
Exploring Emerging Properties in Multimodal Neural Networks

Pranav Singh
ps4364@nyu.edu

Akash Mishra
am11533@nyu.edu

Mayank Poddar
mp6021@nyu.edu

1 Objective

There has been a sharp increase in the number of people using social media recently. Because of this, there has been an increase in the amount of content posted online. Manual moderation of content online is extremely difficult. To overcome this barrier, many automated moderation techniques have been designed. These techniques work with a singular modality (like text, audio, images, etc.) but suffer in a multi-modal environment where one modality complements one or more modalities implying subtle hate/aggressive speech. To address this problem, various unimodal, as well as multimodal approaches have been proposed. In general vision, only models perform worse than natural language models; with this study, we aim to address this by using novel, unimodal self-supervised techniques and further studying the scope of their expansion to make them fully multimodal. These multimodal self-supervised techniques are able to match the accuracy of fully supervised multimodally finetuned models.

2 Dataset

The 2020 hateful meme dataset Kiela et al. [2020] has confounding examples, i.e., counter-intuitive examples based on text or visual features alone. This makes learning from a single modality impossible. Furthermore, even within the purview of multimodal approaches, this is exceptionally challenging as vision and language embeddings extracted from two images may be the same. Still, the downstream labels for the two images would be different. Hence Embedding collapse is an issue.

For the hateful meme challenge dataset we have 8500 images in the training set and a further 2000 images in the testing set. We perform an 80/20 split to get the validation set on the training set. Each entry in the dataset has an image, label, and extracted caption.

3 Approach

3.1 Unimodal and Multimodal Learning

For our approach, we pretrained unimodal and multimodal self-supervised techniques on the training set for 100 epochs. Once trained, we then perform end-to-end finetuning on the self-supervised models. We conducted these experiments over five different seed values to account for variability.

3.2 Graph Neural Network

The second approach that we explored was using GNNs (Graph neural networks). Here, we tried to generate subgraphs for each meme using the image and textual features. Since we want to create dependence between two different kinds of data, i.e., image and text, we try to associate text embedding with image features. Once we have generated this initial subgraph, we pass this into our GNN model, namely GCN Kipf and Welling [2016], Sage, and GAT Veličković et al. [2017]. Our intuition is that this graph will learn the dependence between the multimodal data and help us identify whether the meme is hateful.

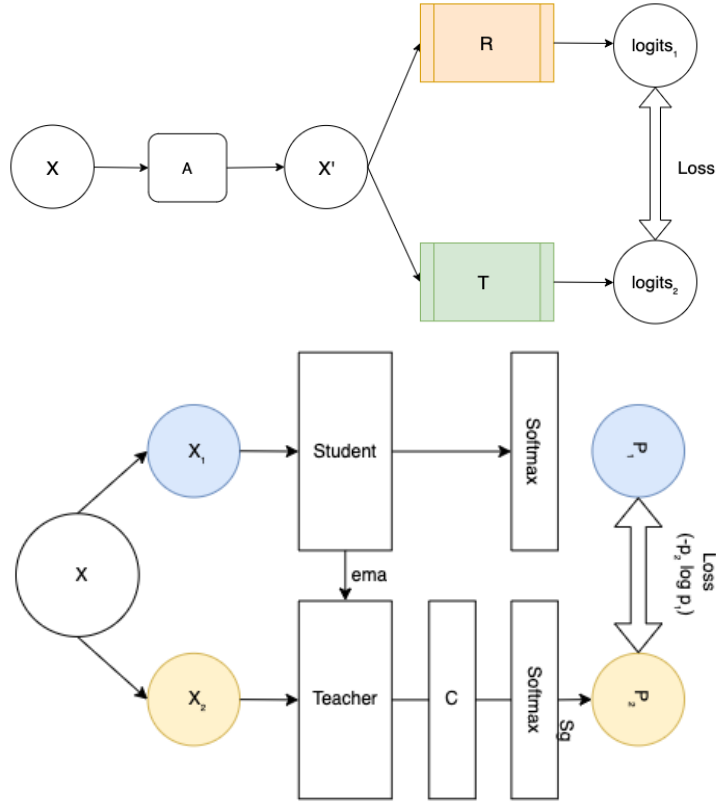


Figure 1: (Top) We show CASS, R represents ResNet-50, a CNN and T in the other box represents the Transformer used (ViT Base/Base-16); X is the input image, which becomes X' after applying augmentations. Note that CASS applies only one set of augmentations to create X'. X' is then passed through both the arms to compute similarity loss. This is different from DINO, that passes different augmentation of the same image through networks with the same architecture but different parameters. The output of the teacher network is centred on a mean computed over a batch. Another key difference is that in CASS, loss is computed over logits meanwhile in DINO it is computed over softmax output.

4 Implementation

For unimodal pre-training, we used CASS Singh et al. [2022] and DINO Caron et al. [2021]. We trained these techniques only on images. Then, during the fine-tuning process, we concatenated text embedding with the image embedding before finally passing through a classifier to perform downstream classification. We expanded upon this by making CASS fully multimodal during the pre-training task. We use ResNet-50 He et al. [2016] and Vision Transformer Dosovitskiy et al. [2020] for obtaining the image embeddings. Similarly, to obtain the text embeddings we use DistilBERT Sanh et al. [2019].

For GNNs, we will use Image and Text Encoder to get the relevant features. We use ResNet-50 pretrained on COCO Dataset for our image encoder and Distilbert pretrained from hugging face. Finally, we use Mask-RCNN He et al. [2017] along with an Image encoder to generate feature embeddings depicting entities in the image. Similarly, we also pass the text through Tokenizer and then through Distilbert to get textual embedding. A graph is constructed by joining the image feature embeddings to textual embedding and the original image. This subgraph will be passed into larger Graph neural networks to get us co-dependent embedding. This newly generated embedding can then be passed through Linear layers and sigmoid function to get us a classification of the Multimodal data.

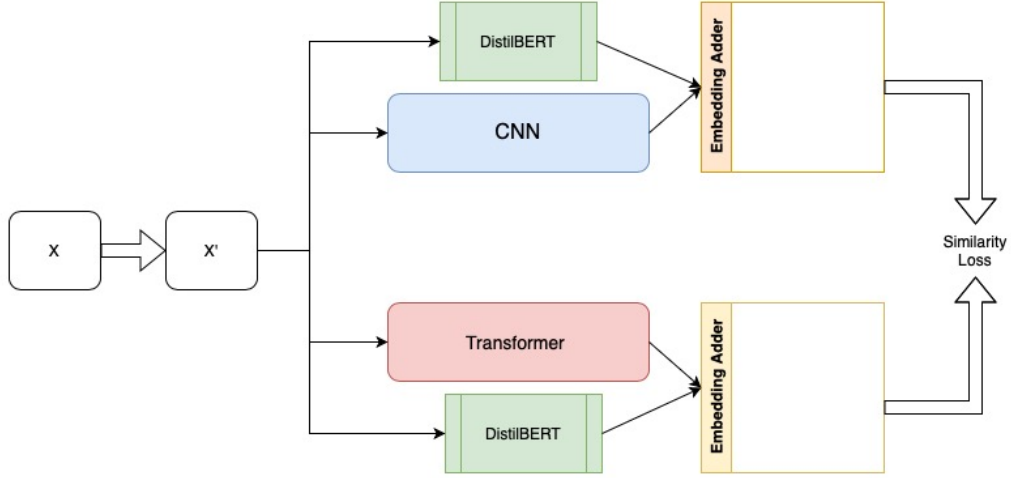


Figure 2: In this figure we show CASS, adopted for multimodal pretraining using DistilBERT and CASS. In this we take an input and apply minimal augmentations for image (resize and channel normalization), and tokenize the text using DistilBERT tokenizer. We then pass this augmented image and tokenized through a combination of vision and text encoder as shown above. The outputs of from the two are then combined in an embedding adder, before calculating counteractive loss during pretraining.

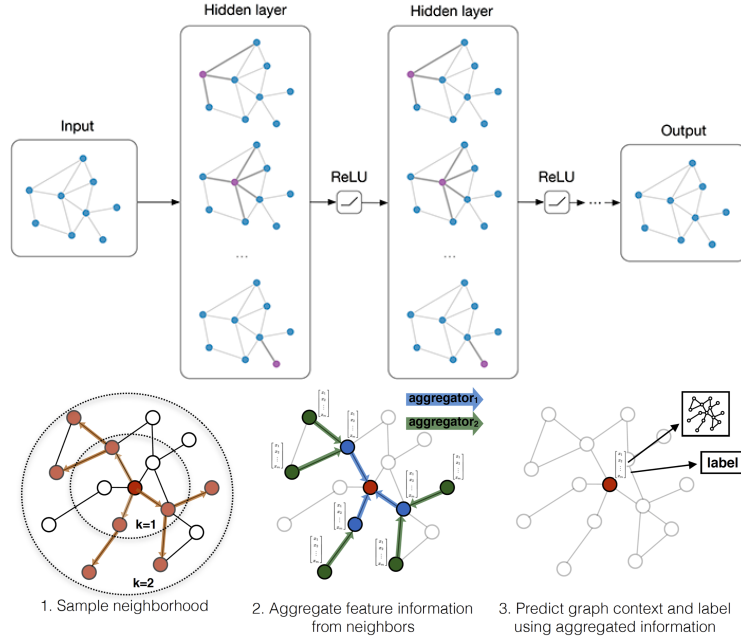


Figure 3: (Top) GCN is based on an efficient variant of convolutional neural networks which operate directly on graphs. Idea comes from getting a localized first-order approximation of spectral graph convolutions. (Bottom) Unlike other GNN, GraphSAGE Hamilton et al. [2017] provides inductive framework to efficiently generate node embeddings for unseen data.

5 Challenges

A major challenge with the Hateful meme dataset is the confounding samples, i.e two memes can have the same image but with different captions, and thereby different labels. Similarly, two memes

can have different images with the same text and hence have different downstream labels. This makes mapping the samples in a latent space extremely difficult. To mitigate this problem we developed multimodal solutions i.e instead of using only one modality we included embedding inputs from the two involved modalities i.e text and image. Furthermore, we have used experimented with graph neural networks (GNN), as they have amazing scene and relationship representation capabilities. We present our results in Table Another limiting factor is the relatively poor performance of vision-only methods as compared to text-only methods.

Lastly, in the case of Graph Neural network, we are facing trouble building a global heterogenous GNN which would better take into account the multimodal nature of our challenge.

Also, since each meme within the dataset is transformed into a subgraph, the scale of data explodes tremendously which limits our batch size and in turn our training speed.

6 Experimentation Results

In Figure 2 we describe and compare the architecture of CASS and DINO.

Modality	Model	AUROC	Accuracy
Unimodal	ResNet-50 224 (Supervised)	0.5157±0.081	0.5175±0.044
	ResNet-50 224 (CASS)	0.5398±0.091	0.5455±0.031
	ResNet-50 224 (DINO)	0.5129±0.009	0.5267±0.055
	ResNet-50 384 (Supervised)	0.5116±0.076	0.511±0.016
	ResNet-50 384 (CASS)	0.5405±0.045	0.5445±0.015
	ResNet-50 384 (DINO)	0.532±0.033	0.543±0.056
Multimodal	ResNet-50 - 384 (CASS)	0.501±0.078	0.551±0.023

Table 1: In this table we compare the performance of unimodally and Multimodally trained CASS and unimodally trained DINO and supervised ResNet-50. We observed that CASS trained Resnet-50 for input size 384, performs the best on the primary metric of comparison, but overall multimodally trained CASS ResNet-50.

Modality	Model	AUROC	Accuracy
Unimodal	VitB16 - 224 (Supervised)	0.51293±0.043	0.5235±0.065
	VitB16 - 224 (DINO)	0.513±0.077	0.5195±0.014
	VitB16 - 224 (CASS)	0.5196±0.051	0.5545±0.01
	VitB16 - 384 (Supervised)	0.5±0.022	0.545±0.087
	VitB16 - 384 (DINO)	0.5002±0.011	0.624±0.071
	VitB16 - 384 (CASS)	0.53386±0.009	0.59±0.017
Multimodal	VitB16 - 384 (CASS)	0.5±0.078	0.625±0.032

Table 2: In this table we compare the performance of unimodally and Multimodally trained CASS and unimodally trained DINO and supervised ViT/Base16. We observed a sinmilar trend as in Table 1, results for input image 384 were considerably better than that for 224. CASS trained models perform generally better than other models show considerable increase in performance over DINO and supervised training.

GNNs	Encoder	AUROC	Accuracy
GCN	Resnet50 + Distilbert	0.514	0.547
SAGE	Resnet50 + Distilbert	0.481	0.625
GAT	Resnet50 + Distilbert	0.541	0.499

Table 3: In this table we compare the performance of various GNN namely GCN, GraphSage and GAT. With the limited scope in GNN experiment, we are not seeing any significant performance improvement here.

7 High Performance Optimization

For High Performance Optimization, we took our best performing model on multimodality - Multimodal Cross Architecture Self-Supervised Learning model and tried to profile and optimize it.

Initially, we looked at profiling our model on CUDA using torch's autograd profiler to get a better understanding of the computational complexity - FLOPs and identifying the bottleneck operations. Here's the output of the profiler:

Two key points to notice here on the operations sorted by cuda time:

Name	Self CPU s	Self CPU %	CPU total s	CPU total %	CPU time avg	Self CUDA s	Self CUDA %	CUDA total s	CUDA time avg	CPU Mem	Self CPU Mem	CUDA Mem	Self CUDA Mem	# of Calls	Total MFLOPs
model inference	2.17s	69.930ms	99.13s	2.499s	2.499s	56.192ms	2.23s	2.520s	2.520s	420 B	-420 B	0 B	-5.33 GB	1	---
aten::conv2d	0.03s	829.000us	79.17s	2.050s	37.220ms	647.000us	0.03s	2.019s	37.170ms	0 B	0 B	509.99 MB	0 B	54	98897.316
aten::convolution	0.04s	1.404ms	79.68s	2.050s	37.205ms	851.000us	0.04s	2.017s	37.340ms	0 B	0 B	509.99 MB	0 B	54	---
aten::convolution	0.04s	1.517ms	79.62s	2.007s	37.174ms	730.000us	0.03s	2.016s	37.143ms	0 B	0 B	509.99 MB	0 B	54	---
aten::convolution	0.04s	1.594ms	79.64s	2.006s	37.148ms	2.013s	79.00s	2.013s	37.109ms	0 B	0 B	509.99 MB	345.29 MB	54	---
aten::conv2d	0.19s	4.780ms	4.31s	108.495ms	2.050ms	4.626ms	0.18s	107.899ms	2.032ms	0 B	0 B	502.98 MB	0 B	53	---
aten::convolution	0.08s	12.084ms	2.94s	74.145ms	378.054ms	10.461ms	0.42s	106.000ms	820.031ms	0 B	0 B	502.98 MB	0 B	53	---
aten::batch_norm	0.03s	775.000us	4.11s	103.553ms	1.958ms	583.000us	0.02s	103.073ms	1.949ms	0 B	0 B	502.98 MB	0 B	53	---
aten::conv2d	3.36s	84.828ms	4.05s	102.435ms	1.933ms	88.200ms	3.50s	102.490ms	1.938ms	0 B	0 B	502.98 MB	0 B	53	---
aten::conv2d	0.22s	5.573ms	1.94s	48.821ms	290.601ms	3.794ms	0.19s	85.812ms	510.786ms	0 B	0 B	1.95 GB	0 B	58	---
aten::relu	2.33s	58.747ms	2.52s	63.662ms	1.299ms	58.480ms	2.32s	63.980ms	1.304ms	0 B	0 B	0 B	0 B	49	---
aten::batch_norm	0.01s	283.000us	2.48s	57.107ms	2.396ms	214.000us	0.01s	61.040ms	2.546ms	0 B	0 B	751.16 MB	751.16 MB	24	---
aten::batch_norm	2.06s	56.880ms	2.26s	57.090ms	2.378ms	60.831ms	2.41s	60.831ms	2.535ms	0 B	0 B	751.16 MB	751.16 MB	24	---
aten::batch_norm	0.33s	8.289ms	0.42s	10.562ms	88.183ms	55.550ms	2.36s	59.530ms	486.120ms	0 B	0 B	840.02 MB	840.02 MB	10	451852.376
aten::batch_norm	1.01s	25.511ms	1.04s	26.155ms	2.615ms	25.480ms	1.02s	25.854ms	2.589ms	0 B	0 B	76.00 MB	76.00 MB	10	24.295
aten::conv2d	0.03s	8.278ms	0.39s	22.444ms	18.134ms	18.134ms	0.01s	18.134ms	30.793ms	420 B	420 B	1.26 GB	1.26 GB	603	---
aten::batch_norm	0.47s	11.817ms	0.53s	13.317ms	70.089ms	15.744ms	0.62s	15.744ms	82.863ms	0 B	0 B	0 B	0 B	100	---
aten::conv2d	0.07s	14.410ms	0.60s	10.015ms	625.625ms	15.494ms	0.61s	15.494ms	645.583ms	0 B	0 B	376.95 MB	376.95 MB	24	---
aten::batch_norm	0.07s	1.770ms	0.78s	19.730ms	122.167ms	1.478ms	0.07s	14.478ms	89.302ms	0 B	0 B	910.73 MB	910.73 MB	1	---
aten::batch_norm	0.12s	3.100ms	0.21s	5.413ms	112.771ms	12.238ms	0.49s	12.238ms	254.959ms	0 B	0 B	846.61 MB	846.61 MB	48	50234.278
aten::batch_norm	0.00s	122.000us	0.00s	11.473ms	109.000ms	9.00s	11.97ms	11.97ms	0 B	0 B	27.00 MB	0 B	1	---	
aten::batch_norm	0.43s	10.891ms	0.47s	11.772ms	11.772ms	11.867ms	0.47s	11.867ms	11.867ms	0 B	0 B	27.00 MB	27.00 MB	1	---
aten::batch_norm	0.44s	11.422ms	0.47s	11.747ms	880.750ms	11.776ms	0.47s	11.776ms	981.333ms	0 B	0 B	6.00 MB	6.00 MB	12	---
aten::batch_norm	0.06s	1.430ms	0.54s	15.686ms	109.472ms	1.846ms	0.05s	8.125ms	63.000ms	0 B	0 B	374.33 MB	374.33 MB	1	---
aten::batch_norm	0.00s	54.000us	0.31s	7.841ms	7.841ms	56.000us	0.00s	7.839ms	7.839ms	0 B	0 B	32.00 MB	32.00 MB	1	---
aten::batch_norm	0.07s	1.730ms	0.46s	11.500ms	106.500ms	1.990ms	0.06s	7.788ms	72.111ms	0 B	0 B	407.98 MB	407.98 MB	2	108
aten::batch_norm	0.31s	7.737ms	0.31s	7.776ms	7.776ms	7.782ms	0.31s	7.782ms	7.782ms	0 B	0 B	32.00 MB	32.00 MB	1	---
aten::batch_norm	0.01s	277.000us	0.26s	6.490ms	649.000ms	336.000ms	0.01s	6.477ms	539.417ms	0 B	0 B	19.34 MB	0 B	12	---
aten::batch_norm	0.02s	540.000us	0.19s	4.683ms	95.571ms	480.000ms	0.02s	5.498ms	112.004ms	0 B	0 B	0 B	0 B	49	---
aten::batch_norm	0.15s	3.890ms	0.18s	4.910ms	86.340ms	5.283ms	0.21s	3.477ms	103.240ms	0 B	0 B	16.00 MB	0 B	53	---
aten::batch_norm	0.18s	4.517ms	0.19s	4.860ms	605.167ms	4.687ms	0.18s	4.672ms	486.000ms	0 B	0 B	0 B	0 B	12	---
aten::batch_norm	0.02s	491.000us	0.27s	6.689ms	131.157ms	380.000ms	0.02s	4.764ms	93.412ms	0 B	0 B	210.00 MB	0 B	51	---
aten::batch_norm	0.10s	2.520ms	0.24s	5.978ms	117.216ms	3.067ms	0.12s	4.384ms	85.961ms	0 B	0 B	110.00 MB	0 B	51	---
aten::conv2d	0.06s	1.395ms	0.09s	2.300ms	21.296ms	3.543ms	0.14s	3.543ms	32.806ms	0 B	0 B	0 B	0 B	108	---
aten::batch_norm	0.02s	302.000us	0.13s	3.243ms	270.250ms	3.312ms	0.13s	3.312ms	276.000ms	0 B	0 B	731.54 MB	731.54 MB	12	191.787
aten::batch_norm	0.09s	2.361ms	0.17s	4.412ms	20.521ms	2.346ms	0.09s	3.164ms	14.716ms	0 B	0 B	0 B	0 B	215	---
aten::batch_norm	0.07s	1.800ms	0.19s	4.748ms	36.477ms	1.520ms	0.06s	3.159ms	28.299ms	0 B	0 B	0 B	0 B	130	---
aten::batch_norm	0.12s	3.004ms	0.16s	4.131ms	6.990ms	2.311ms	0.09s	2.311ms	3.910ms	0 B	0 B	0 B	0 B	391	---
aten::batch_norm	0.11s	2.831ms	0.08s	14.745ms	58.512ms	1.933ms	0.06s	2.011ms	7.980ms	0 B	0 B	0 B	0 B	262	---
aten::batch_norm	0.05s	1.323ms	0.07s	1.720ms	39.765ms	1.820ms	0.07s	1.820ms	35.484ms	0 B	0 B	208.59 MB	208.59 MB	51	51.342
aten::conv2d	0.04s	1.424ms	0.14s	3.471ms	29.168ms	1.936ms	0.05s	1.724ms	14.487ms	0 B	0 B	0 B	0 B	119	---
aten::batch_norm	0.04s	4.610ms	0.09s	2.283ms	3.794ms	1.357ms	0.09s	1.357ms	0 B	0 B	0 B	0 B	0 B	394	---
aten::conv2d	0.00s	114.000us	0.04s	1.070ms	91.083ms	139.000ms	0.01s	1.040ms	87.333ms	0 B	0 B	14.31 MB	0 B	12	---
aten::conv2d	0.01s	195.000us	0.04s	899.000ms	249.750ms	270.000ms	0.01s	895.000ms	248.250ms	0 B	0 B	4.32 MB	0 B	4	---
aten::batch_norm	0.02s	432.000ms	0.04s	951.000ms	79.250ms	937.000ms	0.04s	481.000ms	78.000ms	0 B	0 B	17.44 MB	17.44 MB	12	---
aten::batch_norm	0.01s	189.000ms	0.02s	548.000ms	137.000ms	341.000ms	0.01s	545.000ms	136.250ms	0 B	0 B	4.32 MB	0 B	4	---
aten::batch_norm	0.00s	89.000ms	0.03s	721.000ms	183.500ms	83.000ms	0.00s	466.000ms	233.000ms	0 B	0 B	4.00 MB	0 B	2	---
aten::batch_norm	0.00s	111.000ms	0.02s	406.000ms	33.833ms	144.000ms	0.01s	390.000ms	32.500ms	0 B	0 B	0 B	0 B	12	---
aten::batch_norm	0.00s	60.000ms	0.02s	518.000ms	172.000ms	40.000ms	0.00s	332.000ms	132.667ms	0 B	0 B	7.14 MB	0 B	3	---
aten::batch_norm	0.01s	175.000ms	0.01s	312.000ms	34.667ms	232.000ms	0.01s	293.000ms	32.556ms	0 B	0 B	0 B	0 B	9	---
aten::batch_norm	0.01s	279.000ms	0.02s	446.000ms	148.667ms	229.000ms	0.01s	292.000ms	97.222ms	0 B	0 B	7.14 MB	0 B	3	---

Figure 4: Torch Autograd Profiler on Multimodal CASS

1. Convolution Operations (highlighted in red) take the majority of the time spent on the forward operation: **80.43%**
2. The arithmetic complexity / total number of FLOPs is higher for linear layers (highlighted in blue, likely arising from transformer's attention mechanisms) and FLOPs for Convolution Operation only accounts for **16.44%** of the total number of floating point operations.

Below are the results of different runs of MultiModal CASS Model on different configurations of NYU HPC.

Num GPU's	Type GPU	DataParallel	cuDNN Benchmarking	GPU Core Utilization	GPU Mem Utilization	Per Epoch Time (After Warmup)
1	RTX8000	No	No	97%	63%	8 min 51 secs
1	RTX8000	Yes	Yes	83%	54%	8 min 50 secs
4	RTX8000	Yes	Yes	44.31%	15.64%	5 min 4.8 secs

Table 4: nvidia-smi profile and time counters for model on varying number of GPUs

We deploy strong scaling on our problem statement and notice speedup and scaling efficiency as follows:

$$\text{speedup} = t_{\text{serial}} / t_{\text{parallel}} = 1.74$$
$$\text{scaling-efficiency} = t_{\text{serial}} / (t_{\text{parallel}} * p) = 43.47\%$$

To address the low GPU Core Utilization and Memory Utilization, we employed large batch.size (32 as opposed to 4 before) and increased the num_workers (14 as opposed to 4 before) to reduce the idle time for GPU and bump up the memory utilization. The speedup increased to **3.05** and the scaling efficiency bumped to **76.32%** between the serial and parallel implementations.

Num GPUs	Type GPU	DataParallel	cuDNN Benchmarking	GPU Core Utilization	GPU Mem Utilization	Per Epoch Time (After Warmup)
1	RTX8000	Yes	Yes	95.01%	58.31%	6 min 46 secs
4	RTX8000	Yes	Yes	76.43%	42.13%	2 min 13 secs

Table 5: Performance Improvement with larger batch_size and increased num_workers

Another small experiment we did is to compare the Tesla V100 GPU node vs Quadro RTX 8000 on single GPU performance.

	Quadro RTX8000	Tesla V100
CUDA Cores	4608	5120
Tensor Cores	576	640
Memory	48GB	32GB
GPU Core Utilization	95.01%	95.42%
GPU Mem Utilization	58.31%	46.29%
Per Epoch Time (after Warmup)	6 min 46 secs	5 min 37 secs

Table 6: Performance contrast between RTX8000 and V100 1-GPU implementations

	Training Time
Graph Convolution Network	5 hour 48min
Graph Attention Network	6 hour 29min
GraphSage	5 hour 46min

Table 7: Total time comparison for Graph Networks training on RTX8000 with Core Utilization at **87.76%** and Memory Utilization at **62.93%**

Surprisingly, even though RTX8000 has 48GB memory vs V100 has 32GB memory, the memory utilization in case of Tesla V100 is lesser. This can be seen in all our experiments.

8 Conclusions

In table 8, we summarize our best performing models. We observed that the use of self-supervised pretraining allows models to learn helpful language and image based priors that then helps them to perform better on downstream classification tasks. Furthermore, we also compared the performance of traditional architectures with GNNs.

Modality	Model	AUROC	Accuracy
Unimodal	ResNet-50 (CASS)	0.5405±0.045	0.5445±0.015
	Vit/Base-16 (CASS)	0.53386±0.009	0.59±0.017
Multimodal	ResNet-50 (CASS)	0.501±0.078	0.551±0.023
	Vit/Base-16 (CASS)	0.5±0.078	0.625±0.032
Multimodal GNN	GCN	0.514	0.547
	GAT	0.564	0.499
	SAGE	0.481	0.625

Table 8: In the above table we summarize our best performing models, based on performance metrics described by the challenge on the test set.

References

Douwe Kiela, Hamed Firooz, Aravind Mohan, Vedanuj Goswami, Amanpreet Singh, Pratik Ring-shia, and Davide Testuggine. The hateful memes challenge: Detecting hate speech in multimodal memes. *Advances in Neural Information Processing Systems*, 33:2611–2624, 2020.

- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2016. URL <https://arxiv.org/abs/1609.02907>.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks, 2017. URL <https://arxiv.org/abs/1710.10903>.
- William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs, 2017. URL <https://arxiv.org/abs/1706.02216>.
- Pranav Singh, Elena Sizikova, and Jacopo Cirrone. Cass: Cross architectural self-supervision for medical image analysis. *arXiv preprint arXiv:2206.04170*, 2022.
- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn, 2017. URL <https://arxiv.org/abs/1703.06870>.