

# **GenAI Perspectives from the MediFor and SemaFor Programs**

---

Matt Turek

October 2023



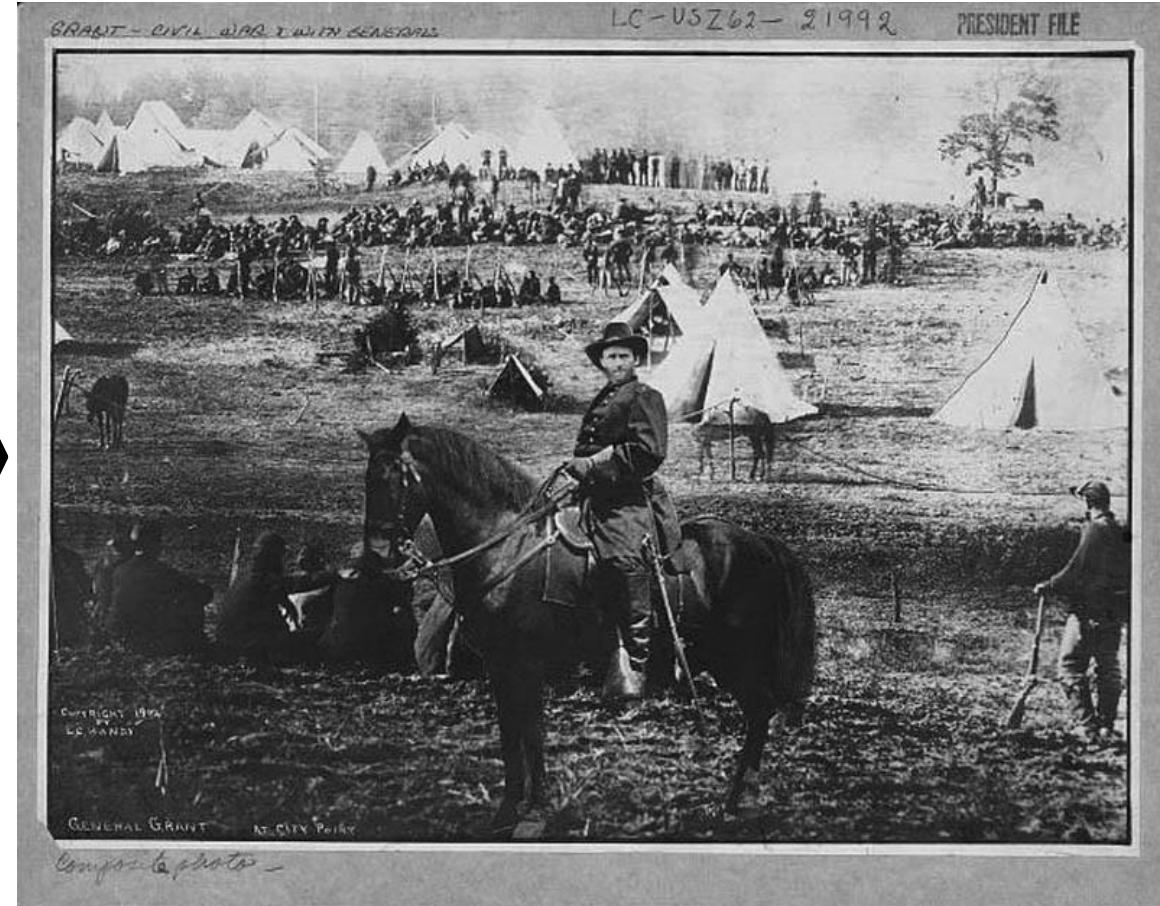
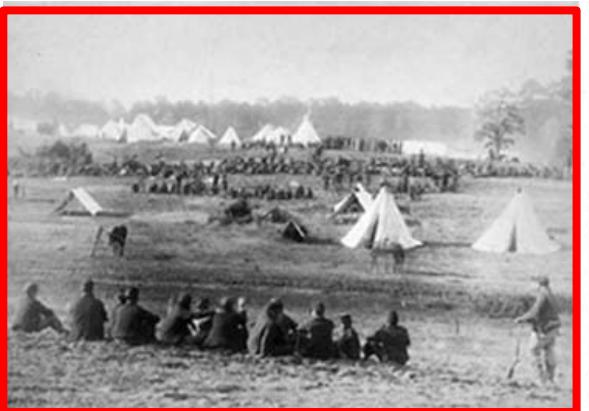
# A brief history of media manipulation technologies



~1865

1860 1910 1920 1930 1940 1950 1960 1970 1980 1990 2000 2010 2020

## A brief history of media manipulation technologies



1902

1860

1910

1920

1930

1940

1950

1960

1970

1980

1990

2000

2010

2020

# A brief history of media manipulation technologies



1917



1860 1910 1920 1930 1940 1950 1960 1970 1980 1990 2000 2010 2020

# A brief history of media manipulation technologies



1939

1860 1910 1920 1930 1940 1950 1960 1970 1980 1990 2000 2010 2020



# A brief history of media manipulation technologies



1860 1910 1920 1930 1940 1950 1960 1970 1980 1990 2000 2010 2020



# A brief history of media manipulation technologies



~80 years



4  
yea  
rs



Text to image

1860 1910 1920 1930 1940 1950 1960 1970 1980 1990 2000 2010 2020

## Targeted Personal Attacks

Peele 2017



AI Multimedia  
Algorithms



Highly realistic video

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

Bricman 2019

AI Multimedia  
Algorithms



Forged  
Evidence

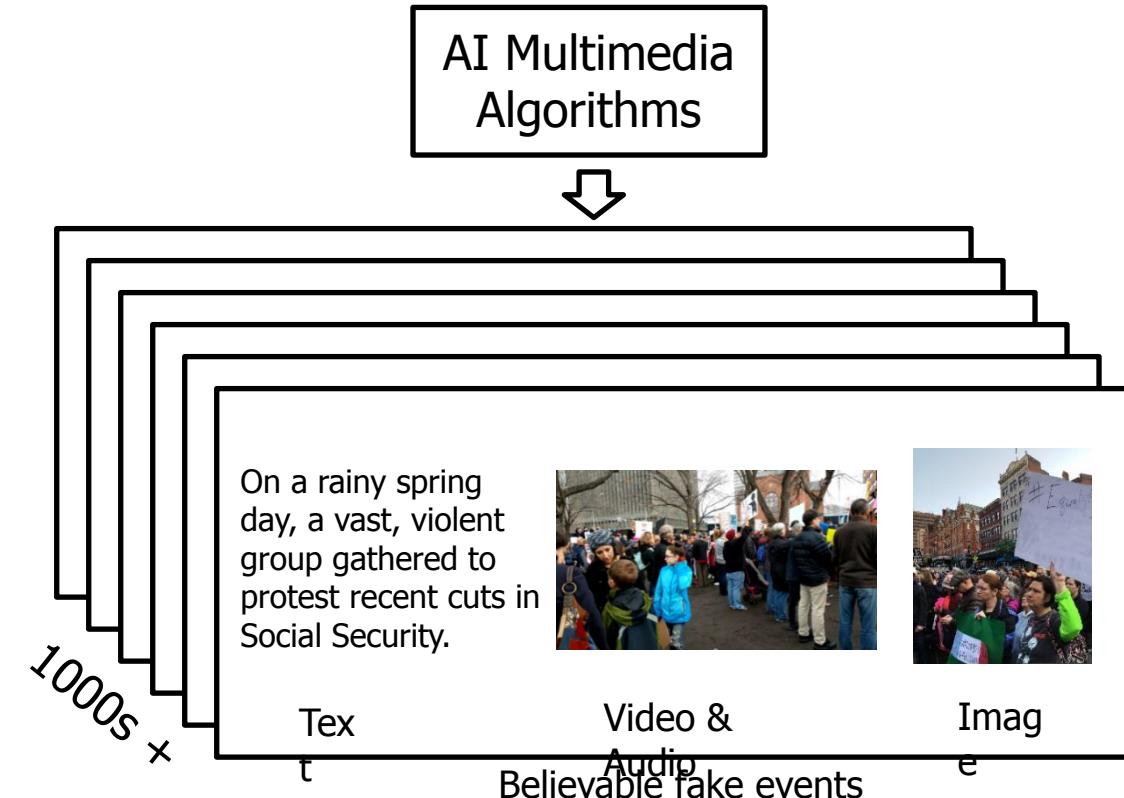


Identity  
Attacks

Examples of possible fakes:

- Substance abuse
- Foreign contacts
- Compromising events
- Social media postings
- Financial inconsistencies
- Forging identity

## Generated Events at Scale



**Undermines key individuals and organizations**



# 2019's speculative falsified media threats are now real

## Targeted Personal Attacks

Peele 2017



Highly realistic video

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

Bricman 2019

AI Multimedia Algorithms



Identity Attacks

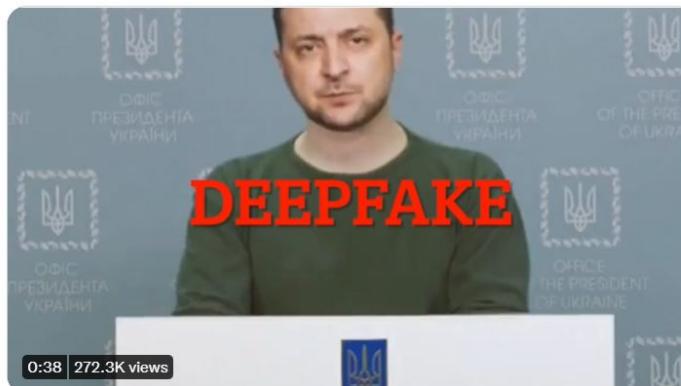
2022



Mikael Thalen  
@MikaelThalen

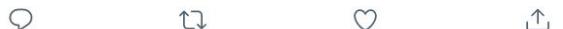
...

A deepfake of Ukrainian President Volodymyr Zelensky calling on his soldiers to lay down their weapons was reportedly uploaded to a hacked Ukrainian news website today, per @Shayan86



11:53 AM · Mar 16, 2022 · Twitter Web App

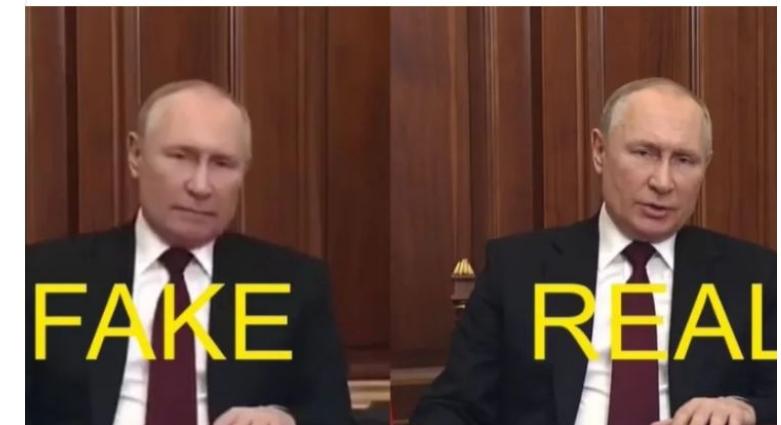
481 Retweets 364 Quote Tweets 870 Likes



## Putin Deepfake Imagines Russian President Announcing Surrender

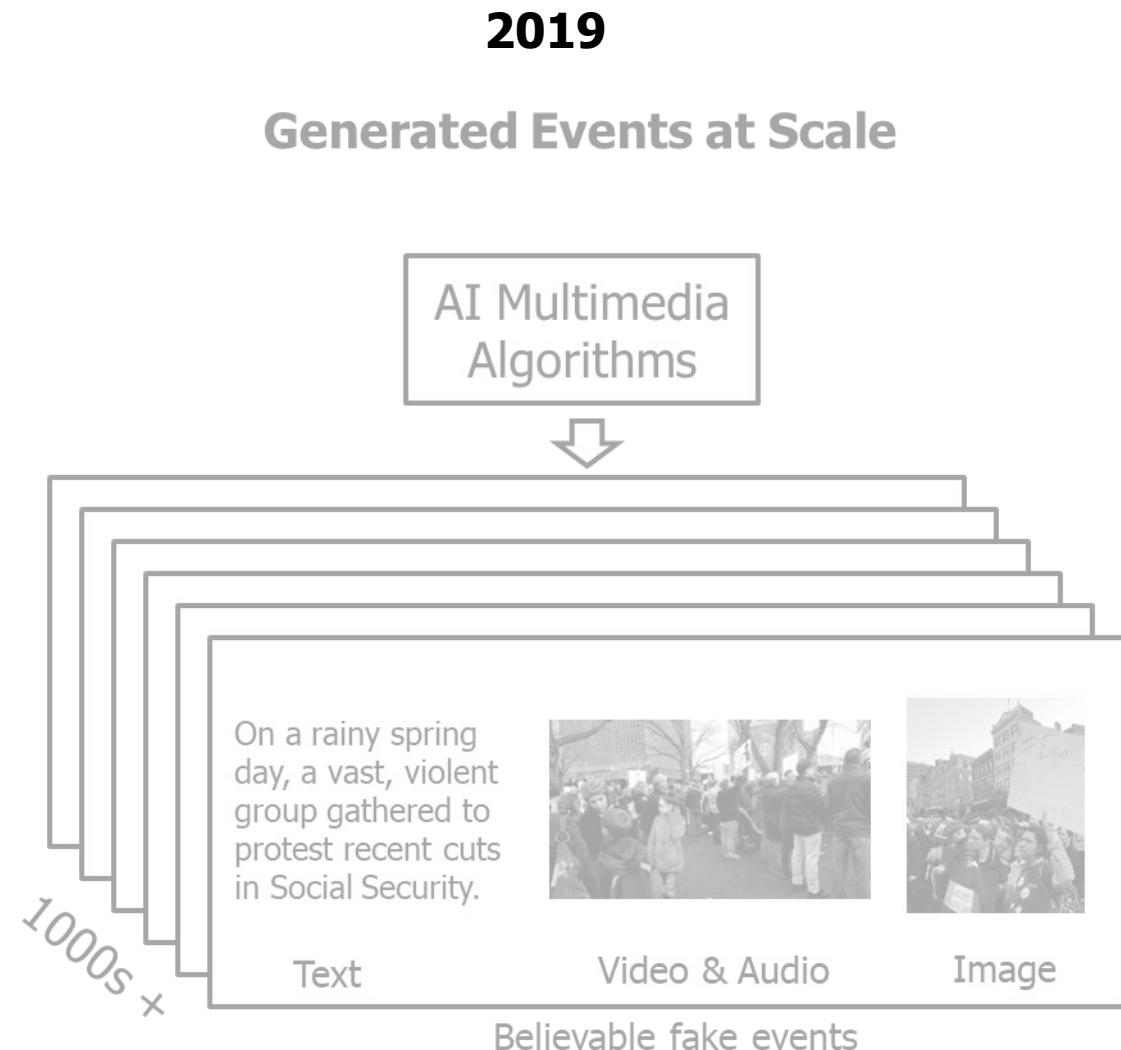
This is likely the first time that deepfake videos have been created and deployed in a war propaganda effort.

By Dan Eron  
Published 18 March 2022





# 2019's speculative falsified media threats are now real





- **Media Forensics (MediFor)**
  - Program Scope

Detect indicators of digital, physical or semantic manipulation in images and video to produce quantitative measures of integrity. Enable integrity assessment of visual media assets at scale.
  - Program kicked off May 2016, ended August 2020
- **Semantic Forensics (SemaFor)**
  - Program Scope

Create rich semantic algorithms that automatically detect, attribute, and characterize falsified multi-modal media to defend against large-scale, automated disinformation attacks
  - Program kicked off August 2020, anticipated completion in September 2024



# Media authentication requires addressing four key problems

---

## Detection

- Determine whether a media asset (image, video, audio, text) is real, manipulated, or AI generated

## Attribution

- Determine the source of the media asset and whether that source is consistent with the purported source
- Attribution to organizations, individuals, tools, techniques
- Answer impacts how the U.S. government responds

## Characterization

- Identify rationale or intent behind the manipulation
  - Target audience
  - Target impact

## Evidence

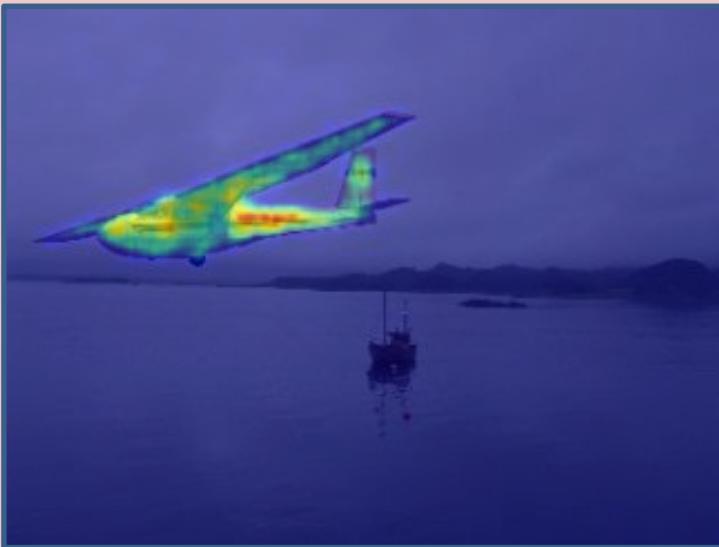
- Not just a determination, but supporting evidence

# Detection & fusion algorithms

**Detection algorithms** analyze media content (pixels) or attributes to determine if manipulation has occurred

**Fusion algorithms** combine information across multiple detectors  
Create a unified score for each media asset

### Digital Integrity



Compression artifacts, PRNU, frequency artifacts

### Physical Integrity



Geometry, lighting

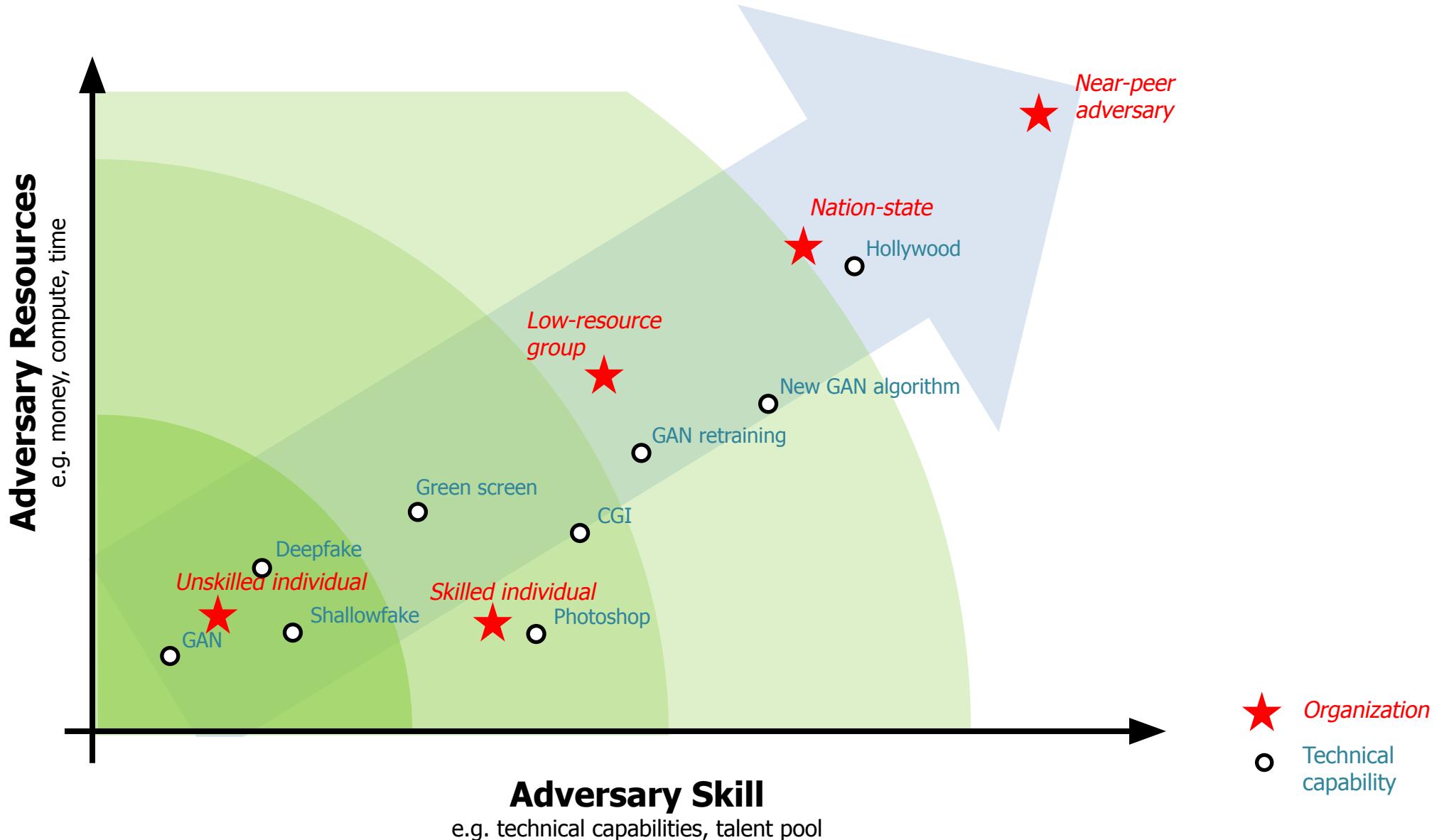
### Semantic Integrity



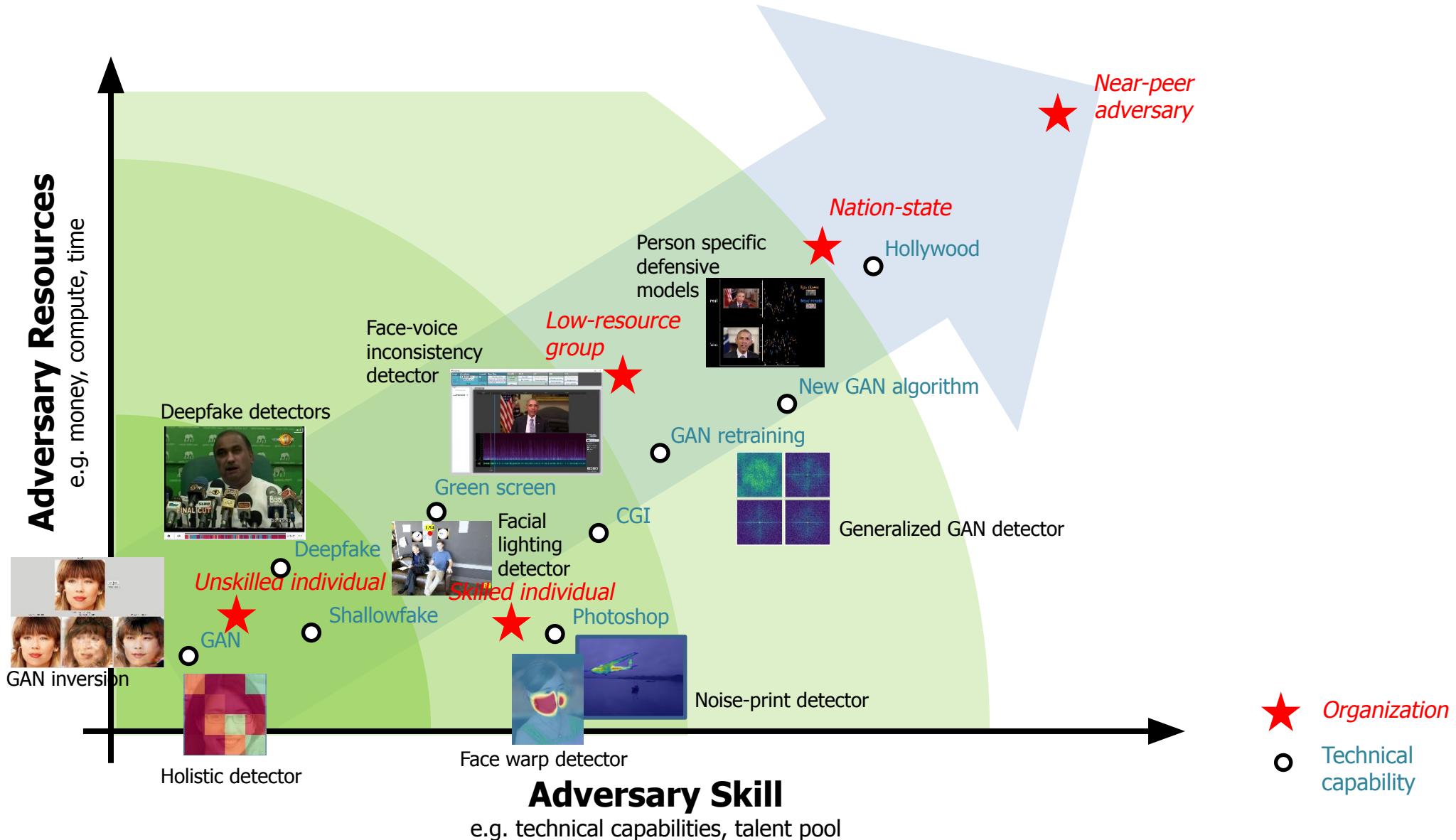
GPS: 43.9196, -78.89419  
Time: 2015/05/29 14:14:30

Comparison to other information

# Adversarial landscape vs. detection algorithms



# Adversarial landscape vs. detection algorithms

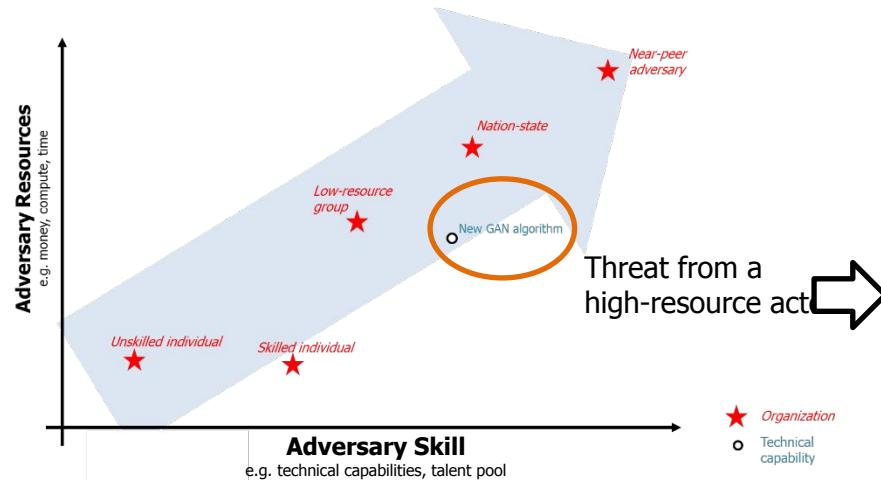


**Challenge:** can detectors identify images from a novel GAN created by a high-resource actor, even without any data from it?

Pro-Chinese Inauthentic Network (2020)

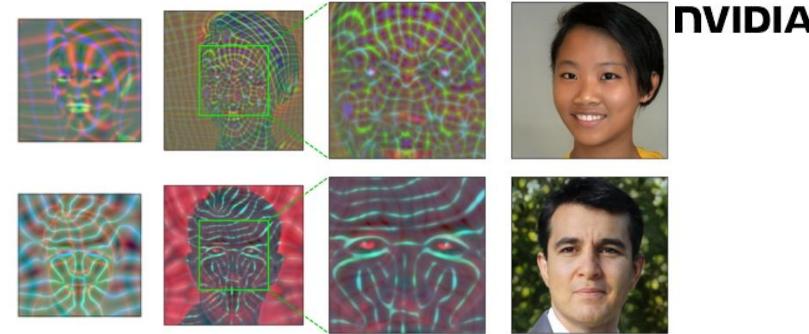


GAN  
Large-scale sock puppet accounts  
Source: Graphika



Alias-Free Generative Adversarial Networks (StyleGAN3)

Official PyTorch implementation of the NeurIPS 2021 paper

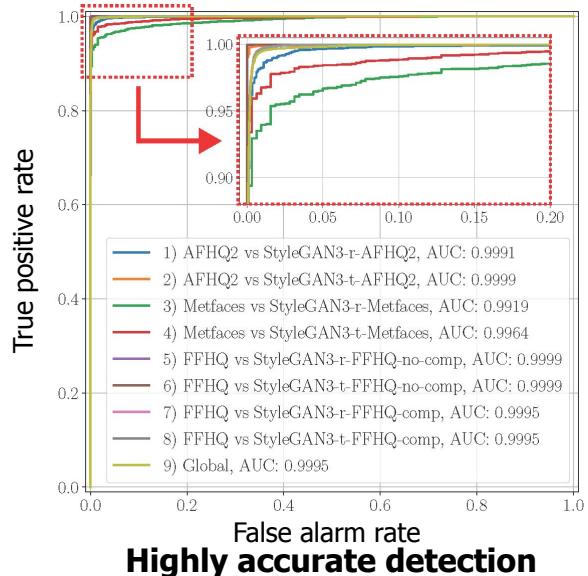


Alias-Free Generative Adversarial Networks  
Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, Timo Aila  
<https://nvlabs.github.io/stylegan3>

Training over semantic categories, augmentation, & many GANs



No training data from StyleGAN3!



Distribution A: Approved for public release: distribution unlimited.

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

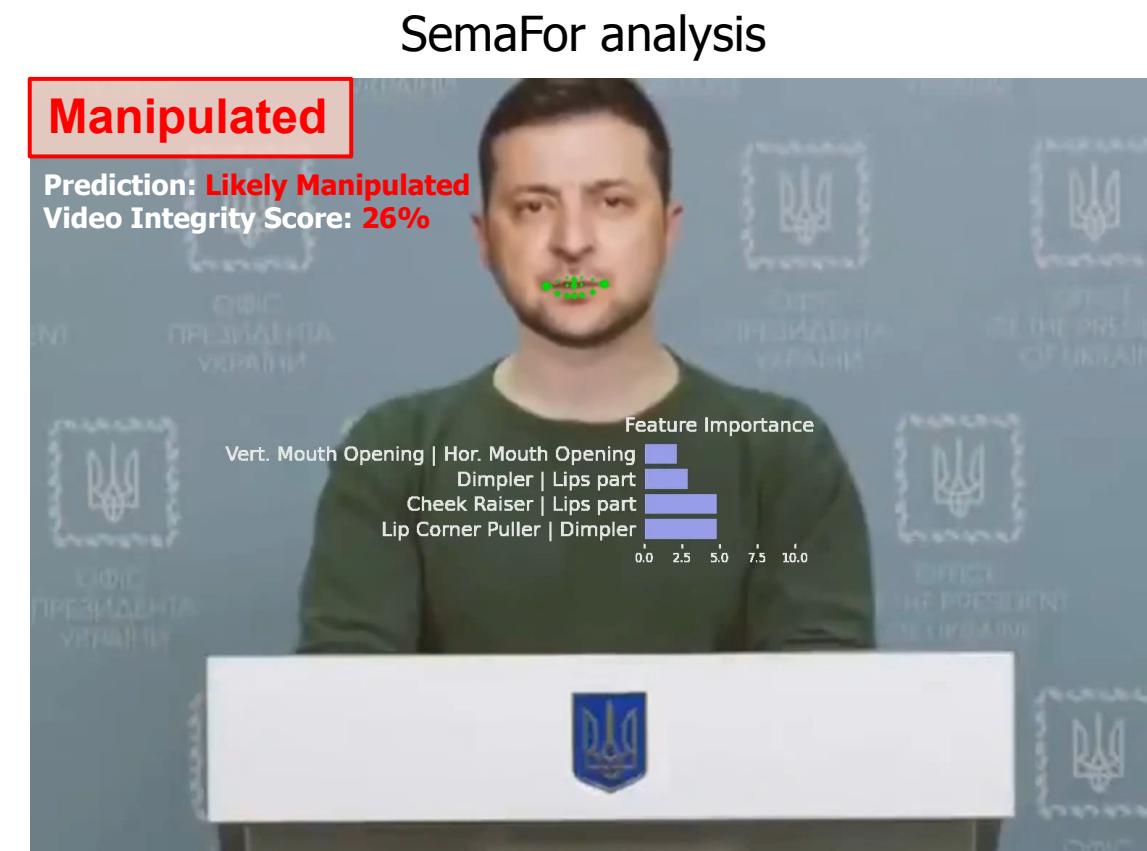
# Person of interest models (attribution to an individual)



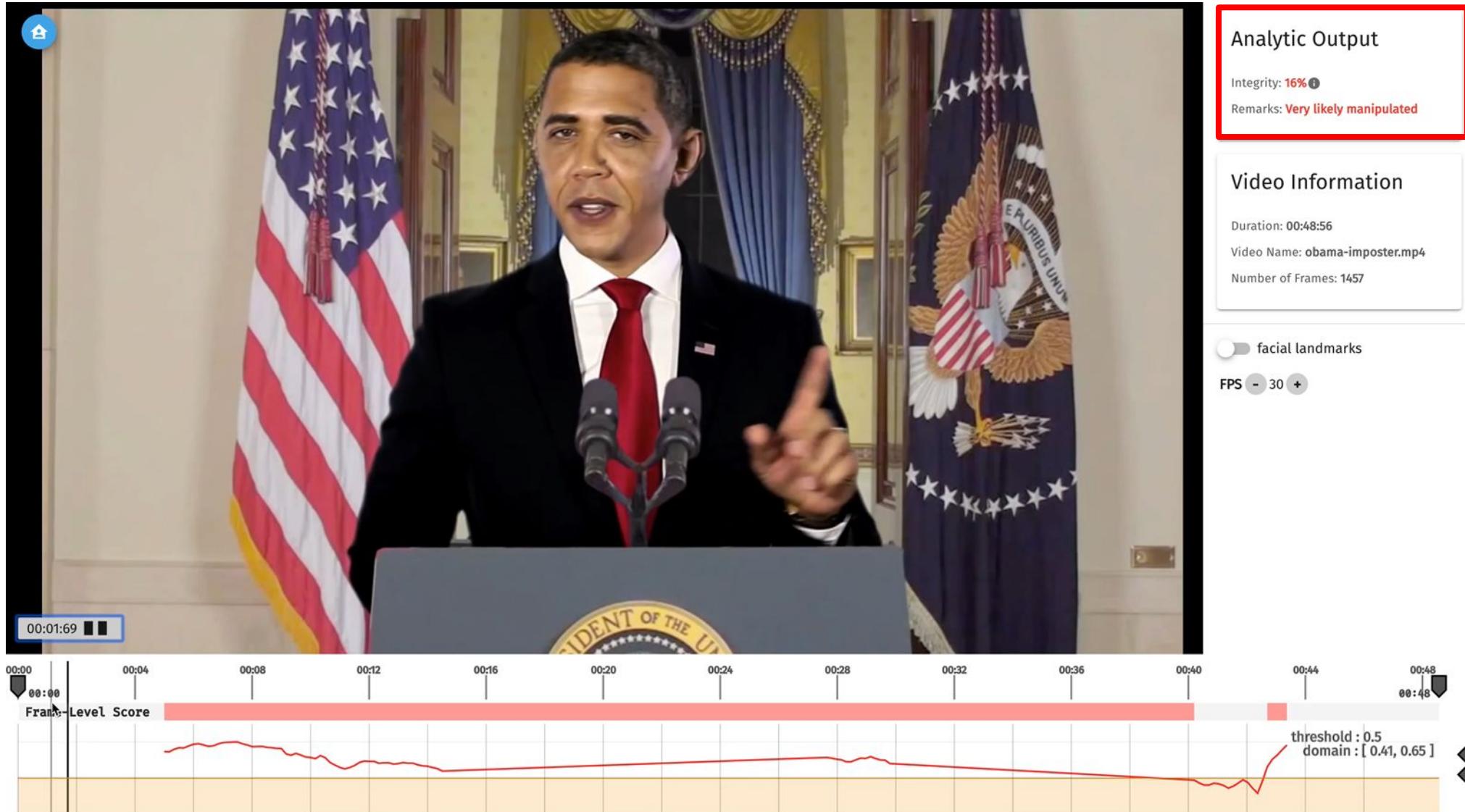
Soft biometric approaches build models of an individual's behavior from known "good" data

Automatically identify facial landmarks, facial action units, and head pose

Machine learning framework creates a face and head movement model for a particular individual



- A video of purported to be President Zelenskyy was posted on several online forums calling for Ukrainian troops to lay down their arms March 16, 2022.
- This video is probably the first use of deepfakes in the context of war.





# SemaFor system: detection, attribution, and characterization analytics

Lockheed Martin

**Pentagon Explosion**

1 Media 1 Needs Review

Uploads started on May 25, 2023 - 18:54 by william.corvey(William Corvey)

WC 5/25/2023

Falsification Likelihood: Certain (circled in red)

Review Status: Need Review (circled in blue)

22 Analytics

Showing 1 - 1 of 1 uploads

**Pentagon Explosion.JPG**

Likely Falsified

Analysis had positive findings and needs to be reviewed. Check 'Mark reviewed' to finish reviewing the result.

Image Analysis Results

Detected: Computer Generated Detection

Needs Review Mark reviewed

Generate Reports

Finding Likelihood: Certain (red), Likely (orange), Maybe (yellow), Unlikely (green), No Chance (blue)

By Analysis Categories (3) By Analytics (32)

Note: The results of this analysis may appear contradictory due to the unique datasets and methods by which each analytic was trained.

Score: 1

Computer Generated Detection (Likely to Certain)

No Detection

Manipulation Detection (Localization)

Opted Out & Failed

Computer Generated Detection (17 Opted Out)

Manipulation Detection

LOCALIZATIONS (1)

Original

Manipulation Detection

Download Media

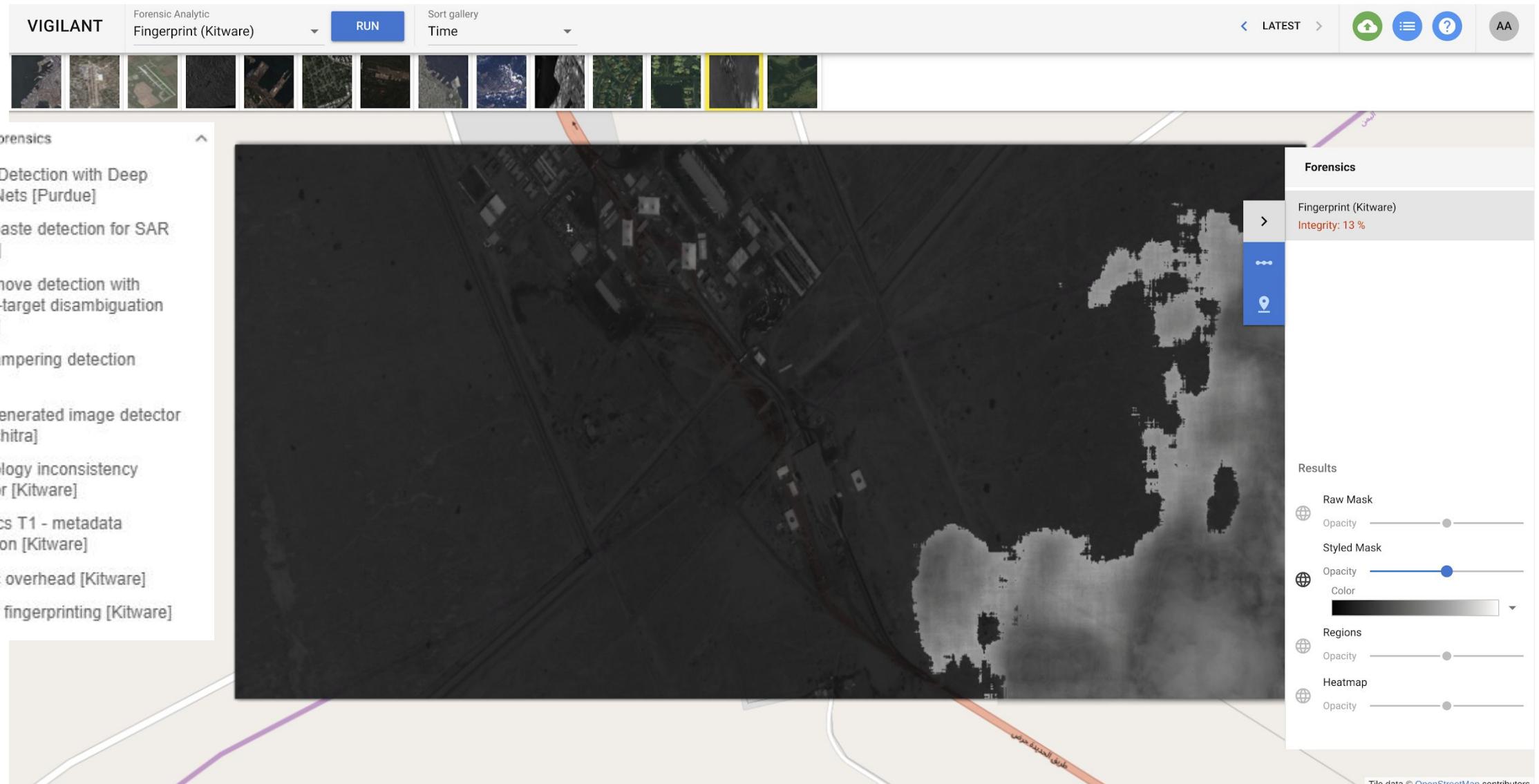
Analysis Date: 05/25/2023, 2:54 PM

Analysis Runtime: 03:55

Mime Type: image/jpeg

Modality: IMAGE

Comments (0)



## Scientific Integrity





WELCOME Rating

Logout

Main Navigation

Upload Image

Upload PDF

View PDFs

Search

Case Report

History

Select All Images

Deselect All Images

Image Analysis

Provenance Graph

Delete Selected Images









figure 2 synthesis of pegylated carbonized...  
hosseinihani\_2005.pdf

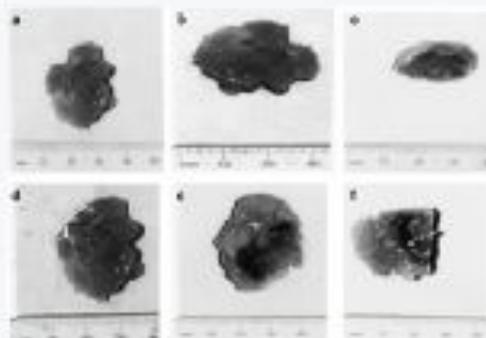


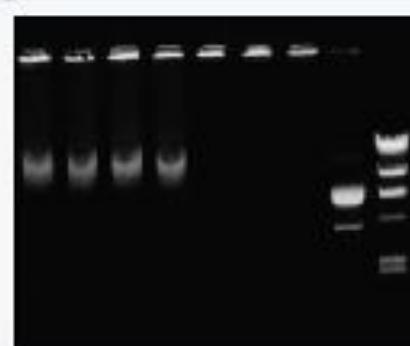
figure 8 shows the survival curve of tumor-be...  
hosseinihani\_2006.pdf



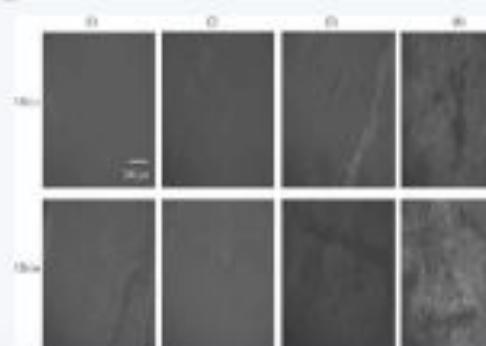
figure 5 time dependency of amount of rnf pro...  
hosseinihani\_2006.pdf



hosseinihani\_2005.pdf



hosseinihani\_2006.pdf



hosseinihani\_2006.pdf



[www.darpa.mil](http://www.darpa.mil)



# SemaFor video analytic

SemaFor Portal Engineering UI Gym

Results New Upload

Q Search Media

478302398.mp4 CHINA CHARACTERIZATION + Tag Generate Report

Manipulation Very Likely

It is **very likely** that this image was manipulated using StyleGAN 3, XYZ, et cetera et cetera.

Analytics

All (23) With Localization (5) **Detection (6)**  
Attribution (4) Characterization (3) Opt-Out (3)

Manipulation Very Likely

No Localization

Very Unlikely      Uncertain      Very Likely

*purdue-unisi-6classesttribution-0-3-8*  
This analytic is used to attribute which image generator created a certain image

Manipulation Very Likely

Very Unlikely      Uncertain      Very Likely

*purdue-unisi-6classesttribution-0-3-8*  
This analytic is used to attribute which image generator created a certain image  
1:32-1:45

Manipulation Likely

Very Unlikely      Uncertain      Very Likely

*purdue-unisi-6classesttribution-0-3-8*  
This analytic is used to attribute which image generator created a certain image  
1:32-1:45

Manipulation Uncertain

Very Unlikely      Uncertain      Very Likely

*purdue-unisi-6classesttribution-0-3-8*

Processing Results... 1 minute ago  
10933751.jpg Batch 21

Very Likely Manipulated 1 day ago  
10933751.jpg Batch 20

Likely Manipulated 1 day ago  
10933751.jpg Batch 20

Very Unlikely Manipulated 3 days ago  
10933751.jpg Batch 19

Very Likely Manipulated 3 days ago  
10933751.jpg Batch 19

Very Likely Manipulated 1 week ago  
4:32 10933751.mp4 Batch 18

Very Unlikely Manipulated 1 week ago  
12:22 10933751.mp4 Batch 18

Very Likely Manipulated 1 week ago  
1:23:22 10933751.mp4 Batch 18

Very Likely Manipulated 1 week ago  
3:23 10933751.mp4 Batch 18

Unlikely Manipulated 2 Weeks Ago  
0:45 10933751.wav Batch 17

Very Unlikely Manipulated 2 Weeks Ago  
0:45 10933751.wav Batch 17

Very Likely Manipulated 2 Weeks Ago  
0:45 10933751.wav Batch 17

Likely Manipulated 2 Weeks Ago  
0:45 10933751.wav Batch 17

Unlikely Manipulated 2 Weeks Ago  
0:45 10933751.wav Batch 16

Very Unlikely Manipulated 2 Weeks Ago  
10933751.jpg Batch 16

Very Unlikely Manipulated 2 Weeks Ago  
10933751.jpg Batch 16

Localization from Analytic Name  
This analytic is used to attribute which image generator created a certain image

Hypothesis  
Lorem ipsum dolor sit amet, consectetur adipiscing lorem ipsum dolor sit amet, consectetur adipiscing

View Original Download

3 of 12

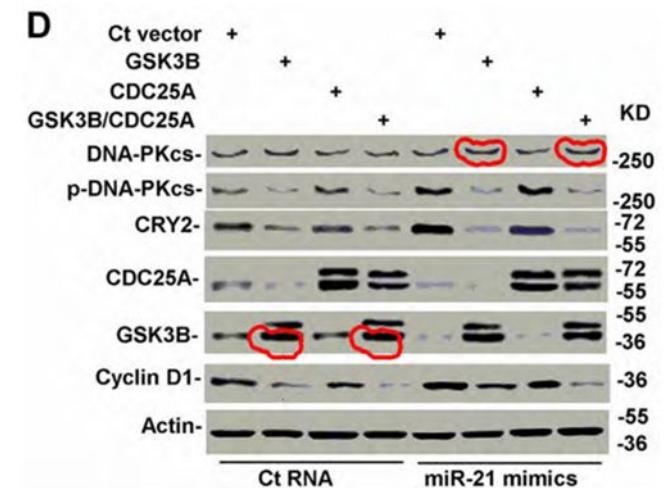
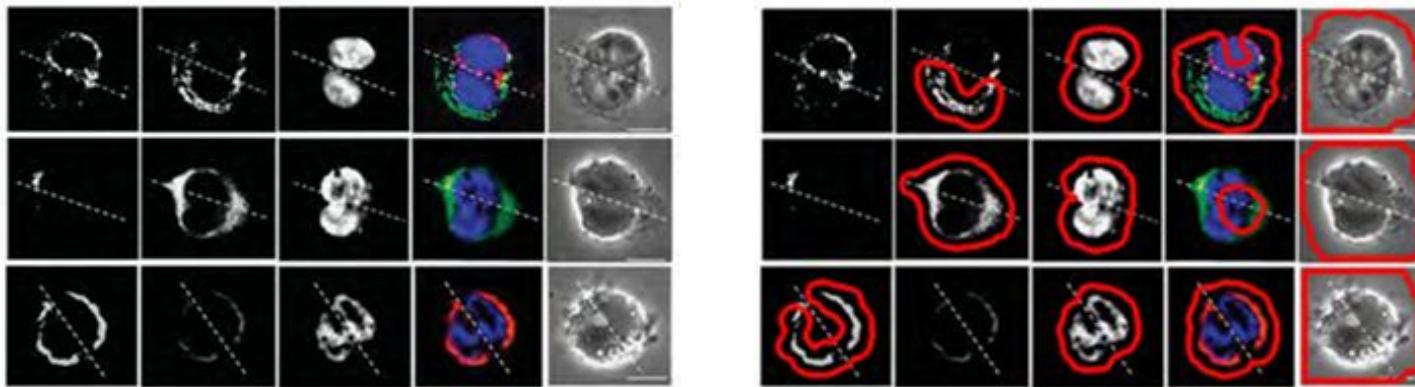
3:25/5:12

\*Mock screenshot



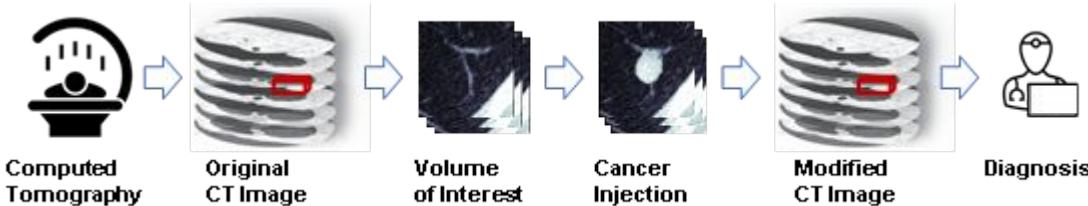
# Misinformation can undermine the scientific process

## Biomedical image manipulation

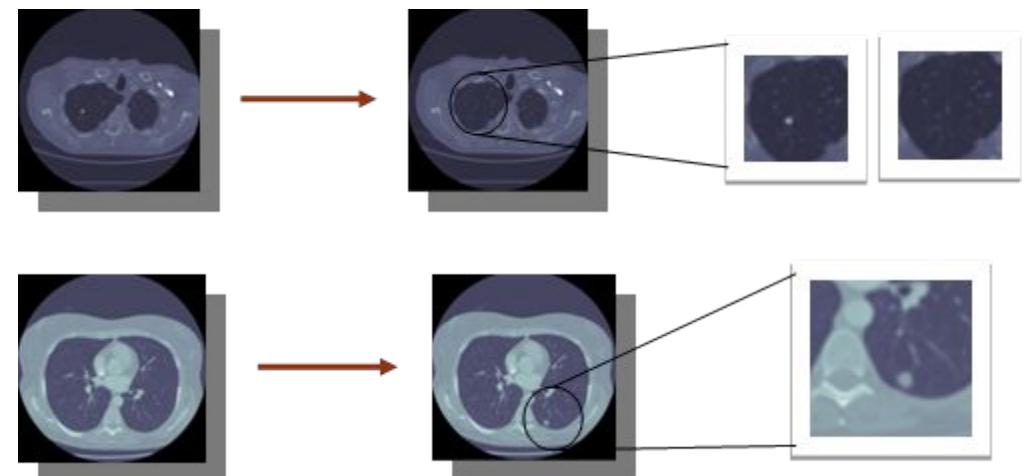


Advanced image manipulation of cancer study results, attempting to show treatment experiment successful. "Copy Move"

## Radiological image manipulation



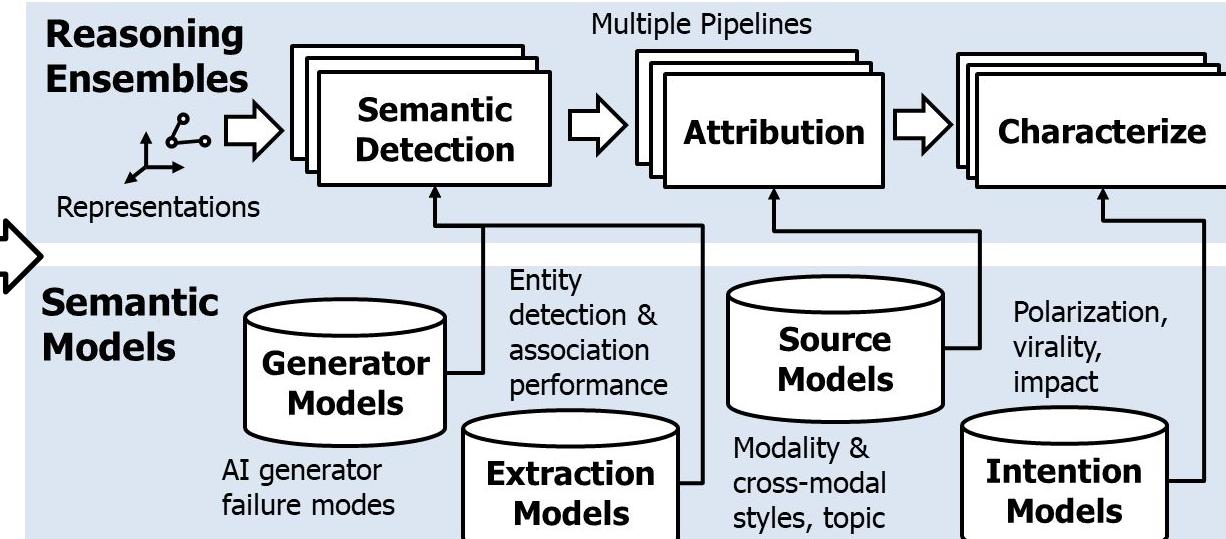
AI generated cancer inserted into 3d Computed Tomography (CT) image fools 95% of radiologists.



# SemaFor technical areas & performers

## TA1: Detection, attribution, characterization

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



**Performers**  
**Kitware**, UIUC, Columbia, ASU, SUNY Albany, UB, UMichigan, Eduworks  
**SRI**, UMd, UB, CMU  
**Purdue**, Notre Dame, Univ Siena, Univ Frederico II of Naples, Politecnico di Milano  
**UC Berkeley**, Boston, University of Washington, UC Davis, Pinscreen

## TA4: Challenge Curation

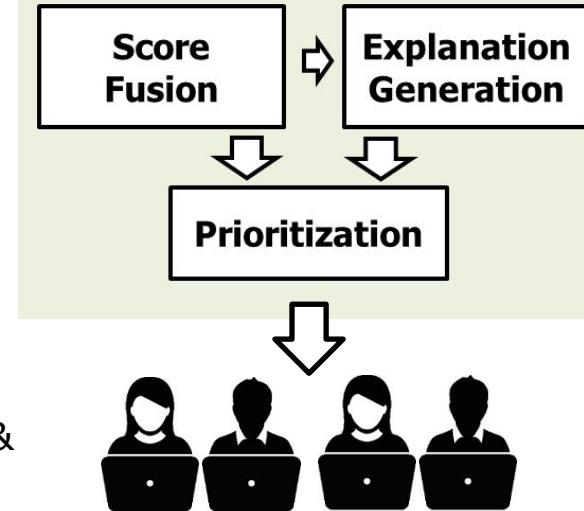
SOTA challenges

Threat modeling

**Performers**  
**STR**, Graphika, Hwang  
**NYU**, Vidrovri, IST Research, Watts  
**Accenture**, Carley  
**NVIDIA**  
**Google**

## TA2: System integration

### Explanation & Integration



**Performers**  
**LM-ATL**,  
 Boston Fusion,  
 CRA

## TA3: Evaluation

Media generation

Evaluations

Metrics

**Performers**  
**PAR**, Aptima, Rank One, Drexel, CU Denver, Syracuse

## Text (Notional)

NewsWire: April 1, 2019, Bob Smith

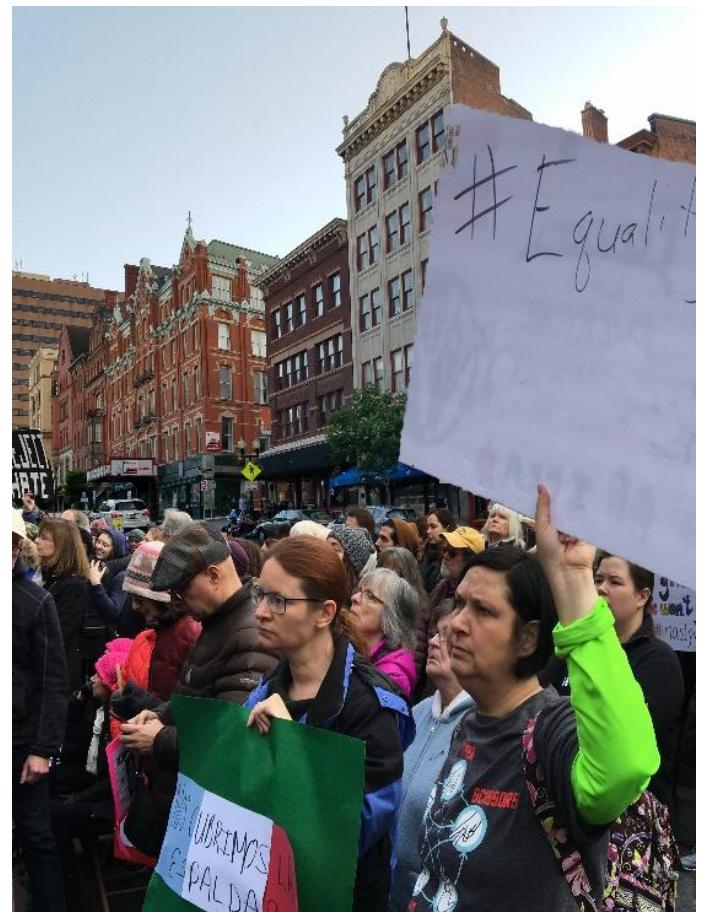
On a rainy spring day, a vast, violent group gathered to protest recent cuts in Social Security.

## 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 to protest recent cuts in Social Security.

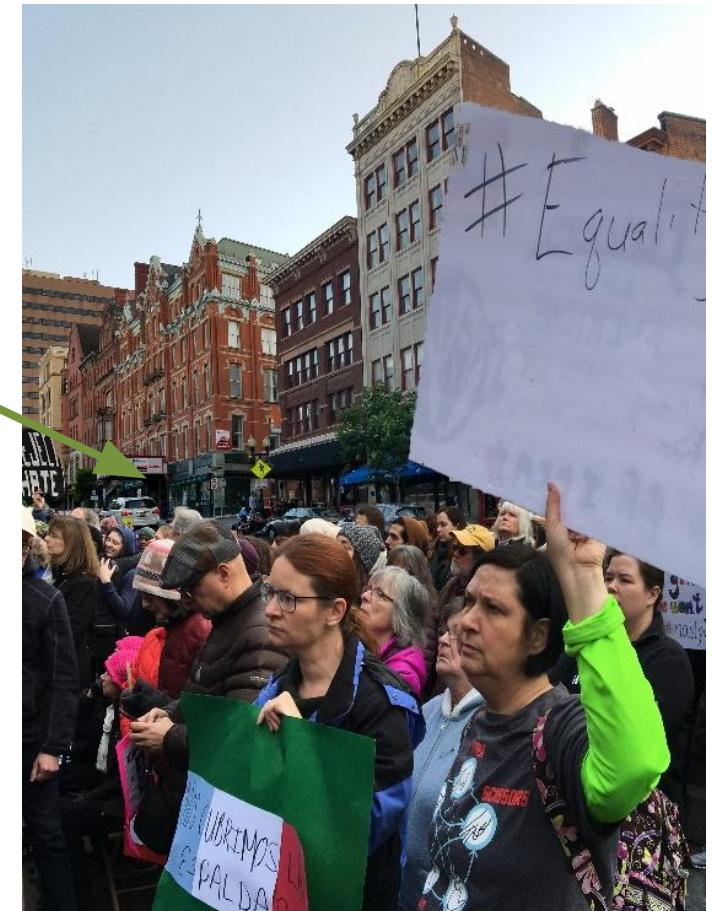
**Audio (Notional)**

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

**Video**

"protest"

Conclusion: Media components consistent across modalities.

**Image**

## Text (Notional)

NewsWire: April 1, 2019, Bob Smith

On a rainy spring day, a vast, violent group gathered 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]..."



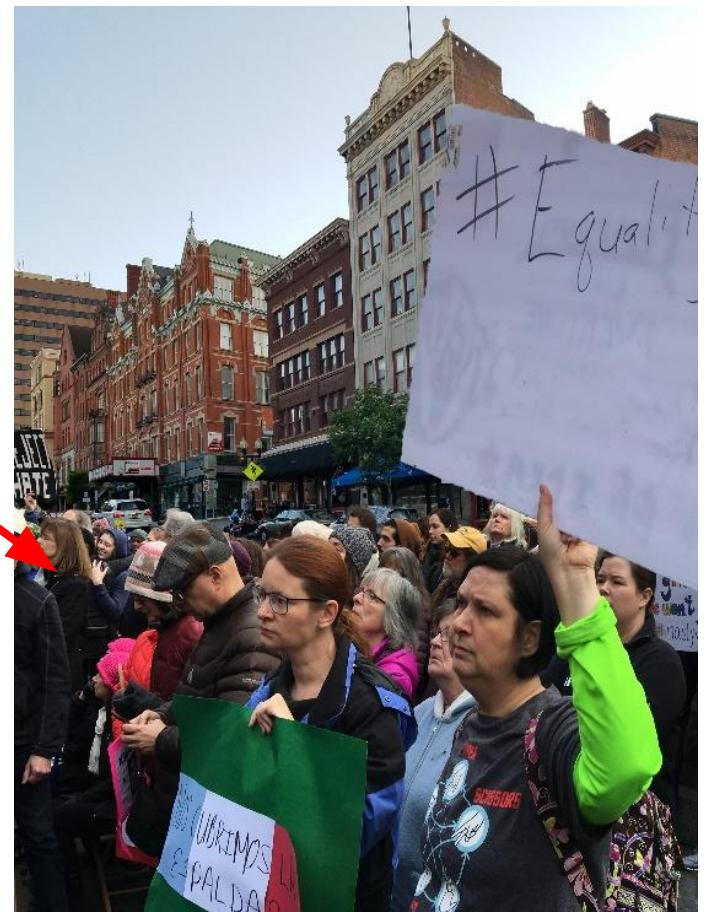
## Video



"violent group"

Conclusion: Media components not consistent across modalities.

## Image



## Text (Notional)

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

### Video



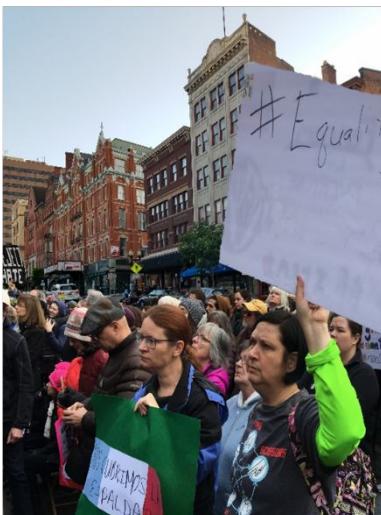
### Audio (Notional)



"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”
  - Topic – “Social Security”
- Use of unsupported term “violent”
- Failed sourcing to high credibility organization (“NewsWire”)

## SUMMARY ATTRIBUTION PERFORMANCE



18 of 20 Attribution Tasks met Program Objectives

Best Performing Tasks	Most Challenging Tasks
Synthetic Media Attribution for Text, Image, Audio, and Video Generators	GPT-j-6b synthetic text generation model with human manipulation



Synthetic image of military equipment generated with Guided diffusion model.



Synthetic image from StyleGAN3.

:: Title :: Cuban technicians fix thousands of SANDF vehicles ::  
:: Slug :: Military Africa ::  
:: By :: Sarah Lesedi ::  
:: Date :: 03/18/2022 00:00 ::  
Cuban technicians fix thousands of SANDF vehicles Published duration 19 July 2018  
image copyright AFP Image caption Rifles and other hardware are transported during their inspection.  
Thousands of pieces of military hardware seized by the South African government are being examined in Havana by technicians from the Republic of Cuba's military-industrial complex.  
South African authorities say thousands of pieces of military equipment are believed to have been looted during the former President's administration.  
The items, seized in 2010, are now being tested to determine the value of them.  
Tens of millions of dollars could change hands if any of the equipment is deemed valuable.  
However, the value of any of the thousands of pieces will likely be less than the cost of transporting the items back to South Africa.  
That is because the state-run Cubadeb utility company is the only company qualified and willing to transport such a large cache of goods.  
The military-industrial complex is responsible for the manufacturing, repair, overhaul, and the maintenance of the equipment.  
It also carries out the transport of the equipment from one end of the Latin American country to another.  
The company has offices in many cities and is in charge of servicing military vehicles in all of them.

Synthetic text article about military capabilities from GPT-j-6b.

## SUMMARY CHARACTERIZATION PERFORMANCE



8 of 8 Characterization Tasks met Program Objectives

Best Performing Tasks	Most Challenging Tasks
Characterization of <i>intents</i> and <i>propaganda</i> tactics (hate speech, bandwagoning, dictat, scapegoating)	<i>Appeal to Fear</i> and <i>Minimization</i> tactics



Utilizing the **appeal to fear tactic** the image is manipulated to exacerbate fear about COVID in a technical information news article.



Supporting the **bandwagoning tactic**, the image is manipulated to insert soldiers to convey and support the tactic in a social media post.