

## Implementing SpokenMedia for the Indian Institute for Human Settlements

Brandon Muramatsu

Andrew McKinney

Peter Wilkins

July 2010

Citation: Muramatsu, B., McKinney, A., Wilkins, P. (2010). Implementing SpokenMedia for the Indian Institute for Human Settlements. Presented at Technology for Education 2010: Mumbai, India, July 1, 2010.

Citation: Muramatsu, B., McKinney, A., Wilkins, P. (2010). Implementing SpokenMedia for the Indian Institute for Human Settlements. Presented at Technology for Education 2010: Mumbai, India, July 1, 2010.

Unless otherwise specified, this work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States License

## Case Study of Using SpokenMedia for IIHS

### ■ Goals

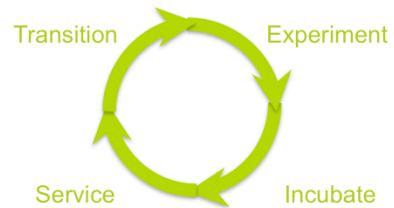
- Demonstrate value of transcripts + translations for IIHS
- Test a process to produce transcripts and translations

### ■ Describe the process and our experiences

- Transcribe -> Edit -> Translate -> Present

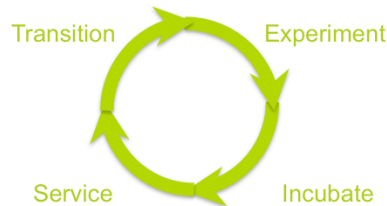
## MIT Office for Educational Innovation and Technology

- Dean for Undergraduate Education
- Support Innovation Cycle
- Novel uses of technology to support teaching and learning
- Academic computing without the LMS/VLE



# SpokenMedia Project

SLS  SPOKEN LANGUAGE SYSTEMS  
MIT Computer Science and Artificial Intelligence Laboratory



*Service:* OEIT Considering  
•What would it do?  
•Is it valuable?

*Experiment:* Spoken Lecture Project  
• Speech recognition research  
• Automatically transcribe video lectures  
• Based on 20+ years of research  
• Built prototype

 SpokenMedia

*Incubate:* SpokenMedia Project  
• Technology transfer from research lab  
• Automatic lecture transcription  
• Search  
• Player  
• Transcript Editor

## Spoken Lecture Project

- Supported by iCampus
- Includes the browser (which was just demo'd) the processor (back end lecture transcription) and a hand workflow to do the processing
- Approximately 400 hours of video indexed

## SpokenMedia

- Technology transfer—get code running outside of original environment
- 4 components: automatic lecture transcription, search, player, transcript editor

- Supported by iCampus

## SpokenMedia as a Service?

- Still investigating

*The Indian Institute for Human Settlements (IIHS) will... “create India’s first independent National Innovation University focused on the challenges and opportunities of its urbanisation.”*

*– Indian Institute for Human Settlements:  
Curriculum Framework Version 3.0  
January 2010*

*"The IIHS Website is our  
commitment to a different  
way of looking at things."*

*– Aromar Revi  
5 January 2010*

*“The Institution will fail or scale  
based on language.”*

– Aromar Revi  
5 January 2010

## What did we do?



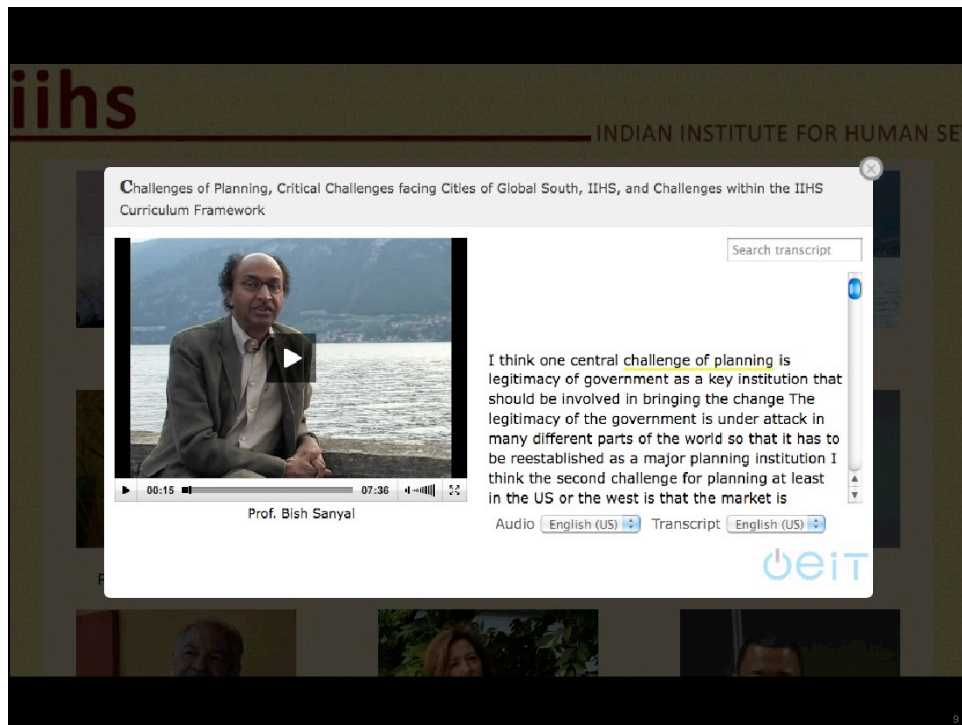
Four step process (simple)

First we used the SpokenMedia to do automatic transcription.

Next we did hand edit and translation steps.

Then we created a player for the presentation of the video and transcripts...



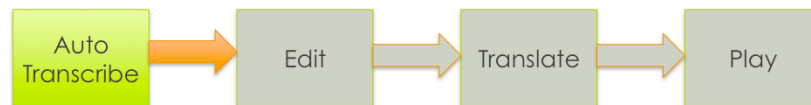


## Demo

In this demo, we see a video interview with Prof. Bish Sanyal from MIT. We'll see three things in this demo:

- As Prof. Sanyal speaks, we will see the text in the transcript highlighted and the highlighting will follow along ("bouncing ball")
- Switching the transcript from English to a hand translation into Hindi that is synchronized with the audio, as the switch occurs the playback is seamless
- Searching within the transcript, after entering the search term and selecting the result, the video and transcript seek to that point in the video and playback continues

## How did we do it?



First we used the SpokenMedia to do automatic transcription.

## How do we do it? Spoken Lecture Research

James Glass  
glass@mit.edu



- Speech recognition & automated transcription of lectures
  - Conversational, spontaneous, starts/stops
  - Different from broadcast news, other types of speech recognition
  - Specialized vocabularies
- Developed recognizer, browser

Supported with iCampus MIT/Microsoft Alliance funding

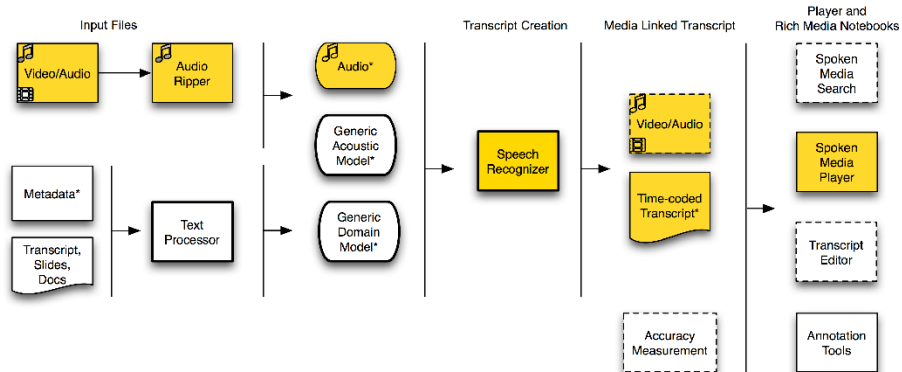
11

Unless otherwise specified this work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0 United States License (<https://creativecommons.org/licenses/by-nc-sa/3.0/us/>)

### Lecture Transcription

- Jim Glass and his group have years of research experience for spoken languages
- Lectures are a different type of spoken language
  - Much of the speech recognition research has focused on real time transcription of news broadcasts, or interactive voice response systems (telephone)
  - Broadcast news has something like 300 unique words in an hour long broadcast
  - Broadcast news is well structured, prepared copy (in studio via teleprompters), clear transitions between speakers, etc.
  - Lectures are conversational and spontaneous
  - Can use highly specialized vocabularies, engineering, physical sciences, mathematics

## SpokenMedia Process



We used a portion of the SpokenMedia process for the demo

12

We only used part of the process due to time constraints.

Audio separation, speech processing, time-coded transcript, and then presentation through a SpokenMedia player.

## How did we do it?



Next we did hand edit and translation steps.

## Edit & Translate: Accuracy

Automatic Transcription	Hand Transcription	Time Adjusted	Translated Hindi
I	I	I	मेरे खयाल से
think	think	think	
once	one	one	नयोजन की एक मुख्य चुनौती है
and	central		
so	challenge	central	
the	of		
challenger	planning	challenge of	
planning	is	planning	
nice	legitimacy	is	
legitimacy	of	legitimacy of	
of	government	government	सरकार की एक ऐसी मुख्य संस्थान के रूप में वैधता
government	as	as	

For this demo, we did computer-based automatic transcription, sent a file to IIHS for “editing” that consisted of performing a hand transcription (due to the format we sent, and the low accuracy of the automatic transcription in this case), a time alignment (though I actually feel that it’s “off” or “slow”) and then a hand translation by IIHS.

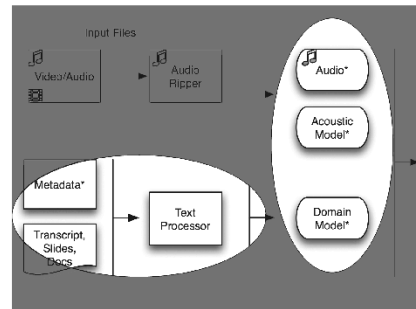
Automatic transcription is in the ~50-55% accuracy range (by way of comparison YouTube Auto Caption for this same video is ~68% accuracy).

## Automatic Speech Recognition Accuracy

### Accuracy

- Domain Model and Speaker Model
- Internal validity measure
- Seed with transcript

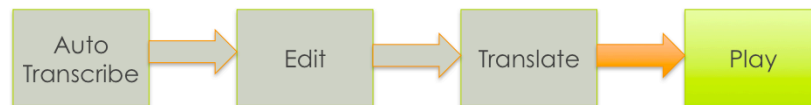
Ongoing research by Jim Glass and his team @ MIT



### Recognizer Accuracy

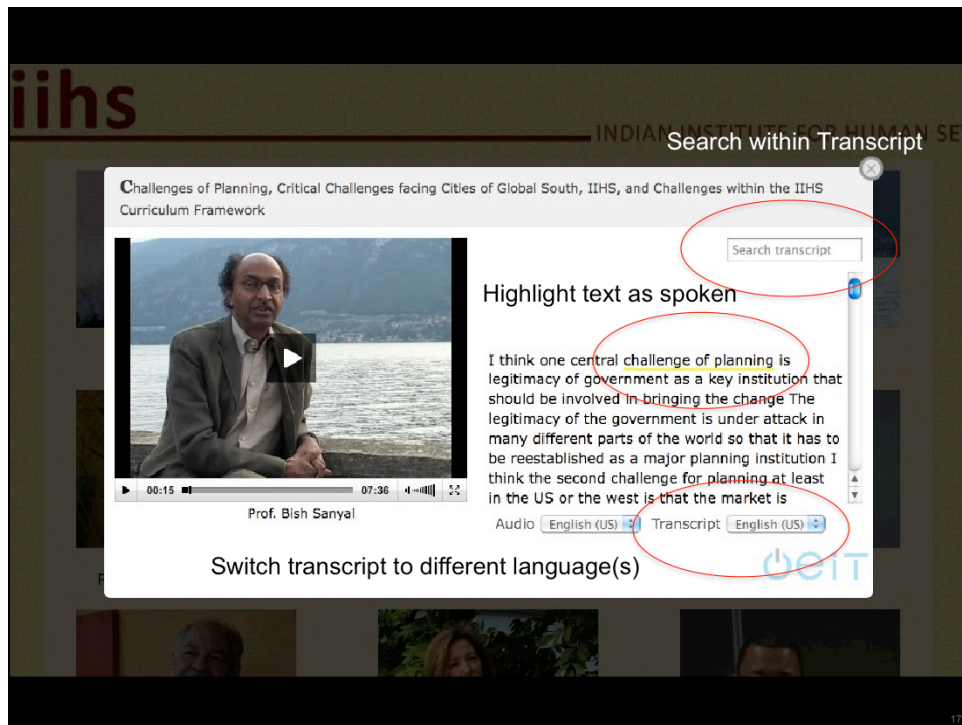
- Base accuracy is approximately 50% (generic domain and speaker models)
- Increase accuracy with speaker model up to 80-85%, and specific domain model
  - This approach is good for courses with multiple lectures by the same speaker
  - Domain models get more useful as more relevant text documents are indexed (keyword/noun phrase extraction)
- Initial results indicate that doing one 99% accurate (by hand/manual) transcript can help immensely for additional lectures by the same speaker
  - Better use of limited resources
- Search accuracy is closer to 90%, searches tend to be for unique words which the processor is better at recognizing

## How did we do it?



Then we created a player for the presentation of the video and transcripts...





### SpokenMedia Player 1.0

- Video-linked transcript
- Highlighted text follows along as the speaker speaks
- Switch transcript to a different transcript track seamlessly during playback
- Search within a transcript

## SpokenMedia (circa January)

### ■ Features

- Automated Lecture Transcription
  - Create a transcript from English lecture videos
- SpokenMedia Player
  - “Bouncing Ball” (underline text) follow along
  - Search within a video
  - Multiple transcript language support

### What we have today

- It's not perfect, but a pretty good start
- Prototype has a number of useful features that demonstrate search interfaces and interaction interfaces

## Challenges

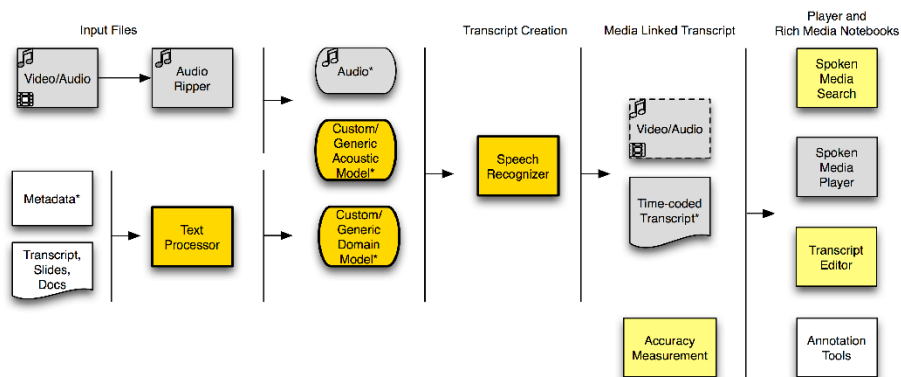
### ■ Accuracy!

- We've seen as high as 90% accuracy, but its not going to be 90% for everyone
- Goal is to do this without special training, and be useful for any lecture video
- Developing an editing tool to correct transcripts by hand

### ■ Applicability of speech recognition

- Indian context and speech require customization of underlying speaker models and speech recognizer

## SpokenMedia (July 2010)



The bright yellow indicates features working in the last two months...

- Text processing to create a custom domain model
- Creation of custom acoustic models in unsupervised mode
- Updated speech recognition software

The gray indicates features we've had working since December 2009 and that were used for IIHS

The light yellow indicates features on which we've just started working.

- Accuracy measurement
- SpokenMedia Search (search across multiple videos)
- Transcript Editor

## Where are we heading?

- **Improved accuracy**
  - Automated accuracy measurement
- **Search across multiple video transcripts**
- **New players with bookmarking, annotation, “paper-based video”**
- **Automate and improve processing**
  - > Starting a lecture transcription service

### Where are we heading?

- Transition from research project to service
- Explore new interactions—what we’re calling Rich Media Notebooks

## Check it out for yourself

- IIHS Demo:  
<http://spokenmedia.mit.edu/demo/iihs/>
- SpokenMedia Website:  
<http://spokenmedia.mit.edu/>

## Thank You!

Brandon Muramatsu, [mura@mit.edu](mailto:mura@mit.edu)

Andrew McKinney, [mckinney@mit.edu](mailto:mckinney@mit.edu)

Peter Wilkins, [pwilkins@mit.edu](mailto:pwilkins@mit.edu)

23

Citation: Muramatsu, B., McKinney, A., Wilkins, P. (2010). Implementing SpokenMedia for the Indian Institute for Human Settlements. Presented at Technology for Education 2010: Mumbai, India, July 1, 2010.

Unless otherwise specified, this work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States License