
INTERN: A New Learning Paradigm Towards General Vision

Jing Shao* Siyu Chen* Yangguang Li* Kun Wang* Zhenfei Yin* Yinan He* Jianing Teng*
 Qinghong Sun* Mengya Gao* Jihao Liu* Gengshi Huang* Guanglu Song Yichao Wu
 Yuming Huang Fenggang Liu Huan Peng Shuo Qin Chengyu Wang Yujie Wang Conghui He
 Ding Liang Yu Liu Fengwei Yu Junjie Yan Dahua Lin Xiaogang Wang Yu Qiao†

Shanghai AI Laboratory SenseTime Research
 The Chinese University of Hong Kong Shanghai Jiao Tong University

Abstract

Enormous waves of technological innovations over the past several years, marked by the advances in AI technologies, are profoundly reshaping the industry and the society. However, down the road, a key challenge awaits us, that is, our capability of meeting rapidly-growing scenario-specific demands is severely limited by the cost of acquiring a commensurate amount of training data. This difficult situation is in essence due to limitations of the mainstream learning paradigm: we need to train a new model for each new scenario, based on a large quantity of well-annotated data and commonly from scratch. In tackling this fundamental problem, we move beyond and develop a new learning paradigm named INTERN. By learning with supervisory signals from multiple sources in multiple stages, the model being trained will develop strong generalizability. We evaluate our model on 26 well-known datasets that cover four categories of tasks in computer vision. In most cases, our models, adapted with only 10% of the training data in the target domain, outperform the counterparts trained with the full set of data, often by a significant margin. This is an important step towards a promising prospect where such a model with general vision capability can dramatically reduce our reliance on data, thus expediting the adoption of AI technologies. Furthermore, revolving around our new paradigm, we also introduce a new data system, a new architecture, and a new benchmark, which, together, form a general vision ecosystem to support its future development in an open and inclusive manner.

1 Introduction

Current state-of-the-art AI models typically over-specialize on a single task, despite remarkable progress in recent years. The consequence is that we develop thousands of models for thousands of tasks or circumstances individually. Each new task requires collecting and annotating a large amount of data, and consumes a giant scale of computational resources. From [2, 3], this presents itself as a significant hurdle in front of AI researches and applications, considering thousands of long-tailed tasks in industries. Alternatively, the artificial general intelligence approach, taking “general intelligence” as a fundamentally distinct property [19], focuses directly on the *generality*, *adaptability*, and *flexibility* of AI models.

Vision and language are two indispensable modalities for artificial general intelligence. For language, impressive progress has been achieved towards general language model (GLM). Recent advances of large-scale pretrained language models such as BERT [15], T5 [31] and GPT-3 [3] have shown potential in developing GLMs that substantially benefit a wide range of language-related downstream tasks by allowing economical task-specific adaptations. Moreover, with

*Equal Contribution

†Corresponding Author: qiaoyu@pjlab.org.cn

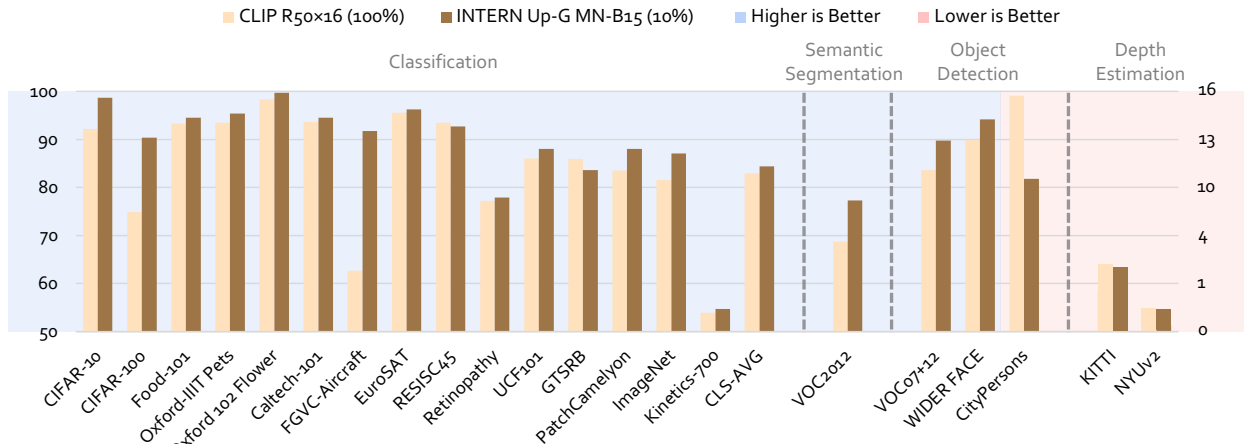


Figure 1: **Comparison of transfer learning performance on diverse tasks.** Our largest pretrained model, Up-G MN-B15, with 90% fewer downstream data, surpasses the best publicly available pretrained model (CLIP-R50 \times 16) on most tasks. All results are obtained with backbone parameters fixed during downstream training. Note that the last three datasets in the pink background use the y -axis on the right, and the lower bar means the better performance.

the advent of task-agnostic training objectives [41, 15], performance gains from pretraining can be steadily improved by scaling up web-crawled data and the model capacity together with computational budgets.

The success of GLMs has inspired new directions for general vision model (GVM) learning. Pioneers working on large-scale supervised [42, 32, 13, 36, 14], self-supervised [10, 21, 8, 11, 1, 7], and cross-modal [23, 30] pretraining have shown certain generality on a limited scope of downstream vision tasks. Nevertheless, it is still challenging to design reliable approaches towards GVMs. Most previous works mainly utilize one source of supervisory signal, *e.g.* ViT-G/14 [42] uses categorical supervision, SEER [20] applies contrastive learning between different augmentation views, and CLIP [30] makes use of paired language descriptions. Pretraining monotonically with an individual supervision is able to produce models performing well in selected scenarios, but cannot offer sufficient competence if we aim at obtaining a “true” GVM that is *generalizable* to a vast set of downstream tasks, even unseen ones. To achieve generality with respect to diverse vision tasks, it is favorable to learn abundant information from miscellaneous types of supervisory signals, including image-level categories, bounding boxes, pixel-wise labels, quantities, as well as natural language.

In this work, we propose a new learning paradigm named INTERN, further pushing forward to successful general vision modeling. Specifically, INTERN introduces a *continuous learning* scheme (see Fig. 2), including a highly extensible upstream pretraining pipeline leveraging large-scale data and various supervisory signals, as well as a flexible downstream adaptation towards diversified tasks.

For an analogy of the upstream pretraining procedure, one may look no further than a typical learning process of an “intern” in real life, which can be roughly divided into the three subsequent stages based on the level of expertise:

- I: An *amateur* with fundamental skill sets who can superficially address encountered problems;
- II: An *expert* who has additionally mastered one particular task with careful supervision;
- III: A *generalist* who is knowledgeable about all known tasks, and adapts fast to unseen tasks.

We demonstrate that it is beneficial to imitate this amateur-to-generalist learning process of an “intern” to ease the non-trivial training of the desired GVM. The resulting multi-stage pretraining scheme not only efficiently absorbs knowledge from broad sources of supervisions, but also is easily scalable with the presence of more data or tasks. Based on this pipeline, our final generalists prove to possess sufficient generality towards a wide range of downstream tasks, and outperform (see Fig. 1) the previous state-of-the-art (CLIP [30]) while only using an order of magnitude fewer (10%) downstream data.

In addition to upstream pretraining, we also introduce a downstream adaptation step. The key challenge is to form downstream task-specific models while largely preserving the merits of an upstream GVM. We show that designing a proxy for flexible knowledge transfer onto various downstream tasks is a promising attempt towards this goal. Moreover, establishing a comprehensive benchmark is essential to propel GVM development. A GVM is expected to not only generalize to different tasks but also have lower requirements on downstream data, achieving few-shot adaptation.

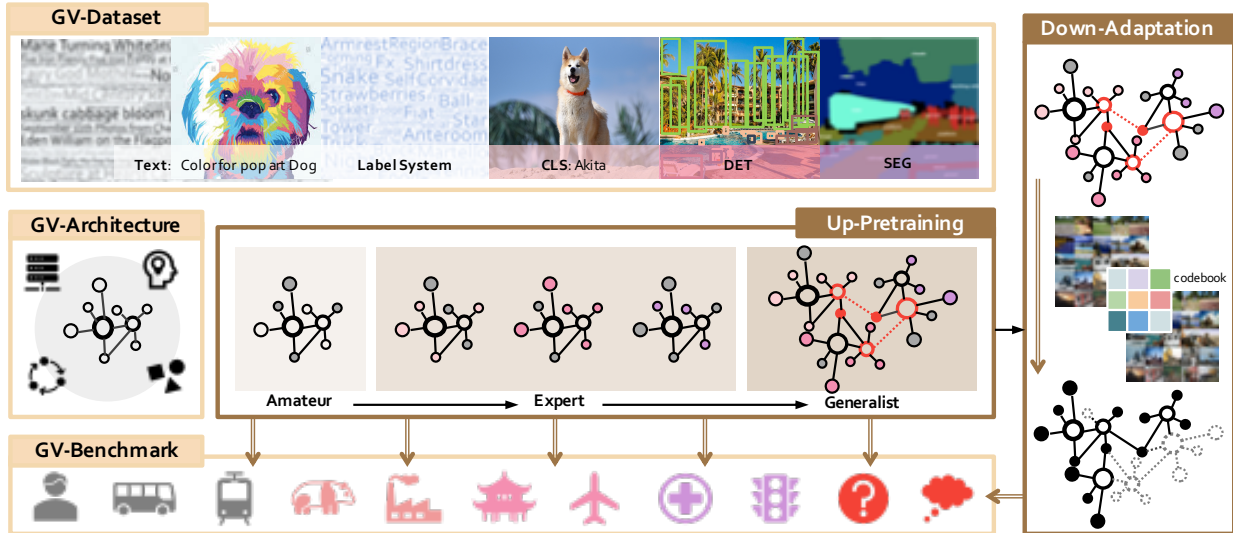


Figure 2: **Overview of INTERN.** Our complete flow of learning and evaluating a general vision model consists of three fundamental bases (*i.e.* GV-Dataset, GV-Architecture, and GV-Benchmark), a three-stage upstream pretraining scheme (*i.e.* Amateur, Expert, and Generalist), and a downstream adaptation algorithm that transfers the up-pretrained models to various downstream tasks in the benchmark. It shows that a general model (*e.g.* Generalist) with a continuous learning process exhibits stronger generalizability even on unseen tasks (shown in a red question mark).

To this end, our general vision benchmark in INTERN contains a wide range of common downstream tasks, and systematically examines important factors relating to GVMs, especially generality and data efficiency.

1.1 Overview of INTERN

As shown in Fig. 2, INTERN consists of seven key components. Three of them serve as fundamental bases: General Vision Dataset is constructed as the database for the upstream step-wise learning process of INTERN. General Vision Architecture is the backbone of INTERN models. General Vision Benchmark consists of a broad range of downstream datasets and evaluation settings to assess the generalization ability of vision models thoroughly. The other four refer to our upstream pretraining scheme with three stages (*i.e.* Upstream-Amateur, Upstream-Expert, and Upstream-Generalist) and Downstream-Adaptation, which provides a refined solution to adapting learned upstream general models towards various types of downstream tasks.

- **General Vision Data (GV-D)** is a super-scale collection of vision datasets with 10 billion samples and various supervisory signals. It presents an extensive label system with 119K visual concepts, covering numerous realms in nature and almost all labels currently studied in computer vision. Guided by this comprehensive label system, we conduct four high-quality general vision datasets actively and continuously. Precisely, GV-D consists of GV-D-10B with multi-modal data, as well as GV-D_c-36M, GV-D_d-3M and GV-D_s-143K with classification, detection, and segmentation annotations, respectively.
- **General Vision Architecture (GV-A)** introduces a set of network architectures with higher modeling capacity, which is constructed from a unified search space with both convolution and transformer operators. We name this automatically assembled and high-performing vision network family as *MetaNet*.
- **General Vision Benchmark (GV-B)** collects 26 downstream tasks consisting of 4 task types, on which models produced by our INTERN paradigm are benchmarked along with publicly released pretrained models for comparison. In addition, GV-B introduces a *percentage-shot* setting where the amount of training data of downstream tasks is shrunk by only taking a portion of the full dataset, such as 10%, 20% and etc. Compared to traditional few-shot settings, our *percentage-shot* setup can well preserve properties like long-tailed distribution of the original dataset and alleviate the sensitivity to sample selection.
- **Upstream-Amateur (Up-A)** is a multi-modal pretraining stage for acquiring the amateur model, which simultaneously uses rich supervisory signals from image-text, image-image, and text-text pairs to train task-agnostic models serving as initialization for the next stage.

- **Upstream-Expert (Up-E)** is the following stage in our pretraining scheme for achieving the expert model, which accumulates specialized knowledge with multi-source supervisions within one of the task types. Each expert only pays attention to its own specialty without interfering with the learning of others.
- **Upstream-Generalist (Up-G)** is a combinational pretraining stage resulting in the generalist model, which integrates knowledge of experts and produces the final form of general representation capable of handling any known or unknown task.
- **Downstream-Adaptation (Down-A)** introduces a transfer learning scheme aimed at enhancing knowledge transfer onto various downstream task types, which can be applied to any upstream pretrained model. It effectively improves the results of adaptation, especially in the low-data regime.

2 Core Results

We first highlight the high performance of models learned with our multi-stage pretraining scheme at downstream transfer learning, especially in the data-efficient (10%) setting. We then evaluate the generality of our paradigm and its extensibility in the face of new tasks. Experiments are conducted on ResNet [22] and the proposed MetaNet. MetaNet-B4 (MN-B4) shares a similar parameter complexity with ResNet-50 (R50), and MetaNet-B15 (MN-B15) is a considerably larger model with $\sim 1\text{B}$ parameters. Transfer learning performances of these models are assessed on GV-B. During the evaluation, we fix pretrained backbone weights and only fine-tune the output head. This setting is previously known as *linear probe* for classification tasks, and we extend this term to other tasks in our benchmark.

2.1 Superior Transfer Learning Performance with Only 10% Training Data

Our goal is to train a general vision model capable of effectively transferring to various downstream tasks with significantly fewer data and annotations. As displayed in Tab. 7, we assess the transferability of our proposed models on a wide range of tasks, covering general and fine-grained classification, object detection, semantic segmentation, and depth estimation. Both low-data (10%) and full-data (100%) settings are considered for appraising data efficiency. We also make comparison with a list of popular large-scale pretraining approaches on our benchmark, including supervised [25, 28], self-supervised [9, 6, 5, 38] and cross-modal [30] approaches with a variety of network architectures [22, 39, 16].

From Tab. 1, we observe that generalist models from our multi-stage pretraining achieve state-of-the-art results on most tasks even with only 10% downstream training data. Take ResNet-50 [22] as an example, in the low-data (10%) regime, Up-G’s performances surpass those of the ImageNet [33] pretrained ResNet-50 by large margins, with +11.5% average accuracy on the classification suite, +18.2% AP on VOC detection, and +8.2% mIoU on semantic segmentation. It is notable that while the pretrained ResNet-50 from CLIP [30] shows a gain of +4.1% on classification, its performance degrades on VOC detection and segmentation tasks. The results of ResNet-50 checkpoints from MoCo v2 [9] and SwAV [6] are roughly similar, both showing unsteady performances across tasks compared to the ImageNet baseline. Therefore, our approach is the only pretraining scheme that leads to superior transfer capabilities on all evaluated tasks. In addition, it is worth mentioning that our ResNet-50 in the 10% setting also steadily outperforms the ImageNet pretrained one with 100% data in terms of all metrics. This further demonstrates the great downstream data efficiency of our models.

For our proposed MetaNet, MN-B4 consistently outperforms its similar-complexity rival, ResNet-50. This observation also reflects INTERN’s good compatibility with different types of backbone networks. Furthermore, our largest generalist model, MN-B15, achieves state-of-the-art performance on our GV-B benchmark in most dimensions, surpassing the best publicly available pretrained model, CLIP-R50 $\times 16$ [30]. It is also worth noting that compared to CLIP-R50 $\times 16$ [30] in the low-data regime, our approach relatively decreases the mean error rate by 40.2% on the classification suite, 52.8% on VOC detection, 45.0% on CityPersons (R), 44.2% on WIDER FACE (M), 34.8% on VOC segmentation, 6.9% on KITTI and 11.9% on NYUv2, respectively.

2.2 Easy Extensibility and Great Generalizability

INTERN is a general-purpose pretraining paradigm suitable for any task or model architecture. In this section, we select the ResNet-50 backbone to demonstrate the extensibility and generalizability of our approach for a fair comparison with previous works.

Effective integration of diverse knowledge. The role of Up-E, *i.e.* expert training stage, is to learn representation for each specialized task type such as image classification, object detection, and semantic segmentation. This design naturally mitigates the learning difficulty caused by task conflicts [12] in multi-task learning setups. The representation learned for each task type performs better than simple ImageNet pretraining when tested on other tasks of the same type [29] due to

Model	Data Setting	CLS-AVG \uparrow	VOC-DET \uparrow	CP (R) \downarrow	WF (M) \uparrow	VOC-SEG \uparrow	KITTI \downarrow	NYUv2 \downarrow
ImageNet [33] R50	10%	62.8	69.5	29.6	84.2	58.0	3.26	0.48
	100%	73.0	79.5	22.7	86.8	66.0	3.09	0.43
MoCo v2 [9] R50	10%	61.1	70.2	23.2	86.3	60.1	3.13	0.46
SwAV [6] R50	10%	64.3	69.2	23.2	86.4	56.8	2.95	0.46
CLIP [30] R50	10%	66.9	68.6	21.3	87.1	55.8	3.15	0.48
CLIP [30] R50 \times 16	10%	75.0	78.4	19.1	89.6	65.2	2.91	0.42
	100%	82.9	83.6	16.2	89.9	68.7	2.83	0.39
Up-G R50	10%	74.3	87.7	14.7	92.2	66.2	2.84	0.39
Up-G MN-B4	10%	78.6	89.1	12.0	92.8	71.4	2.94	0.40
Up-G MN-B15	10%	84.4	89.8	10.5	94.2	77.3	2.71	0.37

Table 1: **Linear probe performance of selected pretrained models.** CLS-AVG denotes average accuracy on 20 classification datasets, CP (R) denotes MR^{-2} of the *reasonable* setup in CityPersons and WF (M) denotes AP50 of the *medium* subset in WIDER FACE. Our Up-G R50 outperforms the ImageNet pretrained version on classification, object detection, semantic segmentation and depth estimation tasks while trained with 90% fewer downstream data.

Pretrain	Data Setting	CLS-AVG \uparrow	VOC-DET \uparrow
ImageNet	100%	73.0	79.5
Up-E (C)	10%	73.7	72.2
Up-E (D)	10%	53.9	87.7
Up-G (C-D)	10%	74.3	87.7

Table 2: **Effectiveness of diverse knowledge integration.** Our classification and detection expert models, denoted by Up-E (C) and Up-E (D) respectively, with only 10% training data, both surpass the ImageNet supervised baseline with full data on their corresponding tasks. Transferability of these two experts are successfully preserved by our generalist, Up-G (C-D), achieving +0.6% accuracy on classification. The red color denotes that the model has been pretrained on the corresponding task.

strengthened knowledge specifically for that task type. Benefiting from GV-D’s large-scale data with rich multi-domain information, Up-E’s quality is improved to a greater extent and it becomes worthy of its name, expert. For example, in Tab. 2, downstream performances on corresponding tasks with only 10% training data of the experts for classification and detection (referred to as Up-E (C) and Up-E (D) respectively) beat those of the ImageNet baseline trained with full data, showing great few-shot capabilities. Next, we consider how to obtain a general representation across different task types. In Up-G, we introduce a simple yet effective *knowledge transfer module* for gluing the experts together. This module lies between any two experts, learning complementary information for the target expert from the source expert. Such a knowledge integration mechanism over all experts results in a universal representation instead of specific ones. In Tab. 2, Up-G (C-D) successfully preserves capabilities of both the classification and detection experts, and we observe an extra gain of +0.6% on average classification accuracy. The results imply that our pretraining approach is able to effectively consolidate diverse knowledge from different task types.

Easily extensible with more experts. To build a robust pretraining system, we always keep extensibility in mind. Our Up-G is a simple multi-task learning approach where experts of new tasks can be easily added. We already show that Up-G (C-D) matches or excels its two experts in terms of transferability on both tasks. As presented in Tab. 3, when introduced with a semantic segmentation expert Up-E (S), we simply extend Up-G (C-D) to Up-G (C-D-S) which further acquires the few-shot ability from the new expert without dropping performances on classification and detection. Moreover, we observe a gain of +1.8% on segmentation compared to Up-E (S), showing again the value of multiple experts in our pretraining system.

Robust representation for generalization to unseen tasks. As a generalist, it is required to be proficient at all known tasks and has a fast adaptation ability to unseen tasks with few samples. We show that pretrained models based on our pipeline meet the goal of possessing sufficient generality towards unknown downstream tasks. As displayed in Tab. 4, with only 10% training data, our Up-G (C-D) achieves 66.2% mIoU on VOC segmentation and 2.84 RMSE on KITTI depth estimation. Both metrics on the two unseen tasks are better than those of the ImageNet supervised model with 100% data. As for Up-G (C-D-S), besides the significant gain on segmentation compared to Up-G (C-D), we see slightly stronger performance on the unseen depth estimation as well. This generalizability is not achieved by any of the single expert models. Specifically, either Up-E (C) or Up-E (D) alone considerably lags behind Up-G (C-D) in

Pretrain	CLS-AVG \uparrow	VOC-DET \uparrow	VOC-SEG \uparrow
Up-E (C)	73.7	72.2	57.7
Up-E (D)	53.9	87.7	62.3
Up-E (S)	47.5	75.0	71.9
Up-G (C-D)	74.3	87.7	66.2
Up-G (C-D-S)	74.3	87.7	73.7

Table 3: **Extensibility of generalist.** Based on Up-G (C-D), Up-G (C-D-S) is easily obtained by additionally linking to Up-E (S). With the new expert added, original performances are not harmed.

Pretrain	Data Setting	VOC-SEG \uparrow	KITTI \downarrow
ImageNet	100%	66.0	3.09
Up-E (C)	10%	57.7	3.21
Up-E (D)	10%	62.3	3.09
Up-G (C-D)	10%	66.2	2.84
Up-G (C-D-S)	10%	73.7	2.80

Table 4: **Generalizability to unseen tasks.** Up-G (C-D) outperforms the ImageNet pretrained ResNet-50 on VOC segmentation and KITTI depth estimation with 90% fewer training data. Up-G (C-D-S) also performs well on the depth estimation task. Up-E (C) or Up-E (D) alone fails to match the performances of the generalist Up-G (C-D).

terms of VOC segmentation and KITTI depth estimation. We attribute this improvement achieved by our generalists to the diverse supervisory signals applied during our multi-stage pretraining. Our pipeline effectively exploits information embedded in multiple tasks, leading to a more robust and generalizable model. Our experiments may further shed light on learning complementary information from a multi-task setting.

2.3 Factors Contributing to General Vision Intelligence

First, *scaling up dataset size, broadening domains, and diversifying supervisory signals matter.* In INTERN, GV-D acts as a foundational database for models to comprehensively “study” in the pretraining step. Most existing datasets are constrained within a narrow computer vision task type, *e.g.* image classification, object detection, etc. We argue that this is insufficient for general vision pretraining which requires wide coverage over various domains. Although some multi-modal or supervised pretraining approaches have used billion-scale data, the performances of resulting models on some structured tasks like object detection and segmentation are unsatisfactory. This indicates that simply increasing the size of a dataset with one single type of supervisory signals, such as paired texts in those multi-modal datasets, is still not good enough because knowledge required by a vast majority of tasks may be missing. In GV-D, as listed in Tab. 5, besides multi-modal data, we include three additional subsets corresponding to three common vision tasks, each of which has the largest scale compared to publicly available datasets for the same type of task. Our label system with 119K categories, which is fully utilized in the classification set GV-D_c-36M, is more than four times larger than that of the ImageNet-21K [14] dataset, covering a much wider range of hierarchically organized visual concepts. It significantly boosts the performance on more fine-grained tasks. We also include billion-level image-text pairs in GV-D-10B, which further expand supervisory signals extensively and reduce negative effects associated with noisy data. In summary, instead of being biased towards one supervision, GV-D provides rich data across multiple task types to achieve pretraining for general vision intelligence. Clear gains are observed in comprehensive evaluation on diverse tasks, which are discussed in Sec. 2.1.

Second, *multi-stage training brings consistent gains.* As shown in Tab. 6, INTERN learns fundamental skills with the *amateur* pretraining stage. Then, after *expert* pretraining with more specific and specialized knowledge, INTERN masters image classification with an average improvement of +2.8% for ResNet-50 and +3.8% for MN-B15 on 20 relevant datasets compared to the *amateur* stage. Finally, based on the multiple expert models from the *expert* stage, the *generalist* stage mines a universal set of skills, obtaining further gains of +0.6% and +0.2%. Meanwhile, the *generalist* model possesses sufficient generality towards more, possibly unseen, tasks and circumstances, which is demonstrated in Sec. 2.2. The consistent gains at successive stages show that our *amateur-expert-generalist* training pipeline effectively learns new knowledge and improves upon the previous stage. These results also imply that our continuous learning paradigm extends the ability of any single-stage pretraining, validating the design choice of combining them appropriately.

Datasets	Concepts	Images	Labels	Open Source
YFCC-100M [37]	-	99M	99M Texts	Yes
WIT [30]	500K Queries	400M	400M Texts	No
ALIGN [24]	-	1.8B	1.8B Texts	No
GV-D-10B	1.65M Queries	10B	10B Texts	Partially
ImageNet-21K [14]	22K Categories	14M	14M Image-Level Labels	Yes
IG-1B [28]	17K Queries	1B	1B Hashtags	No
JFT-3B [42]	30K Categories	3B	3B Noisy Image-Level Labels	No
GV-D_c-36M	119K Categories	36M	36M Image-Level Labels	Yes
COCO [27]	80 Categories	118K	1M Bounding Boxes	Yes
Object365 [34]	365 Categories	609K	10M Bounding Boxes	Yes
OpenImages [26]	600 Categories	2M	15M Bounding Boxes	Yes
GV-D_a-3M	809 Categories	3M	25M Bounding Boxes	Yes
COCO-Stuff [4]	182 Categories	118K	Segmentation Masks	Yes
GV-D_s-143K	334 Categories	143K	Segmentation Masks	Yes

Table 5: **Summary of GV-D and other large-scale datasets for visual pretraining.** A large-scale database is a fundamental component of general vision pretraining. YFCC-100M, ImageNet-21K, COCO, Object365, OpenImages are instances of commonly used public datasets, while IG-1B, WIT, JFT-3B are proprietary ones that cannot be accessed by the community. We construct a novel data system **GV-D** with four subsets: 1) **GV-D-10B** consisting of 10 billion image-text pairs collected with 1.65 million queries; 2) **GV-D_c-36M** contains 36 million images with classification labels from our label system of 119K categories. Although GV-D_c-36M has fewer images than JFT-3B, it has the most manual and clean labels. 3) **GV-D_a-3M** composed of 3 million images with 35 million bounding boxes of 809 categories; 4) **GV-D_s-143K** with 143 thousand images and corresponding semantic segmentation masks of 334 categories.

Model	Data Setting	Up-A	Up-E	Up-G
ResNet-50	10%	70.9	73.7	74.3
MN-B15		80.4	84.2	84.4

Table 6: **Average downstream classification accuracies of different upstream stages.** Classification precision monotonically increases when moving to later stages. Results are consistent at different model scales.

3 Conclusion

INTERN proves to be an effective paradigm towards general vision intelligence. All implementation details and part of our data system will be released. We hope our proposed INTERN will serve as a firm foundation for further general vision studies in the community. To the industry and society, strongly generalizable models from INTERN are expected to dramatically reduce data demands, thus facilitating applications of AI technologies. That being said, there are still many opportunities for improvement in the future. It is possible to incorporate new modalities, including but not limited to videos and audio. The final generalist model could potentially generalize to a far broader range of visual tasks beyond currently considered well-solved ones. We are also interested in boosting the data efficiency of upstream pretraining, making the whole learning process more economical and affordable. Moreover, future researches may leverage learned strong visual representation to approach more complex cognitive understanding tasks.

References

- [1] Hangbo Bao, Li Dong, and Furu Wei. Beit: Bert pre-training of image transformers. *arXiv preprint arXiv:2106.08254*, 2021.
- [2] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, pages 1877–1901, 2020.
- [4] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Coco-stuff: Thing and stuff classes in context. In *CVPR*, pages 1209–1218, 2018.
- [5] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *ECCV*, pages 132–149, 2018.
- [6] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.
- [7] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 1–21, 2021.
- [8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, pages 1597–1607, 2020.
- [9] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [10] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *CVPR*, pages 15750–15758, 2021.
- [11] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. *arXiv preprint arXiv:2104.02057*, 2021.
- [12] Michael Crawshaw. Multi-task learning with deep neural networks: A survey. *arXiv preprint arXiv:2009.09796*, 2020.
- [13] Zihang Dai, Hanxiao Liu, Quoc V Le, and Mingxing Tan. Coatnet: Marrying convolution and attention for all data sizes. *arXiv preprint arXiv:2106.04803*, 2021.
- [14] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.
- [15] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [16] 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*, 2020.
- [17] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *IJCV*, 88(2):303–338, 2010.
- [18] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, pages 3354–3361, 2012.
- [19] Ben Goertzel. Artificial general intelligence: concept, state of the art, and future prospects. *JAGI*, 5(1):1, 2014.
- [20] Priya Goyal, Mathilde Caron, Benjamin Lefaudeaux, Min Xu, Pengchao Wang, Vivek Pai, Mannat Singh, Vitaliy Liptchinsky, Ishan Misra, Armand Joulin, et al. Self-supervised pretraining of visual features in the wild. *arXiv preprint arXiv:2103.01988*, 2021.
- [21] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, Bilal Piot, koray kavukcuoglu, Remi Munos, and Michal Valko. Bootstrap your own latent - a new approach to self-supervised learning. In *NeurIPS*, pages 21271–21284, 2020.
- [22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [23] Yuqi Huo, Manli Zhang, Guangzhen Liu, Haoyu Lu, Yizhao Gao, Guoxing Yang, Jingyuan Wen, Heng Zhang, Baogui Xu, Weihao Zheng, et al. Wenlan: Bridging vision and language by large-scale multi-modal pre-training. *arXiv preprint arXiv:2103.06561*, 2021.
- [24] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc V Le, Yunhsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. *arXiv preprint arXiv:2102.05918*, 2021.
- [25] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. In *ECCV*, pages 491–507, 2020.
- [26] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, et al. The open images dataset v4. *IJCV*, 128(7):1956–1981, 2020.
- [27] 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*, pages 740–755, 2014.
- [28] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, pages 181–196, 2018.
- [29] Thomas Mensink, Jasper Uijlings, Alina Kuznetsova, Michael Gygli, and Vittorio Ferrari. Factors of influence for transfer learning across diverse appearance domains and task types. *arXiv preprint arXiv:2103.13318*, 2021.
- [30] 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 supervision. *arXiv preprint arXiv:2103.00020*, 2021.

- [31] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
- [32] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *arXiv preprint arXiv:2106.05974*, 2021.
- [33] 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.
- [34] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *ICCV*, pages 8430–8439, 2019.
- [35] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, pages 746–760, 2012.
- [36] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, pages 843–852, 2017.
- [37] Bart Thomee, David A Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. Yfcc100m: The new data in multimedia research. *Communications of the ACM*, 59(2):64–73, 2016.
- [38] Enze Xie, Jian Ding, Wenhai Wang, Xiahang Zhan, Hang Xu, Peize Sun, Zhenguo Li, and Ping Luo. Detco: Unsupervised contrastive learning for object detection. In *ICCV*, pages 8392–8401, 2021.
- [39] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *CVPR*, pages 1492–1500, 2017.
- [40] Shuo Yang, Ping Luo, Chen-Change Loy, and Xiaoou Tang. Wider face: A face detection benchmark. In *CVPR*, pages 5525–5533, 2016.
- [41] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *NeurIPS*, 32, 2019.
- [42] Xiaohua Zhai, Alexander Kolesnikov, Neil Houlsby, and Lucas Beyer. Scaling vision transformers. *arXiv preprint arXiv:2106.04560*, 2021.
- [43] Shanshan Zhang, Rodrigo Benenson, and Bernt Schiele. Citypersons: A diverse dataset for pedestrian detection. In *CVPR*, pages 3213–3221, 2017.