# Estimating sound similarities

Michael Cysouw, Jelena Prokić, Steve Moran

# Sound Similarities

- Manually specified
  (Kondrak 2002; Heeringa 2004)

- Hidden Markov Models
  (Bhargava & Kondrak 2009)

- Regular multi-alignment
  (Prokić 2010; Steiner, Stadler & Cysouw 2011)

- Bayesian inference
  (Prokić 2010)

- Investigating almost identical words
  (Holman, Brown & Wichmann 2011)

# Graphemic Normalization

- Widespread idea:
  "Convert everything into IPA"

- IPA is just another orthography !
  (only approximation of sound)

- Still: sound-based normalization is practical
  (but there are strong differences !)

- But: can we do without ?

# Estimating sound similarities

Michael Cysouw, Jelena Prokić, Steve Moran

# Estimating symbol similarities

Michael Cysouw, Jelena Prokić, Steve Moran

# Graphemic parsing

- **Unicode normalization**
  ő vs. o ˜ ´

- **Orthographic parse**
  (separate orthographic units as used in the source: "graphemes")

- **Orthographic normalization**
  (research specific!)

# Graphemic parsing

- **Code points** (7) t s ʰ o ˜ ´ :
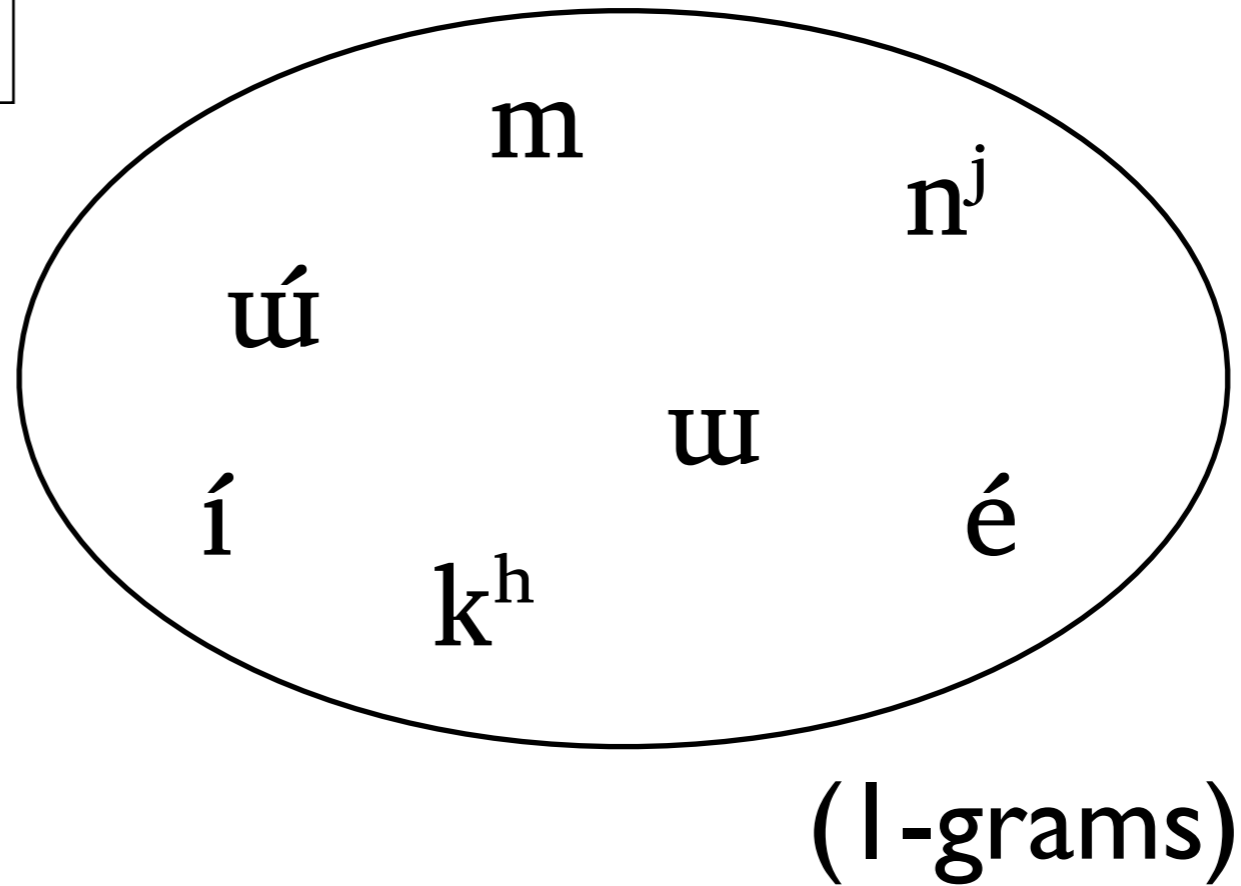
- **Characters** (4) t sʰ ő :

- **Graphemes** (2) tsʰ ő:

# Cross-script mapping

| E | R | freq | dice |
|---|---|------|------|
| r | р | 184 | 0.88874745 |
| n | н | 115 | 0.8461936 |
| l | л | 104 | 0.79646295 |
| s | с | 114 | 0.7927922 |
| t | т | 165 | 0.7701921 |
| m | м | 47 | 0.7699933 |
| o | о | 184 | 0.7510106 |
| k | ть | 21 | 0.74458015 |
| p | п | 50 | 0.7388723 |
| i | и | 102 | 0.7034591 |
| a | а | 221 | 0.6866478 |
| u | у | 40 | 0.6449104 |
| c | к | 77 | 0.6251676 |
| e | е | 219 | 0.59066784 |
| b | б | 32 | 0.525643 |
| w | в | 46 | 0.46787763 |
| d | д | 42 | 0.381996 |
| ⋮ | ⋮ | ⋮ | ⋮ |

Cysouw, Michael & Hagen Jung. 2007. Cognate identification and alignment using practical orthographies. Proceedings of Ninth Meeting of the ACL Special Interest Group in Computational Morphology and Phonology, 109-116.

# 'bag of symbol" approach

mín$^j$éék$^h$ɰú



m

n$^j$

ú

ɰ

í

é

k$^h$

(1-grams)

# 'bag of symbol" approach

mínʲéékʰɯú

mí

nʲé

kʰɯ

ɯú

ínʲ

éé

ékʰ

(2-grams)

# Cross-script mapping

- Ignore linear structure of words "bag of symbols" approach

- Use parallel wordlist to estimate co-occurrences of n-grams

- N-grams that have a high probability of co-occurrence in parallel meaning are interested for historical linguistics

## Chocó[†]
DR  hɨwá hã
CT  huwá hã
CM  húa hã
TD  húa hã
EP  húa hã
BA  húa hã
WM  hua hã

## Chibcha
IK  gúnnɨ
KO  káu-xəldə
DM  gúla
CL  háttakɾaʔ
TN  átkara
BI  ʔakdurə, ʔaxdurə

## Barbacoa
PA  kuse
GU  tasik
TR  ʼkwal
AW  tʃtɨ
TP  tẹde
CH  tʲaapa

## Kamsá
KS  kukuá-tçe

## Quechua
IN  máki

## Arawak
WY  wahápɨ[†]
AC  -káahɨ
CR  no-kapi
PP  wa-káapi
YC  wa-játeʔela
TO  paséwa[†]
CA  kaapi[†]
BN  -káapi
RE  -ké

## Tucano
TC  õbṍ-kã
WN  wãʼbõ bãʼkã
PY  ə̃ʼbõ-kã
WA  ắʼbṍ
BR  ã́ʼbõ
TY  wãʼbṍ-pãbã
YR  ʼwắbõ
DE  bõhõ̃ʼtṍõ
SR  bõʼhṍ
TA  wãʼbṍ
CP  wãʼbṍ
MA  ãbõ
BS  ãbõ
TM  pita-ʼka
CU  pɨʼrɨ́
KG  ʼhɨ̃tʰɨ
SI  ɨ̃́ʼtɨ saʼda[†]
SE  ʼhɨ̃́ʼtɨ
OR  hɨ̃tɨ

## Carib
CJ  eeʼnʲari
YK  óma

## Guahibo
PL  pe-kóbe
GH  pe-kóbe:
CI  pé-kobe, pe-kóbe
JT  pe-kó
GY  peh keʔé

## Sáliba-Piaroa
SL

## Macú-Puinave
PU  mo lap
NK  teiʔ3
KK  tejʔ2-ja4
JU  depũ̃j

## Witoto
MR  ono-ʤɯ, onoɯ
MN  ónoɯ
NP  ónoɯ
OC  oŋõõ(po)
MU  úse
BO  (mé)-ʔóxtsʰɨ́ɨ
MÑ  óhtsʰɨ́ɨ, íʔóhtsʰɨ́ɨ

| | Bora | Muinane |
|---|---|---|
| down | tʃîⁿnʲe, paári | báari, gíino |
| bee | íímúʔóexpʰi, téʔtsʰipa | nîɨbɨri, mîɨbɨriʔɨ |
| sharp | tsʰúʔxɨβáne | síîxéβano |
| … | … | … |

| Bora | Muinane | Bora | Muinane | Bora | Muinane |
|------|---------|------|---------|------|---------|
| #k | #kʰ | #i | #i | #n | #n |
| kɨ | kʰɯ | #a | #a | #m | #m |
| se | tsʰɨ | di | ti | mɨ | mɯ |
| xe | xɨ | du | to | nɨ | nɯ |
| ga | kʷa | #d | #t | us | tsʰɨ |
| ba | pa | #s | #tsʰ | #t | #tʰ |
| #b | #p | gi | tʃi | ɨg | ɯkʷ |
| e# | ɨ# | ni | ni | #ɸ | #pʰ |

# Using bigram matching

- Bora 'two':  mínʲéékʰɯú

- Muinane 'two':  míínokɨ

|      | #m | mi | ii | in | no | ok | kɫ | ɨ# |
|------|----|----|----|----|----|----|----|----|
| #m   | 22 | 3  | 2  | 2  | 2  | 2  | 2  | 2  |
| mi   | 4  | 12 | 2  | 2  | 5  | 1  | 1  | 1  |
| inʲ  | 2  | 1  | 5  | 9  | 3  | 1  | 1  | 2  |
| nʲe  | 1  | 1  | 5  | 5  | 4  | 1  | 1  | 2  |
| ee   | 3  | 3  | 3  | 3  | 6  | 2  | 2  | 2  |
| ekʰ  | 1  | 2  | 1  | 1  | 4  | 2  | 3  | 2  |
| kʰɯ  | 2  | 2  | 2  | 2  | 2  | 1  | 23 | 2  |
| ɯɯ   | 2  | 2  | 3  | 3  | 2  | 2  | 4  | 4  |
| ɯ#   | 2  | 2  | 3  | 2  | 3  | 1  | 3  | 4  |

| | #m | mi | ii | in | no | ok | kɨ | ɨ# |
|---|---|---|---|---|---|---|---|---|
| #m | 22 | 3 | 2 | 2 | 2 | 2 | 2 | 2 |
| mi | 4 | 12 | 2 | 2 | 5 | 1 | 1 | 1 |
| inʲ | 2 | 1 | 5 | 9 | 3 | 1 | 1 | 2 |
| nʲe | 1 | 1 | 5 | 5 | 4 | 1 | 1 | 2 |
| ee | 3 | 3 | 3 | 3 | 6 | 2 | 2 | 2 |
| ekʰ | 1 | 2 | 1 | 1 | 4 | 2 | 3 | 2 |
| kʰɯ | 2 | 2 | 2 | 2 | 2 | 1 | 23 | 2 |
| ɯɯ | 2 | 2 | 3 | 3 | 2 | 2 | 4 | 4 |
| ɯ# | 2 | 2 | 3 | 2 | 3 | 1 | 3 | 4 |

|  | Ocaina | Witoto Murui |
|---|---|---|
| HAND | on̥õõ | onodʒɯ |
| WE | xoxo | koko |
| HOUSE | ɸoo | ɸo |
| DOG | hõʔxo | hɯko |
| JAGUAR | hõʔxo | hɯko |
| FATHER | mõõ | moo |
| HUMMINGBIRD | ɸaʔtííʔtʲo | ɸiθido |
| TREE | am̥ũũɲa | amena |
| STICK | am̥ũũɲa | amena |
| WHO | bõ | bu |
| SLEEP | ɯ́ɯ́nõ | ɯnɯ |
| AGOUTI | ɸɯ́ɯ́tʲo | ɸɯdo |
| THIS | bĩi | bie |
| THIS | baʔi | bie |
| NAME | maam̥ɯ | mamekɯ |
| DAY | mooɲa | aremona |
| BOW | tsipóxatʲa | θɯkuira |
| HEAR | xaaxa | kakade |
| DAY | moɲamó | aremona |

|  | Ocaina | Witoto Murui |
|---|---|---|
| GREASE | ɸaɦĩi | ɸare |
| YOU (PLURAL) | mõʔ | omoɯ |
| THIS | bɯ | bie |
| ARROW | oɯdʲáátʲa | dukɯraθɯ |
| SPEAR | oɯdʲáátʲa | dukɯraθɯ |
| LIP | ɸaʔóóʔko | ɸue igoɯ |
| GREEN | moxóóso | mokorede |
| I | xõ | kue |
| ONE | tʲa | dahe |
| WE | xo | kaɯ |
| TOOTH | aʔtiiʔtʲo | iθido |
| MOUTH | ɸooɯ | ɸue |
| BELLY | gááho | hebe |
| FATHER | mõõhõ | moo |
| YOU (PLURAL) | mõʔxo | omɯko |
| SWAMP | xonɯ́ɯ́βaga | kɯnere |
| RAT | mɯɲõõko | miɲɯe |
| PATH, TRAIL | naahõ | naɯθo |
| OWL | mõõn̥õhõ | monuiθɯ |