# Data Mining and Machine Learning

# Applicati r ASR:
# Feature re f speech

Peter Jančovič

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Objectives

- Front-end analysis for ASR – feature representation of speech
  - To understand motivation and stages for 'typical' parameteris ed for ASR
  - Mel Freque (MFCCs)

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# What is "Front-End Analysis"

- First stage in any speech recognition system
- Goal is to convert the raw acoustic speech waveform into a form which is suitable (or even optimal) for                                    cognition
- In general p                                            ms, front-end analysis is feature extrac
- Where do we start?

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# The Human Auditory System

taken from J N Holmes, "Speech Synthesis and Recognition", Van Nostrand Reinhold (1988)

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# The Basilar Membrane

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Australian National University –
http://online.anu.edu.au/IT
A/ACAT/drw/PPofM/heari
ng/hearing3.html

Data Mining and Machine Learning

# Frequency response of the basilar membrane

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

School for advanced studies, Triste, Italy –

http::/poirot.sissa.it/multidisc/cochlea/utils/basilar.htm

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Lessons from Psycho-Acoustics

- Human speech perception begins with frequency analysis on the basilar membrane

- Frequency is not perceived on a linear scale – hence use of non-linear perc ~~Assignment Project Exam Help~~ **mel** scale, **bark** scale,…

- Individual point on the basilar ~~can be~~ modelled as band-pass filter – a **critical ba** ~~licit~~ bandwidth of such an 'auditory filter'

- Loudness perceived on logarithmic scale

- Phase of limited significance for speech recognition

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

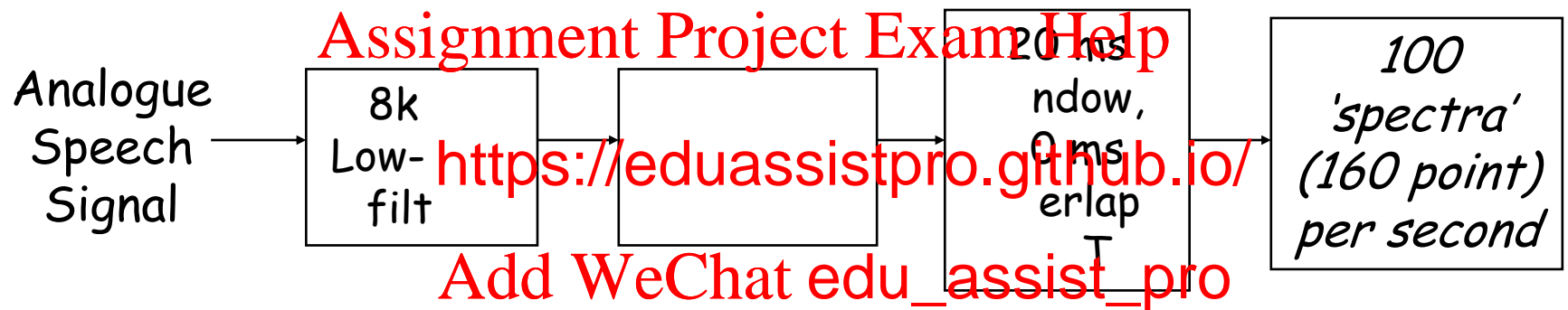# Front-end analysis for ASR

- Speech waveform typically low-pass filtered at 4kHz to 8kHz

- Sampled 8,000 to 16,000 samples per second

- Frequency ana

  - 20 ms anal

  - 10 ms overlap between

  - Hamming window

  - Discrete Fourier Transform

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Frequency analysis for ASR

Analogue Speech Signal →

| 8k Low-filt |

| 20 ms window, 0 ms overlap T |

| 100 'spectra' (160 point) per second |

## Example: 8kHz bandwidth system

UNIVERSITY OF BIRMINGHAM

# Log Power Spectrum

- Phase ignored by taking the **modulus** of the complex spectrum

- Logarithm applied <span style="color:red">Assignment Project Exam Help</span>

  - For consi <span style="color:red">https://eduassistpro.github.io/</span> oustic results

  - To compr

    <span style="color:red">Add WeChat edu_assist_pro</span>

| *160 point short-time spectrum* | → | *modulus & logarithm* | → | *160 point short-time log-power-spectrum* |

Data Mining and Machine Learning
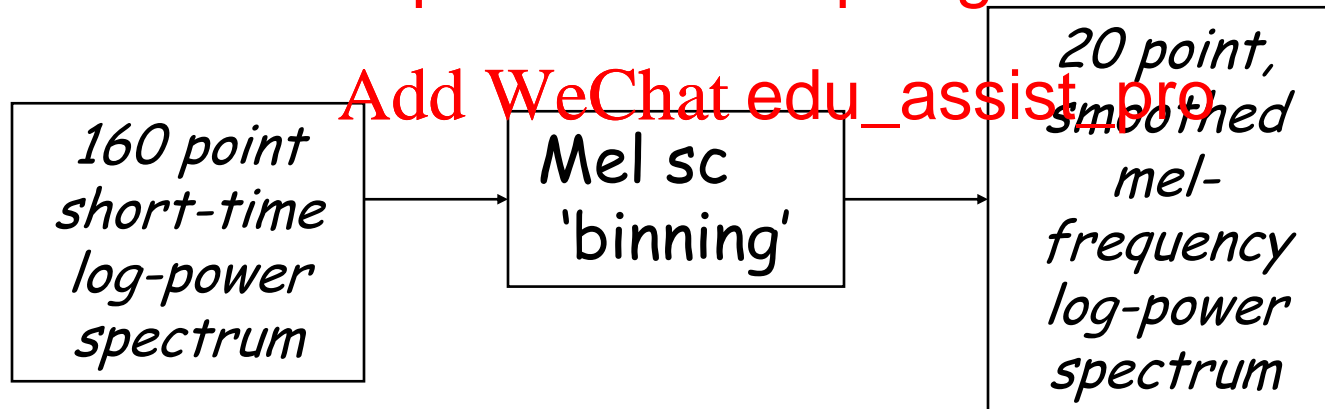
UNIVERSITY OF BIRMINGHAM

# Mel-scale & smoothing

- The **mel spectrum** can be computed by **averaging** the short-time Fourier spectrum over 'bins' whose width depends on frequency…

- …or by using band-pass filters with appropriate, frequency-dependent, ba

```
┌─────────────┐      ┌──────────┐      ┌──────────────┐
│ 160 point   │      │ Mel sc   │      │ 20 point,    │
│ short-time  │ ───> │ 'binning'│ ───> │ smoothed     │
│ log-power   │      │          │      │ mel-         │
│ spectrum    │      │          │      │ frequency    │
│             │      │          │      │ log-power    │
│             │      │          │      │ spectrum     │
└─────────────┘      └──────────┘      └──────────────┘
```

Data Mining and Machine Learning

**UNIVERSITY**OF
**BIRMINGHAM**

# Mel Scale Filterbank

From Steve Young, "The HTK Book", Cambridge University Engineering Department

Data Mining and Machine Learning

UNIVERSITY OF
BIRMINGHAM

# Cepstrum

- Cosine transform applied to remove correlation between components of mel-scale log power spectrum
  - Mel Cepstrum: MFCC = <u>M</u>el <u>F</u>requency <u>C</u>epstral <u>C</u>oefficients

  <span style="color:red">Assignment Project Exam Help</span>

  - Mathemat<span style="color:red">https://eduassistpro.github.io/</span>

<span style="color:red">Add WeChat edu_assist_pro</span>

```
┌──────────────┐        ┌──────────────┐        ┌────────────────────┐
│ 20 point mel │        │    Cosine    │        │     20 MFCCs       │
│  scale log   │ ─────► │  Transform   │ ─────► │ (use only first 12)│
│    power     │        │              │        │                    │
│   spectrum   │        └──────────────┘        └────────────────────┘
└──────────────┘
```

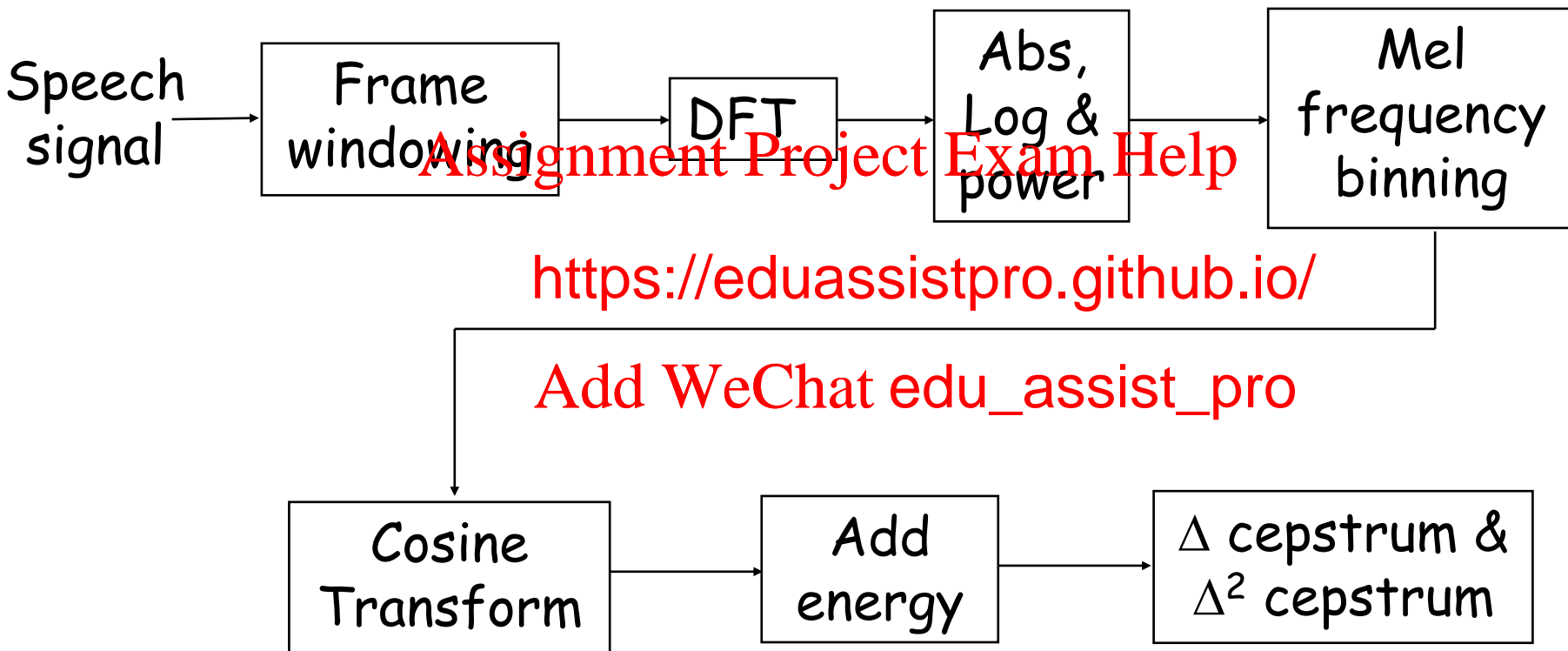Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Energy & Delta Coefficients

- Add energy as 13$^{th}$ parameter

- Compute estimate of time-derivative of each parameter – <span style="color:red">Assignment Project Exam Help</span> pstrum)

  <span style="color:red">https://eduassistpro.github.io/</span>

- Compute est tion of parameter – delta$^2$ cepstru <span style="color:red">Add WeChat edu_assist_pro</span> pstrum)

- Cepstum + $\Delta$ Cepstrum + $\Delta^2$ Cepstrum = 'standard' 39 dimensional representation (e.g. in HTK)

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM

# Front-end analysis – summary

Speech signal → Frame windowing → DFT → Abs, Log & power → Mel frequency binning

Cosine Transform → Add energy → $\Delta$ cepstrum & $\Delta^2$ cepstrum

UNIVERSITY OF BIRMINGHAM

Data Mining and Machine Learning

# Summary

- Introduction to front-end speech processing for ASR
  - Motivations from human hearing
  - Description of 'typical' front-end representation
    - Short-ti
    - Mel scal
    - Cosine transform
    - $\Delta$ and $\Delta^2$ parameters

Assignment Project Exam Help

https://eduassistpro.github.io/

Add WeChat edu_assist_pro

Data Mining and Machine Learning

UNIVERSITY OF BIRMINGHAM