



# Перенесення адаптивної анонімізації в задачі класифікації зображень

---

## Актуальність роботи



Алгоритми генеративного штучного інтелекту в наш час перебувають в апогеї свого розвитку і в свою чергу вимагають великої кількості зображень для коректної роботи.

# Негативний ефект анонімізації зображень на роботу моделі

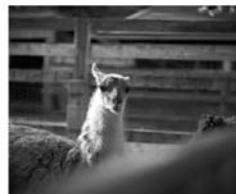
Однак виконання цієї норми шляхом використання наївних методів анонімізації часто призводить до значного погіршення результатів роботи класифікаційних нейронних мереж.



bicycle-built-for-two  
41.39% confidence



unicycle  
54.72% confidence



llama  
99.50% confidence



fountain  
26.30% confidence



ringlet  
92.24% confidence



acorn  
57.85% confidence

# Методи анонімізації зображень



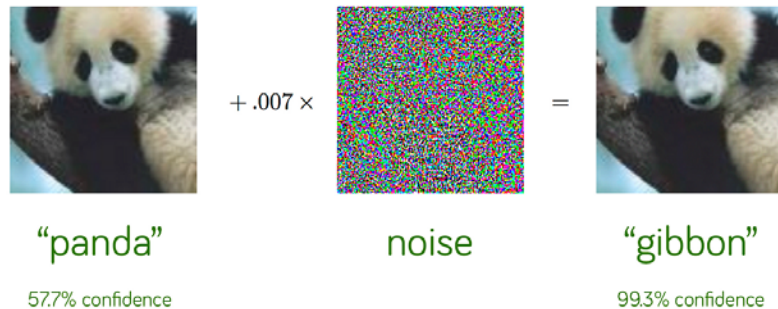
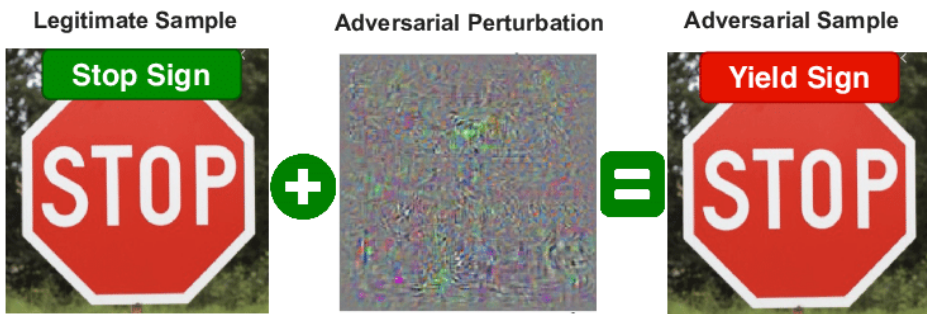
PREDICTED CONCEPT	PROBABILITY
woman	0.991
girl	0.978
portrait	0.974
one	0.973
adult	0.961



PREDICTED CONCEPT	PROBABILITY
woman	0.992
girl	0.986
portrait	0.984
one	0.978
adult	0.963

- Використання Генеративно Змагальних Мереж для анонімізації облич.
- Використання підходу “Зворотня Супер Роздільна Здатність” (ISR). Така роздільна здатність, безумовно, дозволяє зберегти приватність даних.
- Градієнтний метод ітеративної модифікації зображення.

# Змагальні атаки



# Явище переносимості змагальних прикладів

Prediction based on AE with original label “cat”



	(1)	(2)	(3)
F1: Source model	X: “dog”		
F2: Target model	✓: “cat”	X: “frog”	X: “dog”

Definitions of adversarial transferability

Non-targeted transferability	Non-transferable	Transferable	
Targeted transferability (target class: “dog”)	Non-transferable		Transferable
<b>Class-aware transferability</b>	Unfooled	Different mistake	Same mistake



## Переносимість адаптивної анонімізації

Так як задача адаптивної анонімізації по суті є протилежною до задачі створення змагальних атак.

Ми вирішили перевірити базовий градієнтний алгоритм на явище переносимості.

Назва моделі	Точність класифікації
MobileNetV2_100 (базова модель)	100%
ResNet52	60.6%
ResNet101	69.1%
ResNet152	71.9%
VGG13	30%
VGG16	29.4%
VGG19	30.8%
MobileNetV2_050	32.9%
MobileNetV2_140	41.8%
MobileNetV3_large_100	19.8%



## Переносимість адаптивної анонізації

Наступним кроком було вирішено, спробувати покращити базовий алгоритм, розширивши його застосування на декілька моделей водночас.

MobileNetV2_100 (базова модель)	100%
MobileNetV2_140	100%
ResNet101	100%
ResNet152	100%
VGG16	100%
VGG19	100%
Xception61	100%





**Дякую за увагу**