

# Trustworthy AI Systems

-- Accountability of AI

Instructor: Guangjing Wang

[guangjingwang@usf.edu](mailto:guangjingwang@usf.edu)

# Last Lecture

- Membership Inference Attacks
- Model Inversion Attacks
- Model Stealing Attacks
- Privacy Protection Methods

# This Lecture

- Accountability
- Detecting AI-generated Content
- Watermarking Techniques
- Evading Watermarking-based Detection

# Accountability of AI

*If an AI system causes harm (e.g., a self-driving car accident), is it the company deploying the system, the developer or the user liable?*

- Companies: Companies deploying AI systems must ensure they follow legal and ethical guidelines, and mitigate risks such as discrimination, misinformation, or unintended consequences.
- Developers: Responsible for ensuring that AI systems are fair, safe, transparent, and do not perpetuate harm or bias.
- Users: Individuals or entities that use AI systems responsibly. They should understand the limitations of AI, follow ethical guidelines in its application, and avoid misuse.

# Key Questions about AI Accountability

- What is accountability?
  - An obligation or willingness to accept responsibility or to account for one's action.
- How to achieve accountability?
  - AI systems need to be **transparent** (people should understand how decisions are made) and **explainable** (the decision-making process should be clear enough).
  - Regulatory frameworks should be provided to ensure accountability.

# Challenges in AI Accountability

- Complexity of AI Systems: Advanced AI systems (e.g., deep learning-based) often function as black boxes, making it difficult to trace specific decision-making steps.
- Shared Responsibility: AI systems often rely on multiple stakeholders (developers, users, etc.), so it can be hard to pinpoint exactly who is responsible when something goes wrong.
- Autonomous Decision-Making: Making decisions without human intervention challenges traditional notions of accountability, as the AI may make choices its creators did not explicitly intend.

# This Lecture

- Accountability
- Detecting AI-generated Content
  - An important aspect of accountability of AI in the era of generative AI
- Watermarking Techniques
- Evading Watermarking-based Detection

# AI-generated Content VS Deepfake

- **AI-generated content** refers to any content created by AI algorithms, e.g. **text** generation, **image** and art creation, **video** and animation, **music** composition, and **voice** synthesis.
- **Deepfakes** specifically refer to synthetic media where **a person in an existing image or video** is replaced with someone else's likeness. It uses deep learning techniques to make these alterations appear as realistic as possible.

# Generative AI for Image

Generative model	Size	Number	Source
ProGAN [19]	256x256	8.0k	LSUN
StyleGAN [20]	256x256	12.0k	LSUN
BigGAN [3]	256x256	4.0k	ImageNet
CycleGAN [52]	256x256	2.6k	ImageNet
StarGAN [6]	256x256	4.0k	CelebA
GauGAN [33]	256x256	10.0k	COCO
Stylegan2 [21]	256x256	15.9k	LSUN
WFIR [49]	1024x1024	2.0k	FFHQ
ADM [9]	256x256	12.0k	ImageNet
Glide [31]	256x256	12.0k	ImageNet
Midjourney [30]	1024x1024	12.0k	ImageNet
SDv1.4 [39]	512x512	12.0k	ImageNet
SDv1.5 [39]	512x512	16.0k	ImageNet
VQDM [15]	256x256	12.0k	ImageNet
Wukong [50]	512x512	12.0k	ImageNet
DALL-E 2 [37]	256x256	2k	ImageNet

<https://www.whichfaceisreal.com/> (StyleGAN)

# Images from StyleGAN



Water-splotches



Eyeglasses



Background problems



Other asymmetries, hair, teeth, Fluorescent bleed

# The techniques are evolving, but the methodologies persist...



Image Generated by GPT4o

A wide image taken with a phone of a glass whiteboard, in a room overlooking the Bay Bridge. The field of view shows a woman writing, sporting a tshirt with a large OpenAI logo. The handwriting looks natural and a bit messy, and we see the photographer's reflection.

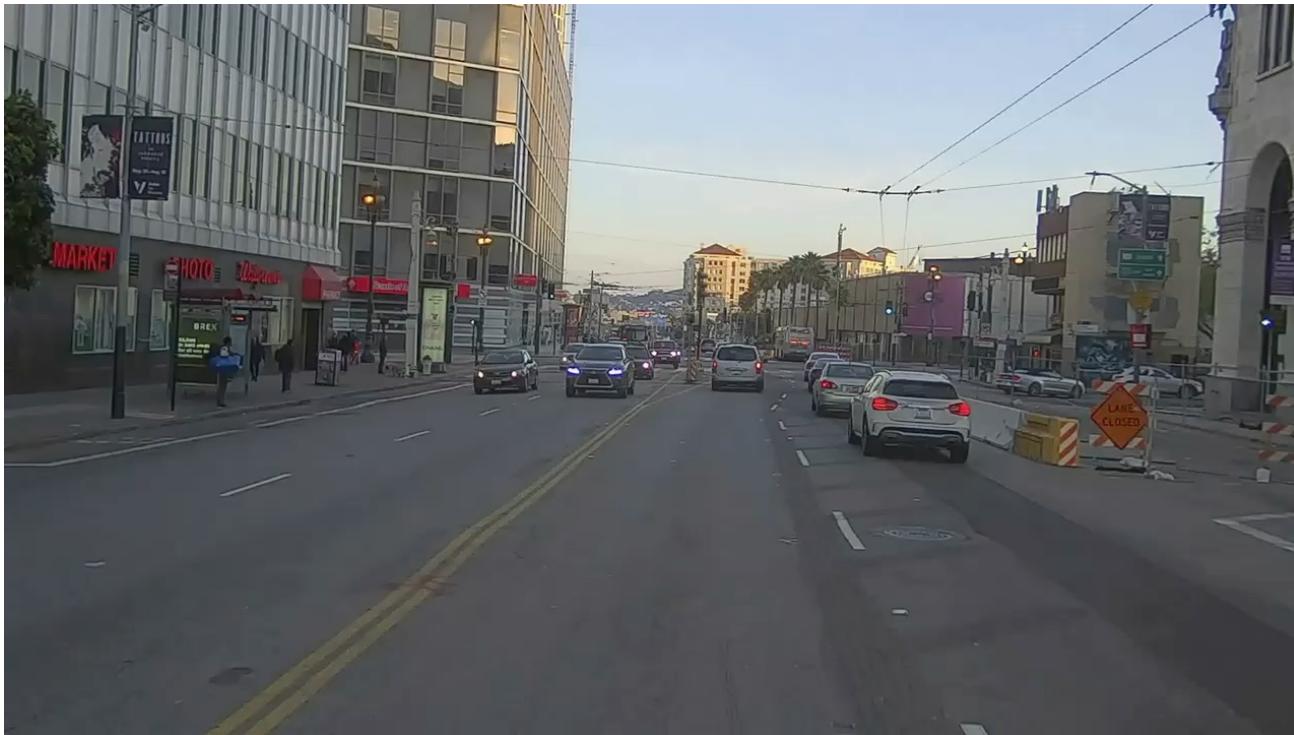
...  
Read more

Transfer between Modalities:  
Suppose we directly model p(text, pixels, sound) with one big autoregressive prior.  
Pros:  
- image generation augmented with world knowledge  
- next-level text rendering  
- native in-context learning  
- unified post-training edits  
Cons:  
- varying bit-rate  
- compute not ada  
Fixes:  
+ model compressed representations  
+ compose autoregressive prior with a powerful decoder  
tokens → [transformer] → [diffusion] pixels

selfie view of the photographer, as she turns around to high five him

Transfer between Modalities:  
Suppose we directly model p(text, pixels, sound) with one big autoregressive prior.  
Pros:  
- image generation augmented with world knowledge  
- next-level text rendering  
- native in-context learning  
- unified post-training edits  
Cons:  
- varying bit-rate  
- compute not ada  
Fixes:  
+ model compressed representations  
+ compose autoregressive prior with a powerful decoder  
tokens → [transformer] → [diffusion] pixels

# Anything in Any Scene Photorealistic Video Object Insertion



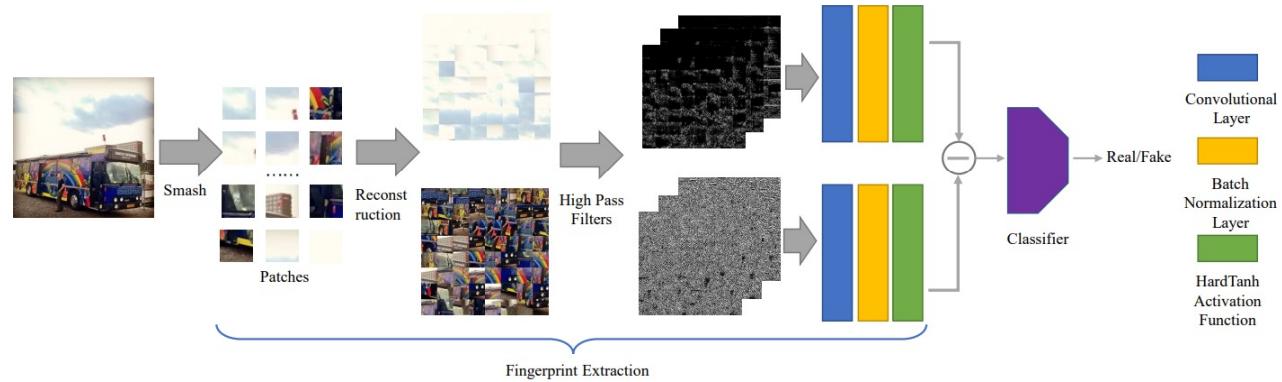
# Detecting AI-generated Content

- Passive Detection
  - Key idea: leverage artifacts in AI-generated content
  - High false positives/negatives
  - Abandoned by OpenAI
- Watermark-based Detection
  - Deployed by Google, Microsoft, OpenAI, Stability AI, etc.

# Passive Detection

- Binary classification: cannot extend to unseen generative models
  - But data augmentation can significantly boost the generalization for the detector.
- Universal fingerprint across different generative models
  - Consistently anomalous behavior of fake images in the frequency domain.
    - Up-sampling operation results in the frequency abnormality of fake images.
  - Image generated by a diffusion model are more likely to be accurately reconstructed by a pre-trained diffusion model.
    - Previous researchers adopt the reconstruction error of an image as the fingerprint to identify diffusion-generated fake images.

# PatchCraft: Exploring Texture Patch for Efficient AI-generated Image Detection (1)



- Inter-pixel correlation is determined by its camera device (CMOS) and ISP (Image Signal Processing).
- The inter-pixel correlations of real images between rich and poor texture regions are very close.

# PatchCraft: Exploring Texture Patch for Efficient AI-generated Image Detection (2)

- Texture diversity measurement

$$l_{div} = \sum_{i=1}^M \sum_{j=1}^{M-1} (|x_{i,j} - x_{i,j+1}|) + \sum_{i=1}^{M-1} \sum_{j=1}^M (|x_{i,j} - x_{i+1,j}|) \\ + \sum_{i=1}^{M-1} \sum_{j=1}^{M-1} (|x_{i,j} - x_{i+1,j+1}|) + \sum_{i=1}^{M-1} \sum_{j=1}^{M-1} (|x_{i+1,j} - x_{i,j+1}|).$$

- High-pass filters are conducive to **suppressing** semantic information and **magnifying** the interpixel correlation

# PatchCraft: AI-generated Image Detection

Generator	CNNSpot	FreDect	Fusing	GramNet	LNP	LGrad	DIRE-G	DIRE-D	UnivFD	Ours
ProGAN	<b>100.00</b>	99.36	<b>100.00</b>	<u>99.99</u>	99.67	99.83	95.19	52.75	99.81	<b>100.00</b>
StyleGAN	90.17	78.02	85.20	87.05	<u>91.75</u>	91.08	83.03	51.31	84.93	<b>92.77</b>
BigGAN	71.17	81.97	77.40	67.33	77.75	85.62	70.12	49.70	<u>95.08</u>	<b>95.80</b>
CycleGAN	<u>87.62</u>	78.77	87.00	86.07	84.10	86.94	74.19	49.58	<b>98.33</b>	70.17
StarGAN	94.60	94.62	97.00	95.05	<u>99.92</u>	99.27	95.47	46.72	95.75	<b>99.97</b>
GauGAN	<u>81.42</u>	80.57	77.00	69.35	75.39	78.46	67.79	51.23	<b>99.47</b>	71.58
Stylegan2	86.91	66.19	83.30	87.28	<b>94.64</b>	85.32	75.31	51.72	74.96	<u>89.55</u>
WFIR	<b>91.65</b>	50.75	66.80	86.80	70.85	55.70	58.05	53.30	<u>86.90</u>	85.80
ADM	60.39	63.42	49.00	58.61	84.73	67.15	75.78	<b>98.25</b>	66.87	<u>82.17</u>
Glide	58.07	54.13	57.20	54.50	80.52	66.11	71.75	<b>92.42</b>	62.46	<u>83.79</u>
Midjourney	51.39	45.87	52.20	50.02	65.55	65.35	58.01	<u>89.45</u>	56.13	<b>90.12</b>
SDv1.4	50.57	38.79	51.00	51.70	85.55	63.02	49.74	<u>91.24</u>	63.66	<b>95.38</b>
SDv1.5	50.53	39.21	51.40	52.16	85.67	63.67	49.83	<u>91.63</u>	63.49	<b>95.30</b>
VQDM	56.46	77.80	55.10	52.86	74.46	72.99	53.68	<b>91.90</b>	85.31	<u>88.91</u>
Wukong	51.03	40.30	51.70	50.76	82.06	59.55	54.46	<u>90.90</u>	70.93	<b>91.07</b>
DALLE2	50.45	34.70	52.80	49.25	88.75	65.45	66.48	<u>92.45</u>	50.75	<b>96.60</b>
Average	70.78	64.03	68.38	68.67	<u>83.84</u>	75.34	68.68	71.53	78.43	<b>89.31</b>

PatchCraft (Ours) <https://fdmas.github.io/AIGCDetect/data/>

# This Lecture

- Accountability
- Detecting AI-generated Content
- Watermarking Techniques
- Evading Watermarking-based Detection

# Watermarking Techniques

- Creators embed identifiable markers or metadata into AI models, datasets, or outputs to trace their origin or prove ownership.
- Watermarking techniques in AI systems are a potential solution for ensuring accountability, particularly in the context of AI-generated content, intellectual property protection, and responsibility.

# AI Model Watermarking

- White-box Watermarking
  - This technique involves embedding the watermark directly into the model's parameters or structure.
  - It requires access to the internal architecture of the model and can be used to verify ownership or identify unauthorized modifications.
- Black-box Watermarking
  - Black-box watermarking embeds a secret, recognizable pattern in the model's outputs rather than its internal structure.
  - It doesn't require access to the internal workings of the AI model and can be verified through observing the model's responses to specific inputs.

# Dataset Watermarking

- Perturbation Watermarking
  - This involves making small, controlled modifications (perturbations) to the data, which do not affect its usability but serve as a watermark to prove dataset ownership or trace responsibility.
- Data Label Watermarking
  - A form of watermarking where specific, unique patterns are embedded in the labels of the dataset.
  - This helps identify if a model has been trained on a specific dataset, even when used by unauthorized entities.

# Watermarking in Large Language Models

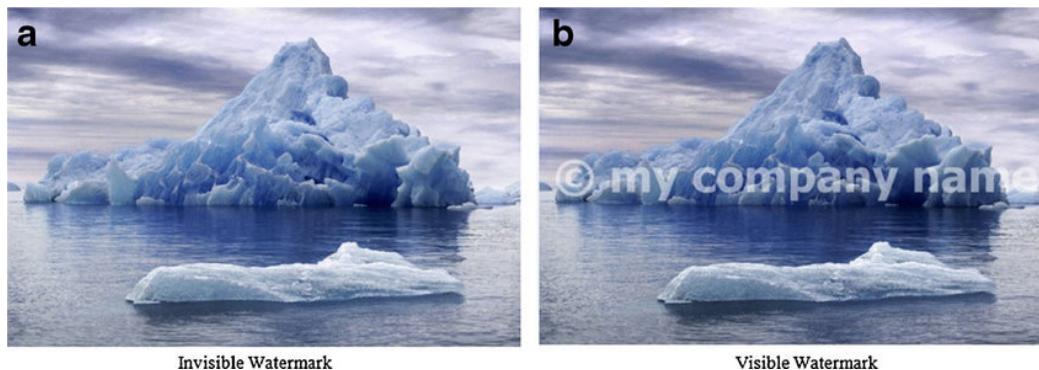
- Token-based Watermarking
  - Inserting specific patterns or sequences of words or tokens into generated text in a way that is hard for a human to detect but easily identified by automated tools.
- Sentence Structure Watermarking
  - Subtly altering sentence structures or punctuation in the generated text to leave an identifiable watermark while maintaining the overall meaning.

# Watermarking Quality

- Quality. The watermark should be imperceptible to humans, meaning the watermarked version appears identical to the original in raw data modality.
- Generalization. The watermark can be applied to any samples, and the watermarking method is compatible with different models.
- Robustness. The watermark cannot be compromised by existing attack methods.

# Digital Watermarking in AI-generated Image

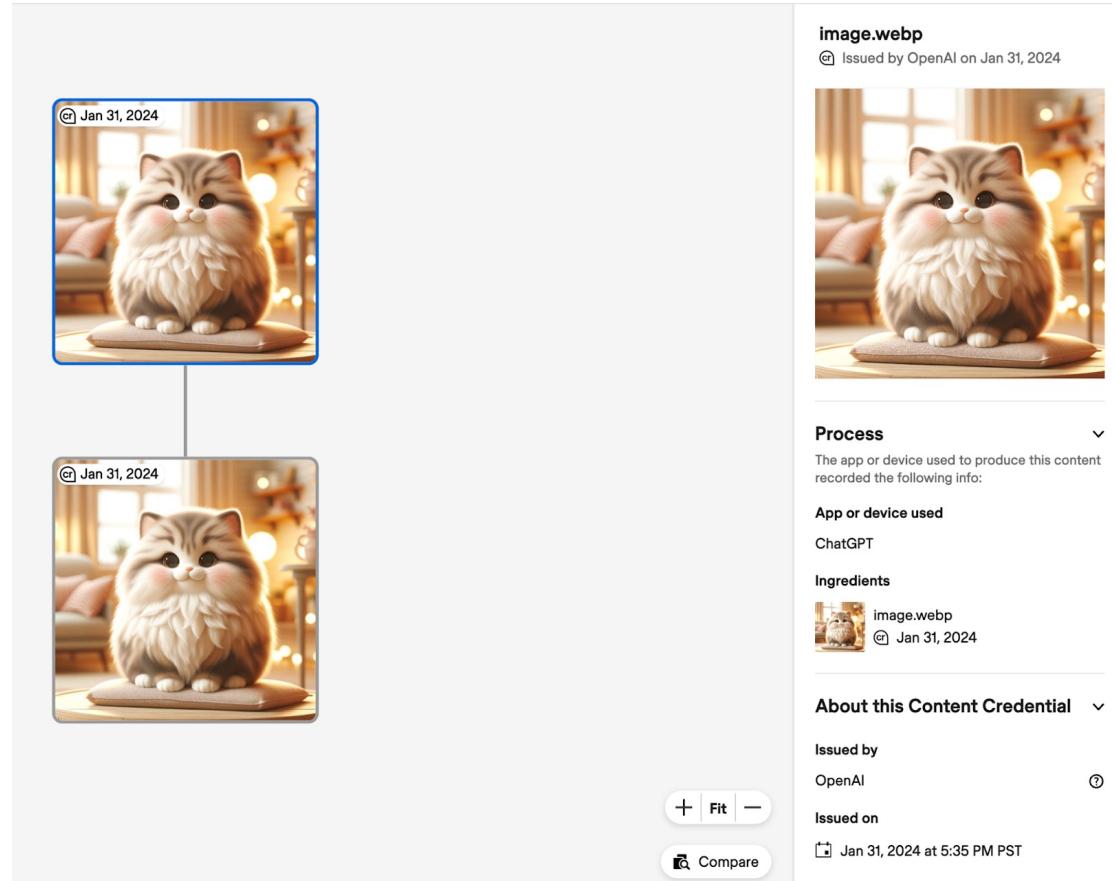
- Invisible Watermarking: embedding markers in a way that is imperceptible to humans but detectable using specific algorithms.
- Visible Watermarking: Placing a visible marker, like a logo or text, onto AI-generated images, videos, or other multimedia content.



<https://link.springer.com/article/10.1007/s11042-013-1515-8>

# C2PA

- Coalition for Content Provenance and Authenticity (C2PA)
- <https://c2pa.org/>
- Taking a screenshot can remove such metadata



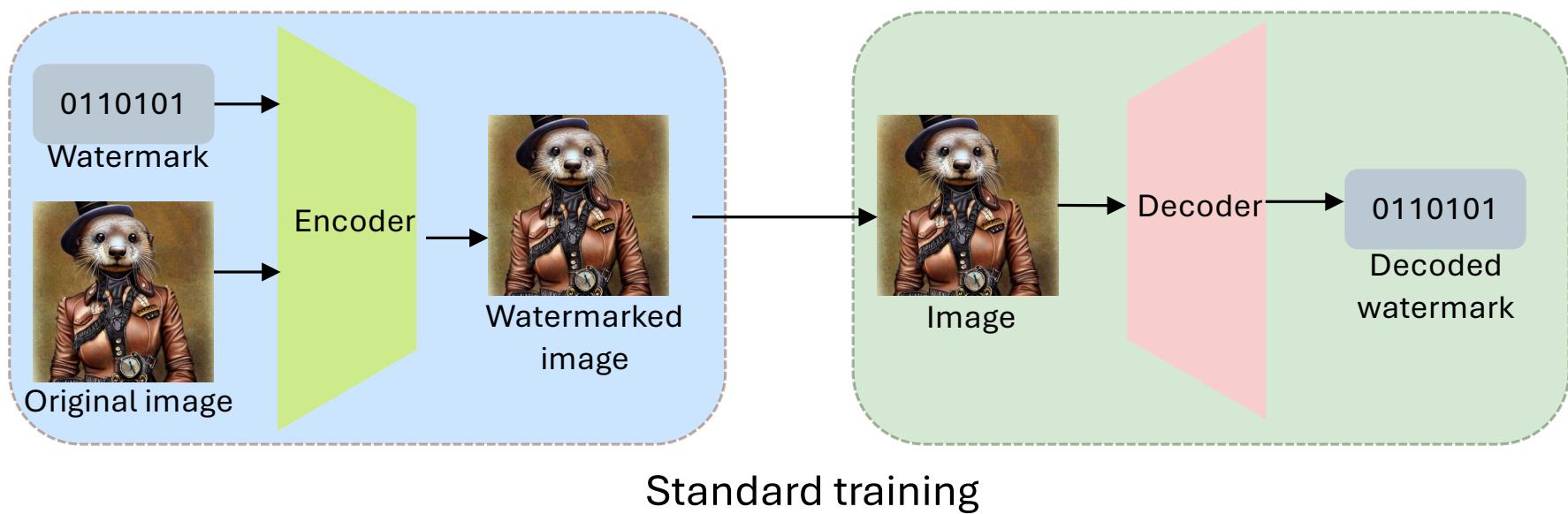
<https://help.openai.com/en/articles/8912793-c2pa-in-chatgpt-images>

# None-learning-based Watermarking Methods

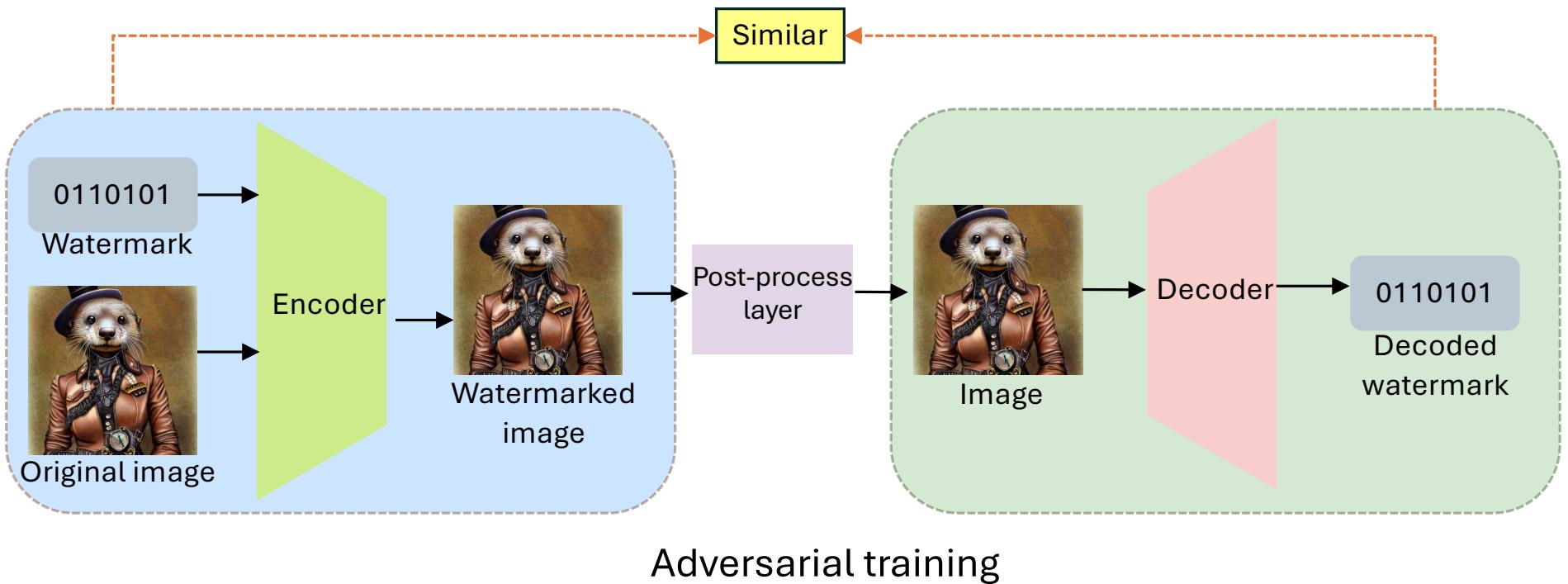
- **dwtDct**: DWT + DCT transform, embed watermark bit into max non-trivial coefficient of block DCT coefficients
- **dwtDctSvd**: DWT + DCT transform, SVD decomposition of each block, embed watermark bit into singular value decomposition
- Discrete wavelet transform (DWT)
- Discrete cosine transform (DCT)
- <https://github.com/ShieldMnt/invisible-watermark>

# Learning-based Watermarking Methods

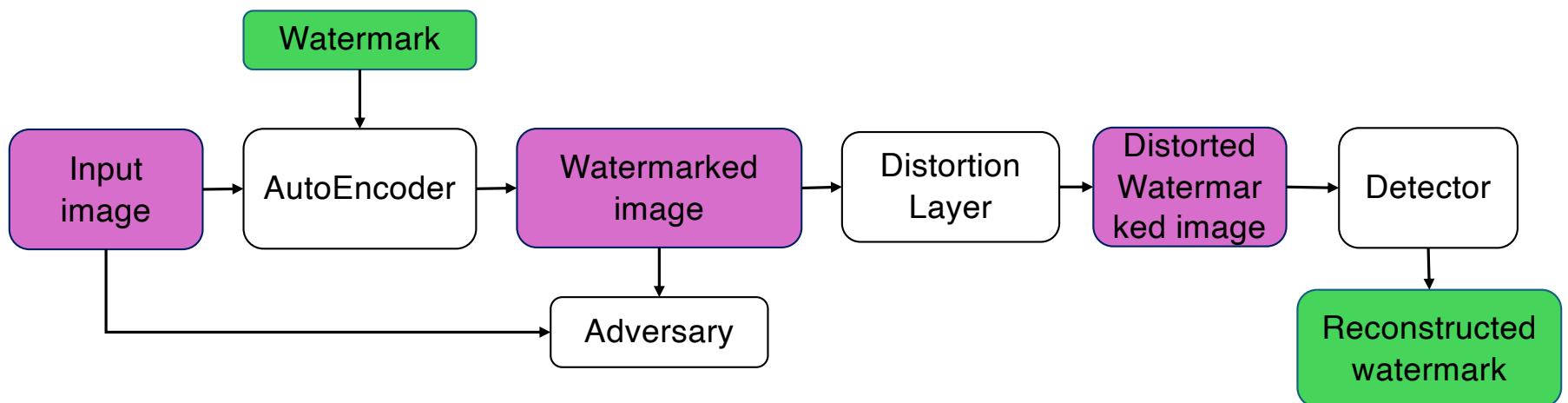
Three main components: Watermark (bitstring), Encoder, Decoder



# Adversarial Training for Watermarking

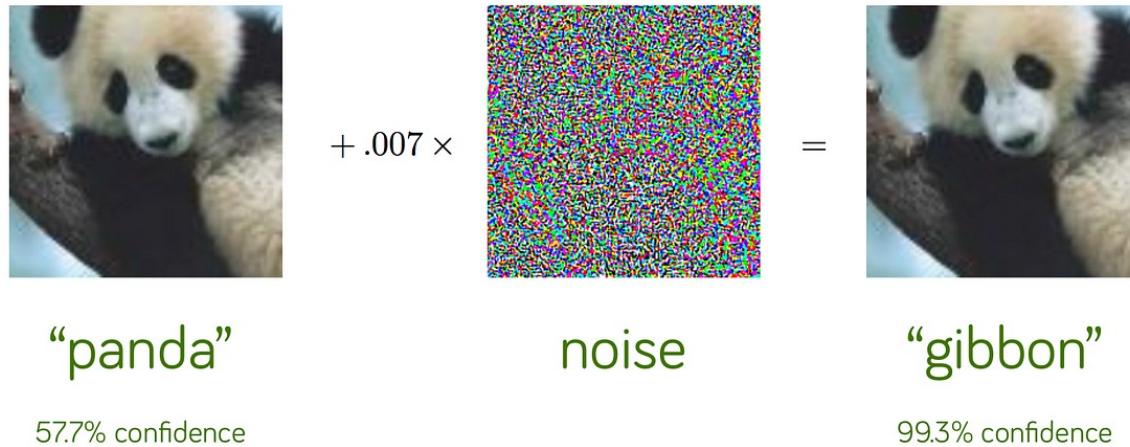


# A General Framework for Training Watermarking Model



- **AutoEncoder**. Embed the watermark into the image, producing an encoded version of it.
- **Adversary**. Improve the quality of the watermarked image and minimize the domain gap between the original and the watermarked image.
- **Detector**. Receive the encoded images or audio and output the extracted message.
- **Distortion Layer**. Simulate potential attack scenarios.

# Watermarking and Adversarial Perturbation



Can DNNs **extract meaningful information** from imperceptible perturbations?

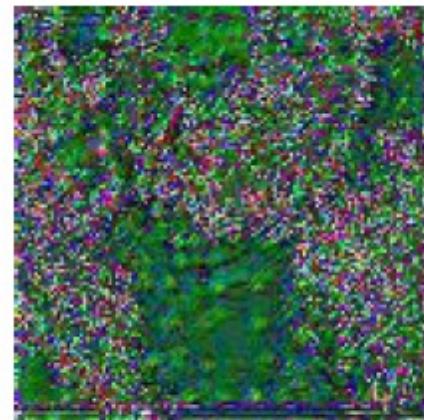
# Watermarking and Adversarial Perturbation

Digital watermarking can be used to identify image ownership



Cover Image

+



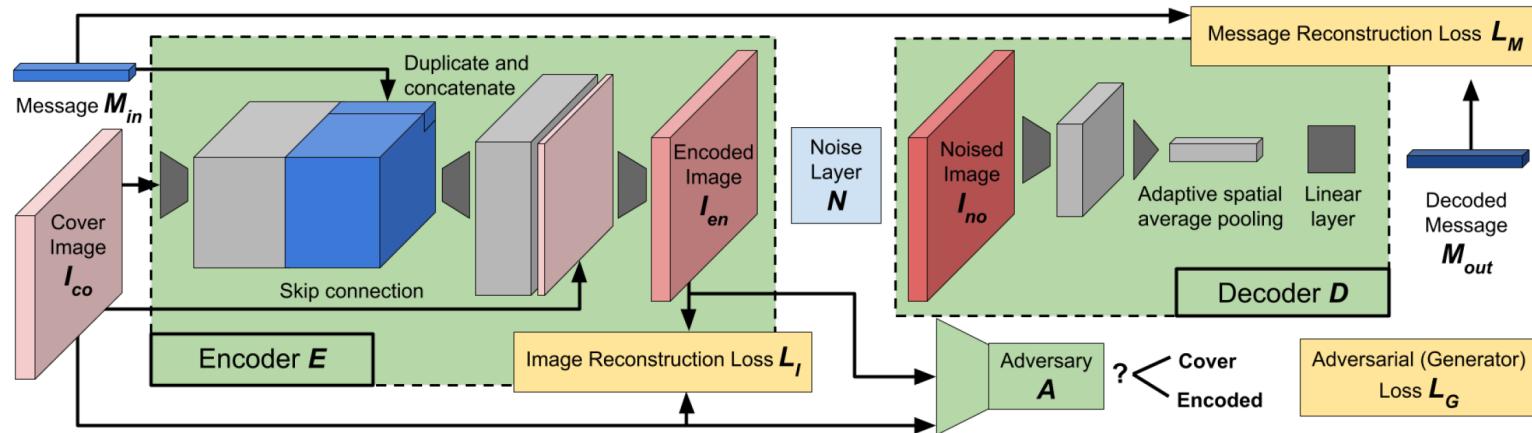
HiDDeN  
Perturbation

=



*“Copyright ID: 1337”*

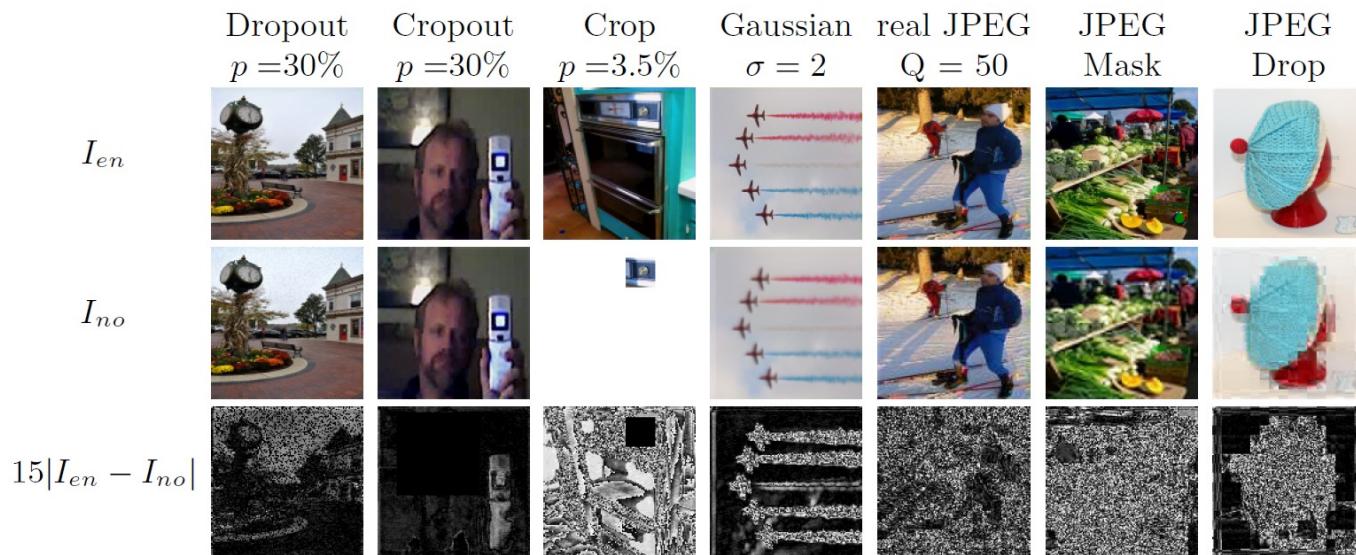
# HiDDeN: Hiding Data With Deep Networks



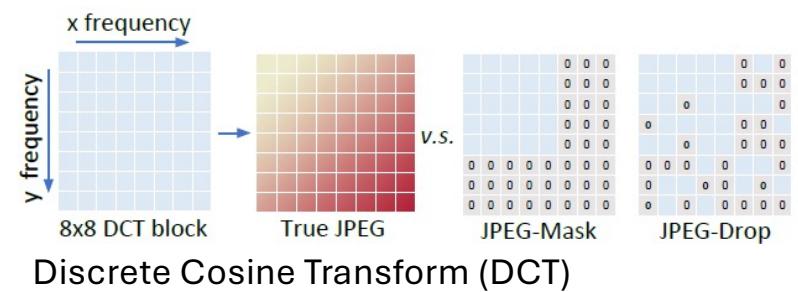
- An **autoencoder** embeds a message into a cover image and outputs an encoded image.
- A **decoder** receives the encoded image and extracts the embedded message.
- The **adversary** predicts whether a given image contains an encoded message.
- The **noisy layer** applies various image transformations to improve robustness against various image distortion scenarios.

<https://arxiv.org/pdf/1807.09937>

# HiDDeN: Distortion Method



- JPEG-Mask zeros a fixed set of high-frequency coefficients
- JPEG-Drop zeros channels with higher drop probabilities for high-frequency coefficients



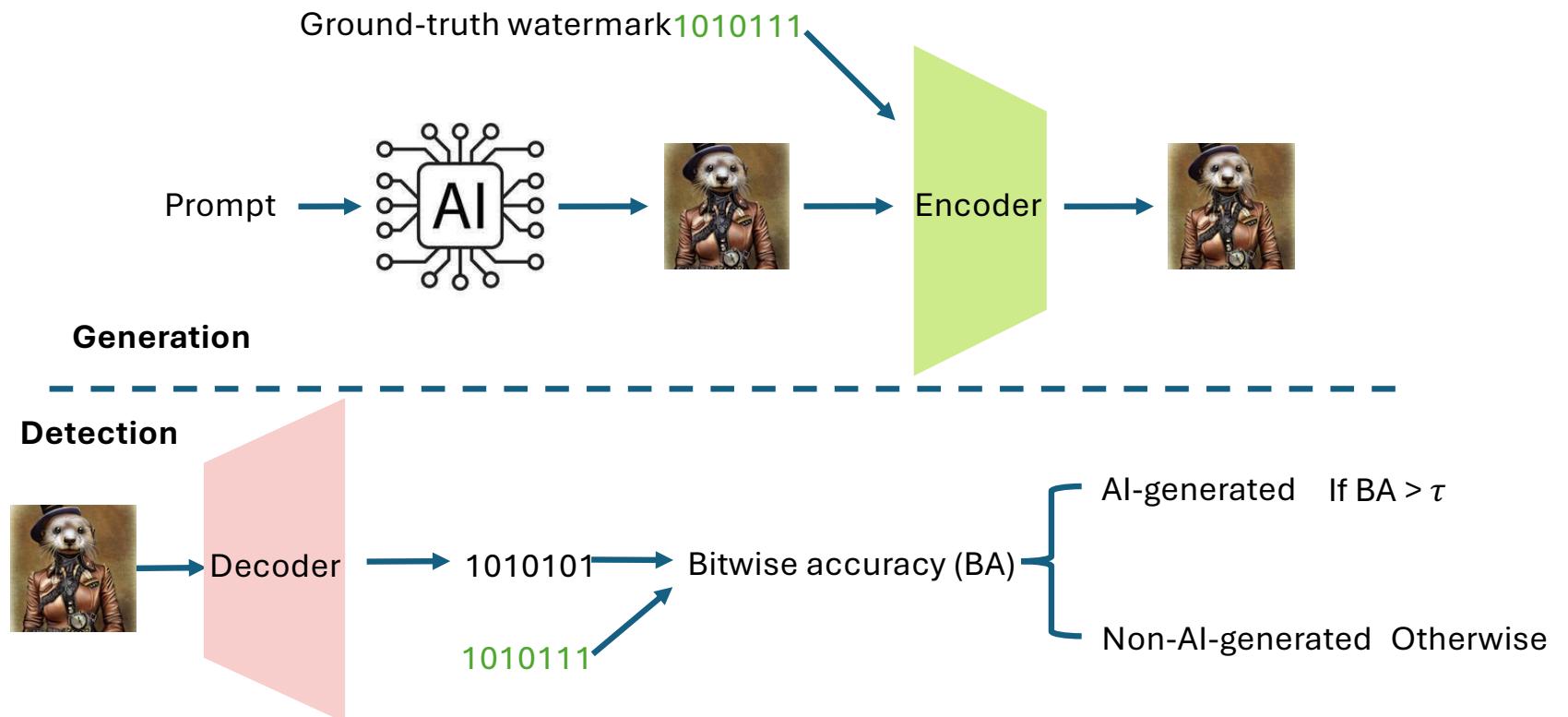
# HiDDeN: Data Quality Results

The image distortion between the cover and stego image using peak signal-to-noise ratio (PSNR)

	Digimarc	Identity	Dropout	Cropout	Crop	Gaussian	JPEG-mask	JPEG-drop	Combined
PSNR(Y)	62.12	44.63	42.52	47.24	35.20	40.55	30.09	28.79	33.55
PSNR(U)	38.23	45.44	38.52	40.97	33.31	41.96	35.33	32.51	38.92
PSNR(V)	52.06	46.90	41.05	41.88	35.86	42.88	36.27	33.42	39.38
<hr/>									
Cover		Trained with Adversary			No Adversary				
Cover	Digimarc	Crop	Gaussian	Combined	Combined	Combined	Combined		
									

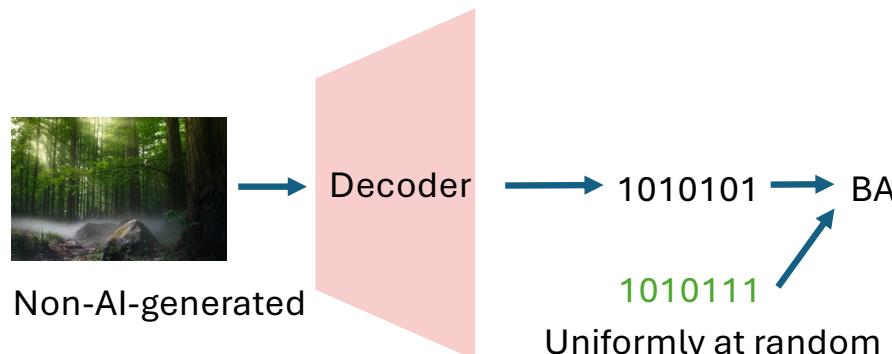
- Encoded images from the HiDDeN model are **visually indistinguishable** from the cover image
- The model trained against an **adversarial discriminator** produces images with **no visible artifacts**

# Watermark-based Detection



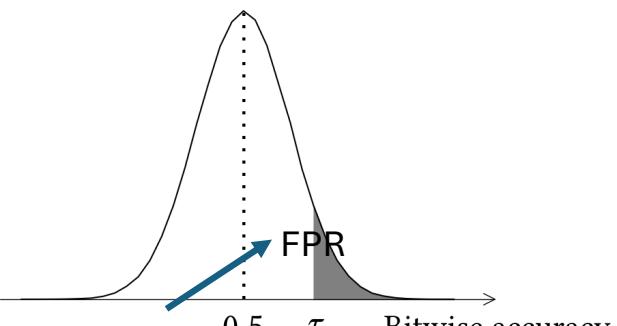
# How to Set Detection Threshold $\tau$ ?

Goal: achieve a desired False Positive Rate (FPR)



n: watermark length (number of experiments)  
0.5: success rate  
BA: Bitwise Accuracy

$$n \cdot BA \sim \text{Binomial}(n, 0.5)$$



$$\begin{aligned} FPR_s(\tau) &= \Pr(BA(D(I_o), w) > \tau) \\ &= \Pr(m > n\tau) = \sum_{k=\lceil n\tau \rceil}^n \binom{n}{k} \frac{1}{2^n}, \end{aligned}$$

$$\begin{array}{ccc} FPR < 10^{-4} & \xrightarrow{\hspace{1cm}} & \tau > 0.83 \\ 30\text{-bit watermark} & & \end{array}$$

<https://arxiv.org/abs/2305.03807>

<https://github.com/zhengyuan-jiang/WEvade>

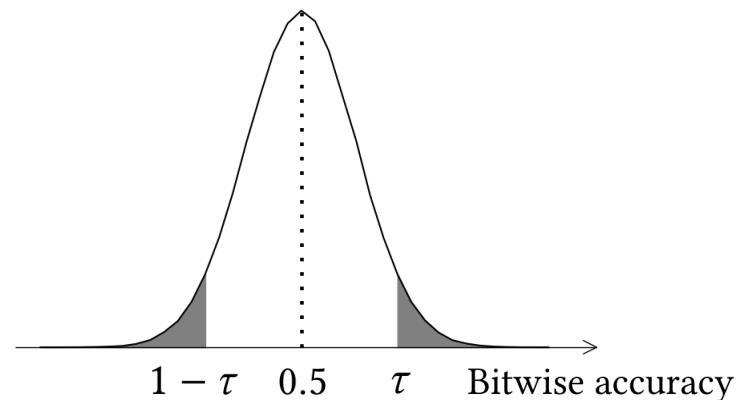
# Problem of Single-tail

- Single-tail detector can be easily evaded by simply extending standard adversarial examples to watermarking.
- Adds perturbation to a watermarked image such that the decoded watermark has a very small bitwise accuracy → 0. This will evade the watermarking-based detection.

# Important Observation

- The watermarks decoded from original images have bitwise accuracy close to 0.5, while those decoded from watermarked images have large bitwise accuracy, e.g., close to 1.
- If the bitwise accuracy of the watermark decoded from an image is significantly smaller than 0.5, it is likely to be an adversarial perturbed image.

To avoid evading detection, using double-tail detector.



AI-generated: If  $BA > \tau$  or  $BA < 1 - \tau$

Non-AI-generated: Otherwise

# This Lecture

- Accountability
- Detecting AI-generated Content
- Watermarking Techniques
- Evading Watermarking-based Detection

# Question: double-tail detector safe to use?

- Intuition: **non-watermarked images have bitwise accuracy  $\approx 0.5$**
- Attack the detector: add minimal perturbation to **make bitwise accuracy  $\approx 0.5$** , making the perturbed watermarked image indistinguishable with non-watermarked ones

$$\begin{aligned} \min_{\delta} \quad & ||\delta||_{\infty} \\ \text{s.t. } & |BA(D(I_w + \delta), w) - 0.5| \leq \epsilon \\ & \text{Decoder} \end{aligned}$$

Ground-truth watermark: Unknown to attacker



Constraints are non-linear, hard to solve!

# Solve the Optimization Problem

$$\min_{\delta} \|\delta\|_{\infty}$$

$$s.t. |BA(D(I_w + \delta), w) - 0.5| \leq \epsilon$$

$$\min_{\delta} l(D(I_w + \delta), w_t)$$

$$s.t. \|\delta\|_{\infty} \leq r,$$

$$BA(D(I_w + \delta), w_t) \geq 1 - \epsilon,$$

Target watermark is a random watermark

Iteratively find the perturbation  $\delta$  that satisfies the constraints (if possible) for a given  $r$ .

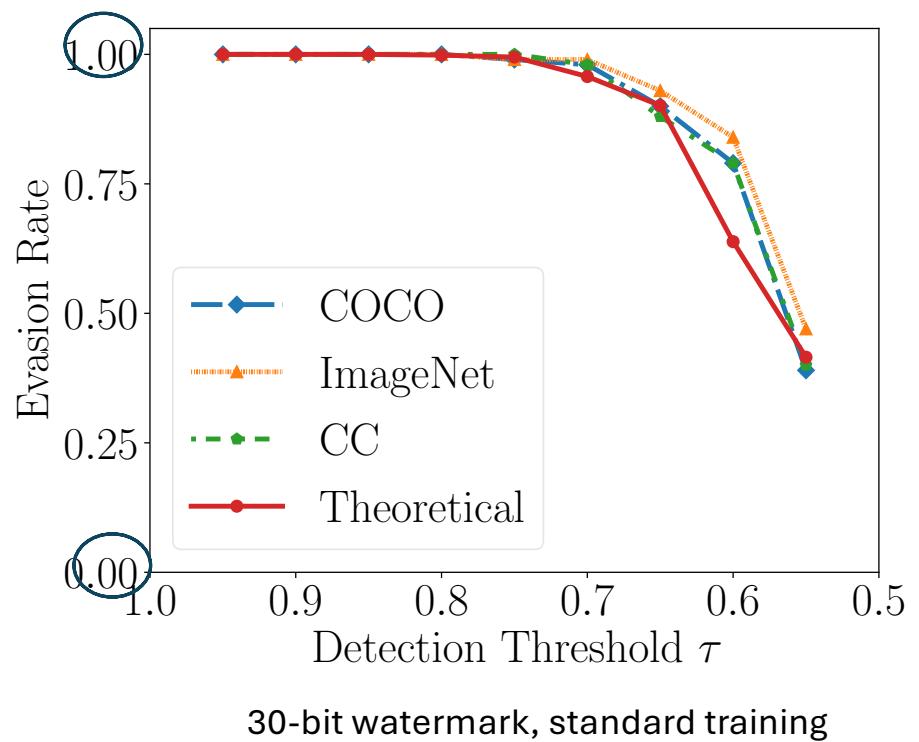
Binary search over  $r$  to find the smallest perturbation  $\delta$  that satisfies the constraints

$[r_a, r_b]$  is initialized such that  $r_a = 0$  and  $r_b$  is a large value

$w_t$  bitwise accuracy close to 0.5 compared to any ground-truth watermark

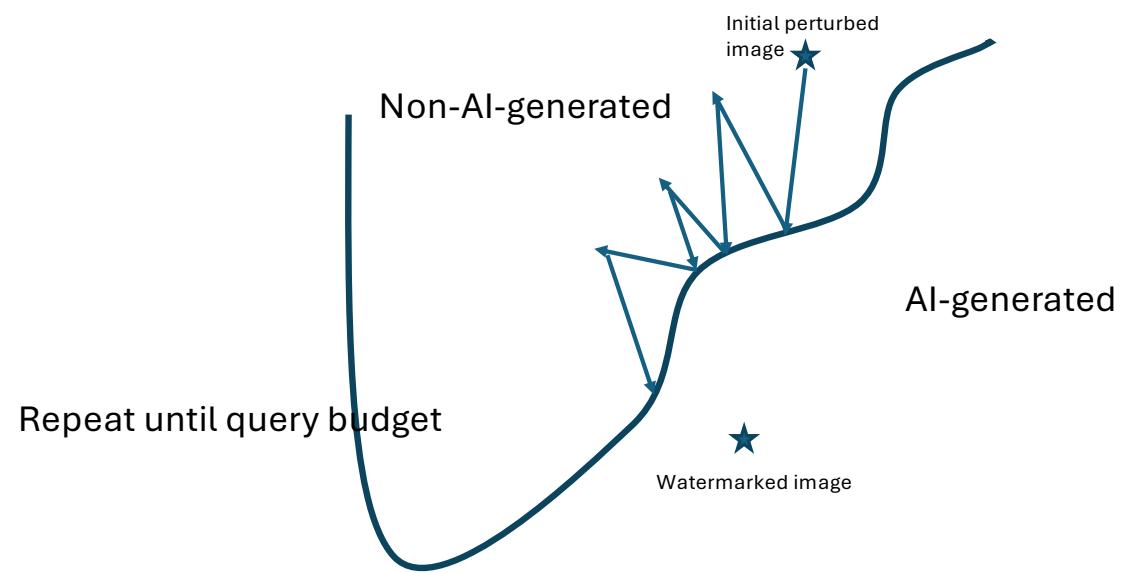
# Empirical Evaluation Results

- Detection Threshold:
  - Fraction of bitwise accuracy  $> \tau$
- Evasion rate: probability that a perturbed watermarked image is detected as non-AI-generated



# Black-box Attack to Evade Watermarking-based Detection

- Surrogate Model
- Query-based



**WEvade-B-Q:** WEvade-B-Q starts from an initial  $I_{pw}$  that evades the target detector. During the iterative process to reduce the perturbation, WEvade-B-Q always guarantees that  $I_{pw}$  evades detection.

# References

- HiDDeN: Hiding Data With Deep Networks  
(<https://arxiv.org/pdf/1807.09937>)
- PatchCraft: Exploring Texture Patch for Efficient AI-generated Image Detection (<https://arxiv.org/pdf/2311.12397>)
- Evading Watermark based Detection of AI-Generated Content  
(<https://arxiv.org/pdf/2305.03807>)