

Guide

- Welcome
- Context
 - The UNL
 - The UNDL Foundation roadmap
 - The UNLversity
 - The workshop
- Morning: First Task (Corpus)
- Afternoon: Second Task (Dictionaries)



The Universal Networking Language (UNL)



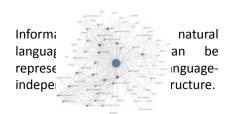
translation knowledge representation

Universal Networking Language





Foundations (I)



Foundations (II)

Relations between languages can be expressed in three layers:

CONCEPTS

= Universal Words (UWs)

CONCEPT MODIFIERS

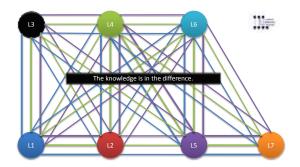
= Attributes

RELATIONS BETWEEN
CONCEPTS

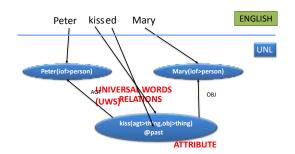
= Relations



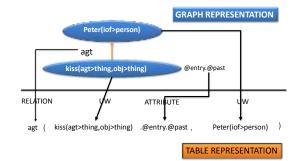
The Universal NETWORKING Language



Natural Language-to-UNL (UNL-ization)



Syntax of UNL



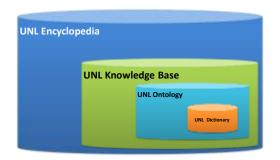
UNL document



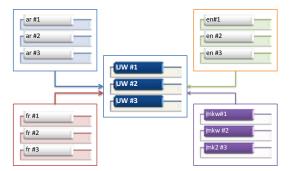
$UNL \ (\text{http://anydomain/anydocument.unl})$



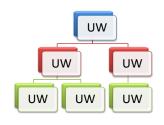
Lexical Databases



UNL Dictionary



UNL Ontology



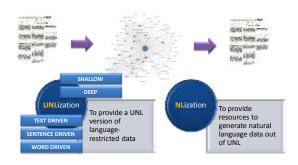
UNL Knowledge Base



UNL Encyclopedia



The UNL System



Is this translation?



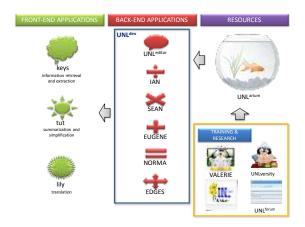
Roles of UNL

representing organizing retrieving extracting inferring generating

monolingual multilingual information



the road map



UNLversity

☐ UNL Symposium

□ A1

- ☐ Internship Program ☐ UNL Grammar Workshop (A1, A2, B1, B2, C1, C2)
 - ☐ European Chapter (Geneva, February 2012)
 - ☐ Indian Chapter (Mumbai, June 2012)☐ Asian Chapter (Macau, October 2012)

 - ☐ Middle East Chapter (Kuwait, February 2013)
 - ☐ African Chapter (Johannesburg, June 2013)
 ☐ American Chapter (Brasilia, October 2013)



Schedule

June 11th, 2012 - Monday 09:00-10:00 10:00-12:00 10:00-12:00 I UNIL-NL diction 14:00-17:00 II UNIL-NL diction June 12th, 2012 - Tuesday 09:00-12:00 III – Morphology 17:00-17:00 IV – NL dictionary I – Corpus II – UNL-NL dictionary June 13th, 2012- Wednesser,
09:00-12:00 V – UNL-NL 9..
14:00-17:00 V – UNL-NL grammar (II)
June 14th, 2012 - Thursday
09:00-12:00 VI – NL-UNL grammar (I)
09:00-17:00 VI – NL-UNL grammar (II)

14:00-17:00 Discussion



Goals

- To build the basic modules of a NL-UNL (analysis) grammar
- To build the basic modules of a UNL-UNL (generation) grammar



Deliverables

During the workshop

- Experimental corpus (50 sentences)
- UNL-NL dictionary (corresponding to the corpus)
- NL dictionary (corresponding to the corpus)
- UNL-NL grammar (corresponding to the corpus)
- NL-UNL grammar (corresponding to the corpus) Before September 15th
- Reference corpus (500 sentences)
- UNL-NL dictionary (corresponding to the corpus)
- NL dictionary (corresponding to the corpus)
- UNL-NL grammar (corresponding to the corpus)
- NL-UNL grammar (corresponding to the corpus)



Warnings

- Doubts are allowed: don't be afraid or shy.
- This is an ongoing initiative: we don't have all the answers yet.
- This is not a competition.

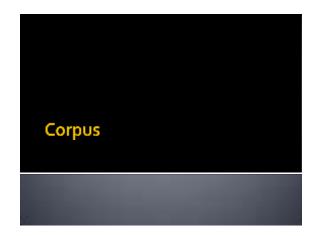


Participants

- Aadil Kak (Kashmiri)
- Arulmozi Selvaraj (Tamil) Balaji Jagan (Tamil)
- Laishram Rishikanta Meitei (Manipuri)
- Navanath Saharia (Assamese)
- Niladri Sekhar Dash (Bengali)
- Parameswarappa S (Kannada)
- Parteek Kumar (Punjabi)
- Pinkey Nainwani (Sindhi)
- Ranjan Das (Oriya)
- Renuka Devi (Telugu) Sachin Pawar (Marathi)
- Shailendra Kumar (Hindi)
- Trupti Nisar (Gujarati)

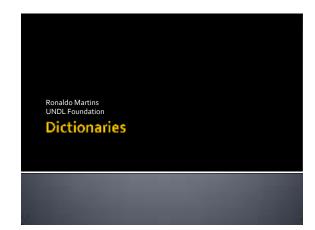






Task #1: Corpus

- Translate the 50 sentences of Corpusto_eng.txt into your native language. Be as close as possible to the original.
- Save the translated text (without the English) original) in a plain text (.txt) file with UTF-8 encoding and upload it to UNLWEB>UNLDEV>PROJECTS>IAN>NL FILES.
- Upload the file <u>Corpus5o_unl.txt</u> to UNLWEB>UNLDEV>PROJECTS>EUGENE>UNL **DOCUMENTS**



Dictionaries

- NL-UNL Dictionary (Analysis)
 - Enumerative (word forms) [table] {2883} "table" (POS=NOU, NUM=SNG) <eng,o,o>; [tables] {2883} "table" (POS=NOU, NUM=PLR) <eng,o,o>;

[foot] {2883} "foot" (POS=NOU,NUM=**SNG**) <eng,o,o>; [feet] {2883} "foot" (POS=NOU,NUM=**PLR**) <eng,o,o>;

- UNL-NL Dictionary (Generation)
 - Generative (base forms)
 - [table] {2883} "table" (POS=NOU,NUM=SNG,PAR=M2) <eng,o,o>;
 - [foot] {2883} "100284665" (POS=NOU,PAR=M1,FLX(PLR:="feet";)) <eng,o,o>;

Building dictionaries



UNLarium



Dictionary Specs

- Dictionary Specs
 - Dictionary structure
 - a plain text file (.txt)
 - one entry per line
 - entries must have the following format:

[NLW] {ID} "UW" (ATTR , ...) < LG , FRE , PRI >; COMMENTS

[NLW]

[NLW] {ID} "UW" (ATTR,...) < LG, FRE, PRI >; COMMENTS

- a multiword expression: [United States of America]
- a compound: [hot-dog]
- a simple word: [happiness]
- a simple morpheme: [happ]
- a complex structure: [[bring] [back]]
- a non-motivated linguistic entity: [g]

$\{ID\}$

[NLW] {ID} "UW" (ATTR,...) < LG, FRE, PRI >; COMMENTS

The unique identifier (primary-key) of the entry.

"UW"

[NLW] {ID} "UW" (ATTR,...) < LG, FRE, PRI >; COMMENTS

 The Universal Word of UNL. This field can be empty if a word does not need a UW.

(ATTR, ...)

[NLW] {ID} "UW" (ATTR,...) < LG, FRE, PRI >; COMMENTS

- The list of features of the NLW.
- Attributes should be separated by ","
- It can be:

 - a list of simple features: (NOU, MCL, SNG)
 a list of attribute-value pairs: (pos=NOU, gen=MCL, num=SNG)
 a list of transformation rules : (plural:="oo"."ee")
 - Replacement
 - <ATTRIBUTE>":="<SOURCE>":"<TARGET>
 - plural:="oo":"ee

 - Left appending
 <attribute>":="<left Deletion>"<"<left Addition>
 - not:=<"un"

 - Right appending
 «ATTRIBUTE»":="<RIGHT ADDITION>">"<RIGHT DELETION>
 - plural:=y>ies

<LG, FRE, PRI>

[NLW] {ID} "UW" (ATTR,...) <LG,FRE,PRI>;

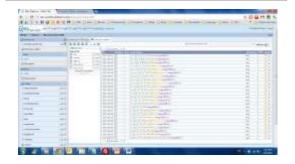
COMMENTS

- FLG
 - The two-character language code according to ISO 639-1.
- FRE
 - . The frequency of NLW in natural texts. Used for natural language analysis (NL-UNL). It can range from o (less frequent) to 255 (most frequent).
- PRI
 - The priority of the NLW. Used for natural language generation (UNL-NL). It can range from 0 to 255.

Examples

[book]{\(\) "book" (\) N, NOU, SNG, M2) < eng, 0, 0 >; [books]{\(\) "book" (\) N, NOU, PLR, M2) < eng, 0, 0 >; [car]{\(\) "car" (\) N, NOU, PLR, M2) < eng, 0, 0 >; [car]{\(\) "car" (\) N, NOU, SNG, M2) < eng, 0, 0 >; [car]{\(\) "car" (\) N, NOU, SNG, M2) < eng, 0, 0 >; [dr.]{\(\) "dr" (\) N, NOU, SNG, M2) < eng, 0, 0 >; [dr.]{\(\) "dr" (\) N, NOU, SNG, M2) < eng, 0, 0 >; [dass]{\(\) "drawd" (\) N, PEN, SNGT, M0) < eng, 0, 0 >; [Edward]{\(\) "Edward]{\(\) "Edward]{\(\) "Edward]{\(\) "Edward]{\(\) Nou, SNGT, M0) < eng, 0, 0 >; [alone]{\(\) "drawd" (\) N, PEN, SNGT, M0) < eng, 0, 0 >; [alone]{\(\) "drawd" (\) "AD, M0) < eng, 0, 0 >; [alone]{\(\) "drawd" (\) "AD, M0) < eng, 0, 0 >; [eatwilt]{\(\) "early" (\) "AAAW (\) Noung, 0, 0 >; [eatwilt]{\(\) "early" (\) "AAAW (\) Noung, 0, 0 >; [eatwilt]{\(\) "early" (\) "AAW (\) Noung, 0, 0 >; [eatwilt]{\(\) "early" (\) "AAW (\) Noung, 0, 0 >; [eatwilt]{\(\) "early" (\) "early" (\) "VER, TST, 1, PS, PS, PS, M7) < eng, 0, 0 >; [arrived]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrived]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrived]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrived]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrived]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1) < eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1, eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1, eng, 0, 0 >; [arrive]{\(\) "arrive" (\) VER, TST, 1, PS, M1, eng, 0, 0 >; [arrive]{\(\) "arrive" (

Using dictionaries (IAN and EUGENE)



Task #2: NL-UNL Dictionary

- Extract the word list (i.e., the set of all distinct word forms) appearing in your translation of the Corpus 50
- Create the NL-UNL dictionary for all the word forms following the English model available at English Analysis Dictionary 50. Use only the tags available at the tagset. For further information on the dictionary structure, see Dictionary Specs.
- Save the NL-UNL dictionary in a plain text (.txt) file with UTF-8 encoding and upload it to UNLWEB>UNLDEV>PROJECTS>IAN>DICTIONARIES.

Inflectional Paradigms

- Inflectional paradigms (regular) x inflectional rules (irregular)
 - book, books => inflectional paradigm (absolutely regular)
 - man, men => inflectional paradigm (irregular, but very common in several lexemes: "service man", "gentleman", "superman", etc.)
 - foot, feet => inflectional rules (irregular and uncommon)

Rule types: simple rules

prefixation

CONDITION := "ADDED" < DELETED;

suffixation

CONDITION := DELETED > "ADDED";

infixation

CONDITION := [REFERENCE] > "ADDED"; CONDITION := "ADDED" < [REFERENCE];

replacement

CONDITION := DELETED : "ADDED"; CONDITION := [INTERVAL] : "ADDED";

Rule types: complex rules

- circumfixation
 - CONDITION := "ADDED" < DELETED , DELETED
 > "ADDED";
- prefixation + infixation
 - CONDITION := "ADDED" < DELETED , DELETED > "ADDED";
- infixation + suffixation
 - CONDITION := DELETED > "ADDED" , "DELETED" > "ADDED";

Prefixation

RULE	BEHAVIOR	BEFORE	AFTER
X:="y"<"z";	if X replace the string "z" by the string "y" in the beginning of the string	z abc	y abc
X:="y"<1;	if X replace the first character of the string by "y"	zabc	yabc
X:="y"<0;	if X add the string "y" to the beginning of the string	zabc	y zabc
X:="y"<;	if X add the string "y" to the beginning of the string (idem previous)	zabc	y zabc
X:="y"<<0;	if X add the string "y" and a blank space to the beginning of the string	zabc	y zabc
X:="y"<<;	if X add the string "y" and a blank space to the beginning of the string (idem previous)	zabc	y zabc

Suffixation

RULE	BEHAVIOR	BEFORE	AFTER
X:="z">"y";	if X replace the string "z" by the string "y" in the end of the string	abcz	abc y
X:=1>"y";	if X replace the last character of the string by "y"	abcz	abc y
X:=o>"y";	if X add the string "y" to the end of the string	abcz	abcz y
X:=>"y";	if X add the string "y" to the end of the string (idem previous)	abcz	abcz y
X:=0>>"y";	if X add a blank space and the string "y" to the end of the string	abcz	abcz y
X:=>>"y";	if X add a blank space and the string "y" to the end of the string (idem previous)	abcz	abcz y

Infixation

RULE	BEHAVIOR	BEFORE	AFTER
X:=[2]>"y";	if X add "y" to the right of the second character	abc	ab y c
X:="y"<[3];	if X add "y" to the left of the third character	abc	ab y c
X:=["b"]>"y";	if X add "y" to the right of "b";	abc	ab y c
X:="y"<["c"];	if X add "y" to the left of "c"	abc	ab y c

Replacement

RULE	BEHAVIOR	BEFO RE	AFTE R
X:="y";	if X replace the whole entry by "y"	Х	у
X:="z":"y";	if X replace the string "z" by "y"	a z bc	a y bc
X:=[2-3]:"y";	if X replace the second to the third character by "z"	a bc z	a y z

Observations (I)

- Rules will only be applied if all conditions are true
 - X:="y"<"z"; ("zabc" changes to "yabc", but "abc" remains "abc" since there is no "z" to be replaced)
- String fields are necessarily continuousX:="aaa"<"xyz"; ("xyzbbb" changes to "aaabbb", but "bxbybz" remains "bxbybz" since there is no continuous string "xyz" to be replaced)
- Each action is applied only once (i.e, rules are not exhaustive)
- PLR:=o>"s"; ("X" becomes "Xs", and not "Xssssss...")
 The replacement rule applies only once to the same
- X:="a":"b"; ("aaa" becomes "baa" and not "bbb")

Observations (II)

In prefixation and suffixation rules, the part to be deleted may be represented by the number of characters (without quotes)

PLR:= "X"<"A"; (ABC becomes XBC) PLR:="X"<1; PLR:= "XY"<"AB"; PLR:="XY"<2; (ABC becomes XYC) PLR:="">"X"; PLR:= 0>"X"; (ABC becomes ABCX) PLR:="C">"X"; = PLR:=1>"X"; (ABC becomes ABX) PLR:="BC">"XY"; = PLR:= 2>"XY"; (ABC becomes AXY)

Observations (III)

In infixation rules, the position of the addition may be made with reference to the end of string by using "-".

	RULE	BEHAVIOR	BEFOR E	AFTER
	X:=[1]>"y";	if X add "y" to the right of the first character	abc	a y bc
I	X:=[-1]>"y";	if X add "y" to the right of the last character	abc	ab y c
ı	X:="y"<[2];	if X add "y" to the left of the second character	abcde	a y bc
	X:="y"<[-2];	if X add "y" to the left of the second character	abcde	abc y de

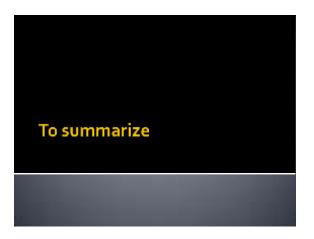
Observations (IV)

- In replacement rules, the part to be deleted may be omitted if
 - the whole string is to be replaced

 PLR:="ABC":"XYZ"; ISTHE SAME AS
- PLR:="XYZ"(ABC becomes XYZ)
 In replacement rules, the part to be deleted may be represented by an interval of characters in the format [beginning-end]
- LR:="R":X";=PLR:=[2-2]:"X";(ABC becomes AXC)
 The symbol "^" is used for negation ("^MCL" means "not MCL")
- "NOU&^MCL:="x":"y"; (If NOU and not MCL then replace "x" by "y")
 "<<" and ">>" add blank spacests!X:="a"<<"br/>b" ("bc" becomes "a bc" and not "abc")
 - X:="a"<<"b" ("bc" becomes "a bc" and not "abc")

Task #3: NL-UNL Dictionary

- Localize the UNL-NL dictionary available at English Generation Dictionary 50. The localized version must reflect the word list of your translated corpus. Use only the tags available at the <u>tagset</u>. For further information on the dictionary structure, see Dictionary
- Save the UNL-NL dictionary in a plain text (.txt) file with UTF-8 encoding and upload it to UNLWEB>UNLDEV>PROJECTS>EUGENE>DICT IONARIES.



Task #1: Corpus

- Translate the 50 sentences of <u>Corpus50 eng.txt</u> into your native language.
 Be as close as possible to the original.
- Save the translated text (without the English original) in a plain text (.txt) file with UTF-8 encoding and upload it to UNLWEB>UNLDEV>PROJECTS>IAN>NL FILES.
- Upload the file <u>Corpus5o_unl.txt</u> to UNLWEB>UNLDEV>PROJECTS>EUGENE>UNL DOCUMENTS without doing any change to it

Task #2: NL-UNL Dictionary

- Extract the word list (i.e., the set of all distinct word forms) appearing in your translation of the Corpus 50
- Create the NL-UNL dictionary for all (and only) the word forms appearing in the corpu. In order to do that, follow the English model available at <u>English</u> <u>Analysis Dictionary 50</u>. Use only the tags available at the <u>tagset</u>. For further information on the dictionary structure, see <u>Dictionary Specs</u>.
- Save the NL-UNL dictionary in a plain text (.txt) file with UTF-8 encoding and upload it to UNLWEB>UNLDEV>PROJECTS>IAN>DICTIONARIES.

Task #3: NL-UNL Dictionary

- Localize the UNL-NL dictionary available at <u>English Generation Dictionary 50</u>. The localized version must reflect the word list of your translated corpus. Use only the tags available at the <u>tagset</u>. For further information on the dictionary structure, see <u>Dictionary</u> <u>Specs</u>.
- Save the UNL-NL dictionary in a plain text (.txt) file with UTF-8 encoding and upload it to UNLWEB>UNLDEV>PROJECTS>EUGENE>DICT IONARIES.