



Natural Language Processing: Dealing with unstructured data

Renata Vieira and Joaquim Neto

FGV Workshop on Data Science
2019





E mudamos nossa forma de pensar, de falar e de
comunicar

Roteiro – dia 1

- Parte 1
 - Complexidade da língua natural
 - Conhecimento linguístico
 - Conhecimento de mundo
 - Tarefas de PLN
 - Língua Portuguesa
- Parte 2
 - Modelos estatísticos da língua
 - Word Embeddings

Roteiro – dia 2

- Parte 1
 - Relações entre ontologias e PLN
 - Pesquisas desenvolvidas nessas áreas
- Parte 2
 - Prática
 - Tarefa: Reconhecimento de entidades nomeadas

Parte 1

- Complexidade da língua natural
- Conhecimento linguístico
- Conhecimento de mundo
- Tarefas de PLN
- Língua Portuguesa

Língua natural



Língua natural





Localizar: puc

A língua natural é ambígua

- Diferentemente das linguagens artificiais
 - Java
 - C
 - HTML
 - Lógica

Comunicação humana e ambiguidade

- Mas como a comunicação humana é eficiente?
 - (pelo menos, às vezes)
- Associações
 - Fala/texto e visão
 - Fala/texto e conhecimento prévio
 - Sistema altamente contextualizado
 - Dito e o não dito



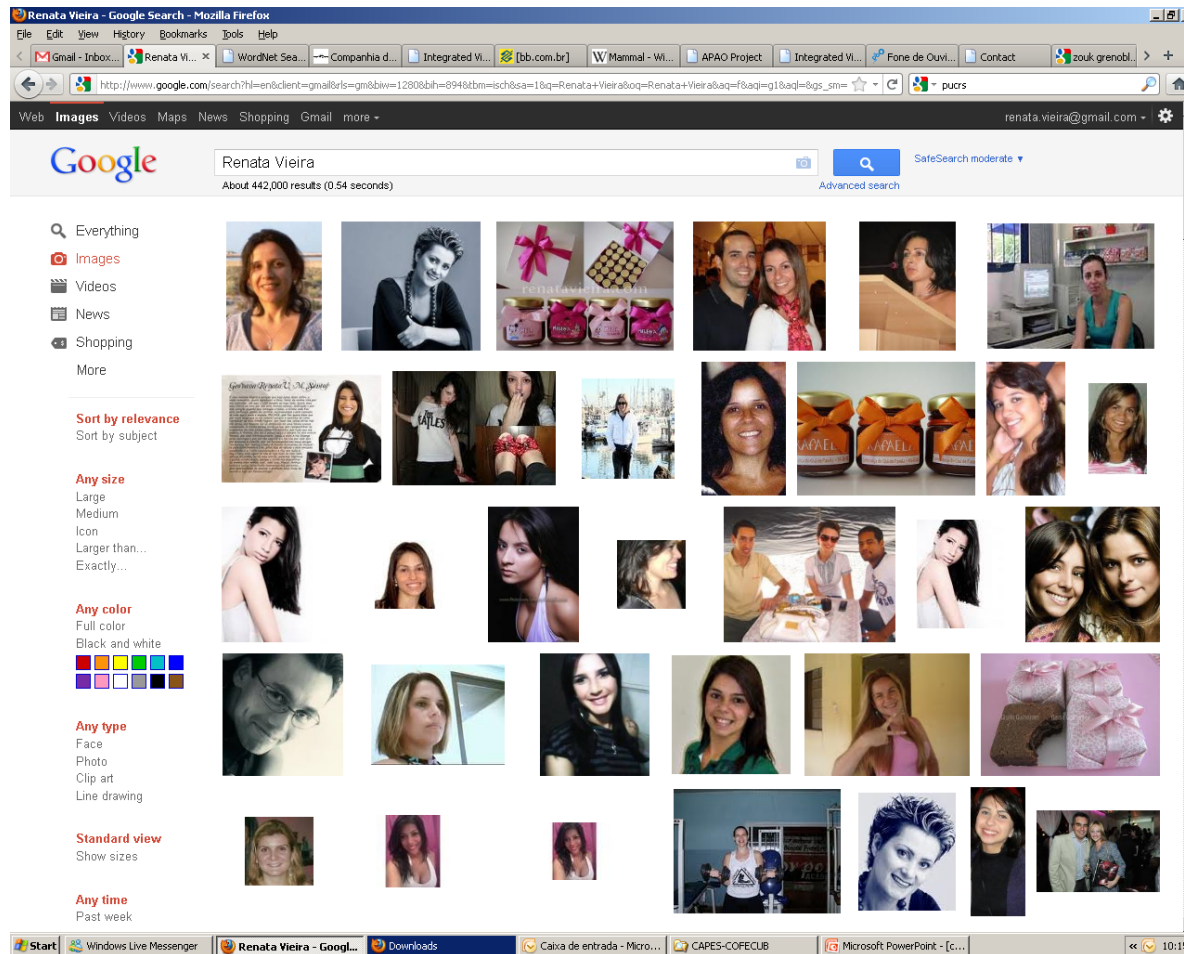
Vamos por partes:
nomes, termos, conceitos,
significados e outras palavras...



Nomes

- Nomes e seus referentes
 - Entidades nomeadas
 - Brasil
 - o país
 - o território
 - time de futebol
 - Renata Vieira
 - pessoa

Renata Vieira



Nomes, **termos**

- Palavras simples e compostas geralmente usadas em contextos específicos
- Termos técnicos usados por exemplo em ciências ou artes específicas
- Terminologia: disciplina lingüística para o estudo científico dos **conceitos** e **termos** usados nas línguas de especialidade”.

Nomes, termos, **conceitos**

- Conceito
 - compreensão que alguém tem de uma palavra; noção, concepção, ideia
 - O termo "**conceito**" tem origem a partir do latim "conceptus" (do verbo concipere) que significa "coisa concebida" ou "**formada na mente**"

Nomes, termos, conceitos, **significados**

- **Significado de Significado**

- Definição atribuída a um termo, palavra, frase, texto; aquilo que alguma coisa quer dizer; sentido.
- Relevância que se dá a algo: sua participação teve muito significado

Nomes, termos, conceitos, **significados**

- **Significado de Significado**
 - [Linguística] Significação; forma representativa e mental que se relaciona com a forma linguística; o que o signo quer significar; a parte do signo linguístico definida pelo conceito.

Nomes, termos, conceitos, significados e **outras palavras**

- Verbos, adjetivos, preposições ...
- **Artigos**

O artigo definido “the”

- B Russell – Mind, 1905.
 - On denoting
- R Vieira - PhD Thesis, 1998
 - Definite Description Resolution in Unrestricted Texts

O artigo definido “the”

- B Abbott - Proceedings of SALT, 2011
 - Support for a Unique Theory of Definiteness
 - There is undoubtedly more to be said on this subject, which seems to continue to excite at least some people today as much as it did Russell at the turn of the century.
- J Cho, International Journal of Bilingualism 2017
 - The acquisition of different types of definite noun phrases in L2-English

Em outras palavras ...

- Somos naturalmente aptos pra lidar com palavras, por que ao estudá-las elas começam a parecer tão complicadas?

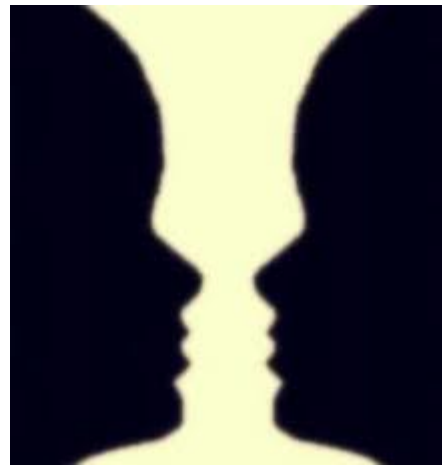
O Sistema Simbólico

- Linguagem:
 - Convenções, associações
 - Nomes próprios
 - Renata Vieira
 - Designadores descritivos
 - A palestrante da tarde de hoje (categorias e propriedades)

Mente

- Conceito
 - "coisa concebida" ou "formada na mente"

A mente e a cunhagem simbólica

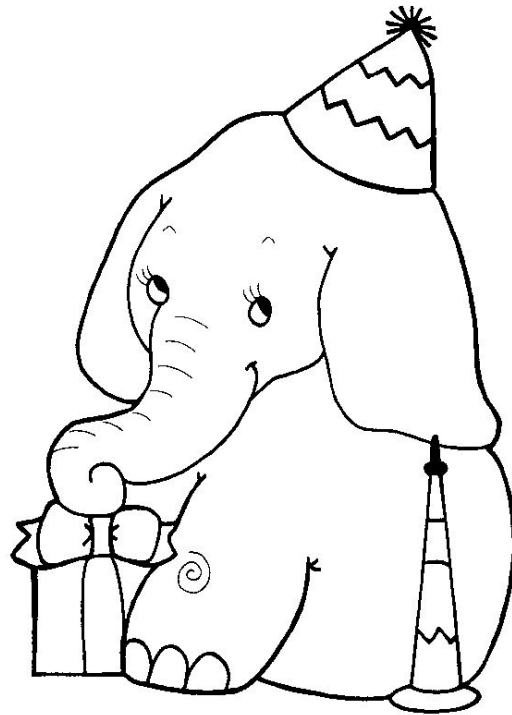


A mente e a cunhagem simbólica

- El____ant

A mente e a cunhagem simbólica

- El____ant



A mente dos outros

- O sucesso da comunicação está na habilidade de reconhecer a intenção comunicativa do falante e inferir sua intenção informativa

Palavras

- Apenas algumas razões das dificuldades
- Outras razões:
 - Mentes
 - Frases, Textos, Posts
 - Sintaxe, semântica e pragmática

Além das palavras

- Frases
 - Parsing ou análise sintática
- Identificar os constituintes de uma frase, verificando as suas partes de acordo com a gramática da língua
 - Sujeito verbo objeto
 - A língua é o único músculo voluntário do corpo humano que não fadiga.

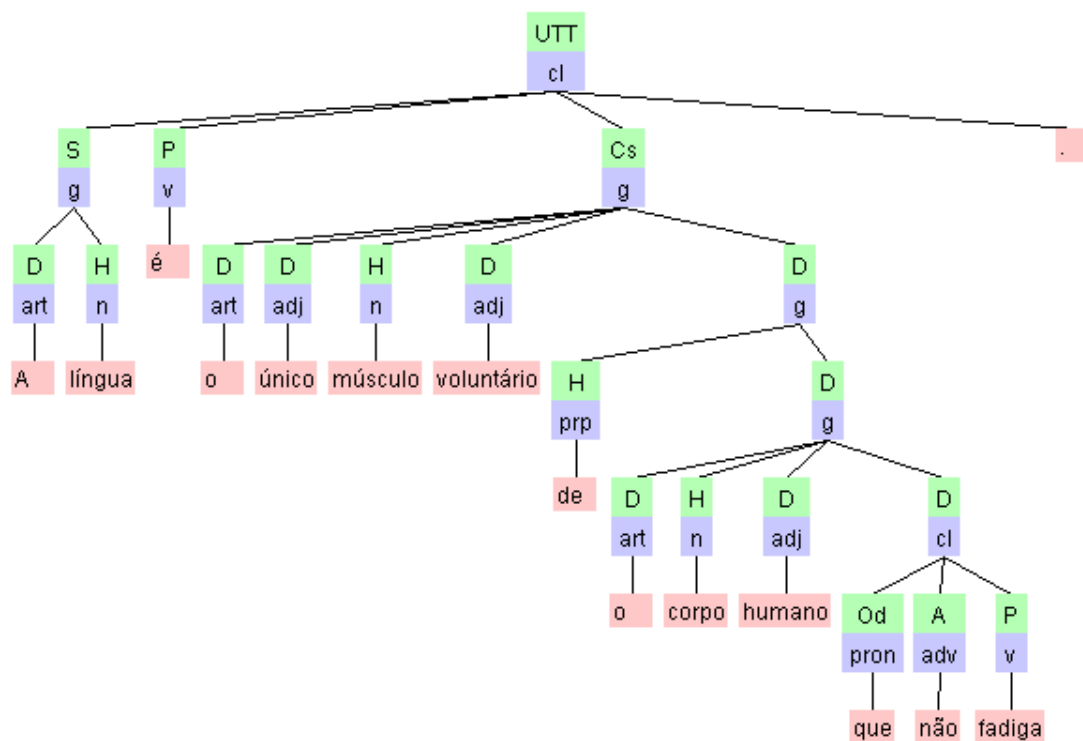
Sentence: A língua é o único músculo voluntário de o corpo humano que não fadiga .

Function: S P Od Oi Op Cs Co fC A SUB CO CJT D H UTT STA QUE COM EXC X

Form: n prop v adj adv art pron prp conj num intj cl par g x

Collapse Tree

Expand Tree



Analysis 1 of 1

Java Applet Window

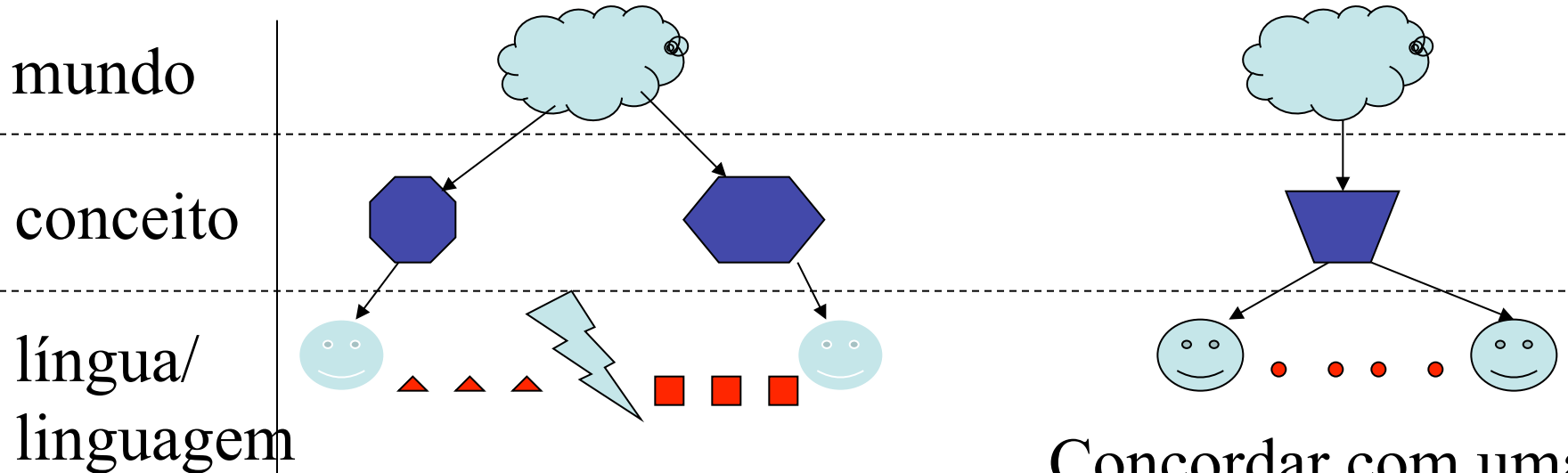
Sintaxe, Semântica, Pragmática

- Forma
- Conteúdo/significado
- Uso
- Conhecimento de mundo

Semântica

- Significado
 - Thesauros
 - Dicionários
 - Ontologias

Ontologias



Desentendimento

Confusão conceitual e
terminológica

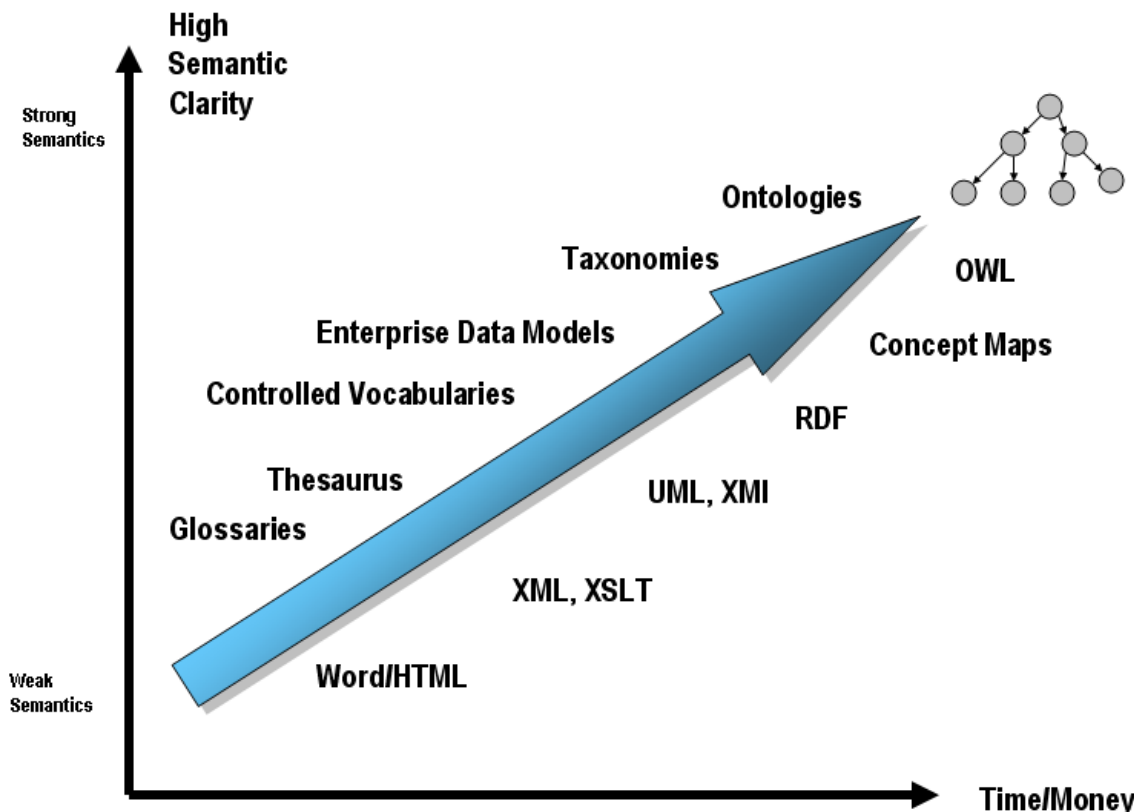
Concordar com uma
conceitualização

Torná-la explícita
em alguma linguagem

Atores: humanos e máquinas

Ontologias

- Semântica explícita
- Significado estruturado
- Estrutura lógica



Ontologias

- Conceito (classe)
 - Extensão x Intensão
 - Conjunto de referentes x conjunto de propriedades
 - Hierarquia de conceitos
 - Relacionamentos

Ontologias

- Relacionamentos
 - É_um (hierarquia)
 - Parte de
 - Same_as

Ontologia da sala

- Conceitos
- cadeira, mesa, porta, auditório, lustres, janelas, paredes
- Hierarquia
- Auditório é um tipo de sala
- Mesa é um tipo de móvel

Ontologia da sala

- Conceitos
 - cadeira, mesa, porta, auditório, lustres, janelas, paredes, teto e piso
- Relacionamentos
 - janelas situam-se nas paredes
 - cadeiras ficam sobre o piso
- Axiomas
 - cadeiras tem 4 pés

Ontologias de topo

- Objetos físicos
- Objetos abstratos
- Continuantes
- Perdurantes

Outros recursos semânticos

- Bases semânticas
 - WordNet
 - ConceptNet
 - BabelNet

WordNet Search - 3.0 - Mozilla Firefox

Arquivo Editar Exibir Histórico Favoritos Ferramentas Ajuda

http://wordnet.princeton.edu/perl/webwn?s=language&sub=Search+Word wordviz

powered by YAHOO! SEARCH Search Web

Nondete... Springer... AAMAS 2008 the Asso... propor2008... Word... Powerse... Microsoft K... wordviz ... Visuwor...

WordNet Search - 3.0 - [WordNet home page](#) - [Glossary](#) - [Help](#)

Word to search for:

Display Options:

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations

Noun

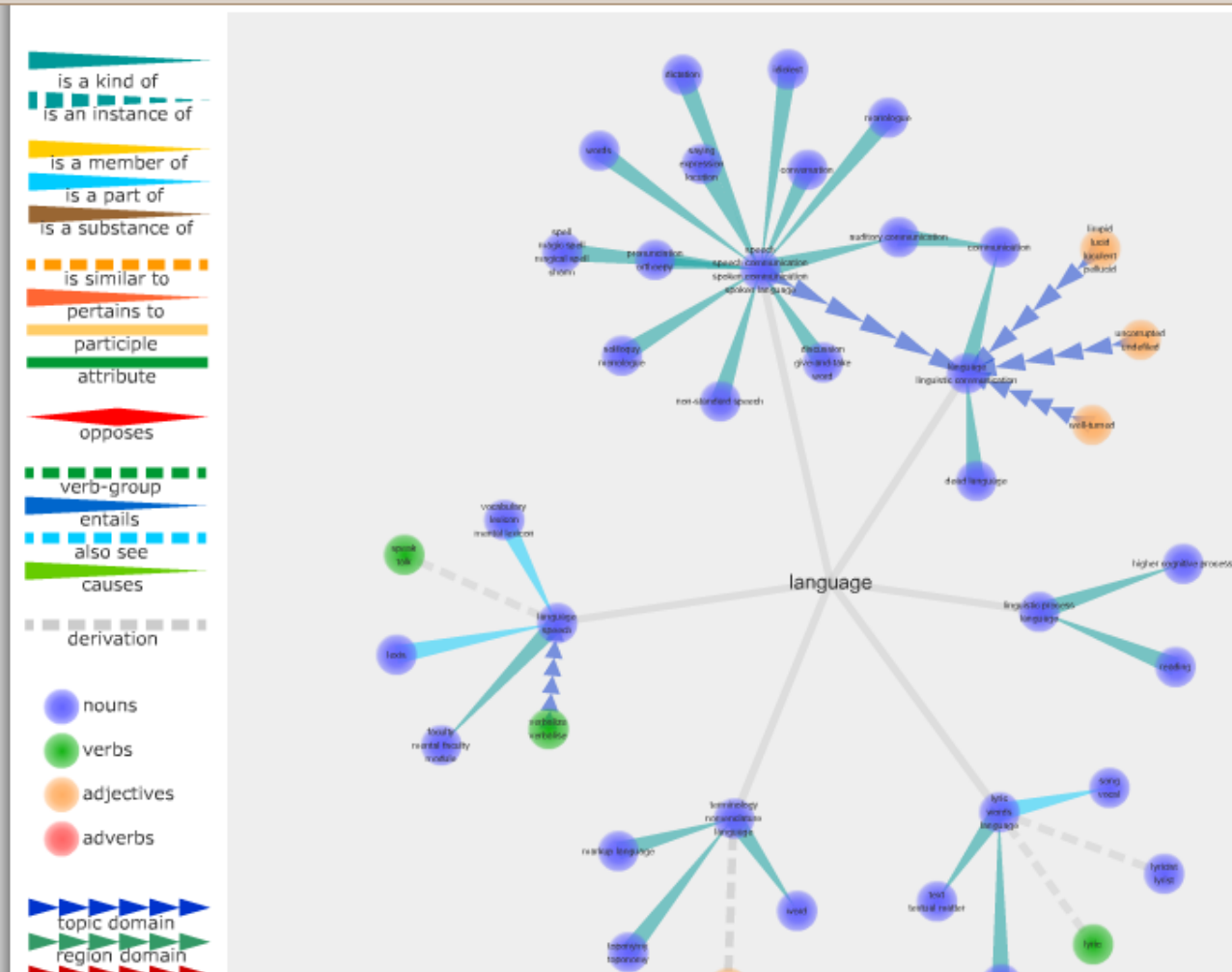
- **S: (n) language, linguistic communication** (a systematic means of communicating by the use of sounds or conventional symbols) *"he taught foreign languages"; "the language introduced is standard throughout the text"; "the speed with which a program can be executed depends on the language in which it is written"*
- **S: (n) speech, speech communication, spoken communication, spoken language, language, voice communication, oral communication** ((language) communication by word of mouth) *"his speech was garbled"; "he uttered harsh language"; "he recorded the spoken language of the streets"*
- **S: (n) lyric, words, language** (the text of a popular song or musical-comedy number) *"his compositions always started with the lyrics"; "he wrote both words and music"; "the song uses colloquial language"*
- **S: (n) linguistic process, language** (the cognitive processes involved in producing and understanding linguistic communication) *"he didn't have the language to express his feelings"*
- **S: (n) language, speech** (the mental faculty or power of vocal communication) *"language sets homo sapiens apart from all other animals"*
- **S: (n) terminology, nomenclature, language** (a system of words used to name things in a particular discipline) *"legal terminology"; "biological nomenclature"; "the language of sociology"*

[WordNet home page](#)

Localizar: ☐ Diferenciar maiúsc./minúsc.

Concluído 5.172s ☒ Images

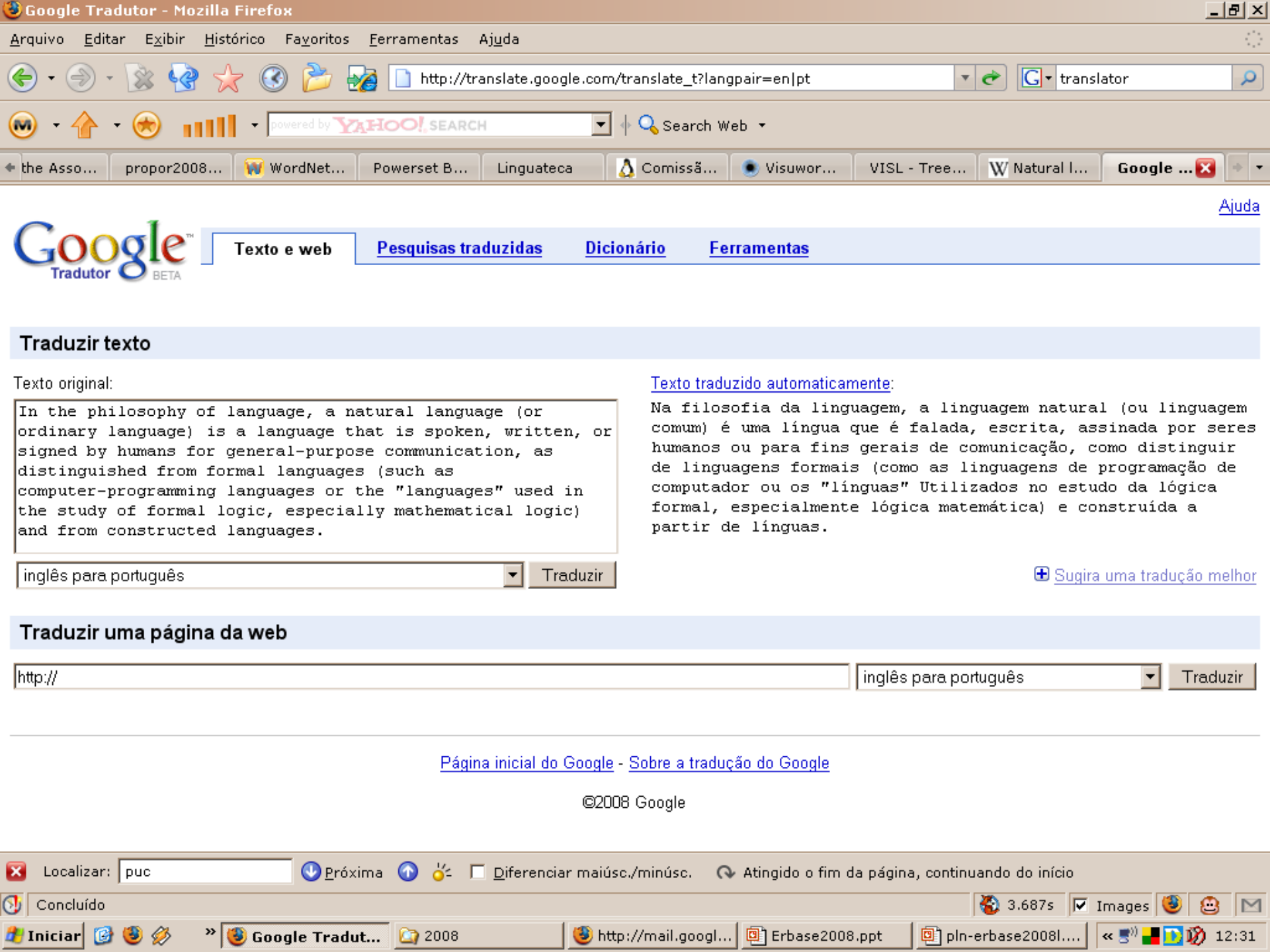
Iniciar WordNet Sear... IA-07 http://mail.goog... Erbase2008.ppt pln-erbase2008l... 11:27



Tarefas PLN

– Alguns exemplos...

- Tradutores
- Corretores ortográficos
- Sumarizadores
- Geração automática de recursos linguísticos e semânticos (dicionários, ontologias)



(2008)

- Na filosofia da linguagem, a linguagem natural (ou linguagem comum) é uma língua que é falada, escrita, **assinada** por seres humanos **ou** para fins gerais de comunicação, **como distinguir de linguagens formais** (como as linguagens de programação de computador ou **os "línguas" Utilizados** no estudo da lógica formal, especialmente lógica matemática) e **construída a partir de línguas**.
- Na filosofia da linguagem, a linguagem natural (ou linguagem comum) é uma língua que é falada, escrita, **ou gestualizada** por seres humanos para fins gerais de comunicação, **e são distintas das linguagens formais** (como as linguagens de programação de computador ou **as linguagens utilizadas** no estudo da lógica formal, especialmente lógica matemática) **e das linguagens construídas**. (Esperanto)

Está evoluindo (2018)

- Na filosofia da linguagem, **uma** linguagem natural (ou linguagem comum) é uma **linguagem** falada, escrita ou **assinada** por seres humanos para comunicação de propósito geral, **como distinta** das linguagens formais (como linguagens de programação de computador ou "línguas" usadas no estudo da lógica formal, especialmente a lógica matemática) e das linguagens construídas.
- Na filosofia da linguagem, **a** linguagem natural (ou linguagem comum) é uma **língua que** é falada, escrita, ou **gestualizada** por seres humanos para fins gerais de comunicação, **e são distintas** das linguagens formais (como as linguagens de programação de computador ou as linguagens utilizadas no estudo da lógica formal, especialmente lógica matemática) e das linguagens construídas.
(esperanto)

Extração de informação de textos

- Reconhecimento de entidades nomeadas
 - Reconhecer a ocorrência de nomes próprios e suas categorias
 - pessoas, organizações, locais, etc...
- Reconhecimento de relações entre entidades
 - Fusões/aquisições de empresas
 - Relações pessoas organizações
 - Etc...
- Análise de sentimentos

Resolução de correferência textual

- Pesquisadores do Museu Nacional do Rio de Janeiro anunciaram a descoberta de **uma nova espécie de dinossauro** no Brasil. **O animal** que na cadeia evolutiva dos dinossauros ocuparia uma posição no grupo Tyrannoraptora, o mesmo do Tyrannosaurus rex, habitou o nordeste brasileiro há 110 milhões de anos, no período Cretáceo. **O fóssil, batizado de Santanaraptor placidus, é o único a ser encontrado no país com tecidos preservados.** Isso pode permitir que os cientistas saibam mais sobre o modo de vida e a evolução dos répteis.

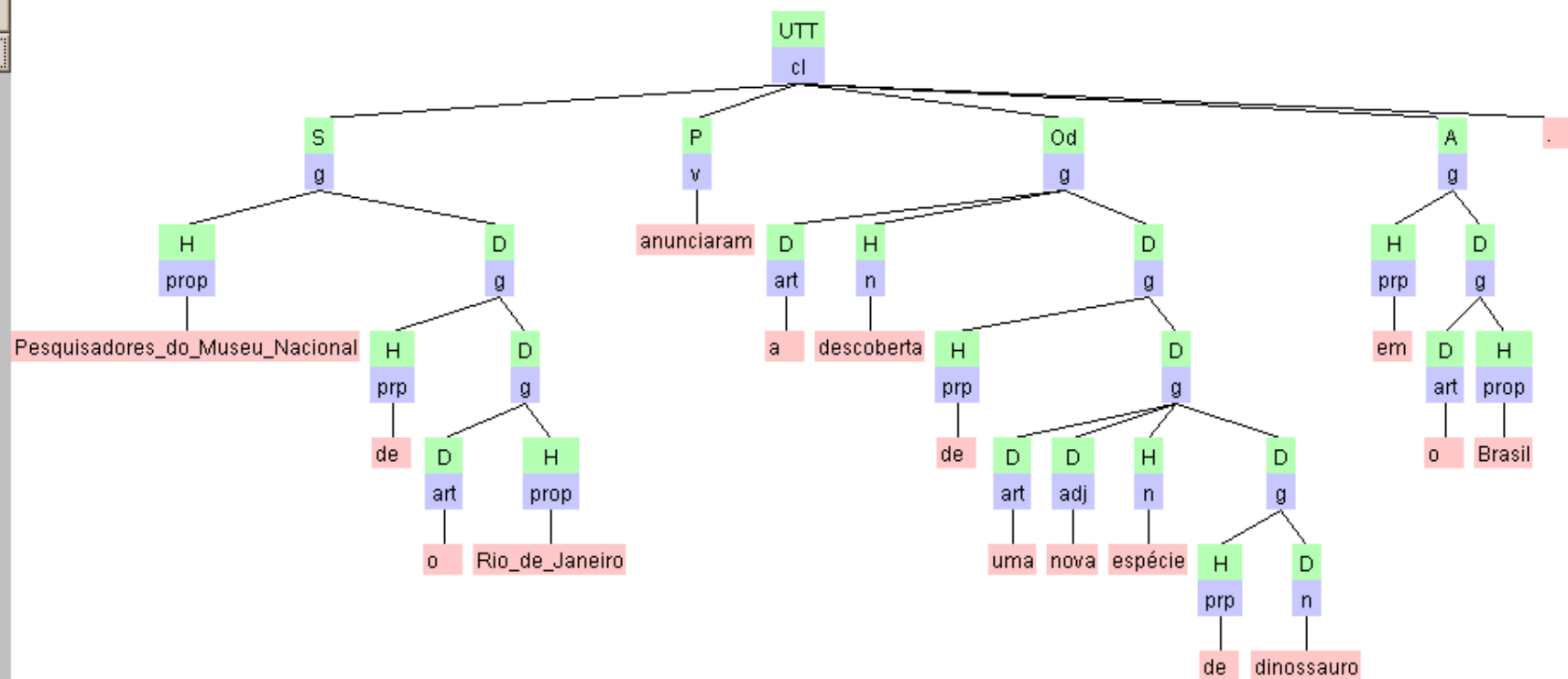
Sentence: Pesquisadores_do_Museu_Nacional de o Rio_de_Janeiro anunciaram a descoberta de uma nova espécie de dinossauro em o Brasil .

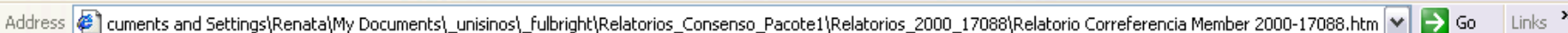
Function: S P Od Oi Op Cs Co fC A SUB CO CJT D H UTT STA QUE COM EXC X

Form: n prop v adj adv art pron prp conj num intj cl par g x

Collapse Tree

Expand Tree





1

Resolução de correferência

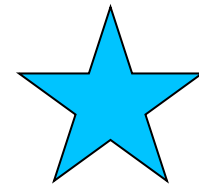
- Aplicações das cadeias de correferência?
 - Sumarização
- Sumário:
 - Segundo ele, pessoas ligadas por computadores tiveram mais contatos pessoais com seus amigos e parentes do que pessoas não-conectadas.

CADEIA : set_10		
word_8..word_9	---new	a internet
word_126..word_127	direct---old	a internet
CADEIA : set_3		
word_13..word_17	---new	os contatos entre as pessoas
word_20	---	os
word_62..word_66	indirect---old	a socialização de as pessoas
CADEIA : set_4		
word_16..word_17	---new	as pessoas
word_65..word_66	direct---old	as pessoas
CADEIA : set_14		
word_32..word_45	---	Barry_Ellman , de o Centro para Estudos_Urbanos e Comunitários de a Universidade_de_Toronto , Canadá
word_48..word_49	indirect---old	o pesquisador
word_73	---	ele
word_134	---	Ellman
word_137	---	ele
word_159..word_160	indirect---old	o pesquisador
CADEIA : set_11		
word_56	---	computadores
word_142	---	computadores
CADEIA : set_9		
word_89..word_95	---associative	as pessoas plugadas em uma rede local
word_108..word_109	---	vizinhos conectados
CADEIA : set_15		

Correferência e sumarização

– Resultado:

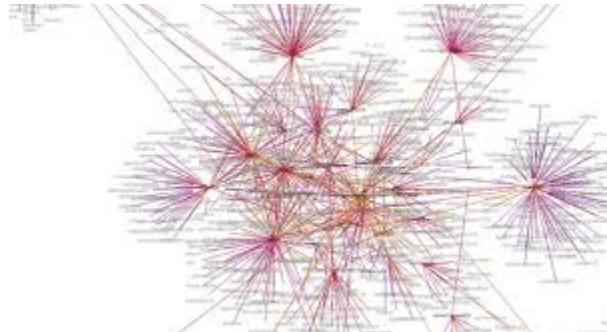
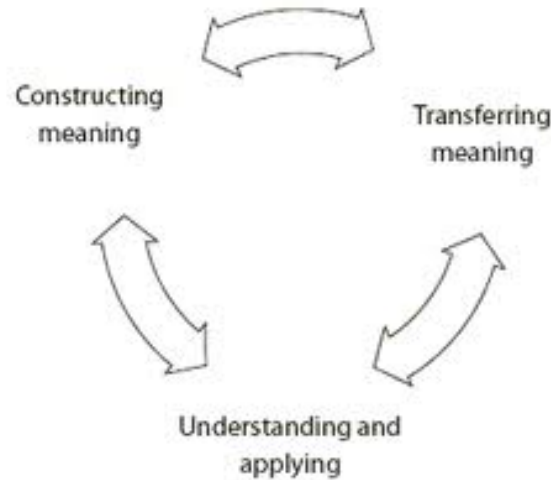
- Segundo Barry Ellman, do Centro para Estudos Urbanos e Comunitários da Universidade de Toronto, Canadá, pessoas ligadas por computadores tiveram mais contatos pessoais com seus amigos e parentes do que pessoas não-conectadas.



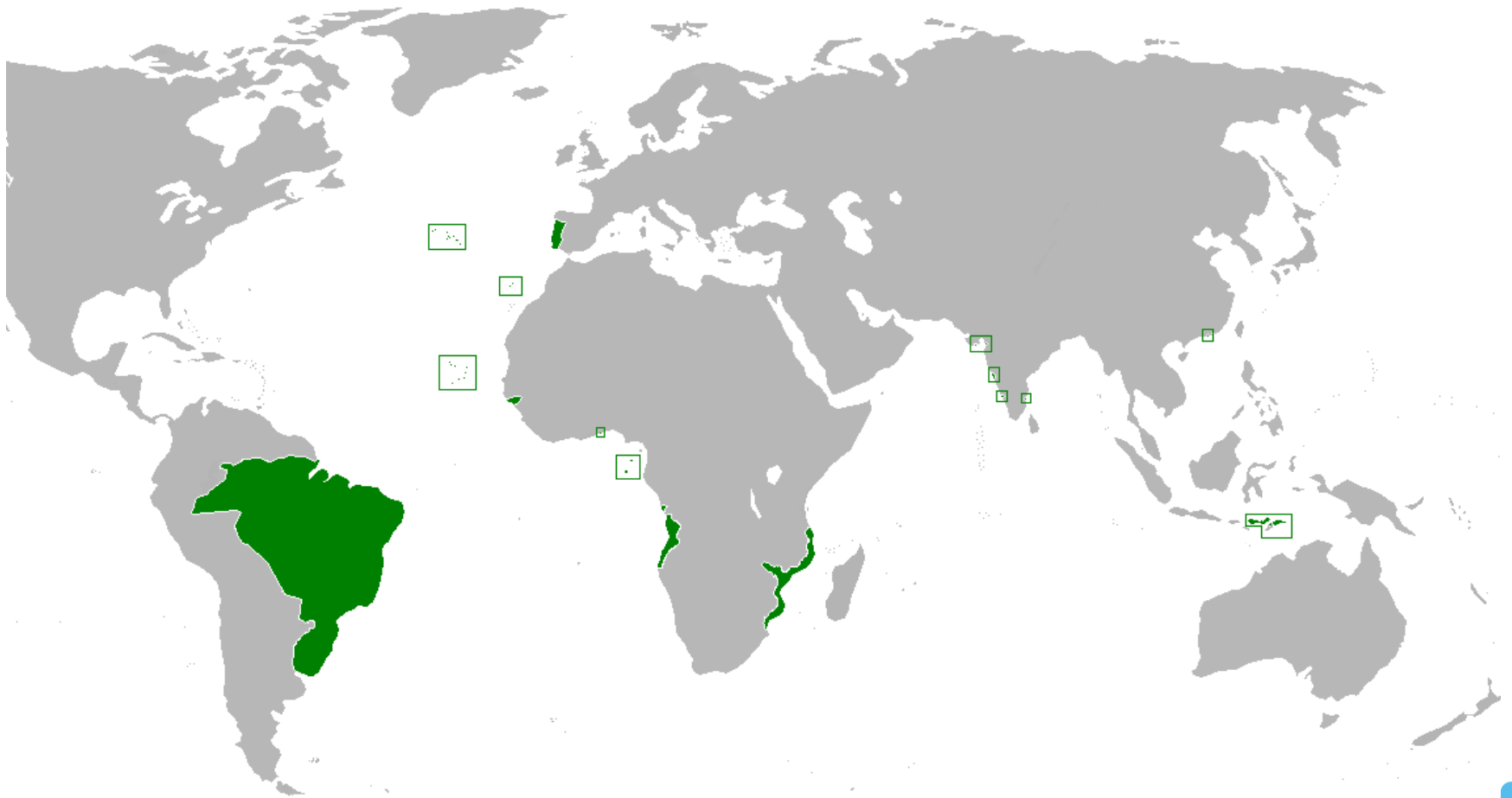
Diálogos

- Sistemas de perguntas e respostas
- Chat bots
- Comunicação Humano-Robô

Muitos desafios



PLN: e a língua portuguesa?



PLN no Brasil

- Necessidade de estudos fundamentais
- Construção de recursos de base
 - Léxicos e bases de dados lexicais
 - Ontologias e vocabulários
 - Gramáticas e *parsers*
 - Grandes coleções anotadas
- PROPOR, STIL
 - Propor 2020 em Évora PT

Arquivo Editar Exibir Histórico Favoritos Ferramentas Ajuda

http://www.nilc.icmc.usp.br/cepln/

powered by YAHOO! SEARCH Search Web

Nondete... Springer... AAMAS 2008 the Asso... propor2008... WordNet... Powerset B... Microsoft K... Comi... Visuwor...

Comissão Especial de Processamento de Linguagem Natural

Principal

Comissão

Regimento

Eventos

Periódicos

Fóruns

Novidades

A criação da Comissão Especial de Processamento de Linguagem Natural (CE-PLN) foi aprovada durante o [XXVII Congresso da Sociedade Brasileira de Computação](#) (realizado no Rio de Janeiro-RJ em Junho/Julho de 2007) por pedido das Profas. Dras. Maria das Graças V. Nunes (da Universidade de São Paulo - [USP/São Carlos](#)), Renata Vieira (da Pontifícia Universidade Católica do Rio Grande do Sul - [PUC-RS](#)) e Vera L. Strube de Lima (da Pontifícia Universidade Católica do Rio Grande do Sul - PUC-RS), que representavam a comunidade de PLN. A comissão reúne associados com interesses comuns na área de PLN.

A área de Processamento da Linguagem Natural (PLN), também denominada Linguística Computacional ou, ainda, Processamento de Línguas Naturais, lida com problemas relacionados à automação da interpretação e da geração da língua humana em aplicações como Tradução Automática, Sumarização Automática de Textos, Ferramentas de Auxílio à Escrita, Perguntas e Respostas, Categorização Textual, Recuperação e Extração de Informação, entre muitas outras, além das tarefas relacionadas de criação e disponibilização de dicionários/léxicos e corpús eletrônicos, desenvolvimento de taxonomias e ontologias, investigações em linguística de corpús, desenvolvimento de esquemas de marcação e anotação de conhecimento linguístico-computacional, resolução anafórica, análise morfossintática automática, análise semântico-discursiva automática, etc.

Em seus processos, e no desenvolvimento de recursos, ferramentas e aplicações, a área tem uma forte interação interdisciplinar, principalmente com as áreas de Linguística e Ciência da Informação, e no Brasil tem suas raízes na área de Inteligência Artificial.

O cenário gerado com a Internet e a demanda por serviços e produtos de Tecnologia da Informação tem ampliado ainda mais o campo de atuação do pesquisador desta área e impulsionado o mercado de trabalho.

O objetivo da CE-PLN é promover e representar a área de PLN no Brasil, apoiando e realizando eventos científicos, propondo e organizando meios de publicação e divulgação para a área e gerenciando listas e fóruns de discussão, dentre outras medidas.



Localizar: prt Próxima Diferenciar maiúsc./minúsc. Atingido o fim da página, continuando do início

Concluído 1.656s Images

Iniciar Comissão Esp... IA-07 http://mail.goog... Erbase2008.ppt pln-erbase2008l... 11:45

Part-Of-Speech and Parsers

COLLOVINI, Sandra et al. Cross-Framework Evaluation for Portuguese POS Taggers and Parsers. In: 19th International Conference on Computational Linguistics and Intelligent Text Processing, CICLing2018, 2018, Hanoi, Vietnam. 19th

POS Taggers

- TreeTagger (Schmid, 1994)
 - <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>
- NLTK (Bird, 2006)
 - <https://www.nltk.org/>
- NLPNet (Fonseca and Rosa, 2013)
 - <http://nilc.icmc.usp.br/nlpnet/>
- UDPipe (Straka et al., 2016)
 - <http://ufal.mff.cuni.cz/udpipe>

Parsers

- Palavras (Bick, 2000)
 - <http://visl.sdu.dk/visl/pt/>
- MaltParser (Nivre et al., 2007)
 - <http://www.maltparser.org/>
- LX-Parser (Silva et al., 2010)
 - <http://lxparser.di.fc.ul>
- Freeling (Padró and Stanilovsky, 2012)
 - <http://nlp.lsi.upc.edu/freeling/demo/demo.php>
- CoGrOO (Silva, 2013)
 - <http://cogroo.sourceforge.net/>

Bases Semânticas

Thiago Machado Lima. Analysing Semantic Resources for Coreference Resolution. In Master's Dissertation. PUCRS, 2019.

- **OntoPT**

- <http://ontopt.dei.uc.pt/>

- **ContoPT**

- <http://ontopt.dei.uc.pt/index.php?sec=contopt>

- **ConceptNet**

- <http://conceptnet.io/>

- **OpenWordNet-PT**

- <https://github.com/own-pt/openWordnet-PT>

- **DBPedia**

- <http://pt.dbpedia.org/>

- **BabelNet**

- <https://babelnet.org/>

Outras Abordagens

- Vamos transformar tudo em números
- Abordagens estatísticas
- Observação da frequência dos vizinhos
 - diga me com quem andas que te direi quem és...
 - word embeddings