



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Speech Communication xxx (2005) xxx–xxx

SPEECH  
COMMUNICATION[www.elsevier.com/locate/specom](http://www.elsevier.com/locate/specom)

## Extraction of pragmatic and semantic salience from spontaneous spoken English

Tong Zhang \*, Mark Hasegawa-Johnson, Stephen E. Levinson

*Department of Electrical and Computer Engineering, Beckman Institute, University of Illinois at Urbana-Champaign,  
405 N. Mathews Avenue, Urbana, IL 61801, USA*

Received 1 January 2005; received in revised form 16 July 2005; accepted 19 July 2005

### 9 Abstract

10 This paper computationalizes two linguistic concepts, *contrast* and *focus*, for the extraction of pragmatic and seman-  
11 tic salience from spontaneous speech. Contrast and focus have been widely investigated in modern linguistics, as cat-  
12 egories that link intonation and information/discourse structure. This paper demonstrates the automatic tagging of  
13 contrast and focus for the purpose of robust spontaneous speech understanding in a tutorial dialogue system. In par-  
14 ticular, we propose two new transcription tasks, and demonstrate automatic replication of human labels in both tasks.  
15 First, we define *focus kernel* to represent those words that contain novel information neither presupposed by the inter-  
16 locutor nor contained in the precedent words of the utterance. We propose detecting the focus kernel based on a word  
17 dissimilarity measure, part-of-speech tagging, and prosodic measurements including duration, pitch, energy, and our  
18 proposed spectral balance cepstral coefficients. In order to measure the word dissimilarity, we test a linear combination  
19 of ontological and statistical dissimilarity measures previously published in the computational linguistics literature. Sec-  
20 ond, we propose identifying *symmetric contrast*, which consists of a set of words that are parallel or symmetric in lin-  
21 guistic structure but distinct or contrastive in meaning. The symmetric contrast identification is performed in a way  
22 similar to the focus kernel detection. The effectiveness of the proposed extraction of symmetric contrast and focus ker-  
23 nel has been tested on a *Wizard-of-Oz* corpus collected in the tutoring dialogue scenario. The corpus consists of 630  
24 non-single word/phrase utterances, containing approximately 5700 words and 48 min of speech. The tests used speech  
25 waveforms together with manual orthographic transcriptions, and yielded an accuracy of 83.8% for focus kernel detec-  
26 tion and 92.8% for symmetric contrast detection. Our tests also demonstrated that the spectral balance cepstral coef-  
27 ficients, the semantic dissimilarity measure, and part-of-speech played important roles in the symmetric contrast and  
28 focus kernel detections.

29 © 2005 Published by Elsevier B.V.

\* Corresponding author. Tel.: +1 217 328 1542; fax: +1 217 244 8371.

E-mail addresses: [tzhang1@ifp.uiuc.edu](mailto:tzhang1@ifp.uiuc.edu) (T. Zhang), [hasegawa@ifp.uiuc.edu](mailto:hasegawa@ifp.uiuc.edu) (M. Hasegawa-Johnson), [sel@ifp.uiuc.edu](mailto:sel@ifp.uiuc.edu) (S.E. Levinson).

30 *Keywords:* Spoken language understanding; Spoken dialogue systems; Computational linguistics; Information extraction

31

## 32 1. Introduction

33 Words are tools; in real speech, every word is  
34 deployed for the purpose of achieving a human  
35 goal. The fields of computational semantics and  
36 pragmatics study quantifiable goal variables—  
37 variables that encode quantifiable aspects of the  
38 goals served by a word in context—and their  
39 semantic and contextual correlates. This paper de-  
40 scribes the computation of two semantic and prag-  
41 matic goal variables, *focus* and *contrast*, from  
42 spontaneous speech.

43 The paper is organized as follows. The remain-  
44 der of Section 1 explains why we are interested in  
45 annotating focus and contrast, defines the aspects  
46 of focus and contrast that are under study with  
47 examples from an intelligent tutoring system  
48 (ITS) corpus, and then puts forward the objectives  
49 of our study in this paper. Section 2 provides some  
50 background in support of our work and describes  
51 related work in modern linguistics and computa-  
52 tional linguistics. Section 3 describes the ITS cor-  
53 pus in detail, with particular attention paid to  
54 annotations and corpus statistics of the proposed  
55 focus and contrast variables. Sections 4 and 5 de-  
56 scribe the algorithms implemented for the purpose  
57 of detecting the proposed focus and contrast vari-  
58 ables: Section 4 describes prosodic analysis, and  
59 Section 5 describes the measurement of word  
60 semantic similarities. Section 6 describes system  
61 integration and results of experimental evaluation  
62 using the ITS corpus. Section 7 discusses and con-  
63 cludes our work.

### 64 1.1. Motivation

65 The motivation of this study is to achieve ro-  
66 bust spontaneous spoken language understanding  
67 (SSLU) in an intelligent tutoring dialogue system.  
68 The system intends to provide a computer-based  
69 environment for education in math and physics,  
70 using the *Lego* construction set, with children of  
71 primary and early middle school ages (9–12 years  
72 old). Due to the characteristics of the dialogue sce-

nario as we describe in Section 3.1, the children  
users' spontaneous utterances are often dysfluent,  
ungrammatical, and even incoherent. Our robust  
speech understanding system design under these  
circumstances basically involves two steps: (1)  
Classification of each utterance into one of a list  
of 30 tutoring events. Similar to *call types* in an  
automatic call center or call router (Gorin et al.,  
2002; Chu-Carroll and Carpenter, 1999), the tutor-  
ing events are used to summarize the content  
meaning of utterances in the tutoring dialogue sce-  
nario in a broad way. For example, the tutoring  
event *AskForPlayInstruction* means that the user  
asks a question requesting the instruction on  
how to play the Legos; *SpinSpeed* means that the  
user is talking about the spinning speed of the  
Lego gears; and *ExplainAction* means that the user  
explains what is being done with the Legos. (2)  
Sometimes the tutoring event itself cannot provide  
sufficient information for the computer to pop up  
proper response. For example, when the tutoring  
event *ArithmeticComputation* is detected, some-  
times the computer needs to know what the type  
of the arithmetic computation is; if it is division,  
then the computer needs to know what the  
dividend and divisor are for proper response. Such  
detailed information needs syntactic/semantic  
structure parsing or named entity recognition  
(Zhang, 2004).

To analyze the content meaning of an utter-  
ance, we are interested in extracting a small set  
of words, from the utterance, that encode prag-  
matically and semantically salient information.  
We investigate the computerization of two linguis-  
tic concepts, focus and contrast, that are assumed  
to be useful for content summarization and struc-  
ture parsing of spoken messages. Both of the con-  
cepts have reasonably clear published definitions.  
We wish to adapt the published definitions as nec-  
essary in order to define a corpus transcription  
experiment, and to train and test algorithms that  
automatically detect these two categories of sal-  
ience based on cues measured in the speech wave-  
form and in its orthographic transcription.

## 117 1.2. Focus

118 The information structure of a sentence can be  
 119 partitioned into *presupposition* and *focus*: presup-  
 120 position is what the interlocutor assumes to be true  
 121 when the sentence is elicited in a conversation, and  
 122 focus is the non-presupposed part of the sentence  
 123 (Chomsky, 1971; Zubizarreta, 1998). For example  
 124 (T represents tutor and U represents user. Focus is  
 125 marked by [ ]<sub>F</sub>),

- 126 (1) T: *What are you exploring there?*  
 U: [*Seeing if the small gears move the big*  
 129 *gears.*]<sub>F</sub>  
 130 (2) T: *How many times does the small gear spin*  
 131 *until they line up again?*  
 U: *I think it goes around [one and a half]<sub>F</sub>*  
 135 *times.*

134  
 137 By definition focus is indicative of pragmati-  
 138 cally and semantically new information not pre-  
 139 supposed by the interlocutor. If focus can be  
 140 reliably detected, it should be possible to use the  
 141 distinction between focus and presupposition to  
 142 detect new information embedded in an utterance.  
 143 Speakers will often signal focus of a sentence by  
 144 the use of *pitch accent* (we use pitch accent to mean  
 145 prosodic prominence marked by F0 extrusion; the  
 146 same word is usually also marked by the other  
 147 acoustic correlates of prominence, including dura-  
 148 tion, energy, and spectral balance). Pitch accent  
 149 marks the constituents within an utterance as high-  
 150 lighted or unexpected; it has been argued that con-  
 151 stituents outside focus are expected, and hence  
 152 tend to be unaccented (Kadmon, 2001; Zubizar-  
 153 reta, 1998; Hedberg and Sosa, 2001). For example  
 154 (pitch accented words are marked with subscript  
 155 *a*),

- 156 (3) T: *Which gear are you counting?*  
 158 U: *I am counting the small<sub>a</sub> gear.*  
 159 (4) T: *Which gear do you think is the strongest?*  
 162 U: *Probably the large<sub>a</sub> gear.*

163 However, the phonological manifestation is not  
 164 straightforward: pitch accent can only approxi-  
 165 mate the location of focus in a sentence. For exam-  
 166 ple, in the sentence *They turn in the opposite<sub>a</sub>*

*direction*, the accented word ‘opposite’ is focus 167  
 for the question *What can you tell me about the* 168  
*directions they turn* but not for the question *What* 169  
*else do you notice?* The latter question requires the 170  
 sentence ‘they turn in the opposite direction’ to be 171  
 focus for interpretation. Such ambiguity in prag- 172  
 matic interpretation of single accent has been 173  
 known traditionally as the *focus projection* phe- 174  
 nomenon, demonstrating that focus expressed by 175  
 a single accent can project to a larger linguistic 176  
 constituent than just the word with pitch accent. 177

Since focus is a syntactic constituent, the 178  
 boundaries of focus need to be determined to iden- 179  
 tify focus. A sentence may have multiple foci, and 180  
 the size of a focus may vary from a single word to 181  
 a phrase or even a sentence. It is difficult to auto- 182  
 matically extract syntactic constituents containing 183  
 novel information without making use of a com- 184  
 plete parse tree for the sentence in question. Even 185  
 with a parse tree available, automatically selecting 186  
 the right constituents would be difficult; for exam- 187  
 ple, in the following exchange, 188

- (5) T: *Which gear are you counting?* 189  
 U: *I am counting the [small]<sub>F</sub> gear [in my* 190  
*hand]<sub>F</sub>.* 193

it would be difficult for an automatic algorithm to 194  
 determine that focus consists of a single word and 195  
 a prepositional phrase; it would be nearly impossi- 196  
 ble without access to a correct parse of the sen- 197  
 tence. It is even harder to extract focus from 198  
 spontaneous speech, since spontaneous speech 199  
 often has loose grammar structure, dysfluency, 200  
 and inconsistency between linguistic segments 201  
 and acoustic segments. For example (‘...’ repre- 202  
 sents silence): 203

- (6) T: *What happens when you spin the left gears?* 204  
 U: *Ahhh ... When you [after it goes around* 205  
*once the other one goes around ... the same* 206  
*... the same I mean it goes around ... you* 207  
*know you only have to spin it around once* 208  
*... and that makes sense basically because* 209  
*they are the same size.]<sub>F</sub>* 212

To robustly understand spontaneous speech, we 213  
 propose labeling individual words containing new 214  
 215

216 information neither presupposed by the interlocu-  
 217 tor nor contained in the preceding part of the  
 218 utterance. Such a word is usually a content word  
 219 because of the information content requirement.  
 220 We hypothesize that words matching this defini-  
 221 tion will typically be the semantically salient part  
 222 of focus. Therefore, we call each of these words a  
 223 *focus kernel*. In the following examples, focus ker-  
 224 nels are marked with **bold**:

(7) *T: What happens to the different gears as you  
 spin the one at the end?*

*U: They **move** with the single gear that I'm  
 spinning.*

(8) *T: Oh, are you having fun?*

*U: **Yeah**, it's kind of interesting.*

(9) *T: How many times would it take for the reds  
 to come back on top?*

*U: It would take **three** times to have the red be  
 back on top.*

238

### 239 1.3. Contrast

240 Contrast is a concept having multiple senses: (1)  
 241 In logic, two propositions are defined to be con-  
 242 trastive if it is impossible for them to be true simul-  
 243 taneously. For example, in the sentence *Bach was*  
 244 *an organ mechanic; Mozart knew little about or-*  
 245 *gans*, the two propositions are not contrastive,  
 246 whereas they become contrastive when 'Mozart'  
 247 is replaced by 'Bach' at the beginning of the second  
 248 sentence (Bosch and van der Sandt, 1999). (2) The  
 249 discourse relation called contrast is induced by  
 250 'but,' and constitutes a pair (or pairs) of contrasted  
 251 alternatives, which can be predicates (e.g., *John*  
 252 *cleaned up the room, but he didn't wash the dishes*),  
 253 individual words (e.g., *John cleaned up the room,*  
 254 *but Bill didn't*), or propositions (e.g., *It is raining,*  
 255 *but we go out for a walk*) (Umbach, 2004). (3) Some  
 256 linguists use contrast to denote the mutually exclu-  
 257 sive disjunction between the words contributing to  
 258 a fact and other alternatives made available by  
 259 context (Vallduví and Vilkuna, 1998). It has been  
 260 argued that focus in general establishes a contrast  
 261 since novel information usually conveys contrast  
 262 between a fact and the potential alternatives  
 263 (Bolinger, 1961; Kruijff-Korbayova and Steedman,  
 264 2003). For example, in the sentence *Last night they*

*had a party*, there is a contrast between the focus  
 'party' and any other alternative activities of the  
 group. (4) Symmetric contrast consists of a set of  
 words that are parallel or symmetric in linguistic  
 structure but mutually exclusive in meaning; the  
 stress on one word is motivated by its distinction  
 from the others, e.g., 'American' and 'Canadian'  
 in *An American farmer was talking to a Canadian*  
*farmer* (Rooth, 1992; Umbach, 2004).

In this study, we seek to make use of the knowl-  
 edge about contrast from the pragmatics and pros-  
 ody literature, for the purpose of detecting pairs of  
 symmetrically contrasted words that are assumed  
 to be useful for spontaneous speech understand-  
 ing. Symmetric contrast can occur within a sen-  
 tence, e.g. (contrasted words are marked with  
**bold**),

(10) *U: The **large** gear has five times as many teeth  
 as the **small** ones.*

(11) *U: How about **small** and **big** and **medium**?*

Topics and/or foci of conjunct phrases or coordi-  
 nated sentences (by 'and', 'but', etc.) can also con-  
 stitute symmetric contrast, e.g.

(12) *T: Where are the gears?*

*U: The **red** gear is on the **bottom** and the **yel-**  
**low** gear is on the **top**.*

(13) *T: How are the gears spinning?*

*U: The two **outside** ones spin in the **same**  
 direction and the **middle** one spins in the **oppo-**  
**site** direction.*

The words participating in a symmetric contrast  
 satisfy semantic parallelism, which has two impli-  
 cations: (a) the conjunct alternatives have to be  
 semantically independent of each other in the sense  
 that neither of them subsumes the other; and (b)  
 there has to be a "common integrator," i.e., a con-  
 cept subsuming both conjunct alternatives  
 (Umbach, 2004).

### 1.4. Study objective

As its primary technical goal, the study intends  
 to test whether the proposed word tags, i.e., focus  
 kernel and symmetric contrast, can be reliably



310 annotated in a spontaneous speech corpus using  
311 both manual and automatic annotation. As part  
312 of this evaluation, this study tests the relationship  
313 of focus kernel and symmetric contrast with the  
314 following prosodic and pragmatic variables: (1)  
315 prosodic prominence—experiments described in  
316 this paper test the reliability of prosodic promi-  
317 nence in the automatic identification of focus ker-  
318 nel and symmetric contrast; (2) novelty and  
319 semantic parallelism, the semantic attributes of fo-  
320 cus kernel and symmetric contrast. Information  
321 theoretic measures of novelty and semantic paral-  
322 lelism are implemented, based on algorithms pro-  
323 posed in the computational linguistics literature.  
324 Implemented algorithms are tested for the purpose  
325 of automatically identifying focus kernel and sym-  
326 metric contrast; and (3) part-of-speech. In addi-  
327 tion, this study discusses the usefulness of focus  
328 kernel and symmetric contrast to spontaneous  
329 speech understanding.

## 330 2. Background and related work

331 Focus and contrast in modern linguistics are  
332 used to “account for the correlation between cer-  
333 tain prosodic patterns and certain pragmatic and  
334 semantic effects” (Kadmon, 2001). Sections 2.1  
335 and 2.2 describe related work on contrast and fo-  
336 cus published in the linguistics and computational  
337 linguistics literature. Section 2.3 describes previous  
338 work on the word dissimilarity measure (given a  
339 pair of words, how much novel information a  
340 word contains with respect to the other word) in  
341 natural language processing (NLP).

### 342 2.1. Focus

343 The information structure of sentences can be  
344 defined in various ways, e.g., presupposition-focus  
345 (Chomsky, 1971; Jackendoff, 1972; Zubizarreta,  
346 1998), topic-comment (Dahl, 1969), theme–rheme  
347 (Firbas, 1964, 1966; Bolinger, 1965; Steedman,  
348 2000), given-new (Halliday, 1967; Kay, 1975),  
349 background-focus (Dahl, 1969; Steedman, 2000),  
350 and background-kontrast (Vallduví and Vilkuna,  
351 1998). Although the information structure defini-  
352 tions are diverse, overlapping and even conflicting,

they can be categorized into two dimensions: (1) 353  
topic-comment or theme–rheme, in which one part 354  
describes what the discourse is talking about, and 355  
the other part advances the discourse; and (2) gi- 356  
ven-new or background-kontrast, in which one 357  
part conveys information that has been known, 358  
and the other part conveys novel information dis- 359  
tinguishing actual occurrence from potential alter- 360  
natives triggered by context (Kruijff-Korbayova 361  
and Steedman, 2003; Lee, 2003). In defining the 362  
information structure of sentences in the given- 363  
new or background-kontrast dimension, we follow 364  
the partition of Chomsky (1971), Jackendoff 365  
(1972) and Zubizarreta (1998), i.e., presupposition 366  
and focus (see Section 1.2). 367

Contrastive focus is a special kind of focus that 368  
expresses *exhaustive identification* of an element 369  
(or subset) given a set of candidates. Exhaustive 370  
identification means that the selection of a candi- 371  
date excludes all other candidates. Ordinary focus 372  
introduces new, non-presupposed information. 373  
Contrastive focus exhaustively selects one or more 374  
candidates from a set of candidates that are pre- 375  
supposed by the interlocutor (Umbach, 2004; 376  
Lee, 2003; Hedberg and Sosa, 2001). For example 377  
*It was [John]<sub>CF</sub> that baked bread for our breakfast.* 378  
When the answer to an alternative question is to 379  
choose a disjunct from the disjunctive alternatives, 380  
the choice is thought to be exhaustive. For exam- 381  
ple, the choice of ‘money’ or ‘pen’ for *Did the baby* 382  
*pick the money first, or did she pick the pen first?* is 383  
contrastive focus (Lee, 2003). Another case of 384  
exhaustiveness is *correction* in dialogues: contra- 385  
stive focus corrects the explicit or implicit assump- 386  
tion made by the interlocutor. 387

Much of the literature studying focus seems to 388  
be motivated by the pitch accent correlate of focus 389  
(e.g., Hedberg and Sosa, 2001; Pierrehumbert and 390  
Hirschberg, 1990): semantic or pragmatic non-pre- 391  
supposed new information is usually signaled by a 392  
language-dependent accentual *F0* contour. Human 393  
detection of focus is very sensitive to the presence 394  
or absence of accenting of the focused item and the 395  
deaccenting of the post-focus items (Welby, 2003; 396  
Xu et al., 2004; Gussenhoven, 2002). Nuclear pitch 397  
accent (the last pitch accent in an intermediate 398  
phrase) and prenuclear pitch accent have different 399  
effects on the listener’s interpretation of focus 400

401 (Welby, 2003). However, as we have described,  
 402 there is no decisive information source to clearly  
 403 mark the focused syntactic constituent. For this  
 404 reason, automatic detection of focus has received  
 405 little study. In one relevant study, Heldner et al.  
 406 (1999) tried to locate *narrow focus* (focus consist-  
 407 ing of a single accented word) in Swedish speech:  
 408 they used energy and high frequency emphasis to  
 409 automatically detect narrow focus within three-  
 410 word phrases. Their automatic focus detector  
 411 was designed based on the assumption that there  
 412 was only one focused word per phrase and the fo-  
 413 cused word was accented. About two thirds of the  
 414 focused words were correctly detected.

415 In general, pitch accent in spoken English  
 416 marks subconstituents within an utterance as high-  
 417 lighted. A central reason for an item to be high-  
 418 lighted is novelty. Focus is marked by pitch  
 419 accent because it expresses novel information.  
 420 However, pitch accent may also signal something  
 421 else; there are a variety of factors affecting the  
 422 placement of pitch accent. As Kadmon (2001)  
 423 pointed, pitch accent can be assigned to a given  
 424 item with special importance, while novelty may  
 425 not be marked with pitch accent. Therefore, pitch  
 426 accent is an important condition, but neither a  
 427 necessary nor a sufficient condition for novelty.

## 428 2.2. Contrast

429 Contrast is widely investigated in modern lin-  
 430 guistics with special attention given to *contrastive*  
 431 *topic* and *contrastive focus* (the latter has been ad-  
 432 dressed in Section 2.1). Contrastive topic can be  
 433 defined as a syntactic constituent that is both to-  
 434 pic-marked and contains a focused item. That  
 435 means, the constituent on one hand forms a topic,  
 436 and on the other hand contains novel information  
 437 in contrast with other alternatives triggered by the  
 438 context (Lee, 1999). In the following example

- 439 (14) *T: What happens to the gears?*  
 440 *U: [The large gear]<sub>CT</sub> only spins once and it*  
 441 *spins slower,*

442  
 443 ‘gears’ is the topic that the speaker is addressing,  
 444 but the choice of ‘large’ is new and in contrast with  
 445 ‘the other gears.’ In addition, the topics of a

sequence of independent contrastive answers to a 447  
 question also form contrastive topics (Krifka, 448  
 1999), e.g. 449

- (15) *T: What about the large gear and the medium* 450  
*gear?* 451  
*U: The large<sub>CT</sub> gear spins left, and the med-*  
*ium<sub>CT</sub> gear spins right.* 455

In (15), ‘large’ and ‘medium’ are identified to be 454  
 contrastive topic and symmetric contrast 456  
 simultaneously. 457  
 458

In the literature, much investigation of contrast 459  
 has focused on the conceptual issue of some prob- 460  
 lematic cases and specification of the types of pitch 461  
 accent correlate (e.g., Gundel and Fretheim, 2001; 462  
 Pierrehumbert and Hirschberg, 1990; Kadmon, 463  
 2001; Lee, 1999). In general, contrast is well and 464  
 clearly defined in linguistics. However, a small 465  
 set of problematic cases receive a great deal of 466  
 attention. For example, authors are divided on 467  
 whether to identify a negated presupposition as 468  
 contrastive focus or contrastive topic (e.g., the 469  
 words ‘anything extraordinary’ in *He did the only* 470  
*thing you could do. He hasn’t done anything* 471  
*extraordinary*). In addition, published studies 472  
 agree that contrast is typically marked with pitch 473  
 accent, but there is considerable disagreement on 474  
 the particular patterns of pitch accent used to 475  
 mark contrastive topic and/or contrastive focus 476  
 (e.g., Lee, 1999; Hedberg and Sosa, 2001). These 477  
 debates are usually framed in terms of the pitch ac- 478  
 cent categories specified by the *ToBI* (tones and 479  
 break indices) prosodic annotation standard 480  
 (Beckman and Ayers, 1994), e.g., whether contras- 481  
 tive topic will be marked by H\*+!H or not. How- 482  
 ever, the disputes in the literature generally do not 483  
 affect the integrity of our task, because: (1) the 484  
 problematic cases requiring conceptual discrimina- 485  
 tion are rare in our corpus; and (2) our study does 486  
 not require the discrimination of pitch accent 487  
 patterns. 488

## 489 2.3. Word semantic similarity

The computational modeling of novelty and 490  
 semantic parallelism involves semantic analysis of 491  
 words, more exactly, the comparison between 492

493 words in terms of semantic meaning. Strictly  
 494 speaking, every word is different from every other  
 495 word. “Absolute *synonymy*, if it exists, is quite  
 496 rare. Usually, words that are close in meaning  
 497 are almost synonyms, but not quite; very similar,  
 498 but not identical, in meaning; not fully intersubsti-  
 499 tutable, but instead varying in their shade of deno-  
 500 tation, connotation, or emphasis” (Edmonds and  
 501 Hirst, 2002). So how can we formally specify simi-  
 502 larity and dissimilarity? A word often has multiple  
 503 senses. What counts as a central trait to compare  
 504 the meaning of a pair of words? Researchers in  
 505 computational linguistics have developed various  
 506 measures to compute the degree of *semantic simi-*  
 507 *larity* between a pair of words. The measures are  
 508 basically divided into two categories: (1) *ontology*  
 509 *hierarchies* (e.g., Lee et al., 1993; Sussna, 1993;  
 510 Resnik, 1995; Jiang and Conrath, 1997). Ontology  
 511 is a structural system of categories or semantic  
 512 types, so that knowledge about a certain domain  
 513 can be organized through the categorization of  
 514 the entities of the domain in terms of the types in  
 515 the ontology. The length of the path between a  
 516 pair of words is a measure of their semantic dis-  
 517 similarity. The well-known *edge-based* method re-  
 518 flects the fact that in a hierarchical semantic  
 519 network, the simplest measure of the distance be-  
 520 tween two elemental concept nodes, *A* and *B*, is  
 521 the shortest path that links *A* and *B*, i.e., the min-  
 522 imum number of edges that separate *A* and *B*  
 523 (Rada et al., 1989). (2) *Corpus statistics* empirically  
 524 model the context-dependence characteristics of  
 525 word meaning in text (e.g., Lin, 1998; Pantel and  
 526 Lin, 2002; Thelen and Riloff, 2002; Terra and  
 527 Clarke, 2003). Three statistics are commonly em-  
 528 ployed to model the similarity of words (Higgins,  
 529 2004):

*Topicality assumption*: similar words tend to have the same neighboring content words.

*Proximity assumption*: similar words tend to occur near each other. Word senses are ultimately grouped according to proximity of meaning.

*Parallelism assumption*: similar words tend to be found in similar grammatical structures.

### 3. Corpus description, annotations and analyses 539

#### 3.1. Tutoring dialogue scenario 540

The intelligent tutoring system helps students 541  
 learn basic math and physics concepts by playing 542  
 with Lego gears, with the objective of helping stu- 543  
 dents develop a physical understanding of abstract 544  
 concepts. For example, one question about the 545  
 relationship between gear size and spinning speed 546  
 is *Line up a 24-tooth gear and a 40-tooth gear. If 547*  
*the 24-tooth gear spins 5 times, then how many 548*  
*times must the 40-tooth gear spin for them to line 549*  
*up again? Why?* Children can answer this question 550  
 by spinning the gears and counting the cycles. Sim- 551  
 ilarly, a physics question about interactive force is 552  
*Put one hand on the 40-tooth gear axle, and put an- 553*  
*other hand on the 8-tooth gear axle. What happens 554*  
*if you hold one of them steady, and try to turn the 555*  
*other one? Why?* Children usually think that the 556  
 big gear is stronger before they do the experiment. 557  
 However, it turns out that the small gear is 558  
 stronger. 559

The complete system has not been finished yet; 560  
 the database used in this study was collected by 561  
 Wizard-of-Oz simulations of the finished system. 562  
 In the experiments, the user and the tutor (human 563  
 wizard) were sitting in separate rooms. The user 564  
 orally communicated with a computerized talking 565  
 head shown on the computer screen ahead of 566  
 him using a head-set microphone. The lip move- 567  
 ment of the talking head was coincident with 568  
 speech synthesized from text that was typed by 569  
 the tutor. The user’s speech was transmitted 570  
 through the microphone to a digital camera placed 571  
 opposite the user, and then was transmitted to the 572  
 earphone of the tutor sitting in another room. 573  
 Both the tutor and the user had Lego gearsets on 574  
 the table in front of them. A video of the tutor’s 575  
 gearset was displayed on the user’s computer, 576  
 and vice versa. The WoZ experiments allowed data 577  
 to be collected that was similar in most respects to 578  
 the data the final spoken language system would 579  
 need to understand: because children felt that they 580  
 were communicating with a computer tutor in- 581  
 stead of a human tutor, they behaved as they 582

583 would in a real computer-interacting environment.  
 584 The number of experimental sessions in which  
 585 each child participated varied from one to three,  
 586 depending on the interest and cooperative attitude  
 587 of the child. The tutor adjusted the content of  
 588 experiments according to the intelligence, coopera-  
 589 tion, and learning progress of each child subject.

### 590 3.2. Characteristics of the ITS corpus

591 Unlike the users of a telephone weather or flight  
 592 ticket system, who are often expert users interested  
 593 in achieving known goals using known tools in the  
 594 shortest time possible, the users of our intelligent  
 595 tutoring system are perpetually naïve with respect  
 596 to the future content of the dialogue. Each child  
 597 participates in at most three experimental sessions  
 598 and each session has different tutorial content, be-  
 599 cause we do not ask children to re-learn knowledge  
 600 that they have mastered. Therefore, although chil-  
 601 dren who participate in more than one experiment  
 602 may gain some expertise in the use of the computer  
 603 interface, the children are never able to memorize  
 604 menu prompts, re-use conversation content, or  
 605 otherwise become expert users of the dialogue  
 606 system.

607 In addition, we encourage children to partici-  
 608 pate in the experiments, and we instigate their  
 609 interest in scientific learning by asking them  
 610 open-ended questions rather than questions with  
 611 an absolute answer. For example, when the child  
 612 is turning gears and the tutor wants to ask the  
 613 child about the motion of the gears, he usually  
 614 does not ask single-choice questions such as *In*  
 615 *which direction are the gears turning?* Instead, he  
 616 would ask questions whose answers are not abso-  
 617 lute, e.g., *What are you noticing?* Compared with  
 618 the single-choice questions, the open-ended ques-  
 619 tions open a wider space for children and arouse  
 620 children's enthusiasm to use their imagination,  
 621 knowledge, and observation to solve problems.

622 Since children are not familiar with the experi-  
 623 ment contents and the answers to open-ended  
 624 questions are usually longer and more complicated  
 625 than those to single-choice questions, their utter-  
 626 ances are even more incoherent and dysfluent than  
 627 is typical in interpersonal conversations. The utter-  
 628 ances usually include loose grammar structure,

fragments, restarts, repairs, meaningless speech  
 (e.g., *That if the...*), and repetitions. Moreover,  
 sentence boundaries in spontaneous speech are  
 ambiguous because of mismatch between acoustic  
 segmentation and linguistic segmentation. The fol-  
 lowing example illustrates the characteristics of the  
 ITS speech data:

- (16) U: *Big gears move in different ways and... 636*  
*uhm...with the first when you push one of 637*  
*the first gears, the other gear the last gear 638*  
*moves you know, and the gear after that 639*  
*moves. 640*

641

### 3.3. Focus kernel annotation 642

Three annotators worked independently on  
 identification of focus kernel, if any, in each utter-  
 ance of the ITS corpus. Annotation was based on  
 perception, text transcriptions, and dialogue con-  
 text. A vast majority of the ITS corpus are ques-  
 tion-answer pairs between the tutor and user, in  
 which the tutor initiates local dialogue topics by  
 asking questions or providing suggestions (e.g.,  
*Can you make them look like this?*). In this case,  
 presupposition of an utterance lies in the question  
 or suggestion of the tutor (Kadmon, 2001). Stu-  
 dents sometimes initiate local dialogue topics by  
 issuing commands, asking questions or simply  
 explaining what they are doing. In this case, pre-  
 suppositions for the students' utterances do not ex-  
 ist. Annotators were given the following criteria to  
 use in their labeling of focus kernels in an  
 utterance:

1. Mark the content words that contain informa- 661  
 tion not already available in presupposition, if 662  
 any, nor in the preceding words of the utter- 663  
 ance. For example, 664  
 665

- (17) T: *What if we now add a gear?*  
 U: *The **third** one **moves** or **spins** the **same** 6*  
***direction** that the **first** one does.*

670

Mark contrastive focus: (a) If the tutor asks an 671  
 alternative disjunctive question and the user 672  
 responds by choosing a single disjunct, then 673



674 the disjunct is characterized by contrastive  
675 focus, and thus, is marked as focus kernel.  
676 For example,

(18) *T: Does the top or the bottom do the  
pushing?*

682 *U: I think it's the top.*

681 (b) If the user's utterance is to correct the  
683 assumption of the tutor, then the corrective  
word is characterized by contrastive focus,  
and thus, is marked as focus kernel. For  
example,

69(19) *T: Take that big gear, please.*

*U: I thought you said **this** gear.*

- 692  
693 2. Focus defined by Chomsky (1971) and Zubizar-  
694 reta (1998) is confined to the rhyme of a sen-  
695 tence (Steedman, 2000). However, we are also  
696 interested in marking the new item in contras-  
697 tive topic, e.g., 'large' in example dialogue  
698 (14), since it conveys novel information. We  
699 label such item as focus kernel although it is  
700 not in the rhyme, because: in our corpus this  
701 case is very rare and focus kernel are generally  
702 in the focus (rhyme), so the inclusion of this  
703 rare case basically does not affect the fact that  
704 we can call the words containing novel informa-  
705 tion as the kernel of focus.
- 706 3. In the case of dysfluency, discourse marker and  
707 repetition should not be marked, while repair  
708 should be marked. For example, in the dysflu-  
709 ent constituent [ ]<sub>dys</sub> of the utterance *If I spin*  
710 *3 times it end up, [all yellows at the bottom,*  
711 *and all yellows at the... I mean, all reds at the*  
712 *bottom,]<sub>dys</sub> *and all yellows at the top,* the dis-  
713 course marker 'I mean' should not be marked.  
714 As correction of 'yellows', 'red' should be  
715 marked.*
- 716 4. Function words are not marked unless they  
717 carry novel information, and are thus stressed  
718 by the speaker, e.g., 'this' in example (19).

719 Each word in an utterance is labeled either fo-  
720 cus kernel or nonfocus kernel. We used the kappa  
721 statistics to evaluate consistency among the three  
722 annotators. The kappa statistics is a chance-cor-

724 rected inter-transcriber agreement rate  
 $\kappa = (P_O - P_C) / (1 - P_C)$ , where  $P_O$  is the rate of in-  
725 ter-transcriber agreement, and  $P_C$  is the average  
726 rate that would be achieved by chance (Flammia,  
727 1998). Comparison of the three transcriptions  
728 yielded a kappa score of 0.79, indicating reason-  
729 ably good agreement among the transcribers. We  
730 used majority voting to resolve the annotation dif-  
731 ference among the three annotators. That is, for  
732 each word, if two or more annotators had the  
733 same focus kernel/nonfocus kernel label, then we  
734 assigned that label to the word as the target label.  
735

### 3.4. Symmetric contrast annotation 736

737 The labeling of symmetric contrast was based  
738 on pairs of words. For example, 'small' and 'big'  
739 are labeled to be a symmetric contrast. If more  
740 than two words were symmetrically contrastive  
741 with one another in an utterance, then we labeled  
742 them by pairs. For example, in *I'm counting the*  
743 *small, medium, and big*, we labeled 'small', 'med-  
744 ium' and 'big' by three symmetric contrasts: 'small'  
745 and 'medium', 'small' and 'big', and 'medium' and  
746 'big'. Moreover, in the case of repair (dysfluency),  
747 e.g., *The small one, no, the big one is turning left*,  
748 the repairing word (e.g., 'big') and the repaired  
749 word (e.g., 'small') were also labelled as a symmet-  
750 ric contrast.

751 Two transcribers worked together to identify  
752 instances of symmetric contrast in the ITS corpus.  
753 The annotation was performed based on speech  
754 perception, transcription, and dialogue context.  
755 Because the transcribers worked together, no in-  
756 ter-transcriber agreement statistics were  
757 computed.

### 3.5. Corpus analysis 758

759 To date 29 WoZ experiments with 17 subjects  
760 have been carried out and transcribed, and we  
761 have collected 11.7 h of audio-visual data. In each  
762 experiment, the child often spent most time silently  
763 playing with the Legos. Therefore, the collected  
764 audio data corpus was small. Recordings of the  
765 students were manually transcribed and segmented  
766 into conversation sides. Each of the user's conver-  
767 sation sides was considered an *utterance*. Some

768 speech data had to be discarded because of Lego  
 769 block noise, heavy breathing, etc. This process re-  
 770 sulted in a total of 714 transcribed user utterances,  
 771 containing a total of approximately 50 minutes of  
 772 relatively clean speech. On average each utterance  
 773 had 4.2 s speech and 8.1 words. The vast majority  
 774 of the utterances contained 1–20 words, while the  
 775 longest utterance had 57 words. Fig. 1 shows the  
 776 histogram of utterance lengths in the corpus.

777 The 714 utterances of the ITS corpus can be  
 778 partitioned into 58 single-word utterances, 26 sin-  
 779 gle-phrase utterances (such as *how many?*), 22  
 780 utterances that merely repeat the words of the tu-  
 781 tor, 19 utterances that are not semantically mean-  
 782 ingful (e.g., *That if the...*), and 589 multi-word,  
 783 multi-phrase utterances containing meaningful  
 784 and novel information. We use the 630 multi-  
 785 word, multi-phrase utterances (589 meaningful  
 786 utterances, 22 repetitions of the tutor, and 19  
 787 semantically void utterances) for experiments.

788 Short utterances tend not to contain symmetric  
 789 contrast. We choose the average number of words  
 790 in a sentence, which is 8, as the threshold to distin-  
 791 guish long utterances from short utterances. Utter-  
 792 ances containing at least 8 words occupy 40% of  
 793 the 630-utterance corpus. Table 1 shows the rela-  
 794 tionship between utterance length and occurrences  
 795 of symmetric contrast. The table shows that 45%  
 796 of the long utterances contain symmetric contrast,  
 797 while 5% of the short utterances contain symmet-  
 798 ric contrast. Phrased in another way, more than  
 799 85% of symmetric contrast instances occur in  
 800 utterances of 8 words or more. Fig. 2 shows the  
 801 histogram of focus kernels per utterance in the

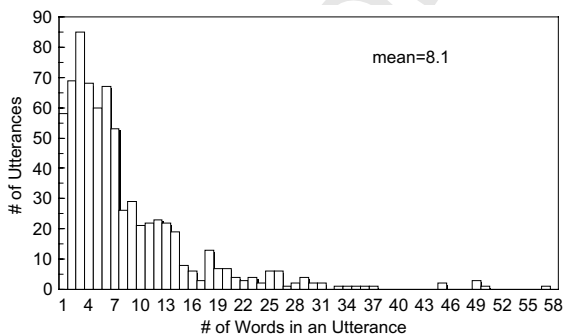


Fig. 1. Histogram of utterance lengths in the ITS corpus.

Table 1  
 Relationship between utterance length and occurrence of symmetric contrast

	Utterance length $\geq 8$	Utterance length $< 8$
# of utterances containing symmetric contrast	113	19
# of utterances not containing symmetric contrast	140	357

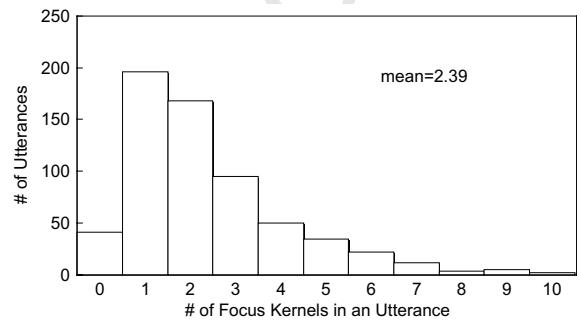


Fig. 2. Histogram of focus kernels per utterance.

630-utterance corpus. An utterance can have maxi-  
 802 mally 10 focus kernels, while on average each  
 803 utterance has 2 focus kernels.  
 804

805 Part-of-speech of all words in the ITS corpus is  
 806 first tagged using an automatic part-of-speech tag-  
 807 ger (Munoz et al., 1999), and then manually  
 808 checked against the tagging standard in Tree-  
 809 bank-3 (Santorini, 1990). The part-of-speech tags  
 810 are used in the automatic detection of symmetric  
 811 contrast and focus kernel, since part-of-speech  
 812 can be used to distinguish content words from  
 813 function words.

#### 4. Prosodic analysis

814  
 815 The literature in both prosody and pragmatics  
 816 reports the pitch accent correlate of contrast and  
 817 focus. Therefore, pitch accent is a reasonable first  
 818 step in the automatic classification of focus kernel  
 819 and symmetric contrast. Because of the man power  
 820 involved in manual labeling of pitch accent in the  
 821 ITS corpus, we try to use an automatic system to  
 822 label pitch accent. To date pitch accent automatic

823 detection concentrates on Radio Speech, in which  
 824 half of all words may be pitch accented (Kim et al.,  
 825 2003), whereas in conversational telephone speech,  
 826 typically about 20% of words are pitch accented  
 827 (Yoon et al., 2004). Pitch accent can be automati-  
 828 cally detected in the Boston Radio Speech Corpus  
 829 with reasonable accuracy: a gender-dependent and  
 830 speaker-independent system achieved a 10% pitch-  
 831 accent detection equal error rate (Kim et al., 2003).  
 832 This section describes the derivation of the pro-  
 833 sodic measurements used for pitch accent labeling.

#### 834 4.1. Duration, pitch and energy

835 English speakers tend to signal prosodic prom-  
 836 inence using extruded pitch, increased energy, and  
 837 longer phoneme durations. In order to match  
 838 durations and energy with individual words, we  
 839 first train a speaker-independent mixture Gaussian  
 840 HMM speech recognizer, and then adapt the rec-  
 841 ognizer to children's speech (Zhang et al., 2004).  
 842 The known word-level transcriptions of each utter-  
 843 ance are expanded into phoneme transcriptions,  
 844 and forcibly aligned to the speech waveform using  
 845 the automatic speech recognizer, resulting in an  
 846 automatic estimate of the duration and time align-  
 847 ment of each word. Pitch and probability of voic-  
 848 ing are extracted using the FORMANT program  
 849 in Entropic XWAVES. Pitch measurements in  
 850 frames with low probability of voicing are dis-  
 851 carded. Extremely high  $F_0$  values (doubling) and  
 852 extremely low  $F_0$  values (halving) are also dis-  
 853 carded. Valid pitch measurements of an utterance  
 854 are then normalized by the highest  $F_0$  in the utter-  
 855 ance to compensate for inter-speaker  $F_0$  differ-  
 856 ences. Kim et al. (2003) found that 95% of the  
 857 pitch accents in their corpus were high- $F_0$  pitch ac-  
 858 cents (in ToBI notation, these were variants of H\*  
 859 or !H\*), and similar statistics have been reported  
 860 for other corpora (e.g., Yoon et al., 2004). There-  
 861 fore, the normalized pitches in the higher half  
 862 pitch region of a word are averaged in order to ob-  
 863 tain an "average peak pitch" measurement for the  
 864 word. The averaging scheme is given by

$$D_m = \frac{1}{N_{\psi_m}} \sum_{f \in \psi_m} F_m[f], \quad (1)$$

where  $D_m$  is the pitch value of word  $W_m$ ,  $F_m[f]$  is  
 the pitch value of frame  $f$  in  $W_m$ ,  $\psi_m$  is a subset  
 of frames in  $W_m$ ,  $N_{\psi_m}$  is the total number of frames  
 in  $\psi_m$ , and

$$\psi_m = \left\{ f | F_m[f] \geq \frac{1}{2} T_m \right\}, \quad (2)$$

$$T_m = \max_{F_m[n] \in W_m} \{F_m[f]\}. \quad (3)$$

Energy of a word is computed using the same  
 strategy as described by Eqs. (1)–(3).

#### 4.2. Spectral balance cepstral coefficients

Spectral balance is intensity emphasis in the  
 higher frequency region. Sluijter et al. (1997)  
 showed that in lexically stressed syllables, speech  
 spectral intensity at higher frequencies (above  
 500 Hz) increased more than the intensity at lower  
 frequencies. Sluijter et al. also found that manipu-  
 lating the high-frequency intensity of speech re-  
 sulted in stronger stress cues than manipulating  
 the entire band. These studies demonstrated that  
 high frequency intensity is a stronger cue for stress  
 than overall intensity. Since pitch accent in spoken  
 sentences normally falls on the lexically stressed  
 syllables, we speculate that intensity of the spec-  
 trum above 500 Hz may be a useful acoustic cue  
 for pitch accent detection.

We propose spectral balance cepstral coeffi-  
 cients (SBCC) to encode the intensity and shape  
 of the speech spectrum in the range between  
 500 Hz and 5000 Hz. The speech waveform is first  
 pre-emphasized and windowed by 30 ms windows  
 with an inter-frame window shift of 10 ms. Speech  
 samples in each frame are decomposed into a series  
 of bands between 0 and 5000 Hz through multires-  
 olution analysis (see Fig. 3). The multiresolution  
 analysis is achieved by iterative application of a  
 filterbank, which consists of quadrature mirror fil-  
 ters (a low-pass smoothing filter and a high-pass  
 differencing filter) with Daubechies-4 orthogonal  
 coefficients (Daubechies, 1990); each time the filter-  
 bank decomposes speech signals into a higher band  
 and a lower band. The Daubechies-4 time domain  
 filterbank has better out-of-band rejection than a  
 filter constructed by adding DFT coefficients (like

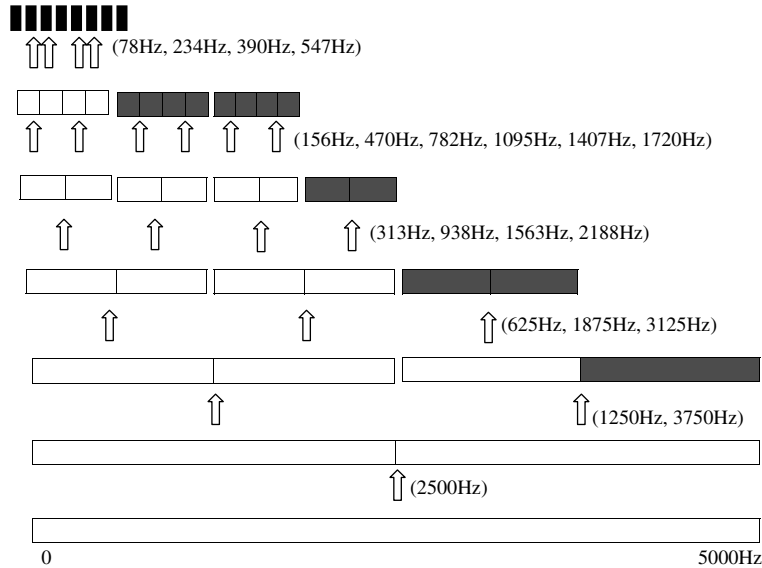


Fig. 3. Multiresolution decomposition of 0–5000 Hz. An arrow splits the band from which it originates into two equal subbands. The frequencies of the splits as arrows indicate are labeled at the right side. The subbands finally obtained from the decomposition are dark marked. The subbands over 500 Hz are the 14 rightmost bands.

914 the filters used in MFCC or PLP analysis), but re-  
 915 tains the desirable properties of MFCC analysis:  
 916 semi-logarithmic frequency scaling, and the flexi-  
 917 bility necessary to select desired sub-bands, e.g.,  
 918 in our case, the sub-band above 500 Hz. After mul-  
 919 tiresolution analysis, each word is characterized by  
 920 its average peak signal intensity in each band, using  
 921 equations similar to those used for pitch:

$$e_m = \frac{1}{N_{\phi_m}} \sum_{n \in \phi_m} y_m^2[n], \quad 1 \leq m \leq M, \quad (4)$$

925 where  $y_m[n]$  is sample  $n$  in band  $B_m$ ,  $e_m$  is the inten-  
 926 sity of band  $B_m$ ,  $M$  is the total number of bands  
 927 spanning from 500 Hz to 5000 Hz,  $N_{\psi_m}$  is the total  
 928 number of speech samples in set  $\Phi_m$ ,  $\Phi_m$  is a subset  
 929 of  $B_m$  and

$$\phi_m = \left\{ n \mid |y_m[n]| \geq \frac{1}{2} T_m \right\}, \quad (5)$$

933 with

$$T_m = \max_{y_m[n] \in B_m} \{|y_m[n]|\}. \quad (6)$$

937 The log energies in each band (logarithms of Eqs.  
 938 (4)–(6)) are transformed using an inverse discrete

cosine transform to compute the cepstral 939  
 coefficients: 940

$$E_l = \sum_{m=1}^M \log(e_m) \cos \left[ \frac{l(m-0.5)\pi}{M} \right], \quad \leq l \leq L, \quad (7)$$

where  $L$  is the desired length of the cepstrum. 943  
 Cepstral mean subtraction is then applied to compen- 944  
 sate for disturbances caused by the transmis- 945  
 sion channel. 946

## 5. Word semantics analysis 947

We use a word dissimilarity measure to model 948  
 the degree of novelty that a word has in compari- 949  
 son with other words. We use  $N_i$  to denote the 950  
 novelty of word  $w_i$  given dialogue context. Accord- 951  
 ing to the definition of focus kernel, we compute 952  
 $N_i$  by the minimum of the dissimilarity between 953  
 $w_i$  and the words in set  $S$ , where  $S$  consists of those 954  
 words appearing in the interlocutor's presupposi- 955  
 tion and those precedent of  $w_i$  in the utterance, i.e., 956

$$N_i = \min_{w_j \in S} \text{dis}(w_i, w_j). \quad (8)$$



960 Semantic parallelism is the semantic attribute  
 961 characterizing symmetric contrast. Given a couple  
 962 of words, the quantification of their semantic par-  
 963 allelism should emphasize: (1) word dissimilarity,  
 964 i.e., a word cannot be parallel to itself or its syno-  
 965 nym; (2) word similarity, e.g., ‘blue’ and ‘white’ are  
 966 more likely than ‘drink’ and ‘pier’ to become a  
 967 symmetric contrast; and (3) non-hypernym rela-  
 968 tion, i.e., a word cannot form symmetric contrast  
 969 with its hypernym. For example, ‘dog’ and ‘cat’  
 970 can form symmetric contrast, while ‘dog’ and ‘ani-  
 971 mal’ cannot.

972 Similarity and dissimilarity of a pair of words  
 973 was computed using measures based on methods  
 974 from the computational linguistics literature. Spe-  
 975 cifically, as suggested in Section 2.3, two methods  
 976 were used: an ontology-based method and a meth-  
 977 od based on corpus statistics. Neither of these two  
 978 methods was found to be an adequate dissimilarity  
 979 measure by itself, but the linear interpolation of  
 980 the two measures was found to be reasonably  
 981 adequate.

982 5.1. Application-oriented ontology

983 Ontology design determines the set of semantic  
 984 categories that properly reflects the particular con-  
 985 ceptual organization of a target domain (Lenci,  
 986 2001). Designing a completely new ontology is  
 987 comparably difficult to the design of a completely  
 988 new dictionary. An attractive solution is to adapt  
 989 a general linguistic resource to the application do-  
 990 main. In this study, we use *WordNet* as the univer-  
 991 sal background linguistic source. WordNet is a  
 992 hierarchical semantic database of English words  
 993 (Miller and Fellbaum, 2002). In WordNet, the  
 994 main relations between words are: synonym, anto-  
 995 nym, and hypernym (representing the ‘is-a’ rela-  
 996 tionship) for nouns, verbs and adjectives; and  
 997 synonym, antonym for adverbs.

998 The ontology in the ITS domain is defined by  
 999  $O = \{C, R, H\}$ , where  $C$  is the set of concepts,  $R$   
 1000 is the set of relations, and  $H$  is the concept hierar-  
 1001 chy. Each content word in the ITS lexicon is a con-  
 1002 cept. Concepts should interconnect with each  
 1003 other in the ontology. The relationship between  
 1004 concepts is represented by links. In this study,  
 1005 the ontology encodes the information of synonym,

1006 antonym, and hypernym for the content words.  
 1007 The hierarchy of word concepts is constructed fol-  
 1008 lowing the procedures in Appendix A. The major  
 1009 semantic categories used in the hierarchical struc-  
 1010 tures are listed in Appendix B. A subset of the  
 1011 ontology is shown schematically in Fig. 4 using a  
 1012 tree structure. In addition, function words are cat-  
 1013 egorized by their part-of-speech under the main  
 1014 semantic class ‘function’ (in contrast with ‘content’,  
 1015 the main semantic class for content words).

1016 We employ the edge-based method, in which an  
 1017 edge represents a direct association between a pair  
 1018 of semantic concepts, to compute the distance be-  
 1019 tween the pair of words. In a more realistic sce-  
 1020 nario, the distances between any two adjacent  
 1021 nodes are not necessarily equal. Generally, the dis-  
 1022 tance shrinks as one descends the hierarchy, since  
 1023 differentiation is based on finer and finer details  
 1024 (Jiang and Conrath, 1997), e.g., the distance be-  
 1025 tween ‘abstraction’ and ‘entity’ (at the top level  
 1026 of the hierarchy) is much bigger than the distance  
 1027 between ‘many’ and ‘some’ (at a lower level). It is  
 1028 therefore necessary to consider that the edge con-  
 1029 necting the two nodes should be weighted by their  
 1030 depth in the hierarchy. We propose the following  
 1031 weighted edge-distance:

$$\text{dis}(w_1, w_2) = \min_{c_1 \in \text{sen}(w_1) c_2 \in \text{sen}(w_2)} \left( \frac{d_{c_1, c_2} + 1}{d_{c_1, c_2}} \right)^\alpha \text{len}(c_1, c_2), \tag{9}$$

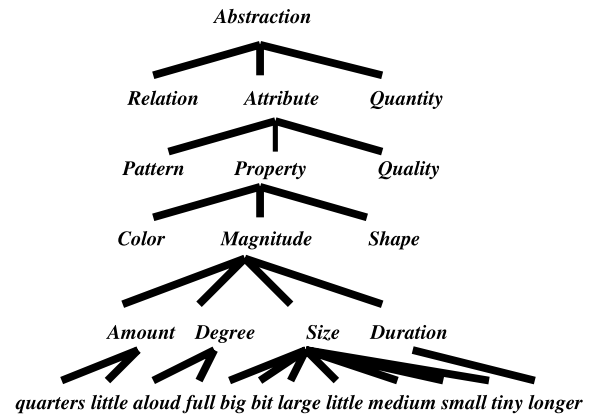


Fig. 4. Partial hierarchy of the application-oriented ontology.

1034 where  $\text{sen}(w)$  denotes the set of possible senses for  
 1035 word  $w$  if case  $w$  has multiple senses;  $d_{c_1, c_2}$  is the  
 1036 mean depth of nodes  $c_1$  and  $c_2$  in the hierarchy;  
 1037  $\text{len}(c_1, c_2)$  is the number of edges connecting  $c_1$   
 1038 and  $c_2$ ; and  $\alpha$  is a constant. Here we choose  
 1039  $\alpha = 3.0$ . In addition, we set the distance score as  
 1040 0 for synonym, and 15.0 for antonym based on  
 1041 the maximum distance between a pair of words  
 1042 in the ontology.

## 1043 5.2. Corpus statistics

1044 Statistical methods attempt to measure depen-  
 1045 dence between words by using statistics taken from  
 1046 a large corpus. Here we use the English *GigaWord*  
 1047 corpus, a billion-word archive of English newswire  
 1048 text distributed by the Linguistic Data Consor-  
 1049 tium. The entire *GigaWord* corpus consists of  
 1050 314 files, nearly 11.7 GB. All texts are presented  
 1051 in SGML form, using a simple markup structure.  
 1052 We remove the SGML tags, leaving only text con-  
 1053 tent. The raw database consists of four different  
 1054 types of documents: *Story*, *Multi*, *Advis*, and  
 1055 *Other*. *Story* has uniform format, each consisting  
 1056 of a few paragraphs describing a same topic. The  
 1057 other three types of documents do not have long  
 1058 paragraphs of text. Therefore, we only choose  
 1059 the story session for experiments; the story session  
 1060 (with paragraph breaks and some other control  
 1061 characters removed) contains 8.57 GB of text.

1062 In this study, similarity between a pair of words  
 1063  $w_1$  and  $w_2$  is measured based on the proximity and  
 1064 topicality assumptions (Section 2.3). First, the  
 1065 proximity assumption tells us that we can model  
 1066 word similarity by the probability of a pair of  
 1067 words occurring together in a sentence, a para-  
 1068 graph, a topic, or even a document. We measure  
 1069 the degree of association between a pair of words  
 1070  $w_1$  and  $w_2$  by mutual information

$$\text{PMI}(w_1, w_2) = \log_2 \frac{p(w_1, w_2)}{p(w_1)p(w_2)}$$

1073 Let the co-occurrence frequency of  $w_1$  and  $w_2$  be  
 1074 denoted by  $f_{w_1, w_2}$ . Because *GigaWord* is a very  
 1075 large text corpus, the co-occurrence frequency is  
 1076 roughly estimated by the number of topics (stories)  
 1077 in which  $w_1$  and  $w_2$  co-occur. Let  $N$  be the size of  
 1078 the corpus in terms of topics. The maximum like-

lihood estimate of the co-occurrence probability 1079  
 is given by  $p(w_1, w_2) = f_{w_1, w_2}/N$ . To compensate 1080  
 for the sparsity of the training data, we compute 1081  
 $p(w_1, w_2)$  by 1082

$$p(w_1, w_2) = \frac{f_{w_1, w_2} + 1}{N + 1} \quad (10)$$

The computations of  $p(w_1)$  and  $p(w_2)$  adopt the 1085  
 same smoothing strategy. Second, the topicality 1086  
 assumption tells us that if we have decided that 1087  
 two words are similar, then we may infer that they 1088  
 have similar mutual information with some other 1089  
 word,  $w$ . Given a context  $C = \{w'_1, w'_2, \dots, w'_n\}$ , 1090  
 $w_1$  and  $w_2$  are considered semantically similar if 1091  
 they are both likely to co-occur with the words 1092  
 in  $C$ . We apply a method (Pantel and Lin, 2002) 1093  
 that computes the cosine distance between the 1094  
 two partial mutual information (PMI) vectors cor- 1095  
 responding to  $w_1$  and  $w_2$ : 1096

$$\begin{aligned} \text{dis}(w_1, w_2) \\ = 1 - \frac{\sum_{w' \in C} \text{PMI}(w', w_1) \text{PMI}(w', w_2)}{\sqrt{\sum_{w' \in C} \text{PMI}(w', w_1)^2} \sqrt{\sum_{w' \in C} \text{PMI}(w', w_2)^2}}, \end{aligned} \quad (11)$$

where  $C = C(w_1) \cup C(w_2)$ ,  $C(w_1)$  and  $C(w_2)$  are con- 1100  
 text of  $w_1$  and  $w_2$ , respectively. Given a word  $w$ , its 1101  
 context  $C(w)$  can be the simple collection of all 1102  
 (content) words appearing within the window of 1103  
 $w$  in the ITS corpus. To more properly represent 1104  
 those words semantically associated with  $w$ , we 1105  
 modify the method proposed by Dagan et al. 1106  
 (1995) to determine  $C(w)$ : 1107

1. Define a pair of words  $(w, v)$  to be strong neigh- 1108  
 bors if  $f(w, v) > t_f$ , where  $f(w, v)$  is the count of 1109  
 $(w, v)$  in a window of size  $d$ ,  $t_f$  is threshold and 1110

$$t_f = \begin{cases} \lambda_1 \times \min(N_w, N_v) & \min(N_w, N_v) > \gamma_1 \\ \tau_1 & \min(N_w, N_v) \leq \gamma_1 \end{cases}, \quad (12)$$

where  $N_w$  and  $N_v$  are the number of occurrences 1113  
 of  $w$  and  $v$ , respectively, and  $\lambda_1$ ,  $\gamma_1$ , and  $\tau_1$  are 1114  
 constant.

2. Collect all the strong neighbors of  $w$  as poten- 1116  
 tial candidates and denote as  $C_1(w)$  1117

$$C_1(w) = \{v | f(w, v) > t_f\}. \quad (13)$$

1120 3. Collect the strong neighbors of all words in  
1121  $C_1(w)$  as potential candidates, and denote as  
1122  $C_2(w)$

$$C_2(w) = \{v|v \in C_1(C_1(w))\}. \quad (14)$$

1125 4. Collect those words which share at least  $t_N$   
1126 strong neighbors with  $w$  in the lexicon and  
1127 denote as  $C_3(w)$

$$C_3(w) = \{v|C_1(w) \cap C_1(v) > t_N\}, \quad (15)$$

1130 where

$$t_N = \begin{cases} \lambda_2 \times \min(C_1(w), C_1(v)) & \min(C_1(w), C_1(v)) > \gamma_2 \\ \tau_2 & \min(C_1(w), C_1(v)) \leq \gamma_2 \end{cases}, \quad (16)$$

1133  $N_w$  and  $N_v$  are the number of occurrences of  $w$   
and  $v$ , respectively,  $\lambda_2$ ,  $\gamma_2$ , and  $\tau_2$  are constant.

1135 5.  $C(w) = C_1(w) \cup C_2(w) \cup C_3(w)$ . (17)  
1139

1140 Examples of corpus statistics are given in  
1141 [Appendix C](#).

### 1142 5.3. Knowledge combination

1143 The computation of word semantic dissimilarity  
1144 is a combination of lexical ontology and corpus  
1145 statistics. First, the manually built pseudo-knowl-  
1146 edge base has advantages in efficient paraphrasing,  
1147 inference and reasoning. The ontology is especially  
1148 useful when the dissimilarities of some words are  
1149 dependent on topic and pragmatic context. For  
1150 example, ‘play’ and ‘build’ are dissimilar to each  
1151 other under a basketball game topic, but similar  
1152 under the ITS topic. However, the top-down orga-  
1153 nization of lexical knowledge cannot incorporate  
1154 the dynamic nature of word meanings. Word  
1155 meanings change dramatically depending on the  
1156 linguistic context. This context-dependent charac-  
1157 teristic is one of the main empirical arguments  
1158 for real text data ([Lenci, 2001](#)). Second, the statis-  
1159 tical model provides computational evidence from  
1160 distributional analysis of corpora data. The statis-  
1161 tical model provides a quantification of the seman-  
1162 tic space. However, corpus statistics has some  
1163 limitations, as it requires the exact correspondence  
1164 between terms (word or character  $n$ -grams). In

unrestricted language many reasonable co-occur-  
ences may fail to occur in the training corpus. 1165  
1166

In order to combine information from both the  
ontology-based dissimilarity measure and the cor-  
pus-based dissimilarity measure, we use linear  
interpolation as 1167  
1168  
1169  
1170

$$\text{dis}(w_1, w_2) = \lambda \text{dis}_1(w_1, w_2) + (1 - \lambda) \text{dis}_2(w_1, w_2), \quad (18)$$

where  $\text{dis}_1$  and  $\text{dis}_2$  are the peak-normalized word  
distances based on the ontological and the statisti-  
cal methods, respectively; and parameters  $\lambda$  is con-  
stant. Here we arbitrarily choose  $\lambda = 0.5$ . The  
prerequisite to use corpus statistics is that there  
must be at least two contextual words for a given  
pair of words ( $w_1, w_2$ ). Otherwise,  $\text{dis}_2$  in Eq. (11)  
will be 1 (# of contextual words is 0) or 0 (# of  
contextual words is 1). If the prerequisite fails to  
be satisfied, then the dissimilarity is computed  
using only the ontology-based dissimilarity mea-  
sure. Examples of feature combination are listed  
in [Appendix D](#). 1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185

## 1186 6. System evaluation

The corpus for focus kernel classification con-  
sisted of 630 multi-word, multi-phrase utterances,  
containing approximately 5700 words and 48 min  
of speech. In the experiments of extracting focus  
kernel and symmetric contrast, training and test  
data included different utterances from the same  
set of talkers, so the experiments were multi-speaker  
speaker-dependent. 1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194

### 1195 6.1. Focus kernel detection

Prosodic observations (duration, average peak  
pitch, average peak energy, and spectral balance  
cepstral coefficients) and part-of-speech tagging  
were integrated with each word’s semantic novelty  
measurement (Eq. (8)) using a time-delay recurrent  
neural network (TDRNN) for the purpose of  
automatically detecting focus kernels. The  
TDRNN is a neural network that encodes and  
integrates dynamic signal context using a combi-  
nation of delayed input nodes (for temporal loca-  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205

1206 tion of important information in the input se-  
 1207 quence) and delayed recurrent nodes (for temporal  
 1208 context information) (Kim, 1998). This study used  
 1209 a modified TDRNN: we added a recurrent layer  
 1210 between the output layer and the hidden layer to  
 1211 feed back the delayed values from the output layer;  
 1212 and the time indices were synchronized to succes-  
 1213 sive words rather than centisecond speech frames.

1214 We have tried using the non/pitch accent auto-  
 1215 matic label as a feature for focus kernel detection:  
 1216 we employed a model trained from Boston Radio  
 1217 News Corpus (Ren et al., 2004a) for our automatic  
 1218 pitch accent labeling task; then the pitch accent/  
 1219 nonpitch accent labels were combined with seman-  
 1220 tic similarity measure to input to the TDRNN.  
 1221 However, the performance was not satisfying,

1222 apparently because of: (1) the application of a  
 1223 model trained on radio news speech to spontane-  
 1224 ous speech; (2) the noisy recording environment;  
 1225 and (3) the fact that erroneous pitch accent detec-  
 1226 tion would lead to erroneous extraction of focus  
 1227 kernel. Therefore, we used the acoustic correlates  
 1228 of pitch accent, rather than pitch accent itself,  
 1229 for the purpose of detecting focus kernel.

1230 We used 90% of the corpus for training and the  
 1231 remaining 10% for test. Our experiments yielded  
 1232 an accuracy of 83.8% on the test set, which con-  
 1233 sisted of 536 words in total. The confusion matri-  
 1234 ces summarizing the classification performance in  
 1235 terms of utterances and words are shown in Tables  
 1236 2 and 3, respectively. In Table 2, *novel* and  
 1237 *nonnovel* denote whether the utterances contain  
 1238 focus kernel or not. Precision, recall, and *F*-score  
 1239 ( $f = \frac{1}{0.5/p+0.5/r}$ ) of non/focus kernel automatic  
 1240 labeling is presented in Table 4. We further present  
 1241 the *F*-scores of non/focus kernel labeling based on  
 1242 different features in Fig. 5. The figure shows that  
 1243 pitch and energy yielded better performance for  
 1244 focus kernel labeling than nonfocus kernel label-  
 1245 ing, while the other features were the visa versa,  
 1246 especially part-of-speech and spectral balance  
 1247 cepstral coefficients.

1248 We compared the efficiency of individual fea-  
 1249 tures using classification accuracy and the *F*-  
 1250 score-based evaluation measures. Classification  
 1251 accuracy is the fraction of the test set that is cor-  
 1252 rectly classified with respect to focus kernel and  
 1253 nonfocus kernel. *F*-score is a harmonic mean of  
 1254 precision and recall: *F*-score is high only when

Table 2

Confusion matrix for utterance classification: novel vs. non-  
 novel; novel = utterance containing focus kernel, non-  
 novel = utterance not containing focus kernel

	Novel	Nonnovel
Novel	53	7
Nonnovel	0	3

Table 3

Confusion matrix for word classification: focus kernel vs.  
 nonfocus kernel

	Focus kernel	Nonfocus kernel
Focus kernel	148	17
Nonfocus kernel	70	301

Table 4

Precision *p*, recall *r*, and *F*-score *f* of non/focus kernel labelling using various features

	Focus kernel			Nonfocus kernel		
	<i>p</i>	<i>r</i>	<i>f</i>	<i>p</i>	<i>r</i>	<i>f</i>
Dur	0.417	0.727	0.530	0.819	0.547	0.656
Egy	0.313	0.885	0.462	0.725	0.135	0.227
Pit	0.312	0.891	0.462	0.723	0.127	0.216
Pos	0.582	0.236	0.336	0.731	0.925	0.817
SBCC	0.373	0.364	0.368	0.720	0.728	0.724
SDM	0.472	0.873	0.613	0.909	0.566	0.698
Dur+Egy+Pit+Pos +SBCC+SDM	0.809	0.770	0.789	0.900	0.919	0.909

Dur = duration, Egy = energy, Pit = pitch, Pos = part-of-speech, SBCC = spectral balance cepstral coefficients, SDM = semantic dissimilarity measure.



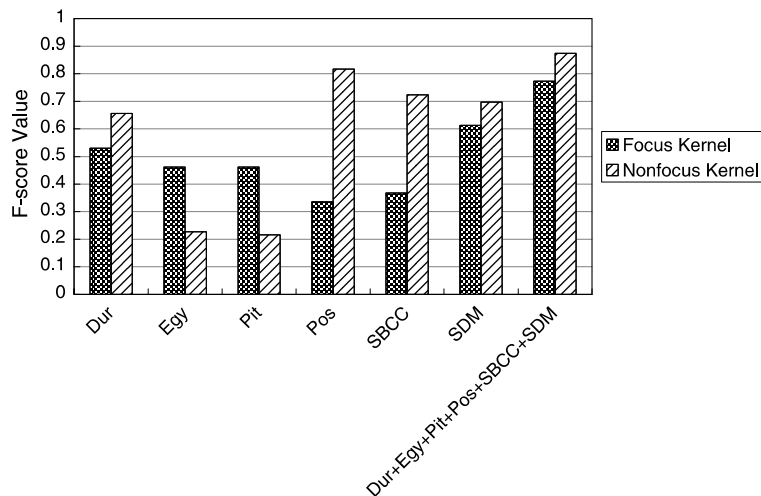


Fig. 5.  $F$ -score values of focus kernel labeling and nonfocus kernel labeling based on various features. Dur = duration, Egy = energy, Pit = pitch, Pos = part-of-speech, SBCC = spectral balance cepstral coefficients, SDM = semantic dissimilarity measure.

1255 both precision and recall are high. Therefore,  
 1256  $F$ -score is a more reliable measure. However, the  
 1257  $F$ -score values that we have derived were class-  
 1258 dependent. We were interested in the global  
 1259 performance measure across the classes, so we de-  
 1260 signed two global  $F$ -score measures: (1) averaged  
 1261 focus kernel  $F$ -score and nonfocus kernel  $F$ -score;  
 1262 and (2)  $f_{ave} = \frac{1}{0.5/p_{ave} + 0.5/r_{ave}}$ , where  $p_{ave}$  was the  
 1263 averaged focus kernel precision and nonfocus  
 1264 kernel precision, and  $r_{ave}$  was the averaged focus  
 1265 kernel recall and nonfocus kernel recall. The com-  
 1266 parison results using classification accuracy and  
 1267 global  $F$ -score measures are shown in Fig. 6. The  
 1268 figure shows that the combination of features  
 1269 much outperformed the individual features. In  
 1270 addition, comparison based on the global  $F$ -score  
 1271 measures shows that word dissimilarity measure  
 1272 was the most important feature, followed by dura-  
 1273 tion, part-of-speech, and SBCC.

1274 As for the acoustic correlates of pitch accent,  
 1275 the test results depicted in Fig. 6 show that dura-  
 1276 tion and SBCC played more important roles than  
 1277 energy and pitch for the detection of focus kernels.  
 1278 Both duration and high-frequency intensity have  
 1279 been shown to be reliable acoustic correlates of  
 1280 lexical stress (Sluijter et al., 1997). Pitch played  
 1281 the least significant role in the focus kernel detec-  
 1282 tion, probably by the frequent occurrence of pitch

1283 tracking inaccuracies in noisy speech data. It was  
 1284 hard to discriminate voiced regions from unvoiced  
 1285 regions in the noisy recording environment. The  
 1286 energy of unvoiced regions carried information  
 1287 irrelevant to pitch estimate, therefore the automati-  
 1288 cally extracted pitch might have contained too  
 1289 many pitch tracking errors to be an efficient fea-  
 1290 ture. Spectral balance cepstral coefficients showed  
 1291 better performance than pitch and energy, possibly  
 1292 because: (1) the band-pass filters eliminated the  
 1293 low frequency noise that adversely affected pitch  
 1294 and energy estimates; and (2) spectral balance  
 1295 may be more representative of vocal effort than  
 1296 was overall energy, because the difference in the  
 1297 shape of glottal waveform tended to cause differ-  
 1298 ences in the higher frequency regions (Sluijter  
 1299 et al., 1997).

## 1300 6.2. Symmetric contrast detection

1301 Symmetric contrast was detected using the algo-  
 1302 rithm shown in Fig. 7: candidate word pairs were  
 1303 pre-filtered using a series of knowledge-based  
 1304 rules, prior to application of a decision tree pro-  
 1305 gram. For the decision tree program, we used  
 1306 *See5*, which is a data mining tool for discovering  
 1307 patterns or relationships in data, assembling them  
 1308 into classifiers that are expressed as decision trees

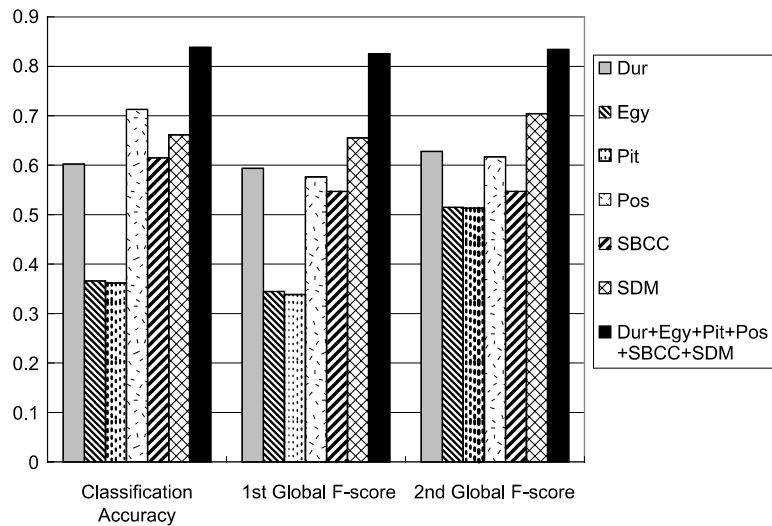


Fig. 6. Feature vector accuracy ranking according to the classification accuracy and two global *F*-score measures: the first measure is the averaged focus kernel *F*-score and nonfocus kernel *F*-score; the second measure is *F*-score computed from the average precision and average recall of the focus kernel and nonfocus kernel labeling. Dur = duration, Egy = energy, Pit = pitch, Pos = part-of-speech, SBCC = spectral balance cepstral coefficients, SDM = semantic dissimilarity measure.

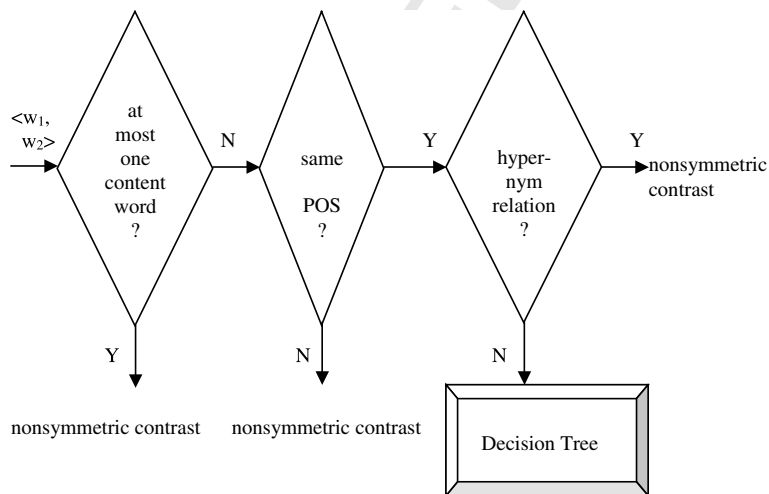


Fig. 7. Architecture of the symmetric contrast detector.

1309 or sets of if-then rules, and using them to make va- 1310  
 1310 lid predictions (Rulequest Research, 2004).

1311 In experimental tests, we labeled all of word 1312  
 1312 pairs in an utterance as either *symmetric contrast* 1313  
 1313 or *nonsymmetric contrast*. We had 265 symmetri- 1314  
 1314 cally contrasted pairs. The corpus study showed 1315  
 1315 that 100% of symmetric contrast cases were pairs

of content words, and 99.2% of contrasted con- 1316  
 1316 tent-word pairs had the same part-of-speech tags. 1317  
 1317 Therefore, we labeled a pair of words  $\langle w_1, w_2 \rangle$  as 1318  
 1318 *nonsymmetric contrast* if at most one of them was 1319  
 1319 a content word, or if they had different part-of- 1320  
 1320 speech tags. We did not model the common inte- 1321  
 1321 grator requirement for semantic parallelism, but 1322

1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322

1323 the semantic independence requirement could be  
 1324 easily represented by requiring that a word may  
 1325 never contrast with its own hypernym.

1326 A pair of words  $\langle w_1, w_2 \rangle$  in the utterance was in-  
 1327 put to the *See5* program, with feature vector being  
 1328 the combination of prosodic measurements and  
 1329 the semantic similarity score, i.e.,  $f = \{\text{duration,}$   
 1330  $\text{pitch, energy, 13 SBCCs, similarity score}\}$ . The  
 1331 similarity score of  $w_1$  and  $w_2$  was

$$\text{sim}(w_1, w_2) = 1 - \text{dis}(w_1, w_2). \quad (19)$$

1334 The prosodic features were the arithmetic mean of  
 1335  $w_1$  and  $w_2$ . The classification target was 1 if  
 1336  $\langle w_1, w_2 \rangle$  was symmetric contrast, and 0 else. The  
 1337 *See5* program was invoked with the options *rule-*  
 1338 *sets* (collections of *if-then* rules), *ignore costs file*  
 1339 (ignore the penalty for misclassification), and *global*  
 1340 *pruning* (a large tree is first grown to fit the data  
 1341 closely and is then pruned by removing parts that  
 1342 are predicted to have a relatively high error rate).

1343 After pre-filtering the word pair candidates  
 1344 using the three knowledge-based rules, we ob-  
 1345 tained 803 word pairs (260 utterances out of the  
 1346 630 utterance corpus) to input to the *See5* pro-  
 1347 gram. Due to the reason that we describe in Sec-  
 1348 tion 6.1, we used the pitch accent acoustic  
 1349 correlates (pitch, energy, duration, and SBCC)  
 1350 rather than pitch accent/nonpitch accent labeling  
 1351 from these acoustic correlates for symmetric con-  
 1352 trast detection. We used 50% of the filtered word

1353 pairs (in 121 utterances) for training the *See5* pro- 1353  
 1354 gram. Then we used the symmetric contrast detec- 1354  
 1355 tor depicted in Fig. 7 to test the other 509 1355  
 1356 (630 – 121 = 509) utterances. Table 5 reports the 1356  
 1357 precision, recall, and *F*-score with which utter- 1357  
 1358 ances were correctly labeled as either containing 1358  
 1359 symmetric contrast (first row in the table) or not 1359  
 1360 containing symmetric contrast (second row in the 1360  
 1361 table). Our test results showed that some individ- 1361  
 1362 ual prosodic features, i.e., pitch, duration and en- 1362  
 1363 ergy, all classified the test word pairs (after pre- 1363  
 1364 filtering) into *nonsymmetric contrast*. In addition, 1364  
 1365 when they were combined with SBCC and seman- 1365  
 1366 tic dissimilarity measure, they played little role in 1366  
 1367 the symmetric contrast detection: almost all the 1367  
 1368 variables in the rule set used for non/symmetric 1368  
 1369 contrast classification were SBCC variables and 1369  
 1370 the semantic measure. Therefore, pitch, energy 1370  
 1371 and duration were not efficient in the symmetric 1371  
 1372 contrast identification. We further compared the 1372  
 1373 performance of SBCC and semantic similarity 1373  
 1374 measure on the pairs of words after the knowledge 1374  
 1375 rule-based pre-filtering, and listed the results in 1375  
 1376 Table 6. The table shows that SBCC outperformed 1376  
 1377 word similarity measure in the efficiency of non/ 1377  
 1378 symmetric contrast classification. We also com- 1378  
 1379 puted the non/symmetric contrast classification 1379  
 1380 accuracy over all of the word pairs in the 509- 1380  
 1381 utterance test corpus, and the accuracy was 92.8%. 1381

Table 5

Precision *p*, recall *r*, and *F*-score *f* for the correct labeling of utterances containing symmetric contrast (row 1) and not containing symmetric contrast (row 2)

	<i>p</i>	<i>r</i>	<i>f</i>
Symmetric Contrast	1.0	0.952	0.975
Nonsymmetric Contrast	0.925	1.0	0.961

Table 6

Precision *p*, recall *r*, and *F*-score *f* of symmetric contrast labelling using various features on the word pairs after the three knowledge rule-based pre-filtering

	Symmetric contrast			Nonsymmetric contrast		
	<i>p</i>	<i>r</i>	<i>f</i>	<i>p</i>	<i>r</i>	<i>f</i>
SBCC	0.789	1.000	0.882	1.000	0.841	0.914
SSM	0.701	1.000	0.824	1.000	0.746	0.855
Dur+Egy+Pit+SBCC+SSM	0.838	1.000	0.912	1.000	0.885	0.939

Dur = duration, Egy = energy, Pit = pitch, SBCC = spectral balance cepstral coefficients, SSM = semantic similarity measure.

### 6.3. Application of focus kernel to a practical SSLU system

1384 Focus kernel has been successfully applied to 1384  
 1385 spontaneous spoken language understanding in 1385  
 1386 the tutoring dialogue scenario by participating in 1386  
 1387 the tutoring event classification (Zhang, 2004). 1387  
 1388 The study corpus was divided into a training set 1388

1389 and a test set. Focus kernels were extracted from  
1390 training utterances. A nonparametric model of  
1391 each tutoring event class was constructed by sim-  
1392 ply listing all of focus kernels used in training  
1393 utterances of that tutoring event. The tutoring  
1394 event classification of test utterances was based  
1395 on: (1) lexical similarity measure between focus  
1396 kernels in the test utterance and focus kernels in  
1397 the nonparametric models of all tutoring event  
1398 candidates; and (2) cognitive state (*confidence*, *puz-*  
1399 *zlement*, or *hesitation*) that reflects the students'  
1400 mental activities during the process of knowledge  
1401 acquisition. Because the lexical similarity measure  
1402 was trained using a task-independent corpus  
1403 (GigaWord) combined with a small task-oriented  
1404 ontology, it was possible to create this SSLU  
1405 system using a comparatively small amount of  
1406 task-specific training data. There were 30 tutoring  
1407 events in total, and the perplexity of the classifica-  
1408 tion task was 22.5. Accuracy of the tutoring  
1409 event classification achieved 75.5% when focus  
1410 kernels and cognitive states were manually anno-  
1411 tated, and reduced by 15.4% relative when focus  
1412 kernels and cognitive states were automatically  
1413 extracted.

1414 Symmetric contrast labels have not yet been  
1415 implemented to our speech understanding task,  
1416 but we expect that automatically labeled symmet-  
1417 ric contrast will be useful as a cue for robust  
1418 semantic parsing.

## 1419 7. Discussion and conclusions

1420 This paper has computationalized two linguistic  
1421 concepts, contrast and focus, that were assumed to  
1422 be useful for robust understanding of spontaneous  
1423 spoken messages in a dialogue system. Standard  
1424 and reasonable linguistic definitions of focus are  
1425 difficult to implement computationally, because  
1426 the scope of focus is dependent on the syntactic  
1427 structure of the utterance and highly variable. In  
1428 order to create a computationally feasible focus  
1429 detector, this paper has defined focus kernel to  
1430 be the word in an utterance containing new infor-  
1431 mation neither presupposed by the interlocutor  
1432 nor contained in the precedent words of the utter-  
1433 ance. This paper has also defined symmetric con-

1434 trast to be a pair of words that are parallel or  
1435 symmetric in linguistic structure but different or  
1436 contrastive in meaning. Symmetric contrast marks  
1437 information about the discourse structure of an  
1438 utterance that may be useful for robust semantic  
1439 parsing.

1440 Novelty detection has been widely investigated  
1441 in natural language processing (NLP) at the sen-  
1442 tence or document level. For example, TREC  
1443 Novelty Track and Topic Detection and Tracking  
1444 intend to extract sentences or documents that dis-  
1445 cuss new development of a document or discourse  
1446 topic. In this study, we have extended the study of  
1447 informative novelty to an analysis of word-level  
1448 novelty in spontaneous speech. In future work, it  
1449 may be possible to use the focus kernel for the pur-  
1450 pose of efficiently summarizing the content of an  
1451 utterance. By definition, the set of words defined  
1452 to be focus kernel is a maximally informative sum-  
1453 mary of the utterance in context. In this sense, the  
1454 use of focus kernels for utterance summarization  
1455 would be similar to the task of sentence-selection  
1456 based topic summarization of a document in  
1457 NLP: topic summarization of a document can be  
1458 realized by seeking a pool of sentences, each of  
1459 which expresses information not contained in its  
1460 precedent sentences in the document (Zechner  
1461 and Waibel, 2000).

1462 The effectiveness of the proposed symmetric  
1463 contrast detection and focus kernel detection sys-  
1464 tems has been evaluated using a transcribed chil-  
1465 dren's spontaneous speech corpus, which was  
1466 collected using Wizard-of-Oz simulations of an  
1467 intelligent tutoring dialogue system. The detection  
1468 of symmetric contrast and focus kernel was based  
1469 on word dissimilarity measures, part-of-speech  
1470 tagging, and measurements of the acoustic corre-  
1471 lates of prosody including duration, pitch, and en-  
1472 ergy. We also propose spectral balance cepstral  
1473 coefficients (intensity and shape of the high-fre-  
1474 quency spectrum) as one more acoustic correlate  
1475 of pitch accent. The word dissimilarity measure  
1476 combines corpus statistics and application-ori-  
1477 ented ontology. Classification achieved accuracies  
1478 of 83.8% for focus kernel classification and  
1479 92.8% for symmetric contrast classification. Our  
1480 tests showed that spectral balance cepstral coeffi-  
1481 cients, the dissimilarity semantic measure, and



1482 part-of-speech played important roles in the focus  
1483 kernel and symmetric contrast detections.

1484 Focus kernel extraction has been applied in a  
1485 spontaneous spoken language understanding sys-  
1486 tem. Symmetric contrast has not yet been applied  
1487 in any spoken language understanding (SLU) sys-  
1488 tem, but seems well suited for use as a cue in ro-  
1489 bust semantic parsing. Our SLU based on the  
1490 detection of focus kernel and symmetric contrast  
1491 differs from most existing SLU systems in that it  
1492 is not keyword-based (does not require the system  
1493 designer to specify, in advance, a list of “case  
1494 markers”), and it also does not require predefined  
1495 grammar (used for syntactic parsing or extraction  
1496 of finite-state-machine-based semantic concepts).  
1497 Focus kernel and symmetric contrast begin SSLU  
1498 by extracting, in a relatively general way, the prag-  
1499 matic and semantic salient information of the  
1500 utterance, thus allowing the interpretation of  
1501 utterance meaning from unconstrained human  
1502 language.

## 1503 8. Uncited reference

1504 Ren et al. (2004b).

## 1505 Acknowledgements

1506 We would like to thank Richard W. Sproat for  
1507 helpful discussions and providing us with the  
1508 GigaWords text corpus. We would also like to  
1509 thank Carla Umbach and Chungmin Lee for their  
1510 comments and suggestions. This work is supported  
1511 by NSF grant number 0085980. Statements in this  
1512 paper reflect the opinions and conclusions of the  
1513 authors, and are not endorsed by the NSF.

## 1514 Appendix A

1515 The ontology construction procedure:

1516 1. Noun, verb, adjective, and adverb are catego-  
1517 rized into the semantic class *Content*, and func-

tion words are categorized into the semantic  
class *Function*.

2. The content words are clustered according to  
their semantic meanings, regardless of their  
part-of-speech:

(1) For each *noun*

a. According to the word meaning in the  
ITS dialogue context, search the Word-  
Net for the most appropriate hypernym  
of the word. Delete those concepts that  
contain redundant information for  
knowledge representation in the ITS sce-  
nario. For example, in WordNet we have  
month → Gregorian month → April (‘→’  
denotes subordination). However, we  
delete ‘Gregorian month’ in constructing  
ontology of the ITS corpus.

b. Check WordNet to see if any word in  
the corpus can be the synonym of the  
noun.

c. If the noun has no hypernym, then it  
becomes a direct subordinate of *Content*.

(2) After the hierarchy of noun words has been  
constructed, we

a. Verify that every member of a concept  
hierarchy is subordinate of the root  
concept.

b. Verify that the direct subordinates of a  
concept are semantically parallel to each  
other.

c. Verify that the hierarchy rooted in a  
concept has distinct meaning from the  
hierarchies rooted in other parallel  
concepts.

(3) For each *verb*, strip tense, and then check  
the WordNet to see whether its hypernyms  
(‘is one way to’ relationship) are identical or  
similar to some noun concepts has been  
defined in the ontology. If yes, then merge  
the associated hierarchy into the existing  
ontology. Otherwise, add a new verb con-  
cept to the ontology following the step  
(1). Some verbs have multiple meanings rel-  
evant to the ITS scenario, e.g., *make* means

1518

1519

1520

1521

1522

---

*Abstraction*  
*Attribute*  
     *Disposition*  
     *Pattern*  
     *Property*  
         *Color*  
         *Magnitude=Amount + Degree + Duration + Size*  
         *Shape=Circles + Crooked + Space*  
         *Quality=Certainty + Correctness + Difference + Excellency + Hardness*  
             + *Importance + Strength*  
     *Quantity*  
         *DefiniteQuantity*  
             *Number=Count + Integer + Fraction*  
         *FundamentalQuantity*  
             *TimePeriod=Unit + Calendar + Time*  
         *IndefiniteQuantity*  
         *Relation=Group + Indication + Link + QuantitativeRelation*  
*Act:*  
     *Change*  
         *Integrity*  
         *Magnitude*  
         *Ownership*  
         *Position=Lined + Move + Rotation + Step*  
         *State*  
         *Do=Exploring + Make + Paint + Put + Show + Try + Use*  
         *Duty*  
         *Interact*  
             *Communication=Approval + Convey*  
             *Contact*  
         *SpeechAct*  
             *Inform*  
             *Request=Ask + Question*  
             *Say=Answer + Repeat*  
     *Be=BeClass + Equal + Have + Stay*  
     *Entity*  
         *Location=Line + Point + Region*  
         *Object*  
             *Artifact=Construction + Device + Representation + Transport*  
         *Thing=Anything + Body + Life*  
     *Phenomenon*  
         *Effect*  
         *PhysicalPhenomenon=Force + Resistance*  
     *PsychologicalFeature*  
         *Cognition*  
             *Awareness*  
             *Information=Example + Experience*  
             *Intelligence*  
             *Perception=Find + See*  
             *Prospect*  
             *Think=Calculation + Figure + Guess*  
         *Feeling*

---

Fig. B1. Major ontological categories in the ITS corpus.

- ‘make somebody do something’ or ‘produce something.’ Decode the multiple meanings of a word, if any.
- (4) Check the ontology structure:
    - a. Search among the ontology for sets of concepts that can be merged together.
    - b. Search for more synonym words in the ITS dialogue scenario. The synonyms of a word can be words of different part-of-speech, e.g., a noun can be a synonym of a verb as long as their semantic meanings are similar, such as *rotation* and *spin*.
    - c. Go to step (2).
  - (5) For each *adjective*, check the WordNet to see whether its hypernyms (‘is a value of relationship’) are identical or similar to a noun concept. If yes, then merge the associated hierarchy into the existing ontology. Otherwise, use the verb or noun derivation (or root) of the adjective to find its hypernym, following step (1). For example, we use the information about *confuse* to find the hypernym of *confused*. Define its synonym and antonym with reference to the WordNet. If the hypernym of the adjective cannot be found, then we use the hypernym of its noun or verb synonyms. For example, we use the hypernym of the verb *agree* to find the hypernym of *alright*. Otherwise, use the synonyms of the adjective as clues to find its hyponyms. Then go to step (4).

- (6) For each *adverb*, WordNet does not define a hypernym. Therefore, we use the stem adjectives as information clues to find the hypernyms. If an adverb has no stem adjective, then use the stem adjectives of its synonyms as information clues. Then go to step (4).
 

1604  
1605  
1606  
1607  
1608  
1609  
1610  
1612
3. Function words typically have little or no semantic content apart from their syntactic use, so we categorize them according to syntactic usage: the words of the same part-of-speech are synonyms to each other, and their part-of-speech is their hypernym.
 

1611

**Appendix B.** Fig. B1. 1613

**Appendix C** 1614

To illustrate the procedure of using corpus statistics for word dissimilarity, consider the utterance “*I lost count let me try again.*” We choose  $\lambda_1 = \lambda_2 = 0.6$ ,  $\gamma_1 = 2$ ,  $\tau_1 = 2$ ,  $\gamma_2 = 1$ ,  $\tau_2 = 1$ , and present the derived context of the content words in Fig. C1. The PMI scores of some pairs of words are shown in Fig. C2. The listed scores show the problem with mutual information is that it is biased towards infrequent words, such as ‘alright.’ The dissimilarity scores are shown in Fig. C3. Therefore as defined by Eq. (8), we have the novelty measures for the content words in the utter-

1615  
1616  
1617  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626

		$C_1$	$C_2$	$C_3$
1	<i>lost</i>	Count	let, loud	$\emptyset$
2	<i>count</i>	let, lost, loud	alright, out, start, try	Try
3	<i>Let</i>	alright, count, start, try	Figure	figure, loud
4	<i>Try</i>	figure, let	alright, count, out, start	count, out
5	<i>again</i>	$\emptyset$	$\emptyset$	$\emptyset$

Fig. C1. Context examples.

1	alright count	2.545	8	count figure	2.102
2	alright let	3.848	9	count lost	1.945
3	alright lost	2.567	10	count loud	1.989
4	alright try	2.982	11	count out	1.873
5	figure lost	0.969	12	count start	1.809
6	figure out	1.237	13	count try	1.840
7	figure try	1.100	14	start try	1.903

Fig. C2. PMI examples.

1	count lost	0.021	6	try let	0.006
2	let lost	0.014	7	again lost	0.002
3	let count	0.037	8	again count	0.003
4	try lost	0.004	9	again let	0.011
5	try count	0.027	10	again try	0.003

Fig. C3. Examples of dissimilarity scores.

1627 ance:  $N_{\text{count}} = 0.021$ ,  $N_{\text{let}} = 0.014$ ,  $N_{\text{try}} = 0.004$ ,  
 1628  $N_{\text{again}} = 0.002$ .

#### 1629 Appendix D.

1630 Fig. D1 lists the dissimilarity scores of some  
 1631 word pairs obtained by the ontological, statistical,

and combinational methods, respectively. The table shows the examples of the four cases in knowledge combination: (1) dissimilar by both ontology and corpus, as in  $\langle \text{any let} \rangle$  and  $\langle \text{center happen} \rangle$ ; (2) dissimilar by ontology but similar by corpus, as in  $\langle \text{already out} \rangle$  and  $\langle \text{take down} \rangle$ ; (3) similar by ontology but dissimilar by corpus, as in  $\langle \text{count figure} \rangle$  and  $\langle \text{rate speed} \rangle$ ; and (4) similar by both ontology

		<i>Ontology</i>	<i>Corpus</i>	<i>Combination</i>
1	<i>any let</i>	14.861	0.052	0.561
2	<i>center happen</i>	13.824	0.125	0.771
3	<i>already out</i>	16.898	0.003	0.456
4	<i>take down</i>	14.857	0.002	0.399
5	<i>count figure</i>	0	0.032	0.104
6	<i>rate speed</i>	1.536	0.028	0.131
7	<i>how when</i>	2.985	0.005	0.095
8	<i>Ones teeth</i>	2.985	0.009	0.108

Fig. D1. Examples of word dissimilarities by combined ontology and statistics.

1640 and corpus, as in *<how when>* and *<ones teeth>*. The  
 1641 combination occurs when the interpolation param-  
 1642 eter  $\lambda = 0.5$ . The examples show that the combina-  
 1643 tion compromises the errors caused by the  
 1644 specification and the generality of ontology and  
 1645 corpus, respectively.

## 1646 References

1647 Beckman, M.E., Ayers, G.M., 1994. Guidelines for ToBI  
 1648 Labeling. Available from: [http://www.ling.ohio-state.edu/](http://www.ling.ohio-state.edu/phonetics/ToBI/main.html)  
 1649 [phonetics/ToBI/main.html](http://www.ling.ohio-state.edu/phonetics/ToBI/main.html)>.  
 1650 Bolinger, D., 1961. Contrastive accent and contrastive stress.  
 1651 *Language* 37, 83–96.  
 1652 Bolinger, D., 1965. *Forms of English*. Harvard University  
 1653 Press, Cambridge, MA.  
 1654 Bosch, P., van der Sandt, R., 1999. *Focus: Linguistic, Cogni-*  
 1655 *tive, and Computational Perspective*. Cambridge University  
 1656 Press, Cambridge, UK.  
 1657 Chu-Carroll, J., Carpenter, B., 1999. Vector-based natural  
 1658 language call routing. *Comput. Linguistics* 25 (3), 361–388.  
 1659 Chomsky, N., 1971. Deep structure, surface structure and  
 1660 semantic interpretation. In: Steinberg, D., Jakobovits, L.  
 1661 (Eds.), *Semantics: An Interdisciplinary Reader in Linguistics,*  
 1662 *Philosophy and Psychology*. Cambridge University  
 1663 Press, Cambridge, UK.  
 1664 Dagan, I., Marcus, S., Markovitch, S., 1995. Contextual word  
 1665 similarity and estimation from sparse data. *Computer*  
 1666 *Speech Language* 9, 123–152.  
 1667 Dahl, D., 1969. *Topic and Focus: a Study in Russian and*  
 1668 *General Transformational Grammar*. Elandres Botryckeri,  
 1669 Göteborg.  
 1670 Daubechies, I., 1990. The wavelet transform, time-frequency  
 1671 localization and signal analysis. *IEEE Trans. Information*  
 1672 *Theory* 36 (5), 961–1005.  
 1673 Edmonds, P., Hirst, G., 2002. Near-synonymy and lexical  
 1674 choice. *Comput. Linguistics* 28 (2), 105–144.  
 1675 Firbas, J., 1964. On defining the theme in functional sentence  
 1676 analysis. *Travaux Linguistiques de Prague* 1, 267–280.  
 1677 Firbas, J., 1966. Non-thematic subjects in contemporary  
 1678 English. *Travaux Linguistiques de Prague* 2, 229–236.  
 1679 Flammia, G., 1998. *Discourse segmentation of spoken dia-*  
 1680 *logue: an empirical approach*. Ph.D. Thesis. MIT.  
 1681 Gorin, A.L., Abella, A., Alonso, T., Riccardi, G., Wright, J.H.,  
 1682 2002. Natural spoken dialog. *IEEE Computer Magazine* 35  
 1683 (4), 51–56.  
 1684 Gundel, J.K., Fretheim, T., 2001. Topic and focus. In: Horn,  
 1685 L., Ward, G. (Eds.), *The Handbook of Pragmatic Theory*.  
 1686 Blackwell Publishers, Malden, MA.  
 1687 Gussenhoven, C. 2002. Intonation and interpretation: phonet-  
 1688 ics and phonology. in: Bel, B., Marlien, I. (Eds.), *Speech*  
 1689 *Prosody 2002, Aix-en-Provence*.  
 1690 Halliday, M., 1967. Notes on transitivity and theme in English.  
 1691 Part II. *J. Linguistics* 3, 199–244.

Hedberg, N., Sosa, J.M., 2001. The prosody of topic and focus  
 1692 in spontaneous English dialogue. *LSA Topic and Focus*  
 1693 *Workshop*.  
 1694  
 1695 Heldner, M., Strangert, E., Deschamps, T., 1999. A focus  
 1696 detector using overall intensity and high frequency empha-  
 1697 sis. *Internat. Congress of Phonetic Science*.  
 1698  
 1699 Higgins, D., 2004. Which statistics reflect semantics? Rethink-  
 1700 ing synonymy and word similarity. *Internat. Conf. on*  
 1701 *Linguistic Evidence*.  
 1702  
 1703 Jackendoff, R., 1972. *Semantic Interpretation in Generative*  
 1704 *Grammar*. MIT Press, Cambridge, MA.  
 1705  
 1706 Jiang, J.J., Conrath, D.W., 1997. Semantic similarity based on  
 1707 corpus statistics and lexical taxonomy. *Proc. Internat. Conf.*  
 1708 *Research on Comput. Linguistics*.  
 1709  
 1710 Kadmon, N., 2001. *Formal Pragmatics*. Blackwell Publishers,  
 1711 Malden, MA.  
 1712  
 1713 Kay, M., 1975. Syntactic processing and functional sentence  
 1714 perspective. In: Schank, R., Nash-Webber, B. (Eds.),  
 1715 *Theoretical Issues in Natural Language Processing*. MIT  
 1716 Press, Cambridge, MA.  
 1717  
 1718 Kim, S.-S., 1998. Time-delay recurrent neural network for  
 1719 temporal correlations and prediction. *Neurocomputing* 20,  
 1720 253–263.  
 1721  
 1722 Kim, S.-S., Hasegawa-Johnson, M., Chen, K., 2003. Automatic  
 1723 recognition of pitch movements using multi-layer percep-  
 1724 tron and time-delay recurrent neural network. *IEEE Signal*  
 1725 *Process. Lett.* 11 (7), 645–648.  
 1726  
 1727 Krifka, M., 1999. Additive particles under stress. *Proc. of*  
 1728 *SALT 8*.  
 1729  
 1730 Kruijff-Korabayova, I., Steedman, M., 2003. Discourse and  
 1731 information structure. *J. Logic Language Inf.* 12, 249–  
 1732 259.  
 1733  
 1734 Lee, C., 1999. Contrastive topic: a locus of the interface. In:  
 1735 Turner et al. (Eds.), *The Semantics/Pragmatics Interface*  
 1736 *from Different Points of View (CRISPI 1)*. Elsevier Science,  
 1737 Amsterdam.  
 1738  
 1739 Lee, C., 2003. Contrastive topic and/or contrastive focus. in:  
 1740 McClure, B. (Ed.), *Japanese/Korean Linguistics 12, CSLI,*  
 1741 *Stanford*.  
 1742  
 1743 Lee, J.H., Kim, M.H., Lee, Y.J., 1993. Information retrieval  
 1744 based on conceptual distance in is-a hierarchies. *J. Docu-*  
 1745 *mentation* 49 (2), 188–207.  
 1746  
 1747 Lenci, A., 2001. Building an ontology for the lexicon: semantic  
 1748 types and word meaning. In: Jensen and Skadhauge (Eds.),  
 1749 *Ontology-based Interpretation of Noun Phrases: Proc. 1st*  
 1750 *Internat. OntoQuery Workshop*.  
 1751  
 1752 Lin, D., 1998. Automatic retrieval and clustering of similar  
 1753 words. *Proc. of COLING-ACL, Montreal, Canada*.  
 1754  
 1755 Miller, G., Fellbaum, C., 2002. *WordNet*. [http://www.cog-](http://www.cogsci.princeton.edu/~wn/)  
 1756 [sci.princeton.edu/~wn/](http://www.cogsci.princeton.edu/~wn/).  
 1757  
 1758 Munoz, M., Punyakanok, V., Roth, D., Zimak, D., 1999. A  
 1759 learning approach to shallow parsing. In *EMNLP-WVLC*  
 1760 '99.  
 1761  
 1762 Pantel, P., Lin, D., 2002. Discovering word senses from text.  
 1763 *ACM SIGKDD Conf. on Knowledge Discovery and Data*  
 1764 *Mining*.  
 1765



- 1748 Pierrehumbert, J., Hirschberg, J., 1990. The meaning of  
1749 intonational contours in the interpretation of discourse.  
1750 In: Cohen, P.R., Morgan, J., Pollack, M.E. (Eds.), *Inten-*  
1751 *tions in Communication*, pp. 271–311.
- 1752 Rada, R., Mili, H., Bicknell, E., Bletner, M., 1989. Develop-  
1753 ment and application of a metric on semantic nets. *IEEE*  
1754 *Trans. Systems Man Cybernet.* 19 (1), 17–30.
- 1755 Ren, Y., Kim, S.-S., Hasegawa-Johnson, M., Cole, J., 2004a.  
1756 Speaker-independent automatic detection of pitch accent.  
1757 *Internat. Conf. on Speech Prosody.*
- 1758 Ren, Y., Kim, S.-S., Hasegawa-Johnson, M., Cole, J., 2004b.  
1759 Speaker-independent automatic detection of pitch accent.  
1760 *Proc. ISCA Internat. Conf. on Speech Prosody.*
- 1761 Resnik, P., 1995. Using information content to evaluate  
1762 semantic similarity in a taxonomy. *Proc. of the 14th*  
1763 *Internat. Joint Conf. on Artificial Intelligence* 1, 448–453.
- 1764 Rooth, M., 1992. A theory of focus interpretation. *Natural*  
1765 *Language Semantics* 1, 75–116.
- 1766 Rulequest Research, 2004. *Data Mining Tools*. Available from:  
1767 <http://www.rulequest.com/see5-info.html>.
- 1768 Santorini, B., 1990. Part-of-speech tagging guidelines for the  
1769 Penn Treebank project. *Linguistic Data Consortium.*
- 1770 Sluijter, A., van Heuven, V.J., Pacilly, J., 1997. Spectral balance  
1771 as a cue in the perception of linguistic stress. *The J. Acoust.*  
1772 *Soc. Amer.* 101 (1), 503–513.
- 1773 Steedman, M., 2000. Information structure and the syntax–  
1774 phonology interface. *Linguistic Inquiry* 31 (4), 649–689.
- 1775 Sussna, M., 1993. Word sense disambiguation for free-text  
1776 indexing using a massive semantic network. *Proc. of the 2nd*  
1777 *Internat. Conf. on Information and Knowledge*  
1778 *Management.*
- Terra, E., Clarke, C.L.A., 2003. Frequency estimates for  
1779 statistical word similarity measures. *Proc. of the HLT-*  
1780 *NAACL.* 1781
- Thelen, M., Riloff, E., 2002. A bootstrapping method for  
1782 learning semantic lexicons using extraction pattern contexts.  
1783 *Proc. of the EMNLP.* 1784
- Umbach, C., 2004. On the notion of contrast in information  
1785 structure and discourse structure. *J. Semantics* 21 (2), 155–  
1786 175. 1787
- Vallduví, E., Vilkuna, M., 1998. On rheme and contrast. In:  
1788 Culicover, P., McNally, L. (Eds.), *Syntax and Semantics,*  
1789 *The Limits of Syntax*, Vol. 29. Academic Press, San Diego,  
1790 CA. 1791
- Welby, P., 2003. Effects of pitch accent position, type, and  
1792 status on focus projection. *Language Speech* 46 (1), 53–81. 1793
- Xu, Y., Xu, C.X., Sun, X., 2004. On the temporal domain of  
1794 focus. *Proc. ISCA Internat. Conf. on Speech Prosody.* 1795
- Yoon, T.-J., Chavarria, S., Cole, J., Hasegawa-Johnson, M.,  
1796 2004. Inter-transcriber reliability of prosodic labeling on  
1797 telephone conversation using ToBI. *Proc. ICSLP.* 1798
- Zechner, K., Waibel, A., 2000. DIASUMM: Flexible summa-  
1799 rization of spontaneous dialogues in unrestricted domains.  
1800 *Proc. of COLING.* 1801
- Zhang, T., 2004. *Spoken language understanding in an intelli-*  
1802 *gent tutoring scenario.* Dissertation, University of Illinois at  
1803 Urbana-Champaign. 1804
- Zhang, T., Hasegawa-Johnson, M., Levinson, S.E., 2004.  
1805 CHILDren’s emotion recognition in an intelligent tutoring  
1806 scenario. *Proc. ICSLP.* 1807
- Zubizarreta, M.L., 1998. *Prosody, Focus, and Word Order.*  
1808 MIT Press, Cambridge, MA. 1809  
1810