

# CS772: Research Project

## Zero Shot Machine Unlearning

Ashutosh Kumar - 210221

Krish Sharma - 210530

Labajyoti Das - 210552

Shubham Patel - 210709

Siddharth Kalra - 211032

May 2, 2024

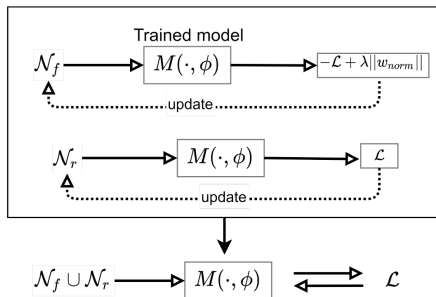
# Problem Statement

- Machine Unlearning
  - Model  $M$ , Data  $D$
  - Request:
    - Forget Data  $D_f \subset D$
    - Retain Data  $D_r = D - D_f$
  - Gold / Retrained Model:  $M^*$
  - Unlearned Model:  $M_u$
  - Aim:  $M_u(x) \approx M^*(x)$

# Problem Statement

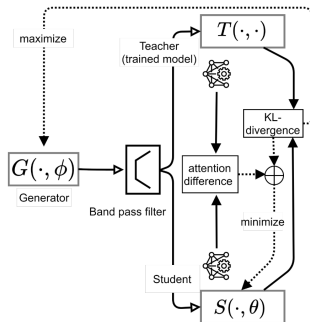
- Machine Unlearning
  - Model  $M$ , Data  $D$
  - Request:
    - Forget Data  $D_f \subset D$
    - Retain Data  $D_r = D - D_f$
  - Gold / Retrained Model:  $M^*$
  - Unlearned Model:  $M_u$
  - Aim:  $M_u(x) \approx M^*(x)$
- Zero-Shot Machine Unlearning
  - No Access to  $D_f$  and  $D_r$
- Proposes two approaches - restricted setting of classification
- Setting
  - Set of Forget  $C_f$  and Retain Classes  $C_r$

# Error Minimization-Maximization Noise



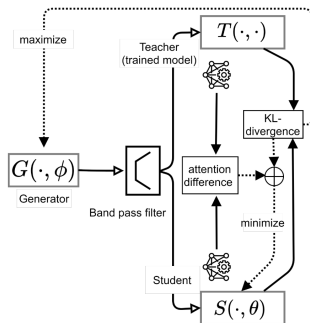
- Anti-Samples  $\mathcal{N}_f$  learnt by maximising loss
- Data representatives  $\mathcal{N}_r$  learnt by minimising loss
- Updates the original model using noise

# Gated Knowledge Transfer



- Knowledge Distillation to train the student from teacher

# Gated Knowledge Transfer



- Knowledge Distillation to train the student from teacher
- Student - Minimise KL
- Attention - Mimic Inner Layers
- Generator:  $\text{Max } D_{KL}(T(x_g)||S(x_g)) = \sum_i^{|C|} t_p^{(i)} \log(t_p^{(i)} / s_p^{(i)})$
- Filter images belonging to  $C_f$

# Entropy Criterion

- Entropy of predictions
  - Reject if  $S(t_p) > \epsilon$
  - Faster Retain Accuracy Restoration
  - Poorer Retain Accuracy
  - Carrying out experimentations

# Deep Inversion

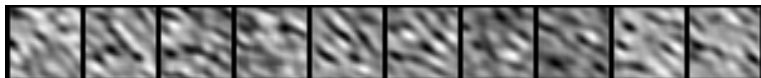
- Difference in  $M^*$  and  $M_u$ 
  - Non-zero Accuracy for  $C_f$
  - Due to Attention implicitly learn for  $C_f$
  - Removing Attention: Impacts Performance
  - Reason: Poor Generated Images
  - Deep Inversion
  - Much Better Images



# Generated Images

MNIST Numbers Dataset  
Images of Digits from 0-9

- GKT\*



- GKT (with entropy criterion)\*



- Deep Inversion



\* Images not in order from 0-9. Images generated by the generator before forget accuracy begin to rise

# Experimental Results

MNIST Numbers Dataset - A11CNN Model

Train: 60,000, Test: 10,000

Retrain Accuracy on Test Set:

- Retrain Model: 99.25 %
- GKT: 97.12 %
- M-M: 10.57 %

# Experimental Results

MNIST Numbers Dataset - AllCNN Model

Train: 60,000, Test: 10,000

Retain Accuracy on Test Set:

- Retrain Model: 99.25 %
- GKT: 97.12 %
- M-M: 10.57 %
- GKT (no attention): <50 %
- Deep Inversion (100 sample/class): 40 - 50 %
- Deep Inversion (6000 sample/class): 80 - 85 %

# Conclusion

- Tackling zero-shot setting
- Non-zero forget class accuracy
- Quality of images generated
- Decent Results

# Learnings

- First research experience
- Ability to read papers
- Tweaking complex machine learning code

Thank You