

Creating a Technical Corpus



**Summer School of Linguistics
PER-FIDE 2013**

Pedro Carvalho

pg22760@alunos.uminho.pt

Agenda



- Introduction
- Thesis Abstracts
 - Sources
 - Vertical Alignment
- ECTS
 - Sources
 - Linking the Pages
- Examples

Introduction



Every day countless technical documents are produced inside our learning institutions walls.

Even a quick inspection of their web pages reveal lots of useful technical information and terminology.

If these resources exist and are open to the public in general why won't we use them?

Thesis Abstracts - Source



The Portuguese Open Access Scientific Repository (RCAAP) contain a great collection of thesis abstracts translated to one or more languages.

When obtaining those abstracts some problems arose, but in this presentation the main focus will be the vertical alignment problem.

Thesis Abstracts - Vertical Alignment



In 30% of the cases the abstracts ended up vertically aligned with no possibility to detect the breaking point between different languages.

(...)

Acredita-se a todas estas actividades que a *Opuntia ficus indica* é muito rica nutricionalmente e daí o seu uso na alimentação. Medicinal plants have been used since Antiquity, by primitive civilizations, to fight illnesses. Nowadays, phytotherapy shows to be more and more sustained due to quality, safety and efficiency requirements, and therefore medicinal plants are being

(...)

Thesis Abstracts - Vertical Alignment



Strategy:

1. Break the chunk of text into phrases.

É um cacto arbustivo constituído por raiz, caule (cladódios) onde se encontram, consoante a época, flores e frutos (doces e suculentos). É uma planta cujos caules jovens (“nopalitos”) e frutos são usados na dieta alimentar dos Mexicanos.



É um cacto arbustivo constituído por raiz, caule (cladódios) onde se encontram, consoante a época, flores e frutos (doces e suculentos).



É uma planta cujos caules jovens (“nopalitos”) e frutos são usados na dieta alimentar dos Mexicanos.

Thesis Abstracts - Vertical Alignment



Strategy:

2. For each phrase:
 - a. Break it in words.

3. For each word:
 - a. Check if it belongs to any dictionaries.

4. In the end choose the dictionary with more words associated.

Thesis Abstracts - Vertical Alignment



É uma planta cujos caules jovens (“nopalitos”) e frutos são usados na dieta alimentar dos Mexicanos.

Is it portuguese?



Yes -> PT + 1

Is it english?



No -> ...

Is it spanish?



Yes -> ES + 1

What should we do in case of a tie?

Thesis Abstracts - Vertical Alignment



Strategy:

3. In case of a tie
 - a. Check the neighborhood.
 - b. Choose the dictionary with more words associated in the global context.

Thesis Abstracts - Vertical Alignment



(...)

Actualmente, a fitoterapia está cada vez mais sustentada por requisitos de eficácia, segurança e qualidade, sendo as plantas medicinais bastante procuradas como recurso terapêutico (Cunha, Silva e Roque, 2003).



PT

A *Opuntia ficus indica* é uma planta originária do México que cresce de forma selvagem em regiões áridas e semi-áridas de todo o mundo (Ennouri et al., 2006b).



PT

É um cacto arbustivo constituído por raiz, caule (cladódios) onde se encontram, consoante a época, flores e frutos (doces e suculentos).



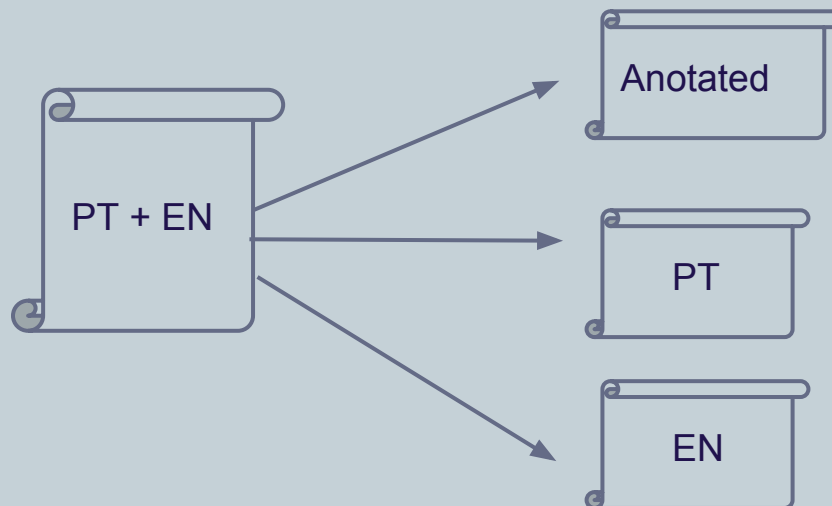
PT

(...)

Thesis Abstracts - Vertical Alignment



Each original document originates a final annotated document and a new document for each detected language.



We randomly chose 50 files that were produced by this tool.

In 1043 phrases that were analyzed we detected 13 errors.

Giving us a success percentage of 98.75%.

ECTS - The Source



In practically any learning institution's web page we can easily find study plans and information about all the courses they offer.

These pages are usually translated at least to english.

The goal was to obtain these pages and somehow link all the original versions to it's counterpart translated versions.

ECTS - Linking the Pages



apps.uc.pt/courses/en/programme/346/2013-2014

UNIVERSIDADE DE COIMBRA

United Nations Educational, Scientific and Cultural Organization

University of Coimbra – Alta and Sofia
inscribed on the World Heritage List in 2013

This uses cookies that do not gather any personal information whatsoever. By using this website, you agree with it's cookie policy. [more info](#)

Courses

[Home](#)

[List all organic units](#)

[< back to course](#)

The work developed by a student during one academic year, according to the indicative plan of the course and on a full totals 60 ECTS.

Common Core

Course unit title	Year	Regime	Type	Subject area	ECTS
Animal Diversity	1	1st Semester	Compulsory	BIO	4.
Biochemistry	1	1st Semester	Compulsory	BQ	6.
Cellular Biology	1	1st Semester	Compulsory	BIO	4.
Chemistry	1	1st Semester	Compulsory	Q	6.
Methods and Techniques in Cytology and Physiology	1	1st Semester	Compulsory	BIO	4.
Zoology	1	1st Semester	Compulsory	BIO	4.
Evolution	1	2nd Semester	Compulsory	BIO	3.
Mathematics	1	2nd Semester	Compulsory	M	6.
Methods and Techniques in Molecular Biology	1	2nd Semester	Compulsory	BIO	6.

ECTS - Linking the Pages



The pages that cannot be directly associated through its url can be linked through mechanization. However each page has its own characteristics.

To detect and possibly correct the alphabetically sorted lists we can use ucts (unambiguos-concept translation sets).

New ideas are welcome!!

Examples



THESIS - ABSTRACT

ECTS

PTD for : PT → EN

glutathiona (65 occurrences)

84.19%		<u>glutathione</u> ✓	77	<input type="button" value="→"/>
3.06%		<u>gsh</u> ✓	60	<input type="button" value="→"/>
1.24%		<u>occasional</u>	35	<input type="button" value="→"/>
1.23%		<u>redox</u> ✓	54	<input type="button" value="→"/>
1.21%		<u>forms</u>	862	<input type="button" value="→"/>
1.16%		<u>namely</u>	2298	<input type="button" value="→"/>
1.15%		<u>widely</u>	394	<input type="button" value="→"/>
1.07%		<u>tolerance</u>	269	<input type="button" value="→"/>

PTD for : PT → EN

ambiental (521 occurrences)

71.81%		<u>environmental</u> ✓	1046	<input type="button" value="→"/>
3.13%		<u>environment</u> ✓	950	<input type="button" value="→"/>
1.66%		<u>management</u>	3991	<input type="button" value="→"/>
0.94%		<u>development</u>	4114	<input type="button" value="→"/>
0.84%		<u>of</u>	101969	<input type="button" value="→"/>
0.63%		<u>noise</u> ✓	122	<input type="button" value="→"/>
0.51%		<u>monitoring</u> ✓	404	<input type="button" value="→"/>
0.50%		<u>transport</u> ✓	444	<input type="button" value="→"/>
0.46%		<u>subjects</u>	900	<input type="button" value="→"/>

Creating a Technical Corpus



**Summer School of Linguistics
PER-FIDE 2013**

Pedro Carvalho

pg22760@alunos.uminho.pt

Thank You :)