# COM3502/4502/6502
# SPEECH PROCESSING

## Lecture 14
## Waveform Processing

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 1*

1

---

# Block Processing

- After capture through a microphone and digitisation through the analogue-to-digital converter, an incoming speech signal becomes a sequence of quantised samples

- Digital signal processing is typically performed on a fixed-length sequence of samples called 'blocks' or 'frames'
  - e.g. in Pure Data the default 'block size' is 64 samples
    (*i.e. 1.45 msecs frame at the default 44.1 kHz sampling rate*)

- Because of the quasi-stationary nature of speech, the frame size is a compromise of …
  - having *sufficient* data in a frame to make the required measurements
  - having *small enough* amount of data that the stationarity assumption is fulfilled

- It is also necessary to ensure that there are a sufficient number of frames to capture the *non-stationary* properties
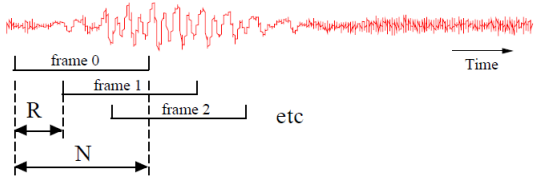
The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 2*

2

## Slide 3

# Block Processing

- To accommodate all these constraints, it is usual to use *overlapping* frames in speech processing
  - 'frame size' ($N$): number of samples per frame
  - 'frame shift' ($R$): number of samples between the start of successive frames

- Frame size is often expressed in time …
  - $NT$ seconds *(where T is the sample period)*

- Frame shift is often expressed as the 'frame rate' …
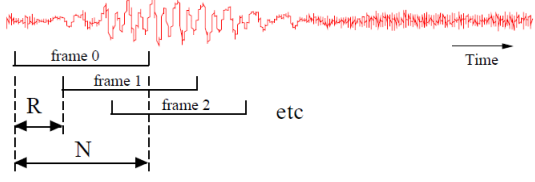  - $f_r = 1/RT$ frames per second (*fps*)



The University Of Sheffield.

3

## Slide 4

# Block Processing

- In speech, it is usual to have …
  - frame length ($NT$) $\approx$ 30 msecs
  - frame rate ($f_r$) $\approx$ 100 fps

- For example …
  - sample rate ($f_s$) = 10 kHz (*10,000 samples/sec*)
  - sample period ($T$) = $1/f_s$ = 100 $\mu$secs/sample
  - frame size ($N$) = $NT/T = 0.03/0.0001$ = 300 samples
  - frame shift ($R$) = $1/f_rT = 1/(100*0.0001)$ = 100 samples
  - frame overlap ($N-R$) = $300-100$ = 20 msecs (*66%*)



The University Of Sheffield.

4

# Block Processing in Pure Data

Signal Block

| A$_1$ | A$_2$ | A$_3$ | A$_4$ |
|---|---|---|---|
| 31.4 | 15.9 | 26.5 | 35.8 |

Wire

| B$_1$ | B$_2$ | B$_3$ | B$_4$ |
|---|---|---|---|
| 97.9 | 42.3 | 84.6 | 26.4 |

Object Box

+~

Inlet

| A$_1$+B$_1$ | A$_2$+B$_2$ | A$_3$+B$_3$ | A$_4$+B$_4$ |
|---|---|---|---|
| 129.3 | 58.2 | 111.1 | 62.2 |

**Taken from: Farnell, A. (2008). *Designing Sound*.**
**London: Applied Scientific Press Limited.**

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 5*

5

---

# Block Processing in Pure Data

- The default settings in Pure Data are …
  - sampling rate = 44.1 kHz (*defined by the sound card*)
  - block size = 64 samples (*1.45 ms*)
  - no overlap

- These can be overridden using the [block~] object
  (*but not in the same patch as the* [adc~] *or* [dac~] *objects*)

- [block~] takes the following parameters …
  - block size (*in samples, power of 2*)
  - overlap (*power of 2*)
  - up/down sampling ratio (*relative to parent window*)

- Only one [block~] object is allowed in a window

- A reasonable setting for speech is [block~ 1024 2 1]
  - 23 msec frame size
  - 86 frames per second
  - 50% frame overlap

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 6*

6

# Block Processing in Pure Data

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 7*

7

---

# Sample-by-Sample Processing in Pd

- Sometimes it is necessary to process data <u>one sample at a time</u>

- Pd provides an object for this: [fexpr~]

- [fexpr~] takes the following arguments …
  - $i#: integer input variable on inlet #
  - $f#: float input variable on inlet #

- Expressions in [fexpr~] are constructed using …
  - $x#[n]: the sample from inlet # indexed by n
  - $y[n]: the output value indexed by n
    (*$x# is shorthand for the current input*)
    (*$y is shorthand for the previous output: $y[-1]*)

- E.g. [fexpr~  $x1+$y] is a simple accumulator

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 8*
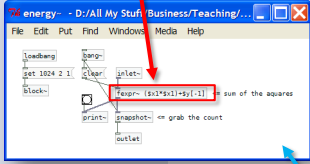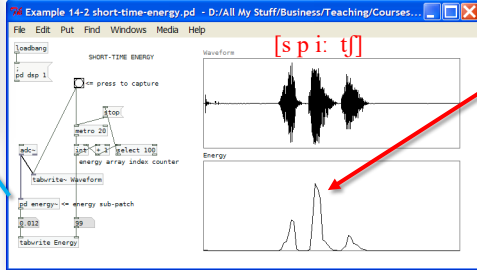
8

© 2019 The University of Sheffield

# Short-Time Energy

'Short-time energy' = sum of the squares of the samples in one frame

Sample by sample computation

$$E = \sum_{i=0}^{N-1} s_i^2$$

Energy is large in voiced speech

[s p i: tʃ]

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 9*

9



© 2019 The University of Sheffield

# Short-Time Energy

energy~ - D:/All My Stuff/Business/Teaching/...

File   Edit   Put   Find   Windows   Media   Help

loadbang        bang~

set 1024 2 1    clear    inlet~

block~                   fexpr~ ($x1*$x1)+$y[-1]   <= sum of the aquares

                print~   snapshot~   <= grab the count

                         outlet

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 10*

10

© 2019  The University of Sheffield

# Zero-Crossing Rate

'ZCR' = number of times the zero axis is crossed in <u>one</u> frame

ZCR is large in <u>unvoiced</u> speech

[s p i: tʃ]

**zerocross~ - D:/All My Stuff/Business/Teachin...**
File  Edit  Put  Find  Windows  Media  Help

loadbang
set 1024 2 1
block~

inlet~
*~ 1e+006   <= multiply by a very large number
clip~ -1 1   <= clip to create a square wave
fexpr~ -$x1*$x1[-1]  <= multiply adjacent samples
              (each '1' indicates a zero-crossing)
clip~ 0 1  <= remove all the '-1's
fexpr~ $x1+$y[-1]  <= count the remaining '1's
snapshot~  <= grab the count
outlet
bang~
clear
print~

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 11*

11

---

© 2019  The University of Sheffield

# Zero-Crossing Rate

**zerocross~  - D:/All My Stuff/Business/Teachin...**

File   Edit   Put   Find   Windows   Media   Help

loadbang
set 1024 2 1
block~

inlet~
*~ 1e+006   <= multiply by a very large number
clip~ -1 1   <= clip to create a square wave
fexpr~ -$x1*$x1[-1]  <= multiply adjacent samples
                (each '1' indicates a zero-crossing)
clip~ 0 1  <= remove all the '-1's
fexpr~ $x1+$y[-1]   <= count the remaining '1's
snapshot~  <= grab the count
outlet
bang~
clear
print~

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 12*

12

# Speech/Non-Speech Detection

- In speech processing it is often useful to be able to detect when someone is speaking

- Accurate speech 'end-point detection' is very difficult

- A simple 'speech/non-speech detector' can be constructed using short-time energy *and* zero-crossing rate
  - energy is high in <u>voiced</u> speech
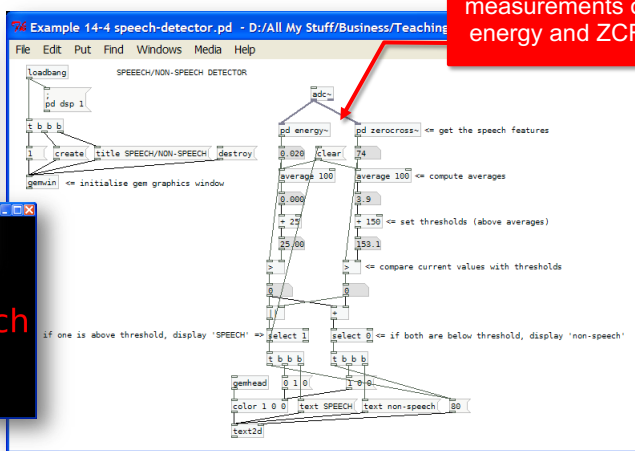  - ZCR is high in <u>unvoiced</u> speech

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 13*

13

---

# Speech/Non-Speech Detection

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 14*

14

# Autocorrelation Function

The 'autocorrelation function' computes the correlation of a signal with itself (*as a function of time*)

$$r_k = \sum_{i=0}^{N-k-1} s_i \cdot s_{i+k}$$



$r_0$ = energy

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 15*

15

---

# Autocorrelation Function

- The autocorrelation function (*ACF*) emphasises periodicity

- ACF is the basis for many spectrum analysis methods

- Short-time ACF (*STACF*) is the basis for many 'pitch detectors' (*fundamental frequency estimators*)

- ACF is fairly expensive to compute (*because there is an inner loop running for every data sample*)

- STACF is often combined with ZCR to construct a 'voiced/unvoiced detector'

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 16*

16

# Correlation

The 'correlation' between two discrete-time signals $s$ and $t$ over an $N$ point interval is …

$$q = \sum_{i=0}^{N-1} s_i . t_i$$

The
University
Of
Sheffield.

19

# Cosine Correlation

For two sinusoids (where *N and T are chosen so the summation is over an <u>integer</u> number of cycles for both signals*) …

$$s[nT] = A.\cos(\omega_s nT)$$
$$t[nT] = \cos(\omega_t nT)$$
$$q = s[0].t[0] + s[T].t[T] + s[2T].t[2T] + \dots$$
$$+ s[(N-1)T.t[(N-1)T]$$

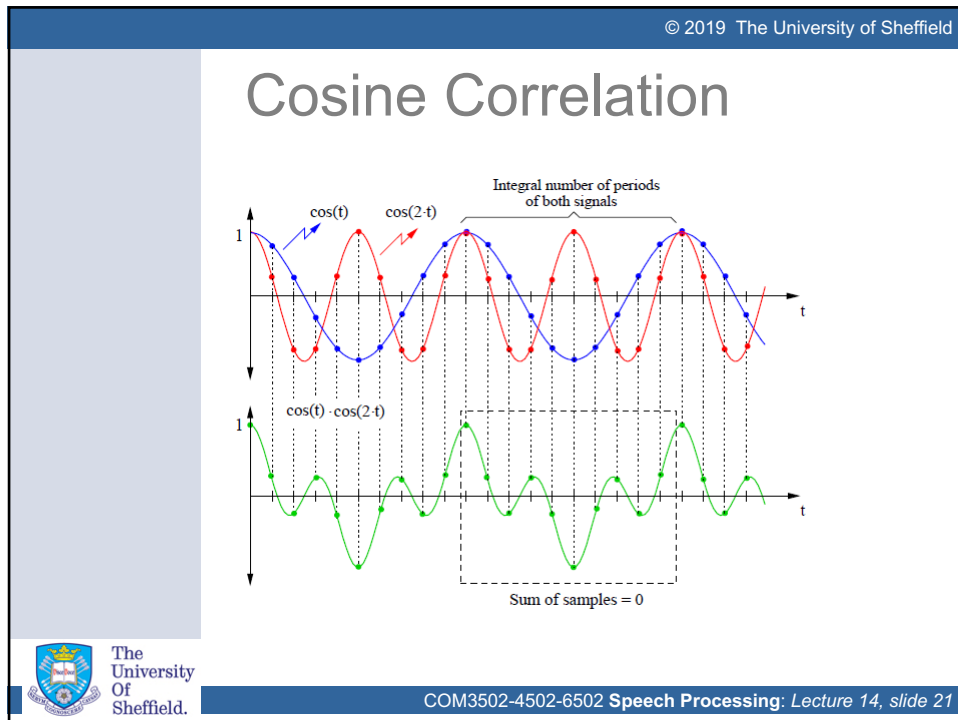It can be shown that …

$$q = \begin{cases} \alpha.A & if \quad \omega_s = \omega_t \\ 0 & otherwise \end{cases}$$

The
University
Of
Sheffield.

20

21



## Cosine Correlation

- The correlation between a test signal $t[NT]$ and a target signal $s[NT]$ is proportional to the amplitude $A$ of the target signal when …

$$cos(\omega_s NT) = cos(\omega_t NT)$$

- So, given that Fourier analysis shows that any signal can be decomposed into sinusoidal waves, cosine correlation can be used as a method to find (*extract*) the cosine components of an arbitrary signal

- This is only possible if the correlation is computed over an integer number of $p$ cycles (*in the test signal*)

$$\omega_t = \frac{2\pi p}{NT} \qquad p = 0, 1, \ldots N-1$$

$$f_t = \frac{p}{NT}$$

22

# Cosine Correlation
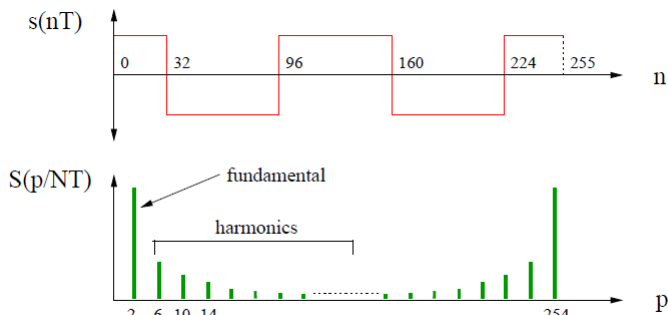
So, the spectrum computed by cosine correlation is …

$$S_p = \sum_{n=0}^{N-1} s_n . \cos\left(\frac{2\pi np}{N}\right) \quad p = 0 \ldots N-1$$

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 23*

23

---

# Cosine Correlation



s(nT)

0  32  96  160  224  255  n

S(p/NT)

fundamental

harmonics

2  6  10  14  254  p

$$S_p = \sum_{n=0}^{N-1} s_n . \cos\left(\frac{2\pi np}{N}\right) \quad p = 0 \ldots N-1$$

The University Of Sheffield.

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 24*

24

## This lecture has covered …

- Processing blocks/frames
- Processing samples
- Short-time energy
- Zero crossing rate
- Speech/non-speech detection
- Autocorrelation
- Cosine correlation

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 26*

© 2019  The University of Sheffield

# Any Questions ?

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 27*

The University Of Sheffield.

27



© 2019  The University of Sheffield

# Next time …

## The Fourier Transform

COM3502-4502-6502 **Speech Processing**: *Lecture 14, slide 28*

The University Of Sheffield.

28