

# REMMA - Reconhecimento de Entidades Mencionadas do MedAlert

Liliana Ferreira e António Teixeira

Encontro do Segundo HAREM  
Setembro 2008  
Aveiro

- 1 Introdução
- 2 Arquitectura
- 3 Participação
- 4 Conclusões



# REMMA - Reconhecimento de Entidades Mencionadas do MedAlert

- Desenvolvido em UIMA - *Unstructured Information Management Architecture*
- Principais recursos:



*almanaques*

+

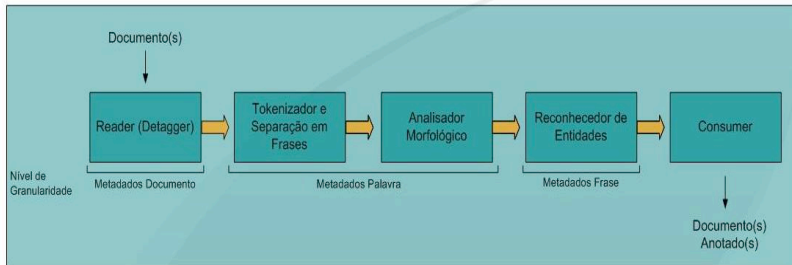


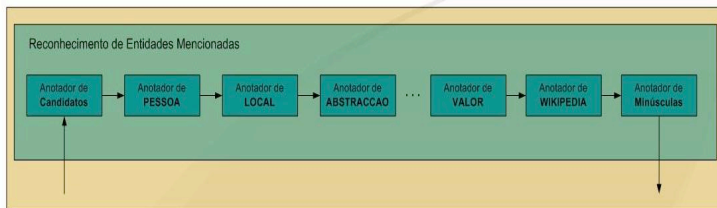
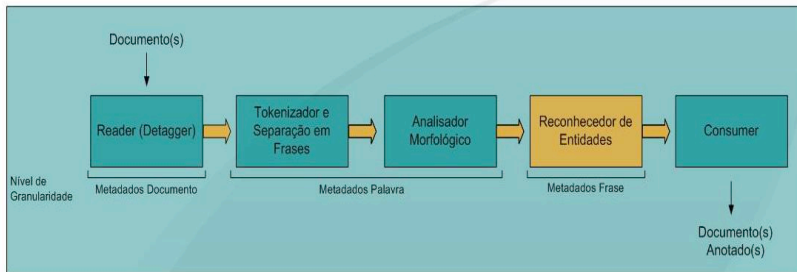
*regras*

+



*wikipedia*





### Problema:

- Almanques ou dicionários de entidades são importantes em REM, no entanto, a sua construção e manutenção é extremamente demorada.

### Solução REMMA:



- A wikipedia disponibiliza os conteúdos da base de dados para cada língua em XML e os ficheiros SQL necessários para construir a BD.

Dumps from any Wikimedia Foundation project:

<http://download.wikimedia.org/>

- Foi utilizada a versão portuguesa de Fevereiro de 2008, contendo cerca de 1 290 836 páginas.

Embora não exista uma convenção, o normal é começar o artigo por uma pequena frase descrevendo o que a entidade é:

### Exemplo

*A "Universidade de Aveiro (UA)" é uma **universidade** pública portuguesa localizada em Aveiro.*

### Estrutura básica:

- Os artigos são identificados por um nome único - concatenação das palavras do título do artigo com " \_";  
*ex.: Universidade de Aveiro - Universidade\_de\_Aveiro*
- Obtemos o artigo correspondente à entidade, seguindo os redireccionamentos caso seja necessário.



### Redirecção:

- Algumas páginas não têm artigo mas são direccionadas para outro artigo com um nome diferente.
- Estes mecanismo é denominado "redirecção" e é identificado por "#REDIRECT [[A B C]]".
- Nestes casos segue-se para o artigo com o título A.B.C.

### Desambiguação:

- Alguns autores desenvolvem uma página de desambiguação para um artigo com nome ambíguo.
- Nós não usamos estas páginas!

Após a obtenção do artigo correspondente, extraímos a primeira frase da qual foi extraída a categoria:

### Exemplo

A "*Universidade de Aveiro (UA)*" é uma *universidade* pública portuguesa localizada em Aveiro.



ORGANIZACAO

- a Não usamos a secção *Categoria* dos artigos uma vez que um artigo pode ter mais que uma *Categoria* e muitas das apresentadas não são claramente hiperónimos da entidade.

## Participação

Cenário *Selectivo 4*: inclui todas as categorias e tipos previstos nas directivas mas não considera os subtipos associados, concretamente os subtipos de LOCAL e de TEMPO.

## Resultados

- 1 Saída I: Regras + Almanagues + Wikipedia  
45 081 entidades em 52m06s
- 2 Saída II: Regras + Almanagues  
37 707 entidades em 32m24s
- 3 Saída III: Wikipedia  
29 949 entidades

Pretendemos demonstrar, com esta experiência, a utilidade de recursos deste tipo na classificação de entidades quando usados em conjunto com outros recursos como almanaques.

- 1 Explorar a *Wikipedia* como fonte de conhecimento externo:
  - a extracção a categoria a partir da primeira frase do artigo.

## Trabalho Futuro:

- 1 Explorar a estrutura dos artigos:
  - a Categoria;
  - b Páginas de desambiguação;
  - c ...
- 2 Aplicar o método a outros recursos semelhantes na área da Medicina!

# REMMA - Reconhecimento de Entidades Mencionadas do MedAlert

Liliana Ferreira e António Teixeira

Encontro do Segundo HAREM  
Setembro 2008  
Aveiro

*Fim! Obrigado.*