

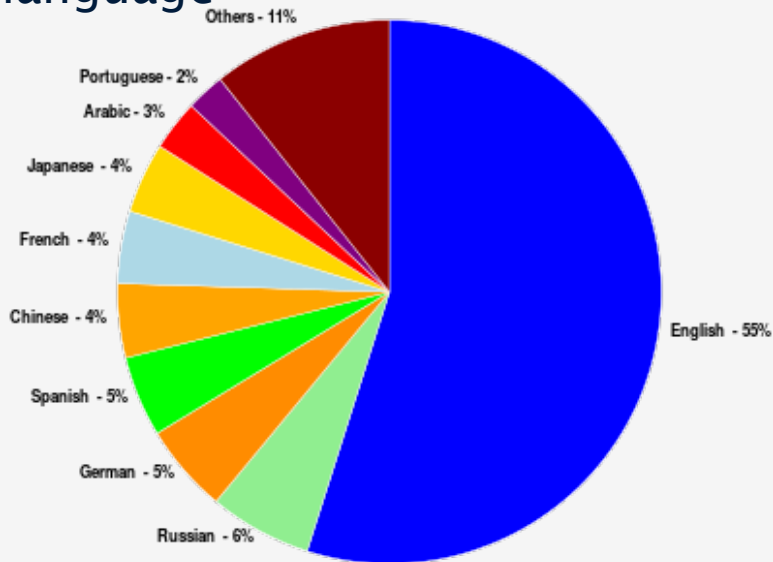
**NLP**

# Introduction to NLP

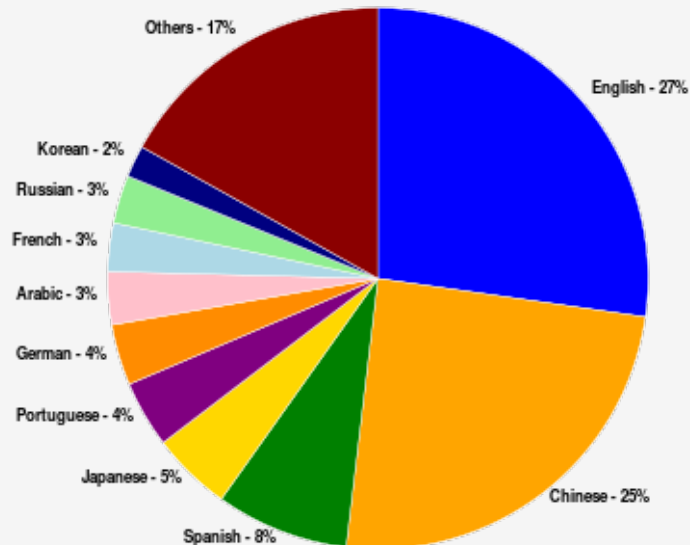
*Machine Translation*

# Multilingual Users

- Content languages for websites language



- Percentage of Internet users by language



April 2013

[http://en.wikipedia.org/wiki/Global\\_Internet\\_usage](http://en.wikipedia.org/wiki/Global_Internet_usage)



[The Tower of Babel, by Pieter Bruegel the Elder, 1563]

# The Rosetta Stone



Carved in 196 BC in Egypt

Deciphered by Champollion in 1822

Mixture of Egyptian (hieroglyphs and Demotic) and Greek

<http://www.ancientegypt.co.uk/writing/rosetta.html>

## NACLO Problem

- <http://nacloweb.org/resources/problems/2012/N2012-C.pdf>
- <http://nacloweb.org/resources/problems/2012/N2012-CS.pdf>
- Problem by Simon Zwarts, based on work by Kevin Knight

# Arcturan Problem – 1/4

*It's hard enough to translate between languages when you understand both languages. It's harder still when you only understand one. But what do computers do? They don't truly understand either language. To illustrate the challenge that computers face, Kevin Knight posed this classic puzzle (Knight 1997): given two equivalent texts in two unknown alien languages, how would you go about translating one to another?*

It is the year 2354 AD. Our scientists have been eavesdropping on messages between two alien civilizations for a very long time, but we have never met either. The closest aliens, the Centauri, have finally begun to communicate with us. Their first message was a message of peace, “Farok crrok hihok yorok klok kantok ok-yurp.”

Now, we know that the Centauri have been in contact for some time with the Arcturan race, who live in another solar system. We have never had contact with the Arcturans, but newly developed technology makes it possible for us to send them a message. We would like to send them, first, a message of peace, but because we do not understand their language, this is not an easy task.

Luckily, we have intercepted communications from the Centauri that include both languages. Here are 12 sentences in Centauri and their 12 translations in Arcturan. Unfortunately, because we have only been eavesdropping, their meaning is unknown. However, we do know that the sentence pairs on each line are translations of each other. We want to use this information to translate the original peace message from the Centauri and then send this to the Arcturans. Your assignment will be to do this translation.



# Arcturan Problem – 2/4

## CENTAURI

ok-voon ororok sprok.

ok-drubel ok-voon anak plok sprok.

erok sprok izok hihok ghrok.

ok-voon anak drok brok jok.

wiwok farok izok stok.

lalok sprok izok jok stok.

lalok farok ororok lalok sprok izok enemok.

lalok brok anak plok nok.

wiwok nok izok kantok ok-yurp.

lalok mok nok yorok ghrok klok.

lalok nok crrok hihok yorok zanzanok.

lalok rarok nok izok hihok mok.

## ARCTURAN

at-voon bichat dat.

at-drubel at-voon pippat rrat dat.

totat dat arrat vat hilat.

at-voon krat pippat sat lat.

totat jjat quat cat.

wat dat krat quat cat.

wat jjat bichat wat dat vat eneat.

iat lat pippat rrat nnat.

totat nnat quat oloat at-yurp.

wat nnat gat mat bat hilat.

wat nnat arrat mat zanzanat.

wat nnat forat arrat vat gat.



# Arcturan Problem – 3/4

**C-1** Let's start with the first Centauri word: "farok". This word occurs in two of our Centauri sentences. Given that these sentences' Arcturan translations only have one word in common with each other, we can assume that this word is the translation for "farok". Which word it is?

farok

[illegible]

**C-2** Do the same thing for “hihok” and “yorok”. For “yorok” you will need to make some assumptions about word ordering.

# hihok

[illegible]

yorok

[illegible]

**C-3** The Centauri word “clók” only occurs once. However, you can figure out its Arcturan translation in another way.

clock

[illegible]

**C-4** Try to use the processes from the previous assignments to complete as much as possible of the following table.

crrok

[illegible]

kantok

[illegible]

ok-yurp

[illegible]

[illegible][illegible]

# Arcturan Solution – 1/3

**C-I** The questions in this assignment are based on examples in Knight (1997). In fact, both Centauri and Arcturan have underlying real world languages, as it turns out Centauri is English and Arcturan is Spanish. The languages are obfuscated to Centauri and Arcturan in order to illustrate how a Statistical Machine Translation (SMT) system must start from scratch, since it has no prior knowledge of how the languages work.

## CENTAURI

Ok-voon ororok sprok.  
Garcia and associates.

Ok-drubel ok-voon anak plok sprok.  
Carlos Garcia has three associates.

Erok sprok izok hihok ghrok.  
His associates are not strong.

Ok-voon anak drok brok jok.  
Garcia has a company also.

Wiwok farok izok stok.  
Its clients are angry.

Lalok sprok izok jok stok.  
The associates are also angry.

## ARCTURAN

At-voon bichat dat.  
Garcia y asociados.

At-drubel at-voon pippat rrat dat.  
Carlos Garcia tiene tres asociados.

Totat dat arrat vat hilat.  
Sus asociados no son fuertes.

At-voon krat pippat sat lat.  
Garcia tambien tiene una empresa.

Totat jjat quat cat.  
Sus clientes están enfadados.

Wat dat krat quat cat.  
Los asociados tambien están enfadados.

## Arcturan Solution – 2/3

Lalok farok ororok lalok sprok izok enemok.  
The clients and the associates are enemies.

Lalok brok anok plok nok.  
The company has three groups.

Wiwok nok izok kantok ok-yurp.  
Its groups are in Europe.

Lalok mok nok yorok ghrok klok.  
The modern groups sell strong pharmaceuticals.

Lalok nok crrrok hihok yorok zanzanok.  
The groups do not sell zanzanine.

Lalok rarok nok izok hihok mok.  
The small groups are not modern.

Wat jjat bichat wat dat vat eneat.  
Los clientes y los asociados son enemigos.

lat lat pippat rrat nnat.  
La empresa tiene tres grupos.

Totat nnat quat oloat at-yurp.  
Sus grupos están en Europa.

Wat nnat gat mat bat hilat.  
Los grupos modernos venden medicinas fuertes.

Wat nnat arrat mat zanzanat.  
Los grupos no venden zanzania.

Wat nnat forat arrat vat gat.  
Los grupos pequeños no son modernos.

# Arcturan Solution – 3/3

The novel sentence which was offered for translation in English is: “clients do not sell pharmaceuticals in Europe.”

## Answers

C-1 jjat

C-2 hihok = arrat, yorok = mat

C-3 We need to use the process of elimination, when mapping all the words between the two sentences two words are unaligned, we assume these are translations of each other. Thus, klok = bat.

C-4 Here are the new matches:

crrok	(empty)
kantok	oloot
ok-yurp	at-yurp

“crrok” does not seem to have a Arcturan equivalent, like in English the word “do” is not translated in “do not sell” which simply becomes “not sells” in Spanish. (Or to put it another way, the Centauri word *crrok* has a translation, but it’s the “empty” word.)

C-5 jjat arrat mat bat oloot at-yurp

Since you cannot deduce with certainty the exact order of the Arcturan sentence, various orders of these words will be accepted.

C-6 Immediately, you are faced with a dilemma: should you translate *totat* as *erok* or *wiwok*? Because *wiwok* occurs more frequently and because you’ve never seen *erok* followed by any of the other words you’re considering, *wiwok* seems more likely. (However, admittedly, this is only a best guess, and *erok* will also be accepted.) Next, you consider various word orders. There appears to be no grammatical path through these words. Suddenly, you remember that curious Centauri word *crrok*, which had no translation. *Crrok*, however, turns out to be a natural bridge between *nok* and *hihok*, giving you the translation:

*wiwok rarok nok crrok hihok yorok klok.*

# Parallel Corpora

- The Rosetta Stone
- The Hansards Corpus
- The Bible

# Hansards Example

- English

- `<s id=960001>` I would like the government and the Postmaster General to agree that we place the union and the Postmaster General under trusteeship so that we can look at his books and records, including those of his management people and all the memos he has received from them, some of which must have shocked him rigid.
- `<s id=960002>` If the minister would like to propose that, I for one would be prepared to support him.

- French

- `<s id=960001>` Je voudrais que le gouvernement et le ministre des Postes conviennent de placer le syndicat et le ministre des Postes sous tutelle afin que nous puissions examiner ses livres et ses dossiers, y compris ceux de ses collaborateurs, et tous les mémoires qu'il a reçus d'eux, dont certains l'ont sidéré.
- `<s id=960002>` Si le ministre voulait proposer cela, je serais pour ma part disposé à l'appuyer.



# English-Cebuano Bible Example

In the beginning God created the heaven and the earth.  
Sa sinugdan gibuhad sa Dios ang mga langit ug ang yuta.

And God called the firmament Heaven.  
Ug gihinganlan sa Dios ang hawan nga Langit.

And God called the dry land Earth  
Ug ang mamala nga dapit gihinganlan sa Dios nga Yuta

- use: co-occurrence, word order, cognates
- corpora are needed
- sentence alignment needs to be done first

[http://en.wikipedia.org/wiki/Bible\\_translations\\_by\\_language](http://en.wikipedia.org/wiki/Bible_translations_by_language)

## NACLO Problem

- <http://nacloweb.org/resources/problems/2012/N2012-D.pdf>
- <http://nacloweb.org/resources/problems/2012/N2012-DS.pdf>
- Problem by Dragomir Radev

Many languages are related to each other for historical reasons. They may have a common ancestor or they may have borrowed words from each other. Linguists group languages into families and branches, based on their common ancestry.

Here is a list of translations of the first article of the Universal Declaration of Human Rights in 17 languages:

Your task is to identify similarities among these languages and group them into seven clusters (groups) of related languages as sketched in the diagram below:

Here is a list of translations of the first article of the Universal Declaration of Human Rights in 17 languages:

A. (English) All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of brotherhood.

B. (Latin) Omnes homines dignitate et iure liberi et pares nascuntur, rationis et conscientiae participes sunt, quibus inter se concordiae studio est agendum.

C. Vsi ljudje se rodijo svobodni in imajo enako dostojanstvo in enake pravice. Obdarjeni so z razumom in vestjo in bi morali ravnati drug z drugim kakor bratje.

D. Dieub ha par en o dellezegezh hag o gwirioù eo ganet an holl dud. Poell ha skiant zo dezho ha dleout a reont bevañ an eil gant egile en ur spered a genvreudeuriezh.

E. Tuots umans naschan libers ed equals in dignità e drets. Els sun dotats cun intellet e conscienza e des-san agir tanter per in uin spiert da fraternità.

F. Toate ființele umane se nasc libere și egale în demnitate și în drepturi. Ele sunt înzestrațe cu rațiune și conștiință și trebuie să se comporte unii față de altele în spiritul fraternității.

G. Genir pawb yn rhydd ac yn gydradd â'i gilydd mewn urddas a hawliau. Fe'u cynysgaeddir â rheswm a chydwybod, a dylai pawb ymddwyn y naill at y llall mewn ysbryd cymodlon.

H. Visi žmonės gimsta laisvi ir lygūs savo orumu ir teisėmis. Jiems suteiktas protas ir sąžinė ir jie turi elgtis vienas kito atžvilgiu kaip broliai.

I. Totu sos èsseres umanos naschint liberos e eguales in dinnidade e in deretos. Issos tenent sa resone e sa cussèntzia e depent operare s'unu cun s'àteru cun ispiritu de fraternidade.

J. Gizon-emakume guztiak aske jaiotzen dira, duintasun eta eskubide berberak dituztela; eta ezaguera eta kontzientzia dutenez gero, elkarren artean senide legez jokatu beharra dute.

K. Kai rahvas roittahes välinny da taza-arvozinnu omas arvos da oigevuksis. Jogahizele heis on annettu mieli da omatundo da heil vältämättä pidäy olla keskenäh, kui vellil.

L. Všetci l'udia sa rodia slobodní a sebe rovní , čo sa týka ich dostôjnosti a práv. Sú obdarení rozumom a majú navzájom jednat' v bratskom duchu.

M. Nascinu tutti l'omi libari è pari di dignità è di diritti. Pussedinu a raghjoni è a cuscenza è li tocca ad agiscia trà elli di modu fraternu.

N. Saoláitear na daoine uile saor agus comhionann ina ndínit agus ina gcearta. Tá baidh an réasúin agus an choinsiasa acu agus dlíd iad féin d'iompar de mheon bhrathreachais i leith a chéile.

O. Visi cilvēki piedzimst brīvi un vienlīdzīgi savā pašcienā un tiesībās. Viņi ir apveltīti ar saprātu un sirdsapziņu, un viņiem jāizturas citam pret citu brālības garā.

P. Kaikki ihmiset syntyvät vapaina ja tasavertaisina arvoltaan ja oikeuksiltaan. Heille on annettu järki ja omatunto, ja heidän on toimittava toisiaan kohtaan veljeyden hengessä.

Q. Wszyscy ludzie rodzą się wolni i równi pod względem swej godności i swych praw. Są oni obdarzeni rozumem i sumieniem i powinni postępować wobec innych w duchu braterstwa.

# Solution

1. CLQ Slavic
2. BEFIM Romance
3. J Basque
4. HO Baltic
5. DGN Celtic
6. KP Finno-Ugric
7. A English

**NLP**