

# Texts, Treebanks, Networks descriptions and analysis of languages

XINYING CHEN

University of Ostrava, Nov. 8, 2017

# Outline

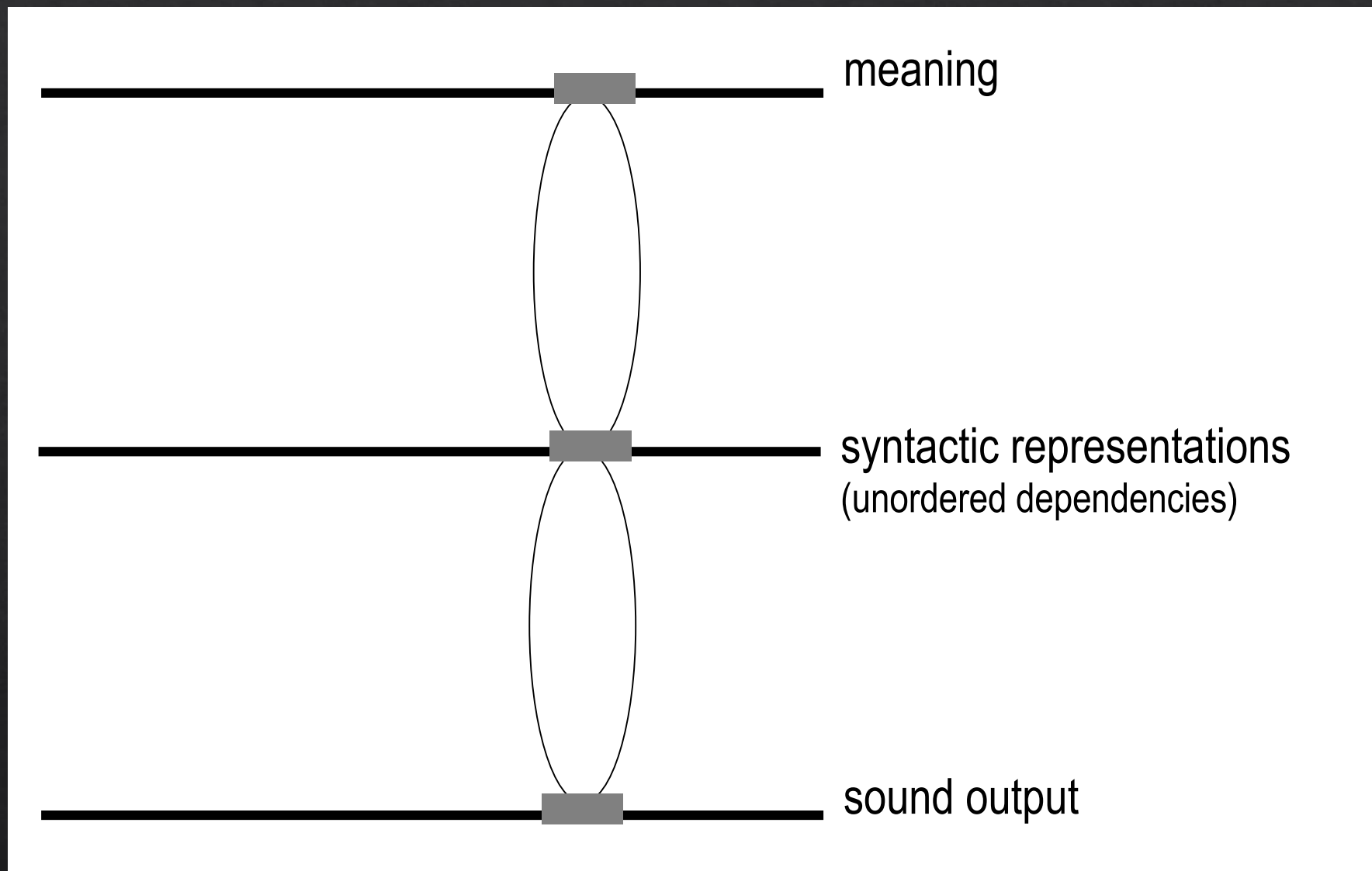
- ◊ Why different language descriptions and analysis
- ◊ Texts
- ◊ Treebanks
- ◊ Networks
- ◊ Questions & comments

# Outline

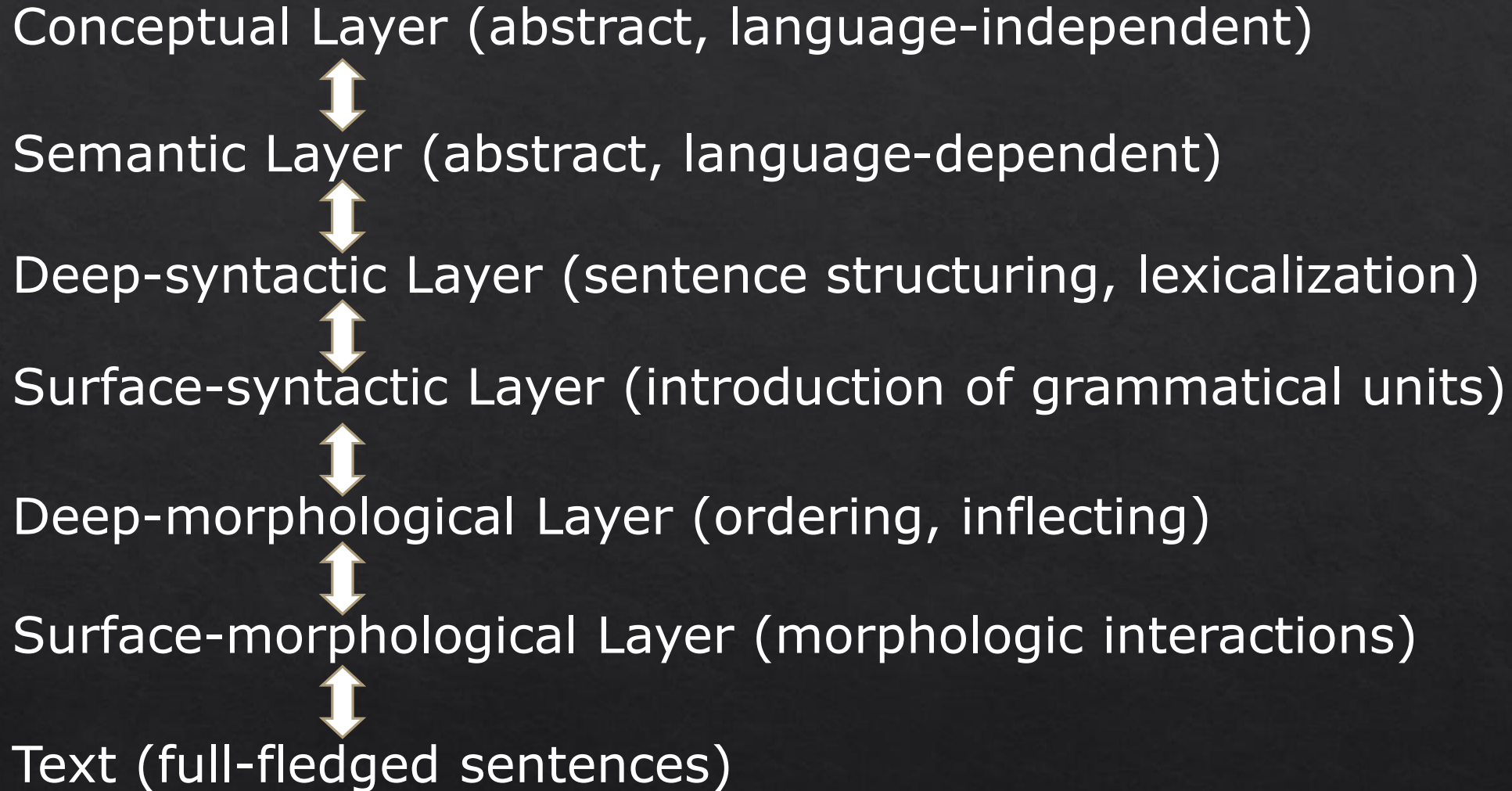
- ◇ Why different language descriptions and analysis

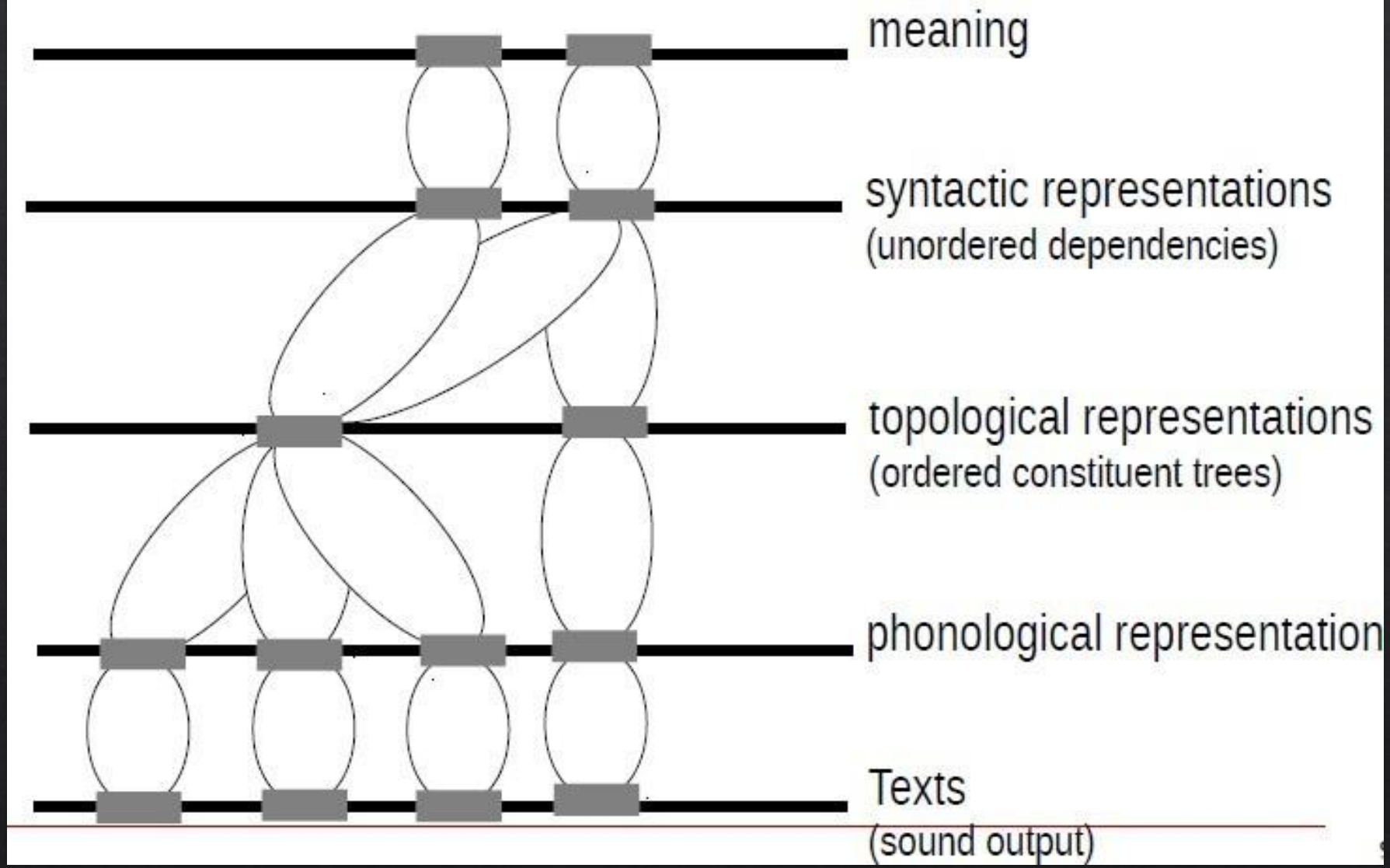
language is a system of signs, each of which is an arbitrary union of sound and meaning... Any given sign, is defined by its relationships with the others. (Saussure 1959)





Saussure's signified and signifier (2011)  
Mel'čuk's meaning and text (1981)  
Chomsky's logical and phonetic structure (2002)





surface-deep hierarchy Mel'čuk 1986

Language understanding and generation require deconstructions or constructions of units on different language layers. Different language descriptions, therefore, yield to this nature of language communications.

# Outline

- ◊ Why different language descriptions and analysis
- ◊ Texts



# Texts

- ◆ You may have heard
  - ◆ word, token, type, lemma, phrase, ngram, sentence, paragraph...
  - ◆ word length, frequency, size...
  - ◆ genre, style, authorship...

# Texts

- ◆ You may have heard
  - ◆ word, token, type, lemma, phrase, ngram, sentence, paragraph...
  - ◆ word length, frequency, size...
  - ◆ genre, style, authorship...
- ◆ **Advantages**
  - ◆ Relatively objective (word vs phrase...)
  - ◆ More available (OCR, Google books...)
  - ◆ Easier to test hypothesis (Zip'f law...)

# Texts

- ◆ You may have heard
  - ◆ word, token, type, lemma, phrase, ngram, sentence, paragraph...
  - ◆ word length, frequency, size...
  - ◆ genre, style, authorship...
- ◆ Advantages
  - ◆ Relatively objective (word vs phrase...)
  - ◆ More available (OCR, Google books...)
  - ◆ Easier to test hypothesis (Zip'f law...)
- ◆ Disadvantages
  - ◆ One dimensional description (linear modals)
  - ◆ Missing details of language understanding or generation

# Structural Complexity of Chinese Characters and Zipf's Law

# Structural Complexity of Chinese Characters and Zipf's Law

## ◇ Hypothesis

- ◇ The relationship between the frequency of Chinese characters and structural complexity of characters should be in accordance with Zipf's law due to the principle of least effort



# Structural Complexity of Chinese Characters and Zipf's Law

## ◇ Hypothesis

- ◇ The relationship between the **frequency** of Chinese characters and **structural complexity** of characters should be in accordance with Zipf's law due to the principle of least effort

# Structural Complexity of Chinese Characters and Zipf's Law

- ◇ Hypothesis

- ◇ The relationship between the frequency of Chinese characters and structural complexity of characters should be in accordance with Zipf's law due to the principle of least effort

- ◇ Frequency

- ◇ Chinese Characters' Frequency Dictionary

- ◇ the most frequent 3,061 different Chinese characters and their frequency in People's Daily, a famous newspaper in China, Corpus. (99.43%)

# Structural Complexity of Chinese Characters and Zipf's Law

## ◇ Hypothesis

- ◇ The relationship between the frequency of Chinese characters and structural complexity of characters should be in accordance with Zipf's law due to the principle of least effort

## ◇ Frequency

### ◇ Chinese Characters' Frequency Dictionary

- ◇ the most frequent 3,061 different Chinese characters and their frequency in People's Daily, a famous newspaper in China, Corpus. (99.43%)

## ◇ Structural Complexity

### ◇ Number of strokes

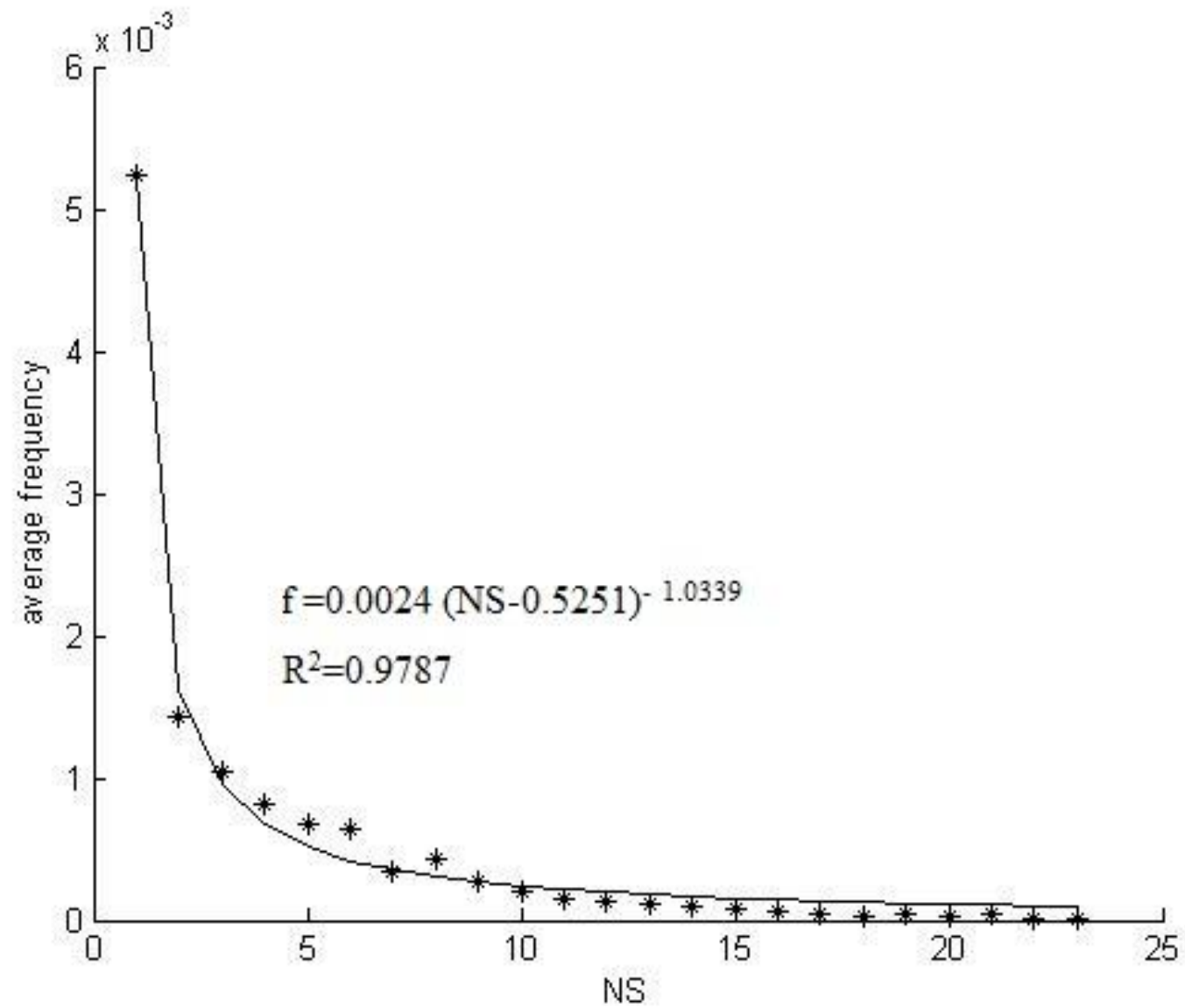
### ◇ Number of components

- ◇ Both counted based on dictionaries and issued formal documents.

睡

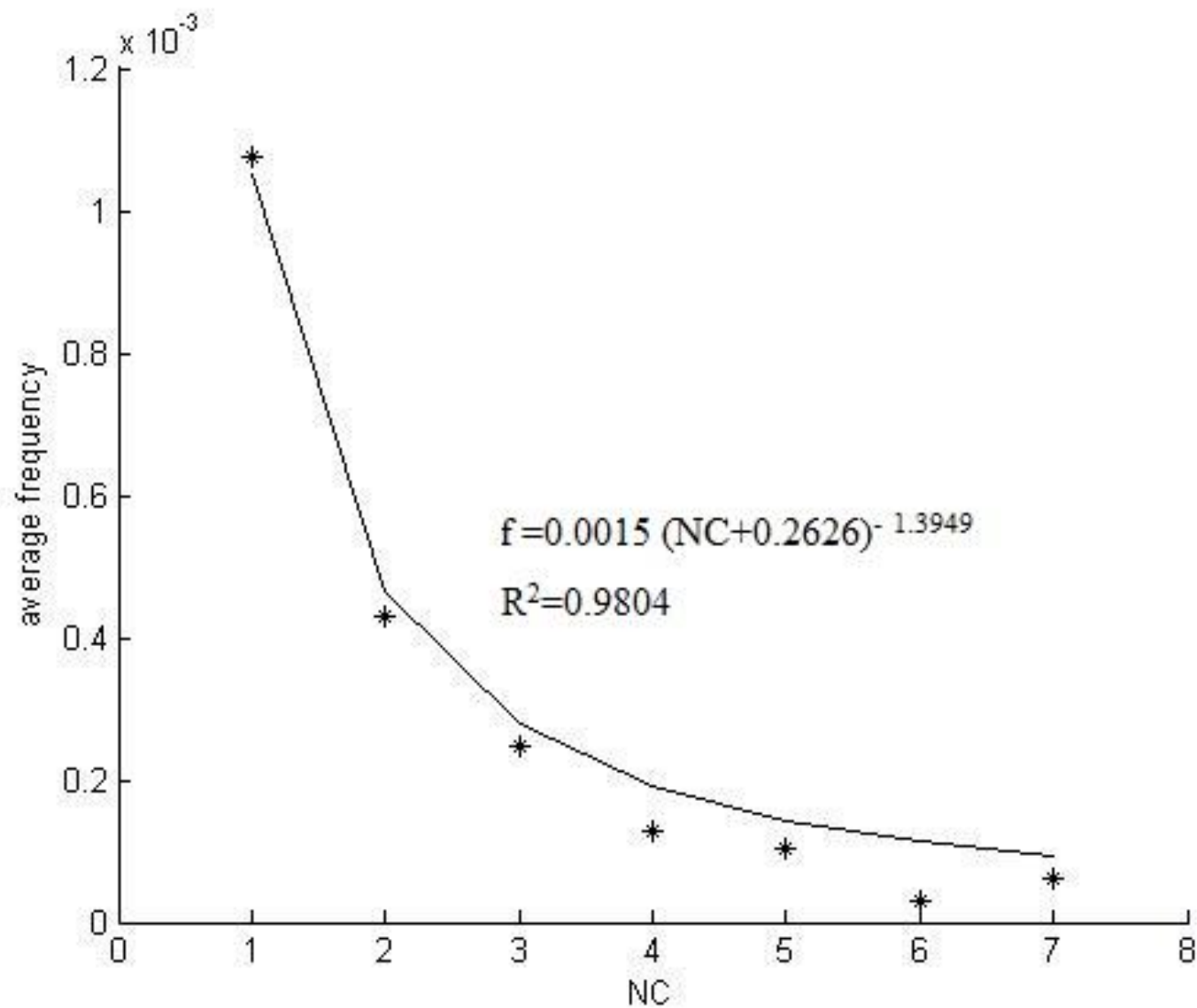
Number of Strokes	Number of Chracters	Examples	Cumulative frequency	Average frequency
1	2	一,乙	0.0104818889	0.0052409444
2	19	人,了,力	0.0270823814	0.0014253885
3	51	大,上,工	0.0535283278	0.0010495751
4	113	中,不,为	0.0924586324	0.0008182180
5	145	发,业,民	0.0989700908	0.0006825524
6	237	在,有,年	0.1514154083	0.0006388836
7	307	作,这,来	0.1035224004	0.0003372065
8	379	的,国,和	0.1611448260	0.0004251842
9	368	是,要,政	0.1006057580	0.0002733852
10	343	家,部,展	0.0675444718	0.0001969227
11	290	理,得,基	0.0437326729	0.0001508023
12	276	就,等,提	0.0366750136	0.0001328805
13	177	新,意,解	0.0215748566	0.0001218918
14	124	道,管,赛	0.0111661176	0.0000900493
15	100	题,增,德	0.0084463093	0.0000844631
16	55	整, 融,器	0.0033965338	0.0000617552
17	33	藏,繁,疑	0.0012677566	0.0000384169
18	11	翻,覆,藤	0.0002244518	0.0000204047
19	13	警,疆,攀	0.0005720697	0.0000440054
20	11	籍,灌,耀	0.0002580577	0.0000234598
21	4	露,霸,髓	0.0001694056	0.0000423514
22	2	囊,镶	0.0000231615	0.0000115808
23	1	罐	0.0000151569	0.0000151569





The relationship between NS and the average frequency of Chinese characters that share the same NS

Frequency of components				
Number of Components	Number of Characters	Examples	Cumulative frequency	Average frequency
1	187	一,中,人	0.2008274347	0.0010739435
2	969	的,国,和	0.4164280467	0.0004297503
3	1186	在,是,发	0.2922298614	0.0002463995
4	534	高,能,说	0.0682191084	0.0001277511
5	150	题,领,解	0.0153812537	0.0001025417
6	32	歌,疑,衡	0.0010030980	0.0000313468
7	3	疆,凝,颤	0.0001869466	0.0000623155



The relationship between NC and the average frequency of Chinese characters that share the same NC

# Diachronic observation of Chinese Words Based on Google Ngram

# Diachronic observation of Chinese Words Based on Google Ngram

- ◇ Observation
  - ◇ the lexicon changes of a language during a time period
- ◇ Hypothesis verification
  - ◇ There is a multi-syllabification trend of Chinese words during Chinese evolution (Chen, H., Liang, J., & Liu, H. 2015 PLoS ONE)



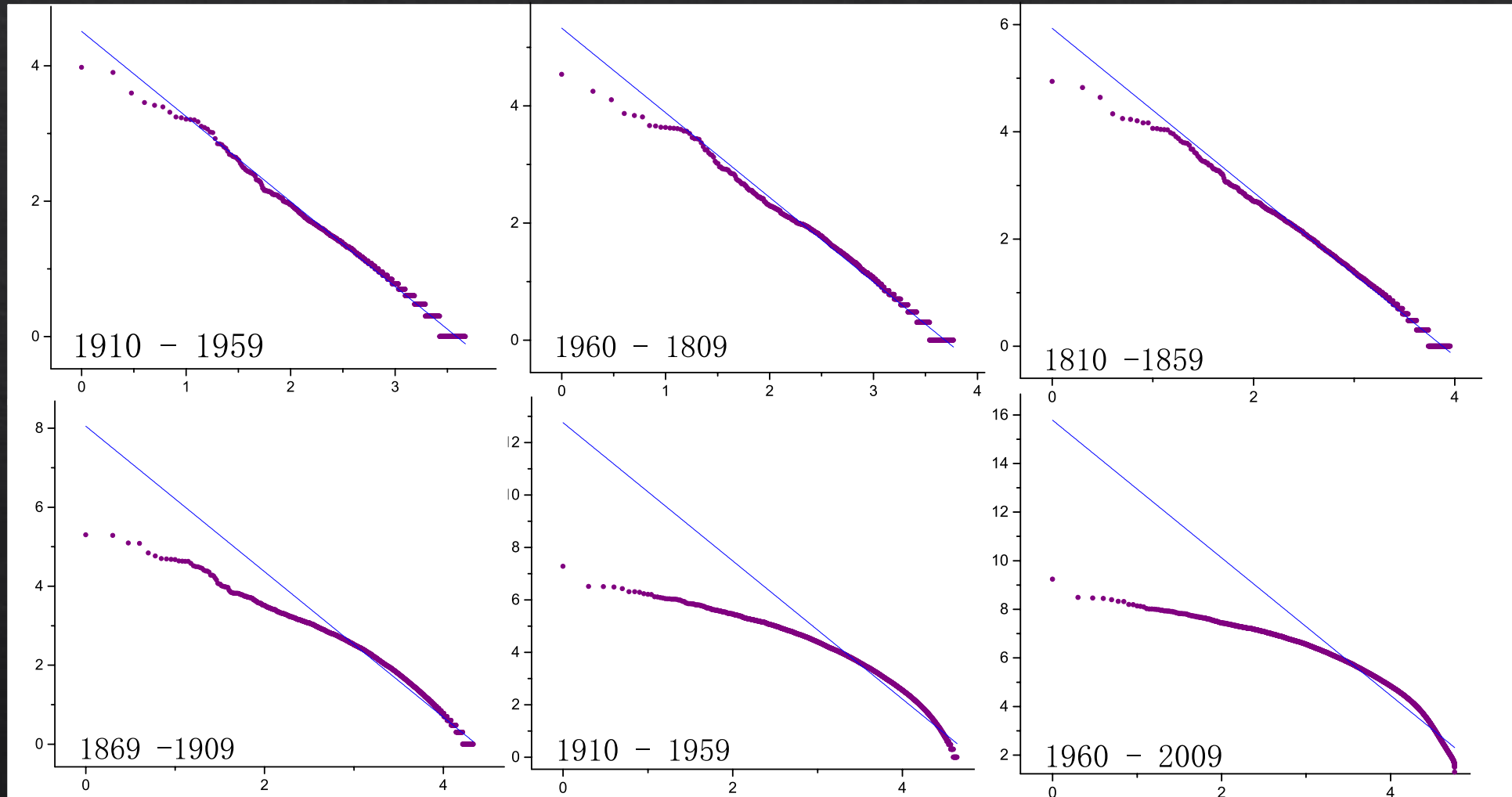
# Diachronic observation of Chinese Words Based on Google Ngram

- ◇ Observation
  - ◇ the lexicon changes of a language during a time period
- ◇ Hypothesis verification
  - ◇ There is a multi-syllabification trend of Chinese words during Chinese evolution (Chen, H., Liang, J., & Liu, H. 2015 PLoS ONE)
- ◇ The Chinese Google 1-gram data (segmented-words frequency)
  - ◇ Symbols, Japanese, Roman letters, Arabic numerals
  - ◇ 1710 ~ 2009 (300 years)
    - ◇ × Before 1710 (sparse, segmentation accuracy)
    - ◇ 6 \* 50 (years)
    - ◇ 1710~59, 1760~1809, 1810~59, 1869~1909, 1910~59, and 1960~2009

# Diachronic observation of Chinese Words Based on Google Ngram

- ◇ Observation
  - ◇ the lexicon changes of a language during a time period
- ◇ Hypothesis verification
  - ◇ There is a multi-syllabification trend of Chinese words during Chinese evolution (Chen, H., Liang, J., & Liu, H. 2015 PLoS ONE)
- ◇ The Chinese Google 1-gram data (segmented-words frequency)
  - ◇ Symbols, Japanese, Roman letters, Arabic numerals
  - ◇ 1710 ~ 2009 (300 years)
    - ◇ × Before 1710 (sparse, segmentation accuracy)
    - ◇ 6 \* 50 (years)
      - ◇ 1710~59, 1760~1809, 1810~59, 1869~1909, 1910~59, and 1960~2009
- ◇ Power law fitting

# Power Law Fitting



# Power Law Fitting: $y = a * x^{-b}$

Time Period	a	b	R <sup>2</sup>
1760-1809	34,791	-0.98122	0.98755
1810-1859	96,530	-0.94701	0.96039
1860-1909	249,946	-0.82656	0.92599
1910-1959	16,579,000	-1.09064	0.91207
1960-2009	1,512,270,000	-1.05617	0.90595

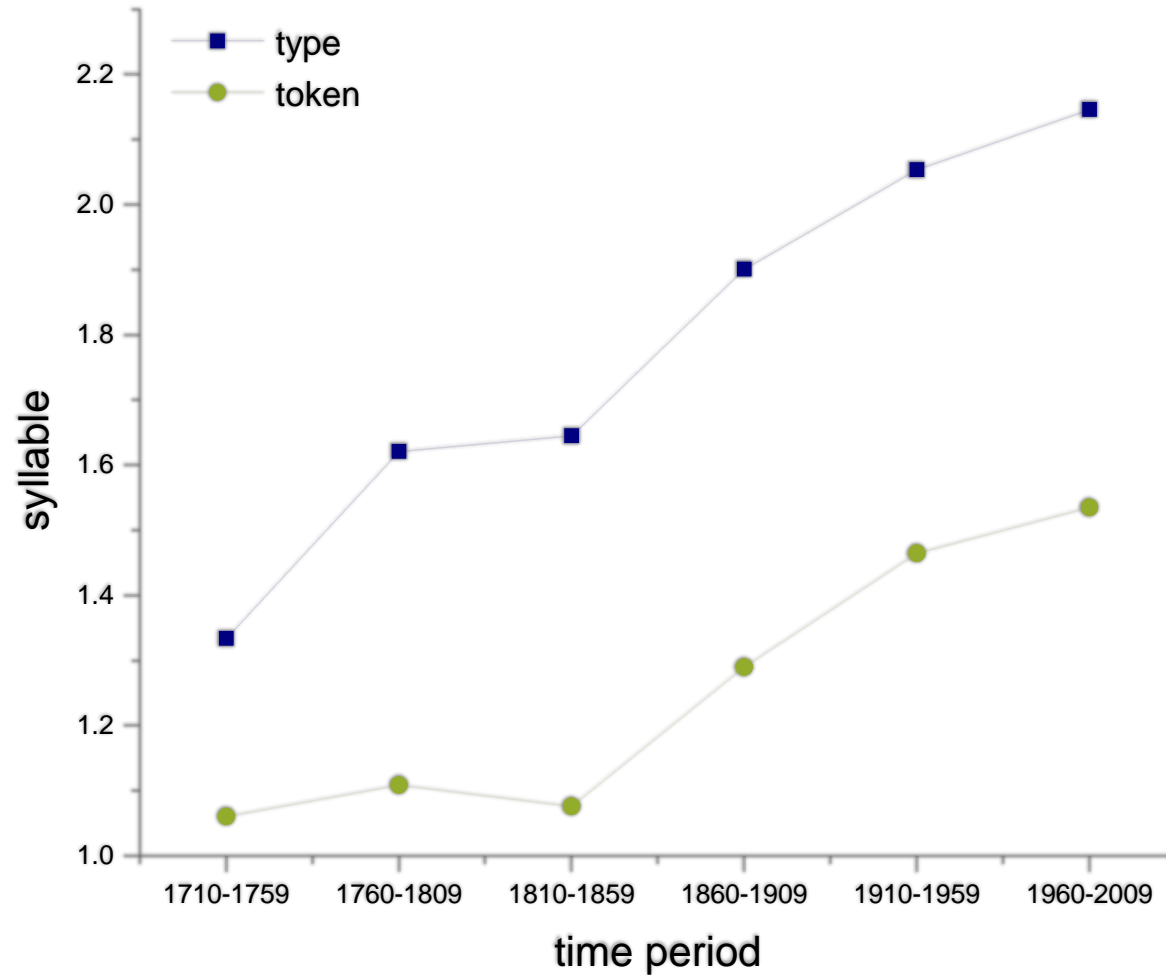
- ♦ all the six frequency lists fit the Power Law well (the minimum R<sup>2</sup> is 0.90595) despite the huge data scale (the largest list includes more than 21 billion tokens)
- ♦ The Google Ngram data should be suitable for doing some linguistic studies

# Diachronic observation of Chinese Words Based on Google Ngram

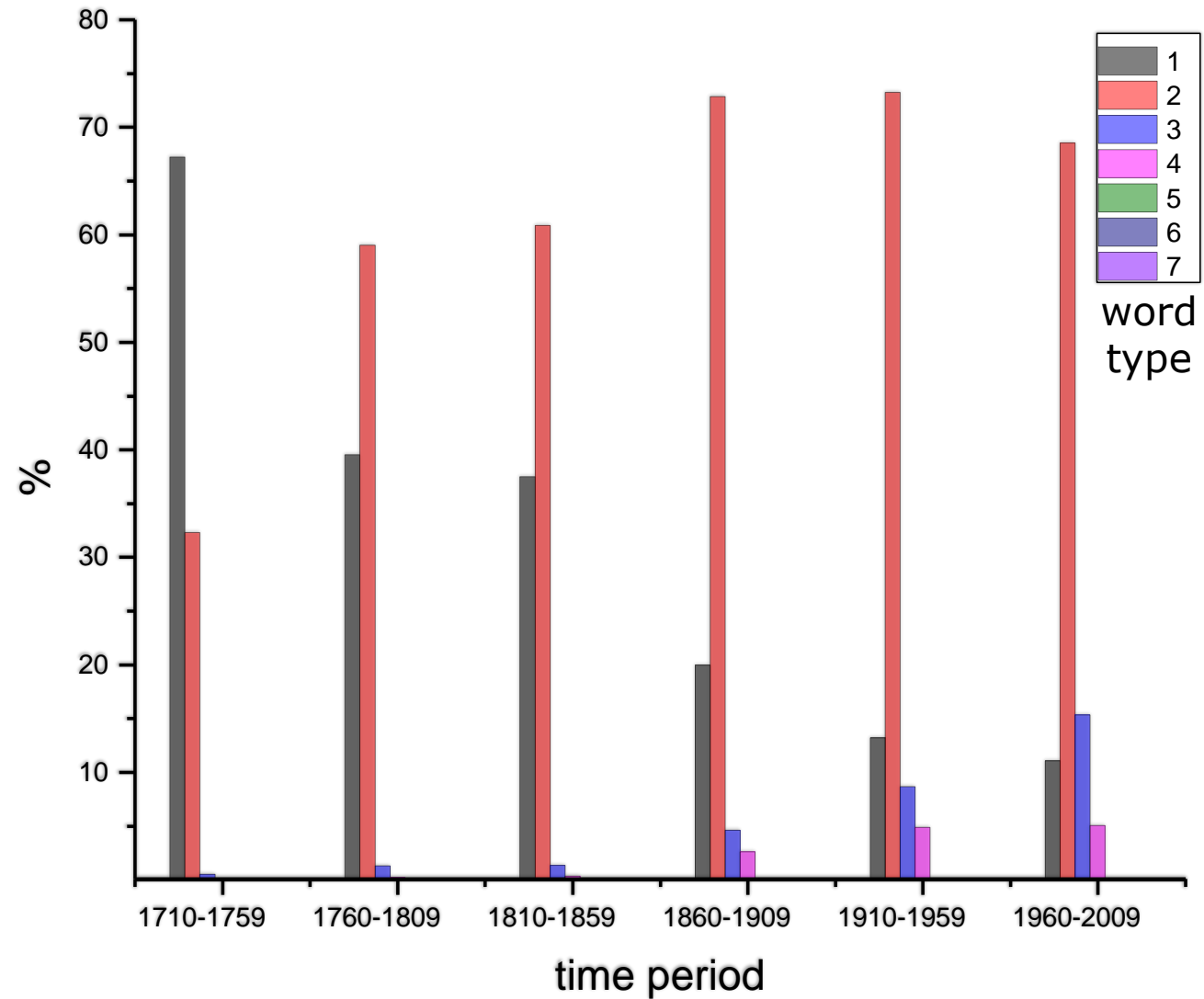
- ◇ Observation
  - ◇ the lexicon changes of a language during a time period
- ◇ Hypothesis verification
  - ◇ There is a **multi-syllabification** trend of Chinese words during Chinese evolution (Chen, H., Liang, J., & Liu, H. 2015 PLoS ONE)
- ◇ The Chinese Google 1-gram data (segmented-words frequency)
  - ◇ Symbols, Japanese, Roman letters, Arabic numerals
  - ◇ 1710 ~ 2009 (300 years)
    - ◇ × Before 1710 (sparse, segmentation accuracy)
    - ◇ 6 \* 50 (years)
    - ◇ 1710~59, 1760~1809, 1810~59, 1869~1909, 1910~59, and 1960~2009
- ◇ Power law fitting
- ◇ **Average word length**



# Average word length



# Word length distribution



# Outline

- ◇ Why different language descriptions and analysis
- ◇ Texts
- ◇ Treebanks

# Treebanks

- ◆ You may have heard
  - ◆ Noam Chomsky
  - ◆ Grammar, syntax, semantics...
  - ◆ Predicate, agent, patient, subject, object, modifier, valence...
  - ◆ Machine translation, natural language processing, natural language understanding...

# Treebanks

- ◆ You may have heard
  - ◆ Noam Chomsky
  - ◆ Grammar, syntax, semantics...
  - ◆ Predicate, agent, patient, subject, object, modifier, valence...
  - ◆ Machine translation, natural language processing, natural language understanding...
- ◆ **Advantages**
  - ◆ Two dimensional description (non-linear modals)
  - ◆ More detailed information of language understanding or generation
  - ◆ Practical uses for language teaching and language processing

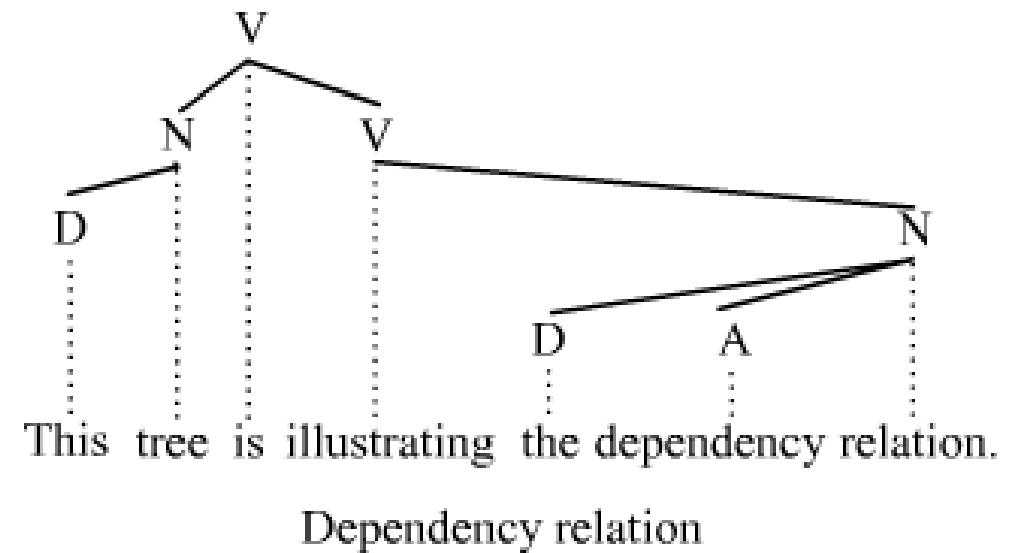
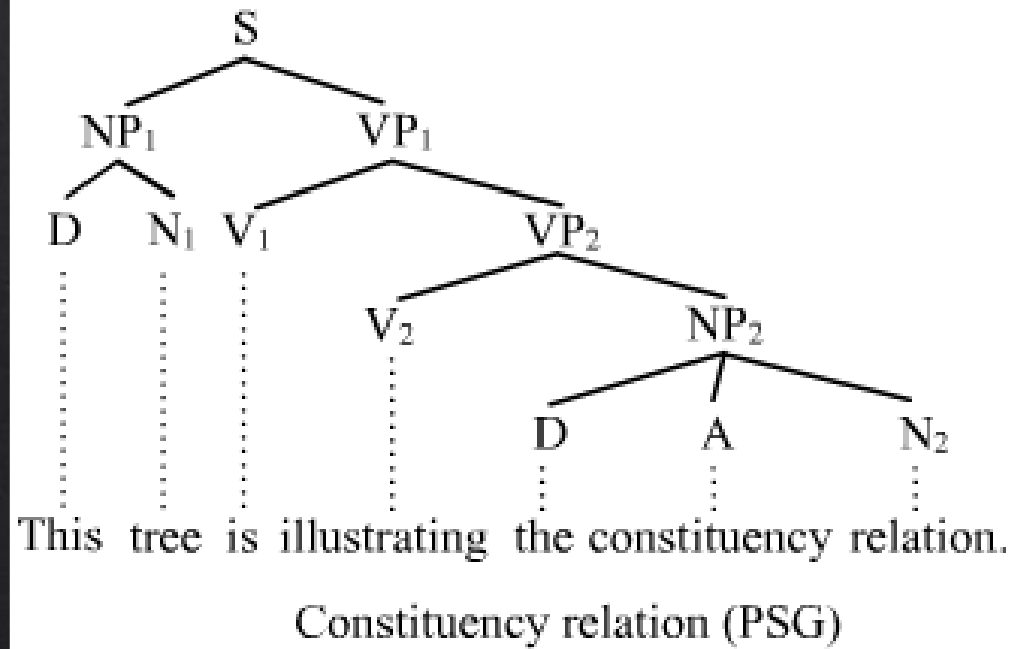


# Treebanks

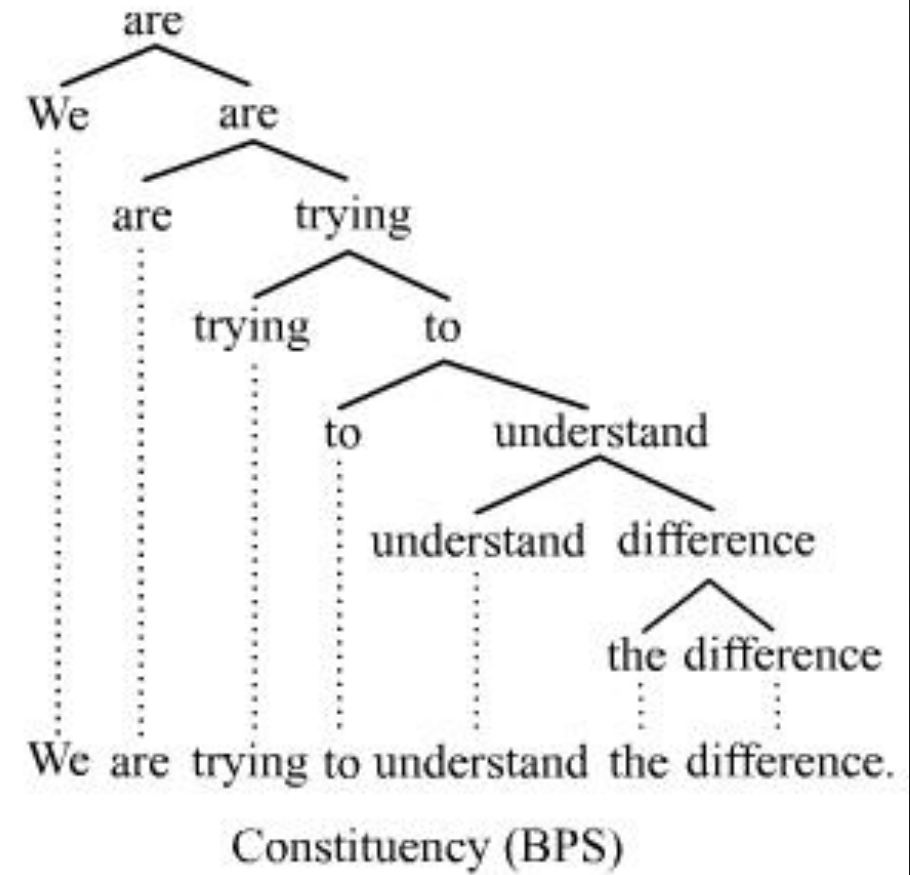
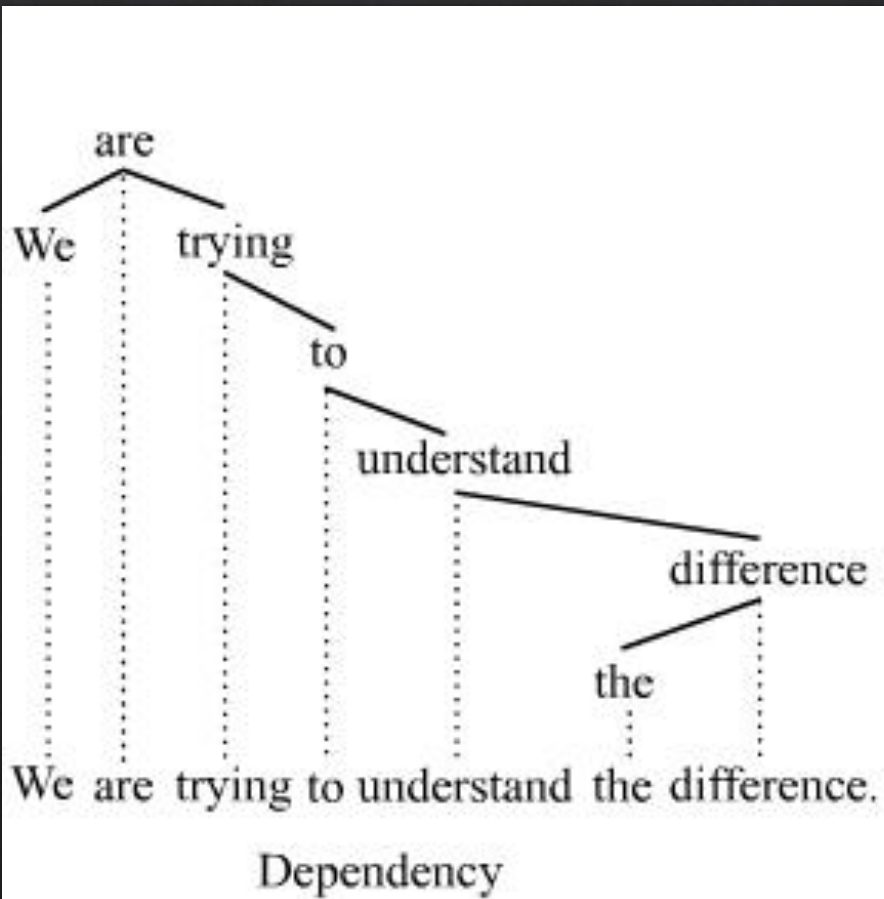
- ◆ You may have heard
  - ◆ Noam Chomsky
  - ◆ Grammar, syntax, semantics...
  - ◆ Predicate, agent, patient, subject, object, modifier, valence...
  - ◆ Machine translation, natural language processing, natural language understanding...
- ◆ Advantages
  - ◆ Two dimensional description (non-linear modals)
  - ◆ More detailed information of language understanding or generation
  - ◆ Practical uses for language teaching and language processing
- ◆ Disadvantages
  - ◆ Disagreements on description frameworks (less objective)
  - ◆ Less available
  - ◆ More difficult to test hypothesis (Null models/hypothesis)

# Treebanks

- ◊ You may have heard
  - ◊ Noam Chomsky
  - ◊ Grammar, syntax, semantics...
  - ◊ Predicate, agent, patient, subject, object, modifier, valence...
  - ◊ Machine translation, natural language processing, natural language understanding...
- ◊ Advantages
  - ◊ Two dimensional description (non-linear modals)
  - ◊ More detailed information of language understanding or generation
  - ◊ Practical uses for language teaching and language processing
- ◊ Disadvantages
  - ◊ Disagreements on description frameworks (less objective)
  - ◊ Less available
  - ◊ More difficult to test hypothesis (Null models/hypothesis)
- ◊ What is treebank

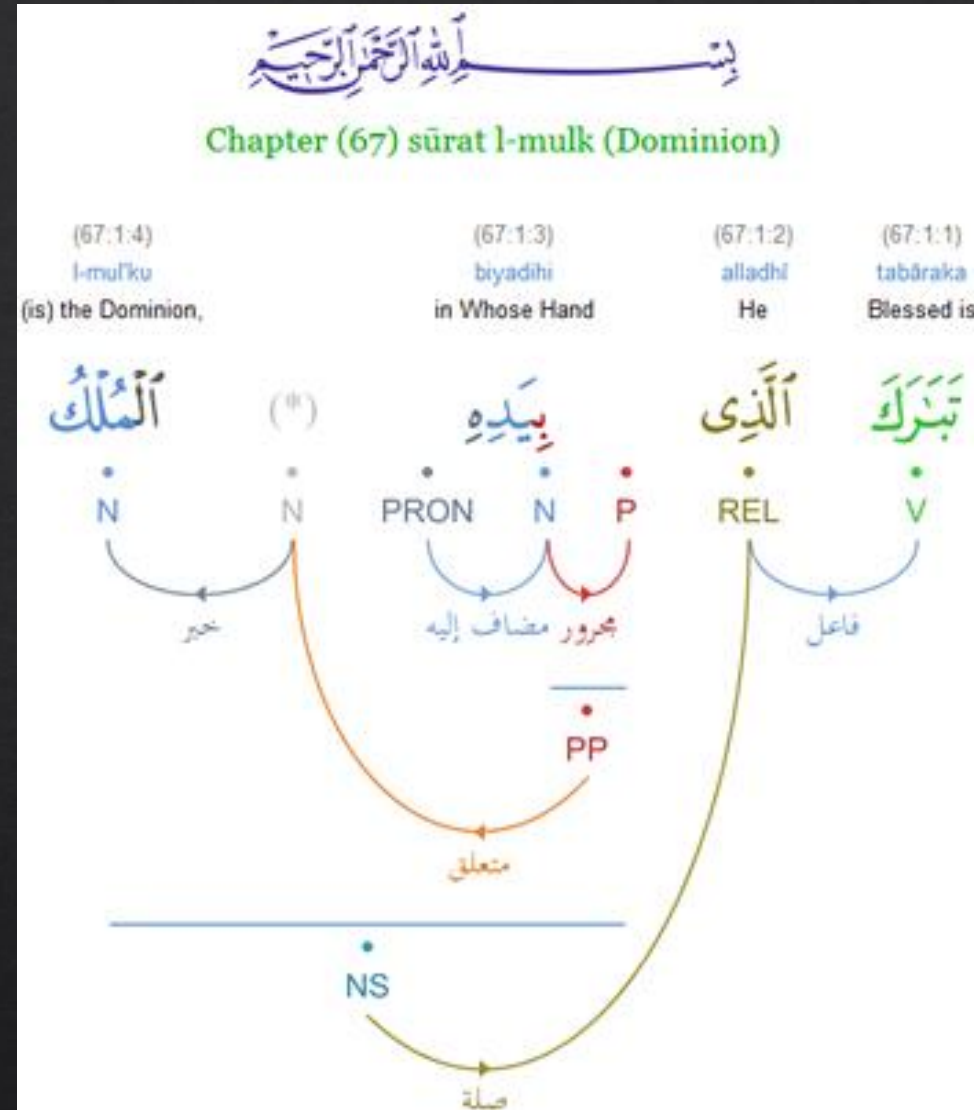


[http://en.wikipedia.org/wiki/Phrase\\_structure\\_grammar](http://en.wikipedia.org/wiki/Phrase_structure_grammar)



[http://en.wikipedia.org/wiki/Dependency\\_grammar](http://en.wikipedia.org/wiki/Dependency_grammar)

Hybrid  
constituency/dependency  
tree from the Quranic  
Arabic Corpus



<http://en.wikipedia.org/wiki/Treebank>



1	Pricing	pricing	NN	—	2	NMOD
2	details	detail	NNS	—	3	SBJ
3	were	be	VBD	—	0	ROOT
4	n't	not	RB	—	3	ADV
5	immediately	immediately	RB	—	6	AMOD
6	available	available	JJ	—	3	PRD
7	.	.	.	—	3	P

1	She	she	PRP	—	2	SBJ
2	bought	buy	VBD	—	0	ROOT
3	a	a	DT	—	4	NMOD
4	car	car	NN	—	2	OBJ
5	.	.	.	—	2	P

## CoNLL-X

## Other format (xml)

```
<?xml-stylesheet type="text/xsl" href="nlp_style.xsl" ?>↓
<xml4nlp>↓
  <doc>↓
    <para id="0">↓
      <sent id="0" cont="迈向充满希望的新世纪——一九九八年新年讲话（附图
片1张）">↓
        <word id="0" cont="迈向" pos="v" parent="-1" relate="HED" />↓
        <word id="1" cont="充满" pos="v" parent="3" relate="DE" />↓
        <word id="2" cont="希望" pos="n" parent="1" relate="VOB" />↓
        <word id="3" cont="的" pos="u" parent="5" relate="ATT" />↓
        <word id="4" cont="新" pos="a" parent="5" relate="ATT" />↓
        <word id="5" cont="世纪" pos="n" parent="0" relate="VOB" />↓
        <word id="6" cont="——" pos="wp" parent="-2" relate="PUN" />↓
        <word id="7" cont="一九九八年" pos="nt" parent="8" relate="ATT"
```

<https://code.google.com/p/clearparser/wiki/DataFormat>

# Classifying Languages by Dependency Structure

## Typologies of Delexicalized Universal Dependency Treebanks

# Classifying Languages by Dependency Structure

## Typologies of Delexicalized Universal Dependency Treebanks

- ◇ Observation
  - ◇ Syntactical differences between different languages
- ◇ Question
  - ◇ Can UD be used for language typology study and reveal the similarity and diversity between language families? If so, then how?

# Classifying Languages by Dependency Structure

## Typologies of Delexicalized Universal Dependency Treebanks

- ◇ Observation
  - ◇ Syntactical differences between different languages
- ◇ Question
  - ◇ Can UD be used for language typology study and reveal the similarity and diversity between language families? If so, then how?
- ◇ **UD treebanks**
  - ◇ UD 2.0 70 treebanks of 50 languages
  - ◇ 63 of which have more than 10,000 tokens.



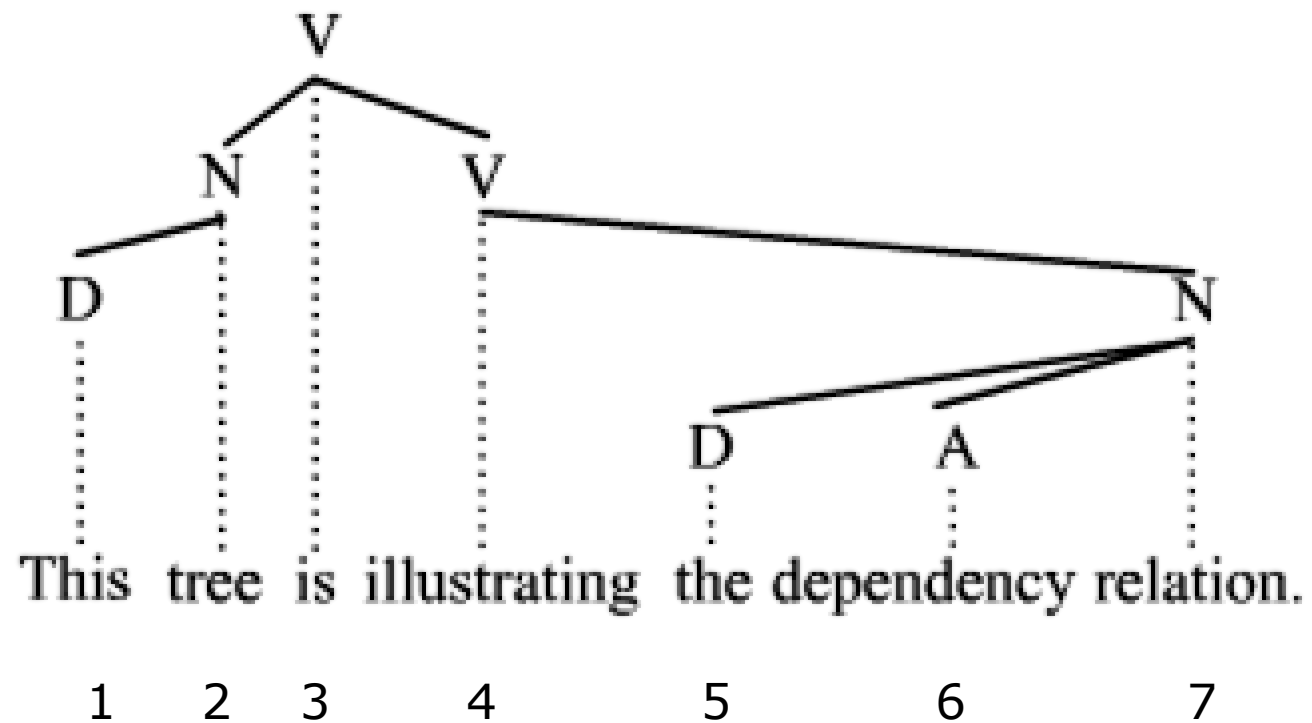
# Classifying Languages by Dependency Structure

## Typologies of Delexicalized Universal Dependency Treebanks

- ◇ Observation
  - ◇ Syntactical differences between different languages
- ◇ Question
  - ◇ Can UD be used for language typology study and reveal the similarity and diversity between language families? If so, then how?
- ◇ UD treebanks
  - ◇ UD 2.0 70 treebanks of 50 languages
  - ◇ 63 of which have more than 10,000 tokens.
- ◇ **Dependency distance & word order (head initial, head final, mixed)**
- ◇ Distributions of frequencies
- ◇ Clustering algorithm



# Dependency Length



$$DL(\text{this-is}) = 3 - 1 = 2$$

$$DL(\text{is-illustrating}) = 3 - 4 = -1$$

# Dependency Length

$$DDD(R) = \frac{\sum_{r \in R} distance(r)}{frequency(R)}$$

R: type of dependencies    r: dependencies

# Classifying Languages by Dependency Structure

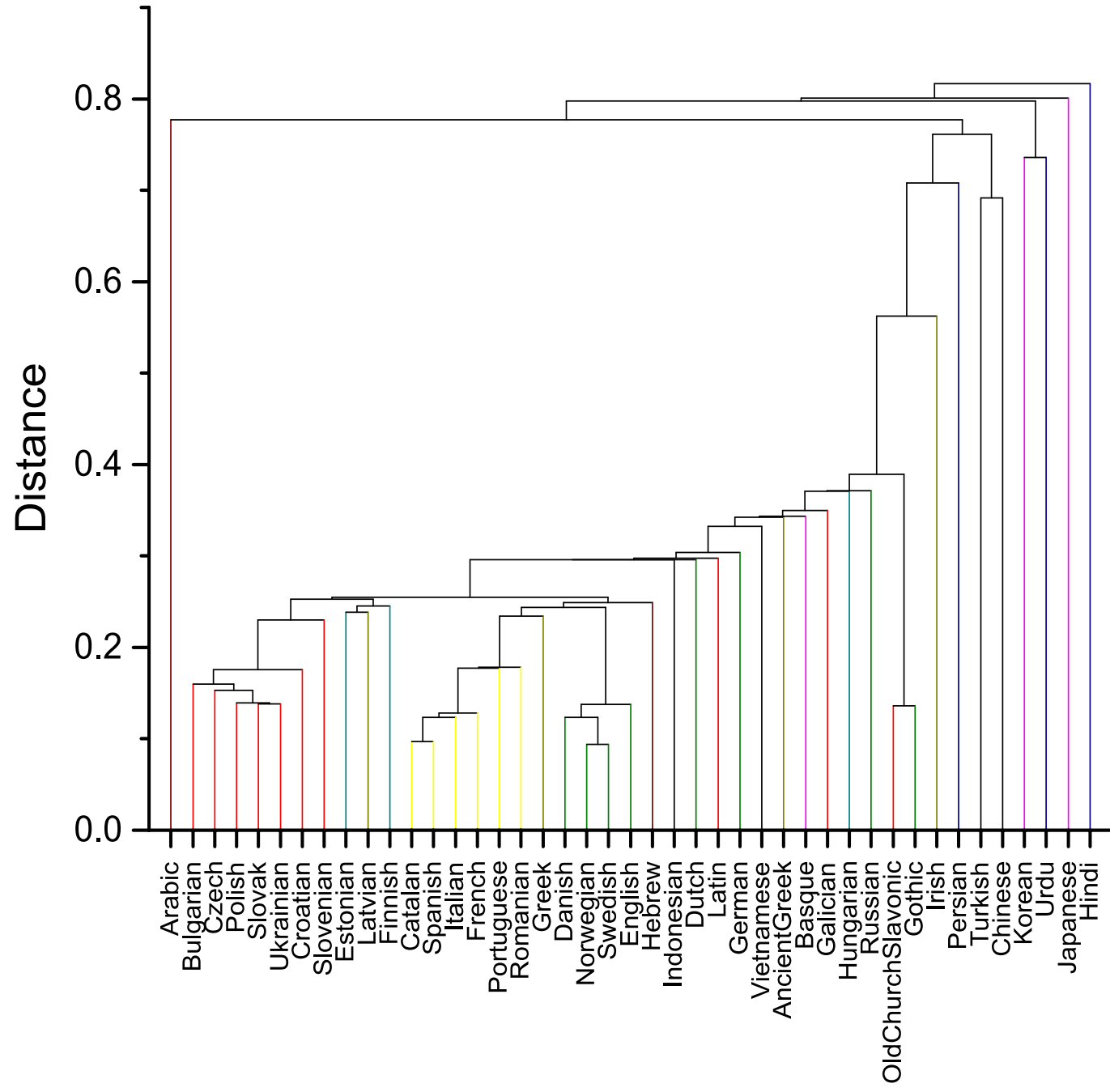
## Typologies of Delexicalized Universal Dependency Treebanks

- ◇ Observation
  - ◇ Syntactical differences between different languages
- ◇ Question
  - ◇ Can UD be used for language typology study and reveal the similarity and diversity between language families? If so, then how?
- ◇ UD treebanks
  - ◇ UD 2.0 70 treebanks of 50 languages
  - ◇ 63 of which have more than 10,000 tokens.
- ◇ Dependency distance & word order (head initial, head final, mixed)
- ◇ **Distributions of frequencies**

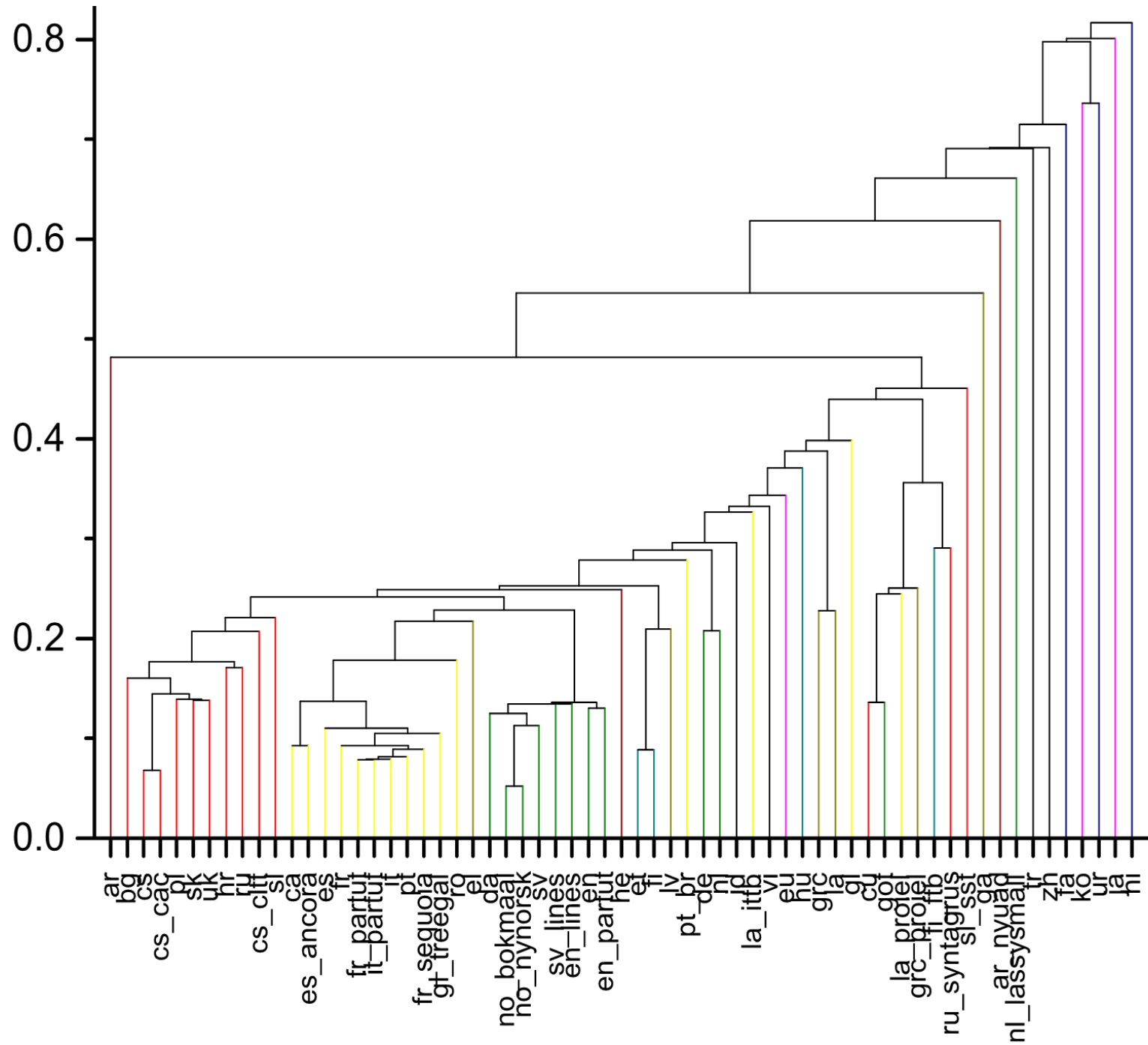
# Classifying Languages by Dependency Structure

## Typologies of Delexicalized Universal Dependency Treebanks

- ◇ Observation
  - ◇ Syntactical differences between different languages
- ◇ Question
  - ◇ Can UD be used for language typology study and reveal the similarity and diversity between language families? If so, then how?
- ◇ UD treebanks
  - ◇ UD 2.0 70 treebanks of 50 languages
  - ◇ 63 of which have more than 10,000 tokens.
- ◇ Dependency distance & word order (head initial, head final, mixed)
- ◇ Distributions of frequencies
- ◇ Clustering algorithm







# Outline

- ◊ Why different language descriptions and analysis
- ◊ Texts
- ◊ Treebanks
- ◊ **Networks**

# Networks

- ◆ You may have heard
  - ◆ six degrees of separation, social network, small world...
  - ◆ graph theory, seven bridges, complex theory...

# Networks

- ◊ You may have heard
  - ◊ six degrees of separation, social network, small world...
  - ◊ graph theory, seven bridges, complex theory...
- ◊ **Advantages**
  - ◊ Three dimensional description
  - ◊ Macro system perspective (break the sentence boundaries)

# Networks

- ◆ You may have heard
  - ◆ six degrees of separation, social network, small world...
  - ◆ graph theory, seven bridges, complex theory...
- ◆ Advantages
  - ◆ Three dimensional description
  - ◆ Macro system perspective (break the sentence boundaries)
- ◆ Disadvantages
  - ◆ Too abstract to be understood or interpreted?
  - ◆ More difficult to test hypothesis (Null models/hypothesis)
  - ◆ Disappeared details?



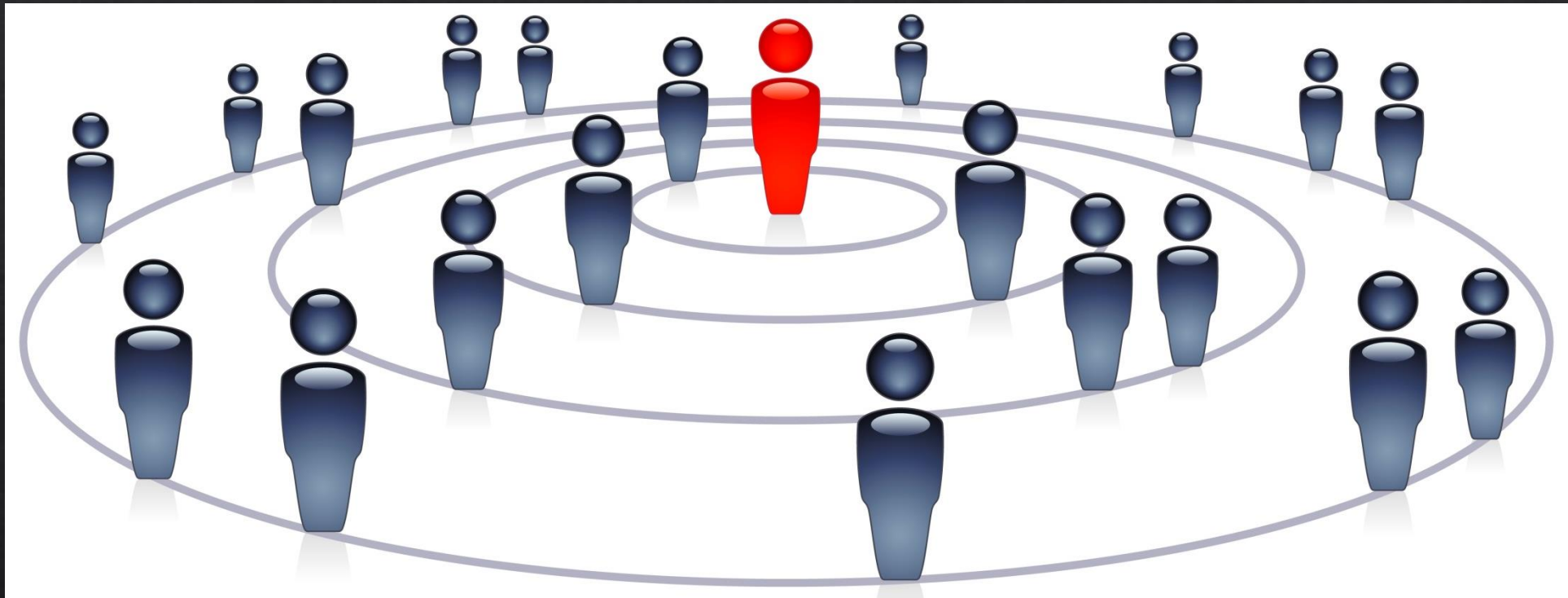
# Networks

- ◆ You may have heard
  - ◆ six degrees of separation, social network, small world...
  - ◆ graph theory, seven bridges, complex theory...
- ◆ Advantages
  - ◆ Three dimensional description
  - ◆ Macro system perspective (break the sentence boundaries)
- ◆ Disadvantages
  - ◆ Too abstract to be understood or interpreted?
  - ◆ More difficult to test hypothesis (Null models/hypothesis)
  - ◆ Disappeared details?
- ◆ What is network

# Origin and development

# Origin and development

## ◇ Six Degrees of Separation



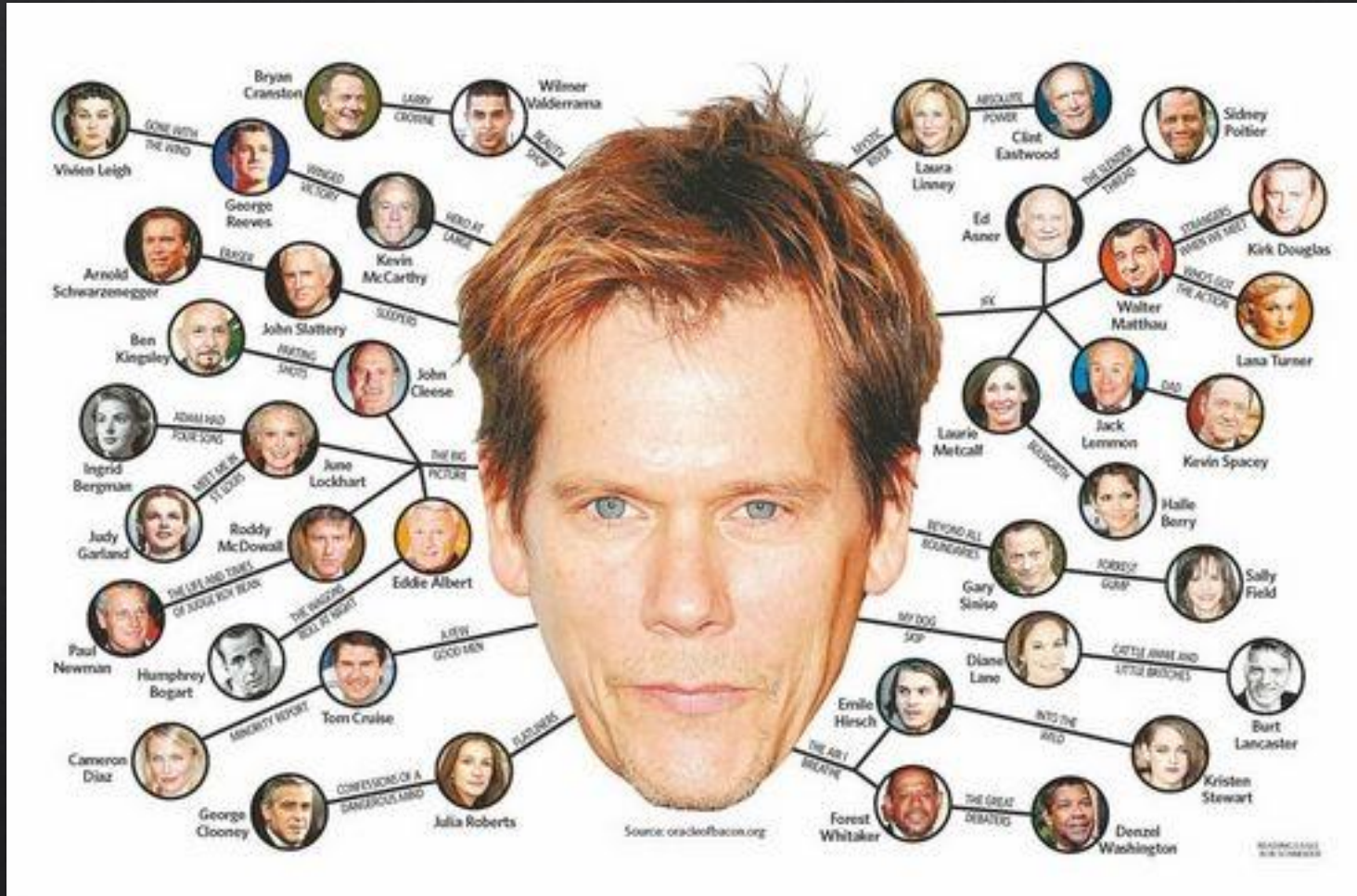
<http://indianvox.com/articles/news/143667.php>

# Origin and development

- ◇ Six Degrees of Separation
  - ◇ Harvard psychologist Stanley Milgram 1967
    - ◇ U.S. / 300 letters / >60 / 6



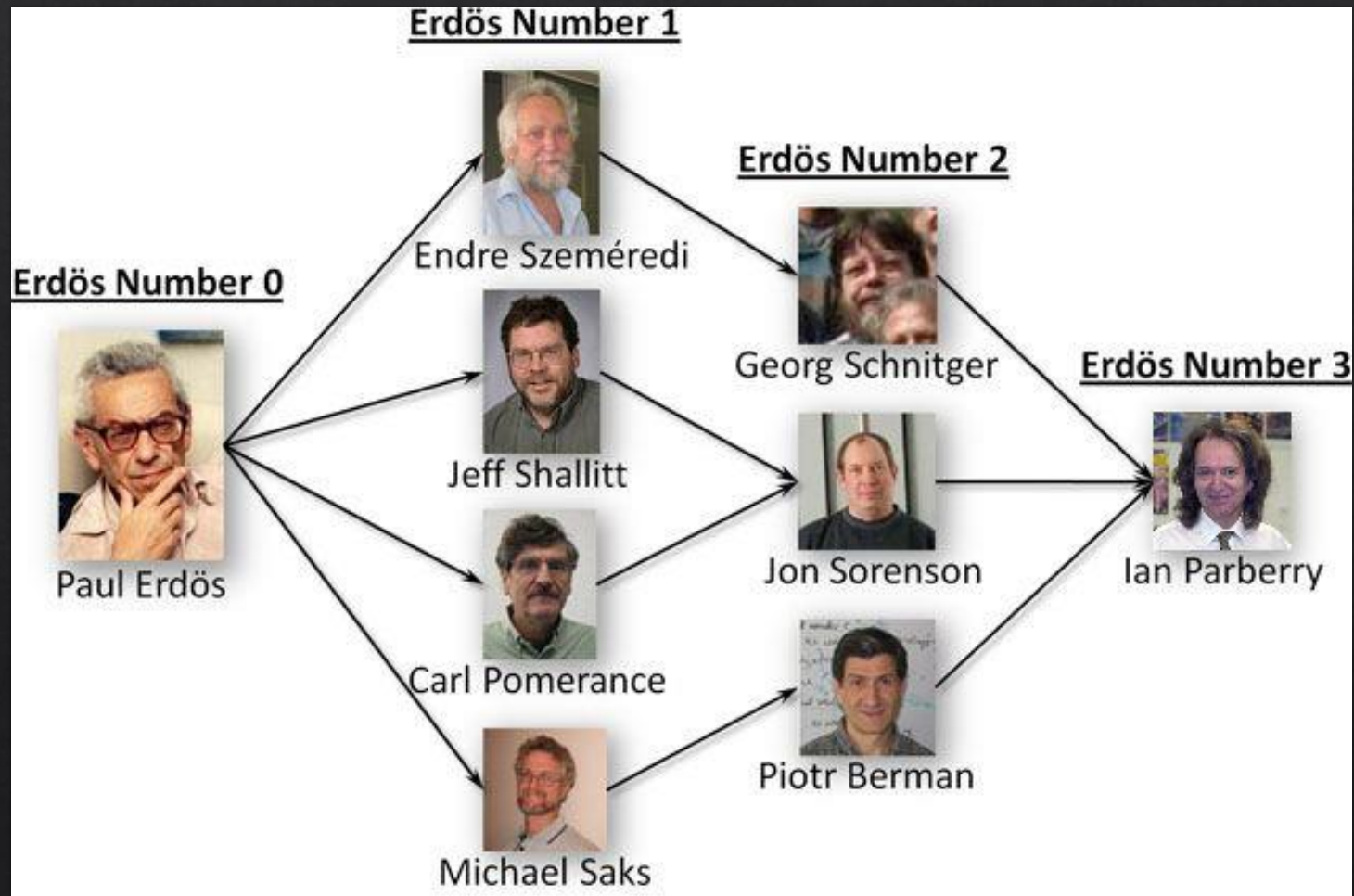
# Bacon number



<http://badgreedraza.blogspot.jp/2015/04/6-degrees-of-kevin-bacon.html>

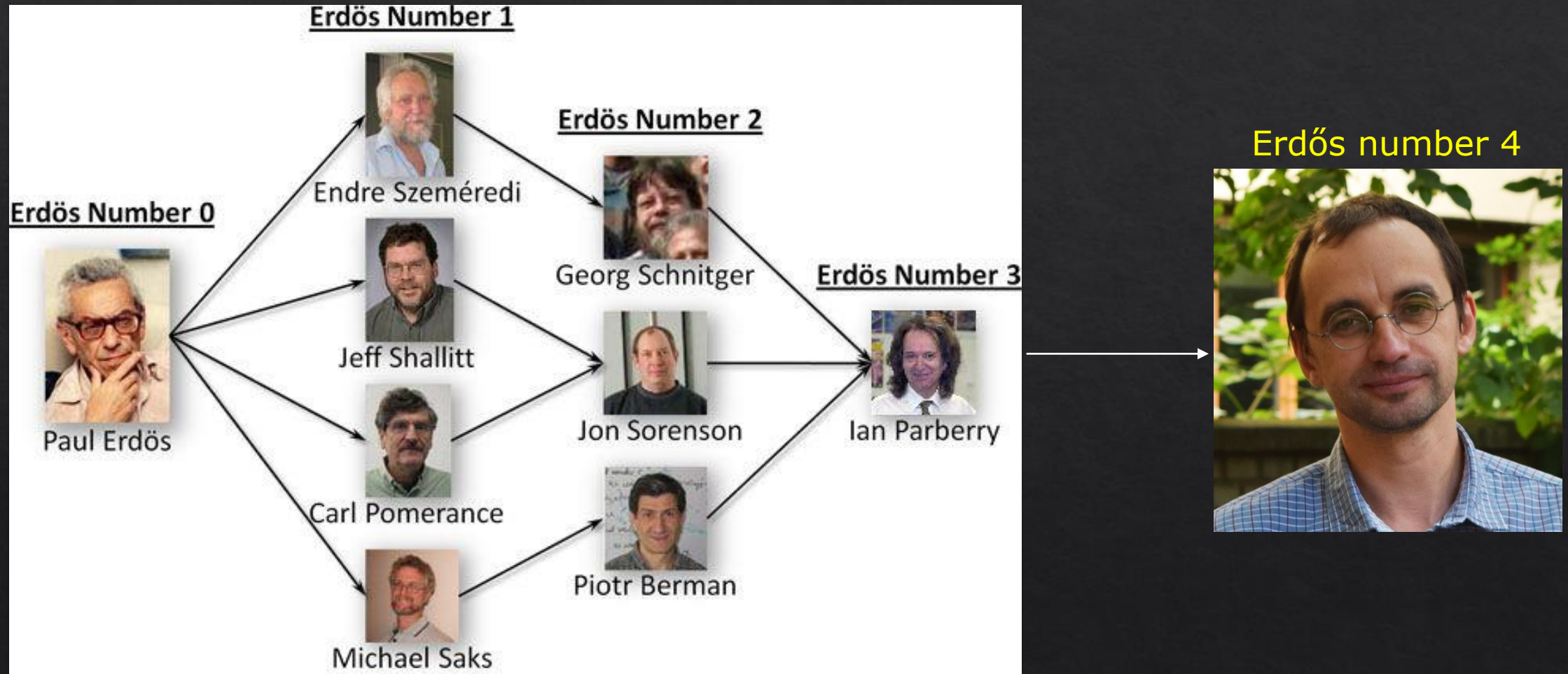


# Erdős number



<https://larc.unt.edu/ian/claimtofame.html>

# Erdős number



<https://larc.unt.edu/ian/claimtofame.html>

# Origin and development

- ◆ Six Degrees of Separation
  - ◆ Harvard psychologist Stanley Milgram 1967
    - ◆ U.S. / 300 letters / >60 / 6
  - ◆ Columbia professor Duncan Watts 2001
    - ◆ 19 email targets (157 countries) / 48,000 senders / 6

[https://en.wikipedia.org/wiki/Six\\_degrees\\_of\\_separation](https://en.wikipedia.org/wiki/Six_degrees_of_separation)



# Origin and development

- ◆ Six Degrees of Separation
  - ◆ Harvard psychologist Stanley Milgram 1967
    - ◆ U.S. / 300 letters / >60 / 6
  - ◆ Columbia professor Duncan Watts 2001
    - ◆ 19 email targets (157 countries) / 48,000 senders / 6
  - ◆ Microsoft researchers Jure Leskovec and Eric Horvitz 2007
    - ◆ 30 billion Skype messages / 240 million people / 6

[https://en.wikipedia.org/wiki/Six\\_degrees\\_of\\_separation](https://en.wikipedia.org/wiki/Six_degrees_of_separation)

# Origin and development

- ◇ Six Degrees of Separation
- ◇ **Big Data Area & Quantitative Method**
  - ◇ UCINET (Borgatti et al. 2002)
  - ◇ PAJEK (Nooy et al. 2005)
  - ◇ NETDRAW (Borgatti 2002)
  - ◇ CYTOSCAPE (Shannon 2003)
  - ◇ ...



# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System

# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
  - ◇ Simplification & Decomposition

# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
- ◇ ~~Simplification & Decomposition~~

# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
  - ◇ Duncan J Watts **Small World**
  - ◇ Albert-László Barabási **Scale Free**

# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
  - ◇ Duncan J Watts **Small World**
  - ◇ Albert-László Barabási **Scale Free**
  - ◇ **Universal Features of Existing Systems**



# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
  - ◇ Universal Features of Existing Systems

# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
  - ◇ Universal Features of Existing Systems
    - ◇ Brain network, Nation network, Social network, Traffic network, Epidemic network...

# Origin and development

- ◇ Six Degrees of Separation
- ◇ Big Data Area & Quantitative Method
- ◇ Complex System
  - ◇ Universal Features of Existing Systems
    - ◇ Brain network, Nation network, Social network, Traffic network, Epidemic network...
    - ◇ Language network

# Origin and development (Language)

# Origin and development (Language)

- ◇ Language is an 'organic' complex system
  - ◇ Ramon Ferrer i Cancho & Ricard Solé



# Origin and development (Language)

- ◇ Language is an 'organic' complex system
  - ◇ Ramon Ferrer i Cancho & Ricard Solé
- ◇ Different language networks
  - ◇ Syntactic dependency networks, language development or language evolution, language clustering and linguistic categorization, manual and machine translation, word sense disambiguation, communication and interaction, semantic networks, phonetics, morphology, parts of speech, knowledge networks, cognitive networks ...

# Origin and development (Language)

- ◊ Why do we need the network approach in linguistic research?

# Origin and development (Language)

- ◇ Why do we need the network approach in linguistic research?
- ◇ **General advantages of the network approach**
  - ◇ Graphic view, intuitive and easy to understand
  - ◇ Overall perspective, not only focusing on the sentence level
  - ◇ General method, easier to communicate and to combine with other disciplines, such as neurology
  - ◇ Tools developed for graphs and other domains can be applied

# Origin and development (Language)

- ◇ Why do we need the network approach in linguistic research?
  - ◇ General advantages of the network approach
    - ◇ Graphic view, intuitive and easy to understand
    - ◇ Overall perspective, not only focusing on the sentence level
    - ◇ General method, easier to communicate and to combine with other disciplines, such as neurology
    - ◇ Tools developed for graphs and other domains can be applied
  - ◇ Limitations of current methods of treebank-based linguistics study
    - ◇ The boundary between sentences



# Origin and development (Language)

- ◇ Why do we need the network approach in linguistic research?
  - ◇ General advantages of the network approach
    - ◇ Graphic view, intuitive and easy to understand
    - ◇ Overall perspective, not only focusing on the sentence level
    - ◇ General method, easier to communicate and to combine with other disciplines, such as neurology
    - ◇ Tools developed for graphs and other domains can be applied
  - ◇ Limitations of current methods of treebank-based linguistics study
    - ◇ The boundary between sentences



1	Pricing	pricing	NN	—	2	NMOD
2	details	detail	NNS	—	3	SBJ
3	were	be	VBD	—	0	ROOT
4	n't	not	RB	—	3	ADV
5	immediately	immediately	RB	—	6	AMOD
6	available	available	JJ	—	3	PRD
7	.	.	.	—	3	P
<hr/>						
1	She	she	PRP	—	2	SBJ
2	bought	buy	VBD	—	0	ROOT
3	a	a	DT	—	4	NMOD
4	car	car	NN	—	2	OBJ
5	.	.	.	—	2	P

1	Pricing	pricing	NN	—	2	NMOD
2	details	detail	NNS	—	3	SBJ
3	were	be	VBD	—	0	ROOT
4	n't	not	RB	—	3	ADV
5	immediately	immediately	RB	—	6	AMOD
6	available	available	JJ	—	3	PRD
7	.	.	.	—	3	P
1	She	she	PRP	—	2	SBJ
2	bought	buy	VBD	—	0	ROOT
3	a	a	DT	—	4	NMOD
4	car	car	NN	—	2	OBJ
5	.	.	.	—	2	P

language is a system of signs, each of which is an arbitrary union of sound and meaning... Any given sign, is defined by its relationships with the others.  
(Saussure 1959)

1	Pricing	pricing	NN	—	2	NMOD
2	details	detail	NNS	—	3	SBJ
3	were	be	VBD	—	0	ROOT
4	n't	not	RB	—	3	ADV
5	immediately	immediately	RB	—	6	AMOD
6	available	available	JJ	—	3	PRD
7	.	.	.	—	3	P
<hr/>						
1	She	she	PRP	—	2	SBJ
2	bought	buy	VBD	—	0	ROOT
3	a	a	DT	—	4	NMOD
4	car	car	NN	—	2	OBJ
5	.	.	.	—	2	P

language is a system of signs, each of which is an arbitrary union of sound and meaning... Any given sign, is defined by its relationships with the others. (Saussure 1959)



language is a **system** of signs, each of which is an arbitrary union of sound and meaning... Any given sign, is defined by its relationships with the others. (Saussure 1959)

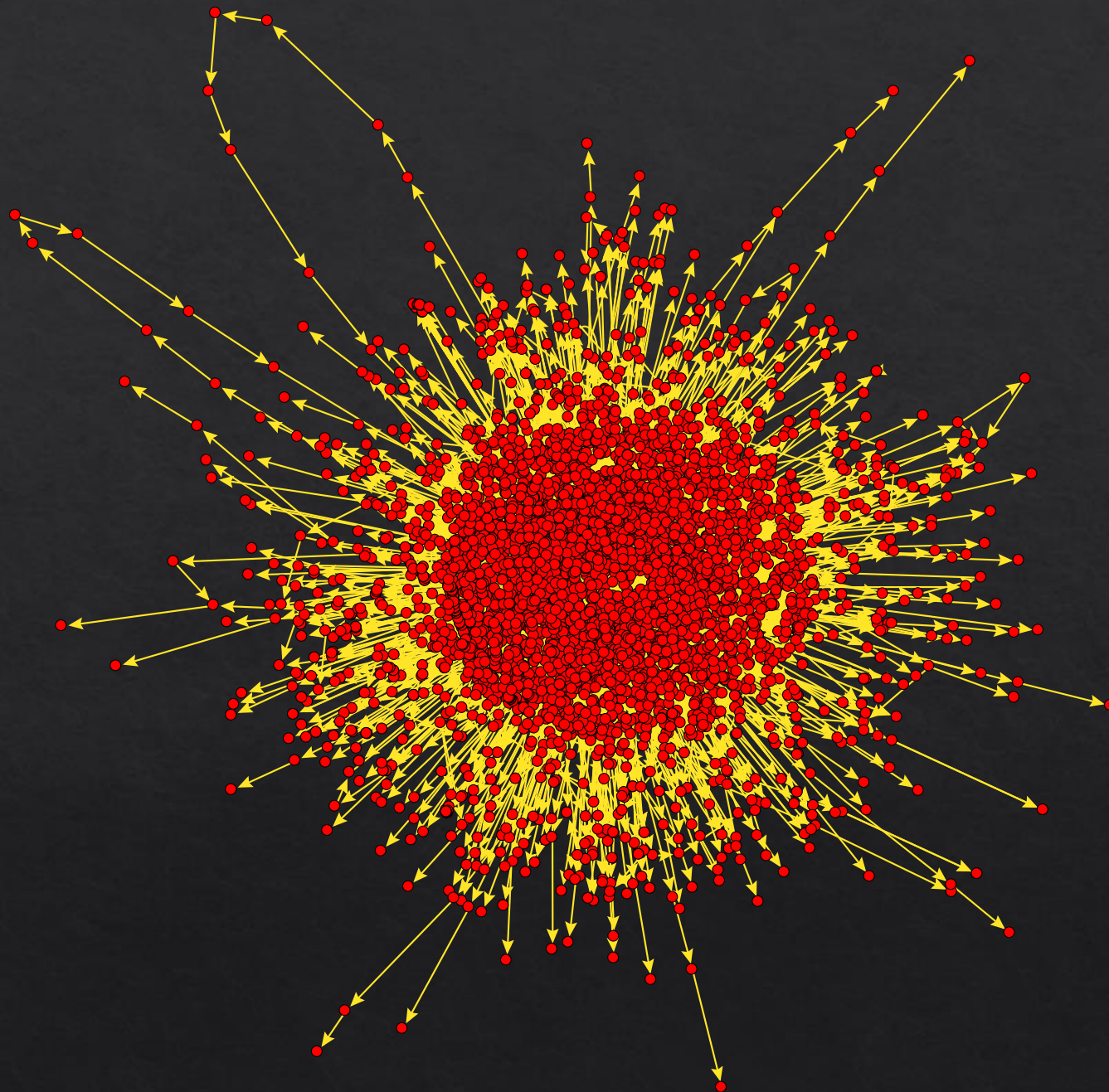


language is a **system** of signs, each of which is an arbitrary union of sound and meaning... Any given sign, is defined by its relationships with the others. (Saussure 1959)

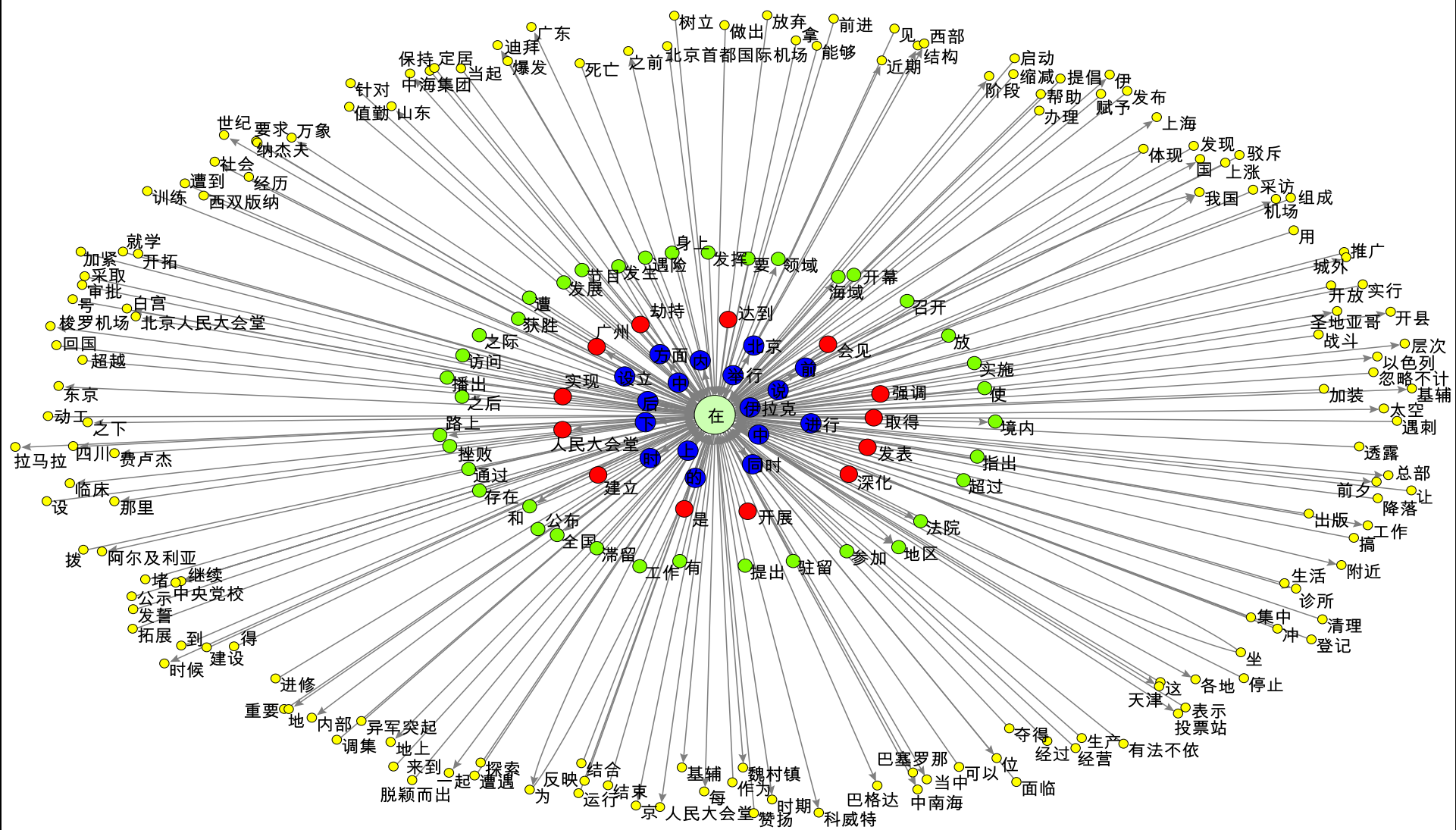
**whole  $\neq$  sum of parts**

# Origin and development (Language)

- ◇ Why do we need the network approach in linguistic research?
  - ◇ General advantages of the network approach
    - ◇ Graphic view, intuitive and easy to understand
    - ◇ Overall perspective, not only focusing on the sentence level
    - ◇ General method, easier to communicate and to combine with other disciplines, such as neurology
    - ◇ Tools developed for graphs and other domains can be applied
  - ◇ Limitations of current methods of treebank-based linguistics study
    - ◇ ~~The boundary between sentences~~

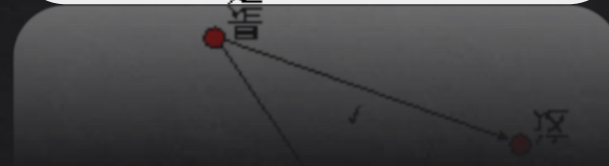
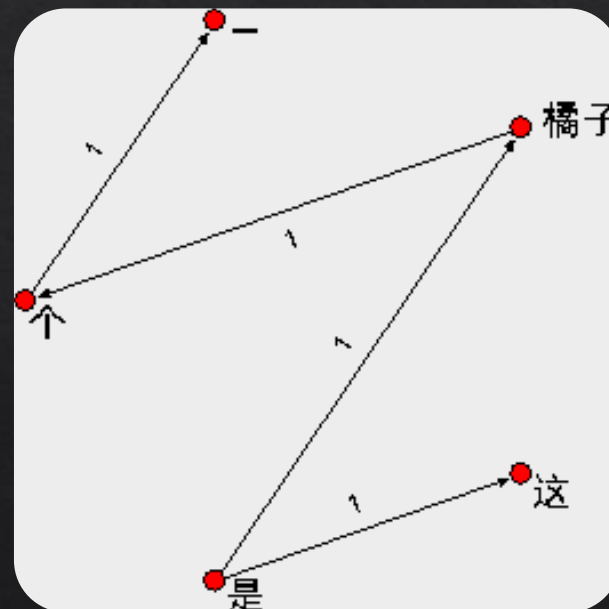
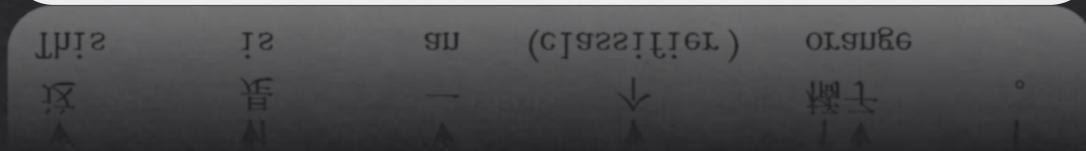
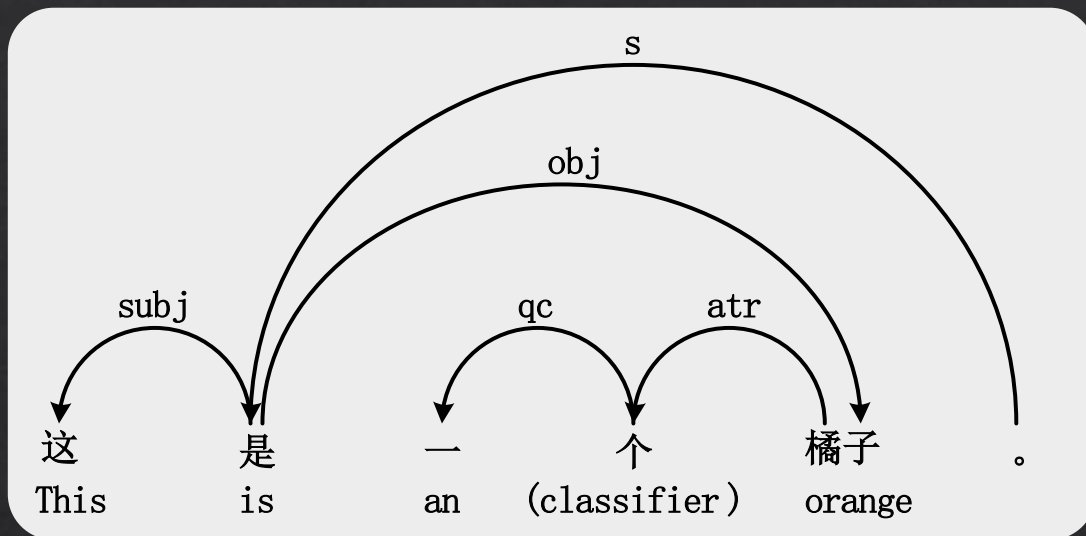






# How to Build the Networks





Annotation of a sample sentence in the **Treebank**.  
 这是一个橘子 *zhe-shi-yi-ge-ju-zi* 'this is an orange'

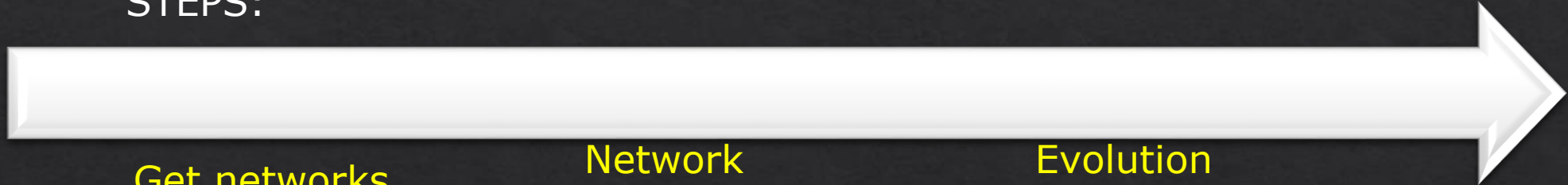
Sente nce Order	Dependent			Governor			Dependency Type
	Order	Charact er	POS	Order	Charact er	POS	
<b>S1</b>	1	这 zhe	pronoun	2	是 shi	verb	Subject
<b>S1</b>	2	是 shi	verb	6	。	punctuati on	main governor
<b>S1</b>	3	一 yi	numeral	4	个 ge	classifier	complement of classifier
<b>S1</b>	4	个 ge	classifier	5	橘子 juzi	noun	Attributer
<b>S1</b>	5	橘子 juzi	noun	2	是 shi	verb	Object
<b>S1</b>	6	。	punctuati on				



# Chinese Function Characters/Words Evolution

# Chinese Function Characters/Words Evolution

STEPS:



## Get networks

- 4 periods
  - Ancient Chinese
  - Middle Ancient Times Chinese
  - Modern Times Chinese
  - Modern Chinese
- 2 single-character words
  - 在 zai 'to exist, be living, to stay or remain, (to be located) in, at, '
  - 人 ren 'people'

## Network characteristics

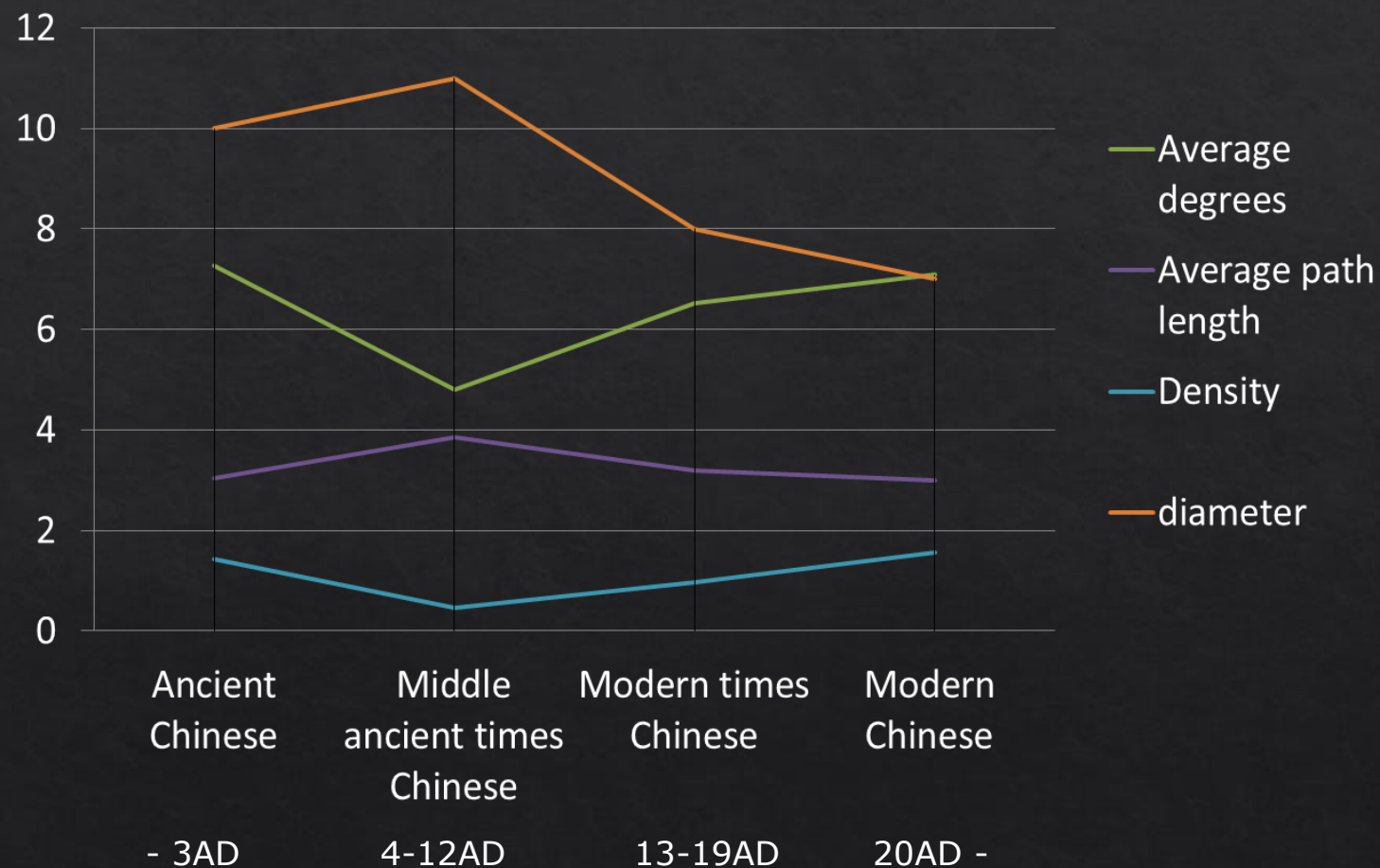
- average degrees
- average path length
- density
- diameter
- degrees of 'zai'
- degrees of 'ren'

## Evolution description

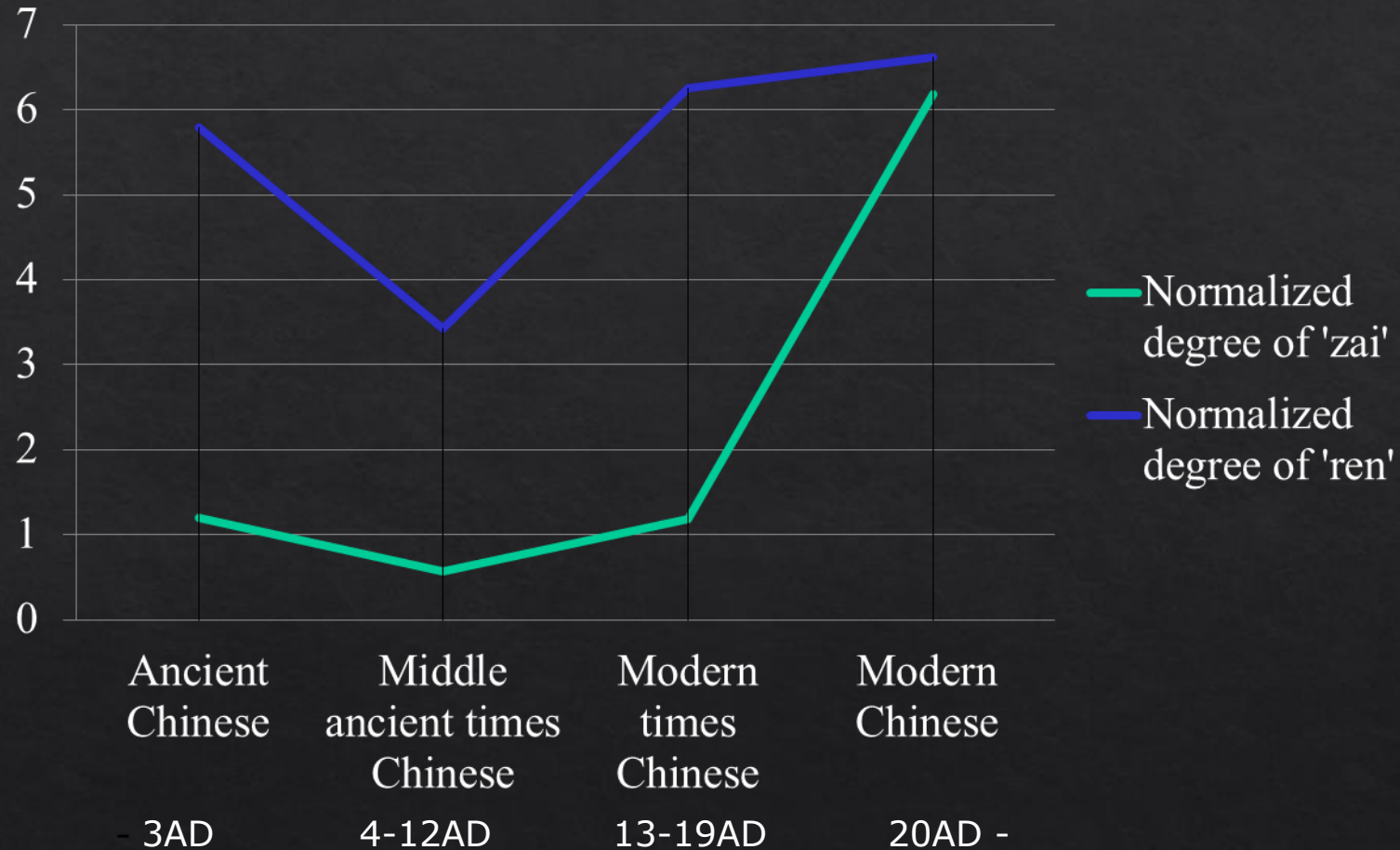
- quantitatively describe the evolution process of Chinese in 4 periods
- quantitatively describe the evolution process of 2 single-character words of Chinese in 4 periods



# Co-occurrence Character Networks



# Comparison of Single-character words



# Comparison of Single-character words

# Comparison of Single-character words

- ◆ Chinese is an isolating language: syntactic structure relies primarily on function words and word order rather than on rich morphological information to encode functional relations between elements (Levy & Manning 2003).

# Comparison of Single-character words

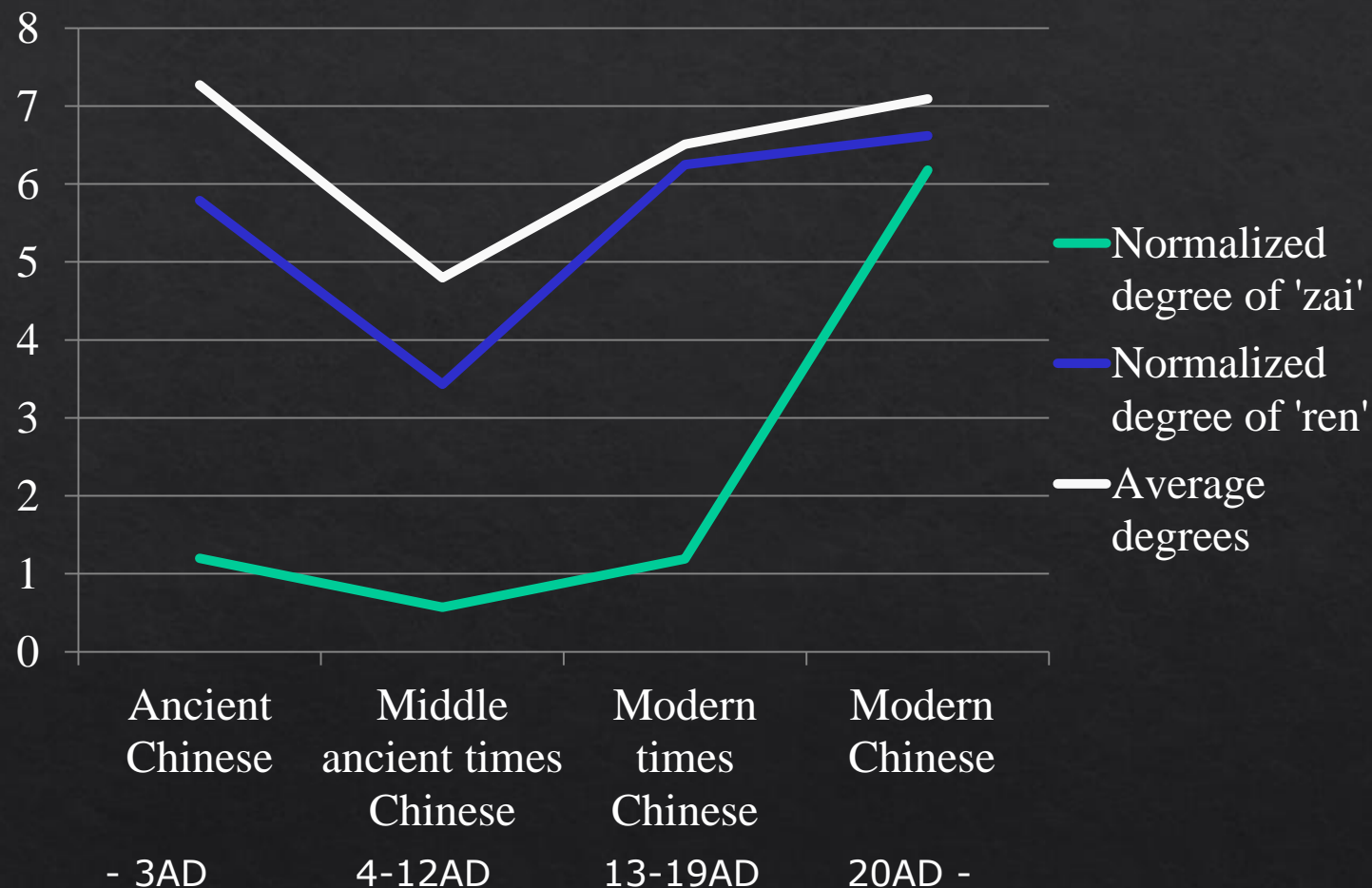
- ◇ Chinese is an isolating language: syntactic structure relies primarily on function words and word order rather than on rich morphological information to encode functional relations between elements (Levy & Manning 2003).
- ◇ 2 single-character words
  - ◇ Frequent characters/words both in the corpus and in general
  - ◇ zai underwent a grammaticalization process, whereas ren remains a content word



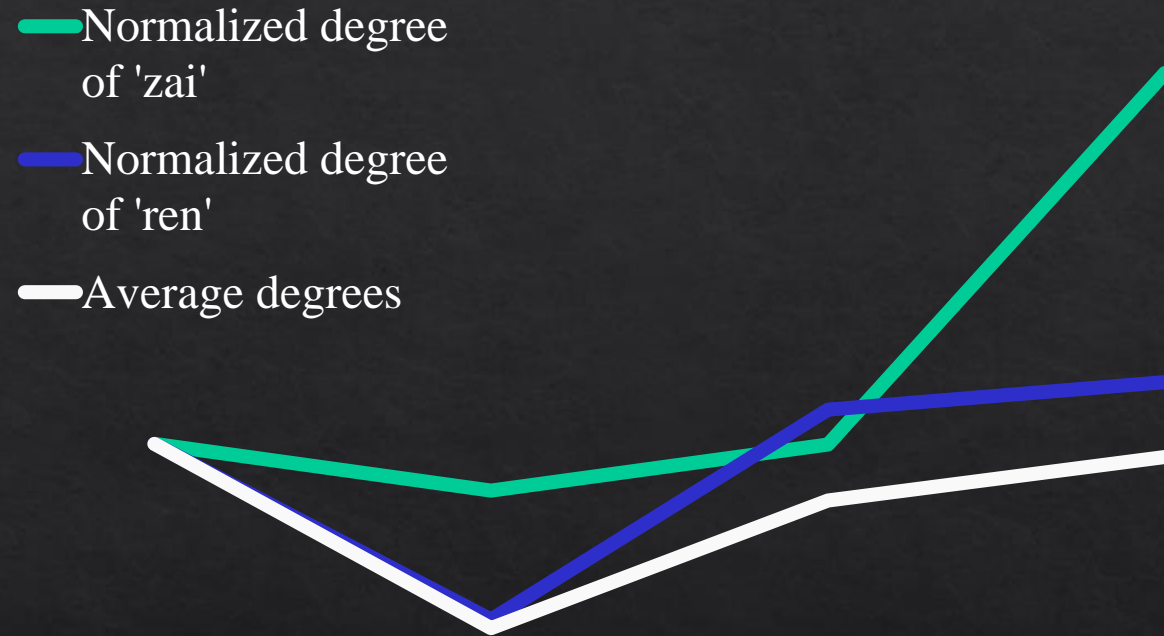
# Comparison of Single-character words

- ◆ Chinese is an isolating language: syntactic structure relies primarily on function words and word order rather than on rich morphological information to encode functional relations between elements (Levy & Manning 2003).
- ◆ 2 single-character words
  - ◆ Frequent characters/words both in the corpus and in general
  - ◆ zai underwent a grammaticalization process, whereas ren remains a content word
- ◆ Hubs could indicate the grammaticalization process and its starting points. Hubs could be functional or potential functional words to undergo future grammaticalization. (Solé et al. 2002)

# Grammaticalization Process in Networks



# Grammaticalization Process in Networks



# Conclusion

- ◈ The network features offer a new source of information to distinguish evolutions of different characters
- ◈ Observing evolution process in networks allows us to consider the evolutions of some specific units and the evolution of the whole system simultaneously.



# Questions & Comments

## Thank you !