# DISTRIBUTION OF SOME LINGUISTIC FEATURES IN SOME TYPES OF DISCOURSE

A survey presented to the Utilika Foundation
by S. M. Colowick

December 2007

*Task: Perform a literature review and current-knowledge summarization about the question: How (i.e. in what proportions) are illocutionary forces, dialog acts, and predicate argument structures distributed in publication and interaction corpora?*

## CONTENTS

# 1 Illocutionary force

## 1.1 Introduction

The 2008 International Conference on Discourse and Grammar will have the theme "Illocutionary force, information structure and subordination between discourse and grammar" (http://liquid.hogent.be/~disgram/). The call for papers includes this sample topic: "What is the exact distribution of illocutionary force in discourse?"

Not only does it seem that no one has found an answer to this question, but it appears to the accidental linguist that no one has even attempted to find it. This lack of exploration could be due to a number of factors, beginning with disagreement as to what exactly the term "illocutionary force" means. For Austin (1962), who first articulated the concept, it referred to the effect that an illocutionary act was meant to have. He first theorized that only some speech acts, the *performatives* (e.g., "I bet" and "I bequeath"), actually *do* anything; he later acknowledged the existence of other types of illocutionary acts, but they all performed some action. According to one linguistics textbook, however, illocutionary force is simply "the intended meaning of an utterance" (O'Grady, Archibald, Aronoff, & Rees-Miller, 2001).

Another possible explanation for the dearth of data is the tendency toward extreme specialization in linguistics. Papers abound on the distribution of specific constructions in specific languages or in highly specific texts (e.g., "The Distribution of the Imperative, Future, and Infinitival Imperative in Dialectal Cretan Inscriptions"), but no one seems to have thought it worthwhile to do a broad-based analysis. The good news is that some people specialize in training computers to classify utterances, and those people's research is the source of most of the data presented here. The bad news is that, because researchers differ in how they classify utterances, any comparisons and aggregation of data are all but impossible. Still, as long as the experts have yet to resolve this issue, we should be happy to get even a rough approximation.

## 1.2 Generalizations from the data

The first five taxonomies presented below share enough tags to allow for fairly easy comparison. The others offer less comparable data, but still contribute somewhat to the following observations:

- Questions of all types constitute less than 5% of the dialog in email and phone interactions, but account for more than 12% of the utterances in instant messaging.

- In the help-desk environment, "Thanking" is the second most frequent type of utterance (after "Statement"), making up about 15% of both email and IM utterances. "Request" is in the top five (10% for email and 6% for IM).

- In general chat situations, greeting and saying good-bye make up about 5-14%of utterances.

- Yes-No questions tend to far outnumber the total of Yes and No answers. In other words, rhetorical and "backchannel" questions are usually grouped with other questions. The exception is SWBD-DAMSL, the only system of the six that seriously attempts to tag utterances with their true illocutionary force.

- In almost all cases, including SWBD-DAMSL, Yes-No Questions seem to occur more frequently than Wh-Questions.

- In almost all cases, the number of Yes answers is close to twice the number of No answers.

## 1.3 The data

### 1.3.1 SWBD-DAMSL

The 43 ttags in the SWBD-DAMSL classification system (Stolke et al., 2000) were adapted from Dialog Act Markup in Several Layers (Allen & Core, 1997). The SWBD-DAMSL set was developed for annotating telephone conversations; it is so extensive and so domain-specific that its utility in annotating written interactions is questionable. However, its attention to the nuances of categorization could be instructive.

| *SWBD-DAMSL Tag* | *Example* | *Count* | *Percent* |
|---|---|---|---|
| **Statement-non-opinion** | *Me, I'm in the legal department.* | 72,824 | 36.0 |
| **Acknowledge (Backchannel)** | *Uh-huh.* | 37,096 | 19.0 |
| **Statement-opinion** | *I think it's great* | 25,197 | 13.0 |
| **Agree/Accept** | *That's exactly it.* | 10,820 | 5.0 |
| **Abandoned or Turn-Exit** | *So, -* | 10,569 | 5.0 |
| **Appreciation** | *I can imagine.* | 4,633 | 2.0 |
| **Yes-No-Question** | *Do you have to have any special training?* | 4,624 | 2.0 |
| **Non-verbal** | *[Laughter], [Throat_clearing]* | 3,548 | 2.0 |
| **Yes answers** | *Yes.* | 2,934 | 1.0 |
| **Conventional-closing** | *Well, it's been nice talking to you.* | 2,486 | 1.0 |
| **Uninterpretable** | *But, uh, yeah* | 2,158 | 1.0 |
| **Wh-Question** | *Well, how old are you?* | 1,911 | 1.0 |
| **No answers** | *No.* | 1,340 | 1.0 |
| **Response Acknowledgement** | *Oh, okay.* | 1,277 | 1.0 |
| **Hedge** | *I don't know if I'm making any sense or not.* | 1,182 | 1.0 |
| **Declarative Yes-No-Question** | *So you can afford to get a house?* | 1,174 | 1.0 |
| **Other** | *Well give me a break, you know.* | 1,074 | 1.0 |
| **Backchannel in question form** | *Is that right?* | 1,019 | 1.0 |
| **Quotation** | *You can't be pregnant and have cats* | 934 | .5 |
| **Summarize/reformulate** | *Oh, you mean you switched schools for the kids.* | 919 | .5 |
| **Affirmative non-yes answers** | *It is.* | 836 | .4 |
| **Action-directive** | *Why don't you go first* | 719 | .4 |
| **Collaborative Completion** | *Who aren't contributing.* | 699 | .4 |
| **Repeat-phrase** | *Oh, fajitas* | 660 | .3 |
| **Open-Question** | *How about you?* | 632 | .3 |
| **Rhetorical-Questions** | *Who would steal a newspaper?* | 557 | .2 |
| **Hold before answer/agreement** | *I'm drawing a blank.* | 540 | .3 |
| **Reject** | *Well, no* | 338 | .2 |
| **Negative non-no answers** | *Uh, not a whole lot.* | 292 | .1 |
| **Signal-non-understanding** | *Excuse me?* | 288 | .1 |
| **Other answers** | *I don't know* | 279 | .1 |
| **Conventional-opening** | *How are you?* | 220 | .1 |
| **Or-Clause** | *or is it more of a company?* | 207 | .1 |
| **Dispreferred answers** | *Well, not so much that.* | 205 | .1 |
| **3rd-party-talk** | *My goodness, Diane, get down from there.* | 115 | .1 |
| **Offers, Options Commits** | *I'll have to check that out* | 109 | .1 |
| **Self-talk** | *What's the word I'm looking for* | 102 | .1 |
| **Downplayer** | *That's all right.* | 100 | .1 |
| **Maybe/Accept-part** | *Something like that* | 98 | <.1 |
| **Tag-Question** | *Right?* | 93 | <.1 |
| **Declarative Wh-Question** | *You are what kind of buff?* | 80 | <.1 |
| **Apology** | *I'm sorry.* | 76 | <.1 |
| **Thanking** | *Hey thanks a lot* | 67 | <.1 |

A puzzling, though insignificant, detail about SWBD-DAMSL is that everyone, including the authors of the coding manual (Jurafsky, Shriberg., & Biasca, 1997), refers to the tag set as containing 42 classifications, when in fact it contains 43. The most likely explanation for the discrepancy is that the Other tag isn't counted, but this is never explicitly stated.

### 1.3.2 Email help-desk study

One set of experiments in automatic sentence classification used a corpus of 160 help-desk email dialogs, totaling 1,486 sentences. (Khoo, Marom, & Albrecht, 2006) The researchers chose to use only the technicians' responses in each dialog, "as they contain well-formed grammatical sentences, as opposed to the customers' emails." (p. 19) The classification system was adapted from DAMSL:

| Khoo tag | Count | Percent |
|---|---|---|
| Statement | 423 | 28.5 |
| Thanking | 228 | 15.3 |
| Request | 146 | 9.8 |
| Salutation | 129 | 8.7 |
| Instruction | 126 | 8.5 |
| Instruction-item | 94 | 6.3 |
| URL | 80 | 5.4 |
| Response-ack | 63 | 4.2 |
| Suggestion | 55 | 3.7 |
| Specification | 41 | 2.8 |
| Signature | 32 | 2.2 |
| Question | 24 | 1.6 |
| Apology | 23 | 1.5 |
| Others | 22 | 1.5 |

### 1.3.3 Instant messaging help-desk study

Another researcher studied a similarly focused corpus, this one consisting of instant messages exchanged in a technical-support environment. The object was to develop an automatic system of identifying distinct utterances within complex messages. The tag set used was a subset of the 42 main SWBD-DAMSL tags. The distribution of the 550 utterances is shown below. (Ivanovic, 2005)

| Ivanovic tag | Example | Percent |
|---|---|---|
| Statement | I am sending you the page now | 36.0 |
| Thanking | Thank you for contacting us | 14.7 |
| Yes-no question | Did you receive the page? | 13.9 |
| Response-ack | Sure | 7.2 |
| Request | Please let me know how I can assist | 5.9 |
| Open-question | how do I use the international version? | 5.3 |
| Yes-answer | yes, yeah | 5.1 |
| Conventional-closing | Bye Bye | 2.9 |
| No-answer | no, nope | 2.5 |
| Conventional-opening | Hello Customer | 2.3 |
| Expressive | haha, :-), grr | 2.3 |
| Downplayer | my pleasure | 1.9 |

### 1.3.4 Chat studies

Two studies of chat dialogs used identical sets of tags for classifying utterances. Because of the similarities in tag set and in the type and size of corpus, it is possible and useful to compare the distributions reported.

The first study was aimed at developing a text-mining tool for analyzing interactions and discussion topics in Internet Relay Chat. The 3,129 postings were labeled with a tag set adapted mainly from two sources; the VerbMobil speech-based machine translation project (Accept, Bye, Clarify, Greet, Reject) and SWBD-DAMSL (Statement, Wh-question, Yes-answer, No-answer, Continuer, Other). The researchers created three tags specific to chat conversations: Emotion, Emphasis, and System. (Wu, Khan, Fisher, Shuler, & Pottenger, 2002) The other study, published five years later, applied this same tag set to 3,507 chat-room posts collected in 2006. The researchers' goal was to build a corpus tagged with lexical, syntactic, and discourse (post classification) information. Only the discourse data are presented here. (Forsyth & Martell, 2007)

As the following table shows, most of the percentages did not vary much between the two corpora. Other than Other (a difficult category to analyze), there are just a few categories with large discrepanices. Two, highlighted in the table, are of particular interest: In the 2002 study, people posted four times as many Agreement utterances as their counterparts did four years later; conversely, the 2006 chatters posted nearly four times as many Reject utterances. These discrepancies could reflect a decline in Internet civility, but they could also be due to differences between the two chatting populations or the topics being discussed (although no particular characteristics or topics were specified by the researchers).

Another interesting change is the three-fold increase in the use of emoticons. This could indicate that people have come to rely less on words to convey their emotions, or again it could be a function of the type of people and topics involved in each study.

| | *Wu et al. (2002)* | | *Forsyth & Martell (2007)* | |
| *Classification* | *Example* | *Percent* | *Example* | *Percent* |
| **Accept** | I agree | 10.0 | yeah it does, they all do | 2.5 |
| **Bye** | See you later | 3.6 | night ya'all. | 1.6 |
| **Clarify** | Wrong spelling | 0.3 | i meant to write the word may..... | 0.3 |
| **Continuer** | And … | 0.4 | and thought I'd share | 3.5 |
| **Emotion** | lol | 3.3 | lol | 11.5 |
| **Emphasis** | I do believe he is right. | 1.5 | Ok I'm gonna put it up ONE MORE TIME | 0.5 |
| **Greet** | Hi, Tom | 5.1 | hiya 10-19-40sUser43 hug | 13.4 |
| **No answer** | No, I'm not. | 0.9 | no I had a roomate who did though | 0.9 |
| **Other** | | 6.7 | 0 | 0.4 |
| **Reject** | I don't think so. | 0.6 | u r not on meds | 2.1 |
| **Statement** | I'll check after class | 42.5 | Yay...democrats have taken the house! | 34.5 |
| **System** | Tom [JADV@11.22.33.44] has left #sacba1 | 9.8 | JOIN | 17.0 |
| **Wh question** | Where are you? | 5.6 | why do you feel that way? | 5.3 |
| **Yes answer** | Yes, I am. | 1.7 | why yes I do | 1.2 |
| **Yes/no question** | Are you still there? | 8.0 | cant we all just get along | 5.2 |

### 1.3.5 Related work

**Effects of machine translation on collaborative work**

Yamashita and Ishida (2006) conducted experiments in which pairs of students whose only shared language was English worked cooperatively on a task. In half the trials, the students communicated in their native languages, using machine translation. In the other trials, they communicated in English. Each participant in the pair was presented with ten tangram figures; the figures were identical for each student, but they were arranged in two different sequences. The subjects' task was to match the arrangements of figures by communicating via a chat interface. The researchers identified six utterance types, according to the primary purpose of the utterance:

| Yamashita tag | Definition |
|---|---|
| Presentation (Description) | A speaker describing a figure: e.g., "Figure 7 looks like a bird flying to the left." "Its neck is long." |
| Presentation (Noun phrase) | A speaker explaining a figure with a noun phrase: e.g., "Figure 5 is a dancing lady." |
| Question or Confirmation | An addressee asking the speaker for clarification, more information, or confirming an understanding: e.g., "Is she wearing a long dress?" |
| Acceptance | An addressee accepting the speaker's presentation: e.g., "Ok," "That's my 5th figure." |
| Not Understood | An addressee telling the speaker that he/she did not understand the message (e.g., "I don't understand."). |
| Others | Utterances that don't belong to any of the above categories. |

The results showed that the subjects "use descriptions when first establishing a common perspective on a referent but shorten these descriptions to noun phrases once common ground is established." (p. 519) This knowledge could be useful when considering panlingual applications in which users would convey information about an image. The researchers did not indicate the proportion of utterances that consisted of either description or noun phrase, but they did compare the distribution of the other four utterance types when the subjects used English versus the use of native languages via machine translation.

| Yamashita tag | English 1st trial | English 2nd trial | MT 1st trial | MT 2nd trial |
|---|---|---|---|---|
| Questions/Confirmation | 25% | 17% | 24% | 20% |
| Acceptance | 66% | 83% | 65% | 78% |
| NotUnderstood | 3% | 0 | 5% | 2% |
| Others | 6% | 0 | 6% | 0 |

**Verbal Response Modes**

The tag set used by Lampert, Dale, and Paris (2006) was derived from Verbal Response Modes, an allegedly exhaustive coding system that was originally developed to analyze client-therapist interactions (Stiles, N.D.). The researchers augmented the corpus of VRM dialogues with additional sentences to ensure a good mix of sentence types, so the distribution is less informative than any of those presented above. What distinguishes it from all the other systems is that it has no "other" category, on the assumption that VRM accommodates all possible utterances.

| Lampert Tag | Count | Percent |
|---|---|---|
| **Disclosure** | 395 | 28.9 |
| **Edification** | 391 | 28.6 |
| **Question** | 218 | 15.9 |
| **Reflection** | 109 | 8.0 |
| **Acknowledgement** | 97 | 7.1 |
| **Advisement** | 73 | 5.3 |
| **Interpretation** | 64 | 4.7 |
| **Confirmation** | 21 | 1.5 |

## ICSI-MR

During the years 2000 to 2002, recordings were made of a series of meetings at the International Computer Science Institute in Berkeley. The tag set used for annotating the 72 hours of dialog was based on the SWBD-DAMSL tags, but had a number of additional, mostly domain-specific tags, such as Topic Change, Floor-holder, and Humorous Material. (Shriberg, Dhillon, Bhagat, Ang, & Carvey, 2004)

Each dialog act could have more than one tag, but each was required to be coded with exactly one general tag from a set of 11. Of those obligatory tags, the second, third, and fifth most frequent were **Backchannel, Floor-holder**, and **Floor-grabber.** These tags and their distribution seem only marginally useful in designing a system that would enable translingual written communication.

## 1.4 Summary

Given the conflicting systems in use, the different domains studied, and the small number of studies found, aggregation of the data presented in this literature review would be both problematic and ill-advised. However, a very rough distribution of utterance types might look like this:

| Utterance | Percent |
|---|---|
| **Statement** | 35 |
| **WH or open question** | 10 |
| **Greeting/opening** | 10 |
| **Agreement/acceptance** | 8 |
| **Emphasis/emotion** | 6 |
| **Yes-No question** | 5 |
| **Goodbye/closing** | 5 |
| **Clarification** | 4 |
| **Thanking** | 4 |
| **Yes answer** | 4 |
| **Rejection/denial** | 3 |
| **Wh-question** | 3 |
| **No answer** | 2 |
| **Apology** | 1 |
| Total | 100 |

The clarification, thanking, and apology categories have been inflated slightly on the assumption that people conversing translingually are more likely than monolingual conversers to be asking for clarification, and multilingual strangers are more likely than monolingual friends to express thanks and apology.

## 2 Predicate-Argument Structure

### 2.1 The Proposition Bank

Palmer, Gildea, & Kingsbury (2005) created the Proposition Bank by adding semantic-role labels to the syntactic structures in the *Wall Street Journal* portion of the Penn Treebank. A major goal of the project was to provide consistent semantic role labels across differing syntactic structures, as in "The window broke" vs. "He broke the window."

Each sense of a verb in the Proposition Bank has a frameset, consisting of all the arguments for that sense. The 4,500 PropBank framesets encompass about 3,300 verbs. Within each frameset the semantic arguments are numbered: Arg0 is generally the Agent, while Arg1 is a Patient or Theme. For the verb *kick*, for example, Arg0 is the kicker while Arg1 is the thing kicked. The higher-numbered arguments tend to be specific to the particular verb. The verb *edge,* for example, has these arguments in addition to Arg0 (causer of motion) and Arg1 (thing moved):

Arg2: distance moved      Arg4: end point
Arg3: start point      Arg5: direction

Additional constituent labels are subtypes of the ArgM tag:

LOC: location      CAU: cause
EXT: extent      TMP: time
DIS: discourse connectives      PNC: purpose
ADV: general-purpose .      MNR: manner
NEG: negation marker      DIR: direction
MOD: modal verb

The following two tables show the frequency statistics for argument roles in PropBank:

Most frequent semantic roles for each syntactic position (percentages)

| position | total | Four most common roles | | | | | | | | other |
|---|---|---|---|---|---|---|---|---|---|---|
| sub | 37364 | Arg0 | 79.0 | Arg1 | 16.8 | Arg2 | 2.4 | TMP | 1.2 | 0.6 |
| obj | 21610 | Arg1 | 84.0 | Arg2 | 9.8 | TMP | 4.6 | Arg3 | 0.8 | 0.8 |
| S | 10110 | Arg1 | 76.0 | ADV | 8.5 | Arg2 | 7.5 | PRP | 2.4 | 5.5 |
| NP | 7755 | Arg2 | 34.3 | Arg1 | 23.6 | Arg4 | 18.9 | Arg3 | 12.9 | 10.4 |
| ADVP | 5920 | TMP | 30.3 | MNR | 22.2 | DIS | 19.8 | ADV | 10.3 | 17.4 |
| MD | 4167 | MOD | 97.4 | ArgM | 2.3 | Arg1 | 0.2 | MNR | 0.0 | 0.0 |
| PP-in | 3134 | LOC | 46.6 | TMP | 35.3 | MNR | 4.6 | DIS | 3.4 | 10.1 |
| SBAR | 2671 | ADV | 36.0 | TMP | 30.4 | Arg1 | 16.8 | PRP | 7.6 | 9.2 |
| RB | 1320 | NEG | 91.4 | ArgM | 3.3 | DIS | 1.6 | DIR | 1.4 | 2.3 |
| PP-at | 824 | EXT | 34.7 | LOC | 27.4 | TMP | 23.2 | MNR | 6.1 | 8.6 |

Most frequent syntactic positions for each semantic role (percentages)

| roles | total | Four most common syntactic positions | | | | | | | | other |
|---|---|---|---|---|---|---|---|---|---|---|
| Arg1 | 35112 | obj | 51.7 | S | 21.9 | subj | 17.9 | NP | 5.2 | 3.4 |
| Arg0 | 30459 | subj | 96.9 | NP | 2.4 | S | 0.2 | obj | 0.2 | 0.2 |
| Arg2 | 7433 | NP | 35.7 | obj | 28.6 | subj | 12.1 | S | 10.2 | 13.4 |
| TMP | 6846 | ADVP | 26.2 | PP-in | 16.2 | obj | 14.6 | SBAR | 11.9 | 31.1 |
| MOD | 4102 | MD | 98.9 | ADVP | 0.8 | NN | 0.1 | RB | 0.0 | 0.1 |
| ADV | 3137 | SBAR | 30.6 | S | 27.4 | ADVP | 19.4 | PP-in | 3.1 | 19.5 |
| LOC | 2469 | PP-in | 59.1 | PP-on | 10.0 | PP-at | 9.2 | ADVP | 6.4 | 15.4 |
| MNR | 2429 | ADVP | 54.2 | PP-by | 9.6 | PP-with | 7.8 | PP-in | 5.9 | 22.5 |
| Arg3 | 1762 | NP | 56.7 | obj | 9.7 | subj | 8.9 | ADJP | 7.8 | 16.9 |
| DIS | 1689 | ADVP | 69.3 | CC | 10.6 | PP-in | 6.2 | PP-for | 5.4 | 8.5 |

The authors made the following observations:

> Arg0, when present, is almost always a syntactic subject, while the subject is Arg0 only 79% of the time. . . . Going from syntactic position to semantic role, the numbered arguments are more predictable than the non-predicate-specific adjunct roles. The two exceptions are the roles of "modal" (MOD) and "negative" (NEG), which . . . are almost always realized as auxiliary verbs and the single adverb (part of speech tag RB) *not,* respectively. [p. 18]

## 2.2 Earlier work

Using the FrameNet tag set and  database (about 50,000 sentences from the British National Corpus), Gildea & Jurafsky (2002) trained statistical classifiers to label the semantic roles of arguments. They provided the following frequency data on parse tree paths, i.e., the path from the target word (VB) to each constituent. Arrows in the path indicate whether the movement is upward or downward in the tree.
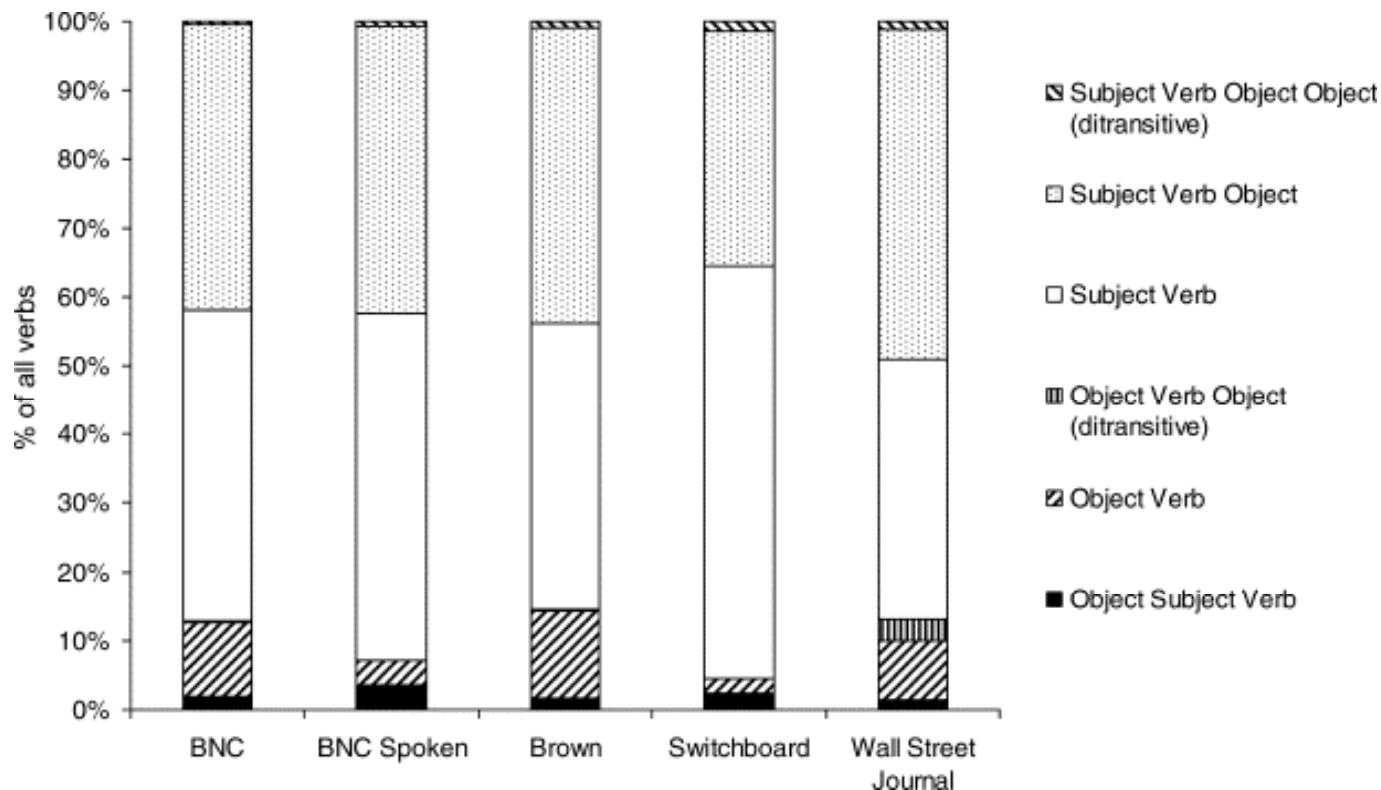
Most frequent values of the path feature in the training data.

| Frequency | Path | Description |
|---|---|---|
| 14.2% | VB↑VP↓PP | PP argument/adjunct |
| 11.8 | VB↑VP↑S↓NP | Subject |
| 10.1 | VB↑VP↓NP | Object |
| 7.9 | VB↑VP↑VP↑S↓NP | Subject (embedded VP) |
| 4.1 | VB↑VP↓ADVP | Adverbial adjunct |
| 3.0 | NN↑NP↑NP↓PP | Prepositional complement of noun |
| 1.7 | VB↑VP↓PRT | Adverbial particle |
| 1.6 | VB↑VP↑VP↑VP↑S↓NP | Subject (embedded VP) |
| 14.2 | | No matching parse constituent |
| 31.4 | Other | |

## 3 Word order

Roland, Dick, and Elman (2007) set out to gather accurate, comprehensive information on structural frequencies in English by analyzing the Brown, *Wall Street Journal,* and Switchboard corpora and both the spoken and written portions of the British National Corpus. One of their conclusions is relevant to any comparison of corpus statistics: They caution that cross-corpus comparisons reveal serious discrepancies in the structural frequencies encountered. Factors that can affect the distribution include discourse type, topics being discussed, and the degree of fluency of the speaker.

Even if accuracy were guaranteed, the distribution of verb subcategorizations, frequency of complementizer omission, and various other phenomena reported in this study do not seem very helpful in designing an interactive communication system. In the end, however, the researchers derived from their data a "global view" of the frequency of six word-order types (in English). The frequencies encountered in the five corpora are shown in the following graph.



These data reflect the fact that English is generally a Subject Verb Object (SVO) type of language. In the written corpora the proportion of Object Verb structures approached 15%, but spoken utterances exhibited closer to 5% frequency of OV order.

Word order could be a significant factor in planning for translingual communication. Over 95% of the world's languages use (in descending order of frequency) SOV, SVO, or VSO as their basic word order. There are a few VOS languages, including Malagasy, and even fewer languages that use OVS (e.g., Hixkaryana) or OSV (e.g., Apuriña). (O'Grady, Archibald, Aronoff, and Rees-Miller, 2001)

As these data illustrate, in the vast majority of English structures the subject appears before the verb. If this consistency is duplicated in other languages, it may be possible and desirable to label the words in a lemmatic/lexemic utterance and have them automatically re-ordered to fit the pattern of the target language.


## 4 Final words

Determining the true illocutionary force or meaning of a speech act requires a deeper analysis than was performed in most of the studies presented here. For the purpose of providing a communication interface, however, shallow knowledge may suffice. After all, in monolingual written exchanges the only information one normally has about the writer's meaning or intention is derived from the sentence type and structure. It is reasonable to expect that people engaging in panlingual communication will also depend on words and syntax alone to convey their meanings. Perhaps some of the data gathered here will prove useful in designing a system that will enable such communication.

# REFERENCES

Allen, J., & Core, M. (1997). Draft of *DAMSL: Dialog Act Markup in Several Layers.*
http://www.cs.rochester.edu/research/speech/damsl/RevisedManual/

Austin, J. L. (1962). *How to do things with words.* Cambridge: Harvard University Press.

Forsyth, E.N., & Martell, C.H. (2007). Lexical and Discourse Analysis of Online Chat Dialog. *International Conference on Semantic Computing* (ICSC 2007), pp. 19-26.

Gildea, D., & Jurafsky, D. (2002) Automatic labeling of semantic roles. *Computational Linguistics, 28*(3), 244-288. http://www.ldc.upenn.edu/acl/J/J02/J02-3001.pdf

Ivanovic, E. (2005) Automatic utterance segmentation in instant messaging dialogue. *Proceedings of the ACL. Student Research Workshop*, pp. 79–84.
http://www.alta.asn.au/events/altw2005/cdrom/pdf/ALTA200533.pdf

Jurafsky, D., Shriberg, L., & Biasca, D. (1997). *Switchboard SWBD-DAMSL shallow-discourse-function annotation coders manual.* http://www.stanford.edu/~jurafsky/ws97/manual.august1.html)

Khoo, A., Marom, Y., & Albrecht, D. (2006) Experiments with sentence classification. *Proceedings of the 2006 Australasian Language Technology Workshop (ALTW2006)*, pp. 18–25.
http://www.alta.asn.au/events/altw2006/proceedings/KhooEtAl.pdf

Lampert, A., Dale, R., & Paris, C. (2006) Classifying speech acts using Verbal Response Modes. *Australasian Language Technology Workshop 2006.*
http://www.ict.csiro.au/staff/Andrew.Lampert/writing/papers/SpeechActsVRM-ALTW2006-Lampert.pdf

O'Grady, W., Archibald, J., Aronoff, M., and Rees-Miller, J. (eds.) (2001). *Contemporary Linguistics: An Introduction* (4th ed). Boston: Bedford/St. Martin's Press.

Palmer, M., Gildea, D., & Kingsbury, P. (2005). The Proposition Bank: An annotated corpus of semantic roles. *Computational Linguistics 31*(1), 71-106. http://verbs.colorado.edu/~mpalmer/papers/prop.pdf

Roland, D., Dick, F., & Elman, J.L. (2007). Frequency of basic English grammatical structures: A corpus analysis. *Journal of Memory and Language, 57*(3), 348-379.

Shriberg, E., Dhillon, R., Bhagat, S., Ang, J., & Carvey, H. (2004) The ICSI Meeting Recorder Dialog Act (MRDA) corpus. *Proceedings of the 5th SIGdial Workshop on Discourse and Dialogue*, pp. 97-100.
http://acl.ldc.upenn.edu/W/W04/W04-2319.pdf

Stiles, W. B. (n.d.). *Verbal response modes coding system.*
http://www.users.muohio.edu/stileswb/verbal_response_modes.htm

Stolcke, A., K. Ries, N. Coccaro, E. Shriberg, R. Bates, D. Jurafsky, et al.(2000). Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational Linguistics, 26*(3), 341-373. http://www.stanford.edu/~jurafsky/ws97/CL-dialog.pdf

Wu, T., Khan, F.M., Fisher, T.A., Shuler, L.A., & Pottenger, W.M. (2002). Posting act tagging using transformation-based learning. *The Proceedings of the Workshop on Foundations of Data Mining and Discovery, IEEE International Conference on Data Mining (ICDM'02).*
http://dimacs.rutgers.edu/~billp/pubs/IEEEICDM02WorkshopTianhao.pdf)

Yamashita, N., & Ishida, T. (2006). Effects of machine translation on collaborative work. *Computer Supported Cooperative Work 2006,* pp. 515-523. http://www.lab7.kuis.kyoto-u.ac.jp/~ishida/pdf/cscw06.pdf