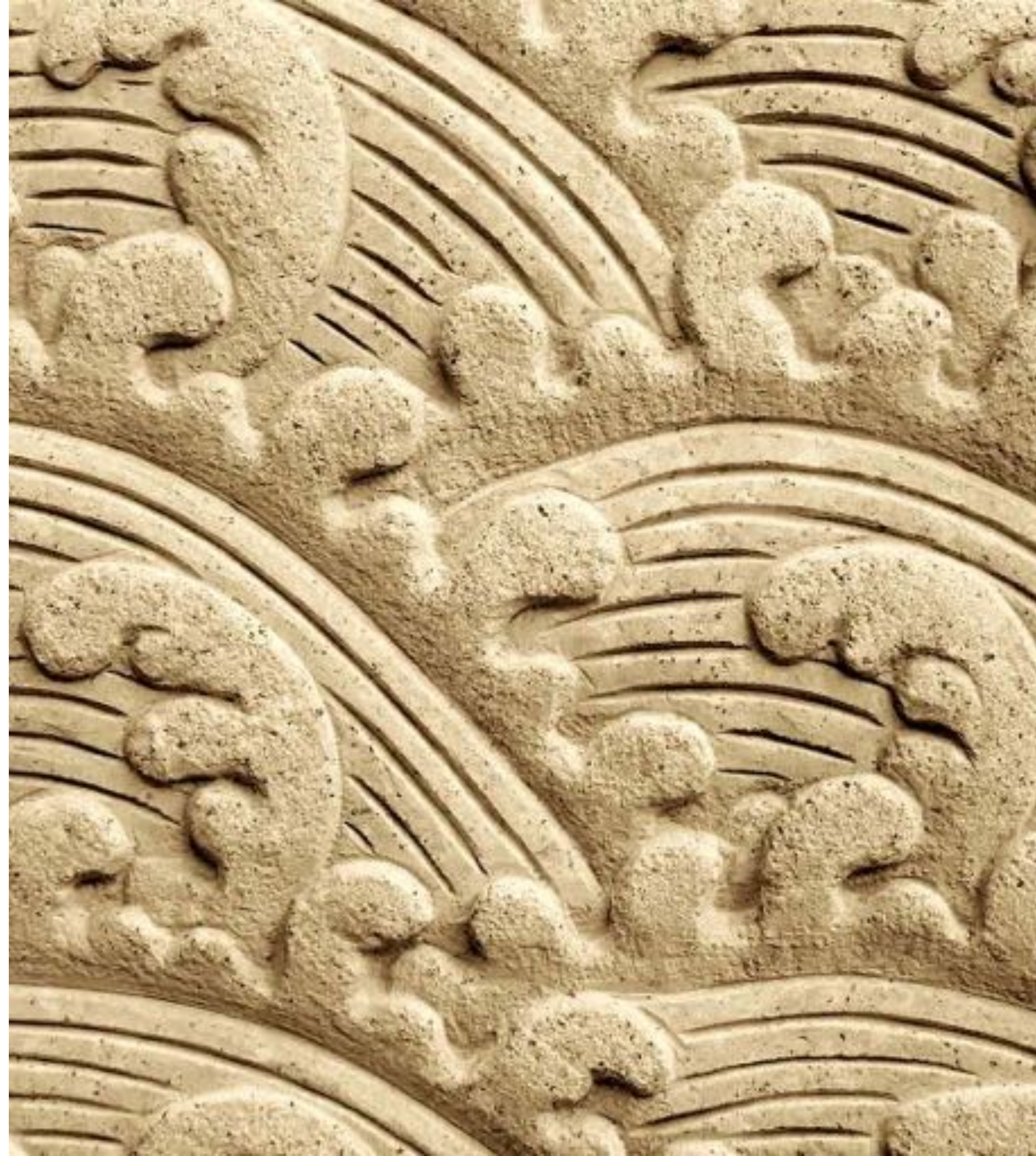# Speech Processing

Simon King
University of Edinburgh

2022-23

# Module 7

---

## Pattern matching

# Orientation

- We're on a journey towards HMMs

  - Pattern **matching**

  - Extracting **features** from speech

  - Probabilistic **generative** modelling

*What we are learning along the way*

Dynamic programming

(in the form of Dynamic Time Warping)

The interaction between
- choice of model
- choice of features

Dynamic programming

(in the form of the Viterbi algorithm)

# What you should already know

- Why the **waveform** is not good for pattern recognition

- Concept of a **feature vector**

- Let's start as simple as possible: whole word templates
  - But we already have to deal with sequences of **different lengths**

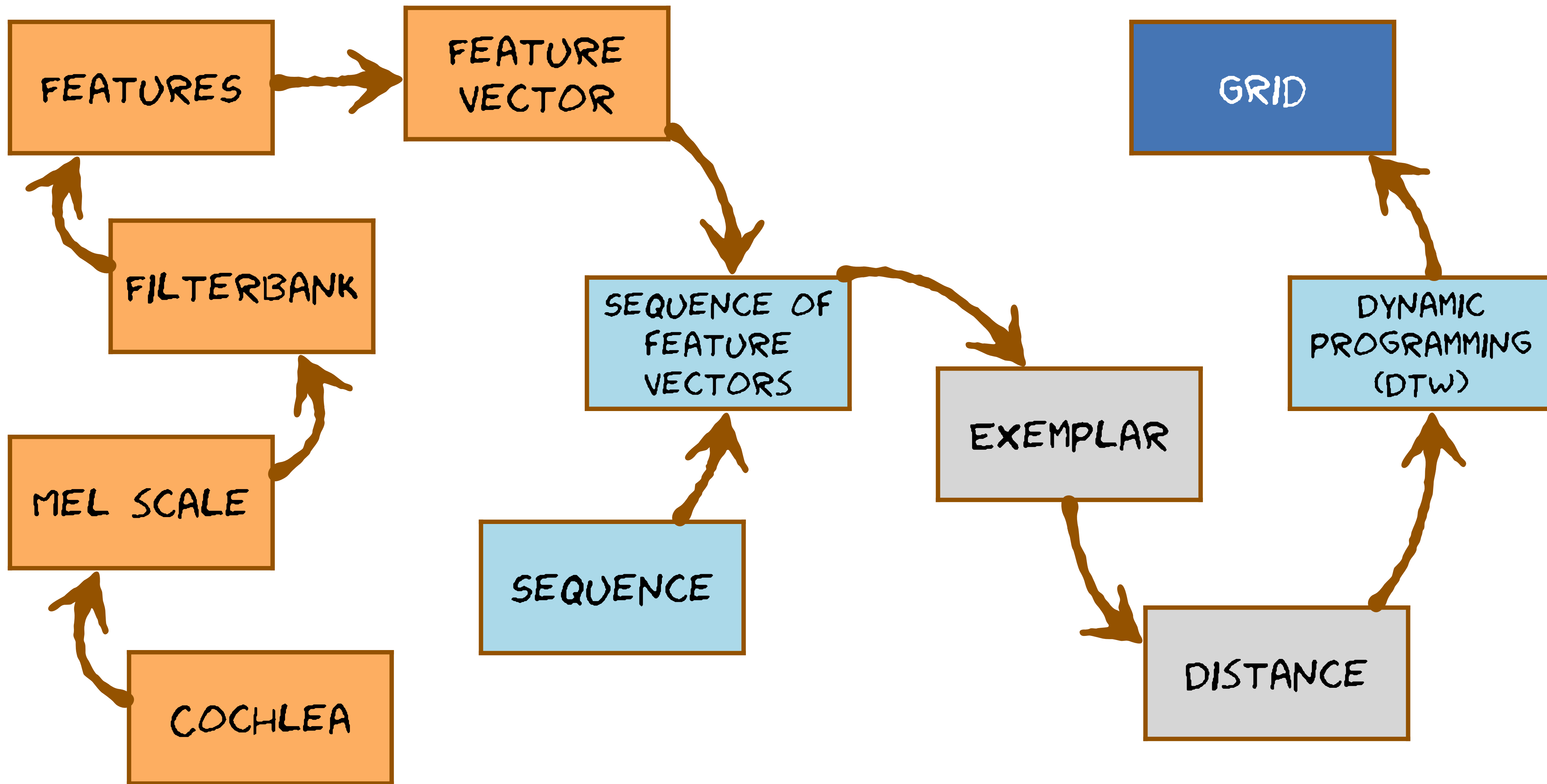Source and filter are combined

But we only want the **filter**

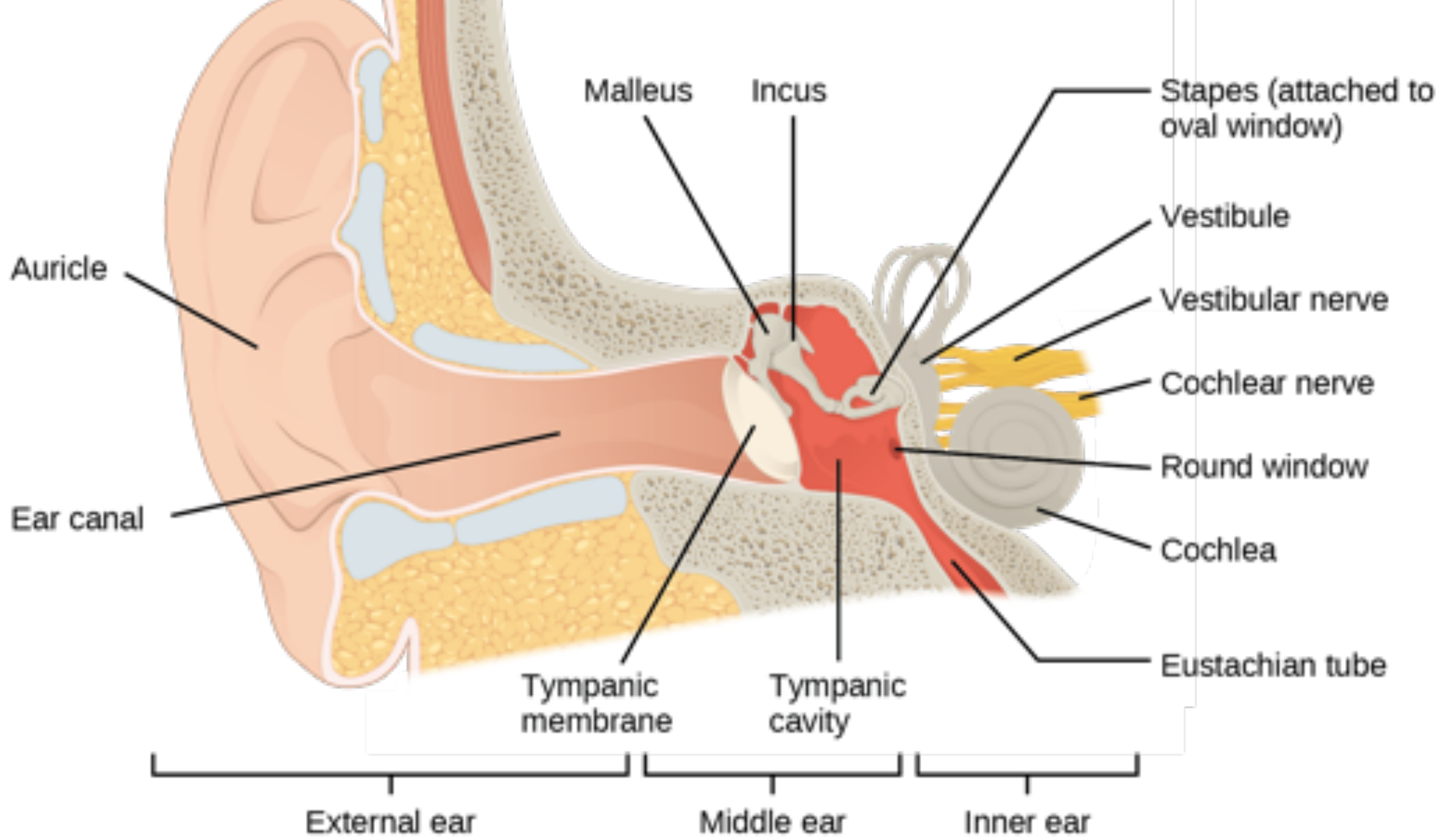Speech waveforms change over time

Use short-term analysis
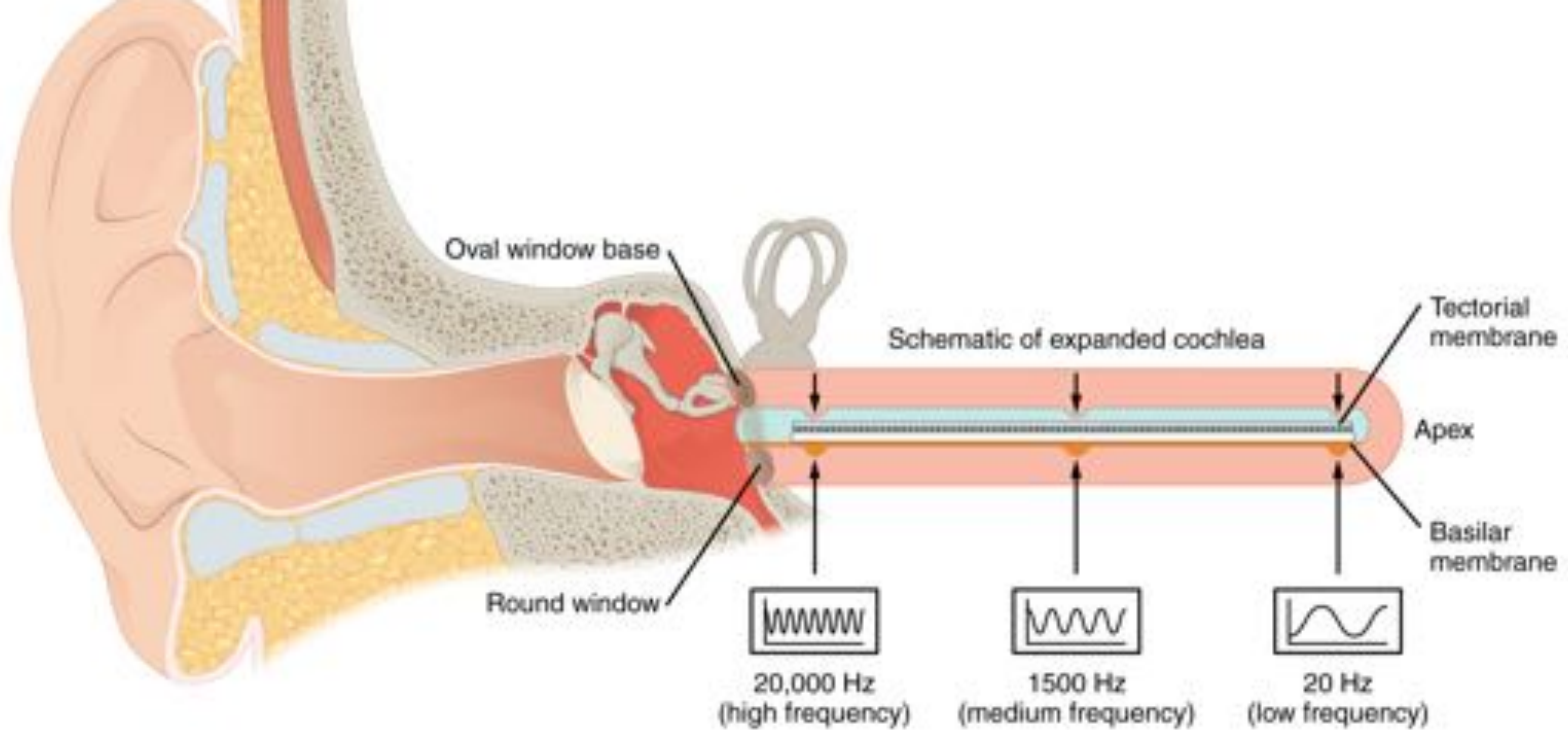
Extract features from frames of speech

Finding an alignment between two sequences

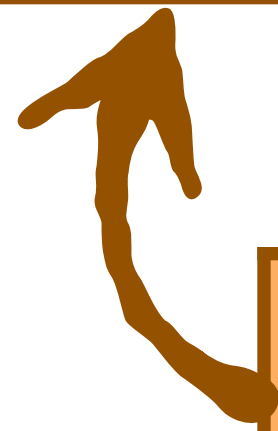- linear time warping

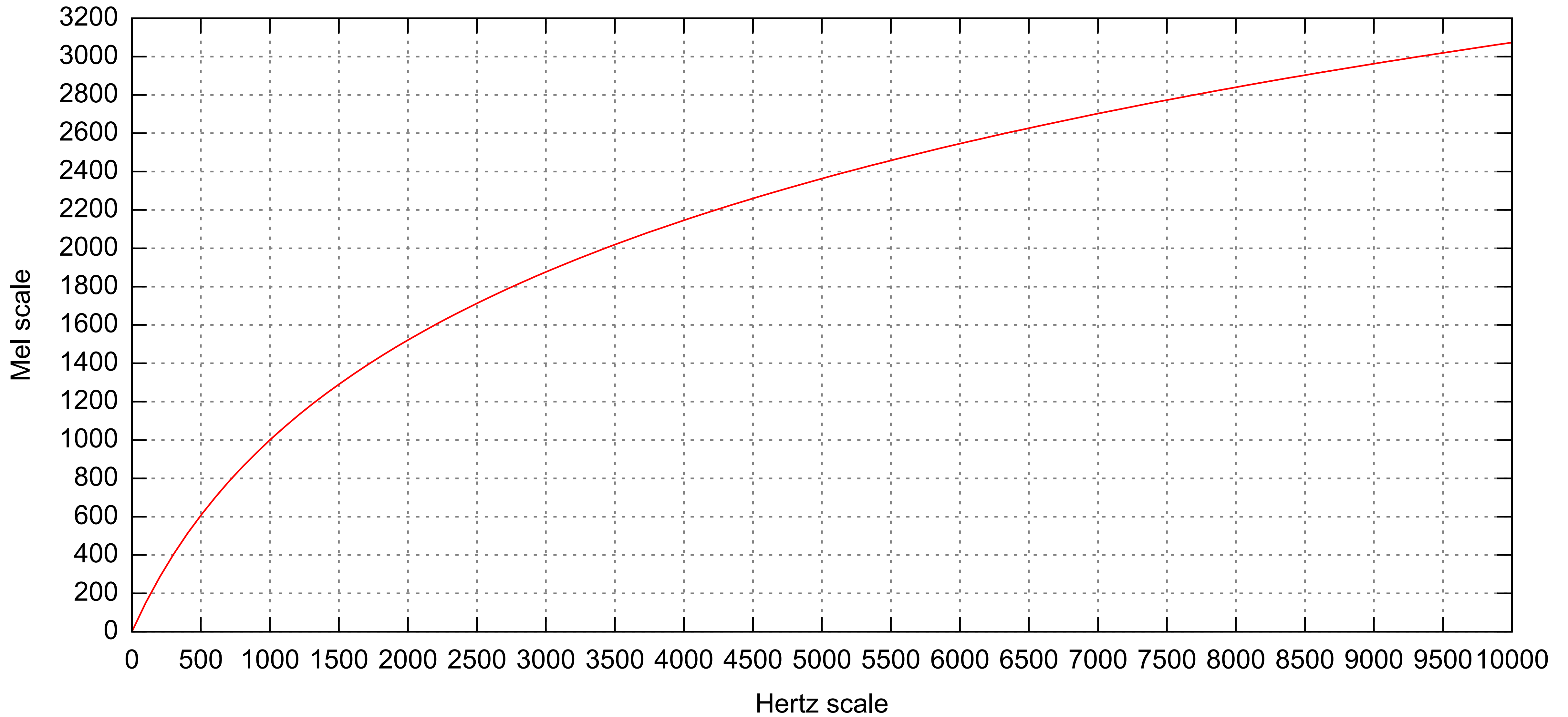- non-linear ('dynamic') time warping

COCHLEA

Malleus  Incus

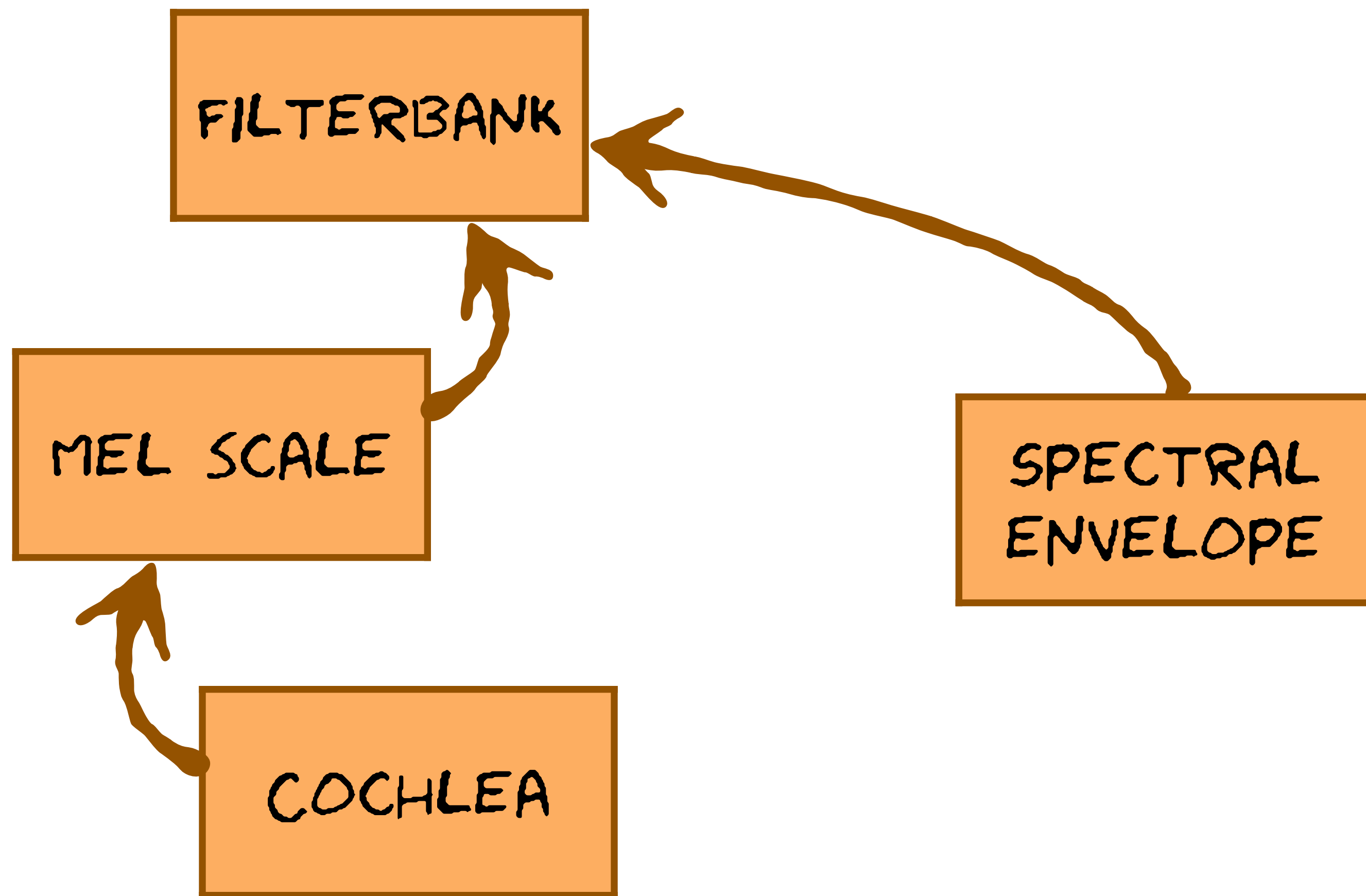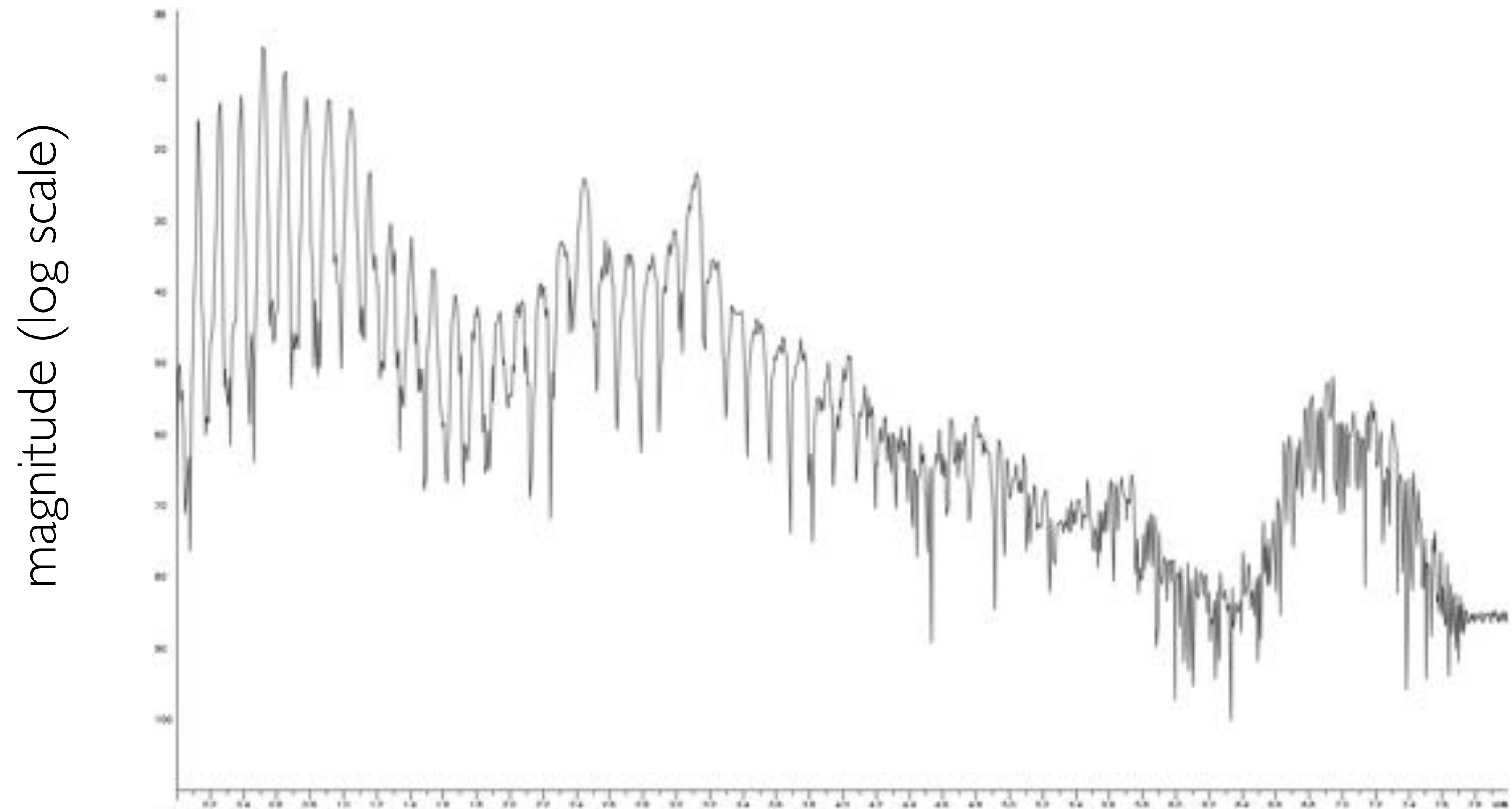Stapes (attached to oval window)

Vestibule

Vestibular nerve

Cochlear nerve

Round window

Cochlea

Eustachian tube

Auricle

Ear canal

Tympanic membrane

Tympanic cavity

External ear

Middle ear

Inner ear

Oval window base

Schematic of expanded cochlea

Tectorial membrane

Apex

Basilar membrane

Round window

20,000 Hz (high frequency)

1500 Hz (medium frequency)

20 Hz (low frequency)

MEL SCALE

COCHLEA

# The auditory system is like a bank of bandpass filters: a "filterbank"



magnitude (log scale)

0                    frequency                    8kHz

FEATURES

FILTERBANK

MEL SCALE

COCHLEA

# Each filter's output is a useful feature for doing Automatic Speech Recognition



magnitude (log scale)

frequency

0

8kHz

# Filterbank features for one frame are speech are stored in a single vector
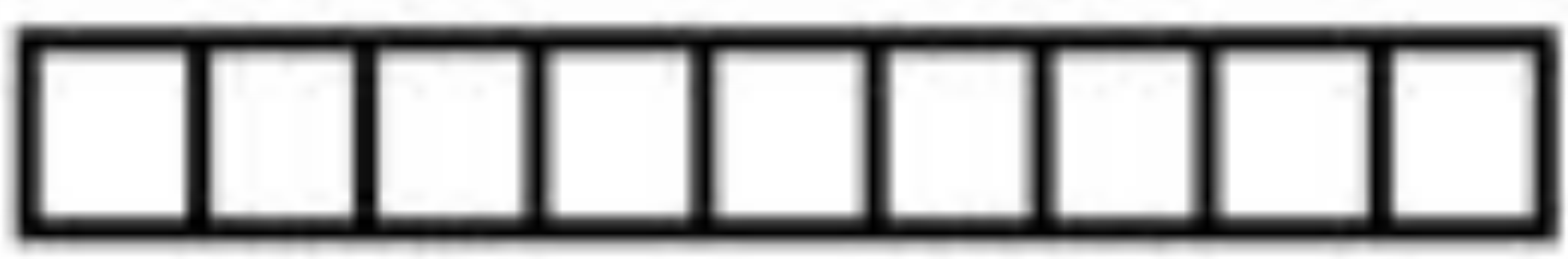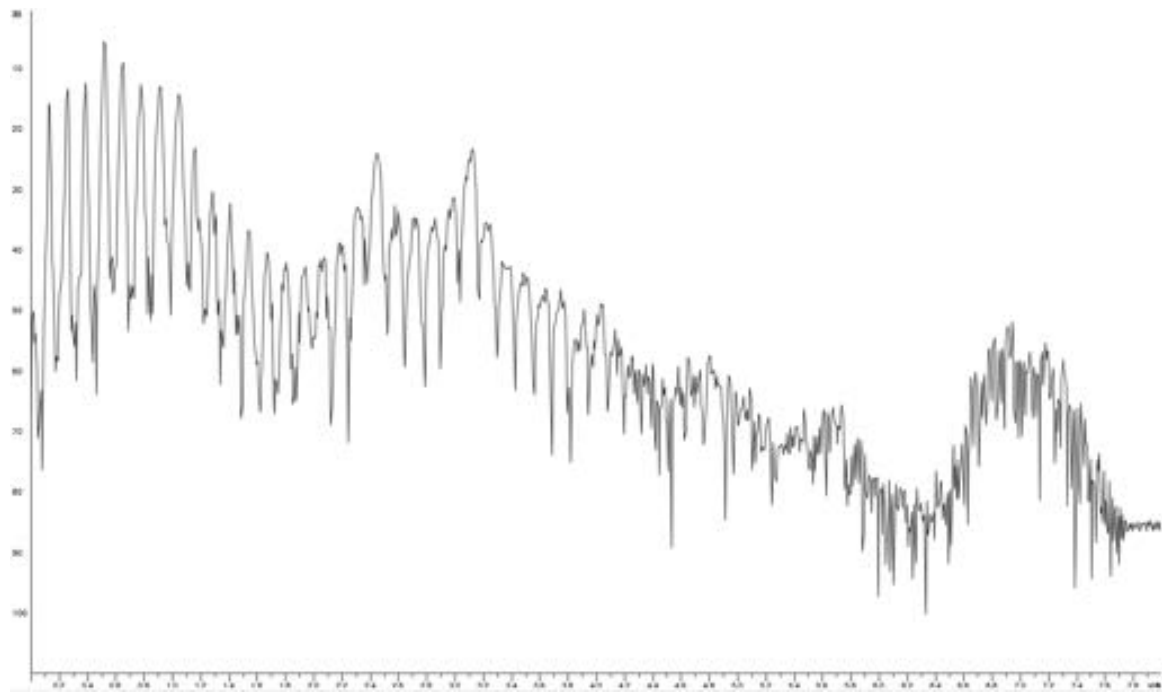


TIME DOMAIN

FEATURES

FFREQUENCY DOMAIN

FEATURE VECTOR

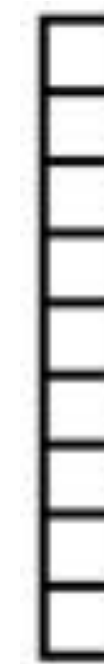# **Sequences** are *everywhere* in language
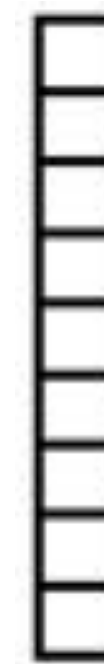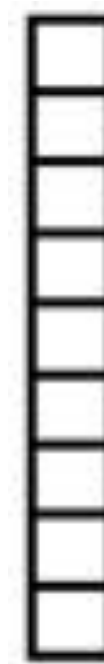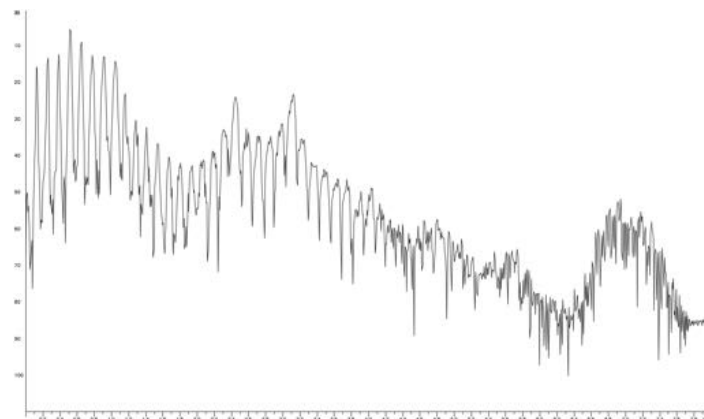
- <u>We've already seen</u>
  - a waveform is a sequence of **samples**
  - a waveform can be analysed as a sequence of overlapping analysis **frames**
  - a sentence is a sequence of **words**
  - a spoken word is a sequence of **phones**
  - a written word is a sequence of **letters**
- <u>Now we have</u>
  - from each frame we extract a feature vector
  - so a waveform becomes a **sequence of feature vectors**

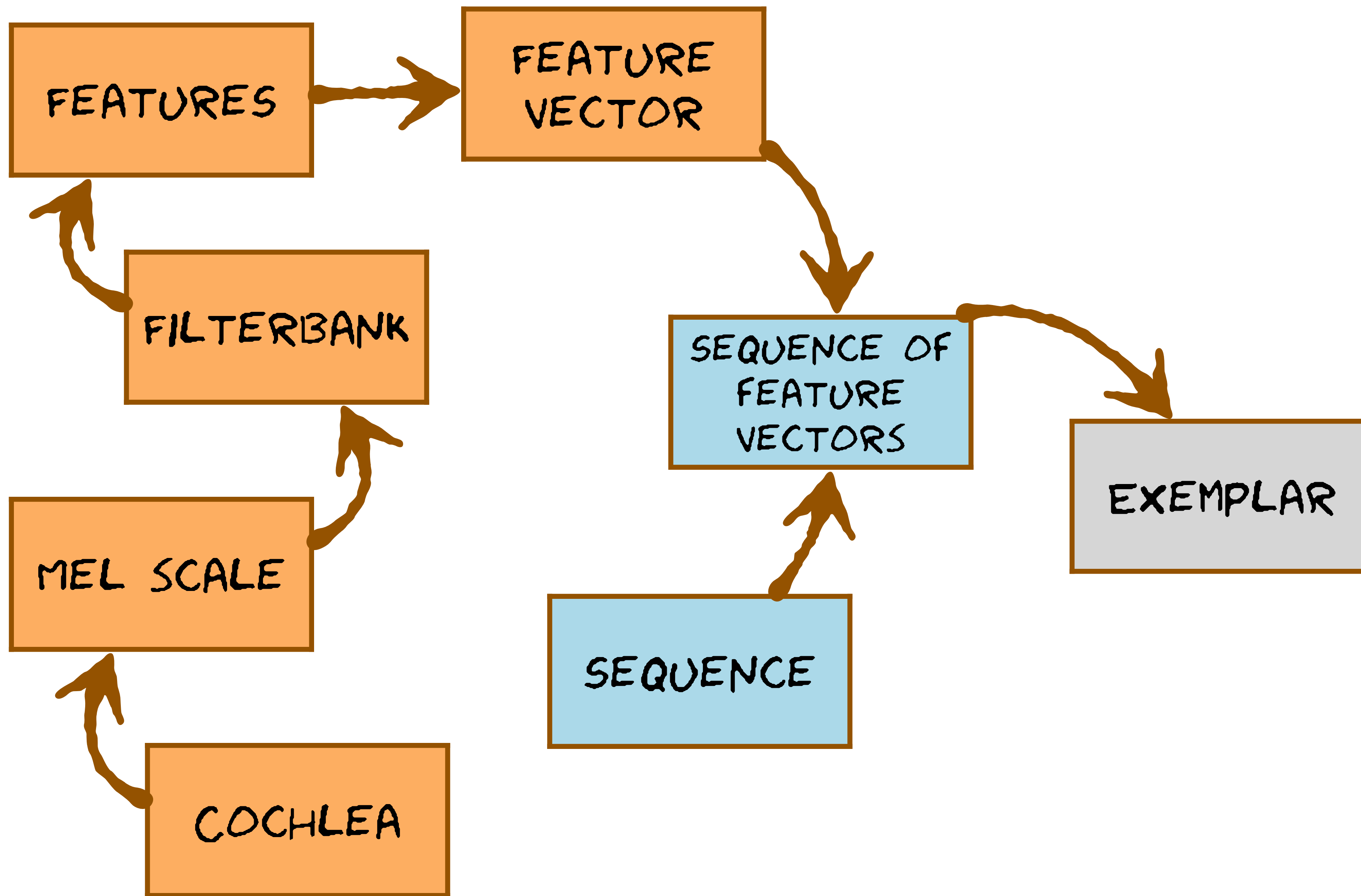# Filterbank features for one frame are speech are stored in a single vector

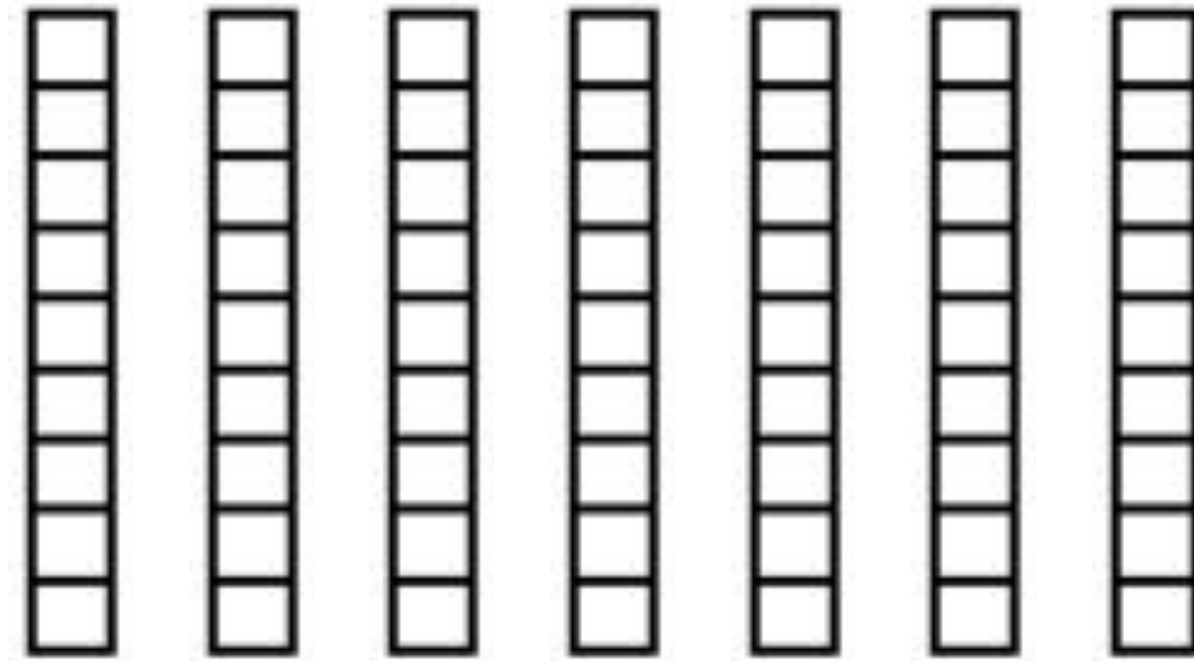# Filterbank features for automatic speech recognition
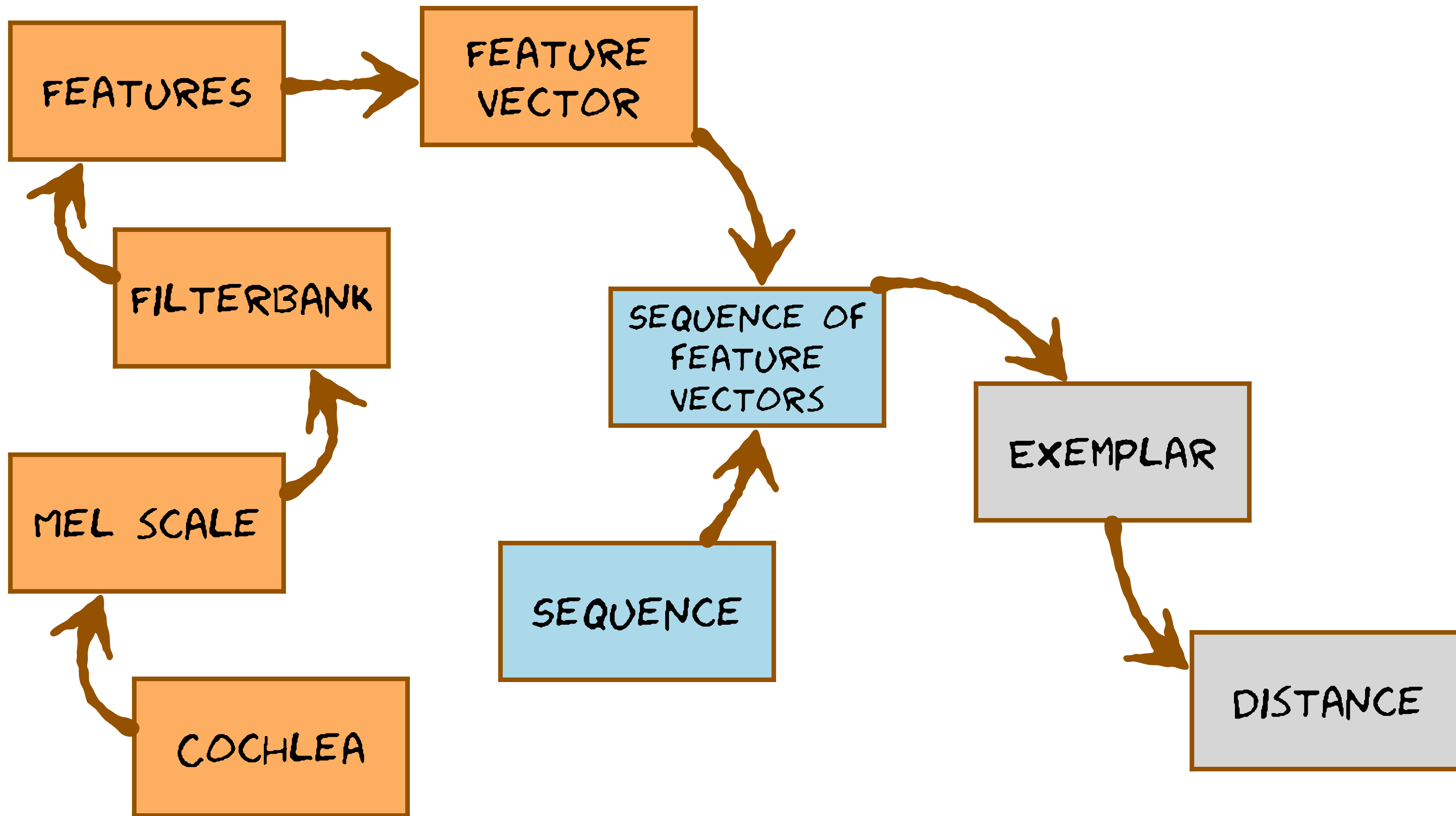


SEQUENCE OF FEATURE VECTORS

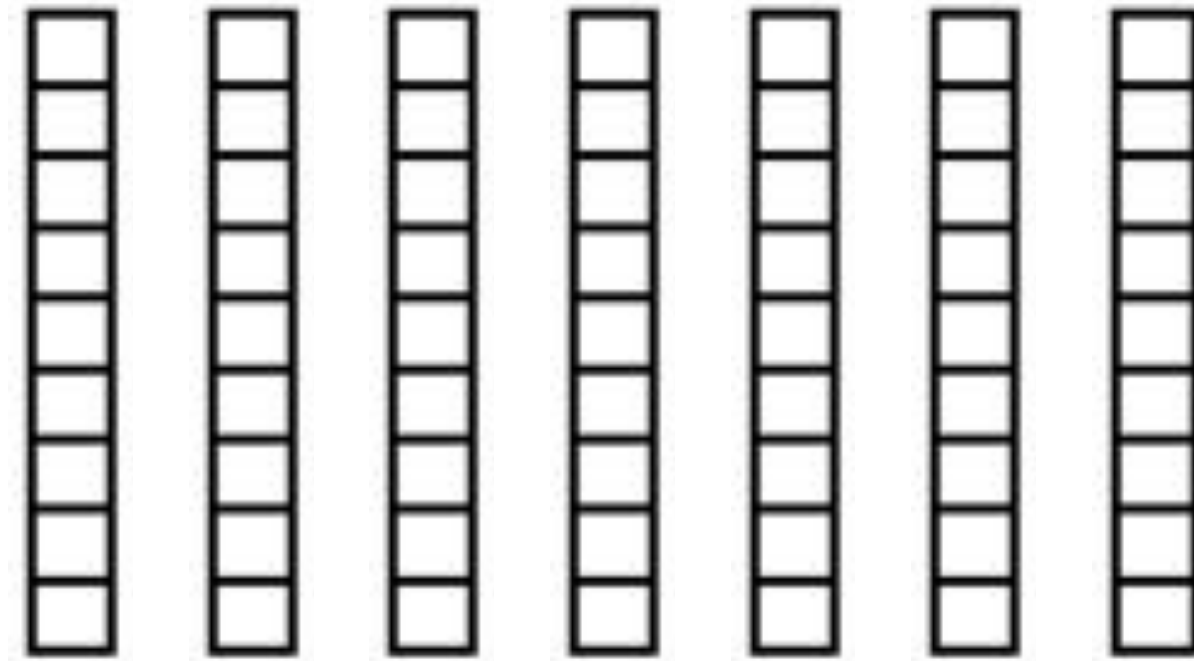# Filterbank features for automatic speech recognition

# "three"

# "three"

"three"

global distance
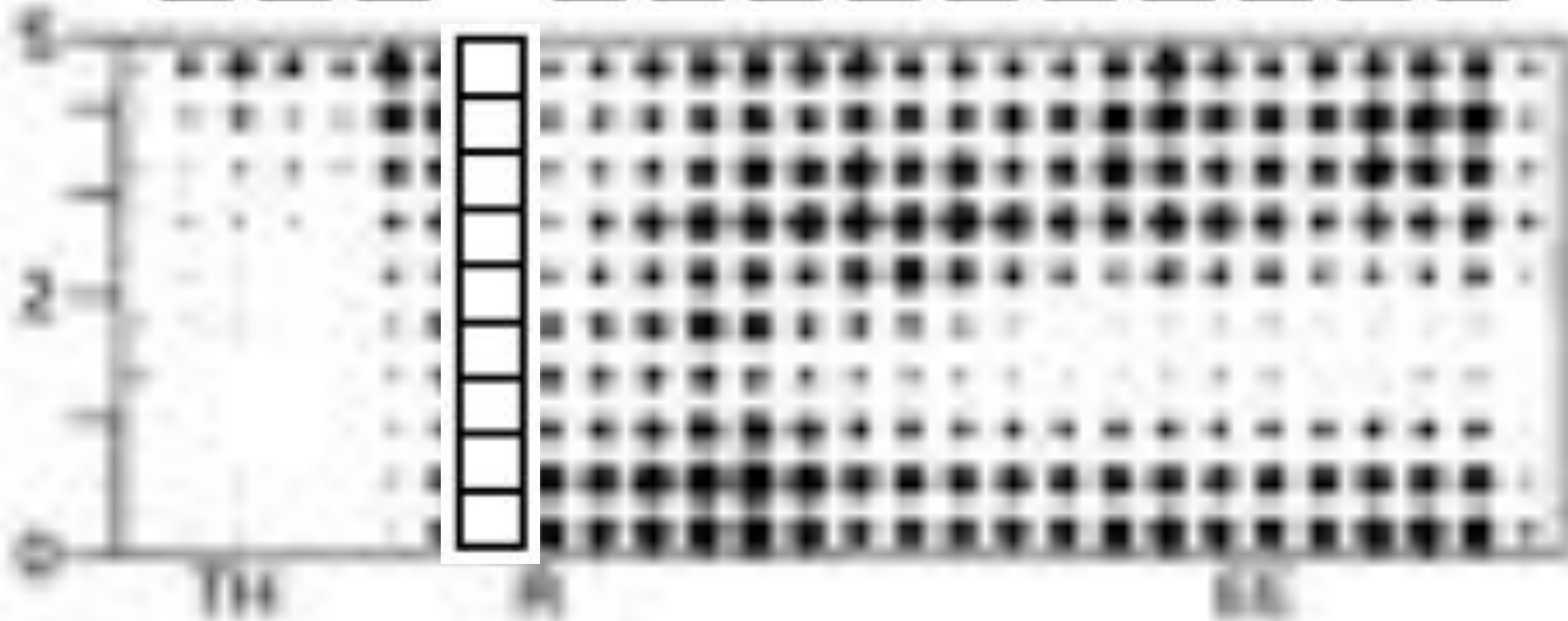
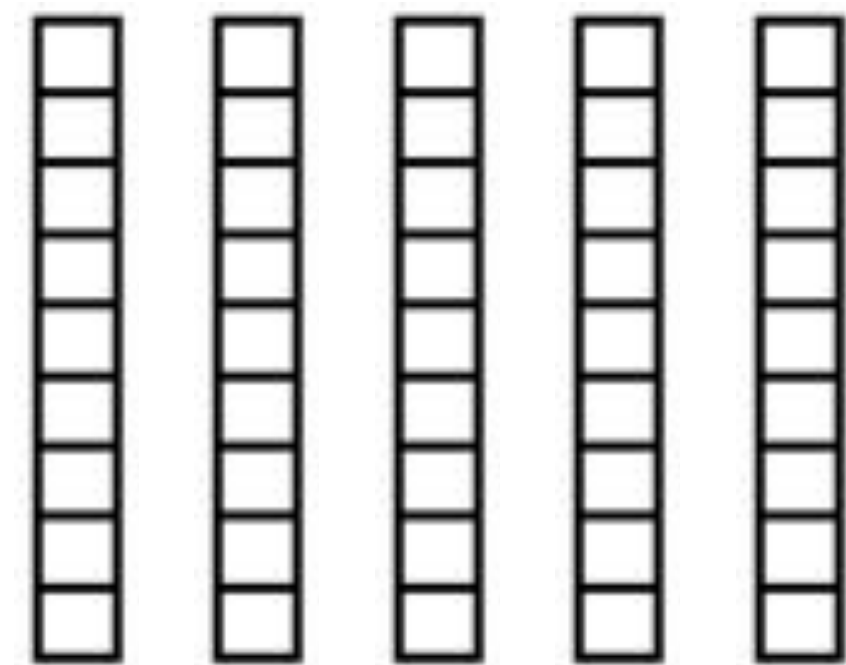$$= \sum \text{local distances}$$

"???"

Image credit: Figure 8.1 from Holmes & Holmes

# Pattern matching by Dynamic Time Warping

template

unknown

| | |
|---|---|
| 1, | 1 |
| 2, | 2 |
| 3, | 3 |
| **4,** | **3** |
| 5, | ? |
| 6, | ? |
| 7, | ? |

| | |
|---|---|
| 1, | 1 |
| 2, | 2 |
| 3, | 2 |
| **4,** | **3** |
| 5, | ? |
| 6, | ? |
| 7, | ? |

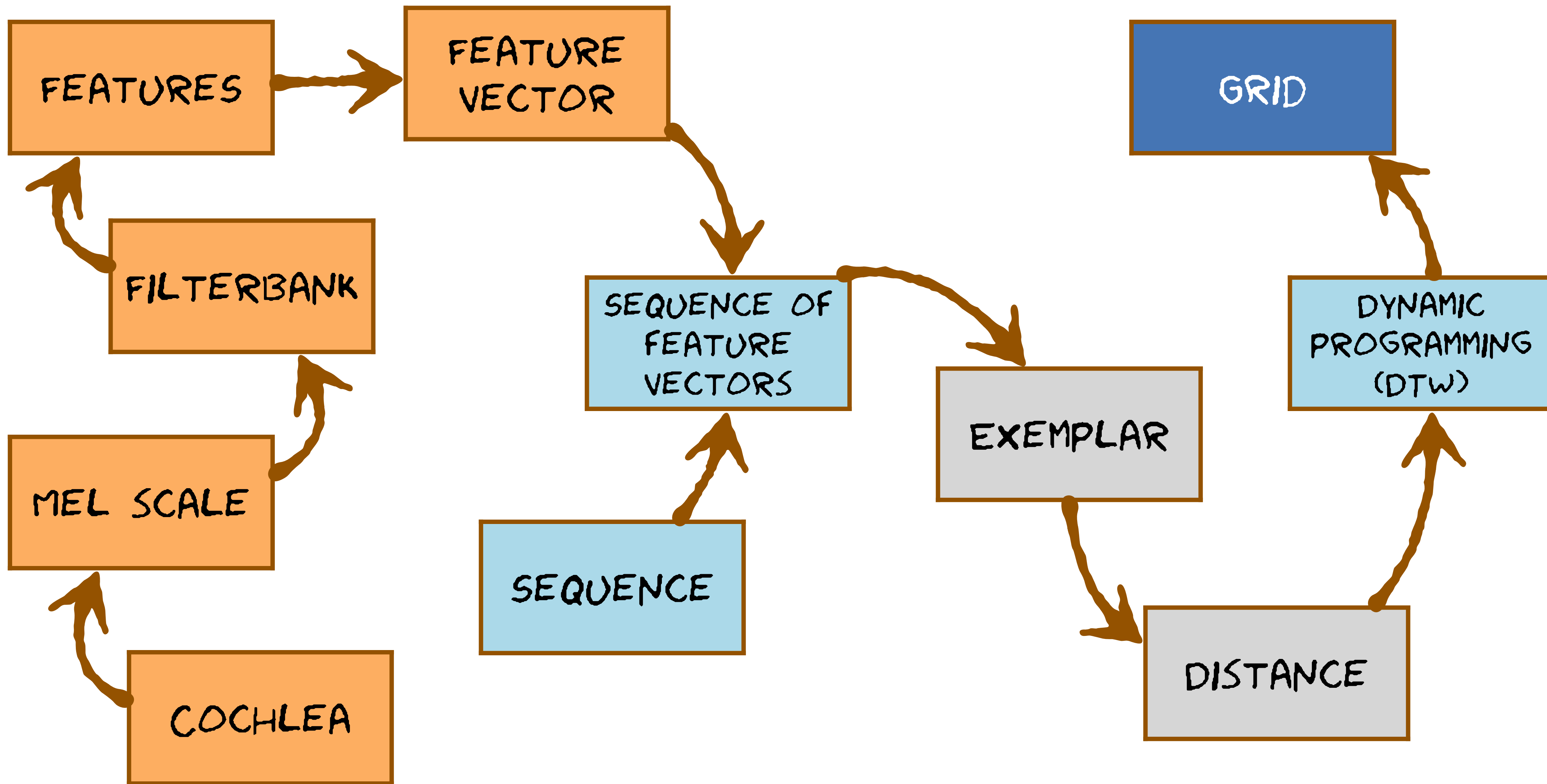| | |
|---|---|
| 1, | 1 |
| 2, | 1 |
| 3, | 2 |
| **4,** | **3** |
| 5, | ? |
| 6, | ? |
| 7, | ? |

# Dynamic Time Warping is a form of Dynamic Programming

- Understanding Dynamic Programming, as an algorithm

**Getting harder**

- Being able to see that Dynamic Programming can be applied to a particular problem

**Really quite difficult**

- Devising a suitable data structure for that problem

**My brain hurts**

# "three"



# "???"

SEQUENCE OF FEATURE VECTORS

GRID

global

DISTANCE

local

SEQUENCE OF FEATURE VECTORS

1 , 1
2 , 2
3 , 3
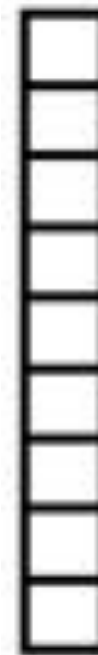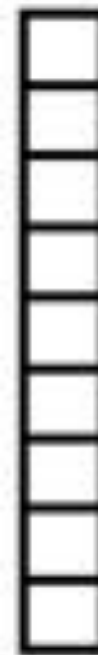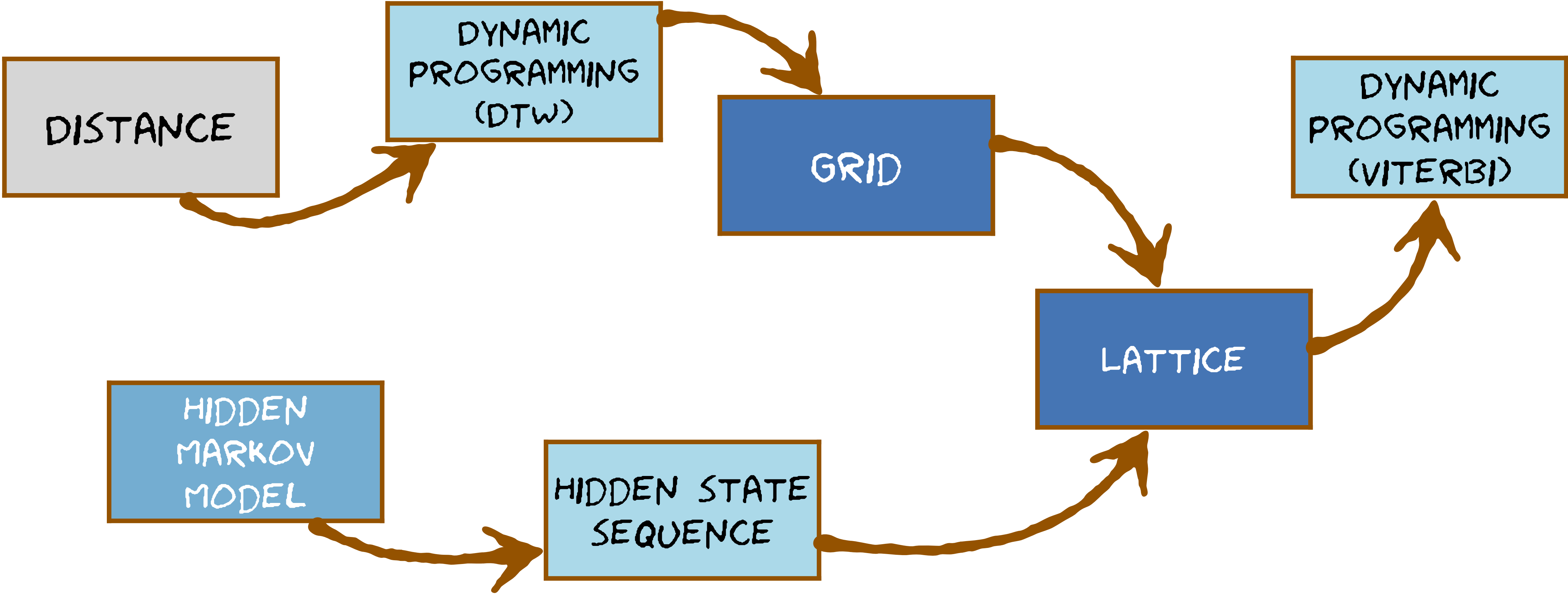**4 , 3**
5 , ?
6 , ?
7 , ?

# What you can learn next

# What next?

- DTW, and especially the local distance measure doesn't account for **variability**

  - so we'll replace it with a **probabilistic model**

- That model will use Gaussian probability density functions

  - to make these simpler, we will first try to remove covariance from our **features**

  - time for some **feature engineering** !

HMMs in Module 9

MFCCs in Module 8