

# REU Project Interim Report

Gregory Koch

July 1, 2011

## Summary

With the exception of problems connecting a microphone to the *Windows XP* VM, steady progress has been made to incorporate speech capacities into Bandit. Added functionality now allows some degree of control over when Bandit talks and estimating the position at which the robot stops talking in the text. A TCP server/client program was written to allow for the transfer of text from *Dragon NaturallySpeaking* running on a *Windows* VM to the Bandit code running natively on Linux (*Ubuntu*). This server/client model was then incorporated into a basic Win32 application designed to be the input controller for all text generated by Dragon. The Linux server was incorporated into Bandit's code to run concurrently. A speech interpreter object was also designed to interface between the TCP server and the rest of Bandit's code, particularly the speech module. The speech capacities of Bandit were also made more modular at this point, with emphasis on encapsulation and reduction in overhead with function calls. A new demo was then made to demonstrate that Bandit's new speech system had been successfully integrated with the existing code by introducing the Repetition Game, the HRI equivalent of an echo server program.

**Speech object enhancements.** Several elements were added to Bandit's `Speech` thread object during June: the functions `stopText`, `resumeText`, and `restartText` provide functionality. The `stopText` function allows any class with access to the `Speech` object to halt any speech being currently spoken by Bandit. Of the implemented functions, the `stopText` function is the most lightweight, because it only requires one boolean to be changed in the `LipSync` object to immediately force all spoken text to be paused. In the original implementation (this function was altered later, see below), the `resumeText` function would resume text from its approximate stopping point or if the robot is already speaking, it would do nothing. The `restartText` function restarts any sentence being currently spoken by determining the current phrase number in the active `sayText` command using the set of lip timings in the `LipSync` class.

The functions displayed near-perfect results in initial testing with a Gettysburg Address demo, however they proved to be insufficient to handle two edge cases: multiple `sayText` commands issued successively and stopping an active `sayText` command, resuming it, then restarting it. To address the first issue, a variable was added to track the current text being spoken. The second case is no longer an issue because `resumeText` has new functionality (see third item listed below).

Later, when being incorporated into the full Bandit demo, where all code was run exactly as it would be run on an actual version of Bandit, several problems were encountered with the implementation of these `Speech` functions:

- **Improper FSM Transitions.** First, a new boolean had to be added to the `Speech` class to explicitly indicate when the text had been paused, rather than just stopped. This was necessary because many FSM transitions in Bandit's code check whether Bandit is still speaking. If a command is issued to stop Bandit's current speech, however, the FSM transitions will immediately take

place. This will stop Bandit from continuing with his current phrase, but in many cases, Bandit immediately begins speaking the next planned phrase from the next state. Therefore, we must check both whether Bandit is still actively speaking *and* if Bandit has been paused by a command.

- **VTML Pause Tags.** The `Speech` object had previously not dealt with `sayText` commands containing VTML pause tags, which destroyed the parsing capacities that had been added. Functions that returned approximate speech position in the current text or restarted the `Speech` object from the beginning of the current sentence ceased to function. In response, functions were enhanced to automatically skip past any such tags by looking for an opening `<` character to delimit the tag open, with the corresponding closing `>` character to delimit the tag close. However, while this immediately solved most problems with the pause tags, the function estimating speech position in the current text was not repaired. The combination of additional time added to the `.wav` file and extra characters added to the `sayText` command ruined the previous method, which was to calculate the percentage of the `.wav` file already played by using a timer, then use that percentage with total string length to determine the approximate character being read.
- **resumeText changes.** Because having Bandit restart from an arbitrary point in a sentence is not a realistic behavior in most circumstances, the resume command was altered to make Bandit resume speaking from the beginning of the current phrase he is currently speaking. This was a convenient change to make since the position estimation function was broken by the VTML Pause Tags (see above) and the previous implementation of `resumeText` is therefore no longer possible. If Bandit is already speaking, nothing will happen. However, the restart command already does nearly the exact same thing, except that it always executes, whether Bandit is speaking or not. It should be questioned whether the resume command is necessary anymore, unless it is desired to resume Bandit from his exact stopping point.

In general, it is clear that the original implementation of Bandit did not account for speech controls; therefore, some changes were made as detailed above, and further changes will have to be made to ensure the new functionality is properly integrated with existing code.

**Win32 Application and TCP Server.** After the *Windows XP* VM was installed and mounted using Virtualbox, it was not difficult to establish communication between the VM and Linux by setting up a basic Win32 application on the *Windows* VM to serve as a client. The application consists of a small window with a text area. The application is preconfigured to send messages to a specific server IP over a predetermined port. Any text entered into the text area is immediately forwarded to the server. The idea is for *Dragon NaturallySpeaking* to eventually input speech-to-text output directly into this application, thus provided an indirect but immediate link between microphone and Bandit. This portion of the project is complete.

**Microphone issues.** Unfortunately, although extensive efforts were made to get a *Bluetooth* USB Microphone (*Plantronics* Calisto Headset) to work with the host machine, no progress was made. No effort to update drivers, use USB filters, or adjust volume settings succeeded in getting audio input with the headset on the *Windows XP* VM. However, no audio input was possible on the native Linux *Ubuntu* operation system install either, so it is possible that a permissions issue. No default sound settings seemed to be inhibiting microphone functionality. Resolving these microphone issues should be considered one of the highest priorities at this time, because the microphone is integral to enabling actual communication between human and robot.

**TCPServer/Interpreter.** The `TCPServer` runs on its own thread, repeatedly checking for input. It was given a queue to place all new text into, but otherwise functions purely as a TCP server. The `SpeechInterpreter` class was created to bridge the gap between the `TCPServer` object and the other classes incorporated in *Bandit*. The `SpeechInterpreter` class also runs on its own thread, continually checking the `TCPServer` object for new additions to its queue. If a message received by the `TCPServer` is a recognized command, this command is immediately executed. Otherwise, the `SpeechInterpreter` object moves the string into a queue of its own, which it maintains. Further, the `SpeechInterpreter` class contains a boolean to track if a break has been requested by the player. This effectively replaces the Wiimote's break functionality, because the `ExerciseInstructor` method that continually checks for a break now simply checks this boolean in the `SpeechInterpreter` class instead of reading in Wiimote button presses.

**Encapsulation.** Since the `Speech` class spawns the `SpeechInterpreter` thread, which in turn spawns the `TCPServer` thread, any object with access to the `Speech` object can access the `TCPServer` object. However, in practice, this access is limited, because all of `TCPServer`'s members are private. Therefore, an arbitrary object could not access client session data, for example. We restrict our attention to the relationship between the `Speech` and `SpeechInterpreter` classes. The `SpeechInterpreter` is constructed with a pointer to the `Speech` object solely to allow it to execute commands as they are received, such as stopping the current speech module from executing.

**Repetition Game demo.** An initial speech demo using the Gettysburg address and several other speeches was initially created with a basic FSM to demonstrate and test `Speech` object functions. The next logical step was to try to integrate these changes, along with the TCP Server, to simulate real scenarios on a fully functioning version of *Bandit*. To do this, several days were spent reviewing various code for different components of *Bandit* to understand its functionality.

The Repetition Game demo is five minutes in total length. When the demo runs, first the introduction module is executed. This has also been modified to now ask the player for their name before beginning the session. As a consequence, this has eliminated any need to supply the robot in advance with the player's name and

can instead be determined as a consequence of human-robot interaction. The name is stored in a variable and can be saved in the database; Bandit will not ask the player's name if it is not the first session. After the introduction is finished, Bandit will immediately begin the Repetition Game.

When the instructions are complete, the human player can then speak anything into the microphone, and is immediately mimicked by Bandit, who copies the input back verbally to the player. Because the microphone is not functioning, this currently functions as copy-pasting text into the Win32 application on the client side. During this time, any of the `Speech` commands can be executed (`sayText`, `stopText`, `resumeText`, `restartText`), although a few issues remain with the functionality of these commands during full demos, due to the complexity of Bandit's FSM.

**Further outlook and conclusion.** Current progress has been adequate to identify many significant implementation issues with integrating active speech functionality. It has proven fairly straightforward to control Bandit's speech module, send text from an external source, and integrate received text via a TCP server and parsing object. The crucial concern for July is certainly getting a microphone online so that *Dragon NaturallySpeaking* can actually be used. Secondary concerns include improving the integration of existing features to ensure that edge cases are smoothly handled and developing a more advanced parsing system in order to provide a more intelligent interaction between human and robot.

It will be necessary to consider the structure of any parsing object that interfaces between the message server and the rest of Bandit's code. The exact specifications of Bandit's eventual speech functions should be immediately decided to influence further implementation choices. Demos like the Repetition Game should serve as a test for Bandit's speech-to-text integration, but should be phased out in favor of real interaction in the context of existing games that serve the basis for Bandit's physical rehabilitation and exercise functions in ongoing experiments. Provided a good solution for the microphone issue is found, this will allow examination of how well two-way vocal communication can work and help to identify further errors that have not been uncovered.