# Combatting with DeepFakes
# New Trends and Thoughts

## Siwei Lyu, Ph.D.

- Department of Computer Science and Engineering,
- Center for Information Integrity,
- Media Forensics Lab,
- University at Buffalo, State University of New York

# May 21st, 2023

# May 21st, 2023



**AI Generated**
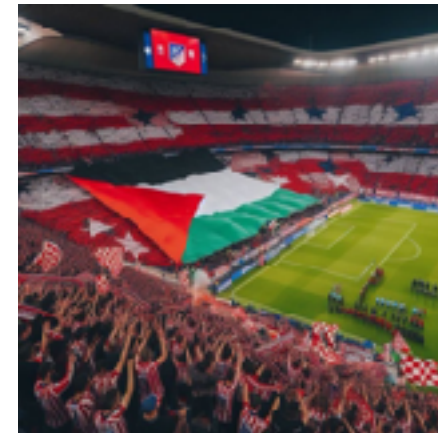
# DeepFakes on the rise



Twitter 11/13/2023

Twitter 11/06/2023
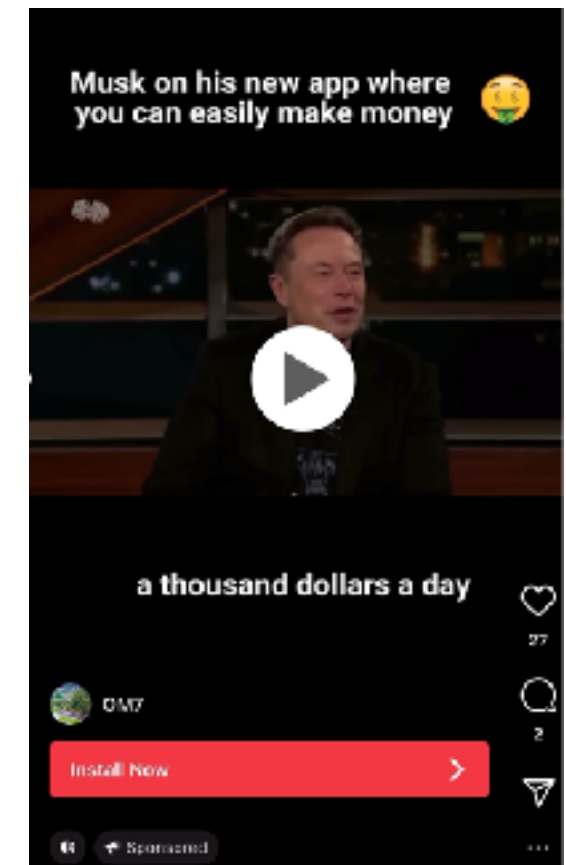
Twitter 10/30/2023

Twitter 10/20/2023

Facebook 11/06/2023

Facebook 10/27/2023

Twitter 10/22/2023

# Rise of AIGC

**Science & engineering**



VAE  Transformer  Autoregression

GAN  Normalized Flow  Diffusion

**Tools & services**



**Data source & spreading channels**



**Computation power**

# White House Executive Order on AI



Administration    Priorities    The Record

Gen AI

Watermarking

CTOBER 30, 2023

resident Bic
r on Safe, Se
rtificial Inte

OM ▸ STATEMENTS AND RELE

enabled fraud and
est practices for d
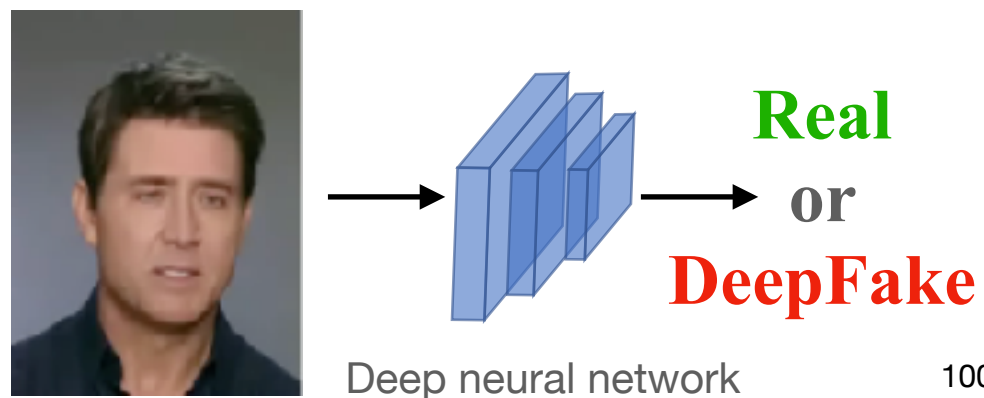fficial content. Th

Commerce will develop guidance for content authentication and watermarking to clearly label AI-generated content. Federal agencies will use these tools to make it easy for Americans to know that the communications they receive from their government are authentic—and set an example for the private sector and governments around the world.
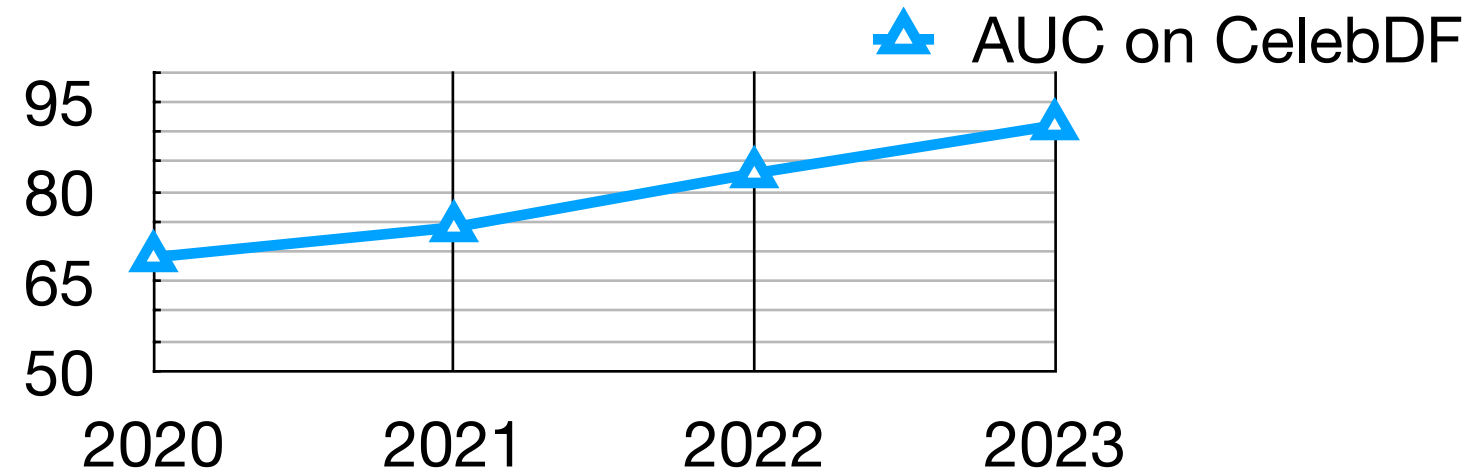
# AIGC detection — semantic based

# AIGC detection — data driven



Deep neural network

Real
or
DeepFake

- More sophisticated and powerful NNs
- Larger and more diverse datasets
- Increasing performance on benchmarks

Resent

VGG

Xception Net

Capsule Net

Efficient Net

HRNet

# of 1,000 frames

10000

20,000

10,000

3,068

4,133

5,369

1000

1,000

100

640

49

10

- UADFV
- DeepFake-TIMIT
- FaceForensics++ (DeepFake)
- Google DeepFake Dataset
- DeepFake Detection Challenge (initial set)
- Celeb-DF
- DeepFake Detection Challenge (full set)
- Deeperforensics-1.0

AUC on CelebDF

95

80

65

50

2020    2021    2022    2023

# AIGC detection — are we there yet?

- Lack clarity in the problem definition

  - Is "real vs. DeepFake" well-defined?

- What is the purpose of detection

  - Use for triaging (overall accuracy/throughput) vs. use for evidence (individual accuracy/explanation)

- No option for "I am not sure"

  - Forcing detector to make decision on every input

- Too much reliance on labeled data and no diagnosis

  - Learning circumstantial attributes — source for lack of robustness to post-processing and low generalizability

- No user-friendly testbeds for performance

  - Commercial closed source systems vs. open source research code on code repositories

# DeepFake-o-meter

- Open-platform of open source AIGC detection methods

- Simple API/containerization to simplify integration

- Easy interface to casual users to run on individual media

- Comprehensive evaluation of in-the-wild performance



https://zinc.cse.buffalo.edu/ubmdfl/deep-o-meter/

# AIGC detection — what's next?
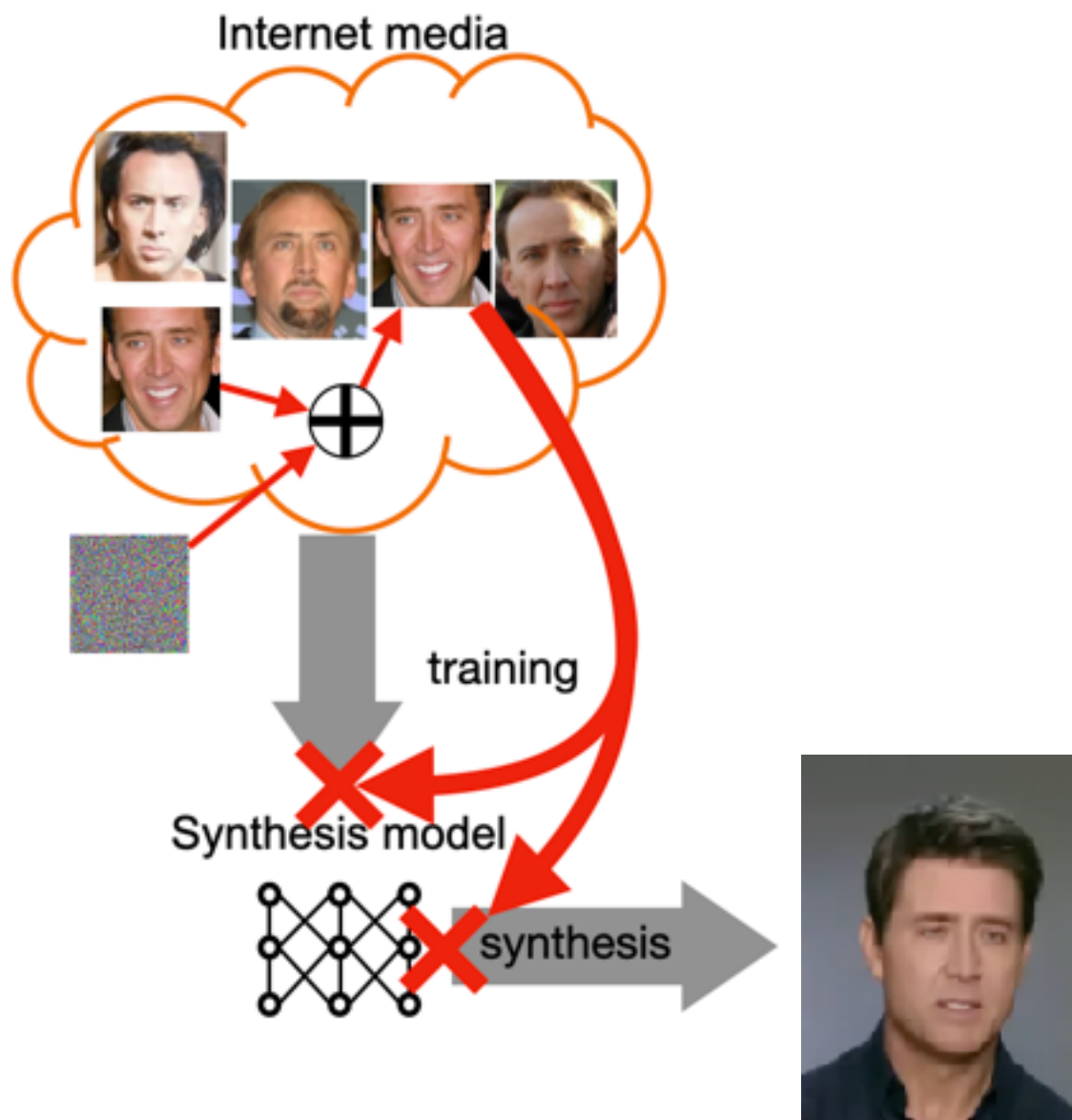
- Detection of AIGC may become irrelevant



Apple is releasing their much-rumored Vision Pro headset next year, complete with an AR/VR version of Facetime. To allow others you're chatting with to see your likeness while you have the Vision Pro on, Apple will prompt users to scan their face using the device's cameras, creating a lifelike representation of your face (a "persona") that moves its lips and emotes as cameras sense your own movements and gestures.
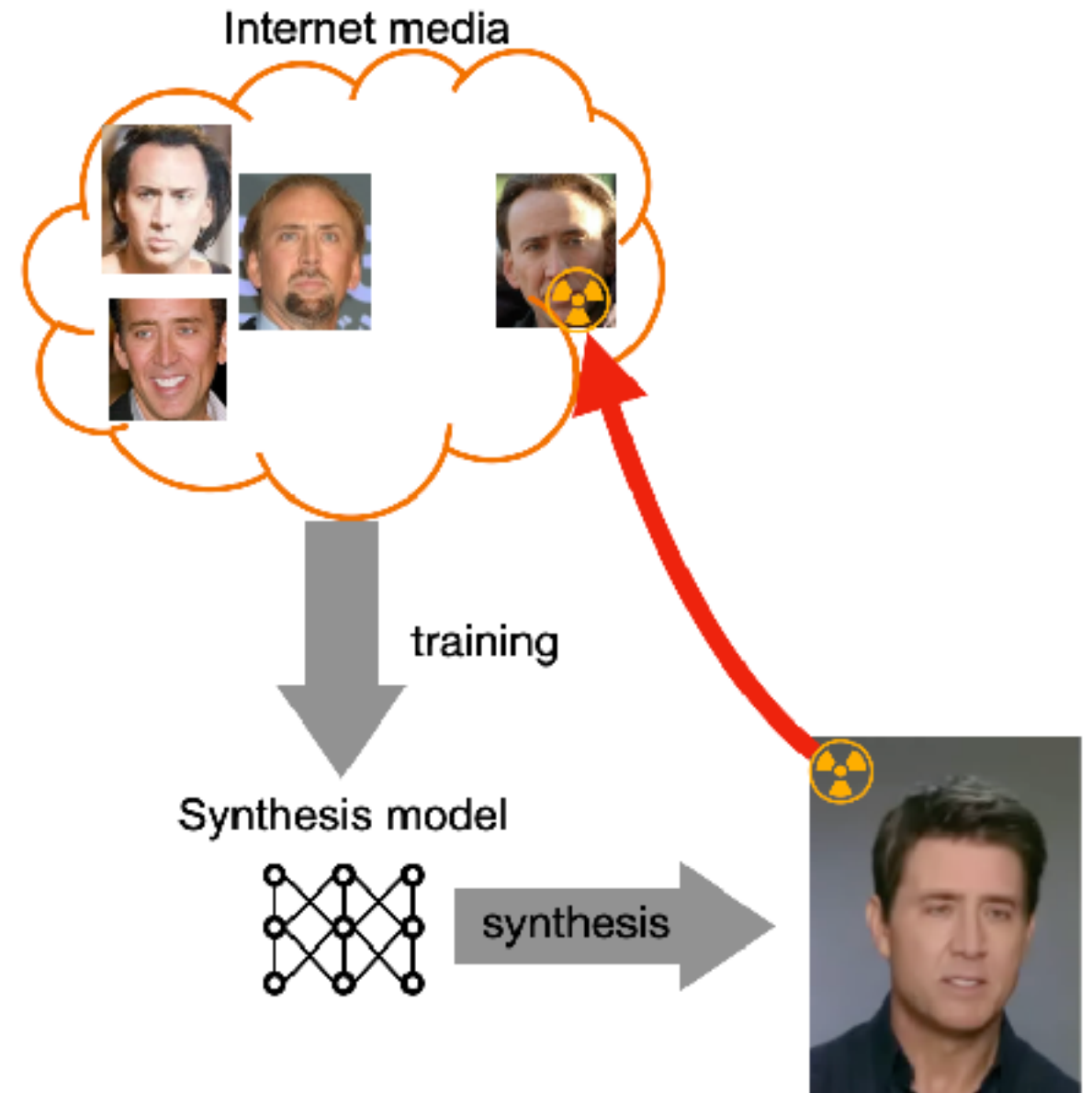
  - Verification of authenticated use of gen AI models

    - Liveness/presence detection for real-time AIGC

- Attribution and provenance of the generation process

- Providing fine-grained analysis — spatial and temporal localization

- Other modalities and cross-modality analysis

- Handle the zero-day attack: making detection generalizable
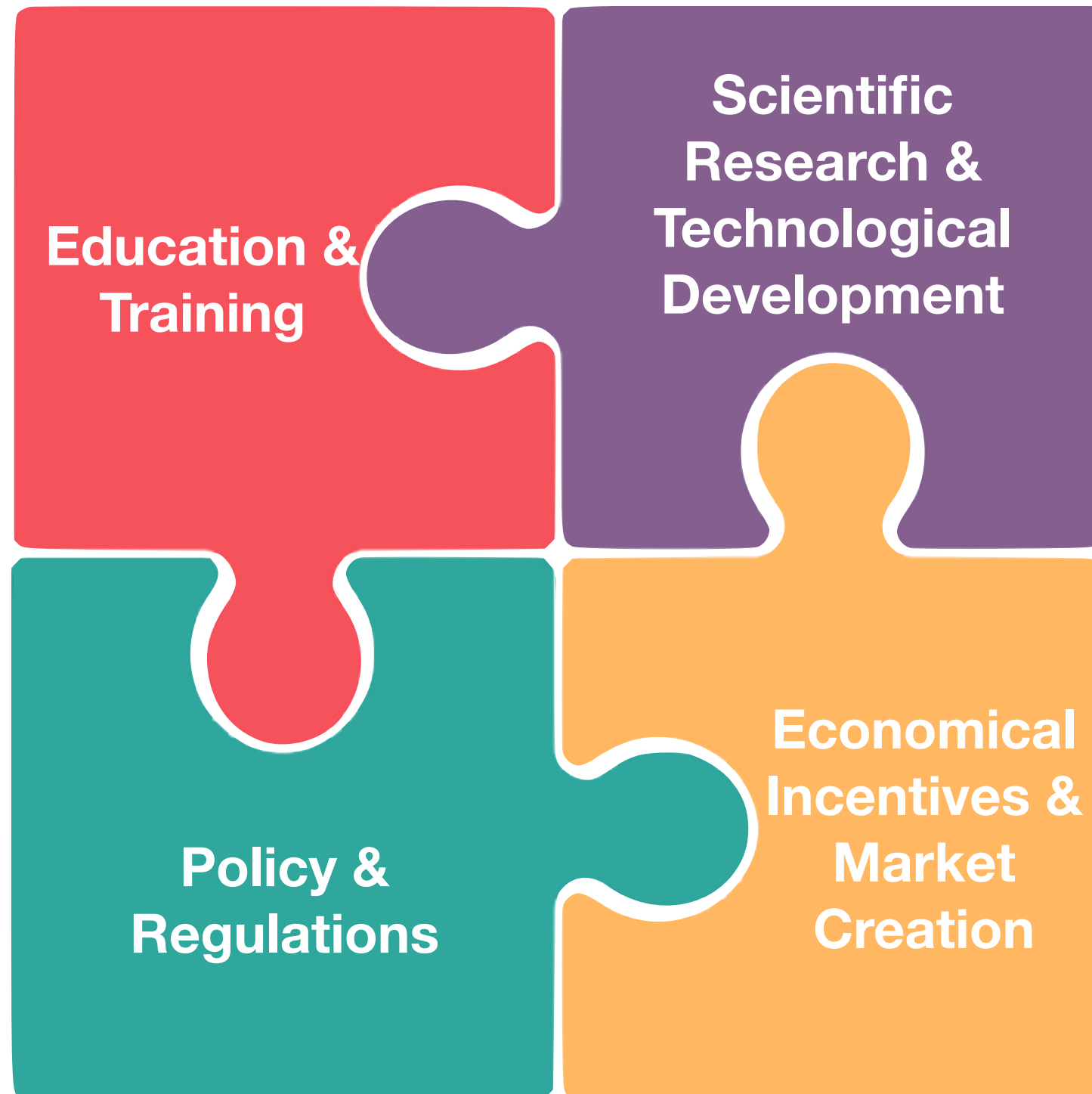
# Beyond detection



**Obstructing DeepFakes**

**Tracing DeepFakes**

# Beyond technologies

DART Learn

DeepCover

# Conclusions

- AIGC continues to have rapid developments

- Misused AIGC (aka DeepFakes) erode our trust to online information

- There is no "silver bullet" solution, we need all

  - Content authentication and provenance

  - Watermarking

  - Reactive forensics — detection, attribution

  - Active forensics — affecting training/synthesis

- Future directions for AIGC detection

  - Improve robustness, transparency, and explainability

  - Handling diversity and complexity of generation

  - Providing more insights beyond binary classification

# Thank you!

Contact: siweilyu@buffalo.edu