# Enhancing Unsupervised Image Representation Learning: A Synthesis of ImageGPT with NLP-Inspired Techniques and Adapter Modules

**Anonymous Authors**[1]

## Abstract

This paper delves into the realm of unsupervised image representation learning, leveraging advancements in Natural Language Processing (NLP) to enhance image processing models. The motivation is drawn from the success of unsupervised pre-training in NLP and the untapped potential in image domain, particularly in direct pixel prediction. Building on the foundation laid by "Generative Pretraining from Pixels," which introduced a novel sequence Transformer architecture for pixel prediction, this work extends the exploration into efficient transfer learning using adapter modules, as proposed in "Parameter-Efficient Transfer Learning for NLP." Our implementation focuses on the practical adaptation of the ImageGPT model, employing linear probing across its layers on the CIFAR-10 dataset. The results demonstrate significant insights into the model's feature representation capabilities, with optimal performance observed at intermediate layers. This study not only reinforces the practicality of leveraging pretrained models for classification tasks but also contributes to the broader understanding of deep learning models' internal representations, highlighting the potential for more parameter-conservative strategies in image processing.

## 1. Introduction

The project is rooted in the growing intersection of image processing and NLP techniques, particularly in the context of unsupervised learning. The motivation for this research stems from the challenges and opportunities presented by images as a more complex modality compared to text for generative modeling. While generative pre-training has shown significant promise in NLP, its application in image processing, especially in direct pixel prediction, is relatively unexplored. This gap in research sets the stage for our study, which aims to leverage the potential of unsupervised learning in image processing, drawing inspiration from the triumphs in NLP.

The papers reviewed for this project lay the foundational concepts and methodologies that guide our implementation. "Generative Pretraining from Pixels" introduces a sequence Transformer architecture for predicting image pixels, a method departing from traditional convolutional neural network approaches. This methodology aligns with our goal of exploring unsupervised learning techniques in image representation. Simultaneously, "Parameter-Efficient Transfer Learning for NLP" introduces adapter modules, which offer a novel approach to efficient transfer learning without the need for extensive retraining or fine-tuning of the entire model. This concept of parameter efficiency is particularly appealing for our project, considering the computational intensity of training models from scratch.

Our implementation focuses on adapting the ImageGPT model, a state-of-the-art model in image processing, for diverse tasks using linear probing. This approach involves appending a new layer or classifier to the model's existing structure and assessing its adaptability across different tasks within the CIFAR-10 dataset. The objective is to explore how effectively a pretrained model, initially designed for a specific domain, can be repurposed for broader applications. This exploration is driven by the hypothesis that mid-network layers, rather than the final layer, may offer more effective representations for downstream tasks, challenging conventional linear probing approaches. The study culminates in an analysis of the model's adaptability and the potential for incorporating NLP-inspired methodologies in enhancing unsupervised image representation learning.

---

[1]Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

## 2. Review of paper to implement or extend

### 2.1. Storyline

#### 2.1.1. HIGH-LEVEL MOTIVATION

The research in "Generative Pretraining from Pixels" is driven by the potential of unsupervised learning in image processing. In an era abundant with data but scarce in labeled datasets, there's a compelling need for models that can derive meaningful insights without explicit supervision. Inspired by the triumphs of unsupervised pre-training in Natural Language Processing, the authors venture into the image domain, aiming to predict pixels without relying on traditional 2D structures. This work hints at a future where models can adeptly understand and generate visual content without the necessity of labeled data, promising significant advancements in fields like computer vision and medical imaging.

#### 2.1.2. PRIOR WORK ON THIS PROBLEM

Unsupervised pre-training has played a crucial role in deep learning, with methods like Deep Belief Networks and Denoising Autoencoders commonly used in neural networks for various applications. While these methods initially showed promise, advancements in neural networks reduced the need for such pre-training. The approach of unsupervised pre-training then gained traction in the domain of Natural Language Processing (NLP), with models like BERT pushing the state of the art forward. In the realm of images, generative models like GANs and VAEs have been explored for representation learning, but most haven't been as competitive as supervised and self-supervised methods.

#### 2.1.3. RESEARCH GAP

Images, being a more complex modality than text, have been traditionally challenging for generative modeling. While generative pre-training has shown promise in NLP, its potential in image processing, especially in direct pixel prediction, remains underexplored compared to recent self-supervised methods.

#### 2.1.4. CONTRIBUTIONS

The paper presents several key contributions to the field of unsupervised representation learning for images. Firstly, the authors re-evaluate generative pre-training on images, demonstrating its competitiveness with other self-supervised approaches. This re-evaluation is particularly significant given the historical context and recent advancements in the domain. Secondly, they introduce a novel methodology that employs a sequence Transformer architecture, traditionally used in NLP, to predict pixels in images. This represents a departure from the conventional use of convolutional neural networks in image processing. Furthermore, their approach

showcases significant improvements in low-resolution unsupervised representation learning settings, pushing the boundaries of what's achievable in this domain.

### 2.2. Proposed solution

The paper introduces a novel method for unsupervised image representation learning using generative pre-training. It utilizes a sequence Transformer architecture, commonly used in Natural Language Processing (NLP), to predict image pixels. This approach comprises two primary stages:

#### 2.2.1. PRE-TRAINING STAGE:

In this phase, the model is trained on a large dataset without using any labels. The aim is to learn general features and representations from the data. Two main objectives are explored during this stage:

**Auto-regressive Objective:** This involves predicting the next pixel in a sequence based on the previous pixels. Mathematically, the model is trained to minimize the negative log-likelihood of the data ($L_{AR}$), where X is an unlabeled dataset consisting of high dimensional data $x = (x_1, ..., x_n)$ and a permutation $\pi$ of the set $[1, n]$ is picked:

$$p(x) = \prod_{i=1}^{n} p(x_{\pi_i} | x_{\pi_1}, \ldots, x_{\pi_{i-1}}, \theta) \tag{1}$$

$$L_{AR} = \mathbb{E}_{x \sim X}[-\log p(x)] \tag{2}$$

**Bert Objective:** Inspired by the BERT model from NLP, this objective involves masking a subset of pixels and then training the model to predict the masked pixels based on the unmasked ones. Specifically, a sub-sequence $M$ is sampled such that each index $i$ has a 15% probability of being "masked" or hidden. The model then predicts the "masked" elements $x_M$ conditioned on the "unmasked" ones $X_{[1,n] \setminus M}$, with the objective being

$$L_{BERT} = \mathbb{E}_{x \sim X} \mathbb{E}_M \left[ -\sum_{i \in M} \log p\left(x_i | x_{[1,n] \setminus M}\right) \right] \tag{3}$$

#### 2.2.2. MEASURING REPRESENTATION QUALITY STAGE:

**Fine-tuning** After pre-training, the model is further refined on a smaller labeled dataset specific to the task at hand. This involves adding a small classification head to the model and optimizing a classification objective.

**Linear probing** Another approach introduced to measure the quality of the representations learned is linear probing. In this method, fixed features are extracted from the model, and a linear classifier is trained on these features. This

technique, known as "linear probing," captures the intuition that good features should be able to linearly separate the classes of transfer tasks. The features are viewed as fixed during linear probing, and a projection is learned to produce class logits from these intermediate features. This projection contains the only trainable weights, so the optimization is only done on $L_{CLF}$

## 2.3. Claims-Evidence

### 2.3.1. CLAIM1

**Claim:** Better generative models learn better representations

**Evidence:** In Figure 1, as the validation loss on the auto-regressive objective decreases throughout training, the accuracy of the linear probe increases. This trend is consistent across various model sizes, with larger models achieving better validation losses.
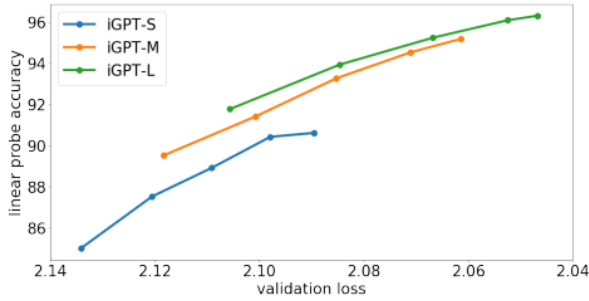


*Figure 1.* Fig 3 in section 4.3 from(Chen et al., 2020) The graph demonstrates that as the generative model improves, so does the quality of the learned representations.

### 2.3.2. CLAIM2

**Claim:** The quality of representations from generative models varies with depth, and the best representations often lie in the middle layers of the network.

**Evidence:** Figure 2 showcases how representation quality first improves as a function of depth and then starts deteriorating around the middle layer, continuing until the penultimate layer. This behavior suggests that generative models might operate in two phases, resembling encoder-decoder architectures but learned within a single architecture via a pre-training objective.

### 2.3.3. CLAIM3

**Claim:** Predicting pixels directly, as done in the paper's approach, leads to state-of-the-art representations for low-resolution datasets.
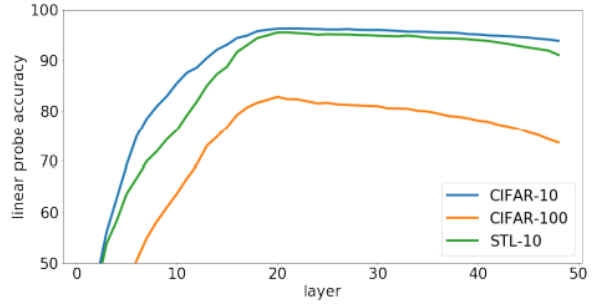


*Figure 2.* Fig 2 in section 4.1 from(Chen et al., 2020) This figure presents a relationship between the validation loss on the auto-regressive objective and the accuracy of the linear probe throughout the generative pre-training. The graph uses dotted markers to denote checkpoints at various steps (e.g., 65K, 131K, 262K, 524K, and 1000K). The positive slope of the graph suggests a link between improved generative performance and improved representation quality..

**Evidence:** The results indicate that the proposed approach of directly predicting pixels achieves state-of-the-art performance in low-resolution settings. In high-resolution scenarios, the approach remains competitive with other self-supervised methods on ImageNet.On CIFAR-10, the proposed approach achieves an accuracy of 99.0%, outperforming other notable methods such as AutoAugment and GPipe.

| Model | Acc | Unsup Transfer | Sup Transfer |
|---|---|---|---|
| **CIFAR-10** | | | |
| AutoAugment | 98.5 | | |
| GPipe | 99.0 | | ✓ |
| iGPT-L | 99.0 | ✓ | |
| **CIFAR-100** | | | |
| iGPT-L | 88.5 | ✓ | |
| AutoAugment | 89.3 | | |
| EfficientNet | 91.7 | | ✓ |

*Figure 3.* Table 3 in section 4.4 from(Chen et al., 2020)

## 2.4. Critique and Discussion

This article presents a novel approach to image learning by drawing inspiration from techniques prevalent in natural language processing (NLP). One of the most captivating aspects of this research is the authors' innovative methodology of downscaling images into a one-dimensional sequence. Traditionally, image learning has been approached with at least 2D or higher-dimensional input data. However, this pa-

| Model | Acc | Unsup Transfer | Sup Transfer |
|---|---|---|---|
| **CIFAR-10** | | | |
| AMDIM-L | 91.2 | ✓ | |
| ResNet-152 | 94 | | ✓ |
| iGPT-L | 96.3 | ✓ | |
| **CIFAR-100** | | | |
| AMDIM-L | 70.2 | ✓ | |
| ResNet-152 | 78 | | ✓ |
| iGPT-L | 82.8 | ✓ | |
| **STL-10** | | | |
| AMDIM-L | 94.2 | ✓ | |
| iGPT-L (IR $32^2 \cdot 3$) | 95.5 | ✓ | |
| iGPT-L (IR $96^2 \cdot 3$) | 97.1 | ✓ | |

*Figure 4.* Table 1 in section 4.3 from(Chen et al., 2020)

per challenges the norm by training a sequence Transformer to autoregressively predict pixels, a technique reminiscent of NLP methodologies.

While the paper delves into the intricacies of the process, the mention of the downsampling technique using VQ-VAE could have been elucidated further. Although the results indicate that this technique performs moderately on extensive training datasets like ImageNet, a deeper dive into its workings would have provided readers with a more comprehensive understanding.

The evidence and claims presented in the article are robust and align well with the assumptions made by the authors. Their experimental setup, especially the choice of datasets and the evaluation metrics, provides a solid foundation for their claims. It's commendable how the authors have ventured away from the traditional convolutional neural network methods, opting instead for sequence Transformers for image processing. This paradigm shift, while innovative, carries inherent risks. However, the results presented in the paper not only validate this approach but also hint at the potential future applications of such methodologies.

In conclusion, the paper offers a fresh and intriguing perspective on image learning, effectively bridging techniques from NLP and image processing. While certain aspects could benefit from further clarity, the overall direction and results of the research are commendable and pave the way for future explorations in this domain.

## 3. Review of second paper

### 3.1. Storyline

#### 3.1.1. HIGH-LEVEL MOTIVATION

Transfer learning has emerged as a powerful technique in the realm of Natural Language Processing (NLP). By leveraging knowledge from pre-trained models, researchers have been able to achieve impressive performance across a myriad of NLP tasks. However, the online setting, where tasks arrive in a stream, poses a challenge. The ideal scenario is to have a system that excels across all tasks without necessitating the training of a new model for each incoming task. This is especially pertinent for applications like cloud services, where models are trained to solve multiple tasks that arrive sequentially. The inefficiency of fine-tuning large pre-trained models for every downstream task is a significant concern in this context.

#### 3.1.2. PRIOR WORK ON THIS PROBLEM

Historically, transfer learning in NLP has been approached in two primary ways: feature-based transfer and fine-tuning. Feature-based transfer typically involves pre-training real-valued embedding vectors at various levels (word, sentence, paragraph) which are then fed to custom downstream models. On the other hand, fine-tuning involves adjusting the weights of a pre-trained network for the downstream task. Notably, BERT, a Transformer network trained on vast text corpora with an unsupervised loss, has set the benchmark by achieving state-of-the-art performance on tasks like text classification and extractive question answering.

#### 3.1.3. RESEARCH GAP

Despite the successes of the aforementioned methods, there exists a gap in achieving parameter efficiency. Fine-tuning, while effective, often requires training a significant number of parameters for each new task. This is not only computationally expensive but also storage-intensive, especially when dealing with a series of tasks. The challenge, therefore, is to devise a method that retains the effectiveness of transfer learning while being parameter-efficient.

#### 3.1.4. CONTRIBUTIONS

This paper introduces a novel approach to address the above gap: adapter-based tuning. The primary contributions are:
1. The proposal of "adapter modules" as an alternative to full fine-tuning. These modules are compact, adding only a minimal number of trainable parameters per task, ensuring a high degree of parameter sharing.
2. A demonstration that, using adapters, it's possible to achieve near state-of-the-art performance while adding only a fraction of parameters per task.
3. An empirical evaluation of the adapter-based tuning

approach, highlighting its efficiency and effectiveness across various NLP tasks and benchmarks.

4. An in-depth analysis and discussion on the robustness and architectural choices of adapter modules, providing insights into their design and performance trade-offs.

## 3.2. Proposed solution

### 3.2.1. ADAPTER MODULES

Adapter modules are designed to facilitate transfer learning within Transformer architectures without fine-tuning the entire model. Despite exploring various designs, the authors found that simpler configurations yielded better performance.

### 3.2.2. TRANSFORMER WITH ADAPTER

**Transformer Layers:** Each layer of a Transformer model typically consists of two primary sub-layers: an attention mechanism and a feed-forward neural network.

**Skip-Connections:** In Transformers, each of these sub-layers is usually connected via a "skip connection" (or residual connection) which helps with gradient flow during training.

**Adapter Placement:** Adapters are placed after both the attention and feed-forward sub-layers. The processed output from an adapter feeds into the skip connection and then proceeds to the next layer.

### 3.2.3. ADAPTER'S FUNCTIONALITY

**Bottleneck Architecture:** The adapter employs a bottleneck design to ensure parameter efficiency. It first reduces the feature dimensions from the original size (denoted as "d") to a much smaller size (denoted as "m"). After processing, it projects the features back to the original dimension "d".

**Parameter Calculation:** The number of parameters introduced by this bottleneck design in each layer is given by the formula

Downward projection weights + Upward projection weights

$\quad$ + Downward biases + Upward biases

$= (d \times m) + (m \times d) + m + d$

$= 2md + d + m$

$$\quad (4)$$

This design ensures that the additional parameters introduced by the adapter are only a small fraction (ranging from 0.5% to 8%) of the original model's parameters, striking a balance between performance and efficiency.

## 3.3. Claims-Evidence

### 3.3.1. CLAIM1

**Claim:** Adapter-based tuning achieves a balance between performance and the number of trained parameters, outperforming traditional fine-tuning in terms of parameter efficiency.

**Evidence:** Figure 4 (from the paper) showcases the validation set accuracy versus the number of trained parameters for various methods, including adapter tuning. The results indicate that adapter tuning, even with smaller adapter sizes, achieves comparable performance to full fine-tuning but with significantly fewer trained parameters.
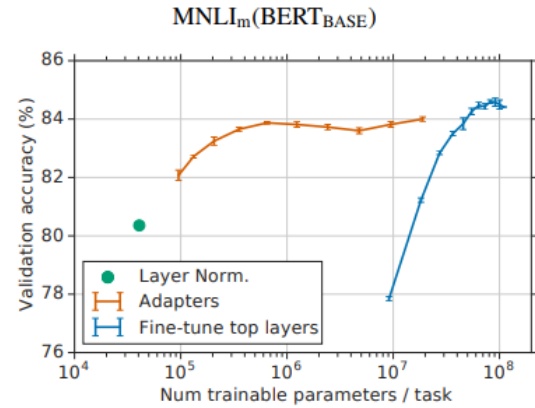


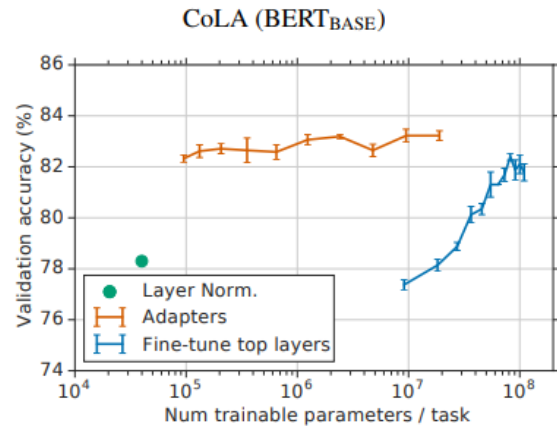*Figure 5.* Fig 4 in section 3.4 from(Houlsby et al., 2019)



*Figure 6.* Fig 4 in section 3.4 from(Houlsby et al., 2019)

### 3.3.2. CLAIM2

**Claim:** Adapter-based tuning is effective not just for classification tasks but also for more complex tasks like extractive question answering.

**Evidence:** The results on the SQuAD v1.1 dataset, as presented in the paper, demonstrate that adapters achieve near state-of-the-art performance. For instance, adapters of size 64 (2% parameters) attain a best F1 of 90.4%, which is very close to the 90.7% achieved by full fine-tuning.
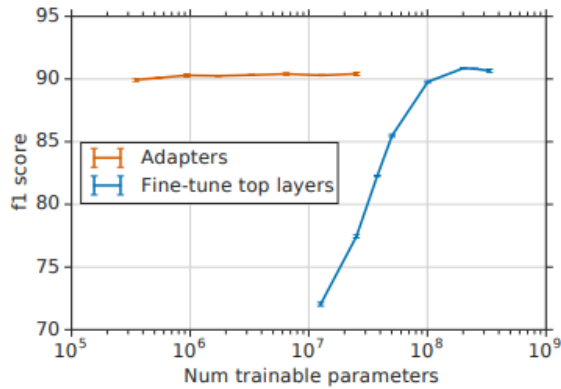


*Figure 7.* Fig 5 in section 3.6 from(Houlsby et al., 2019)

### 3.3.3. CLAIM3

**Claim:** The performance of adapter modules is robust to variations in initialization scale and the number of neurons.

**Evidence:** The paper presents an analysis where the initialization scale of the adapter module weights is varied. The results, summarized in Figure 6 (from the paper), indicate that performance remains consistent for standard deviations below $10^{-2}$. Additionally, the paper's experimental data from Section 3.2 suggests that the performance of adapters remains robust across different adapter sizes, further emphasizing their versatility.

### 3.4. Critique and Discussion

The paper introduces the innovative concept of adapter modules to address transfer learning inefficiencies, a notable achievement given the reduced parameter count while maintaining performance. However, some architectural details, particularly the bottleneck design and initialization strategies, lacked depth. While the experiments, especially on the GLUE benchmark, substantiate the paper's claims, a broader comparison with other transfer methods would enrich the findings. The chosen benchmarks are appropriate, but the
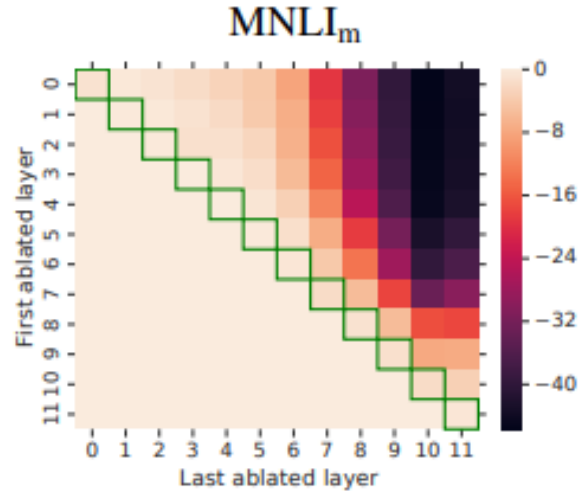


*Figure 8.* Fig 6 in section 3.6 from(Houlsby et al., 2019)

paper might have delved deeper into adapter variations. Additionally, the assumption of tasks arriving sequentially may not always align with real-world scenarios. Overall, the paper offers a fresh take on NLP transfer learning, but certain areas could benefit from enhanced clarity and exploration.

**How might this paper related to the paper I'm implementing(Generative Pretraining from Pixels)** Both the first paper and the second paper underscore the significance of unsupervised pretraining, albeit in different domains. While the latter delves deep into the intricacies of natural language processing using adapter modules, the former innovatively employs a sequence Transformer to autoregressively predict pixels.

The most compelling takeaway from the NLP-focused paper, when juxtaposed with "Generative Pretraining from Pixels," is the strategy for achieving efficient transfer learning. The iGPT model, as highlighted in the first paper, showcases remarkable performance across diverse datasets. Yet, the "Full Fine-tuning" methodology it employs can be resource-intensive. This is where the brilliance of the "Adapter tuning" technique from the second paper shines. By minimizing the parameters introduced during the fine-tuning phase, it offers a promising avenue for enhancing the efficiency of the iGPT model's application across a broader spectrum of datasets.

In essence, the insights from "Parameter-Efficient Transfer Learning for NLP" could pave the way for more parameter-conservative transfer learning strategies in the context of "Generative Pretraining from Pixels," broadening its applicability and efficiency.
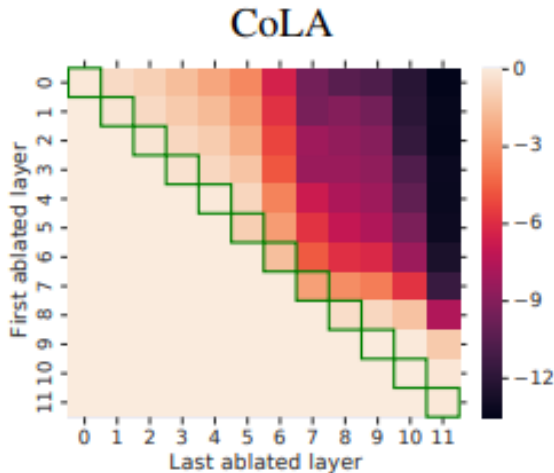
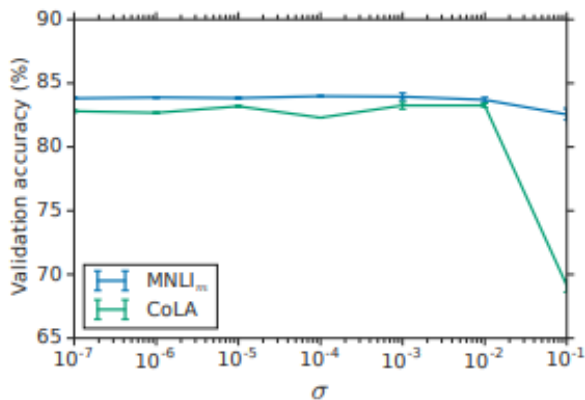*Figure 9.* Fig 6 in section 3.6 from(Houlsby et al., 2019)



*Figure 10.* Fig 6 in section 3.6 from(Houlsby et al., 2019)

## 4. Review of third paper

### 4.1. Storyline

#### 4.1.1. HIGH-LEVEL MOTIVATION

The realm of computer vision has traditionally relied on models pre-trained on crowd-labeled datasets, such as ImageNet. However, the potential of leveraging the vast amount of text available on the internet to train computer vision systems offers a broader source of supervision compared to traditional methods that rely on specific object categories. The overarching vision of the research is to harness the power of natural language supervision to revolutionize computer vision in a manner similar to how pre-training methods transformed NLP.

#### 4.1.2. PRIOR WORK ON THIS PROBLEM

Historically, there have been attempts to enhance content-based image retrieval by training models to predict textual elements in documents paired with images. For instance, Mori et al. explored this approach over two decades ago. Quattoni et al. demonstrated the potential of learning data-efficient image representations by predicting words in captions associated with images. More recent works, like that of Joulin et al., have shown that CNNs trained to predict words in image captions can learn valuable image representations. These efforts, while promising, have not fully harnessed the potential of natural language supervision at scale.

#### 4.1.3. RESEARCH GAP

While prior research has made strides in integrating text and image data, there remains a gap in efficiently and scalably harnessing the potential of natural language supervision for computer vision. Traditional methods have been limited in their scope, often not leveraging the vast amount of available text on the internet. Additionally, there's a need for models that can perform zero-shot transfer across a wide range of tasks without requiring dataset-specific training.

#### 4.1.4. CONTRIBUTIONS

The paper introduces CLIP, a method that demonstrates the efficacy of predicting which caption aligns with which image as a scalable way to learn state-of-the-art image representations from scratch. This approach is based on a dataset of 400 million (image, text) pairs collected from the internet. After pre-training, CLIP uses natural language to reference or describe visual concepts, enabling zero-shot transfer to various downstream tasks. The research benchmarks CLIP on over 30 different existing computer vision datasets, showcasing its versatility and robustness. Notably, CLIP matches the accuracy of the original ResNet-50 on ImageNet zero-shot, without utilizing any of its training examples.

### 4.2. Proposed solution

The paper introduces CLIP (Contrastive Language–Image Pre-training) as a novel approach to bridge the research gap in leveraging natural language supervision for computer vision. The central idea behind CLIP is to train vision models using a contrastive objective in a joint vision-language embedding space. Rather than predicting specific words in the text accompanying an image, CLIP is designed to discern which text snippet from a set corresponds to which image. It achieves this by maximizing the similarity between an image and its paired text while minimizing the similarity with other texts in the batch.

To implement this, the paper describes a training process where an image and a text snippet are encoded into a shared embedding space. A contrastive loss is then applied to pull the embeddings of matching image-text pairs closer and push apart the embeddings of non-matching pairs. The objective is for the model to correctly pair images and texts across a large batch of negative examples.

The paper provides a numpy-like pseudocode for the core of the CLIP implementation: Where $I_e$ and $T_e$ are the em-

```
# image_encoder - ResNet or Vision Transformer
# text_encoder  - CBOW or Text Transformer
# I[n, h, w, c] - minibatch of aligned images
# T[n, l]       - minibatch of aligned texts
# W_i[d_i, d_e] - learned proj of image to embed
# W_t[d_t, d_e] - learned proj of text to embed
# t             - learned temperature parameter

# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T)  #[n, d_t]

# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)

# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)

# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss   = (loss_i + loss_t)/2
```

*Figure 11.* Fig 3 in section 2.2 from(Radford et al., 2021)

beddings for images and texts, respectively, and loss is the symmetric loss function.

By connecting visual concepts with natural language in this manner, CLIP is capable of performing zero-shot transfer across a multitude of tasks without the need for dataset-specific training. The paper further benchmarks CLIP on over 30 different existing computer vision datasets, emphasizing its potential as a robust and scalable method for computer vision.

### 4.3. Claims-Evidence

#### 4.3.1. CLAIM1

**Claim:** CLIP significantly improves performance on various datasets, demonstrating its capability as a flexible and practical zero-shot computer vision classifier.

**Evidence:** In the paper, Table 1 compares Visual N-Grams to CLIP. The best CLIP model improves accuracy on ImageNet from 11.5% (achieved by Visual N-Grams) to 76.2%. This matches the performance of the original ResNet-50

without using any of the 1.28 million crowd-labeled training examples available for ImageNet. The top-5 accuracy of CLIP models is also notably high.

| | aYahoo | ImageNet | SUN |
|---|---|---|---|
| Visual N-Grams | 72.4 | 11.5 | 23.0 |
| CLIP | **98.4** | **76.2** | **58.5** |

*Figure 12.* table 1 in section 3.1 from(Radford et al., 2021)

#### 4.3.2. CLAIM2

**Claim:** Zero-shot CLIP demonstrates competitive performance when compared to few-shot methods, often matching or even surpassing the performance of models trained with multiple shots.

**Evidence:** As visualized in Figure 5, zero-shot CLIP matches the average performance of a 4-shot linear classifier trained on the same feature space and nearly matches the best results of a 16-shot linear classifier across publicly available models.
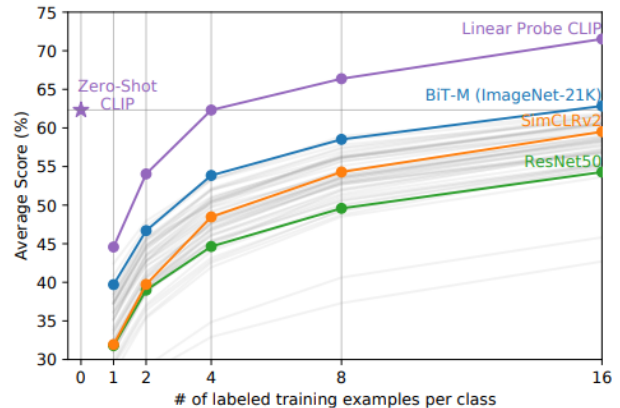


*Figure 13.* Fig 5 in section 3.2 from(Radford et al., 2021)

#### 4.3.3. CLAIM3

**Claim:** CLIP models demonstrate superior representation learning capabilities, outperforming state-of-the-art (SOTA) computer vision models in terms of both overall score and compute efficiency.

**Evidence:** Figure 6 showcases the performance of CLIP models against other SOTA computer vision models, with

8

scores averaged over datasets studied by Kornblith et al. (2019) and a broader set of 27 datasets. Notably, the largest CLIP model slightly edges out the performance of the renowned Noisy Student EfficientNet-L2. Furthermore, the standout model from the paper, ViT-L/14@336px, consistently surpasses the best existing models across the evaluation suite, achieving an average performance lead of 2.6%.
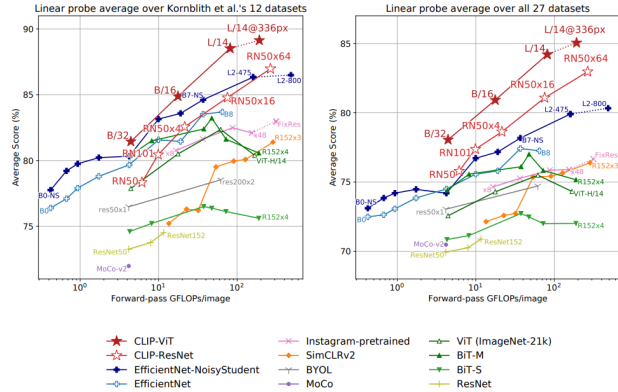


*Figure 14.* Fig 6 in section 3.3 from(Radford et al., 2021)

### 4.4. Critique and Discussion

The paper introduces the innovative CLIP method, showcasing the potential of natural language supervision in computer vision. The contrastive objective in a joint vision-language space is a highlight, offering a fresh approach to training vision models. However, while the paper provides compelling evidence of CLIP's performance, a deeper dive into the model architecture, ethical considerations, and potential biases would have added depth. Additionally, while the creation of the WIT dataset is commendable, more transparency on data collection and preprocessing would be beneficial. Overall, the paper presents a promising direction, but certain areas could benefit from further clarity and exploration.

**How might this paper related to the paper I'm implementing(Generative Pretraining from Pixels)** The third paper and the first paper both delve deep into enhancing the capabilities of domain-specific models using unsupervised learning techniques. The latter focuses on the potential of natural language supervision in the field of computer vision, drawing inspiration from the success of unsupervised pretraining in natural language processing (NLP). The former, on the other hand, is driven by the potential of unsupervised learning in image processing and aims to predict pixels without relying on traditional 2D structures. Both of these papers are inspired by the success of unsupervised pretraining in NLP. In the third paper, the CLIP method is designed to

harness the power of natural language to train visual models, with a particular focus on zero-shot capabilities. On the other hand, the first paper reevaluates generative pretraining for images and introduces a novel approach that employs a sequence Transformer architecture, traditionally used in NLP, for pixel prediction in images.

## 5. Implementation

### 5.1. Implementation Motivation

The inspiration for this implementation stems from the intriguing application of natural language processing (NLP) techniques in image learning, as detailed in the referenced article. The original study's ability to leverage generative image modeling for unsupervised image representation, particularly the adaptation of ImageGPT models, resonated with my research interests. However, the complexity and computational intensity of training such models from scratch present substantial challenges. Consequently, my implementation pivots towards exploiting pretrained models for downstream tasks, aligning with methods outlined in the article.

Two primary methods are under consideration: linear probing and entire model fine-tuning. Initially, I gravitated towards model fine-tuning, inspired by the "Parameter-Efficient Transfer Learning for NLP" article, which promises computational efficiency. This method's potential to achieve effective transfer learning with reduced resource demand was a compelling aspect, as discussed in the "Critique and Discussion" section. Unfortunately, preliminary experiments with this fine-tuning approach did not yield the anticipated results, prompting a reassessment of the strategy.

Given this outcome, the focus has shifted towards linear probing. This approach involves adding a new layer or classifier to the pretrained ImageGPT model's last layer, thereby adapting it to various tasks. The objective is to explore how effectively the ImageGPT model, trained through linear probing, can adapt to diverse tasks. This exploration aims to uncover the extent to which a pretrained model, initially designed for a specific domain (image representation), can be repurposed for a broader spectrum of applications. The experimental results from this approach are expected to provide insights into the versatility and adaptability of pretrained models in machine learning.

### 5.2. Implementation setup and plan

The primary objective of my implementation is to utilize the pretrained ImageGPT model for a series of linear probing experiments on the CIFAR-10 dataset. Recognized as a benchmark in image classification, CIFAR-10's variety of images is ideal for assessing the ImageGPT model's adaptability across different image recognition tasks.

### 5.2.1. CODE BASE AND DATASET:

My approach will harness the 'transformers' library from Hugging Face to access the ImageGPT model and employ the 'datasets' library for the CIFAR-10 dataset. These libraries are chosen for their comprehensive APIs, which are conducive to effective and streamlined model training and evaluation. What's more, Building on the foundation provided by Hugging Face's transformer documentation and github repository from(Wolf et al., 2020), I will develop and train various models and classifiers to enhance the ImageGPT model's application to a broader spectrum of downstream tasks.

### 5.2.2. METHODOLOGY:

The implementation will encompass the following key stages:

**1. Data Loading and Preprocessing:** The CIFAR-10 dataset will be loaded and preprocessed using the 'datasets' library to conform to the input specifications of the ImageGPT model.

**Feature Extraction:** Leveraging the pretrained ImageGPT model, I will extract features from the images in CIFAR-10, focusing on capturing the hidden states across the model's layers.

**Linear Probing:** A logistic regression classifier will be trained on the extracted features. This classifier will be applied to features from each layer of ImageGPT to identify the layer offering the most effective image classification representation.

**Evaluation Metrics:** Accuracy will be the primary metric for performance evaluation, calculated for both the training and testing subsets of CIFAR-10. Additionally, t-SNE visualizations will be employed to illustrate the differences in accuracy between the most effective and the final layers.

**Additional Classification Models:** Expanding the scope, I plan to incorporate other classifiers like k-Nearest Neighbors (k-NN) and neural networks, and compare their performance against the logistic regression model.

### 5.2.3. PRIORITIZATION AND ALIGNMENT WITH MOTIVATION

My focus will primarily be on exploring the efficacy of linear probing with ImageGPT, resonating with my intent to investigate the transferability of NLP techniques to the realm of image learning. The outcomes of this research will offer valuable insights into the broader applicability and versatility of pretrained models within the field of machine learning.

### 5.3. Implementation details

My implementation primarily utilized the Hugging Face transformers library for accessing the pretrained ImageGPT model and the datasets library for the CIFAR-10 dataset. The feature extraction process, tailored to extract hidden states from each layer of ImageGPT, was custom-implemented, while the logistic regression classifier for linear probing was adapted from scikit-learn.

Key modifications included adapting the logistic regression training process to accommodate the format of extracted features. Standard scikit-learn and matplotlib methods were used for additional classifiers and t-SNE visualizations, respectively.

In summary, the implementation combined established libraries for model and dataset access, custom feature extraction, and adaptation of conventional classification and visualization techniques. This approach was chosen to explore the adaptability of the ImageGPT model to image classification tasks, aligning with the project's research goals.

### 5.4. Results and interpretation

In my study, I sought to enhance the ImageGPT model's adaptability for new tasks by employing linear probing across its layers. My results align well with my initial expectations and reveal significant insights into the model's feature representation capabilities.

From my logistic regression analysis (Figure 15), I observed optimal classification performance at the 13th layer, with accuracy reaching 78.22%. This finding is consistent with my hypothesis that intermediate layers of the ImageGPT model harbor a balanced representation of features conducive to classification tasks, particularly within the CIFAR-10 dataset. The substantial drop in accuracy at the final layer, plummeting to 48.34%, suggests a diminishing utility of deeper features for direct classification, possibly due to over-specialization towards the model's pretraining tasks.

The t-SNE visualizations (Figures 16 and 17) provide compelling visual evidence supporting my quantitative results. At the 13th layer, the training and test set features display clear class delineation, underscoring the discriminative strength of the representations. Conversely, the last layer's visualization illustrates a muddled feature space, indicating less class separability.

The k-NN classifier performance (Figure 18) peaked at a much lower accuracy compared to the neural network and logistic regression, maxing out around 32%. This could imply that k-NN, which relies on proximity in feature space, may not be as effective in leveraging the abstract representations

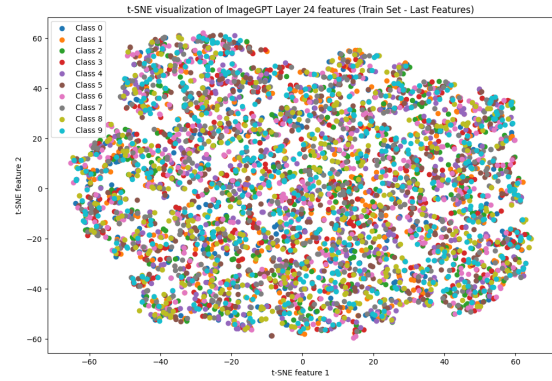*Figure 15.* Logistic Regression Model Accuracy



*Figure 17.* t-SNE visualization of ImageGPT Layer 24 features
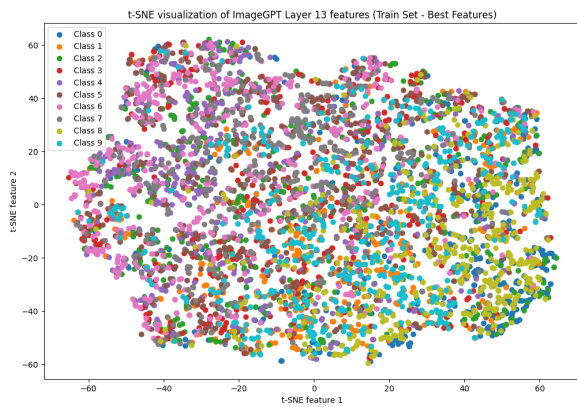


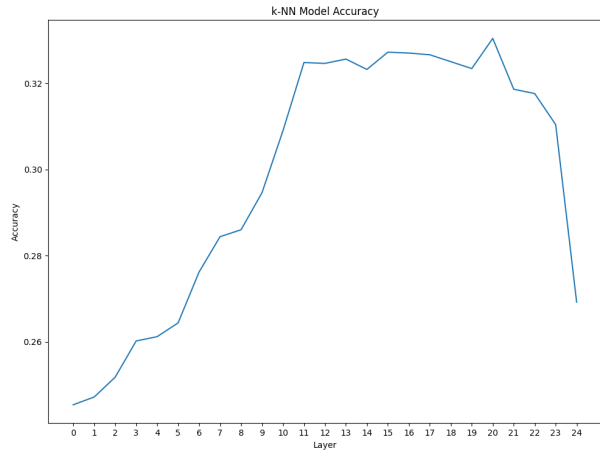*Figure 16.* t-SNE visualization of ImageGPT Layer 13 features



*Figure 18.* K-Nearest Neighbors Model Accuracy

learned by ImageGPT.

Remarkably, my neural network classifier (Figure 19) surpassed all other models, nearly achieving an 80% accuracy peak. This superior performance accentuates the neural network's capability to exploit the complex and non-linear associations in the data, which the ImageGPT model encodes. A recurrent theme across my classifiers was the suboptimal performance of the last layer, challenging the conventional approach of linear probing that typically focuses on the final layer's output. My study advocates for a paradigm shift towards utilizing mid-network layers for downstream tasks, which could yield more effective transfer learning strategies.

Despite computational constraints limiting the study to three classifiers, the evidence suggests that the ImageGPT model is robust and versatile. Its ability to adapt to diverse downstream tasks through linear probing is clear, with mid-layer features offering the most promising results. My study

not only reinforces the practicality of leveraging pretrained models for classification tasks but also contributes to the broader understanding of deep learning models' internal representations.

My findings have practical implications for those in the field of machine learning, providing a strategy to unlock the full potential of pretrained models without the need for extensive retraining. The results should be of interest to researchers and practitioners aiming to apply deep learning models efficiently across various domains.

## 6. Conclusion and Discussion

In this project, we ventured into integrating NLP methodologies into unsupervised image representation learning, primarily focusing on enhancing ImageGPT model capabilities using linear probing and NLP-inspired adapter modules. This approach aimed to address the challenge of efficiently
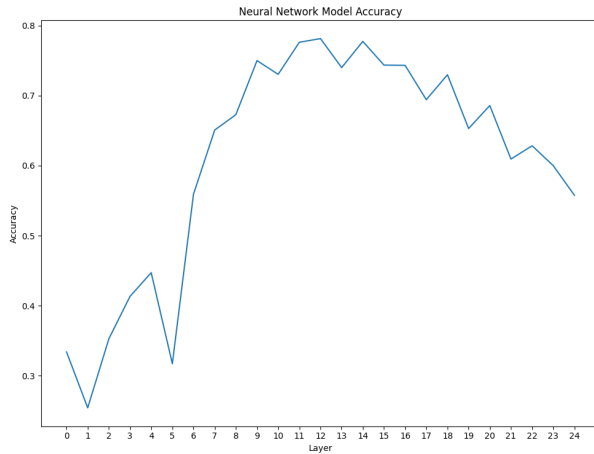
11

*Figure 19.* Neural Network Model Accuracy

processing visual content without relying heavily on labeled data, a significant gap in fields like computer vision and medical imaging.

**Key Insights**  Our findings challenged traditional approaches by revealing that intermediate layers of ImageGPT, rather than the final layers, provided more effective feature representations for image classification tasks. This insight suggests a shift in transfer learning strategies, emphasizing the potential of mid-network layers.

**Limitations and Future Directions**  The study faced computational complexity challenges, common in large-scale models like ImageGPT. Future research could further explore adapter modules to balance performance with computational efficiency. Additionally, expanding the range of datasets for evaluation would provide a more comprehensive understanding of the model's versatility. Tailoring linear probing techniques to specific model and dataset characteristics could also enhance efficacy.

**Broader Implications**  This research contributes to the overarching goal of advancing unsupervised learning in image processing. By addressing the identified challenges, future work can further integrate NLP techniques in image processing, potentially leading to more sophisticated, efficient image analysis methods, crucial for fields relying on visual data.

## References

Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., and Sutskever, I. Generative pretraining from pixels. In III, H. D. and Singh, A. (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1691–1703. PMLR, 13–18 Jul 2020. URL https://proceedings.mlr.press/v119/chen20s.html.

Houlsby, N., Giurgiu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., and Gelly, S. Parameter-efficient transfer learning for NLP. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2790–2799. PMLR, 09–15 Jun 2019. URL https://proceedings.mlr.press/v97/houlsby19a.html.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., and Sutskever, I. Learning transferable visual models from natural language supervision. In Meila, M. and Zhang, T. (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR, 18–24 Jul 2021. URL https://proceedings.mlr.press/v139/radford21a.html.

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Le Scao, T., Gugger, S., Drame, M., Lhoest, Q., and Rush, A. M. Transformers: State-of-the-Art Natural Language Processing. pp. 38–45. Association for Computational Linguistics, October 2020. URL https://github.com/huggingface/transformers/tree/d1a00f9dd0b851245e4a54cbd70816a80e781ec2/src/transformers/models/imagegpt.