

ViTamin: Designing Scalable Vision Models in the Vision-Language Era

Jieneng Chen^{1*} Qihang Yu^{2*} Xiaohui Shen² Alan Yuille¹ Liang-Chieh Chen²
¹Johns Hopkins University ²ByteDance *equal contribution
<https://beckschen.github.io/vitamin.html>

Abstract

Recent breakthroughs in vision-language models (VLMs) start a new page in the vision community. The VLMs provide stronger and more generalizable feature embeddings compared to those from ImageNet-pretrained models, thanks to the training on the large-scale Internet image-text pairs. However, despite the amazing achievement from the VLMs, vanilla Vision Transformers (ViTs) remain the default choice for the image encoder. Although pure transformer proves its effectiveness in the text encoding area, it remains questionable whether it is also the case for image encoding, especially considering that various types of networks are proposed on the ImageNet benchmark, which, unfortunately, are rarely studied in VLMs. Due to small data/model scale, the original conclusions of model design on ImageNet can be limited and biased. In this paper, we aim at building an evaluation protocol of vision models in the vision-language era under the contrastive language-image pretraining (CLIP) framework. We provide a comprehensive way to benchmark different vision models, covering their zero-shot performance and scalability in both model and training data sizes. To this end, we introduce ViTamin, a new vision models tailored for VLMs. ViTamin-L significantly outperforms ViT-L by 2.0% ImageNet zero-shot accuracy, when using the same publicly available DataComp-1B dataset and the same OpenCLIP training scheme. ViTamin-L presents promising results on 60 diverse benchmarks, including classification, retrieval, open-vocabulary detection and segmentation, and large multi-modal models. When further scaling up the model size, our ViTamin-XL with only 436M parameters attains 82.9% ImageNet zero-shot accuracy, surpassing 82.0% achieved by EVA-E that has ten times more parameters (4.4B).

1. Introduction

The past decades have witnessed significant progress in computer vision, like visual recognition tasks. The advent of AlexNet [72] marked a significant milestone, catalyzing the extensive evolution and dominance of Convolutional

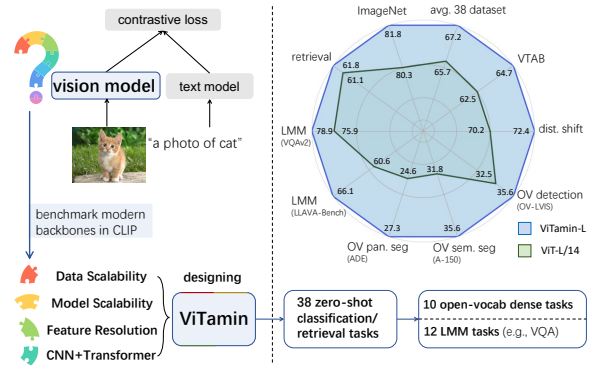


Figure 1. **Practices of designing scalable vision models in the vision-language era.** We benchmark modern vision models with various model and data scales under CLIP setting using DataComp-1B [41], leading to findings about data and model scalability, feature resolution, and hybrid architecture, which motivate us to develop ViTamin for VLM. ViTamin-L achieves superior zero-shot performance over ViT-L/14 [80] on ImageNet [114] and average 38 datasets [41], and advances a suite of 22 downstream tasks for Open-Vocabulary (OV) detection [143] and segmentation [159], and Large Multi-modal Model (LMM) tasks [88].

Neural Networks (ConvNets) [12, 44, 50, 51, 61, 75, 93, 94] in computer vision. More recently, with the debut of Vision Transformer [32, 133], a growing number of transformer-based architectures [25, 92, 132, 139, 151, 156] have shown great potential to surpass the prior ConvNet counterparts.

The rapid advancement of neural network design in computer vision can be attributed to a combination of factors. Among them, an important factor is the well-established benchmarks, allowing the community to examine the developments in a standardized way. Particularly, ImageNet [114] has become the de facto testing ground for new vision models. It not only sets a standard benchmark for visual recognition, but also serves as a mature pre-training dataset for transferring the network backbone to a variety of downstream tasks (e.g., detection and segmentation) [13, 14, 19, 51, 68, 87, 94, 124, 137, 152, 157, 158].

Recently, the emergence of vision-language models (VLMs) [65, 108] has changed the paradigm by leveraging the pre-training schedule on the extremely large scale noisy

Internet data up to billions of image-text pairs [117], much larger than the ImageNet scale. VLMs not only produce strong and generalizable features [65, 108], but also excel in zero-shot downstream tasks [43, 73, 89, 99, 113, 159, 175]. However, unlike the ImageNet benchmark, where many types of neural networks are designed and blossomed [50, 58, 61, 72, 119, 126], the existing VLMs mostly employ the vanilla Vision Transformer (ViT) architecture [32]^{*}, and the recent benchmark DataComp [41] focuses on the data curation under the common (yet unverified) belief that ViTs scale much better than any other architectures in this vision-language era [27, 81] and thus ViT is all we need.

The current trend can be characterized by several key observations: (1) The high computational demand, requiring extensive resources [63] for months, is a significant barrier for advancing VLMs [108], limiting exploring diverse vision models. (2) Traditional vision models are mainly optimized for the ImageNet benchmark, which may not scale well for larger datasets [41, 117], unlike purely transformer-based architectures [133] that have proven scalable in language tasks [105, 131] and are now being adopted for VLMs as image encoders [32]. (3) Current VLM benchmarks focus on zero-shot classification/retrieval tasks [41], with a notable lack of downstream tasks involving open-vocabulary dense prediction [29, 30, 43, 73, 149, 150, 159, 165, 172], as well as a gap in assessing Large Multi-modal Models (LMMs) [77, 88, 89, 175].

In this paper, we aim to address the aforementioned issues with practices as shown in Fig. 1. To begin with, we establish a new test bed for designing vision models under the CLIP framework [65, 108] using the DataComp-1B dataset [41], which is one of the largest publicly available datasets with high quality. Specifically, we employ two training protocols: *short schedule* for fast benchmarking vision models across model and data scales, and *long schedule* for training best performing vision models. With the *short schedule*, we re-benchmark state-of-the-art vision models found on ImageNet settings for VLMs. Particularly, we select ViT [32], ConvNeXt [93], CoAtNet [25], as representatives for pure Transformer, pure ConvNet, and hybrid architecture, respectively. We combine various model scales and data scales to provide a comprehensive evaluation towards different architectures, revealing several critical findings. First, increasing data scales improves all vision models across all model sizes, while ViT scale slightly better than others in terms of model parameters. Second, the final resolution of the extracted features affects prediction performance. Third, CoAtNet performs better than ViT and ConvNeXt in general, though it is hard to scale up CoAtNet-4 to billions of data due to computational constraints.

Those findings motivate us to develop a new vision model, named ViTamin tailored for VLM. ViTamin is a 3-

stage hybrid architectures, combining two stages of MB-Conv blocks with a final stage of Transformer blocks. This hybrid design leverages its Transformer stage to enhance data and model scalability, along with output stride of 16 to enjoy high feature resolution. As a result, ViTamin-L outshines its ViT-L/14 counterpart [41] by +2.0% zero-shot imageNet accuracy in identical OpenCLIP training scheme and identical 256 token length. When increasing feature resolution to 576 patches, ViTamin-L further attains 81.8% zero-shot imageNet accuracy, surpassing the prior art ViT-L/14 CLIPA-v2 [80] by +1.5%. In average performance across 38 datasets, it not only exceeds ViT-L/14 counterpart [80] by +1.5%, but also outperforms the larger ViT-H/14 model [80] by +0.4% while having only half parameters. When further scaling up the model size, our ViTamin-XL with only 436M parameters attains 82.9% ImageNet zero-shot accuracy, surpassing 82.0% achieved by EVA-E (*i.e.*, EVA-02-CLIP-E/14 [123]) that has ten times more parameters (4.4B). Furthermore, we introduce an effective training scheme Locked-Text Tuning (LTT), which guides the training of vision backbone with a frozen pretrained text encoder. It enhances the small variant by +4.0% and the base variant by +4.9% without any extra cost.

Our another intriguing observation is the prevailing emphasis on data filtering over vision architecture design in VLM. For instance, while the best DataComp challenge solution [154] achieved only a +2.3% gain, our ViTamin with LTT largely improves performance by +23.3% on the same dataset size, without intensive data filtering. Finally, we introduce a suite of downstream tasks, including open-vocabulary detection and segmentation, and LMMs, for evaluating VLM-specific vision models. ViTamin outperforms the ViT-L model, enhancing detector by +3.1% on OV-LVIS and segmentor by +2.6% on average 8 datasets, and excelling across 12 LMM benchmarks. Notably, ViTamin sets a new state-of-the-art on 7 benchmarks for open-vocabulary segmentation.

We aim for our findings to encourage a reevaluation of the current limitations in VLM designs and hope that our extensive benchmarking and evaluations will drive the development of more advanced vision models for VLMs.

2. Related Work

Vision Backbone: On the ImageNet benchmark [114], ConvNets [50, 61, 72, 93, 116, 119, 125–127, 147, 163] have been the dominant networks choice since the advent of AlexNet [72]. Recently, the vision community has witnessed impressive emergence of the Transformer architecture [133], a trend that began with the widespread adoption of the ViT [32] and its subsequent developments [36, 78, 82, 92, 109, 128, 136, 139, 161, 173]. Among these, hybrid architectures [18, 25, 34, 46, 55, 83, 98, 121, 132, 140, 146, 151] combine Transformer self-attention with convo-

^{*}with only a few exceptions, *e.g.*, ConvNeXt [93] by OpenCLIP [63].

lution, where CoAtNet [25] particularly obtains impressive results on ImageNet. Notably, MaX-DeepLab [137], emerged as early as 2020, successfully developed a hybrid network backbone for dense pixel predictions, where the first two stages utilize residual bottleneck blocks [50], followed by two subsequent stages employing axial attention [136]. More recently, by leveraging the design practices of a Vision Transformer, a ResNet [50] can be modernized to ConvNeXt [93], competing favorably with ViT. Along the same direction, but not limited to the ImageNet scale, our work aims to develop a novel vision model for training with billions of data [41] in the vision-language era.

Language-Image Pre-training: Language-image pre-training has seen significant advancements [1, 3, 65, 77, 89, 108, 155] with the emergence of LLMs [11, 105, 130]. The huge progress can be attributed to the pre-training on an immense scale of noisy web-collected image-text data [41, 117], much larger than the ImageNet. Notably, CLIP [65, 108] generates strong image features and excels in zero-shot transfer learning [43, 73, 89, 99, 113, 159, 175], which make it an essential role in large multi-modal model [17, 77, 88, 89]. CLIP has been improved by advanced training strategies including self-supervised learning [85, 101], efficient tuning [90, 168] and training [79, 81, 123, 141, 169]. These studies predominantly employ ViT [32] as the only vision model. As a result, the architectural design for the CLIP vision model has not been thoroughly investigated. Thus, we attempt to bridge the gap by developing a novel vision model for VLMs.

3. Method

In the section, we revisit the problem definition of CLIP and propose two training protocols (*short* and *long* schedules) on DataComp-1B (Sec. 3.1). With *short schedule*, we re-benchmark modern vision models found on ImageNet under the CLIP setting (Sec. 3.2). We then introduce the proposed ViTamin architecture design, motivated by the discoveries in the re-benchmarking results (Sec. 3.3).

3.1. CLIP and Training Protocols

CLIP Framework: Given a batch of N image-text pairs $\{(I_1, T_1), \dots, (I_N, T_N)\}$ (where I_i and T_i denote image and text for i_{th} pair), the objective of CLIP [108] learns to align the image embeddings \mathbf{x}_i and text embeddings \mathbf{y}_i for each pair. Formally, the loss function is defined as follows:

$$-\frac{1}{2N} \sum_{i=1}^N \left(\underbrace{\log \frac{e^{\mathbf{x}_i^T \mathbf{y}_i / \tau}}{\sum_{j=1}^N e^{\mathbf{x}_i^T \mathbf{y}_j / \tau}}}_{\text{image to text}} + \log \frac{e^{\mathbf{y}_i^T \mathbf{x}_i / \tau}}{\sum_{j=1}^N e^{\mathbf{y}_i^T \mathbf{x}_j / \tau}} \right), \quad (1)$$

where $\mathbf{x}_i = \frac{f(I_i)}{\|f(I_i)\|_2}$, $\mathbf{y}_i = \frac{g(T_i)}{\|g(T_i)\|_2}$, and τ is a temperature variable. A vision model $f(\cdot)$ and a text model $g(\cdot)$ are

trained to minimize the loss function. We focus on vision model design and use the text models from OpenCLIP [63].

Training Protocols: We employ two training protocols: *short schedule* and *long schedule*. The *short schedule* is designed for efficiently benchmarking vision models up to 1 training epoch on DataComp-1B [41] (*i.e.*, 1.28B seen samples). As detailed in Tab. 2, given a descent amount of resources (*e.g.*, 32 A100 GPUs), it takes less than two days to train a small (~ 25 M parameters) model variant. The *long schedule* is designed for training the best performing models with up to 40B seen samples.

3.2. Benchmarking Vision Models in CLIP Setting

The *short schedule* allows us to efficiently re-benchmark state-of-the-art vision models found on ImageNet under the CLIP setting using DataComp-1B. The experimented models are ViT [32] (a pure Transformer), ConvNeXt [93] (a pure ConvNet), and CoAtNet [25] (a hybrid model). We examine their scalability in terms of both model scales and data sizes. Each vision model has sizes varying from small (~ 25 M parameters), base (~ 85 M) to large (~ 300 M), while the data sizes range from 128M, 512M to 1.28B training seen samples (1 epoch is equal to 1.28B seen samples). The metric is zero-shot accuracy on ImageNet, supplemented by the results on the 38 tasks following DataComp [41]. As shown in Fig. 2, we analyze the benchmarked results from four aspects, including data scalability, model scalability, feature resolution, and hybrid architecture. For simplicity, we use “X@Y” to denote the vision model X trained with Y seen samples. See appendix for numerical results.

Data Scalability: When training seen samples increase from 128M to 1.28B, we observe a consistent trend of improvements across all model sizes and all vision models (a1-a5). Interestingly, ViT-S/16@512M (22M parameter) attains the zero-shot performance of 53.8% on ImageNet, which is better than 45.8% by ViT-B/16@128M (86M parameter). It shows the effectiveness of training large scale data that quadrupling training seen samples can be more impactful than quadrupling the number of model parameters. Additionally, ViT-B/16@512M & @1.28B significantly boost ViT-B/16@128M from 45.8% to 60.0% (+14.2%) and 65.6% (+19.8%).

→ As the training seen samples increase, the performances consistently improves in all cases.

Model Scalability: When the model sizes increase, the performances of all vision models are also boosted (b1-b3). However, we observe a different gain among them (b4). For example, ConvNeXt-XL@128M brings only +1.4% gain over ConvNeXt-B@128M, while ViT-L/16@128M enhances ViT-B/16@128M by +3.6%. Given plenty of data, ViT still shows a better model scalability, especially scaling from base to large (*e.g.*, +6.4% for ViT vs. +3.6% for both CoAtNet and ConvNeXt at 512M samples; +6.3% for ViT

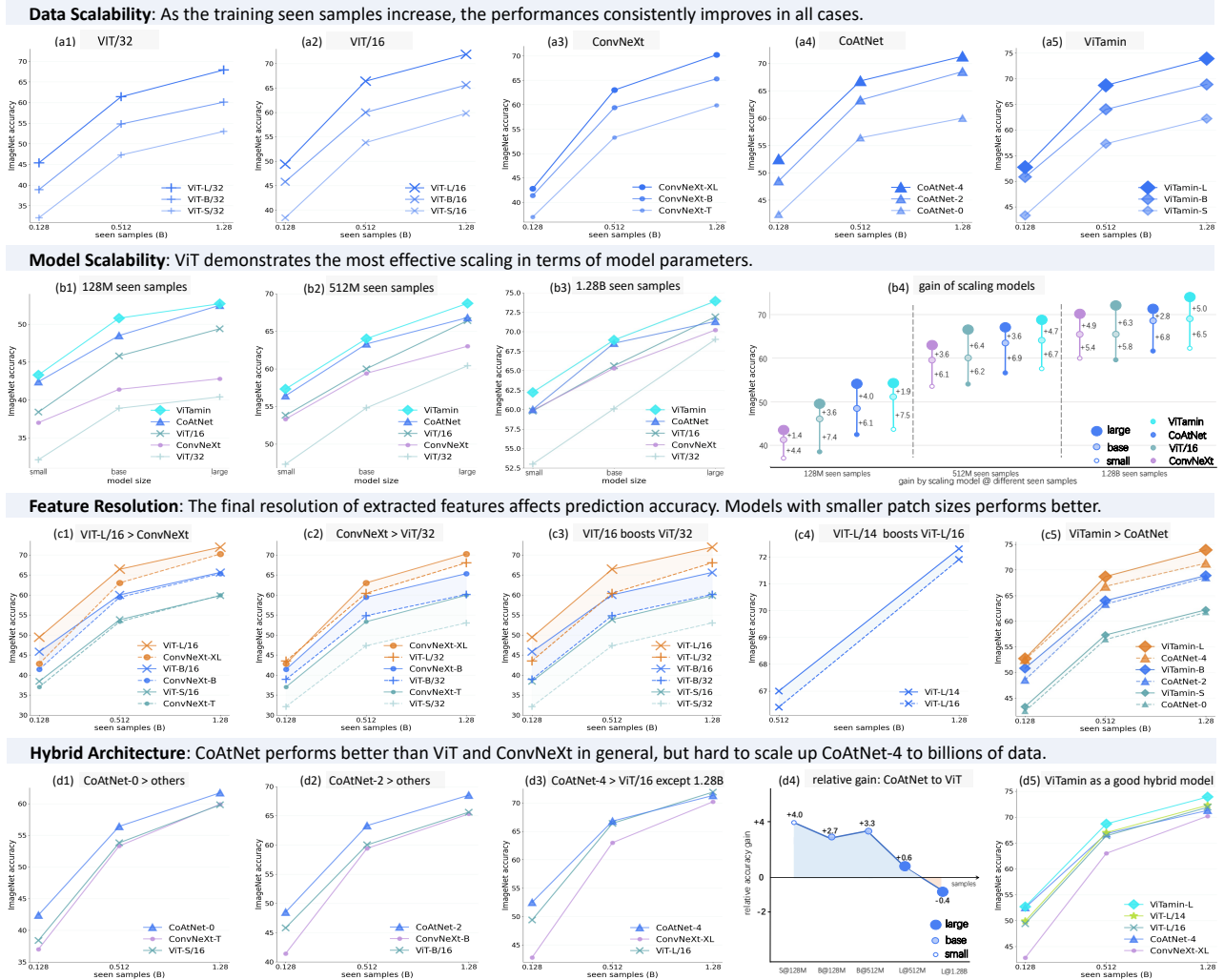


Figure 2. **Benchmarking vision models under CLIP setting on DataComp-1B**, including ViT (a pure Transformer), ConvNeXt (a pure ConvNet), and CoAtNet (a hybrid model). We examine their scalability in terms of both data sizes (1st row) and model scales (2nd row), and further analyze the results from the aspects of feature resolution (3rd row) and hybrid architecture (4th row).

vs. +2.8% for CoAtNet and + 4.9% for ConvNeXt at 1.28B samples). As a result, ViT shows the best scalability.

→ ViT demonstrates the most effective scaling in terms of model parameters.

Feature Resolution: Across all model scales and data sizes, ConvNeXt performs better than ViT/32 but loses its advantage to ViT/16 (c1 & c2). This trend deviates significantly from what is observed in ImageNet era, where ConvNeXt consistently outperforms ViT/16 (also see our ImageNet-scale VLM experiments in appendix). We hypothesize that, comparing to ImageNet’s object class label, the text in CLIP captures broader area of information, and thus is beneficial from higher feature resolution. Besides, ViT also benefits from using smaller patch sizes (thus high feature resolution) over larger path sizes (c3 & c4).

→ The final resolution of extracted features affects the prediction performance. ViT with smaller patch sizes out-

performs ViT with larger patch sizes and ConvNeXt.

Hybrid Architecture: We observe that ConvNeXt consistently lags behind ViT- $\{S,B\}/16$ and particularly ViT-L/14, suggesting that a pure ConvNet has limited capacity under the CLIP setting when presented with abundant of data (d1-d3). By contrast, CoAtNet significantly surpasses both ViT and ConvNeXt (e.g., CoAtNet-2@1.28B has a remarkable +2.9% and +3.2% gain over ViT-B/16@1.28B and ConvNeXt-B@1.28B, respectively), indicating the effectiveness of hybrid models. However, CoAtNet requires the most GPU memory; we can only train CoAtNet-4 with batch size 8k on 64 A100 GPUs, while all the other large models are trained with batch size 16k on 32 A100 GPUs. This affects CoAtNet’s scalability in large variant.

→ CoAtNet surpasses ViT and ConvNeXt in general, yet it is hard to scale up CoAtNet-4 to billions of data.

3.3. Novel Vision Transformer for Vision-Language

Herein, we distill from the aforementioned observations, culminating in the proposed vision model, ViTamin (**V**ision **T**rAnsfor**M**er for **v**ision-l**a**Ng**u**age), which notably takes the lead in the benchmarking results across all settings in Fig. 2. To introduce ViTamin, we commence by its macro-level network design (Sec. 3.3.1), followed by the micro-level block design (Sec. 3.3.2). Finally, we develop a vision model family with a simple scaling rule (Sec. 3.3.3).

3.3.1 Macro-level Network Design

Overview: The macro-level network design of ViTamin is inspired by the ViT and CoAtNet. Specifically, on top of a simple convolutional stem (*i.e.*, two 3×3 convolutions) [25], we adopt a 3-stage network architecture, where the first two stages employ the Mobile Convolution Blocks (MBConv) [57, 116] and the third stage uses the Transformer Blocks (TFB) [32, 133]. Fig. 3 shows the overview of ViTamin. We detail the design principles below, based on the discoveries from the re-benchmarking results

Data and Model Scalability: ViT demonstrates the best scalability in terms of both model scales and data sizes. We thus opt for using Transformer Blocks in our last stage, and we stack most blocks here across different model sizes.

Feature Resolution: We tailor the network to generate high resolution feature maps in the end. Our 3-stage network design thus yields a feature map with output stride 16 (*i.e.*, a downsampling factor of 16).

Hybrid Architecture: Similar to CoAtNet, we employ MBConv in the first two stages, resulting in a hybrid model. However, unlike CoAtNet that is constrained by its large memory usage, we propose a light-weight design of stage 1 and 2, which contain only two and four MBConv blocks.

Given the macro-level network design, we then move on to further improve the micro-level block design below.

3.3.2 Micro-level Block Design

Overview: The proposed ViTamin depends on two types of blocks: Mobile Convolution Blocks (MBConv) and Transformer Blocks (TFB). We refine each block in our model.

MBConv-LN: The Mobile Convolution Block (MBConv) [116] employs the “inverted bottleneck” design [50], starting with a first 1×1 convolution to expand the channel size, followed by a 3×3 depthwise convolution [58] for spatial interaction, and ending with another 1×1 convolution to revert to the original channel size. Modern MBConv, as in MobileNetv3 [57], adds numerous batch normalization (BN) [64] layers and squeeze-and-excitation (SE) [59]. We adopt a simple modification by removing all BN layers and SE, and just using a single layer normalization (LN) [4] as the first layer in our block, akin to the pre-norm layer in the Transformer block, resulting in the proposed MBConv-LN. Ablation (in appendix) shows that MBConv-LN enjoys

block	stride	ViTamin-S		ViTamin-B		ViTamin-L		ViTamin-XL	
		B	C	B	C	B	C	B	C
conv-stem	2	2	64	2	128	2	160	2	192
MBConv-LN	4	2	64	2	128	2	160	2	192
MBConv-LN	8	4	128	4	256	4	320	4	384
TFB-GeGLU	16	14	384	14	768	31	1024	32	1152

Table 1. **ViTamin model variants.** ViTamin variants differ in the number of blocks B and number of channels C in each stage.

a simple design while attaining a similar performance to the original MBConv-BN-SE in MobileNetv3.

TFB-GeGLU: The Transformer Block (TFB) [133] contains two residual blocks: one with self-attention and the other with feed-forward network (FFN). We empirically discover that substituting the first linear layer with GeGLU [118], an enhanced version of the Gated Linear Unit [26] that has a $2 \times$ expansion rate, can enhance accuracy in the FFN. We denote the Transformer Block with the updated FFN as TFB-GeGLU. Ablation (in appendix) shows that TFB-GeGLU requires 12% fewer parameters than TFB due to half expansion ratio, allowing us to stack additional Transformer blocks towards deeper architectures [125, 129, 173].

3.3.3 Meta Architecture and Scaling Rule

Meta Architecture: After introducing our macro-level network and micro-level block designs, we now put everything together to form the meta architecture of ViTamin. Specifically, ViTamin is a hybrid architecture that contains only three stages, built on top of a simple convolutional stem (*i.e.*, two 3×3 convolutions). The first two stages are composed of MBConv-LN, where we stack two and four of them for stage 1 and 2, respectively. The third stage are obtained by stacking N_B TFB-GeGLU blocks. With the meta architecture in mind, we are ready to discuss the scaling rule to generate a family of ViTamin with different model sizes.

Scaling Rule: Our scaling rule is extremely simple and straightforward, controlled by two hyper-parameters: width (*i.e.*, the channel sizes of those three stages) and depth (*i.e.*, N_B , the number of TFB-GeGLU blocks in stage 3). Note that our convolutional stem has the same channel size as the first stage. We define four model sizes: Small, Base, Large, and X-Large (S, B, and L variants have a similar amount of model parameters to ViT [32, 167]). We use the same channel size as ViT in our 3rd stage for each model variant. Specifically, we set the channel sizes of our three stages as $(C, 2C, 6C)$, where $6C = \{384, 768, 1024, 1152\}$ for Small, Base, Large and X-Large model variant, respectively[†]. Subsequently, given the target model parameter, the value of N_B (*i.e.*, the number of TFB-GeGLU blocks in stage 3) can be easily found. We show the family of ViTamin- $\{S, B, L, XL\}$ in Tab. 1.

[†]We calculate the channel size for stage 1 as $1/6C$, rounding to the nearest value that is divisible by 32.

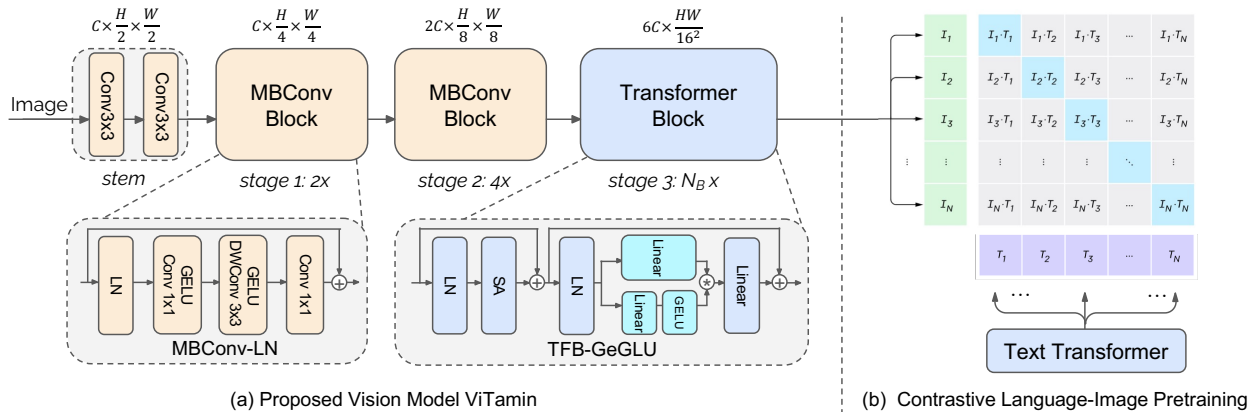


Figure 3. **Overview of ViTamin architecture.** (a) ViTamin begins with a convolutional stem, followed by Mobile Convolution Blocks (MBConv) in stage 1 and 2, and Transformer Blocks (TFB) in stage 3. The 2D input to the stage 3 is flattened to 1D. For the *macro-level* designs, the three-stage layout generates the final feature map with output stride 16, similar to ViT/16 [32]. We set channels sizes for the three stages to be $(C, 2C, 6C)$. For the *micro-level* designs, the employed MBConv-LN modifies MBConv [116] by using a single LayerNorm [4]. TFB-GeGLU upgrades TFB’s FFNs [133] (Feed-Forward Networks) with GELU Gated Linear Units [118]. (b) In the CLIP framework, given N image-text pairs, the vision model’s output I_i is learned to align with its corresponding text Transformer’s output T_i . Our text Transformers are the same as OpenCLIP [63]. +: Addition. *: Multiplication.

	short schedule for benchmarking			long schedule		
	ViTamin-S	ViTamin-B	ViTamin-L	ViTamin-L	ViTamin-XL	
batch size	8k	8k	16k	90k	90k	90k
image size	224	224	224	224	256	256
# A100 GPUs	32	32	32	184	312	312
# epochs	1	1	1	10	10	30
seen samples (B)	1.28	1.28	1.28	12.8	12.8	40.0
training days	1.8	3.3	5.6	11	15	46

Table 2. **Short and long training schedules on DataComp-1B.**

Locked-Text Tuning for CLIP: Besides model design, we propose Locked-Text Tuning (LTT) to exploits a pre-trained frozen text encoder. In light of the aligned image and text embeddings in CLIP, we leverage the pre-trained text encoder from a large VLM to guide the training of image encoders of smaller VLMs. Specifically, when training other ViTamin variants (*e.g.*, ViTamin-S and ViTamin-B), we initialize their text encoder with the one from a pre-trained ViTamin-L. The text encoder is then frozen, used as a teacher to guide the training of the randomly initialized image encoder. This scheme can be considered as a way to distill the knowledge [56] from a pre-trained frozen text encoder to a randomly initialized image encoder.

4. Experimental Results

In this section, we detail the implementation in Sec. 4.1, compare with the state-of-the-arts in Sec. 4.2, and deploy ViTamin to downstream tasks, including open-vocabulary detection/segmentation, and large multi-modal models in Sec. 4.3. See appendix for ablation studies.

4.1. Implementation Details

Training Strategy: We train the VLMs using OpenCLIP [63] on the public dataset DataComp-1B [41]. Tab. 2 summarizes the settings for our training schedules and

model variants. We use the short schedule to benchmark vision models and conduct our ablation studies, and long schedule to train our best ViTamin-L. We closely follow the training hyper-parameter settings in OpenCLIP [41, 63]. The training and fine-tuning details are in the appendix.

Evaluation Strategy: We follow DataComp [41] to zero-shot evaluate VLMs with a testbed of 38 tasks, including ImageNet [114], 6 distribution shift tasks [7, 53, 54, 111, 135], VTAB tasks [166], WILDS tasks [69, 115], and 3 retrieval tasks [9, 15, 153].

Other Downstream Tasks: We evaluate the trained VLM in downstream tasks. For open-vocabulary detection, we exploit the F-ViT framework [143], while for open-vocabulary segmentation, we adopt the FC-CLIP framework [159] and zero-shot evaluate on multiple segmentation datasets. Finally, we evaluate VLMs in LLaVA-1.5 [88] for LMMs across multiple benchmarks. In all the cases, F-ViT, FC-CLIP, and LLaVA employ the frozen VLM backbone to effectively ablate different pre-trained VLMs.

4.2. Main Results

Comparison with other State-of-the-arts: Tab. 3 summarizes the comparison between ViTamin-L and other state-of-the-art models, which exclusively employ the ViT backbone [32] but use different training schemes and datasets. For a fair comparison, we focus on the methods that use the same training data DataComp-1B [41], but still list other methods in the table for reference. For simplicity, we use “X@Z” to denote the vision model X trained with input size Z[‡]. ImageNet zero-shot accuracy is our main metric; other results are still reported in the table. As shown

[‡]Notation @ here is slightly abused to denote the training seen samples.

image encoder	image size	num patches	text encoder depth/width	seen samples (B)	training scheme	training dataset	trainable params Image+Text (M)	MACs Image+Text (G)	ImageNet Acc.	avg. 38 datasets	ImageNet dist. shift.	VTAB	retrieval
ViT-L/14 [41]	224	256	12 / 768	12.8	OpenCLIP	DataComp-1B	304.0 + 123.7	77.8 + 6.6	79.2	66.3	67.9	65.2	60.8
ViT-L/14 [80]	224	256	12 / 768	12.8 + 0.5	CLIPA-v2	DataComp-1B	304.0 + 110.3	77.8 + 2.7	79.7	65.4	68.6	62.9	60.6
ViT-L/14 [80]	336	576	12 / 768	12.8 + 0.5 + 0.1	CLIPA-v2	DataComp-1B	304.3 + 110.3	174.7 + 2.7	80.3	65.7	70.2	62.5	61.1
ViTamin-L	224	196	12 / 768	12.8	OpenCLIP	DataComp-1B	333.3 + 123.7	72.6 + 6.6	80.8	66.7	69.8	65.3	60.3
ViTamin-L	256 [†]	256	12 / 768	12.8 + 0.2	OpenCLIP	DataComp-1B	333.4 + 123.7	94.8 + 6.6	81.2	67.0	71.1	65.3	61.2
ViTamin-L	336	441	12 / 768	12.8 + 0.2	OpenCLIP	DataComp-1B	333.6 + 123.7	163.4 + 6.6	81.6	67.0	72.1	64.4	61.6
ViTamin-L	384 [†]	576	12 / 768	12.8 + 0.2	OpenCLIP	DataComp-1B	333.7 + 123.7	213.4 + 6.6	81.8	67.2	72.4	64.7	61.8
ViTamin-L2	224	196	24 / 1024	12.8	OpenCLIP	DataComp-1B	333.6 + 354.0	72.6 + 23.3	80.9	66.4	70.6	63.4	61.5
ViTamin-L2	256 [†]	256	24 / 1024	12.8 + 0.5	OpenCLIP	DataComp-1B	333.6 + 354.0	94.8 + 23.3	81.5	67.4	71.9	64.1	63.1
ViTamin-L2	336	441	24 / 1024	12.8 + 0.5	OpenCLIP	DataComp-1B	333.8 + 354.0	163.4 + 23.3	81.8	67.8	73.0	64.5	63.6
ViTamin-L2	384 [†]	576	24 / 1024	12.8 + 0.5	OpenCLIP	DataComp-1B	334.0 + 354.0	213.4 + 23.3	82.1	68.1	73.4	64.8	63.7
ViTamin-XL	256 [†]	256	27 / 1152	12.8 + 0.5	OpenCLIP	DataComp-1B	436.1 + 488.7	125.3 + 33.1	82.1	67.6	72.3	65.4	62.7
ViTamin-XL	384 [†]	576	27 / 1152	12.8 + 0.5	OpenCLIP	DataComp-1B	436.1 + 488.7	281.9 + 33.1	82.6	68.1	73.6	65.6	63.8
ViTamin-XL	256 [†]	256	27 / 1152	40.0	OpenCLIP	DataComp-1B	436.1 + 488.7	125.3 + 33.1	82.3	67.5	72.8	64.0	62.1
ViTamin-XL	336 [†]	441	27 / 1152	40.0 + 1.0	OpenCLIP	DataComp-1B	436.1 + 488.7	215.9 + 33.1	82.7	68.0	73.9	64.1	62.6
ViTamin-XL	384 [†]	576	27 / 1152	40.0 + 1.0	OpenCLIP	DataComp-1B	436.1 + 488.7	281.9 + 33.1	82.9	68.1	74.1	64.0	62.5
ViT-L/14 [123]	224	256	12 / 768	4.0	EVA-CLIP	Merged-2B	333.3 + 123.7	72.6 + 6.6	79.8	64.9	68.9	62.8	63.3
ViT-L/14 [123]	336	576	12 / 768	6.0	EVA-CLIP	Merged-2B	333.3 + 123.7	72.6 + 6.6	80.4	65.8	70.9	63.2	63.5
ViT-L/16 [169]	256	256	24 / 1024	40.0	SigLIP	WebLI	316.0 + 336.2	78.1 + 19.3	80.5	65.6	70.2	62.5	61.1
ViT-L/16 [169]	384	576	24 / 1024	40.0 + 5.0	SigLIP	WebLI	316.3 + 336.2	175.8 + 19.3	82.1	66.8	70.9	63.1	68.7
ViT-G/14 [63]	224	256	32 / 1280	39.0	OpenCLIP	LAION-2B	1844.9 + 694.7	473.4 + 48.5	80.1	66.7	69.1	64.6	63.5
ViT-H/14 [80]	336	576	24 / 1024	12.8 + 0.5 + 0.1	CLIPA-v2	DataComp-1B	632.5 + 354.0	363.7 + 9.7	81.8	66.8	72.4	63.7	62.6
ViT-E/14 [123]	224	256	24 / 1024	4.0	EVA-CLIP	LAION-2B	4350.6 + 354.0	1117.3 + 23.3	82.0	66.9	72.0	63.6	62.8
ViT-G/14 [80]	336	576	32 / 1280	12.8 + 0.5 + 0.1	CLIPA-v2	DataComp-1B	1845.4 + 672.3	1062.9 + 20.2	83.1	68.4	74.0	64.5	63.1
SoViT-400M/14 [2]	224	256	27 / 1152	40.0	SigLIP	WebLI	428.2 + 449.7	106.2 + 6.6	82.0	68.1	69.5	64.8	66.8
SoViT-400M/14 [2]	384	729	27 / 1152	40.0 + 5.0	SigLIP	WebLI	428.2 + 449.7	302.3 + 26.3	83.1	69.2	72.4	64.6	69.8
ViT-H/14 [37]	224	256	24 / 1024	39.0	OpenCLIP	DFN-5B	632.1 + 354.0	162.0 + 23.3	83.4	69.6	69.9	67.5	68.3
ViT-H/14 [37]	378	729	24 / 1024	39.0 + 5.0	OpenCLIP	DFN-5B	632.7 + 354.0	460.1 + 23.3	84.4	70.8	72.8	68.5	69.5

Table 3. **Comparison with state-of-the-art models.** Our models are only trained on the publicly available DataComp-1B [41]. CLIPA-v2 [80] uses an advanced progressive training scheme (from smaller images to larger ones) than the original OpenCLIP [41, 63] scheme that we follow. Other methods that use different settings are marked in gray for reference. Specifically, EVA-CLIP [123] uses EVA weights [38], better training scheme FLIP [85], and different training datasets [38, 117]. SigLIP [169] employs better sigmoid loss, stronger text encoders, and an extremely long schedule on the proprietary WebLI dataset [16] (40B for training and another 5B seen samples for fine-tuning). †: ViT-L/14 benefits from more image tokens by using a smaller output stride 14 than 16 that we use. To have the same image tokens, we slightly enlarge the image size (*e.g.*, $224/14 = 256/16$ and $336/14 = 384/16$). We note that all compared results are from the **OpenCLIP-results** that are evaluated under the same setting to ensure a fair comparison.

in the table, ViTamin-L@224 outperforms ViT-L/14@224 OpenCLIP [63] by +1.6%. However, ViT-L/14 benefits from more image tokens by using a smaller output stride 14 than 16 that we use (as benchmarked in the appendix). To have the same image tokens, we slightly enlarge the image size. As a result, our ViTamin-L@256 surpasses ViT-L/14@224 OpenCLIP [63] and CLIPA-v2 [80] by 2.0% and 1.5%, respectively. After fine-tuning on larger input sizes, ViTamin-L@384 and ViTamin-L@336 still performs better than ViT-L/14@336 CLIPA-v2 [80] by +1.5% and +1.3%, respectively. Impressively, with only half the parameters, our ViTamin-L attains an average of 67.2% performance across 38 datasets, exceeding the larger ViT-H/14 CLIPA-v2 model’s performance by +0.4%. Scaling up the text encoder to match the model size of the image encoder (specifically, ViTamin-L2) notably increases zero-shot ImageNet accuracy to 82.1% and average 38 datasets performance to 68.1%. Further scaling up the model parameters (*i.e.*, ViTamin-XL) and 40 billion seen samples reaches 82.9% zero-shot ImageNet accuracy.

Locked-Text Tuning: Fig. 4 shows that our LTT improves our ViTamin-S/-B by a large margin, especially when data sizes are small. Notably, LTT lifts ViTamin-B to the next scale of model performance, surpassing ViT-L/16 by +14% in 128M samples and +1.1% in 512M seen samples. Interestingly, LTT can save 10% training budget for ViTamin-B as the text tower is completely frozen.

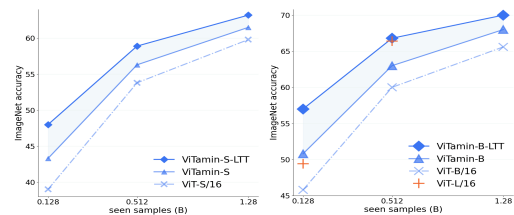


Figure 4. **Locked-text tuning (LTT).** LTT exploits a pretrained frozen text encoder, and effectively boosts the model performance.

Data Quality vs. Model Capacity: The DataComp challenge [41] underscores the role of data filtering for VLM, however, using a fixed ViT model. As shown in Tab. 4, the leading solution [154] of DataComp challenge in ICCV 2023 employed a complicated 24 filtering rules to improve the dataset quality, resulting in +2.3% gain. Surprisingly, our ViTamin-B improves the performance by a healthy margin of +12.8% accuracy, and locked-text tuning can lift the gain to +23.3%. The result highlights the importance to co-design the vision-language dataset and model.

4.3. New Suite of Downstream Tasks

The evaluations so far are mostly on classification/retrieval-based task, highlighting a lack of downstream tasks similar to those employed in the ImageNet era. Yet, in contrast to ImageNet-based vision models where downstream tasks mainly involve transfer learning for conventional de-

image encoder	data filtering	dataset size	seen samp.	IN acc. (%)	avg. 38 datasets (%)
leaderboard					
ViT-B/32	DataComp [41]	14M	128M	29.7	32.8
ViT-B/32	SIEVE [96]	24M	128M	30.3 (+0.6)	35.3 (+2.5)
ViT-B/32	Top-1 Solution [154]	23M	128M	32.0 (+2.3)	37.1 (+4.3)
our experiments					
ViT-B/32	DataComp [41]	14M	128M	29.4	31.5
ViT-B/16	DataComp [41]	14M	128M	35.8 (+6.4)	34.6 (+3.1)
ViTamin-B	DataComp [41]	14M	128M	42.2 (+12.8)	38.3 (+6.8)
ViT-B/16-LTT	DataComp [41]	14M	128M	43.6 (+14.2)	41.1 (+9.6)
ViTamin-B-LTT	DataComp [41]	14M	128M	52.7 (+23.3)	47.2 (+15.7)

Table 4. **Data quality vs. model capacity.** The leaderboard results are from ICCV 2023 DataComp challenge medium filtering track.

image encoder	pretraining		OV-COCO [165] (AP ₅₀ ^{novel})	OV-LVIS [47] (AP _r)
	dataset	scheme		
ViT-L/14	DataComp-1B	CLIPA-v2	36.1	32.5
ConvNeXt-L	LAION-2B	OpenCLIP	36.4	29.1
ViTamin-L	DataComp-1B	OpenCLIP	37.5	35.6

Table 5. **Open-vocabulary detection.** Different image encoders (ViT-L/14 by [80], ConvNeXt-L by [63]) are using the F-ViT framework [143] in a sliding window manner [160], trained on OV-COCO [165] and OV-LVIS [47]. ConvNeXt-L is marked in gray due to different pretrained dataset.

image encoder	pretraining		panoptic dataset (PQ)			semantic dataset (mIoU)				
	dataset	scheme	ADE [171]	Cityscapes [24]	MV [103]	A-150 [171]	A-847 [171]	PC-459 [100]	PC-59 [100]	PAS-21 [35]
ViT-L/14	DataComp-1B	CLIPA-v2	24.6	40.7	16.5	31.8	14.3	18.3	55.1	81.5
ConvNeXt-L	LAION-2B	OpenCLIP	26.8	44.0	18.3	34.1	14.8	18.2	58.4	81.8
ViTamin-L	DataComp-1B	OpenCLIP	27.3	44.0	18.2	35.6	16.1	20.4	58.4	83.4

Table 6. **Open-vocabulary segmentation.** Different image encoders (ViT-L/14 by [80], ConvNeXt-L by [63]) are using the FC-CLIP framework [159] in a sliding window manner [160], trained on COCO [87] and zero-shot evaluated on the other datasets. ConvNeXt-L is marked in gray due to different pretrained dataset.

tection and segmentation, VLMs excel with zero-shot capability and provides feature embeddings that are well-aligned across the vision-language domain. In light of this, we introduce a novel suite of downstream tasks aimed at the holistic evaluation of VLMs, including open-vocabulary detection and segmentation and multi-modal LLM.

Open-Vocabulary Detection and Segmentation: To examine how well the trained VLMs can adapt to downstream tasks, we consider two simple yet effective frameworks F-ViT [143] and FC-CLIP [159] which utilize a frozen CLIP backbone for open-vocabulary detection and segmentation, respectively. Specifically, we consider different VLMs as plug-in frozen backbones to these frameworks, while for ViT and ViTamin that may not easily generalize to high resolution input, we extract the feature in a sliding window manner [160], with window size equal to the pre-train image size, resulting in Sliding F-ViT and Sliding FC-CLIP, respectively. Tab. 5 illustrates that ViTamin-L serves as a stronger image encoder for open-vocabulary detection, surpassing its ViT-L/14 counterpart by 1.4% and 3.1% on OV-COCO and OV-LVIS. Tab. 6 shows that ViTamin-L outperforms ViT-L/14 by 2.6% on average 3 panoptic datasets and by 2.6% on average 5 semantic datasets. Notably, surpassing prior art, ViTamin-L sets a new state-of-the-art

image encoder	training scheme	VQAv2	GQA	VizWiz	SQA	T-VQA	POPE	MME	MMBench	MMB ^{C/N}	SEED	LLaVA ^V	MM-Vet
		[45]	[62]	[49]	[95]	[120]	[40]	[84]	[91]	[91]	[76]	[89]	[162]
ViT-L/14	OpenAI	78.5	62.0	50.0	66.8	58.2	85.9	1511	64.3	58.3	58.6	65.4	31.1
ViT-L/14	CLIPA-v2	75.9	60.3	48.8	65.6	55.0	84.9	1396	60.8	54.6	54.6	60.6	28.6
ViTamin-L	OpenCLIP	78.4	61.6	51.1	66.9	58.7	84.6	1421	65.4	58.4	57.7	64.5	33.6
ViTamin-L [†]	OpenCLIP	78.9	61.6	55.4	67.6	59.8	85.5	1447	64.5	58.3	57.9	66.1	33.6

Table 7. **Large Multi-modal Model (LMM) performance with different VLMs.** The results in 1st row originate from LLaVA-1.5 paper [88] and are marked in gray due to pretraining on OpenAI WIT dataset [108] unlike DataComp-1B [41] used by other rows. All listed models are trained following the same settings in LLaVA-1.5 [88] with Vicuna-V1.5-7B [20], for a fair comparison. †: image size of 384 rather than the default 336.

performance across seven benchmarks for open-vocabulary panoptic segmentation and semantic segmentation.

Large Multi-modal Models: Another key application of VLMs lies in their role as vision encoders within LMMs [77, 89, 175], as image features in VLMs that is well-aligned with text, thereby bridging the visual comprehension gap for LLMs. Specifically, we consider LLaVA-1.5 [88] as the evaluated framework. We follow [88] for all experimental settings, where the image is processed through a frozen CLIP model and a MLP projector, retaining the image as visual tokens, which are prepended to a text sequence and fed into a frozen Vicuna-v1.5-7B [20]. We run evaluation on 12 LMM benchmarks following [88], with results in Tab. 7. It should be noted that while OpenAI-trained ViT-L/14 underperforms CLIPAv2-trained counterpart by -3.7% ImageNet accuracy, it excels remarkably in LLaVA (+4.4% on VQAv2 and +4.3% on VizWiz). This highlights the need for incorporating a variety of downstream tasks to ensure a comprehensive evaluation. Surprisingly, simply replacing LLaVA’s image encoder to ViTamin-L can achieve new state-of-the-art across various benchmarks.

5. Conclusion

In this work, we build an evaluation protocols of modern vision models in VLM and re-benchmark them under CLIP setting. We examine vision models from four aspects of data scalability, model scalability, feature resolution and hybrid architecture. The four pillars motivate us to propose ViTamin, which not only competes favorably with ViT in zero-shot ImageNet accuracy and average 38 dataset accuracy, but also achieves the state-of-the-art on 22 downstream tasks covering open-vocabulary detection and segmentation and large multi-modal models. We hope that our design practices will drive the development of more advanced vision models for VLMs.

Acknowledgement: We thank Haichao Yu and Zhanpeng Zeng for the discussions about DataComp challenge and micro-level block design, respectively. This work was supported in part by ONR N00014-23-1-2641.

Appendix

In the supplementary materials, we provide additional information, as listed below.

- **Sec. A:** The ablation studies on the ViTamin macro-level network and micro-level block designs.
- **Sec. B:** ViTamin sets new SoTA in open-vocabulary dense prediction tasks including the OV-LVIS detection benchmark and 6 segmentation benchmarks.
- **Sec. C:** The results of using the proposed Locked-Text Tuning (LTT) training scheme.
- **Sec. D:** The results of benchmarking vision models under CLIP setting with an ImageNet-22K data scale.
- **Sec. E:** The numerical results of benchmarking vision models under CLIP setting with DataComp-1B.
- **Sec. F:** Detailed results of 38 datasets for different VLMs.
- **Sec. G:** The training hyper-parameter settings for short/long schedules and high-resolution input fine-tuning.

A. Ablation Studies

We conduct ablation studies on ViTamin design from two aspects: macro-level network and micro-level block. At the macro-level network design, we ablate the hybrid architecture and channel sizes of our three-stage network. At the micro-level block design, we ablate the design choices of convolution blocks and feed-forward network. In the tables, ‘IN acc.’ and ‘avg. 38’ denote the ImageNet accuracy (%) and the average accuracy (%) of 38 datasets, respectively. The ImageNet accuracy is used as the main metric. For simplicity, all the ablation studies are performed using base model variants with 128M seen samples.

Hybrid Architecture: In Tab. 8, we ablate design choices of hybrid architectures. Specifically, the compared architectures include ViT-B/16 (pure transformer with TFB or TFB-GeGLU blocks in stage 3), a new MBCConvNet-B (pure ConvNet with MBCConv-LN blocks in all three stages), and our ViTamin-B (MBCConv-LN in stage 1 and 2, and TFB-GeGLU in stage 3). The ablated models may differ in depth but share a similar number of parameters. As shown in the table, our ViTamin-B outperforms both the pure Transformer ViT-B/16 and the pure ConvNet MBCConvNet-B by more than +4.7%.

Channel Sizes of ViTamin: We ablate the effect of varying channel sizes within our ViTamin. The channel sizes (x_1C , x_2C , x_3C) denote the channel sizes of stage 1, 2, and 3, respectively. We set the channel size multipliers x_1 and x_2 to be 1 and 2 (commonly used in the literature for ImageNet). We ablate different values for x_3 in Tab. 9. Our final setting of (C , $2C$, $6C$) improves over (C , $2C$, $8C$) by +1.1%, and is on par with (C , $2C$, $4C$) but uses fewer parameters and MACs.

Design Choice of Convolution Blocks: In Tab. 10, we

model	# block type		depth stage 3	params (M)	IN acc. (%)	avg. 38 datasets
	stage 1 & 2	stage 3				
ViT-B/16	-	TFB	12	86.2	45.8	41.0
ViT-B/16	-	TFB-GeGLU	14	84.2	45.4	40.9
ViT-B/16	-	TFB-GeGLU	15	90.2	46.1	41.2
MBCConvNet-B	MBCConv-LN	MBCConv-LN	18	87.3	45.8	41.7
ViTamin-B	MBCConv-LN	TFB-GeGLU	14	87.5	50.8	44.6

Table 8. **Ablation study for hybrid architecture.** MBCConv-LN: Mobile Convolution with LayerNorm. TFB-GeGLU: Transformer Block with GeGLU. In this ablation study, we ablate TFB and TFB-GeGLU in ViT-B/16, and design a pure ConvNet using only MBCConv-LN across all three stages (called MBCConvNet-B in the table). Our final setting is marked in blue.

channel size	params (M)	MACs (G)	IN acc.	avg. 38
(C, 2C, 8C)	86.0	19.5	49.7	44.8
(C, 2C, 6C)	87.5	21.8	50.8	44.6
(C, 2C, 4C)	91.5	28.5	51.0	44.8

Table 9. **Ablation study on the channel sizes.** The channel sizes (x_1C , x_2C , x_3C) denote the channel sizes of stage 1, 2, and 3, respectively, w.r.t. a constant C (e.g., $(x_1, x_2, x_3) = (1, 2, 6)$ and $C = 128$ for ViTamin-B). Our final setting is marked in blue.

block type	params (M)	MACs (G)	IN acc.	avg. 38
ConvNeXt	88.0	21.0	49.8	44.9
MBCConv-BN	87.5	21.9	50.5	44.9
MBCConv-BN-SE	88.5	21.9	50.9	45.0
MBCConv-LN	87.5	21.8	50.8	44.6

Table 10. **Ablation study for design choice of convolutional blocks.** BN: BatchNorm. SE: Squeeze-and-Excitation. LN: LayerNorm. Our final setting is marked in blue.

ablate the design choices of convolution blocks in stage 1 and 2. The design choices include ConvNeXt, MBCConv-BN, MBCConv-BN-SE, and our MBCConv-LN. MBCConv-BN block is the original MBCConv block used in MobileNetv2 [116] with three BatchNorm layers [64], while the MBCConv-BN-SE block, proposed by MobileNetv3 [57], augments MBCConv-BN with the Squeeze-and-Excitation layer [59]. Each of the MBCConv variants demonstrates a superior performance to the ConvNeXt block [93]. Our MBCConv-LN, which employs a single Layer Normalization [4], outperforms the MBCConv-BN block, and achieves a similar result to MBCConv-BN-SE while requiring fewer parameters.

Design Choice of Feed-Forward Network: In Tab. 11, we study the effectiveness of GeGLU [118] in a Transformer Block (TFB) [133]. We experiment with ViT-B/16 and our ViTamin-B, and ablate on the effect of using the original TFB vs. the adopted TFB-GeGLU [118]. Remarkably, with the same depth of 12 blocks, ViTamin-B with GeGLU can achieve 49.9% accuracy and surpass the plain ViT-B/16 by a significant +4.1% margin and requires 13% fewer parameters. Adding two more blocks to align the parameters with ViT-B/16, our ViTamin-B boosts its performance to 50.8%, which not only improves over the GeGLU-

image encoder	GeGLU [118]	depth	params (M)	IN acc.	avg. 38
ViT-B/16		12	86.2	45.8	41.0
ViTamin-B		12	89.8	50.3	44.0
ViTamin-B	✓	12	75.7	49.9	43.5
ViTamin-B	✓	14	87.5	50.8	44.6
ViTamin-B		14	104.0	50.4	44.5

Table 11. **Ablation study for design choice of FFN.** Our final setting is marked in blue.

absent ViTamin-B counterpart (last row) by +0.4% but also maintains a reduced parameters by 26%.

B. Open-Vocabulary Dense Prediction

Frozen Feature Extraction via Sliding Window: We tested the transferability of VLMs to open-vocabulary detection tasks using F-ViT [143] and open-vocabulary segmentation tasks using FC-CLIP [159] frameworks, which both rely on a frozen CLIP backbone. The image size (*e.g.*, 1344×1344) for dense prediction tasks is usually larger than that of upstream VLM pre-training (*e.g.*, 224×224). To employ a frozen transformer-based architecture in these framework, we did not use any distillation [143] or convolutional backbone [159], while we find that a simple sliding window strategy [160] for frozen image feature extraction is effective enough to obtain reasonable performance on downstream tasks requiring high resolution input. The window size is the same as the input image size used during its VLM pretraining. We denote the slightly modified frameworks as *Sliding F-ViT* and *Sliding FC-CLIP*. We follow [143, 159] and use 896×896 and 1344×1344 input size for the open-vocabulary detection and segmentation tasks, respectively.

B.1. Open-Vocabulary Detection

In Tab.5 of main paper, ViTamin has been validated to be effective for open-vocabulary object detection on the OV-COCO dataset. In this section, we supplement the results on an additional benchmark OV-LVIS, where ViTamin sets a new state-of-the-art performance.

Experimental Setting: The open-vocabulary LVIS (OV-LVIS), introduced in ViLD [47], redefines the 337 rare categories from the LVIS v1.0 [48] dataset as novel categories. We strictly follow the F-ViT [143] framework to perform the open-vocabulary detection tasks, excepting the frozen image features are extracted in a sliding-window manner [160] (denoted as *Sliding F-ViT* in Tab. 13). The effectiveness of VLMs is validated through simply replacing the frozen backbone of F-ViT [143] framework. For evaluation, we follow previous works to use the mean mask AP on rare categories (AP_r) as the metric on OV-LVIS.

Results Analysis: Tab. 12 demonstrates that ViTamin-L is a stronger image encoder for open-vocabulary detector, surpassing its ViT-L/14 counterpart by 3.1% on OV-LVIS

image encoder	pretraining		OV-LVIS [47] (mAP_r)
	dataset	scheme	
ViT-L/14	DataComp-1B	CLIPA-v2	32.5
ConvNeXt-L	LAION-2B	OpenCLIP	29.1
ViTamin-L	DataComp-1B	OpenCLIP	35.6

Table 12. **Open-vocabulary detection.** Different image encoders (ViT-L/14 by [80] and ConvNeXt-L by [63]) are deployed using the F-ViT framework [143] in a sliding window manner [160], trained on OV-LVIS dataset [47]. ConvNeXt-L is marked in gray due to different pretrained dataset.

detector	image encoder	OV-LVIS (AP_r)	OV-COCO (AP_{50}^{novel})
ViLD [47]	RN50	16.6	27.6
OV-DETR [164]	RN50	17.4	29.4
DetPro [33]	RN50	19.8	-
OC-OVD [6]	RN50	21.1	36.6
OADP [138]	RN50	21.7	-
RegionCLIP [170]	RN50x4	22.0	-
CORA [144]	RN50x4	22.2	41.7
BARON-KD [142]	RN50	22.6	34.0
VLDet [86]	SwinB	26.3	-
F-VLM [73]	RN50x64	32.8	28.0
Detic [174]	SwinB	33.8	-
RO-ViT [67]	ViT-L/16	32.4	33.0
RO-ViT [67]	ViT-H/16	34.1	-
F-ViT [143]	ViT-L/14	24.2	24.7
F-ViT+CLIPSelf [143]	ViT-L/14	34.9	44.3
Sliding F-ViT	ViTamin-L	35.6	37.5

Table 13. **Comparison with prior arts** on open-vocabulary detection on OV-LVIS [47] and OV-COCO [165]. The last row (Sliding F-ViT) shows the result of employing our ViTamin-L using the F-ViT framework [143] in a sliding window manner [160].

dataset [47].

Comparison with Prior Arts: As shown in Tab. 13, ViTamin consistently outperforms all previous methods in the open-vocabulary detection task on OV-LVIS, setting a new state-of-the-art performance of 35.6% AP_r . Notably, our approach surpasses not only the distillation-based backbone (*e.g.*, CLIPSelf [143]) but also larger backbone (*e.g.*, ViT-H/16 in RO-ViT [67]).

B.2. Open-Vocabulary Segmentation

In Tab.6 of main paper, ViTamin has been validated to be effective for open-vocabulary panoptic and semantic segmentation on 8 dataset. We strictly follow the FC-CLIP framework [159] to perform the open-vocabulary segmentation tasks, excepting the frozen image features are extracted in a sliding-window manner [160] (denoted as *Sliding FC-CLIP* in Tab. 12). Following prior works [159], the *Sliding FC-CLIP* is trained on COCO [87] and zero-shot evaluated on the other datasets. In this section, we compare ViTamin with previous state-of-the-art methods.

Comparison with Prior Arts: In Tab. 14, our approach consistently outperforms all previous open-vocabulary seg-

method	image encoder	panoptic dataset (PQ)			semantic dataset (mIoU)				
		ADE	Cityscapes	MV	A-150	A-847	PC-459	PC-59	PAS-21
		[171]	[24]	[103]	[171]	[171]	[100]	[100]	[35]
FreeSeg [107]	-	16.3	-	-	-	-	-	-	-
OpenSeg [43]	-	-	-	-	21.1	6.3	9.0	42.1	-
GroupViT [148]	ViT-S/16	-	-	-	10.6	6.3	9.0	42.1	-
MaskCLIP [31]	ViT-B/16	15.1	-	-	23.7	8.2	10.0	45.9	-
ODISE [149]	-	22.2	23.9	14.2	29.9	11.1	14.5	57.3	84.6
FC-CLIP [159]	ConvNeXt-L	26.8	44.0	18.3	34.1	14.8	18.2	58.4	81.8
Sliding FC-CLIP	ViTamin-L	27.3	44.0	18.2	35.6	16.1	20.4	58.4	83.4

Table 14. **Comparison with prior arts** on open-vocabulary segmentation. ViTamin sets a new state-of-the-art result on various panoptic and semantic segmentation datasets. The last row (Sliding FC-CLIP) shows the result of employing our ViTamin-L using the FC-CLIP framework [159] in a sliding window manner [160].

mentation methods in 2 panoptic dataset and 4 semantic benchmarks, setting a new state-of-the-art. Notably, ViTamin surpasses the prior art by 0.5% PQ on ADE panoptic dataset and 1.5% mIoU on A-150 semantic dataset.

C. Locked-Text Tuning

Tab. 15 summarizes the detailed results of using the proposed new training scheme, Locked-Text Tuning (LTT). Specifically, when using the LTT training scheme, we employ the text encoder pretrained from ViTamin-L, and use it to guide the training of image encoders of ViTamin-S and ViTamin-B. As shown in the table, we consistently observe the improvements of using LTT. Compared to other distillation-based CLIP training schemes (See the rows marked in grey), our models achieve higher classification and retrieval accuracy in similar model parameters. Practically, despite being adopted from the larger model, the text encoder is much lighter compared to the image encoder (6.6 vs 21.8 GMACs), resulting in only a 14% increase in overall model MACs. Interestingly, using LTT results in a 10% savings in training costs for ViTamin-B, due to the text encoder being fully frozen.

D. Benchmarking Vision Models in CLIP with ImageNet-22K Data Scale

Tab. 16 summarizes the results of benchmarking vision models under CLIP setting with ImageNet-22K data scale. Specifically, we mimic the ImageNet-22K data scale by randomly selecting 14.2M data samples from DataComp-1B, and set the training epochs to 90, a standard training setting on ImageNet-22K. Similar to the findings on ImageNet-22K in the literature [93], under such a small data scale (14.2M data samples), ConvNeXt-T consistently outperforms ViT-S/32 and ViT-S/16. However, when the data scales up to 128M, or even 1.28B, the results are totally different, where ViT/16 shows a superior performance to ConvNeXt by a large margin, across all model sizes (see Tab. 17). We note that hybrid models, such as CoAtNet-0 and our ViTamin-S, still demonstrate the best

training scheme	training dataset	image encoder	params (M)	seen samp.	IN acc. (%)	avg. 38 (%)	retrieval COCO (%)
<i>models on private/other dataset, for reference</i>							
LiT [168]	Private-4B	ViT-B/32	86.2	0.9B	68.8	-	36.1
TinyCLIP [141]	LAION+YFCC	ViT-45M/32	45.0	1.6B	62.1	-	45.4
TinyCLIP [141]	LAION+YFCC	ViT-63M/32	63.0	1.6B	64.5	-	47.7
<i>our experiments</i>							
OpenCLIP	DataComp-1B	ViTamin-S	22.0	128M	43.3	40.8	25.8
OpenCLIP	DataComp-1B	ViTamin-S	22.0	512M	57.3	49.6	36.6
OpenCLIP	DataComp-1B	ViTamin-S	22.0	1.28B	62.2	53.2	40.2
OpenCLIP	DataComp-1B	ViTamin-B	87.5	128M	50.8	44.6	31.2
OpenCLIP	DataComp-1B	ViTamin-B	87.5	512M	64.0	53.9	41.7
OpenCLIP	DataComp-1B	ViTamin-B	87.5	1.28B	68.9	57.7	44.9
LTT (ours)	DataComp-1B	ViTamin-S	22.0	128M	47.5	44.8	33.4
LTT (ours)	DataComp-1B	ViTamin-S	22.0	512M	58.9	52.0	41.6
LTT (ours)	DataComp-1B	ViTamin-S	22.0	1.28B	63.4	54.6	45.0
LTT (ours)	DataComp-1B	ViTamin-B	87.5	128M	56.7	50.5	39.8
LTT (ours)	DataComp-1B	ViTamin-B	87.5	512M	66.8	57.3	47.1
LTT (ours)	DataComp-1B	ViTamin-B	87.5	1.28B	70.8	59.4	50.0

Table 15. **Locked-Text Tuning (LTT) training scheme.** We use the pretrained text encoder from ViTamin-L and train the image encoders of ViTamin-{S,B}. Due to the use of private or other filtered/merged dataset, the results borrowed from LiT [168] and TinyCLIP [141] are just for reference, and LiT [168] reports retrieval on COCO only. †: a filtered subset of WebLI dataset [16].

image encoder	data size (M)	epoch	#params (M)	MACs (G)	ImageNet Acc.(%)	avg. 38 datasets (%)
<i>ImageNet-22K scale</i>						
ViT-S/32	14.2	90	21.81	1.12	39.4	36.7
ViT-S/16	14.2	90	21.81	4.25	45.7	38.7
ConvNeXt-T	14.2	90	28.61	4.47	45.9	39.3
CoAtNet-0	14.2	90	24.56	4.43	49.1	41.4
ViTamin-S	14.2	90	22.03	5.50	50.3	41.3

Table 16. **Benchmarking vision models under CLIP setting with an ImageNet-22K data scale.** We mimic the ImageNet-22K data scale with 14.2M data size and 90 training epochs (standard training setting on ImageNet-22K). The benchmarked vision models include ViT (pure transformer), ConvNeXt (pure convolution), CoAtNet (hybrid model), and our proposed ViTamin.

performances under this small data scale, showing that the hybrid design works well across all data sizes.

E. Numerical Results of Benchmarking Vision Models with DataComp-1B

In Fig.2 of the main paper, we provide the analysis of benchmarked results from various aspects. In this section, we further supplement the numerical results of benchmarking vision models (including ViT, ConvNeXt, CoAtNet, and our ViTamin) across different model scales and data sizes in Tab. 17. As shown in the table, the proposed ViTamin consistently outperforms all the other vision models in almost all settings.

F. Results of 38 dataset for different VLMs.

Tab. 18 demonstrates the detailed results for VLMs with different large-variant image encoders. This table is associated with Tab. 3 of the main paper.

image encoder	image size	text encoder depth / width	seen samples	#params (M) image+text	MACs (G) image+text	ImageNet Acc. (%)	avg. 38 datasets	ImageNet dist. shift.	VTAB	Retrieval
<i>small model variants</i>										
ViT-S/32	224	12 / 384	128M	21.81 + 40.44	1.12 + 1.64	32.1	34.1	25.3	35.8	27.2
ViT-S/16	224	12 / 384	128M	21.81 + 40.44	4.25 + 1.64	38.4	36.7	29.3	38.5	31.3
ConvNeXt-T	224	12 / 384	128M	28.61 + 40.44	4.47 + 1.64	37.0	35.8	30.1	37.5	30.8
CoAtNet-0	224	12 / 384	128M	24.56 + 40.44	4.43 + 1.64	42.4	38.9	33.5	39.9	34.2
ViTamin-S	224	12 / 384	128M	22.03 + 40.44	5.50 + 1.64	43.3	40.8	35.6	41.0	35.2
ViT-S/32	224	12 / 384	512M	21.81 + 40.44	1.12 + 1.64	47.3	44.3	36.7	46.9	36.9
ViT-S/16	224	12 / 384	512M	21.81 + 40.44	4.25 + 1.64	53.8	46.5	41.9	46.7	42.1
ConvNeXt-T	224	12 / 384	512M	28.61 + 40.44	4.47 + 1.64	53.3	46.1	42.4	46.4	42.2
CoAtNet-0	224	12 / 384	512M	24.56 + 40.44	4.43 + 1.64	56.4	49.0	45.2	49.0	45.0
ViTamin-S	224	12 / 384	512M	22.03 + 40.44	5.50 + 1.64	57.3	49.6	46.9	48.8	45.4
ViT-S/32	224	12 / 384	1.28B	21.81 + 40.44	1.12 + 1.64	53.0	47.3	41.3	48.9	41.5
ViT-S/16	224	12 / 384	1.28B	21.81 + 40.44	4.25 + 1.64	59.8	50.9	47.2	51.3	47.5
ConvNeXt-T	224	12 / 384	1.28B	28.61 + 40.44	4.47 + 1.64	59.9	51.3	47.8	52.7	48.3
CoAtNet-0	224	12 / 384	1.28B	24.56 + 40.44	4.43 + 1.64	61.7	51.6	50.1	51.3	48.9
ViTamin-S	224	12 / 384	1.28B	22.03 + 40.44	5.50 + 1.64	62.2	53.2	51.3	51.7	50.0
<i>base model variants</i>										
ViT-B/32	224	12 / 512	128M	86.19 + 63.43	4.37 + 2.91	38.9	38.0	30.6	40.6	30.7
ViT-B/16	224	12 / 512	128M	86.19 + 63.43	16.87 + 2.91	45.8	41.0	35.8	42.1	36.2
ConvNeXt-B	224	12 / 512	128M	88.09 + 63.43	15.38 + 2.91	41.4	39.7	33.5	41.2	34.1
CoAtNet-2	224	12 / 512	128M	74.18 + 63.43	15.94 + 2.91	48.5	43.5	38.9	43.8	39.1
ViTamin-B	224	12 / 512	128M	87.53 + 63.43	21.84 + 2.91	50.8	44.6	41.3	45.1	40.8
ViT-B/32	224	12 / 512	512M	86.19 + 63.43	4.37 + 2.91	54.8	48.3	42.7	50.1	42.4
ViT-B/16	224	12 / 512	512M	86.19 + 63.43	16.87 + 2.91	60.0	51.0	48.2	51.4	47.5
ConvNeXt-B	224	12 / 512	512M	88.09 + 63.43	15.38 + 2.91	59.4	50.3	47.9	49.9	47.2
CoAtNet-2	224	12 / 512	512M	74.18 + 63.43	15.94 + 2.91	63.3	52.4	52.4	51.0	49.7
ViTamin-B	224	12 / 512	512M	87.53 + 63.43	21.84 + 2.91	64.0	53.9	53.3	53.7	50.8
ViT-B/32	224	12 / 512	1.28B	86.19 + 63.43	4.37 + 2.91	60.1	52.5	47.4	53.6	47.5
ViT-B/16	224	12 / 512	1.28B	86.19 + 63.43	16.87 + 2.91	65.6	55.6	53.1	55.3	51.7
ConvNeXt-B	224	12 / 512	1.28B	88.09 + 63.43	15.38 + 2.91	65.3	54.7	54.0	54.2	51.7
CoAtNet-2	224	12 / 512	1.28B	74.18 + 63.43	15.94 + 2.91	68.5	56.8	57.2	56.0	53.4
ViTamin-B	224	12 / 512	1.28B	87.53 + 63.43	21.84 + 2.91	68.9	57.7	58.3	56.4	54.1
<i>large model variants</i>										
ViT-L/32	224	12 / 768	128M	303.97 + 123.65	15.27 + 6.55	43.5	40.8	34.0	42.7	34.2
ViT-L/16	224	12 / 768	128M	303.97 + 123.65	59.70 + 6.55	49.4	43.8	38.7	44.3	38.9
ViT-L/14	224	12 / 768	128M	303.97 + 123.65	77.83 + 6.55	49.9	43.8	39.4	44.5	39.3
ConvNeXt-XL	224	12 / 768	128M	350.25 + 123.65	79.65 + 6.55	42.8	38.4	33.3	38.4	35.0
CoAtNet-4	224	12 / 768	128M	275.07 + 123.65	60.81 + 6.55	52.5	45.2	42.0	45.2	41.1
ViTamin-L	224	12 / 768	128M	333.32 + 123.65	72.60 + 6.55	52.7	44.8	42.4	44.6	41.8
ViT-L/32	224	12 / 768	512M	303.97 + 123.65	15.27 + 6.55	60.4	51.8	47.4	52.7	47.3
ViT-L/16	224	12 / 768	512M	303.97 + 123.65	59.70 + 6.55	66.4	55.6	53.6	55.5	52.2
ViT-L/14	224	12 / 768	512M	303.97 + 123.65	77.83 + 6.55	67.0	55.4	54.8	54.2	52.0
ConvNeXt-XL	224	12 / 768	512M	350.25 + 123.65	79.65 + 6.55	63.0	52.5	51.1	51.8	49.4
CoAtNet-4	224	12 / 768	512M	275.07 + 123.65	60.81 + 6.55	66.8	56.1	56.4	56.5	50.4
ViTamin-L	224	12 / 768	512M	333.32 + 123.65	72.60 + 6.55	68.7	56.6	56.8	56.5	53.2
ViT-L/32	224	12 / 768	1.28B	303.97 + 123.65	15.27 + 6.55	67.5	57.0	54.1	57.9	51.9
ViT-L/16	224	12 / 768	1.28B	303.97 + 123.65	59.70 + 6.55	71.9	60.1	59.9	59.9	56.0
ViT-L/14	224	12 / 768	1.28B	303.97 + 123.65	77.83 + 6.55	72.3	60.7	60.5	60.0	56.0
ConvNeXt-XL	224	12 / 768	1.28B	350.25 + 123.65	79.65 + 6.55	70.2	58.3	59.1	57.0	55.5
CoAtNet-4	224	12 / 768	1.28B	275.07 + 123.65	60.81 + 6.55	71.3	59.4	61.4	59.1	53.4
ViTamin-L	224	12 / 768	1.28B	333.32 + 123.65	72.60 + 6.55	73.9	62.0	62.9	61.4	56.6

Table 17. **Benchmarking vision backbones on Datacomp-1B under CLIP setting (contrastive language-image pretraining).** We benchmark popular vision backbones, including ViT [32] (pure transformer model), ConvNeXt [93] (pure convolution model), CoAtNet [25] (hybrid convolution-transformer model), and our proposed ViTamin, under different model parameters and training seen samples.

image encoder	training scheme	avg. 38	ImageNet 1k [28]	Caltech-101 [39]	CIFAR-10 [71]	CIFAR-100 [71]	CLEVR Counts [66]	CLEVR Distance [71]	Country211 [108]	Describable Textures [22]	EuroSAT [52]	FGVC Aircraft [97]	Food-101 [10]	GTSRB [122]	ImageNet Sketch [135]	ImageNet v2 [111]	ImageNet-A [54]	ImageNet-O [54]	ImageNet-R [53]	KFTT Vehicle Distance [42]	MINIST [74]	ObjectNet [7]	Oxford Flowers-102 [104]	Oxford-IIT Pet [106]	Pascal VOC 2007 [35]	PatchCamelyon [134]	Rendered SST2 [166]	RESISC45 [166]	Stanford Cars [70]	STL-10 [23]	SUN397 [145]	SVHN [102]	Fllickr [153]	MSCOCO [15]	WinoGAVIL [9]	iWildCam [8]	Camelyon17 [5]	FMoW [21]	Dollar Street [112]	GeoDE [110]
ViT-L/14 [63]		66.3	79.2	94.7	98.2	87.3	35.6	24.4	31.6	66.5	71.2	47.5	94.5	58.5	68.0	72.1	69.6	32.6	90.8	27.9	86.6	74.3	82.6	95.1	82.5	51.2	61.0	69.4	93.1	99.3	74.3	67.7	81.2	54.5	46.7	16.1	50.9	24.0	66.2	91.5
ViT-L/14 [80]		65.4	79.6	94.5	98.7	88.5	18.6	24.5	29.4	69.6	60.4	43.0	94.2	59.1	70.6	73.0	71.2	33.7	92.9	19.3	73.7	69.9	81.0	95.0	80.7	59.2	53.9	68.4	93.7	99.2	75.3	63.9	81.9	56.0	43.9	17.2	67.6	24.6	66.5	91.5
ViT-L/14 † [80]		65.7	80.3	94.4	98.6	88.3	15.7	24.4	30.7	68.6	58.1	42.8	94.6	57.0	70.9	73.5	77.7	32.9	93.3	20.0	76.7	73.2	81.0	95.0	79.8	60.3	53.2	68.8	94.1	99.3	75.6	62.9	82.5	56.4	44.5	19.4	67.8	25.0	67.5	92.4
ViTamin-L [63]		66.7	80.8	95.1	98.5	88.0	35.4	24.1	33.0	68.1	66.6	49.2	94.9	61.9	71.5	73.6	72.4	32.0	93.0	23.6	81.8	76.1	85.0	95.6	80.2	52.2	61.1	73.0	94.8	99.3	75.8	66.1	81.1	54.7	45.1	17.2	58.0	16.4	68.1	91.5
ViTamin-L † [63]		67.2	81.8	95.6	98.5	87.8	31.7	24.1	36.0	69.2	64.7	49.7	95.8	63.4	72.1	75.2	81.7	31.1	93.8	21.2	81.3	80.7	84.7	95.8	82.0	50.1	60.5	73.1	95.1	99.5	76.3	66.9	82.5	55.7	47.3	19.1	49.5	17.5	70.7	92.2
ViTamin-XL † [63]		68.1	82.6	95.7	98.7	88.8	19.1	20.0	37.8	71.5	75.6	53.7	96.0	53.2	73.1	76.3	83.1	33.0	94.2	17.4	88.9	81.9	85.6	95.9	83.8	56.2	61.9	75.9	94.8	99.4	76.2	74.0	84.7	58.7	47.9	21.4	46.0	22.5	68.6	92.5

Table 18. **Detailed results of 38 dataset for different VLMs.** The compared models are trained with the scheme of either OpenCLIP [63] or CLIPA-v2 [80]. All models are trained on DataComp-1B [41] dataset with similar seen samples for a fair comparison. †: using larger number of patches of 576 (*i.e.*, image size of 336 for row 3 and 384 for row 5, respectively).

training config	<i>short schedule</i>	<i>long schedule</i>
	ViTamin-S/B/L 224 ²	ViTamin-L/L2/XL/XL 224 ² /224 ² /256 ² /256 ²
batch size	8k/8k/16k	90k
seen samples	1.28B	12.8B/12.8B/12.8B/40B
optimizer	AdamW	AdamW
base learning rate	5e-4	2e-3
weight decay	0.02	0.02
optimizer momentum β_1	0.9	0.9
optimizer momentum β_2	0.98/0.98/0.95	0.95
learning rate schedule	cosine decay	cosine decay
warmup steps	500	782/4436/4436/9981
warmup schedule	linear	linear
random crop ratio	none	[0.4, 1.0]
stochastic depth [60]	0.1	0.1
precision	amp bfloat16	amp bfloat16

Table 19. **Short/Long schedule training settings for ViTamin variants.**

pre-training config	ViTamin-L 224 ²	ViTamin-L2 224 ²	ViTamin-XL 256 ²
fine-tuning config	256 ² /336 ² /384 ²	256 ² /336 ² /384 ²	256 ² /384 ²
batch size	90k	90k	90k
seen samples	0.2B	0.5B	0.5B
optimizer	AdamW	AdamW	AdamW
base learning rate	1e-5	1e-5	1e-5
weight decay	0	0	0
optimizer momentum β_1	0.9	0.9	0.9
optimizer momentum β_2	0.95	0.95	0.95
learning rate schedule	constant	constant	constant
warmup steps	0	0	0
random crop ratio	none	none	none
stochastic depth [60]	0.1	0.1	0.1
precision	amp bfloat16	amp bfloat16	amp bfloat16

Table 20. **Fine-tuning setting for high resolution.** The models are pre-trained with *long schedule* and then fine-tuned on the target resolution.

G. Training Hyper-parameter Settings

Tab. 19 and Tab. 20 provide our details of training hyper-parameter settings for short/long schedules and fine-tuning for high resolution, respectively. The short schedule is used to benchmark several vision models on DataComp-1B,

along with our ablation studies, while the long schedule is used to train our ViTamin-L for better performances. When fine-tuning the trained model on larger input resolution, we fine-tune with only 200M seen samples and a small constant learning rate.

References

- [1] Gpt-4v(ision) system card. 2023. 3
- [2] Ibrahim Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. Getting vit in shape: Scaling laws for compute-optimal model design. *NeurIPS*, 2023. 7
- [3] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *NeurIPS*, 2022. 3
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 5, 6, 9
- [5] Peter Bandi, Oscar Geessink, Quirine Manson, Marcory Van Dijk, Maschenka Balkenhol, Meyke Hermsen, Babak Ehteshami Bejnordi, Byungjae Lee, Kyunghyun Paeng, Aoxiao Zhong, et al. From detection of individual metastases to classification of lymph node status at the patient level: the camelyon17 challenge. *IEEE transactions on medical imaging*, 38(2):550–560, 2018. 13
- [6] Hanoona Bangalath, Muhammad Maaz, Muhammad Uzair Khattak, Salman H Khan, and Fahad Shahbaz Khan. Bridging the gap between object and image-level representations for open-vocabulary detection. *NeurIPS*, 2022. 10
- [7] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *NeurIPS*, 2019. 6, 13
- [8] Sara Beery, Arushi Agarwal, Elijah Cole, and Vighnesh Birodkar. The iwildcam 2021 competition dataset. *arXiv preprint arXiv:2105.03494*, 2021. 13
- [9] Yonatan Bitton, Nitzan Bitton Guetta, Ron Yosef, Yuval Elovici, Mohit Bansal, Gabriel Stanovsky, and Roy Schwartz. Winogavil: Gamified association benchmark to

- challenge vision-and-language models. *NeurIPS*, 2022. [6](#), [13](#)
- [10] Lukas Bossard, Matthieu Guillaumin, and Luc Van Gool. Food-101—mining discriminative components with random forests. In *ECCV*, 2014. [13](#)
- [11] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *NeurIPS*, 2020. [3](#)
- [12] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *ICLR*, 2015. [1](#)
- [13] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *TPAMI*, 40(4):834–848, 2017. [1](#)
- [14] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. [1](#)
- [15] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. [6](#), [13](#)
- [16] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. *arXiv preprint arXiv:2209.06794*, 2022. [7](#), [11](#)
- [17] Xi Chen, Xiao Wang, Lucas Beyer, Alexander Kolesnikov, Jialin Wu, Paul Voigtlaender, Basil Mustafa, Sebastian Goodman, Ibrahim Alabdulmohsin, Piotr Padlewski, et al. Pali-3 vision language models: Smaller, faster, stronger. *arXiv preprint arXiv:2310.09199*, 2023. [3](#)
- [18] Yinpeng Chen, Xiyang Dai, Dongdong Chen, Mengchen Liu, Xiaoyi Dong, Lu Yuan, and Zicheng Liu. Mobileformer: Bridging mobilenet and transformer. In *CVPR*, 2022. [2](#)
- [19] Bowen Cheng, Maxwell D Collins, Yukun Zhu, Ting Liu, Thomas S Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-DeepLab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *CVPR*, 2020. [1](#)
- [20] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, 2023. [8](#)
- [21] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In *CVPR*, 2018. [13](#)
- [22] Mircea Cimpoi, Subhansu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *CVPR*, 2014. [13](#)
- [23] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *AISTATS*, 2011. [13](#)
- [24] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *CVPR*, 2016. [8](#), [11](#)
- [25] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. *NeurIPS*, 2021. [1](#), [2](#), [3](#), [5](#), [12](#)
- [26] Yann N Dauphin, Angela Fan, Michael Auli, and David Grangier. Language modeling with gated convolutional networks. In *ICML*, 2017. [5](#)
- [27] Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer, Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling vision transformers to 22 billion parameters. In *ICML*, 2023. [2](#)
- [28] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. [13](#)
- [29] Xueqing Deng, Qihang Yu, Peng Wang, Xiaohui Shen, and Liang-Chieh Chen. Coconut: Modernizing coco segmentation. In *CVPR*, 2024. [2](#)
- [30] Zheng Ding, Jieke Wang, and Zhuowen Tu. Open-vocabulary universal image segmentation with maskclip. In *ICML*, 2023. [2](#)
- [31] Xiaoyi Dong, Jianmin Bao, Yinglin Zheng, Ting Zhang, Dongdong Chen, Hao Yang, Ming Zeng, Weiming Zhang, Lu Yuan, Dong Chen, et al. Maskclip: Masked self-distillation advances contrastive language-image pretraining. In *CVPR*, 2023. [11](#)
- [32] 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. In *ICLR*, 2021. [1](#), [2](#), [3](#), [5](#), [6](#), [12](#)
- [33] Yu Du, Fangyun Wei, Zihe Zhang, Miaoqing Shi, Yue Gao, and Guoqi Li. Learning to prompt for open-vocabulary object detection with vision-language model. In *CVPR*, 2022. [10](#)
- [34] Stéphane d’Ascoli, Hugo Touvron, Matthew L Leavitt, Ari S Morcos, Giulio Biroli, and Levent Sagun. Convit: Improving vision transformers with soft convolutional inductive biases. In *ICML*, 2021. [2](#)
- [35] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88:303–338, 2010. [8](#), [11](#), [13](#)
- [36] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021. [2](#)
- [37] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023. [7](#)

- [38] Yuxin Fang, Wen Wang, Binhui Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *CVPR*, 2023. 7
- [39] Li Fei-Fei, Rob Fergus, and Pietro Perona. Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. In *CVPRW*, 2004. 13
- [40] Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Zhenyu Qiu, Wei Lin, Jinrui Yang, Xiawu Zheng, et al. Mme: A comprehensive evaluation benchmark for multimodal large language models. *arXiv preprint arXiv:2306.13394*, 2023. 8
- [41] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv preprint arXiv:2304.14108*, 2023. 1, 2, 3, 6, 7, 8, 13
- [42] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 13
- [43] Golnaz Ghiasi, Xiuye Gu, Yin Cui, and Tsung-Yi Lin. Scaling open-vocabulary image segmentation with image-level labels. In *ECCV*, 2022. 2, 3, 11
- [44] Ross Girshick. Fast r-cnn. In *CVPR*, 2015. 1
- [45] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In *CVPR*, 2017. 8
- [46] Benjamin Graham, Alaaeldin El-Nouby, Hugo Touvron, Pierre Stock, Armand Joulin, Hervé Jégou, and Matthijs Douze. Levit: a vision transformer in convnet’s clothing for faster inference. In *CVPR*, 2021. 2
- [47] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. In *ICLR*, 2022. 8, 10
- [48] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 10
- [49] Danna Gurari, Qing Li, Abigale J Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and Jeffrey P Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In *CVPR*, 2018. 8
- [50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 2, 3, 5
- [51] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *CVPR*, 2017. 1
- [52] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2019. 13
- [53] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In *ICCV*, 2021. 6, 13
- [54] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In *CVPR*, 2021. 6, 13
- [55] Byeongho Heo, Sangdoo Yun, Dongyoon Han, Sanghyuk Chun, Junsuk Choe, and Seong Joon Oh. Rethinking spatial dimensions of vision transformers. In *ICCV*, 2021. 2
- [56] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 6
- [57] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *ICCV*, 2019. 5, 9
- [58] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. 2, 5
- [59] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 5, 9
- [60] Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016. 13
- [61] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *CVPR*, 2017. 1, 2
- [62] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019. 8
- [63] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hananeh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. 2, 3, 6, 7, 8, 10, 13
- [64] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5, 9
- [65] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *ICML*, 2021. 1, 2, 3
- [66] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, 2017. 13
- [67] Dahun Kim, Anelia Angelova, and Weicheng Kuo. Region-aware pretraining for open-vocabulary object detection with vision transformers. In *CVPR*, 2023. 10
- [68] Alexander Kirillov, Kaiming He, Ross Girshick, Carsten Rother, and Piotr Dollár. Panoptic segmentation. In *CVPR*, 2019. 1

- [69] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *ICML*, 2021. 6
- [70] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In *ICCVW*, 2013. 13
- [71] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 13
- [72] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *NeurIPS*, 2012. 1, 2
- [73] Weicheng Kuo, Yin Cui, Xiuye Gu, AJ Piergiovanni, and Anelia Angelova. F-vm: Open-vocabulary object detection upon frozen vision and language models. In *ICLR*, 2023. 2, 3, 10
- [74] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998. 13
- [75] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 1
- [76] Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Benchmarking multimodal llms with generative comprehension. *arXiv preprint arXiv:2307.16125*, 2023. 8
- [77] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022. 2, 3, 8
- [78] Kunchang Li, Yali Wang, Junhao Zhang, Peng Gao, Guanglu Song, Yu Liu, Hongsheng Li, and Yu Qiao. Uniformer: Unifying convolution and self-attention for visual recognition. *TPAMI*, 2023. 2
- [79] Runze Li, Dahun Kim, Bir Bhanu, and Weicheng Kuo. Reclip: Resource-efficient clip by training with small images. *arXiv preprint arXiv:2304.06028*, 2023. 3
- [80] Xianhang Li, Zeyu Wang, and Cihang Xie. Scaling clip training with 81.1% zero-shot imagenet accuracy within a \$10,000 budget; an extra \$4,000 unlocks 81.8% accuracy. *arXiv preprint arXiv:2306.15658*, 2023. 1, 2, 7, 8, 10, 13
- [81] Xianhang Li, Zeyu Wang, and Cihang Xie. An inverse scaling law for clip training. *NeurIPS*, 2023. 2, 3
- [82] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Mvitv2: Improved multiscale vision transformers for classification and detection. In *CVPR*, 2022. 2
- [83] Yanyu Li, Geng Yuan, Yang Wen, Ju Hu, Georgios Evangelidis, Sergey Tulyakov, Yanzhi Wang, and Jian Ren. Efficientformer: Vision transformers at mobilenet speed. *NeurIPS*, 2022. 2
- [84] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*, 2023. 8
- [85] Yanghao Li, Haoqi Fan, Ronghang Hu, Christoph Feichtenhofer, and Kaiming He. Scaling language-image pre-training via masking. In *CVPR*, 2023. 3, 7
- [86] Chuang Lin, Peize Sun, Yi Jiang, Ping Luo, Lizhen Qu, Gholamreza Haffari, Zehuan Yuan, and Jianfei Cai. Learning object-language alignments for open-vocabulary object detection. *arXiv preprint arXiv:2211.14843*, 2022. 10
- [87] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 1, 8, 10
- [88] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 1, 2, 3, 6, 8
- [89] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *NeurIPS*, 2023. 2, 3, 8
- [90] Haotian Liu, Kilho Son, Jianwei Yang, Ce Liu, Jianfeng Gao, Yong Jae Lee, and Chunyuan Li. Learning customized visual models with retrieval-augmented knowledge. In *CVPR*, 2023. 3
- [91] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multimodal model an all-around player? *arXiv preprint arXiv:2307.06281*, 2023. 8
- [92] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, 2021. 1, 2
- [93] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. In *CVPR*, 2022. 1, 2, 3, 9, 11, 12
- [94] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1
- [95] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *NeurIPS*, 2022. 8
- [96] Anas Mahmoud, Mostafa Elhoushi, Amro Abbas, Yu Yang, Newsha Ardalani, Hugh Leather, and Ari S Morcos. Sieve: Multimodal dataset pruning using image captioning models. *arXiv preprint arXiv:2310.02110*, 2023. 8
- [97] Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013. 13
- [98] Sachin Mehta and Mohammad Rastegari. Mobilevit: lightweight, general-purpose, and mobile-friendly vision transformer. *arXiv preprint arXiv:2110.02178*, 2021. 2
- [99] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple open-vocabulary object detection. In *ECCV*, 2022. 2, 3

- [100] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, 2014. 8, 11
- [101] Norman Mu, Alexander Kirillov, David Wagner, and Saining Xie. Slip: Self-supervision meets language-image pre-training. In *ECCV*, 2022. 3
- [102] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bis-sacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 13
- [103] Gerhard Neuhof, Tobias Ollmann, Samuel Rota Bulo, and Peter Kontschieder. The mapillary vistas dataset for semantic understanding of street scenes. In *ICCV*, 2017. 8, 11
- [104] Maria-Elena Nilsback and Andrew Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian conference on computer vision, graphics & image processing*, pages 722–729. IEEE, 2008. 13
- [105] OpenAI. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 2, 3
- [106] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In *CVPR*, 2012. 13
- [107] Jie Qin, Jie Wu, Pengxiang Yan, Ming Li, Ren Yuxi, Xue-feng Xiao, Yitong Wang, Rui Wang, Shilei Wen, Xin Pan, et al. Freeseq: Unified, universal and open-vocabulary image segmentation. In *CVPR*, 2023. 11
- [108] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language super-vision. In *ICML*, 2021. 1, 2, 3, 8, 13
- [109] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jon Shlens. Stand-alone self-attention in vision models. *NeurIPS*, 2019. 2
- [110] Vikram V Ramaswamy, Sing Yu Lin, Dora Zhao, Aaron B Adcock, Laurens van der Maaten, Deepti Ghadiyaram, and Olga Russakovsky. Beyond web-scraping: Crowd-sourcing a geographically diverse image dataset. *arXiv preprint arXiv:2301.02560*, 2023. 13
- [111] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to im-agenet? In *ICML*, 2019. 6, 13
- [112] William A Gaviria Rojas, Sudnya Diamos, Keertan Ranjan Kini, David Kanter, Vijay Janapa Reddi, and Cody Cole-man. The dollar street dataset: Images representing the geo-graphic and socioeconomic diversity of the world. *NeurIPS*, 2022. 13
- [113] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 3
- [114] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115 (3):211–252, 2015. 1, 2, 6
- [115] Shiori Sagawa, Pang Wei Koh, Tony Lee, Irena Gao, Sang Michael Xie, Kendrick Shen, Ananya Kumar, Weihua Hu, Michihiro Yasunaga, Henrik Marklund, et al. Extending the wilds benchmark for unsupervised adaptation. In *ICLR*, 2022. 6
- [116] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018. 2, 5, 6, 9
- [117] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Worts-man, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *NeurIPS*, 2022. 2, 3, 7
- [118] Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020. 5, 6, 9, 10
- [119] Karen Simonyan and Andrew Zisserman. Very deep convo-lutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2
- [120] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *CVPR*, 2019. 8
- [121] Aravind Srinivas, Tsung-Yi Lin, Niki Parmar, Jonathon Shlens, Pieter Abbeel, and Ashish Vaswani. Bottleneck transformers for visual recognition. In *CVPR*, 2021. 2
- [122] Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The german traffic sign recognition bench-mark: a multi-class classification competition. In *IJCNN*, 2011. 13
- [123] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023. 2, 3, 7
- [124] Shuyang Sun, Weijun Wang, Andrew Howard, Qihang Yu, Philip Torr, and Liang-Chieh Chen. Remax: Relaxing for better training on efficient panoptic segmentation. *NeurIPS*, 2024. 1
- [125] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 2, 5
- [126] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019. 2
- [127] Mingxing Tan and Quoc Le. Efficientnetv2: Smaller mod-els and faster training. In *ICML*, 2021. 2
- [128] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 2
- [129] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles, Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *ICCV*, 2021. 5
- [130] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Bap-tiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language mod-els. *arXiv preprint arXiv:2302.13971*, 2023. 3

- [131] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. **2**
- [132] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. Maxvit: Multi-axis vision transformer. In *ECCV*, 2022. **1, 2**
- [133] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017. **1, 2, 5, 6, 9**
- [134] Bastiaan S Veeling, Jasper Linmans, Jim Winkens, Taco Cohen, and Max Welling. Rotation equivariant cnns for digital pathology. In *MICCAI*, 2018. **13**
- [135] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *NeurIPS*, 2019. **6, 13**
- [136] Huiyu Wang, Yukun Zhu, Bradley Green, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Axial-DeepLab: Stand-Alone Axial-Attention for Panoptic Segmentation. In *ECCV*, 2020. **2, 3**
- [137] Huiyu Wang, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Max-deeplab: End-to-end panoptic segmentation with mask transformers. In *CVPR*, 2021. **1, 3**
- [138] Luting Wang, Yi Liu, Penghui Du, Zihan Ding, Yue Liao, Qiaosong Qi, Bialong Chen, and Si Liu. Object-aware distillation pyramid for open-vocabulary object detection. In *CVPR*, 2023. **10**
- [139] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*, 2021. **1, 2**
- [140] Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, Xiyang Dai, Lu Yuan, and Lei Zhang. Cvt: Introducing convolutions to vision transformers. In *ICCV*, 2021. **2**
- [141] Kan Wu, Houwen Peng, Zhenghong Zhou, Bin Xiao, Mengchen Liu, Lu Yuan, Hong Xuan, Michael Valenzuela, Xi Stephen Chen, Xinggang Wang, et al. Tinyclip: Clip distillation via affinity mimicking and weight inheritance. In *ICCV*, 2023. **3, 11**
- [142] Size Wu, Wenwei Zhang, Sheng Jin, Wentao Liu, and Chen Change Loy. Aligning bag of regions for open-vocabulary object detection. In *CVPR*, 2023. **10**
- [143] Size Wu, Wenwei Zhang, Lumin Xu, Sheng Jin, Xiangtai Li, Wentao Liu, and Chen Change Loy. Clipself: Vision transformer distills itself for open-vocabulary dense prediction. *arXiv preprint arXiv:2310.01403*, 2023. **1, 6, 8, 10**
- [144] Xiaoshi Wu, Feng Zhu, Rui Zhao, and Hongsheng Li. Cora: Adapting clip for open-vocabulary detection with region prompting and anchor pre-matching. In *CVPR*, 2023. **10**
- [145] Jianxiong Xiao, Krista A Ehinger, James Hays, Antonio Torralba, and Aude Oliva. Sun database: Exploring a large collection of scene categories. *IJCV*, 119:3–22, 2016. **13**
- [146] Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. *NeurIPS*, 2021. **2**
- [147] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *ICCV*, 2017. **2**
- [148] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *CVPR*, 2022. **11**
- [149] Jiarui Xu, Sifei Liu, Arash Vahdat, Wonmin Byeon, Xiaolong Wang, and Shalini De Mello. Open-vocabulary panoptic segmentation with text-to-image diffusion models. In *CVPR*, 2023. **2, 11**
- [150] Mengde Xu, Zheng Zhang, Fangyun Wei, Yutong Lin, Yue Cao, Han Hu, and Xiang Bai. A simple baseline for zero-shot semantic segmentation with pre-trained vision-language model. In *ECCV*, 2022. **2**
- [151] Chenglin Yang, Siyuan Qiao, Qihang Yu, Xiaoding Yuan, Yukun Zhu, Alan Yuille, Hartwig Adam, and Liang-Chieh Chen. Moat: Alternating mobile convolution and attention brings strong vision models. In *ICLR*, 2023. **1, 2**
- [152] Xuan Yang, Liangzhe Yuan, Kimberly Wilber, Astuti Sharma, Xiuye Gu, Siyuan Qiao, Stephanie Debats, Huisheng Wang, Hartwig Adam, Mikhail Sirotenko, and Liang-Chieh Chen. Polymax: General dense prediction with mask transformer. *arXiv preprint arXiv:2311.05770*, 2024. **1**
- [153] Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. *Transactions of the Association for Computational Linguistics*, 2:67–78, 2014. **6, 13**
- [154] Haichao Yu, Yu Tian, Sateesh Kumar, Linjie Yang, and Heng Wang. The devil is in the details: A deep dive into the rabbit hole of data filtering. *arXiv preprint arXiv:2309.15954*, 2023. **2, 7, 8**
- [155] Jiahui Yu, Zirui Wang, Vijay Vasudevan, Legg Yeung, Mojtaba Seyedhosseini, and Yonghui Wu. Coca: Contrastive captioners are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022. **3**
- [156] Qihang Yu, Yingda Xia, Yutong Bai, Yongyi Lu, Alan L Yuille, and Wei Shen. Glance-and-gaze vision transformer. *NeurIPS*, 2021. **1**
- [157] Qihang Yu, Huiyu Wang, Dahun Kim, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. Cmt-deeplab: Clustering mask transformers for panoptic segmentation. In *CVPR*, 2022. **1**
- [158] Qihang Yu, Huiyu Wang, Siyuan Qiao, Maxwell Collins, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. k-means Mask Transformer. In *ECCV*, 2022. **1**
- [159] Qihang Yu, Ju He, Xueqing Deng, Xiaohui Shen, and Liang-Chieh Chen. Convolutions die hard: Open-vocabulary segmentation with single frozen convolutional clip. *NeurIPS*, 2023. **1, 2, 3, 6, 8, 10, 11**
- [160] Qihang Yu, Xiaohui Shen, and Liang-Chieh Chen. Towards open-ended visual recognition with large language model. *arXiv preprint arXiv:2311.08400*, 2023. **8, 10, 11**
- [161] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan.

- Metaformer is actually what you need for vision. In *CVPR*, 2022. 2
- [162] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023. 8
- [163] Xiaowei Yu, Yao Xue, Lu Zhang, Li Wang, Tianming Liu, and Dajiang Zhu. Exploring the influence of information entropy change in learning systems. *arXiv preprint arXiv:2309.10625*, 2023. 2
- [164] Yuhang Zang, Wei Li, Kaiyang Zhou, Chen Huang, and Chen Change Loy. Open-vocabulary detr with conditional matching. In *ECCV*, 2022. 10
- [165] Alireza Zareian, Kevin Dela Rosa, Derek Hao Hu, and Shih-Fu Chang. Open-vocabulary object detection using captions. In *CVPR*, 2021. 2, 8, 10
- [166] Xiaohua Zhai, Joan Puigcerver, Alexander Kolesnikov, Pierre Ruysen, Carlos Riquelme, Mario Lucic, Josip Djolonga, Andre Susano Pinto, Maxim Neumann, Alexey Dosovitskiy, et al. The visual task adaptation benchmark. 2019. 6, 13
- [167] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. In *CVPR*, 2022. 5
- [168] Xiaohua Zhai, Xiao Wang, Basil Mustafa, Andreas Steiner, Daniel Keysers, Alexander Kolesnikov, and Lucas Beyer. Lit: Zero-shot transfer with locked-image text tuning. In *CVPR*, 2022. 3, 11
- [169] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, 2023. 3, 7
- [170] Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li, Luwei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. Regionclip: Region-based language-image pretraining. In *CVPR*, 2022. 10
- [171] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *CVPR*, 2017. 8, 11
- [172] Chong Zhou, Chen Change Loy, and Bo Dai. Extract free dense labels from clip. In *ECCV*, 2022. 2
- [173] Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. Deepvit: Towards deeper vision transformer. *arXiv preprint arXiv:2103.11886*, 2021. 2, 5
- [174] Xingyi Zhou, Rohit Girdhar, Armand Joulin, Philipp Krähenbühl, and Ishan Misra. Detecting twenty-thousand classes using image-level supervision. In *ECCV*, 2022. 10
- [175] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 2, 3, 8