



Motores de busca

João Pedro Gonçalves
jp@sapo.pt
Abril 2010





A Pesquisa do SAPO

- 50 Queries por Segundo (QPS).
- Acordo com Microsoft – Bing
- Indexador FAST com cerca de 100 milhões de páginas indexadas.



[Web](#) | [Imagens](#) | [Notícias](#) | [Blogs](#) | [PAi](#) | [PBi](#) | [Directório](#) | [Produtos](#) | [mais...](#)

Pesquisar: páginas de Portugal em língua portuguesa toda a web

[Portal SAPO](#) | [Registe o seu site](#) | [Anuncie na Pesquisa](#) | [Blog](#) | [Ajuda?](#)





A Pesquisa do SAPO

- Frontend em mod_perl
- Perl, Python, C e C++
- Equipa de 14 programadores nas áreas de Pesquisa, Anúncios contextualizados e Directório



[Web](#) | [Imagens](#) | [Notícias](#) | [Blogs](#) | [PAi](#) | [PBi](#) | [Directório](#) | [Produtos](#) | [mais...](#)

[Portal SAPO](#) | [Registe o seu site](#) | [Anuncie na Pesquisa](#) | [Blog](#) | [Ajuda?](#)

Directório ✕

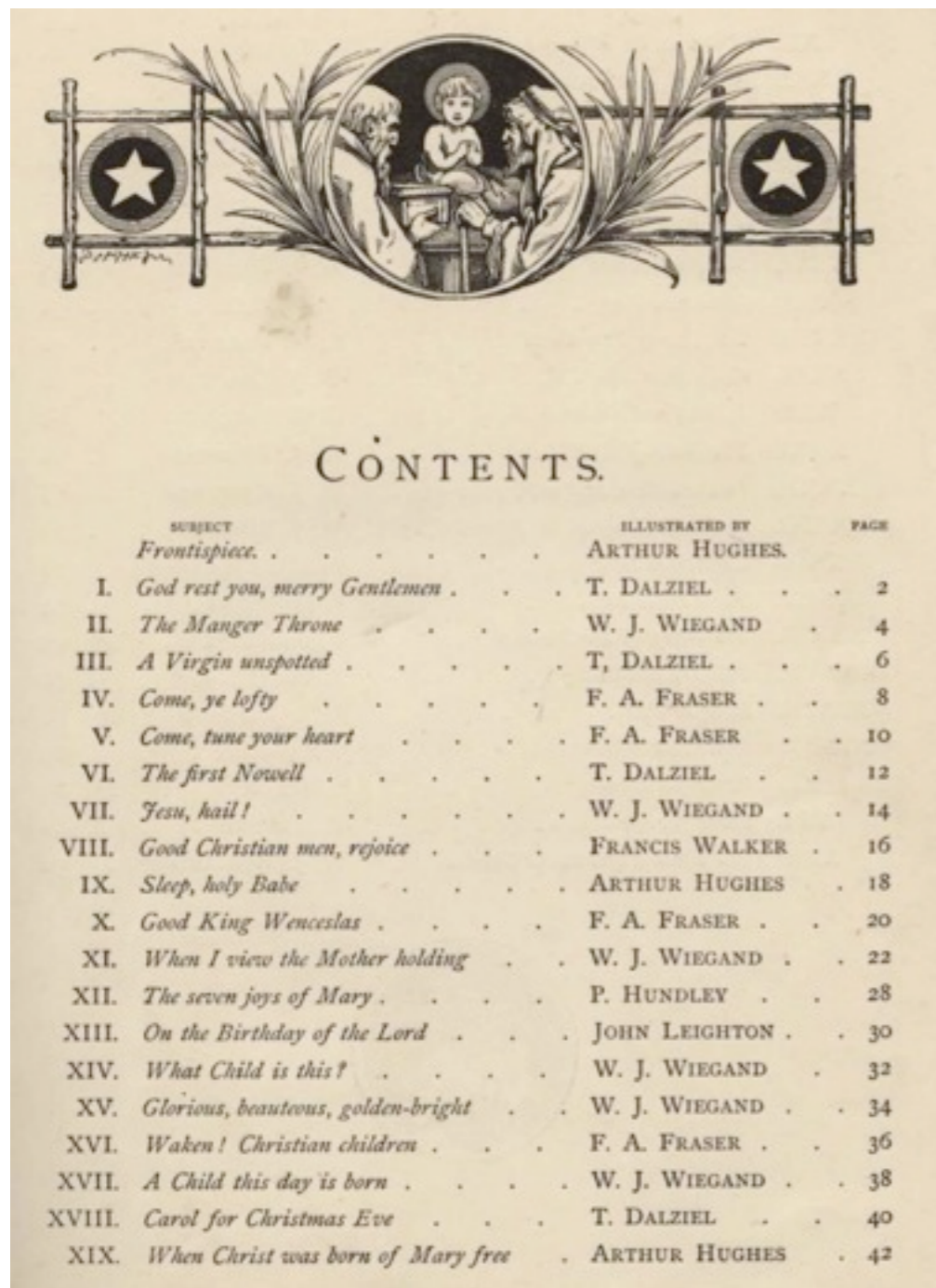
<ul style="list-style-type: none">• Arte e cultura Cinema, Teatro, Fotografia• Ciência Astronomia, Gestão, Medicina• Comunicação Social Jornais, Televisões, Revistas• Desporto Futebol, Motorizados, Aquáticos	<ul style="list-style-type: none">• Economia e Negócios Imobiliário, Emprego, Comércio• Educação Superior, Profissional, E-learning• Entretenimento e Lazer Cinema, Desporto, Humor• Estado e Administração Governo, ONGs, UE	<ul style="list-style-type: none">• Nacional e Regional Lisboa, Porto, Faro• Saúde Hospitais, Farmácias, Crianças• Sociedade Política, Religião, Minorias• Tecnologia e Internet Mail, Software, Biotecnologia
--	--	---

PT.COM © 2007



Directórios

por categorias



[Web](#) | [Imagens](#) | [Notícias](#) | [Blogs](#) | [PAi](#) | [PBi](#) | [Directório](#) | [Produtos](#) | [mais...](#)

[Portal SAPO](#) | [Registe o seu site](#) | [Anuncie na Pesquisa](#) | [Blog](#) | [Ajuda?](#)

✕

Directório

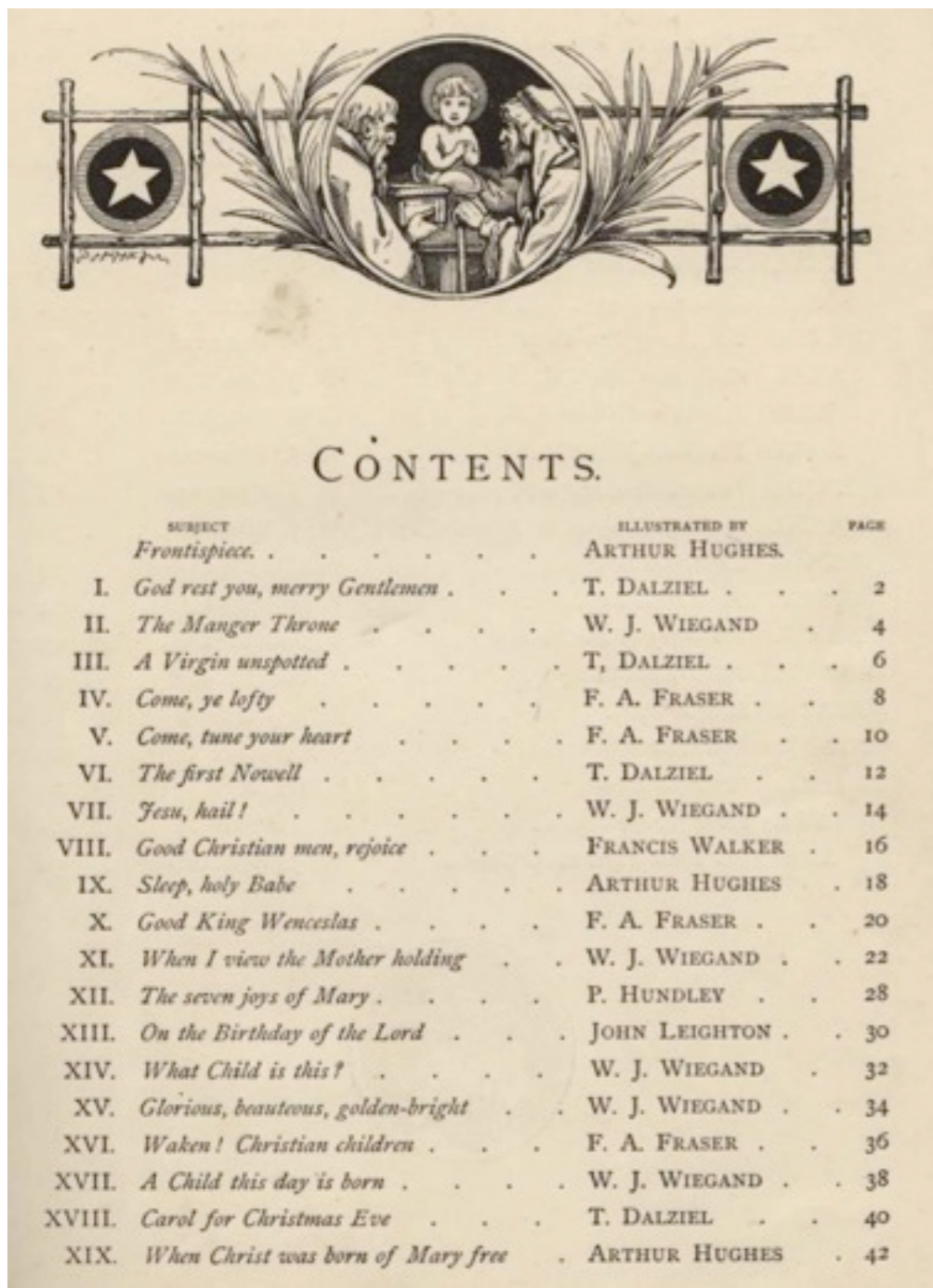
- [Arte e cultura](#)
[Cinema](#), [Teatro](#), [Fotografia](#)
- [Economia e Negócios](#)
[Imobiliário](#), [Emprego](#), [Comércio](#)
- [Nacional e Regional](#)
[Lisboa](#), [Porto](#), [Faro](#)
- [Ciência](#)
[Astronomia](#), [Gestão](#), [Medicina](#)
- [Educação](#)
[Superior](#), [Profissional](#), [E-learning](#)
- [Saúde](#)
[Hospitais](#), [Farmácias](#), [Crianças](#)
- [Comunicação Social](#)
[Jornais](#), [Televisões](#), [Revistas](#)
- [Entretenimento e Lazer](#)
[Cinema](#), [Desporto](#), [Humor](#)
- [Sociedade](#)
[Política](#), [Religião](#), [Minorias](#)
- [Desporto](#)
[Futebol](#), [Motorizados](#), [Aquáticos](#)
- [Estado e Administração](#)
[Governo](#), [ONGs](#), [UE](#)
- [Tecnologia e Internet](#)
[Mail](#), [Software](#), [Biotecnologia](#)

PT.COM © 2007



Directórios

por categorias



pesquisa

Web | Imagens | Notícias | Blogs | PAi | PBi | **Directório** | Produtos | mais...

Pesquisar

Portal SAPO | Registe o seu site | Anuncie na Pesquisa | Blog | Ajuda?

Directório

- [Arte e cultura](#)
Cinema, Teatro, Fotografia
- [Ciência](#)
Astronomia, Gestão, Medicina
- [Comunicação Social](#)
Jornais, Televisões, Revistas
- [Desporto](#)
Futebol, Motorizados, Aquáticos
- [Economia e Negócios](#)
Imobiliário, Emprego, Comércio
- [Educação](#)
Superior, Profissional, E-learning
- [Entretenimento e Lazer](#)
Cinema, Desporto, Humor
- [Estado e Administração](#)
Governo, ONGs, UE
- [Nacional e Regional](#)
Lisboa, Porto, Faro
- [Saúde](#)
Hospitais, Farmácias, Crianças
- [Sociedade](#)
Política, Religião, Minorias
- [Tecnologia e Internet](#)
Mail, Software, Biotecnologia

PT.COM © 2007

portugal gdp

Input interpretation:

Mathematica form

Portugal GDP

Result:

\$243.5 billion per year (US dollars per year) (2008 estimate)

Local currency conversion:

More

€ 178.1 billion per year (euros per year) (at current quoted rate)

GDP history:

Linear scale



(from 1970 to 2008)
(in billions of US dollars per year)
(log scale)

Economic properties:

More

GDP at exchange rate	\$243.5 billion per year (2008) (world rank: 38 th)
GDP at parity	\$237.3 billion per year (2008) (world rank: 47 th)
GDP in local currency	€ 166.2 billion (2008)
GDP per capita	\$22 805.48 per person (2008) (world rank: 55 th)
GDP real growth	-0.04466% per year
inflation rate	1.914% per year
unemployment rate	7.6% (2008) (world rank: 106 th)

Units »

Computed by: Wolfram Mathematica

Source information »

Download as: PDF | Live Mathematica

portugal gdp vs italy vs greece vs spain

Input interpretation:

Mathematica form

Portugal

Italy

Greece

Spain

GDP

Results:

Portugal	\$243.5 billion per year (US dollars per year) (2008 estimate)
Italy	\$2.303 trillion per year (US dollars per year) (2008 estimate)
Greece	\$355.9 billion per year (US dollars per year) (2008 estimate)
Spain	\$1.604 trillion per year (US dollars per year) (2008 estimate)

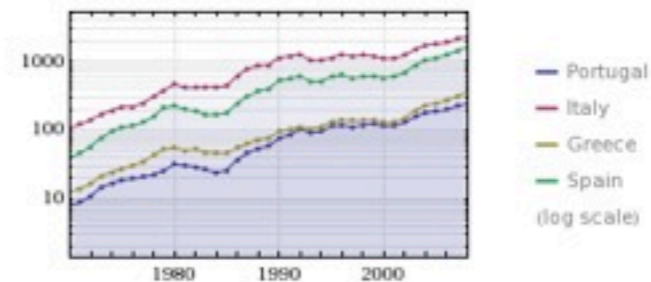
Ranked values:

Reverse

		visual	ratios	
1	Italy		9.458	1
2	Spain		6.588	0.6966
3	Greece		1.462	0.1545
4	Portugal		1	0.1057

GDP history:

Linear scale



(from 1970 to 2008) (in billions of US dollars per year)

Computed by: Wolfram Mathematica

Source information »

Download as: PDF | Live Mathematica

100,000 characters in Portuguese



Input interpretation:

100 000 characters (in Portuguese)

Approximate word count:

16 260 words

(assuming 5.15 characters per word and 1 space between each word)

Printed length:

[More](#) | [Weight, thickness](#) | [Show assumptions](#)

single-spaced document	33 pages	1476 lines
double-spaced document	66 pages	1476 lines

Times:

[More](#)

typical typing	5 hours
typical speaking	110 minutes
silent reading	59 minutes

Data size:

100 kB (kilobytes) | 800 kb (kilobits)
(assuming 8-bit encoding)

Language translation lengths:

English	94 000 characters
French	100 000 characters
German	100 000 characters
Spanish	100 000 characters
Italian	110 000 characters
Mandarin Chinese	50 000 characters

(estimates based on typical translations)

Computed by: [Wolfram Mathematica](#) | [Source information](#) | [Download as: PDF](#) | [Live Mathematica](#)

C major 7th chord



Input interpretation:

C major seventh chord

Music notation:

[Play notes](#)



Note names:

C E G B

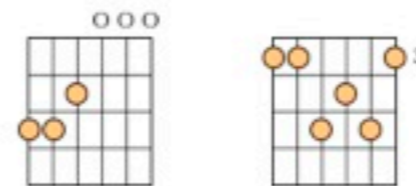
Keyboard display:

[Play chord](#)



Guitar chord voicings:

[More](#)

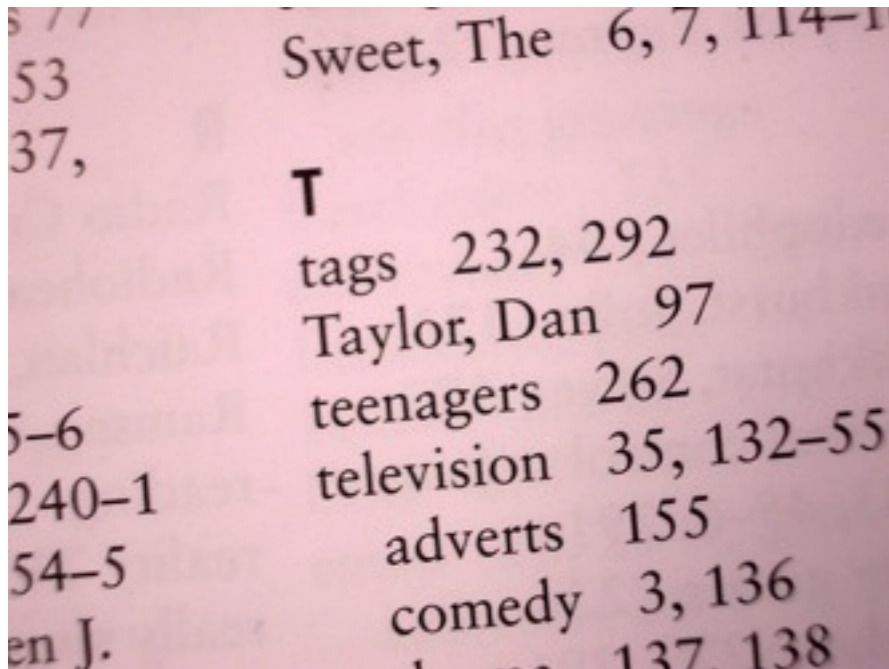


Computed by: [Wolfram Mathematica](#)

[Download as: PDF](#) | [Live Mathematica](#)

Pesquisar

por índice invertido



[Web](#) | [Imagens](#) | [Notícias](#) | [Blogs](#) | [PAi](#) | [PBi](#) | [Directório](#) | [Produtos](#) | [mais...](#)

[Portal SAPO](#) | [Registe o seu site](#) | [Anuncie na Pesquisa](#) | [Blog](#) | [Ajuda?](#)

Directório

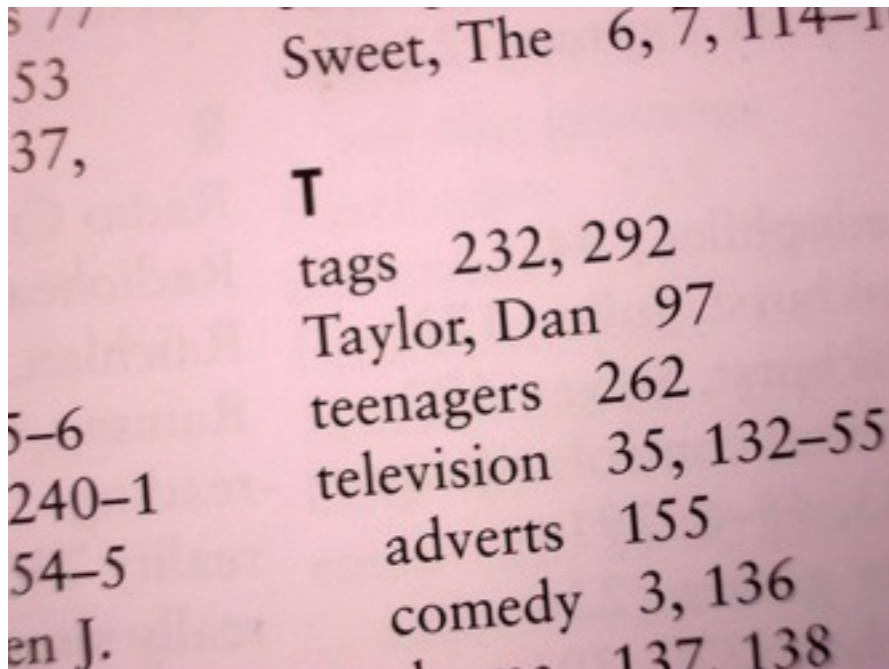
- [Arte e cultura](#)
Cinema, Teatro, Fotografia
- [Economia e Negócios](#)
Imobiliário, Emprego, Comércio
- [Nacional e Regional](#)
Lisboa, Porto, Faro
- [Ciência](#)
Astronomia, Gestão, Medicina
- [Educação](#)
Superior, Profissional, E-learning
- [Saúde](#)
Hospitais, Farmácias, Crianças
- [Comunicação Social](#)
Jornais, Televisões, Revistas
- [Entretenimento e Lazer](#)
Cinema, Desporto, Humor
- [Sociedade](#)
Política, Religião, Minorias
- [Desporto](#)
Futebol, Motorizados, Aquáticos
- [Estado e Administração](#)
Governo, ONGs, UE
- [Tecnologia e Internet](#)
Mail, Software, Biotecnologia

PT.COM © 2007



Pesquisar

por índice invertido



[Web](#) | [Imagens](#) | [Notícias](#) | [Blogs](#) | [PAi](#) | [PBi](#) | [Directório](#) | [Produtos](#) | [mais...](#)

[Portal SAPO](#) | [Registe o seu site](#) | [Anuncie na Pesquisa](#) | [Blog](#) | [Ajuda?](#)

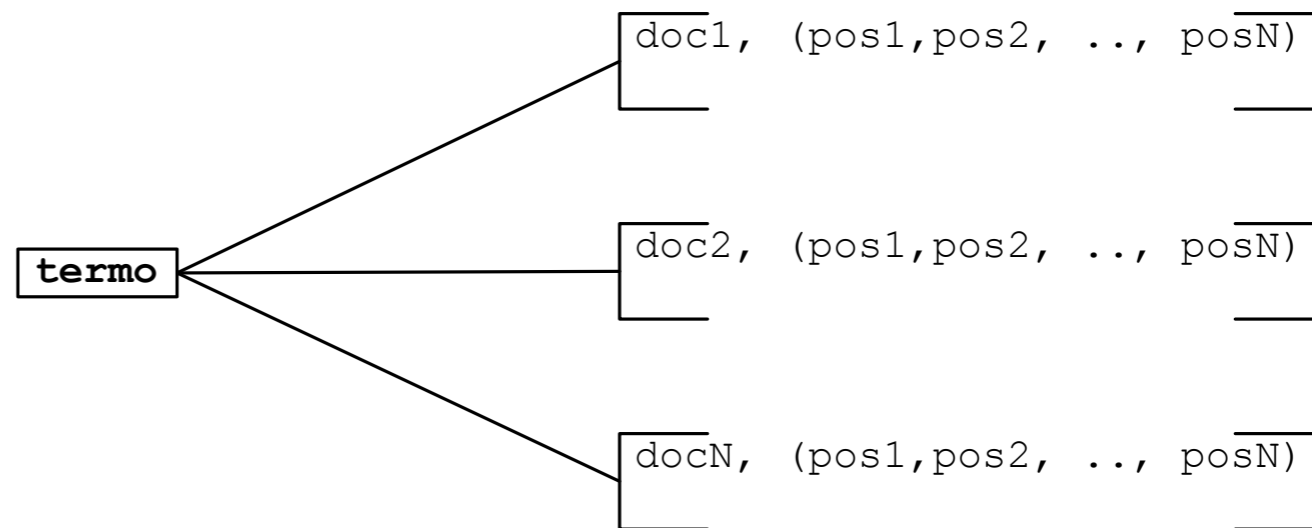
Directório ✕

- [Arte e cultura](#)
Cinema, Teatro, Fotografia
- [Economia e Negócios](#)
Imobiliário, Emprego, Comércio
- [Nacional e Regional](#)
Lisboa, Porto, Faro
- [Ciência](#)
Astronomia, Gestão, Medicina
- [Educação](#)
Superior, Profissional, E-learning
- [Saúde](#)
Hospitais, Farmácias, Crianças
- [Comunicação Social](#)
Jornais, Televisões, Revistas
- [Entretenimento e Lazer](#)
Cinema, Desporto, Humor
- [Sociedade](#)
Política, Religião, Minorias
- [Desporto](#)
Futebol, Motorizados, Aquáticos
- [Estado e Administração](#)
Governo, ONGs, UE
- [Tecnologia e Internet](#)
Mail, Software, Biotecnologia

PT.COM © 2007

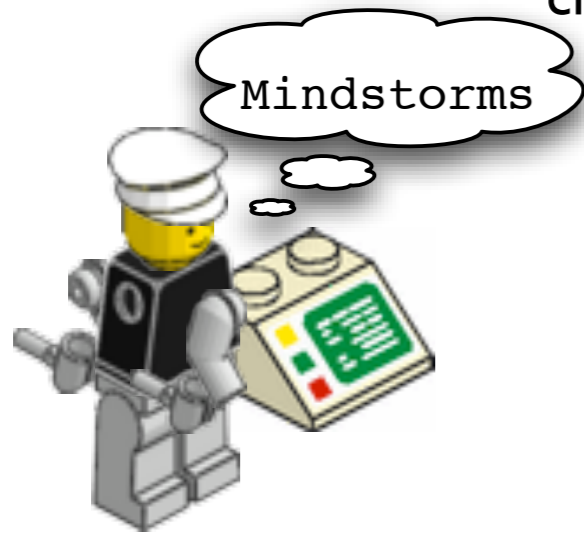


53
 37,
 Sweet, The 6, 7, 114-115
T
 tags 232, 292
 Taylor, Dan 97
 teenagers 262
 television 35, 132-55
 adverts 155
 comedy 3, 136
 137, 138



Encontrar

ciclo Search/Find/Obtain

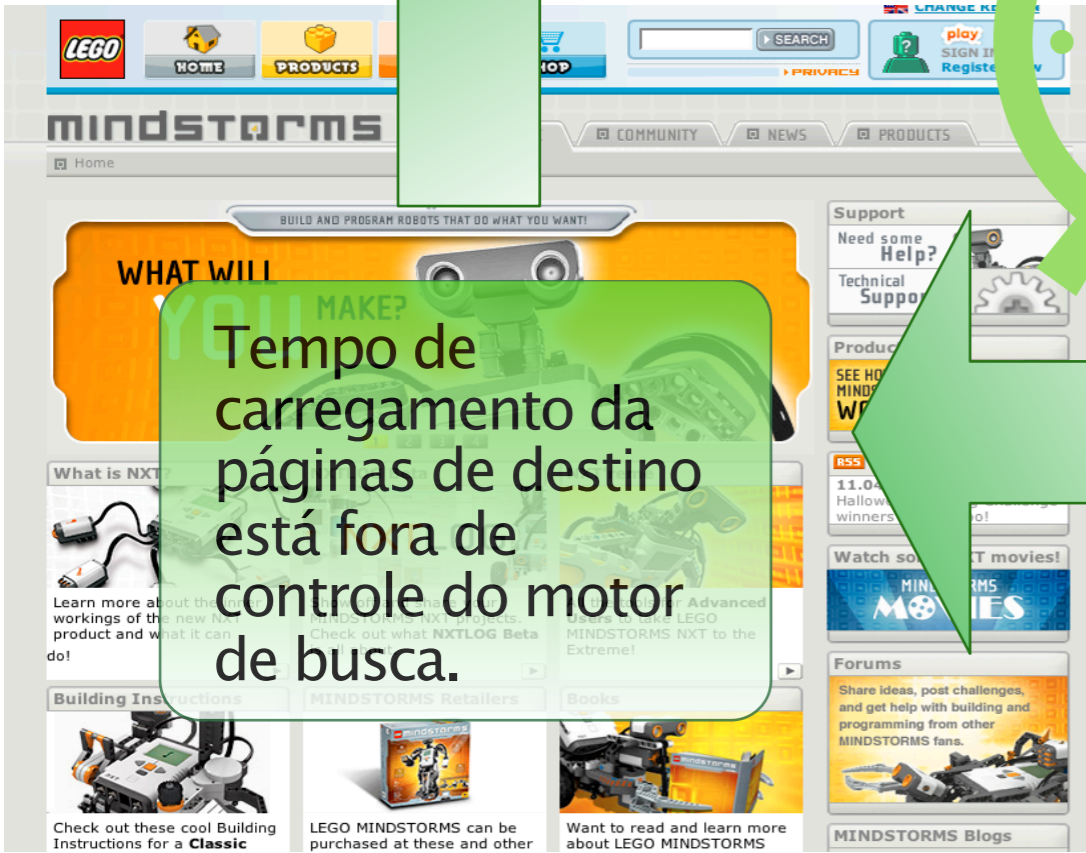
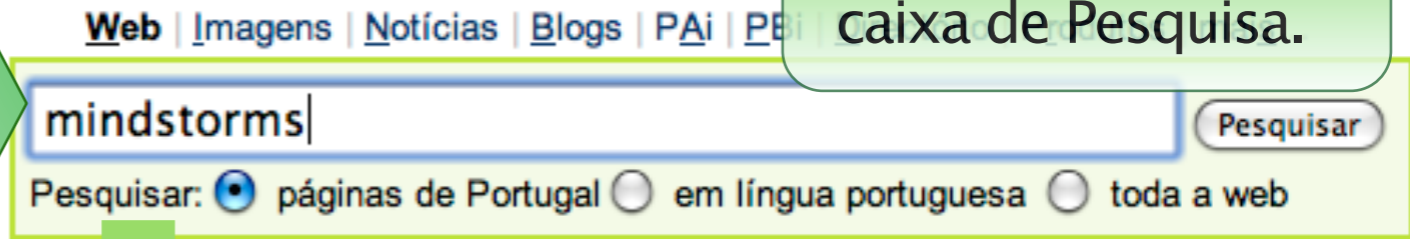


Encontrar

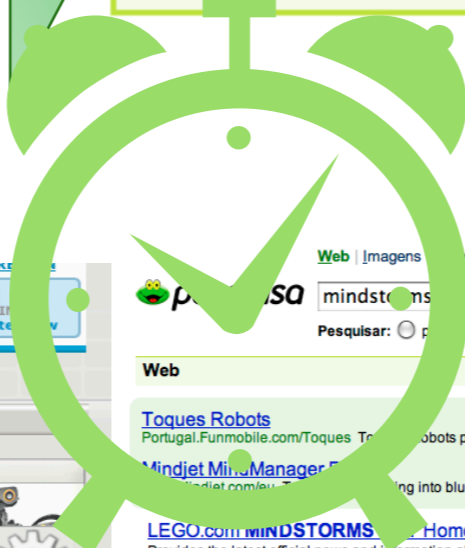
ciclo Search/Find/Obtain



Acesso rápido à caixa de Pesquisa.



Tempo de carregamento da páginas de destino está fora de controle do motor de busca.



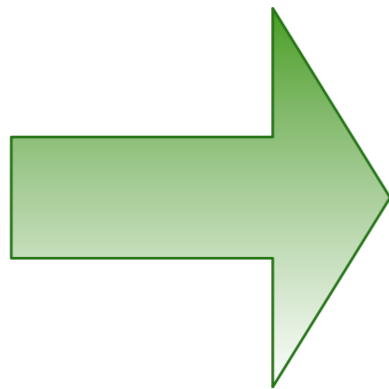
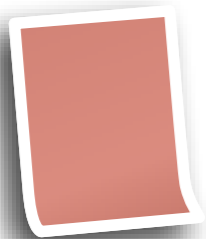
Resultado pretendido rápido de encontrar na página.



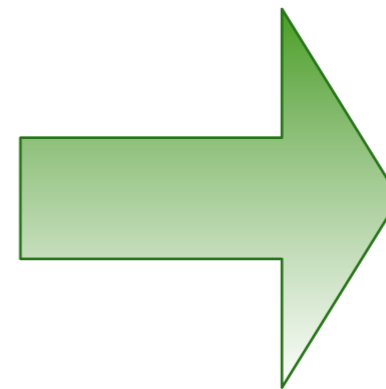
Volume da Informação – Desafios

Volume da Informação – Desafios

~15 KB
por documento



100
milhões de páginas



1,4 Tb de Índice de
Pesquisa



No caso Português, garantir 50 Pesquisas por segundo, cada uma em menos de 500ms por Pesquisa, com um Índice de 1,4 Terabyte.

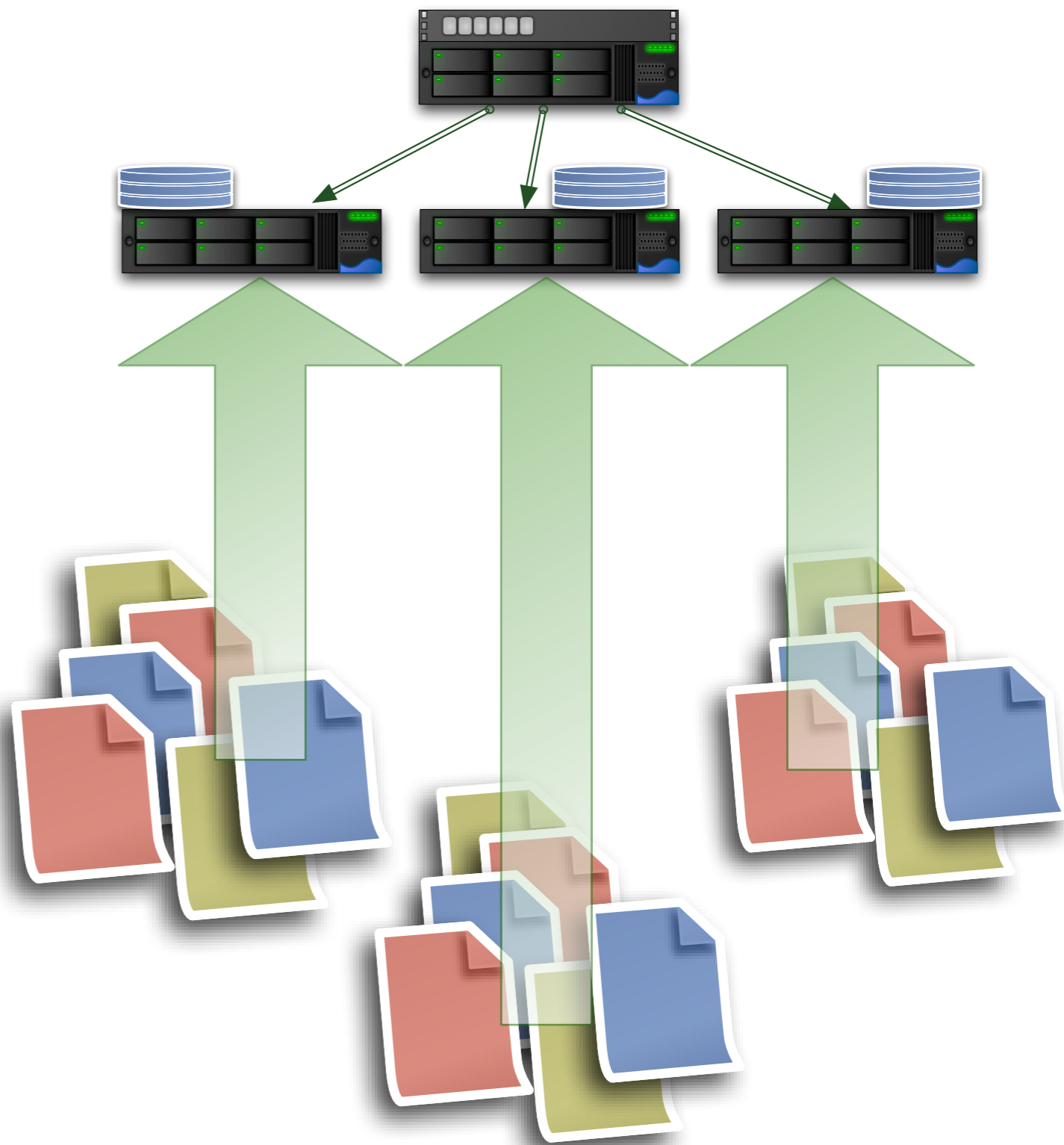


Crawlers



Crawlers

- Percorrem as páginas que irão ser indexadas.
- Limite de pedidos por site.
- Vários sites em simultâneo, optimizado para a largura de banda disponível.
- Crawling distribuído.
- Master distribui tarefas, cada nó é responsável pelo seu armazenamento.





Document Pipeline



Document Pipeline

- Conteúdo dos crawlers é enviado para a Document Pipeline.
- Envio em batch dos documentos.
- Cada url é um objecto Document.
- No início do processo, apenas sabemos a url e o conteúdo do documento não processado.

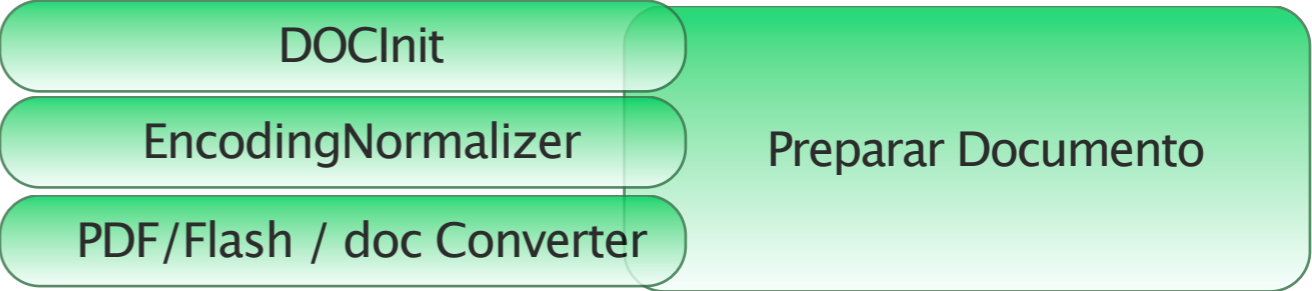




Document Pipeline

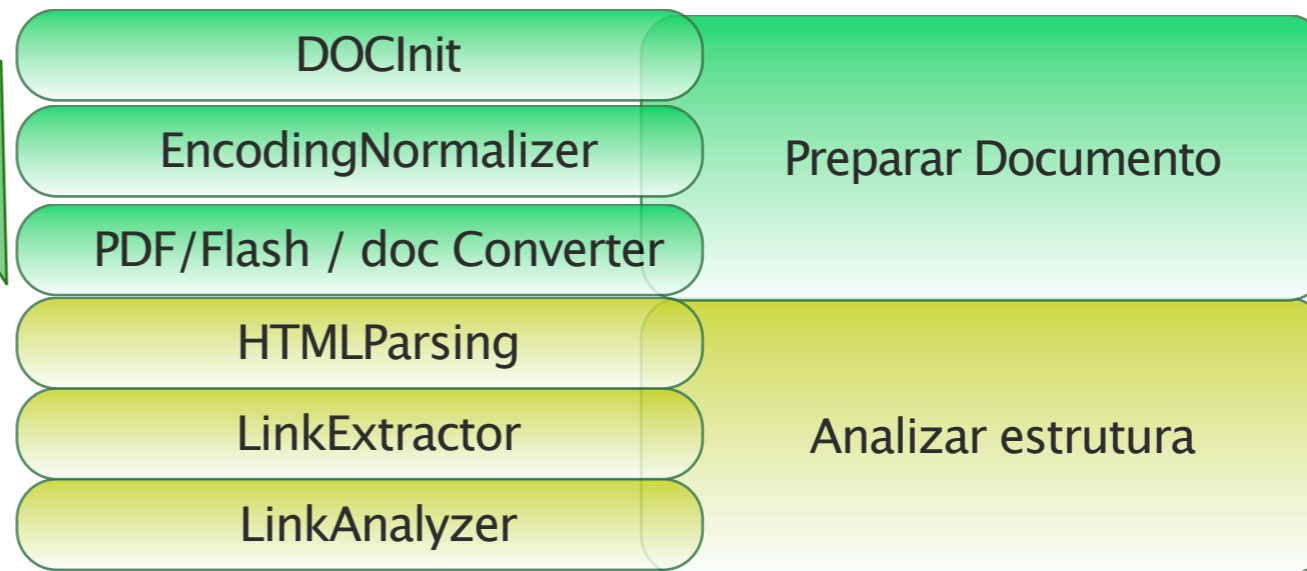


Document Pipeline



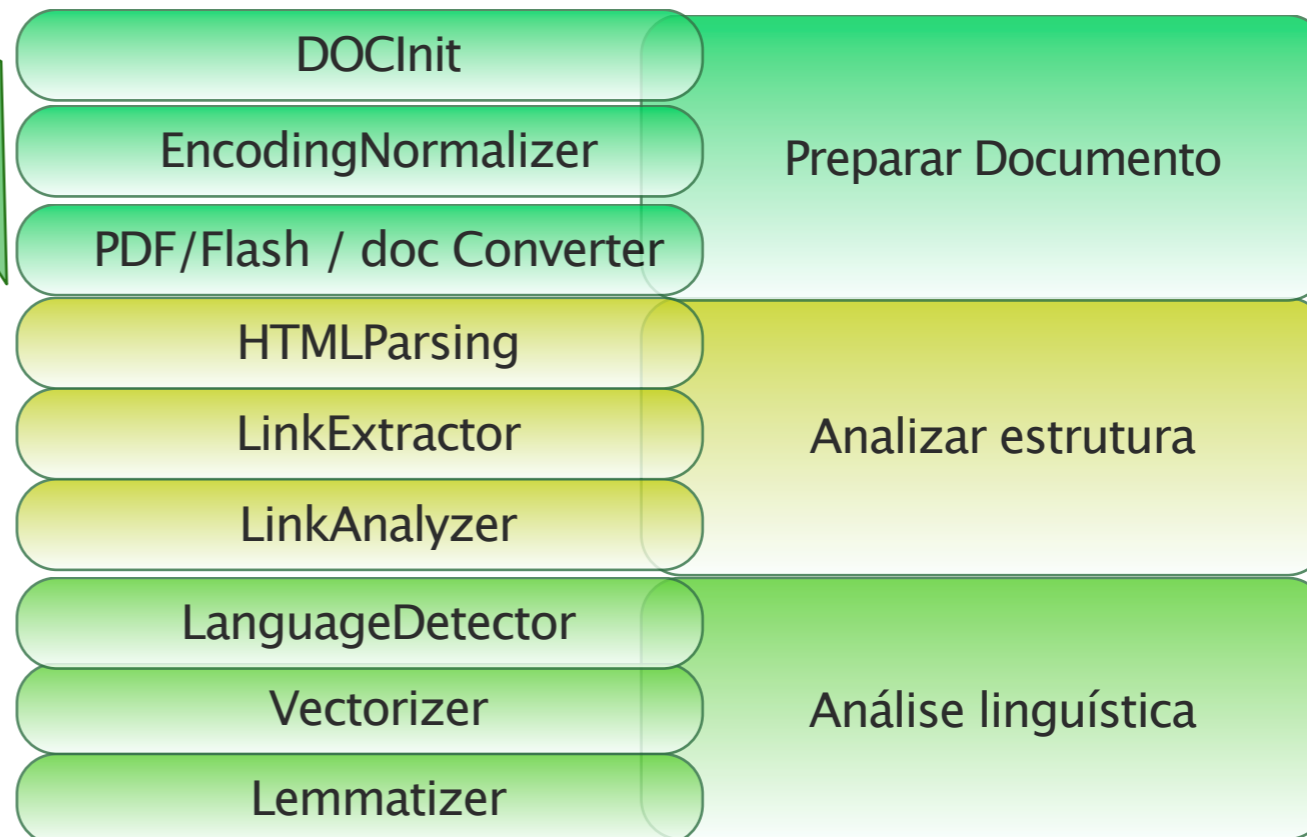


Document Pipeline



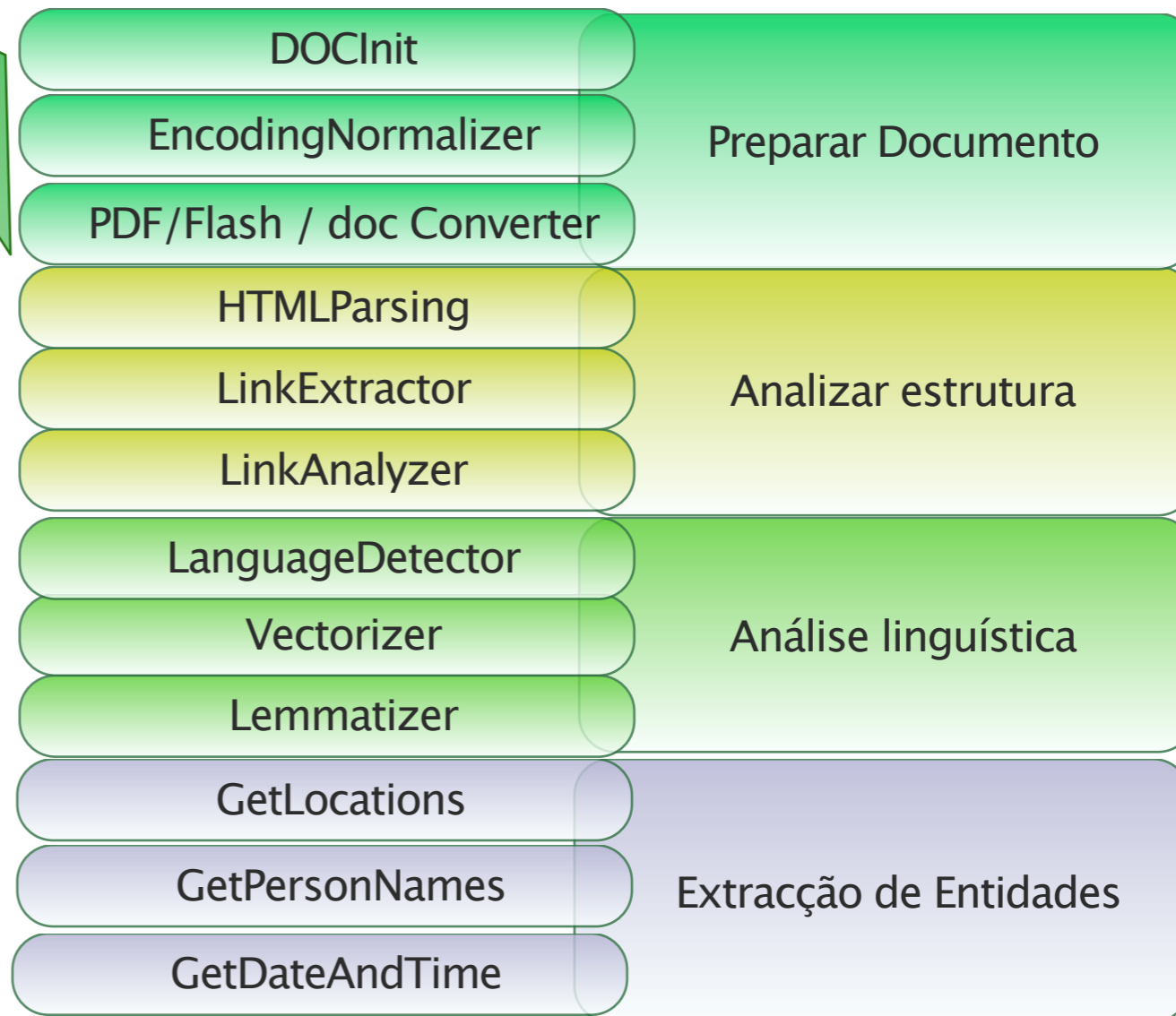


Document Pipeline





Document Pipeline



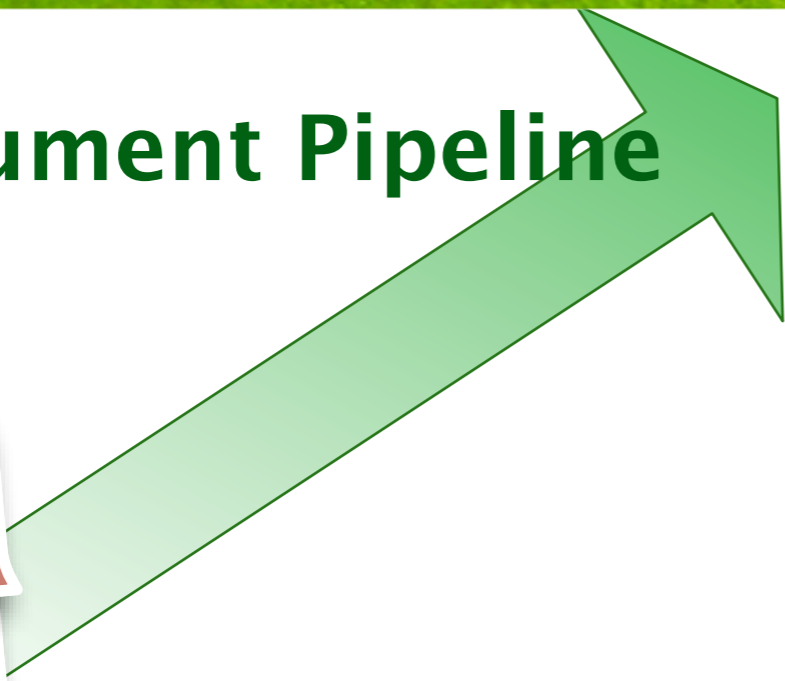


Document Pipeline

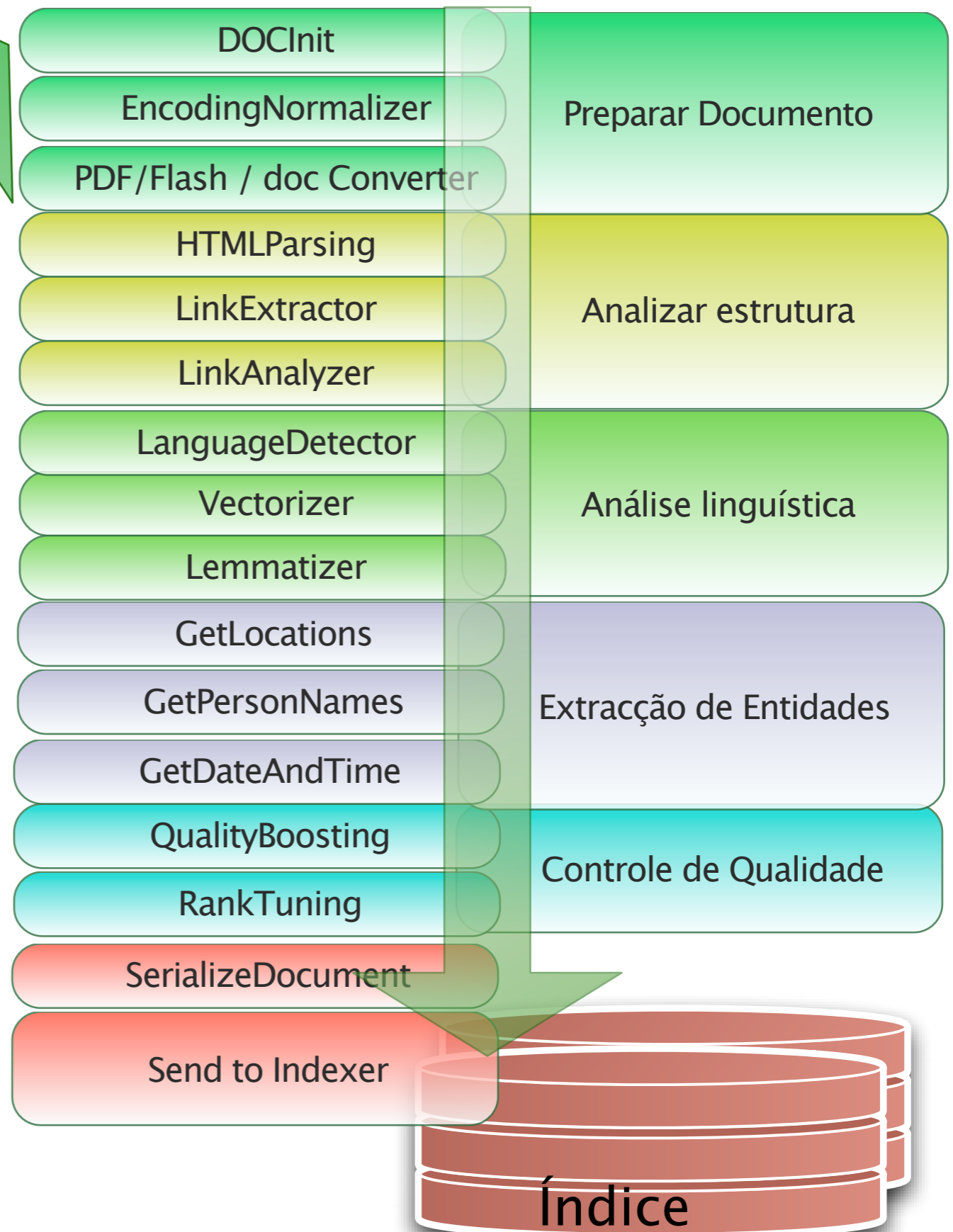




Document Pipeline



- Cada componente na cadeia é um módulo isolado que altera o objecto Document.
- Módulos desenvolvidos em Python.





Mais tecnologia em:

NOTÍCIAS ANÁLISES OPINIÃO MULTIMÉDIA EXTRAS

Pesquisar no TeK

OK

Últimas | Computadores | Negócios | Internet | Telecomunicações

Subscriver RSS Subscriver Newsletter

TeK > Notícias > Negócios

Nokia, TMN e Sapo continuam a liderar preferências

Publicado por Casa dos Bits há 20 horas e 27 minutos | 0 comentários



Nokia, TMN, Sapo, HP e Canon lideram, cada uma na sua categoria, as preferências dos portugueses, segundo o estudo "Marcas de Confiança" de 2010, das [Seleções do Reader's Digest](#), relativamente à área das tecnologias.

Na décima edição da análise, a **Nokia** mantém-se em primeiro (68%), em clara vantagem face às marcas que ocupam a segunda e terceira posições, a Samsung (13%) e a Sony Ericsson (9%). Relativamente a 2009, as marcas de confiança são as mesmas, mas as percentagens eram outras - 73, 8,5 e seis por cento, respectivamente.

A **TMN** continua a ser o operador de rede móvel mais popular, com um primeiro lugar reforçado, já que 49 por cento dos inquiridos votaram na marca, face aos 45 por cento que o tinham feito em 2009. Segue-se a Vodafone com 36 por cento e a Optimus, com 13 por cento.

Na categoria de Empresas de Serviço de Internet, o **Sapo** conseguiu superar os resultados do ano anterior (38%), somando 45 por cento dos votos. A Zon Multimédia perdeu três por cento relativamente a 2009, apresentando agora 14 por cento, enquanto o MEO ganhou o terceiro lugar, com oito por cento, destronando a Cabovisão.

Na área dos computadores pessoais, a **Hewlett-Packard** viu a confiança na marca reforçada, consolidando o primeiro lugar com 35 por cento, mais cinco por cento do que em 2009. A Toshiba, com 17 por cento, e a Asus, com 12 por cento, mantêm, respectivamente, a segunda e a terceira posições nesta categoria.

A **Canon** (34%), a Sony (24%) e a Olympus (12%) são, por esta ordem, as três marcas de confiança dos portugueses na categoria Máquinas Fotográficas.

A análise "Marcas de Confiança" das Seleções do Reader's Digest foi promovida em 16 países, entre os quais Alemanha, Áustria, Bélgica, Espanha, Finlândia, França, Holanda, Hungria e Polónia. Em Portugal, este estudo foi realizado através de questionário postal endereçado a 12.200 assinantes da revista Seleções do Reader's Digest.

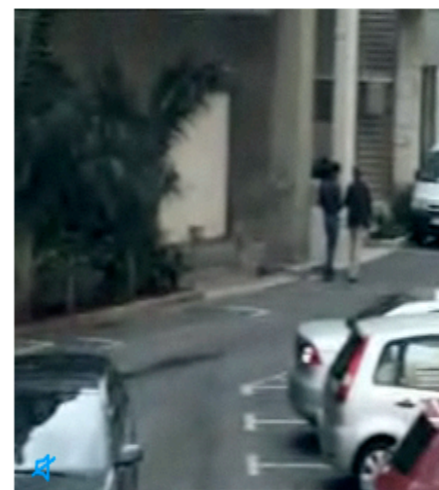
[Comentar este artigo »](#)

Quer ganhar um Nokia N97?

Simples e fácil... Clica aqui para saber como!
www.funphone.com.pt

Nokia na Vodafone

Grandes oportunidades online. Veja todos os modelos. Compre já!
Vodafone.pt/nokia

Anúncios **sapo** **mobile** [saiba mais >](#)
Notícias de tecnologia no telemóvelSugestões [Montra](#)

Sugestão TeK: Vai uma boleia para chegar ao destino?



Poupar dinheiro, ter companhia e uma atitude ambientalmente mais responsável são as consequências mais directas de partilhar o carro com terceiros. Mas a alternativa também pode ser uma saída em dias de greve.

[Ler artigo](#)[Outras sugestões »](#)



Mais tecnologia em: PPLWARE

NOTÍCIAS ANÁLISES OPINIÃO MULTIMÉDIA EXTRAS

Pesquisar no TeK OK

Últimas | Computadores | Negócios | Internet | Telecomunicações

Subscrever RSS Subscrever Newsletter

Title

TeK > Notícias > Negócios

Nokia, TMN e Sapo continuam a liderar preferências

Author, Date

Publicado por Casa dos Bits há 20 horas e 27 minutos | 0 comentários



Keywords

Nokia, TMN, Sapo, HP e Canon lideram, cada uma na sua categoria, as preferências dos portugueses, segundo o estudo "Marcas de Confiança" de 2010, das Selecções do Reader's Digest, relativamente à área das tecnologias.

Na décima edição da análise, a Nokia mantém-se em primeiro (68%), em clara vantagem face às marcas que ocupam a segunda e terceira posições, a Samsung (13%) e a Sony Ericsson (9%). Relativamente a 2009, as marcas de confiança são as mesmas, mas as percentagens eram outras - 73, 8,5 e seis por cento, respectivamente.

A TMN continua a ser o operador de rede móvel mais popular, com um primeiro lugar reforçado, já que 49 por cento dos inquiridos votaram na marca, face aos 45 por cento que o tinham feito em 2009. Segue-se a Vodafone com 36 por cento e a Optimus, com 13 por cento.

Na categoria de Empresas de Serviço de Internet, o Sapo conseguiu superar os resultados do ano anterior (38%), somando 45 por cento dos votos. A Zon Multimédia perdeu três por cento relativamente a 2009, apresentando agora 14 por cento, enquanto o MEO ganhou o terceiro lugar, com oito por cento, destronando a Cabovisão.

Na área dos computadores pessoais, a Hewlett-Packard viu a confiança na marca reforçada, consolidando o primeiro lugar com 35 por cento, mais cinco por cento do que em 2009. A Toshiba, com 17 por cento, e a Asus, com 12 por cento, mantêm, respectivamente, a segunda e a terceira posições nesta categoria.

A Canon (34%), a Sony (24%) e a Olympus (12%) são, por esta ordem, as três marcas de confiança dos portugueses na categoria Máquinas Fotográficas.

A análise "Marcas de Confiança" das Selecções do Reader's Digest foi promovida em 16 países, entre os quais Alemanha, Áustria, Bélgica, Espanha, Finlândia, França, Holanda, Hungria e Polónia. Em Portugal, este estudo foi realizado através de questionário postal endereçado a 12.200 assinantes da revista Selecções do Reader's Digest.

Comentar este artigo »

Quer ganhar um Nokia N97?

Simples e fácil... Clica aqui para saber como!
www.funphone.com.pt

Nokia na Vodafone

Grandes oportunidades online. Veja todos os modelos. Compre já!
Vodafone.pt/nokia

Anúncios sapo

Últimas Em destaque Mais comentados

- > Vídeos do MyWay na Reuters
- > Portugal em 6º em vírus nos PCs
- > Android chega às 50 mil aplicações
- > Anti-vírus comparados
- > 5 perguntas a... Paulo Costa
- > TomTom anuncia novidades



mobile [saiba mais >](#)
Notícias de tecnologia no telemóvel

Sugestões Montra

Sugestão TeK: Vai uma boleia para chegar ao destino?



Poupar dinheiro, ter companhia e uma atitude ambientalmente mais responsável são as consequências mais directas de partilhar o carro com terceiros. Mas a alternativa também pode ser uma saída em dias de greve.

Ler artigo

[Outras sugestões »](#)



Indexação

```
<field name="title" fullsort="yes" tokenize="auto" lemmatize="yes" >
  <vectorize default="10:0" />
</field>
<field name="body" tokenize="auto" max-result-size="1024"
  fallback-ref="teaser" result="dynamic" index="no" lemmatize="yes">
  <vectorize default="5:5" alternative="{ja,ko,zh,szh,tzh}:5:0" />
</field>
<field name="teaser" index="no" />
<field name="description" element-name="meta_description" result="no" tokenize="auto" />
<field name="anchortext" result="no" tokenize="auto" />
<field name="keywords" element-name="meta_keywords_off" result="no" tokenize="auto" />
```

- Cada índice é gerado localmente.
- Perfil do Índice permite definir quais os campos mais relevantes e qual o peso.



Indexação

```
<context weight="50">  
  <field-weight field-ref="body" value="5" />  
  <field-weight field-ref="description" value="30" />  
  <field-weight field-ref="urlkeywords" value="40" />  
  <field-weight field-ref="keywords" value="50" />  
  <field-weight field-ref="title" value="60" />  
</context>
```

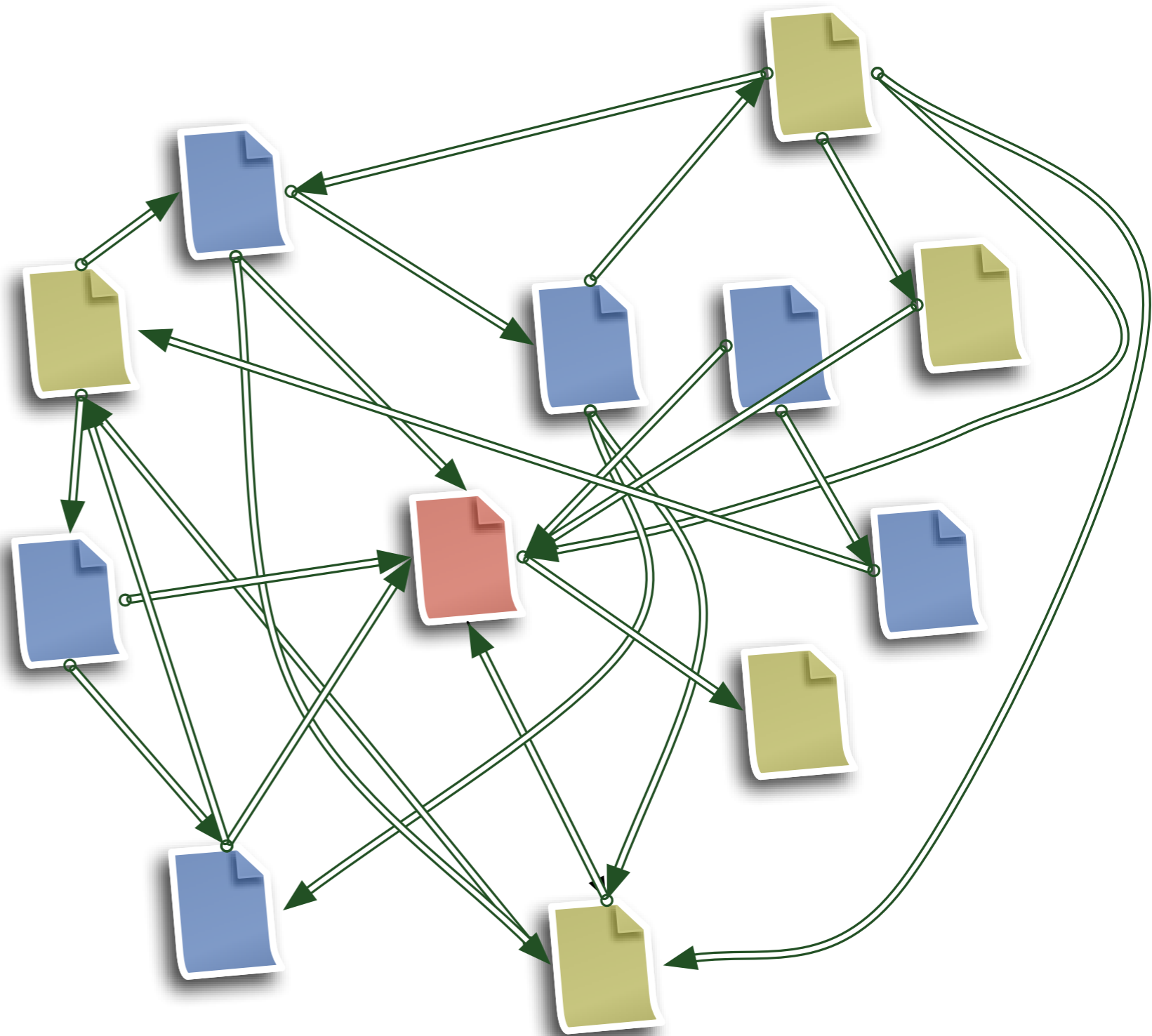
- Cada índice é gerado localmente.
- Perfil do Índice permite definir quais os campos mais relevantes e qual o peso.

Confiar

PageRank

$$\sum_{i=1}^N \ell(p_i, p_j) = 1$$

- Análise estrutural da Web.
- Autoridade nem sempre qualidade
- Comportamento dos utilizadores ajuda no ranking.
- Problemas com Spam.

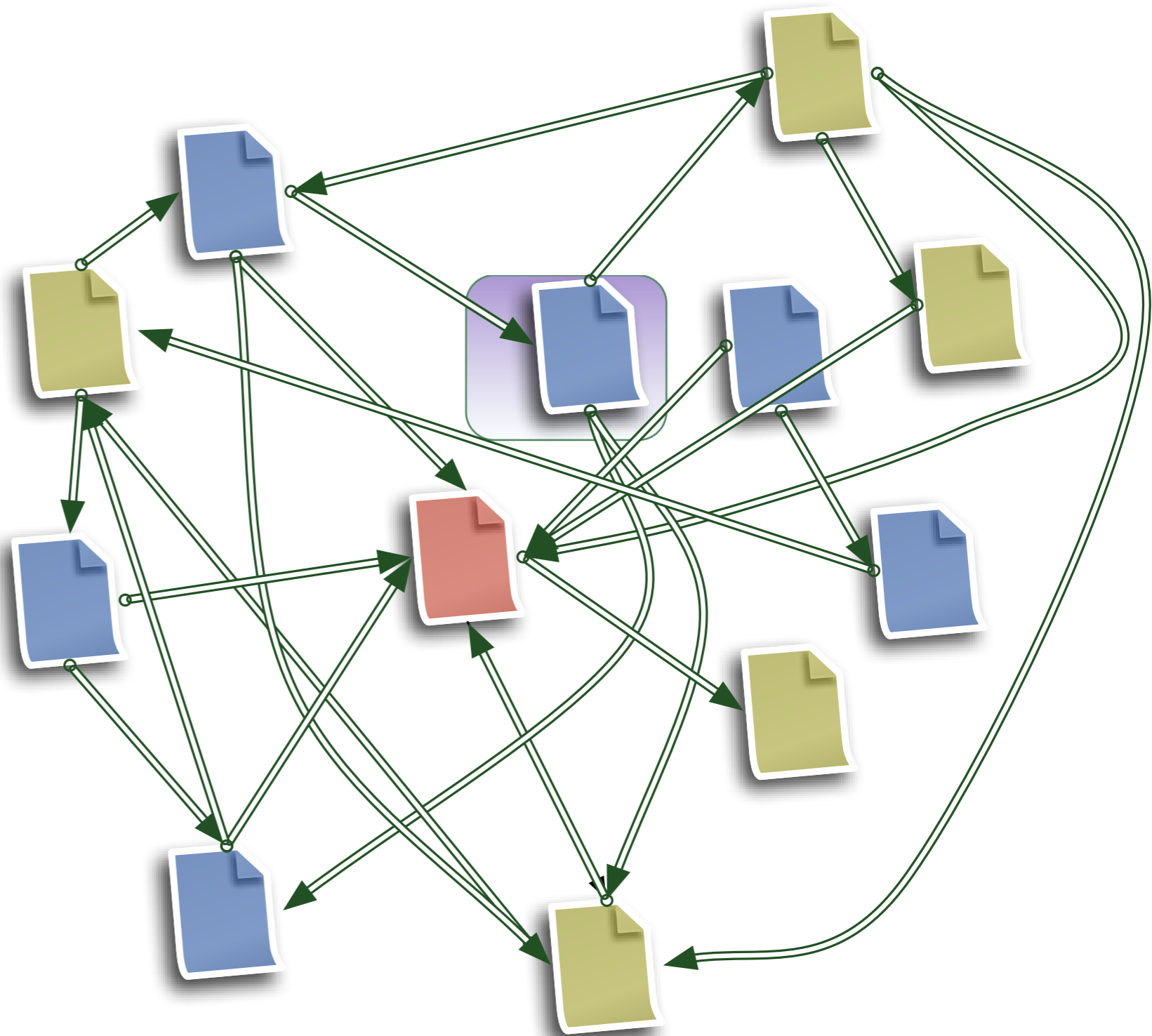


Confiar

PageRank

$$\sum_{i=1}^N \ell(p_i, p_j) = 1$$

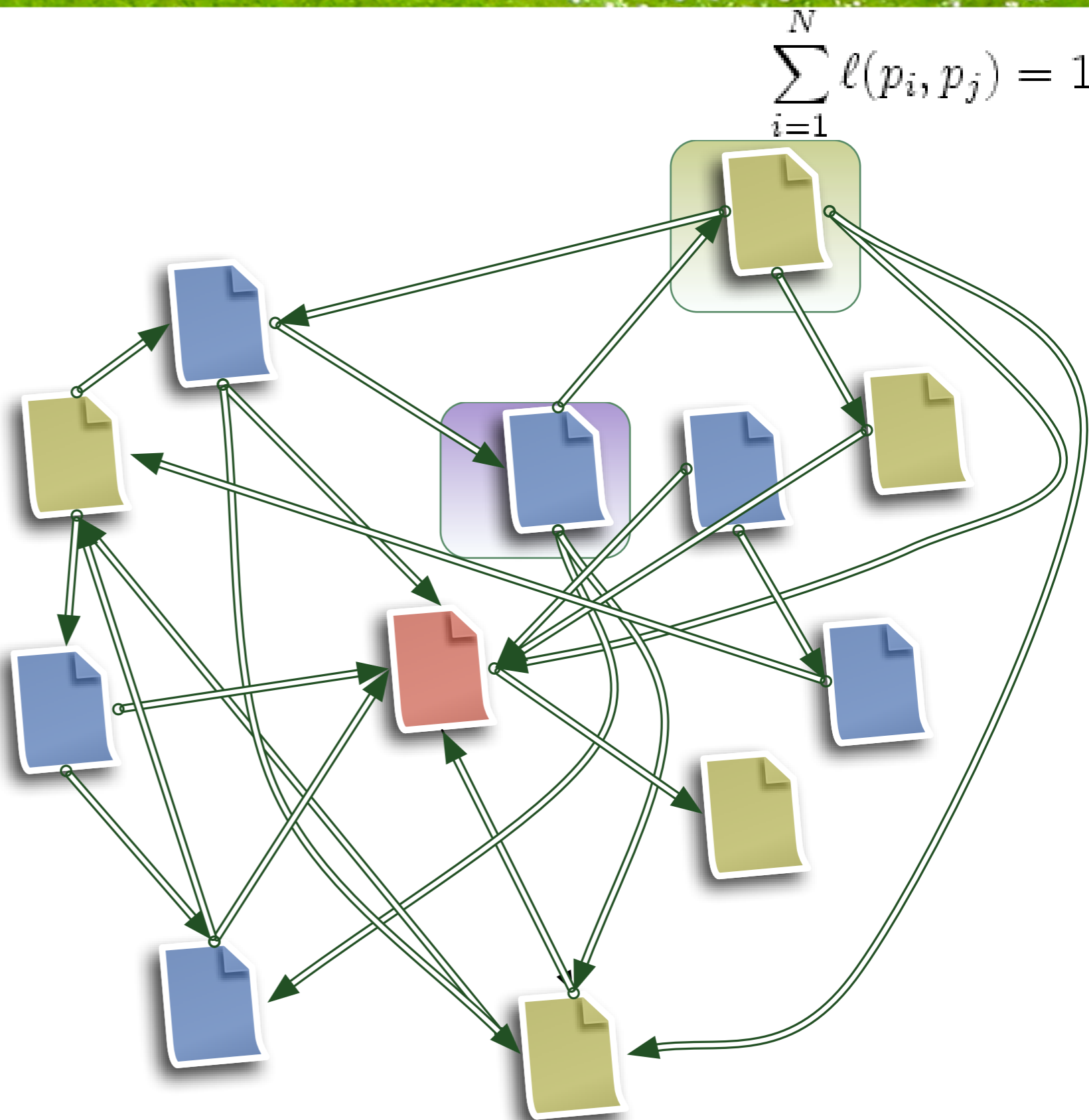
- Análise estrutural da Web.
- Autoridade nem sempre qualidade
- Comportamento dos utilizadores ajuda no ranking.
- Problemas com Spam.



Confiar

PageRank

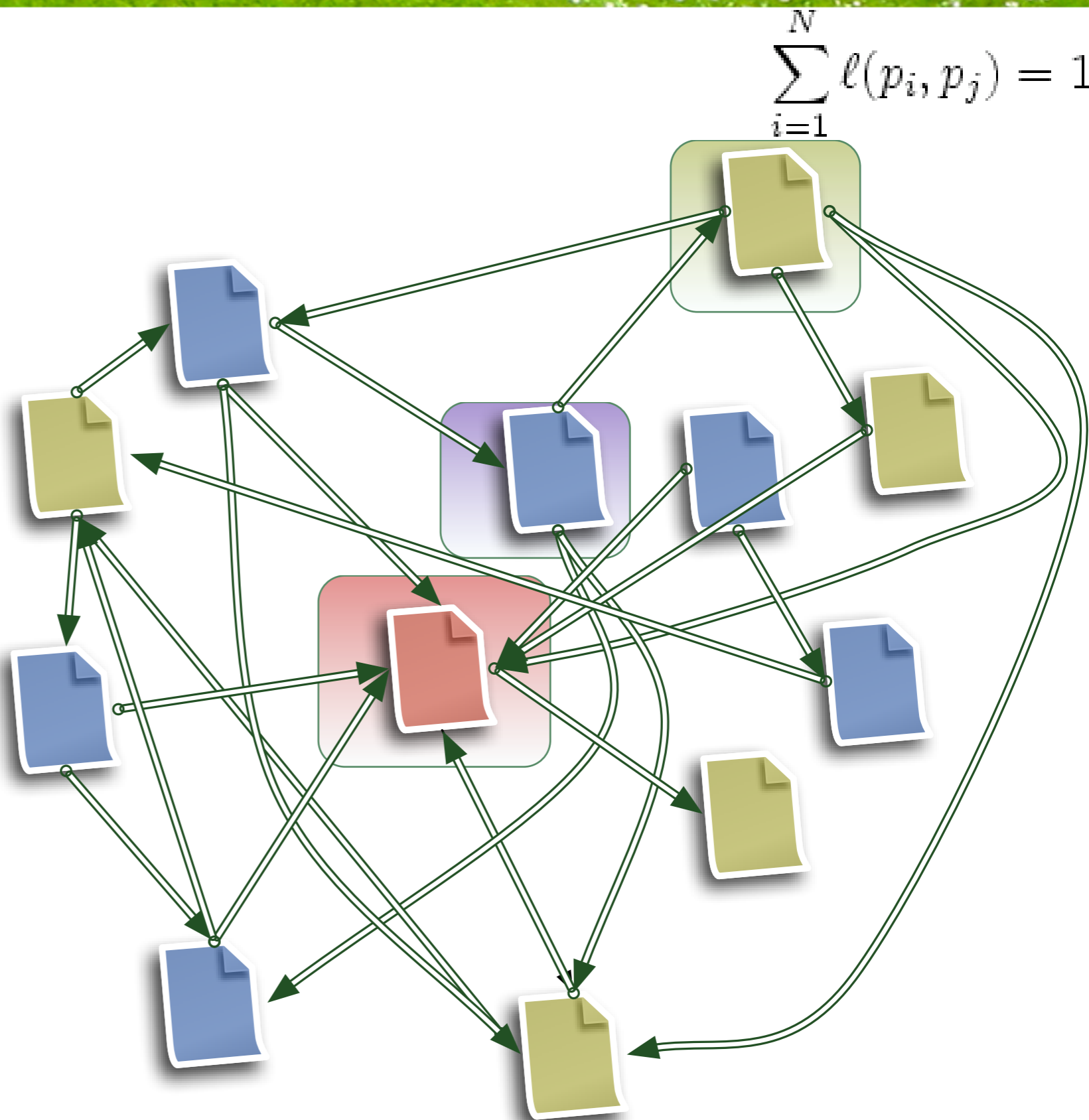
- Análise estrutural da Web.
- Autoridade nem sempre qualidade
- Comportamento dos utilizadores ajuda no ranking.
- Problemas com Spam.



Confiar

PageRank

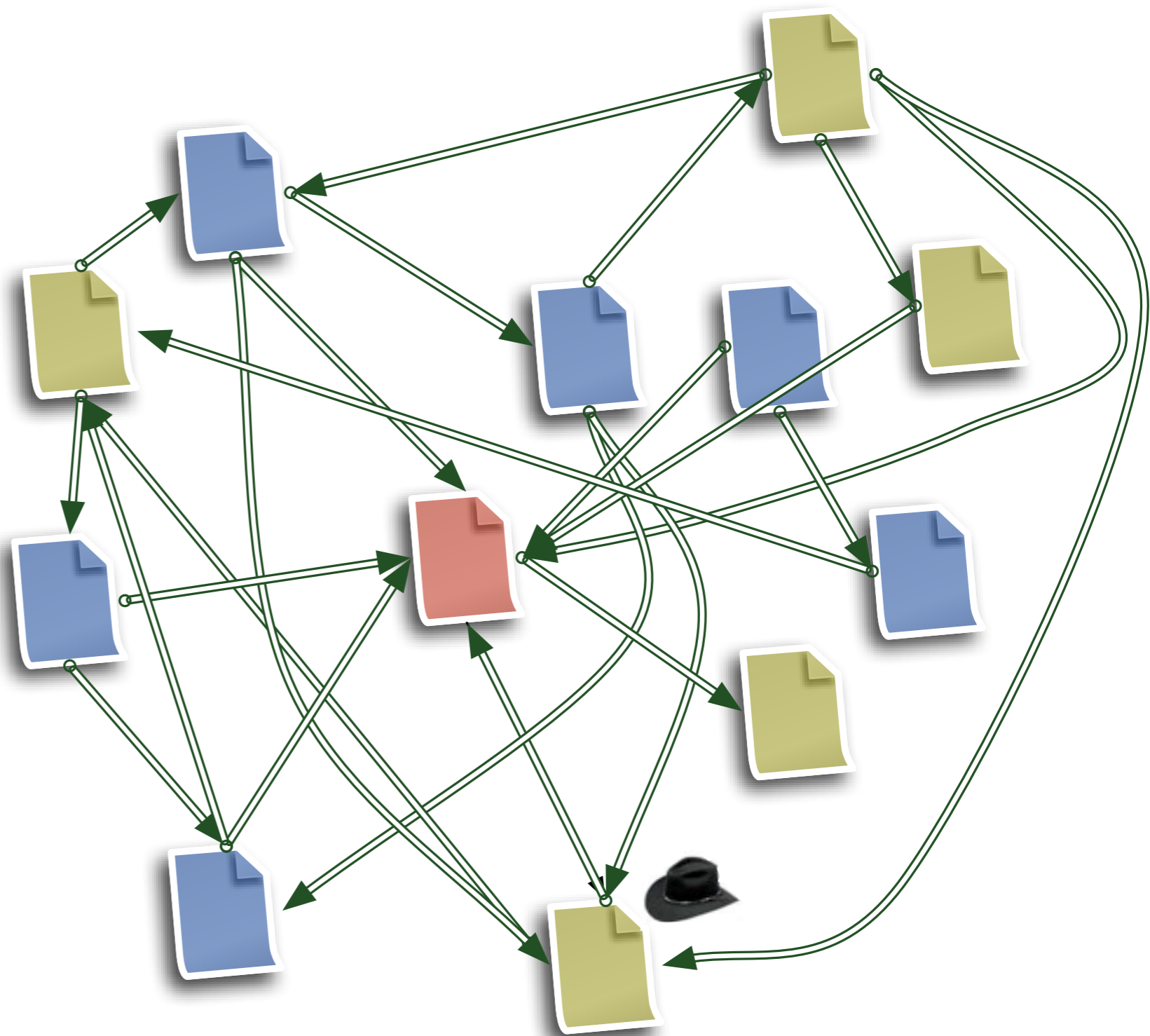
- Análise estrutural da Web.
- Autoridade nem sempre qualidade
- Comportamento dos utilizadores ajuda no ranking.
- Problemas com Spam.



Confiar

SEO

- Como manipular a confiança a favor do meu site?
- White hat
 - Trazer resultados mais relevantes, seguindo as regras.
- Black hat
 - Trazer mais tráfego, com recursos a meios menos legítimos.



Obrigado!



João Pedro Gonçalves
jp@sapo.pt
Abril 2010

