

**Why Laboratory Syntax?**

The traditional ways of data collection are deeply flawed:

- statistically unreliable
- highly subjective
- fraught with confounding factors that may have nothing to do with language

—Maria Polinsky. 2005. Linguistic Typology and Grammar Construction. LSA Workshop on Typology in American Linguistics: An Appraisal of the Field, January 9, 2005, Oakland.

# I. Problems with introspective data

## Example 1: ‘The Affectedness Constraint’

the city’s destruction

the boy’s removal

the picture’s defacement

\*the event’s recollection

\*the problem’s perception

\*the picture’s observation

—A. Giorgi and G. Longobardi. 1991. *The Syntax of Noun Phrases: Configuration, Parameters and Empty Categories*, pp. 140ff. Cambridge University Press.  
(and many others)

Taylor: possessors have to be *topical* and *informative* relative to the possessed.

‘Concerning those events, their recollection still frightens me.’

‘Concerning that problem, its perception varies from person to person.’

‘Concerning that picture, its careful observation will reveal many interesting details.’

—J. R. Taylor, 1994. “Subjective” and “Objective” readings of possessor nominals. In *Cognitive Linguistics* 5: 201–242.

from the WWW:

Certainly, between the presentation of information to the senses and **its recollection**, various cognitive processes take place.

Lesson 2: Sound Properties and **Their Perception**.

But the standard idea that an event is inseparable from **its observation** is just scientific silliness.

Lesson:

Introspective judgments about decontextualized examples may underestimate the space of grammatical possibility

## Example 2: Verbal Subcategorization

Verbs taking APs but not participles:

Kim turned out political.

\*Kim turned out doing all the work.

Kim ended up political.

\*Kim ended up sent more and more leaflets.

—Pollard and Sag (1994: 105–108) (and many others)



Usage in *New York Times* contradicts these observations:

But it turned out having a greater impact than any of us dreamed.

On the big night, Horatio ended up flattened on the ground like a fried egg with the yolk broken.

—Chris Manning. 2003. Probabilistic syntax. In *Probabilistic Linguistics*, ed. by R. Bod, J. Hay, and S. Jannedy. 289-341.

Manning (2003: 300f):

“What is going on here? Pollard and Sag’s judgments seem reasonable when looking at the somewhat stilted “linguists’ sentences.” But with richer content and context, the *New York Times* examples sound (in my opinion) in the range between quite good and perfect. None of them would make me choke on my morning granola. They in all likelihood made it past a copy editor.”

Manning:

“... the subcategorization frames that Pollard and Sag do not recognize are extremely rare, whereas the ones they give encompass the common subcategorization frames of the verbs in question.”

“Corpus frequencies can be used to quantify linguistic intuitions and lexical generalizations such as Levin’s (1993) semantic classification...” (BNC data)

–Maria Lapata. 1999. Acquiring lexical generalizations from corpora: A case study for diathesis alternations. In *Proceedings of the 37th Meeting of the North American Chapter of the Association for Computational Linguistics*, 397–404. College Park, Maryland.

Lesson:

Introspective judgments about constructed examples may reflect relative frequency within the space of grammatical possibility.

### Example 3: the dative alternation

*That movie gave me the creeps.*

*\*That movie gave the creeps to me.*

*The lighting here gives me a headache.*

*\*The lighting here gives a headache to me.*

—Oehrle 1976 and many linguists thereafter; recently in *Linguistic Inquiry* (2001: 261)

Joan Bresnan and Tatiana Nikitina. 2003. "On the  
Gradience of the Dative Alternation".

<http://www-lfg.stanford.edu/bresnan/download.html>

tried Google

GIVE THE CREEPS TO



many examples like these:

This life-sized prop will **give the creeps to just about anyone!** Guess he wasn't quite dead when we buried him!

...Stories like these must **give the creeps to people whose idea of heaven is a world without religion...**

**GIVE A HEADACHE TO**

many examples like these:

She found it hard to look at the Sage's form for long. The spells that protected her identity also **gave a headache to anyone trying to determine even her size**, the constant bulging and rippling of her form gave Sarah vertigo.

Design? Well, unless you take pride in **giving a headache to your visitors** with a flashing background? no.

Compare:

\*That movie gave the creeps to me.

...Stories like these must **give the creeps to people whose idea of heaven is a world without religion...**

??Stories like these must **give people whose idea of heaven is a world without religion the creeps...**

That movie gave me the creeps.

Bresnan, Cueni, Nikitina, and Baayen (2005):

The longer phrase is placed at the end — *the principle of end weight*. (Behaghel 1910, Wasow 2002)

Idioms like *give the creeps* have a strong bias toward the double object construction, but the principle of end weight overrides it.

Linguistics textbook data:

Ted denied Kim the opportunity to march.

\*Ted denied the opportunity to march to Kim.

The brass refused Tony the promotion.

\*The brass refused the promotion to Tony.

Georgia Green. 1971. Some implications of an interaction among constraints. CLS 7. 85-100.

\*Ted gave Joey permission to march, but he denied Kim it.

Ted gave Joey permission to march, but he denied it to Kim.

\*The brass gave Martin permission to sit, but they denied Tony it.

The brass gave Martin permission to sit, but they denied it to Tony.

Lesson:

Introspective judgments about constructed examples may fail to reflect the interactions of multiple conflicting constraints, including processing constraints.



Further examples: the benefactive alternation

Chris baked/bought/decorated/sliced a cake  
for Kim.

Chris baked/bought Kim a cake

\*Chris decorated/sliced Kim a cake.

See Christiane Fellbaum. 2005. Examining the constraints on the benefactive alternation by using the World Wide Web as a corpus. In *Evidence in Linguistics: Empirical, Theoretical, and Computational Perspectives*, ed. by M. Reis and S. Kepser. Mouton de Gruyter.

## Further examples: *wh*- questions

Hofmeister, Philip, T. Florian Jaeger, Ivan A. Sag, Inbal Arnon, and Neal Snider. Locality and Accessibility in Wh-Questions. To appear in the proceedings of the conference: Linguistic Evidence: Empirical, Theoretical, and Computational Perspectives, Tbingen, 2-4 February 2006 (hosted by SFB441: 'Linguistic Data Structures'), University of Tbingen, Germany. On-line, Stanford: <http://lingo.stanford.edu/sag/publications.html>.

Sag, Ivan A., Inbal Arnon, Bruno Estigarribia, Philip Hofmeister, T. Florian Jaeger, Jeanette Pettibone, and Neal Snider. Processing Accounts for Superiority Effects. Under Review. On-line, Stanford: <http://lingo.stanford.edu/sag/publications.html>

## Summary:

Introspective judgments about decontextualized, constructed examples...

- may underestimate the space of grammatical possibility because of absence of context
- may reflect relative frequency within the space of grammatical possibility
- may fail to reflect the interactions of multiple conflicting constraints, including processing constraints

## II. Problems with corpus data

Corpus studies of English have found that various properties of the recipient and theme have a quantitative influence on dative syntax (Thompson 1990, Collins 1995, Snyder 2003, Gries 2003, ao):

discourse accessibility

relative length

pronominality

definiteness

animacy

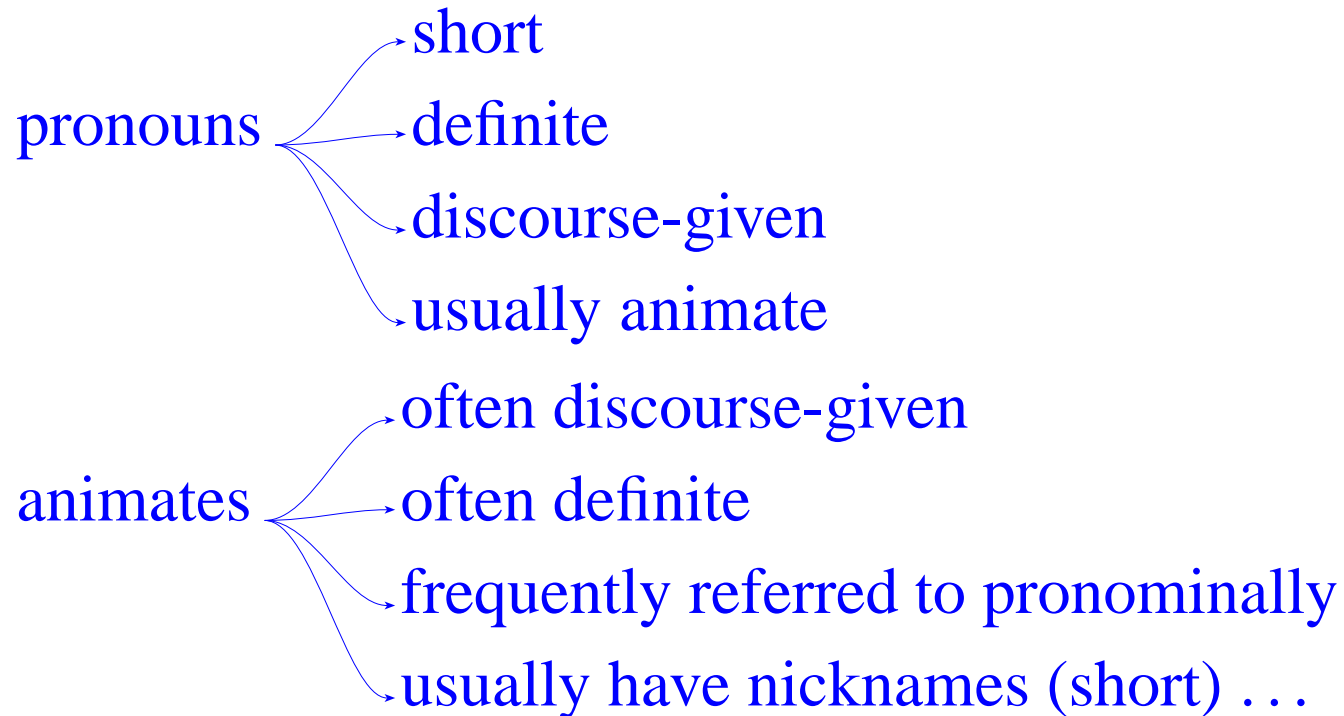
⇒

dative construction choice

Yet what really drives the English dative alternation remains unclear...

# 1. The problem of confounds

What really drives the dative alternation remains unclear because of pervasive correlations in the data:



Correlations tempt us into reductive theories that explain effects in terms of just one or two variables (e.g. Hawkins 1994, Snyder 2003)



A beautifully simple theory:

1. **Givenness correlates with shorter**, less complex expressions (less description needed to identify)
2. **Shorter expressions occur earlier** in order to facilitate parsing (more complex after less)

Apparent effects of givenness (and correlated properties like animacy) could reduce to the preference to process syntactically complex phrases later than simple ones (Hawkins 1994).

## 2. The problem of pooling different speakers' data

## A persistent question about corpus studies of grammar ...

in Newmeyer's (2003: 696) words:

“The Switchboard Corpus explicitly encompasses conversations from a wide variety of speech communities. But how could usage facts from a speech community to which one does not belong have any relevance whatsoever to the nature of one's grammar? **There is no way that one can draw conclusions about the grammar of an individual from usage facts about communities, particularly communities from which the individual receives no speech input.**”

### 3. The problem of lexical biases

Observations of NP properties are not independent of their verbs. Just as the pooling of data from different speakers introduces unknown dependencies among the observations, so does the pooling of NP observations from different verbs.

The properties of recipients and themes depend on the verbs which describe the transfer events they are participating in. For example:

*bring* is nearly three times more likely to have a given recipient than *take*

*take* is over seven times more likely to have a nongiven recipient than *bring*.

(The goal of bringing is usually located near the speaker, the goal of taking is usually located away from the speaker)

## 4. The problem of cross-corpus differences

*Does it make sense to relate frequencies of usage to grammar? (Keller and Asudeh 2002: 240)*

*After all, unlike the grammaticality of a linguistic form, which is an idealization over usage, the actual frequency of usage of a form is a function of both grammatical structure and extra-grammatical factors such as memory limitations, processing load, and the context.*



In fact it is true that

the frequencies of double-object constructions in the Switchboard collection of recordings of telephone conversations  $\neq$

frequencies in the Treebank Wall Street Journal collection of news and financial reportage

V NP NP's = 79% of total Switchboard datives  
( $n = 2360$ )

V NP NP's = 62% of total Wall Street Journal datives ( $n = 905$ )

## Summary:

Corpus data are problematic because...

- correlated variables introduce confounds
- pooled data from different speakers may invalidate grammatical inference
- observed properties of nouns may depend on verbs
- cross-corpus differences appear to undermine the relevance of corpus studies to grammatical theory

III. Are controlled psycholinguistic experiments the answer?

Psycholinguistic experiments usually involve constructed sentences isolated from connected discourse. Experimental data lack discourse cohesion and subjects resort to default referents. The prompts used in experiments also have significant effects.

—Douglas Roland and Daniel Jurafsky. 2002. Verb sense and verb subcategorization probabilities. In *The Lexical Basis of Sentence Processing. Formal, Computational and Experimental Issues*, ed. by P. Merlo and S. Stevenson, 325–345. Amsterdam: Benjamins.

“Factorial designs are commonly used where regression is more appropriate. . . . Psycholinguists are generally very reluctant to include covariates in their analyses. . . . When relevant covariates are not taken into account, the conclusions suggested by one’s model may be unwarranted.”

—R. Harald Baayen. 2004. Statistics in psycholinguistics: A critique of some current gold standards.

<http://www.mpi.nl/world/persons/private/baayen/publications.html>

**Why Laboratory Syntax?**

These problems can be solved using modern statistical modeling techniques on both corpus and psycholinguistic data (Roland and Jurafsky 2002, Baayen 2004, Bresnan et al. 2005)

For Thursday:

Joan Bresnan. 2005. “A Few Lessons from Typology”. Comments from the LSA Workshop: Typology in American Linguistics. An Appraisal of the Field. LSA 79th Annual Meeting, Oakland, January 9, 2005. 11 pages.

<http://www.stanford.edu/~bresnan/download.html>



- Read and evaluate “A Few Lessons from Typology”.
- Test the thesis of the typology lessons paper by Googling data from a syntactic or semantic domain that interests you.

## Example from last night:

From: K.P. Mohanan <ellkpmoh@nus.edu.sg>

To: bresnan@stanford.edu

Date: Sep 25, 2006 8:29 PM

Dear Joan,

Would you happen to know how to explain the asymmetry between (1d) and the rest? I haven't seen any instances of proximity agreement when the conjunction is "and" instead of "or"

- 1a. Three boys and one girl are/\*is in the room.
- 1b. One girl and three boys are/\*is in the room
- 1c. There are/is three boys and one girl in the room.
- 1d. There is/\*are one girl and three boys in the room.

Mohanan

Google “there are one \* and three”:

Above  $4 \times 10^{19} eV$ , there are one triplet and three doublets within separation angle of  $2.5^\circ$  and the probability of observing these clusters by a chance ...

The Adsorption and Reaction of a Titanate Coupling Reagent on the ... There are one isopropoxy and three organic long-chains in the structural formula of CA7.

Within the Kurdistan Regional Government, there are one Turkmen and three Christian cabinet ministers.

There are one homeowner and three renters on the block, in addition to the church, and many of the properties have been vacant for as long as 20 years, ...

Google “there is one \* and three”:

There is one female and three males: all of them are pure black except for the

A doctrine of Christianity that there is one God and three divine persons in the one God: the Father, the Son (Jesus), and the Holy Spirit. ...

Figure A-4 shows a configuration in which there is one SSM and three SSCs.

There is one nurse and three technicians for every nine patients.

On Thursday:

- You will INSTALL R on your own computer.

Go to

<http://lib.stat.cmu.edu/R/CRAN/>

download the binaries for your platform, and read the directions on installation.