

# Natural Language Processing Research

USC CS Annual Research Review  
March 23, 2010

Kevin Knight

Research Associate Professor  
Computer Science Department  
USC

Senior Research Scientist and Fellow  
Information Sciences Institute (ISI)  
USC

# What is Natural Language Processing?

- Challenges:
  - Understand what people say!
  - Generate utterances that make sense!
- Valuable applications:
  - Language Translation
  - Dictation
  - Question answering
  - Dialogue
  - Summarization
  - Creative language generation

# Natural Language Research at USC

Computer Science Faculty

Institute for Creative Technologies (ICT)

Information Sciences Institute (ISI)

Electrical Engineering Department (EE)

**M. Arbib**  
**D. Chiang**  
**A. Gordon**  
**J. Gratch**  
**J. Hobbs**  
**E. Hovy**  
**K. Knight**  
**A. Leuski**  
**D. Marcu**  
**S. Marsella**  
**S. Narayanan**  
**P. Pantel**  
**F. Sha**  
**W. Swartout**  
**D. Traum**



# Natural Language Research at USC

# Computer Science Faculty

## Text and gesture (Stacy Marsella)

## Origins of language (Michael Arbib)

New research faculty in 2010  
(**Liang Huang** and **Kenji Sagae**)

**M. Arbib**  
**D. Chiang**  
**A. Gordon**  
**J. Gratch**  
**J. Hobbs**  
**E. Hovy**  
**K. Knight**  
**A. Leuski**  
**D. Marcu**  
**S. Marsella**  
**S. Narayanan**  
**P. Pantel**  
**F. Sha**  
**W. Swartout**  
**D. Traum**



# Natural Language Research at USC

Computer Science Faculty

See **Andrew Gordon's** text mining research in last week's Economist magazine.

M. Arbib  
D. Chiang  
A. Gordon  
J. Gratch  
J. Hobbs  
E. Hovy  
K. Knight  
A. Leuski  
D. Marcu  
S. Marsella  
S. Narayanan  
P. Pantel  
F. Sha  
W. Swartout  
D. Traum





# Natural Language Research at USC

Computer Science Faculty

M. Arbib  
D. Chiang  
A. Gordon  
J. Gratch  
J. Hobbs  
E. Hovy  
K. Knight  
A. Leuski  
D. Marcu  
S. Marsella  
S. Narayanan  
P. Pantel  
F. Sha  
W. Swartout  
D. Traum

Shri Narayanan won a 2009 Best Paper Award for  
“Toward Detecting Emotions in Spoken Dialogs.”

David Chiang won a 2009 Best Paper Award  
for “11,001 New Features for Statistical  
Machine Translation.”



# Natural Language Research at USC

## Computer Science Faculty

**M. Arbib**  
**D. Chiang**  
**A. Gordon**  
**J. Gratch**  
**J. Hobbs**  
**E. Hovy**  
**K. Knight**  
**A. Leuski**  
**D. Marcu**  
**S. Marsella**  
**S. Narayanan**  
**P. Pantel**  
**F. Sha**  
**W. Swartout**  
**D. Traum**

**Daniel Marcu** co-founded an language translation software company that employs 110 people.



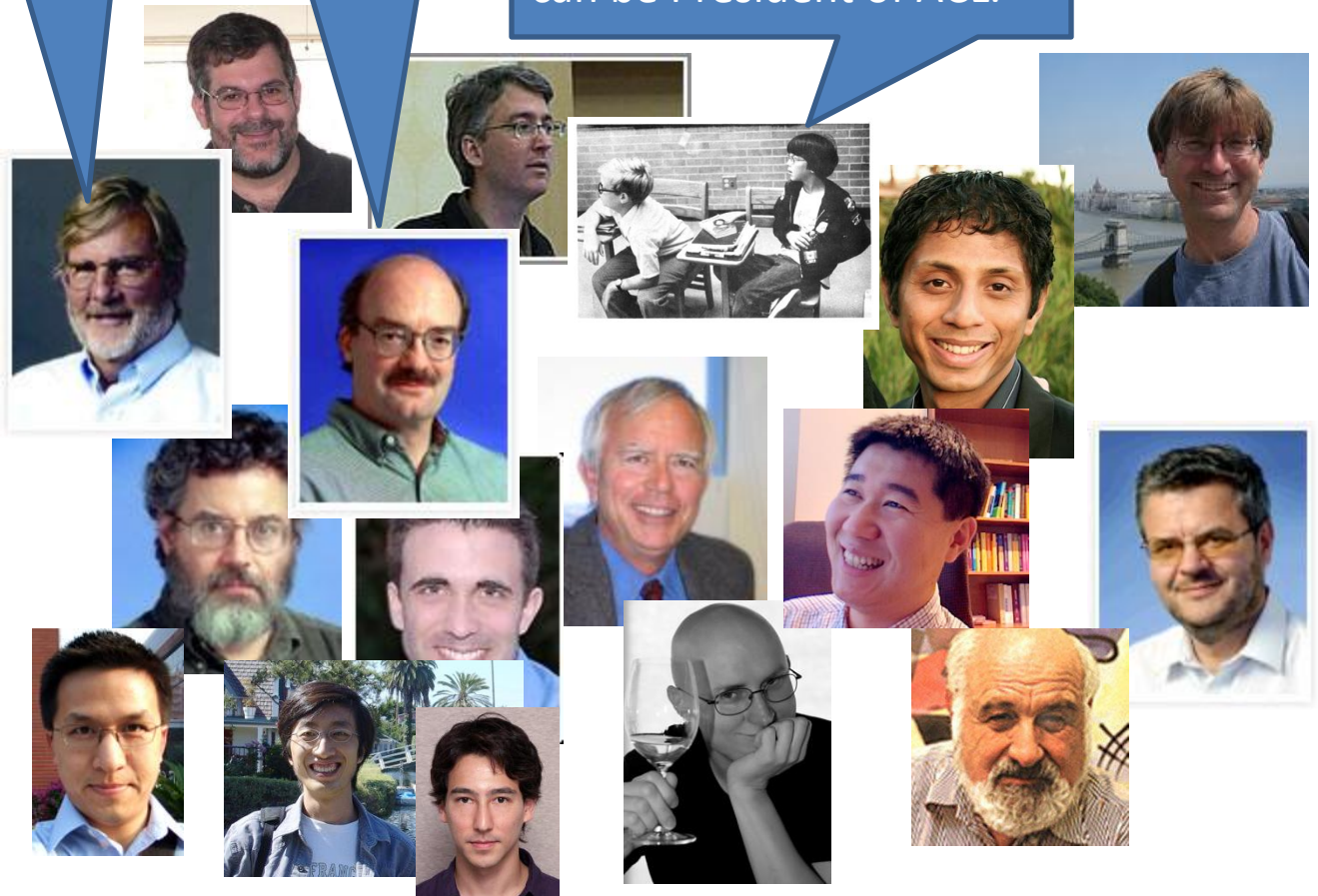
**Jerry Hobbs** was President of the Association for Computational Linguistics (ACL).

**Eduard Hovy** was also President of the ACL.

Even a **boy** from Louisiana can be President of ACL.

Computer Science Faculty

**M. Arbib**  
**D. Chiang**  
**A. Gordon**  
**J. Gratch**  
**J. Hobbs**  
**E. Hovy**  
**K. Knight**  
**A. Leuski**  
**D. Marcu**  
**S. Marsella**  
**S. Narayanan**  
**P. Pantel**  
**F. Sha**  
**W. Swartout**  
**D. Traum**

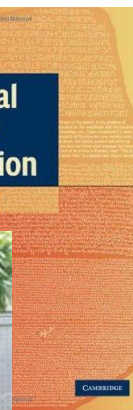




## PhD Alumni

Prof. Philipp Koehn  
University of Edinburgh

Statistical  
Machine  
Translation  
Philipp Koehn



## Computer Science Faculty

M. Arbib  
D. Chiang  
A. Gordon  
J. Gratch  
J. Hobbs  
E. Hovy  
K. Knight  
A. Leuski  
D. Marcu  
S. Marsella  
S. Narayanan  
P. Pantel  
F. Sha  
W. Swartout  
D. Traum



# Institute for Creative Technologies (ICT)

- Virtual humans
- Spoken dialogue
  - Understanding
  - Generation
  - Dialogue tracking
  - Emotion



# Electrical Engineering (EE)

Prof. Shri Narayanan lab

- Virtual Sick Call
- Enriched **Speech Translation** Systems
- User-Centric Mixed-Initiative Spoken **Dialog Systems**
- Modeling **Emotive Improvisation** in Theater Performance
- Virtual Human Museum Guides
- Exploring Emotional Vocal Productions with **MRI**
- Human-Robot Interaction and **Socially Assistive Robotics**
- Social Communication Training in Children with **Autism**
- Early Assessment of Academic Standards
- Human-Like Speech Processing
- Dynamics of **Vocal Tract Shaping**
- Prosody and **Articulatory Dynamics** in Spoken Language

DIALOGUE

EMOTION

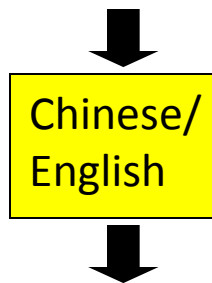
EDUCATION

DYNAMICS

# Information Sciences Institute (ISI)

- Automatic Language 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.



# Why People Get into Automatic Language Translation

- Passion about understanding how human language works
  - What makes one sequence of words grammatical, and another not?
- Interest in languages
  - What's the difference between English and Chinese?
- Desire to change the world
  - How will the world be different when the language barrier disappears?

# Why It's Hard

- Each word has tons of meanings
  - I'll **get** a cup of coffee → ?
  - I didn't **get** that joke → ?
  - I **get** up at 8am → ?
  - I **get** nervous → ?
  - Yeah, I **get** around ... → ?
- Each word has zillions of contexts
- Word order is different

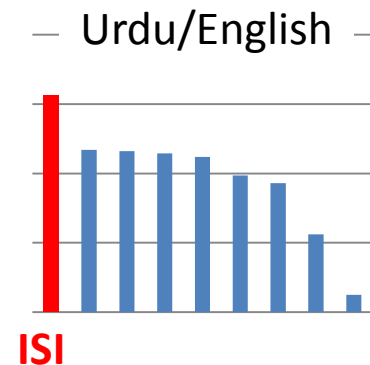
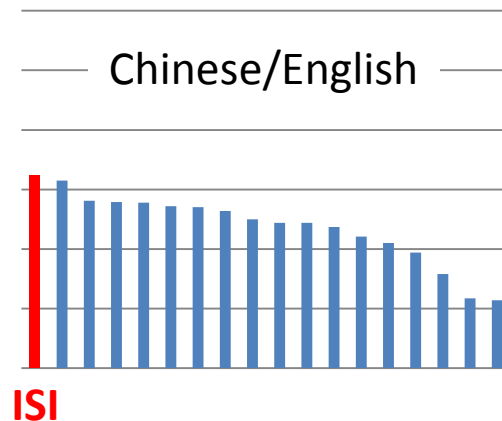
# Why It's Hard

- Output must be a grammatical, sensible, never-before-uttered sentence!
- Computers **consume** lots of human language
  - Google, Yahoo, Bing ...
  - Speech recognizers ...
- More challenging to also **produce** human language
  - What makes one sequence of words grammatical, and another not?

# Information Sciences Institute (ISI)

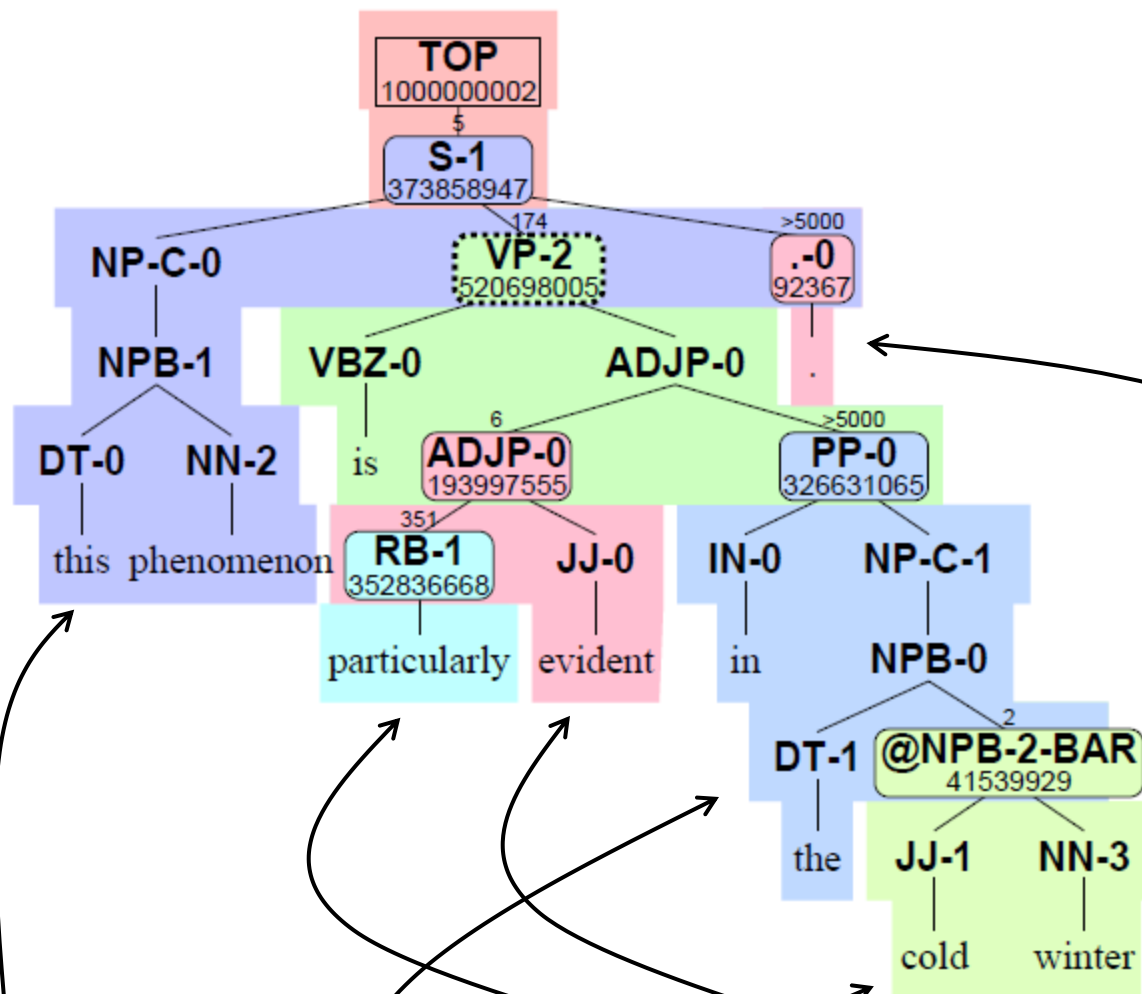
- Automatic Language Translation
  - Grand challenge machine learning problem
  - Billions of words of human translations to learn from
  - Marry language structure and text statistics

Annual  
international  
bake-off  
(2009)

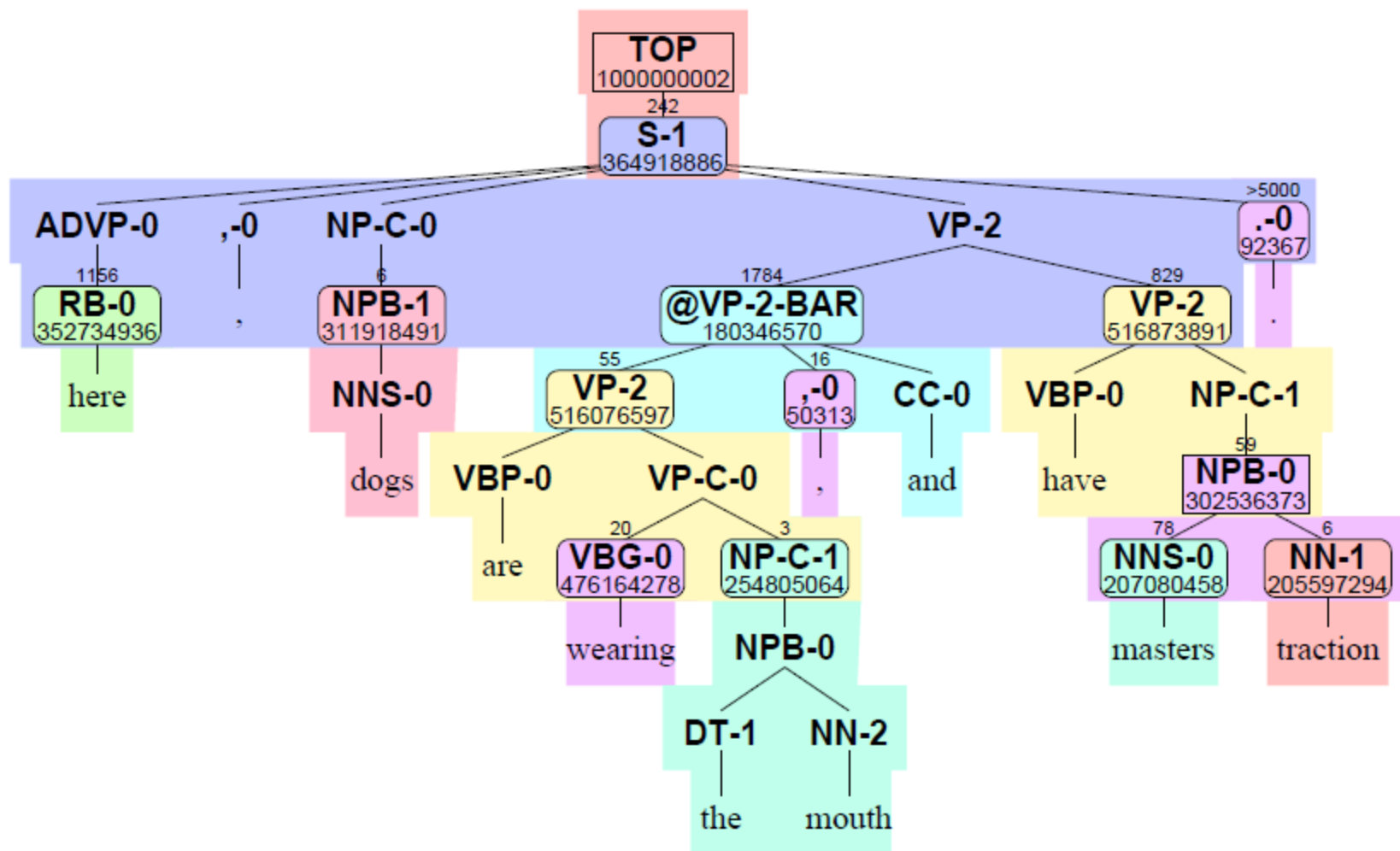




这种现象在寒冷的冬季尤其明显。



这种现象在寒冷的冬季尤其明显。



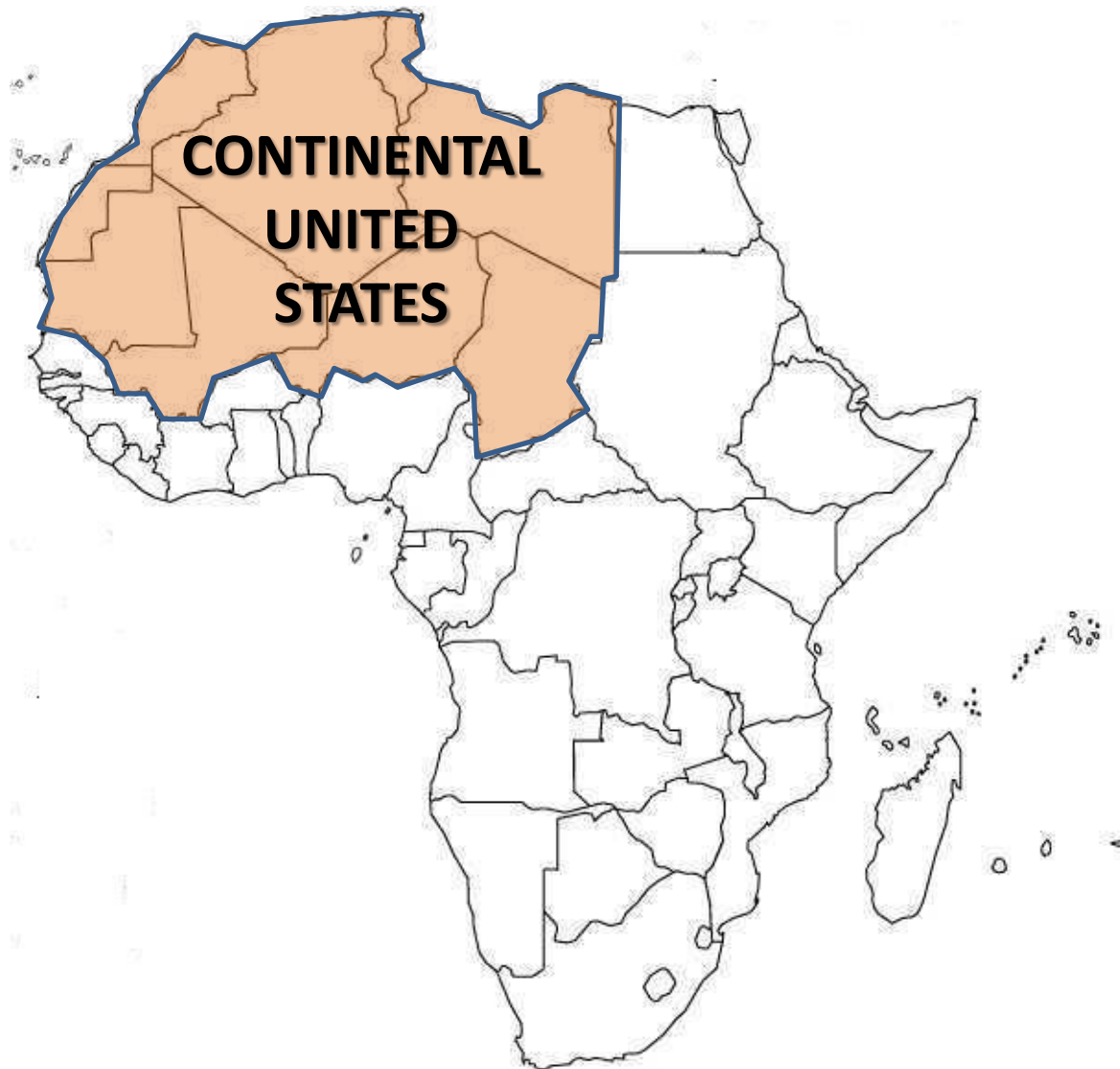
在这里，狗都配戴嘴套，并有主人牵引。

# African Languages

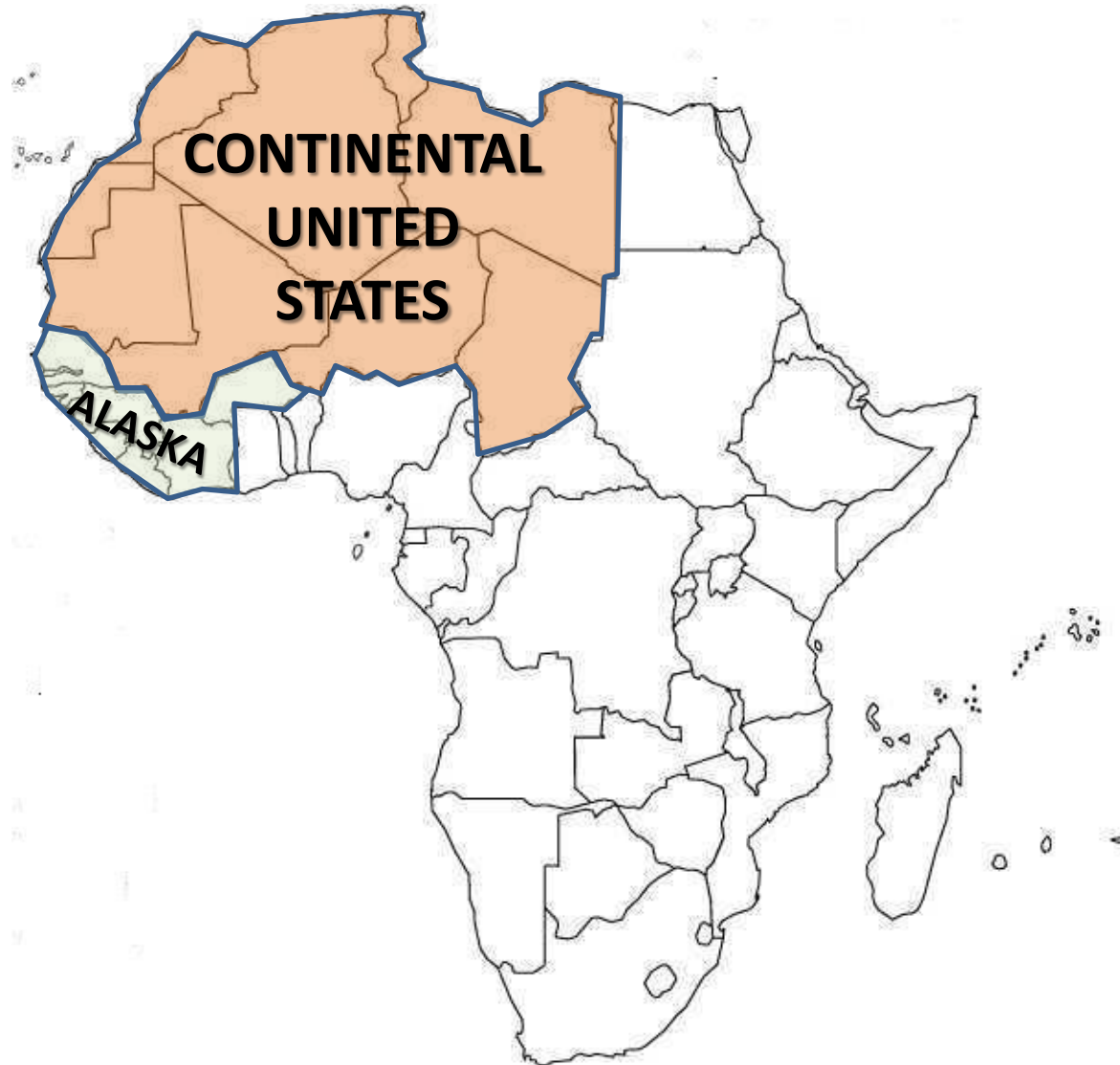




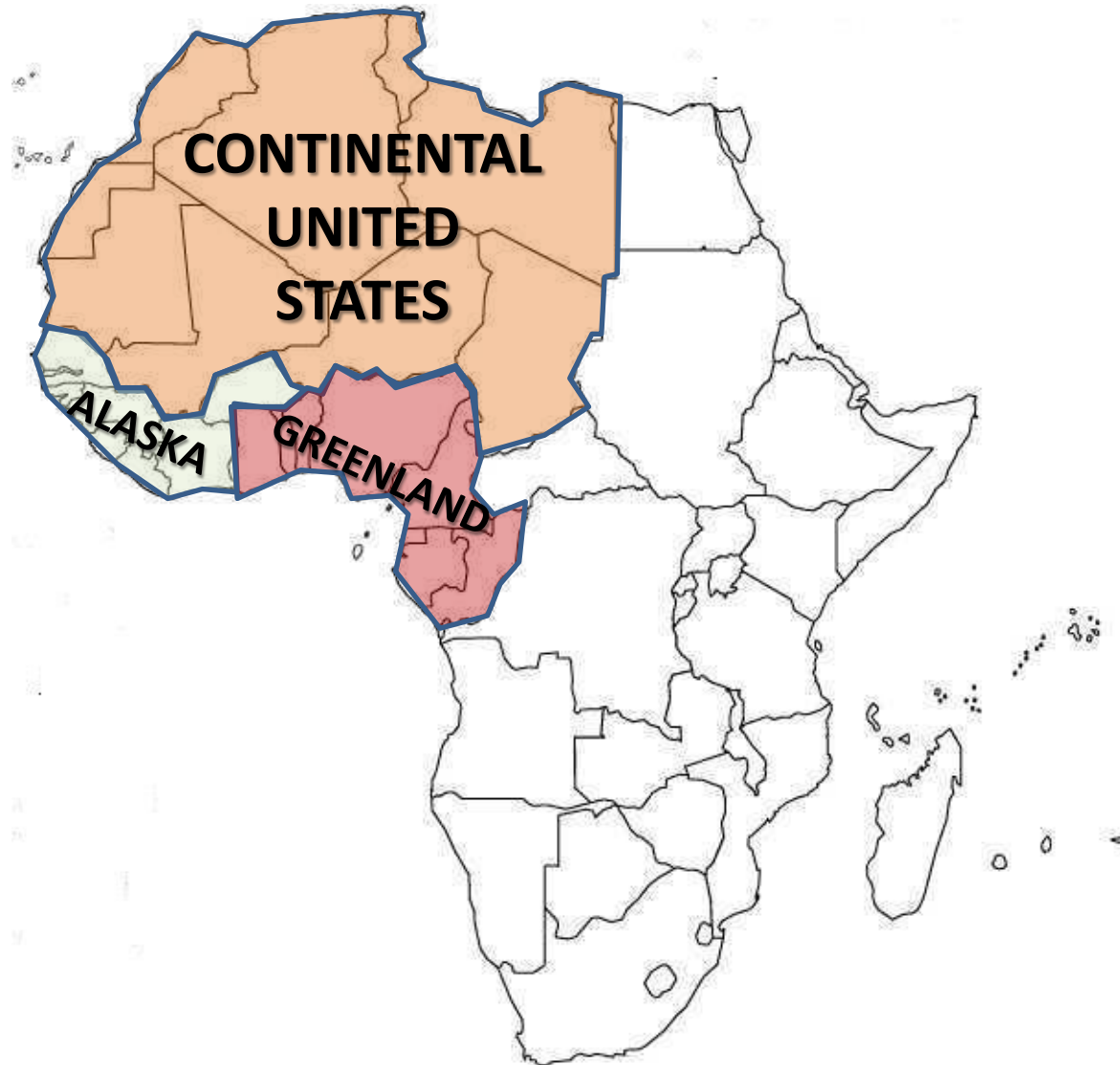
# African Languages



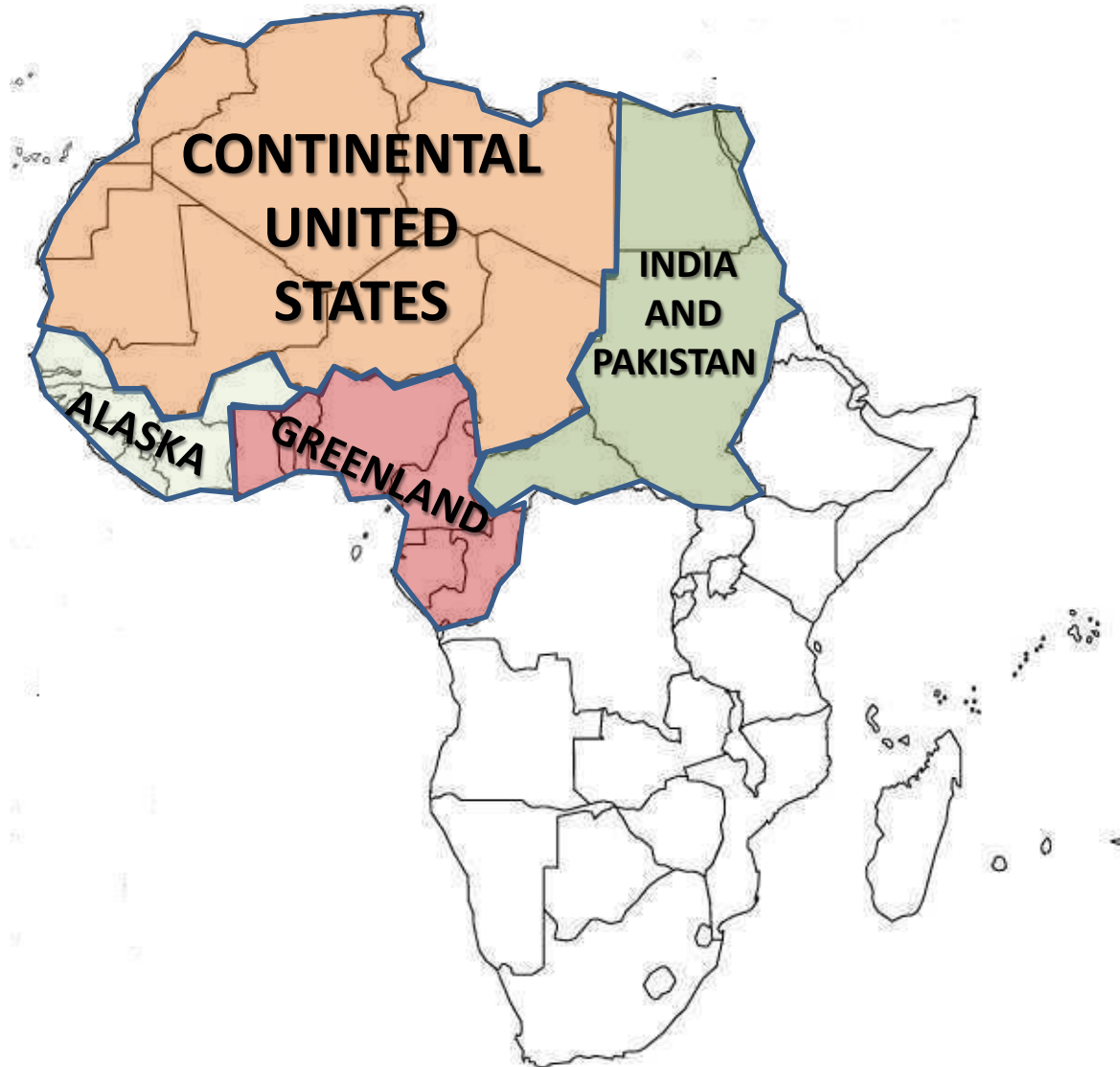
# African Languages



# African Languages

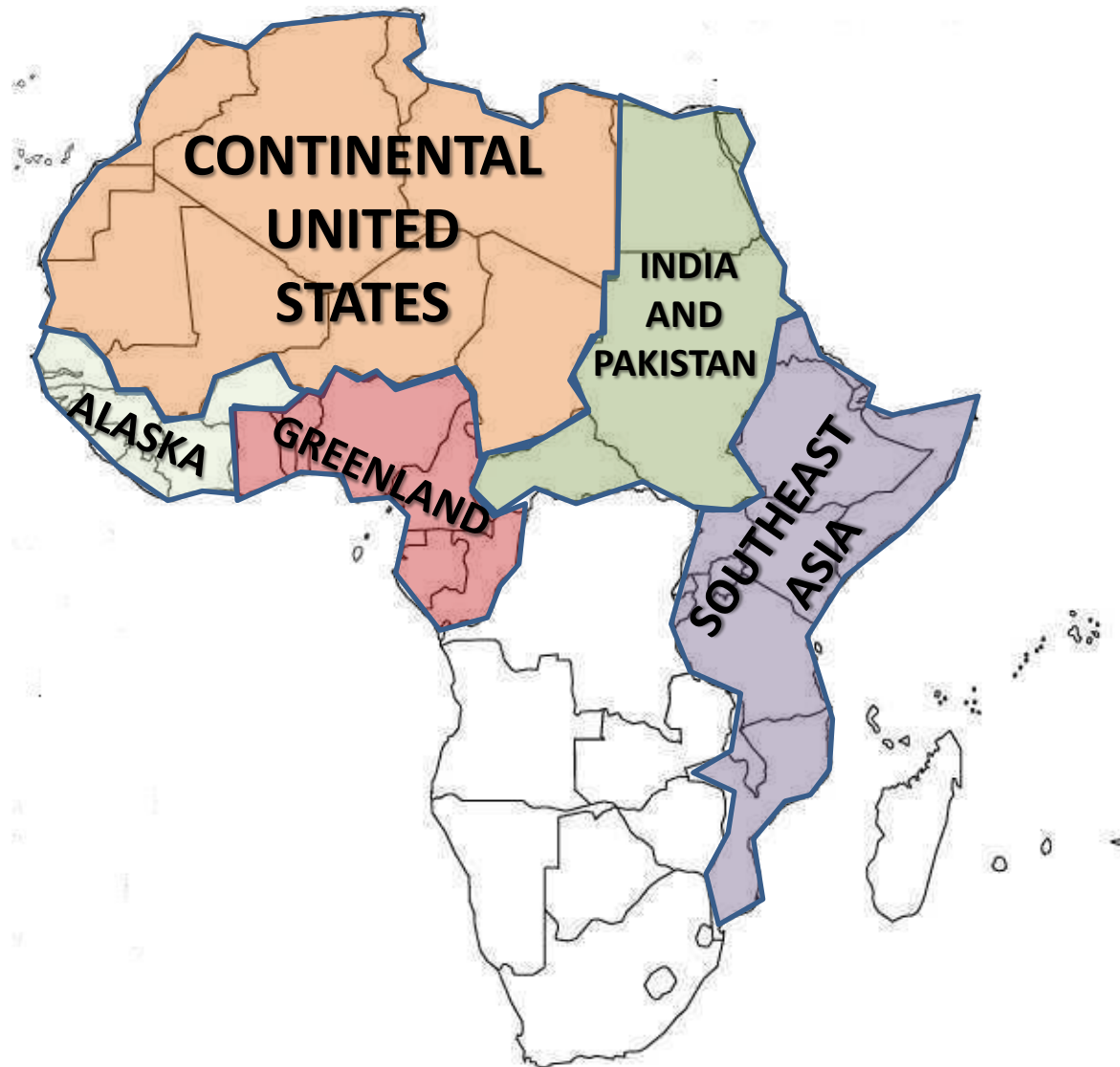


# African Languages

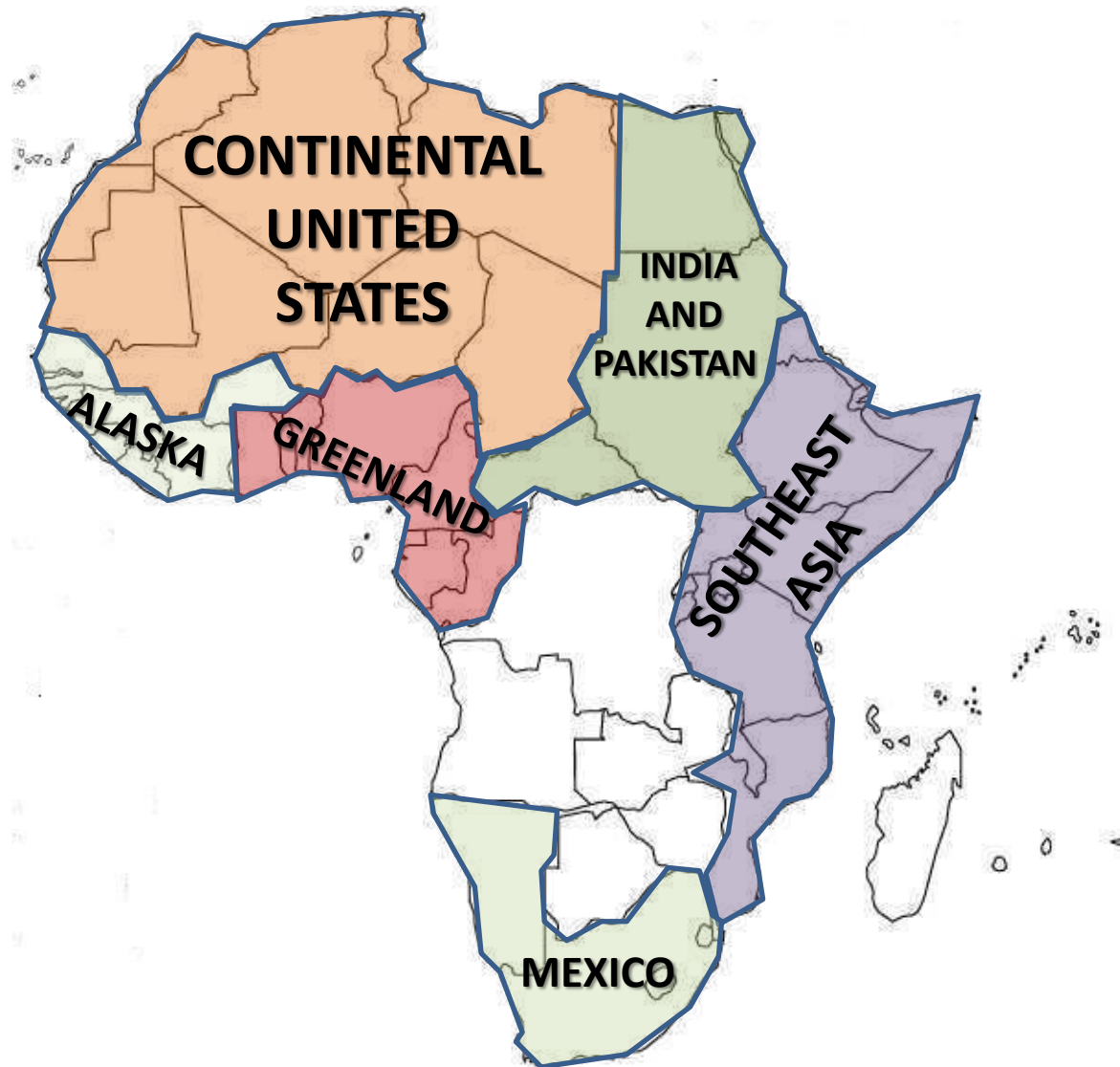




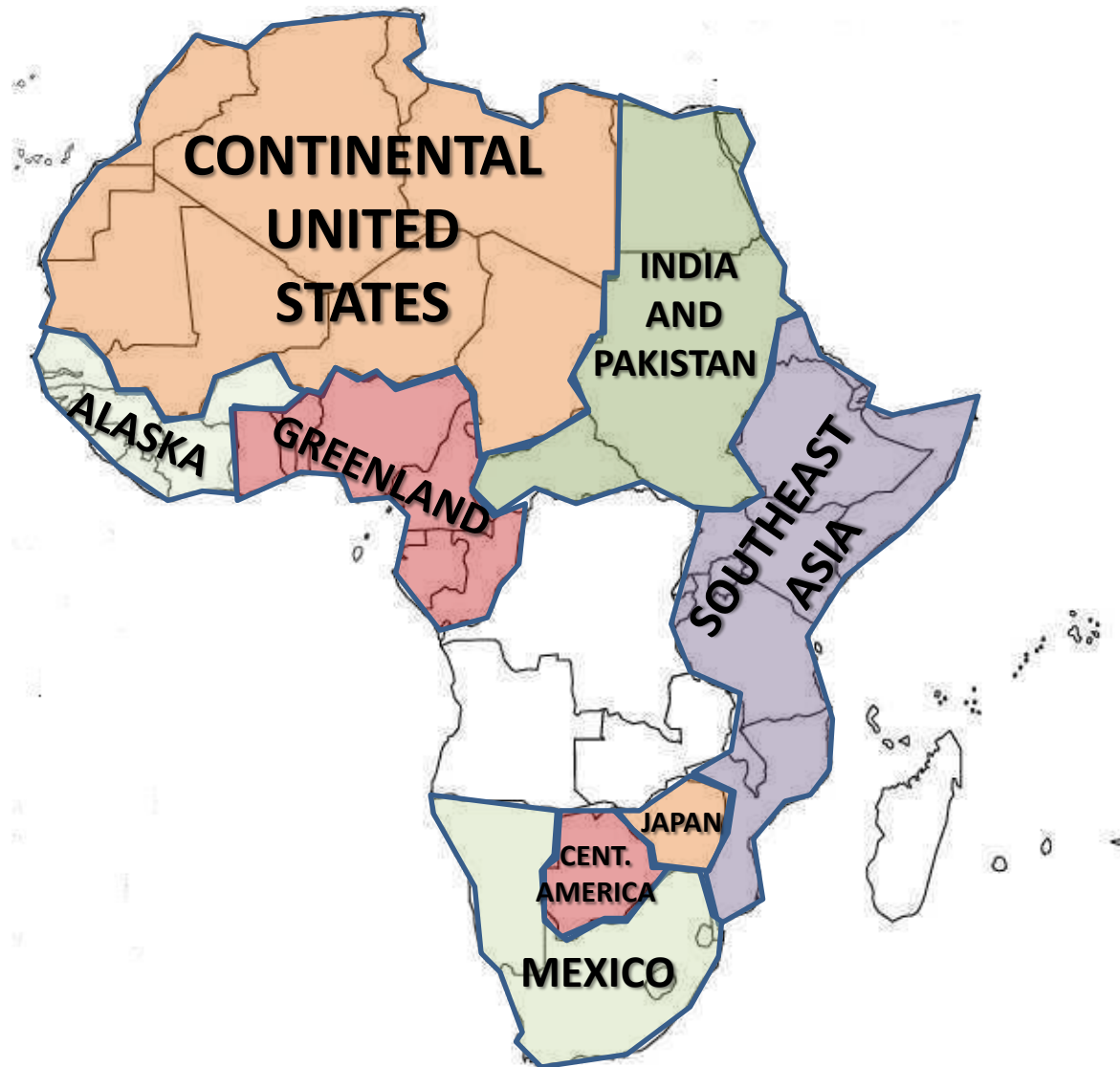
# African Languages



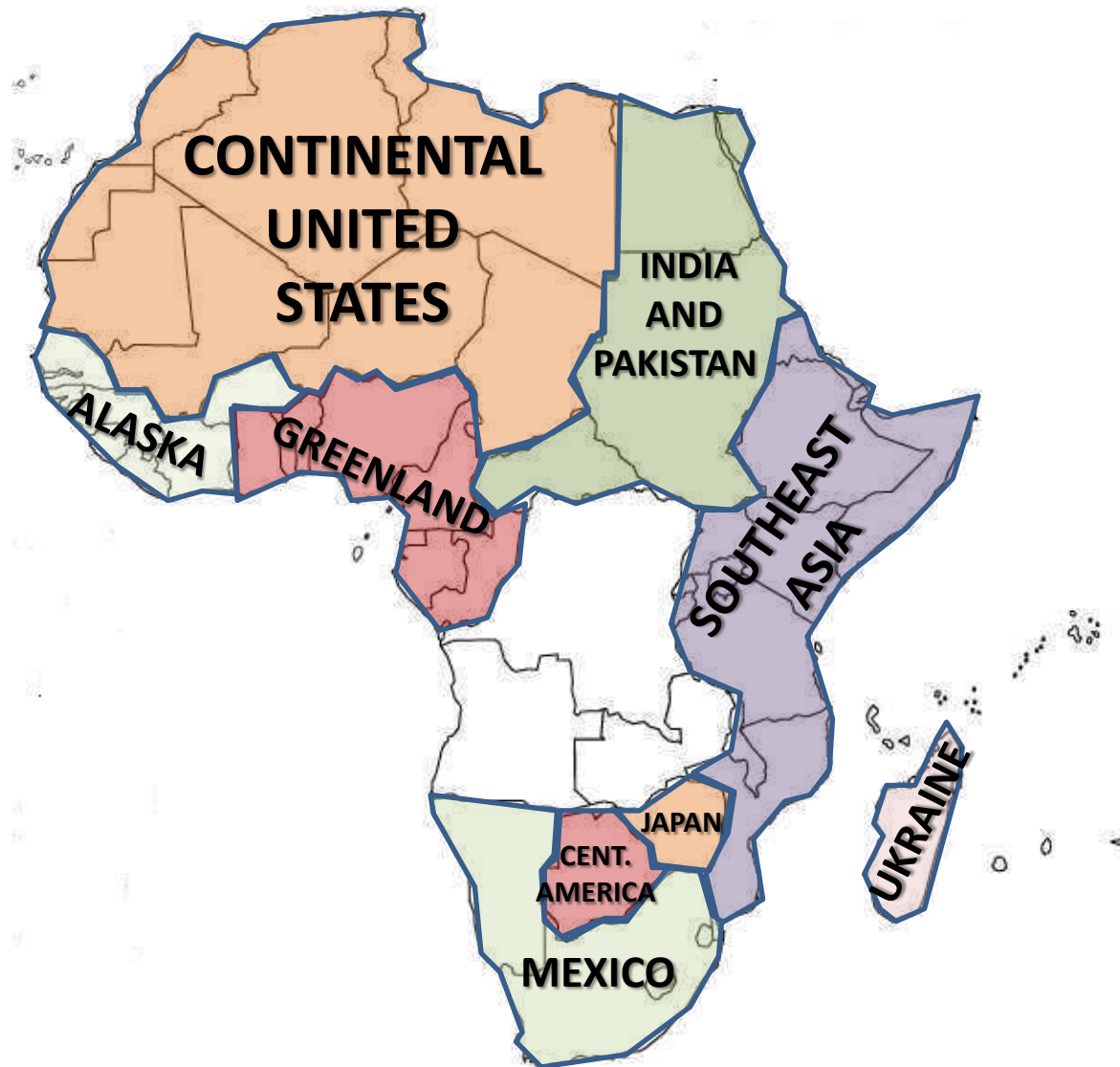
# African Languages



# African Languages



# African Languages



# African Languages



# Africa



**almost  
as big as  
the Moon**



# Africa



Zero languages spoken

# Africa



Zero languages  
spoken

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

# Information Sciences Institute (ISI)

- Learning by Reading
  - Can a machine learn about a topic by reading articles?
- Biomedical Text Extraction
  - How to turn a scientific article into a formal representation of the experiment?
- Phylogeny of Languages
- Automatic Decipherment

# Decipherment

Δ ρ / Z / U B □ X O R X 9 X X B  
W V + 3 6 Y F O Δ H P □ K E ρ Y 3  
M J Y Λ U I X Δ ρ T L N G Y ρ □  
S □ / Δ □ B P O R A U □ 7 R J ρ E  
X A L M Z T G R \ 9 F H V W 3 Δ Y  
□ + □ □ Δ K I □ □ □ X A □ □ S □  
R N I I Y E J O Δ ρ G B T G S □ B  
L O / P □ B □ X ρ E H M U Λ R R X ...

Zodiac 408-letter cipher  
(solved in 1969)

H E R > 9 J Λ V P X I □ L T G □ □  
N 9 + B □ □ O □ D W Y · < □ K 7 □  
B X □ □ M + u z G W □ □ L □ □ H J  
S 9 9 Δ Λ J Δ □ V O 9 O + + R K □  
□ Δ M + □ L T □ I □ F P + P □ X /  
9 Δ R Λ F J O - □ □ C □ F > □ D □  
□ □ + K □ □ □ □ □ X G V · □ L I  
□ G □ J 7 T □ O + □ N Y □ + □ L Δ  
□ < M + 8 + Z R □ F B □ Y A □ □ K  
- □ J u v + Λ J + O 9 Δ < F B Y -  
U + R / □ L E I D Y B 9 8 T M K □  
□ < □ J R J I □ □ T □ M · + P B F  
□ □ Δ S Y □ + N I □ F B □ □ □ R  
J G F N Λ 7 □ □ □ □ □ □ V □ L + +  
Y B X □ □ □ □ □ □ □ > V u z - +  
I □ · □ □ B K □ □ □ 9 Λ · 7 M □ □  
R □ T + L □ □ □ < + F J W B I □ L  
+ + □ W C □ □ P O S H T / □ □ 9  
I F X □ W < Δ L □ □ Y O B □ - □ □  
> M D H N 9 X S □ Z O Δ A I K □ +

Zodiac 340-letter cipher  
(undeciphered)



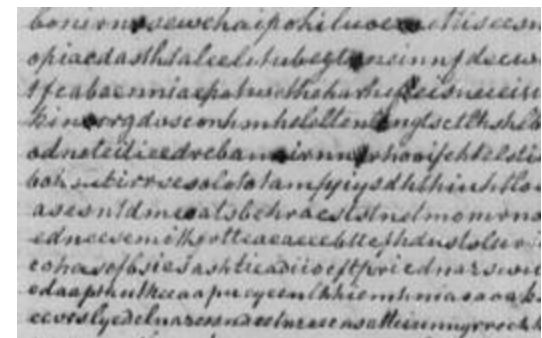
Ugaritic tablets  
(ISI with MIT)



Foreign language as a code  
for English (Weaver 1949)



Voynich Manuscript, 240pp,  
(undeciphered)



Thomas Jefferson cipher  
(solved, see WSJ 7/6/09)

# Zodiac Serial Killer

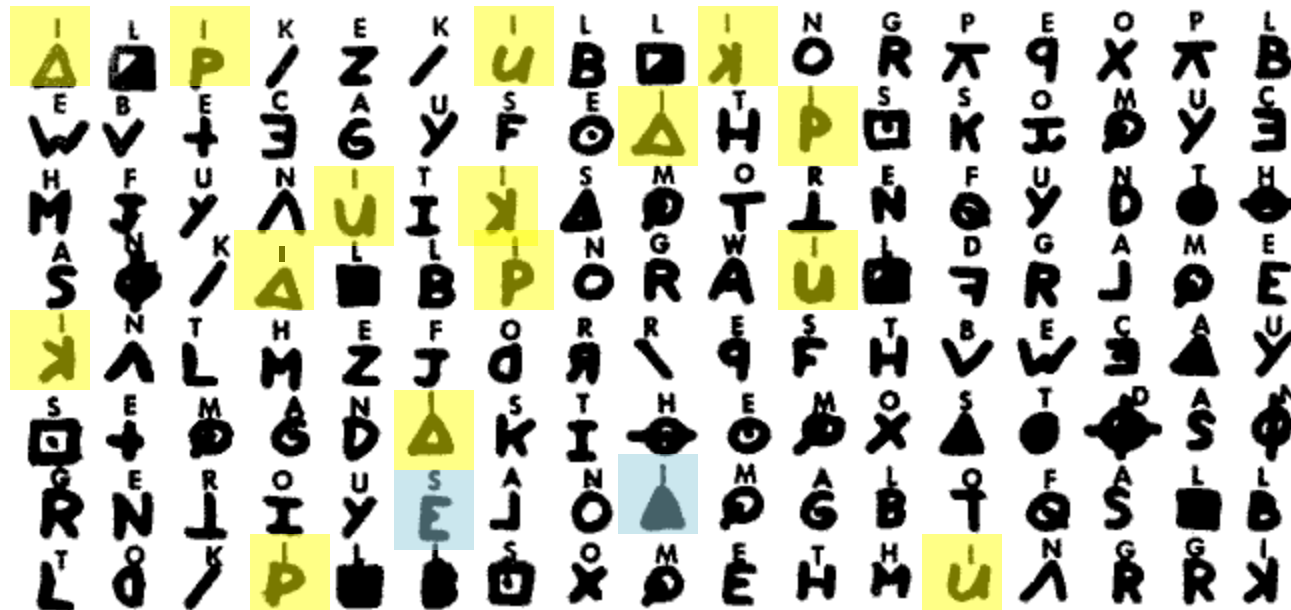
408-letter cipher, solved by husband and wife over breakfast:

Δ ▣ P / Z / U B ▣ X O R π 9 X π B  
W V + 3 6 Y F ⊙ Δ H P ▣ K π ρ Y 3  
M J Y Λ U I X Δ ρ T ⊥ N ⊙ Y D ⊙ ⊙  
S ⊙ / Δ ▣ B P O R A U ▣ 7 R J ρ E  
X Λ L M Z T ⊙ R \ 9 F H V W 3 Δ Y  
▣ + ⊙ ⊙ ⊙ Δ K I ⊙ ⊙ ρ X Δ ⊙ ⊙ S ⊙  
R N I I Y F J O Δ ρ G B T ⊙ S ▣ B  
L ⊙ / P ▣ B ▣ X ρ E H M U Λ R R X

(plus two more sections)

# Zodiac Serial Killer

408-letter cipher, solved by husband and wife over breakfast:



(plus two more sections)

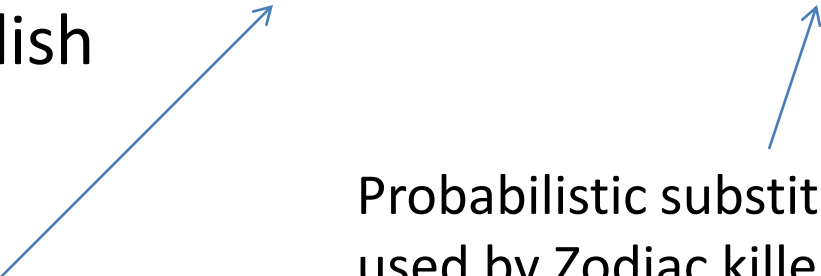
Homophonic Cipher

Plain	Cipher
A	1 3 G 1 S
B	V
C	e
D	@ f
E	+ E I N P W Z
F	J Q
G	R
H	= M
I	2 k P U
J	
K	/
L	4 7 B
M	q
N	^ \$ D O
O	; d T X
P	:
Q	
R	! \ r
S	1 3 6 F K
T	* H I L N
U	Y
V	l
W	A 2
X	t
Y	5
Z	

‘We felt that the word “KILL” or “KILLING” would appear in his code, and the word “I”, because he had an ego.’ -- Donald Harden



# Probabilistic Approach

$$P(\text{cipher}) = \sum_{\text{english}} P(\text{english}) * P(\text{cipher} \mid \text{english})$$


Program assigns a score to any sequence of English letters.

High score = good English!

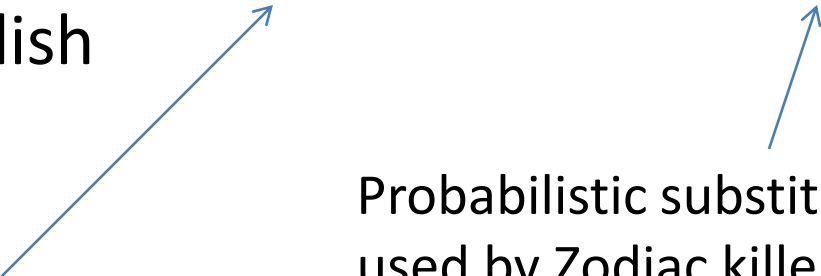
**We have to supply this.**

Probabilistic substitution table used by Zodiac killer.

If plaintext letter is “E”, probability of enciphering as “q” is 0.043.

**Hidden! Learning program sets this to maximize  $P(\text{cipher})$ .**

# Probabilistic Approach

$$P(\text{cipher}) = \sum_{\text{english}} P(\text{english}) * P(\text{cipher} \mid \text{english})$$


Program assigns a score to any sequence of English letters.

High score = good English!

**We have to supply this.**

Probabilistic substitution table used by Zodiac killer.

If plaintext letter is “E”, probability of enciphering as “q” is 0.043.

**Hidden! Learning program sets this to maximize  $P(\text{cipher})$ .**

Does maximizing  $P(\text{cipher})$  yield more accurate substitution tables?

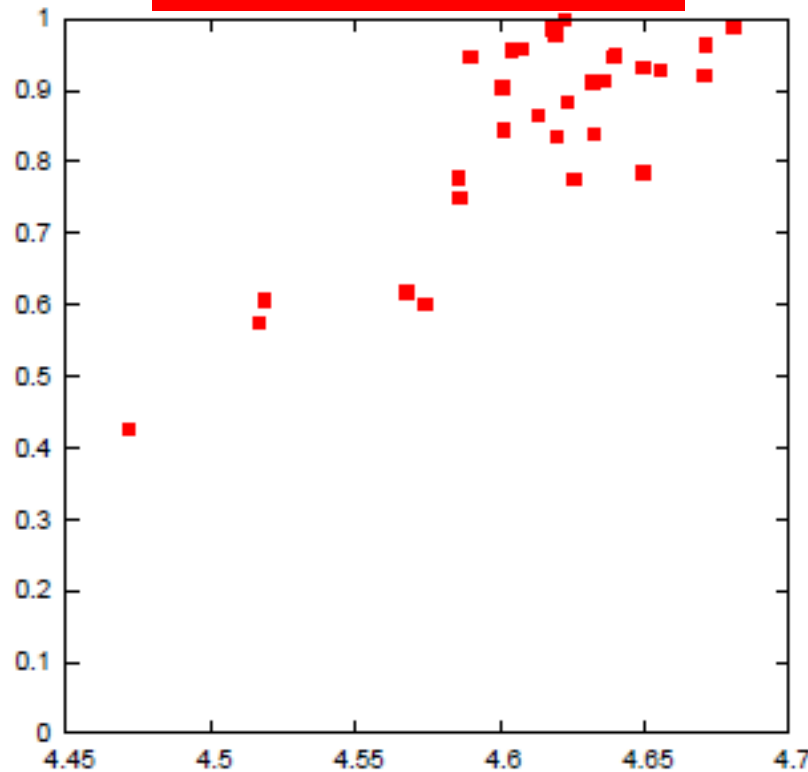
How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

How much English do you have to know to solve a 408-letter homophonic cipher?

Does maximizing  $P(\text{cipher})$  yield more accurate substitution tables?

Each square is a different substitution table

Fortuitous result:  
**Minimize  
Decipherment  
(Task)  
Error**



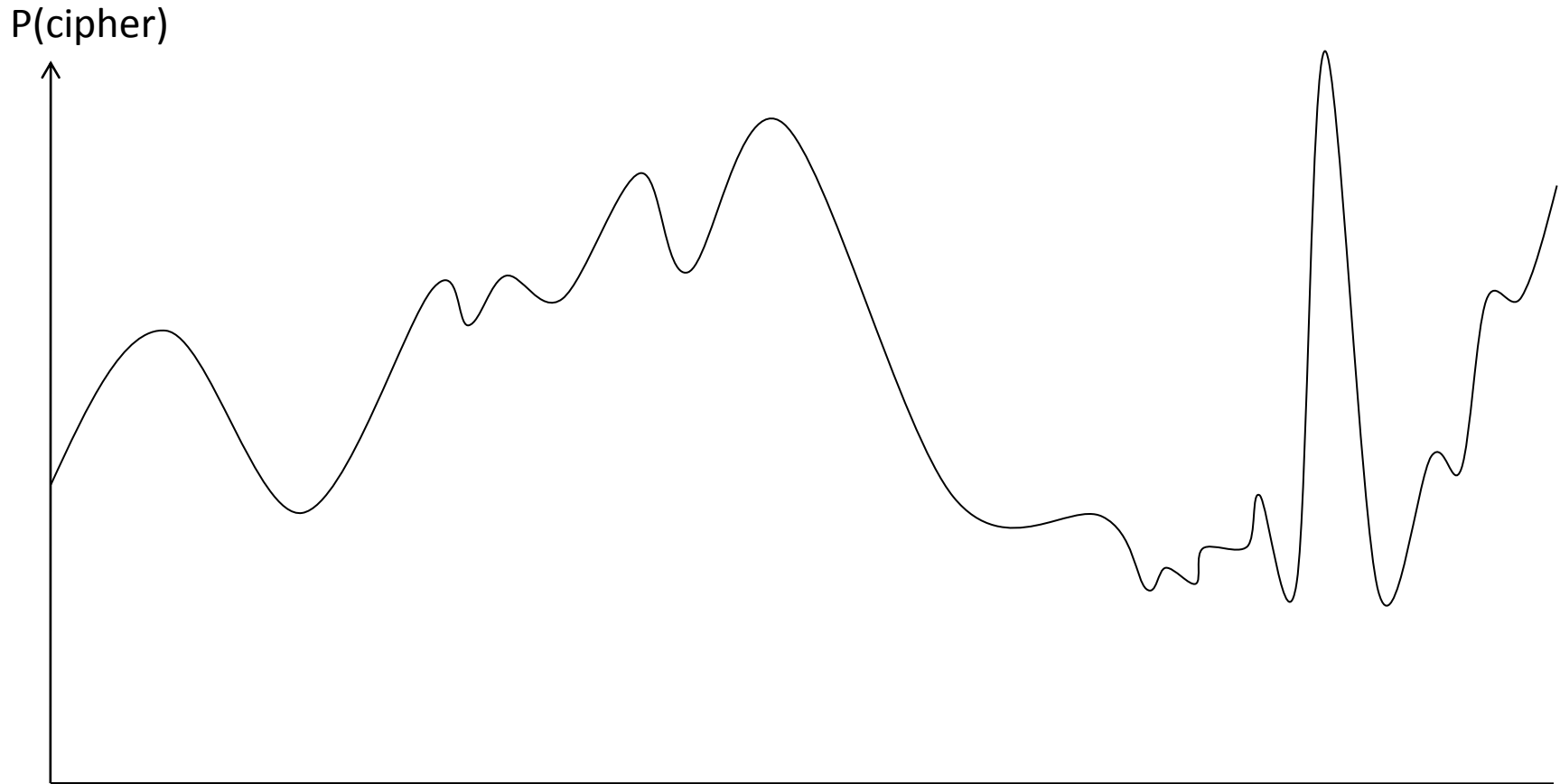
What computer does:  
**Maximize  $P(\text{cipher})$**



all  
hail  
maximum  
likelihood

How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

# Search Space

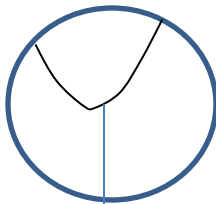


All substitution tables  
(fancifully laid out in one dimension instead of  $26 \times 64$ )

How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

# Search Space

$P(\text{cipher})$

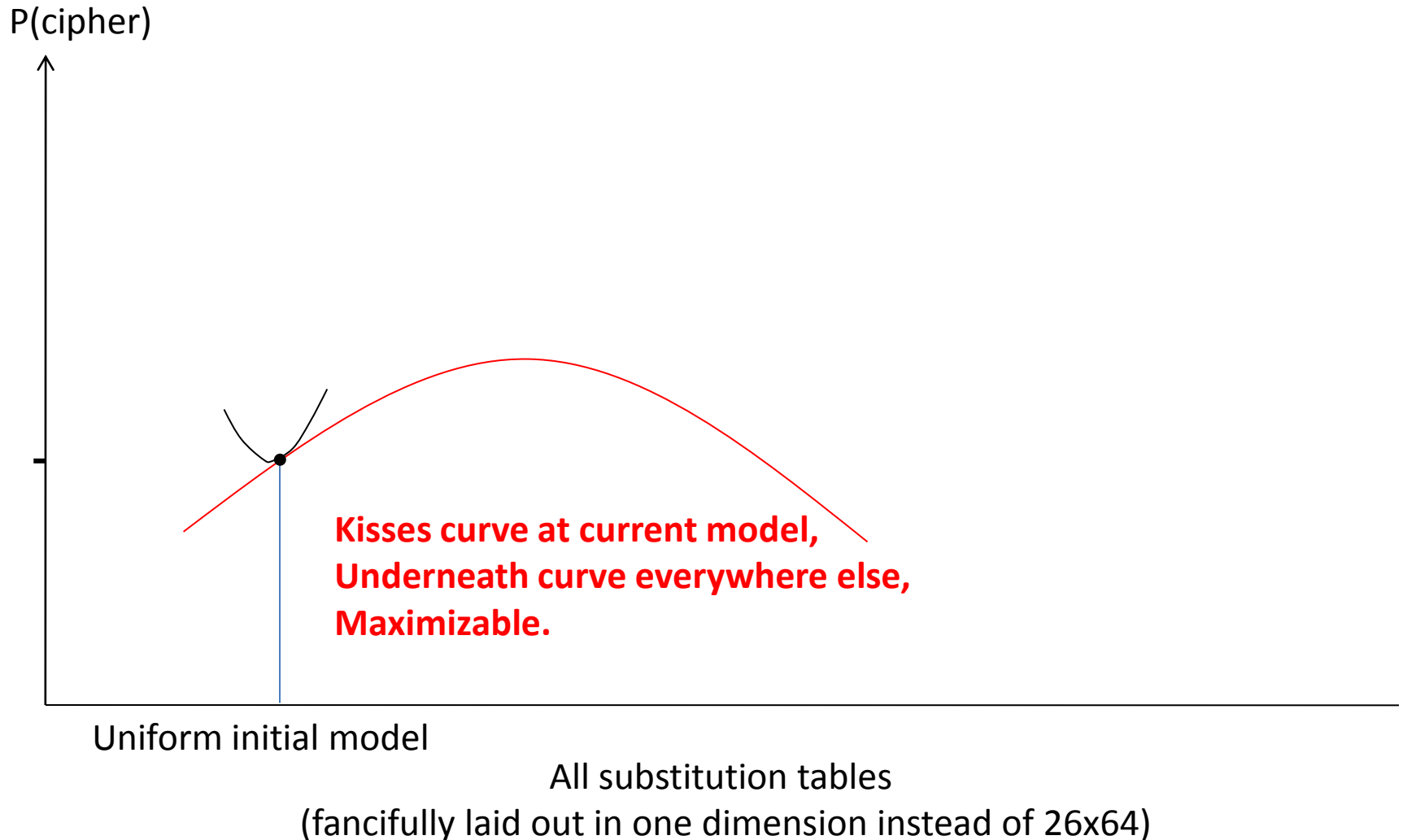


Uniform initial model

All substitution tables  
(fancifully laid out in one dimension instead of  $26 \times 64$ )

How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

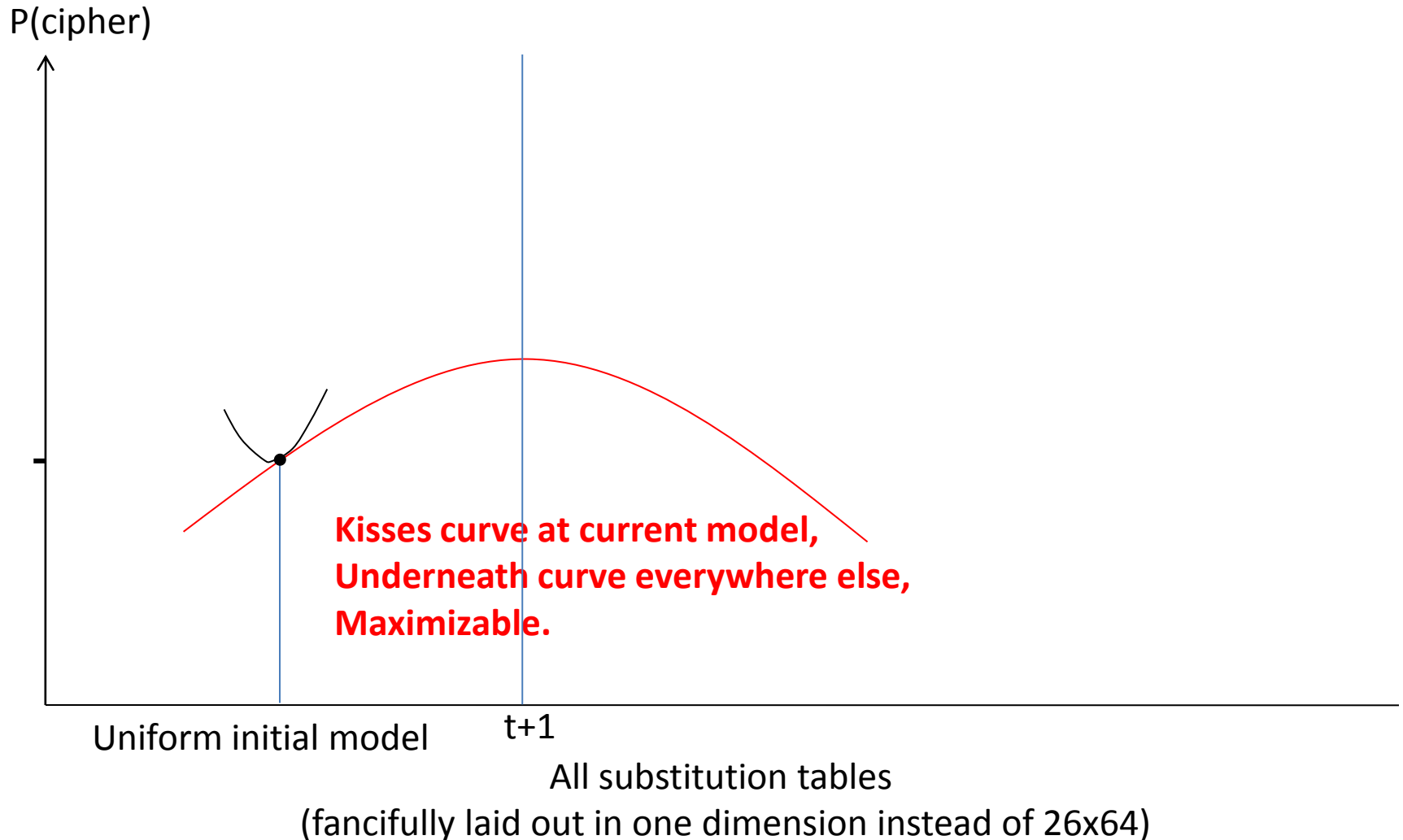
# EM Algorithm





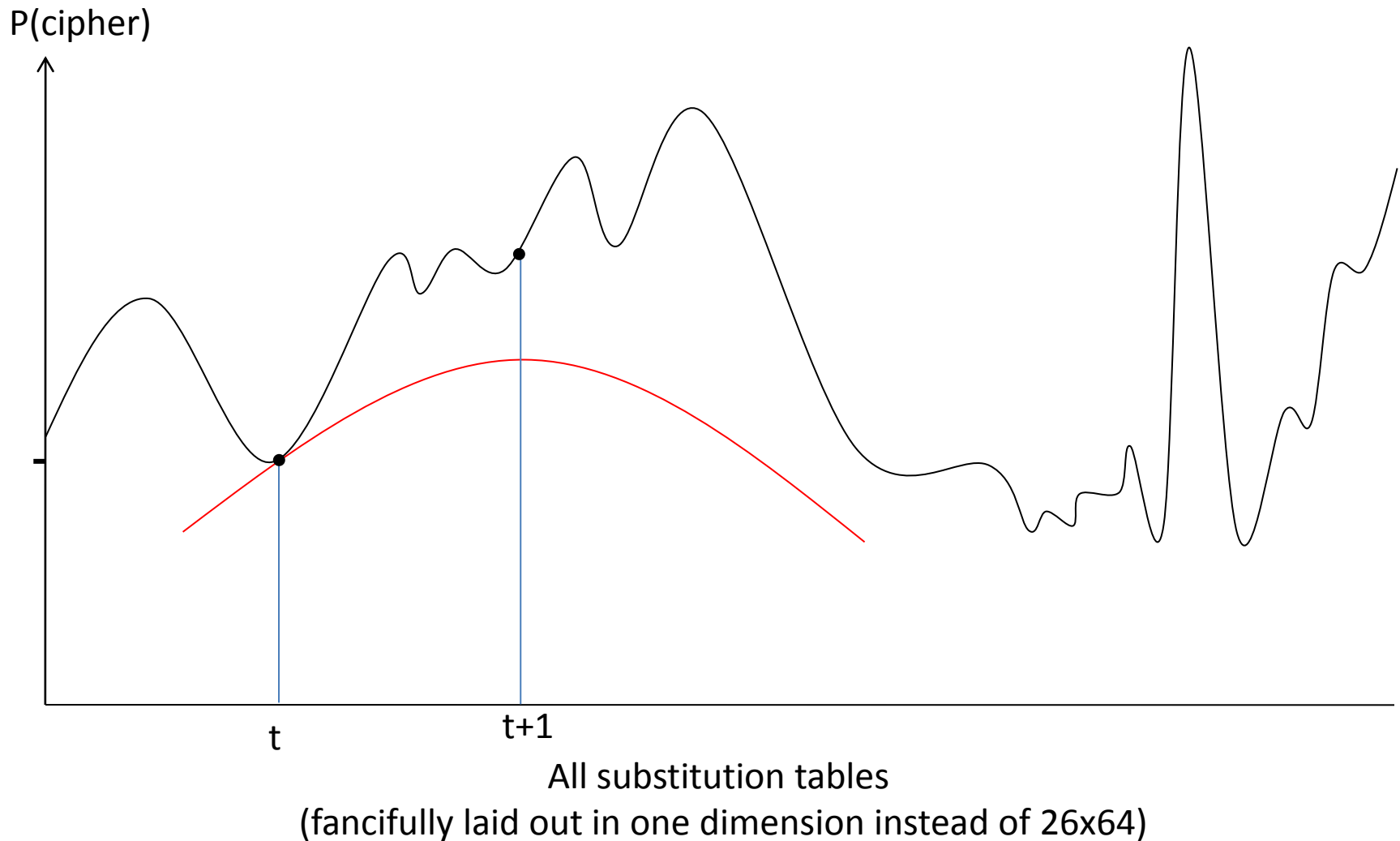
How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

# EM Algorithm



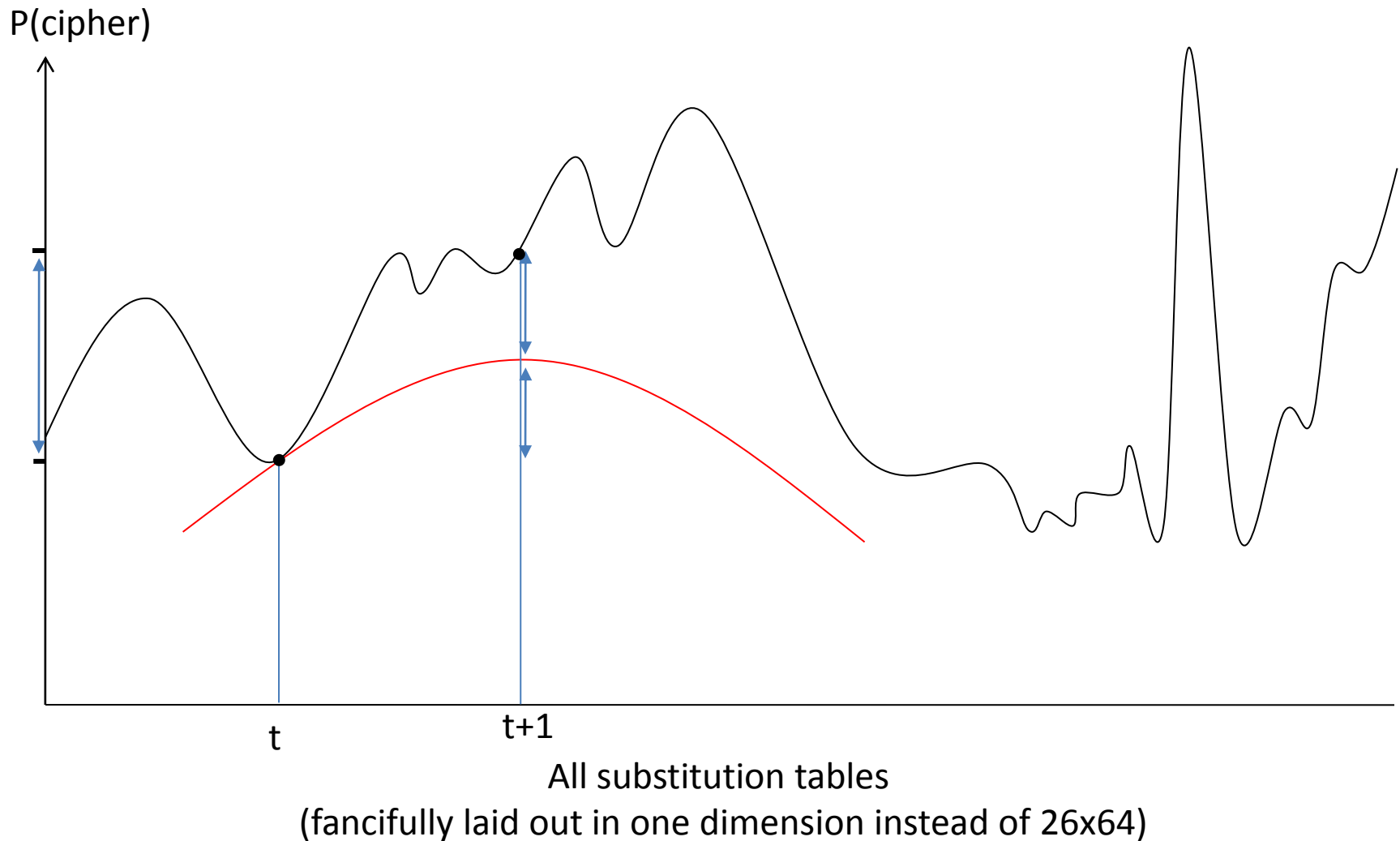
How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

# EM Algorithm



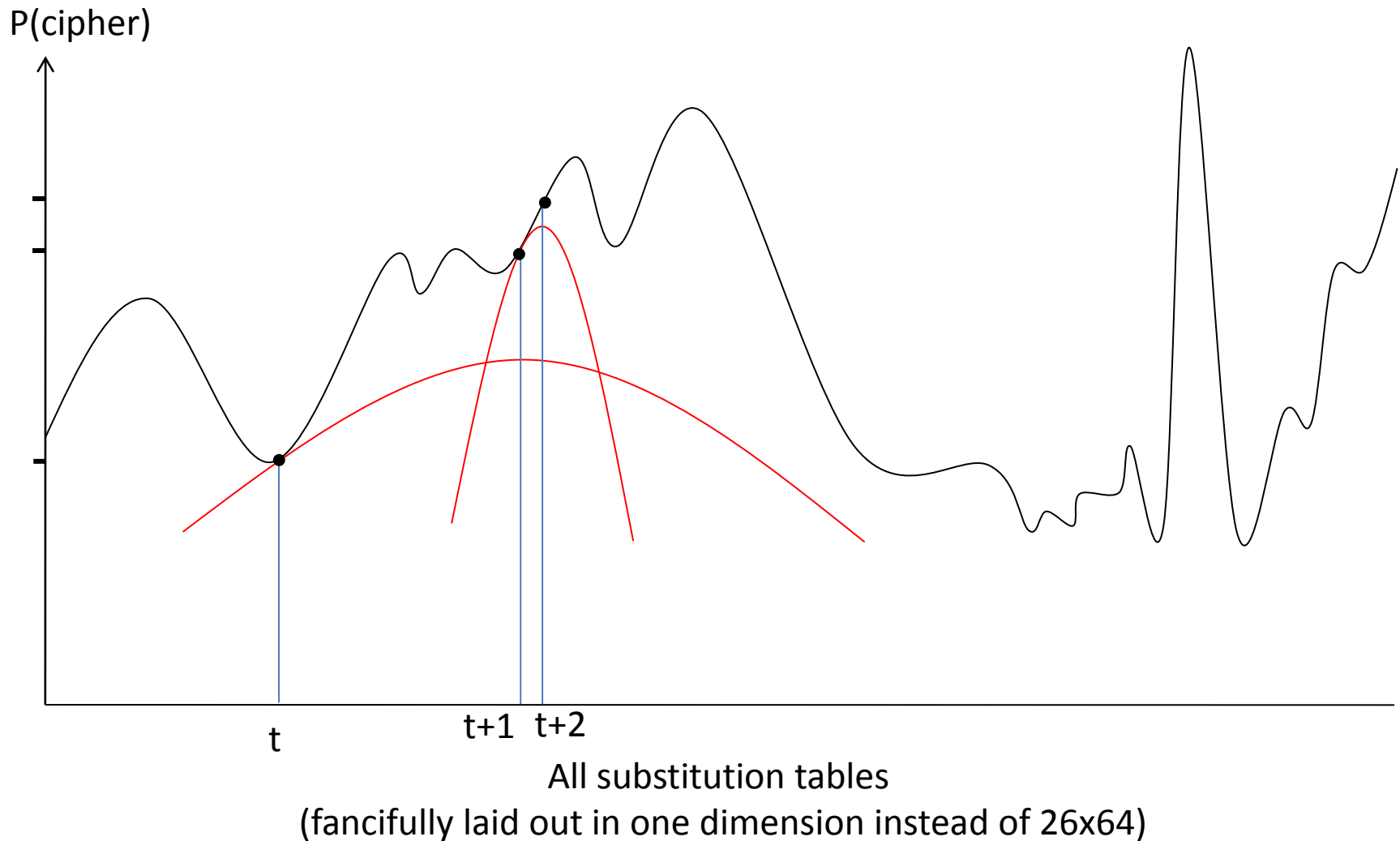
How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

# EM Algorithm



How to efficiently search for the table that maximizes  $P(\text{cipher})$ ?

# EM Algorithm



How much English do you have to know  
to solve a 408-letter homophonic cipher?



# Experiments on Zodiac 408

English knowledge	EM starts	Crib?	Decipherment error
3-gram letter statistics $P(s \mid t \ t)$	1	-	99.7 – everything wrong!
	100	-	71.2
	100	KILLING	24.6
1-gram word statistics $P(\text{pittsburgh})$	1	-	98.4
	20	-	65.6
	20	ILIKEKILLINGPEOPLE	4.1 – almost perfect

Secret weapons:

- EM
- Bayesian reasoning
- Finite-state transducer toolkits
- Modeling

# Natural Language Research at USC

Computer Science Faculty

Institute for Creative Technologies (ICT)

Information Sciences Institute (ISI)

Electrical Engineering Department (EE)

**M. Arbib**  
**D. Chiang**  
**A. Gordon**  
**J. Gratch**  
**J. Hobbs**  
**E. Hovy**  
**K. Knight**  
**A. Leuski**  
**D. Marcu**  
**S. Marsella**  
**S. Narayanan**  
**P. Pantel**  
**F. Sha**  
**W. Swartout**  
**D. Traum**





thanks