# VIPriors 3: Visual Inductive Priors for Data-Efficient Deep Learning Challenges

Robert-Jan Bruintjes, Attila Lengyel, Marcos Baptista Rios, Osman Semih Kayhan, Davide Zambrano, Nergis Tomen and Jan van Gemert

**Abstract**—The third edition of the "VIPriors: Visual Inductive Priors for Data-Efficient Deep Learning" workshop featured four data-impaired challenges, focusing on addressing the limitations of data availability in training deep learning models for computer vision tasks. The challenges comprised of four distinct data-impaired tasks, where participants were required to train models from scratch using a reduced number of training samples. The primary objective was to encourage novel approaches that incorporate relevant inductive biases to enhance the data efficiency of deep learning models. To foster creativity and exploration, participants were strictly prohibited from utilizing pre-trained checkpoints and other transfer learning techniques. Significant advancements were made compared to the provided baselines, where winning solutions surpassed the baselines by a considerable margin in all four tasks. These achievements were primarily attributed to the effective utilization of extensive data augmentation policies, model ensembling techniques, and the implementation of data-efficient training methods, including self-supervised representation learning. This report highlights the key aspects of the challenges and their outcomes.

**Index Terms**—Visual inductive priors, challenge, image classification, object detection, instance segmentation, action recognition.

✦

## 1 INTRODUCTION

DATA is fueling deep learning, yet obtaining high quality annotations is often costly. In recent years, extensive research has been dedicated to exploring ways to utilize large quantities of data to train comprehensive foundation models for vision and language [5], and to combine multiple modalities for weak supervision [43]. While these approaches have demonstrated impressive results, self-supervision is not yet the holy grail. Training on massive datasets still requires a significant amount of energy, contributing to carbon emissions. Furthermore, only a handful of deep learning behemoths have access to billions of data points and expensive deep learning hardware. In addition, large quantities of data may simply not be available in certain domains. The Visual Inductive Priors for Data-Efficient Deep Learning workshop (VIPriors) therefore aims to encourage research on learning efficiently from few data samples by combining the power of deep learning with hard-won knowledge priors from various fields. We focus on data efficiency through visual inductive priors.

The Visual Inductive Priors for Data-Efficient Deep Learning workshop has now been organized for the third year in a row, with the latest 2022 edition taking place at ECCV in Tel Aviv, Israel. In order to stimulate research in data-efficient computer vision, the workshop includes challenges in which participants train computer vision models on small subsets of (publicly available) datasets. We challenge participants to submit solutions that are able to learn an effective representation of the dataset without access to

the large quantities of data that is used to train state-of-the-art deep computer vision models.

In this report, we present the outcomes of the third edition of these challenges. We discuss specific details and top-ranking solutions of each challenge. It was observed that the top competitors in all challenges heavily relied on model ensembling and data augmentation to improve the data efficiency of their solutions. Additionally, many of the participants' solutions utilized a limited number of backbones and baseline methods, which seem to possess properties conducive to learning from small data. To recognize submissions that introduce innovative methods, a jury prize for each challenge was awarded to the most noteworthy submission.

## 2 CHALLENGES

The workshop accommodates four common computer vision challenges in which the number of training samples are reduced to a small fraction of the full set:

**Image classification**: We use a subset of Imagenet [15]. The subset contains 50 images from 1,000 classes for training, validation and testing.

**Object detection**: DelftBikes [30] dataset is used for the object detection challenge. The dataset includes 8,000 bike images for training and 2,000 images for testing (Fig. 1). Each image contains 22 different bike parts that are annotated as bounding box, class and object state labels.

**Instance segmentation**: The main objective of the challenge is to segment basketball players and the ball on images recorded of a basketball court. The dataset is provided by SynergySports[1] and contains a train, validation and test set of basketball games recorded at different courts with instance labels.

- R.J. Bruintjes, A. Lengyel, N. Tomen and J. van Gemert are with Delft University of Technology.
  E-mail: r.bruintjes@tudelft.nl
- M. Baptista Rios is with Gradiant.
- O. S. Kayhan is with Bosch Security Systems B.V.
- D. Zambrano is with Synergy Sports.

1. https://synergysports.com

TABLE 1
Overview of challenge submissions. J indicates jury prize. Bold-faced methods are contributions by the competitors.

| Ranking | Teams | Encoder architectures | Data augmentation | Methods | Main metric |
|---|---|---|---|---|---|
| **Classification** | | | | | |
| 1 | **Ma et al.** | SE+PyramidNet, ResNeSt200e, ReXNet,EfficientNet-B8, ConvNeXt-XL* | CutMix, AutoAugment [13], Stubborn Image Augmentation(SIA) | Label smoothing, AdvProp, Random image cropping and patching (RICP), extra training on stubborn images, hard fusion | **78.7** |
| 2 & J | Lu et al. | HorNet, ConvNeXt | Automix [42] | Cross-decoupled knowledge distillation [74], label smoothing | 77.9 |
| 3 | Zuo et al. | CoAtNet, TResNet, Resnet50, Resnext50, EdgeNeXt | CutMix, Random erasing, MixUp, AutoAugment | Knowledge distillation between encoders | 77.7 |
| 4 | Wang et al. | ResNeSt, TResNet, SE-ResNet, ReXNet, ECA-NFNet, ResNet-RS [1], Inception-ResNet, RegNet, EfficientNet, MixNet | AutoAugment, MixUp, CutMix, padding | label smoothing, train on larger images, data resampling | 76.8 |
| 5 | She et al. | ResNeSt, Res2Net, Xception, DPN [11], EfficientNet, SENet | AutoAugment, MixUp | label smoothing, train on larger images, hard negative resampling | 75.4 |
| 6 | Chen et al. | ResNeSt, EfficientNet, ReXNet, RegNetY | AutoAugment, MixUp, CutMix, ColorJitter | label smoothing, train on larger images, Exponential Moving Average on network parameters | 70.8 |
| **Object detection** | | | | | |
| 1 | **Lu et al.** | YOLOv4 [3], YOLOv7 [54], YOLOR [55], CBNetv2 [34] | Mosaic [3], mix-up [72], copy-paste [32] | Weighted Boxes Fusion [45], TTA, Model Soups [60], Image Uncertainty Weighted | **33.0** |
| 2 | Xu et al. | Cascade RCNN [7], Swin T. [38], ConvNext [41], ResNext [63] | AutoAugment [13], random flip, multi-scale augmentations [37] | MoCoV3 [10], MoBY [64], Soft-NMS [4], FPN [36], SSFPN [22], non-maximum weighted (NMW) [75] | 32.9 |
| 3 | J. Zhao et al. | Cascade RCNN [7], Swin T. [38], Convnext [41] | Albu, MixUp [72], AutoAugment [13] | Stochastic Weight Averaging [24], Hard classes retraining, FPN [36], Soft-NMS [4], pseudo labeling | 32.4 |
| 4 & J | P. Zhao et al. | Cascade RCNN [7], Swin T. [38], Pyramid ViT [58] | Mosaic [3], MixUp [72] | SimMIM [65], GIOU loss [52], Soft-NMS [4] | 30.9 |
| **Instance segmentation** | | | | | |
| 1 | **Yan et al.** [66] | CBSwin-T [35] | **TS-DA**, **TS-IP** [66] Random scaling, cropping AutoAugment [13] | Hybrid Task Cascade [8], CBFPN [35] | **53.1** |
| 2 (shared) | Leng et al. | Swin Transformer-Large [38] | ImgAug [26], Copy-Paste [18], Horizontal Flip and Multi-scale Training | CBNetV2 [35] | 50.6 |
| | Lu et al. | CBSwin-T [35] ResNet [21] ConvNeXt [41] Swinv2 [38] CBNetv2 [35] | MixUp [72], Mosaic Task-Specific Copy-Paste [69] Color and geometric transformations | Hybrid Task Cascade [8] CBFPN [35] Group Normalization [61] | 50.6 |
| 3 | Zhang et al. | CBSwin-T [35] | Location-aware MixUp, RandAugment [14], GridMask [9], Random scaling, Copy-paste [18], Multi-scale augmentation, TTA | Hybrid Task Cascade [8] Seesaw Loss [56] SWA [25] | 49.8 |
| 4 | Cheng et al. | CBSwin-T [35] | RandAugment Copy-Paste [18] GridMask | Hybrid Task Cascade [8] Mask Transfiner [31] | 18.5 |
| 5 & J | Cheng et al. [12] | ResNet [21] | Random flip and scale jitter | Sparse Instance Activation for Real-Time Instance Segmentation [12] | 18.5 |
| **Action recognition** | | | | | |
| 1 | **Song et al.** | R(2+1)D [49], SlowFast [16], CSN [50], X3D [17], TANet [40], Timesformer [2] | Random flipping, TenCrop | Soft voting | 0.71 |
| 2 | He et al. | SlowFast [16], Timesformer [2], TIN [46], TPN [67], X3D [17], Video Swin Transformers [39], R(2+1)D [49], DirecFormer [51]. | AutoAugment [13], CutMix [68] random flip, grayscale, jitter, temporal aug., TenCrop, test-time aug. | Label smoothing | 0.69 |
| 3 & J | Tan et al. | TSN [57], TANet [40], TPN [67], SlowFast [16], CSN [50], Video MAE [48] | MixUp [72], CutMix [68] | MoCo [19], TVL-1 [70] | 0.59 |

**Action recognition**: For this challenge we have provided Kinetics400ViPriors, which is an adaptation of the well-known Kinetics400 dataset [27]. The training set consists of approximately 40k clips, while the validation and test sets contain about 10k and 20k clips, respectively.

We provide a toolkit[2] which consists of guidelines, baseline models and datasets for each challenge. The competitions are hosted on the Codalab platform. Each participating team submits their predictions computed over a test set of samples for which labels are withheld from competitors.

The challenges include certain rules to follow:

- Models ought to train from scratch with only the given dataset.
- The usage of other data rather than the provided training data, pretraining the models and transfer learning methods are prohibited.
- The participating teams need to write a technical report about their methodology and experiments.

**Shared rankings.** Due to confusion around the exact deadline of the competitions, we have merged rankings of two different moments. This has resulted in shared places in some of the rankings of the individual challenges.

## 2.1 Classification

Image classification serves as an important benchmark for the progress of deep computer vision research. In particular, the ImageNet dataset [15] has been the go-to benchmark for image classification research. ImageNet gained popularity because of its significantly larger scale than those of existing benchmarks. Ever since, even larger datasets have been used to improve computer vision, such as the Google-owned JFT-300M [47]. However, we anticipate that relying on the increasing scale of datasets is problematic, as increased data collection is expensive and can clash with privacy interests of the subjects. In addition, for domains like medical imaging, the amount of labeled data is limited and the collection and annotation of such data relies on domain expertise. Therefore, we posit that the design of data efficient methods for deep computer vision is crucial.

As in the two earlier editions of this workshop, in our image classification challenge we provide a subset [29] of the Imagenet dataset [15] consisting of 50 images per class for each of the train, validation and test splits. The classification challenge had 14 participating teams, of which six teams submitted a report. The final ranking and the results can be seen in Table 2.

### 2.1.1 First place

The team from Xidian University led by Tianzhi Ma uses an ensemble of SE+PyramidNet, ResNeSt200e, ReXNet, EfficientNet-B8 and ConvNeXt-XL* models. They apply diverse data augmentation strategies to increase the diversity in the data, and include several other optimization tricks like RICP and hard fusion. Ultimately, they were able to improve their top-1 accuracy from last year's challenge (68.6) to a winning top-1 accuracy of 79%.

2. https://github.com/VIPriors/vipriors-challenges-toolkit

### 2.1.2 Second place & jury prize

The team from Xidian University led by Xiaoqiang Lu uses only two models in their ensemble, and instead gain performance by using cross-decoupled knowledge distillation [74]. Other than this, only Automix [42] and label smoothing are required to secure second place. For this minimal yet effective solution we award this team the jury prize.

### 2.1.3 Third place

The team from Xidian University lead by Yi Zou uses five different encoder architectures. They exhaustively apply knowledge distillation from all encoders to all other encoders to train twenty models, which are all ensembled for the final model. All models are trained with severe data augmentation: CutMix, random erasing, MixUp, AutoAugment.

### 2.1.4 Conclusion

As in previous editions [6], [33], the crucial components of a competitive submission to the image classification competition are ensembling of many different classification architectures, as well as combining multiple different augmentation policies. Aside from label smoothing and training with larger image sizes, knowledge distillation gained in popularity among the methods used to train the networks.

## 2.2 Object Detection

Similar to the object detection challenge last year [33], we also use DelftBikes [30] dataset this year (Fig. 1). Each image in the DelftBikes contains 22 labeled bike parts as class and bounding box labels of each part. In addition, the dataset includes extra object state labels as intact, missing, broken or occluded. The dataset has 10k bike images in total and 2k of the images are used for only testing purposes. The dataset contains different object sizes, and contextual and location biases that can cause false positive detections [28], [30]. Note that, some of the object boxes are noisy which introduces more challenges to detect object parts.



Fig. 1. Some images from the DelftBikes dataset. Each image has a single bike with 22 labeled parts.

We provide a baseline detector as a Faster RCNN with a Resnet-50 FPN [44] backbone from scratch for 16 epochs. This baseline network is trained with the original image size without any data augmentation. It performs 25.8% AP score on the test set. Note that, the evaluation is done on available parts which are intact, damaged and occluded parts.

TABLE 2
Final rankings of the Image Classification challenge.

| Ranking | Teams | Top-1 Accuracy |
|---|---|---|
| 1 | **Tianzhi Ma, Zihan Gao, Wenxin He, Licheng Jiao** <br> *School of Artificial Intelligence, Xidian University.* | **78.7** |
| 2 & J | Xiaoqiang Lu, Chao Li, Chenghui Li, Xiao Tan, Zhongjian Huang, Yuting Yang <br> *School of Artificial Intelligence, Xidian University.* | 77.9 |
| 3 | Yi Zuo, Zitao Wang, Xiaowen Zhang, Licheng Jiao <br> *School of Artificial Intelligence, Xidian University.* | 77.7 |
| 4 | Jiahao Wang, Hao Wang, Hua Yang, Fang liu, Lichang Jiao <br> *School of Artificial Intelligence, Xidian University.* | 76.8 |
| 5 | Wenxuan She, Mengjia Wang, Zixiao Zhang, Fang Liu, Licheng Jiao <br> *School of Artificial Intelligence, Xidian University.* | 75.4 |
| 6 | Baoliang Chen, Yuxuan Zhao, Fang Liu, Licheng Jiao <br> *School of Artificial Intelligence, Xidian University.* | 70.8 |

TABLE 3
Final rankings of the Image Object Detection challenge.

| Ranking | Teams | AP @ 0.5:0.95 |
|---|---|---|
| 1 | **Xiaoqiang Lu, Yuting Yang, Zhongjian Huang, Xiao Tan, Chenghui Li.** <br> *School of Artificial Intelligence, Xidian University* | **33** |
| 2 | Bocheng Xu, Rui Zhang, and Yanyi Feng. <br> *Department of AI R&D, Terminus Technologies.* | 32.9 |
| 3 | Jiawei Zhao, Zhaolin Cui, Xuede Li, Xingyue Chen, Junfeng Luo, and Xiaolin Wei. <br> *Vision Intelligence Department (VID).* | 32.1 |
| 4 & J | Ping Zhao, Xinyan Zhang, Weijian Sun, and Xin Zhang. <br> *Huawei Technologies Co., Ltd. and Tongji University* | 30.9 |

The detection challenge had 41 participant teams. The team from Xidian University obtained first place by 33% AP scores. Terminus Technologies and Vision Intelligence Department from Meituan followed them by 32.9% AP and 32.1% AP respectively. The team from Huawei Technologies and Tongji University won the jury prize for their 'coarse-to-fine' idea.

### 2.2.1 First place

Lu et al. employ an ensemble of various YOLO detectors [3], [54], [55] and CBNetv2 [34] (Fig. 2). They design two-stage training: (i) pre-training by using weak data augmentation and (ii) fine-tuning by using strong data augmentation such as mosaic [3], mix-up [72], and copy-paste [32] and a weighted training strategy based on image uncertainty. The authors further improved the results by weighted boxes fusion (WBF) [45] and TTA strategies and obtain 33 % AP on the test set.

### 2.2.2 Second place

Xu et al. train Cascade RCNN [7] using Swin Transformer [38], ConvNext [41] and ResNext [63] as backbone architectures. These backbones are pretrained by using self-supervised methods such as MoCoV3 [10] and MoBY [64]. In addition, they use AutoAugment [13], random flip and multi-scale augmentation methods [37] to improve the detection performance. Finally, the non-maximum weighted (NMW) [75] method, Soft-NMS [4] and model ensemble methods are used on the test set. The method obtained 32.9% AP on the test set.
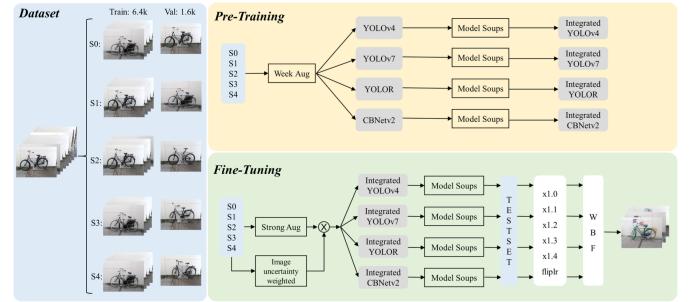


Fig. 2. Bag of Freebies for training detector [73]. They train different models during the pretraining and fine-tuning phases with different types of data augmentation methods. They also use image uncertainty to improve object detection performance.

### 2.2.3 Third place

Zhao et al. initially train Cascade RCNN [7] detector using ConvNext backbone [41]. Then, they create a synthetic dataset (Fig. 3) from the training set and obtain pseudo labels on this dataset with the initial trained model. Afterwards, they train the same model only on the pseudo labels with smaller-resolution images. In the end, they retrain the pseudo-label pretrained network with the original train set and select some hard classes to improve the detector performance on them. During training phases, they also use various data augmentation methods as colour jittering and RGB shifting, mix-up [72] and AutoAugment [13]. The method obtains 32.1% AP detection performance.
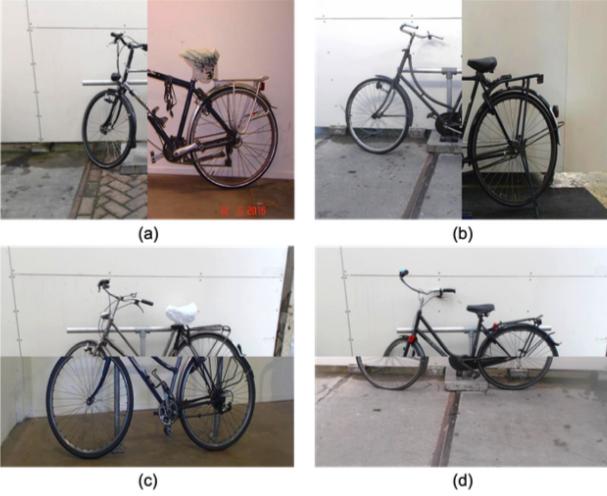
Fig. 3. Synthetic images generated for backbone pretraining.



Fig. 5. Data augmentation policy of the first-place instance segmentation submission by Yan et al. [66].

### 2.2.4  Jury prize

Method of Zhao et al. has two phases: pretraining and adaptation phases (Fig. 4). In the pretraining phase, they utilize mosaic [3] and mix-up [72] data augmentations on object and image level features and train SimMIM [65]. In the adaptation phase, a pretrained encoder of SimMIM is used to initialize the backbone of Cascade RCNN [7]. In 'coarse detection', the model detects bike objects. In the 'fine detection' phase, the fine detection module runs on the cropped bike object from the previous phase and tries to detect relevant bike parts. The final model obtains 30.94% AP. The team earned the jury prize because of their 'coarse-to-fine' idea, well-written article and discussion of strategies that did not work.
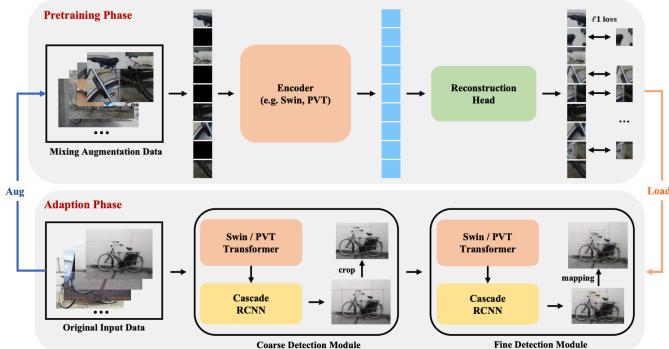


Fig. 4. Method pipeline. First, the backbone is pretrained by using SimMIM [65] to obtain strong features. In the adaptation phase, coarse to fine detection strategy improves detection.

## 2.3  Instance Segmentation

Instance segmentation is the task of detecting and segmenting specific objects and instances in an image. With applications ranging from autonomous driving, surveillance, remote sensing to sport analysis, it is a fundamental computer vision problem. Similarly to last year, our challenge is based

on the basketball dataset provided by SynergySports [53], consisting of images recorded during various basketball games played on different courts. The goal is to detect and predict segmentation masks of all players and ball objects in the images. With a mere 184, 62, and 64 samples for the train, validation and test splits, respectively, the dataset is considered very small. The test labels are withheld from the challenge participants and final performance on the test set is evaluated on an online server. The main metric used is the Average Precision (AP) @ 0.50:0.95. Our baseline method is based on the Detectron2 [62] implementation of Mask-RCNN [20].

Twelve teams submitted solutions to the evaluation server, of which six teams submitted a report to qualify their submission to the challenge. The final rankings are shown in Table 4.

### 2.3.1  First place

The method of Yan et al. [66] introduces a task-specific data augmentation (TS-DA) strategy to generate additional training data, and a task-specific inference processing (TS-IP) which is applied at test time. TS-DA employs Copy-Paste [18] augmentations with constraints on the absolute locations of the synthetic players and ball objects to ensure all objects are placed inside the court, and their relative locations to mimic player-ball interactions. Subsequently, geometric and photometric augmentations are applied to the image to further increase the variety in their appearance. During inference, random scaling and cropping is applied to the images, and additional filtering employed to the predictions to ensure only one basketball of reasonable dimensions is present on the court. The complete data augmentation policy is illustrated in Figure 5.

The segmentation model is based on the Hybrid Task Cascade (HTC) detector [8] and the CBSwin-T backbone with CBFPN [35]. Mask Scoring R-CNN [23] is used to further improve segmentation quality. After training the model, it is further finetuned using the SWA [71] strategy.

TABLE 4
Final rankings of the Instance Segmentation challenge. J indicates jury prize.

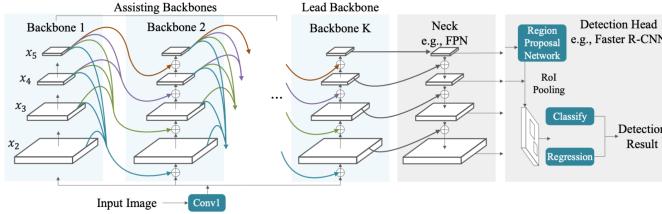| Ranking | Teams | % AP @ 0.50:0.95 |
|---|---|---|
| 1 | **Bo Yan, Xingran Zhao, Yadong Li, Hongbin Wang.** *Ant Group, China.* | **53.1** |
| 2 (shared) | Fuxing Leng, Jinghua Yan, Peibin Chen, Chenglong Yi. *ByteDance, Huazhong University of Science and Technology.* | 50.6 |
| | Xiaoqiang Lu, Yuting Yang, Zhongjian Huang. *School of Artificial Intelligence, Xidian University, Xi'an, China.* | 50.6 |
| 3 | Junpei Zhang, Kexin Zhang, Rui Peng, Yanbiao Ma, Licheng Jiao Fang Liu. *Team Yanbiao_Ma.* | 49.8 |
| 4 | Yi Cheng, ShuHan Wang, Yifei Chen, Zhongjian Huang. *School of Artificial Intelligence, Xidian University, Xi'an, China.* | 47.6 |
| 5 & J | Tianheng Cheng, Xinggang Wang, Shaoyu Chen, Qian Zhang, Chang Huang, Zhaoxiang Zhang, Wenqiang Zhang, Wenyu Liu. *(1) School of EIC, Huazhong University of Science & Technology; (2) Horizon Robotics; (3) Institute of Automation, Chinese Academy of Sciences (CASIA)* | 34.0 |



Fig. 6. Instance segmentation model architecture and training pipeline of the method of Leng et al.



Fig. 7. Overview of instance segmentation method by the third place competitor, Zhang et al.

### 2.3.2 Shared second place - A

Leng et al. demonstrate that a straightforward combination of well-proven methods can yield near-SoTA performance. The approach uses a Swin Transformer-Large [38] as the backbone, and the pipeline is based on CBNetV2 [35], as shown in Figure 6. In terms of data augmentations the method relies on a combination of AutoAugment [13], ImgAug [26] and Copy-Paste [18].

### 2.3.3 Shared second place - B

Lu et al. make use of the popular HTC detector [8] with CBSwin-T [35] backbone with CBFPN [35] using group normalization, Mosaic, test-time augmentations and the Task-Specific Copy-Paste Data Augmentation Method [18] from a previous edition of the VIPriors Instance Segmentation challenge. Moreover, different backbones, namely ResNet [21], ConvNeXt [41], Swinv2 [38] and CBNetv2 [35] are trained and combined using Model Soups [59]. Multiple predictions are combined together using mask voting.

### 2.3.4 Third place

The method of Zhang et al. employs Location-aware MixUp, RandAugment, GridMask, Random scaling, Copy-paste, Multi-scale augmentation, and test-time augmentation in terms of data augmentation techniques. The model used is the popular HTC detector [8] and soft non-maxima suppression is applied on the predicted target boxes. The overall framework is depicted in figure 7.
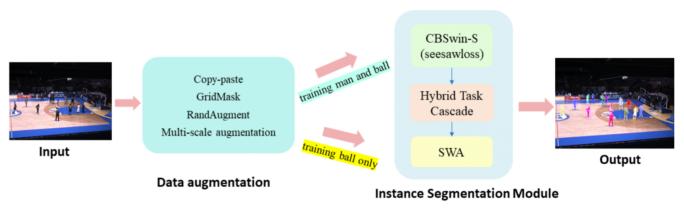
### 2.3.5 Jury prize

This year's jury prize is awarded to Sparse Instance Activation for Real-Time Instance Segmentation [12] by Cheng et al. The paper presents a method for instance segmentation using a novel representation of instance activation maps. These maps highlight informative regions for each object, which are then used to obtain instance-level features for recognition and segmentation. The method avoids the need for non-maximum suppression in post-processing by predicting objects in a one-to-one style using bipartite matching.

## 2.4 Action Recognition

Many of the popular Action Recognition models are deep networks that require a large amount of data for training, which can be challenging when there is limited data availability or insufficient compute resources. In line with the workshop's goals, we present the Kinetics400ViPriors dataset, which is a modified version of the well-known Kinetics400 dataset. We have created a smaller variant with 40k, 10k, and 20k clips for the train, validation, and test sets, respectively, while preserving the original number of action classes. Our aim is to motivate researchers in Action Recognition to develop efficient models that can leverage visual prior knowledge from the data.

For evaluation, we use the average classification accuracy across all classes on the test set. The accuracy for a single class is calculated as $\text{Acc} = \frac{P}{N}$, where $P$ represents the number of correct predictions for the evaluated class and
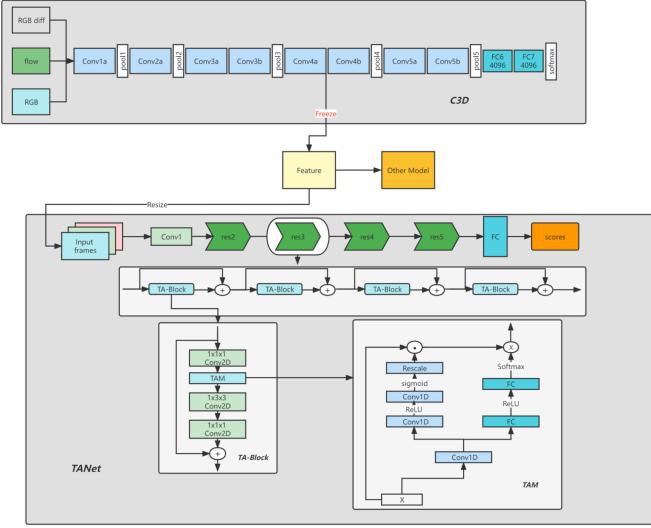
Fig. 8. Method proposed by first place Song et al. from Xidian University.



Fig. 9. Schematic diagram of the method (inference mode) proposed by second place He et al. from Xidian University.

$N$ is the total number of samples in that class. The average accuracy is determined by taking the mean of accuracies for all classes.

9 teams submitted solutions to the evaluation server, of which 3 teams submitted a report to qualify their submission to the challenge. The final rankings are shown in Table 5.

### 2.4.1 First place

The authors train a selection of models, including R(2+1)D [49], SlowFast [16], CSN [50], X3D [17], TANet [40] and Timesformer [2], and apply a model fusion approach by assigning different weights to the models and using the soft voting method to combine their results. In terms of data augmentation, frames were extracted from videos and a subset was selected by choosing every second frame. The videos were resized and noise was added through random flipping. During testing, TenCrop was used as a test-time enhancement. The evaluation involved ten-fold cross-validation, where the training dataset was combined with the validation dataset. An overview of the method is provided in Fig. 8.

### 2.4.2 Second place

The method proposes a multi-network dynamic fusion model combining a variety of backbones, including Slow-Fast [16], Timesformer [2], TIN [46], TPN [67], Video Swin Transformers [39], R(2+1)D [49], X3D [17], DirecFormer [51]. Model predictions are combined as a weighted average by the prediction score of each model. Test-time augmentation with majority voting is used, as well as AutoAugment [13], CutMix [68], and a variety of other spatial, photometric and temporal augmentations during training. An overview of the method is provided in Fig. 9.

### 2.4.3 Third place & Jury prize

The method combines self-supervised pre-training of various backbone models, optical flow estimation and model ensembling to train a d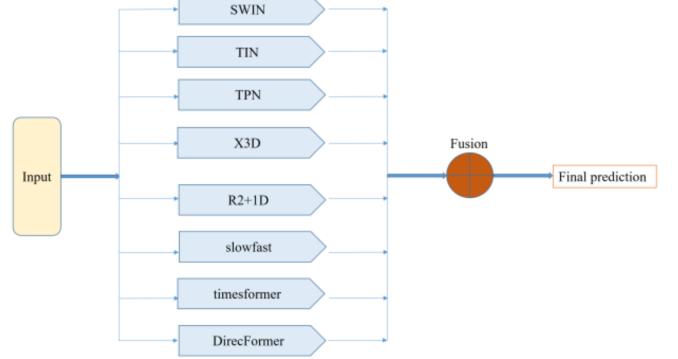ata efficient video classification model. First, the 2D model encoders are pre-trained using the MoCo [19] self-supervised representation learning framework on image data using the individual frames of the provided dataset. Next, optical flow features are extracted using the TVL-1 [70] method. To correct for camera movement, consecutive image frames are aligned by calculcating the transformation matrix based on extracted SIFT features. Finally, a range of models including TSN [57], TANet [40], TPN [67], SlowFast [16], CSN [50] and Video MAE [48] are trained on the training data and pre-extracted optical flow features. MixUp [72] and CutMix [68] data augmentation is employed. Model ensembling is performed by concatenating the features of all models and training a single linear classifier layer after normalization. An ablation study is performed to show that self-supervised pre-training improves model performance.

## 3 CONCLUSION

We have summarized all solutions in Table 1 in terms of the encoder architecture, data augmentation techniques and main methods used.

Organizing the same challenges for the third year in a row gives a unique perspective on trends: which methods and/or architectures prevail over time, and which are replaced? The use of combining large numbers of models in ensembles and heavy data augmentation have been unchanging throughout the VIPiors challenge series. The models used in the ensembles are a mix of CNNs and Vision Transformers for the tasks of object detection and instance segmentation, whereas for image classification and action recognition Vision Transformers are not seeing use in our challenges. As for data augmentation, AutoAugment, MixUp and CutMix are unchanging constants in the training regimes of our competitors, regardless of the task.

Though we did not explicitly perform the required analysis, we cannot escape the impression that the simplicity of model ensembling is hard to beat with task-specific or domain-specific knowledge, especially when considering the effort required in design and implementation. Winning methods tend to use ensembling, and the bigger the ensemble, the better, as is shown in the ablation studies of some of the competitors' reports. If one is to follow this approach, we speculate that choosing models with a variety of inductive

TABLE 5
Final rankings of the Action Recognition challenge. J indicates Jury prize.

| Ranking | Teams | Acc |
|---|---|---|
| 1 | **Xinran Song, Chengyuan Yang, Chang Liu, Yang Liu, Fang Liu, Licheng Jiao** *School of Artificial Intelligence, Xidian University, Xi'an, China.* | **0.71** |
| 2 | Wenxin He, Zihan Gao, Tianzhi Ma , Licheng Jiao *School of Artificial Intelligence, Xidian University, Xi'an, China.* | 0.69 |
| 3 & J | Bo Tan, Yang Xiao, Wenzheng Zeng, Xingyu Tong, Zhiguo Cao, Joey Tianyi Zhou *Huazhong University of Science and Technology (China) and CFAR (Singapore)* | 0.59 |

biases (e.g. CNNs and Vision Transformers) could make the ensemble more effective. However, such heavy use of ensembles may just be possible in our challenges because of the limited size of the datasets, which makes training many models feasible.

# REFERENCES

[1] Bello, I., Fedus, W., Du, X., Cubuk, E.D., Srinivas, A., Lin, T.Y., Shlens, J., Zoph, B.: Revisiting resnets: Improved training and scaling strategies. Advances in Neural Information Processing Systems **34**, 22614–22627 (2021) 2

[2] Bertasius, G., Wang, H., Torresani, L.: Is space-time attention all you need for video understanding? arXiv preprint arXiv:2102.05095 (2021) 2, 7

[3] Bochkovskiy, A., Wang, C.Y., Liao, H.Y.M.: Yolov4: Optimal speed and accuracy of object detection (2020) 2, 4, 5

[4] Bodla, N., Singh, B., Chellappa, R., Davis, L.S.: Soft-nms–improving object detection with one line of code. In: Proceedings of the IEEE international conference on computer vision. pp. 5561–5569 (2017) 2, 4

[5] Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M.S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N.S., Chen, A.S., Creel, K.A., Davis, J., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L.E., Goel, K., Goodman, N.D., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D.E., Hong, J., Hsu, K., Huang, J., Icard, T.F., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P.W., Krass, M.S., Krishna, R., Kuditipudi, R., Kumar, A., Ladhak, F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X.L., Li, X., Ma, T., Malik, A., Manning, C.D., Mirchandani, S., Mitchell, E., Munyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B., Nie, A., Niebles, J.C., Nilforoshan, H., Nyarko, J.F., Ogut, G., Orr, L.J., Papadimitriou, I., Park, J.S., Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y.H., Ruiz, C., Ryan, J., R'e, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K.P., Tamkin, A., Taori, R., Thomas, A.W., Tramèr, F., Wang, R.E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie, S.M., Yasunaga, M., You, J., Zaharia, M.A., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L., Zhou, K., Liang, P.: On the opportunities and risks of foundation models. ArXiv **abs/2108.07258** (2021) 1

[6] Bruintjes, R.J., Lengyel, A., Rios, M.B., Kayhan, O.S., van Gemert, J.: Vipriors 1: Visual inductive priors for data-efficient deep learning challenges. arXiv preprint arXiv:2103.03768 (2021) 3

[7] Cai, Z., Vasconcelos, N.: Cascade r-cnn: Delving into high quality object detection. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 6154–6162 (2018) 2, 4, 5

[8] Chen, K., Pang, J., Wang, J., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Shi, J., Ouyang, W., Loy, C.C., Lin, D.: Hybrid task cascade for instance segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019) 2, 5, 6

[9] Chen, P., Liu, S., Zhao, H., Jia, J.: Gridmask data augmentation (2020) 2

[10] Chen, X., Xie, S., He, K.: An empirical study of training self-supervised vision transformers. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9640–9649 (2021) 2, 4

[11] Chen, Y., Li, J., Xiao, H., Jin, X., Yan, S., Feng, J.: Dual path networks. Advances in neural information processing systems **30** (2017) 2

[12] Cheng, T., Wang, X., Chen, S., Zhang, W., Zhang, Q., Huang, C., Zhang, Z., Liu, W.: Sparse instance activation for real-time instance segmentation. In: Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR) (2022) 2, 6

[13] Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: Autoaugment: Learning augmentation policies from data. arXiv preprint arXiv:1805.09501 (2018) 2, 4, 6, 7

[14] Cubuk, E.D., Zoph, B., Shlens, J., Le, Q.V.: Randaugment: Practical automated data augmentation with a reduced search space. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops. pp. 702–703 (2020) 2

[15] Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009) 1, 3

[16] Feichtenhofer, C., Fan, H., Malik, J., He, K.: Slowfast networks for video recognition. In: 2019 IEEE/CVF International Conference on Computer Vision (ICCV). pp. 6201–6210 (2019). https://doi.org/10.1109/ICCV.2019.00630 2, 7

[17] Feichtenhofer, C.: X3d: Expanding architectures for efficient video recognition. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 200–210 (2020). https://doi.org/10.1109/CVPR42600.2020.00028 2, 7

[18] Ghiasi, G., Cui, Y., Srinivas, A., Qian, R., Lin, T.Y., Cubuk, E.D., Le, Q.V., Zoph, B.: Simple copy-paste is a strong data augmentation method for instance segmentation. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 2917–2927 (2021). https://doi.org/10.1109/CVPR46437.2021.00294 2, 5, 6

[19] He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9726–9735 (2020). https://doi.org/10.1109/CVPR42600.2020.00975 2, 7

[20] He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: 2017 IEEE International Conference on Computer Vision (ICCV). pp. 2980–2988 (2017). https://doi.org/10.1109/ICCV.2017.322 5

[21] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition (2015) 2, 6

[22] Hong, M., Li, S., Yang, Y., Zhu, F., Zhao, Q., Lu, L.: Sspnet: Scale selection pyramid network for tiny person detection from uav images. IEEE Geoscience and Remote Sensing Letters **19**, 1–5 (2021) 2

[23] Huang, Z., Huang, L., Gong, Y., Huang, C., Wang, X.: Mask scoring r-cnn. In: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 6402–6411 (2019). https://doi.org/10.1109/CVPR.2019.00657 5

[24] Izmailov, P., Podoprikhin, D., Garipov, T., Vetrov, D., Wilson, A.G.: Averaging weights leads to wider optima and better generalization. arXiv preprint arXiv:1803.05407 (2018) 2

[25] Izmailov, P., Podoprikhin, D., Garipov, T., Vetrov, D.P., Wilson, A.G.: Averaging weights leads to wider optima and better generalization. CoRR **abs/1803.05407** (2018), http://arxiv.org/abs/1803.05407 2

[26] Jung, A.B.: imgaug. https://github.com/aleju/imgaug (2018), [Online; accessed 30-Oct-2018] 2, 6

[27] Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., Suleyman, M., Zisserman, A.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017) 3

[28] Kayhan, O.S., van Gemert, J.C.: Evaluating context for deep object detectors. arXiv preprint arXiv:2205.02887 (2022) 3

[29] Kayhan, O.S., Gemert, J.C.v.: On translation invariance in cnns: Convolutional layers can exploit absolute spatial location. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14274–14285 (2020) 3

[30] Kayhan, O.S., Vredebregt, B., van Gemert, J.C.: Hallucination in object detection–a study in visual part verification. arXiv preprint arXiv:2106.02523 (2021) 1, 3

[31] Ke, L., Danelljan, M., Li, X., Tai, Y.W., Tang, C.K., Yu, F.: Mask transfiner for high-quality instance segmentation. In: CVPR (2022) 2

[32] Kisantal, M., Wojna, Z., Murawski, J., Naruniec, J., Cho, K.: Augmentation for small object detection. arXiv preprint arXiv:1902.07296 (2019) 2, 4

[33] Lengyel, A., Bruintjes, R.J., Rios, M.B., Kayhan, O.S., Zambrano, D., Tomen, N., van Gemert, J.: Vipriors 2: visual inductive priors for data-efficient deep learning challenges. arXiv preprint arXiv:2201.08625 (2022) 3

[34] Liang, T., Chu, X., Liu, Y., Wang, Y., Tang, Z., Chu, W., Chen, J., Ling, H.: Cbnetv2: A composite backbone network architecture for object detection. arXiv preprint arXiv:2107.00420 (2021) 2, 4

[35] Liang, T., Chu, X., Liu, Y., Wang, Y., Tang, Z., Chu, W., Chen, J., Ling, H.: Cbnetv2: A composite backbone network architecture for object detection. arXiv preprint arXiv:2107.00420 (2021) 2, 5, 6

[36] Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B., Belongie, S.: Feature pyramid networks for object detection (2017) 2

[37] Liu, W.K., Hao, S., Belytschko, T., Li, S., Chang, C.T.: Multiscale methods. International Journal for Numerical Methods in Engineering 47(7), 1343–1361 (2000) 2, 4

[38] Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. International Conference on Computer Vision (ICCV) (2021) 2, 4, 6

[39] Liu, Z., Ning, J., Cao, Y., Wei, Y., Zhang, Z., Lin, S., Hu, H.: Video swin transformer. arXiv preprint arXiv:2106.13230 (2021) 2, 7

[40] Liu, Z., Zhao, X., Huang, T., Hu, R., Zhou, Y., Bai, X.: Tanet: Robust 3d object detection from point clouds with triple attention. AAAI (2020), https://arxiv.org/pdf/1912.05163.pdf 2, 7

[41] Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A convnet for the 2020s. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022) 2, 4, 6

[42] Liu, Z., Li, S., Wu, D., Chen, Z., Wu, L., Guo, J., Li, S.Z.: Unveiling the power of mixup for stronger classifiers. arXiv preprint arXiv:2103.13027 (2021) 2, 3

[43] Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I.: Learning transferable visual models from natural language supervision. In: Meila, M., Zhang, T. (eds.) Proceedings of the 38th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 139, pp. 8748–8763. PMLR (18–24 Jul 2021), https://proceedings.mlr.press/v139/radford21a.html 1

[44] Ren, S., He, K., Girshick, R., Sun, J.: Faster r-cnn: Towards real-time object detection with region proposal networks (2016) 3

[45] Roman Solovyev, W.W., Gabruseva, T.: Weighted boxes fusion: Ensembling boxes from different object detection models. Image and Vision Computing 107, 104117 (2021) 2, 4

[46] Shao, H., Qian, S., Liu, Y.: Temporal interlacing network. AAAI (2020) 2, 7

[47] Sun, C., Shrivastava, A., Singh, S., Gupta, A.: Revisiting unreasonable effectiveness of data in deep learning era. In: Proceedings of the IEEE international conference on computer vision. pp. 843–852 (2017) 3

[48] Tong, Z., Song, Y., Wang, J., Wang, L.: VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In: Advances in Neural Information Processing Systems (2022) 2, 7

[49] Tran, D., Wang, H., Torresani, L., Ray, J., LeCun, Y., Paluri, M.: A closer look at spatiotemporal convolutions for action recognition. In: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6450–6459 (2018). https://doi.org/10.1109/CVPR.2018.00675 2, 7

[50] Tran, D., Wang, H., Torresani, L., Feiszli, M.: Video classification with channel-separated convolutional networks. CoRR abs/1904.02811 (2019), http://arxiv.org/abs/1904.02811 2, 7

[51] Truong, T.D., Bui, Q.H., Duong, C.N., Seo, H.S., Phung, S.L., Li, X., Luu, K.: Direcformer: A directed attention in transformer approach to robust action recognition. In: Computer Vision and Pattern Recognition (2022) 2, 7

[52] Union, G.I.O.: A metric and a loss for bounding box regression. In: Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, S. Savarese//IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), Long Beach, CA, USA. pp. 658–666 (2019) 2

[53] Van Zandycke, G., Somers, V., Istasse, M., Don, C.D., Zambrano, D.: Deepsportradar-v1: Computer vision dataset for sports understanding with high quality annotations. In: Proceedings of the 5th International ACM Workshop on Multimedia Content Analysis in Sports. p. 1–8. MMSports '22, Association for Computing Machinery, New York, NY, USA (2022). https://doi.org/10.1145/3552437.3555699, https://doi.org/10.1145/3552437.3555699 5

[54] Wang, C.Y., Bochkovskiy, A., Liao, H.Y.M.: Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. arXiv preprint arXiv:2207.02696 (2022) 2, 4

[55] Wang, C.Y., Yeh, I.H., Liao, H.Y.M.: You only learn one representation: Unified network for multiple tasks. arXiv preprint arXiv:2105.04206 (2021) 2, 4

[56] Wang, J., Zhang, W., Zang, Y., Cao, Y., Pang, J., Gong, T., Chen, K., Liu, Z., Loy, C.C., Lin, D.: Seesaw loss for long-tailed instance segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 9695–9704 (June 2021) 2

[57] Wang, L., Xiong, Y., Wang, Z., Qiao, Y., Lin, D., Tang, X., Van Gool, L.: Temporal segment networks: Towards good practices for deep action recognition. vol. 9912 (10 2016) 2, 7

[58] Wang, W., Xie, E., Li, X., Fan, D.P., Song, K., Liang, D., Lu, T., Luo, P., Shao, L.: Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 568–578 (2021) 2

[59] Wortsman, M., Ilharco, G., Gadre, S.Y., Roelofs, R., Gontijo-Lopes, R., Morcos, A.S., Namkoong, H., Farhadi, A., Carmon, Y., Kornblith, S., Schmidt, L.: Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In: Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., Sabato, S. (eds.) Proceedings of the 39th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 162, pp. 23965–23998. PMLR (17–23 Jul 2022), https://proceedings.mlr.press/v162/wortsman22a.html 6

[60] Wortsman, M., Ilharco, G., Gadre, S.Y., Roelofs, R., Gontijo-Lopes, R., Morcos, A.S., Namkoong, H., Farhadi, A., Carmon, Y., Kornblith, S., et al.: Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In: International Conference on Machine Learning. pp. 23965–23998. PMLR (2022) 2

[61] Wu, Y., He, K.: Group normalization. In: ECCV (2018) 2

[62] Wu, Y., Kirillov, A., Massa, F., Lo, W.Y., Girshick, R.: Detectron2. https://github.com/facebookresearch/detectron2 (2019) 5

[63] Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K.: Aggregated residual transformations for deep neural networks. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1492–1500 (2017) 2, 4

[64] Xie, Z., Lin, Y., Yao, Z., Zhang, Z., Dai, Q., Cao, Y., Hu, H.: Self-supervised learning with swin transformers. arXiv preprint arXiv:2105.04553 (2021) 2, 4

[65] Xie, Z., Zhang, Z., Cao, Y., Lin, Y., Bao, J., Yao, Z., Dai, Q., Hu, H.: Simmim: A simple framework for masked image modeling. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 9653–9663 (2022) 2, 5

[66] Yan, B., Zhao, X., Li, Y., Wang, H.: Task-specific data augmentation and inference processing for vipriors instance segmentation challenge (2022), https://arxiv.org/abs/2211.11282 2, 5

[67] Yang, C., Xu, Y., Shi, J., Dai, B., Zhou, B.: Temporal pyramid network for action recognition. In: Proceedings of the IEEE Confer-

ence on Computer Vision and Pattern Recognition (CVPR) (2020) 2, 7

[68] Yun, S., Han, D., Oh, S.J., Chun, S., Choe, J., Yoo, Y.: Cutmix: Regularization strategy to train strong classifiers with localizable features. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 6023–6032 (2019) 2, 7

[69] Yunusov, J., Rakhmatov, S., Namozov, A., Gaybulayev, A., Kim, T.H.: Instance segmentation challenge track technical report, vipriors workshop at iccv 2021: Task-specific copy-paste data augmentation method for instance segmentation (2021) 2

[70] Zach, C., Pock, T., Bischof, H.: A duality based approach for realtime tv-l1 optical flow. In: Hamprecht, F.A., Schnörr, C., Jähne, B. (eds.) Pattern Recognition. pp. 214–223. Springer Berlin Heidelberg, Berlin, Heidelberg (2007) 2, 7

[71] Zhang, H., Wang, Y., Dayoub, F., Sünderhauf, N.: Swa object detection. arXiv preprint arXiv:2012.12645 (2020) 5

[72] Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. In: International Conference on Learning Representations (2018), https://openreview.net/forum?id=r1Ddp1-Rb 2, 4, 5, 7

[73] Zhang, Z., He, T., Zhang, H., Zhang, Z., Xie, J., Li, M.: Bag of freebies for training object detection neural networks (2019), http://arxiv.org/abs/1902.04103, cite arxiv:1902.04103 4

[74] Zhao, B., Cui, Q., Song, R., Qiu, Y., Liang, J.: Decoupled knowledge distillation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11953–11962 (2022) 2, 3

[75] Zhou, H., Li, Z., Ning, C., Tang, J.: Cad: Scale invariant framework for real-time object detection. In: Proceedings of the IEEE international conference on computer vision workshops. pp. 760–768 (2017) 2, 4