

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: ABCDEFGHIJKLMN**OP**QRSTUVWXYZ
CIPHER: P**L**OKMIJNUHBYGV**TFC**RD**XES**Z**AQ**W
- Plaintext: **HELLO KITTY . . .**
- Ciphertext: **NMYYT BUXXQ . . .**
- 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."

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
L 10
M 6
N 1
O
P 1
Q 10
R 3
S
T 7
U
V
W 1
X 5
Y 7
Z 2

.
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.
HYD FKXC, FQ MKX RLQIQ 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
L 10
M 6
N 1 .
O
P 1 .
Q 10
R 3 .
S
T 7
U
V
W 1 .
X 5
Y 7
Z 2 .

.
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.
HYD FKXC, FQ MKX RLQIQ 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	.a	.a	.	.	
KDCY	LQZKTLJKX	CY	MBCYJQL:	"TR	
.	.a	.	a	. . .	
HYD	FKXC,	FQ	MKX	RLQIQ	HYDL
aa	
MKL	DXCTW	RDCDLQ	JQMNKXTMB		
.	.	a	.a.	.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 .

a	e.a	.a	.	e	.
KDCY	LQZKTLJKX	CY	MBCYJQL:	"TR	
.	.a	.e	a	. ee.e	.
HYD	FKXC,	FQ	MKX	RLQIQ	HYDL
a	.	.	e	.e	.a
MKL	DXCTW	RDCDLQ	JQMNKXTMB		
.	.e	a	.a.	e.a	
PTBMYEQL	K	FKH	CY	LQZKTL	TC."

didn't create "ae"

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.ao .a	.e	o.		
KDCY	LQZKTLJKX	CY	MBCYJQL:	"TR	
.	.a	.e a	. ee.e	.	
HYD	FKXC,	FQ	MKX	RLQIQ	HYDL
a	o.	. e	.e .a	o	
MKL	DXCTW	RDCDLQ	JQMNKXTMB		
.o	.e a	.a.	e.ao	o	
PTBMYEQL	K	FKH	CY	LQZKTL	TC."
don't like "ao" – back up!					

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	o e.a	.a	o	o.e	.
KDCY	LQZKTLJKX	CY	MBCYJQL:	"TR	
.o	.a	.e a	. ee.e	.o	
HYD	FKXC,	FQ	MKX	RLQIQ	HYDL
a	.	. e	.e .a		
MKL	DXCTW	RDCDLQ	JQMNKXTMB		
.	o.e	a .a.	o 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	6 ### V
Z	2 .

a o re.a r.a o o.e f
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.o .a .e a freeze .o r
HYD FKXC, FQ MKX RLQIQ HYDL

ar . f re .e .a
MKL DXCTW RDCDLQ JQMNKXTMB

. o.er a .a. o re.a r
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 6 #### V
 Z 2 .

a o re.a r.a o o.e f
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.o .a .e a freeze .o r
HYD FKXC, FQ MKX RLQIQ HYDL

ar . f re .e .a
MKL DXCTW RDCDLQ JQMNKXTMB

. o.er a .a. o re.a r
PTBMYEQL K FKH CY LQZKTL TC."

frequent cipher letters: ~~Q~~ ~~Z~~ ~~X~~ C D T M ~~X~~ X
 frequent English letters: ~~e~~ t ~~r~~ a n i ~~s~~ h

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 6 #### V
 Z 2 .

a no re.air.a no no.e if
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.o .a n .e a freeze .o r
HYD FKXC, FQ MKX RLQQIQ HYDL

ar ni. f n re .e .a i
MKL DXCTW RDCDLQ JQMNKXTMB

.i o.er a .a. no re.air in
PTBMYEQL K FKH CY LQZKTL TC."

frequent cipher letters: ~~Q~~ ~~Z~~ ~~X~~ C D T M ~~X~~ X
 frequent English letters: ~~e~~ ~~t~~ ~~a~~ ~~n~~ ~~i~~ ~~s~~ h

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 6 #### V
 Z 2 .

a to re.air.a to to.e if
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.o .a t .e a freeze .o r
HYD FKXC, FQ MKX RLQQIQ HYDL

ar ti. f t re .e .a i
MKL DXCTW RDCDLQ JQMNKXTMB

.i o.er a .a. to re.air it
PTBMYEQL K FKH CY LQZKTL TC."

frequent cipher letters: ~~Q~~ ~~Z~~ ~~X~~ C D ~~T~~ M ~~X~~ X
 frequent English letters: ~~e~~ ~~t~~ ~~a~~ ~~n~~ ~~i~~ ~~s~~ h

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 6 #### V
 Z 2 .

a to repair.a to to.e if
KDCY LQZKTLJKX CY MDBCYJQL: "TR

.o .a t .e a freeze .o r
HYD FKXC, FQ MKX RLQQIQ HYDL

ar ti. f t re .e .a i
MKL DXCTW RDCDLQ JQMNKXTMB

.i o.er a .a. to repair it
PTBMYEQL K FKH CY LQZKTL TC."

frequent cipher letters: ~~Q~~ ~~Z~~ ~~X~~ ~~C~~ D ~~F~~ M ~~X~~ X
 frequent English letters: ~~e~~ ~~t~~ ~~r~~ a n ~~i~~ s h

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 6 #### V
 Z 2 .

auto repairman to customer: if
KDCY LQZKTLJKX CY MDBCYJQL: "TR

you wait we can freeze your
HYD FKXC, FQ MKX RLQQIQ HYDL

car until future mechanics
MKL DXCTW RDCDLQ JQMNKXTMB

discover a way to repair it
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 6 #### V
 Z 2 .

Pattern word dictionaries

KDCY **LQZKTLJKX** CY MDBCYJQL: "TR
abcdeafdg

HYD FKXC **abnegated** MKX **RLQQIQ** HYDL
abccdc
 abnegates
 advocator
 airedales
 alienages
 alienated
 alienates
 amperages
 cadencies
 capricorn
 cogencies
 escapeway
 healthily
 imbeciles
 imperiled
 incurious
 inherited
 injurious
 landslide
 octagonal
 oklahoman
 overboard
 repairman
 sacristry
 unrebuked
 unsecured

MKL DXCT **CDLQ** JQ **abnegated** TMB **abcdefghijklm**
 basses
 bassos
 bosses
 breeze
 budded
 ...
 cheese
 cusses
 dossen
 finnan
 fleece
 fosses
 freeze
 ...
 terror
 tosses
 tweeze
 wadded
 wheeze

PTBMYEQI **KH** CY L **L TC."**
 consumptively
 copyrightable
 documentarily
 lycanthropies
 musicotherapy
 semivoluntary
 subordinately
 unpredictably

OR, NORWEGIAN!
 filmprodusent
 kurspamelding
 publikasjon
 stylemarginpx
 uproblematisk

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 = ȝ i l y †

B = û

C = ô ñ

D = ƞ

E = ˆ ƒ Δ ƶ f † ˆ ˆ ˆ

F = p

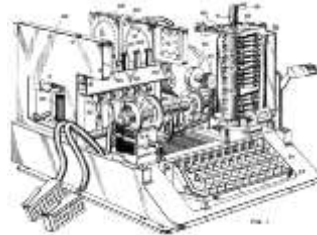
G = ˆ ˆ ...

Cipher Types

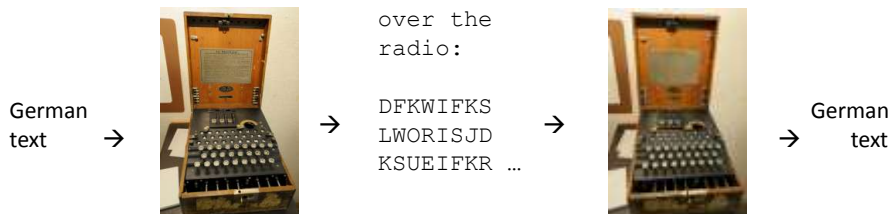
- http://cryptogram.org/cipher_types.html
 - documents ~70 types
- E.g., RUNNING KEY cipher
 - key = agreed-upon standard English text
 - $\text{ciphertext}(i) = [\text{plaintext}(i) + \text{key}(i)] \bmod 26$
 - effectively uses 26 substitution keys
 - breakable!
 - we search for a key and (resulting)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

How to Make Things Efficient?

- Mechanical encryption/decryption devices



“First NLP Task Ever” (1930s-40s) Breaking the German Enigma Cipher



input (intercepted ciphertext): DFKWIFKSLWORISJDKSUEIFKR ...
output (plaintext): VASISTDASHERRCAPITANRICH ...

“First NLP Task Ever” (1930s-40s) Breaking the German Enigma Cipher

Substitution system

$N \rightarrow J$

Substitution table **changes**
with every keystroke:

$NNN \rightarrow JTE$

Flattens out ciphertext
letter distributions.

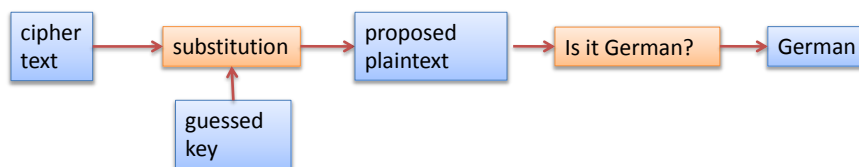


Secret key =
initial rotor
ordering and
settings

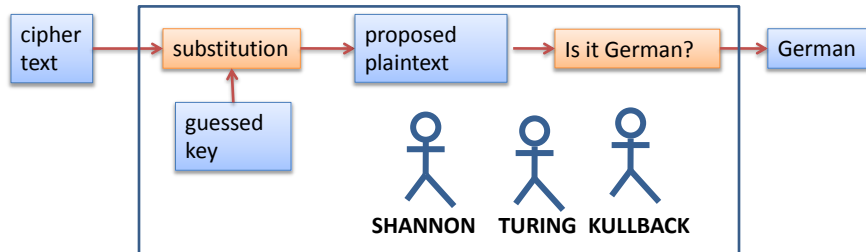
Reversible behavior

$NNN \rightarrow JTE \rightarrow NNN$

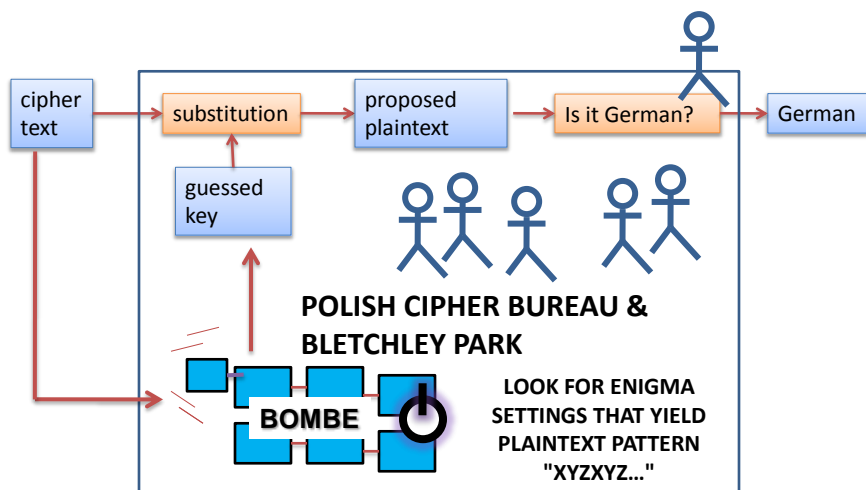
Breaking Enigma

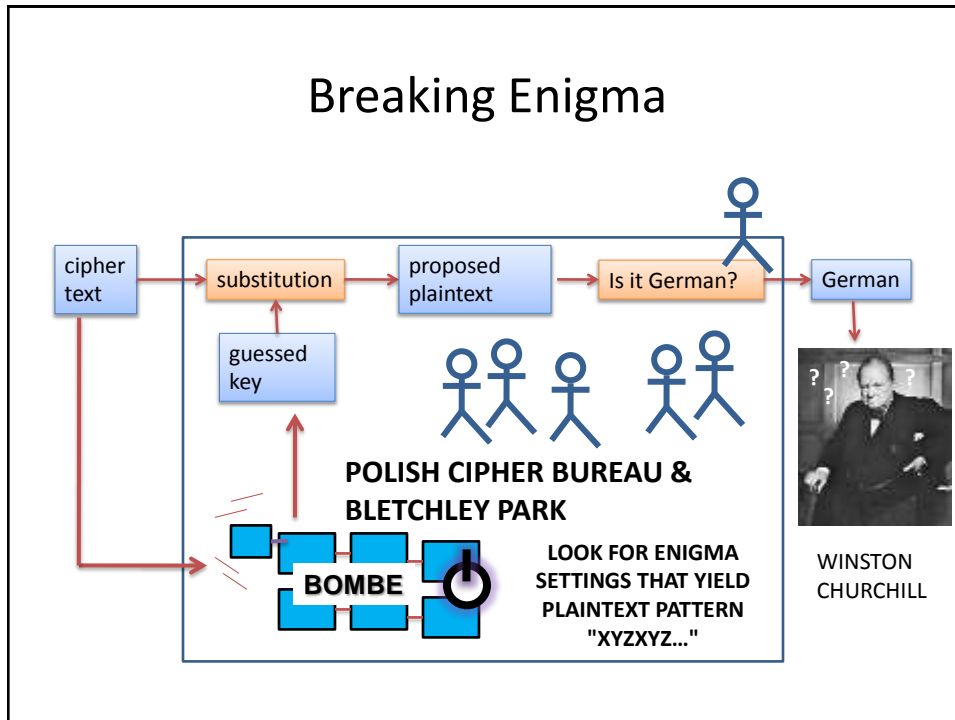


Breaking Enigma



Breaking Enigma





Enigma

- Mathematical breakthroughs:
 - Log-odds for weight of evidence [Good, Turing]
 - Smoothing with prior [Good, Turing]
 - Information theory [Shannon]
- elegant,
powerful,
widely-applicable
mathematics
- 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)

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 setting up inverse rods for the middle wheel. The rods are paired off according to R.H.W. couplings, i.e. M.W. output. This has been done for the couplings ku, fx, ep which arose in the DANZIGVON crib in Fig 55, assuming the red wheel in the middle. The pairs in each column of this set up give possible M.W. couplings. We have now to find out whether these couplings are possible. 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 sufficient to have 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.



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 setting up inverse rods for the middle wheel. The rods are paired off according to R.H.W. couplings, i.e. M.W. output. This has been done for the couplings ku, fx, ep which arose in the DANZIGVON crib in Fig 55, assuming the red wheel in the middle. The pairs in each column of this set up give possible M.W. couplings. We have now to find out whether these couplings are possible. 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 sufficient to have 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
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 trying all possible inverse rods for the middle wheel. The rods are paired off to give possible R.H.W. couplings, i.e. M.W. output. This has been done for the pair **fx, ep** which arose in the DANZIGVON crib in Fig 55, assuming **ep** in the middle. The pairs in each column of this set up give possible M.W. couplings. We have now to find out whether these couplings are possible. 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 sufficient to have 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



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 trying all possible inverse rods for the middle wheel. The rods are paired off to give possible R.H.W. couplings, i.e. M.W. output. This has been done for the pair **fx, ep** which arose in the DANZIGVON crib in Fig 55, assuming **ep** in the middle. The pairs in each column of this set up give possible M.W. couplings. We have now to find out whether these couplings are possible. 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 sufficient to have 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

possible M.W.
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 cryptanalyst. 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 season 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.

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 cryptanalyst. 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.

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 cryptanalyst. 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.

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, and a 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 cryptanalyst. 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)

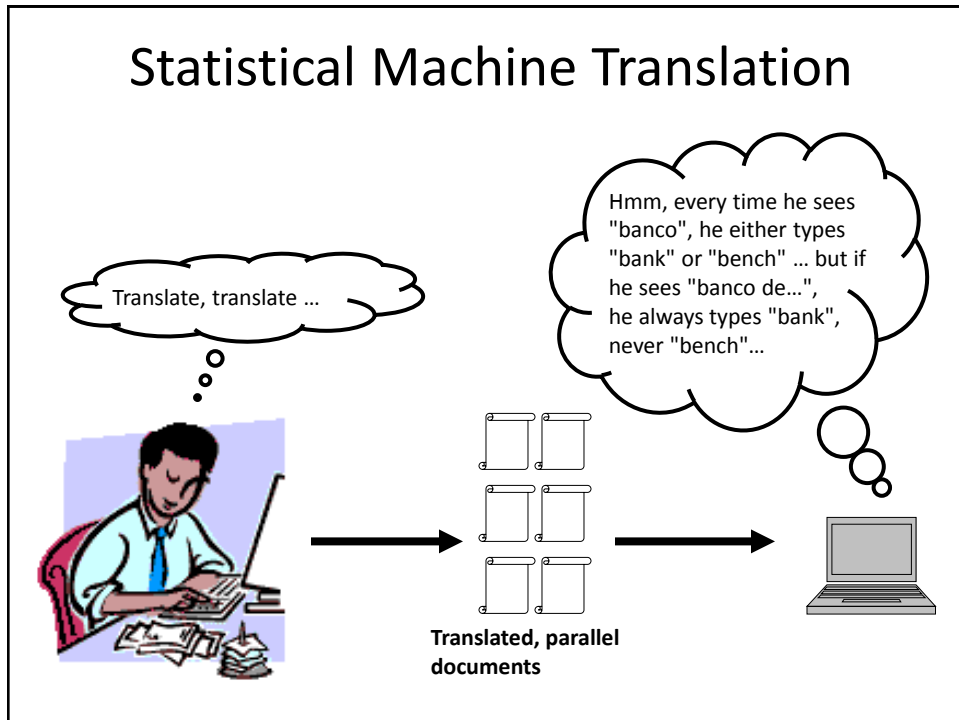


OUR HERO

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 clients 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 zenzanine . 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 .	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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok 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 crrrok hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok 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 crrrok hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok 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 crrrok hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok 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** **errrok** hihok yorok klok 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok hihok yorok zanzanak .
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** **errrok** **hihok** yorok klok 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok errrok hihok yorok zanzanak .
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** crrrok **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 anak plok sprok .	8a. lalok brok anak 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 anak 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 crrrok hihok yorok zanzanak .
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** crrrok **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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok clock .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok hihok yorok zanzanak .
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** **clock** 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok clock . ???
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** **clock** 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok clock .
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 crrrok hihok yorok klok** 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 anak plok sprok .	8a. lalok brok anak 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 anak drok brok jok .	10a. lalok mok nok yorok ghirok klok .
4b. at-voon krat pippat sat lat .	10b. wat nnat gat mat bat hilat .
5a. wiwok farok izok stok .	11a. lalok nok crrrok 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 .

Learn Translation Knowledge from Non-Parallel Text?

English/Albanian
Parallel text



Translation model

English text Albanian text



Translation model

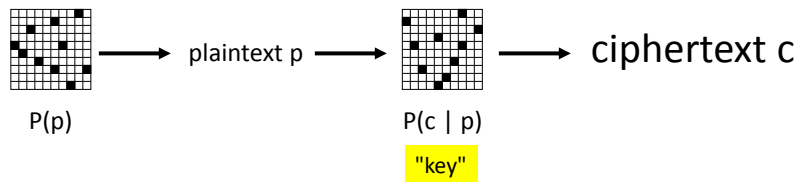
Is this what Weaver had in mind?
We'll come back to this idea.

Automatic decipherment

Letter Substitution Cipher

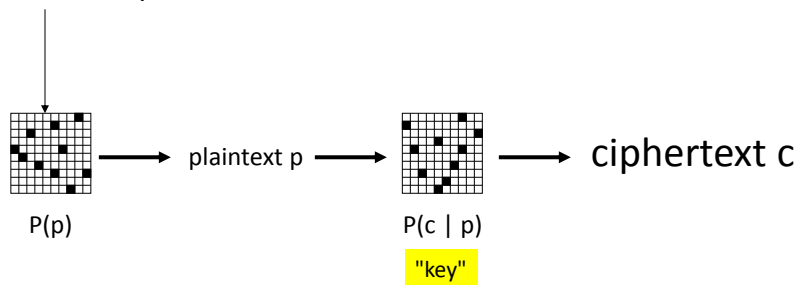
ciphertext c

Letter Substitution Cipher



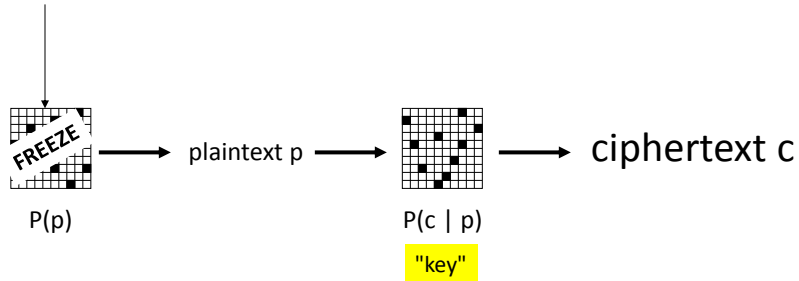
Letter Substitution Cipher

plaintext samples,
unrelated to ciphertext

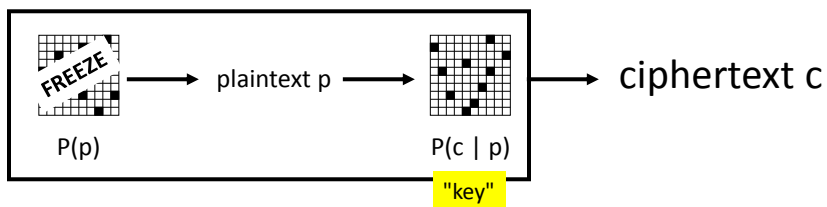


Letter Substitution Cipher

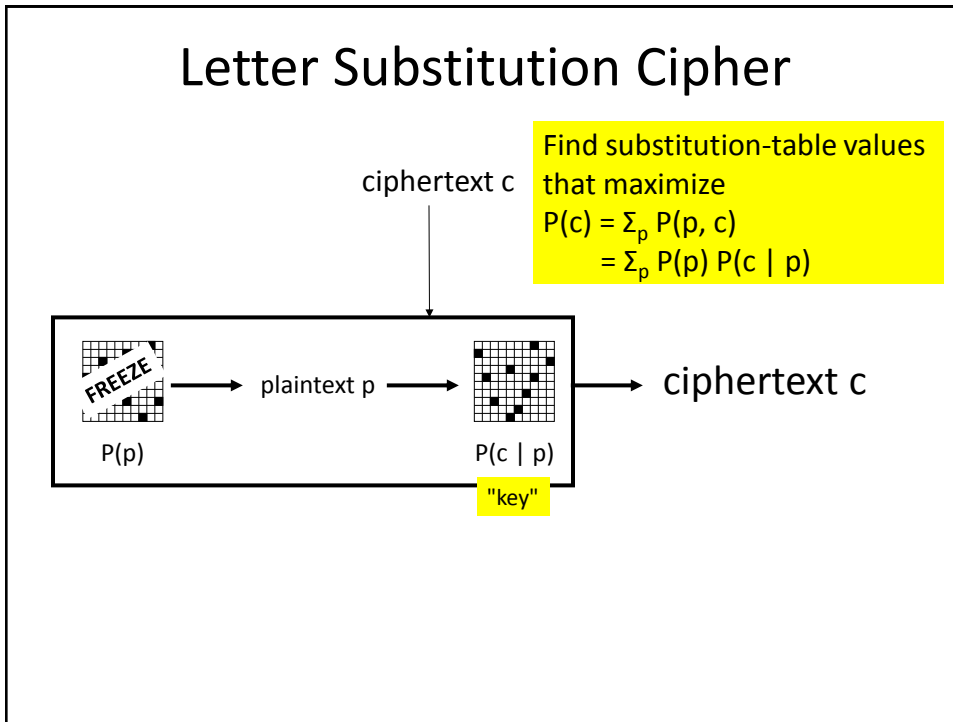
plaintext samples,
unrelated to ciphertext



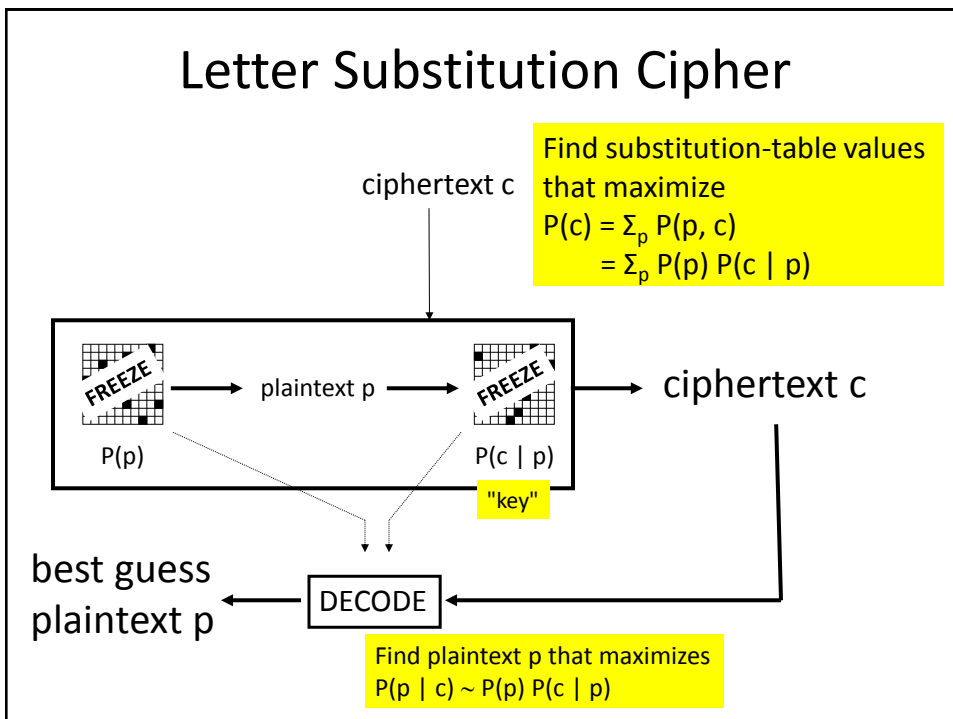
Letter Substitution Cipher

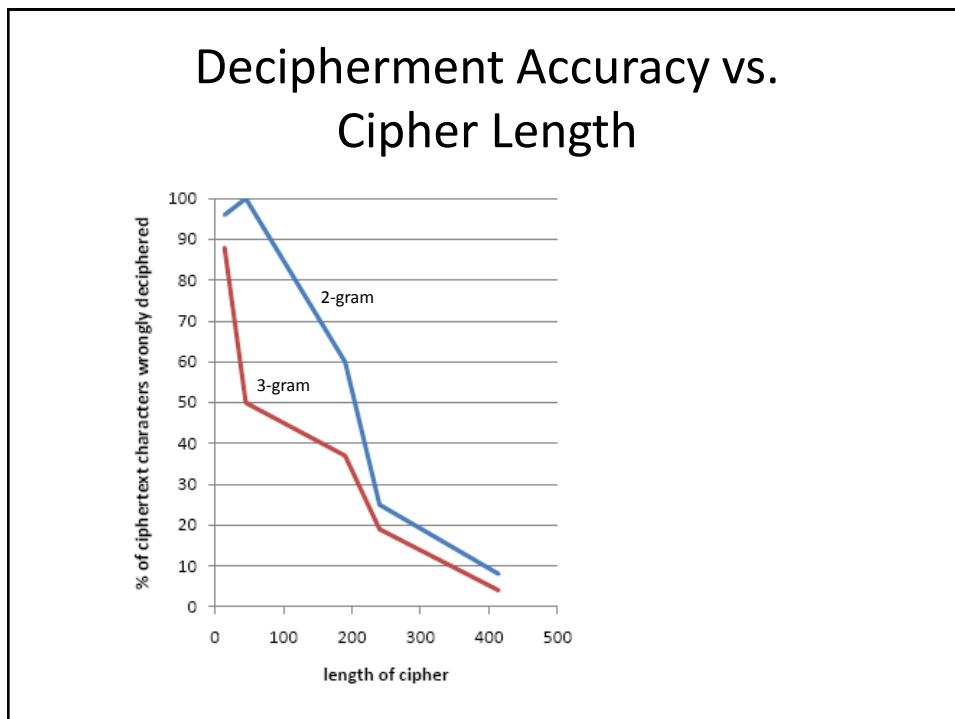
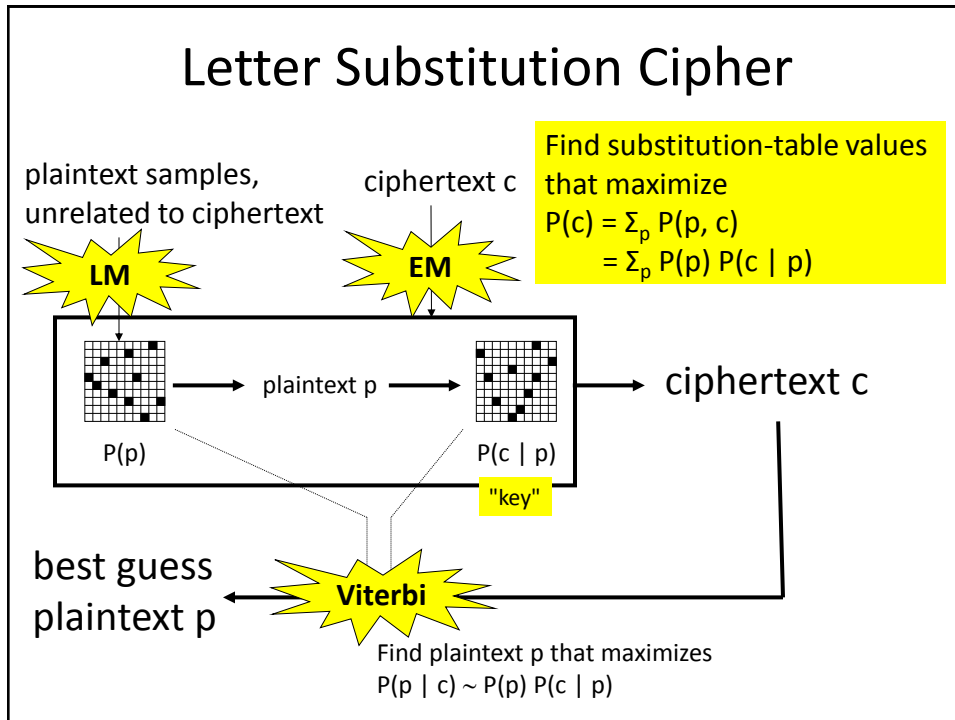


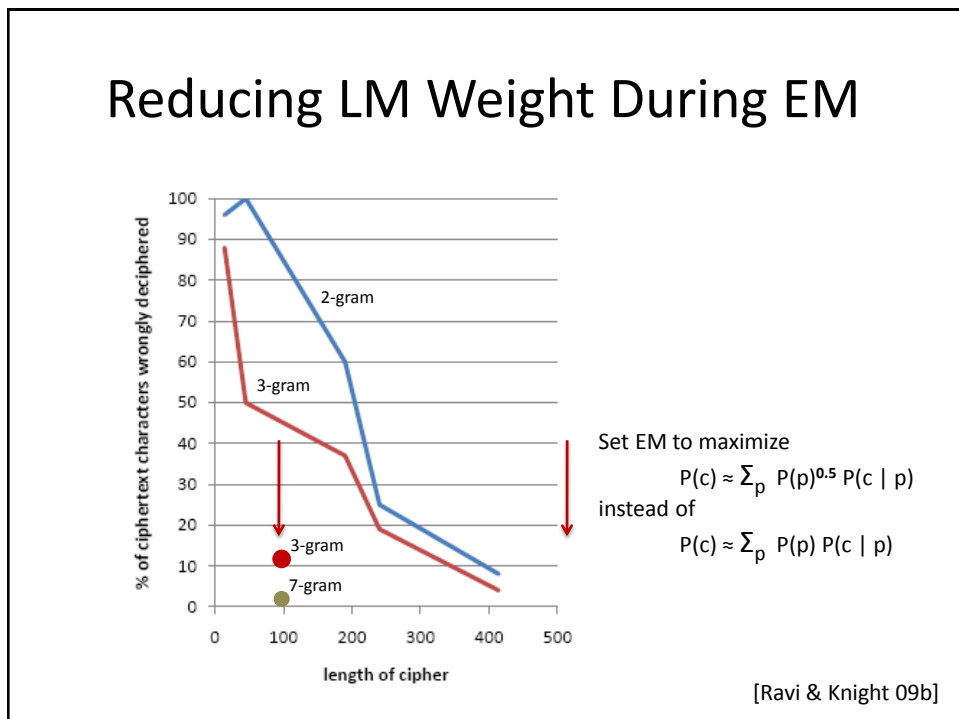
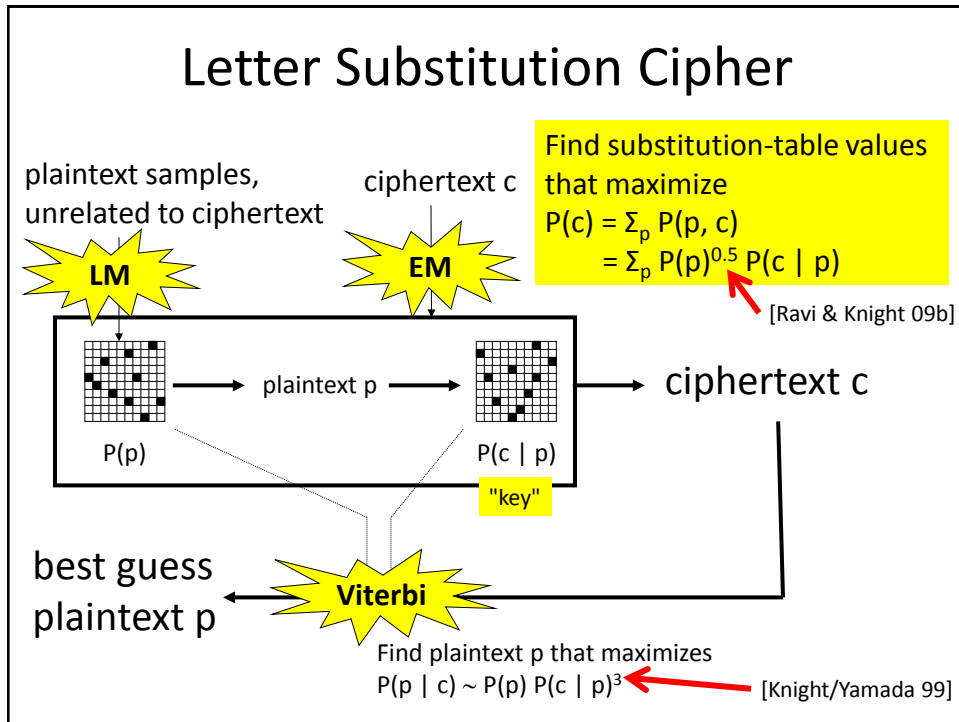
Letter Substitution Cipher



Letter Substitution Cipher

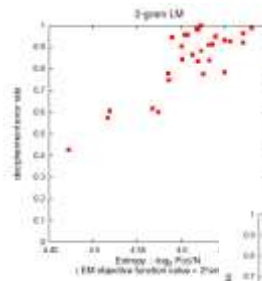
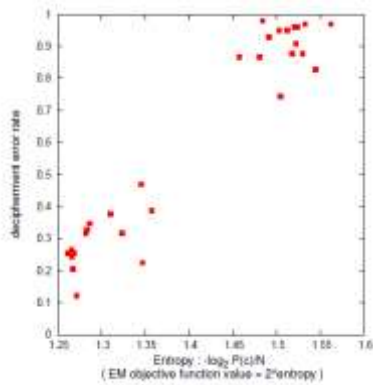




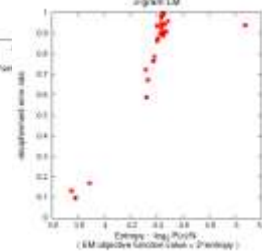


Random Restarts are Critical

English 98-letter cipher, 3-gram LM



Japanese
syllable
cipher

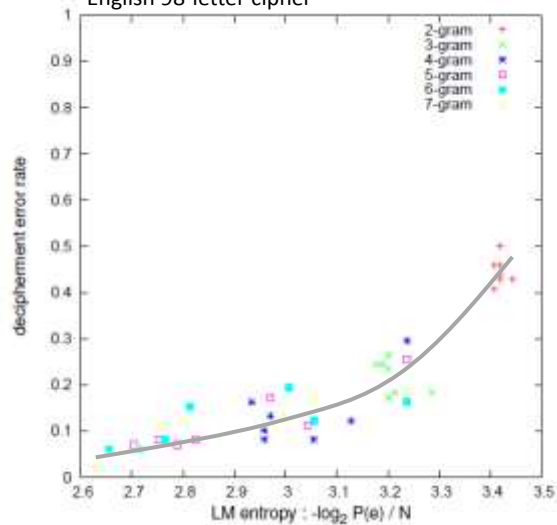


even people do restarts!

[Ravi & Knight 09b]

Good Language Models are Critical

English 98-letter cipher



[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

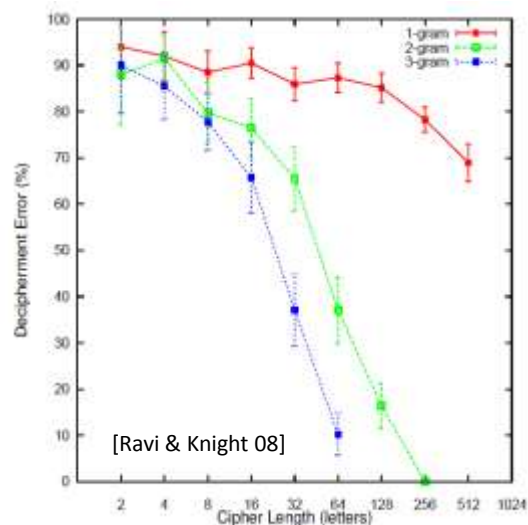
Deterministic Keys

* Use ILP to search only deterministic keys.

* Exact, no restarts.

Cipher Length	EM error	ILP error
52	85 %	21 %
98	45 %	12 %
414	10 %	0.5 %

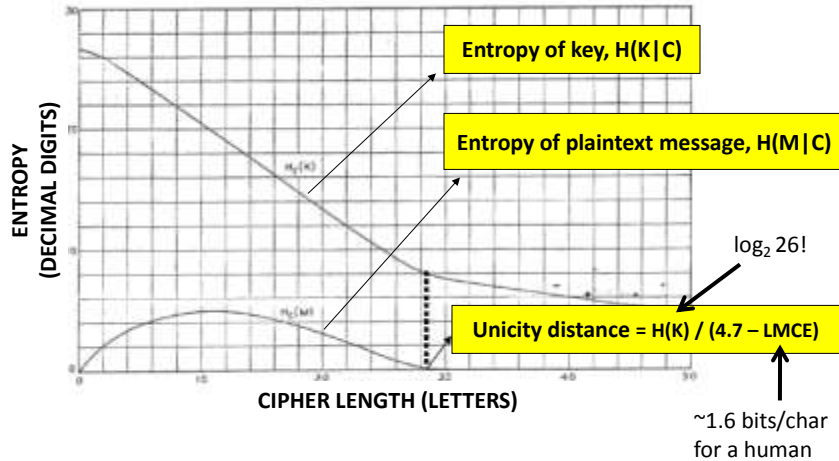
Using 2-gram letter-based LM



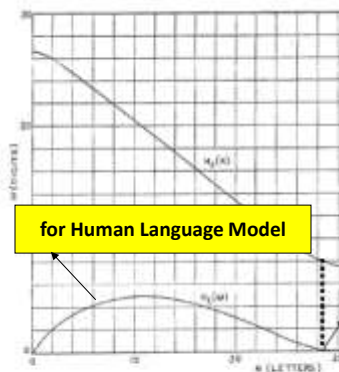
[Shannon 46, 49]

"Communication Theory of Secrecy Systems"

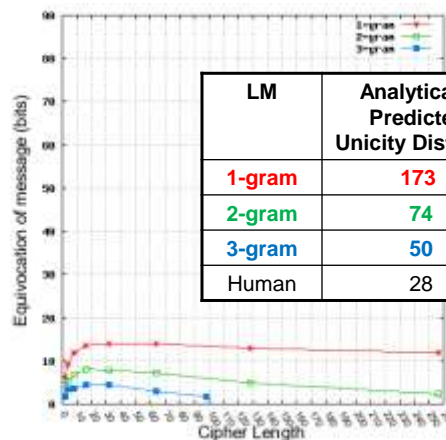
- Shannon analytically predicted uncertainty about key and message
- Graphed it for a human-level language model



Verifying Shannon's Prediction of Plaintext Message Uncertainty



ANALYTIC CURVES (Shannon)



ACTUAL CURVES

[Ravi & 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



[Knight, Megyesi, Schaefer 11]

Copiale Cipher

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

Section headers

Lines ≈ equal length

Paragraphs and section titles always begin with **capitalized Roman letters**.

Non-enciphered inscriptions: **Copiales 3 and Philipp 1866**

Some scratch-outs, rare

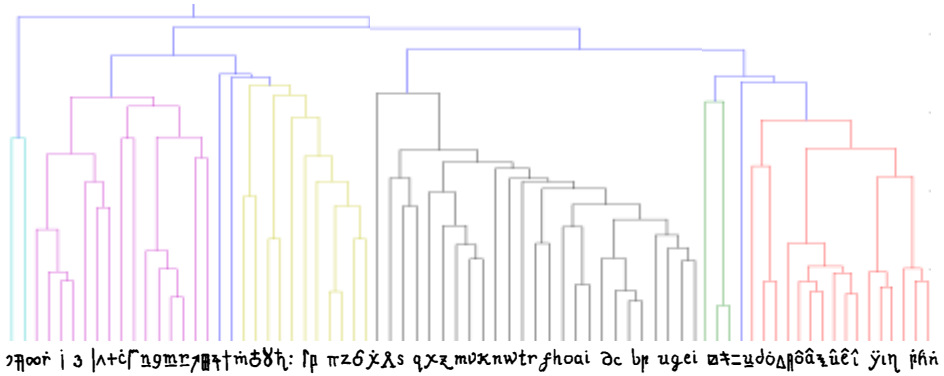
Preview text fragments ("catchwords")

Letter Frequencies

digraphs:	trigraphs:	tendencies:
ɔ ħ 99	ɔ ħ ^ 47	â, ê, î, ô, û followed by ʒ and j
ç : 66	ç : ɷ 23	â, ê, î, ô, û preceded by z and π
ħ ^ 49	η ɔ ħ 22	
: ɷ 48	ÿ ɔ ħ 18	
z ʘ 44	ʁ ç 17	

Clustering of Cipher Letters

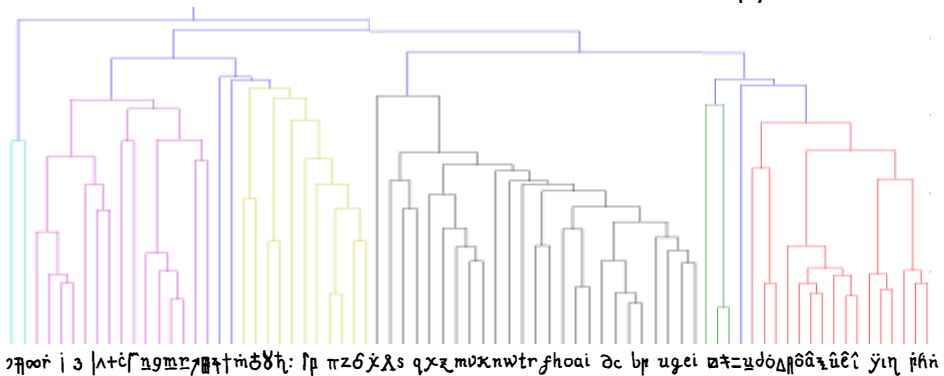
letters grouped if they have similar contexts (L/R neighbors)
Scipy software



thanks Jon Graehl

Clustering of Cipher Letters

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



underlined
letters

unaccented
Roman letters

circumflexed
vowels

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

κ m û r: p z i ô f | ÿ ʷ ħ ê ĭ ħ ɛ ι λ π π â ɔ b Δ g z =
 j ʎ l z ɹ ɸ ç λ â r ɔ g κ λ ħ ʃ r ħ ʃ λ ɛ i n r î | Δ r ô ʀ λ a
 = g z w π ÿ ê c Δ r ô Δ + b z η r i ʃ ÿ î r z ɹ f λ z
 π ʃ i ʎ ð r = ʃ | û λ s ʃ m δ m | ɛ η g â | κ ħ = λ ħ | l ʃ ô
 ʋ ʃ : r î λ b i f u m ʏ j z ʋ z â i x ʎ r m π i z ħ λ c ô ð g g
 z û + ʀ ʃ ʃ n ʃ ê ɔ ʒ h λ ħ ʃ i ħ ʃ λ î ɔ t n ʃ | r ô ʏ m â
 + h ħ r z ô ɔ n b s η + : ɛ r κ ʀ r ô ʀ ħ Δ c û λ g = n z κ
 ʀ ɛ ∞ n z ʀ ɔ ʃ ħ ĩ r π z ê ʏ n g = r π g ô Δ ɔ n z ʀ | κ
 π η i ʃ ÿ r î ʒ π ɹ λ b g λ i t b ô û ʃ ʃ ħ ʎ ô λ e z η ô
 | û r c ʃ ɔ ʒ ô λ b η ħ m ð

κ f n g l x n a c b f z m x
 l b u v c g h t r h b k g n x n
 f g g n x b g b e c b ...

Decipher against 80 plaintext languages.

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

κ m û r: p z i ô f | ÿ ʷ ħ ê ĭ ħ ɛ ι λ π π â ɔ b Δ g z =
 j ʎ l z ɹ ɸ ç λ â r ɔ g κ λ ħ ʃ r ħ ʃ λ ɛ i n r î | Δ r ô ʀ λ a
 = g z w π ÿ ê c Δ r ô Δ + b z η r i ʃ ÿ î r z ɹ f λ z
 π ʃ i ʎ ð r = ʃ | û λ s ʃ m δ m | ɛ η g â | κ ħ = λ ħ | l ʃ ô
 ʋ ʃ : r î λ b i f u m ʏ j z ʋ z â i x ʎ r m π i z ħ λ c ô ð g g
 z û + ʀ ʃ ʃ n ʃ ê ɔ ʒ h λ ħ ʃ i ħ ʃ λ î ɔ t n ʃ | r ô ʏ m â
 + h ħ r z ô ɔ n b s η + : ɛ r κ ʀ r ô ʀ ħ Δ c û λ g = n z κ
 ʀ ɛ ∞ n z ʀ ɔ ʃ ħ ĩ r π z ê ʏ n g = r π g ô Δ ɔ n z ʀ | κ
 π η i ʃ ÿ r î ʒ π ɹ λ b g λ i t b ô û ʃ ʃ ħ ʎ ô λ e z η ô
 | û r c ʃ ɔ ʒ ô λ b η ħ m ð

κ f n g l x n a c b f z m x
 l b u v c g h t r h b k g n x n
 f g g n x b g b e c b ...

Decipher against 80 **FAIL** plaintext languages.

Second Decipherment Approach

Homophonic cipher,
e.g.:

A = 8 | l y r

B = û

C = ó ñ

D = ʈ

E = ʃ Δ ʂ f î ʒ

F = ʀ

G = ʝ

etc.



κ m û r : p z i ô f | ʝ ʈ ê j ʈ ʒ i λ n π â z b Δ g z =
j ʀ l z u p φ c λ â r g κ λ ʈ ʈ r ʈ ʈ λ ʒ i n r î | Δ r ô ʀ λ a
= g z w π ʝ ê c Δ r ô Δ + b z η r i ʃ ʝ î r z u f λ z
π ʂ j ʈ ʈ r = ʀ | û λ s ʂ m ô m | ʒ η g â | κ ʈ = λ ʈ | l ʃ ô
ø ʀ : r î λ b i f u m ʝ j z v z â j x ʈ r m π i z h λ c ô ô g g
z û + ʀ ʂ ʈ ʈ ʃ ê z g h λ ʈ ʈ | ʈ ʈ λ î z t n ʂ | r ô ʝ m â
+ h h r z ô z ñ b s η + : ʒ r κ ʀ r ô ʀ ʈ Δ c û λ g = ñ z κ
ʀ ʒ o o n z ʀ z ʃ ʈ î r π z ê ʝ ñ g = r π g ô Δ z n z ʀ | κ
π η j i ʃ ʝ r î g π u λ b g λ i t b ô û ʀ ʈ ʈ z ô λ e z η ô
| û r c ʀ z φ ô λ b η ʈ m ô

Homophonic Cipher

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

FAIL

But, slight numerical preference for
German

Cipher Characteristics

digraphs:

ɔ ħ 99

ć : 66

ħ ^ 49

: ȷ 48

z ʀ 44

trigraphs:

ɔ ħ ^ 47

ć : ȷ 23

η ɔ ħ 22

ÿ ɔ ħ 18

ħ ć | 17

tendencies:

â, ê, î, ô, û followed by ʒ and j

â, ê, î, ô, û preceded by z and π



?



?

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 ħ

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

So guess: ɔ = C, ħ = H

One Thing Leads to Another

$$\eta = CH \rightarrow \eta \wedge = CHT \rightarrow \wedge = T ?$$

Each step is guesswork.

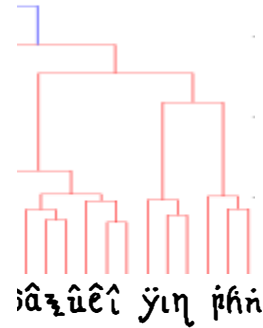
Must be willing to retract.

Weird task, not knowing German.

No longer care what the book says.

Cluster diagram crucial:

$$\ddot{y} = | \rightarrow \iota = | , \eta = |$$



Spring Break
2011



Spring Break 2011

German letters

Cipher letters, in groups

Grid

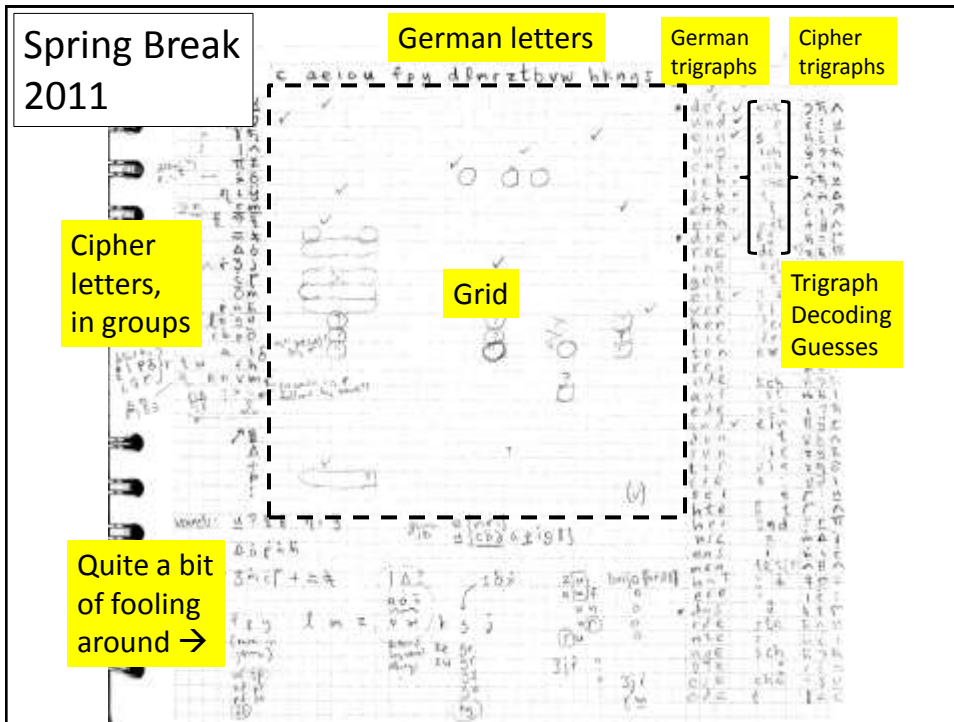
Spring Break 2011

German letters

Cipher letters, in groups

Grid

Quite a bit of fooling around →



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

κ μ̂ ρ: ρ z i o f | y o h e i h ε ι λ η π α z b Δ g z =
 j l z u p q c l a r g κ λ h η r h η λ ε i n r i | Δ r δ ρ λ a
 = g z w π y e c Δ r δ Δ + b z η r i x y i r z u f λ z
 π κ j r δ r = f | u λ s κ m δ m | ε η g â | κ h = λ h | l x o
 ρ f : r i λ b i f u m y j z v z â j x r m π i z h λ c δ o g g
 z ũ + p κ η η x e z o g h λ h η i h η λ i z t n κ | r o y m â
 + h h r z o z n b s η + : ε r κ r p δ ρ h Δ c ũ λ g = n z κ
 p ε o o n z ρ z f h i r p z e y n g = r p g δ Δ z n z p | κ
 π η j x y r i g π u λ b g λ i t b δ ũ f η h z o λ e z η o
 | ũ r c f z o δ λ b η h m o

Actually, those are space bars

Copiale Decipherment

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



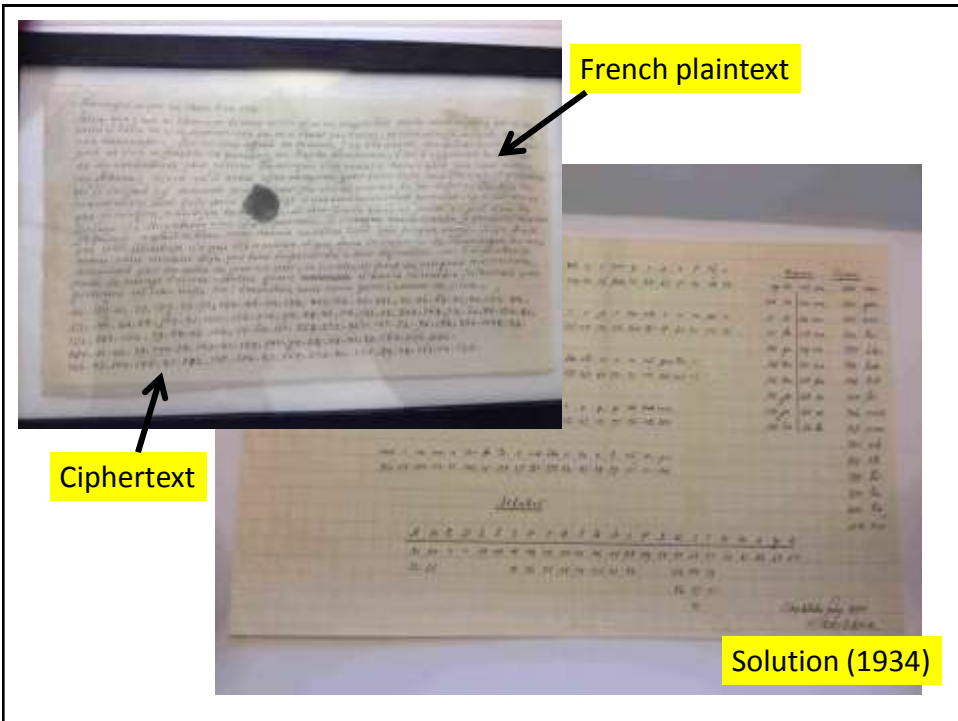
gesetz buchs
 der hochehleuchte ◊ e ⊙
 geheimer theil.
 erster abschnitt
 geheimer unterricht vor die gesellen.
 erster titul.
 ceremonien der aufnahme.
 wenn die sicherheit der Δ durch den ältern
 thürhüter besorget und die Δ vom dirigirenden Λ
 durch aufsetzung seines huths geöffnet ist wird der
 candidat von dem jüngern thürhüter aus einem andern
 zimmer abgeholt und bey der hand ein und vor des
 dirigirenden Λ tisch geführt dieser frägt ihn:
 erstlich ob er begehre ◊ zu werden
 zweytens denen verordnungen der ⊙ sich
 unterwerffen und ohne widerspenstigkeit die lehrzeit
 ausstehen wolle.
 drittens die Δ der ⊙ zu verschweigen und dazu
 auf das verbindlichste sich anheischig zu machen
 gesinnet sey.
 der candidat antwortet ja.

Copiale Decipherment

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



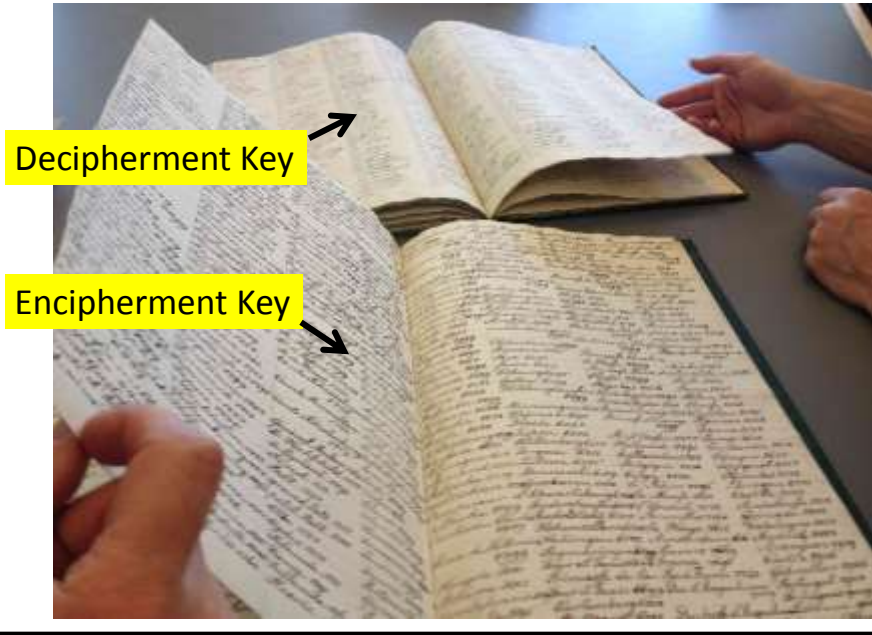
First lawbook
 of the ◊ e ⊙
 Secret part.
 First section
 Secret teachings for apprentices.
 First title.
 Initiation rite.
 If the safety of the Δ is guaranteed, and the Δ is
 opened by the chief Λ, by putting on his hat, the
 candidate is fetched from another room by the
 younger doorman and by the hand is led in and to the
 table of the chief Λ, who asks him:
 First, if he desires to become ◊.
 Secondly, if he submits to the rules of the ⊙ and
 without rebelliousness suffer through the time of
 apprenticeship.
 Thirdly, be silent about the Δ of the ⊙ and
 furthermore be willing to offer himself to volunteer
 in the most committed way.
 The candidate answers yes.



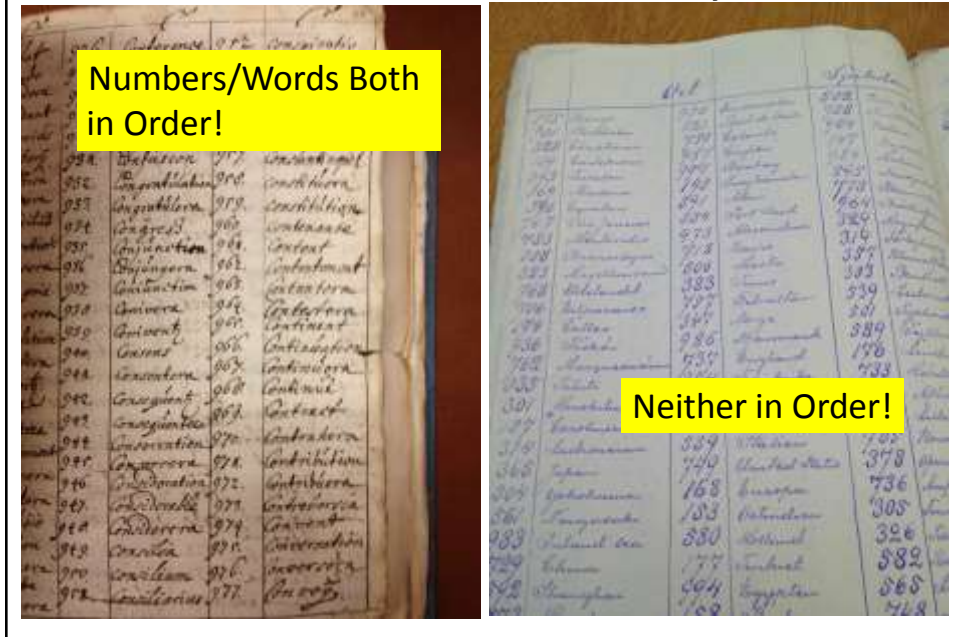
Word Substitution Encipherment Key



Word Substitution Keys



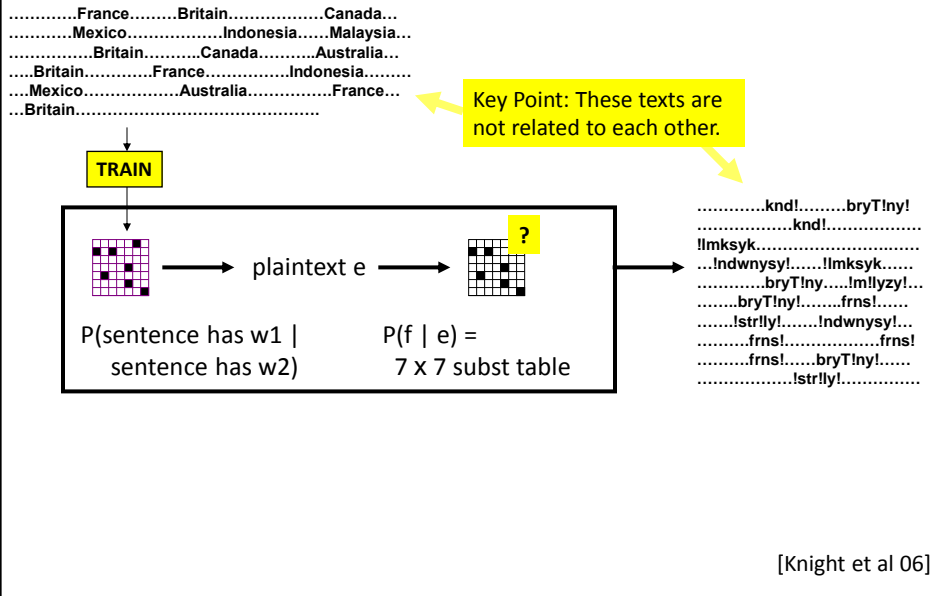
Word Substitution Keys



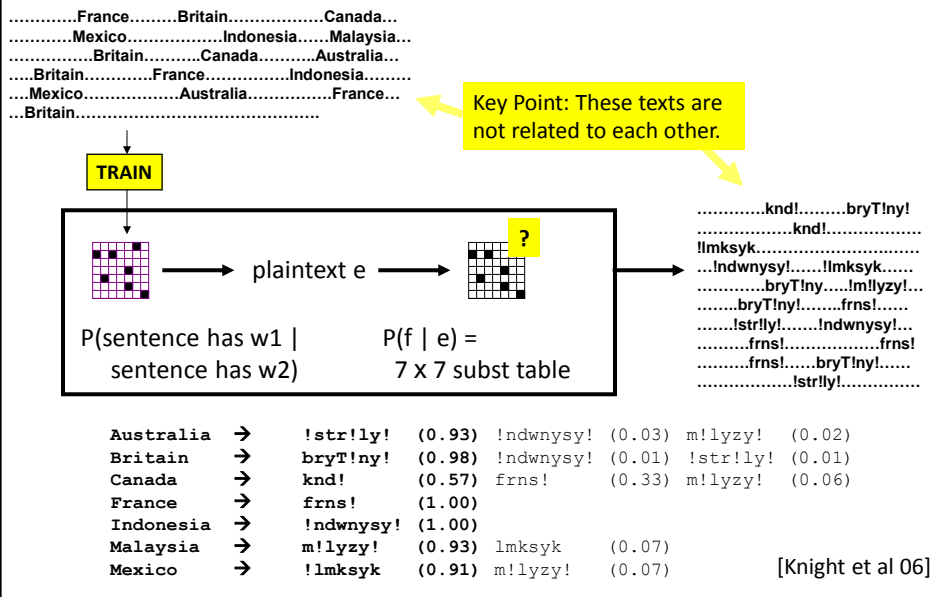
Word Substitution

- 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)



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 x 100,000

Word Substitution (Giga-scale)

- Gibbs sampling fixes memory problem

Cipher: 24234 1899 39902 5716 29948 ...

Plain: the man is car are ...

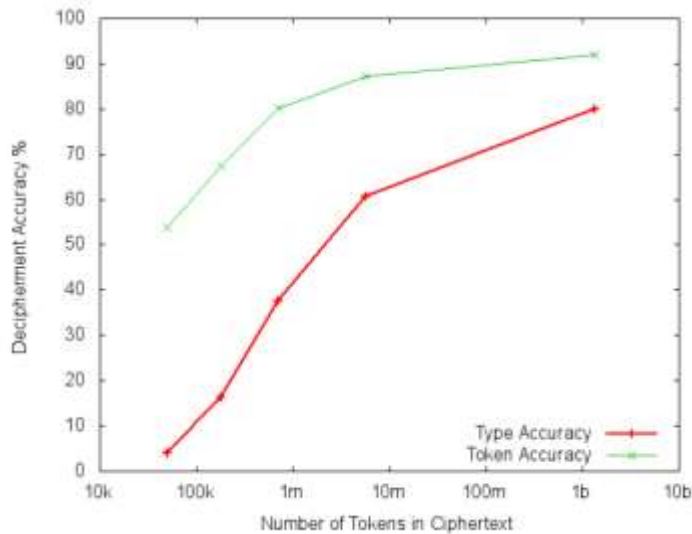
Resample:

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)



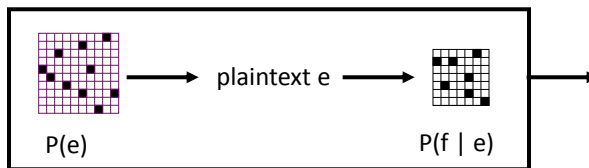
(Dou & Knight 2011)

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



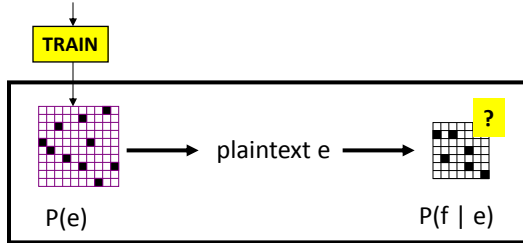
رفعت رئيس السلطة الفلسطينية محمود عباس مجدداً تصويحات وزير الخارجية الإسرائيلي سولمان شالوم التي قال فيها انه يتعين على إسرائيل إعادة النظر في اسمائها من غزة، المقرر ان يتم المصوب، المعلق اذا فازت حركة المقاومة الإسلامية حماس في الانتخابات التشريعية وقال عباس في مؤتمر صحفي على هامش مشاركته في القمة العربية-اللاتينية الأولى انه يتعين على إسرائيل احترام خيار الشعب الفلسطيني حتى لو فازت حماس بالانتخابات، وأضاف "اذا نجحت حماس أو فتح سيكون هذا خيار الشعب الفلسطيني، وعلى الجميع قبول هذا الخيار بكل ترحاب".

من جانبه شجب رئيس الحكومة الفلسطينية أحمد قريع الطابع الأيديولوجي الحزبي للانحياز الإسرائيلي من غزة، وقد ان إسرائيل تريد معاقبة هذا الأراضي لتعزيز سيطرتها على الضفة الغربية وقال قريع في كلمة له خلال مؤتمر نظمته وزارة الأوقاف في رام الله السبتينحون من غزة ولكننا لا نعرف ما هو شكل هذا الانحياز وماذا سيتركون، وما هو مصير المعابر والحدود، وكل ذلك غامض لأنه قرار أحادي الجانب

Foreign Language as a Cipher

BAGHDAD, Iraq (CNN) -- Six bombings 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 Hawija when a suicide bomber detonated explosives hidden under his clothing, Iraqi officials said. That attack in the town about 130 miles (209 kilometers) north of Baghdad also wounded 30 Iraqis, said Iraqi army Lt. Col. Khaliil al-Zawbai. A car bombing in Saddam Hussein's ancestral homeland of Tikrit also killed 30 Iraqis and wounded another 40, Iraqi officials said. The Tikrit explosion...

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



وقد زعم رئيس السلطة الفلسطينية محمود عباس مجدداً تسميحات وزير الخارجية الإسرائيلي سيلفان شالوم التي قال فيها إنه يتعين على إسرائيل إعادة النظر في سحبها من غزة، المقرر أن يتم السبت المقبل إذا فارت حركة المقاومة الإسلامية حماس في الانتخابات التشريعية وقال عباس في مؤتمر صحفي على هامش مشاركته في القمة العربية-اللاتينية الأولى إنه يتعين على إسرائيل احترام خيار الشعب الفلسطيني حتى أو فارت حماس بالانتخابات، وأضاف "إذا نجحت حماس أو فتح سيكون هذا خيار الشعب الفلسطيني، وعلى الجميع قبول هذا الخيار بكل ترحاب".

من جانبه شجب رئيس الحكومة الفلسطينية أحمد قريع الطابع الأحادي الجانب للاتسحاب الإسرائيلي من غزة، وأكد أن إسرائيل تريد مغادرة هذه الأراضي لتعزيز سيطرتها على الضفة الغربية.

وقال قريع في كلمة له خلال مؤتمر لطلته ووزارة الأوقاف في رام الله "سينسحبون من غزة ولكننا لا نعرف ما هو شكل هذا الاتسحاب وماذا سيتركون، وما هو مصير المعابر والحدود، وكل ذلك غامض لأنه قرار أحادي الجانب".

Time Expressions

!!@!m
!lywm
!lth!ny&
!!@!m !l!m!Dy
Sfr
@!m
th!ny&
@!m 1992
@!m 1993
ywm
!!!sbw@ !l!m!Dy
fy !ldqyq&
!lsn& !l!ry&
!lsn&
!lsh=hr !l!m!Dy
!lsh=hr !l!ry
snw!t
sn&
=hdh! !!@!m
s!@&
!!@Sr
@!m 1991

@!m 1990
w!lth!ny&
fy !lywm
mn !lsh=hr !l!ry
!lqrn
!'y!m
@!m!aN
!s!@&
17 shb!T 1994
th!lth snw!t
dqyq&
=hdh=h !lsn&
ywmyn
mn !!@!m !l!m!Dy
!lsn& !lmqbl&
fy !lsn&
kl ywm
fy !!@!m !l!m!Dy

!!@Swr
=hdh! !lsh=hr
fy ywm
nys!n
!sbw@
=hdh=h !!!'y!m
qbl !'y!m
fy !!@Sr
mn !lsn&
!lsnw!t
b@d ywm
!!'y!m
13 nys!n 1994
!lth!ny& @shr&
th!lth& !y!m
qbl !sbw@yn
fy !lywm !t!ly
sh@b!n
tmwz
3 dhw !lHj& 1414
fy shb!T !l!m!Dy
qbl ywmyn

Time Expressions

!!@!m
!lywm
!lth!ny&
!!@!m !l!m!Dy
Sfr
@!m
th!ny&
@!m 1992
@!m 1993
ywm
!!sbw@ !l!m!Dy
fy !ldqyq&
!lsn& !lj!ry&
!lsn&
!lsh=hr !l!m!Dy
!lsh=hr !lj!ry
snw!t
sn&
=hdh! !l!m
s!@&
!!@Sr
@!m 1991

@!m 1990
w!lth!ny&
fy !lywm
mn !lsh=hr !lj!ry
!lqrn
!y!m
@!m!aN
!s!@&
17 shb!T 1994
!l!th snw!t
dqyq&
=hdh=h !lsn&
ywmy
mn !l!m !l!m!Dy
!lsn& !lmqbl&
fy !lsn&
kl ywm
fy !l!m !l!m!Dy

!!@Swr
=hdh! !lsh=hr
fy ywm
nys!n
!sbw@
=hdh=h !!y!m
qbl !y!m
fy !!@Sr
mn !lsn&
!lsnw!t
b@d ywm
!ly!m
13 nys!n 1994
!lth!ny& @ch&
th!lth& !y!m
qbl !sbw@yn
fy !lywm !l!ly
sh@b!n
tmwz
3 dhw !lHj& 1414
fy shb!T !l!m!Dy
qbl ywmy

Time Expressions

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

9 Hzyr!n 1942	27 tmwz 1993	21 Hzyr!n 1967
8 tshryn !!wl 1990	26 tmwz 1953	20 !y!r 1990
7 k!nwn !!wl 1993	26 shb!T 1993	20 tshryn !'wl 1983
6 !y!r 1993	26 k!nwn !!wl 1994	20 tshryn !!'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 tshryn !!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 Expr

<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

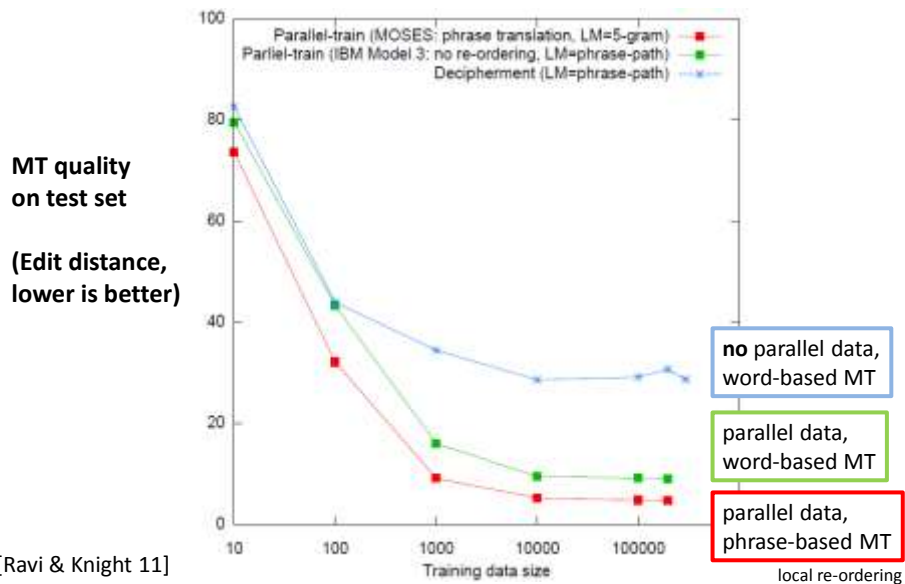
229	fy Hzyr!n !l!m!Dy	16	n=h!y& Hzyr!n !l!m!Dy
207	fy Hzyr!n	16	fy Hzyr!n 1990
75	fy Hzyr!n !lmqbl	15	sh=hr Hzyr!n
61	fy Hzyr!n 1993	15	fy sh=hr Hzyr!n !l!m!Dy
31	fy Hzyr!n 1992	15	fy Hzyr!n 1994
27	!l!r!b@ mn Hzyr!n	14	mn 17 Hzyr!n
27	fy Hzyr!n 1967	14	fy Hzyr!n 1996
19	fy 30 Hzyr!n !l!m!Dy	14	fy 30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

Time Expressions

Hzyr!n

229	fy Hzyr!n !!m!Dy	16	n=h!y& Hzyr!n !!m!Dy
207	fy Hzyr!n	16	fy Hzyr!n 1990
75	fy Hzyr!n !!mqbl	15	sh=hr Hzyr!n
61	fy Hzyr!n 1993	15	fy sh=hr Hzyr!n !!m!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 !!m!Dy	14	fy 30 Hzyr!n
18	fy n=h!y& Hzyr!n !!m!Dy	13	fy sh=hr Hzyr!n
18	fy Hzyr!n 1991	13	fy 20 Hzyr!n !!m!Dy
17	mn Hzyr!n	13	4 Hzyr!n 1967
17	mdh Hzyr!n !!m!Dy	12	n=h!y& Hzyr!n
17	4 Hzyr!n	12	!lr!b@ mn Hzyr!n 1967

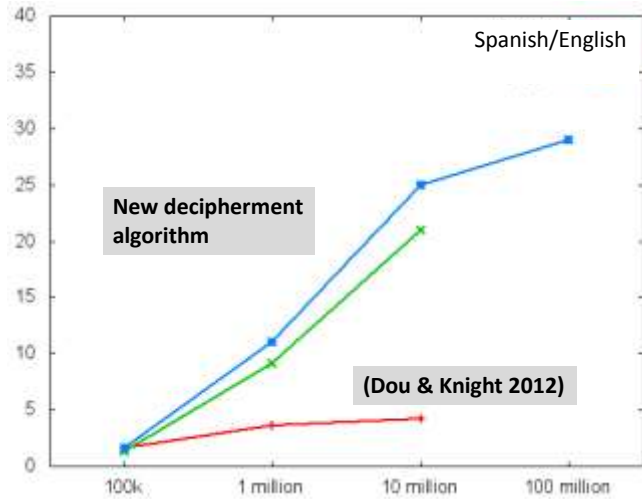
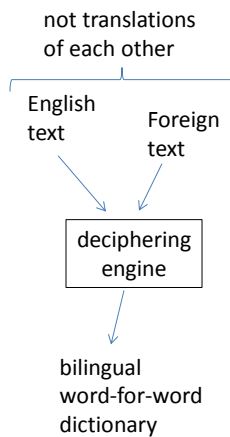
Deciphering Spanish Time Expressions



Deciphering Foreign Language at Giga-Scale

(Dou & Knight subm.)

Accuracy of learned
bilingual dictionary



How much foreign text (running words)

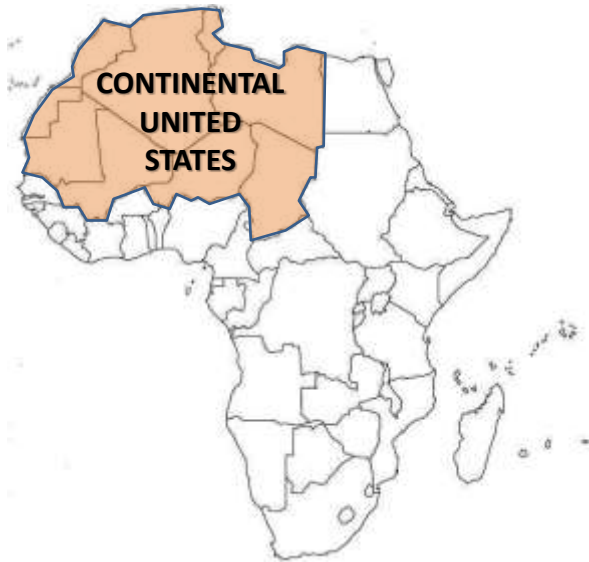
Practical Value

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

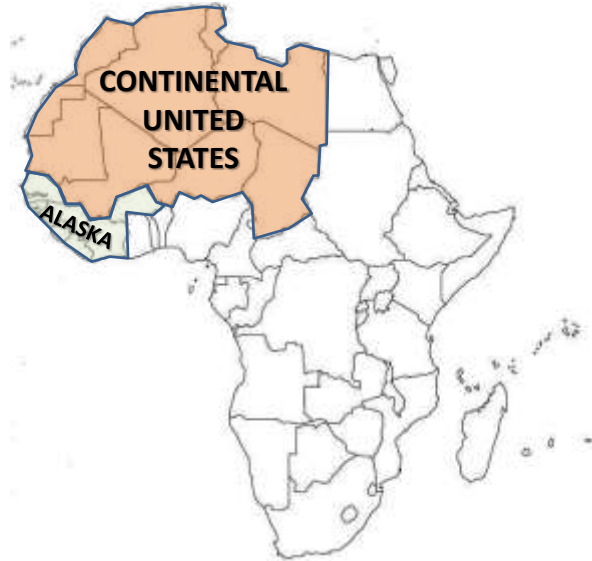
African Languages



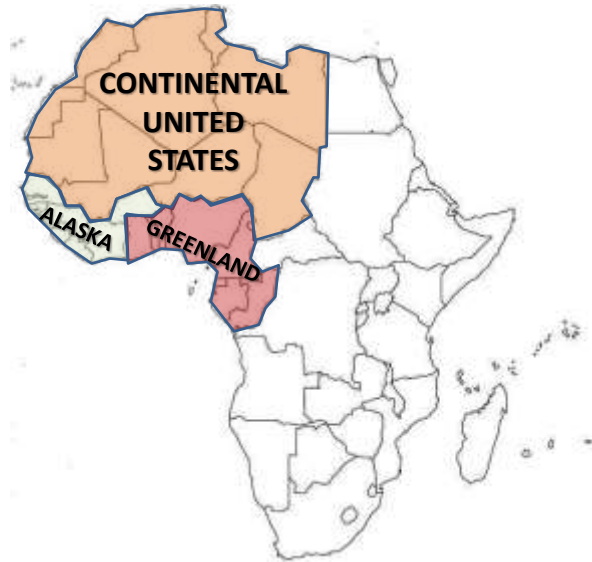
African Languages



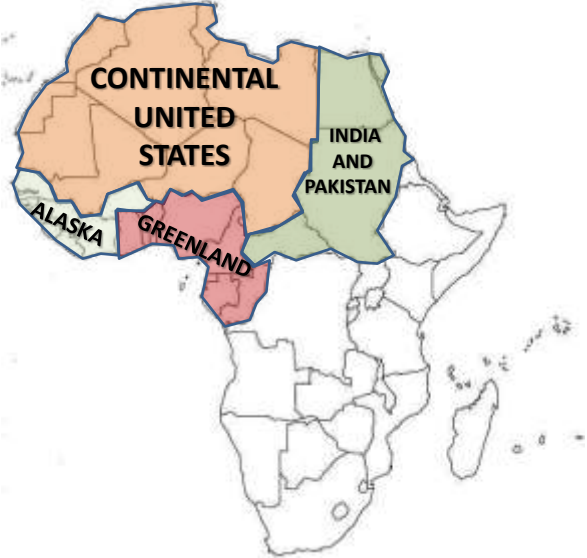
African Languages



African Languages



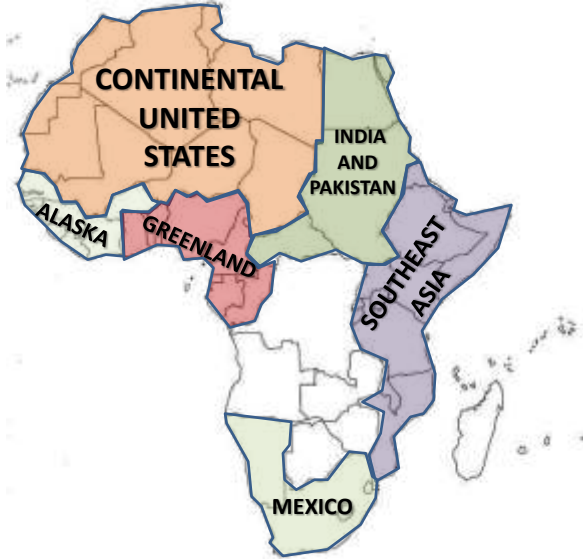
African Languages



African Languages



African Languages



African Languages

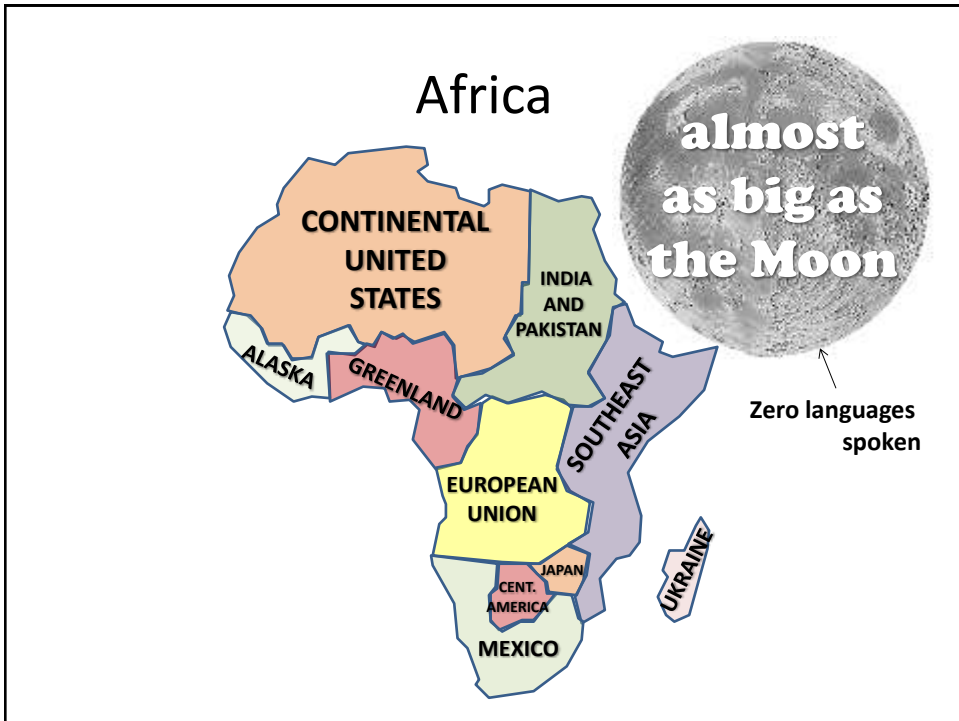


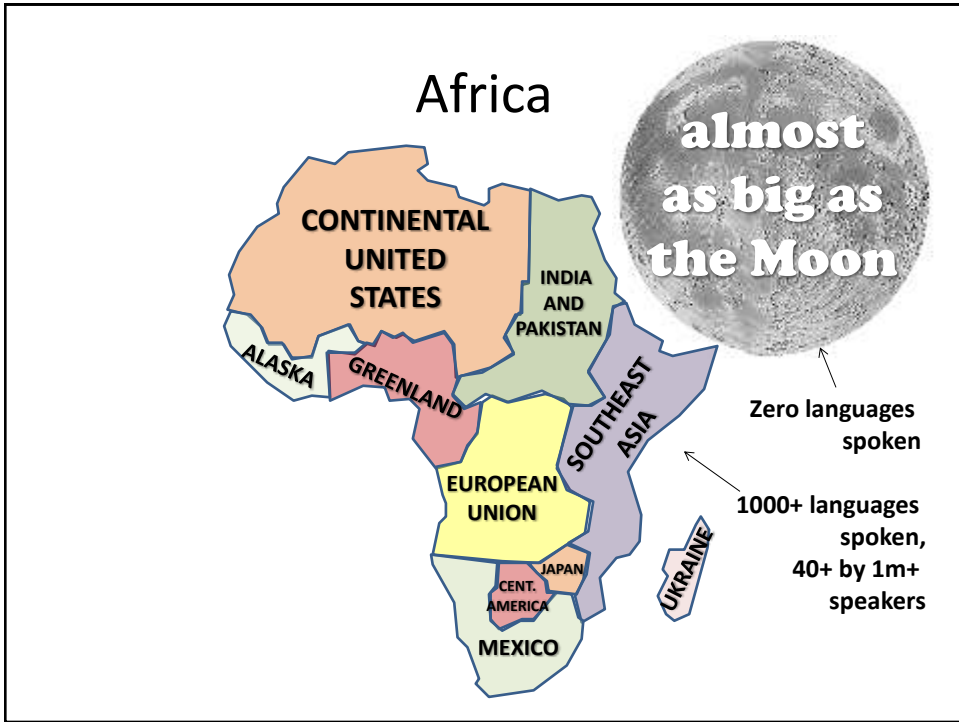
African Languages



African Languages







Unsolved ciphers

Voynich Manuscript (VMS)



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



Voynich Manuscript (VMS)



72289 90119 90119 90119 90119 90119 90119 90119 90119 90119
 72289 90119 90119 90119 90119 90119 90119 90119 90119 90119
 82289 90119 90119 90119 90119 90119 90119 90119 90119 90119

PTC80X 0FFC9 40X 1FC89 40FFC9 40FF9 7C77C9 4077C9 91A0 0FFA29
 72289 40FFC9 40FFA0 289 40FFC9 7C89 40FFC9 40FFC9 97C89 90FF9
 82289 40FFC9 8A77C89 40FFC9 40FF9

BSC8AE OPCC9 4OE FCC89 4OFCC9 4OP9 SCBS9 4OBSC9 EFAM OPAE29
 2ZC9 4OFCC9 4OFAM Z89 4OFCC9 SC89 4OFCC9 4OFCC9 ESC89 EOP9
 8ZC9 4OPCCC9 8ARSC89 4OFC9 4OP9



“Herbal” section



Many pictures look like grafting.



Sunflower? Would date VMS as post-1492.

“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 Tepenez signs VMS in Bohemian court	1865	WV born in Lithuania
1630s	George Baresch owns VMS 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 Tepenez signature
		1919	WV writes to Bohemian State Archvs
		1921	WV presents VMS + inserted Marci letter mentioning Francis Bacon, asks \$160k
		1921	Newbold & WV announce decipherment
		1930	WV dies. VMS placed in vault, \$100k
		1931	VMS appraised at \$19,400
		1960	Ethel dies, VMS to secretary Ann Nill "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	R	S	T	U	V	X	Z
A	V	Z	B	F	G	L	M	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											N									
L																				
M												A								
N																				
O																				
P											N									
Q																				
R																				
S																				
T																				
U																				
V																				
X																				
Z																				

22x22 table

Encoding:

A → CC, OM, ...
 B → ...
 ...
 N → HA, MI, DO, NU ...
 ...
 Z → ...

Decoding:

...
 DO → N
 ...

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

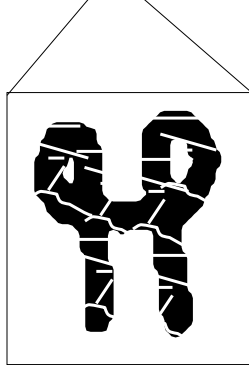
Hmm, by the method, both plaintext **and ciphertext** should be in Latin letters...

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



William Newbold,
Polymath, PhD UPenn

... 409cc89 ...

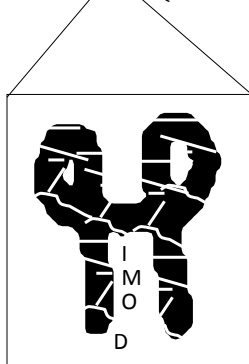


apparent
ciphertext




William Newbold,
Polymath, PhD UPenn

... 409cc89 ...



apparent
ciphertext

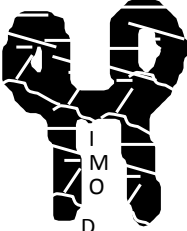
real
ciphertext:
DOMI...



Let's Decipher with Newbold !

Ɔcc89 ...

apparent ciphertext



real ciphertext

DOMI...

DO OM MI ...

OM DO MI ...

a n n ...


o n n ...

doubling

non-deterministic anagramming

lookup in 22² table

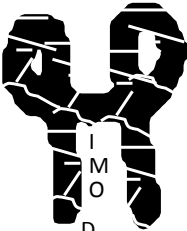
non-deterministic mapping from 11 Latin letters to full 22



Let's Decipher with Newbold !

Ɔcc89 ...

apparent ciphertext



real ciphertext

DOMI...

DO OM MI ...

OM DO MI ...

a n n ...

o n n ...

doubling

non-deterministic anagramming

lookup in 22² table

non-deterministic mapping from 11 Latin letters to full 22

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

22x22 table
(values guessed)

Alphabet: Currier/D'Imperio Transcription

c c c C S Z	P P P P P F B V	Q Q Q Q Q X W Y
J a x r o \ v J A E R O I D	6 7 8 9 4 ? 6 7 8 9 4 2	
G H 1 G H 1	T U 0 T U 0	N M 3 N M 3
K L 5 K L 5		

VMS Letters

count	letter	count	letter	count	letter
25468	O	o	2886	2	?
20227	C	c	1752	N	U
17655	9	9	1413	B	6
14281	A	a	1046	J	Y
12973	8	8	950	Q	K
11008	S	s	908	X	G
10471	E	e	591	T	L
10026	F	f	524	*	H
6716	R	r	431	V	1
5994	P	p	316	I	5
5423	4	4	217	W	0
4501	Z	z	157	D	
4076	M	m	156	3	

Total
 63k character tokens

VMS Words

count	word		count	word		count	word	
863	8AM	8aW	212	OFAM	oIfaW	140	OPCC9	oIfcc9
537	OE	oR	211	8AN	8aW	138	OFAE	oIfaR
501	SC89	rc89	191	4OFAE	4oIfaR	130	ZO	co
469	AM	aW	186	ZOE	coR	129	OFAR	oIfaR
426	ZC89	co89	177	OFCC9	oIfcc9	119	ESC89	rc89
396	SOE	coR	174	SCC9	rc89	118	OFC89	oIfcc9
363	OR	oR	172	SCOE	rc89			
350	AR	oR	155	S9	cc9	+ many more!		
344	SC9	aR	155	OPC89	oIfcc9			
318	8AR	rc89	154	OPAM	oIfaW			
308	4OFCC9	8aR	152	4OFAR	4oIfaR			
305	4OFCC89	4oIfcc9	151	9	9			
283	ZC9	4oIfcc89	151	4OE	4oR			
279	4OFAN	co89	150	S89	cc89			
272	4OFC89	4oIfaW	147	4OF9	4oIf9			
270	89	4oIfcc89	144	ZCC9	co89			
262	4OFAM	89	144	OFAN	oIfaW			
260	AE	4oIfaW	144	2AM	caW			
253	8AE	89	143	OPAE	oIfaR			
243	Z	oR	141	OPAR	oIfaR			
219	SOR	8aR	140	SX9	rc89			
		?						
		coR						

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
... 4oIfcc89 4oIfcc89 ...
- 8 distinct words appear tripled
... 4oIfcc89 4oIfcc89 4oIfcc89 ...
... coR coR coR ...
... co89 co89 co89 ...
... oIfaW oIfaW oIfaW ...
... oR oR oR ...
... 9IfaW 9IfaW 9IfaW ...
... 8aW 8aW 8aW ...
... 4oIfcc89 4oIfcc89 4oIfcc89 ...

Substitution Cipher?

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

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

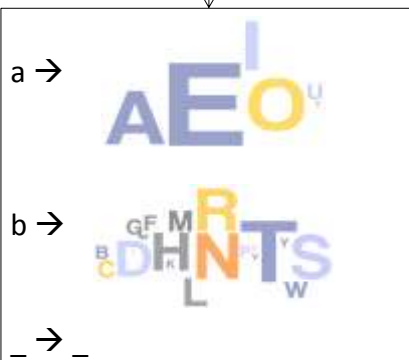


- Tried 80+ languages written without vowels.

Letter Clustering

Trigram model over {a, b, _}

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

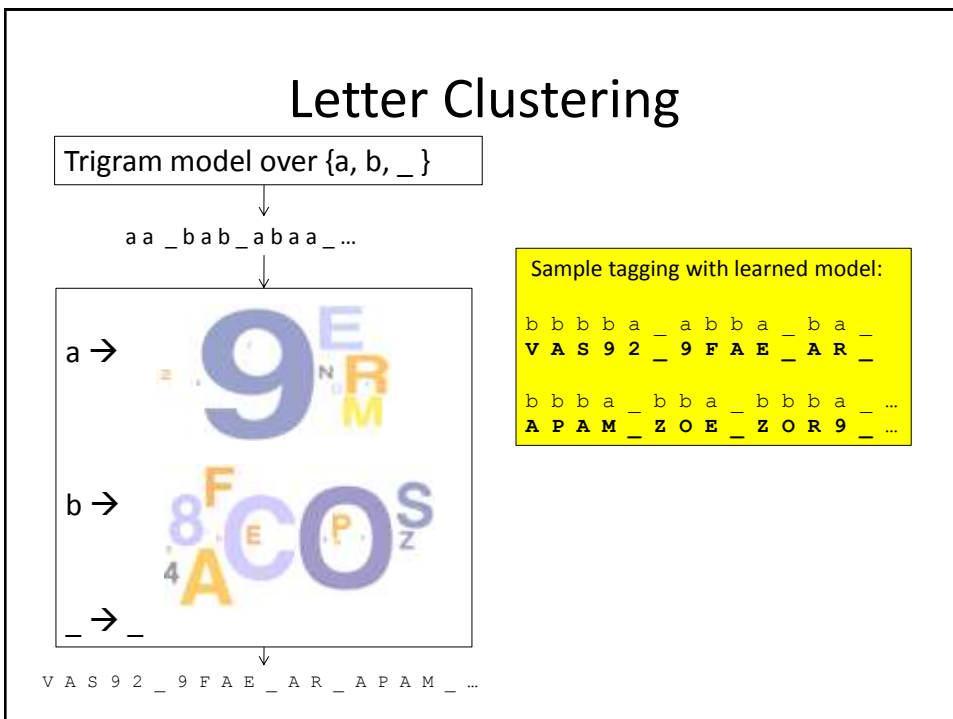
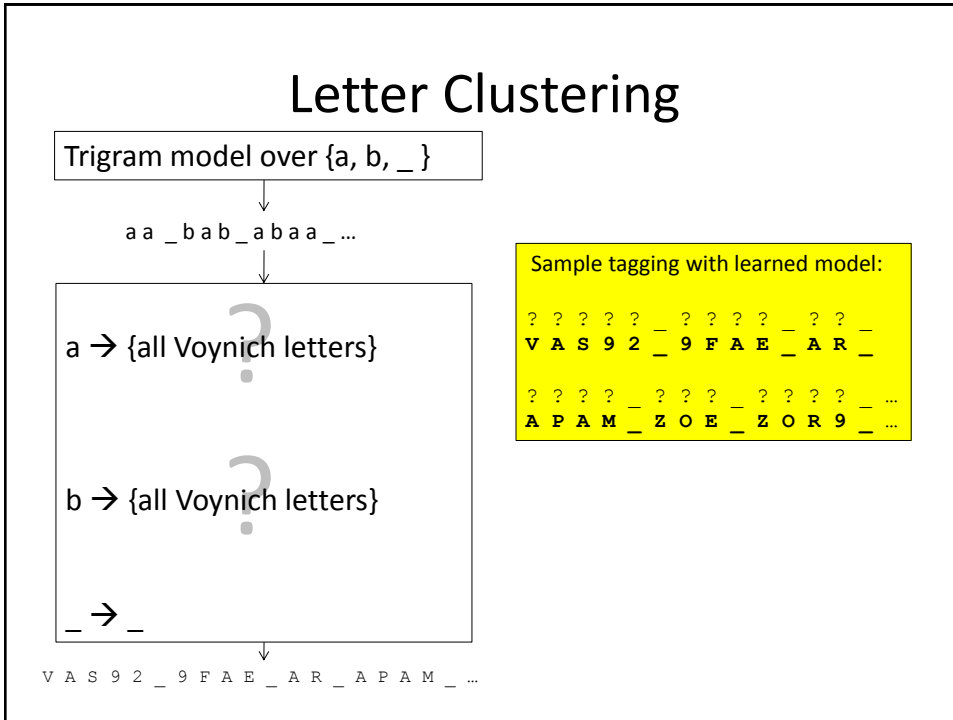


i n _ t h e _ t o w n _ w h e r e _ i _ w a s ...

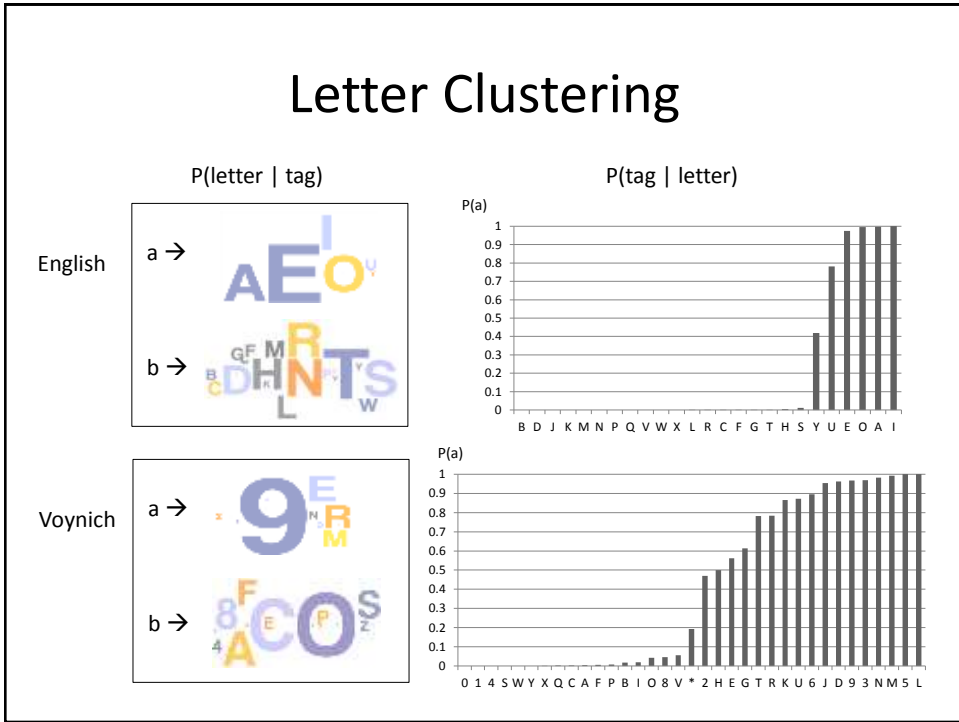
Sample tagging with learned model:

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



Word Clustering

Bigram model over {a, b}

a a b a b a b a a ...



VAS92 9FAE AR APAM ZOE ZOR9 QRC2 9 ...

Sample tagging with learned model:

a a a a a

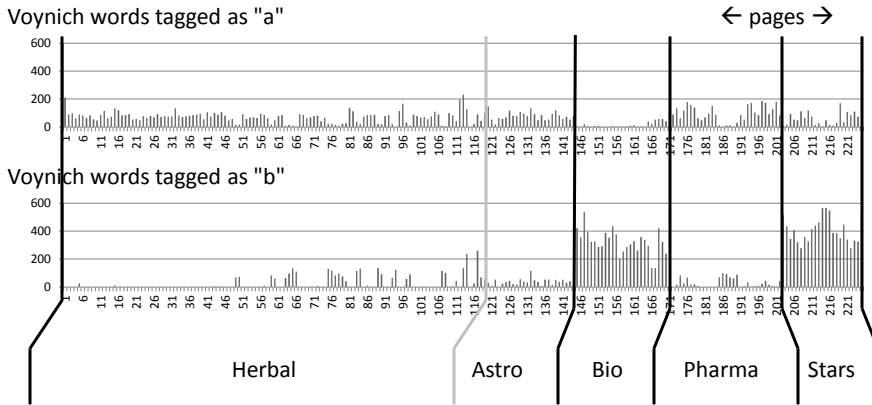
VAS92 9FAE AR APAM ZOE

a a a a a ...

ZOR9 QRC2 9 FOR ZOE89 ...



Word Clustering



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

Approved for Release by NSA on 08-03-2009, E.O. 13526, Case # 58742

An Application of PTAH to the Voynich Manuscript (U)

BY MARY E. D'IMPERIO

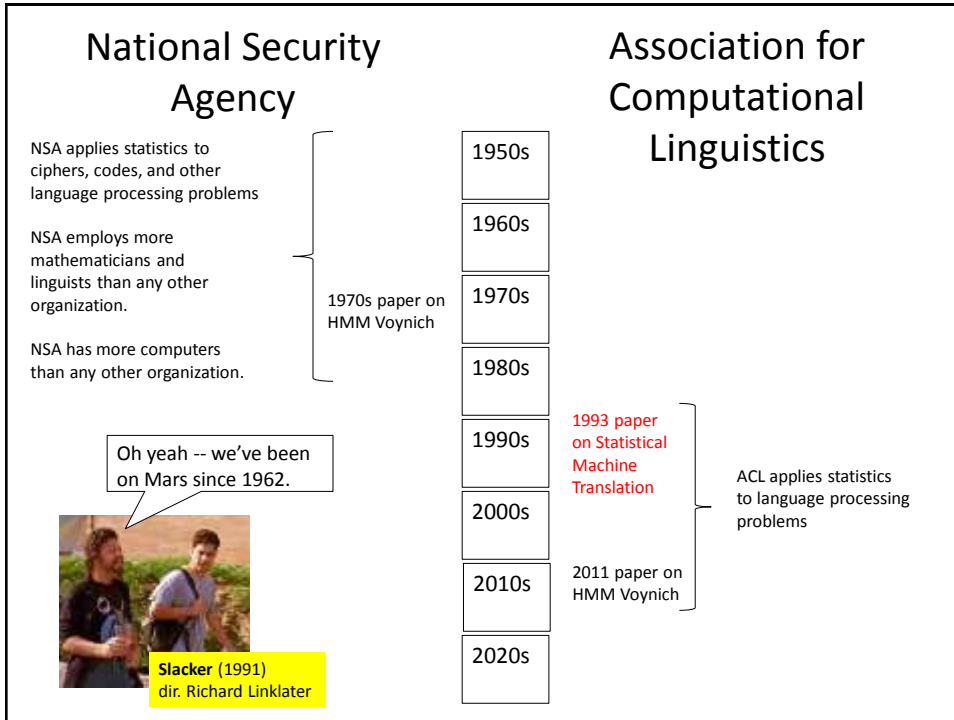
~~Top Secret Umbra~~

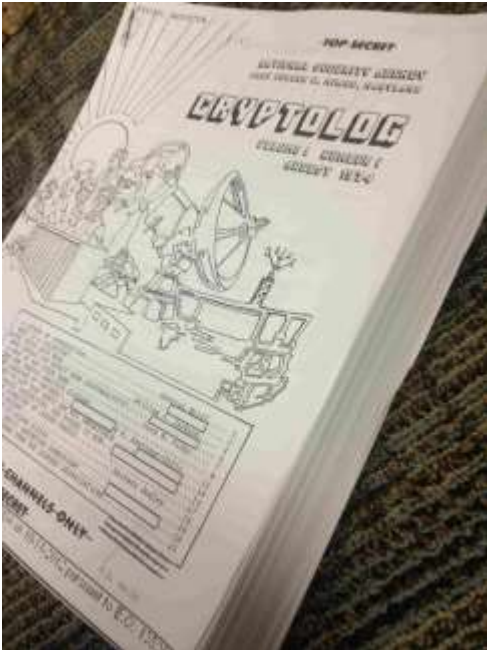
(U) This article is the second in a series of studies applying some modern statistical techniques to the problems posed by the Voynich manuscript. This study attempts to discover and demonstrate regularities of patterning in the Voynich text subjectively noted by many earlier students of the manuscript. 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.

... is a general Analyses), paper in the times. Mr. ling on his program. He was struck by the passage "immenso Prah 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 [redacted] of R51, the name is pronounced "however you like" [8]. The technique itself and its uses are classified Top Secret Codeword.

[redacted] 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 [redacted]

1970s National Security Agency report recently declassified!





CRYPTOLOG

NSA newsletter declassified in 2013.

4400 pages (1974-1997).
238 Mbyte PDF file.

Covers intelligence gathering, linguistics, military, cryptography, office space, pay grades, human factors, etc.

Heavily redacted.

CRYPTOLOG: Voynich

DOCID: 4009723

UNCLASSIFIED

The Voynich Manuscript

When a newspaper editor sends a letter to our Chicago full back on the back page...
...the Voynich manuscript "myst"? No, is it a hoax? No. What is it, then? A...
...That was my reaction the first time I looked at it closely, but faced with all the profound...
...However, a recent reworking of Elizabeth Friedman's article in the Washington Post (Sept 8, 1942) and of Benjamin Titman's paper in the NSA Technical Journal (Summer 1947), plus some observations I have seen in the magazine...
...The history of the manuscript, which has been detailed in other places, needs only passing mention since it does not throw any light on the content. During the 18th and 19th centuries, it was said by James Huxley, author of the "History of the University of Prague," to have belonged to one John de Sponcor (died 1111) (1376-1412). Huxley writes in 1866 to the Rev. John...
...The author of the manuscript, the author of which, he had heard from another source, was the great natural scholar Roger Bacon. (This...
...Huxley himself withheld judgment on the attribution, but an Irish poet scholar about 1840...
...The information in this paragraph and the preceding paragraph was taken from Johnston, January 1945 (NSA, V. 10, 5).

FOUO: 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 Revisited P16

The Voynich Manuscript, an object of interest...
...The history of the manuscript, which has been detailed in other places, needs only passing mention since it does not throw any light on the content. During the 18th and 19th centuries, it was said by James Huxley, author of the "History of the University of Prague," to have belonged to one John de Sponcor (died 1111) (1376-1412). Huxley writes in 1866 to the Rev. John...
...The author of the manuscript, the author of which, he had heard from another source, was the great natural scholar Roger Bacon. (This...
...Huxley himself withheld judgment on the attribution, but an Irish poet scholar about 1840...
...The information in this paragraph and the preceding paragraph was taken from Johnston, January 1945 (NSA, V. 10, 5).

...the Voynich manuscript...
...The information in this paragraph and the preceding paragraph was taken from Johnston, January 1945 (NSA, V. 10, 5).

April 78 • CRYPTOLOG • Page 21

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 wasn't as easy and rapid as had been anticipated, however, and 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. CILT (Chinese University Language Translator) in Hong Kong translates two periodicals into English. And a system was developed by a U.S. company for FIP 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 post-editing 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, save without extensive pre-editing or post-editing. 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...
...P...
...M...
...P...

DOCID: 4510113
TOP SECRET COMINT
MACHINE TRANSLATION: What can it do for us?
...
TOP SECRET COMINT

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 Beary Tetrault and myself, with many valuable suggestions from our colleagues on Professional Qualification Examination (PQE) Committees 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 inv languages, mainly Indo-Europe other families. The results aging enough in both instance send its use in the PQE Handb

CRYPTOLOG: Linguists

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

At the ...
foreign lan
their

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 incessant complaints about lack of recognition, few specific suggestions were made regarding how that reco

TEACHING COMPUTER SCIENCE TO LINGUISTS

by [redacted] PI6

He entered ...
with foreign ...
or perhaps to ...
amounts of ...
of ambiguous ...
ble rules of ...
erent in the ...
into another ...
rign or the ...
the physical ...
m. He felt at ...
it is so foreign ...
SA.

12. PUBLICATIONS (24% total; 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 Services

... col sep
ling sys
see inc
ord
axi
pro
fal
ext

Back to Word Clustering

Bigram model over {a, b}

a a b a b a b a a ...



VAS92 9FAE AR APAM ZOE ZOR9 QRC2 9 ...

Let's try 10 clusters.

Let's limit ourselves to the more homogenous Bio + Stars sections.

10-Class Word Clustering: English



etc

etc

etc

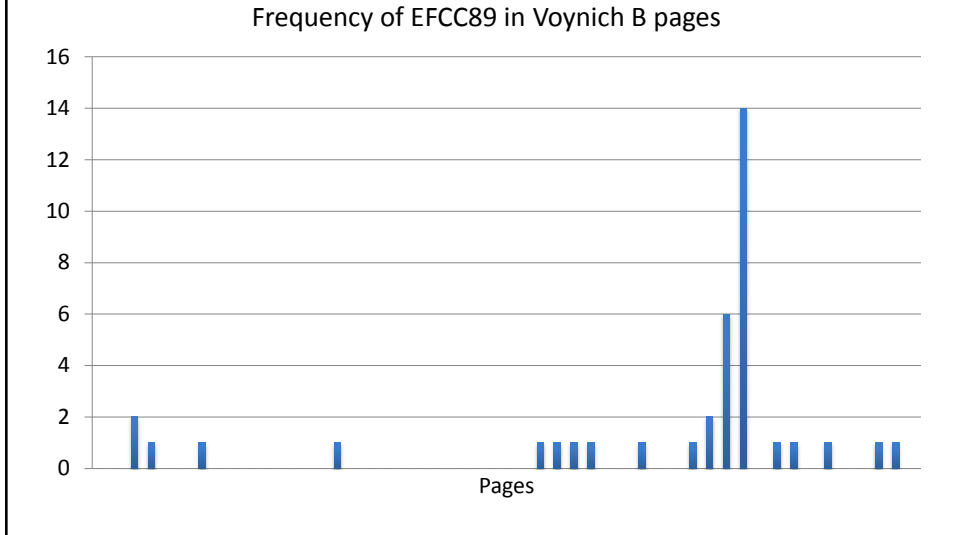


etc

etc

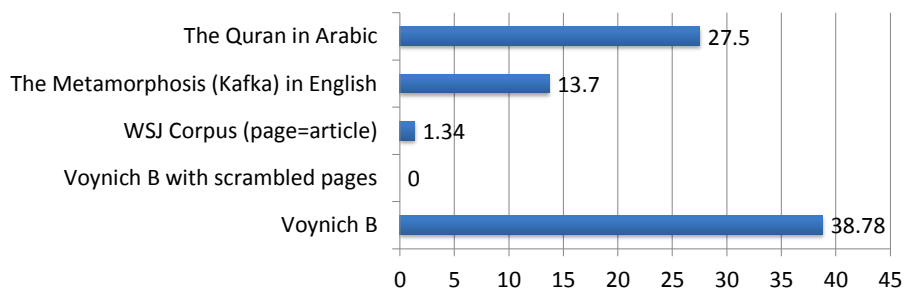
etc

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!

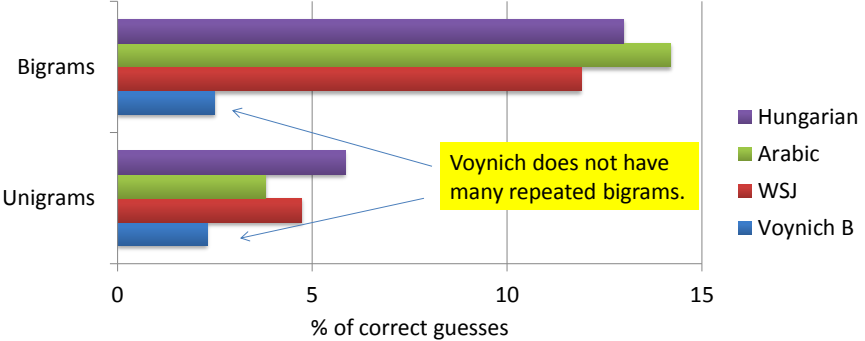
The line is a functional entity...



Prescott H. Cutler

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
Unigrams	~6%	~4%	~5%	~3%
Bigrams	~13%	~14%	~12%	~3%

Zodiac Killer Ciphers

Zodiac 408 (solved, 1969)



Zodiac 340 (still unsolved)



Zodiac Serial Killer

408-letter cipher (solved):

I Δ L I P K E Z / U B O R K P X K B
 W V + E G Y F D Δ H P Q K X P Y E
 M Z / Λ U I Z F Δ B T L N G Y D Q E
 S Φ / Δ M B P O R A U Q T R J P E
 X Λ L M Z J O R \ / Q F H V W Z Δ Y
 Q + P G D Δ K I → Φ P X Δ E Φ S Φ
 R N I Z Y E J Q Δ P G B T Q S B
 L O / P B B Q X P E H H U A R R X
 C Z H P S I → W P I A → L H R Δ F
 B P D R + T X Q → N F Z E U H X F
 Z > Q O V W J → + I L → J A R O H
 I Δ D R R O T Y H \ / O X J T G A
 P → M A R U L Δ L M J N A → Z Φ P
 R I Z J X Δ Δ B V W \ + V T L O P
 Φ U Q S A Δ B V W \ + V T L O P
 Λ Π S H J F U Z → O Δ D → G G Q I H
 N X → Z Z E / Δ B Q Z T A P B B V
 Q Z X P W Q P F B Q Y J O Δ A Δ B
 Q O T B R U C + Q O Y O O L A B W
 V Z E G Y K E P T Y A Δ B B L L O
 H E F B X Δ Φ X A P C \ Δ L I T C P
 Q J O B B O S Φ P O R X Q P F B S
 Z O J T L P O Δ T T + J B P G W X
 V I E X R Δ W I Q P E F E

(plus two more sections)

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23
A	J	G	A	S																		
R	V																					
E	B																					
7	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
8	4	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
M	+	W	Q	N	Z	O	B															
F	T	O																				
S	R																					
N	S																					
M	G																					

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 DANGEROUE ANAMAL OF ALL TO KILL SOMETHING GIVES ME THE MOST THRILLING EXPERENCE 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 ”
EBEORIEEMETHHPITI

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	3.3 / 2.6

[Ravi & Knight 11]

See also Malte Nuhn's paper at ACL 2013!

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



[Knight & Yamada 99]

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..."

ciphertext

primera parte
del ingenioso
hidalgo don ...

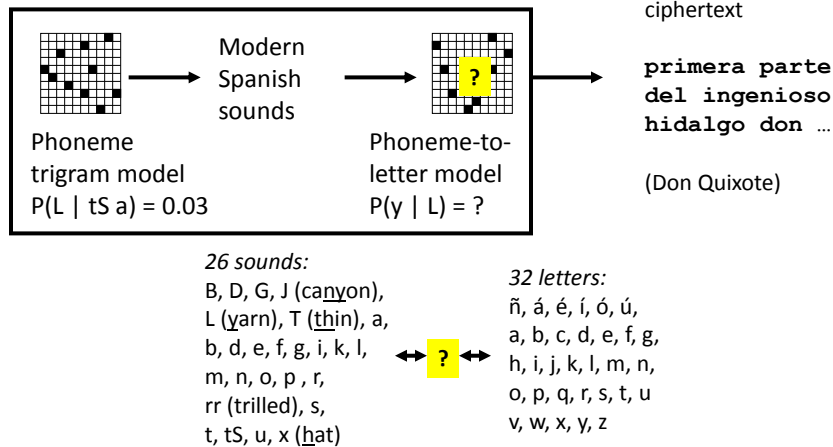
(Don Quixote)



OUR HERO

[Knight & Yamada 99]

Archaeological Decipherment



EM approach = 93% accurate phonetic decipherment

[Knight & 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]

Undeciphered Writing Systems

Undeciphered writing systems

Indus Valley
Script
(3300BC)



Linear A
(1900BC)



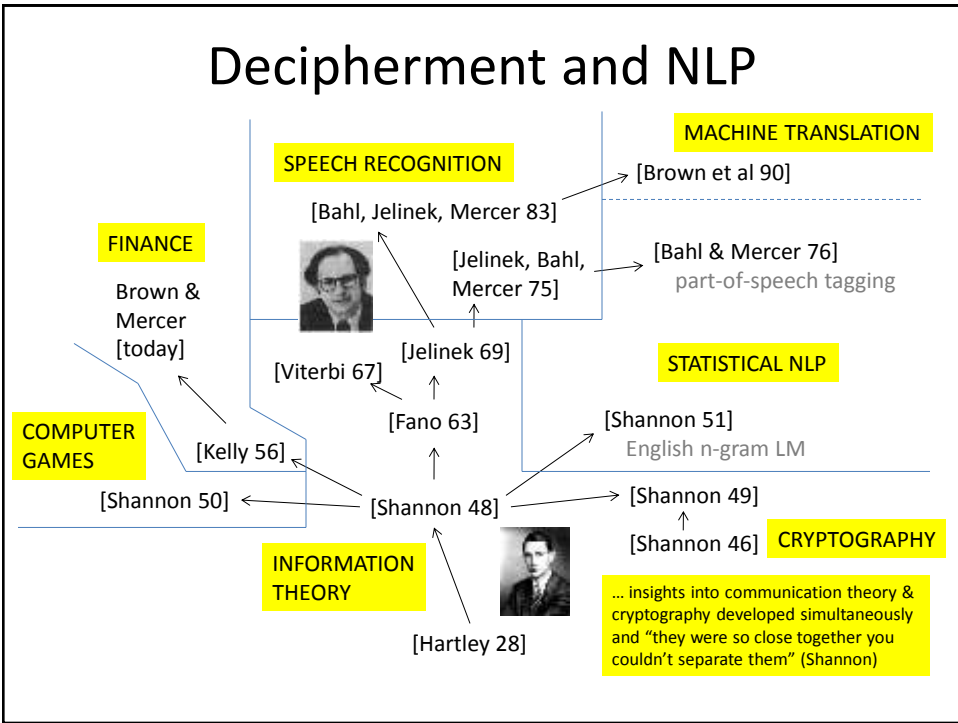
Rongorongo (1800s?)







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