

Recent developments in the TransSearch bilingual concordancer

Guy Lapalme RALI - Université de Montréal

joint work with

Fabrizio Gotti, Philippe Langlais, Julien Bourdaillet, Stéphane Huet, Jacques Steinlin Industrial contributor:Terminotix





university NLP laboratory in Canada.

Research Projects	System demos
 <u>Translation</u> <u>Information extraction</u> <u>Automatic summarization</u> <u>Information retrieval</u> <u>Judicial texts processing</u> <u>Environmental information</u> 	 <u>Translation</u> <u>TransSearch</u>: bilingual concordancer <u>TransType</u>: animation of an interactive translation session <u>Reacc</u>: automatic French accentuation <u>SILC</u>: language and coding detection <u>Lexiqum</u>: Québec text concordancer
	Information
 <u>Publications</u> <u>Textual Resources</u> <u>Contacts</u> <u>Members</u> <u>Collaborators</u> 	 <u>RALI OLST-MITACS seminars (mostly in French)</u> Courses by professors of the RALI Fall 2011: <u>IFT6810: Traitement Statistique des Langues Naturelles</u> Fall 2011: <u>IFT3335: Intelligence Artificielle</u> <u>RALI in the press</u>
A selection of international conferences in NI	P

RALI - Today

- 3 professors
 - Philippe Langlais : machine translation
 - Jian-Yun Nie : information retrieval
 - Guy Lapalme : summarization, generation, info extraction
- Adjunct professor : Atefeh Farzindar
- Students
 - I post-doc
 - 8 Ph.D.
 - 5 M.Sc.
- I Research associate
 - Fabrizio Gotti



Goals of RALI

- Applied NLP using symbolic and statistical methods
 - Tools for translators
 - Writing aids
- Collaborative work
 - Linguistics department
 - Industrial partners
 - Governmental partners



http://www.tsrali.com

USER: la	apalme	QUERIES M	Y ACCOUNT	PREFERENCES	CONTACT	HELP	QUIT
Personaliz	zed Favorite / Bookmark : TransSearch (w	hat is this?)				Bilingua	al query
	Document collection Expression	on : House of Comm on: take+ ride	ons Hansard (19	986-2011)	\$		
1	Ils peuvent se faire passer des « p'tit	es vites ».	They can	be taken for a rid	e.		
2	Voyez ce qui leur arrive lorsqu'on pro	ofite d'elles.	Look what	at happens when the	ey are taken fo	r a ride.	
3	C'est trop embêtant et on ne soupço vouloir nous embobiner.	nne pas les compagnie	s de It becom us for a	es too cumbersome ride.	and we think no	obody is t a	iking
4	On le prend pour aller se promener e endroit ou dans un autre.	t on le laisse dans un	They tak	e it for a ride and	leave it somewh	ere else.	
5	Ils ont été abordés par des gens qu'i des amis, en qui ils avaient confiance	ls considéraient comme e, et ils se sont fait rou	e They had ler. friends ar	l been approached b nd trusted and they	oy people who the the second sec	hey consid a ride.	ered
6	Il serait donc important de connaître revenus, si on ne veut pas être les di	toute l'information su ndons de la farce.	les That is w about rev	why it is important to venue if we do not v	have full access vant to be take	s to inform n for a rid	ation l e .

Timeline of TransSearch

- 1992 Michel Simard (CITI) on Sun machines
- 1995 First Web version (free)
- 1997 Web at RALI
- 1999 Paying subscribers TSRali.com
- 2003 Transfer to Terminotix
- 2008 NSERC CRD project for new version
- 2011 TransSearch with translation spotting in beta

Julien Bourdaillet, Stéphane Huet, Philippe Langlais and Guy Lapalme. TransSearch: from a Bilingual Concordancer to a Translation Finder. *Machine Translation*, vol. 24, number. 3-4, p. 241-271, Dec 2010



TRANSSEAR



TRANSSEAR



TransBases (1986-2011)

Corpus	docs	K pairs	M words
Hansard	2 854	10 500,0	339,6
Senate	I 025	I 200,0	47,6
Canadian courts	13 068	3 200,0	142,6
Meteorological alerts	2 398	7,2	0,2
	19 345	14 907,2	530,0



TransBase Creation

- Finding texts and their translation
- Cleaning and normalization
- Sentence alignment
- Indexing
- Web based querying









Transpotting



- Word alignment combining HMM and IBM models
 - combine alignments from both directions
 - add a contiguity constraint
- Phrase-based transpotting
 - use the phrase table produced by MOSES



Postprocessing

on behalf of

Transpot	Frequency
au nom de	1424
au nom du	763
au nom des	683
de*	136
•••	
dans l'intérêt des	15
de la part de	13
dans*	13
part de	13
• • •	
pour l'ensemble de	1
parler pour	1
loin que*	1
le bien*	1



Postprocessing

- Must deal with transpotting errors
 - between 20% and 40%
- User will focus on the first ten transpots
 - as many correct
 - as many diversified



Postprocessing

- Filtering bad transpots using supervised classification
- Merging variants
 - inflectional form of canonical words
 - ignore some grammatical words
 - start merging by considering the most frequent ones
- Pseudo-relevance feedback
 - search for more frequent transpots in sentences with less frequent ones
 - develop a new transfer table from this specialized corpus and find new alignments





Evaluation corpora

- Hansard (8.3 million Fr-En sentence pairs)
- 5000 most frequent English queries to TransSearch between 2001 and 2007
- 5000 sentence pairs for each query
- Bilingual phrase lexicon for testing
 - 2 358 entries (284 DEV, 2074 TEST)
 - average number of translations (3.6 DEV, 3.9 TEST)



Evaluation process

• Transpotting

- only the ones with the source containing the query and the target containing a translation from the dictionary (180 000 DEV, 1.4M TEST)
- Translation
 - comparison with bilingual phrase lexicon
- Classification
 - hand annotation of 531 queries (12 144 examples)



Results

- Two best methods: PBM and C-HMM-bi
- Benefits from filtering and merging
- Limited impact of pseudo-relevance feedback
- But in the production version of TS3
 - at most 900 results for each query
 - transpotting
 - IBM-bi
 - filtering of transpots
 - 75% of grammatical words
 - only auxiliairies or numbers



Non web uses of TransSearch

- Web service
 - application embedding
 - call from a text processor (e.g. Word)
- Educational use
- Integration à la TransType



Other bilingual concordancers

- With transpotting
 - TotalRecall (Wu, 2003)
 - Chinese-English word alignment
 - LinearB (Callison-Burch, 2005)
 - basic translation spotting
 - Linguee.com
 - many language pairs of web text (free with ads)
- Without transpotting
 - Opus
 - multiple parallel corpora
 - WeBiText (CNRC/Terminotix)
 - meta crawler for bilingual texts



Why is TransSearch so popular ?

- Complements a terminology bank
- Not a bilingual dictionary
 87% of queries have more than one word
- Ready-made solution to translation problems
 - idiomatic expressions
 - non compositional translations
 - non literal translations



Questions ?

7 traductions de questions from the floor dans 7 occurrences

questions à partir du 1 parquet	l'auditoire à poser des questions	1
répondu à leurs questions 1 les députés lui posaient 1 des questions	Each group gave presentations which lasted between five and ten minutes, and at the end of all the presentations, I invited <u>questions from the floor</u> .	Chaque groupe a fait une allocution variant entre cinq et dix minutes et, à la fin de tous les exposés, j'ai invité <u>l'auditoire à poser des questions</u> .
questions1poser des questions1questions des spectateurs1l'auditoire à poser des1questions1		



Guy Lapalme : Mon bureau | Liste de mes cours | Mon compte utilisateur | Mes messages | Quitter

Anglais

EN001 - Guy Deville

CoBRA Demo > EN001 > Lecture de textes > Leçon

Lecture de textes

2

Retire later to live longer

Early **retirement**, as is **commonly** believed, does not help **retirees** to live longer and it may ever shorten one's life. This is the conclusion of a study published on October 21 by the British Medic Journal. The research involved **tracking** more than 3,500 **employees** working for Shell Oil in Texa over a 26-year period. The workers retired at 55, 60 or 65 and were **monitored** to see what effer their age at **retirement** had on their **lifespan**. Researchers considered factors such as gender ar **socio-economic status** in **ascertaining** whether retiring early is associated with better survival. seems **the findings** have **displaced** the myth that spending our golden years **at a leisurely par** away from **the daily grind** of the nine to five will increase our longevity. It appears that retiring lat **provides for** a longer life.

The results were **astonishing**. The life expectancy of employees who retired at 55 was significantly reduced compared with those who retired at the age of 65. The researchers conclude that: "Retiring early at 55 or 60 was not associated with better survival than retiring at 65 a cohort of past employees of the petrochemical industry. Mortality was higher in employees were tretired at 55 than in those who continued working." Leader of the research team Shan Tsai sai "Although some workers retired at 55 because of failing health, these results clearly show th early retirement is not associated with increased survival. On the contrary, mortality improved wir increasing age at retirement for people from both high and low socio-economic groups."

retirement : nom commun	(plur. retirements) 🤍		
He opposes pensioned retirement.	Il s'oppose à la retraite pensionnée.		
with other senators in wishing him well in his retirement, and a Merry Christmas!	sénateurs pour lui souhaiter la meilleure des retraites et un très joyeux Noël!		
records kept of refusals, retirements, moves, etc. As well, replacement questionnaires	indiquant les refus, les retraites , les déménagements, etc. De même, des questionnaires		
Thus, the gradual raising of the female retirement age, initiated several years	graduel de l'âge de la retraite des femmes, décidé il y a plusieurs années, continue		
	retirement : nom commun He opposes pensioned retirement. with other senators in wishing him well in his retirement, and a Merry Christmas! records kept of refusals, retirements, moves, etc. As well, replacement questionnaires Thus, the gradual raising of the female retirement age, initiated several years		

E-Lex

Deville et al. Université de Namur



TransType

File Edit Go Options Product overview Aperçu de la machine: The machine is controlled from a liquid crystal color touch screen where you can view your operation and application settings. Aperçu de la machine: La machine est contrôlée à partir d'un écran tactile à cristaux liquides vous permettant de visualiser vos paramètres d'application et d'opération.	őhour
Product overviewAperçu de la machine:The machine is controlled from a liquid crystal color touch screen where you can view your operation and application settings.La machine est contrôlée à partir d'un écran tactile à cristaux liquides vous permettant de visualiser vos paramètres d'application et d'opération.	ADOU
The machine is controlled from a liquid crystal color touch screen where you can view your operation and application settings.	ĺ
Printer settings can be pre-programmed for specific production job types and when such a job type is selected, the printer is set up automatically for the paper type and application.	r

