

# **HAREM e MiniHAREM: Uma análise comparativa**

Nuno Cardoso

Encontro HAREM

15 de Julho de 2006

FLUP, Porto

*Linguateca*

# Participação no HAREM

## **HAREM (14-2-2005)**

- 10 participantes
- 6 países
  - PT, BR, MX, DK, ES e FR
- 18 saídas (3 não-oficiais)

## **MiniHAREM (3-4-2006)**

- 5 participantes
- 2 países
  - PT e BR
- 20 saídas
  - Só para os participantes do primeiro HAREM

# HAREM numa casca de noz

## Colecção de textos HAREM

Eça de Queirós nasceu na  
Póvoa de Varzim em 1845.

# HAREM numa casca de noz

Colecção de textos HAREM

Eça de Queirós nasceu na  
Póvoa de Varzim em 1845.

Participante

Sistema REM  
participante

Etiquetagem  
automática

Saída do Participante

```
<PESSOA TIPO="INDIVIDUAL" MOREF="M,S">
Eça de Queirós</PESSOA> nasceu na
<PESSOA TIPO="INDIVIDUAL" MOREF="M,S">
Póvoa</PESSOA> de Varzim em 1845.
```

# HAREM numa casca de noz

## Colecção de textos HAREM

Eça de Queirós nasceu na  
Póvoa de Varzim em 1845.

## Participante

Sistema REM  
participante

Etiquetagem  
automática

## Saída do Participante

<PESSOA TIPO="INDIVIDUAL" MORF="M,S">  
Eça de Queirós</PESSOA> nasceu na  
<PESSOA TIPO="INDIVIDUAL" MORF="M,S">  
Póvoa</PESSOA> de Varzim em 1845.

Avaliação  
HAREM

## Colecção Dourada

<PESSOA TIPO="INDIVIDUAL" MORF="M,S">  
Eça de Queirós</PESSOA> nasceu na  
<LOCAL TIPO="ADMINISTRATIVO" MORF="M,S">  
Póvoa de Varzim</LOCAL> em <TEMPO  
TIPO="DATA">1845</TEMPO>.

# HAREM numa casca de noz

## Colecção de textos HAREM

Eça de Queirós nasceu na Póvoa de Varzim em 1845.

## Participante

Sistema REM participante

Etiquetagem automática

## Saída do Participante

<PESSOA TIPO="INDIVIDUAL" MORF="M,S">  
Eça de Queirós</PESSOA> nasceu na  
<PESSOA TIPO="INDIVIDUAL" MORF="M,S">  
Póvoa</PESSOA> de Varzim em 1845.

Avaliação  
HAREM

## Colecção Dourada

<PESSOA TIPO="INDIVIDUAL" MORF="M,S">  
Eça de Queirós</PESSOA> nasceu na  
<LOCAL TIPO="ADMINISTRATIVO" MORF="M,S">  
Póvoa de Varzim</LOCAL> em <TEMPO  
TIPO="DATA">1845</TEMPO>.

## Pontuações

### Identificação:

Eça de Queirós: **Correcto**  
Póvoa de Varzim:  
**Parcialmente Correcto**  
1845: **EmFalta**

### Classificação Morfológica:

Eça de Queirós: **Correcto**  
Póvoa de Varzim: **Errado**  
**no Género**

### Classificação Semântica:

Eça de Queirós: **Correcto**  
Póvoa de Varzim:  
**EmFalta LOCAL**  
**Espúrio PESSOA**  
1845:  
**EmFalta TEMPO**

# Cenários selectivos em 2005

## • HAREM

Sistema	Saída	PESSOA	ORGAN.	LOCAL	TEMPO	VALOR	ACONT.	ABSTR.	COISA	OBRA	VARIADO
CaGE	1			■							
	2	■	■	■		■					
	3			■							
Cortex	1	■	■	■	■	■					
	2	■	■	■	■	■					
	3	■	■	■	■	■	■	■	■	■	■
ELLE			■			■					
Malinche		■	■	■	■	■	■	■	■	■	■
Nerua	1	■	■	■	■	■	■	■	■	■	■
	2	■	■	■	■	■	■	■	■	■	■
	3	■	■	■	■	■	■	■	■	■	■
RSN-NILC		■	■	■	■	■	■	■	■	■	■
PALAVRAS-NER		■	■	■	■	■	■	■	■	■	■
RENA		■	■	■	■	■	■	■	■	■	■
SIEMÊS	1	■	■	■	■	■	■	■	■	■	■
	2	■	■	■	■	■	■	■	■	■	■



= alguns tipos



= todos os tipos

# Cenários selectivos em 2006

- MiniHAREM

Sistema Saída	PESSOA	ORGAN.	LOCAL	TEMPO	VALOR	ACONT.	ABSTR.	COISA	OBRA	VARIADO
CaGE			■							
Siemês 2	■	■	■	■	■	■	■	■	■	■
Cortex	■	■	■	■	■	■				
SMELL	■	■	■	■	■	■				
Stencil-Nooj	■	■	■	■	■					

■ = alguns tipos

■ = todos os tipos



## Colecção Dourada (CD)



# Propósito da Coleção Dourada

- Incluir todas as EMs relevantes em texto português, obtendo uma marcação “ideal” de EMs no texto
- Categorias das EMs criadas empiricamente a partir de análise do texto. Categorização feita em dois níveis, categorias e tipos.
- Não representa o que os sistemas REM devem obter hoje, mas permitem:
  - avaliar a dificuldade da tarefa REM
  - estabelecer um limite superior
- Há muito mais em REM do que pessoas, organizações, locais e números...

# Colecções Douradas usadas

- HAREM: Colecção Dourada de 2005 (CD 2005)
- MiniHAREM: Colecção Dourada de 2006 (CD 2006)
- Ambas as CDs foram retiradas da Colecção HAREM

Tamanhos	Colecção HAREM	CD 2005	CD 2006	Ambas
Palavras	600 086	92 830	62 461	155 291
Documentos	1 202	129	128	257
EMs	~ 40 000	5 270	3 858	9 128
EMs vagas (class.)	~ 1 000	133	142	275
EMs vagas (ident.)	~ 500	71	56	127

# Categorias e Tipos em 2005

- **ABSTRACCAO**
  - DISCIPLINA
  - ESTADO
  - ESCOLA
  - OBRA
  - PLANO
  - IDEIA
  - NOME
- **OBRA**
  - ARTE
  - REPRODUZIDA
  - PRODUTO
  - PUBLICACAO
- **ACONTECIMENTO**
  - EFEMERIDE
  - ORGANIZADO
  - EVENTO
- **COISA**
  - OBJECTO
  - SUBSTANCIA
  - CLASSE
- **TEMPO**
  - DATA
  - HORA
  - PERIODO
  - CICLICO
- **ORGANIZACAO**
  - INSTITUICAO
  - ADMINISTRACAO
  - EMPRESA
  - SUB
- **PESSOA**
  - INDIVIDUAL
  - GRUPOIND
  - CARGO
  - GRUPOCARGO
  - MEMBRO
  - GRUPOMEMBRO
- **LOCAL**
  - GEOGRAFICO
  - ADMINISTRATIVO
  - VIRTUAL
  - ALARGADO
  - CORREIO
- **VARIADO**
  - OUTRO
- **VALOR**
  - MOEDA
  - CLASSIFICACAO
  - QUANTIDADE

# Categorias e Tipos em 2006

- **ABSTRACCAO**

- DISCIPLINA
- ESTADO
- ESCOLA
- OBRA
- PLANO
- IDEIA
- NOME

- **PESSOA**

- INDIVIDUAL
- GRUPOIND
- CARGO
- GRUPOCARGO
- MEMBRO
- GRUPOMEMBRO

- **OBRA**

- ARTE
- REPRODUZIDA
- **PRODUTO**
- PUBLICACAO

- **ORGANIZACAO**

- INSTITUICAO
- ADMINISTRACAO
- EMPRESA
- SUB

- **LOCAL**

- GEOGRAFICO
- ADMINISTRATIVO
- VIRTUAL
- ALARGADO
- CORREIO

- **ACONTECIMENTO**

- EFEMERIDE
- ORGANIZADO
- EVENTO

- **TEMPO**

- DATA
- HORA
- PERIODO
- CICLICO

- **COISA**

- OBJECTO
- SUBSTANCIA
- CLASSE
- **MEMBROCLASSE**

- **VARIADO**

- OUTRO

- **VALOR**

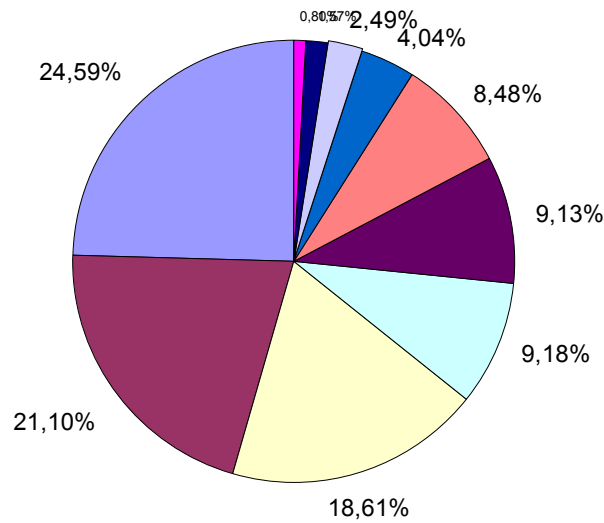
- MOEDA
- CLASSIFICACAO
- QUANTIDADE

# Principais alterações HAREM => MiniHAREM

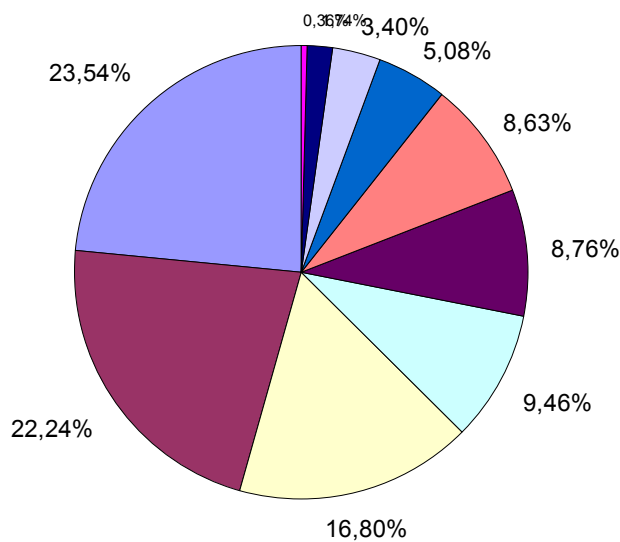
- OBRA TIPO="PRODUTO" eliminado
- COISA TIPO="MEMBROCLASSE" criado
- Emails e URLs deixam de ser marcados
- Referências anafóricas mantêm significado semântico
  - ex: Revolução de 1830 (...) e a de 1832...
- PESSOA TIPO="MEMBRO" para referências a pessoas a partir de uma organização

# CD: Distribuição por Categorias

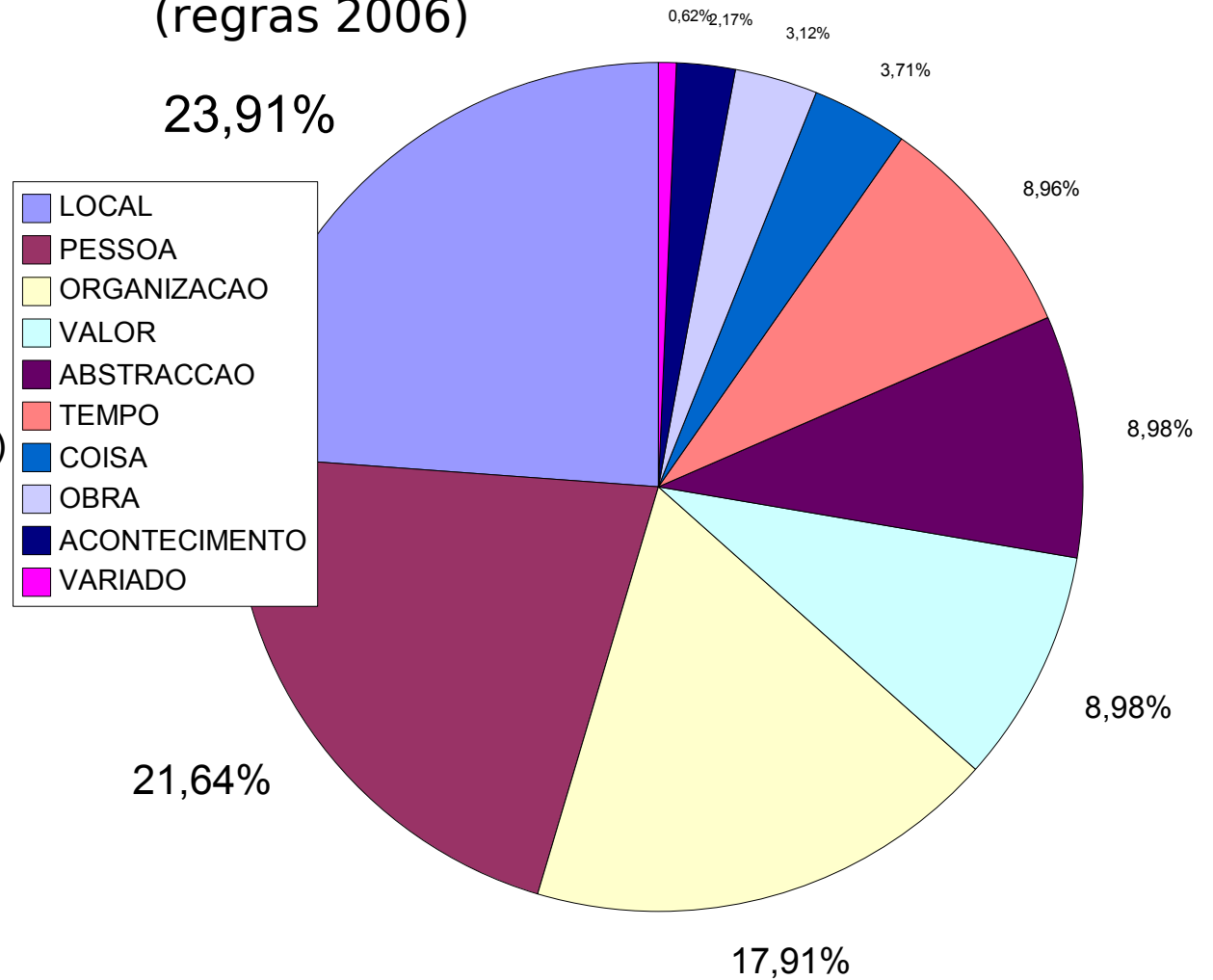
Colecção Dourada de 2005 (regras 2005)



Colecção Dourada de 2006 (regras 2006)



## Ambas as Colecções Douradas (regras 2006)

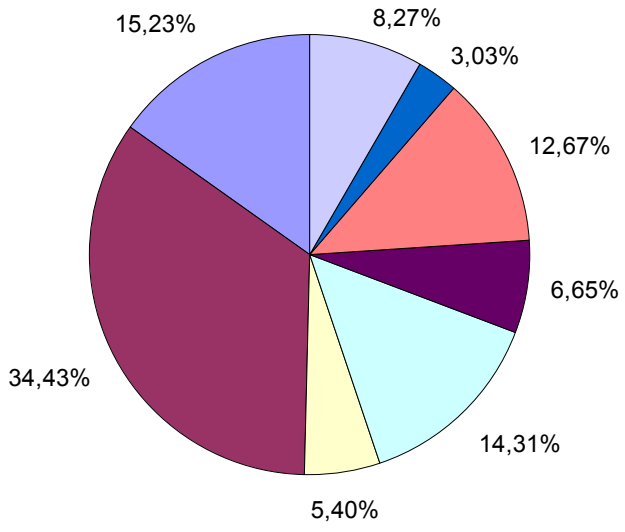


- LOCAL
- PESSOA
- ORGANIZACAO
- VALOR
- ABSTRACAO
- TEMPO
- COISA
- OBRA
- ACONTECIMENTO
- VARIADO

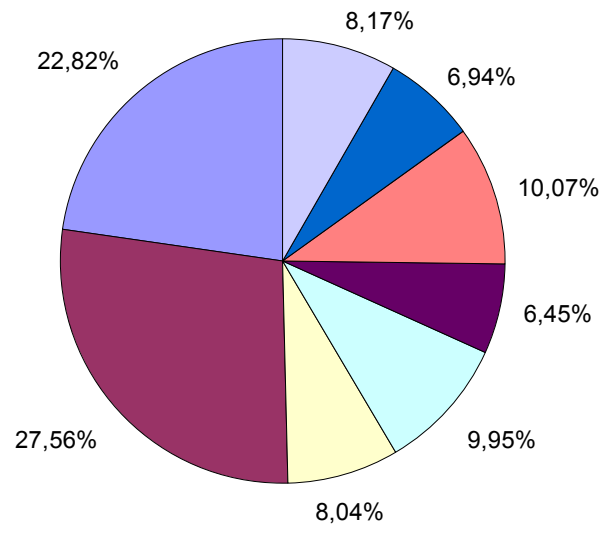
# CD: Distribuição por Género Textual (nº de palavras)

## Ambas as Colecções Douradas

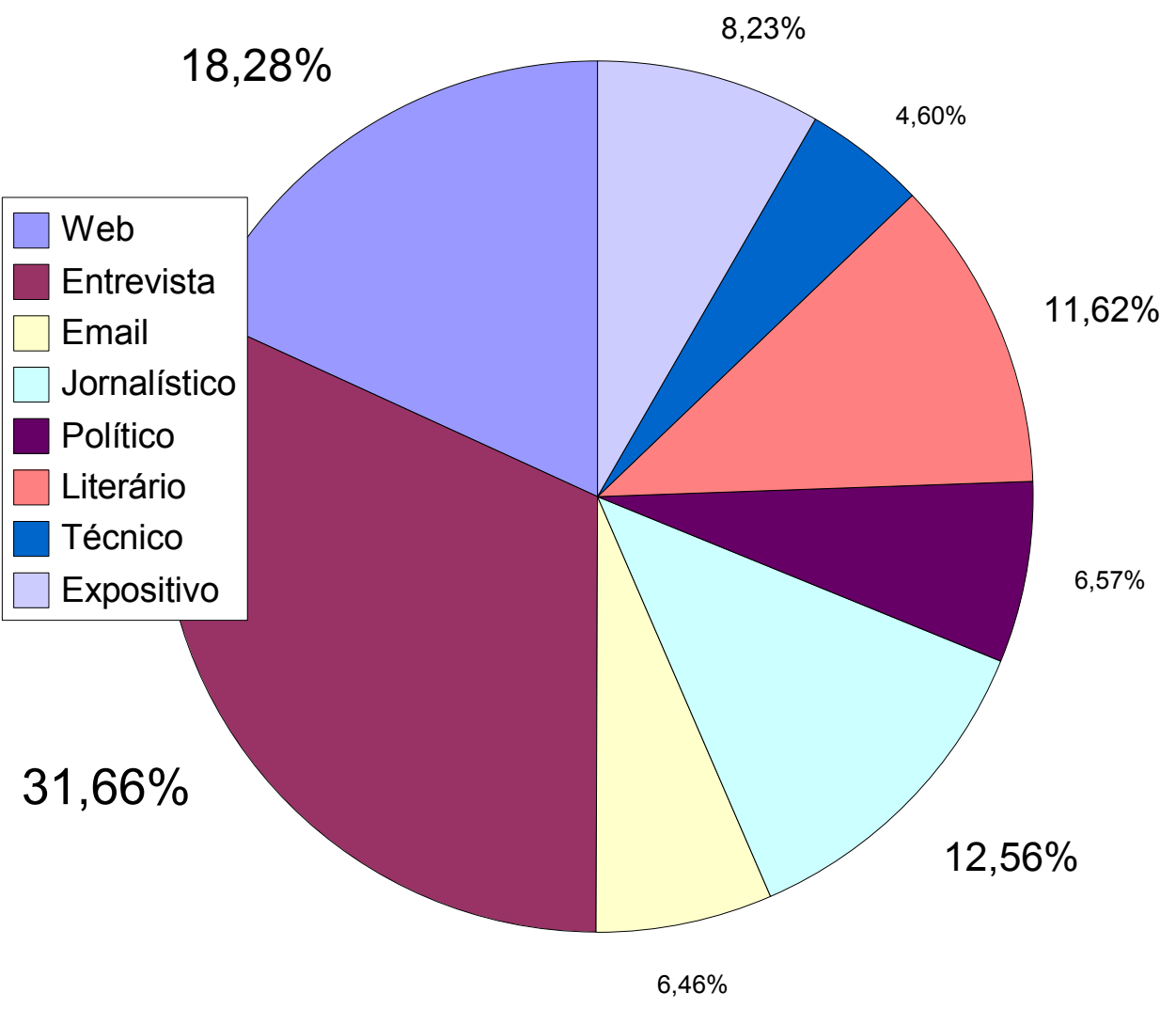
Colecção Dourada de 2005



Colecção Dourada de 2006



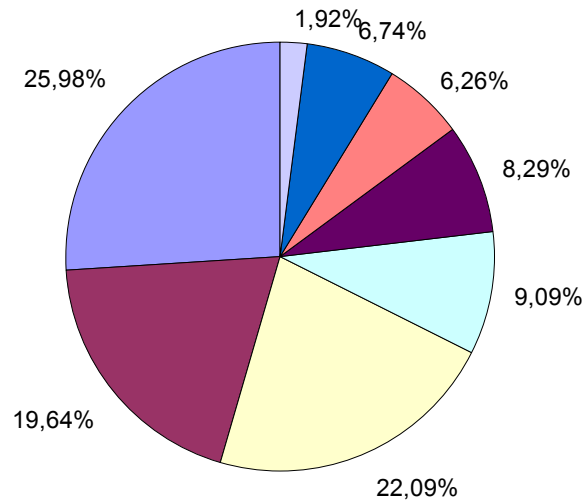
- Web
- Entrevista
- Email
- Jornalístico
- Político
- Literário
- Técnico
- Expositivo



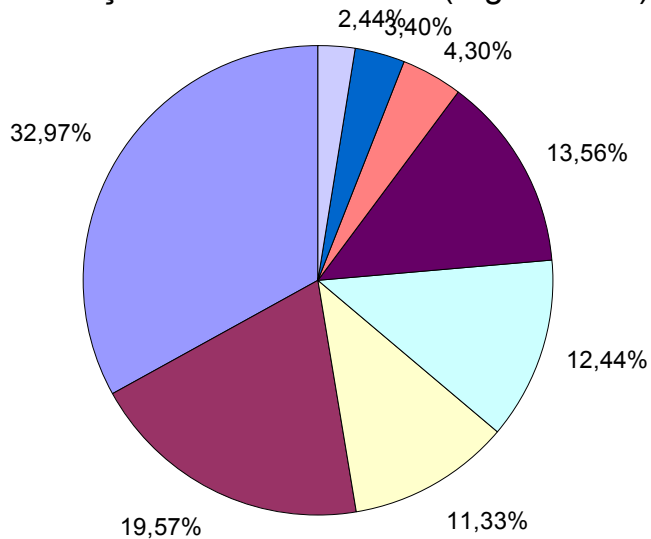


# CD: Distribuição por Género Textual (nº de EMs)

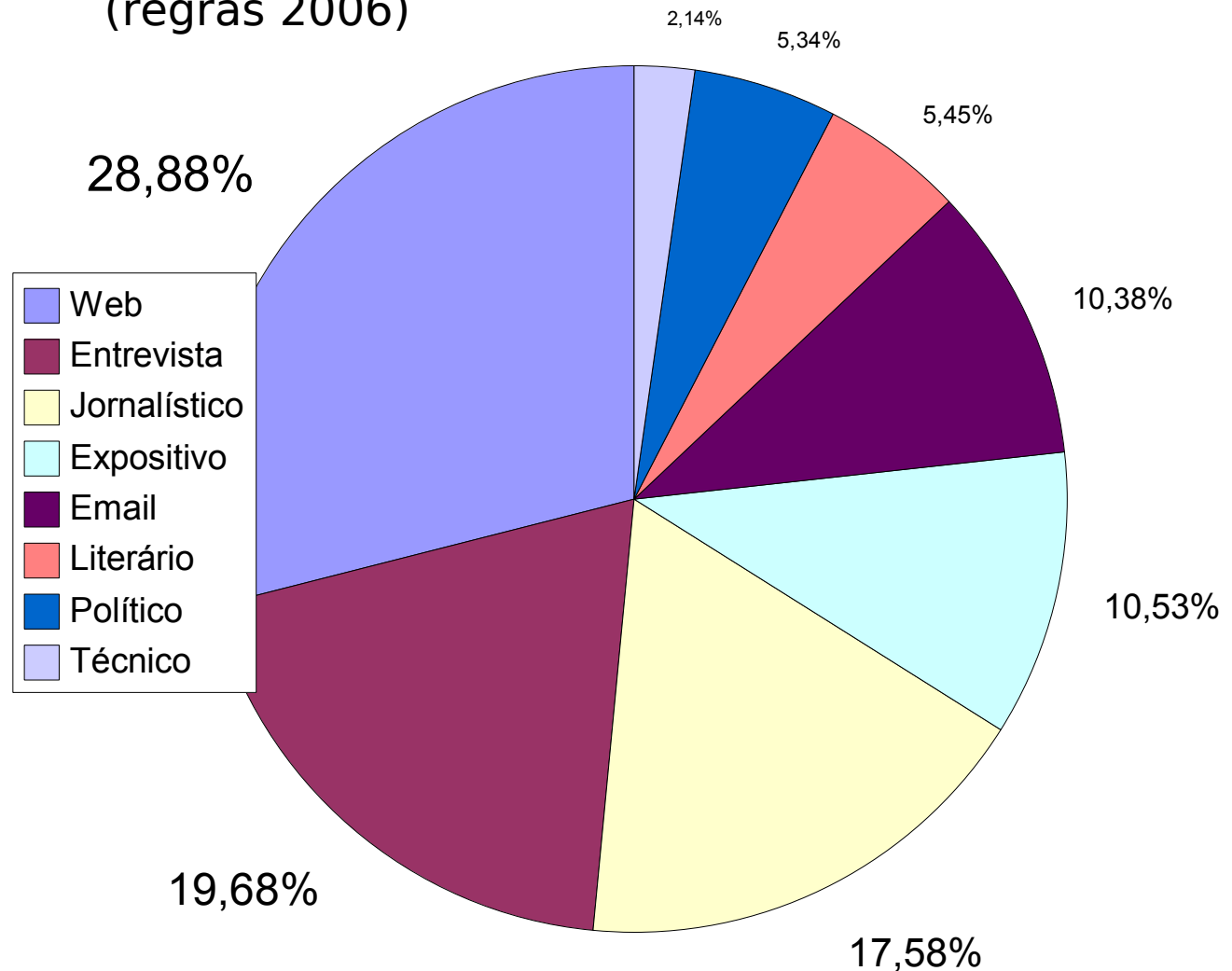
Colecção Dourada de 2005 (regras 2005)



Colecção Dourada de 2006 (regras 2006)

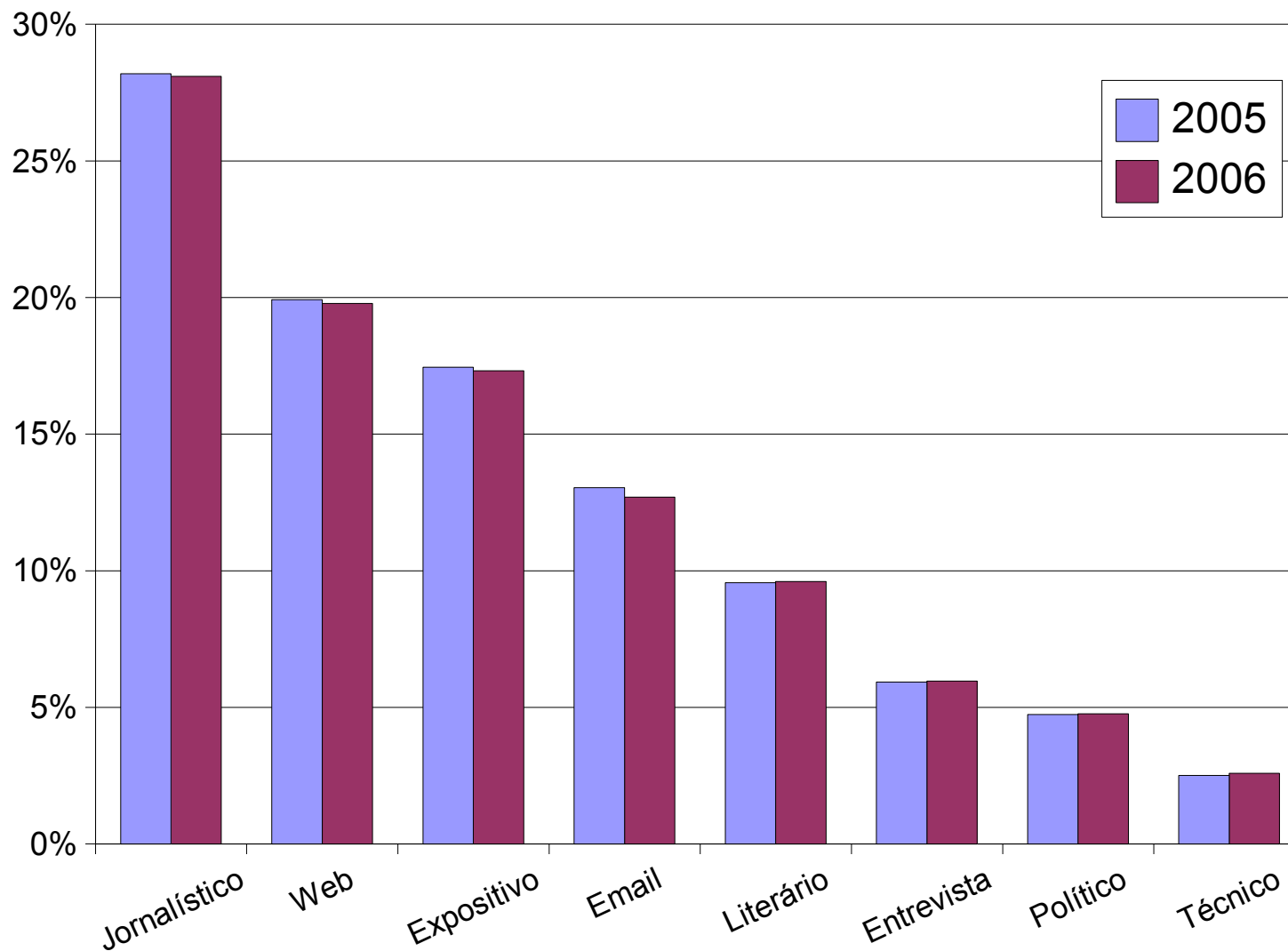


## Ambas as Colecções Douradas (regras 2006)



# Densidade de EMs, por Género Textual

$$\text{Densidade EM}_G = \frac{\text{N}^\circ \text{ palavras que pertencem e EMs, para o género G}}{\text{N}^\circ \text{ total de palavras, para o género G}}$$



# Nº de palavras das EMs, por Categoria

Categorias	2005			2006		
	Média	Mediana	Desv.Pad.	Média	Mediana	Desv.Pad.
ACONTECIMENTO	3,34	3	2,94	3,76	3	3,16
OBRA	3,26	2	2,89	3,5	3	3,19
VARIADO	2,25	1	2,51	2,23	1	2,48
ABSTRACCAO	2,19	1	2,44	2,21	1	2,01
ORGANIZACAO	2,19	1	1,96	2,21	1	2,45
PESSOA	1,9	2	1,12	1,9	2	1,10
TEMPO	1,81	1	1,34	1,82	1	1,34
VALOR	1,75	2	0,90	1,75	2	0,91
LOCAL	1,65	1	1,43	1,66	1	1,46
COISA	1,45	1	0,83	1,54	1	0,88
<b>TOTAL</b>	<b>1,97</b>	<b>1</b>	<b>1,73</b>	<b>1,98</b>	<b>1</b>	<b>1,76</b>

# Teor de Sobreposição de EMs

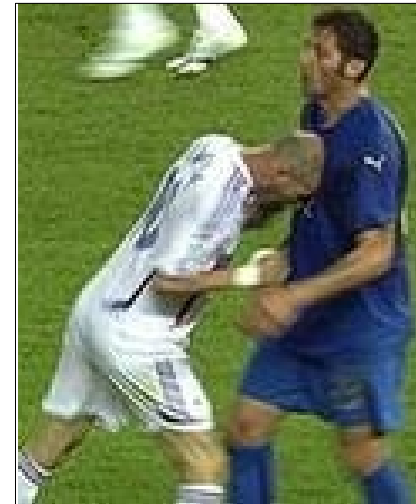
Quantas EMs aparecem em ambas as Coleções Douradas?

	CD 2005	CD 2006
Número Total de EMs	5132	3712
Número de EMs Distintas	3060	2434
Rácio	59,63%	65,57%
<hr/>		
Nº Total de EMs comuns	<b>623</b>	
Nº EMs Distintas comuns	<b>380</b>	
Rácio de total EMs	12,14%	16,78%
Rácio de EMs distintas	12,42%	15,61%

# Análise Estatística ao HAREM e MiniHAREM



VS



# Objectivos da análise estatística

- Determinar o nível de confiança dos resultados do HAREM e MiniHAREM
- Distinguir os sistemas / estratégias realmente diferentes
- Verificar se o tamanho da CD é suficiente
- **Validar os eventos de avaliação**

# Requisitos dos testes estatísticos


- Não paramétrico -- distribuição de EMs (quase) impossível de determinar --
- Comparação directa entre pares de saídas
- Simples e robusto
- Nível de confiança calculado facilmente

# Teste estatístico escolhido: permutação

- Escolha: Testes de permutação
  - *Approximate Randomization*
- Sinopse do teste de permutação:
  - Permutar aleatoriamente observações entre duas saídas
  - Se a diferença inicial diminuir, é provável que as duas saídas sejam diferentes
  - Se a diferença inicial se mantiver, é provável que seja obra do acaso





# Teste de Permutação - exemplo

Saída A   $P_A = 16/20 = 80\%$

Saída B   $P_B = 6/20 = 30\%$

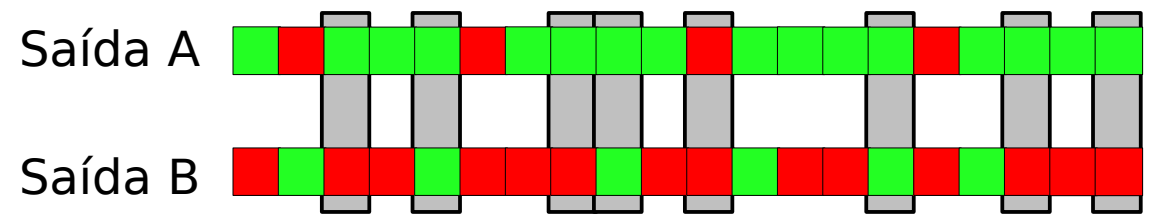
# Teste de Permutação - exemplo

Saída A   $P_A = 16/20 = 80\%$


Saída B   $P_B = 6/20 = 30\%$

## Fase 1: Permutação

- Troca de observações, com 1/2 de probabilidade



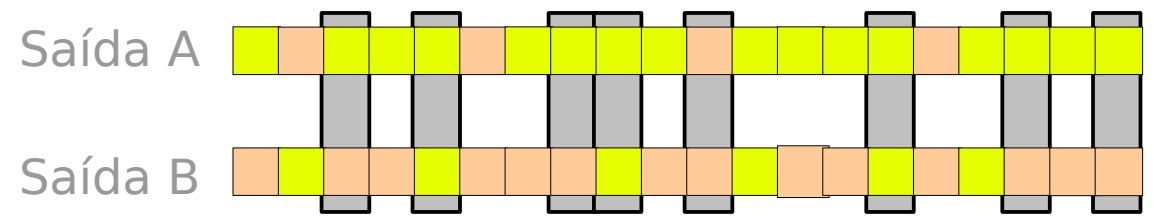
# Teste de Permutação - exemplo

Saída A   $P_A = 16/20 = 80\%$

Saída B   $P_B = 6/20 = 30\%$


## Fase 1: Permutação

- Troca de observações, com 1/2 de probabilidade




## Fase 2: Cálculo da nova diferença

Pseudo-Saída A   $P_A^* = 12/20 = 60\%$

Pseudo-Saída B   $P_B^* = 10/20 = 50\%$

# Teste de Permutação - exemplo

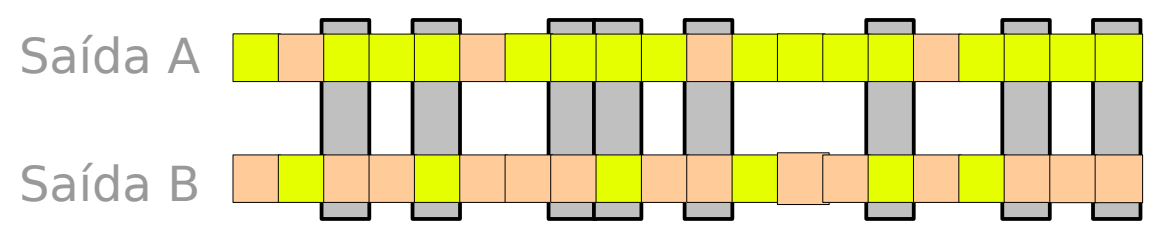
Saída A   $P_A = 16/20 = 80\%$

Saída B   $P_B = 6/20 = 30\%$


50%


## Fase 1: Permutação

- Troca de observações, com 1/2 de probabilidade



## Fase 2: Cálculo da nova diferença

Pseudo-Saída A   $P_A^* = 12/20 = 60\%$

Pseudo-Saída B   $P_B^* = 10/20 = 50\%$

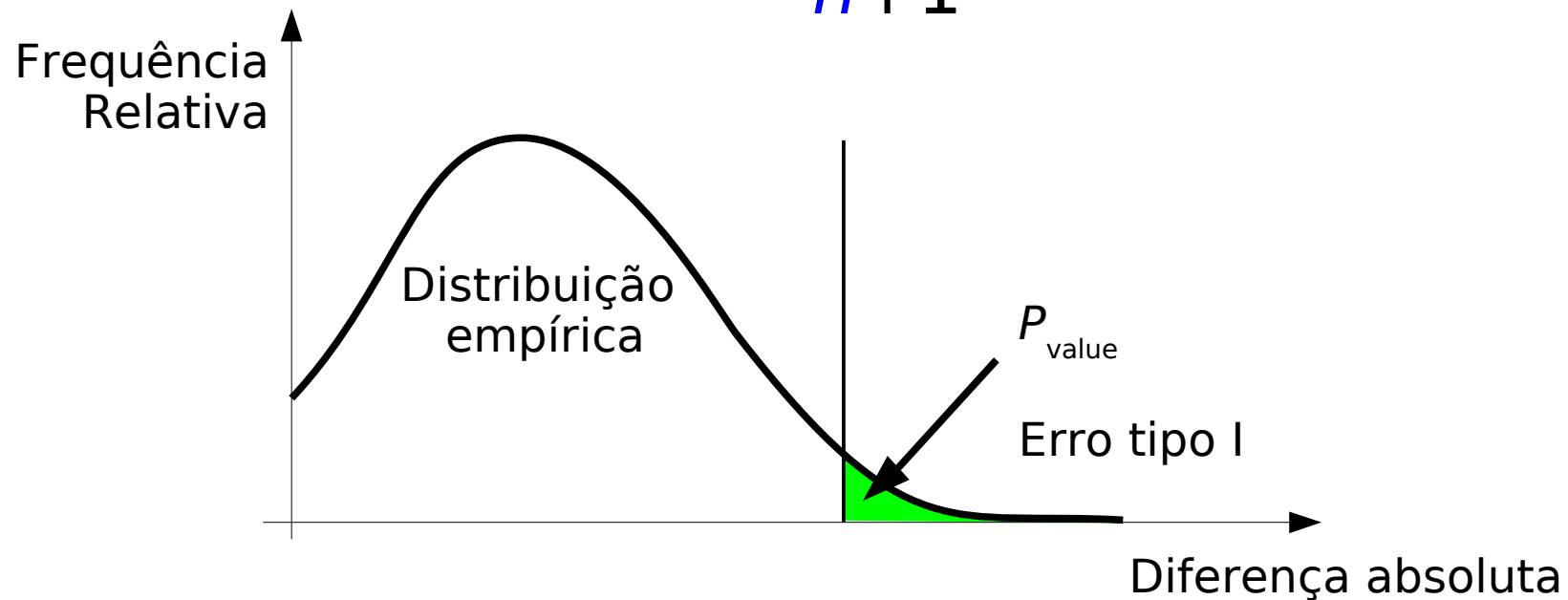
10%



# Teste de Permutação - exemplo

- Repetir para  $n$  iterações (ou seja, gerar *pseudo-saídas*)
  - Contar o nº de vezes  $m$  que a diferença aumentou em vez de diminuir

- Valor de P:  $P_{\text{value}} = \frac{m+1}{n+1}$



# Mas... como permutar no HAREM?

- Observações das saídas são dependentes!

Saídas	Texto / EMs
<b>A</b>	<p>Segundo<sup>①</sup> o presidente da Fundação<sup>②</sup> para o Desenvolvimento da Produção, Costa e Silva, ...<sup>③</sup></p>
<b>CD</b>	<p>Segundo o presidente da Fundação<sup>①</sup> para o Desenvolvimento da Produção, Costa e Silva, ...<sup>②</sup></p>
<b>B</b>	<p>Segundo o presidente<sup>①</sup> da Fundação<sup>②</sup> para o Desenvolvimento da Produção<sup>③</sup>, Costa<sup>④</sup> e Silva<sup>⑤</sup>, ...</p>

- Como permutar “Costa”? E se “Silva” não for permutado?

# Permutação por blocos

- Solução: agrupar EMs em blocos independentes

Saídas	Texto / EMs
A	
B	

- Mantém pontuações
- Resolve diferenças dadas pelo <ALT>
- Ver a permutação como a troca de desempenho entre saídas, para uma determinada frase ou pedaço de texto





# Testes estatísticos com tamanho da CD

<i>Saída A</i>			<i>Saída B</i>			<i>Diferença</i>		
Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>0.728</b>	<b>0.696</b>	<b>0.712</b>	<b>0.798</b>	<b>0.870</b>	<b>0.832</b>	<b>0.069</b>	<b>0.174</b>	<b>0.121</b>

# Testes estatísticos com tamanho da CD

<i>Saída A</i>			<i>Saída B</i>			<i>Diferença</i>		
Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>0.728</b>	<b>0.696</b>	<b>0.712</b>	<b>0.798</b>	<b>0.870</b>	<b>0.832</b>	<b>0.069</b>	<b>0.174</b>	<b>0.121</b>

- Utilizando cada vez menos blocos no teste...

# Blocos	<i>n</i> iterações = 9999			<i>PSEUDO-SAÍDAS de A</i>						<i>PSEUDO-SAÍDAS de B</i>					
	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão		
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>Todos</b>	0.0001	0.0001	0.0001	0.765	0.783	0.774	0.003	0.004	0.003	0.765	0.783	0.774	0.003	0.004	0.003
2000	0.0001	0.0001	0.0001	0.765	0.772	0.769	0.008	0.009	0.007	0.765	0.772	0.769	0.008	0.009	0.007
1000	0.0001	0.0001	0.0001	0.765	0.765	0.765	0.013	0.015	0.012	0.765	0.765	0.765	0.013	0.014	0.012
500	0.0004	0.0001	0.0001	0.766	0.761	0.763	0.019	0.022	0.017	0.765	0.761	0.763	0.018	0.021	0.017
250	<b>0.0181</b>	0.0001	0.0001	0.765	0.759	0.762	0.027	0.031	0.025	0.766	0.760	0.763	0.027	0.031	0.025
200	<b>0.0351</b>	0.0001	0.0001	0.765	0.759	0.762	0.030	0.034	0.028	0.765	0.759	0.762	0.030	0.035	0.028
100	<b>0.1391</b>	0.0009	0.0047	0.766	0.759	0.761	0.043	0.049	0.040	0.765	0.758	0.761	0.044	0.049	0.041
75	<b>0.1912</b>	0.0034	<b>0.0123</b>	0.767	0.759	0.762	0.050	0.057	0.047	0.767	0.759	0.762	0.050	0.057	0.047
50	<b>0.2900</b>	<b>0.0181</b>	<b>0.0453</b>	0.766	0.759	0.761	0.062	0.069	0.058	0.766	0.760	0.761	0.062	0.068	0.057
25	<b>0.4488</b>	<b>0.0843</b>	<b>0.1505</b>	0.767	0.762	0.762	0.087	0.093	0.079	0.766	0.760	0.760	0.086	0.093	0.079

- ... aumenta o desvio padrão das pseudo-saídas...

# Testes estatísticos com tamanho da CD

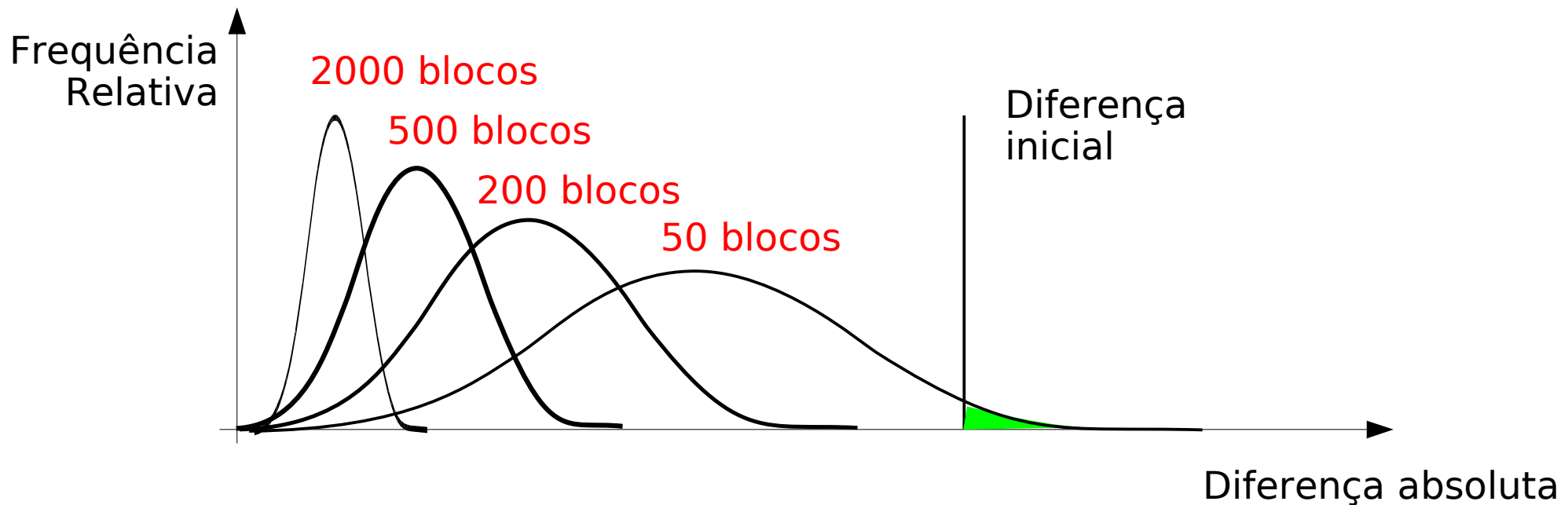
<i>Saída A</i>			<i>Saída B</i>			<i>Diferença</i>		
Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>0.728</b>	<b>0.696</b>	<b>0.712</b>	<b>0.798</b>	<b>0.870</b>	<b>0.832</b>	<b>0.069</b>	<b>0.174</b>	<b>0.121</b>

- ... e aumenta a média e desvio padrão das diferenças entre pseudo-saídas!

*PSEUDO-DIFERENÇAS*

# Blocos	Valor de P			Média			Desvio Padrão		
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>Todos</b>	0.0001	0.0001	0.0001	0.006	0.006	0.005	0.004	0.005	0.004
2000	0.0001	0.0001	0.0001	0.008	0.009	0.008	0.006	0.007	0.006
1000	0.0001	0.0001	0.0001	0.012	0.013	0.011	0.009	0.010	0.008
500	0.0004	0.0001	0.0001	0.017	0.018	0.015	0.013	0.014	0.011
250	<b>0.0181</b>	0.0001	0.0001	0.024	0.026	0.021	0.018	0.020	0.016
200	<b>0.0351</b>	0.0001	0.0001	0.026	0.029	0.024	0.020	0.022	0.018
100	<b>0.1391</b>	0.0009	0.0047	0.037	0.041	0.034	0.028	0.031	0.026
75	<b>0.1912</b>	0.0034	<b>0.0123</b>	0.043	0.048	0.039	0.032	0.036	0.029
50	<b>0.2900</b>	<b>0.0181</b>	<b>0.0453</b>	0.053	0.058	0.048	0.040	0.045	0.036
25	<b>0.4488</b>	<b>0.0843</b>	<b>0.1505</b>	0.073	0.081	0.066	0.056	0.061	0.051

# Variação do nº de blocos vs valor de P



**PSEUDO-DIFERENÇAS**

# Blocos	Valor de P			Média			Desvio Padrão		
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>Todos</b>	0.0001	0.0001	0.0001	0.006	0.006	0.005	0.004	0.005	0.004
2000	0.0001	0.0001	0.0001	0.008	0.009	0.008	0.006	0.007	0.006
1000	0.0001	0.0001	0.0001	0.012	0.013	0.011	0.009	0.010	0.008
500	0.0004	0.0001	0.0001	0.017	0.018	0.015	0.013	0.014	0.011
250	<b>0.0181</b>	0.0001	0.0001	0.024	0.026	0.021	0.018	0.020	0.016
200	<b>0.0351</b>	0.0001	0.0001	0.026	0.029	0.024	0.020	0.022	0.018
100	<b>0.1391</b>	0.0009	0.0047	0.037	0.041	0.034	0.028	0.031	0.026
75	<b>0.1912</b>	0.0034	<b>0.0123</b>	0.043	0.048	0.039	0.032	0.036	0.029
50	<b>0.2900</b>	<b>0.0181</b>	<b>0.0453</b>	0.053	0.058	0.048	0.040	0.045	0.036
25	<b>0.4488</b>	<b>0.0843</b>	<b>0.1505</b>	0.073	0.081	0.066	0.056	0.061	0.051

# Prova dos nove...

<b>Saída A</b>			<b>Saída B</b>			<b>Diferença</b>		
Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>0.728</b>	<b>0.696</b>	<b>0.712</b>	<b>0.798</b>	<b>0.870</b>	<b>0.832</b>	<b>0.069</b>	<b>0.174</b>	<b>0.121</b>

## PSEUDO-DIFERENÇAS

# Blocos	Valor de P			Média			Desvio Padrão		
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF
<b>Todos</b>	0.0001	0.0001	0.0001	0.006	0.006	0.005	0.004	0.005	0.004
2000	0.0001	0.0001	0.0001	0.008	0.009	0.008	0.006	0.007	0.006
1000	0.0001	0.0001	0.0001	0.012	0.013	0.011	0.009	0.010	0.008
500	0.0004	0.0001	0.0001	0.017	0.018	0.015	0.013	0.014	0.011
250	<b>0.0181</b>	0.0001	0.0001	0.024	0.026	0.021	0.018	0.020	0.016
200	<b>0.0351</b>	0.0001	0.0001	0.026	0.029	0.024	0.020	0.022	0.018
100	<b>0.1391</b>	0.0009	0.0047	0.037	0.041	0.034	0.028	0.031	0.026
75	<b>0.1912</b>	0.0034	<b>0.0123</b>	0.043	0.048	0.039	0.032	0.036	0.029
50	<b>0.2900</b>	<b>0.0181</b>	<b>0.0453</b>	0.053	0.058	0.048	0.040	0.045	0.036
25	<b>0.4488</b>	<b>0.0843</b>	<b>0.1505</b>	0.073	0.081	0.066	0.056	0.061	0.051

<b>Teste T-Student (99%) = 2,58</b>					
Média - t * desvPad			Média + t * desvPad		
-0.005	-0.006	-0.005	0.016	0.019	0.015
-0.008	-0.009	-0.007	0.024	0.028	0.022
-0.011	-0.013	-0.01	0.035	<b>0.039</b>	0.032
-0.016	-0.017	-0.014	0.050	0.054	0.044
-0.022	-0.025	-0.02	<b>0.069</b>	0.078	0.063
-0.025	-0.028	-0.023	<b>0.078</b>	0.086	0.070
-0.036	-0.039	-0.033	<b>0.110</b>	0.122	0.100
-0.041	-0.045	-0.037	<b>0.126</b>	0.140	0.114
-0.05	-0.057	-0.046	<b>0.155</b>	<b>0.173</b>	<b>0.142</b>
<b>-0.072</b>	-0.078	-0.065	<b>0.218</b>	<b>0.239</b>	<b>0.198</b>

# Influencia do nº de iterações no valor de P

<b>n blocos = 2000</b>				<b>PSEUDO-SAÍDAS de A</b>									<b>PSEUDO-SAÍDAS de B</b>									<b>PSEUDO-DIFERENÇAS</b>								
# Iter	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão			Média			Desvio Padrão											
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF									
<b>9999</b>	0.0001	0.0001	0.0001	0.765	0.772	0.769	0.008	0.009	0.007	0.765	0.772	0.769	0.008	0.009	0.007	0.00831	0.00937	0.00750	0.00626	0.00707	0.00570									
999	0.0010	0.0010	0.0010	0.765	0.772	0.768	0.008	0.009	0.007	0.765	0.771	0.768	0.008	0.009	0.007	0.00836	0.00930	0.00744	0.00609	0.00692	0.00562									
99	0.0100	0.0100	0.0100	0.764	0.771	0.768	0.009	0.009	0.007	0.766	0.772	0.769	0.009	0.009	0.008	0.00984	0.00952	0.00855	0.00750	0.00642	0.00609									
<b>n blocos = 200</b>				<b>PSEUDO-SAÍDAS de A</b>									<b>PSEUDO-SAÍDAS de B</b>									<b>PSEUDO-DIFERENÇAS</b>								
# Iter	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão			Média			Desvio Padrão											
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF									
<b>9999</b>	<b>0.0351</b>	0.0001	0.0001	0.765	0.759	0.762	0.030	0.034	0.028	0.765	0.759	0.762	0.030	0.035	0.028	0.02650	0.02904	0.02379	0.01991	0.02212	0.01804									
999	<b>0.0290</b>	0.0010	0.0010	0.767	0.758	0.762	0.030	0.036	0.029	0.767	0.759	0.762	0.030	0.035	0.028	0.02516	0.02896	0.02360	0.01866	0.02212	0.01753									
99	<b>0.0500</b>	0.0100	0.0100	0.772	0.760	0.766	0.032	0.035	0.029	0.767	0.760	0.763	0.032	0.039	0.032	0.02613	0.03271	0.02526	0.02005	0.02390	0.01957									
<b>n blocos = 25</b>				<b>PSEUDO-SAÍDAS de A</b>									<b>PSEUDO-SAÍDAS de B</b>									<b>PSEUDO-DIFERENÇAS</b>								
# Iter	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão			Média			Desvio Padrão											
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF									
<b>9999</b>	<b>0.4488</b>	<b>0.0843</b>	<b>0.1505</b>	0.767	0.762	0.762	0.087	0.093	0.079	0.766	0.760	0.760	0.086	0.093	0.079	0.07302	0.08053	0.06637	0.05618	0.06136	0.05085									
999	<b>0.4330</b>	<b>0.0930</b>	<b>0.1580</b>	0.766	0.761	0.761	0.086	0.094	0.079	0.764	0.761	0.760	0.089	0.096	0.082	0.07373	0.08299	0.06646	0.05823	0.06107	0.05393									
99	<b>0.4800</b>	<b>0.0800</b>	<b>0.1200</b>	0.762	0.764	0.760	0.090	0.088	0.077	0.765	0.763	0.762	0.088	0.096	0.082	0.07849	0.08180	0.06811	0.05604	0.06272	0.04987									

# Uma experiência com n<sup>o</sup> iterações

n blocos = 2000				PSEUDO-SAÍDAS de A									PSEUDO-SAÍDAS de B									PSEUDO-DIFERENÇAS								
# Iter	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão			Média			Desvio Padrão											
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF									
9999	0.0001	0.0001	0.0001	0.765	0.772	0.769	0.008	0.009	0.007	0.765	0.772	0.769	0.008	0.009	0.007	0.00831	0.00937	0.00750	0.00626	0.00707	0.00570									
999	0.0010	0.0010	0.0010	0.765	0.772	0.768	0.008	0.009	0.007	0.765	0.771	0.768	0.008	0.009	0.007	0.00836	0.00930	0.00744	0.00609	0.00692	0.00562									
99	0.0100	0.0100	0.0100	0.764	0.771	0.768	0.009	0.009	0.007	0.766	0.772	0.769	0.009	0.009	0.008	0.00984	0.00952	0.00855	0.00750	0.00642	0.00609									

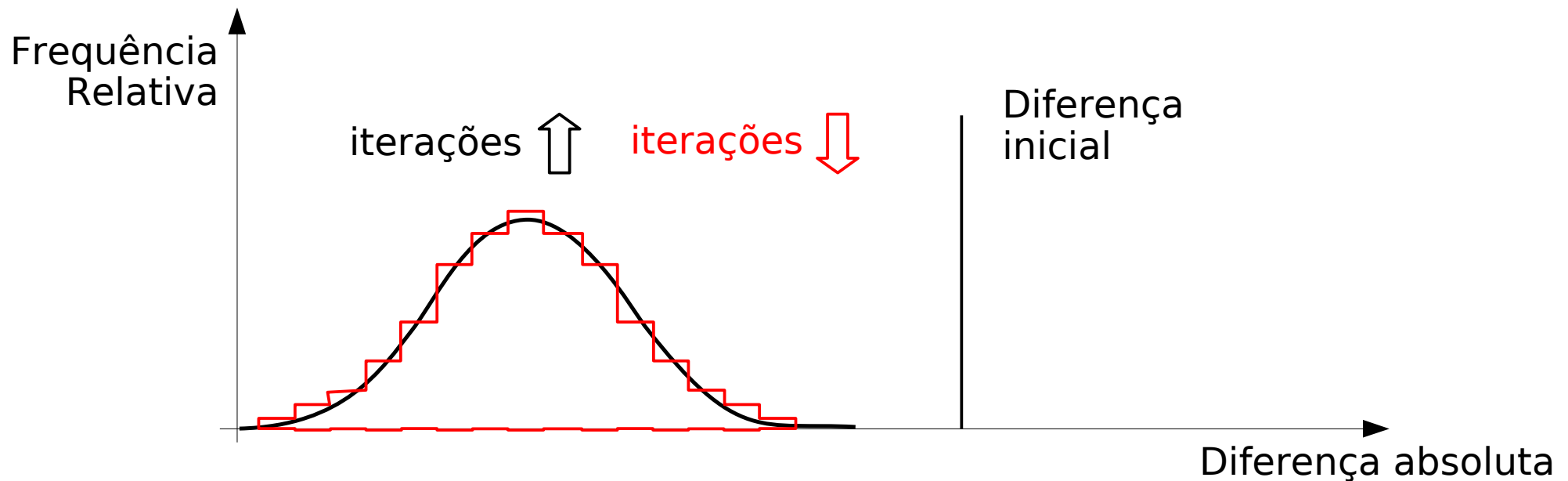
  

n blocos = 200				PSEUDO-SAÍDAS de A									PSEUDO-SAÍDAS de B									PSEUDO-DIFERENÇAS								
# Iter	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão			Média			Desvio Padrão											
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF									
9999	<b>0.0351</b>	0.0001	0.0001	0.765	0.759	0.762	0.030	0.034	0.028	0.765	0.759	0.762	0.030	0.035	0.028	0.02650	0.02904	0.02379	0.01991	0.02212	0.01804									
999	<b>0.0290</b>	0.0010	0.0010	0.767	0.758	0.762	0.030	0.036	0.029	0.767	0.759	0.762	0.030	0.035	0.028	0.02516	0.02896	0.02360	0.01866	0.02212	0.01753									
99	<b>0.0500</b>	0.0100	0.0100	0.772	0.760	0.766	0.032	0.035	0.029	0.767	0.760	0.763	0.032	0.039	0.032	0.02613	0.03271	0.02526	0.02005	0.02390	0.01957									

n blocos = 25				PSEUDO-SAÍDAS de A									PSEUDO-SAÍDAS de B									PSEUDO-DIFERENÇAS								
# Iter	Valor de P			Média			Desvio Padrão			Média			Desvio Padrão			Média			Desvio Padrão											
	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF	Prec	Abr	MedF									
9999	<b>0.4488</b>	<b>0.0843</b>	<b>0.1505</b>	0.767	0.762	0.762	0.087	0.093	0.079	0.766	0.760	0.760	0.086	0.093	0.079	0.07302	0.08053	0.06637	0.05618	0.06136	0.05085									
999	<b>0.4330</b>	<b>0.0930</b>	<b>0.1580</b>	0.766	0.761	0.761	0.086	0.094	0.079	0.764	0.761	0.760	0.089	0.096	0.082	0.07373	0.08299	0.06646	0.05823	0.06107	0.05393									
99	<b>0.4800</b>	<b>0.0800</b>	<b>0.1200</b>	0.762	0.764	0.760	0.090	0.088	0.077	0.765	0.763	0.762	0.088	0.096	0.082	0.07849	0.08180	0.06811	0.05604	0.06272	0.04987									

- Sem grandes diferenças. Média e desvio padrão na mesma.

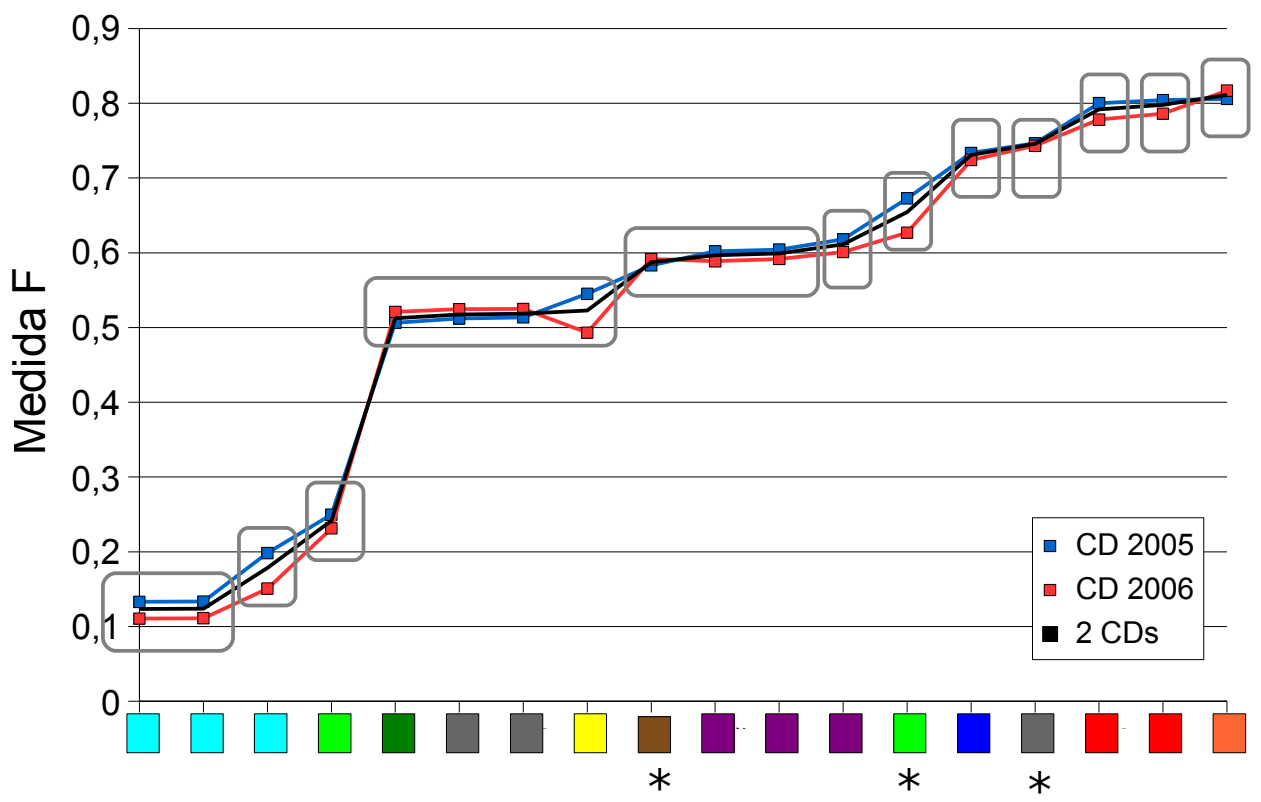
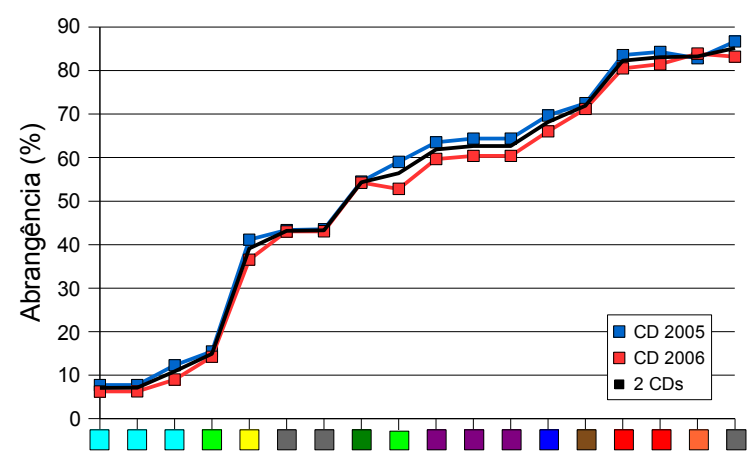
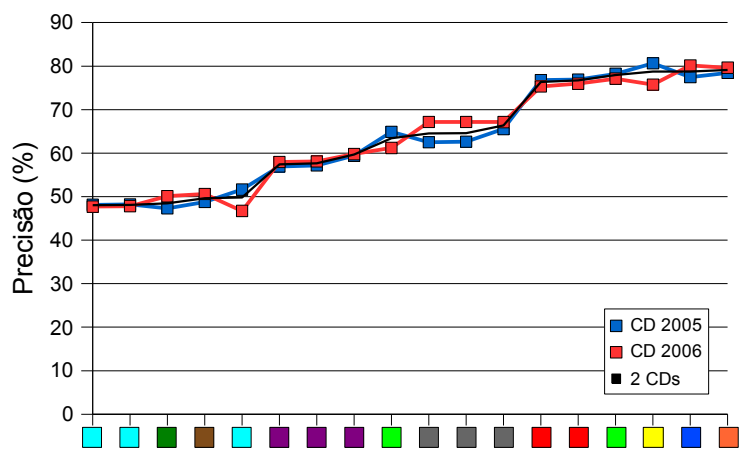


# Resultados finais



# HAREM 2005

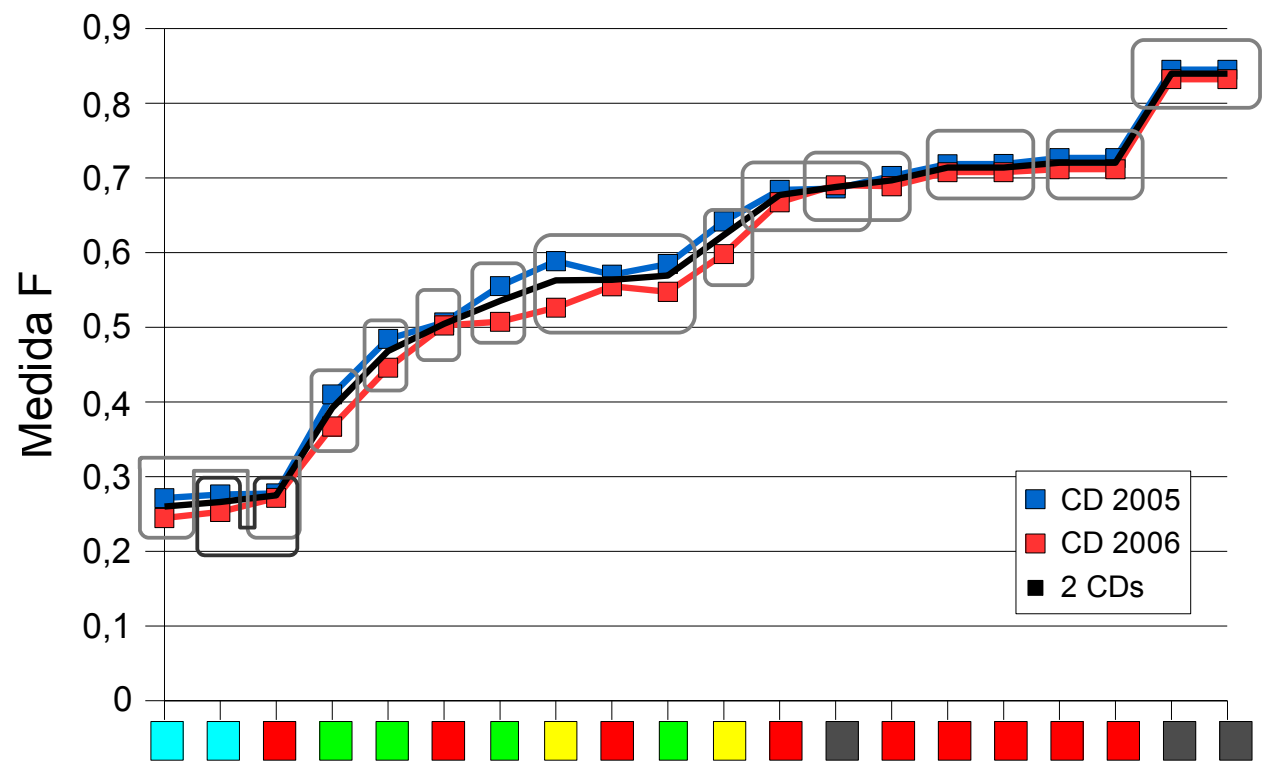
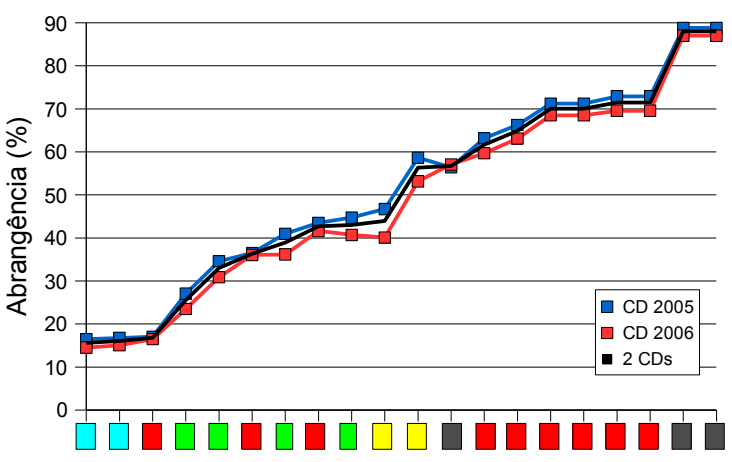
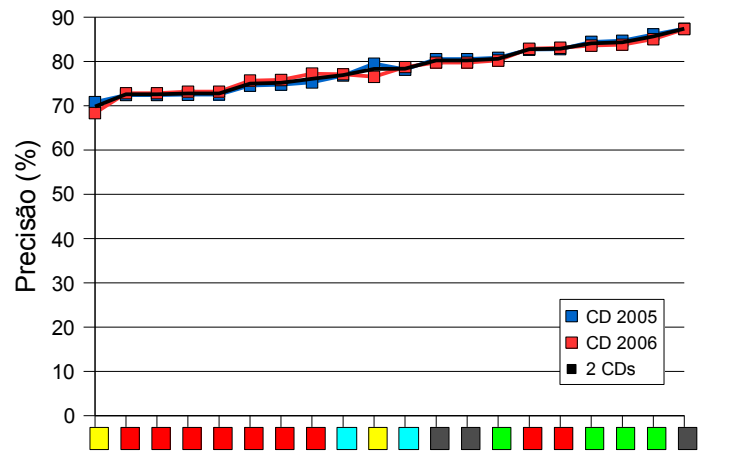
## Tarefa de Identificação



\* - Saídas não oficiais

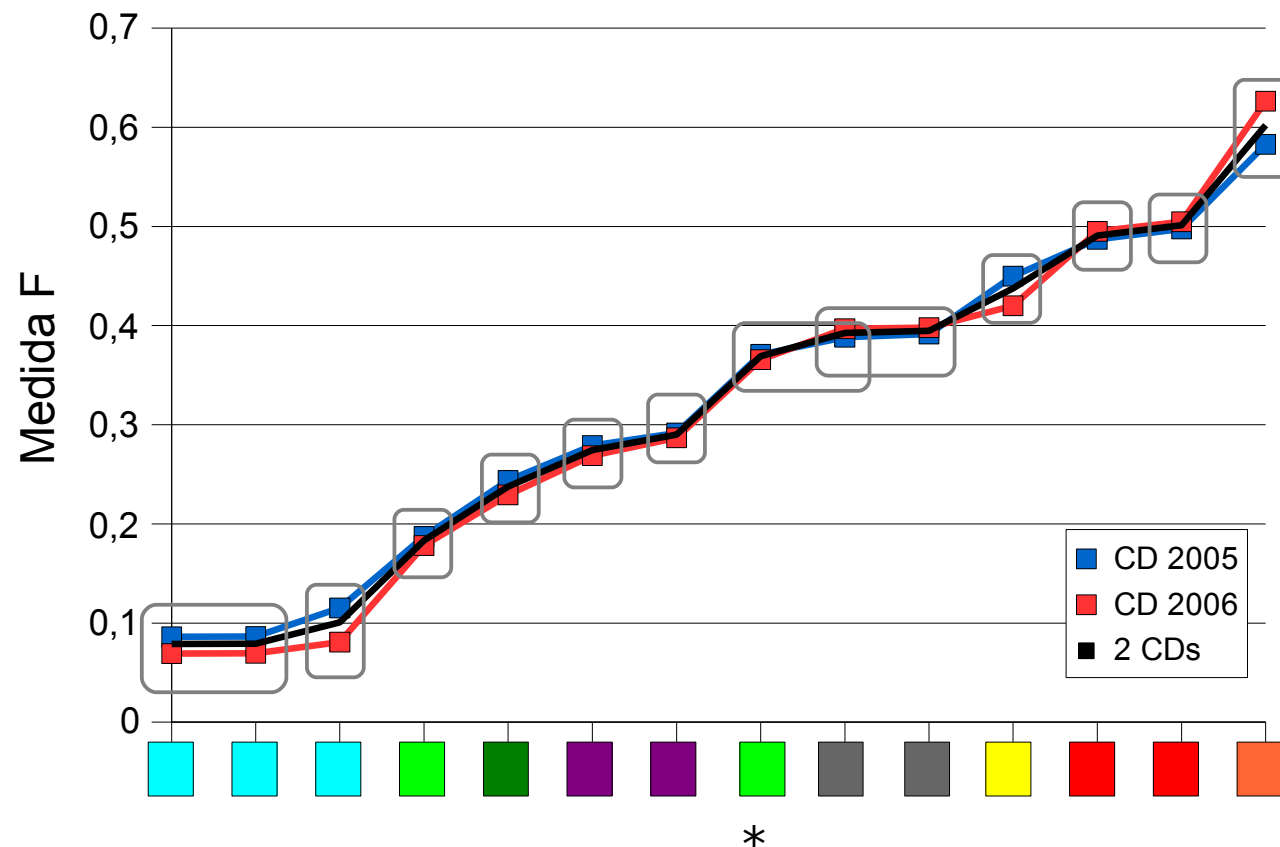
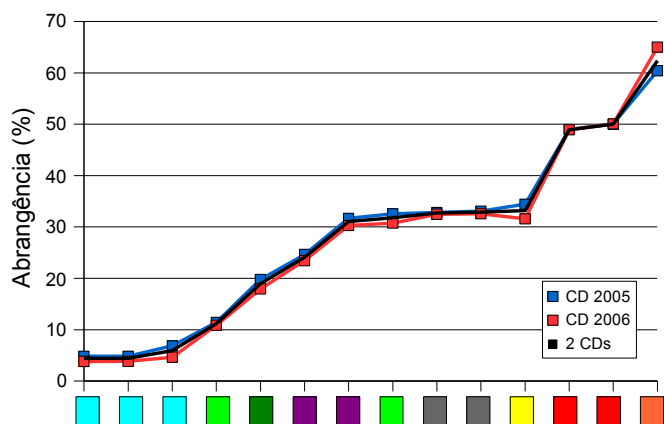
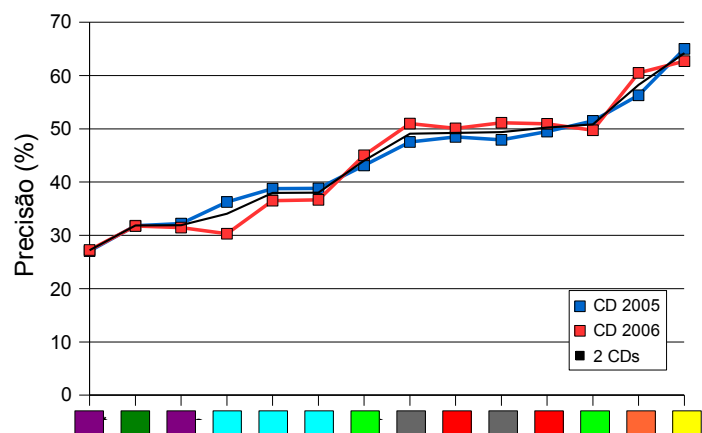
# MiniHAREM 2006

## Tarefa de Identificação



# HAREM 2005

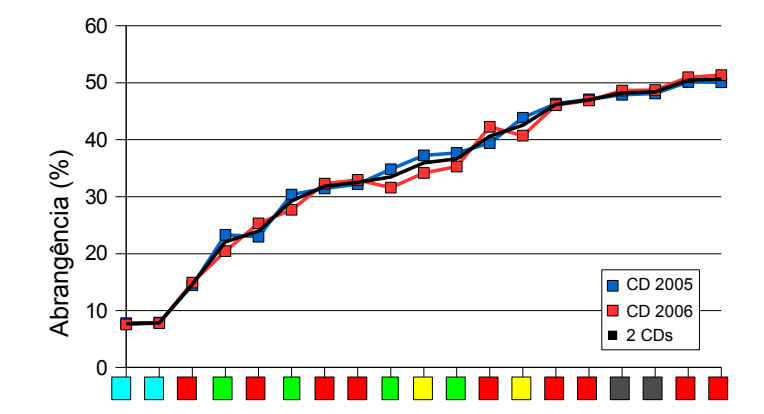
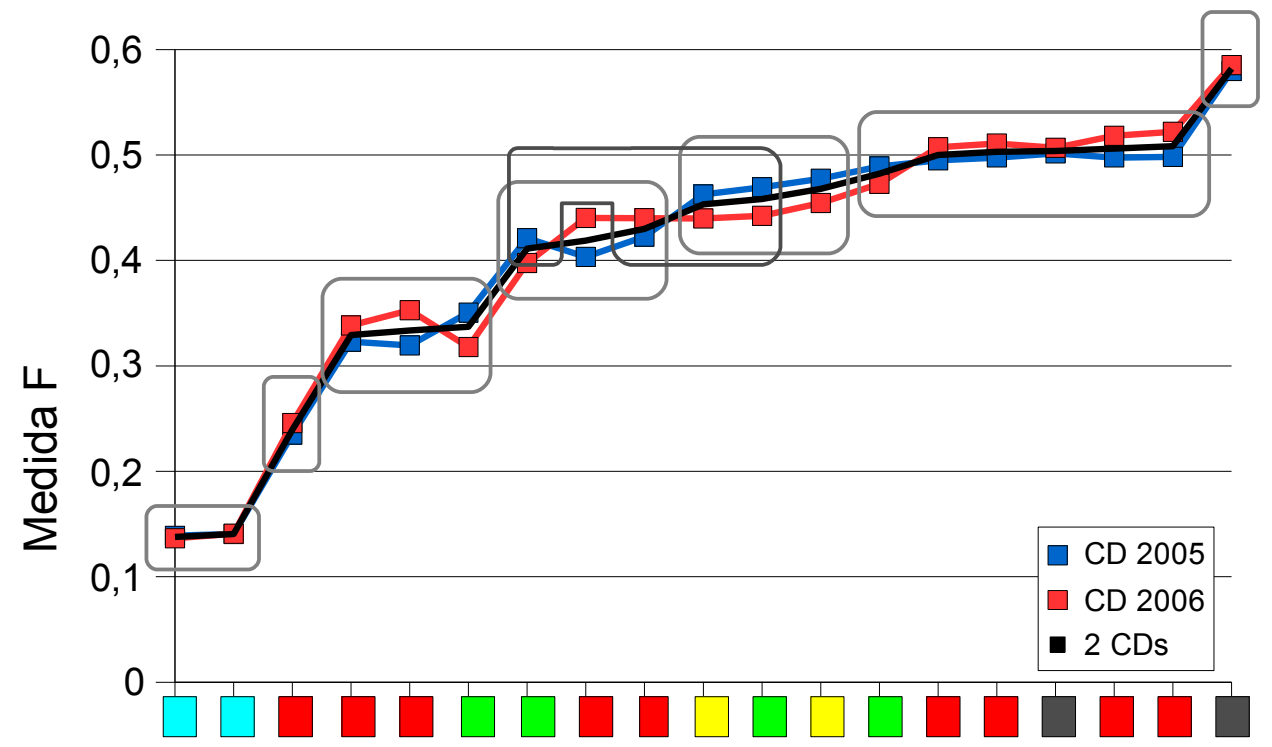
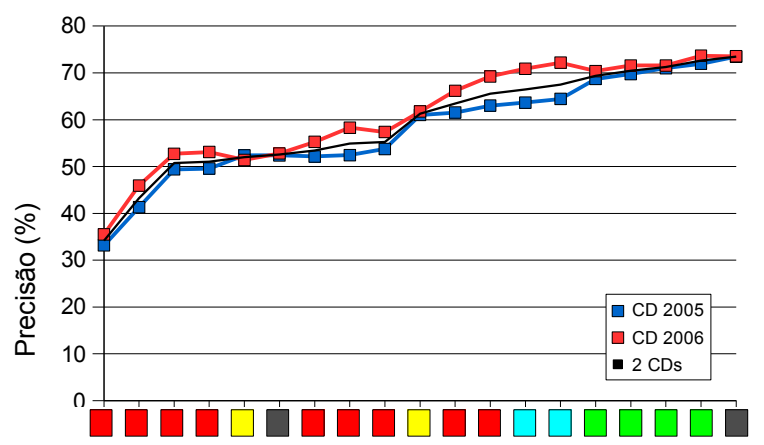
## Tarefa de **Classificação Semântica** (Combinada)



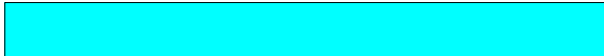


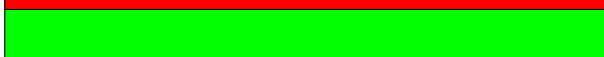

\* - Saída não-oficial

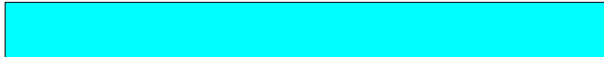


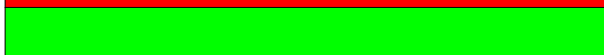

# MiniHAREM 2006

## Tarefa de Classificação Semântica (Combinada)



# Evolução dos Sistemas, em 1 ano?

Identificação	2005	2006	%
	0,178	0,266	49,4%
	0,523	0,624	19,2%
	0,798	0,720	-9,7%
	0,655	0,569	-13,0%
	0,746	<b>0,839</b>	12,6%

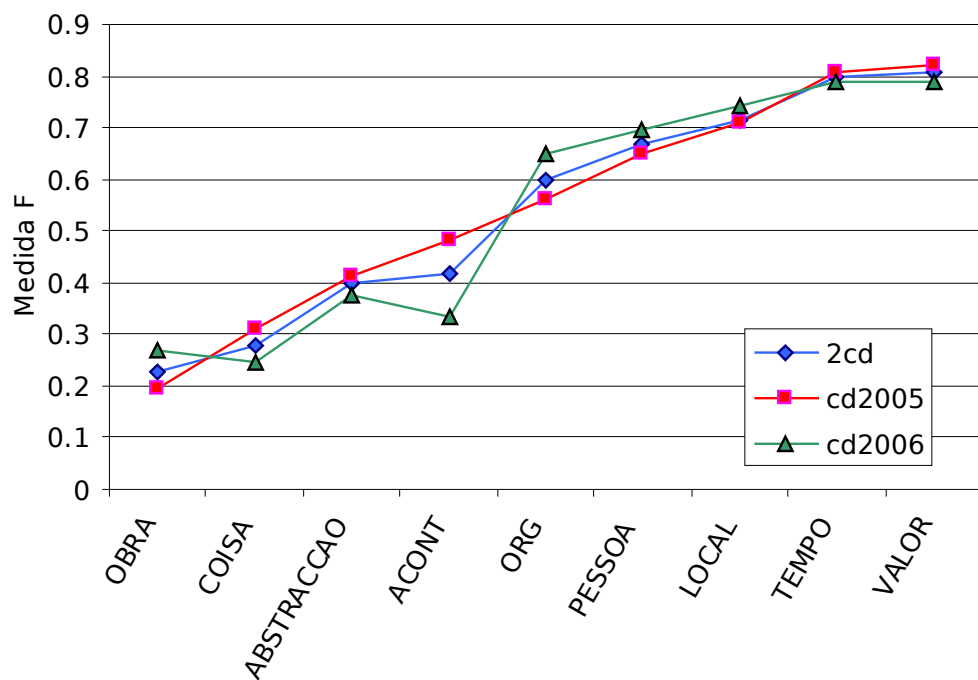
Semântica (CSC)	2005	2006	%
	0,101	0,141	39,8%
	0,438	0,468	6,9%
	0,501	0,508	1,4%
	0,369	0,482	30,7%
	0,395	<b>0,582</b>	47,5%

# Estado da Arte em REM: Categorias

# Panorama de Identificação, por Categorias

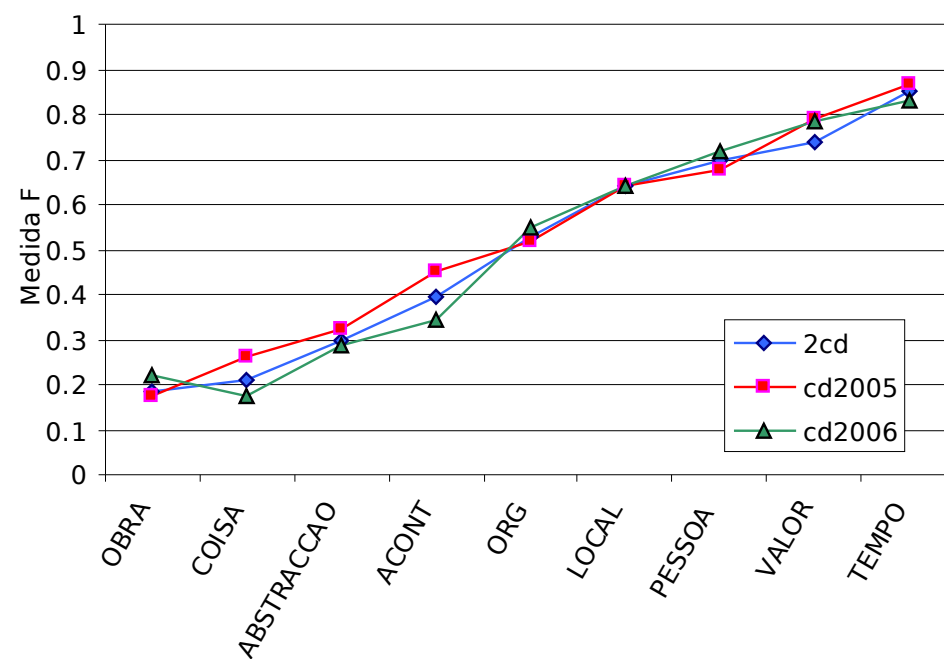
## 2005

Identificação por Categorias (HAREM)



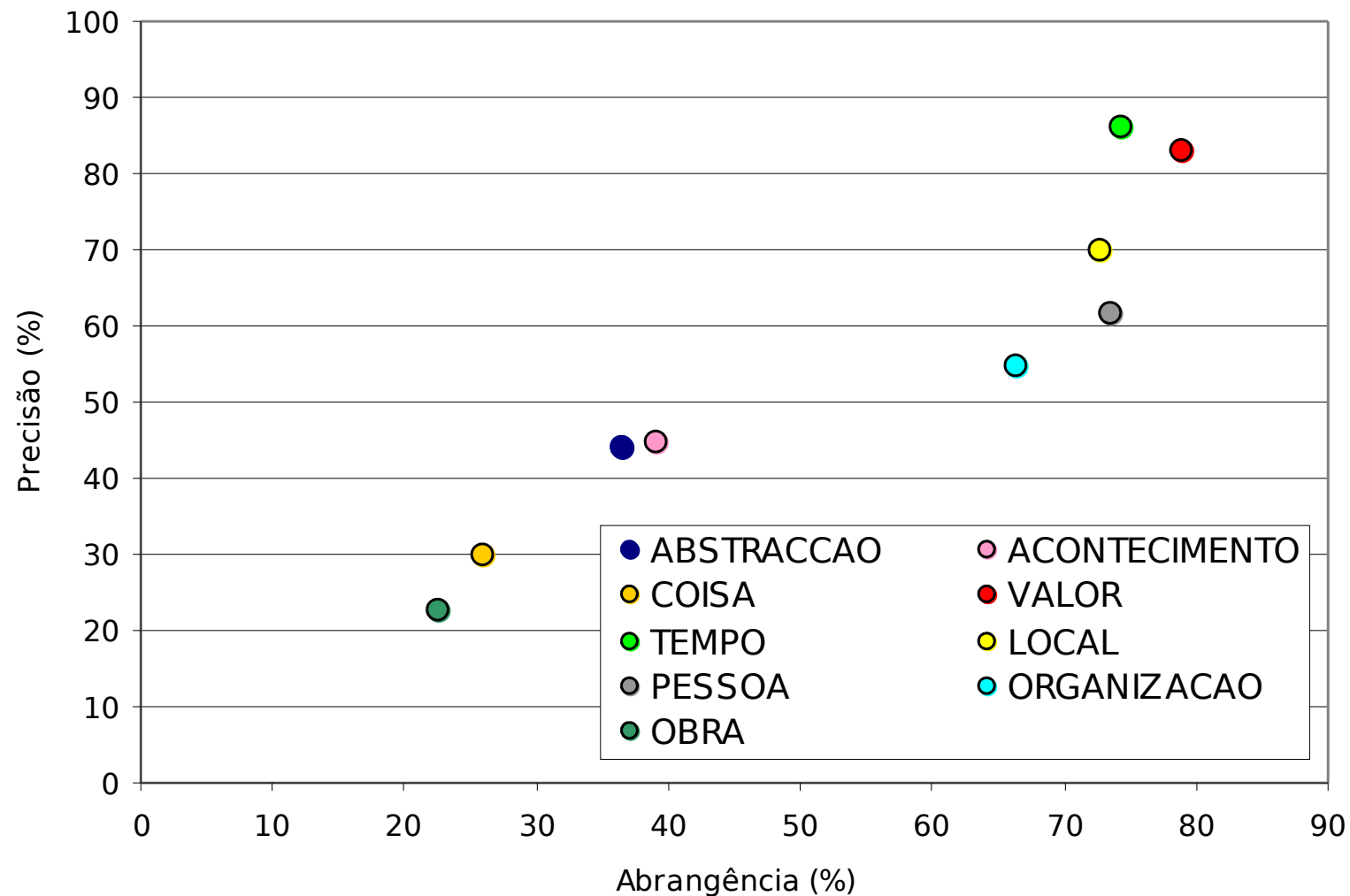
## 2006

Identificação por Categorias (MiniHAREM)



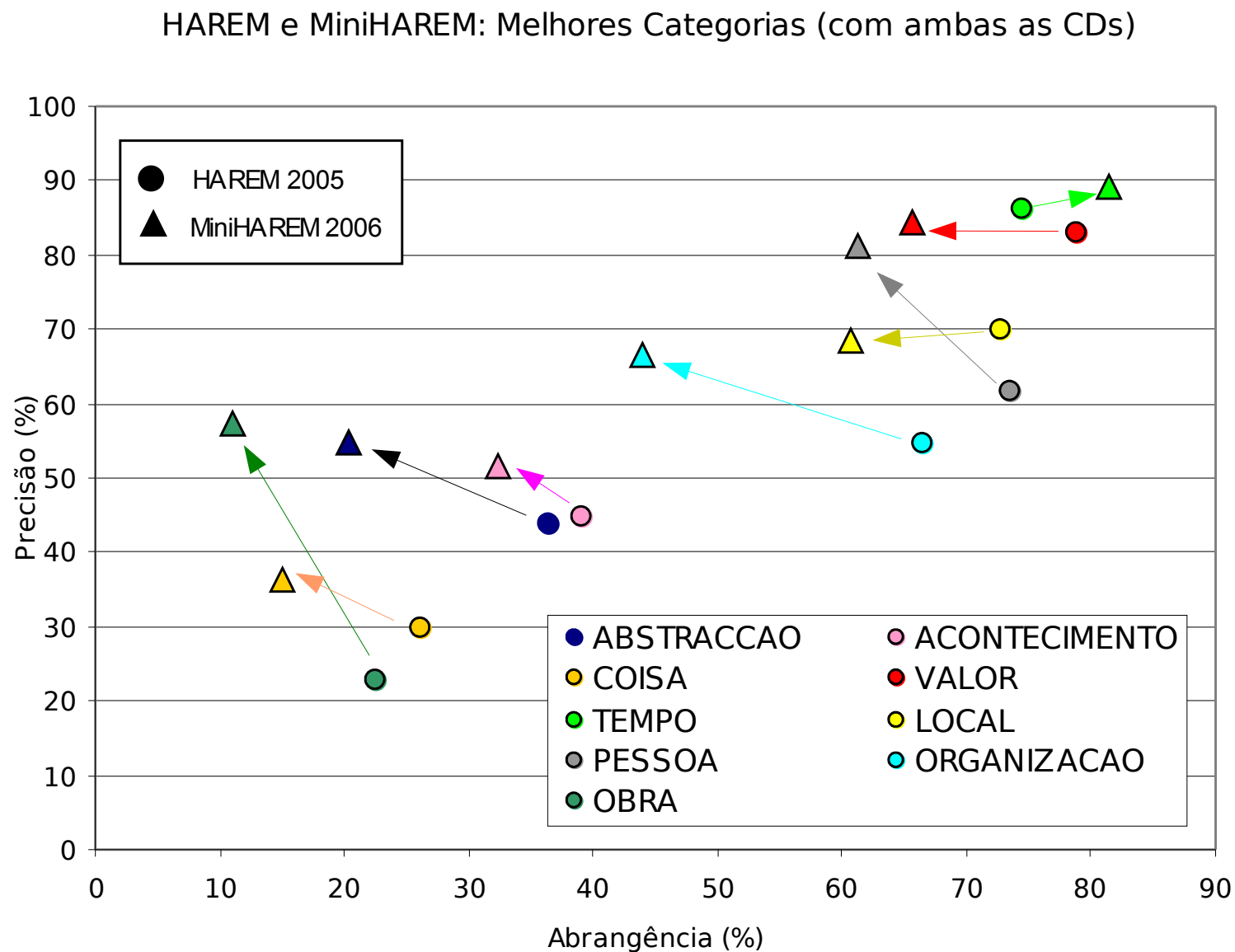
# Panorama de Identificação, por Categorias

HAREM: Melhores Categorias (com ambas as CDs)





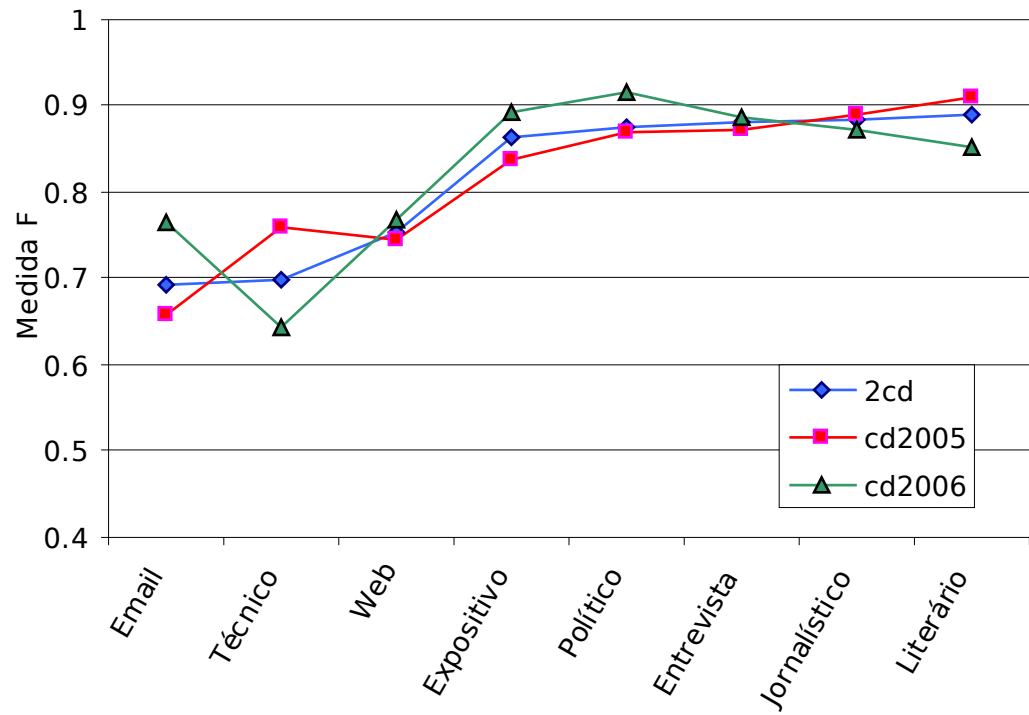
# Panorama de Identificação, por Categorias



# Panorama de Identificação, por Género

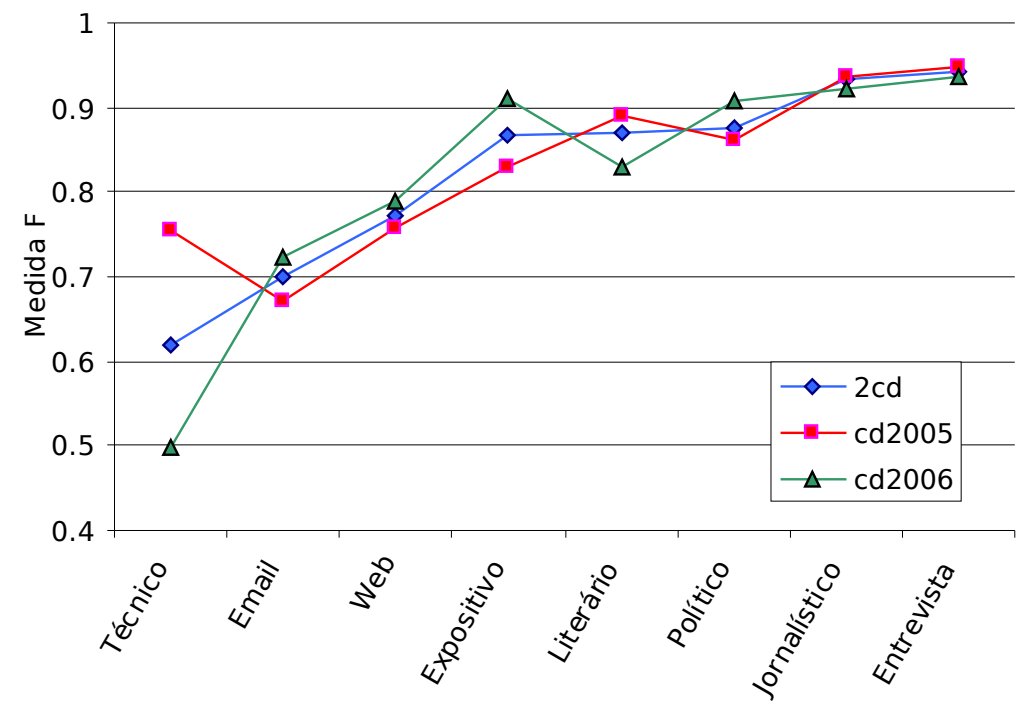
## 2005

Identificação por Género (HAREM)



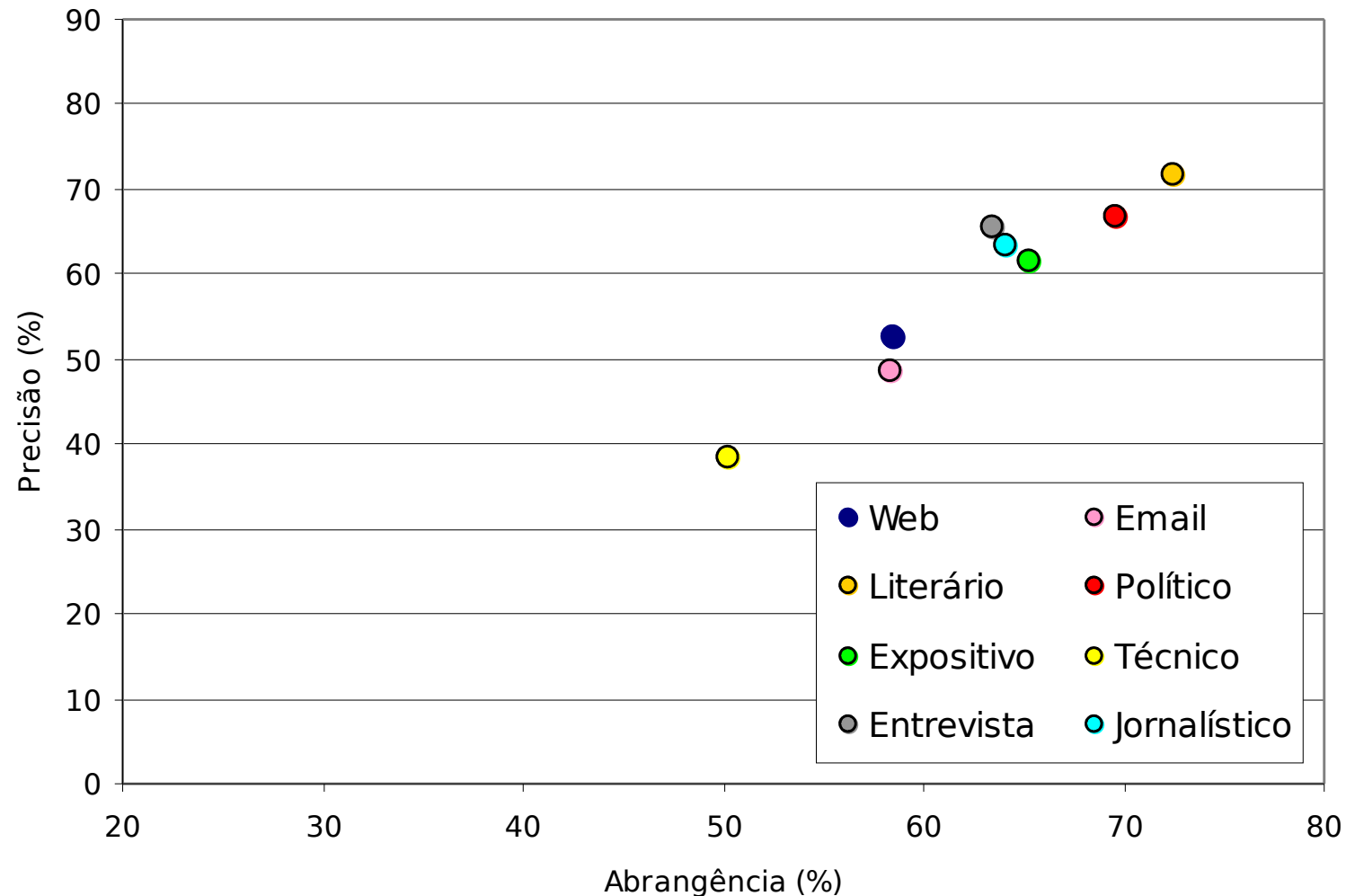
## 2006

Identificação por Género (HAREM)

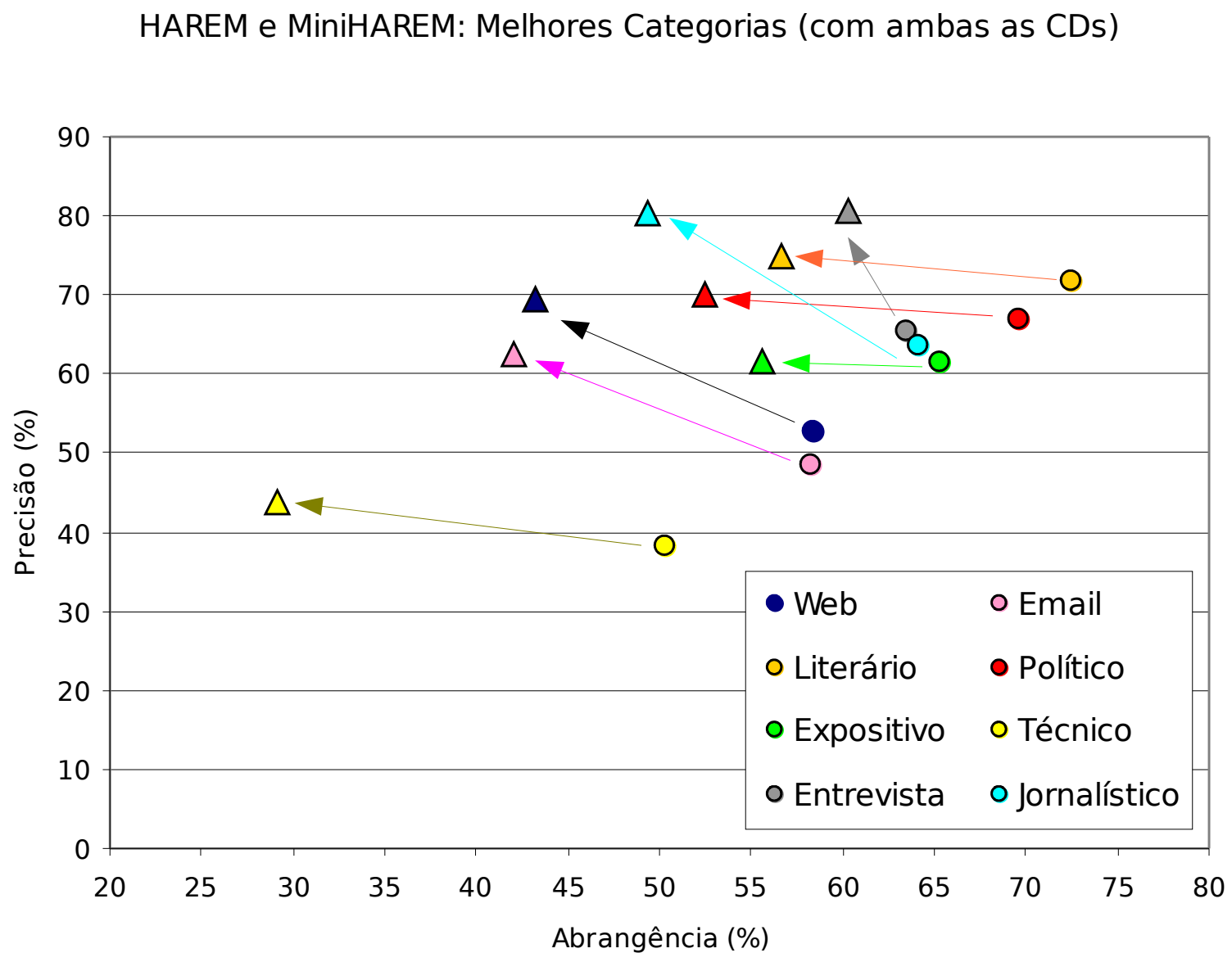


# Panorama de Identificação, por Género

HAREM: Melhores Géneros (com ambas as CDs)



# Panorama de Identificação, por Género



# Conclusões

- HAREM e MiniHAREM: São comparáveis
  - CDs são semelhantes
  - Produzem desempenhos semelhantes
  - Tamanho da CD mais que adequado
- A avaliação HAREM está validada.
- As melhores estratégias para REM foram aferidas e comparadas
- Os sistemas REM melhoraram em 1 ano!
- REM em PT com futuro risonho.

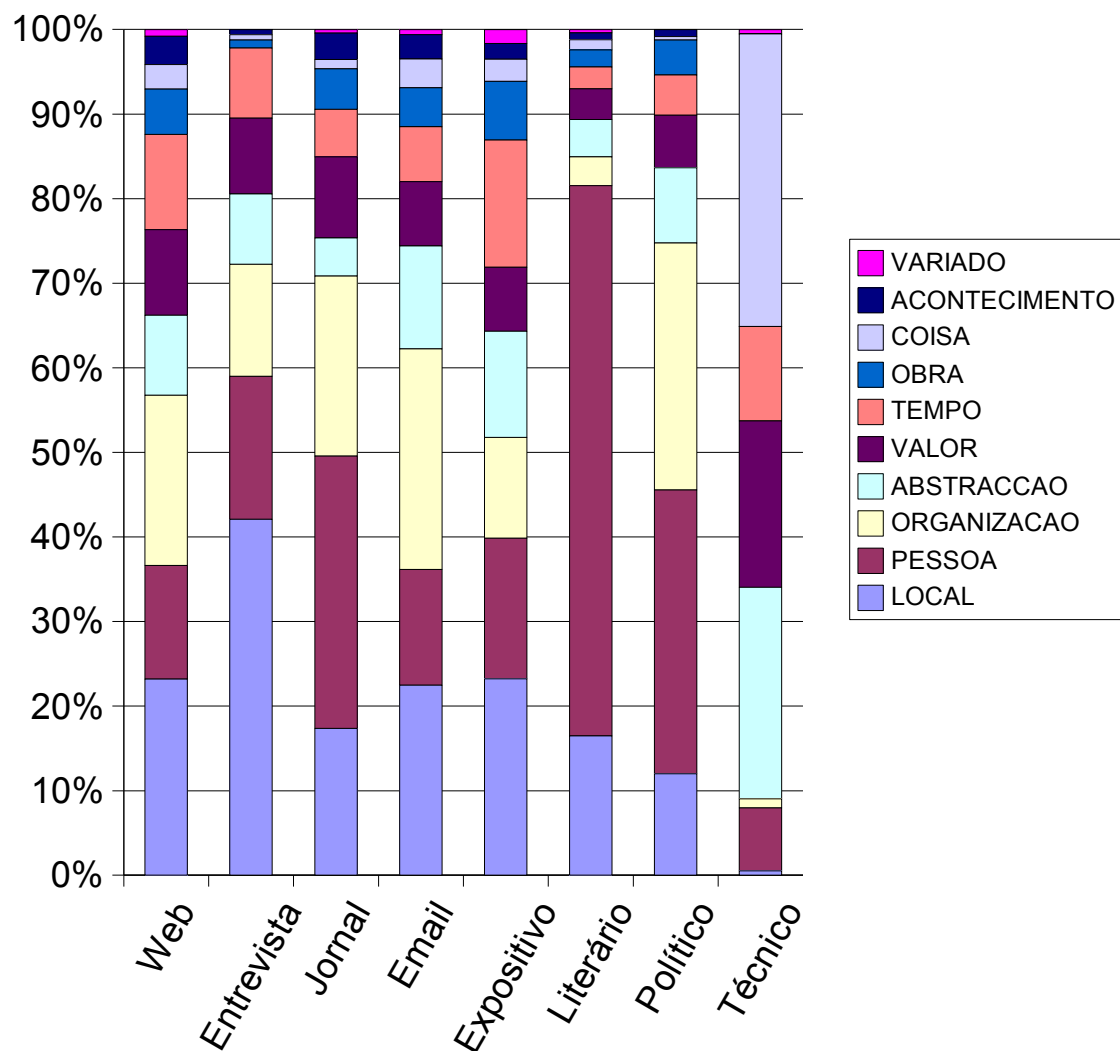
# Fim



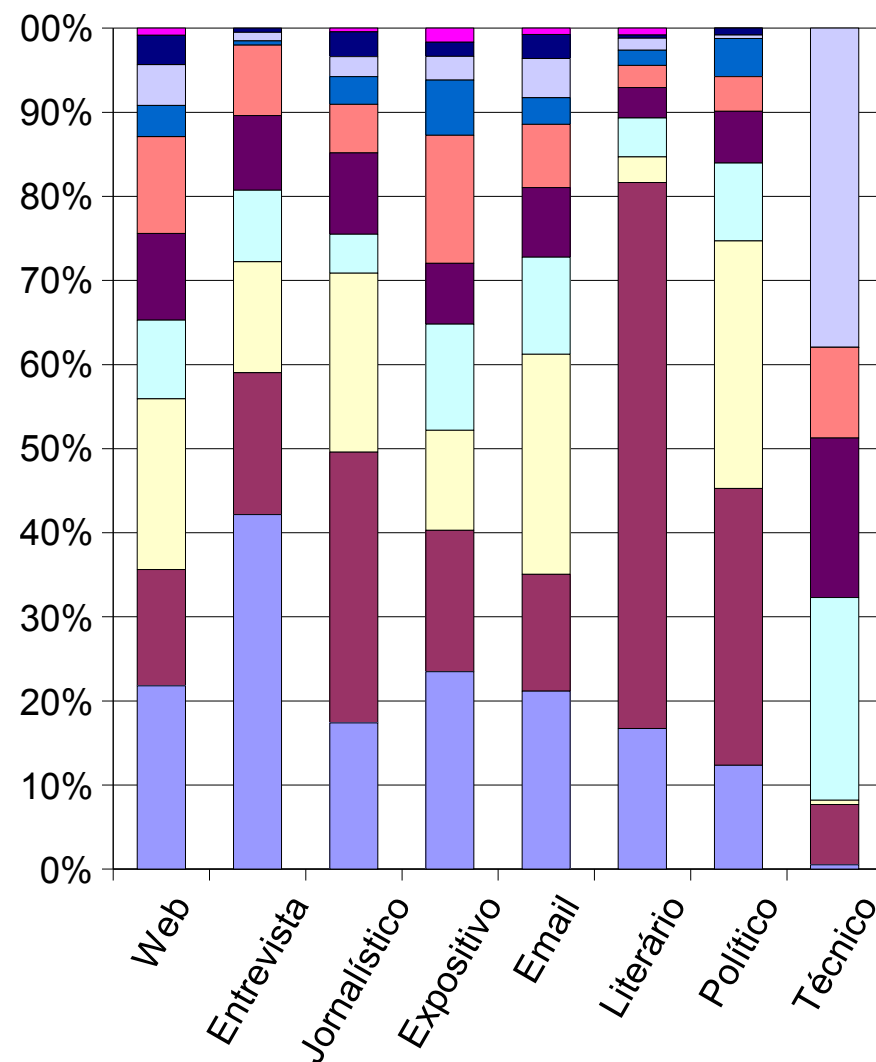
Obrigado pela atenção.

# Distribuição de Categorias, por Género Textual

## Regras de 2005

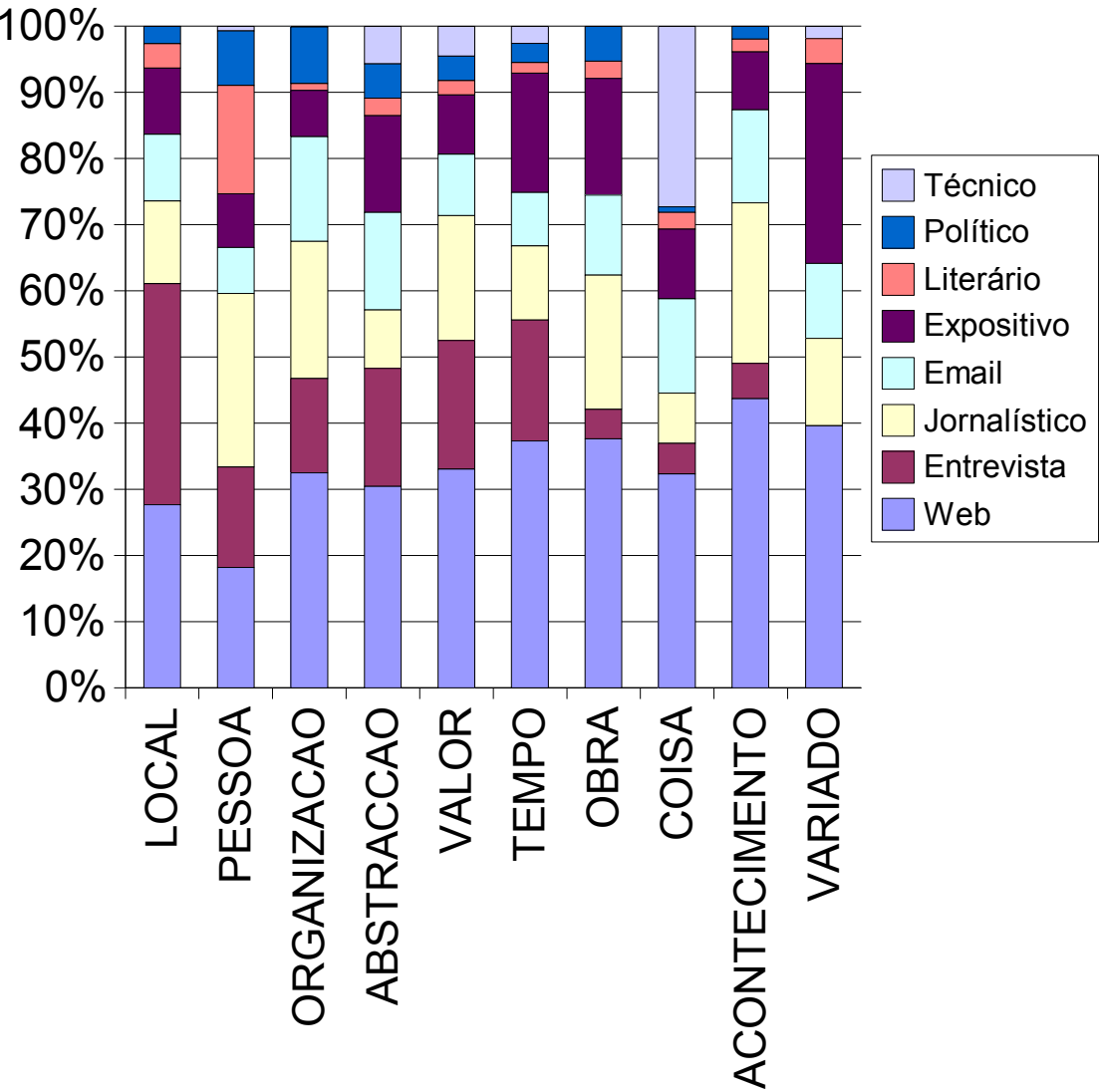


## Regras de 2006



# Distribuição de Género Textual, por Categoria

### Regras de 2005



### Regras de 2006

