



# Special Module on Media Processing and Communication

**Dayalbagh Educational Institute  
(DEI)  
Dayalbagh Agra**

**PHM 961**

**Indian Institute of Technology Delhi  
(IITD)  
New Delhi**

**SIV 864**



# Text-to-Audiovisual Speech Synthesizer





# Background

- Two approaches

Model based

- Flexible
- Lacks video realism

Image based



# Background

## Model based



3D Polygonal  
Surface

Muscle Simulation  
RFFD



Minimum Perceptible  
Actions (MPA)



Expressions  
Phonemes



Emotions and  
Sentences

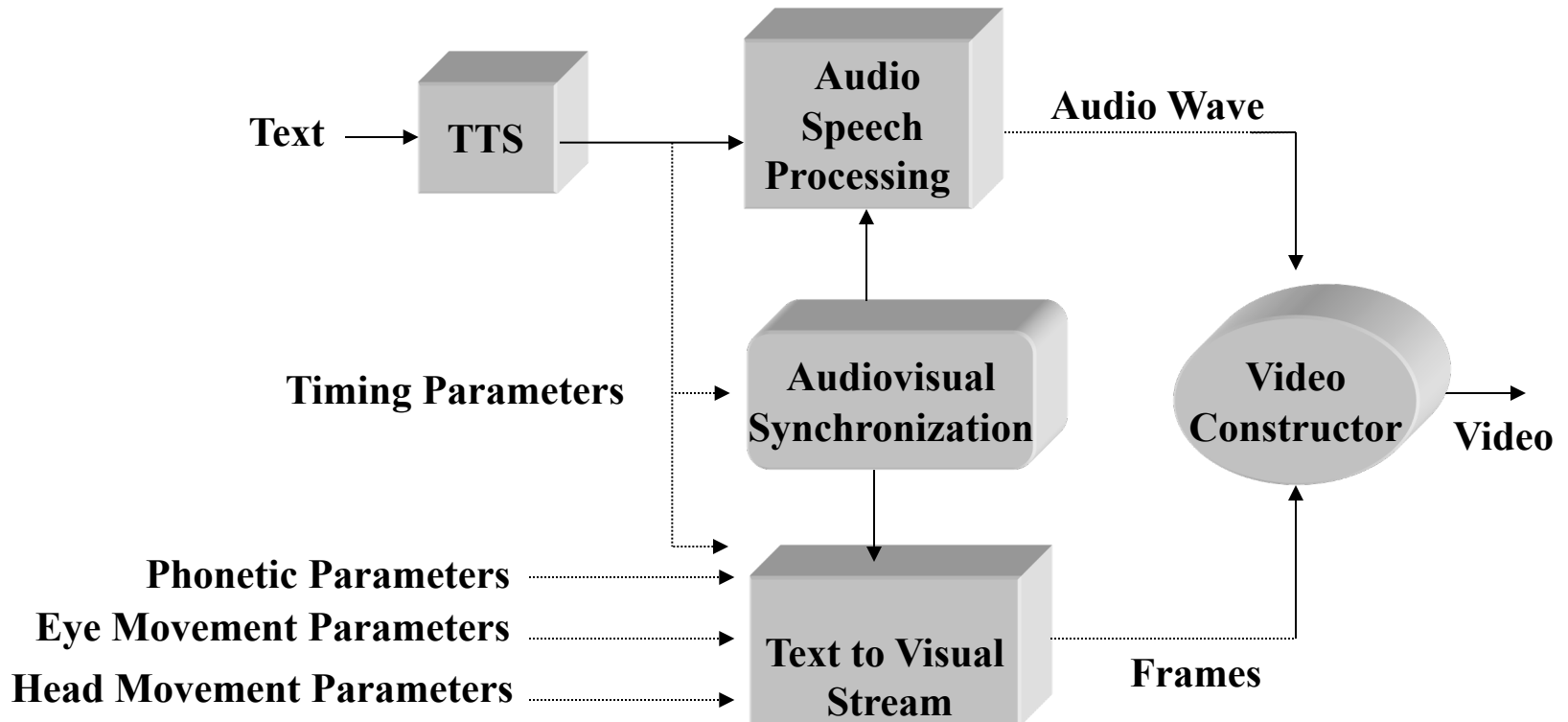


# Approach

- Basic idea
  - Morphing visemes (Ezzat and Poggio)
- Add ons
  - Eye movements
  - Head movements
  - Co-articulation*



# Overview





# Text-to-Audiovisual Stream

- Four sub-tasks
  - Viseme Extraction
  - Morphing
  - Morph Concatenation
  - Synchronization





# Text-to-Audiovisual Stream

- Viseme Extraction

Phoneme counterparts

Many to one mapping

- 16 snapshots

Set of keywords covering all the phonemes



/p, b, m/



/oo/

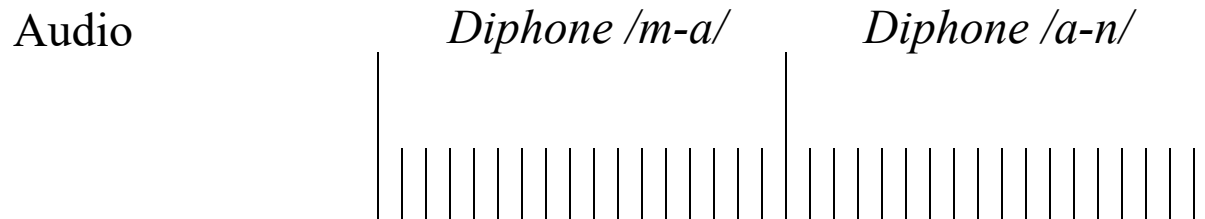


# Text-to-Audiovisual Stream

- Morphing
  - Optical flow (automatic correspondence)
  - No other moving part/object (assumption)
- Morph Concatenation
  - Simple concatenation of viseme morphs

# Text-to-Audiovisual Stream

- Audiovisual Synchronization



Viseme Transition



$$\text{Morph parameter} = \frac{s(F_k) - s(V_j)}{l(V_j)}$$



# Text-to-Audiovisual Stream

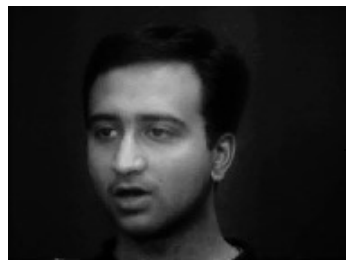
- Examples



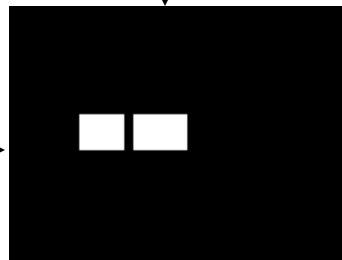
# Eye Movement

- Mask based approach

Basic image with closed eyes



Viseme /aa/



Mask

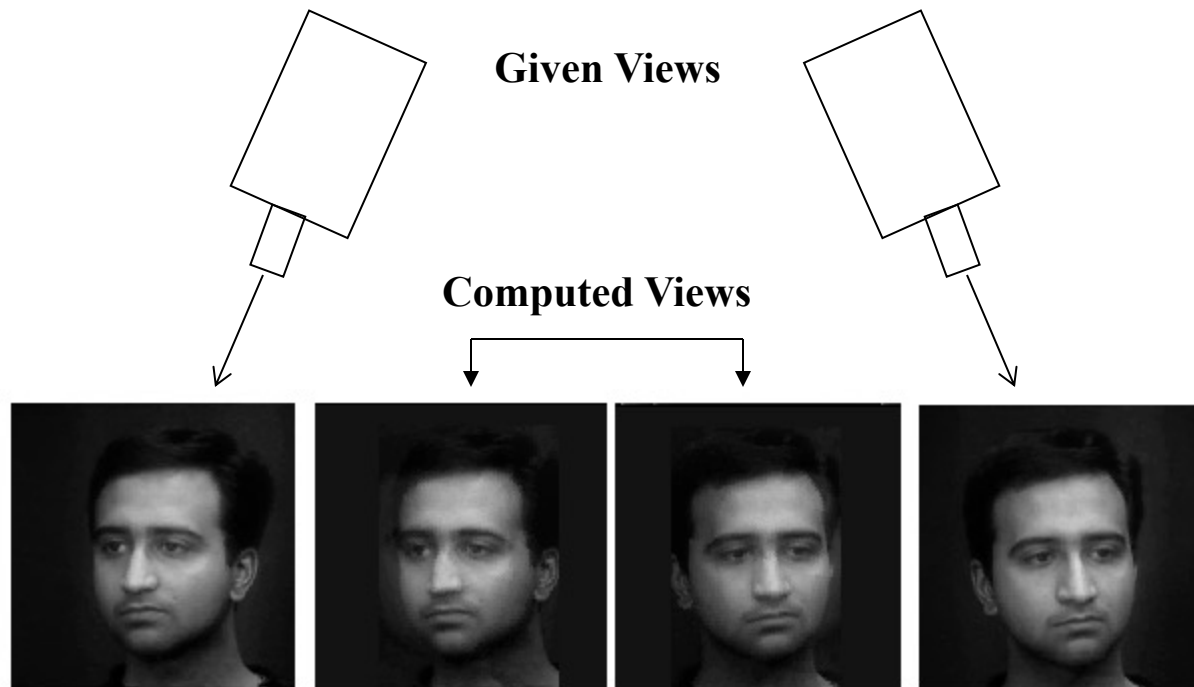


Viseme /aa/ with closed eyes



# Head Movement

- View Morphing (Seitz and Dyer 1996)



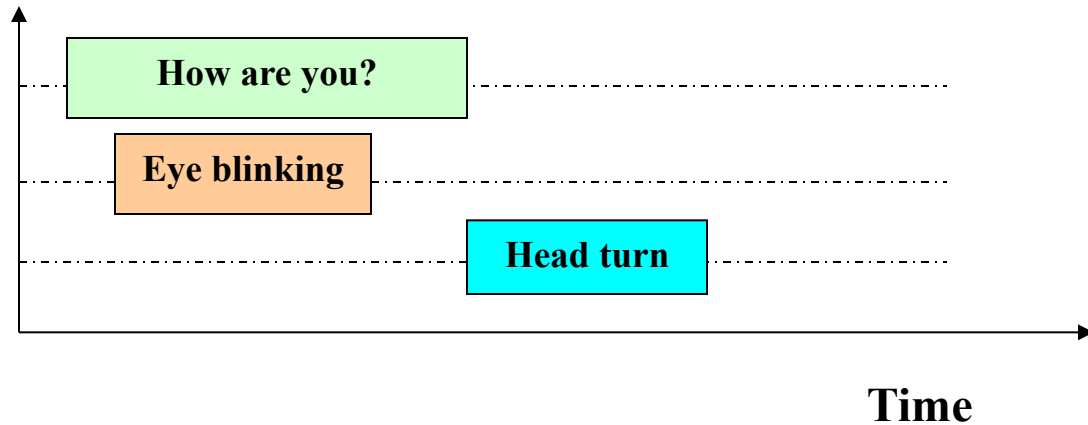


# Integration

Channel 1 (text)

Channel 2 (expression)

Channel 3  
(head movement)





# Results

**Demo 1**



**“I miss you”**

**Demo 2**



**“I am fine, thank you”**





# Co-articulation

- Problem

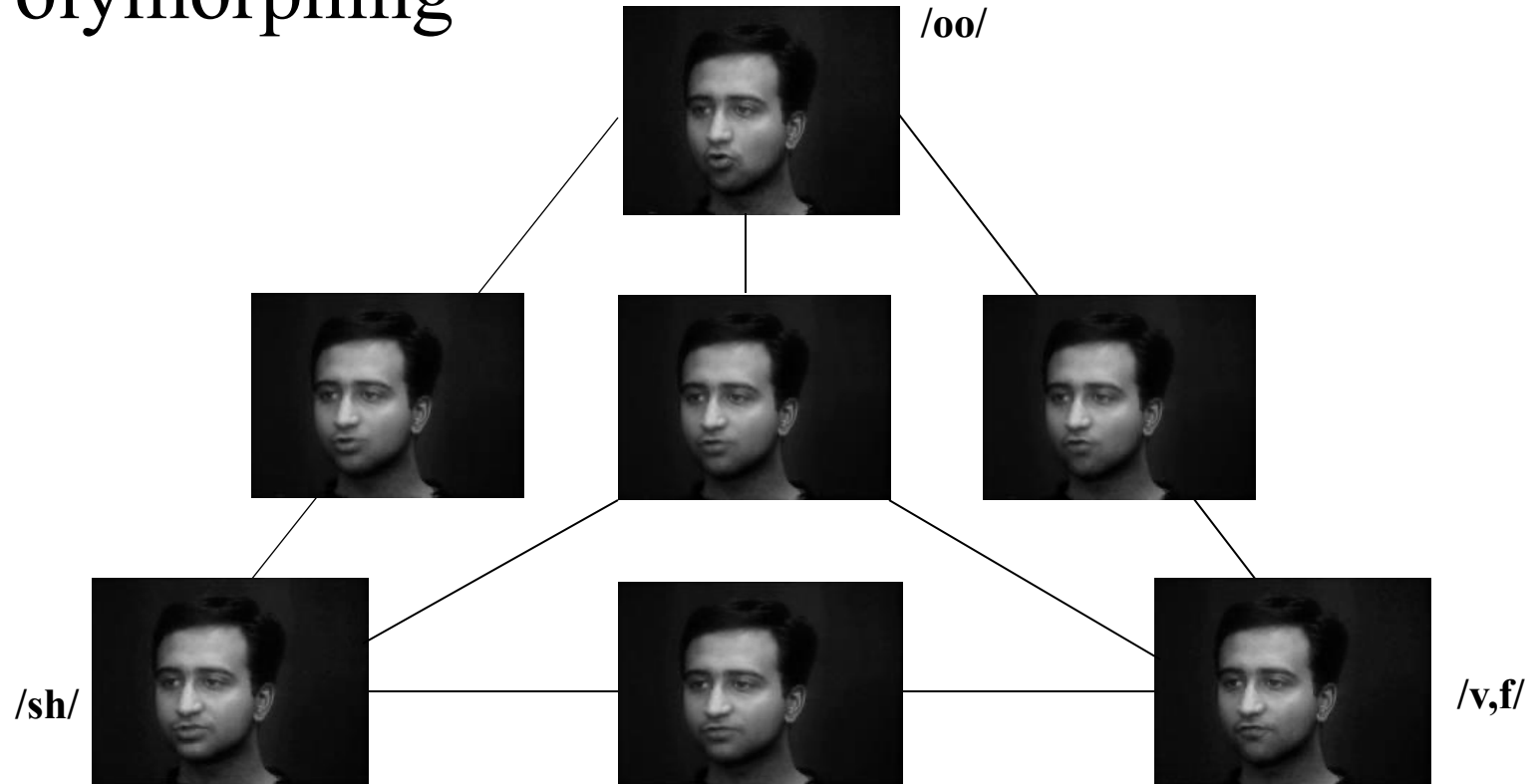
*There is an overlap in the production of syllables and phonemes*

- Approach

*Polymorphing*

# Co-articulation

## Polymorphing





# Co-articulation

- Some Results



**Without co-articulation**



**With co-articulation**

**“Tea twenty two temporary food stew”**



# Conclusion

- A text-to-audiovisual speech synthesizer
- Non verbal communication (expressions and head movement)
- Co-articulation
- Learning based methods