

# MULTIMEDIA CLASS 2016 TA CLASS

SLIDE 3:

INTRODUCTION TO SPEECH PROCESSING USING MATLAB

# PREFACE, WHAT DO WE LEARN?

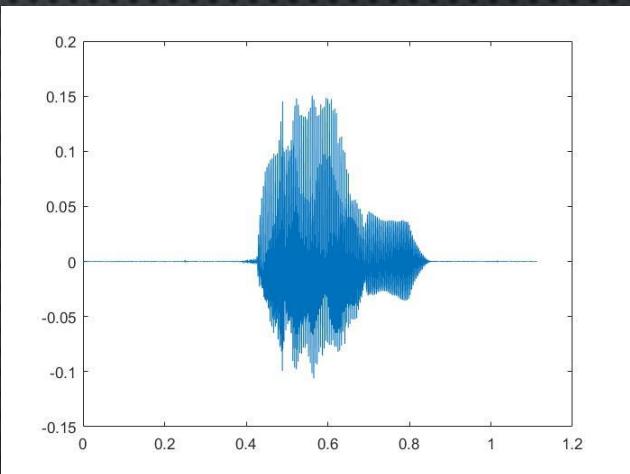
- IN THIS CHAPTER WE WILL LEARN A BIT MORE ABOUT SPEECH PROCESSING IN MATLAB
- YOU KNOW THE BASICS FROM YOUR COURSE
- SOME THE THINGS WE WILL LEARN IN THIS SESSION:
  - ALL ABOUT MFCC, WHAT IT IS AND HOW IT CAN BE USED
  - A BIT MORE ABOUT ADVANCED AUTOMATIC SPEECH RECOGNITION (ASR) AND ROBUST SPEECH PROCESSING (RSP)
  - TAKE IT DOWN A NOTCH: A SIMPLE 2-CLASS SPEECH RECOGNITION + IMPLEMENTATIONAL DETAILS (GUESS WHAT, NEW HOMEWORK IS UP!)
  - A BIT MORE ABOUT GMMS (GAUSSIAN MIXTURE MODELS), AND HOW WE CAN USE IT

# WHAT DO WE WANT TO ACHIEVE?

- WE SHOULD BE ABLE TO INTERACT WITH COMPUTER USING ONLY OUR BEAUTIFUL VOICES
- SIMPLEST FORM: WE SHOULD BE ABLE TO ANSWER THE COMPUTER WITH YES/NO
  - LIKE: WOULD YOU LIKE ME TO ORDER YOU PIZZA FOR LUNCH, YES, OK ORDER HAS BEEN PLACES
- HOW CAN COMPUTER UNDERSTAND WHAT WE JUST SAID?
  - JUST LIKE A BABY, USING PROBABILISTIC FRAMEWORKS
  - SHOW IT A COUPLE HUNDREDS OF EXAMPLES, IT WILL FIND THE COMMON DENOMINATOR
  - AT LEAST IT LEARN FASTER THAN YOU! (HOW OLD WERE YOU WHEN YOU LEARNED TO UNDERSTAND WHAT OTHER PEOPLE SAID?)

# WHAT IS AUDIO SIGNAL IN MATLAB?

- WHAT HAPPENS WHEN WE READ AN AUDIO FILE IN MATLAB?
  - THEY BECOME AND ARRAY (1-D FOR MONO AND 2-D FOR STEREO) OF DOUBLES
  - EACH NUMBER IS RELATIVE TO SOUND PRESSURE OF MICROPHONE AT THE TIME OF RECORDING
  - WE CAN PLOT IT EASILY, BUT IT IS A LOT OF DATA
  - YOU HAVE TO PERFORM RECOGNITION BASED ON AN 8908-D FEATURE DIMENSION
    - HUGE AMOUNT OF DATA FOR PROBABILISTIC APPROACHES



```
>> whos y
  Name      Size            Bytes  Class       Attributes
  y        8908x1           71264  double
```

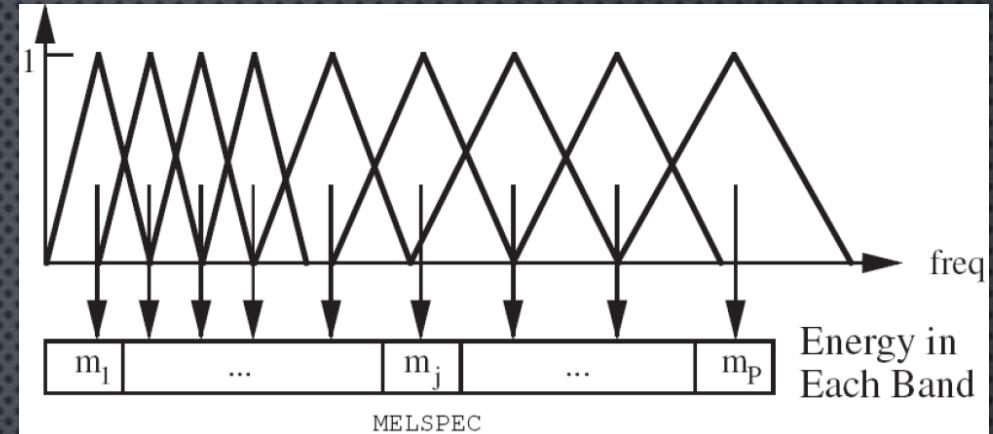
# FEATURE EXTRACTION (FE)

- OK, SO WE SAID THAT 8K+ DIMENSION ARRAY WITH VARIABLE LENGTH IS REALLY HARD TO PERFORM RECOGNITION ON
- SO, HOW DOES GOOGLE DO MILLIONS OF MILLIONS OF THEM ON A DAILY BASIS?
- FEATURE EXTRACTION IS THE ART OF EXTRACTING WHAT REALLY MATTERS FROM ALL THE DATA YOU HAVE
- IT WAS DONE UNTIL RECENTLY BASED ON HAND-CRAFTED FEATURES THAT MADE SENSE
- LATELY BY THE EVOLUTION OF CNNs, IT IS DESIGNED AUTOMATICALLY FOR YOU (IT IS CALLED END-TO-END LEARNING, JUST GOOGLE IT)

# FE, CONTINUED

- HERE COMES THE FE WE ARE GOING TO USE: MFCC
- IT STANDS FOR MEL-FREQUENCY CEPSTRAL COEFFICIENTS
  - WAIT A MINUTE, MEL SOUND FAMILIAR, WHERE HAVE WE HEARD IT?
  - YOU HAVE HEARD ALL ABOUT IT FROM DR. SEYEDIN, IN CHAPTER3: SPEECH PROCESSING
  - IN CASE YOU NEED TO REMEMBER IT AGAIN, TAKE A LOOK AT YOUR NOTES OR HANDBOOKS

# MFCC



- WHAT IS IT?
  - WELL IT SAYS HOW MUCH ENERGY IS THERE IS THERE IN EACH GIVEN FREQUENCY BAND
  - BANDS ARE DERIVED FROM HUMAN PERCEPTION OF AUDIO
  - IN CASE YOU FORGOT THE FORMULA FOR BANDS, GET BACK TO Dr. SEYEDIN PRESENTATIONS
- HOW IT IS DONE (PERFORMED ON EACH SIGNAL WINDOW) [WIKIPEDIA]:
  - 1- TAKE THE FOURIER OF THE WINDOWED SIGNAL
  - 2- USE TRIANGULAR OVERLAPPING WINDOWS TO MAP THE FREQUENCIES TO MEL-SCALE
  - 3- TAKE THE LOGS OF POWERS
  - 4- TAKE DCT OF THE LOGS OF POWERS (DON'T WORRY, YOU WILL READ ALL ABOUT IT IN DSP)
  - 5- MFCC IS THE AMPLITUDE OF THE RESULTING SPECTRUM
- WAIT A MINUTE, YOU SAID POWER??? WHAT ABOUT THE PHASE???
  - DON'T WORRY, AS I SAID BEFORE PHASE IS NOT VERY IMPORTANT IN SPEECH PROCESSING

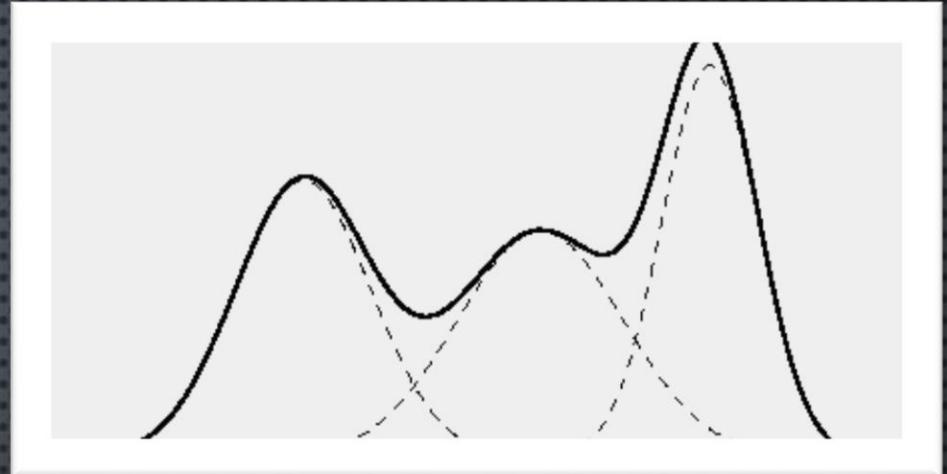
# MFCC, CONTINUED

- OK, I DON'T REALLY UNDERSTAND, YOU SAID SIMPLE!
  - IT IS OK, THAT'S WHERE LIBRARIES COME USEFUL
  - RASTAMAT IS A REALLY USEFUL LIBRARY MADE FOR SPEECH PROCESSING IN MATLAB
    - CHECK THE END OF THIS TUTORIAL FOR LINKS
  - WE WILL LET RASTAMAT DO THE HEAVY LIFTING FOR US, WE WILL JUST USE ITS RESULTS
- SO WHAT DO YOU NEED TO DO?
  - BE PATIENT, GMM IS UP

# MFCC IN MATLAB

- HERE IS THE CODE:
- `[MM,~] = MELFCC(Y, Fs, 'MAXFREQ', Fs/2, 'NUMCEP', 13 ...`
- `, 'NBANDS', 13, 'FBTYPE', 'FCMEL', 'DCTTYPE', 1, ...`
- `'USECMP', 1, 'WINTIME', 0.025, 'HOPTIME', 0.010, ...`
- `'PREEMPH', 0, 'DITHER', 1);`
- NOW YOU CAN ASK ABOUT THE PARAMETERS...

# NEXT UP, GMM



- SO WHAT IS GMM?
  - GAUSSIAN MIXTURE MODEL IS A WAY TO EASILY MODEL A PDF, BY SAYING IT IS THE ADDITION OF SOME GAUSSIANS
  - LOOK AT THE PICTURE ABOVE, SOLID LINE IS THE ORIGINAL PDF, AND DOTTED LINES ARE THE GAUSSIANS ADDED TO BUILD IT
  - IT IS VERY USEFUL IN AUDIO AND SPEECH PROCESSING, SINCE IT MODELS THE DATA WELL
  - A COMPLETE PDF RESOLVES TO SOME MEANS AND VARIANCES, SO IT BECOMES MUCH EASIER TO KEEP TRACK OF THOUSAND OF MODEL FOR DIFFERENT WORDS
- I ADDED A TUTORIAL IN PDF TO YOUR GIT, TAKE A LOOK AT IT

# COMPLETE SIMPLE SPEECH RECOGNITION, TRAINING

1. SEGMENT EACH AUDIO FILE INTO WINDOWS
2. PERFORM MFCC FEATURE EXTRACTION ON EACH WINDOW
  1. IN THIS SECTION YOU HAVE LOTS OF LOTS OF 13-D ARRAYS FOR EACH WORD MODEL
3. FIT A GMM MODEL (3-4 MIXTURES) TO EACH CLASS
4. SAVE YOUR MODELS AND DONE!

# COMPLETE SIMPLE SPEECH RECOGNITION, TESTING

PERFORM THE FOLLOWING FOR EACH AUDIO FILE:

1. PERFORM STEPS 1 AND 2 IN TRAINING TO GET THE MFCCS
2. GET THE LIKELIHOOD OF BELONGING TO EACH CLASS FOR EACH WINDOW (VALUE IN [0-1])
3. TAKE THE LOG OF LIKELIHOODS
4. SUM THEM UP FOR EACH CLASS (RESULTS IN 2 NUMBERS, 1 FOR CLASS “YES” AND ANOTHER FOR CLASS “No”)
5. THE WORD WITH THE HIGHEST LIKELIHOOD IS YOUR ESTIMATION

# FINISHING

- DO WHAT YOU LEARNED FOR ALL THE AUDIO FILES IN THE TESTING SEQUENCE, REPORT YOUR RECOGNITION RATE
- TUNE THE PARAMETERS TO GET BETTER RESULTS
- TRY TO WRITE A FAST CODE
- OPTIONAL FEATURES:
  - VAD: VOICE ACTIVITY DETECTOR (REMOVE NOISE AT THE BEGINNING AND END OF EACH FILE)
  - FEATURE NORMALIZATION (CMVN: CEPSTRAL MEAN AND VARIANCE NORMALIZATION)
  - ADD DERIVATIVES OF FEATURES TO FE.M AND MAKE YOUR RESULTS GET BETTER
- OK, LET'S SEE WHAT WE'VE GOT IN MATLAB

# ACKNOWLEDGEMENTS

- LIBRARY FOR FEATURE EXTRACTION (MFCC AND ALL):
  - RASTAMAT LIBRARY BY DANIEL ELLIS
    - AVAILABLE AT: <HTTP://WWW.EE.COLUMBIA.EDU/LN/ROSA/MATLAB/RASTAMAT/>
  - HMM TOOLKIT BY KEVIN MURPHY
    - AVAILABLE AT: <HTTP://WWW.AI.MIT.EDU/~MURPHYK/SOFTWARE/HMM.HTML>
- SOME OF THE AUDIO FILES BELONG THE AURA DATASET, THEY HAVE BEEN GIVEN TO ME BY SPRL (SIGNAL AND SPEECH PROCESSING RESEARCH LAB) UNDER SUPERVISION OF DR. AHADI AND DR. SEYEDIN
- OTHER STUFF (USES OF THE FOLLOWING LIBRARY) HAVE BEEN WRITTEN BY ME A COUPLE OF YEARS AGO