### A HYBRID APPROACH TO SEGMENTATION OF SPEECH USING GROUP DELAY PROCESSING AND HMM BASED S. Aswin Shanmugam, Hema A. Murthy

- forced alignment.





## Importance of Phase $0.5 \pi$ $0.5 \pi$ $0.5 \pi$ Angular Frequency -> Angular Frequency Angular Frequency -

Figure 2: High resolution property of minimum-phase group delay [1]

The authors would like to thank the Department of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" The authors would like to thank the Department of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under the project "CSE/1112/129/DITX/HEMA" of Communications & IT, Government of India, for funding this research under DON Lab, Department of Computer Science and Engineering - IIT Madras - India

Department of Computer Science and Engineering, Indian Institute of Technology Madras, Chennai, Tamil Nadu, India



Figure 3: Syllable boundaries given by HMM based segmentation and GD based segmentation with WSF=10 & WSF=30

- With WSF "10", the number of spurious boundaries is large (resolution is high) but correct boundaries are not misplaced
- or a semi-vowel, boundary corrections are performed

### Hybrid Segmentation Algorithm

Build HSMM monophone mode
HSMM-based forc
Build HMM monophone models a
Syllable boundary correction us
Splice waveforms at the syllable I training on the syllable waveforms
HMM-based force
Syllable boundary correction us to obtain the final syllable
Splice waveforms at the syllable I training on syllable waveforms to
HMM-based forced alignmen to obtain phoneme bound
Combine phoneme level alig constituting an utterance to segmentation for the
Figure 4: Steps involved in the

► If the syllable does not end with a nasal, fricative or if it's not followed by a nasal, fricative, affricate

dels using flat start training	
ced alignment	
and perform forced alignment	
using group delay algorithm	
level and perform embedded s to refine monophone models	
ced alignment	
using group delay algorithm le level segmentation	
level and perform embedded to refine monophone models	
ent on syllable waveforms daries within syllables	
ignment within syllables o obtain phoneme level e entire utterance	
e proposed hybrid method	

# Segmentation Accuracy Increase in acoustic likelihood during forced alignment. Method Hybric HMM 0.2 Time (s) Quality of TTS HMM based speech synthesis systems (HTS) [4] built for Tamil. Both Phone HTS and Syllable HTS [5] show improved WER Order independent pair comparison tests also show improvement. Acknowledgement "CSE/1112/129/DITX/HEMA". Conclusion accurate on the particular. processing and machine learning to obtain accurate segmentation. References

- Sadhana, vol. 36, no. 5, pp. 745-782, Nov 2011."

- to speech system," in EUSIPCO, 2013.

Average log probability per frame							
Nasals	Fricatives	Semi-	Vowels	Stop	Overall		
		Vowels		Consonants			
-72.22	-79.23	-73.36	-70.77	-82.15	-73.17		
-73.74	-81.72	-75.46	-70.99	-85.04	-74.16		
Table 1: Average log probability per frame							

The syllable "ta" marked in the spectrogram is identified correctly only by the proposed method.



Figure 5: Syllable level segmentation given by the proposed hybrid method compared to HMM, HSMM and HSMM followed by HMM

System	WER
HTS - Syllable (HMM Segmentation)	11.11%
HTS - Syllable (Hybrid Segmentation)	7.07%
<b>HTS - Phone (HMM Segmentation)</b>	4.04%
<b>HTS - Phone (Hybrid Segmentation)</b>	1.01%
Table 2. Word error rate (WFR)	1

HTS - Syllable			HTS - Phone		
A-B	<b>B-A</b>	A-B+B-A	A-B	<b>B-A</b>	A-B+B-A
75	20	77.5	70	15	77.5
Т	able 3	B: Pair con	npari	ison	tests

The authors would like to thank the Department of Electronics & Information Technology (DeitY), Ministry of Communications & IT, Government of India, for funding this research under the project

Machine learning techniques are robust on the average, while signal processing techniques are

An attempt has been made to synergize the benefits of knowledge-based domain specific signal

[1] Hema A. Murthy and B. Yegnanarayana, "Group delay functions and its application to speech processing,"

[2] S. Young and P. Woodland, "HTK: Speech recognition toolkit," http://htk.eng.cam.ac.uk/

[3] T. Nagarajan, V. Kamakshi Prasad, and Hema Murthy, "Minimum phase signal derived from root cepstrum," *Electronics Letters*, vol. 39, no. 12, pp. 941-942, Jun 2003.

[4] "HMM-based speech synthesis system (HTS)," http://hts.sp.nitech.ac.jp/

[5] A. Pradhan, S. Aswin Shanmugam, A. Prakash, V. Kamakoti, and Hema A. Murthy, "A syllable based statistical text

WWW:http://www.cse.iitm.ac.in, http://lantana.tenet.res.in