### Objectives

$\square$ be able to build simple voices within the system	$\square$ be able to explain how the system works	$\square$ be able to adapt it for your applications	$\square$ be able to influence the way it does it	$\square$ be able to make Festival speak what you want	After the course, you will	$\Box$ become familiar with Festival's architecture and know what it can and can't do.	$\square$ understand relative complexity of implementing solutions to problems	□ understand basic processes in speech synthesis
--	---	---	---	--	----------------------------	--	--	--

## Text to speech

Four major topics in speech synthesis

- $\square$  Architecture
- objects and processes required
- □ Text processing
- from text to tokens to utterances to words.
- ☐ Linguistic/prosodic processing
- lexicons, phrasing, intonation, duration
- □ Waveform synthesis
- diphone, unit selection, concatenation and modification

## Course outline

- $\square$  March:
- history, basic Festival use
- TTS, Utterance structure, processes
- Text analysis, Lexicons, LTS

- Prosody: phrasing, intonation, duration and power

- Prosody modelling techniques
- | April:
- Large projects:
- Waveform synthesis: diphones, unit selection
- Building a new voices in new languages
- Limited domain synthesis
- Project time
- Concept to speech, Voice Transformation
- Speech Synthesis evaluation and diagnosis

## Course evaluation

- $\square$  (approx) weekly small exercises:
- best 4 contribute to grade
- □ 1 "large" project:
- set beginning of April
- e.g. build a new voice
- Requires (appropriate) writeup.

 $\square$  No exam

## Important weblinks

- □ Course notes, slides, exercises
- http://www.cs.cmu.edu/~awb/11752.html
- $\square$  Building Voices in Festival
- http://www.festvox.org

audio examples from

From Dennis Klatt's 1987 JASA 82, 737-793 article. http://www.festvox.org/history/klatt.html (Sampled history, this isn't eveything).

### □ 1930s:

- o 1936: UK Telephone Company's speaking clock, used catenation optical storage, with phrases, words and part-word con-
- o (1) 1939: Bell Lab's Homer Dudley developed VODER mechanical (organ like) device that can "play" speech.
- o Dudley's VOCODER, devised by decomposing speech signal into source/filter model

### 1940s, 50s

- o Analogue based models
- o (4) OVE formant synthesis (Gunnar Fant) 1953
- o (11) The DAVO articulatory synthesizer developed by George Rosen at M.I.T., 1958

- o Digital models and text-to-speech, prosodic rules
- o (17) Elegant rule program for British English by John Holmes, Ignatius Mattingly, and John Shearne, 1964
- o (19) Rules to control a low-dimensionality articulatory model, by Cecil Coker, 1968
- o (20) First prosodic synthesis by rule, by Ignatius Mattingly, 1968

### ☐ **1970s**:

- o TTS, products, and diphones
- o (24) The first full text-to-speech system, done in Japan by Noriko Umeda et al., 1968
- o (21) Sentence-level phonology incorporated in rules by Dennis Klatt, 1976
- o (22) Concatenation of linear-prediction diphones, by Joe Olive, 1977.
- o (28) The inexpensive Votrax Type-n-Talk system, by Richard Gagnon, 1978.
- o (30) The M.I.T. MITalk system by Jonathan Allen, Sheri Hunnicut, and Dennis Klatt, 1979.

### □ 1980s:

- o concatenative speech, larger systems
- o (34) The AT&T Bell Laboratories text-to-speech system, 1985.
- o (35) Several of the DECtalk voices.
- o (36) DECtalk speaking at about 300 words/minute.

### $\square$ 1990s:

- o products, multi-lingual, unit selection
- o Japanese NUU-talk system, Sagisaka, Iwahashi, ATR, 1992
- o General unit-selection CHATR 1994 (diphones/unit selection example from Festival 1996)
- o Free software and everyone joins in, MBROLA 1995, Festival 1996.
- o Unit selection in AT&T's NextGen system 1999.

#### But ...

- □ Commercial systems still use recorded prompts
- $\square$  and word concatenation (cf. 1936)
- □ Because it sounds better.

### Festival: a generic speech synthesis system □ multi-lingual text to speech $\hfill\Box$ synthesis for language systems □ synthesis development environment

# Festival Speech Synthesis System

- □ Open Source speech synthesizer system:
- designed for development and runtime use
- □ Used in many commercial and academic systems:
- from large (AT&T) to small (various startups)
- distributed with RedHat 8.x
- hundreds of thousands users
- □ Multi-lingual
- no built-in language
- designed to allow addition of new languages
- □ Additional tools for rapid voice development:
- guidelines and documentation
- autoaligner for labelling recorded speech
- statistical learning tools
- scripts for building models

## Festival as software

$\square$ Free software	□ Full sources in distribution	□ Multiplatform (Unix/Windows NT)	$\square$ New languages without new C++ code	□ No fixed theories	o intonation analysis (F0, Tilt), signal processing, CART building, n-gram, SCFG, WFST, OLS	□ General tools:	o lexicons, LTS, duration, intonation, phrasing, POS tagging, to-kenizing, diphone/unit selection, signal processing	□ General replaceable modules:	$\Box$ C/C++ code with Scheme scripting language	□ General system for multi-lingual TTS	http://festvox.org/festival/	
11-752. LTI. Carnegie Mello					build-		ng, to-					

# The CMU FestVox project

<ul> <li>text analysers</li> <li>prosodic modelling</li> <li>lexicon building</li> </ul>	<ul> <li>diphone</li> <li>unit selection</li> <li>limited domain</li> </ul> Other support:	☐ Support for English and other languages ☐ Support for different waveform synthesis methods:	<ul> <li>Discussion and examples for building voices</li> <li>Example voice databases</li> <li>Step by step walktrhoughs of processes</li> </ul>	☐ Festvox: building synthetic voices:  — Tools, scripts documentation	http://festvox.org
11-752, LTI, Carnegie Mellon					

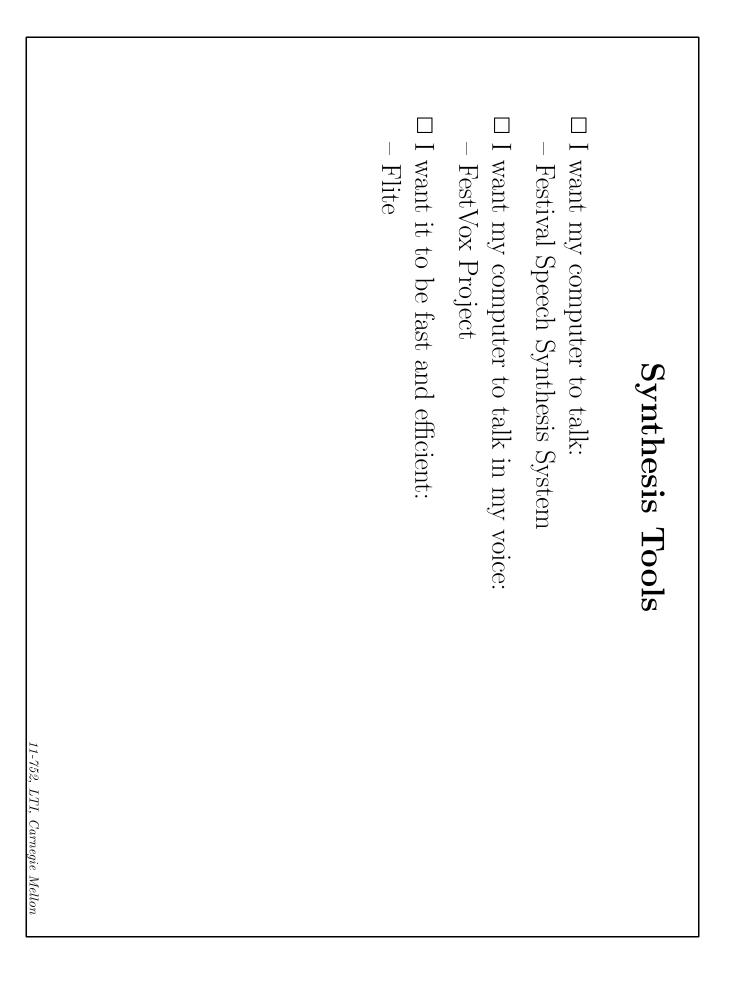
## The CMU Flite project

http://cmuflite.org FLITE a fast, small portable run-time synthesis engine

- $\square$  C based (no loadable files)
- □ Basic FestVox voices compiled into C/data
- □ Thread-safe
- Suitable for embedded systems and multi-client servers
- Ipaq, Linux, WinCE, etc
- □ Scalable:
- quality/size/speed trade offs
- frequency based lexicon pruning
- quality/speed db compression

### $\square$ Sizes:

- less than 4Meg footprint (code+data+runtime RAM)
- less that 0.025 secs time to speak (streaming synthesis)



### □ Basic Festival commands □ Scheme (Festival's scripting language) □ How to get Festival to talk □ Exercises Using Festival 11-752, LTI, Carnegie Mellon

## Getting it to talk

- $\square$  Say a file
- festival --tts file.txt
- □ From Emacs
- say region, say buffer
- □ Command line interpreter
   festival> (SayText "Hello")

# Scheme – Festival's scripting language

- $\square$  Why:
- too many options
- need flexibility
- easy to add functionality
- □ Why Scheme:
- very simplevery powerfulwell established
- authors are familiar with it

## Bluffer's guide to Scheme

Scheme is a dialect of Lisp

- $\square$  **expressions** are
- atoms or

lists

- a bcd "hello world" 12.3
- (a b c) (a (1 2) seven)
- $\square$  Interpter evaluates expressions
- atoms evaluate as variables

lists as function calls

festival\_version

3.14

(+23)

# Bluffer's guide to Scheme

```
□ Setting variables
- (set! a 3.14)
□ defining functions
- (define (timestwo n) (* 2 n))
(timestwo a)
6.28
```

### Scheme: lists

```
apples
                                                              festival> (length alist)
festival> (length (append alist blist))
                                                                                             bananas)
                                                                                                                                                           festival> (append alist blist)
                                                                                                                                                                                                                         festival> (set! blist (cons 'oranges alist))
                                                                                                                                                                                                                                                       (pears bananas)
                                                                                                                                                                                                                                                                                          festival> (cdr alist)
                                                                                                                                                                                                                                                                                                                                                         festival> (car alist)
                                                                                                                                                                                                                                                                                                                                                                                                                       festival> (set! alist '(apples pears bananas))
                                                                                                                                                                                         (oranges apples pears bananas)
                                                                                                                                                                                                                                                                                                                                                                                      (apples pears bananas)
                                                                                                                             (apples pears bananas oranges apples pears
```

## Scheme: speech

```
<UTT 983277>
                    festival> (SayText "This
                                                                                                                                                     festival> (utt.synth utt1)
                                                                                                                                                                           Synthesize an utterance
                                                                                   festival> (utt.play utt1)
                                                                                                                                                                                                                      festival> (set! utt1 (Utterance Text "Hello"))
                                                                 <UTT 982345>
                                                                                                       Play waveform
                                                                                                                                <UTT 982345>
                                                                                                                                                                                                  <UTT 982345>
                                                                                                                                                                                                                                           Make an utterance of type text
                                        Do all together
                        ը.
                   an example")
```



□ alt-m goto manual page  $\square$  alt-s to speak help  $\square$  alt-h on current symbol short help □ Online manual "http://festvox.org/docs/manual-1.4.3/"

 $\square$  Use TAB key for completion

# Lexicons and Lexical entries

Festival will make mistakes in pronunciation. You can explicitly give pronunciations for words. □ '("cepstra" n (((k eh p) 1) ((s t r aa) 0)))) □ Each language/dialect has its own lexicon And entry consists of You can add entries to the current lexicon You can lookup words with a Syllable is ( ( PHONE0 PHONE1 ...) STRESS ) ( WORD POS ( SYL0 SYL1 ...)) (lex.lookup WORD PartOfSpeech) (lex.add.entry NEWENTRY)

### Exercises

by noon March 10th to antoine@cs.cmu.edu

- 1. \* Make Festival say your name, and then everyones name in the class. Add explicit pronunciations to the lexicon.
- 2. Use Festival to say selected pieces of text. Find ten things Festival fails on.
- 3. How long does it take for Festival to say Alice's Adventures in Wonderland?

### lints

```
3. use
                                                                                                                               2. See Festival manual for instructions for installation
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        l. to test pronunciations, use
$SPPPDIR/src/festtut/examples/books/alice29.txt
                                                                                        "http://festvox/docs/manual-1.4.3/"
                                                                                                                                                                           on similarly pronounced words
                                                                                                                                                                                                                                                                                                                                                                                                              to set new entries in the lexicon use
                                                                                                                                                                                                                                                         To find out phone names use
                                                                                                                                                                                                                (lex.lookup "word" nil)
                                                                                                                                                                                                                                                                                                                                                                                                                                                   (SayText "My name is John Smith.")
                                                                                                                                                                                                                                                                                            '("edinburgh" n (((eh d) 1) ((ah n) 0) ((b ax ) 0) ((r ow) 0))))
                                                                                                                                                                                                                                                                                                                                                                    (lex.add.entry
```

an *estimate* of the time is fine