# Semantic Forensics (SemaFor)

William Corvey, PhD
DARPA I2O

Deepfakes Day 2022

8/30/2022

# Image & video manipulation technology evolution



**Image Synthesis**

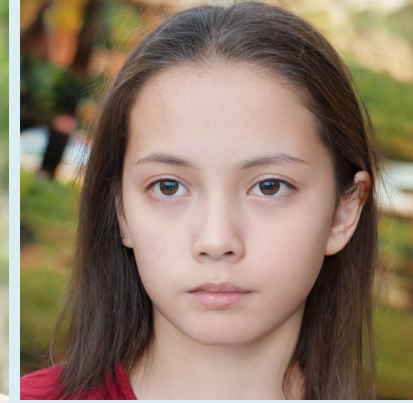**2014** Goodfellow et al.     **2015** Radford et al.     **2016** Liu and Tuzel     **2017** Karras et al.     **2018** Karras et al.     **2019** Karras et al.
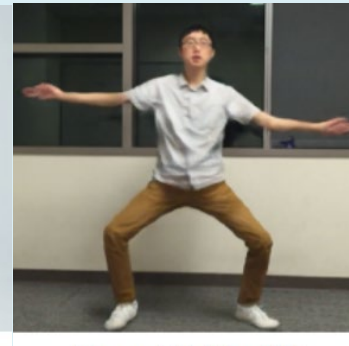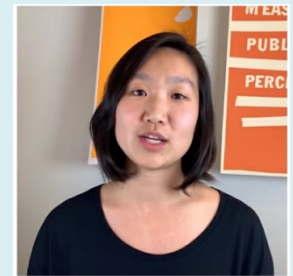
**Video Manipulation**

Target     Source

Source Subject     Target Subject 2

..stock price closed at one hundred ~~ninety one point four~~ eighty two point two five dollars

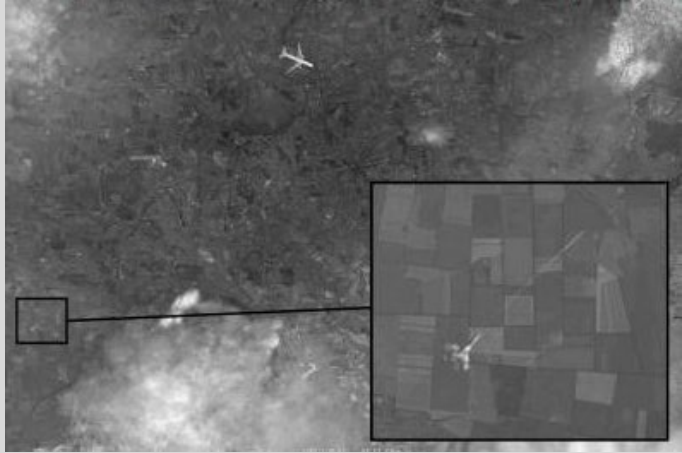**2016** Thies et al.     **2017** Deepfakes     **2018** Chan et al.     **2019** Fried et al.

# Image & video manipulations we've already seen

**Manual manipulation**

## MH17 Downing (2014)



*Photoshop*　　Source: Bellingcat

## Catalan Independence (2017)



*Photoshop*　Source: El Pais

## Chinese government social media (2020)



*Photoshop*　Source: Twitter

https://twitter.com/zlj517/status/1333214766880888448

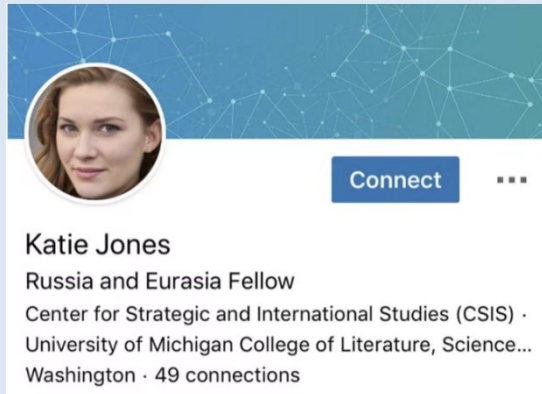**Automated manipulation**

## Jordan Peele's Obama (2018)



*Deepfake*　　Source: YouTube

## Katie Jones – LinkedIn (2019)



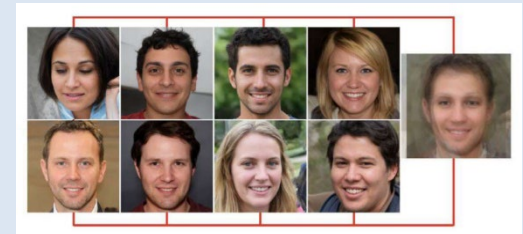*GAN*　　Source: CNET

## Nixon Moon Speech (2020)



*Multiple*　Source: moondisaster.org

## Pro-Chinese Inauthentic Network (2020)



*GAN*　　Source: Graphika

# Future threats from falsified multi-modal media

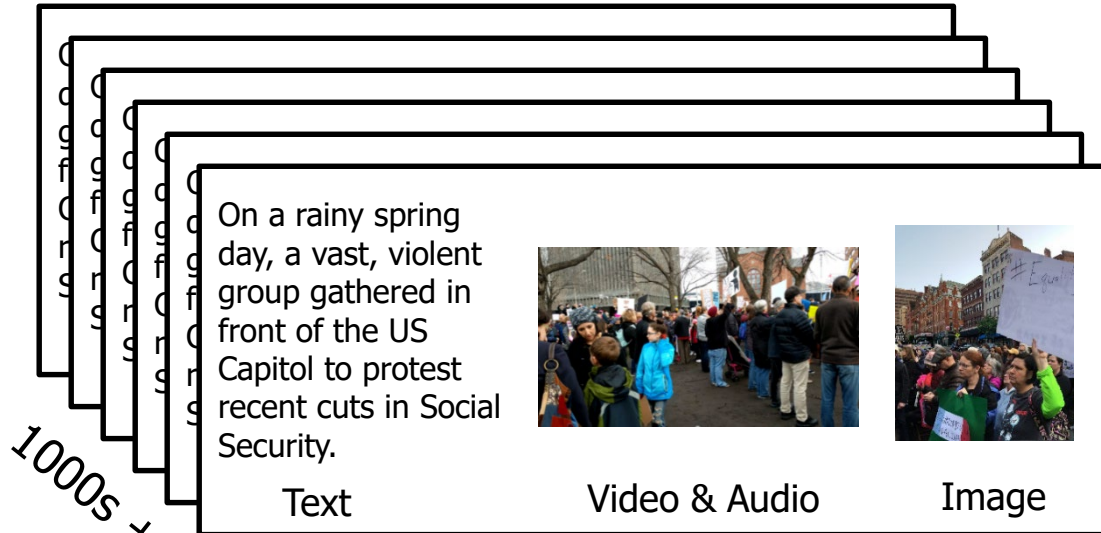## Targeted Personal Attacks

Peele 2017



AI Multimedia Algorithms

⬇



Highly realistic video

## Generated Events at Scale

AI Multimedia Algorithms

⬇

1000s ×

On a rainy spring day, a vast, violent group gathered in front of the US Capitol to protest recent cuts in Social Security.

Text

Video & Audio

Image

Believable fake events

## Ransomfake concept: Identity Attacks as a service (IAaaS)

Bricman 2019

AI Multimedia Algorithms

⬇

Forged Evidence

⬇

Identity Attacks

Examples of possible fakes:
• Substance abuse
• Foreign contacts
• Compromising events
• Social media postings
• Financial inconsistencies
• Forging identity

## Undermines key individuals and organizations

Create rich semantic algorithms that automatically detect, attribute, and characterize falsified multi-modal media to defend against large-scale, automated disinformation attacks

# Synthetic Media Detection, Attribution, and Characterization Capabilities

| | Desired Capability | Today | SemaFor |
|---|---|---|---|
| **Detection** | Automatically detect semantic generation/manipulation errors | Limited | Yes |
| | Detect manipulations across multiple modalities and assets | Limited | Yes |
| | Robust to many manipulation algorithms | Fragile | Highly robust |
| | Increased adversary effort needed to fool detection algorithms | Some | Significant |
| **Attribution** | Automatically confirm source or author | Limited | Yes |
| | Automatically identify unique source fingerprints | No | Yes |
| | Explain authorship inconsistencies | No | Yes |
| **Characterization** | Automatically characterize manipulation intent or impact | No | Yes |
| | Provide evidence and explanation for manipulation intent | No | Yes |
| | Correctly prioritize generated/manipulated media for review | No | Yes |

# Semantic Detection

**Text (Notional)**

*NewsWire: April 1, 2019, Bob Smith*
On a rainy spring day, a vast, violent group gathered in front of the US Capitol to protest recent cuts in Social Security.

**Audio (Notional)**

"We'd like to welcome you here on this beautiful spring day. Thank you all for coming out [cheering]…"

**Image**



**Video**

**Text (Notional)**

*NewsWire: April 1, 2019, Bob Smith*
On a rainy spring day, a vast, violent group gathered in front of the US Capitol to protest recent cuts in Social Security.
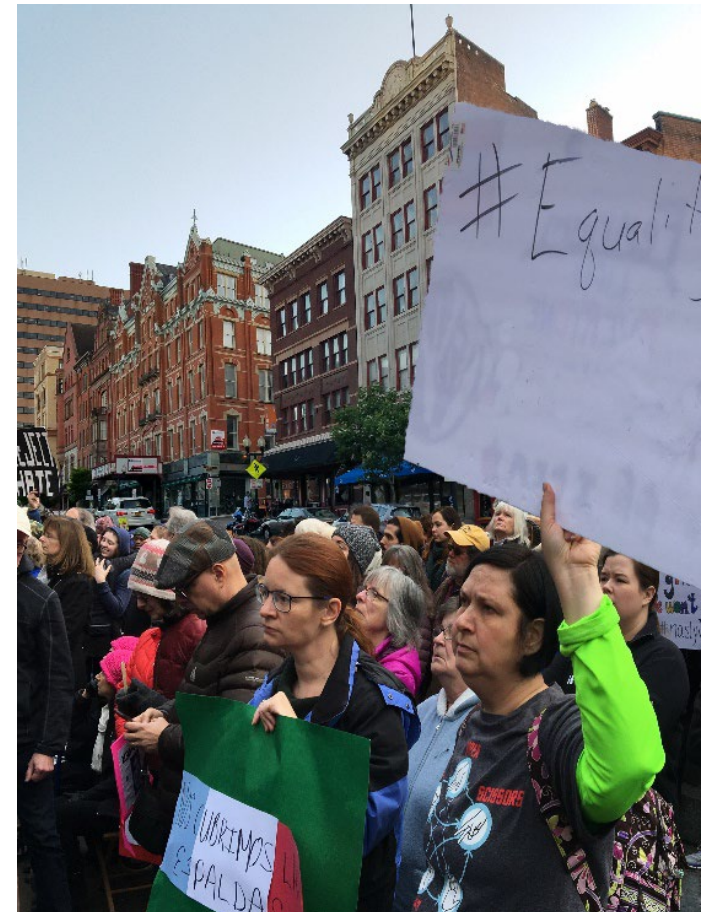
**Audio (Notional)**

"We'd like to welcome you here on this beautiful spring day. Thank you all for coming out [cheering]…"

Conclusion: Media components consistent across modalities.

**Image**

**Video**

"protest"

**Text (Notional)**

*NewsWire: April 1, 2019, Bob Smith*
On a rainy spring day, a vast, violent group gathered in front of the US Capitol to protest recent cuts in Social Security.

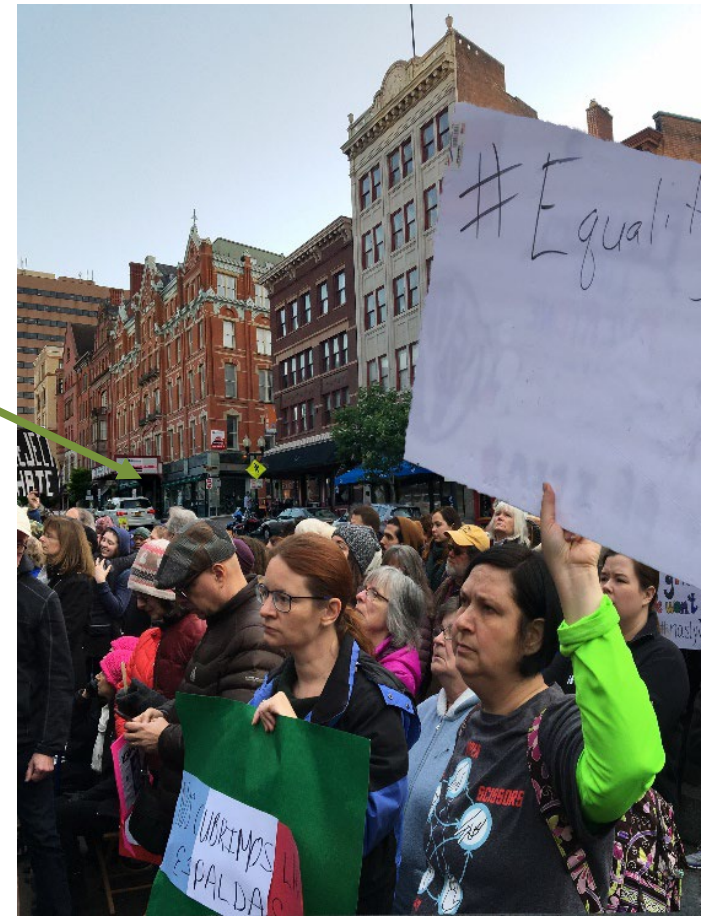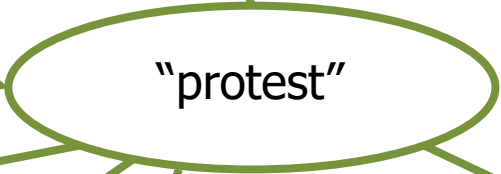**Audio (Notional)**

"We'd like to welcome you here on this beautiful spring day. Thank you all for coming out [cheering]…"

Conclusion: Media components not consistent across modalities.

**Image**

**Video**

"violent group"

# Semantic Attribution & Characterization

## Text (Notional)

*NewsWire: April 1, 2019, Bob Smith*

On a rainy spring day, a vast, violent group gathered in front of the US Capitol to protest recent cuts in Social Security.

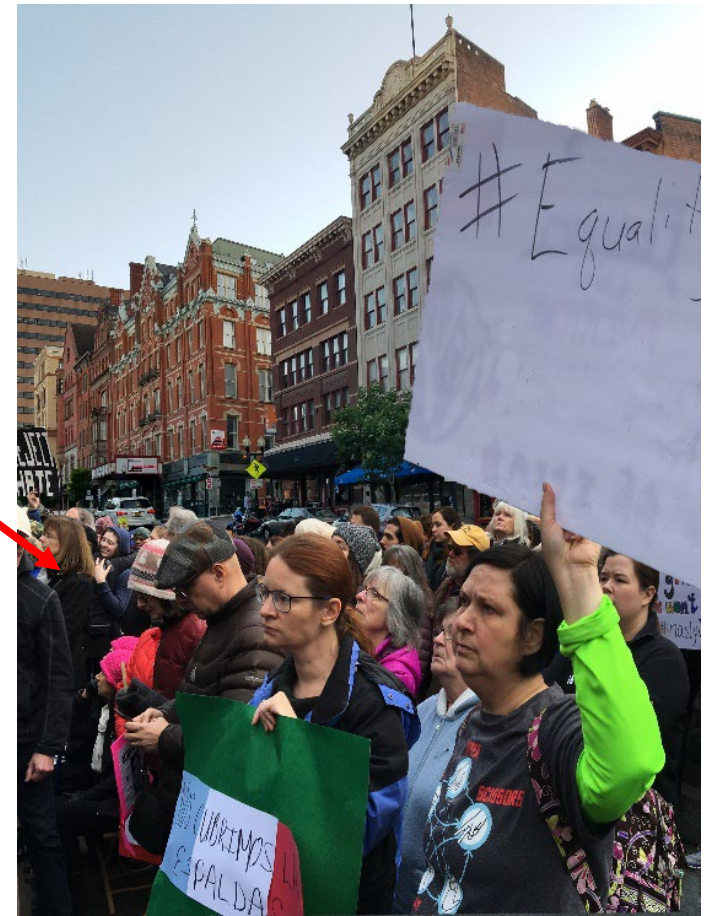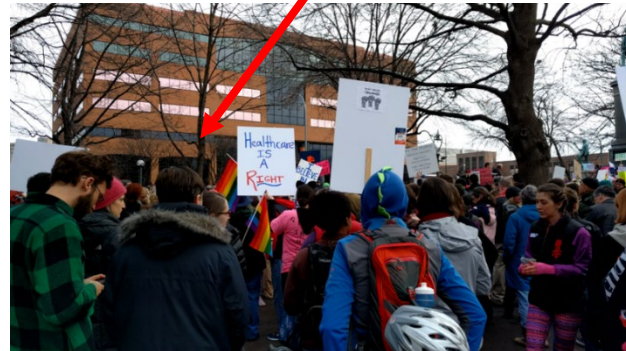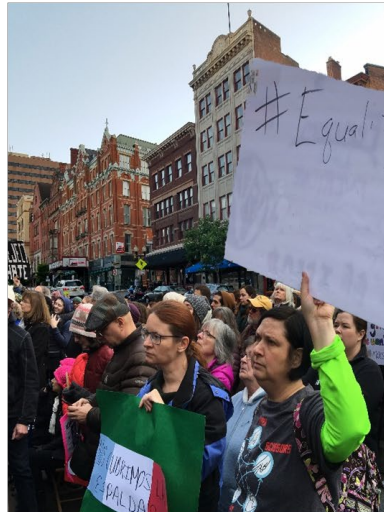**Video**

**Audio (Notional)**

"We'd like to welcome you here on this beautiful spring day. Thank you all for coming out [cheering]…"

**Image**



## Attribution: Incorrect
- Bob Smith is a tech reporter, doesn't report on social events
- Vocabulary indicates different author
- NewsWire has a different style for use of images in news article
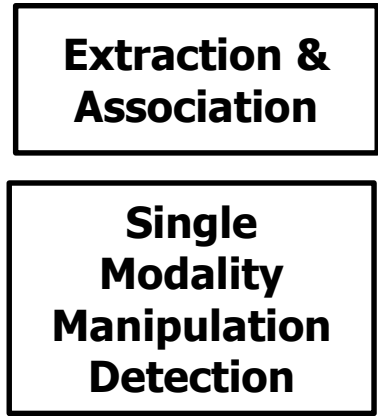
## Characterization: Malicious
- Large number of inconsistencies across media
  - Environment – "rainy spring day"
  - Behavior – "violent group"
  - Location – "US Capitol"
  - Topic – "Social Security"
- Use of unsupported term "violent"
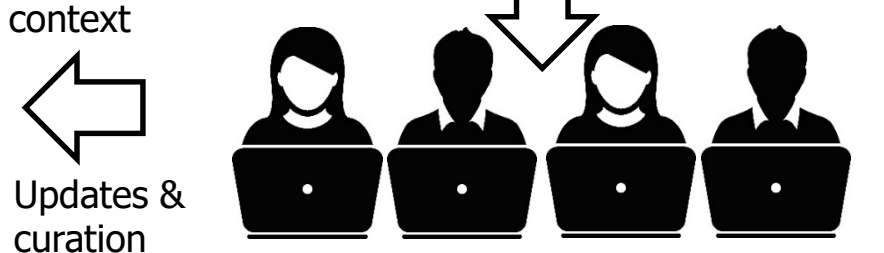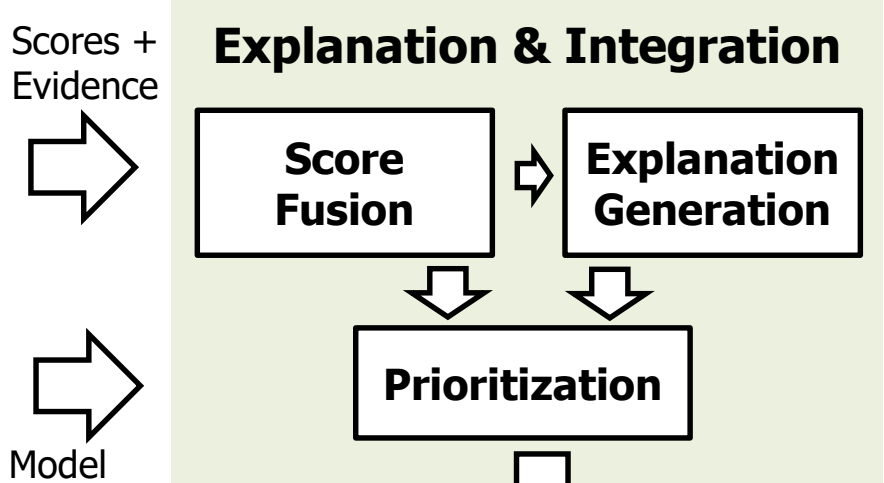- Failed sourcing to high credibility organization ("*NewsWire*")

# SemaFor System

**Multimedia:**
Text,
Audio,
Images,
Video,
Source metadata

**Extraction & Association**

**Single Modality Manipulation Detection**

**Multimodal Representations**

Attributed Graphs    Semantic Embeddings    Hybrid Representations    ...

**Reasoning Ensembles**

Representations

Multiple Pipelines

**Semantic Detection**

**Attribution**

**Characterize**

Scores + Evidence

**Explanation & Integration**

**Score Fusion**

**Explanation Generation**

**Prioritization**

**Semantic Models**

**Generator Models**

AI generator failure modes

Entity detection & association performance

**Extraction Models**

**Source Models**

Modality & cross-modal styles, topic models

Polarization, virality, impact

**Intention Models**

Model context

Updates & curation

# Technical Areas

**TA1: Detection, Attribution, Characterization**

**Extraction & Association**

**Single Modality Manipulation Detection**

**Multimodal Representations**

Attributed Graphs   Semantic Embeddings   Hybrid Representations   ...

Multiple Pipelines

**Reasoning Ensembles**

Representations

**Semantic Detection** → **Attribution** → **Characterize**

**Semantic Models**

Generator Models

AI generator failure modes

Entity detection & association performance

Extraction Models

Source Models

Modality & cross-modal styles, topic models

Polarization, virality, impact

Intention Models

**TA2: Explanation & Integration**

Score Fusion → Explanation Generation

Prioritization

**TA3: Evaluation**

Media generation   Evaluations

Metrics

**TA4: Challenge Curation**

SOTA challenges   Threat modeling

**Adversary Resources**
e.g. money, compute, time

**Adversary Skill**
e.g. technical capabilities, talent pool

Nation-state, Corporations

Criminal Enterprise

Skilled Networked activists

Low-skill Networked activists

Skilled individual

Unskilled Individual

Hollywood

New GAN algorithm
Eval 2: Puppeteering, Synthetic Audio

GAN retraining

Green screen

Eval 2: CGI Generated Text and Images

Eval Photoshop 1

Deepfake

Eval 2: Attribution and Cheap Fakes

Shallowfake

GAN

★ Organization

○ Technical capability

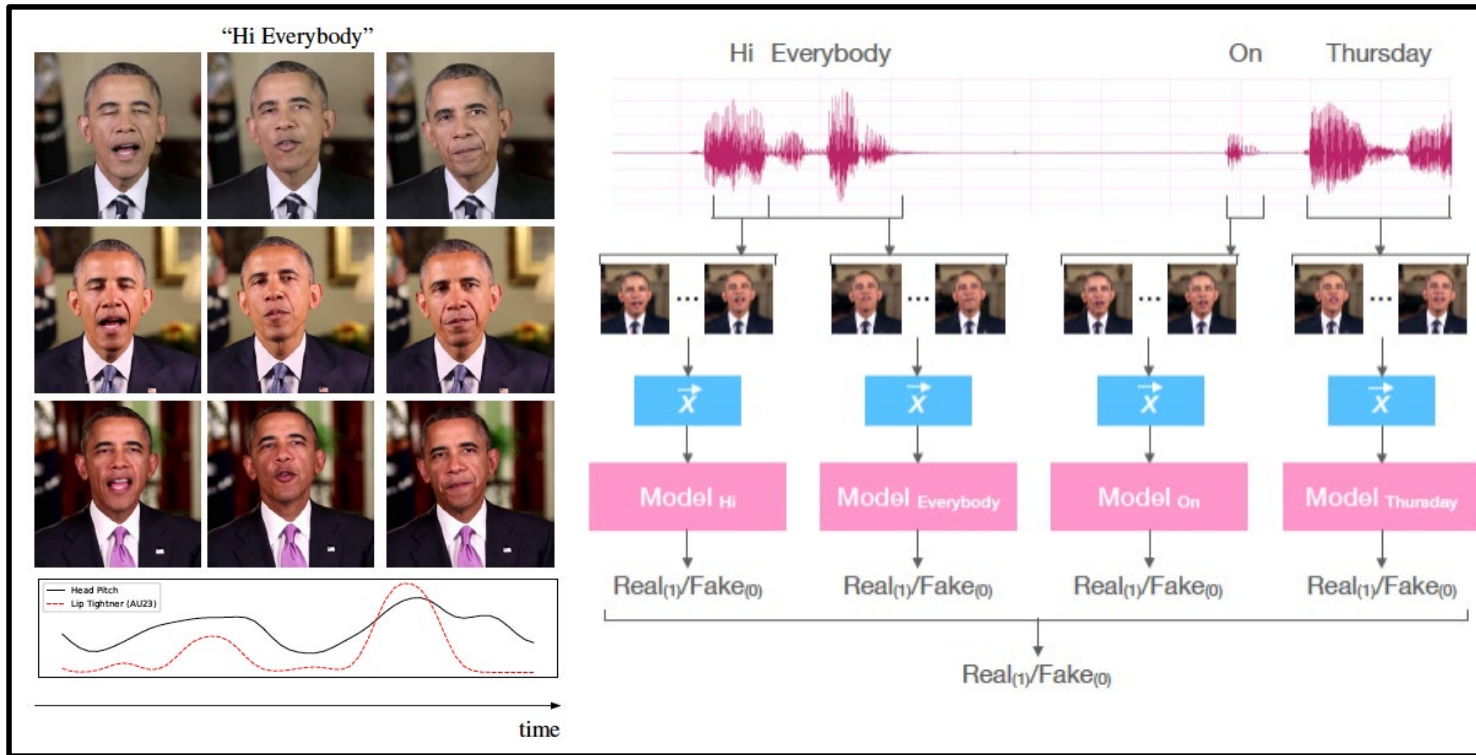Evaluation Mapping

Development of first-of-kind dataset of actual context manipulations, including human annotations of in-the-wild manipulations.

Collaborate with USG disinformation and misinformation community to identify topics and campaigns of interest.

Threat Landscape: A multi-dimensional representation to quantify and categorize the "who, what, why, where and how" of threats using manipulated media.

UCBerkeley
Pinscreen



**AUC on 10s video clips & 5 falsifications**

| | Audio Dubbing | Wav2Lip | Impersonator | FaceSwap | in-the-wild |
|---|---|---|---|---|---|
| Obama | 1.00 | 1.00 | 0.95 | 0.90 | 0.98 |
| Trump | 0.99 | 0.99 | 0.89 | 0.92 | 0.98 |
| Biden | 0.89 | 0.93 | 0.98 | 0.73 | 0.95 |
| Harris | 0.89 | 0.92 | 0.82 | 0.93 | - |
| O' Brien | 0.92 | 0.88 | 0.90 | 0.84 | - |
| Oliver | 0.95 | 0.92 | 0.86 | 0.87 | - |
| Avg | 0.94 | 0.94 | 0.90 | 0.87 | 0.97 |

**Future integration & deployment**



- Biometric-based (person-specific) forensic approach learns a semantic model of an individual's movement & speech
- Can detect both deepfakes and cheapfakes (such as audio dubbing or impersonator)
- Generalizes to falsification methods unseen during training

## Detect misinformation from a cluster of topically related documents

**News 1**

… Maduro was not targeted by the drones, the prime minister said, but state security services reported that the drones were meant for him. "The explosion was caused by two machine guns," Maduro said, adding that there were no injuries. …

**News 2**

… Venezuela's president, Nicolás Maduro, has survived an apparent assassination attempt after what officials described as drones armed with explosives detonated overhead during a speech he was making at a military event. …

**Cross-document KG**

**Within-document KG**

Nicolas Maduro — targeted — guns — explosion — News 1

drones

**Within-document KG**

Nicolás Maduro — detonated — News 2

Venezuela — explosives

Event cluster

**Detector**

GNN → Event features → Real / Fake

Detection results

GNN → Doc features → Real / Fake

**Event-level detector**    **Doc-level detector**

### Detection Accuracy (%)

| Document-level | IED | TL17 | Crisis |
|---|---|---|---|
| HDSF | 78.42 | 80.62 | 82.14 |
| GROVER-median | 79.06 | 79.40 | 86.84 |
| GROVER-mega | 82.90 | 90.00 | 87.13 |
| Ours | **86.76** | **90.21** | **93.89** |

| Event | IED | | TL17 | | Crisis | |
|---|---|---|---|---|---|---|
| | F1 | AUC | F1 | AUC | F1 | AUC |
| Random | 16.31 | 50.44 | 19.44 | 49.65 | 21.70 | 50.41 |
| LR | 31.26 | 77.87 | 29.14 | 68.19 | 31.67 | 68.17 |
| Ours | 44.86 | **88.46** | **41.56** | **82.59** | **48.48** | **85.60** |
| w/o cross-doc | **45.00** | **88.54** | **41.66** | 82.28 | 47.78 | 85.17 |

*Real news*

...She noted that at least 1% of people who catch coronavirus die of it. "Another 10-20% are hospitalized. Another 30% or more have long lasting symptoms. The vaccine is far safer, with only minor temporary side effects," Ranney said on Twitter. …In both Pfizer and Moderna's vaccine trials, no worrying side-effects were seen. …

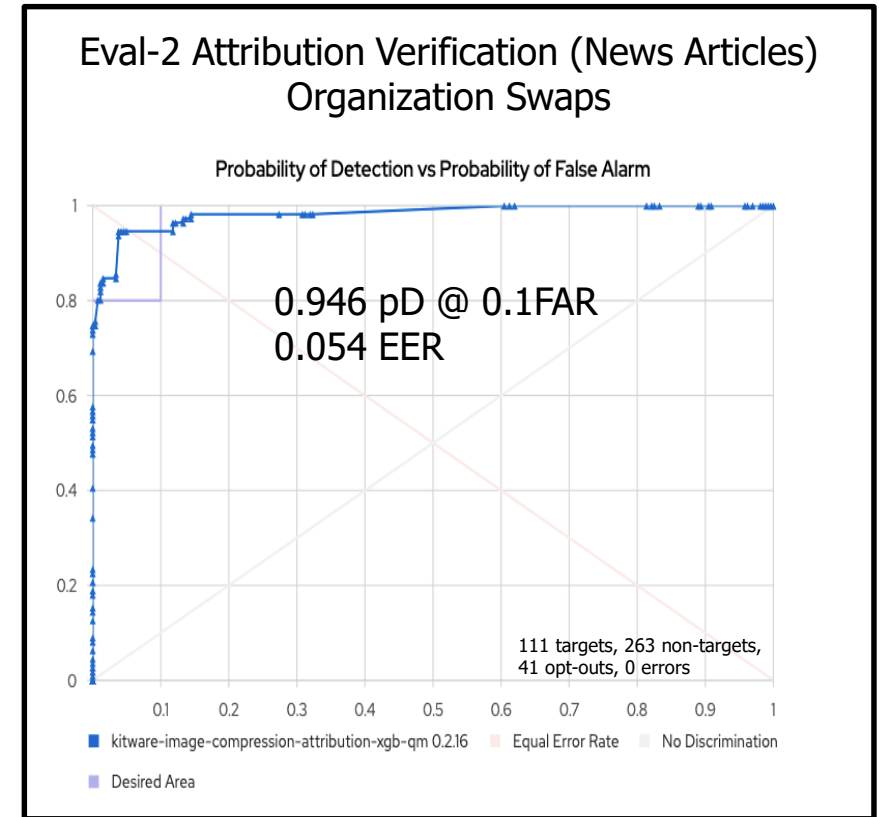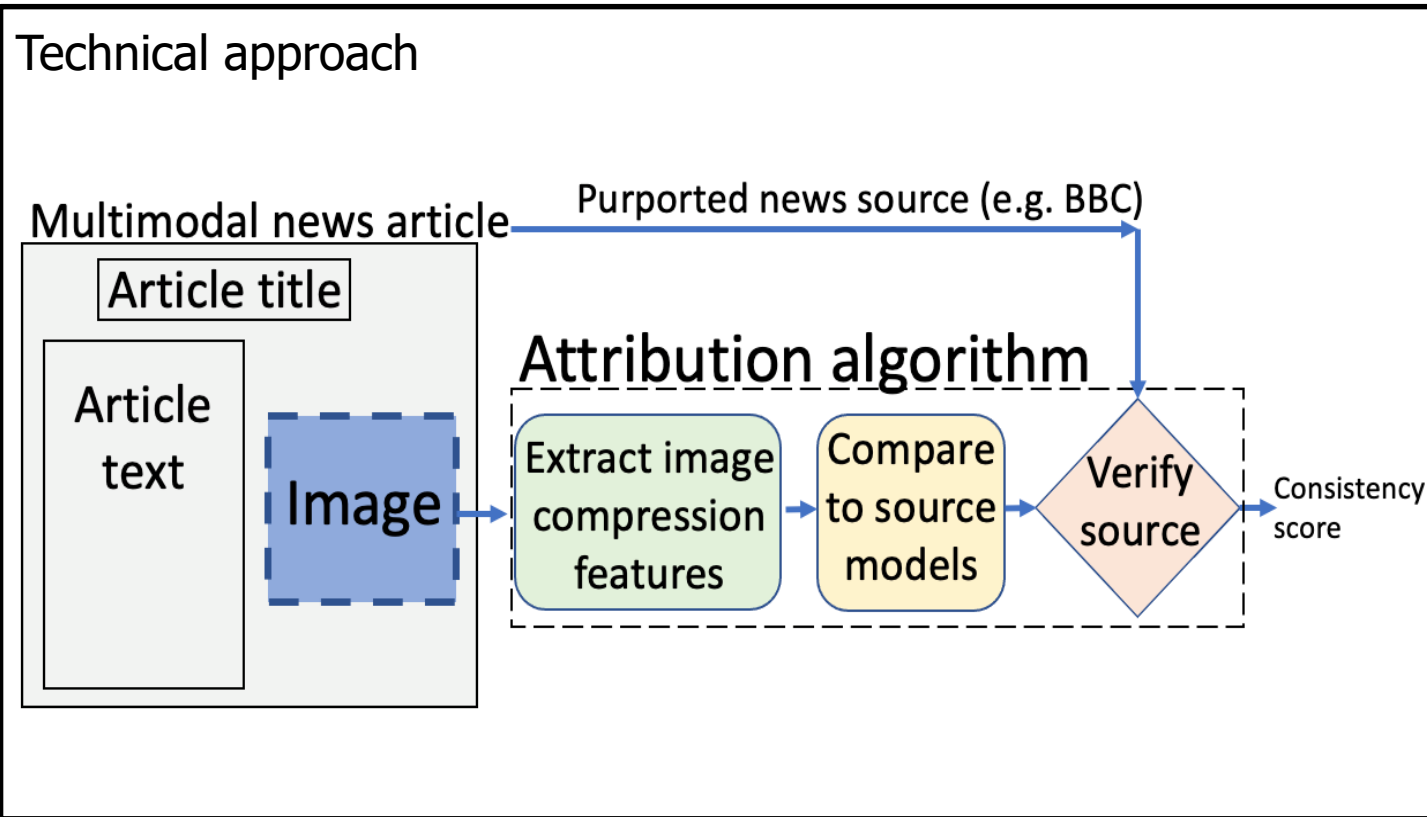**Claim**: COVID vaccine is safer than COVID

⟷ contradictory

*Fake news*

As many as 45,000 people may have died from the mRNA shots being given to halt COVID, according to prominent physician Dr. Peter McCullough. And teenagers — especially boys — are more at risk from being hospitalized from the vaccine than they are for COVID, he said. The culprit is myocarditis, inflammation of the heart. …
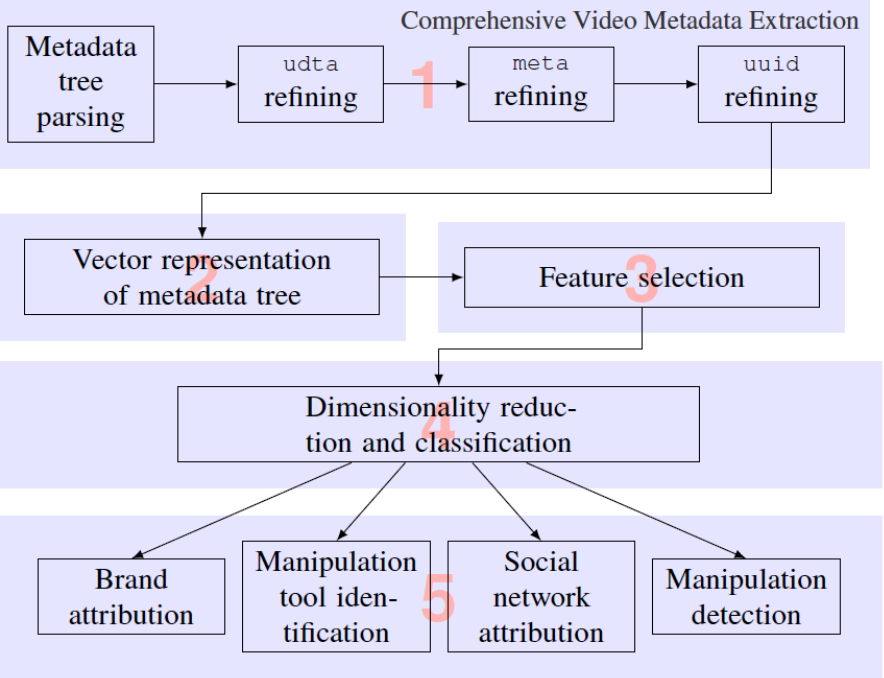
**Claim**: COVID vaccine is more dangerous than COVID

## Unique single-document to cross-document event coreference resolution and reasoning
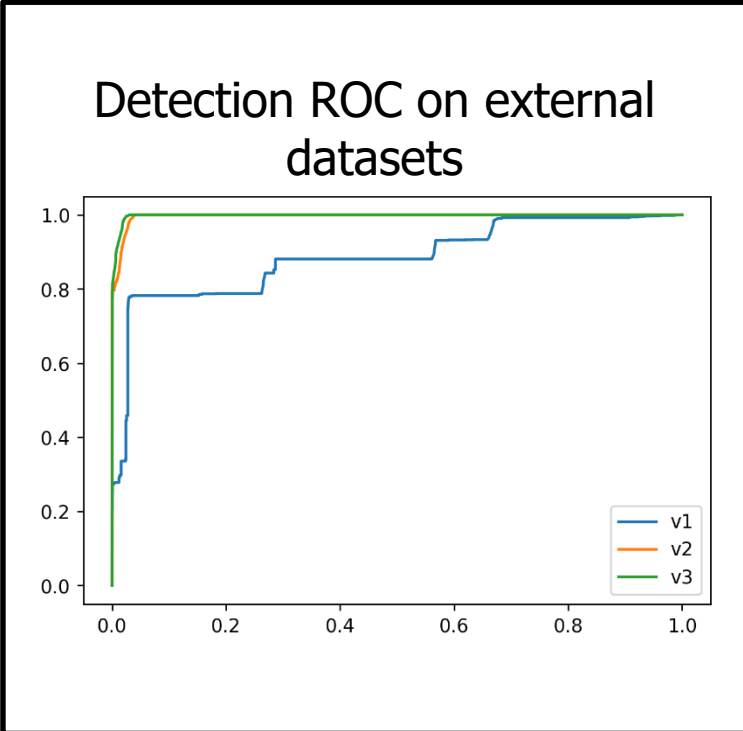
16

# News organization attribution based on compression settings

## Technical approach

Multimodal news article → Purported news source (e.g. BBC)

**Article title**

**Article text** | **Image**

## Attribution algorithm

Extract image compression features → Compare to source models → Verify source → Consistency score

### Eval-2 Attribution Verification (News Articles) Organization Swaps

Probability of Detection vs Probability of False Alarm

0.946 pD @ 0.1FAR
0.054 EER

111 targets, 263 non-targets, 41 opt-outs, 0 errors

kitware-image-compression-attribution-xgb-qm 0.2.16    Equal Error Rate    No Discrimination

Desired Area

- News organizations have highly curated pipelines for media production that leave signatures in the media
- Analytic verifies whether images originate from a purported (news) source or not, by comparing each image's compression settings to known image compression patterns of the source
- Strong performance in SemaFor Eval 2: Attribution Verification (News Articles) Organization Swaps
- Novel forensic analysis of image file compression with low computational cost

Comprehensive Video Metadata Extraction

Detection ROC on external datasets

- Analytics use metadata in the video files to determine if a video is manipulated
- Difficult to hide manipulation traces in metadata
- Method is significantly faster than pixel-based manipulation detection techniques
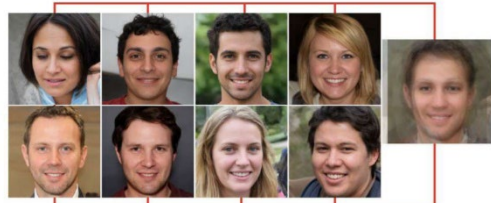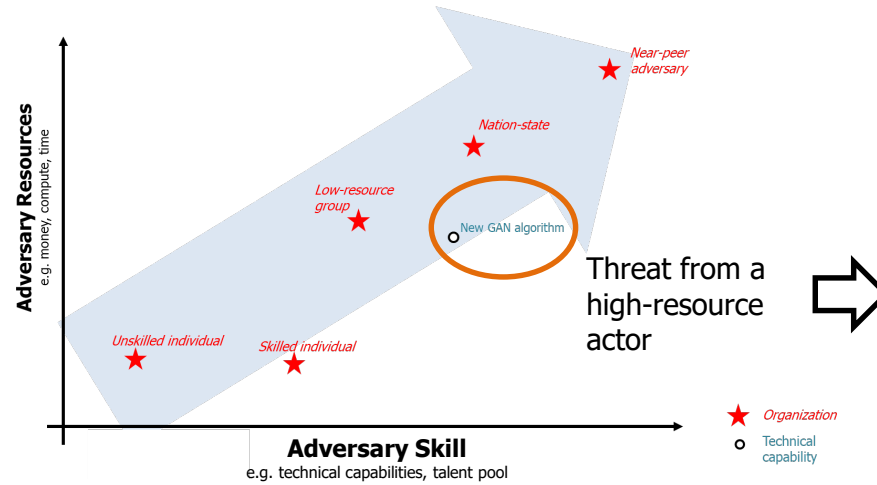
**Challenge:** can detectors identify images from a novel GAN created by a high-resource actor, even without any data from it?
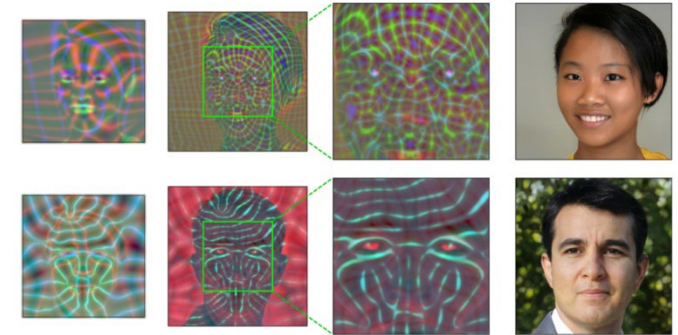
Pro-Chinese Inauthentic Network (2020)



*GAN*
Large-scale sock puppet accounts

Source: Graphika



**Adversary Resources**
e.g. money, compute, time

*Near-peer adversary*

*Nation-state*

*Low-resource group*

New GAN algorithm

*Unskilled individual*

*Skilled individual*

**Adversary Skill**
e.g. technical capabilities, talent pool

★ *Organization*
○ *Technical capability*

Threat from a high-resource actor

Alias-Free Generative Adversarial Networks (StyleGAN3)
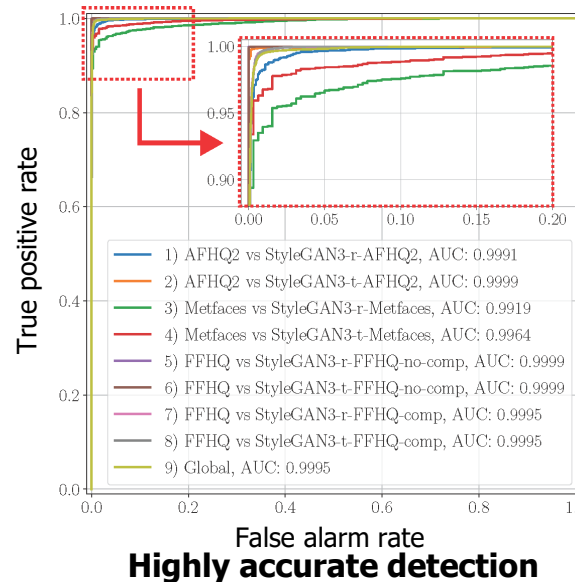Official PyTorch implementation of the NeurIPS 2021 paper



Alias-Free Generative Adversarial Networks
Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, Timo Aila
https://nvlabs.github.io/stylegan3

Training over semantic categories, augmentation, & many GANs



Real

Fake

**No training data from StyleGAN3!**



True positive rate

False alarm rate

1) AFHQ2 vs StyleGAN3-r-AFHQ2, AUC: 0.9991
2) AFHQ2 vs StyleGAN3-t-AFHQ2, AUC: 0.9999
3) Metfaces vs StyleGAN3-r-Metfaces, AUC: 0.9919
4) Metfaces vs StyleGAN3-t-Metfaces, AUC: 0.9964
5) FFHQ vs StyleGAN3-r-FFHQ-no-comp, AUC: 0.9999
6) FFHQ vs StyleGAN3-t-FFHQ-no-comp, AUC: 0.9999
7) FFHQ vs StyleGAN3-r-FFHQ-comp, AUC: 0.9995
8) FFHQ vs StyleGAN3-t-FFHQ-comp, AUC: 0.9995
9) Global, AUC: 0.9995

**Highly accurate detection**



NVlabs / stylegan3-detector  Public

StyleGAN3 Synthetic Image Detection

**Overview**

While new generator models, such as StyleGAN3, enable new media synthesis capabilities, they may also present a new challenge for AI forensics algorithms for detection, attribution, and characterization of synthetic media.

As part of DARPA's Semantic Forensics (SemaFor, for short) program, NVIDIA has been collaborating with digital forensics experts and researchers to help advance the capabilities to verify the authenticity and provenance of synthetic media.

**NVIDIA delayed releasing the GAN software & published the detectors alongside StyleGAN3**

# Detecting & localizing synthetic audio



Human · Synth · Ambiance

## Example of humans + synth voices + background

**Shawn**: Has everyone had a chance to review the data?
**Synth 1**: It's quite disturbing. And fascinating!
**Dave**: It doesn't mean the vaccine isn't safe—
**Synth 2**: Are you really sure about that?
**Dave**: As sure as I am of anything.
**Synth 3**: Well, well, I'm not convinced—
**Shawn**: Nothing will convince you!
**Dave**: OK, we all have to calm down.
**Synth 2**: You want me to calm down? *Of course* I'm already calm!
**Shawn**: This can't be allowed to go public. If it does…
**Synth 1**: [Chuckle-whispers] Our scheme will be found out.
**Dave**: There's still time to deal with the fallout.
**Synth 2**: There's not much time left, though!
**Shawn**: Have you spoken to our friends overseas?
**Dave**: They're as worried as we are.
**Synth 1**: They have the most to lose. Quite possibly even more.

**All the Feels: NVIDIA Shares Expressive Speech Synthesis Research at Interspeech**

Developers and creators can access state-of-the-art conversational AI models for expressive speech synthesis to generate voices for characters, virtual assistants and personalized avatars.

August 31, 2021 by ISHA SALIAN

https://blogs.nvidia.com/blog/2021/08/31/conversational-ai-research-speech-synthesis-interspeech/

https://openreview.net/pdf?id=0NQwnnwAORi

## Hackathon results (SRI)

| Subset | Task 1 DETECT % EER ResNet | Task 2 LOCALIZE % EER ResNet |
|--------|------|------|
| ALL | 26.4 | 33.3 |
| GRP1 | 14.1 | 24.0 |
| GRP2 | 14.1 | 21.5 |
| GRP3 | 18.8 | 38.0 |
| GRP4 | 4.7 | 5.1 |
| GRP5 | 10.9 | 12.2 |
| GRP6 | 53.1 | 86.9 |
| GRP7 | 26.6 | 27.0 |
| GRP8 | 21.9 | 39.1 |
| GRP9 | 21.9 | 32.6 |
| GRP10 | 26.6 | 38.1 |
| GRP11 | 12.5 | 28.5 |
| GRP12 | 23.4 | 26.5 |
| GRP13 | 20.3 | 41.5 |
| GRP14 | 57.8 | 89.7 |
| GRP15 | 10.9 | 30.9 |
| GRP16 | 20.3 | 39.8 |

Different audio categories

# Prototype HMI

### Media asset

### Algorithm results



### Evidence representation

www.darpa.mil