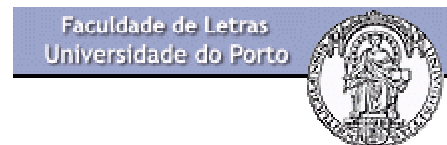


ParaMT: A Paraphraser for Machine Translation

Anabela Barreiro

barreiro_anabela@hotmail.com

FLUP & CLUP-Linguatca
New York University



Motivation

1. MT is a highly desired NLP application
2. MT has eluded researchers and developers for over 50 years
3. MT is increasingly more used and useful

BUT...

4. MT is still far from perfect!

Subject and Motivation

Support verb construction = predicate noun construction

is a multiword expression containing a verb with weak semantic value and a noun which is the predicate of the sentence.

Predicate nouns can be:

➤ **morphologically related to a verb**

fazer uma apresentação de = apresentar

to pay a visit to = to visit

➤ **autonomous**

fazer um mestrado - *mestrar

to have fun - *to fun

Subject and Motivation

Francisco Vieira adianta ainda que está a fazer um esforço no sentido de tomar uma decisão ainda esta semana, definindo se avança ou não para uma candidatura à RTLRS. [CdP]

Translation Engine	Translation Results
FreeTranslation	Francisco Scallop advances even if is it do an effort in the sense of <i>take a decision</i> still this week, defined advances or not for a candidacy to the RTLRS.
WorldLingo	advances despite he is to make an effort in the direction to still <i>take a decision</i> this week, defining if he advances or he does not stop a candidacy to the RTLRS.

I can't make a decision about anything these days. [Compara]

Translation Engine	Translation Results
Google	Eu não posso <i>fazer a uma decisão</i> sobre qualquer coisa estes dias.
Amikai	que eu não posso <i>fazer para uma decisão</i> sobre qualquer coisa estes dias.
FreeTranslation	Eu não posso <i>tomar uma decisão</i> sobre algo estes dias.
Babelfish	Eu não posso <i>fazer a uma decisão</i> sobre qualquer coisa estes dias.
WorldLingo	Eu no posso <i>fazer a uma deciso</i> sobre qualquer coisa estes dias.
E-Translation Server	Não posso <i>tomar uma decisão</i> sobre qualquer coisa estes dias.

Main Objectives

1. Build a body of lexical, syntactic and semantic knowledge around support verb constructions
2. Apply this linguistic knowledge to paraphrasing
3. Improve machine translation

Outcome

➤ **Port4NooJ**

- an open source, ontology driven Portuguese linguistic system, which integrates a bilingual extension for Portuguese-English machine translation

➤ **ReWriter**

- a monolingual paraphraser to pre-edit texts

➤ **ParaMT**

- a bilingual/multilingual paraphraser to be integrated in machine translation systems

Resources

- Port4NooJ - Publicly available at:
 - <http://www.nooj4nlp.net>
 - <http://www.linguateca.pt/Repositorio/Port4Nooj/>
- Based on:
 - NooJ linguistic environment (<http://www.nooj4nlp.net/>)
 - OpenLogos English-Portuguese dictionary (<http://logos-os.dfki.de/>)
OpenLogos is an open-source derivative of the Logos Machine Translation System
- Data Used
 - COMPARA (<http://www.linguateca.pt/COMPARA>)
 - METRA (<http://www.linguateca.pt/metra>)
 - Other corpora

Port4NooJ Dictionaries

Part of Speech Inflectional Paradigm Syntactic-Semantic Attributes English Transfer

Lemma

mesa, N+FLX=CASA+CO+surf+EN=table
 cair, V+FLX=ATRAIR+INMO+IntoType+EN=fall
 holandês, A+FLX=INGLÊS+AN+lang+EN=Dutch
 actualmente, ADV+FLX=FACILMENTE+TEMP+punc+pres+EN=nowadays
 alguém, PRO+IMPERS+INDEF+EN=somebody
 porque, RELINT+why+EN=why
 e, CONJ+JOIN+EN=and
 durante, PREP+TEMP+EN=during
 cada, DET+IMPERS+INDEF+SG+EN=each
 terceiro+NUM+ord+EN=one third

General dictionary sample representing all PoS, variable and invariable forms

a curto prazo, ADV+TEMP+EN=in the short run
 a favor de, PREP+CAUS+EN=in favor of
 cada um, PRO+INDEF+SG+EN=each one
 de quem, INT+ThatType+EN=whose
 quem quer que seja, REL+WhateverType+EN=whoever
 além disso, CONJ+COOR+EN=besides
 um quarto, NUM+frac+EN=one fourth

Sample of invariable compounds in the general dictionary

HIV, N+FLX=PORTUGAL+EN=HIV
 doença maniaco-depressiva, N+FLX=CASA+EN=manic-depressive disorder
 doença bipolar, N+FLX=CASA+EN=bipolar disorder
 asma, N+FLX=CASA+EN=asthma

Sample of the dictionary of Biomedical Terms

adro da igreja, N+FLX=MENINO+PL+end+EN=churchyard
 cabo de vassoura, N+FLX=MENINO+COTool+EN=broomstick
 bebida alcoólica, N+FLX=CASA+MA+liqu+EN=alcoholic drink+UNAMB
 bebida alcoólica, N+FLX=CASA+MA+liqu+EN=booze+slang
 cor de laranja, A+NAV+Apred+EN=orange
 sul-americano, A+FLX=ALTO+AN+des+EN=South American
 a curto prazo, ADV+LocTime+TEMP+EN=in the short run
 fora de serviço, ADV+STAT+phr+EN=out of order
 há muito tempo, ADV+LocTime+TEMP+puncpast+EN=a long time ago
 isto é, CONJ+COOR+EN=i.e.
 já não, CONJ+COOR+EN=no longer
 mesmo assim, CONJ+SUB+EN=even so
 juntamente com, PREP+ASSOC+EN=along
 à direita de, PREP+Loc+AT+EN=at the right of
 em conformidade com, PREP+ALOG+EN=in congruence with

Sample of the dictionary of Terms and Multiword Expressions

DicTUM

Amsterdão, N+PL+EN=Amsterdam
 Estados Unidos da América, N+PL+EN=United States of America
 África, N+PL+EN=Africa
 Extremo Oriente, N+PL+EN=Far East
 Mediterrâneo, N+FLX=ANO+PL+EN=Mediterranean
 Alpes Peninos, N+FLX=ALPES+PL+EN=Pennine Alps
 ONU, N+AN+EN=UN

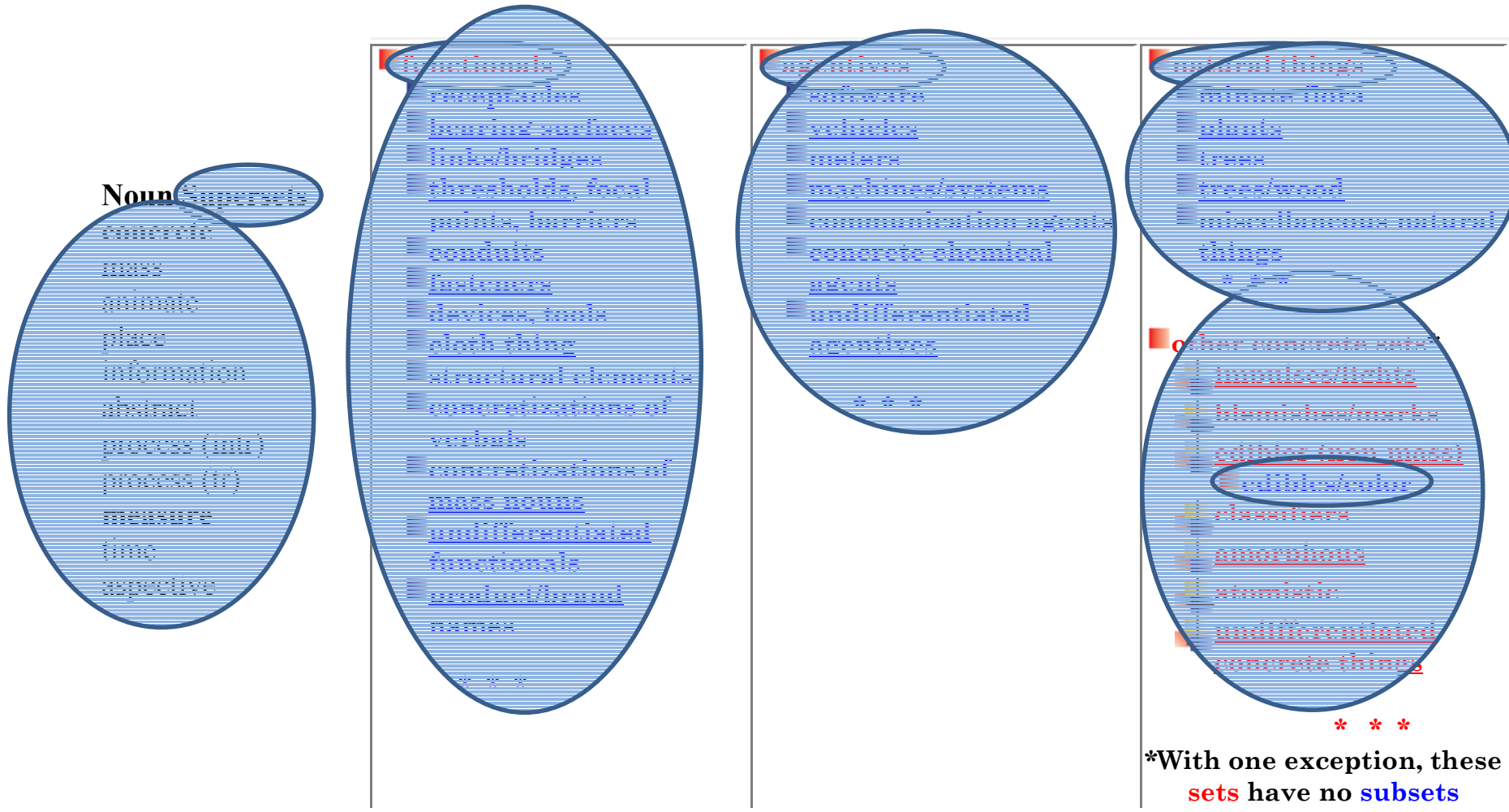
Sample of the dictionary of Proper Names

Syntactic-Semantic Ontology

- Representation abstract language
- Hierarchical taxonomy (sets, supersets and (sometimes) subsets)
- Based on Logos SAL ontology
- Integrated in the dictionary
- It represents both meaning (semantics), and structure (syntax)
- Over 1,000 categories

Syntactic-Semantic Ontology

Sets and **Subsets** of the **CONCRETE Noun Superset**
 Click on **CONCRETE Superset**, **sets** and **subsets** for explanations



Syntactic-Semantic Ontology

Category

agentives
 software
 concrete chemical agents
 machines/systems
 vehicles
 meters
 communication agents
 functionals
 devices/tools
 fasteners
 bearing surfaces
 receptacles
 conduits
 thresholds/focal points/barriers
 links/bridges
 cloth things
 structural elements
 concretizations of verbals
 concretizations of mass nouns
 product/brand names
 natural things
 minute flora
 plants
 trees
 trees/wood
 misc. natural things
 edibles (non-mass)
 edibles/color
 impulses/lights
 blemishes/marks
 classifiers
 amorphous
 atomistic

Subcategory

CO+undagt
 CO+soft
 CO+chem
 CO+mach
 CO+vehic
 CO+meter
 CO+comm
 CO+undfunc
 CO+tool
 CO+fast
 CO+surf
 CO+recp
 CO+cond
 CO+barr
 CO+link
 CO+cloth
 CO+struc
 CO+verb
 CO+mass
 CO+brand
 CO+nat
 CO+flora
 CO+plant
 CO+tree
 CO+trwd
 CO+mnat
 CO+ednm
 CO+edcol
 Col+ight
 CO+blem
 CO+class
 CO+amor
 CO+atom

Examples in English

See subsets
routine
catalyst, warhead
battery, camera
truck, ship
clock, gauge
radio, radar
trinket, ornament
pliers
nail, tendon
table, shelf
bottle, barrel
chute, artery
wall, door
circuit, nerve
shirt, blanket
spar, bone
threading
acid lining
Windows NT
 See subsets
algae, spore
rose, weed
apple, willow
oak, maple
pebble, iceberg
pork chop
orange, cherry
lamp, beam
scratch, freckle
element
breeze, tide
electron, atom

Examples in Portuguese

See subsets
rotina, ficheiro
ácido sulfúrico
máquina fotográfica
automóvel
manómetro
rádio
ornamento
alicate
prego
mesa
garrafa
artéria
porta
circuito
camisola
osso

Windows NT
 See subsets
alga
erva
macieira
carvalho
iceberg
costoleta
laranja
lâmpada
sarda
elemento
brisa
átomo

Categories of
CONCRETE nouns

Syntactic-Semantic Ontology

ME - MEASURE Noun Sets and Subsets		
Sets and Subsets	Mnemonics (= SynSem)	Examples
abstract concepts measured by unit	ME+abs	<i>humidity, length</i>
discrete measurable concepts	ME+dis	<i>sum, increment</i>
units of measure	ME+unit	See subsets
units of weight	ME+unit+wt	<i>ounce, pound</i>
units of velocity	ME+unit+vel	<i>mph, megahertz</i>
units of volume measure	ME+unit+vol	<i>gallon, liter</i>
units of temperature	ME+unit+temp	<i>degrees celsius</i>
units of energy/force	ME+unit+ener	<i>watt, horsepower</i>
measurement systems	ME+unit+sys	<i>fahrenheit, kelvin</i>
units of duration	ME+unit+dur	<i>hour, minute, year</i>
specialized units of measure	ME+unit+spec	<i>oersted, ohm, phon</i>
units of money/value	ME+unit+value	<i>dollar, euro, forint</i>
units of linear/area measure	ME+unit+lin	<i>inch, yard, mile</i>
general undifferentiated measure	ME+undif	<i>degree, gross, share</i>

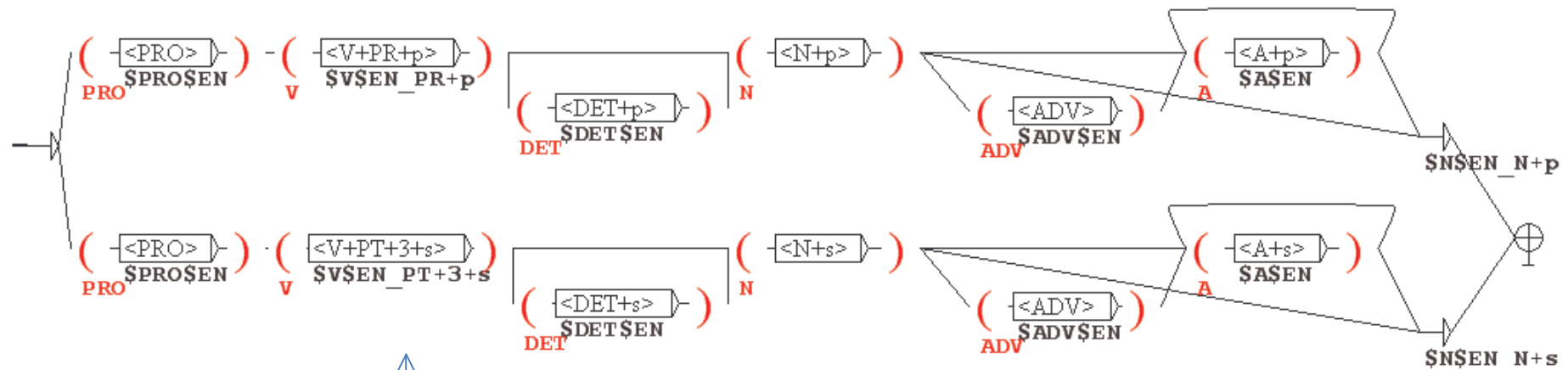
Categories of
MEASURE nouns

Inflectional and Derivational Description

The screenshot displays the NooJ Community Edition interface with several windows open, each showing a different type of linguistic rule. Blue boxes with arrows point to specific windows:

- Noun Inflectional Paradigm**: Points to the 'Verbs.flx' window, which contains rules for verb inflection. A blue box highlights the rule for plural adding -s: `# plural adding -s` and `# Logos PAT 99`.
- Verb Inflectional Paradigm**: Points to the 'Verbs.flx' window, which contains rules for verb inflection. A blue box highlights the rule for regular ending in -ar (1st conjugation), like falar: `# regular ending in -ar (1st conjugation), like falar`.
- Adjective Inflectional Paradigm**: Points to the 'Verbs.flx' window, which contains rules for verb inflection. A blue box highlights the rule for ending in -o, like alto: `# ending in -o, like alto`.
- Adverb Inflectional Paradigm**: Points to the 'Adv.s.flx' window, which contains rules for adverb inflection. A blue box highlights the rule for ending in -velmente, like amigavelmente: `# ending in -velmente; like amigavelmente`.
- Determiner Inflectional Paradigm**: Points to the 'Dets.flx' window, which contains rules for determiner inflection. A blue box highlights the rule for definite article, like o: `# definite article, like o`.
- Pronoun Inflectional Paradigm**: Points to the 'Prons.flx' window, which contains rules for pronoun inflection. A blue box highlights the rule for demonstrative pronoun, like este/esse/aquele: `# demonstrative pronoun, like este/esse/aquele`.
- Interrogative Pronoun Inflectional Paradigm**: Points to the 'OtherProns.flx' window, which contains rules for other pronoun inflection. A blue box highlights the rule for interrogative pronoun, like qual: `# interrogative pronoun, like qual`.
- Nominalization Derivational Paradigm**: Points to the 'Noms.flx' window, which contains rules for nominalization. A blue box highlights the rule for 3são: `# 3são:`.

Paraphrasing and Translation Grammars



Graph to translate simple sentences

Translation and bilingual paraphrasing of simple sentences

Concordance for Text PTSimpleSentences-3pl.not

Clear Concordance | 4 characters before, and 5 after. Display: Inputs Outputs

Text	Before	Seq.	After
	Elas visitam cidades.	Elas são homens honestos/they are honest men	. Elas são mulheres simpáticas. Eles
	Elas visitam cidades.	Elas são homens honestos/they are honest men	. Elas são mulheres simpáticas. Eles
	comida. Eles vendem sapatos.	Ela compra casas/the purchase houses	. Elas sonham acordados. Eles sonham
	comida. Eles vendem sapatos.	Elas compram casas/they purchase homes	. Elas sonham acordados. Eles sonham
	comida. Eles vendem sapatos.	Elas compram casas/the buy houses	. Elas sonham acordados. Eles sonham
	comida. Eles vendem sapatos.	Elas compram casas/they buy houses	. Elas sonham acordados. Eles sonham

Query 6/29

Explicit Marking of Derivation and Support Verb

Verb entries:

- Identification of derivational paradigms for **nominalizations** (annotation *NDRV*) and **predicate adjectives** (annotation *ADRV*)
- Link to the derived noun's **support verbs** and to the adjective's **copula verbs** (annotation *VSUP* and annotation *VCOP*)

adaptar,V+FLX=FALAR+Aux=1+INOP57+Subset=132+EN=adapt+**VSUP=fazer**+**DRV=NDRV00**:CANÇÃO

azedar,V+FLX=LIMPAR+Aux=1+OBJTRundif98+Subset=740+EN=sour+**VCOP=estar**+**DRV=ADRV00**:ALTO

Explicit Marking of Derivation and Semantic Verb Association

Adjective entries:

- Identification of derivational paradigms for **adverbializations** (annotation *AVDRV*)

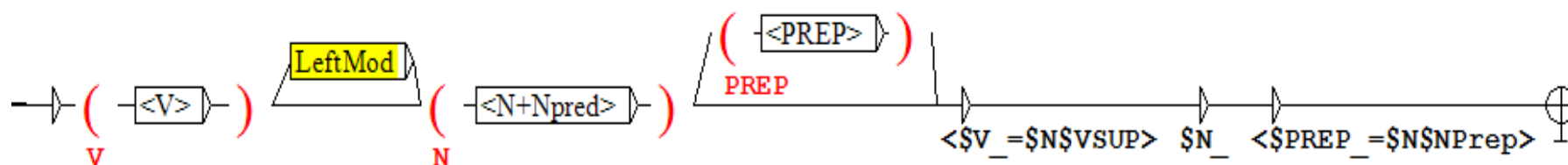
literal,A+FLX=PRINCIPAL+IN+symb+EN=literal+**DRV=AVDRV00**:LITERALMENTE

Autonomous predicate nouns:

- Identification of **autonomous predicate nouns** (annotation *Npred*)
- Identification of a semantically related verb

curso,N+FLX=ANO+**Npred**+IN+inst+EN=course+VSUP=tirar+**VRB=estudar**+NPrep=de+Det=um

ReWriter: a Monolingual Paraphraser



gosto de ver o comboio a	fazer corridas /correr	à velocidade máxima ao longo
io de cheque especial para	fazer doações /doar	às entidades que escolher. A
cores e, quando é preciso ir	fazer filmagens/filmar	fora do estúdio, às vezes fic
o que queria trocar de pares e	fazer um jogo /jogar	ao melhor de três sets , mas
o dra deu-me um papel para	fazer uma lista de/listar	todas as coisas boas que ex
res foram à caracterização	fazer uns retoques/retocar	, outros estão a descansar n

Recognition and monolingual paraphrasing
of support verb constructions
(support verb construction / morphologically related lexical verb)

ReWriter: Examples

o cirurgião Faivre, ao	fazer uma amputação	uma amputação	
o cirurgião Faivre, ao	fazer uma amputação	uma amputação	uma amputação
o cirurgião Faivre, ao	fazer uma amputação	uma amputação	uma amputação
omista britânico, conseguiu	fazer uma transfusão de sangue	uma transfusão de sangue	de um cão para outro. A juiza
omista britânico, conseguiu	fazer uma transfusão de sangue	uma transfusão de sangue	de um cão para outro. A juiza
os pacientes que precisam	fazer uma transfusão de sangue	uma transfusão de sangue	colaborando também com a c
os pacientes que precisam	fazer uma transfusão de sangue	uma transfusão de sangue	colaborando também com a c
os pacientes que precisam	fazer uma transfusão de sangue	uma transfusão de sangue	colaborando também com a c
os pacientes que precisam	fazer uma transfusão de sangue	uma transfusão de sangue	colaborando também com a c

Elementary SVC > Lexical Verb

Elementary SVC > non-elementary SVC

realizar/effectuar

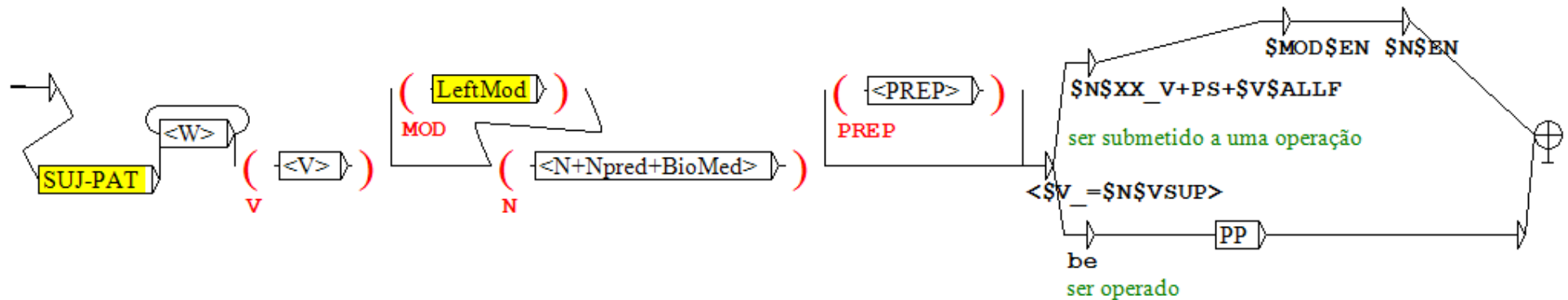
Elementary SVC > *sujeitar-se a*

submeter-se a

ONLY if the SUBJECT is a patient

Recognition and paraphrasing of elementary support verb constructions co-occurring with predicate nouns of the biomedical field

(support verb construction / lexical verb or stylistic variant / non-elementary support verb construction)



ReWriter: Example of Possible Application

1. Toda a pessoa tem o direito de tomar parte na direcção dos negócios, públicos do seu país, quer directamente, quer por intermédio de representantes livremente escolhidos.

2. Toda a pessoa tem direito de acesso, em condições de igualdade, às funções públicas do seu país.

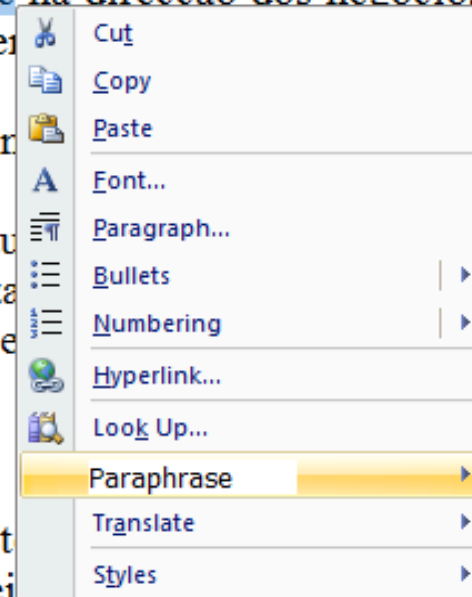
3. A vontade do povo é o fundamento da autoridade dos poderes públicos: e deve exprimir-se através de eleições honestas a realizar-se em sufrágio universal e igual, com voto secreto ou segundo processo equivalente e sufrágio directo e a liberdade de voto.

Artigo 22º

Toda a pessoa, como membro da sociedade, tem direito de participar livremente na direcção dos negócios, públicos do seu país, quer directamente, quer por intermédio de representantes livremente escolhidos. Todo o cidadão tem direito de acesso, em condições de igualdade, às funções públicas do seu país. A vontade do povo é o fundamento da autoridade dos poderes públicos: e deve exprimir-se através de eleições honestas a realizar-se em sufrágio universal e igual, com voto secreto ou segundo processo equivalente e sufrágio directo e a liberdade de voto.

Artigo 23º

1. Toda a pessoa tem direito ao trabalho, à livre escolha do trabalho,



Interactive ReWriter
for word processing applications
such as text editing

ParaMT: a Bilingual/Multilingual Paraphraser

a fazer um estágio para	dar aulas de/teach	religião, mas não se import
m -- os filhos -- juntos e	fizeram a mudança para/change	Johannesburg, e ensinaram
. Necessitava apenas de	ter a certeza de/know	que não escapara à sua
ente hipotética. -- Deves	ter alguma ideia/know	. Dorothy andava a fazer um
. não podemos deixar de	ter cautela/beware	. Pobre Caro, pensou Lync
ra dos chinelos, antes de	ter chance de/can	mudar de idéia. Como pos
ope a Jean, esta pareceu	ter dificuldade em/avoid	olhá-lo nos olhos. Deixou
ao Kiss dela. Apesar de	ter falta de/lack	amor-póprio, isso não sign
igos e imprensa estava a	ter lugar /occur	numa longa galeria com car
guiu ter filhos. -- Tens de	ter mão /control	nessa confusão toda. Sam
spondi, minha mãe deve	ter medo de/fear	cobras. Eu disse no Gabin
da loja antes de ele	ter tempo de/could	chamar a brigada de narcó
a triste aventura havia de	ter um fim/finish	.
Ela ouvira a tia Velma	ter uma discussão com/argue	Jack acerca de mostarda r
de olhos fechados para	ter uma ideia de/know	como seria ser cego e
ter paciência.» «Voltei a	ter uma imensa vontade de/want	viver. A conversa parecia :

Recognition and bilingual paraphrasing of support verb constructions
(Portuguese support verb construction / corresponding English verb)

Preliminary Quantitative Results

500 sentences

100 for each elementary support verb

	SVC Recognition Precision	SVC Recognition Recall	SVC Paraphrasing Precision
Pôr	73/73 - 100%	73/100 - 73%	72/73 - 98.6%
Tomar	75/75 - 100%	75/100 - 75%	68/73 - 93.1%
Ter	65/65 - 100%	65/100 - 65%	59/65 - 90.7%
Dar	57/60 - 95%	57/100 - 57%	46/51 - 90.1%
Fazer	43/45 - 95.5%	43/100 - 43%	40/45 - 88.8%
Average	62.6/63.6 - 98.4%	62.6/100 - 62.6%	57/61 - 93.4%

Evaluation of recognition and paraphrasing
of support verb constructions

Conclusions

- Linguistic knowledge applied to a machine translation system improves its output quality.
- Effective results from linguistically based research on paraphrases can save substantial effort and resources employed by statistically based machine translation systems

Acknowledgements

Thank you for your attention!

This work was partly supported by grant SFRH/BD/14076/2003 from *Fundação para a Ciência e a Tecnologia*, co-financed by POSI and partly by *Fundação para a Computação Científica Nacional*.



ParaMT: A Paraphraser for Machine Translation

Anabela Barreiro

barreiro_anabela@hotmail.com

FLUP & CLUP-Linguatca
New York University

Faculdade de Letras
Universidade do Porto

