

Para quê gastar 3 dias em trabalho manual quando podemos fazer (em 3 meses) uma ferramenta que o automatize?

ou

Lutando por um **NLP operating system**

J.João Dias de Almeida (jj@di.uminho.pt)

11 de Setembro de 2008



Missão

Melhorar a qualidade do PLN do Português, através de atenuar as dificuldade com que se deparam que os investigadores e desenvolvedores desta área.

- Disponibilizando recursos e ferramentas que permitam processamento sofisticado do Português.
- Catalogando e monitorizando a area
- Organizando actividades de avaliação (conjunta)



Missão

... de acordo com as nossa preocupações e crenças:

- Language processing techniques:
 - Domain specific languages (DSL)
 - grammar based
 - rule based
- Formal approach applied to informal, noisy domains
- Scripting and fast prototype
- Deal with real size resources

Inserido na realidade académica (U.Minho / DI)

e integrando coisas divertidas



com1 < a > b

b = f(a)

com1 < a | com2 | com3 > b

b = com3(com2(com1(a)))

estruturas mais complexas

grafo de cálculo -- Makefile

texto = seq(linha)

-- texto = frases, palavras, lemas

grep

-- nlgrep

yacc

-- nlyacc



- **TerminUM**: Terminology from bilingual resources:
 - **NATools**: set of tools for parallel corpora:
 - aligners
 - probabilistic translation dictionary extractor
 - parallel corpora client-server
 - bilingual terminology extractor
 - **Parguess**: find bibtexes in the WEB



- Recursos e ferramentas de língua:
 - **JSpell**: analisador morfológico
 - **Chuveiro de Dicionários**: um sistema para geração de dicionários de correctores ortográficos Firefox, Thunderbird, OpenOffice, aspell, etc para PT
 - **DPL**: Dictionary programming language
 - **DAC**
 - **Lingua::StarDict::Gen**: geração de dicionários StarDict



- **T₂O**: Ontology algebra:
 - **Biblio::Thesaurus**
 - **TabularThesaurus**
- General Support tools:
 - **Makefile::Parallel**: DSL for large scale computing using cluster
 - **Text::RewriteRules**: DSL for Textual rewriting systems in Perl context
 - **XML::TMX**: processing Translation Memories
 - **XML::DT**: processing HTML and XML files







