

A New Dataset for Tonal Segmental Dialectometry

From the Yue- and Pinghua-Speaking Area

Matthew Sung, Jelena Prokic & Yiya Chen | 6th SIGTYP Workshop (Malta)



**Universiteit
Leiden**
The Netherlands

Discover the world at Leiden University

Motivation

Traditional Dialectology

- Originated in Europe
- Predominantly been focusing on sounds (segments)

Motivation

Traditional Dialectology

- Originated in Europe
- Predominantly been focusing on sounds (segments)

Dialectometry/ Quantitative Dialectology

- Similar to traditional dialectology
- Mostly on European languages, but have been expanding beyond Europe recently

Motivation

Traditional Dialectology

- Originated in Europe
- Predominantly been focusing on sounds (segments)

Dialectometry/ Quantitative Dialectology

- Similar to traditional dialectology
- Mostly on European languages, but have been expanding beyond Europe recently

Tone languages

- More than half of the world's languages are tonal (Yip 2002)
- Not many datasets with more than 20 dialects (of the same language) are publically available which contains both tones and segments

New Dataset!

Sinitic languages

- Yue and Pinghua
- Spoken in Southern China
- Tone languages

New Dataset!

Sinitic languages

- Yue and Pinghua
- Spoken in Southern China
- Tone languages

It consists of:

- 104 varieties
- 130 monosyllabic words
 - E.g. Body parts, numbers, animals, geographical features etc.
- Both segments and tones

Data sources

Dialect Surveys and Homonymic Syllabaries

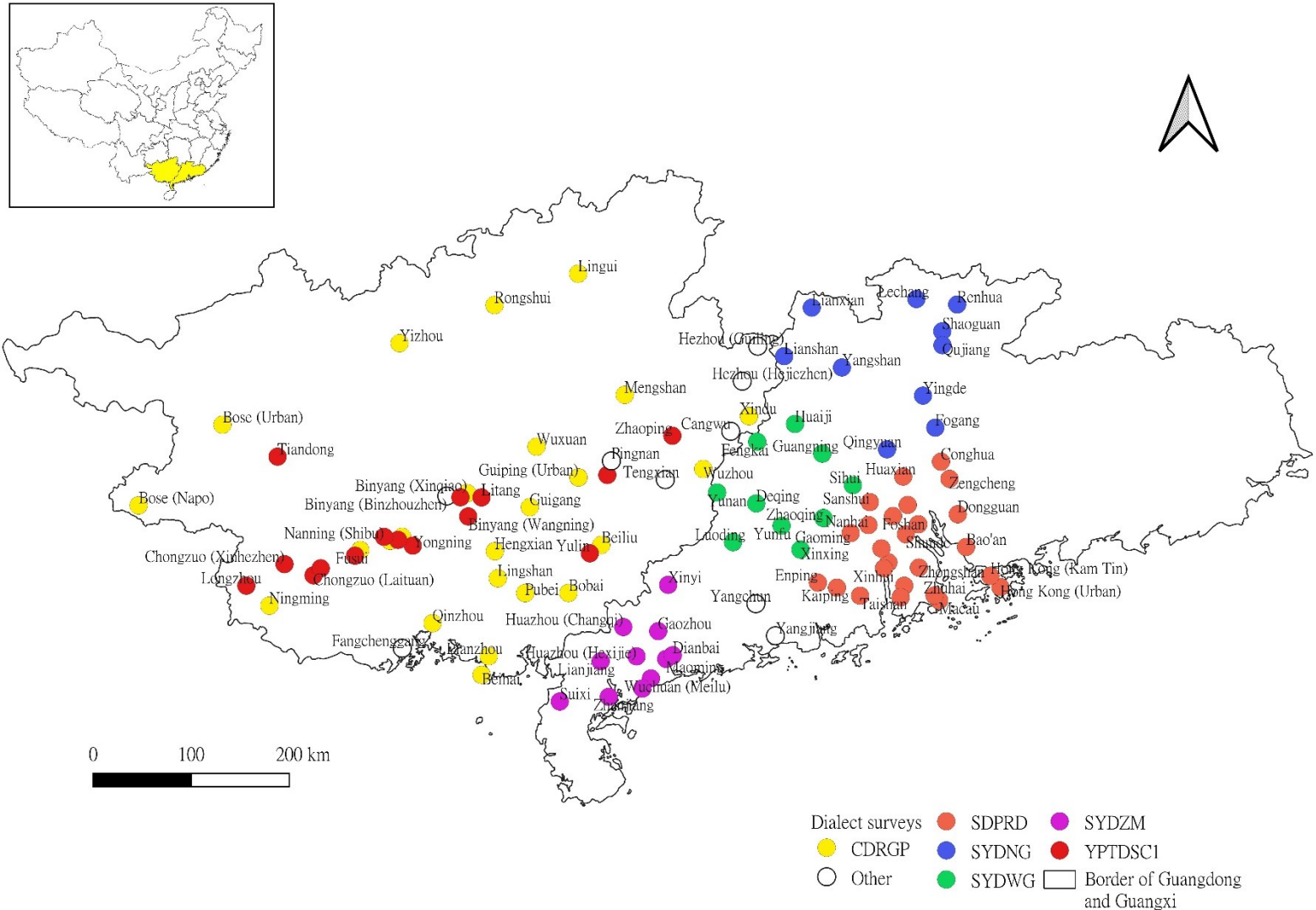
韻攝	1	2	3	4	5	6	7	8	9	10
韻攝	多	拖	他	駝	駝	駝	大	駝	椰	哪
韻攝	果開一	果開一	果開一	果開一	果開一	果開一	果開一	果開一	果開一	果開一
韻攝	平歌端	平歌透	平歌定	平歌定	平歌定	上聲定	去聲定	去聲定	去聲定	去聲定
韻攝	tuə ⁵⁵	tʰuə ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	huə ⁵⁵	na ⁵⁵
廣州(市區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
香港(市區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
汕頭(新區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
廈門(市區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
番禺(市區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
花縣(花山)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
從化(城內)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
增城(縣城)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
佛山(市區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
南海(沙頭)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
順德(大良)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
三水(西樵)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
高明(明城)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
中山(石岐)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
珠海(前山)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
斗門(上橫水上鄉)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
斗門(斗門鎮)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
江門(白沙)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
新會(會城)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
台山(台城)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
開平(赤坎)	u ⁵⁵	hu ⁵⁵	ha ⁵⁵	hu ⁵⁵	hu ⁵⁵	hu ⁵⁵	hu ⁵⁵	hu ⁵⁵	no ⁵⁵	na ⁵⁵
恩平(牛江)	tuə ⁵⁵	huə ⁵⁵	ha ⁵⁵	huə ⁵⁵	huə ⁵⁵	huə ⁵⁵	huə ⁵⁵	huə ⁵⁵	no ⁵⁵	na ⁵⁵
鶴山(鶴城)	ɔu ⁵⁵	hu ⁵⁵	ha ⁵⁵	hu ⁵⁵	hu ⁵⁵	hu ⁵⁵	hu ⁵⁵	hu ⁵⁵	no ⁵⁵	na ⁵⁵
東莞(莞城)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
東莞(沙井)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
惠州(市區)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
東莞(清遠)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
從化(良田)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
中山(龍山)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵
中山(龍山)	tɔ ⁵⁵	tʰɔ ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	tʰa ⁵⁵	no ⁵⁵	na ⁵⁵

Dialect Survey (Survey of Dialects in the Pearl River Delta)

- ɛ [ɿ]惹□願□njɛpɿ-:蜻蜓 [ɿ]嘢东西
- o
- p [ɿ]□~老:大舅母 [ɿ]婆袍大棉~:棉袄 [ɿ]抱□泡沫 [ɿ]拖解暴~露□--:一串
- 6 [ɿ]煲 [ɿ]□像蝙蝠的一种鸟 [ɿ]保宝堡碉~ [ɿ]报
- m [ɿ]□脸肿 [ɿ]毛 [ɿ]冒帽
- f [ɿ]搔骚臊 [ɿ]曹槽槽 [ɿ]嫂 [ɿ]造 [ɿ]扫~地|~把 [ɿ]□早~晚~:早稻晚稻
- t [ɿ]糟遭 [ɿ]掏滔涛桃陶淘逃萄葡~绸捆扎 [ɿ]早枣蚤 [ɿ]导祷祈~ [ɿ]灶~头:灶 [ɿ]道稻盗
- tʰ [ɿ]操 [ɿ]讨草 [ɿ]套澡糙燥
- d [ɿ]刀叨 [ɿ]岛捣倒打~|~茶 [ɿ]到
- n [ɿ]恼脑 [ɿ]□瞪眼睛
- l [ɿ]劳牢痨 [ɿ]佬 [ɿ]老 [ɿ]涝
- ʃ [ɿ]□形容密而多 [ɿ]□做~:做什么 [ɿ]傻
- k [ɿ]高糕篙膏羔蒿蒿~ [ɿ]稿 [ɿ]个告筒~Gokɿ:这里
- ŋ [ɿ]我 [ɿ]傲 [ɿ]熬
- h [ɿ]好~人 [ɿ]耗好喜~
- θ [ɿ]豪毫壕 [ɿ]澳~门 [ɿ]浩号蚝~油
- w [ɿ]浣~lɛwɿ:髒
- i
- p [ɿ]皮肥脾琵琶~髻枇~把 [ɿ]婢被棉~|~迫 [ɿ]备鼻避

Homonymic Syllabary (The Phonology of Cangwu Local Vernacular in Guangxi)

Localities and their respective sources



Segmental Data



Universiteit
Leiden
The Netherlands

Discover the world at Leiden University

Modifications to the original transcriptions

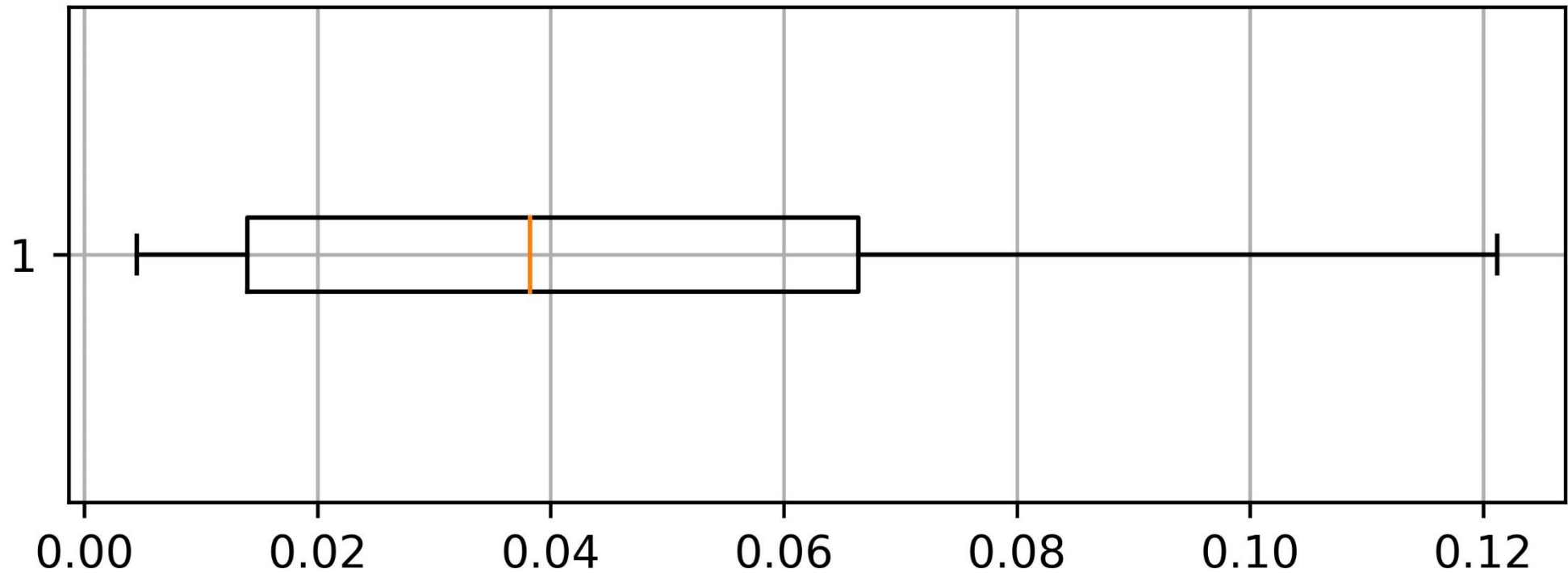
- Transcribers' differences
 - More than 1 transcriber
 - Different conventions/ habits of transcription for the same sound
 - Can influence the dialectometric analysis
- Normalisation/ cleaning required

Modifications

1. Comparison with existing recordings
 - Guangzhou 'water' sœy → søy
2. Maintaining contrasts
3. Removal of redundant characters
 - Gaoyang dialects 'person' niɛn → nɛn
4. Simplification of overly detailed transcriptions
 - HK (Kam Tin) a̱ → a
5. Consistency of onsets
6. Conversion from Chinese IPA to Standard IPA
 - Guangzhou 'skin' p'ei → p^hei
7. Phonetic alignment

Levenshtein Distance between Raw vs. Cleaned Transcriptions

Boxplot of Distances between Raw and Cleaned Transcriptions



Tonal Data



Universiteit
Leiden
The Netherlands

Discover the world at Leiden University

Tonal data

- The existing digital tonal datasets generally consist of about 20-30 dialects
- Our dataset
 - Same number of words as the segmental data
 - Same items (130)
 - Same dialects (104)
 - Allows us to compare the segmental and tonal variation of the same area

Tone notations

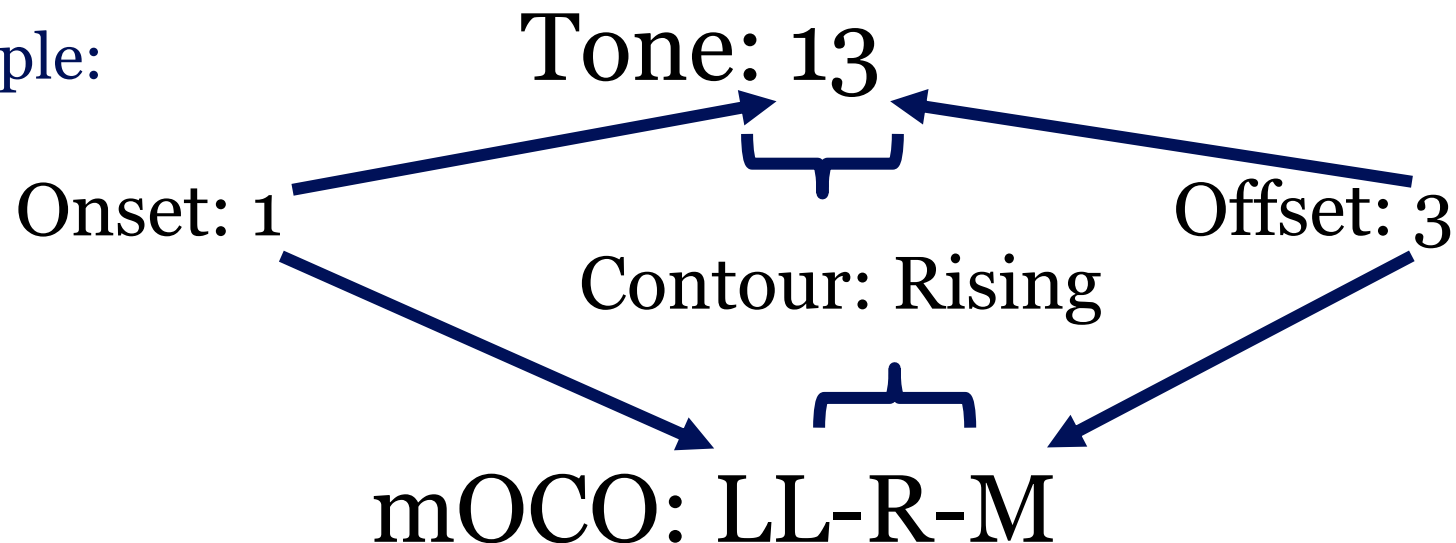
- Chao's (1930) tone letters
- Contour levels: 1, 2, 3, 4, 5
 - 1 is the lower pitch level, 5 is the highest
- Combination of pitch levels to represent a tone contour
 - Level tones: 11, 22, 33, 44, 55
 - Rising tones: e.g. 12, 24, 25, 45
 - Falling tones: e.g. 51, 42, 31, 43
 - Concave tones: e.g. 413, 512
 - Convex tones: e.g. 232, 253

Tone representation in dialectometry

- Chao's (1930) transcription system cannot directly be used for dialectometry (Sung et al. Forthcoming)
- Requires further conversion to yield meaningful tone distances (when using tools like *Gabmap* or *LED-A.org*)
 - Tone-to-string (Tang 2009)
 - Onset-Contour-Offset (Yang and Castro 2008)
 - **Modified Onset-Contour-Offset**

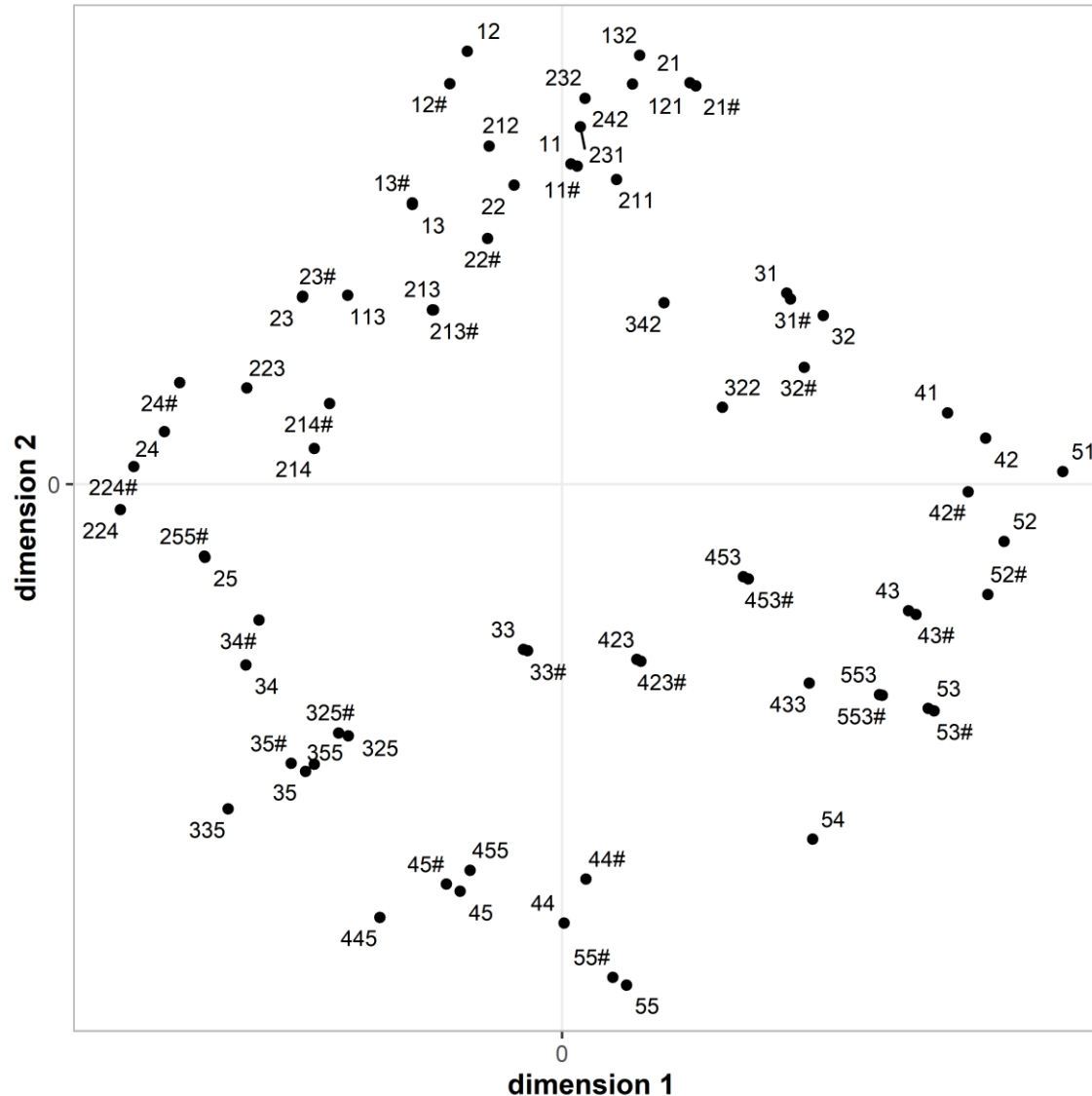
Modified Onset-Contour-Offset (mOCO)

- Modified from Yang and Castro's (2008) representation
- Converts Chao's (1930) transcription system into 3 parts:
 - Onset (starting pitch of the tone)
 - Contour (shape of the tone)
 - Offset (ending pitch of the tone)
- Here is an example:



Tone distances calculated with mOCO

- Tone distances: applying Levenshtein Distance (Levenshtein 1966, Heeringa 2004) on the mOCO representation
- This representation corresponds to the perceptual dimensions (Gandour and Harshman 1978)
- It can differentiate 72 out of 73 of the tones in the data



Preliminary analysis of tonal variation



Universiteit
Leiden
The Netherlands

Discover the world at Leiden University

Tonal variation

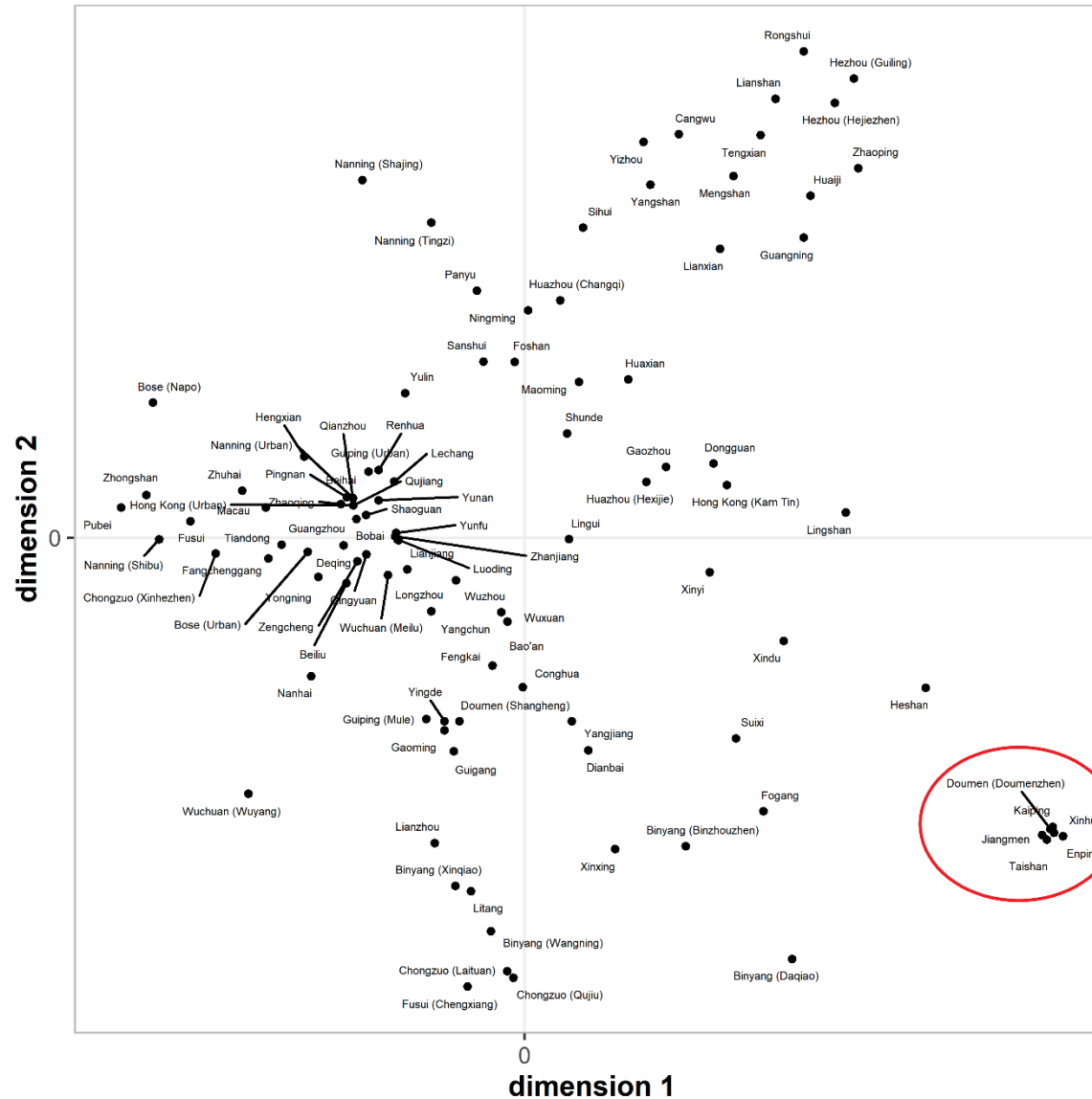
- In Chinese dialectology, very often studies on tones are descriptive
 - Correspondences between Middle Chinese tone categories and present-day pronunciation
- Existing dialectometric studies of tonal languages
 - Sometimes tones were neglected
 - Typically with a small dataset
 - Often combined with segments
 - The tone distance metrics cannot differentiate a lot of the tones in the data (Sung et al. Forthcoming)

Tone distances between dialects

- Between each pair of dialects in the data
 - Calculate the Levenshtein distance between the tones of each word
 - Sum the distances and divide it by the number of words compared (normalisation)

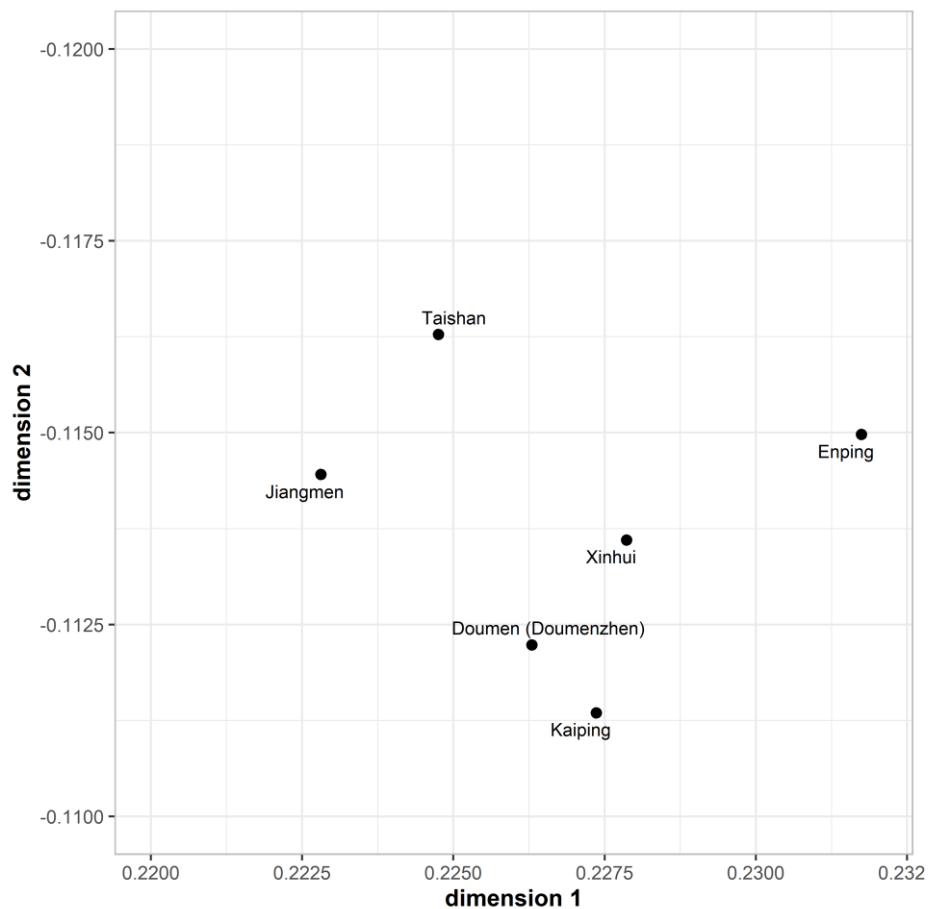
Dialect distances of Yue (tones only)

- MDS plot of 104 Yue dialects
- General pattern: continuum-like
- There is a clear cluster which seem to be isolated from the broader continuum
 - Siyi dialects (red circle)



Dialect distances of Siyi dialects (tones only)

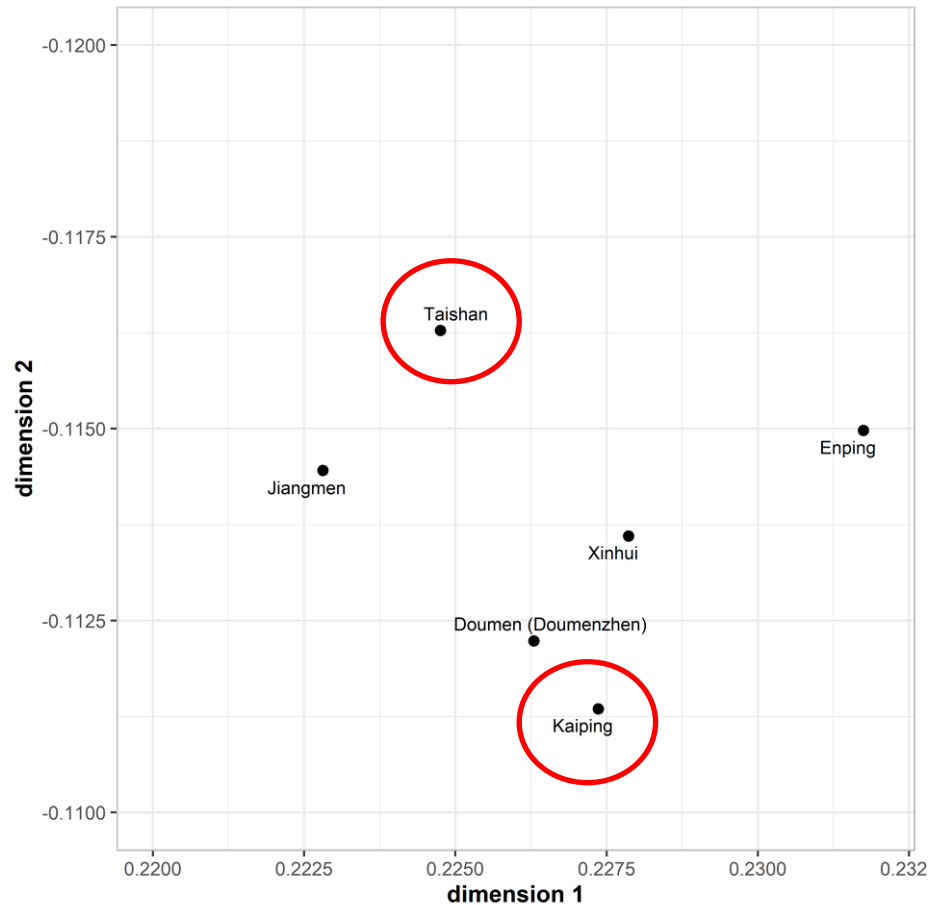
- Siyi dialects
- We can see 3 patterns here:
 - Group 1: Doumen, Taishan, Kaiping
 - Group 2: Enping
 - Group 3: Jiangman, Xinhui



Tone Categories	Doumen	Taishan	Kaiping	Enping	Jiangmen	Xinhui
Yin Ping	33	33	33	33	23	23
Yang Ping	22	22	22	22	22	22
Yin Shang	55	55	55	55	45	45
Yang Shang	21	21	21	31	21	21
Qu	31	31	31		31	31
Yin Ru1	55#	55#	55#	55#	55#	55#
Yin Ru2	33#	33#	33#	33#	33#	33#
Yang Ru	21#	21#	21#	21#	21#	21#

Dialect distances of Siyi dialects (tones only)

- Group 1 dialects
- Same tonemic inventory
- Why do they not overlap on the MDS plot?
- How do we explain this pattern?



Tone Categories	Doumen	Taishan	Kaiping	Enping	Jiangmen	Xinhui
Yin Ping	33	33	33	33	23	23
Yang Ping	22	22	22	22	22	22
Yin Shang	55	55	55	55	45	45
Yang Shang	21	21	21	31	21	21
Qu	31	31	31	31	31	31
Yin Ru1	55#	55#	55#	55#	55#	55#
Yin Ru2	33#	33#	33#	33#	33#	33#
Yang Ru	21#	21#	21#	21#	21#	21#

Tone correspondences

- Tone distances has to be accompanied with a tone correspondence table to get further insights
- We see *Lexical Distribution* differences (Wells 1982), or ‘exceptions’ in grey
- These are detected and reflected in the aggregate tone distances

Correspondences	No. of Items
11# : 21#	1
21:21	4
21:31	2
21# : 21#	12
21# : 33#	1
22:22	21
22:55	1
31:31	10
33:21	2
33:33	33
33:55	1
33# : 21#	1
33# : 33#	3
35:21	1
55:55	26
55# : 55#	11

Correspondence Table of Tones between Taishan (left) and Kaiping (right) Dialects (irregular correspondences in gray)

Conclusion

New Dataset

- 104 Yue-Pinghua dialects
- 130 words
- Segments and tones

New tone distance metric

- mOCO is able to differentiate 72/73 tones in the dataset

New discovery on tonal variation

- There seems to be a dialect continuum on the tonal level in the Yue-speaking area
- Tones can vary on the tonetic, tonemic as well as lexical distribution

Data availability

OSF repository:

<https://osf.io/j5gxz/>

- Contains the data and
Tone conversion scripts (Python)

Full datasets are under embargo at the moment



Contact info:

Matthew Sung:

h.w.m.sung@hum.leidenuniv.nl

Thank you for your attention!



Universiteit
Leiden
The Netherlands

Contact info:

Matthew Sung:

h.w.m.sung@hum.leidenuniv.nl

Personal website:

<https://sites.google.com/view/matthew-sung-dialectologist/home>