

Sinc Net for speaker Recognition

Jagabandhu Mishra

IIT Dharwad

outline:-

1. SPEAKER recognition ✓
2. Earlier approaches ✓
3. End-to-End approaches
4. Proposed approach
5. Experimental setup
6. Result & discussion

1. SPEAKER RECOGNITION

• SPEAKER RECOGNITION

Speaker Identification (SI)

Speaker Verification (SV)

THIS APPROACH

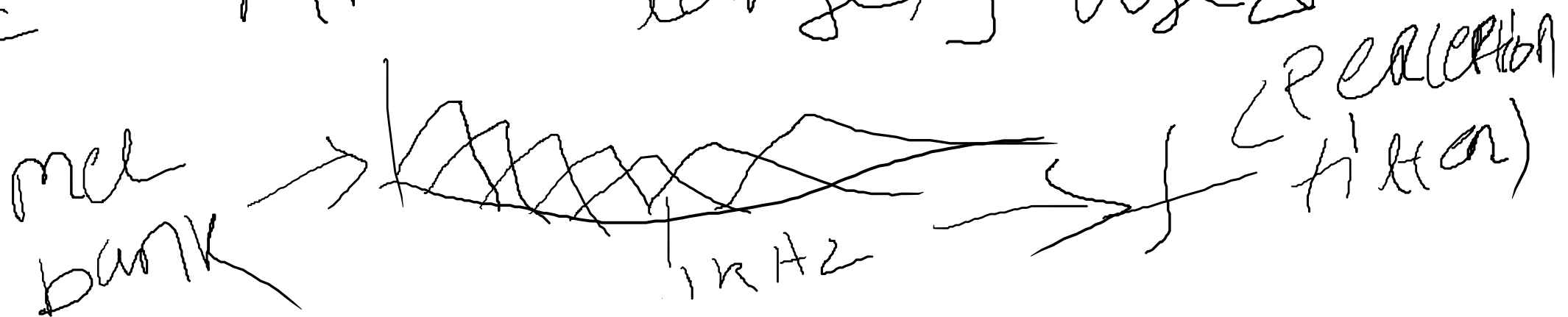
• SI :- Identify from N speakers

• SV :- Verify if genuine speaker or not (true/false)

2. Earlier approaches

• speech signal \rightarrow Feature extraction (FE) \rightarrow Modeling
 \downarrow
decision

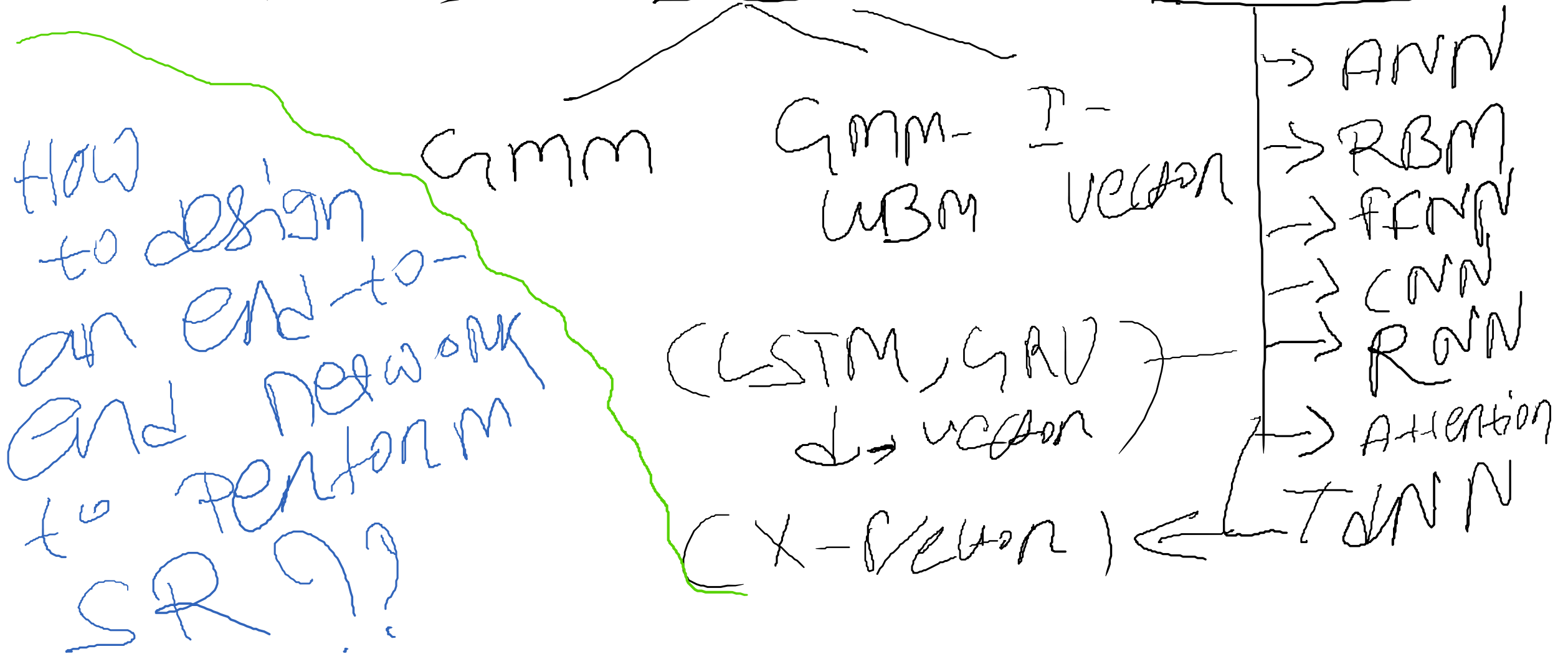
• FE = MFCC largely used



\Rightarrow Why it is good for speaker Recognition ??

2. Further Approaches (continue)

• modeling: Statistical Neural Networks



2: End-to-end approaches

- ⊗ Filterbank \rightarrow DNN
- ⊗ Spectrogram \rightarrow CNN - (H-U)
- ⊗ RAW Speech \rightarrow CNN - (don't know what layers)
- ⊗ will all this captures appropriate info ?? (for SR). (ANS:- we don't know)
- ⊗ Intuition:- Can we put some constraint on the CNN filter.??

4. PROPOSED approach:

- ⊗ Objective! - capture task specific info. from data.
- ⊗ See fig-1 of Ref. Paper.
- ⊗ 1st layer! - used as a data driven F.E.
↓ Putting constraints on
learn filters
only learn higher
and lower (w/af frequency)
- ⊗ See Eq. 1 to 4 (Ref Paper)

continue!

• eqn 1 :- $G(f, f_2) = \text{rect}\left(\frac{f}{2f_L}\right) - \text{rect}\left(\frac{f}{2f_1}\right)$

⊗ will this be feasible to set (finite sensor)

⊗ $f_2 > f_1$ (how)

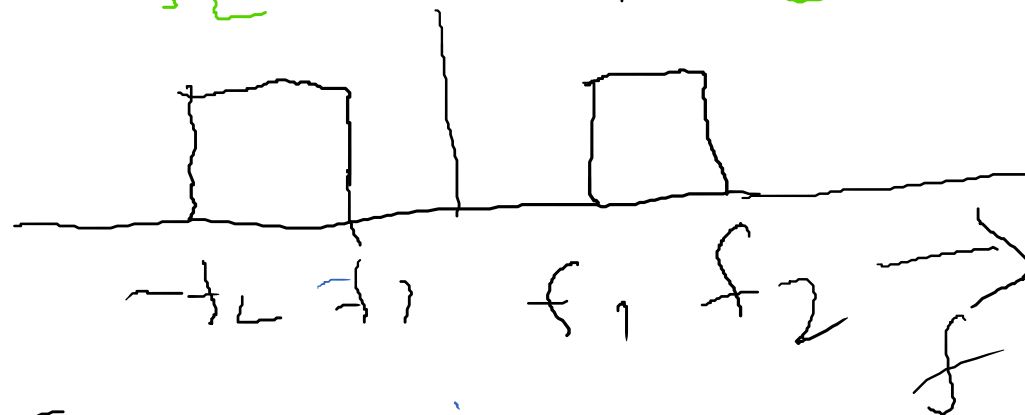
⊗ Initialization?

IFT

=



=



✓

$$g(f, f_2) = 2f_2 \sin((2\pi + \pi)) - 2f_1 \sin((2\pi + \pi))$$

~~$\sin \infty$~~

Continue to -

① $f_1 > f_2 \rightarrow e^h$ 526

② finite length: e^h 728

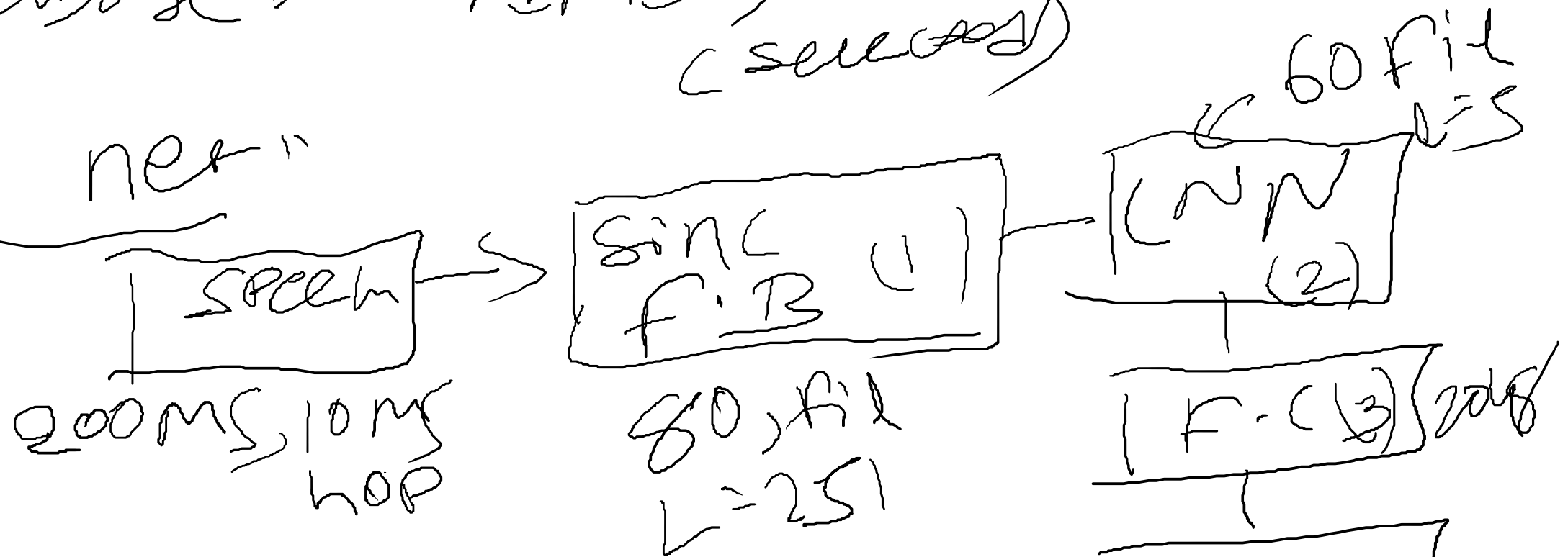
③ Initializer: —
Random ($0 - f_s/2$)
mel filter

④ Expected outcome: \rightarrow Interpretability
Fast convergence Few Parameters
(FL vs 2F)

5. Experimental Setup:

① Database - TIMIT, Librispeech
(several)

② Sinc net



③ Initial - Glorot

④ Optimizer - RMS PROP, $\text{LR} = 0.001$
 $\alpha = 0.95$, $\epsilon = 10^{-7}$, 128 mini batch size

6. Result & Discussion 1

① see fig 2, 3 & 4 (filter analysis)

② S.I & S.V. Result & Discussion

Section - 5

③ Future direction :-

- can take other filter (input of rect)
- improve model's after F.E
- other application (speech based)

Thank you

② Retention! -

Speaker Recognition from raw waveform
with SINCNET

by
M. Ravanelli, Joshua BenSD

Pioneer of
DL, along with
Prof. Hinton