

***WRATH*: Turning Watermark Robustness Against Itself via a Watermark-Agnostic Black-Box Invalidation Attack**

*Nan Jiang, *Juan Hu, *Bangjie Sun, *Terence Sim, †Jun Han

*National University of Singapore, †KAIST

Contact: jiangnan@u.nus.edu



YOU CAN DO HARD THINGS

BE KIND ONLINE

THINK BEFORE YOU CLICK

KEYBOARD SHORTCUTS
Ctrl + C Copy
Ctrl + V Paste
Ctrl + X Cut
Ctrl + Z Undo
Ctrl + S Save
Ctrl + P Print
Alt + Tab Switch Windows
Ctrl + F Find

ChatGPT

ChatGPT

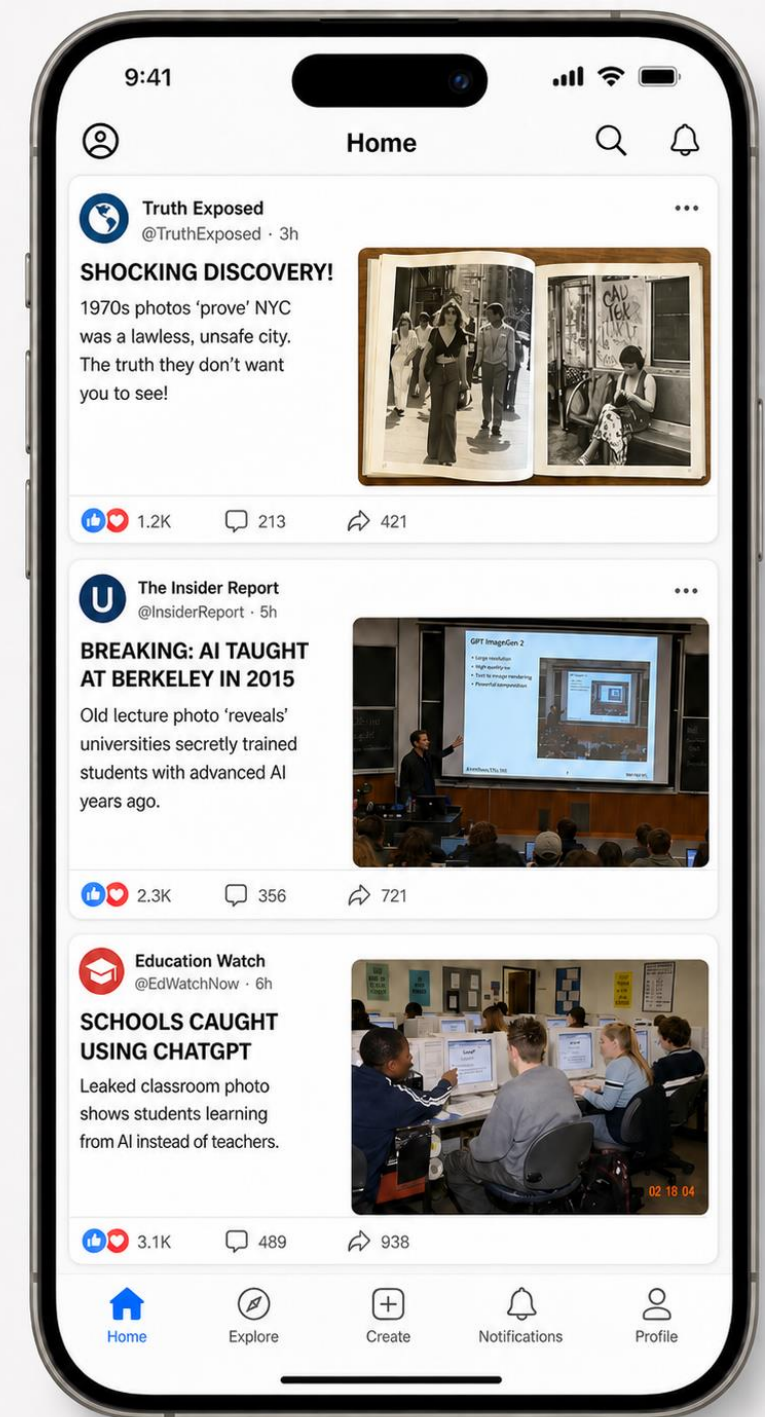
ChatGPT
Ask anything
How do glaciers move?
Glaciers move slowly due to gravity. The ice deforms and slides over rock and boulders.

ChatGPT
Ask anything
What are binary stars?
Binary stars are two stars that orbit around their common center of mass. They are held together by gravity.

02 18 04

Photo-realistic AI-generated Images Raise Concerns

- Spread misleading information on social media
- Fabricate evidence of events that never happened
- **Undermine trust in online visual content**



Watermarking as a Trust Mechanism

- Many providers already use watermarks to identify AI-generated images



Amazon Bedrock



ChatGPT



Background: Watermarking in GenAI

- **Generation:** **Invisible** watermarks are embedded into AI-generated images

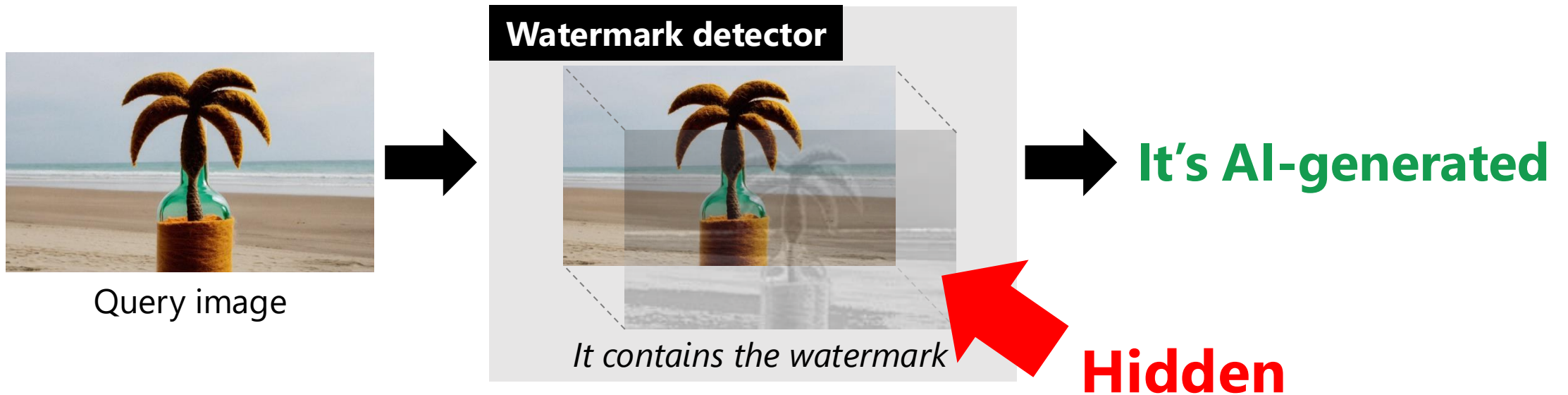
Hidden



Black-box

Background: Watermarking in GenAI

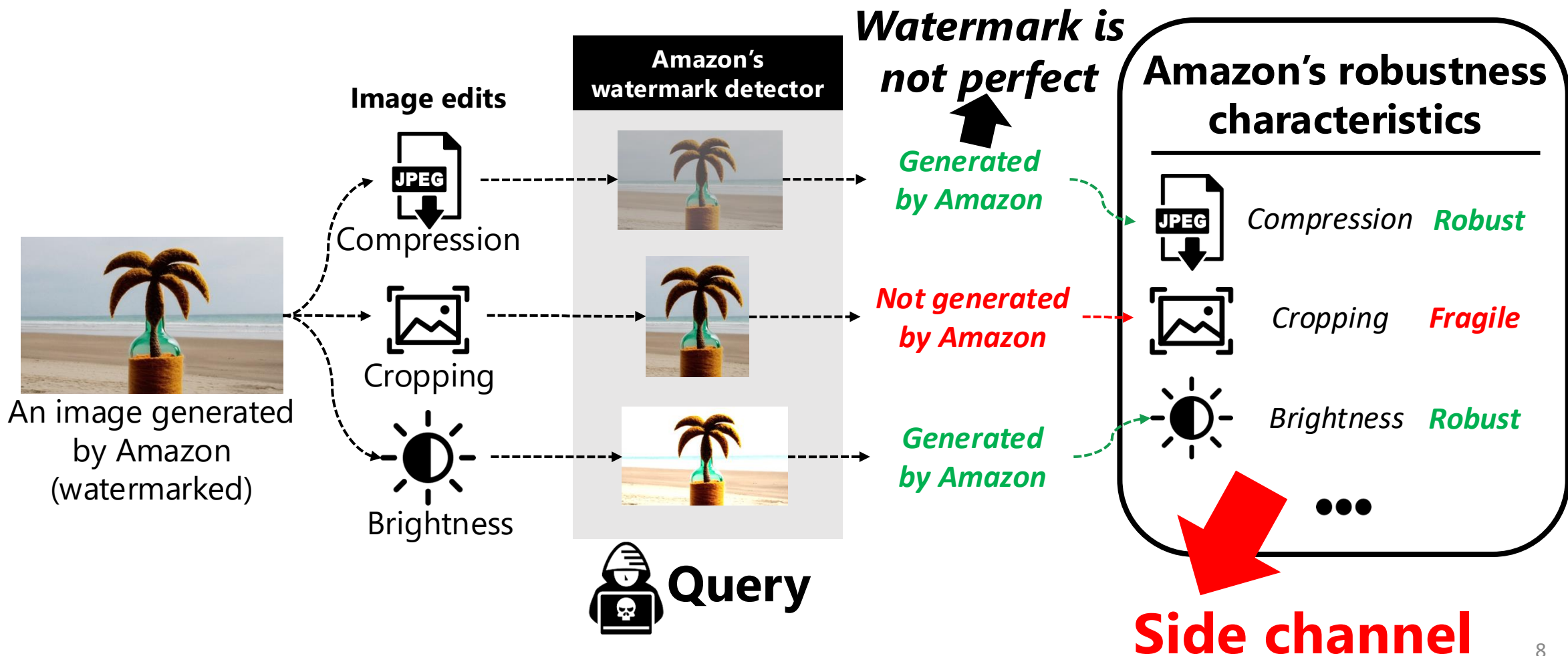
- **Generation:** **Invisible** watermarks are embedded into AI-generated images
- **Detection:** Images are deemed as AI-generated if **detected** with a watermark



Does it really reveal nothing about the watermark pattern?

Our Work: Watermark Robustness as Side Channel

- Watermark's **robustness characteristics** can reveal **watermark pattern**



How Does Robustness Reveal Watermark?

- Watermark patterns are **preserved** by **robust** edits but **disrupted** by **fragile** edits

Amazon's robustness characteristics



Compression **Robust**



Cropping **Fragile**



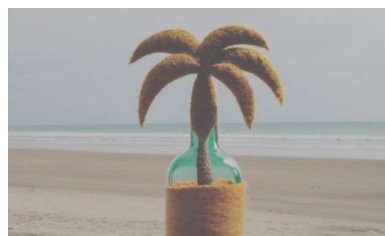
Brightness **Robust**



Robust edits



vs.



Original

Compression

Brightness

Watermark pattern \subseteq Common **preserved** patterns

Fragile edits



vs.



Original

Rotation

Cropping

Watermark pattern \subseteq Common **disrupted** patterns

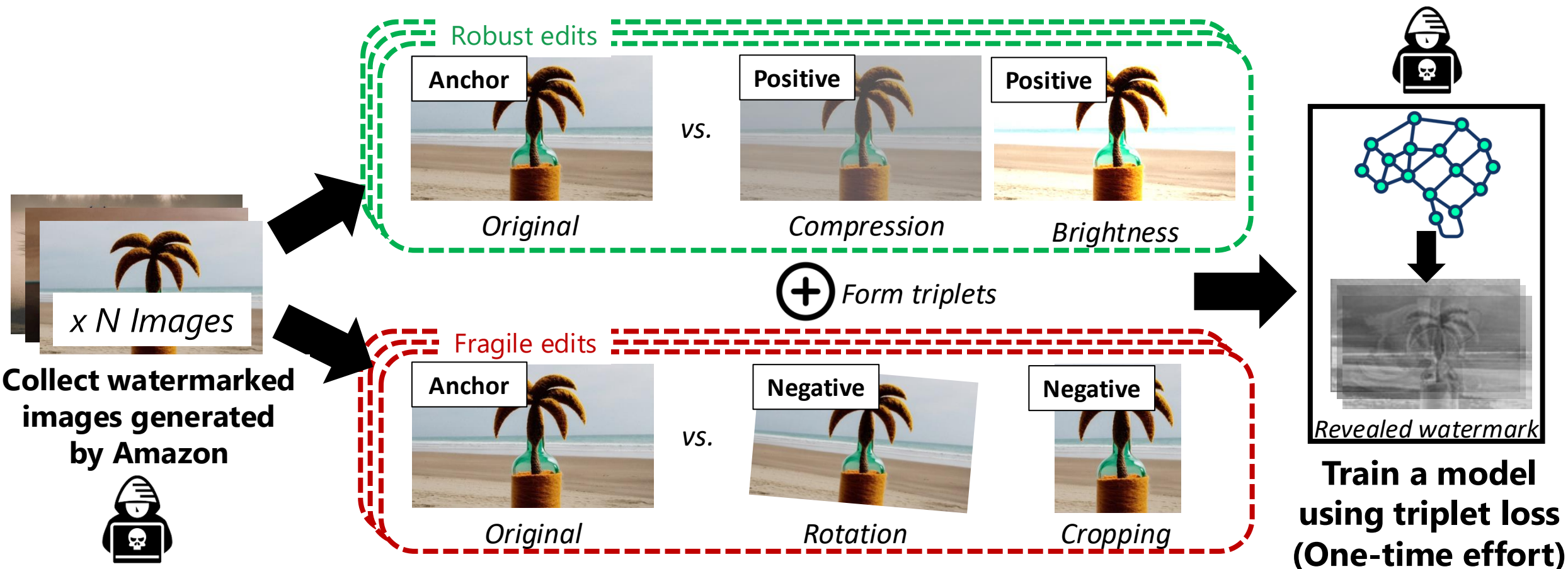
Aggregate evidence



Revealed watermark

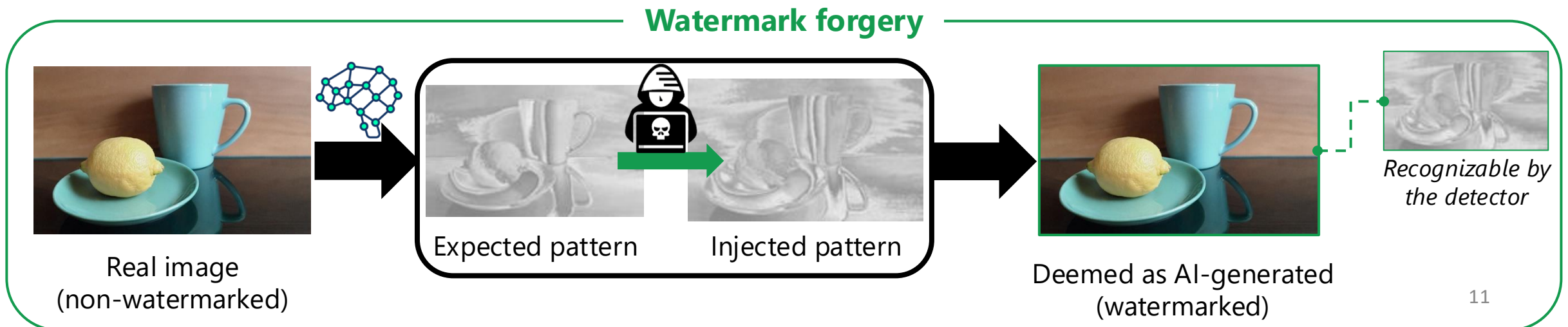
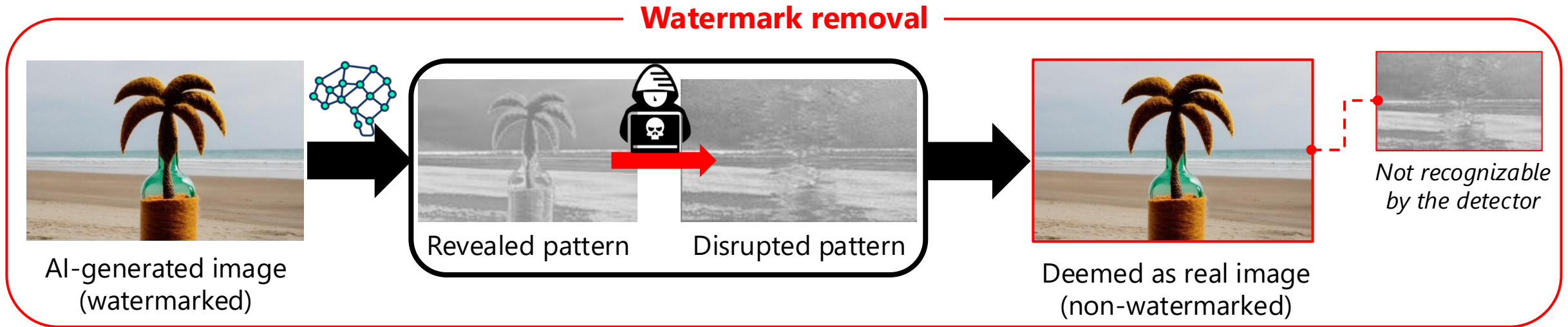
How Does Robustness Reveal Watermark?

- We aggregate evidence across images



Revealed Watermark Underpins Digital Trust

- Revealed watermark pattern enables **removal** and **forgery** attacks



Evaluation: Performance on **Removal** and **Forgery**

- **Setup:**

- **Six** representative academic works and **one** real-world system (Amazon)
- **< 100** queries to the detector; **1000** watermarked images for training

Evaluation: Performance on **Removal** and **Forgery**

- **Results:**

- **Removal:** Up to **95%** success rate while preserving perceptual quality (**LPIPS < 0.1**)
- **Forgery:** Up to **100%** success rate while preserving perceptual quality (**LPIPS < 0.1**)

Watermarked (Amazon)



Removed



Non-watermarked



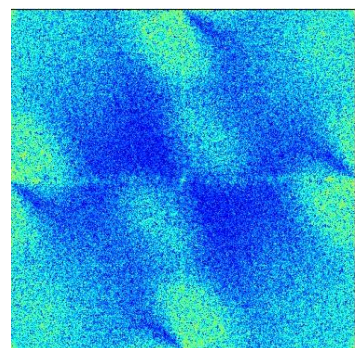
Forged (Amazon)



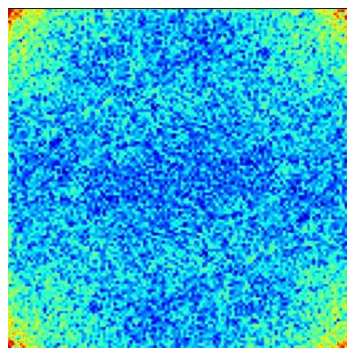
Root Cause and Countermeasure

- **Root cause:** **consistent** watermark patterns across images

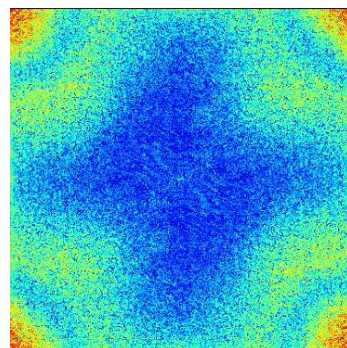
Examples in frequency domain



DWT-DCT

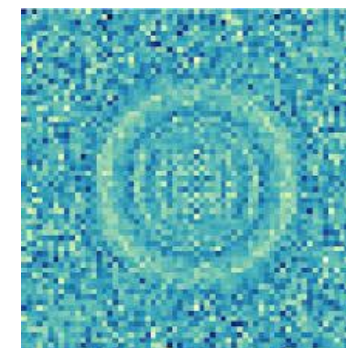


HiDDeN



StegaStamp

Example in latent domain



Tree-Ring




Allows us to **train a model** to reveal common watermark pattern

- **Countermeasure:** **randomize** watermark patterns using **cryptographic keys**

Conclusion

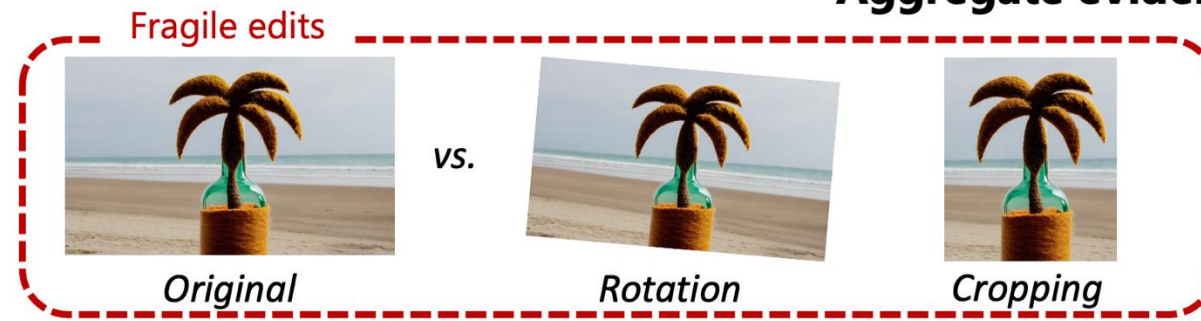
- Robustness enables watermark **revelation**, **removal**, and **forgery**
- Future watermark should **avoid consistent** watermark pattern

Amazon's robustness characteristics

	Compression	Robust
	Cropping	Fragile
	Brightness	Robust
...		



Watermark pattern \subseteq Common preserved patterns



Watermark pattern \subseteq Common disrupted patterns

Aggregate evidence



Revealed watermark