# PanLex
## A Panlingual Lexicon

## Jonathan Pool & Susan Colowick
### Utilika Foundation

## DELPH-IN Summit
## 26 June 2011

# Outline

- Purpose
- Related work
- Construction
- Size
- Applications
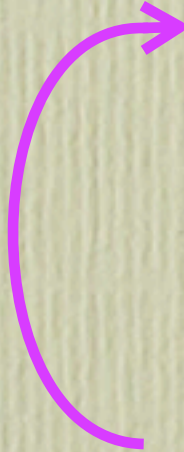- Current work
- PanLex and grammar engineering
- Team

# Purpose

PanLex aims to become a **panlingual lexical translation resource.**

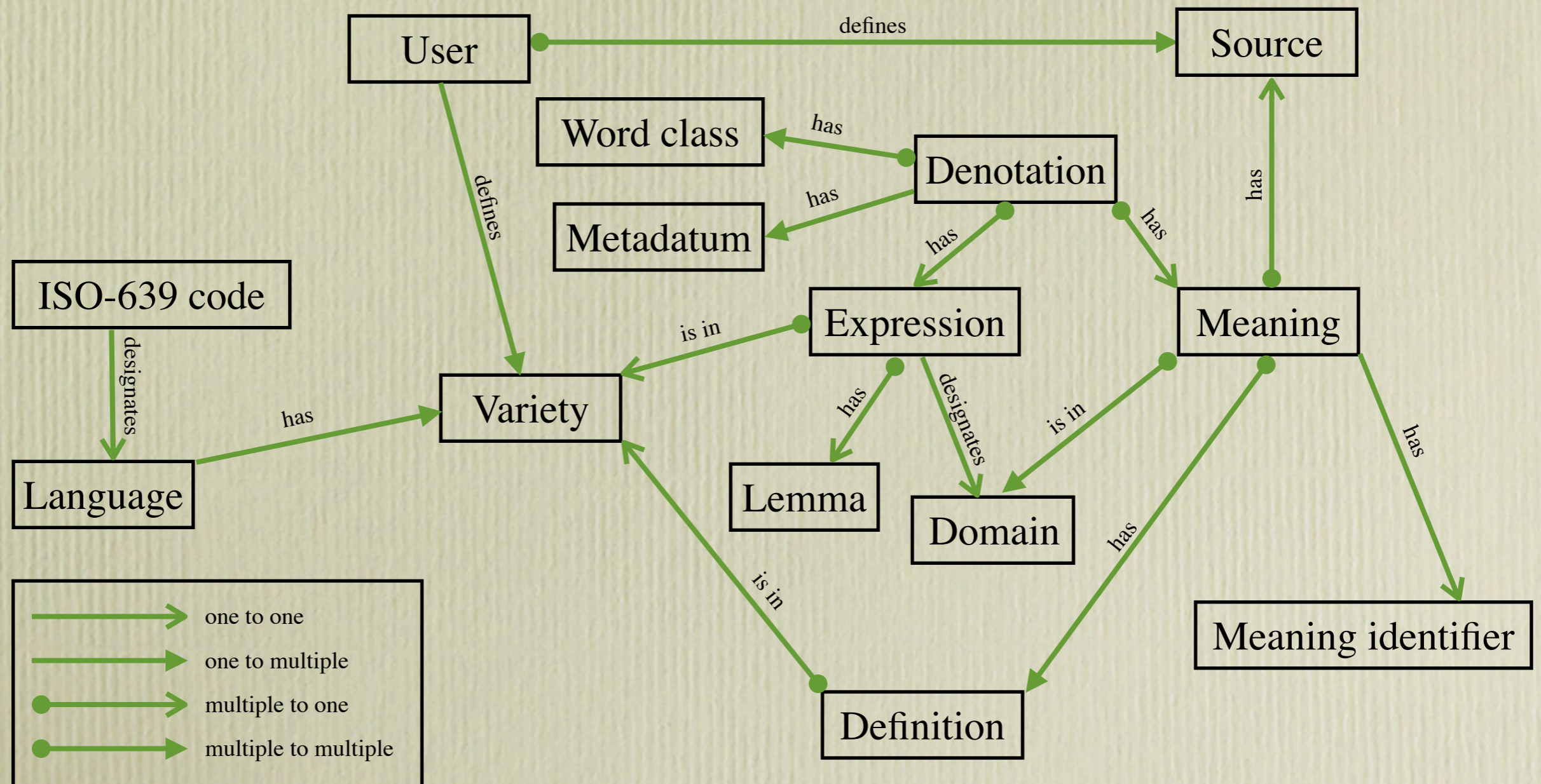| | | | | | |
|---|---|---|---|---|---|
| eng-000 | English | five | fay-000 | Fars | pänj |
| epo-000 | Esperanto | kvin | fer-000 | Feroge | wi |
| erg-000 | Sye | sukrim | ffm-000 | Maasina | joyi |
| erh-000 | Eruwa | íì-sòrī` | fie-000 | Fyer | háwá |
| erk-000 | Fate | lim | fij-000 | vosa Vakaviti | e lima na |
| ero-000 | Horpa | gwai | fin-000 | suomi | viisi |
| ers-000 | Ersu | ŋuàr | fin-000 | suomi | viisi |
| ers-001 | Thochu | wa-re | fip-000 | Fipa | visaano |
| ers-002 | Lyusu | ŋâ | fli-000 | Fali | kè~ɽ è~w |
| ers-003 | Menia | nga | fmp-000 | Fe'fe' | tiǀì |
| ers-004 | Muli | ngo | fng-000 | Fanagalo | fayif |
| erw-000 | Erokwanas | rim | fni-000 | Fanya | luñe |
| ese-000 | Ese Ejja | me-oe-xi | fra-000 | français | cinq |
| ese-001 | Huarayo | iamatamata | fry-000 | Frysk | fiif |
| esi-000 | Iñupiat | tallimat | gil-000 | taetae ni Kiribati | nimaua |
| esk-000 | Iñupiatun | tallimat | gla-000 | Gàidhlig na h-Alba | còig |
| esq-000 | Huelel | pemaxala | gle-000 | Gaeilge | cúig |
| ess-000 | Chaplino | taɬimat | glv-000 | chengey Vannin | queig |
| est-000 | Estonian | viis | gug-000 | avañe'ẽ | po |
| esu-000 | Central Yupik | taɬiman | heb-000 | עברית | חָמֵשׁ |
| etr-000 | Edolo | bi | heb-000 | עברית | חֲמִשָּׁה |
| ett-000 | mechl Rasnal | mach | hin-000 | हिंदी | पाँच |
| etu-000 | Ejagham | é-rôn | | | |
| etx-000 | Aten | wiǀ | hrv-000 | hrvatski | pet |
| etx-001 | Niten | wéé | hun-000 | magyar | öt |
| eus-000 | euskara | bost | hye-000 | արեւմտահայերէն | հինգ |
| eus-000 | euskara | bost | ido-000 | Ido | kin |
| eus-001 | Aitzineuskara | *bortz | ina-000 | interlingua | cinque |
| eus-001 | Aitzineuskara | *bortze | ind-000 | bahasa Indonesia | lima |
| eve-000 | эвэды торэн | tunŋᵉn | isl-000 | íslenska | fimm |
| evh-000 | Uvbie | i-siorī | ita-000 | italiano | cinque |
| evn-000 | орочон турэн | tunŋa | jbo-000 | la lojban. | mu |
| evn-001 | Solon | tongnga | jpn-000 | 日本語 | 五 |
| ewe-000 | Ɛʋɛgbɛ | ató~ | jpn-000 | 日本語 | 五つ |
| ewo-000 | Ewondo | tán | | | |
| eya-000 | Eyak | tcó~i | kal-000 | kalaallisut | tallimat |
| eyo-000 | Keiyo | mú:t | kat-000 | ქართული | ხუთი |
| faf-000 | Fagani | rima | kmr-000 | Kurmancî | pênc |
| fag-000 | Finungwa | yale yale kobok | kor-000 | 한국어 | 다섯 |
| fan-000 | Pahouin | tan | kor-000 | 한국어 | 오 |
| fao-000 | føroyskt | fimm | kpv-000 | коми кыв | вит |

# Related work

- Dicts.info
- FreeDict
- Freelang
- Global WordNet Association
- Langtolang
- Logos Foundation Dictionary
- OmegaWiki

- OneLook
- Open Dictionary
- The Rosetta Project
- Slovnik
- La Vortaro
- Webster's Online Dictionary
- Wiktionary
- Тезаурус

# Construction

- Step 0. Design the schema.

- Step 1. Acquire resources.

- Step 2. Obtain and normalize facts from resources.
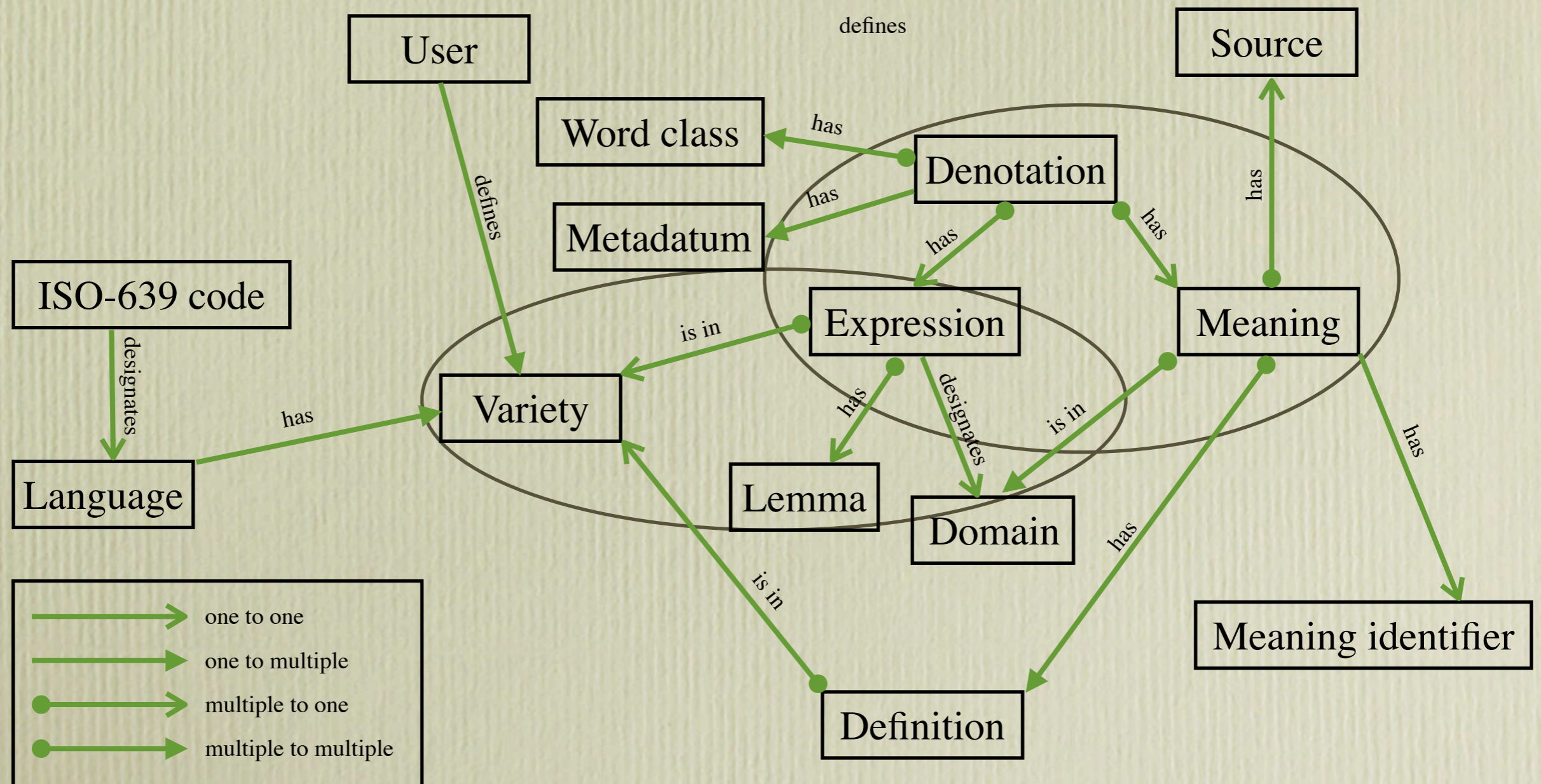
- Step 3. Add those facts to PanLex.
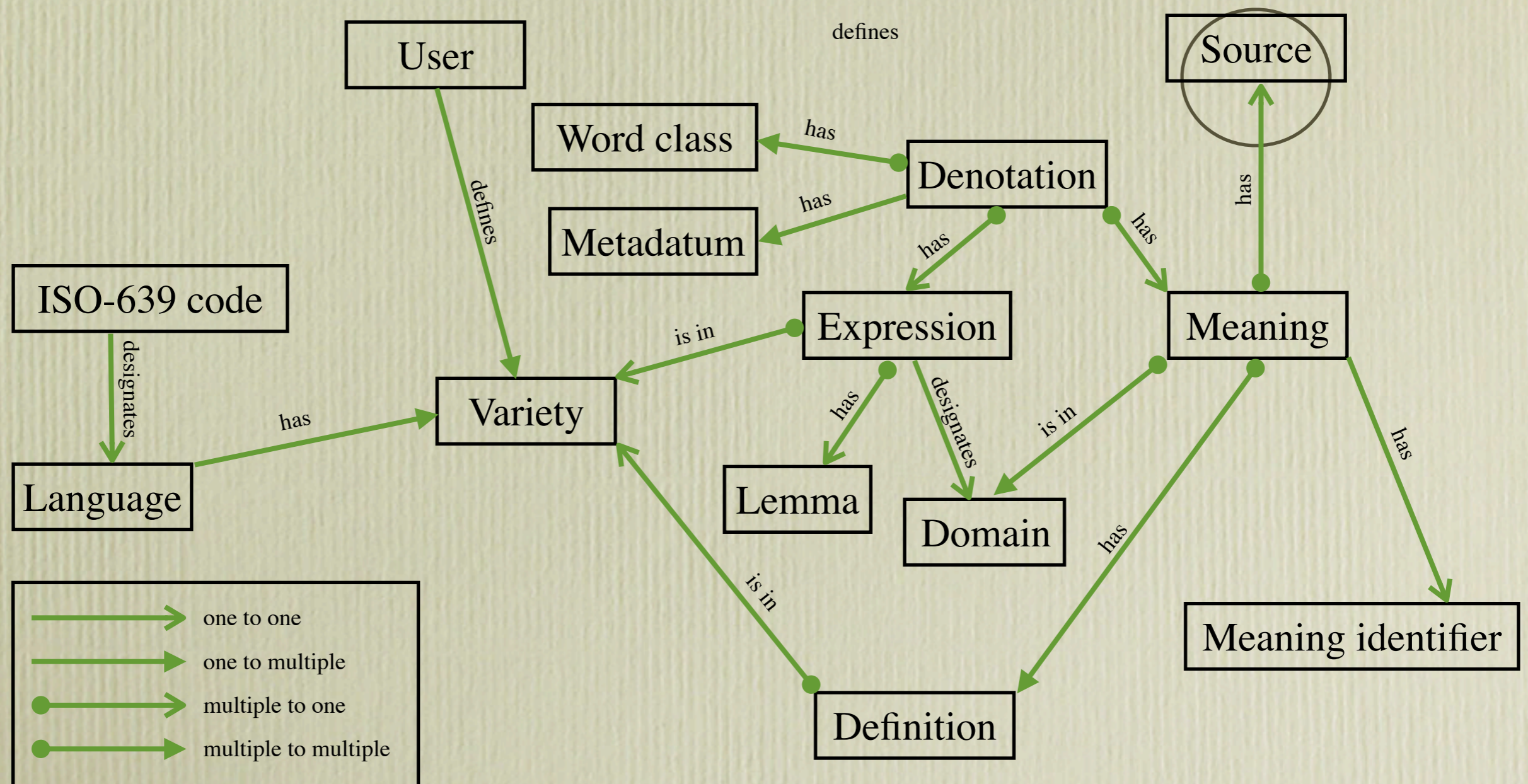
# Construction

Step 0. Design the schema.

# Construction

## Step 0. Design the schema.

# Construction
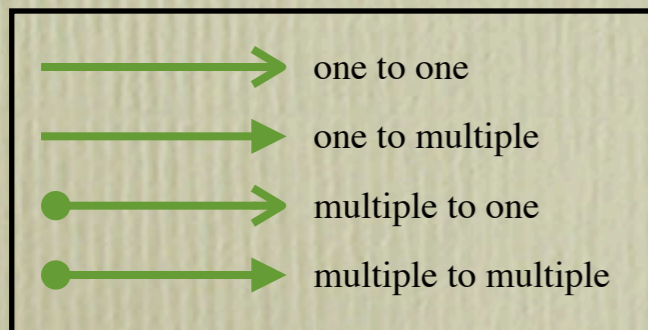
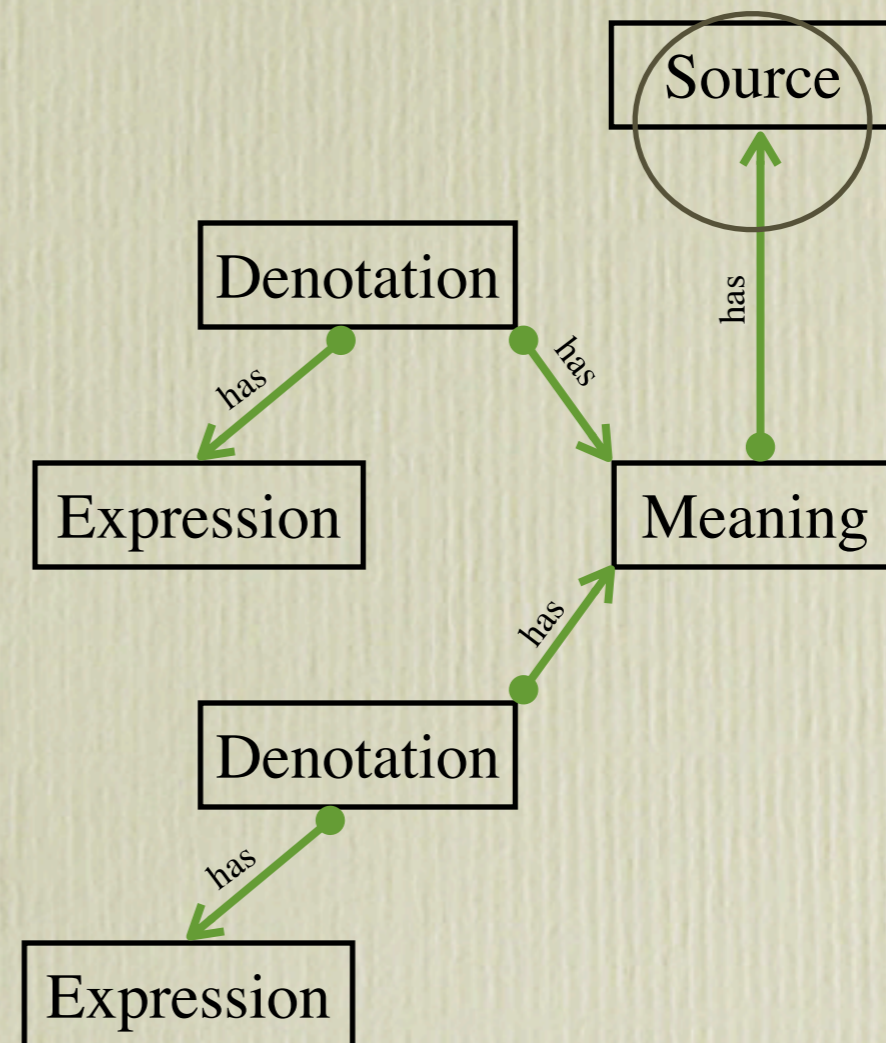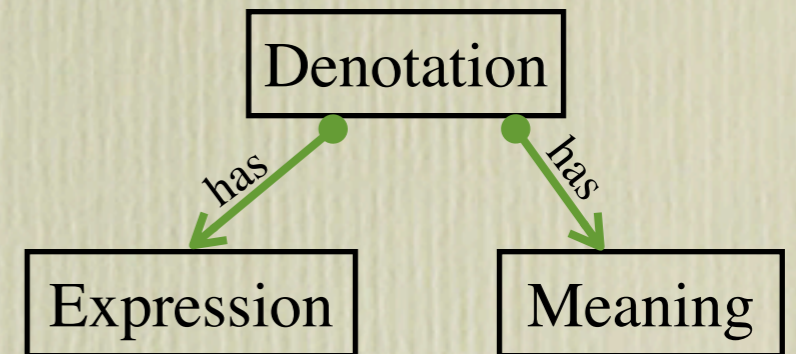## Step 0. Design the schema.

# Construction

Step 0. Design the schema.

Translation

Source

Denotation

has

has

has

Expression    Meaning

has

Denotation

has

Expression

→ one to one

→ one to multiple

●→ multiple to one

●→ multiple to multiple

# Construction

## Step 0. Design the schema.

Denotation
has          has
Expression          Meaning

```
                 Table "public.dn"
 Column | Type    | Modifiers | Storage | Description
--------+---------+-----------+---------+-------------
 dn     | integer | not null  | plain   | ID
 mn     | integer | not null  | plain   | meaning
 ex     | integer | not null  | plain   | expression
Indexes:
    "dn_pkey" PRIMARY KEY, btree (dn)
    "dn_mn_ex_key" UNIQUE, btree (mn, ex) CLUSTER
    "dn_ex_idx" btree (ex)
    "dn_mn_idx" btree (mn)
Foreign-key constraints:
    "dn_ex_fkey" FOREIGN KEY (ex) REFERENCES ex(ex)
    "dn_mn_fkey" FOREIGN KEY (mn) REFERENCES mn(mn)
Referenced by:
    TABLE "md" CONSTRAINT "md_dn_fkey" FOREIGN KEY (dn) REFERENCES dn(dn)
    TABLE "pl0" CONSTRAINT "pl0_mn_fkey" FOREIGN KEY (mn, ex) REFERENCES dn(mn, ex)
        ON UPDATE CASCADE ON DELETE CASCADE
    TABLE "pl1" CONSTRAINT "pl1_mn_fkey1" FOREIGN KEY (mn, ex) REFERENCES dn(mn, ex)
        ON UPDATE CASCADE ON DELETE CASCADE
    TABLE "wc" CONSTRAINT "wc_dn_fkey" FOREIGN KEY (dn) REFERENCES dn(dn)
Triggers:
    dnexap AFTER INSERT OR DELETE OR UPDATE ON dn FOR EACH ROW EXECUTE PROCEDURE exap()
```

# Construction

Step 1. Acquire resources.

- Monolingual dictionaries
- Bilingual dictionaries
- Multilingual dictionaries
- Wiktionaries
- Glossaries
- Standards
- Terminologies
- Wordnets
- Thesauri
- Vocabulary databases
- Locale databases

VikiSözlük
Özgür Sözlük

**Arrest: ∩ᒧᔓᐅᓂᖅ: Tigujauniq: 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: ∆ᑭ∩ᑦ∩ᓂᖅ: Ikitittiniq: Crime d'incendie**

The crime of deliberately setting fire to property. 1996,

*γλώσσα για ειδικούς σκοπούς*
**MT** (70.20)
**Da:** fagsprog
**De:** Fachsprache
**En:** language for special purposes
**Es:** lenguaje especializado
**Fi:** kieli tiettyihin tarkoituksiin
**Fr:** langage spécialisé
**He:** שפה למטרות מיוחדות
**Hu:** szaknyelv
**It:** lingua speciale
**Nl:** vaktaal
**Sv:** fackspråk
**BT** γλώσσες

| **SE (English: Sweden)** | |
|---|---|
| **An tSualainn** | ·ga· |
| **isveç** | ·az· |
| **İsveç** | ·tr· |
| **Iswidhan** | ·so· |
| **Rootsi** | ·et· |
| **Ruotsi** | ·fi· |
| **Ruotta** | *se·* |
| **Schweden** | ·de· |
| **Schweede** | ·gsw· |

# Construction

Step 2. Obtain and normalize
facts from resources.



```
:
2
rus-000
epo-000

ex
гавиал
ex
gavialo

ex
гагара
ex
grebo

ex
гагара
ex
kolimbo

ex
гадюка
ex
vipero
ex
vipuro
```

# Construction

Step 2. Obtain and normalize facts from resources.

**adial** *n.* 1) obsidian; a volcanic glass-like substance. *glas.*
2) fish species; a name designating a bottlefish, unicorn surgeonfish, or some other species of unicornfish (also known as **gelenga**). *botolpis. Naso annulatus; Naso tuberosus; Naso unicornis; Naso brevirostris; Naso brachycentron.*
**adial itna** *n.* fish species; Bannerfish. *Heniochus diphreutes.*
**adiuol** *n.* a species of vine.

```
:
1
bch-000

ex
adial
wc
noun
ex
eng-000
obsidian
df
eng-000
a volcanic glass-like substance
ex
tpi-000
glas
```

# Construction

Step 2. Obtain and normalize facts from resources.

| English | Tagalog | Ilocano | S. Kalinga |
|---------|---------|---------|------------|
| sky | lángit | lángit | langit |
| cloud | alapáap | úlep | lifuu |
| rainbow | bahaghári | bullaláyaw | afungar |
| star | bituín | bitwén | fituwon |

```
.
1
eng-000

sky
tgl-000
lángit
ilo-000
lángit
ksc-000
langit

cloud
tgl-000
alapáap
ilo-000
úlep
ksc-000
lifuu
```

# Construction

Step 2. Obtain and normalize
facts from resources.

- Character recognition
- Character recoding
- Compositional normalization
- Punctuation standardization
- Lemma standardization

- Word-class standardization
- Entry-structure classification
- Language-variety identification
- Duplicate removal

Jonathan Pool, "Processing LEGO Data for PanLex",
http://www.panlex.org/dev/panlex-lego-prep.html.

Timothy Baldwin, Jonathan Pool, and Susan M. Colowick,
"PanLex and LEXTRACT", Coling 2010.
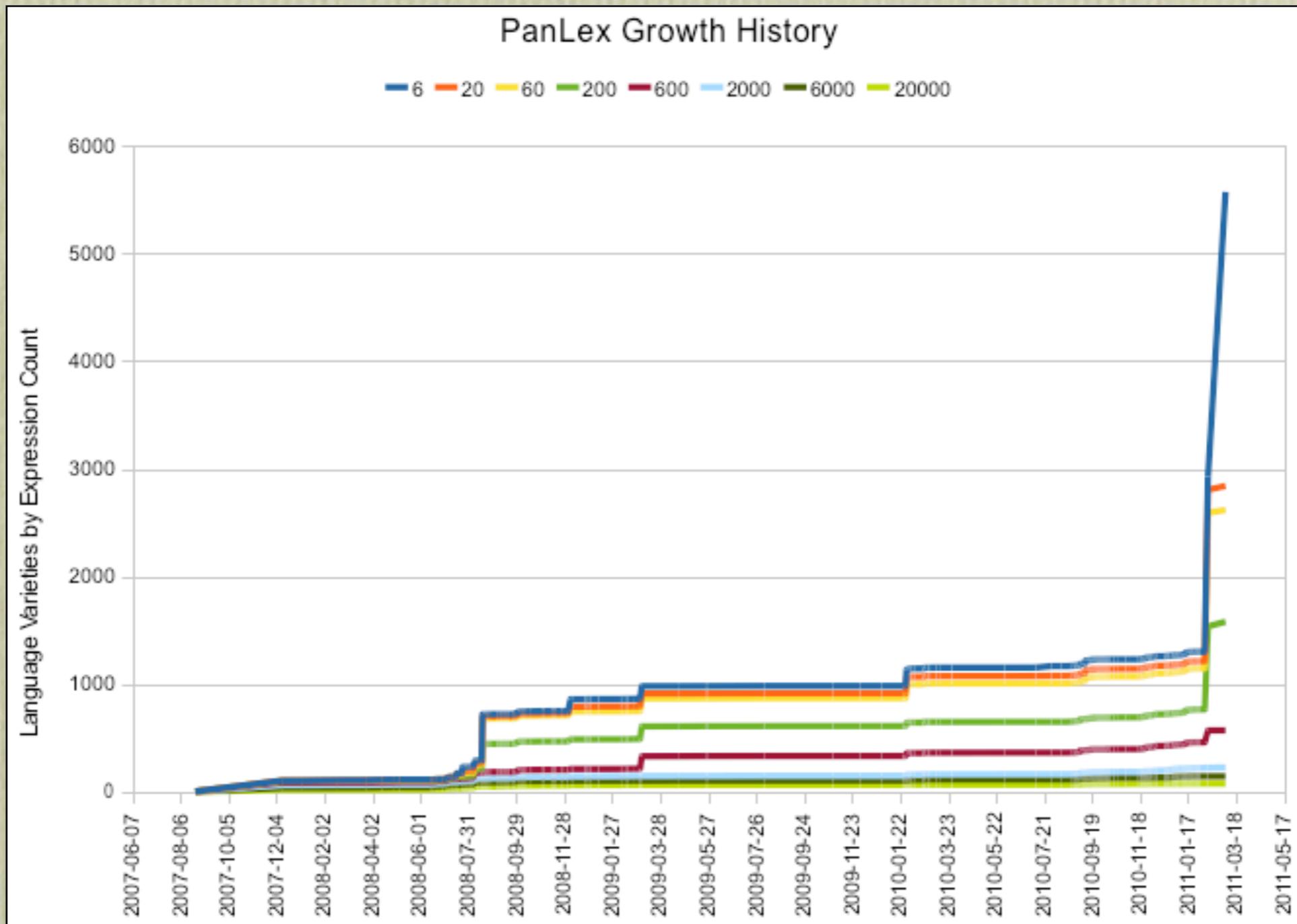
# Construction

Step 3. Add those facts to PanLex.

# Size

Current metrics:

- 17,621,880 expressions

- 6,662 language varieties

- 446,906,568 pairwise translations

- 1,121 sources processed

- 2,202 sources awaiting processing

# Size

# Applications

## UIs for PanLex

**PanLex 2.8**

translation

through

| language | | expression | | see |
|---|---|---|---|---|
| bul-000 | български | 62252 | лист | ⬭ |
| deu-000 | Deutsch | 91652 | Blatt | ⬭ |
| eng-000 | English | 463039 | leaf | ⬭ |
| epo-000 | Esperanto | 633660 | folio | ⬭ |
| ita-000 | italiano | 69185 | foglia | ⬭ |
| ita-000 | italiano | 69189 | foglio | ⬭ |
| ita-000 | italiano | 72631 | pagina | ⬭ |
| pol-000 | polski | 732673 | liść | ⬭ |
| rus-000 | русский | 756911 | лист | ⬭ |
| slv-000 | slovenščina | 1329204 | list | ⬭ |

**from**

| language | | expression | |
|---|---|---|---|
| lit-000 | lietuvių | 1307318 | lapas |

**into**

| language | | expression | |
|---|---|---|---|
| tgl-000 | Tagalog | 1051479 | dahon |

# Applications

## UIs for PanLex

# Applications

## UIs for PanLex



**TeraDict**

syntax

**TeraDict can translate this into:**

| | |
|---|---|
| als-000 | toskërishte |
| arb-000 | العربية |
| asm-000 | অসমীয়া ভাষা |
| bel-000 | беларуская |
| ben-000 | বাংলা |
| bul-000 | български |
| cat-000 | català |
| ces-000 | čeština |
| cmn-000 | 简体字 |
| cmn-001 | 繁體中文 |
| cym-000 | Cymraeg |
| dan-000 | dansk |
| deu-000 | Deutsch |

**TümSöz**

söz

| fra-000 | français |
|---|---|

**Çeviriler:**

| |
|---|
| c'est dit |
| discours |
| mot |
| mots |
| parlé |
| parole |
| propos |
| verbe |

Yeni söz veya deyim:

# Applications

## Search



Janara Christensen, Mausam, and Oren Etzioni, "A Rose is a Roos is a Ruusu: Querying Translatins for Web Image Search", ACL-IJCNLP 2009.

# Applications

## Translation



Stephen Soderland, Christopher Lim, Mausam, Bo Qin, Oren Etzioni, and Jonathan Pool, "Lemmatic Machine Translation", Proceedings of Machine Translation Summit XII, 2009.

# Applications

## Lemmatic Communication



Katherine Everitt, Christopher Lim, Oren Etzioni, Jonathan Pool, Susan Colowick, Stephen Soderland, "Evaluating Lemmatic Communication", *trans-kom*, 3, 2010, 70–84.

# Applications

## Polling?

| 48. | Do you consume fast food? | ⦿ Never / don't know |
| --- | --- | --- |
| | | ○ Once per week or less |
| | | ○ 2-3 times per week |
| | | ○ 4-7 times per week |
| | | ○ A lot / more than once daily |
| 49. | Aside from fast food, how often do you consume deep-fried foods? | ⦿ Never / don't know |
| | | ○ Once per week or less |

## Messaging?

jonathan | PanMail   PanImages   Help   Settings   Sign Out

**Panlingual Mail** Beta

**Email in Any Language!**

[Search Mail]

**Compose Message**

Inbox

Sent Mail

Drafts

Trash

[ Delete ] [ Mark Unread ]

Source Language: [            ]
Target Language: [            ]
To : [            ]
From : [ jonathan ]
Subject : [            ]
Message : No Translations...
[            ]

[ Delete ] [ Mark Unread ]

Panlingual Mail is provided by the Turing Center

# Applications

Games?

Social networking?

# Applications

Subtitles?



cinstit preferinţă nu ingrijorare



รวบรวม ปกติ ต่างๆ คลางแคลง มนษย

# Current work

- Process 2,000+ acquired sources

- Acquire new digital sources

- Digitize and process printed sources

- Cultivate partnerships and volunteers

http://utilika.org/info/intern2011.html

# Current work

## Internship projects, 2011–2012

- Lexical acquisition:
  - ✓ Optical character recognition
  - ✓ Lexicographic parsing
  - ✓ Game development

- Applications:
  - ✓ Image search
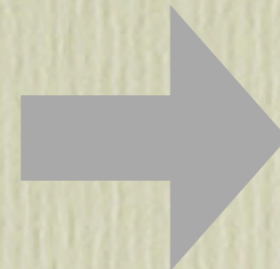  - ✓ Mobile app
  - ✓ Social app

- Infrastructure:
  - ✓ Translation inference
  - ✓ Grammar engineering
  - ✓ Graph visualization

http://utilika.org/info/intern2011.html

# PanLex and Grammar Engineering

## Lexical acquisition **from** grammars

```
noun3_det=opt
  noun3_stem1_orth=mazzita
  noun3_stem1_pred=_blutwurst_n_rel
  noun3_stem2_orth=ittra
  noun3_stem2_pred=_letter_n_rel
  noun3_stem3_orth=universita
  noun3_stem3_pred=_university_n_rel
```

```
:
2
mlt-000
eng-000

…

ex
ittra
wc
noun
md
det
opt
ex
letter

…
```

# PanLex and Grammar Engineering

## Lexemes and predicates **for** grammars

**Noun Types**

▼ lemmatic (noun1)

ⓧ **Noun type 1**:

Type name: `lemmatic`

Features:

( Add a Feature )

For nouns of this type, a determiner is ○ obligatory ⦿ optional ○ impossible

Stems:

ⓧ | Spelling: `кальме`    Predicate: `_glue_n_rel`

ⓧ | Spelling: `леф`    Predicate: `_lion_n_rel`

ⓧ | Spelling: `ошеряй`    Predicate: `_citizen_n_rel`

ⓧ | Spelling: `тевонь пула`    Predicate: `_chain reaction_n_rel`

ⓧ | Spelling: `шнамань кельгома`    Predicate: `_vanity_n_rel`

# PanLex and Grammar Engineering

## Instantly generate 6,000 minimal grammars

version=6

section=language
language=Bandjalang

section=word-order
word-order=free
has-dets=no
has-aux=no

section=number

section=person
person=none

---

section=gender

section=other-features
section=case
case-marking=none

section=sentential-negation

section=coordination

section=matrix-yes-no

Cf. David Wax,
Matrix ODIN Mash-up

---

section=lexicon
noun1_orth=bayaɲ
noun1_pred=_bayaɲ_n_rel
noun1_det=imp
noun2_orth=ɟibali
noun2_pred=_ɟibali_n_rel
noun2_det=imp
noun3_orth=wumar wumar
noun3_pred=_wumar wumar_n_rel
noun3_det=imp
…

section=test-sentences

# PanLex and Grammar Engineering

## Then incrementally differentiate and relate the grammars

version=6

section=language
language=Bandjalang

section=word-order
word-order=sov
has-dets=no
has-aux=no

section=number

section=person
person=none

section=gender

section=other-features
section=case
case-marking=none

section=sentential-negation

section=coordination

section=matrix-yes-no

WALS          PanLex

section=lexicon
advb1_orth=bayaŋ
advb1_pred=_today_n_rel

noun2_orth=ɹibali
noun2_pred=_ɹibali_n_rel
noun2_det=opt
noun3_orth=wumar wumar
noun3_pred=_wumar wumar_n_rel
noun3_det=imp
…

section=test-sentences

Documentation, crowdsourcing

# Team

## Turing Center, University of Washington
### http://www.turing.washington.edu

- Oren Etzioni
- Katherine Everett
- Christopher Lim
- Mausam
- Kobi Reiter

- Marcus Sammer
- Michael Schmitz
- Michael Skinner
- Stephen Soderland

**Turing Center**

Investigating problems at the crossroads of natural language processing, data mining, Web search, and the Semantic Web.

# Team

## Utilika Foundation
http://utilika.org

- Susan Colowick
- Jonathan Pool
- Miranda Taylor

# Team

## Collaborators

- Academic
  - LEGO Project (Timothy Usher, Jeff Good, Helen Aristar-Dry)
  - Timothy Baldwin, University of Melbourne
  - Larry Hyman, UC Berkeley
  - Interns from Univ. of Washington, Univ. of Illinois, Univ. of Maryland, UC Davis, Univ. of Melbourne, CMU, and Ohio State
- Nonprofit
  - Long Now Foundation / Rosetta Project (Laura Welcher)
- For-profit
  - CJK Dictionary Institute (Jack Halpern)
  - Digital Sonata (Vadim Berman)
  - PostgreSQL Experts (Josh Berkus, Quinn Weaver)

# Comments/Questions

Info:

http://panlex.org