

# BIB<sub>TEX</sub> Bibliography Index Maker: Informe del Projecte

Ramon Xuriguera  
rxuriguera@gmail.com

Març 2010

## 1 Objectius del projecte

*BIB<sub>TEX</sub> Index Maker* és una eina d'ajuda a la creació d'índexos bibliogràfics pensada com un complement a aplicacions de maneig de referències ja existents com poden ser *JabRef*<sup>1</sup> o *Mendeley*<sup>2</sup>.

La principal funcionalitat consisteix en escanejar un directori que conté articles científics en PDF i generar un índex bibliogràfic en BIB<sub>TEX</sub> amb les referències d'aquests fitxers. Aquest índex es pot importar des de les aplicacions esmentades o bé pot ser referenciat directament des d'un nou document T<sub>EX</sub>.

## 2 Treball existent

Actualment existeixen nombroses aplicacions dedicades al maneig de referències. Algunes d'elles utilitzen les metadades dels fitxers per tal de trobar informació com ara el títol o l'autor, però cap de les eines que hem trobat llegeix el contingut dels documents.

Empreses com ara *Google* o *Microsoft* extreuen informació de documents PDF per oferir serveis com ara *Scholar* o *Academic Search*, però no ofereixen el codi font per tal que altres el puguin utilitzar i per tant, no sabem com funcionen. *CiteSeer* és un projecte *open source* de característiques similars, però que també té limitacions. El sistema funciona extraient la bibliografia dels articles, però té problemes per obtenir els camps de la capçalera del propi fitxer.

## 3 Descripció detallada

El format PDF no permet extreure el contingut de forma fàcil i, per tant, les eines d'extracció de text disponibles tenen bastants problemes. En especial hi ha problemes amb el flux del text dins el fitxer, els caràcters Unicode, les lligadures de caràcters (e.g. *f* es representa amb un sol caràcter), superíndexos, etc. A més, les capçaleres dels documents varien molt depenent de la revista a la qual han estat publicats.

---

<sup>1</sup><http://jabref.sourceforge.net/>

<sup>2</sup><http://www.mendeley.com/>

Per evitar aquests problemes, fem ús de la informació de qualsevol de les biblioteques digitals accessibles des d'Internet, per exemple *SpringerLink* o *ACM*. Per poder extreure la referència corresponent a un article seguirem el següent procediment que mostra la figura 1. Bàsicament, per cada document PDF, l'aplicació:

1. Extreu el contingut en forma de text.
2. Genera un llistat de consultes a partir del text. Per exemple, agafa frases de l'*abstract*.
3. Obté els resultats d'aquestes consultes en algun cercador (Google, Bing, etc.).
4. Mira d'extreure la referència de les pàgines retornades pel cercador.
5. Comprova que les dades de la referència corresponen al fitxer.

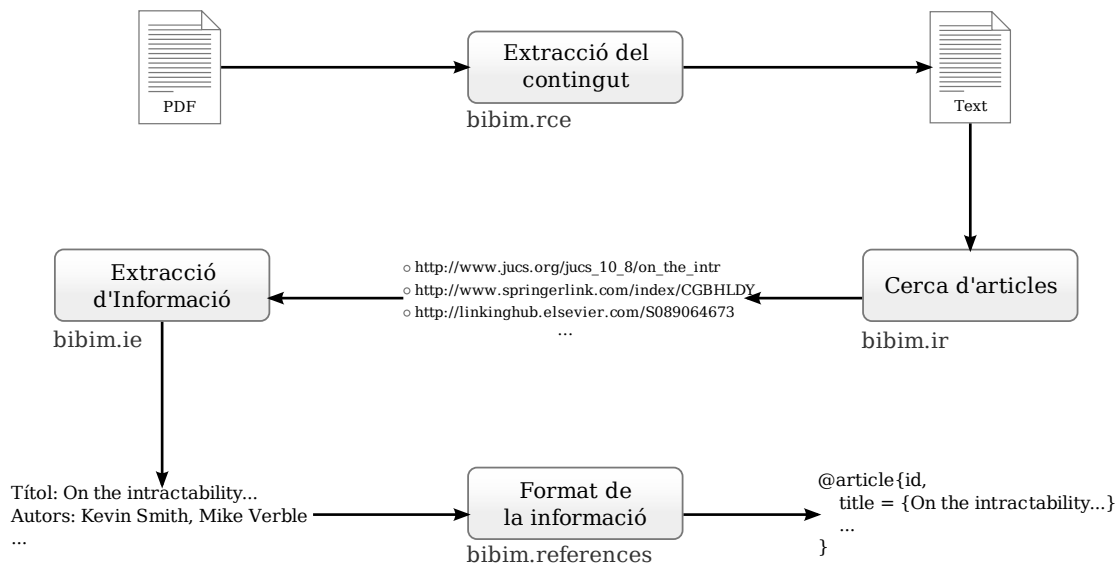


Figura 1: Procediment per obtenir la referència a un article

El pas més interessant d'aquest procés és l'extracció de la referència a partir de les pàgines resultants de fer la cerca. Com és lògic, els resultats obtinguts d'Internet corresponen a pàgines completament diferents i el nostre programa ha de ser capaç de detectar on es troba cada fragment d'informació i extreure'l. A primera vista podria semblar que si les pàgines són tant diferents, el fet de buscar l'article a Internet no ens ha solucionat res. La principal diferència, però, és que la informació d'aquestes pàgines està estructurada.

Una primera idea per fer-ho és tenir un conjunt de *wrappers* de manera que cadascun d'ells sap com extreure la informació de cada pàgina. Aquesta solució dista de ser perfecta, però pot arribar a funcionar. Hem pogut comprovar que els articles d'un mateix camp de coneixement acostumen a aparèixer a les mateixes biblioteques. Això implica que, a la pràctica, amb un número no molt gran de *wrappers* es poden obtenir resultats prou bons. El principal problema d'aquesta tècnica és que quan alguna de les pàgines de les quals es depèn canvia la seva estructura, cal modificar també el *wrapper* corresponent.

Una alternativa és la possibilitat de generar aquests *wrappers* de forma automàtica a partir d'un número reduït d'exemples. Això no només facilita la possibilitat d'extreure informació de noves pàgines sinó que quan l'aplicació detecta que els resultats comencen a empitjorar es pot reaprendre les regles d'extracció. Per poder fer això necessitem tenir les dades de les extraccions anteriors emmagatzemades a una base de dades d'exemples on. Per cada exemple hem de tenir, com a mínim, la URL de la pàgina i les dades que volem extreure. El procediment és el següent:

1. Cerquem cada un dels camps que volem extreure dins de la pàgina.
2. Establim la posició de l'etiqueta HTML que conté la informació dins de la pàgina.
3. Trobem la posició del tros d'informació dins de l'etiqueta HTML
4. Generem la regla d'extracció
5. Validem la regla amb la resta d'exemples

## 4 Paquets de treball

La taula següent mostra el llistat de paquets de treball acompanyats del temps estimat que ha calgut o caldrà per finalitzar-los i l'estat en què es troben. Cal comentar que les tasques amb un estat  $\geq 95\%$  es consideren tasques completes, però susceptibles a correccions.

Tasca	Temps estimat	Estat
Recerca i proves de concepte	20 dies	100%
Infraestructura (repositori, <i>build tracker</i> , etc.)	4 dies	100%
Extracció del text dels PDF ( <i>bibim.rce</i> )	4 dies	95%
Obtenció pàgines amb la referència ( <i>bibim.ir</i> )	9 dies	
Obtenir consultes del text	1 dia	95%
<i>Browser</i> per obtenir consultes	2 dies	95%
Cerca de les consultes (múltiples cercadors)	5 dies	95%
Ordenar resultats de la cerca	1 dia	80%
Extracció d'informació ( <i>bibim.ie</i> )	17 dies	
<i>Wrappers</i> a mà	5 dies	60%
Generació automàtica de <i>wrappers</i>	12 dies	
Extracció d'exemples	1 dia	50%
Generació de regles	8 dies	10%
<i>Wrapper</i> genèric	3 dies	15%
Tractament de referències ( <i>bibim.references</i> )	5 dies	
<i>Parsing</i>	2 dies	95%
Generació	2 dies	95%
Crear índexs	1 dia	90%
Base de dades ( <i>bibim.db</i> )	3 dies	95%
Mòdul principal ( <i>bibim.main</i> )	8 dies	
Escaneig de fitxers	1 dia	95%
Codi principal de l'aplicació	3 dia	95%
<i>Threading</i>	2 dies	95%
Validació de referències	2 dies	85%
Interfície d'usuari	4 dies	15%
Memòria	-	$\approx 10\%$

Alguns aspectes que cal comentar:

- La implementació de *wrappers* a mà està parada. Aquests *wrappers* seran substituïts si s'obtenen bons resultats amb la generació automàtica.
- El *parsing* de referències és necessari per poder validar els camps de la referència en aquells casos en que se n'extreu el codi complet Bib<sub>T</sub><sub>E</sub>X.
- Els *threads* s'utilitzen per poder processar més dun fitxer a la vegada. Per exemple, mentre s'estan esperant els resultats d'una cerca a Internet, es pot anar extraient el text d'un altre PDF.