

# The Matrix: Future Directions

## Wrap up

---

Ling 567

March 8, 2010

# Overview

---

- Wrap up/reflections
- Matrix: Future directions
- MT Lab questions

# Goals: Of Grammar Engineering

---

- Build useful, usable resources
- Test linguistic hypotheses
- Represent grammaticality/minimize ambiguity
- Build modular systems: maintenance, reuse

# Goals: Of this course

---

- Mastery of tfs formalism
- Hands-on experience with grammar engineering
- A different perspective on natural language syntax
- Practice building (and debugging!) extensible system
- Contribute to on-going research in multilingual grammar engineering

# Reflections

---

- Where have the analyses provided by the Matrix (or suggested by the labs) seemed like a good fit?
- Where have they been awkward?
- What have you learned in this class about syntax?
- ... about knowledge engineering for NLP?
- ... about computational linguistics in general?
- ... about linguistics in general?

# Feedback: Pair projects

---

- How did you divide the work?
- In what ways was having a partner helpful?
- Would you have learned more working on your own?

# Future directions overview

---

- More libraries (and semantic harmonization)
- How this class will evolve
- Auto-generated transfer rules
- Lexical acquisition
- Ontological annotation
- Matrix-ODIN Mash-up
- Typological seeding of statistical NLP

# More libraries

---

- New this year: Argument optionality, updates to word order
- Demonstratives
- Extensions/retrofits to negation, questions, coordination
- (more) extensions to word order
- Non-verbal predicates
- Intersective modifiers
- Numeral classifiers
- More verb subcategorization
- Embedded clauses
- Marking of information structure

# Evolution of 567

---

- New phenomena: Wh-questions, relative clauses, while-clauses ...?
- Ever bigger jump start --- reaching the limit on this one?
  - Would working in groups of three make it possible to get to even bigger grammar fragments?
- Partnership with field linguists
- Work with small corpora

# Autogenerated transfer rules

---

- Identify “grammaticalized” differences in MRSs
- “Publish” choices along these dimensions for each grammar
- Create a library of transfer rules from property to property:
  - pro-drop to pronouns (and vice versa)
  - Mismatches in demonstrative distinctions
  - can <> the possibility exists

# Autogenerated transfer rules

---

- Use language-specific pred values
- Create transfer rules on the basis of PanDictionary or other lexical resources
- Measure the extent of translation divergence
- Use bitexts and statistical methods to detect word pairs requiring more than straight pred-mapping transfer rules

# Lexical acquisition

---

- How can we import lexical entries from other linguistic resources (e.g., FIELD lexicons, ODIN)?
- How big do the grammars have to get before we can embark on (semi-)automated lexical acquisition?
- To what extent do the lexical properties of translational equivalents predict lexical properties in another language?
- How can we most effectively leverage human effort?
- How do we know when we're missing an appropriate type?

# Ontological annotation

---

- Annotate grammars with links to GOLD (Farrar & Langendoen 2003)
  - Locate which constraints contribute to which phenomena
  - Index analyses for discovery in grammars and treebanks
- Annotations in Matrix core
- Annotations in customization system
- Support for user annotation

# Matrix-ODIN Mash-up

---

- ODIN: Online Database of INterlinear glossed text (Lewis 2006)
- Lewis & Xia 2007 explore learning typological properties from ODIN data
- Next steps:
  - Answer Matrix customization system questionnaire automatically
  - Including lexical information
  - ... and information about affixes
- Step 3: Automatically generated precision grammars!

# Seeding statistical NLP with typological knowledge

---

- Haghghi & Klein 2006: Unsupervised parsing (“Prototype Driven Grammar Induction”)
- Syntax-based statistical MT is finally coming into its own (e.g., Callison-Burch’s talk last week)
- Matrix Customization system-generated starter grammars represent a middle ground between broad-coverage precision grammars and coarse-grained typological information (as in WALS).
  - Testable over hand-constructed test suites
  - Usable to create prototype trees or even translation pairs

# Overview

---

- Wrap up/reflections
- Matrix: Future directions
- MT Lab questions