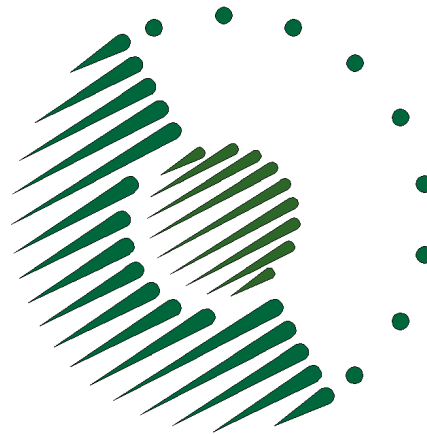


Panlingual Localization

Jonathan Pool

Utilika Foundation

pool@utilika.org



LISA @ Berkeley

5 August 2009

The Localization Industry Standards Association

Panlingual Localization

1. Goal
2. Strategy
3. Method
4. Demonstration
5. Challenges
6. Solutions
7. Discussion



1. Goal: Panlingual Localization

Make any information instantly and freely accessible in any language.

Google: 127

~~Universal Declaration of Human Rights: 360~~

Wikipedia: 235



2. Strategy: Lemmatic Communication

Все счастливые семьи
похожи друг на друга,
каждая несчастливая
семья несчастлива по-
своему.



все	счастливый	семья	похожий
все	несчастливый	семья	особенный



Les familles heureuses
se ressemblent toutes;
les familles
malheureuses sont
malheureuses chacune
à leur façon.



tout	heureux	famille	similaire
tout	malheureux	famille	particulier

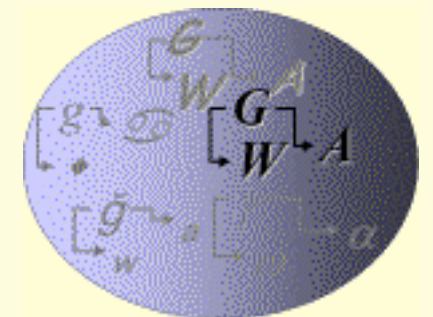
3. Method: Mine Diverse Lexical Resources

Thesauri

Dictionaries



WordNets



Glossaries

Arrest: $\cap \mathbb{J} \triangleright \sigma^{\mathfrak{b}}$: Tigujaunig: Arrestation

The act of placing a person in custody, according to law. The powers of ordinary citizens and peace officers to arrest a person are set out in the *Criminal Code*, 1996, Part XVI.

Arson: ΔΡΟΚΟσ^ς: Ikitittiniq: Crime d'incendie

The crime of deliberately setting fire to property.
Criminal Code, 1996, sections 433-436.

Locale Repositories

SE (English: Sweden)

An tSualainn ga

isveç az .

İsveç ·tr·

Iswidhan 'so'

Rootsi ·et·

Ruotsi ·f₁·

Ruotta *se* ·

Schweden ·de·

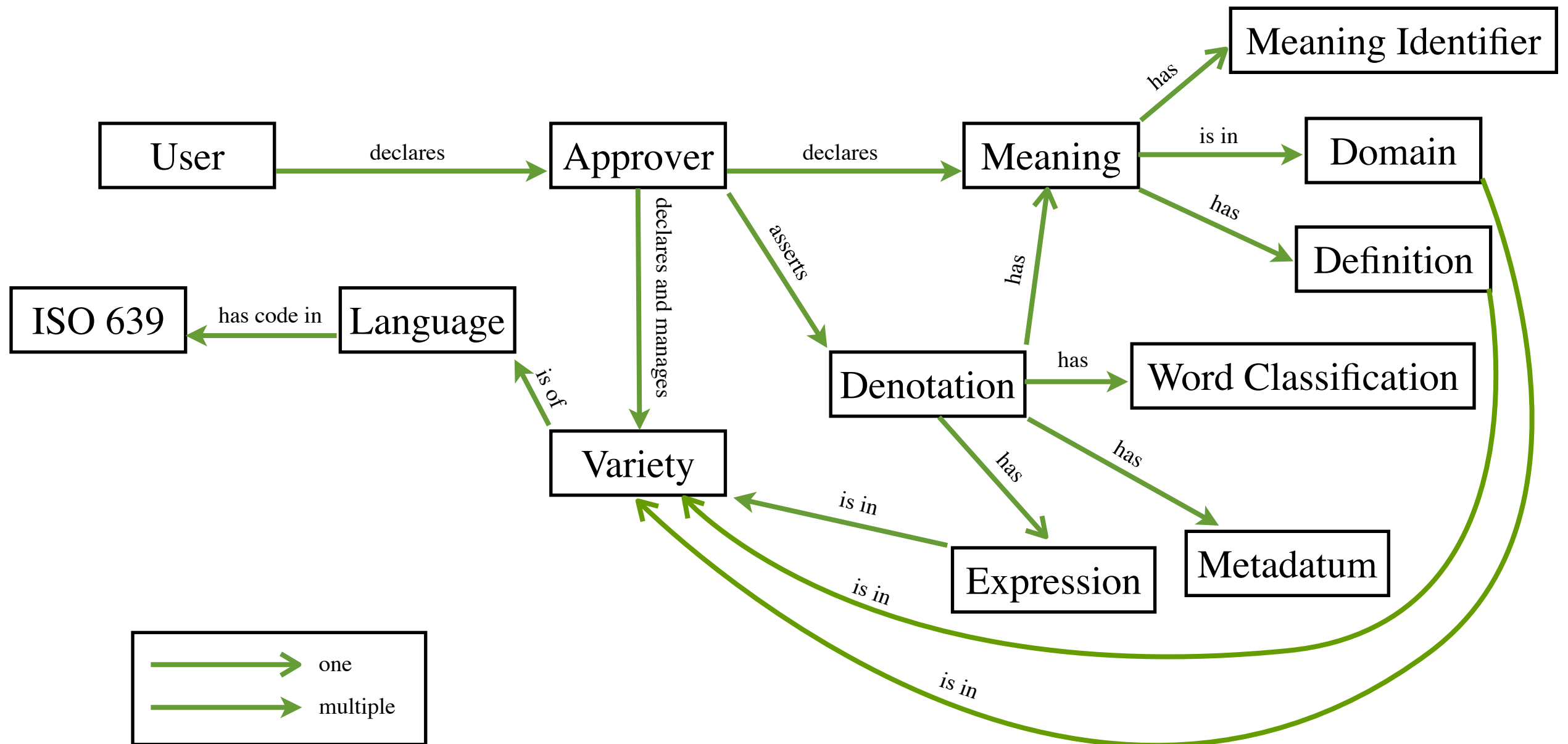
Schweede ·gsw·

Wiktionaries

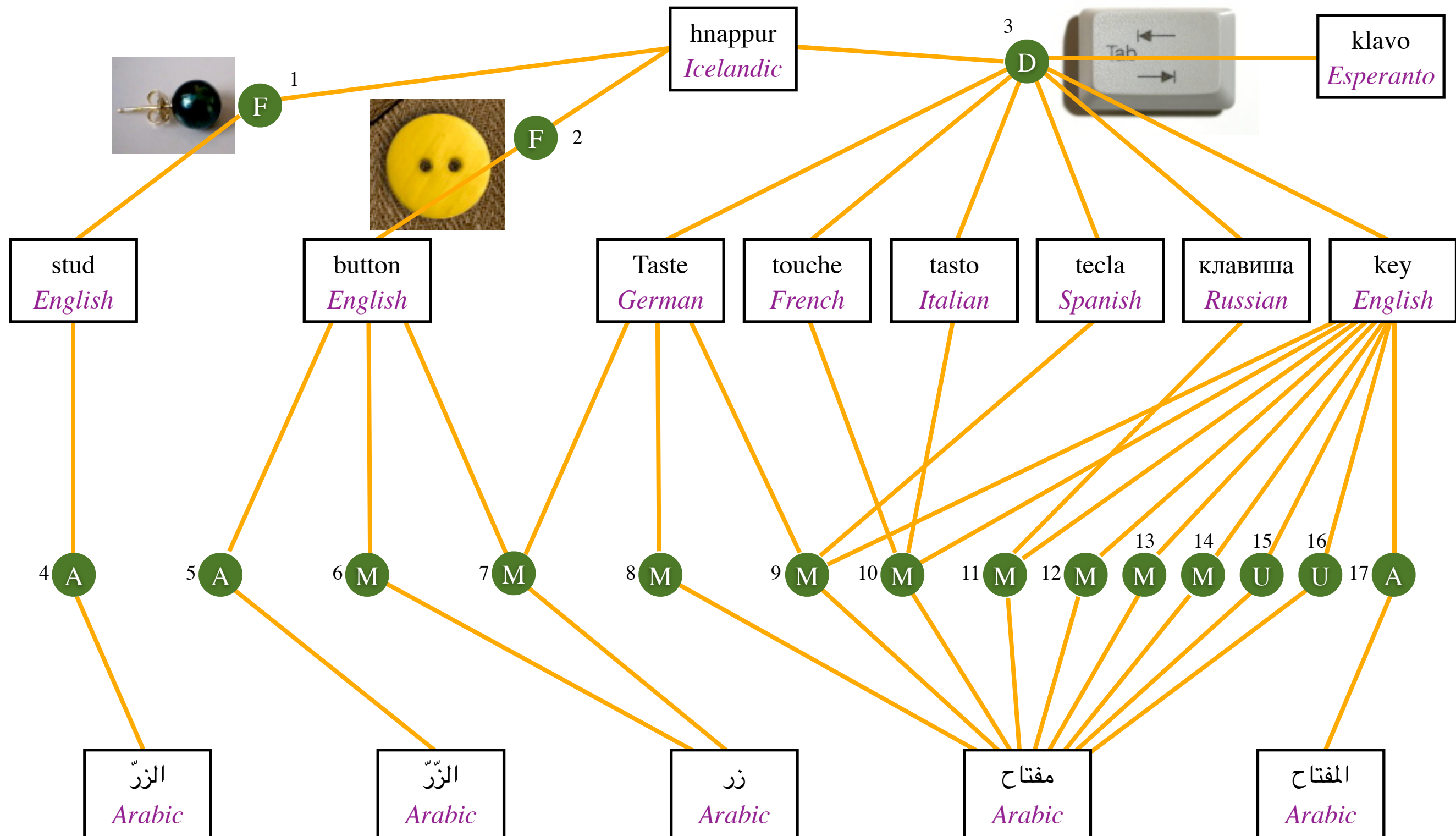


VikiSözlük
Özgür Sözlük

3. Method: Unify Heterogeneous Data



3. Method: Infer New Results



4. Demonstration

Lemmatic UI

PanLex 2.1

translation	expression	meaning
<input type="radio"/> see <input type="radio"/> edit	<input type="radio"/> see <input type="radio"/> change	<input type="radio"/> see
language	source	person
<input type="radio"/> see <input type="radio"/> edit <input type="radio"/> new	<input type="radio"/> see <input type="radio"/> edit <input type="radio"/> new <input type="radio"/> definition <input type="radio"/> file — submit — definition <input type="radio"/> file — submit — example <input type="radio"/> file — submit <input type="radio"/> file — propose <input type="radio"/> file — get	<input type="radio"/> see <input type="radio"/> edit <input type="radio"/> new

English

<http://panlex.org/u>



Utilika Foundation

4. Demonstration

Lematic UI

TümSöz 2.1

çeviri	anlatım	anlam
<input type="checkbox"/> görmek	<input type="checkbox"/> görmek	<input type="checkbox"/> görmek
<input type="checkbox"/> düzenlemek	<input type="checkbox"/> değiştirmek	

dil	kaynak	kişi
<input type="checkbox"/> görmek	<input type="checkbox"/> görmek	<input type="checkbox"/> görmek
<input type="checkbox"/> düzenlemek	<input type="checkbox"/> düzenlemek	<input type="checkbox"/> düzenlemek
<input type="checkbox"/> yeni	<input type="checkbox"/> yeni	<input type="checkbox"/> yeni
	<input type="checkbox"/> tanım	
	<input type="checkbox"/> dosya — göndermek — tanım	
	<input type="checkbox"/> dosya — göndermek — örnek	
	<input type="checkbox"/> dosya — göndermek	
	<input type="checkbox"/> dosya — önermek	
	<input type="checkbox"/> dosya — almak	

Turkish

<http://panlex.org/u>



5. Challenges: Incomplete Data

PanLex 2.1

превод <input type="checkbox"/> поглеждам <input type="checkbox"/> редактирам	източник <input type="checkbox"/> поглеждам <input type="checkbox"/> редактирам <input type="checkbox"/> нов <input type="checkbox"/> <i>dfn</i> <input type="checkbox"/> досие — публикувам — <i>dfn</i> <input type="checkbox"/> досие — публикувам — пример <input type="checkbox"/> досие — публикувам <input type="checkbox"/> досие — предлага <input type="checkbox"/> досие — получа	значение <input type="checkbox"/> поглеждам
език <input type="checkbox"/> поглеждам <input type="checkbox"/> редактирам <input type="checkbox"/> нов	източник <input type="checkbox"/> поглеждам <input type="checkbox"/> редактирам <input type="checkbox"/> нов <input type="checkbox"/> <i>dfn</i> <input type="checkbox"/> досие — публикувам — <i>dfn</i> <input type="checkbox"/> досие — публикувам — пример <input type="checkbox"/> досие — публикувам <input type="checkbox"/> досие — предлага <input type="checkbox"/> досие — получа	човек <input type="checkbox"/> поглеждам <input type="checkbox"/> редактирам <input type="checkbox"/> нов

No translation found

Bulgarian



Utilika Foundation

5. Challenges: Mistranslation

PanLex 2.1

troidigezh	tro-lavar	ster
<input type="checkbox"/> gwelet	<input type="checkbox"/> gwelet	<input type="checkbox"/> gwelet
<input type="checkbox"/> skrid	<input type="checkbox"/> trubardiñ	

teod	erienenn eien	den
<input type="checkbox"/> gwelet	<input type="checkbox"/> gwelet	<input type="checkbox"/> gwelet
<input type="checkbox"/> skrid	<input type="checkbox"/> skrid	<input type="checkbox"/> skrid
<input type="checkbox"/> nevez	<input type="checkbox"/> nevez	<input type="checkbox"/> nevez
	<input type="checkbox"/> termenadur	
	<input type="checkbox"/> fichennaoueg — bazhyevañ — termenadur	
	<input type="checkbox"/> fichennaoueg — bazhyevañ — skouer	
	<input type="checkbox"/> fichennaoueg — bazhyevañ	
	<input type="checkbox"/> fichennaoueg — kinnig	
	<input type="checkbox"/> fichennaoueg — degemer	

Means “betray”
instead of “change”

Breton



Utilika Foundation

5. Challenges: Inconsistency

PanLex 2.1

tarjima	epr	mazmun
<input type="checkbox"/> кўрмоқ <input type="checkbox"/> таҳрир қилмоқ	<input type="checkbox"/> кўрмоқ <input type="checkbox"/> ўзгартирмоқ	<input type="checkbox"/> кўрмоқ
til	payvandlash	shaxs
<input type="checkbox"/> кўрмоқ <input type="checkbox"/> таҳрир қилмоқ <input type="checkbox"/> yangi	<input type="checkbox"/> кўрмоқ <input type="checkbox"/> таҳрир қилмоқ <input type="checkbox"/> yangi <input type="checkbox"/> <i>dfn</i> <input type="checkbox"/> shaxsiy hujjatlar — jo'natmoq — <i>dfn</i> <input type="checkbox"/> shaxsiy hujjatlar — jo'natmoq — o'xshash <input type="checkbox"/> shaxsiy hujjatlar — jo'natmoq <input type="checkbox"/> shaxsiy hujjatlar — буюрмоқ <input type="checkbox"/> shaxsiy hujjatlar — олмоқ	<input type="checkbox"/> кўрмоқ <input type="checkbox"/> таҳрир қилмоқ <input type="checkbox"/> yangi

Scripts vary

Uzbek



5. Challenges: Intelligibility

PanLex 2.1

language

number	10
code	ain
variety	0
meaning — expression — many	yes
expression — meaning — many	yes
language — variety — name	アイヌ イタク
source	mul:pool
expression — count	549

☐ expression — text — character — count — see

Intent: “Let me see the counts of the characters in the texts of the expressions (in Ainu)”: Clear?

6. Solutions: Patch Holes

PanLex 2.1

from — expression

language		expression	
art-001	PanLex	60684	mod

променям
↓

into — expression

language		expression		choose
bul-000	български	62984	променям	<input type="radio"/>

☐ expression — new

Add individual
translations



6. Solutions: Patch Holes

PanLex 2.1

превод	<i>epi</i>	значение
<input type="checkbox"/> поглеждам	<input type="checkbox"/> поглеждам	<input type="checkbox"/> поглеждам
<input type="checkbox"/> редактирам	<input type="checkbox"/> променям	

език	ИЗТОЧНИК	ЧОВЕК
<input type="checkbox"/> поглеждам	<input type="checkbox"/> поглеждам	<input type="checkbox"/> поглеждам
<input type="checkbox"/> редактирам	<input type="checkbox"/> редактирам	<input type="checkbox"/> редактирам
<input type="checkbox"/> нов	<input type="checkbox"/> нов	<input type="checkbox"/> нов
	<input type="checkbox"/> <i>dfn</i>	
	<input type="checkbox"/> досие — публикувам — <i>dfn</i>	
	<input type="checkbox"/> досие — публикувам — пример	
	<input type="checkbox"/> досие — публикувам	
	<input type="checkbox"/> досие — предлага	
	<input type="checkbox"/> досие — получа	

New translation used



6. Solutions: Add Massive Content

Indicator	Now	Goal
Resources	600	10000
Language varieties	1300	7000
Expressions	12,000,000	350,000,000
Expression-meaning pairs	27,000,000	1,000,000,000
Translation/synonym pairs	102,000,000	1,000,000,000

Content Sources

Machine-readable lexical resources

Printed lexical resources

Human computation (games, volunteering, piece work)

Structured and unstructured text



6. Solutions: Add Power

To lemmatic communication:

User training

Interfaces to computational grammars and text analyzers

More intelligent inference

Interactive inference

Beyond lemmatic communication:

Morphological data

Syntactic data

Language-independent relational representations



7. Discussion

Is panlingual localization realistic?

Who would be good collaborators?

What applications could benefit?



Principal Collaborators

University of Washington Turing Center:

Oren Etzioni

Mausam

Stephen Soderland

Marcus Sammer

Kobi Reiter

Katherine Everitt

Utilika Foundation:

Susan Colowick

PostgreSQL Experts, Inc.:

Quinn Weaver



Utilika Foundation

More Information and Publications

Utilika Foundation:

<http://utilika.org>

PanLex project:

<http://panlex.org>

PanLex UI prototype:

<http://panlex.org/u>

Turing Center:

<http://www.turing.washington.edu>

