

Language Translation and Code-Breaking

Kevin Knight

Information Sciences Institute & Computer Science Department
University of Southern California



joint work with:

Sujith Ravi, Qing Dou (USC), Beáta Megyesi, Christiane Schaefer (Uppsala)
Regina Barzilay, Ben Snyder (MIT), Sravana Reddy (Chicago)

Stanford University November 4, 2014

Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。

Kowane mutum na da hakkin ya sami yancin yin tunani da na sanin yakamata da na bin addini; saboda haka yana da yancin sake addini ko ra'ayin da ya bada gaskiya gare shi, da kuma yancin nuna addininsa ko ra'ayinsa, shi daya ko a cikin taro kuma a fili ko a boye ta hanyar koyarwa ko yin ibada, ko bauta wa abin da ya bada gaskiya gare shi da yin abubuwan da abin da yake bauta wa din ya nuna masa.

Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。



The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.

Kowane mutum na da hakkin ya sami yancin yin tunani da na sanin yakamata da na bin addini; saboda haka yana da yancin sake addini ko ra'ayin da ya bada gaskiya gare shi, da kuma yancin nuna addininsa ko ra'ayinsa, shi daya ko a cikin taro kuma a fili ko a boye ta hanyar koyarwa ko yin ibada, ko bauta wa abin da ya bada gaskiya gare shi da yin abubuwan da abin da yake bauta wa din ya nuna masa.



Everyone has the right to freedom of thought, conscience and religion; this right includes freedom to change his religion or belief, and freedom, either alone or in community with others and in public or private, to manifest his religion or belief in teaching, practice, worship and observance.

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)

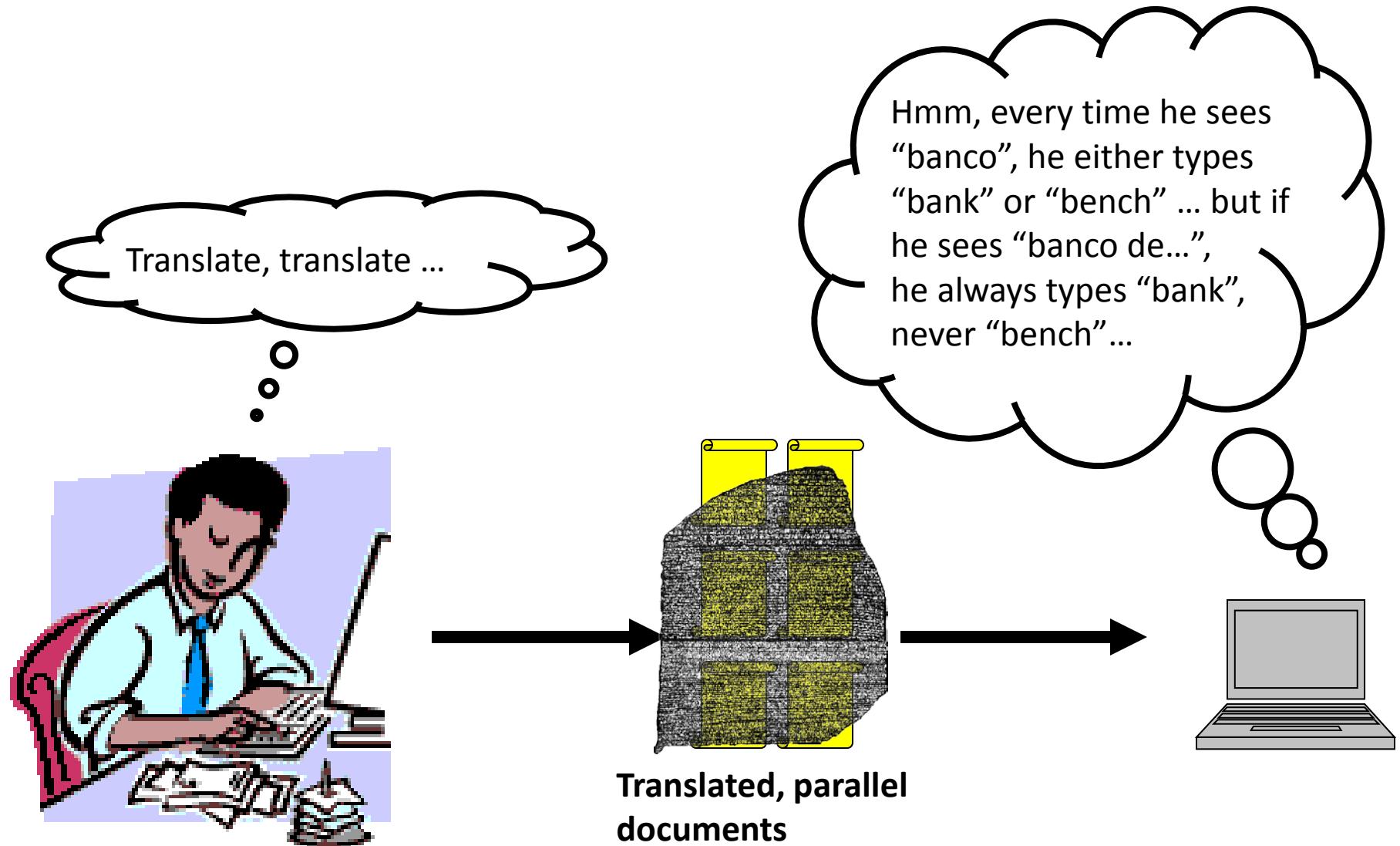


Weaver saw a colleague decoding intercepts into Turkish, without “knowing” Turkish.

... maybe a computer could translate into English without “knowing” English?

OUR HERO

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 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 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 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 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 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 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 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 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 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 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 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 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 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** 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 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 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 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 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 quat oloat at-yurp . totat nnat quat oloat at-yurp . 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 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 . X / /
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 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 . X / •••••• process of 10b. wat nnat gat mat •••••• elimination
4b. at-voon krat pippat sat lat .	11a. lalok nok crrrok hihok yorok zanzanok . / / /
5a. wiwok farok izok stok .	11b. wat nnat arrat mat zanzanat .
5b. totat jjat quat cat .	12a. lalok rarok nok izok hihok mok . / / / /
6a. lalok sprok izok jok stok .	12b. wat nnat forat arrat vat gat .
6b. wat dat krat quat cat .	

Statistical Machine Translation

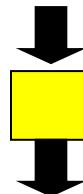
1999

- one page per day, low quality
- limited domains, languages
- no commercial offerings

2014

- fast, 50+ languages
- quality much improved
- products for business, intelligence, end users

أعلن الرئيس الصومالي شريف شيخ أحمد أثناء زيارة لمواقع لقوات الحكومية والأفريقية في حي هودن في مدينتي أن الحملة العسكرية الحالية "لن تتوقف حتى تتحرر الصومال من الشباب والقاعدة".



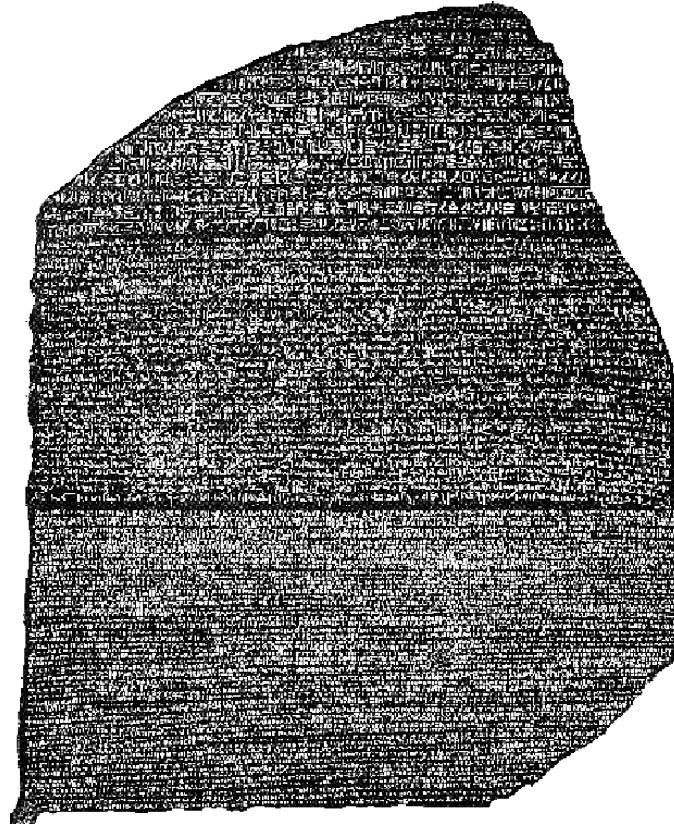
Somali President Sharif Sheikh Ahmed during a visit to sites of government forces and African in the district of Howden in Mogadishu that the current military campaign "will not stop until Somalia is liberated from the youth and al-Qaeda. "

Learning Translation Knowledge from Parallel Text

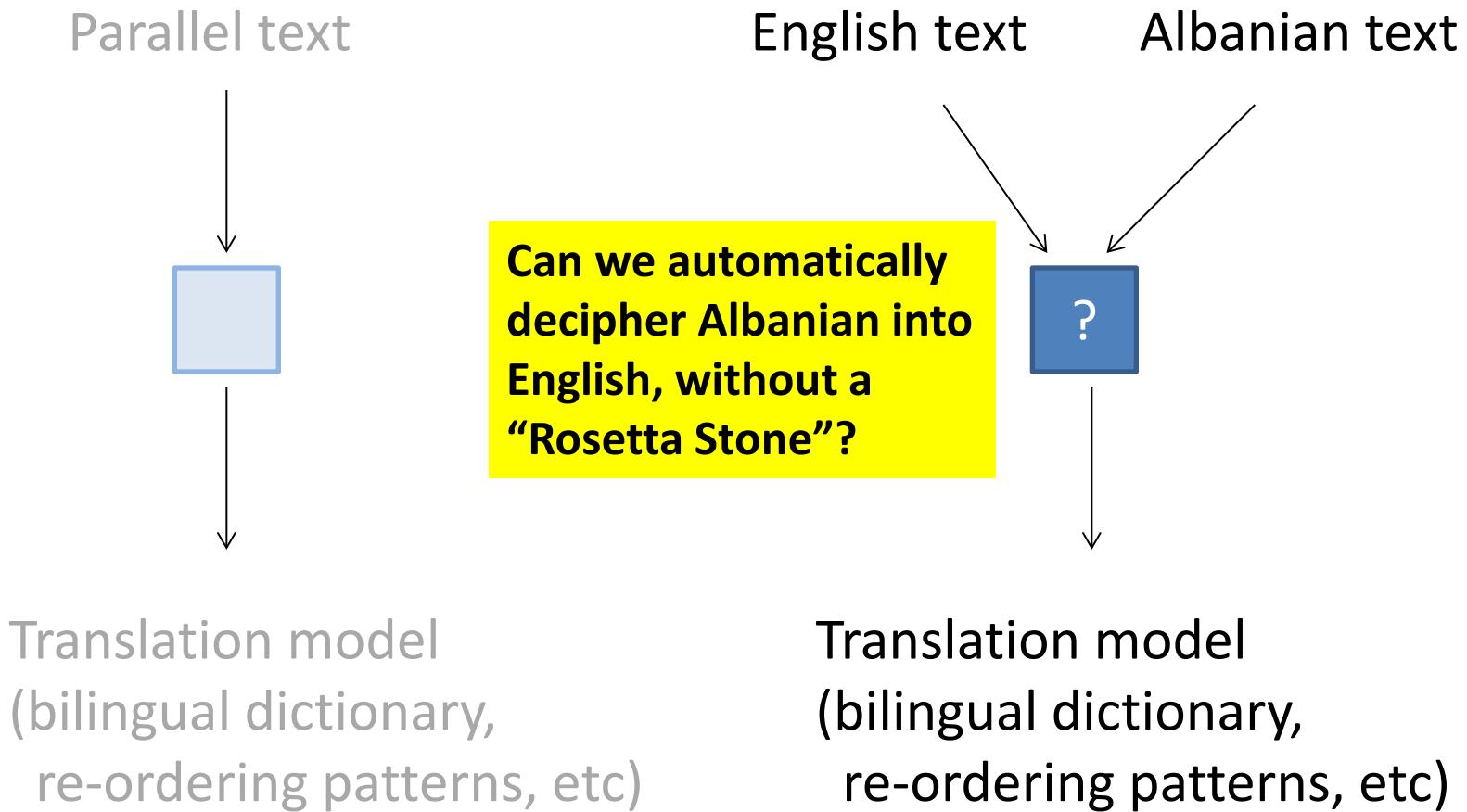
Parallel text



Translation model
(bilingual dictionary,
re-ordering patterns, etc)



Learning Translation Knowledge from Non-Parallel Text?



African Languages

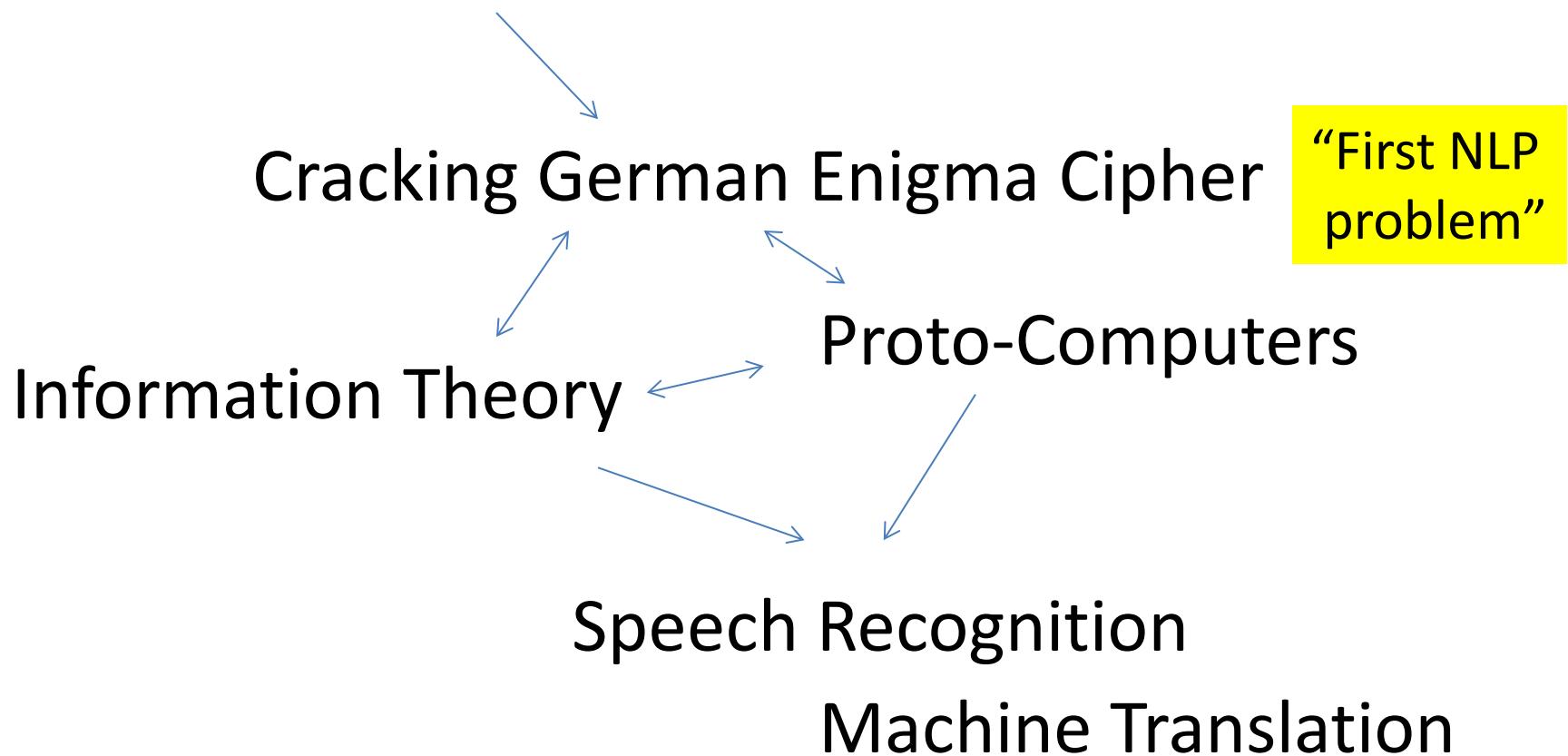


Zero languages spoken

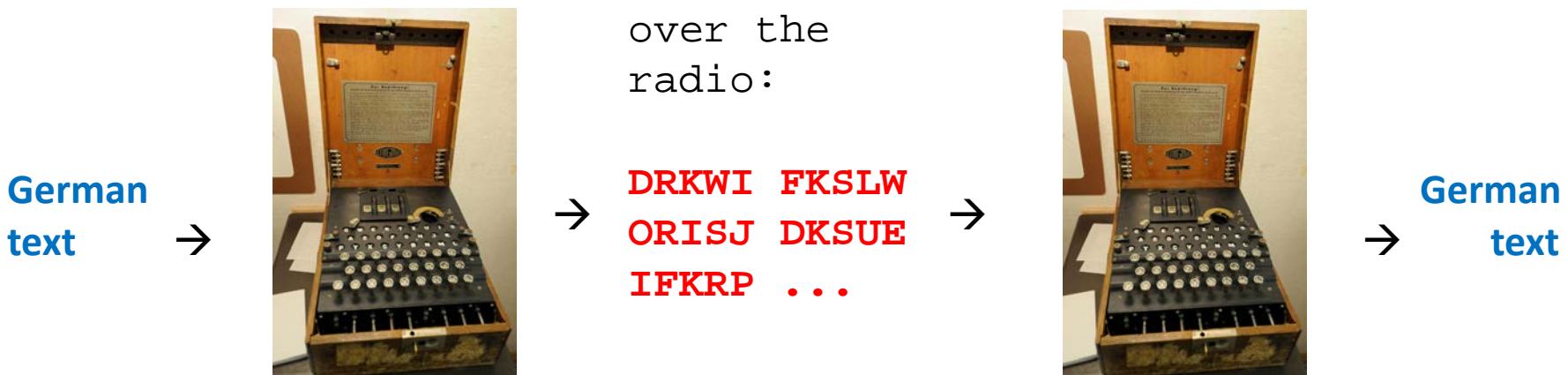
1000+ languages spoken,
40+ by 1m+ speakers

Translation as Decipherment? Our Cryptographic Roots

Automating Cryptography (1930s)



German Enigma Cipher Machine (1920s-1940s)



input (intercepted ciphertext):
output (plaintext):

DRKWI FKSLW ORISJ DKSUE IFK ...
VASIS TDASH ERRCA PITAN RIC ...

German Enigma Cipher Machine

Substitution system

$N \rightarrow J$

Substitution table **changes** with every keystroke:

$NNN \rightarrow JTE$

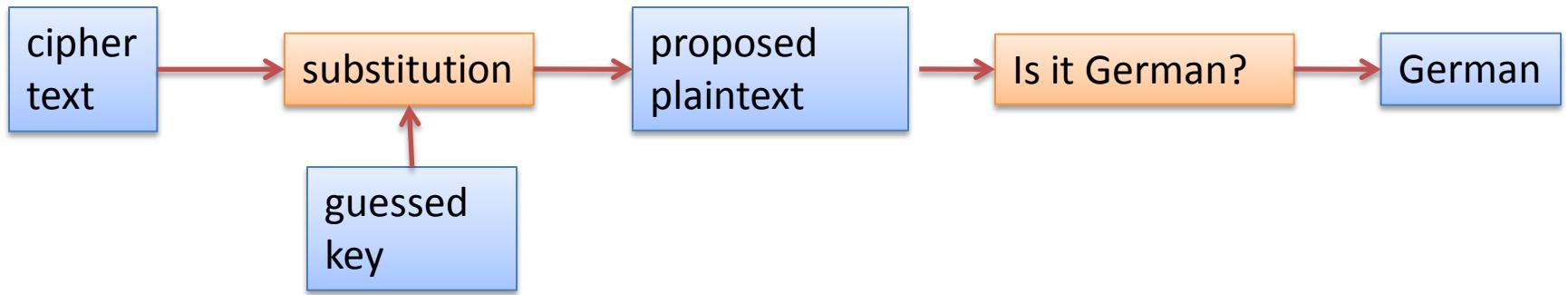
Secret key =
initial rotor
ordering and
settings

Reversible behavior

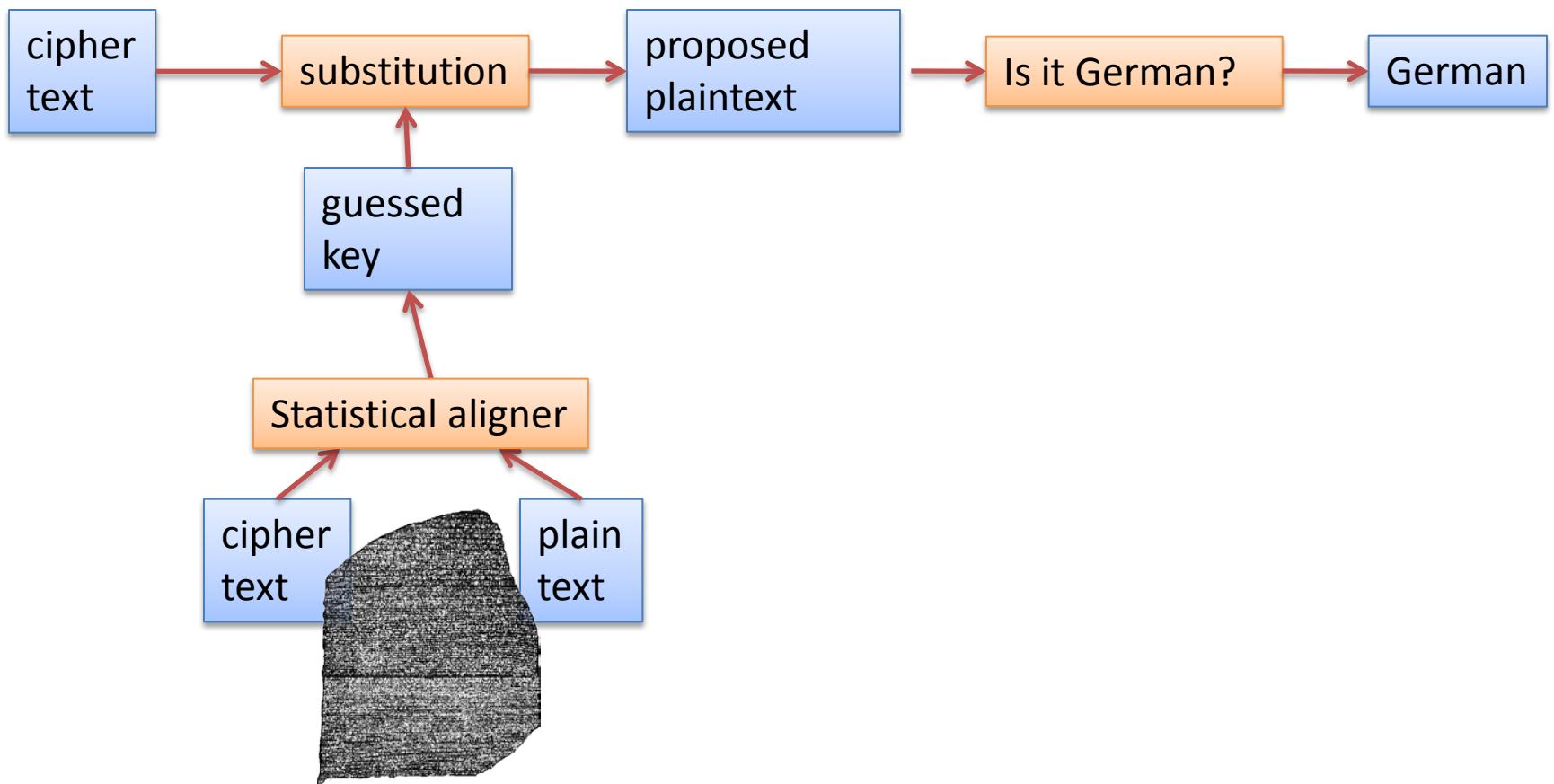
$NNN \rightarrow JTE \rightarrow NNN$



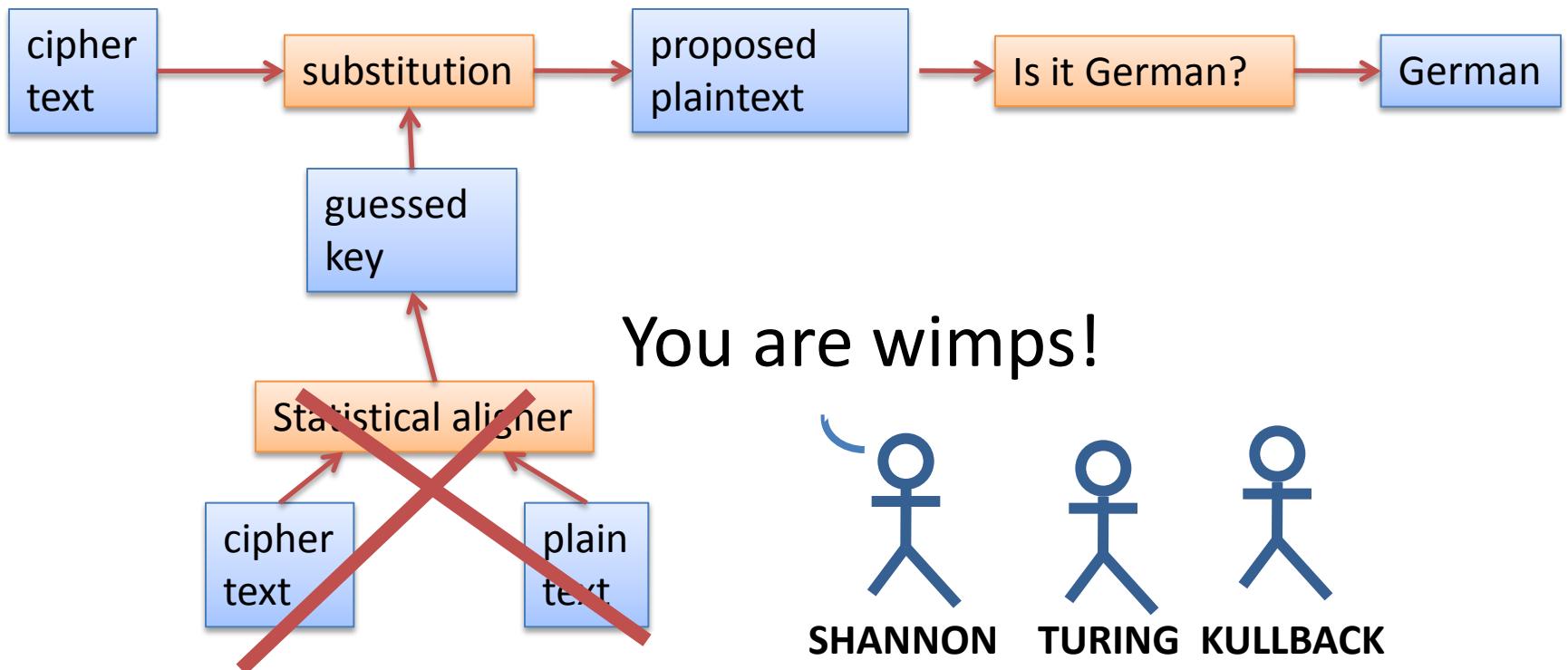
Breaking Enigma



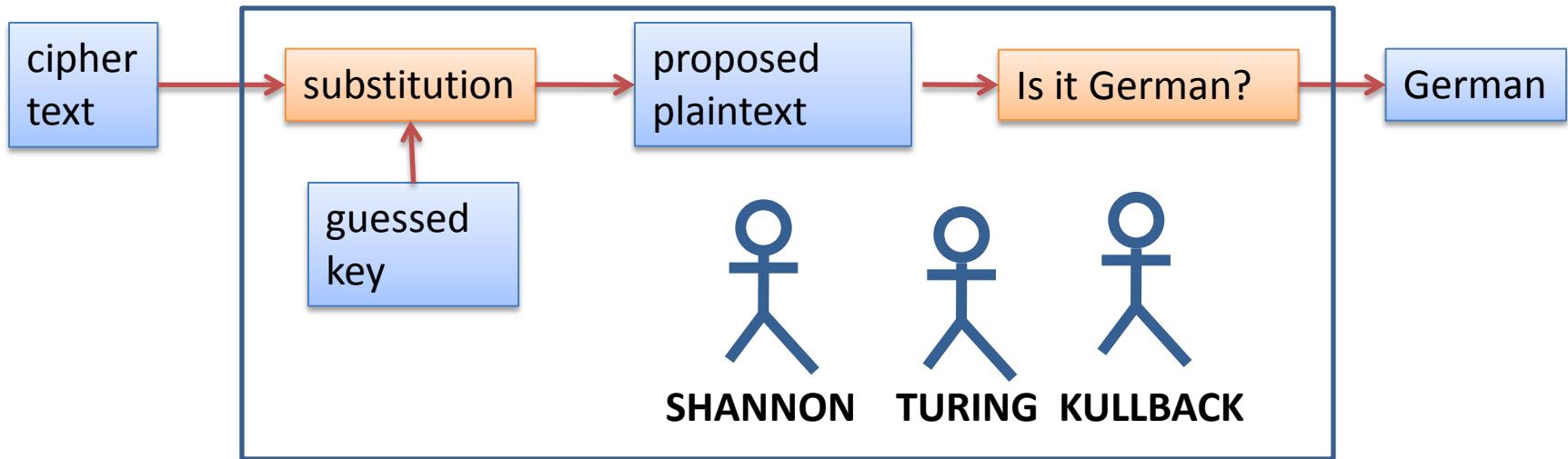
Breaking Enigma



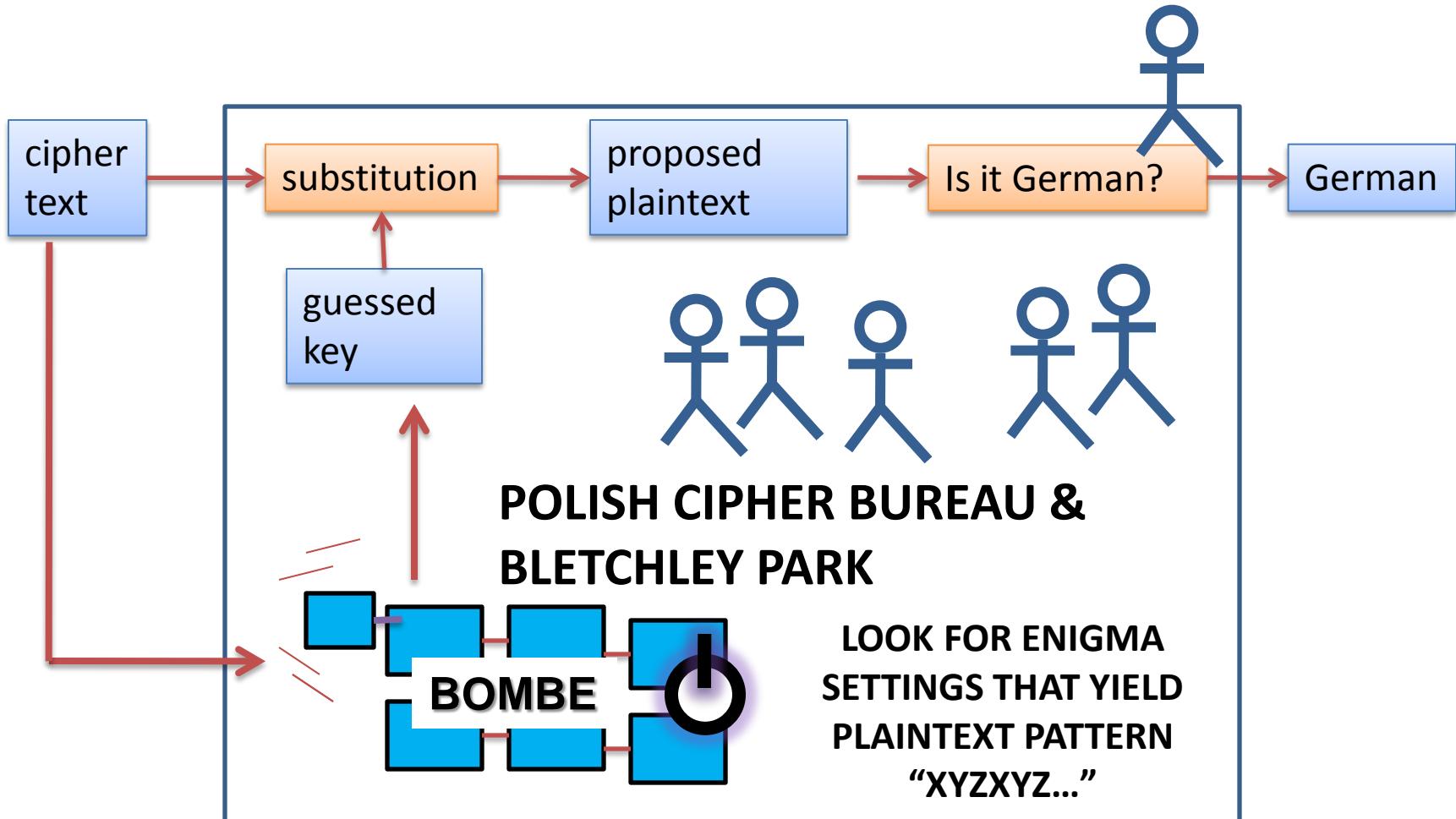
Breaking Enigma



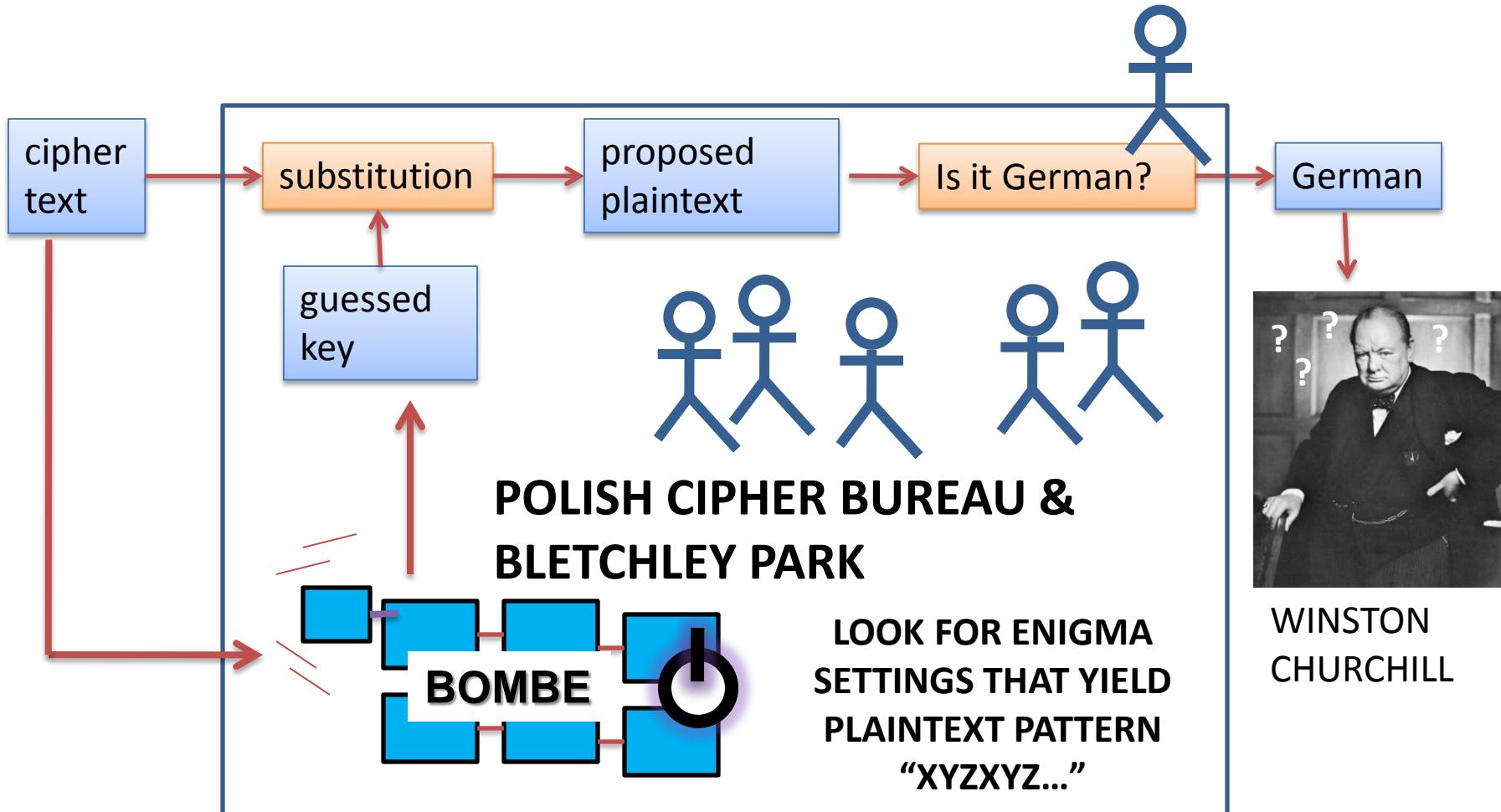
Breaking Enigma



Breaking Enigma



Breaking Enigma

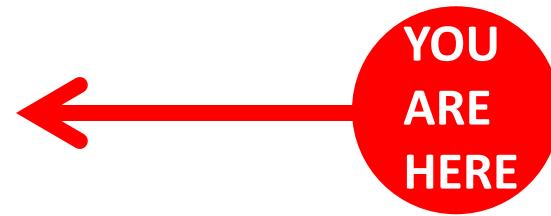


Plan for This Talk

- Break a series of codes
 - Letter substitution
 - Phonetic substitution
 - archaeology
 - transliteration
 - Word substitution
 - Foreign language as cipher
- Bonus
 - Two historical ciphers
 - Final thought on translation and cryptography

Plan for This Talk

- Break a series of codes
 - Letter substitution
 - Phonetic substitution
 - archaeology
 - transliteration
 - Word substitution
 - Foreign language as cipher
- Bonus
 - Two historical ciphers
 - Final thought on translation and cryptography



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."

KDCY LQZKTLJKX CY MDBCYJQL: "TR

HYD FKXC, FQ MKX RLQQIQ HYDL

MKL DXCTW RDSDLQ 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 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 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 . a . a

KDCY LQZKTLJKX CY MDBCYJQL: "TR

. . a . a . . .

HYD FKXC, FQ MKX RLQQIQ HYDL

a a

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 MDBCYJQL: "TR

. .a .e a . ee.e .

HYD FKXC, FQ MKX RLQQIQ HYDL

a . . e .e .a

MKL DXCTW RDSDLQ 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 MDBCYJQL: "TR

. .a .e a . ee.e .

HYD FKXC, FQ MKX RLQQIQ HYDL

a o. . e .e .a o

MKL DXCTW RDSDLQ 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 MDBCYJQL: "TR

.o .a .e a . ee.e .o

HYD FKXC, FQ MKX RLQQIQ HYDL

a . . e .e .a

MKL DXCTW RDSDLQ 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 RLQQIQ HYDL

ar . f re .e .a

MKL DXCTW RDSDLQ 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 RLQQIQ HYDL

ar . f re .e .a

MKL DXCTW RDSDLQ JQMNKXTMB

. o.er a .a. o re.a r

PTBMYEQL K FKH CY LQZKTL TC."

frequent cipher letters: Q L K C D T M Y X

frequent English letters: e t o n i r 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	##### V
R	3 .
S	
T	### 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 L K C D T M Y X~~

frequent English letters: ~~e t o a n i r 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 L K C D T M Y X~~

frequent English letters: ~~e t o a n i r 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 L K C D T M Y X~~

frequent English letters: ~~e t o a n i r 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 .

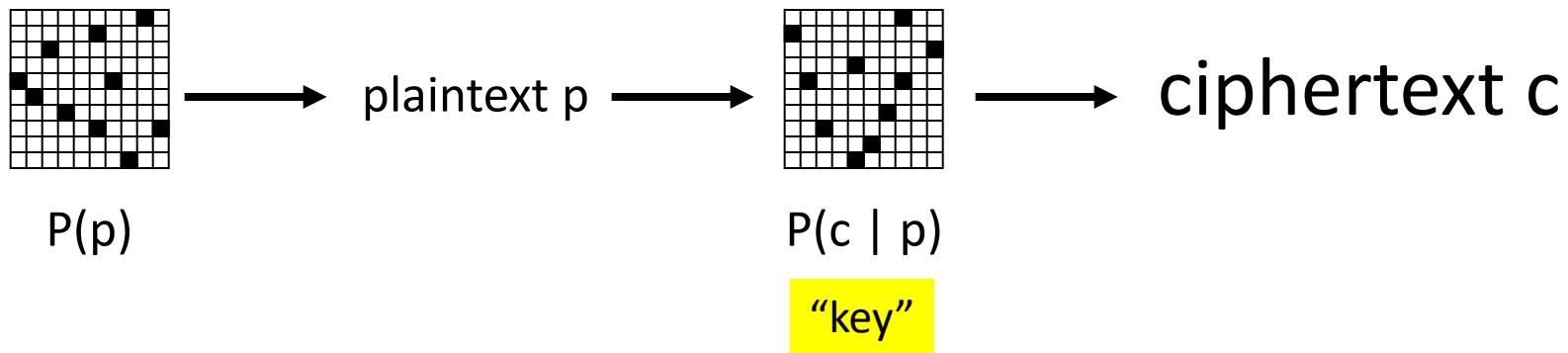
auto repairman to customer: if
KDCY LQZKTLJKX CY MDBCYJQL: "TR
you want we can freeze your
HYD FKXC, FQ MKX RLQQIQ HYDL
car until future mechanics
MKL DXCTW RDCDLQ JQMNKXTMB
discover a way to repair it
PTBMYEQL K FKH CY LQZKTL TC."

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	.

Letter Substitution Cipher

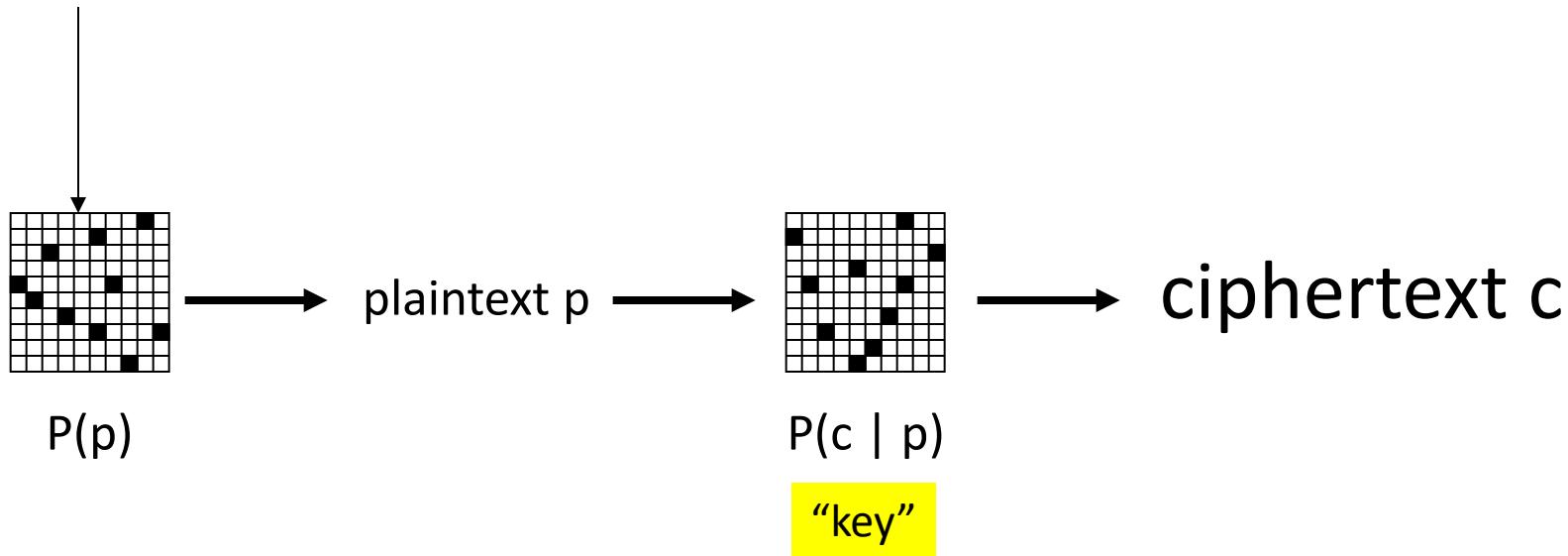
ciphertext c

Letter Substitution Cipher



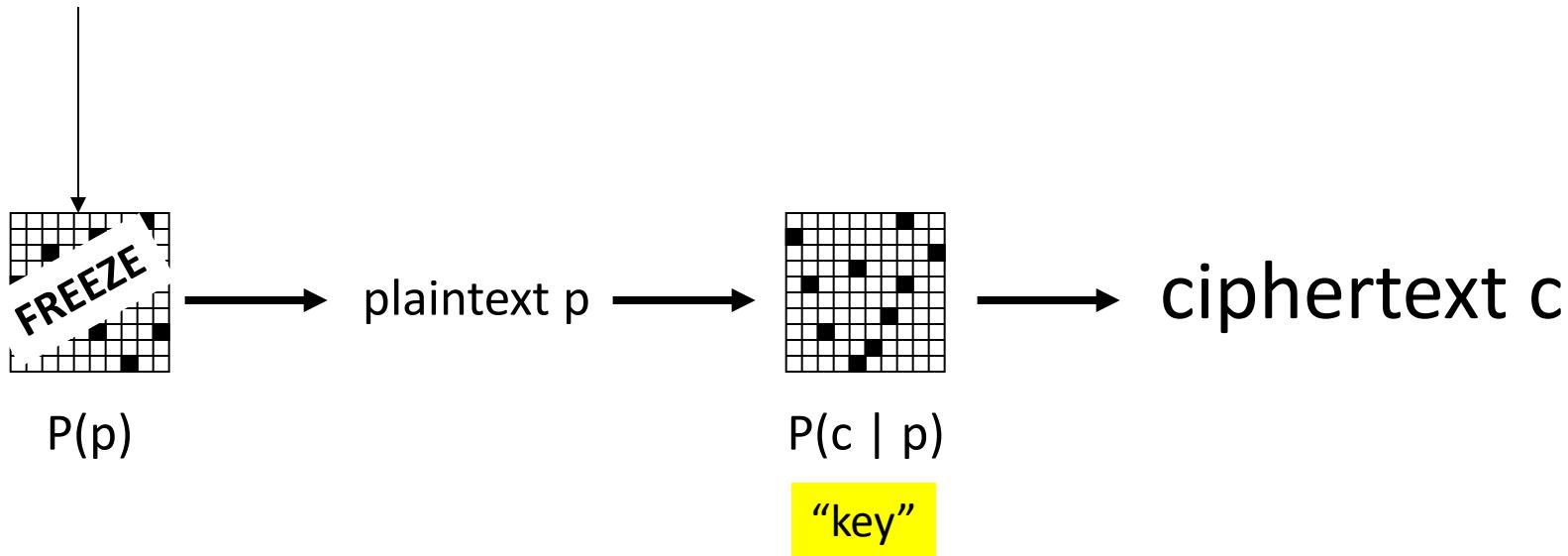
Letter Substitution Cipher

plaintext samples,
unrelated to ciphertext

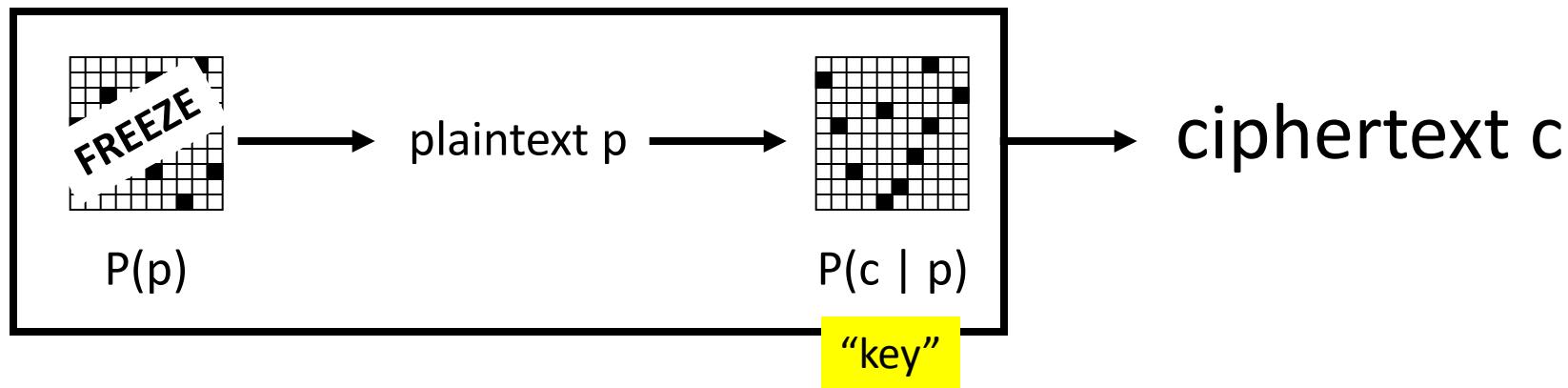


Letter Substitution Cipher

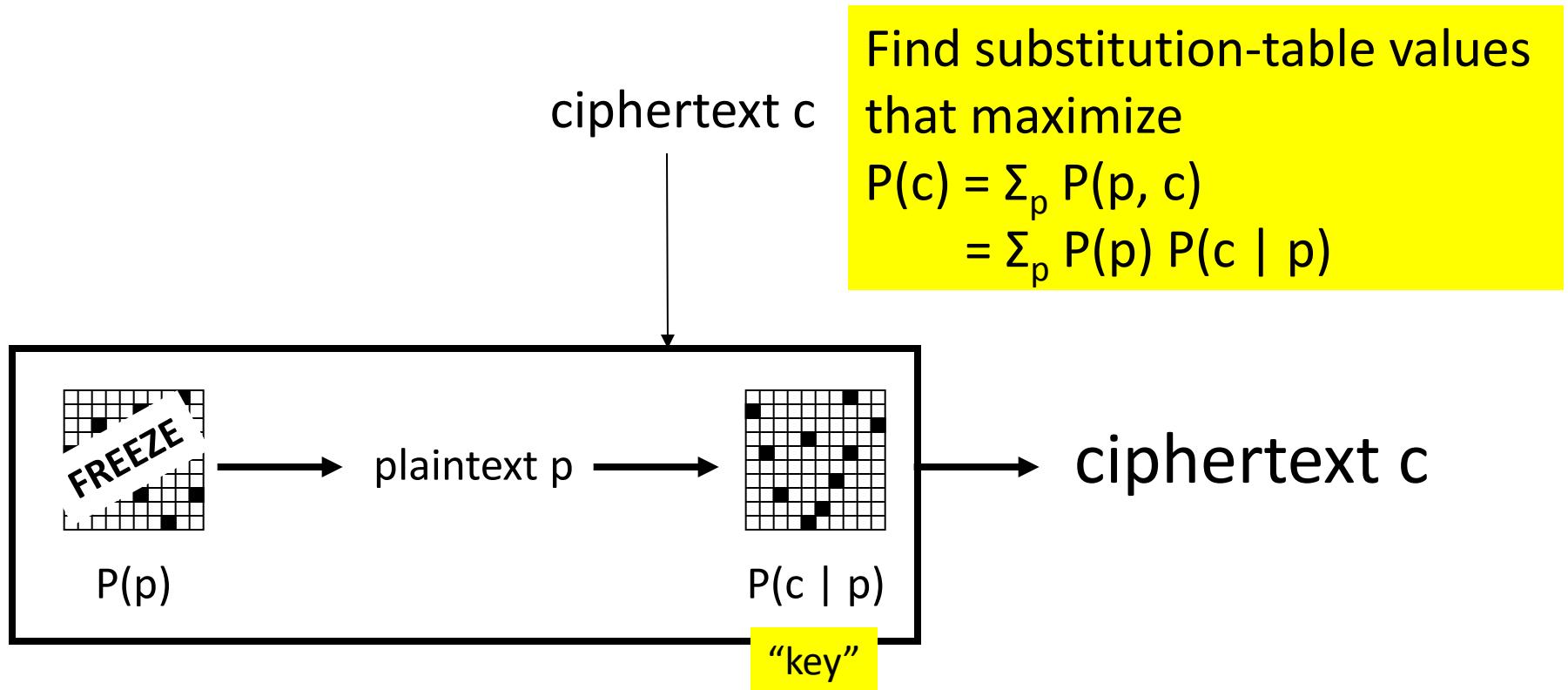
plaintext samples,
unrelated to ciphertext



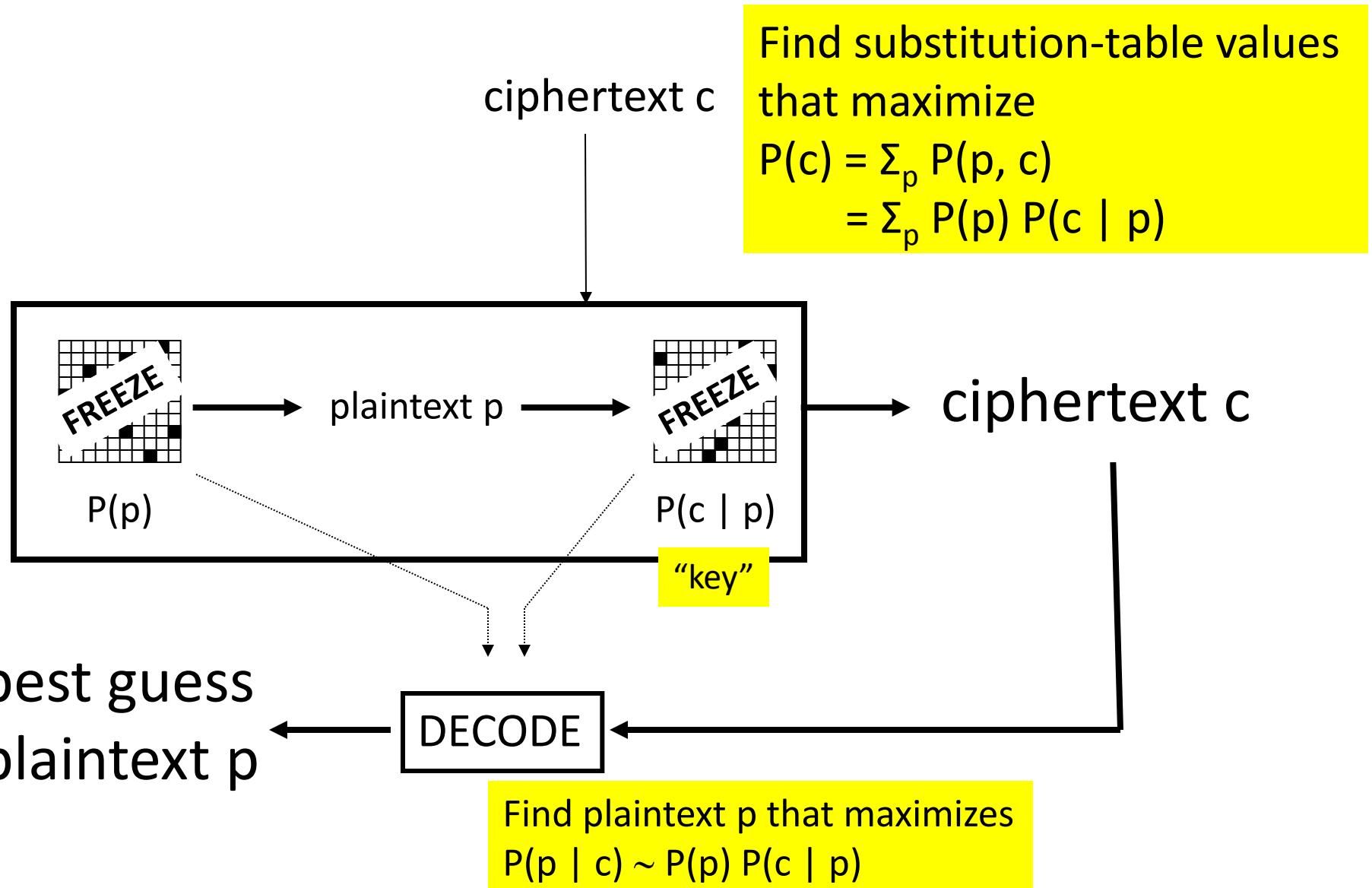
Letter Substitution Cipher



Letter Substitution Cipher



Letter Substitution Cipher

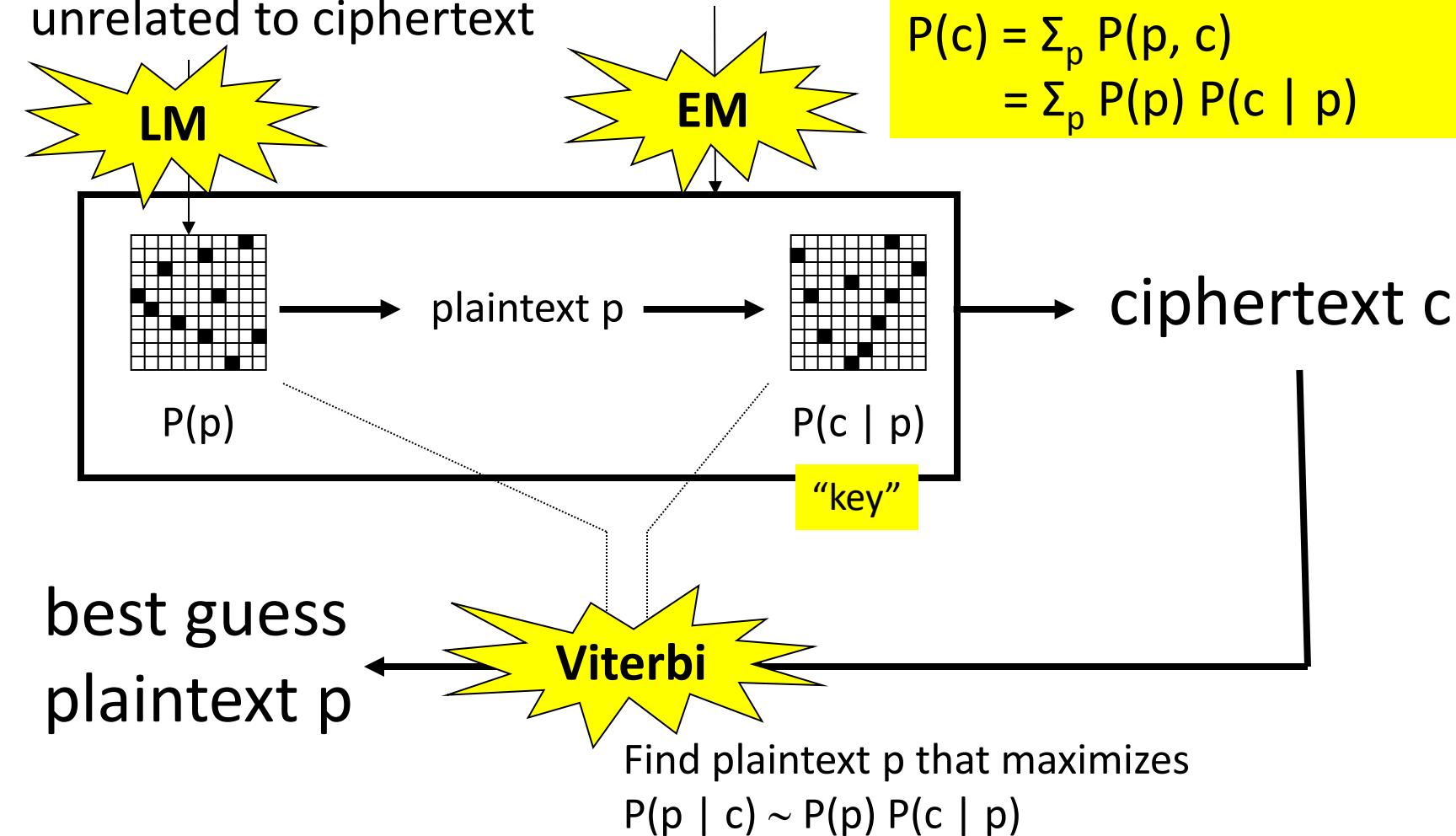


Letter Substitution Cipher

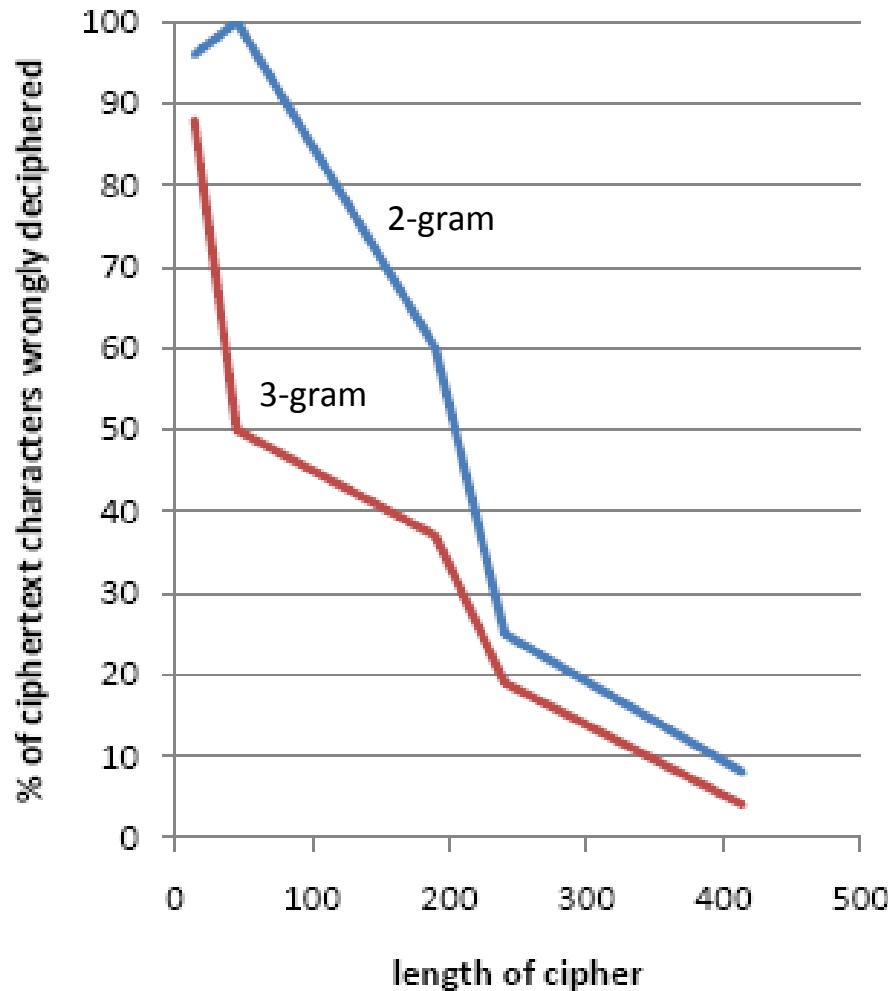
plaintext samples,
unrelated to ciphertext

ciphertext c

Find substitution-table values
that maximize
 $P(c) = \sum_p P(p, c)$
 $= \sum_p P(p) P(c | p)$



Decipherment Accuracy vs. Cipher Length



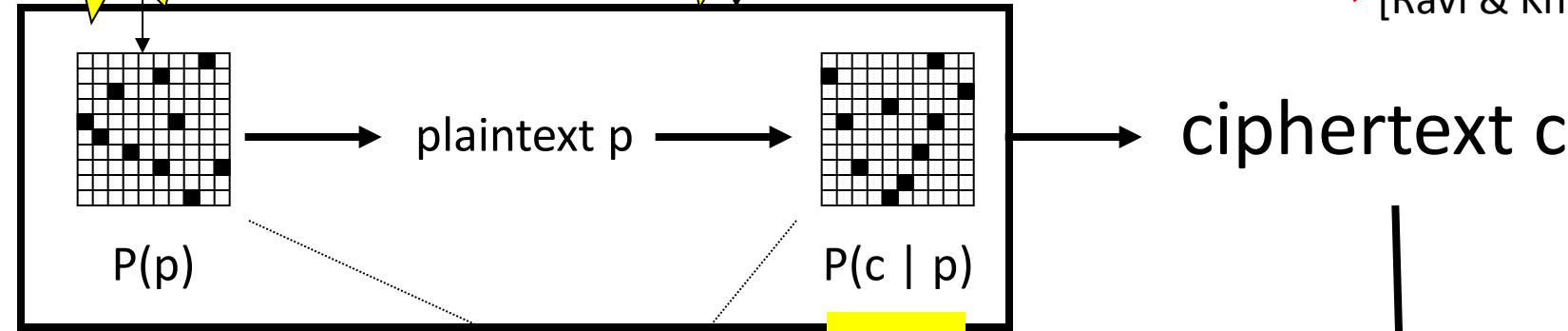
Letter Substitution Cipher

plaintext samples,
unrelated to ciphertext

ciphertext c

Find substitution-table values
that maximize
 $P(c) = \sum_p P(p, c)$
 $= \sum_p P(p)^{0.5} P(c | p)$

[Ravi & Knight 09b]



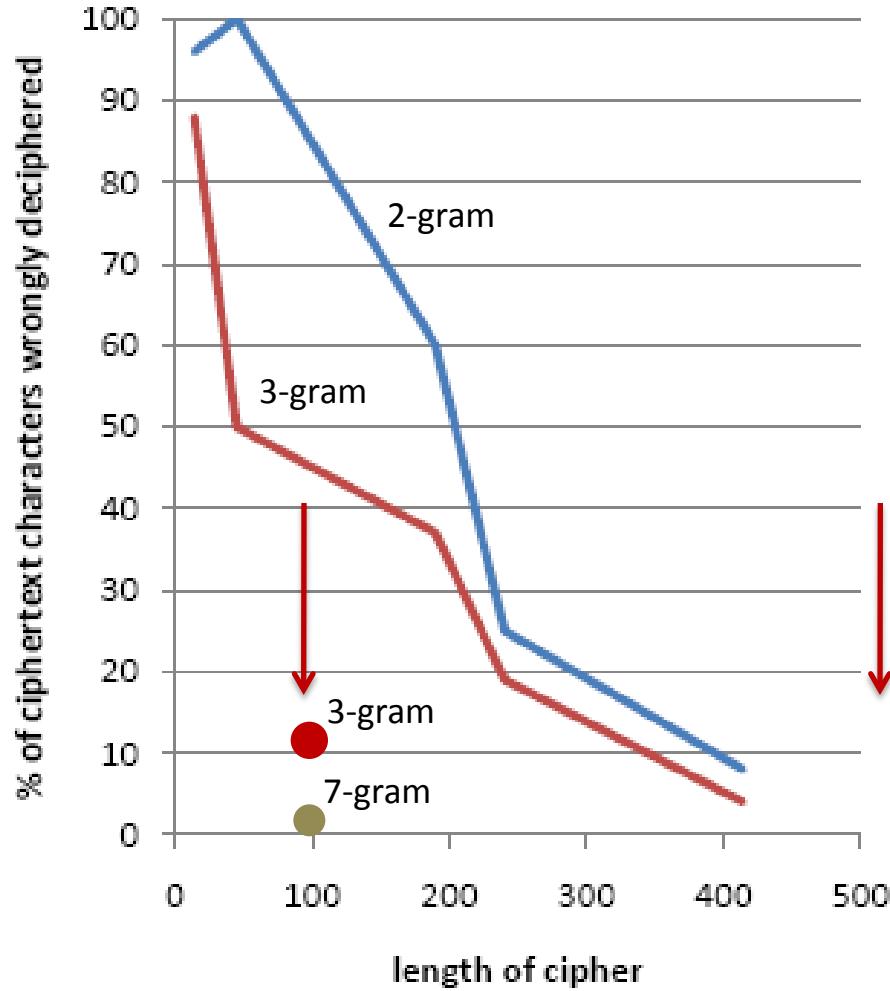
best guess
plaintext p

Viterbi

Find plaintext p that maximizes
 $P(p | c) \sim P(p) P(c | p)^3$

[Knight/Yamada 99]

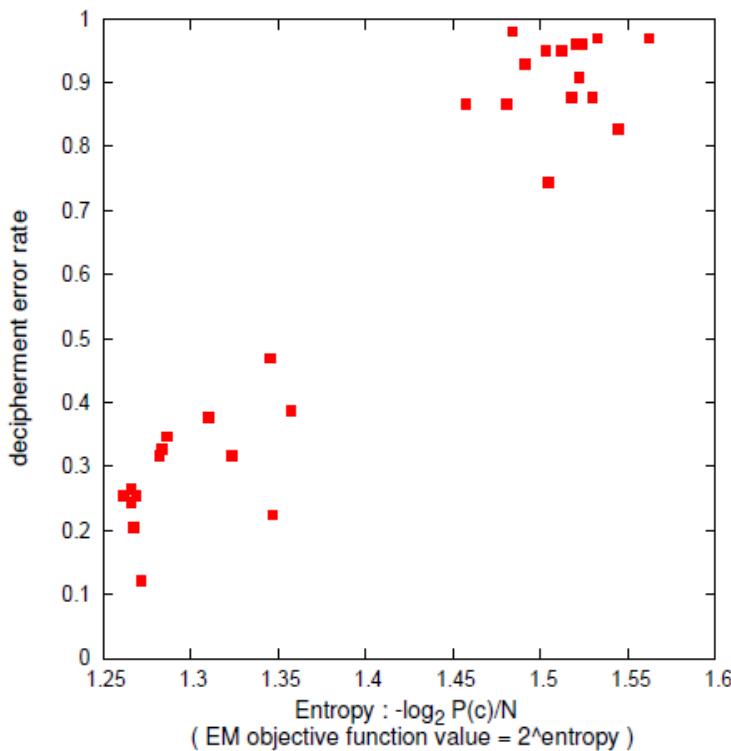
Reducing LM Weight During EM



Set EM to maximize
 $P(c) \approx \sum_p P(p)^{0.5} P(c | p)$
instead of
 $P(c) \approx \sum_p P(p) P(c | p)$

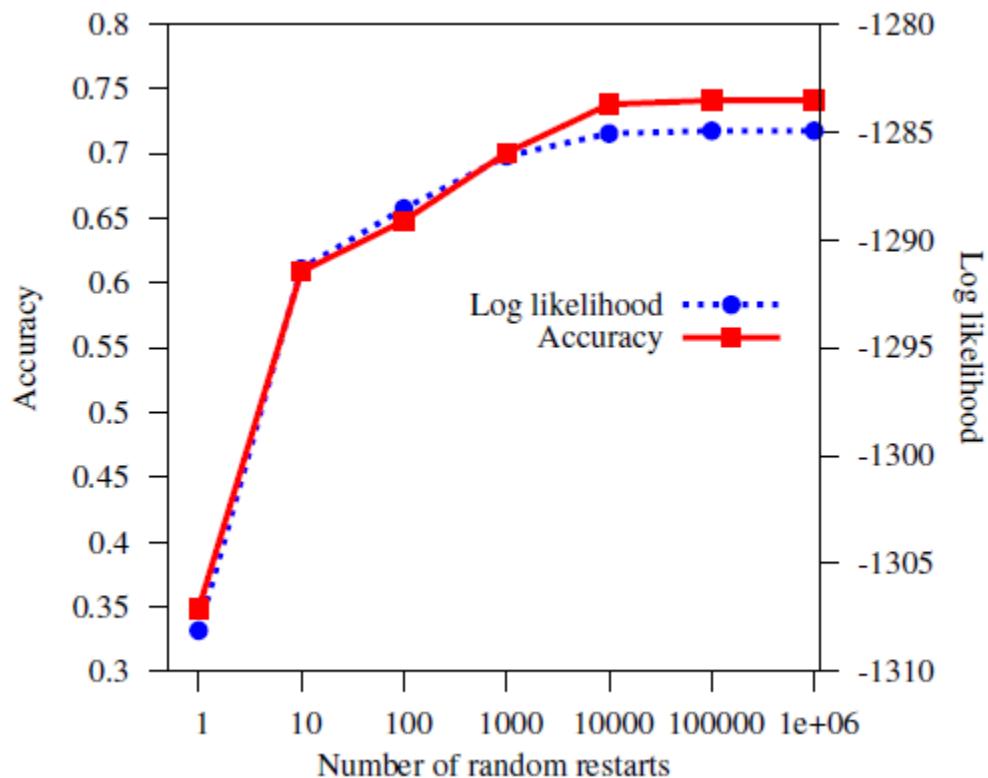
Random Restarts are Critical

English 98-letter cipher, 3-gram LM



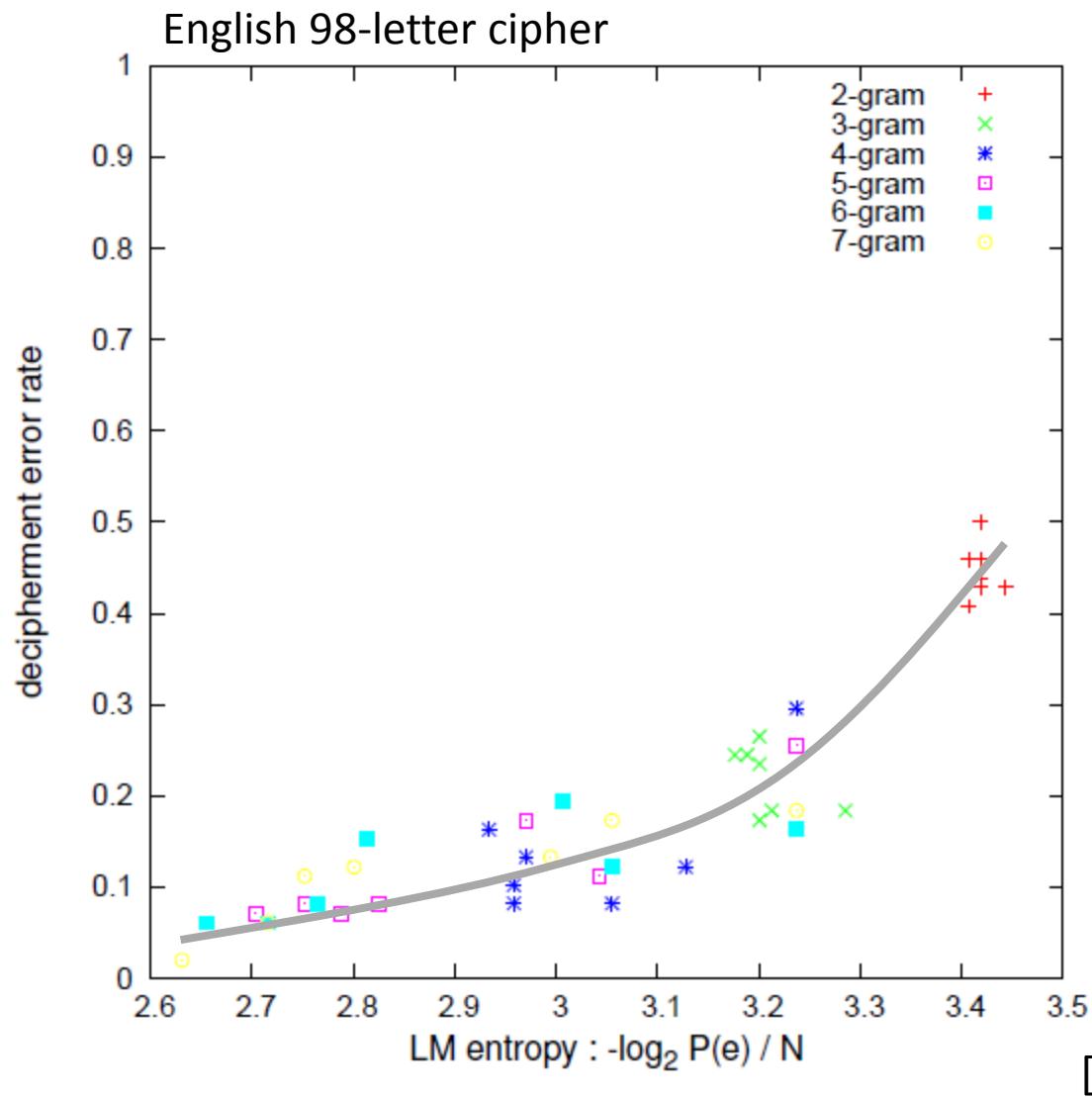
[Ravi & Knight 09b]

Zodiac killer 340 cipher



[Berg-Kirkpatrick & Klein 13]

Good Language Models are Critical



Deterministic Substitution Constraint

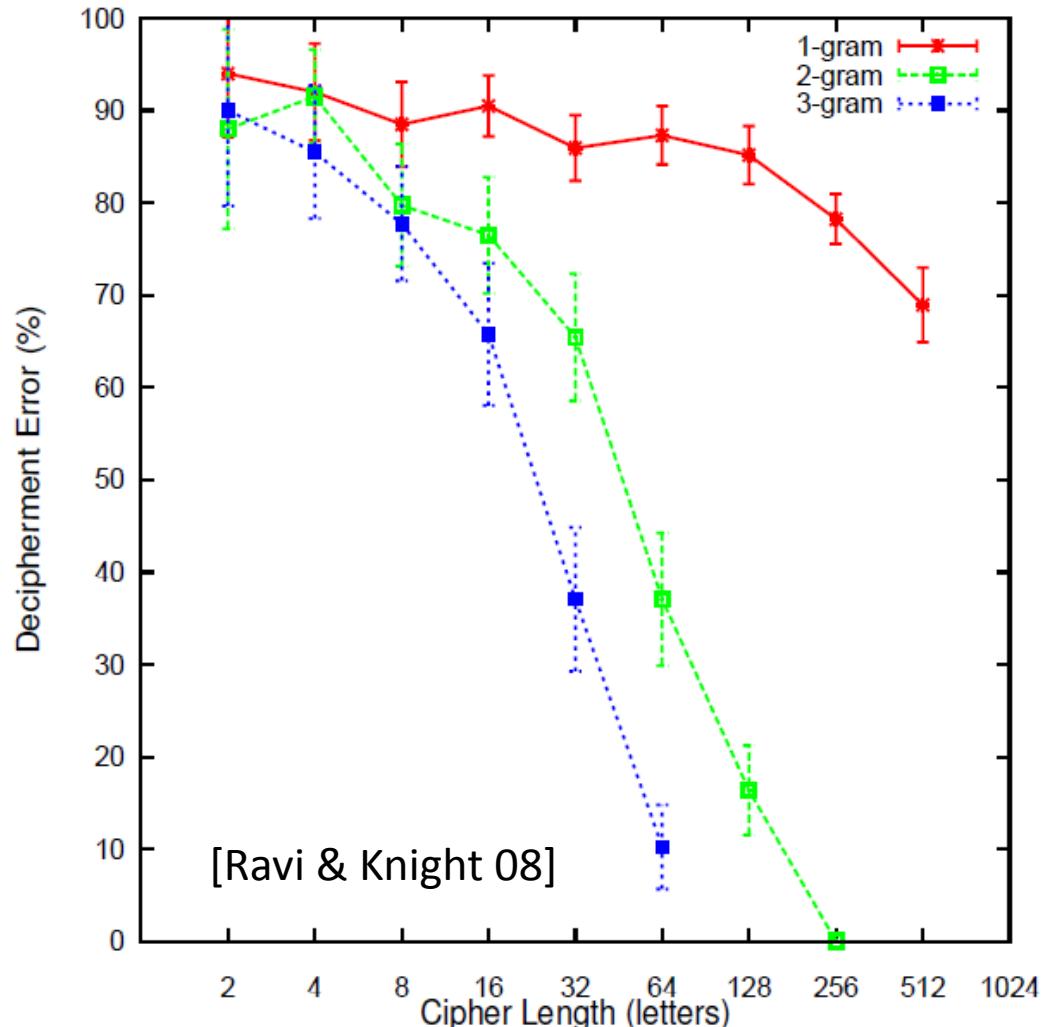
Using ILP instead of EM

- * Search only over 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

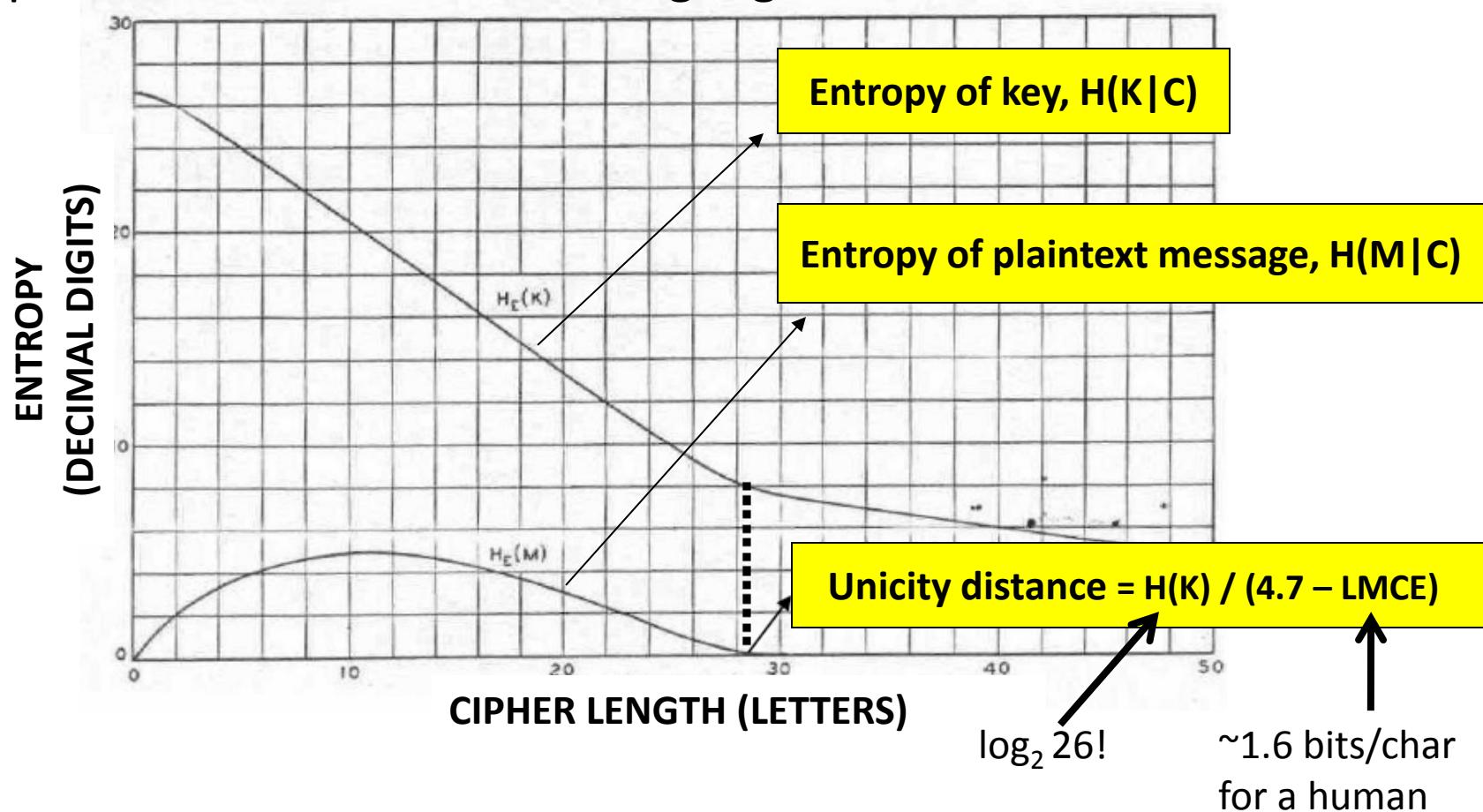


see also [Nuhn, Ney, Schamper] on beam search

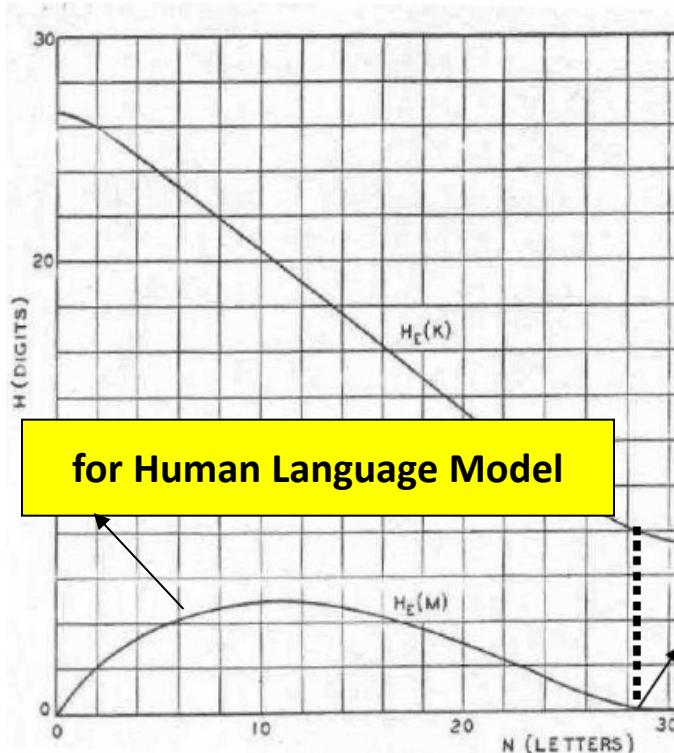
[Shannon 46, 49]

“Communication Theory of Secrecy Systems”

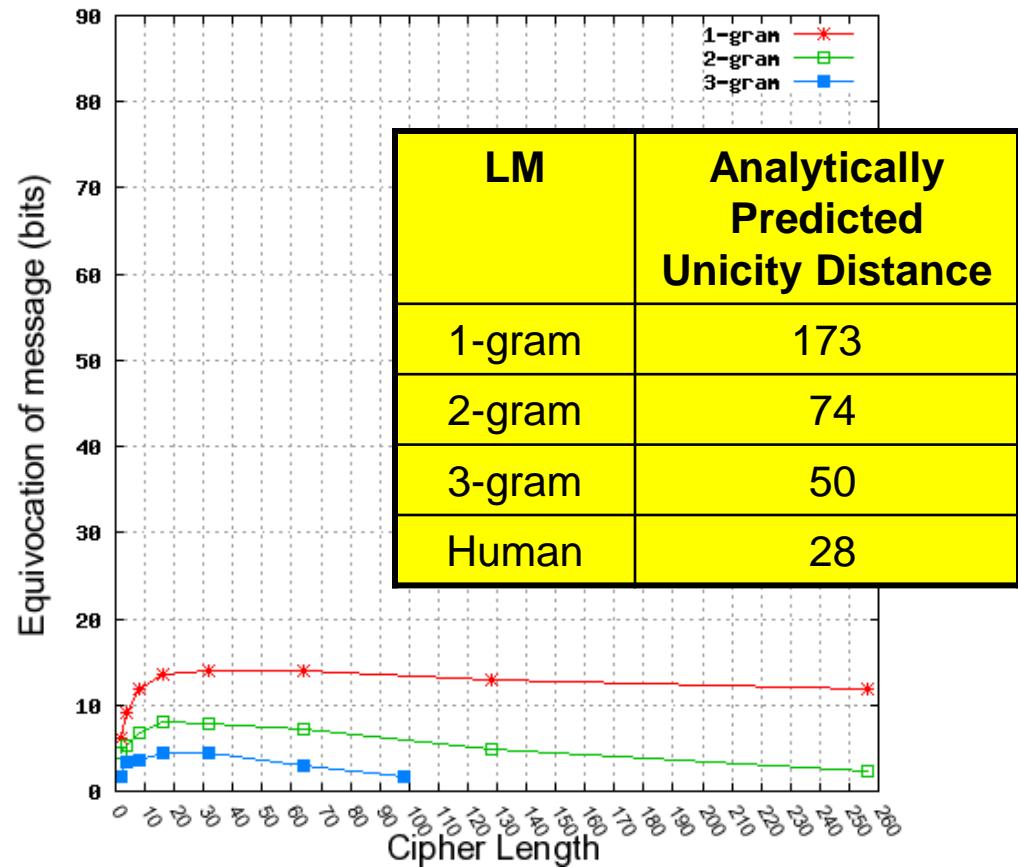
- Shannon analytically predicted uncertainty about key and message
- Graphed it for a human-level language model



Verifying Shannon's Prediction of Plaintext Message Uncertainty



ANALYTIC CURVES
(Shannon's)

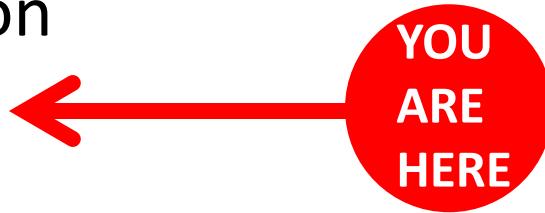


ACTUAL CURVES
(ours)

[Ravi & Knight 08]

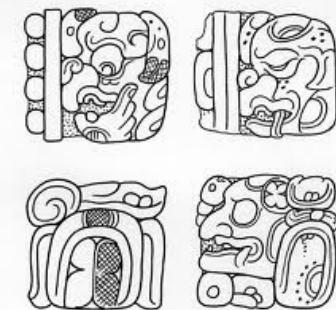
Plan for This Talk

- Break a series of codes
 - Simple letter substitution
 - Phonetic substitution
 - archaeology
 - transliteration
 - Word substitution
 - Foreign language as cipher
- Bonus
 - Two historical ciphers
 - Final thought on translation and cryptography



Phonetic Decipherment

ciphertext



Phonetic Decipherment

ciphertext

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

Phonetic 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...”



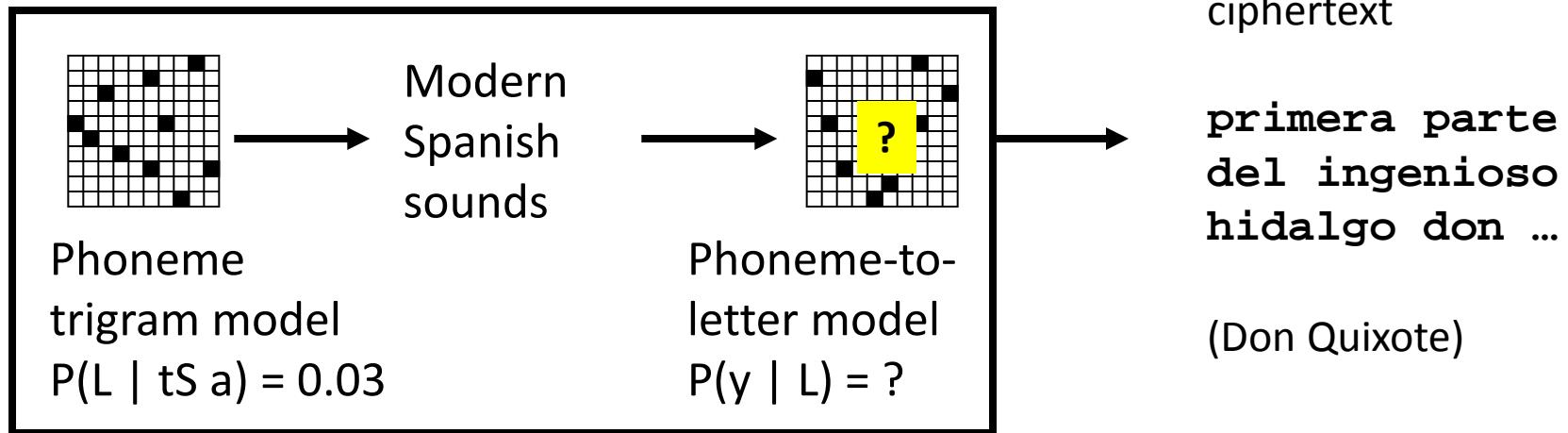
OUR HERO

ciphertext

*primera parte
del ingenioso
hidalgo don ...*

(Don Quixote)

Phonetic 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

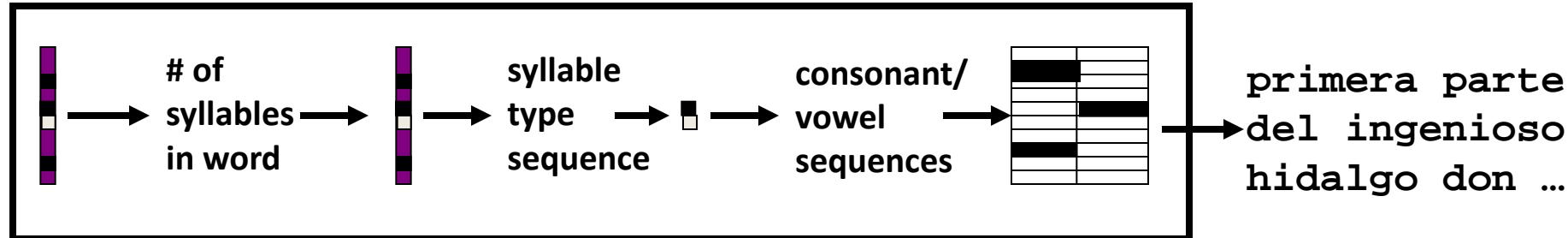
EM approach = 93% accurate phonetic decipherment

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



$$P(1) = ?$$

$$P(2) = ?$$

etc.

$$P(CV) = ?$$

$$P(V) = ?$$

$$P(CVC) = ?$$

+ 7 others

$$P(V | V) = ?$$

$$P(VV | V) = ?$$

$$P(a | V) = ?$$

$$P(a | C) = ?$$

etc.

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

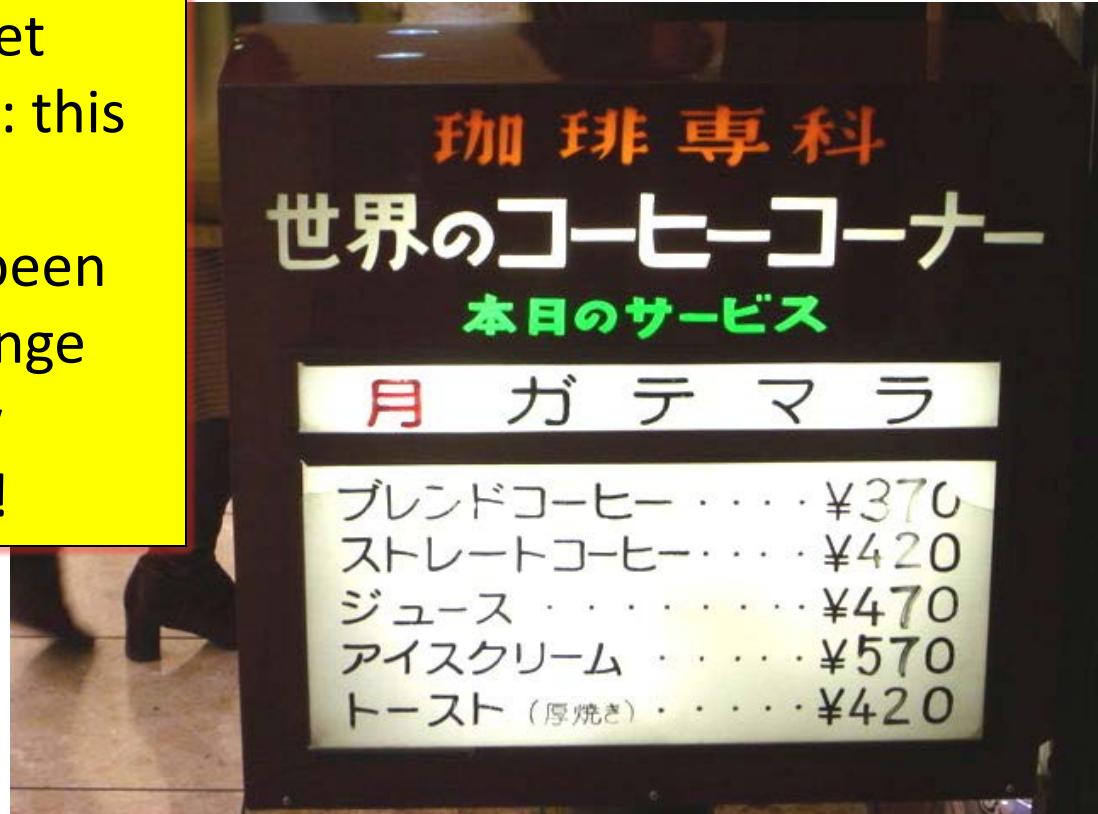
[Knight et al 06]
see also
[Kim & Snyder 13]

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!



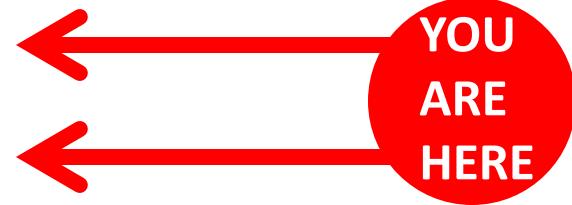
OUR HERO



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

Plan for This Talk

- Break a series of codes
 - Simple letter substitution
 - Phonetic substitution
 - archaeology
 - transliteration
 - Word substitution
 - Foreign language as cipher
- Bonus
 - Two historical ciphers
 - Final thought on translation and cryptography



Word Substitution

Berlin le. n. d. Mars 1783.

Monsieur

Each code number represents
a plaintext **word**, not letter

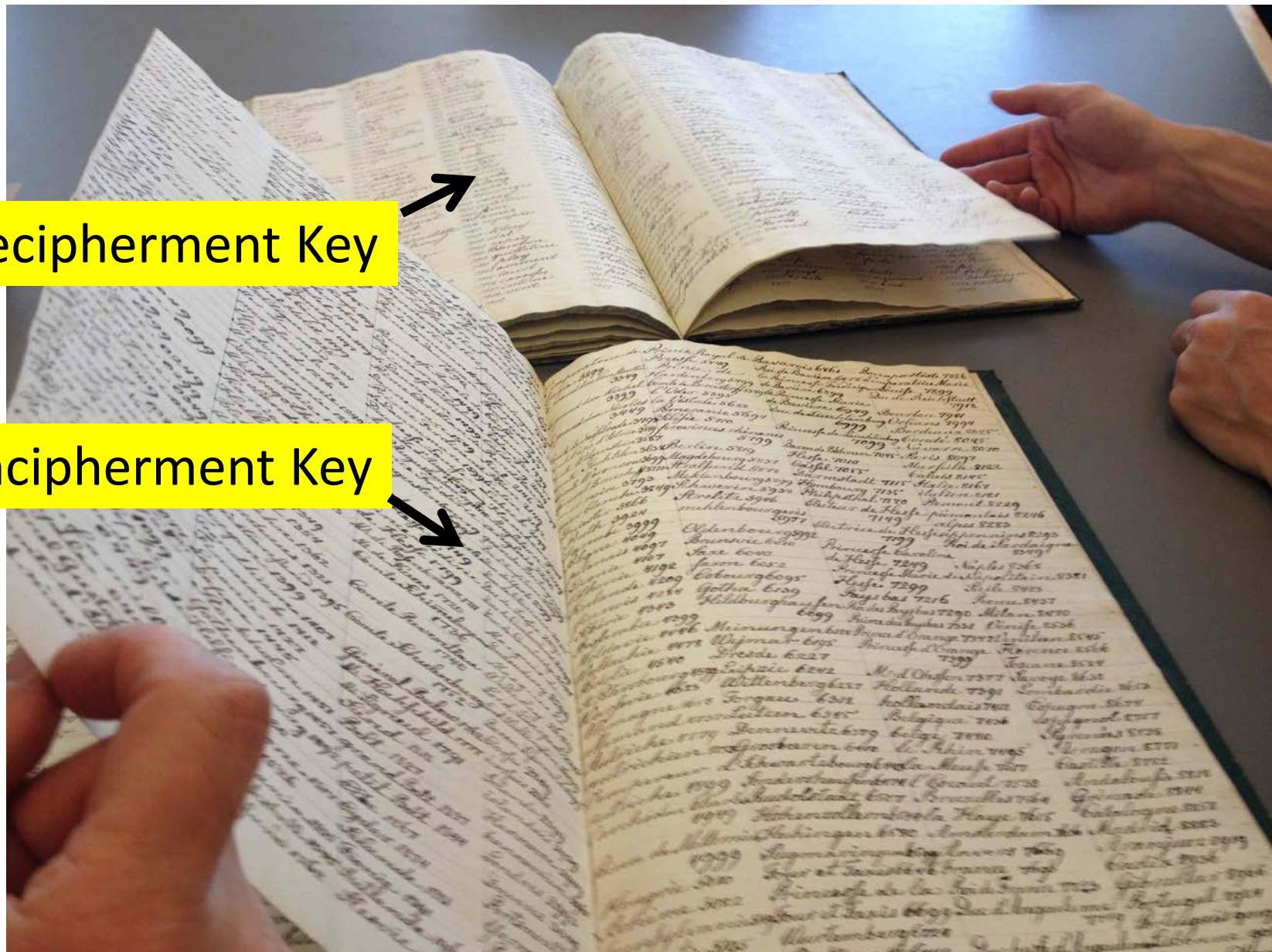
1200. 3660. 2300. 2012. 5519. 601. 1268. 3228. 53. 3039. 2017. 3595. 690. 280. 2059. 4116.
1496. 4102. 3346. 1A87. 56. 2362. 2665. 3054. 2097. 1456. 4133. 760. 3079. 2181. 3179. 512.
3480. 0969. 3754. 165. 1826. 3101. 2743. 4385. 4317. 1497. 1199. 416. 120. 601. 2003. 2006. 2015.
863. 1067. 2637. 3794. 1496. 2254. 3818. 979. 2722. 3A22 3661. 4020. 433. 1630. 1480. 1969.
2587. 518. 286. 452. 2362. 2A10. 2987. 615. 1031. 2524. 380. 2006. 5793. 2819. 3849. 330.
11021. 1A8. 3468. 2446. 451. 1224. 3915. 2279. 3503. 55. 1521. 0A36. 1692. 2012. 1393. 615. 1129.

Word Substitution Keys

Decipherment Key



Encipherment Key



Foreign Language as a Cipher

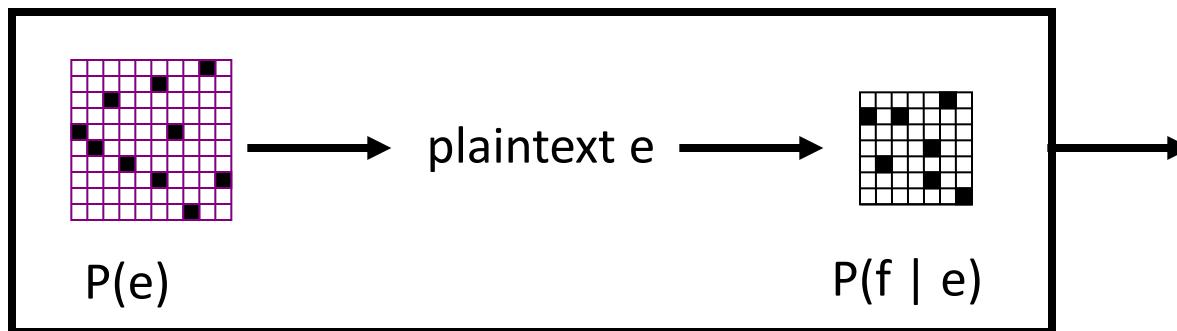
رفض رئيس السلطة الفلسطينية محمود عباس مجددا تصريحات وزير الخارجية الإسرائيلي سيلفان شالوم التي قال فيها إنه يتبعن على إسرائيل إعادة النظر في انسحابها من غزة، المقرر أن يتم الصيف المقبل إذا فازت حركة المقاومة الإسلامية حماس في الانتخابات التشريعية وقال عباس في مؤتمر صحفي على هامش مشاركته في القمة العربية-اللاتينية الأولى إنه يتبعن على إسرائيل احترام خيار الشعب الفلسطيني حتى لو فازت 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 الطبع الأحادي الجانبي لانسحاب الإسرائيلى من غزة، وأكد أن إسرائيل تريد مغادرة هذه الأرضى لتعزيز سيطرتها على الضفة الغربية.

وقال قريع في كلمة له خلال مؤتمر نظمته وزارة الأوقاف في رام الله "سينسحبون من غزة ولكننا لا نعرف ما هو شكل هذا الانسحاب وماذا سيتركون، وما هو مصير المعابر والحدود، وكل ذلك غامض لأنه قرار أحادي الجنان

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!!!”



رفض رئيس السلطة الفلسطينية محمود عباس مجددا تصريحات وزير الخارجية الإسرائيلي سيلفان شالوم التي

قال فيها إنه يتمنى على إسرائيل إعادة النظر في انسحابها من غزة، المقرر أن يتم الصيف المقبل إذا فازت حركة المقاومة الإسلامية حماس في الانتخابات التشريعية وقال عباس في مؤتمر صحفي على هامش مشاركته في القمة العربية-اللاتينية الأولى إنه يتمنى على إسرائيل احترام خيار الشعب الفلسطيني حتى لو فازت 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 لانسحاب الإسرائيلي من غزة، وأكد أن إسرائيل تريد مغادرة هذه الأرضي لتعزيز سيطرتها على الضفة الغربية.

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

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

Time Expressions

@!m 1990
w!lth!ny&
fy !lywm
mn !lsh=hr !lj!ry
!lqrn
!y!m
@!m!aN
!ls!@&
17 shb!T 1994
th!th snw!t
dqyq&
=hdh=h !lsn&
ywmyn
mn !!@!m !lm!Dy
!lsn& !lmqbl&
fy !lsn&
kl ywm
fy !!@!m !lm!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!th& !y!m
qbl !sbw@yn
fy !lywm !lt!ly
sh@b!n
tmwz
3 dhw !!Hj& 1414
fy shb!T !lm!Dy
qbl ywmyn

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

Time Expressions

@!m 1990
w!lth!ny&
fy !lywm
mn !lsh=hr !lj!ry
!lqrn
!y!m
@!m!aN
!!s!@&
17 shb!T 1994
th!lth snw!t
dqyq&
=hdh=h !lsn&
ywmyn
mn !!@!m !lm!Dy
!lsn& !lmqbl&
fy !lsn&
kl ywm
fy !!@!m !lm!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& @sh!&
th!lth& ly!m
qbl !sbw@yn
fy !lywm !lt!ly
sh@b!n
tmwz
3 dhw !!Hj& 1414
fy shb!T !lm!Dy
qbl ywmyn

Time Expressions

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

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

<nb> Hzyr!n <nb>

| | | | |
|----|--------------------|---|-------------------|
| 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 !lmbql | 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 |

Time Expressions

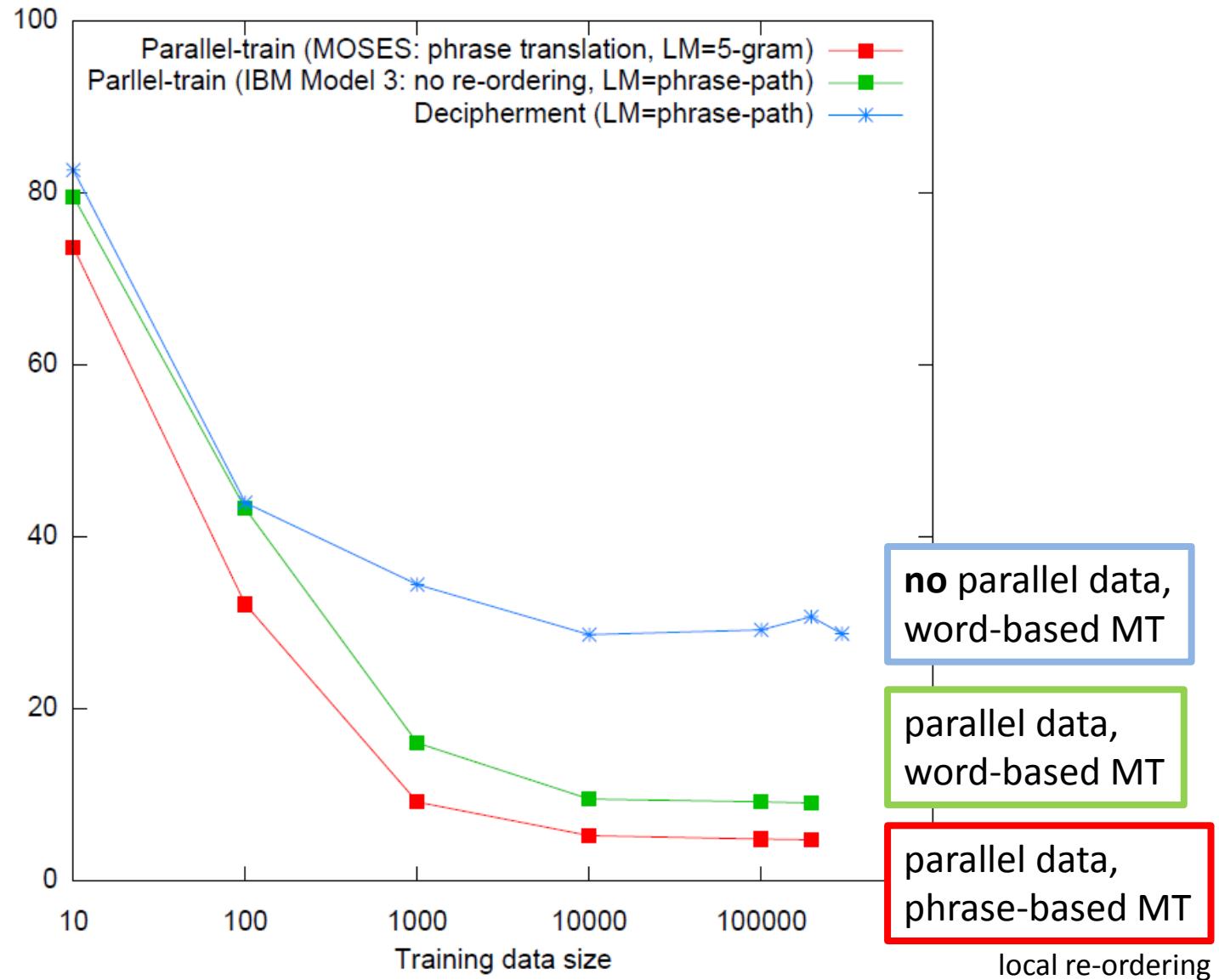
Hzyr!n

| | | | |
|-----|-------------------------|----|------------------------|
| 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 | ty 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 |

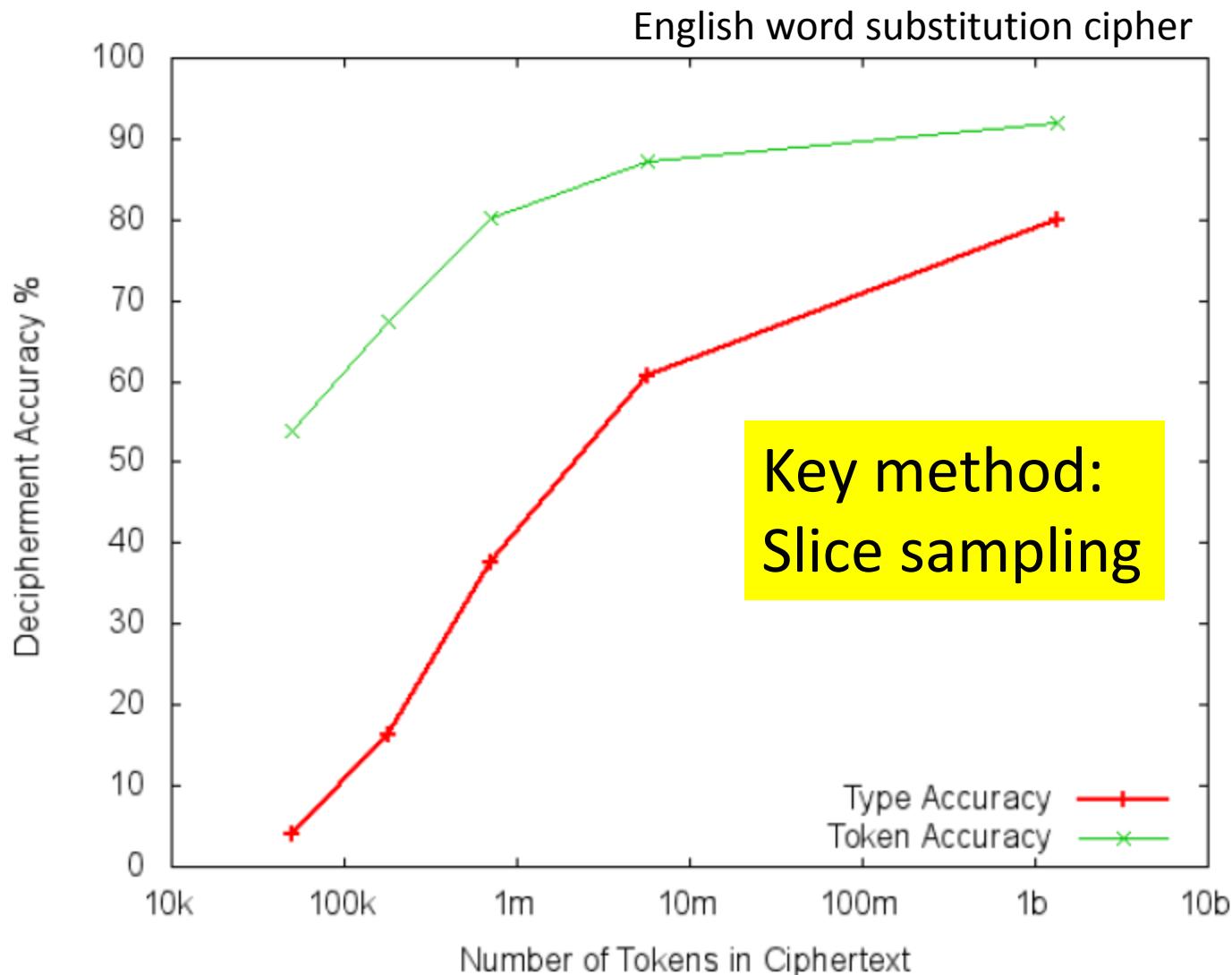
Deciphering Spanish Time Expressions

MT quality
on test set

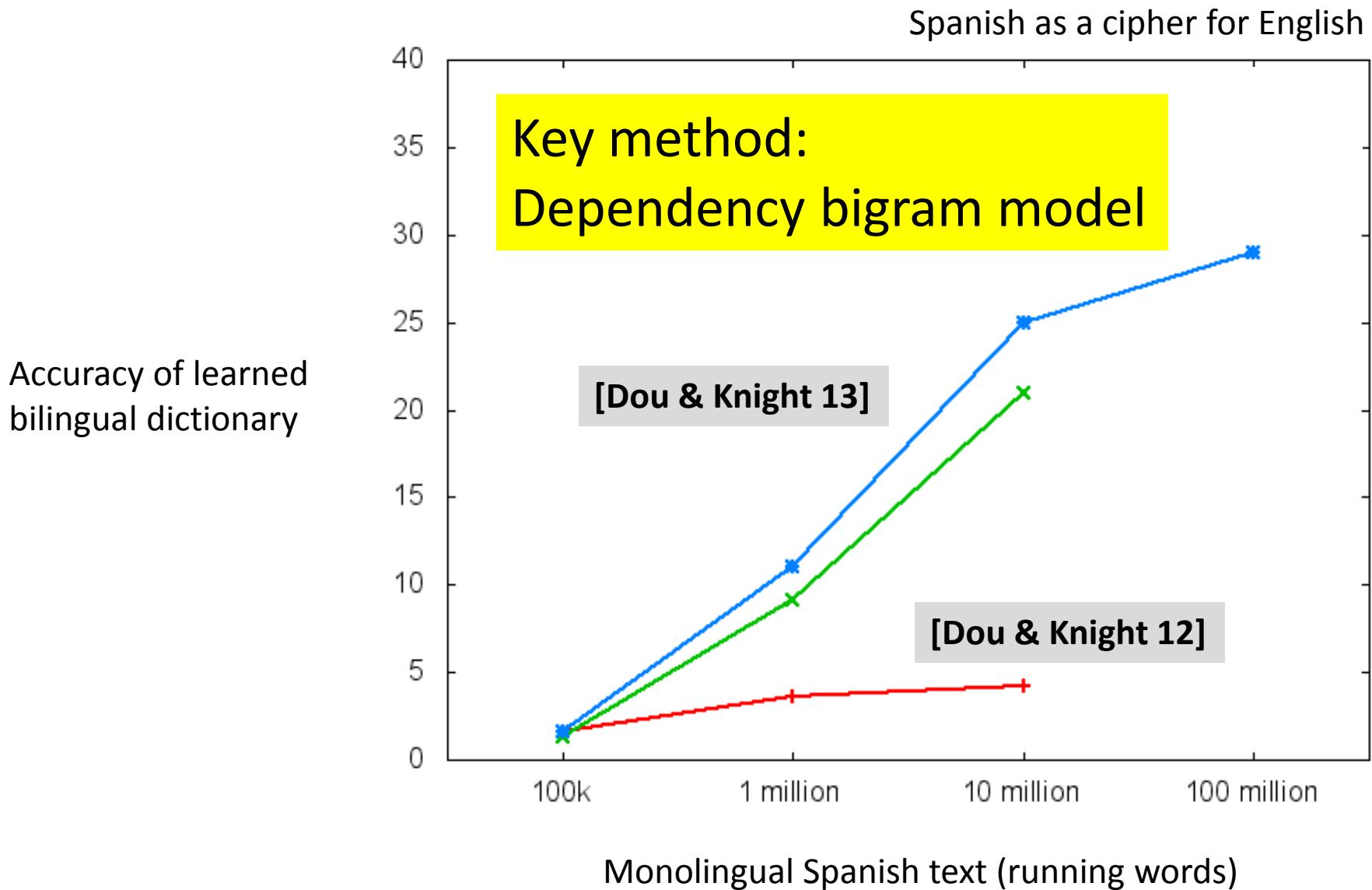
(Edit distance,
lower is better)



Scaling Up: Word Substitution

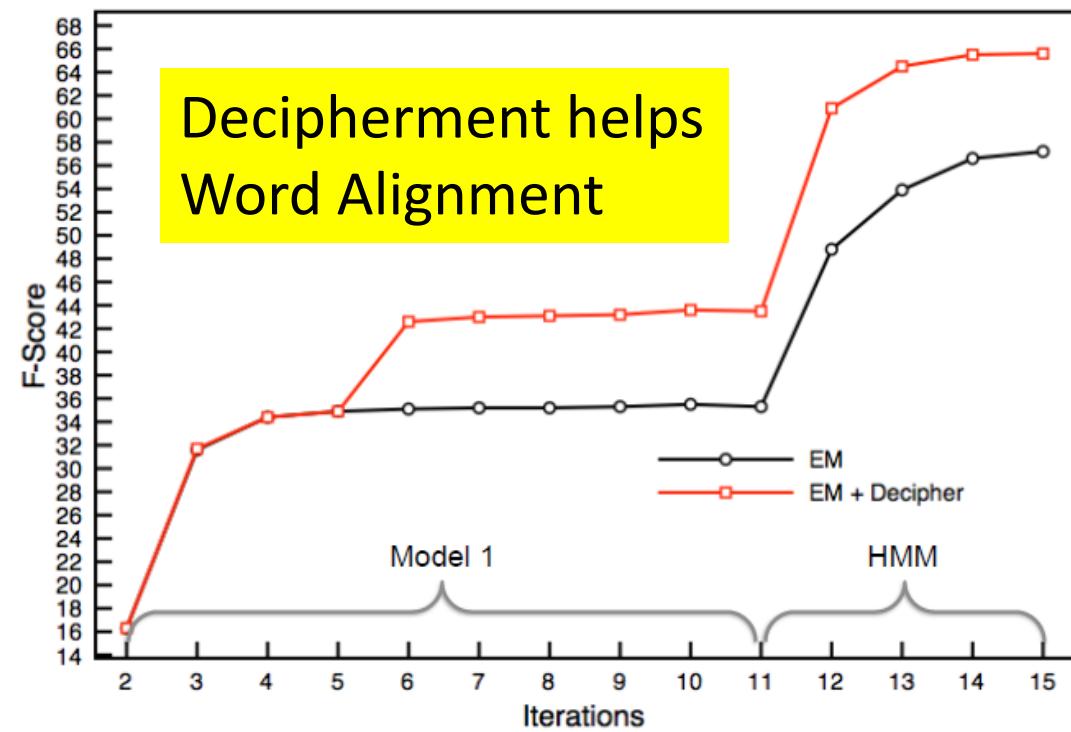


Scaling Up: Translation Dictionaries

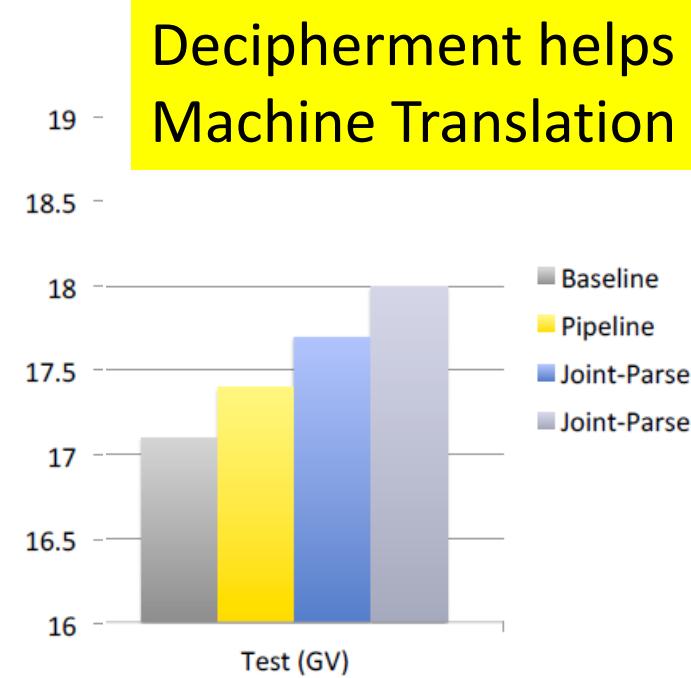


Scaling Up: Malagasy Translation

Small bilingual Malagasy/English text $\xleftrightarrow{\text{joint}}$ Large Malagasy monolingual text
(need to **align** words [Brown et al 93]) (need to **decipher** [Dou & Knight 13])



Decipherment helps
Word Alignment



Decipherment helps
Machine Translation

Plan for This Talk

- Break a series of codes
 - Simple letter substitution
 - Phonetic substitution
 - archaeology
 - transliteration
 - Word substitution
 - Foreign language as cipher
- Bonus 
 - Two historical ciphers
 - Final thought on translation and cryptography



YOU
ARE
HERE

Copiale Cipher

341

Αθηναῖς καὶ πάσῃ τῇ πόλει τὸν οὐρανὸν τοῦτον τοποθετεῖσθαι
τοῦτον τὸν οὐρανὸν τοποθετεῖσθαι τοποθετεῖσθαι τοποθετεῖσθαι
τοποθετεῖσθαι τοποθετεῖσθαι τοποθετεῖσθαι τοποθετεῖσθαι τοποθετεῖσθαι

'Néco+lin.

Caprodēiνοντάντας Ανάγκην προβογόνησιν

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

Copiale Cipher

Some scratch-outs, rare

Preview text fragments ("catchwords")

Δάλοι οι τίτλοι σε παραγράφους είναι πάντα με μεγάλες γράμματα.
Σύμφωνα με την αρχή της συγγραφής, η πρώτη λέξη της πρώτης σειράς παραγράφου πρέπει να έχει μεγάλη πρώτη γραμμή.

Επίσημη παραγράφηση:

Αθηναϊκός παραγράφος:

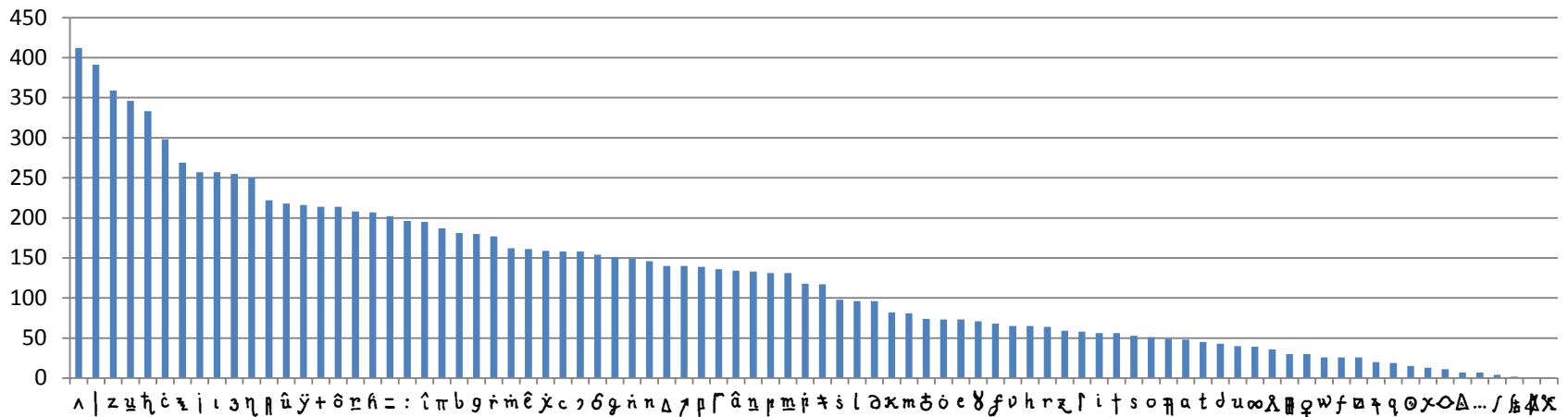
Ο παραγράφος είναι η μεγάλη παραγράφηση της συγγραφής, η οποία περιλαμβάνει την παραγράφηση της πρώτης σειράς παραγράφου και την παραγράφηση της δεύτερης σειράς παραγράφου.

Ο παραγράφος είναι η μεγάλη παραγράφηση της συγγραφής, η οποία περιλαμβάνει την παραγράφηση της πρώτης σειράς παραγράφου και την παραγράφηση της δεύτερης σειράς παραγράφου.

Paragraphs and section titles
always begin with
capitalized Roman letters.

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

Letter Frequencies



digraphs:

, 99

č : 66

† 49

: ፳፻፲፭ ፪፭፷፭

z R 44

trigraphs:

, 11 47

č : ፲፻ 23

η , ተ 22

ÿ, þ 18

h c | 17

tendencies:

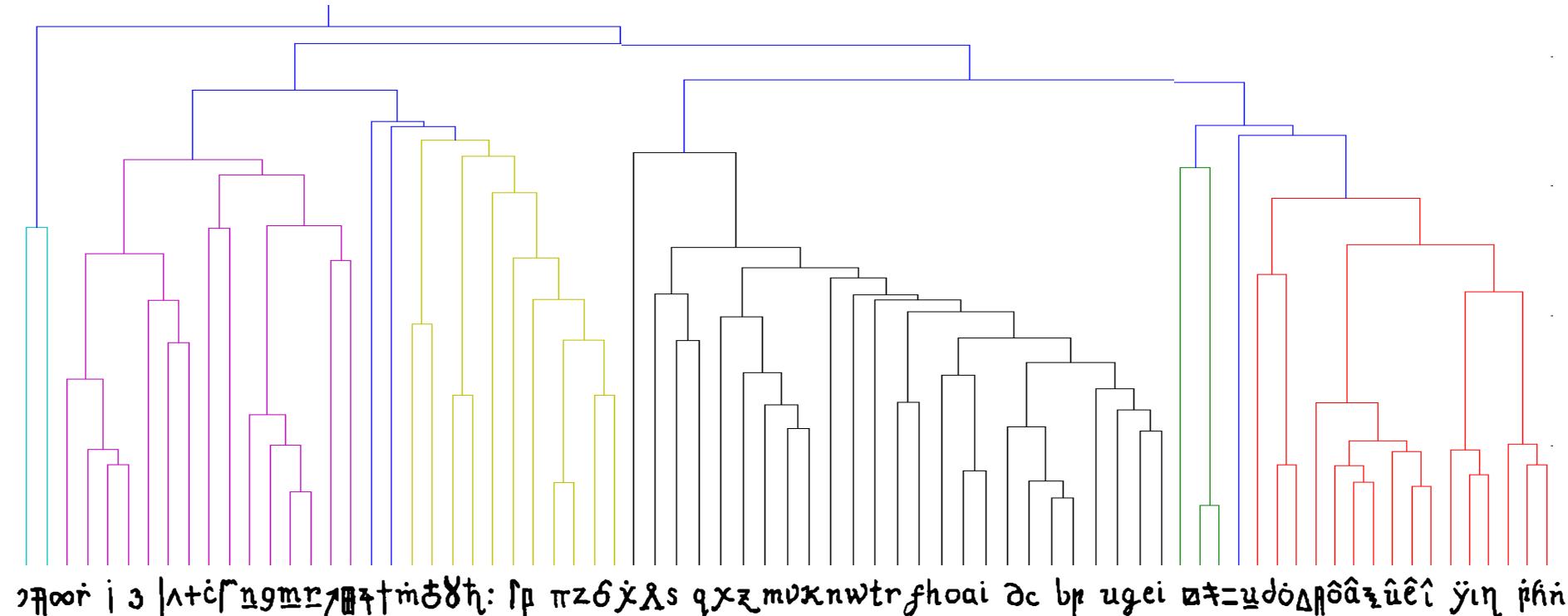
â, ê, î, ô, û followed by ʒ and j

\hat{a} , \hat{e} , \hat{i} , \hat{o} , \hat{u} preceded by z and π

Clustering of Cipher Letters

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

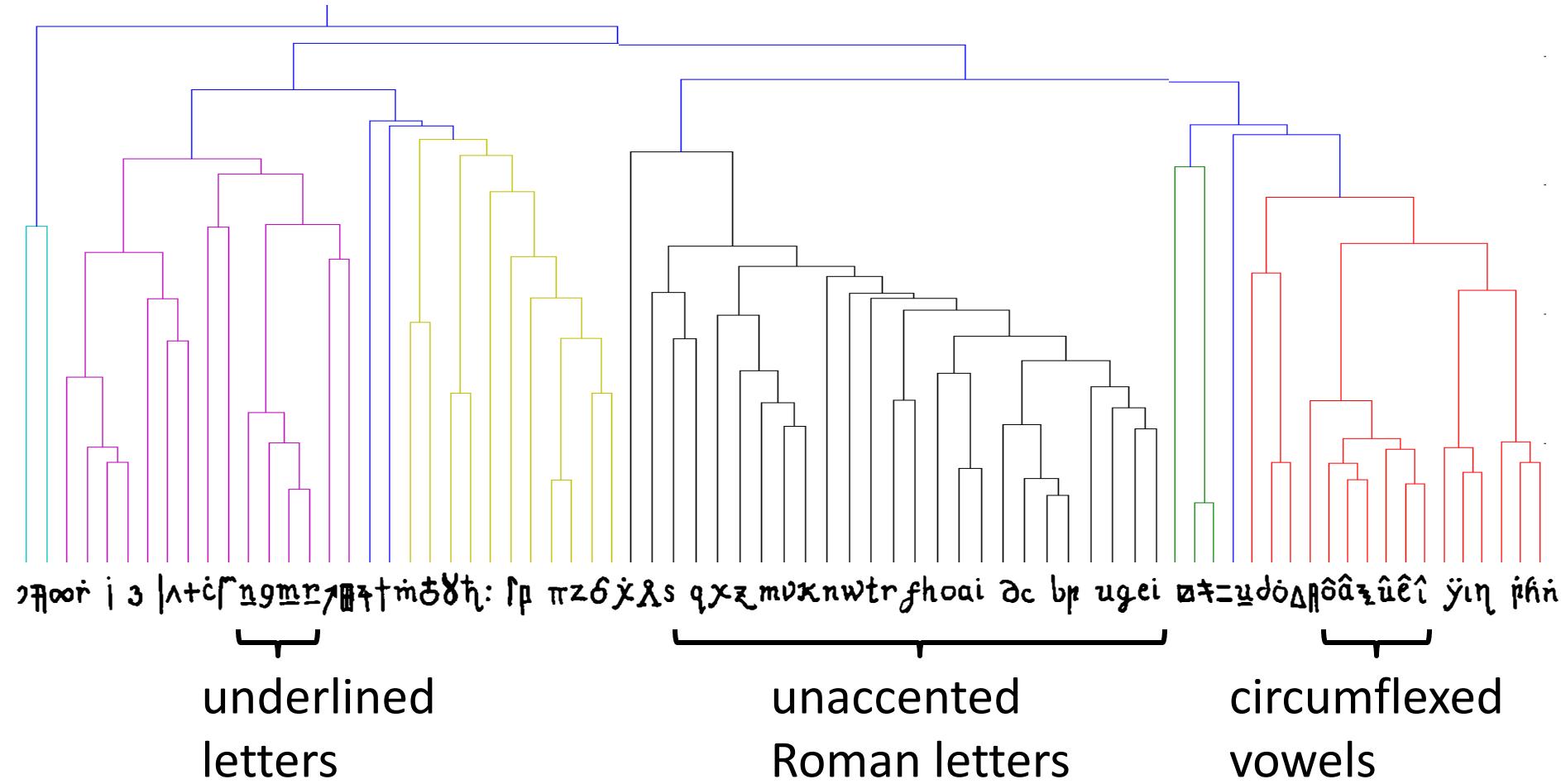
Scipy software



Clustering of Cipher Letters

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

Scipy software



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

**XfnglknaCbfzmk
lbuvCghtrhbkgnkn
fggnkbgbeCb ...**

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

Kmûr: pziô f| ý, hêi hêi x̄i l n pâz bAg z =
i xlz u pç c lâ r gKl h p i h p i l z i n p i | d r bR l a
= g z w p i y e c A r t d + b z ñ r i x y i l r z u f l z
p i x i d p i = f | u l s x m d m | x ñ g a | K h = l h | l x ô
o f : r i l b i l u m y j z v z â i x , p i m p i l h l c t ò o g g
z û t p k k n x e z g h l h p i h p i l i z t n x | r ô y m â
+ h h i r z ô z n b s ñ + : x r K p b R h d c û l g = n z K
p x o o n z p i f h i n r p z ê y ñ g = r p g t d a n z p | K
p ñ i l x y r i g p u l b g l i t b d û f p h i z o l e z ñ ô
f u r c f z q b n l b ñ h m a

*zfnqlknacbfxmk
lbuvcghtrhbkgnkn
fggnkbgbeeb ...*

D
8

80

FAIL

rst
.. languages.

Second Decipherment Approach

Homophonic cipher, e.g.:

A = Ȧ i l y r
B = Ȣ
C = ȶ n
D = ȴ
E = ȫ f ȭ Ȯ f Ȫ ȳ ȵ ȷ ȸ
F = ȣ
G = Ȭ



etc.

Homophonic Cipher

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

FAIL

But, slight numerical preference for
German

Cipher Characteristics

digraphs:

, \ddot{h} 99

\ddot{c} : 66

\ddot{h} ^ 49

: \ddot{u} 48

z R 44

trigraphs:

, \ddot{h} ^ 47

\ddot{c} : \ddot{u} 23

\ddot{h} , 22

\ddot{y} , \ddot{h} 18

\ddot{h} c | 17

tendencies:

$\hat{a}, \hat{e}, \hat{i}, \hat{o}, \hat{u}$ followed by \ddot{z} and $\ddot{\pi}$

$\hat{a}, \hat{e}, \hat{i}, \hat{o}, \hat{u}$ preceded by \ddot{z} and $\ddot{\pi}$



should appear
adjacent in German text

Make full digraph table for cipher and for German

Key Observation #1

In Copiale, \mathfrak{C} almost always followed by \mathfrak{H}

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

So guess: $\mathfrak{C} = C, \mathfrak{H} = H$

One Thing Leads to Another

ſt̄ = CH → ſt̄Λ = CHT → Λ = T ?

Each step is guesswork.

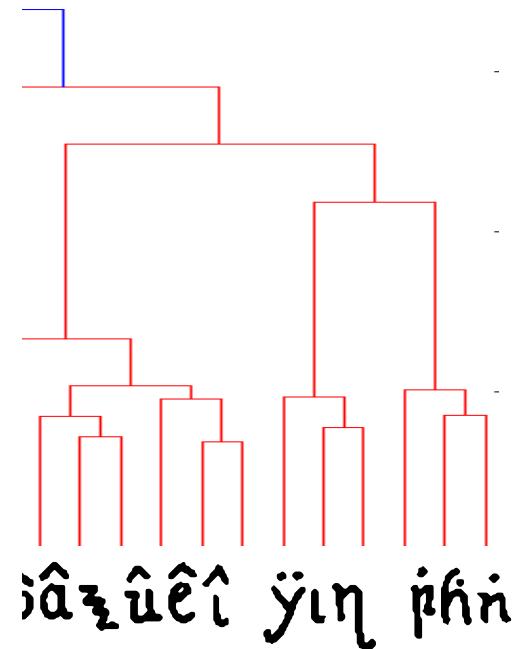
Must be willing to retract.

Weird task, not knowing German.

No longer care what the book says.

Cluster diagram crucial:

ÿ = I → ȳ = I , ȳ = I



Spring Break 2011

c aeiou fpy dlmrztbvw hkngs j gernin' ub

Spring Break 2011

Cipher
letters,
in groups

German letters

Other letters, groups

followed by

- 1 P t w
- 2 K N V M
- 3 T J P X S Z

not preceded by r

- a i g
- e h
- o u
- u e

preceded by r followed by vowels

- 1 A E I O U
- 2 E I O U
- 3 A E I O U

Grid

KN VM
 preceding vowels
 followed by vowels
 P A T V
 V ?
 Vowels: u ð ð R n t y - - - - -
 J D e {nri} - - - - -
 u {c b d} a p i g l l
 Döpikh
 " 3 m c f + = ≠ | Δ? : 8 x zu[u]f (svi) o {nrdl} o
 need: f g y l m z v w k s j
 (rare in
german)
 affricate
 lateral
 affricate
 (f) (t)
 u f sp
 nf pr
 nf pe
 (ff)

| | der | v | et | 2th |
|---|------|---|-----|---------|
| * | und | ✓ | e | é : u |
| | ein | ✓ | s i | hé l |
| | ung | | ich | g ? n |
| | ent | ✓ | ich | y n ? n |
| | rech | ✓ | che | 7th y |
| | sche | ✓ | t | ám ñ |
| | ech | ✓ | i | c l |
| * | die | ✓ | ett | + y = p |
| | rec | | sau | ? e k |
| | ne | | ds | l n |
| | gen | | eil | + c l n |
| | erit | ✓ | t e | l i w |
| | ver | | ti | z h 3 |
| | hen | | ter | z h f |
| | lic | | ler | z a t |
| | ten | | ew | u a t |
| | rei | | sch | á s |
| | nke | | st | h ? h |
| | anf | | sh | l ? h |
| | ede | | int | ly z |
| | and | ✓ | ein | g b z |
| | den | | rie | z n o |
| | run | | die | z n o |
| | ter | | de | yd l |
| | ere | | tu | r m n |
| | sei | | nd | = r n |
| | h te | | z | m Á l |
| | hei | | i | h i c |
| | nsg | | z | h i c |
| | ens | | i | h i c |
| | ment | | z | h i c |
| | h ne | | o | h i c |
| | ere | | ö | h i c |
| * | das | | te | k t p |
| | ide | | s | h e : |
| | ntie | | sch | k ? n |
| | nge | | a | c s |
| | gle | | che | z d u |
| | erk | | t | l a c |
| | ede | | | |

Spring Break 2011

Cipher
letters,
in groups

Quite a bit
of fooling
around →

German letters

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

| | | | |
|---|-----------|-------|---------|
| * | d e r v | c h e | z h a |
| * | u n d v | e | é : u |
| * | e i n v | s i | h i c t |
| * | u n g v | i c h | g y z n |
| * | c h t v | i c h | u n z n |
| * | r e h v | c h e | z h u |
| * | s i c h v | t | h m A |
| * | e c h v | i | c l A |
| * | e c h | e i t | + y h |
| * | d i e v | s a | h = f |
| * | r e c | d s | z z h M |
| * | u n f | e i t | t i n |
| * | g e n | t | c l n |
| * | e i t v | r e | t i n |
| * | v e r | t i | l i z |
| * | h e n | f e r | z R 3 |
| * | l i c | l e e | z R P |
| * | t e n | e w | z a t |
| * | r e i | a c i | z h z |
| * | n t e | s c h | h g z |
| * | a n f | s t | m h i |
| * | e d e | i c h | i z h |
| * | a n d v | e i n | h y E |
| * | d e n | t | o b A |
| * | r u n | i e | z z z |
| * | t e r | d i e | z n o |
| * | e f e | e | y d i |
| * | s e i | e r | z n u |
| * | h t e | h i | r : n |
| * | h e i | u n d | z = c T |
| * | n s g | z | m A i |
| * | e n s | i | k i c |
| * | m e n | t e t | z H A |
| * | b e n | z | z m z |
| * | e r e | i g | l C : |
| * | a d s | z g | k t m |
| * | r d e | z t e | h A n |
| * | n t e | z h e | h e i |
| * | n g e | s c h | k z n |
| * | l t e | z c s | z z z |
| * | a r e | z v u | z v u |
| * | e d e | t | l A C |

Grid

Spring Break 2011

Cipher
letters,
in groups

Quite a bit
of fooling
around →

German letters

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

German trigraphs

Cipher trigraphs

Grid

| | | | | |
|---|------|---|-----|------|
| * | der | v | cht | ztn |
| | Und | v | e | é:u |
| | ein | v | s | hct |
| | ung | v | ich | y?n |
| | cht | v | ich | n?n |
| | reh | v | che | ?n |
| | sich | v | t | má |
| | chev | v | i | c> |
| | ech | v | et | +g> |
| * | die | v | sa | h= |
| | rek | v | ds | ?zhm |
| | nde | v | ei | |
| | gen | v | t | |
| | eit | v | r | |
| | ver | v | a | |
| | hen | v | te | |
| | lic | v | de | |
| | ten | v | ew | |
| | rei | v | | |
| | nte | v | | |
| | ant | v | | |
| | ede | v | | |
| | and | v | | |
| | den | v | | |
| | run | v | | |
| | ter | v | | |
| | elle | v | | |
| | sei | v | | |
| | hte | v | | |
| | hei | v | | |
| | nsg | v | | |
| | ens | v | | |
| | men | v | | |
| | hnt | v | | |
| | ere | v | | |
| | ere | v | | |
| | das | v | | |
| | nde | v | | |
| | nte | v | | |
| | nge | v | | |
| | ite | v | | |
| | sch | v | | |
| | hzn | v | | |
| | z | v | | |
| | má | v | | |
| | hlc | v | | |
| | z | v | | |
| | hac | v | | |

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

Kmûr:rzlôf|y, hêi hziln pâzba g z= iplz u kôc lârg k l h p r h p l z i n p u l d r h l a = g z w p y ê c A r t ð a + b z q r i x y j i r z u f l z p t i p d r i = f | u l s t p l d m | z n g â | k h = l h | l x ô o f : r i l b i f u t y j z v z â j x , r p l p i z h l c t ð o g g z û t p f p n x e z g h l h p i h p l i z t h t | r ô y m â + h h r z ô z n b s n t : z r k p d h h d c n l g = n z k p z o o n z p z f h i n r p z e y n g = r p g t ð a z n z p k p n j i x y r i g p u l b g l i t b d n f p h h z ô l e z n ô f i n r c f z d n l b n h h m d

Actually, those are space bars

Copiale Decipherment

lit:mz||bl
v̄x̄|̄l̄s̄k̄r̄t̄|w̄n

ποιηταρεωνα=γλυκ̄ουρθ̄

δημι+σηματηηγ̄f.

cūl̄ēt̄p̄t̄ḡl̄:k̄

δχ̄ηη+τ̄ρ̄t̄l̄:ȳγ̄l̄t̄ōq̄z̄īx̄āj̄ēc̄:ūd̄.
fūr̄f̄k̄l̄l̄:c̄.

μηρ̄t̄+Δḡȳūx̄z̄ōm̄n̄f̄γ̄h̄h̄+f̄.

κην̄r̄p̄z̄īōf̄j̄ȳt̄h̄ēīh̄īl̄ōp̄āz̄b̄Δḡz̄īl̄z̄ūp̄f̄c̄l̄āt̄ḡk̄l̄t̄
π̄t̄h̄f̄l̄īn̄p̄īl̄d̄r̄f̄l̄āz̄ḡw̄p̄ȳēc̄Δr̄t̄d̄+b̄z̄r̄t̄x̄ȳīr̄z̄ūf̄λ̄z̄
π̄t̄īḡd̄ī=γ̄l̄n̄l̄s̄t̄l̄n̄|̄l̄ḡā|̄x̄t̄=l̄t̄|̄l̄x̄d̄w̄:r̄īl̄b̄īl̄ūm̄ȳj̄z̄
z̄āl̄x̄r̄t̄p̄īz̄h̄l̄c̄t̄ōḡz̄ū+r̄t̄p̄t̄x̄ēz̄ḡh̄l̄t̄f̄īh̄l̄l̄īz̄t̄n̄f̄r̄ȳ
m̄āt̄h̄h̄r̄z̄ōn̄b̄s̄h̄t̄:r̄t̄k̄r̄p̄f̄h̄d̄c̄ūl̄ḡ=ūz̄k̄r̄z̄ōn̄z̄f̄h̄r̄p̄
z̄ē̄ȳnḡ=r̄p̄ḡt̄d̄ēn̄z̄p̄|̄k̄p̄h̄īx̄ȳr̄īḡp̄ūλ̄b̄ḡl̄īt̄b̄d̄n̄f̄f̄h̄z̄ōlē
z̄n̄l̄īr̄c̄z̄f̄n̄l̄b̄h̄t̄d̄:

n̄īr̄c̄īḡēōp̄h̄ōr̄d̄p̄z̄d̄n̄h̄z̄āk̄Ōḡs̄=d̄m̄ēj̄z̄ūl̄

ḡs̄m̄n̄ōl̄d̄ūp̄ēḡūd̄t̄r̄d̄īz̄ḡt̄r̄x̄ūp̄n̄ēk̄Θ̄n̄p̄h̄b̄=n̄
λ̄ēt̄m̄ōḡ:ūā=r̄z̄l̄d̄h̄r̄p̄m̄h̄āz̄d̄j̄l̄r̄n̄ī6̄ḡȳl̄ūz̄īēḡc̄ūh̄
r̄s̄ēn̄λ̄
c̄p̄=f̄h̄ūb̄t̄d̄c̄:z̄d̄.

hz̄j̄l̄:ō̄ḡēz̄n̄h̄āz̄Ān̄p̄āz̄Θ̄p̄s̄t̄t̄īz̄t̄m̄ȳb̄ūl̄ī=n̄p̄
z̄p̄s̄=n̄h̄f̄r̄z̄p̄t̄z̄īp̄ūp̄c̄ȳh̄f̄b̄l̄īx̄n̄ḡh̄īȳt̄īb̄s̄=b̄t̄r̄ūd̄
j̄d̄l̄r̄:n̄l̄ūn̄ōn̄.
l̄p̄ōz̄f̄īh̄ḡz̄ȳp̄r̄l̄īt̄l̄t̄d̄r̄l̄l̄ōēz̄īr̄.

gesetz buchs

der hoherleuchte ◊ e ◎

geheimer theil.

erster abschnitt

geheimer unterricht vor die gesellen.
erster titul.

ceremonien der aufnahme.

wenn die sicherheit der Δ durch den ältern

thürheter 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 geführet dieser frägt ihn:

erstlich ob er begehre ◊ zu werden

zweytens denen verordnungen der Θ sich

unterwerffen und ohne wiederspenstigkeit die lehrzeit
ausstehen wolle.

drittens die Δ der Θ gu verschweigen und dazu
auf das verbindlichste sich anheischig zu machen
gesinnet sey.

der candidat antwortet ja.

Copiale Decipherment

lit:mzplbl
vix{t̪as̪kɒŋg̪wɪn

ποίησέ γέγονα = πλήρης θεός

ԾԱՀԿԱՏՈՂՊՐԼԻՆԸ. ՀԱՅԻՆ
ԸՆԴՀԵՅՏԻՐՐԴԳՈՒԱԿ

Ճշնհեղ+քրիտլնչ: յրաւէօլզւ իշալիք: Ա. Գործութաւունչ:

መግለጫ+Δግሂያዥዘዕዢሙናቸውንከተማቸው.

κατηγορίας οι οποίες πρέπει να γίνουν στην πλατφόρμα της Διεύθυνσης Αστυνομίας Κρήτης. Η πλατφόρμα αυτή θα παρέχει στην αρχή της διάρθρωσης την ενημέρωση για την προσέλευση των αστυνομικών στην περιοχή, την επιβολή της αναγνωρίσεως και την επιβολή της απόδοσης της αστυνομικής δύναμης. Στην πλατφόρμα θα μπορεί να γίνεται η απόδοση της αστυνομικής δύναμης στην περιοχή, η επιβολή της αναγνωρίσεως και η επιβολή της απόδοσης της αστυνομικής δύναμης.

First lawbook of the e

Secret part.

First section

Secret teachings for apprentices.

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 

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.

Copiale Decipherment

lit:mzplbl
v̄x̄z̄/̄l̄s̄k̄p̄z̄/w̄n http://
π̄ōīn̄t̄h̄Δ̄ḡēz̄c̄â̄=̄ḡl̄ūb̄ Ω̄ūr̄Ω̄z̄
δ̄p̄t̄z̄ī+̄ōīn̄p̄r̄l̄t̄h̄īn̄c̄ f̄.
c̄ūīf̄ēz̄t̄īp̄t̄ḡn̄l̄:κ̄
δ̄x̄ūt̄ēn̄t̄+̄z̄r̄p̄t̄m̄l̄īz̄:̄ȳp̄l̄t̄s̄ōīq̄z̄īx̄â̄j̄ēc̄:̄ūd̄.
f̄ūr̄p̄īκ̄l̄ūl̄z̄c̄.
m̄p̄r̄â̄+̄d̄ḡȳūx̄z̄ōz̄m̄n̄t̄Γ̄ūh̄t̄+̄īf̄.

Documentation at
<http://stp.lingfil.uu.se/~bea/copiale>
Google: “copiale”

First lawbook of the e

Secret part.

First section

Secret teachings for apprentices.

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 

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.

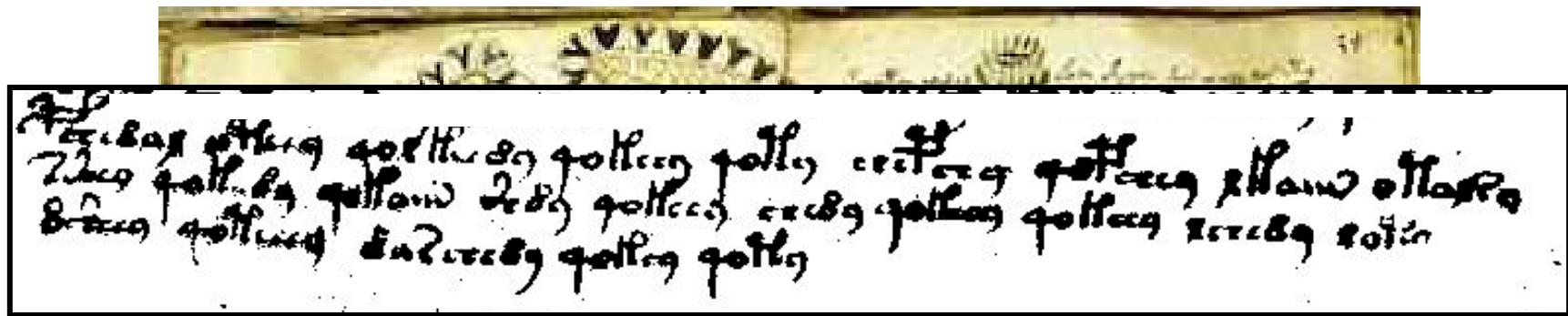
The candidate answers yes.

Voynich Manuscript



- Medieval illustrated manuscript
- 235 pages, 6 sections, 38k word tokens, 35 letter types
- Undeciphered

Voynich Manuscript



¶cc8ae 0¶cc9 40¶ FCC89 40¶cc9 40¶9 cc¶cc9 40¶cc9 8¶rall 0¶a¶?9
?cc9 40¶cc89 40¶rall 2¶89 40¶cc9 cc¶cc9 40¶cc9 40¶cc9 8cc89 20¶9
8cc9 40¶ccc9 8a2cc89 40¶cc9 40¶9

BSC8AE OPCC9 4OE FCC89 40FCC9 4OP9 SCBS9 4OBSC9 EFAM OPAE29
2ZC9 4OFC89 4OFAM Z89 4OFCC9 SC89 4OFCC9 4OFCC9 ESC89 EOP9
8ZC9 4OPCCC9 8ARSC89 4OFC9 4OP9

Voynich Letter Substitution

Latin letter trigram model



quo_vade_bre ...



a → {all Voynich letters}

b → {all Voynich letters}

c → {all Voynich letters}

...

z → {all Voynich letters}

_ → _



v A S 9 2 _ 9 F A E _ A R _ A P A M _ ...

| Input | Decipherment |
|---------------|---------------|
| VAS92 9FAE AR | quiss squm is |
| APAM ZOE ZOR9 | onum pom |
| QOR92 9 FOR | quss hates s |
| ZOE89 ... | qum hatis ... |

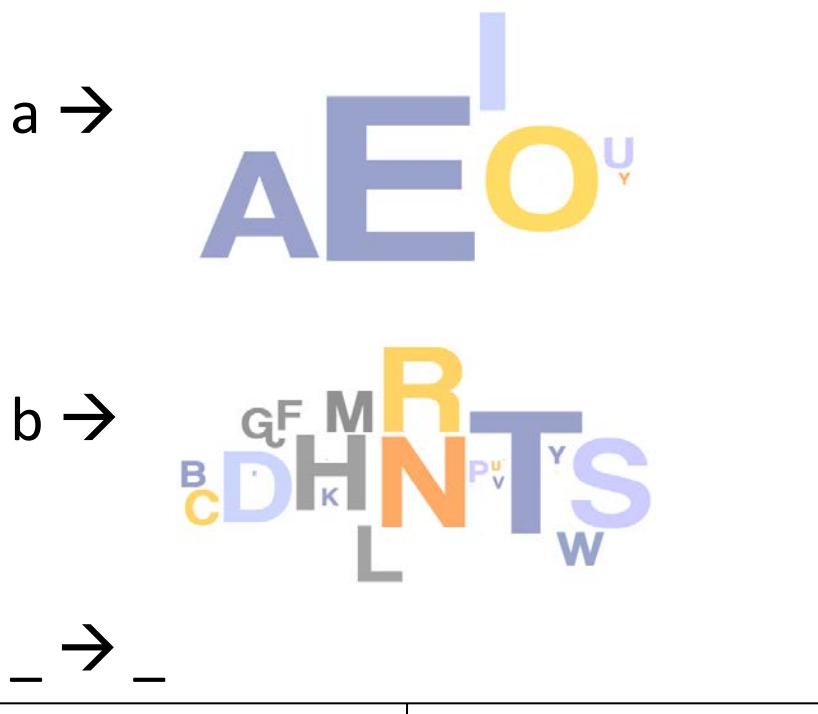


Letter Clustering

Trigram model over {a, b, _ }



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



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}



_ → _



V A S 9 2 _ 9 F A E _ A R _ A P A M _ ...

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 _ ...

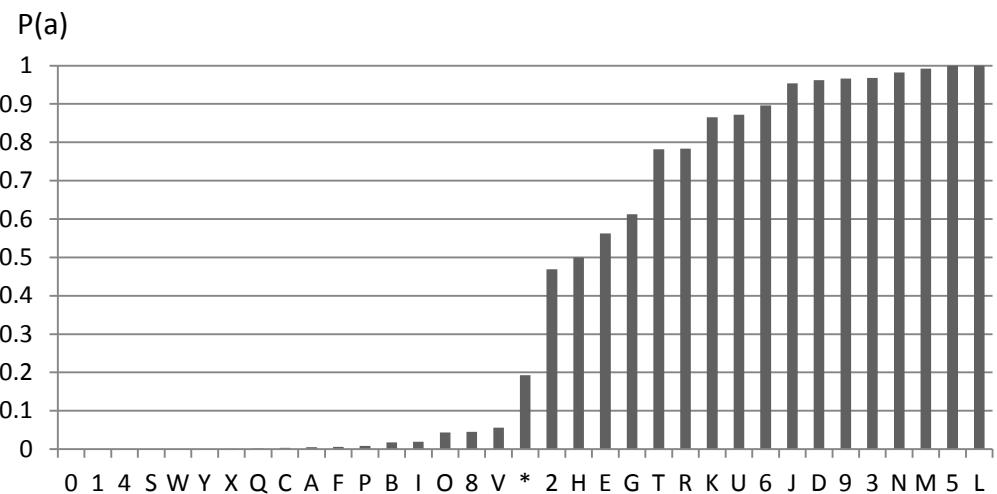
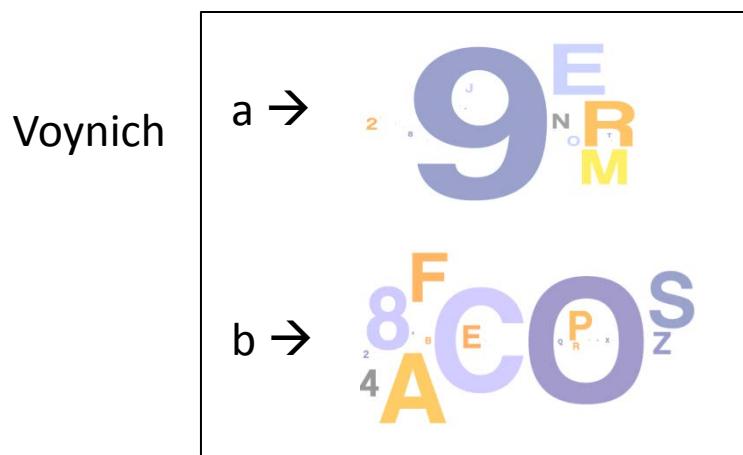
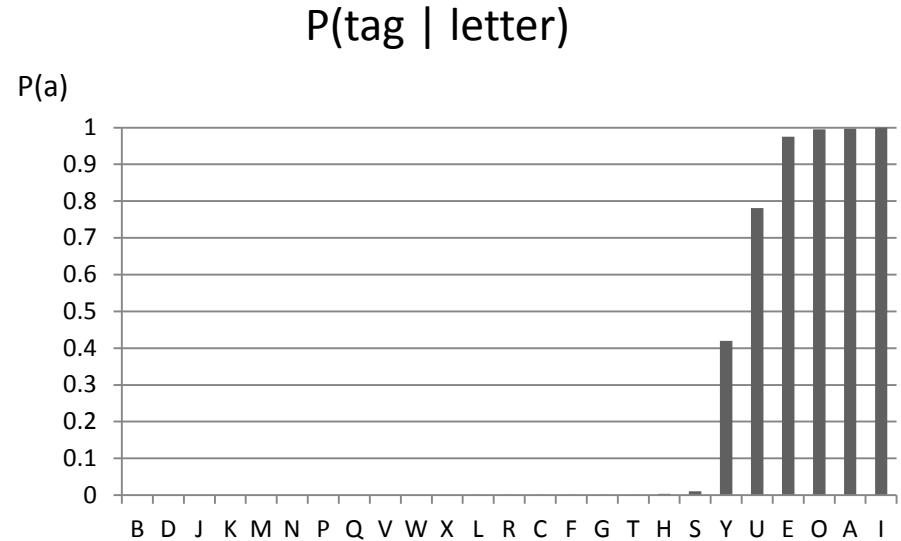
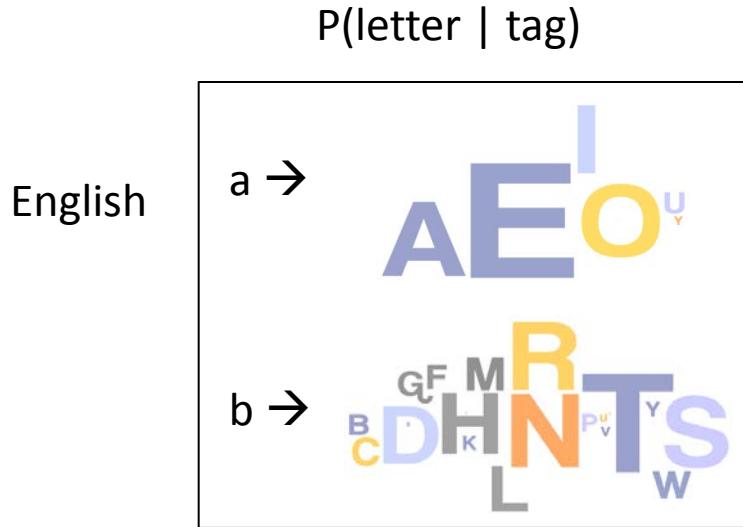


Sample tagging with learned model:

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

Bigram model over {a, b}

1

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

1

a →



b →



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

Sample tagging with learned model:

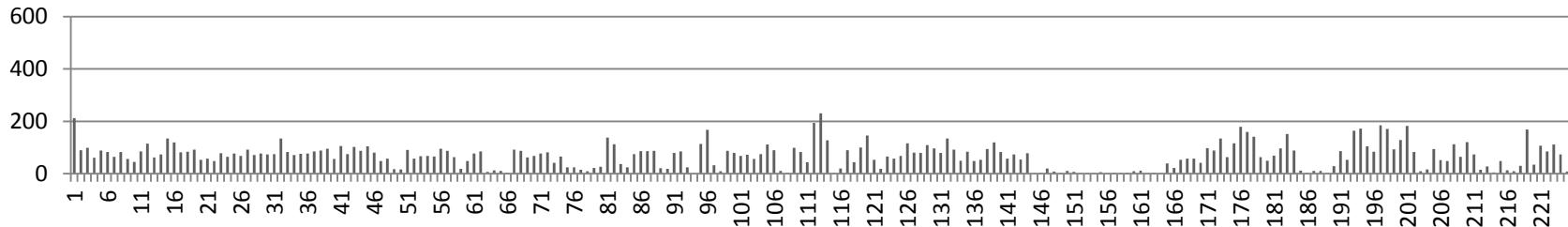
a a a a a
VAS92 9FAE AR APAM ZOE

a a a a a ...
ZOR9 ORC2 9 FOR ZOE89 ...



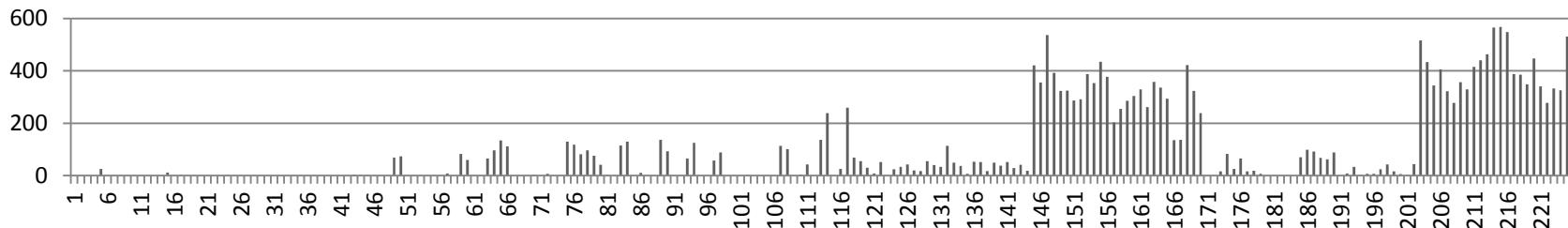
Word Clustering

Voynich words tagged as “a”



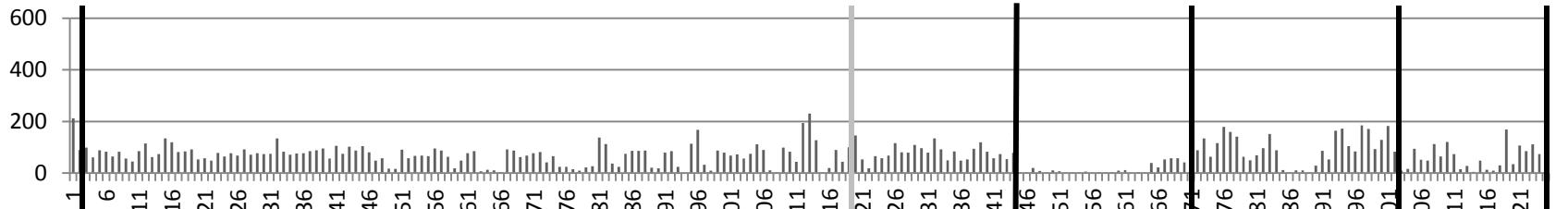
← pages →

Voynich words tagged as “b”



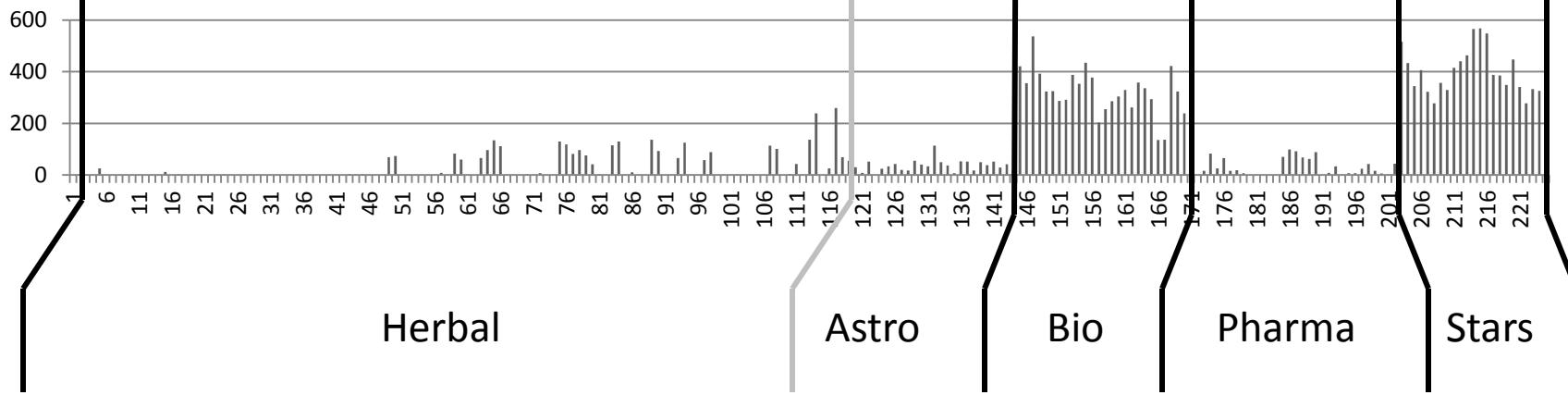
Word Clustering

Voynich words tagged as “a”



← pages →

Voynich words tagged as “b”



Voynich sections, per drawings observed.
Captain Currier’s “two languages” (1976).

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.

s a general Analyses), paper in the mmer. Mr. King on his

1970s NSA 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 [redacted] of RSI, the name is pronounced "however you like" [8]. The technique itself and its uses are classified Top Secret Codeword.

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]

National Security Agency

NSA applies statistics to ciphers, codes, and other language processing problems

NSA employs more mathematicians and linguists than any other organization.

NSA has more computers than any other organization.



Oh yeah -- we've been
on Mars since 1962!

“Slacker”

Association for Computational Linguistics

1950s

1960s

1970s

1980s

1990s

2000s

2010s

2020s

??

1970s paper on
HMM Voynich

1993 paper
on Statistical
Machine
Translation

2011 paper on
HMM Voynich

ACL applies statistics
to language processing
problems

10-Class Word Clustering: English

The image shows a large purple bar chart with white text labels on the bars. The chart consists of four vertical bars of increasing height from left to right. The first bar has the word 'MY' in large grey letters at the top, with 'ITS' in smaller blue letters above it. The second bar has 'THIS' in yellow at the top, with 'your' in smaller blue letters below it. The third bar has 'OUR' in white at the bottom. The fourth bar has 'AN' in blue at the top, with 'HIS' in white below it. The background is a light purple color.

etc

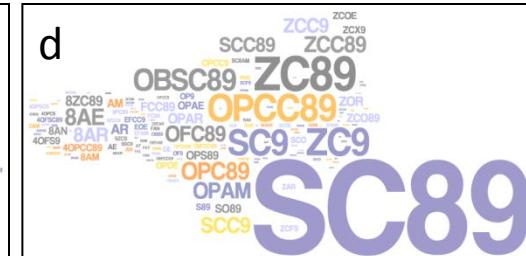
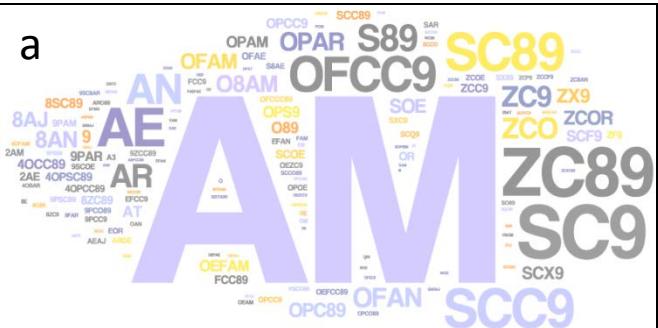
etc

etc

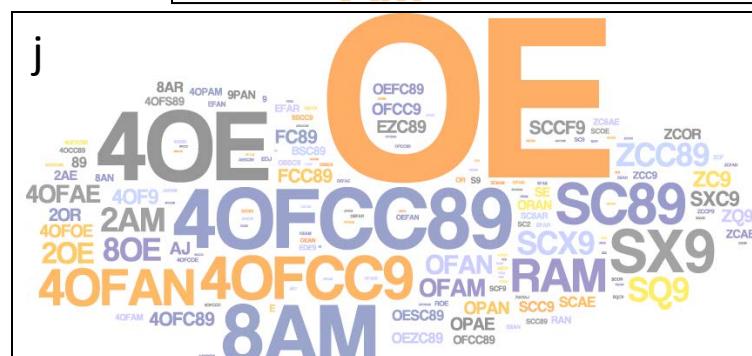
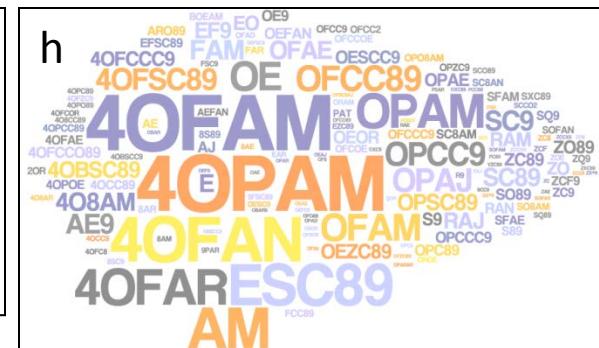
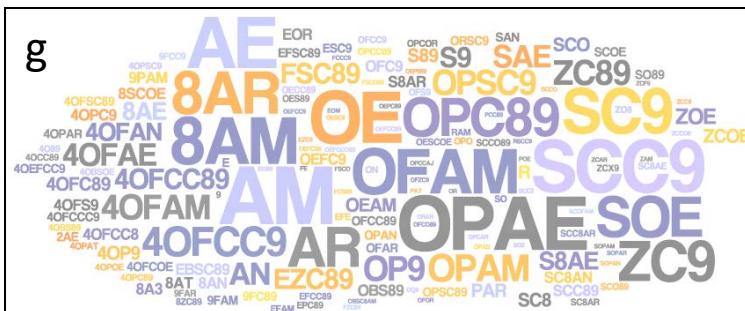
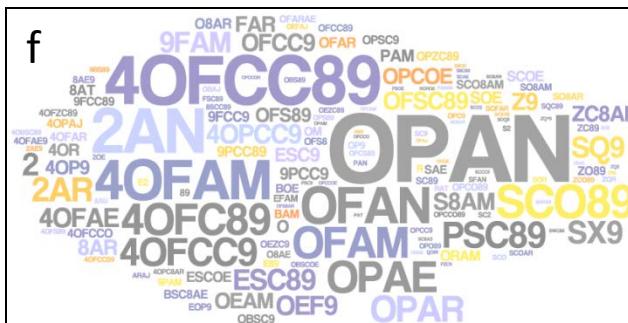
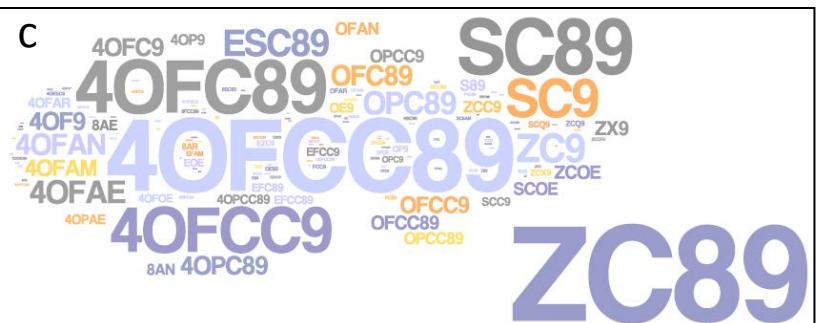
etc

etc

etc



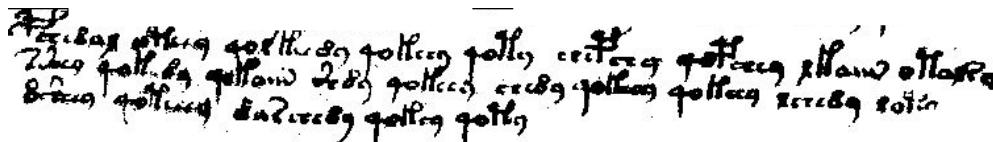
10 Classes: Voynich-B



Other Unsolved Ciphers

Voynich Manuscript (1400)

Reddy &
Knight 11

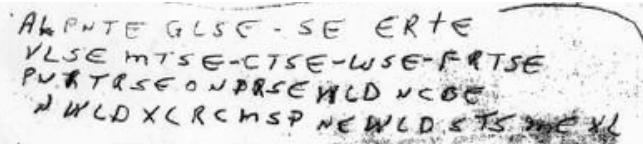
A block of handwritten text in the Voynich script, which is believed to be a constructed language. The text is written in two columns of approximately 15 characters each.

Zodiac Killer (1967)

Ravi & Knight 11

| | | | | | | | | | | | | | | | | | |
|---|---|-----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Δ | □ | P | / | Z | / | U | ▀ | B | ▀ | Λ | Ο | R | Λ | Ψ | X | Λ | Β |
| W | V | + E | G | Y | F | Ω | Δ | H | R | □ | K | I | Ψ | Y | E | | |
| M | J | Y | Λ | U | I | K | Δ | Ρ | T | Τ | N | Θ | Y | D | Ω | Ω | |

FBI cipher (1999)

A block of handwritten text in a cursive script, appearing to be a cipher message. The text is written in several lines and includes some smudges and corrections.

Kryptos (1990)

OBKR

UOXOGHULBSOLIFBBWFLRVQQPRNGKSSO
TWTQSJQSSEKZZWATJKLUDIAWINFBNYP
VTTMZFPKGDKZXTJCDIGKUHUAUEKCAR

Lost Language Decipherment



Snyder, Barzilay, Knight 10

Unsupervised Training for NLP

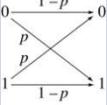
Machine Translation ...

Plan for This Talk

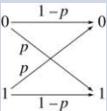
- Break a series of codes
 - Simple letter substitution
 - Phonetic substitution
 - archaeology
 - transliteration
 - Word substitution
 - Foreign language as cipher
- Bonus
 - Two historical ciphers
 - Final thought on translation and cryptography



Future Prospects for Translation

| | Cryptography | Translation |
|---|--|---|
| Manual |  | Manual encoding |
| Mechanical |  | 1920s Mechanical encoding;
intuition-based decryption |
| Mathematical |  | 1950s Computer decryption,
based on information theory |
| Higher math,
deeper
understanding |  | 1980s Public-key systems,
based on number theory |

Future Prospects for Translation

| | 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
??? ??? ??? |

end

Zodiac Killer Ciphers

Zodiac 408 (solved, 1969)

Δ □ P / Z / 4 B □ K O R A X 9 X B
W V + E G Y F O D H R □ K I R Y E
M J U L A N I K A R T N Q Y D S
S F / Δ □ B R O R A U □ F R A R E
K A L M Z D J Z R \ 9 F H V W E A Y
□ + @ G D A K I - O R X A @ - S F
R N I T Y E L O □ R G B T Q S □ B
L D / P □ B □ X R E H M U A R R K

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



Zodiac 340 (still unsolved)

H E R > 9 L A V R K I O L T G O D
N P + B F □ O D W Y < □ K F □
B Y I M + M C I Y S 9 A D V 0 + + R K O
D M + T M D I D T + F P + P O K /
P A R A F > O L F - O C F > O D F
■ O + K Q □ I O C X G V . L I
Φ G O T Z J D O + D N Y F + O L A
D < M + 8 + Z R O F B C X A O O K
- O L V + L J + O 9 A < F B Y -
U + R / O T E I D Y B 9 8 T M K O
O < J R J I □ O T O M . + P B F
Φ O D S Y □ + N I O F B C F I A R
J G F N V F O B O C . B O C .
Y B X O E O D C E > V U Z O - +
I O C I O F B K F O 9 L . F M Q G O
R C T + L E O C < + F J W B I O L
+ O W C F W C P O S H T / F O P
I F W < D A B T D O Y O B
C - C S K P N H M >
+ I A I K I

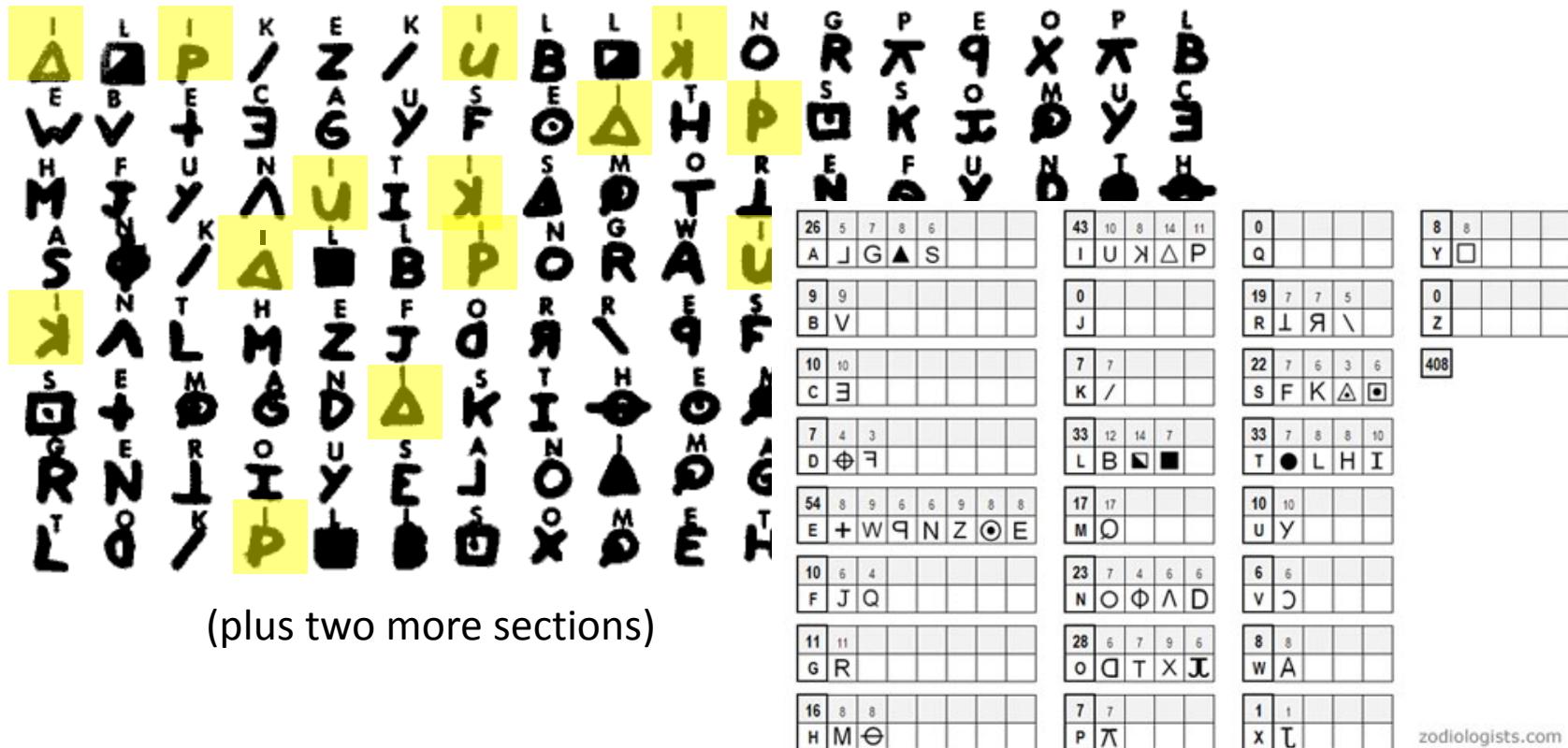
A circle with a crosshair symbol drawn below the cipher text.

COOP-SFPO
1546-78
7-14-78 GML
7-17-78

#2 11-9-69

Zodiac Serial Killer

408-letter cipher (solved):



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 DANGEROUUE 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 THAЕ 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”

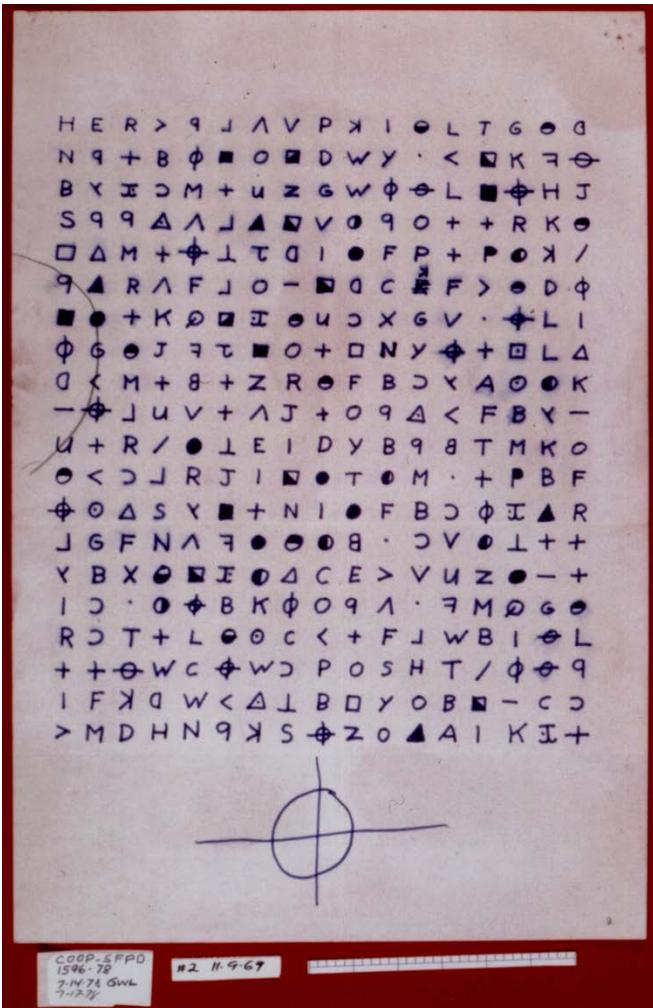
Deciphering Zodiac 408

Bayesian models

Extended Carmel finite-state toolkit to do Bayesian inference.
[Chiang et al 10]

| Language Model | Initial Sample | Decipherment Error |
|--|-----------------|--------------------|
| 3-gram | Random | 62.3 / 48.5 / 47.4 |
| 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 |

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.

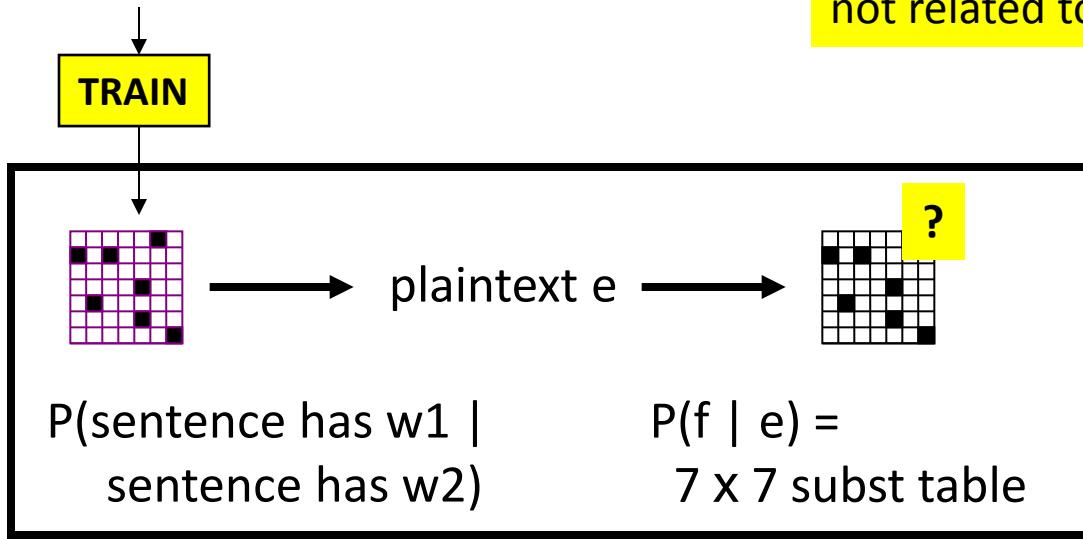
Decipherment Papers by ACL-ers

- "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.

Word Substitution Cipher

.....France.....Britain.....Canada...
.....Mexico.....Indonesia.....Malaysia...
.....Britain.....Canada.....Australia...
.....Britain.....France.....Indonesia.....
....Mexico.....Australia.....France...
...Britain.....

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

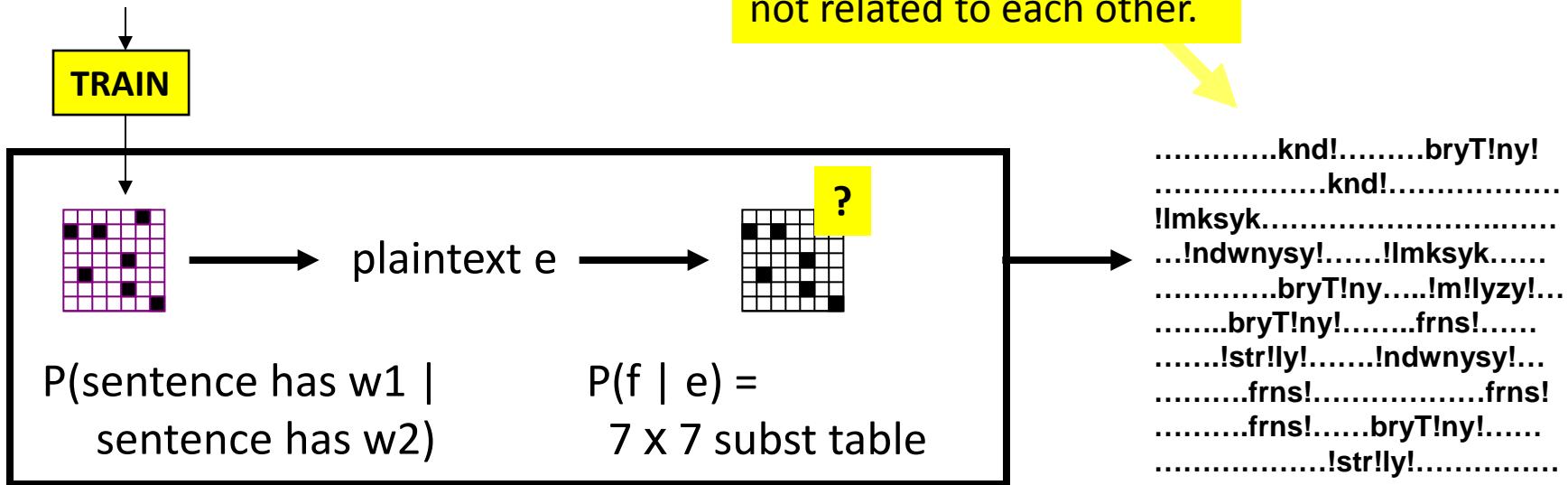


.....knd!.....bryT!ny!
.....knd!.....
!lmksyk.....
...!ndwnysy!.....!lmksyk.....
.....bryT!ny.....!m!lyzy!...
.....bryT!ny!.....frns!.....
.....!str!ly!.....!ndwnysy!...
.....frns!.....frns!.....
.....frns!.....bryT!ny!.....
.....!str!ly!.....

Word Substitution Cipher

.....France.....Britain.....Canada...
.....Mexico.....Indonesia.....Malaysia...
.....Britain.....Canada.....Australia...
.....Britain.....France.....Indonesia.....
....Mexico.....Australia.....France...
...Britain.....

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



| | | | | | | | |
|-----------|---|-----------|--------|-----------|--------|----------|--------|
| Australia | → | !str!ly! | (0.93) | !ndwnysy! | (0.03) | m!lyzy! | (0.02) |
| Britain | → | bryT!ny! | (0.98) | !ndwnysy! | (0.01) | !str!ly! | (0.01) |
| Canada | → | knd! | (0.57) | frns! | (0.33) | m!lyzy! | (0.06) |
| France | → | frns! | (1.00) | | | | |
| Indonesia | → | !ndwnysy! | (1.00) | | | | |
| Malaysia | → | m!lyzy! | (0.93) | lmksyk | (0.07) | | |
| Mexico | → | !lmksyk | (0.91) | m!lyzy! | (0.07) | | |

[Knight et al 06]