# Matthew Lee: Who am I?
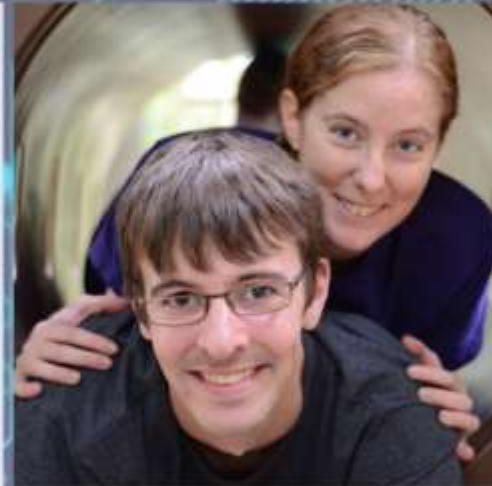
- Language Technology Specialist
  Member of SIL Cameroon

- Pursuing Master's Degree in Descriptive Linguistics

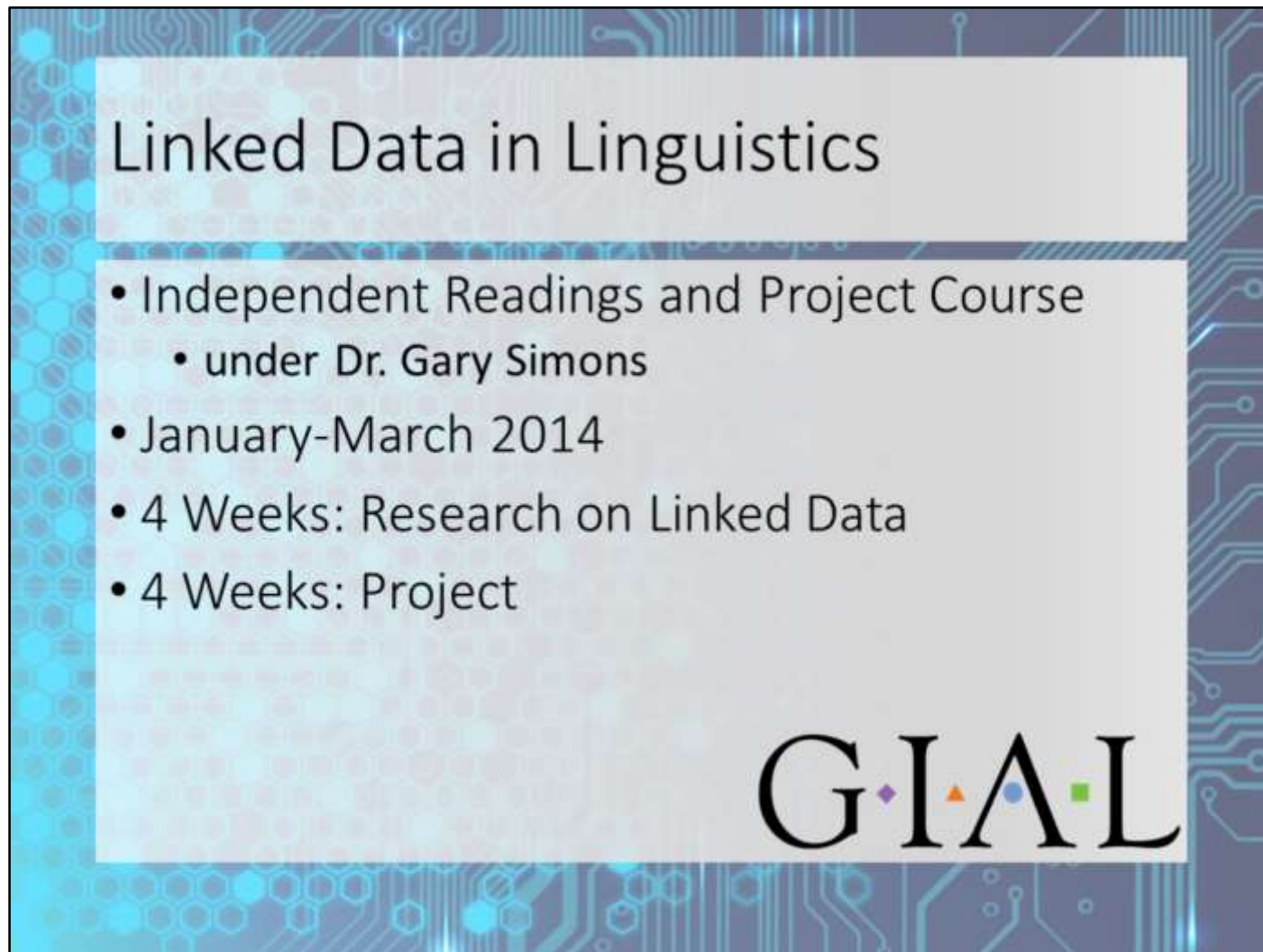- BS: Integrated Science & Technology
- BA: Philosophy

Matthew & Teresa Lee

Independent Study

Textbook: Linked Data in Linguistics

Chiarcos, Christian, Sebastian Nordhoff & Sebastian Hellman (eds.). 2012. *Linked Data in Linguistics*. Springer.

Full of articles detailing ways to make linguistic data smarter and to share, connect, and analyze it in new ways.

## Course Goals:

- Learn about Linked Data and how it can be used in linguistics.
- Learn the technology necessary (XSLT) to transform XML Data.
- Transform XML lexical databases into an interoperable Linked Data format (RDF).
- Demonstrate interesting cross-linguistic searching across those databases.

# Interoperability

*noun*

1977 : ability of a system to work with or use the parts or equipment of another system interoperable \-ˈä-p(ə-)rə-bəl\ *adjective*

Merriam-Webster, I. (2003). Merriam-Websters collegiate dictionary. Springfield, MA: Merriam-Webster, Inc.
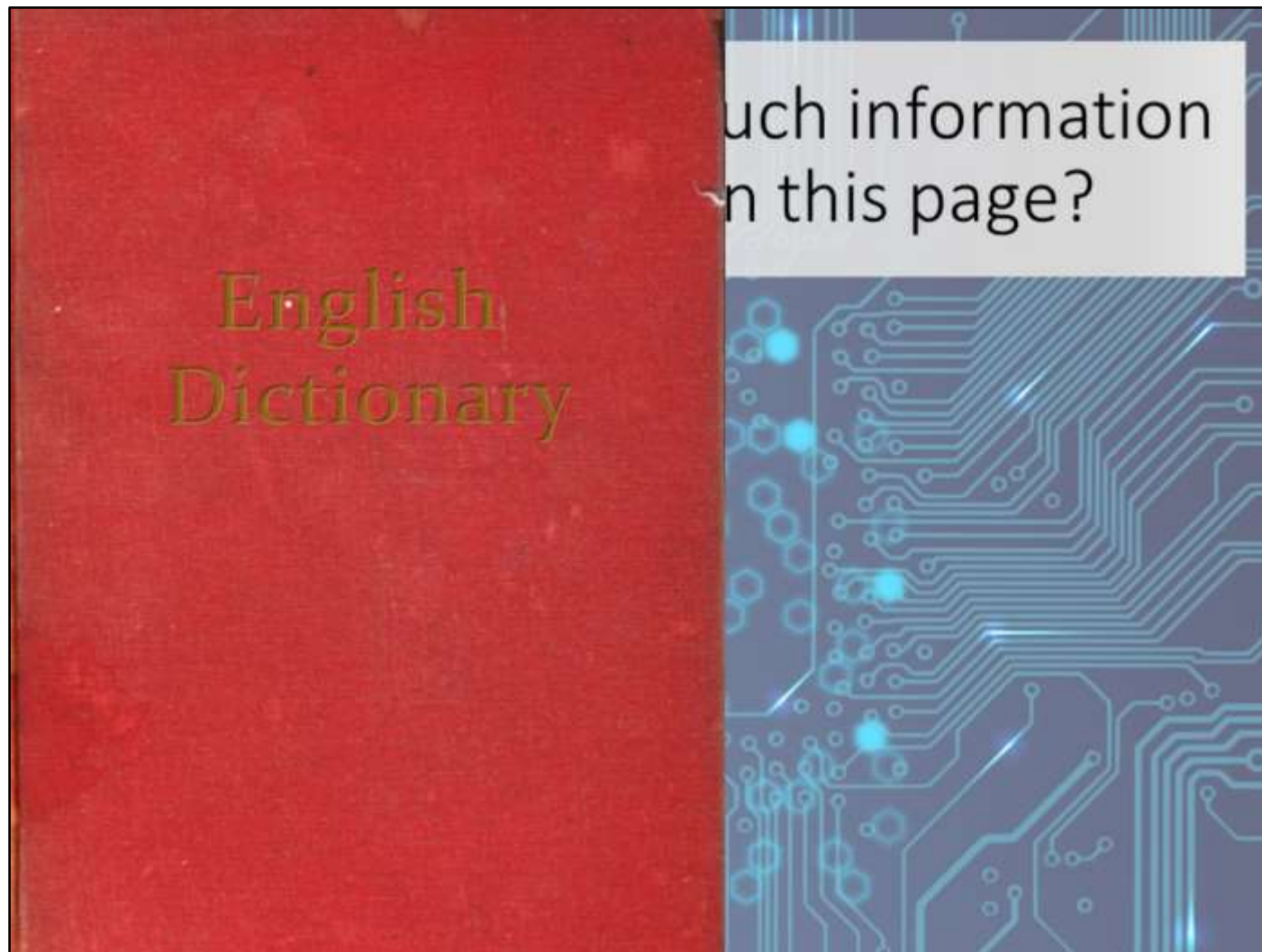
# Previous Work

- Interoperable Lexicons are not a new idea.
- Gary proposed this in 2005 (Simons), but at that point there were not many robust XML Lexicons in existence.
- Thanks to the advances and adoption of XML-based Lexical tools like **Fieldworks Language Explorer**, **WeSay**, and **Lexus** we can now move forward into that reality.

Simons, Gary F. 2005. Beyond the Brink: Realizing Interoperation through an RDF Database. Linguistic Ontologies and Data Categories for Linguistic Resources. Cambridge, MA. http://emeld.org/workshop/2005/papers/simons-paper.pdf (8 December, 2013).

# Previous Work

- Helen Arister-Dry and a team at the Max Planck Institute has attempted around 2011 to create a bridge between FLEx Lexicons and LMF lexicons (a competing standard from MPI).

- Some lexicons and wordlists were converted and uploaded to http://lego.linguistlist.org/ but there is not much interoperability other than a simple search.

- It seems that there was never an attempt to do cross-linguistic comparison.

Arister-Dry, Helen. LIFTing LEGO with RELISH: Lexicon Interchange FormaT in Use. Paper presented at the Institute for Language Information and Technology. http://www.mpi.nl/departments/other-research/research-projects/language-archiving-technology/events/relish-workshop/program/LexiconStandardsandtheLEGOProject.pptx.

# Dictionaries as Datapoints: Matthew Lee



How much information is on this page?

How much information is on this page?

- Dictionary
- Language
- Words
- Phrases
- Pronunciations
- Meaning
- Relationships
- Categories

Your brain knows how to interpret this data, it sees: [Click]

# Dictionaries as Datapoints: Matthew Lee



Let's move this information into a Desktop Publishing program….for example Microsoft Word.
Just so you know I'm not picking on word…this could be [Click]

Microsoft Publisher

Dictionaries as Datapoints: Matthew Lee



Libreoffice or OpenOffice

Adobe InDesign

Dictionaries as Datapoints: Matthew Lee



Or Apple's Pages

# Dictionaries as Datapoints: Matthew Lee



Back to Word…these programs are designed for creating formatted documents, they lack critical features for organizing data. What does Word "understand" about this document?

# Dictionaries as Datapoints: Matthew Lee



Though it can reliably print you document, Word knows almost nothing about the content of a document. How many facts can the computer access? Can you ask it for a list of lexemes? Can Word help you to create a reversal for your dictionary?

This is the reason why...if you want to create a document in another form, for printing on A4 or a different audience...you now have two documents and you must be careful to make consistent changes to both.

Formatted Dictionary

Computer sees:
- Text
- Formatting
- (That's all folks!)

Brain sees:
- Dictionary
- Language
- Words
- Phrases
- Meaning
- Relationships
- Categories
- Connections

Lexical Databases: A Step Forward

Dictionaries as Datapoints: Matthew Lee

```
\lx Sirius
\ph sir'i-us
\ps n
\sn 1
\de the dog-star.

\lx sirloin
\ph sẽr'loin,
\ps n
\sn 1
\de the loin, or upper part of the loin, of beef

\lx sirocco
\ph si-rok'ō
\ps n
\pl siroccos (sirok'ōz)
\sn 1
\de a hot, relaxing wind, from the Libyan deserts

\lx sirrah
\ph sir'a
\ps n
\sn 1
```

I have converted the previous document into another format. How many of you recognize this lexicon format?

```
\lx Sirius
\ph sir'i-us
\ps n                Standard Format Ma
\sn 1
\de the dog-star.

\lx sirloin
\ph sẽr'loin,
\ps n
\sn 1
\de the loin, or upper part of the loin, of beef

\lx sirocco
\ph si-rok'ō
\ps n
\pl siroccos (sirok'ōz)
\sn 1
\de a hot, relaxing wind, from the Libyan deserts

\lx sirrah
\ph sir'a
\ps n
\sn 1
```

This is the Standard Format, the file used by Shoebox and Toolbox. I count 282 "Facts". Now I can search and filter the data if I put it into Toolbox.

People love Toolbox!  The best part of Toolbox (for working) is that it doesn't constrain you to any consistent structure.  The worst part of toolbox (for publishing) it that doesn't constrain you to any consistent structure. If you want to add a custom column that keeps track of which words you've taught your pet parrot…it will let you put that anywhere. As a result, Toolbox Lexicons tend to mix and match languages, abbreviations, order, and notations over the years. They also tend to have references to non-existent entries. There is no safety net.

IN SFM, the content is still a mystery to the computer, but at least the program "knows" that each line represents a specific type of data. With the right configuration, Toolbox is designed to keep track of structured data in many languages. Though there are published standards, the structure here is really only understood by convention (Lexeme, part of speech, etc…), and you often need the original author to explain what some obscure items mean.

Publishing a Toolbox Dictionary

Image Source: http://www.tipsquirrel.com/photoshop-3d-whats-an-extrusion/

At best, publishing a toolbox Lexicon is like this:

and at it's worst…

SFM Dictionary (Toolbox/Shoebox)

Computer sees:
- Text
- ~~Formatting~~
- Data Categories
- Language

Brain sees:
- Dictionary
- Language
- Words
- Phrases
- Meaning
- Relationships
- Categories
- Connections

XML: Structured Data

Dictionaries as Datapoints: Matthew Lee



Structured Data

A Lexicon has many Entries. Typically, each entry contains one or more senses and a pronunciation. Each sense can have zero or more example sentences. Each sense can have zero or more Subsenses.  Even if you have a custom field…[click]…it has a place in the structure.

What does XML data look like?

Dictionaries as Datapoints: Matthew Lee

```xml
<entry dateCreated="2014-05-02T15:49:00Z" dateModified="2014-05-02T15:49:00Z"
id="sisal-grass_0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d"
guid="0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d">
  <lexical-unit>
    <form lang="en">
      <text>sisal-grass</text>
    </form>
  </lexical-unit>
  <trait name="morph-type" value="stem" />
  <pronunciation>
    <form lang="en">
      <text>sis'al-gras</text>
    </form>
  </pronunciation>
  <sense id="382513ad-f2e9-44fd-b9f0-9e3f185be051">
    <grammatical-info value="Noun"></grammatical-info>
    <definition>
      <form lang="en">
        <text>the prepared fiber of the American aloe, used for cordage</text>
      </form>
    </definition>
  </sense>
</entry>
```

LIFT XML looks like this…which may seem overwhelming at first, but it is self describing.

Dictionaries as Datapoints: Matthew Lee

```xml
<entry dateCreated="2014-05-02T15:49:00Z" dateModified="2014-05-02T15:49:00Z"
id="sisal-grass_0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d"
guid="0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d">
  <lexical-unit>
    <form lang="en">
      <text>sisal-grass</text>
    </form>
  </lexical-unit>
  <trait name="morph-type" value="stem" />
  <pronunciation>
    <form lang="en">
      <text>sis'al-gras</text>
    </form>
  </pronunciation>
  <sense id="382513ad-f2e9-44fd-b9f0-9e3f185be051">
    <grammatical-info value="Noun"></grammatical-info>
    <definition>
      <form lang="en">
        <text>the prepared fiber of the American aloe, used for cordage</text>
      </form>
    </definition>
  </sense>
</entry>
```

Entry

LIFT XML

35

Dictionaries as Datapoints: Matthew Lee



```xml
<entry dateCreated="2014-05-02T15:49:00Z" dateModified="2014-05-02T15:49:00Z"
id="sisal-grass_0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d"
guid="0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d">
  <lexical-unit>
    <form lang="en">
      <text>sisal-grass</text>
    </form>
  </lexical-unit>
  <trait name="morph-type" value="stem" />
  <pronunciation>
    <form lang="en">
      <text>sis'al-gras</text>
    </form>
  </pronunciation>
  <sense id="382513ad-f2e9-44fd-b9f0-9e3f185be051">
    <grammatical-info value="Noun"></grammatical-info>
    <definition>
      <form lang="en">
        <text>the prepared fiber of the American aloe, used for cordage</text>
      </form>
    </definition>
  </sense>
</entry>
```
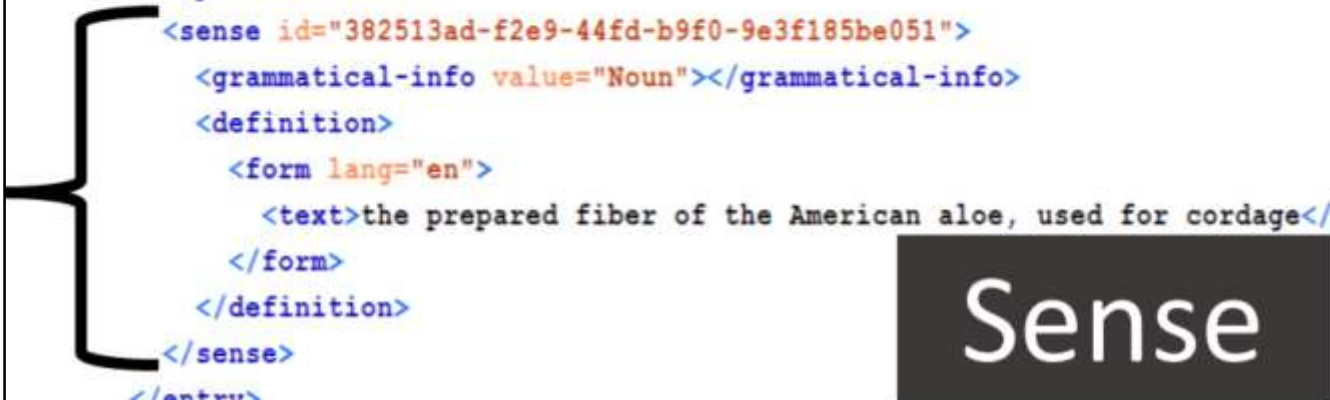
Sense

LIFT XML

36

```xml
<entry dateCreated="2014-05-02T15:49:00Z" dateModified="2014-05-02T15:49:00Z"
id="sisal-grass_0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d"
guid="0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d">
  <lexical-unit>
    <form lang="en">
      <text>sisal-grass</text>
    </form>
  </lexical-unit>
  <trait name="morph-type" value="stem" />
  <pronunciation>
    <form lang="en">
      <text>sis'al-gras</text>
    </form>
  </pronunciation>
  <sense id="382513ad-f2e9-44fd-b9f0-9e3f185be051">
    <grammatical-info value="Noun"></grammatical-info>
    <definition>
      <form lang="en">
        <text>the prepared fiber of the American aloe, used for cordage</text>
      </form>
    </definition>
  </sense>
</entry>
```
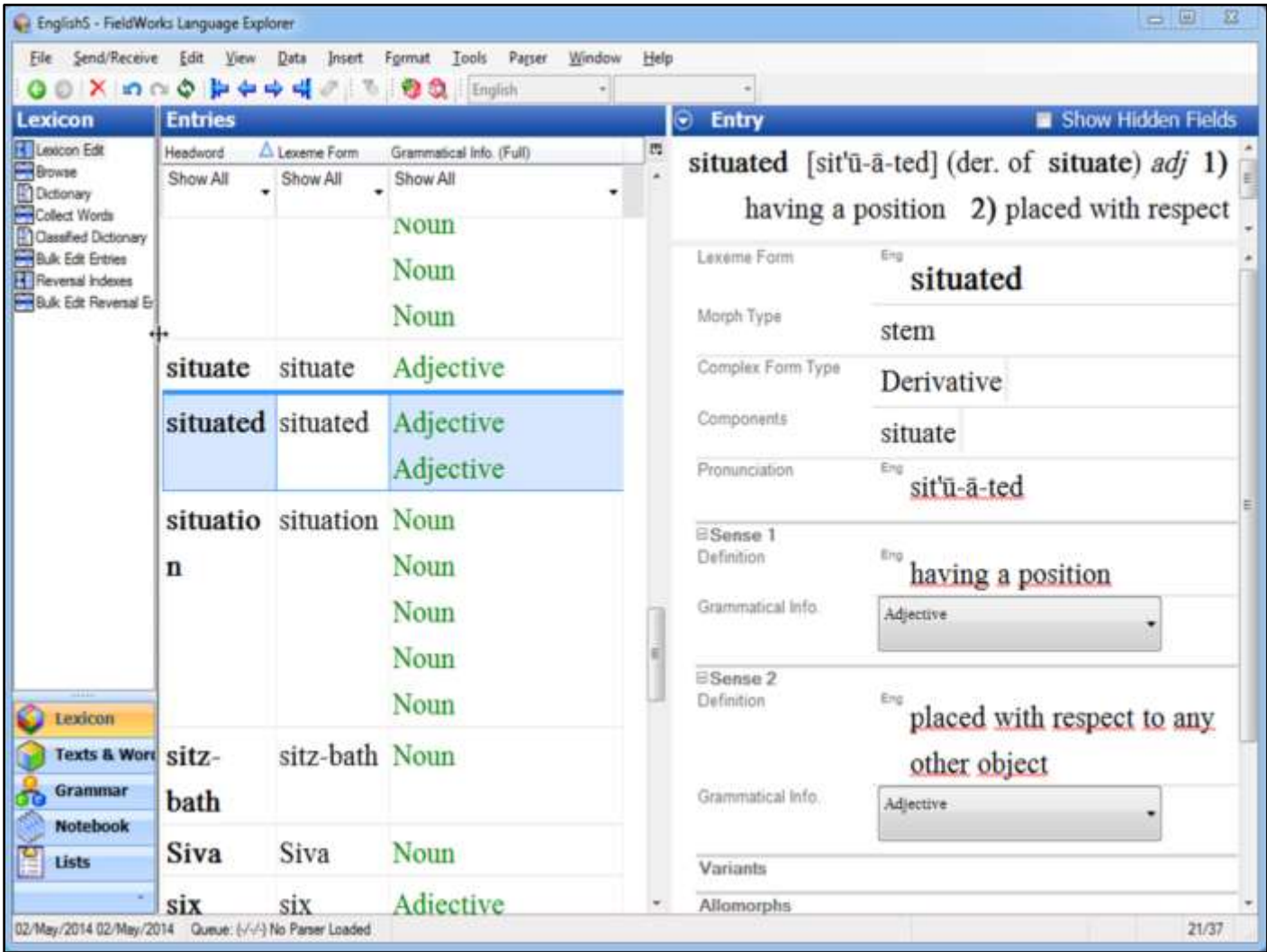
Pronunciation

LIFT XML

**Transformation Technologies**

- XSLT
- CSS
- XSL-FO
- LaTeX

These are all technologies already being harnessed for publishing dictionaries and Scripture.

Because the data is structured, it is relatively easy for the computer to reorganize the data and create a new form.

# Dictionaries as Datapoints: Matthew Lee



This is SIL's Fieldworks Language Explorer. A new generation of linguists has been trained on FLEx, and many are transitioning from Toolbox to FLEx. FLEx uses XML to store it's data and offers fields for most of the categories that you will need, with the option of adding custom ones into the structure. **XML Data can look like this…[click]**

**situate** [sit'ū-āt] *adj* placed. der. **situated**

**situated** [sit'ū-ā-ted] (der. of **situate**) *adj* 1) having a position 2) placed with respect to any other object

**situation** [sit-ū-ā'shun] *n* 1) position 2) locality 3) circumstances 4) office 5) employment

**sitz-bath** [sits'bath] *n* a bath for bathing in a sitting posture

**Siva** [sē'va] *cf:* **Brahma**. *n* a god in the Hindu triad, appearing with Brahma and Vishnu. Siva is the destroying god, and his emblem is a bull.

**six** [siks] *adj* 1) one more than five:\ps n 2) the number greater by one than five 3) the symbol representing 6 units

**sixfold** [siks'fōld] *adj* six times as many or as much

**sixpence** [siks'pens] *n* a small British silver coin, value six pennies, or 12½ cents

**sixpenny** [siks'pen-i] *adj* worth six pence

**sixscore** [siks'skōr] 1) *n* six times twenty 2) *adj* six times twenty

**six-shooter** [siks-shōōt'er] *n* a six-chambered revolver

Or this

40

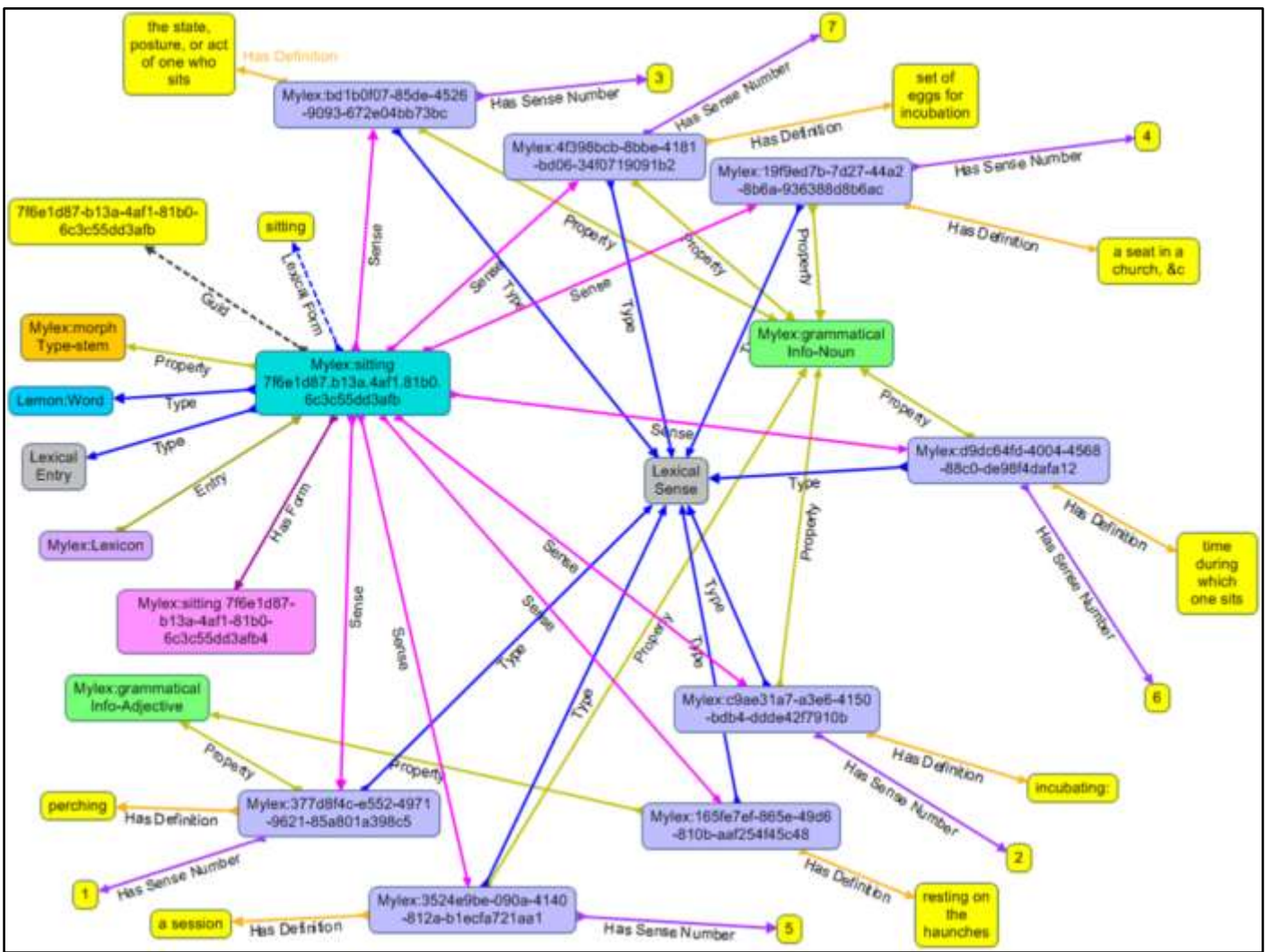Dictionaries as Datapoints: Matthew Lee



Or this (Bonngi Webonary)

Dictionaries as Datapoints: Matthew Lee



Or this (Anki Flashcards)

# Dictionaries as Datapoints: Matthew Lee
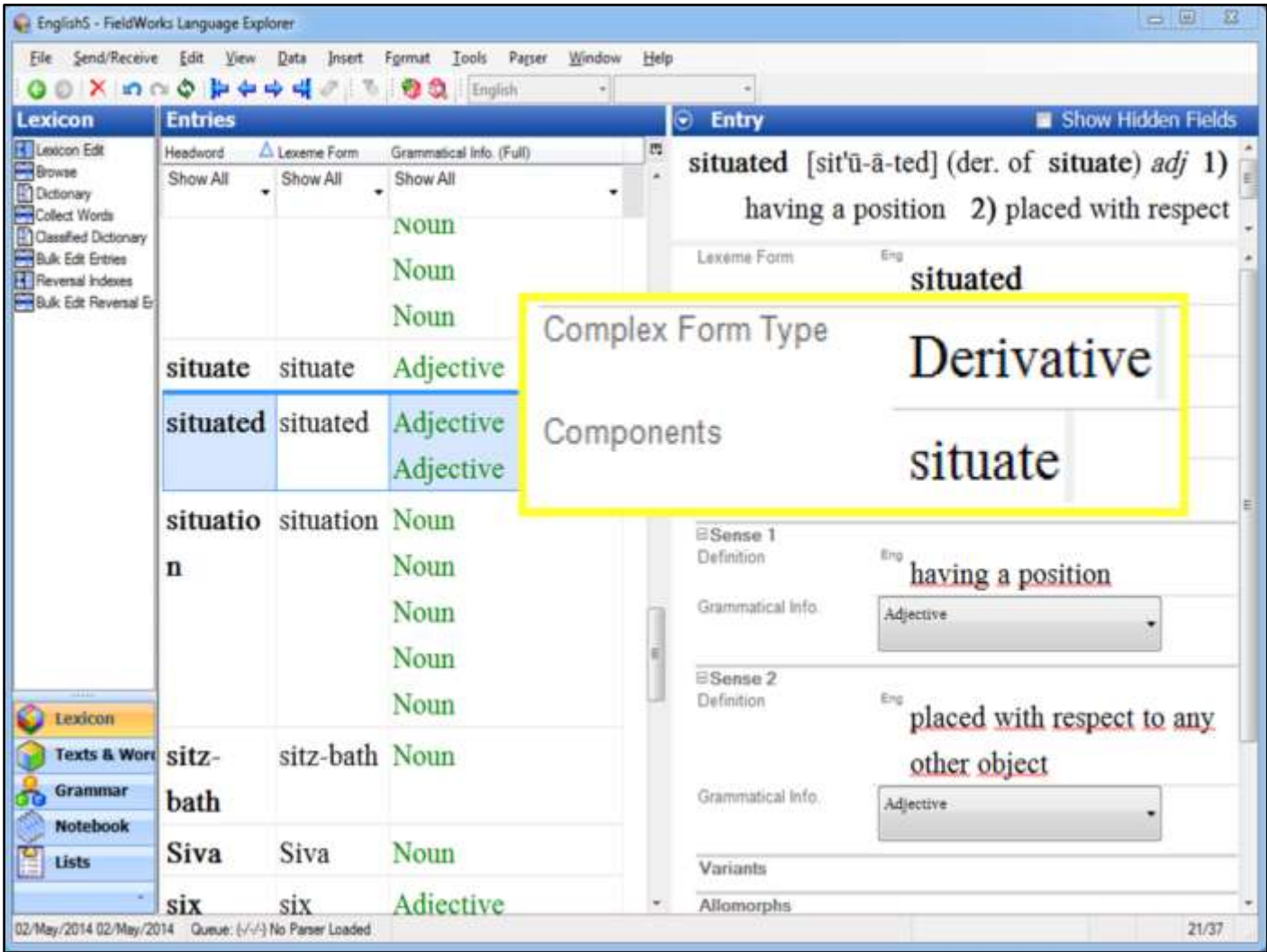


I'll come back to this one..

Or of course it can be printed…

# Dictionaries as Datapoints: Matthew Lee



XML data can also be linked…this shows that the adjective situate is related to "situate".  The tool makes sure that the target entry exists to prevent dead links.
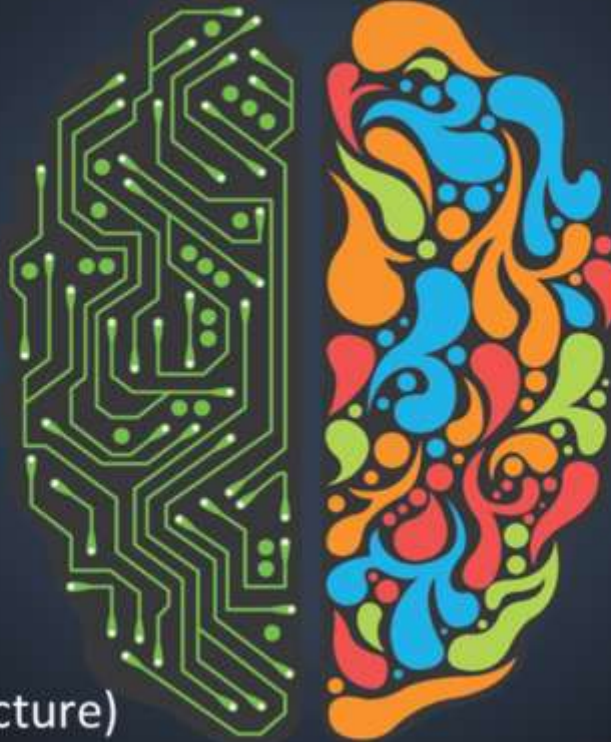
## Publishing is Automagic (Mostly)

- Because of Structure, exports of XML data are designed to be both configurable and automatic, and often they are.
  - Configure once, export as needed.
- If an automatic export doesn't work for your context, you will need to work with a technician to tweak that template, but after that, you can generate the it automagically.

## Standing on the Shoulders of Giants:

- XML is the ideal working format, where the linguist works.

- It contains enough structured information to transform into other interesting formats.

- Anytime I want to create a new presentation form, I can export a new copy from the current lexicon and update my system automagically.

- The linguist can keep working! I'm working on a snapshot.

| Data Format | Searching | Filtering | Formatting | Publishing | Derived Forms | Linking | Checking Consistency |
|---|---|---|---|---|---|---|---|
| Typed or Handwritten | | | | | | | |
| Standard Format | ✔ | ✔ | | | | | |
| Desktop Publishing | ✔ | | ✔ | ✔ | | | |
| XML | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ | ✔ |

This chart shows what these types of data are good at, though it could be argued that each of these things is possible.  If I had room, I'd add "Leveling your Coffee Table" so that the handwritten dictionary could have a point in its favor.

The Technology that I used to link the databases is called Resource Description Format. Here's a video about a little company that is exploring this technology.

Google, Inc. 2012. *Introducing the Knowledge Graph.*
http://www.youtube.com/watch?v=mmQl6VGvX-
c&feature=youtube_gdata_player (3 May, 2014).

All Knowledge can be represented by nodes and relationships.
The more relationships…the more you know about a concept. If you don't believe that everything can be defined by relationships…

gondola
elevator car
cable car
railway car
railroad car
railcar
auto
automobile
motorcar machine
https://www.visualthesaurus.com/

If you don't believe that everything can be defined by relationships… Take a Thesaurus… Each line represents "synonym". Even if you don't know the word you could learn something about it by studying the related words.  The computer can help us understand things by showing those relationships.

Unique Identifiers (URIs)

RDF relies heavily on Unique Resource Identifiers known as URIs. When you link to something, you want to know that you are linking to one specific page and that that page will be there for the forseeable future.

This requires disambiguation.
Many names for one "Thing"

Many things with the same name.

## What could I link to?

There are many organizations on the Internet making comprehensive lists of things.

- SIL manages the ISO 639-3 (Ethnologue Codes) for the world's languages.
- IMDB (and Freebase) have lists of movies, actors and music.
- Lexvo is a Massively Multilingual Lexicon
- Encyclopedia of Life (eol.org) for Plants and Animals

To get around this, you can link to the specific thing on an authoritative source.

What could I link to?

What's the web's largest list of specific things?

I know that Wikipedia is a **Bad Word** in the academic world. You are not saying that you agree with everything on the page, but by linking to Wikipedia you can say: I'm talking about this specific thing that exists in the real world!

Dictionaries as Datapoints: Matthew Lee



Going back to our example, here we see a Wikipedia Disambiguation Page for All the Famous people named Dean Martin.

Dictionaries as Datapoints: Matthew Lee



And here's the one we want.

# Dictionaries as Datapoints: Matthew Lee

| | | |
|---|---|---|
| abstract | dateOfBirth | is artist of |
| activeYearsEndYear | dateOfDeath | is associatedBand of |
| activeYearsStartYear | deathDate | is associatedMusicalArtist of |
| alias | deathPlace | is creator of |
| background | genre | is influencedBy of |
| birthDate | hasPhotoCollection | is musicalArtist of |
| birthPlace | label | is musicalBand of |
| deathDate | name | is presenter of |
| deathPlace | occupation | is starring of |
| genre | placeOfBirth | is wikiPageDisambiguates of |
| occupation | placeOfDeath | is wikiPageRedirects of |
| recordLabel | shortDescription | is artist of |
| thumbnail | yearsActive | is associatedActs of |
| wikiPageExternalLink | description | is creator of |
| wikiPageID | subject | is extra of |
| wikiPageInLinkCount | type | is guests of |
| wikiPageOutLinkCount | comment | is influences of |
| wikiPageRevisionID | label | is mainCharTeam of |
| alias | sameAs | is presenter of |
| alternativeNames | wasDerivedFrom | is recordedBy of |
| background | depiction | is starring of |
| birthDate | givenName | is title of |
| birthName | isPrimaryTopicOf | is writer of |
| birthPlace | name | is sameAs of |
| caption | surname | is primaryTopic of |

Just By linking

Dictionaries as Datapoints: Matthew Lee



"Linking Open Data cloud diagram, by Richard Cyganiak and Anja Jentzsch. http://lod-cloud.net/"

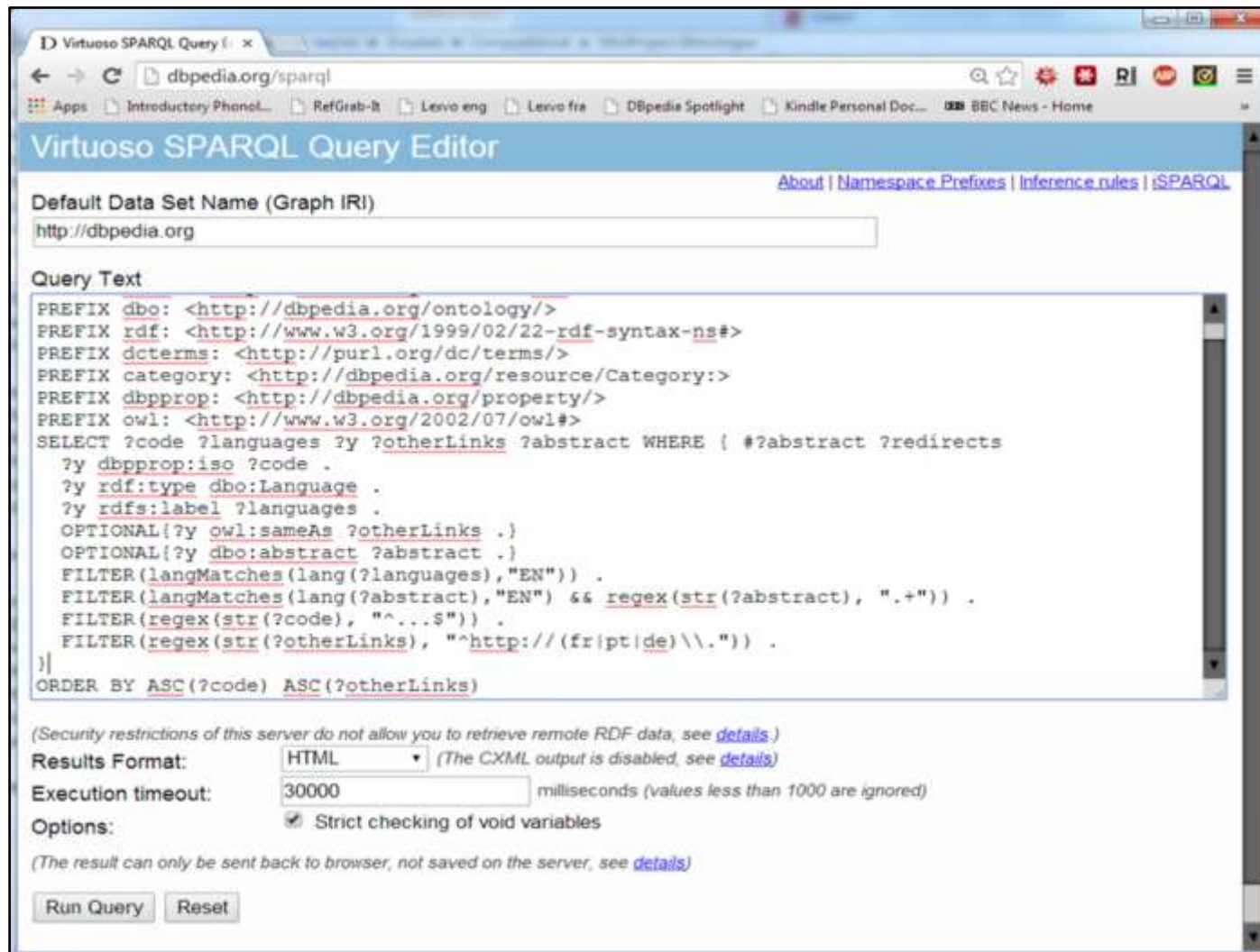Dictionaries as Datapoints: Matthew Lee

# What can you ask of RDF?

- Gary wanted to know which languages had Wikipedia (DBpedia) pages that were more than just stubs.

Query in prose:

Show me the DBPedia data linked to each 3-letter code. Include the Abstract if available and links to other language Wikipedia pages about the same language.

# Dictionaries as Datapoints: Matthew Lee



Query in SPARQL

Dictionaries as Datapoints: Matthew Lee

New Query:

- Query in Prose:

Download all facts and links on DBpedia from English pages about every language, then sort them by language code.

Dictionaries as Datapoints: Matthew Lee

Dictionaries as Datapoints: Matthew Lee

Links for Lexicons

## What should a Dictionary link to?

- ISO 639-3 (Ethnologue) Codes for a language
- GOLD is a structured web of Linguistic Terms and Concepts.
- Ideally, one could link to Wiktionary Entries for specific senses.
- Any classification system, like Semantic Domains, is a viable candidate.

## "Time Zones"

- Each language has specific classes of features that behave in a specific way.
- Linguists love to show off what is DIFFERENT about the language they study.

...so that they can become rich and famous "in the linguistic world".

## "Time Zones"

- Let's say a linguist defines these categories:
  - Recent Past
  - Distant Past
  - Not-so-recent past
  - Yesterday
- What do they have in common?
- They are ALL past tense!

## "Time Zones"

FLEx has pre-populated lists to choose from:
- Morpheme Types
- Parts of Speech
- Inflectional Features
- Relationship Types (Synonym/Antonym)
- Semantic Domains
- Anthropological Categories

I went through these lists in FLEx, pre-linking them to their exact feature (i.e. Recent Past) in GOLD if available, and also to their parent features (i.e. "Past Tense" and "Tense").

Dictionaries as Datapoints: Matthew Lee

"Time Zones"

- If you have defined "Recent Tense" as inflectional features in your language, we can now ask for lists of:
  - Morphemes marking Tense
  - Morphemes marking Past Tense
  - Morphemes marking Recent Past
- This is very similar to the phonological categories in WALS.

Dryer, Matthew S. & Martin Haspelmath (eds.). 2013. *WALS Online*. Leipzig: Max Planck Institute for Evolutionary Anthropology. http://wals.info/ (7 May, 2014).

McCrae, John, Dennis Spohr & Philipp Cimiano. 2011. Linking Lexical Resources and Ontologies on the Semantic Web with lemon. *The Semantic Web: Research and Applications*, 245–259. http://pub.uni-bielefeld.de/luur/download?func=downloadFile&recordOId=2278403&fileOId=2526034.

Lemon is a proposed architecture for Lexicons in RDF and I used their model as a skeleton.

## Future Work: Glosses

- Each entry could be linked to Wiktionary or Wikipedia pages in the analysis language for their referents.

- The linguist would have to do and verify this work, so this is not likely to happen until linguists see the value of it.

- These pages are linked to entries for the same referents in other languages.

There are tools out there to help with making these links…

- OpenRefine. http://openrefine.org
- Dbpedia Spotlight http://dbpedia-spotlight.github.io/demo/

- But none of them are adapted (yet) for this purpose.

# Dictionaries as Datapoints: Matthew Lee



I take a Flex database, like this example database…and Export the parts that are important to me as LIFT XML.

Dictionaries as Datapoints: Matthew Lee

```xml
<entry dateCreated="2014-05-02T15:49:00Z" dateModified="2014-05-02T15:49:00Z"
id="sisal-grass_0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d"
guid="0047e7a0-5dbe-40e8-9fd6-a244f8aceb0d">
  <lexical-unit>
    <form lang="en">
      <text>sisal-grass</text>
    </form>
  </lexical-unit>
  <trait name="morph-type" value="stem" />
  <pronunciation>
    <form lang="en">
      <text>sis'al-gras</text>
    </form>
  </pronunciation>
  <sense id="382513ad-f2e9-44fd-b9f0-9e3f185be051">
    <grammatical-info value="Noun"></grammatical-info>
    <definition>
      <form lang="en">
        <text>the prepared fiber of the American aloe, used for cordage</text>
      </form>
    </definition>
  </sense>
</entry>
```

LIFT XML

84

This is where I spent most of my time on the project, tweaking the transformations to pull out the important structure and data.

## Mapping to URIs Prefix:Reference

- gold:Noun
  - http://purl.org/linguistics/gold/Noun
- lemon:sense
  - http://lemon-model.net/lemon#sense
  - Lemon is a proposed model of lexical information in RDF.
- lexsil:GrammaticalInfo

  Point to the definition of Grammatical info on my server.
- Semantic Domains
  - These are already structured and in-depth.
- mylex:sixteenmo_fd24c341.3793.48ab.b978.9ef84 8fa8f89
  - Point to the specific lexical entry with this name in this lexicon.

The Prefix, or part before the colon, is a shortcut pointing to an authority that catalogs such things, their catalog is called an Ontology.

Dictionaries as Datapoints: Matthew Lee

| Subject | Predicate | Object |
|---|---|---|
| mylex:d85533bb-0326-42c2-b2d8-626aa8c95a74 | lexsil:hasDefinition | "sexto decimo" |
| mylex:grammaticalInfo-Noun | rdf:type | lexsil:GrammaticalInfo |
| mylex:d85533bb-0326-42c2-b2d8-626aa8c95a74 | lemon:property | mylex:grammaticalInfo-Noun |
| mylex:d85533bb-0326-42c2-b2d8-626aa8c95a74 | rdf:type | lemon:LexicalSense |
| mylex:sixteenmo_fd24c341.3793.48ab.b978.9ef848fa8f89 | lemon:sense | mylex:d85533bb-0326-42c2-b2d8-626aa8c95a74 |
| mylex:sixteenmo_fd24c341-3793-48ab-b978-9ef848fa8f893 | gold:acousticRealization | "siks`tēn-mō" |
| mylex:sixteenmo_fd24c341-3793-48ab-b978-9ef848fa8f893 | rdf:type | lexsil:SpokenRepresentation |
| mylex:sixteenmo_fd24c34 | | mylex:sixteenmo_fd24c |

# Dictionaries as Datapoints: Matthew Lee



This is a single entry from the example database from before.

Dictionaries as Datapoints: Matthew Lee



This is all of the data from the dictionary page before. Once the relationships are defined, we have 884 Facts.

This is now an automated conversion...so I just have to point the process to an exported lexicon and push play.

# Test Databases for the Project

| Large Databases | Small (Student) Databases |
|---|---|
| | *Field Methods/Data Mgmt.* |
| • **Badwee** (Cameroon) | • Fe?fe? (Cameroon) |
| Keith Beavon | • Hebrew (Israel) |
| • **Bonggi** (Malaysia) | • Japanese (Japan) |
| Dr. Michael Boutin | • Laari (Central Afr. Rep.) |
| • **Ewondo** (Cameroon) | • Mandarin (China) |
| Dr. Essono | |
| • **Marwari** (India) | |
| Jonathan Dailey | |

Results:

Results: Small Lexicons

- The 5 Student Databases that I transformed:
  - Laari
  - Japanese
  - Hebrew
  - Fe?fe?
  - Mandarin

- Total: 146,373 Triples (pieces of information).

## Results: Large Lexicons

- The 4 Large Lexicons that I transformed:
  - Badwee
  - Bonggi
  - Ewondo
  - Marwari

- Total: 2,164,469 Triples.

- An import of Webster's Dictionary contained over 5 million triples by itself.

Ontology

Now that the data is linked, we can start to ask questions…

Dictionaries as Datapoints: Matthew Lee



The primary database I used, Allegrograph, has two interfaces: this one for searching, and [click]

Dictionaries as Datapoints: Matthew Lee



this one for Browsing

# Search Syntax

- The primary tool used for searching RDF is a language called SPARQL.
- This language is related to SQL and is very powerful, but as you'll see, it has a learning curve.

## Simple Queries

Which words have a gloss of "monkey"?

```
select Distinct ?lexeme ?gloss {?s lemon:sense ?sense.
            ?s lemon:lexicalForm ?lexeme.
            ?sense lexsil:hasGloss ?gloss .
              ?sense skos:narrowMatch ?sd
              FILTER regex(?gloss,'monkey').
                FILTER(langMatches(lang(?gloss), "EN"))}
```

Simple Query (monkey)

| lexeme | gloss | Entry Name |
|--------|-------|------------|
| "hoʌtsəˌ\" | "monkey" | mandarin:hoʌtsəˌ_a204ba47.a0b2.4b20.b762.e53c01da23c5 |
| "hoʌtsəˌ˥" | "monkey" | mandarin:hoʌtsəˌ_a204ba47.a0b2.4b20.b762.e53c01da23c5 |
| "hoʌtsəˌ˥" | "monkey" | mandarin:hoʌtsəˌ_a204ba47.a0b2.4b20.b762.e53c01da23c5 |
| "saru" | "monkey" | japanese:ˈsɑ.ru_ad2af1e1.3892.435d.bdbf.b2d11ca1e604 |
| "ˈsa.ru" | "monkey" | japanese:ˈsɑ.ru_ad2af1e1.3892.435d.bdbf.b2d11ca1e604 |

## Simple Queries

**What are all the words that have a gloss or definition including "moon"?**

```
select Distinct ?lx ?def ?gloss ?s where {
{ ?ent lemon:lexicalForm  ?lx .
  ?ent lemon:sense ?s.
  ?s lexsil:hasGloss ?gloss
FILTER regex(?gloss,'^(lune|moon)$') }
UNION
{ ?ent lemon:lexicalForm ?lx .
  ?ent lemon:sense ?s.
  ?s lexsil:hasDefinition ?def
FILTER regex(?def,'(lune|moon)( |$)')  } }
Order by ?s
```

Simple Query (moon/lune)

| lx | def | gloss | Lang |
|---|---|---|---|
| "lomo" | | "lune" | Badwee |
| "buaidn" | | "moon" | Bonggi |
| "mata-buaidn" | | "moon" | Bonggi |
| "lugut" | "clouds covering moon or stars" | | Bonggi |
| "kelomon" | "very dark; no moon" | | Bonggi |
| "ǹkos" | "Bâton-talisman, bâtonnet magique, baguette de fée utilisée pour rendre un culte à la nouvelle lune. " | | Ewondo |
| "məmua" | "Pleine lune" | | Ewondo |

## Simple Queries

Select all words in the Semantic Domain "5.2.3.1 Food from Plants"

```
select Distinct ?lexeme ?gloss ?sd ?s{?s lemon:sense ?sense.
        ?s lemon:lexicalForm ?lexeme.
        ?sense lexsil:hasGloss ?gloss .
            ?sense skos:broadMatch <lexsil:Ddp4-5.2.3.1>.
            ?sense skos:narrowMatch ?sd.
            #FILTER regex(?gloss,'monkey').
                    #FILTER regex(?sd,'Ddp').
                FILTER(langMatches(lang(?gloss), "EN"))}
order by ?s
```

Dictionaries as Datapoints: Matthew Lee



## Simple Query (5.2.3.1 Food from Plants)

| lexeme | gloss | sd | Lx name |
|---|---|---|---|
| "anggur" | "grape" | lexsil:Ddp4-5.2.3.1.2 | bonggi:anggur_a539 8537.7494.4366.8e 1e.42e530f023e6 |
| "bembangan" | "type fruit" | lexsil:Ddp4-5.2.3.1.2 | bonggi:bembangan_ a30d03d6.075a.455 a.b89b.783cde2f80e d |
| "biabas" | "guava" | lexsil:Ddp4-5.2.3.1.2 | bonggi:biabas_21da 9228.d57f.4034.a34 e.54e2c155ba0f |
| "ə̃ŋgur" | "grape" | lexsil:Ddp4-5.2.3.1.2 | marwari:अंगूर_6093 8c51.406c.4167.9b4 4.66588b751505 |
| "अखरोट" | "nut" | lexsil:Ddp4-5.2.3.1.1 | marwari:अखरोट_24 08f60b.35cd.4cbb.8 8d1 a79abdf0d47f |

Notice that I have asked for "broad" matches. Semantic domains are hierarchical, so anything in category 5.2.3.1.2 is also contained one step up in 5.3.2.1.

I didn't have any related lexicons, and the lexicons I had didn't have much information on borrowing...so I thought that I wouldn't be able to show this with my dataset.  But…

Dictionaries as Datapoints: Matthew Lee



Take a look at the words for grapes from the previous search… This is part of the fun of cross-linguistic study.

How much data can an RDF database hold?

# RDF Scalability: DBPedia (Wikipedia)

- Dbpedia contains 119 languages
- Together **24.9 million** things
  - 24.6 million images
  - 27.6 million links to external web pages
  - 67.0 million links to Wikipedia categories
- Total: **2.46 billion** pieces of information (RDF triples)

Sahnwaldt, Christopher. 2013. DBpedia: About. http://dbpedia.org/About (7 May, 2014).

## Digital Stewardship

- Publishing digital data is a step toward SIL's End C:
- SIL exists to the end that:

*Individuals and communities benefit from our contribution to an increasing body of knowledge regarding the world's languages and cultures, and to the academic and professional disciplines related to our work in language development.*

- Could offer this as a service
  - Single requests for users (Web interface)
  - Internal Access
  - APIs and Subscriptions for Academic Organizations
- Make Structured Information Available as a service.

## Beyond the Proof of Concept

- My Current Test is small-scale.
- 7 million pieces of information across 9 Lexicons
- I used Allegrograph as a database server.
  - Offers Value-added features for work with RDF.
  - Free to use up to 5 million triples per graph.
- Running on a rented server.
- Access is SPARQL and Gruff-only (not yet user friendly).

## Build a More User-Friendly Interface

To be used by members of the language community and linguists, the database needs:

- A Simple web interface for asking easy questions and setting limits.

- Web-based graphical ways to view and manipulate graphs.

- Program interfaces (APIs) to allow high-power computation.

and the Interesting Questions...

- How can this be expanded to include other elements in FLEx (Grammar, Phonology, Texts, Discourse)?

- How can we help the user with Disambiguation to link to other external URIs?

- Can this become my Master's thesis?

## Major References

Chiarcos, Christian, Sebastian Nordhoff & Sebastian Hellman (eds.). 2012. *Linked Data in Linguistics*. Springer.

McCrae, John, Dennis Spohr & Philipp Cimiano. 2011. Linking Lexical Resources and Ontologies on the Semantic Web with lemon. *The Semantic Web: Research and Applications*, 245–259. http://pub.uni-bielefeld.de/luur/download?func=downloadFile&recordOId=2278403&fileOId=2526034 .

Simons, Gary F. 2005. Beyond the Brink: Realizing Interoperation through an RDF Database. *Linguistic Ontologies and Data Categories for Linguistic Resources*. Cambridge, MA. http://emeld.org/workshop/2005/papers/simons-paper.pdf (8 December, 2013).

Simons, Gary F. & H. Andrew Black. 2008. Third Wave Writing and Publishing. http://www-01.sil.org/silepubs/Pubs/52287/SILForum2009-005.pdf