

# *Linguateca's* infrastructure for Portuguese... and how it allows the detailed study of language varieties

Diana Santos



# A map of the talk

- Brief introduction of *Linguateca*
  - An infrastructure for Portuguese language technology
  - Short history
- The linguistic analysis of running text
  - Corpus projects for Portuguese
  - Three *Linguateca* projects: AC/DC, Floresta Sintáctica, and CorTrad
- Studying variation and varieties with the AC/DC cluster
  - Data
  - Formal variational linguistics support
  - New capabilities

# Never heard about *Linguateca*?

- It is a government funded initiative to significantly raise the quality and availability of resources for the **computational processing of Portuguese**
- After an initial plan for discussion by the community (white paper) a network was launched, headed by a small group (Linguateca's Oslo node) at SINTEF ICT (formerly SINTEF Tele og Data)
- This network has had as main goal to guarantee that
  - Information was provided and gathered at one place on the Web
  - Resources were made public, maintained, and further developed in connection with the scientific community
  - Evaluation initiatives were launched

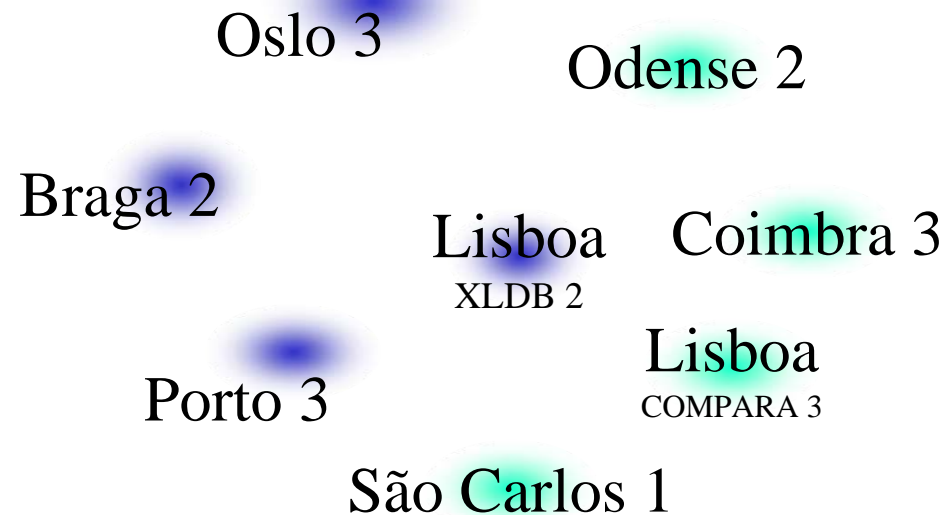
# Linguatca, a project for Portuguese

- A distributed resource center for Portuguese language technology

## IRE model

- Information
- Resources
- Evaluation

[www.linguatca.pt](http://www.linguatca.pt)



# *Linguatca* highlights, [www.linguatca.pt](http://www.linguatca.pt)

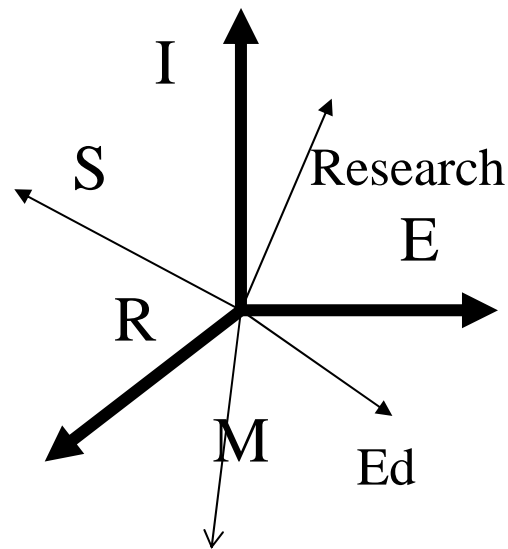
- *> 2000 links* More than 7,000,000 visits to the Web site
- [AC/DC](#), [CETEMPúblico](#), [COMPARA](#) ... Considerable resources for processing the Portuguese language
- *Morfolimpíadas* The first evaluation contest for Portuguese, followed by CLEF and HAREM
  
- Public resources
- Foster research and collaboration
- Formal measuring and comparison
- One language, many cultures
- Cooperation using the Internet
- Do not adapt applications from English

# Linguateca's premises: not a research project

- a project whose aim is to considerably improve the conditions of the community who deals with the computational processing of the Portuguese language
- Is processing of Portuguese = NLP specialized to Portuguese? **NO**
- Does one build a community just by financing individual research projects? **NO**
- One has to **build a research infrastructure** and actively foster collaboration and joint evaluation

# The IRE model and its evolution

- First: Information, Resources and Evaluation
- But then
  - (resource) Maintenance:
  - Support
  - Research (PhDs)
  - Education



# A document to discuss the future of the area

## ■ Main points: in 1998

- There was hardly anything publicly available
- People were alone doing the same things without knowledge of each other
- No evaluation whatsoever

## ■ Main need: an umbrella service

- Maintaining and making resources available cannot be considered research
- The sharing spirit for a common goal: open source philosophy
- No separation of commercial/industrial and academic venues



# At this moment, LINGUATECA is or has (produced)...

- Probably the largest repository on one language (computational processing) in the world (on the Web): kept at FCCN premises
- Well-known in the national communities (Portugal and Brazil) and in the international community (?)
- A set of reusable tools and resources that can be put to use by other researchers
- A set of studies on Portuguese and Portuguese processing (IR, GIR, MT, automatic terminology extraction, QA)
- A set of documents that enrich the area and can be used pedagogically
- A sizeable group of people trained in this area, a lot of others with some exposure to these activities through contact

# Linguateca's achievements

- A lot of publicly available resources
- Several evaluation contests which advanced the state of the art
- Information, dissemination, gathering of relevant data and a team who answers
- The first evaluation contest for Portuguese
- The first treebank for Portuguese
- The first Web-based corpus service for Portuguese
- The first QA system for Portuguese
- The largest revised and annotated parallel corpus in the world
- The first national Web snapshot available

# International impact

- Resources created by Linguateca available from the (Pennsylvania-based) Linguistic Data Consortium (LDC)
- Portuguese as one of the major languages in CLEF (more than 100 research groups worldwide participate in the largest evaluation forum for European languages and crosslingual information retrieval)
  - Linguateca belongs to the steering committee
  - Innovative pilots have been suggested by Linguateca, who has helped shaping the future
- The Portuguese treebank has often been used by third parties as example or resource in international venues, such as CoNLL or LREC
- According to Bernardo Magnini, Linguateca was the main inspiration for EVALITA, evaluation for Italian

# Evaluation contests (*avaliação conjunta*)

Model: DARPA and NIST eval. cont.

- Jointly agree on a task and discuss the details together
- Create an evaluation setup
  - measures
  - resources
  - procedure
- Compare the performance of the several systems and get a state of the art
- Make public both resources, programs and systems' outputs for
  - external validation
  - research on both the task and the evaluation methodology
  - organization of future evaluation contests
  - training of newcomers

# Linguistic analysis of running text

- Researchers on Portuguese needed support for computer-based empirical studies that were replicable and based on the same materials, available for extended periods of time, and that did not require physical access to specific premises
- Web-based services are the obvious answer, if they serve material that is curated and properly documented, and if they can be freely used
- AC/DC: providing access, making access possible
  - AC/DC cluster: a set of corpus projects, all inheriting from AC/DC, but with additional capabilities or features
  - Parallel corpora: COMPARA, CorTrad
  - Human revision: Floresta, COMPARA, ...

# A brief history of Portuguese corpus linguistics

In the 1970s, oral corpora were collected

- *Português Fundamental* (inspired by the *Français Fondamental*)
- Projeto NURC (Labov-inspired)

Both in Portugal and Brazil, continuation of corpus studies

- VARSUL, Variação Lingüística Urbana do Sul do País (1982- )
- CRPC, Corpus de Referência do Português Contemporâneo (1988- )

In the 1990s, due to better computer facilities, a renewal/revival

- 1994 - CIPM, Corpus informatizado do português medieval
- 1998 - Tycho Brahe, *Padrões rítmicos, domínios prosódicos ...*
- Projecto Natura, INESC, Corpus NILC/São Carlos, ...

# A brief history of Portuguese corpus linguistics (ct)

- Banco de português (199x-)
- CORDIAL-SIN...DUPLEX (1998-)
- Português Falado - Variedades Geográficas e Sociais (1995-97)

## International projects involving Portuguese

- CHILDES
- ENPC
- Borba-Ramsay corpus, ECI
- PORTEXT (1988-?)
- VISL (1994-)
- MLCC Multilingual and Parallel Corpora, *Official Journal of the EC*

# Portuguese corpora during LINGUATECA's lifetime

- Lácio-Web (2002-)
- C-ORAL-ROM (2001-2004)
- COMET (2005-)
- Corpus do português (2006-)
- etc.

- EuroParl
- Turigal
- JRC-Acquis

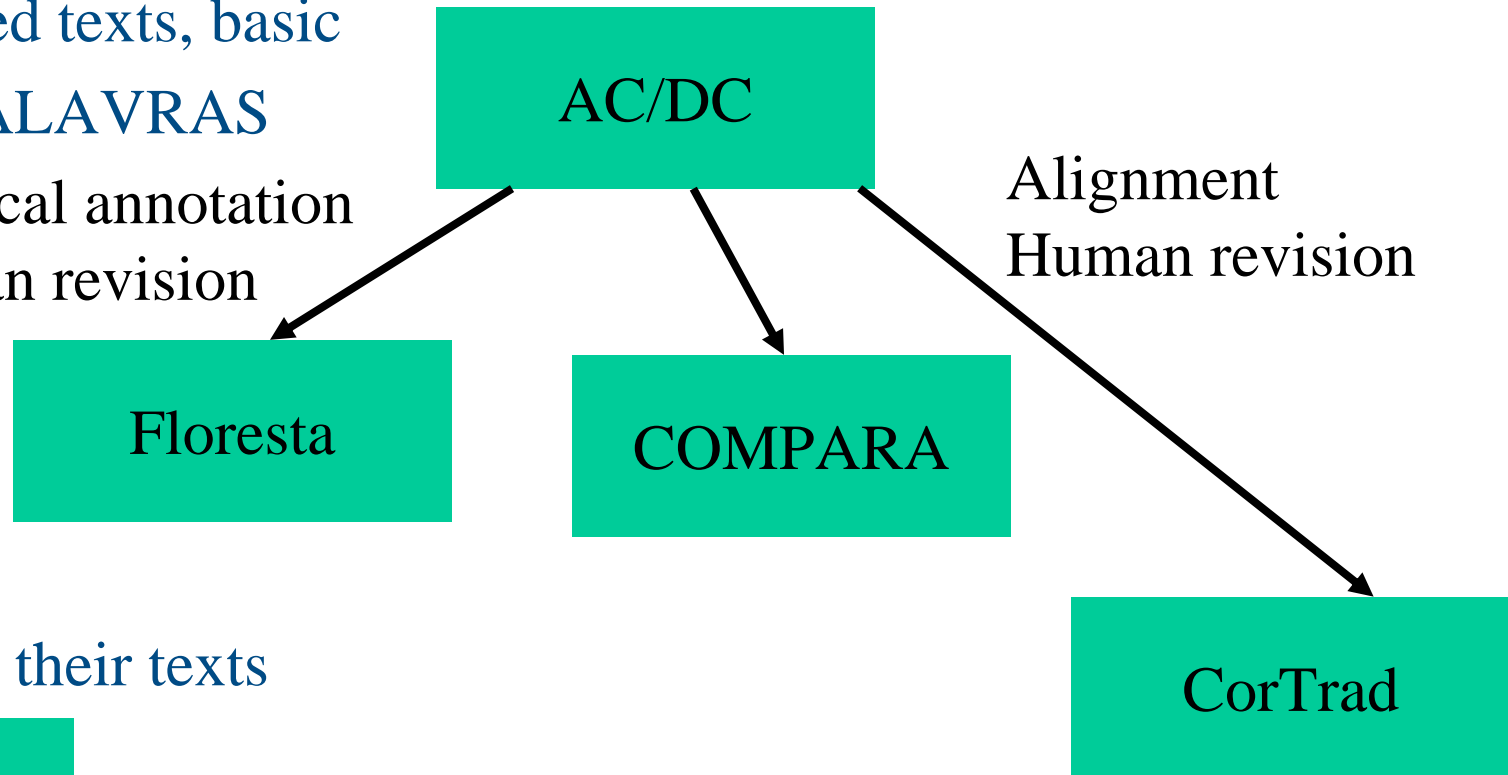
See also the ELC (*Encontros de linguística de corpus*) series in Brazil since 1999



# Similarities and differences in LINGUATECA corpora

- A set of closed texts, basic parsing from PALAVRAS

Hierarchical annotation  
Human revision



- Users choose their texts

Corpógrafo

# Corpus gallery in the AC/DC cluster

## ■ General newspapers

- CETEMPúblico
- CETENFolha (→ São Carlos)
- CHAVE
- Notícias de Moçambique



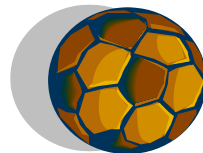
## ■ Regional newspapers

- NatMinho
- DiaCLAV
- Diário Gaúcho



## ■ Specific newspapers

- Sports : CONDIVport
- Political: Avante!
- Fashion: CONDIVport
- Health: CONDIVport
- Science: CorTradjorn



## ■ Literary

- Vercial
- ClassLPPE
- ENPCpub
- COMPARA
- CorTradlit



Adapted from Rocha (2007)

# Corpus gallery in the AC/DC cluster (cont.)

## ■ Oral documents

- Museu da Pessoa
- ECI-EBR falado
- Selva falado



## ■ Technical

- CorTradtec
- ECI-EE
- NILC/São Carlos tec
- Selva Ciência



## ■ Evaluation resources

- CDHAREM
- AmostRA
- FrasesPP

## ■ Email

- Listas: ANCIB
- SPAM: CoNE



## ■ Web

- Amazônia

## ■ “Historical”

- CETEMPúblico (primeiro milhão)
- NatPublico

Adapted from Rocha (2007)

# Brief description of AC/DC

- **A**cesso a **C**orpora / **D**isponibilização de **C**orpora
- Ca. 20 different corpora
- Ca. 360 million words, 16 million sentences
- Portuguese and Brazilian varieties, a few other texts from others
- Different genres, mainly contemporary
  
- Perl interface to the IMS (Open) CWB (corpus workbench)
- Common tokenization
- Use of the PALAVRAS parser (Bick, 2000) for linguistic annotation
- (Semi-automatic) annotation of selected semantic features

**Linguatca**

[Estrutura](#)  
[Equipa](#)

[Apresentação](#)  
[Acesso a recursos](#)

- [AC/DC](#)
- [- Procura](#)
- [- Corpora](#)
- [- Anotação](#)
- [- Exemplos](#)
- [CETEMPúblico](#)
- [CETENFolha](#)
- [CHAVE](#)
- [COMPARA](#)
- [Corpógrafo](#)
- [Esfinge](#)
- [Floresta Sintá\(ótica\)](#)
- [METRA](#)
- [PAPEL](#)
- [REPENTINO](#)
- [Repositório](#)
- [WebJspell](#)
- [WPT03](#)

[Catálogo de recursos](#)  
[Catálogo de ferramentas](#)  
[Catálogo de actores](#)  
[Catálogo de publicações](#)  
[Informação interessante](#)  
[Fórum](#)

# Projecto AC/DC: corpus Museu da Pessoa

[AC/DC : Linguatca](#)

O corpus **Museu da Pessoa** é um corpus de 109 entrevistas transcritas pelo [Núcleo Português do Museu da Pessoa](#) no âmbito dos seus projectos.

Procurar:

### Resultado:

- Concordância
- Distribuição das formas
- Distribuição dos lemas
- Distribuição da categoria gramatical (PoS)
- Distribuição do tempo verbal e/ou do caso pronominal
- Distribuição de pessoa e/ou número
- Distribuição do género
- Distribuição da função sintáctica
- Distribuição por entrevista

### Opções

- Resultados por ordem alfabética (só distribuições)

### Estrutura do corpus

Marcadores estruturais: **ent** [entrevista], **p** [parágrafo], **s** [frase], **resposta**, **pergunta**,

Veja um [excerto do corpus e informação adicional](#).

Tipo	Entrevistas
Variante(s)	PT BR
Tamanho (unidades)	456 mil
Tamanho (palavras)	315 mil

### [Página principal](#)

#### Procure noutros corpora:

- [AmostrA-NILC](#) [ANCIB](#) [Avante!](#) [CD HAREM](#)
- [CETEMPúblico](#)
- [CETEMPúblico \(primeiro milhão\)](#) [CHAVE](#)
- [Clássicos LP/Porto Editora](#) [CONDIVport](#)
- [CoNE](#) [DiaCLAV](#) [ECI-EBR](#) [ECI-EE](#)
- [ENPCPUB \(parte portuguesa\)](#) [FrasesPB](#)
- [FrasesPP](#) [Museu da Pessoa](#) [Natura/Minho](#)
- [Natura/Público](#) [NILC/São Carlos](#) [Vercial](#)

# Linguatca

## Estrutura

## Equipa

## Apresentação

## Acesso a recursos

- [AC/DC](#)
- [- Procura](#)
- [- Corpora](#)
- [- Anotação](#)
- [- Exemplos](#)
- [CETEMPúblico](#)
- [CETENFolha](#)
- [CHAVE](#)
- [COMPARA](#)
- [Corpógrafo](#)
- [Esfinge](#)
- [Floresta Sintáctica](#)
- [METRA](#)
- [PAPEL](#)
- [REPENTINO](#)
- [Repositório](#)
- [WebJspell](#)
- [WPT03](#)

## Catálogo de recursos

## Catálogo de ferramentas

## Catálogo de actores


## Catálogo de publicações

## Informação interessante

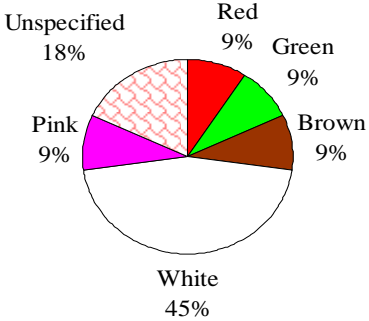
## Fórum

Corpus	Lemas								
	N	ADJ	ADV	V	NUM	GRAM	PROP	todos	todos/pos
AmostRA	<a href="#">5088</a>	<a href="#">1903</a>	<a href="#">324</a>	<a href="#">2351</a>	<a href="#">323</a>	<a href="#">184</a>	<a href="#">1387</a>	<a href="#">11197</a>	<a href="#">11560</a>
ANCIB	<a href="#">14957</a>	<a href="#">4116</a>	<a href="#">702</a>	<a href="#">4201</a>	<a href="#">4964</a>	<a href="#">356</a>	<a href="#">29717</a>	<a href="#">57993</a>	<a href="#">59013</a>
Avante!	<a href="#">20512</a>	<a href="#">8996</a>	<a href="#">1707</a>	<a href="#">9419</a>	<a href="#">6264</a>	<a href="#">851</a>	<a href="#">47500</a>	<a href="#">93014</a>	<a href="#">95249</a>
CDHAREM	<a href="#">5432</a>	<a href="#">1970</a>	<a href="#">391</a>	<a href="#">2237</a>	<a href="#">685</a>	<a href="#">293</a>	<a href="#">3915</a>	<a href="#">14453</a>	<a href="#">14923</a>
CETEMPúblico	<a href="#">160824</a>	<a href="#">47512</a>	<a href="#">5475</a>	<a href="#">54926</a>	<a href="#">109796</a>	<a href="#">2047</a>	<a href="#">1070220</a>	<a href="#">1430678</a>	<a href="#">1450800</a>
CETEMPúblico (primeiro milhão)	<a href="#">13510</a>	<a href="#">4765</a>	<a href="#">845</a>	<a href="#">5437</a>	<a href="#">2389</a>	<a href="#">284</a>	<a href="#">23758</a>	<a href="#">49517</a>	<a href="#">50988</a>
CHAVE	<a href="#">113806</a>	<a href="#">39865</a>	<a href="#">4715</a>	<a href="#">40533</a>	<a href="#">91180</a>	<a href="#">1383</a>	<a href="#">711584</a>	<a href="#">988949</a>	<a href="#">1003066</a>
Clássicos da Literatura Portuguesa/Porto Editora	<a href="#">12818</a>	<a href="#">5109</a>	<a href="#">1116</a>	<a href="#">8961</a>	<a href="#">267</a>	<a href="#">355</a>	<a href="#">4393</a>	<a href="#">31726</a>	<a href="#">33019</a>
ConDIVport	<a href="#">15162</a>	<a href="#">8009</a>	<a href="#">1507</a>	<a href="#">11619</a>	<a href="#">2184</a>	<a href="#">580</a>	<a href="#">31571</a>	<a href="#">60790</a>	<a href="#">63123</a>
ConE	<a href="#">10163</a>	<a href="#">2727</a>	<a href="#">435</a>	<a href="#">2744</a>	<a href="#">4295</a>	<a href="#">324</a>	<a href="#">18504</a>	<a href="#">38527</a>	<a href="#">39192</a>
DiaCLAV	<a href="#">20854</a>	<a href="#">6809</a>	<a href="#">1174</a>	<a href="#">8924</a>	<a href="#">5754</a>	<a href="#">416</a>	<a href="#">64786</a>	<a href="#">105800</a>	<a href="#">108717</a>
ECI-EBR	<a href="#">13909</a>	<a href="#">5773</a>	<a href="#">936</a>	<a href="#">6192</a>	<a href="#">925</a>	<a href="#">353</a>	<a href="#">8987</a>	<a href="#">35815</a>	<a href="#">37075</a>
ECI-EE	<a href="#">1043</a>	<a href="#">515</a>	<a href="#">183</a>	<a href="#">575</a>	<a href="#">226</a>	<a href="#">126</a>	<a href="#">173</a>	<a href="#">2727</a>	<a href="#">2841</a>
ENPC (parte pública)	<a href="#">3555</a>	<a href="#">1375</a>	<a href="#">364</a>	<a href="#">1881</a>	<a href="#">138</a>	<a href="#">183</a>	<a href="#">793</a>	<a href="#">8023</a>	<a href="#">8289</a>
FrasesPB	<a href="#">2156</a>	<a href="#">749</a>	<a href="#">187</a>	<a href="#">888</a>	<a href="#">60</a>	<a href="#">129</a>	<a href="#">215</a>	<a href="#">4255</a>	<a href="#">4384</a>
FrasesPP	<a href="#">1698</a>	<a href="#">689</a>	<a href="#">183</a>	<a href="#">774</a>	<a href="#">70</a>	<a href="#">131</a>	<a href="#">197</a>	<a href="#">3640</a>	<a href="#">3742</a>
Museu da Pessoa	<a href="#">5221</a>	<a href="#">1378</a>	<a href="#">320</a>	<a href="#">2641</a>	<a href="#">353</a>	<a href="#">227</a>	<a href="#">2106</a>	<a href="#">11808</a>	<a href="#">12246</a>
Natura/Mínho	<a href="#">12829</a>	<a href="#">5339</a>	<a href="#">851</a>	<a href="#">5199</a>	<a href="#">4431</a>	<a href="#">425</a>	<a href="#">30354</a>	<a href="#">58141</a>	<a href="#">59428</a>
Natura/Público	<a href="#">35717</a>	<a href="#">12121</a>	<a href="#">1534</a>	<a href="#">12170</a>	<a href="#">9723</a>	<a href="#">837</a>	<a href="#">83606</a>	<a href="#">152574</a>	<a href="#">155708</a>
NILC/São Carlos	<a href="#">64650</a>	<a href="#">22874</a>	<a href="#">2769</a>	<a href="#">25444</a>	<a href="#">60630</a>	<a href="#">815</a>	<a href="#">302016</a>	<a href="#">472123</a>	<a href="#">479198</a>
Vercial	<a href="#">35918</a>	<a href="#">12240</a>	<a href="#">2629</a>	<a href="#">27351</a>	<a href="#">2710</a>	<a href="#">628</a>	<a href="#">42890</a>	<a href="#">116376</a>	<a href="#">124366</a>


# COMPARA: (EN) Author with highest % of colour:



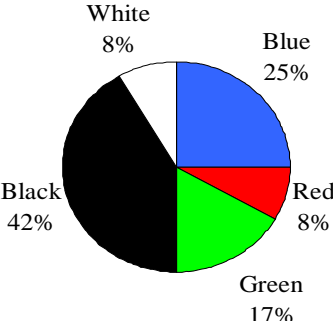
**Lewis Carrol**




Color	Percentage
White	45%
Unspecified	18%
Pink	9%
Red	9%
Green	9%
Brown	9%



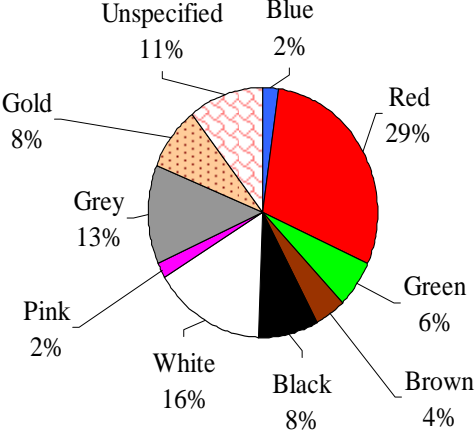
**Mary Shelley**




Color	Percentage
Black	42%
Blue	25%
Green	17%
Red	8%
White	8%



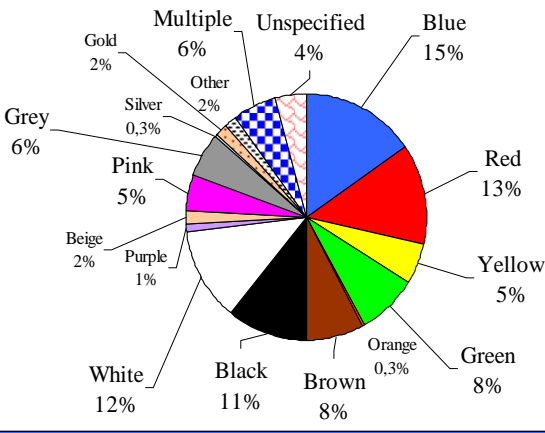
**Henry James**



Color	Percentage
Red	29%
Unspecified	11%
Grey	13%
White	16%
Black	8%
Green	6%
Brown	4%
Gold	8%
Blue	2%
Pink	2%




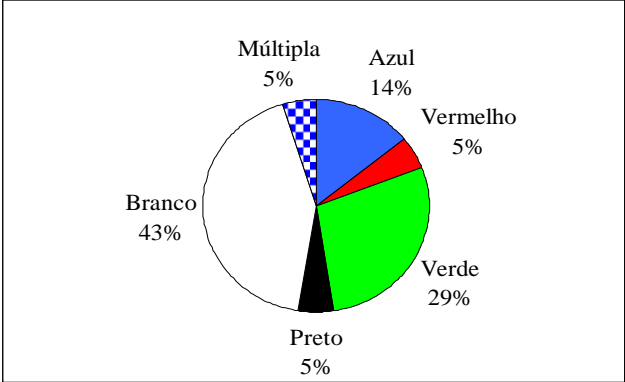
**Joanna Trollope**



Color	Percentage
Blue	15%
Red	13%
Yellow	5%
Green	8%
Orange	0.3%
Brown	8%
Black	11%
White	12%
Grey	6%
Purple	1%
Beige	2%
Pink	5%
Silver	0.3%
Gold	2%
Multiple	6%
Unspecified	4%

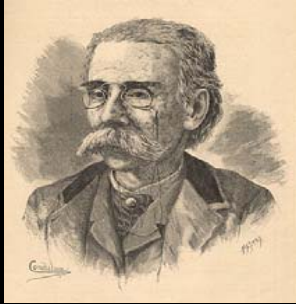
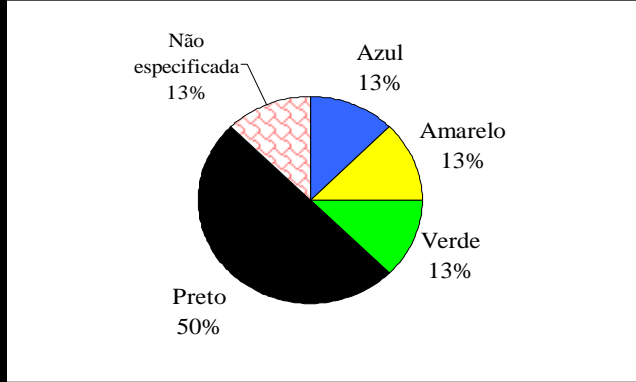
Silva, Inácio & Santos (2008)

# COMPARA: (PT) author with highest % of colour:

**José de Alencar**

Color	Percentage
Branco	43%
Verde	29%
Preto	5%
Azul	14%
Vermelho	5%
Múltipla	5%


**Camilo Castelo Branco**

Color	Percentage
Preto	50%
Não especificada	13%
Verde	13%
Amarelo	13%
Azul	13%



**Mia Couto**

**26%**




**José Eduardo Agualusa**

**31%**



**Jorge de Sena**

**24%**



**Marcos Rey**

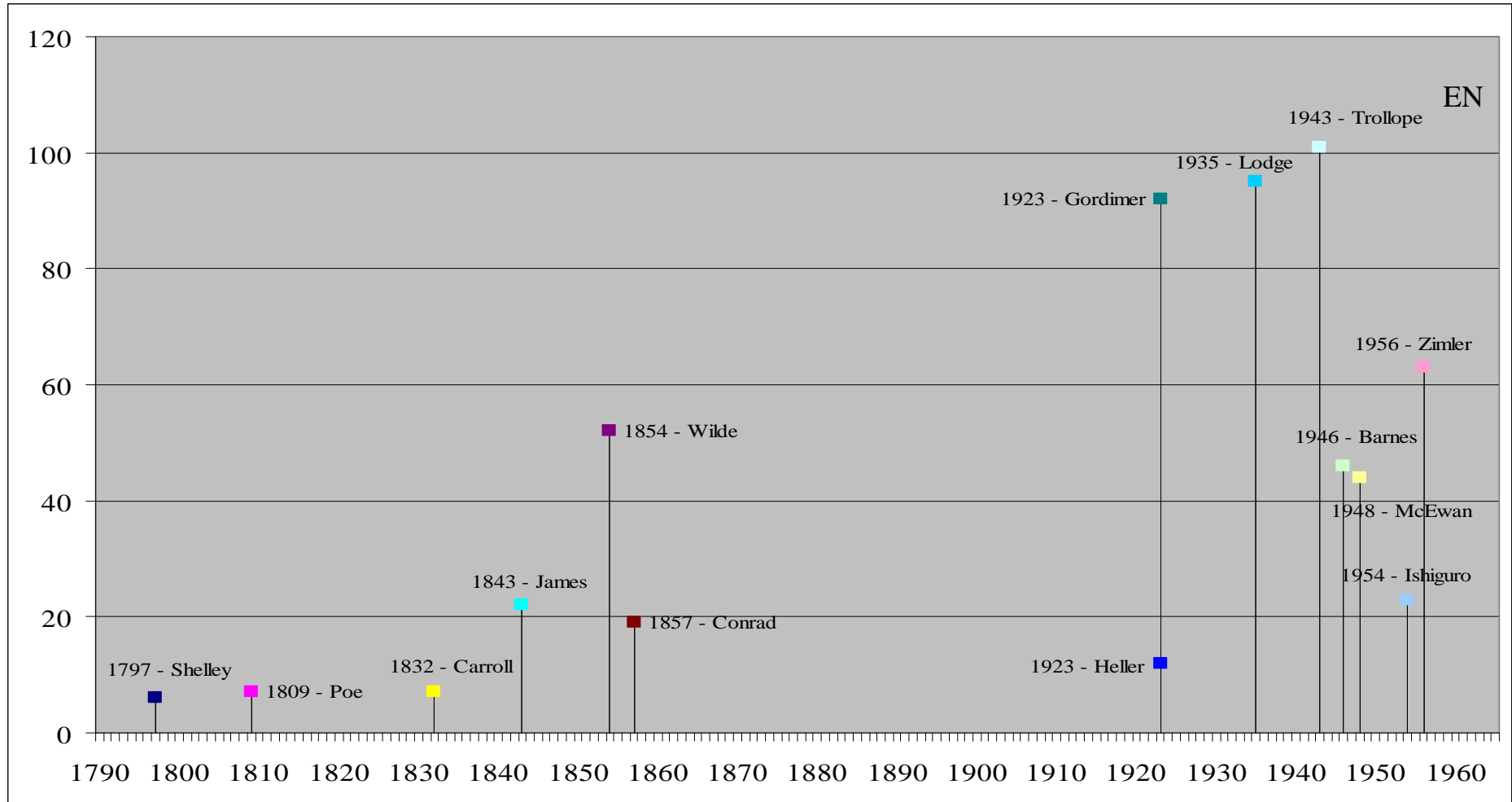
**44%**

Silva, Inácio & Santos (2008)



# COMPARA: Does colour quantity change with time?

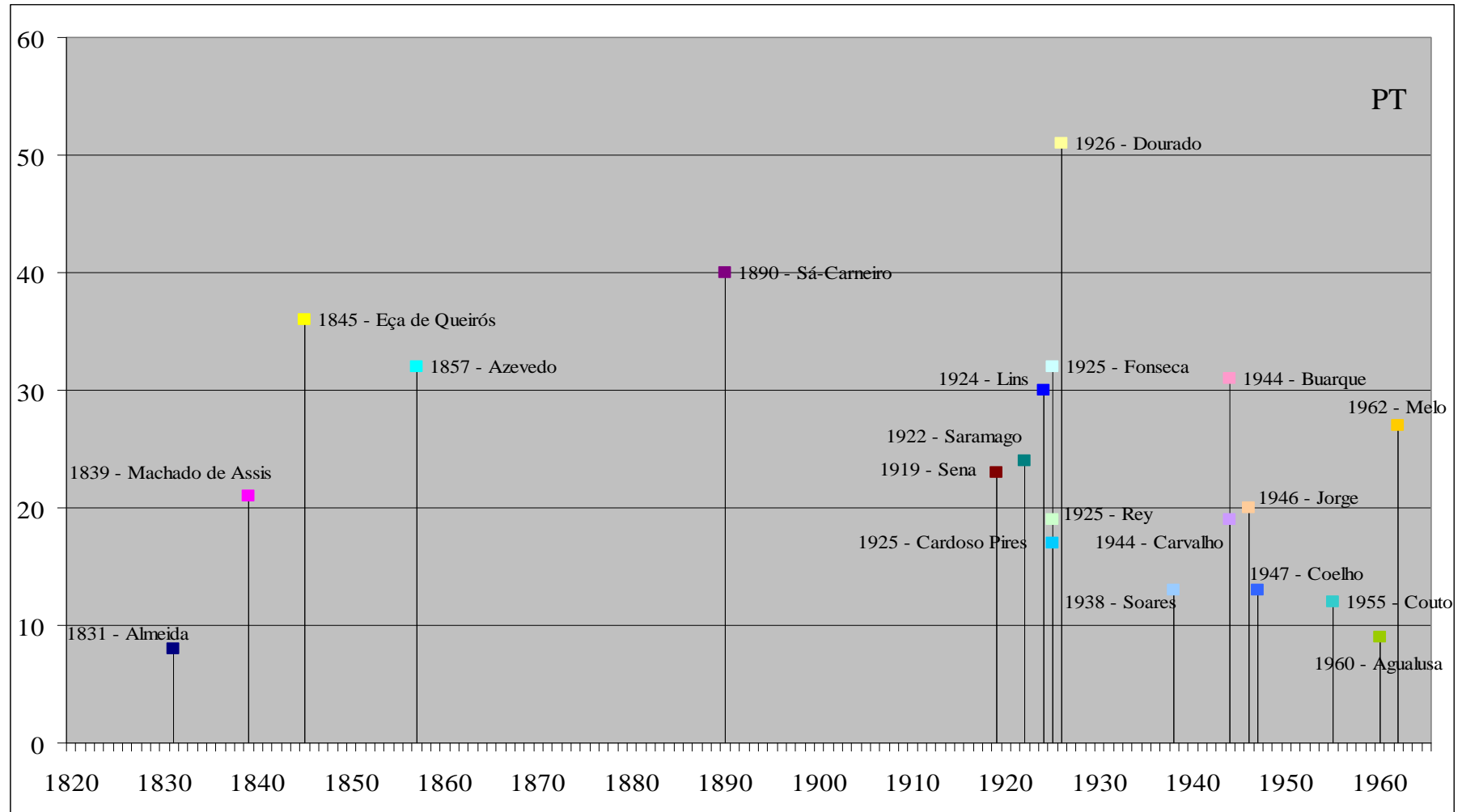
Number of colour types per authors' birth date (English-speaking authors)



Silva, Inácio & Santos (2008)

# COMPARA: Does colour quantity change with time?

Number of colour types per authors' birth date (Portuguese-speaking authors)



Silva, Inácio & Santos (2008)

# The Floresta Sintáctica treebank: history

- **2000** Formal cooperation between VISL and Linguateca started
- **2000-2001** Root planting: three linguistic scholarships at Odense, active preparation of tools and workflow at Oslo-Odense, launching of the basic resource and project philosophy (3-4 years work)
- **2002-2004** Partial work, incremental versions, stable but slow work
- **2005** Work on format validation at the Braga node

Sleeping forest...

- **2007-2008** New team, at Coimbra node: new material, new tools

Sleeping forest...

There is support and answer to questions, but not actual development

# Floresta's “international” impact

- Used by Sabine Buchholz & Darren Green at their LREC 2006 article to illustrate treebanks' maintenance problems
- Used by Jason Balridge to infer a Portuguese grammar
- Used by CoNLL-X 2006, *ConLL-X shared task on multilingual dependency parsing* for Portuguese
- Integrated by Steven Bird in NLTK, *Natural Language Toolkit*, since September 2007
- Other explicit uses
  - John Hopkins University
  - Essex University
- Floresta anonymous downloads: (from our logs)

# Some numbers from 2004

Clauses	21,931
Finite clauses	15,566
Infinite clauses	5,602
Averbal clauses	763
Noun phrases*	43,096
Prepositional phrases*	32,210
Adjectival phrases*	1,780
Adverbial phrases*	833
Coordinated items	5,448
Trees	9,431

\* phrase = *more than one word*

# Lessons from Floresta

- Very hard to gather a community: much easier to create own treebanks
- Very hard to have any feedback from theoretical linguists
  - They were not invited in the first place!
  - Who is this German guy anyway?
  - Why do engineers talk about syntax?
- Very hard to have consensus on any subject whatsoever of linguistics:
  - What is a word? What is a phrase? What is a multiword expression?
  - What is an argument? What is an object? What is a phrase?
  - What is a head? What is a noun? What is a noun phrase?
- Ahead of our time? Impossible aspiration? Users will come in later?
  - Merge with current (independent) projects?

# Brief presentation of CorTrad

## ■ New material

- Portuguese-to-English translation
- Non-native translation vs. Native translation
- Technical translation

## ■ Multiversion

- Study the changes from initial to published translation
- Varieties might also be construed as rough translation to more idiomatic one

## ■ Tailored for specific genres

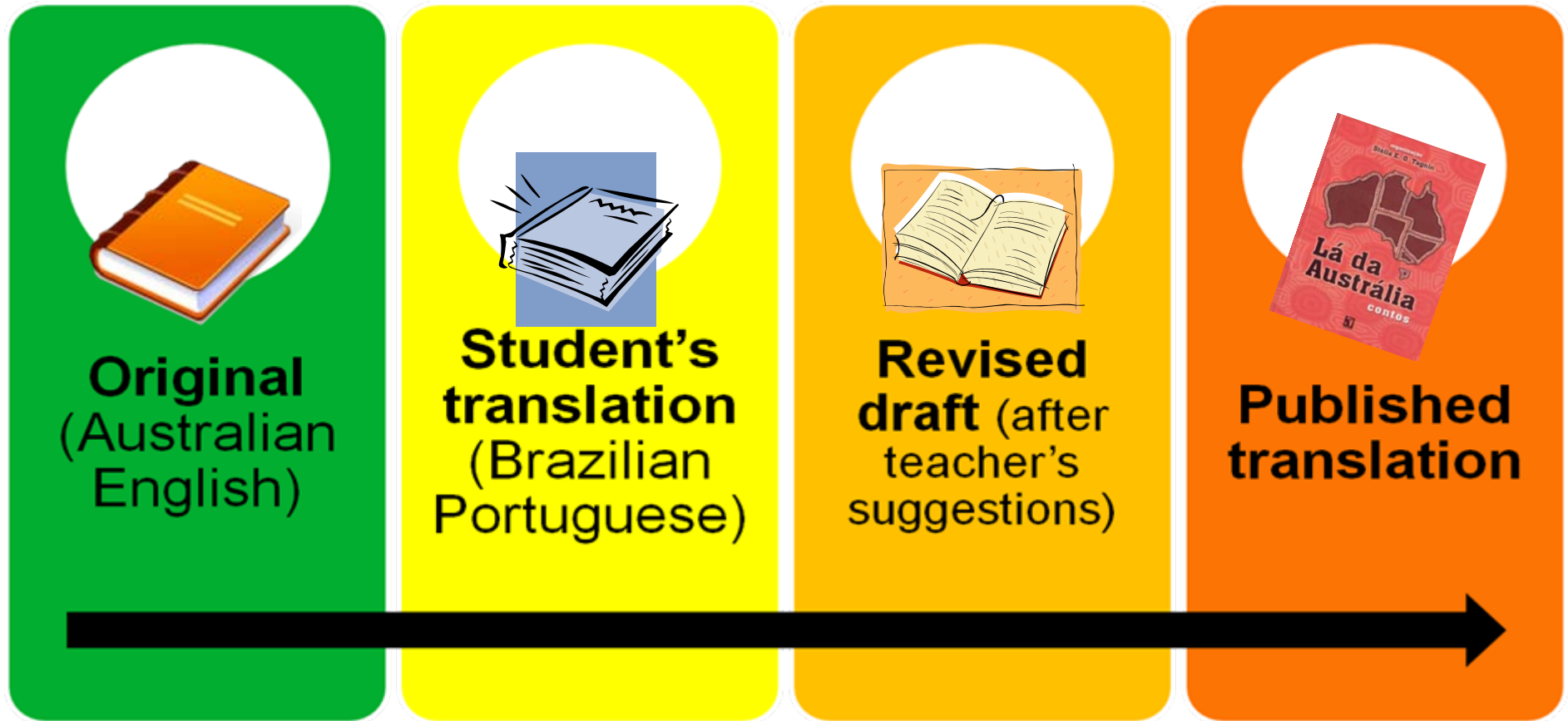
- Cookbook
- Short scientific news

CorTrad is the parallel subcorpus of COMET, encoded with DISPARA, Linguateca's environment to make corpora available on the Web. A joint project of Univ. of São Paulo, Linguateca and NILC.

# Literary CorTrad:

## Australian short stories

(\*learner corpus)



Tagnin, Santos & Teixeira (2009)





- Principal
- O Projeto
- Equipe
- CorTec
- CoMAprend
- CorTrad
- Artigos, etc.
- Links
- Informativo
- Contato
- Site FFLCH
- Site USP

I have a new <b>house</b> .	Tenho uma casa nova .	Tenho uma casa nova .	Tenho uma casa nova .
I lock the <b>house</b> and hide the key and walk up the street , past neat little gardens of lawns and roses .	Tranco a porta , escondo a chave e subo a rua passando por lindos jardins gramados com roseiras .	É meio dia e meia , hora de ir . <s> Tranco a porta , escondo a chave e subo a rua passando por lindos jardins gramados com roseiras . Passo por casas desabitadas e por desertas ruas de bairro .	É meio dia e meia , hora de ir . <s> Tranco a porta , escondo a chave e subo a rua passando por lindos jardins gramados com roseiras . Passo por casas desabitadas e por desertas ruas de bairro .
The <b>house</b> where Lady Weare was to stay for two nights of her visit was old and had a friendly garden .	A casa onde Lady Weare ficaria por duas noites em sua visita a Melbourne era antiga e tinha um belo jardim .	A casa onde Lady Weare ficaria por duas noites em sua visita a Melbourne era antiga e tinha um belo jardim .	A casa onde Lady Weare se hospedaria durante as duas noites em Melbourne era antiga e tinha um jardim acolhedor .
The charming <b>house</b> , her friends , welcomed her .	Na casa charmosa , seus amigos lhe deram boas-vindas .	Na casa charmosa , seus amigos lhe deram boas-vindas .	De volta à casa encantadora , os amigos lhe deram boas vindas .
At times , because of her own unhappiness , she had the gift of making everyone unhappy.Mr Berrington lived alone in a <b>house</b> which was bigger and better than those in the surrounding streets and he died alone .	Às vezes , devido à sua própria infelicidade , possuía o dom de deixar todos infelizes.O Sr.Berrington vivia sozinho em uma casa que era maior e melhor do que as outras na vizinhança e morreu sozinho .	Às vezes , devido à sua própria infelicidade , possuía o dom de deixar todos infelizes.O Sr.Berrington vivia sozinho em uma casa que era maior e melhor do que as outras na vizinhança e morreu sozinho .	Às vezes , em função da sua própria infelicidade , possuía o dom de deixar a todos infelizes . O Sr.Berrington vivia sozinho em uma casa que era maior e melhor do que as outras na vizinhança e morreu sozinho .
He had been dead for about a week when the police , responding to a call from a neighbour , broke into his <b>house</b> .	Estava morto há uma semana quando a polícia , atendendo ao chamado de uma vizinha , arrombou a casa .	Estava morto há uma semana quando a polícia , atendendo ao chamado de uma vizinha , arrombou a casa .	Estava morto há uma semana quando a polícia , atendendo ao chamado de uma vizinha , arrombou a casa .
She had not attempted to visit him when he had not come to the <b>house</b> as usual .	Ela não havia tentado ir visitá-lo , quando ele não aparecera na sua casa , como de costume .	Ela não havia tentado ir visitá-lo , quando ele não aparecera na sua casa , como de costume . Teria morrido rapidamente , durante o sono ?	Ela não tentara ir visitá-lo quando ele não aparecera na sua casa , como de costume .

Tagnin, Santos & Teixeira (2009)

# The detailed study of language varieties... with the AC/DC cluster?

- Three moments: what is the material and how is it marked up?
  - Variety (**country**, province, social class, age, ...)
  - Time of publication (decade, year, semester, day, ...), time of writing
  - Genre, register, publication channel, author, ...
  - Original/translated (from...)/transcribed
  - Revised at all?
  - Coherent or discontinuous?
- How comparable it is? How do intra-variety and inter-variety correlate?
  - Corpus homogeneity, corpus signature, or maximum quantity as the ideal good?

# Support for formal variational linguistics

- Inspired by the Quantitative Lexicology and Variational Linguistics group <http://wwwling.arts.kuleuven.be/qlvl/> at the Catholic University at Leuven, and its Portuguese counterpart, CONDIVport, created by Augusto Soares da Silva and his team at the Catholic Univ. of Braga, and made available through AC/DC, we started to provide support for this kind of studies as a merge with our semantic annotation efforts
- CONDIVport developed a set of onomasiological profiles for the themes of football and fashion (health is underway)
- Linguateca did the same for colour, and revised annotation in context
- Both fashion and colour profiles were reused and improved and all AC/DC corpora were automatically annotated with them

# Profiling...

- Profile names (fashion): *blusa* or *blusão* or *calças curtas*
- Profile names (colours): *vermelho* or *branco* or *creme*
- *blusão*: *blazer, blusão, camurça, casaco de pele, colete, etc.*
- *calças curtas*: *bermudas, calças à corsário, calças ¾, calções, shorts, etc.*
- *vermelho*: *cor de carmim, cor de cereja, cor de chama, cor de colorau, cor de fogo alaranjado, cor de lagosta, cor de lagosta de viveiro, cor de morango, cor de morango esborrachado, encarniçado, escarlate, grená, magenta, ruborizar-se, rubro, vermelho-Benfica, vermelho-bordeaux, etc.*
- *creme*: *aperolada, bege, bege África, bege-areia, marfim, cor de pele, etc.*

# Comparing profile-based measures (Geeraerts & Grondelaers, 99)

$$A_{K,Z}(Y) = \sum_{i=1}^n F_{Z,Y}(x_i) \cdot W_{x_i}$$

- $A_{K,Z}(Y)$  is the ratio of terms with a feature K in the onomasiological profile for concept Z in dataset Y
- K = set of terms with a particular feature (for example FRENCH)
- Z = concept (for example VERDE, or VEST, or BLUSÃO)
- $F_{Z,Y}$  relative frequency of x for concept Z in Y

$$A_K(Y) = 1/n * \sum_{i=1}^n A_{K,Z_i}(Y)$$

- $A_K(Y)$  is the global proportion of the subset K in dataset Y
- Comparing values of **relevant features** for different “datasets” (decades, varieties) convergence or divergence can be investigated

# Current data about AC/DC profiles

- Number of different colour terms (lemmas) in the set of all corpora: **1672**, in 23 profiles (colour groups)  
(not counting proper names or other cases deriving from colour)
- Number of different clothing terms in the set of all corpora: **318**, in 28 profiles

	Colour tokens	Colour types	CT per 10 <sup>4</sup>
Condiv	20,380	547	47.5
CHAVE	85,506	526	5.47
CLASSLPPE	3,214	145	49.9
museudapessoa	167	24	20.0

# Future capabilities in AC/DC

- Search helped/enhanced by semantic relations (synonyms, hypernyms, antonyms...) and syntactic relations (adj-past participle, clitics, contractions, nominalizations, ...)
- Reuse of complex queries
  - previous queries listed and explained (documentation effort, enriched FAQ)
  - programming of macros to make them more compact
- Automatic contrast between two different searches
  - Parallel results
  - Similarities
  - Contrasts

# A preview of what may come...

- Compare X with Y according to...
  - Pure frequency
  - Distribution of lemmas
  - Distribution of passive/active, tense...
  - Presence of postmodifiers
  - Kind of subject, or of verb...
  - ...
- X and Y can be lexical items, or constructions, or any search whatsoever, restricted to whatever subsets
  - Compare ADJ N with N ADJ
  - Compare SEE pron VGER with SEE THAT finite
  - Compare *castanho* (PT) with *marrom* (BR), or in translation vs. original



# Feedback highly appreciated!

We are open to

- Comments
- Questions
- Doubts
- Suggestions
- Ideas
- Critical remarks
- Cooperation proposals