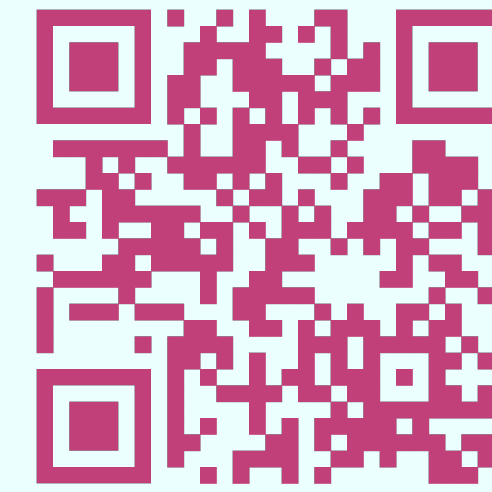


PEKING  
UNIVERSITY



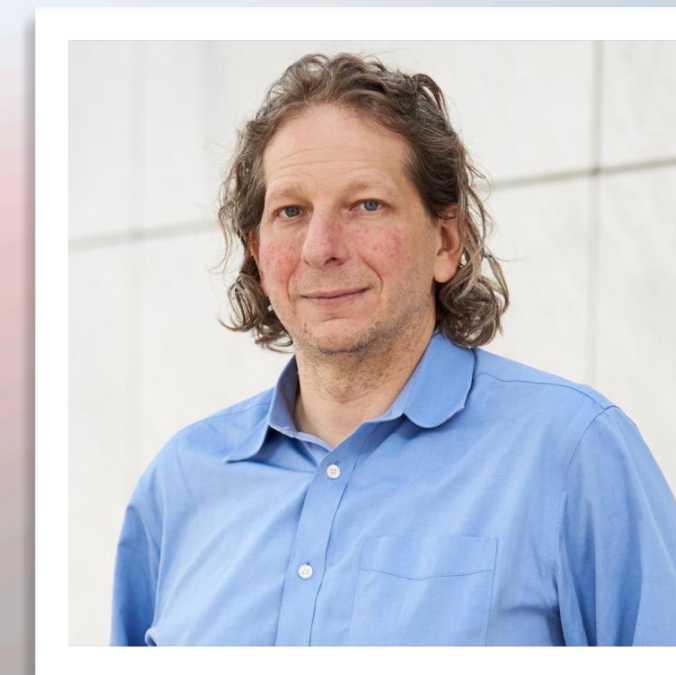
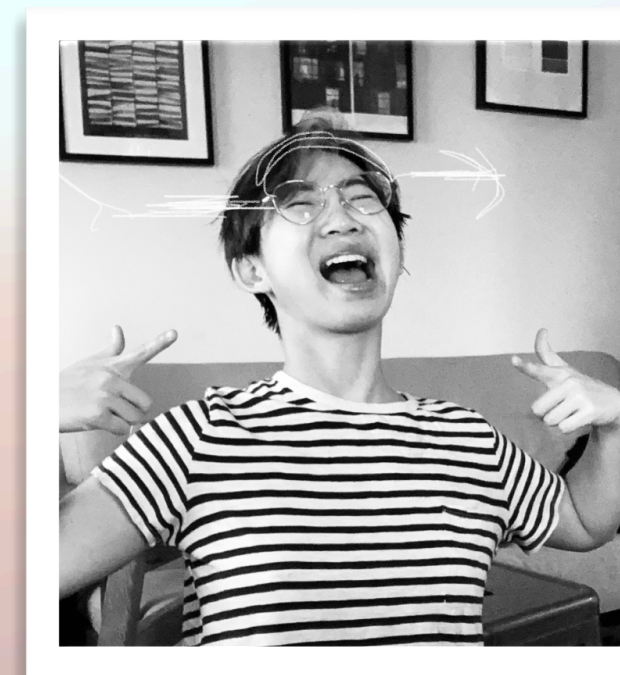
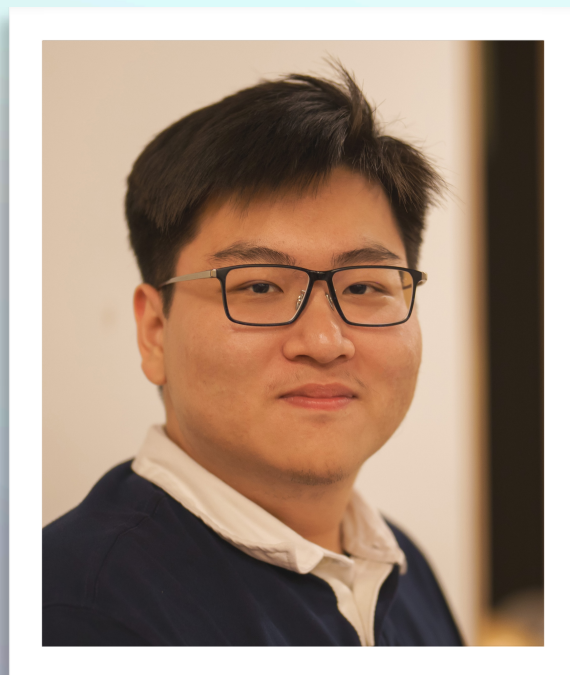
Full Paper



Sayan Gul Award  
Best Undergrad Paper

# Finding structure in logographic writing with library learning

Guangyuan Jiang, Matthias Hofer, Jiayuan Mao, Lionel Wong, Josh Tenenbaum, and Roger Levy



# Human language is deeply structured

- Human language is deeply structured — a **universal** trait of human communication systems.

# Human language is deeply structured

- Human language is deeply structured — a **universal** trait of human communication systems.

$$\left[ \begin{array}{l} + \text{stop} \\ + \text{consonantal} \\ + \text{alveolar} \end{array} \right] \rightarrow [r] \ / \ \left[ \begin{array}{l} + \text{vowel} \\ + \text{stressed} \end{array} \right] \text{ \_\_\_\_ } \left[ \begin{array}{l} + \text{vowel} \\ - \text{stressed} \end{array} \right]$$

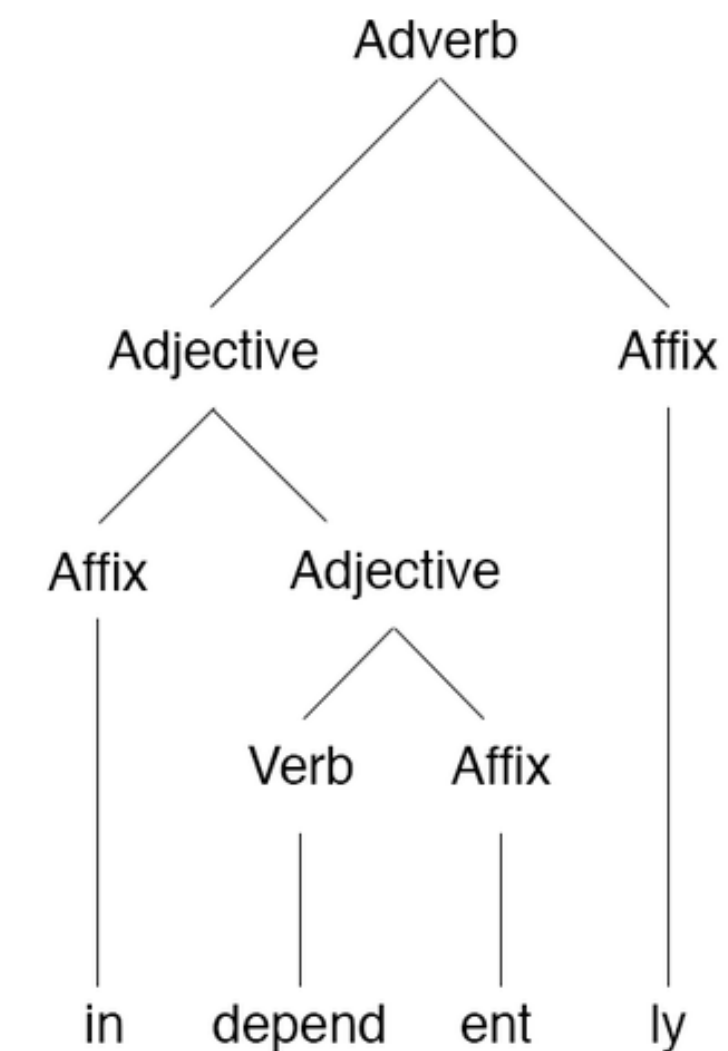
phonetics and phonology

# Human language is deeply structured

- Human language is deeply structured — a **universal** trait of human communication systems.

$\left[ \begin{array}{l} + \text{ stop} \\ + \text{ consonantal} \\ + \text{ alveolar} \end{array} \right] \rightarrow [r] / \left[ \begin{array}{l} + \text{ vowel} \\ + \text{ stressed} \end{array} \right] \text{ \_\_\_\_ } \left[ \begin{array}{l} + \text{ vowel} \\ - \text{ stressed} \end{array} \right]$

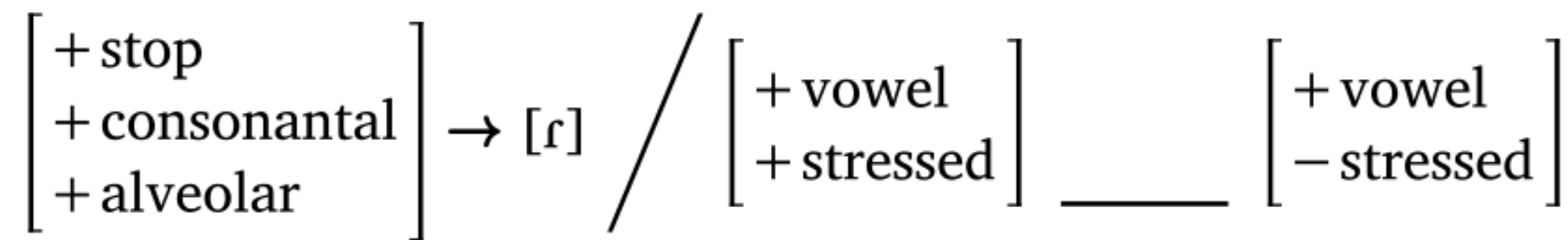
phonetics and phonology



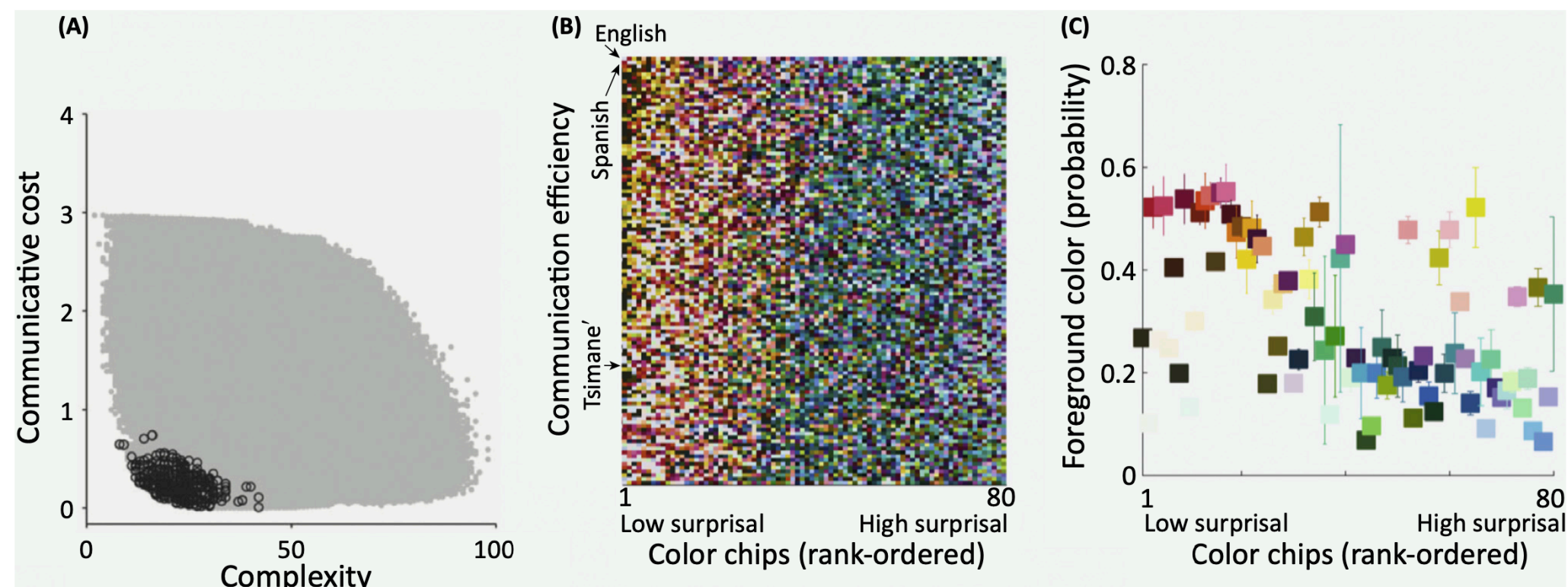
morphology and syntax

# Human language is deeply structured

- Human language is deeply structured — a **universal** trait of human communication systems.

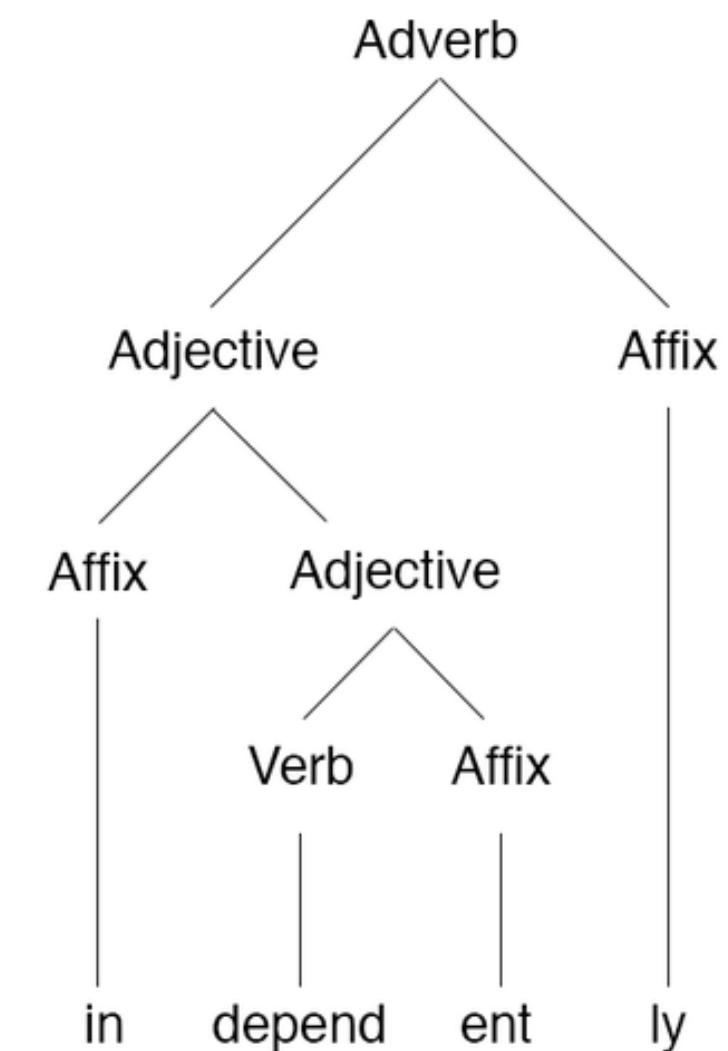


phonetics and phonology



semantic space: kinship, color

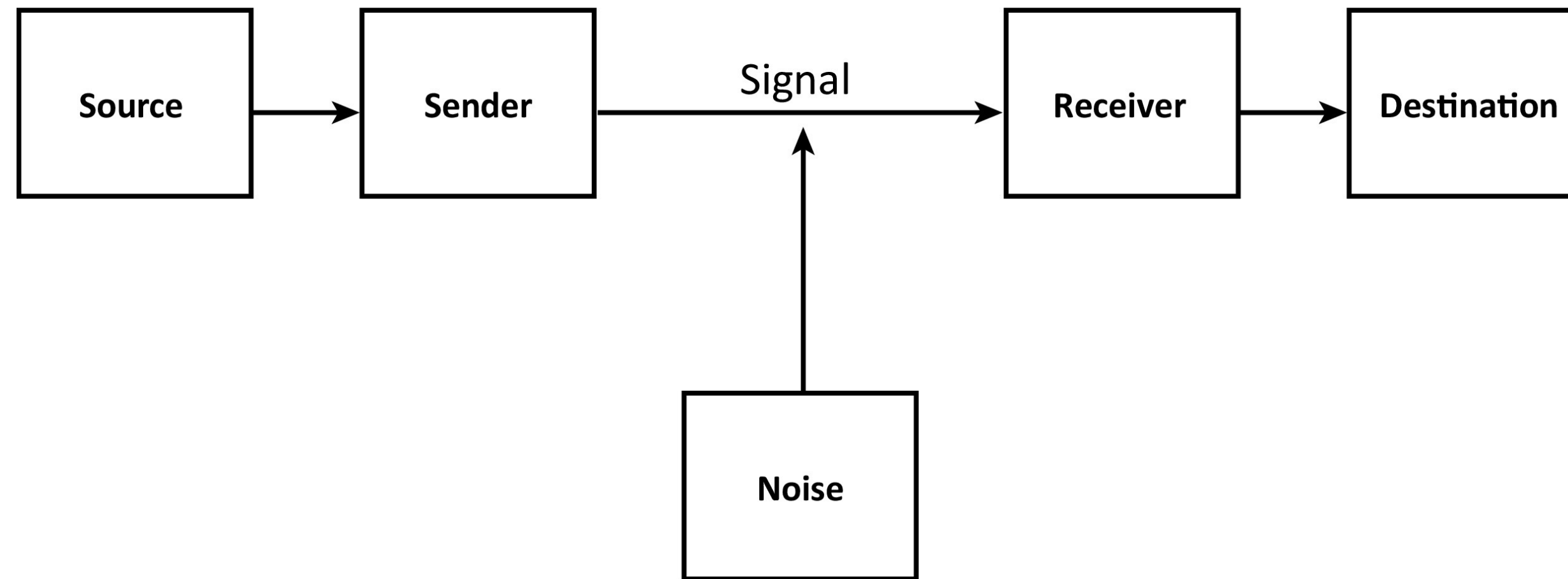
(Kemp & Regier, 2012; Zaslavsky et al., 2018)



morphology and syntax

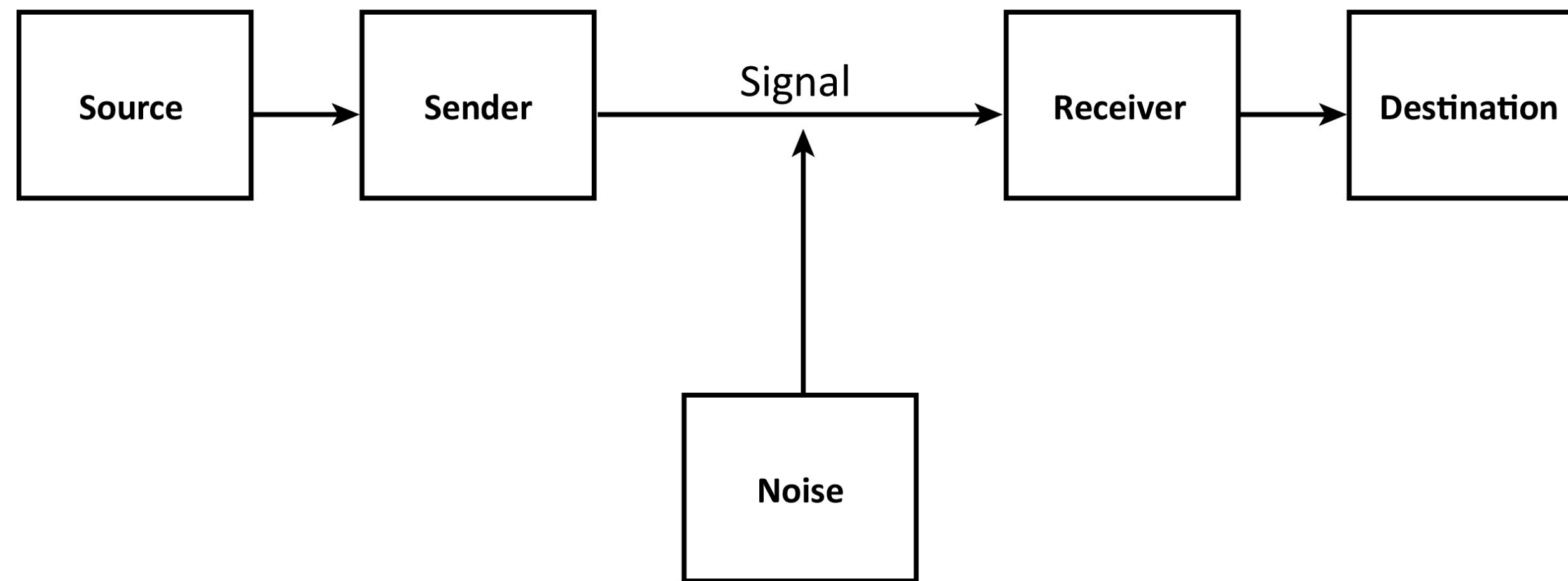
# Why structured? Shaped by efficiency

# Why structured? Shaped by efficiency

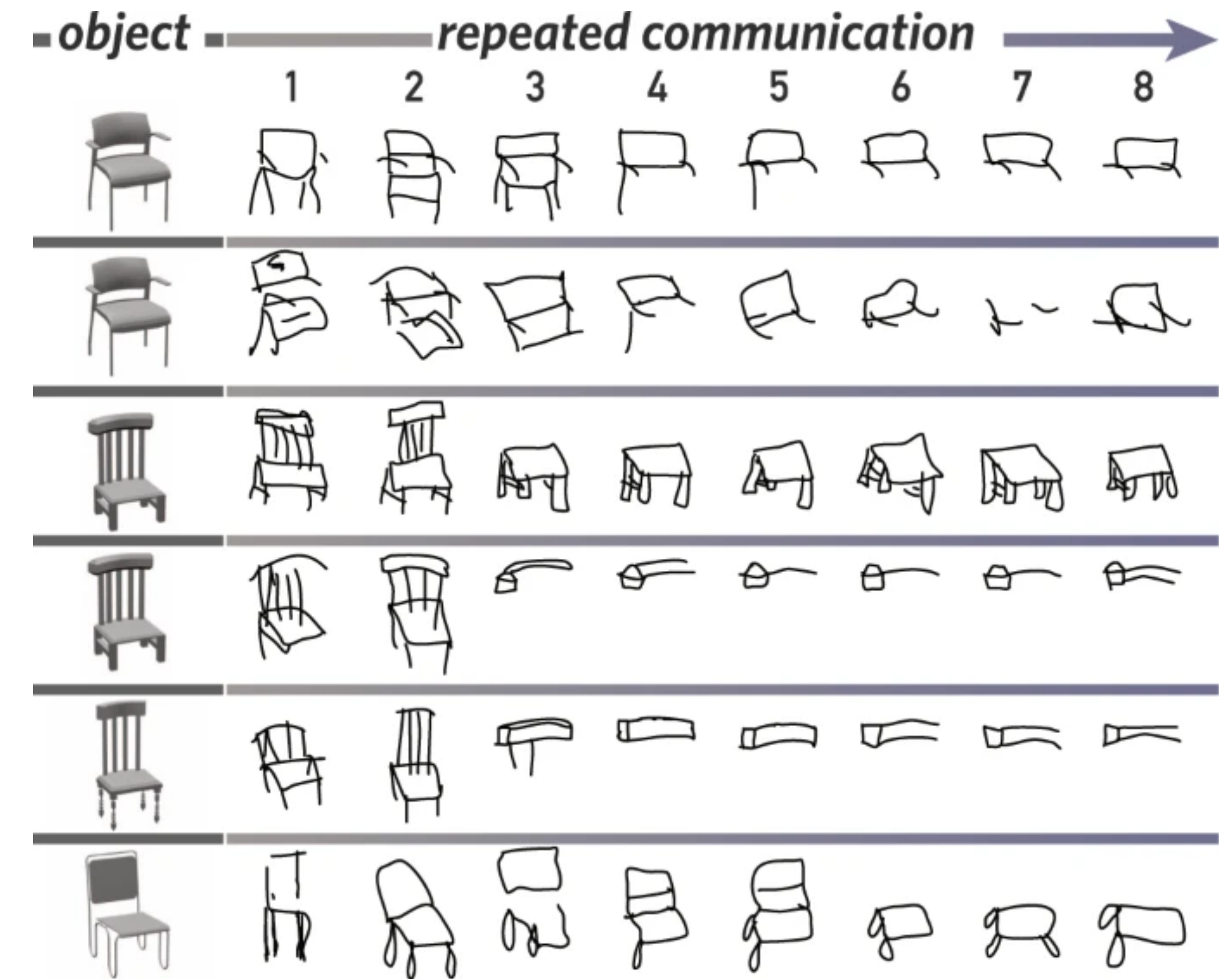


communicative efficiency

# Why structured? Shaped by efficiency



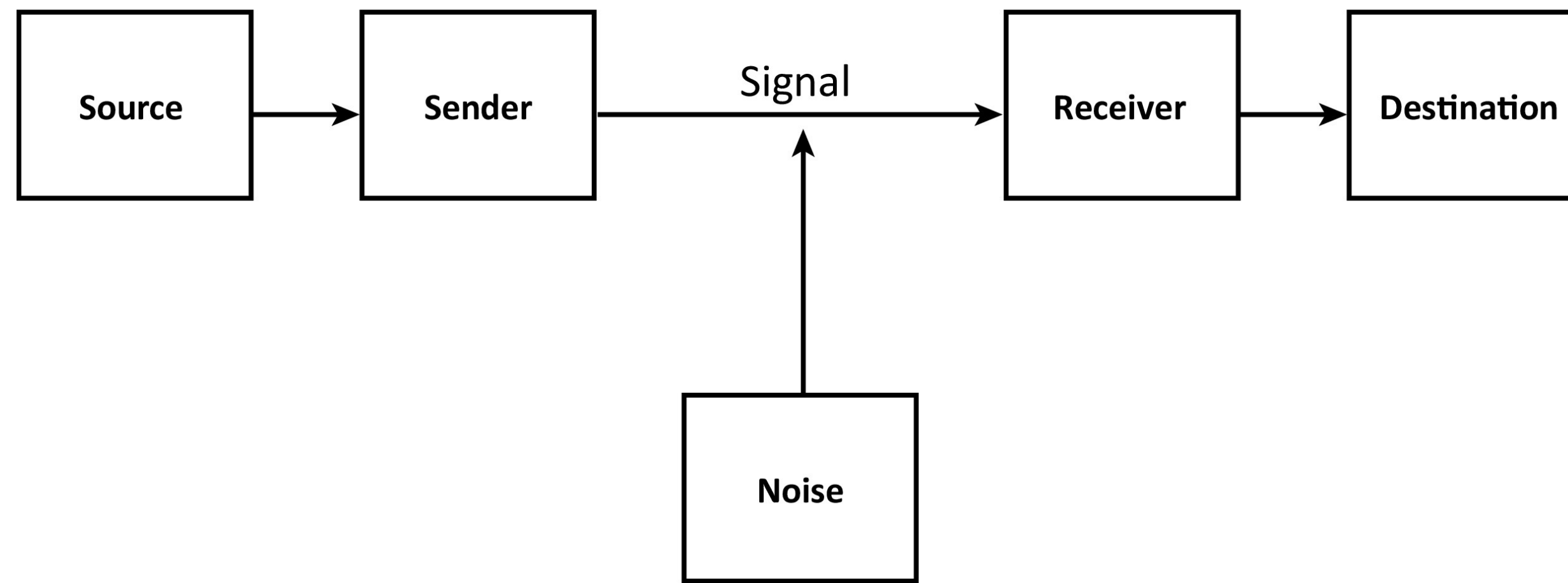
communicative efficiency



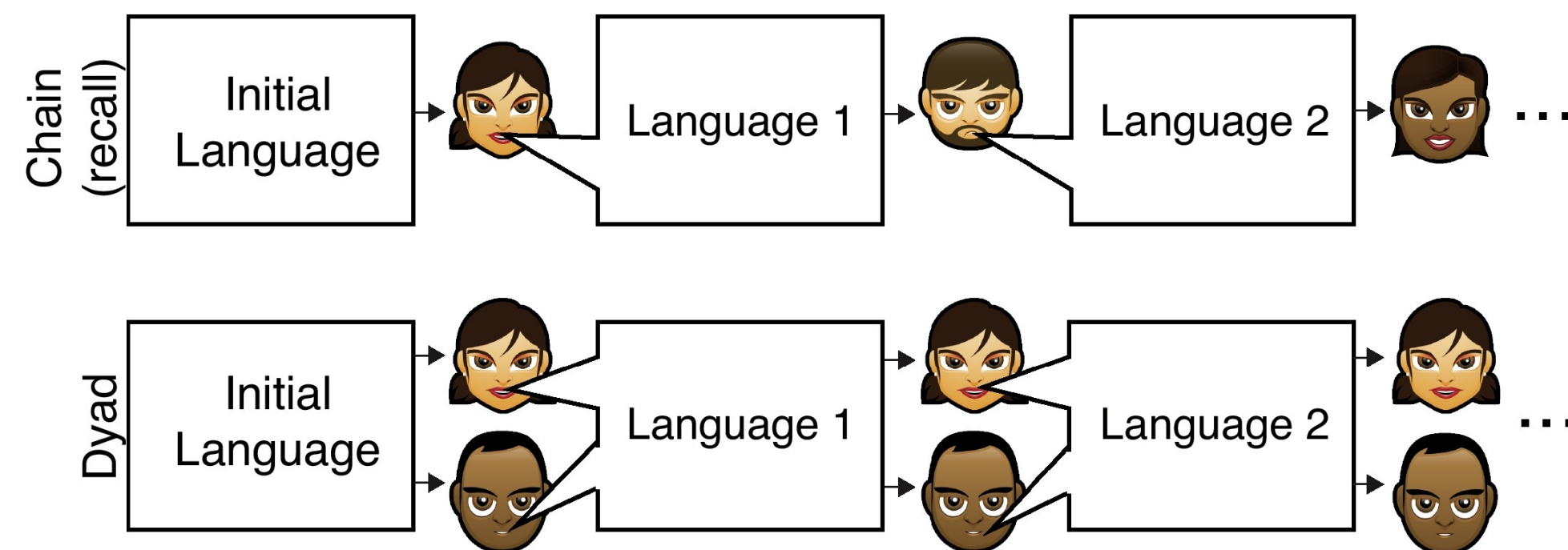
representational efficiency



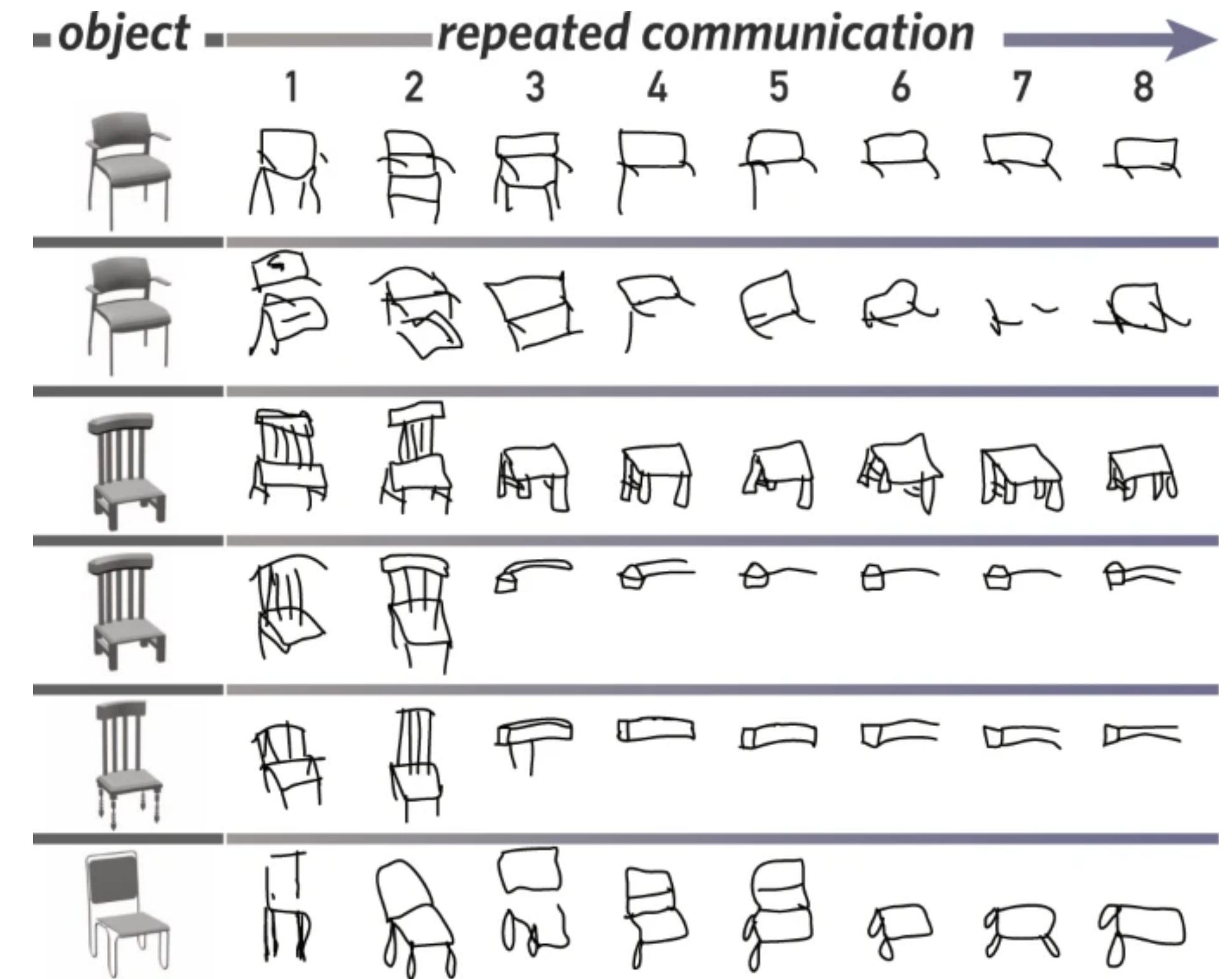
# Why structured? Shaped by efficiency



communicative efficiency



learnability



representational efficiency

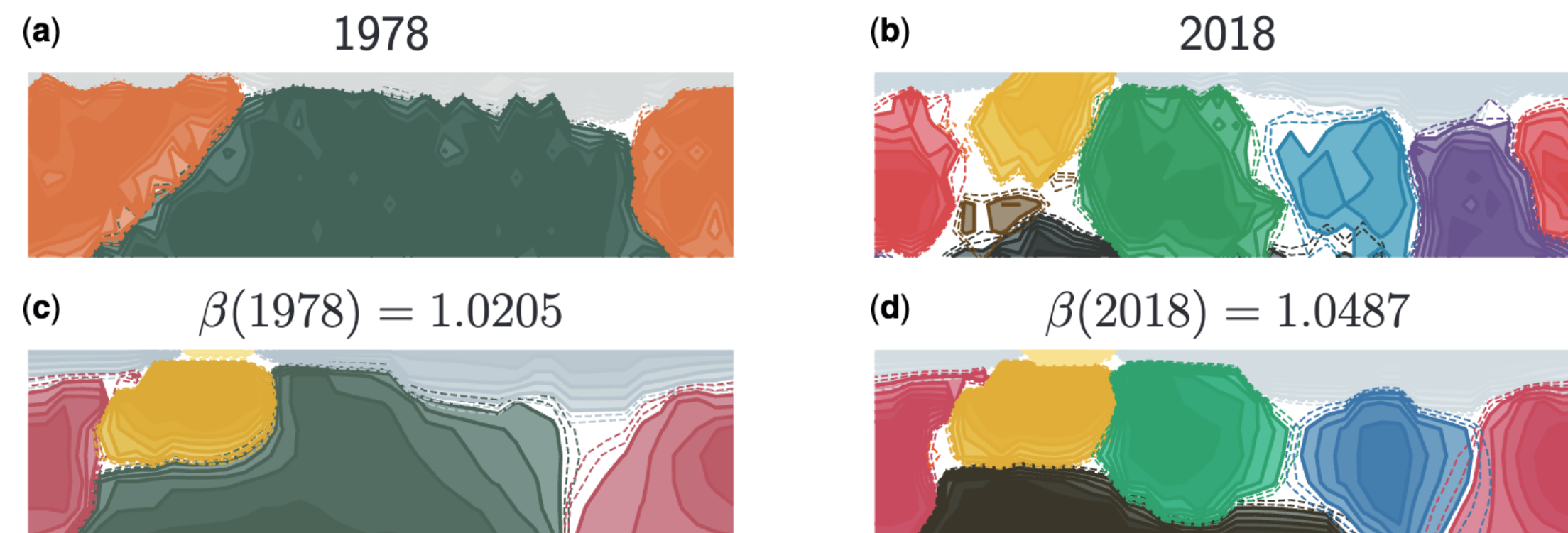
# **Human language shaped by efficiency**

**...and how languages evolve over time**

# Human language shaped by efficiency

## ...and how languages evolve over time

In the wild



(Zaslavsky et al., 2022)

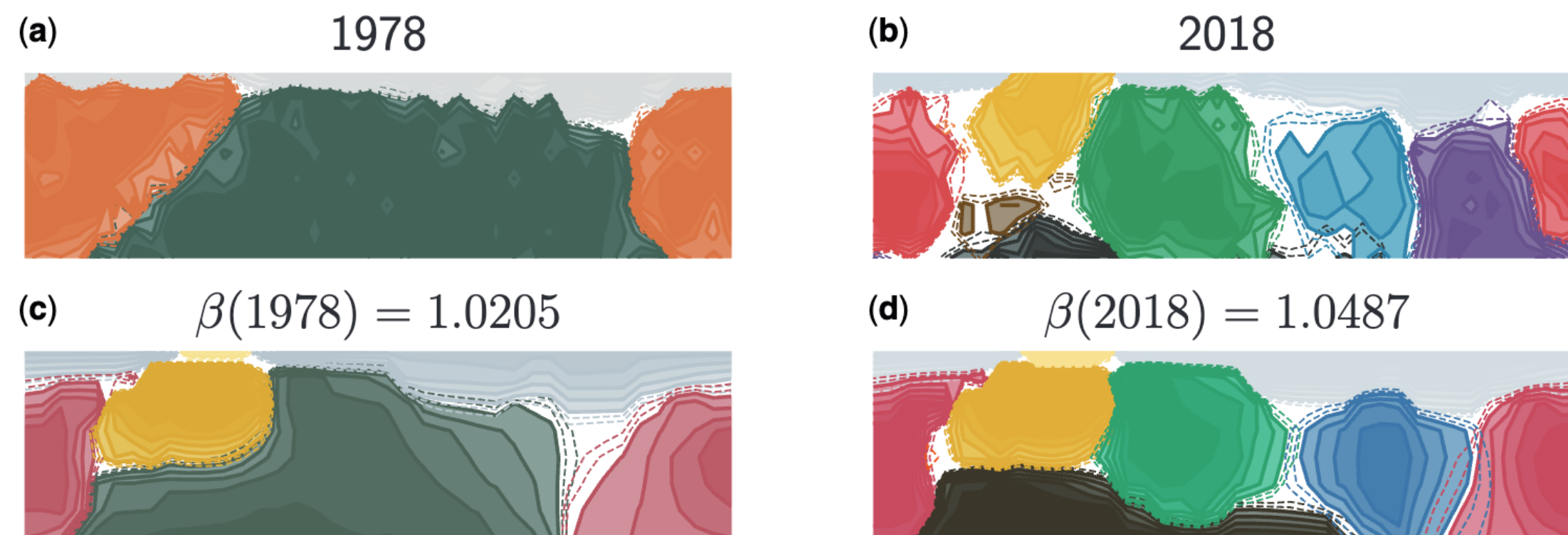


Chinese characters

# Human language shaped by efficiency

## ...and how languages evolve over time

In the wild



(Zaslavsky et al., 2022)



Chinese characters

In the lab

→	n-ere-ki	l-ere-ki	renana	□
	n-ehe-ki	l-aho-ki	r-ene-ki	○
	n-eke-ki	l-ake-ki	r-ahe-ki	△
↗↘	n-ere-plo	l-ane-plo	r-e-plo	□
	n-eho-plo	l-aho-plo	r-eho-plo	○
	n-eki-plo	l-aki-plo	r-aho-plo	△
↻	n-e-pilu	l-ane-pilu	r-e-pilu	□
	n-eho-pilu	l-aho-pilu	r-eho-pilu	○
	n-eki-pilu	l-aki-pilu	r-aho-pilu	△










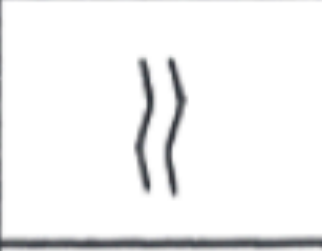










(Kirby, Cornish, & Smith, 2008)

# **Human language shaped by efficiency**

**Languages' earliest records – writing**

# Human language shaped by efficiency










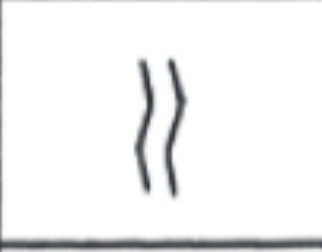










## Languages' earliest records – writing

	3100 BC	3000 BC	2400 BC	1000 BC
head				
mouth/speak				
water				
drink				
go/stand/bring				

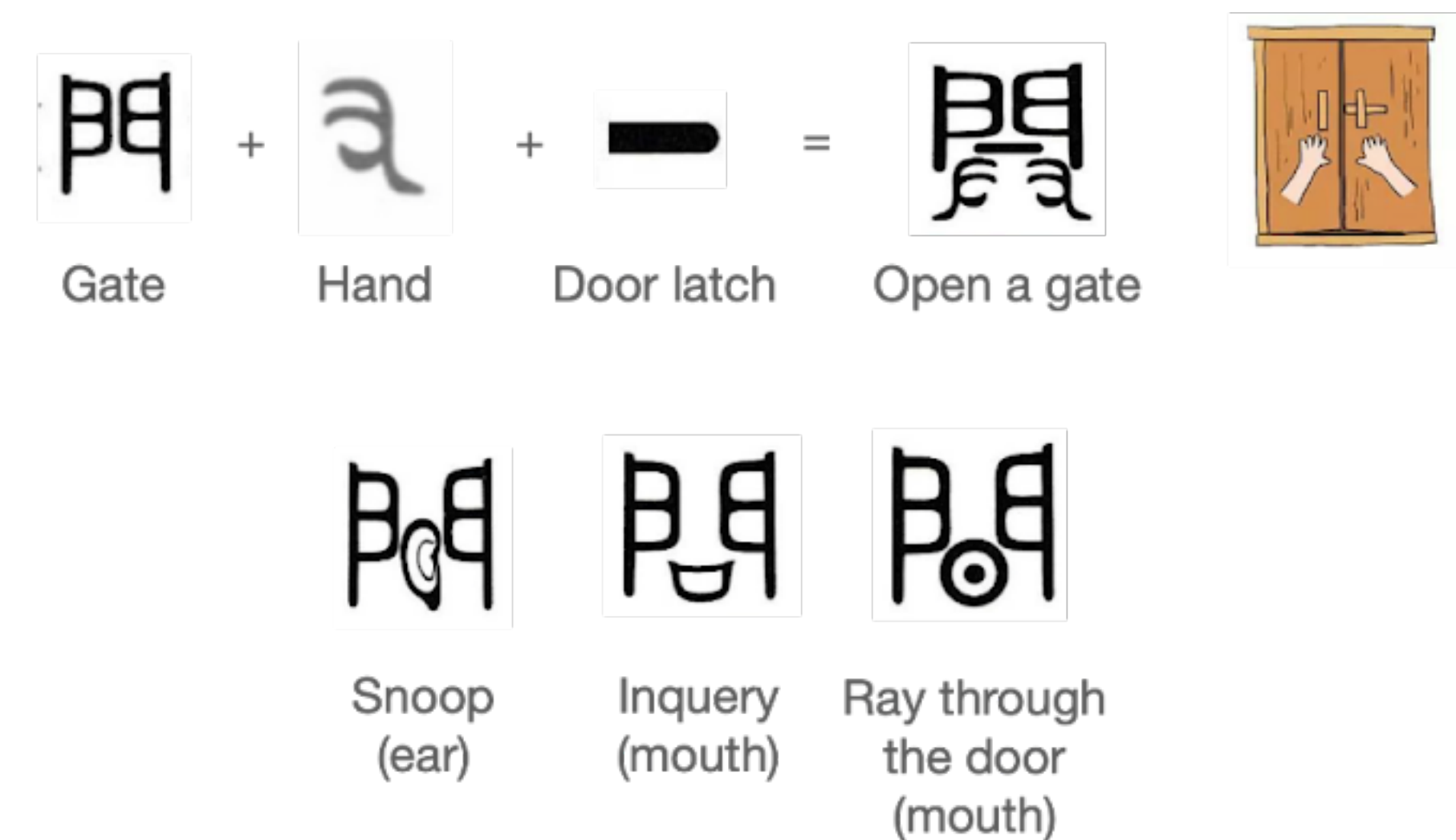
Sumerian Cuneiform (Sampson, 1985)

# Human language shaped by efficiency

## Languages' earliest records – writing

	3100 BC	3000 BC	2400 BC	1000 BC
head				
mouth/speak				
water				
drink				
go/stand/bring				

Sumerian Cuneiform (Sampson, 1985)



Chinese oracle bone scripts

# Efficiency-based structure discovery

## Morphology discovery (Goldsmith, 2001)

*laughed*

*laughing*

*laughs*

*walked*

*walking*

*walks*

*jumped*

*jumping*

*jumps*

total letter count: 57



# Efficiency-based structure discovery

## Morphology discovery (Goldsmith, 2001)

*laugh***ed**

*laugh***ing**

*laugh***s**

*walk***ed**

*walk***ing**

*walk***s**

*jump***ed**

*jump***ing**

*jump***s**

total letter count: **57**

# Efficiency-based structure discovery

## Morphology discovery (Goldsmith, 2001)

*laugh***ed**

*laugh***ing**

*laugh***s**

*walk***ed**

*walk***ing**

*walk***s**

*jump***ed**

*jump***ing**

*jump***s**



total letter count: **57**

# Efficiency-based structure discovery

## Morphology discovery (Goldsmith, 2001)

*laugh***ed**  
*laugh***ing**  
*laugh***s**  
*walk***ed**  
*walk***ing**  
*walk***s**  
*jump***ed**  
*jump***ing**  
*jump***s**

total letter count: **57**



*laugh*  
*walk*  
*jump*

**ed**  
**ing**  
**s**

total letter count: **19**

# Efficiency-based structure discovery

## Morphology discovery (Goldsmith, 2001)

*laugh***ed**  
*laugh***ing**  
*laugh***s**  
*walk***ed**  
*walk***ing**  
*walk***s**  
*jump***ed**  
*jump***ing**  
*jump***s**

total letter count: **57**



*laugh* } { **ed**  
*walk* } { **ing**  
*jump* } { **s**


*Efficient representation leads to  
morphological structure*

total letter count: **19**

# Efficiency-based structure discovery

Learning syntax/morphological rules (Kim, Dyer, & Rush, 2019; Ellis et al., 2022)

PCFG Rule	DMV parameter		<u>masculine</u>	<u>feminine</u>
$S \rightarrow Y_h$	$P_{root}(h)$			
$Y_h \rightarrow L_h^0 R_h^0$	1	<i>rich</i>	bogat	bogata
$L_h^0 \rightarrow h_l$	$P_{stop}(\text{Stop} h, \leftarrow, \text{no\_dep})$	<i>mild</i>	blag	blaga
$L_h^0 \rightarrow L'_h$	$P_{stop}(\neg\text{Stop} h, \leftarrow, \text{no\_dep})$	<i>green</i>	zelen	zelena
$L'_h \rightarrow Y_d L_h$	$P_{choose}(d h, \leftarrow)$			
$L_h \rightarrow h_l$	$P_{stop}(\text{Stop} h, \leftarrow, \text{one\_dep})$			
$L_h \rightarrow L'_h$	$P_{stop}(\neg\text{Stop} h, \leftarrow, \text{one\_dep})$			



add /a/ to feminine

learning syntax/grammar  
(Kim, Dyer, & Rush, 2019)

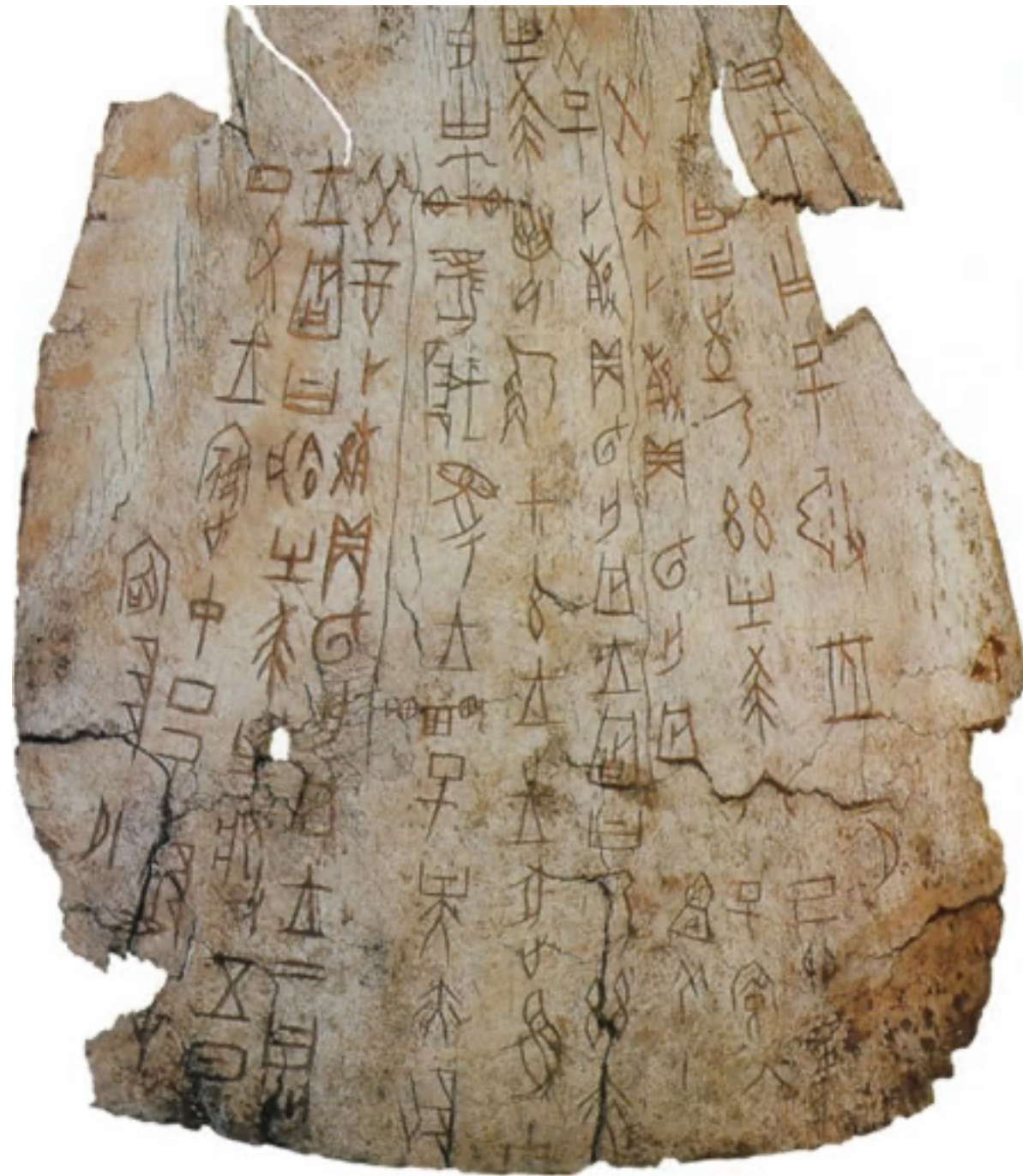
morpho-phonology discovery  
(Ellis et al., 2022)

# Our hypotheses

- If language has been shaped by pressure for efficiency,
  - then search for representational efficiency should recover its combinatorial elements.
- Furthermore, we should see an increase in efficiency over time.

# Our domain: Chinese Orthography

Earliest records — oracle bone scripts (~1500-1050 BC)



人 rén person	男 nán man	女 nǚ woman	子 zi child	夫 fu husband	妻 qī wife	王 wáng king	口 kǒu mouth
目 mù eye	耳 ěr ear	心 xīn heart	日 rì sun	月 yuè moon	山 shān mountain	雨 yǔ rain	田 tián field
土 tǔ earth	水 shuǐ water	火 huǒ fire	貝 bèi cowrie shell	大 dà big	小 xiǎo small	上 shàng above	下 xià below
力 lì strength	中 zhōng middle	先 xiān first	光 guāng bright	肉 ròu meat	出 chū to go out	刀 dāo knife	南 nán south

# Our domain: Chinese Orthography

oracle

seal

traditional

simplified

*sink*

*float*

*color*

*insect*

*orange*

*peace*

1500 BC

1050 BC

200 AD

1950 AD



# Our domain: Chinese Orthography

oracle

seal

traditional

simplified

*sink*

*float*

*color*

*insect*

*orange*

*peace*

1500 BC

1050 BC

200 AD

1950 AD

# Chinese 101

Long evolutionary history (of over 3,000 years).

	oracle	seal	traditional	simplified
<i>sink</i>				
<i>float</i>				
<i>color</i>				
<i>insect</i>				
<i>orange</i>				
<i>peace</i>				

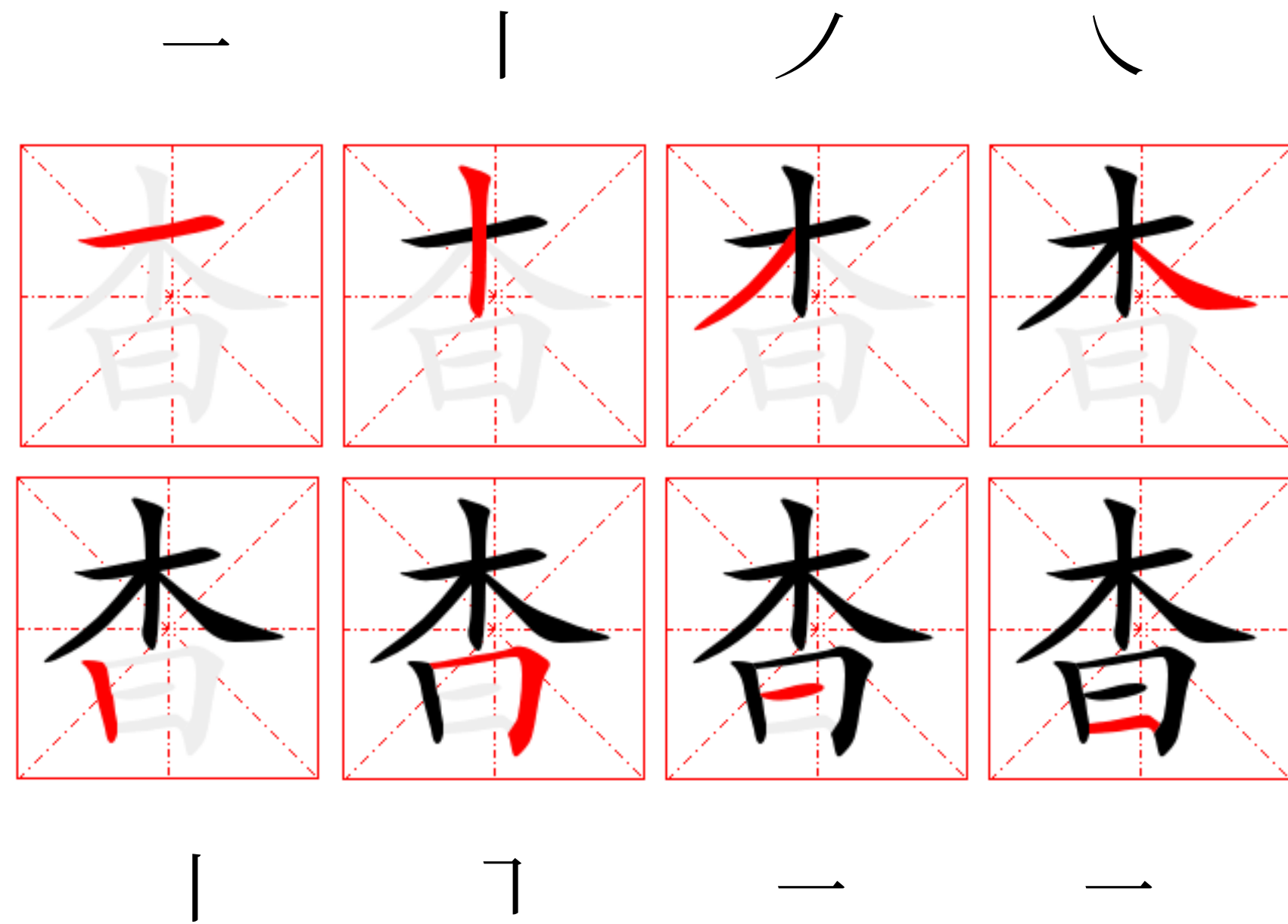
1500 BC
1050 BC
200 AD
1950 AD

- **Frequent reuse** of graphical elements **within** individual **characters** and **across** the **writing system**.
- **Unique opportunity** for studying combinatorial structure.



# Chinese 101

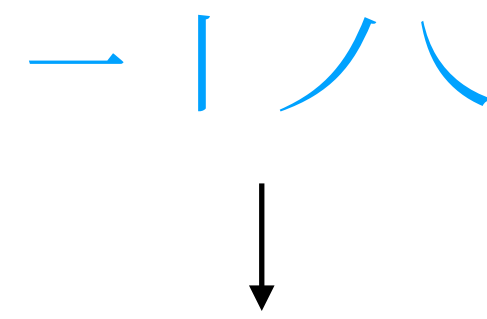
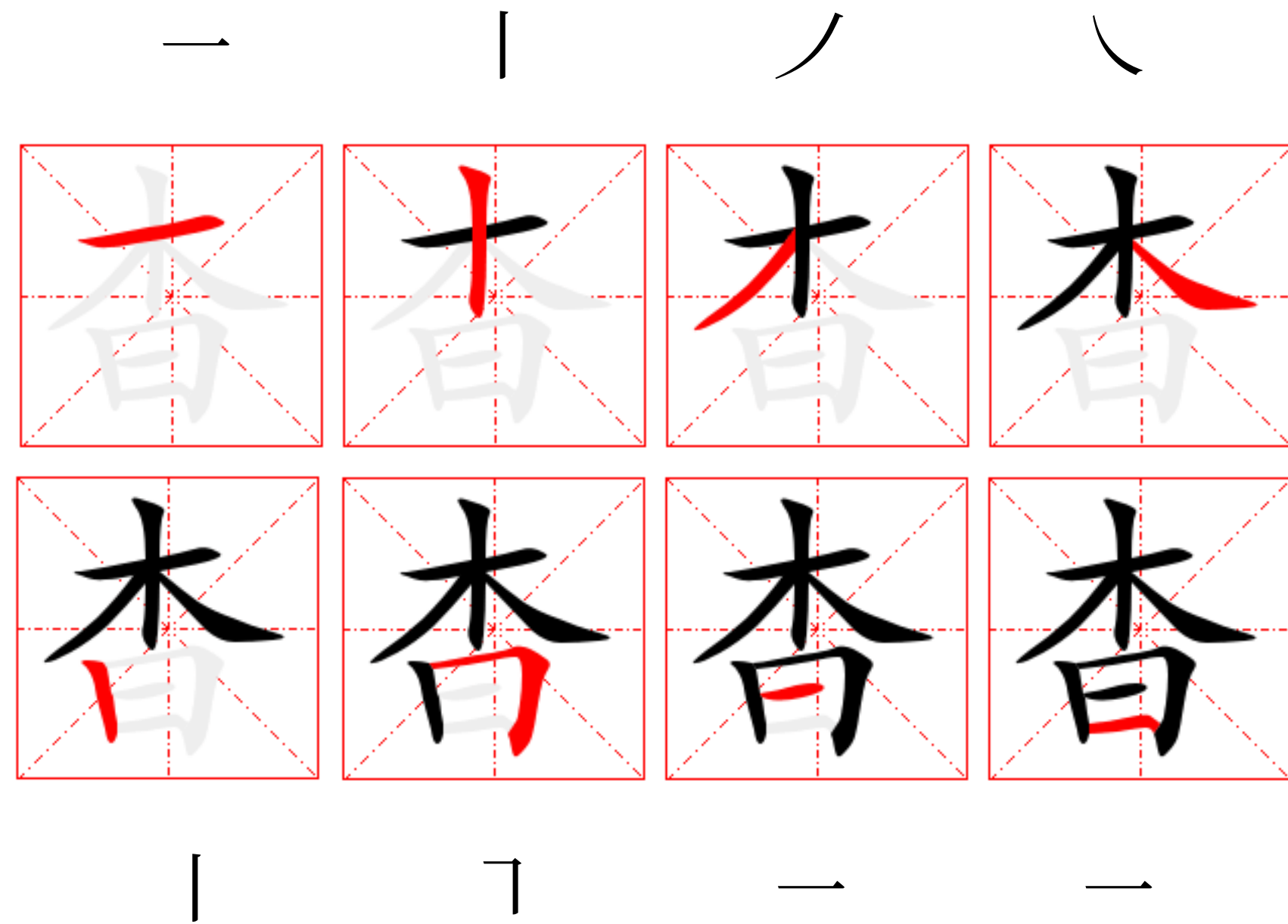
## Characters, radicals, and strokes



*Chinese characters are made up of strokes*

# Chinese 101

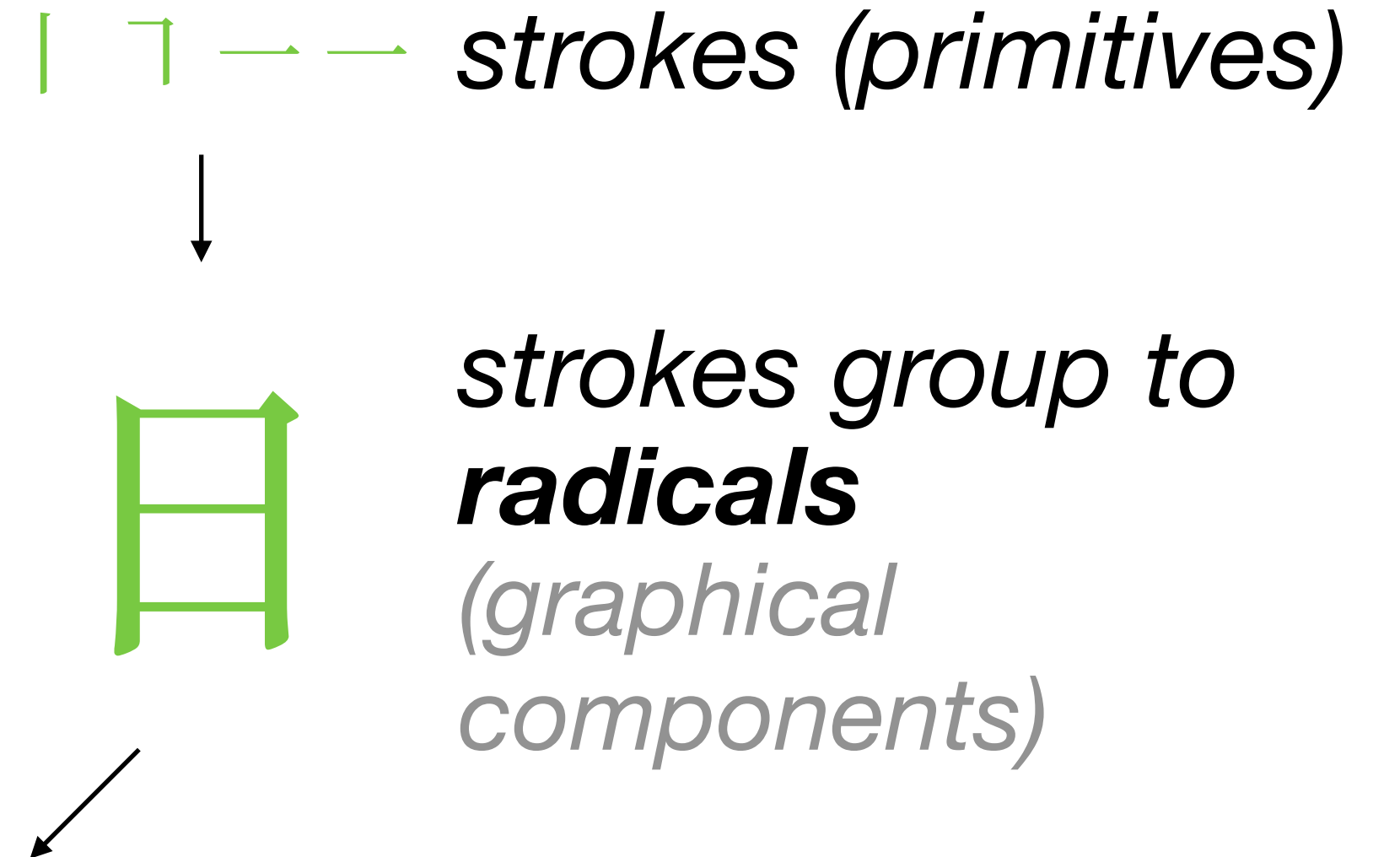
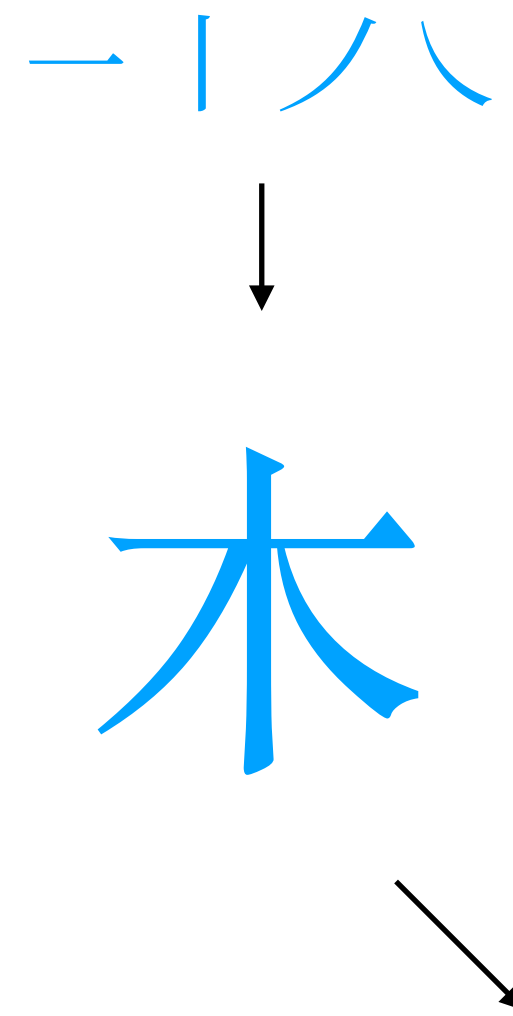
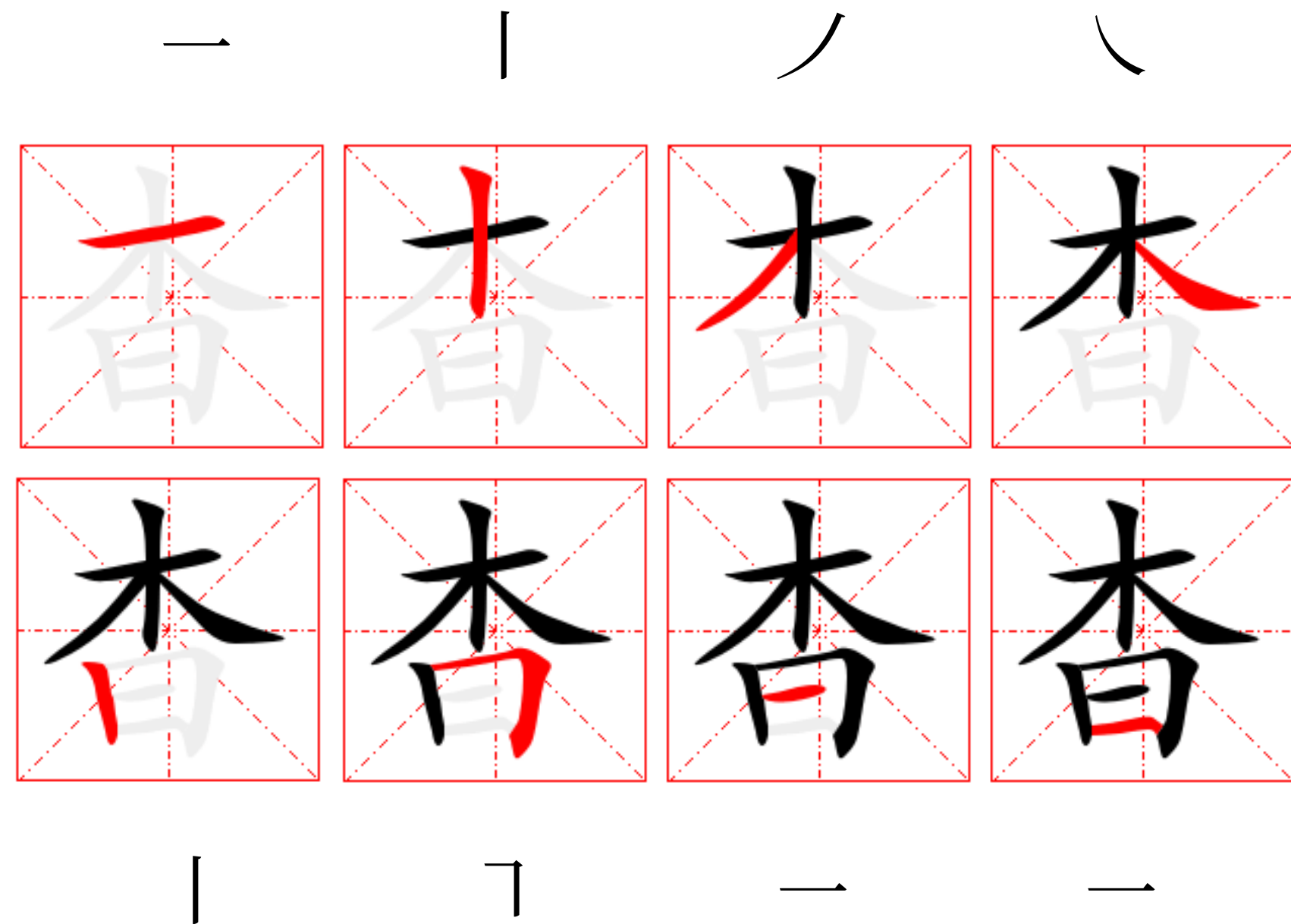
## Characters, radicals, and strokes



*Chinese characters are made up of strokes*

# Chinese 101

## Characters, radicals, and strokes



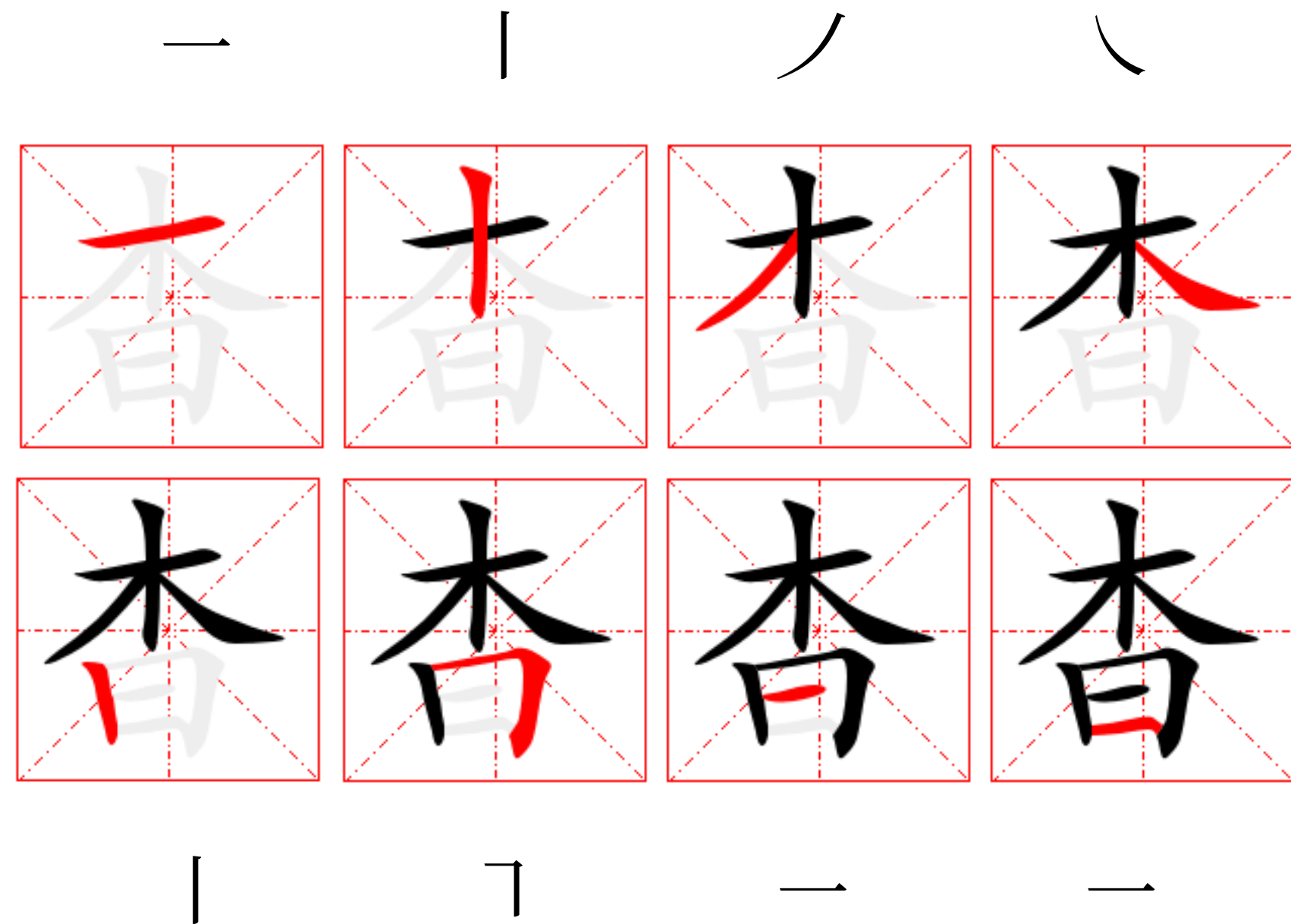
*strokes (primitives)*

*strokes group to  
**radicals**  
(graphical  
components)*

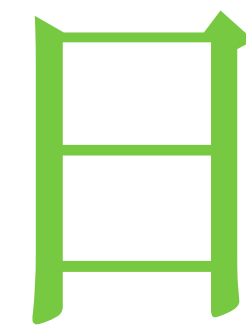
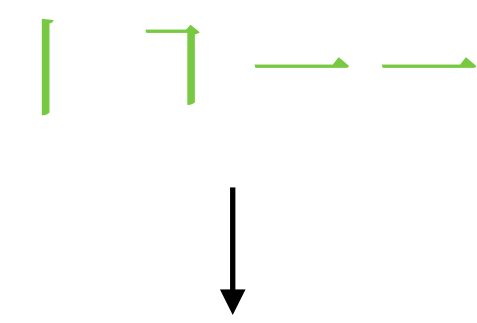
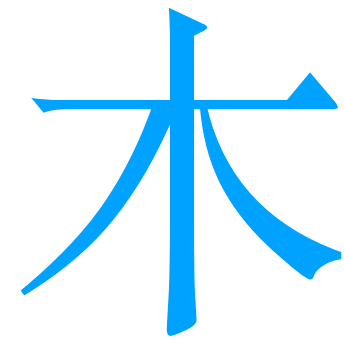
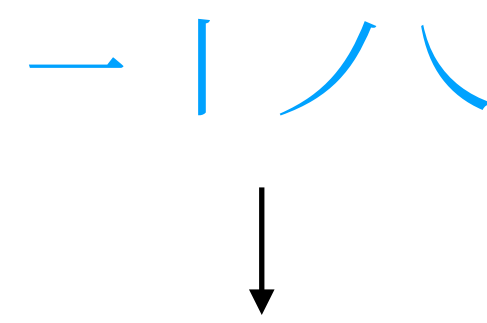
*Chinese characters are  
made up of strokes*

# Chinese 101

## Characters, radicals, and strokes



*Chinese characters are made up of strokes*



*strokes (primitives)*

*strokes group to radicals  
(graphical components)*

*radicals group to characters*

# Chinese 101

## Characters, radicals, and strokes

*empirically,  
radicals can be  
combinatorially  
reused*

木

木

目

木

目

# Chinese 101

## Characters, radicals, and strokes

*empirically,  
radicals can be  
combinatorially  
reused*

木 x 3 =

木 + 日 =

日 x 3 =

木 + 月 =

日 + 月 =



# Chinese 101

## Characters, radicals, and strokes

*empirically,  
radicals can be  
combinatorially  
reused*

木 x 3 = 森

木 + 日 = 杳

日 x 3 = 晶

木 + 月 = 朙

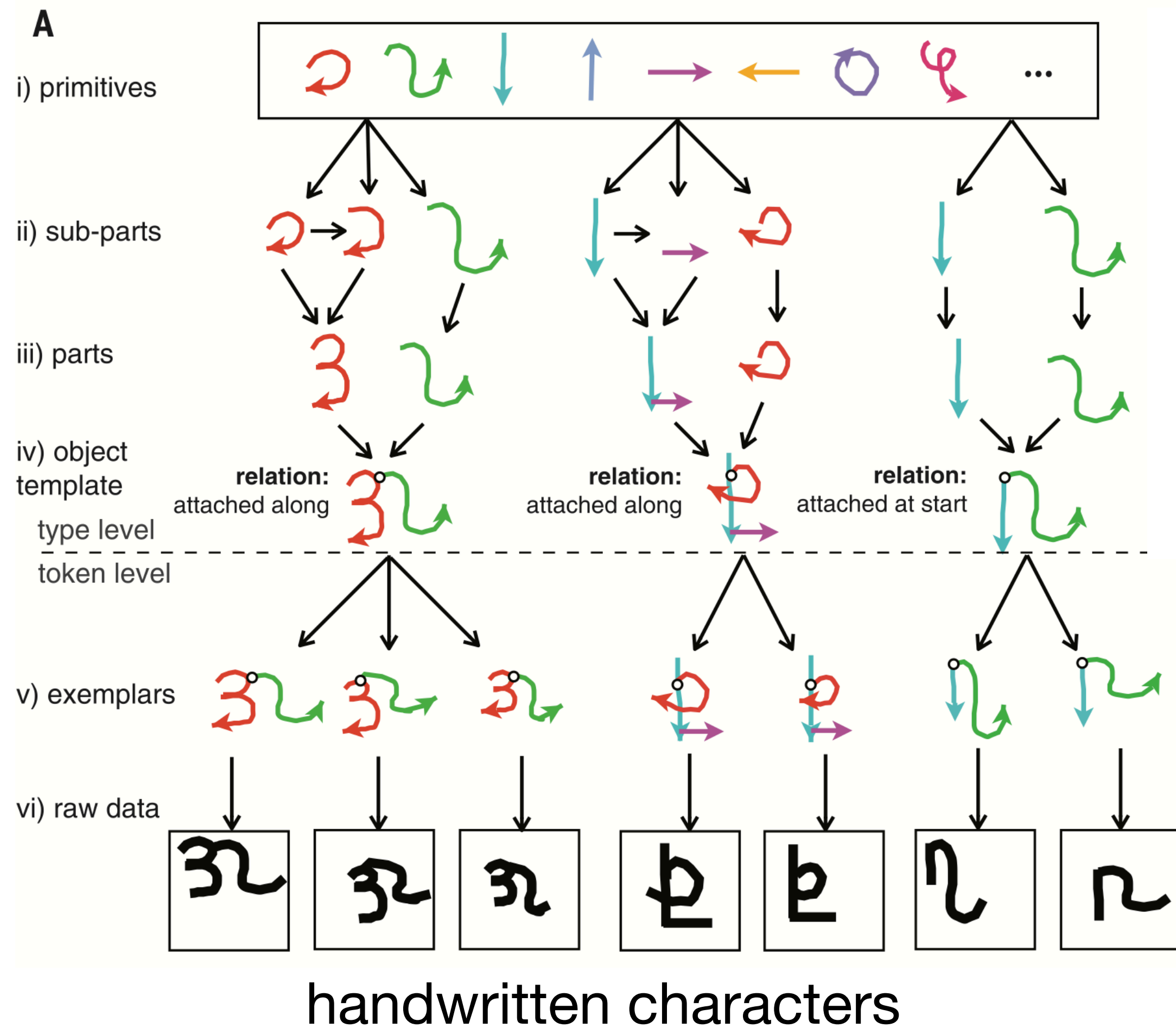
日 + 月 = 明

# Measuring representational efficiency

Inferring motor programs from images with bayesian inference

# Measuring representational efficiency

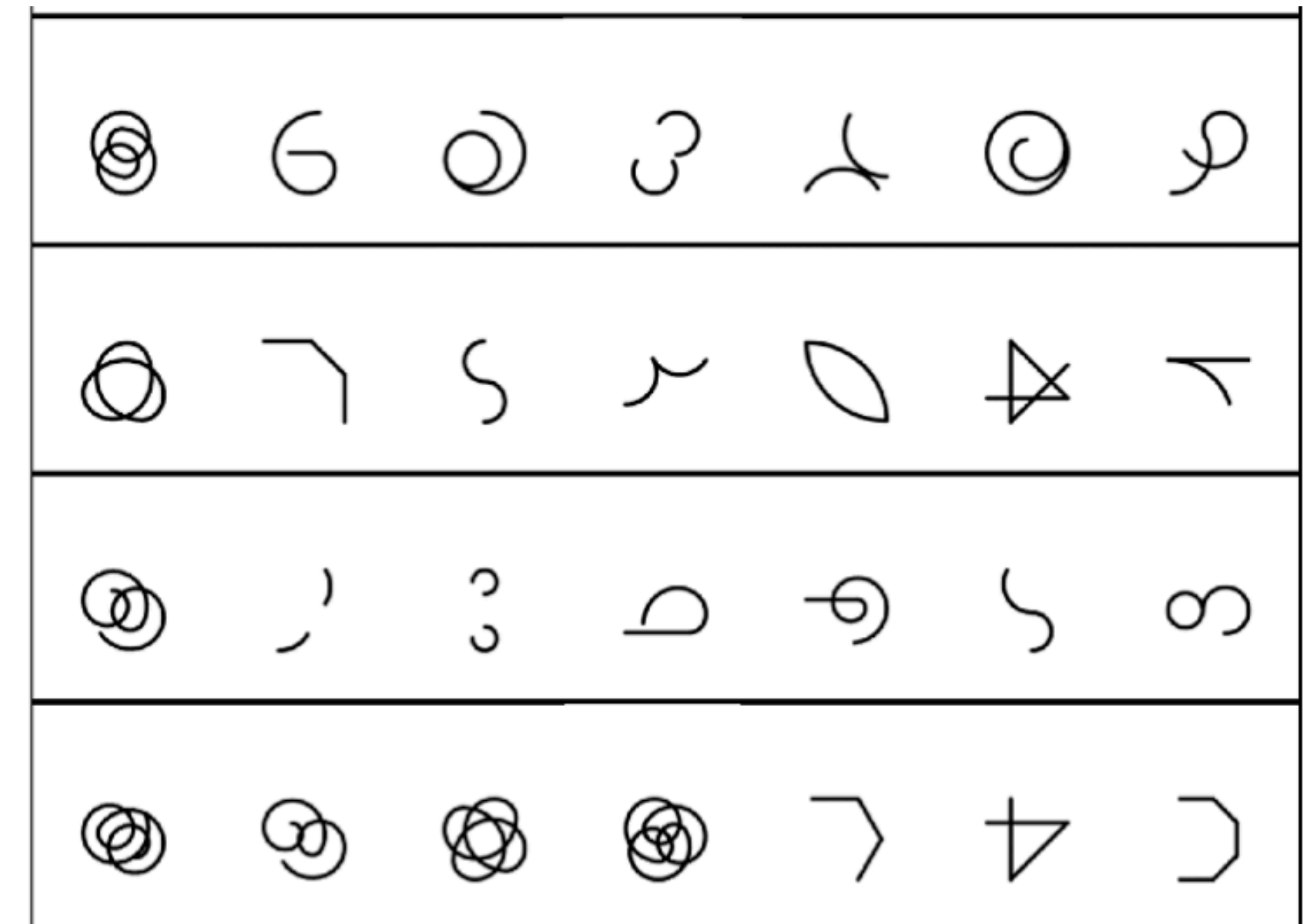
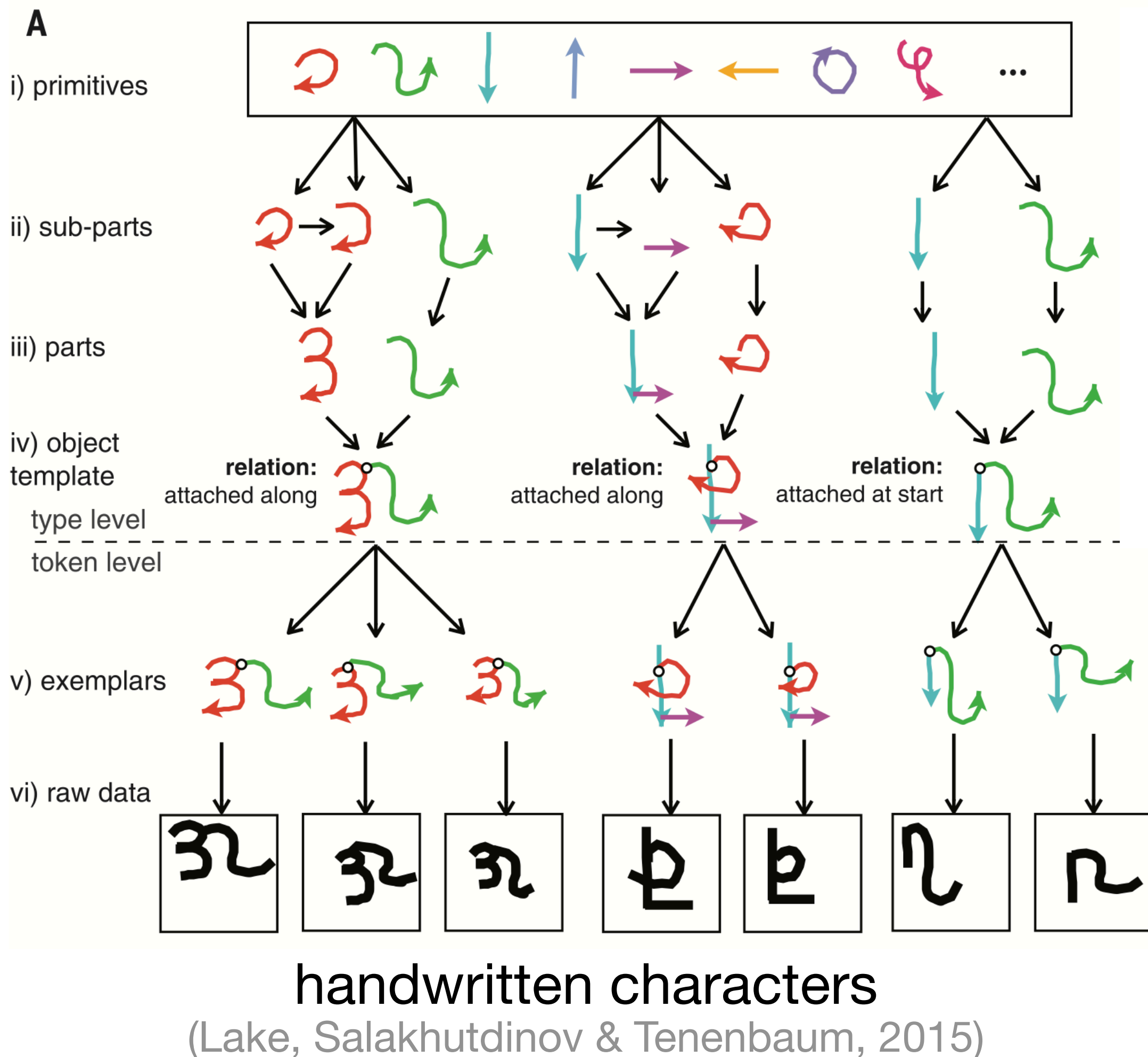
## Inferring motor programs from images with bayesian inference



(Lake, Salakhutdinov & Tenenbaum, 2015)

# Measuring representational efficiency

## Inferring motor programs from images with bayesian inference



geometric shapes

(Sablé-Meyer et al., 2022)

Shape perception formulated as searching programs with MDL.

# Contributions

# Contributions

- We **develop a library learning model** that allows us to jointly discover an underlying inventory of **higher order graphical forms** and evaluate the **MDL** of writing system represented with that set of abstract components.

# Contributions

- We **develop a library learning model** that allows us to jointly discover an underlying inventory of **higher order graphical forms** and evaluate the **MDL** of writing system represented with that set of abstract components.
- Our findings:

# Contributions

- We **develop a library learning model** that allows us to jointly discover an underlying inventory of **higher order graphical forms** and evaluate the **MDL** of writing system represented with that set of abstract components.
- Our findings:
  - Our model **rediscovers widely recognized theories** of combinatorial structure in the Chinese orthography.



# Contributions

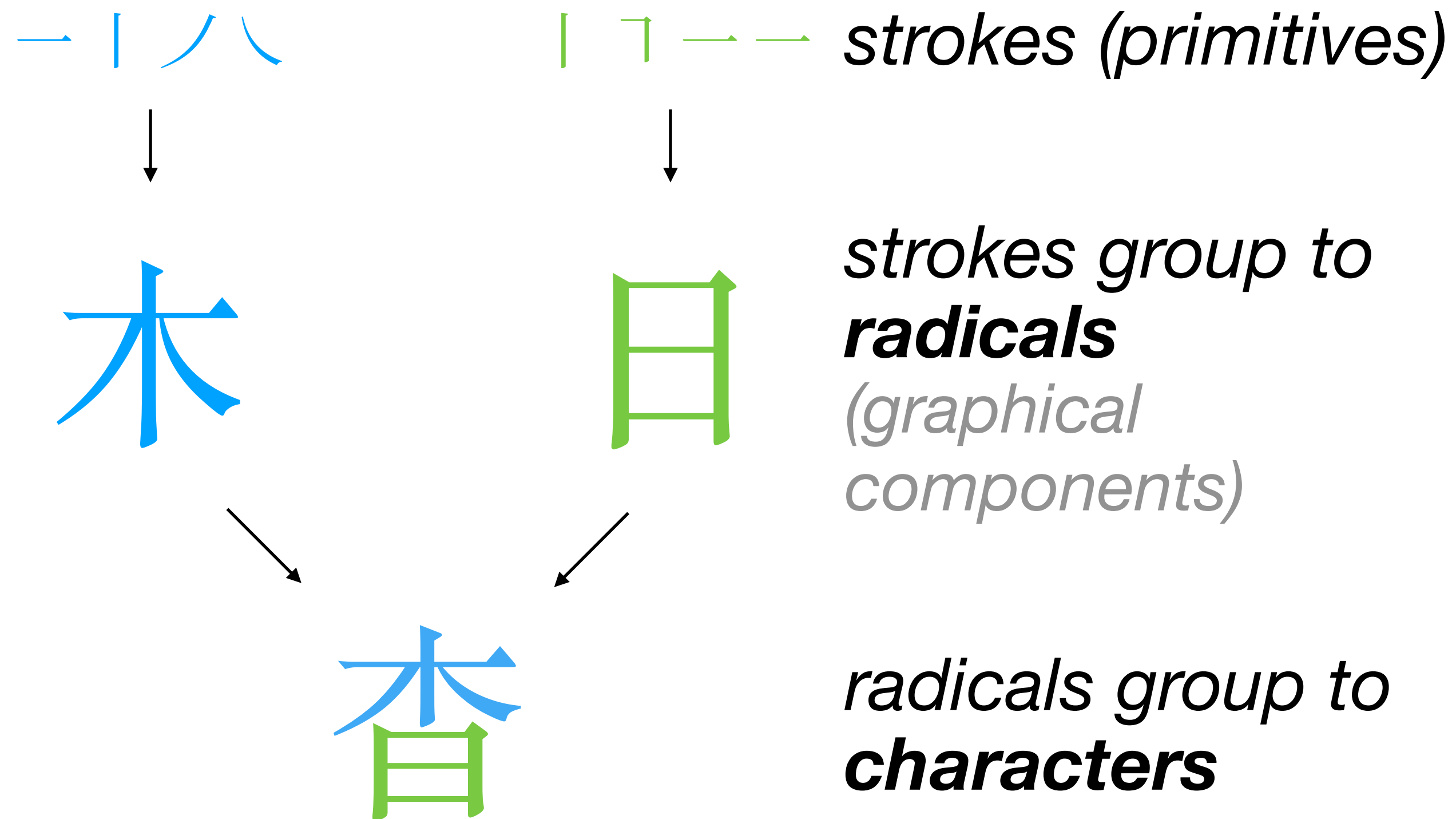
- We **develop a library learning model** that allows us to jointly discover an underlying inventory of **higher order graphical forms** and evaluate the **MDL** of writing system represented with that set of abstract components.
- Our findings:
  - Our model **rediscovers widely recognized theories** of combinatorial structure in the Chinese orthography.
  - We extend this analysis **diachronically**, investigating the **evolution of Chinese scripts** over several representative historical stages.

# Contributions

- We **develop a library learning model** that allows us to jointly discover an underlying inventory of **higher order graphical forms** and evaluate the **MDL** of writing system represented with that set of abstract components.
- Our findings:
  - Our model **rediscovers widely recognized theories** of combinatorial structure in the Chinese orthography.
  - We extend this analysis **diachronically**, investigating the **evolution of Chinese scripts** over several representative historical stages.
  - And yield an interesting diachronic finding about the relationship **between two modern Chinese scripts**.

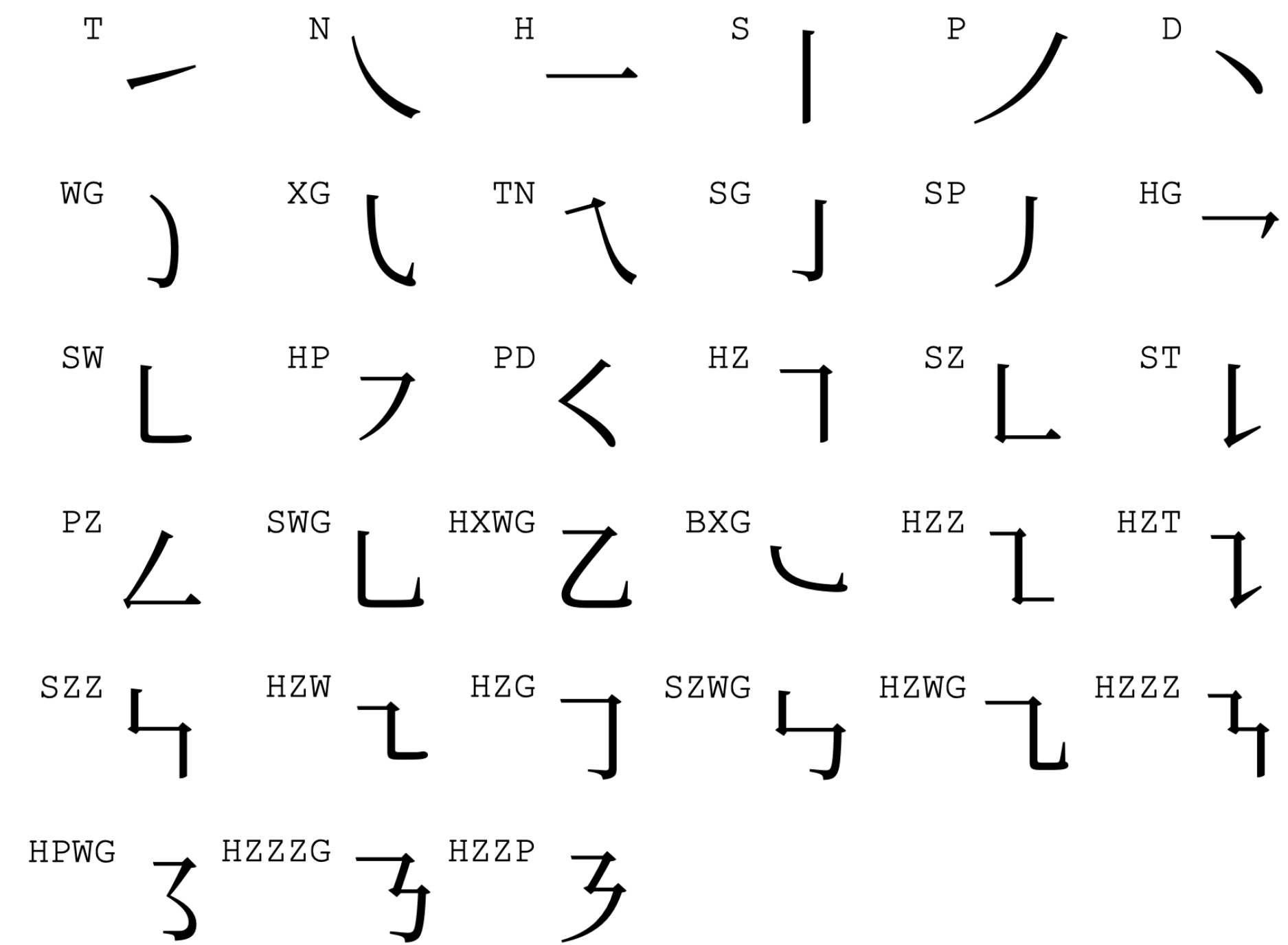
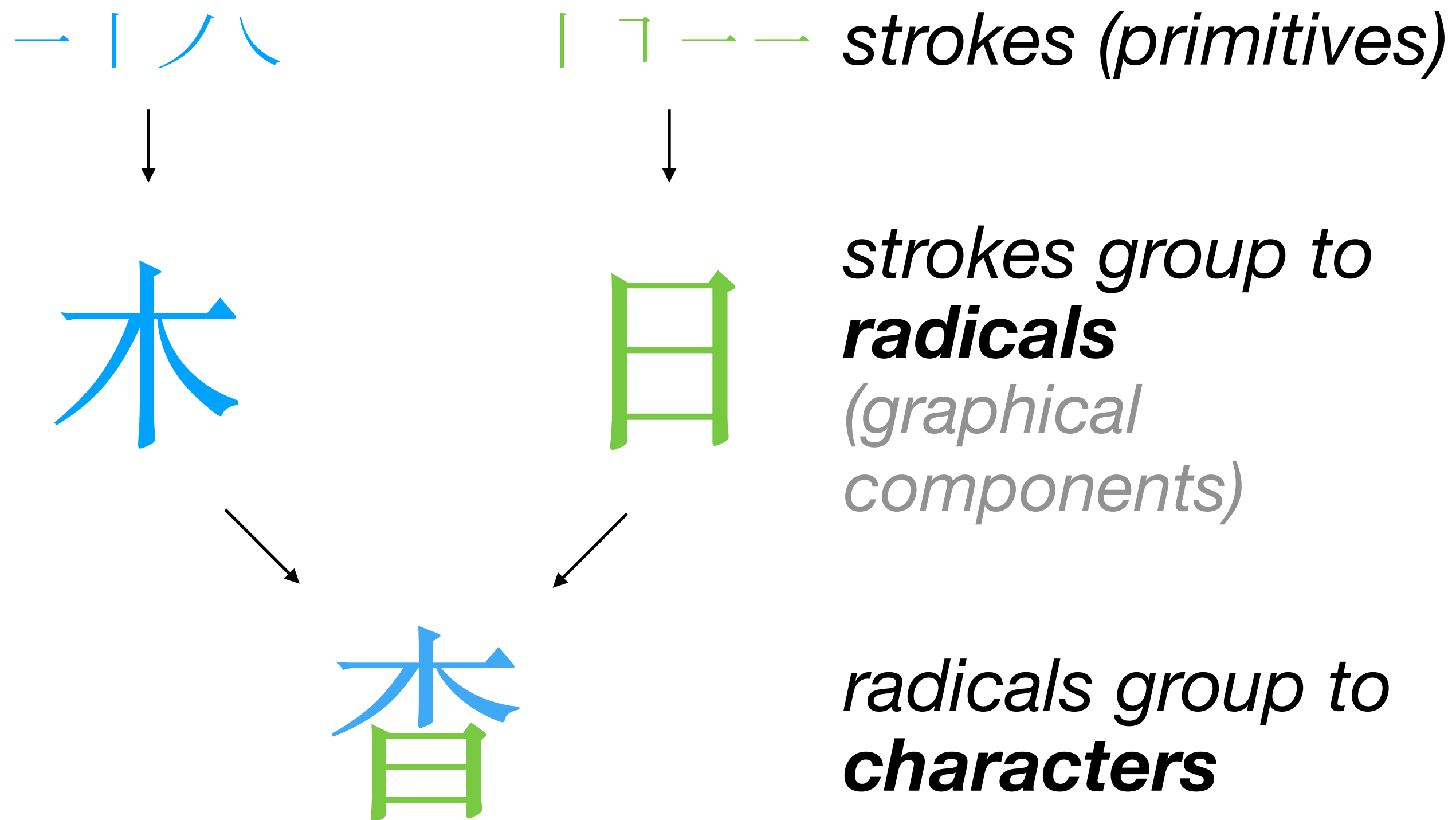
# Our Approach: Strokes

## Structure discovery with library learning



# Our Approach: Strokes

## Structure discovery with library learning



base strokes

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杏 晶 栲 明 朋

$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}base}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杳 晶 朧 明 朋  
三 丨 八 一 丨 八 丨 三 一 丨 八 丨 三 丨 三 二 丨 三 二

$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}base}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杳 晶 朧 明 朋

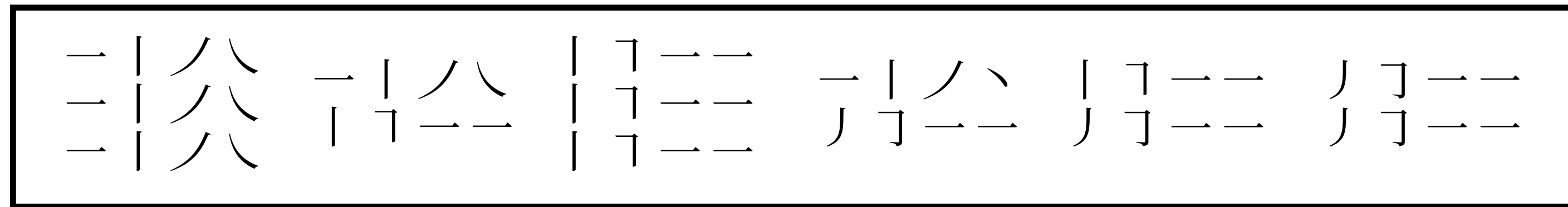
三 | 八 一 | 八 一 | 三 | 三 一 | 三 一 | 三 一 | 三 一 | 三 一 | 三 一

$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}base}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杏 晶 朧 明 朋



$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}base}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

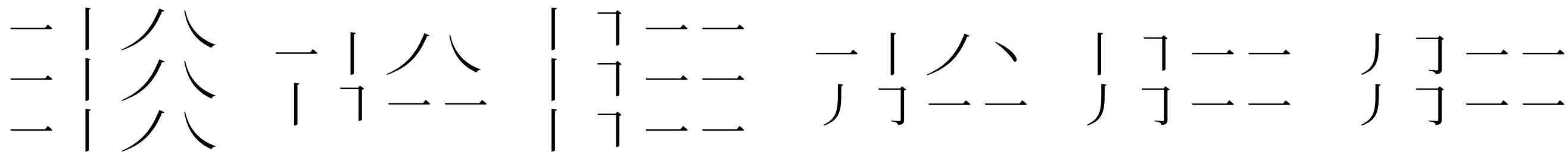
literal length = 56  
(#strokes used)



# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杏 晶 栒 明 朋



Raw DL = 56 description length of the rewritten characters

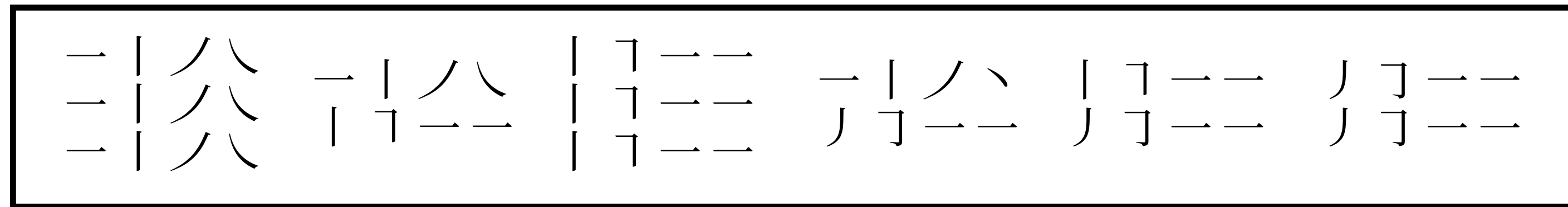
$$DL_{\mathcal{L}}(\mathcal{W}) = \sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} \overbrace{DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

literal length = 56  
(#strokes used)

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杏 晶 栊 明 朋



Raw DL = 56 description length of the rewritten characters

$$DL_{\mathcal{L}}(\mathcal{W}) = \sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))$$

literal length = 56  
(#strokes used)

Our goal: utilize a library (or a vocabulary) of patterns to efficiently represent the stroke sequences.

# Structure discovery with library learning

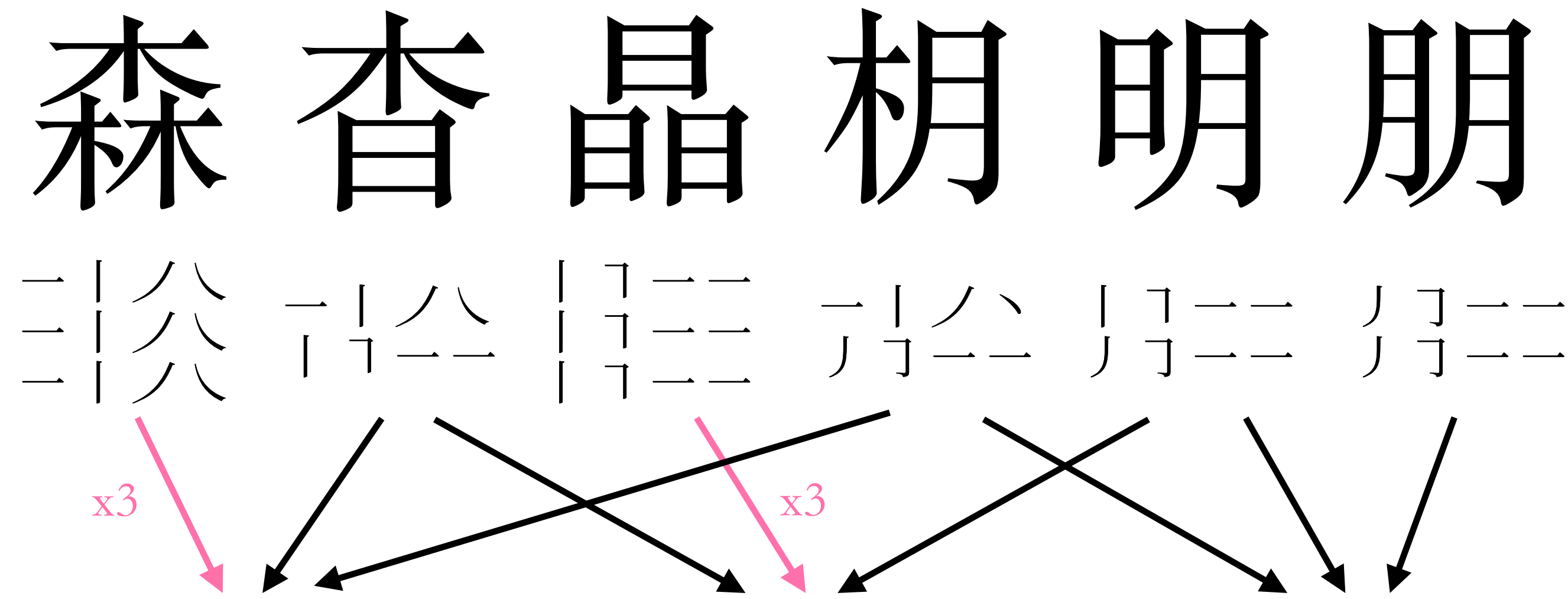
Library learning as finding MDL (minimum description length)

森 杳 晶 朧 明 朋  
三 | 八 一 | 八 一 | 三 | 三 一 | 二 二 一 | 二 二 一 | 二 二 一

$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}base}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)



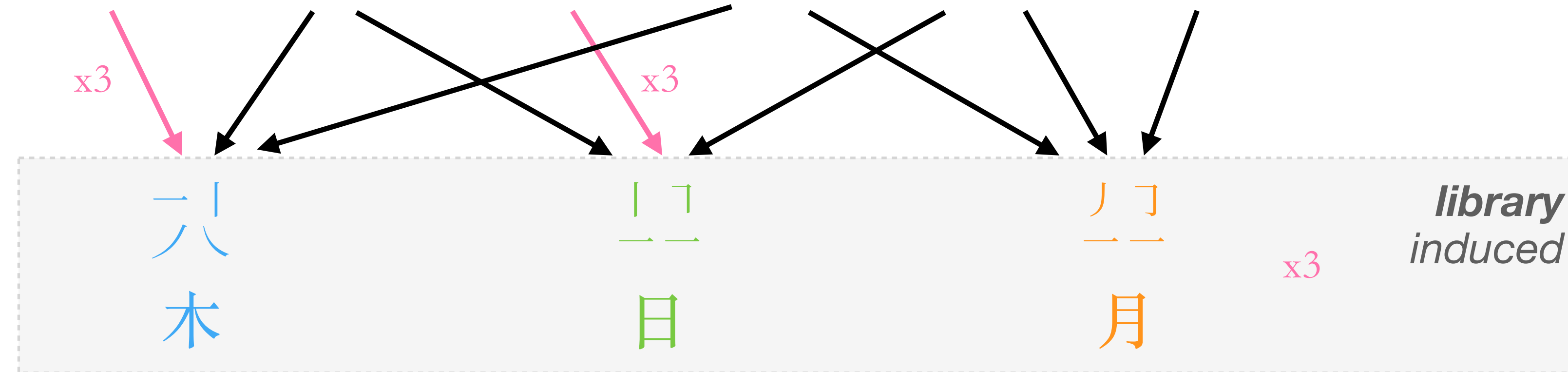
$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)

森 杏 晶 栒 明 朋

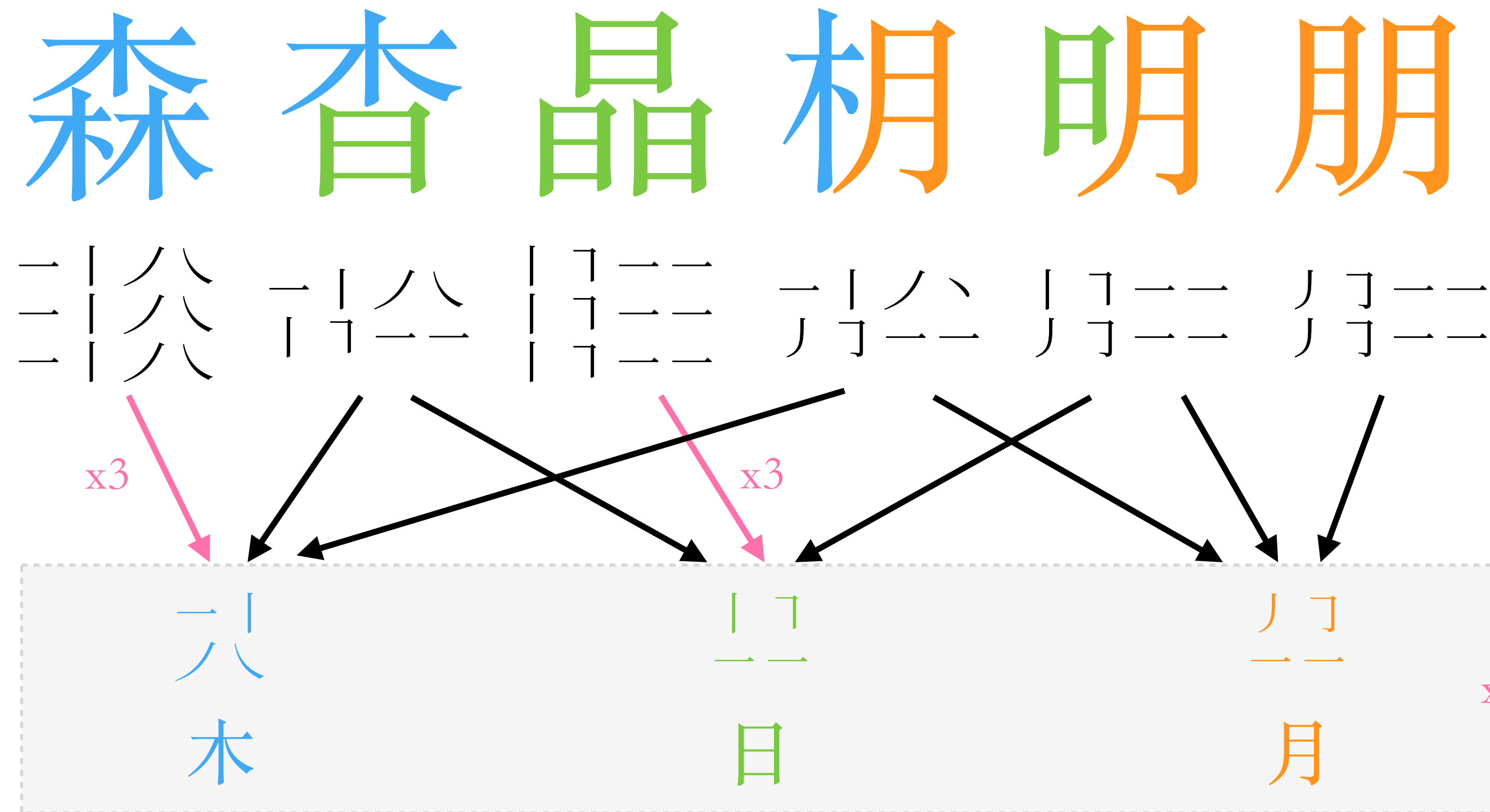
三 丨 欠 丨 丨 欠 丨 丨 三 丨 丨 欠 丨 丨 欠 丨 丨 二 丨 丨 二 丨 丨 二



$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

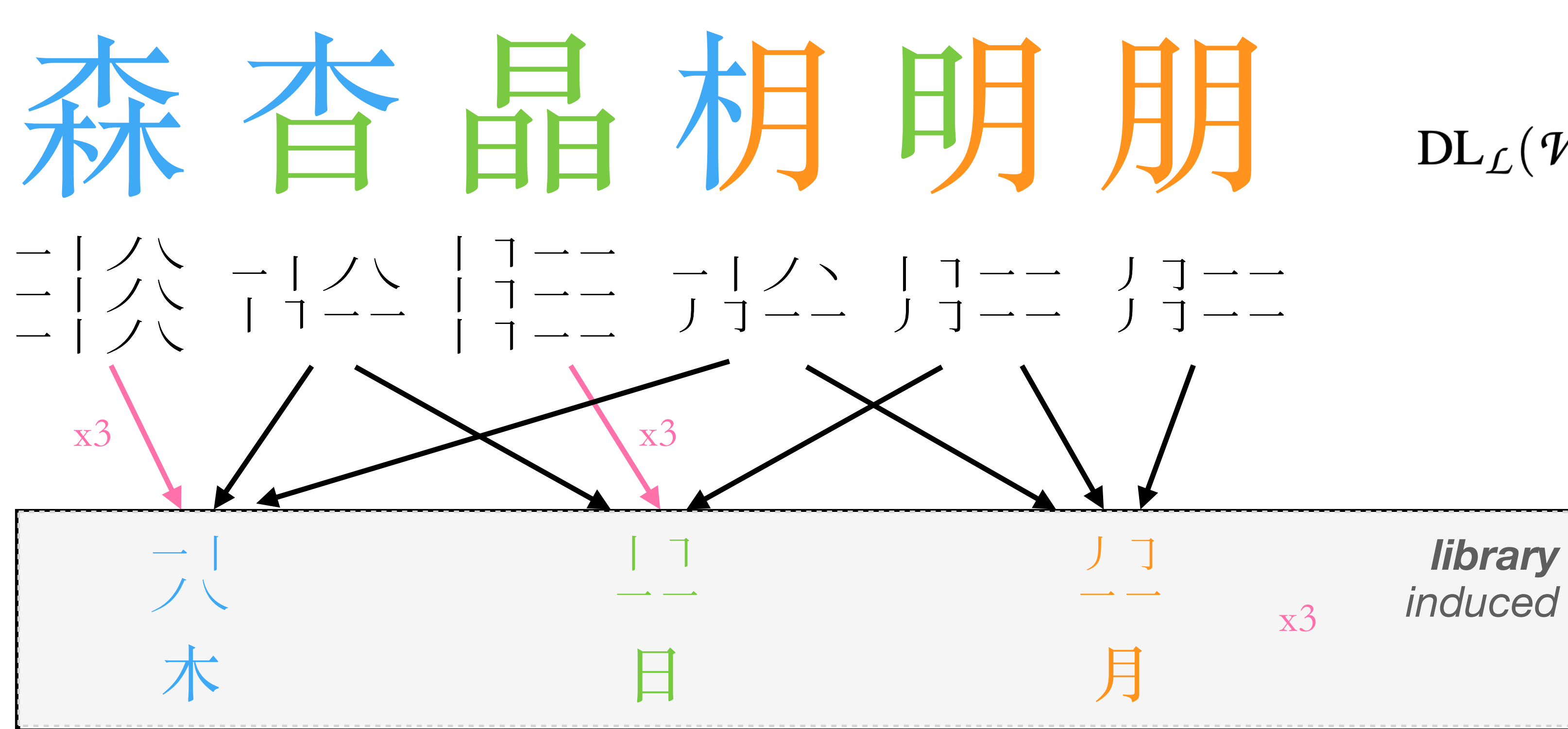
Library learning as finding MDL (minimum description length)



$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)



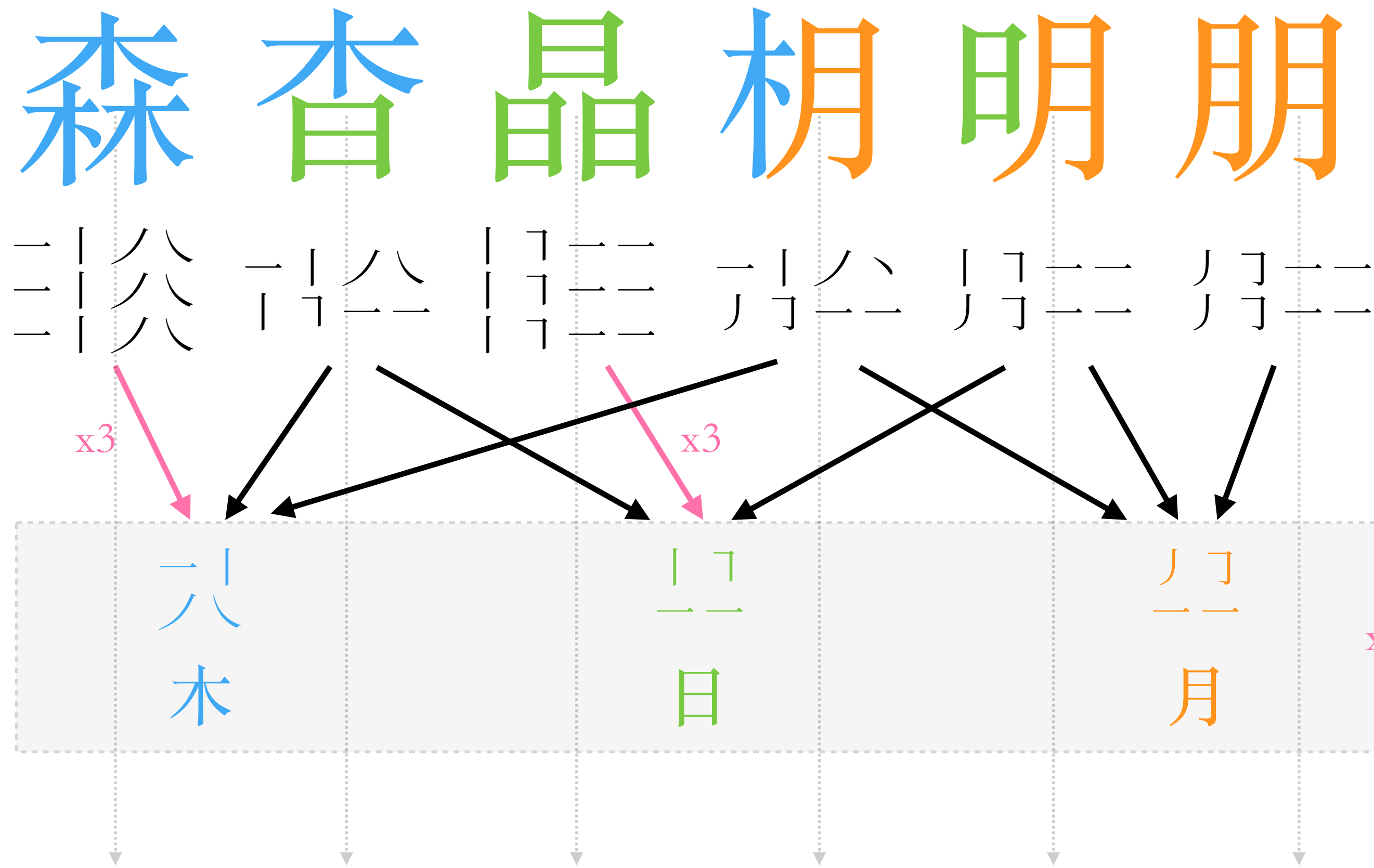
$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

A library of recurring patterns discovered.

How do we know it is a good library of patterns?

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)



$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

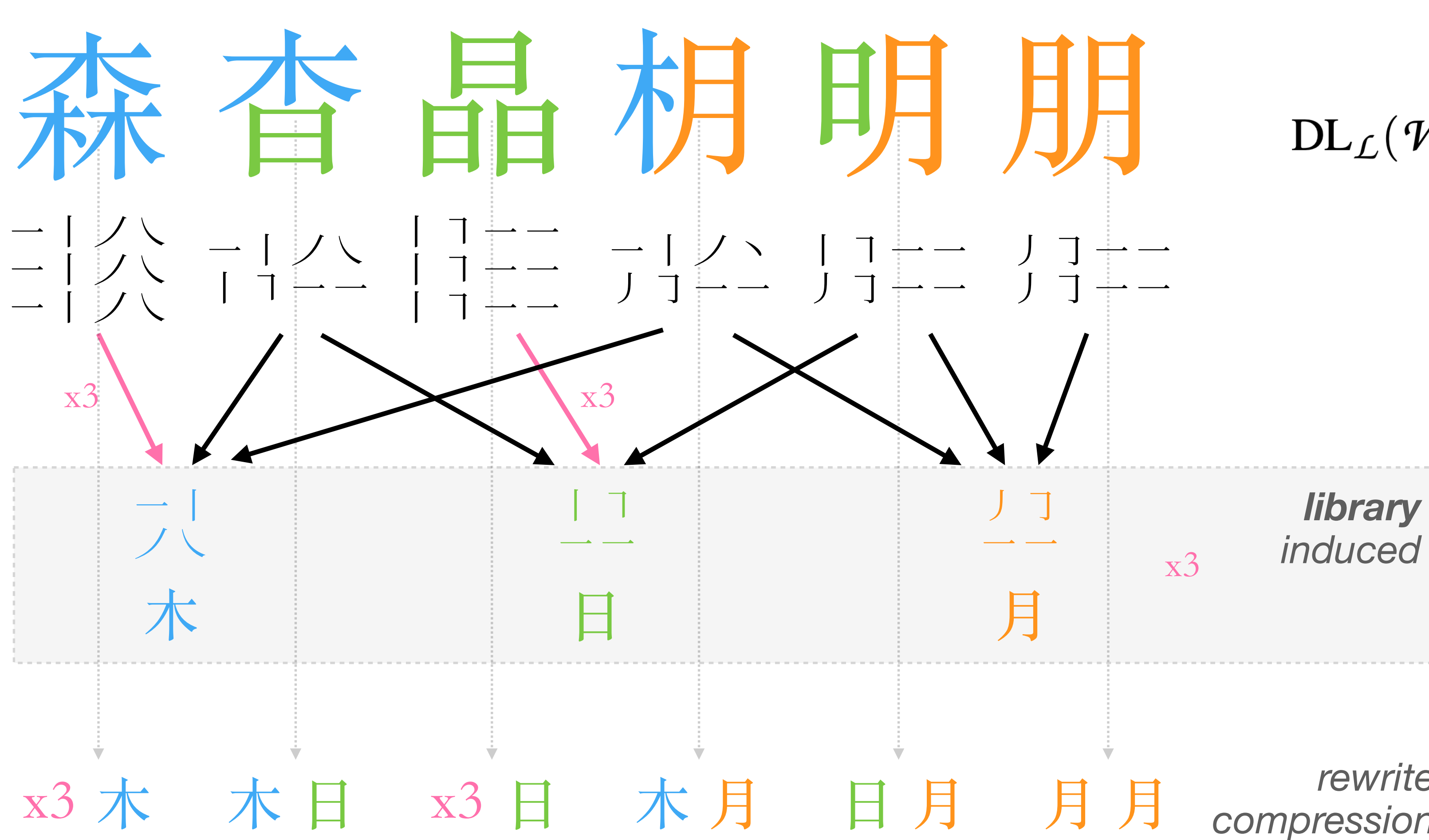
$$+ \overbrace{DL(\mathcal{L})}^{\text{description length of the library}}$$

$$DL(\mathcal{L}) = \sum_{fn \in \mathcal{L}} DL(\text{BODY}(fn))$$



# Structure discovery with library learning

Library learning as finding MDL (minimum description length)



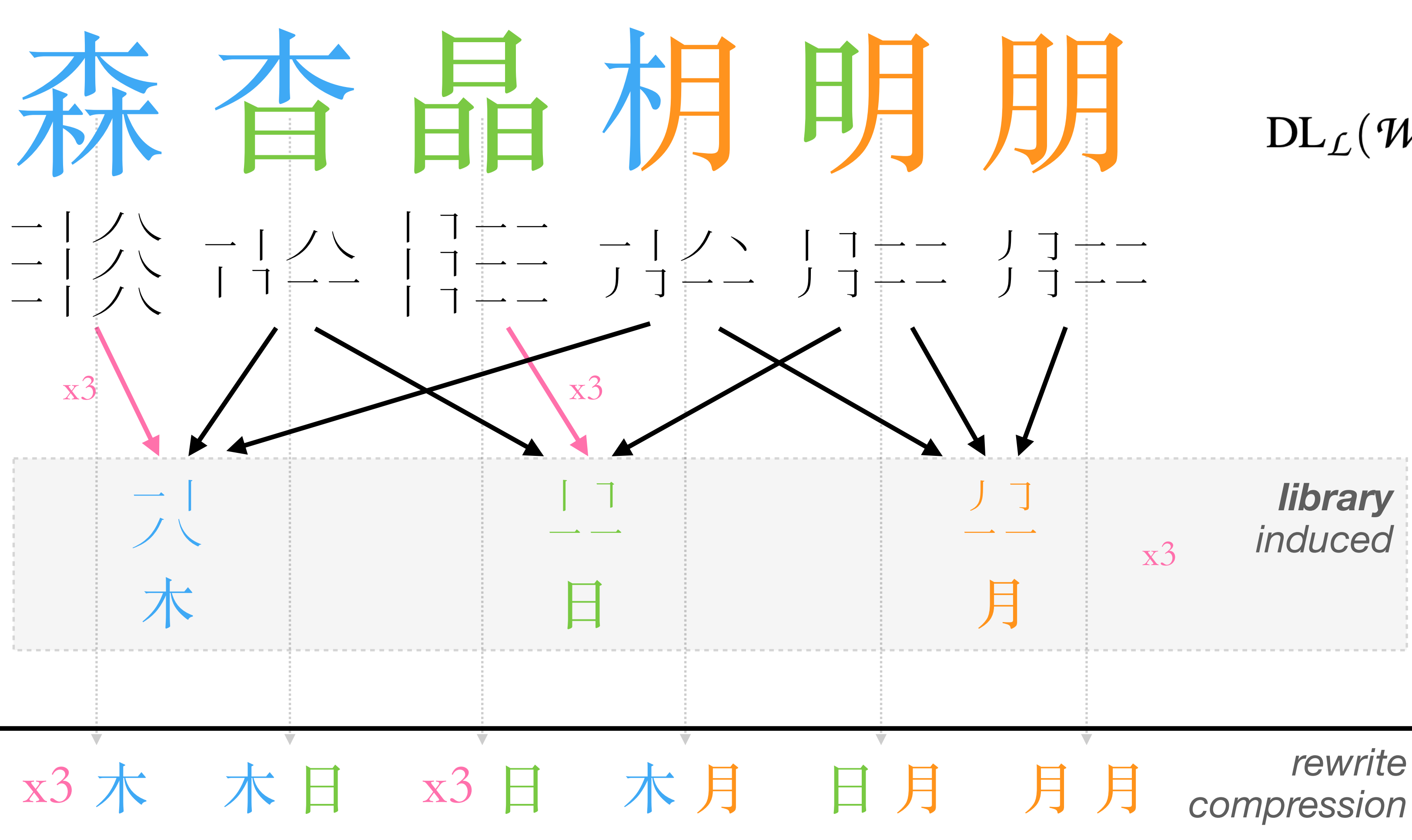
$$DL_{\mathcal{L}}(\mathcal{W}) = \overbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}^{\text{description length of the rewritten characters}}$$

$$+ \overbrace{DL(\mathcal{L})}^{\text{description length of the library}}$$

$$DL(\mathcal{L}) = \sum_{fn \in \mathcal{L}} DL(\text{BODY}(fn))$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)



$$DL_{\mathcal{L}}(\mathcal{W}) = \underbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}_{\text{description length of the rewritten characters}}$$

**compressed length = 12**  
(#primitives used)

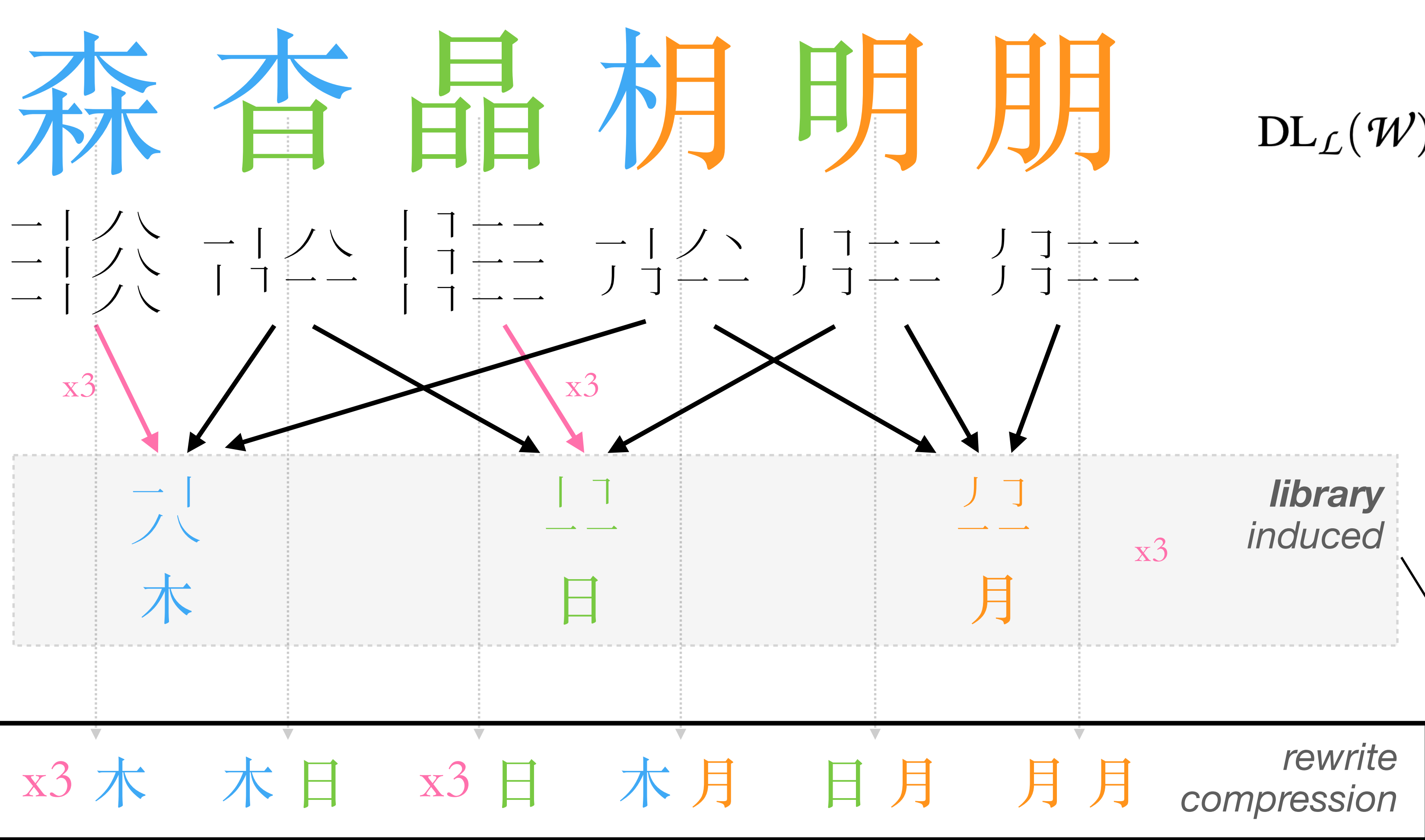
description length of the library

$$+ \overbrace{DL(\mathcal{L})}$$

$$DL(\mathcal{L}) = \sum_{fn \in \mathcal{L}} DL(\text{BODY}(fn))$$

# Structure discovery with library learning

Library learning as finding MDL (minimum description length)



$$DL_{\mathcal{L}}(\mathcal{W}) = \underbrace{\sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))}_{\text{description length of the rewritten characters}}$$

**compressed length = 12**  
(#primitives used)

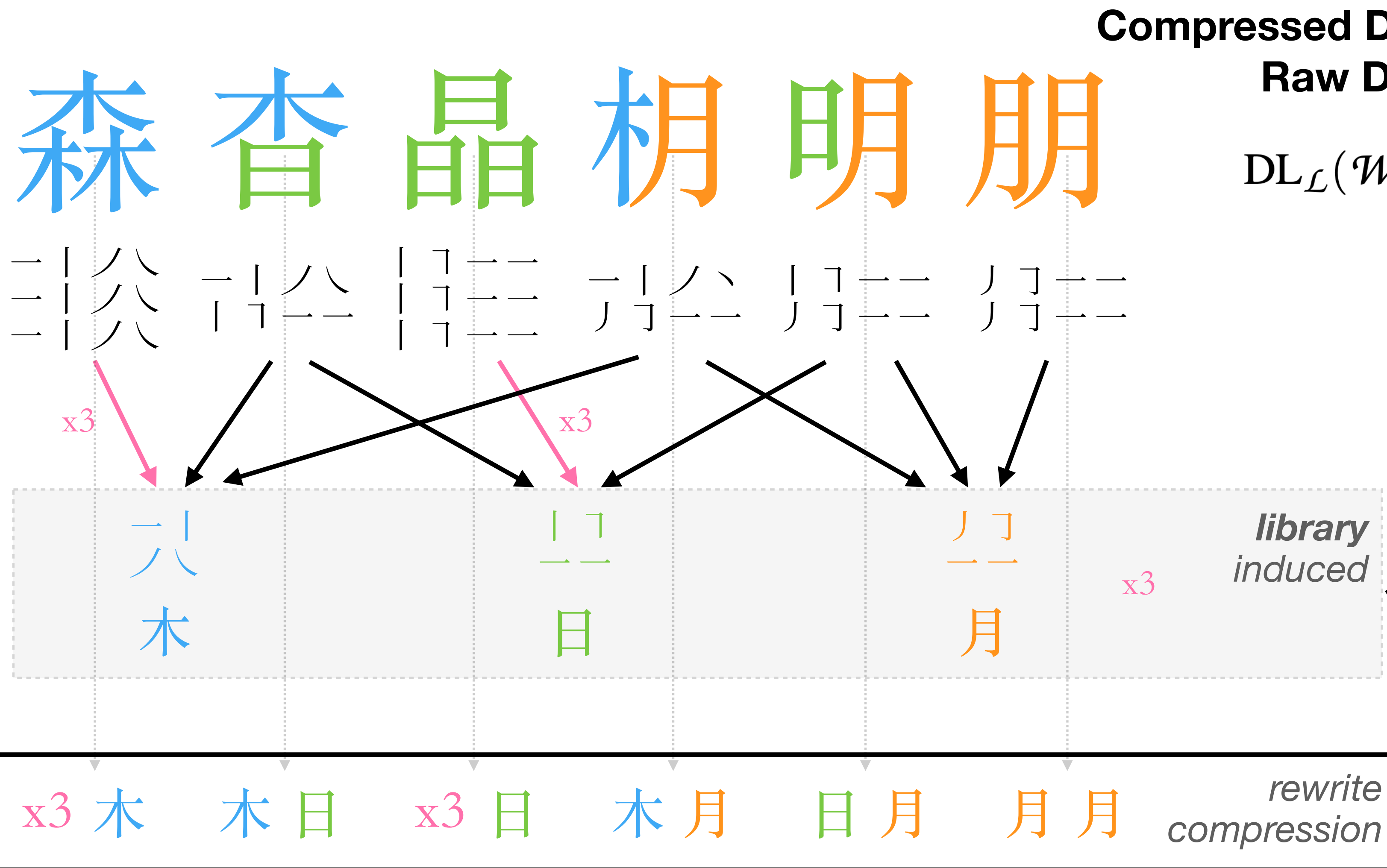
$$\text{description length of the library} + \overbrace{DL(\mathcal{L})}$$

$$DL(\mathcal{L}) = \sum_{fn \in \mathcal{L}} DL(\text{BODY}(fn))$$

**library overhead = 13**  
(sum of #primitives used in every abstraction)

# Structure discovery with library learning

## Library learning as finding MDL (minimum description length)



Compressed DL = **25** (**12** + **13**)  
 Raw DL = 56  
 description length of the rewritten characters

$$DL_{\mathcal{L}}(\mathcal{W}) = \sum_{p \in P_{\mathcal{L}_{base}}(\mathcal{W})} DL(\text{REWRITE}(p, \mathcal{L}))$$

compressed length = **12**  
 (#primitives used)

description length of the library  
 +  $DL(\mathcal{L})$

$$DL(\mathcal{L}) = \sum_{fn \in \mathcal{L}} DL(\text{BODY}(fn))$$

library overhead = **13**  
 (sum of #primitives used in every abstraction)

# At scale analysis of simplified Chinese script

oracle

seal

traditional

simplified

*sink*

*float*

*color*

*insect*

*orange*

*peace*

1500 BC

1050 BC

200 AD

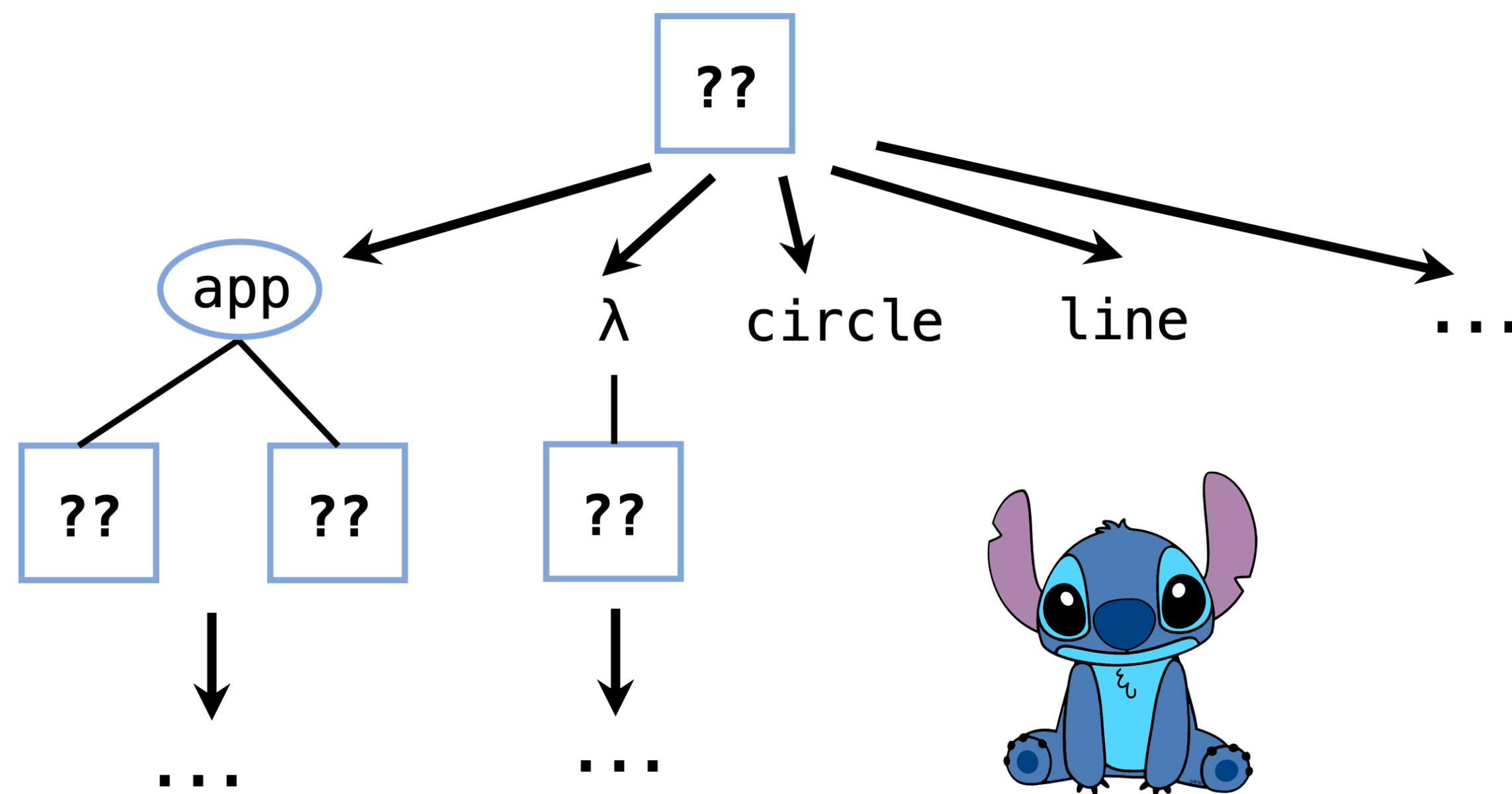
1950 AD

# At scale analysis of simplified Chinese script



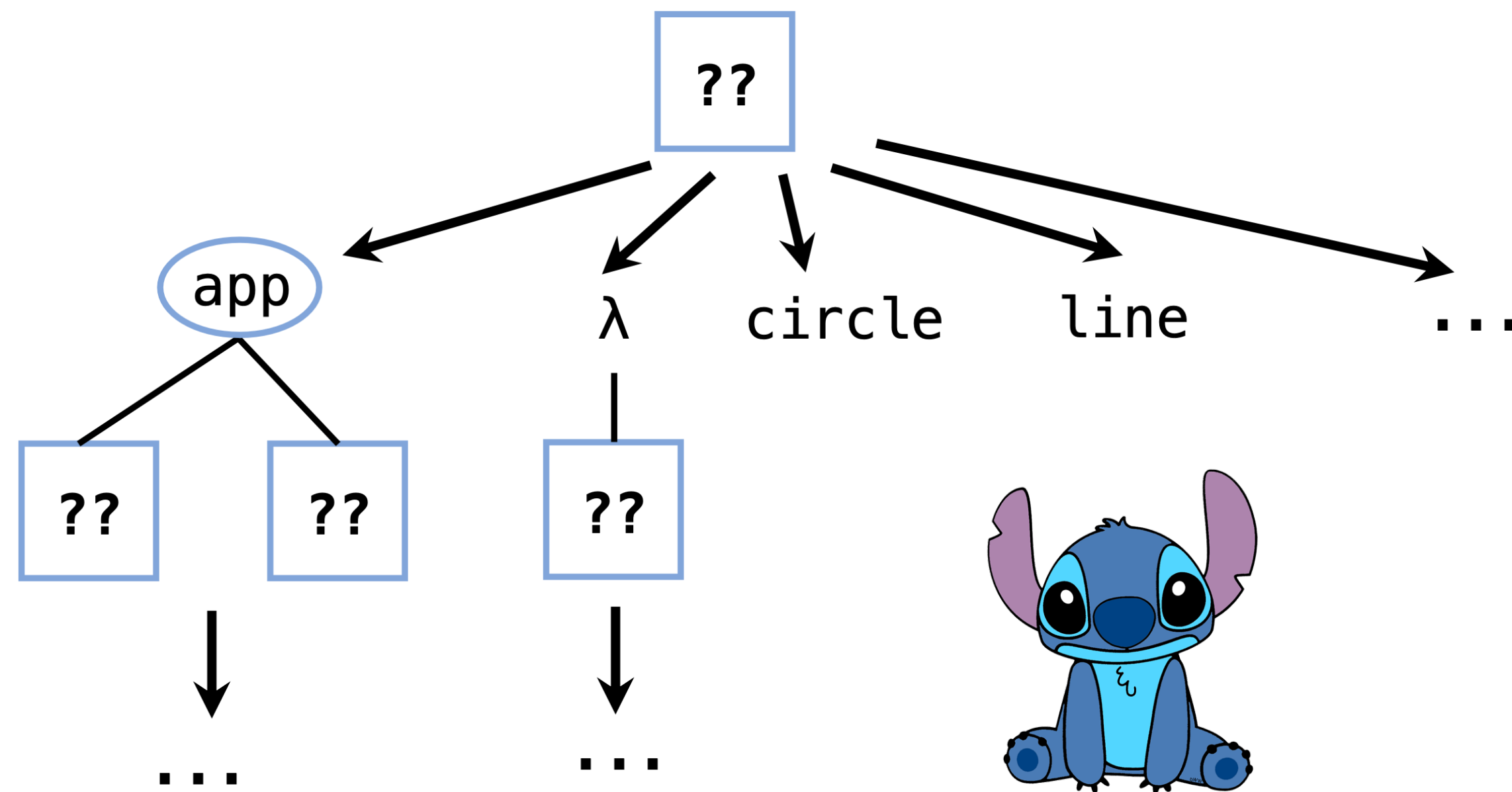
# Scaling up to the simplified Chinese script

- Leverage the **Stitch** (Bowers et al., 2023; also see DreamCoder) for efficiently discovering library functions.



# Scaling up to the simplified Chinese script

- Leverage the **Stitch** (Bowers et al., 2023; also see DreamCoder) for efficiently discovering library functions.



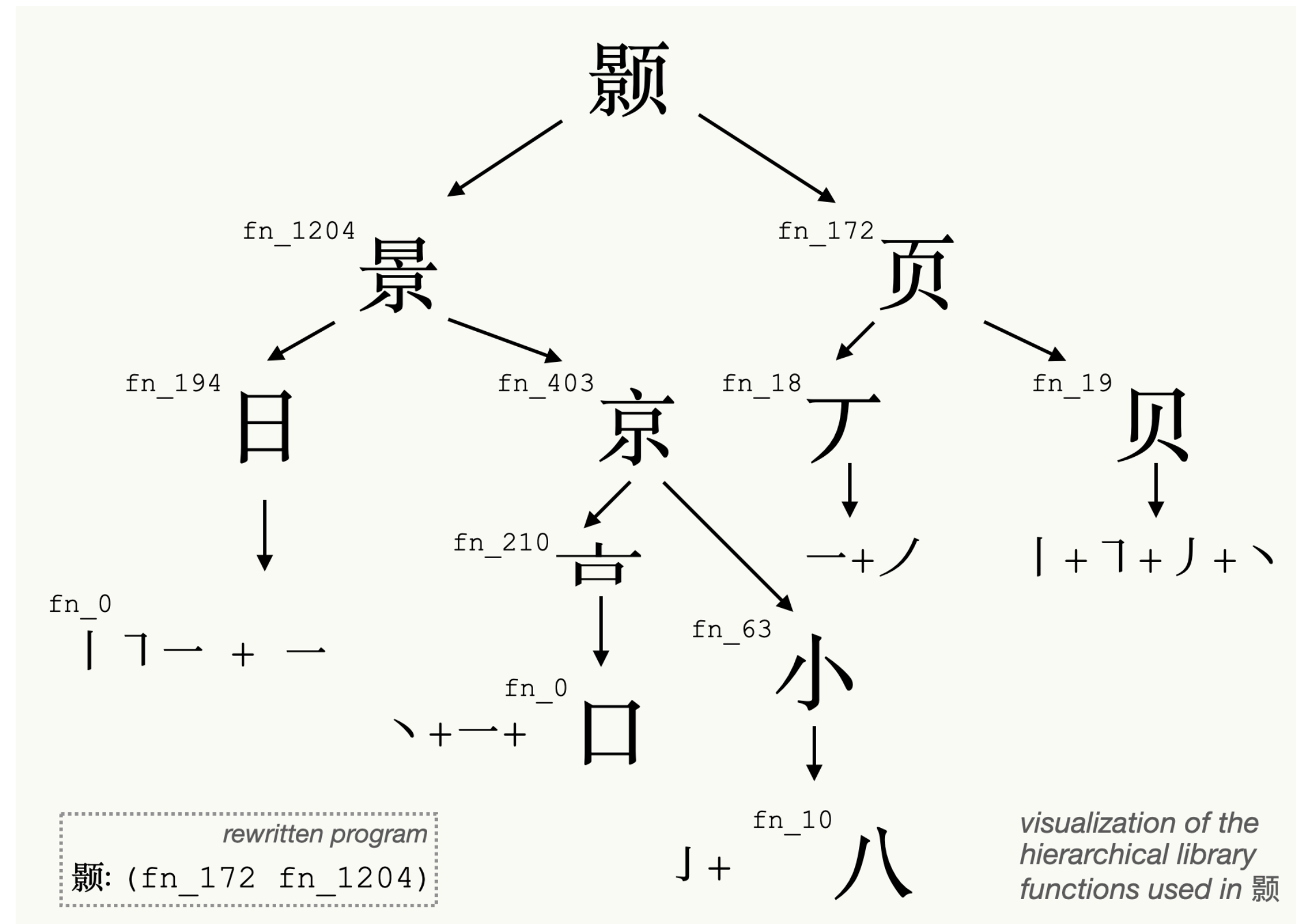
- 6,596 simplified Chinese characters. Represented as **programs**.



**Our model rediscovers widely recognized theories of combinatorial structure in the Chinese orthography.**

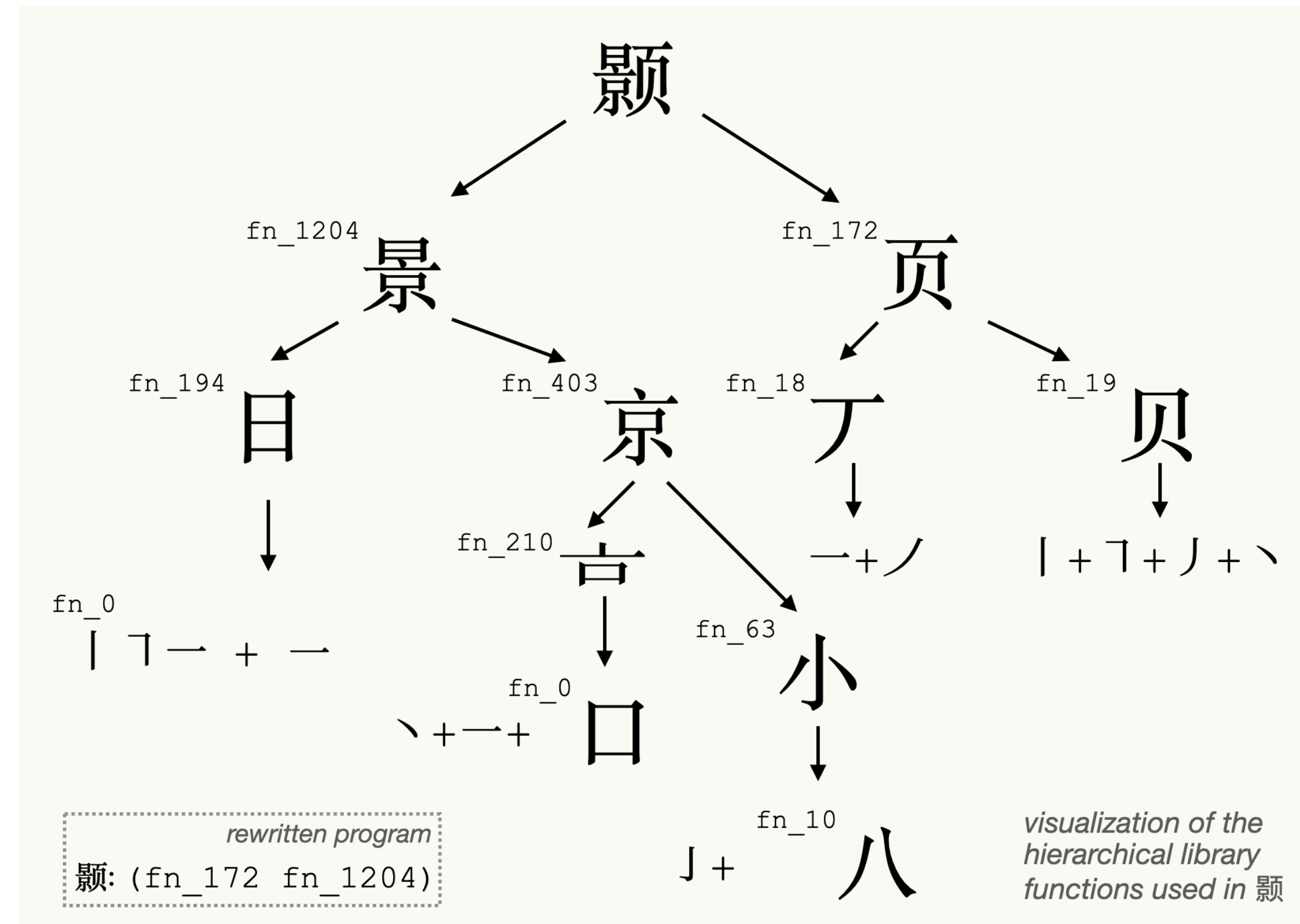
# What are the library functions learned?

hierarchically  
defined  
**graphical**  
**components**



# What are the library functions learned?

hierarchically  
defined  
**graphical  
components**



combinatorial  
**patterns** and  
**templates**

`fn_1136(#0) := (#0 (#0 (#0 list))) repeat #0 3 times`

`fn_577(#0, #1) := (lambda (#0 (#0 (#0 #1)))) repeat #0 3 times + one radical`

`fn_1135(#0, #1) := (lambda ($0 ($1 \ 丿) 一 | 丿 \)) fixed part + two radicals`

# Library functions resemble expert-defined radicals

Can library learning models uncover the structural theories underlying the Chinese language?

## Discovered and aligned radicals (187 / 201)

H	S	SP	D	HG	1	8	161	24	29	10	23	120	66	57	89
一	丨	丿	丶	㇇	十	厂	匚	卜	冂	八	人	勹	儿	匕	几
82	259	31	36	163	999	41	25	33	190	17	6	64	37	106	304
宀	冫	冫	凵	冂	刀	力	又	厶	廴	干	工	土	艹	寸	井
52	126	721	63	0	497	124	100	262	164	224	107	414	104	95	80
大	尢	弋	小	口	口	山	巾	彳	彡	夕	冬	斗	广	门	宀
49	71	77	160	173	67	169	40	219	112	295	32	378	666	27	507
辶	ヨ	尸	己	弓	子	巾	女	马	幺	《	王	无	韦	木	支
208	475	78	953	433	399	930	264	144	43	19	149	205	225	453	235
犬	歹	车	牙	戈	比	瓦	止	支	日	贝	水	见	牛	手	气
320	462	396	139	1582	582	22	782	154	679	155	510	189	125	561	371
毛	长	片	斤	爪	父	月	氏	欠	风	彡	文	方	火	斗	户
56	1423	167	178	58	309	1099	34	69	410	159	715	695	275	116	162
心	毋	示	甘	石	龙	业	目	田	罍	皿	生	矢	禾	白	鸟
42	445	242	294	177	1268	1579	372	389	88	1654	276	271	172	216	81
疒	立	穴	疋	皮	夂	矛	耒	老	耳	臣	西	而	页	至	声
113	519	581	51	330	283	995	535	1069	196	328	110	457	117	28	247
虫	缶	舌	竹	白	自	血	舟	色	衣	羊	米	聿	艮	羽	糸
1004	407	789	175	1018	86	222	675	1307	704	1354	290	590	701	293	265
麦	走	豆	酉	辰	豕	里	足	邑	身	采	谷	豸	角	言	辛
489	1192	258	228	597	1061	35	463	46	303	751	959	137	1663	783	906
青	卓	雨	非	齿	龟	隹	金	鱼	革	面	韭	骨	香	鬼	食
614	1254	129	943	468	586	287	543	1705	1252	877					
音	首	髟	高	黄	麻	鹿	黑	黍	鼓	鼻					

## Radicals failed to discover (14 / 201)

飞 瓜 肉 齐 赤 鹵 龟 阜 隶 鬲 鬥 鼎 鼠 龠

# Library functions resemble expert-defined radicals

Can library learning models uncover the structural theories underlying the Chinese language?

## Discovered and aligned radicals (187 / 201)

H	S	SP	D	HG	1	8	161	24	29	10	23	120	66	57	89
一	丨	丿	丶	フ	十	厂	匚	卜	冂	八	人	勹	儿	匕	几
82	259	31	36	163	999	41	25	33	190	17	6	64	37	106	304
宀	冫	冫	凵	冂	刀	力	又	厶	廴	干	工	土	艹	寸	井
52	126	721	63	0	497	124	100	262	164	224	107	414	104	95	80
大	尢	弋	小	口	口	山	巾	彳	彡	夕	冬	斗	广	门	宀
49	71	77	160	173	67	169	40	219	112	295	32	378	666	27	507
辶	ヨ	尸	己	弓	子	巾	女	马	幺	《	王	无	韦	木	支
208	475	78	953	433	399	930	264	144	43	19	149	205	225	453	235
犬	歹	车	牙	戈	比	瓦	止	支	日	贝	水	见	牛	手	气
320	462	396	139	1582	582	22	782	154	679	155	510	189	125	561	371
毛	长	片	斤	爪	父	月	氏	欠	风	彡	文	方	火	斗	户
56	1423	167	178	58	309	1099	34	69	410	159	715	695	275	116	162
心	毋	示	甘	石	龙	业	目	田	罍	皿	生	矢	禾	白	鸟
42	445	242	294	177	1268	1579	372	389	88	1654	276	271	172	216	81
疒	立	穴	疋	皮	夂	矛	耒	老	耳	臣	西	而	页	至	声
113	519	581	51	330	283	995	535	1069	196	328	110	457	117	28	247
虫	缶	舌	竹	白	自	血	舟	色	衣	羊	米	聿	艮	羽	糸
1004	407	789	175	1018	86	222	675	1307	704	1354	290	590	701	293	265
麦	走	豆	酉	辰	豕	里	足	邑	身	采	谷	豸	角	言	辛
489	1192	258	228	597	1061	35	463	46	303	751	959	137	1663	783	906
青	卓	雨	非	齿	龟	隹	金	鱼	革	面	韭	骨	香	鬼	食
614	1254	129	943	468	586	287	543	1705	1252	877					
音	首	髟	高	黄	麻	鹿	黑	黍	鼓	鼻					

## Radicals failed to discover (14 / 201)

飞 瓜 肉 齐 赤 鹵 龟 阜 隶 鬲 鬥 鼎 鼠 龠

- Our model discovered 187 (93.0%) radicals defined by experts.

# Library functions resemble expert-defined radicals

Can library learning models uncover the structural theories underlying the Chinese language?

## Discovered and aligned radicals (187 / 201)

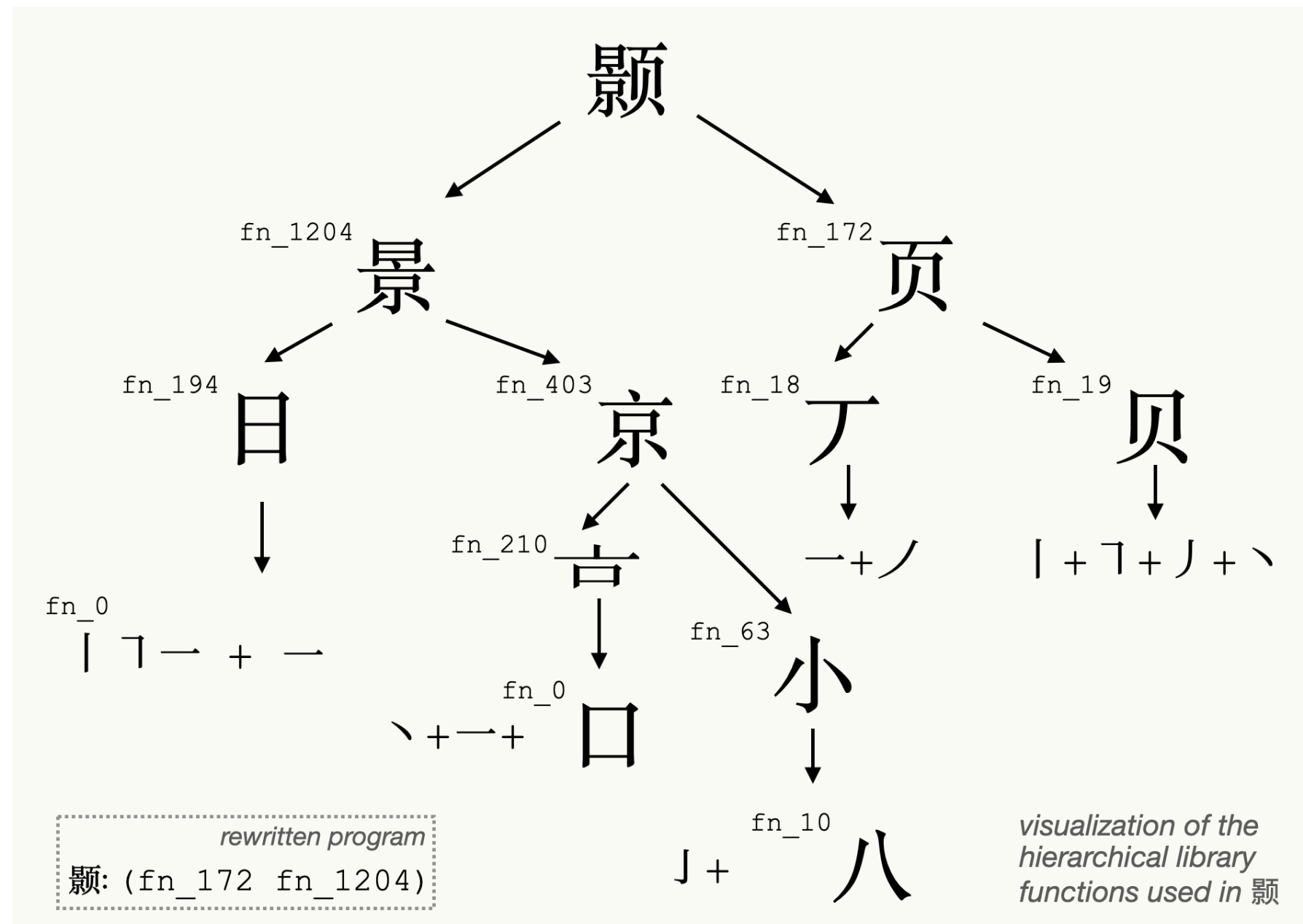
H	S	SP	D	HG	1	8	161	24	29	10	23	120	66	57	89
一	丨	丿	丶	フ	十	厂	匚	卜	冂	八	人	勹	儿	匕	几
82	259	31	36	163	999	41	25	33	190	17	6	64	37	106	304
二	冫	冫	凵	冂	刀	力	又	厶	廴	干	工	土	艹	寸	井
52	126	721	63	0	497	124	100	262	164	224	107	414	104	95	80
大	尢	弋	小	口	口	山	巾	彳	彡	夕	冬	斗	广	门	宀
49	71	77	160	173	67	169	40	219	112	295	32	378	666	27	507
辶	ヨ	尸	己	弓	子	巾	女	马	幺	《	王	无	韦	木	支
208	475	78	953	433	399	930	264	144	43	19	149	205	225	453	235
犬	歹	车	牙	戈	比	瓦	止	支	日	贝	水	见	牛	手	气
320	462	396	139	1582	582	22	782	154	679	155	510	189	125	561	371
毛	长	片	斤	爪	父	月	氏	欠	风	彡	文	方	火	斗	户
56	1423	167	178	58	309	1099	34	69	410	159	715	695	275	116	162
心	毋	示	甘	石	龙	业	目	田	罍	皿	生	矢	禾	白	鸟
42	445	242	294	177	1268	1579	372	389	88	1654	276	271	172	216	81
疒	立	穴	疋	皮	夂	矛	耒	老	耳	臣	西	而	页	至	声
113	519	581	51	330	283	995	535	1069	196	328	110	457	117	28	247
虫	缶	舌	竹	白	自	血	舟	色	衣	羊	米	聿	艮	羽	糸
1004	407	789	175	1018	86	222	675	1307	704	1354	290	590	701	293	265
麦	走	豆	酉	辰	豕	里	足	邑	身	采	谷	豸	角	言	辛
489	1192	258	228	597	1061	35	463	46	303	751	959	137	1663	783	906
青	卓	雨	非	齿	龟	隹	金	鱼	革	面	韭	骨	香	鬼	食
614	1254	129	943	468	586	287	543	1705	1252	877					
音	首	髟	高	黄	麻	鹿	黑	黍	鼓	鼻					

## Radicals failed to discover (14 / 201)

飞 瓜 肉 齐 赤 鹵 龟 阜 隶 鬲 鬥 鼎 鼠 龠

- Our model **discovered 187 (93.0%) radicals defined by experts.**
- Recap: radicals are
- **Graphical components discovered by experts that frequently occur in Chinese characters.**

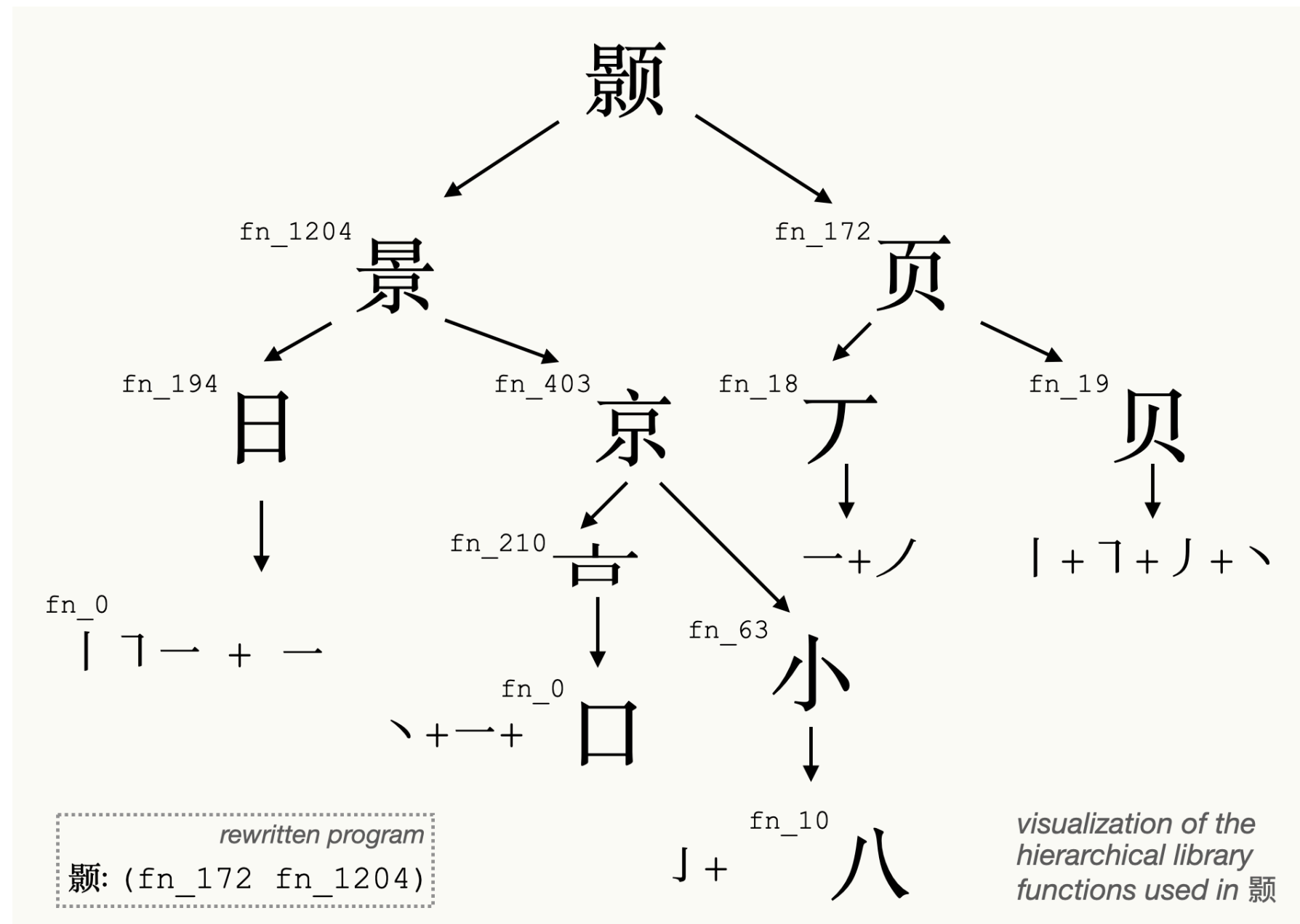
# Learned library captures the hierarchical organization of simplified Chinese characters



Model	F <sub>1</sub>
Library learning	61.6
Baselines	
– Balanced binary tree	34.4
– Random binary tree	28.5
– Left-branching tree	30.8
– Right-branching tree	36.0

- **Compared to gold standard.**
- Scores calculated over spans of the parsed trees.

# Learned library captures the hierarchical organization of simplified Chinese characters



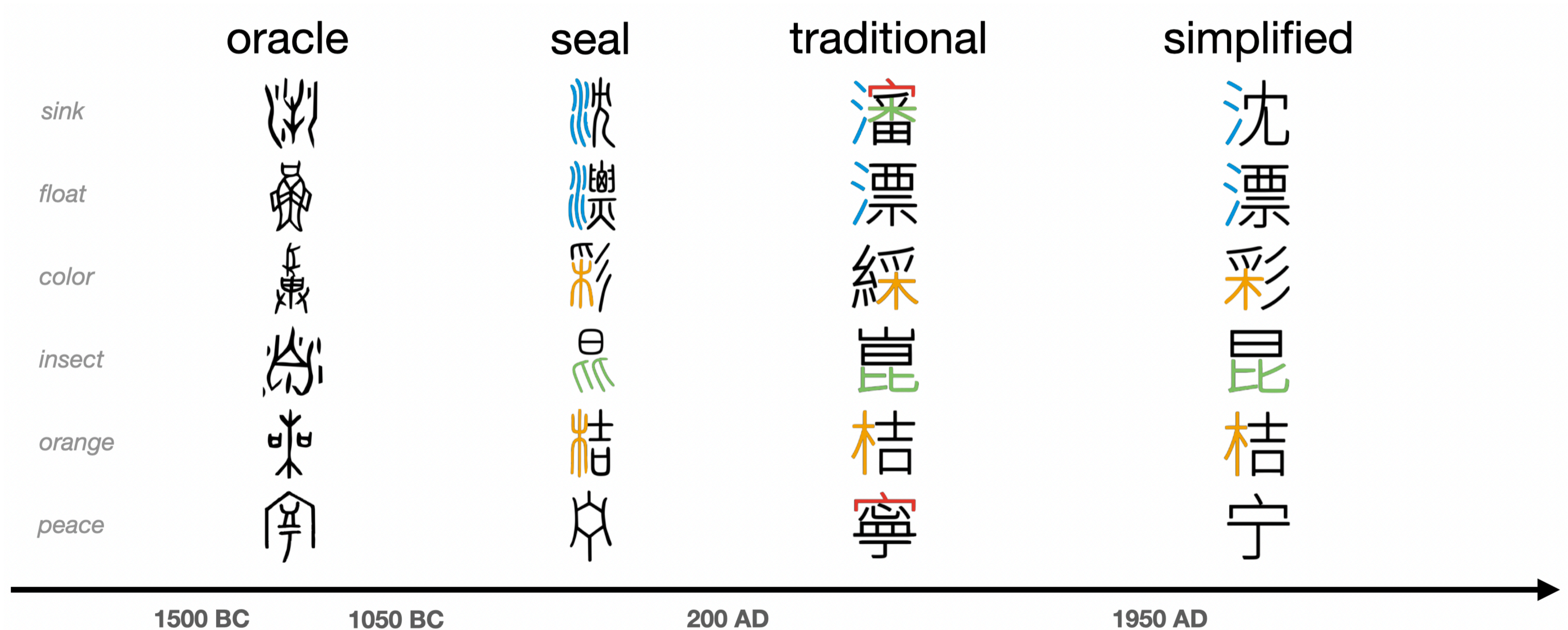
Model	F <sub>1</sub>
Library learning	61.6
Baselines	
– Balanced binary tree	34.4
– Random binary tree	28.5
– Left-branching tree	30.8
– Right-branching tree	36.0

- Compared to gold standard.
- Scores calculated over spans of the parsed trees.



# Diachronic analysis:

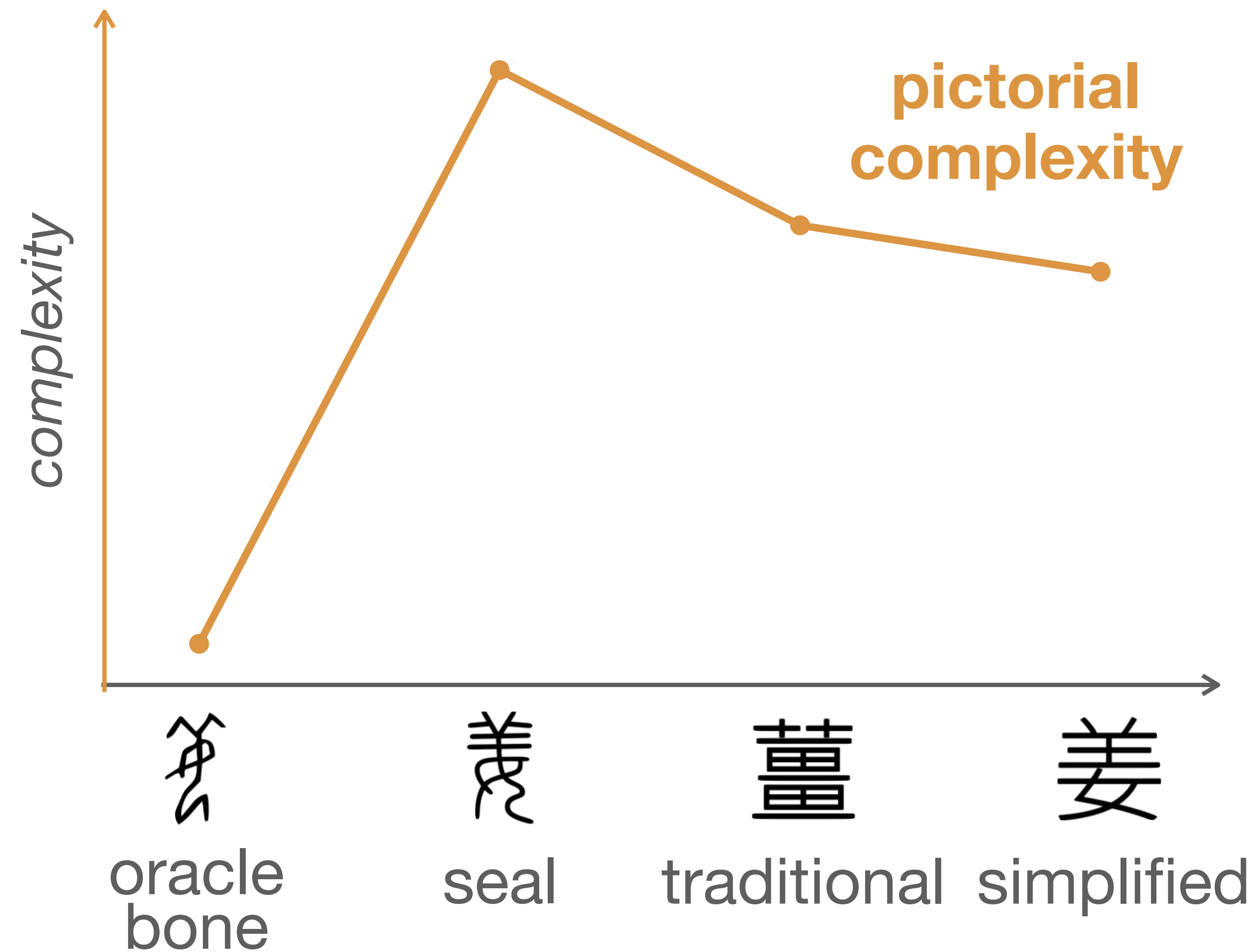
## How has the Chinese writing system evolved?



- A simplification? More efficient? Visually more complex?

# Diachronic analysis:

How has the Chinese writing system evolved?

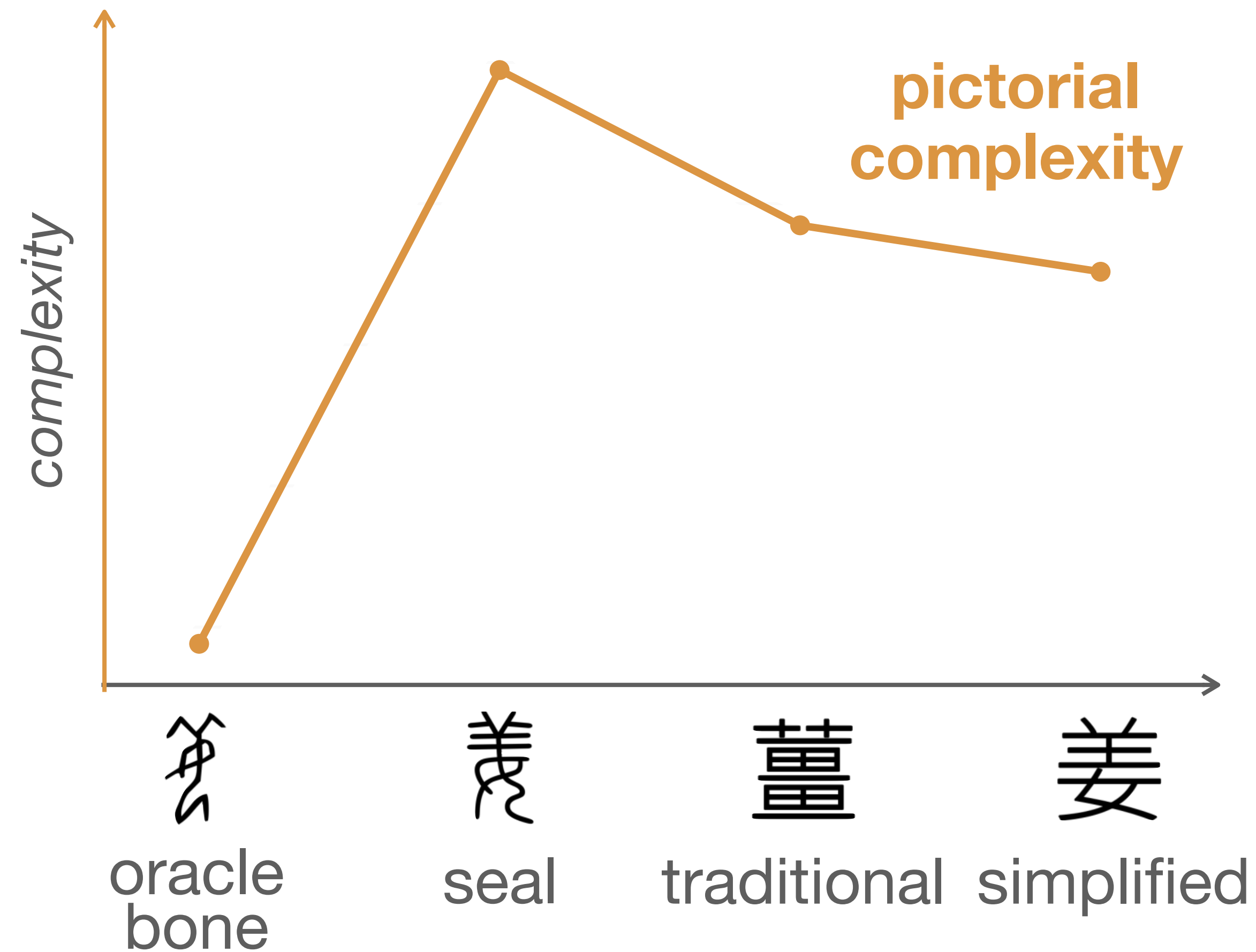


*pictorial  
(perimetric)  
complexity:*  $C = \frac{p^2}{4\pi A}$

Widely used for drawings and simple shapes!

# Diachronic analysis:

How has the Chinese writing system evolved?



*pictorial  
(perimetric)  
complexity:*  $C = \frac{p^2}{4\pi A}$

Widely used for drawings and simple shapes!

Previous results based on **pictorial complexity** did not show a gradual simplification over time (Han et al., 2022).

# Is pictorial complexity enough?

- It is good for **simple drawings**.
- But is **not** capable of **capturing complexity and reuse** at a **system** level.

# Is pictorial complexity enough?

- It is good for **simple drawings**.
- But is **not** capable of **capturing complexity and reuse** at a **system** level.

駟鞠

# Is pictorial complexity enough?

- It is good for **simple drawings**.
- But is **not** capable of **capturing complexity and reuse** at a **system** level.

馬鞴

真直

# Is pictorial complexity enough?

- It is good for **simple drawings**.
- But is **not** capable of **capturing complexity and reuse** at a **system** level.

駟

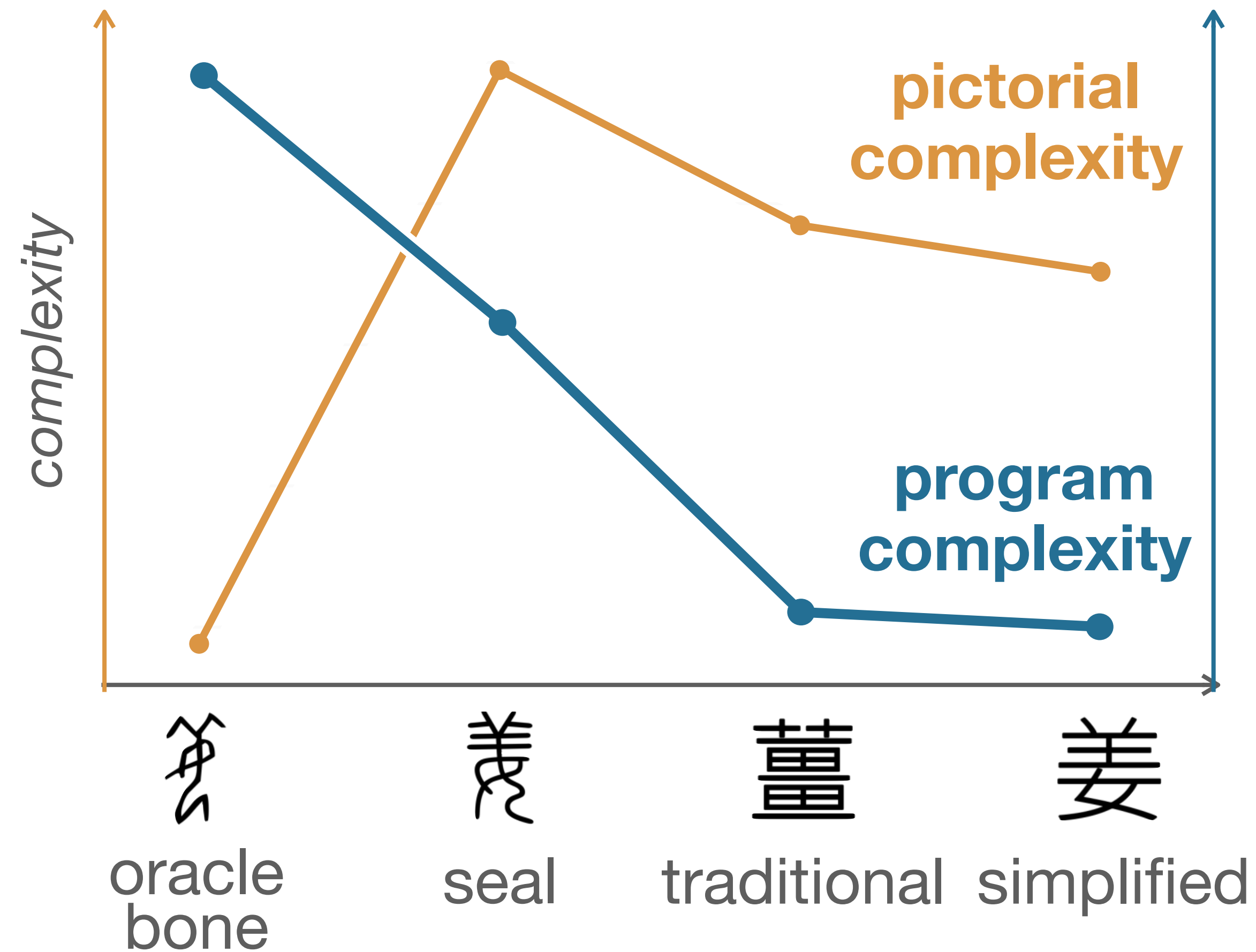
轟

These two characters have comparable pictorial complexities.  
However, considering reuse and patterns, 轟 is much simpler.

as 轟 = 直 × 3

# Diachronic analysis:

## How has the Chinese writing system evolved?

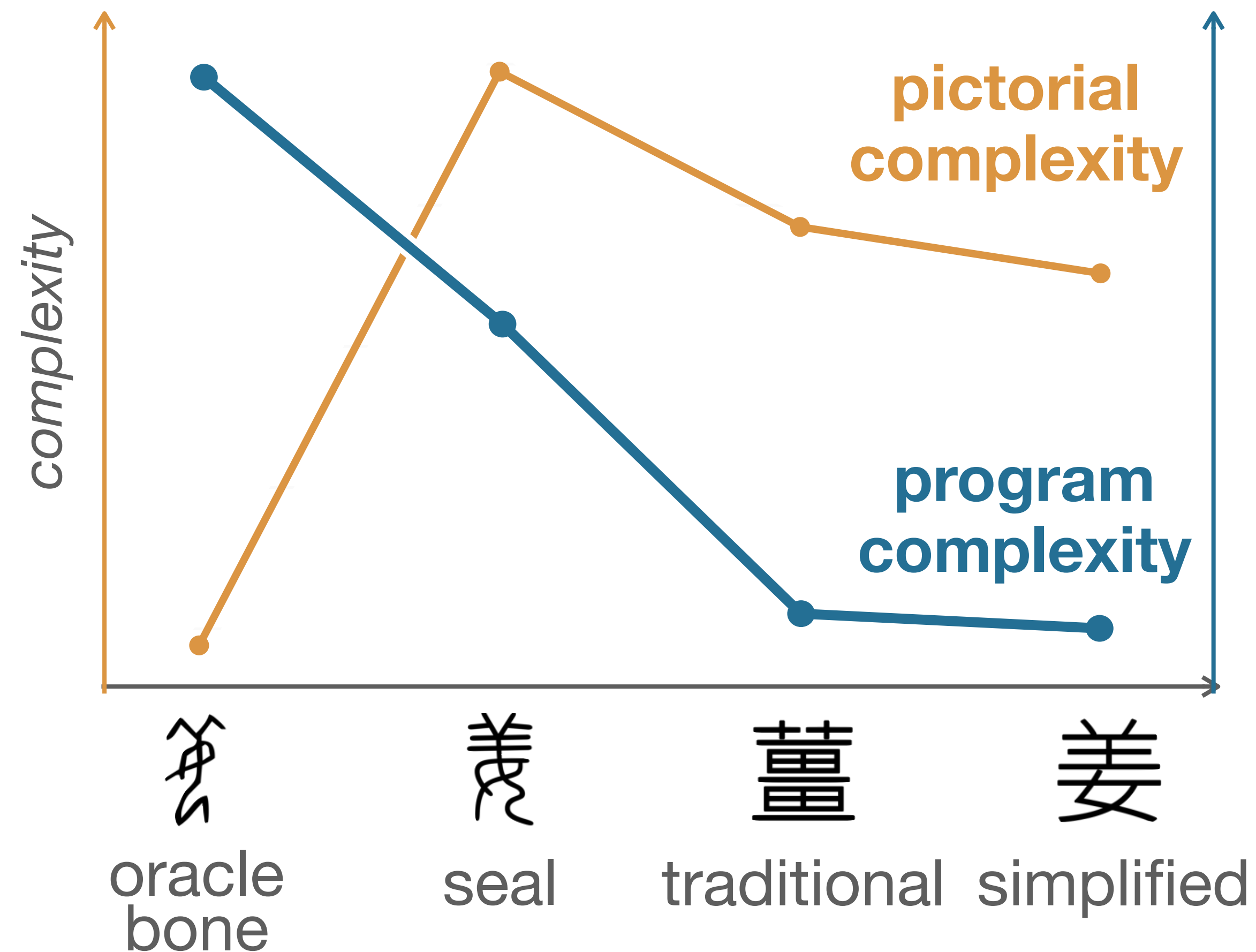


Our prediction:  
the **library learning model**  
should **reveal a gradual  
simplification** as systems  
**adapt to these biases** in  
cultural evolution.



# Diachronic analysis:

## How has the Chinese writing system evolved?



Our prediction:

the **library learning model** should **reveal a gradual simplification** as systems **adapt to these biases** in cultural evolution.

Result:

**Program complexity**  $C(\mathcal{W})$  has shown a **monotonic decrease** across time, confirming earlier empirical arguments.

# More on recent changes ~1950s

## Traditional Chinese -> Simplified Chinese

oracle

seal

traditional

simplified

*sink*

*float*

*color*

*insect*

*orange*

*peace*

1500 BC

1050 BC

200 AD

1950 AD

# More on recent changes ~1950s

Traditional Chinese -> Simplified Chinese

traditional

simplified

*sink*

沈

沈

*float*

漂

漂

*color*

綵

彩

*insect*

昆

昆

*orange*

桔

桔

*peace*

寧

宁

1500 BC

1050 BC

200 AD

1950 AD

# More on recent changes ~1950s

## Traditional Chinese -> Simplified Chinese

observation  
(traditional ⇒ simplified)

開 ⇒ 开

鑄 ⇒ 铸

盧 ⇒ 卢

爐 ⇒ 炉

將 ⇒ 将

獎 ⇒ 奖

(non-systematic  
rules observed)

inconsistent  
one-to-multiple  
mapping

inconsistent  
simplification

- **A real simplification?**

- This process may have **disrupted established systematicity** and lead to a loss of established semantic-phonetic and graphic patterns (Handel, 2013; Zhao & Baldauf, 2011).

# More on recent changes ~1950s

## Traditional Chinese -> Simplified Chinese

observation  
(traditional  $\Rightarrow$  simplified)

開  $\Rightarrow$  开

開  $\Rightarrow$  开

盧  $\Rightarrow$  卢

爐  $\Rightarrow$  炉

將  $\Rightarrow$  将

將  $\Rightarrow$  将

(non-systematic  
rules observed)

simplification rules

1 開  $\Rightarrow$  开 / 开

2 盧  $\Rightarrow$  卢 / 户

3 將  $\Rightarrow$  将 / 将

inconsistent  
one-to-multiple  
mapping

inconsistent  
simplification

# More on recent changes ~1950s

## Traditional Chinese -> Simplified Chinese

observation  
(traditional  $\Rightarrow$  simplified)

開  $\Rightarrow$  开

鑄  $\Rightarrow$  铸

盧  $\Rightarrow$  卢

爐  $\Rightarrow$  炉

將  $\Rightarrow$  将

獎  $\Rightarrow$  奖

(non-systematic  
rules observed)

inconsistent  
one-to-multiple  
mapping

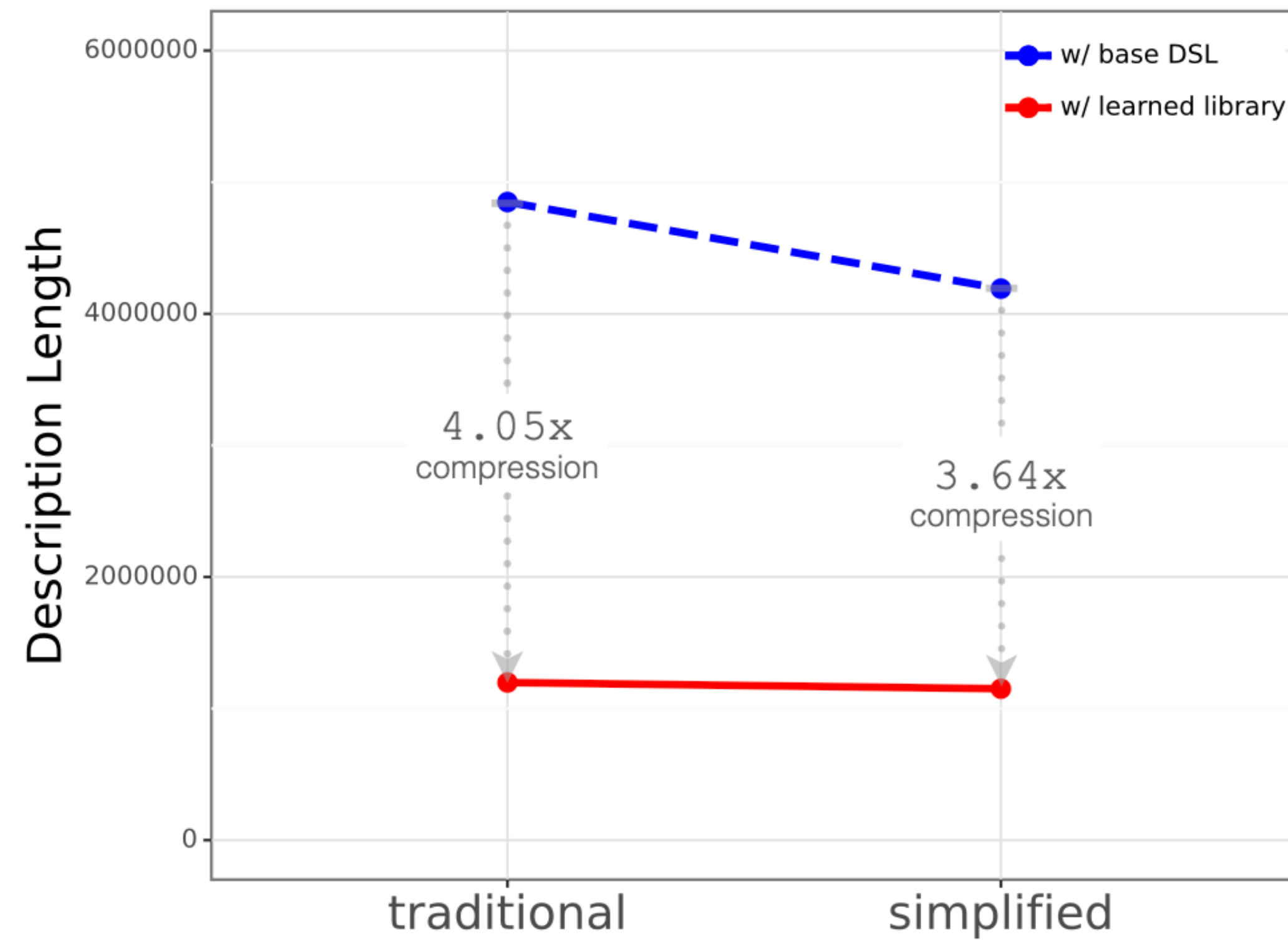
inconsistent  
simplification

- Can our computational model provide concrete evidence of the loss of systematicity?

- Our prediction:

- **Systematic scripts should be more compressive.**

# Simplified Chinese is simpler but less systematic compared to traditional Chinese



- Compression ratio (Raw DL / Compressed DL):
  - **Traditional > Simplified**
  - Suggesting the simplification process **did break** part of the systematicity.

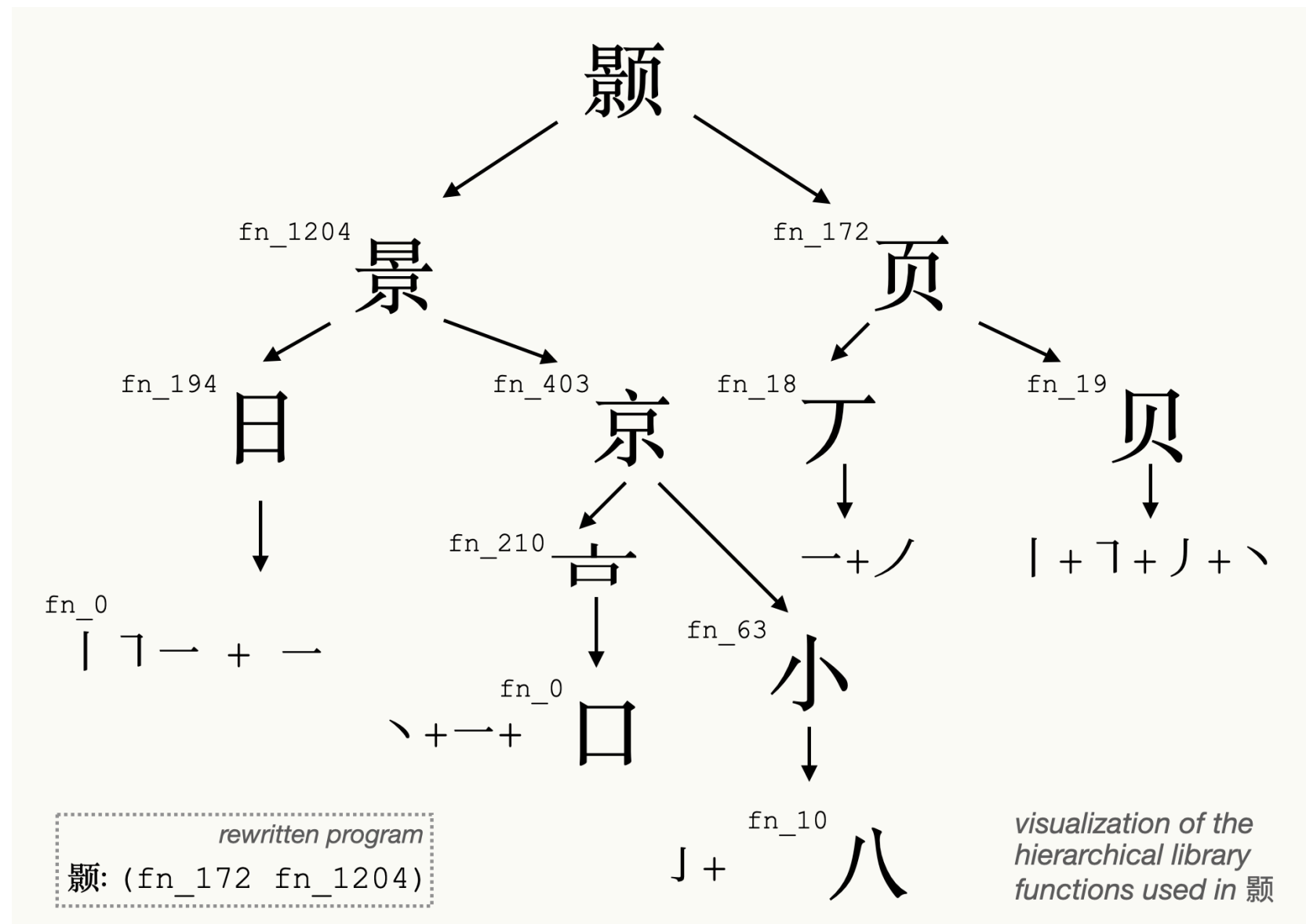
# Conclusions & Takeaway

- **A library learning-based** computational model can reveal the **inductive biases** behind the **emergence and evolution of combinatorial structures** in human language.
- **Combinatoriality**
  - **Develops** from a **MDL** perspective of representational efficiency
  - **By discovering inventories of reusable parts**
  - **By compressing the language**



# Conclusions & Takeaway

← Synchronic

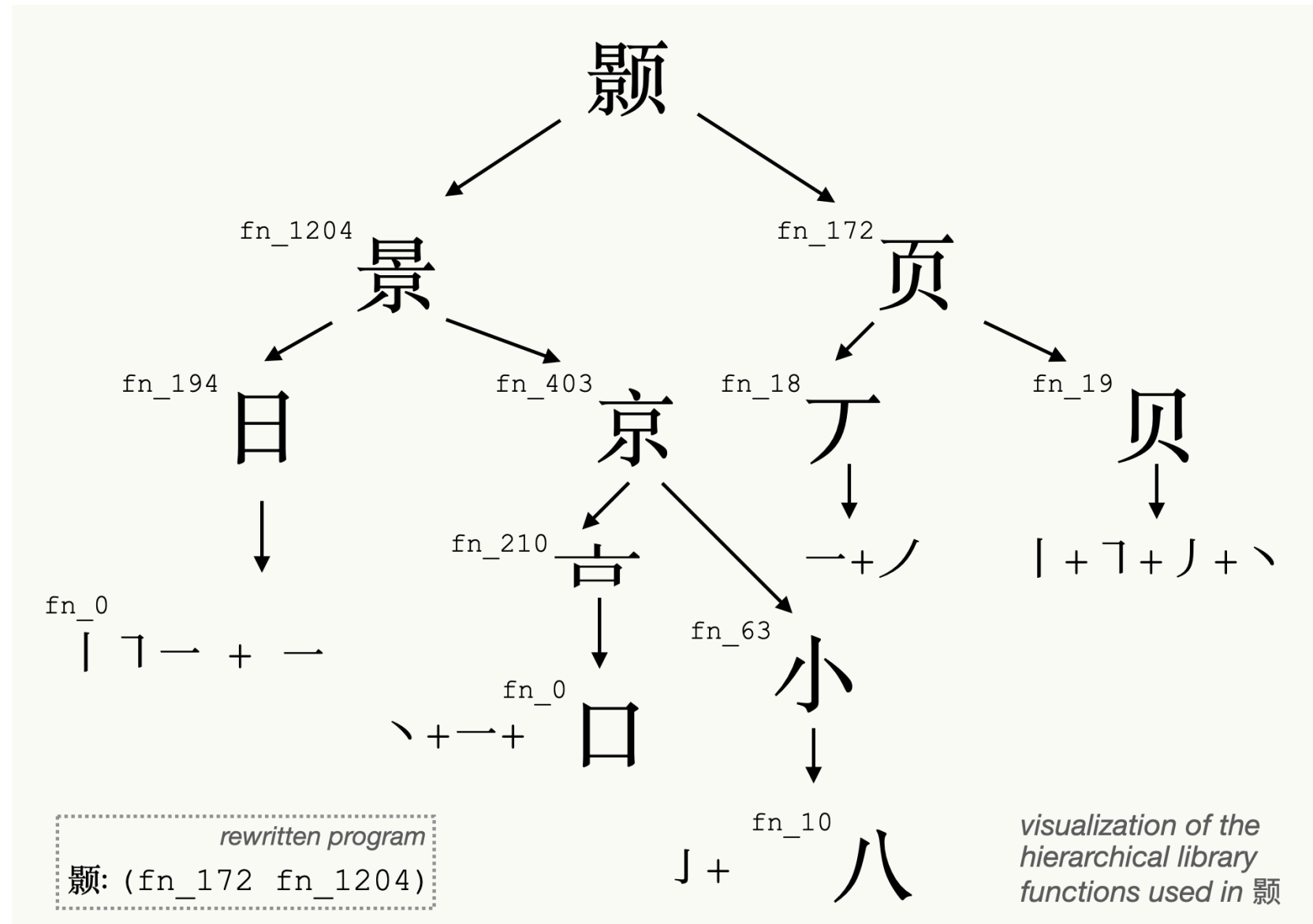


## Discovered and aligned radicals (187 / 201)

H	一	丨	丿	丶	冫	十	厂	匚	卜	冂	八	人	勹	儿	匕	几
82	259	31	36	163	999	41	25	33	190	17	6	64	37	106	304	
52	126	721	63	0	497	124	100	262	164	224	107	414	104	95	80	
49	71	77	160	173	67	169	40	219	112	295	32	378	666	27	507	
208	475	78	953	433	399	930	264	144	43	19	149	205	225	453	235	
320	462	396	139	1582	582	22	782	154	679	155	510	189	125	561	371	
56	1423	167	178	58	309	1099	34	69	410	159	715	695	275	116	162	
42	445	242	294	177	1268	1579	372	389	88	1654	276	271	172	216	81	
113	519	581	51	330	283	995	535	1069	196	328	110	457	117	28	247	
1004	407	789	175	1018	86	222	675	1307	704	1354	290	590	701	293	265	
489	1192	258	228	597	1061	35	463	46	303	751	959	137	1663	783	906	
614	1254	129	943	468	586	287	543	1705	1252	877						
	音	首	影	高	黄	麻	鹿	黑	黍	鼓	鼻					

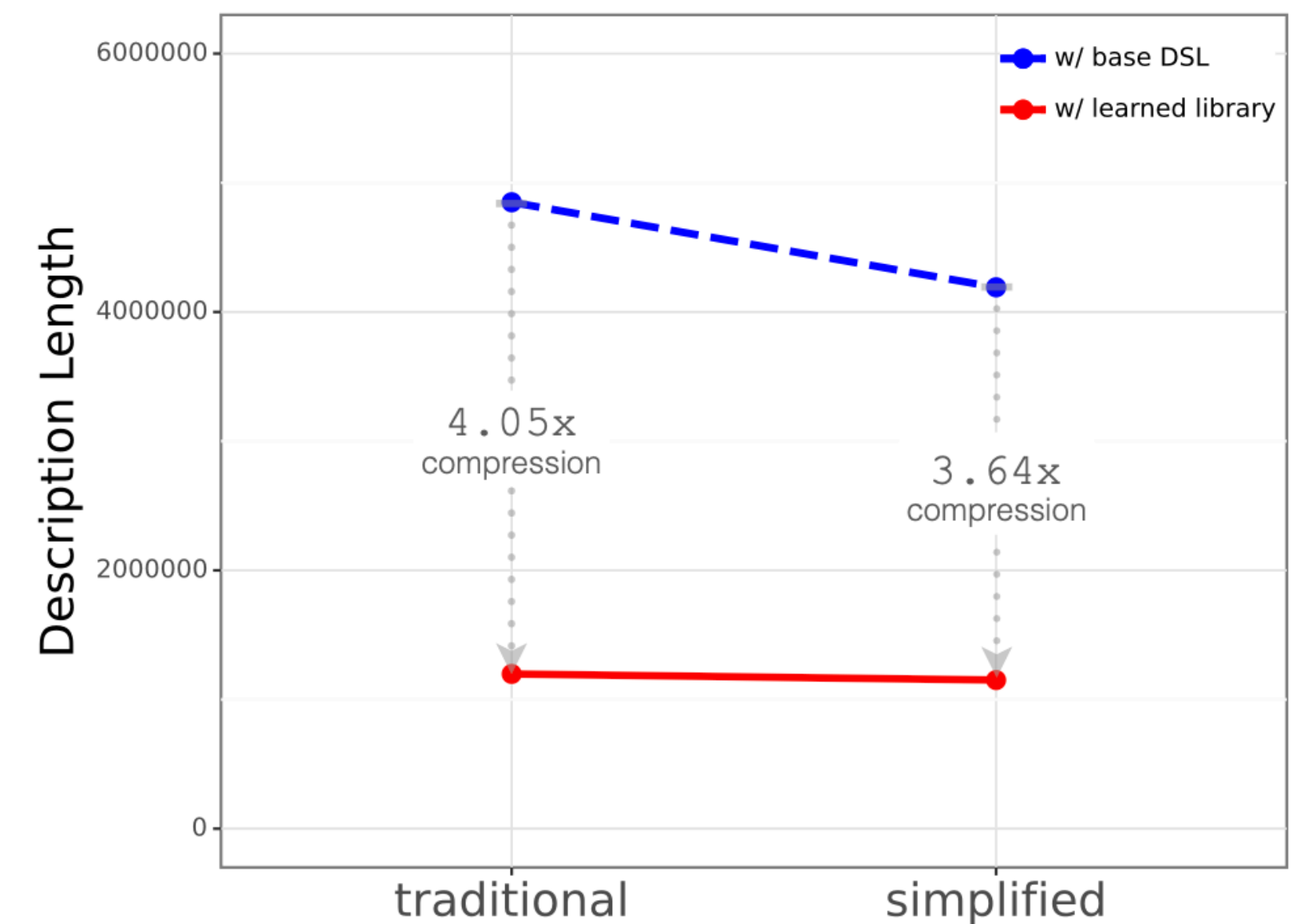
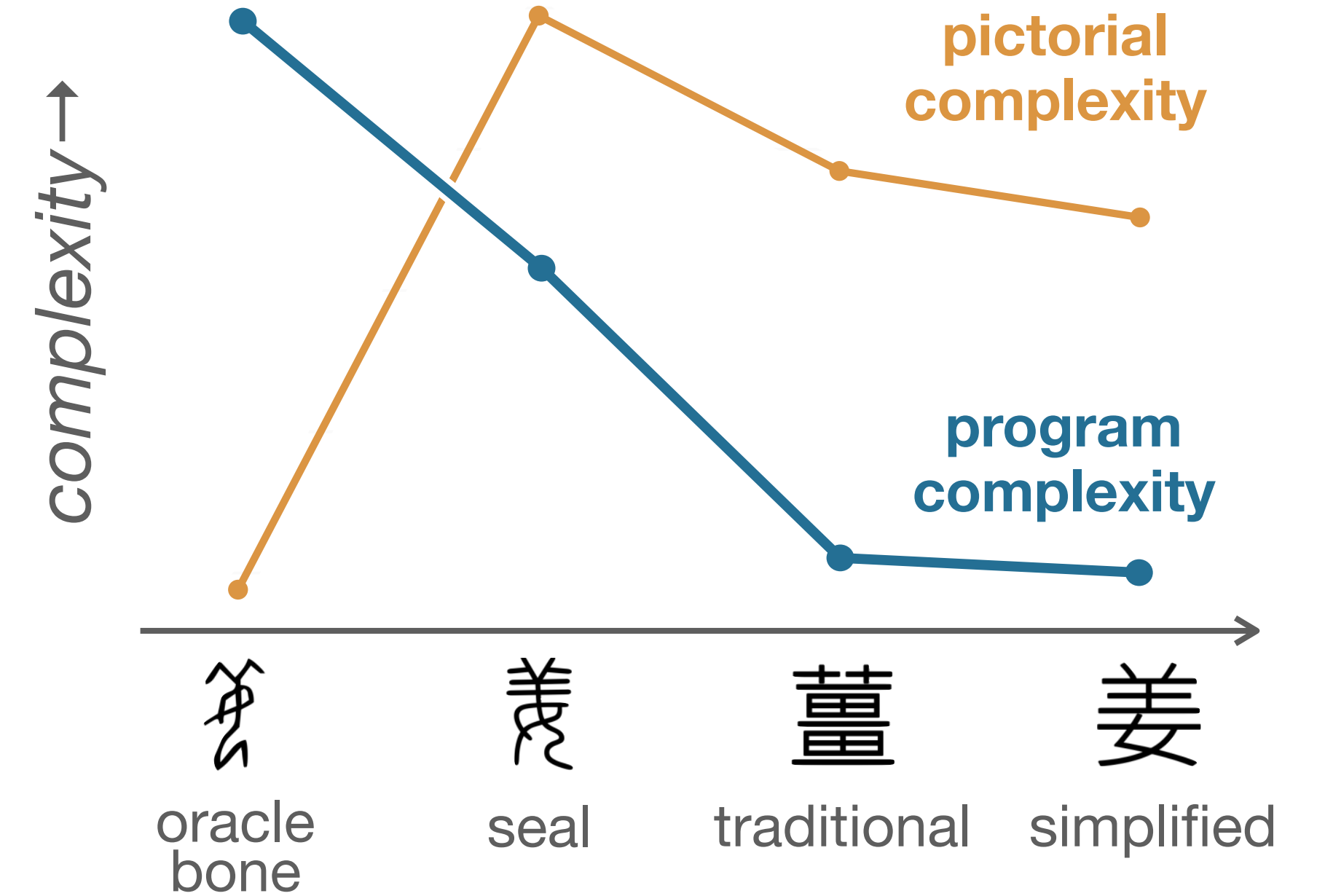
# Conclusions & Takeaway

← Synchronic



## Discovered and aligned radicals (187 / 201)

H	S	SP	D	HG	1	8	161	24	29	10	23	120	66	57	89
一	丨	ノ	丶	冫	十	厂	匚	卜	冂	八	人	勹	儿	匕	几
82	259	31	36	163	999	41	25	33	190	17	6	64	37	106	304
宀	冫	冂	冂	冂	刀	力	又	彡	彡	干	工	土	艹	寸	井
52	126	721	63	0	497	124	100	262	164	224	107	414	104	95	80
大	尢	弋	小	口	口	山	巾	彡	彡	夕	夕	斗	广	门	宀
49	71	77	160	173	67	169	40	219	112	295	32	378	666	27	507
辶	冫	尸	己	弓	子	巾	女	马	幺	《	王	无	韦	木	支
208	475	78	953	433	399	930	264	144	43	19	149	205	225	453	235
犬	歹	车	牙	戈	比	瓦	止	支	日	贝	水	见	牛	手	气
320	462	396	139	1582	582	22	782	154	679	155	510	189	125	561	371
毛	长	片	斤	爪	父	月	氏	欠	风	彡	文	方	火	斗	户
56	1423	167	178	58	309	1099	34	69	410	159	715	695	275	116	162
心	毋	示	甘	石	龙	业	目	田	𠂔	皿	生	矢	禾	白	鸟
42	445	242	294	177	1268	1579	372	389	88	1654	276	271	172	216	81
疒	立	穴	足	皮	彡	矛	耒	老	耳	臣	西	而	页	至	声
113	519	581	51	330	283	995	535	1069	196	328	110	457	117	28	247
虫	缶	舌	竹	白	自	血	舟	色	衣	羊	米	聿	艮	羽	糸
1004	407	789	175	1018	86	222	675	1307	704	1354	290	590	701	293	265
麦	走	豆	酉	辰	豕	里	足	邑	身	采	谷	豸	角	言	辛
489	1192	258	228	597	1061	35	463	46	303	751	959	137	1663	783	906
青	卓	雨	非	齿	鬲	佳	金	鱼	革	面	韭	骨	香	鬼	食
614	1254	129	943	468	586	287	543	1705	1252	877					
音	首	影	高	黄	麻	鹿	黑	黍	鼓	鼻					



Diachronic →

# Future work

- Extend to meaning compositionality:
  - Logographics captures the **multi-level** structure: not only combinatoriality in **forms**, but also compositionality in **meanings**.
- A wider range of logographic languages:
  - Cuneiform, Vai script...
- Consider more factors: frequency, motor cost, iconicity, visual complexity, etc.

# Thanks!

I'm actively seeking PhD positions starting  
25fall :)



Guangyuan



Paper