



**HAL**  
open science

## Modelling Chinese dialect evolution

Johann-Mattis List, Shijulal Nelson-Sathi, Dagan Tal

► **To cite this version:**

Johann-Mattis List, Shijulal Nelson-Sathi, Dagan Tal. Modelling Chinese dialect evolution. Beyond Phylogeny: Quantitative diachronic explanations of language diversity (SLE Workshop), Aug 2012, Stockholm, Sweden. hprints-00758549

**HAL Id: hprints-00758549**

**<https://hal-hprints.archives-ouvertes.fr/hprints-00758549v1>**

Submitted on 28 Nov 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Modelling Chinese Dialect Evolution

Johann-Mattis List\*, Shijulal Nelson-Sathi<sup>+</sup>, and Tal Dagan<sup>+</sup>

\*Institute for Romance Languages and Literature

<sup>+</sup>Institute for Genomic Microbiology  
Heinrich Heine University Düsseldorf

2012/08/31

# Structure of the Talk

- 1 Languages
  - Languages
  - Diasystems
  - Change
- 2 Modelling Language History
  - Trees
  - Waves
  - Networks
- 3 Modelling Chinese Dialect History
  - Data
  - Analysis
  - Results

语言



ЯЗЫК



Languages

*language*



*språk*



# Languages and Dialects

Norwegian, Danish, and Swedish are different languages.

Beijing-Chinese, Shanghai-Chinese, and Hakka-Chinese are dialects of the same Chinese language.

# Languages and Dialects

Beijing Chinese	1	iou <sup>21</sup>	i <sup>55</sup>	xuei <sup>35</sup>	pei <sup>21</sup> fəŋ <sup>55</sup>	kən <sup>55</sup>	t <sup>h</sup> ai <sup>51</sup> iaŋ <sup>11</sup>	t͡səŋ <sup>55</sup>	tsai <sup>53</sup>	naə <sup>51</sup>	t͡səŋ <sup>55</sup> luən <sup>51</sup>
Hakka Chinese	1	iu <sup>33</sup>	it <sup>55</sup>	pai <sup>33</sup> a <sup>11</sup>	pet <sup>33</sup> fuŋ <sup>33</sup>	t <sup>h</sup> uŋ <sup>11</sup>	nit <sup>11</sup> t <sup>h</sup> eu <sup>11</sup>	hək <sup>33</sup>	e <sup>53</sup>		au <sup>55</sup>
Shanghai Chinese	1	fi <sup>22</sup>		t <sup>h</sup> a <sup>55</sup> tsɿ <sup>21</sup>	poʔ <sup>3</sup> foŋ <sup>44</sup>	taʔ <sup>5</sup>	t <sup>h</sup> a <sup>33</sup> fiä <sup>44</sup>	tsəŋ <sup>33</sup>	hɔ <sup>44</sup>		ləʔ <sup>1</sup> lə <sup>23</sup> tsa <sup>53</sup>
Beijing Chinese	2	ʂei <sup>35</sup>		də <sup>55</sup>		pən <sup>35</sup> liŋ <sup>21</sup>	ta <sup>51</sup>				
Hakka Chinese	2	man <sup>33</sup>	ɲin <sup>11</sup>		k <sup>w</sup> ɔ <sup>55</sup>	vɔi <sup>53</sup>					
Shanghai Chinese	2	sa <sup>33</sup>	ɲiŋ <sup>55</sup>	fiəʔ <sup>21</sup>		pəŋ <sup>33</sup> zɿ <sup>44</sup>	du <sup>13</sup>				
Norwegian	1	nu:ravɪn·ŋ	ɔ	su:lŋ						kranɿlɔt	ɔm
Swedish	1	nu:ɖanvɪndən	ɔ	su:lən		tyɪstadə	ən gɔŋ				ɔm
Danish	1	noʌnvenʔŋ	ʌ	so:lʔn	k <sup>h</sup> ʌm		enɕəŋ	i sɖɛiðʔ			ʌmʔ
Norwegian	2	vem	a	dem	sŋ	va:	ɖŋ	stærkəstə			
Swedish	2	vem	av	dɔm	sɔm	va		stærkast			
Danish	2	vemʔ	a	bŋ	ɖ	va	ɖŋ	sɖærkəstə			

# Languages and Dialects

From the perspective of the lexicon and the sound system, the Chinese **dialects** are at least equally if not more different than the Scandinavian **languages**.

# Language as a Diasystem

Languages are complex aggregates of different linguistic systems that 'coexist and influence each other' (Cosseriu 1973: 40, my translation).

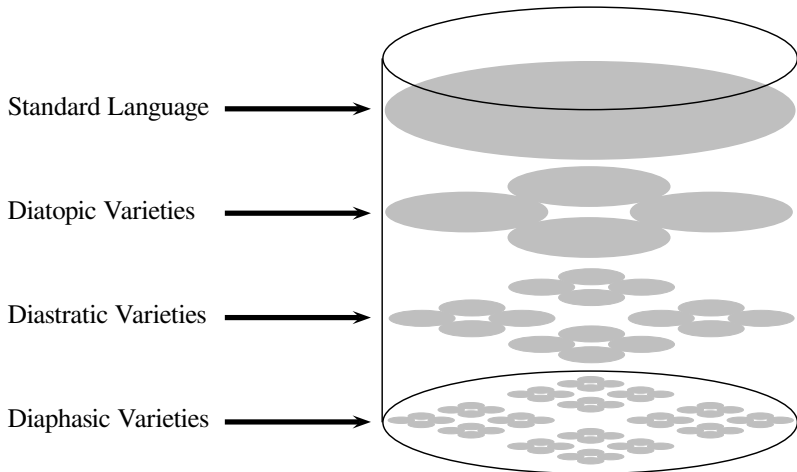


# Language as a Diasystem

Languages are complex aggregates of different linguistic systems that 'coexist and influence each other' (Cosseriu 1973: 40, my translation).

A linguistic diasystem requires a "roof language" (Goossens 1973:11), i.e. a linguistic variety that serves as a standard for interdialectal communication.

# Language as a Diasystem



# Change



# Change



expected

Mandarin

[ma<sub>55</sub>po<sub>21</sub>lou]

# Change



expected

Mandarin [ma<sub>55</sub>po<sub>21</sub>lou]

attested

Mandarin [wan<sub>51</sub>paw<sub>21</sub>lu<sub>51</sub>]

# Change



expected

Mandarin [ma<sub>55</sub>po<sub>21</sub>lou]

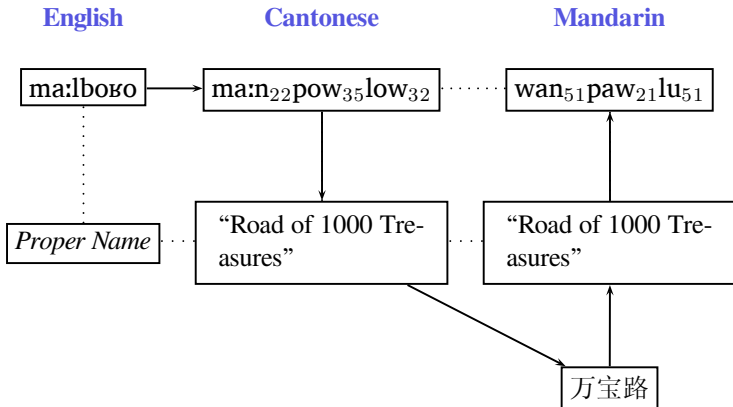
attested

Mandarin [wan<sub>51</sub>paw<sub>21</sub>lu<sub>51</sub>]

explanation

Cantonese [ma:n<sub>22</sub>pow<sub>35</sub>low<sub>32</sub>]

# Change







# Dendrophilia

August Schleicher  
(1821-1868)



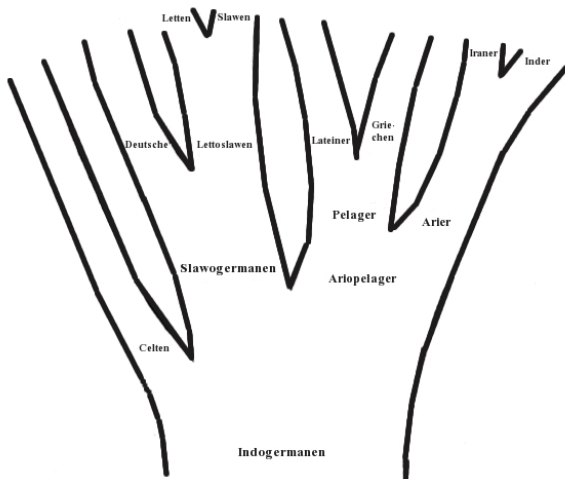
# Dendrophilia

*These assumptions that logically follow from the results of our research can be best illustrated with help of a branching tree. (Schleicher 1853: 787, my translation)*



August Schleicher  
(1821-1868)

# Dendrophilia



**Schleicher (1853)**

# Dendrophobia



Johannes Schmidt  
(1843-1901)

## Dendrophobia



*No matter how we look at it, as long as we stick to the assumption that today's languages originated from their common proto-language via multiple furcation, we will never be able to explain all facts in a scientifically adequate way. (Schmidt 1872: 17, my translation)*

Johannes Schmidt  
(1843-1901)

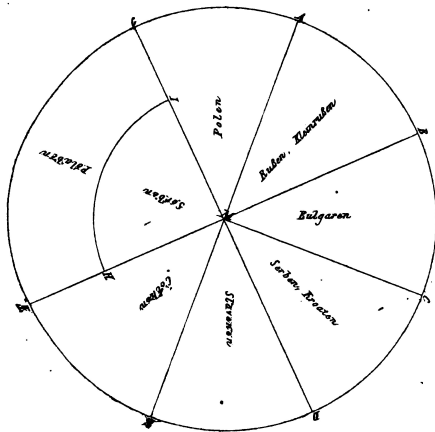
# Dendrophobia



*I want to replace [the tree] by the image of a wave that spreads out from the center in concentric circles becoming weaker and weaker the farther they get away from the center.*  
(Schmidt 1872: 27, my translation)

Johannes Schmidt  
(1843-1901)

# Dendrophobia

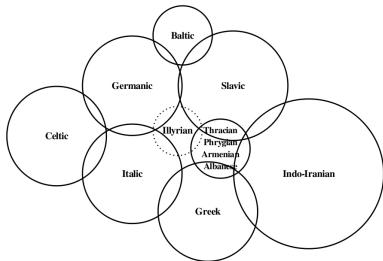


Schmidt (1875)

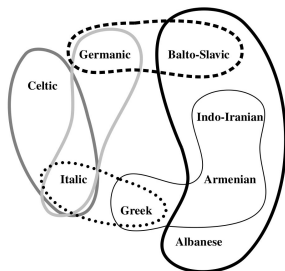
# Dendrophobia



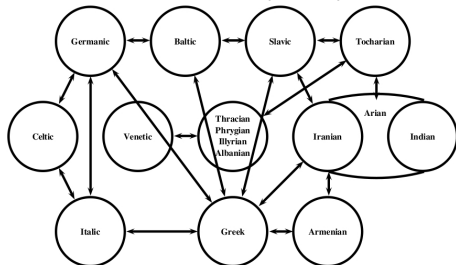
Meillet (1908)



Hirt (1905)



Bloomfield (1933)



Bonfante (1931)



# Phylogenetic Networks

Trees are bad because

# Phylogenetic Networks

## Trees are bad because

- they are difficult to reconstruct

# Phylogenetic Networks

## Trees are bad because

- they are difficult to reconstruct
- languages do not separate in split processes

# Phylogenetic Networks

## Trees are bad because

- they are difficult to reconstruct
- languages do not separate in split processes
- they are boring, since they only capture certain aspects of language history, namely the vertical relations

# Phylogenetic Networks

## Trees are bad because

- they are difficult to reconstruct
- languages do not separate in split processes
- they are boring, since they only capture certain aspects of language history, namely the vertical relations

## Waves are bad because

- nobody knows how to reconstruct them

# Phylogenetic Networks

## Trees are bad because

- they are difficult to reconstruct
- languages do not separate in split processes
- they are boring, since they only capture certain aspects of language history, namely the vertical relations

## Waves are bad because

- nobody knows how to reconstruct them
- languages still separate, even if not in split processes

# Phylogenetic Networks

## Trees are bad because

- they are difficult to reconstruct
- languages do not separate in split processes
- they are boring, since they only capture certain aspects of language history, namely the vertical relations

## Waves are bad because

- nobody knows how to reconstruct them
- languages still separate, even if not in split processes
- they are boring, since they only capture certain aspects of language history, namely, the horizontal relations

# Phylogenetic Networks



Hugo Schuchardt  
(1842-1927)



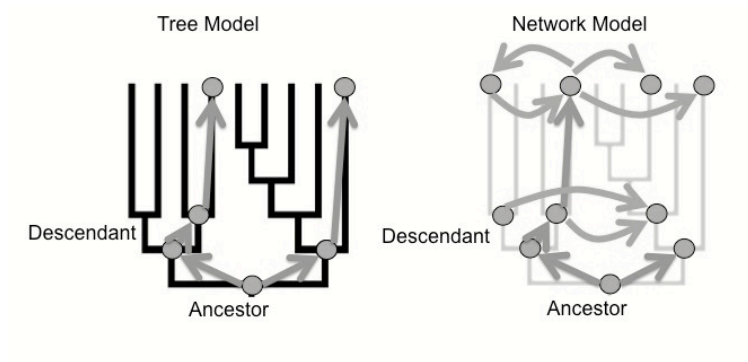
# Phylogenetic Networks



*We connect the branches and twigs of the tree with countless horizontal lines and it ceases to be a tree*  
(Schuchardt 1870 [1900]: 11)

Hugo Schuchardt  
(1842-1927)

# Phylogenetic Networks



魚 𩺰 𩺰 ?

## Modelling Chinese Dialect History

首 𩺰 𩺰 𩺰

# Data

# Data

- The data for this study was taken from the Xiàndài Hànyǔ Fāngyán Yīnkù (Hou 2004).

# Data

- The data for this study was taken from the Xiàndài Hànyǔ Fāngyán Yīnkù (Hou 2004).
- It consists of 180 items (“meanings”) translated into 40 contemporary Chinese dialects.

# Data

- The data for this study was taken from the Xiàndài Hànyǔ Fāngyán Yīnkù (Hou 2004).
- It consists of 180 items (“meanings”) translated into 40 contemporary Chinese dialects.
- The data is available on a CD in RTF format along with recordings for all dialect entries.

# Data

- The data for this study was taken from the Xiàndài Hànyǔ Fāngyán Yīnkù (Hou 2004).
- It consists of 180 items (“meanings”) translated into 40 contemporary Chinese dialects.
- The data is available on a CD in RTF format along with recordings for all dialect entries.
- For this study, the transcriptions in RTF were converted to Unicode.



# Data

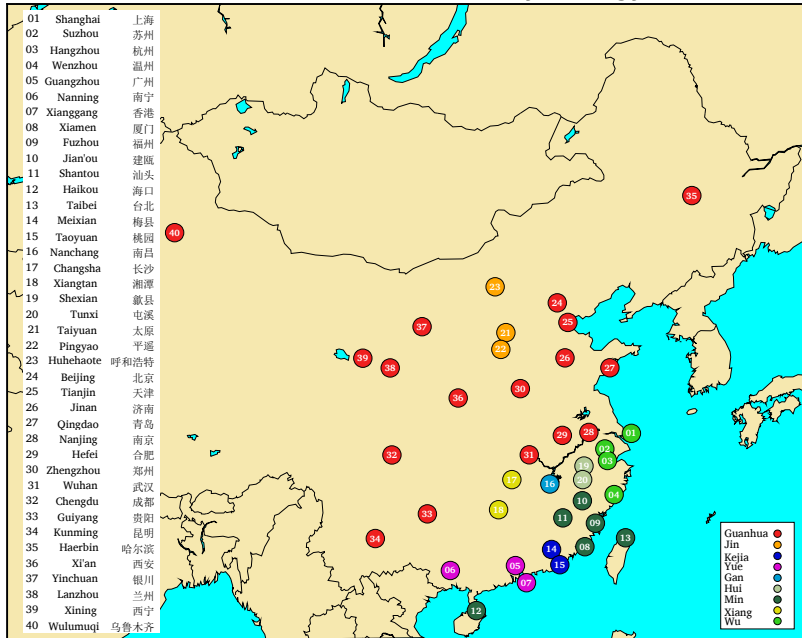
- The data for this study was taken from the Xiàndài Hànyǔ Fāngyán Yīnkù (Hou 2004).
- It consists of 180 items (“meanings”) translated into 40 contemporary Chinese dialects.
- The data is available on a CD in RTF format along with recordings for all dialect entries.
- For this study, the transcriptions in RTF were converted to Unicode.
- Every word was compared with the recordings in order to minimize errors resulting from the extraction process and the original encoding itself.

## Data

ITEM 太阳 *tàiyáng* "sun"

Dialect	Pronunciation	Characters	Cognacy
Shanghai	t <sup>h</sup> a <sup>34-33</sup> fiã <sup>13-44</sup>	太阳	1
Shanghai	ŋjɪŋ <sup>1-11</sup> dʏ <sup>13-23</sup>	日头	2
Wenzhou	t <sup>h</sup> a <sup>42-22</sup> ji	太阳	1
Wenzhou	ŋi <sup>213-22</sup> dʏu	日头	2
Guangzhou	jit <sup>2</sup> t <sup>h</sup> ɐu <sup>21-35</sup>	热头	3
Guangzhou	t <sup>h</sup> ai <sup>33</sup> joɛŋ <sup>21</sup>	太阳	1
Haikou	zit <sup>3</sup> hau <sup>31</sup>	日头	2
Beijing	t <sup>h</sup> ai <sup>51</sup> iaŋ <sup>1</sup>	太阳	1

# Dialect Locations in the Xiàndài Hànyǔ Fāngyán Yīnkù



# Analysis

## Analysis

- The data was analyzed with help of Dagan and Martin's (2008) method for phylogenetic network reconstruction, that was applied to linguistic data before (Nelson-Sathi et al. 2011).

## Analysis

- The data was analyzed with help of Dagan and Martin's (2008) method for phylogenetic network reconstruction, that was applied to linguistic data before (Nelson-Sathi et al. 2011).
- Given a binary reference tree reflecting the vertical history of a language family and a list of homologs ("cognates") distributed over the languages, the method reconstructs horizontal relations between the languages and the internal nodes of the tree.

## Analysis

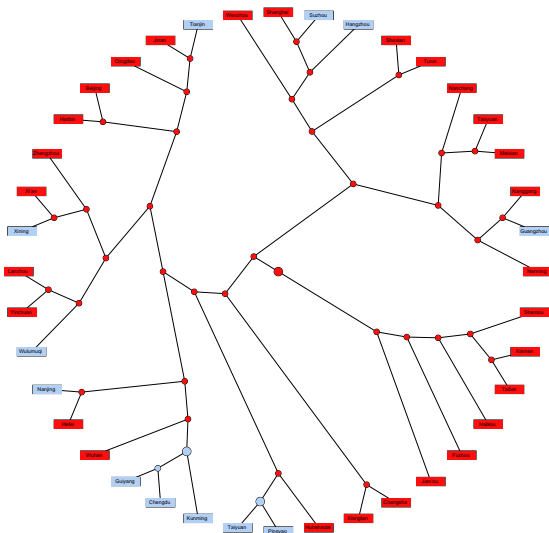
- The data was analyzed with help of Dagan and Martin's (2008) method for phylogenetic network reconstruction, that was applied to linguistic data before (Nelson-Sathi et al. 2011).
- Given a binary reference tree reflecting the vertical history of a language family and a list of homologs ("cognates") distributed over the languages, the method reconstructs horizontal relations between the languages and the internal nodes of the tree.
- The reconstruction of horizontal relations is done by seeking specific evolutionary models (loss and gain of characters) that fit the given distribution best.

## Analysis

- The data was analyzed with help of Dagan and Martin's (2008) method for phylogenetic network reconstruction, that was applied to linguistic data before (Nelson-Sathi et al. 2011).
- Given a binary reference tree reflecting the vertical history of a language family and a list of homologs ("cognates") distributed over the languages, the method reconstructs horizontal relations between the languages and the internal nodes of the tree.
- The reconstruction of horizontal relations is done by seeking specific evolutionary models (loss and gain of characters) that fit the given distribution best.
- The main criterion by which the fitness of the distributions is evaluated is the "vocabulary size", i.e. the distribution of word forms over a set of meanings. Comparing the vocabulary sizes of different models that infer different amounts of lateral events, the model that comes closest to the vocabulary sizes of the contemporary languages is chosen.

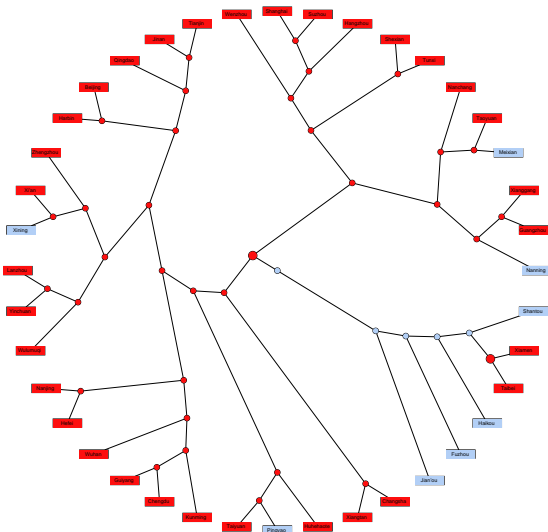


## Analysis



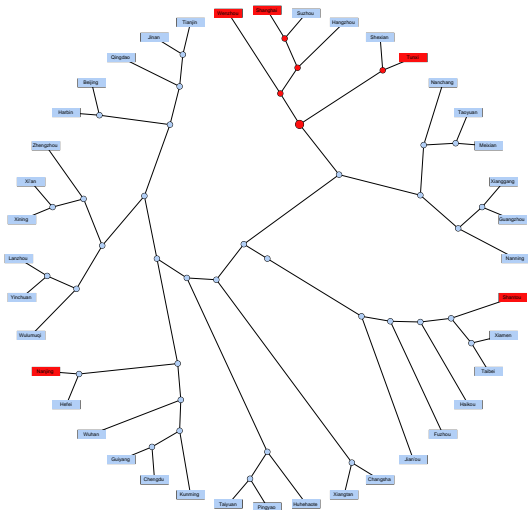
“sun” 日头 *ritou*

# Analysis



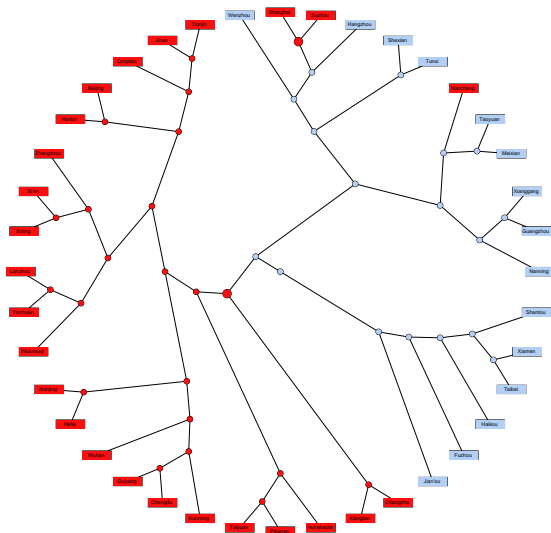
“sun” 太阳 *tàiyáng*

## Analysis



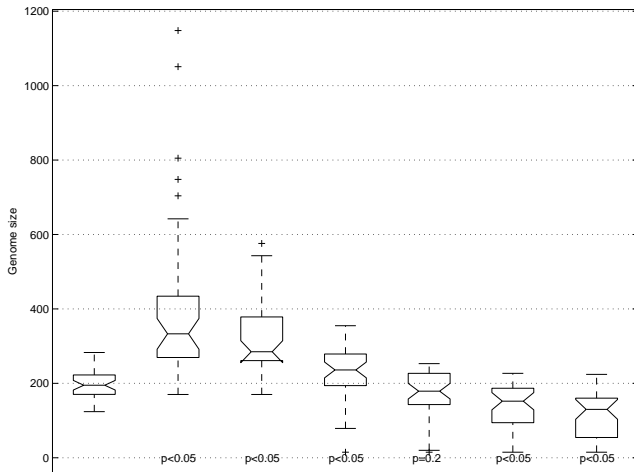
“become sick” 生病 *shēngbìng*

# Analysis



“aubergine” 茄子 *qiézi*

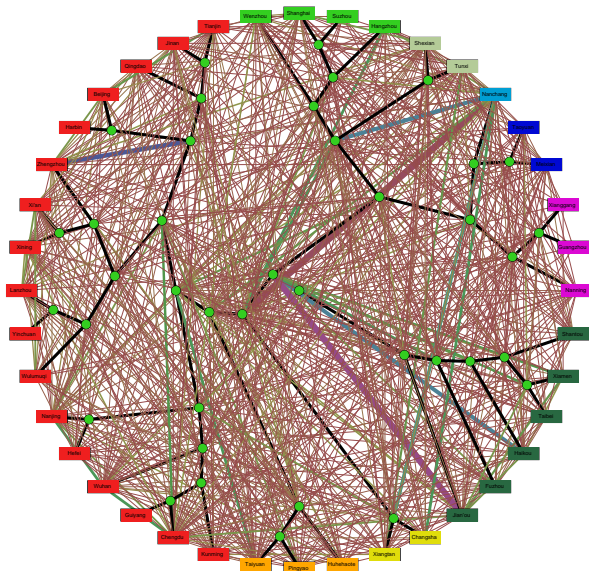
# Results



## Results

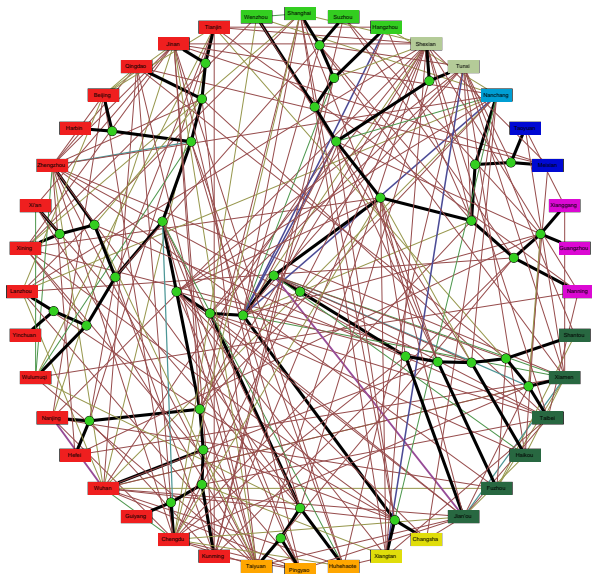
- The BOR3-model fits the distribution best. It allows up to three lateral connections per homolog.
- Out of 1152 homologs distributed over the Chinese dialects, 264 are monophyletic, 328 require one, 355 two, and 177 three lateral links in order to explain the distribution neatly.
- This corresponds to a borrowing rate of 0.5286 borrowing events per homolog per lifetime.
- For 78 percent of all homologs in the dataset the method reconstructs lateral links and therefore suggests that these have been involved in borrowing events during their history.
- Surprisingly, the 48 homologs that correspond to basic vocabulary concepts in the dataset do not show significant differences in their borrowing rates compared to the non-basic items.

# Results: General Results



Whole Dataset

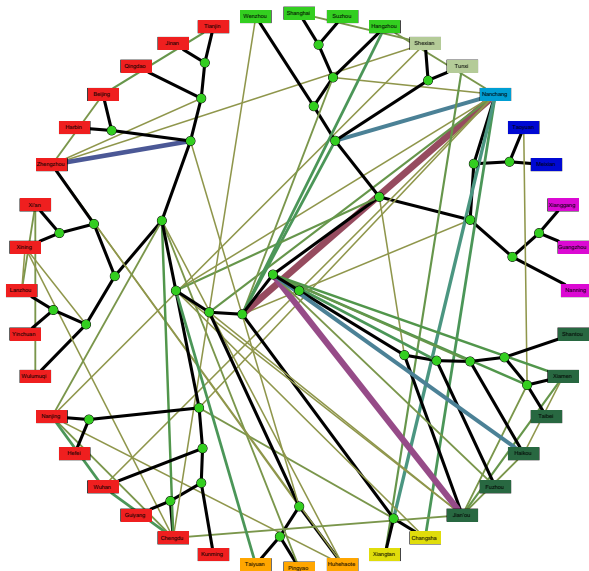
# Results: General Results



Swadesh Subset

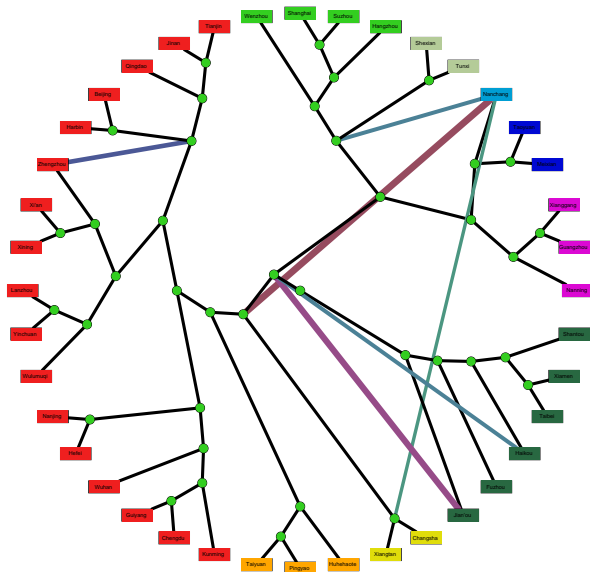


# Results: General Results



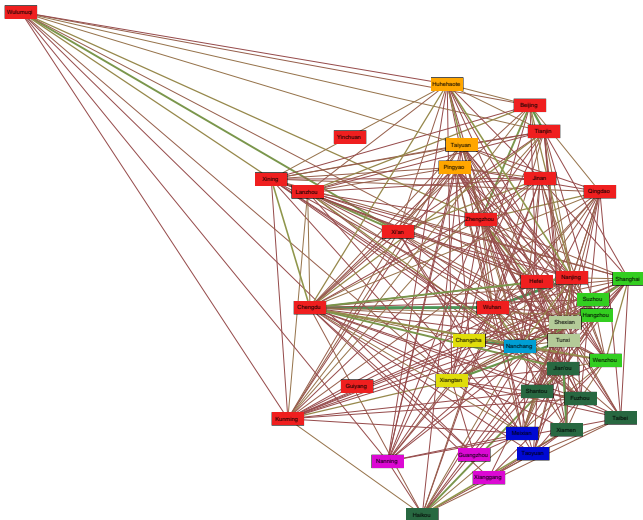
Whole Dataset (Cutoff 5)

## Results: General Results



Whole Dataset (Cutoff 10)

# Results: Chengdu



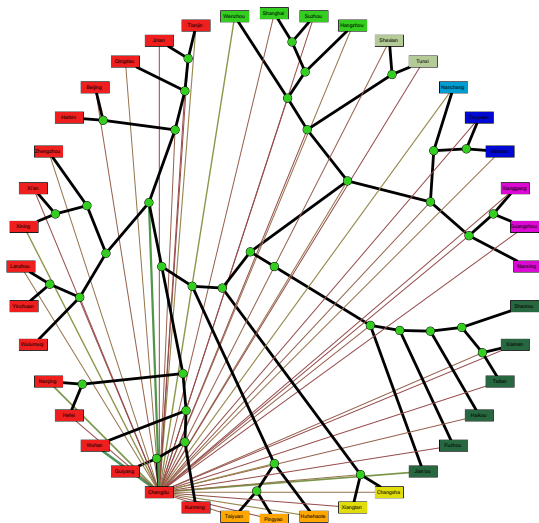
Contemporary Links Mapped to Coordinates

## Results: Chengdu



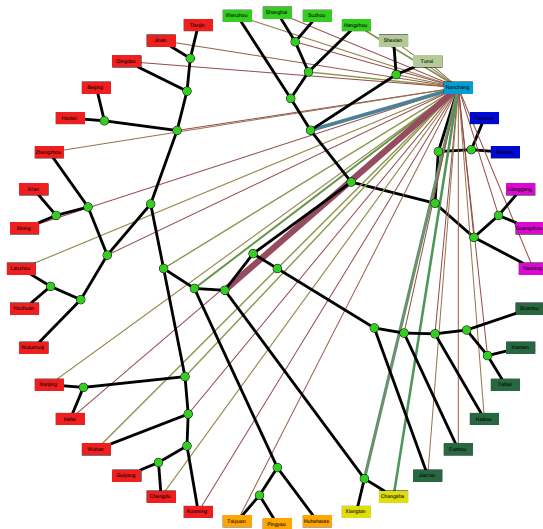
Contemporary Links of Chengdu

# Results: Chengdu



Links of Chengdu

# Results: Nanchang

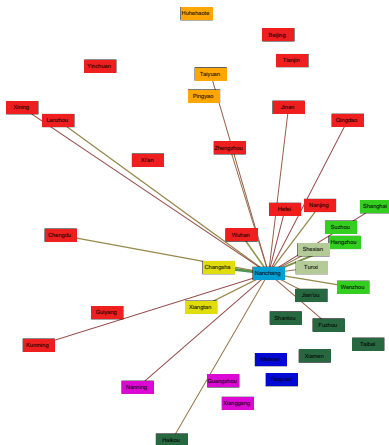


Links of Nanchang

## Results: Nanchang

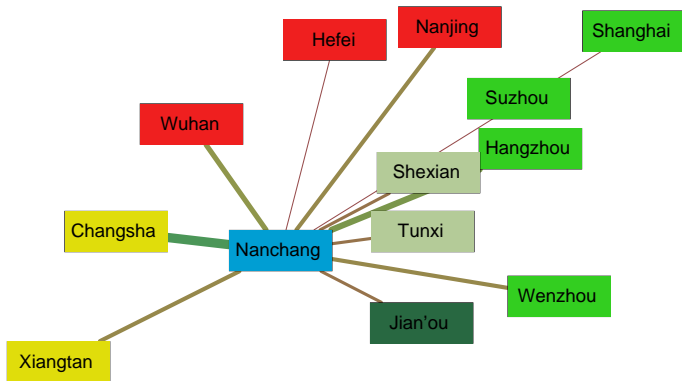
Hubness

Hubs



Contemporary Links of Nanchang

# Results: Nanchang



Links between Nanchang and its Neighbors



## Concluding Remarks

- Phylogenetic networks look nice.
- Phylogenetic networks are – if properly reconstructed – a valid alternative to both the tree and the wave model.
- We need to test the method by Dagan and Martin (2008) on more data and in more detail in order to be able to give an account on its full potential and its limits.

## Concluding Remarks

谢谢大家！

## Concluding Remarks

Thank you!