# Simulating the Early Evolution of Writing

Richard Sproat Google, Inc

Signs of Writing Beijing June 25-27

# The problem

- There have been only four (more or less unequivocal) cases of the independent discovery of writing:
  - o Mesopotamia
  - Egypt
  - China
  - Mesoamerica
- With only four data points, it is hard to make any generalizations about the conditions that might favor (or disfavor) the discovery.

# **Computational simulation**

- Computational simulation has been used in a number of areas of linguistics to model phenomena that are hard to test in the laboratory or in the field:
  - Spread of linguistic features in social networks (e.g. Steels, 2012)
  - Historical change (e.g. Niyogi, 2006)
  - Emergence of linguistic properties (e.g. Kirby, 1999 and much subsequent work)
- The present research seeks to model the emergence of writing from non-linguistic symbol systems

# **Phenomena of interest**

- The non-linguistic symbol systems in use in the culture, and the existence of combinatorial systems where symbols occur in "texts".
- Linguistic properties favoring the discovery of writing, and favoring a particular kind of writing system over another (e.g. a consonantary versus a syllabary).
- Economic or other factors that would encourage the development of better means of record keeping.
- The development of lightweight materials encouraging the wider use of writing (Farmer et al. 2002).

# **Phenomena of interest**

- The non-linguistic symbol systems in use in the culture, and the existence of combinatorial systems where symbols occur in "texts".
- Linguistic properties favoring the discovery of writing, and favoring a particular kind of writing system over another (e.g. a consonantary versus a syllabary).
- Economic or other factors that would encourage the development of better means of record keeping.
- The development of lightweight materials encouraging the wider use of writing (Farmer et al. 2002).

### **Summary of Chicago Presentation**

# Goals

- Generate a writing system for a language from a nonlinguistic symbol system. Symbols are assigned to morphemes, spreading to new morphemes by:
  - Semantic similarity
  - Phonetic similarity
- Determine how many spellings for morphemes are purely semantic, purely phonetic, or mixed
  Compare these with the situation in true ancient writing systems
- Determine which linguistic properties make phonetic spread easiest.

# **Parameters: phonotactics**

### **Basic phonotactics for morphemes**

Monosyllabic language	σ = C? V R? C?	[R= sonorant]
Disyllabic language	σσ?	

Randomly generate ≈1K morphemes from these templates

# **Parameters: what's phonetically close?**

### • CLOSE:

- obstruents match on place of articulation
- vowels must match
- sonorants optional
- Ex: daNg matches tak
- CLOSE\_RHYME:
  - as above, but "rhyming" is sufficient
  - Ex: daNg matches bak

### • CLOSE\_V\_FREE:

- $\circ$  same as CLOSE, but vowels don't need to match, and V matches Ø
- Ex: donek matches dink

### • STRICT:

• Exact match required

# **Parameters: semantics**

### 100 basic "concepts":

PERSON, MAN, WOMAN, HOUSE, BRONZE, SWORD, MEAT, SHEEP, OX, GOAT, FISH, TREE, BARLEY, WHEAT, WATER, STONE, CLAY, THREAD, CLOTHING, FIELD, TEMPLE, GOD, AXE, SCYTHE, DOG, LION, WOLF, DEMON, SNAKE, TURTLE, FRUIT, HILL, CAVE, TOWN, ENCLOSURE, FLOWER, RAIN, THUNDER, CLOUD, SUN, MOON, HEART, LUNG, LEG, ARM, FINGER, HEAD, TONGUE, EYE, EAR, NOSE, GUTS, PENIS, VAGINA, HAIR, SKIN, SHELL, BONE, BLOOD, LIVER, FARM, LOCUST, STICK, STAR, EARTH, ASS, DEATH, BIRTH, WOMB, MILK, COAL, SEED, LEAF, CHILD, ANTELOPE, BEAR, BEE, MOUSE, DUNG, PLOUGH, SPROUT, ICE, DAY, NIGHT, WINTER, SUMMER, AUTUMN, SPRING, KING, GOOSE, PRIEST, ROAD, CART, GRASS, FIRE, WIND, NAIL, BREAST, BOWL, CUP

Randomly associate these with "morphemes" and for remaining morphemes associate those with **random combinations of up to 3 of these** 

# **Parameters: symbols**

Each basic concept is associated with a symbol. We used dingbats and other mostly non-script symbols from the Unicode Basic Multilingual Plane.\* Some examples:



\*Expanded in the current system to include other sets such as Yi symbols

# Parameters: "languages"

- Initial parameters:
  - MONOSYLLABIC / DISYLLABIC
  - CLOSE / CLOSE\_RHYME / CLOSE\_V\_FREE / STRICT
- For each pair of the above
  - Generate a "language" 5 times: a language consists of about 1,000 morphemes generated from the grammar
  - Run 5 simulations for each of these

# **Steps in simulation**

- 1. Assign simplex concepts and symbols to morphemes
  - a. From a set of morphemes associated with a symbol, pick *one* morpheme to associate to the symbol
- 2. Assign complex concepts to remaining morphemes
- 3. Assign spellings to morphemes that do not have one
  - a. Symbols associated with shared semantic components, and/or
  - b. Symbols associated with similar sound ⇐ "eureka" moment
- 4. Extend the phonetic coverage with *telescoping*:
  - a.  $ba + ad \rightarrow bad$
- 5. Result is a mix of semantic and phonetic components

# **MONOSYLLABIC, CLOSE**

toN	PERSON	S(emantic){&}
kerg	PERSON	S{發}+P(honetic){☆]
gib	PERSON	S{ॐ}+P{□}
kak	MAN	S{ d }
gab	WOMAN	S{"}
pet	WOMAN	S{□}+P{\_}
urb	HOUSE	S{ <b>▲</b> }
kuNt	HOUSE	S{ <b>▲</b> }+P{⊡}
а	BRONZE	S{\\$
dep	BRONZE	S{\$\$}+P{\$\$}
keg	BRONZE	S{☞}+P{☆}

## **Results from previous simulations**

	# entries	w/ spelling	w/out spelling	phon	$\operatorname{sem+phon}$	$\mathbf{purephon}$	$\operatorname{sem}$
DI, CLOSE	661	462 (0.70)	198(0.30)	47(0.07)	35(0.05)	11 (0.02)	451 (0.68)
DI, CLOSE-RHYME	648	508(0.79)	139(0.21)	127(0.20)	100(0.15)	26(0.04)	482(0.74)
DI, CLOSE-V-FREE	644	640(0.99)	3(0.01)	324(0.50)	262(0.41)	61(0.10)	579(0.90)
DI, STRICT-V-FREE	649	504(0.78)	144 (0.22)	112(0.17)	87(0.14)	24 (0.04)	481(0.74)
MONO, CLOSE	643	580(0.90)	62(0.10)	239 (0.37)	193(0.30)	45(0.07)	534(0.83)
MONO, CLOSE-RHYME	651	630(0.97)	20(0.03)	310(0.48)	253 (0.39)	56(0.09)	574(0.88)
MONO, STRICT	640	487(0.76)	152(0.24)	95(0.15)	71(0.11)	23(0.04)	464(0.73)

	Table 5 St	ructural Classifi	cation of Charac	ters
Principle	Oracle Bones (Shang dynasty)	Xu Shen (2nd century)	Zheng Qiao (12th century)	Kang Xi (18th centuty)
Pictographic	227 (23%)	364 (4%)	608 (3%)	
Simple indicative	20(2%)	- 125 (1%)	107(1%)	\$ ±1,500(3%)
Compound indicative	396 (41%)	1,167 (13%)	740(3%)	
Semantic- phonetic	334 (34%)	7,697 (82%)	21,810(93%)	47,141 (97%)
Total	977	9,353	23,265	48,641
	Principle Pictographic Simple indicative Compound indicative Semantic- phonetic Total	Principle Oracle Bones (Shang dynasty)   Pictographic 227 (23%)   Simple indicative 20 (2%)   Compound indicative 396 (41%)   Semantic- phonetic 334 (34%)   Total 977	Diable yStructural ClassifieOracle Bones (Shang dynasty)Xu Shen (2nd century)Principledynasty)century)Pictographic227 (23%)364 (4%)Simple indicative20 (2%)125 (1%)Compound indicative396 (41%)1,167 (13%)Semantic- phonetic334 (34%)7,697 (82%)Total9779,353	Diracie     Structural Classification of Charace       Oracle     Zheng       Bones     (2nd       (Shang     (2nd       (I2th       dynasty)     century)       Principle     227 (23%)       364 (4%)     608 (3%)       Simple     indicative       indicative     20 (2%)       125 (1%)     107 (1%)       Compound     1,167 (13%)       indicative     396 (41%)       Semantic-     334 (34%)       phonetic     334 (34%)       7,697 (82%)     21,810 (93%)       Total     977     9,353       23,265     23,265

#### Table 3 Structural Classification of Characters

# **Take away points**

- Strict matches are not very effective.
- Close consonant matches useful for largely "monosyllabic" languages; less so for "disyllabic" languages.
- Allowing vowels to be free is useful in "disyllabic" languages.
- The rate of semantic-phonetic compounds in "monosyllabic" languages with close consonant matches is broadly similar to the rate of semantic-phonetic compounds in Oracle Bone texts.
- But we don't really provide much evaluation of how well the system works

# Note: some details of the simulation

- Currently about 1100 lines of Python code
- Uses the *pyfst* interface to OpenFst
- Grammars written in *Thrax* finitestate grammar compiler (http://www.openfst.

org/twiki/bin/view/GRM/Thrax

### **Evaluating the Similarity to Known Systems**

# **Problems with the previous work**

- Letting a symbol inherit the phonology from just one morpheme is not a good model of how some ancient systems apparently worked.
- We need some way of quantifying how good a fit the models are to (the few) documented cases of *ex nihilo* evolution of writing.

# Similarity within a phonetic class

- For a set of ancient writing systems, we want to measure how phonologically similar the morphemes are within a phonetic class
- For example in Chinese, those morphemes that share the same phonetic radical

# **Ancient Chinese**

- Source: Baxter-Sagart list. (<u>http://en.wiktionary.org/wiki/Appendix:Baxter-Sagart\_Old\_Chinese\_reconstruction</u>)
- Intersect with characters with semantic-phonetic decomposition from <a href="http://www.zhongwen.com/">http://www.zhongwen.com/</a>.
- Results in 902 entries. Some examples:

丘	丘	qiu1	khjuw	k <sup>wh</sup> ə
蚯	臣	qiu1	khjuw	k <sup>wh</sup> ə
虛	丘	xu1	khjo	q <sup>h</sup> a
虛	臣	xu1	xjo	qʰa
툢	臣	yue4	ngæwk	ŋ <sup>ç</sup> rok

# **Edit distance (ED)**

	а	b	С	d
С	а		е	d



- Efficient to compute:  $\mathcal{O}(\text{Length}_1 \times \text{Length}_2)$
- Operations could be weighted by phonetic distance.

# **Ancient Chinese**

- 156 equivalence classes (out of 198 79%) with >1 member
  - o qwak gwak gwa?
  - $\circ$  g<sup>c</sup>aŋ g<sup>c</sup>raŋ
  - mu? tha? thsa? tsa?
  - $\circ \quad \text{lap } m^{\varsigma} \text{ron } r^{\varsigma} \text{on}$
- Compute mean NED for each equivalence class with >1 member; average over all equivalence classes. Smaller is better.

	Mean edit distance
Ancient Chinese	0.54
Middle Chinese	0.60
Modern Mandarin	0.57

# **Sumerian**

• Source: Electronic Text Corpus of Sumerian Literature sign list (<u>http://etcsl.orinst.ox.ac.uk/edition2/signlist.php</u>)

Sign names	Signs	ETCSL values
Α	Ϋ́	a, dur <sub>5</sub> , duru <sub>5</sub>
A.AN	ĭ₩ <del>X</del>	am <sub>3</sub> , em <sub>x</sub> , šeĝ <sub>3</sub>
A.EDIN.LAL	Ţ¥₩Ĵ∕Q¢<Ţ <sup>►</sup>	ummud
A.HA.TAR.DU	TATA	girim <sub>3</sub>
A.IGI	T¥4⊢	er <sub>2</sub> , še <sub>x</sub>
А.КА	ĭ₩ <sup>™</sup>	ugu <sub>2</sub>
A.KAL	Ĩ <b>ŧ</b> ⊑∔	illu

• Kept all lower-case (phonetic) values, without subscripts.

# **Sumerian**

- 212 equivalence classes (out of 617 34%) with >1 member
- Same procedure as with Chinese

	Mean edit distance
Ancient Chinese	0.54
Middle Chinese	0.60
Modern Mandarin	0.57
Sumerian	0.89

# **Middle Egyptian**

• 21 equivalence classes (out of 256 — 8%) with more than one member from Gardiner's sign list (via <a href="http://www.http://www.sign.com">http://www.sign.com</a>

//wwwegyptianhieroglyphsnet/gardiners-sign-list/low-narrow-signs/

	Mean edit distance
Ancient Chinese	0.54
Middle Chinese	0.60
Modern Mandarin	0.57
Middle Egyptian	0.60
Sumerian	0.89



Image from:

Dehaene, Stanislas. 2010. *Reading in the Brain.* New York, Penguin. Fig 2.2. p. 63

### Language and speech

**Visual processing** 





Symbol  $\sigma$  associated with concept and thence a set of related meanings/morphemes  $\mu_k$ 

Via the morphemes becomes associated to set of pronunciations  $\phi_k$ 





Symbol  $\sigma$  associated with concept and thence a set of related meanings/morphemes  $\mu_k$ 

One morpheme  $\mu_1$ becomes most strongly associated with the symbol.

Via this morpheme becomes associated to pronunciation  $\phi_1$ 



Symbol  $\sigma$  associated with concept and thence a set of related meanings/morphemes  $\mu_k$ 

One morpheme  $\mu_1$ becomes most strongly associated with the symbol.

Via this morpheme becomes associated to pronunciation  $\phi_1$ 



Symbol  $\sigma$  associated with concept and thence a set of related meanings/morphemes  $\mu_k$ 

One morpheme  $\mu_1$ becomes most strongly associated with the symbol.

Via this morpheme becomes associated to pronunciation  $\phi_1$ 

- The first effectively "fossilizes" the non-linguistic origin of the sign, preserving it through to multiple phonetic functions. **This is like Sumerian.**
- The second treats the sign as linguistic earlier by associating it to a particular morpheme and thence to a particular sound. This is like Chinese or Egyptian.
  - The second seems to reflect a more advanced stage: the inventors of the system realize that a sign can stand for a particular abstract linguistic unit.

# **Simulation**

Parameterize with two options:

- As in previous simulations, pick just one morpheme per semantic group and base phonetics on that morpheme (Chinese, Egyptian)
- Use all morphemes in a semantic group (Sumerian)

# **Coverage results**

	# entries	w/ spelling	w/out spelling	$\mathbf{phon}$	sem+phon	$\mathbf{purephon}$	$\operatorname{sem}$
DI, CLOSE	661	462(0.70)	198(0.30)	47 (0.07)	35 (0.05)	11(0.02)	451 (0.68)
DI, CLOSE-RHYME	648	508(0.79)	139(0.21)	127(0.20)	100(0.15)	26(0.04)	482(0.74)
DI, CLOSE-V-FREE	644	640(0.99)	3(0.01)	324(0.50)	262(0.41)	61(0.10)	579(0.90)
DI. STRICT-V-FREE	649	504(0.78)	144(0.22)	112(0.17)	87 (0.14)	24(0.04)	481 (0.74)
MONO, CLOSE	643	580 (0.90)	62(0.10)	239(0.37)	193 (0.30)	45(0.07)	534(0.83)
MONO, CLOSE-RHYME	651	630(0.97)	20(0.03)	310(0.48)	253(0.39)	56(0.09)	574(0.88)
MONO, STRICT	640	487(0.76)	152(0.24)	95(0.15)	71(0.11)	23(0.04)	464(0.73)

	# entries	w/ spell	w/out sp	phon	sem+phon	purephon	sem
MONO CLOSE 1 morph	638	587 (0.92)	50 (0.08)**	235 (0.40)**	203 (0.32)*	50 (0.08)**	537 (0.84)**
MONO CLOSE all morphs	645	636 (0.99)	8 (0.01)**	222 (0.34)**	187 (0.29)*	34 (0.05)**	601 (0.93)**

# **Comparison with extant systems**

	Mean edit distance	Standard deviation
Ancient Chinese	0.54	
Middle Chinese	0.60	
Modern Mandarin	0.57	
Middle Egyptian	0.60	
MONOSYLLABIC_CLOSE: 1 morph*	0.52	0.026
MONOSYLLABIC_CLOSE: all morphs*	0.78	0.026
Sumerian	0.89	

\*In both cases the percentage of signs with multiple pronunciations is about 64% — less than Chinese but more than Sumerian or Egyptian

# **Simulation using OC syllables**

404 Old Chinese reconstructed syllables:
ba, baŋ, ben, but, b<sup>c</sup>a, b<sup>c</sup>awk, b<sup>c</sup>aŋ,

b<sup>ç</sup>rak...

- 198 equivalence classes used to define "CLOSE"
- One morpheme per group

- IFa? ENCLOSURE,SKIN,EYE S{♥, ◎}+P{\*
- lek NIGHT,HOUSE,LEAF S{♣, ℂ, △}
- pa? NIGHT,HOUSE,LEAF S{♣, ℂ }
- d<sup>c</sup>ak CUP,LEAF S{♣}+P{�}
- lek LEAF,BIRTH S{♣,♂}
- Irə LEAF S{\*\*}



# **Example:** >

- r<sup>c</sup>ew WOMB,ARM,MOON S{ **〕**, **♥**}
- nij COAL S{∎}+P{ )}
- nij CLOTHING,SUMMER,THUNDER  $S\{5\}+P\{\)$
- nij SILVER,FISH S{§}+P{ )}

# **Coverage results**

	# entries	w/ spelling	w/out spelling	$\mathbf{phon}$	$\operatorname{sem+phon}$	purephon	$\operatorname{sem}$
DI, CLOSE	661	462(0.70)	198(0.30)	47(0.07)	35 (0.05)	11 (0.02)	451 (0.68)
DI, CLOSE-RHYME	648	508(0.79)	139(0.21)	127(0.20)	100(0.15)	26(0.04)	482(0.74)
DI, CLOSE-V-FREE	<b>644</b>	640(0.99)	3(0.01)	324(0.50)	262(0.41)	61(0.10)	579(0.90)
DI. STRICT-V-FREE	649	504(0.78)	144(0.22)	112(0.17)	87 (0.14)	24(0.04)	481(0.74)
MONO, CLOSE	643	580 (0.90)	62 (0.10)	239 (0.37)	193 (0.30)	45 (0.07)	534 (0.83)
MONO, CLOSE-RHYME	651	630(0.97)	20(0.03)	310(0.48)	253(0.39)	56(0.09)	574(0.88)
MONO, STRICT	640	487(0.76)	152(0.24)	95(0.15)	71(0.11)	23(0.04)	464(0.73)

	# entries	w/ spell	w/out sp	phon	sem+phon	purephon	sem
MONO CLOSE 1 morph	638	587 (0.92)	50 (0.08)	235 (0.4)	203 (0.32)	50 (0.08)	537 (0.84)
MONO CLOSE all morphs	645	636 (0.99)	8 (0.01)	222 (0.34)	187 (0.29)	34 (0.05)	601 (0.93)
"Ancient Chinese" 1 morph	650	533 (0.82)	116 (0.25)	163 (0.25)	130 (0.20)	32 (0.05)	500 (0.77)

# **Comparison with simulations**

	Mean edit distance	Standard deviation
Ancient Chinese	0.54	
Middle Chinese	0.60	
Modern Mandarin	0.57	
Middle Egyptian	0.60	
MONOSYLLABIC_CLOSE: 1 morph	0.52	0.03
"Ancient Chinese"-based: 1 morph	0.67	0.08
MONOSYLLABIC_CLOSE: all morphs	0.78	0.03
Sumerian	0.89	

# **Synopsis and conclusions**

- A model where pronunciations spread from all morphemes associated with a concept fits Sumerian better.
- A model where one morpheme is picked as *the* denotation of the symbol fits Chinese (or Egyptian) better.
- Modeling based on "actual" Ancient Chinese syllables and limiting ourselves to reconstructed phonetic equivalence classes does not work well — but that is probably because the phonetic equivalence classes are too limiting.
- Proportion of signs with multiple pronunciations also differs significantly from known systems. This could partly reflect standardization in real systems, which we are not modeling.

# **Synopsis and conclusions**

- Does this reflect a difference in the evolution of the scripts?
- Sumerian developed from a raw "ideographic" system.
- But Chinese had symbols already associated with specific morphemes — a later phase in the evolution of writing? Either:
  - Phonetics were standardized from an earlier system...
  - or maybe Chinese got the idea of writing from elsewhere...

# **Further work**

- Link consonantal systems to ablaut-like processes
  - We already have some results but no time to report here.
- Simulate writing full sentences and affixal morphology
  - We already have some results but no time to report here.
- Simulate a wider range of phonetic shapes for morphemes, and a wider range of phonetic closeness
- Provide a more plausible model of lexical statistics
- Model of standardization processes
- Set of symbols is currently static: but new symbols are invented in real writing systems