



清華大學
Tsinghua University



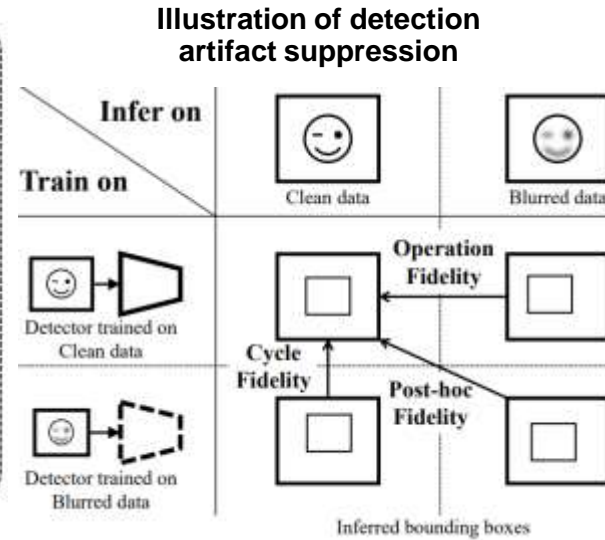
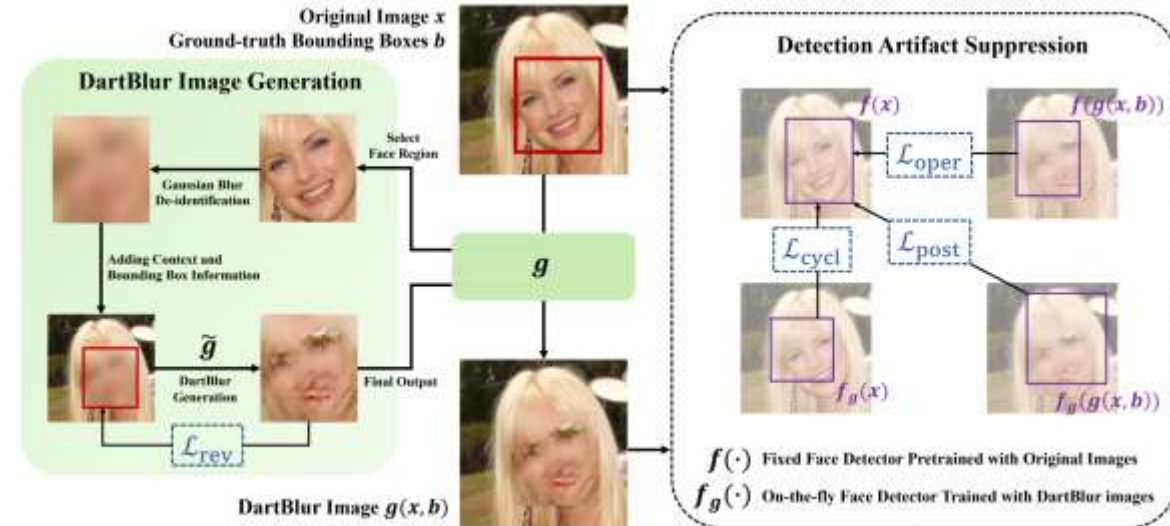
DartBlur: Privacy Preservation with Detection Artifact Suppression

Baowei Jiang*, Bing Bai*, Haozhe Lin*, Yu Wang,
Yuchen Guo, Lu Fang

Tsinghua University

*authors contributed equally

Summary of highlights



Algorithm 1: Training Algorithm for DartBlur

Input: Dataset \mathcal{D} , pretrained face detector f
Hyper-parameter: Step size α and β , threshold ϵ_{rev}
Output: Optimized parameters θ_g^*

- 1 Randomly initialize θ_g and θ_{f_g} ;
- 2 **while not converge do**
- 3 Sample a batch of data $\mathbf{x}, \mathbf{b} \sim \mathcal{D}$;
- 4 // Optimize g with f_g fixed
- 5 Update $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{rev}}(g, \mathbf{x}, \mathbf{b}, \epsilon_{\text{rev}})$;
- 6 Update $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{det}}(f(g(\mathbf{x}, \mathbf{b})), f(\mathbf{x}))$;
- 7 Update $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{det}}(f_g(g(\mathbf{x}, \mathbf{b})), f(\mathbf{x}))$;
- 8 // Optimize g considering second-order effects
- 9 Compute adapted parameters with gradient descent:
 $\theta'_{f_g} = \theta_{f_g} - \alpha \nabla_{\theta_{f_g}} \mathcal{L}_{\text{det}}(f_g(g(\mathbf{x}, \mathbf{b})), \mathbf{b})$;
- 10 Update $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{det}}(f_g(\mathbf{x}; \theta'_{f_g}), f(\mathbf{x}))$;
- 11 // Optimize f_g with g fixed
- 12 Update $\theta_{f_g} \leftarrow \theta_{f_g} - \beta \nabla_{\theta_{f_g}} \mathcal{L}_{\text{det}}(f_g(g(\mathbf{x}, \mathbf{b})), \mathbf{b})$;
- 13 **end**

- Propose a new blur-based privacy preservation model DartBlur by taking into account the actual accessibility of the model, review convenience, and detection artifact suppression simultaneously.
- Introduce four novel training objectives that each directly addresses the desired properties and design an adversarial training strategy with a second-order optimization for model training.
- Effectively protect personal privacy while suppressing detection artifacts on various benchmarks.



清華大學
Tsinghua University



DartBlur: Privacy Preservation with Detection Artifact Suppression

Baowei Jiang*, Bing Bai*, Haozhe Lin*, Yu Wang,
Yuchen Guo, Lu Fang

Tsinghua University

*authors contributed equally

Motivation

- Blur-based methods are already preferred to protect private information in real-world data, e.g. ImageNet and Ego4D.



ImageNet



Ego4D

Motivation

- Images with blurred faces produce artifacts for object detection task, and affect the performance of the face detector

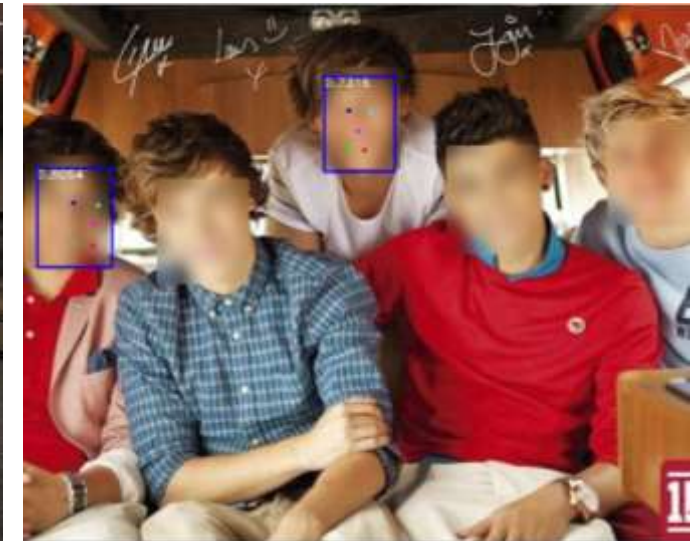
‘Blurred’ RetinaFace
on clean faces

‘Clean’ RetinaFace on
blurred faces



ImageNet

Ego4D



Motivation

- Based on visualization and evaluation, both Gaussian blur and Pixelation produce artifacts on the face detector

easy
medium
hard



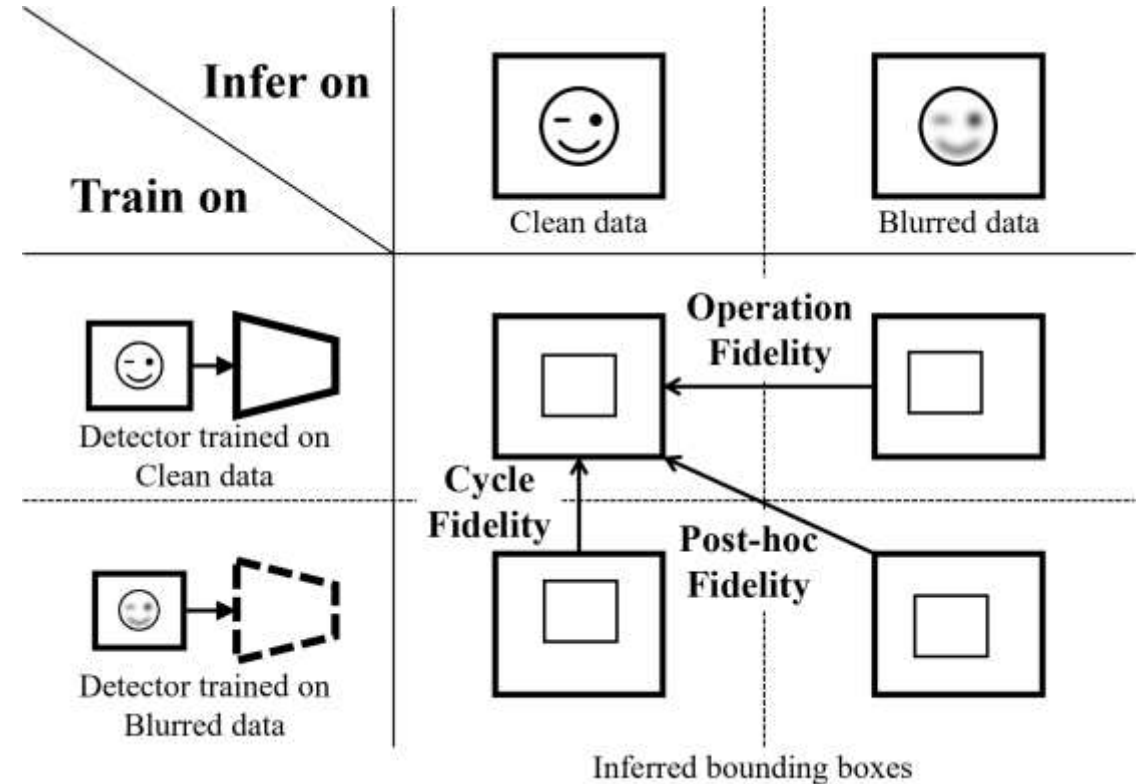
RetinaFace training and evaluation

Train \ val dataset	Origin	Gaussian Blur	Pixelation
Origin	0.921005 0.901564 0.745564	0.784570 0.790573 0.651209	0.764316 0.708425 0.544112
Gaussian Blur	0.000000 0.000037 0.011646	0.996980 0.993339 0.897396	
Pixelation	0.000000 0.000076 0.000121		0.996482 0.995469 0.952062

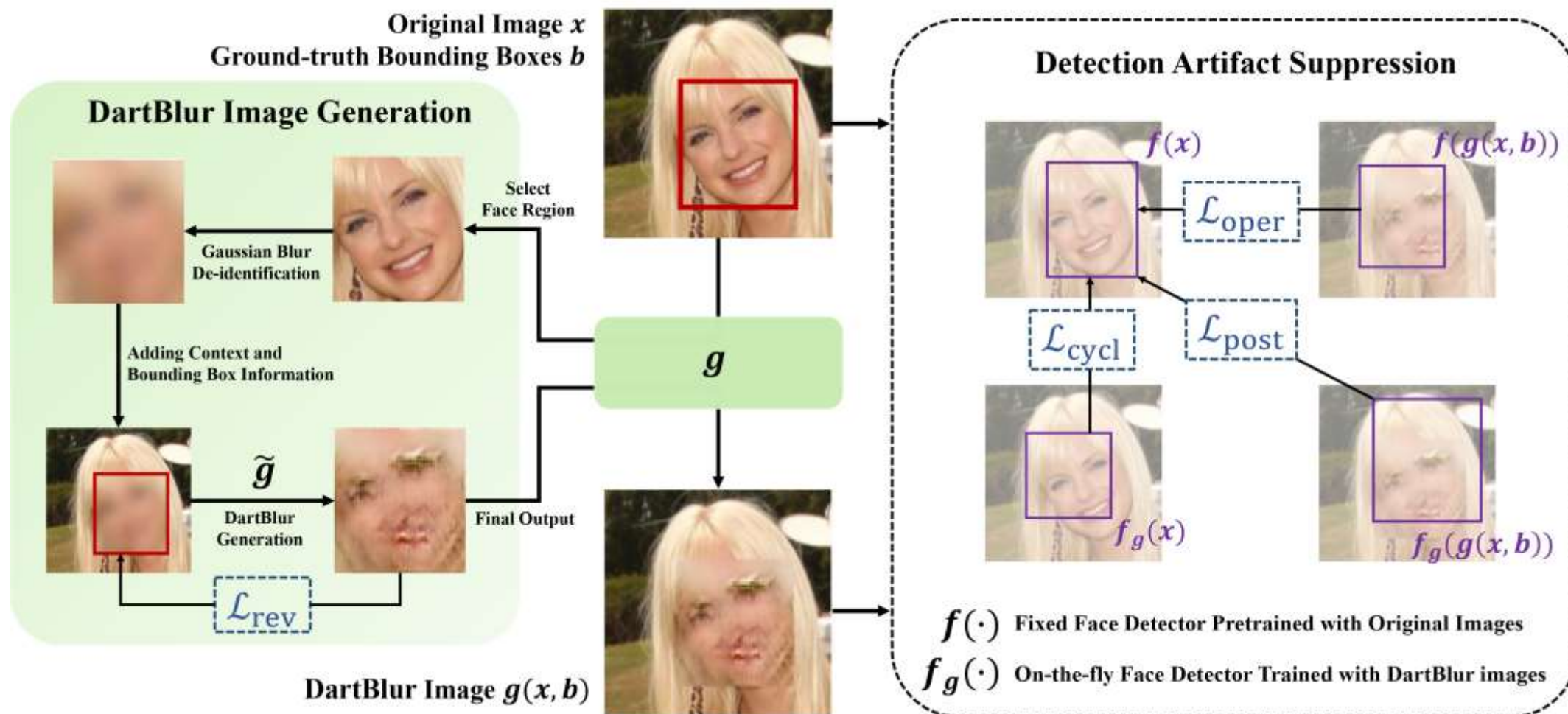
Overview

We design a novel blur-based privacy protection paradigm with the following goals:

- Accessibility.
- Review convenience.
- Detection artifact suppression.



Pipeline



Optimization

- Input image \mathbf{x} and bbox \mathbf{b} , Unet g
- Fixed \mathbf{f}_g
 - Solve $\min_{\mathbf{g}} \|\mathbf{g}(\mathbf{x}) - \mathbf{x}\|$
 - Solve $\min_{\mathbf{g}} \|\mathbf{f}(\mathbf{x}) - \mathbf{f}(\mathbf{g}(\mathbf{x}))\|$
 - Solve $\min_{\mathbf{g}} \|\mathbf{f}_g(\mathbf{x}) - \mathbf{f}(\mathbf{g}(\mathbf{x}))\|$
 - Solve $\min_{\mathbf{g}} \|\mathbf{f}'_g(\mathbf{x}) - \mathbf{f}(\mathbf{x})\|$, $\mathbf{f}_g \xrightarrow{\mathbf{b}} \mathbf{f}'_g$
- Fixed g
 - Train \mathbf{f}_g on $\mathbf{g}(\mathbf{x}, \mathbf{b})$ dataset

Algorithm 1: Training Algorithm for DartBlur

Input: Dataset \mathcal{D} , pretrained face detector f

Hyper-parameter: Step size α and β , threshold ϵ_{rev}

Output: Optimized parameters θ_g^*

```
1 Randomly initialize  $\theta_g$  and  $\theta_{f_g}$ ;
2 while not converge do
3   Sample a batch of data  $\mathbf{x}, \mathbf{b} \sim \mathcal{D}$ ;
4   // Optimize  $g$  with  $f_g$  fixed
5   Update  $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{rev}}(g, \mathbf{x}, \mathbf{b}, \epsilon_{\text{rev}})$ ;
6   Update  $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{det}}(f(g(\mathbf{x}, \mathbf{b})), f(\mathbf{x}))$ ;
7   Update  $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{det}}(f_g(g(\mathbf{x}, \mathbf{b})), f(\mathbf{x}))$ ;
8   // Optimize  $g$  considering second-order effects
9   Compute adapted parameters with gradient descent:
10   $\theta'_{f_g} = \theta_{f_g} - \alpha \nabla_{\theta_{f_g}} \mathcal{L}_{\text{det}}(f_g(g(\mathbf{x}, \mathbf{b})), \mathbf{b})$ ;
11  Update  $\theta_g \leftarrow \theta_g - \beta \nabla_{\theta_g} \mathcal{L}_{\text{det}}(f_g(\mathbf{x}; \theta'_{f_g}), f(\mathbf{x}))$ ;
12  // Optimize  $f_g$  with  $g$  fixed
13  Update  $\theta_{f_g} \leftarrow \theta_{f_g} - \beta \nabla_{\theta_{f_g}} \mathcal{L}_{\text{det}}(f_g(g(\mathbf{x}, \mathbf{b})), \mathbf{b})$ ;
14 end
```

Evaluation

- DartBlur trained with WIDER FACE successfully suppressed detection artifacts, and was generalizable across different datasets, and could be well generalized across different datasets and architectures for detection.

Dataset	Fidelity	Block	Blur	Pixel.	DartBlur
WIDER FACE	Oper. Fid.	19.20	<u>83.10</u>	46.35	96.76
	Post-hoc Fid.	77.77	<u>84.04</u>	80.08	84.70
	Cycle Fid.	0.10	1.97	<u>7.64</u>	47.22
FDDB	Oper. Fid.	4.63	<u>89.00</u>	4.82	98.06
	Post-hoc Fid.	85.49	<u>93.41</u>	88.34	94.64
	Cycle Fid.	0.00	<u>0.05</u>	0.02	41.93
CrowdHuman	Oper. Fid.	16.69	<u>75.74</u>	45.26	81.64
	Post-hoc Fid.	58.23	<u>57.51</u>	<u>59.60</u>	62.52
	Cycle Fid.	0.83	<u>36.72</u>	28.75	45.11

Evaluation results of detection artifact suppression and cross-dataset transferability

Architecture	Fidelity	Block	Blur	Pixel.	DartBlur
PyramidBox	Oper. Fid.	21.34	<u>84.60</u>	30.55	95.18
	Post-hoc Fid.	<u>75.04</u>	70.59	65.18	75.16
	Cycle Fid.	0.01	1.70	<u>10.98</u>	24.68
YOLOv5	Oper. Fid.	35.93	<u>84.84</u>	35.22	96.17
	Post-hoc Fid.	85.68	<u>87.44</u>	86.01	91.72
	Cycle Fid.	0.21	<u>10.00</u>	0.36	37.15

Cross-architecture transferability evaluation on WIDER FACE

Privacy recovering

An indomain model is trained to reconstruct the original images from DartBlur images. As illustrated in the following equation, the function \bar{g} acts only on the faces region,

$$\bar{g}(g(x, b), b) = x \odot (1 - b) + \bar{g}(g(x, b)) \odot b,$$

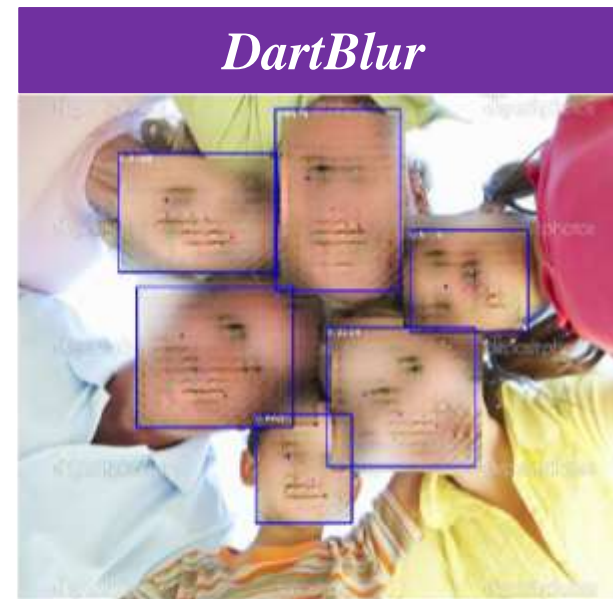
where \bar{g} maintain the same structure as \tilde{g} . The figure shows 3 cases and demonstrates that erased privacy information cannot be trivially restored.



Dartblur-Faces



Detection results comparison



Thanks!