

Bits 6 Dcripts 00000101



# Tutorial Encountering Writing

### Sitabhra Sinha

# How can you tell...

### This is writing...



#### ....but this is not !



# How can you tell...

### This is writing...



#### an Achameneid cuneiform tablet

#### ....but this is not !



#### a Nan e-Barbari (Persian flat bread)

# Writing and Language are not equivalent!

What is writing? a system of communicating – typically using a language - by means of conventional visible marks

What is language ? A system of syntactic communication capable of encoding ideas of arbitrary complexity

Syntax = Compositionality + Recursion ("words/letters") (embedding)

A writing system can be used to write different languages e.g., Roman alphabet

Look for syntactic structure in a symbolic string

The same writing system can be used to write several different languages...

## Examples

### Latin alphabet

used to write many modern European languages, including English and many modern Asian languages, including Malay, Turkish and Indonesian

#### ✤ Arabic script

used to write texts in Arabic, Persian (Farsi and Dari), Malay (Jawi), Cham (Akhar Srak), Uyghur, Kurdish, Punjabi (Shahmukhi), Sindhi, Balti, Balochi, Pashto, Lurish, Urdu, Kashmiri, Rohingya, Somali, among others

## A B C D E F G H I K L M N O P Q R S T V X Y Z



## ... and the same language can be written using different writing systems Example

Turkic Languages have been written using

#### Orkhon script

(8<sup>th</sup>-10<sup>th</sup> century CE) written from right to left

### Ottoman Turkish alphabet

(10<sup>th</sup>-20<sup>th</sup> century) A form of the Perso-Arabic script written from right to left

### Latin Turkish alphabet

(1928 onwards) written from left to right d#↓L\$----: \$1T?1: \$6¼{%1: \$1: F4I1: d+: 6{M1: 6{: ?TX1: \$F1: k{F: 6{M1: %h: 4#%: k\*(: {LF\$↓\$; 447: k47: 6{M1: 6{%CF: ----4---: kL\${\$ ?3X1: k#¥@\$: 1%----: ----41TY1: ?----:

٢	Ż	۲	6	د	Ŷ	Ÿ	ш	Ņ	۴
d	ĥ	ķ	Ç	С	ş	t	р	b	)
ä	b	Ġ	0	Ŵ	w	5	ÿ	J	Ż
Ż	ţ	ż	Ş	ş	s	ğ	Z	r	₫
Ú	٢	J	٩	ڪ	5	قر	ė	Ė	٤
n	m	Ι	ŋ	g	k	ķ	f	ġ	¢

Aa Bb Cc Çç Dd Ee Ff Gg Ğğ Hh Iı İi Jj Kk LI Mm Nn Oo Öö Pp Rr Ss Şş Tt Uu Üü Vv Yy Zz

# The spectrum of writing

All known writing systems involve a mix of

- Phonetic elements: signs with sound values
- Logographic elements: semantic signs



(Robinson 2002)

#### Classifying writing systems based on number of characters

Alphabetic  $(\sim 25 \text{ signs})$ 



Syllabic (~100 signs)

あa	621	うu	え・	お・
かka	き ki	< ku	け ke	こ ko
දී sa	L shi	す su	せ se	₹ so
たta	ち chi	つ tsu	て te	と <b>to</b>
な na	にni	& nu	ね ne	の no
は ha	ひ hi	یک، fu	🔨 he	ほ ho
ま ma	みmi	む mu	Ø me	も mo
やya		<i>И</i> Ф уи		よyo
ら ra	りri	3 ru	th re	ろ <b>ro</b>
わwa				を (w)o
٨n				

E.g., Latin

E.g., Japanese Kana

Alphabets themselves further distinguished into Pure Alphabets: distinct letters for consonants & vowels

Abugida: vowels modify characters for consonants Abjad: vowels are omitted as they are implied rather than being explicit

Ideographic (>50000 signs)

是		我	的	我	我	意	敬	你
我	概	出	様	求	們	能	願	在
心	是	惡	兒	你	需	殼	你	天
所	你	爲	不	觅	用	成	的	上
願	的	的	要	我	的	就	國	願
的	直	是	由	的	粮	在	圖	你
啊	到	那	我	債	食	地	降	的
	世	國	入	照	求	如	臨	名
	世	權	迷	我	你	在	願	兒
	代	勢	願	免	4	天	你	被
	代.	榮	你	人	日	_	的	人
	這	耀	敥	債	給	様	旨	尊

#### E.g., Chinese

#### Logo-syllabic (~900 signs)



E.g., Sumerian cuneiform

Arthur Conan Doyle, The Adventure of the Dancing Men, in Return of SherlockHolmes (1903)



Granada Television 1984

长光大大大大大大大大

Criminal's second message

汉父父子子子父父子

Criminal's third message

XYXX XXXYX

Elsie's response

XX4XF

Criminal's final message

Sherlock Holmes' message to the criminal

ATTY YXYXX TAAX

Criminal's second message



Criminal's third message

えぞき X XXXXX

Elsie's response

Criminal's final message

机小水水水水水水水水水水水水

Sherlock Holmes' message to the criminal えれまだ エンドズ エム イエン

### Frequency of English Characters



#### Samuel Morse 1791-1872



To increase the efficiency of encoding, Morse code was originally designed so that the length of each symbol is approximately inverse to the frequency of occurrence of the character that it represents in text of the English language.

But Morse's character frequency counts are based on pieces of text considers the commonest words multiple times



If one considers all distinct words so that word frequency does not produce artefacts, e.g., by considering words from a dictionary, the highest frequency characters are

#### EARIOT NSLCUD

In either case, E is the commonest letter

Criminal's second message

Criminal's third message

*X*<sub>E</sub> ⊢<sup>1</sup> E <sup>\*</sup>/<sub>F</sub>

Elsie's response

Criminal's final message

Sherlock Holmes' message to the criminal えれまれ インドンド エム イエー

えええ E E えよく E

ÅE⊢ĴE ⋡

 Image: Second message

XX E ダチダイ E 炎 Criminal's third message

Elsie's response

There are no word boundaries – could the flags indicate word endings ?

Criminal's final message

Sherlock Holmes' message to the criminal えれまだ ナキドガ また オキュズ

为为王子E为人E 为为王人 Criminal's second message 光XE汽子之YE装 Criminal's third message えええ E E Z J Y Y E Elsie's response E E No flags – single word! Criminal's final message Sherlock Holmes' message to the criminal 对我我父女孩子孩孩 我我我父

Criminal's second message 光式 E 浅子文Y E 芸 **Possible English** words of the form \* E \* E \* are Criminal's third message えええ E E Z J Y Y E **SEVER** Elsie's response LEVER EHEK NEVER Criminal's final message Sherlock Holmes' message to the criminal 对我我父女孩子孩孩 我我我父

 Image: Second message

 Image: Second message

光式 E 汽子气入 E 芸

Criminal's third message

Elsie's response

えぞき E E きょうう N E V E R Assuming that this is a conversation between two people, very likely that it is NEVER

Ε

Criminal's final message

Sherlock Holmes' message to the criminal えれまれ インドンドン イン・イン・





























## Substitution ciphers in real life

A typical example of coded messages found in prisons in California sent to Stanford's Statistics Department for decryption

 $\frac{A}{A} = \frac{A}{A} = \frac{A}$ 

How to find the unknown function f such f : {code space}  $\rightarrow$  {English alphabet}

## Markov chain Monte Carlo

A standard approach to deciphering is to use the statistics of written English to guess at probable choices for f, try these out, and see if the decrypted messages make sense.

The statistics is obtained from a text corpus. The proportion of consecutive text symbols from x to y gives a matrix M(x; y) of  $I^{st}$  order transitions from one symbol to another. Then the plausibility of the mapping f

$$PI(f) = \prod_{i} M(f(s_{i}), f(s_{i+1}))$$

where s<sub>i</sub> runs over consecutive symbols in the coded message.

Mappings f which have high values of PI(f) are good candidates for deciphering. To maximize the plausibility we can run the following Markov Chain Monte Carlo algorithm

- Start with a preliminary guess, say f.
- Compute  $\operatorname{Pl}(f)$ .
- Change to  $f_*$  by making a random transposition of the values f assigns to two symbols.
- Compute  $Pl(f_*)$ ; if this is larger than Pl(f), accept  $f_*$ .
- If not, flip a  $Pl(f_*)/Pl(f)$  coin; if it comes up heads, accept  $f_*$ .
- If the coin toss comes up tails, stay at f.

# "Reps do it"

The space of possible f s is extremely large – can the Monte Carlo guided random walk achieve the "correct" f in reasonable time, or at all?

ENTER HAMLET HAM TO BE OR NOT TO BE THAT IS THE QUESTION WHETHER TIS

#### The text is scrambled at random and the Monte Carlo algorithm run

100 ER ENOHDLAE OHDLO UOZEOUNORU O UOZEO HD OITO HEOQSET IUROFHE HENO ITORUZAEN 200 ES ELOHRNDE OHRNO UOVEOULOSU O UOVEO HR OITO HEOQAET IUSOPHE HELO ITOSUVDEL 300 ES ELOHANDE OHANO UOVEOULOSU O UOVEO HA OITO HEOQRET IUSOFHE HELO ITOSUVDEL 400 ES ELOHINME OHINO UOVEOULOSU O UOVEO HI OATO HEOQRET AUSOWHE HELO ATOSUVMEL 500 ES ELOHINME OHINO UODEOULOSU O UODEO HI OATO HEOQRET AUSOWHE HELO ATOSUDMEL 600 ES ELOHINME OHINO UODEOULOSU O UODEO HI OATO HEOQRET AUSOWHE HELO ATOSUDMEL 900 ES ELOHANME OHANO UODEOULOSU O UODEO HA OITO HEOQRET IUSOWHE HELO ITOSUDMEL 1000 IS ILOHANMI OHANO RODIORLOSR O RODIO HA OETO HIOQUIT ERSOWHI HILO ETOSRDMIL. 1100 ISTILOHANMITOHANOT ODIO LOS TOT ODIOTHATOEROTHIOQUIRTE SOWHITHILOTEROS DMIL 1200 ISTILOHANMITOHANOT ODIO LOS TOT ODIOTHATOEROTHIOQUIRTE SOWHITHILOTEROS DMIL 1300 ISTILOHARMITOHAROT ODIO LOS TOT ODIOTHATOENOTHIOQUINTE SOWHITHILOTENOS DMIL 1400 ISTILOHAMRITOHAMOT OFIO LOS TOT OFIOTHATOENOTHIOQUINTE SOWHITHILOTENOS FRIL 1600 ESTEL HAMRET HAM TO CE OL SOT TO CE THAT IN THE QUENTIOS WHETHEL TIN SOCREL 1700 ESTEL HAMRET HAM TO BE OL SOT TO BE THAT IN THE QUENTIOS WHETHEL TIN SOBREL 1800 ESTER HAMLET HAM TO BE OR SOT TO BE THAT IN THE QUENTIOS WHETHER TIN SOBLER 1900 ENTER HAMLET HAM TO BE OR NOT TO BE THAT IS THE QUESTION WHETHER TIS NOBLER 2000 ENTER HAMLET HAM TO BE OR NOT TO BE THAT IS THE QUESTION WHETHER TIS NOBLER

The simple optimization process converges to original text in a few thousand steps

## Substitution ciphers in real life

Decipherment of the typical example of coded messages found in prisons in California sent to Stanford's Statistics Department for decryption

to bat-rb. con todo mi respeto. i was sitting down playing chess with danny de emf and boxer de el centro was sitting next to us. boxer was making loud and loud voices so i tell him por favor can you kick back homie cause im playing chess a minute later the vato starts back up again so this time i tell him con respecto homie can you kick back. the vato stop for a minute and he starts up again so i tell him check this out shut the f\*\*k up cause im tired of your voice and if you got a problem with it we can go to celda and handle it. i really felt disrespected thats why i told him. anyways after i tell him that the next thing I know that vato slashes me and leaves. dy the time i figure im hit i try to get away but the c.o. is walking in my direction and he gets me right dy a celda. so i go to the hole. when im in the hole my home boys hit doxer so now "b" is also in the hole. while im in the hole im getting schoold wrong and

# Phaistos Disk



Found in 1908 by an Italian archaeologist Louis Pernier in connection with the excavations of Phaistos in Crete.

features 241 tokens, comprising 45 distinct signs, apparently made by pressing hieroglyphic "seals" into the soft clay - spiraling clockwise toward the center



## Phaistos Disk sign statistics

Sign			Frequency																
No. <b>\$</b>	Glyph	Font	Name	A \$	B \$	A+B ▼	33	33		TUNNY	2	2 4 6		03	Ð	TATTOOED HEAD	2	0	2
02	ET.		PLUMED HEAD	14	5	19	40	$\Omega$		OX BACK	3	3	6	09	ß)	TIARA	0	2	2
07	$\cap$		HELMET	3	15	18	45			WAVY BAND	2	4	6	14	B	MANACLES	1	1	2
								Ш						16	V	SAW	0	2	2
12	63		SHIELD	15	2	17	08	P		GAUNTLET	1	4	5	20	6	DOLIUM	0	þ	2
27	ſ		HIDE	10	5	15	22	R		SLING	0	5	5	21		СОМВ	2	0	2
18	$\ll$		BOOMERANG	6	6	12	31	e (je)		EAGLE	5	0	5	28	J	BULLS LEG	2	0	2
01	Ř		PEDESTRIAN	6	5	11	06	135		WOMAN	2	2	4	41	Ĩ	FLUTE	2	0	2
23	8			5	6	11	10			ARROW	4	0	4	04		CAPTIVE	1	0	1
20	U				0		36	1000		VINE	0	4	4	05	ŝ	CHILD	0	1	1
29	é l		CAT	3	8	11	07	∆ %g		DA DVDU O	0	0	4	11	¢	BOW	1	0	1
35	SFE		PLANE TREE	5	6	11	37	t, see		PAPTRUS	2	2	4	15	÷	MATTOCK	0	1	1
	۹ ۲				_		38	£€ €		ROSETTE	3	1	4	17	٩)	LID	1	0	1
25	Ç		SHIP	2	5	7	39	Y		LILY	1	3	4	30	~ ~3	RAM	0	1	1
13	CIIIID		CLUB	3	3	6	19	$\mathbb{V}$		CARPENTRY PLANE	3	0	3	42		GRATER	0	1	1
24			BEEHIVE	1	5	6		ึง						12	<u></u>		0	1	1
	000						32	Ð		DOVE	2	1	3	43	¥	STRAINER	0		1
26	8		HORN	5	1	6	34	$\mathbb{C}$		BEE	1	2	3	44	<3	SMALL AXE	1	0	1

# Inferring the signary size

The frequency of occurrence of symbols on the Phaistos Disc. (Reference number of symbol in Evans's list)

	44																			
	43	41																		
	42	28																		
	30	21			4	10														
	17	20			4	15														
	15	16	38	39	36 3	33					35									
	11	14	34	37	31 2	26					29									
	5	9	32	10	22 2	24					23				46					
	4	3	19	6	8	13	25				1	18			27		12	7	2	
t=0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20

We make the general assertion that, in a small sample of an alphabetic or syllabic writing system consisting of  $L_1$  characters of  $M_1$  different kinds, the probable number of symbols in the alphabet or syllabary is, subject to various restrictions discussed below, given approximately by  $L_1^2/(L_1-M_1) - L_1$ .

 $\Rightarrow$  If the Phaistos disk is genuine (the writing system follows the same general frequency distribution as other known alphabetic/syllabic systems), about 10 more symbols in the signary



Alan L Mackay

