



# Decipherment

Kevin Knight  
Information Sciences Institute  
University of Southern California

includes joint work with:

**S. Ravi** (USC/ISI, now Google), **Q. Dou**, **K. Yamada** (USC/ISI)

**B. Megyesi, C. Schaefer** (Uppsala Univ.)

**R. Barzilay, B. Snyder** (MIT)

**S. Reddy** (Univ. Chicago, now Dartmouth)

ACL Tutorial

August 2013

## Why Decipherment?

- It's fun and cool
  - ancient languages
  - secret societies
- Breaking codes was the first application of NLP
- Intellectual root of NLP
  - language models, log-odds ratios, smoothing
  - ASR and MT use "decoders"
- View foreign language as a code for English

# Decipherment Papers by ACL-ers

- "Unsupervised Analysis for Decipherment Problems," (K. Knight, A. Nair, N. Rathod, and K. Yamada), Proc. ACL-COLING, 2006.  
(Rejected four times previously, but OK!)
- "Attacking Decipherment Problems Optimally with Low-Order N-gram Models," (S. Ravi and K. Knight), *Cryptologia*, 2009.
- "Probabilistic Methods for a Japanese Syllable Cipher," (S. Ravi and K. Knight), Proc. ICCPOL, 2009.
- "A Statistical Model for Lost Language Decipherment," (B. Snyder, R. Barzilay, and K. Knight), Proc. ACL, 2010.
- "An Exact A\* Method for Deciphering Letter-Substitution Ciphers," (E. Corlett and G. Penn), Proc. ACL, 2010.
- "Deciphering Foreign Language," (S. Ravi and K. Knight), Proc. ACL, 2011.
- "The Copiale Cipher," (K. Knight, B. Megyesi, and C. Schaefer), Proc. ACL BUCC, 2011.
- "Bayesian Inference for Zodiac and Other Homophonic Ciphers," (S. Ravi and K. Knight), Proc. ACL, 2011.
- "What We Know About the Voynich Manuscript," (S. Reddy and K. Knight), Proc. ACL LaTECH, 2011.
- "Simple Effective Decipherment via Combinatorial Optimization," (T. Berg-Kirkpatrick and D. Klein), Proc. EMNLP, 2011.
- "Decoding Running Key Ciphers," (S. Reddy and K. Knight), Proc. ACL, 2012.
- "Large Scale Decipherment for Out-of-Domain Machine Translation," (Q. Dou and K. Knight), Proc. EMNLP, 2012.
- "Deciphering Foreign Language by Combining Language Models and Context Vectors," (M. Nuhn, A. Mauser, and H. Ney), Proc. ACL, 2012.
- "Decipherment Complexity in 1:1 Substitution Ciphers," (M. Nuhn, and H. Ney), Proc. ACL, 2013.
- "Beam Search for Solving Substitution Ciphers," (M. Nuhn, J. Schamper, and H. Ney), Proc. ACL, 2013.
- "Scalable decipherment for machine translation via hash sampling," (S. Ravi), Proc. ACL, 2013.
- "Unsupervised Consonant-Vowel Prediction over Hundreds of Languages," (Y. Kim and B. Snyder), Proc. ACL, 2013.

## Outline

- Classical military/diplomatic ciphers (15 mins)
- Foreign language as a code (10 mins)
- Automatic decipherment (55 mins)
- **Break (30 mins)**
- Unsolved ciphers (40 mins)
- Writing as a code for speech (20 mins)
- Undeciphered writing systems (15 mins)
- Conclusions (15 mins)

## Classical military/diplomatic ciphers

### Letter Substitution Cipher

- Encipherment key:

PLAIN : ABCDEFGHIJKLMNOPQRSTUVWXYZ

CIPHER : PLOKMIJNUHBYGVTFCRDXESZAQW

- Plaintext: **HELLO WORLD . . .**

- Ciphertext: **NMYYT ZTRYK . . .**

- Key itself doesn't change: "simple substitution"

- What key, if applied to the ciphertext, would yield sensible plaintext?

KDCY LQZKTLJKX CY MDBCYJQL: "TR

HYD FKXC, FQ MKX RLQQIQ HYDL

MKL DXCTW RDCDLQ JQMNKXTMB

PTBMYEQL K FKH CY LQZKTL TC."

```

A
B 3
C 8
D 7 #
E 1 .
F 3 .
G
H 3 .
I 1 .
J 3 .
K 10 ##### V
L 10 ##
M 6 #
N 1 .
O
P 1 .
Q 10 ##### V
R 3 .
S
T 7 ### V
U
V
W 1 .
X 5
Y 7 ### V
Z 2 .

```

a e.a .a .e .

KDCY LQZKTLJKX CY MDBCYJQL: "TR

. .a .e a . ee.e .

HYD FKXC, FQ MKX RLQQIQ HYDL

a . . e .e .a

MKL DXCTW RDCDLQ JQMNKXTMB

. .e a .a. e.a

PTBMYEQL K FKH CY LQZKTL TC."

```

A
B 3
C 8
D 7 #
E 1 .
F 3 .
G
H 3 .
I 1 .
J 3 .
K 10 ##### V
L 10 ##
M 6 #
N 1 .
O
P 1 .
Q 10 ##### V
R 3 .
S
T 7 ### V
U
V
W 1 .
X 5
Y 7 ### V
Z 2 .

```

didn't create "ae"

a e.ao .a .e o.  
**KDCY LQZKTLJKX CY MDBCYJQL: "TR**  
. .a .e a . ee.e .  
**HYD FKXC, FQ MKX RLQQIQ HYDL**  
a o. . e .e .a o  
**MKL DXCTW RDCDLQ JQMNKXTMB**  
. o .e a .a . e.ao o  
**PTBMYEQL K FKH CY LQZKTL TC."**

A	
B	3
C	8
D	7 #
E	1 .
F	3 .
G	
H	3 .
I	1 .
J	3 .
K	10 ##### V
L	10 ##
M	6 #
N	1 .
O	
P	1 .
Q	10 ##### V
R	3 .
S	
T	7 ### V
U	
V	
W	1 .
X	5
Y	7 ### V
Z	2 .

don't like "ao" – back up!

### Pattern word dictionaries

**KDCY (LQZKTLJKX) CY MDBCYJQL: "TR**

**abcdeafdg**

**HYD FKXC (MKX (RLQQIQ) HYDL**

**abccdc**

**MKL DXCTW (CDLQ JQ (XTMB**

**abcdefghijklm**

**PTBMYEQL (KH CY L (L TC."**

abnegated  
abnegates  
advocator  
airedales  
alienages  
alienated  
alienates  
amperages  
cadencies  
capricorn  
cogencies  
escapeway  
healthily  
imbeciles  
imperiled  
incurious  
inherited  
injurious  
landslide  
octagonal  
oklahoman  
overboard  
repairman  
sacristy  
unrebuked  
unsecured

basses  
bassos  
bosses  
breeze  
budded  
...  
cheese  
cusses  
dosses  
finnan  
fleece  
fosses  
freeze  
...  
terror  
tosses  
tweeze  
wadded  
wheeze

consumptively  
copyrightable  
documentarily  
lycanthropies  
musicotherapy  
semivoluntary  
subordinately  
unpredictably

OR, NORWEGIAN!  
filmprodusent  
kurspamelding  
publikasjoner  
stylemarginpx  
upproblematisk

## Fundamental Questions

- How much English does a system need to know to break a cipher?
- How long does the cipher need to be, to admit a unique solution?
- How much computational effort is required to decipher?

and...

## How to Make Things Harder?

- Homophonic cipher
  - ciphertext values from 00 to 99
    - A → 02, 14, 16, 22, 49, 51, 58, 90
    - B → 04, 76
    - C → 15, 56, 71
    - etc
  - flattens out ciphertext distribution
    - “a cab...” becomes “22 56 14 04...”
  - still deterministic in the deciphering direction
- Polyalphabetic ciphers
  - the secret key changes at each plaintext letter token
    - e.g., rotate through 10 different keys
- Transposition ciphers

or perhaps:

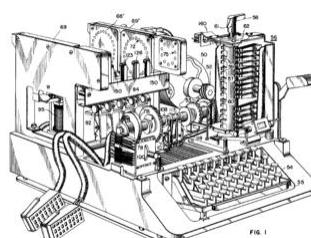
A = ȳ ȷ ȴ ȵ
B = ȫ
C = ȶ ȸ
D = ȷ
E = Ȣ ȣ Ȥ ȥ Ȧ ȧ Ȩ ȩ
F = ȫ
G = ȳ ...

## Cipher Types

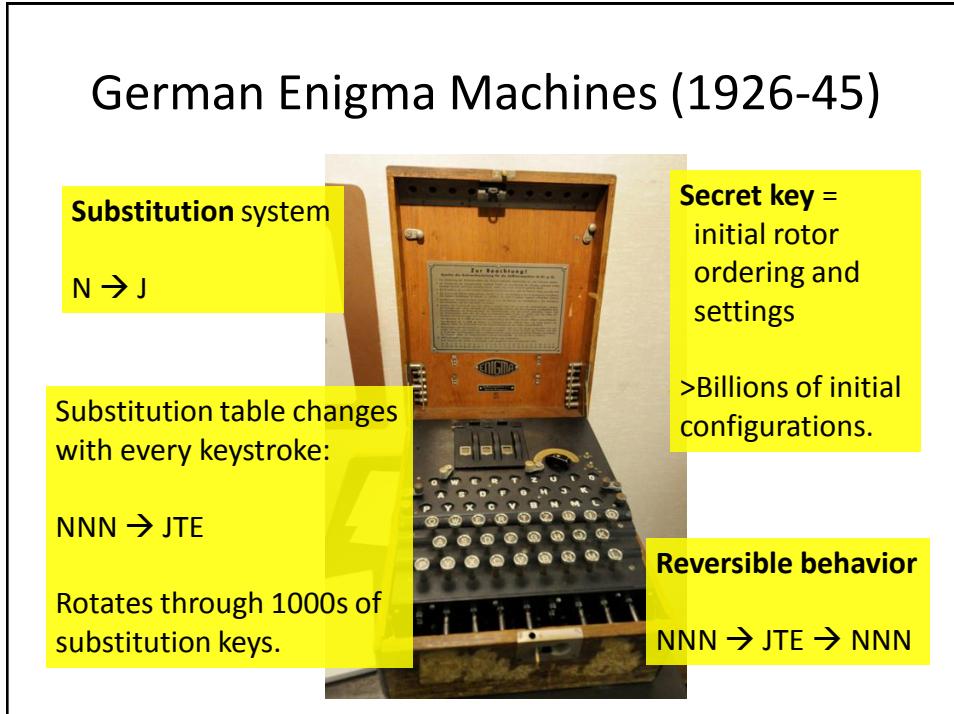
- [http://cryptogram.org/cipher\\_types.html](http://cryptogram.org/cipher_types.html)
  - documents ~70 types
- E.g., RUNNING KEY cipher
  - key = agreed-upon standard English text
  - ciphertext(i) = [ plaintext(i) + key(i) ] mod 26
  - effectively uses 26 substitution keys
  - breakable!
  - we search for a key and (resulting) plaintext that are both natural language

## How to Make Things Efficient?

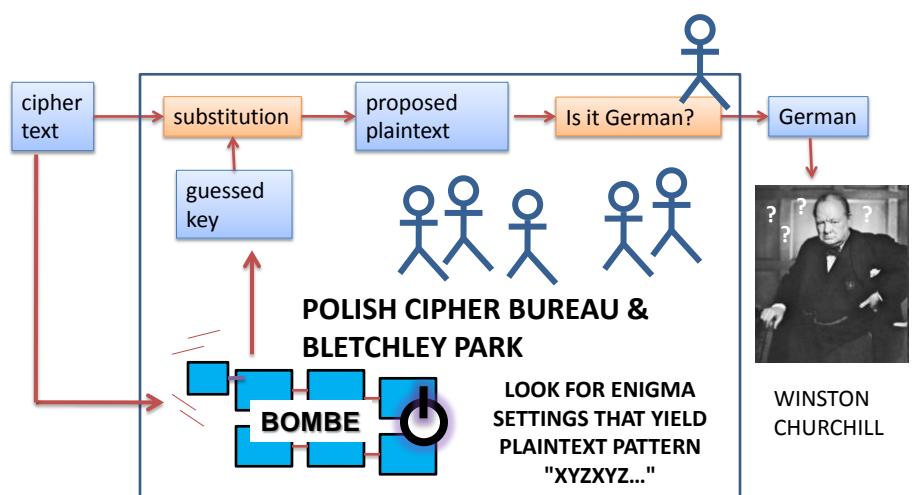
- Mechanical encryption/decryption devices



## German Enigma Machines (1926-45)



## Breaking Enigma



# Enigma

- Mathematical breakthroughs:
  - Log-odds for weight of evidence [Good, Turing]
  - Smoothing with prior [Good, Turing]
  - Information theory [Shannon]
  
- 1945: War ends
- 1973: Wartime Enigma decipherment leaked
- 1975: Last surplus Enigma given to developing countries
- 1996: One Turing Enigma treatise declassified
- 2012: Another declassified (but have to go to England)

elegant,  
powerful,  
widely-applicable  
mathematics

## Turing Enigma Treatise

(aka NR 964, Box 201, RG 457, aka "The Prof's Book")

140pp (written sometime between 1939 and 1942)

One method is to try independently all the possible positions for the middle wheel. We shall want to know the middle wheel couplings which are consequences of these various assumptions. This can be done by using inverse rods for the middle wheel. The rods are paired off in pairs of R.H.W. couplings, i.e. M.W. output. This has been done for the case of fx, ep which arose in the DANZIGVON crib in Fig 55, assuming that the middle wheel does not rotate. The pairs in each column of this set up give possible M.W. couplings. We have now to find out whether these couplings are good or bad. Our procedure is rather different according as the U.K.W. does or does not rotate. In the case that the U.K.W. does not rotate it will be necessary to use a Foss sheet (the rows and columns lettered preferably with the diagonal alphabet) in which, in the RW square are entered the positions of the left hand wheel at which the RW is one of the pairs in the L.H.W. output alphabet Fig 51. This is known as the 'short catalogue' for this wheel.

elegant,  
powerful, war-winning  
widely-applicable-  
mathematics

if we worked this  
hard on machine  
translation ...



# Foreign language as a code

## Alan Turing, on Thinking Machines

Instead we propose to try and see what can be done with a 'brain' which is more or less without a body, providing at most, organs of sight speech and hearing. We are then faced with the problem of finding suitable branches of thought for the machine to exercise its powers in. The following fields appear to me to have advantages:-

- (i) Various games e.g. chess, noughts and crosses, bridge, poker.
- (ii) The learning of languages.
- (iii) Translation of languages.
- (iv) Cryptography.
- (v) Mathematics.



of these (i), (iv), and to a lesser extent (iii) and (v) are good in that they require little contact with the outside world. For instance in order that the machine should be able to play chess its only organs need be 'eyes' capable of distinguishing the various positions on a specially made board, and means for announcing its own moves. Mathematics should preferably be restricted to branches where diagrams are not much used. Of the above possible fields the learning of languages would be the most impressive, since it is the most human of these activities. This field seems however to depend rather too much on sense organs and locomotion to be feasible.

The field of cryptography will perhaps be the most rewarding. There is a remarkably close parallel between the problems of the physicist and those of the cryptographer. The system on which a message is enciphered corresponds to the laws of the universe, the intercepted messages to the evidence available, the keys for a day or a message to important constants which have to be determined. The correspondence is very close, but the subject matter of cryptography is very easily dealt with by discrete machinery, physics not so easily.

# Statistical Machine Translation

"When I look at an article in Russian, I say: this is really written in English, but it has been coded in some strange symbols. I will now proceed to decode." -- Warren Weaver (1947)

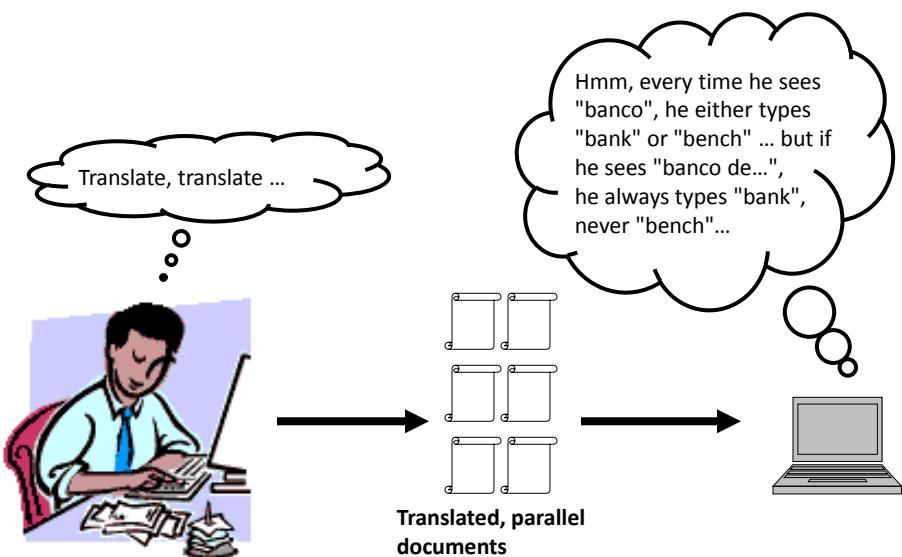


OUR HERO

Weaver saw a colleague decoding intercepts into Turkish, without "knowing" Turkish.

... maybe a computer could translate into English without "knowing" English?

# Statistical Machine Translation



# Parallel Corpus

12 English sentences in English and Spanish.

1a. Garcia and associates . 1b. Garcia y asociados .	7a. the clients and the associates are enemies . 7b. los clientes y los asociados son enemigos .
2a. Carlos Garcia has three associates . 2b. Carlos Garcia tiene tres asociados .	8a. the company has three groups . 8b. la empresa tiene tres grupos .
3a. his associates are not strong . 3b. sus asociados no son fuertes .	9a. its groups are in Europe . 9b. sus grupos estan en Europa .
4a. Garcia has a company also . 4b. Garcia tambien tiene una empresa .	10a. the modern groups sell strong pharmaceuticals . 10b. los grupos modernos venden medicinas fuertes .
5a. its clients are angry . 5b. sus clientes estan enfadados .	11a. the groups do not sell zanzanine . 11b. los grupos no venden zanzanina .
6a. the associates are also angry . 6b. los asociados tambien estan enfadados .	12a. the small groups are not modern . 12b. los grupos pequenos no son modernos .

# Parallel Corpus

12 English sentences in Centauri and Arcturan.

1a. ok-voon ororok sprok . 1b. at-voon bichat dat .	7a. lalok farok ororok lalok sprok izok enemok . 7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anok plok sprok . 2b. at-drubel at-voon pippat rrat dat .	8a. lalok brok anok plok nok . 8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok . 3b. totat dat arrat vat hilat .	9a. wiwok nok izok kantok ok-yurp . 9b. totat nnat quat oloaat at-yurp .
4a. ok-voon anok drok brok jok . 4b. at-voon krat pippat sat lat .	10a. lalok mok nok yorok ghirok clok . 10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok . 5b. totat jjat quat cat .	11a. lalok nok errrok hihok yorok zanzanok . 11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok . 6b. wat dat krat quat cat .	12a. lalok rarok nok izok hihok mok . 12b. wat nnat forat arrat vat gat .

## Centauri/Arcturan

Your assignment, translate this to Arcturan: farok crrok hihok yorok clok kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anok plok sprok .	8a. lalok brok anok plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anok drok brok jok .	10a. lalok mok nok yorok ghirok clok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrok hihok yorok zanzanok .
5b. totat jjat quat cat .	11b. wat nnat arrat mat zanzanat .
6a. lalok sprok izok jok stok .	12a. lalok rarok nok izok hihok mok .
6b. wat dat krat quat cat .	12b. wat nnat forat arrat vat gat .

## Centauri/Arcturan

Your assignment, translate this to Arcturan: farok crrok hihok yorok **clok** kantok ok-yurp

1a. ok-voon ororok sprok .	7a. lalok farok ororok lalok sprok izok enemok .   /
1b. at-voon bichat dat .	7b. wat jjat bichat wat dat vat eneat .
2a. ok-drubel ok-voon anok plok sprok .	8a. lalok brok anok plok nok .
2b. at-drubel at-voon pippat rrat dat .	8b. iat lat pippat rrat nnat .
3a. erok sprok izok hihok ghirok .	9a. wiwok nok izok kantok ok-yurp .
3b. totat dat arrat vat hilat .	9b. totat nnat quat oloat at-yurp .
4a. ok-voon anok drok brok jok .	10a. lalok mok nok yorok ghirok <b>clok</b> .   X / / ••• process of 4b. at-voon krat pippat sat lat .
5a. wiwok farok izok stok .	10b. wat nnat gat mat bat hilat . elimination
5b. totat jjat quat cat .	11a. lalok nok crrok hihok yorok zanzanok .   /
6a. lalok sprok izok jok stok .	11b. wat nnat arrat mat zanzanat .   / / /
6b. wat dat krat quat cat .	12a. lalok rarok nok izok hihok mok .   / / /
	12b. wat nnat forat arrat vat gat .

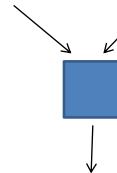
## Learn Translation Knowledge from Non-Parallel Text?

English/Albanian  
Parallel text



Translation model

English text      Albanian text



Is this what Weaver had in mind?  
We'll come back to this idea.

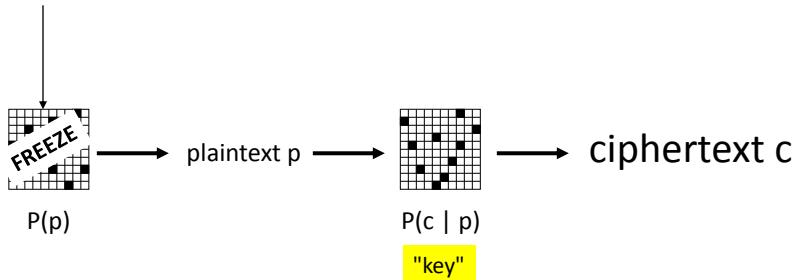
Automatic decipherment

# Letter Substitution Cipher

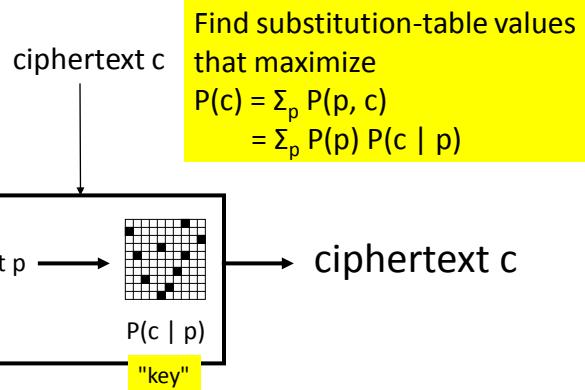
ciphertext c

# Letter Substitution Cipher

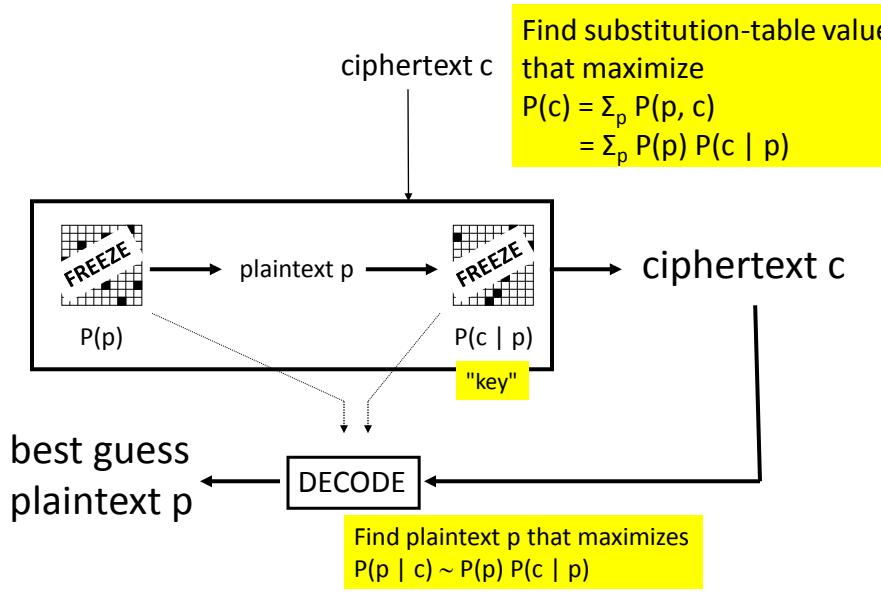
plaintext samples,  
unrelated to ciphertext



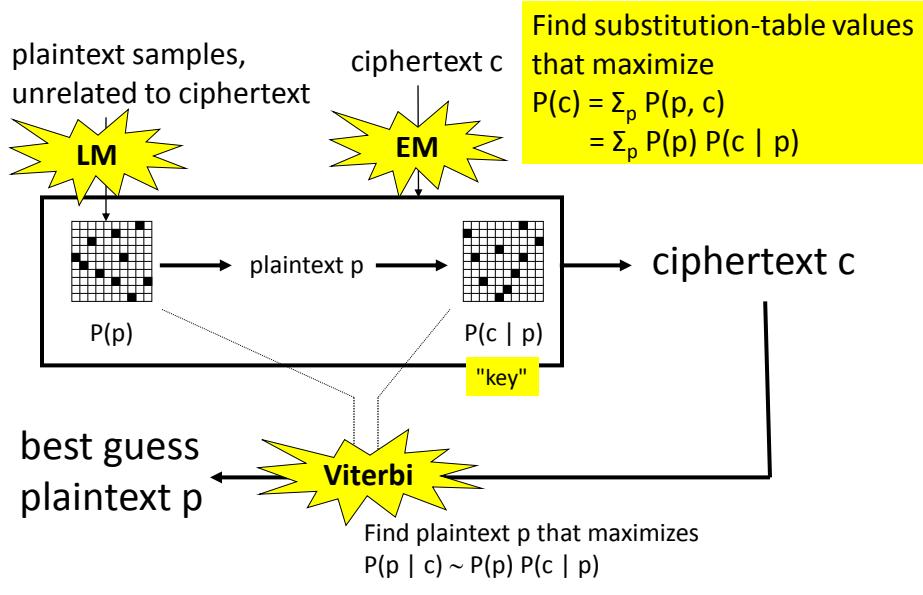
## Letter Substitution Cipher



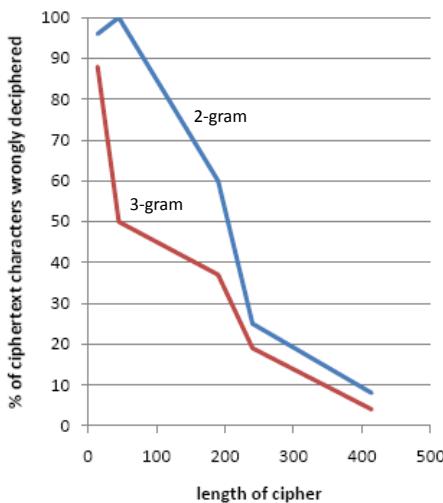
## Letter Substitution Cipher



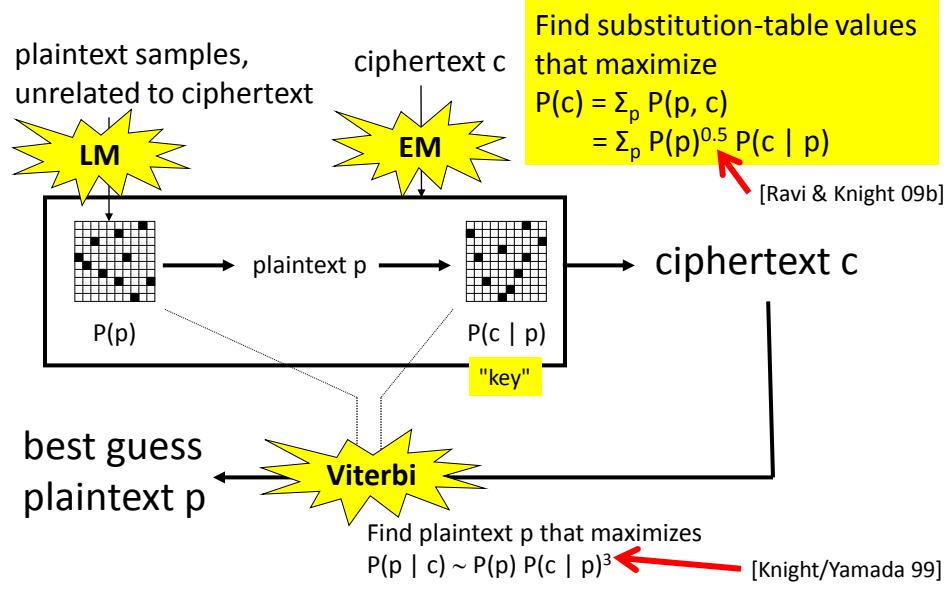
## Letter Substitution Cipher



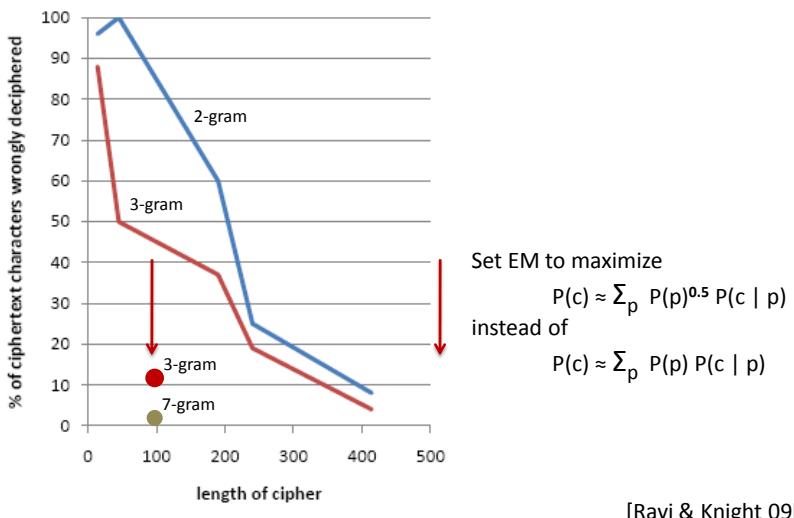
## Decipherment Accuracy vs. Cipher Length



## Letter Substitution Cipher

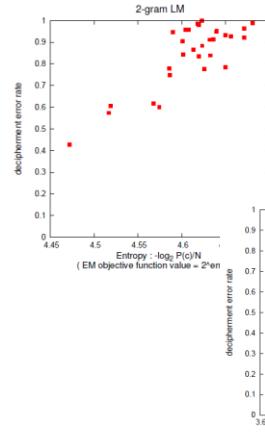
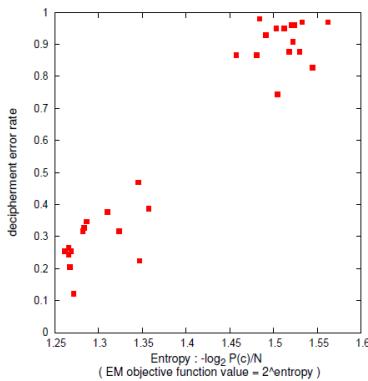


## Reducing LM Weight During EM



## Random Restarts are Critical

English 98-letter cipher, 3-gram LM

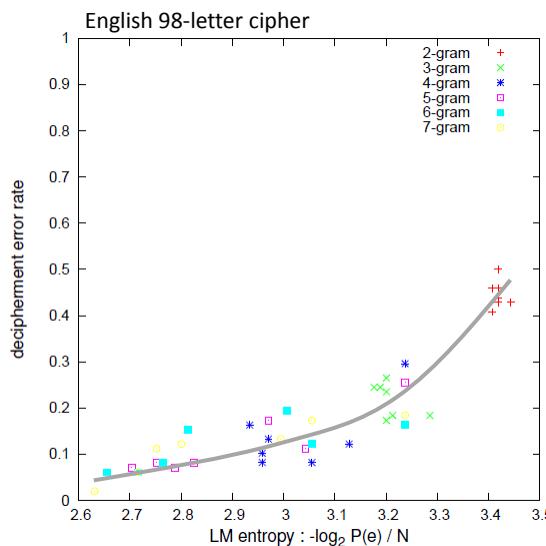


Japanese  
syllable  
cipher

even people do restarts!

[Ravi & Knight 09b]

## Good Language Models are Critical



[Ravi & Knight 09b]

## Searching for Deterministic Keys

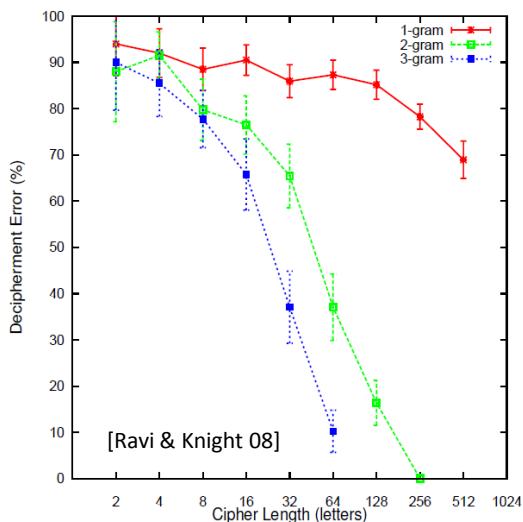
- Peleg & Rosenfeld, 1979
  - relaxation search
  - ...
- Ravi & Knight, 2008
  - ILP, exact search
- Corlett & Penn, 2010
  - A\* exact search
- Nuhn, Schamper, and Ney, 2013
  - beam search

## Deterministic Keys

- \* Use ILP to search only deterministic keys.
- \* Exact, no restarts.

Using 2-gram letter-based LM

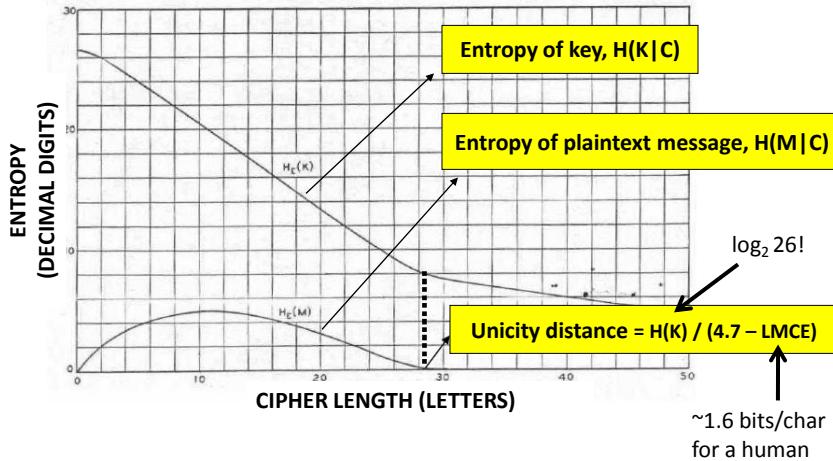
Cipher Length	EM error	ILP error
52	85 %	<b>21 %</b>
98	45 %	<b>12 %</b>
414	10 %	<b>0.5 %</b>



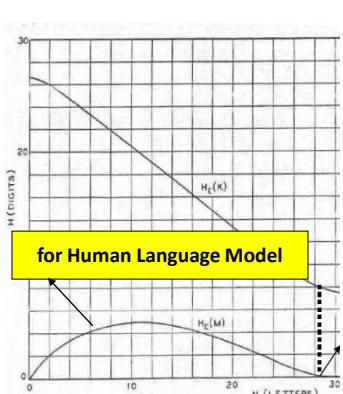
## [Shannon 46, 49]

### "Communication Theory of Secrecy Systems"

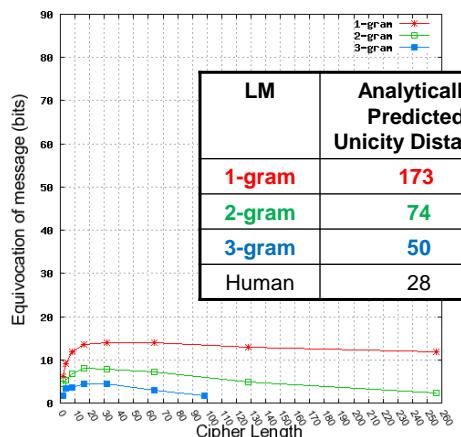
- Shannon analytically predicted uncertainty about key and message
- Graphed it for a human-level language model



### Verifying Shannon's Prediction of Plaintext Message Uncertainty



ANALYTIC CURVES (Shannon)



ACTUAL CURVES

[Ravi &amp; Knight 08]

# Some Recent Historical Decipherments

- Jefferson cipher (L. Smithline)
    - <http://online.wsj.com/article/SB124648494429082661.html>
    - For more than 200 years, buried deep within Thomas Jefferson's correspondence and papers, there lay a mysterious cipher -- a coded message that appears to have remained unsolved. Until now.
  - Civil War ciphers (K. Boklan)
    - Cryptologia, 30:340–345
    - We study a previously undeciphered Civil War cryptogram, limiting ourselves to pencil and paper, and discover not only a missive of military importance, but in the process identify a new Confederate codeword. Our methods rely not only upon cryptanalysis of the encryption method but also on the exploitation of an elementary mistake.
  - German Naval Enigma
    - <http://www.enigma.hoerenberg.com>
    - The "Breaking German Navy Ciphers" Project was founded in 2012. The goal is to break original radio messages, which were encoded with the famous German ENIGMA cipher machine. Up to now, we've succeeded in deciphering 53 original World War II Enigma M4 messages. Many of these messages had never been broken before, so you can read them for the first time in history.

# Copiale Cipher

*Capročēipvōyrār & Anāyxuacypbg=aghia*

[Knight, Megyesi, Schaefer 11]

## Copiale Cipher

105 pages, 75000 letter tokens, no word spacing, no illustrations.

Some scratch-outs, rare

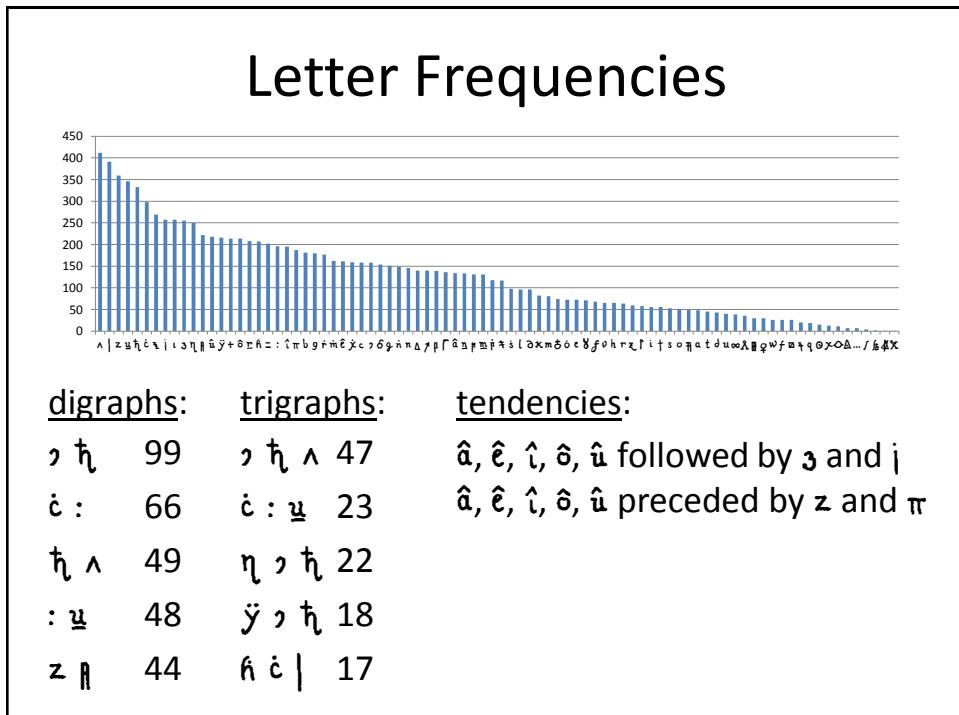
Preview text fragments ("catchwords")

Section headers

Lines ≈ equal length

Paragraphs and section titles always begin with capitalized Roman letters.

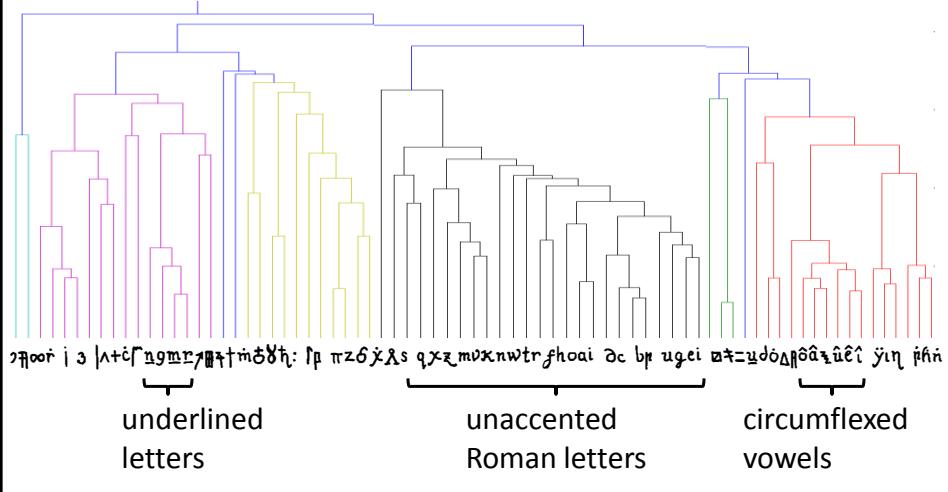
Non-enciphered inscriptions:  
Copiales 3 and Philipp 1866



# Clustering of Cipher Letters

letters grouped if they have similar contexts (L/R neighbors)

## Scipy software



thanks Jon Graehl

# First Decipherment Approach

unaccented Roman letters that cluster:

a b c d e f g h i  
k l m n o p q r s  
t u v w x y z

most common letter = 12%  
least common = very small

x fnglxnacbfxmk  
lbuvcghtrhbkgnxn  
f ggnxkbgbecb ...

Decipher against  
80 plaintext languages.

# Second Decipherment Approach

Homophonic cipher,  
e.g.:



A = γ i l y r  
B = ū  
C = ö ü  
D = ñ  
E = ġ f A ī f ī ī 3  
F = p  
G = ÿ

etc.

# Homophonic Cipher

Result of computer attack on Copiale, using  
80 possible plaintext languages?

**FAIL**

But, slight numerical preference for German

## Cipher Characteristics

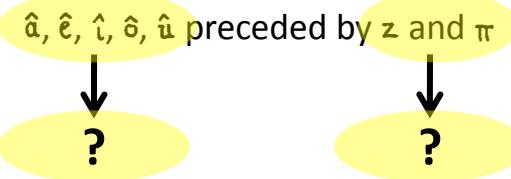
digraphs:

, ĥ	99	, ĥ ḥ	47
č :	66	č : ü	23
ḥ ḥ	49	η , ḥ	22
: ü	48	ÿ , ḥ	18
z ḥ	44	h č	17

trigraphs:

tendencies:

â, ê, î, ô, û followed by ɔ and ɔ



should appear  
adjacent in German text

Make full digraph table for cipher and for German

## Key Observation #1

In Copiale, ɔ almost always followed by ḥ

In German, C almost always followed by H  
(German CH is like English QU)

So guess: ɔ = C, ḥ = H

# One Thing Leads to Another

$$\partial \hat{h} = CH \quad \rightarrow \quad \partial \hat{h} \wedge = CHT \quad \rightarrow \quad \wedge = T ?$$

Each step is guesswork.

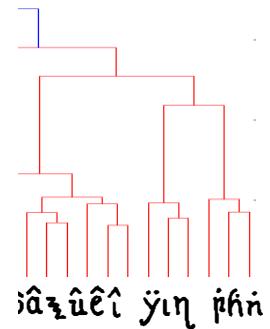
Must be willing to retract.

# Weird task, not knowing German.

No longer care what the book says.

# Cluster diagram crucial:

$$\ddot{y} = l \quad \rightarrow \quad \iota = l, \eta = l$$



# Spring Break 2011

Cipher  
letters,  
in groups

Quite a bit  
of fooling  
around →

## German letters

## German trigraphs

Cipher  
trigraphs

✓ Grid

# Trigraph Decoding Guesses

## Key Observation #2

unaccented Roman letters that cluster:

a b c d e f g h i  
k l m n o p q r s  
t u v w x y z

κτημάριον φύσης ή ζωής πατέρας ή  
ιχθυόφρεστης ή θαλάσσης πατέρας ή  
πατέρας της γης ή θεού της γης ή  
πατέρας της θάλασσης ή θεού της θάλασσης  
πατέρας της γης ή θεού της γης ή  
πατέρας της θάλασσης ή θεού της θάλασσης

# Actually, those are space bars

# Copiale Decipherment

→

weil die sterblichen der  $\Delta$  durch den thürheter besorget und die  $\Delta$  vom dirigirenden  $\Lambda$  durch aufsetzung seines huths geöffnet ist wird der candidat von dem jüngern thürhiter aus einem andern zimmer abgeholet und bey der hand ein und vor des dirigirenden  $\Lambda$  tisch geführet dieser frägt ihn:

erstlich ob er begehrte  zu werden

zweytens denen verordnungen der **O** sich

unterwerffen und ohne Widerspenstigkeit die Lehrzeit ausstehen wolle.

drittens die **A** der **O** gu verschweigen und dazu auf das verbindlichste sich anheischig zu machen gesinnet sey.

der candidat antwortet ja.

# Copiale Decipherment

1

First lawbook

of the  e 

### Secret part.

## First section

## Teachings for April

## First title. Initiation rite

If the safety of the **A** is guaranteed, and the **A** is opened by the chief **A**, by putting on his hat, the candidate is fetched from another room by the younger doorman and by the hand is led in and to the table of the chief **A**, who asks him:

First, if he desires to become **O**.

Secondly, if he submits to the rules of the **O** and without rebelliousness suffer through the time of apprenticeship.

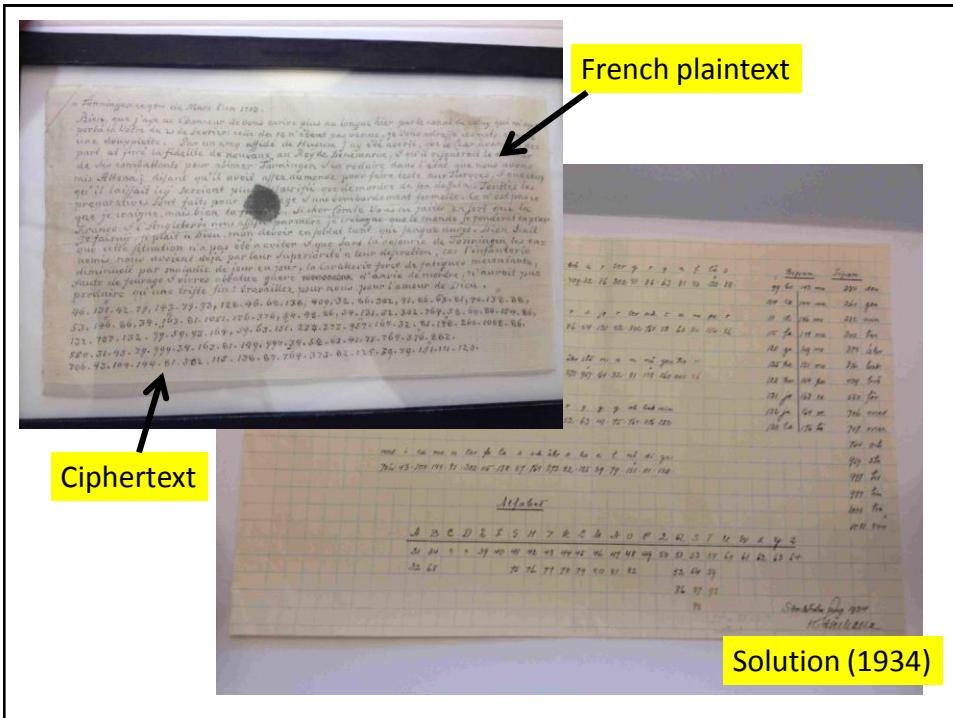
Thirdly, be silent about the ~~A~~ of the ~~O~~ and furthermore be willing to offer himself to volunteer in the most committed way.

The candidate answers yes.



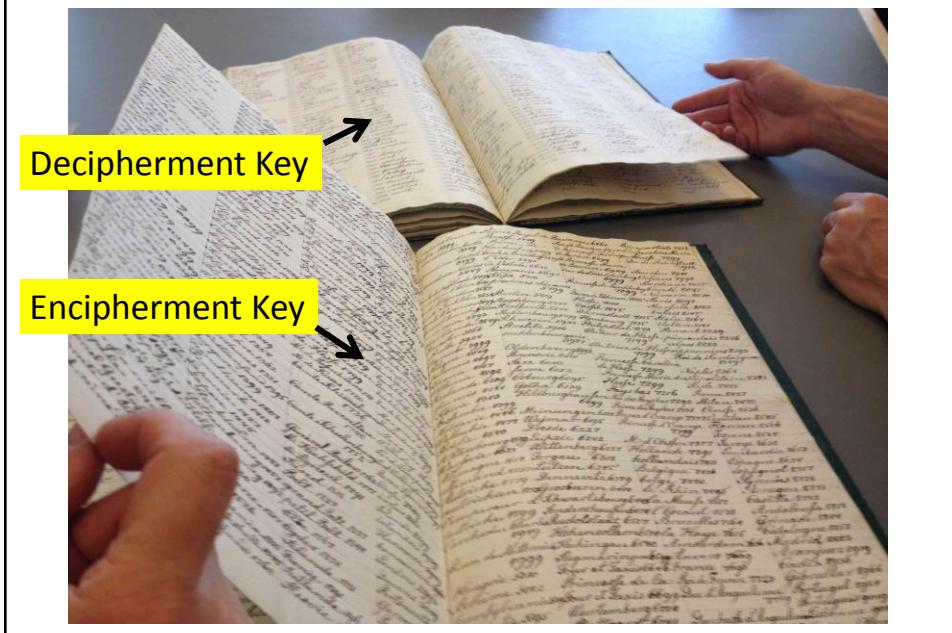
# Historical Archives





# Word Substitution Encipherment Key

## Word Substitution Keys



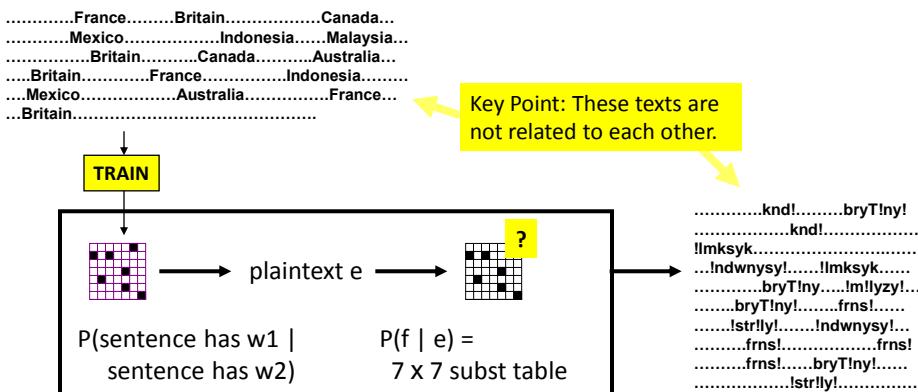
## Word Substitution Keys

<p><b>Numbers/Words Both in Order!</b></p> <table border="1"> <thead> <tr> <th>Number</th> <th>Word</th> </tr> </thead> <tbody> <tr><td>951.</td><td>Bafascon</td></tr> <tr><td>952.</td><td>Congratulation</td></tr> <tr><td>953.</td><td>Congratulera</td></tr> <tr><td>954.</td><td>Congress</td></tr> <tr><td>955.</td><td>Conjunction</td></tr> <tr><td>956.</td><td>Content</td></tr> <tr><td>957.</td><td>Contentam</td></tr> <tr><td>958.</td><td>Contentam</td></tr> <tr><td>959.</td><td>Convere</td></tr> <tr><td>960.</td><td>Conivent</td></tr> <tr><td>961.</td><td>Convens</td></tr> <tr><td>962.</td><td>Conventore</td></tr> <tr><td>963.</td><td>Consequent</td></tr> <tr><td>964.</td><td>Contract</td></tr> <tr><td>965.</td><td>Consequentes</td></tr> <tr><td>966.</td><td>Conservation</td></tr> <tr><td>967.</td><td>Contribution</td></tr> <tr><td>968.</td><td>Converra</td></tr> <tr><td>969.</td><td>Contribution</td></tr> <tr><td>970.</td><td>Contribution</td></tr> <tr><td>971.</td><td>Contributio</td></tr> <tr><td>972.</td><td>Contribuera</td></tr> <tr><td>973.</td><td>Contributoria</td></tr> <tr><td>974.</td><td>Convenant</td></tr> <tr><td>975.</td><td>Constituta</td></tr> <tr><td>976.</td><td>Constitutio</td></tr> <tr><td>977.</td><td>Convoys</td></tr> </tbody> </table>	Number	Word	951.	Bafascon	952.	Congratulation	953.	Congratulera	954.	Congress	955.	Conjunction	956.	Content	957.	Contentam	958.	Contentam	959.	Convere	960.	Conivent	961.	Convens	962.	Conventore	963.	Consequent	964.	Contract	965.	Consequentes	966.	Conservation	967.	Contribution	968.	Converra	969.	Contribution	970.	Contribution	971.	Contributio	972.	Contribuera	973.	Contributoria	974.	Convenant	975.	Constituta	976.	Constitutio	977.	Convoys	<p><b>Neither in Order!</b></p> <table border="1"> <thead> <tr> <th>Number</th> <th>Word</th> </tr> </thead> <tbody> <tr><td>175</td><td>Foreign</td></tr> <tr><td>901</td><td>Verchader</td></tr> <tr><td>328</td><td>Verchaderen</td></tr> <tr><td>09</td><td>Verleben</td></tr> <tr><td>343</td><td>Verleben</td></tr> <tr><td>369</td><td>Verdere</td></tr> <tr><td>376</td><td>Versterken</td></tr> <tr><td>347</td><td>Vie Janers</td></tr> <tr><td>583</td><td>Versterke</td></tr> <tr><td>338</td><td>Versterkings</td></tr> <tr><td>323</td><td>Versterkende</td></tr> <tr><td>768</td><td>Versterkend</td></tr> <tr><td>706</td><td>Versterke</td></tr> <tr><td>194</td><td>Vellos</td></tr> <tr><td>936</td><td>Velhos</td></tr> <tr><td>762</td><td>Verstecasium</td></tr> <tr><td>535</td><td>Vahite</td></tr> <tr><td>301</td><td>+ Honshu</td></tr> <tr><td>187</td><td>Varshava</td></tr> <tr><td>344</td><td>Varshonewa</td></tr> <tr><td>365</td><td>Japan</td></tr> <tr><td>304</td><td>Yokohama</td></tr> <tr><td>361</td><td>Yungasati</td></tr> <tr><td>983</td><td>Yuland sea</td></tr> <tr><td>729</td><td>China</td></tr> <tr><td>792</td><td>Shanghai</td></tr> <tr><td>669</td><td>10.</td></tr> <tr><td>502</td><td>Sincomata</td></tr> <tr><td>528</td><td>Spal de Sula</td></tr> <tr><td>904</td><td>Tolomeo</td></tr> <tr><td>147</td><td>Ungar</td></tr> <tr><td>454</td><td>Ungary</td></tr> <tr><td>545</td><td>Ungary</td></tr> <tr><td>778</td><td>Ungary</td></tr> <tr><td>773</td><td>Ungary</td></tr> <tr><td>664</td><td>Ungary</td></tr> <tr><td>329</td><td>Ungary</td></tr> <tr><td>314</td><td>Ungary</td></tr> <tr><td>357</td><td>Almeria</td></tr> <tr><td>303</td><td>Marta</td></tr> <tr><td>539</td><td>Tunis</td></tr> <tr><td>501</td><td>Sebartar</td></tr> <tr><td>589</td><td>Sorge</td></tr> <tr><td>196</td><td>Barwmark</td></tr> <tr><td>733</td><td>England</td></tr> <tr><td>705</td><td>Italiens</td></tr> <tr><td>378</td><td>United States</td></tr> <tr><td>736</td><td>Europa</td></tr> <tr><td>305</td><td>Ungary</td></tr> <tr><td>326</td><td>Ungary</td></tr> <tr><td>582</td><td>Turkiet</td></tr> <tr><td>565</td><td>Egypten</td></tr> <tr><td>102</td><td>MUR</td></tr> </tbody> </table>	Number	Word	175	Foreign	901	Verchader	328	Verchaderen	09	Verleben	343	Verleben	369	Verdere	376	Versterken	347	Vie Janers	583	Versterke	338	Versterkings	323	Versterkende	768	Versterkend	706	Versterke	194	Vellos	936	Velhos	762	Verstecasium	535	Vahite	301	+ Honshu	187	Varshava	344	Varshonewa	365	Japan	304	Yokohama	361	Yungasati	983	Yuland sea	729	China	792	Shanghai	669	10.	502	Sincomata	528	Spal de Sula	904	Tolomeo	147	Ungar	454	Ungary	545	Ungary	778	Ungary	773	Ungary	664	Ungary	329	Ungary	314	Ungary	357	Almeria	303	Marta	539	Tunis	501	Sebartar	589	Sorge	196	Barwmark	733	England	705	Italiens	378	United States	736	Europa	305	Ungary	326	Ungary	582	Turkiet	565	Egypten	102	MUR
Number	Word																																																																																																																																																																				
951.	Bafascon																																																																																																																																																																				
952.	Congratulation																																																																																																																																																																				
953.	Congratulera																																																																																																																																																																				
954.	Congress																																																																																																																																																																				
955.	Conjunction																																																																																																																																																																				
956.	Content																																																																																																																																																																				
957.	Contentam																																																																																																																																																																				
958.	Contentam																																																																																																																																																																				
959.	Convere																																																																																																																																																																				
960.	Conivent																																																																																																																																																																				
961.	Convens																																																																																																																																																																				
962.	Conventore																																																																																																																																																																				
963.	Consequent																																																																																																																																																																				
964.	Contract																																																																																																																																																																				
965.	Consequentes																																																																																																																																																																				
966.	Conservation																																																																																																																																																																				
967.	Contribution																																																																																																																																																																				
968.	Converra																																																																																																																																																																				
969.	Contribution																																																																																																																																																																				
970.	Contribution																																																																																																																																																																				
971.	Contributio																																																																																																																																																																				
972.	Contribuera																																																																																																																																																																				
973.	Contributoria																																																																																																																																																																				
974.	Convenant																																																																																																																																																																				
975.	Constituta																																																																																																																																																																				
976.	Constitutio																																																																																																																																																																				
977.	Convoys																																																																																																																																																																				
Number	Word																																																																																																																																																																				
175	Foreign																																																																																																																																																																				
901	Verchader																																																																																																																																																																				
328	Verchaderen																																																																																																																																																																				
09	Verleben																																																																																																																																																																				
343	Verleben																																																																																																																																																																				
369	Verdere																																																																																																																																																																				
376	Versterken																																																																																																																																																																				
347	Vie Janers																																																																																																																																																																				
583	Versterke																																																																																																																																																																				
338	Versterkings																																																																																																																																																																				
323	Versterkende																																																																																																																																																																				
768	Versterkend																																																																																																																																																																				
706	Versterke																																																																																																																																																																				
194	Vellos																																																																																																																																																																				
936	Velhos																																																																																																																																																																				
762	Verstecasium																																																																																																																																																																				
535	Vahite																																																																																																																																																																				
301	+ Honshu																																																																																																																																																																				
187	Varshava																																																																																																																																																																				
344	Varshonewa																																																																																																																																																																				
365	Japan																																																																																																																																																																				
304	Yokohama																																																																																																																																																																				
361	Yungasati																																																																																																																																																																				
983	Yuland sea																																																																																																																																																																				
729	China																																																																																																																																																																				
792	Shanghai																																																																																																																																																																				
669	10.																																																																																																																																																																				
502	Sincomata																																																																																																																																																																				
528	Spal de Sula																																																																																																																																																																				
904	Tolomeo																																																																																																																																																																				
147	Ungar																																																																																																																																																																				
454	Ungary																																																																																																																																																																				
545	Ungary																																																																																																																																																																				
778	Ungary																																																																																																																																																																				
773	Ungary																																																																																																																																																																				
664	Ungary																																																																																																																																																																				
329	Ungary																																																																																																																																																																				
314	Ungary																																																																																																																																																																				
357	Almeria																																																																																																																																																																				
303	Marta																																																																																																																																																																				
539	Tunis																																																																																																																																																																				
501	Sebartar																																																																																																																																																																				
589	Sorge																																																																																																																																																																				
196	Barwmark																																																																																																																																																																				
733	England																																																																																																																																																																				
705	Italiens																																																																																																																																																																				
378	United States																																																																																																																																																																				
736	Europa																																																																																																																																																																				
305	Ungary																																																																																																																																																																				
326	Ungary																																																																																																																																																																				
582	Turkiet																																																																																																																																																																				
565	Egypten																																																																																																																																																																				
102	MUR																																																																																																																																																																				

## Word Substitution

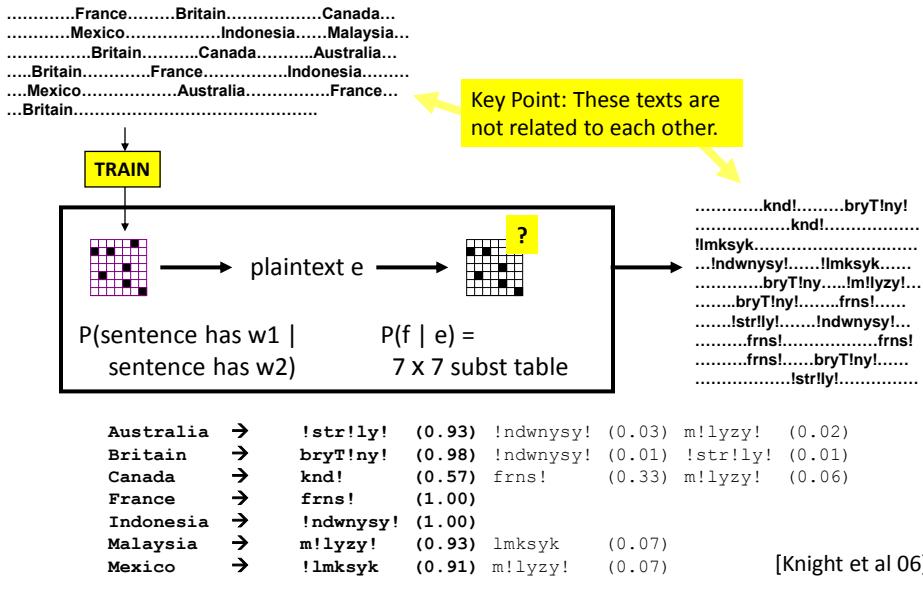
- Interesting for NLP
- Language translation can be viewed as word substitution (and transposition)
- Certainly, that is how IBM models 1-5 view it

### Word Substitution (Small-scale)



[Knight et al 06]

## Word Substitution (Small-scale)



## Word Substitution (Giga-scale)

- Suppose I replace each English word on your hard drive with some integer.
- Can you recover your texts?
- In principle, apply the same techniques we used for letter substitution.
  - English word-bigram LM drives decipherment
  - But for EM, initially-uniform substitution table is too big!
  - $100,000 \times 100,000$

## Word Substitution (Giga-scale)

- Gibbs sampling fixes memory problem

Cipher: 24234 1899 39902 5716 29948 ...

Plain: the man is car are ...

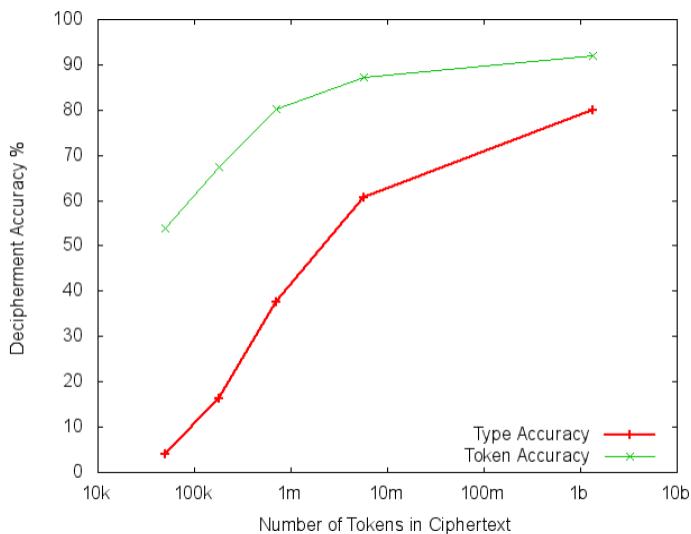
Resample:

a
an
apple
...
man
zoo

Still need to sample 100,000 alternatives at each cipher token, for each epoch.

- Slice sampling (Dou & Knight 12) fixes speed problem

## Word Substitution (Giga-scale)



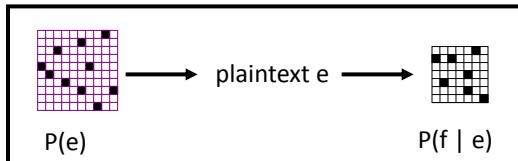
(Dou & Knight 2011)

# Foreign Language as a Cipher

"When I look at this giant corpus of Arabic, I say to myself, this is really English, but it has been encoded in some strange symbols!!! Let's decode!!!!"



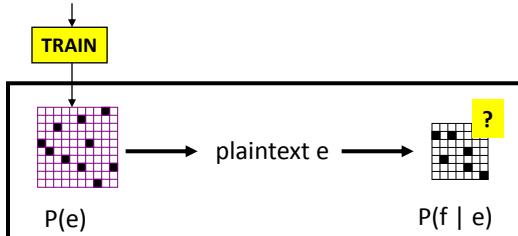
## OUR HERO



# Foreign Language as a Cipher

**BAGHDAD, Iraq (CNN) —** Six bombs killed at least 54 Iraqis and wounded 96 others Wednesday, including 20 civilians who died as they lined up to join the Iraqi army in Hawijah when a suicide bomber detonated explosives hidden under his clothing, Iraqi officials said. That attack in the town about 130 miles (200 kilometers) north of Baghdad also wounded 30 Iraqis, said Iraqi army Lt. Col. Tahar al-Zawabili. A car bombing in Saddam Hussein's ancestral homeland of Tikrit also killed 30 Iraqis and wounded another 40, Iraqi officials said. The Tiktrit explosion...

**Key Point:** These texts are not related to each other.



رسالتهم ورؤى حلقة الارض الارمني سلسلة نشرات التي  
لاقت اهتماماً عالياً بين اوساط المفكرين والعلماء في العالم العربي.  
المقدمة الافتتاحية تناولت في المقدمة الافتتاحية ملخصاً في الاحداث التأريخية  
وقد اشارت الى مفهوم سفينة على متنها مشاركاً في  
القمة العربية في الدار البيضاء التي انعقدت في 1986  
اخراج فلم "سفينة العروبة" الذي اعده وينون من اسراف  
على اتفاقيات، اضاف<sup>[15]</sup> "نحو حجج احادي وبحقون"  
هذا يشير الى سفينة العروبة التي اعادت احياء وعول قبول  
الخطاب العرقي في العالم العربي.  
[15] "في الواقع، لا يزال تجاذب  
من جانب شعب يعيش بمنطقة الانتفاضة المديدة احدى  
المناطق الاهلياء بالاسلاميين من عرق  
واكذب ان اسرائيل تدرك هذه الاعراض  
بسيلانها على اهلها، ويسهلونها على  
رفق قوي في كلمة الله، ومن ثم يصرّون وراء  
ما هي تشكيل "العناد" -بعضهم من عرق  
العناد- وهذا ينطبق على اهلنا، وهو  
مسير العساكر والجنود، وكل ذلك ينبع من افراد  
احدى الجان

!@!m		!@Swr
!ywm		=hdh! !lsh=hr
!lth!ny&		fy ywm
!@!m !lm!Dy		nys!n
Sfr	@!m 1990	!sbw@
@!m	w!lth!ny&	=hdh=h !!!'y!m
th!ny&	fy !ywm	qbl !y!m
@!m 1992	mn !lsh=hr !lj!ry	fy !!@Sr
@!m 1993	!lqrn	mn !lsn&
ywm	!y!m	!lsnw!t
!!!sbw@ !lm!Dy	@!m!aN	b@d ywm
fy !!dqyq&	!!:@&	!!!'m
!lsn& !lj!ry&	17 shb!T 1994	13 nys!n 1994
!lsn&	th!lth snw!t	!lth!ny& @ch!&
!lsh=hr !lm!Dy	dqyq&	thl!th& !y!m
!lsh=hr !lj!ry	=hdh=h !lsn&	qbl !sbw@yn
snw!t	ywmyn	fy !ywm !!!ly
sn&	mn !!@!m !lm!Dy	sh@b!n
=hdh! !!@!m	!lsn& !lmqbl&	tmwz
s!@&	fy !lsn&	3 dhw !!Hj& 1414
!!@Sr	kl ywm	fy shb!T !lm!Dy
@!m 1991	fy !!@!m !lm!Dy	qbl ywmyn

## Time Expressions

<n><n>\* ??? 19<n><n>

9 Hzyr!n 1942	27 tmwz 1993	21 Hzyr!n 1967
8 tshry!n !!!wl 1990	26 tmwz 1953	20 !y!r 1990
7 k!nwn !!!wl 1993	26 shb!T 1993	20 tshry!n !wl 1983
6 !y!r 1993	26 k!nwn !!!wl 1994	20 tshry!n !!!wl 1921
6 !~Adh!r 1991	25 !ylwl 1926	1 !y!r 1994
5 shb!T 1950	24 !~Adh!r 1993	17 Hzyr!n 1972
4 Hzyr!n 1989	22 !ylwl 1957	16 !ylwl 1919
30 !~Adh!r 1944	22 tshry!n !!!wl 1948	16 Hzyr!n 1984
29 !y!r 1945	22 tmwz 1952	16 !~Ab 1929
29 !~Adh!r 1993	21 !y!r 1994	
28 k!nwn !!!wl 1994	21 k!nwn !!!wl 1988	

## Time Expressions

<n> Hzyr!n <n>

13	4 Hzyr!n 1967	2	fy 30 Hzyr!n 1995
12	fy 12 Hzyr!n 1993	2	fy 18 Hzyr!n 1994
7	5 Hzyr!n 1967	2	fy 14 Hzyr!n 1993
6	fy 30 Hzyr!n 1989	2	fy 14 Hzyr!n 1991
6	30 Hzyr!n 1989	2	fy 12 Hzyr!n 1990
4	fy 30 Hzyr!n 1994	2	7 Hzyr!n 1994
4	fy 30 Hzyr!n 1993	2	6 Hzyr!n 1941
3	fy 19 Hzyr!n 1967	2	26 Hzyr!n 1994
2	ywm 30 Hzyr!n 1989	2	21 Hzyr!n 1994
2	w 6 Hzyr!n 1994	2	1 Hzyr!n 1994
2	qbl 5 Hzyr!n 1967	2	19 Hzyr!n 1965
2	fy 9 Hzyr!n 1967	2	18 Hzyr!n 1994
2	fy 7 Hzyr!n 1981	2	18 Hzyr!n 1940
2	fy 6 Hzyr!n 1994	2	12 Hzyr!n 1993
2	fy 5 Hzyr!n 1967	2	11 Hzyr!n 1994

## Time Expressions

<n> Hzyr!n <n>

13	4 Hzyr!n 1967	2	fy 30 Hzyr!n 1995
12	fy 12 Hzyr!n 1993	2	fy 18 Hzyr!n 1994
7	5 Hzyr!n 1967	2	fy 14 Hzyr!n 1993
6	fy 30 Hzyr!n 1989	2	fy 14 Hzyr!n 1991
6	30 Hzyr!n 1989	2	fy 12 Hzyr!n 1990
4	fy 30 Hzyr!n 1994	2	7 Hzyr!n 1994
4	fy 30 Hzyr!n 1993	2	6 Hzyr!n 1941
3	fy 19 Hzyr!n 1967	2	26 Hzyr!n 1994
2	ywm 30 Hzyr!n 1989	2	21 Hzyr!n 1994
2	w 6 Hzyr!n 1994	2	1 Hzyr!n 1994
2	qbl 5 Hzyr!n 1967	2	19 Hzyr!n 1965
2	fy 9 Hzyr!n 1967	2	18 Hzyr!n 1994
2	fy 7 Hzyr!n 1981	2	18 Hzyr!n 1940
2	fy 6 Hzyr!n 1994	2	12 Hzyr!n 1993
2	fy 5 Hzyr!n 1967	2	11 Hzyr!n 1994

## Time Expressions

<n> Hzyr!n <n>

- 13      4 Hzyr!n 1967
- 12      fy 12 Hzyr!n 1993
- 7        5 Hzyr!n 1967
- 6        fy 30 Hzyr!n 1989
- 6        30 Hzyr!n 1989
- 4        fy 30 Hzyr!n 1994
- 4        fy 30 Hzyr!n 1993
- 3        fy 19 Hzyr!n 1967
- 2        ywm 30 Hzyr!n 1989
- 2        w 6 Hzyr!n 1994
- 2        qbl 5 Hzyr!n 1967
- 2        fy 9 Hzyr!n 1967
- 2        fy 7 Hzyr!n 1981
- 2        fy 6 Hzyr!n 1994
- 2        fy 5 Hzyr!n 1967

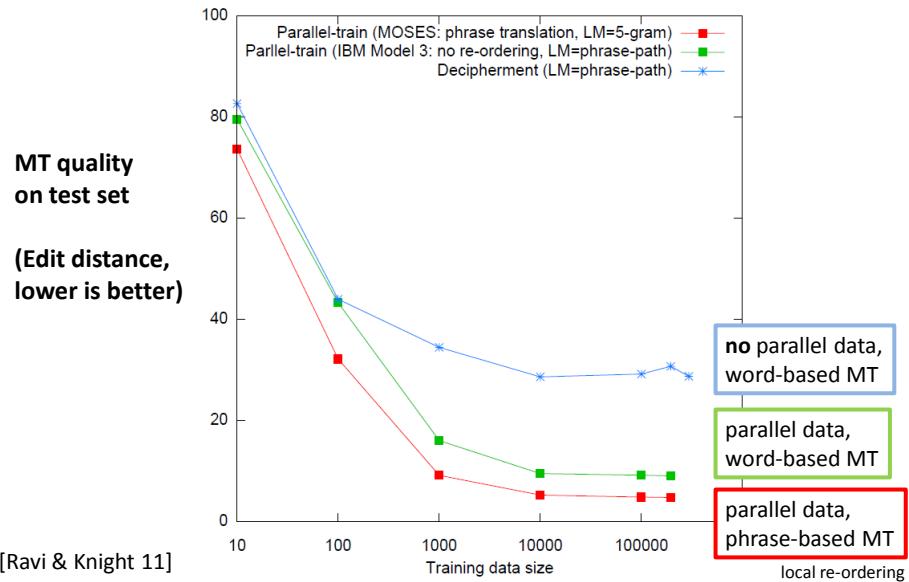
Search query	Documents
January 4, 1967	8040
February 4, 1967	9270
March 4, 1967	10700
April 4, 1967	21800
May 4, 1967	14000
June 4, 1967	39300
July 4, 1967	12600
August 4, 1967	7970
September 4, 1967	7390
October 4, 1967	8800
November 4, 1967	6560
December 4, 1967	9770

## Time Expressions

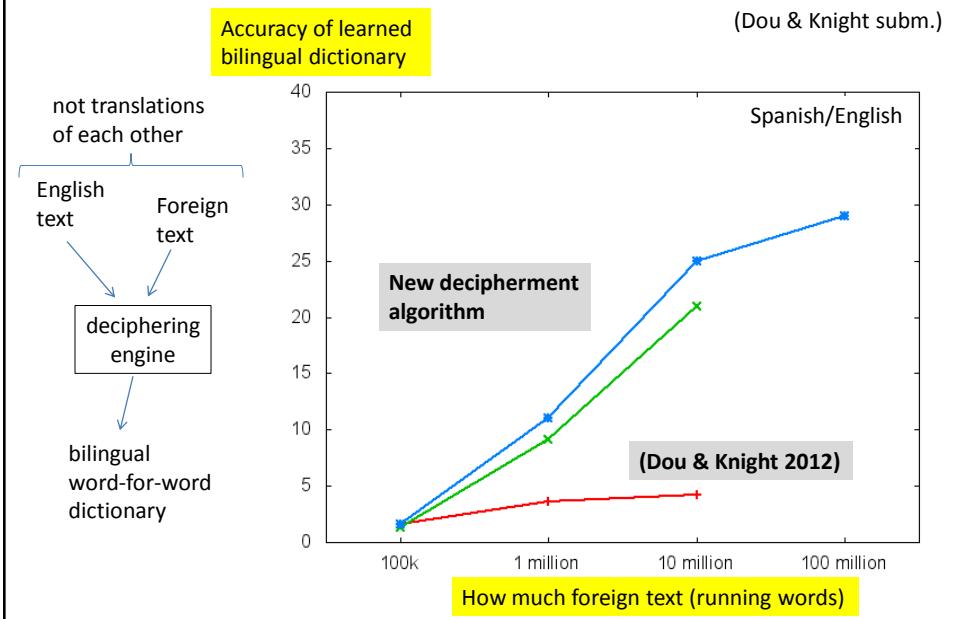
Hzyr!n

- |                                   |                                |
|-----------------------------------|--------------------------------|
| 229      fy Hzyr!n !lm!Dy         | 16      n=h!y& Hzyr!n !lm!Dy   |
| 207      fy Hzyr!n                | 16      fy Hzyr!n 1990         |
| 75        fy Hzyr!n !lmqbl        | 15      sh=hr Hzyr!n           |
| 61        fy Hzyr!n 1993          | 15      fy sh=hr Hzyr!n !lm!Dy |
| 31        fy Hzyr!n 1992          | 15      fy Hzyr!n 1994         |
| 27        !lr!b@ mn Hzyr!n        | 14      mn 17 Hzyr!n           |
| 27        fy Hzyr!n 1967          | 14      fy Hzyr!n 1996         |
| 19        fy 30 Hzyr!n !lm!Dy     | 14      fy 30 Hzyr!n           |
| 18        fy n=h!y& Hzyr!n !lm!Dy | 13      fy sh=hr Hzyr!n        |
| 18        fy Hzyr!n 1991          | 13      fy 20 Hzyr!n !lm!Dy    |
| 17        mn Hzyr!n               | 13      4 Hzyr!n 1967          |
| 17        mndh Hzyr!n !lm!Dy      | 12      n=h!y& Hzyr!n          |
| 17        4 Hzyr!n                | 12      !lr!b@ mn Hzyr!n 1967  |

## Deciphering Spanish Time Expressions

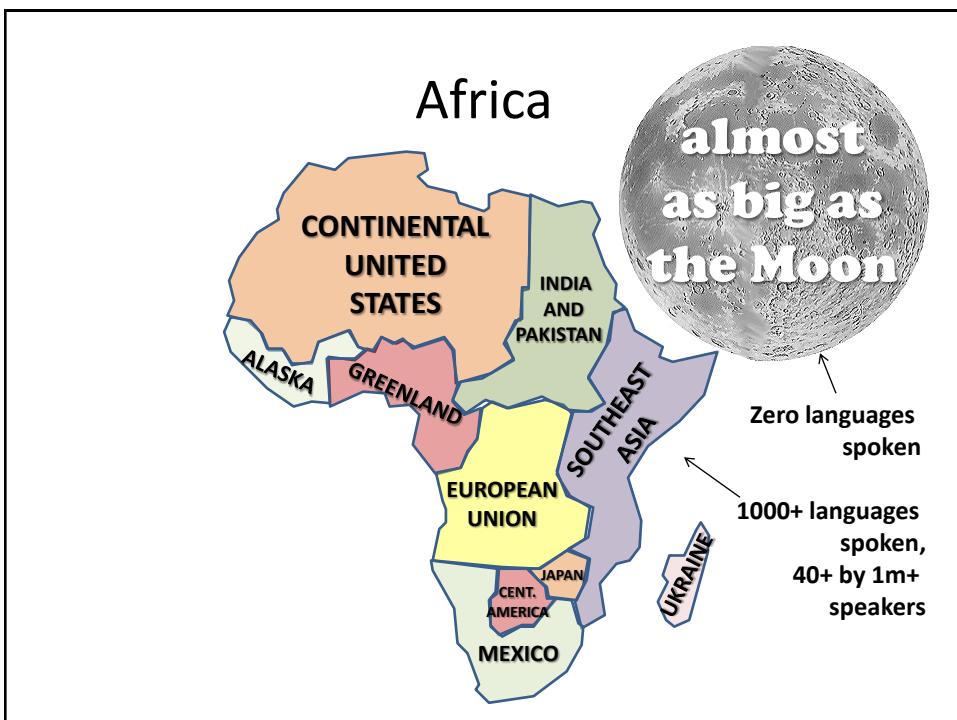


## Deciphering Foreign Language at Giga-Scale



## Practical Value

- Scenarios where in-domain parallel data is scarce.
- Decipher large monolingual in-domain corpora to improve systems trained on small amounts of parallel text



## Unsolved ciphers

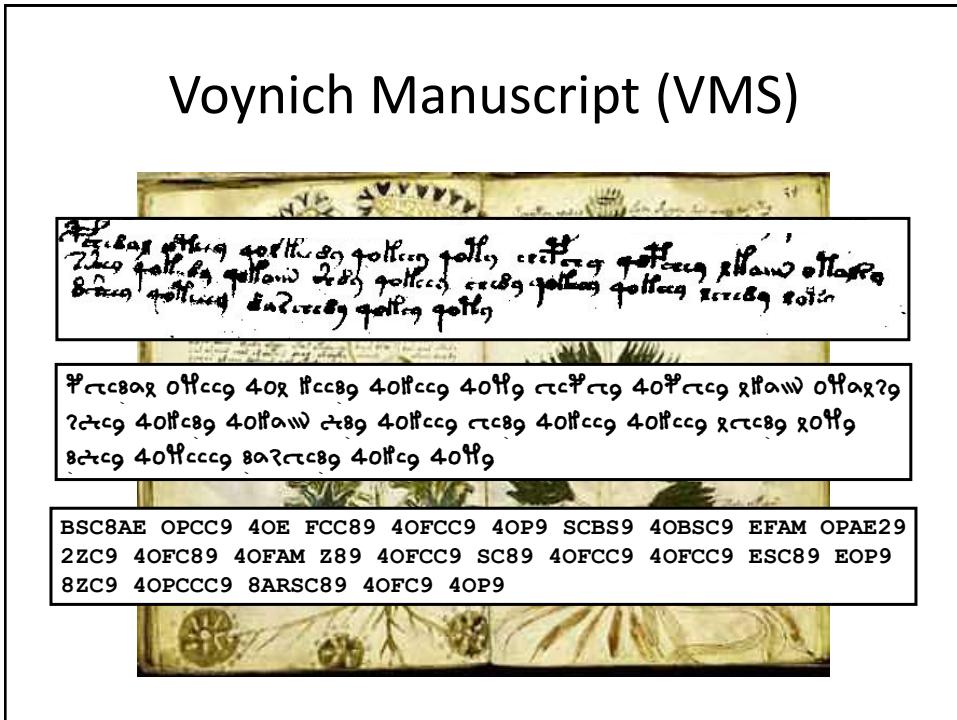
### Voynich Manuscript (VMS)



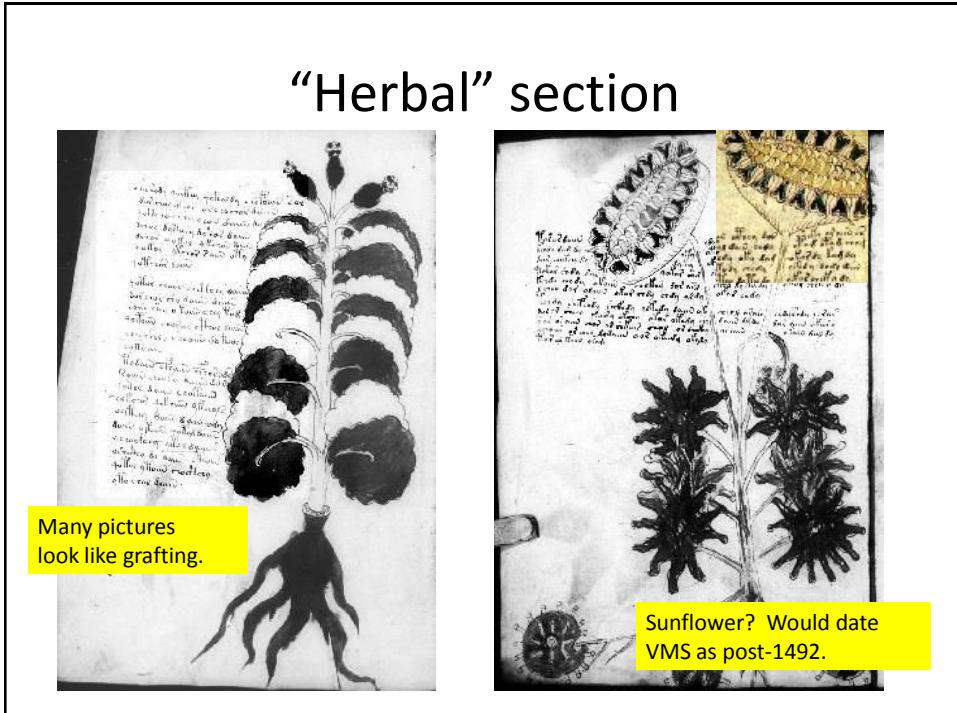
- Medieval illustrated manuscript (early 1400s)
- 235 pages, 6 sections, 38k word tokens, 35 letter types
- Undeciphered



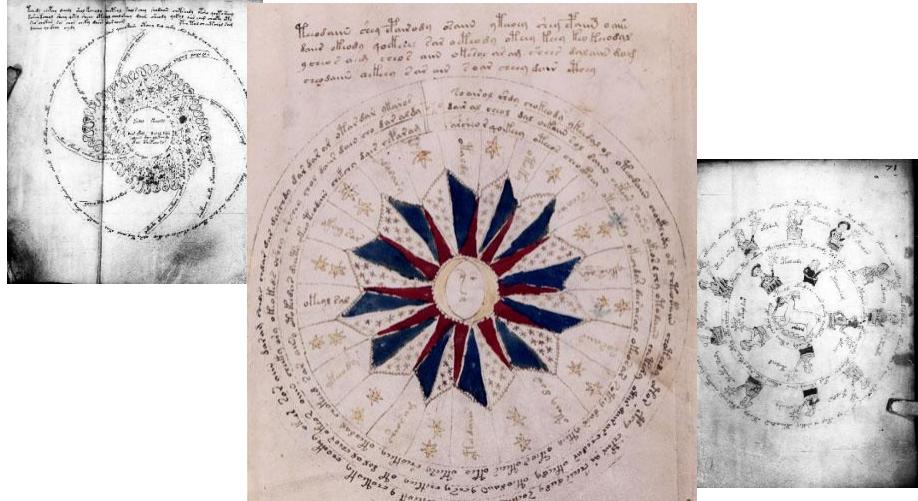
# Voynich Manuscript (VMS)



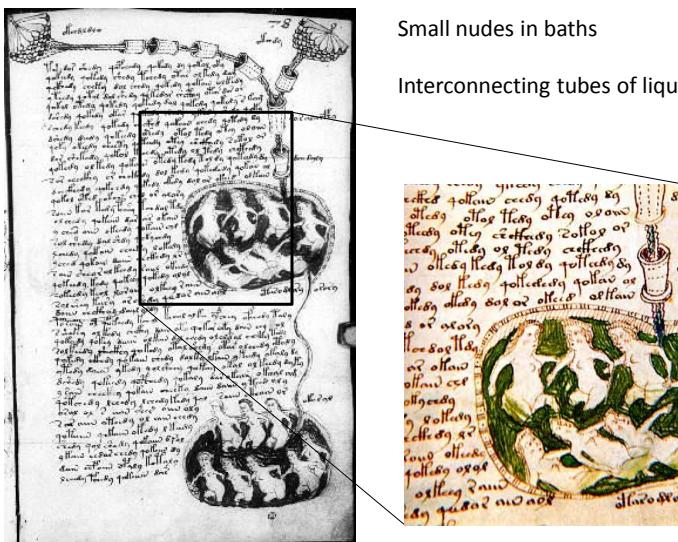
## “Herbal” section



## “Astrological” section



## “Biological” section



Small nudes in baths

Interconnecting tubes of liquids



## “Pharmacological” section

### History of Voynich Manuscript (VMS)

- |           |  |      |  |
|-----------|--|------|--|
| 1576-1612 | Rudolf II purchases VMS                            | 1864 | Ethel Boole born in England  |
| 1608-1622 | J. de Tepenecz signs VMS in Bohemian court         | 1865 | WV born in Lithuania   |
| 1630s     | George Baresch owns VMS<br>sends letter to Kircher | 1885 | WV imprisoned, Polish nationalist  |
| 1639      | GB writes Kircher again                            | 1890 | WV & EB meet, marry in 1902  |
| 16xx      | Marci inherits VMS from GB                         | 1898 | WV publishes first book list   |
| 1665      | Marci sends VMS to Kircher with letter             | 1912 | WV acquires VMS in “ancient castle”  |
| 1665-80   | Kircher owns VMS                                   | 1914 | WV moves to USA, opens bookshop  |
| 1680      | Kircher dies                                       | 1919 | WV sends photostatic copies of VMS   |
|           |  | 1919 | Copying reveals de Tepenecz signature  |
|           |  | 1919 | WV writes to Bohemian State Archvs   |
|           |  | 1921 | <b>WV presents VMS + inserted Marci letter mentioning Francis Bacon, asks \$160k</b> |
|           |  | 1921 | <b>Newbold &amp; WV announce decipherment</b>  |
|           |  | 1930 | WV dies. VMS placed in vault, \$100k   |
|           |  | 1931 | VMS appraised at \$19,400  |
|           |  | 1960 | Ethel dies, VMS to secretary Ann Nill<br>“Castle” revealed as Villa Mondragone       |
|           |  | 1961 | NY dealer Hans Kraus buys for \$24,500   |
|           |  | 1969 | Kraus donates VMS to Yale  |
|           |  | 1972 | Brumbaugh finds WV letters in BSA  |
|           |  | 200x | Zandbergen finds 1639 Baresch letter in newly online Kircher archive                 |



## Newbold Decipherment

Marci letter → Bacon → Cabala → “letter doubling” cipher

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

22x22 table

### Encoding:

A → CC, OM, ...

B → ...

### Decoding:

...

DO → N

...

N → HA, MI, DO, NU ...

...

Z → ...

Encoder has freedom to devise  
a “cover text” to hide real message.

### Example:

a n n ... → DO MI NU ... → DOMINU ...

## Newbold System

- Too hard to assemble good “cover” text!
- So, make cipher letter-pairs overlap:  
a n n ... → AD DB BR ... → ADBR ...
- Then, employ anagramming:  
a n n ... → OM DO MI ... → DO OM MI ... → DOMI ...
- Now can construct a plausible looking “cover” text in Latin for our secret message (also in Latin)
- An ingenious system, to be sure!

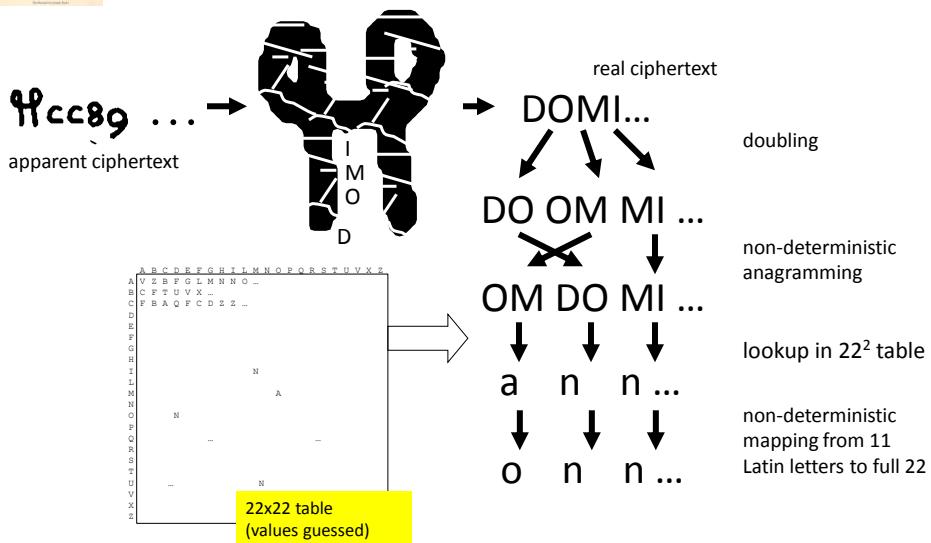
# Newbold Decipherment

Hmm, by the method, both plaintext **and ciphertext** should be in Latin letters...

But the VMS doesn't have Latin letters...



## Let's Decipher with Newbold !





# Newbold's Results

1300 real ciphertext “letters” in first 3 lines

Decipherment of those first lines:  
“I, Roger Bacon, have written this...”  
(in Latin)

Anagramming sets of 55 letters is sometimes required.

Slow but steady progress... Andromeda galaxy, ovaries ... so ... Roger Bacon must have had a microscope & telescope, hundreds of years before they were invented ... !

# VMS Transcription

Філософія 09:00-10:30  
Історія 10:45-12:00  
Література 12:00-13:30  
Математика 13:30-15:00  
Підготовка до вступу 15:00-16:00

BSC8AE OPCC9 4OE FCC89 40FCC9 4OP9 SCBS9 4OBSC9 EFAM OPAAE29  
2ZC9 40FC89 40FAM Z89 40FCC9 SC89 40FCC9 4OFCC9 ESC89 EOP9  
8ZC9 4OPCCC9 8ARSC89 40FC9 4OP9

*last paragraph, f103r*

## Alphabet: Currier/D'Imperio Transcription

ც ც ც  
C S Z

ბ ბ ბ ბ  
P F B V

ტ ტ ტ ტ  
Q X W Y

ჯ ა გ რ ი ძ  
J A E R O I D

გ ჸ ჸ გ 4 ?  
6 7 8 9 4 2

ღ ღ ღ  
G H 1

რ რ რ  
T U O

ლ ლ ლ  
N M 3

ლ ლ ლ  
K L 5

## VMS Letters

count letter

25468 O 0  
20227 C c  
17655 9 9  
14281 A ა  
12973 8 8  
11008 S ს  
10471 E ე  
10026 F ֆ  
6716 R რ  
5994 P պ  
5423 4 4  
4501 Z շ  
4076 M մ

count letter

2886 2 ?  
1752 N ნ  
1413 B բ  
1046 J յ  
950 Q პ  
908 X ხ  
591 T ტ  
524 \* \*

count letter

148 U უ  
96 6 ჭ  
74 Y ყ  
52 K კ  
31 G ღ  
17 L ლ  
14 H հ  
2 1 մ  
1 5 լ  
1 0 մ

Total  
63k character tokens

## VMS Words

count	word	count	word	count	word
863	8AM	8AM	212	OFAM	0FAM
537	OE	OE	211	8AN	8AN
501	SC89	SC89	191	4OFAE	40F0R
469	AM	AM	186	ZOE	ZOR
426	ZC89	ZC89	177	OFC9	0FCC9
396	SOE	SOE	174	SCC9	SCC9
363	OR	OR	172	SCOE	SC0R
350	AR	AR	155	S9	CT9
344	SC9	SC9	155	OPC89	0FCC9
318	8AR	8AR	154	OPAM	0FAM
308	4OFC9	4OFC9	152	4OFAR	40F0R
305	4OFC89	4OFC89	151	9	9
283	ZC9	ZC9	151	4OE	40R
279	4OFAN	4OFAN	150	S89	CT89
272	4OFC89	4OFC89	147	4OF9	40F9
270	89	89	144	ZCC9	CTCC9
262	4OFAF	4OFAF	144	OFAN	0FAM
260	AE	AE	144	2AM	2AM
253	8AE	8AE	143	OPAE	0F0R
243	2	2	141	OPAR	0F0R
219	SOR	SOR	140	SX9	CTDK9
					+ many more!
					Total: 8116 distinct words

## VMS Word Bigrams

- Very few repeated bigrams: **Extremely troubling!**  
Nothing like “of the” in English.
- 115 (out of 8116) distinct words appear doubled  
... 40FCC89 40FCC89 ...
- 8 distinct words appear tripled
  - ... 40FCC89 40FCC89 40FCC89 ...
  - ... CT0R CT0R CT0R ...
  - ... CT0R CT0R CT0R ...
  - ... 0FAM 0FAM 0FAM ...
  - ... OR OR OR ...
  - ... 9FAM 9FAM 9FAM ...
  - ... 8AM 8AM 8AM ...
  - ... 40FCC89 40FCC89 40FCC89 ...

## Substitution Cipher?

- Nope.
- Tried 80+ languages.
- For example, if we decipher assuming Latin plaintext:

quiss squm is onum pom  
quss hates s qum hatis ...



- Tried 80+ languages written without vowels.

## Letter Clustering

Trigram model over {a, b, \_ }

a a \_ b a b \_ a b a a \_ ...

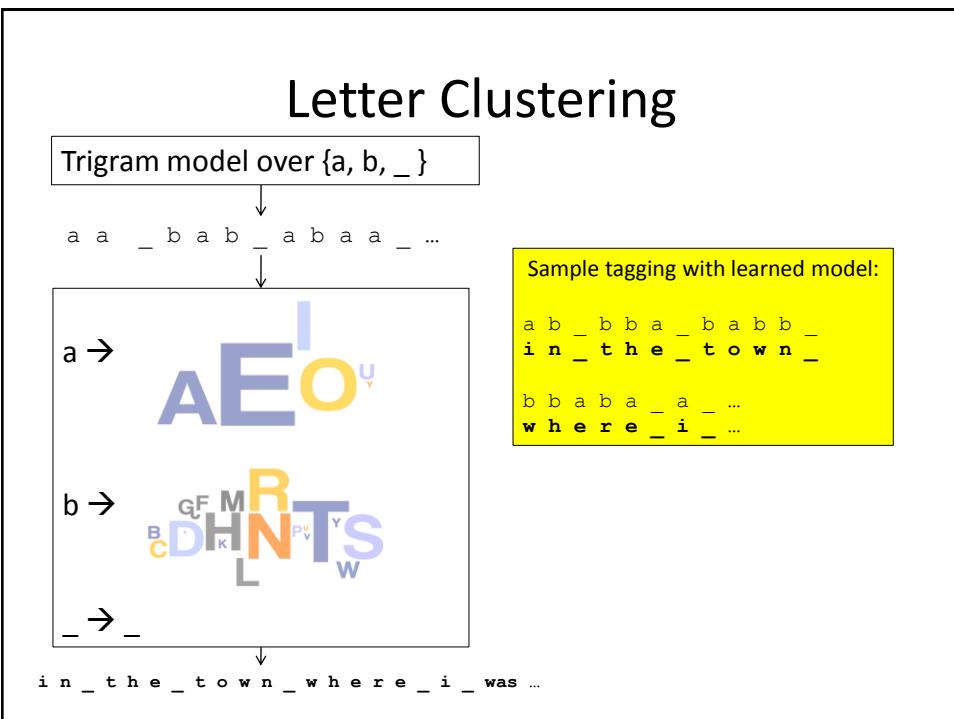
a → AEI

b → BDHNTS

\_ → \_

Sample tagging with learned model:

a b \_ b b a \_ b a b b \_  
i n \_ t h e \_ t o w n \_  
b b a b a \_ a \_ ...  
w h e r e \_ i \_ ...



# Letter Clustering

Trigram model over {a, b, \_}

a a \_ b a b \_ a b a a \_ ...

a → {all Voynich letters}

b → {all Voynich letters}

\_ → \_

Sample tagging with learned model:

? ? ? ? ? \_ ? ? ? ? ? \_ ? ? ?  
V A S 9 2 \_ 9 F A E \_ A R \_

? ? ? ? ? \_ ? ? ? ? ? \_ ? ? ? ?  
A P A M \_ Z O E \_ Z O R 9 \_ ...

# Letter Clustering

Trigram model over {a, b, \_}

a a \_ b a b \_ a b a a \_ ...

a → 

b → 

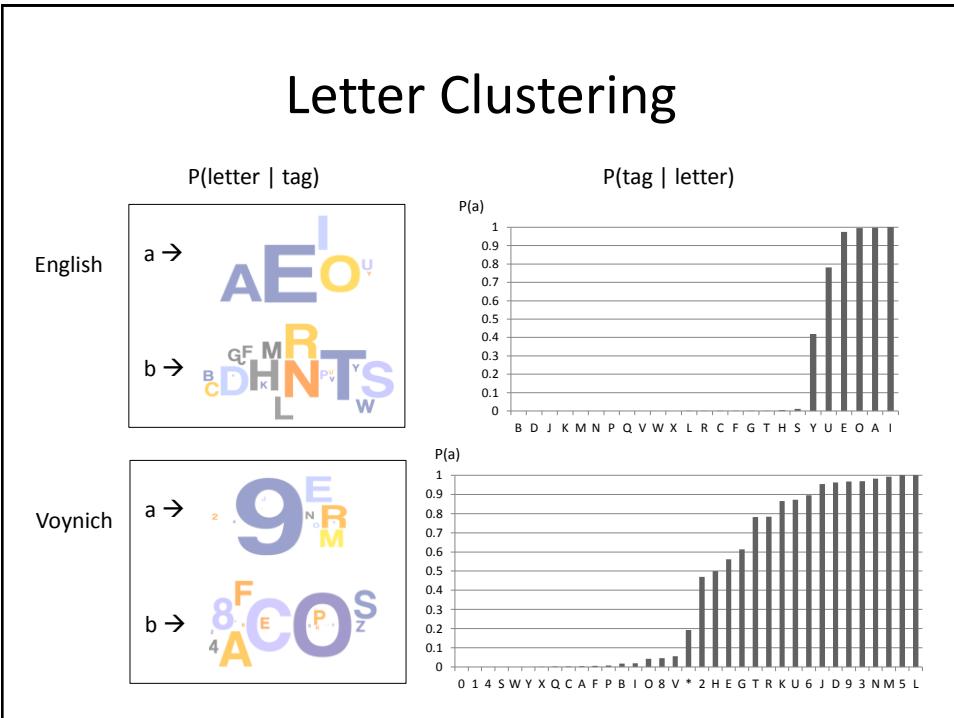
\_ → \_

Sample tagging with learned model:

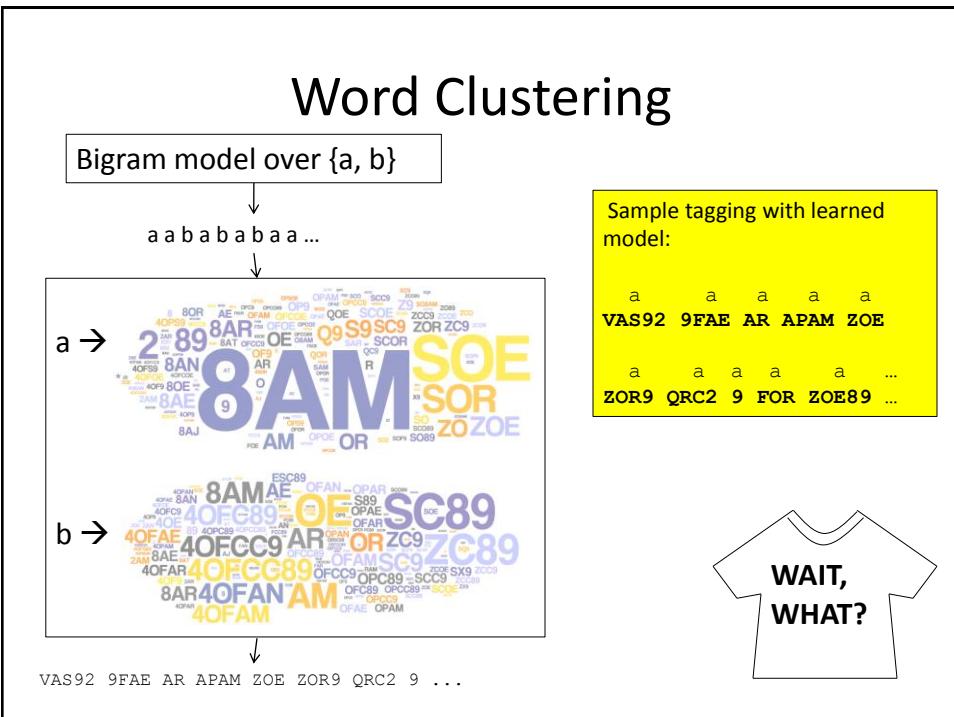
b b b b a \_ a b b a \_ b a \_  
V A S 9 2 \_ 9 F A E \_ A R \_

b b b a \_ b b a \_ b b b a \_ ...  
A P A M \_ Z O E \_ Z O R 9 \_ ...

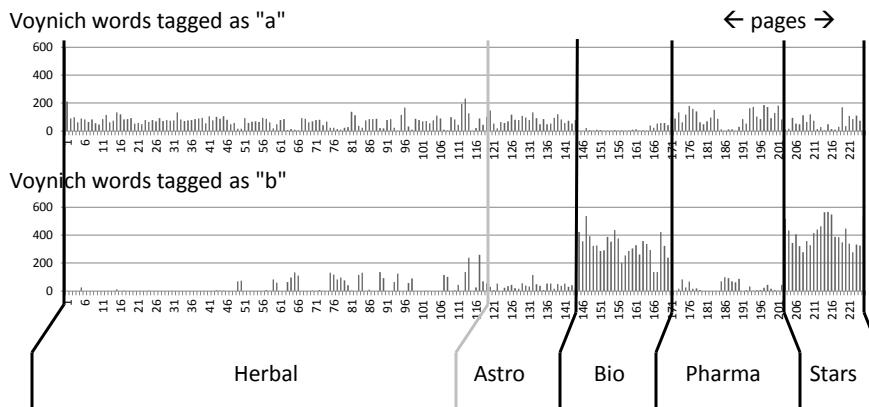
## Letter Clustering



## Word Clustering



# Word Clustering



Voynich sections, per drawings observed.  
Captain Currier's "two languages" (1976).

## An Application of PTAH to the Voynich Manuscript (U)

Approved for Release by NSA or  
06-03-2009, FOIA Case # 58742

BY MARY E. D'IMPERIO

~~-Top Secret Umbra~~

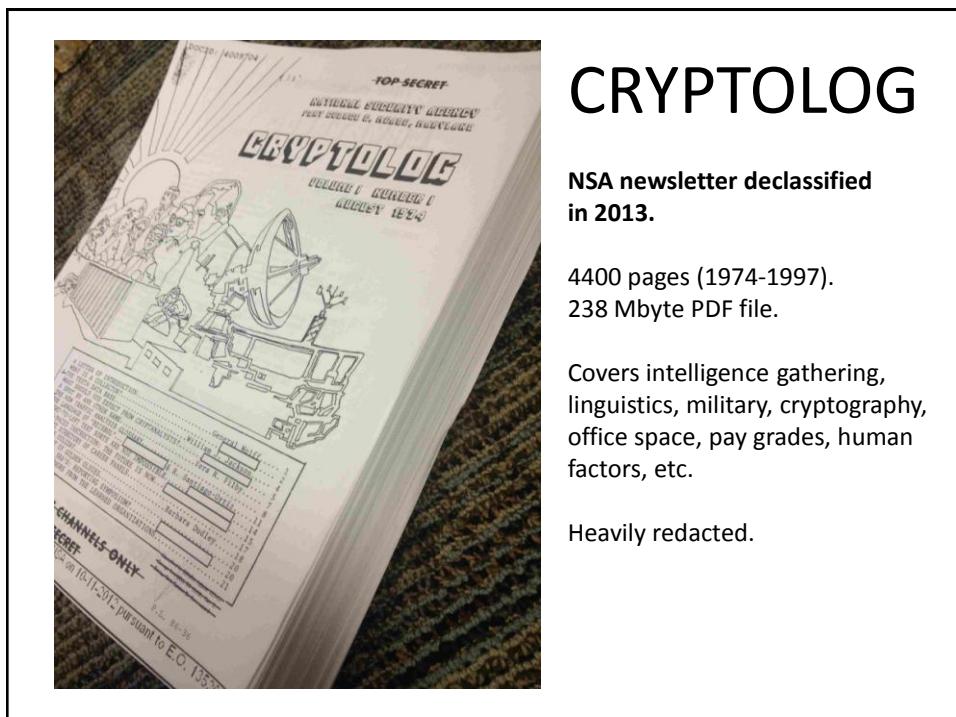
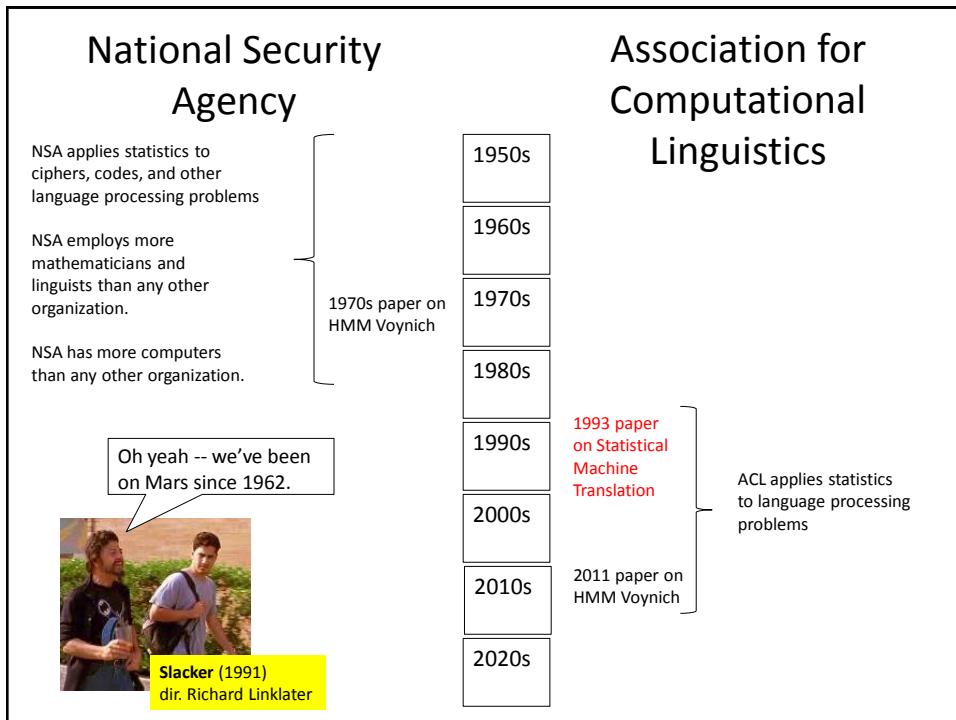
(U) This article is the second in a series of studies applying some modern statistical techniques to the problems posed by the Voynich manuscript. This study attempts to discover and demonstrate regularities of patterning in the Voynich text subjectively noted by many earlier students of the manuscript. Three separate PTAH studies are described, attacking the Voynich text at three levels: single symbols, whole "words," and a carefully chosen set of substrings within "words." These analyses are applied to samples of text from the "Biological B" section of the manuscript, in Currier's transcription. A brief general characterization of PTAH is provided, with an explanation of how it is used in the present application.

# 1970s National Security Agency report recently declassified!

program. He was struck by the passage "immenso Ptah noi invociam," and named his program after the Egyptian god. The name was ultimately extended from this program, implementing a particular application of the method, to the method and its mathematical theory as well [2, p 85]. According to [ ] of RSI, the name is pronounced "however you like" [8]. The technique itself and its uses are classified Top Secret Codeword.

Digitized by srujanika@gmail.com

I chose PTAH for the present study for two main reasons: first, because of the applications of PTAH to book codes, and second, because I wished to learn more about PTAH itself.



# CRYPTOLOG: Voynich

DOCID: 4009723 UNCLASSIFIED

**The Voynich Manuscript**

When a newspaper editor needs a filler, he always falls back on the Toch Ness Monster or the Abominable Snowman. By the editor of a cryptologic magazine the obvious device is another blurb on the subject here discussed. So, evidently, thought a former editor, among whose efforts the following paragraphs were found:

Is the Voynich manuscript "real"? No. Is it a hoax? No. What is it then? A make-believe language, ingeniously produced purely for the satisfaction of the maker.

That was my reaction the first time I looked at it closely, but faced with all the profound theories about it I lacked the courage to say so. However, a recent rereading of Elizabeth Foss' article in *Scientific American* (August 5, 1962) and of Brigadier Tilman's paper in the NSA Technical Journal (Summer 1967), plus some phenomena I have seen in the meantime, have emboldened me to give the world the benefit of my thoughts.

## The Voynich Manuscript Revisited

P16

**FOCUS** An example of all of these problems is the Voynich manuscript, a unique European manuscript thought to date most probably from the 15th or 16th century, which has resisted solution, not only by philologists early in this century, but by NSA cryptanalysts as well.

The Voynich Manuscript, an object of interest off and on since the seventeenth century, contains over 200 pages written in a partially deciphered script which is completely incomprehensible. Equally enigmatic are the large number of drawings -- of plants, few of which are identifiable, and of naked women sitting in tubs or drawing water from pipes (one author has called the latter a "plumber's nightmare").

The history of the manuscript, which has been detailed in other places, needs only passing mention since it does not add any light on the question of its origin. About 1500, it was given to Joannes Marci, mathematician and orientalist at the University of Prague, to have belonged at that time to the Emperor Maximilian I (1519-1521). Marci writes in 1666 to the Bishop of Altötting, Kircher, in Rome, that he was making a present to the latter of the manuscript, the author of which he could not identify. But before Kircher, the great medieval scholar Roger Bacon, (how Marci came into possession of it, I do not know.)

Marci himself withheld judgment on the attribution, but at least one scholar since his time has come to the same conclusion. But the authorship, Professor William Newbold of the University of Pennsylvania was convinced that it

was an enciphered text prepared by Bacon and he worked on this assumption from 1919 until his death in 1948. Although he never deciphered some of it, including an occurrence of "eg. Bacon" on the last page\*. His "solution" has been convincingly refuted by other scholars, who however have not offered anything better.

I now rush in where angels fear to tread. Although I am no specialist in Old Norse, I am convinced that the manuscript is a text in fourteenth century Danish or Norwegian -- perhaps a mixture of the two. This suggestion has also been suggested. For reasons too complex to go into here, I have tentatively identified Old Norse (that is, Old Swedish) and rejected altogether the second choice of Old West Norse, Old Icelandic. The reasoning which suggested Dano-Norwegian is given below.

Most of the manuscript has a depressing number of intercrossed words and phrases, of little help unless collateral information is available, suggesting that these are prayers, incantations, or formulas of a specific character. This is

\*The information in this paragraph and the preceding paragraph was taken from *Herculanum*, January 1963 (Vol. V, No. 3). (UNCLASSIFIED)

April 76 • CRYPTOLOG • Page 11

# CRYPTOLOG: Machine Translation

is machine translation. Machine translation is actually having a computer prepare a translation. There was to have been no difference in quality or style between a translation done by a machine and one done by a person. Georgetown University was very active in the field for some time. Progress was as easy and rapid as had been anticipated. However, in 1966 the Automatic Language Processing Advisory Committee published a report recommending that research along machine-translation lines be cut back. This report sharply curtailed federal funding. There is still, however, research being done both here and abroad, and there are several machine-translation systems that claim to be operational. One is the METEO project in Canada, which developed a system that translates weather reports from English into French. CULT (Chinese University Language Translator) in Hong Kong translates two periodicals into English. And a system was developed by a U.S. company for FID and was adapted for use by NASA during the Apollo-Soyuz Test Project. These systems differ a great deal in their approach and in the amount of pre-editing and postediting that is necessary, but all are true machine-translation efforts.

At present, NSA has a rather limited machine-translation effort.

As machine translation stands today, we haven't reached the stage where we can feed a "source" (foreign-language) text into a computer and produce a text in the "target" (in our case, English) language which is as good as the human product, not without extensive pre-editing or postediting. But in the science and technology world, current machine translation has a place. Some scientists prefer it to the

**Partial Machine Translation: A Final Report U**

P16  
and  
P16

DOCID: 4010113 TOP SECRET UMBRA CRYPTOLOG

MACHINE TRANSLATION:  
What can it do for us?

DOCID: 4010113 TOP SECRET UMBRA CRYPTOLOG

ED 1-4, (1)  
PL 86-36

# CRYPTOLOG: Evaluating Translations

## An Objective Approach to SCORING TRANSLATIONS

Reprinted from *QRL (Quarterly Review for Linguists)*, November 1973

*Author's note: The philosophy underlying the translation grading system described in this paper has been developed and applied by Emery Tetrault and myself, with many valuable suggestions from our colleagues on Professional Qualification Examination (PQE) Committee, and from other Agency linguists. My use of the pronoun "we" reflects this collaboration. I personally take full responsibility for presenting our findings here.*

Translation as an intellectual activity has been practiced since antiquity for practical as

tuitive judgments across lang in source language-to-English

Over the past 2 or 3 years I have developed a way to sco which may obviate this proble tent even though our results been far from perfect (total grading any kind of connected impossible). Our first large system, which I will describe the Russian PQE. We have sub in a number of other PQEs in languages, mainly Indo-European other families. The results aging enough in both instance mend its use in the PQE Handb

# CRYPTOLOG: Linguists

## LET'S GIVE LINGUISTS A BIGGER PIECE OF THE PIE!

### • Recognition

Most linguists specified that they want recognition above all else. A number felt that lack of recognition of the worth of linguists is evident in the inability of Agency linguists to compete successfully with managers or others for promotion. Despite almost unanimous complaints about lack of recognition, few specific suggestions were made regarding how that recogni

### TEACHING COMPUTER SCIENCE TO LINGUISTS

by [redacted] P16

### 12. PUBLICATIONS (List titles; do not confuse this with reports prepared as a regular part of the job)

## SOME TIPS ON GETTING PROMOTED

Article based on talk given in April 1978 to WIN (Women in NSA)

Promotion. The word inevitably stirs response of some kind in every red-blooded NSA employee: hope, pleasure, challenge; despair, frustration, disappointment; even inertia, resentment, resignation. Despite disparate views on promotion,

serving on the Agency Grade 14 my experience there has simply held impressions and reinforced the critical importance of the covered in this article.

*Personnel Summaries*

[redacted]

cou  
sep

ling

sys

see

inc

dog

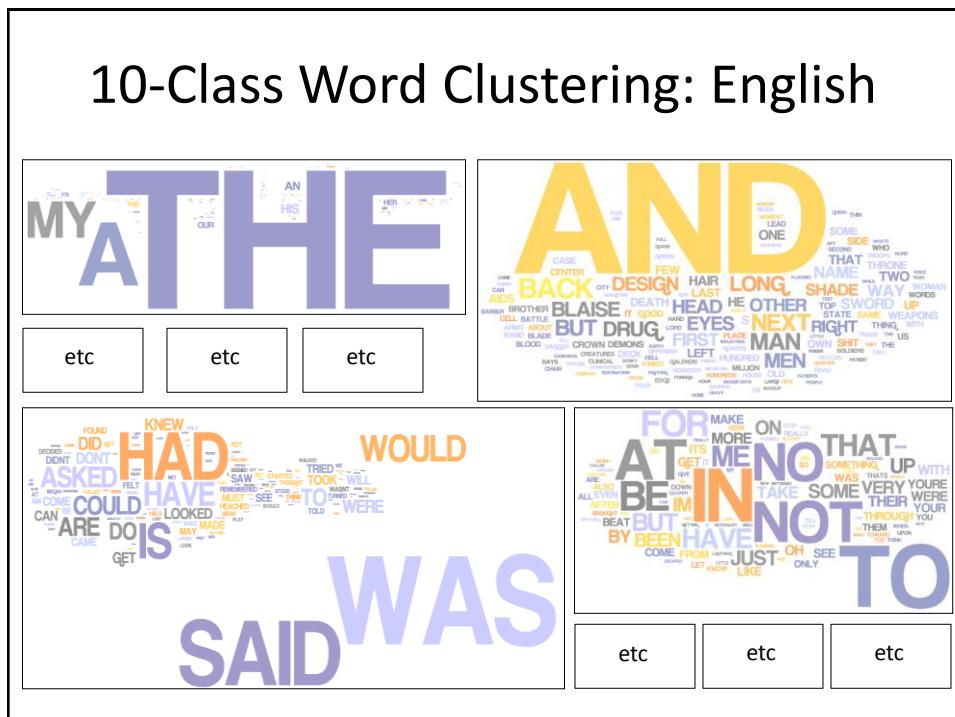
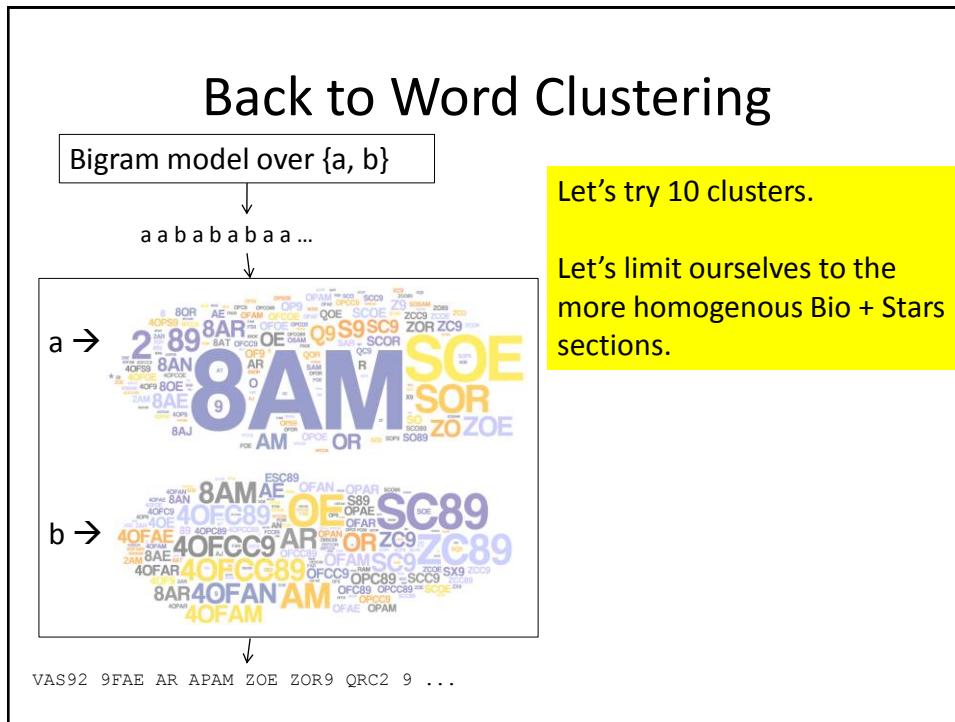
ord

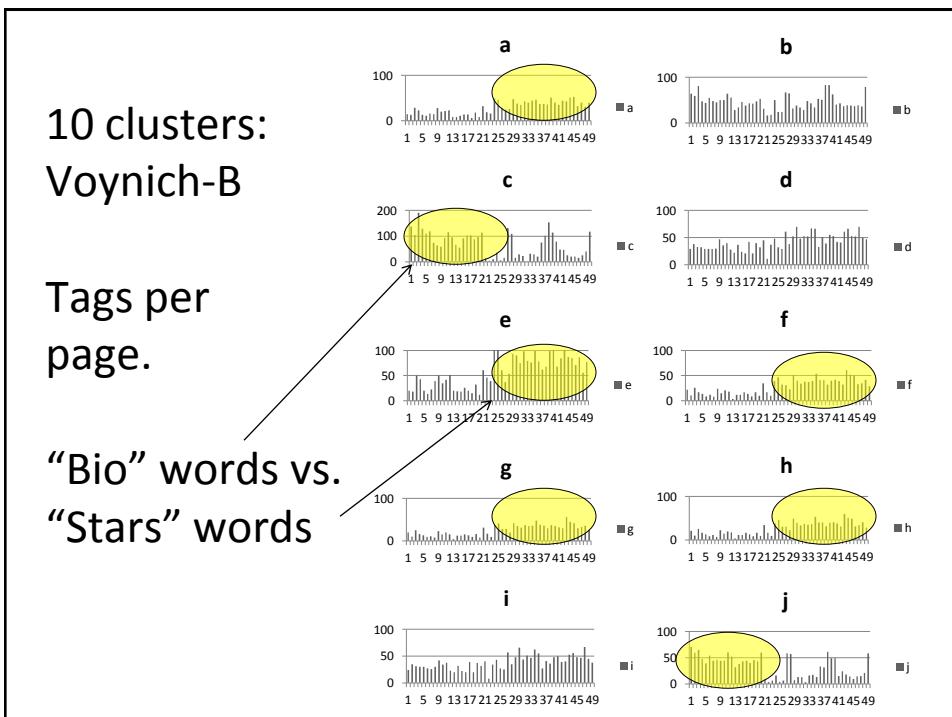
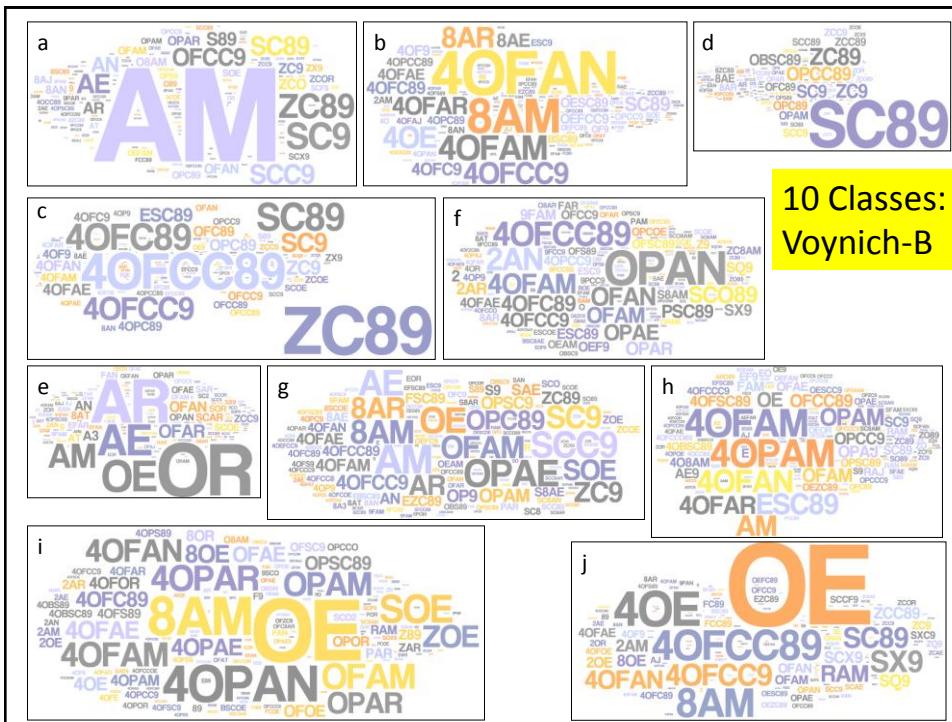
axis

pro

fai

ext





# Does VMS Have Content Words?

Measure the saliency of a word in a page  
with TF-IDF

$$\text{TF-IDF}(w, d) = \text{TF}(w, d) \times \log \frac{N}{\text{DF}(w)}$$

# times that word w  
occurs in page d

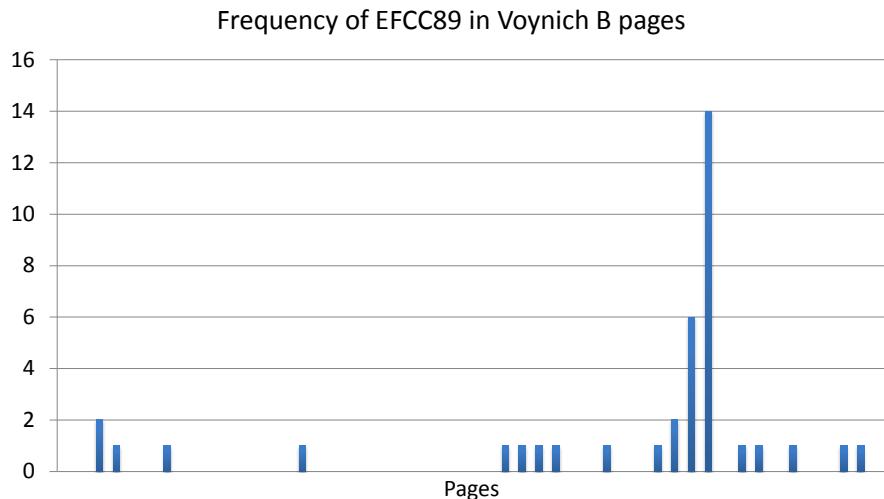
# pages that  
contain word w

(Reddy & Knight, 2011)

# Does VMS Have Content Words?

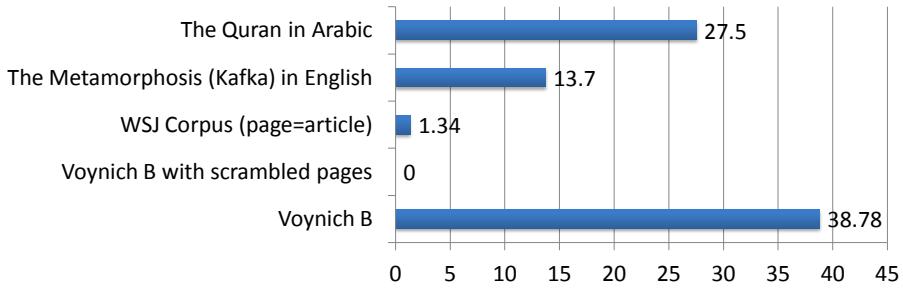
OFCC9 ZC89 8AN ZC9 SCFC89 R0R 4OFAN SAE FAN ZC09 EPAN OR 8FC09 EZC9 FAE 4OFCC89 FAR SC9 R AN ZC89 40FCC9 4OF9 FCCOR SAR OFC9 SQ9 BSC89 SX9 EZC89  
QFC89 40FC89 QESC89 2AE S89 OPC9 ESC89 40FAE DEAR FCC89 2AN OPAR 2AM 0PZC89 BSC8AR BAE DEAN SAR AR ESC9 EDR AJ 40FCC9 SCQ9 ZX9 0PC89 8AT ZC9 E  
SFAE 40FC9 8CF9 SQC89 EPC89 0FAE 0P9 40FC89 OPCCC89 4OPC9 40FC89 DFC89 0PAJ 2ZC89 2 O 8AM EFAR 8E E889 FC9 SDC89 8AR9 ZDOR 40FCOR AM ZD9I SE FCC9  
4OPAN 4OPB89 OPC9 8X9 EF C89 40FC00 AE SC8AN **EFCC89** ZC89 SOE OEAZ ZC8AE 4OFAM 9 8CDE EFC9 4OBSC9 OFS9 **SCAE** AEOE ZCAR 8RA  
40XC9 RAI OFC8AE SC0C9 40PC2 ZC8AJ OFCCC89 40XC89 SCX89 0P0DE 2AEFCC89 PC89 SX89 **PCC9** ZCFC89 ZCC0E ZCAE ON 40FCAR SR 8AJ  
EAJ 8ZCC89 8SC0 ZAR SQ9 E0FCC89 ER SXAE 4OPC89 ESCOE SC8AR IRAN 40PSC9 AT OIF\*9 40SC9 DEFAE 8ZC9 PAR EFO **EFCC9** ZC0 AEOJ FCC0E  
40FSC89 8AN ZFAE **8SC8** 40VSC89 S89 40F0DE ESC8 SFCC9 8SC8 E0FAJ 40F OFCC0E SCBSC8 80M ZCX89 PC9 0PCBAR OPC2 40FC89 OPC089  
SC089 40FC089 OPC89 AK EV89 SAM PAT CCC2 OJ EFAJ **FCCC9** SC02 8AK ZCCF E0P9 ZCFAN 40FCC89 OFCAE EFCC089 ZC08 ZC8AN BSAE OPCAE  
OPCC0I 40AN 40B 8C0 SC0 SCFC89 AEAJ **OFCCO** EFCS89 EFC889 SCCFAN OPCC89 EFAT E8AE OPCCOE SO FCS89 0EFC889 40FC08  
40SC89 OPCO ZCC08AR SC0J OPARAE OAM OEFCC0 EFC0E EFAE EFCC9 PC8AJ OFC8AN ZOF **40FCC0E** 40FCC02 0PCC8AN EFCC89 8CAJ  
SCX89 8SC9 FSOES8AR A1B SCFAE89 4CCAE 40FCCA2 SC0F89 ZCCFZ9 SOFSC9 SX9 9SCC08AN 9FCC8M RCCC9 CEA3 AF9 ZFAM 8ZCC0 SCFC89  
OPCOEAT 40ZC0 8SC8AR 9FCC02 BSC8C9 OPC080 BS8AJ OFC8AN ZCP ZCOPA1 ZPAR OESCO89 ZC0C9 ESCCCF AJ 40FCC080R SCC89 40PC8E  
EOC89 SC8C9 EFAK O\*OR F98CC89 ZCCP SCOP889 2ZC0 PSC8 FCQCC89 40FCZ8 FCCZ089 **OFCC89** ESR 20 AAE8 O'AR 20AM OPC08AM  
40PC88AM PCC8AN ZC0FAR RF9 SPAR 40TAN ZCOPSC89 ZFC9 **40FCAN** ZFC089 8ZCCOPCC9 PCAR SOEFC89 ESCS89 40CCCO SC8A FCC08AE  
PCO OFC0J 40FC08 EFC8EFCS **PCC8** 9SC8E BOEAE B9FCOR SCCV9 OBSC8AE EVSC89 ESCCOE OPCON SCAJAR 40FCOF89 OPCCOEFC89 EAN  
40FCCAE 40FCC E8SC8P COFAN EFC8AR EFC8AN OFCAJ PC0EFC8AN EPCCAE U OFCC2C9 OESAE 4CCAR BOCOFCC9 PC8AN SCBSAJ 40FC0FAN  
FCC8 BOEFC8 COFCAN I'AR \*AN 9ZC OFZ89 ZFCC9 SXAM

## Do Content Words Indicate Topics?



## Are VMS Pages in Order?

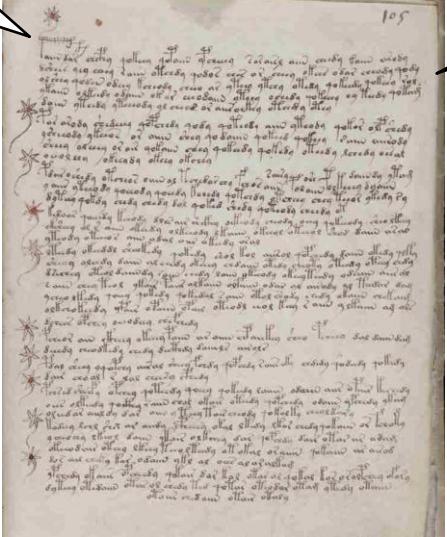
- Measure similarity between a pair of pages using cosine similarity (with bag-of-words)
- Count the % of pages  $P$  where the most similar page to  $P$  is adjacent to it



# Is VMS Prose?

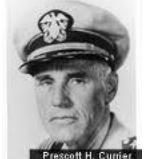
Special ligatures at beginning of "paragraphs"

Looks like paragraph structure



BUT:  
Lines begin and end disproportionately often with certain characters!

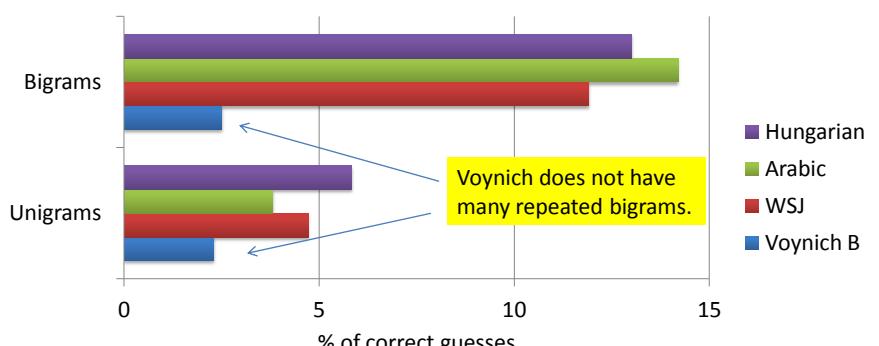
The line is a functional entity...



Prescott H. Currier

## Are VMS Word Sequences Predictable?

- Guess most likely word to follow current word
- Simulate game from bigram probabilities  
90-10 train-test splits



Category	Hungarian	Arabic	WSJ	Voynich B
Bigrams	~13.5%	~14.5%	~12.5%	~2.5%
Unigrams	~6.0%	~4.0%	~5.0%	~2.5%

# Zodiac Killer Ciphers

**Zodiac 408** (solved, 1969)

A grid of letters and symbols representing the solved Zodiac 408 cipher. The grid is approximately 20 columns by 20 rows. It contains various letters (A-Z), numbers, and special characters. A small portion of the grid is highlighted with red boxes, showing the word "KILLER" and other specific sequences.



**Zodiac 340** (still unsolved)

A grid of letters and symbols representing the unsolved Zodiac 340 cipher. The grid is approximately 20 columns by 20 rows. It contains various letters (A-Z), numbers, and special characters. A small portion of the grid is highlighted with red boxes, showing the word "KILLER" and other specific sequences.

# Zodiac Serial Killer

408-letter cipher (solved):

A grid of letters and symbols representing the solved 408-letter cipher. The grid is approximately 20 columns by 20 rows. It contains various letters (A-Z), numbers, and special characters. Some sections of the grid are highlighted in yellow, indicating specific patterns or solved segments.

(plus two more sections)

28	5	7	3	6	43	10	8	14	11	0	8	8
E	B	P	E	C	I	Z	/	U	B	D	K	I
W	V	+	E	3	G	Y	F	O	H	T	P	S
H	F	U	N	I	T	I	S	M	O	R	E	U
M	J	Y	L	U	I	K	N	G	W	T	D	A
A	N	K	I	A	L	B	P	R	T	I	U	M
S	I	N	T	H	E	F	O	R	S	H	J	E
R	E	M	G	D	A	D	K	I	T	B	R	Y
N	I	L	I	Y	S	J	O	N	O	V	W	U
L	O	J	P	U	P	S	X	M	E	H	Y	Y

[zodiologists.com](http://zodiologists.com)

# Zodiac Serial Killer

## Plaintext solution

“ I LIKE KILLING PEOPLE BECAUSE IT IS SO MUCH FUN IT IS MORE FUN THAN KILLING WILD GAME IN THE FORREST BECAUSE MAN IS THE MOST DANGEROUCE ANAMAL OF ALL TO KILL SOMETHING GIVES ME THE MOST THRILLING EXPERIENCE IT IS EVEN BETTER THAN GETTING YOUR ROCKS OFF WITH A GIRL THE BEST PART OF IT IS THAE WHEN I DIE I WILL BE REBORN IN PARADICE AND THEI HAVE KILLED WILL BECOME MY SLAVES I WILL NOT GIVE YOU MY NAME BECAUSE YOU WILL TRY TO SLOI DOWN OR ATOP MY COLLECTIOG OF SLAVES FOR MY AFTERLIFE EBEORIETEMETHHPITI ”

Plaintext has many misspellings

Final 18 plaintext characters of 408 are "junk"

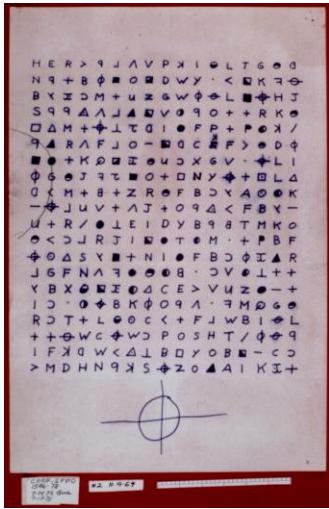
## Deciphering Zodiac 408 Bayesian models & Gibbs sampling

Language Model	Initial Sample	Decipherment Error
3-gram	Random	62.3
5-gram	Random	all wrong!
”	3-gram solution	42.6
Word 1-gram	Random	all wrong!
<i>Interpolated</i> 5-gram and word 1-gram	Random	79.2
”	5-gram solution	<b>3.3 / 2.6</b>

[Ravi & Knight 11]

See also Malte Nuhn's paper at ACL 2013!

## Unsolved Zodiac 340



Has no obvious reading order bias:

% cipher bigram types that repeat (freq > 1)	Left/ Right order	Up/ Down order	Diag. North- East	Diag. South- East
Zodiac 408 (solved)	13 %	5	7	5
Zodiac 340 (unsolved)	7	6	8	5

Could be nonsense ... or maybe  
bigrams are smoothed out via  
more careful substitutions.

## Other Unsolved Ciphers

### Beale (1885)

11, 95, 94, 241, 975,  
604, 230, 436, 664, 582, 150, 251, 284, 308, 211, 124, 211, 486, 225, 401, 370,  
11, 101, 305, 139, 189, 17, 33, 88, 208, 193, 145, 1, 94, 73, 416, 918, 263, 28, 500,  
138, 210, 118, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113,  
118, 320, 138, 36, 416, 280, 15, 71, 224, 991, 44, 16, 401, 39, 88, 61, 304, 12, 21,  
24, 283, 134, 92, 63, 249, 486, 682, 7, 219, 184, 369, 780, 18, 64, 463, 474, 131,  
106, 97, 103, 862, 70, 60, 1317, 471, 540, 208, 121, 890, 346, 36, 150, 59, 546,  
614, 13, 120, 63, 219, 812, 2160, 1780, 99, 35, 18, 21, 136, 872, 15, 25, 170, 88, 4,  
50, 106, 301, 13, 408, 680, 93, 86, 116, 530, 82, 568, 8, 102, 38, 416, 89, 71, 216,  
728, 965, 818, 2, 38, 121, 195, 14, 326, 148, 234, 18, 95, 131, 234, 361, 824, 5,  
21, 219, 324, 824, 431, 64, 326, 19, 48, 122, 89, 216, 284, 919, 861, 326, 985,  
233, 64, 68, 232, 431, 960, 50, 29, 81, 216, 323, 603, 14, 612, 81, 360, 36, 51, 62,  
124, 232, 233, 234, 235, 236, 237, 238, 239, 239, 240, 241, 242, 243, 244, 245, 246,  
10, 6, 66, 119, 38, 41, 49, 602, 423, 962, 392, 294, 875, 78, 14, 23, 111, 111, 111,  
31, 501, 823, 218, 280, 34, 24, 150, 1060, 162, 286, 19, 21, 17, 340, 19, 242, 31,  
25, 26, 27, 28, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29, 29,  
548, 96, 11, 203, 77, 364, 218, 65, 667, 890, 236, 154, 211, 10, 98, 34, 119, 56,  
216, 119, 71, 218, 1164, 1496, 1817, 51, 39, 210, 36, 3, 10, 546, 232, 22, 141, 617,  
24, 54, 24, 54, 24, 54, 24, 54, 24, 54, 24, 54, 24, 54, 24, 54, 24, 54, 24, 54, 24, 54,  
212, 416, 127, 931, 19, 4, 63, 96, 12, 101, 418, 16, 140, 230, 460, 538, 19, 27, 88,  
612, 1431, 90, 716, 275, 74, 83, 11, 426, 89, 72, 84, 1300, 1700, 814, 221, 132,  
40, 23,  
447, 55, 86, 34, 43, 212, 107, 96, 314, 264, 1055, 323, 428, 601, 203, 124, 95, 216,  
814, 290, 654, 820, 2, 301, 112, 176, 213, 71, 87, 96, 202, 35, 10, 2, 41, 17, 84,  
221, 730, 626, 214, 11, 86, 766.

### Taman Shud (1948)



### FBI (1999)



### Kryptos (1990)

OXOOGHULBSOLIFBBWFLRVQQPRNGKSSO  
TWTQSJQSEKZZWATJKLUDIAWINFB**NYP**  
**VTTMZFPKGDKZXTJCDIGKUHUAEKCAR**

**NYPVTT** = BERLIN (2011 clue)

A K P N T E G L S E - S E E R T E  
V L S E M T S E - C T S E - W S E - F R T S E  
P U T T R S E O N P R S E N C D V C C  
N W L D X C R C H S P N E W L D S

(2 pp total)

# Collected Ciphers

Mina aen gevar, af perer-oli,  
ni lava obet pera, meni mina  
lurus grön var all in ..... der  
ujor. Je an dora  
af salta SONORAM  
nek o oli new.

—.5. WYN 4-a E-TN-3 CYA 4-a 87 ZBM R-N  
6M Y-R-E7A ER YM BWTAS NMEN 3TA E71 NED2  
—.5 37 ZBM NY M2R-E7A, SWY-A VBA ANKURAN

Digitized by srujanika@gmail.com

+ many more!

1866, 20 July

Ziffernblatt vom 20.07.1866  
Nr. 201

Schleswig-Holsteinische Zeitung  
aus dem Hause der Meinung. Bei C. G. A. und  
H. C. H. 1866. In englischer Sprache abgedruckt.  
Gebundene Heft 44 V. 11. 1866. D. 173 N. M.

# Writing as a code for speech

# Archaeological Decipherment

ciphertext



Mayan  
glyphs

# Archaeological Decipherment

Thinks Mayan decipherment should be based on ideographic rather than linguistic principles.

Resists notion that the glyphs have a phonetic component.



J. Eric S. Thompson

It's phonetic.



Yuri Knorozov

ciphertext



Mayan  
glyphs

## Archaeological Decipherment

- Mayan glyphs
  - Egyptian glyphs (Rosetta Stone)
  - Linear B
- etc

Computer did not play much of a role in these successful decipherments

## Archaeological Decipherment

ciphertext

primera parte  
del ingenioso  
hidalgo don ...

[Knight & Yamada 99]

## Archaeological Decipherment

"When I look at these squiggles, I say to myself, this is **really a sequence of Spanish phonemes**, but it has been encoded in some strange symbols..."



OUR HERO

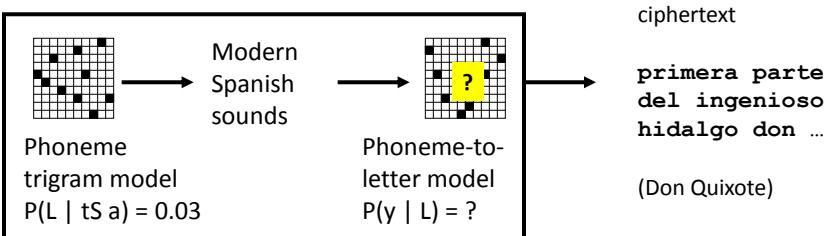
ciphertext

**primera parte  
del ingenioso  
hidalgo don ...**

(Don Quixote)

[Knight &amp; Yamada 99]

## Archaeological Decipherment



26 sounds:

B, D, G, J (canyon),  
L (yarn), T (thin), a,  
b, d, e, f, g, i, k, l,  
m, n, o, p, r,  
rr (trilled), s,  
t, tS, u, x (hat)

32 letters:

ñ, á, é, í, ó, ú,  
a, b, c, d, e, f, g,  
h, i, j, k, l, m, n,  
o, p, q, r, s, t, u  
v, w, x, y, z

EM approach = 93% accurate phonetic decipherment

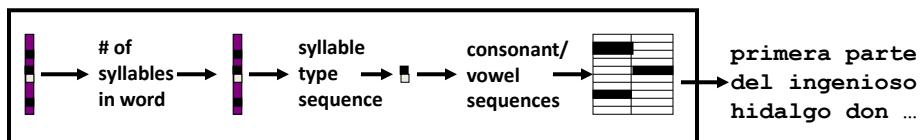
[Knight &amp; Yamada 99]

## What if Spoken Language Behind Script is Unknown?

- Build a universal model  $P(p)$  of human phoneme sequence production
  - human might generally say: K AH N AH R IY
  - human won't generally say: R T R K L K
- Find a  $P(c | p)$  table
  - such that there is a decoding with a good universal  $P(p)$  score
- Phoneme & syllable inventory
  - if z, then s
  - all have CV syllables; if VCC, then also VC
- Syllable sonority structure
  - dram, lomp, ? rdam, ? lopm
- Physiological preference constraints
  - tomp, tont, ? tomk, ? tonp

[Knight et al 06]

## Unknown Source Language



$$\begin{array}{llll}
 P(1) = ? & P(CV) = ? & P(V | V) = ? & P(a | V) = ? \\
 P(2) = ? & P(V) = ? & P(VV | V) = ? & P(a | C) = ? \\
 \text{etc.} & P(CVC) = ? & & \text{etc.} \\
 & + 7 \text{ others} & &
 \end{array}$$

Input: primera parte del ingenioso ...  
 Output: NSV.NV.NV NVS.NV NVS VS.NV.SV.V.NV ...

**S** = sonorous consonant phoneme  
**N** = non-sonorous consonant phoneme  
**V** = vowel phoneme

[Knight et al 06]

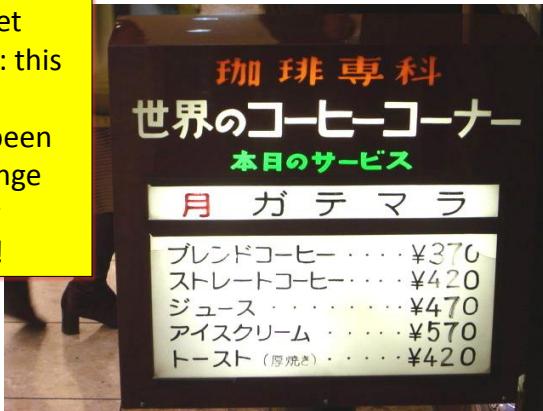
See Y. Kim & B. Snyder's  
 ACL 2013 paper addressing  
 100s of human languages!

## Practical Detour: Phoneme Substitution Ciphers

When I look at street signs in Tokyo, I say: this is **really written in English**, but it has been coded in some strange symbols. I will now proceed to decode!



OUR HERO



Parallel data: [Knight & Graehl 97]  
Non-parallel data: [Ravi & Knight 09a]

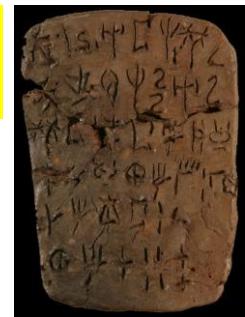
## Undeciphered Writing Systems

## Undeciphered writing systems

Indus Valley  
Script  
(3300BC)



Linear A  
(1900BC)



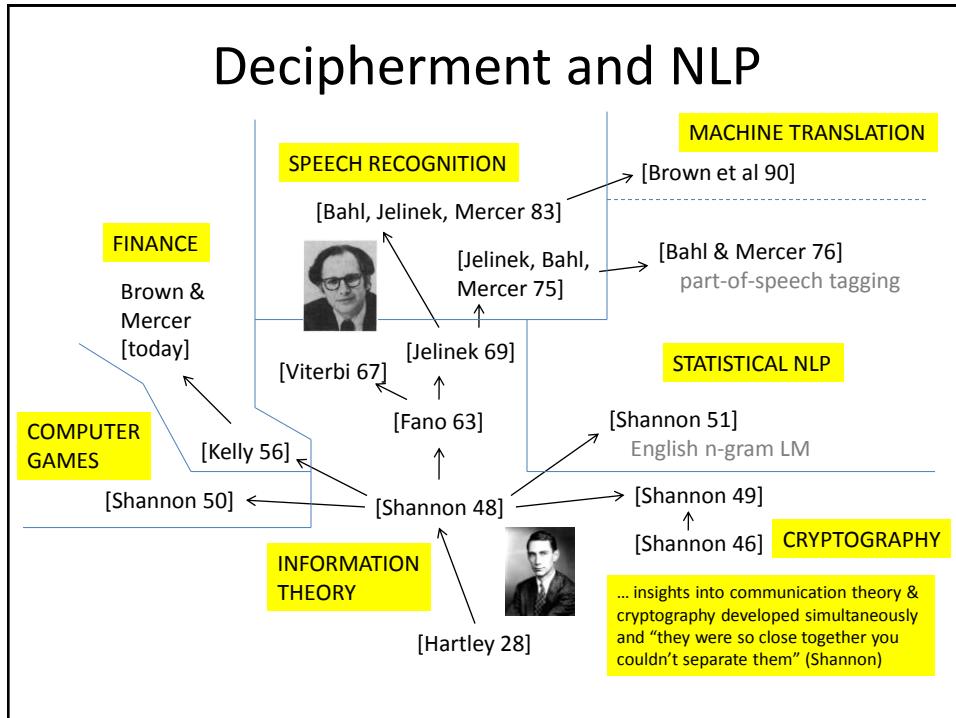
Rongorongo (1800s?)



Phaistos Disc (1700BC?)



## Conclusions



# Decipherment and NLP

	Cryptography	Translation
Manual	Manual encoding	Human translation
Mechanical	1920s Mechanical encoding; intuition-based decryption	1960s Rule-based MT
Mathematical	1950s Computer decryption, based on information theory	1990s Statistical MT
Higher math, deeper understanding	1980s Public-key systems, based on number theory	2020s ??? ??? ???



thanks