

# The Electronic Data and Retrieval of the Secret History of the Mongols

Di Jiang



Institute of Ethnology & Anthropology  
Chinese Academy of Social Sciences



# Outline

---

- The background of SHM
- The original format of Chinese-transliterated document
- Make the electronic data for SHM
- Design a data retrieval system for SHM
- What may the electronic data tell us?

# What is SHM?

- SHM is a great classical historical work in 13<sup>th</sup> century. The name of the book is called *Monqol-un Nihuča Tobčiyān* in ancient Mongol (The Secret History of the Mongols, SHM)
- The book tells a historical story about Mongols and Genghis Khan, the conqueror of Central Asia and southern Russia, the founder and pioneer of Yuan Dynasty.





# SHM of the world

---

- The work SHM is so famous that it has become a branch of learning: SHK-ology all over the world.
  - There are more than ten language translations of SHM in the world
  - There are more than twenty Latin transliterations of SHM from Chinese characters
  - There are more than ten thousands of papers and books on SHM
  - There are more than one hundred experts of SHM in the world
  - .....





# Secrets of the Secret History

---

- Why does SHM arouse so much interests of scholars?  
Its secrets or riddles?
- People say: you may find the tomb of Genghis Khan, however, you may have no way to find the original SHM.
  - The original SHM has been lost. Can it be found again?
  - Who is the author of the book?
  - Why was the book handed down to generations in Chinese characters?
  - How can the book be restored to original Mongol from Chinese-transliterated characters?
  - What is concealed among the Chinese characters of SHM?
  - .....



# The format of Chinese-transliterated document

---

- the length of the book: 300 thousand Chinese characters (in guess before statistics).
- There are 12 volumes and 282 chapters by transliterators
- Format:
  - layouts: characters written vertically
  - the original shape of the archaic handwriting form with three lines representing one content
    - The first line is word-for-word Chinese characters (interlinearize)
    - The second line is Chinese-transliterated characters from Mongol
    - The third line is initials of the pronunciation of Mongol words
    - There are also endings indicating pronunciation of Mongol words within the second lines



## 元朝秘史卷五

太祖皇帝

那裏

種一行

擄着

種

成吉思哈罕

田迭

泰亦赤兀的

倒兀里周

泰亦赤兀台

骨頭有的

人名

勇士

人名

人名

牙速禿克兀泥

阿兀出把阿禿覘

豁團幹覘昌

忽都兀答覘

種每行

子孫的

子孫行

直到

灰飛一般

泰亦赤兀的

兀魯渾

兀魯哈

古覘帖列

忽捏速額覘

盡殺了

國

百姓行

他的教勸着

客亦思堅乞都罷

伯原作

兀魯思

亦覘格泥

阿訥

歌多格周

元秘史五

來着

太祖皇帝

地名

住冬了

亦列周

成吉思哈罕

惚巴恰牙

兀不者罷

伯原作

Chinese translation

成吉思將泰亦赤兀的阿兀出把阿禿覘等子孫殺盡將百姓起來至惚巴恰牙地面住冬了。

種

的

人名

老人

人名

人名

你出古

巴阿哩訥

失覘古額禿

額不堅

阿刺

納牙阿

兒子每

一同自的行

種的

官人

人名

可兀

魯額邊

泰亦赤兀敦

那顏

塔兒

忽台乞鄰禿

入林着

有的行

驚有的人

有來麼道

上馬

不

鬼刺周

不灰直

幹禿古溫

不列額客延

抹灑

兀祿

initial

volume

page

ending

ending

transliteration



# Can we make an electronic SHM?

---

- SHM is an ancient handwriting document
- SHM is a typical complicated document
- The carrier of information is unintegrated
- The characters and texts of SHM are non-standard and irregular





# Principles of making Electronic data for SHM

---

- Preserve the original format of SHM
  - Version, direction of handwriting, and annotation of pronunciation
- Keep all the information of characters
  - Simplified Characters, traditional Characters, and variant characters
- Keep all the information of volumes, chapters, and pages
- Keep all the information of segmentation and word-for-word format

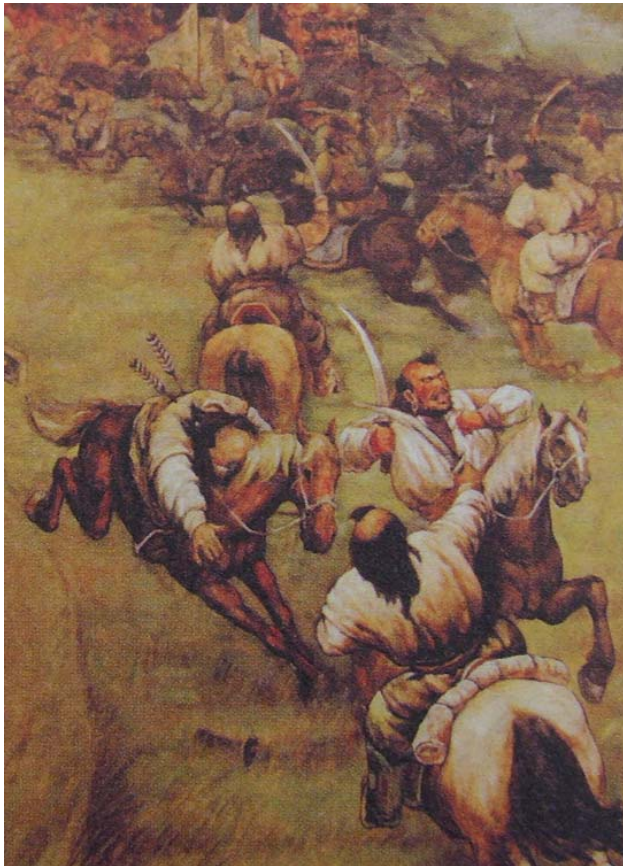


# The information of the handwriting document

- Layouts: three lines for one content
  - direction of handwriting: from vertical lines to lines sideways
  - Initials stand at one line
  - Word-for-word translation stands at one line

成吉思 太祖	中 合罕 皇帝	田迭 那裏	泰亦赤兀 種	的 行	倒兀里周 擄着	泰亦赤兀台 種	—initial —transliterate —word-for-word trans
-----------	---------------	----------	-----------	--------	------------	------------	--

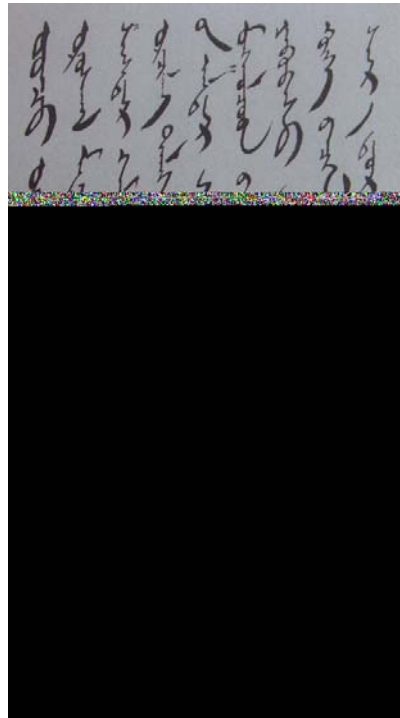
# The information of the handwriting document



- Volumes:
  - add a volume number before each volume
- Chapters:
  - keep the numbers of chapters
- Chinese translation
  - Be separated from the main texts

# The information of the handwriting document

- Characters:
  - keep all the forms of characters
  - “斷”/“断”
  - “驢”/“鐵”
  - “幾”/“几”
  - “仇”/“讎”/“讐”
  - .....



[11]	[10]
<p>都鑿鑿兒 阿合亦訥 桑兒遠可元亮 不列額</p> <p>人名 兄他的 子 看米</p> <p>元秘史一</p> <p>七</p> <p>帖堆阿塔刺</p>	<p>阿闐阿 桑奔兒干 途兒亦列周</p> <p>人名 人名 名均 有米</p> <p>行 來着</p> <p>西箇子 生丁</p> <p>不古訥台 別古訥台 捏列田不列額</p> <p>河闐阿。桑奔兒干取了為妻的後頭。生二子。一箇名不古訥台。一箇名別古訥台。</p> <p>帖堆阿塔刺</p>

page

# The format of Electronic version (deep structure)

chapter

transliteration

initial

ending

align

元朝秘史 卷一 忙豁侖紐察脫察安

1.1

#1

成吉思[名] 合罕訥[皇帝的] 忽札兀兒[根源]

迭額列{舌}[上]騰格理{舌}[天] 額扯[處] 札牙阿禿[命有的] 脫列{舌}<克>先[生子的] 孛兒帖[蒼色] 赤那[狼] 阿主兀[有] 格兒該[妻]亦訥[他的] 豁{中}埃[慄白色]馬闌{舌}<勒>[鹿] 阿只埃[有來] 騰汲思客[水名]禿<勒>周[渡着]亦列{舌}罷[來了] 斡難[河名] 沐漣{舌}訥帖里{舌}兀捏[河的源行] 不峒罕哈<勒>敦[山名]納[行]嫩兀<黑>刺周[營盤做着] 脫列{舌}<克>先[生子的] 巴塔赤罕[人名] 阿主兀[有來]

//當初元朝的人祖。是天生一箇蒼色的狼。與一箇慄白色的鹿相配了。同渡過騰吉思名字的水來。到於斡難名字的河源頭。不兒罕名字的山前住着。產了一箇人名字喚作巴塔赤罕。//

translation

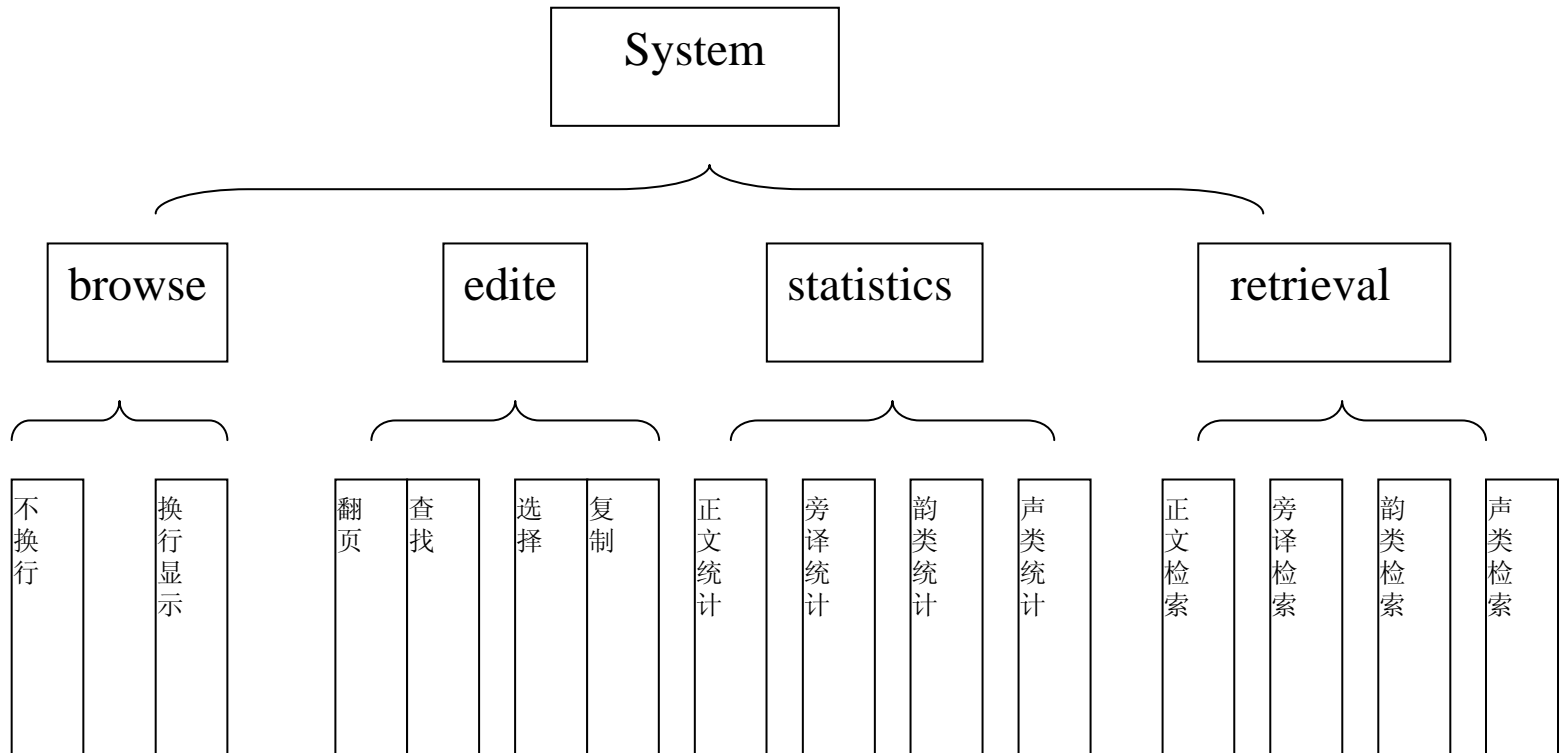


# Design for the retrieval system

---

- Functions of Browsing and seeking
  - Any characters or words
- Functions of statistics
  - Transliteration, Translation, Initials, endings
- Functions of retrieval
  - Transliteration, Translation, Initials, endings
- Functions of data output
- .....

# Design for the retrieval system





# Browse the original format

元朝秘史

文件 (F) 浏览 (V) 统计 (S) 检索 (I) 帮助 (H)

元朝秘史 卷十  
10.1  
#230

皇帝	說	雲有的	夜	天窓有的房	我的	圍臥着	
成吉思合罕	鳴話列論	額兀列台雪泥	額兀列台雪泥	幹魯格台格兒	米訥	額額連客 (卜) 帖周 幹	
中	舌			舌	舌	舌	
老的每	宿衛	我的	星有的	夜	宮 室我的	周圍臥着	被窠內
脫古思	客 (卜) 帖兀 (勤)	米訥	豁都台雪泥	幹兒朶格兒米訥	豁兒臣客帖周	幹樂朶脫刺	舌
				舌	舌	舌	舌
的	高位子裏	教到了		移動有的風雪行			
訥	溫突兒幹樂突兒	古兒格罷 (原作別)	失 (勤) 只鄰不恢孛羅幹納	失 (勤) 古 (惕)			
	舌	舌	舌	舌	舌		
歇息	不曾做立着	心	教安了的至誠心有的每	宿衛	我的		
只林兀祿勤擺亦周	只魯格	阿木兀魯 (黑) 三誠薛 (惕) 乞 (勤) 田	客 (卜) 帖兀 (勤) 米訥				
舌	舌						

# Retrieval by alignment

- To realize one content with three lines, the technology of interlinearization (alignment) has been used, which judges the browsing position of the three chunks according to what line is longer between 1<sup>st</sup> lines, 2<sup>nd</sup> lines, and 3<sup>rd</sup> lines.



# The result of aligning retrieval

longer

equal length

shorter

检索-正文

文件(F)

输入待检索字符: 必孫

检索

字符	卷	页	章	段	位置	显示
必孫 1	1.49	68	1	97	落後的每 魯<#>撒<#>	弟每自的行 迭兀捏里顏 舌
必孫 2	2.5	74	3	20	婦人名 行 訶額侖兀真 泥	寡婦 必孫 行 子每 小的每 泥 可兀<#> 兀出格<#>
必孫 12	12.24	272	1	186	格<#>	弟姪每自的行 迭兀捏里顏 舌

寡婦 必孫 嫂嫂行 別兒格泥顏 阿撒刺  
寡婦 必孫 行 子每 小的每  
媳婦兒自的行 別里顏 別魯迭

总共找到个数: 3

# Retrieval of translation by concordance

- A method of **concordance** may tell the contexts of what you seek.

检索-汉译

文件 (F)

输入待检索字符：

字符	卷	页	章	段	位置	显示
野獸	3	3.44	123	2	63	。并好馬都將來與你。野獸
野獸	3	3.44	123	2	77	圍呵。俺首先出去圍將來與你。如廝殺時違了野獸
野獸	5	5.38	164	1	179	處剿捕時。一同剿捕。野獸
野獸	6	6.19	175	2	97	成吉思止當不從。因趕野獸
野獸	7	7.4	187	2	102	搶得財物。打獵時得的野獸
野獸	9	9.27	219	1	53	得的財物。圍獵時得的野獸
野獸	12	12.51	279	1	140	川地面先因無水。止有野獸
野獸	12	12.58	281	1	123	陰害了。一件將天生的野獸

总共找到个数：8



# What may the electronic data tell us? (for historians)

---

- Who appears in SHM?
- Where did the named persons go?
- What events happened to them?
- How can you find the secrets of original texts from Chinese characters?
  - Seek names of persons
  - Seek names of places
  - Seek some event keywords

输入待检索字符：

斡難河

检索

字符	卷	页	章	段	位置	显示
斡難河	1	1. 16	24	1	63	梁瘡秃尾子的馬。順着
斡難河	1	1. 19	30	1	18	不忽合塔吉。後來順着
斡難河	1	1. 31	50	4	61	一子。名不里孛闊。於
斡難河	1	1. 34	54	3	15	父。也速該把阿秃兒在
斡難河	1	1. 36	55	2	160	了。即便打着馬。逆着
斡難河	1	1. 38	56	2	116	怎生般艱難。哭的聲將
斡難河	1	1. 41	59	1	49	兒的妻訶額侖正懷孕於
斡難河	2	2. 7	75	2	38	親上頭。將針做鉤兒於
斡難河	2	2. 18	81	1	52	十六日。泰亦赤兀每於
斡難河	2	2. 18	81	1	114	於頭上打倒走了。走到
斡難河	2	2. 18	81	1	129	內臥着恐怕人見。又入
斡難河	2	2. 19	82	2	45	來着。白日般月明裏。
斡難河	2	2. 22	84	2	116	。必救了我。所以順着
斡難河	3	3. 9	106	2	107	軍。共二萬軍上馬逆着
斡難河	3	3. 13	108	1	24	字那裏相合了。起去往
斡難河	3	3. 27	116	1	66	時。帖木真十一歲。於
斡難河	4	4. 5	129	2	134	被札木合推動。退著於
斡難河	4	4. 7	130	2	81	百姓來了。喜歡。着於
斡難河	4	4. 38	144	2	89	落。將百姓起了。渡過
斡難河	8	8. 27	202	2	21	捕了。至是虎兒年。於
斡難河	9	9. 6	211	2	31	風匣自不峒罕山來。於

总共找到个数：21

# A sample (for grammarians)

检索-汉译

文件(E)

输入待检索字符：

被

检索

字符	卷	页	章	段	位置	显示
被	1	1.6	9	2	30	敦地面貂鼠青鼠野物。
被	1	1.34	53	1	5	親自送去。
被	1	1.34	53	1	97	人。為親送女兒上頭。
被	1	1.38	56	2	88	說。我的丈夫頭髮不曾
被	1	1.39	57	2	6	因俺巴孩合罕
被	1	1.48	67	2	58	乞顏來了。因記起舊日
被	1	1.49	68	2	52	帖木真去做女婿回時。
被	2	2.3	71	1	38	俺巴孩皇帝死了麼道。
被	2	2.5	73	2	5	察刺合老人
被	2	2.5	73	2	36	收的并俺衆人的百姓。
被	2	2.5	73	2	47	他將去。因勤他的時分
被	2	2.8	76	2	80	。我釣得一箇金色魚。
被	2	2.10	77	2	31	我昨前射得箇雀兒。也
被	2	2.10	77	2	45	。今遍釣得箇魚兒。又
被	2	2.10	77	2	115	。將箭抽着要射他時。
被	2	2.14	79	2	174	。上馬走入山林裏去。
被	2	2.16	80	2	189	割開。牽着馬下山來。
被	2	2.24	85	1	41	箇名字的兒子說。雀兒
被	2	2.30	90	2	14	真的慘白驢馬八正在家
被	2	2.31	90	1	25	來。帖木真說。我的馬
被	2	2.51	103	2	99	躲得過。我的小性命。
被	3	3.3	104	2	43	王罕處去到了說。不想
被	3	3.6	105	3	39	。教對他說。我的妻子
被	3	3.7	105	1	34	也聽得帖木真安蒼的妻
被	3	3.19	111	2	56	弟赤列都。他妻訶額命
被	3	3.29	117	1	135	下。做了筵席。夜晚共
被	4	4.5	129	2	124	巴主地面對陣。成吉思
被	4	4.8	130	1	76	坤太子。死了的上頭。
被	4	4.10	131	1	3	韉繩。





# A sample (for corpus linguists)

- Statistics for Chinese characters in SHM

Char. types	Amount	Token	Rate
1, Characters of translation	1669	47907	21.3354
Initials	2	1302	
Endings	9	608	
Character	1665	40674	
Full stop	1	5323	
2, Characters of alignment	1567	61893	27.5641
3, Annotative characters	10	1582	0.7045
4, Characters of transliteration	546	113160	50.0396
Single words	540	90457	
Word initials	2	15533	
Word endings	15	7170	
Total Chinese characters	2099	224542	100.00





# Thank you for your attention!

---

- The electronic data of SHM will be useful for
  - Historians
  - Philologists
  - Linguists
  - Geographers
  - Military scientists
  - Students



# Types of characters in SHM

---

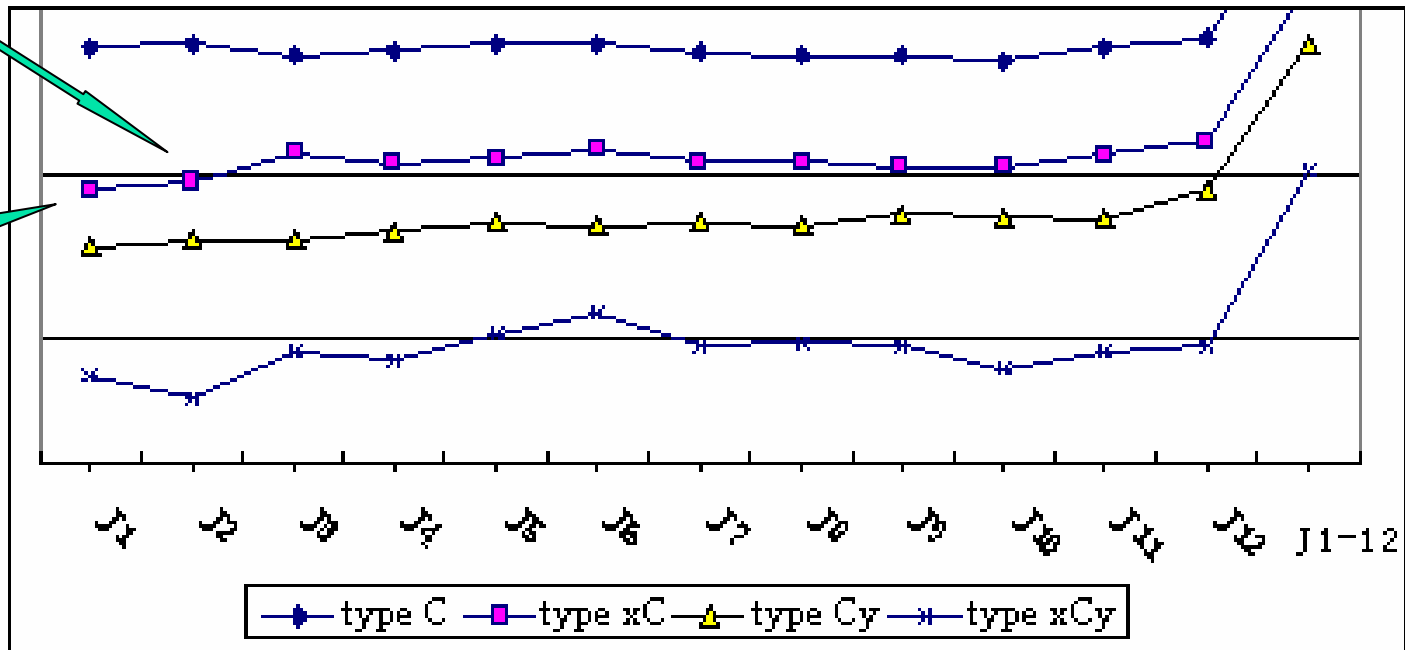
- Type one (C): single character.
  - “兕”
- Type two (xC): character with initials.
  - “舌兕”
- Type three (Cy): character with endings.
  - “阿<sub>勒</sub>”
- Type four (xCy): character with initials and endings.
  - “舌魯<sub>黑</sub>”

# Statistics for types of characters in SHM

- Pay attention to tokens of Type xC in V.1 & v.2

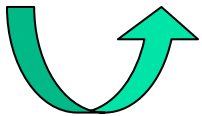
	C	token	rate	xC	token	rate	Cy	token	rate	xCy	token	rate
V.1	250	5949	6.5766	37	<u>820</u>	<u>0.9065</u>	111	367	0.4057	19	<u>62</u>	0.0685
V.2	236	6241	6.8994	34	<u>907</u>	<u>1.0026</u>	113	396	0.4377	13	<u>45</u>	0.0497
V.3	253	5309	5.8690	37	1354	1.4968	115	406	0.4488	21	87	0.0961
V.4	236	5670	6.2681	32	1175	1.2989	114	465	0.5140	25	76	0.0840
V.5	224	6088	6.7302	30	1252	1.3840	114	505	0.5582	24	109	0.1205
V.6	251	6220	6.8761	42	<u>1417</u>	<u>1.5665</u>	115	495	0.5472	27	<u>144</u>	<u>0.1592</u>
V.7	231	5460	6.0360	28	1192	1.3177	134	517	0.5715	22	94	0.1039
V.8	230	5217	5.7674	33	1216	1.3443	128	491	0.5428	24	97	0.1072
V.9	208	5183	5.7298	27	1098	1.2138	110	590	0.6522	22	93	0.1028
V.10	219	4765	5.2677	32	1096	1.2116	115	551	0.6091	23	68	0.0752
V.11	240	5931	6.5567	32	1337	1.4781	117	538	0.5947	23	86	0.0951
V.12	235	<u>6777</u>	<u>7.4919</u>	33	<u>1613</u>	<u>1.7832</u>	119	<u>793</u>	<u>0.8766</u>	25	95	0.1050
Total	510	68810	76.069	83	14477	16.004	317	6114	6.7590	52	1056	1.1674

# The reason: An interesting case



# See how the character types influence the whole text

- “兕”(and “舌兕”) is a character in the most high frequency in the whole document.
- There are two forms of “兕”, one is single “兕”, another is “舌兕” with initials.
  - “兕” in vol. 1 is 362, “舌兕” is 0
  - “兕” in vol. 2 is 478, “舌兕” is 8
  - “兕” in vol. 3 is 17, “舌兕” is 391
  - “兕” in vol. 4 is 57, “舌兕” is 298
  - “兕” in vol. 5 is 16, “舌兕” is 376
  - .....
  - “兕” in vol. 1~12 is 1116, “舌兕” is 3707



- Now we know the transliterators made some errors in vol.1 and 2. All forms of the character should be “舌兕”.