

DEEP LEARNING FOR COMPUTER VISION

Summer School at UPC TelecomBCN Barcelona. June 28-July 4, 2018



Instructors



Organized by



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



Supported by



Co-funded by the
Erasmus+ Programme
of the European Union

GitHub Education

Google Cloud Platform

+ info: <http://bit.ly/dlcv2018>

<http://bit.ly/dlcv2018>



#DLUPC

Day 4 Lecture 6

Speech and Vision



Xavier Giro-i-Nieto

xavier.giro@upc.edu

Associate Professor

Universitat Politècnica de Catalunya
Technical University of Catalonia



Acknowledgments



Antonio
Bonafonte



Santiago
Pascual



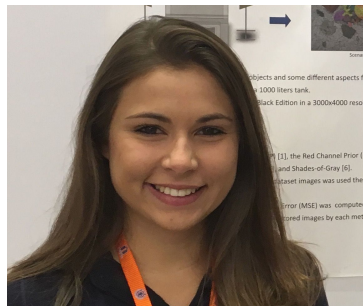
Marta R.
Costa-jussà



Jose A.
Rodríguez Fonollosa



Acknowledgments



Amanda
Duarte



Janna
Escur



Alejandro
Woodward



Fran
Roldan

Janna Escur, "[Exploring Automatic Speech Recognition with Tensorflow](#)" UPC ETSETB 2018.

Roldan F, Pascual S, Salvador A, McGuinness K, Giro-i-Nieto X, "Speech-conditioned Face Generation with Deep Adversarial Networks" (under progress)

Acknowledgments

Speech Recognition



**DEEP LEARNING
FOR SPEECH AND LANGUAGE**

Winter School at UPC TelecomBCN Barcelona. 24-30 January 2018.

Instructors

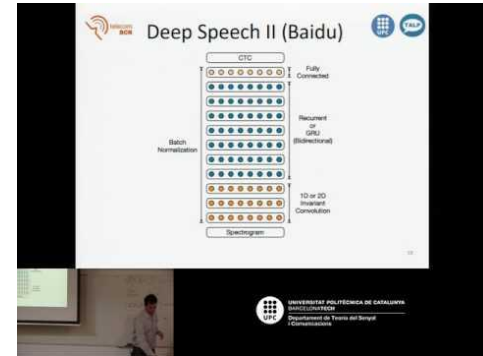
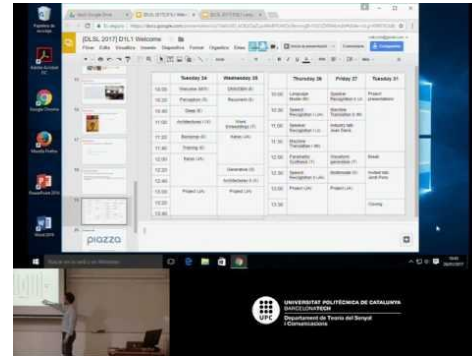


Marta R. Costa-jussa, José A. R. Fonollosa, Santiago Pascual, Javier Hernando, Antonio Bonafonte, Xavier Giró-i-Nieto

Organized by  **Supported by**    **GitHub Education**

+ info: <https://telecombcn-dl.github.io/2018-dlsl/>

- [1st edition](#) (2017)
- [2nd edition](#) (2018)



Speech Synthesis



**Day 3 Lecture 5
Parametric Speech
Synthesis**

Antonio Bonafonte



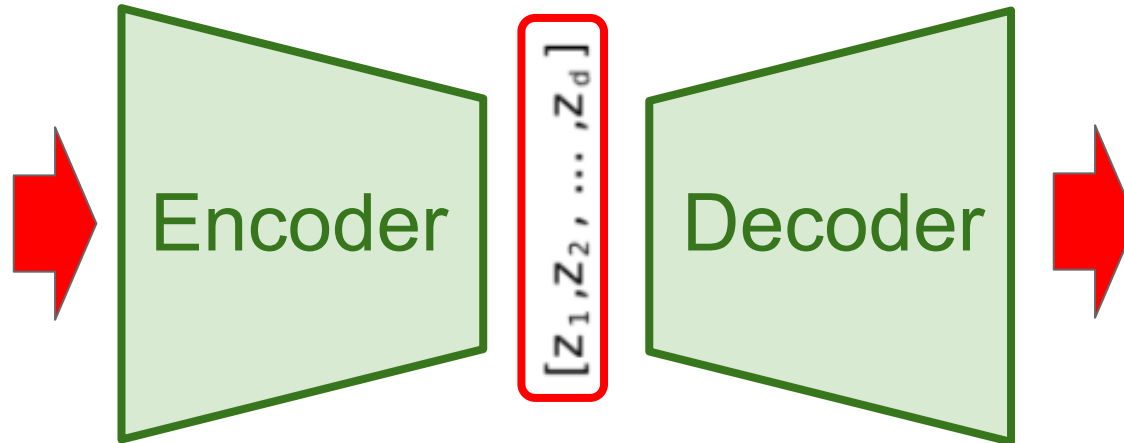
**Day 4 Lecture 4
Speech Synthesis
WaveNet**

Antonio Bonafonte

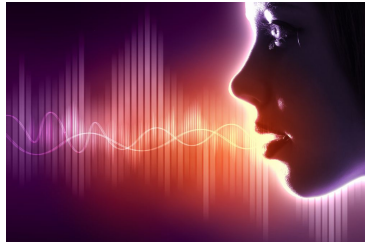


Representation or Embedding

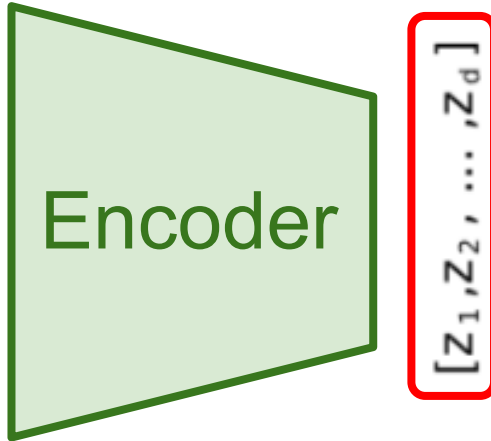




Representation or Embedding

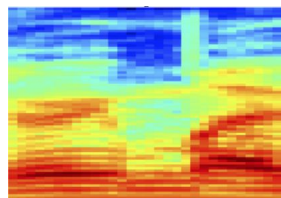


Speech

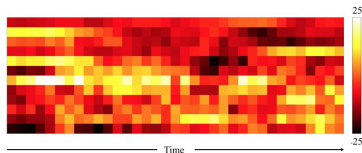
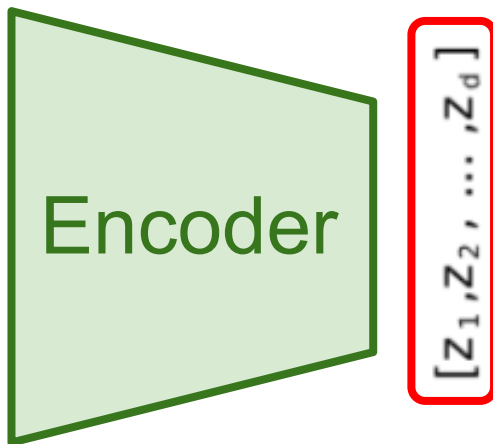


Representation or Embedding

Raw



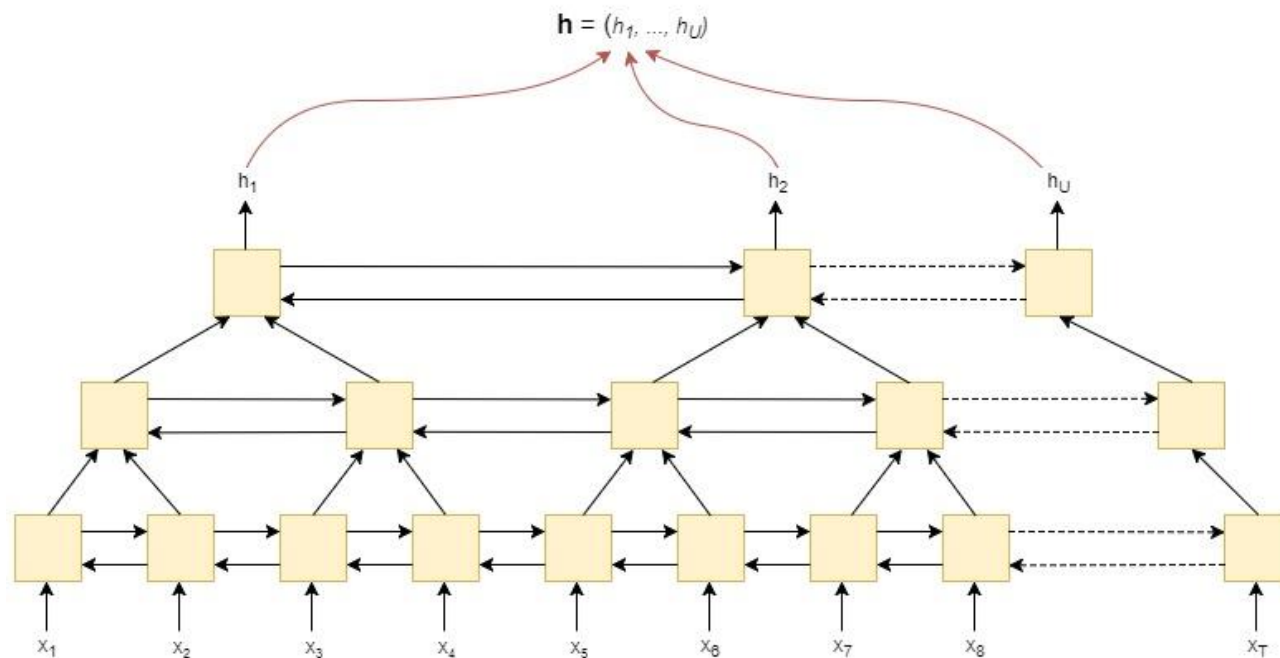
Mel spectrum



MFCC

More details: [Slides](#) by Kishore Prahallad (CMU)

LAS Speech Encoder

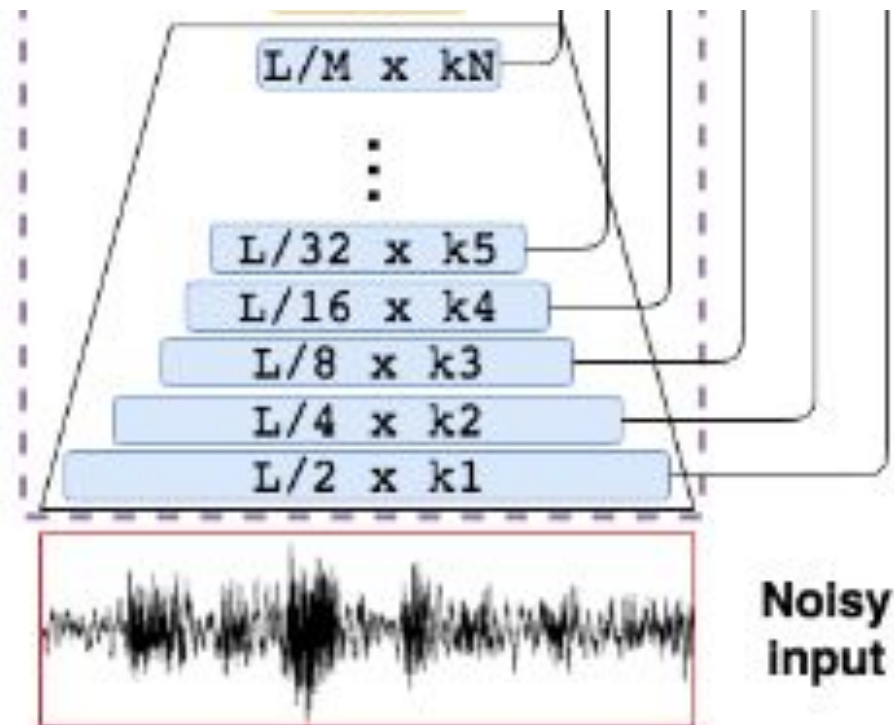


Pyramidal BLSTM

Chan, William, Navdeep Jaitly, Quoc Le, and Oriol Vinyals. ["Listen, attend and spell: A neural network for large vocabulary conversational speech recognition."](#) ICASSP 2016.

Janna Eскур, ["Exploring Automatic Speech Recognition with Tensorflow"](#) UPC ETSETB 2018.

SEGAN Speech Encoder



Pascual, Santiago, Antonio Bonafonte, and Joan Serra. "[SEGAN: Speech enhancement generative adversarial network.](#)" Interspeech 2017.



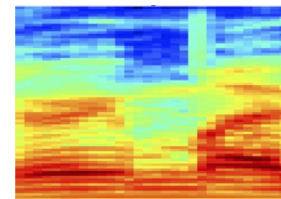
Representation or Embedding

$[z_1, z_2, \dots, z_d]$

Decoder

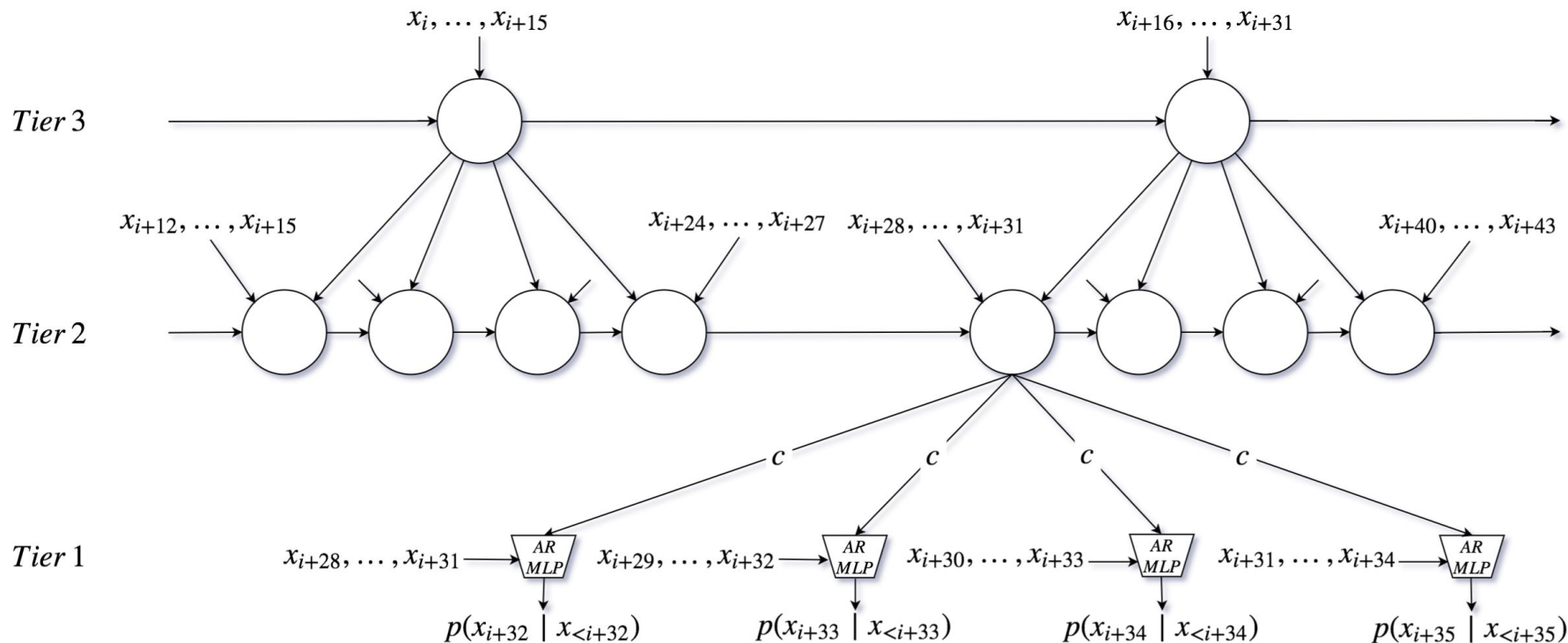


Raw

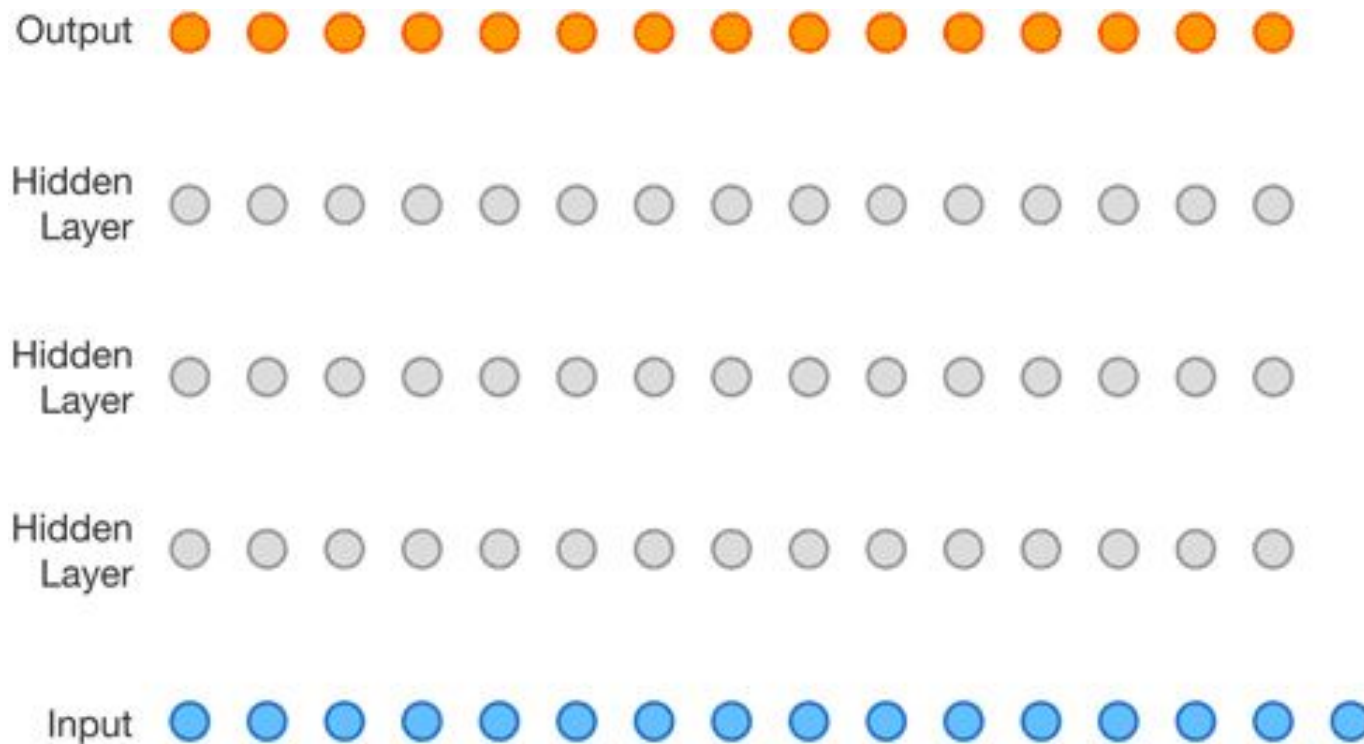


Mel spectrum

SampleRNN Speech Decoder

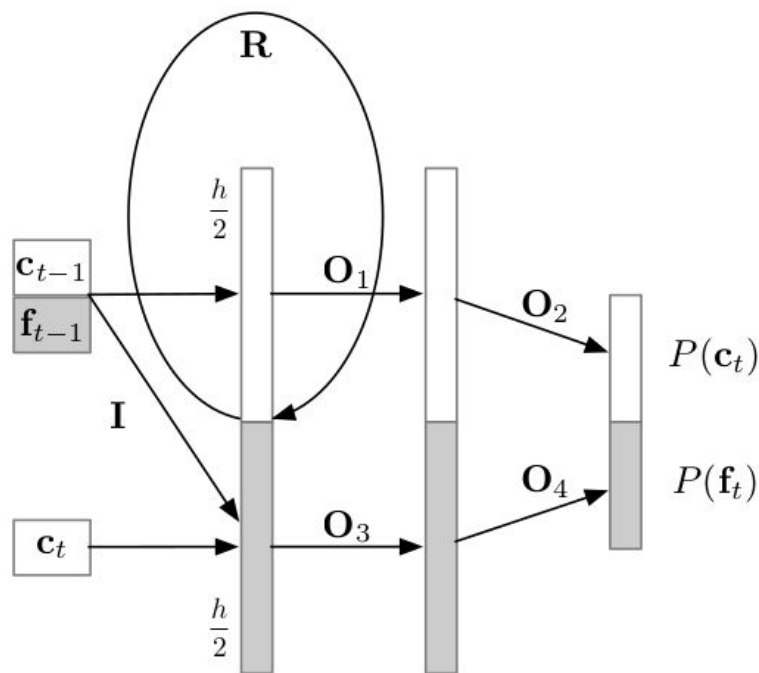


WaveNet Speech Decoder



Oord, Aaron van den, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. ["Wavenet: A generative model for raw audio."](https://arxiv.org/abs/1609.03499) arXiv preprint arXiv:1609.03499 (2016).

WaveRNN Speech Decoder



Kalchbrenner, Nal, Erich Elsen, Karen Simonyan, Seb Noury, Norman Casagrande, Edward Lockhart, Florian Stimberg, Aaron van den Oord, Sander Dieleman, and Koray Kavukcuoglu. ["Efficient Neural Audio Synthesis."](#) arXiv preprint arXiv:1802.08435 (2018).

WaveGAN Speech Decoder

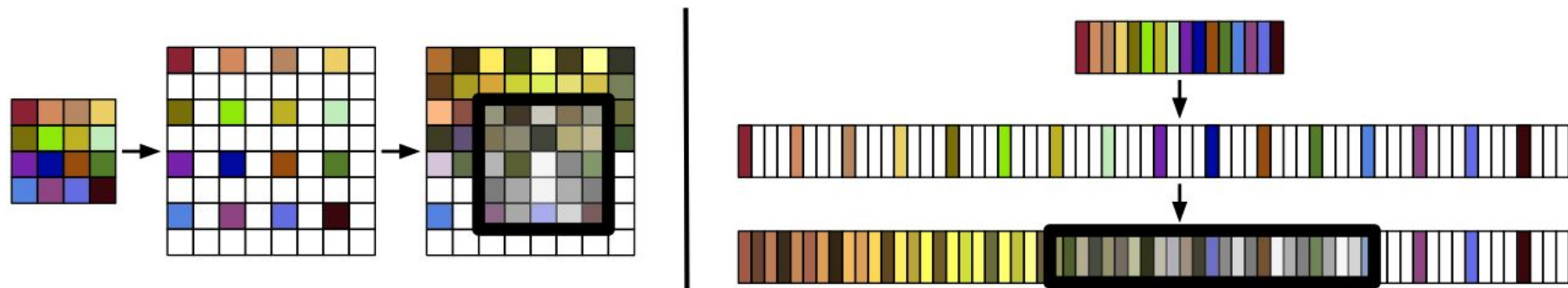
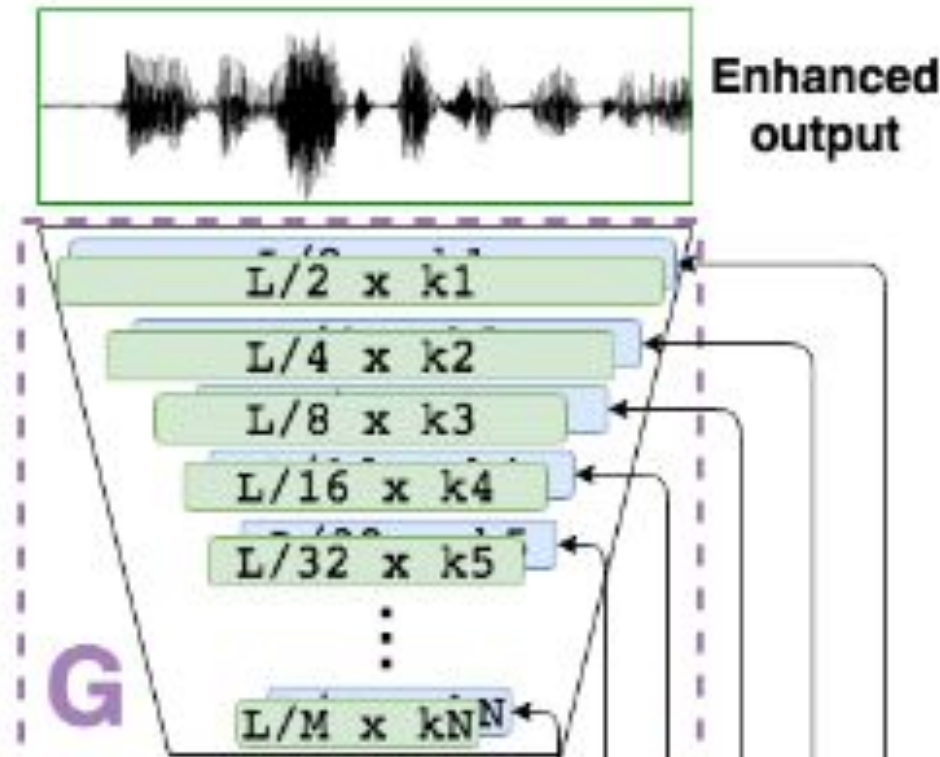


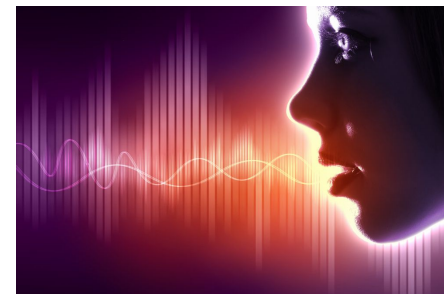
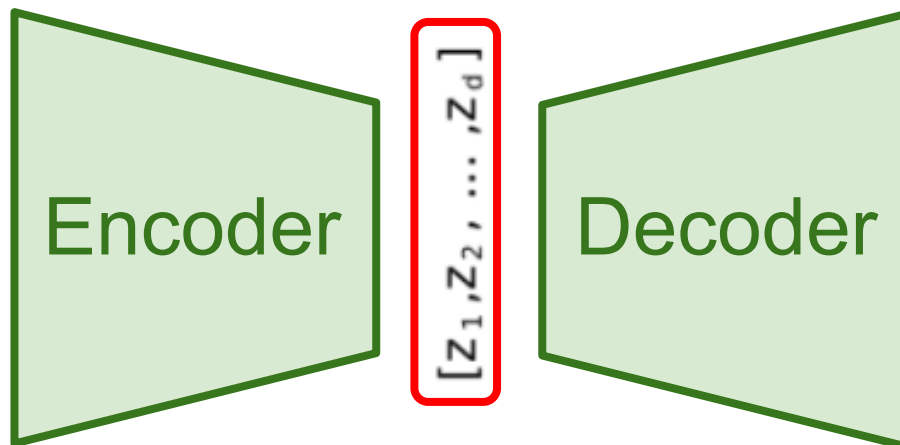
Figure 2. Depiction of the transposed convolution operation for the first layers of the DCGAN (Radford et al., 2016) (**left**) and WaveGAN (**right**) generators. DCGAN uses small (5x5), two-dimensional filters while WaveGAN uses longer (length-25), one-dimensional filters and a larger upsampling factor. The two operations have the same number of parameters and numerical operations.

SEGAN Speech Decoder



Pascual, Santiago, Antonio Bonafonte, and Joan Serra. ["SEGAN: Speech enhancement generative adversarial network."](#) Interspeech 2017.

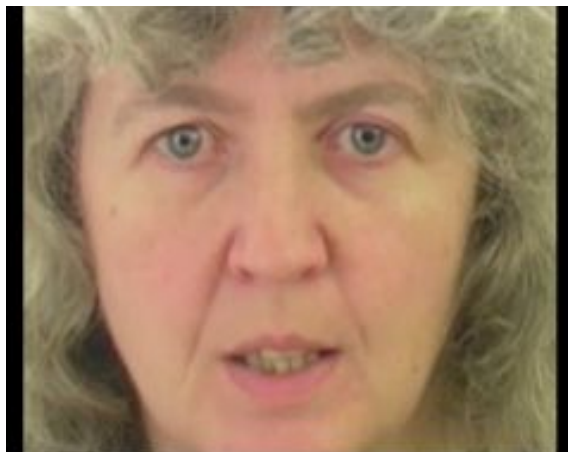
Representation or Embedding



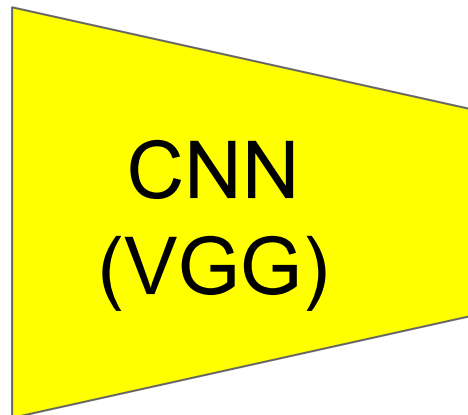


Ephrat, Ariel, Tavi Halperin, and Shmuel Peleg. "Improved speech reconstruction from silent video." In ICCV 2017 Workshop on Computer Vision for Audio-Visual Media. 2017.

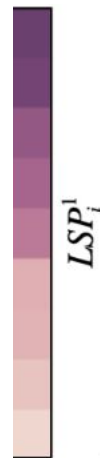
Speech Reconstruction from Video



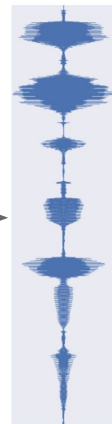
Frame from a
silent video



CNN
(VGG)



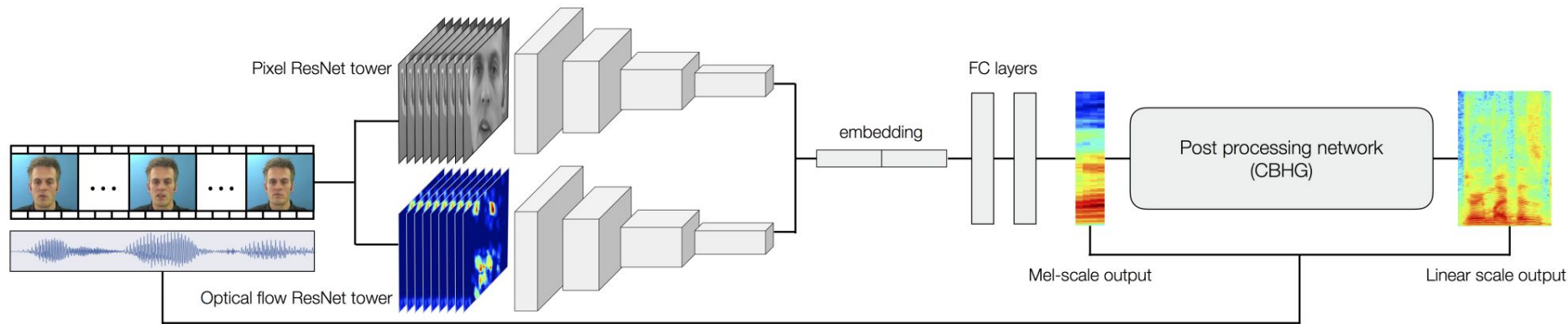
Post-hoc
synthesis



Audio feature

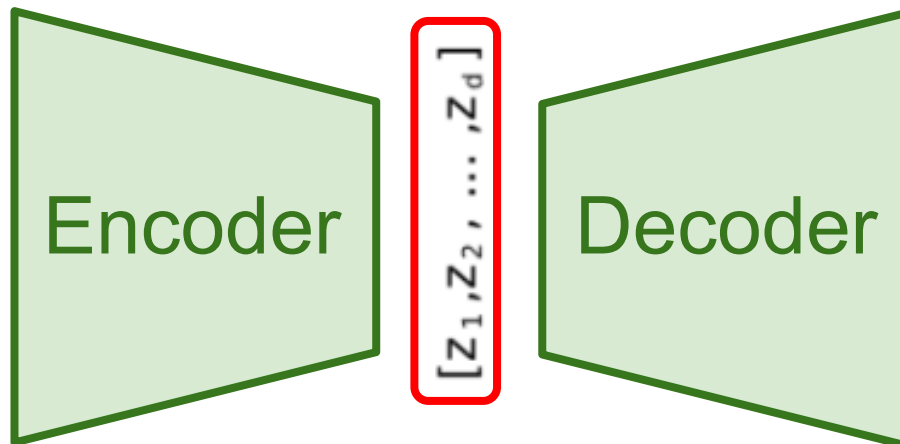
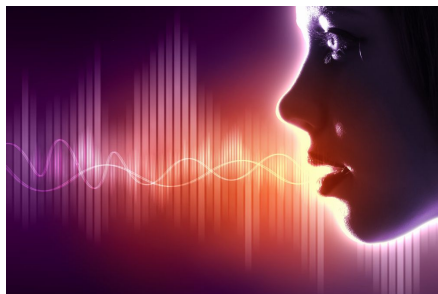
Speech Reconstruction from Video

bit.ly/DLCV2018
#DLUPC



Ephrat, Ariel, Tavi Halperin, and Shmuel Peleg. "Improved speech reconstruction from silent video." In ICCV 2017 Workshop on Computer Vision for Audio-Visual Media. 2017.

Representation or Embedding

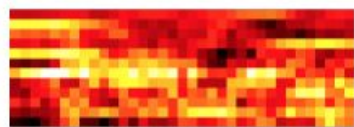


Speech to Frame Synthesis (face pixels)

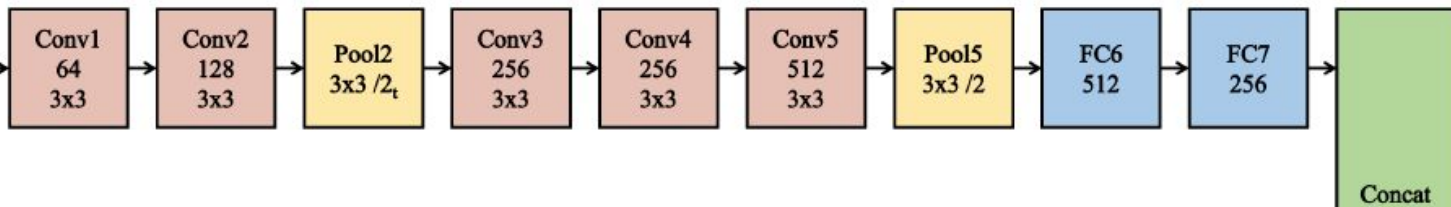
bit.ly/DLCV2018
#DLUPC



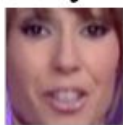
Audio Encoder



12x35x1



Identity Encoder



112x112x3

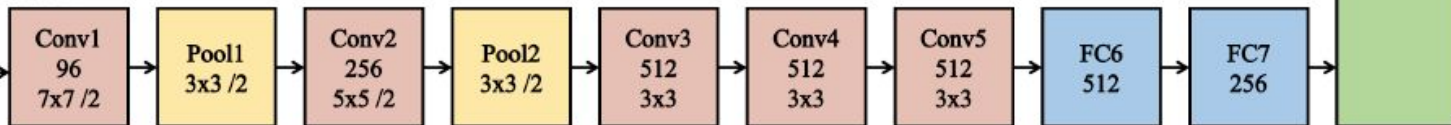
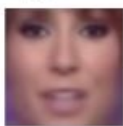


Image Decoder



109x109x3

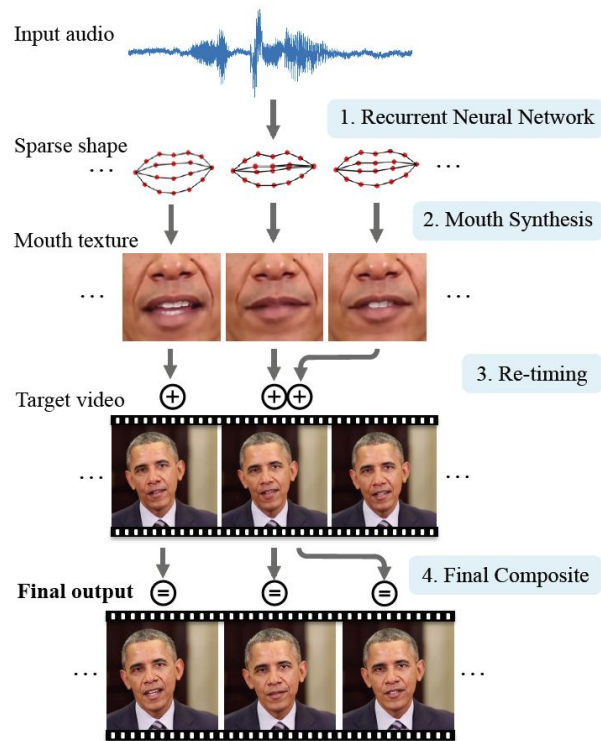




Speech to Video Synthesis (lip keypoints)

bit.ly/DLCV2018

#DLUPC



Suwajanakorn, Supasorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman. ["Synthesizing Obama: learning lip sync from audio."](#) SIGGRAPH 2017.



Without Re-timing



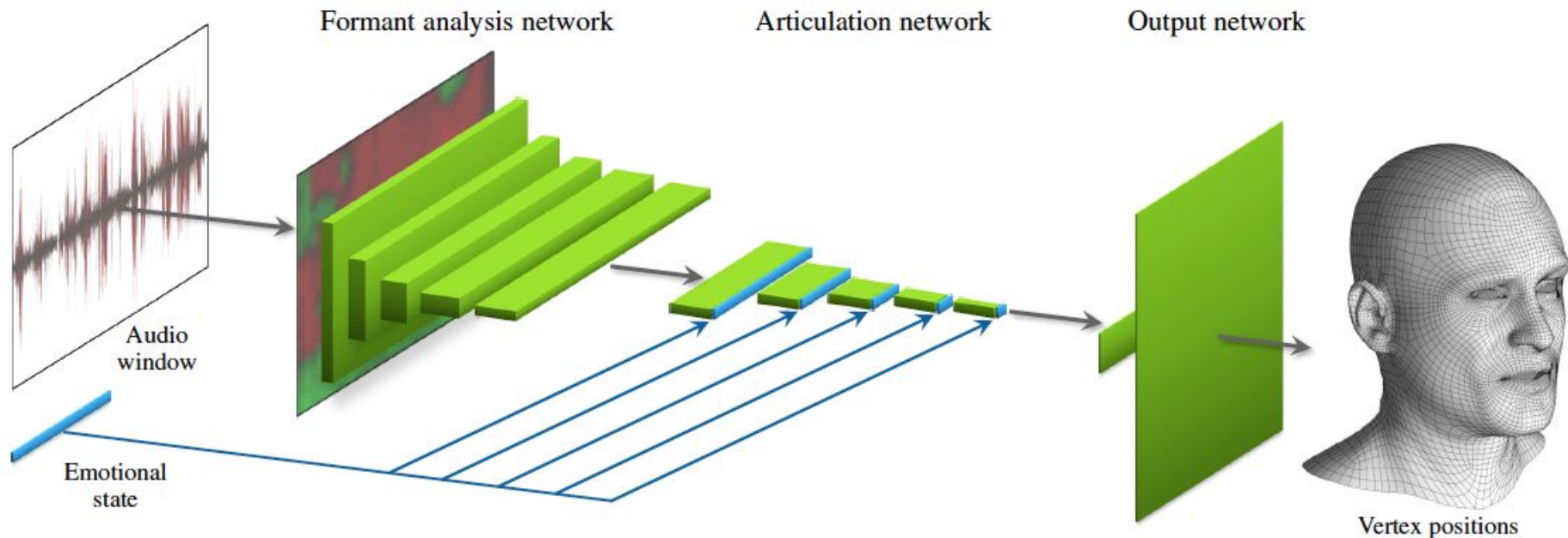
With Re-timing
(Our Result)

Karras, Tero, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. "Audio-driven facial animation by joint end-to-end learning of pose and emotion." SIGGRAPH 2017

Speech to Video Synthesis (vertex positions)

bit.ly/DLCV2018

#DLUPC



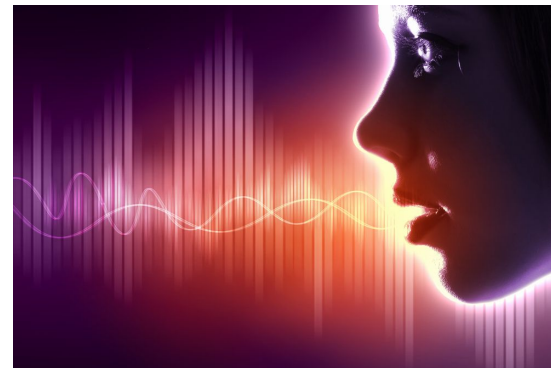
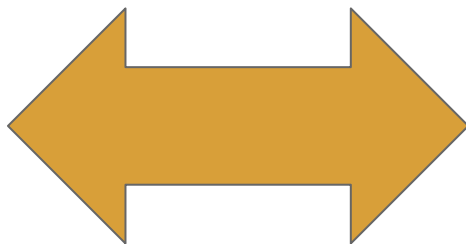
Karras, Tero, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. ["Audio-driven facial animation by joint end-to-end learning of pose and emotion."](#) SIGGRAPH 2017



Karras, Tero, Timo Aila, Samuli Laine, Antti Herva, and Jaakko Lehtinen. "Audio-driven facial animation by joint end-to-end learning of pose and emotion." SIGGRAPH 2017



Vision



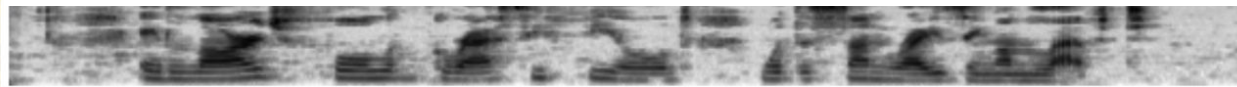
Speech

Matching speech to images



Humans understand speech much earlier than text, could computers do the same ?

Large dataset (120,000) of speech description of images from Places dataset.

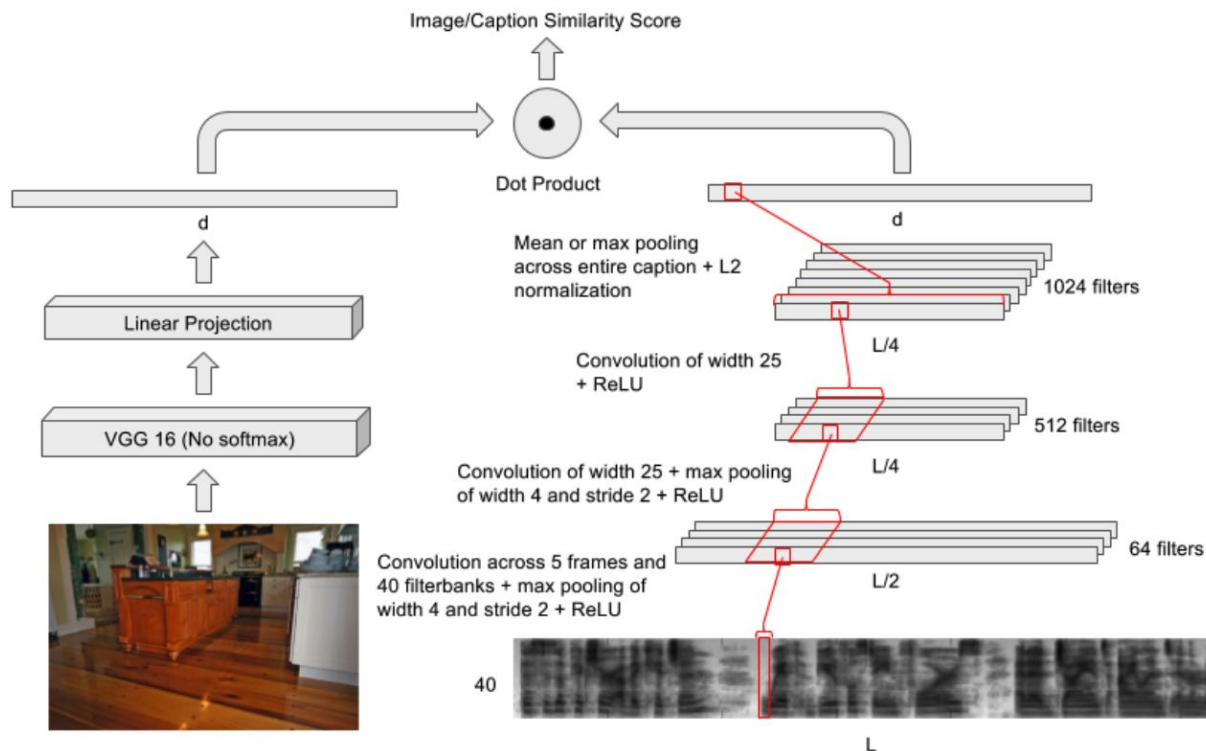


Harwath, David, **Antonio Torralba**, and James Glass. "[Unsupervised learning of spoken language with visual context.](#)" NIPS 2016. [\[talk\]](#)

Matching speech to images



Train a visual & speech networks with pairs of (non-)corresponding images & speech.



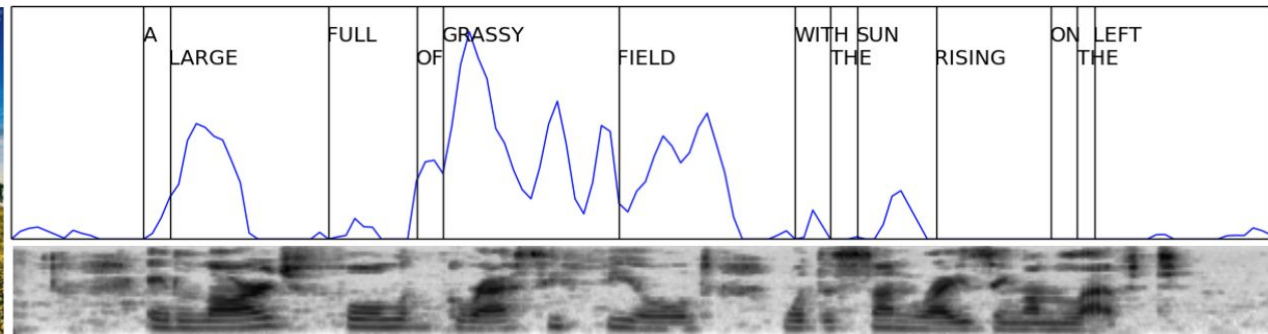
Harwath, David, Antonio Torralba, and James Glass. ["Unsupervised learning of spoken language with visual context."](#) NIPS 2016. [\[talk\]](#)

Matching speech to images



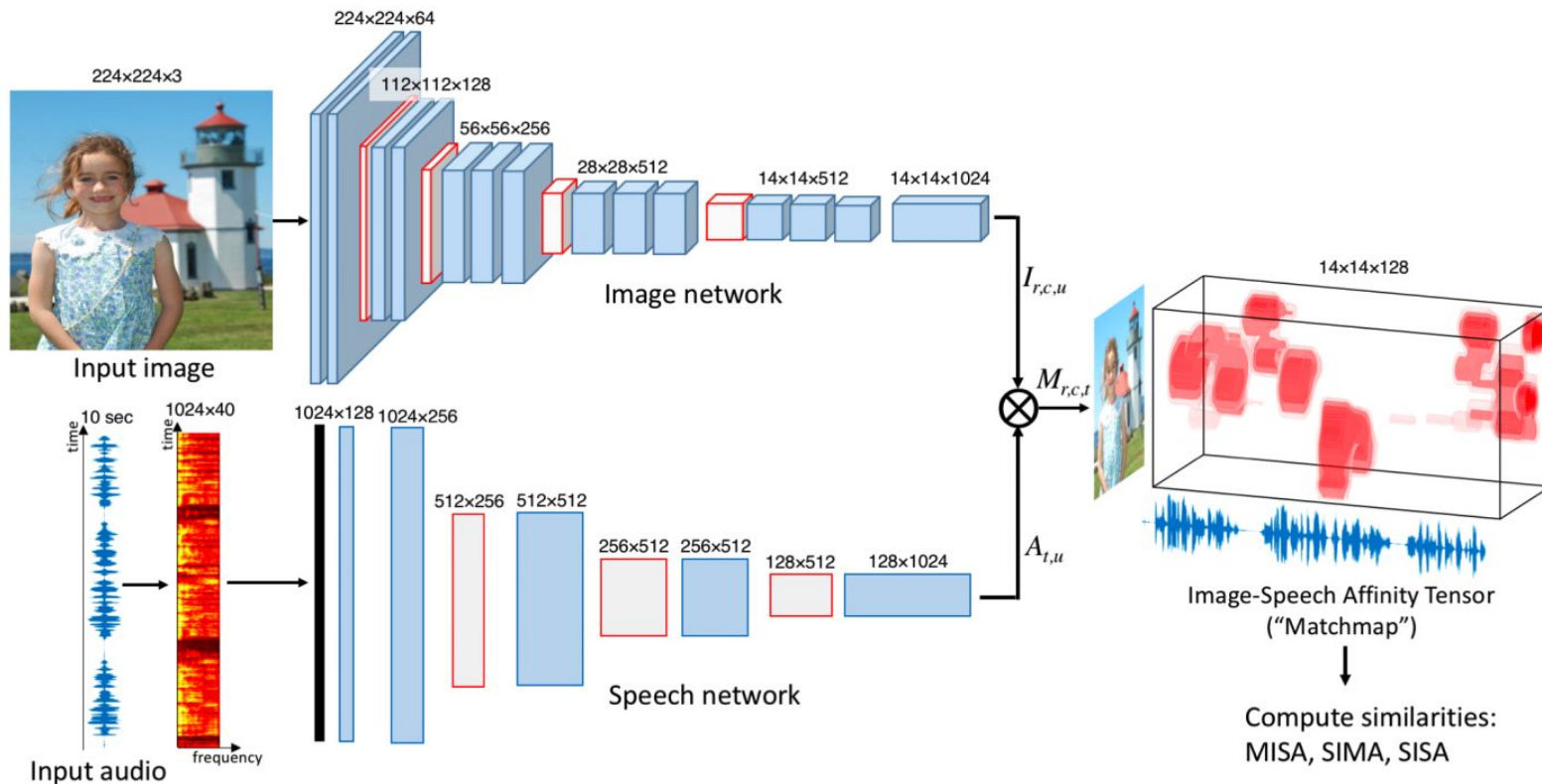
Similarity curve show which regions of the spectrogram are relevant for the image.

Important: no text transcriptions used during the training !!



Matching speech to objects (heatmap)

bit.ly/DLCV2018
#DLUPC



Harwath, David, Adrià Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass. ["Jointly Discovering Visual Objects and Spoken Words from Raw Sensory Input."](#) arXiv preprint arXiv:1804.01452 (2018).

Matching speech to objects (heatmap)

bit.ly/DLCV2018
#DLUPC

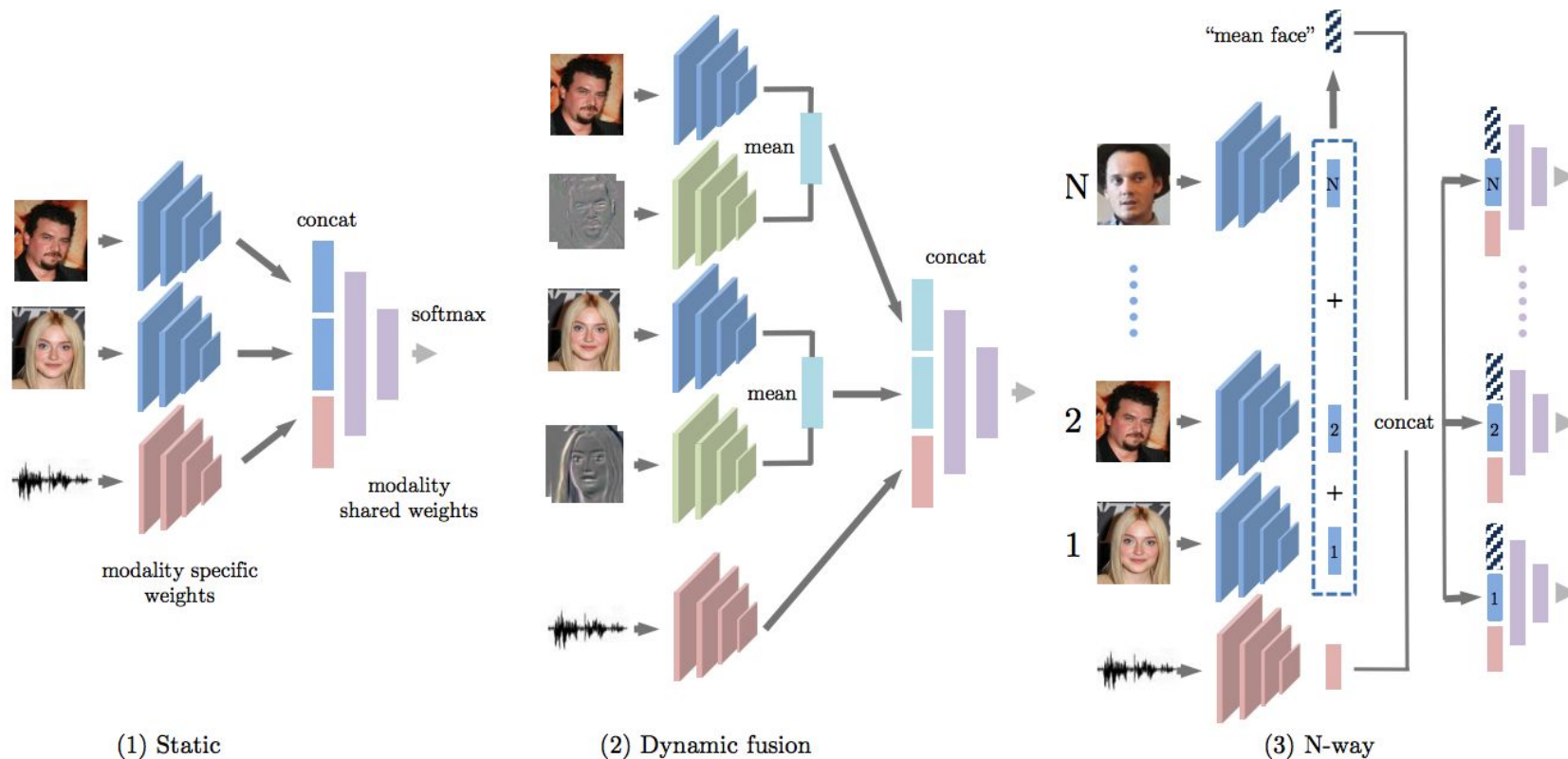


Regions matching the spoken word "WOMAN":



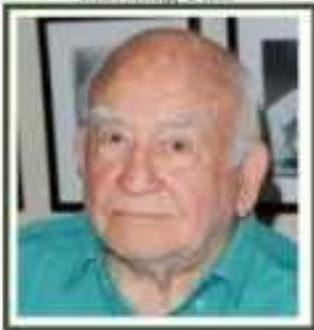
Harwath, David, Adrià Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass. ["Jointly Discovering Visual Objects and Spoken Words from Raw Sensory Input."](#) arXiv preprint arXiv:1804.01452 (2018)

Matching speech to objects (faces)





Matching Face

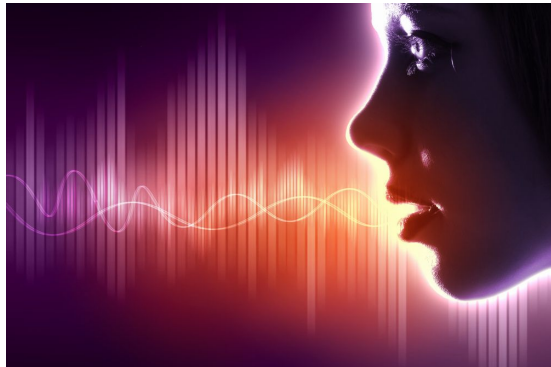


Audio Source

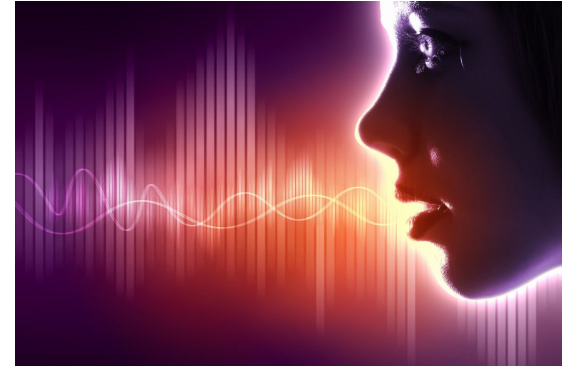
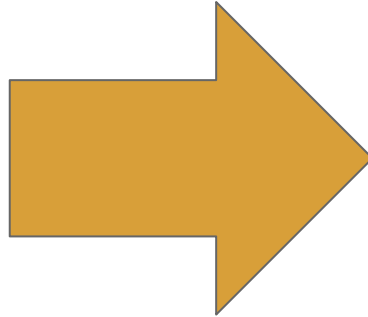




Vision

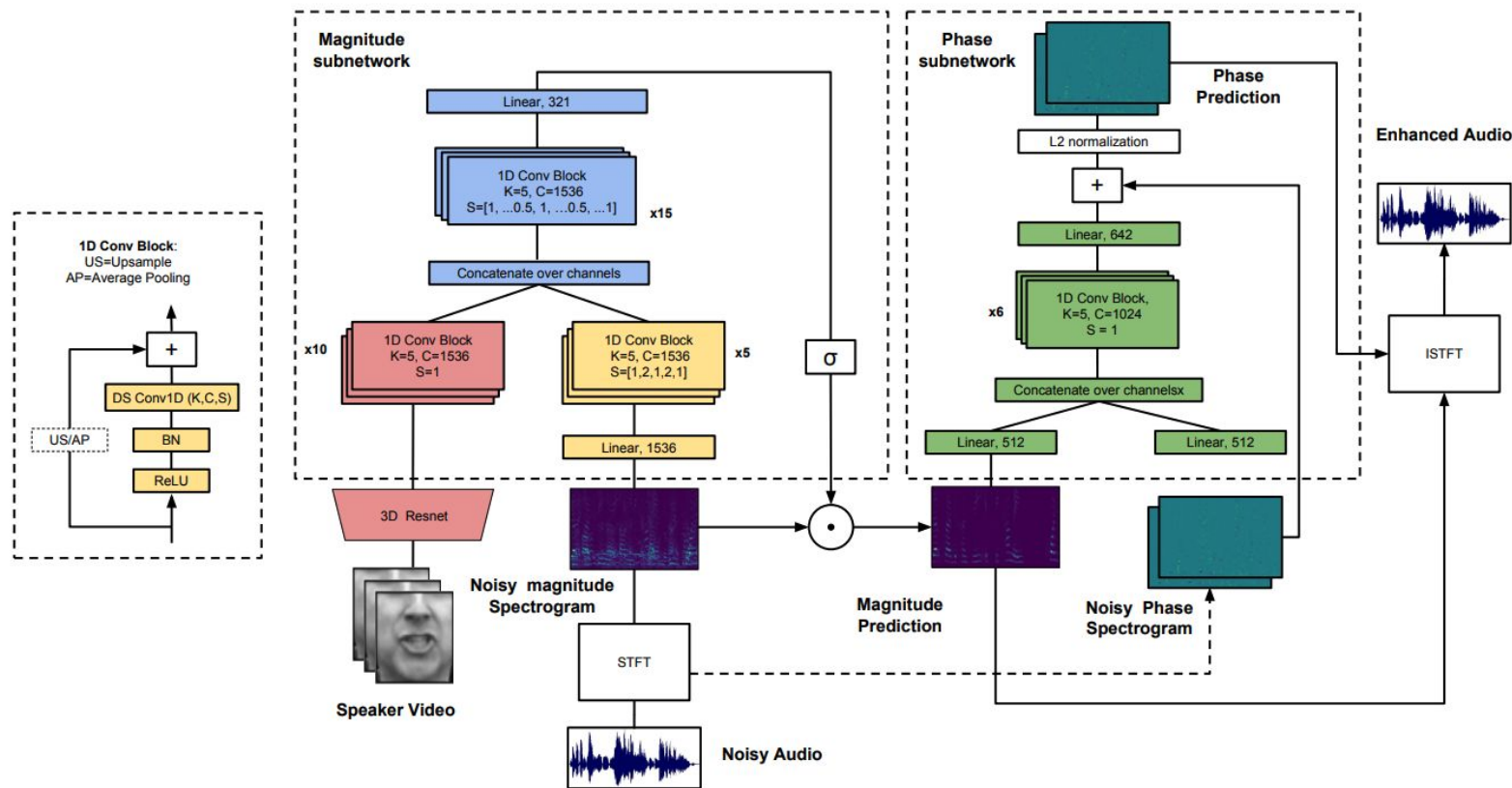


Speech



Speech

Speech Separation with Vision (lips)

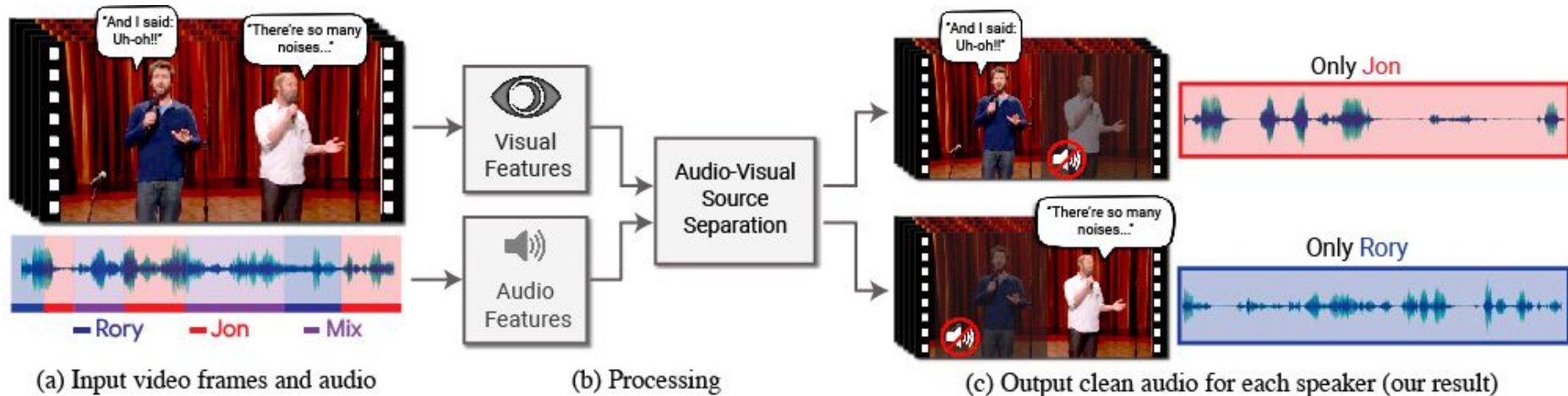




Afouras, Triantafyllos, Joon Son Chung, and Andrew Zisserman. "The Conversation: Deep Audio-Visual Speech Enhancement." Interspeech 2018..

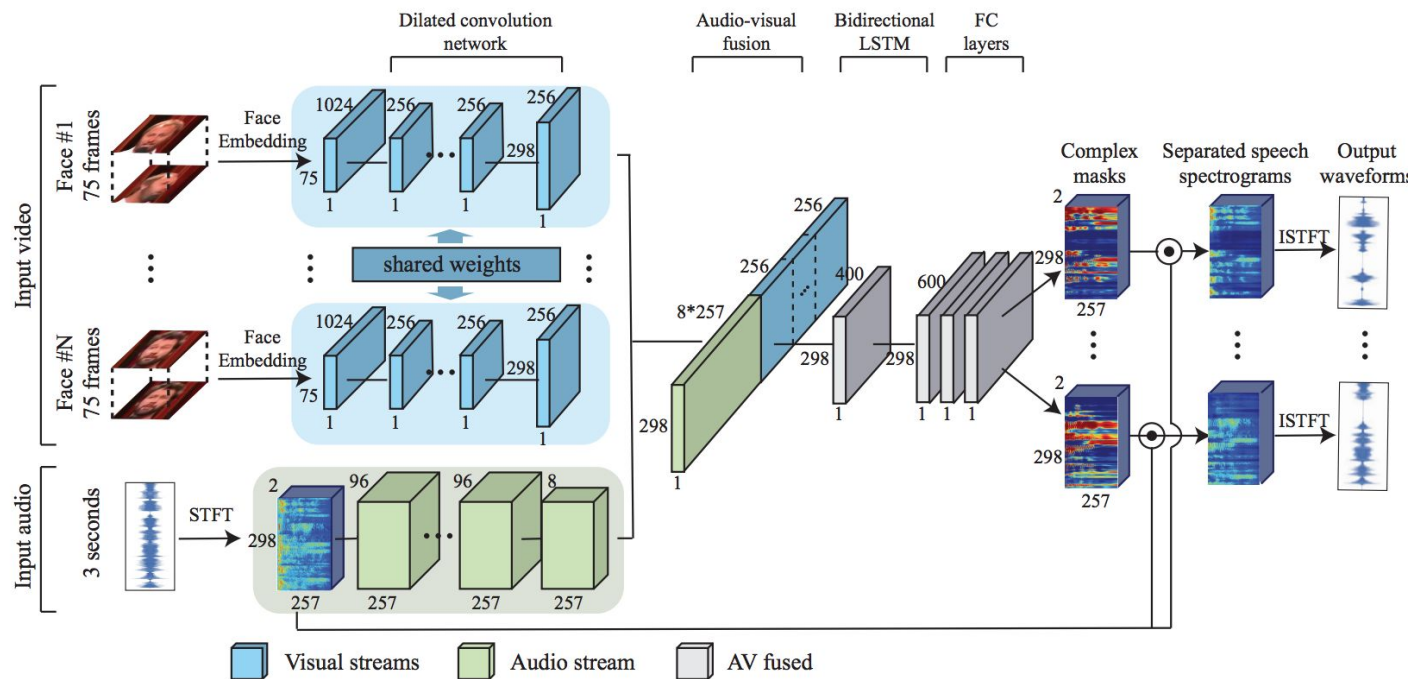
Speech Separation with Vision (faces)

bit.ly/DLCV2018
#DLUPC



Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T. Freeman and Michael Rubinstein, ["Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation"](#) SIGGRAPH 2018.

Speech Separation with Vision (faces)



Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T. Freeman and Michael Rubinstein, "[Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation](#)" SIGGRAPH 2018.



Video source: Team Coco, <https://www.youtube.com/watch?v=UT7h4nRcWJU>

Audio-Visual Speech Separation Results

Comparison with Audio-Only

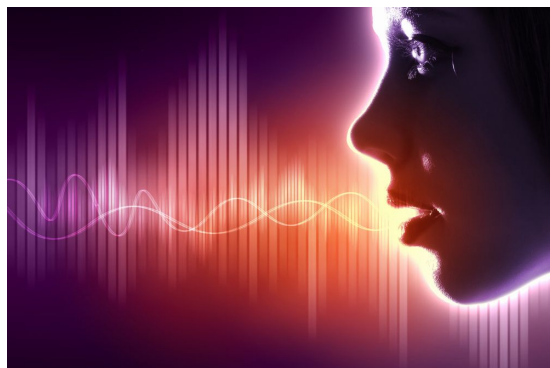
Comparison with Audio-Visual Methods

Application to Video Transcription

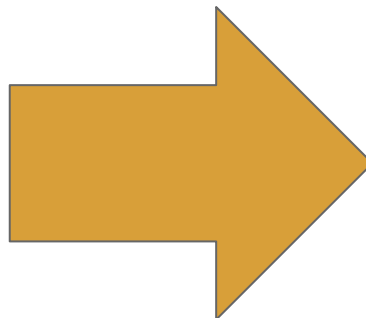
Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T. Freeman and Michael Rubinstein, "Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation" SIGGRAPH 2018.



Vision



Speech



con-
eius-
re et
n ve-
umco
con-
hen-
re eu
ccae-
ulpa
labo-
ndip-
icidi-
a. Ut
exer-
ex ea
olor
e cil-
cep-
lent,
sollit
n sit

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt

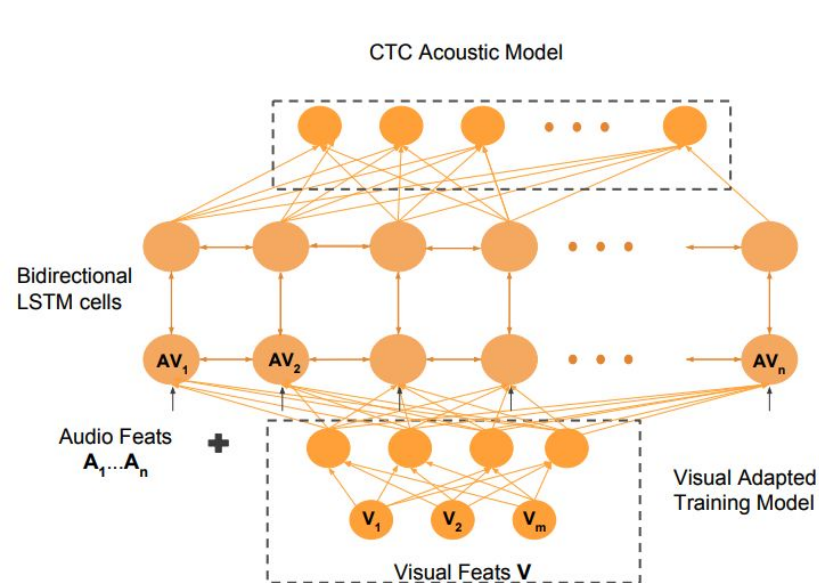
eiusmod tempor incididunt ut labore et do-
lore magna aliqua. Ut enim ad minim ve-
niam, quis nostrud exercitation ullamco
laboris nisi ut aliquip ex ea commodo con-
sequat. Duis aute irure dolor in reprehen-
derit in voluptate velit esse cillum dolore eu
fugiat nulla pariatur. Excepteur sint occae-
cat cupidatat non proident, sunt in culpa
qui officia deserunt mollit anim id est labo-
rum. Lorem ipsum sit amet, consectetur
adipiscing elit, sed do eiusmod tempor in-
cididunt ut labore et dolore magna aliqua.
Ut enim ad minim veniam, quis nostrud ex-
ercitation ullamco laboris nisi ut aliquip ex
ea commodo consequat. Duis aute irure

rept
sont
niffo
con-
was
blat
whi
wait
the
sint
in c
est
sect
tem
na z
nos
aliq
aut
lupt
null
con
tem
na a

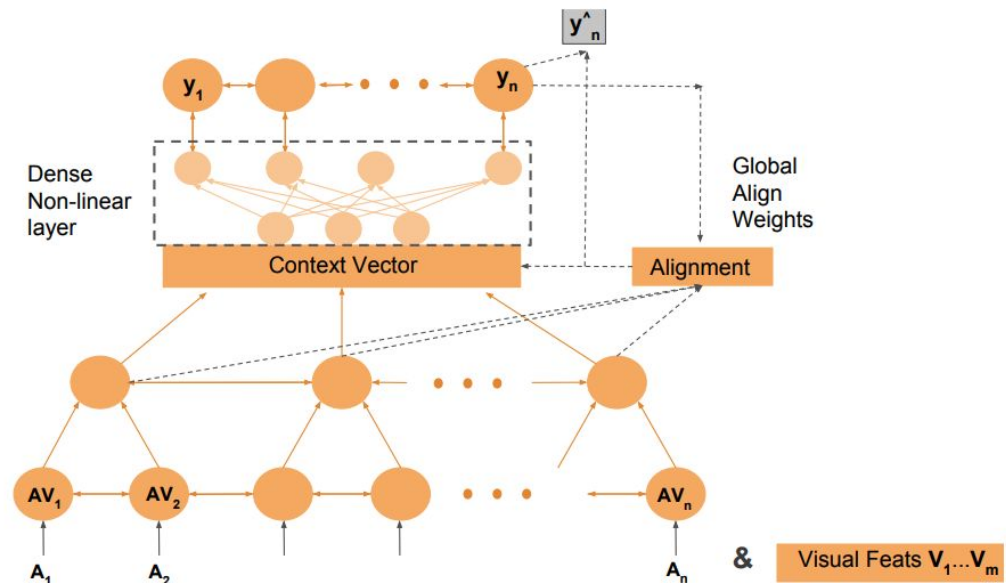
aliquip ex ea commodo consequat.
Duis aute irure dolor in reprehenderit in
voluptate velit esse cillum dolore eu fugiat
nulla pariatur. Excepteur sint occaecat cup-
idatat non proident, sunt in culpa qui offi-
cia deserunt mollit anim id est laborum.
Lorem ipsum sit amet, consectetur adipi-
scing elit, sed do eiusmod tempor incidi-
dunt ut labore et dolore magna aliqua. Ut
enim ad minim veniam, quis nostrud exer-
citation ullamco laboris nisi ut aliquip ex ea
commodo consequat. Duis aute irure dolor
in reprehenderit in voluptate velit esse cil-
lum dolore eu fugiat nulla pariatur. Lorem
ipsum dolor sit amet, consectetur adiniscine

Text

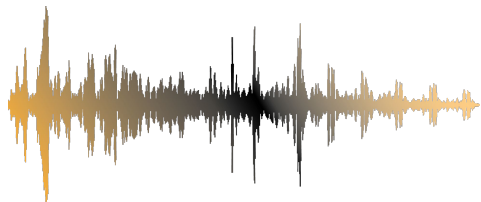
Speech Recognition with vision



(a) CTC model architecture with adaptation



(b) S2S model architecture with global attention and adaptation



Audio



Vision

con-
etis-
re et
nive-
unco
con-
hen-
re eu
cuae-
culpa
labo-

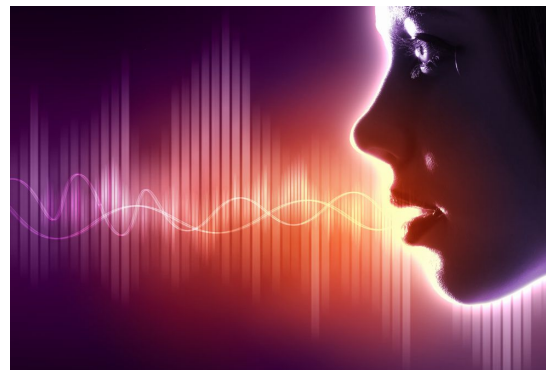
Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt

eiusmod tempor incididunt ut labore et do-
lore magna aliqua. Ut enim ad minim ve-
niam, quis nostrud exercitation ullamco
laboris nisi ut aliquip ex ea commodo con-
sequat. Duis aute irure dolor in reprehen-
derit in voluptate velit esse cillum dolore eu
fugiat nulla pariatur. Excepteur sint occae-
cat cupidatat non proident, sunt in culpa
qui officia deserunt mollit anim id est labo-
rum. Lorem ipsum sit amet, consectetur
adipiscing elit, sed do eiusmod tempor in-
cididunt ut labore et dolore magna aliqua.
Ut enim ad minim veniam, quis nostrud ex-
ercitation ullamco laboris nisi ut aliquip ex
ea commodo consequat. Duis aute irure

repi-
son-
nifi-
cou-
was
blar
whi
wait
the
sint
in c
est
sect
tem
na i
nos
aliq
aut
lupt
null
con
tem
na a

aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum. Lorem ipsum sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Lorem ipsum dolor sit amet, consectetur adipisicing

Text



Speech



Vision

Speech

The diagram illustrates a neural network architecture for multimodal learning. It features three input modalities: Audio, Text, and Vision. Each modality is processed by a set of weights ($w_0, w_1, w_2, \dots, w_n$) and an input vector (x_1, x_2, \dots, x_n). The weighted inputs are summed (Σ) and passed through a sigmoid function to produce the final output.

eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Lorem ipsum sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Lorem ipsum sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

reputations of the
son of a famous
nobleman
could be
was
blame
which
wait
the
since
in c
est
sect
tem
na a
nos
alig
auto
fup
nul
con
tem
na

Speech



Undergradese

What undergrads ask vs. what they're REALLY asking

"Is it going to be an open book exam?"

Translation: "I don't have to actually memorize anything, do I?"

"Hmm, what do you mean by that?"

Translation: "What's the answer so we can all go home."

"Are you going to have office hours today?"

Translation: "Can I do my homework in your office?"

"Can i get an extension?"

Translation: "Can you re-arrange your life around mine?"

"Is this going to be on the test?"

Translation: "Tell us what's going to be on the test."

"Is grading going to be curved?"

Translation: "Can I do a mediocre job and still get an A?"

JORGE CHAM © 2008

