

Attacks against Neural Networks

Student: Martin Matak

Supervisor: Georg Weissenbacher



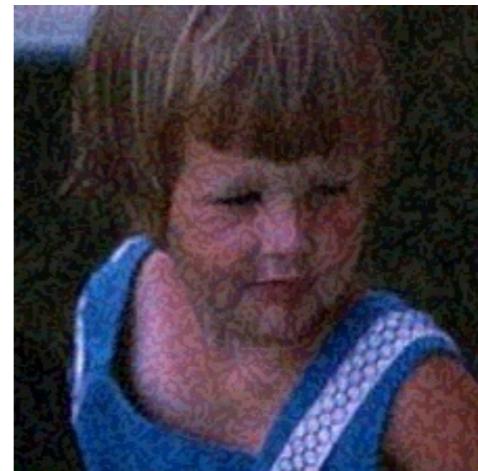
 How-Old.net
How old do I look? #HowOldRobot



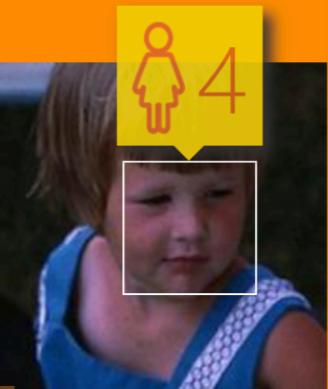
Sorry if we didn't quite get it right - [we are still improving this feature..](#)

[Try Another Photo!](#)

 Microsoft



 How-Old.net
How old do I look? #HowOldRobot

 A photograph of the same young girl, with a yellow speech bubble overlay containing a female icon and the number '4'. A white rectangular box highlights the area around her eyes and nose, indicating the facial recognition detection zone.

Sorry if we didn't quite get it right - [we are still improving this feature..](#)

[Try Another Photo!](#)

 Microsoft



How-Old.net

How old do I look? #HowOldRobot

Sorry if we didn't quite get it right - [we are still improving this feature..](#)

Try Another Photo!

Microsoft

A screenshot of the How-Old.net website. The interface features a large orange header with the site's logo and name. Below the header, a photograph of a young girl is displayed with a yellow speech bubble containing a red icon of a person and the number '4'. A white rectangular box highlights a specific area of the girl's face. At the bottom of the page, there is a message about the AI's learning process and a button to try another photo. The Microsoft logo is at the very bottom.

How-Old.net

How old do I look? #HowOldRobot

Sorry if we didn't quite get it right - [we are still improving this feature..](#)

Try Another Photo!

Microsoft

A screenshot of the How-Old.net website, identical in layout to the first one but showing a different result. The photograph of the young girl now has a yellow speech bubble containing a red icon of a person and the number '79'. A white rectangular box highlights a different area of the girl's face. The message at the bottom remains the same, along with the Microsoft logo.

Outline:

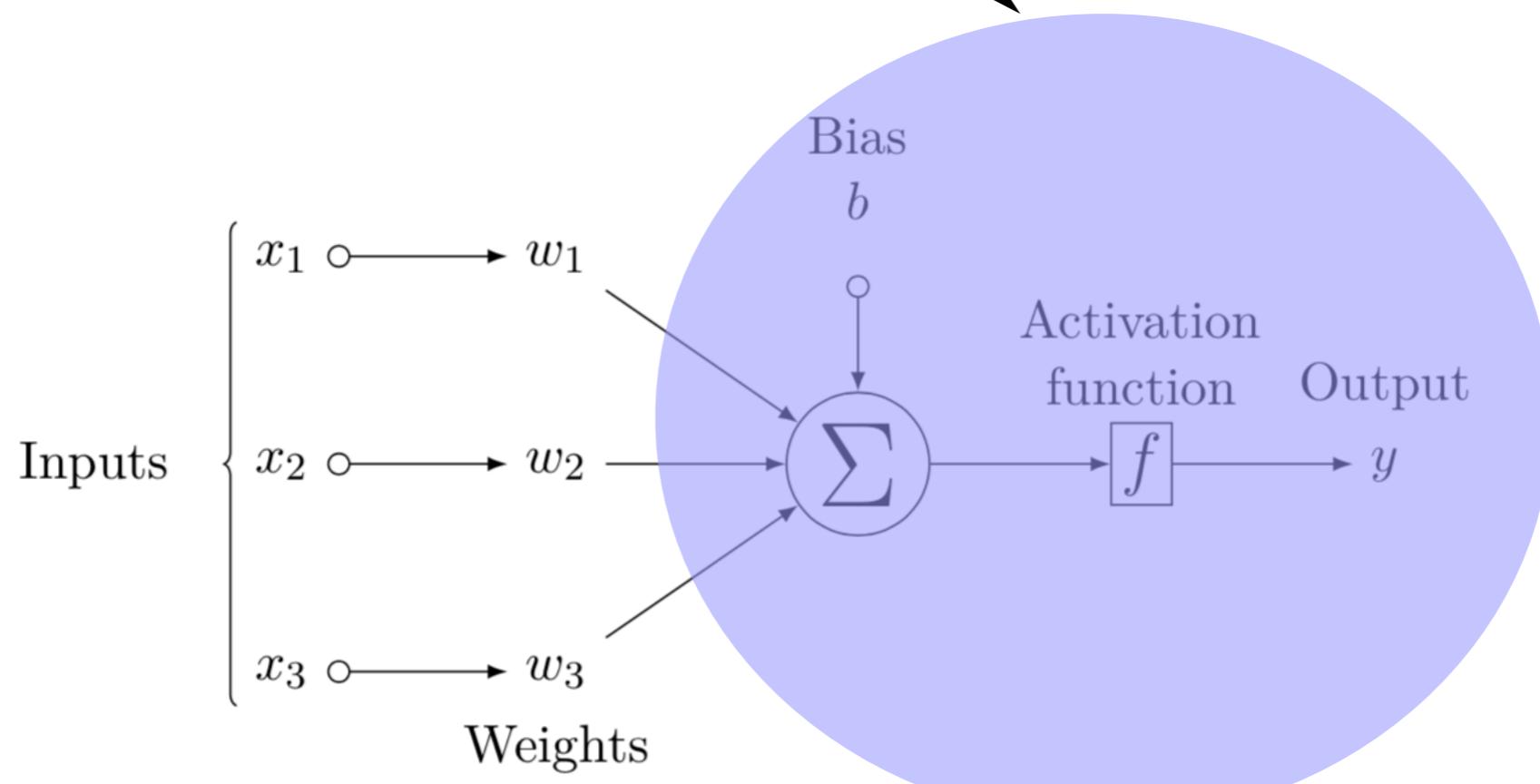
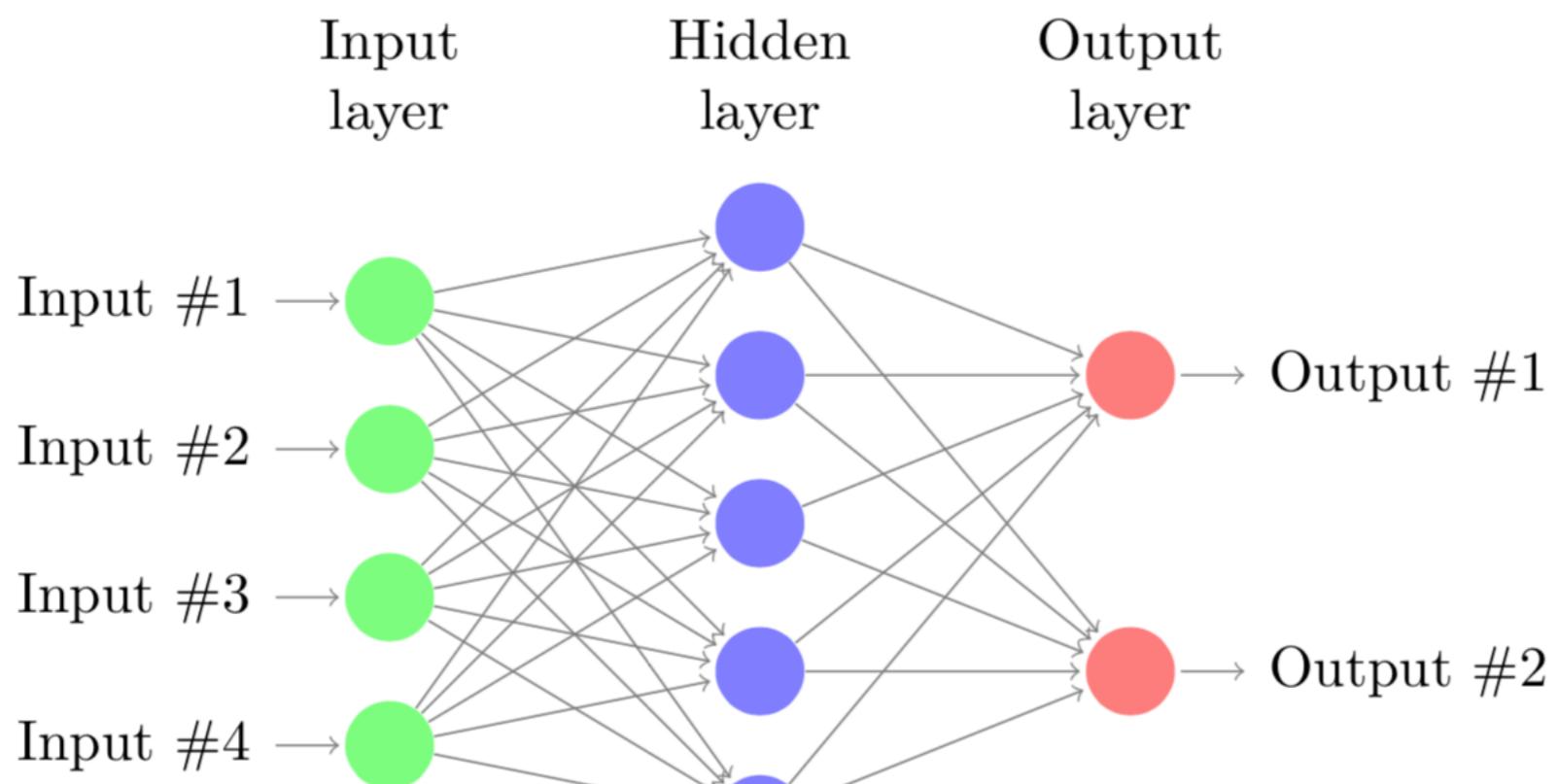
1. Train neural networks for age estimation

2. Attack the networks

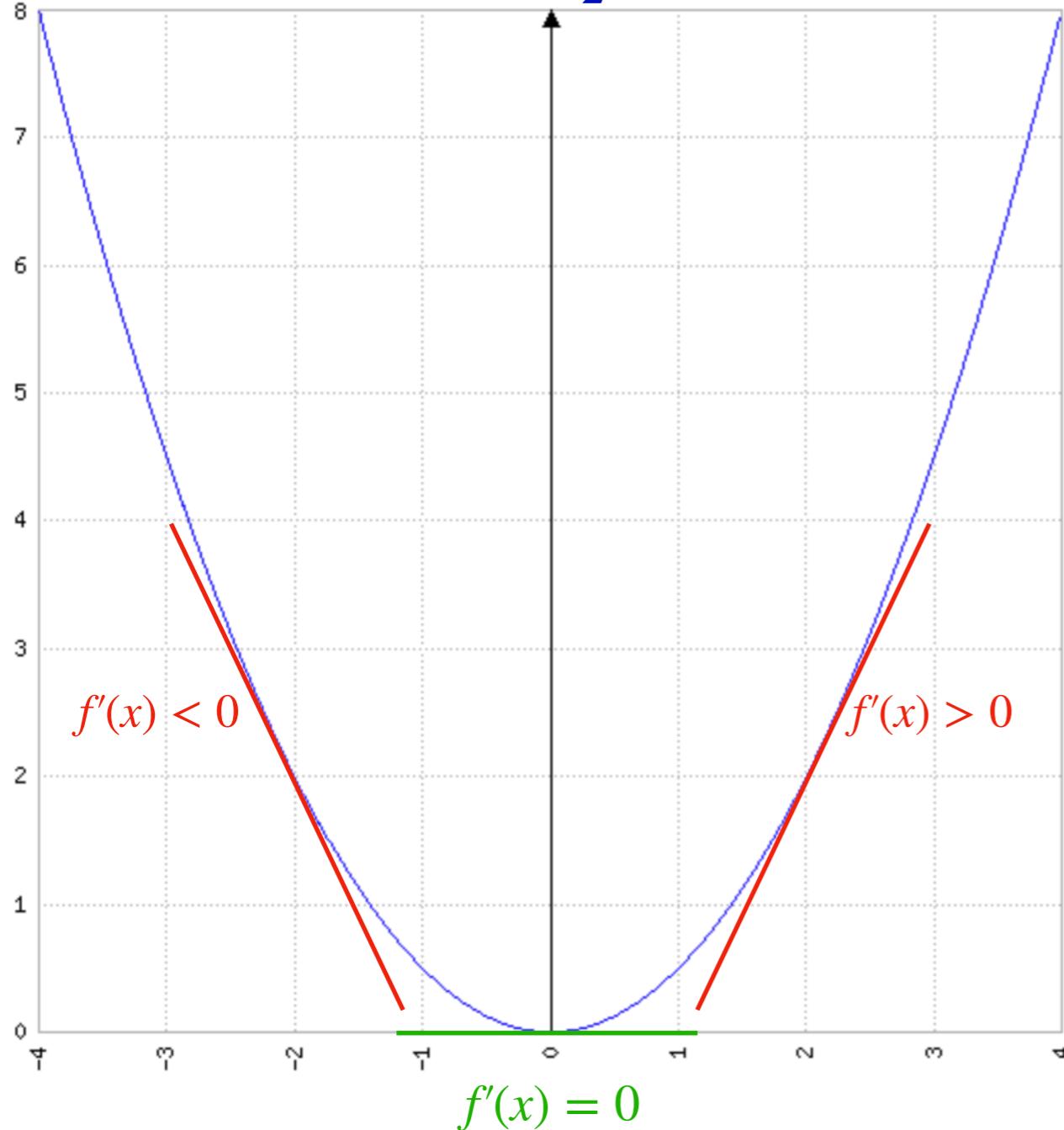
- **White-box environment**
- **Black-box environment**

3. The Semi-targeted Approach

Background



$$f(x) = \frac{1}{2}x^2$$



Gradient descent algorithm

input: α, p

while true:

gradient := $f'(p)$

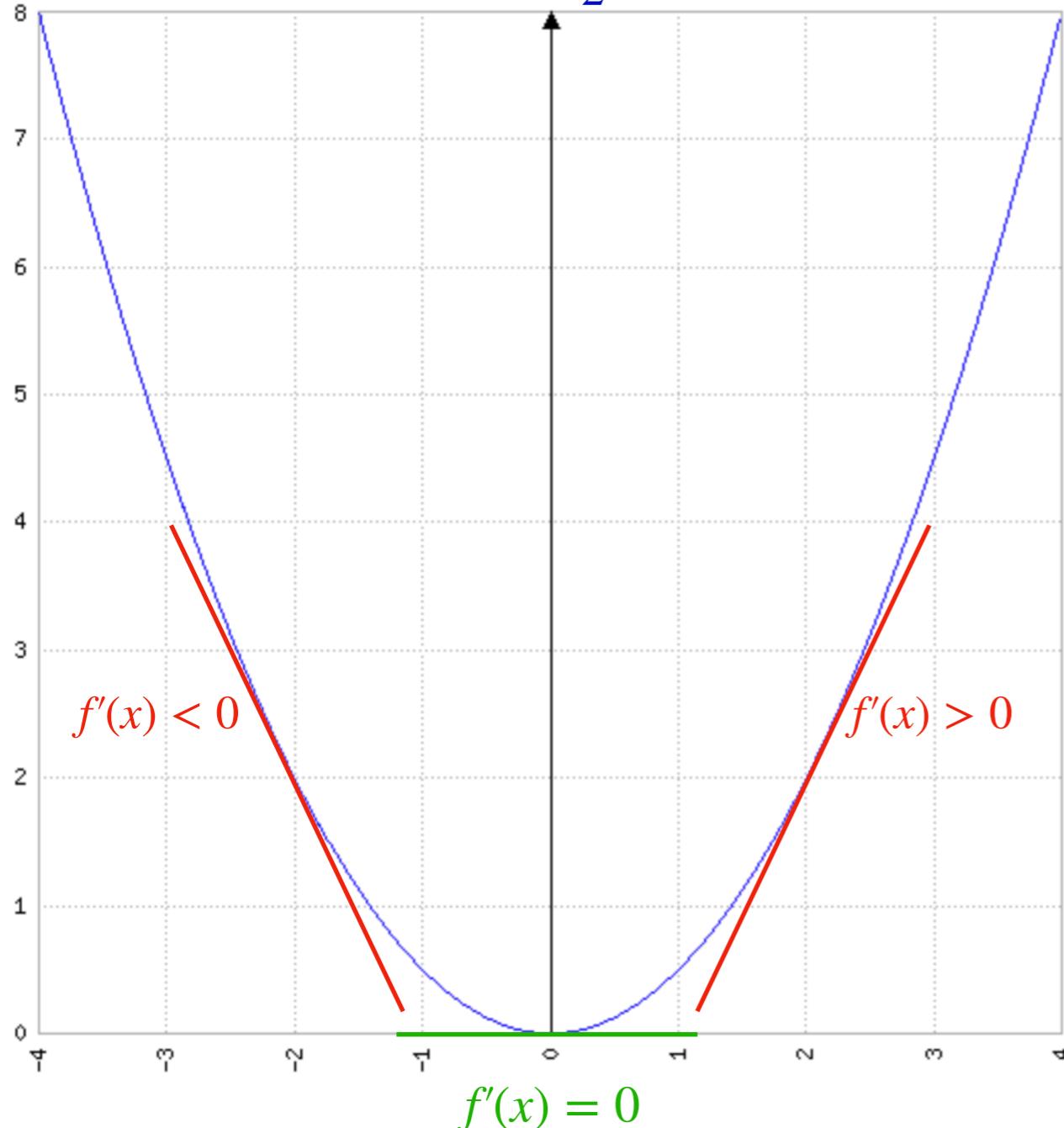
if gradient is small enough:

break

else:

$p := p - \alpha * \text{gradient}$

$$f(x) = \frac{1}{2}x^2$$



Gradient descent algorithm

input: α, p

while true:

gradient := $f'(p)$

if gradient is small enough:

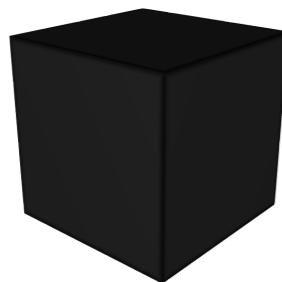
break

else:

$p := p - \alpha * \text{gradient}$

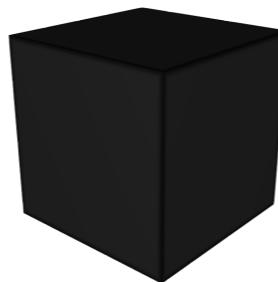
$$L(weights) = \frac{1}{S} \sum_{s=1}^S (y_s - \text{neural_network}(x_s; weights))^2$$

Current classification

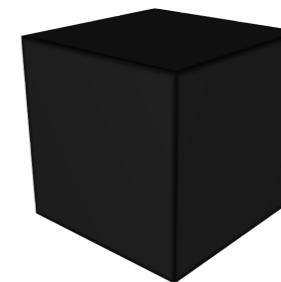
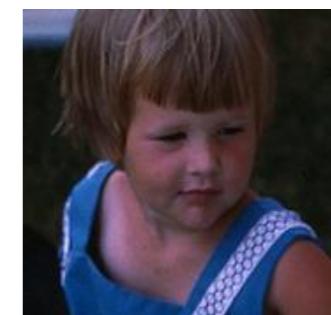


4

Misclassification



Targeted 90 misclassification



90

Trained models

Dataset: 30,843 images

	Architecture	Optimizer	Validation MAE
model 1	ResNet-50	SGD	5.15
model 2	ResNet-50	Adam	6.77
model 3	InceptionResNetV2	SGD	4.50
model 4	InceptionResNetV2	Adam	3.92

White-box Attacks

Fast gradient sign method (FGSM)

Fast gradient sign method (FGSM)

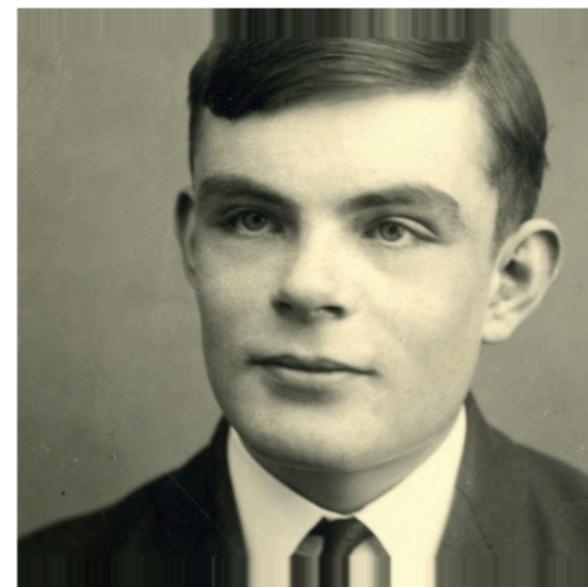
$$L(\text{weights}, x, \text{target_y}) = (\text{target_y} - \text{neural_network}(x; \text{weights}))^2$$

Fast gradient sign method (FGSM)

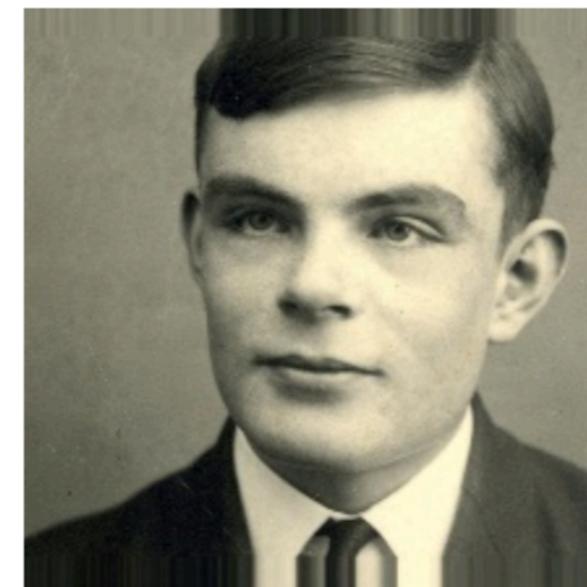
$$L(\text{weights}, x, \text{target_y}) = (\text{target_y} - \text{neural_network}(x; \text{weights}))^2$$

input: x , target_y , α , neural network (white-box)

```
gradient := L'(weights, x, target_y)
x := x - α * sign(gradient)
```



28



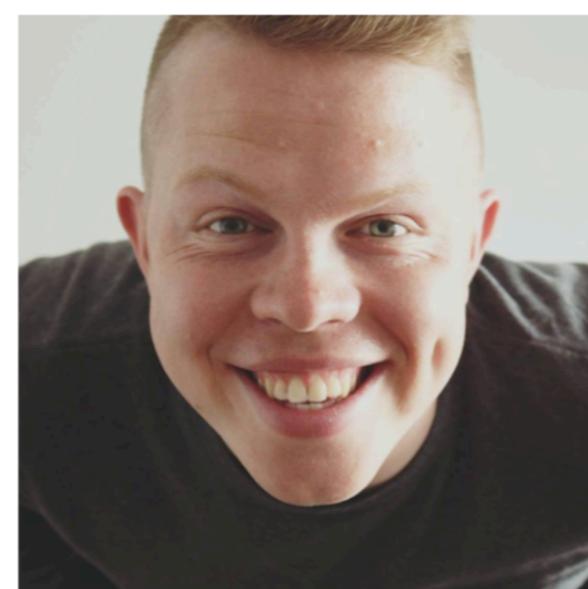
59



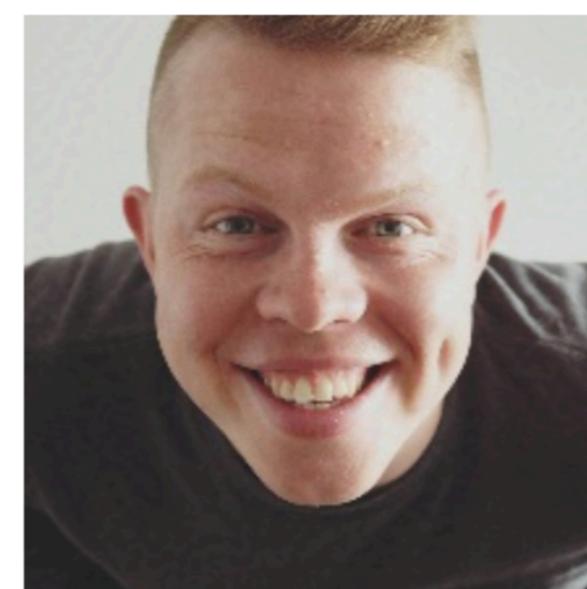
59



28



28



64

Carlini and Wagner (CW)

Carlini and Wagner (CW)

1. *Minimize* added perturbation
2. Must be adversarial (target class)
3. Must be a valid image

Carlini and Wagner (CW)

1. *Minimize* added perturbation

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

2. Must be adversarial (target class)

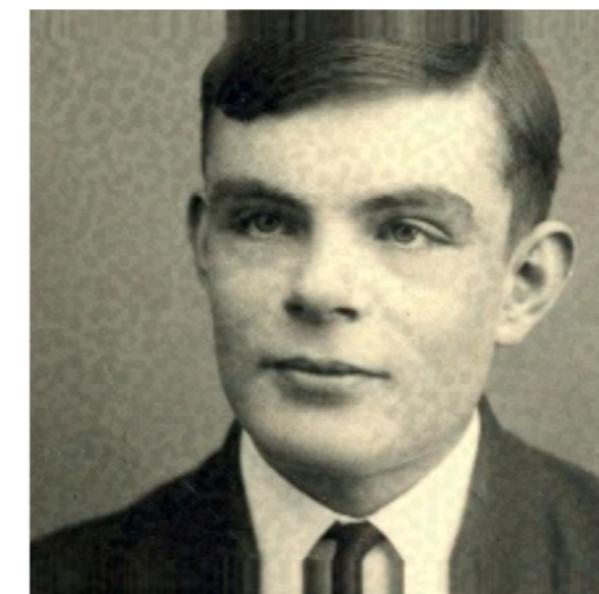
$$\text{such that } \mathcal{C}(\mathbf{x} + \boldsymbol{\delta}) = t$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

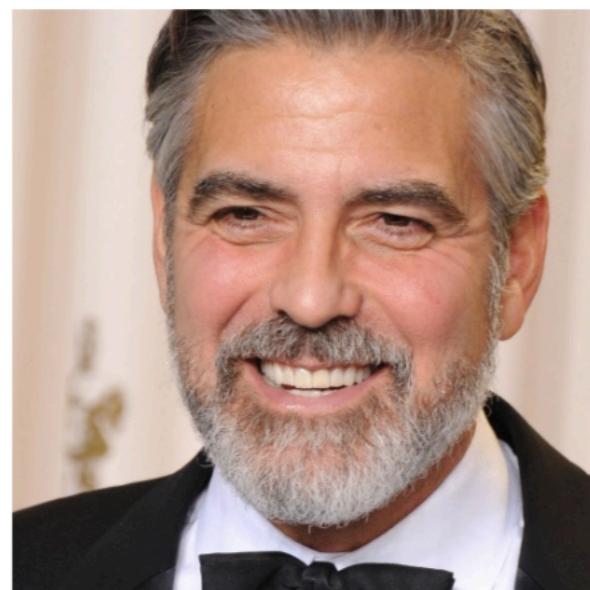
3. Must be a valid image



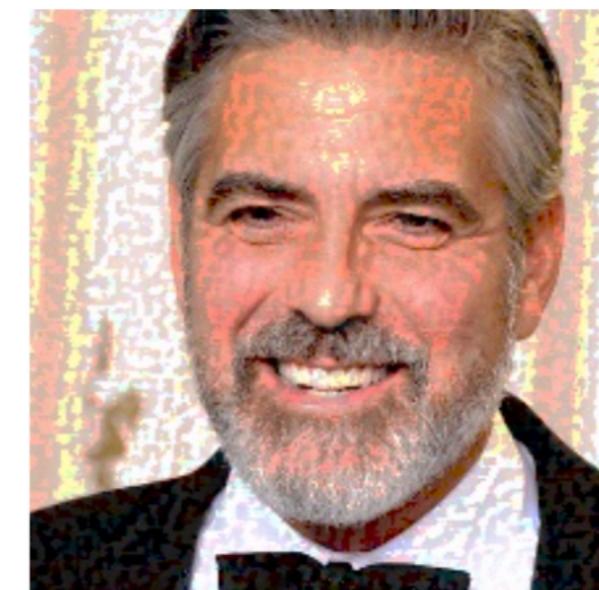
36



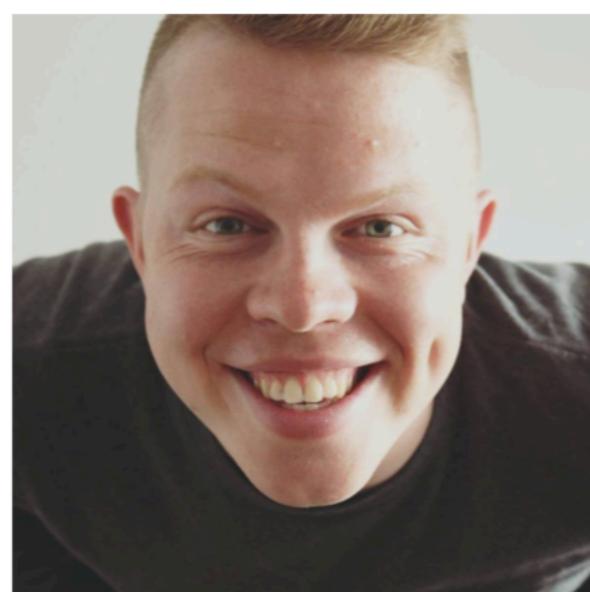
79



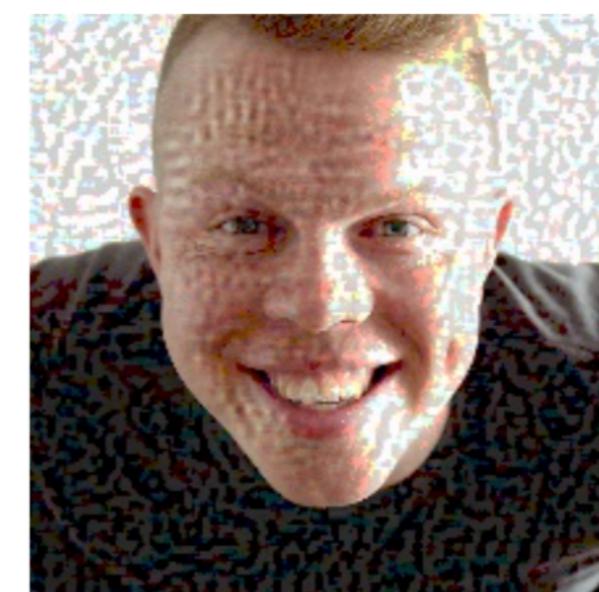
55



10



35



84

White-box: results

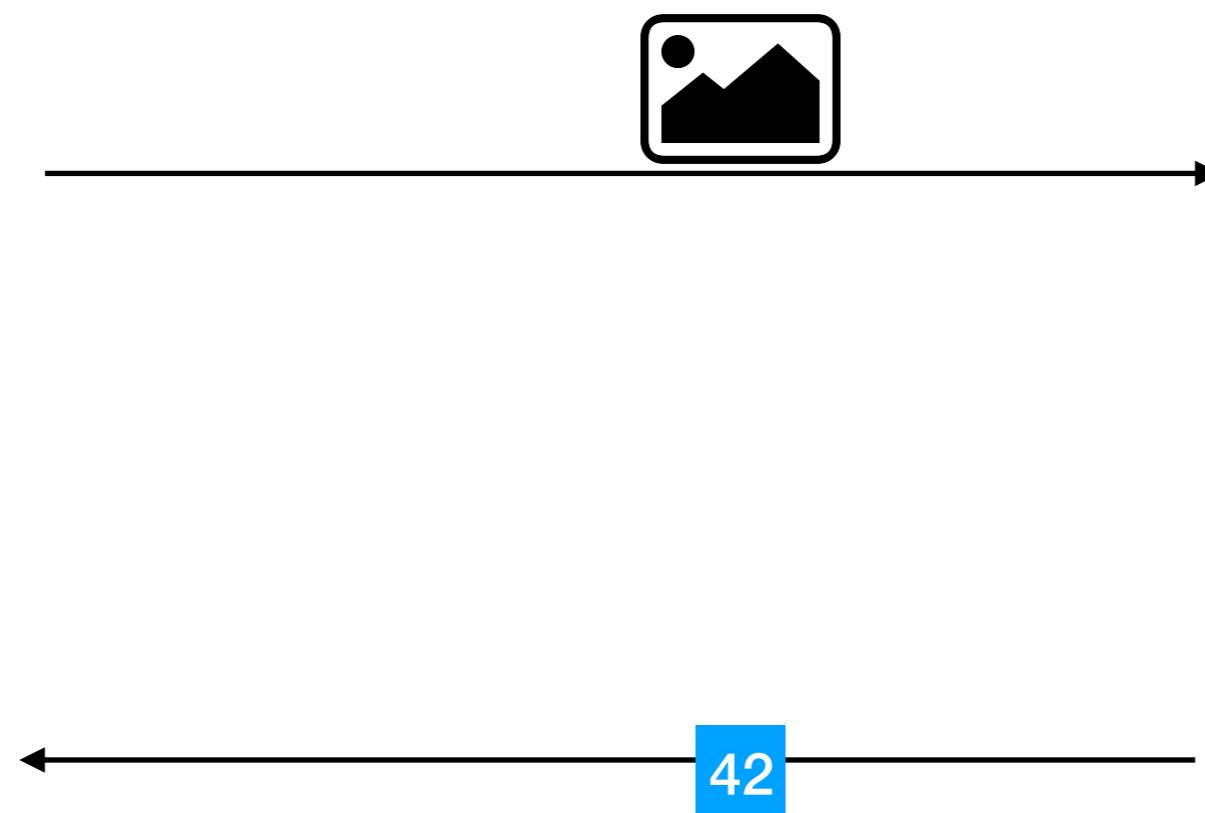
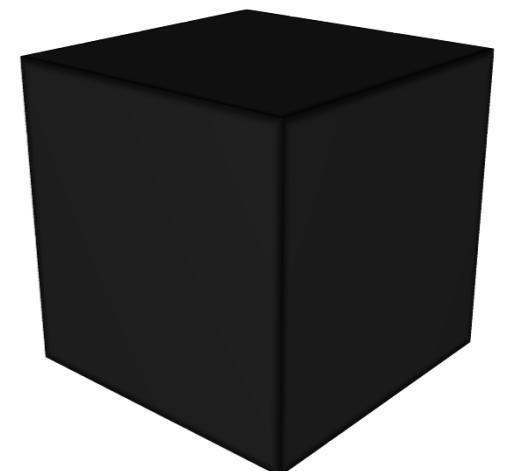
	clean samples MAE	FGSM MAE	FGSM AVG L2	CW MAE	CW AVG L2
model 1	6.69	32.25	1879.63	29.35	2977.28
model 2	9.50	47.79	1879.49	9.50	0.00
model 3	5.40	19.32	2510.92	7.51	1307.51
model 4	4.62	15.03	2511.10	6.79	1104.78

Black-box attacks

Adversary



Black-box Neural Network



Boundary attack

The starting image:



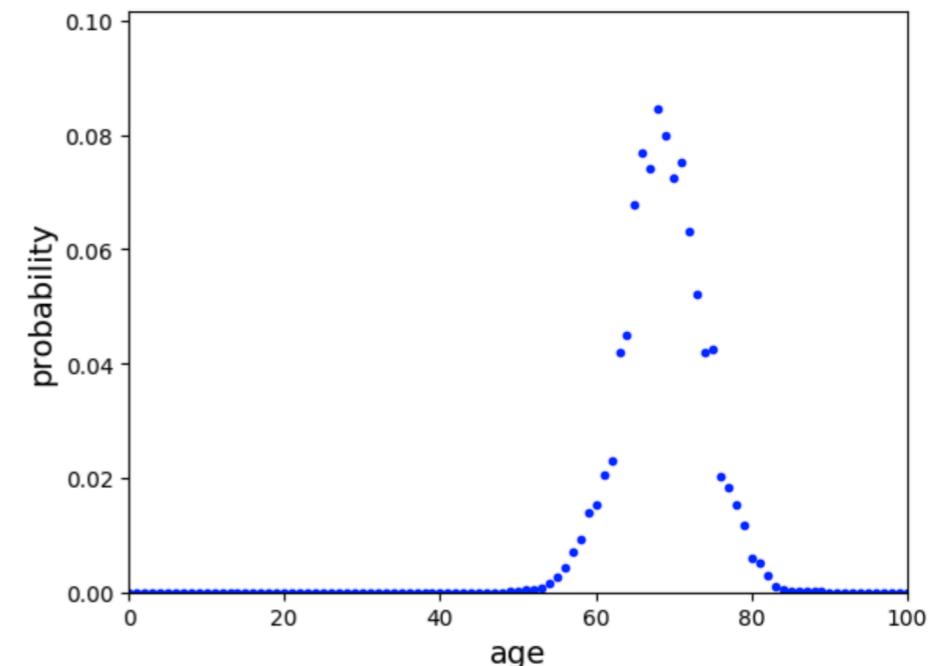
68

The target image

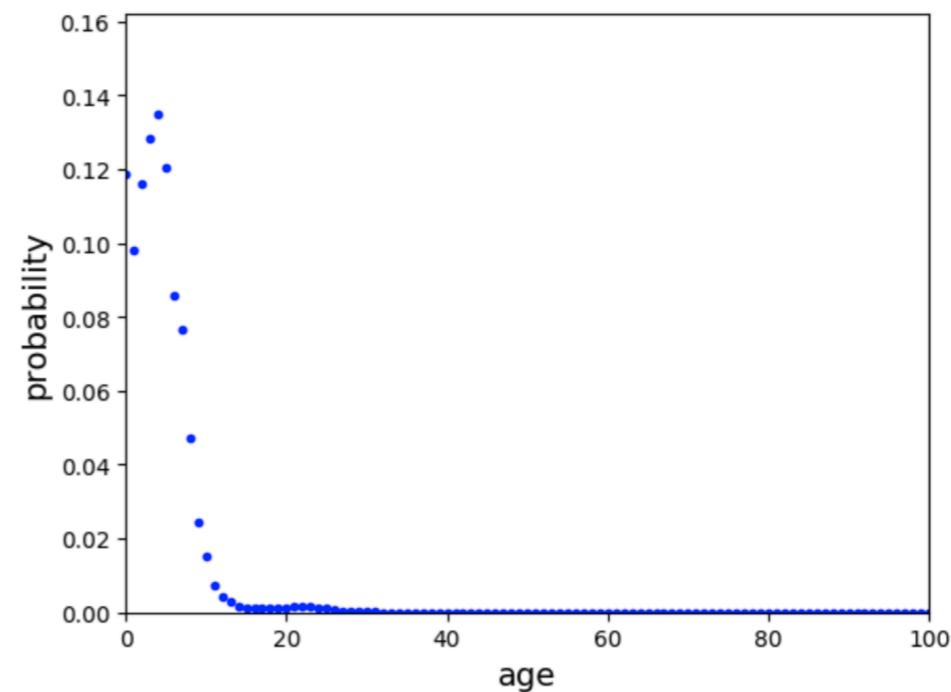


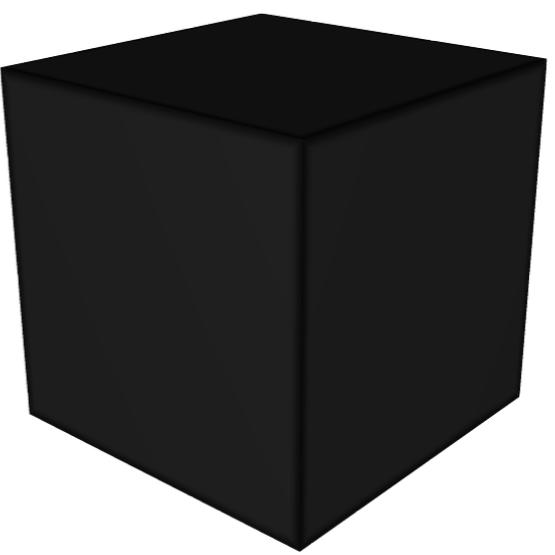
5

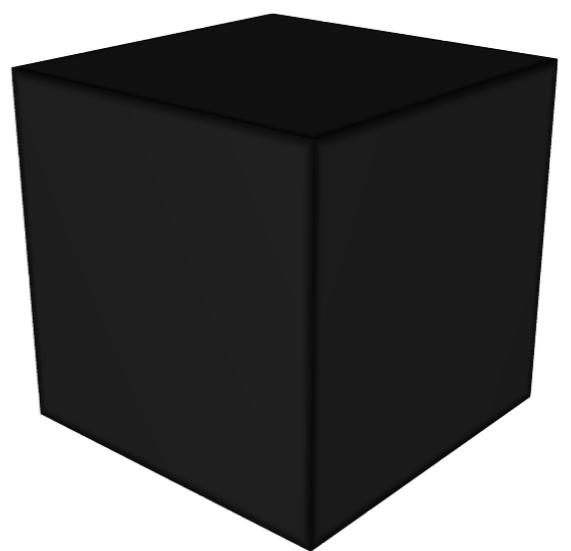
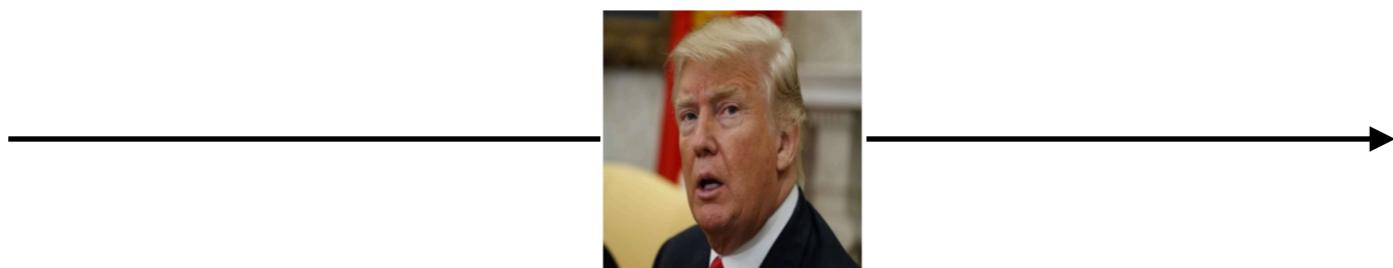
Confidence graph

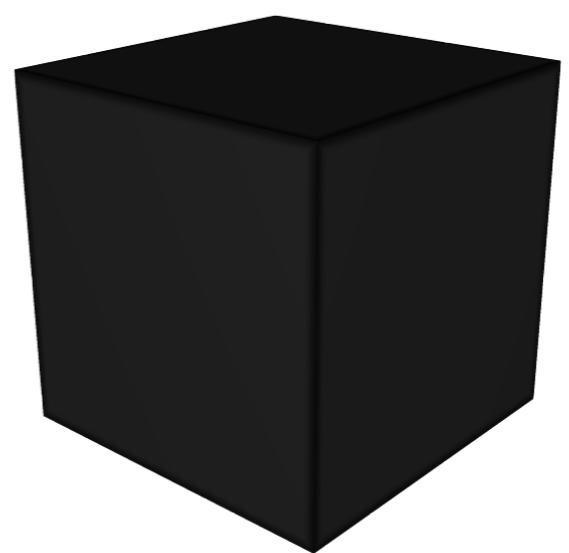
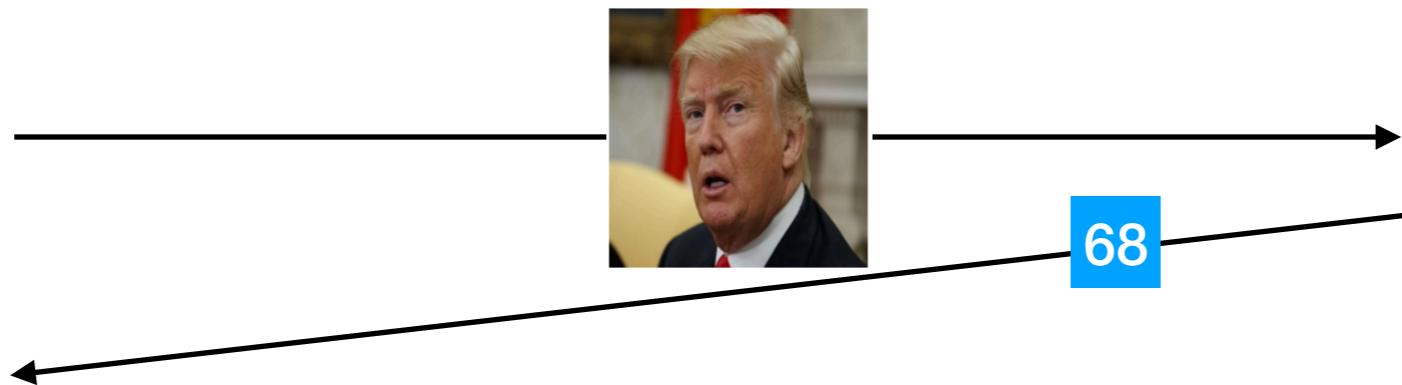


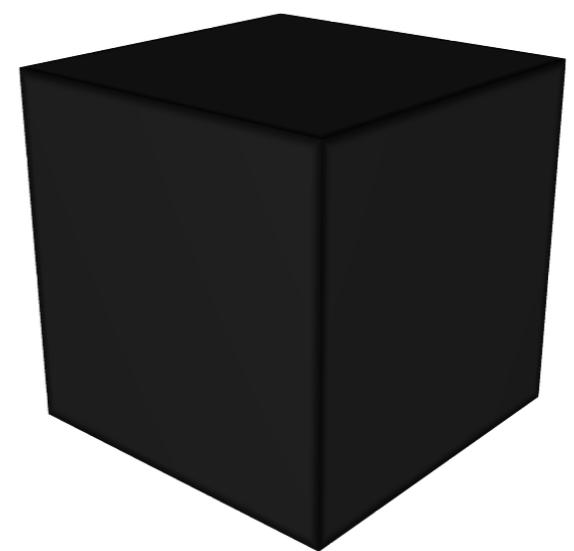
Confidence graph

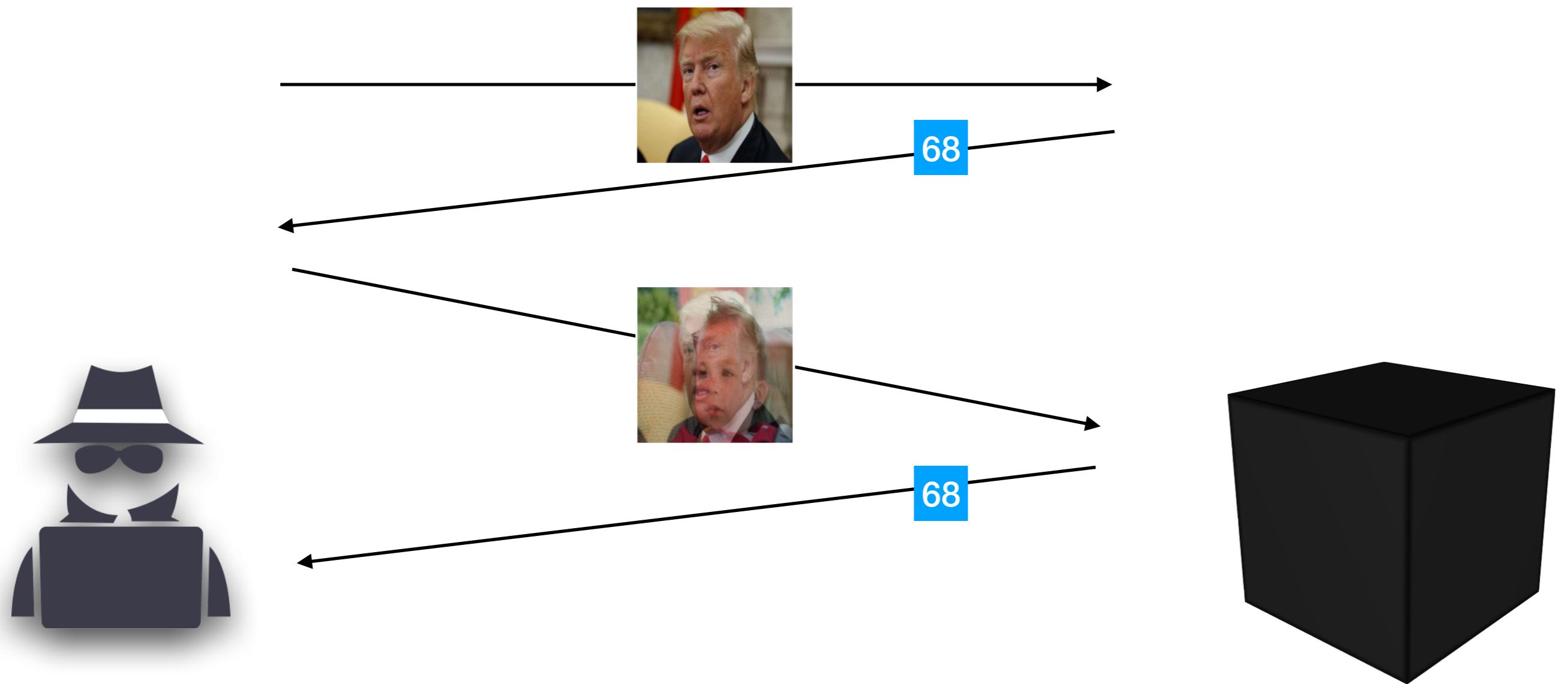


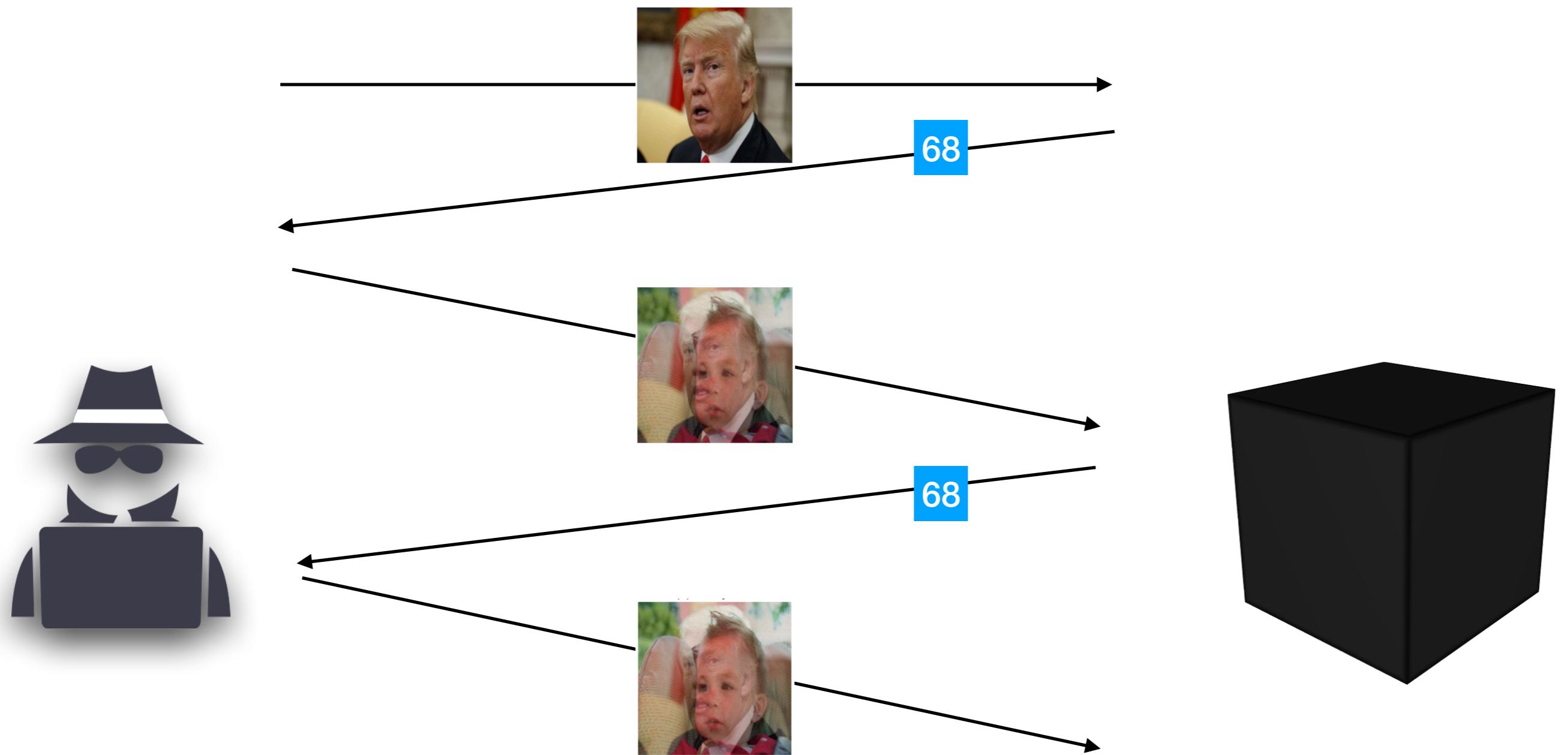


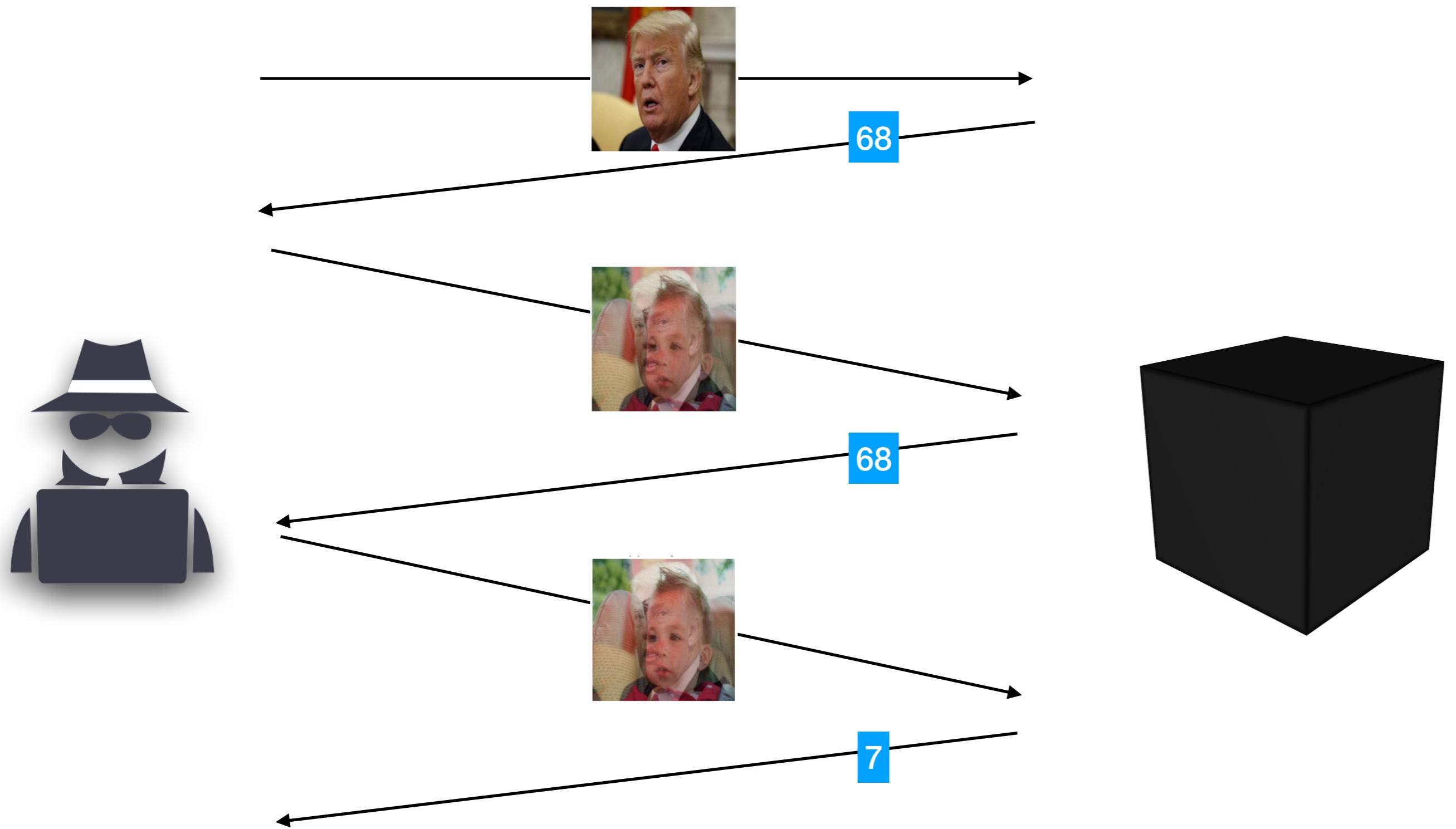


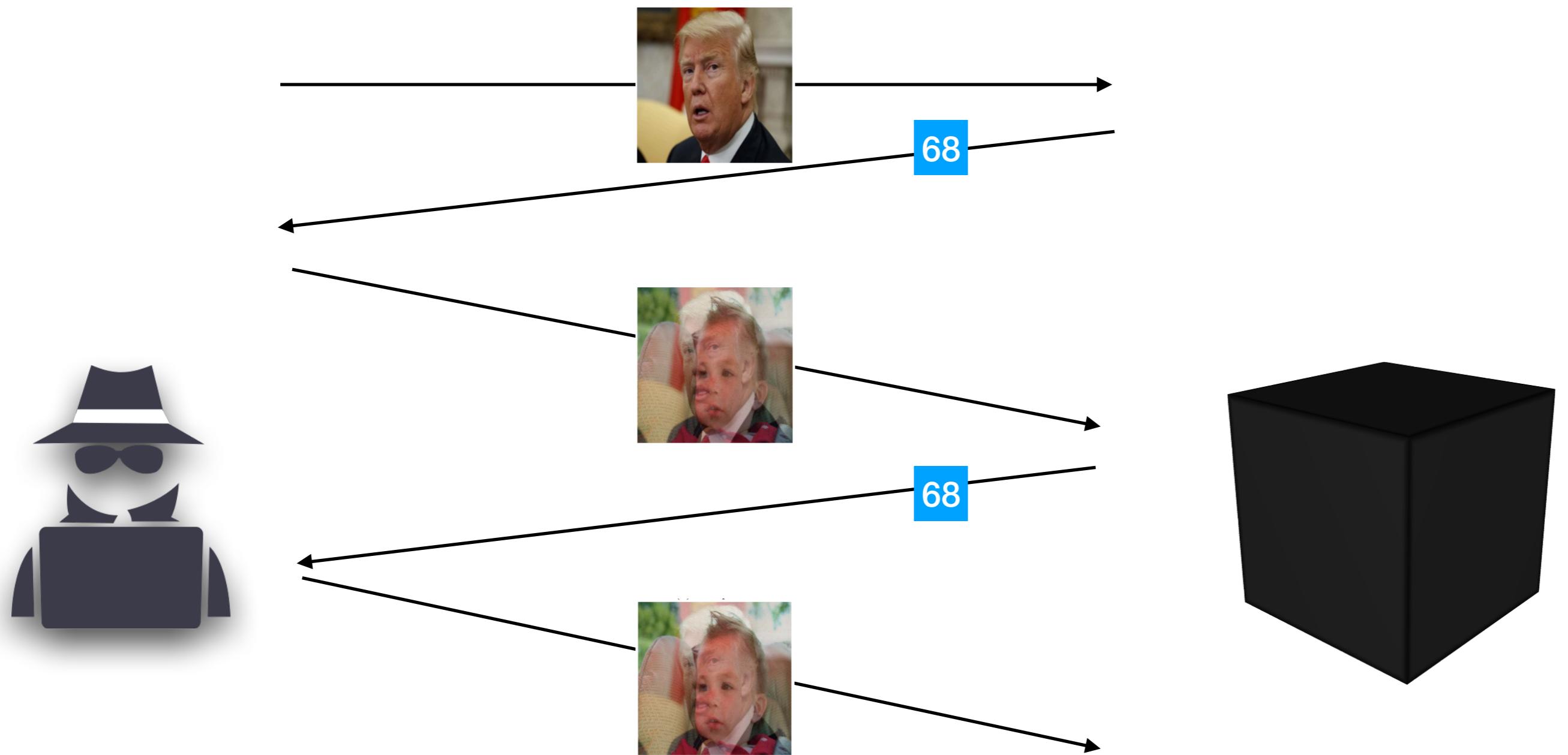


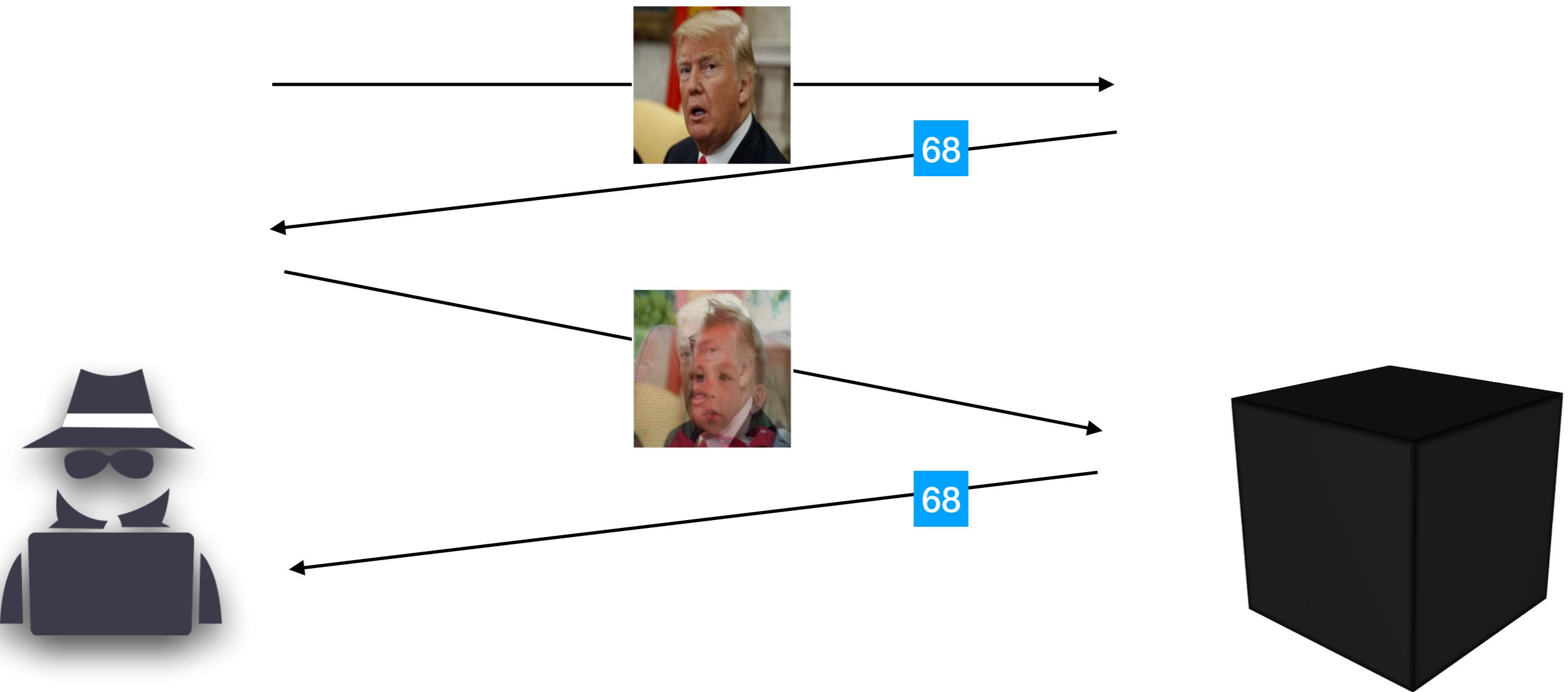


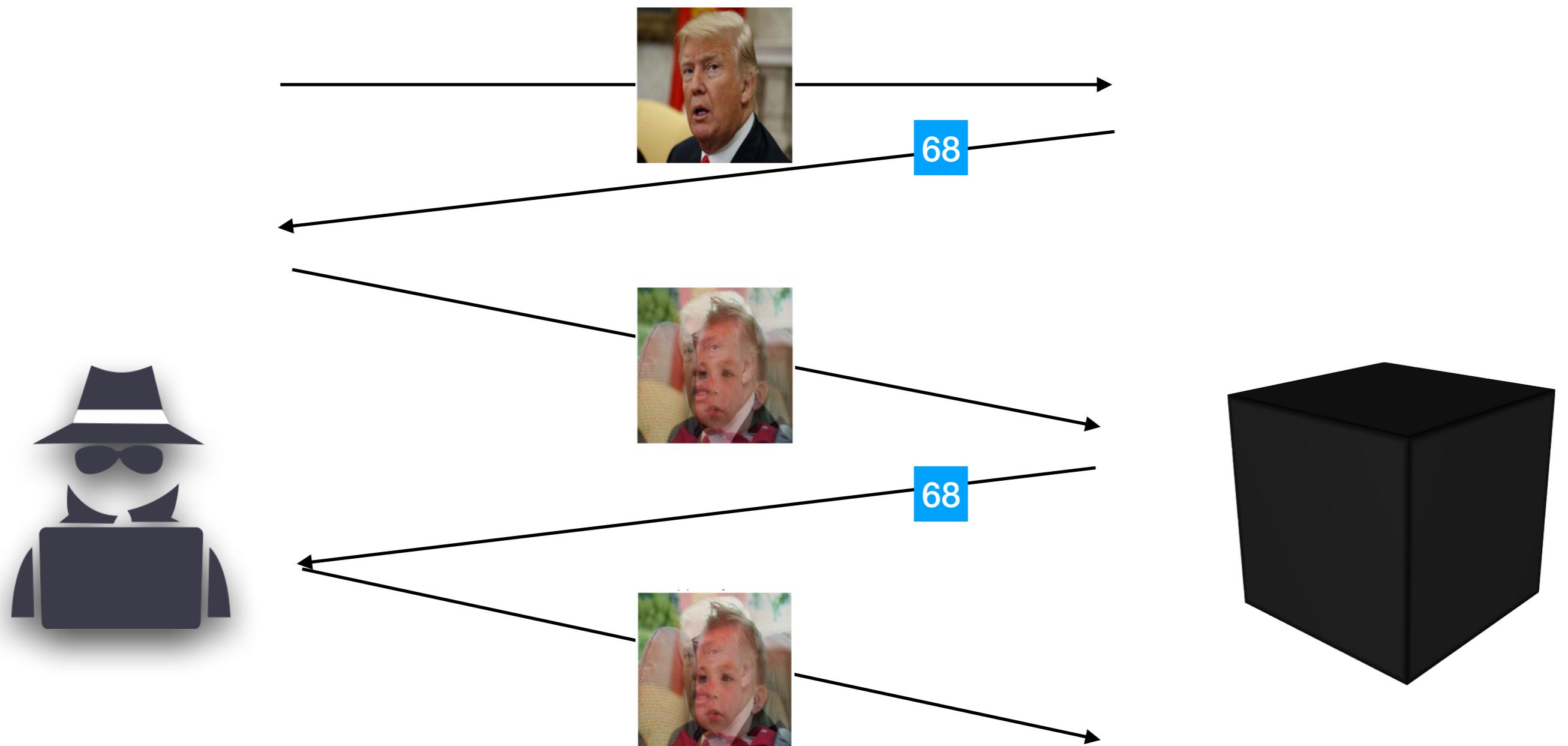


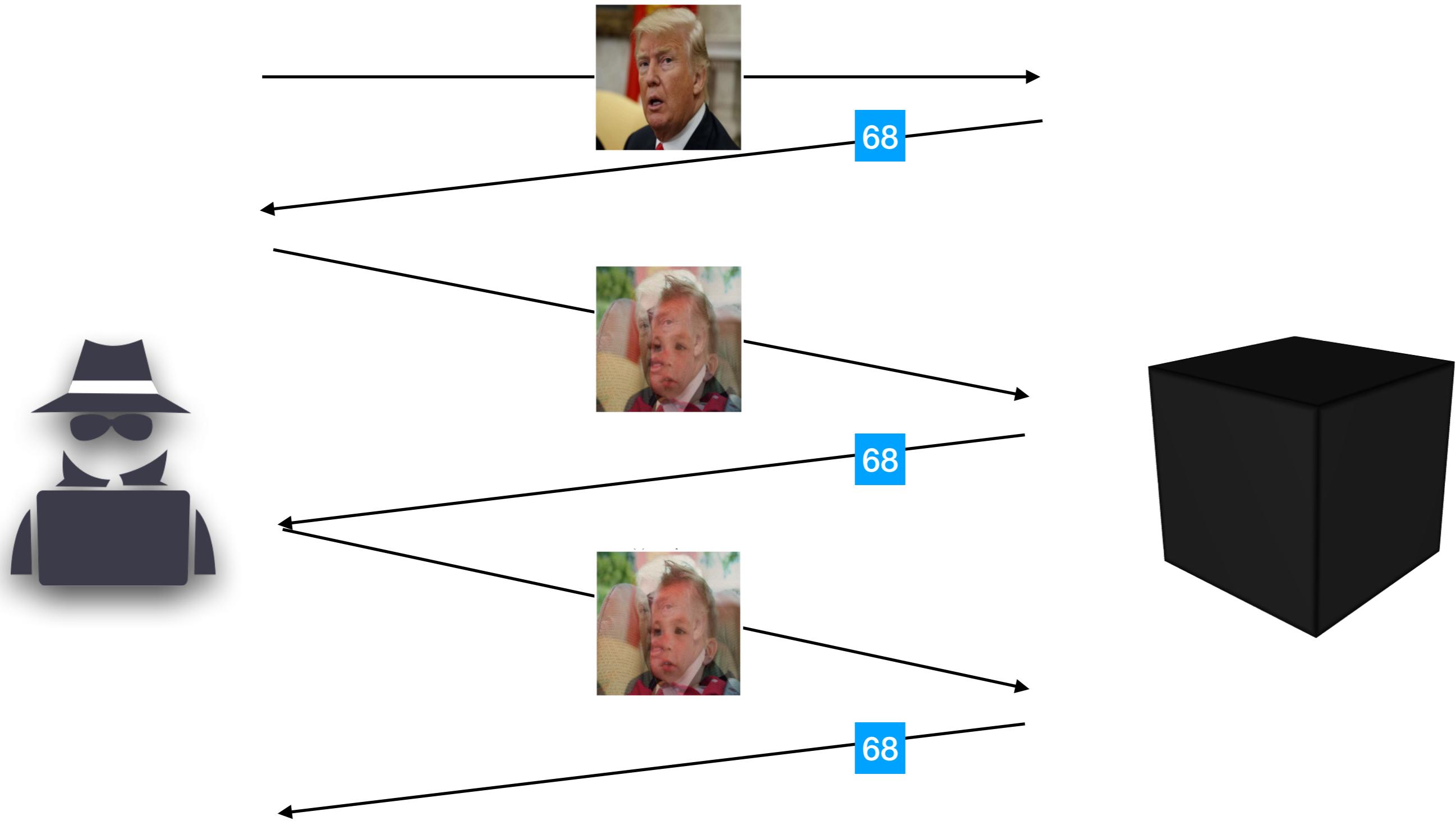


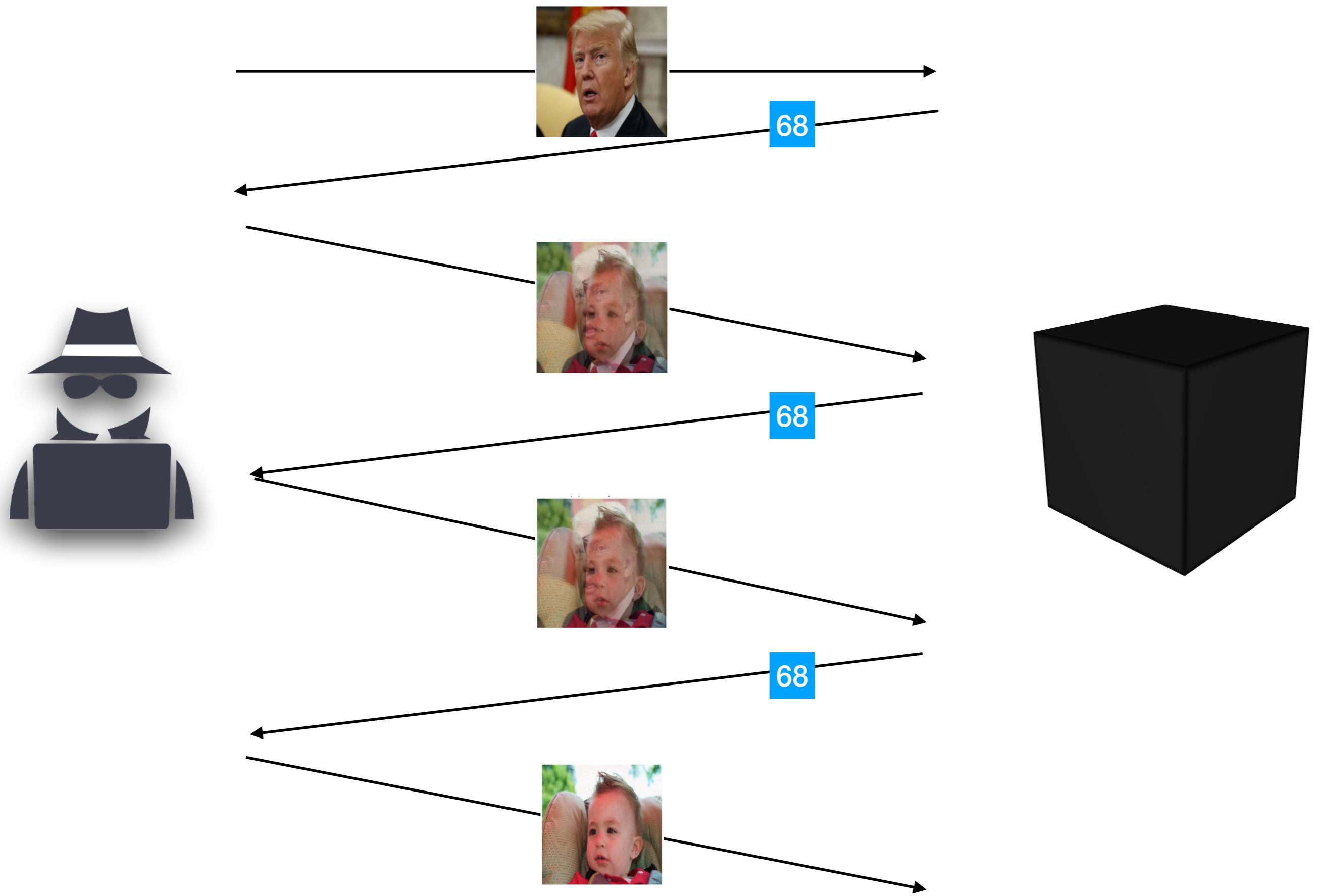


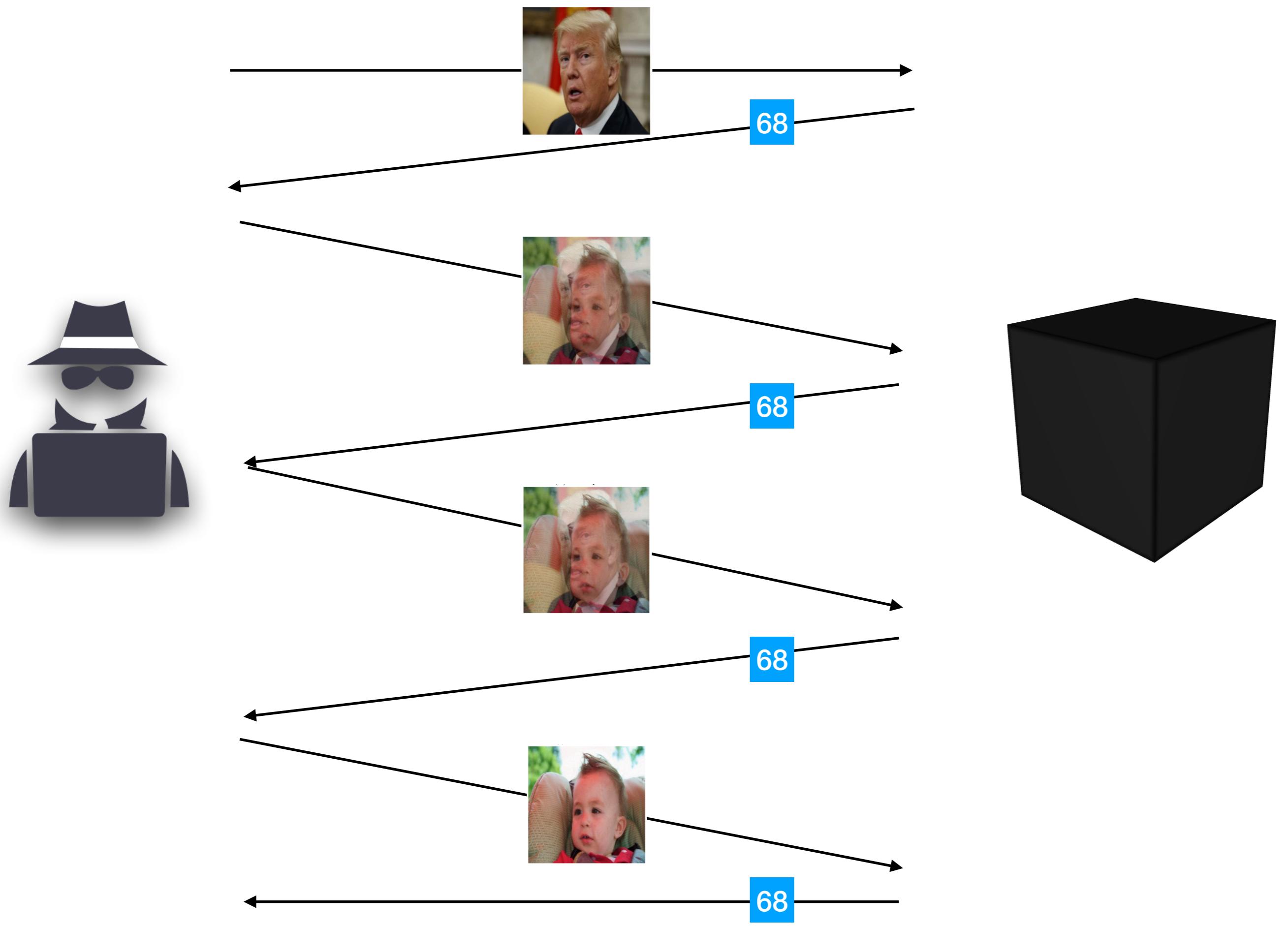














(a) The starting image



(c) 2000 queries



(e) 8000 queries



(g) 16000 queries



(b) 1000 queries



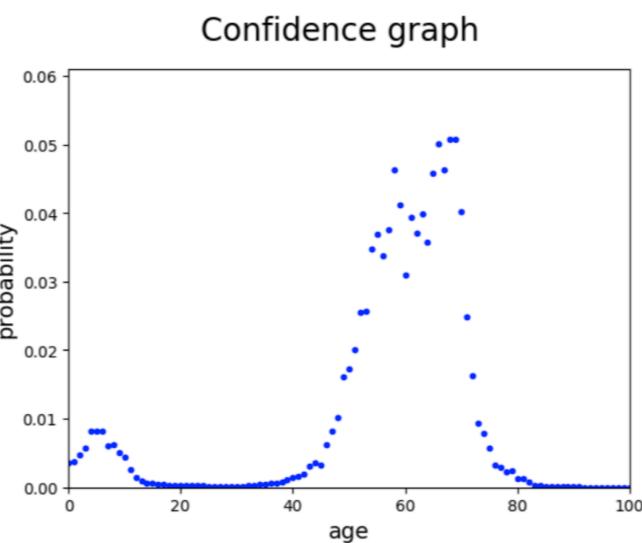
(d) 4000 queries



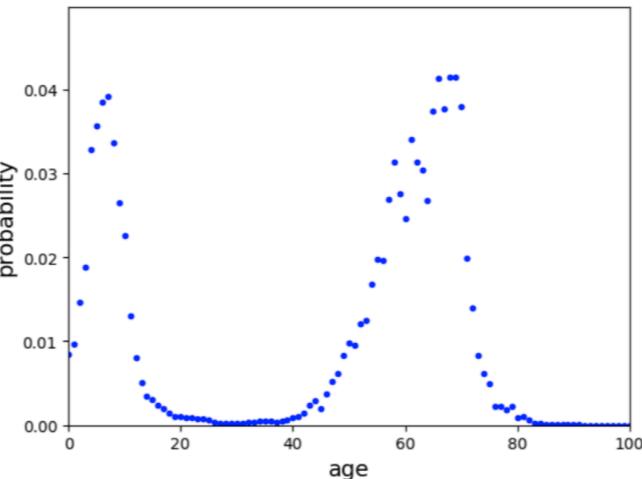
(f) 12000 queries



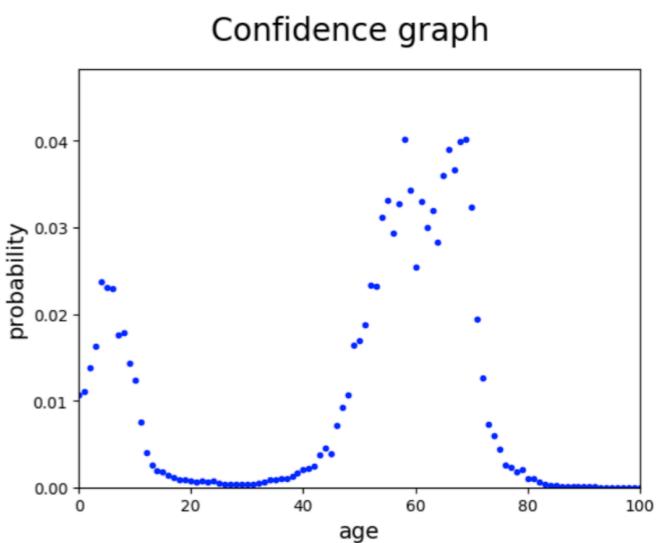
(h) Final adversarial sample



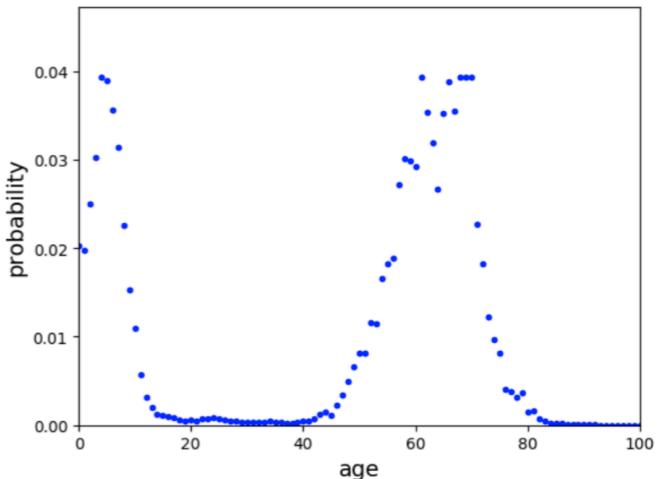
(a) 3000 queries
Confidence graph



(c) 8000 queries

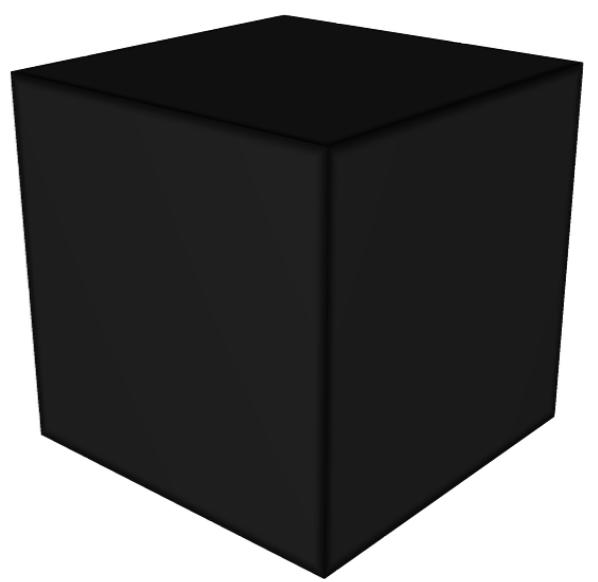


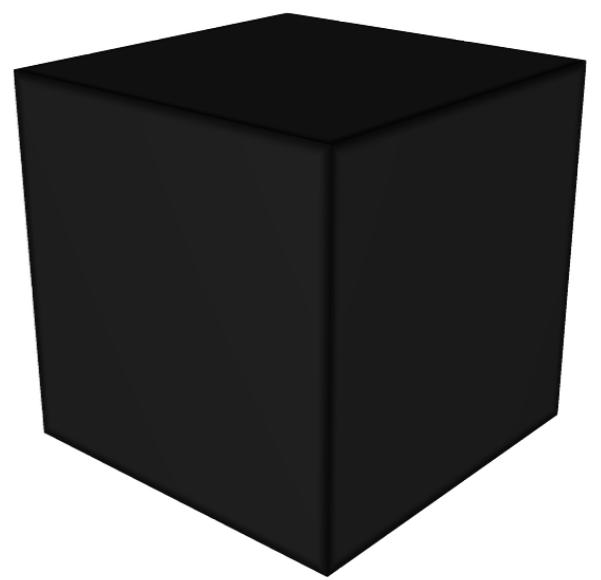
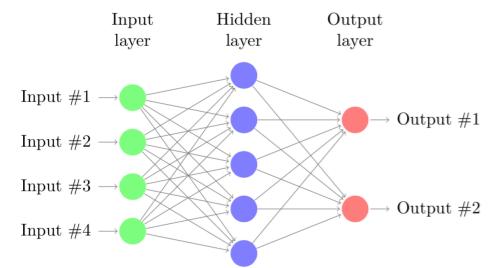
(b) 4000 queries
Confidence graph

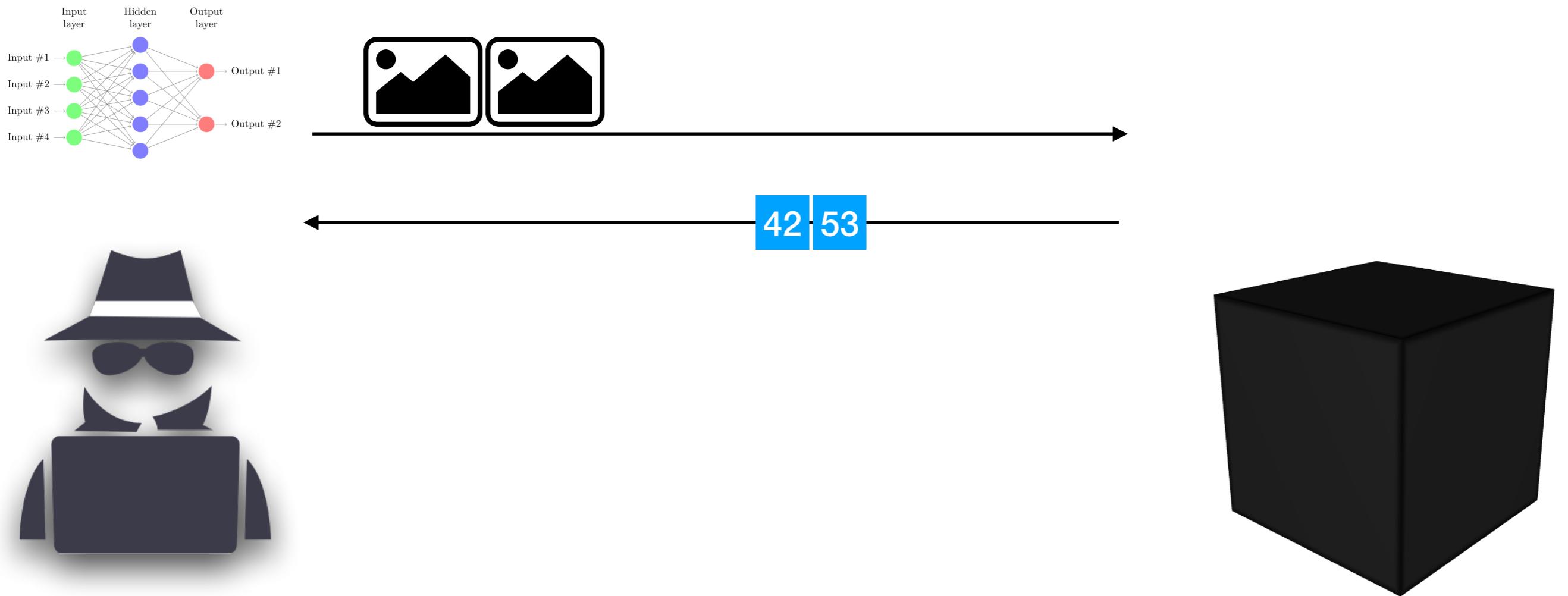


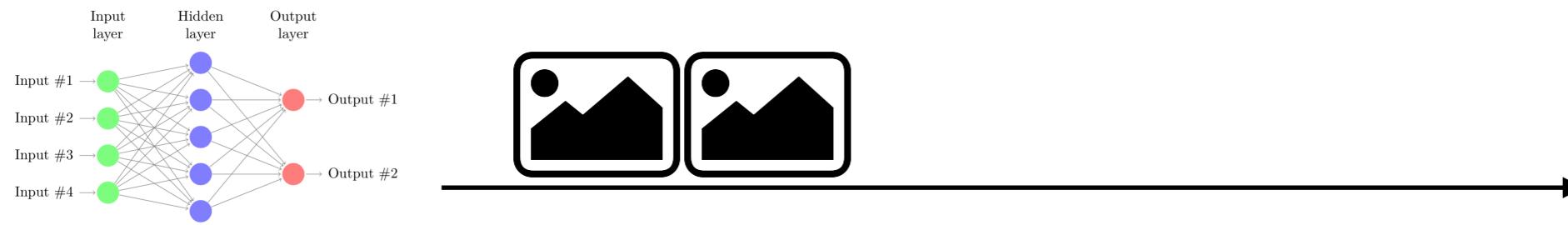
(d) The final adversarial sample (77 000 queries)

Transfer Based Approach





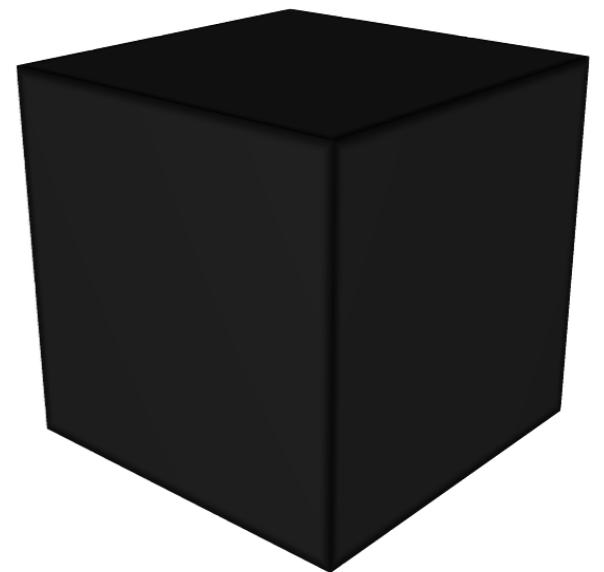


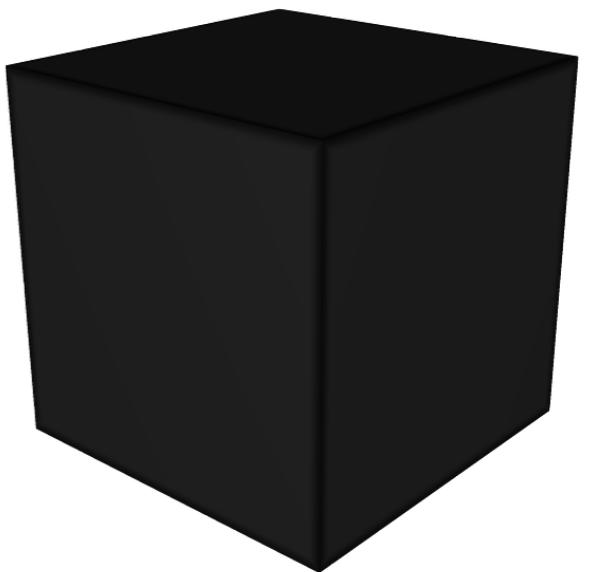
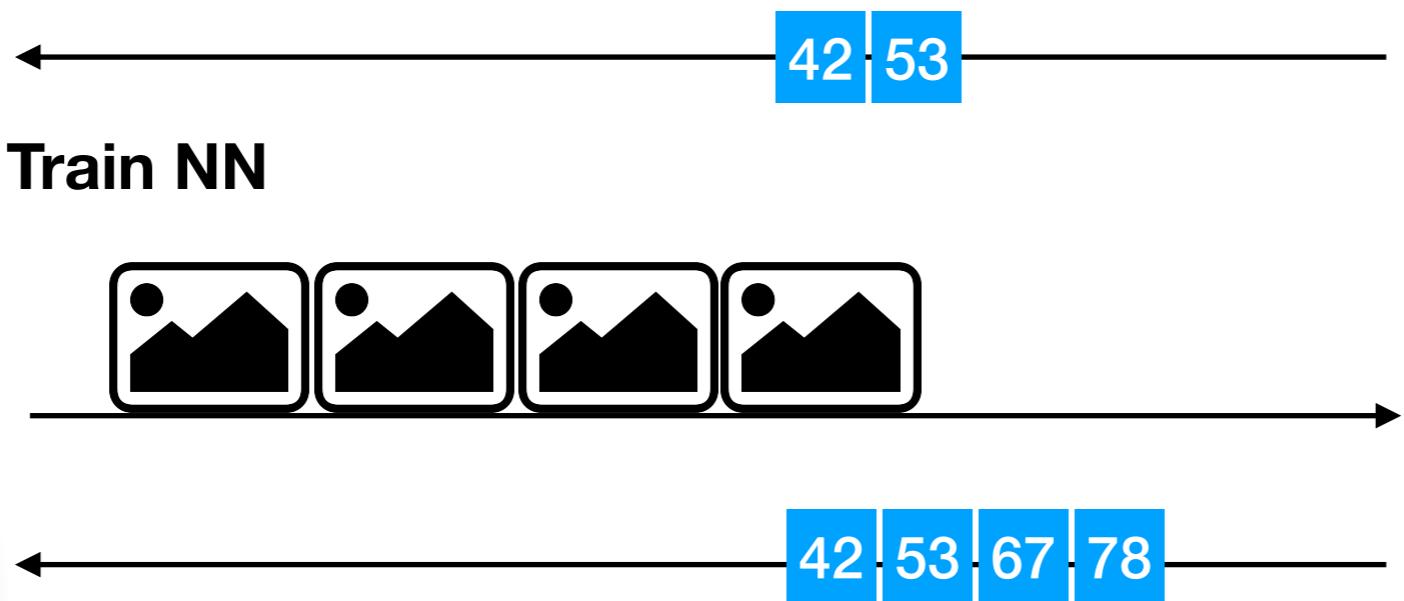
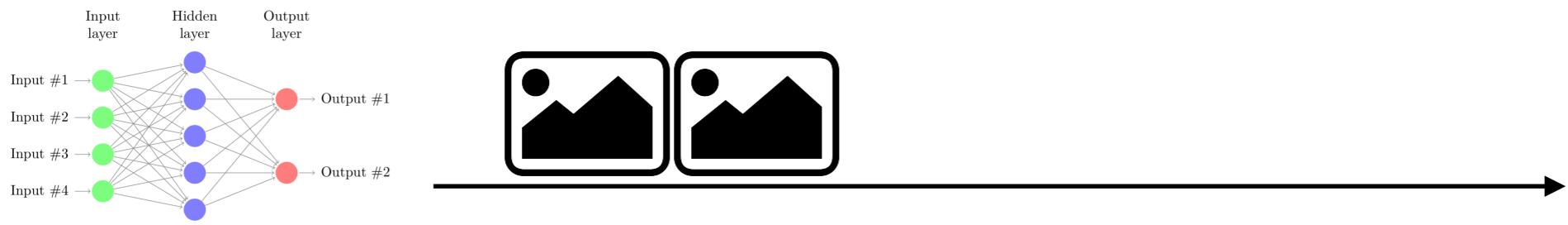


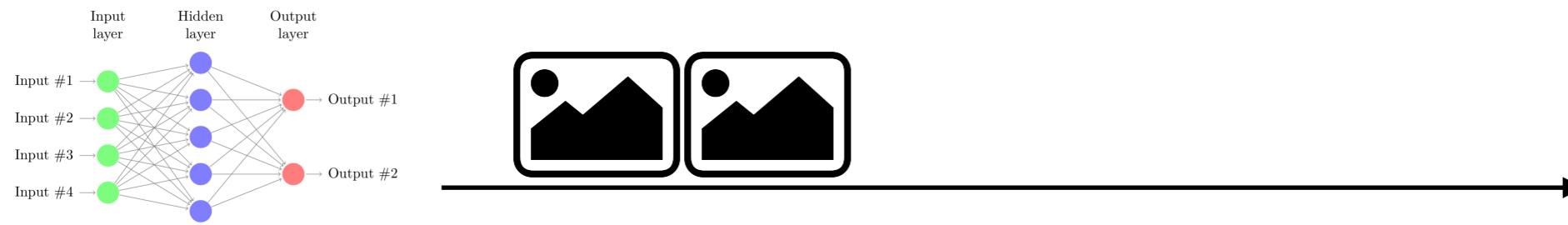
Train NN



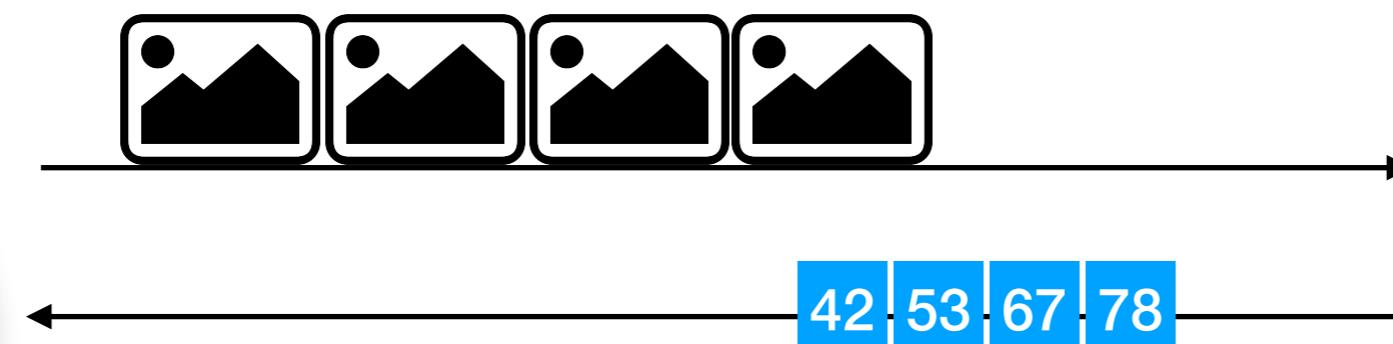
42 | 53



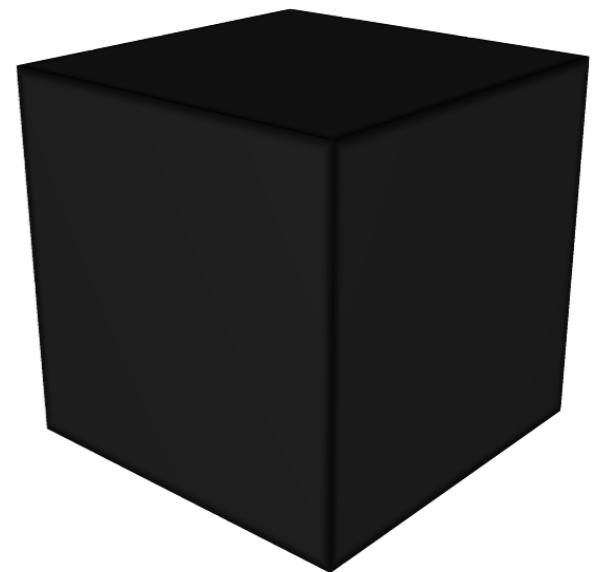


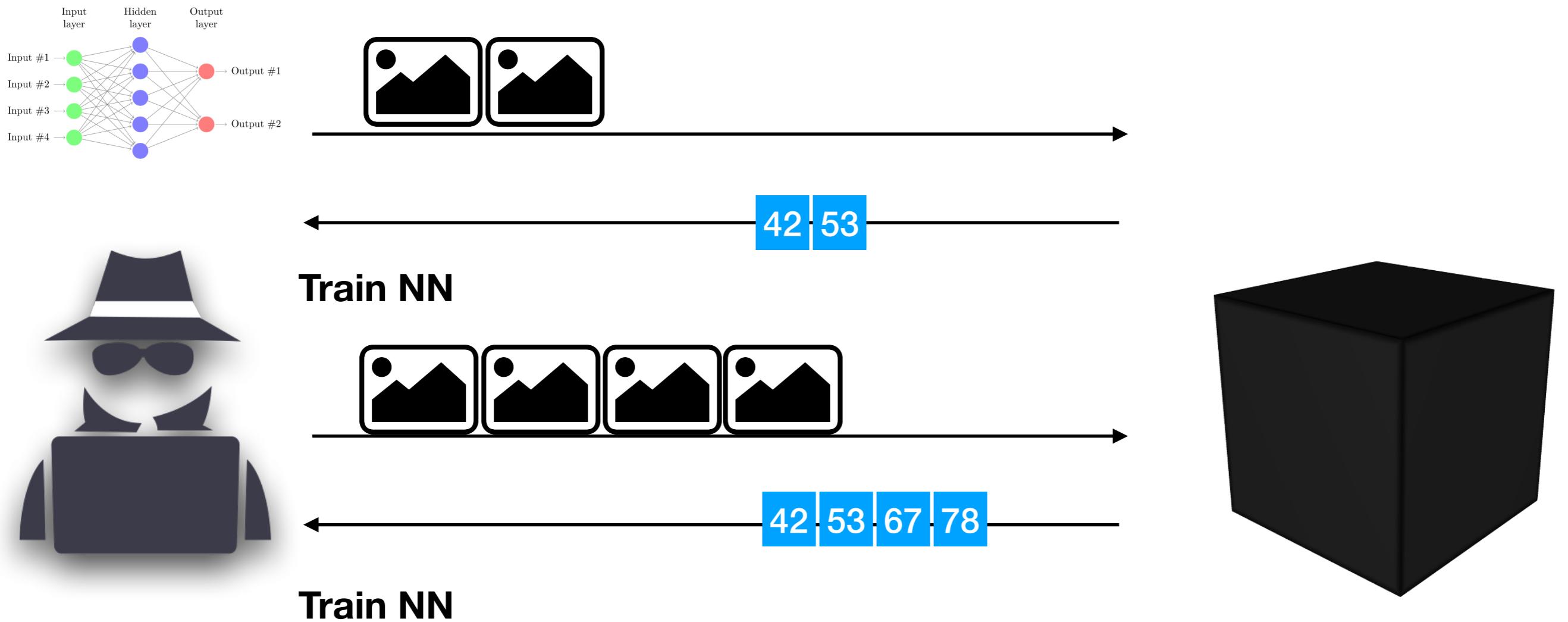


Train NN

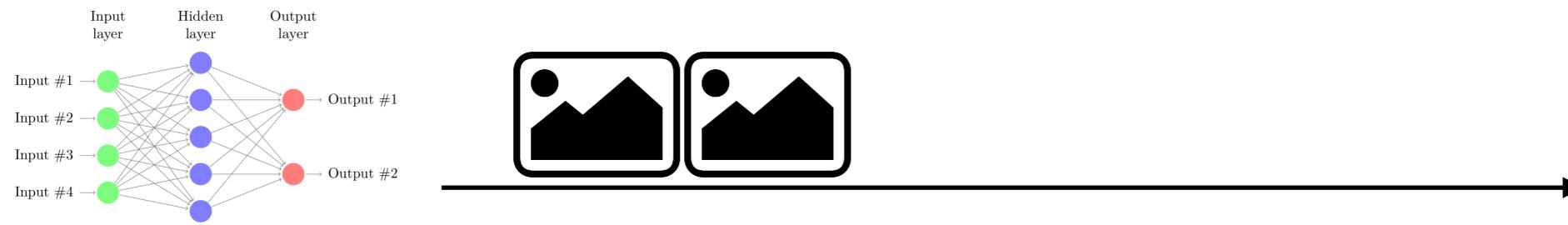


Train NN

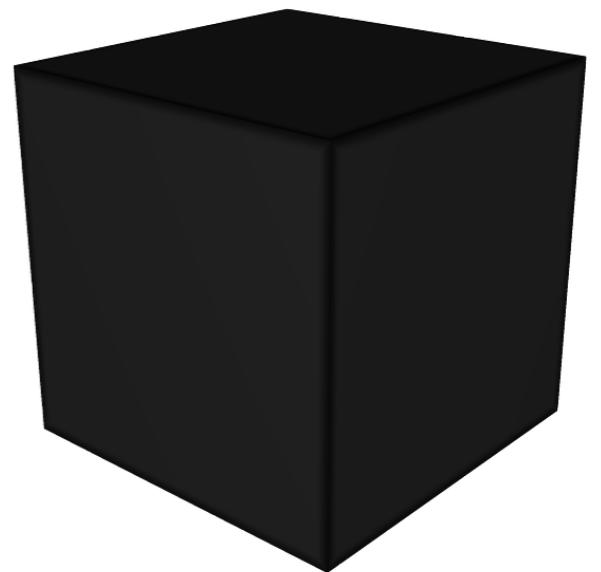
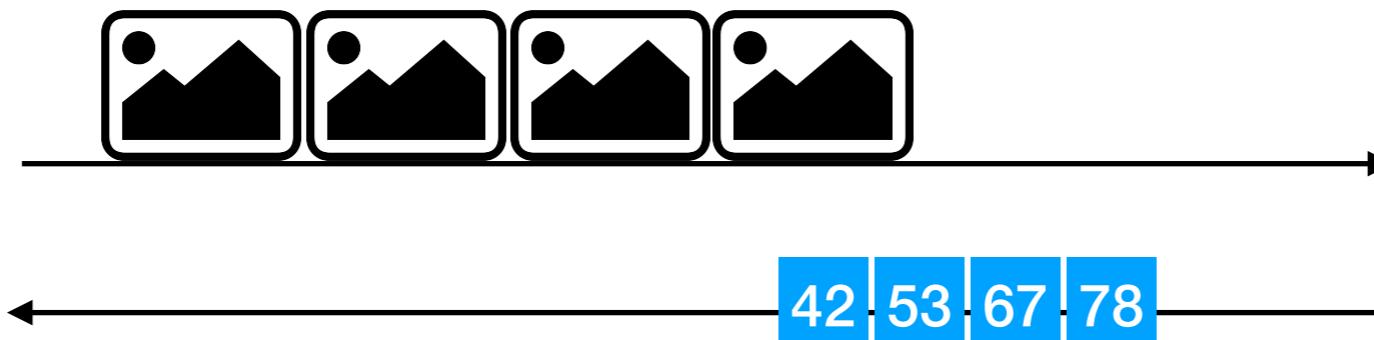




Attack white-box NN: FGSM or CW



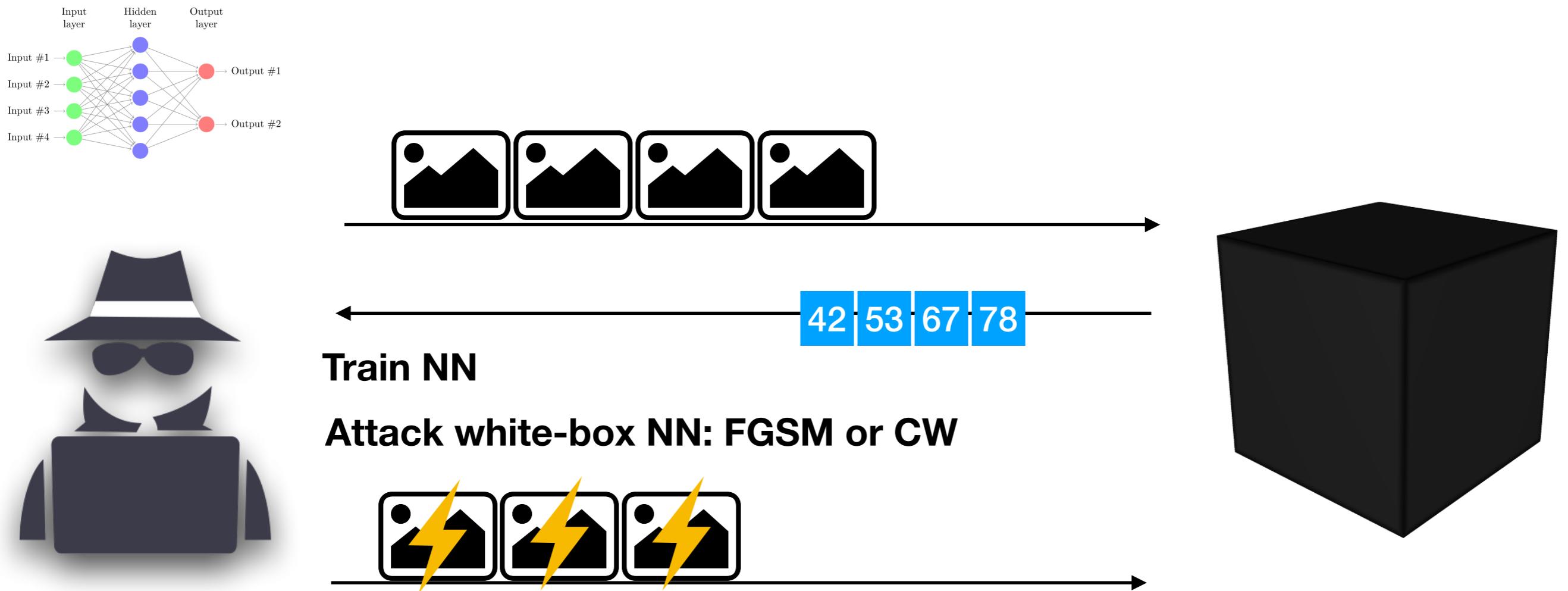
Train NN

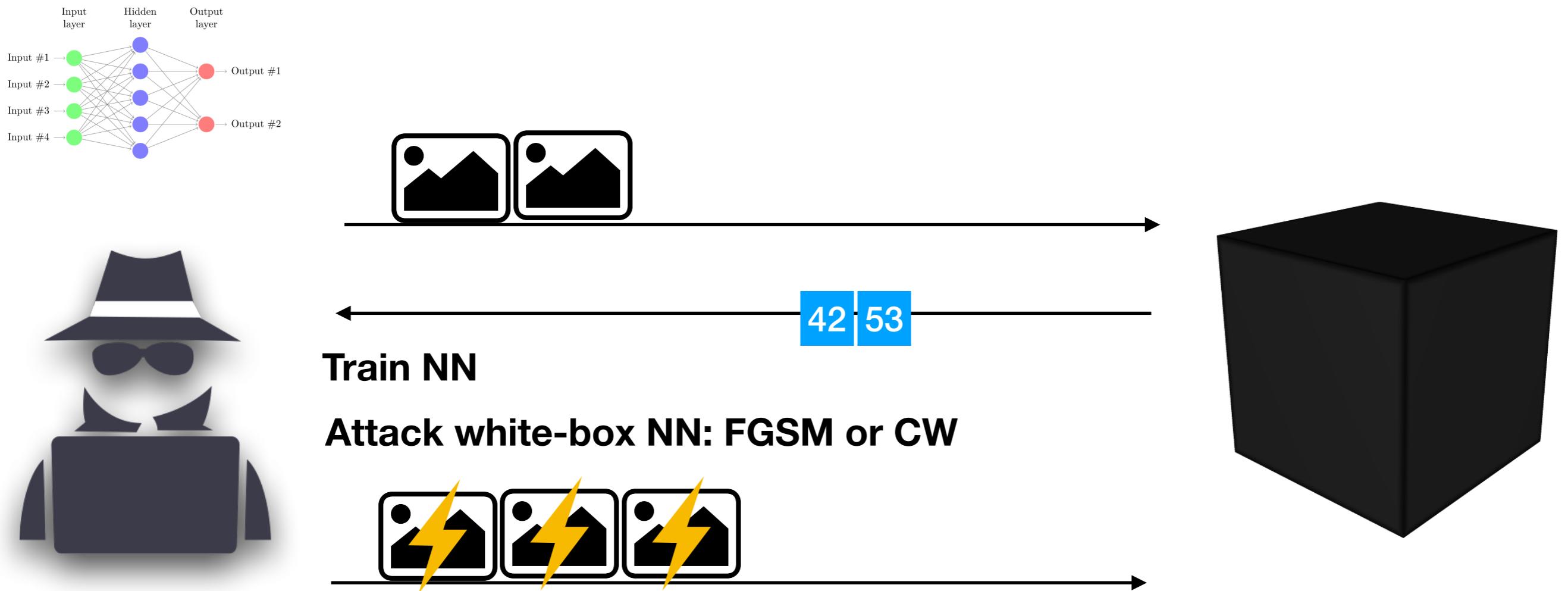


Train NN

Attack white-box NN: FGSM or CW







1. High memory expectation for the dataset expansion
2. Poor results without the dataset expansion

1. High memory expectation for the dataset expansion
2. Poor results without the dataset expansion
3. Misclassification transfers better than targeted misclassification*

* Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song
Delving into transferable adversarial examples and black-box attacks
CoRR, abs/1611.02770, 2016.

1. High memory expectation for the dataset expansion
2. Poor results without the dataset expansion
3. Misclassification transfers better than targeted misclassification*

The Semi-targeted Approach

* Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song
Delving into transferable adversarial examples and black-box attacks
CoRR, abs/1611.02770, 2016.

1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

.....

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

Classical transfer based approach

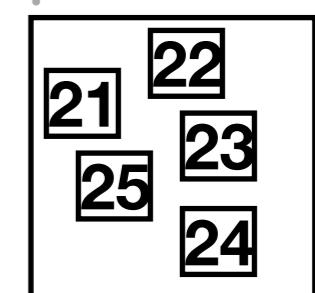
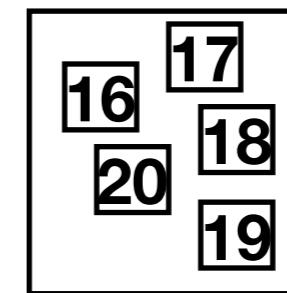
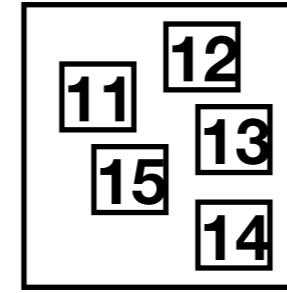
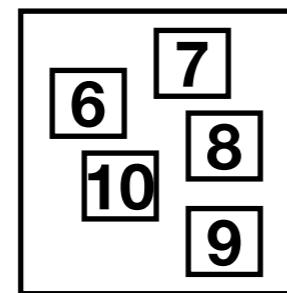
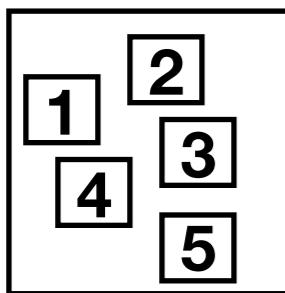
1 **2** **3** **4** **5** **6** **7** **8** **9** **10** **11** **12** **13** **14** **15** **16** **17** **18** **19** **20** **21** **22** **23** **24** **25**

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

Classical transfer based approach

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25

The Semi-targeted approach



The semi-targeted approach: results

- 4 "buckets": 0-25, 26-50, 51-75, 76-100
- FGSM
- **Managed to execute 3 iterations of dataset expansion**

	clean samples MAE	transfer based approach MAE	semi-targeted approach MAE
model 1	6.08	7.57	6.00
model 2	6.99	10.66	9.25
model 3	5.29	5.57	7.75
model 4	3.75	4.25	4.83

Contributions

- A framework for the white-box and the black-box attacks is developed
- Adversarial algorithms are evaluated in different environments
- The new semi-targeted black-box attack is introduced

Contributions

- A framework for the white-box and the black-box attacks is developed
- Adversarial algorithms are evaluated in different environments
- The new semi-targeted black-box attack is introduced

Thank you!

Appendix A

CW

Carlini and Wagner (CW)

Carlini and Wagner (CW)

1. *Minimize* added perturbation
2. Must be adversarial (target class)
3. Must be a valid image

Carlini and Wagner (CW)

1. *Minimize* added perturbation

$$\begin{aligned} & \text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta}) \\ & \text{such that } \mathcal{C}(\mathbf{x} + \boldsymbol{\delta}) = t \\ & \mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n \end{aligned}$$

2. Must be adversarial (target class)

3. Must be a valid image

Carlini and Wagner (CW)

1. *Minimize* added perturbation

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } \mathcal{C}(\mathbf{x} + \boldsymbol{\delta}) = t$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

2. Must be adversarial (target class)

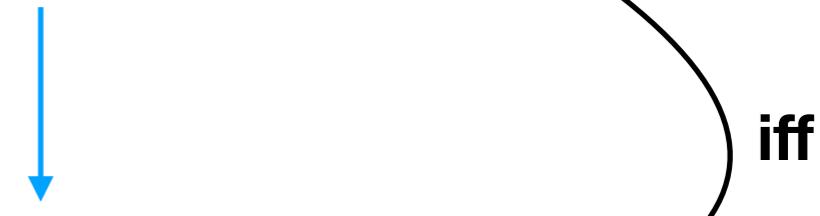
$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } f(\mathbf{x} + \boldsymbol{\delta}) \leq 0$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

3. Must be a valid image

iff



Carlini and Wagner (CW)

1. *Minimize added perturbation*

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } \mathcal{C}(\mathbf{x} + \boldsymbol{\delta}) = t$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

2. Must be adversarial (target class)

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } f(\mathbf{x} + \boldsymbol{\delta}) \leq 0$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

3. Must be a valid image

$$c > 0$$

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta}) + c \cdot f(\mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } \mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

iff

Carlini and Wagner (CW)

1. *Minimize added perturbation*

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } \mathcal{C}(\mathbf{x} + \boldsymbol{\delta}) = t$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

2. Must be adversarial (target class)

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } f(\mathbf{x} + \boldsymbol{\delta}) \leq 0$$

$$\mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

3. Must be a valid image

$$c > 0$$

$$\text{minimize } \mathcal{D}(\mathbf{x}, \mathbf{x} + \boldsymbol{\delta}) + c \cdot f(\mathbf{x} + \boldsymbol{\delta})$$

$$\text{such that } \mathbf{x} + \boldsymbol{\delta} \in [0, 1]^n$$

$$f(\mathbf{x} + \boldsymbol{\delta}) = \max(\max(\{\mathbf{Z}(\mathbf{x} + \boldsymbol{\delta})_i : i \neq t\}) - \mathbf{Z}(\mathbf{x} + \boldsymbol{\delta})_t, 0)$$

Appendix B

dataset augmentation

$$\nabla \mathbf{F}(\mathbf{X}) = \left[\frac{\partial \mathbf{F}(\mathbf{X})}{\partial x_1}, \frac{\partial \mathbf{F}(\mathbf{X})}{\partial x_2} \right]$$

$$S_{\rho+1} = \{\vec{x} + \lambda \cdot \text{sgn}(J_F[\tilde{O}(\vec{x})]) : \vec{x} \in S_\rho\} \cup S_\rho$$

Appendix C

Detailed results

The semi-targeted approach: results

- 4 "buckets": 0-25, 26-50, 51-75, 76-100
- FGSM
- **Managed to execute 3 iterations of dataset expansion**

blackbox model id	clean samples		adversarial samples					
	blackbox MAE	substitute accuracy	blackbox MAE	substitute accuracy	avg L2	std dev L2	min L2	max L2
1	6.08	0.67	6.00	0.37	2525.36	76.90	2102.57	2589.41
2	6.99	0.68	9.25	0.30	2525.20	76.96	2103.83	2589.41
3	5.29	0.69	7.75	0.35	2525.43	76.61	2108.77	2589.41
4	3.75	0.74	4.83	0.34	2525.48	76.57	2105.55	2589.41

Classical transfer based approach:

blackbox model id	clean samples		adversarial samples					
	blackbox MAE	substitute MAE	blackbox MAE	substitute MAE	avg L2	std dev L2	min L2	max L2
1	6.08	10.20	7.57	13.63	2524.84	77.62	2097.67	2589.41
2	6.99	11.09	10.66	18.59	2525.10	76.92	2105.15	2589.41
3	5.29	21.02	5.57	25.09	2527.32	73.56	2143.89	2589.41
4	3.75	8.25	4.25	14.75	2525.40	76.66	2111.09	2589.41

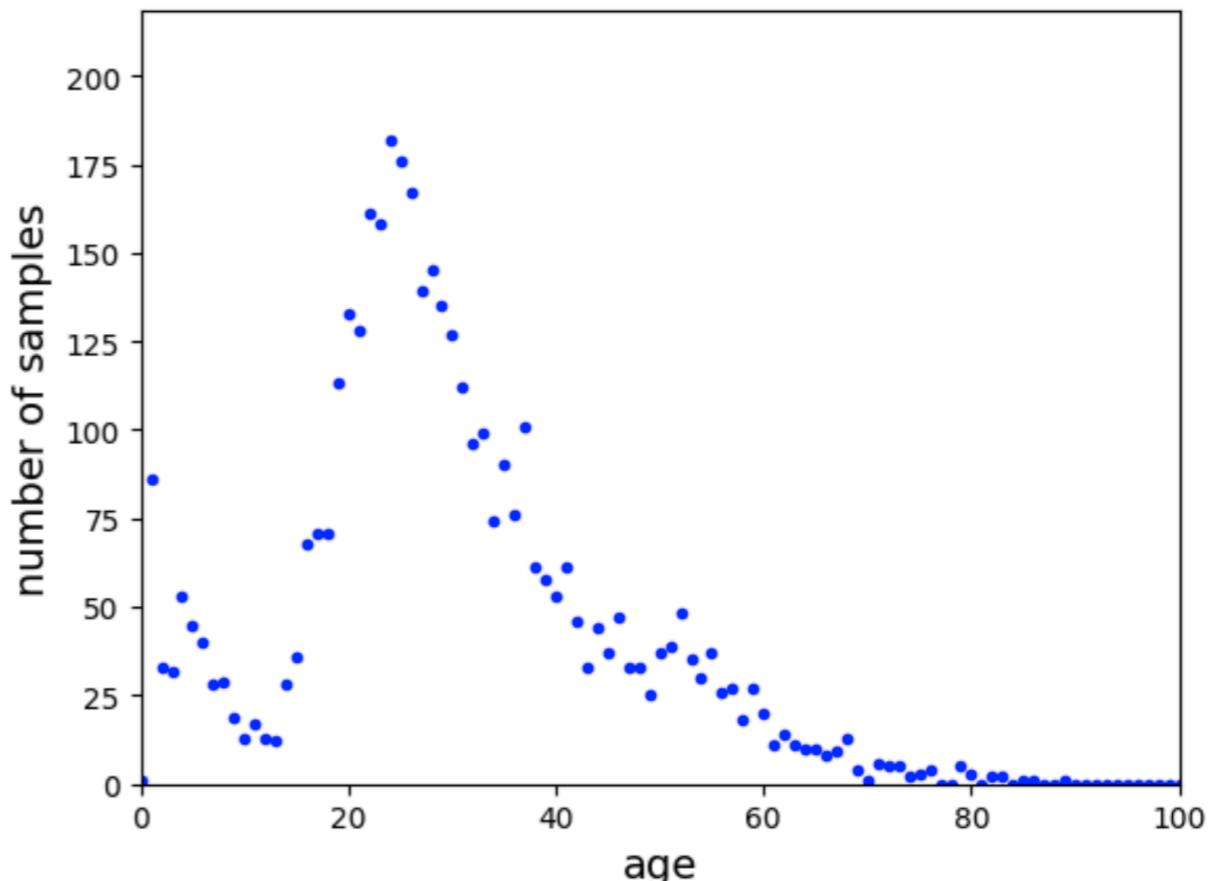
Appendix D

Dataset

The APPA REAL Training Dataset

The APPA-REAL dataset:

- total 7,591 images
 - 4113 training images
 - 1500 validation images
 - 1978 test images



The UTK Face Dataset

The UTK-Face dataset:

- total 23,252 images

