

Phases of speech recognition: (This is like compiler phases - ASR is complex process, hence divided into many stages for processing).

Computer  
↓  
Micro-processor  
Computer  
(Audacity)  
to record  
sound

word probability  
queen .8

coin .9  
so, sentence is →

1. Signal processing (input is speech signal & o/p is feature vector) , phones = 46
2. Phonetic recognition (i/p is feature vector along with database of acoustic models that help in generation of phone lattice)
3. Word recognition (i/p is phone lattice & a database of lexicon, i.e. dictionary that comprises word models) o/p is the word
4. Task recognition (it has i/p the word lattice, and the database of grammar  
→ I have killed your coin. (so, correct is "queen".

-- output of this last stage is final sentence.

Feature vector: -

sample collection



Actual text  
which  
was  
spoken

- The speech signal is sampled  
every 25 msec

- its properties are sample

(39 nos)

Only these  
parameters go to  
next stages

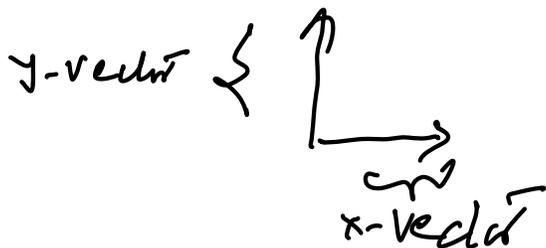
- frequency
- Amplitude
- Power
- ...

These parameters are called feature vectors

e.g.  $(x, y) = (5, 7)$  <sup>two vectors x, y</sup>

$(x, y, z) = (5, 7, 10)$  <sup>three vectors x, y, z.</sup>

→ 39 parameters means 39 vectors



The vectors take far less space, compared to the original signal; hence the vectors effectively provide compression (of data).

### Search Algorithms:

- for phonemes: Easy search, as only 46 phonemes. Search space is
- " Lexicon: medium search (1000s)
- " Sentence: difficult search as it requires to search out 7 millions of sentences

Let the sentence was:

"Recognize speech", is pronounced like,

like "urp ck a nite beesh" ✓

When this happens, we need to search  
a correct sentence. This is language

prob. ∴ we require a model of  
language

To recognize words  
individually, which is  
possible if words are  
sufficiently isolated

(a word does not affect next word)

We can simply search the words in lexicon

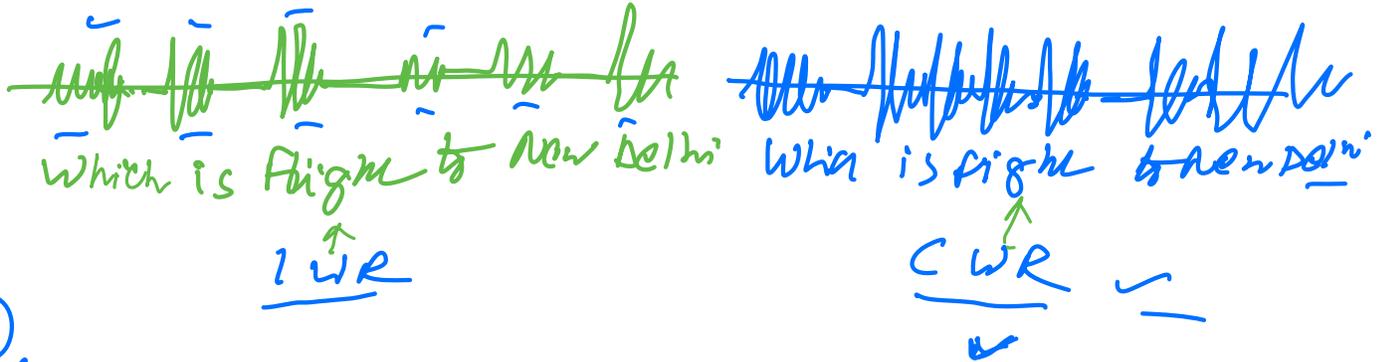
When words are isolated (slow speed) so that no word affects the pronunciation of any word the search space is less.

This is called - IWR - Isolated word recognition.

- Its search algorithm is simply the lexicon search.

When entire sentence is spoken at speed, the situation becomes that of "Recognize speech" is recognized as "wreck a nice beach". This requires language model based algorithm.

and the algorithm is called CWR - Continuous Word Recognition. This is time consuming also.

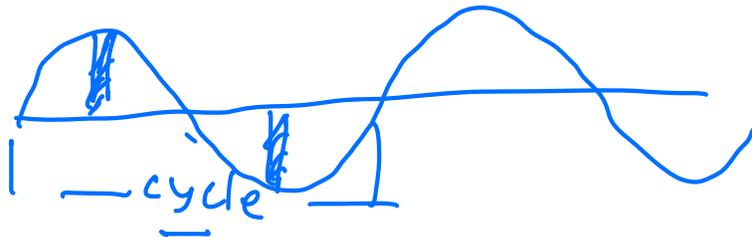


$$\frac{1}{4 \times 10^3} = 0.25 \text{ msec}$$

Domestic Supply: (India)  $\approx$  cycles/sec

$$2 \text{ kHz}$$

$$\frac{1}{2 \times 10^3 \times 2}$$



∴ sampling frequency is 2x signal frequency

