

000
001
002
003
004
005
006
007
008
009
010
011
012
013

054

055

056

Cross-Language Speech Dependent Lip-Synchronization

057

058

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

079

080

081

082

083

084

085

086

087

088

089

090

091

092

093

094

095

096

097

098

099

100

101

102

103

104

105

106

107

Anonymous WACV submission

Paper ID 141

Abstract

Speech videos such as movie dialogues, public speech, online courses are a great source of infotainment. These videos are often limited by the linguistic constraint of audiences being from different demographics. Vernacular backgrounds that are non-native to the accent or the language of the content producer often makes it difficult for a listener to comprehend the full essence of the content. Such videos are often supplemented with foreign language subtitles which hamper viewing experience. Otherwise, simple audio dubbing in a different language makes the video appear unnatural due to unsynchronized lip motion.

In this paper, we try to address this issue by proposing two lip synchronization methods: 1) Lip-synchronization for change in accent of the same language audio dubbing, and 2) Cross-language lip-synchronization for speech videos dubbed in a different language. We describe an automated pipeline to synchronize the lip movements of the speaker conditioned upon the audio in both cases. Our quantitative evaluation shows high SSIM index between generated cross-language lip-synchronized videos and the original videos. With the help of a user-based study, we verify that our method is preferred over unsynchronized videos.

1. Introduction

Speech videos are an effective way of story-telling. Many socio-political changes are caused by public speeches, movies depict the cultural aspect of societies. Similarly, online instructional videos, especially Massive Open Online Courses (MOOCs), are prime examples of how education can help skill development beyond the boundaries of conventional classrooms. Moreover, the Internet has made it possible for a global inter-cultural exchange of ideas where any information is just a click away. Yet we find limited penetration for these speech videos when they cross international boundaries. For instance, the retention rates in MOOC courses can be as low as 10% [13]. One of the major reasons for this is a cultural gap between the linguistics of the audience and the content pro-

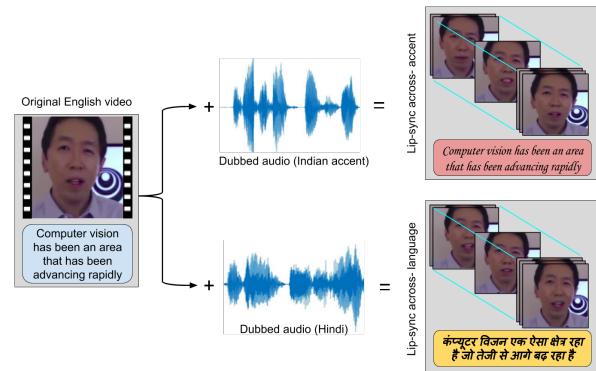


Figure 1. Lip synchronization on Andrew Ng Machine learning tutorial video based on dubbed audio: (top-right) shows Dynamic Programming to synchronize lip-motion of the original English video (left) into Indian English accent, (bottom-right) cross-language lip-sync to synchronize into different language (Hindi).

ducer. Students from different parts of the world often find it difficult to understand the accent and language of the instructors, owing to their non-familiarity with them. This results in slow learning curves as well as dropouts from such online courses. Subtitles in different languages do not lend enough help since they divert the attention of the audience. A quick-fix solution to this would be to dub speech videos in the accent or language of the audience. However, dubbing without lip synchronization makes the video appear unnatural [17].

This problem can be solved by synchronizing the lip motion of the speaker in the target video to be coherent with the dubbed audio. A similar problem exists in the field of computer animation, where lip-motion of the animated characters are constrained upon the textual script of the character. This usually required a human in the loop to manually lay the visemes, hence such a system cannot be scaled for photo-realistic lip-synchronization (lip-sync).

With recent developments in deep neural architectures like Generative Adversarial Networks (GANs) [10, 14], we are now able to generate photo-realistic images conditioned upon an input prior. Similarly, Recurrent neural networks [12] have given way to better sequence learning methods

108 which has become the de facto in time-series data modeling
109 such as speech videos. Both of these methods are
110 quintessentially data-driven approaches, and require large
111 amount of training data.

112 In this paper, we propose ‘Visual Dubbing’ for synchronizing
113 lip motion in speech videos according to the language it is dubbed in.
114 Our main ideas and contributions are two-fold: 1) we propose a model for lip-syncing a target
115 video with audio dubbing in a different accent of the same language, such as Indian English accent or French
116 English accent, as shown in Figure 1 (top); 2) we propose another model to lip-sync the speech video based on audio
117 dubbing in a different language, for instance English video with Hindi audio dubbing, as shown in Figure 1 (bottom).

118 The input to both our models is speech video where the
119 lip-motion of a speaker is clearly visible, and the dubbed audio.
120 The output is a video generated with synchronized lip
121 motion. We also propose a scalable pipeline for dataset creation,
122 which will later be used to train our models. Unlike
123 audio-dubbing which requires professional dubbing artists
124 to give their voices, visual-dubbing does not depend on hu-
125 man visual input for lip-synchronization of a target video.
126 Figure 1 shows the visual dubbing for a video clip of an
127 Andrew Ng MOOC video.

128 We evaluate our generative model based on the struc-
129 tural similarity (SSIM) index of the lip-synchronized videos
130 with respect to the original English videos. Lastly, with the
131 help of a user-based study we show that lip-synchronization
132 makes the speech video more engaging while preserving
133 photorealism.

134 2. Related Work

135 The current process of dubbing involves *translation* to
136 the target language to resemble the lip motion of the source
137 language as much as possible, *recording* of the dubbed content
138 in pace with the original performance, and *editing* of the
139 dubbed soundtrack and lip motion to be temporally close.
140 This process is performed by production companies, and is
141 both time-consuming and expensive. Even after such effort,
142 the result of dubbing is not visually pleasing to the viewer
143 because of clear visual discrepancy between the lip motion
144 and the audio track.[25]

145 This discrepancy between the speech information of the
146 dubbed track and the facial motion of the video track is due
147 to differences in correspondence of phoneme sequences and
148 lip motions [25]. It causes a strong discomfort for viewers,
149 and is also a huge distraction for those who are hearing-
150 impaired, as they rely significantly on lip reading [23, 21].
151 (This difference is one of the reasons why people dislike
152 watching dubbed content [17], because it alters the sound
153 perceived by the observer [22].

154 These findings are the motivation for our work to solve
155 the problem of lip-synchronization.

156 2.1. Lip synchronization

157 The essential component of speech perception in videos
158 and animations is visual cues, such as visemes [26] and
159 phonemes. In particular, this is of importance for hearing-
160 impaired people because they rely on these visual cues to
161 understand videos. [21]. So, for dubbing of videos to an-
162 other language or even another accent of the same language,
163 it is necessary to learn the viseme-phonemic relation be-
164 tween both the accents or languages, so as to synchronize
165 the lip motion with the original performance in the video.

166 The earliest work related to ours could be animation of
167 facial movements in avatars modeled either from audio [15]
168 or text [29, 30]. These work mostly used HMM [9] for se-
169 quential lip trajectory generation [31, 30]. One of the first
170 systems for animating a virtual avatar’s face directly from
171 speech was proposed by [2], which models the joint distri-
172 bution of acoustic and visual speech.

173 The work by Bregler et.al.[3] learns a mapping between
174 the visemes and phonemes for one specific actor and lan-
175 guage, and synthesizes new lip movements through image
176 warping. However, the results of this method fail to dub be-
177 tween different languages and different individuals. Some
178 of the recent work focus on synthesizing photo-realistic lip-
179 motions and facial expressions. Face2Face [28] morphs the
180 facial landmarks of a person in a target video based on those
181 of another actor. However, these kind of models require a
182 human in the loop, which can be quite expensive to scale,
183 and prone to error.

184 The advent of Recurrent Neural Networks, especially
185 LSTMs [12], gave way to better sequence learning, which
186 made generation of features from speech more efficient.
187 Pham et al [24] used CNN followed by an LSTM to gen-
188 erate face parameters from input audio waveform. Karras
189 et al [16] proposed a network consisting of a spatial convolu-
190 tion layer followed by a temporal convolution network on
191 top of fully connected layers to convert speech audio into
192 facial expressions. Chung et al [6] proposed an encoder-
193 decoder convolutional network to jointly embed face and
194 audio. The encoder network consisted of an audio stream
195 and a video stream which merges in a bottleneck represen-
196 tation. This representation is used by the decoder to synthe-
197 size lips-motion video frames.

198 Most similar to our work are [27, 20] which use speech
199 audio represented as MFCC features [27] and text [20]
200 to train an LSTM to produce a sequence of lip landmark
201 points. The lip landmarks are then used to generate mouth
202 texture. Finally this mouth texture is merged with the face
203 in the original frame. Our work is different from [27, 20]
204 in that our method synchronizes lip motion across two
205 different languages, in contrast to just English-to-English.
206 Hence, our challenges include learning higher-level viseme-
207 phonemic relations across languages.

216

3. Method

Instructional videos provide a controlled framework for this problem, since the speakers usually speak scripted dialogues in good lighting facing the camera. The challenge is to model the lips, and generate new lip movements for the same speaker, given the dubbed audio. We consider dubbing to be of two types: a) in the same language as the original speech video but with a different accent, to alleviate the nuances related to unfamiliar accent; and b) where the video is re-dubbed into an entirely different language, to cater to non-native audiences. Both the dubbing types pose different challenges in lip synchronization. For English to non-native English, words remain the same while the pitch, tone and timing of the spoken word change across accents. On the other hand, for cross-language dubbing both words and phonemes change. In this section we will discuss two different methods to address these challenges.

234

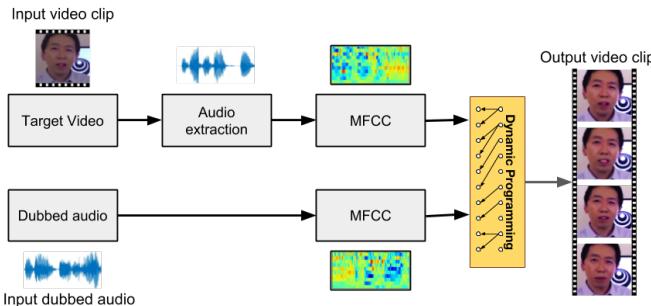


Figure 2. Pipeline for Dynamic Programming: Inputs are the target video and the dubbed audio in Indian accent, which are dynamically synced based on the similarity between the MFCC features of the voice of the native English speaker and the dubbed audio.

249

3.1. Cross-accent lip-sync

While dubbing for different accent, since the same words are spoken by the instructor, all the required viseme sequences are already present in the original video. However, the time instances of the words being spoken might change. Hence, this problem of cross-accent lip-sync can be reduced to a non-linear alignment between original audio and the dubbed audio. This can be done by creating a mapping between different segments of the two audio clips.

In our setup, we have Andrew Ng's Machine Learning audio-visual clip, along with a dubbed audio clip of the same dialogues in Indian-accented English. Spoken words can be broken down into sequences of phonemes. Therefore, we densely segment these audio clips into 25 millisecond (ms) clips. We then use Dynamic Programming [8] to create a dynamic map between the segments of the two audio clips in the feature space. Mel frequency Cepstral Coefficients (MFCC) have been widely used by the speech community, as they provide an optimal encoding of those audio bands which are most relevant to spoken speech. There-

fore our mapping is based on finding the nearest neighbor of a segment from one audio clip in another. This is then converted into a mapping between the corresponding video frames. This procedure is illustrated in Figure 2, where we show the overall pipeline for non-linear alignment between original English audio and non-native English accent dubbed audio using dynamic programming.

Each segment of the original video clip is assigned a new time-stamp based on the mapping. This creates a non-uniform separation between adjacent segments of the original video clip. We render a new lip-synced video by interpolating frames in the original clip to fill the voids, and down-sampling to remove excess frames. Here we propose to use a Dynamic Programming algorithm (Algorithm 1), where x and y are the two inputs (in this case, the MFCCs of original and dubbed audio), L_x and L_y are the lengths of the respective audio segments. C represents the cumulative cost of the dynamic mapping with time, and M is used to find the optimal path for the mapping between input and output frames. $xTOy$ is the variable that records the mapping from x to y . Similarly, to find $yTOx$, the following can be appended at line 18: $yTOx[q - 1] \leftarrow p - 1$.

Algorithm 1 Dynamic Programming

```

1:  $init \leftarrow 5$ 
2:  $C[i][j] \leftarrow 0, i = 0 \text{to} L_x, j = 0 \text{to} L_y$ 
3: for  $i = 1 \text{to} L_x$  do  $C[i][0] \leftarrow i * init$ 
4: for  $j = 1 \text{to} L_y$  do  $C[0][j] \leftarrow j * init$ 
5:  $M[i][j] \leftarrow 0, i = 0 \text{to} (L_x - 1), j = 0 \text{to} (L_y - 1)$ 
6: for  $i = 1 \text{to} L_x$  do
7:   for  $j = 1 \text{to} L_y$  do
8:      $min1 \leftarrow C[i - 1][j - 1] + \text{cost}(input_x[i - 1], input_y[j - 1])$ 
9:      $min2 \leftarrow C[i - 1][j] + init$ 
10:     $min3 \leftarrow C[i][j - 1] + init$ 
11:     $C[i][j] \leftarrow \min(min1, min2, min3)$ 
12:    if  $C[i][j] = min1$  then  $M[i - 1][j - 1] \leftarrow 1$ 
13:    if  $C[i][j] = min2$  then  $M[i - 1][j - 1] \leftarrow 2$ 
14:    if  $C[i][j] = min3$  then  $M[i - 1][j - 1] \leftarrow 3$ 
15:  $p \leftarrow L_x, q \leftarrow L_y$ 
16: while  $p \neq 0$  and  $q \neq 0$  do
17:   if  $M[p - 1][q - 1] = 1$  then
18:      $xTOy[p - 1] \leftarrow q - 1, p \leftarrow p - 1, q \leftarrow q - 1$ 
19:   else if  $M[p - 1][q - 1] = 2$  then  $p \leftarrow p - 1$ 
20:   else if  $M[p - 1][q - 1] = 3$  then  $q \leftarrow q - 1$ 
  
```

3.2. Cross-language lip-sync

Major challenges in lip-syncing audio of a foreign language (e.g. Hindi) on video of original language (e.g. English) are the differences in their grammatical structure and set of phonemes. One way is to directly generate lip-images

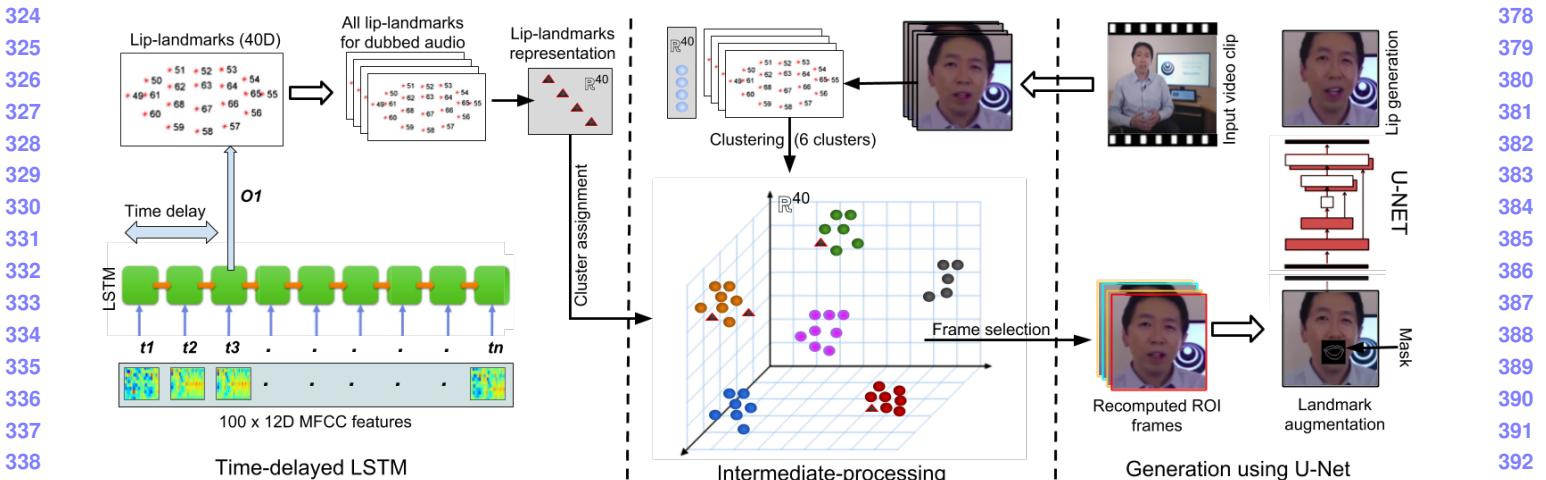


Figure 3. Cross-language lip-sync: (left) Pipeline for training LSTM with Hindi speech and lip landmarks, (center) shows reassignment of frames for each predicted lip-landmark using intermediate-processing step, (right) pipeline for inferring using U-Net on frames from English video.

conditioned upon the foreign language audio and target video. But such end-to-end systems require large amount of training data to learn the complex audio-visual relation between the two modalities [6]. At the same time, recent developments in generative networks [10, 14] have yielded impressive results. Considering these factors, we first learn an embedding between Hindi audio and lip-landmarks. This allows us to predict lip-landmarks from a relatively smaller speech corpus. From these predicted lip-landmarks, we generate mouth images over the original English video to match the Hindi audio. This entire two-step pipeline can be seen in Figure 3 (which also includes an intermediate processing step, discussed in Section 4.3).

3.2.1 Audio to Lip Landmarks

The first step is to encode audio into lip-landmarks. For each phoneme there exists a viseme, and the lip-motion responsible for the transition between two different visemes depends on its location in the viseme-sequence that constitutes the spoken word. This makes audio to lip-landmarks a sequence modeling problem. Hence, similar to [27, 20], we use an LSTM [12] to encode audio. This LSTM takes audio MFCC features at each time step as input, and predicts lip-landmarks at time step 't' after each input, and is therefore called Time-Delayed LSTM (TD-LSTM) [11]. During training, the TD-LSTM is trained on Hindi speech audio-visual. The input is MFCC for every 25ms audio segment at 10ms time step, and the output is the lip-landmark of the Hindi speaker at the 200ms delayed frame, as shown in Figure 3 (left). During inference, given a new audio sample, we use the predicted lip-landmarks as the prior to generate the mouth (lip-region) in the second step.



Figure 4. Generation results for the U-Net: (top) input to U-Net, (middle) generated images, (bottom) ground truth

3.2.2 Lip Landmarks to Generated Faces

Once lip landmarks are predicted from the audio in foreign language, in the second step the lips of the speakers in the original video must be modified to match these landmarks. To solve this problem, we use a U-Net similar to [20] to generate mouth of the speaker based on an encoded prior. During training, the input to the network is the face image of the speaker, with the mouth masked by a black box of constant size, and the original landmarks in the face drawn as a white polygon, see Figure 3 (right). The output of the network is the original face image. This allows the network to learn to generate actual face of the speaker with the lip-region conditioned upon the lip-polygon on the masked face.

As L1 loss is commonly used while generating images, we use this to train the U-Net. In addition, since our main focus is on correctly generating mouth region of the speaker, we add another loss term to penalize wrongly predicted pixels in that region. Considering the mean of the black mask as the center of the mouth region, we add a Gaussian weight kernel G_{loss} to the L1 loss such that the weight of this loss decreases radially from the center of the mouth to the face extremities. Formally, for a ground truth

432 \hat{y} and the predicted face frame y , where any pixel location
 433 is represented by (i, j) , our loss is defined as:
 434

$$L = L1 * (1 + G_{loss}) \quad (1)$$

437 where;

$$L1 = \sum_{i,j} \|\hat{y}_{ij} - y_{ij}\| \quad (2)$$

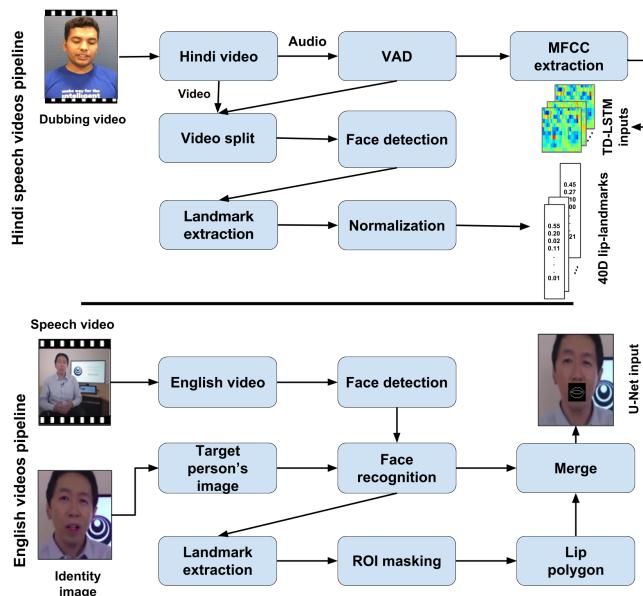
$$G_{loss} = \sum_{i,j} c * \exp \frac{(i - u_i) * (j - u_j)}{v_{ij}} \quad (3)$$

444 In equation 3, c is a normalization constant, u_i and u_j
 445 represent the mean pixel location of the black mask (mouth
 446 region), and v_{ij} represents the covariance.
 447

448 During inference, the mouth in every frame of the video
 449 is replaced by a constant black square, and a white polygon
 450 of the lip landmarks predicted by the LSTM network from
 451 the previous step. Thus, the U-Net will then generate faces
 452 according to the Hindi dubbing audio. Unlike [27, 20], we
 453 train our network on multiple sources, allowing the network
 454 to generalize over multiple speakers, as shown in Figure 4.
 455

3.3. Dataset

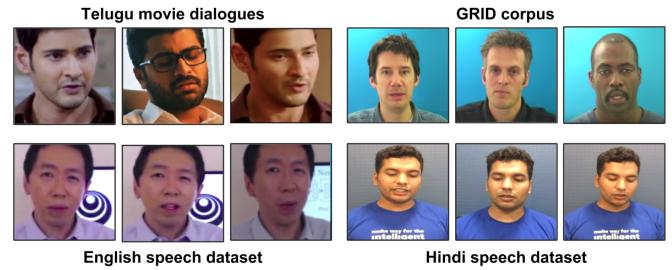
456 In this section, we discuss our dataset curation pipeline
 457 used in our cross-language lip-synchronization method. We
 458 require two different datasets, one for each language: (i)
 459 Hindi speech dataset, for training time-delayed LSTM, and
 460 (ii) English videos dataset, to train U-Net for lip-generation.
 461



481 Figure 5. Dataset curation: (top) shows pipeline to create dataset
 482 from Hindi speech audio to give MFCC features and lip landmarks
 483 for training Time delayed LSTM (TD-LSTM); (bottom) shows
 484 pipeline to create dataset from Andrew Ng (and other) videos in
 485 English to give masked frames, for training U-Net.

3.3.1 Hindi speech dataset

486 As we wish to learn an encoding from Hindi audio to lip
 487 landmarks, we require a dataset consisting of Hindi audio
 488 to train a time-delayed LSTM. Since parallel audio speech
 489 corpus is difficult to find, and because dubbing is mostly a
 490 post-production phenomenon, we record 5 hours of audio-
 491 visual data of a native Hindi speaker narrating articles from
 492 Hindi newspapers and stories. Using voice activity detec-
 493 tion [1], the video clips are segmented to give continuous
 494 segments of speech clips. For 5 hours of speech data, we
 495 get 5000 video clips of average length 2 seconds. Each such
 496 video clip is then sampled at 25 frames per second (fps). To
 497 obtain landmarks, we use a HOG-based face detector using
 498 dlib [18] to find the speaker's face in the clip, and predict
 499 68 face landmarks using dlib. We then choose the land-
 500 marks corresponding to mouth region (landmarks 49 to 68)
 501 a.k.a. lip-landmarks, and normalize them. For each video
 502 clip segmented by voice activity detection, these normalized
 503 lip-landmarks are saved. Similarly, for each video clip we
 504 extract audio and sample it at 100Hz. We extract MFCCs
 505 for each sampled segment of the extracted audio clip. The
 506 training set consisted of 90% of total dataset, validation was
 507 done on rest of 10% of the dataset. The dataset curation
 508 pipeline for Hindi Speech can be seen in Figure 5 (top).



511 Figure 6. Random frames from Telugu movie dialogue clips (top-
 512 left), GRID corpus (top-right), English speech dataset(bottom-
 513 left), Hindi speech dataset (bottom-right).
 514

3.3.2 English speech dataset

515 As our aim is to generate lip-synced Andrew Ng's machine
 516 learning videos with Hindi dubbing, we use 20 Andrew Ng
 517 videos to create a dataset of English speech videos. The
 518 input to our U-Net is frames from instructional video clips
 519 from the English speech dataset, where the pixels in the
 520 mouth region of the instructor are masked with the wire-
 521 frame structure of the lip-landmarks. For each frame where
 522 the face of the speaker has been detected, the face region
 523 is noted as the bounding box of the 68 landmarks, similar
 524 to that in Hindi speech dataset. The square region around
 525 the face with 1.5 times the face width is extracted. This
 526 results in images with full visibility of the instructor's face.
 527 The mouth region of each face is considered as the bound-
 528 ing box around the mean of the mouth landmarks (49 to
 529 68).
 530

540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
68), and of width 0.25 times the width of the face region.
It is then replaced with a black mask and a white polygon
connecting the lip landmarks, and resized to the input shape
of the U-Net. The output of the U-Net is the original face
image. The training set consisted of 10 video clips while
the validation was done on remaining 10 video clips. This
pipeline is shown in Figure 5 (bottom).

548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
It was important to not let the network overfit on the
input images. We therefore used multiple sources of images
as the training dataset for the U-Net — frames extracted
from 1) Telugu movies, 2) videos of Andrew Ng’s deep
learning.ai lectures, 3) GRID corpus [7], 4) Hindi Speech
dataset. Table 1 details the number of frames from each
source. Figure 6 shows some randomly sampled frames
from these datasets.

Source	Train	Validation
Telugu movies	37130 frames	4159 frames
English speech	24359 frames	16035 frames
GRID	13350 frames	1500 frames
Hindi speech	37790 frames	3714 frames

Table 1. Number of images in Train and Validation sets for training U-Net.

562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589
For each source of images, we detected faces and pre-
dicted landmarks using [4].

3.4. Representation

590
591
592
593
Audio: The input to dynamic programming algorithm
and time-delayed LSTM is raw audio from the Hindi speech
video clips represented as Mel-frequency cepstral coefficients
(MFCC). The audio was sampled at 100Hz with
each sample being of length 25ms. We extract 13 coefficient
MFCC feature for each sampled segment, but only use
12 coefficient as the feature representation, discounting the
first feature. For TD-LSTM, length of each input training
sample is 100 (or total 1 second), with shape 100 x 12D.

594
595
596
597
598
599
600
598
Lip-landmarks: The output of the TD-LSTM are lip-
landmarks of the speaker of Hindi speech video clips i.e. 40
dimensional (D) normalized lip-landmarks for each frame
(20 lip-landmarks x 2D). The output is computed at a delay
of 200ms, and is represented as 1 x 40D vector.

601
602
603
604
605
606
607
608
609
Images: The input to our U-Net is a masked face image
of size 256 x 256 x 3D similarly output is a 256 x 256 x 3D
image consisting original face of the speaker.

4. Implementation

4.1. TD-LSTM

610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
Our proposed TD-LSTM model consists of a single layer
LSTM with 60 neurons in the hidden layer, followed by
a 40D dense layer. We also up-sample each video clip at
100Hz to compute lip-landmarks. In each forward pass,

637
638
639
640
641
642
643
644
645
646
647
the network takes 100 time-steps MFCC features (100 x
12D) and predicts the 20th up-sampled lip-landmark frame
(1x40D). This is done densely for each Hindi audio-visual
clip. This results in an offset in the prediction of lip-
landmarks of 200ms at the beginning, and 800ms at the
end of the video. We compensate for this by replicating
the first and the last frame’s predicted landmarks respec-
tively for the appropriate number of frames. We also imple-
mented TD-LSTM with 500ms and 800ms time-delays. But
we found very little perceptual difference between the re-
sults, and therefore chose 200ms delay. Using Bidirectional
TD-LSTM also did not perceptually affect the results.

648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
689
We implemented the network in Keras deep learning
framework [5], with a batch size of 64, mean square loss,
and Adam [19] as the choice of optimizer. We trained
our network for 20 epochs, with a total time of around
4 hours on Nvidia GTX 1080 Ti, when the loss started
plateauing. Pre-training with 10% videos randomly sam-
pled from GRID corpus resulted in faster saturation of loss.
As the TD-LSTM learns to encode dubbed audio into lip-
landmarks, we observed that it captures the artifacts cor-
responding to the dubbing artist, such as thickness of lips and
span of mouth. Hence we found pre-training with GRID
corpus also helps in generalization, in case of different dub-
bing artists.

4.2. U-Net

690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
769
We use U-Net architecture similar to [20] We trained the
U-Net on 4 NVIDIA TitanX GPUs, using a batch size of
16, counting 4 batches per iteration, until \approx 5000 iterations.
This took \approx 2 seconds per iteration, and occupied
 \approx 3.3GB of memory including model weights and images
kept in the buffer. As U-Net has been trained on Andrew
Ng’s lip-landmarks to predict original frame, it also learns
an undesirable mapping between jaw location and the shape
of lip in the output. This badly affects the generation of lip
in the instances where the predicted landmarks correspond
to a closed lip while the target frame has mouth open, or
vice-versa. To overcome this, we introduce an *intermediate-
processing* step between TD-LSTM and U-Net.

4.3. Intermediate processing

770
771
772
773
774
775
776
777
778
779
779
We normalize the lip-landmarks from all the frames in
the target instructional video in English, and group them
into 6 clusters. The frames in the target video nearest to
the cluster centroids are chosen to represent the clusters,
with their landmarks as their new centroids. All the lip-
landmarks predicted by TD-LSTM are then assigned to a
cluster based in their distances from the new centroids. This
allows the predicted lip-landmarks to be assigned appropri-
ate face frames. Only these 6 frames are then fed to the
U-Net (after masking the mouth region). This results in a
set of generated frames consisting of lip-synced mouth re-

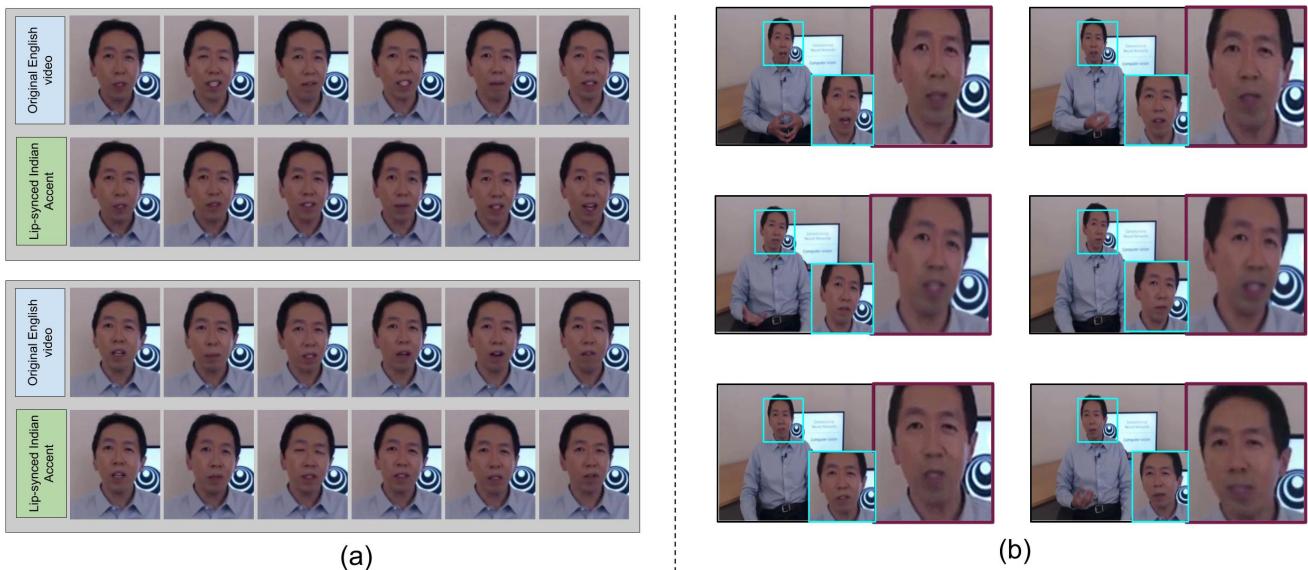


Figure 7. Qualitative results for cross-accent and cross-language lip-sync: (a) In both the examples, frames are sampled at 3 fps from the original instructional video in English (top) and cross-accent lip-synced video (bottom). (b) Each of the 6 images depicts original English video (left) along with its enlarged ROI, (right) shows our generated Hindi lip-synced video.

gions but in only 6 distinct facial poses corresponding to the new centroids. This results in a jittery face video with proper lip-synchronization.

4.4. Homograph computation

The frames generated from the U-Net are slightly blurred, therefore we use a pre-trained CNN deblur network as in [6], trained on facial images for sharpening. We then compute the pairwise homography between each generated video frame and that of the original instructional video clip using all the 3D face-landmarks predicted using [4], except those corresponding to eyes and lip region. This gives a transformation matrix between the frame pairs. We then crop a rectangular mouth region in the generated video frames, with its center at the mean of lip-landmarks, and its length twice of that between the mean of lip-landmarks and the landmark corresponding to the tip of the nose. This region of interest (ROI) is augmented over the original video using the the computed transformation matrix. All these ROI augmented frames along with Hindi dubbing are then used to create the final Hindi lip-synced video.

4.5. Evaluation metric

Audio-visual data generation is done for human consumption, hence its evaluation is subjective. Hence, we conducted a user-based study to evaluate the outputs of our cross-accent and cross-language methods. This allowed us to get an average subjective evaluation.

Structural similarity (SSIM) index [32] is a widely used metric to evaluate the quality of generated images and videos. We compare the SSIM index to evaluate the gen-

eration of the U-Net.

5. Results

5.1. User-based study

To check the quality of our proposed models, we asked 20 different people to rate the lip-unsynced video and the generated lip-synced video pairs, for each of 10 Andrew Ng's ML video clips, for each of the two cases of cross-accent and cross-language. The 10 video clip-pairs were of upto 1 minute duration each.

5.1.1 Cross-accent lip-sync

For each of the pairs of un-synced dubbed video where the Indian English-accented dubbed audio is naively overlaid on the original English video, and our dynamically lip-synced videos, we randomly selected 5 video pairs for each subject and asked them to rank the videos between 1 (Hard to understand) to 5 (Easy to understand) based on their comfort. As shown in Table 2, users preferred our dynamically programmed lip-synchronized videos.

	US(N)	S(N)	US(F)	S(F)
Dynamic	3.0	4.6	1.9	3.2

Table 2. Mean scores for Dynamic Programming (Dynamic) on Indian-English: for un-synced speech overlay ('US'), and lip-synced version 'S', by naive (N) and familiar (F) users.

756

5.1.2 Cross-language lip-sync

We perform a similar user-based experiment with cross-language lip-synchronized videos — we showed 10 videos, with Hindi audio naively overlaid (un-synced), and Hindi lip-synced video to 20 users. Since comfort is subjective and ill-defined, we asked users to rate the percentage of lip-synchronization perceived by them for each pair. As shown in Table 3, the means of the comfort score and percentage lip-synchronization were higher for our cross-language lip-synced videos. The average comfort rating across users for each video pair can be seen in Figure 8 (a), where as average percentage lip-synchronization can be seen in Figure 8 (b). We also show qualitative results of cross-accent and cross-language lip-synchronization in Figure 7 (a) and (b) respectively.

	C-US	C-S	LS%-US	LS%-S
Mean	2.51	3.1	23.86	45.95
Std-dev	1.07	0.6	25.9	24.1

Table 3. Mean scores and standard deviation for Cross-language lip-sync on Hindi: (C) comfort level for (US) un-synced speech overlay, and (S) lip-synced version; (LS%) Lip-Sync percentage for (US) un-synced and (S) lip-synced versions.

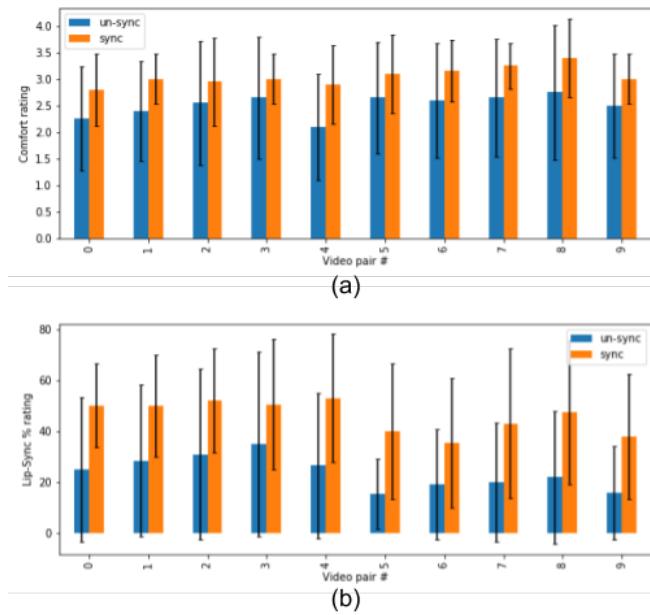


Figure 8. User feedback for cross-language lip-sync corresponding to 10 video pairs - (a) shows average comfort rating and its standard deviation for lip-unsynced (blue) and lip-synced videos (orange), (b) shows average percentage of perceived lip-synchronization and its standard deviation for lip-unsynced (blue) and lip-synced videos (orange).

5.2. Quality of generation

We compare SSIM index to evaluate the frame quality in original English video and its Hindi lip-synced version for each of the 10 video pairs. The average SSIM index across all the generated frames w.r.t the frames in the original videos was 0.98, with an overall standard deviation of 0.01. To evaluate the generation output of our cross-language model, we also computed SSIM scores for the 4 datasets we used to train U-Net. The average SSIM score for each of these dataset can be seen in Figure 9, with mean average SSIM score for all the dataset to be 0.58 with the overall standard deviation of 0.05.

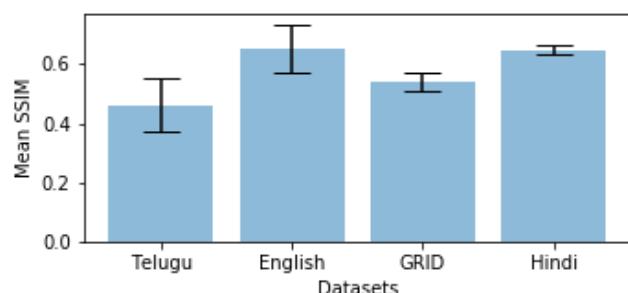


Figure 9. Mean and standard deviation of SSIM scores for various datasets used to train U-Net

6. Discussion

In this work, we assume the availability of dubbed audio, which can be automated using Machine translation (MT) systems and text-to-speech (TTS) synthesizers. However, imperfections in MT translations and lack of personality in the TTS-synthesized speech could make them unsuitable for instructional videos. Furthermore, handling multiple speakers, extreme head poses, and robust key point tracking present future scope of improvement. Lastly, we believe this work can help expand the reach of instructional videos across diverse linguistic groups.

7. Conclusion

We propose two different lip-synchronization methods for educational videos for same language with different accents, and two different languages to improve instructor-student engagement during online video lectures. We detail our pipelines for dataset creation for audio-to-lip-landmarks as well as lip-landmarks-to-mouth-generation. Our user-based-study shows that lip-synchronization can improve effectiveness of content delivery through dubbed speech videos.

References

- [1] WEBRTC, VAD. <https://webrtc.org/>.

- | | | |
|-----|--|-----|
| 864 | [2] M. Brand. Voice puppetry. In <i>Proceedings of the 26th annual conference on Computer graphics and interactive techniques</i> , pages 21–28. ACM Press/Addison-Wesley Publishing Co., 1999. | 918 |
| 865 | | 919 |
| 866 | | 920 |
| 867 | | 921 |
| 868 | [3] C. Bregler, M. Covell, and M. Slaney. Video rewrite: Driving visual speech with audio. In <i>Proceedings of the 24th annual conference on Computer graphics and interactive techniques</i> , pages 353–360. ACM Press/Addison-Wesley Publishing Co., 1997. | 922 |
| 869 | | 923 |
| 870 | [4] A. Bulat and G. Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In <i>International Conference on Computer Vision</i> , 2017. | 924 |
| 871 | | 925 |
| 872 | | 926 |
| 873 | [5] F. Chollet et al. Keras, 2015. | 927 |
| 874 | | 928 |
| 875 | [6] J. S. Chung, A. Jamaludin, and A. Zisserman. You said that? In <i>British Machine Vision Conference</i> , 2017. | 929 |
| 876 | | 930 |
| 877 | [7] M. Cooke, J. Barker, S. Cunningham, and X. Shao. An audio-visual corpus for speech perception and automatic speech recognition. <i>The Journal of the Acoustical Society of America</i> , 120(5):2421–2424, 2006. | 931 |
| 878 | | 932 |
| 879 | [8] T. H. Cormen. <i>Introduction to algorithms</i> . 2009. | 933 |
| 880 | | 934 |
| 881 | [9] R. Durbin, S. R. Eddy, A. Krogh, and G. Mitchison. <i>Biological sequence analysis: probabilistic models of proteins and nucleic acids</i> . Cambridge university press, 1998. | 935 |
| 882 | | 936 |
| 883 | [10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In <i>Advances in neural information processing systems</i> , pages 2672–2680, 2014. | 937 |
| 884 | | 938 |
| 885 | [11] A. Graves and J. Schmidhuber. Framewise phoneme classification with bidirectional lstm and other neural network architectures. <i>NEURAL NETWORKS</i> , pages 5–6, 2005. | 939 |
| 886 | | 940 |
| 887 | [12] S. Hochreiter and J. Schmidhuber. Long short-term memory. 9:1735–80, 12 1997. | 941 |
| 888 | | 942 |
| 889 | [13] K. S. Hone and G. R. El Said. Exploring the factors affecting mooc retention: A survey study. <i>Computers & Education</i> , 98:157–168, 2016. | 943 |
| 890 | | 944 |
| 891 | [14] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. <i>arxiv</i> , 2016. | 945 |
| 892 | | 946 |
| 893 | [15] P. Kakumanu, R. Gutierrez-Osuna, A. Esposito, R. Bryll, A. Goshtasby, and O. Garcia. Speech driven facial animation. In <i>Proceedings of the 2001 workshop on Perceptive user interfaces</i> , pages 1–5. ACM, 2001. | 947 |
| 894 | | 948 |
| 895 | [16] T. Karras, T. Aila, S. Laine, A. Herva, and J. Lehtinen. Audio-driven facial animation by joint end-to-end learning of pose and emotion. <i>ACM Transactions on Graphics (TOG)</i> , 36(4):94, 2017. | 949 |
| 896 | | 950 |
| 897 | [17] R. Kilborn. Speak my language': current attitudes to television subtitling and dubbing. <i>Media, culture & society</i> , 15(4):641–660, 1993. | 951 |
| 898 | | 952 |
| 899 | [18] D. E. King. Dlib-ml: A machine learning toolkit. <i>Journal of Machine Learning Research</i> , 10(Jul):1755–1758, 2009. | 953 |
| 900 | | 954 |
| 901 | [19] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. <i>CoRR</i> , abs/1412.6980, 2014. | 955 |
| 902 | | 956 |
| 903 | [20] R. Kumar, J. Sotelo, K. Kumar, A. de Brébisson, and Y. Ben-gio. Obamanet: Photo-realistic lip-sync from text. | 957 |
| 904 | | 958 |
| 905 | | 959 |
| 906 | | 960 |
| 907 | | 961 |
| 908 | | 962 |
| 909 | | 963 |
| 910 | | 964 |
| 911 | | 965 |
| 912 | | 966 |
| 913 | | 967 |
| 914 | | 968 |
| 915 | | 969 |
| 916 | | 970 |
| 917 | | 971 |