

An Integrated Speech-to-Text Parser for a Humanoid Robot

Gregory Koch¹ and Juan Fasola²
Austin College¹ University of Southern California²



Funded by the National Science Foundation (NSF) via
Research Experiences for Undergraduates (REU) Grant #0755534

Abstract:

Current work with test beds like the Bandit II humanoid robot platform has indicated early success with the efficacy of robots in physical rehabilitation therapy and autism spectrum disorder in early childhood development. However, few researchers have taken advantage of modern natural language processing systems to develop versatile dialog systems for robots that are designed for everyday use. The principal aim of this work is to integrate a general speech-to-text system using natural language processing to facilitate better human-robot interaction in the home environment. Further work will extend the robot's capacities to promote learning within the context of active dialog with other entities in real world situations.

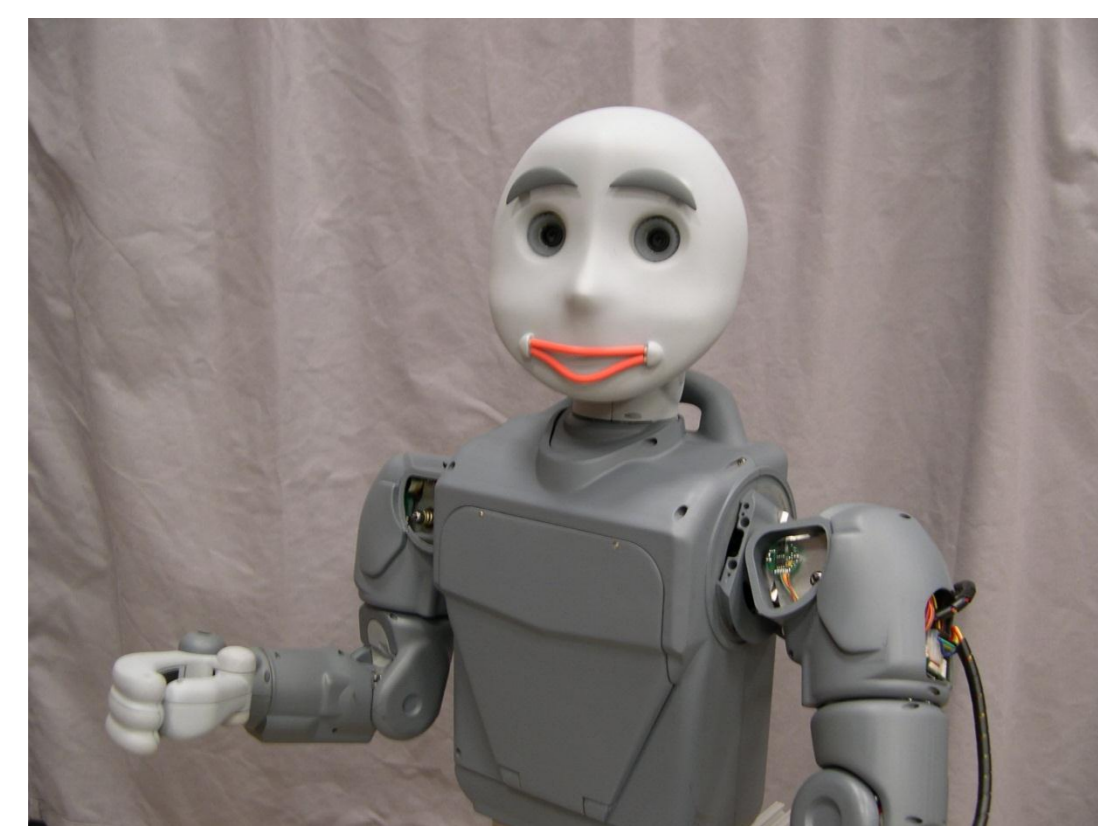


Figure 1: The Bandit II humanoid robot platform (Bluesky Robotics)

Stanford Natural Language Parser:

To begin to develop an advanced system using semantics to assign meaning to arbitrary phrases received as input data, we chose to use the **Stanford Natural Language Processing (SNLP) Parser**, a prevalent tool for extracting part-of-speech information.

By default, the SNLP parser runs in **Java**, so the major portion of integrating it into the rest of the code base involving writing a Java application capable of receiving microphone data and turning that into a usable data structure with the parser output contained in it. This data structure combines the three types of default output data that may be obtained from the parser: **part-of-speech tags** for individual words, **typed dependencies** representing the relationships between words (directed graph), and the **parse tree**, showing the connections between different parts of the entire phrase.

```
Your query
Bandit, go to the green table.

Tagging
Bandit/UH ,/, go/VB to/TO the/DT green/JJ table/NM ./

Parse
(ROOT
(S
(INTJ (UH Bandit))
(, ,)
(VP (VB go)
  (PP (TO to)
    (NP (DT the) (JJ green) (NN table))))
(. .)))

Typed dependencies
dep(go-3, Bandit-1)
prep(go-3, to-4)
det(table-7, the-5)
amod(table-7, green-6)
pobj(to-4, table-7)
```

Figure 3: Sample output NLP data obtained from the Stanford NLP Parser online demo, available at: <http://nlp.stanford.edu/software/lex-parser.shtml>

Architecture:

We offer a simple architecture emphasizing the modularity of each functionality of the program. Where possible, logic handling the treatment of data is separated from the modules responsible for extracting the microphone or NLP data itself.

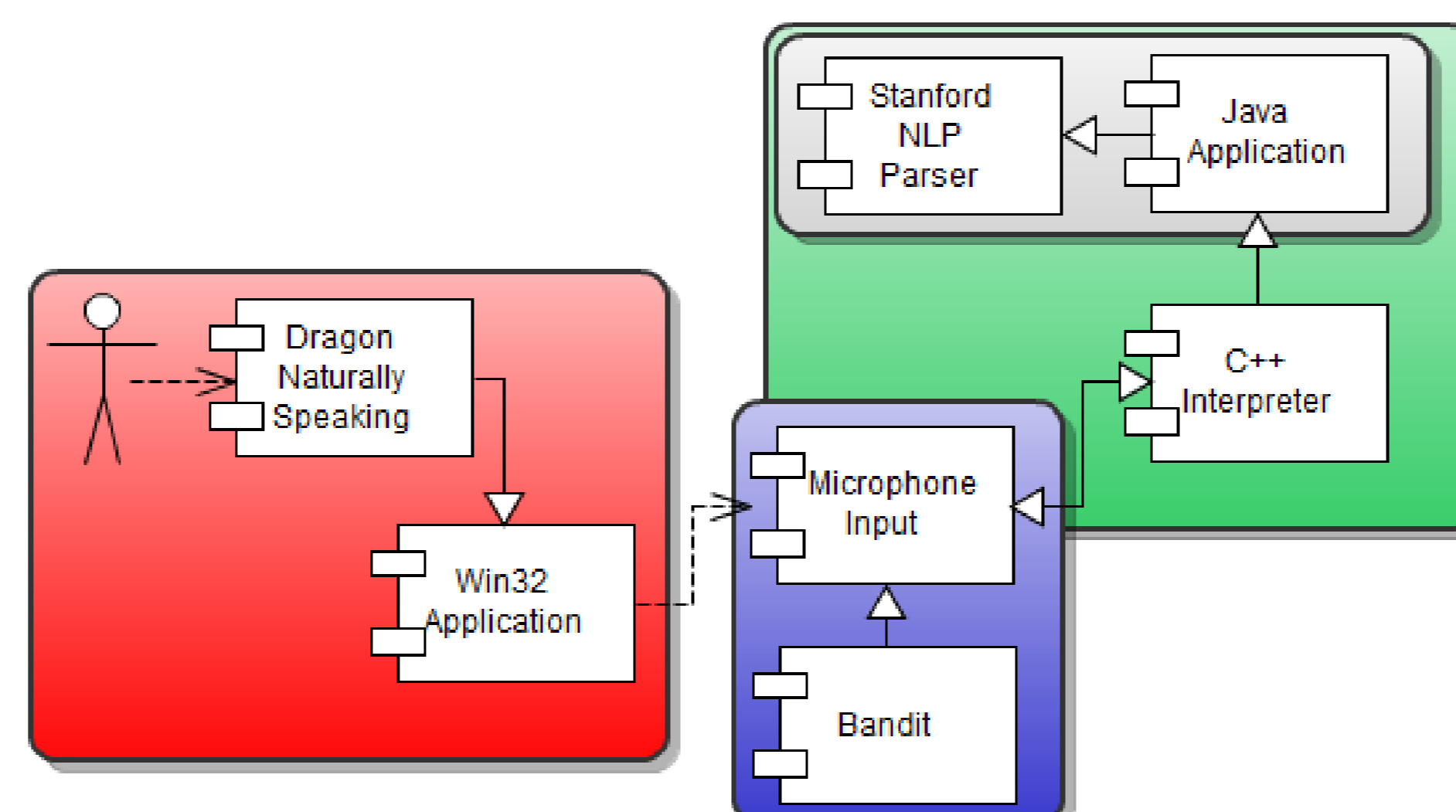


Figure 4: UML component diagram, providing an abstract view of system architecture with principal components of each module

Learning Game Demo:

The **Learning Game** demo encompasses all work that has been completed thus far for the project. It uses speech-to-text and NLP data to make real-time decisions and build up a concrete knowledge base of **verbs**, **nouns**, and **adjectives**, inferring which words correspond to which parts of speech.

The demo consists of a continuous series of prompts which attempt to augment the robot's knowledge base by asking to be taught commands. All commands consist of achieving some arm position with either the left or the right arm. Once the command is processed, the robot will ask about any words that are not in its knowledge base, then perform the corresponding action once it has made the semantic associations.

For nouns and adjectives, it will attempt to relate new words to ones it already knows. The base adjectives in the knowledge base are "left" and "right". The single base noun is "arm". In the case of verbs, however, the user has a choice: either relate the unknown verb to one that is already known ("raise" and "lower" are the base verbs that Bandit has been taught) or physically demonstrate an action by taking advantage of Bandit's vision processing capabilities.



Figure 5: Bandit learning a verb during a sample run of the Learning Game

As Bandit's knowledge base grows, the variety of commands it already knows becomes greater, so it becomes possible to relate new words to other new words via synonyms and use verbs outside of their original context.

Future Work:

With a working system in place that is able to independently capture speech data from a microphone and also retrieve natural language processing information for that data, the future of this project is flexible.

However, the direct extension of current work may be oriented towards learning commands to be used in a dialog setting in the context of human-robot interaction. By employing another higher-level parsing algorithm coupled with machine learning techniques, progress may be made towards a flexible system for adaptation in a home or clinical setting.

With minor modification, the code could be extended to any similar robot platform to the Bandit II humanoid robot; the many existing modules for Bandit serve as precedent, providing a practical test bed for testing speech functionality. Therefore, future work may also be implemented on the PR-2, a larger multifunctional robotics platform built by Willow Garage.

Further details available at: <http://robotics.usc.edu/~gkoch>, the website supporting this project.

Overview:

The project consisted of four primary phases. In the first phase, Dragon's speech controls were improved to allow for more dynamic behavior, such as stopping and restarting speech. This was necessary to build a flexible speech system, but will not be discussed further. A speech-to-text engine was then integrated to capture microphone data from an external microphone. Using this data with a natural language parser, we were able to obtain contextual information about each input phrase. Finally, we sent the data to the Robot Operation System (ROS) running on Bandit, where a higher-level interpreter used the data to learn new commands within the context of human-robot interaction.

Dragon Speech-to-Text Engine:

To implement speech-to-text functionality, we used **Dragon NaturallySpeaking**, a proprietary speech recognition available for Microsoft Windows. By installing Dragon on a Windows VM running concurrently with the Linux code for Bandit, we were able to forward the input data from the microphone to a TCP client that connects to ROS on Linux.

With this system in place, we are able to access current speech data from anywhere in the robot with relative ease.

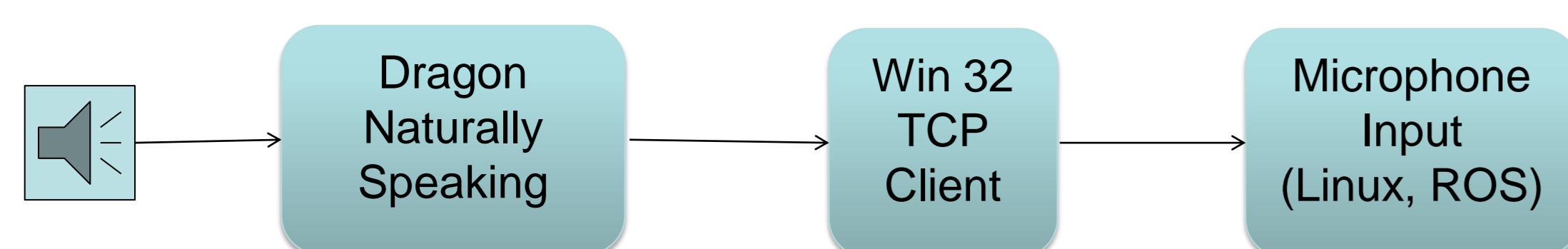


Figure 2: Microphone data is forwarded from the user, processed by Dragon, and then sent as a string to a ROS module on Linux.