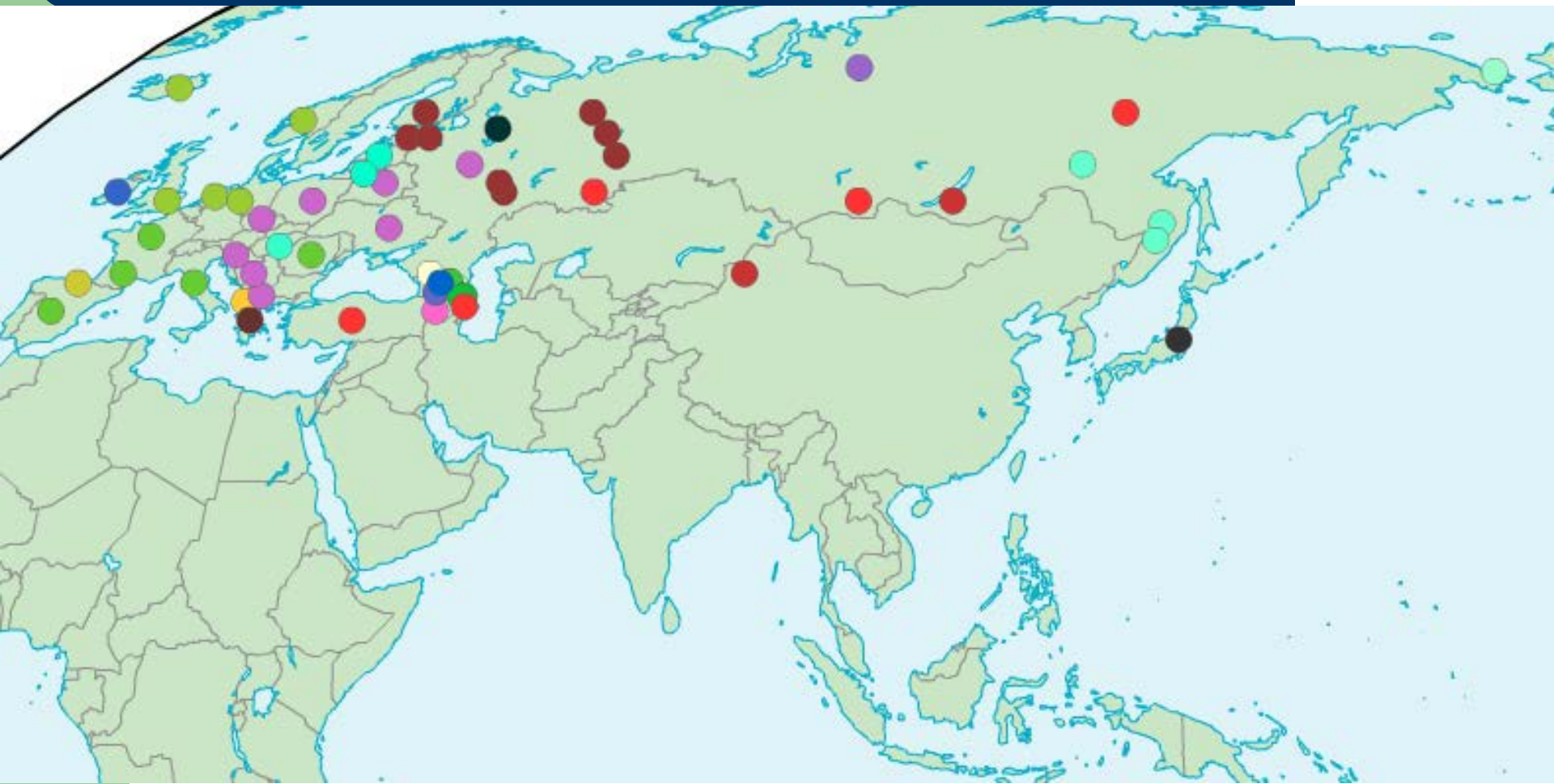| #  | Predicate | Translation | Valency Class |
|----|-----------|-------------|---------------|
| … |  |  |  |
| 21 | reach | *Petros-ə hasav ap'-i-n*<br>Petros[**NOM**]-DEF reach:AOR:**3SG** bank-**DAT**-DEF<br>'P. reached the bank' | **NOM_DAT** |
| 22 | touch | *Petros-ə dipav pat-i-n*<br>Petros[**NOM**]-DEF touch:AOR:**3SG** wall-**DAT**-DEF<br>'Petros touched the wall' | **NOM_DAT** |
| 53 | attack | *Arĵ-ə harjakvec' jknors-i vra*<br>bear[**NOM**]-DEF attack:AOR:**3SG** fisherman-**DAT on**<br>'A bear attacked a fisherman' | **NOM_DATvra** |

# Data collection: questionnaire

- One class in each language was identified as **transitive**
  - in the sense of e.g. [Haspelmath 2011]: the class encompassing 'break'
- The number of valency classses: from 7 (Modern Greek) to 33 (Abaza)

# Data collection: sample

- 57 languages of Northern Eurasia
    - roughly, to the North of 35°N
    - including two extinct languages: Latin and Ancient Greek
    - 9 families (following WALS)
    - 24 genera (following WALS)
- Total datapoints: 6799
    - = 7410 (57 lgs x 130 predicates) − 611 gaps (≈11 per language)

# Data collection: sample
## (coloured by genera)

# Languages and language experts

| Lg | Family | Genus | Expert |
|---|---|---|---|
| Abaza | NorthwestCaucasian | NorthwestC | Peter Arkadiev |
| Albanian | IndoEuropean | Albanian | Varvara Diveeva |
| ArmenianEastern | IndoEuropean | Armenian | Vasilisa Krylova |
| Azerbaijani | Altaic | Turkic | Lejla Kurbanova |
| Bagwalal | NakhDaghestanian | AvarAndicT | Dmitry Gerasimov |
| Bashkir | Altaic | Turkic | Sergey Say |
| Basque | Basque | Basque | Natalia Zaika |
| Belarusian | IndoEuropean | Slavic | Olga Gorickaja |
| Buriat | Altaic | Mongolic | Mikhail Knazev |
| Chukchee | ChukotkoKamchatkan | NorthernCh | Maria Pupynina |
| Czech | IndoEuropean | Slavic | Anastasija Makarova |
| Dutch | IndoEuropean | Germanic | Mikhail Knazev |
| Enets | Uralic | Samoyedic | Maria Ovsjannikova |
| English | IndoEuropean | Germanic | Dmitry Nikolaev |
| ErzyaMordvin | Uralic | Finnic | Ksenia Shagal |
| Estonian | Uralic | Finnic | Irina Külmoja |
| Evenki | Altaic | Tungusic | Nadezhda Bulatova, Elena Perekhvalskaja |
| Finnish | Uralic | Finnic | Ksenia Shagal |
| French | IndoEuropean | Romance | Elena Kordi |
| Gascon | IndoEuropean | Romance | Natalia Zaika |
| Georgian | Kartvelian | Kartvelian | Alexander Rostovtsev- |
| German | IndoEuropean | Germanic | Sandra Birzer |
| GreekAncient | IndoEuropean | Greek | Ildar Ibragimov |
| GreekModern | IndoEuropean | Greek | Ekaterina Zheltova |
| HillMari | Uralic | Finnic | Ksenia Studenikina |
| Hungarian | Uralic | Ugric | Vasilisa Zhigulskaja |
| Icelandic | IndoEuropean | Germanic | Ingunn Hreinberg Indr |
| IngrianFinnish | Uralic | Finnic | Daria Mischenko |

| Lg | Family | Genus | Expert |
|---|---|---|---|
| Ingush | NakhDaghestanian | Nakh | Johanna Nichols |
| Irish | IndoEuropean | Celtic | Dmitry Nikolaev |
| Italian | IndoEuropean | Romance | Anna Alexandrova |
| Japanese | Japanese | Japanese | Yukari Konuma |
| Kalmyk | Altaic | Mongolic | Sergey Say |
| KomiPermyak | Uralic | Finnic | Ekaterina Sergeeva |
| KomiZyrian | Uralic | Finnic | Ekaterina Sergeeva |
| Latin | IndoEuropean | Romance | Inna Popova |
| Latvian | IndoEuropean | Baltic | Natalia Perkova |
| Lezgi | NakhDaghestanian | Lezgic | Ramazan Mamedshax |
| Lithuanian | IndoEuropean | Baltic | Natalia Zaika |
| Macedonian | IndoEuropean | Slavic | Vladimir Fedorov |
| MokshaMordvin | Uralic | Finnic | Maria Kholodilova |
| Nanai | Altaic | Tungusic | Daria Mischenko |
| NorwegianBokmal | IndoEuropean | Germanic | Olga Kuznecova |
| Ossetic | IndoEuropean | Iranian | Arsenij Vydrin |
| Polish | IndoEuropean | Slavic | Georgij Moroz |
| Romanian | IndoEuropean | Romance | Daria Suetina |
| RomaniKalderash | IndoEuropean | Indic | Kirill Kozhanov |
| Russian | IndoEuropean | Slavic | Sergey Say |
| Rutul | NakhDaghestanian | Lezgic | Anastasia Vasilisina, S |
| Serbian | IndoEuropean | Slavic | A.Makarova |
| Slovene | IndoEuropean | Slavic | Andreja Žele, Mladen |
| Spanish | IndoEuropean | Romance | Elena Gorbova |
| Turkish | Altaic | Turkic | Maria Ovsjannikova |
| Tuvan | Altaic | Turkic | Arzhaana Syuryun |
| Udihe | Altaic | Tungusic | Elena Perkhvalskaja |
| Ukrainian | IndoEuropean | Slavic | Natalia Zaika |
| Yakut | Altaic | Turkic | Ajtalina Nogovitsyna |

# Structure of the talk

- Background and aims
- Data collection
- Distance metrics
- Results
- Conclusions

# Distance metrics

For each pair of languages
- Genetic distance
- Areal distance
- Structural distances

# Distance metrics: genetic

- Three levels, based on WALS:
  - 1: same genus
  - 2: same family, different genera
  - 3: different families

  E.g.: DistGenetic (Eastern Armenian, Azerbaijani) = 3

# Distance metrics: geographic

- Calculated as the geographic distance (in kilometers) between the two points associated with individual languages
- Coordinates are taken from WALS
- The distance is calculated using `distCosine()` from the R package `geosphere` [Hijmans 2016]
- NB: this is a very coarse metric for languages spoken over vast areas
- For statistical purposes, the decimal logarithm of the distance is used, e.g.

DistGeo (Eastern Armenian, Azerbaijani) = 277 km

LogDistGeo (Eastern Armenian, Azerbaijani) = 2.44

# Distance metrics: structural

- Structural distances:
  - **DistTrRat**: measures (dis)similarity in transivity prominence
  - **DistTrProf**: measures (dis)similarity in transivity profiles
  - **DistValPat**: measures (dis)similarity between systems of valency classes

# Distance metrics, structural (1)

- Transitivity Ratio (TrRatio): the number of transitive verbs divided by the total number of verbs, cf. [Haspelmath 2015]

E.g. TrRatio (Azerbaijani) = 0.48 (58 transitive verbs / 121 total)

- DistTrRatio is the absolute value of the difference between transitivity prominence in the two languages

DistTrRatio (Azerbaijani, Eastern Armenian) = |0.48 − 0.50| = 0.02

# Distance metrics, structural (2)

- Transitivity profile of a language: sets of +/- transitive verbs

- DistTrProf measures (dis)similarity between "transitivity profiles"

- The relative Hamming distance: the ratio of predicates that are transitive in one language and intransitive in the other

# Distance metrics, structural (2)

| | Eastern Armenian | Azerbaijani |
|---|---|---|
| win | TR | INTR |
| be_afraid | INTR | INTR |
| believe | INTR | INTR |
| see | TR | TR |
| reach | INTR | INTR |
| touch | INTR | INTR |
| forget | INTR | TR |
| wait | TR | TR |
| know | TR | TR |
| avoid | INTR | INTR |
| … | | |

# Distance metrics, structural (2)

|                |   | Azerbaijani | |
|----------------|---|-------------|------|
|                |   | t           | i    |
| Eastern        | t | 53          | 8    |
| Armenian       | i | 5           | 53   |

DistrTrProf (Eastern Armenian, Azerbaijani) = (5+8)/(53+8+5+53) = 13 / 119 = 0.109

- Low DistTrProf entails low DistTrRat, but not vice versa.

# Distance metrics, structural (2)

- DistrTrProf (Eastern Armenian, Azerbaijani) = 0.109
Is this a big difference or a small difference?

- Standardization: **z-scores**
- Mean value of DistTrProf among all pairs = 0.209, and $\sigma$ = 0.07

z (DistrTrProf (Eastern Armenian, Azerbaijani)) = $\frac{0.109 - 0.209}{0.07}$ = $-1.43$

NB: Negative z-scores signal more similarity between languages!

# Distance metrics, structural (3)

- Cross-linguistic identification of minor minor valency classes (cf. "ablative verbs"?, "instrumental verbs"?) is not feasible

- Measuring (dis)similarity in valency class systems is the biggest challenge

- I propose **DistValPat**, a metric based on entropy and MI (mutual information)

- Entropy ≈ the amount of information (conveyed by the valency class assignment)
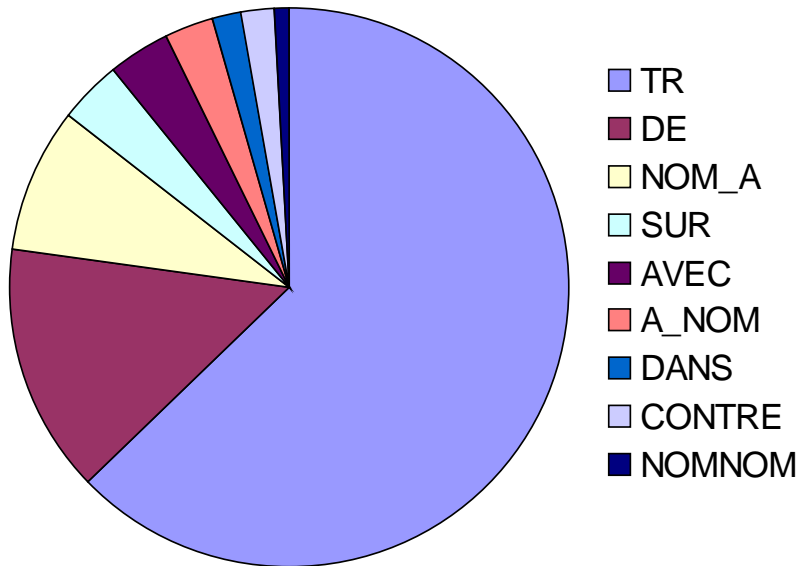
# Distance metrics, structural (3)

$$H(x) = -\sum_{i=1}^{k} p(x_i) \cdot \log(p(x_i))$$

**Hypothetical Language 1:
All verbs belong to the same
class**   $H = 0$

**Hypothetical Language 2:
130 verb classes**

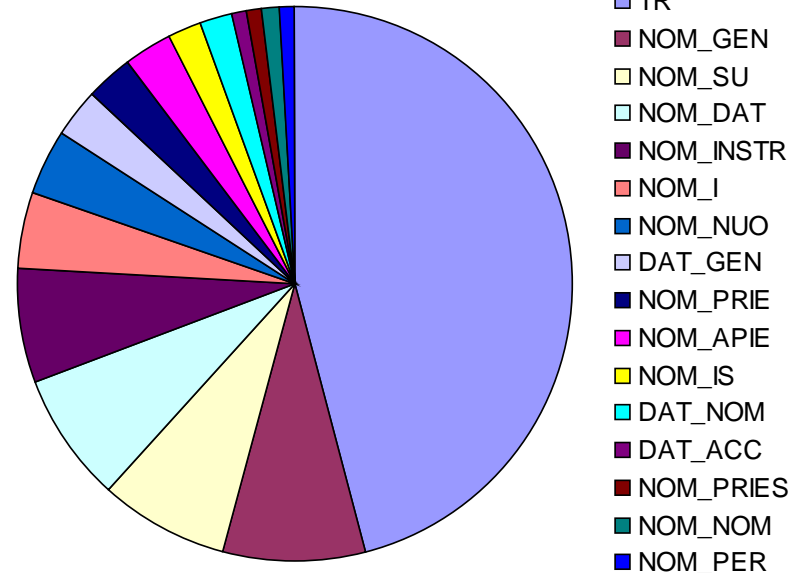$$H = -\log\left(\frac{1}{130}\right) \approx 4{,}87$$

# Entropy

### French



Legend:
- TR
- DE
- NOM_A
- SUR
- AVEC
- A_NOM
- DANS
- CONTRE
- NOMNOM

**H ≈ 1.31**

### Lithuanian



Legend:
- TR
- NOM_GEN
- NOM_SU
- NOM_DAT
- NOM_INSTR
- NOM_I
- NOM_NUO
- DAT_GEN
- NOM_PRIE
- NOM_APIE
- NOM_IS
- DAT_NOM
- DAT_ACC
- NOM_PRIES
- NOM_NOM
- NOM_PER

**H ≈ 2.02**

|  | Armenian | Azerbaijani |  |
|---|---|---|---|
| take | TR | TR | |
| see | TR | TR | |
| influence | NOMvra | NOMDAT | |
| encounter | TR | NOMCOM | |
| enter | NOMNOM | NOMCOM | |
| win | TR | NOMDAT | |
| go_out | NOMABL | NOMABL | |
| drive | TR | TR | |
| bend | TR | TR | |
| tell | NOMDAT | TR | |
| hold | TR | TR | |
| catch_up | NOMDAT | NOMDAT | |
| milk | TR | TR | |
| reach | NOMDAT | NOMDAT | |
| touch | NOMDAT | NOMDAT | |
| fight | NOMhet | NOMCOM | |
| be_friends | NOMhet | NOMCOM | |
| think | NOMmasin | NOMABL | |
| … | | | |
| **H (Entropy)** | **1.658** | **1.462** | |

| | Armenian | Azerbaijani | Joint Distribution |
|---|---|---|---|
| take | TR | TR | TR_TR |
| see | TR | TR | TR_TR |
| influence | NOMvra | NOMDAT | NOMvra_NOMDAT |
| encounter | TR | NOMCOM | TR_NOMCOM |
| enter | NOMNOM | NOMCOM | NOMNOM_NOMCOM |
| win | TR | NOMDAT | TR_NOMDAT |
| go_out | NOMABL | NOMABL | NOMABL_NOMABL |
| drive | TR | TR | TR_TR |
| bend | TR | TR | TR_TR |
| tell | NOMDAT | TR | NOMDAT_TR |
| hold | TR | TR | TR_TR |
| catch_up | NOMDAT | NOMDAT | NOMDAT_NOMDAT |
| milk | TR | TR | TR_TR |
| reach | NOMDAT | NOMDAT | NOMDAT_NOMDAT |
| touch | NOMDAT | NOMDAT | NOMDAT_NOMDAT |
| fight | NOMhet | NOMCOM | NOMhet_NOMCOM |
| be_friends | NOMhet | NOMCOM | NOMhet_NOMCOM |
| think | NOMmasin | NOMABL | NOMmasin_NOMABL |
| … | | | |
| **H (Entropy)** | **1.658** | **1.462** | **2.196** |

# Distance metrics, structural (3)

- MI (Mutual Information) = H (X) + H (Y) − H (X, Y)
- MI (Armenian, Azerbaijani) = 1.658 + 1.462 − 2.196 = 0.924
- Higher MI values reflect higher similarity between valency class systems in the two languages
- MI was calculated using R package `infotheo` [Meyer 2014]

# Distance metrics, structural (3)

- Converting MI into a distance metric

$$\text{DistValPat (L1, L2)} = 1 - \frac{\frac{MI\,(L1,L2)}{H\,(L1)} + \frac{MI(L1,L2)}{H(L2)}}{2}$$

- DistValPat is high if the joint entropy is high relative to individual entropies
- DistValPaI is higher if valency class systems are divergent

33

# Distance metrics, structural (3)

- DistValPal (Armenian, Azerbaijani) = 0.405

- z (DistValPat(Armenian, Azerbaijani)) = $\frac{0.405 - 0.499}{0.096} = -0.97$

- This means that valency class assignment in Armenian and Azerbaijani is rather similar: the distance between the two languages is almost one standard deviation below the mean

# Distance metrics: summary

- Pairs of languages: 1596 = (57*56)/2
- 5 distance metrics for each pair:
  - genetic
  - geographical
  - 3 structural

# Structure of the talk

- Background and aims
- Data collection
- Distance metrics
- Results
- Conclusions

# Results

- All the three structural distance metrics correlate positively with both the genetic and areal distance

# Results

- All the three structural distance metrics correlate positively with both the genetic and areal distance.

    => Expected

- But the devil is in the detail

# Results: transitivity prominence

- Transitivity / intransitivity prominence is primarily an areal phenomenon with **subcontinental** degree of granularity

  – Transitivity peaks are in Central Western Europe and in the Far East

  – Intransitivity peaks are in the Caucasus and in the Eastern Europe

# Results: transitivity prominence

## The ratio of intransitive verbs

# Results: transitivity prominence

- Genera are relatively homogeneous in terms of transitivity prominence: DistTrRatio's are low
- No traceable family-size effects, e.g. both Indo-European and Uralic languages are very diverse

**DistTrRat = Difference in Transitivity Ratio**

# Results: transitivity profiles

- DistTrProf: significant genetic signal not only on the level of individual genera, but also on the family-size level

- Also visible on the MDS (Multidimensional scaling) plot
  – However, Uralic languages are somewhat distorted

43

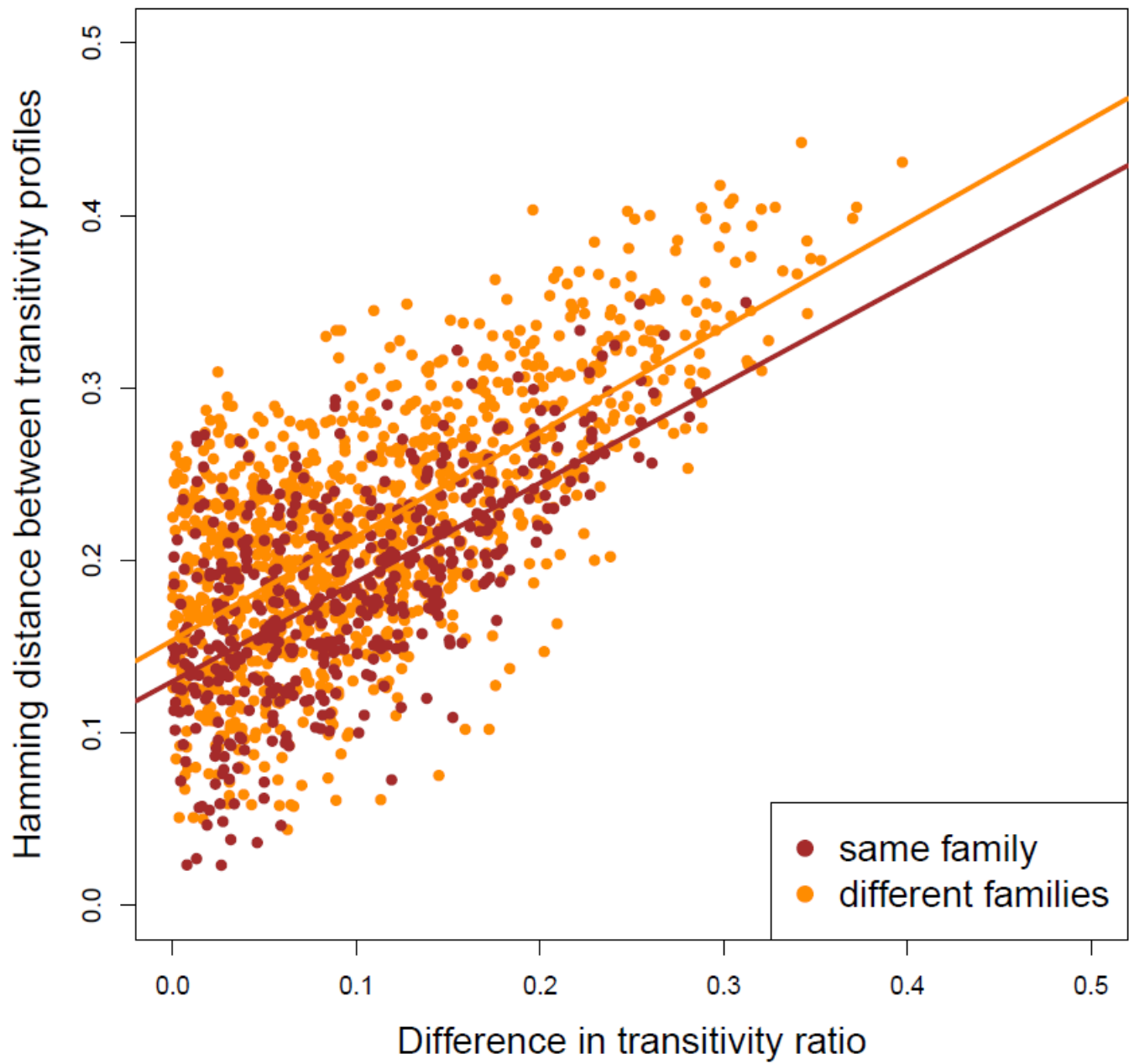**DistTrProf = Hamming distance between transitivity profiles**

**MDS plot for
DistTrProf = Hamming distance between transitivity profiles**

# Results: transitivity profiles

- Given a certain level of DistTrRat, genetically related languages show lower DistTrProf

**DistTrRat, DistTrProf & Genetic relatedness**

Difference in transitivity ratio (x-axis)

Hamming distance between transitivity profiles (y-axis)

Legend:
- same family
- different families

# Results: transitivity profiles

- This would not be expected if the transitivity-prominence scale of verbs were universal

- Probably, verb hierarchies of transitivity prominence are family-specific, e.g.:

  – Experiential predicates ('see', 'know', 'love', 'want') are especially prone to be intransitive in Nakh-Daghestanian

  – Verbs of contact ('follow', 'reach', 'touch', 'kiss', 'attack') are especially prone to be intransitive in Uralic (though not Hungarian)
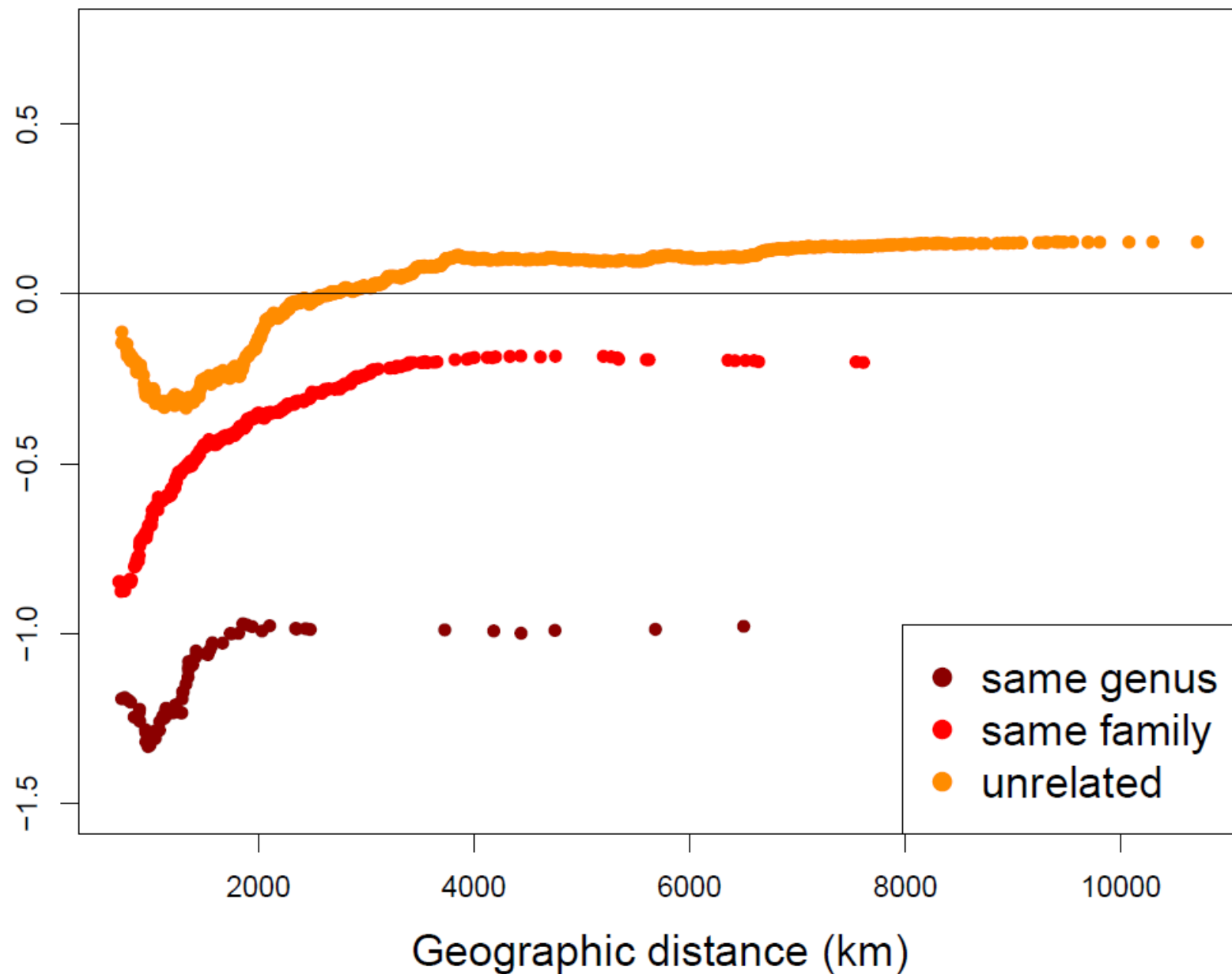
  – etc.

# Results: transitivity profiles

- Next slide: the role of geographic distance
  - X-axis: geographic distance in kilometers
  - Y-axis: mean DistTrProf for pairs of languages spoken closer than N kilometers tp each other (cumulative mean)
  - separately for three levels of genetic distance

  This method is inspired by [Wichmann & Holman 2009: 75 ff.]

**DistTrProf = Hamming distance between transitivity profiles**

Geographic distance (km)

Cumulative average z-scores: DistTrProf

- same genus
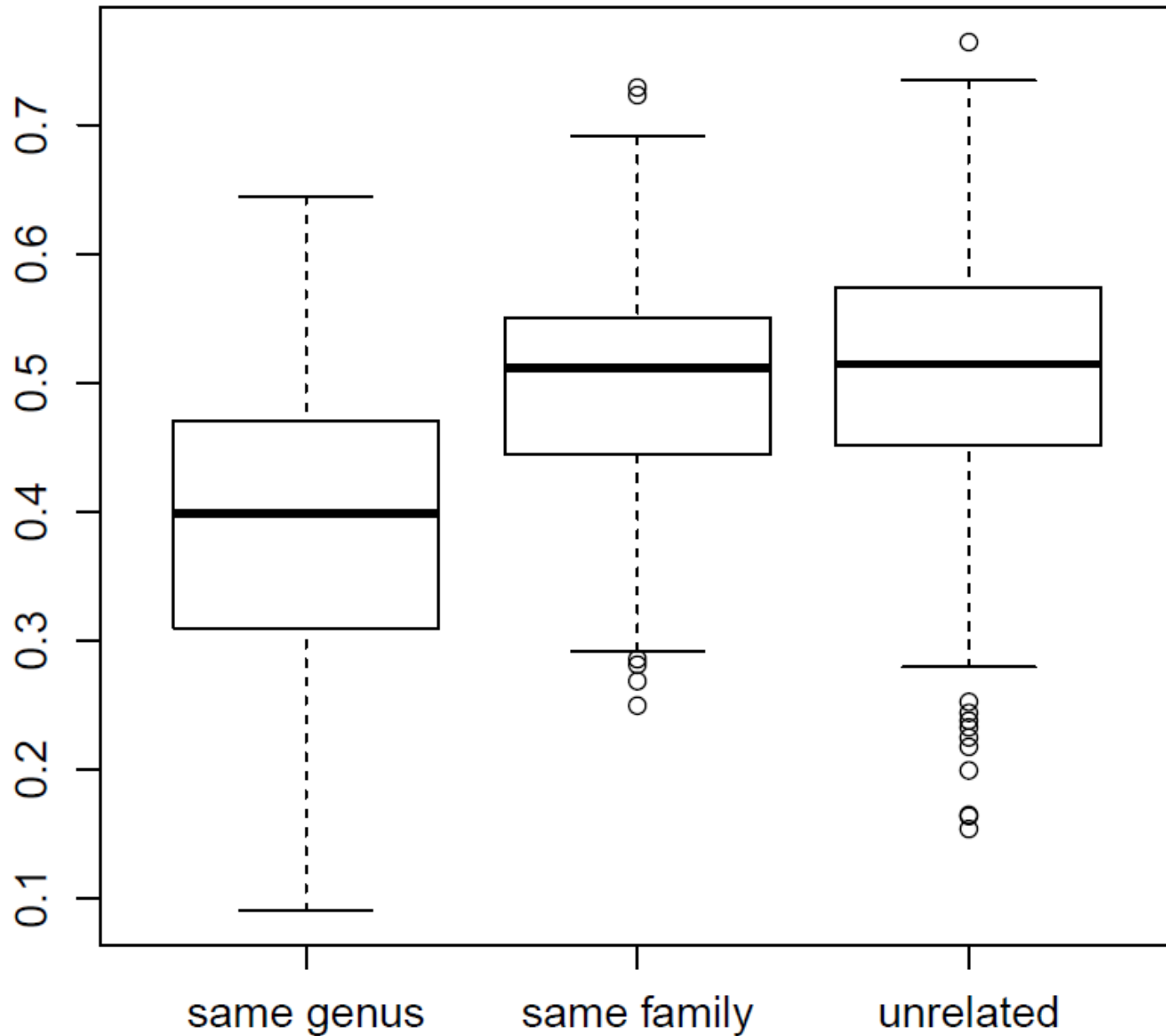- same family
- unrelated

# Results: transitivity profiles

- Robust genetic signal: three curves are very different
- If genetic factor is levelled out, the role of geographic proximity rapidly fades away after ≈2000 km

# Results: valency class systems

- DistValPat displays no family-level effects, only genus-level effects
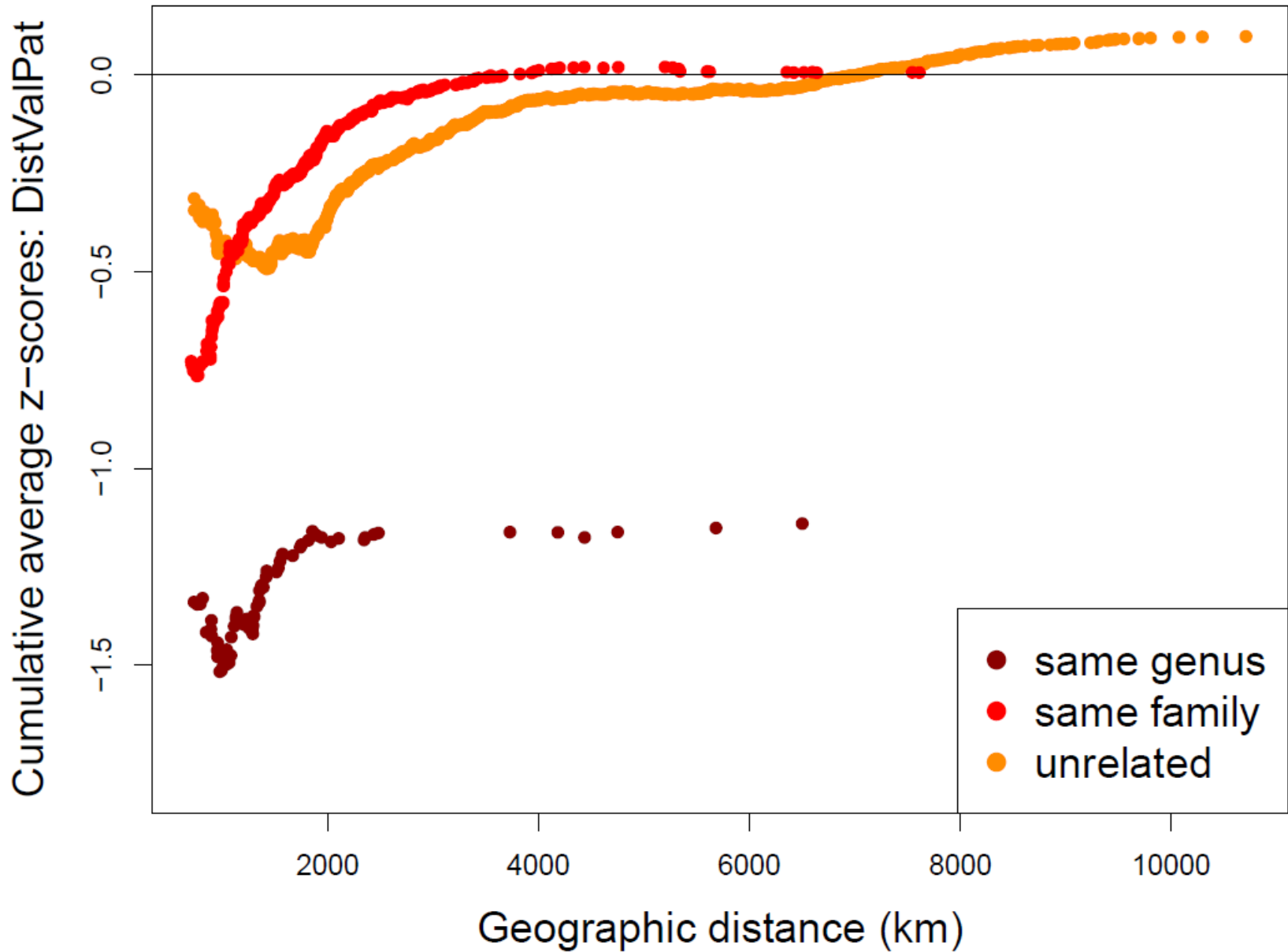
# DistValPat = Entropy-based distance between valency class systems

# Results: valency class systems

- DistValPat: geographical effects (next slide)
  - The curves for languages from same vs. different families show no consistent effect for distances > 1000 km
  - DistValPat shows the strongest areal signal for both genetically related and unrelated languages
  - Caucasus is an exception: many pairs of geographically proximate languages with huge DistValPat; this accounts for the anomaly on the left margin of the orange curve
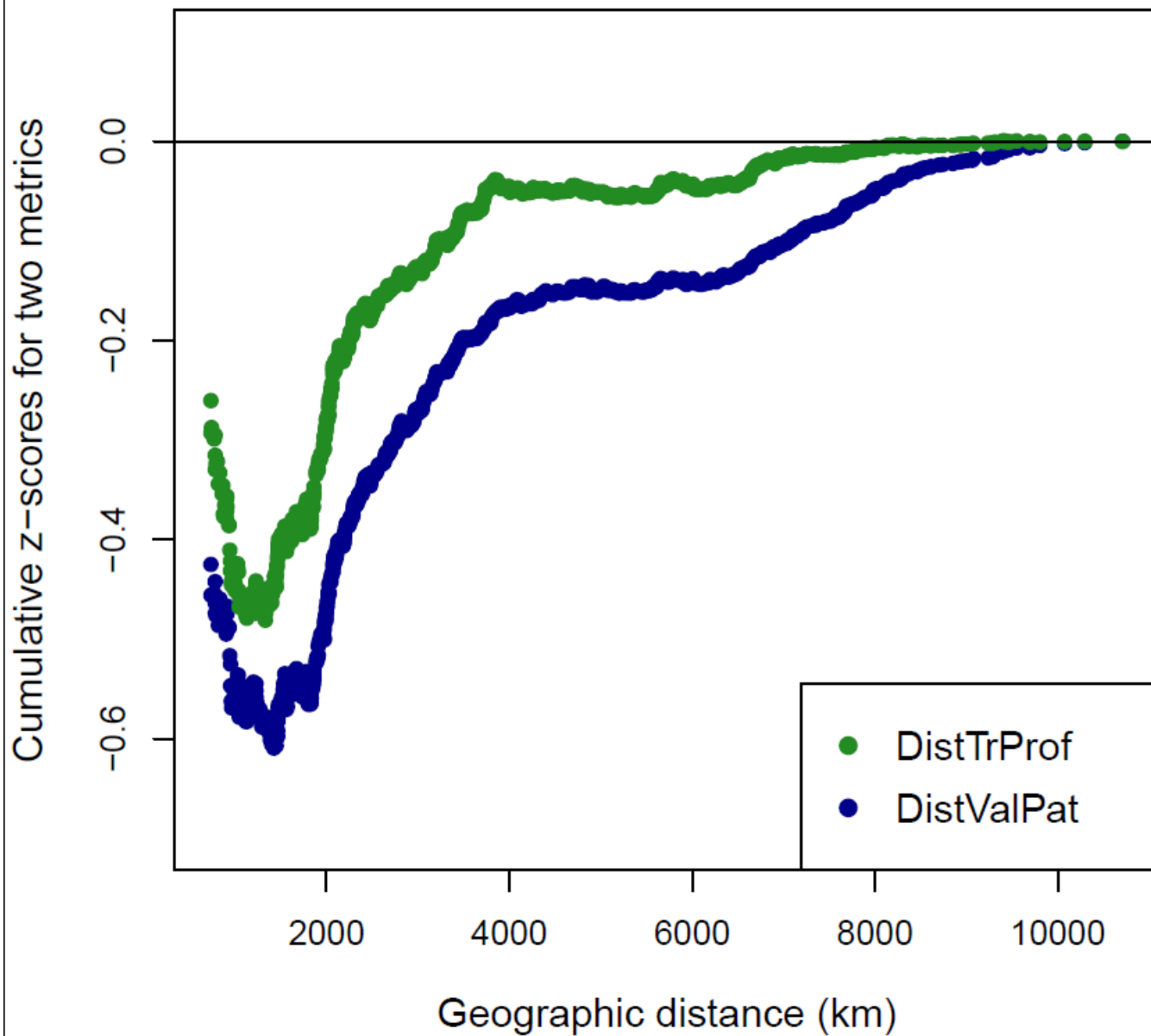
**DistValPat = Entropy−based distance between valency class systems**

Cumulative average z−scores: DistValPat vs. Geographic distance (km)

- same genus
- same family
- unrelated

# Results: valency class systems

- DisValPat displays a stronger and more lasting effect of geographic distance than DistTrProf
  - See the next slide: pairs of genetically related languages are disregarded, z-scores are re-calculated
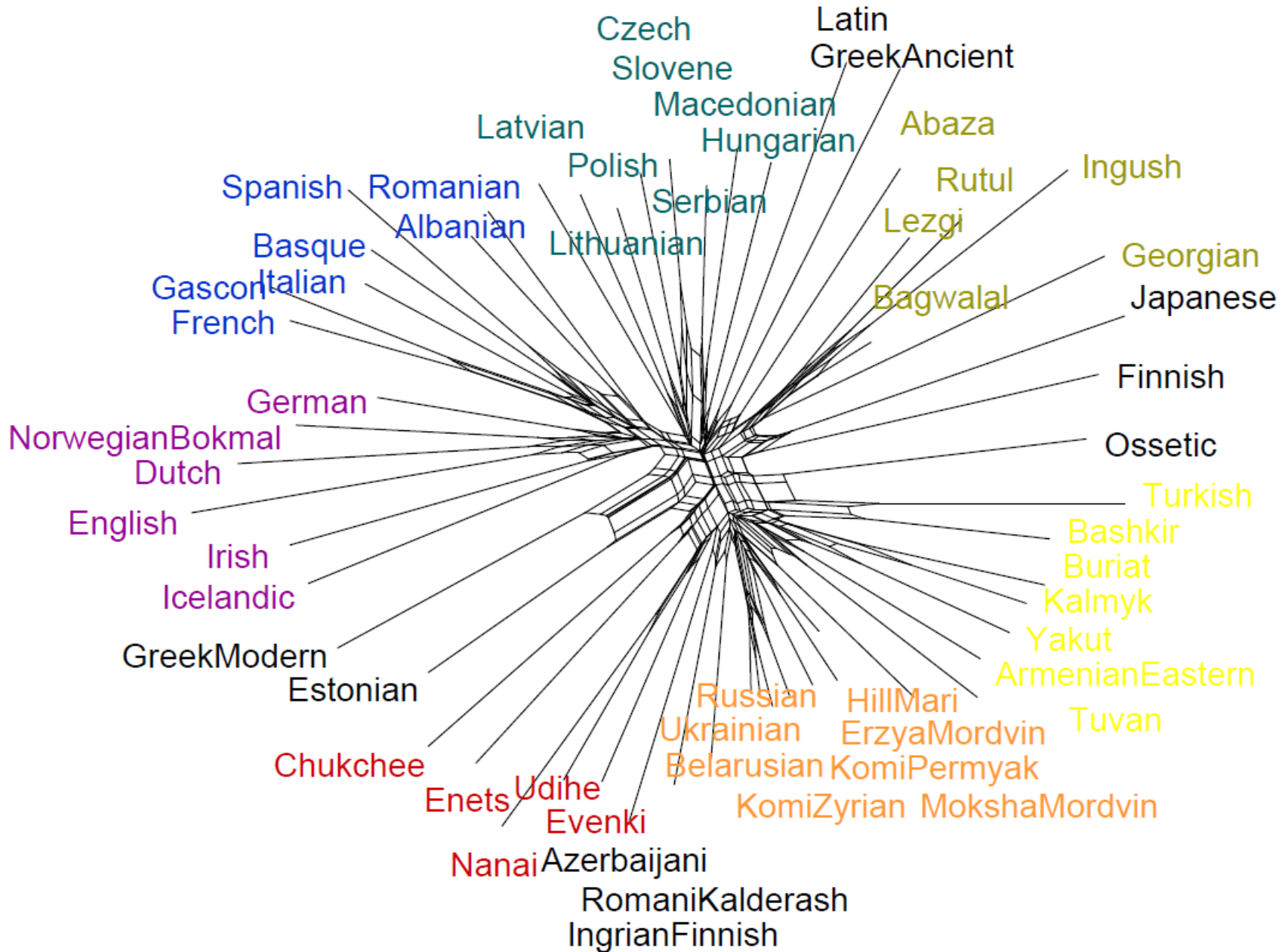
**Two metrics for genetically unrelated languages**

# Results

- Areal effects are clearly visible if the distance matrix is visualized using the NeighborNet algorithm
  - implemented in the SplitsTree software [Huson & Bryant 2006]

# NeighborNet, DistValPat: (dis)similarity in bivalent valency class systems

# Structure of the talk

- Background and aims
- Data collection
- Distance metrics
- Results
- Conclusions

# Conclusions

- Transitivity prominence is an areal phenomenon with subcontinental granularity
- Similarities in transitivity profiles: strong genetic effects, no large-scale geographic effects
- Similarities in valency class organization, including minor classes: no family-level genetic effects, strong areal effects

# Conclusions

- Plausible explanation
  - valency patterns of individual verbs change relatively fast and are easily transferable in language contact
  - languages are relatively stable in terms of those semantic features that are relevant for the assignment of the [+/-] transitivity values to individual verbs
  - and transitivity hierarchies of verb meanings can be family-specific

# Acknowledgements

- Language experts (see above)
- My colleagues from the Institute for linguistic studies, RAS, who participated in research projects supported by Russian Foundation for Humanities (2009-2011; 2011-2013)
- And especially **Maria Ovsjannikova** who created the plots in R

# Thank you!

# Selected references

- Aikhenvald A.Y., Dixon R.M.W., Onishi M. (eds). 2001. Non-canonical marking of subjects and objects. Amsterdam/Philadelphia: John Benjamins.
- Bhaskararao P., Subbarao K.V. (eds). 2004. Non-nominative Subjects. 2 vols. Amsterdam/Philadelphia: John Benjamins.
- Bossong, Georg. 1998. Le marquage de l'expérient dans les langues d'Europe'. In: Feuillet (ed.). Actance et valence dans les langues de l'Europe. Berlin: Mouton de Gruyter. 259–94.
- Dixon, R.M.W., Aikhenvald, A.Y. (eds). 2000: Changing valency: case studies in transitivity. Cambridge.
- Dowty D. 1991. Thematic proto-roles and argument selection. Language 67. 547-619.
- Haspelmath, Martin. 1993. More on the typology of inchoative / causative verb alternations. In: Comrie, Bernard & Maria Polinsky (eds.) *Causatives and Transitivity.* Amsterdam: Benjamins. 87–120.
- Haspelmath, Martin. 2001. Non-canonical marking of core arguments in European languages. In: Aikhenvald et al. eds. 53-83.
- Haspelmath, Martin. 2011. On S, A, P, T, and R as comparative concepts for alignment typology. *Lingustic Typology* 15(3). 535–567.
- Haspelmath, Martin. 2015. Transitivity prominence. In: Malchukov & Comrie (eds.), 131-147.
- Hijmans, Robert J. 2016. geosphere: Spherical Trigonometry. R package version 1.5-5. https://CRAN.R-project.org/package=geosphere

# Selected references

- Hopper, P.J., Thompson, S.A. 1980. Transitivity in grammar and discourse. Language. 1980, 56. (2). P. 251–299.

- Kittilä, S. 2002. Transitivity: towards a comprehensive typology. Turku, 2002.

- Levin, Beth. 1993. English Verb Classes and Alternations. Chicago: University of Chicago Press.

- Malchukov, A. 2006. Transitivity parameters and transitivity alternations: constraining co-variation. In: Case, valency and transitivity, ed. by L. Kulikov, A. Malchukov, P. de Swart. Amsterdam, Philadelphia. P. 175–190.

- Meyer, Patrick E. 2014. infotheo: Information-Theoretic Measures. R package version 1.2.0. https://CRAN.R-project.org/package=infotheo

- Næss, Å. 2007. Prototypical Transitivity. Amsterdam, Philadelphia.

- Nedjalkov, V.P. 1969. Nekotorye verojatnostnye universalii v glagolnom slovoobrazovanii. In: F. Vardul' (ed.). Jazykovye universalii i lingvisticheskaja tipologija. Moscow: Nauka. 106-114.

- Nichols, Johanna. 1992. *Linguistic Diversity in Space and Time*. Chicago: University of Chicago Press.

- Nichols, Johanna. 2008. Why are stative-active languages rare in Eurasia? Typological perspective on split subject marking. In Mark Donohue and Søren Wichmann (eds). The Typology of Semantic Alignment Systems, 121-139. Oxford: Oxford University Press.

- Nichols, Johanna, David A. Peterson & Jonathan Barnes. 2004. Transitivizing and detransitivizing languages. *Linguistic Typology* 8: 149–211.

# Selected references

- Say, S. 2014. Bivalent Verb Classes in the Languages of Europe: A Quantitative Typological Study. *Language dynamics and change*, 4 (2014), 116–166.
- Tsunoda, T. 1981. Split case-marking patterns in verb-types and tense / aspect / mood // Linguistics. Vol. 19. P. 389–438.
- Tsunoda, Tasaku. 1985. Remarks on transitivity. *Journal of Linguistics* 21. 385–396.
- WATP: The World Atlas of Transitivity Pairs [http://watp.ninjal.ac.jp/en/]
- Wichmann, S. & Holman, E. 2009. Assessing temporal stability for linguistic typological features. Munchen: LINCOM Europa.
- Wichmann, Søren. 2015. Statistical observations on implicational (verb) hierarchies. In: Malchukov & Comrie (eds.), 155-181.