

Voice Recognition

By

Luigi Penaloza



Voice recognition, automatic/computer speech recognition is the process of converting acoustic signal captured by microphone into set of words.

Voice recognition:

- Recognize the voice of the individual
- Identifies what you are saying
- Performance of required task by computer

Voice Recognition History

Humans

vs

Computers

When we articulate, sound waves are produced and ear conveys to brain for processing.

- Digitization
- Acoustic Analysis of speech signal
- Linguistic interpretation

1940's and 1950's- US Department of Defense and MIT largest contributors to the voice recognition field.

1960's and 1970's- Technology attributed by educational institutions.

Voice Recognition now

- Built into our phones, games consoles and smart watches.
- Used to automate our house.
- Integrated into personal AI assistants, that allows us to shop online and make appointments.

All possible because deep learning has made voice recognition accurate enough.

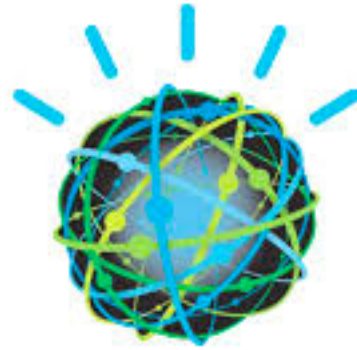
Right now we have about 95% word accuracy for the English language.



Famous Voice Assistants

- Siri (Apple)
- Duer (Baidu)
- Cortana (Microsoft)
- Alexa (Amazon)
- Google Now (Google)
- Watson (IBM)

amazon alexa



IBM WATSON

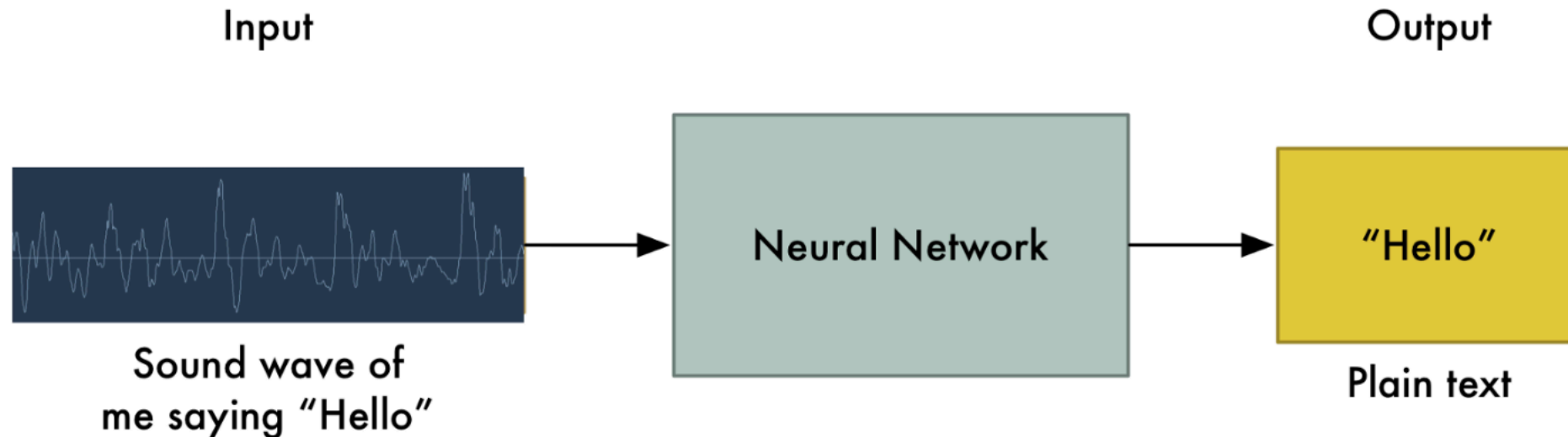
 Cortana



Hey Siri

How does it work?

- The user speaks some words by invoking voice recognition on a mobile app.
- The spoken words are processed by the recognition software and converted to text.
- The converted text is then provided as input to the search mechanism, which returns the results.



Difficulties with speech recognition

Main problem is that speech varies in speed.

For example: one person might say “hello Siri!”, and another might say “heeeelllllooooo Siiiiiriii!”

Both files should be recognized as same text, “hello Siri!” so we need to align audio files of different lengths to a fixed length piece.

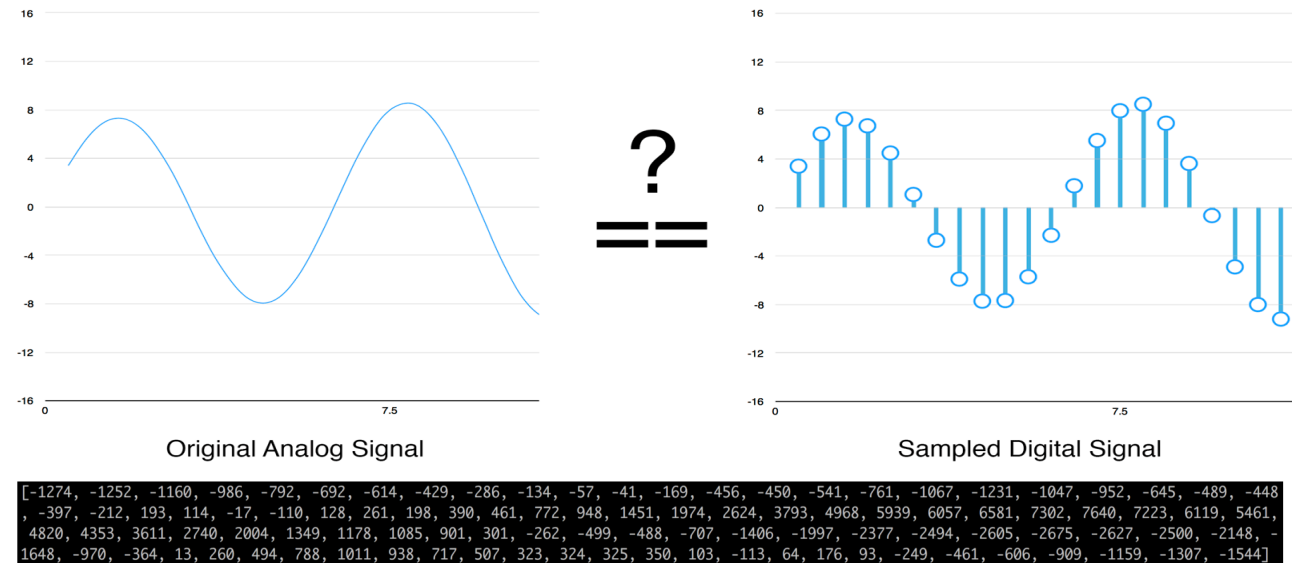
"Hey, Siri"



Turning sound into bits

Sampling

- We feed sound waves into a computer. Since sound waves are one dimensional, we can obtain a single value based on height of wave.
- We record the height of wave at equally spaced points.
- We can take a reading a thousand times of a second and record wavelength height number.
- 16,000 samples per second is enough to cover the frequency range of human speech.



Each number represents the amplitude of the sound wave at 1/16000th of a second intervals

Pre-processing sampled data

- We have an array of number representing the sound wave's amplitude.

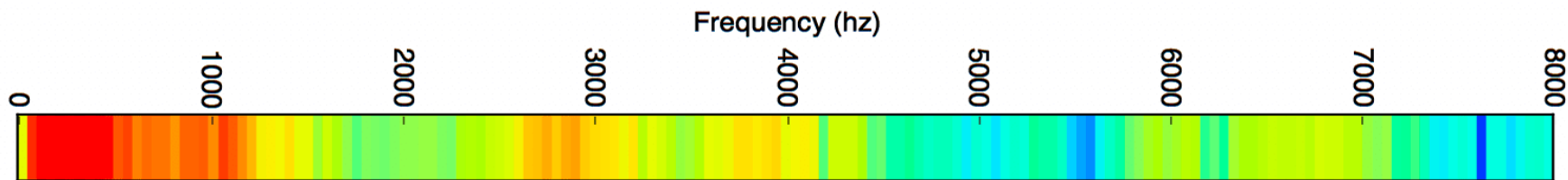
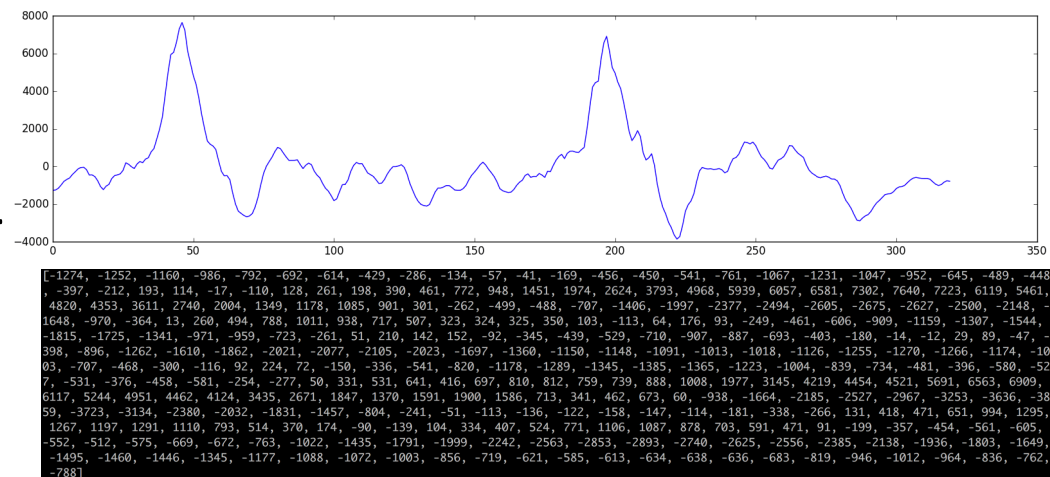
in our case $1/16,000^{\text{th}}$ of second intervals.

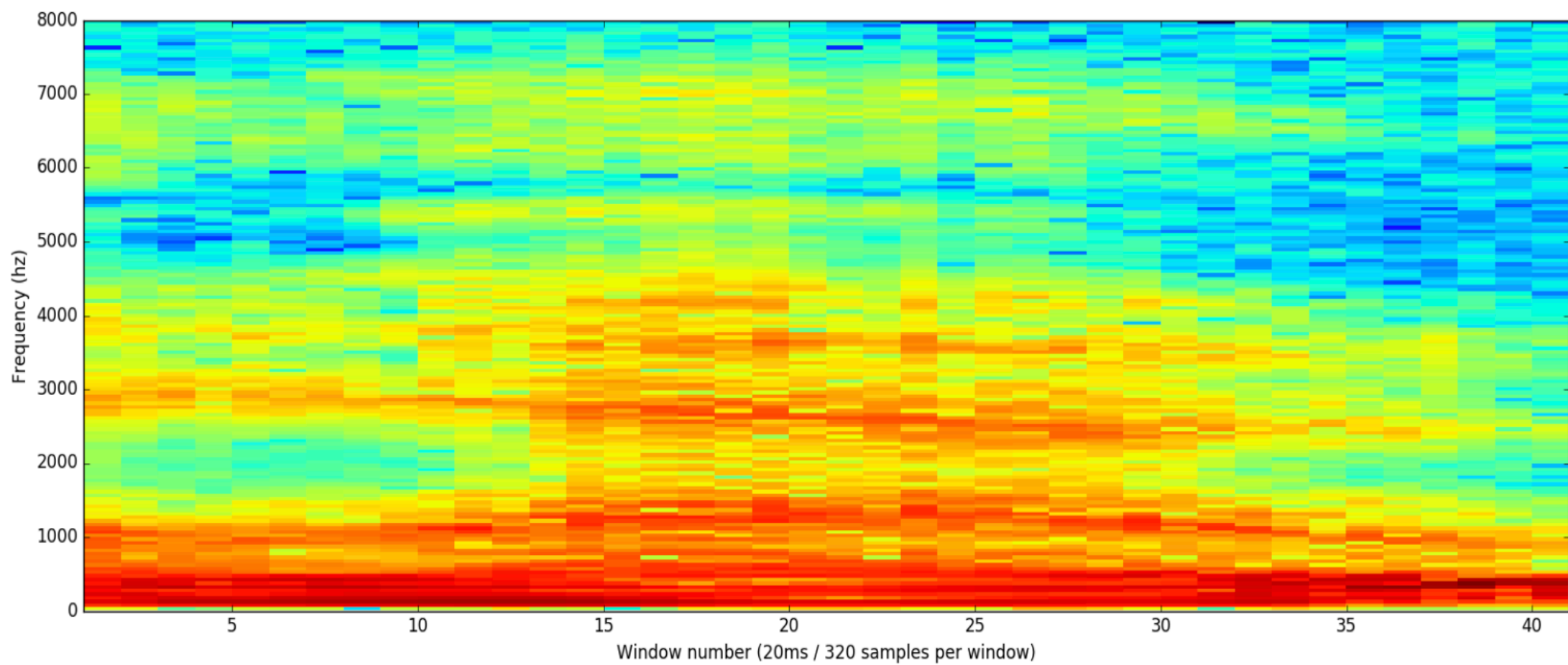
- We can preprocess the data and then feed numbers to neural network.

- First we group sampled audio into chunks
- Recordings usually mix of different frequencies, low and high pitched.

- We can use Fourier Transform which

breaks apart complex sound waves into simple sound waves. Then we add how much energy is contained in each wave.

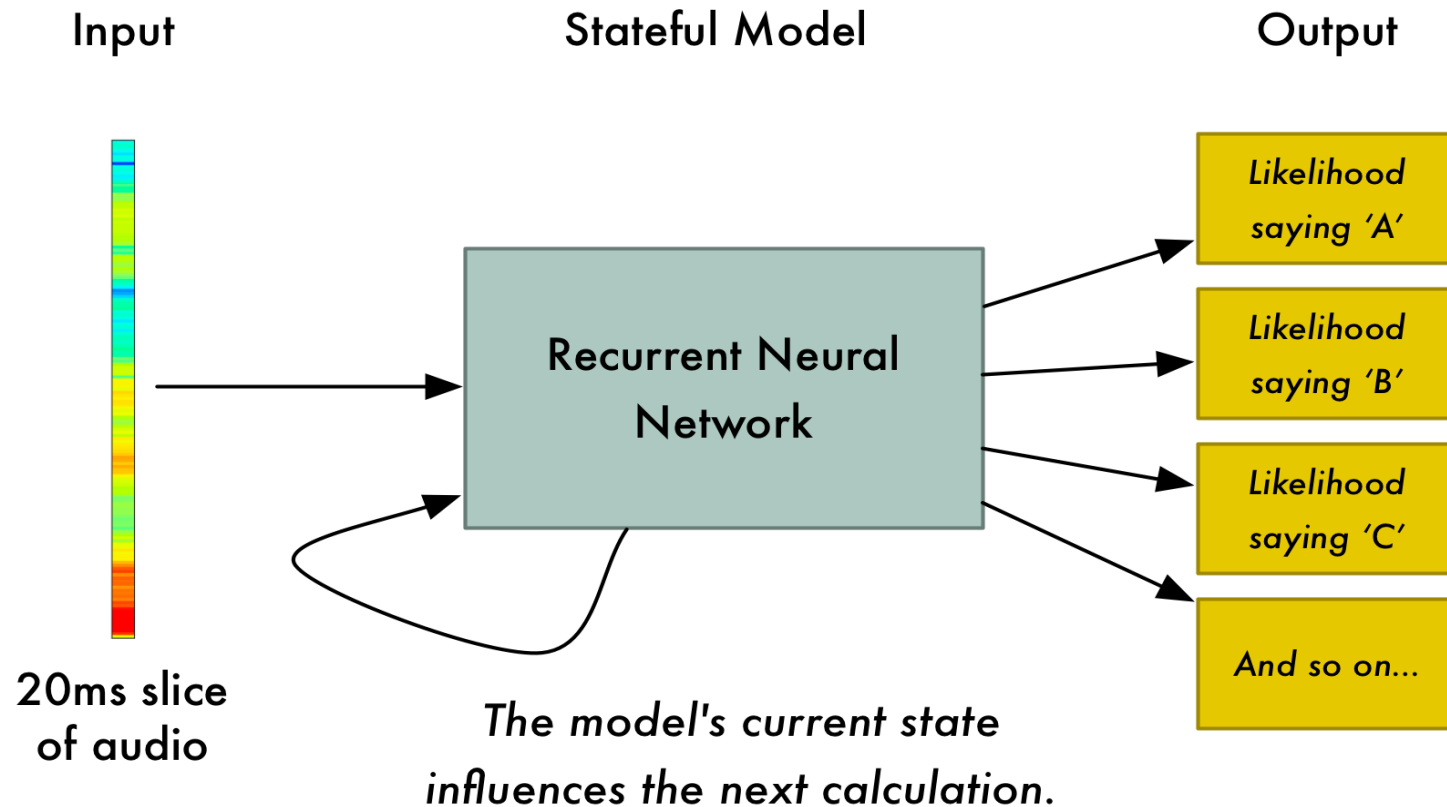




The full spectrogram of the "hello" sound clip

Recognizing characters from short sounds

- We'll feed audio in previous format to neural network
- Neural network input = 20 millisecond audio chunks
- For each audio, model will try to figure out the letter that corresponds to sound being spoken.





Most likely letter:
(per 20 milliseconds)

□ □ H H H E E _ L L _ _ L L L O O O □ □ □ □

- Neural network thinks we said
- “HHHEE_LL_LLLOOO” or “AAAUU_LL_LLLOOO”
- With cleaning:
- “HHHEE_LL_LLLOOO”-> “HE_L_LO”
- “AAAUU_LL_LLLOOO”-> “AU_L_LO”
- “HELLO” or “AULLO”
- **“HELLO”**

Future of Voice Recognition

Competitiveness between voice assistants, and mobile app development companies will drive voice recognition advancement ahead.





SOURCES:

<http://research.baidu.com>

<http://web.mit.edu>

<https://medium.com>

www.sas.upenn.edu

www.voicerecognition.com.au



Luigi Penaloza