

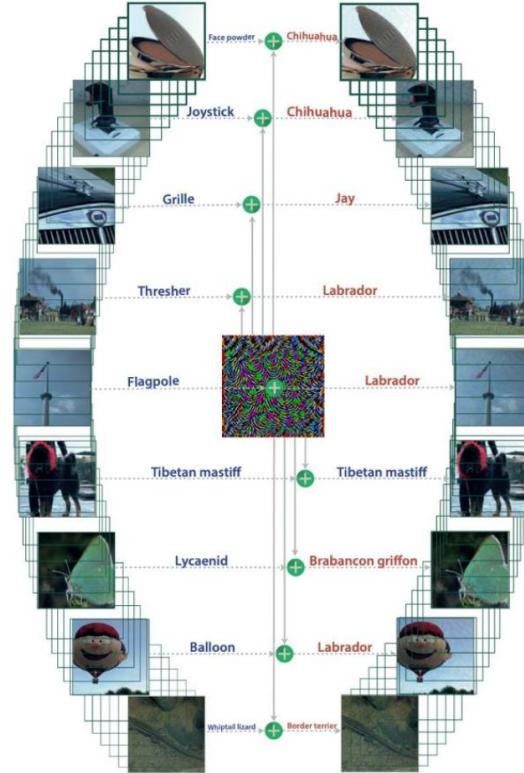
Universal Adversarial Perturbations - Lab

(Machine Learning Security - Fall 2025)

Dorjan Hitaj
Fabio De Gaspari

Universal adversarial perturbations

There exists a universal and small perturbation vector that causes images to be misclassified with high probability

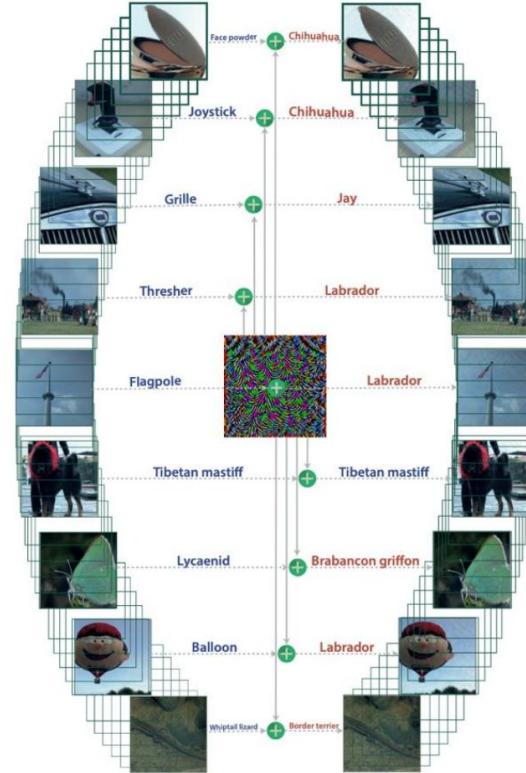


Universal adversarial perturbations

There exists a universal and small perturbation vector that causes images to be misclassified with high probability.

In principle, same as before:

- Optimise the perturbation over different images until you get one that works on all intended classes



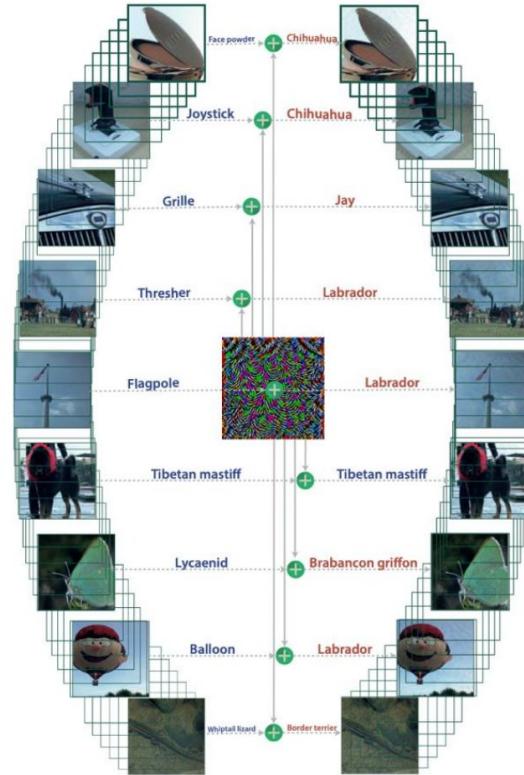
Universal adversarial perturbations

Algorithm 1 Computation of universal perturbations.

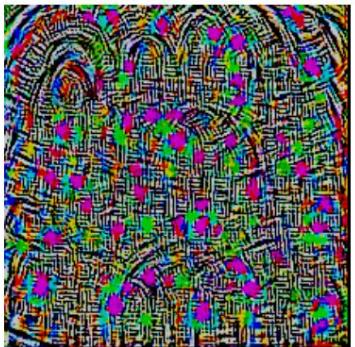
```
1: input: Data points  $X$ , classifier  $\hat{k}$ , desired  $\ell_p$  norm of  
the perturbation  $\xi$ , desired accuracy on perturbed sam-  
ples  $\delta$ .  
2: output: Universal perturbation vector  $v$ .  
3: Initialize  $v \leftarrow 0$ .  
4: while Err( $X_v$ )  $\leq 1 - \delta$  do  
5:   for each datapoint  $x_i \in X$  do  
6:     if  $\hat{k}(x_i + v) = \hat{k}(x_i)$  then  
7:       Compute the minimal perturbation that  
sends  $x_i + v$  to the decision boundary:  

$$\Delta v_i \leftarrow \arg \min_r \|r\|_2 \text{ s.t. } \hat{k}(x_i + v + r) \neq \hat{k}(x_i).$$
  
8:       Update the perturbation:  

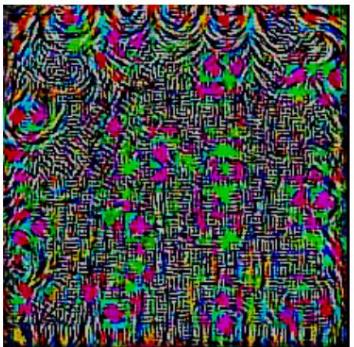
$$v \leftarrow \mathcal{P}_{p,\xi}(v + \Delta v_i).$$
  
9:     end if  
10:    end for  
11:  end while
```



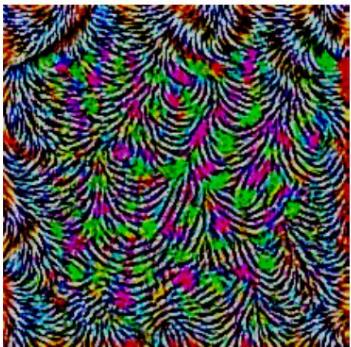
Universal adversarial perturbations



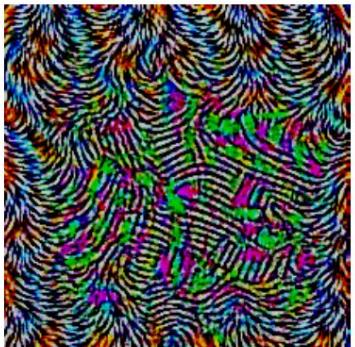
(a) CaffeNet



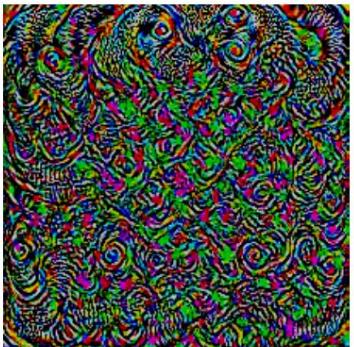
(b) VGG-F



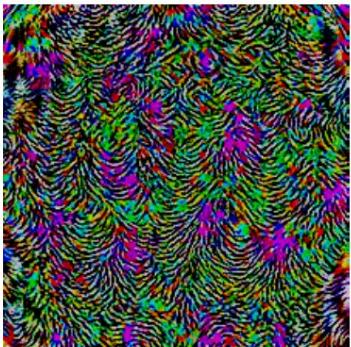
(c) VGG-16



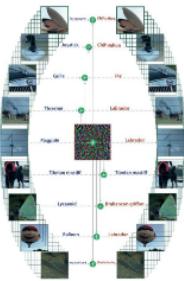
(d) VGG-19



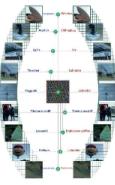
(e) GoogLeNet



(f) ResNet-152



Universal adversarial perturbations



- A bit too slow and sometimes takes time to/does not converge

Algorithm 1 Computation of universal perturbations.

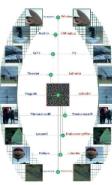
- 1: **input:** Data points X , classifier \hat{k} , desired ℓ_p norm of the perturbation ξ , desired accuracy on perturbed samples δ .
- 2: **output:** Universal perturbation vector v .
- 3: Initialize $v \leftarrow 0$.
- 4: **while** $\text{Err}(X_v) \leq 1 - \delta$ **do**
- 5: **for** each datapoint $x_i \in X$ **do**
- 6: **if** $\hat{k}(x_i + v) = \hat{k}(x_i)$ **then**
- 7: Compute the *minimal* perturbation that sends $x_i + v$ to the decision boundary:

$$\Delta v_i \leftarrow \arg \min_r \|r\|_2 \text{ s.t. } \hat{k}(x_i + v + r) \neq \hat{k}(x_i).$$

- 8: Update the perturbation:

$$v \leftarrow \mathcal{P}_{p,\xi}(v + \Delta v_i).$$

- 9: **end if**
- 10: **end for**
- 11: **end while**



Universal adversarial perturbations - a faster variant

Algorithm 1 Computation of universal perturbations.

1: **input:** Data points X , classifier \hat{k} , desired ℓ_p norm of
 perturbation δ , number of epochs N_{ep} , number of
 embed samples n_e .

Algorithm 2 Stochastic gradient for universal perturbation

for epoch = 1 . . . N_{ep} **do**

for minibatch $B \subset X$ **do**

 Update δ with gradient variant $\delta \leftarrow \delta + g$

 Project δ to ℓ_p ball

end for

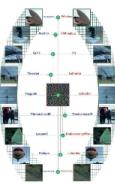
end for

9: **end if**

10: **end for**

11: **end while**

Universal adversarial perturbations - a faster variant



The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI-20)

Universal Adversarial Training

Ali Shafahi,* Mahyar Najibi,* Zheng Xu,* John Dickerson, Larry S. Davis, Tom Goldstein

Department of Computer Science

University of Maryland

College Park, Maryland 20742

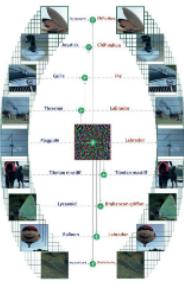
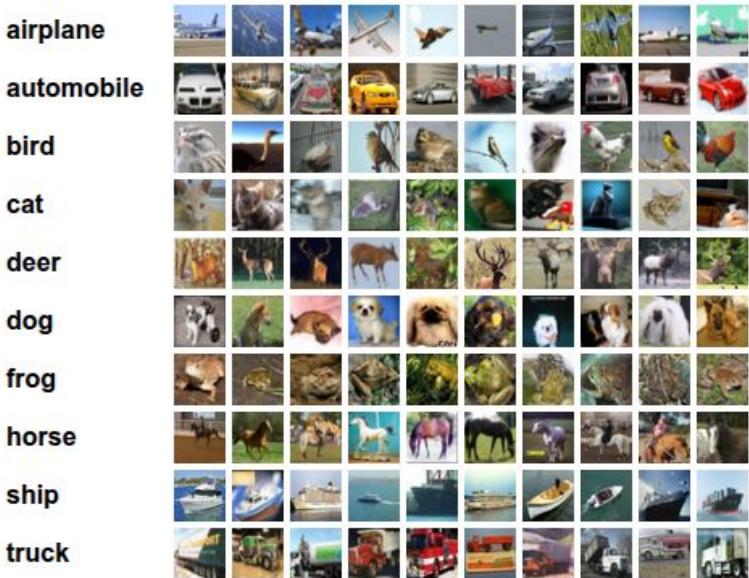
{ashafahi, najibi, xuzh, john, lsd,tomg}@cs.umd.edu

*Attached to this Lab materials

Task at hand

CIFAR 10

- The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.



Task at hand - Loading the data

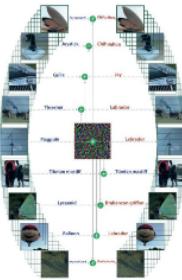
CIFAR 10

- The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.



```
trainset = torchvision.datasets.CIFAR10(root='./data', train=True, download=True, transform=transform_train)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=128, shuffle=True, num_workers=2)
```

```
testset = torchvision.datasets.CIFAR10(root='./data', train=False, download=True, transform=transform_test)
testloader = torch.utils.data.DataLoader(testset, batch_size=100, shuffle=False, num_workers=2)
```



Task at hand - Loading the data

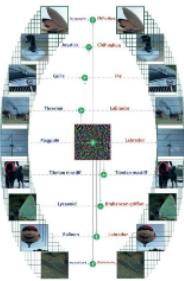
CIFAR 10

- The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.



```
trainset = torchvision.datasets.CIFAR10(root='./data', train=True, download=True, transform=transform_train)
trainloader = torch.utils.data.DataLoader(trainset, batch_size=128, shuffle=True, num_workers=2)
```

```
testset = torchvision.datasets.CIFAR10(root='./data', train=False, download=True, transform=transform_test)
testloader = torch.utils.data.DataLoader(testset, batch_size=100, shuffle=False, num_workers=2)
```



The overall code structure

```
✓ checkpoint
    ckpt.pth
✓ data
    cifar-10-batches-py
    cifar10
    cifar-10-python.tar.gz
✓ models
    __init__.py
    densenet.py
    resnet.py
    vgg.py
✓ uaps
    uap.png
    uap_sgd.pth
    run_attack.py
    train_models.py
    uap_attack.py
    utilities.py
```

- **models/** - 3 different DNN architectures already implemented
- **data/** - the folder where cifar 10 data will be stored, automatically by the torch dataloader
- **uaps/** - the folder where you will store the different perturbations that you will generate for further use in evaluation.

The overall code structure

```
✓ checkpoint
    ? ckpt.pth
✓ data
    > cifar-10-batches-py
    > cifar10
    ? cifar-10-python.tar.gz
✓ models
    ? __init__.py
    ? densenet.py
    ? resnet.py
    ? vgg.py
✓ uaps
    ? uap.png
    ? uap_sgd.pth
    ? run_attack.py
    ? train_models.py
    ? uap_attack.py
    ? utilities.py
```

- **run_attack.py** - short snippet of code to evaluate uaps
- **train_models.py** - contains the code to train different type of models on cifar dataset
- **utilities.py** - bunch of useful functions and definitions
- **uap_attack.py** - contains a method definition that you need to fill out

What you will need to do in UAP method?

- Expand current universal perturbation to the batch.
- Compute the loss of the model on perturbed inputs (with per-sample losses clamped).
- Compute gradients of that loss w.r.t. the per-sample perturbations.
- Average the gradients across the batch, take the sign, and update the universal perturbation by a small step.
- Clip the universal perturbation to the allowed Linf ball (the bounds [-eps, eps]).

What you will need to do in UAP method?

```
def uap_sgd2( 2 usages
    model: nn.Module,
    loader: torch.utils.data.DataLoader,
    nb_epoch: int,
    eps: float,
    beta: float = 12.0,
    step_decay: float = 0.8,
    y_target: Optional[int] = None,
    loss_fn: Optional[nn.Module] = None,
    uap_init: Optional[torch.Tensor] = None,
    device: Optional[torch.device] = None,
) -> Tuple[torch.Tensor, List[float]]:
    """
    Universal Adversarial Perturbation (UAP) via sign-of-mean-grad updates.

    Args:
        model: pretrained classifier (expects inputs in [0,1]).
        loader: DataLoader yielding (images, labels).
        nb_epoch: number of epochs over the dataset.
        eps: Linf bound for the universal perturbation (e.g., 8/255).
        beta: clamp value for per-sample loss (scalar).
        step_decay: multiplicative decay factor applied per epoch: step = eps * (step_decay ** epoch).
        y_target: if provided, generate a targeted UAP towards this class (int).
        loss_fn: optional loss function (default CrossEntropyLoss(reduction='none')).
        uap_init: optional initial perturbation tensor of shape (C,H,W).
        device: torch.device or None (will use model.device if None).

    Returns:
        delta: tensor (C,H,W) -- learned universal perturbation (detached).
        losses: list of scalar loss values recorded (one per batch).
    """

```

Pseudocode

Input:

Pretrained model $f(\cdot)$

Dataset $D = \{(x_i, y_i)\}$

Epochs: N

Perturbation limit: ϵ

Step decay factor: α

Clamping constant: β

(Optional) Target label y_{target}

(Optional) Loss function $L(\cdot)$ — default: Cross-Entropy

Output:

Universal perturbation δ

Pseudocode

Set $\delta \leftarrow 0$ (or provided initialization)

Define clamped loss:

$L_{\text{clamp}}(x, y) = \text{mean}(\min(L(f(x)), y), \beta))$

For each epoch $e = 1 \dots N$:

Set step size: $\eta \leftarrow \varepsilon \times (\alpha)^e$

For each mini-batch $(x_b, y_b) \in D$:

If targeted \rightarrow replace $y_b \leftarrow y_{\text{target}}$

Compute perturbed inputs:

$x_{\text{pert}} \leftarrow \text{clip}(x_b + \delta, 0, 1)$

Forward pass: $p \leftarrow f(x_{\text{pert}})$

Compute loss:

$l \leftarrow L_{\text{clamp}}(p, y_b)$

If targeted $\rightarrow l \leftarrow -l$

Compute gradient of l w.r.t δ : $g \leftarrow \nabla_{\delta} l$

Average gradients and take sign:

$g_{\text{sign}} \leftarrow \text{sign}(\text{mean}(g))$

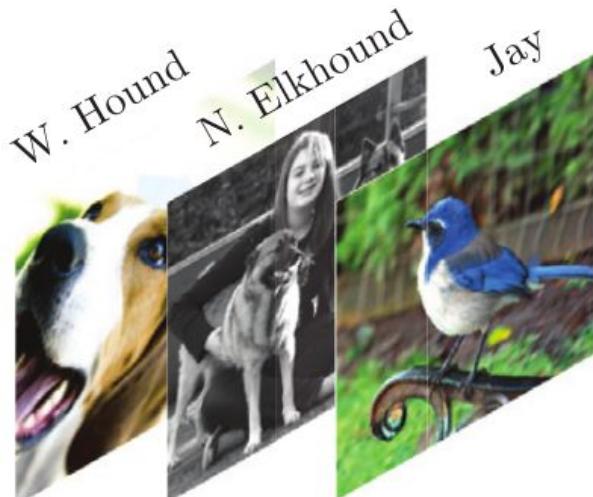
Update perturbation:

$\delta \leftarrow \delta + \eta \times g_{\text{sign}}$

Project δ to valid range:

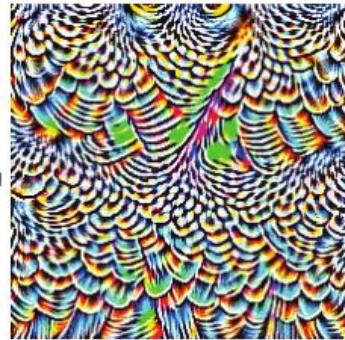
$\delta \leftarrow \text{clip}(\delta, -\varepsilon, \varepsilon)$

Where you will need to end up?



+

Universal
Perturbation



≡



```
    delta: tensor (C,H,W) -- learned universal perturbation (detached).  
    losses: list of scalar loss values recorded (one per batch).  
    ...
```

Up next...

- Explainable AI and Guaranteed robustness towards adversarial examples
[Next Lecture]
- Lab on explainable AI approaches for image classifiers and malware detectors.