# Statistical NLP
# Spring 2010
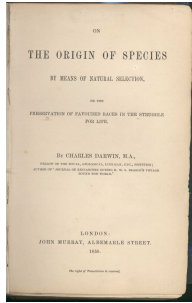
University of California
C A L
N L P
**Berkeley**

Lecture 25: Diachronics

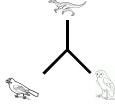Dan Klein – UC Berkeley

---

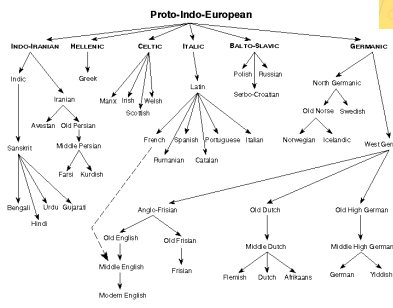## Evolution: Main Phenomena
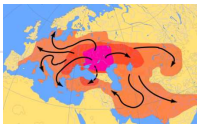
**Mutations of sequences**



Time

**Speciation**

Time

---

## Tree of Languages

**Challenge: identify the phylogeny**

- Much work in biology, e.g. work by Warnow, Felsenstein, Steele…
- Also in linguistics, e.g. Warnow et al., Gray and Atkinson…

http://andromeda.rutgers.edu/~jlynch/language.html



---

## Statistical Inference Tasks

**Inputs**

Modern Text

FR  IT  PT  ES

Phylogeny

**Outputs**

focus
fuego   feu
Ancestral Word Forms

fuego   oeuf
huevo   feu
Cognate Groups / Translations

les faits sont très clairs
Grammatical Inference

---

## Outline

focus
fuego   feu
Ancestral Word Forms

fuego   oeuf
huevo   feu
Cognate Groups / Translations

les faits sont très clairs
Grammatical Inference

---

## Language Evolution: Sound Change

Latin     camera /kamera/

Deletion: /e/

Change: /k/ .. /tʃ/ .. /ʃ/

Insertion: /b/

French    chambre /ʃambʁ/

Eng. camera from Latin, "camera obscura"

Eng. chamber from Old Fr. before the initial /t/ dropped

---

## Diachronic Evidence

### Yahoo! Answers [2009]

**Resolved Question**    Show me another »

**Which is correct....tonight or tonite?**
11 months ago

[!] Report Abuse

**Best Answer** - Chosen by Voters

"Tonight" is the traditional version.

If you'll observe, "tonite" is listed as a misspelling by the system here.
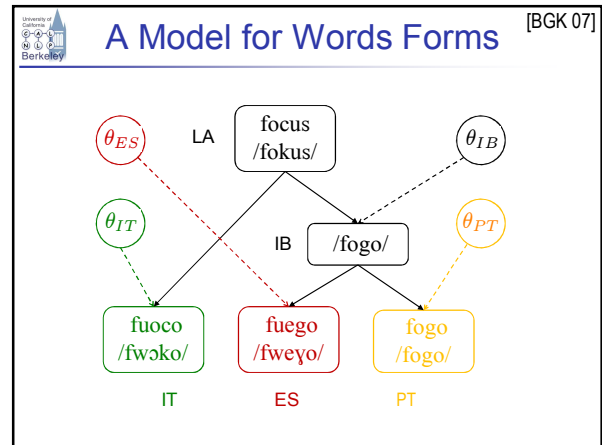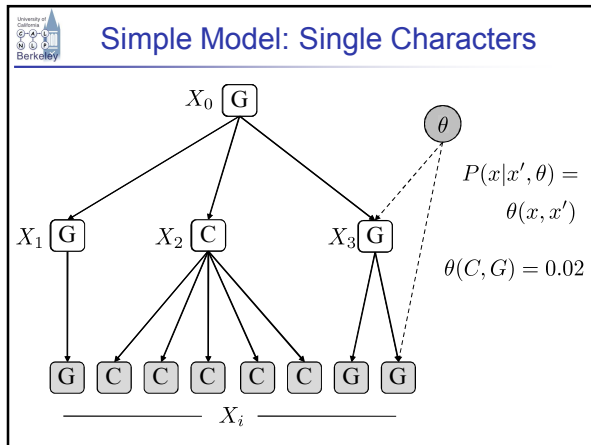
The use of "tonite" can probably be traced to the way that people make mistakes and they stick with a small group and then the use of it expands, making it become a use that people accept
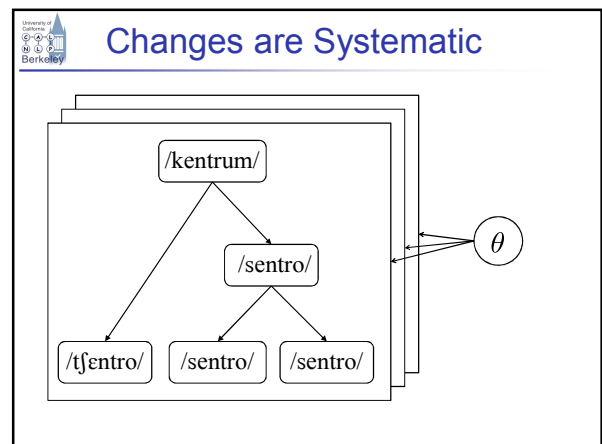
11 months ago

tonight not tonite

### Appendix Probi [ca 300]

tonitru non tonotru

---

## Synchronic (Comparative) Evidence

| Gloss | Latin | Italian | Spanish | Portuguese |
|---|---|---|---|---|
| Word/verb | verbum | verbo | verbo | verbu |
| Fruit | fructus | frutta | fruta | fruta |
| Laugh | ridere | ridere | reir | rir |
| Center | centrum | centro | centro | centro |
| August | augustus | agosto | agosto | agosto |
| Swim | natare | nuotare | nadar | nadar |

---

## Simple Model: Single Characters

$X_0$ [G]       [$\theta$]

$$P(x|x',\theta) = \theta(x,x')$$

$$\theta(C,G) = 0.02$$

$X_1$ [G]   $X_2$ [C]   $X_3$ [G]

[G] [C] [C] [C] [C] [C] [G] [G]

$$\underline{\qquad\qquad X_i \qquad\qquad}$$

---

## A Model for Words Forms    [BGK 07]

$\theta_{ES}$   LA   focus /fokus/   $\theta_{IB}$

$\theta_{IT}$   IB   /fogo/   $\theta_{PT}$

fuoco /fwɔko/   fuego /fweɣo/   fogo /fogo/

IT   ES   PT

---

## Contextual Changes

/fokus/

[#] [f] [o]

[#] [f] [w] [ɔ]   . . .

$\theta_{IT}$

/fwɔko/

$$P(w,a|w',\theta_\ell) =$$

$$\prod_k P(w_k, a_k | w_{k-1}, w', \theta_\ell) \propto$$

$$\exp\left(\theta_\ell^\top f(w_k, w_{k-1}, w'_{a_{k-1}}, w'_{a_k}, w'_{a_{k+1}})\right)$$

---

## Changes are Systematic

/kentrum/

/sentro/

$\theta$

/tʃentro/   /sentro/   /sentro/

## Experimental Setup
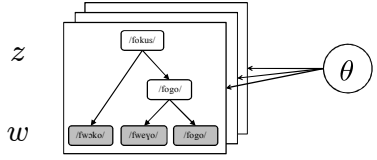
University of California
C A L
N L P
Berkeley

- Data sets
  - Small: Romance
    - French, Italian, Portuguese, Spanish
    - 2344 words
    - Complete cognate sets
    - Target: (Vulgar) Latin

FR   IT   PT   ES

  - Large: Oceanic
    - 661 languages
    - 140K words
    - Incomplete cognate sets
    - Target: Proto-Oceanic [Blust, 1993]

## Data: Romance

University of California
C A L
N L P
Berkeley

| Gloss | Latin | Italian | Spanish | Portuguese |
|-------|-------|---------|---------|------------|
| Word/verb | verbum | verbo | verbo | verbu |
| Fruit | fructus | frutta | fruta | fruta |
| Laugh | ridere | ridere | reir | rir |
| Center | centrum | centro | centro | centro |
| August | augustus | agosto | agosto | agosto |
| Swim | natare | nuotare | nadar | nadar |

## Learning: Objective

University of California
C A L
N L P
Berkeley

$z$

$w$

/fokus/
/fogo/
/twoko/   /fweyo/   /fogo/
$\theta$

$$\max_{\theta} P(\theta|w_1 \dots w_L)$$

$$\max_{\theta, z} P(\theta, z|w_1 \dots w_L)$$

## Learning: EM

University of California
C A L
N L P
Berkeley

/fokus/
/fogo/
/twoko/   /fweyo/   /fogo/
$\theta$

/fokus/
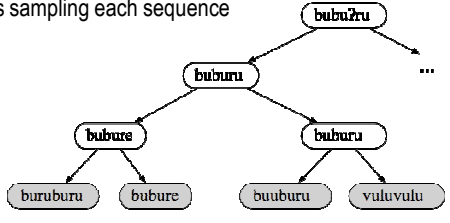/fogo/
/twoko/   /fweyo/   /fogo/
$\theta$

- M-Step
  - Find parameters which fit (expected) sound change counts
  - Easy: gradient ascent on theta

- E-Step
  - Find (expected) change counts given parameters
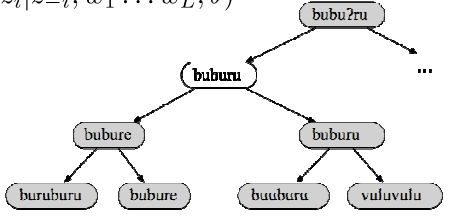  - Hard: variables are string-valued

## Computing Expectations

University of California
C A L
N L P
Berkeley

[Holmes 01, BGK 07]

Standard approach, e.g. [Holmes 2001]:
Gibbs sampling each sequence

bubu?ru
buburu
bubure          buburu
buruburu   bubure   buuburu   vuluvulu

...

'grass'

## A Gibbs Sampler

University of California
C A L
N L P
Berkeley

$$P(z_i|z_{-i}, w_1 \dots w_L, \theta)$$

bubu?ru
buburu
bubure          buburu
buruburu   bubure   buuburu   vuluvulu

...

'grass'

## A Gibbs Sampler



'grass'

## A Gibbs Sampler



'grass'

## Getting Stuck



How to jump to a state where the liquids
/r/ and /l/ have a common ancestor?

## Getting Stuck



## Solution: Vertical Slices

[BGK 08]

Single
Sequence
Resampling

Ancestry
Resampling



## Details: Defining "Slices"

The sampling domains (kernels) are indexed by contiguous
subsequences (*anchors*) of the observed leaf sequences



anchor

Correct construction section($G$) is
non-trivial but very efficient

## Results: Alignment Efficiency

Is ancestry resampling faster than basic Gibbs?

Hypothesis: Larger gains for deeper trees

Setup: Fixed wall time

Synthetic data, same parameters



Depth of the phylogenetic tree

---

## Results: Romance

| Gloss | Latin | Italian | Spanish | Portuguese |
|---|---|---|---|---|
| Word/verb | verbum | verbo | verbo | verbu |
| Fruit | fructus | frutta | fruta | fruta |
| Laugh | ridere | ridere | reir | rir |
| Center | centrum | centro | centro | centro |
| August | augustus | agosto | agosto | agosto |
| Swim | natare | nuotare | nadar | nadar |

---

## Learned Rules / Mutations



/werbum/ (la)

m →
u → o
w → v

m → / _ #
u → o / _
w → v / many environments
...

/verbo/ (vl)

r → f        e → ε

coluber    non colober
passim     non passi

---

## Learned Rules / Mutations



u → o / many environments
v → b / init. or intervocal.
t → t e / ALV _ #
...

/verbo/ (ib)

v → b        u → o

/berbo/ (es)    /verbu/ (pt)

r → f

---

## Comparison to Other Methods

▪ Evaluation metric: edit distance from a reconstruction made by a linguist (lower is better)

▪ Comparison to system from [Oakes, 2000]

   ▪ Uses exact inference and deterministic rules

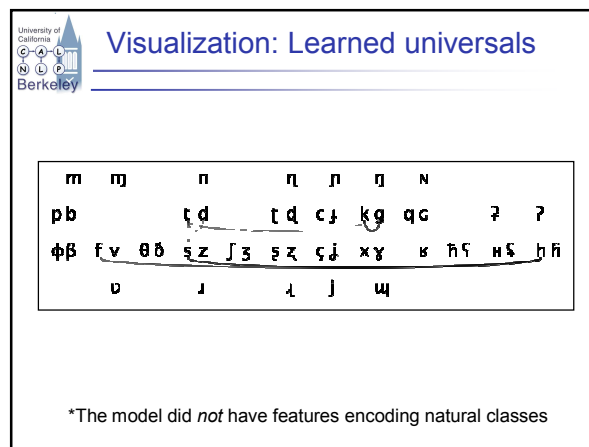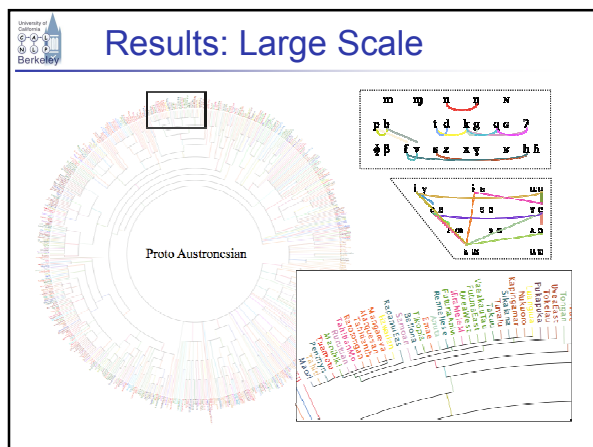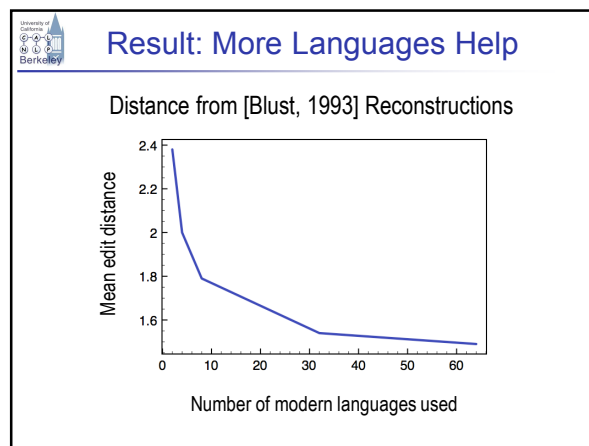   ▪ Reconstruction of Proto-Malayo-Javanic cf [Nothefer, 1975]



2.20
1.65
1.10
0.55
0

Oakes    Us

---

## Data: Oceanic



Proto-Oceanic

---

## Data: Oceanic

| Gloss | Hawai'ian | Maori | Samoan | Tongan |
|-------|-----------|-------|--------|--------|
| 'break' | haki | whati | fati | fasi |
| 'house' | hale | whare | fale | fale |
| 'yam' | uhi | uhi | ufi | ufi |
| 'woman' | wahine | wahine | fafine | fefine |
| 'moon' | mahina | mahina | masina | mahina |

http://language.psy.auckland.ac.nz/austronesian/research.php

## Result: More Languages Help

### Distance from [Blust, 1993] Reconstructions



## Results: Large Scale



## Visualization: Learned universals



*The model did *not* have features encoding natural classes

## Regularity and Functional Load

In a language, some pairs of sounds are more contrastive than others (higher functional load)

**Example:** English "p"/"b" versus "t"/"th"

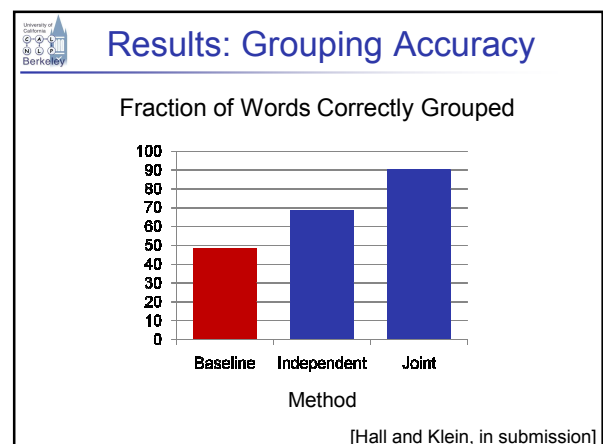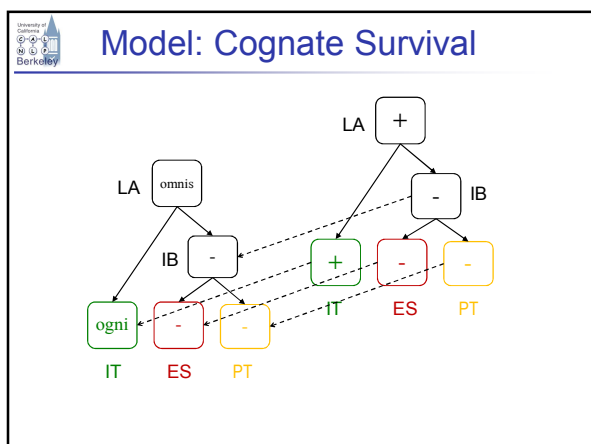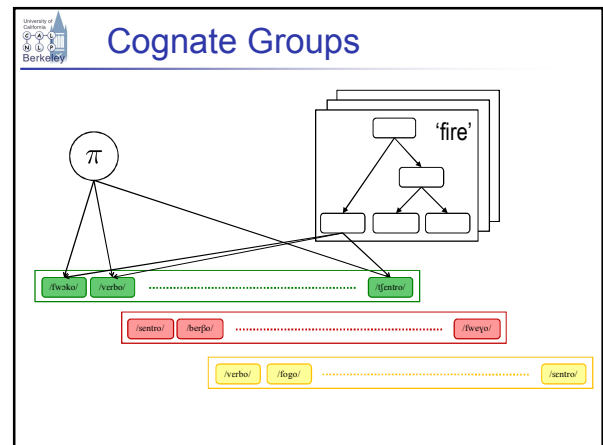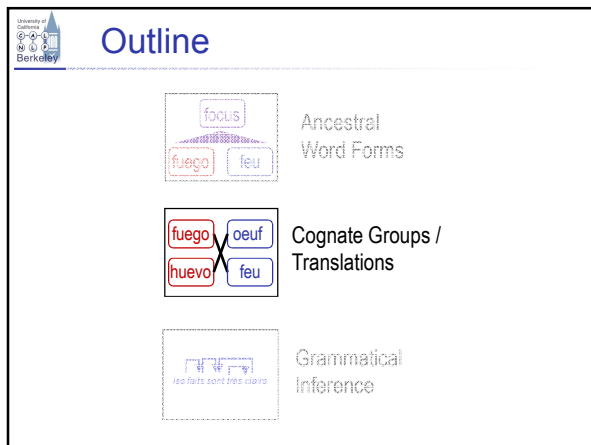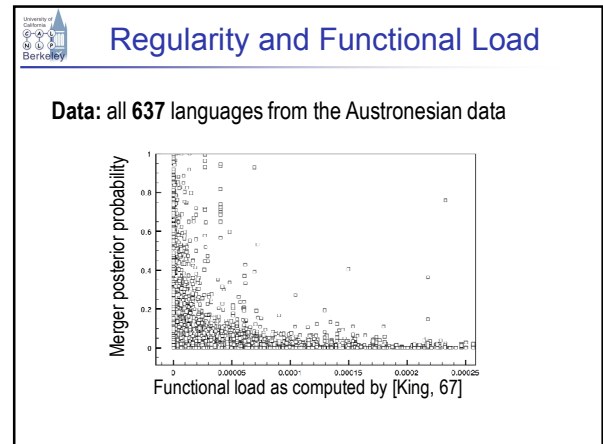"p"/"b": pot/dot, pin/din, dress/press, pew/dew, ...

"t"/"th": thin/tin

## Functional Load: Timeline

**Functional Load Hypothesis (FLH):** sounds changes are less frequent when they merge phonemes with high functional load    [Martinet, 55]

**Previous research within linguistics**: "FLH does not seem to be supported by the data" [King, 67]

**Caveat**: only four languages were used in King's study [Hocket 67; Surandran et al., 06]

**Our work:** we reexamined the question with two orders of magnitude more data [BGK, *under review*]

## Regularity and Functional Load

**Data:** only **4** languages from the Austronesian data



*Each dot is a sound change identified by the system*

Y-axis: Merger posterior probability
X-axis: Functional load as computed by [King, 67]

## Regularity and Functional Load

**Data:** all **637** languages from the Austronesian data



Y-axis: Merger posterior probability
X-axis: Functional load as computed by [King, 67]

## Outline



focus / fuego / feu — Ancestral Word Forms

fuego / oeuf / huevo / feu — Cognate Groups / Translations

les faits sont très clairs — Grammatical Inference

## Cognate Groups



'fire'

π

/fwoko/ /verbo/ ............ /tʃentro/

/sentro/ /berʃo/ ............ /fweɣo/

/verbo/ /fogo/ ............ /sentro/

## Model: Cognate Survival



LA +
LA omnis
IB -
IB -
+ - -
ogni - -
IT ES PT
IT ES PT

## Results: Grouping Accuracy

Fraction of Words Correctly Grouped



Method: Baseline, Independent, Joint

[Hall and Klein, in submission]

## Semantics: Matching Meanings

EN | day

Occurs with:
"night"
"sun"
"week"

tag | EN

Occurs with:
"name"
"label"
"along"

DE
tag

Occurs with:
"nacht"
"sonne"
"woche"

## Outline

Ancestral Word Forms
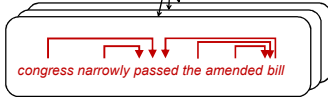
focus
fuego | feu

Cognate Groups / Translations

fuego | oeuf
huevo | feu

Grammatical Inference

les faits sont très clairs

## Grammar Induction

**Task:** Given sentences, infer grammar (and parse tree structures)

$\theta_{EN}$

congress narrowly passed the amended bill

$\theta_{FR}$

les faits sont très clairs

## Shared Prior

$$P(\theta) = N(0, \sigma^2 I)$$

$\theta$

$$P(\theta_\ell|\theta) = N(\theta, \sigma^2 I)$$

$\theta_{EN}$

$\theta_{FR}$

congress narrowly passed the amended bill

les faits sont très clairs

## Results: Phylogenetic Prior

Avg rel gain: 29%

GL
IE
G
WG  NG
RM

English, Dutch, Danish, Swedish, Spanish, Portuguese, Slovene, Chinese

70 60 50 40 30 20 10 0

## Conclusion

- Phylogeny-structured models can:
  - Accurately reconstruct ancestral words
  - Give evidence to open linguistic debates
  - Detect translations from form and context
  - Improve language learning algorithms

- Lots of questions still open:
  - Can we get better phylogenies using these high-res models?
  - What do these models have to say about the very earliest languages? Proto-world?

Thank you!

University of
California

Berkeley

nlp.cs.berkeley.edu