

# FAKE NEWS

## Current State of the art in Disinformation Generation and Detection



**Alexander Schindler**

Scientist

Information Management

Center for Digital Safety & Security

AIT Austrian Institute of Technology GmbH

**Mina Schütz**

Intern

Information Management

Center for Digital Safety & Security

AIT Austrian Institute of Technology GmbH



# Agenda

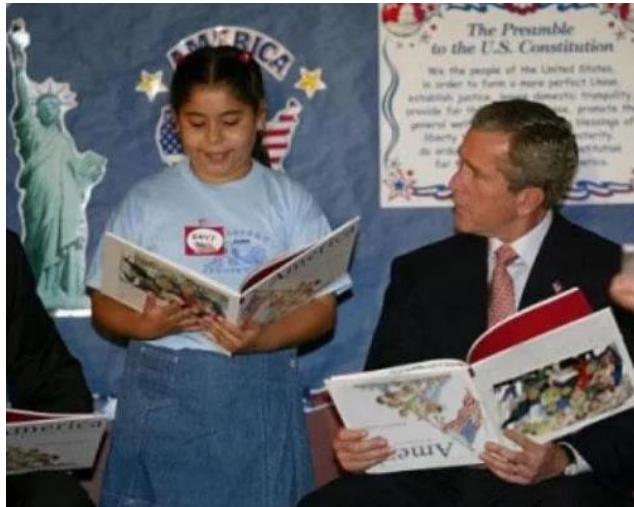
- Introduction to Fake News
- Generative Processes
  - GAN / VAE
  - Sequence to Sequence
- Fake News Detection approaches
  - Visual
  - Accoustic
  - Textual
- Discussion and Conclusion

# Who shares Fake News?

- Do you / have you shared **Fake News**?
- Are you likely to share **Fake News**?
- Can you identify **Fake News**?

# George Bush Reading Upside-Down

- George W. Bush (US Pres.)
- Rick Perry (Gov. Texas)
- George Sanchez Charter School (Houston)
- 2002



Fake Image



Original Image



Flipping Error

26.11.2019



Additional Proof

# Ideological Priors

or

Why do we fall for Fake News?

- **Naive Realism**

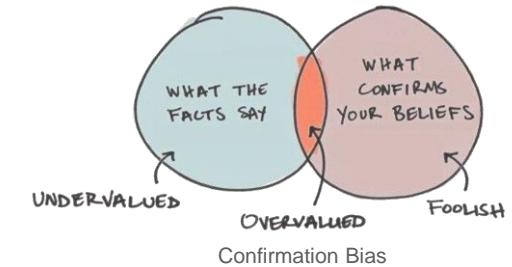
- Believe information that is aligned with your views

- **Confirmation Bias**

- Seek information that confirms your existing views

- **Normative Influence Theory**

- Consume/Share socially safe options → for social acceptance, affirmation



→ Individual level of Fake News



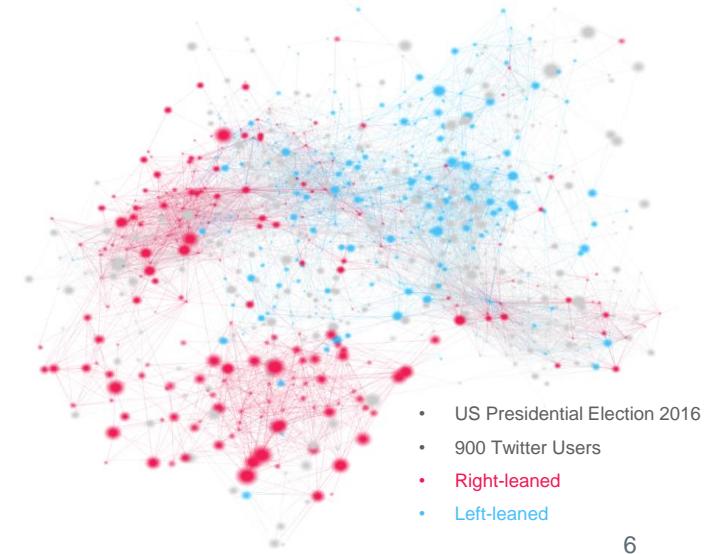
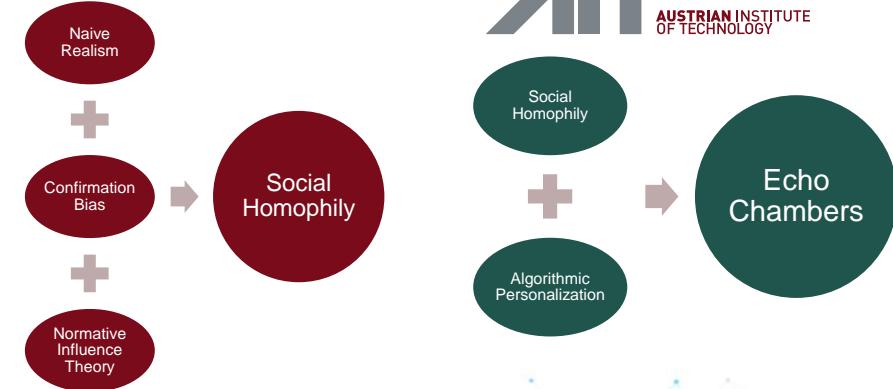
Teilung im Nationalrat



EU-Ausgleich für Massenzug

# Nature / Characteristics

- **Social Level**
- **Echo Chambers / Filter Bubbles**
  - Social homophily
    - Form connections with ideologically similar individuals
  - Algorithmic personalization
    - Read content
    - Follow / befriend persons
- Consequences
  - Less exposure to conflicting viewpoints
  - Isolation in own filter bubble
  - Improve survival / spread of fake news

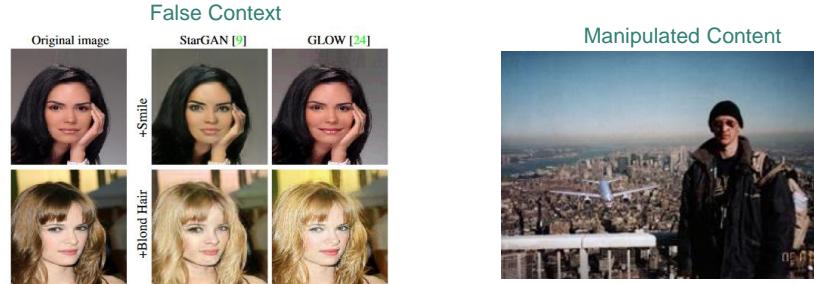


# Types of Fake News

- **Fabricated content**
  - Completely false
- **Misleading content**
  - Misleading use / framing of issue
- **Imposter content**
  - Genuine source impersonated with false sources
- **Manipulated content**
  - For deception (e.g. images)
- **False connection**
  - Headlines, visuals do not support content
- **False context**
  - Genuine content shared with false context information

Pope Francis Shocks World, Endorses Donald Trump for President, Releases Statement

TOPICS: Pope Francis Endorses Donald Trump



Cozzolino, D., Thies, J., Rössler, A., Riess, C., Nießner, M., & Verdoliva, L. (2018). Forensictransfer: Weakly-supervised domain adaptation for forgery detection. arXiv preprint arXiv:1812.02510.





VOLUME 28 ISSUE 17

NUMBER ONE IN NEWS

12-16 DECEMBER 1998

## Congress Hires Drummer

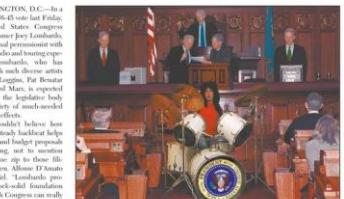


See Onion, page 28

### THE DOW

World Sport Technology Entertainment Style Travel Money

Hoax



### Bang the Drum Slowly

Here's where the absence of a skilled percussionist

## Hoax



Singapore dismisses Lee Kuan Yew death report as hoax  
By Jason Hanna, CNN  
Updated 1458 GMT (2258 HKT) March 18, 2015



**CNN** A top government spokesman dismissed as a hoax. We Singapore's founding father had died.

Former Prime Minister Lee Kuan Yew is alive, said Farah Rahim, Ministry of Communications and Information. The 91-year-old is

Lee Kuan Yew, Singapore's founding father and first prime minister, dies at 91, government website says.



Singapore dismisses Lee Kuan Yew death report as hoax  
<https://edition.cnn.com/2015/03/18/world/singapore-lee-kuan-yew/index.html>

# Sepcial Types of Fake News

- **Satire**

- *Excluded from definition*

- **Hoax**

- *False story – masquerade truth*

- **Rumors**

- *Unverified sources*
- *Not nescessarily false*
- *May be verified later as true / false*

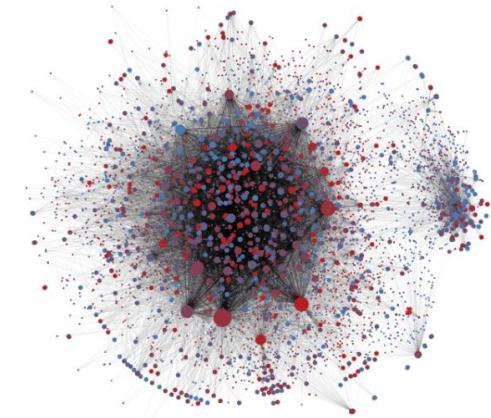
# FAKE NEWS DETECTION



# Primary Characteristics

- **Source / Promoters**
  - Who posts fake news? Who shares it?
  - Bots
- **Information Content**
  - Content
  - Linguistic style
- **User Responses**
  - Reactions to articles
  - Positive/Negative/Neutral

Bots spreading false information in social media



Kumar, S., & Shah, N. (2018). False information on web and social media: A survey. arXiv preprint arXiv:1804.08559.

# Challenges

- High stakes and multiple players
- Adversarial intent
- Lack of awareness
- Propagation dynamics
- Constant change

# General Remarks

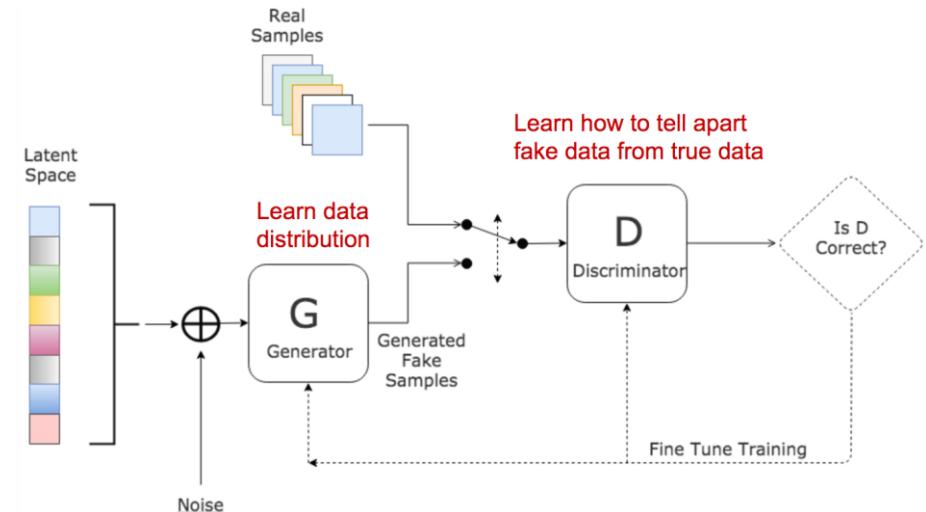
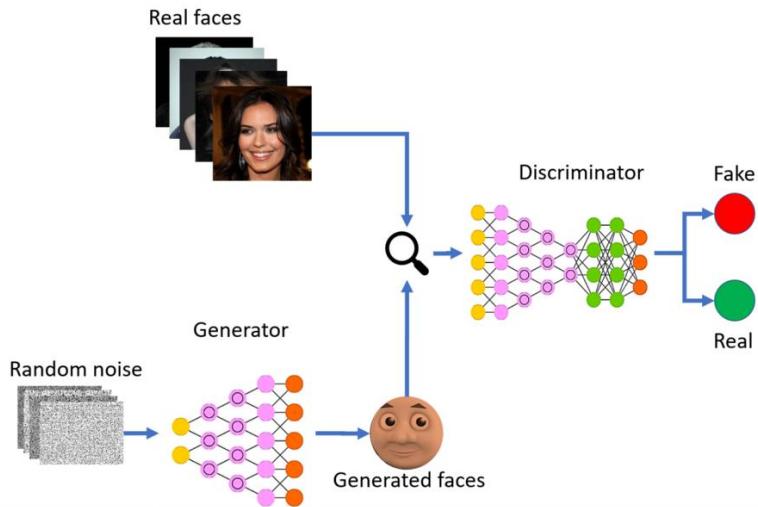
- No general Fake News Detection Tool available
- Modality dependant approaches
  - Text Domain
  - Visual Domain
  - Accoustic Domain
  - Network Domain
- Hybrid / Multi-Modal approaches

# GENERATIVE PROCESSES

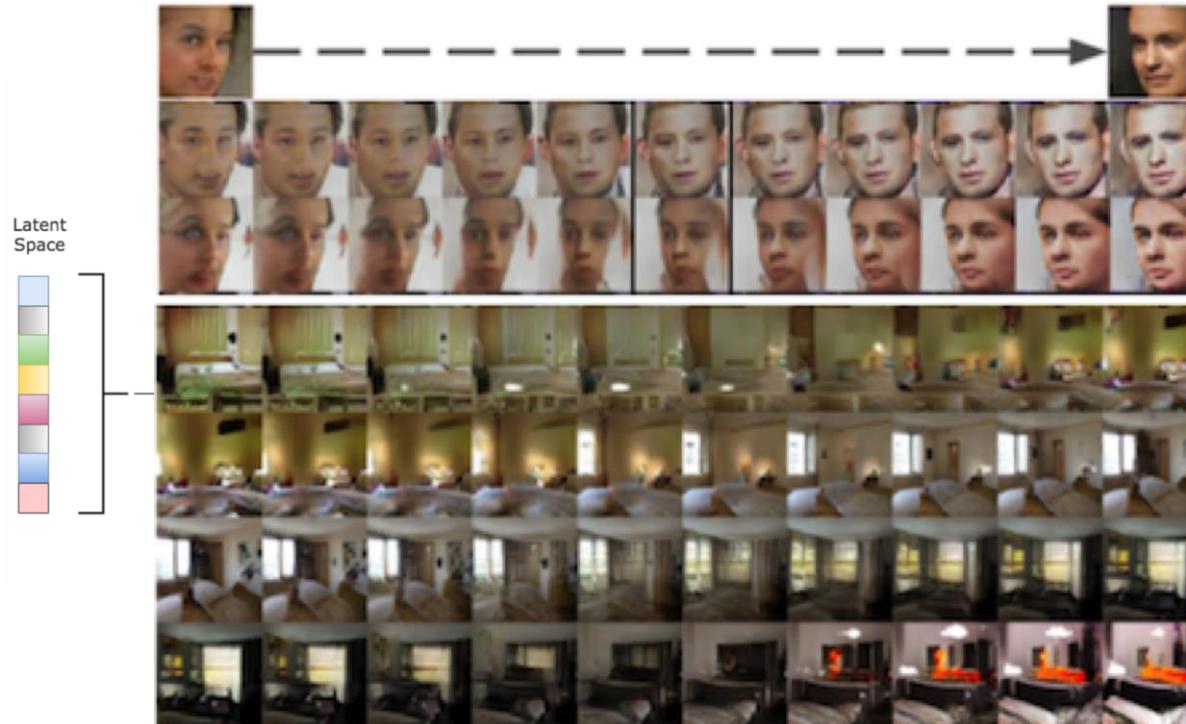
## Content Synthesization and Manipulation



# Generative Adversarial Networks (GAN)



# DCGAN: Interpolation in Latent Space



# DCGAN: Vector Arithmetic for Visual Concepts



**Generative Adversarial Networks:** One network transforms random noise into images, a second one tries to distinguish generated from real images. Both are trained at the same time.



Photographs of bed rooms that do not actually exist

# PROGRESSIVE GROWING OF GANs FOR IMPROVED QUALITY, STABILITY, AND VARIATION

Tero Karras

NVIDIA

Timo Aila

NVIDIA

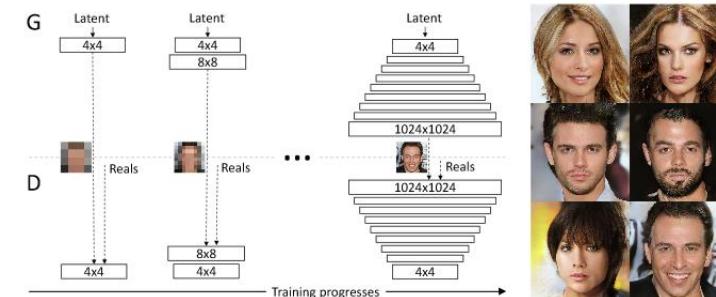
Samuli Laine

NVIDIA

Jaakko Lehtinen

NVIDIA and Aalto University

{tkarras, taila, slaine, jlehtinen}@nvidia.com



# A Style-Based Generator Architecture for Generative Adversarial Networks

Tero Karras  
 NVIDIA

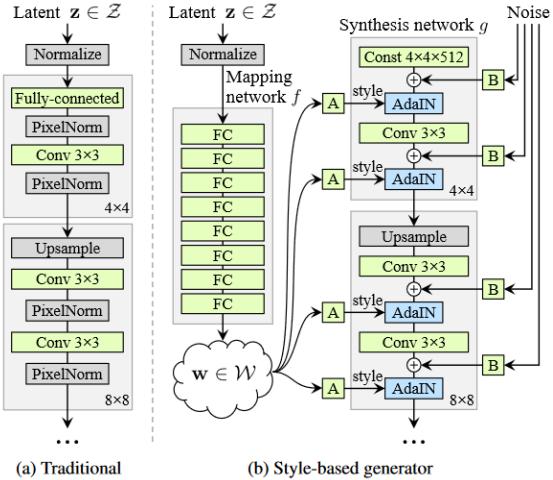
tkarras@nvidia.com

Samuli Laine  
 NVIDIA

slaine@nvidia.com

Timo Aila  
 NVIDIA

taila@nvidia.com



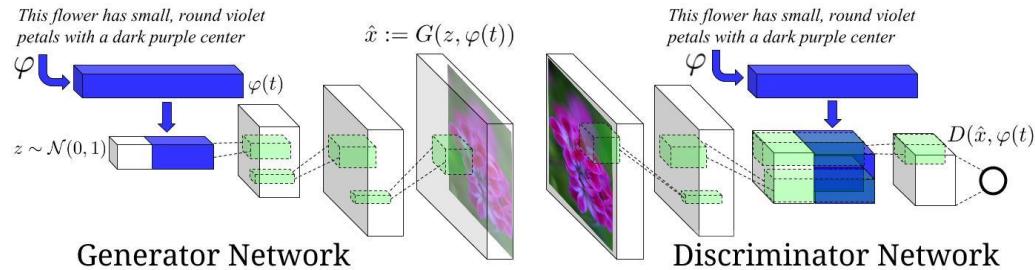
# Generative Adversarial Text to Image Synthesis

Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran  
Bernt Schiele, Honglak Lee

REEDSCOT<sup>1</sup>, AKATA<sup>2</sup>, XCYAN<sup>1</sup>, LLAJAN<sup>1</sup>  
SCHIELE<sup>2</sup>, HONGLAK<sup>1</sup>

<sup>1</sup> University of Michigan, Ann Arbor, MI, USA (UMICH.EDU)

<sup>2</sup> Max Planck Institute for Informatics, Saarbrücken, Germany (MPI-INF.MPG.DE)



this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen

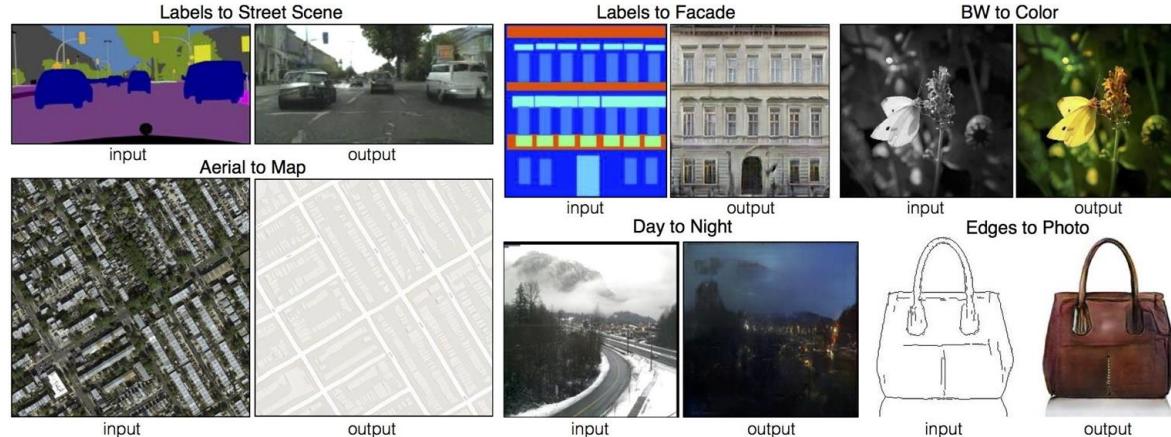


- Paper: <http://arxiv.org/abs/1605.05396>
- Github: <https://github.com/paarthneekhara/text-to-image>

Figure 1. Examples of generated images from text descriptions. Left: captions are from zero-shot (held out) categories, unseen text. Right: captions are from the training set.

# CycleGAN

- Based on pix2pix: Fully-convolutional architecture to transform an image into another



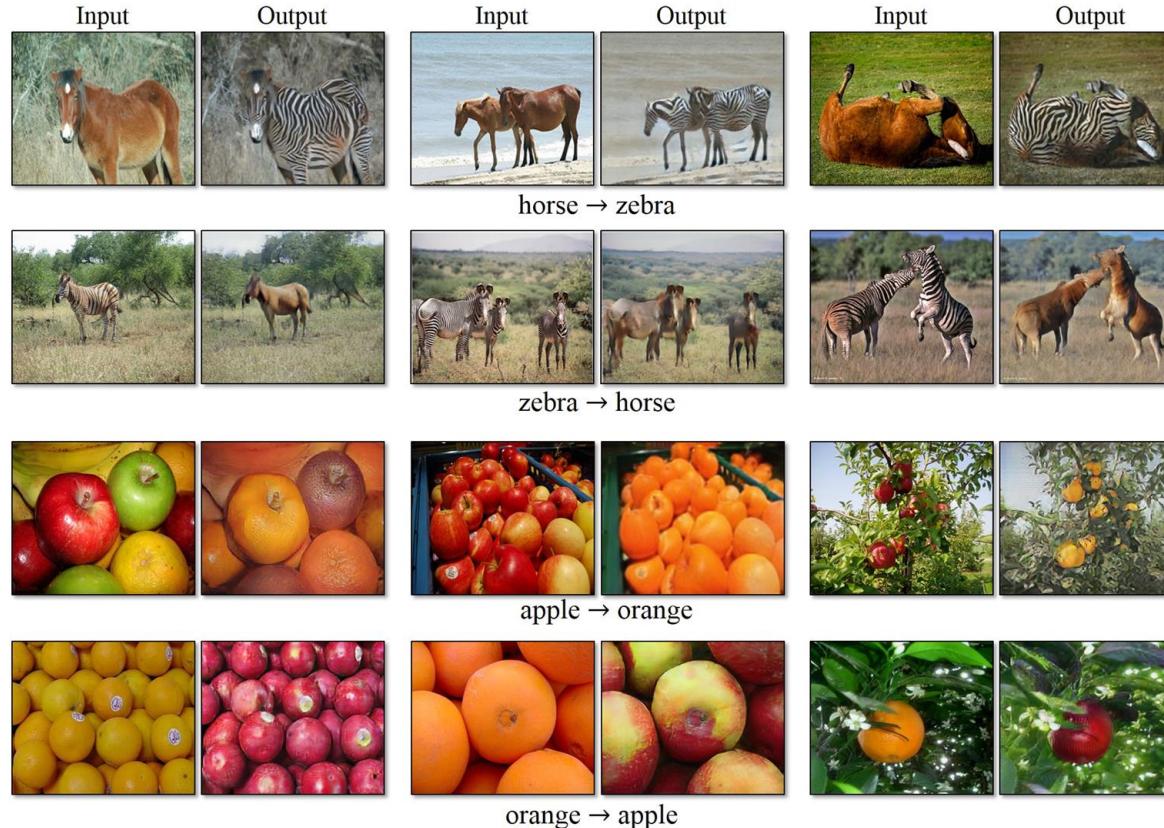
- Interactive demo: <https://affinelayer.com/pixsrv/>

# Image super-resolution through deep learning

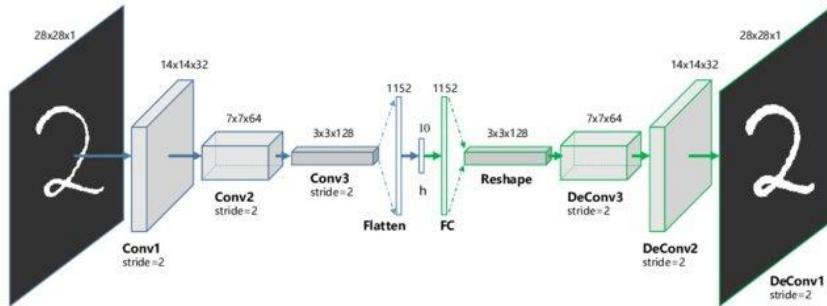
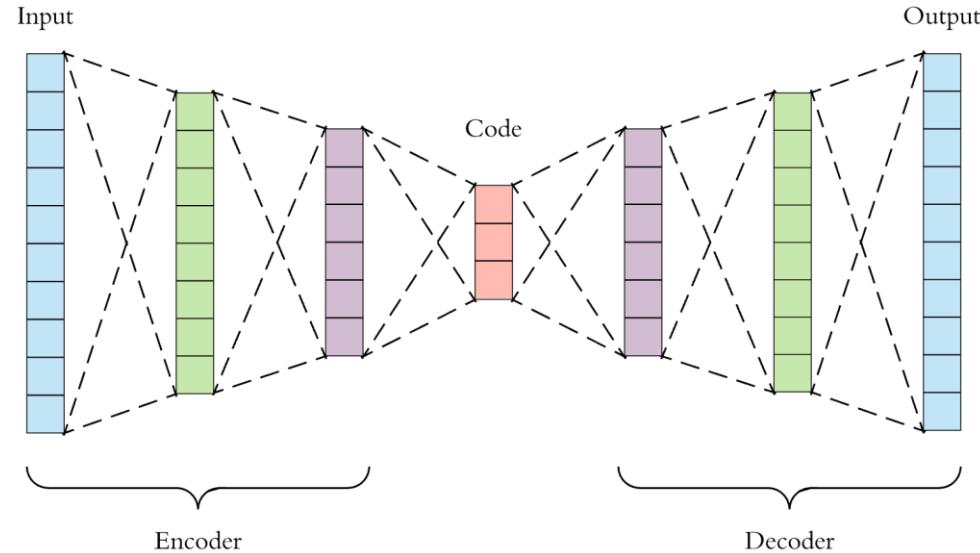
- Upscale 16x16 images by factor 4
- Deep Convolutional Generative Adversarial Network (DCGAN)
  - Image as input instead of gaussian noise
  - Loss function measures difference between input and scaled version
  - Generator uses Residual Network (ResNet) modules
- Code on Github
  - <https://github.com/david-gpu/srez>



# CycleGAN

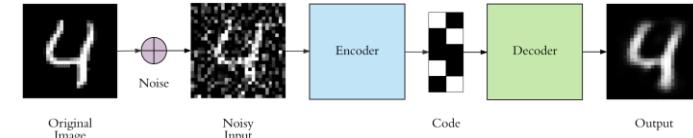


# Auto Encoders



Guo, Xifeng, et al. "Deep clustering with convolutional autoencoders." International Conference on Neural Information Processing. Springer, Cham, 2017.

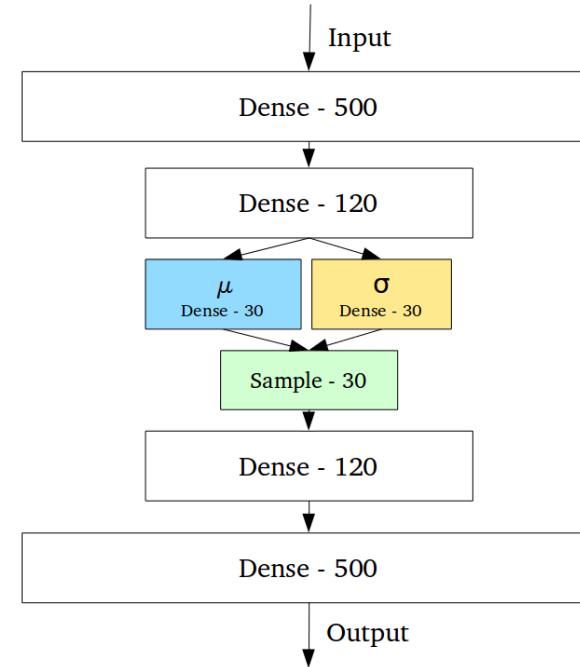
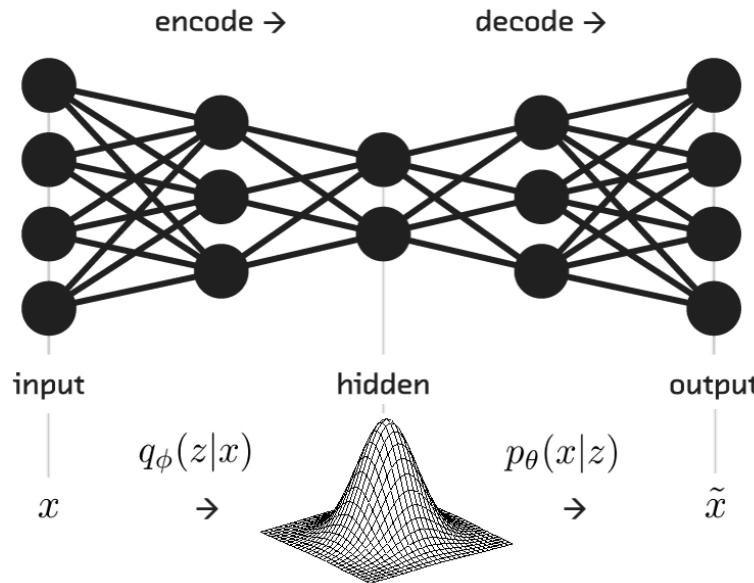
26.11.2019



Applied Deep Learning - Part 3: Autoencoders  
<https://towardsdatascience.com/applied-deep-learning-part-3-autoencoders-1c083af4d798>

24

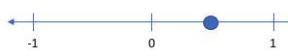
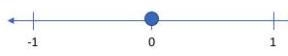
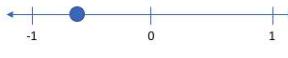
# Variational Auto Encoders



# Probability Distribution of Latent Factors

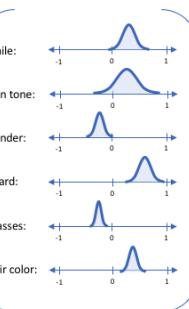
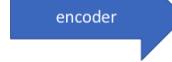
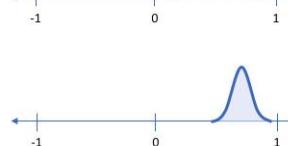
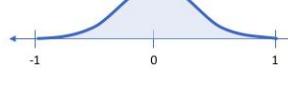
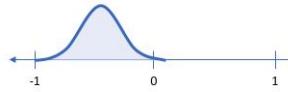


Smile (discrete value)



vs.

Smile (probability distribution)

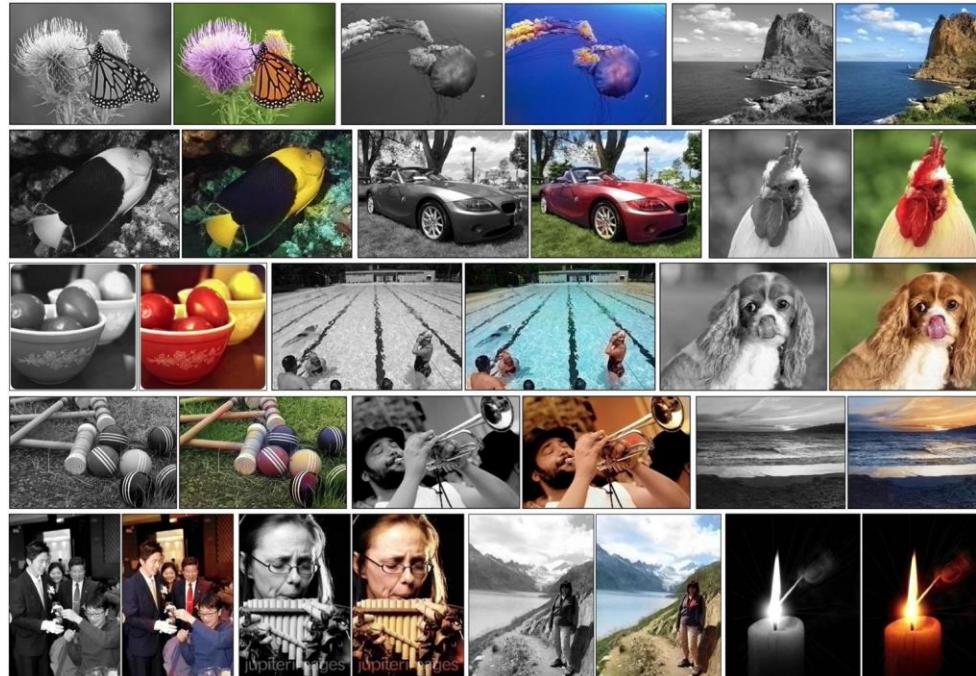


Latent attributes



# Image Colorization

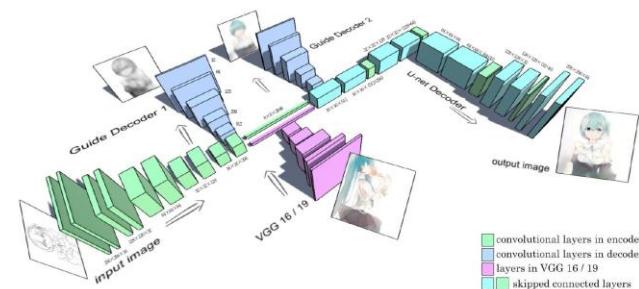
**Input:** Gray-scale image. **Output:** Colored image.  
**Method:** fully-convolutional network



# Style Transfer for Anime Sketches with Enhanced Residual U-net and Auxiliary Classifier GAN



- Apply the style of a painting to a grayscale sketch
- Residual U-Net
  - Image segmentation
  - Large hints from VGG 19 outputs
- with Auxiliary Classifier Generative Adversarial Network (AC-GAN)



<https://arxiv.org/abs/1706.03319>

ConvNet that can **predict next step of time sequence**, using a clever architecture for processing a large temporal context (about 3000-6000 past time steps)

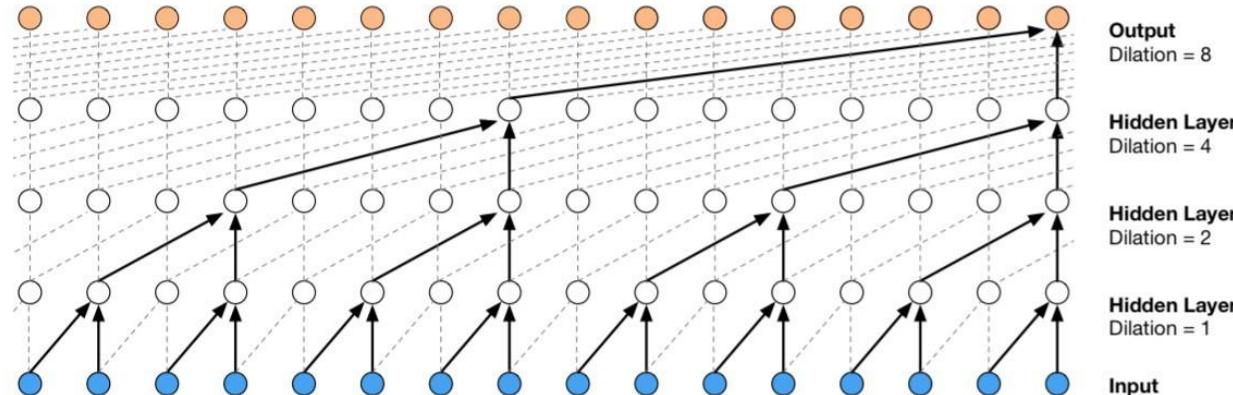
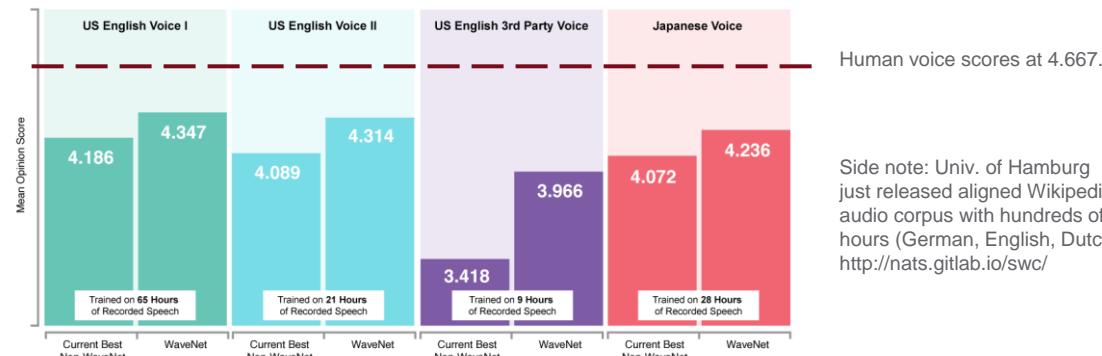


Figure 3: Visualization of a stack of *dilated* causal convolutional layers.

# Text-to-Speech WaveNet in production

- Previously discussed:
  - WaveNet (DeepMind, Sep 2016, arXiv 1609.03499)
  - Deep Voice 1 (Baidu, Feb 2017, arXiv 1702.07825)
  - Tacotron (Google, Mar 2017, arXiv 1703.10135)
  - Deep Voice 2 (Baidu, May 2017, arXiv 1705.08947)
- Oct 2017: Improved WaveNet used for Google Assistant
  - Generate 20 s of audio in 1 s (old WaveNet: 20 ms in 1 s)
  - 24 kHz, 16 bit (old WaveNet: 16 kHz, 8 bit)

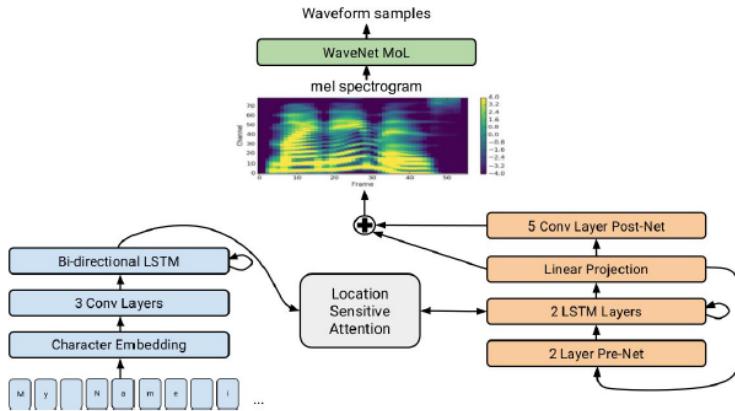
Mean Opinion Scores



Side note: Univ. of Hamburg just released aligned Wikipedia audio corpus with hundreds of hours (German, English, Dutch): <http://nats.gitlab.io/swc/>

# Tacotron 2

- Text-to-speech directly from characters, learned end-to-end
- Improvement over Tacotron: mel spectrogram + WaveNet instead of linear spectrogram + Griffin-Lim
- Improvement over some earlier WaveNet-based systems: no need to explicitly predict phones, timing and f0
- Comparison to Deep Voice 3: higher MOS, more mispronounciations, fewer repetitions or skipped words



Name	MOS
Parametric	$3.492 \pm 0.096$
Tacotron (Griffin-Lim)	$4.001 \pm 0.087$
Concatenative	$4.166 \pm 0.091$
WaveNet (Linguistic)	$4.341 \pm 0.051$
Ground Truth	$4.582 \pm 0.053$
Tacotron 2 (this paper)	<b><math>4.526 \pm 0.066</math></b>

# Tacotron-2 Examples

- Synthetic or Real?

“That girl did a video about Star Wars lipstick.”



“She earned a doctorate in sociology at Columbia University.”



“George Washington was the first President of the United States.”



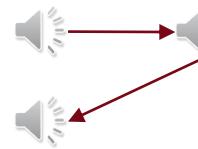
Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., ... & Saurous, R. A. (2018, April). Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 4779-4783). IEEE.

- Stress and intonation
- Questions
- Prosody
  - Intonation, rhythm, tone

#### Style / Reference

**Reference text:** Alice was not much surprised at this, she was getting so used to queer things happening.

#### Base Voice Model



#### Result

**Perturbed text:** Eric was not much surprised at this, he was getting so used to TensorFlow breaking.

#### Singing



Skerry-Ryan, R. J., Battenberg, E., Xiao, Y., Wang, Y., Stanton, D., Shor, J., ... & Saurous, R. A. (2018). Towards end-to-end prosody transfer for expressive speech synthesis with tacotron. arXiv preprint arXiv:1803.09047.

# VISUAL CONTENT

## Finally, Deep Fakes!



# Steps to create a Deep Fake with Deep Face Lab

1. Define Source and Destination Videos
2. Extract faces
3. Remove all but destination faces (other persons)
4. Clean source faces (false detections, occlusions, etc.)
5. Align faces
6. Train model
7. Apply model → Convert destination to source faces
8. Assemble video from all frames



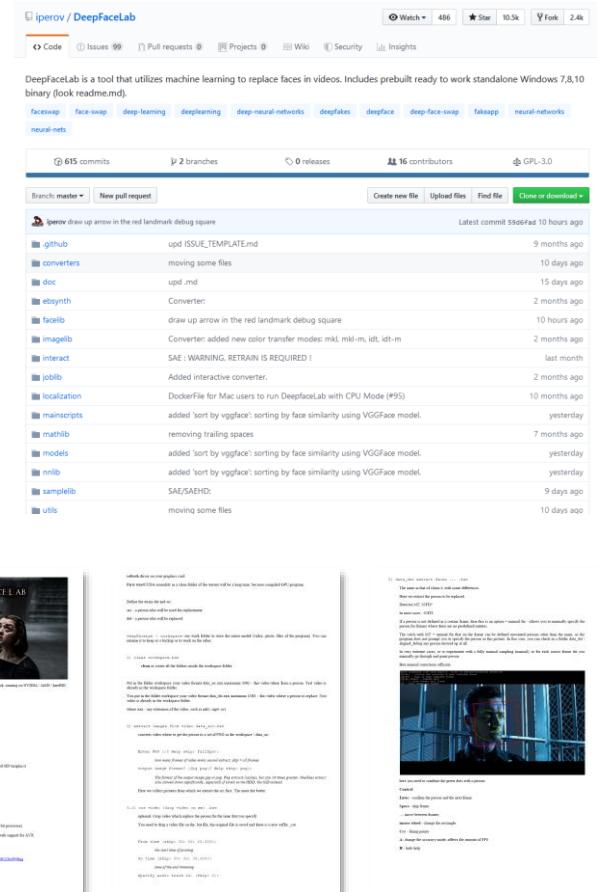
## Effort to create this Deep Fake

- ~ 10min manual interaction (for everything)
- ~ 3h computation time – pre-processing
- ~ 1h computation time – model training (10K epochs / 10% of normal)
- ~ 1h computation time – conversion
- ~ 10min computation time – re-assembling video
- ~ 6h to create this video!

# Deep Face Lab

- Available on Github
- Step-wise documentation
- Prepared, enumerated scripts for Windows

- 1) clear workspace.bat
- 2) extract images from video data\_src.bat
- 3.1) cut video (drop video on me).bat
- 3.2) extract images from video data\_dst FULL FPS.bat
- 3.other) denoise extracted data\_dst.bat
- 4) data\_src extract faces MANUAL.bat
- 4) data\_src extract faces MT all GPU debug.bat
- 4) data\_src extract faces MT all GPU.b



The screenshot shows the GitHub repository for DeepFaceLab. At the top, there's a header with 'iperov / DeepFaceLab' and various repository statistics: Watch (486), Star (10.5k), Fork (2.4k). Below the header are tabs for Code, Issues (99), Pull requests (0), Projects (0), Wiki, Security, and Insights.

The main content area displays the repository's code structure. It includes sections for faceapp, face-swap, deep-learning, deeplearning, deep-neural-networks, deepfaces, deepface, deep-face-swap, faceapp, and neural-networks. Key files listed include .gitignore, .github, converters, moving some files, upd.jnd, upd.jnd, draw up arrow in the red landmark debug square, Converter, draw up arrow in the red landmark debug square, Converter: added new color transfer modes: mkl, mkl-m, idt, idt-m, interact, SAE : WARNING. RETRAIN IS REQUIRED !, Added interactive converter, localization, Dockerfile for Mac users to run DeepFaceLab with CPU Mode (#95), mainscripts, added 'sort by vggface': sorting by face similarity using VGGFace model, matlib, removing trailing spaces, models, added 'sort by vggface': sorting by face similarity using VGGFace model, nnlib, added 'sort by vggface': sorting by face similarity using VGGFace model, samplelib, SAE/SAEHD, and utils, moving some files.

Below the code listing is a commit history showing the last 10 commits. The commits are dated from 9 months ago to yesterday. The latest commit was made 10 hours ago.

At the bottom of the screenshot, there are two examples of the DeepFaceLab software interface. The left example shows a woman's face being processed with a progress bar at 100%. The right example shows a man's face being processed with a progress bar at 100%.

# Corridor Crew – Deep Fake Demonstrations

We Made The  
Best Deepfake on  
The Internet

<https://www.youtube.com/watch?v=3vHvOyZ0GbY&t=886s>



How We Faked  
Keanu Reeves  
Stopping a  
Robbery

<https://www.youtube.com/watch?v=lzEFnbZ0Zd4>



# Never hold a sign into a camera

- No AI required!



RACISTS  
RAPISTS

RAP

Fake – Das Mädchen mit dem Schild „Will trade racists for rapists“,  
<https://www.mimikama.at/allgemein/fake-schild-will-trade-racists-for-rapists/>

# DETECTION APPROACHES

## Visual tampering / forgery



# Forgery Detection

- Identify GAN generated / forged images
- Encoder-Decoder Network
- Learn disentangled partitioned latent spaces
  - Real image embedding space
  - Fake image embedding space
  - Difference used to detect fakes
- Disentangled loss
- VGG-like architecture
- Class Activation Mapping (CAM)
  - Visualize / explain faked image regions

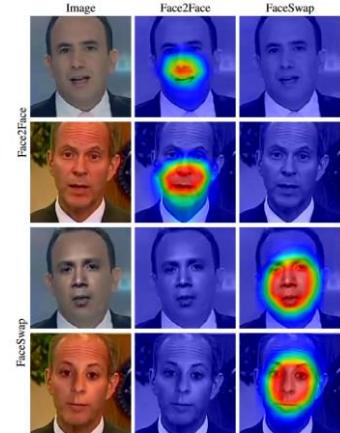
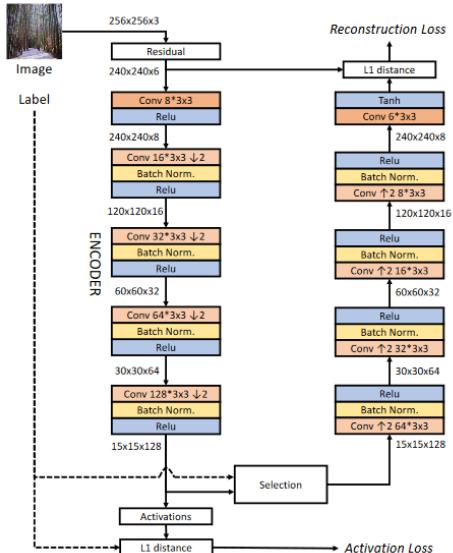


Figure 2: Two examples of images manipulated with Face2Face [39] and FaceSwap [2] (left) and their corresponding class activation maps, when the network (XceptionNet [10]) is trained on Face2Face forgeries (middle) and when it is trained on FaceSwap ones (right).

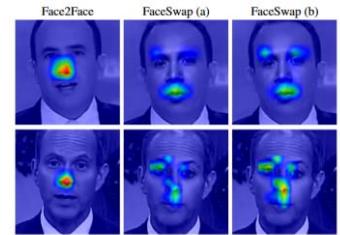


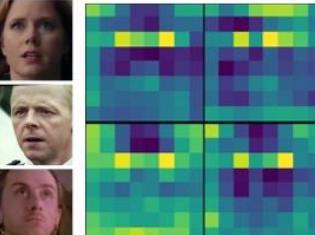
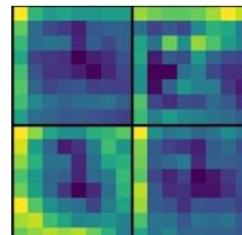
Figure 3: Class activation maps for our method when it is trained on Face2Face [39] and tested on Face2Face forgeries (left) or FaceSwap [3] (middle), and finally trained on Face2Face but fine-tuned using only four images manipulated with FaceSwap and tested on FaceSwap (right).

# Detecting GAN generated Faces

- Identify GAN generated faces
- VGG-like architecture / Inception variant
- Train on DeepFake Dataset
- High accuracy with simple approach
- Observation
  - Activations
    - Fake → Background
    - Real → Eyes



mean layer output of 100 *deepfake* faces



mean layer output of 100 *real* faces

Marra, F., Gragniello, D., Cozzolino, D., & Verdoliva, L. (2018, April). Detection of GAN-generated fake images over social networks. In 2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR) (pp. 384-389). IEEE.

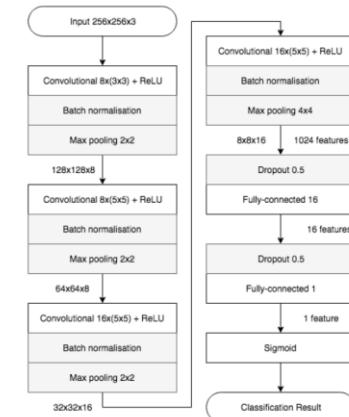
Network	Deepfake classification score		
Class	forged	real	total
Meso-4	0.882	0.901	0.891
MesoInception-4	0.934	0.900	0.917

Table 3. Classification scores of several networks on the *Deepfake* dataset, considering each frame independently.

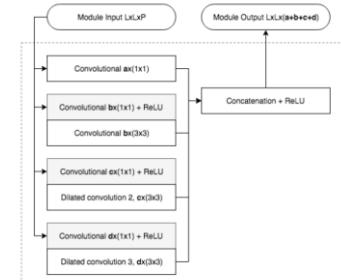
Network	Aggregation score	
Dataset	Deepfake	Face2Face (23)
Meso-4	<b>0.969</b>	<b>0.953</b>
MesoInception-4	<b>0.984</b>	<b>0.953</b>

Table 5. Video classification scores on the two dataset using image aggregation, with the *Face2Face* dataset compressed at rate 23.

Model Architecture



Inception Variant

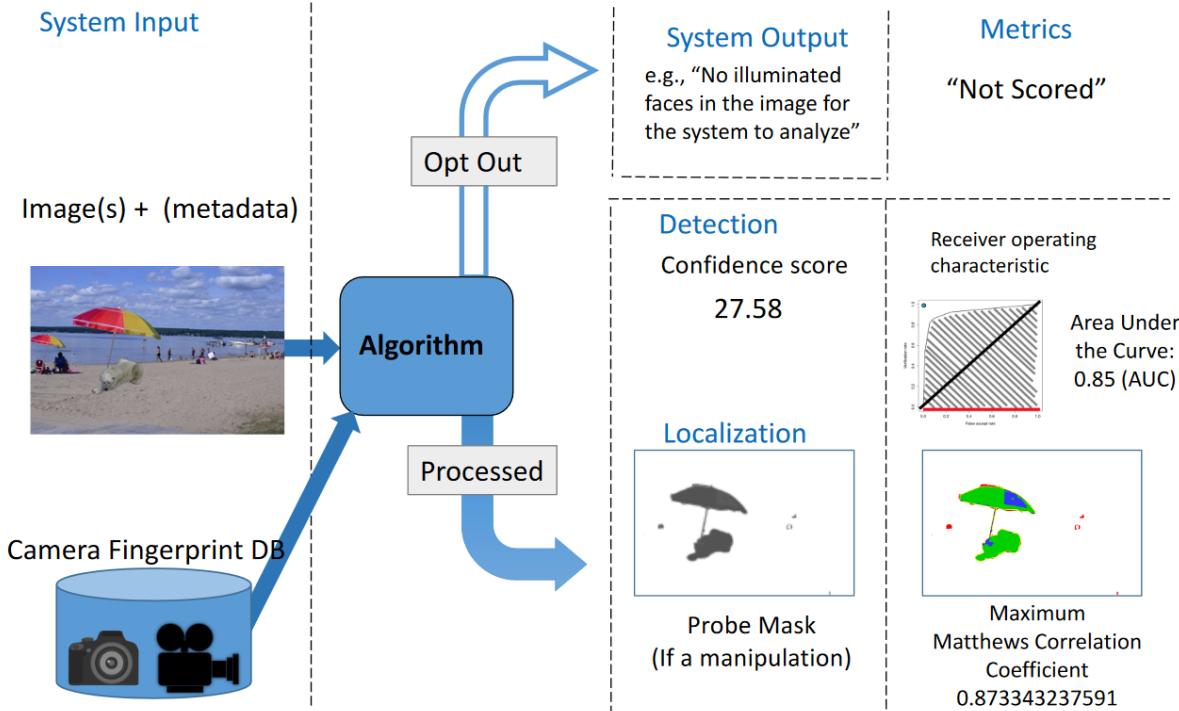


# Nimble Challenge

- Motivation
  - Media Forensic Technology Development
  - Develop Tools, Evaluation Tasks, Datasets
- 4 Tasks
  - Manipulation Detection and Localization
  - Splice Detection and Localization
  - Provenance Filtering
  - Provenance Graph Building
- Hosted in 2017, 2018
- Successor
  - Media Forensics Challenge 2018, 2019

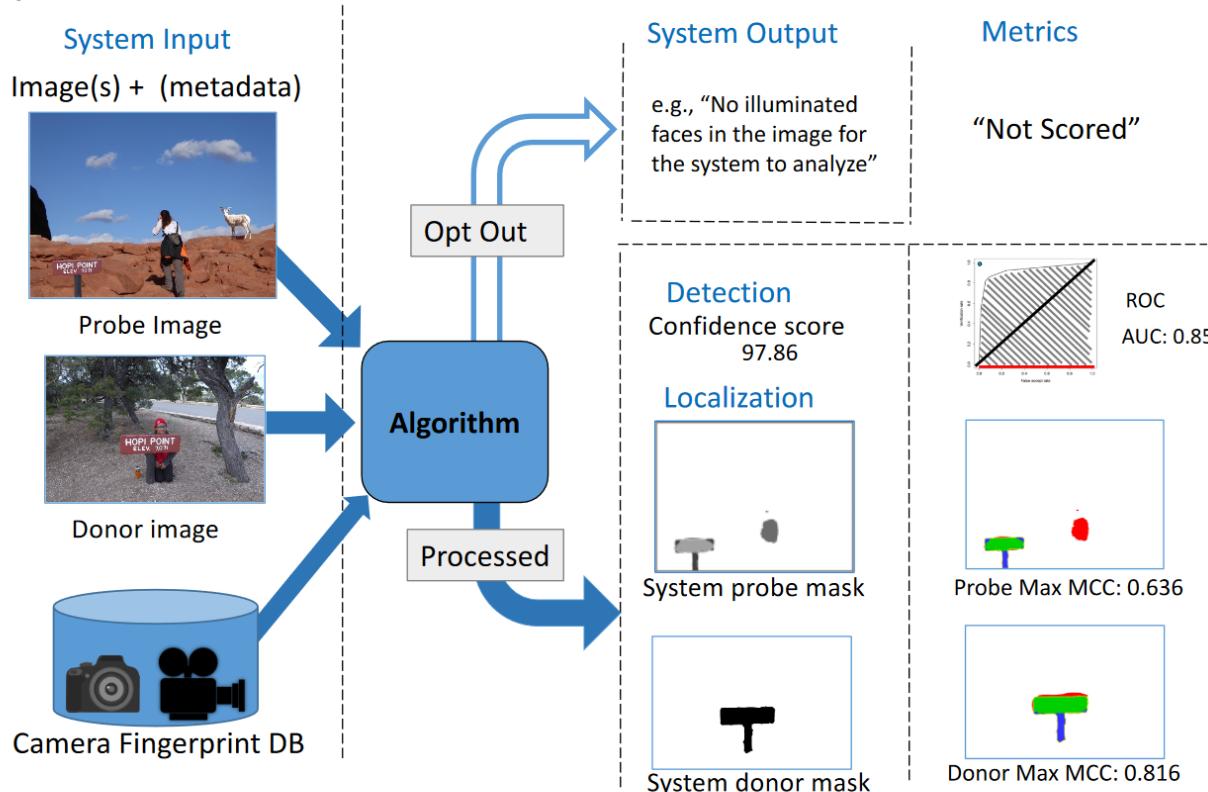


## Manipulation Detection and Localization Evaluation Task



Nimble Challenge 2017 Evaluation - Slides: for  
the 1 year PI meeting  
[https://www.nist.gov/sites/default/files/documents/2017/09/05/nist\\_medinfor\\_july17pimeeting-merge\\_v12.pdf](https://www.nist.gov/sites/default/files/documents/2017/09/05/nist_medinfor_july17pimeeting-merge_v12.pdf)

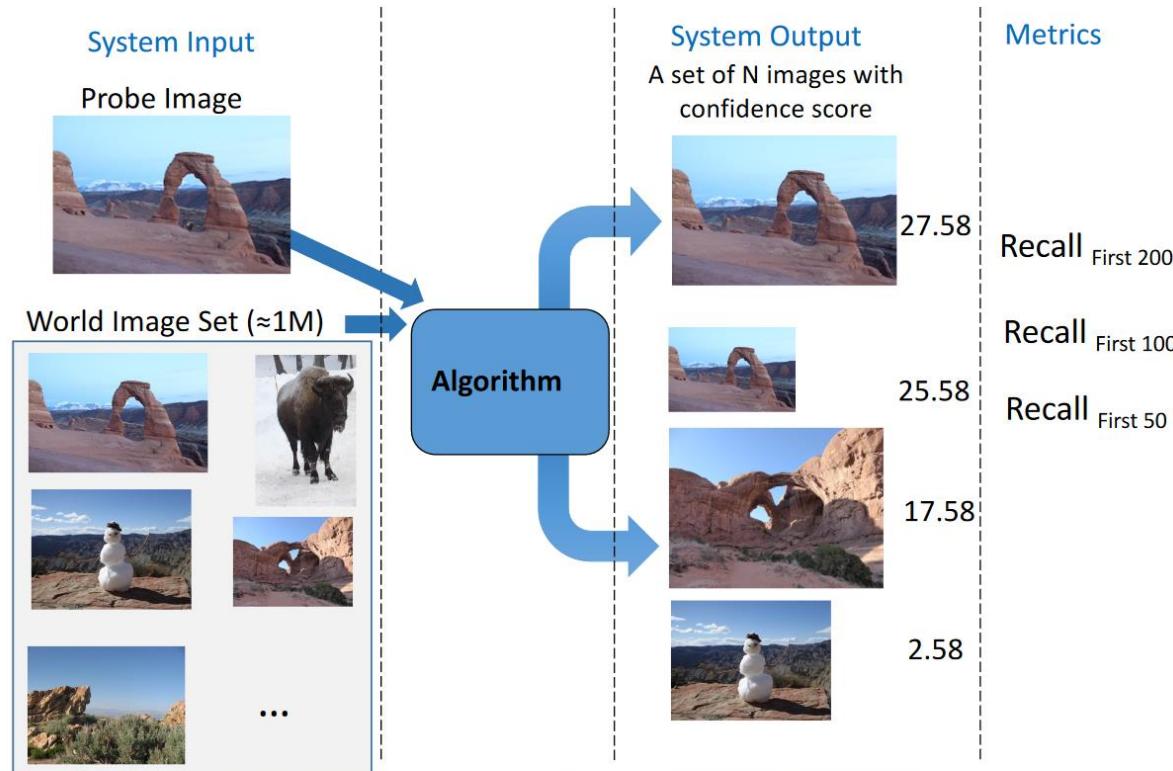
## Splice Detection and Localization Evaluation Task



Nimble Challenge 2017 Evaluation - Slides: for  
the 1 year PI meeting  
[https://www.nist.gov/sites/default/files/documents/2017/09/05/nist\\_medinfor\\_july17pimeeting-merge\\_v12.pdf](https://www.nist.gov/sites/default/files/documents/2017/09/05/nist_medinfor_july17pimeeting-merge_v12.pdf)

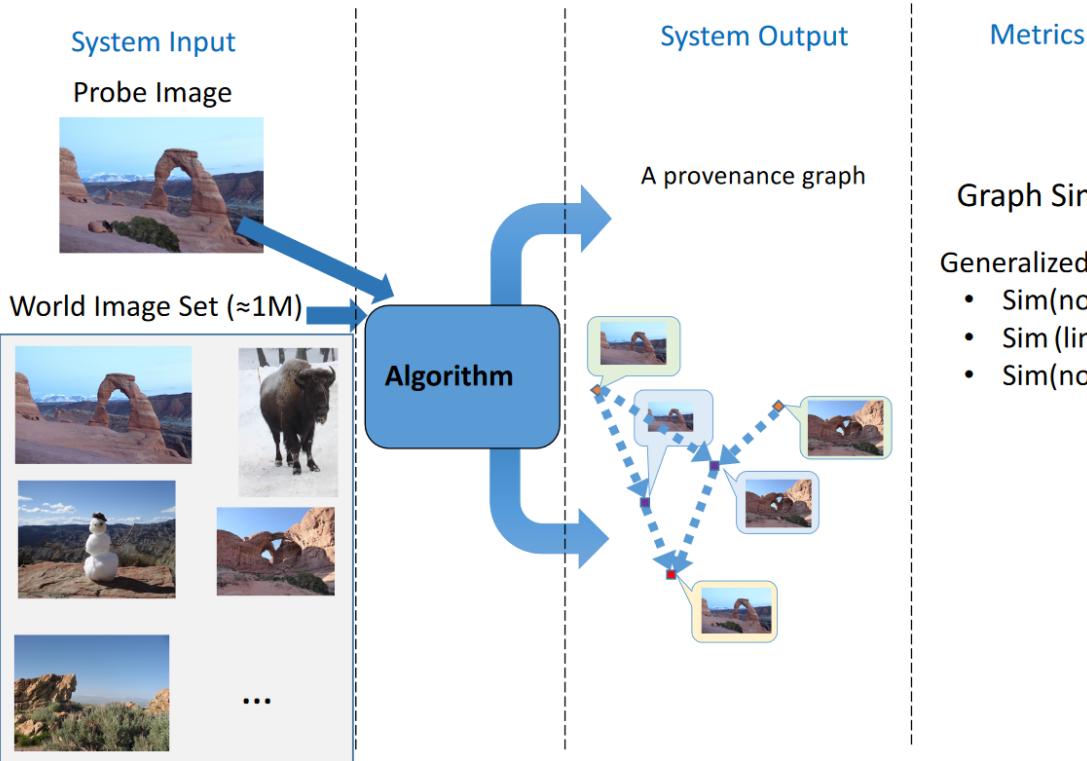
# Nimble Challenge 2017

## Provenance Filtering Evaluation Task



Nimble Challenge 2017 Evaluation - Slides: for  
the 1 year PI meeting  
[https://www.nist.gov/sites/default/files/documents/2017/09/05/nist\\_medinfor\\_july17pimeeting-merge\\_v12.pdf](https://www.nist.gov/sites/default/files/documents/2017/09/05/nist_medinfor_july17pimeeting-merge_v12.pdf)

## Provenance Graph Building Evaluation Task



Nimble Challenge 2017 Evaluation - Slides: for  
the 1 year PI meeting  
[https://www.nist.gov/sites/default/files/documents/2017/09/05/nist\\_medinfor\\_july17pimeeting-merge\\_v12.pdf](https://www.nist.gov/sites/default/files/documents/2017/09/05/nist_medinfor_july17pimeeting-merge_v12.pdf)



# AUDIO TAMPERING



# Fake Audio

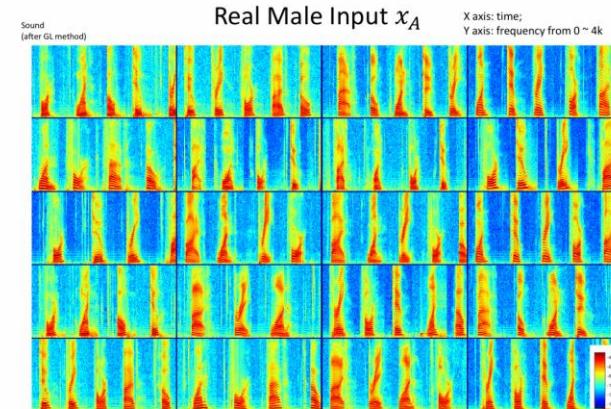
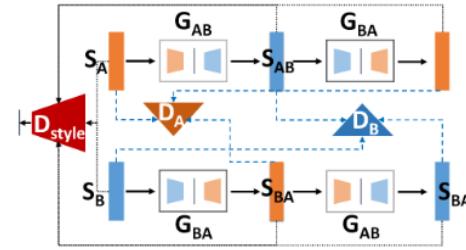
- **No Deep Fake for Audio!!!**
  - Deep Fakes with voice require
    - Traditional audio manipulation
    - Impressionist
- **Approaches to style-transfer in audio**
  - Not promising



Deep Fake VFX - Pity the poor impressionist by Jim Meskimen <https://www.youtube.com/watch?v=Wm3squcz7Aw>

# Voice-GAN

- Spectrogram-based-GAN
- Examples
  - Female
  - Female → Male
  - Female → Male → Female

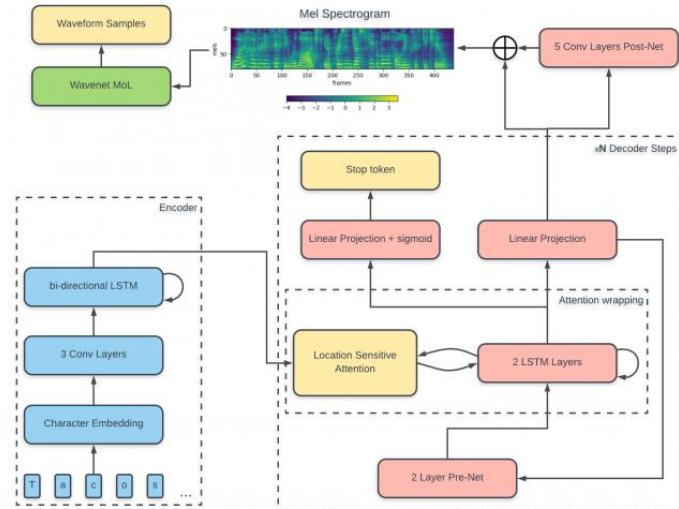


Gao, Y., Singh, R., & Raj, B. (2018, April). Voice impersonation using generative adversarial networks. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 2506-2510). IEEE.

# Speech Synthesis

- **Text to Speech**  
→ different approach
- **Tacotron-2**
  - Conditioned Wavenets

## Model Architecture:



Tacotron-2 <https://github.com/Rayhane-mamah/Tacotron-2>

# Tacotron-2 Examples

- Synthetic or Real?

“That girl did a video about Star Wars lipstick.”



“She earned a doctorate in sociology at Columbia University.”



“George Washington was the first President of the United States.”



Shen, J., Pang, R., Weiss, R. J., Schuster, M., Jaitly, N., Yang, Z., ... & Saurous, R. A. (2018, April). Natural tts synthesis by conditioning wavenet on mel spectrogram predictions. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 4779-4783). IEEE.

- Stress and intonation
- Questions
- Prosody
  - Intonation, rhythm, tone

## Style / Reference

**Reference text:** Alice was not much surprised at this, she was getting so used to queer things happening.



## Result

**Perturbed text:** Eric was not much surprised at this, he was getting so used to TensorFlow breaking.



## Singing



Skerry-Ryan, R. J., Battenberg, E., Xiao, Y., Wang, Y., Stanton, D., Shor, J., ... & Saurous, R. A. (2018). Towards end-to-end prosody transfer for expressive speech synthesis with tacotron. arXiv preprint arXiv:1803.09047.

# Simple but effective / No AI



Pelosi videos manipulated to make her appear drunk are being shared on social media  
<https://www.youtube.com/watch?v=sDOo5nDjwgA>

**FFmpeg**

Anmelden Einstellungen | Hilfe/Anleitung | Über Trac | Register  
 Wiki Journal Tickets anzeigen Suche Tags

Wiki: **How to speed up / slow down a video**

**Speeding up/slowing down video**

You can change the speed of a video stream using the `>setpts` video filter. Note that in the following examples, the audio stream is not changed, so it should ideally be disabled with `-an`.

To double the speed of the video, you can use:

```
ffmpeg -i input.mkv -filter:v "setpts=0.5*PTS" output.mkv
```

The filter works by changing the presentation timestamp (PTS) of each video frame. For example, if there are two successive frames shown at timestamps 1 and 2, and you want to speed up the video, those timestamps need to become 0.5 and 1, respectively. Thus, we have to multiply them by 0.5.

Note that this method will drop frames to achieve the desired speed. You can avoid dropped frames by specifying a higher output frame rate than the input. For example, to go from an input of 4 FPS to one that is sped up to 4x that (16 FPS):

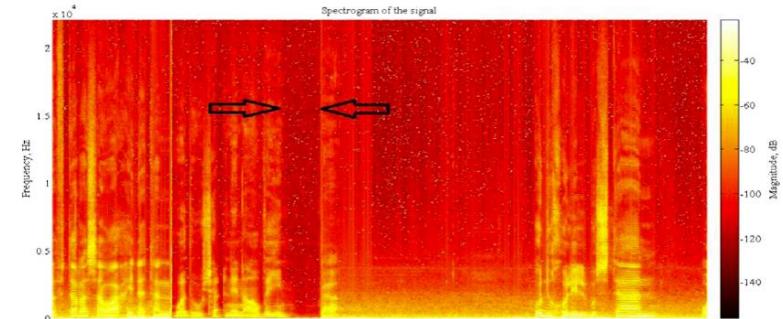
```
ffmpeg -i input.mkv -r 16 -filter:v "setpts=0.25*PTS" output.mkv
```

To slow down your video, you have to use a multiplier greater than 1:

```
ffmpeg -i input.mkv -filter:v "setpts=2.0*PTS" output.mkv
```

# Audio tampering research

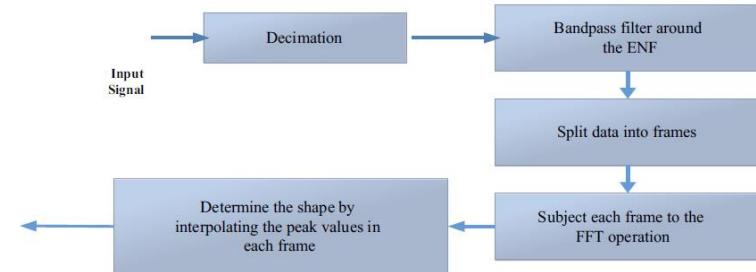
- Motivation
  - Identify modified or fabricated content
  - E.g. Forged / faked evidence in court cases
- Common Modifications
  - Splicing, copying, moving, insertions
- Common Countermeasures
  - Local noise level estimation: Splicing → different noise levels
  - Exploring pitch similarity: Copy-move
  - Electric Network Frequency (ENF)
- Problems
  - Many traces get lost in compression (e.g. mp3)



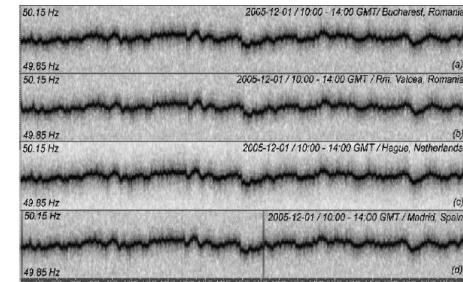
Zakariah, M., Khan, M. K., & Malik, H. (2018). Digital multimedia audio forensics: past, present and future. *Multimedia tools and applications*, 77(1), 1009-1040.

# Electronic Network Frequency (ENF)

- Traces of ENF in the recording
  - ENF deviates from 50 to 60Hz
  - Distinct pattern
  - Similar on different networks
  - Captured in recording
- Compares
  - Estimated ENF signature of recording
  - Reference frequency database by power supply company
- Max offset for cross correlation (MOCC)
  - For query ENF and reference ENF



Zakariah, M., Khan, M. K., & Malik, H. (2018). Digital multimedia audio forensics: past, present and future. *Multimedia tools and applications*, 77(1), 1009-1040.



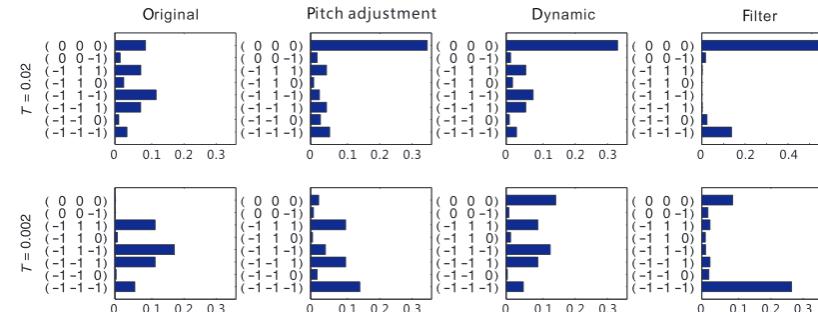
Grigoras, C. (2007). Applications of ENF criterion in forensic audio, video, computer and telecommunication analysis. *Forensic science international*, 167(2-3), 136-145.

# Audio Postprocessing Detection

- Amplitude Cooccurrence Vectors (ACV)
- Gray Level coocurrence matrix (GLCM)
  - Joint probability of two pixels
  - Texture analysis / characterize coocurrence patterns
  - Adapted to audio
- Distinguish tampered audio
- Identify type of modification



<https://www.youtube.com/watch?v=sDOo5nDjwgA>



Luo, D., Sun, M., & Huang, J. (2016). Audio postprocessing detection based on amplitude cooccurrence vector feature. IEEE Signal Processing Letters, 23(5), 688-692.

# TEXT BASED FAKE NEWS DETECTION

## TEXT CONTENT ANALYSIS

Mina Schütz



# FAKE NEWS DETECTION: TEXT CONTENT

- **Classification task**
  - Binary or multi-classification
- **Regression task**
  - Output as a numeric score of truthfulness
- **Similar NLP concepts**
  - Fake product reviews
  - Online resumes
  - Opinion spamming
  - Fake profiles
  - Spamming and phishing

# CONTENT-BASED APPRoAChES

- Headline & body text
- **Linguistic cues and patterns**
  - Character, word, sentence or document level
  - Linguistic Inquiry and Word Count (LIWC)
- **Style-based**
  - Hashtags, mentions, punctuation marks, sentiment
  - Topics languages, domains

Attribute Type	Feature
Quantity	Character count
	Word count
	Noun count
	Verb count
	Number of noun phrases
	Sentence count
	Paragraph count
Complexity	Number of modifiers (e.g., adjectives and adverbs)
	Average number of clauses per sentence
	Average number of words per sentence
	Average number of characters per word
Uncertainty	Average number of punctuations per sentence
	Percentage of modal verbs
	Percentage of certainty terms
	Percentage of generalizing terms
	Percentage of tentative terms
	Percentage of numbers and quantifiers
Subjectivity	Number of question marks
	Percentage of subjective verbs
	Percentage of report verbs
	Percentage of factive verbs
	Percentage of imperative commands

# Knowledge-based approaches

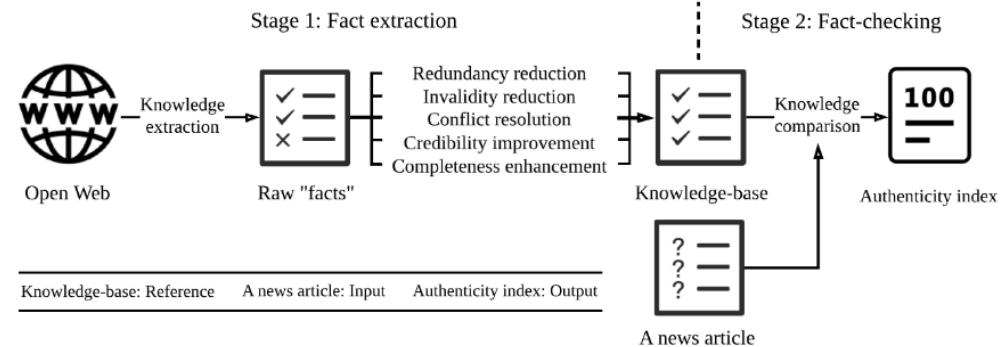
- Check for truthfulness of claims

- **Manual fact-checking**

- Expert-based
- Crowdsources
- politifact.com, snopes.com

- **Automatic fact-checking**

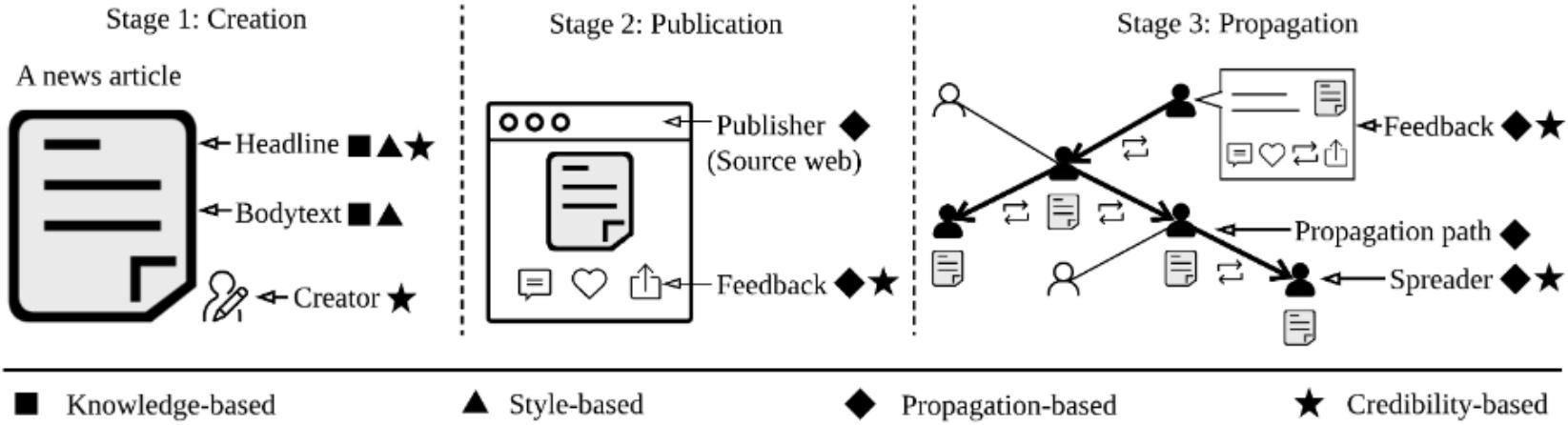
- Linked Open Data (i.e. DBpedia)
- Fact-extraction (Knowledge base construction)
- Fact-checking (Knowledge comparison)



# SOCIAL-CONTEXT Based APPROACHES

- **Stance-based**
  - Stance of body text relative to the headline claim
  - Viewpoint of user
  - Infer validity of original article
  - Support or refute claim
- **User-based**
  - Registration age
  - Numer of followers / followees
- **Propagation-based**
  - Propagation networks
  - i.e. Twitter shares / retweets
  - Likes
- **Credibility-based**
  - Headlines (clickbaits)
  - News source, spreaders, author

# HYBRID APPROACHES



# preprocessing

- Tokenization
- Stemming
- Lemmatization
- Part-of-Speech Tagging
- Generalization
- Vector representation
  - Bag-of-Words, N-Gram, LDA
  - TF, TF-IDF
  - Word2Vec, GloVe
  - WordPiece Model
  - FastText
- **Datasets**
  - FEVER
  - LIAR
  - CREDBANK
  - BuzzFace / BuzzFeedNews
  - PHEME
  - Some-like-it-hoax
  - Kaggle datasets (WSDM)
  - FakeNewsNet
  - Fake News Challenge (FNC-1)
  - GermanFakeNC

# Classification MODELS & Evaluation

- **Supervised**
  - Naïve Bayes Classifier
  - Logistic Regression
  - Decision Trees / Random Forest
  - Support Vector Machines
- **Evaluation Metrics**
  - Precision
  - Recall
  - F1
  - Accuracy
  - ROC
  - ROC AUC
  - Cross-fold-validation
- **Unsupervised**
  - Recurrent Neural Network
  - Long Short-Term Memory
  - Convolutional Neural Networks
  - Uni- or bidirectional Transformers

# CHALLENGES

- Datasets are biased, too small or not available anymore
  - Domain / dataset specific
  - Many features to choose from
  - Model learns writing style of a news source (not if it is fake news)
  - Often only linguistic cues or wordcount without context
- Almost no knowledge-based approaches
- No common fake news definition
  - No standardized approach available (Accuracy from ~ 20% to ~ 95%)
  - Early fake news identification

# FAKE NEWS PROBLEM

Overestimated?

Technical Problem?

Solvable?



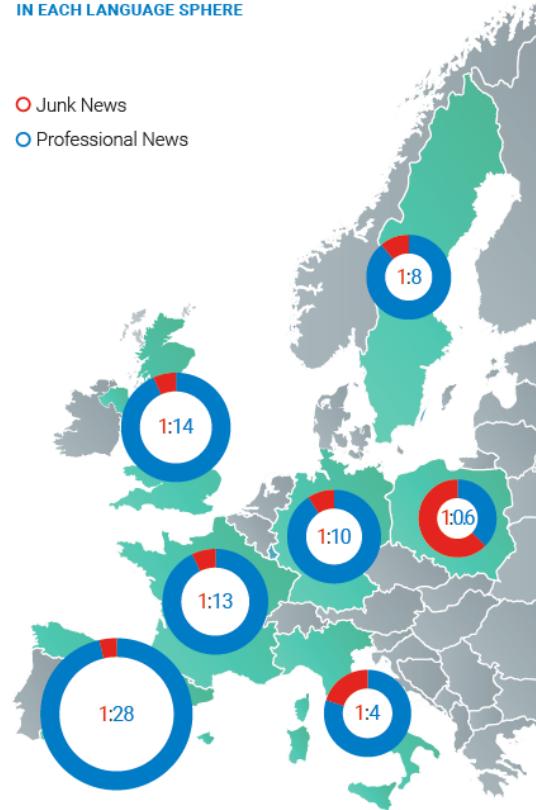
# Conclusions

- **General**
  - Fake News detection currently a hyped topic
  - Fear mongering to create business opportunities
- **„Traditional“ attacks**
  - Media Forensic
    - traditional research fields
    - Many well established approaches
    - Ongoing evaluation campaigns
- **AI-Threat**
  - Low-Res GAN – no threat
  - GAN content identifiable (incl. Deep Fakes)
  - Challenge: future Hi-Res GANs
- **Biggest Challenge**
  - Text content analysis
  - Claim / Fact checking

# Extent of Fake News in EU Parliament Elections

- Oxford Institute study
  - 7 different languages
- **Only 4% source articles → Fake**
  - ~600M Tweets
  - ~140K Users
  - ~6K unique articles

Figure 1 - RATIO OF JUNK TO PROFESSIONAL NEWS  
IN EACH LANGUAGE SPHERE



# Conclusions

- **General Discussion about Fake News**
  - Technological Problem?
  - Sociological Problem?
  - Political Problem?
- **Who is responsible?**
  - Tech-Companies (Twitter, Facebook, etc.)
- **Most cited solutions**
  - Renaissance of quality journalism
  - Promote media literacy

# RECOMMENDED LITERATURE



# Surveys

- Kumar, S., & Shah, N. (2018). **False information on web and social media: A survey**. *arXiv preprint arXiv:1804.08559*.
- Sharma, K., Qian, F., Jiang, H., Ruchansky, N., Zhang, M., & Liu, Y. (2019). **Combating fake news: A survey on identification and mitigation techniques**. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 10(3), 21.

arXiv:1901.06479v1 [cs.IJG] 18 Jan 2019

**Combating Fake News: A Survey on Identification and Mitigation Techniques**

KARISHMA SHARMA, University of Southern California  
 FENG QIAN, University of Southern California  
 HE JIANG, University of Southern California  
 NATHALIA RUCHANSKY, University of Southern California  
 MING ZHANG, Peking University  
 YAN LIU, University of Southern California

The proliferation of fake news on social media has opened up new directions of research for timely identification and containment of fake news, and mitigation of its widespread impact on public opinion. While much of the work on fake news has focused on identifying fake news from user-generated content, with the rise of fake news engagements with the news on social media, there has been a rising interest in proactive intervention strategies to counter the spread of misinformation and its impact on society. In this survey, we describe the modern-day challenges of combating fake news and propose a taxonomy of methods for combating fake news. We also discuss the various existing methods and techniques applicable to both identification and mitigation, with a focus on the significant advances in each method and their advantages and limitations. In addition, research has often been limited to specific domains such as politics, health, and science. In this survey, we aim to provide a more comprehensive and summarative characterization of available datasets. Furthermore, we outline new directions of research to facilitate future development of effective and interdisciplinary solutions.

CCS Concepts: Information systems → Social networking sites; Data mining; Computing methodologies → Machine learning

Additional Key Words and Phrases: AI, fake news detection, rumor detection, misinformation

ACM Reference Format:

Sharma, Karishma, Qian, Feng, He, Jiang, Nathalia Ruchansky, Ming Zhang, and Yan Liu. 2019. Combating Fake News: A Survey on Identification and Mitigation Techniques. *ACM Trans. Intell. Syst. Technol.* 10, 4, Article 111 (August 2018), 41 pages. <https://doi.org/10.1145/322455.322456>

**1 INTRODUCTION**

In recent years, the topic of *fake news* has experienced a resurgence of interest in society. The increased attention stems largely from growing concerns around the widespread impact of fake news on society and politics. In January 2017, a spokesman for the German government stated that they “are dealing with a phenomenon of a dimension that [they] have never seen before”.

Authors' addresses: Karishma Sharma, Information Systems, University of Southern California; Feng Qian, University of Southern California, iis.usc.edu; He Jiang, Nathalia Ruchansky, University of Southern California, jiang@usc.edu; Ming Zhang, Peking University, mzhang@pku.edu; Yan Liu, University of Southern California, liuyan@usc.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions.acm.org](http://permissions.acm.org).

© 2019 Association for Computing Machinery.  
<https://doi.org/10.1145/322455.322456>

arXiv:1804.08559v1 [cs.IJG] 23 Apr 2018

## False Information on Web and Social Media: A Survey

SRIJAN KUMAR, Computer Science, Stanford University, USA  
 NEIL SHAH\*, Computer Science, Carnegie Mellon University, USA

Fake information can be created and spread easily through the web and social media platforms, resulting in widespread real world impact. Characterizing how fake information proliferates on social platforms and why it succeeds in deceiving readers are critical to develop efficient detection algorithms and tools for early detection. A recent surge of research in this area has led to the development of various methods for combating fake news, such as detection, identification, and mitigation. Majority of the research has primarily focused on two broad categories of fake information: opinion-based (e.g., fake reviews, rumors) and news-based (e.g., fake news). This survey aims to provide a comprehensive overview of the survey spanning diverse aspects of fake information, namely (i) the actors involved in spreading fake information, (ii) rationale behind successful deceiving readers, (iii) quantifying the impact of fake information, (iv) measuring its characteristics across different dimensions, (v) detection and mitigation techniques, and (vi) challenges and opportunities in building unified frameworks to describe these recent methods and highlight a number of important directions for future research.<sup>1</sup>

Additional Key Words and Phrases: misinformation, fake news, fake reviews, rumors, houses, web, internet, social media, social networks, fake news, fake, propaganda, conspiracy, knowledge bases, e-commerce, disinformation, impact, mechanism, rationale, detection, prediction

ACM Reference Format:  
 Kumar, Srijan, and Shah, Neil. 2018. False Information on Web and Social Media: A Survey. 1, 1 (April 2018), 35 pages. <https://doi.org/10.1145/3196403>

## 1 INTRODUCTION

The web provides a highly interconnected world-wide platform for everyone to spread information to millions of people in a matter of few minutes, at little to no cost [1]. While it has led to groundbreaking phenomena such as the growth of social media, it has also led to the propagation of fake news, rumors, and false information [2]. False information on the web and social media has affected stock markets [3], slowed responses during disasters [4], and terrorist attacks [27, 30]. Recent surveys have alarmingly shown that people believe fake news more than real news [5, 6]. Thus, combating fake news has become one of the most important importance to combat fake information on such platforms. With primary motives of influencing opinions and earning money [1, 46, 56, 57], the wide impact of fake information makes it one of the modern dangers to society,民族, and economy [1, 46, 56, 57]. Therefore, it is important to understand the nature of fake news and how it is created to proactively detect it and mitigate its impact. In this survey, we review the state-of-the-art scientific literature on fake information on the web and social media to give a comprehensive description of its mechanisms, characteristics, impact, and detection. While recent surveys have focused on fake

<sup>1</sup>This is a new at best approximation of the paper will appear in the book titled *Social Media Analytics: Advances and Applications*. © ACM, in 2018.

Authors' addresses: Srijan Kumar, Computer Science, Stanford University, USA, [wjgjwqz@cs.stanford.edu](mailto:wjgjwqz@cs.stanford.edu); Neil Shah, Computer Science, Carnegie Mellon University, USA, [neilshah@cs.cmu.edu](mailto:neilshah@cs.cmu.edu).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions.acm.org](http://permissions.acm.org).

© 2019 Association for Computing Machinery.  
<https://doi.org/10.1145/3196403>

## Bibliography

- Graves, Lucas (2018). Understanding the Promise and Limits of Automated Fact-Checking. *Factsheet February 2018*, Reuters Institute, 1 – 8. <https://reutersinstitute.politics.ox.ac.uk/our-research/understanding-promise-and-limits-automated-fact-checking>.
- Figueira, Álvaro; Oliveira, Luciana (2017). The current state of fake news: challenges and opportunities. *Procedia Computer Science* Vol. 121, 817 – 825. <https://www.sciencedirect.com/science/article/pii/S1877050917323086>. 10.1016/j.procs.2017.11.106.
- Mahid, Zaitul Iradah; Manickam, Selvakumar; Karuppiah, Shankar (2018). Fake News on Social Media: Brief Review on Detection Techniques. *2018 Fourth International Conference on Advances in Computing, Communication & Automation (ICACCA)*, Subang Jaya, Malaysia, 1 – 5. <https://ieeexplore.ieee.org/document/8776689>. 10.1109/ICACCAF.2018.8776689.
- Oshikawa, Ray; Qian, Jing; Wang, William Yang (2018). A Survey on Natural Language Processing for Fake News Detection. 1 – 11. <https://arxiv.org/abs/1811.00770>.

# Thank you!

Alexander Schindler and Mina Schütz

25.11.2019

