

ECML/PKDD 2020

 **WATERMARK**

Attacking Optical Character Recognition (OCR) Systems with Adversarial Watermarks



清華大學
Tsinghua University



交叉信息研究院
Institute for Interdisciplinary
Information Sciences

Lu Chen¹, Jiao Sun², Wei Xu¹

¹IIS, Tsinghua University

²CS, University of Southern California



Optical Character Recognition

Introduction

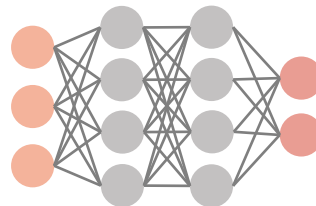
Optical Character Recognition (OCR) is a widely adopted application for converting printed or handwritten images to text, which becomes a critical preprocessing component in text analysis pipelines, such as document retrieval and summarization. OCR has been significantly improved in recent years thanks to the wide adoption of the deep neural network (DNN), and thus deployed in many critical applications where OCR's quality is vital. For example, photo-based ID recognition depends on OCR's quality to automatically structure information into databases, and automatic trading sometimes relies on OCR to read certain news articles for determining the sentiment of news.

Unfortunately, OCR also inherits all counter-intuitive security problems of the DNNs. Especially, the OCR model is also vulnerable to *adversarial examples*, which are crafted by making human-imperceptible perturbations on original images with the intent of misleading the model. The wide adoption of OCR in real pipelines gives more incentives for adversaries to game the OCR, such as causing fake ID information, incorrect readings of metrics or instructions, etc. Figure 2 and 3 in the evaluation section illustrate two real-world examples with attacking the ID number and financial

Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



Deep Neural Network



Introduction

Optical Character Recognition (OCR) is a widely adopted application for converting printed or handwritten images to text, which becomes a critical preprocessing component in text analysis pipelines, such as document retrieval and summarization. OCR has been significantly improved in recent years thanks to the wide adoption of the deep neural network (DNN), and thus deployed in many critical applications where OCR's quality is vital. For example, photo-based ID recognition depends on OCR's quality to automatically structure information into databases, and automatic trading sometimes relies on OCR to read certain news articles for determining the sentiment of news.

Unfortunately, OCR also inherits all counter-intuitive security problems of the DNNs. Especially, the OCR model is also vulnerable to adversarial examples, which are crafted by making human-imperceptible perturbations on original images with the intent of misleading the model. The wide adoption of OCR in real pipelines gives more incentives for adversaries to game the OCR, such as causing fake ID information, incorrect readings of metrics or instructions, etc. Figure 2 and 3 in the evaluation section illustrate two realworld examples with attacking the ID number and financial

Copyright c 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



PDF
FILE



IMAGE
FILE



SCANNED
DOCUMENT



SEARCHABLE
TEXT DOCUMENT



Optical Character Recognition

Introduction

Optical Character Recognition (OCR) is a widely adopted application for converting printed or handwritten images to text, which becomes a critical preprocessing component in text analysis pipelines, such as document retrieval and summarization. OCR has been significantly improved in recent years thanks to the wide adoption of the deep neural network (DNN), and thus deployed in many critical applications where OCR's quality is vital. For example, photo-based ID recognition depends on OCR's quality to automatically structure information into databases, and automatic trading sometimes relies on OCR to read certain news articles for determining the sentiment of news.

Unfortunately, OCR also inherits all counter-intuitive security problems of the DNNs. Especially, the OCR model is also vulnerable to *adversarial examples*, which are crafted by making human-imperceptible perturbations on original images with the intent of misleading the model. The wide adoption of OCR in real pipelines gives more incentives for adversaries to game the OCR, such as causing fake ID information, incorrect readings of metrics or instructions, etc. Figure 2 and 3 in the evaluation section illustrate two real-world examples with attacking the ID number and financial

Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



Introduction

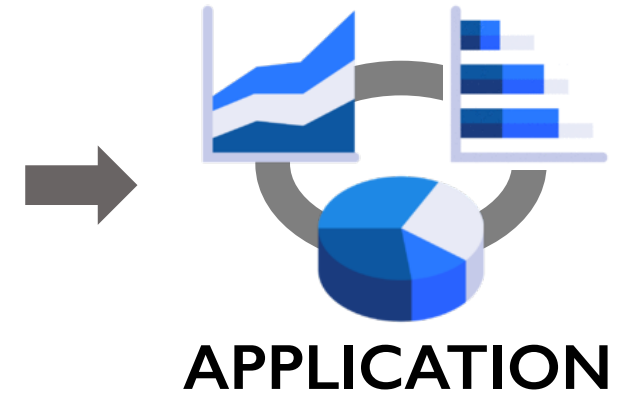
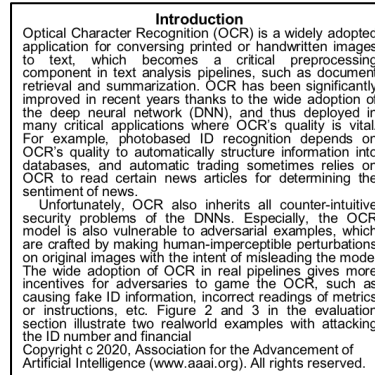
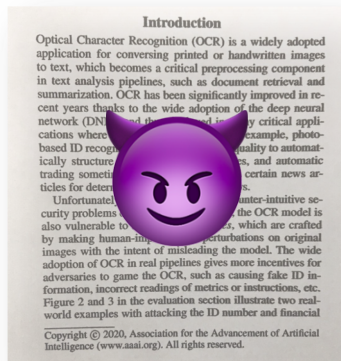
Optical Character Recognition (OCR) is a widely adopted application for converting printed or handwritten images to text, which becomes a critical preprocessing component in text analysis pipelines, such as document retrieval and summarization. OCR has been significantly improved in recent years thanks to the wide adoption of the deep neural network (DNN), and thus deployed in many critical applications where OCR's quality is vital. For example, photo-based ID recognition depends on OCR's quality to automatically structure information into databases, and automatic trading sometimes relies on OCR to read certain news articles for determining the sentiment of news.

Unfortunately, OCR also inherits all counter-intuitive security problems of the DNNs. Especially, the OCR model is also vulnerable to *adversarial examples*, which are crafted by making human-imperceptible perturbations on original images with the intent of misleading the model. The wide adoption of OCR in real pipelines gives more incentives for adversaries to game the OCR, such as causing fake ID information, incorrect readings of metrics or instructions, etc. Figure 2 and 3 in the evaluation section illustrate two realworld examples with attacking the ID number and financial

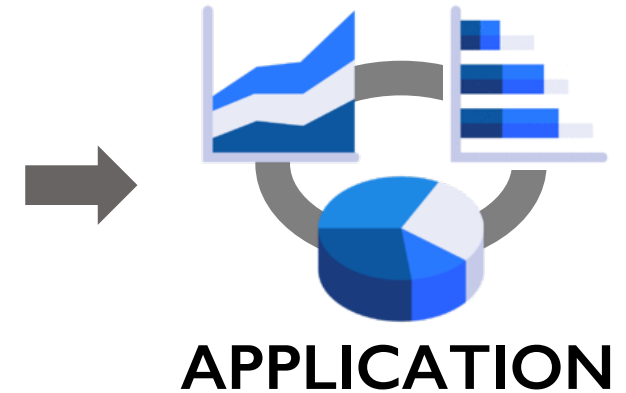
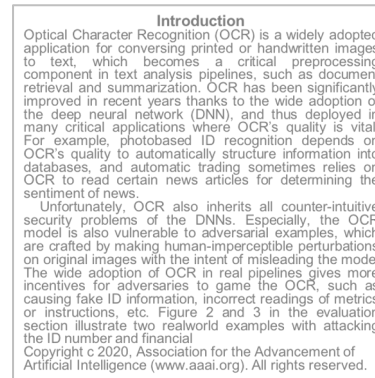
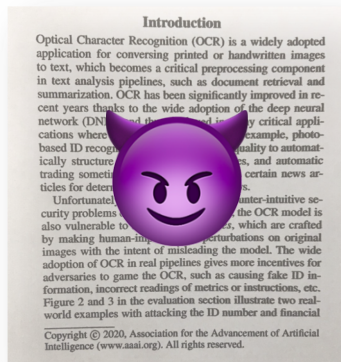
Copyright c 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.



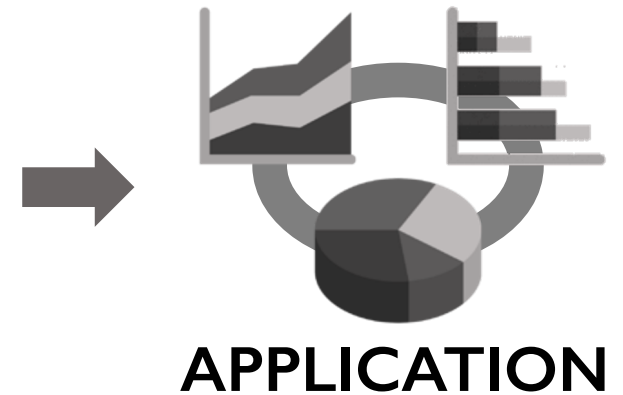
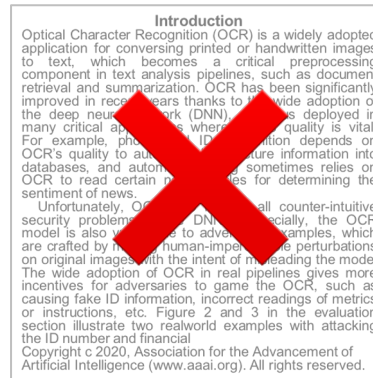
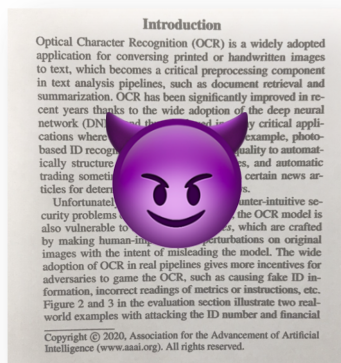
Optical Character Recognition



Optical Character Recognition



Optical Character Recognition



Traditional Attacks

[Goodfellow, 2015]



Original image

noise

Adversarial example

“Panda” 57.7%

“nematode” 8.2%

“gibbon” 99.3%

Visually imperceptible

[Brown, 2017]



Original image

Patch

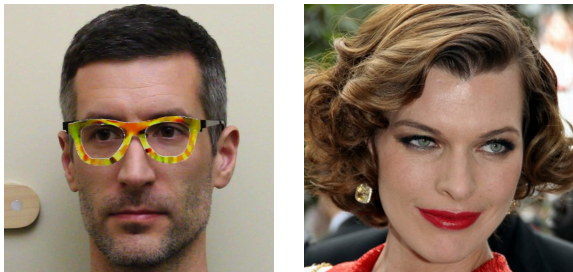
Adversarial example

“Banana” 97%

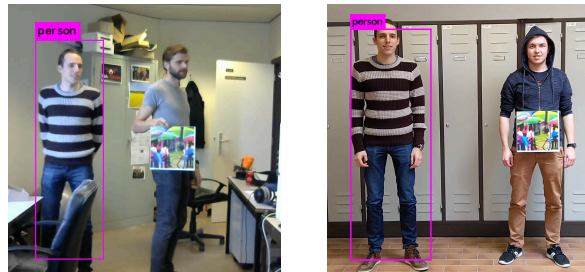
“toaster” 99%

Adversarial patch

[Sharif, 2016]



[Thys, 2019]



[Eykholt, 2018]

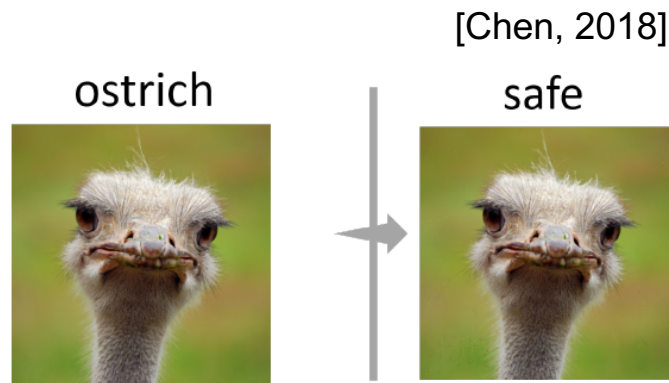


Adversarial attacks in the real world

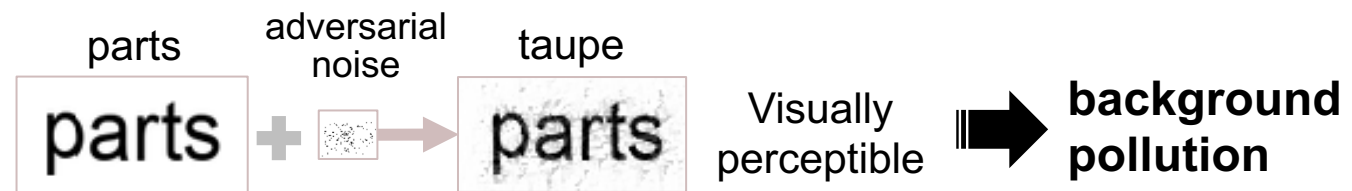
Traditional Attack vs. OCR Attack:

1. Image Backgrounds

- **colorful** background vs. **white** background



Colorful Background
(Photos)

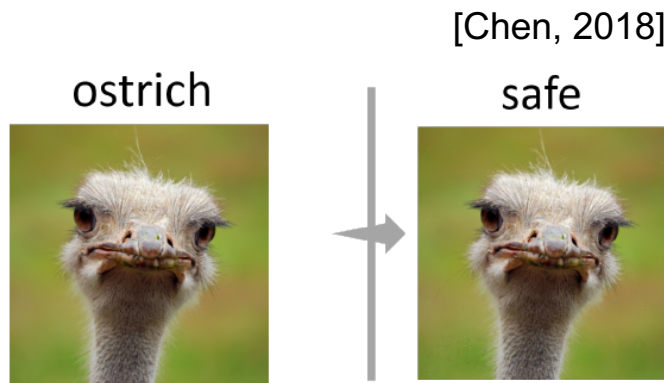


White Background
(Printed-text Images)

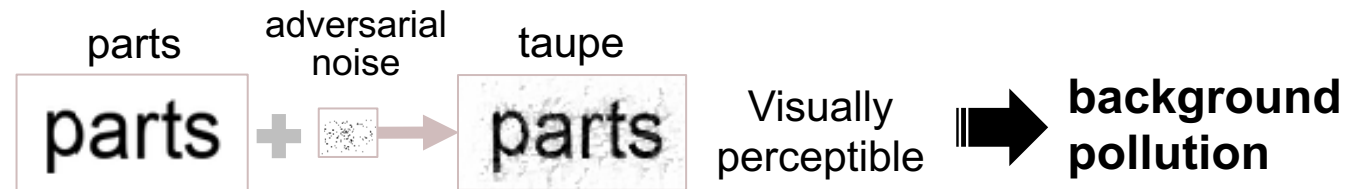
Traditional Attack vs. OCR Attack:

1. Image Backgrounds

■ **colorful background** vs. **white background**



Colorful Background
(Photos)

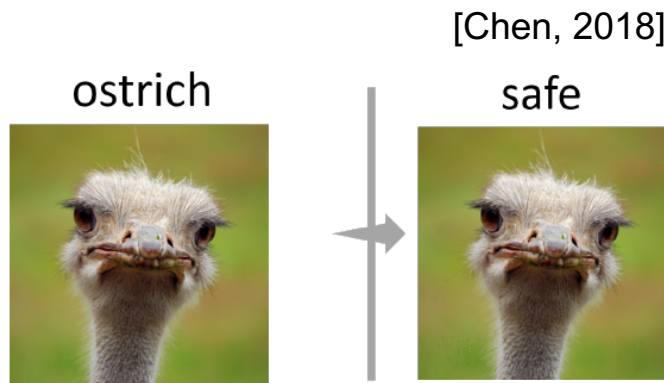


White Background
(Printed-text Images)

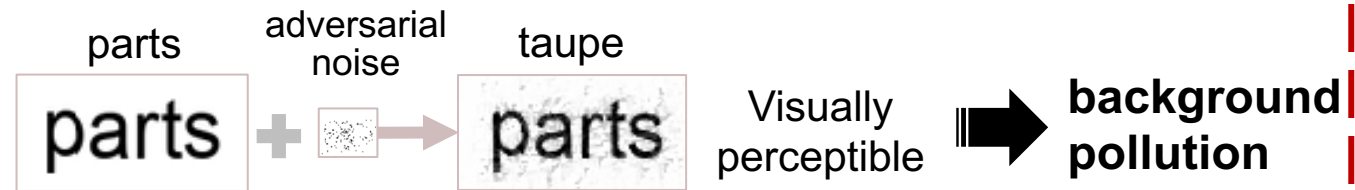
Traditional Attack vs. OCR Attack:

1. Image Backgrounds

- **colorful background** vs. **white background**



Colorful Background
(Photos)



White Background
(Printed-text Images)

Traditional Attack vs. OCR Attack:

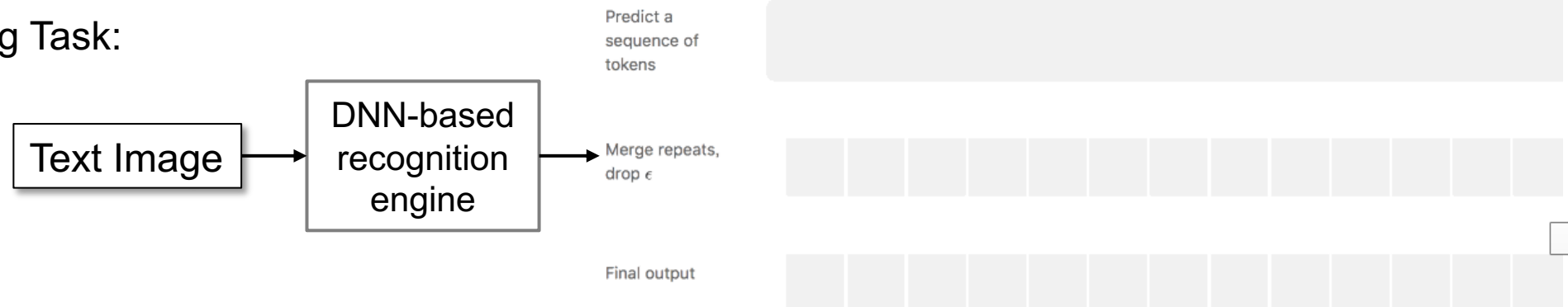
2. Model Tasks

- Traditional models ---- **image classification task** ➡ cross entropy loss
- DNN-based OCR models ---- **sequential labeling task** ➡ CTC loss

Image Classification Task:



Sequential Labeling Task:



FAWA: Fast Adversarial Watermark Attack

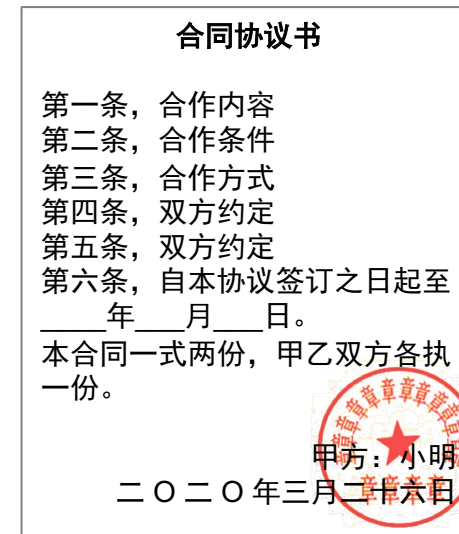
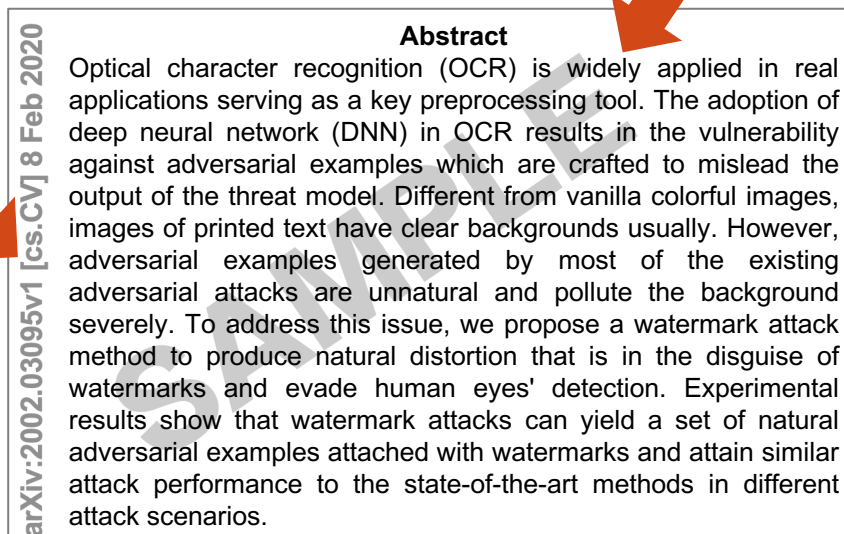


background pollution
sequential labeling task



Basic idea:

Making use of the popularity of watermark (**WM**) in the documents, we hide noise in watermarks.



FAWA: Fast Adversarial Watermark Attack

💡 Basic idea:

Making use of the popularity of watermark (**WM**) in the documents, we hide noise in watermarks.



FAWA

1. Natural

targeted text: *taupe*

traditional attack



FAWA



2. Fast **100%** attack success rate

78% fewer iterations

3. Low Perturbation Level

60% less noise

FAWA

1. Colored Watermark

targeted text: randem



2. OCR Model of Other Language

9月17日止当周，美国原油每月进口量增加31.3万桶。小方一点下场，左侧配置资产。

Target Output :

(Chinese) 9 月 **12** 日**上**当周，美国原油每**且**进口量增加 31.3万桶。**孙**方**二**点**丕**下场，**右**侧配置资产。

3. OCR accuracy-enhancing mechanism



FAWA



FAWA:

Attack Settings

White-box Model

- Attackers have perfect knowledge of the DNN architecture and parameters.

Targeted Attack

- Attackers aim to generate specific recognition texts.

Attack Pipeline



harmless image



adversarial image

intermediate status



clean image x

$f(x)$: random

random



adversarial image x'

t : random

random

Attack Pipeline:

1. Find Positions

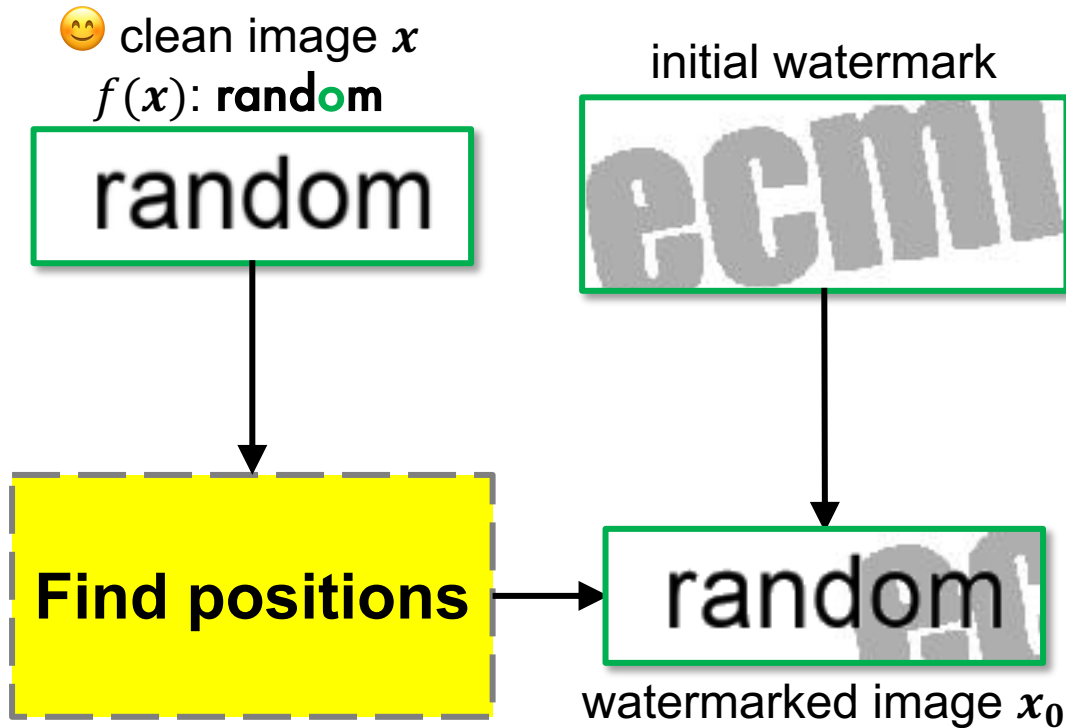


harmless image



adversarial image

intermediate status



adversarial image x'
 t : random



Attack Pipeline:

2. Watermark Attack

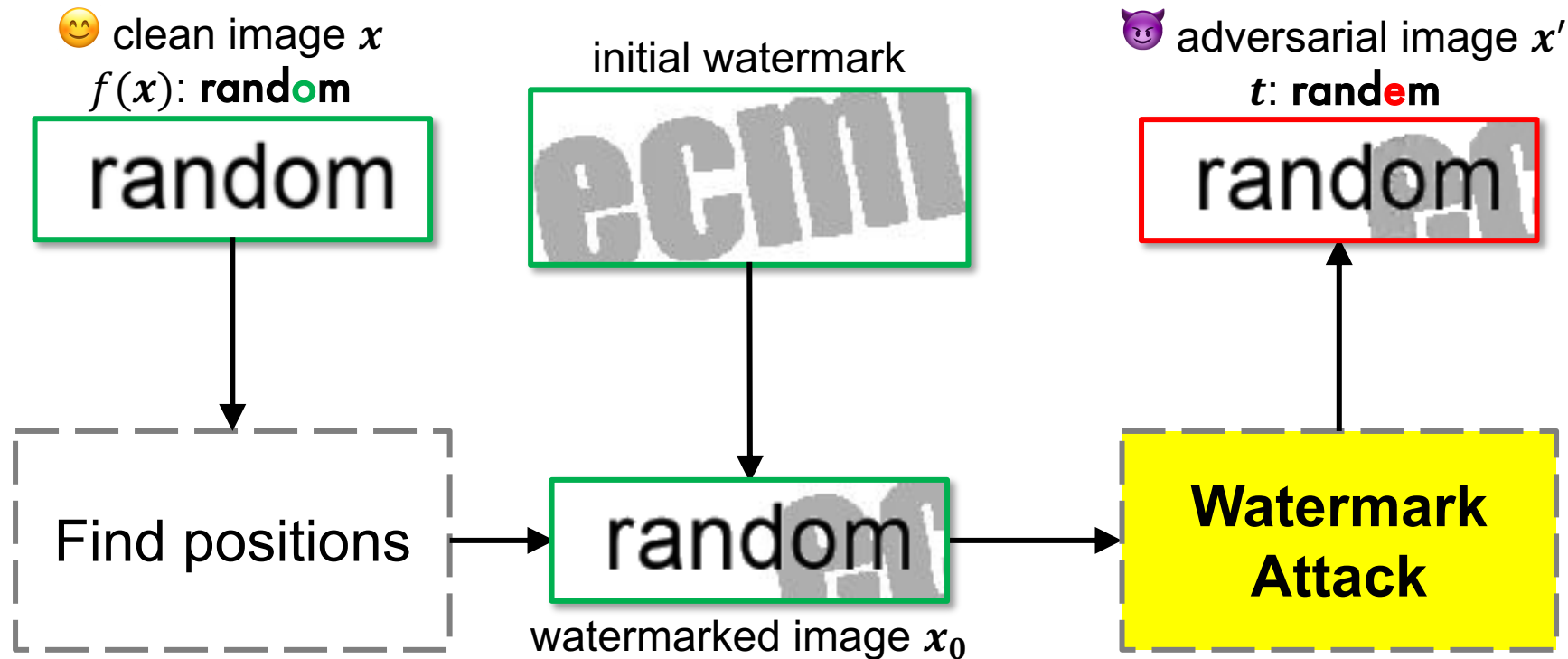


harmless image



adversarial image

intermediate status



Attack Pipeline: 3*. Full-Color Conversion

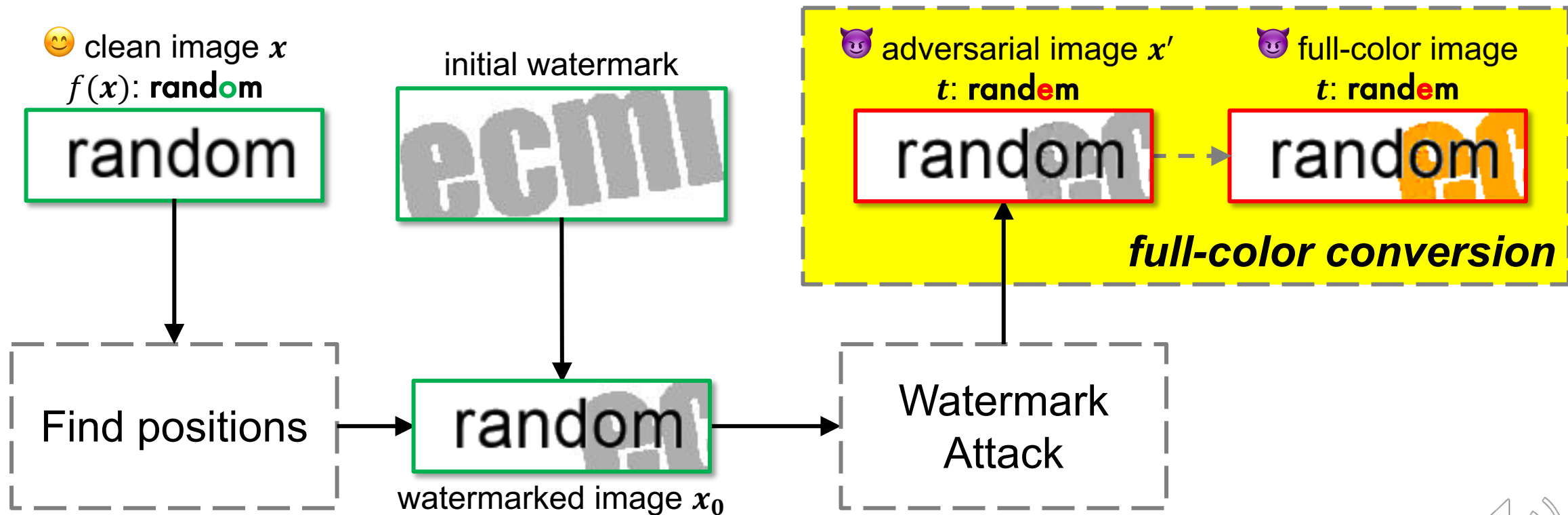


harmless image

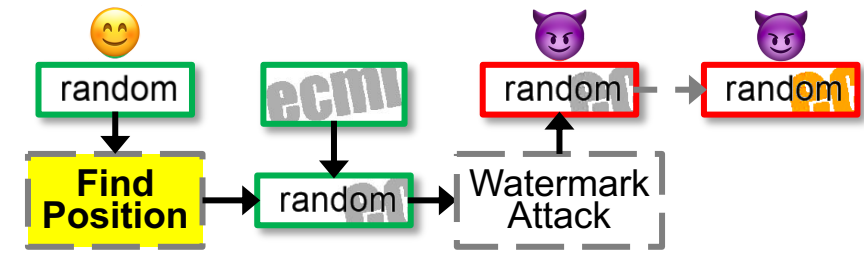


adversarial image

intermediate status



Critical Preprocess: Find Positions



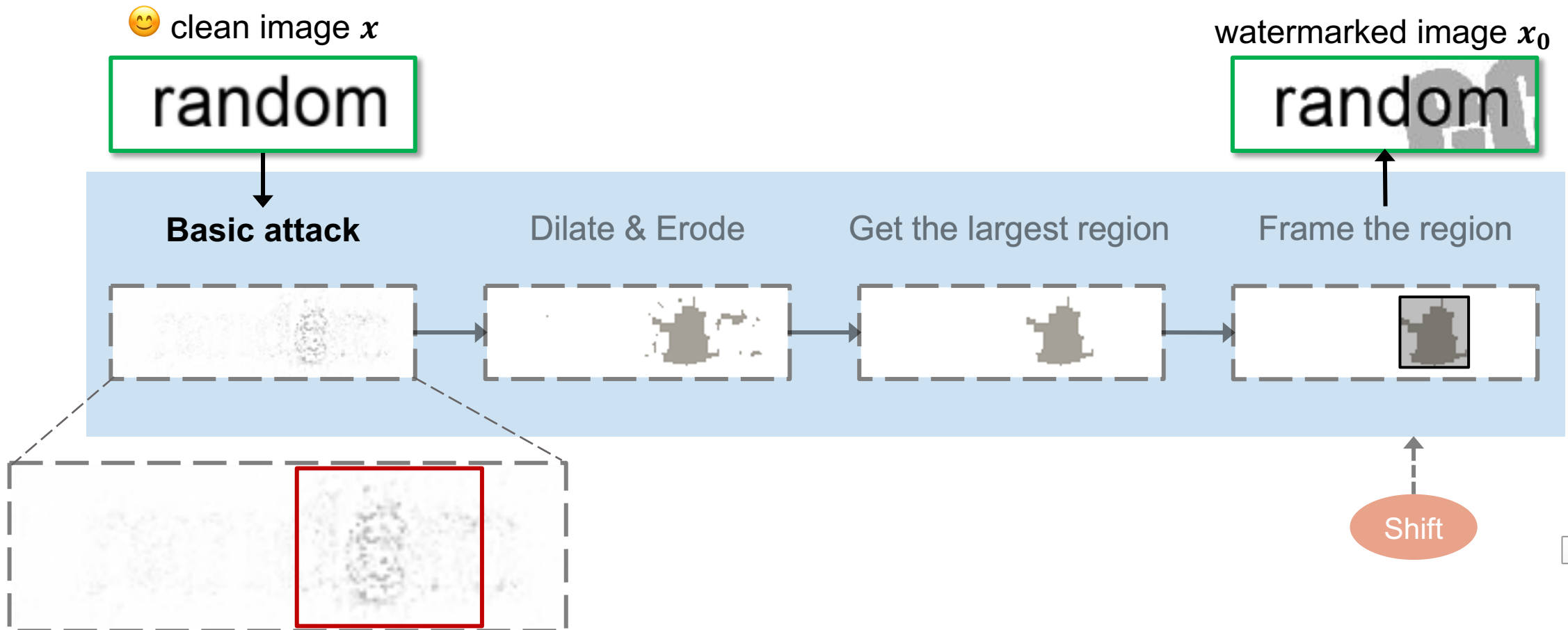
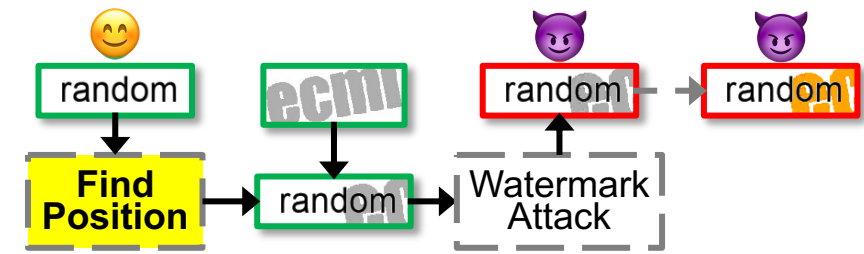
- **Basic Attack**

- We find positions based on the noise of the basic attack.
- We use *Momentum Iterative Method (MIM)* as the basic attack to find positions.



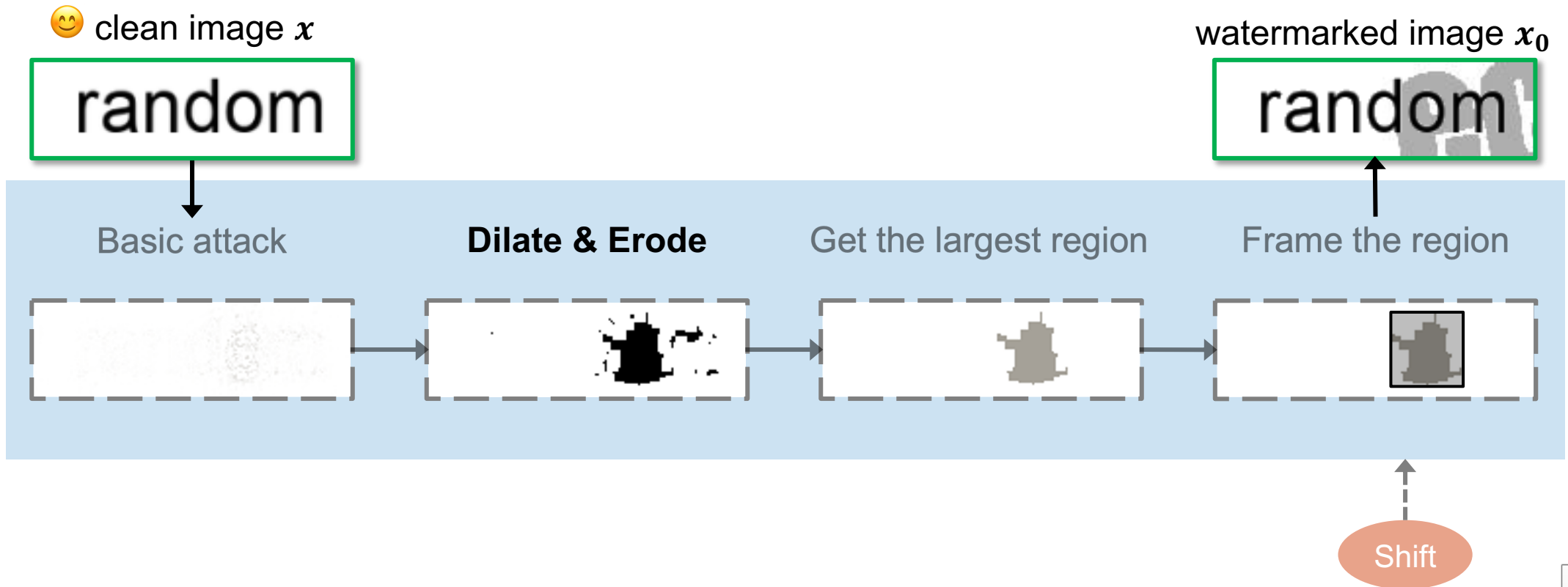
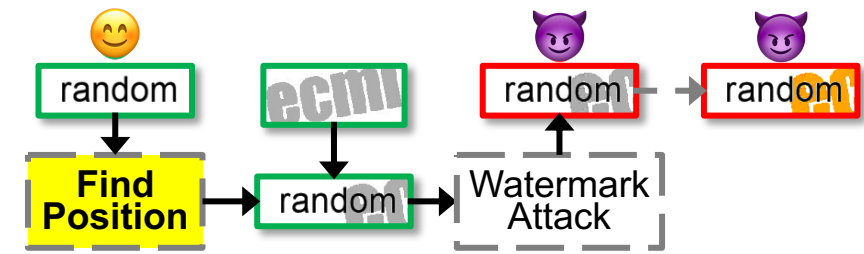
Find Positions:

1. Basic Attack



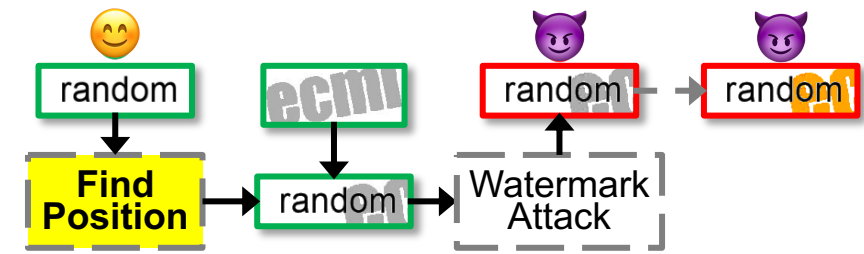
Find Positions:

2. Dilate & Erode



Find Positions:

3. Get the largest region



😊 clean image x

random

Basic attack



Dilate & Erode



Get the largest region



Frame the region



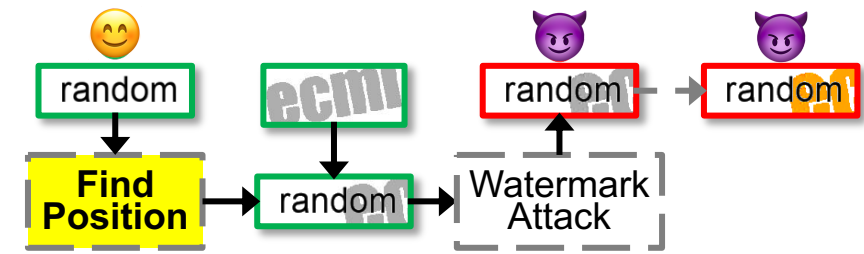
watermarked image x_0

random

Shift

Find Positions:

4. Frame the region



😊 clean image x

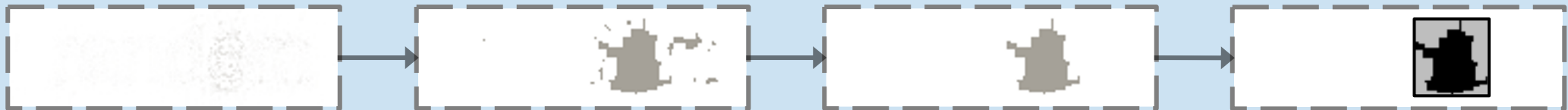
random

Basic attack

Dilate & Erode

Get the largest region

Frame the region

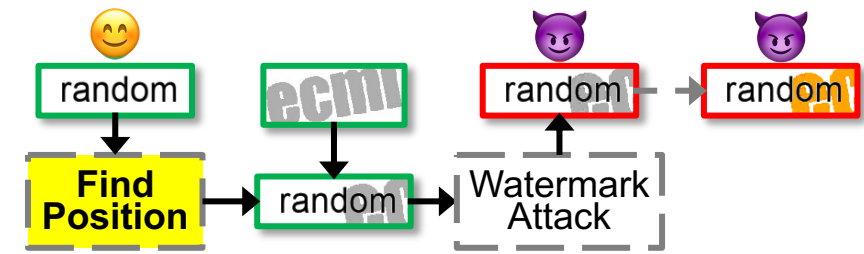


watermarked image x_0

random

Shift

Find Positions: 4*. Add a Shift



😊 clean image x

random

Basic attack

Dilate & Erode

Get the largest region

Frame the region



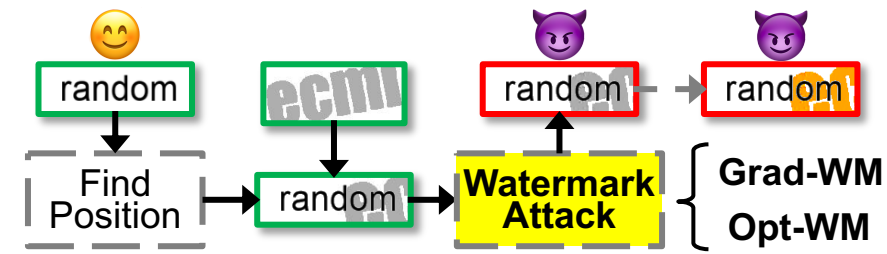
Shift

watermarked image x_0

random

Watermark Attack:

WM Attack = Basic Attack + WM

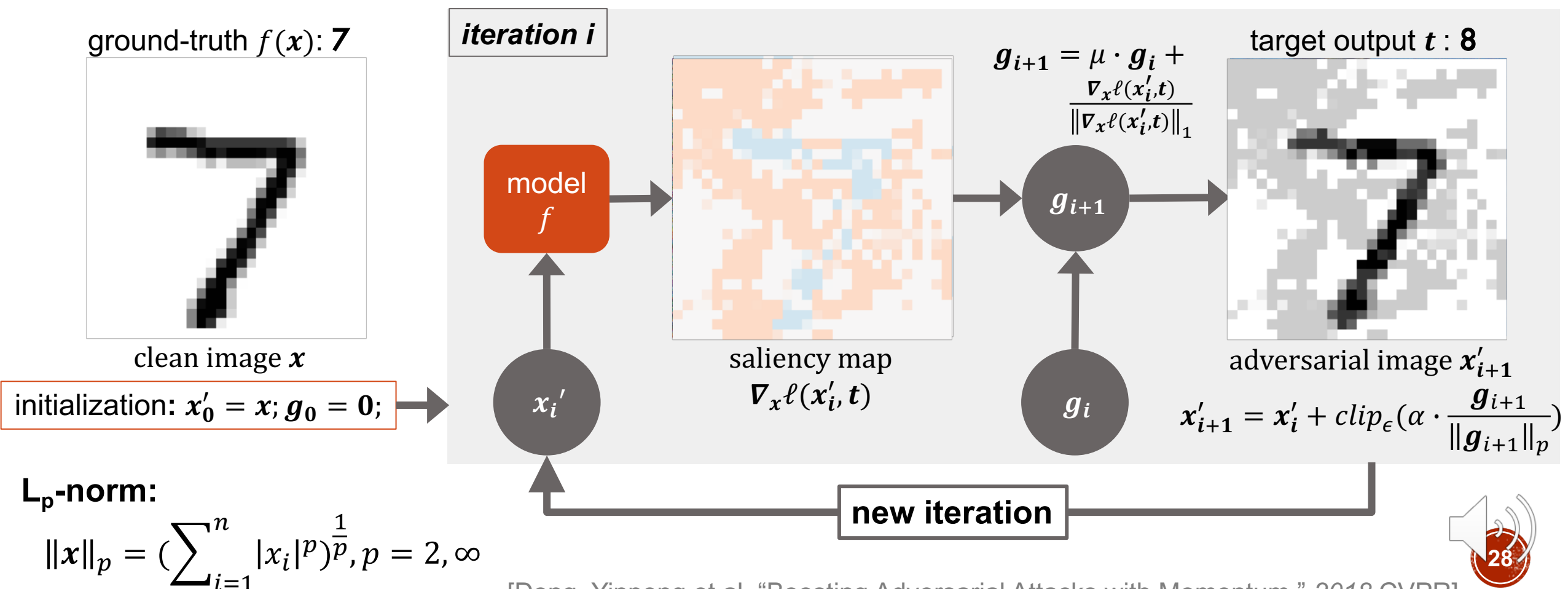


- **Basic Attacks:** (traditional attacks)
 - **Grad-Basic:** Gradient-based Basic Attack
 - Momentum Iterative Method (MIM) [Dong et al. 2018]
 - **Opt-Basic:** Optimization-based Basic Attack
 - OCR Attack [Song et al. 2018]
- **WM:** Watermark region allowed to add noise
- **Watermark Attacks:** (our attacks)
 - **Grad-WM:** Gradient-based Watermark Attack
 - **Grad-WM = Grad-Basic + WM**
 - **Opt-WM:** Optimization-based Watermark Attack
 - **Opt-WM = Opt-Basic + WM**

Gradient-based Basic Attack: Grad-Basic

✗ background pollution
✗ sequential labeling task

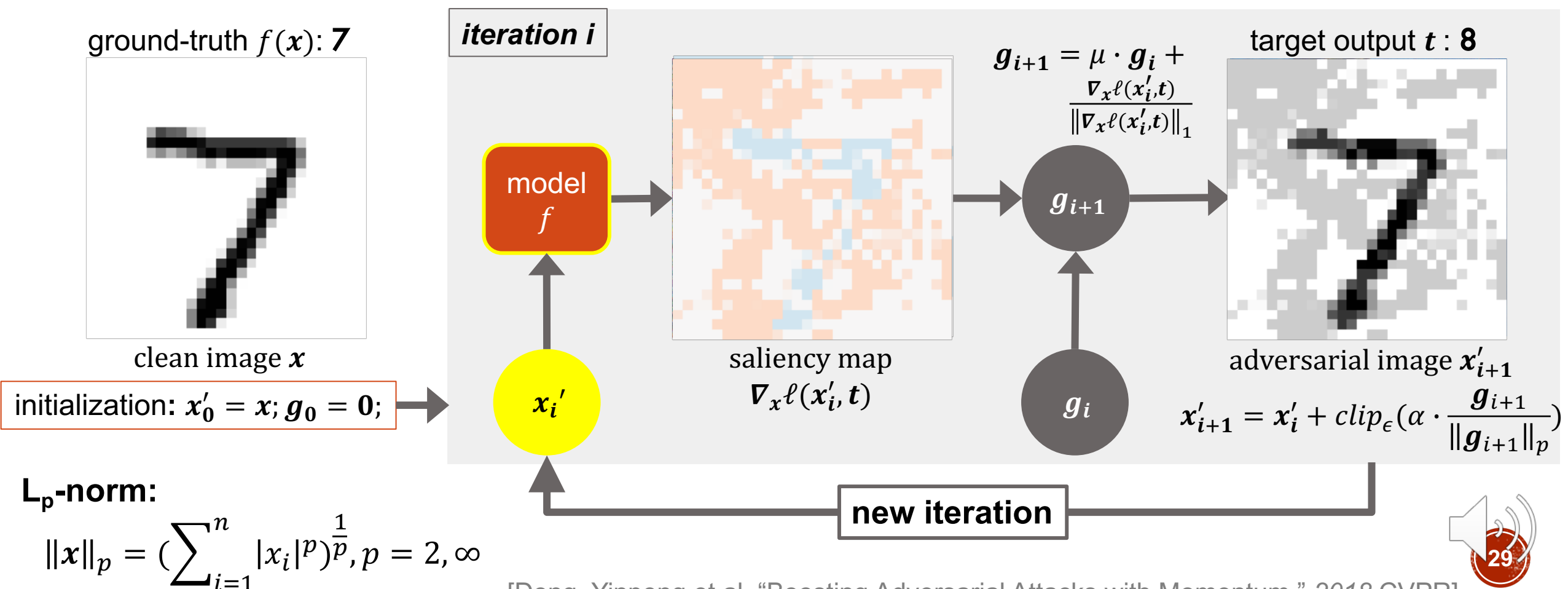
- cross entropy ℓ
- ϵ -bounded noise
- step size α
- decay factor μ



Gradient-based Basic Attack: Grad-Basic

✗ background pollution
✗ sequential labeling task

- cross entropy ℓ
- ϵ -bounded noise
- step size α
- decay factor μ

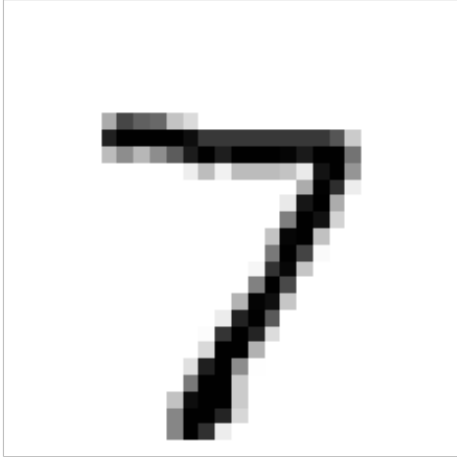


Gradient-based Basic Attack: Grad-Basic

background pollution
sequential labeling task

- cross entropy ℓ
- ϵ -bounded noise
- step size α
- decay factor μ

ground-truth $f(x): 7$



clean image x

iteration i

model
 f



saliency map

$$\nabla_x \ell(x'_i, t)$$

$$g_{i+1} = \mu \cdot g_i + \frac{\nabla_x \ell(x'_i, t)}{\|\nabla_x \ell(x'_i, t)\|_1}$$

g_{i+1}

target output $t: 8$



adversarial image x'_{i+1}

$$x'_{i+1} = x'_i + \text{clip}_\epsilon(\alpha \cdot \frac{g_{i+1}}{\|g_{i+1}\|_p})$$

initialization: $x'_0 = x; g_0 = 0;$

x'_i

g_i

new iteration

L_p -norm:

$$\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}, p = 2, \infty$$

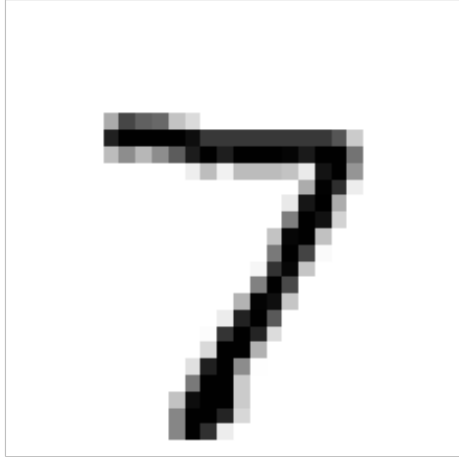


Gradient-based Basic Attack: Grad-Basic

background pollution
sequential labeling task

- cross entropy ℓ
- ϵ -bounded noise
- step size α
- decay factor μ

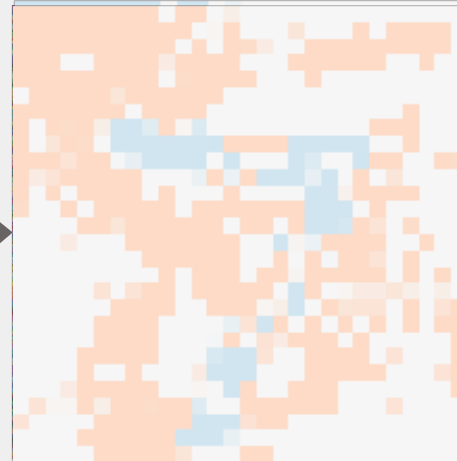
ground-truth $f(x): 7$



clean image x

iteration i

model
 f



saliency map
 $\nabla_x \ell(x'_i, t)$

$$g_{i+1} = \mu \cdot g_i + \frac{\nabla_x \ell(x'_i, t)}{\|\nabla_x \ell(x'_i, t)\|_1}$$

g_{i+1}

g_i

target output $t: 8$



adversarial image x'_{i+1}

$$x'_{i+1} = x'_i + \text{clip}_\epsilon(\alpha \cdot \frac{g_{i+1}}{\|g_{i+1}\|_p})$$

new iteration

initialization: $x'_0 = x; g_0 = 0;$

L_p -norm:

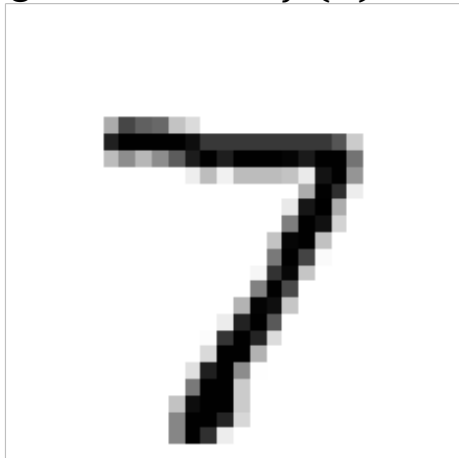
$$\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}, p = 2, \infty$$

Gradient-based Basic Attack: Grad-Basic

background pollution
sequential labeling task

- cross entropy ℓ
- ϵ -bounded noise
- step size α
- decay factor μ

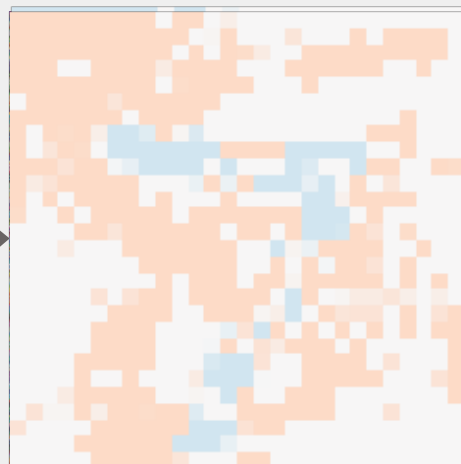
ground-truth $f(x): 7$



clean image x

iteration i

model
 f



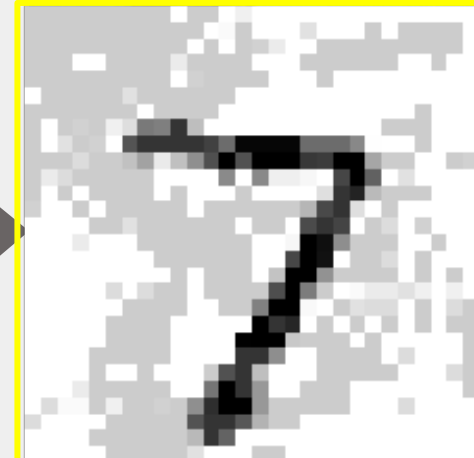
saliency map
 $\nabla_x \ell(x'_i, t)$

$$g_{i+1} = \mu \cdot g_i + \frac{\nabla_x \ell(x'_i, t)}{\|\nabla_x \ell(x'_i, t)\|_1}$$

g_{i+1}

g_i

target output $t: 8$



adversarial image x'_{i+1}

$$x'_{i+1} = x'_i + \text{clip}_\epsilon(\alpha \cdot \frac{g_{i+1}}{\|g_{i+1}\|_p})$$

initialization: $x'_0 = x; g_0 = 0;$

x'_i

new iteration

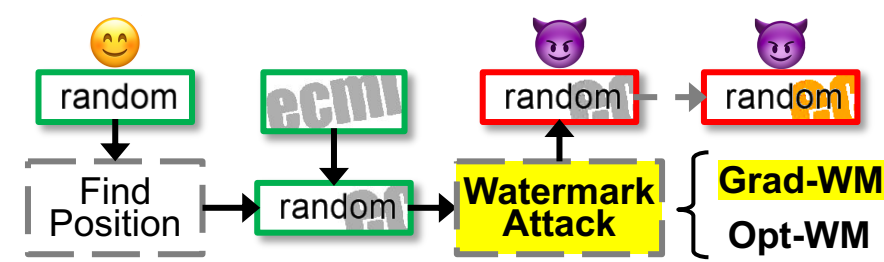
L_p -norm:

$$\|x\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}, p = 2, \infty$$



Watermark Attack:

Grad-WM=Grad-Basic+WM



sequential labeling task
background pollution

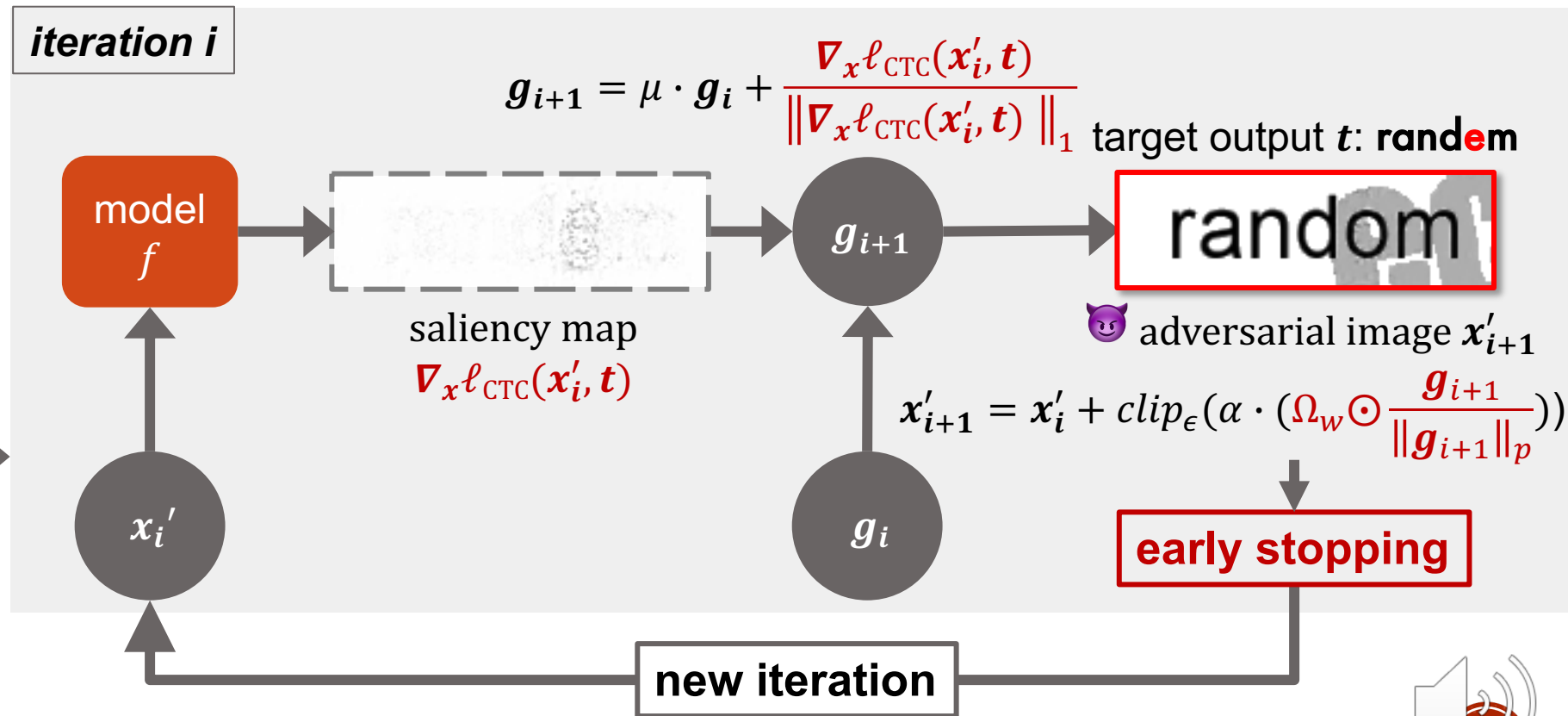
- grayscale β
- watermark mask Ω_w
- text mask Ω_t

ground-truth $f(x)$: random

random

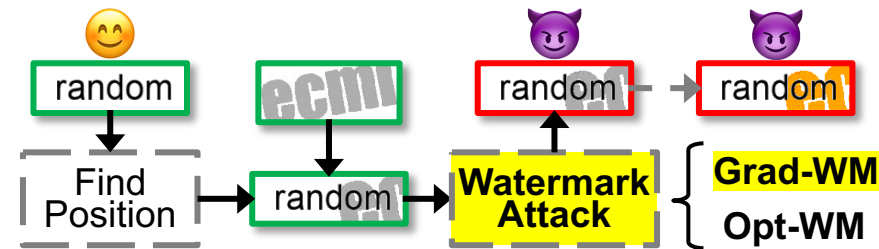
😊 clean image x

initialization: $g_0 = 0$;
 $x'_0 = (1 - \Omega_w \odot \overline{\Omega_t}) \odot x$
 $+ \beta \cdot \Omega_w \odot \overline{\Omega_t}$;



Watermark Attack:

Grad-WM=Grad-Basic+WM



- ✓ sequential labeling task
- ✓ background pollution

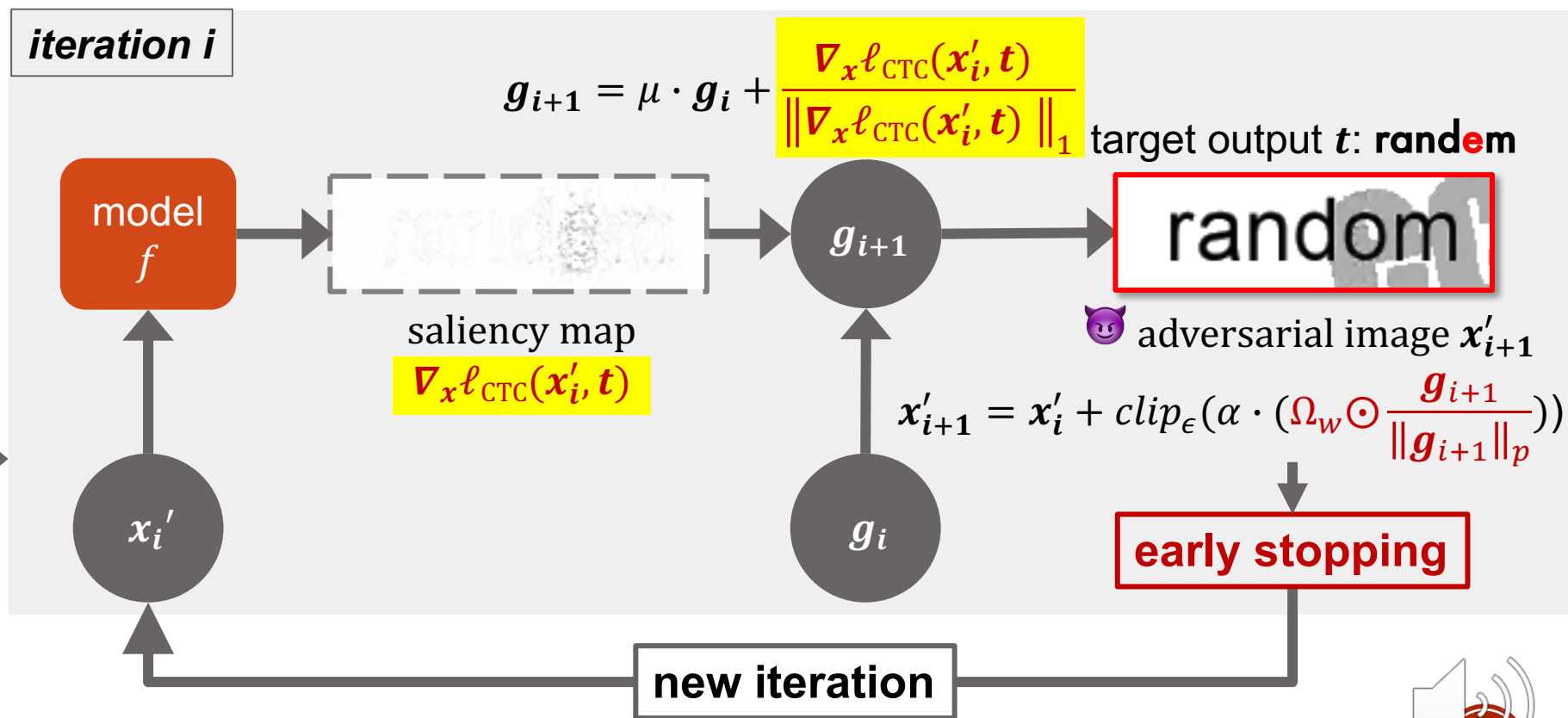
- grayscale β
- watermark mask Ω_w
- text mask Ω_t

ground-truth $f(x)$: random

random

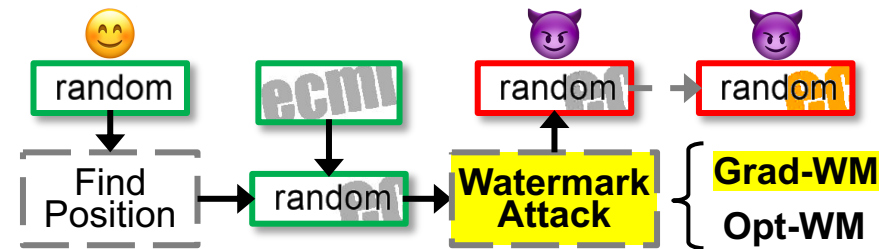
😊 clean image x

initialization: $g_0 = 0$;
 $x'_0 = (1 - \Omega_w \odot \overline{\Omega_t}) \odot x$
 $+ \beta \cdot \Omega_w \odot \overline{\Omega_t}$



Watermark Attack:

Grad-WM=Grad-Basic+WM



sequential labeling task
background pollution

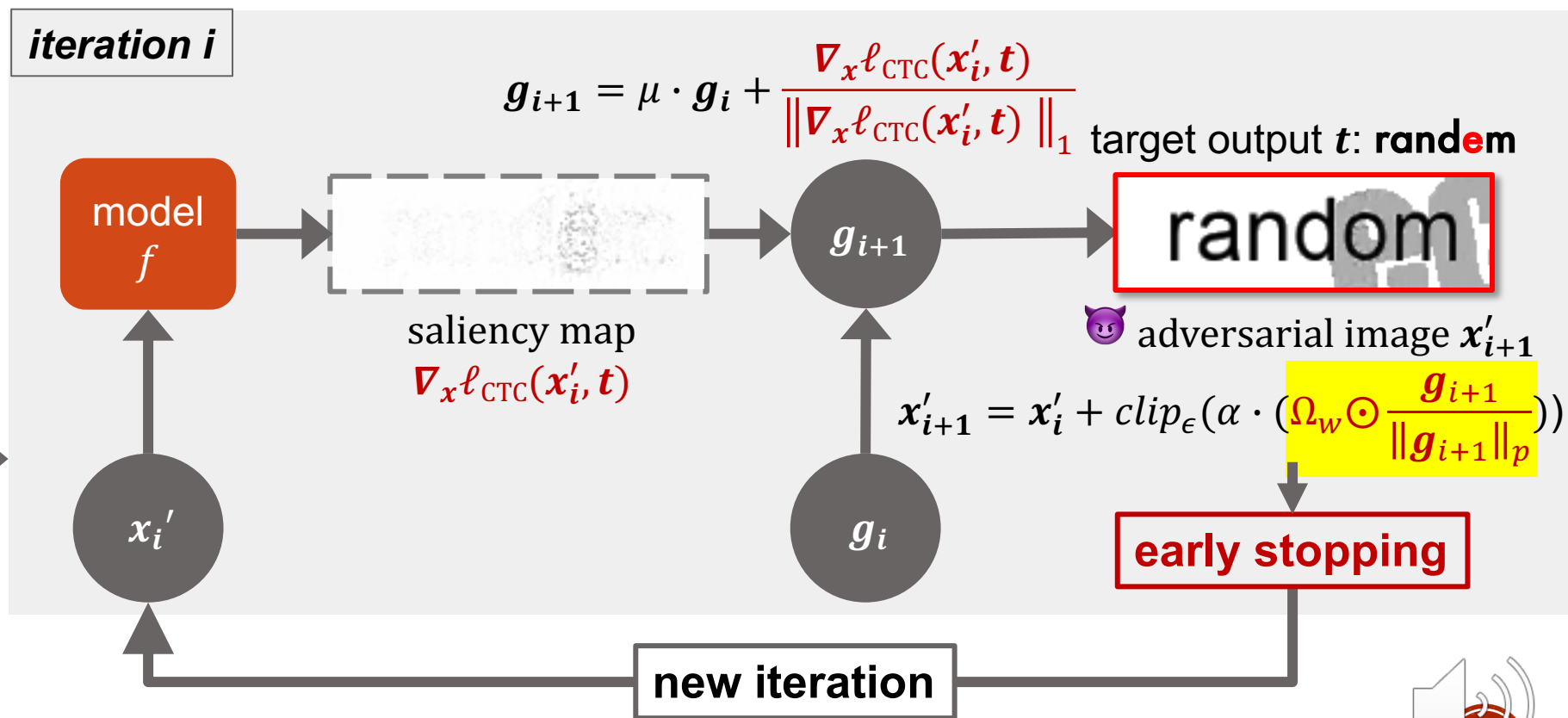
- grayscale β
- watermark mask Ω_w
- text mask Ω_t

ground-truth $f(x)$: random

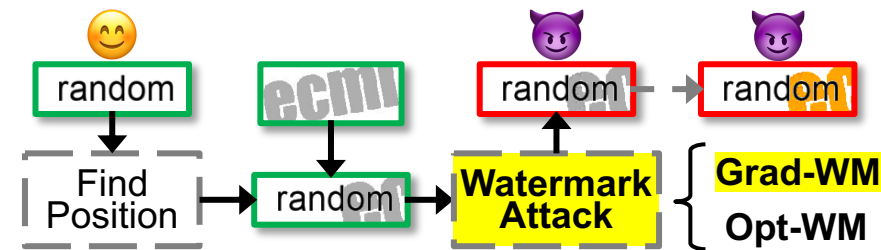
random

😊 clean image x

initialization: $g_0 = 0$;
 $x'_0 = (1 - \Omega_w \odot \overline{\Omega_t}) \odot x$
 $+ \beta \cdot \Omega_w \odot \overline{\Omega_t}$



Watermark Attack: Grad-WM=Grad-Basic+WM



sequential labeling task
background pollution

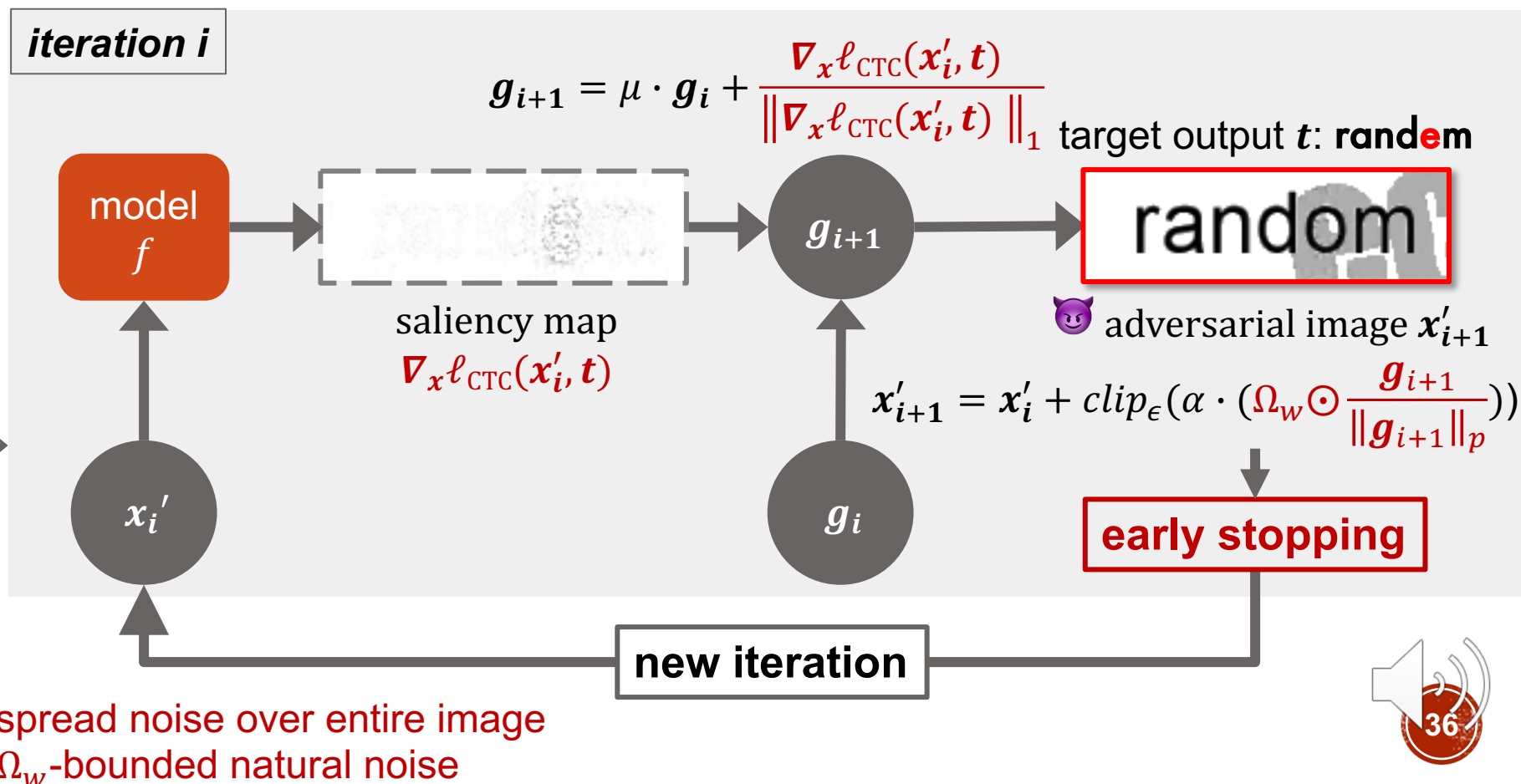
- grayscale β
- watermark mask Ω_w
- text mask Ω_t

ground-truth $f(x)$: random

random

😊 clean image x

initialization: $g_0 = 0$;
 $x'_0 = (1 - \Omega_w \odot \overline{\Omega_t}) \odot x$
 $+ \beta \cdot \Omega_w \odot \overline{\Omega_t}$



amazon

amazon

tennis

vs.

tennis

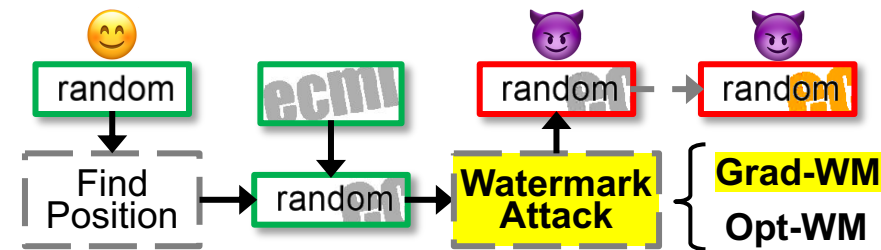
Grad-Basic

Grad-WM



Watermark Attack:

Grad-WM=Grad-Basic+WM



sequential labeling task
background pollution

- grayscale β
- watermark mask Ω_w
- text mask Ω_t

ground-truth $f(x)$: random

random

😊 clean image x

initialization: $g_0 = 0$;

$$x'_0 = (1 - \Omega_w \odot \overline{\Omega_t}) \odot x + \beta \cdot \Omega_w \odot \overline{\Omega_t};$$

readability ↗

model f

saliency map
 $\nabla_x \ell_{CTC}(x'_i, t)$

$$g_{i+1} = \mu \cdot g_i + \frac{\nabla_x \ell_{CTC}(x'_i, t)}{\|\nabla_x \ell_{CTC}(x'_i, t)\|_1} \text{ target output } t: \text{random}$$

g_{i+1}

random

😈 adversarial image x'_{i+1}

$$x'_{i+1} = x'_i + \text{clip}_\epsilon(\alpha \cdot (\Omega_w \odot \frac{g_{i+1}}{\|g_{i+1}\|_p}))$$

g_i

early stopping

new iteration

ecm!

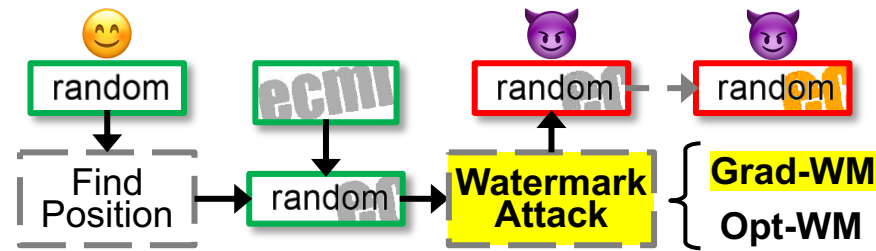
random

watermark mask Ω_w

text mask Ω_t

Watermark Attack:

Grad-WM=Grad-Basic+WM



sequential labeling task
background pollution

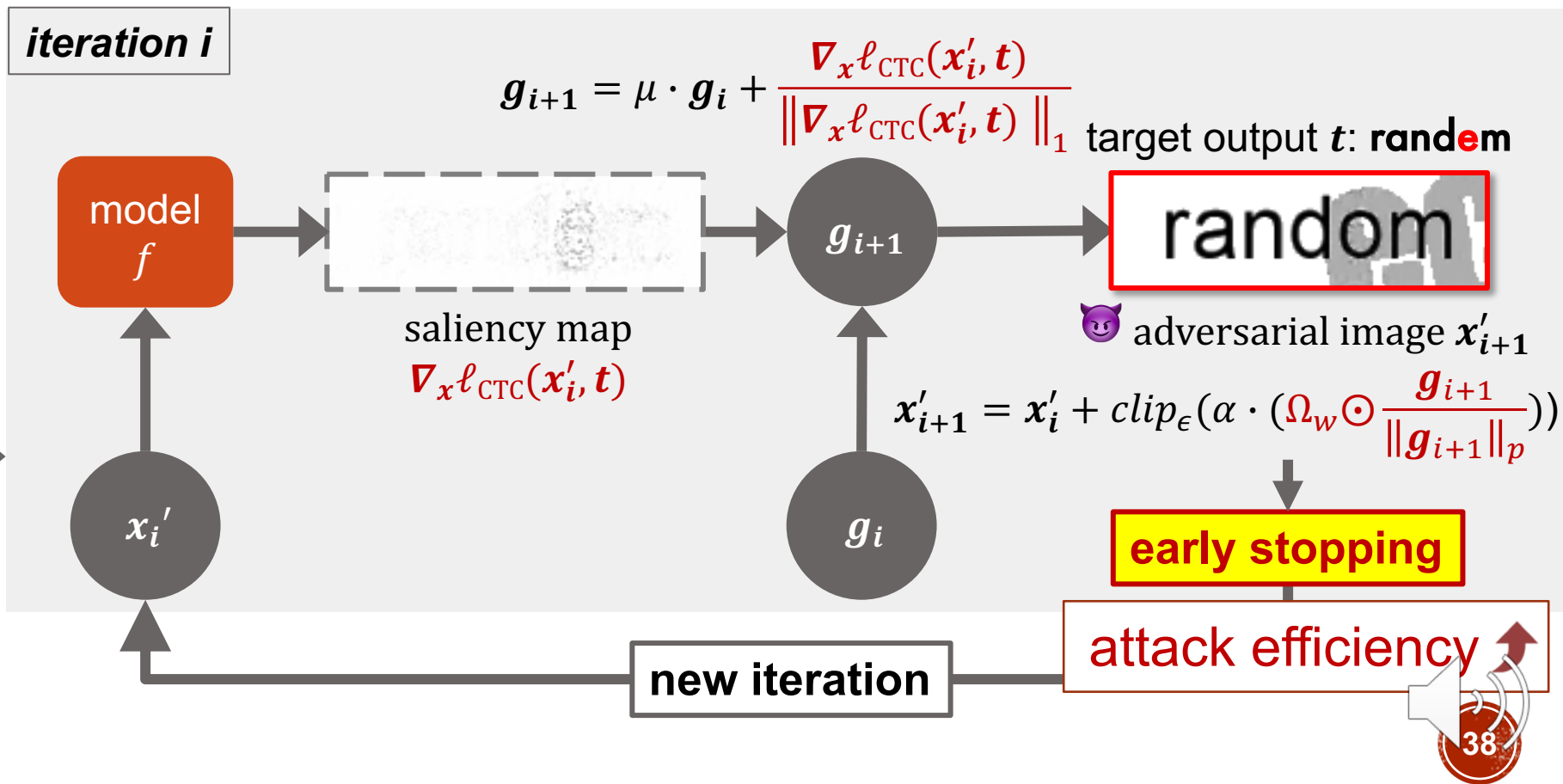
- grayscale β
- watermark mask Ω_w
- text mask Ω_t

ground-truth $f(x)$: random

random

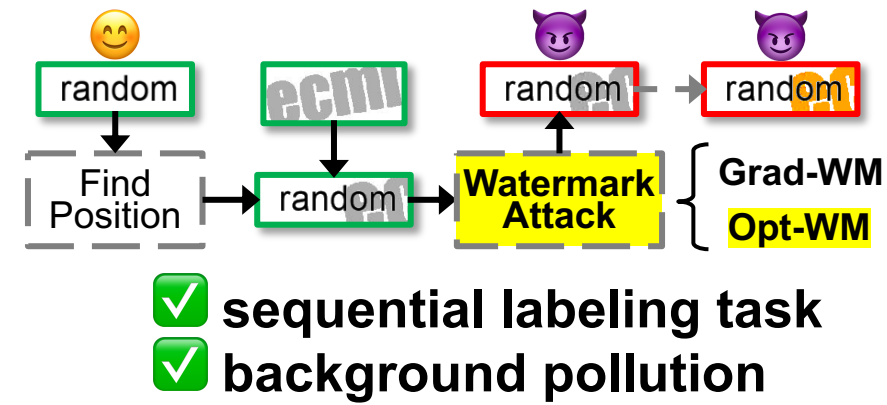
😊 clean image x

initialization: $g_0 = 0$;
 $x'_0 = (1 - \Omega_w \odot \overline{\Omega_t}) \odot x$
 $+ \beta \cdot \Omega_w \odot \overline{\Omega_t}$



Watermark Attack:

Opt-WM=Opt-Basic+WM



Opt-Basic:

$$\min_w c \cdot \ell_{CTC} \left(\frac{\tanh(\mathbf{w}) + 1}{2}, \mathbf{t} \right) + \left\| \frac{\tanh(\mathbf{w}) + 1}{2} - \mathbf{x} \right\|_2^2$$

Opt-WM:

- 1. Separate the *perturbation term* \mathbf{w} :

$$\min_w c \cdot \ell_{CTC} \left(\frac{\tanh(\mathbf{w} + \mathbf{x}) + 1}{2}, \mathbf{t} \right) + \left\| \frac{\tanh(\mathbf{w} + \mathbf{x}) + 1}{2} - \mathbf{x} \right\|_2^2$$

- 2. Introduce the *watermark mask* Ω_w :

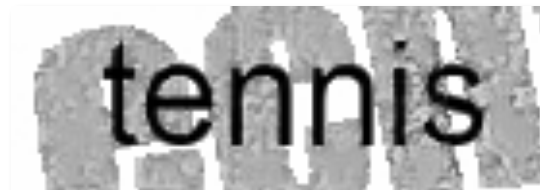
$$\min_w c \cdot \ell_{CTC} \left(\frac{\tanh(\Omega_w \odot \mathbf{w} + \mathbf{x}) + 1}{2}, \mathbf{t} \right) + \left\| \frac{\tanh(\Omega_w \odot \mathbf{w} + \mathbf{x}) + 1}{2} - \mathbf{x} \right\|_2^2$$

amazon



Opt-Basic (traditional)

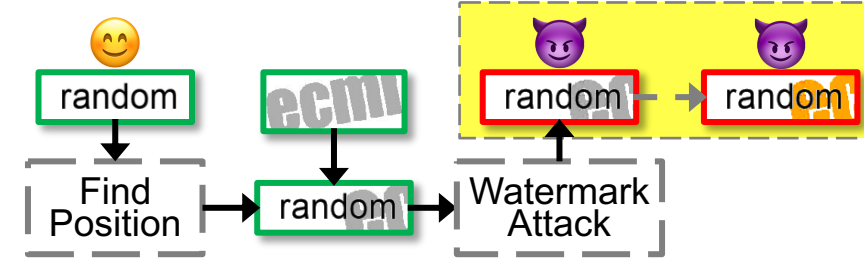
vs



Opt-WM (ours)

amazon

Improving Readability: Full-Color Conversion



- Grayscale watermark 🖱️ Full-color watermark
- Given grayscale value $Gray$, fix R value and B value, we can calculate the left G value by the *ITU-R 601-2 luma transform*:

$$Gray = R * 0.299 + G * 0.587 + B * 0.114$$

Experiment Settings:

Threat Model

- **Calamari OCR**
 - a open-source OCR model
 - 2 convolutional layers, 2 pooling layers, a LSTM layer
 - trained by CTC in Tensorflow



Experiment Settings:

Data Generation ---- IMDB

- Printed-text images (100% accuracy)
 - 5 fonts: *Courier*, *Georgia*, *Helvetica*, *Times*, *Arial* (font size:32 px)

parts parts parts parts parts

- 1092 word images
- 1479 sentence images
- 97 paragraph images

Experiment Settings:

Data Generation ---- IMDB

- Attack pairs

- Letter-Level Attacks (word images)

- Difficulty: **Easy** Case / **Random** Case / **Hard** Case (Replace)

parts pants parts pacts parts pasts

- Operation: **Replace** Case / **Insert** Case / **Delete** Case

parts pants parts partis parts pars (parts)

- Word-Level Attacks (word / sentence / paragraph images) (Replace)

parts taupe This one did exactly that. Tale one did exactly that.

Evaluation Metrics

▪ Perturbation Level

image quality

★ **MSE**: mean-square error

$$\text{MSE} = \frac{1}{|x|} (x - x')^2$$

▪ PSNR: peak-signal-to-noise ratio

$$\text{PSNR} = 10 \log\left(\frac{D^2}{\text{MSE}}\right)$$

▪ SSIM: structural similarity index

▪ Success Rate

▪ ASR: targeted attack success rate 100

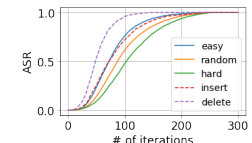
$$\text{ASR} = \frac{\#(f(x')=t)}{\#(x)}$$

▪ Attack Efficiency ⚡

★ **I_{avg}**: average iterations of successful attacks

$$I_{\text{avg}} = \frac{\sum_i^n \#(\text{iteration}_i)}{n}$$

★ **ASR-Iter**: plot the change of ASR along with increasing iterations



Basic Attack vs. FAWA: Grad-Basic vs. Grad-WM

- MSE: mean-square error
- I_{avg} : average iterations
of successful attacks
- ASR: targeted attack
success rate 100

		replacement						insertion		deletion	
		easy		random		hard					
		MSE	I_{avg}	MSE	I_{avg}	MSE	I_{avg}	MSE	I_{avg}	MSE	I_{avg}
Grad-Basic	Courier	10.5	59	14.0	74	17.0	70	11.6	50	3.2	21
	Georgia	27.4	43	32.8	99	37.3	104	22.1	83	17.3	55
	Helvetica	27.0	51	33.6	113	38.6	113	23.0	70	16.7	43
	Times	26.4	62	31.5	85	35.8	109	20.3	98	17.2	68
	Arial	29.8	51	36.7	73	42.5	66	24.3	88	19.2	59
example		parts		parts		parts		parts		parts	
Grad-WM	Courier	2.8	30	3.6	18	4.3	27	3.6	21	0.7	8
	Georgia	7.8	15	8.9	33	9.8	30	5.1	39	3.5	21
	Helvetica	8.4	9	10.0	52	11.2	52	6.3	23	3.7	19
	Times	7.3	15	8.3	20	9.3	34	4.5	7	3.4	21
	Arial	9.4	13	11.1	14	12.7	25	6.2	33	4.4	20
example		parts		parts		parts		parts		parts	
target output		pants		pacts		pasts		partis		pars	

■ **ASR: 100% in default**

letter-level attack in word images



FAWA: Lower Perturbation Level

- MSE: mean-square error
- I_{avg} : average iterations of successful attacks
- ASR: targeted attack success rate 100

		replacement						insertion		deletion	
		easy		random		hard					
		MSE	I _{avg}	MSE	I _{avg}	MSE	I _{avg}	MSE	I _{avg}	MSE	I _{avg}
Grad-Basic	Courier	10.5	59	14.0	74	17.0	70	11.6	50	3.2	21
	Georgia	27.4	43	32.8	99	37.3	104	22.1	83	17.3	55
	Helvetica	27.0	51	33.6	113	38.6	113	23.0	70	16.7	43
	Times	26.4	62	31.5	85	35.8	109	20.3	98	17.2	68
	Arial	29.8	51	36.7	73	42.5	66	24.3	88	19.2	59
example		parts		parts		parts		parts		parts	
Grad-WM	Courier	2.8	↘ 30	3.6	↘ 18	4.3	↘ 27	3.6	↘ 21	0.7	↘ 8
	Georgia	7.8	↘ 15	8.9	↘ 33	9.8	↘ 30	5.1	↘ 39	3.5	↘ 21
	Helvetica	8.4	↘ 9	10.0	↘ 52	11.2	↘ 52	6.3	↘ 23	3.7	↘ 19
	Times	7.3	↘ 15	8.3	↘ 20	9.3	↘ 34	4.5	↘ 7	3.4	↘ 21
	Arial	9.4	↘ 13	11.1	↘ 14	12.7	↘ 25	6.2	↘ 33	4.4	↘ 20
example		parts		parts		parts		parts		parts	
target output		pants		pacts		pasts		partis		pars	

■ **74% less noise** (MSE)
on average

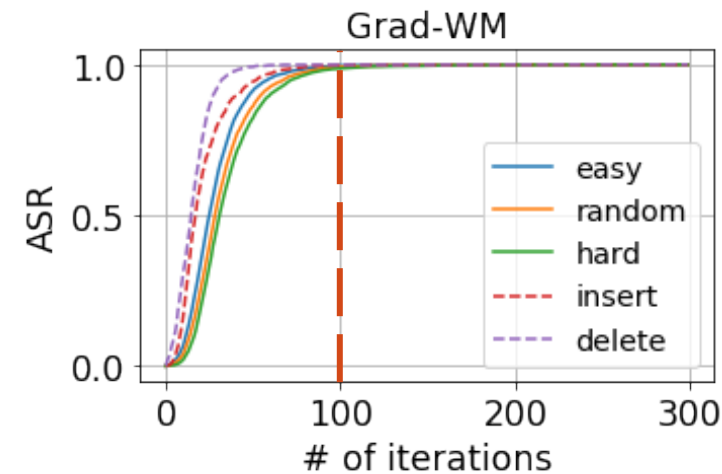
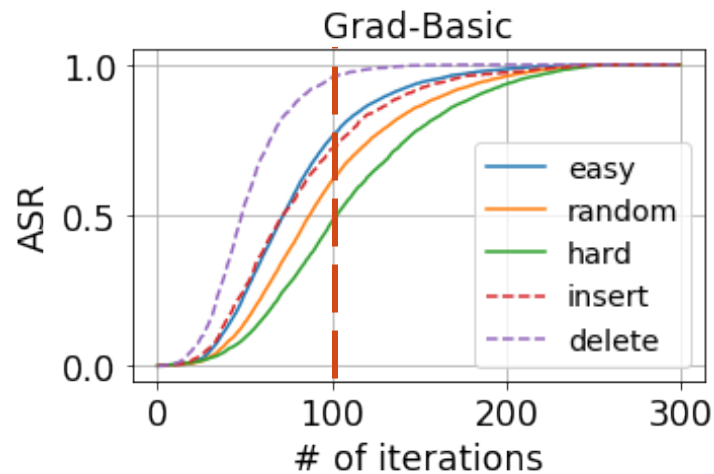
letter-level attack in word images



FAWA: Faster Attack Speed

- MSE: mean-square error
- I_{avg} : average iterations of successful attacks
- ASR: targeted attack success rate 100

- **67% fewer iterations** (I_{avg}) on average
- A sharper slope indicates faster attack speed in the figure.



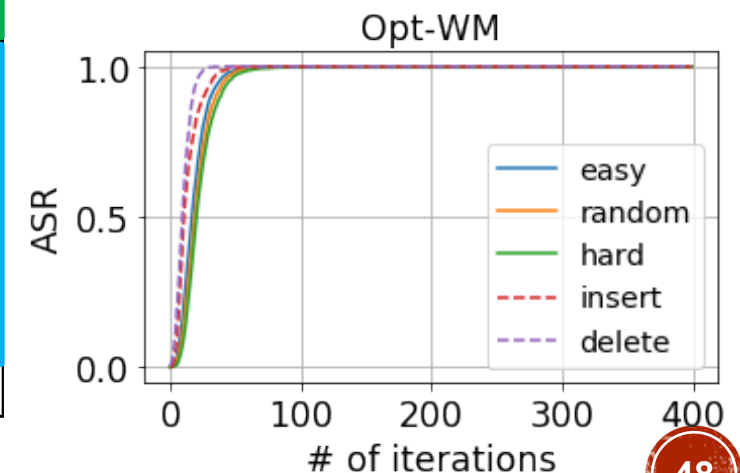
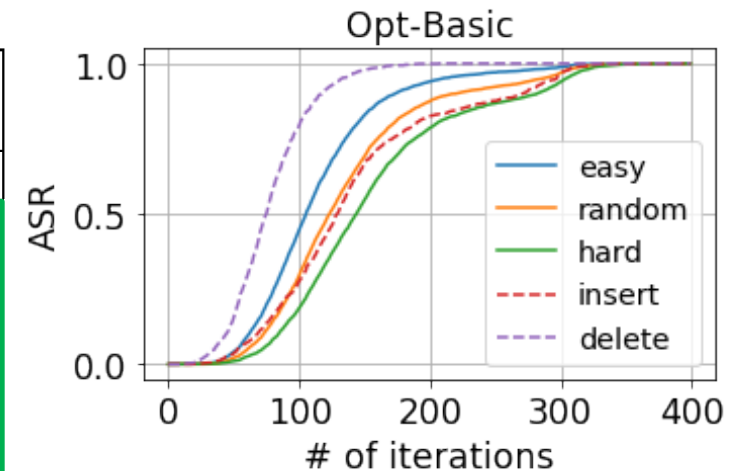
Basic Attack vs. FAWA: Opt-Basic vs. Opt-WM

- MSE: mean-square error
- I_{avg} : average iterations of successful attacks

- 44% less noise
- 88% fewer iterations

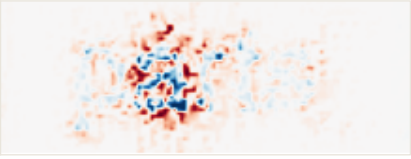
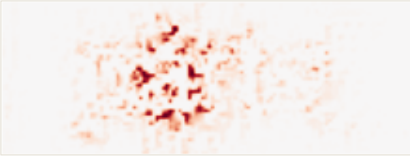

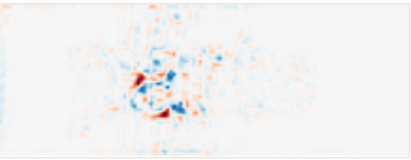
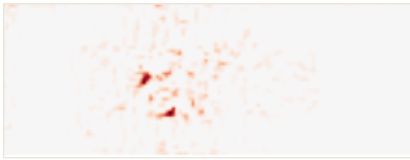

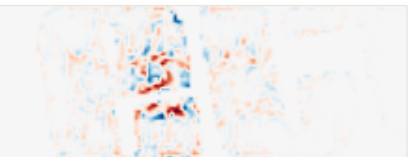
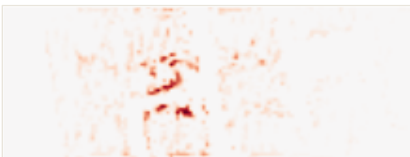

		replacement						insertion		deletion	
		easy		random		hard					
		MSE	I _{avg}	MSE	I _{avg}	MSE	I _{avg}	MSE	I _{avg}	MSE	I _{avg}
Opt-Basic	Courier	25.4	266	30.3	313	36.7	321	25.4	309	13.6	43
	Georgia	52.0	292	59.4	318	67.5	328	41.6	337	45.0	169
	Helvetica	52.1	301	60.2	328	68.2	340	47.1	321	45.0	178
	Times	49.9	294	56.1	324	61.6	345	41.7	314	44.3	172
	Arial	56.3	304	65.3	327	73.8	341	48.3	324	51.0	176
example		parts		parts		parts		parts		parts	
Opt-WM	Courier	16.7	116	20.1	96	20.4	95	31.1	29	3.2	13
	Georgia	31.6	30	35.1	32	38.3	37	21.7	12	16.2	9
	Helvetica	33.3	31	37.0	42	38.8	53	25.1	13	16.5	9
	Times	30.3	22	33.9	26	35.9	36	19.2	11	15.4	8
	Arial	37.2	30	40.4	45	43.6	50	25.4	16	19.4	10
example		parts		parts		parts		parts		parts	
target output		pants		pacts		pasts		partis		pars	

letter-level attack in word images



Visualization Explanation: Saliency Map

- The first to the last line are **clean**, **gray** and **watermark** backgrounds.
- **Reduced contrast** is beneficial to reduce noise.

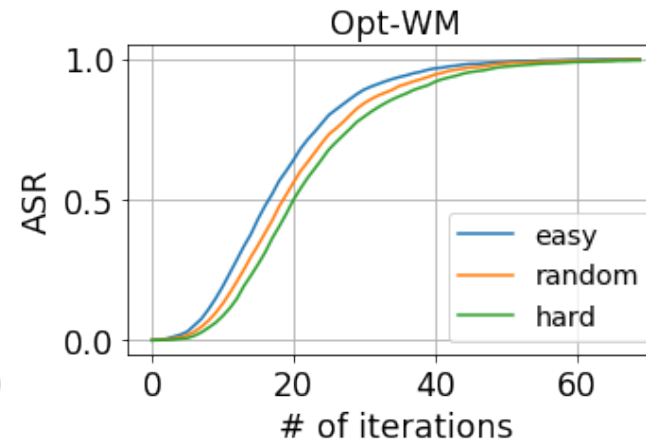
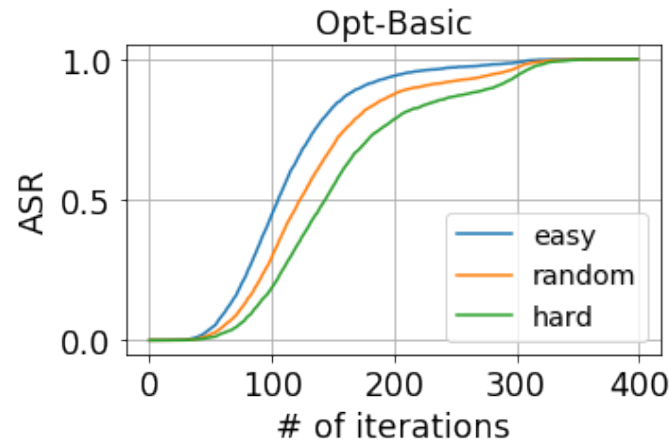
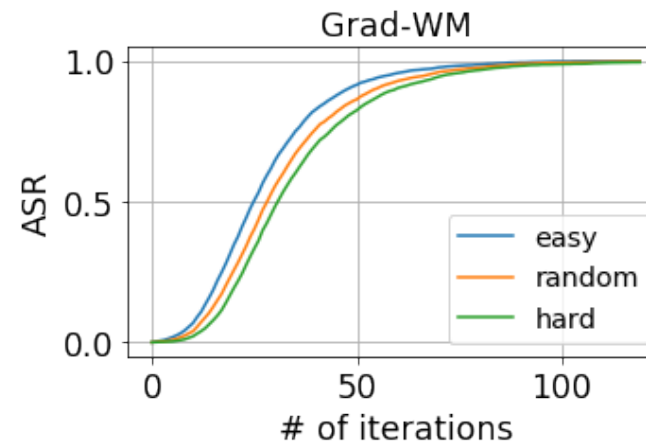
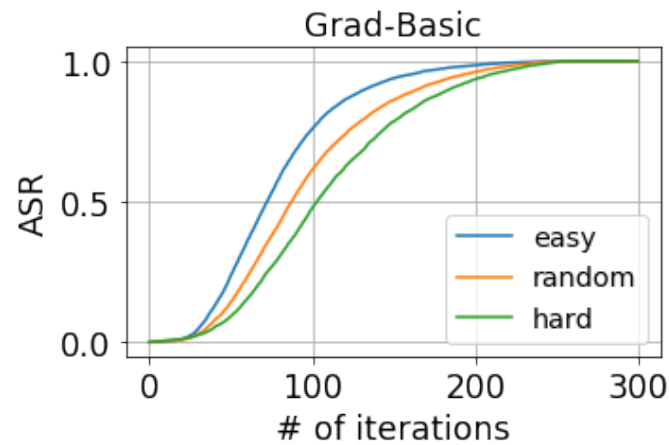
clean images	adversarial images	saliency map	MSE	saliency map+	MSE+	saliency map-	MSE-
parts	parts		55.82		34.77		21.05
parts	parts		15.47		6.92		8.55
parts	parts		24.97		14.13		10.84

The target output is "ports".



Attack Difficulty:

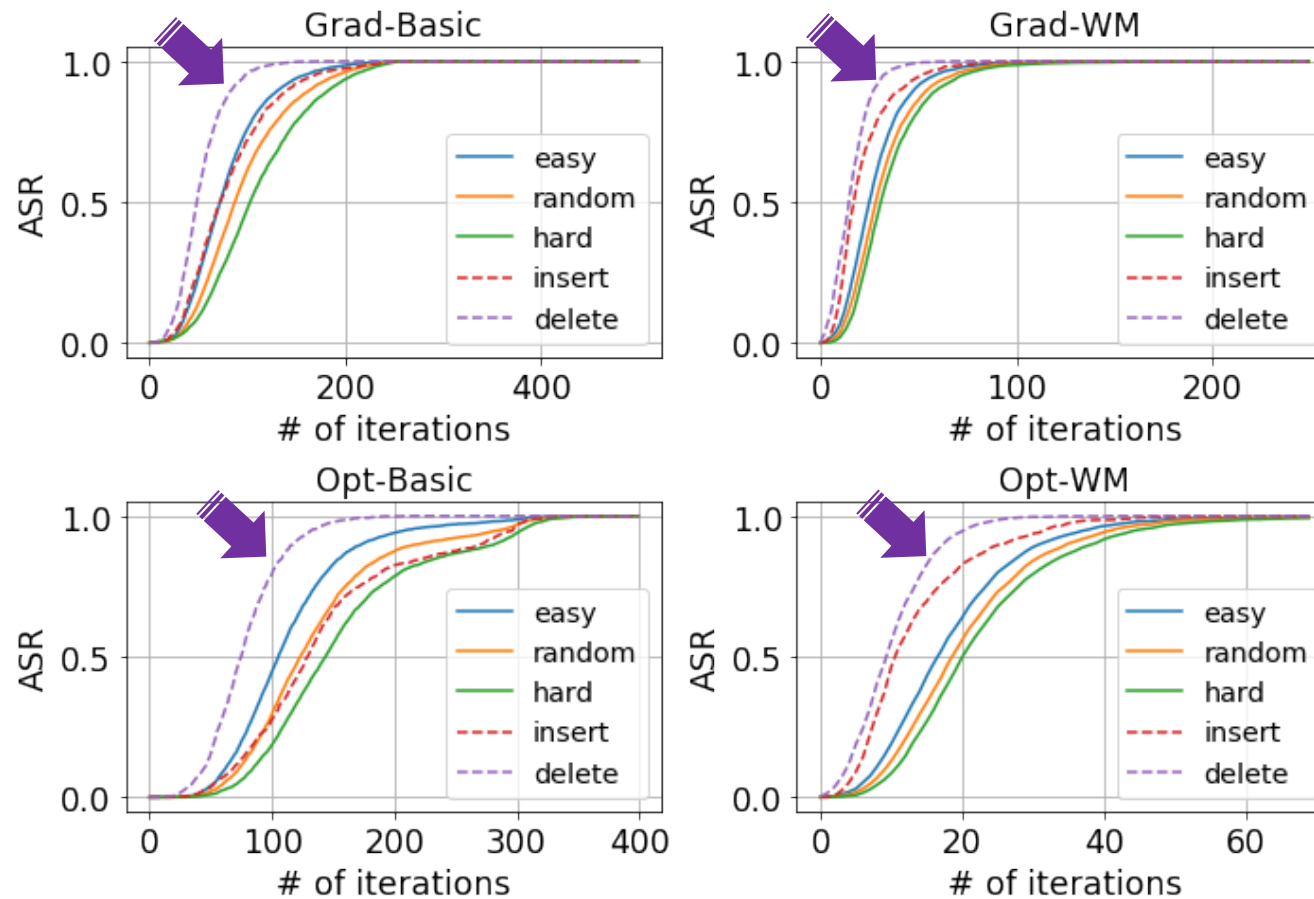
Easy < **Random** < **Hard**



letter-level attack in word images with Arial font

Attack Difficulty:

Delete < Insert & Replace

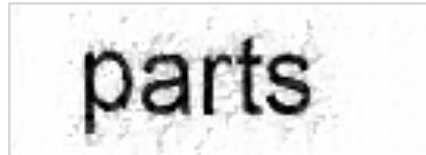


letter-level attack in word images with Arial font

Word-Level Attacks: WM is More Natural

taupe

Grad-Basic:



Grad-WM:

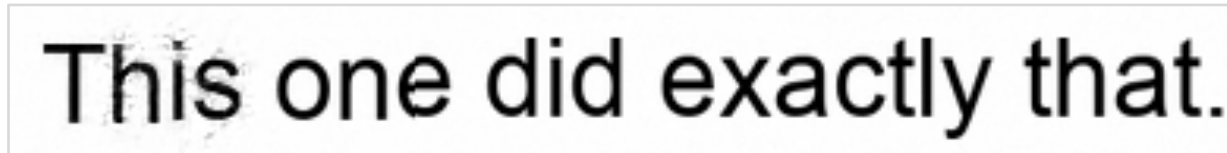


WM attack:

- 56% lower noise
- 50% less iterations

Tale one did exactly that.

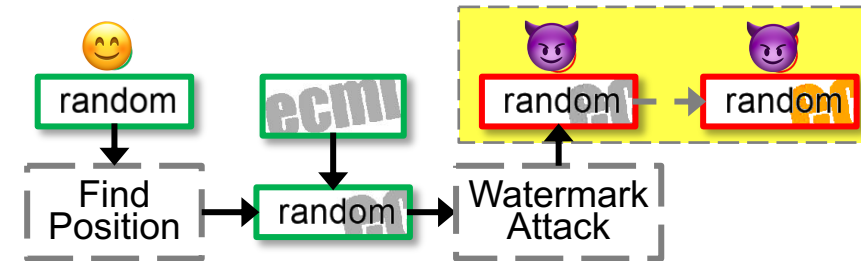
Grad-Basic:



Grad-WM:



Improve Readability: Full-Color Watermarks



Input Image: ● positive

This is one of the funniest movies I have ever seen. This, in my opinion, is Rob Lowe at his best. I'm not quite sure why this film has gotten such a low rating. I guess you either love it or hate it, but if nothing else, it is definitely worth a rental.

OCR Output: ● negative

This is one of the scariest movies I have ever seen. This, in my opinion, is Rob Lowe at his worst. I'm not quite sure why this film has gotten such a high rating. I guess you either love it or hate it, but if nothing else, it is not definitely worth a rental.

Other Colored Adversarial Examples:

easy	type	toy	run	where	pan
	tyre	foy	rum	there	pen
hard	wash	humour	expect	good	shaft
	cash	rumour	expert	hood	shift
word	rob	warm	sing	pill	off
	tow	heir	blue	dirt	add

Conclusion

- We propose **fast adversarial watermark attacks (FAWA)** on sequence-based **OCR** models.
 - Sequential labeling task——CTC loss
 - Background pollution——watermark
 - **Natural** watermark-style noise.
 - **Lower** perturbation level.
 - **Faster** attack speed.

WATERMARK



THANK YOU

Lu Chen¹, Jiao Sun², Wei Xu¹

¹IIS, Tsinghua University

²CS, University of Southern California



清华大学
Tsinghua University



交叉信息研究院
Institute for Interdisciplinary
Information Sciences

