# Modelling Chinese dialect evolution

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

# Modelling Chinese Dialect Evolution

Johann-Mattis List[*], Shijulal Nelson-Sathi[+], and Tal Dagan[+]

[*]Institute for Romance Languages and Literature
[+]Institute for Genomic Microbiology
Heinrich Heine University Düsseldorf

2012/08/31

# Structure of the Talk

Languages

# 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 | iou²¹ i⁵⁵ | xuei³⁵ | pei²¹fəŋ⁵⁵ | kən⁵⁵ | tʰai⁵¹iaŋ¹¹ | tʂəŋ⁵⁵tsai⁵³ | naɚ⁵¹ | t͡ʂəŋ⁵⁵luən⁵¹ |
| **Hakka Chinese** | 1 | iu³³ it⁵⁵ | pai³³a¹¹ | pet³³fuŋ³³ | tʰuŋ¹¹ | ɲit¹¹tʰeu¹¹ | hɔk³³ | e⁵³ | au⁵⁵ |
| **Shanghai Chinese** | 1 | fii²² | tʰã⁵⁵tsɿ²¹ | poʔ³foŋ⁴⁴ | taʔ⁵ | tʰa³³ɦiã⁴⁴ | tsəŋ³³hɔ⁴⁴ | | ləʔ²¹lə²³tsa⁵³ |

| | | | | | | |
|---|---|---|---|---|---|---|
| **Beijing Chinese** | 2 | ʂei³⁵ | də⁵⁵ | pən³⁵ liŋ²¹ | ta⁵¹ | |
| **Hakka Chinese** | 2 | man³³ | ɲin¹¹ | kʷɔ⁵⁵ | vɔi⁵³ | |
| **Shanghai Chinese** | 2 | sa³³ | ɲiŋ⁵⁵ | fiəʔ²¹ | pəŋ³³ zʅ⁴⁴ | du¹³ |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Norwegian** | 1 | nuːɾɑʋinˑn̩ | ɔ | suːln̩ | | | kɾɑŋlət | əm |
| **Swedish** | 1 | nuːɖɑnvɪndən | ɔ | suːlən | tyɪstadə | ən gɔŋ | | əm |
| **Danish** | 1 | noʌʌnvenˀn̩ | ʌ | soːlˀn̩ | kʰʌm | eŋɡ̊aŋ | i sd̥ʁɪðˀ | ʌmˀ |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Norwegian** | 2 | ʋem | ɑ | dem | sm̩ | ʋɑː | dn̩ | stæɾkəstə |
| **Swedish** | 2 | ʋɛm | ɑv | dɔm | sɔm | ʋɑ | | stɑɹkast |
| **Danish** | 2 | ʋɛmˀ | a | b̥m̩ | d̥ | ʋɑ | d̥n̩ | sd̥æʌɡ̊əsd̥ə |

# 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**.
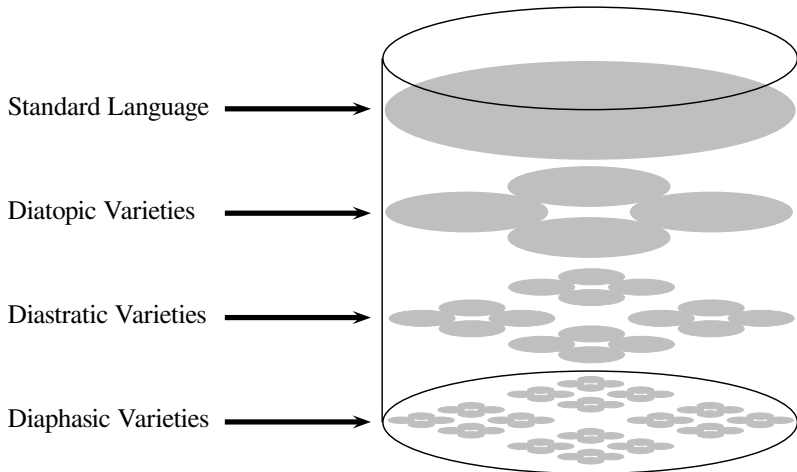
# Language as a Diasystem

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

# Language as a Diasystem

Languages are complex aggregates of different linguistic systems that 'coexist and influence each other' (Coseriu 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



Standard Language

Diatopic Varieties

Diastratic Varieties

Diaphasic Varieties

# Change

# Change



expected        Mandarin        $[ma_{55}po_{21}lou]$

# Change



expected    Mandarin    $[ma_{55}po_{21}lou]$

attested    Mandarin    $[wan_{51}paw_{21}lu_{51}]$

# Change



| expected | Mandarin | $[ma_{55}po_{21}lou]$ |
| attested | Mandarin | $[wan_{51}paw_{21}lu_{51}]$ |
| explanation | Cantonese | $[maːn_{22}pow_{35}low_{32}]$ |

# Change



English        Cantonese        Mandarin

| ma:lboʁo | → | ma:n$_{22}$pow$_{35}$low$_{32}$ | ······ | wan$_{51}$paw$_{21}$lu$_{51}$ |

*Proper Name* ····· "Road of 1000 Treasures" ····· "Road of 1000 Treasures"
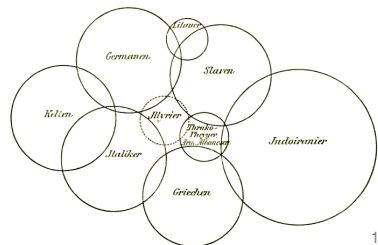
万宝路

# Modelling Language History

# Dendrophilia
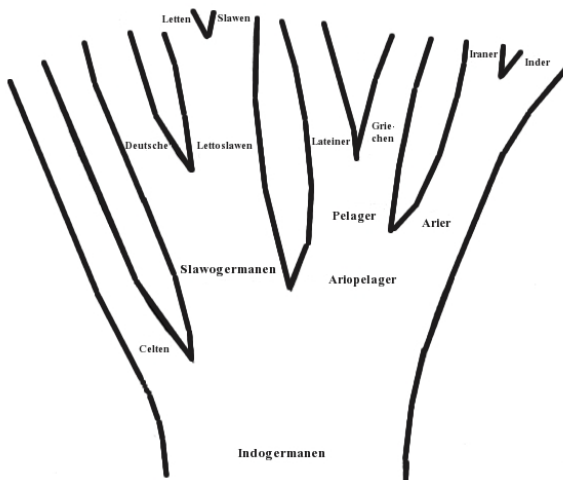


August Schleicher
(1821-1868)

# 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)
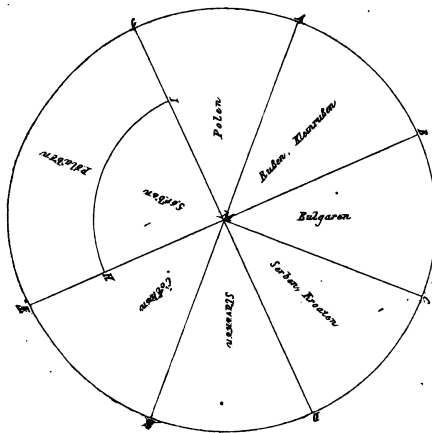
Johannes Schmidt
(1843-1901)

# 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)
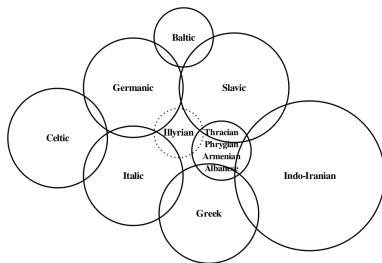
Johannes Schmidt
(1843-1901)

# Dendrophobia
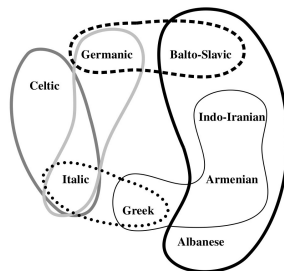


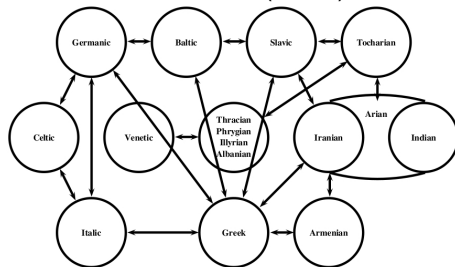**Schmidt (1875)**

# Dendrophobia



Meillet (1908)



Bloomfield (1933)



Hirt (1905)



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
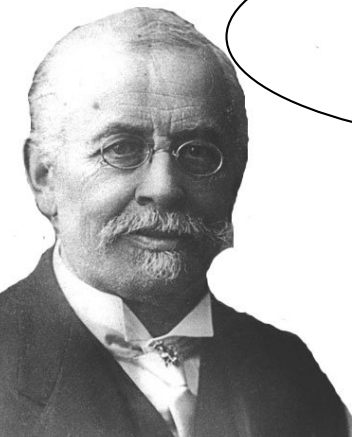
## Waves are bad because

- 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
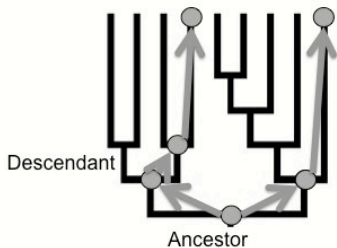(1842-1927)

# 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

**Modelling Chinese Dialect History**

# 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.
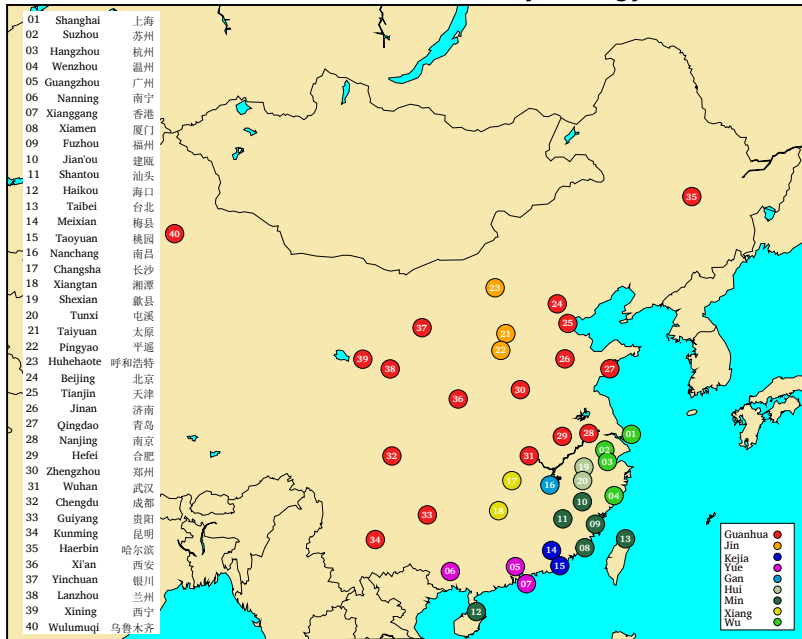- For this study, the transcriptions in RTF were converted to Unicode.

# 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 | $t^ha^{34\text{-}33}ɦiã^{13\text{-}44}$ | 太阳 | 1 |
| Shanghai | $ȵjɿʔ^{1\text{-}11}dɤ^{13\text{-}23}$ | 日头 | 2 |
| Wenzhou | $t^ha^{42\text{-}22}ji$ | 太阳 | 1 |
| Wenzhou | $ȵi^{213\text{-}22}dɤu$ | 日头 | 2 |
| Guangzhou | $jit^2t^hɐu^{21\text{-}35}$ | 热头 | 3 |
| Guangzhou | $t^hai^{33}jœŋ^{21}$ | 太阳 | 1 |
| Haikou | $zit^3hau^{31}$ | 日头 | 2 |
| Beijing | $t^hai^{51}iɑŋ^1$ | 太阳 | 1 |

# Dialect Locations in the Xiàndài Hànyǔ Fāngyán Yīnkù

| 01 | Shanghai | 上海 |
| 02 | Suzhou | 苏州 |
| 03 | Hangzhou | 杭州 |
| 04 | Wenzhou | 温州 |
| 05 | Guangzhou | 广州 |
| 06 | Nanning | 南宁 |
| 07 | Xianggang | 香港 |
| 08 | Xiamen | 厦门 |
| 09 | Fuzhou | 福州 |
| 10 | Jian'ou | 建瓯 |
| 11 | Shantou | 汕头 |
| 12 | Haikou | 海口 |
| 13 | Taibei | 台北 |
| 14 | Meixian | 梅县 |
| 15 | Taoyuan | 桃园 |
| 16 | Nanchang | 南昌 |
| 17 | Changsha | 长沙 |
| 18 | Xiangtan | 湘潭 |
| 19 | Shexian | 歙县 |
| 20 | Tunxi | 屯溪 |
| 21 | Taiyuan | 太原 |
| 22 | Pingyao | 平遥 |
| 23 | Huhehaote | 呼和浩特 |
| 24 | Beijing | 北京 |
| 25 | Tianjin | 天津 |
| 26 | Jinan | 济南 |
| 27 | Qingdao | 青岛 |
| 28 | Nanjing | 南京 |
| 29 | Hefei | 合肥 |
| 30 | Zhengzhou | 郑州 |
| 31 | Wuhan | 武汉 |
| 32 | Chengdu | 成都 |
| 33 | Guiyang | 贵阳 |
| 34 | Kunming | 昆明 |
| 35 | Haerbin | 哈尔滨 |
| 36 | Xi'an | 西安 |
| 37 | Yinchuan | 银川 |
| 38 | Lanzhou | 兰州 |
| 39 | Xining | 西宁 |
| 40 | Wulumuqi | 乌鲁木齐 |

Legend:
- Guanhua
- Jin
- Kejia
- Yue
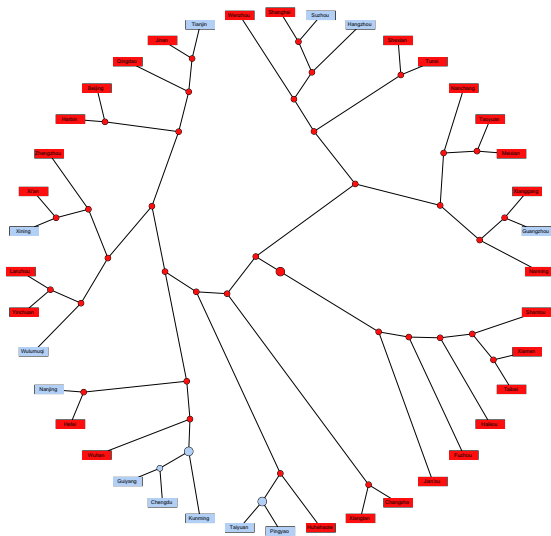- Gan
- Hui
- Min
- Xiang
- Wu

# 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).
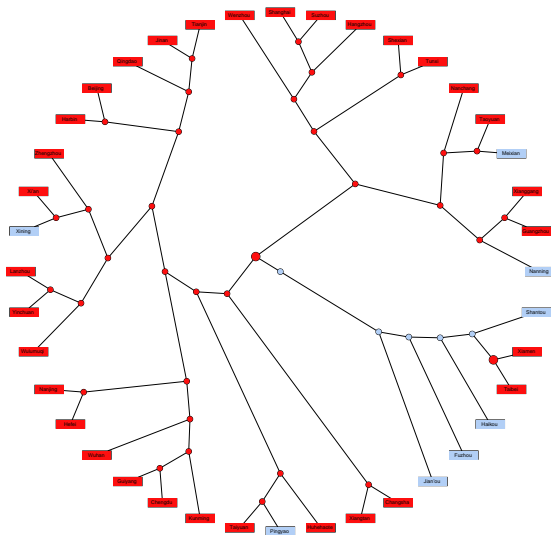
## Analysis

- 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
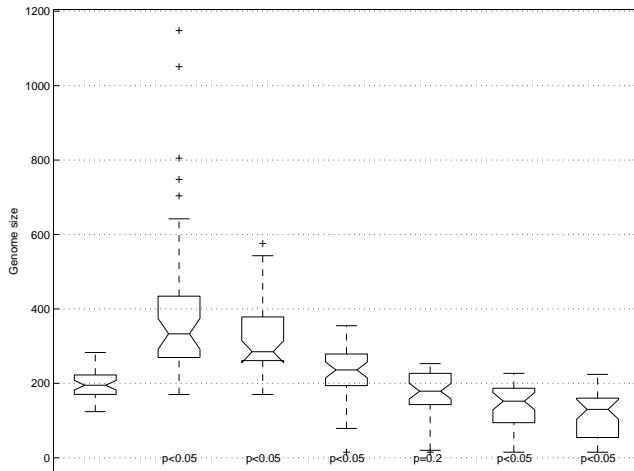


"sun" 太阳 *tàiyáng*

# Analysis



"become sick" 生病 *shēngbìng*
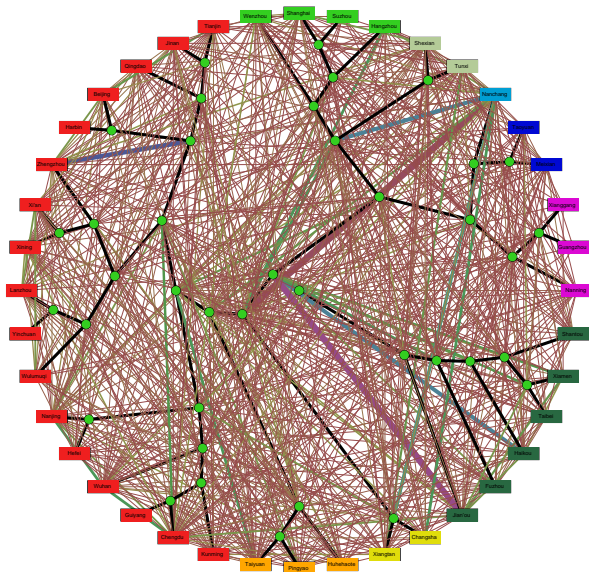
# Analysis



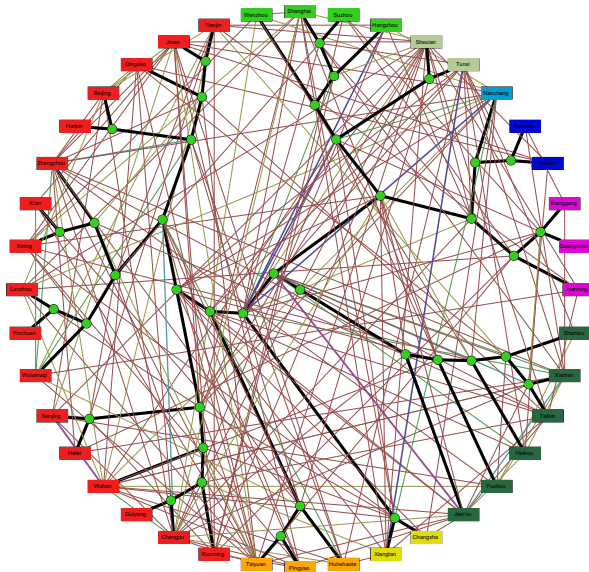"aubergine" 茄子 *qiézi*

# 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.

- Suprisingly, 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.
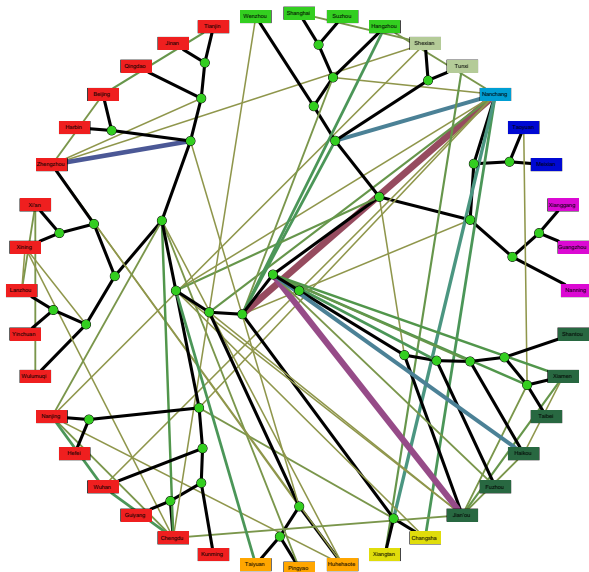
# Results: General Results
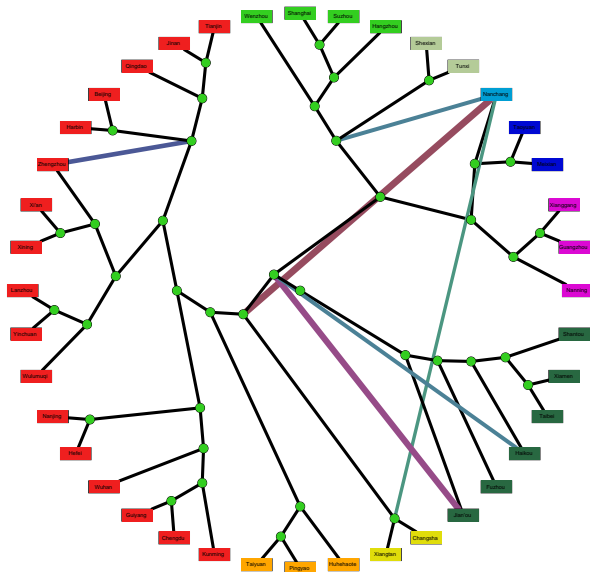


Whole Dataset

# Results: General Results


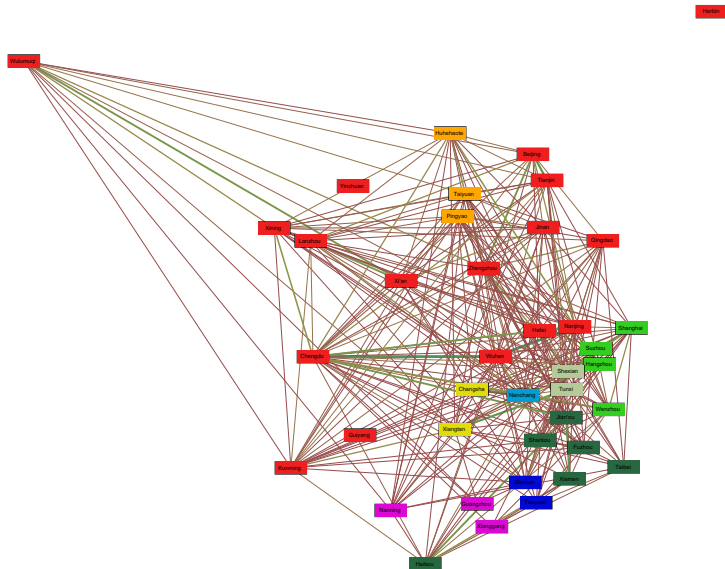
Swadesh Subset

# Results: General Results



Whole Dataset (Cutoff 5)
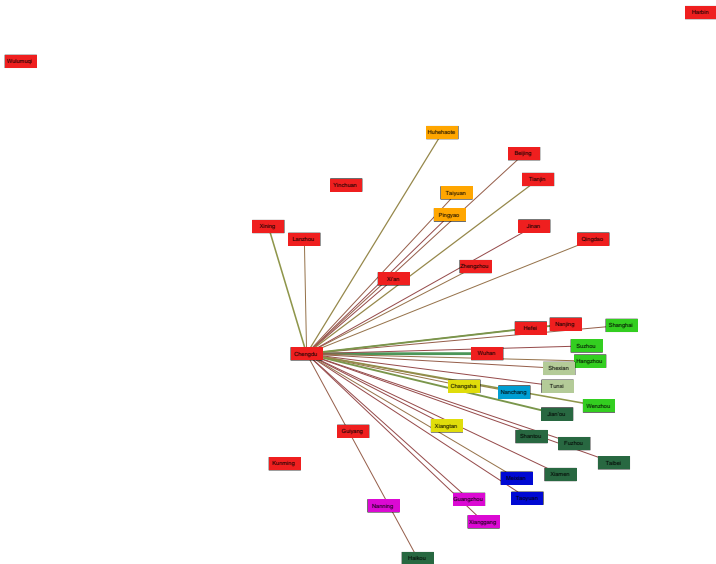
# Results: General Results



Whole Dataset (Cutoff 10)
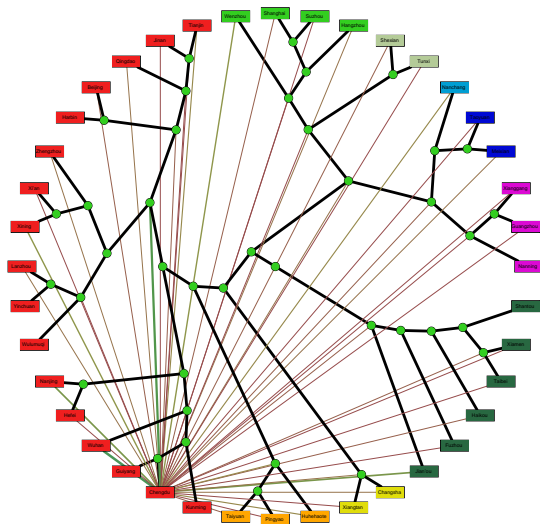
# Results: Chengdu



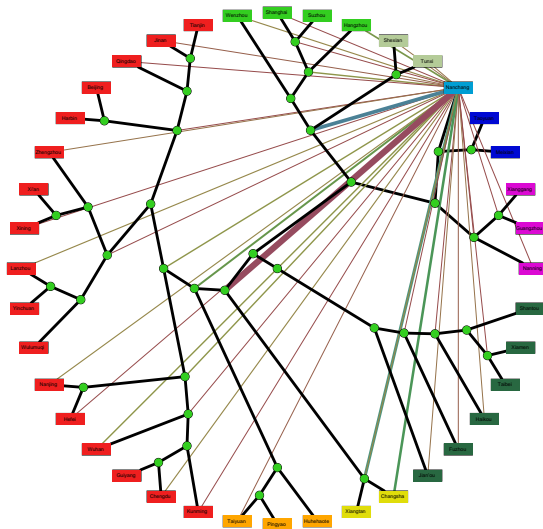Contemporary Links Mapped to Coordinates

# Results: Chengdu



Contemporary Links of Chengdu
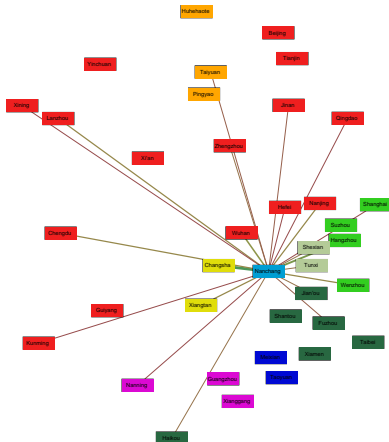
# Results: Chengdu



Links of Chengdu
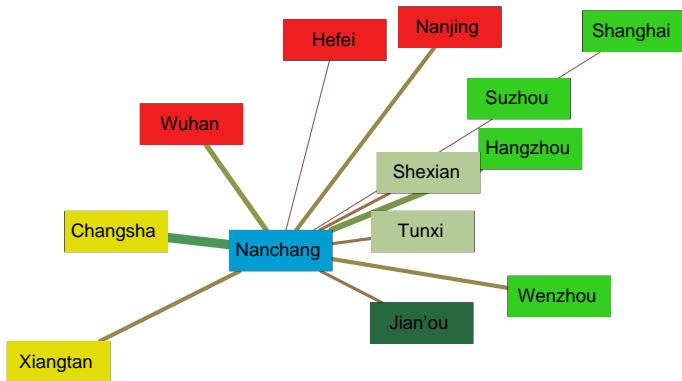
# Results: Nanchang



Links of Nanchang

# Results: Nanchang



Contemporary Links of Nanchang

# Results: Nanchang



Hefei

Nanjing

Shanghai

Wuhan

Suzhou

Hangzhou

Shexian

Changsha

Tunxi

Nanchang

Wenzhou

Xiangtan

Jian'ou

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！

# Thank you!