

SYNERGY AND DIVERSITY IN CLIP: ENHANCING PERFORMANCE THROUGH ADAPTIVE BACKBONE ENSEMBLING

Anonymous authors

Paper under double-blind review

ABSTRACT

Contrastive Language-Image Pretraining (CLIP) stands out as a prominent method for image representation learning. Various architectures, from vision transformers (ViTs) to convolutional networks (ResNets) have been trained with CLIP to serve as general solutions to diverse vision tasks. This paper explores the differences across various CLIP-trained vision backbones. Despite using the same data and training objective, we find that these architectures have notably different representations, different classification performances across datasets, and different robustness properties to certain types of image perturbations. Our findings indicate a remarkable possible synergy across backbones by leveraging their respective strengths. In principle, classification accuracy could be improved by over 40 percent with an informed selection of the optimal backbone per test example. Using this insight, we develop a straightforward yet powerful approach to adaptively ensemble multiple backbones. The approach uses as few as one labeled example per class to tune the adaptive combination of backbones. On a large collection of datasets, the method achieves a remarkable increase in accuracy of up to 39.1% over the best single backbone, well beyond traditional ensembles.

1 INTRODUCTION

Large pre-trained models are transforming machine learning, computer vision (Radford et al., 2021; He et al., 2022; Kirillov et al., 2023; Liu et al., 2023), and natural language processing (Devlin et al., 2018; Brown et al., 2020; Touvron et al., 2023). These models are typically trained with self-supervised objectives, eschewing the need for manual annotations and enabling the use of very large datasets (Jia et al., 2021; Schuhmann et al., 2022).

Contrastive Language–Image Pre-training (CLIP) (Radford et al., 2021) is one such approach that enables various downstream applications involving vision and language. CLIP learns to align text and image representations when training a pair of vision and language encoders, which can be implemented with various architectures. These encoders are then used in downstream applications to obtain representations of images and text, whose similarity can be simply evaluated with a dot product. This enables e.g. zero-shot classification (Zhai et al., 2022; Li et al., 2022a; Jia et al., 2021) and cross-modal retrieval (Li et al., 2020a;b; Yu et al., 2022).

Despite extensive prior work on CLIP, there is still a gap in comparing the representations from different vision backbones trained with this paradigm. Existing work has compared backbones for their generalization capabilities (Goldblum et al., 2023; Li et al., 2022a; Zhang et al., 2021; Gao et al., 2021), finding that larger architectures generally perform better. This paper contributes to this area with an empirical study that compares CLIP-trained backbones. We present new observations on how performance and robustness (e.g. invariance to image transformations) vary substantially across architectures. Within the same backbone family (e.g. different ViTs), different models present different patterns of performance across datasets (Figure 4). The relation with model scale is also more complicated than suggested in prior work (Figure 5). Our empirical findings suggest that CLIP’s effectiveness for image classification could be improved by combining the strengths of different backbones. We thus present an ensemble method that combines backbones adaptively to each test

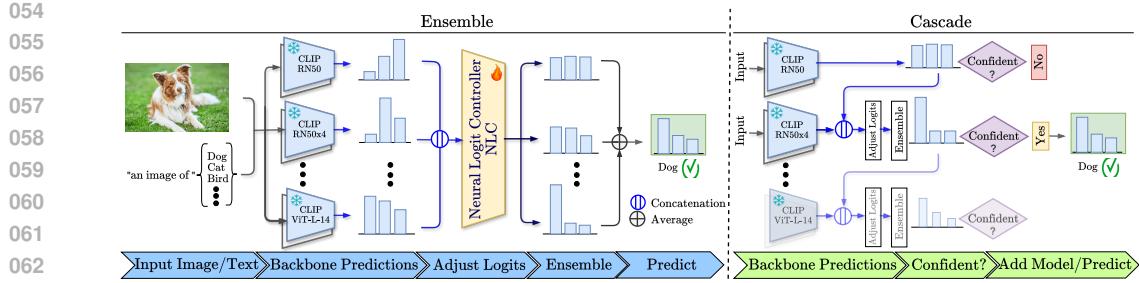


Figure 1: We propose a method to improve CLIP’s effectiveness for image classification by combining the strengths of different backbones. **(Left)** For a given test image, the logits from different backbones are combined with a temperature scaling that weights their contribution to the final prediction. The scaling is implemented in the Neural Logit Controller (NLC, a small MLP) that is learned from as little as one labeled example per class. **(Right)** To reduce the computational load, our method can be combined with the Cascade framework (Wang et al., 2022a).

example, using a temperature scaling that weights each backbone’s contribution to the final prediction condition on the input image (see Figure 1).

A naive combination of backbones with traditional ensembling techniques (Lakshminarayanan et al., 2017; Dietterich, 2000) (i.e. averaging the outputs) does not consistently enhance generalization performance. These approaches focus on prediction agreements rather than leveraging diversity. In contrast, our method uses a temperature scaling based on the idea of calibrating each backbone’s confidence. This intuition is further supported by the “modality gap” and variation in performance observed in prior work when adjusting CLIP’s logits (Shi et al., 2023). We compare our approach to more advanced ensembles, such as SuperLearner Cheng Ju and van der Laan (2018), which also combines the logits of multiple models by learning temperature scaling. However, unlike our method, SuperLearner does not adaptively adjust these temperatures based on the input features, implying that SuperLearner cannot exploit the diversity of predictions Tab. C.1.

Our proposed method, named the Neural Logit Controller (NLC), uses a few labeled examples (as little as one per class) to tune the combination of backbones. Hence, we benchmark it against the state-of-the-art methods in a few-shot setting (Zhang et al., 2021; Radford et al., 2021; Zhou et al., 2022; Gao et al., 2021). NLC shows remarkable performance and consistently outperforms other approaches. Furthermore, we show that NLC complements existing methods like Tip-Adapter (Zhang et al., 2021). Combining it with NLC in its original experimental framework yields improvements in accuracy of over 15%.

Finally, in-depth experiments reveal a clear correlation between NLC’s improvement over the best backbone and the diversity of correct predictions

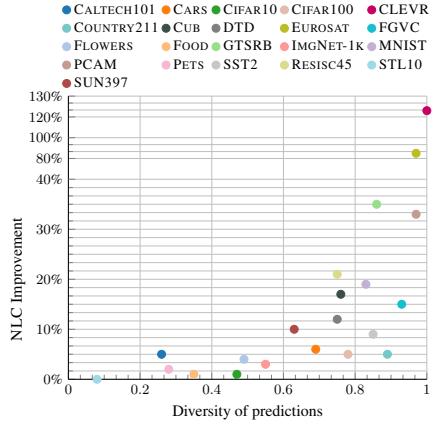


Figure 2: **(Y axis)** Relative improvement of NLC over best backbone vs. **(X axis)** predictions diversity¹. NLC always improves over the best backbone. Moreover, the clear correlation shows that higher diversity in predictions tends to result in greater improvements with NLC.

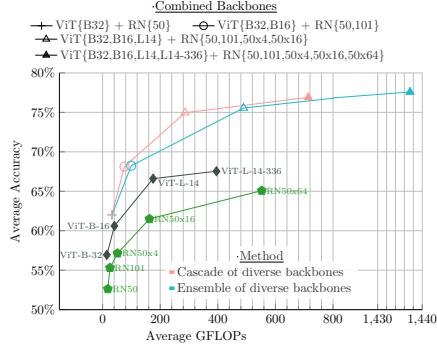


Figure 3: Average accuracy across 21 datasets for zero-shot ResNets and ViTs backbones, NLC Ensemble, and Cascade using 2 to 9 backbones, plotted against the Average GFLOPs. Demonstrates that ensembles can surpass the best zero-shot backbone with fewer GFLOPs.

¹The diversity of predictions is measured as: $1 - \frac{\# \text{samples correctly predicted by all backbones}}{\# \text{samples correctly predicted by any backbone}}$.

in each dataset (see Figure 2). It achieves over 120% of relative improvement when the diversity of predictions is close to one, and never degrades the performance of the best backbone. In terms of efficiency, NLC outperforms the best backbone by combining only the top-four most efficient backbones, while using approximately 300 fewer GFLOPs (see Figure 3). We also show how to use the Cascade framework (Wang et al., 2022a; Varshney and Baral, 2022) with NLC to enhance performance while maintaining computational requirements within the bounds of the original backbones.

Our contributions are summarized as follows.

- We perform extensive experiments across 21 datasets that reveal diverse predictions across CLIP backbones and distinct robustness properties to image transformations.
- We evaluate the complementarity of backbones by measuring the potential improvement in classification accuracy (up to 43.5%) of an optimal oracle selection of backbone per test example.
- To leverage this complementarity, we propose an adaptive ensembling method (NLC) that uses temperature scaling and requires a single labeled example per class. It improves accuracy over the best backbone by 9.1% on average, surpassing previous ensembling frameworks.
- We demonstrate significant advantages in computational efficiency. Combining the top four most efficient backbones enables NLC to outperform the best backbone with ~ 300 fewer GFLOPs. Integrating NLC with the Cascade approach (Wang et al., 2022a) maintains computational efficiency.
- We demonstrate that NLC consistently outperforms state-of-the-art few-shot methods. Moreover, integrating NLC with Tip-Adapter (Zhang et al., 2021) gives the latter a performance boost of over 15%.

2 INVESTIGATING DIFFERENCES ACROSS VISION BACKBONES

CLIP for zero-shot image classification. CLIP (Radford et al., 2021) is a general technique to train a pair of vision and text encoders ϕ_l and ϕ_v on paired image/text data. The method uses a self-supervised objective to align the encoded representations of matching images and text descriptions. Any image and text can then be mapped to a unified semantic space. This enables zero-shot image classification by computing a compatibility score between an image and a class description (“prompt”) as $\text{score}(x, l) = \phi_v(x)^\top \phi_l(l)$. We follow previous work (Menon and Vondrick, 2023) for the generating of prompts e.g. an image of $\{\text{label}\}$, which is a $\{\text{concept}\}$. Subsequently, we calculate the compatibility score between the image feature $\phi_v(x)$ and the prompt representations $\phi_l(y)$, where y represents a label within the class set. The prompt with the highest similarity is then chosen as the label for the given image.

2.1 COMPLEMENTARITY OF DIFFERENT BACKBONES

The vision encoder or “backbone” trained with CLIP can be implemented using various architectures. We want to identify whether different backbones have distinctive behaviours. We analyze the output of 9 of them in a zero-shot classification setting on 21 standard benchmark datasets (see Section 4 for details). We propose two approaches to assess the differences in their zero-shot predictions and their possible complementary behaviour. First, we define an ORACLE prediction that combines the output of the 9 models on a per-image basis, using a correct prediction if any of the models is correct, any other otherwise. This simulates a scenario where an informed choice of the optimal backbone could be made for each test image to obtain the highest classification accuracy. Second, we use overlap

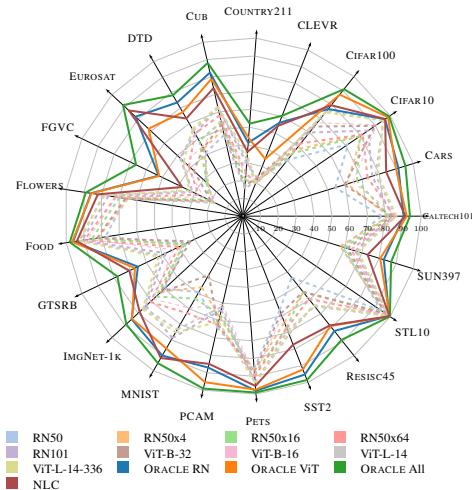


Figure 4: Zero-shot classification accuracy of various CLIP models on 21 datasets, and of the upper-bound “ORACLE” combination of ResNets (RN), ViTs, and all backbones.

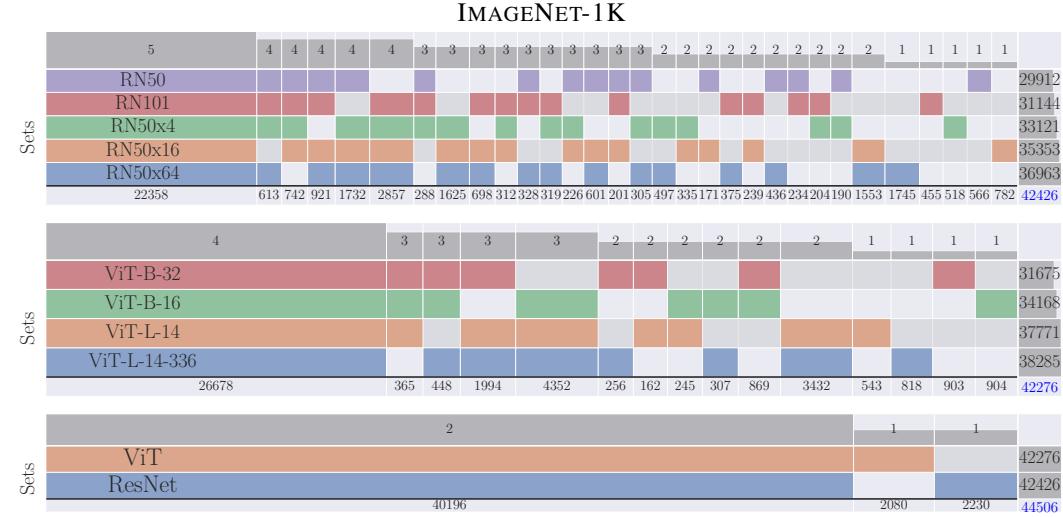


Figure 5: Linear Venn diagrams showing the overlap of test images from ImageNet-1k correctly classified by different backbones (rows). Each column represents a subset of images correctly classified by a specific group of backbones (group size in column header). Row/column sums indicate the number of correct predictions per backbone/subset. We observe that (1) different backbones agree on a large part of the data. (2) They also make additional correct predictions on different subsets. (3) Accuracy usually grows with architectures size, but even within a same family (ViTs, ResNets), different models show different patterns of (in)correct predictions.

diagrams to visualize the patterns of (in)correct predictions across backbones and their different strengths on different subsets of images.

Figure 4 and Table B.1 summarize the performance of the different backbones. Note that they are pre-trained using the exact same dataset and objectives, as provided in the OpenClip project (Ilharco et al., 2021). Our results show that backbones from different families (ViTs vs. ResNets) show significant differences while controlling for model size (number of parameters). For example, ViT-B-32 outperforms ResNet-50x4 in the CIFAR10 and CIFAR100 benchmarks. However, this trend is not consistent across all datasets, e.g. on IMGNET-1k. Additionally, the results indicate that backbones with more parameters do not always outperform smaller ones.

We then assess the ORACLE predictions using backbones from the same family (ORACLE RN for ResNets, ORACLE ViT for ViTs, and ORACLE All for a combination of all backbones). The ORACLE All obviously outperforms the best individual backbone in every case. For instance, on EUROSAT, CLEVR, CUB, DTD, CIFAR100, and CARS, the ORACLE All shows improvements of 43.5%, 36.0%, 25.1%, 21.8%, 16.6% and 16.0%, respectively. **The large magnitude of these improvements highlights the complementarity of backbones. In other words, the correctly-classified data is not simply growing as one considers larger and better models. Instead, different backbones perform well on different subsets of the data.** This insight is the key motivation for the ensemble method we propose in Section 3. We also found that similar findings have been made in domains other than ours Roth et al. (2024); Zhong et al. (2021); Ramé et al. (2022).

We visualize in Figure 5 the overlap between different subsets of IMGNET-1k correctly predicted by different backbones. We observe comparable accuracy between the largest ViT (ViT-L-14-336, correct predictions) and ResNet (RN50x64) with 38,285 vs. 36,963 correct predictions. However, there are significant differences within each family: for example, among ResNets, only 22,358 images are correctly classified by all ResNets models, out of the 42,426 ones correctly classified by *any* ResNet.¹

We further observe that every model has unique strengths by examining the performance of the oracle models. Specifically, ORACLE RN, ORACLE ViT, and ORACLE All increase the number of correct predictions to 42,426, 42,276, and 44,506 out of 50,000. Each backbone exhibits correct predictions unique to itself. In the ORACLE RN, RN50, RN101, RN50x4, RN50x16, and RN50x64 contribute 566, 455, 518, 782, and 1,745 exclusive predictions, respectively. In the ORACLE ViT, B-32, B-16, L-14, and L-14-336 present 903, 904, 543, and 818 unique predictions, respectively.

Finally, when all backbones are used, ResNet contributes 2,230 exclusive predictions, and ViT contributes 2,080. It is important to note that correct predictions of any given model are not simply a subset of those by the best model in its family.

2.2 OTHER SOURCES OF DIVERSITY

To better understand the effects of backbone complementarity, and contrast these effects against the role of variables that are known to drive performance in ensembles such random parameter initialization and learning via stochastic gradient descent, we propose to evaluate the diversity of correct predictions when combining multiple backbones, which as seen in our experiments above (Figure 2) strongly correlates with ensemble performance gains.

Concretely, we isolate the effects of three variables on prediction diversity: (1) backbone architecture, (2) pre-training dataset, and (3) number of training steps. For (1), we evaluate the complementarity of ViT-L-14, ViT-B-32, and ViT-B-16 backbones using the same dataset. For (2), we analyze each ViT’s performance using different pre-training datasets (LAION-400M vs. OpenCLIP). For (3), we consider ViTs pretrained on LAION-400M using the checkpoints at two different epochs (31 and 32).

Since, to the best of our knowledge, there are no publicly released CLIP models that differ only on their initialization, studies on this variable would, in principle, require us to train CLIP backbones from scratch, which is extremely compute-intensive.

Our diversity metric considers the aggregate of test examples correctly classified and it is the ratio between the instances correctly classified by all backbones versus the instances correctly classified by any backbone¹. It can be thought of as $(1 - \text{IoU})$, where IoU is the "intersection over union" of n sets of correctly predicted samples, where n is the number of methods considered. If every sample in this aggregate is correctly classified by all backbones (i.e. they all agree while being correct), $\text{diversity}=0$. Otherwise, if every sample can only be correctly classified by a single backbone, $\text{diversity}=1$. This allows us to evaluate the complementarity of models.

Figure 6 (Tab. I.2) demonstrates that combining multiple backbones achieves higher average diversity than using the same backbone across different datasets or nearby solutions in the loss landscape. This indicates that training stochasticity and pretraining dataset are not the primary factors. Instead, the most meaningful complementarity comes from diverse backbones, as they excel on different data subsets, as confirmed by the ORACLE results in Table I.1.

2.3 ROBUSTNESS TO IMAGE TRANSFORMATIONS

To explore qualitative differences in behaviour of different backbones, we selected a subset of each benchmark’s test set and applied specific image transformations (resize, flip, grayscale, color jitter, and Gaussian blur). We then evaluate the zero-shot classification performance of each backbone on these corrupted images. Figure 7 shows that ResNet 50 and ViT-B-16 are the most resilient to flip, RN50x64 to resize, and RN50x16 to grayscale. ViT-L-14 and L-14-336 are the most robust to color jitter and Gaussian blur. These findings suggest

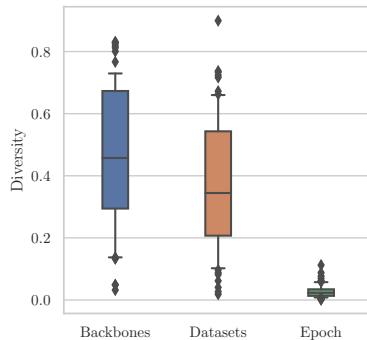


Figure 6: Diversity of predictions of (1) different **Backbones** with same pre-trained dataset, (2) same backbone with different pre-trained **Datasets**, and (3) same backbone and same pre-trained dataset in two different **Epochs**. Results show complementarity is higher when we combine different backbones.

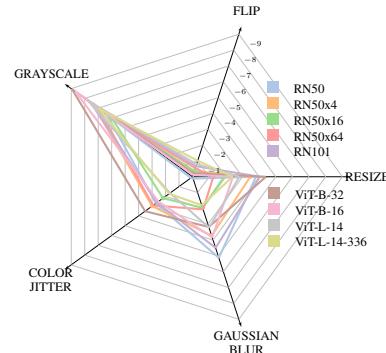


Figure 7: Impact on test classification accuracy of image transformations. Smaller values (closer to center) mean higher robustness, suggesting models are particularly robust to specific transformations.

270 that the backbones' ability to deal with out-of-distribution images varies widely. Certain models
 271 prove particularly robust to specific transformations.
 272

273 3 PROPOSED ENSEMBLING METHOD

276 The previous section showed that different backbones have different strengths and perform well on
 277 different types of data. We now use this insight to improve classification performance by combining
 278 multiple CLIP backbones.

280 **Combining models by temperature scaling.** Our approach is technically inspired by prior work
 281 on network calibration, which we apply to an ensemble of models. Specifically, our technique can
 282 be seen as a variation of **Platt scaling** (Platt et al., 1999). This classical method uses the logits of a
 283 model as features for a logistic regression, which is trained on the validation set to produce calibrated
 284 probabilities. More precisely, given logit scores z_i for an example i , Platt scaling learns two scalars
 285 a and b and produces calibrated probabilities as $\hat{q}_i = \sigma(az_i + b)$. The parameters a and b can be
 286 optimized using the negative log likelihood (NLL) loss over the validation set, with the model weights
 287 being frozen during this process.

288 A common use of Platt scaling is known as **temperature scaling** (Guo et al., 2017). In this approach,
 289 a single scalar parameter $t > 0$ is used for all classes of a given model. The new, calibrated confidence
 290 prediction is given by $\hat{q}_i = \max_k \text{softmax}(z_i/t)^{(k)}$, where t is called the temperature, z_i are the
 291 logits for example i returned by the uncalibrated model. Usually, t is optimized with respect to the
 292 NLL on the validation set aiming to reduce the overconfidence of the model on its predictions and
 293 to produce more reliable predictions, but because the parameter t does not change the maximum of
 294 the softmax function, the class prediction \hat{y}_i remains unchanged, meaning that the performance of a
 295 given model remains the same.

296 **Learning temperature coefficients.** In our method, we aim to jointly optimize a set of temperature
 297 parameters t_b with $b \in [1, \dots, B]$ for a set of B backbones. The aim is to combine their predictions
 298 by adjusting each backbone's confidence depending on the confidence of the others, and the input
 299 example. We learn the temperatures t_b that weigh the logit z_i^b for a backbone b and example x using
 300 the cross-entropy loss, [and then we combined the logits by a weighted sum](#).

301 Concretely, we train a one-layer MLP (our Neural Logit Controller) to predict the set of temperatures
 302 that best calibrate the backbone mixture. the MLP takes as input the concatenated representations
 303 obtained by passing the images through the encoder ϕ_v of each backbone $b \in \mathbb{B}$. As depicted in
 304 Figure 1, The MLP directly produces a vector of temperatures $t \in \mathbb{R}^B$ and is trained on a holdout set
 305 for the training set of each target dataset using the cross-entropy loss between final predictions and
 306 ground truth labels.

308 4 EXPERIMENTS

310 **Backbones.** We consider the original selection of backbones described by Radford et al. (2021).
 311 This includes the ResNets (He et al., 2016) RN50, RN50x4, RN50x16, RN50x64 and RN101, and
 312 the ViTs (Dosovitskiy et al., 2021) B-16, B-32, L-14 and L-14-336. All models are obtained through
 313 the open-source project OpenCLIP (Ilharco et al., 2021).

315 **Datasets.** We use a selection of 21 popular image classification datasets: CALTECH101 (Li et al.,
 316 2022b), CARS (Krause et al., 2013), CIFAR10 (Krizhevsky et al., 2009), CIFAR100 (Krizhevsky et al.,
 317 2009), CLEVR (Johnson et al., 2017), CUB (Wah et al., 2011), DTD (Cimpoi et al., 2014), EUROSAT
 318 (Helber et al., 2018), FGVC (Maji et al., 2013), FLOWERS (Nilsback and Zisserman, 2008), FOOD
 319 (Bossard et al., 2014), GTSRB (Houben et al., 2013), IMAGENET-1K (Deng et al., 2009) MNIST (Deng,
 320 2012), PCAM (Veeling et al., 2018), PETS (Parkhi et al., 2012), RenderedSST2 (Socher et al., 2013),
 321 RESISC45 (Cheng et al., 2017), STL10 (Coates et al., 2011) and SUN397 (Xiao et al., 2010). They
 322 include diverse images of nature, animals, places, medical scans, satellite images, and man-made
 323 objects. The evaluation metric is simply the accuracy of each dataset's test split. The aggregated
 performance is the average accuracy across all datasets weighted equally.

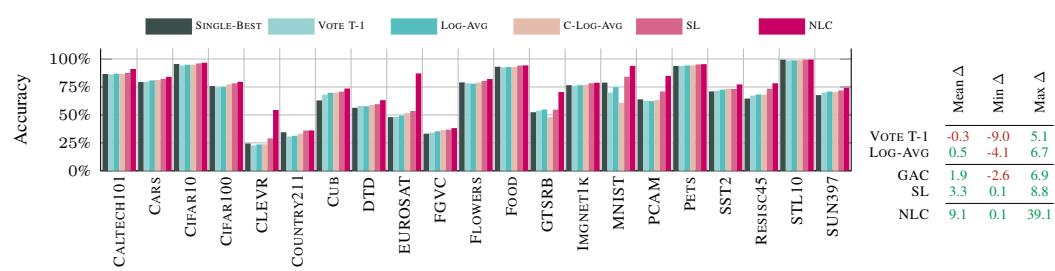


Figure 8: Comparison of the top-performing single backbone (SINGLE-BEST) with NLC and other ensembling strategies on top of zero-shot CLIP backbones. In the table, Mean, Max and Min Δ summarize the difference in performance across datasets with respect to the SINGLE-BEST backbone.

Baselines. Our experiments consider several ensembling baselines. Non-parametric approaches:

- **Logit averaging** (LOG-AVG) simply takes the average of the logits (scores) of the models.
- **Voting** (VOTE T-1) uses a majority vote. Each model casts a vote for its max-scored class. The final prediction is the class with the highest number of votes. Compared to LOG-AVG, this does not use the *soft* score values and the uncertainty reflected therein.

Parametric approaches:

- **Calibrated logit averaging** (C-LOG-AVG) first calibrates each model independently using temperature scaling. They are then combined by simply averaging their logits.
- **Super Learner** (SL) (Cheng Ju and van der Laan, 2018) a flexible ensemble learning framework that optimally combines predictions from multiple base learners. In contrast to our method, SuperLearner learns temperature scaling factors without considering the input. Thus, finding the best scaling for the total validation-set.
- **Mixture of experts** (MOE) is a popular method to improve performance by combining specialized submodels of experts (Chen et al., 2022; Lepikhin et al., 2020; Rau, 2019; Zoph et al., 2022), each one of which focuses on a specific input region. A gating network determines which expert is relevant for a given input. We use the Sparse MoE implementation (Zoph et al., 2022) by Rau (2019), where the MoE layer’s input is the concatenation of vision features x_b from ϕ_b . We use standard hyperparameters: 9 experts trained for classification with cross-entropy loss with Adam (Kingma and Ba, 2014) and a learning rate of $2e-5$ for 300 epochs.

See Appendix C a discussion of additional baselines. The appendix also contains additional results on the combination of backbones and how they perform under distribution shifts on the IMGNET-1K dataset (Section E). Finally, Section H presents dataset-specific results including overlap diagrams, an analysis of possible model combinations using NLC and our ORACLE, as well as the dataset-specific learned temperature values.

4.1 COMPARISON OF BACKBONE-ENSEMBLING METHODS

In this section, we evaluate the proposed NLC approach along with various non-parametric and parametric ensembling techniques. The non-parametric ones are “static” while the parametric ones use the **training split** of each target dataset to adjust the ensemble adaptively. The following results will show the importance of this adaptive setting to leverage the backbones’ respective strengths, and the superiority of our method over baseline parametric techniques.

Figure 8 and Table C.1 show the performance of each technique. Among **non-parametric baselines**, LOG-AVG fails to enhance overall performance beyond the best backbone. One notable exception is the COUNTRY dataset where LOG-AVG shows a substantial improvement over the best backbone. Overall, though, simple score averaging does not seem to exploit the complementarity of the backbones, with an improvement in accuracy of only +1.3%. When the backbones are calibrated with C-LOG-AVG, performance improves compared to the non-calibrated version, across all datasets except GTSRB, MNIST, RESISC45, and SUN397. **Vote T-1** is mostly ineffective and does not significantly improve over the best backbone. Among **parametric methods**, the proposed NLC shows the most substantial improvement, up to +39.1% on the EUROSAT dataset. On average, we obtain a commendable improvement of +9.1% over the SINGLE-BEST backbone across all datasets.



Figure 9: Comparison of the top-performing single backbone (SINGLE-BEST) with NLC and ensemble strategies on top of CLIP backbones linear classifier, where Mean, Max and Min Δ summarize the difference in performance across datasets with respect to the SINGLE-BEST.

4.2 INVESTIGATING BENEFITS BEYOND DATASET-SPECIFIC LINEAR CLASSIFIERS

We now design an additional experiment to investigate the benefits of the proposed NLC. A standard approach to adapt CLIP is to train a dataset-specific linear classifier on frozen visual features. In this section, we first train such linear classifiers, using each backbone’s frozen features and each dataset’s training split. Whereas NLC normally acts on raw CLIP scores, we now train it to act on the logits from these classifiers instead. Since the linear classifiers already are strong dataset-specific models, one might expect that only small additional improvements are possible. Yet, the following results will show that the NLC obtains further performance gains thanks to the *instance*-level adaptation. In contrast to alternative approaches, it also never reduces the performance compared to the best single linear classifier.

Setup. We follow the setup of Li et al. (2022a). We initialize the weights of the linear classifiers using language weights, which was shown to be more stable than a random initialization. The output of each backbone ϕ_b is L2-normalized before being passed to each linear classifier, which is trained using 90% of the target dataset’s training split. The remaining 10% are used to train NLC.

Results. Figure 9 and Table D.2 present the results combining linear classifiers with both non-parametric and parametric approaches. Similar to Section 4.1, the proposed NLC consistently improves performance across all datasets, with an average improvement of +1.6% over the SINGLE-BEST linear classifier. As expected, the improvement is reduced compared to the original NLC since some of the gains are now realised by the dataset-specific linear classifier. Yet Table D.1 clearly shows that additional gains can be made over the linear classifiers, i.e. that there remains an exploitable complementarity across backbones. By examining the ORACLE performance in this setting, we can see that the potential benefits of combining multiple backbones remain. Interestingly, the MOE approach does not reach the performance level of the best backbone on several datasets. It shows a degradation sometimes down to -20.9%. This suggests that MoEs may struggle to partition the input space effectively into distinct clusters, specific to certain experts, potentially limiting its functionality.

Table 1: Accuracy (%) in a few-shot setting

Few-shot	1	2	4	8	16
Linear-probe CLIP Radford et al. (2021)	22.2	31.9	41.2	49.5	56.1
CoOP Zhou et al. (2022)	47.6	50.9	56.2	59.9	63.0
CLIP-Adapter Gao et al. (2021)	61.2	61.5	61.8	62.7	63.6
Tip-Adapter Zhang et al. (2021)	60.7	61.0	61.0	61.5	62.0
Tip-Adapter-F Zhang et al. (2021)	61.3	61.7	62.5	64.0	65.5
NLC	78.2	78.1	78.2	78.3	78.4

4.3 COMPARISON WITH FEW-SHOT ADAPTER METHODS

We now examine the performance of the proposed NLC in a few-shot setting. We also show that it complements existing few-shot adapter methods and improves their accuracy across various datasets.

We compare our method against Tip-Adapter, Tip-Adapter-F (Zhang et al., 2021), few-shot Linear-probe CLIP (Radford et al., 2021), CoOP (Zhou et al., 2022), and CLIP-Adapter (Gao et al., 2021). Linear-probe CLIP fine-tunes a linear classifier on a few-shot training set on the frozen CLIP backbones. CoOP (Zhou et al., 2022) generates different prompt designs to make prompts learnable. CLIP-Adapter (Gao et al., 2021) enhances few-shot classification by introducing a feature adapter on CLIP’s visual and textual encoders. Tip-Adapter (Zhang et al., 2021) achieves comparable performance to CLIP-Adapter without requiring training, using a key-value cache model from the

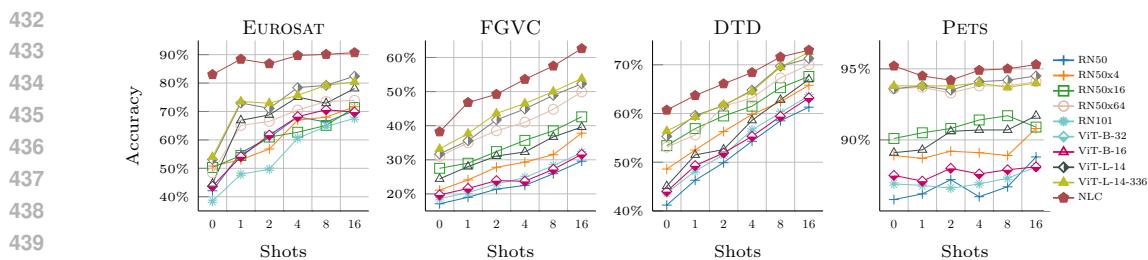


Figure 10: Peformance of TiP-Adapter applied to different zero-shot CLIP backbones, and the combination of these with NLC, showing how our method provides complementary benefits.

few-shot training set. Tip-Adapter-F can further enhance performance by fine-tuning the cache. As shown in Table 1, our NLC shows an outstanding performance over compared methods and consistently surpasses other few-shot methods by learning how to combine multiple backbones, which is a complementary approach to any of these adapters. In the appendix, Table F.2 shows that current few-shot adapter methods also have performance differences with various backbones, and still, our NLC surpasses the best backbone reported by these existing methods.

Finally, we integrate Tip-Adapter (Zhang et al., 2021) with our NLC and conduct experiments on EUROSAT, FGVC, DTD, and PETS. Initially, we apply Tip-Adapter and Tip-Adapter-F independently on each backbone, following the protocol from Zhang et al. (2021) with 1, 2, 4, 8, and 16 shots. Subsequently, we employ our ensembling mechanism to fuse the adapted backbones. To train NLC, we follow their setting and use the same validation split as used by Tip-Adapter to linearly combine their logits with CLIP.

See Figure 10 and Tables F.3–F.6 for a comparison of the performance of Tip-Adapter over different zero-shot CLIP backbones, and the combination of all these Tip-Adapter versions with NLC. Across the four evaluated datasets, NLC improves over each version of Tip-Adapter in the few-shot setting. Notably, NLC improve the performance of the best Tip-Adapter backbone (L-14-336) of up to 10% for EUROSAT using 16 shots. Moreover, when we compare NLC and SL in the Tip-Adapter setting, NLC shows better performance, Table F.3–F.6.

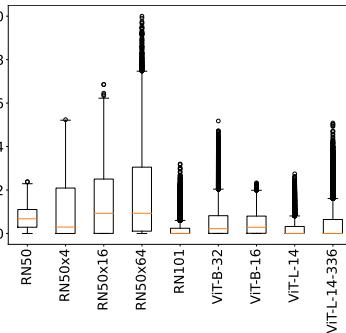


Figure 11: Distribution of max-normalized temperature values learned by our NLC method on the CLEVR dataset.

4.4 EXAMINING THE LEARNED TEMPERATURE VALUES

We visualise the learned temperature values to understand how the NLC method assigns weights to different backbones, revealing its strategy for leveraging the strengths of diverse architectures. Figure 11 shows the learned temperature values for NLC on the CLEVR dataset (see Section H for other datasets). By examining the learned temperatures, we can gain insight into the contribution of each backbone. We note that the most weighted backbone is not consistently the deepest one such as ResNet-101. The uniform distribution of learned temperatures across backbones indicates that NLC effectively utilizes the strengths of multiple backbones, resulting in a more balanced and accurate labeling process, complementing our analysis in Figure 2. These findings highlight that NLC benefits from the diversity in model predictions, reinforcing the advantage of combining various backbones.

5 RELATED WORK

Foundation models in computer vision. The trend in computer vision over the past decade has been to train larger and larger models on increasingly diverse tasks and increasingly larger datasets (Gadre et al., 2023; Schuhmann et al., 2022). Scale and generality have proved to be key enablers to the performance and robustness of these models (Fang et al., 2022; Gan et al., 2022; Mayilvahanan et al., 2023; Santurkar et al., 2022; Tu et al., 2024). Various methods have been proposed to pre-train such models (He et al., 2022; Radford et al., 2021; Yu et al., 2022). CLIP (Radford et al., 2021) is a

notable one, specifically designed to exploit noisy image-text pairs scraped from the internet. CLIP’s significance lies in its scalability and its ability to generate meaningful alignments with prompts, facilitating zero-shot classification gao2021clip,zhang2021tip,li2022elevater. CLIP’s popularity is due to its scalability and ability to handle diverse text prompts that enable a variety of downstream applications (Gao et al., 2021; Zhang et al., 2021; Li et al., 2022a). Various versions of CLIP have been trained on different datasets and vision backbones such as ResNets (He et al., 2016) and ViTs (Dosovitskiy et al., 2021).

Comparing vision backbones. Goldblum et al. (2023) recently conducted a comprehensive comparison of pretrained backbones across various tasks including image classification, object detection, segmentation, and image retrieval. In comparison, this paper presents complementary observations of the strengths of different backbones on different datasets and types of data. Moreover, we leverage these findings with a novel ensembling approach this complementarity of different backbones. Other works investigating differences across vision architectures include (Angarano et al., 2022; Pinto et al., 2022; 2021; Wang et al., 2022b) and others specific to CNNs (Abello et al., 2021; Hermann et al., 2020) and ViTs (Naseer et al., 2021). Finally, we find that there are works analysing the diversity of models for knowledge distillation Roth et al. (2024), complementarity on models trained don different subset of data Ramé et al. (2022) and Zhong et al. (2021) exploring the diversity on LLMs.

Model ensembling. Combining multiple machine learning models is a classical approach for improving predictive performance. Simply averaging the outputs of several models is a simple, effective technique (Bauer and Kohavi, 1999; Breiman, 1996; Dietterich, 2000; Lakshminarayanan et al., 2017) with studies dating back more than three decades ago. The approach has also been applied to deep neural networks (deep ensembles (Lakshminarayanan et al., 2017)), which has shown benefits in various domains including higher accuracy under distribution shift (Ovadia et al., 2019; Teney et al., 2018). The diversity of the combined models (in terms of uncorrelated errors) has also been shown to be critical to these improvements (Hao et al., 2024; Wortsman et al., 2022). The closest to our work is the SuperLearner framework (Cheng Ju and van der Laan, 2018), which similarly investigates multiple classical ensembles and introduces a new ensemble that learns temperature scaling. However, our study not only evaluates the performance of various classical ensembles but also explores the complementarity of CLIP backbones across 21 datasets. Furthermore, unlike SuperLearner, our method adapts its temperature scaling factors based on the input, which plays a crucial role in ensembling CLIP backbones.

6 CONCLUSIONS

This paper presents an analysis of vision backbones in the CLIP framework, focusing on the task of image classification. Unlike prior studies that span various downstream tasks (e.g. Goldblum et al. (2023)) our emphasis lies in identifying the unique strengths of various backbones. Our experiments revealed a distinctive complementarity across architectures and an avenue for enhancing CLIP’s performance by synergistic combination. We proposed an ensemble approach that reweights the logits from each backbone condition on the input data to yield more accurate predictions in image classification across a variety of datasets.

Limitations. First, our evaluation focused mostly on image classification, even though CLIP can be used for a variety of downstream tasks. Although we test our method on out-of-distribution Sec. E and also obtain the upper bound performance on Image-Text retrieval Sec. K, more work needs to be done on how to combine the representations. Second, our approach relies on a late fusion of backbones by adaptively adjusting their logits. Alternative approaches that leverage the same initial motivation use knowledge distillation, which requires training the student backbone again Roth et al. (2024) Third, although we employed the Cascade method (Wang et al., 2022a) to reduce the computations, our method still needs to pass a test image multiple times through vision encoders. This adds complexity and computational overhead, which can be important if a large number of backbones are considered. A possible solution would involve fusing backbones at the training stage to require only a single forward pass at test time.

Future work. We could emulate the behaviour of a given backbone by learning their representations using the main branch’s early layers. Other efficient methods could also be developed to combine predictions from multiple backbones to enhance the scalability of the idea.

540 REFERENCES
541

- 542 A. A. Abello, R. Hirata, and Z. Wang. Dissecting the high-frequency bias in convolutional neural net-
543 works. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
544 pages 863–871, 2021.
- 545 S. Angarano, M. Martini, F. Salvetti, V. Mazzia, and M. Chiaberge. Back-to-bones: Rediscovering
546 the role of backbones in domain generalization. *arXiv preprint arXiv:2209.01121*, 2022.
- 547 E. Bauer and R. Kohavi. An empirical comparison of voting classification algorithms: Bagging,
548 boosting, and variants. *Machine learning*, 36:105–139, 1999.
- 549 L. Bossard, M. Guillaumin, and L. Van Gool. Food-101-mining discriminative components with
550 random forests. In *Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland,*
551 *September 6–12, 2014, Proceedings, Part VI 13*, pages 446–461. Springer, 2014.
- 552 L. Breiman. Bagging predictors. *Machine learning*, 24:123–140, 1996.
- 553 T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan,
554 P. Shyam, G. Sastry, A. Askell, et al. Language models are few-shot learners. *arXiv preprint*
555 *arXiv:2005.14165*, 2020.
- 556 Z. Chen, Y. Deng, Y. Wu, Q. Gu, and Y. Li. Towards understanding the mixture-of-experts layer in
557 deep learning. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors,
558 *Advances in Neural Information Processing Systems*, volume 35, pages 23049–23062. Curran
559 Associates, Inc., 2022.
- 560 G. Cheng, J. Han, and X. Lu. Remote sensing image scene classification: Benchmark and state of the
561 art. *Proceedings of the IEEE*, 105(10):1865–1883, Oct 2017. ISSN 1558-2256.
- 562 A. B. Cheng Ju and M. van der Laan. The relative performance of ensemble methods with deep
563 convolutional neural networks for image classification. *Journal of Applied Statistics*, 45(15):2800–
564 2818, 2018. doi: 10.1080/02664763.2018.1441383. URL <https://doi.org/10.1080/02664763.2018.1441383>. PMID: 31631918.
- 565 M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing textures in the wild. In
566 *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- 567 A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning.
568 In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*,
569 pages 215–223. JMLR Workshop and Conference Proceedings, 2011.
- 570 J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical
571 image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages
572 248–255. Ieee, 2009.
- 573 L. Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal
574 Processing Magazine*, 29(6):141–142, 2012.
- 575 J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. Bert: Pre-training of deep bidirectional
576 transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 577 T. G. Dietterich. Ensemble methods in machine learning. In *International workshop on multiple
578 classifier systems*, pages 1–15. Springer, 2000.
- 579 A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani,
580 M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby. An image is worth 16x16 words:
581 Transformers for image recognition at scale. *ICLR*, 2021.
- 582 A. Fang, G. Ilharco, M. Wortsman, Y. Wan, V. Shankar, A. Dave, and L. Schmidt. Data deter-
583 mines distributional robustness in contrastive language image pre-training (clip). In *International
584 Conference on Machine Learning*, pages 6216–6234. PMLR, 2022.

- 594 S. Y. Gadre, G. Ilharco, A. Fang, J. Hayase, G. Smyrnis, T. Nguyen, R. Marten, M. Wortsman,
 595 D. Ghosh, J. Zhang, E. Orgad, R. Entezari, G. Daras, S. Pratt, V. Ramanujan, Y. Bitton, K. Marathe,
 596 S. Mussmann, R. Vencu, M. Cherti, R. Krishna, P. W. Koh, O. Saukh, A. Ratner, S. Song,
 597 H. Hajishirzi, A. Farhadi, R. Beaumont, S. Oh, A. Dimakis, J. Jitsev, Y. Carmon, V. Shankar, and
 598 L. Schmidt. Datacomp: In search of the next generation of multimodal datasets, 2023.
- 599 Z. Gan, L. Li, C. Li, L. Wang, Z. Liu, J. Gao, et al. Vision-language pre-training: Basics, recent
 600 advances, and future trends. *Foundations and Trends® in Computer Graphics and Vision*, 14(3–4):
 601 163–352, 2022.
- 602 P. Gao, S. Geng, R. Zhang, T. Ma, R. Fang, Y. Zhang, H. Li, and Y. Qiao. Clip-adapter: Better
 603 vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021.
- 604 M. Goldblum, H. Souri, R. Ni, M. Shu, V. Prabhu, G. Somepalli, P. Chattopadhyay, M. Ibrahim,
 605 A. Bardes, J. Hoffman, R. Chellappa, A. G. Wilson, and T. Goldstein. Battle of the backbones: A
 606 large-scale comparison of pretrained models across computer vision tasks, 2023.
- 607 C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger. On Calibration of Modern Neural Networks. In
 608 *Proceedings of the 34th International Conference on Machine Learning*, pages 1321–1330. PMLR,
 609 July 2017.
- 610 Y. Hao, Y. Lin, D. Zou, and T. Zhang. On the benefits of over-parameterization for out-of-distribution
 611 generalization. *arXiv preprint arXiv:2403.17592*, 2024.
- 612 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings
 613 of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- 614 K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick. Masked autoencoders are scalable vision
 615 learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
 616 pages 16000–16009, 2022.
- 617 P. Helber, B. Bischke, A. Dengel, and D. Borth. Introducing eurosat: A novel dataset and deep
 618 learning benchmark for land use and land cover classification. In *IGARSS 2018-2018 IEEE
 619 International Geoscience and Remote Sensing Symposium*, pages 204–207. IEEE, 2018.
- 620 D. Hendrycks, S. Basart, N. Mu, S. Kadavath, F. Wang, E. Dorundo, R. Desai, T. Zhu, S. Parajuli,
 621 M. Guo, D. Song, J. Steinhardt, and J. Gilmer. The many faces of robustness: A critical analysis of
 622 out-of-distribution generalization. *ICCV*, 2021a.
- 623 D. Hendrycks, K. Zhao, S. Basart, J. Steinhardt, and D. Song. Natural adversarial examples. *CVPR*,
 624 2021b.
- 625 K. Hermann, T. Chen, and S. Kornblith. The origins and prevalence of texture bias in convolutional
 626 neural networks. *Advances in Neural Information Processing Systems*, 33:19000–19015, 2020.
- 627 S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel. Detection of traffic signs in real-world
 628 images: The German Traffic Sign Detection Benchmark. In *International Joint Conference on
 629 Neural Networks*, 2013.
- 630 G. Ilharco, M. Wortsman, R. Wightman, C. Gordon, N. Carlini, R. Taori, A. Dave, V. Shankar,
 631 H. Namkoong, J. Miller, H. Hajishirzi, A. Farhadi, and L. Schmidt. Openclip. Technical report,
 632 July 2021. If you use this software, please cite it as below.
- 633 C. Jia, Y. Yang, Y. Xia, Y.-T. Chen, Z. Parekh, H. Pham, Q. Le, Y.-H. Sung, Z. Li, and T. Duerig.
 634 Scaling up visual and vision-language representation learning with noisy text supervision. In
 635 *International conference on machine learning*, pages 4904–4916. PMLR, 2021.
- 636 J. Johnson, B. Hariharan, L. Van Der Maaten, L. Fei-Fei, C. Lawrence Zitnick, and R. Girshick. Clevr:
 637 A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings
 638 of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910, 2017.
- 639 D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*,
 640 2014.

- 648 A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C.
649 Berg, W.-Y. Lo, P. Dollár, and R. Girshick. Segment anything, 2023.
- 650
- 651 J. Krause, J. Deng, M. Stark, and L. Fei-Fei. Collecting a large-scale dataset of fine-grained cars.
652 Technical report, 2013.
- 653 A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. Technical
654 report, 2009.
- 655
- 656 B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty
657 estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- 658 D. Lepikhin, H. Lee, Y. Xu, D. Chen, O. Firat, Y. Huang, M. Krikun, N. Shazeer, and Z. Chen.
659 Gshard: Scaling giant models with conditional computation and automatic sharding. *arXiv preprint*
660 *arXiv:2006.16668*, 2020.
- 661
- 662 C. Li, H. Liu, L. H. Li, P. Zhang, J. Aneja, J. Yang, P. Jin, H. Hu, Z. Liu, Y. J. Lee, and J. Gao.
663 Elevater: A benchmark and toolkit for evaluating language-augmented visual models. *Neural
664 Information Processing Systems*, 2022a.
- 665 F.-F. Li, M. Andreetto, M. Ranzato, and P. Perona. Caltech 101, Apr 2022b.
- 666
- 667 G. Li, N. Duan, Y. Fang, M. Gong, and D. Jiang. Unicoder-vl: A universal encoder for vision
668 and language by cross-modal pre-training. In *Proceedings of the AAAI conference on artificial
669 intelligence*, 2020a.
- 670 X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei, et al. Oscar:
671 Object-semantics aligned pre-training for vision-language tasks. In *Computer Vision–ECCV 2020:
672 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXX 16*, pages
673 121–137. Springer, 2020b.
- 674 H. Liu, C. Li, Y. Li, and Y. J. Lee. Improved baselines with visual instruction tuning, 2023.
- 675
- 676 S. Maji, J. Kannala, E. Rahtu, M. Blaschko, and A. Vedaldi. Fine-grained visual classification of
677 aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- 678 P. Mayilvahanan, T. Wiedemer, E. Rusak, M. Bethge, and W. Brendel. Does clip’s generalization
679 performance mainly stem from high train-test similarity? *arXiv preprint arXiv:2310.09562*, 2023.
- 680
- 681 S. Menon and C. Vondrick. Visual classification via description from large language models. *ICLR*,
682 2023.
- 683 M. M. Naseer, K. Ranasinghe, S. H. Khan, M. Hayat, F. Shahbaz Khan, and M.-H. Yang. Intriguing
684 properties of vision transformers. *Advances in Neural Information Processing Systems*, 34:23296–
685 23308, 2021.
- 686
- 687 M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In
688 *2008 Sixth Indian conference on computer vision, graphics & image processing*, pages 722–729.
689 IEEE, 2008.
- 690 Y. Ovadia, E. Fertig, J. Ren, Z. Nado, D. Sculley, S. Nowozin, J. Dillon, B. Lakshminarayanan, and
691 J. Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset
692 shift. In *Advances in Neural Information Processing Systems*, pages 13991–14002, 2019.
- 693 O. M. Parkhi, A. Vedaldi, A. Zisserman, and C. Jawahar. Cats and dogs. In *2012 IEEE conference on
694 computer vision and pattern recognition*, pages 3498–3505. IEEE, 2012.
- 695
- 696 F. Pinto, P. Torr, and P. K. Dokania. Are vision transformers always more robust than convolutional
697 neural networks? In *Advances in Neural Information Processing Systems*, 2021.
- 698 F. Pinto, P. H. Torr, and P. K. Dokania. An impartial take to the cnn vs transformer robustness contest.
699 In *European Conference on Computer Vision*, pages 466–480. Springer, 2022.
- 700
- 701 J. Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized
likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.

- 702 A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin,
 703 J. Clark, et al. Learning transferable visual models from natural language supervision. In
 704 *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- 705 A. Ramé, K. Ahuja, J. Zhang, M. Cord, L. Bottou, and D. Lopez-Paz. Model ratatouille: Recycling
 706 diverse models for out-of-distribution generalization. *arXiv preprint arXiv:2212.10445*, 2022.
- 707 D. Rau. Sparsely-gated mixture-of-experts pytorch implementation, 2019.
- 708 B. Recht, R. Roelofs, L. Schmidt, and V. Shankar. Do ImageNet classifiers generalize to ImageNet?
 709 In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference
 710 on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5389–5400.
 711 PMLR, 09–15 Jun 2019.
- 712 K. Roth, L. Thede, A. S. Koepke, O. Vinyals, O. J. Henaff, and Z. Akata. Fantastic gains and
 713 where to find them: On the existence and prospect of general knowledge transfer between any
 714 pretrained model. In *The Twelfth International Conference on Learning Representations*, 2024.
 715 URL <https://openreview.net/forum?id=m50eKHcttz>.
- 716 S. Santurkar, Y. Dubois, R. Taori, P. Liang, and T. Hashimoto. Is a caption worth a thousand
 717 images? a study on representation learning. In *The Eleventh International Conference on Learning
 718 Representations*, 2022.
- 719 C. Schuhmann, R. Beaumont, R. Vencu, C. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta,
 720 C. Mullis, M. Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation
 721 image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022.
- 722 P. Shi, M. C. Welle, M. Björkman, and D. Krägic. Towards understanding the modality gap in CLIP.
 723 In *ICLR 2023 Workshop on Multimodal Representation Learning: Perks and Pitfalls*, 2023.
- 724 R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts. Recursive deep
 725 models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013
 726 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle,
 727 Washington, USA, Oct. 2013. Association for Computational Linguistics.
- 728 D. Teney, P. Anderson, X. He, and A. Van Den Hengel. Tips and tricks for visual question answering:
 729 Learnings from the 2017 challenge. In *Proceedings of the IEEE conference on computer vision
 730 and pattern recognition*, pages 4223–4232, 2018.
- 731 H. Touvron, T. Lavigil, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal,
 732 E. Hambro, F. Azhar, A. Rodriguez, A. Joulin, E. Grave, and G. Lample. Llama: Open and efficient
 733 foundation language models. *ArXiv*, abs/2302.13971, 2023.
- 734 W. Tu, W. Deng, and T. Gedeon. A closer look at the robustness of contrastive language-image
 735 pre-training (clip). *Advances in Neural Information Processing Systems*, 36, 2024.
- 736 N. Varshney and C. Baral. Model cascading: Towards jointly improving efficiency and accuracy of
 737 NLP systems. In Y. Goldberg, Z. Kozareva, and Y. Zhang, editors, *Proceedings of the 2022 Confer-
 738 ference on Empirical Methods in Natural Language Processing*, pages 11007–11021, Abu Dhabi,
 739 United Arab Emirates, Dec. 2022. Association for Computational Linguistics. doi: 10.18653/v1/
 740 2022.emnlp-main.756. URL <https://aclanthology.org/2022.emnlp-main.756>.
- 741 B. S. Veeling, J. Linmans, J. Winkens, T. Cohen, and M. Welling. Rotation equivariant CNNs for
 742 digital pathology. *arXiv preprint arXiv:1806.03962*, June 2018.
- 743 C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. Caltech-ucsd birds-200-2011. Technical
 744 Report CNS-TR-2011-001, California Institute of Technology, 2011.
- 745 J. Wang, K. Markert, M. Everingham, et al. Learning models for object recognition from natural
 746 language descriptions. In *BMVC*, 2009.
- 747 X. Wang, D. Kondratyuk, E. Christiansen, K. M. Kitani, Y. Movshovitz-Attias, and E. Eban. Wisdom
 748 of committees: An overlooked approach to faster and more accurate models. In *International
 749 Conference on Learning Representations*, 2022a. URL <https://openreview.net/forum?id=MvO2t0vbs4->.

- 756 Z. Wang, Y. Bai, Y. Zhou, and C. Xie. Can CNNs be more robust than transformers? *arXiv preprint*
 757 *arXiv:2206.03452*, 2022b.
- 758
- 759 M. Wortsman, G. Ilharco, S. Y. Gadre, R. Roelofs, R. Gontijo-Lopes, A. S. Morcos, H. Namkoong,
 760 A. Farhadi, Y. Carmon, S. Kornblith, et al. Model soups: averaging weights of multiple fine-tuned
 761 models improves accuracy without increasing inference time. In *International Conference on*
 762 *Machine Learning*, pages 23965–23998. PMLR, 2022.
- 763 J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene
 764 recognition from abbey to zoo. In *2010 IEEE Computer Society Conference on Computer Vision*
 765 and *Pattern Recognition*, pages 3485–3492, June 2010.
- 766
- 767 J. Yu, Z. Wang, V. Vasudevan, L. Yeung, M. Seyedhosseini, and Y. Wu. Coca: Contrastive captioners
 768 are image-text foundation models. *arXiv preprint arXiv:2205.01917*, 2022.
- 769
- 770 X. Zhai, X. Wang, B. Mustafa, A. Steiner, D. Keysers, A. Kolesnikov, and L. Beyer. Lit: Zero-shot
 771 transfer with locked-image text tuning. In *Proceedings of the IEEE/CVF Conference on Computer*
 772 *Vision and Pattern Recognition (CVPR)*, pages 18123–18133, June 2022.
- 773
- 774 R. Zhang, R. Fang, P. Gao, W. Zhang, K. Li, J. Dai, Y. Qiao, and H. Li. Tip-adapter: Training-free
 775 clip-adapter for better vision-language modeling. *arXiv preprint arXiv:2111.03930*, 2021.
- 776
- 777 R. Zhong, D. Ghosh, D. Klein, and J. Steinhardt. Are larger pretrained language models uniformly
 778 better? comparing performance at the instance level. In C. Zong, F. Xia, W. Li, and R. Navigli,
 779 editors, *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages
 3813–3827, Online, Aug. 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.
 findings-acl.334. URL <https://aclanthology.org/2021.findings-acl.334>.
- 780
- 781 K. Zhou, J. Yang, C. C. Loy, and Z. Liu. Learning to prompt for vision-language models. *International*
 782 *Journal of Computer Vision*, 130(9):2337–2348, 2022.
- 783
- 784 B. Zoph, I. Bello, S. Kumar, N. Du, Y. Huang, J. Dean, N. Shazeer, and W. Fedus. St-moe: Designing
 785 stable and transferable sparse expert models. *arXiv preprint arXiv:2202.08906*, 2022.
- 786
- 787
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- 795
- 796
- 797
- 798
- 799
- 800
- 801
- 802
- 803
- 804
- 805
- 806
- 807
- 808
- 809

972 **Table F.1: Classification accuracy of models un-**
 973 **der few-shot settings.**

Few-shot	1	2	4	8	16
Linear-probe CLIP Radford et al. (2021)	22.2	31.9	41.2	49.5	56.1
CoOPZhou et al. (2022)	47.6	50.9	56.2	59.9	63.0
CLIP-AdapterGao et al. (2021)	61.2	61.5	61.8	62.7	63.6
Tip-AdapterZhang et al. (2021)	60.7	61.0	61.0	61.5	62.0
Tip-Adapter-FZhang et al. (2021)	61.3	61.7	62.5	64.0	65.5
NLC	78.2	78.1	78.2	78.3	78.4

974 **Table F.2: Classification accuracy of models**
 975 **on various vision backbones using 16-shots.**

Models	ResNet			ViT	
	50	101	50x16	B-32	B-16
Zero-shot CLIP Radford et al. (2021)	60.3	62.5	70.9	63.8	68.7
CoOPZhou et al. (2022)	47.6	50.9	56.2	59.9	63.0
CoOP	63.0	66.6	-	66.9	71.9
CLIP-AdapterGao et al. (2021)	63.6	65.4	-	66.2	71.1
Tip-AdapterZhang et al. (2021)	62.0	64.8	73.0	65.6	70.8
Tip-Adapter-FZhang et al. (2021)	65.5	68.6	75.8	68.7	73.7
NLC				78.4	

980 E ROBUSTNESS UNDER DISTRIBUTION SHIFT OF IMGNET-1K

981
 982
 983 To assess how well the combination of backbones performs under varying conditions, we test its
 984 robustness using natural distribution shifts in the ImageNet dataset. We evaluate its performance on
 985 four datasets representing different distribution shifts: ImageNet-V2 Recht et al. (2019), ImageNet
 986 Adversarial Hendrycks et al. (2021b), ImageNet Rendition Hendrycks et al. (2021a) and ImageNet
 987 Sketch Wang et al. (2009). Specifically, we employ a NLC trained on the ZeroShot CLIP backbones
 988 using the original IMGNET-1K data Deng et al. (2009). This allows us to determine whether the
 989 learned combination of backbones can maintain its performance across these distribution shift datasets.
 990 In Table E.1, we observe that the learned combination of backbones, denoted as NLC, enhances
 991 performance in 3 out of 4 selected benchmarks, with improvements ranging from 0.3% to 2.1%.
 992 However, in the case of ImageNet Adversarial, the performance of the combination of backbones
 993 appears to suffer, possibly due to a more complex decision boundary.

994 F THE COMBINATION OF NLC WITH TIP-ADAPTER

995
 996 We integrate the Tip-Adapter Zhang et al. (2021) with our NLC ensembling mechanism and conduct
 997 experiments on EUROSAT, FGVC, DTD, and PETS datasets. Initially, we apply Tip-Adapter and
 998 Tip-Adapter-F independently on each backbone, following their protocol with 1, 2, 4, 8, and 16 shots.
 999 Subsequently, we employ our ensembling mechanism to fuse the adapted backbones. For training
 1000 NLC, we utilize the validation set used by Tip-Adapter to combine their logits with CLIP linearly.

1001
 1002 In Table F.1 and F.2, we compare the performance of different few-shot adapters of CLIP on the
 1003 ImageNet dataset. It shows that current few-shot adapter methods also have performance differences
 1004 with various backbones; still, our NLC surpassed the best backbone reported by their method.

1005
 1006 Table F.3, F.4, F.5, and F.6 show the zero-shot performance of the CLIP used by Tip-Adapter. It also
 1007 shows the Tip-Adapter and Tip-Adapter-F performance on each backbone and when we combine all
 1008 the Tip-Adapter versions and backbones with NLC. Across the four datasets used for this experiment,
 1009 NLC improve the performance of each version of Tip-Adapter. Notably, NLC for EUROSAT obtained
 1010 an improvement of up to 15% with respect to the best Tip-Adapter backbone (L-14-336) using 1 shot.

1026
1027
1028
1029
1030
1031
1032
1033
1034

1035 Table F.3: Tip-Adapter and Tip-Adapter-F fused Table F.4: Tip-Adapter and Tip-Adapter-F fused
1036 with our NLC and applied to EUROSAT dataset with our NLC and applied to FGVC dataset

1037

ZeroShot													
ResNet					ViT								
50	50x4	50x16	50x64	101	B-32	B-16	L-14	L-14-336					
37.5	32.0	40.3	49.4	32.5	45.2	47.6	58.1	63.5					
TiP-Adapter													
TiP-Adapter-F													
Shots													
Model		1	2	4	8	16	Model		1	2	4	8	16
ResNet	50	55.5	60.8	68.1	66.3	70.3	ResNet	50	19.0	21.4	22.5	25.9	29.6
	50x4	53.0	56.8	67.0	68.0	71.7		50x4	24.1	27.8	29.3	31.5	37.8
	50x16	54.5	61.0	62.7	65.1	71.4		50x16	29.0	32.4	35.7	38.5	42.6
	50x64	65.1	66.5	70.5	73.5	73.9		50x64	35.0	38.5	41.0	44.8	49.9
	101	47.9	49.7	60.6	64.7	67.6		101	20.4	22.9	24.8	28.3	32.1
ViT	B-32	54.1	61.7	68.2	70.6	69.8	ViT	B-32	21.6	23.9	23.7	27.2	31.6
	B-16	66.9	68.8	75.1	72.9	78.1		B-16	28.0	31.1	32.3	36.7	39.6
	L-14	73.0	71.1	78.4	79.2	82.4		L-14	35.6	41.7	44.9	48.9	52.3
	L-14-336	73.5	72.9	75.7	79.2	80.4		L-14-336	37.7	43.6	46.5	50.0	53.8
	With SL	77.8	81.5	85.6	86.4	85.4		With SL	39.1	44.2	47.7	51.2	54.8
With NLC		88.4	86.8	89.7	90.1	90.7		With NLC	40.9	44.5	47.9	51.5	55.0
TiP-Adapter-F													
Shots													
Model		1	2	4	8	16	Model		1	2	4	8	16
ResNet	50	60.7	64.4	73.3	77.7	84.9	ResNet	50	20.5	22.9	26.5	30.4	35.5
	50x4	59.5	61.9	76.2	81.9	84.9		50x4	26.3	29.1	32.8	36.8	42.2
	50x16	61.4	68.7	75.9	80.8	83.7		50x16	31.2	36.9	38.3	43.6	49.4
	50x64	71.2	69.7	78.7	81.4	86.6		50x64	37.0	40.8	45.0	49.8	54.8
	101	62.7	57.7	75.4	78.8	83.6		101	21.7	24.0	26.9	32.0	38.0
ViT	B-32	59.4	70.1	76.5	79.9	84.9	ViT	B-32	22.7	25.3	27.5	32.9	36.9
	B-16	66.6	71.0	79.1	83.8	88.9		B-16	30.2	34.1	36.1	40.9	44.6
	L-14	74.9	75.0	86.1	86.4	90.6		L-14	38.6	44.1	48.5	51.9	57.4
	L-14-336	72.5	76.4	86.1	85.3	91.0		L-14-336	40.9	45.2	49.6	52.7	59.0
	With SL	84.5	86.7	89.9	89.0	93.1		With SL	46.7	49.1	53.4	57.4	62.4
With NLC		89.7	88.3	91.4	91.1	93.8		With NLC	46.8	49.2	53.6	57.5	62.6

1071
1072
1073
1074
1075
1076
1077
1078
1079

1134 Table G.1: Ablation of NLC performance when changing the number of samples used to train.
 1135 We use $\text{NLC}(n)$ to denote NLC with n samples per class. and showcase the improvement in
 1136 performance compared with SINGLE-BEST backbone in a zero-shot setting.

	CALTECH101																		Mean Δ	Min Δ	Max Δ			
	CARS	CIFAR10	CIFAR100	CLEVR	COUNTRY211	CUB	DTD	EUROSAT	FGVC	FLOWERS	FOOD	GTSRB	IMGNET-1K	MNIST	PCAM	PETS	SST2	RESISC45	STL10	SUN397				
SINGLE-BEST	86.6	79.4	95.6	75.8	24.4	34.5	63.0	56.4	48.0	33.2	79.1	93.1	52.4	76.6	78.9	63.9	93.8	71.0	64.6	99.4	67.7			
NNC(1)	87.3	82.3	93.6	78.0	25.1	36.0	71.0	59.3	55.0	37.1	79.4	93.6	54.8	77.6	83.3	62.5	94.3	73.8	70.0	98.9	71.1	2.2	-2.0	8.0
NNC(2)	87.3	81.4	93.0	78.2	25.1	33.2	70.0	59.9	52.3	37.2	79.5	93.5	55.3	77.4	84.4	62.5	94.5	73.1	70.2	98.9	71.1	1.9	-2.6	7.0
NNC(4)	87.1	81.3	95.3	75.9	26.9	35.4	70.2	59.9	54.0	35.9	78.2	93.1	55.2	78.1	81.6	62.7	94.5	72.7	68.8	99.1	71.2	1.9	-1.3	7.2
NNC(8)	87.1	81.6	95.6	76.2	26.9	36.1	70.3	58.1	54.0	37.0	80.3	93.1	54.9	78.3	82.5	67.0	94.2	72.4	68.9	99.5	71.2	2.3	0.0	7.3
NNC(16)	87.1	82.8	95.8	77.2	27.0	36.2	71.3	60.4	53.5	37.0	80.4	93.5	55.2	78.4	83.3	69.6	94.2	73.6	69.6	99.4	72.2	2.9	0.0	8.3
NNC(32)	87.3	82.9	96.0	78.2	27.2	36.1	71.2	60.4	51.6	37.2	80.4	94.2	55.1	78.1	85.3	75.5	94.4	74.1	70.1	99.4	72.4	3.3	0.0	11.6

G TRAINING NLC UNDER LIMITED SAMPLES

In our exploration of the effectiveness of the NLC approach, we extend our analysis to a scenario where we limit the samples to combine the zero-shot CLIPs. This experiment allows us to assess the adaptability and performance of our proposed method under limited training data conditions. Table G.1 presents the performance of NLC by means of a limited number of samples. Although the performance of NLC overall improves when it has more data available to combine the backbones, in most cases, just using one sample NLC(1) per class is enough to improve its performance. Notably, there is a stop in performance degradation when we use more than 8 samples NLC(8) in all the benchmarks.

Moreover, we run five different seeds to train NLC on the different few-shots settings. Results show that the standard deviation of our obtained performance is stable and low across all of our studied datasets, suggesting that our approach is not sensitive to the effects of different samples. It varies with the number of samples used with a mean standard deviation of 0.62 when we used 1 shot and 0.48 when we used 32 shots.

H LEARNED TEMPERATURE VALUES

Figure H.1 and H.2 present the distribution of the temperature values using a box plot for the NLC method, normalized by their maximum value. Notably, there is a dominance in the temperature values for the best backbones in each family, *i.e.*, RN50x64 and ViT-L-14 or ViT-L-14-336. Particularly prominent in STL10, PETS, DTD, FGVC and FLOWERS, making these backbones especially relevant for NLC. Interestingly, it is observed that the most weighted backbone within the ResNet family is not consistently ResNet-101, despite its deeper architecture. This observation is evident in datasets such as PETS, CUB, FGVC and SUN397, where the mean value of the temperature corresponding to ResNet-50 surpasses that of ResNet-101.

Furthermore, the distribution of the temperature values across backbones for datasets such as CLEVR, EUROSAT, and GTSRB is more uniform compared to other datasets. This suggests that the NLC method is effectively leveraging the strengths of each backbone to arrive at accurate labels for each sample, resulting in a more balanced distribution of weights across different backbones, complementing our analysis on Fig. 2.

I DIVERSITY AND ORACLE PERFORMANCE OF DIFFERENT COMBINATIONS

Table I.1 shows the Oracle performance by combinaning different options of diversity. Table I.2 present the diversity of different source of complementarity.

J CASCADING AND ENSEMBLE

Figure J.3 shows the performance of multiple ensemble and cascading methods of 2 to 9 backbones, notice that NLC obtains the best performance, and when we add cascade it maintains the computational requirements

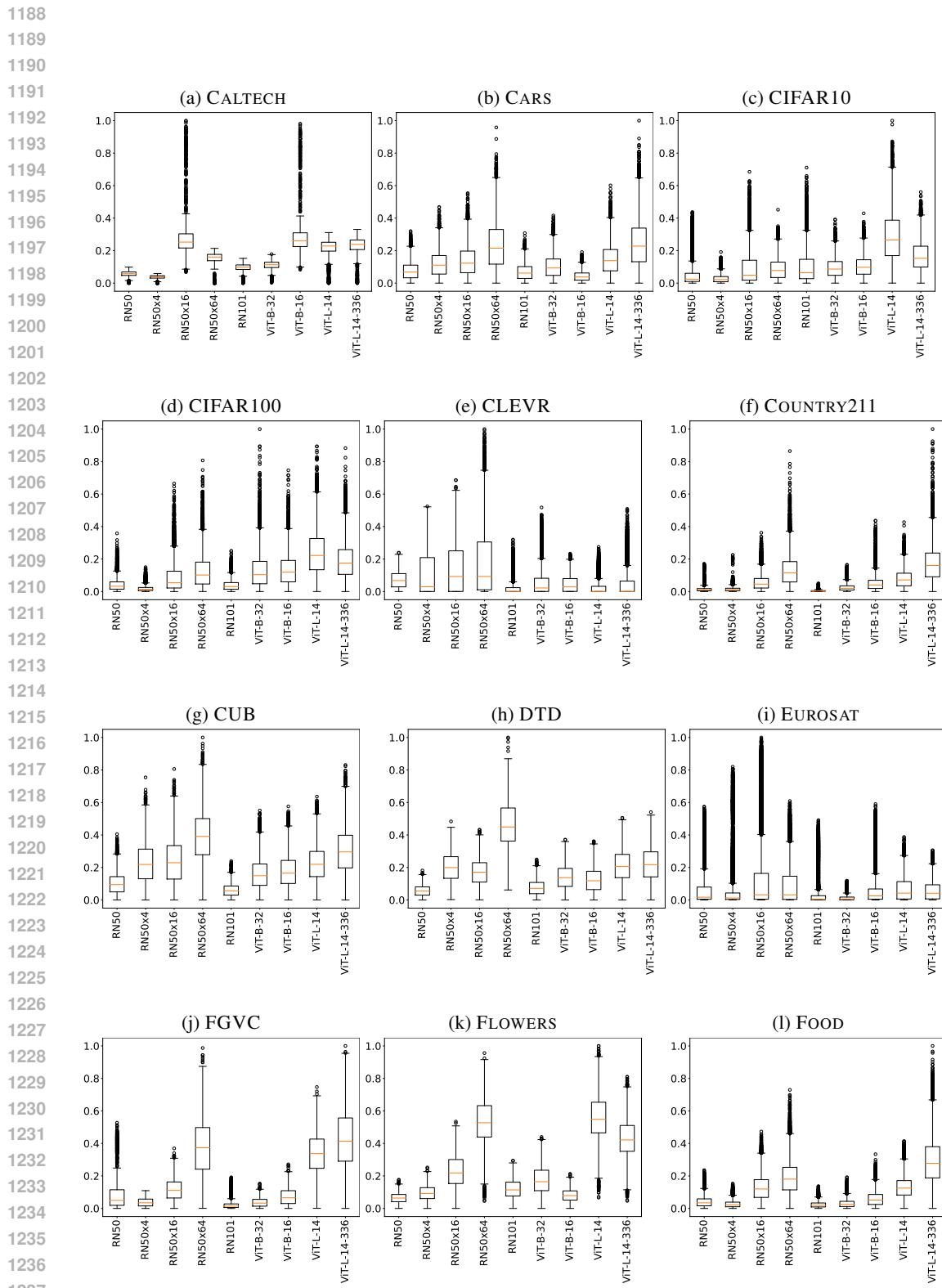


Figure H.1: Alpha values for each dataset using NLC

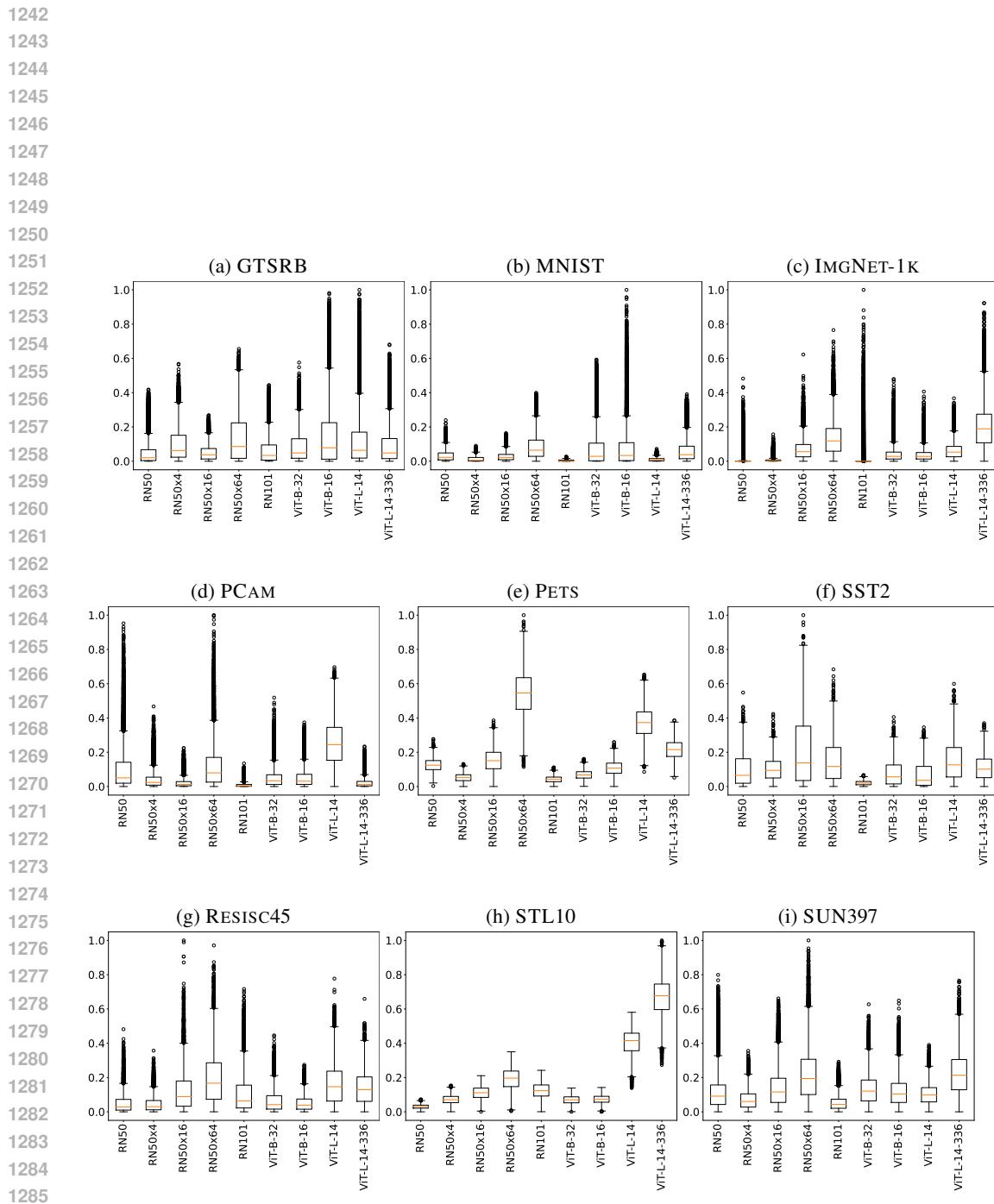
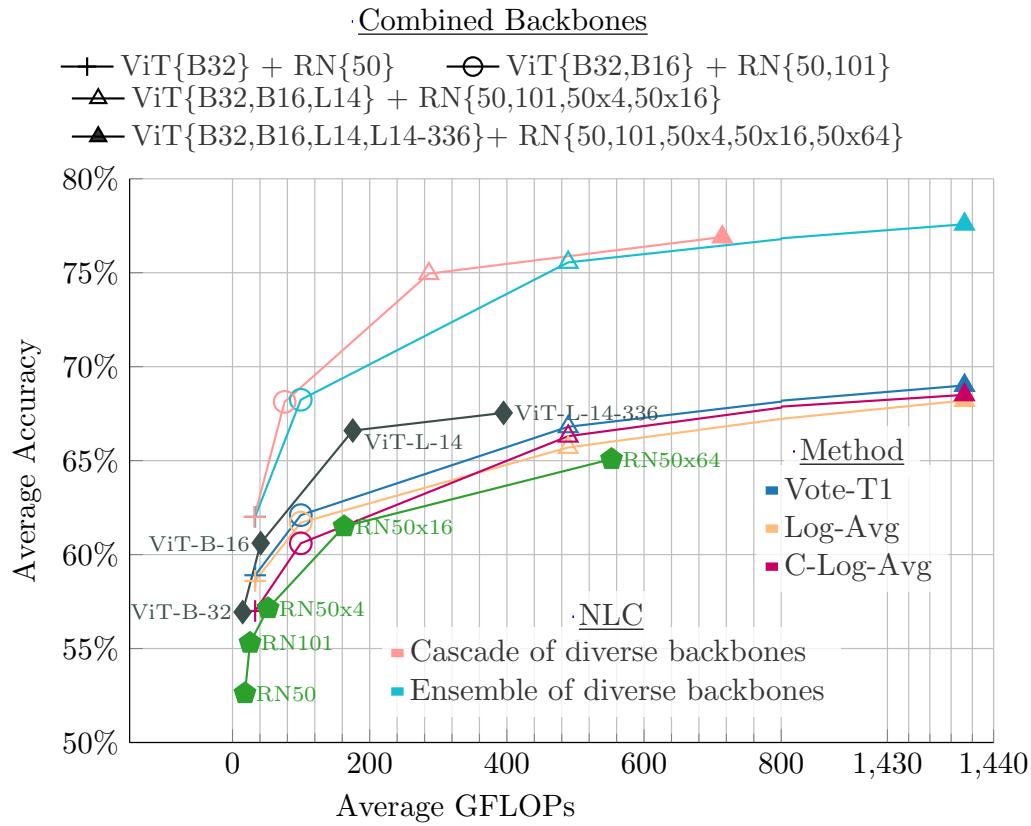


Figure H.2: Alpha values for each dataset using NLC



1388
1389 Figure J.3: Average accuracy across 21 datasets for zero-shot ResNets and ViTs backbones, LOG-
1390 AVG Ensemble, VOTE T-1 Ensemble, C-LOG-AVG Ensemble, NLC Ensemble, and NLC Cascade
1391 using 2 to 9 backbones, plotted against the Average GFLOPs. Demonstrates that our NLC ensembles
1392 can surpass the best zero-shot backbone, and the NLC cascade can also surpass the best zero-shot
1393 backbone with fewer GFLOPs. Moreover, our method surpasses the standard ensembling techniques.
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

1404 L OVERLAP DIAGRAMS FOR OTHER DATASETS. 1405

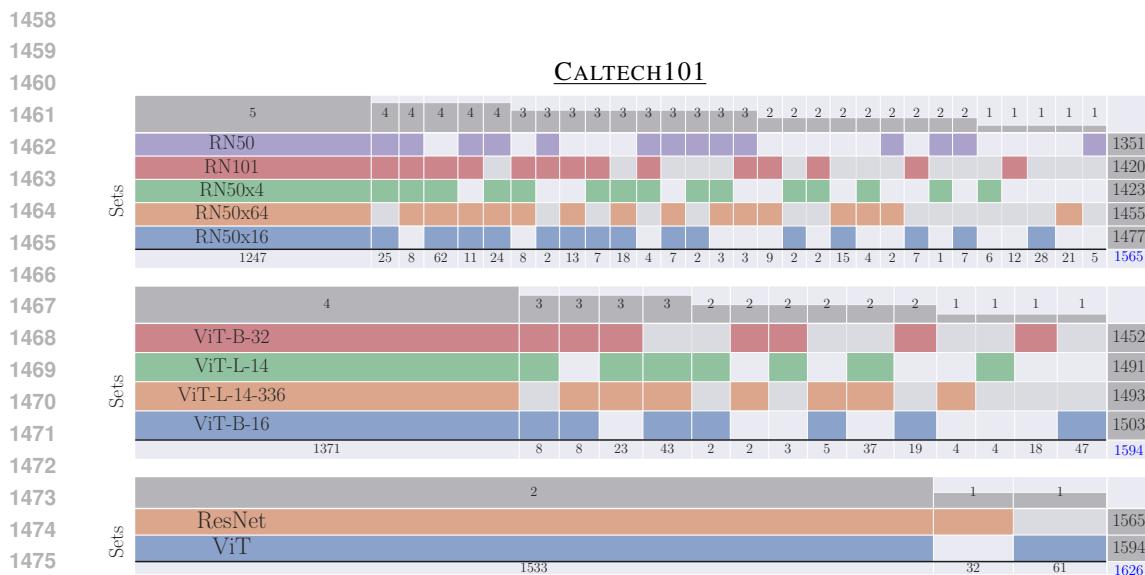
1406 In this section, we present the linear Venn Diagrams for each of the other datasets used in the
1407 experiment section CALTECH (Figure L.4), CARS (Figure L.5), CUB (Figure L.10), CIFAR10 (Figure
1408 L.6), CIFAR100 (Figure L.7), CLEVR (Figure L.8), COUNTRY211 (Figure L.9), CUB (Figure L.10),
1409 DTD (Figure L.11), EUROSAT(Figure L.12), FGVC (Figure L.13), FLOWERS (Figure L.14), FOOD
1410 (Figure L.15), GTSRB (Figure L.16), MNIST (Figure L.17), PCAM (Figure L.18), PETS (Figure L.19),
1411 RenderedSST2 (Figure L.20), RESISC45 (Figure L.21), STL10 (Figure L.22), and SUN397 (Figure
1412 L.23). We can see that in each dataset, the CLIP backbones present possible complementarities that
1413 could be exploited.

1414 M POSSIBLE COMBINATIONS 1415

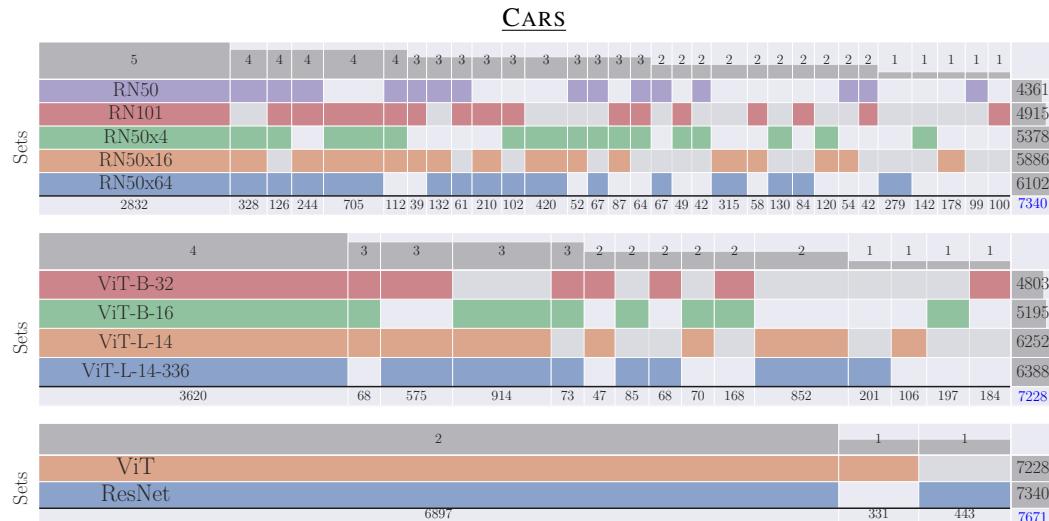
1416 Tables M.4, M.5, M.6, M.7, M.8, M.9, M.10, and M.11 present the results of possible combinations
1417 of backbones using the non-parametric and parametric approaches proposed in the paper. Notably,
1418 the performance of NLC consistently emerges as the best across various backbone combinations and
1419 datasets when compared to other methods.

1420 Notably, instances exist where the combination of specific backbones yields a more substantial
1421 performance boost than utilizing all backbones together. For instance, in the PETS dataset, combining
1422 ResNet 50, 101, and ViT-B-32 results in a delta improvement of 2.37%, surpassing the 0.99%
1423 improvement achieved by using the five backbones selected for this experiment. This phenomenon
1424 is consistent across datasets with different backbone combinations. In CARS, there is a boost of
1425 5.71% when combining ResNet-101 and ViT-B-32, compared to the 2.55% boost when using all
1426 five different backbones. While the best delta improvement among backbones may not necessarily
1427 come from combining all backbones, the best overall accuracy is consistently obtained when using
1428 the combination of all backbones.
1429

1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457



1477 Figure L.4: CALTECH101 Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.



1504 Figure L.5: Overlap diagrams for the CARS dataset with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

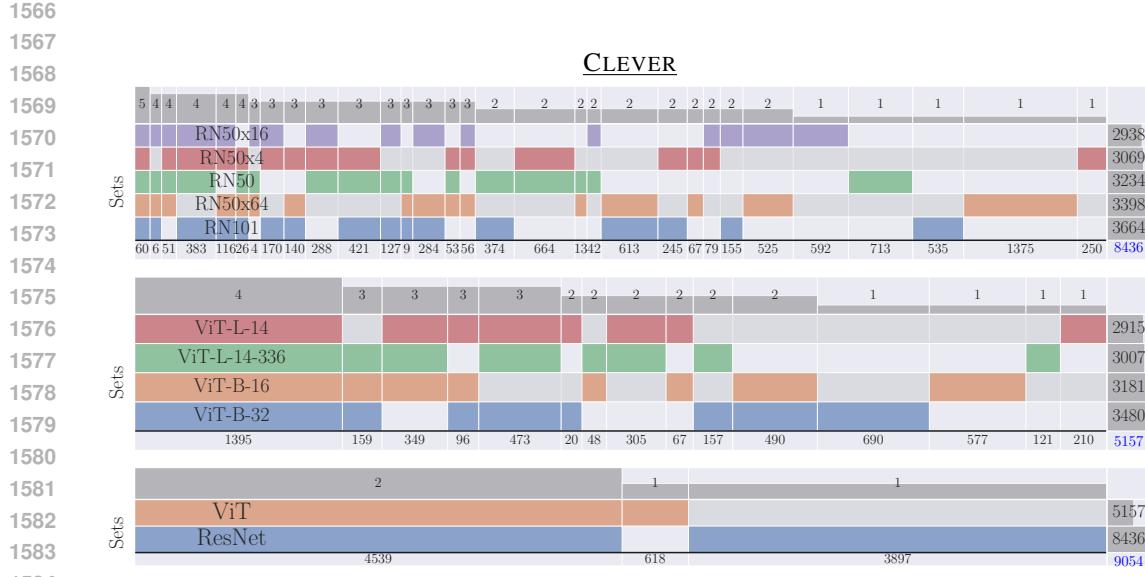


Figure L.8: CLEVER Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

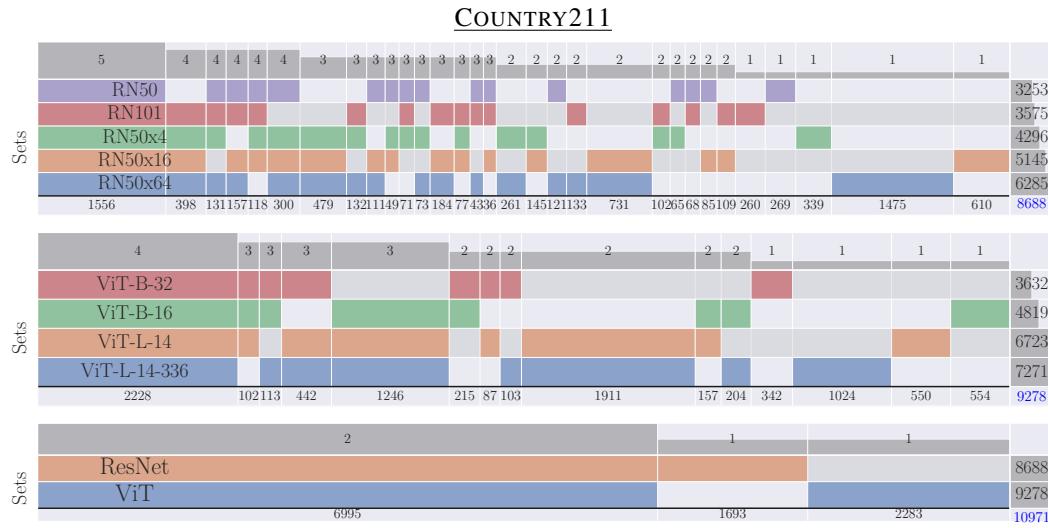


Figure L.9: COUNTRY211 Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

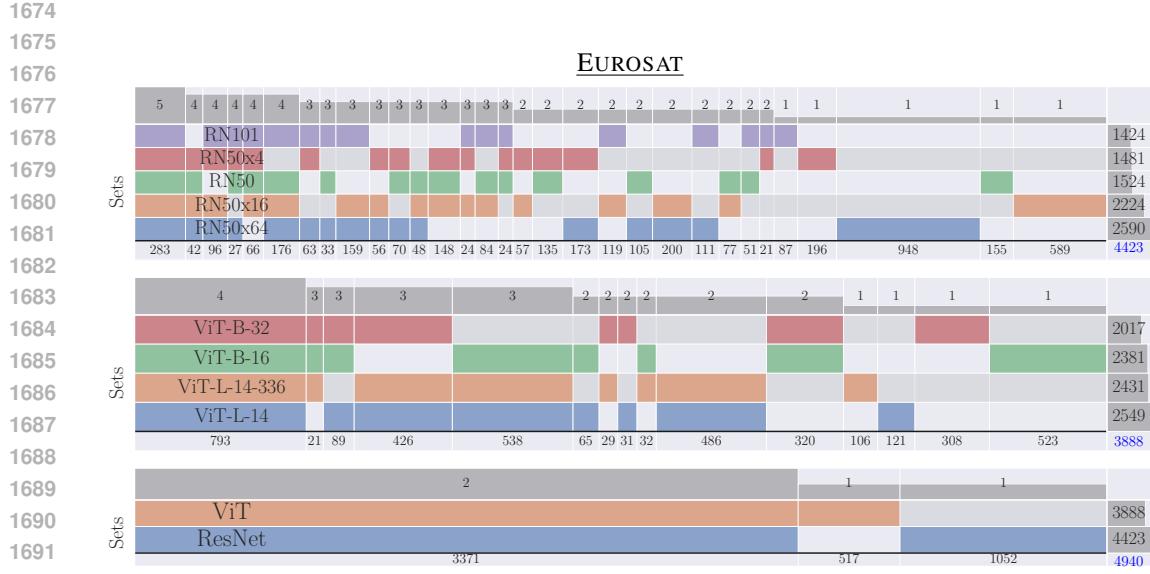


Figure L.12: EUROSAT Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

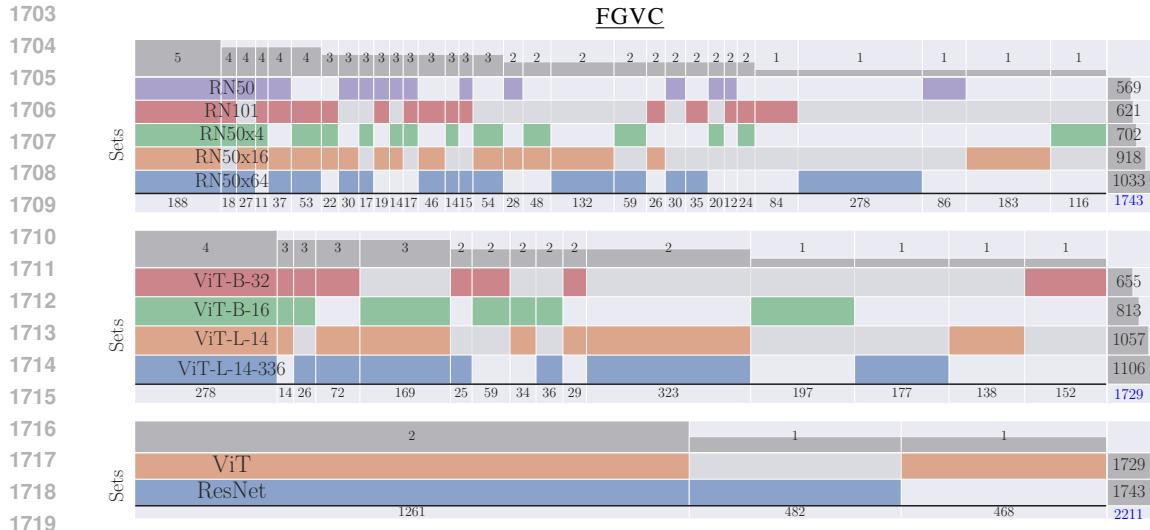


Figure L.13: FGVC Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

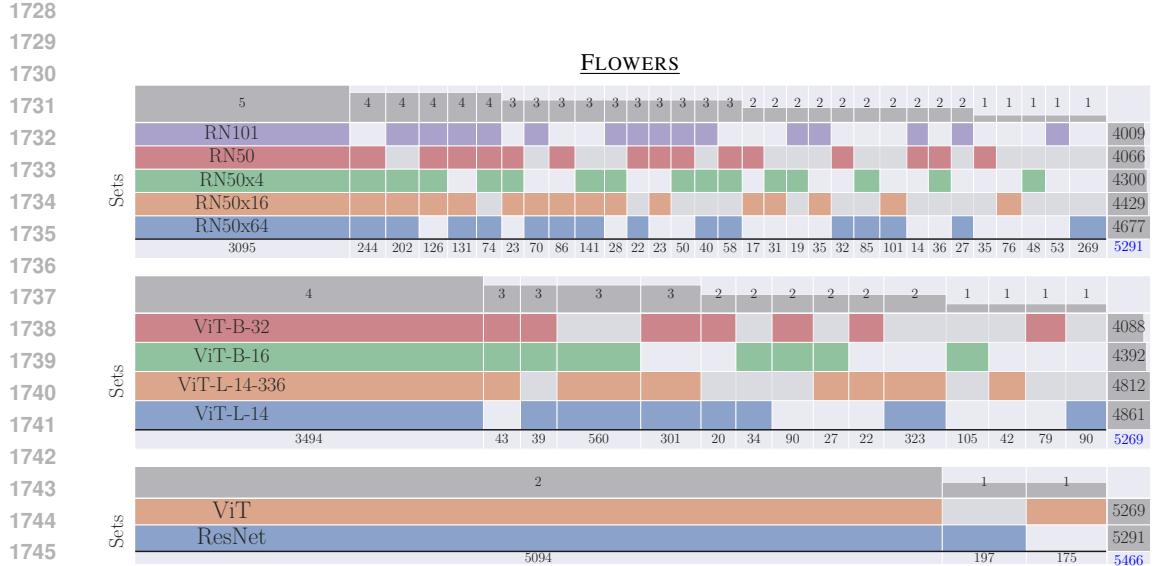


Figure L.14: FLOWERS Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

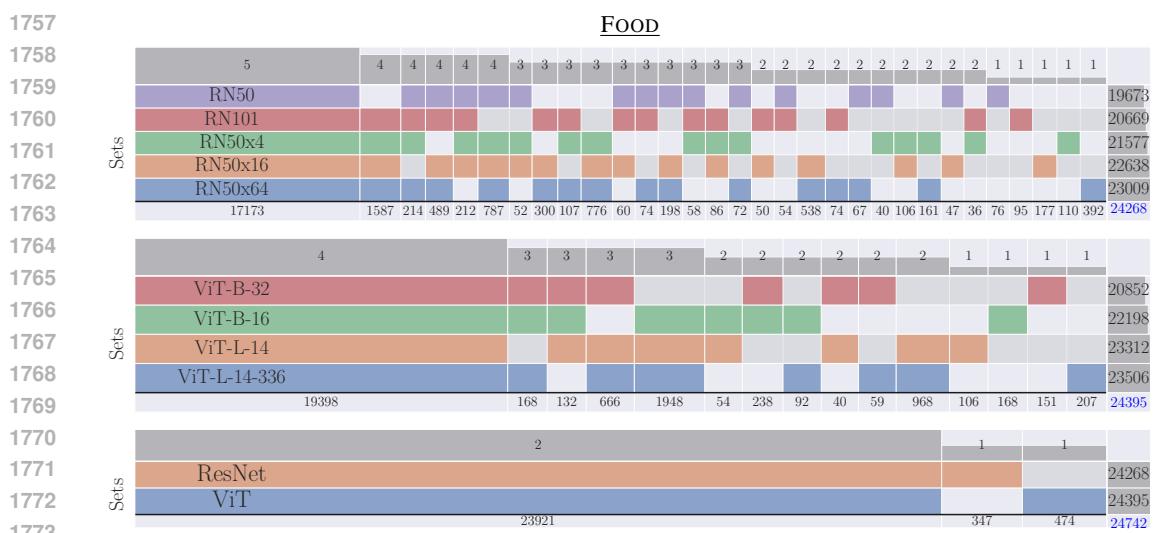


Figure L.15: FOOD Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

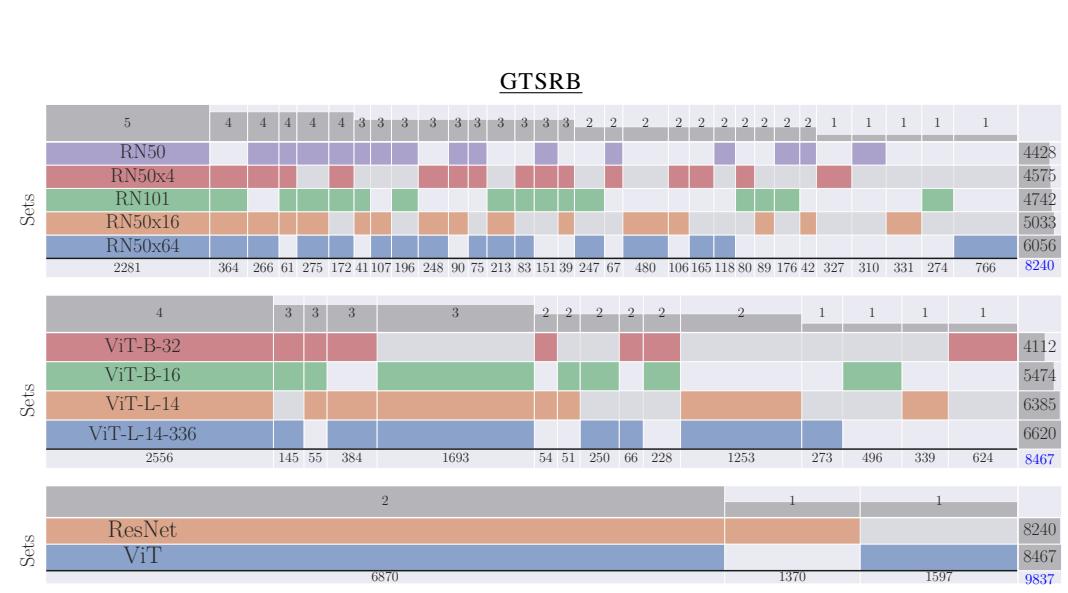


Figure L.16: GTSRB Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

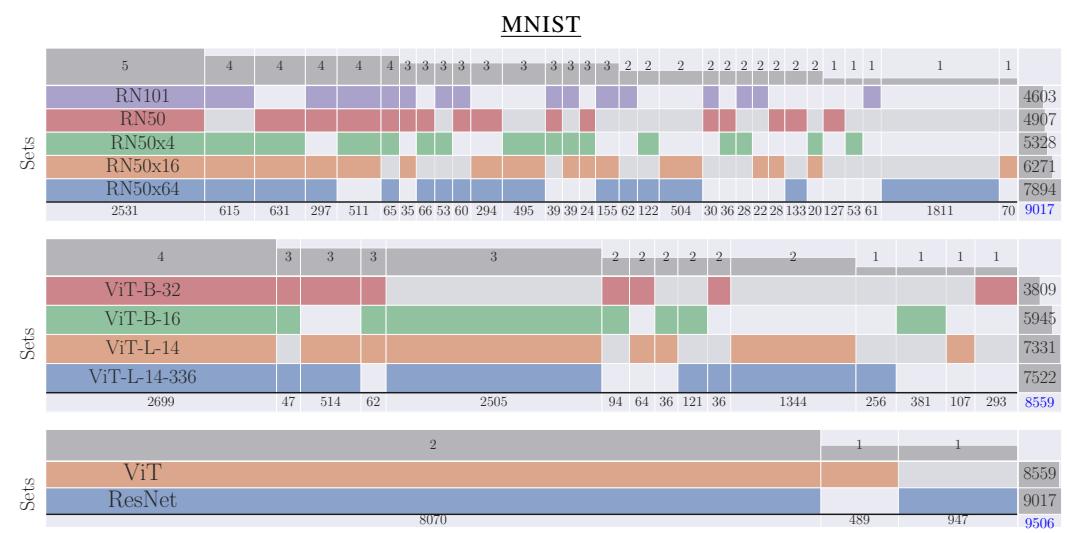


Figure L.17: MNIST Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

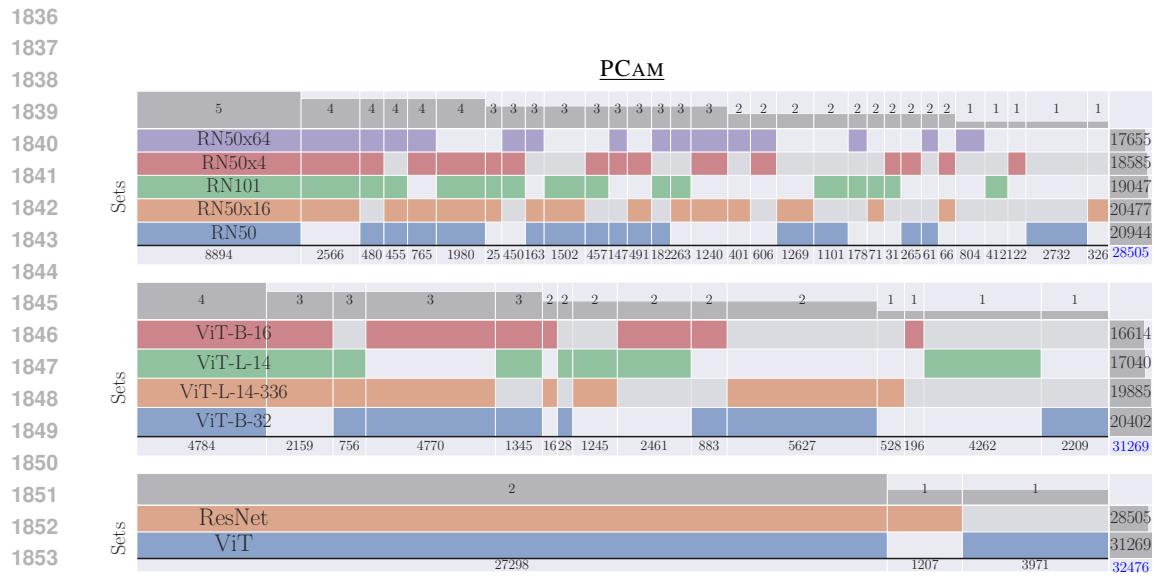


Figure L.18: PCAM Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

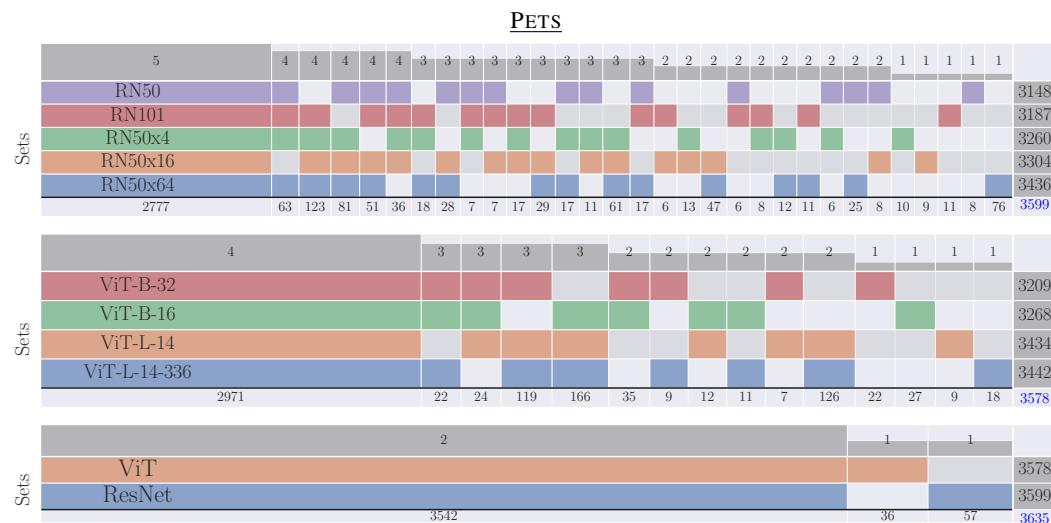


Figure L.19: PETS Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

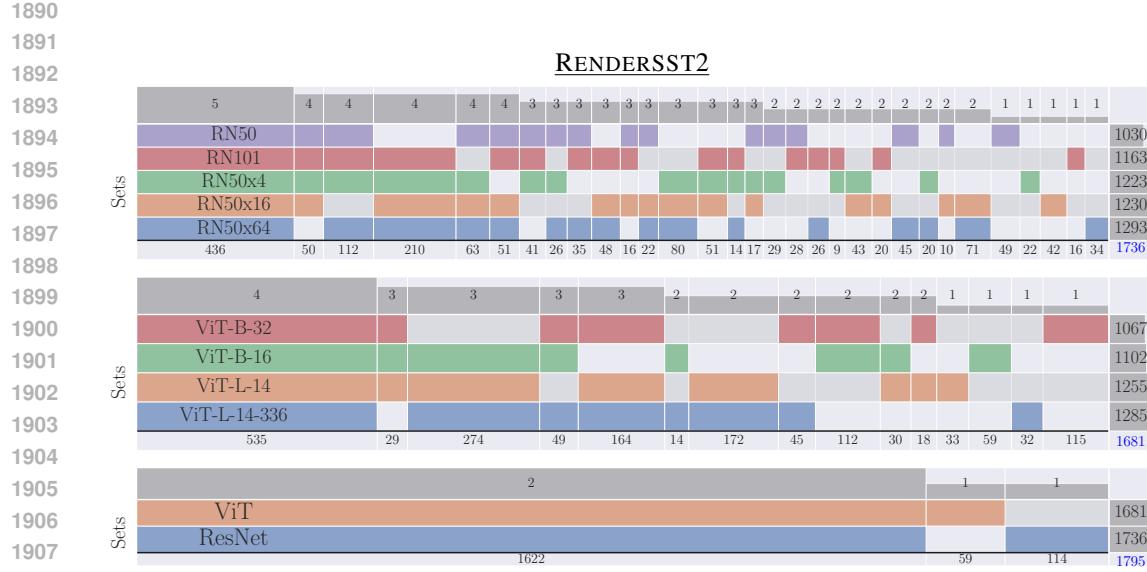


Figure L.20: RENDERSST2 Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

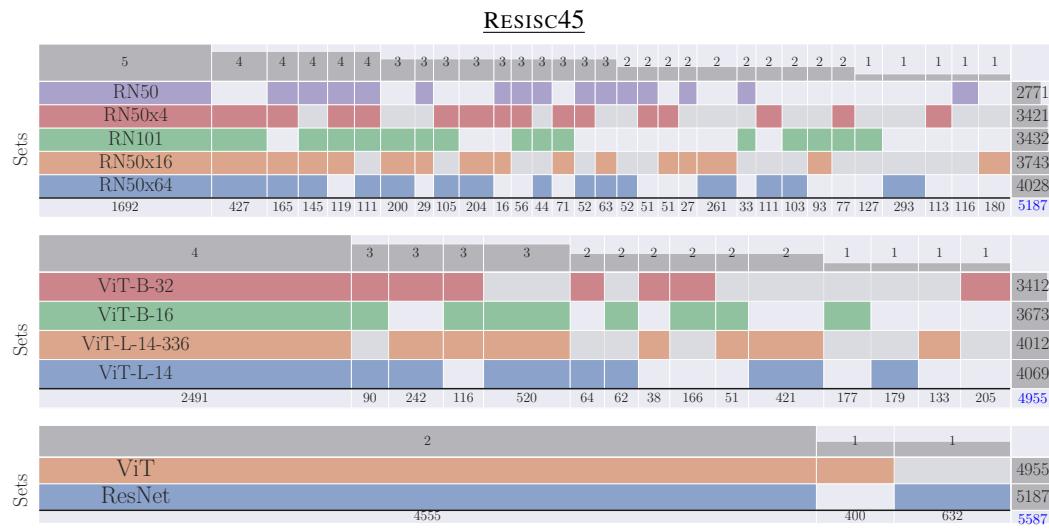


Figure L.21: RESISC45 Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

1943

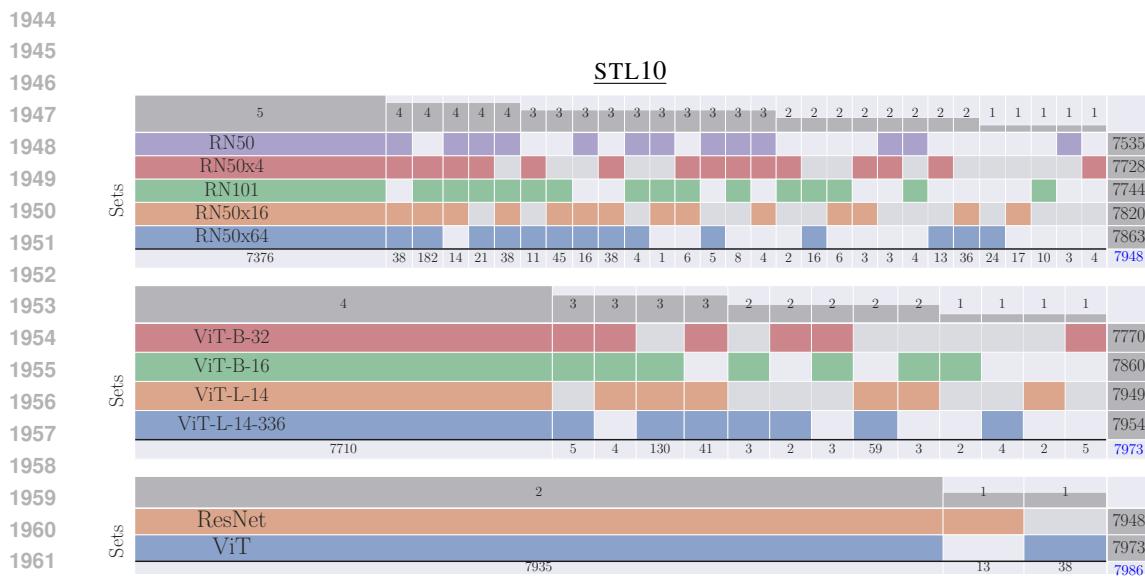


Figure L.22: STL10 Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

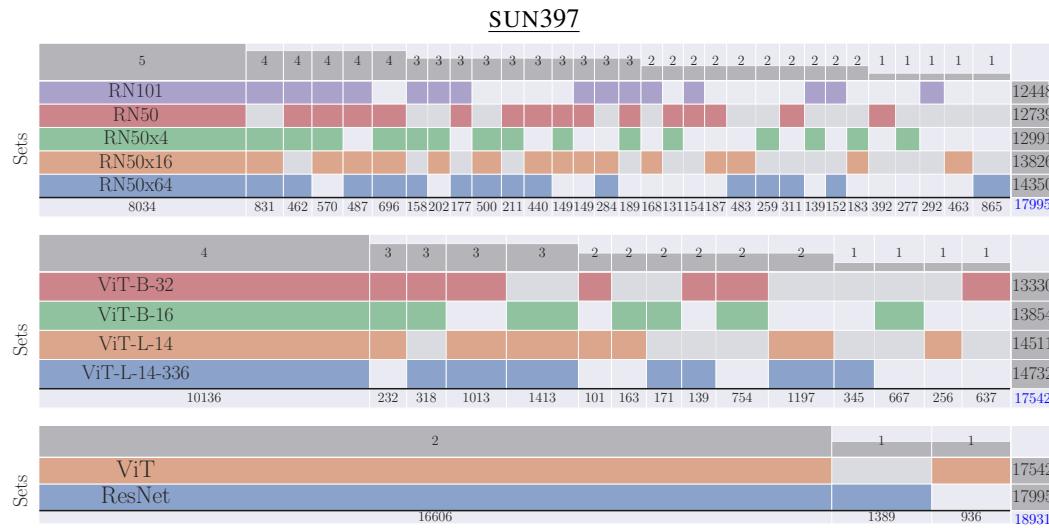


Figure L.23: SUN397 Overlap diagrams with the correct prediction of each backbone. The Top part of the Overlap diagram shows the number of backbones that are predicting correctly a set of images. Each column represents a set of image instances that are predicted correctly by some group of backbones. Each row in the diagram shows in colour the backbone that correctly predicts a certain set of image instances, in grey when the backbone is not correctly predicting those instances. The bottom part of the Overlap diagram shows the number of images in a certain set. The right part is the total amount of correctly predicted images per backbone.

2376

2377

2378

2379

2380

2381

2382

2383

2384

2385

2386

Table M.11: Our results on IMGNET-1K dataset for all the possible combinations of combining the zero-shot predictions of CLIP backbones, which we group into non-parametric and parametric techniques. Also, the best-performing single backbone (SINGLE-BEST) and the ORACLE performance. We present, for each combination of backbones, the improvement, constancy – and deterioration of accuracy performance for each method when we compare it against the SINGLE-BEST backbone. Mean, Max, and Min Δ summarize the difference in performance across methods and backbone combinations.

2393

IMGNET-1K															
ResNet	ViT				SINGLE-BEST	Non-Parametric				Parametric			ORACLE		
	50	101	B-32	B-16	L-14	VOTE T-1	VOTE T-3	CONF	LOG-AVG	C-CONF	C-LOG-AVG	GAC	NLC		
✓										59.84	-0.00				
✓	✓									62.28	-0.00				
✓		✓								63.35	-0.00				
✓		✓	✓							68.34	-0.00				
✓	✓					62.30	63.76 1.46	64.16 1.86	59.05 -3.25	64.61 2.31	63.22 0.91	64.47 2.17	64.71 2.41	65.14 2.84	70.75 8.45
✓	✓					63.36	65.00 1.65	65.37 2.01	64.40 1.04	65.78 2.42	64.39 1.03	65.54 2.19	65.86 2.51	66.07 2.71	71.76 8.41
✓		✓				68.34	68.47 0.12	68.72 0.37	68.15 -0.19	68.89 0.55	68.09 -0.26	68.85 0.51	69.60 1.26	69.86 1.52	74.51 6.17
✓		✓	✓			75.53	75.13 -0.4	75.27 -0.27	60.23 -15.3	75.09 -0.45	74.51 -1.02	74.97 -0.56	75.97 0.44	76.17 0.64	80.13 4.6
✓	✓					63.36	66.03 2.67	66.41 3.05	60.29 -3.07	66.99 3.63	65.45 2.09	66.70 3.35	67.00 3.65	67.35 3.99	72.96 9.61
✓	✓					68.34	69.09 0.75	69.29 0.95	68.54 0.2	69.58 1.24	68.45 0.11	69.40 1.05	69.90 1.56	70.29 1.95	75.15 6.8
✓		✓				75.53	75.18 -0.35	75.35 -0.18	62.53 -13.0	75.16 -0.38	74.56 -0.97	75.10 -0.44	76.01 0.48	76.39 0.86	80.28 4.75
✓	✓					68.34	68.90 0.56	69.17 0.82	63.17 -5.18	69.41 1.07	68.48 0.13	69.23 0.89	69.76 1.42	70.04 1.7	74.97 6.63
✓	✓					75.53	75.37 -0.16	75.40 -0.13	63.53 -12.01	75.19 -0.35	74.79 -0.74	75.20 -0.33	75.98 0.45	76.31 0.78	80.49 4.96
✓		✓				75.53	75.47 -0.06	75.61 0.08	68.42 -7.12	75.63 0.1	75.04 -0.5	75.50 -0.03	75.98 0.44	76.26 0.73	80.60 5.06
✓	✓	✓				63.36	66.15 2.79	66.76 3.4	60.27 -3.08	67.33 3.97	65.51 2.16	67.08 3.72	67.47 4.11	67.87 4.51	76.39 13.03
✓	✓	✓				68.34	68.22 -0.12	68.94 0.6	66.82 -1.53	69.23 0.89	68.17 -0.18	69.18 0.83	70.11 1.77	70.71 2.37	77.95 9.61
✓	✓	✓	✓			75.53	72.15 -3.38	74.09 -1.45	62.70 -12.83	74.22 -1.31	74.03 -1.5	74.08 -1.45	76.07 0.54	76.42 0.88	82.41 6.87
✓	✓	✓	✓			68.34	68.58 0.24	69.20 0.85	63.72 -4.62	69.45 1.11	68.35 -0.0	69.45 1.11	70.18 1.84	70.60 2.26	78.06 9.72
✓	✓	✓	✓			75.53	72.80 -2.73	74.35 -1.18	63.26 -12.28	74.36 -1.17	74.24 -1.3	74.32 -1.21	76.14 0.61	76.35 0.82	82.61 7.08
✓	✓	✓	✓			75.53	73.95 -1.58	74.93 -0.61	67.33 -8.21	74.98 -0.55	74.58 -0.95	74.90 -0.63	76.16 0.63	76.54 1.01	82.76 7.23
✓	✓	✓	✓			68.34	69.15 0.81	69.73 1.39	62.46 -5.88	70.02 1.68	68.66 0.32	69.95 1.61	70.43 2.09	70.87 2.53	78.58 10.24
✓	✓	✓	✓			75.53	73.32 -2.21	74.61 -0.92	63.05 -12.49	74.63 0.9	74.27 -1.27	74.58 -0.95	76.10 0.57	76.60 1.07	82.82 7.29
✓	✓	✓	✓			75.53	74.13 -1.4	75.04 -0.49	67.39 -8.14	75.11 -0.42	74.59 -0.94	75.14 -0.39	76.09 0.56	76.73 1.2	82.89 7.35
✓	✓	✓	✓			75.53	74.01 -1.52	74.90 -0.63	68.56 -6.97	75.02 -0.51	74.64 -0.9	74.98 -0.55	76.07 0.54	76.55 1.01	82.91 7.38
✓	✓	✓	✓			68.34	69.10 0.75	69.51 1.17	62.51 -5.83	69.66 1.32	68.43 0.08	69.74 1.4	70.52 2.18	71.09 2.75	80.31 11.97
✓	✓	✓	✓			75.53	72.37 -3.16	73.71 -1.82	63.11 -12.42	73.67 -1.86	73.94 -1.6	73.62 -1.92	76.12 0.58	76.69 1.16	84.08 8.55
✓	✓	✓	✓	✓		75.53	73.47 -2.06	74.46 -1.07	67.47 -8.06	74.46 -1.07	74.21 -1.32	74.41 -1.12	76.10 0.57	76.63 1.09	84.17 8.64
✓	✓	✓	✓	✓		75.53	73.54 -1.99	74.40 -1.13	67.54 -7.99	74.39 -1.14	74.32 -1.21	74.39 -1.14	76.14 0.61	76.65 1.11	84.26 8.73
✓	✓	✓	✓	✓		75.53	73.78 -1.75	74.64 -0.89	67.46 -8.07	74.58 -0.95	74.32 -1.21	74.59 -0.94	76.19 0.66	76.77 1.23	84.42 8.89
✓	✓	✓	✓	✓		75.53	72.67 -2.86	73.95 -1.58	67.46 -8.07	73.89 -1.64	74.06 -1.47	73.88 -1.65	76.22 0.69	76.59 1.06	85.29 9.76
Mean Δ -0.54 0.16 -7.09 0.29 -0.40 0.21 1.28 1.68 7.99 Max Δ 2.79 3.40 1.04 3.97 2.16 3.72 4.11 4.51 13.03 Min Δ -3.38 -1.82 -15.30 -1.86 -1.60 -1.92 0.44 0.64 4.60															

2410

2411

2412

2413

2414

2415

2416

2417

2418

2419

2420

2421

2422

2423

2424

2425

2426

2427

2428

2429