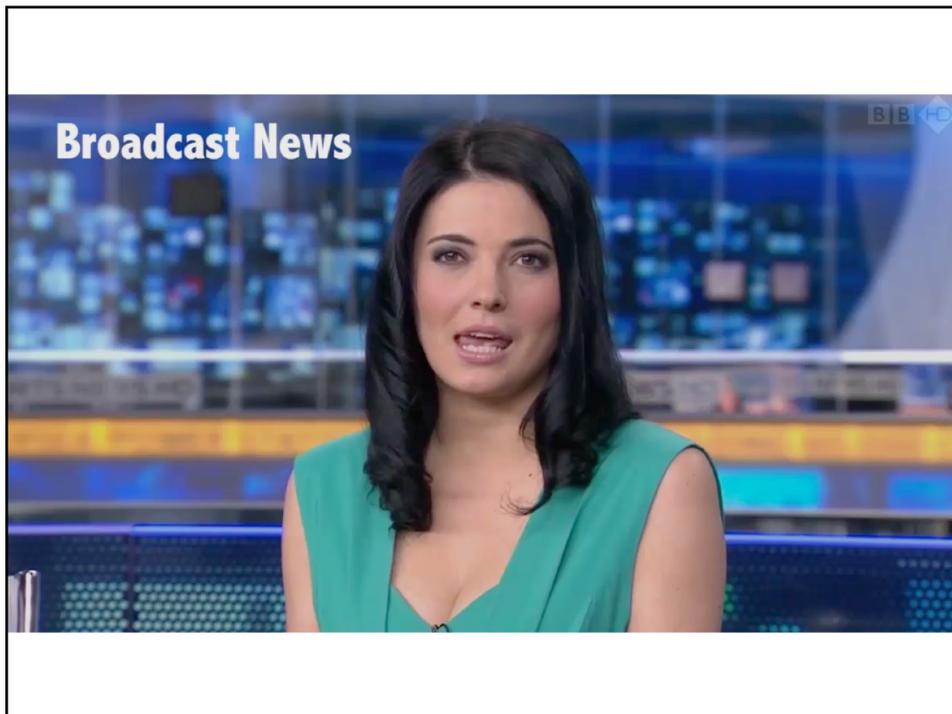
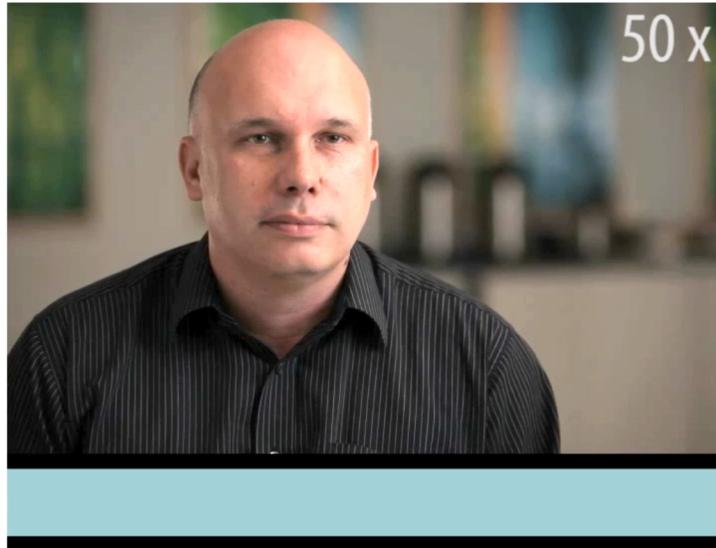


# Transcript-Based Manipulation and Browsing

**Tools for Placing Cuts and Transitions in Interview Video.** Floraine Berthouzo, Wilmot Li and  
Maneesh Agrawala, SIGGRAPH 2012.



## Raw Footage: Hours of Material



## Editing Process

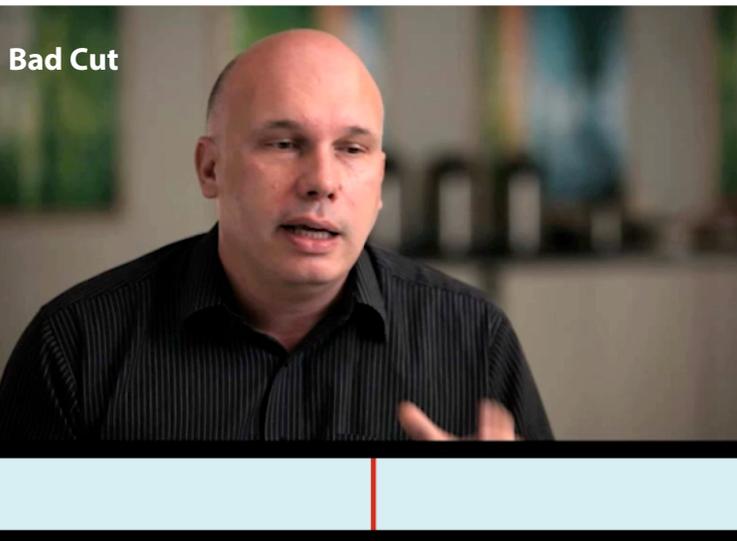
## Step 1: Index Footage



0:40 Some of our customers that have been working in print for many years ...  
1:15 So what we're trying to do ... (Noise)  
1:28 So what we're trying to do is to bring to our customers the tools that will allow them to express sophisticated layout ...

3:45 Now you'll be able to have much more creative layouts ... you can have as many regions as you want ... and the text will be able to flow across all of these regions ...  
1:02 If the orientation of your device changes from landscape to portrait ... the layout and text will reflow across all these different regions ...  
7:23 ... we have a fantastic engineering team ...

## Step 2: Find Good Cut Locations



Not gesturing  
Not in the middle of saying a word

## Step 3: Insert Transition



Jump Cut

Fade-to-Black

Zoom-In

### Visible Transition

**Jump/Fade** seen as time passing (short/long)

**Zoom** seen as continuous action from multiple cameras

## Step 3: Insert Transition



Optical Flow Interpolation

### Hidden Transition

Seen as single continuous shot

Difficult to smoothly interpolate large changes

## Current Problems

### 1. Index footage

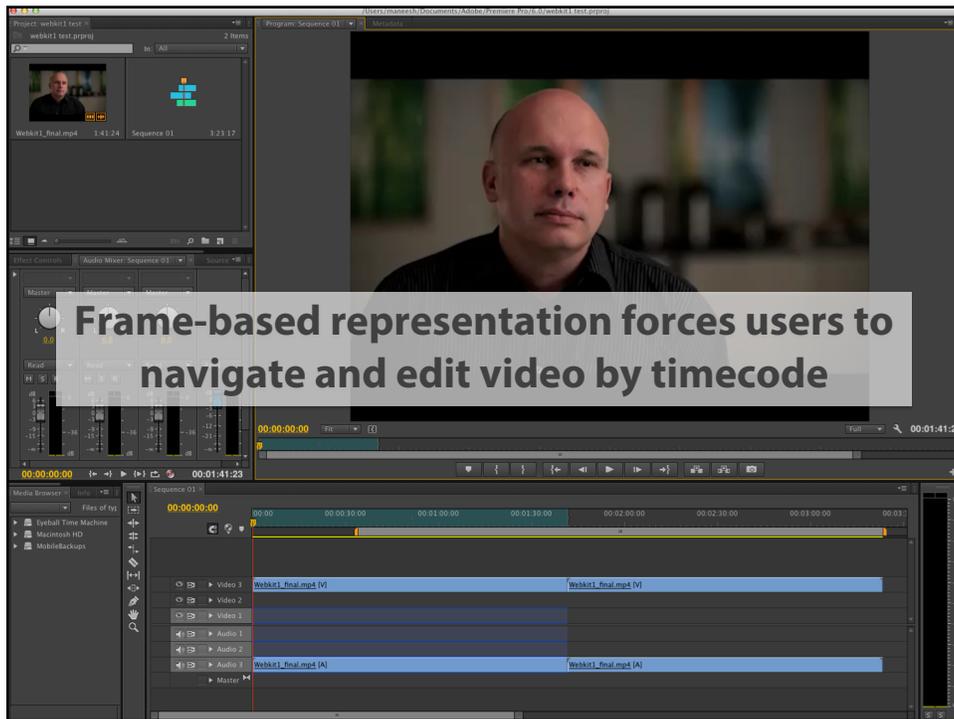
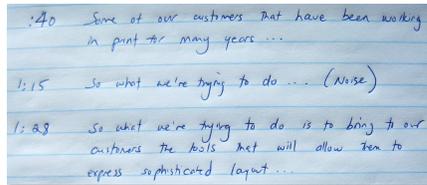
Must listen many times (very slow)

### 2. Find good cut locations

Requires additional scrubbing

### 3. Hidden transition

Interpolation introduces artifacts



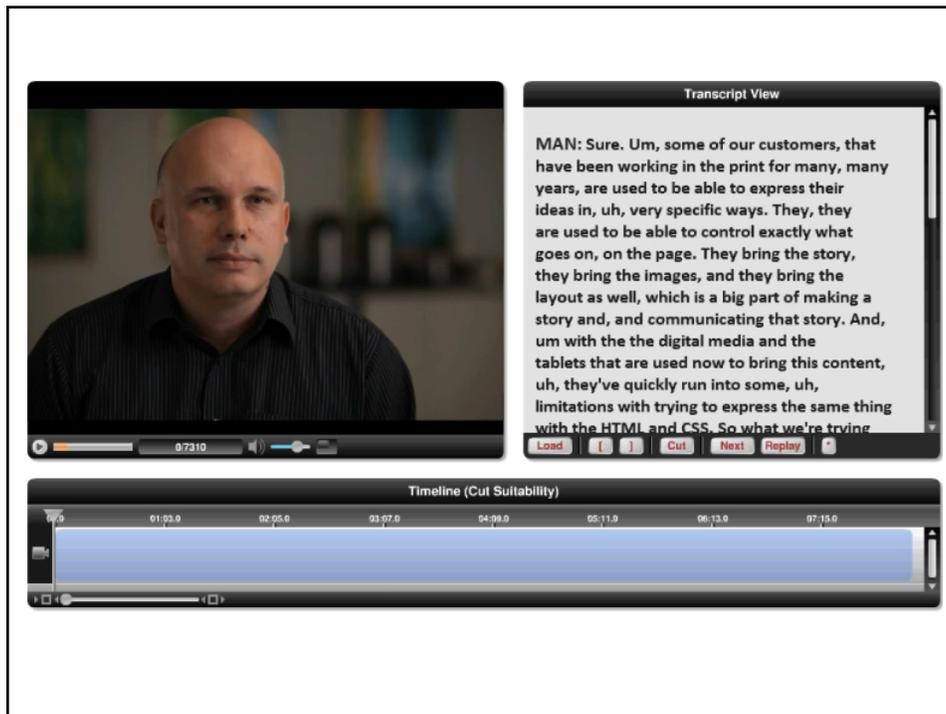
## Conceptual Model

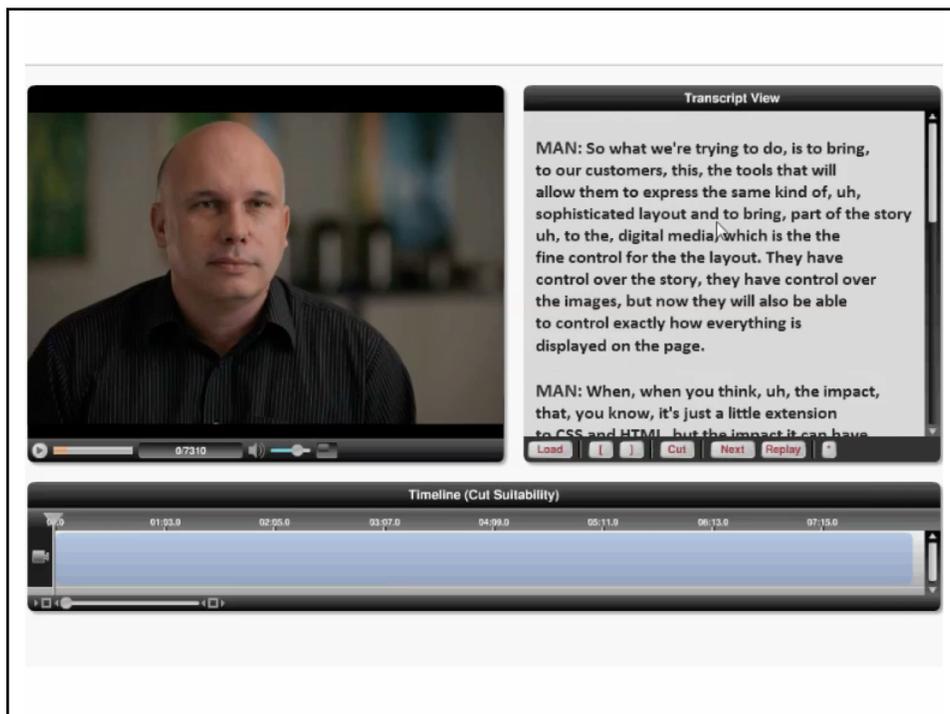
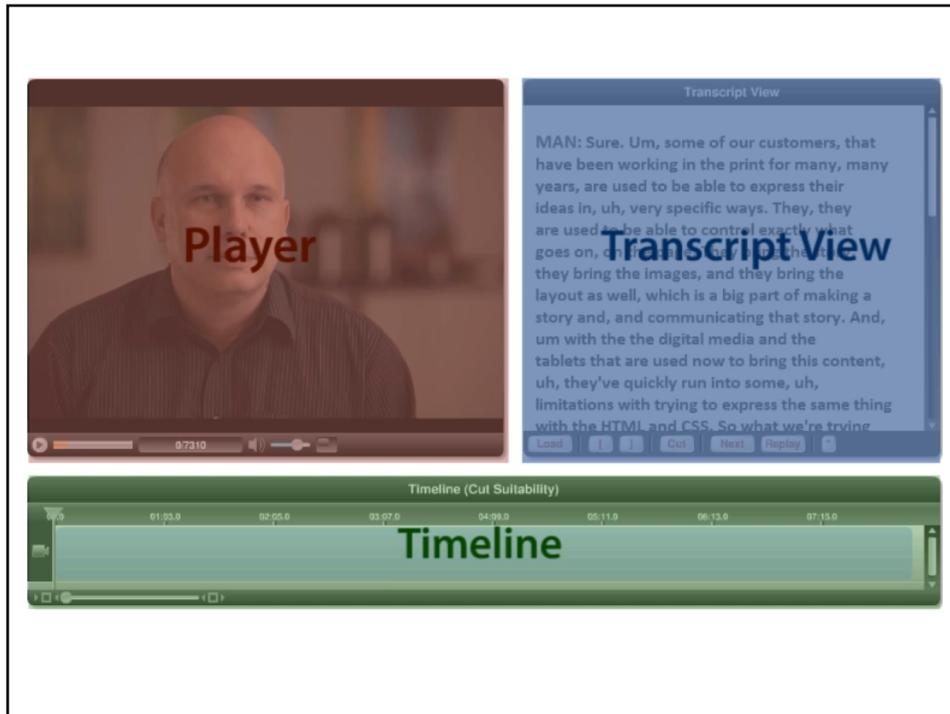
**Man:** Sure. Um, some of our customers, that have been working in the print for many, many years, are used to being able to express their ideas in, uh, very specific ways. They, they are used to being able to control exactly what goes on, on the page. They bring the story, they bring the images and they bring the layout as well, which is a big part of making a story and, and communicating that story.

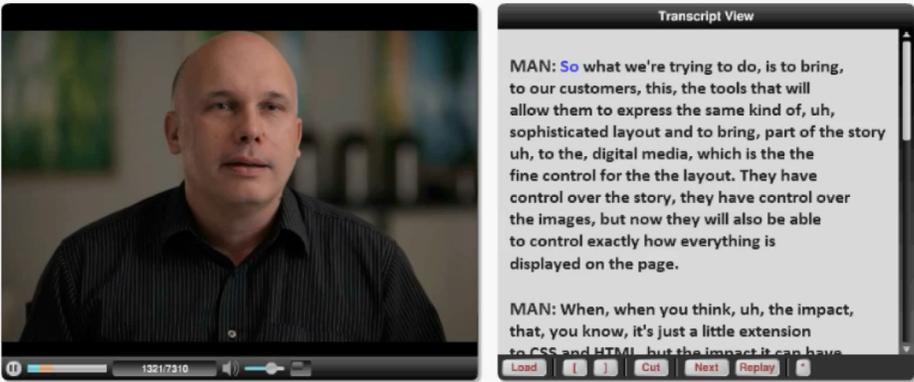
And, um, with the, the digital media and the tablets, that are used now to bring this content, uh, they've quickly run into some, uh, limitations with trying to express the same thing with the HTML and CSS.

**People think of interviews in terms of content**

**Goal: Let users navigate and edit video using this higher-level transcript-based representation**







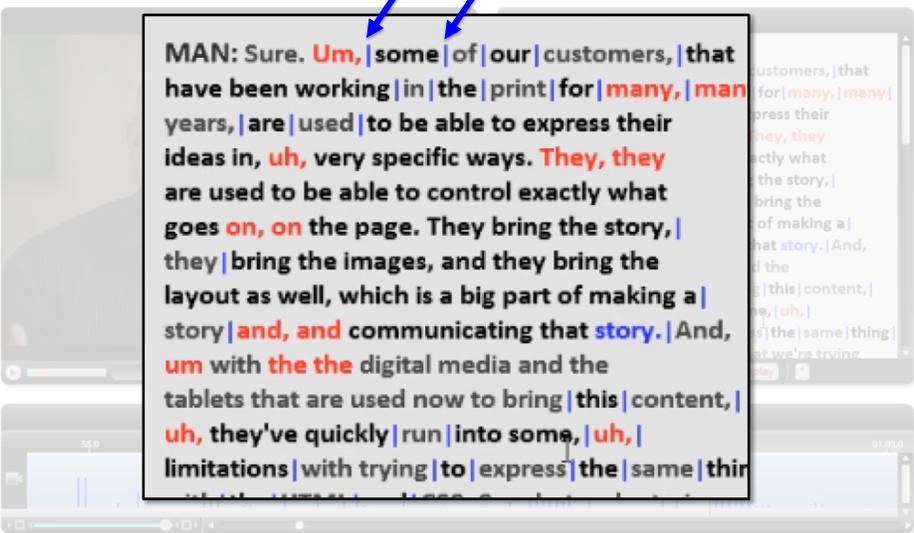
Transcript View

MAN: So what we're trying to do, is to bring, to our customers, this, the tools that will allow them to express the same kind of, uh, sophisticated layout and to bring, part of the story uh, to the, digital media, which is the the fine control for the the layout. They have control over the story, they have control over the images, but now they will also be able to control exactly how everything is displayed on the page.

MAN: When, when you think, uh, the impact, that, you know, it's just a little extension to CSS and HTML but the impact it can have

MAN: So what we're trying to do, is to bring, to our customers, this, the tools that will allow them to express the same kind of, uh,

Transcript View Zoom In



MAN: Sure. Um, some of our customers, that have been working in the print for many, many years, are used to be able to express their ideas in, uh, very specific ways. They, they are used to be able to control exactly what goes on, on the page. They bring the story, they bring the images, and they bring the layout as well, which is a big part of making a story and, and communicating that story. And, um with the the digital media and the tablets that are used now to bring this content, uh, they've quickly run into some, uh, limitations with trying to express the same thing

Transcript View Zoom In

The screenshot displays a video player interface. On the left is a video frame showing a man speaking. On the right is a 'Transcript View' panel with the following text:
 

MAN: Sure. Um, some of our customers, that have been working in the print for many, many years, are used to be able to express their ideas in, uh, very specific ways. They, they are used to be able to control exactly what goes on, on the page. They bring the story, they bring the images, and they bring the layout as well, which is a big part of making a story and, and communicating that story. And, um with the the digital media and the tablets that are used now to bring this content, uh, they've quickly run into some, uh, limitations with trying to express the same thing with the HTML and CSS. So what we're trying

 Below the transcript is a 'Timeline (Cut Suitability)' graph showing a blue bar chart representing suitability levels over time. The video progress bar at the bottom shows a timestamp of 9:18:7310.

This screenshot is similar to the one above, showing the same video player interface. The transcript text is identical, but the video progress bar now shows a timestamp of 9:18:7285. A vertical line is drawn on the transcript at the end of the sentence 'So what we're trying to do... is we're trying to...'. Below the video player, the text 'Hidden Transition' is centered.

The screenshot displays a video editing software interface. On the left is a video preview window showing a man speaking. On the right is a 'Transcript View' window with the following text:
 

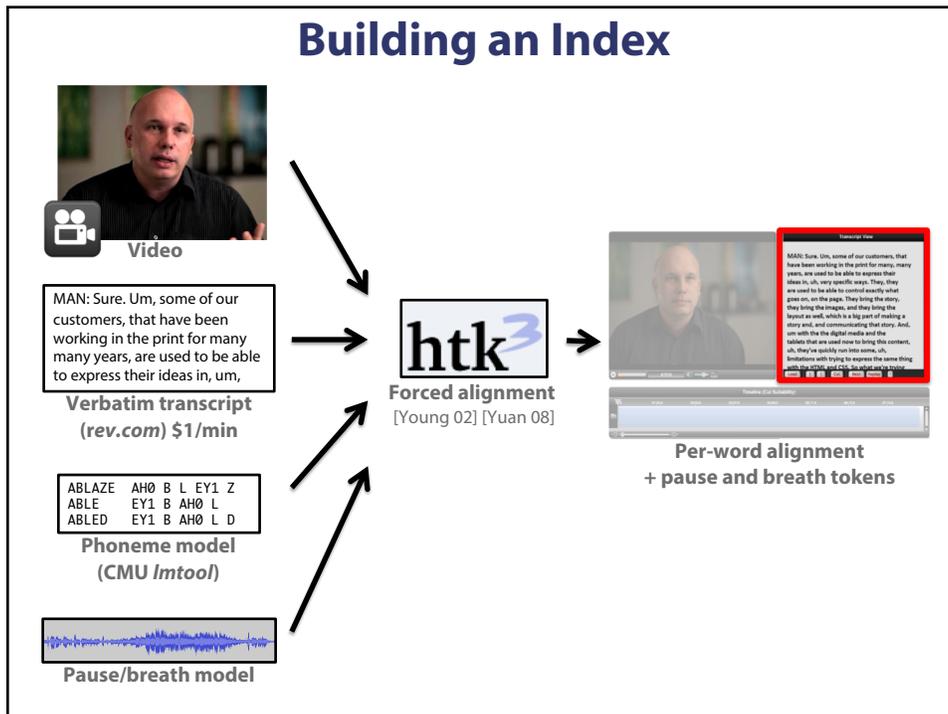
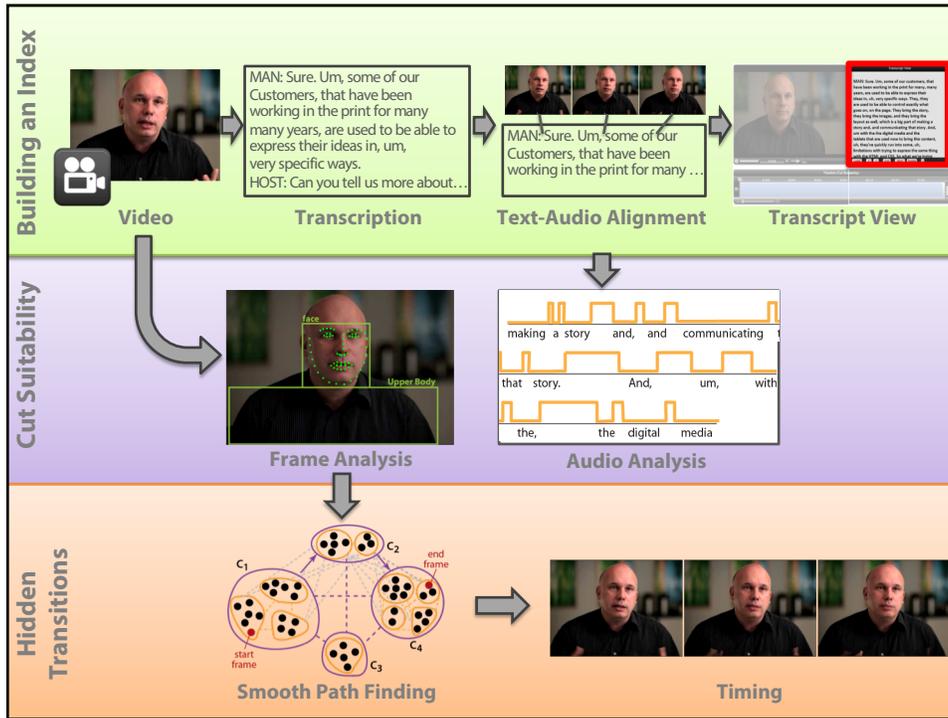
MAN: Sure. Um, some of our customers, that have been working in the print for many, many years, are used to be able to express their ideas in, uh, very specific ways. They, they are used to be able to control exactly what goes on, on the page. They bring the story, they bring the images, and they bring the layout as well, which is a big part of making a story. And, um with the the digital media and the tablets that are used now to bring this content, uh, they've quickly run into some, uh, limitations with trying to express the same thing with the HTML and CSS. So what we're trying to do, is we're trying to...

 Below the transcript is a 'Timeline (Cut Suitability)' window showing a blue bar graph representing cut suitability over time. A vertical line is positioned at approximately 01:01:00. The video player at the bottom shows a progress bar at 00:09:7283.

Jump Cut

The screenshot displays the same video editing software interface as above. The video preview window shows the man speaking. The 'Transcript View' window shows the same text as above. The 'Timeline (Cut Suitability)' window shows the same blue bar graph. A vertical line is positioned at approximately 01:01:00. The video player at the bottom shows a progress bar at 01:08:7285.

Hidden Transition



## Cut Suitability Score

MAN: Sure. **Um**, some of our customers, that have been working in the print for **many, many** years, are used to be able to express their ideas in, **uh**, very specific ways. **They, they** are used to be able to control exactly what goes **on, on** the page. They bring the story, they bring the images, and they bring the layout as well, which is a big part of making a story **and, and** communicating that **story**. **And**, **um** with **the the** digital media and the tablets that are used now to bring **this** content, **uh**, they've quickly run into some, **uh**, limitations with trying to express the same thing

## Scoring

Not gesturing

Not in the middle of saying a word



$$Score(i) = S_v(i)S_a(i)$$

Visual Score

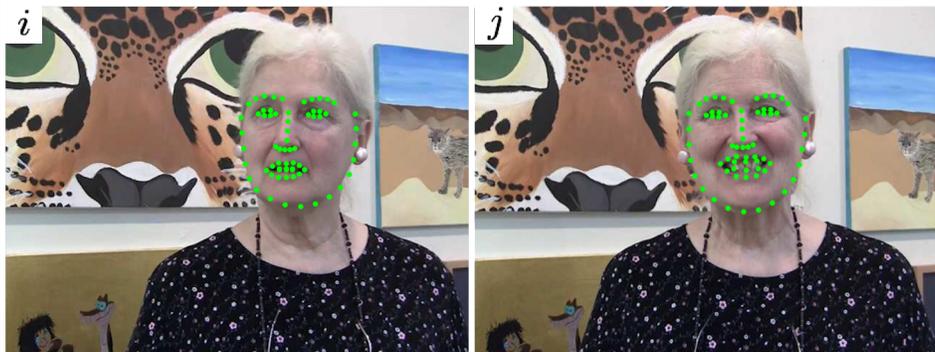
Audio Score

## Visual Frame Distance



$$Distance(i, j) = Spatial\ Distance + Appearance\ Distance$$

## Spatial Distance



Facial Feature Tracking [Saragih 09]

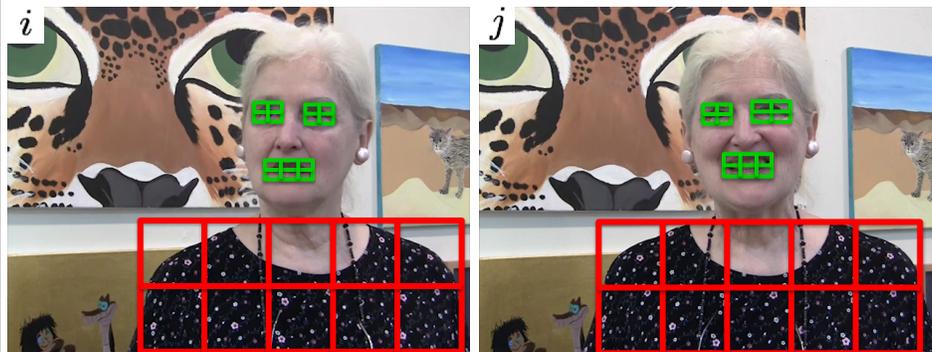
## Spatial Distance



Facial Feature Tracking [Saragih 09]

$$\text{Spatial Distance}(i, j) = \|p_i - p_j\|$$

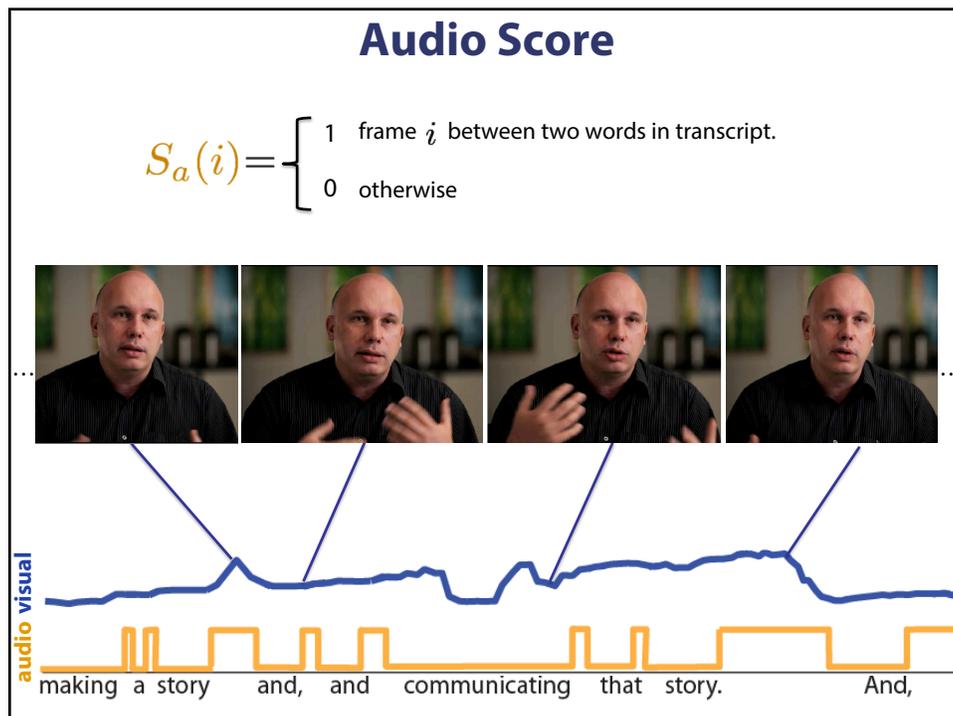
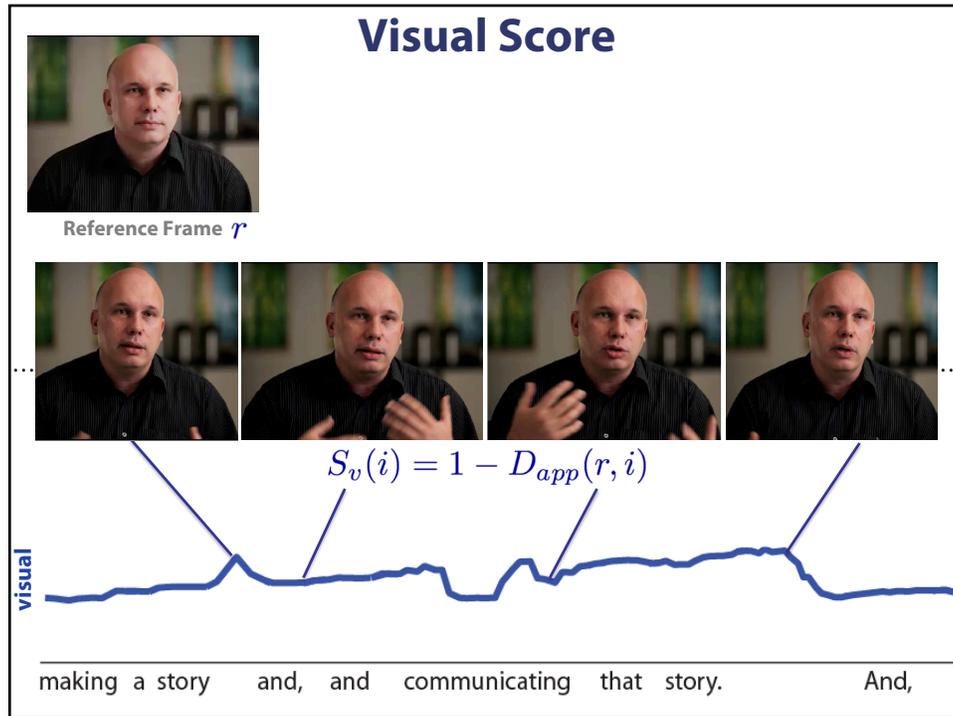
## Appearance Distance



Compute histogram of oriented gradients (HOG) feature for each region

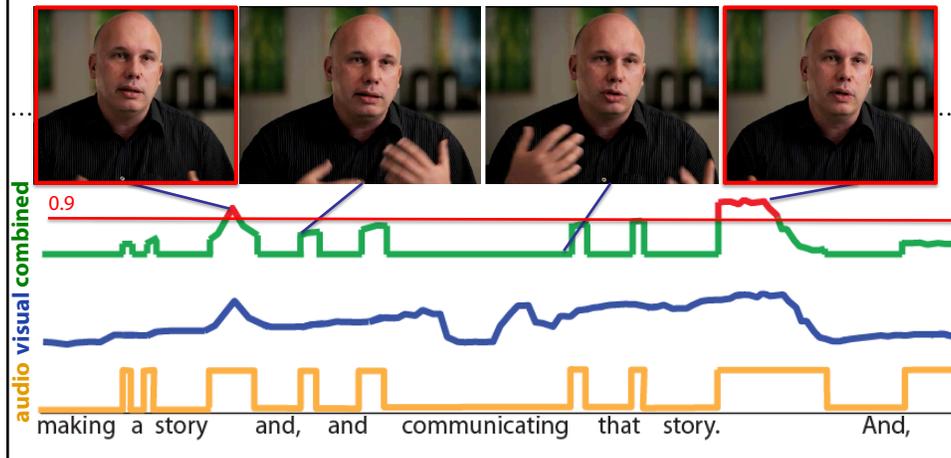
$$\text{Appearance Distance}(i, j) \simeq \chi_{i,j}^2(\text{eyes}) \times \chi_{i,j}^2(\text{mouth}) \times \chi_{i,j}^2(\text{body})$$

Based on [Kemelmacher-Shlizerman 11]



## Cut Suitability Score

$$Score(i) = S_v(i)S_a(i) > 0.9$$



MAN: Sure. **Um**, **some** of our customers, that have been working in the print for **many**, **many** years, are used to be able to express their ideas in, **uh**, very specific ways. **They**, **they** are used to be able to control exactly what goes **on**, **on** the page. They bring the story, they bring the images, and they bring the layout as well, which is a big part of making a story **and**, **and** communicating that **story**. **And**, **um** with **the** **the** digital media and the tablets that are used now to bring **this** content, **uh**, they've quickly run into some, **uh**, limitations with trying to express the same thing

## Hidden Transitions

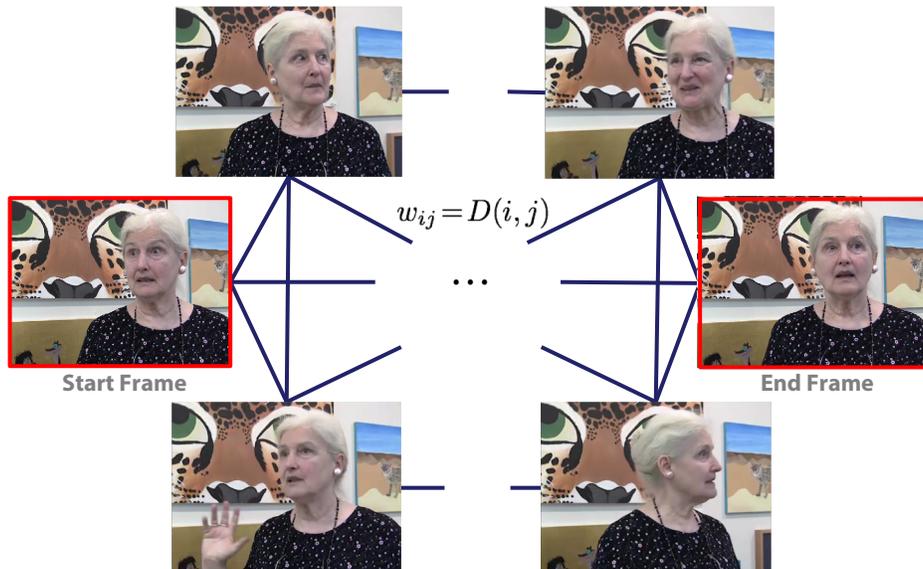


Jump Cut

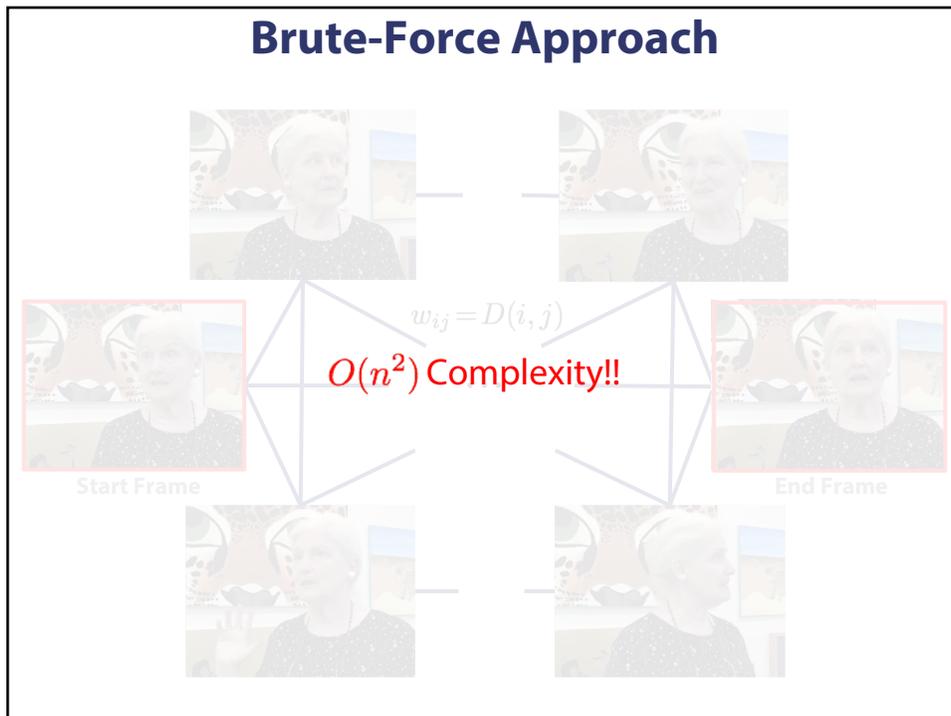


Hidden Transition

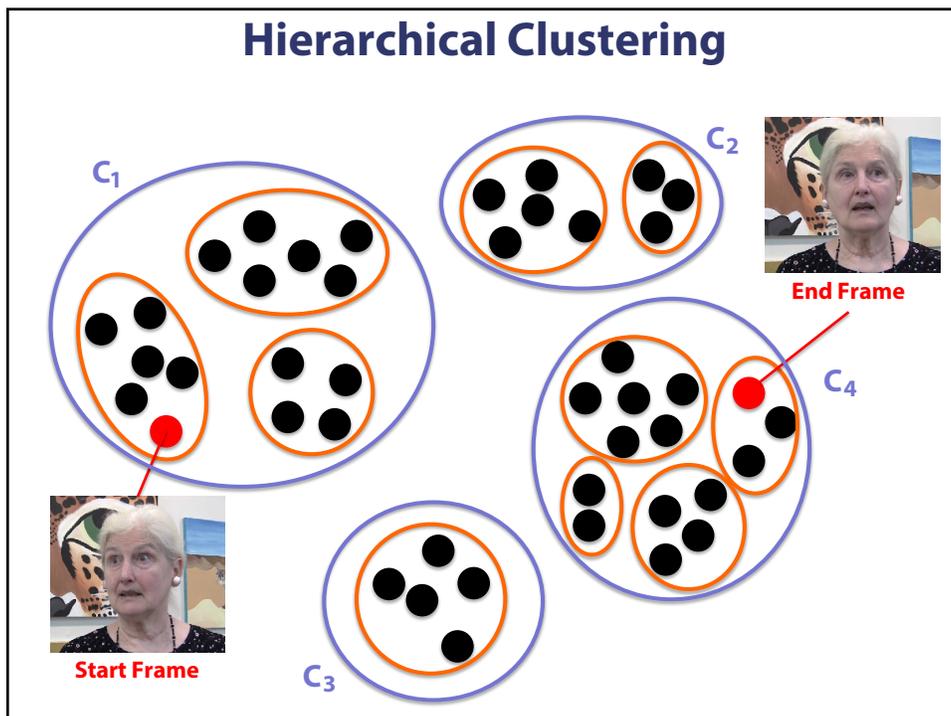
## Brute-Force Approach

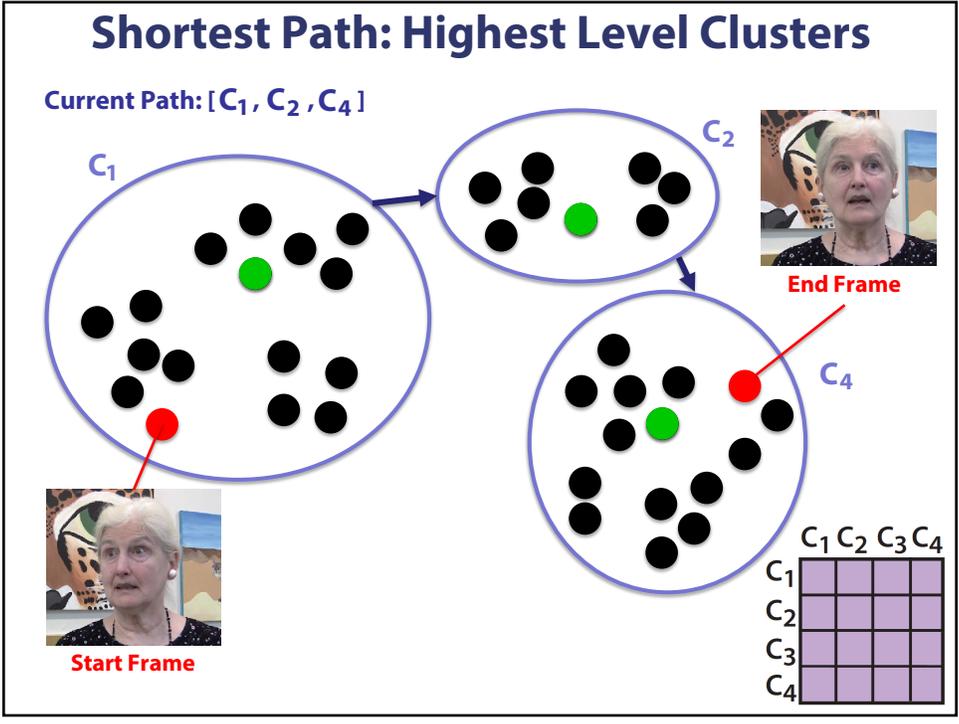
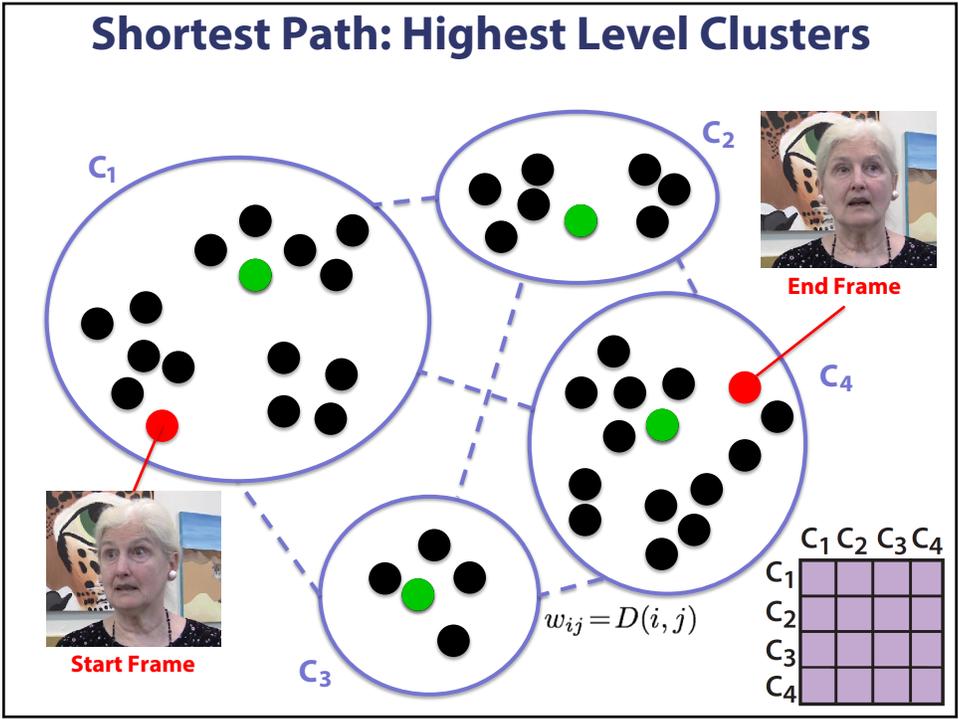


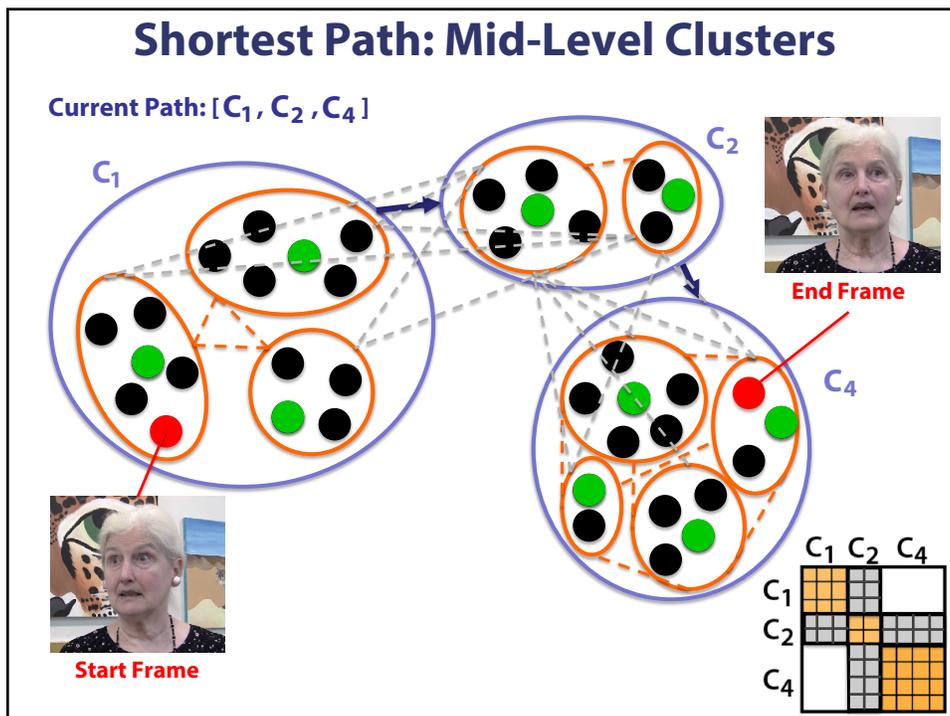
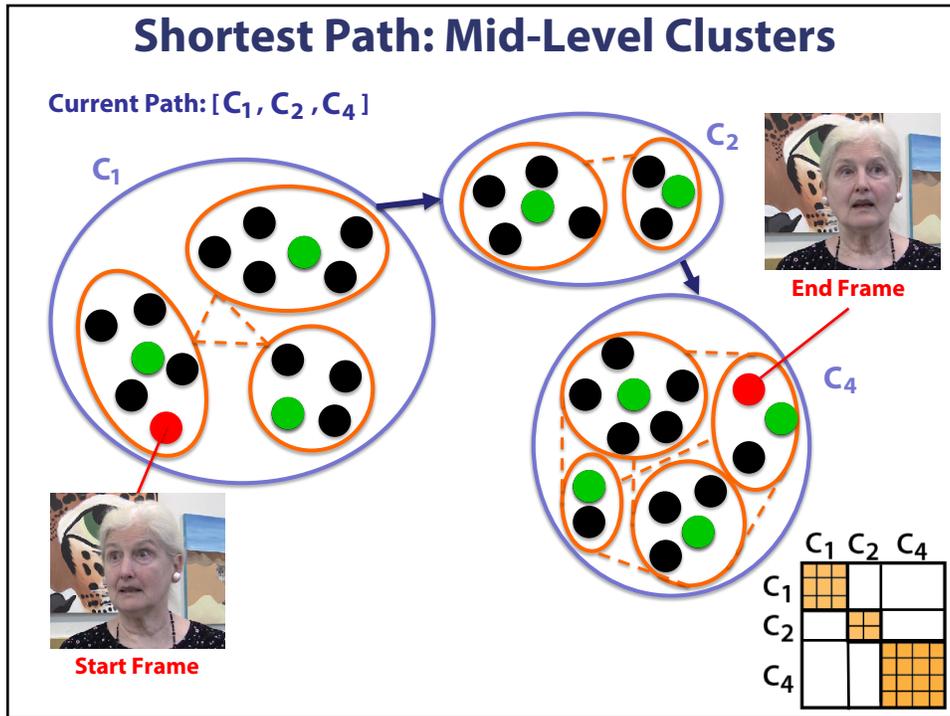
## Brute-Force Approach



## Hierarchical Clustering







### Jump-Cut



### Shortest Path Transition



## Shortest Path

Shortest Path In-Between Frames



## Optical Flow Interpolation

Shortest Path In-Between Frames



Interpolation Frames [ Brox 04 ]

## Data-Driven Timing

Raw Footage Frames



$$N(i, j) = 2$$

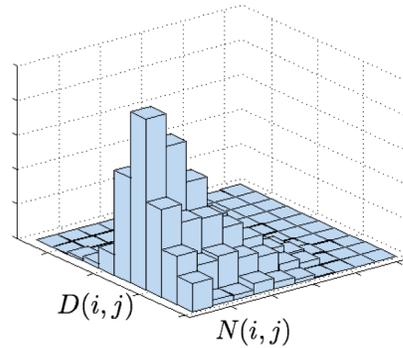
$$D(i, j) = 0.02$$

Number of intermediate frames

$$N(i, j) = |i - j|$$

Frame distance

$$D(i, j)$$

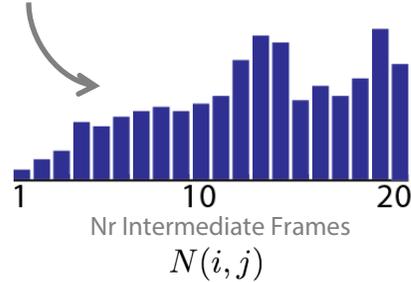


## Setting Number of Frames

Shortest Path In-Between Frames



$$D(i, j) = 3.6$$





## Hidden Transition



## Jump cuts (in red)



**Our result** hidden transitions in blue  
pauses in green



Comp. time:  
clusters 52m  
hidden 8s  
pauses 16s

**Jump cuts (in red)**



**Our result** hidden transitions in blue  
pauses in green



Comp. time:  
clusters 72m  
hidden 9s  
pauses 14s



**Jump cuts (in red)**



**Our result** hidden transitions in blue  
pauses in green



Comp. time:  
clusters 22m  
hidden 5s  
pauses 9s

## Limitations



Poor hidden transitions when endpoints very different  
Assumptions: Still background, fixed lighting and camera position

## User Feedback

### Professional Users

4 Video editors  
5 Journalists

### Interface

Would fit in existing workflow

***Transcript view extremely useful***

Easy to find good cut locations  
Could quickly test different cut possibilities

### Hidden Transitions

Useful alternative to visible transitions  
Good for Web

It begins with reporting.

UC BERKELEY  
GRADUATE SCHOOL  
OF JOURNALISM



## Takeaways

### Choose the *right* representation

Transcript much better than frames for navigation and editing of interview video

### Principle of congruence

For effective interfaces, structure of external representation should match structure of mental representation

[Tversky 02] [Norman 86, 88]

### What is appropriate representation for other types of video?