

DLP - Week - 8

L1 Recap of Previous Discussions

- speech is much better than text.
- chatgpt is cd multi-modal LLMs.
- language identifiers
- speaker diarization
- automatic speech recognition

↳ this is the half part. What if I want to listen to result? That is TTS (Text To Speech)

L2 Intro to TTS

- Text to Speech conversion
- text → TTS → output
↳ is speech
↳ synthesized audio → natural and intelligible speech.
- Applications
 - 1. speech based technologies
 - 2. help people with literacy difficulties
 - 3. aid the visually challenged.

- Frameworks .

1. unit select " synthesis (USS)
2. hidden markov model based (HTS)
3. neural network based (conventional)
4. end to end (E2E)

- Phases

1. Training phase
2. synthesis / testing phase

L3. History of TTS

- long history
 - in 1990, \rightarrow first female voice synthesizer (successful)
 - now very powerful to achieve TTS .
 - speech has phonemes (they were considered as states) \rightarrow to approx. the acoustic vectors .
 - DB
- <text , audio> pairs \rightarrow continuous speech.

L4. Analogous of ASR & TTS

Text

How are you?

characters

Speech

~~How are you~~

phonemes

in spoken language
gets concatenated
by them.

liy | → ſ

l ih | → ſ

ASR
↓
editable
text

Today, we talk about
the reverse problem.

Text



TTS



How
Speech

maybe gender

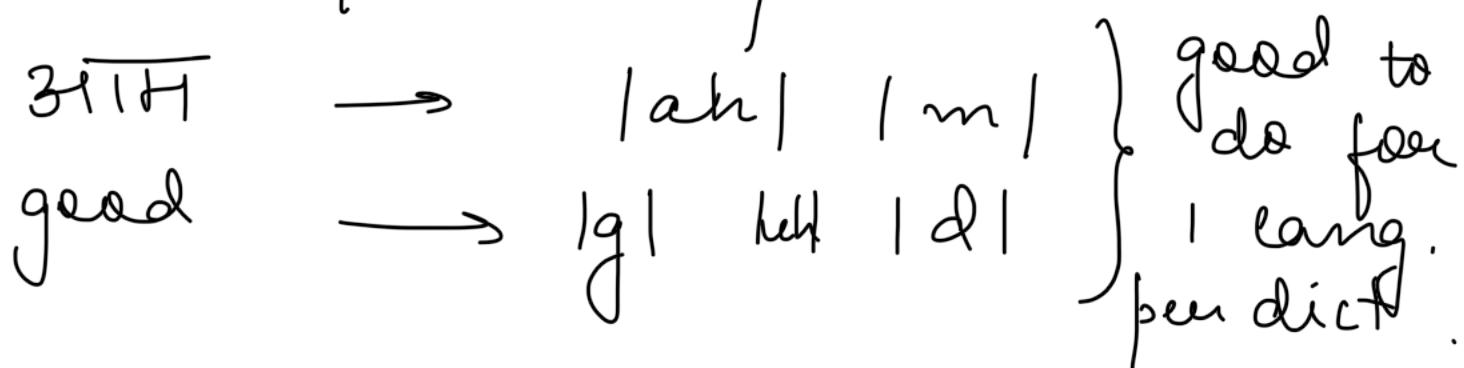
big actors

deepfake (ethnic
issues)

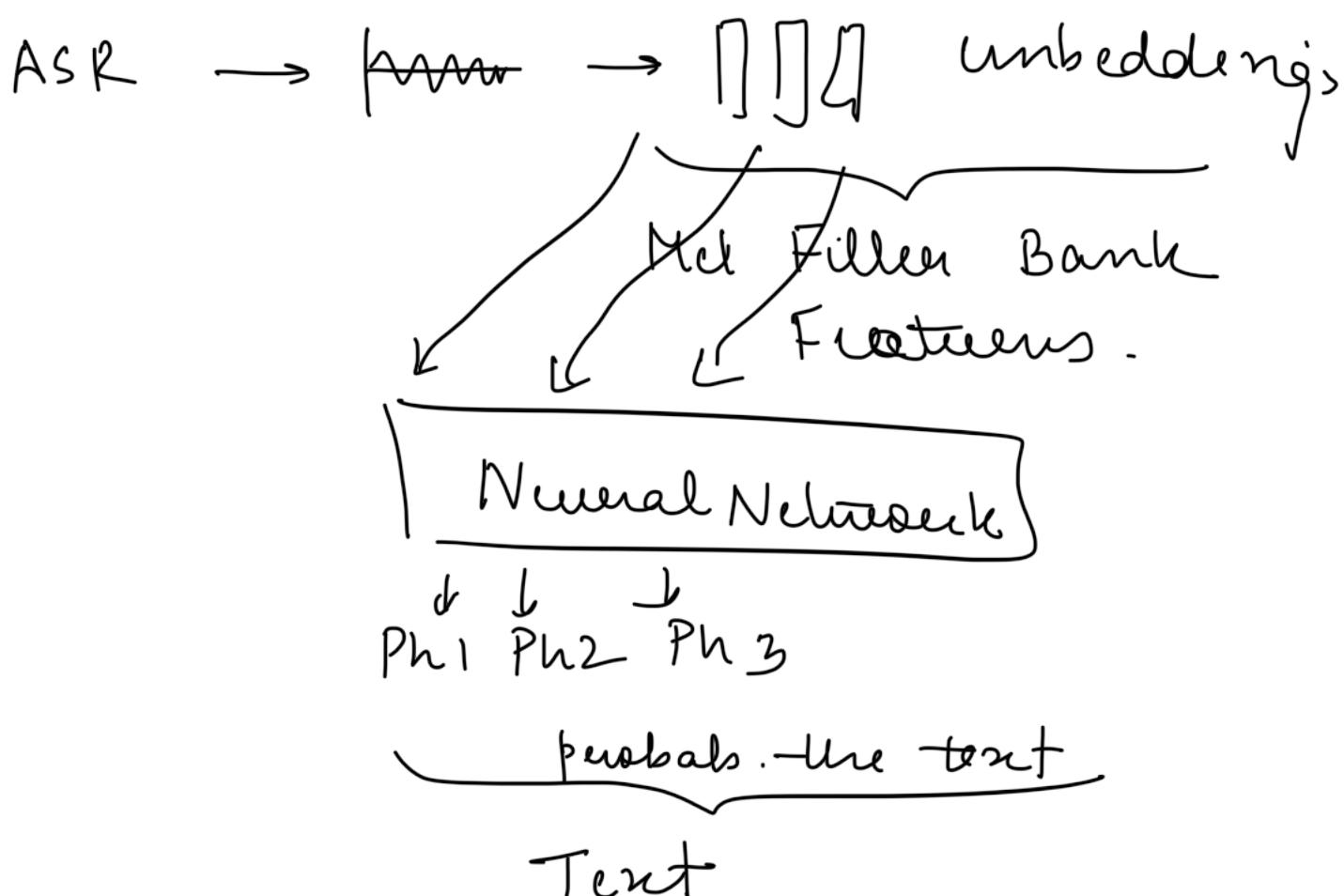
dubbing applications.

L5 Phoneme To Text

Lexicon / dictionary



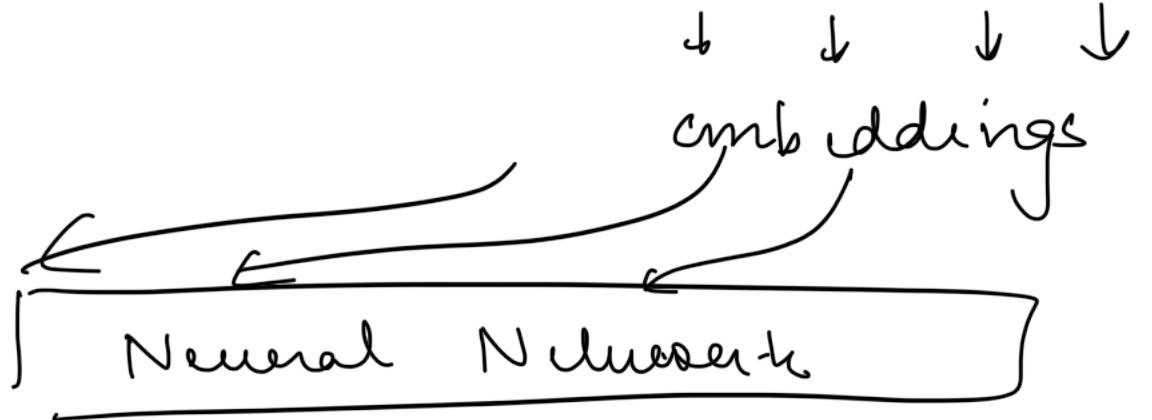
Text → sig. of phonemes
 lexicon



- In text to speech

How are you

lexicon → |h| |ow| |ə| |u|



Mel
Spectrogram

Vocoder

↳ this synthesizes the speech.

+ Add. speaker (speaker embedding using speaker vector representation)

- LG to 18 Demo