# "Did You Mean...?" and Dictionary Repair: from Science to Engineering

Michael Maxwell, presented by Petra Bradley
(mmaxwell@casl.umd.edu, pbradley@casl.umd.edu)
University of Maryland, College Park MD 20742 USA

AMTA 2016

29 October 2016

# Two problems with electronic (XML) dictionaries

- How can you look up a word? Especially when…
  - You can't spell gud
  - There's noise in the environment
  - You can't hear well
  - It's in a foreign language
- What happens when the dictionary has an error?
  - Missing information, information in the wrong field, wrongly structured fields
  - Errors range from typos to missing pages' worth of data
  - CASL has worked with dozens of XML dictionaries, mostly from government sources; no dictionary has been without errors
  - Errors are pervasive: ~10,000 errors in one Urdu dictionary (with ~50,000 entries)

# CASL's solutions

- "Did You Mean…?" (DYM)
  Smart fuzzy lookup. Three versions:
  - ‣ Basic DYM: Spell correction only
  - ‣ Morphologically Aware DYM (MADYM): Spell correction + morphological parsing
  - ‣ Cross-language DYMs: Basic or MADYM + lookup in multiple dialect/ language/ script dictionaries
    - - Waziri Pashto: lookup in 1902 Waziri dictionary, plus backoff lookup in modern Pashto dictionary via sound changes
    - - French-in-Arabic: lookup of Arabic script queries in Moroccan Arabic dictionary *and* in French dictionary (with Arabic morphological parsing)
- Dictionary repair tools
  - ‣ ADALT discovers anomalies: rare structures, data that doesn't "fit the mold"
  - ‣ VELMA is an editor specialized for dictionaries and similar documents

ADALT + VELMA and DYM are *tools*; neither replaces human labor

- J.C.R. Licklider, 1960 *Man-Computer Symbiosis*: "Computing machines can do readily, well, and rapidly many things that are difficult or impossible for man, and men can do readily and well, though not rapidly, many things that are difficult or impossible for computers. That suggests that a symbiotic cooperation, if successful in integrating the positive characteristics of men and computers, would be of great value."

- ADALT and VELMA complement human abilities in error-finding

- DYM enables humans to better understand critical foreign language texts by doing smart fuzzy lookup

# Two problems with electronic (XML) dictionaries

- How can you look up a word? Especially when…
  - ‣ You can't spell gud
  - ‣ There's noise in the environment
  - ‣ You can't hear well
  - ‣ It's in a foreign language
  - ‣ *"Did You Mean...?"*
- What happens when the dictionary has an error?
  - ‣ Missing information, information in the wrong field, wrongly structured fields
  - ‣ Errors range from typos to missing pages' worth of data
  - ‣ CASL has worked with dozens of XML dictionaries, mostly from government sources; no dictionary has been without errors
  - ‣ Errors are pervasive: ~10,000 errors in one Urdu dictionary (with ~50,000 entries)

# How we got there: DYMs

- 2007: Dr. Anton Rytting (CASL) suggests need for spell correction for dictionary lookup by Arabic language analysts
- 2008: First "Did You Mean…?" for Georgetown Iraqi Arabic dictionary.
- 2009: Second DYM: Urdu, overwhelmingly positive response.
- Later: Additional Arabic dictionaries, gazetteers, 10,001 Arabic names; Pashto, Russian, Ukrainian, Somali, Persian (Farsi), Swahili…
- 2012-2013: Dr. Corey Miller adds Persian morphological parser to Persian DYM to produce first Morphologically Aware DYM (MADYM)
- Later: MADYMs for Arabic, Somali, Swahili, Korean; DYMs for Chinese (pinyin), Dhivehi, Punjabi, Portuguese (Open Street Maps for Rio de Janeiro)
- 2016:
  - ‣ DYM Toolkit for building (part basic of) DYMs for new languages
  - ‣ DYM Platforms for deploying basic DYMs and MADYMs

# Electronic Dictionary Lookup

**Ordinary dictionary lookup:** type in a word, get back definition(s) *if* you spelled it correctly, and you chose the dictionary citation form

**Wildcard lookup:** For each letter you're uncertain about, type in a '*'
If you hear German [rat], type *ra**.
But this will give you unwanted words: *Rah, Rap, Rank, rar, rau, Rat,* and *Rad*–whereas only the last two are likely to be what you want.

**Regular expression lookup:** If you think it might end in a 'd' or a 't', type *ra[d|t]*. Works if you understand regular expressions, and if you know where the mistakes are likely to be, and if you remember what the possibilities are.

**DYM:** Builds in the knowledge of regular expressions, likely mistakes, possible substitutions (and deletions and insertions), and *how likely* a particular mistake is; and automagically applies this knowledge everywhere in the word where a mistake is possible. Type in *rat,* get back *Rat* and *Rad* (but not the others).

# How does a DYM work? outside the box

Urdu Romanization. th = aspirated dental, t = unaspirated dental, Th = aspirated retroflex.
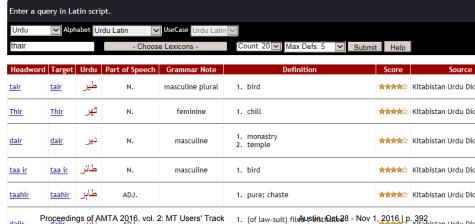


## Did You Mean...?
[Version 2.0]

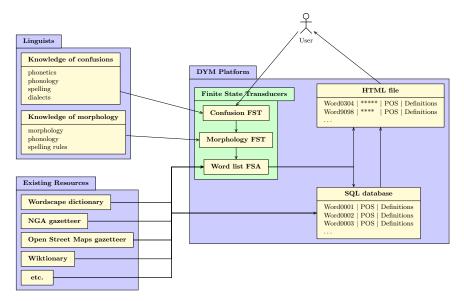CENTER FOR ADVA
STUDY OF LANG

Enter a query in Latin script.

| Urdu ▼ | Alphabet | Urdu Latin ▼ | UseCase | Urdu Latin ▼ | | |
| thair | | - Choose Lexicons - | | Count: 20 ▼ | Max Defs: 5 ▼ | Submit | Help |

| Headword | Target | Urdu | Part of Speech | Grammar Note | Definition | Score | Source |
|----------|--------|------|----------------|--------------|------------|-------|--------|
| tair | tair | طیر | N. | masculine plural | 1. bird | ★★★★☆ | Kitabistan Urdu Dic |
| Thir | Thir | ٹھہر | N. | feminine | 1. chill | ★★★★☆ | Kitabistan Urdu Dic |
| dair | dair | دیر | N. | masculine | 1. monastry 2. temple | ★★★★☆ | Kitabistan Urdu Dic |
| taa ir | taa ir | طائر | N. | masculine | 1. bird | ★★★★☆ | Kitabistan Urdu Dic |
| taahir | taahir | طاہر | ADJ. | | 1. pure; chaste | ★★★★☆ | Kitabistan Urdu Dic |
| daiir | daiir | دائر | ADJ. | | 1. (of law-suit) filed; instituted 2. circling ; whirling | ★★★★☆ | Kitabistan Urdu Dic |

# How does a DYM work? inside the box

# But that's so complicated!

CASL tools to build DYMs:



- DYM Toolkit to incorporate the linguists' knowledge of confusions
  - ‣ DYM Toolkit enables end users (e.g. language analysts) to build confusion matrices for their confusions.
  - ‣ The toolkit can also be used to build confusion matrices for dialects, non-standard spelling systems, native speaker errors.
    - In conjunction with a transliterator, this can be used for Romanized scripts (Arabeze, Pinglish/ Farslish,…)
- Dictionaries, gazetteers often exist (but stay tuned for part 2 of this talk)
- DYM Platform takes care of the software infrastructure

# But that's so complicated!

Morphological parsing...



...is still rocket science (and optional).

# DYM software

- Platforms 1 and 2, with DYMs, submitted as deliverable (includes documentation and about ten DYMs); USG has free use.
  - Also available with consulting/ training.
  - DYM Toolkit included.
- What do I need to build a new DYM?
  - Dictionary or other lexical resource (e.g. gazetteer)
    If it needs cleanup…stay tuned!
  - Confusion matrix
    - Available for non-native listener confusions for ten languages
    - …or language analysts can build their own with the provided Toolkit
  - …or talk to us about building a new one.

# DYM software (cont'd)

Planned work

- Cross-language DYMs (like Waziri Pashto, French-in-Arabic)
  - ‣ CASL has built such DYMs, but current Platforms do not support them
    …future work.
- Stand-alone DYMs (without network connection/ server)
  - ‣ This could also be used to rapidly field stand-alone dictionaries, even without a 'real' confusion matrix.
- Coming soon to your cellphone?

# Two problems with electronic (XML) dictionaries

- How can you look up a word? Especially when…
  - You can't spell gud
  - There's noise in the environment
  - You can't hear well
  - It's in a foreign language
- What happens when the dictionary has an error?
  - Missing information, information in the wrong field, wrongly structured fields
  - Errors range from typos to missing pages' worth of data
  - CASL has worked with dozens of XML dictionaries, mostly from government sources; no dictionary has been without errors
  - Errors are pervasive: ~10,000 errors in one Urdu dictionary (with ~50,000 entries)
  - *ADALT and VELMA*

# VELMA: dictionary editor

- VELMA = Visual Environment for Lexicography and MAchine learning
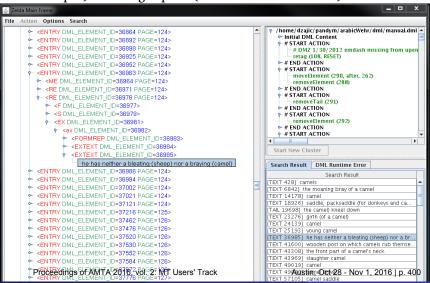




New version = Zelda

- VELMA is designed for editing of existing dictionaries (not building new ones)
  - ‣ Allows easy manipulation of existing XML nodes and text…
  - ‣ …meaning Ordinary Working Lexicographers (not computer scientists)* can modify dictionaries.
  - ‣ Integrated with ADALT.
  - ‣ Uses ADALT output to suggest a workflow for linguists/ lexicographers.
    *Ok, you need to learn about XML…

# VELMA in pictures

Dictionary file on left, actions performed upper right pane; search results (or ADALT output) lower right pane (here: search for 'camel').

# VELMA in pictures

Supports typical editor operations: search and replace, drag-and-drop, change text (or XML tags), undo (of *any* action, not just most recent)

# So what is this ADALT stuff?

How do you find errors in the electronic equivalent of hundreds of pages of a dictionary?

- Errors come in two types:
  - ‣ Frequent: Some common structure is mis-represented throughout
    - Easy for humans to notice at least some of them.
    - …but finding all instances of an error type can be hard!
  - ‣ Rare: One-offs, typos
    - Needle-in-a-haystack
    - Hard for humans to find these…
    - …especially when you don't know what kinds of needles there are!

# Finding those needles in a haystack

You need…

# CASL's metal detector

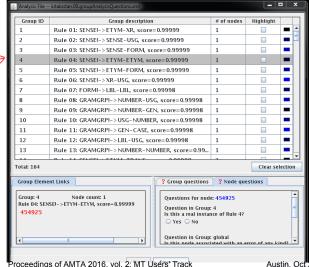**ADALT** (Automatic Detection of Anomalies in Lexicographic Text)

- ADALT uses machine learning (AI) to find rare structures
- These are anomalies; some anomalies are errors, some are simply rare structures.
- Examples:
  - ‣ (Spanish, but based on actual anomaly in an Urdu dictionary):
    **raro** *adj.* (Usage: rare)
    Missing definition!
  - ‣ (English, but similar anomalies found in Urdu dictionary):
    **colo(u)r** ......
    Prevents lookup: user will never search for this spelling!
  - ‣ (Illustrative; actual errors would be in XML)
- ADALT builds models of structures found in a particular dictionary, and uses heuristics to find structures which are rare within those models.
- It therefore doesn't know (and doesn't need to be told) what is 'correct'.
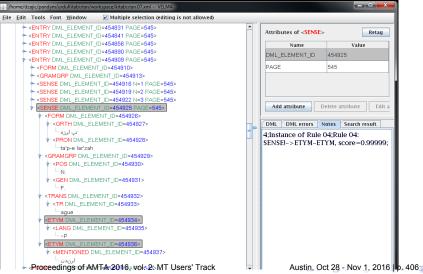
# Example of ADALT usage: the metal detector

The metal detector has found anomalies; user chooses one to work on (a single lexeme that has two etymologies):
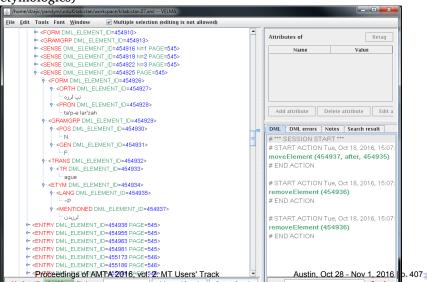
# Example of ADALT usage: The needle is found!

VELMA jumps to that anomaly; the first 'etymology' contains just a language tag (Persian), the second lacks a language tag.

# Example of ADALT usage: grab the needle

The user uses VELMA to remove the needle from the haystack (merge the two etymologies)

# …about those frequent errors

Latest version of VELMA (Zelda!) provides "query-by-example" (QBE) and "edit-by-example" (EBE)

- If the user notices multiple of instances of some structural error…
  - ‣ Dozens might be evident on inspection
  - ‣ …and hundreds more might be lurking.
- …then use QBE to find all instances of that structure (or similar structures).
- Once the user has verified that QBE finds the desired erroneous structures, correct one, then use EBE to correct others in the same way.

# ADALT and VELMA software

- Mature: CASL has used this (or previous versions) on about a dozen dictionaries.
- VELMA is preferable to text editors/ Perl/ version control systems because...
  - ‣ Users are less likely to make mistakes
  - ‣ ...and if you do make a mistake, it can be undone six months from now without messing up everything done between now and then.
  - ‣ Provides audit trail of changes.
  - ‣ Users report working faster.
- Submitted as deliverable (includes documentation); USG has free use.
  - ‣ Also available with consulting/ training.
- What can I do with this?
  - ‣ Import/ clean up dictionaries
  - ‣ Import/ clean up gazetteer data
  - ‣ Potentially useful for other kinds of semi-structured text data

# Summary

- DYMs and dictionary repair started as CASL research projects less than ten years ago.
- Many useful tools and resources (cleaned up dictionaries, DYMs) came out of the work along the way.
- Dictionary repair and building DYMs are now (mostly) engineering.
  - Fielding a new DYM can be done in a few days.
  - Repairing a dictionary can take a few months (rather than a year or more)
    ...and the output is of higher quality than CASL's early dictionaries.
- These tools help–not replace–human beings.
- Contact information:
  - Dr. Mike Maxwell (Technical Director for HLT): mmaxwell@casl.umd.edu 301-356-2639
  - Dr. David Zajic (Research Scientist): dzajic@casl.umd.edu 301-356-8995
  - Dr. Mike Bunting (CASL Executive Director): mbunting@casl.umd.edu 301-356-8894

# Current DYMs

- Arabic (Iraqi, MSA, Sudanese, Levantine…; dictionaries, gazetteers, 10,001 Arabic names); *French in Arabic (code switching)
- Urdu, Pashto, *Waziri Pashto, Punjabi
- Russian, Ukrainian
- Somali, Maay (related to Somali), Swahili, Chimwiini (related to Swahili)
- Persian (Farsi, Dari)
- Korean
- Mandarin Chinese (pinyin search)
- Dhivehi

*Not available in curent Platforms