

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/271764570>

# Natural Language Processing: Past, Present and Future

Chapter · January 2013

DOI: 10.1007/978-1-4614-6018-3\_4

CITATIONS

3

READS

2,094

1 author:



[Deborah A. Dahl](#)

Conversational Technologies

81 PUBLICATIONS 784 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



W3C Multimodal Interaction Working Group [View project](#)



Autonomy Operating System for UAVs [View project](#)

# Natural Language Processing: Past, Present and Future

**Deborah A. Dahl**

Conversational Technologies

## **Abstract**

This chapter provides a broad discussion of the history of natural language understanding for both speech and text. It includes a survey of the general approaches that have been and are currently being applied to the goals of extracting the user's meaning from human-language inputs and performing useful tasks based on that analysis. The discussion utilizes examples from a wide variety of applications, including mobile personal assistants, Interactive Voice Response (IVR) applications, and question answering.

## **1. Introduction**

Enabling a computer to understand everyday human speech or ordinary written language, and do something useful based on that understanding has been a scientific goal at least since Alan Turing proposed that the ability to carry on a believable conversation could serve as a test of a truly intelligent machine in 1950 [1]. The difficulty of doing this task in its full generality has been consistently underestimated throughout the history of the field. However, in the past fifteen years, (and accelerating at an even more rapid pace during the last couple of years) significant progress has been made towards making natural language understanding (NLU) practical and useful.

This progress has not been based on any fundamental, new insights in how human language works. Instead, I would argue, the progress made in NLU is based on factors having to do with the engineering aspects of natural language processing, as opposed to scientific ones. Specifically, 1) recognizing the need for robust processing in the face of uncertain input; 2) identifying tractable tasks that are less ambitious than full, human-quality NLU; 3) network capabilities that allow systems to leverage the full power of distant servers 4) vast amounts of real data accessible over the Internet; and 4) Moore's Law which allows algorithms that were once impractically resource-intensive to be tested and put into practice.

Today, we have powerful personal assistants, like Apple's Siri, that respond to everyday types of natural language requests like checking the weather, setting up meetings, setting reminders and answering general knowledge questions, very much in the way that early researchers imagined so many years ago. These assistants are far from perfect – Siri makes plenty of mistakes – but they are good enough to be practical, and they are getting better. Let's look at the history of the technology behind these applications to see how this technology has made possible mobile personal-assistants, Interactive Voice Response (IVR) systems and other modern-day uses of NLU.

## 2. Beginnings

Natural language processing has been a topic of interest since the earliest days of computing. Early publications such as Claude Shannon's 1948 paper on information theory [2] proposed a statistical theory of communication that considered communication as a statistically-based process involving decoding of a signal; and in fact, some early work was done in the field following this model. However, Noam Chomsky's influential 1957 book *Syntactic Structures* [3] changed the fundamental direction of natural language processing research with its claim that the structure of natural language is inherently incapable of being captured by statistical processes. The following thirty years of work followed a path based on formal languages as the primary tool for addressing the problem of NLU. However, in the early 1990's statistics again came to the forefront of NLU. This was at least partly due to breakthroughs in speech-recognition systems, enabled through the use of statistics, such as [4], as well as the efforts to bring speech recognition and NLU together in programs such as the DARPA Spoken Language Program (1989-1994)<sup>1</sup>

## 3. The Process of Natural Language Understanding

All natural language processing systems take some form of language -- whether it's a spoken dialog, a typed input, or a text -- and extract its meaning. Some natural language processing systems go directly from words to meanings while others perform one or more levels of intermediate analysis. Systems that go directly from words to meanings are typically used to identify fairly coarse-grained meanings, such as classifying requests into categories or web search. A system that is retrieving the web pages that are relevant to a specific search doesn't need to do a detailed analysis of the query or the documents themselves. On the other hand, finer-

---

<sup>1</sup> See [5] for an introduction to the first workshop.

grained analysis with more levels of processing is usually used for tasks where the system has to understand exactly what the user has in mind. If I say “I need a flight from Boston to Denver on July 22 that arrives before 10 a.m.,” the system needs to extract the exact date, time and cities to provide an answer that satisfies the user.

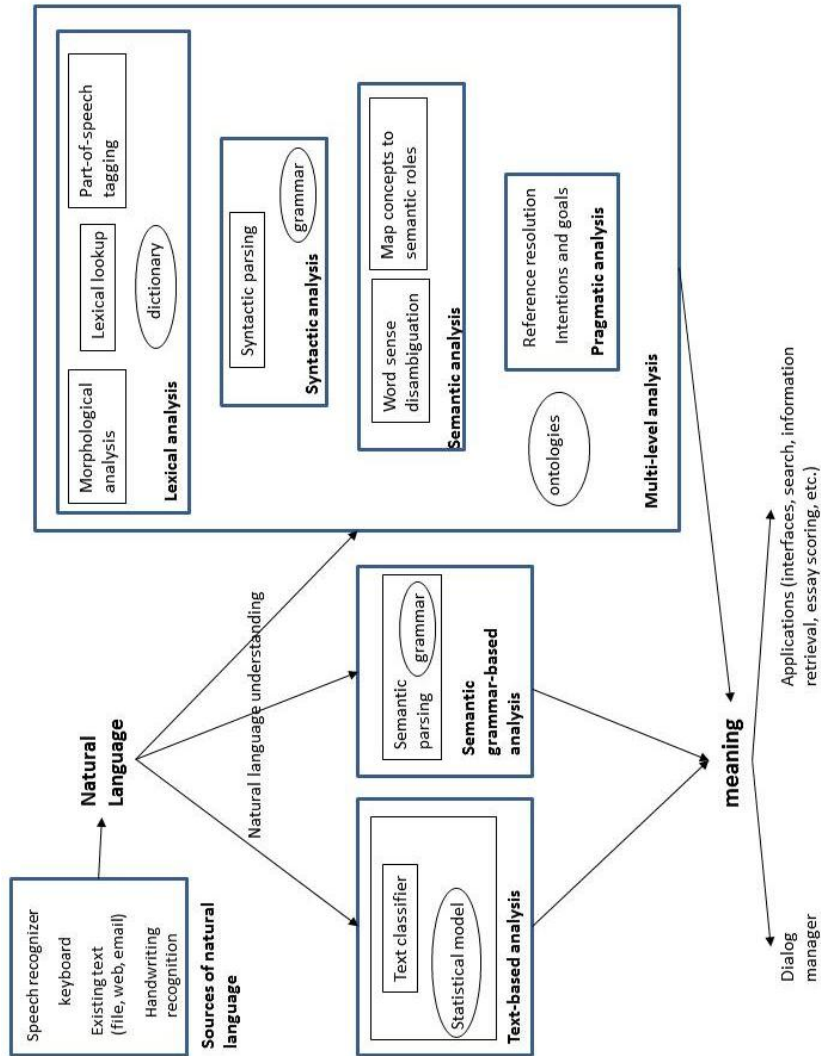
Figure 1 provides a broad perspective on the three main general approaches to NLU.

Natural language can come from many sources as shown in the upper left-hand corner of Figure 1 – speech recognition, a keyboard, handwriting recognition, or existing text, such as a file or web page. The goal is to find out what the meaning of that language is, where “meaning” is very broadly understood as some representation of the content of the language that is relevant to a particular application. Figure 1 shows three approaches to NLU, labeled *text-based analysis*, *semantic grammar-based analysis*, and *multi-level analysis*, which we will explore in detail in this chapter.

In text-based analysis, the basic unit of analysis is the text itself. Statistical models based on information such as the proximity of different words to each other in the text, the relative frequency of the words in that text and other texts, and how often the words co-occur in other texts are used to perform such tasks as web search and document classification.

In contrast, the other two approaches, semantic grammars and multi-level analyses, both attempt to define some kind of structure or organization of the text to pull out specific information that is of interest to an application. Semantic grammars look directly for a structure that can be used by an application; whereas multi-level analysis looks for multiple levels of intermediate structure that eventually result in a representation of the meaning of the text in a form that is useful to an application. (It should be noted that these are idealized systems; most actual systems contain elements of different approaches.)

As shown in the bottom of Figure 1, after the meaning is produced it can be used by other software, such as a dialog manager or another application. The rest of this chapter will discuss these components in detail, and will conclude with a discussion of integrating natural language processing with other technologies.



**Figure 1 Approaches to NLU**

### 3.1. Multi-level Analysis

We'll start by looking at the multi-level analysis approach.

Natural language processing systems which do a detailed analysis of their inputs traditionally have components that are based on subfields of linguistics such as the lexicon, syntax, semantics, and pragmatics. The relative importance of these components in processing the language often depends on the language. For example, analyzing written text in languages that don't have spaces between words, such as Chinese or Japanese, often includes an extra process for detecting word boundaries. Processing can be done sequentially or in parallel, depending on the architecture of the system. Many implemented systems also include some aspect of probability. That is, how to analyze an input may be uncertain when the input is analyzed, but if one of the analyses is more likely, the less likely analyses can be either eliminated or explored at a lower priority. For example, "bank" in the sense of a financial institution is a more likely meaning in most contexts than the verb "bank" in the sense of piling up a substance against something else.

#### 3.1.1. Lexical lookup

Starting from either a written input or the output of a speech recognizer, lexical lookup describes information about a word in the input. It may include a step of *morphological analysis* where words are taken apart into their components. For example, the English word "books" can be analyzed as "book" + "plural." This is especially important for languages where words have many forms depending on their use in a sentence (highly inflected languages). Spanish, for example, has more different word forms than English, and a word like "hablaremos" would be analyzed as "speak" + "future" + "first person" + "plural" or "we will speak." There are many other languages that are much more complicated than Spanish, and it would be very impractical to list each possible word in a dictionary for these languages. So these words need to be broken into their components.

Related to morphological analysis is a process called *part of speech tagging*, which identifies a word as a noun, verb, adjective, or other part of speech (see [5] for an example). This process provides extremely useful information, especially for words that can be used in many different contexts.

The English word "like" is a good example of a word that can occur as at least six different parts of speech, as shown in Table 1.

**Table 1 Parts of Speech for "like"**

Example of usage	Part of Speech
I like that	Verb

Her likes and dislikes are a mystery	Noun
He was, like, eight feet tall	Interjection
People like that drive me crazy	Preposition
They are of like minds about that	Adjective
Your food is cooked like you wanted	Conjunction

- 1.1.1. Automatically identifying the part of speech of a word is helpful for later stages of processing, such as parsing and word-sense disambiguation, which we will discuss below, because it eliminates some analysis options. If the system knows that “like” is a verb in a particular sentence, then it can rule out any other possible analysis that uses “like” as a noun.

### 3.1.2. Parsing

Parsing is a stage in natural language processing which breaks down a sentence into its components and shows how they’re related to each other. Parsing can have the goal of finding either syntactic or semantic relationships within an utterance. Syntactic parsing is the older approach, and has been explored in a large body of research since early papers such as [6] [7, 8].

#### 3.1.2.1. Syntactic Parsing

Syntactic components include parts of speech and phrases, but not the meanings of those words or phrases. Rather, syntactic analysis is based on a set of rules defining the structure of the language. This set of rules is *called a syntactic grammar*. Figure 2 shows an example of a simple syntactic grammar that could analyze English sentences like “the cat sleeps on the chair.”<sup>2</sup>The first rule states that a sentence consists of a noun phrase (NP) followed by a verb phrase (VP), and the following rules describe how noun phrases and verb phrases are built, until we get to the actual words (dog, cat, table, etc.) or terminal symbols in the grammar. In this example parenthesized components are optional and alternatives are indicated by “|.” Full syntactic grammars for actual human languages are obviously much more complex.

---

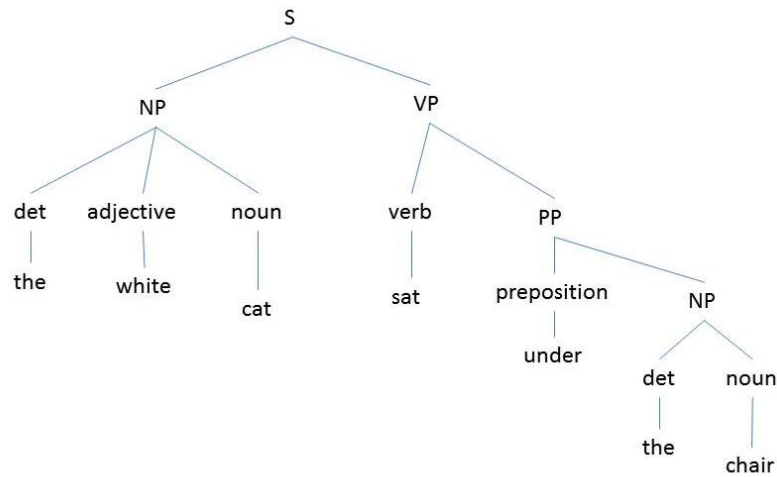
<sup>2</sup> For simplicity, this example is a context free grammar (CFG) in the terminology of formal languages, although normally a syntactic grammar of a natural language would be at least as powerful as a context-sensitive grammar. A commonly used syntax for writing context-free grammars is Backus-Naur Form or Backus Normal Form (BNF), invented by John Backus.

**Figure 2 A simple syntactic grammar**

```
Sentence → NP VP
NP → Det (Adj) N (PrepP)
VP → V (NP) (PrepP)
Det → the | a
Adj → blue | white | green | black
PrepP → Prep NP
Prep → under | on | in | behind
N → dog | cat | table | chair
V → sleeps | sits
```

Figure 3 shows a syntactic analysis for “the white cat sat under the chair.” Because a syntactic analysis doesn’t take into account the meanings of the words, we would get the same syntactic analysis for sentences like “the black dog slept behind the sofa,” as we would for “the white cat sat under the chair.” This is because both sentences have the same syntactic structure, even though they have entirely different meanings. The advantage of an approach that uses syntactic parsing is that the process of syntactic analysis can be decoupled from the meanings of the words that were spoken or even from the domain of the application. Consequently, the syntactic grammar of a language can be reused for many applications. On the other hand, one disadvantage of syntactic analysis is that it is based on a general, domain-independent grammar of a language. Such grammars require a great deal of work to put together even with modern machine learning techniques. However, most applications can be useful without a general grammar. Therefore, other techniques have been developed, most notably parsing approaches that do take into account the semantics of the domain, and whose output is based on semantic relationships among the parts of the utterance.



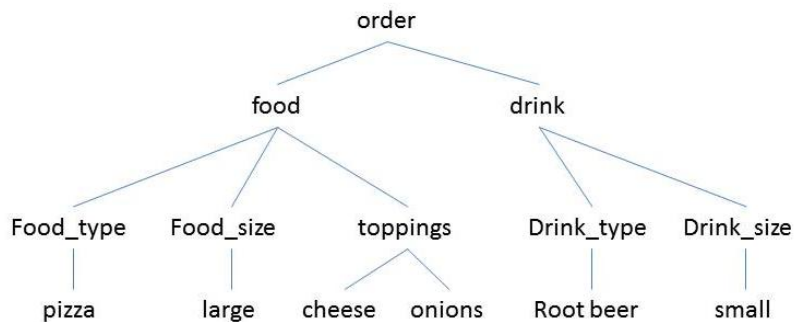


**Figure 3 Syntactic analysis for "the white cat sat under the table"**

### 3.1.2.2. Semantic Parsing

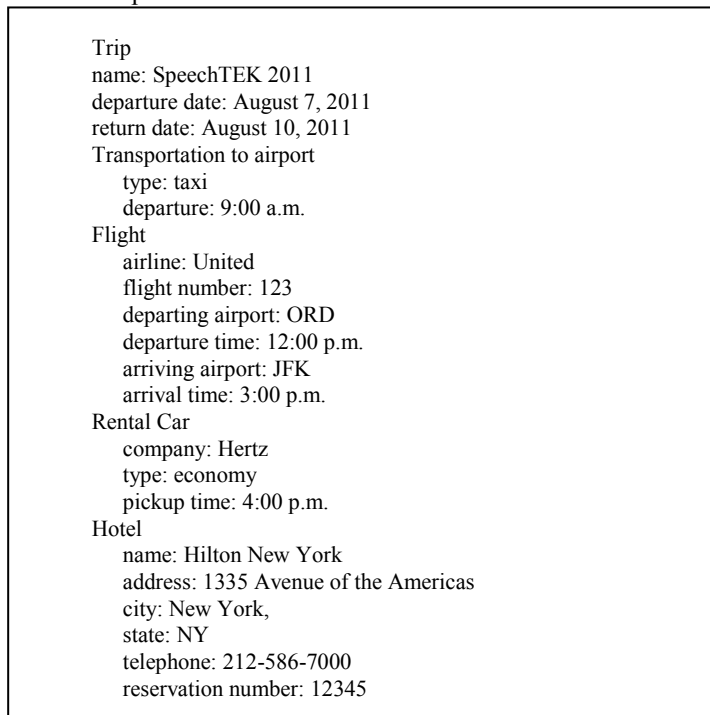
Semantic parsing analyzes utterances in the context of a specific application: ordering fast food. Figure 4 represents the categories in that application: the type of food, the type of drink, toppings to be put on the pizza, and so on. There is no syntactic information, such as the fact that “pizza” is a noun, or that “large” is an adjective. This is also a less general approach than syntactic parsing in that every new application needs a new grammar. It is, however, generally much faster to develop a one-time, semantic grammar for a single application than to develop a general syntactic grammar for an entire language.

Semantic parsing is also popular for speech systems because the grammar can be used to constrain the recognizer by ruling out unlikely recognition results. Using a grammar to constrain speech recognition supports the fast processing required by the real-time nature of speech recognition. In practice, this means that a grammar used to constrain speech recognizers has to be more computationally tractable than grammars used to analyze text. Speech grammars are always either finite state grammars (FSG) or context free (CFG), as defined in [9].



**Figure 4 Semantic parse for "I want a large cheese pizza with onions and a small root beer"**

The result of semantic parsing is a semantic frame, a structured way of representing related information which is popular in artificial intelligence. [10]. Figure 5 shows a complex semantic frame for travel information.



**Figure 5 A semantic frame for travel information**

**VoiceXML [11], a widely used language for defining IVR applications, uses semantic frames (or *forms*, in VoiceXML terminology).**

Figure 6 shows a VoiceXML form with *fields* (which correspond to the slots of a semantic frame) which will be filled by the information that the user provides for a card number and expiration date.

```

<form>
  <prompt>Welcome to the electronic payment system.</prompt>
  <field name="card_number">
    <prompt> Please enter your credit card number? </prompt>
    <grammar
src="http://www.ajax.com/credit_card_number.grxml"/>
  </field>
  <field name="date">
    <prompt>Please enter your expiration date </prompt>
    <grammar
src="http://www.ajax.com/credit_card_date.grxml"/>
  </field>
</form>

```

**Figure 6 A VoiceXML form with "card\_number" and "date" fields**

Semantic parsing first became popular in the early 1990's as a relatively quick way to get a speech system running. Examples of this approach include [12], [13], and [14].

All current commercial grammar-based, speech-recognition systems use semantic parsing.

As speech-recognition systems began to mature during the 1990's, the need for standard ways to write grammars became apparent. Initially, every recognizer had its own format, which made it extremely difficult to use a different recognizer in an existing system. Tools like the Unisys Natural Language Speech Assistant were developed to allow grammars to be authored in a recognizer-independent fashion with a graphical tool that would generate multiple grammars in the various formats. There were also a number of efforts to develop open grammar formats that could be used by multiple recognizers. These included Microsoft's Speech Application Programming Interface (SAPI) grammar format and Sun's Java Speech Grammar Format (JSGF) format. The JSGF format was contributed to the World Wide Web Consortium's (W3C's) Voice Browser Working Group in 2000 and became the basis of the ABNF format of the W3C's Speech Recognition Grammar Format (SRGS) specification [15], which became a formal standard in 2004. Because Extensible Markup Language (XML) [16] was rapidly increasing in popularity at this time, the SRGS specification also defines an XML version of the grammar standard. While the ABNF format is more compact than the XML for-

mat, the XML format is much more amenable to machine processing since there are many tools available for editing and validating XML documents.

Figure 7 shows an XML SRGS grammar rule for a fast-food order that would enable a recognizer to recognize sentences like “I would like a coke and a pizza with onions.” The `<ruleref>` tags point to other rules that aren’t shown here that recognize the different ways of asking for a drink (“#drink” and the different ways of describing a pizza (“#pizza”).

**Figure 7 An SRGS rule for a fast food order**

```
<rule id="order">
  I would like a
  <ruleref uri="#drink"/>
  and
  <ruleref uri="#pizza"/>
</rule>
```

The existence of a standard grammar format for speech recognizers made it possible to use grammars to constrain recognition in a vendor-independent way, but that didn’t solve the problem of representing the meaning of the utterance. To address that need, SRGS provides for inserting semantic tags into a grammar that would do things, for example, like expressing the fact that whatever was parsed in the “drink” rule should be labeled as a drink. However, SRGS doesn’t define a format for the tags. Another W3C standard, Semantic Interpretation for Speech Recognition (SISR) [17] defines a standard format for semantic tags that can be used within an SRGS grammar. Figure 8 shows the rule from Figure 7 with semantic tags. The tags are written in ECMAScript 237 [18], a standardized version of Javascript. This rule is essentially building a semantic frame that includes “drink” and “pizza” slots. The “drink” slot in turn has slots for the liquid and size of the drink. So, the reference to “out.drink.liquid,” for example, means that the “liquid” value of the “drink” frame will be filled by whatever matched the drink in the user’s utterance. If the user said “Coke” that value would be “Coke,” if the user said “lemonade,” the value would be “lemonade,” and so on.

**Figure 8 SRGS rule with SISR semantic tags**

```

<rule id="order">
  I would like a
  <ruleref uri="#drink"/>
  <tag>out.drink = new Object();
    out.drink.liquid=rules.drink.type;
    out.drink.drinksiz=rules.drink.drinksiz;
  </tag>
  and
  <ruleref uri="#pizza"/>
  <tag>out.pizza=rules.pizza;</tag>
</rule>

```

The semantic frame that is generated by this rule is the final result of NLU in the semantic-parsing paradigm. It is ready to be acted upon by an application to perform a task such as an interaction with an IVR (e.g. ordering fast food or making travel plans) or a web search. The W3C EMMA (Extensible MultiModal Annotation) specification [19] provides a standard way of representing the output semantic frame as well as other important annotations, such as the time of the utterance and the processor's confidence in the result.

We've brought the utterance through speech recognition and semantic analysis, to a final representation of a meaning that can be used by an application. (How natural language results can be used in an application is something we'll address in a later section.)

, At the beginning of the parsing section we described another approach to parsing: syntactic parsing. Looking back at Figure 3, it is clear that a syntactic analysis is not at all ready to be used by an application. So let's return to the syntactic parse in Figure 3 and talk about what other steps need to be taken to finish getting the meaning from the utterance once the syntactic parsing has been accomplished. Once we have a syntactic analysis of the input, the next step is semantic interpretation.

### 3.1.3. Semantic Analysis and Representation

The process of semantic interpretation provides a representation of the meaning of an utterance. In the semantic-parsing approach discussed above, the processes of looking at the structural relationships among words and deciding the overall meaning of the utterance were not differentiated. This can be efficient, especially for simpler applications, and as we have said, this is the way that all current grammar-based, speech-recognition applications work. However, it is also possible to separate syntactic analysis from semantic interpretation. This has been done

in research contexts, in some earlier commercial systems [20], as well in some very new systems such as IBM’s Watson [21].

### 3.1.3.1. Representation

We start with the goal of semantic analysis: obtaining the meaning of an utterance or text. We know what texts and utterances are, but what does a “meaning” look like? We saw one example in Figure 5, a semantic frame with slots and fillers (or attribute/value pairs) like “destination: New York.” This is still a very common type of representation. However, many other types of semantic representations have been explored in the past. There have been a number of approaches based on formal logic, for example the research system described in [22] and the commercial system described in [23]. In these systems meanings are expressed as logical expressions. For example, “the flight from Philadelphia to Denver has been cancelled” might be expressed as the following

$$\exists(x) (\text{flight}(x) \wedge \text{from}(\text{Philadelphia}, x) \wedge \text{to}(\text{Denver}, x) \wedge \text{cancelled}(x))$$

This is read “There is something, “x,” which is a flight, and which is from Philadelphia and is to Denver and is cancelled.”

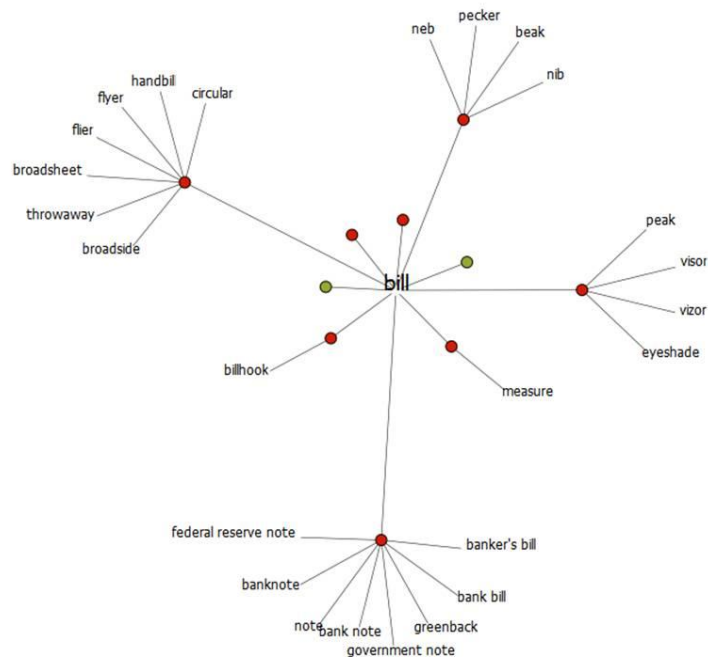
Another interesting type of semantic representation is similar to semantic frames, except that the slot names are application independent. These are often called *case frames*. For example, in a sentence like “send an email to Richard” the subject of the verb “send,” that is, the understood subject of the command, is classified as an “agent” slot because the subject is acting. The email is classified as a “theme,” and the recipient is assigned to the slot “goal.” The idea of case frames for semantic representation originated in Fillmore’s work [24] and was later elaborated in the work of Levin [25], which presents a detailed analysis of hundreds of English verbs. This approach supports very generic, application-independent systems because the case frames themselves are application independent. [20, 26, 27] are examples of systems that used this approach. On the other hand, the disadvantage of this approach is that, because the slots are application-independent, they still need to be associated with application-specific slots before they can be integrated with an application.

### 3.1.3.2. Word-Sense Disambiguation

Another important aspect of semantic processing is word-sense disambiguation (WSD). Many words have more than one meaning, a phenomenon that is called “polysemy.” For example, “bill” can refer to something you pay, or a bird’s beak. A “tie” can refer to something men wear around their necks or to a game where both teams have the same score, and so on. In order to unambiguously represent the meaning of an utterance or text, the correct senses of words need to be assigned. WSD is especially important in machine translation, because words that

are polysemous in one language usually must be translated into two different words. For example, the word “hot” in English can refer to temperature or spiciness, but Spanish uses two different words, “caliente” and “picante,” for the two concepts. The strategy for WSD is to examine the context around the polysemous words, rule out senses that are impossible in that context, and then select the sense that is most probable in that context from the remaining senses.

Word sense disambiguation often relies on a resource called an *ontology*. An ontology is a structured representation of concepts and their relationships, and defines what the senses are for a particular word. A well-known example of an ontology is WordNet [28], developed at Princeton University. WordNet was originally developed for English, but WordNets for a number of other languages have been developed. Figure 9 shows the senses in WordNet for the word “bill.” As Figure 9 shows, “bill” has six senses in WordNet: a bird’s beak, a handbill, the brim of a hat, a banknote, a billhook, or legislation (not counting the proper name “Bill.”) A WSD component would have the task of assigning one of those senses to a particular occurrence of “bill” in an input.<sup>3</sup>



**Figure 9 The WordNet senses for "bill"**

Fortunately, the general problem of sense disambiguation can be avoided in many applications. This is because either the topics are not broad enough to in-

<sup>3</sup> The WordNet visualization shown in Figure 9 was created using the Google Code project “Synonym”.

clude multiple senses, or the number of polysemous words in the application is small enough that the appropriate contextual information can be hand coded. For example, I could ask my mobile personal-assistant a question like “what are the times of my next two appointments” or “what is five times three.” The word “times” has two different senses in those two requests, but the senses can be distinguished because they occur in different contexts. In a personal-assistant application, the application is specific enough that the contexts can be hand coded. In this example the developer could simply specify that if the words on either side of the word “times” are numbers, “times” means “multiplication.” Work on the more general problem of WSD in broader domains such as newspapers, translations, or broadcast news relies heavily on automated ways of acquiring the necessary contextual information.

To sum up, at the end of the semantic processing phase, we have an exact description of the meaning of the input. It might be represented in any of a number of ways: as an application-specific semantic frame, a logical form, or an application-independent set of case roles, among others, and the senses of polysemous words have been disambiguated.

The next stage of processing is pragmatic processing. As with syntactic analysis and semantic interpretation, not all systems perform pragmatic analysis as a distinct step.

### 3.1.4. Pragmatic Analysis and Representation

Pragmatics is the subfield of linguistics that deals with the relationship of language to its context. By “context” we mean both the linguistic context (i.e., what has been said before in a dialog or text) as well as the non-linguistic context, which includes the relationship of language to the physical or social world. If I point to something, and say “I like that,” the understanding of both “I” and “that” depends on the state of affairs in the world: who is speaking and what they’re pointing to. Moreover, use of the present tense of the verb “like” ties the speech to the current time, another aspect of the non-linguistic context.

#### 3.1.4.1. Reference Resolution

An important and unsolved problem in NLU is a general solution to understanding so-called *referring expressions*. Referring expressions include pronouns such as “I” and “that;” *one*-anaphora, as in “the blue one;” and definite noun phrases, such as “the house.” This task is called *reference resolution*. Reference resolution is the task of associating a referring expression (“I,” “he,” “the blue one” or “the house”) with a *referent*, or the thing that’s being referred to. Reference resolution is difficult because understanding references can require complex, open-ended knowledge. As in other areas of NLU, the need for a general solution to reference resolution has been finessed in practice by addressing simpler, less



general, but nevertheless useful problems. For example, in an IVR application, the system never really has to interpret “I,” even though it is used all the time (“I want to fly to Philadelphia”) because there is never more than one human in the conversation at a time. Other pronouns are rarely used in IVR applications. You could imagine something like “I want to fly to Philadelphia. My husband is coming, too, and he needs a vegetarian meal.” If the user did say something like that the IVR would need to figure out what “he” means, and that one passenger on this reservation needs a vegetarian meal. However, in practice, speech directed at an IVR is much simpler and consequently the system rarely has to address interpreting pronouns.

Ontologies, as discussed above, also provide useful information for pragmatic analysis because they represent conceptual hierarchies. Some references can be interpreted if we know what kind of thing a word refers to. For example, knowing that “Boston” is a city provides the information needed to know that “the city” in “If I fly into Boston, what’s the best way to get into the city?” refers to Boston.

Pragmatic analyses in commercial systems are normally represented in semantic frames where any context-dependent references have been resolved. For example, a user might say “I want to schedule an appointment for tomorrow” instead of a specific date. Because “tomorrow” is a word that must be interpreted with information from the non-linguistic context, pragmatic processing has to identify the actual date that “tomorrow” refers to. The final semantic frame would then include the specific date for the appointment, rather than just the word “tomorrow.”

#### 3.1.4.2. Named Entity Recognition

Another example of tying language to the world is in the task of *named entity recognition*, or identifying references to people, organizations or locations through textual descriptions. Named entity recognition is a type of reference resolution where the referent is an actual individual, place or organization. The descriptions can be extremely diverse, but if an application needs to associate events and activities to an individual, it’s important to identify the individual, no matter how the reference is expressed. For example, someone might refer to Barack Obama as “the President” (assuming that we know we’re talking about the United States and we’re talking about the current president), “the Commander in Chief,” “Mr. Obama,” “he,” or more indirectly, as in “the winner of the 2008 presidential election,” or “the author of *Dreams from my Father*.” This is a very active research area, and researchers are looking at a number of interesting questions, such as how to recognize named entities in tweets [29].

#### 3.1.4.3. Sentiment Analysis

Sentiment analysis is a new and important application of natural language processing that looks at an aspect other the literal meaning, or *propositional content*,

of an utterance or text. All of the types of processing that we've talked about so far have addressed the goal of extracting the literal meaning from natural language. In contrast, the goal of sentiment analysis is to characterize the speaker or writer's attitude toward the topic of the text. As reviews of products and businesses proliferate on the Web, companies that are interested in monitoring public attitudes about their products and services are increasingly turning to automated techniques such as sentiment analysis to help them identify potential problems. Sentiment analysis tries to classify texts as expressing positive, negative, or neutral sentiments, and can also look at the strength of the expressed sentiment. Sentiment analysis of written texts is technically a type of text classification, which will be discussed in the next section in detail. However, in sentiment analysis, the classification categories have to do with attitudes rather than specific topics. Initial work on sentiment analysis in text is described in [30]. Sentiment analysis can also be done using spoken input, using information such as prosody, which is not available in texts. For example, [31] describes using prosody to detect sentiments in spoken interactions.

### ***3.2. Text Classification***

Looking back at Figure 1, we note that we haven't really touched on the text-based approaches to NLU. As we said in the discussion of Figure 1, text-based approaches map inputs fairly directly to meaning, without going through the levels of intermediate analyses that the semantically based or the multi-level approaches perform.

One way to think about text classification is that the goal is to take some text and classify it into one of a set of categories, or bins. Ordinary web search is a kind of text classification. In the case of web search, the bins are just "relevant to my search query" or "not relevant to my search query." The classification result is assigned a score (used internally) so that higher scoring, and presumably more relevant, web pages are seen first by the user. There are many text-classification techniques available, primarily based on machine-learning methods. For example, Naïve Bayes, vector-space classifiers, and support-vector machines are used in text classification, to name only a few. This area is a very active field of research.

Text classification, in combination with statistical speech recognition based on statistical language models (SLM's), has become very popular in the last ten years as a tool that enables IVR systems to accept more open-ended input than is typically possible with hand-constructed, semantic grammars. Unlike semantic grammar-based systems, the speakers' utterances do not have to match exactly anything that was directly coded during system development. This is because the matching of text to bin is not all or none, but statistical. Combining text classification with statistical-language models of speech was first proposed in [32] and has become very successful. Users are typically much more satisfied with systems that allow them to express themselves in their own words.

As these systems have been deployed in IVR systems and other spoken-dialog systems, a number of refinements in best practices have been learned. For example, users have more success if the system's opening prompt is not as open ended as "How may I help you?" because users may not understand how to respond to this kind of very open prompt. A more constrained prompt, such as "Please tell me the reason for your call today" is usually more effective.

Because these are statistically-based systems, a drawback to SLM systems is that they require collection of significant numbers of the utterances that are used to train the system, up to tens of thousands in some cases. Moreover, not only must these training utterances be collected, but they must also be manually classified into their appropriate categories by human annotators. This is because the system develops the statistical preferences that it will use to categorize future utterances on the basis of human-annotated data. Training based on human annotation, or *supervised training*, is an expensive procedure. For this reason, training with little or no attention from human annotators, called *unsupervised training*, or *weakly supervised training*, is an important goal of work in this area, although the problem of effective unsupervised training is far from solved.

However, once trained, these systems can be very accurate. The expense of human annotation can be cost-effective in some larger-scale applications, if the alternative is sending the caller to a human agent. Figure 10 shows an example of how accurate these systems can be [33], even on very indirect requests.

**Figure 10 Correct processing of an open-ended user request in an IVR**

User: "I've been on in and out of the hospital and I know I'm late on it and I'm... I'm... I'm wondering, I'm out of the hospital now and they finally took my cast off, but I still can't work and I can't walk and I'm wondering...."

Classified as "Caller would like to get an extension on paying his utility bill"

Commercial systems based on this technology are often referred to as "natural language systems," because they can effectively process users' unconstrained, natural language, inputs. However, as we have seen in this chapter, natural language systems are much more general than this specific technology.

#### 4. Summary of approaches

We have reviewed three general approaches to NLU: multi-level approaches semantic parsing approaches and text-based approaches.

1. The multi-level approaches include several levels of linguistically-based analysis, each building on the previous level. These include lex-

ical analysis, syntactic parsing, semantic analysis, and pragmatic analysis. The claim of these systems is that by developing a set of application-independent resources (dictionaries, grammars, semantic information and ontologies), the task of developing new applications can be greatly simplified. In practice, however, the application-independent resources on which these systems are premised have proven to be extremely expensive and time consuming to develop. Organizations with extensive development capabilities can still create these kinds of systems. For example, the Watson Jeopardy-playing system implemented by IBM is a multi-level system [21]. Unlike the Watson project, most natural language processing application-development efforts have constrained budgets and cannot afford to develop these resources on their own. In a few cases government funding has enabled the creation of shared resources. Comlex (Common Lexicon) [34] and WordNet [28] are notable examples. They are exceptions because in general the required resources are not widely available and must be constructed by each organization.

2. Semantic-grammar based approaches were particularly useful for early speech applications, through the 1990's and early 2000's, because semantic grammars (an example can be seen in Figure 8) are sufficient to process the utterances that were found in limited domains, such as banking or air travel planning. In addition, the semantic grammar serves a useful role in constraining the speech recognizer so that it will only recognize utterances that are appropriate to the application. This significantly improves the accuracy of speech recognition. Semantic grammars are, however, difficult to maintain, especially as the complexity of the application increases. Nevertheless, the vast majority of current IVR applications use this approach. Fortunately, most IVR applications do not require complicated grammars, making this approach highly effective for IVR applications.
3. Text-based approaches became popular in speech applications in the early 2000's, as developers realized that more natural input to IVR's was highly desirable. It was impossible to create semantic grammars broad enough to recognize this more open input, so the text-based approaches came into general use. The large amount of annotated training data that these systems require makes them expensive to build and maintain. This is particularly true if the data changes dynamically, which is the case for seasonal applications. A seasonal retail application, for example, needs new data for each new product added to the application because new products introduce new words for users to say

Clearly, no single approach is ideal. Each application has its own goals and requirements, making some approaches better for other applications than others. Limited applications like IVR's do well with semantic grammar approaches. Multi-level systems are a good approach for very broad question-answering systems

that require a detailed analysis of the questions, like IBM's Watson. Text-based systems are good for classification tasks that require only a general understanding of the input.

Many current systems are hybrids, and incorporate techniques drawn from several of the generic approaches. Mobile personal assistants like Apple's Siri, for example, make use of multiple techniques. Text-processing techniques enable mobile personal-assistants to work with wide-ranging input on unpredictable topics such as web searches from millions of different users. On the other hand, in many cases the inputs to mobile personal-assistants require more detailed understanding of the user's request. A request that includes a specific date or time needs to be analyzed in detail so that the date or time is handled correctly. A semantic grammar that parses dates and times is the perfect tool for this. Even simple word-spotting can be used, although sometimes that produces incorrect results. For example comment to Siri such as "I need \$100" gets the response "Ok, I set up your meeting for tomorrow at 1 p.m." Clearly Siri must only be paying attention to the word "one" in that query. Clearly these mobile personal-assistants use an eclectic mix of techniques because of the many different types of conversations they have with their users.

## 5. Methodology -- Getting data

All natural language based systems are based on data. In a multi-level system the data may be in the form of dictionaries or syntactic grammars. In a semantics-based system the data may take the form of a semantic grammar. Text-classification systems rely on associations between texts and their classifications (training data) which allow them to classify new texts based on their resemblance to the training texts. Similarly, any kind of system that makes use of probability will derive its probabilities from training data. Early systems used data hand coded by experts, which was time consuming and expensive. As machine learning became more sophisticated, many systems began to use training data that was annotated with the correct analysis by humans without using data that was explicitly hand coded by experts. Human annotators, while expensive, are much less expensive (and more available) than grammar experts. For example, an extensive annotation effort at the University of Pennsylvania, Treebank [35], provided a large set of syntactic parses prepared by humans which were intended to be used in machine learning of parsing techniques. A similar effort, PropBank [36], added semantic case-frames to the Treebank data. Treebank and PropBank represent the supervised approach to annotation. As discussed earlier in the section on text-based approaches, unsupervised approaches require less attention from human annotators but much research needs to be done before unsupervised techniques are good enough for widespread use. At this point, the general problem of data acquisition has not yet been solved.

## 6. “Frequently Bought With”

Natural language processing can be part of many other types of systems and often serve as only one component of a complete system. Here we review some of the other components that are often combined with natural language processing. We will focus on interactive dialog systems, like mobile personal-assistants.

Figure 11 shows the complete architecture for a typical spoken-dialog system. As Figure 11 shows, the natural language processing component is only one part of the larger system.

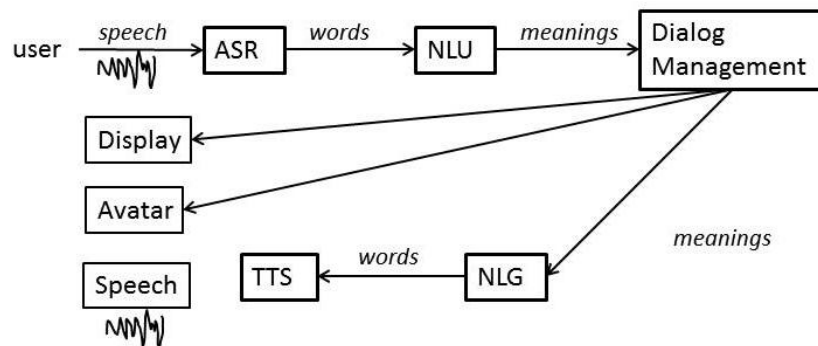


Figure 11 A generic interactive spoken-dialog system

### 6.1. Speech recognition

Early work on dialog systems with spoken input was part of the DARPA Speech Understanding Program (SUR) of the 1970's [37-40]. These were strictly research projects, since the speech recognition of the time was too slow and inaccurate for practical applications.

In the early 1990's speech recognition started to improve dramatically. This improvement was stimulated by two factors. One factor was a technical breakthrough: the use of Hidden Markov Models [41]. The second was the series of formal evaluations of recognition accuracy conducted by NIST [42], which helped researchers understand how each specific algorithmic improvement contributed to overall recognition accuracy. These improvements made possible the development of speech-enabled IVR applications, which continue to be very successful.

At the same time, the formerly favored multi-level approaches [43, 44] were being replaced by the less resource-intensive, semantic-parsing approach [12]. As discussed earlier, semantic parsing-based methods work best in limited applications, such as checking on banking information. This is because in these systems, every possible input has to be anticipated by the developers. So-called *out of grammar* or *out of vocabulary* utterances cannot be processed. If the user says something that the developer had not anticipated, the system has to engage the user in a tedious dialog to try to get the user to say something it knows how to process.

In order to support more general applications, such as personal assistants, speech recognition has to be able to accept much less constrained inputs. Fortunately, while parsing-based IVR applications were spreading, the technology needed for recognizing less constrained inputs was being developed independently in the context of speech recognizers used for dictation (e.g., Dragon Dictate and Dragon Naturally Speaking). Dictation systems use Statistical Language Models (SLM's) to define the expected possibilities of words in a user's utterance, rather than grammars. These possibilities are expressed as word pairs (*bigrams*); triples (*trigrams*); or more generally, as *N-grams*. SLM's contain information, such as the fact that the sequence "the cat" is more probable than the sequence "the it." Because this information is probabilistic rather than absolute, recognizers using the SLM approach are more flexible than grammar-based recognizers for recognizing unexpected input.

Dictation technology has continued to improve as it is applied to tasks like web search, which allows for the collection of vast amounts of data from millions of users. The result is now that dictation speech recognition works reasonably well in the context of spoken-dialog systems such as mobile personal-assistants, although factors like noise and accents still affect recognition accuracy.

## 6.2. *Multimodal inputs*

Devices that include a display, keyboard, touchscreen and/or mouse enable the user to interact with the device in ways other than voice. This style of computer-human interaction is called *multimodal interaction*. Multimodal interaction has been a research topic for many years (see [45-47] for early work). However, several factors prevented this early research work from being widely used in commercial systems. One major factor was that speech recognition was not as accurate as needed to support seamless multimodal interaction (error correction was a constant distraction from the user's goals); another was that, for many years, spoken input was limited to a few specific situations. For example, in telephone-based, IVR applications, the alternative is touchtones, which are even more cumbersome than speech.

Another type of application where even error-prone speech recognition made sense was where the user was, for some reason, unable to use a mouse or key-

board. This included users who used speech as an assistive device or users in hands-busy situations. Now, we have both much better speech recognition as well as powerful small devices with keyboards that are difficult to use. The combination of better speech recognition with small keyboards makes spoken and multimodal interaction much more appealing than in the past. In a mobile application like Siri, users can either speak a request, or in some cases, interact with Siri using a touch alternative. The touch alternative is especially useful for tasks such as confirming or canceling a request. AT&T's Speak4it multimodal mobile assistant also supports simultaneous speech and drawing input. For example, Speak4it allows a user to draw a circle on a map while saying "Show me Italian restaurants around here."

In addition to enabling input combining spoken interaction with touchscreens, today's mobile devices routinely include other capabilities that provide additional opportunities for multimodal interaction. These include cameras, accelerometers, and GPS technology. There are also special-purpose sensors that can be added to mobile devices, such as glucose meters or blood pressure meters. These special-purpose sensors provide even more opportunities for multimodal interaction.

Because the number of different modalities continues to increase, it is important to have generic, modality-independent ways of representing inputs from a wide range of modalities. The W3C's Extensible Multimodal Interaction (EMMA) specification [19] provides a way to manage inputs from an open-ended set of modalities. In order to do this, EMMA defines a uniform standard for representing inputs from any modality, whether it is speech, keyboard, touchscreen, accelerometer, camera, or even future modalities. The meaning of the input is represented in the same way, independent of the modality. For example, if the user says "where are some Italian restaurants around here" to her mobile device, the meaning as represented in EMMA would look the same as if the user typed the same request. The modality (speech or keyboard) would be represented as a property of the meaning, but the interpretation itself would be the same.

### ***6.3. Dialog processing***

Looking back at Figure 11, we see that the meaning resulting from NLU process can be sent to application components or to a dialog manager. For interactive applications, a dialog manager is very important, since it is the component of the overall system that decides how to act on the user's request. It does so either by reacting to the user with a system response or by taking action to accomplish a task, or both.

**The most commonly-used tool for dialog management in commercial systems is VoiceXML [11, 48]. As discussed earlier, VoiceXML is an XML language that defines a set of slots (called a "form" in VoiceXML) along with system prompts associated with each slot and speech-recognition grammars**



**that are used to process the user’s speech and extract the user’s meaning from the utterance.**

Figure 6 shows an example of a VoiceXML form. Originally, the grammar associated with a VoiceXML form was always a semantic grammar in SRGS [15] format; however, with the popularity of statistical natural language processing based on SLM’s and text classification, today the URL for a VoiceXML grammar often points to a statistical SLM recognizer.

There is also a considerable research literature on dialog management, particularly task-oriented dialog management [49]. Major approaches include systems based on planning [50, 51], information states [52], and agents [53]. [54] provides an excellent overview of commercial and academic approaches.

#### ***6.4. Text to Speech Technology***

Spoken output from a system can be provided by audio recordings, as it is in most IVR systems. Synthesized speech can also be used, and is required when it is impossible to pre-record every possible system response. The technology for synthesizing speech from text is called Text to Speech (TTS). There are two general approaches to TTS: *Formant-based* synthesis creates speech from rules describing how to generate the speech sounds; *concatenative synthesis* creates speech by piecing together snippets of prerecorded speech. Concatenative TTS is generally considered to sound better, but formant-based synthesis has a much smaller memory footprint because it doesn’t require a large database of prerecorded speech. It is therefore very practical to run formant-based synthesis locally on devices, which is important for minimizing latency.

#### ***6.5. Application Integration***

NLU is not very useful unless it’s able to accomplish tasks through interfaces to other software. Certainly, there are applications that just have a conversation, such as the very early program ELIZA [55] or more modern programs called “chatbots,” but most practical applications need an interface to other software that actually does something. For a personal-assistant program like Siri, this includes being able to access programs running on the device, like the user’s calendar and contacts, as well as being able to access external software such as Wolfram Alpha [56]. Siri can also access some device hardware, such as the GPS system. GPS information enables it to answer questions such as “Where am I” (although not similar questions such as “What is my exact location?”). However, Siri cannot access other hardware, such as the camera. Surprisingly little work has been done on the principles of integrating language with external systems; however, Ball, et. al., [57] and Norton, et. al. [58] describe a rule-based system for integrating natural

language processing results with other software and Dahl et. al., [59] discuss an XML interface to external services.

## 7. Summary

This chapter has reviewed the history of natural language processing and discussed the most common general approaches: multi-level analysis, semantic approaches, and approaches based on statistical text-classification -- using examples from such applications as IVR applications and mobile personal assistants. The chapter also places natural language processing in the context of larger systems for spoken and multimodal dialog interaction. In addition, it reviews related technologies, including speech recognition, dialog management, and text to speech. Today's NLU applications are extremely impressive, and the pace of their improvement is accelerating. Looking to the future, it is clear that these applications will become even more capable. These improvements will be driven by such factors as the dramatic increases in the power of devices, the development of new techniques for exploiting the vast amounts of data available on the World Wide Web, and improvements in related technologies such as speech recognition. All these factors are creating a synergy that will make the next generation of natural language applications ubiquitous and indispensable parts of our lives.

## 8. References

- [1] A. Turing, "Computing machinery and intelligence," *Mind*, vol. 59, pp. 433-60, 1950.
- [2] C. Shannon, "A Mathematical Theory of Communication," *Bell System Technical Journal*, vol. 27, pp. pp. 379-423, 623-656, 1948.
- [3] N. Chomsky, *Syntactic structures*. The Hague: Mouton, 1957.
- [4] K.-F. Lee, *Automatic Speech Recognition: The development of the SPHINX system*. Norwell, Massachusetts, USA: Kluwer Academic Publishers, 1989.
- [5] L. Hirschman, "Overview of the DARPA Speech and Natural Language Workshop," in *Speech and Natural Language: Proceedings of a Workshop Held at Philadelphia, Pennsylvania*, Philadelphia, PA, USA, 1989.
- [6] V. H. A. Yngve, "A model and a hypothesis for language structure," in *Proceedings of the American Philosophical Society*. vol. 104, ed, 1960, pp. 444-466.
- [7] M. Marcus, *A Theory of Syntactic Recognition for Natural Language*. Cambridge, MA: The M.I.T. Press, 1980.

- [8] W. A. Woods, "Transition Network Grammars for Natural Language Analysis," *Communications of the Association for Computing Machinery*, vol. 13, 1970.
- [9] J. E. Hopcroft and J. D. Ullman, "Introduction to Automata Theory, Languages, and Computation," ed: Addison-Wesley Publishing Company, 1987, pp. 13--45.
- [10] M. Minsky, "A framework for representing knowledge," in *Theoretical Issues in Natural Language Processing I*, ed, 1975.
- [11] S. McGlashan, *et al.*, "Voice Extensible Markup Language (VoiceXML 2.0)," ed: W3C, 2004.
- [12] W. Ward, "Understanding Spontaneous Speech," in *DARPA Speech and Language Workshop*, 1989.
- [13] S. Seneff, "TINA: A Natural Language System for Spoken Language Applications," *Computational Linguistics*, vol. 18, pp. 61-86, March 1992.
- [14] E. Jackson, *et al.*, "A Template Matcher for Robust NL Interpretation " in *Speech and Natural Language: Proceedings of a Workshop Held at Pacific Grove, California, February 19-22, 1991*, Pacific Grove, CA, 1991.
- [15] A. Hunt and S. McGlashan, "W3C Speech Recognition Grammar Specification (SRGS)," ed: W3C, 2004.
- [16] T. Bray, *et al.* (2004). *Extensible Markup Language (XML) 1.0 (Third Edition)*. Available: <http://www.w3.org/TR/2004/REC-xml-20040204/>
- [17] L. Van Tichelen and D. Burke. (2007). *Semantic Interpretation for Speech Recognition*. Available: <http://www.w3.org/TR/semantic-interpretation/>
- [18] (2001). *Standard ECMA-327 ECMAScript 3rd Edition Compact Profile*. Available: <http://www.ecma-international.org/publications/standards/Ecma-327.htm>
- [19] M. Johnston, *et al.* (2009). *EMMA: Extensible MultiModal Annotation markup language*. Available: <http://www.w3.org/TR/emma/>
- [20] D. A. Dahl, *et al.* (2000, November) Commercialization of natural language processing technology. *Communications of the ACM (electronic edition)*.
- [21] A. Moschitti, *et al.*, "Using Syntactic and Semantic Structural Kernels for Classifying Definition Questions in Jeopardy!," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, Edinburgh, UK, 2011.
- [22] H. Alshawi and J. van Eijck, "Logical Forms in the Core Language Engine," in *27th Annual Meeting of the Association for Computational Linguistics*, Vancouver, BC, Canada, 1989.
- [23] P. Clark and P. Harrison, "Boeing's NLP System and the Challenges of Semantic Representation," in *Semantics in Text Processing. STEP 2008 Conference Proceedings*, Venice, Italy, 2008.

- [24] *Telephony Voice User Interface Conference*. Available: <http://www.tmaa.com/>
- [25] B. Levin, *English Verb Classes and Alternations*. Chicago: The University of Chicago Press, 1993.
- [26] L. M. Norton, *et al.*, "Augmented role filling capabilities for semantic interpretation of natural language," in *Proceedings of the DARPA Speech and Language Workshop*, Pacific Grove, CA, USA, 1991.
- [27] M. Palmer, *et al.*, "The Kernel Text Understanding System," *Artificial Intelligence*, vol. 63, pp. 17-68, 1993.
- [28] C. Fellbaum, Ed., *WordNet: An Electronic Lexical Database*. Cambridge, Massachusetts: MIT Press, 1998, p.^pp. Pages.
- [29] X. Liu, *et al.*, "Recognizing Named Entities in Tweets," presented at the Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics, Portland, Oregon, 2011.
- [30] P. Turney, "Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews," in *Proceedings of the Association for Computational Linguistics*, 2002, pp. 417-424.
- [31] S. Crouch and R. Khosla, "Sentiment analysis of speech prosody for dialogue adaptation in a diet suggestion program," *SIGHIT Rec.*, vol. 2, pp. 8-8, 2012.
- [32] J. Chu-Carroll and B. Carpenter, "Dialog management in vector-based call routing," in *36th ACL/COLING*, Montreal, Quebec, Canada, 1998, pp. 256-267.
- [33] D. A. Dahl. (2006, September) Natural Language Processing: The next steps. *Speech Technology Magazine*.
- [34] R. Grishman, *et al.*, "Complex Syntax: Building a Computational Lexicon," in *COLING*, Kyoto, Japan, 1994, pp. 268-272.
- [35] M. Marcus, *et al.*, "Building a large annotated corpus of English: The Penn Treebank," *Computational Linguistics* vol. 19, pp. 313-330, 1993.
- [36] M. Palmer, *et al.*, "The Proposition Bank: An Annotated Corpus of Semantic Roles," *Computational Linguistics*, vol. 31, pp. 71-105, March 2005.
- [37] W. Woods, *et al.*, "The Lunar Sciences Natural Language Information System: Final Report," Bolt, Beranek and Newman 1972.
- [38] L. D. Erman, *et al.*, "The HEARSAY-II Speech Understanding System: Integrating Knowledge to Resolve Uncertainty," *Computing Surveys*, vol. 12, pp. 213--253, 1980.
- [39] J. A. Barnett, *et al.*, "The SDC Speech Understanding System," in *Trends in Speech Recognition*, W. A. Lea, Ed., ed Englewood Cliffs, NJ: Prentice-Hall, 1980, pp. 272-293.
- [40] J. J. Wolf and W. Woods, "The HWIM Speech Understanding System," in *Trends in Speech Recognition*, W. Lea, A. , Ed., ed Englewood Cliffs, New Jersey: Prentice-Hall, 1980, pp. 316-339.

- [41] L. R. Rabiner, "A tutorial on hidden Markov models and selective applications in speech recognition," *Proceedings of the IEEE*, vol. 77, February 1989.
- [42] D. A. Dahl, *et al.*, "Expanding the scope of the ATIS task: the ATIS-3 corpus," in *ARPA Human Language Technology Workshop*, Princeton, NJ, USA, 1994.
- [43] L. M. Norton, *et al.*, "Recent improvements and benchmark results for the Paramax ATIS system," in *Proceedings of the DARPA Speech and Language Workshop*, Harriman, New York, 1992.
- [44] S. Austin, *et al.*, "Proceedings of the Speech and Natural Language Workshop," Pacific Grove, CA, USA, 1991.
- [45] R. Bolt, "Put-That-There: Voice and gesture at the graphics interface.," *Computer Graphics*, vol. 14, pp. 262-270., 1980.
- [46] M. M. Taylor, *et al.*, Eds., *The Structure of Multimodal Dialogue* (Human Factors in Information Technology. North-Holland, 1989, p.^pp. Pages.
- [47] A. I. Rudnicky and A. G. Hauptmann, "Chapter 10: Multimodal interaction in speech systems," in *Multimedia interface design*, M. M. Blattner and R. B. Dannenberg, Eds., ed New York, NY: ACM Press, 1992, pp. 147-171.
- [48] M. Oshry, *et al.* (2007). *Voice Extensible Markup Language (VoiceXML) 2.1*. Available: <http://www.w3.org/TR/voicexml21/>
- [49] J. Allen, *et al.*, "An architecture for a generic dialogue shell," *Nat. Lang. Eng.*, vol. 6, pp. 213-228, 2000.
- [50] D. Bohus and A. I. Rudnicky, "RavenClaw: Dialog Management Using Hierarchical Task Decomposition and an Expectation Agenda," in *Eurospeech*, Geneva, Switzerland, 2003.
- [51] C. L. Sidner, "Building Spoken-Language Collaborative Interface Agents," in *Practical Spoken Dialog Systems*, D. A. Dahl, Ed., ed, 2004.
- [52] S. Larrson and D. Traum, "Information state and dialog management in the TRINDI dialog move engine toolkit," *Natural Language Engineering*, vol. 6, pp. 323-340, September 2000.
- [53] A. Nguyen and W. Wobcke, "An agent-based approach to dialogue management in personal assistants," presented at the Proceedings of the 10th international conference on Intelligent user interfaces, San Diego, California, USA, 2005.
- [54] K. Jokinen and M. McTear, *Spoken Dialog Systems*: Morgan & Claypool, 2010.
- [55] J. Weizenbaum, "ELIZA--A Computer Program For the Study of Natural Language Communication Between Man and Machine," *Communications of the ACM*, vol. 9, January, 1966 1966.
- [56] T. Claburn. (2009, March 10) Stephen Wolfram's Answer To Google. *Information Week*. Available: <http://www.informationweek.com/news/internet/search/215801388?pgno=1>

- [57] C. N. Ball, *et al.*, "Answers and Questions: Processing Messages and Queries," in *Speech and Natural Language: Proceedings of a Workshop Held at Philadelphia, Pennsylvania*, Philadelphia, PA, USA, 1989.
- [58] L. M. Norton, *et al.*, "Methodology for application development for spoken language systems," in *International Conference on Spoken Language Processing*, Philadelphia, PA, USA, 1996, pp. 662-664.
- [59] D. A. Dahl, *et al.*, "A Conversational Personal Assistant for Senior Users," in *Conversational Agents and Natural Language Interaction: Techniques and Effective Practices*, D. Perez-Marin and I. Pascual-Nieto, Eds., ed Hershey, PA, USA: IGI Global, 2011.