# Voice Commands to Control Recording Sessions

## John "Marty" Goddard

*Thesis in partial fulfillment of the degree*

*Master of Science in Computer Engineering*

## Syracuse University

### Fall 2011

9/8/2011

1

# Introduction: Motivation

## Goal

- Provide hands-free recording operation
- Allow musicians to record themselves

## Historical

- "Put That There" gestures and voice 1981
- Voice Navigator for Macintosh Musicians 1989
- Guitar Guru music teaching software, slogan "Keep your hands on your axe"

# Recording Workflow for Overdubs:
## (adding instruments and voices to existing recording)

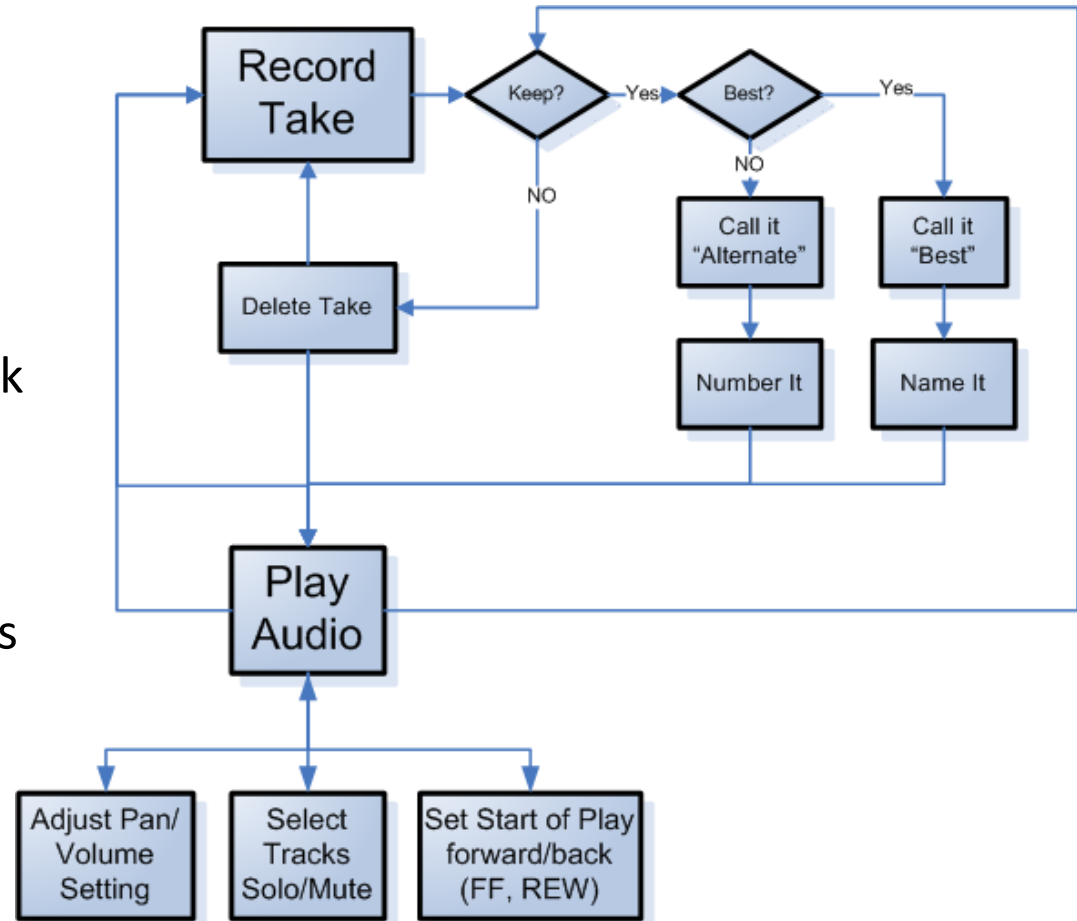Record Take: Begin Recording, perform music, stop recording

Delete take if not salvageable

Play Audio: Audition, or playback recorded tracks and listen for performance problems.

Solo: Selects track to hear, others are temporarily turned off

Mute: Temporarily turn off track

Pan: shift left or right in stereo mix

# Experimental Apparatus

- <u>Audacity</u> (audio recorder/editor)
  - Open-source: No cost, can modify
  - Scripting command interface
    - Allows control from other applications/programs

- <u>SayPlay</u> (created by me for this research)
  - Accepts voice commands, sends to Audacity
  - Utilizes Windows Speech Recognition
    - Built into Windows (.NET 3.5 and later)
    - Supports custom grammars, Documented API

9/8/2011

# Audacity audio recorder/editor

# SayPlay: handles voice command events, sends formatted commands to Audacity

- Accepts Voice input

- Interacts with Speech
  Recognition Engine
  – Creates Grammar structures
  – Handles recognition events

- Sends commands to Audacity via Script Interface



- Logs command events and results to a text file

# Voice Commands
(categorized by function)

- <u>Single word or phrase</u>:
  - Play, Stop, Record, Save, "Skip Start", "Play One Second", etc.
  - There are ~100 such commands, hard to remember them.

- <u>Naming things:</u> (tracks; in future: takes, sections)
  Ex: "Name this track 'kazoo' "

- <u>Adjusting things</u>: Refer to tracks by *<u>assigned name</u>*, to control mute, solo, pan, gain, selection settings (etc.)
  Ex: "Play the kazoo track",  "Mute the 'violin' track"

- <u>System Administration</u>:  Comments, context switching commands, refresh the list of names.

# Grammar Structure for Naming Tracks

Creating the Name

| Name |
|------|

Trapezoid containing: **This**, **The Recorded**

Trapezoid containing: **Take**, **Track**

Trapezoid containing: **Wildcard (Dictation)**, **Pre-Defined Track Names**

- Trapezoid represents choices
- { } Indicates optional word(s)

# Grammar Structure for Named Tracks

## Using the Name



- Trapezoid represents choices
- { } Indicates optional word(s)

\* Pan command appended with "Medium Left", "Hard Right" Etc.

# Learning Performance

63 = "Vocal Scat"

63

55 = "Vocal Scat"

47 = "Violins"

21 = "Vocal Scat"

26 = "Tambourine", "Alto Sax"

Alternating success/failure = "Pass"
misrecognized as "Pause"

7 = "Banjo", "Guitar", "Saxophone"

**Legend:**
- Perfect
- A
- B
- C
- D
- E

# Types of Recognition Failures

- **<u>Wrong Name</u>:** track name is misrecognized
- **<u>Wrong phrasing</u>:** User error, incorrect wording
- **<u>Low Confidence</u>:** Correctly recognized but value of confidence (returned by WSR) is below threshold
- **<u>Timeout</u>:** Long pause in speaking truncates phrase
- **<u>False positive</u>:** Utterance misinterpreted as a command, with confidence above threshold
- **<u>True rejection</u>:** A non-command is misrecognized, but confidence is below threshold, so it is justly rejected
- **<u>"Breath After"</u>:** WSR misrecognizes a dysfluency following a recognized word. Possible when an *optional additional* word is allowed.  Ex: "Wrong <u>*and*</u>"

# Problem With Dictation Speech Recognition: Names from "out of the blue"

- Random names (even made up ones), require large vocabulary <u>dictation speech recognition</u>

- Even after a name is correctly assigned, it can still be misrecognized when using the name

- Techniques were developed to improve accuracy <u>assigning names</u> and <u>using names</u>

# 3 Techniques were developed to Improve Assigning Names

- **1. Elaboration:** Use "Like" or "As In" to add words to a phrase about, or containing the desired name.   Ex: <u>"Name this track 'bass' as in bass guitar"</u>.

- **2. Quotation:** Enable a quoted phrase to be the name.  Ex: <u>"Name this track quote 'scream like a banshee' unquote"</u>.

- **3. Spell it:** <u>"Name this track **spelled** W, O, W"</u>.

# 3 Techniques to Correctly Refer to a Named Item, Once the Name is Assigned

(in addition to methods for Assigning Names)

1. Load track names into the Grammar loaded into the Speech Recognition Engine (see next slide)

2. Add a name to the Windows Speech Dictionary

3. Prevent recognition of confused names

9/8/2011

14

# Loading Names into Grammar

## Moves Recognition Duty…



From Here

To Here

# Loading Names into Grammar

– Get list of track names from Audacity

– Unload TrackCommands Grammar from Speech Recognition engine

– Construct a new grammar with new track names

– Reload the new TrackCommands Grammar

9/8/2011

# Steps to Load Names into Grammar

| User | Windows Speech Recognition | SayPlay | Audacity |
|------|---------------------------|---------|----------|

User says "Computer, Please Update the Session"

Recognition Event: "Computer, Please Update the Session"

Is Confidence > 0.93?

YES

Send Command to Get all track names

Return all track names

Ask WSR for a chance to Interrupt the running recognizer

When finished with Current task, make callback

Unload TrackCommand Grammar

Rebuild TrackCommand Grammar including All Track Names

Reload TrackCommand Grammar into Speech Recognition Engine

When finished Loading Grammar, make callback

Print out the new TrackCommand Grammar to status log

# Experiments with Homophones

Pairs of words that sound the same …
… but are Spelled Differently

- Example: "Pairs" and "Pears"
- I selected from over 1000 homophone pairs only those having to do with sound or music
- Find which is usually recognized (the default)
- Then try elaboration to assign the non-default
- And measure the effectiveness of loading the non-default spelling into the loaded grammar

# elaboration success rates



Horizontal bar chart of elaboration success rates (0% to 70%):

| Word | Rate |
|------|------|
| hurts | ~67% |
| wood | ~67% |
| Presence | ~62% |
| cord | ~50% |
| Bass | ~30% |
| rap | ~24% |
| Beet | ~22% |
| Wrap | ~21% |
| weekly | ~18% |
| four | ~17% |
| auricle | ~17% |
| tick | ~14% |
| hymn | ~10% |
| bated | ~6% |

# Interpretation of Elaboration Results

- Elaboration works only sometimes
- There seems to be a skew toward assigning Proper Names to named entities
  - Wei instead of way
  - Ryan or Orion, instead of rhyme
  - Kazue, instead of kazoo
  - Other confused names: Sarandon, Shanti, Peres,…
- Conclude that the command phrase "Name this track Theremin", has greater influence over outcome than by adding "like the Russian Inventor"
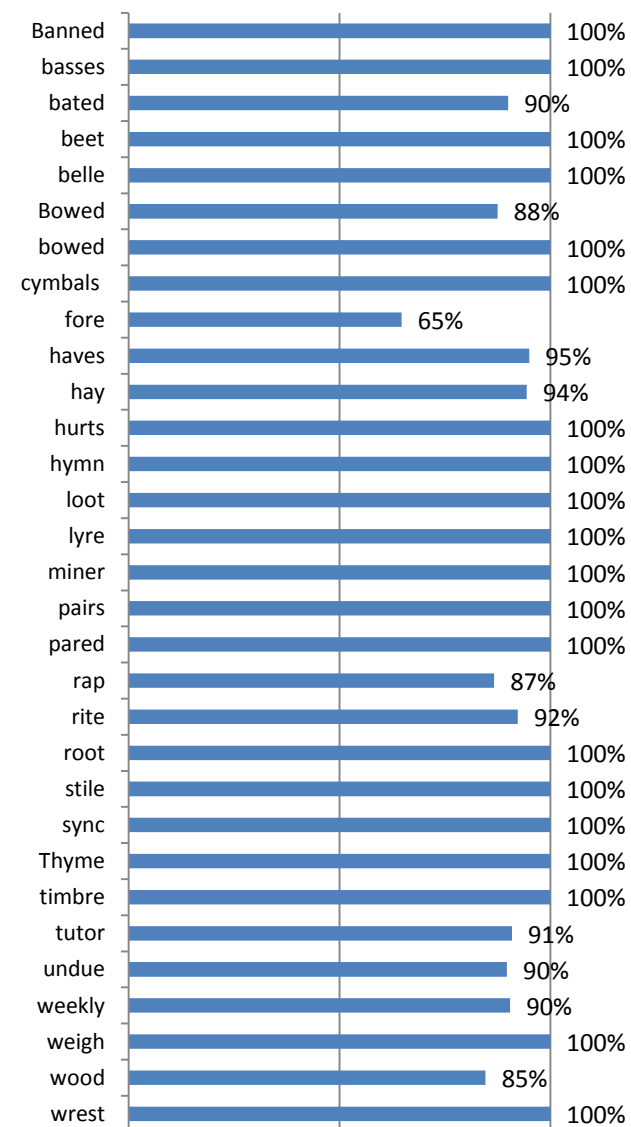
# After Loading into Grammar

- The non-default spelling of the Homophone pair was recognized with very high accuracy (even though it was rarely, if ever, recognized before loading the name into the grammar)

**Recognition Accuracy**
**(after loaded into Grammar)**

| Word | Accuracy |
|------|----------|
| Banned | 100% |
| basses | 100% |
| bated | 90% |
| beet | 100% |
| belle | 100% |
| Bowed | 88% |
| bowed | 100% |
| cymbals | 100% |
| fore | 65% |
| haves | 95% |
| hay | 94% |
| hurts | 100% |
| hymn | 100% |
| loot | 100% |
| lyre | 100% |
| miner | 100% |
| pairs | 100% |
| pared | 100% |
| rap | 87% |
| rite | 92% |
| root | 100% |
| stile | 100% |
| sync | 100% |
| Thyme | 100% |
| timbre | 100% |
| tutor | 91% |
| undue | 90% |
| weekly | 90% |
| weigh | 100% |
| wood | 85% |
| wrest | 100% |

9/8/2011

21

# Interpretation

- Using dictation speech recognition for names already assigned, does perform reasonably well (>90%) on many ordinary track names

- Loading names into the loaded grammar works even better than relying on dictation recognition

- So, why not *always* load the names into the grammar? Because the names must first be acknowledged as correctly spelled.

- Why not always load session names into grammar *when first opening a session* (when all names are known good)? Future work.

# Future Work on Voice Commands to Control Recording Sessions

- Name song sections for navigating in time
  - "Name this section 'chorus', or 'verse', or 'hook'"
  - "Play the Chorus section" or "jump to the chorus section"

- Command Macros using "This Means That"
  - "Computer, please rename the 'Play' command 'Audition'"

- Reference multiple objects in one command
  - "Mute the piano, bass and backing vocal tracks"

# Other Possible Future Work

- Naming tasks, processes, or searches for voice commands to complete user-defined multi-step commands.

- Transpose this work to
  - Personal note taker/assistant
  - Video recording/playback
  - Other domains of process control

# Questions?

- Thank you for attending