

FAKE NEWS

From Shallow to Deep. How to create, detect and fight it.



Alexander Schindler

Scientist

Information Management

Center for Digital Safety & Security

AIT Austrian Institute of Technology GmbH



Preface / Motivation

- **Why we talk about Fake News @VDLM?**
 - Recent Hype Topic
 - Many scary talks about the biggest threat to democracy



How the Obama / Jordan Peele DEEPCODEAK actually works | Ian Hislop's Fake News – BBC
<https://www.youtube.com/watch?v=g5wLaJYBAm4>



The Guardian https://www.theguardian.com/us-news/2019/jun/06/fake-news-how-misinformation-became-the-new-front-in-us-political-warfare

The Guardian https://www.theguardian.com/us-news/2019/jun/06/fake-news-how-misinformation-became-the-new-front-in-us-political-warfare



Trend Micro White paper on Fake News
https://documents.trendmicro.com/assets/white_papers/wp-fake-news-machine-how-propagandists-abuse-the-internet.pdf

Preface / Motivation

- **Why we talk about Fake News @VDLM?**
 - Recent Hype Topic
 - Many scary talks about the biggest threat to democracy
- **Why do I talk about Fake News?**
 - Fake News Team
 - AIT / Center for Digital Safety and Security
 - My research contribution
 - Multi-modal content retrieval / classification
 - Supervision of MSc students, internships
- **I'm not an expert on Fake News!**
 - Introduction to the problem domain

Agenda

- Introduction to Fake News
- Fake News Detection approaches
 - Visual
 - Accoustic
 - Textual
- Discussion and Conclusion

Who shares **Fake News**?

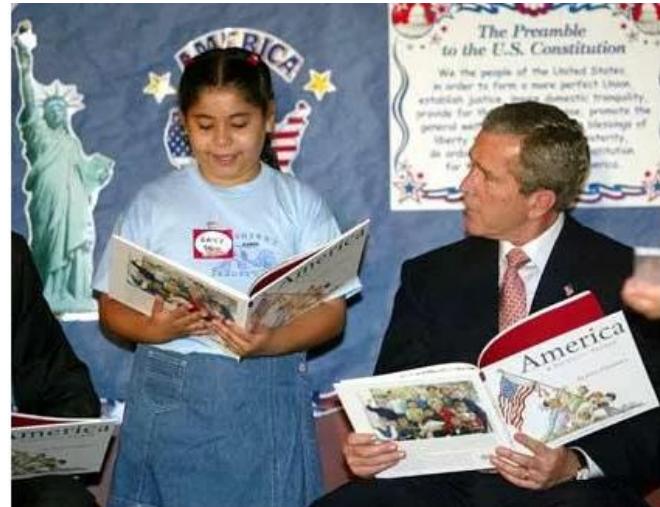
- Do you / have you shared **Fake News**?
- Are you likely to share **Fake News**?
- Can you identify **Fake News**?

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

30.10.2019



Additional Proof

Ideological Priors

or

Why do we fall for Fake News?

- **Naive Realism**

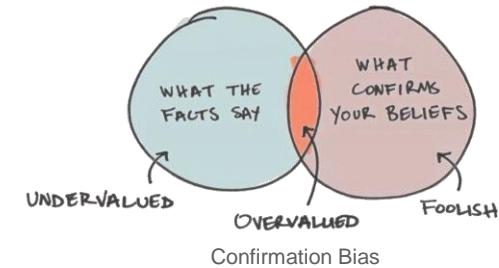
- Believe information that is aligned with your views

- **Confirmation Bias**

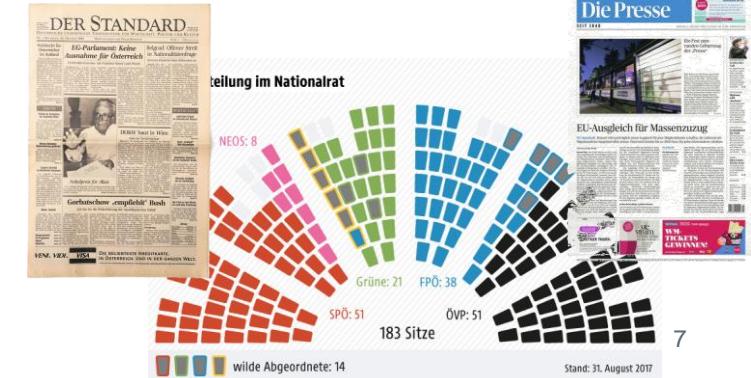
- Seek information that confirms your existing views

- **Normative Influence Theory**

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

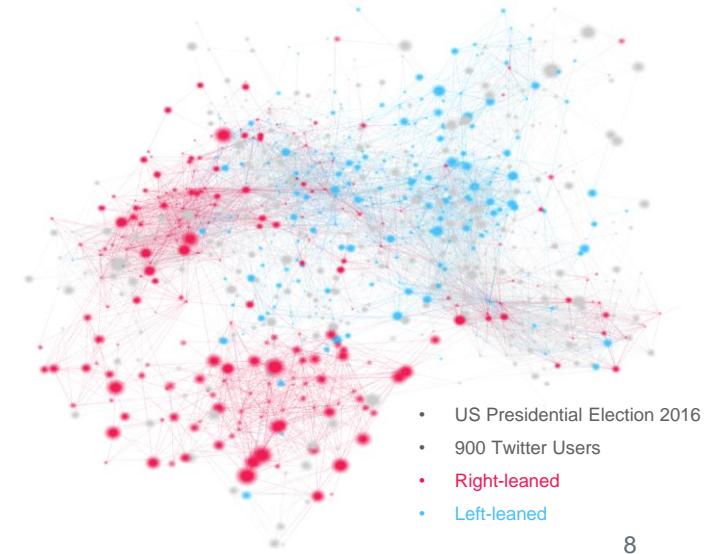
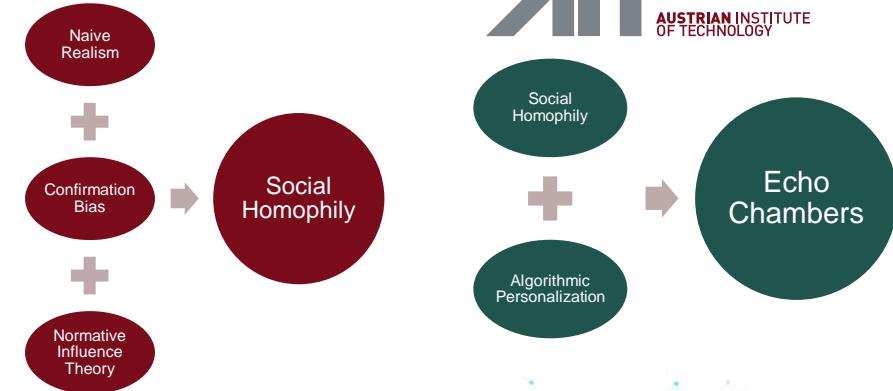


→ Individual level of Fake News



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



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

TOPICS: Pope Francis Endorses Donald Trump



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)

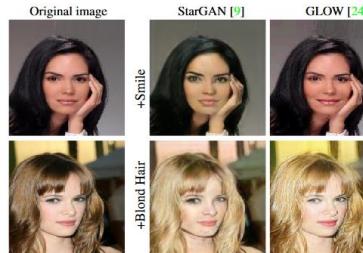
• False connection

- Headlines, visuals do not support content

• False context

- Genuine content shared with false context information

False Context



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.

Manipulated Content



Imposter Content



How Your Audience Will Believe Anything: The Psychology Behind the Fake News
<https://www.click.co.uk/blog/how-your-audience-will-believe-anything-the-psychology-behind-the-fake-news-bas-van-den-belds-benchmark-2018-talk-review/>



VOLUME 28 ISSUE 17

NUMBER ONE IN NEWS

12-16 DECEMBER 1998

Congress Hires Drummer

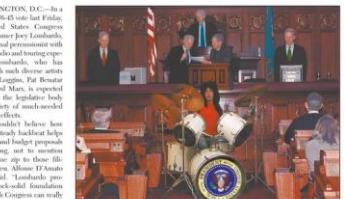


See Onion, page 22

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 Wednesday 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 necessarily false*
- *May be verified later as true / false*

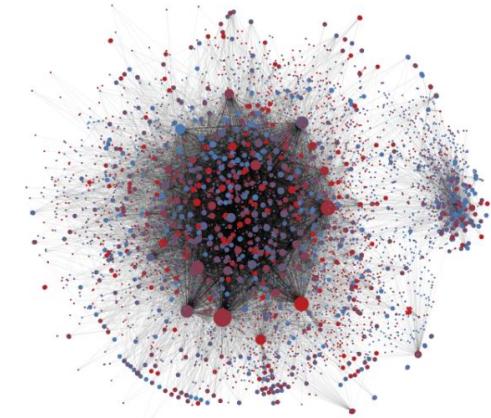
FAKE NEWS DETECTION



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

VISUAL CONTENT

Finally, Deep Fakes!



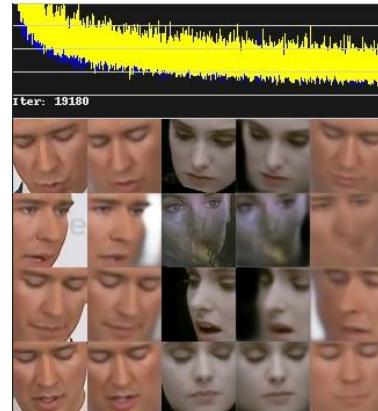


Steps to create a Deep Fake with Deep Face Lab

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



Destination Video



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

Screenshot of the DeepFaceLab GitHub repository page:

- Code**: 615 commits, 2 branches, 0 releases, 16 contributors, GPL 3.0
- Branch: master**, **New pull request**
- Issues**: 99, **Pull requests**: 0, **Projects**: 0, **Wiki**, **Security**, **Insights**
- Description**: DeepFaceLab is a tool that utilizes machine learning to replace faces in videos. Includes prebuilt ready to work standalone Windows 7,8,10 binary (look [readme.md](#)).
- Tags**: faceswap, face-swap, deep-learning, deeplearning, deep-neural-networks, deepfakes, deepface, deep-face-swap, fakesapp, neural-networks
- Contributors**: iperov, converters, doc, ebgynth, facelib, imagelib, interact, joblib, localization, mainscripts, matlib, models, nnlib, samplelib, utils
- Commits** (Recent):
 - draw up arrow in the red landmark debug square
 - update ISSUE_TEMPLATE.md
 - moving some files
 - upd .md
 - Converter: draw up arrow in the red landmark debug square
 - Converter: added new color transfer modes: mkl, mkl-m, idt, idt-m
 - SAE : WARNING: RETRAIN IS REQUIRED !
 - Added interactive converter.
 - Dockerfile for Mac users to run DeepFaceLab with CPU Mode (#95)
 - added 'sort by vggface': sorting by face similarity using VGGFace model.
 - removing trailing spaces
 - added 'sort by vggface': sorting by face similarity using VGGFace model.
 - added 'sort by vggface': sorting by face similarity using VGGFace model.
 - SAE/SAEHD:
 - moving some files

Two screenshots of the DeepFaceLab software interface showing the process of extracting faces from a video:

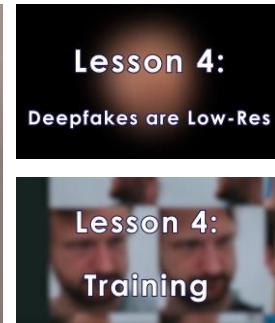
- Screenshot 1**: Shows the main interface with a video frame and various settings like "DEEP FACE LAB", "Video Input", "Video Output", and "Face Detection".
- Screenshot 2**: Shows a close-up view of a person's face with a bounding box and processing options.

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>



30.10.2019

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

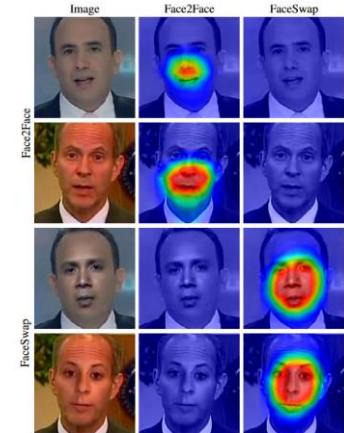
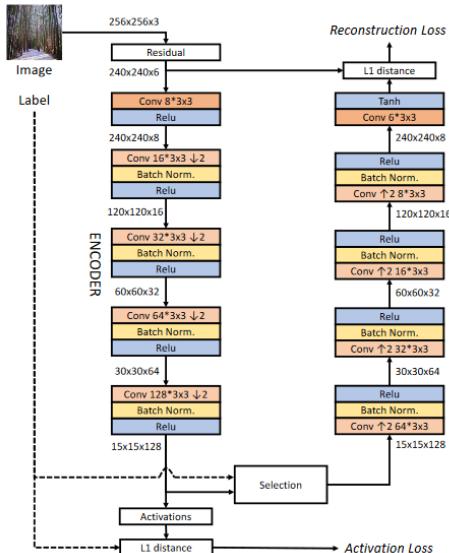
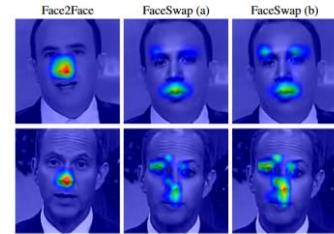


Figure 2: Two examples of images manipulated with Face2Face [39] and FaceSwap [2] (left) and their corresponding class activation maps, when the network (XceptionNet [10]) is trained on Face2Face forgeries (middle) and when it is trained on FaceSwap ones (right).



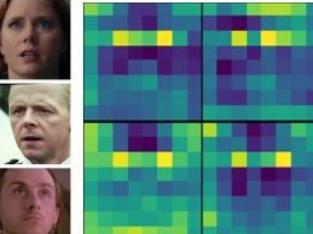
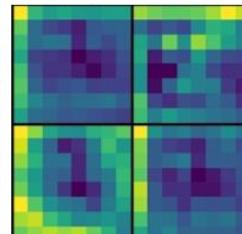
Cozzolino, D., Thies, J., Rössler, A., Riess, C., Nießner, M., & Verdoliva, L. (2018). Forensictransfer: Weakly-supervised domain adaptation for forgery detection. arXiv preprint arXiv:1812.02510.

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.

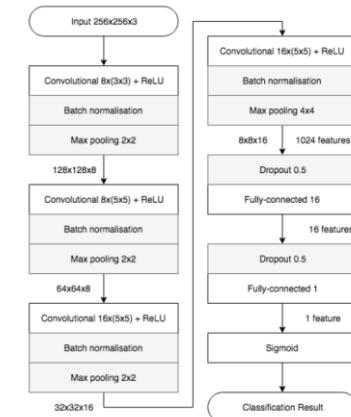
Network	Deepfake classification score		
Class	forged	real	total
Meso-4	0.882	0.901	0.891
MesoInception-4	0.934	0.900	0.917

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

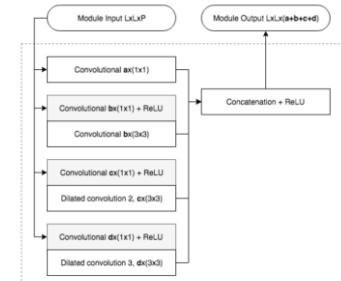
Network	Aggregation score	
Dataset	Deepfake	Face2Face (23)
Meso-4	0.969	0.953
MesoInception-4	0.984	0.953

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

Model Architecture



Inception Variant



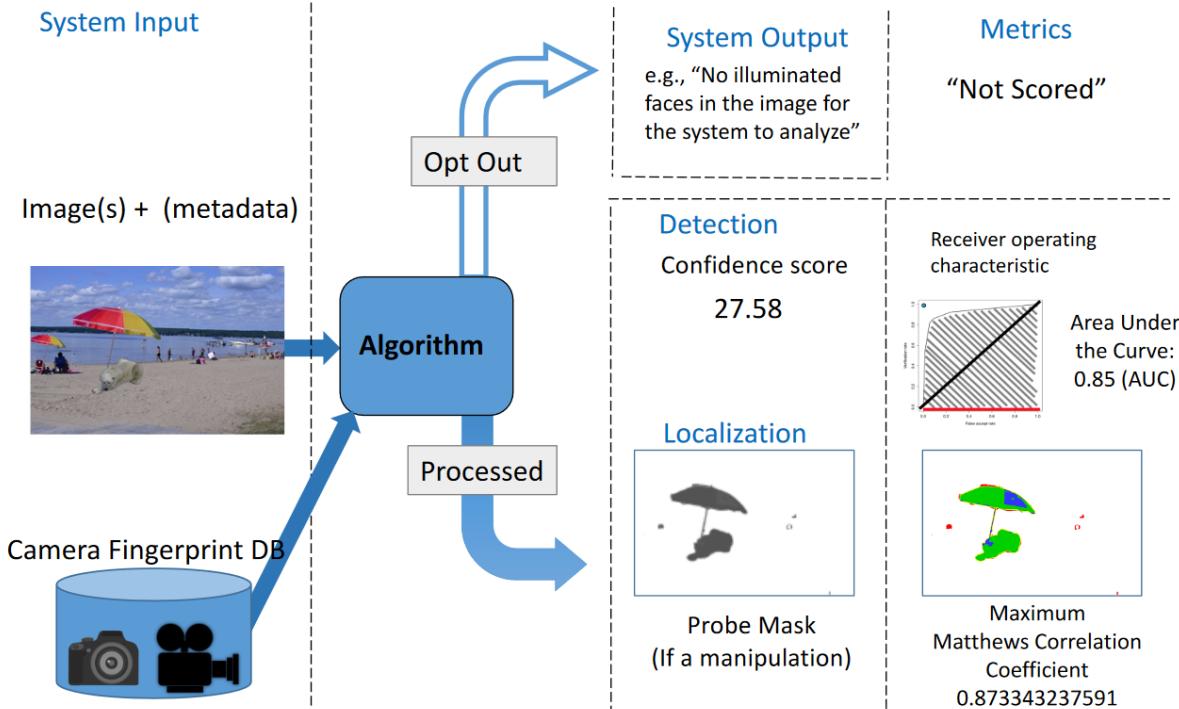
Nimble Challenge

- 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



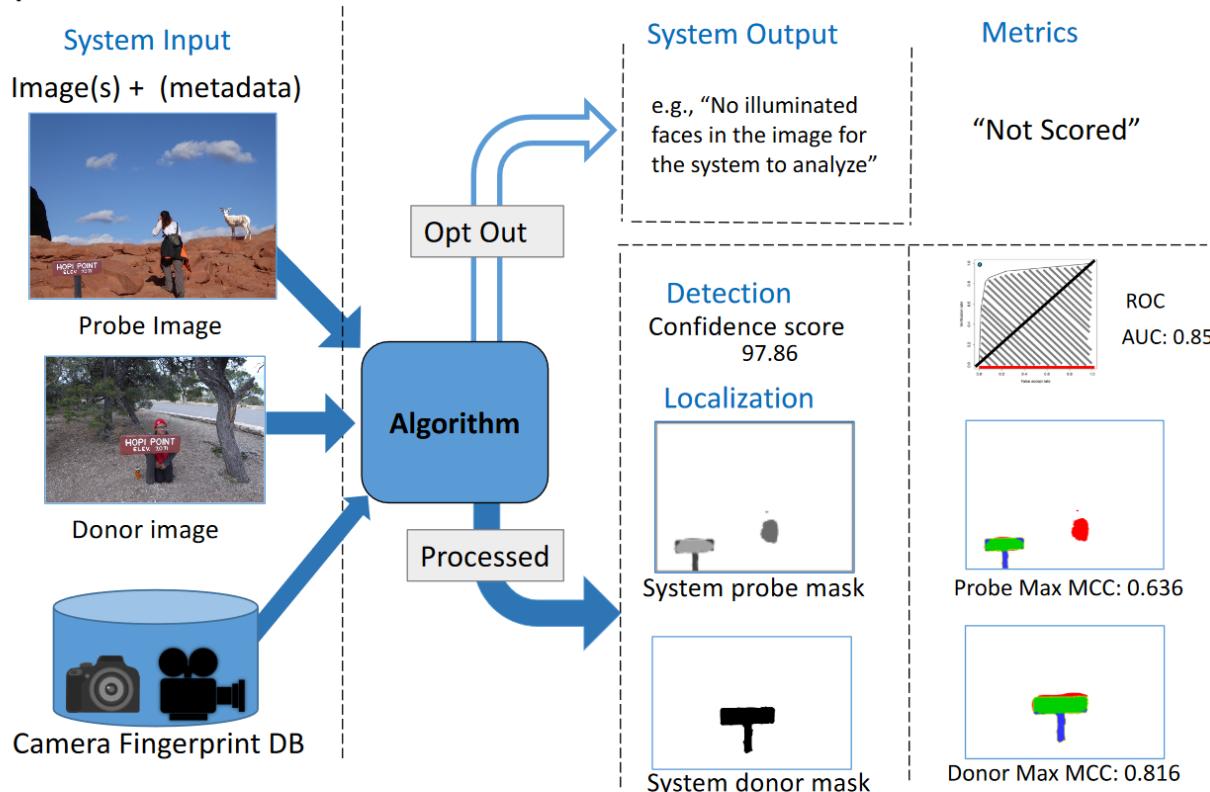
Nimble Challenge 2017 Evaluation - Slides: for the 1 year PI meeting
https://www.nist.gov/sites/default/files/documents/2017/09/05/25_medifor_july17pimeeting-merge_v12.pdf

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

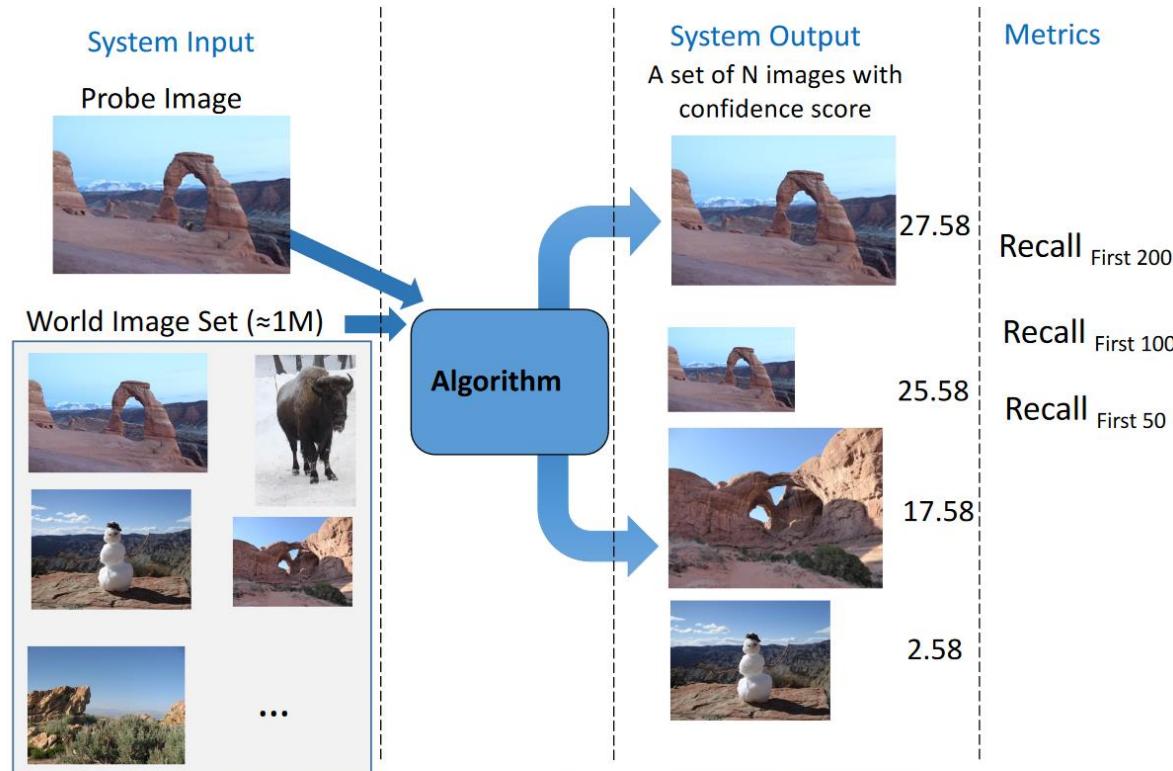
Splice Detection and Localization Evaluation Task



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

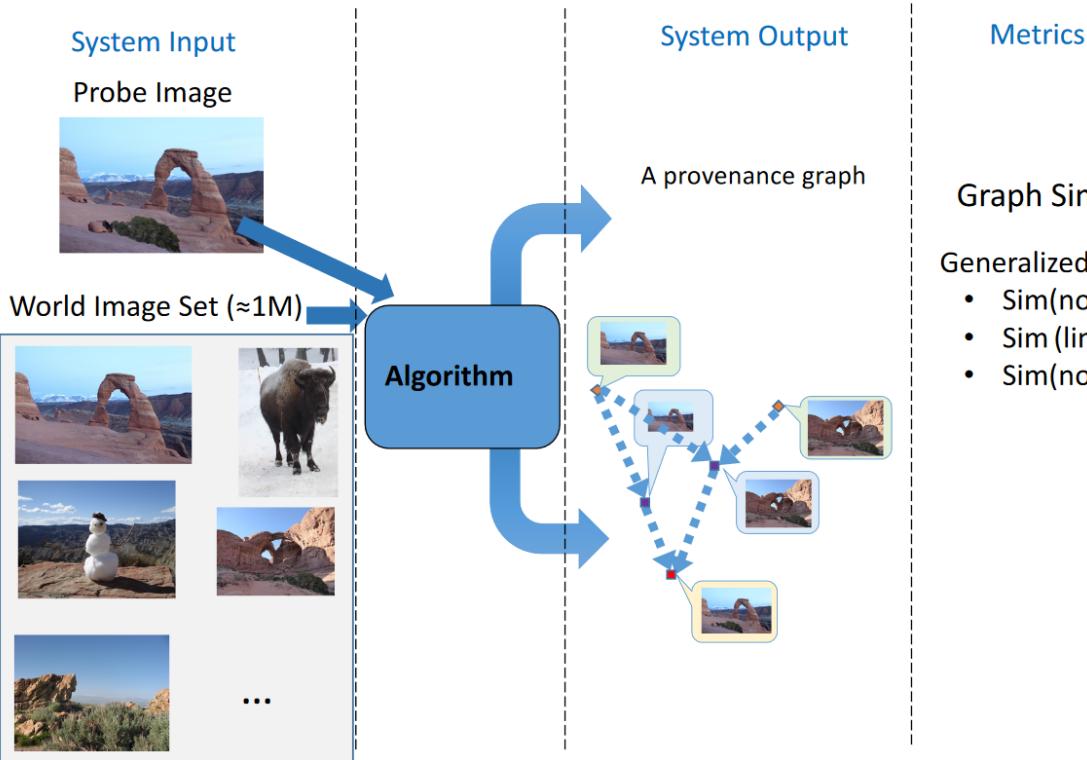
Nimble Challenge 2017

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_medinor_july17pimeeting-merge_v12.pdf

Provenance Graph Building Evaluation Task



Graph Similarity

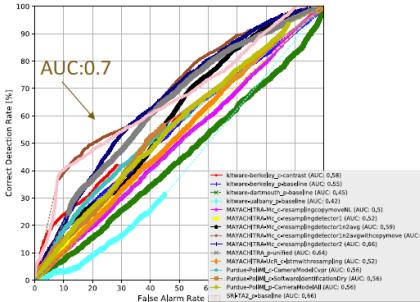
Generalized F-measure:

- Sim(nodes)
- Sim(links)
- Sim(nodes+links)

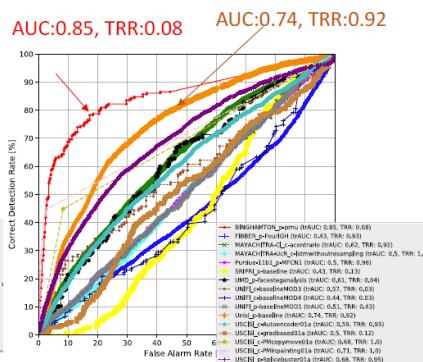
Results

NC17 Image Manipulation Detection Results

Systems That Processed All Probes



Systems that "Opted In" to Process Some of the Probes



7/28/17



National Institute of Standards and Technology / U.S. Department of Commerce

16

Provenance Graph Evaluation Example: ND-Purdue, Baseline System

	Mean Similarity			
	Mean Node Recall	Node Overlap	Link Overlap	
	0.778	0.778	0.375	0.588



Image Legend

- Wide Green image border - The Probe image.
- Green image border - Correctly included image.
- Red image border - False alarm image.
- Grey image border - Omitted provenance image (missed detection).

Link Legend

- Green link - Correctly linked images.
- Red link - False alarm link.
- Grey link - Omitted link.

7/28/17



National Institute of Standards and Technology / U.S. Department of Commerce

29

AUDIO TAMPERING



Fake Audio

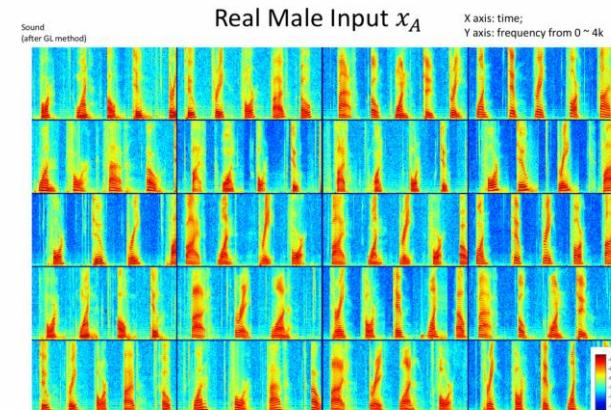
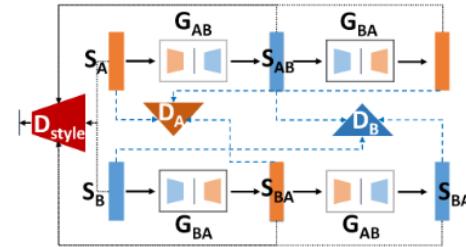
- **No Deep Fake for Audio!!!**
 - Deep Fakes with voice require
 - Traditional audio manipulation
 - Impressionist
- **Approaches to style-transfer in audio**
 - Not promising



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

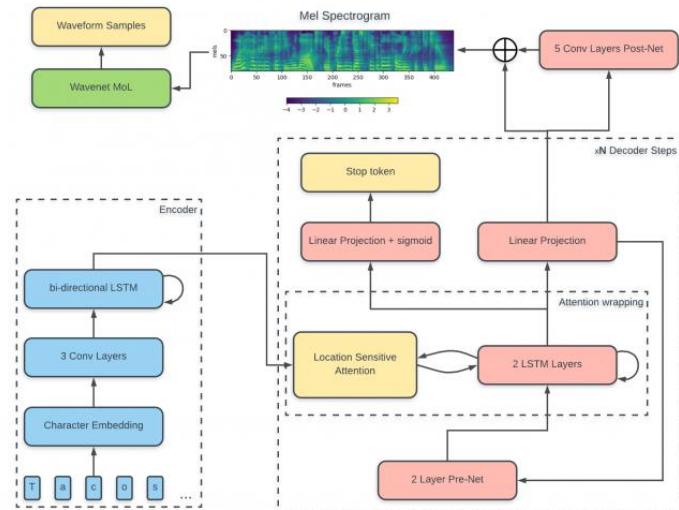


Gao, Y., Singh, R., & Raj, B. (2018, April). Voice impersonation using generative adversarial networks. In 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP) (pp. 2506-2510). IEEE.

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



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):

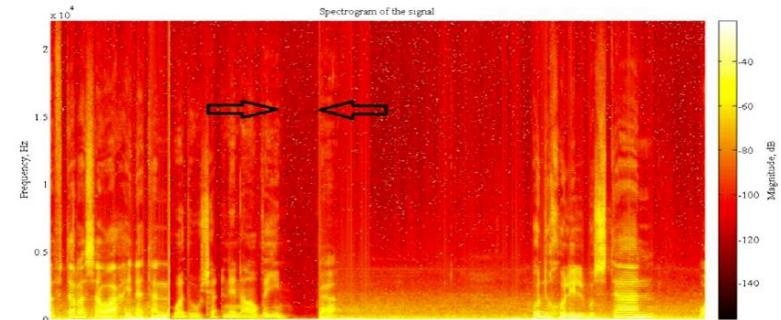
```
ffmpeg -i input.mkv -r 16 -filter:v "setpts=0.25*PTS" output.mkv
```

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

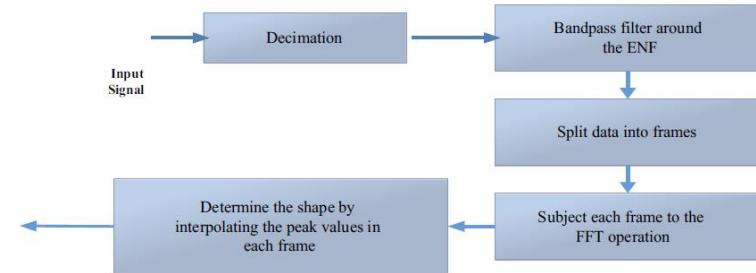
- Motivation
 - Identify modified or fabricated content
 - E.g. Forged / faked evidence in court cases
- Common Modifications
 - Splicing, copying, moving, insertions
- Common Countermeasures
 - Local noise level estimation: Splicing → different noise levels
 - Exploring pitch similarity: Copy-move
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



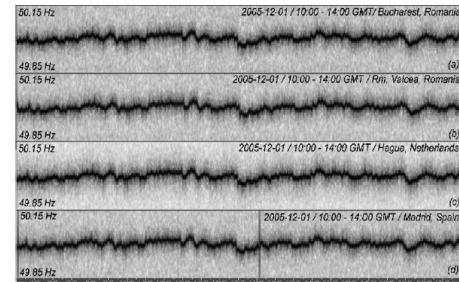
Zakariah, M., Khan, M. K., & Malik, H. (2018). Digital multimedia audio forensics: past, present and future. *Multimedia tools and applications*, 77(1), 1009-1040.

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.



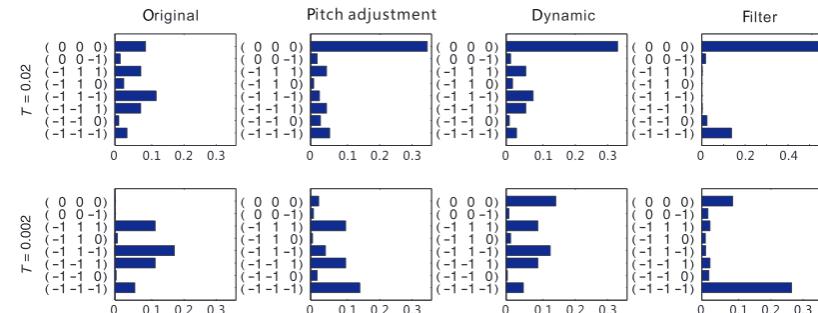
Grigoras, C. (2007). Applications of ENF criterion in forensic audio, video, computer and telecommunication analysis. *Forensic science international*, 167(2-3), 136-145.

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



Content based approaches

- **Text representation**
 - Word frequency, Syntax, punctuation marks
 - Linguistic cues
 - Lie: you, he, she, ...
 - Truth: I ...
- N-gram based approaches: promising
- Deep Learning based approaches: more promising (e.g. Bert)

Knowledge Graph based approaches

- **Fact Checking**
 - Use Knowledge-Graphs (e.g. Dbpedia)
 - Most reliable approach
- **Challenges**
 - Fact Identification, verification, correction
 - Knowledge-base construction
 - Manual fact-checking expensive

BBC NEWS
UK
BBC News (UK)
@BBCNewsUK

BREAKING: Buckingham Palace announces the death of Queen Elizabeth II at the age of 90. Circumstances are unknown. More to follow.



RETWEETS 236 LIKES 59
10:47 AM - 29 Dec 2016
46 236 59 ***

What is 'fake news,' and how can you spot it? Try our quiz
<https://www.theglobeandmail.com/community/digital-lab/fake-news-quiz-how-to-spot/article33821986/>

FAKE NEWS PROBLEM

Overestimated?

Technical Problem?

Solvable?



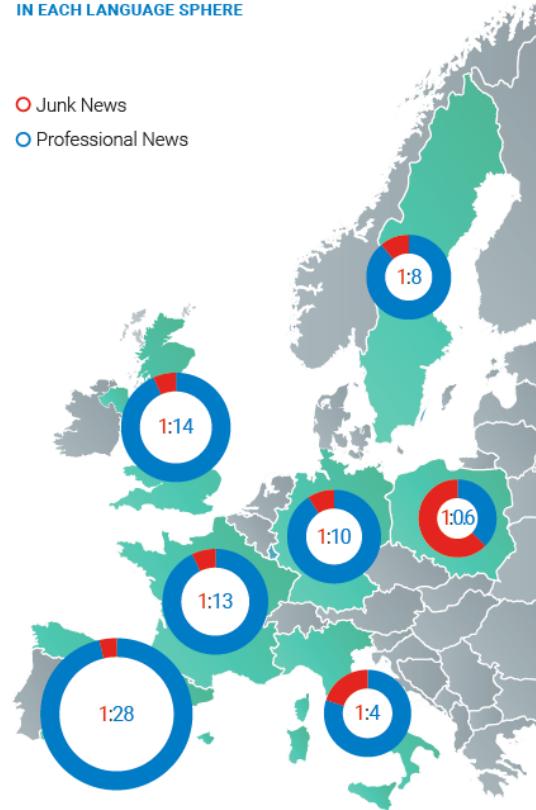
Conclusions

- **General**
 - Fake News detection currently a hyped topic
 - Fear mongering to create business opportunities
- **„Traditional“ attacks**
 - Image / Video / audio 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

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

Media Literacy for Citizenships

- NGO
- Research and projects
- Promote media literacy



EAVI is a non-profit organisation based in Brussels. We work in Europe (and beyond) to empower individuals to be active, engaged citizens in today's increasingly challenging media environment.



30.10.2019

BEYOND 'FAKE NEWS'

10 TYPES OF MISLEADING NEWS

propaganda	partisan	IMPACT
adopted by governments, corporations and non-profits to manage attitudes, values and knowledge appeals to emotions can be beneficial or harmful	ideological and includes interpretation of facts but may claim to be impartial privileges facts that conform to the narrative whilst forging others emotional and passionate language	neutral low medium high
clickbait	conspiracy theory	MOTIVATION
eye catching, sensational headlines designed to distract often misleading and content may not reflect headline drives ad revenue	tries to explain simply complex realities as response to fear or uncertainty not falsifiable and evidence that refutes the conspiracy is regarded as further proof of the conspiracy rejects experts and authority	money politics/power humour/fun passion (mis)inform
sponsored content	pseudoscience	
advertising made to look like editorial potential conflict of interest for genuine news organisations consumers might not identify content as advertising if it is not clearly labeled	purveyors of greenwashing, miracle cures, anti-vaccination and climate change denial misrepresents real scientific studies with exaggerated or false claims often contradicts experts	
satire and hoax	misinformation	
social commentary or humour varies widely in quality and intended meaning may not be apparent can embarrass people who confuse the content as true	includes a mix of factual, false or partly-false content intention can be to inform but author may not be aware the content is false false attributions, doctored content and misleading headlines	
error	bogus	
established news organisations sometimes make mistakes mistakes can hurt the brand, offend or result in litigation reputable orgs publish apologies	entirely fabricated content spread intentionally to disinform guerrilla marketing tactics; bots, comments and counterfeit branding motivated by ad revenue, political influence or both	
DIG DEEPER...		
false attribution	authentic images, video or quotes are attributed to the wrong events or person	
counterfeit	websites and Twitter accounts that pose as a well-known brand or person	content does not represent what the headline and captions suggest
	doctored content	content, such as statistics, graphs, photos and video have been modified or doctored

N.B. The impact and motivation assignments are not definitive and should just be used as a guide for discussion

eavi
MEDIA LITERACY
for CITIZENSHIP
www.eavi.eu



Conclusions

- **Technological Solutions**

- **Requested:**

- Tools to identify/score news items
- Tools to handle waves of fake news (e.g. during elections)

- **Promoted:**

- End-to-end tools
- Filters
- One Tool to rule them all

- **Suggested:**

- Tools to assist journalists in
 - fact-checking
 - Research
 - Quality journalism

Social Responsibility

```
if model.predict(article) == "regime critical":
    delete(article)
    report(user)
```



DER STANDARD · Web

Meinungsfreiheit unter Erdogan

01.09.2019, 16:56 Uhr

Türkei: Nun sind auch Streaming-Dienste und Internetmedien abgekoppelt

DER STANDARD · Web > Netzpolitik

ZENSUR

Türkei: Neue Regelung ermöglicht Zensur von Onlineinhalten

Onlinemedien können durch die Regelung künftig censuriert werden. Besonders Medien, die im Internet eine Plattform gefunden haben, sind betroffen

2. August 2019, 15:20 | 10 Postings

ZEIT ONLINE

Politik Gesellschaft Wirtschaft Kultur • Wissen Digital Campus • Arbeit Entdecken Sport ZEITmagazin Podcasts mehr •

Suche

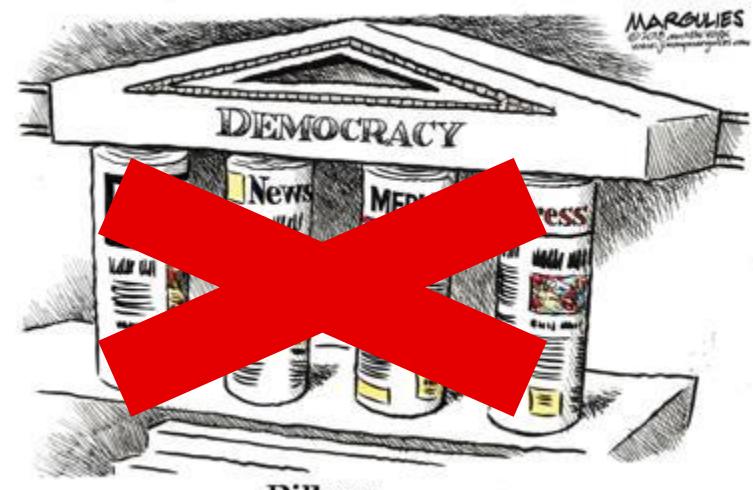
Internationale

Türkei sperrt oppositionelle Websites und Internetkonten

Ein Gericht hat die Schließung regierungskritischer Inhalte im Internet angeordnet. Zusätzlich wurde eine neue Regelung zur Zensur von Onlineplattformen eingeführt.

6. August 2019, 10:02 Uhr / Quelle: ZEIT ONLINE, dpa, AFP und / 207 Kommentare

Flage mit einem Abbild von Präsident Recep Tayyip Erdogan © Chris McGrath/Getty Images



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:1901.06437v1 [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 research has focused on the detection of fake news, there has also been significant work on mitigating 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 present a comprehensive review of the state-of-the-art methods for three main 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 summarize characteristic features 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, Feng Qian, He Jiang, Natalia 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 not been faced with before.”

Authors’ addresses: Karishma Sharma, ksharma@usc.edu, University of Southern California; Feng Qian, University of Southern California, fqliang@usc.edu; He Jiang, jianghe@usc.edu; Natalia Ruchansky, University of Southern California, nruhansky@usc.edu; Ming Zhang, Peking University, mzhang@pku.edu.cn; Yan Liu, University of Southern California, yanliu@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.

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

111

ACM Trans. Intell. Syst. Technol., Vol. 10, No. 4, Article 111. Publication date: August 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 focused on the identification and mitigation of fake news on social media platforms for news recommendation modeling. Majority of the research has primarily focused on two broad categories of fake information: opinion-based (e.g., fake reviews, fake news) and fact-based (e.g., fake news, fake reviews). This survey aims to provide a comprehensive overview 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) mitigation techniques, and (vi) detection methods. We also propose a unified framework to describe these recent methods and highlight a number of important directions for future research.

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

ACM Reference Format:
 Kumar, S., and Shah, N. 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 ground-breaking phenomena such as the rise of social media, it has also led to the rise of fake news, fake reviews, and fake 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 increasingly believe fake news [5, 6, 7]. The rise of fake news has also led to significant economic losses and the increasing 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, second only to terrorism [27]. Therefore, it is important to understand the nature of fake information and the damage it creates 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

“fake news” at the top level, this survey will cover the look below.

Authors’ addresses: Srijan Kumar, Computer Science, Stanford University, USA, wjnjnq@cs.stanford.edu; Neil Shah, Computer Science, Carnegie Mellon University, USA, neilshah@cs.cmu.edu

Permissions 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.

© 2018 Association for Computing Machinery.
<https://doi.org/10.1145/3196403>

Thank you!

Alexander Schindler, 29.10.2019

