In Appendix A, we provide training details for each of our fine-tuning run. In Appendix B, we provide details on how to compute the concatenated feature vector $\Phi_\theta(\mathbf{x})$ from a ViT-B/32 CLIP model for LDIFS. In Appendix C, we provide additional results on **i)** observing how fine-tuning can lead to concept forgetting, **ii)** our investigations on why LP-init-L2SP is a better baseline compared to others in avoiding concept forgetting and **ii)** performance of our proposed LDIFS regularizer.

## A  TRAINING DETAILS

**Training datasets and hyper-parameters:** We fine-tune the CLIP ViT-B/32 model on each of the 9 image classification datasets: **a)** Stanford Cars, **b)** CIFAR-10, **c)** CIFAR-100, **d)** DTD, **e)** EuroSAT, **f)** GTSRB, **g)** MNIST, **h)** RESISC45, **i)** SVHN. For each dataset, we have a separate fine-tune run for each of the baselines discussed in §2. For each run, we train using the AdamW Loshchilov & Hutter (2017) optimizer, with an initial learning rate of $1e-5$, a weight decay of $0.1$ and a cosine learning rate scheduler with a warmup length of $500$. For all the runs, we use a batch-size of $128$. Following the code of Ilharco et al. (2022a), we use the following number of epochs to fine-tune each dataset: **a)** Stanford Cars: 35 epochs, **b)** CIFAR-10/100: 10 epochs, **c)** DTD: 76 epochs, **d)** EuroSAT: 12 epochs, **e)** GTSRB: 11 epochs, **f)** MNIST: 10 epochs, **g)** RESISC45: 15 epochs and **h)** SVHN: 10 epochs. For each dataset, we keep a minimum of 10 epochs for fine-tuning.

**Compute and number of experiments:** Each of our fine-tuning run is done on a single NVIDIA A100 GPU. As there are a total of 9 classification datasets and a 8 fine-tuning baselines (including our proposed LDIFS), we perform a total of 72 fine-tune training runs from a pre-trained CLIP ViT-B/32. Furthermore, in order to produce the plots capturing $\mathcal{A}_{\text{ZS}}$, $\mathcal{A}_{\text{LP}}$, as well as $\ell_2$ distance in parameter and feature space, we store intermediate checkpoints over the course of fine-tuning. Particularly, for each run, we evaluate checkpoints after every $20\%$ of training completion. Finally, in order to obtain the full set of results in Table 6, for the 72 fine-tuned models, we evaluate each of them on 9 datasets (i.e., including the test set of the fine-tuned task), thereby having a total of $648$ evaluations. Similar to training, each evaluation is also performed on a single NVIDIA A100 GPU.

**Choosing $\lambda_{\text{LDIFS}}$:** One hyper-parameter which the LDIFS regularizer introduces is $\lambda_{\text{LDIFS}}$. A higher value of $\lambda_{\text{LDIFS}}$ encourages the model to preserve features of the original foundation model and vice versa. For each classification task, we performed a grid search over $\lambda_{\text{LDIFS}} \in \{0.01, 0.05, 0.1, 0.5, 1, 10, 100\}$ and cross-validated this hyper-parameter on a held-out validation set, choosing the value which produces the best performance on the validation set. We found $\lambda_{\text{LDIFS}} = 10$ to produce the best performance over datasets in general, so all the results we present in this paper are with $\lambda_{\text{LDIFS}}$ set to 10.

## B  COMPUTING $\Phi_\theta(\mathbf{x})$

In this section, we discuss how we compute the concatenated feature vector $\Phi_\theta(\mathbf{x})$ given input $\mathbf{x}$ and model parameters $\theta$, specifically for a ViT-B/32 model. This is used for computing LDIFS, our proposed regularizer for fine-tuning.

Let the feature output from layer $l$ in the network for input $\mathbf{x}$ be $\Phi_{\theta(l)}(\mathbf{x}) \in \mathbb{R}^l$. Then the normalized feature vector for layer $l$ can be represented as $\frac{\Phi_{\theta(l)}(\mathbf{x})}{||\Phi_{\theta(l)}(\mathbf{x})||}$. In order to form the concatenated feature vector, one can take technically take features from every intermediate layer of the network. However, storing all the features is memory intensive. Thus, we follow a similar approach to the LPIPS Zhang et al. (2018) implementation for VGG and AlexNet, and choose 5 intermediate points in the ViT-B/32 architecture to collect features from. For a single input image $\mathbf{x} \in \mathbb{R}^{3 \times 224 \times 224}$, each of the intermediate feature representations has a dimension of $\mathbb{R}^{50 \times 768}$ with 50 tokens and 768 dimensional representation for each token. When normalizing the feature vector, we flatten this vector out to a single 38400 dimensional vector. Thus the full concatenated feature vector $\Phi_\theta(\mathbf{x})$ has dimensionality $5 \times 38400 = 192000$. However, note that there can be other ablations to this design and we leave that for future exploration.

**LDIFS vs LPIPS Zhang et al. (2018)**: One can find similarities between our proposed LDIFS regulariser and the LPIPS metric Zhang et al. (2018) used for measuring perceptual similarity between images. However, while LPIPS uses feature space distance on a *pre-trained, frozen* model to find perceptual similarity between pairs or sets of images, LDIFS instead feature space distance

between a pair of pre-trained and fine-tuned model for the same image, to preserve the input-output behaviour of the pre-trained model.

## C    ADDITIONAL RESULTS

In this section, we present additional empirical results to supplement our observations and conclusions in the main paper.

### C.1    CRIPPLING EFFECT OF FINE-TUNING

In §3, we observed how most existing fine-tuning methods, while gaining state-of-the-art performance on fine-tuned tasks, can lead to concept forgetting in models. In this section, we provide additional empirical results to further strengthen those observations.

In Figure 7 and Figure 15, we present the $\mathcal{A}_{\text{ZS}}$ and $\mathcal{A}_{\text{LP}}$ accuracy for models fine-tuned on each of the 9 classification tasks, on their respective test sets. This shows expected behaviour as models broadly achieve very high test set accuracy, which increases over the course of fine-tuning, on their respective fine-tuned tasks. Note that LP-init baselines seem to have relatively lower performance on $\mathcal{A}_{\text{ZS}}$ accuracy compared to their ZS-init counter-parts. However, this reduced performance is only limited to $\mathcal{A}_{\text{ZS}}$ as this does not translate to a reduced performance in $\mathcal{A}_{\text{LP}}$. Moreover, L2SP baselines (both ZS and LP-init) obtain relatively lower fine-tuned accuracy both in case of $\mathcal{A}_{\text{ZS}}$ and $\mathcal{A}_{\text{LP}}$.

In Figure 8 and Figure 16, we present the $\mathcal{A}_{\text{ZS}}$ and $\mathcal{A}_{\text{LP}}$ for models fine-tuned on EuroSAT (first row), GTSRB (second row) and SVHN (third row) as captured on 8 classification datasets different from their respective fine-tuned set. Broadly, the ZS and LP performance for all models drops on other datasets as they are fine-tuned, thereby capturing concept forgetting in the fine-tuned models. Among the baselines, we consistently observe LP-init-L2SP to perform better than others in avoiding concept forgetting. This is evident through the distinctly higher $\mathcal{A}_{\text{ZS}}$ and $\mathcal{A}_{\text{LP}}$ accuracies over fine-tuning for the LP-init-L2SP baselines.

### C.2    ANALYSING L2SP

In this section, we provide additional results to supplement our observations related to investigating the L2SP baseline in §4 of the main paper. In Figure 19, we present the $\ell_2$ distance in parameter space $||\theta - \theta_0||_2^2$ between the pre-trained and current models captured over fine-tuning. For all the datasets, we observe that L2SP baselines while initially diverging slightly in the parameter space, converge back to a model having low $\ell_2$ parameter space distance from the original foundation model. Other baselines on the other hand, completely diverge away from the original model in the parameter space.

To further investigate the change in input-output behaviour of the model over fine-tuning, we measure the distance in feature space (see Equation (3)) over fine-tuning. In Figure 20, we present the $\ell_2$ distance in feature space captured for models fine-tuned on EuroSAT. Again, consistent with our previous observation, we find that unlike other baselines, L2SP first diverges away from the original foundation model and then converges back to the original input-output behaviour, as is indicative through a decreasing L2 feature space distance in the later stages of fine-tuning. This observation is consistent on the EuroSAT train set, the EuroSAT test set as well as on other datasets, thereby providing our motivation for the LDIFS regularizer.

### C.3    PERFORMANCE OF LDIFS

In this section, we provide additional results to supplement observations related to the performance of the proposed LDIFS regularizer.

**Analysing LDIFS on parameter space and feature space distance:** In Figure 9 and Figure 10, we plot the $\ell_2$ distance in parameter/weight and feature space respectively (same as Figure 3 and Figure 4 in the main paper), but with the LDIFS baselines added in. In the weight space, LDIFS, while lower than other fine-tuning methods, has a relatively higher $\ell_2$ distance from the pre-trained model compared to L2SP. On the contrary, in the feature space, LDIFS consistently gets the lowest distance from the pre-trained model. This is expected and was indeed the purpose behind LDIFS's
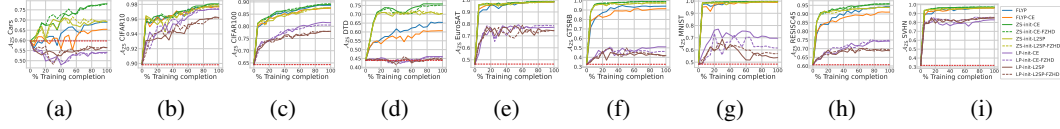
Figure 7: **Test set ZS Accuracy $\mathcal{A}_{ZS}$ for different fine-tuning methods** on 9 image classification tasks: (a) Stanford Cars, (b) CIFAR-10, (c) CIFAR-100, (d) DTD, (e) EuroSAT, (f) GTSRB, (g) MNIST (h) RESISC45 and (i) SVHN. $\mathcal{A}_{ZS}$ generally rises over the course of fine-tuning. However, for the ZS-init-L2SP and LP-init-L2SP baselines, the gain in $\mathcal{A}_{ZS}$ is relatively lower. Furthermore, for LP-init baselines, the performance is consistently lower compared to other baselines



Figure 8: **ZS Accuracy $\mathcal{A}_{ZS}$ for models fine-tuned on EuroSAT (first row), GTSRB (second row), and SVHN (third row) and evaluated on 8 different datasets different from their fine-tuning dataset.** Most fine-tuning methods show a drop in $\mathcal{A}_{ZS}$ performance over the course of fine-tuning indicating a reduction in the model's transferability.



Figure 9: $\ell_2$ **distance in weight space** $||\theta_{v(t)} - \theta_{v(0)}||_2^2$ between the image encoder fine-tuned to the current time-step $f_{\theta_{v(t)}}$ and the pre-trained image encoder $f_{\theta_{v(0)}}$ over the course of fine-tuning.



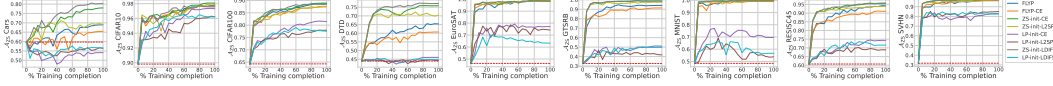Figure 10: $\ell_2$ **distance in feature space** $d(\theta_{v(t)}, \theta_{v(0)}, \mathcal{D})$, between image encoders, $f_{\theta_{v(t)}}$ and $f_{\theta_{v(0)}}$ computed over different fine-tuning methods for models fine-tuned on EuroSAT.

design which is to minimize the difference in input-output behaviour, captured through feature space distance between the fine-tuned and pre-trained models and not necessarily constrain the parameter space of the fine-tuned model to lie close to the pre-trained model. Furthermore, note from Figure 10, that although we are applying the LDIFS regularizer only on the EuroSAT samples during fine-tuning, the preservation of features extends even to other datasets like CIFAR10, CIFAR100 and GTSRB, as is indicated by the low feature space distance even on these datasets.

**Results and observations:** Firstly, in Figure 13 and Figure 17, we present the test set accuracy of all fine-tuning methods including LDIFS on all 9 classification tasks. We see that LDIFS performs competitively with all other baselines and improves on L2SP consistently. In Figure 14 and Figure 18, we present the $\mathcal{A}_{ZS}$ and $\mathcal{A}_{LP}$ accuracy respectively, for models fine-tuned on EuroSAT (first row),

Figure 11: $\mathcal{A}_{\text{LP}}$ **for models trained with LDIFS on EuroSAT using the full concatenated feature vector vs just the last layer (LL) feature vectors.** Full feature vectors with earlier features is crucial for LDIFS's performance.

GTSRB (second row) and SVHN (third row) over the course of fine-tuning. Furthermore, we also provide the $\Delta_{\text{ZS}}$, $\Delta_{\text{LP}}$ and the $\mathcal{A}_{\text{LP}}$ values for fully fine-tuned models on all 9 classification tasks and all baselines in Table 6.

As is crystal clear from our results, LP-init-LDIFS consistently minimizes concept forgetting without compromising on performance on the fine-tuned task. This is evident from its high $\mathcal{A}_{\text{ZS}}$, $\mathcal{A}_{\text{LP}}$ (see Figure 14 and Figure 18) and high $\Delta_{\text{ZS}}$ and $\Delta_{\text{LP}}$ (see Table 6) when evaluated on tasks other than the fine-tuning task at hand. Note from Table 6 that even when LP-init-LDIFS is not achieving the best performance on a certain dataset pair, it is often a close second with very little difference compared to the best performing baseline. In addition to this advantage, its $\mathcal{A}_{\text{LP}}$ accuracy on the fine-tuned task itself is very competitive with the top scores obtained by other fine-tuning baselines and provides a consistent improvement on L2SP.

Finally, in Table 7, we present the full results of the 2 task sequence setup for continual fine-tuning. Even in this setting, we find LP-initLDIFS to not only preserve performance of the first fine-tuning task during the second fine-tuning stage but also to preserve performance of all other tasks during the entire sequence of end-to-end fine-tuning.

**Ablation with last layer features:** In Figure 11, we present an ablation where we fine-tune on EuroSAT using LDIFS but using just the last layer features as opposed to the full concatenated feature vector. Clearly, including earlier features when computing the LDIFS regularizer is crucial to its performance as preserving just the last layer features leads to a significant performance drop compared to using the full feature vector.

## C.4 COMPARISON WITH ADAPTERS AND PROMPT TUNING

In Table 5, we present results of fine-tuning a CLIP ViT-B/32 model using linear probing, CoOp, CLIP-Adapter, Tip-Adapter and full end-to-end fine-tuning using ZS-init-CE loss on 9 downstream image classification tasks. We consistently observe end-to-end fine-tuning to produce the best test set accuracy across all 9 tasks. Furthermore, we also note that in case of tasks where the gap between linear probe performance and end-to-end fine-tuning performance is significant (e.g. SVHN), adapter and prompt tuning approaches don't cover this performance gap. Hence, if the pre-trained encoder is not performant on a downstream task, the only way to achieve state-of-the-art results on the task requires end-to-end fine-tuning the model on the task itself. This observation thus necessitates study into better end-to-end fine-tuning methods for foundation models.

## C.5 COMPARISON WITH WISE-FT

In this section, we perform an ablation to see the effect of Wise-FT Wortsman et al. (2022b) on preventing concept forgetting. Wise-FT is a weight interpolation method which forms a linear combination between the pre-trained $\theta_0$ and fine-tuned $\theta_f$ model parameters:

$$\theta_{\text{wse}} = \alpha\theta_0 + (1 - \alpha)\theta_f \tag{6}$$

Since this can work with any fine-tuned model, we can apply it on all the end-to-end fine-tuning baselines in this work. Specifically, we tune the hyperparameter $\alpha$ on a held-out validation set of the downstream task at hand to maximum validation accuracy. We evaluate concept forgetting on 3 downstream tasks: CIFAR-10, EuroSAT and SVHN using the $\Delta_{\text{LP}}$ metric and report the results in Table 4. For each dataset, we evaluate the mean $\Delta_{\text{LP}}$ on 5 other downstream tasks: Cars, DTD, GTSRB, RESISC45 and MNIST. Our observations show a consistent reduction in concept forgetting across all fine-tuning baselines. However, the order of performance between the fine-tuning baselines does not change and LDIFS still maintains superior performance over other baselines both pre and

Figure 12: **Visualizing test LP accuracy of fine-tuning baselines on different datasets.** For each dataset, we plot its test set LP accuracy on the x-axis and the average test set LP accuracy of the remaining datasets on the y-axis.

| Dataset | Fine-tuning baselines | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | FLYP | | FLYP-CE | | ZS-init-CE | | LP-init-CE (LP-FT) | | ZS-init-L2SP | | LP-init-L2SP | | ZS-init-LDIFS (Ours) | | LP-init-LDIFS (Ours) | |
| | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ | $\Delta_{LP}$ | $\Delta_{LP}(+wse)$ |
| CIFAR10 | −5.14 | −0.11 | −1.33 | −0.16 | −1.67 | −0.1 | −1.49 | −0.01 | 1.58 | 1.64 | 1.49 | 1.96 | 1.53 | 1.99 | **1.72** | **2.01** |
| EuroSAT | −6.5 | −1.32 | −5.53 | −1.04 | −5.62 | −1.23 | −4.08 | −0.76 | −1.77 | −0.12 | −0.87 | 0.08 | 0.83 | 1.2 | **1.24** | **1.96** |
| SVHN | −10.82 | −6.24 | −10.18 | −6.37 | −10.98 | −7.2 | −9.2 | −4.11 | −2.65 | −0.66 | −2.13 | −0.23 | −0.57 | 0.12 | **−0.18** | **0.33** |

Table 4: $\Delta_{LP}$ without and with Wise-FT Wortsman et al. (2022b) for all fine-tuning baselines on CIFAR-10, EuroSAT and SVHN.



Figure 13: **Comparing test set ZS Accuracy $\mathcal{A}_{ZS}$ for different fine-tuning methods with LDIFS** on 9 image classification tasks. While LP-init baselines generally underperform on $\mathcal{A}_{ZS}$, we find ZS-init-LDIFS to be competitive with the best baselines.

post application of Wise-FT. Thus, the results show that Wise-FT can be combined with any E2E fine-tuning method to further minimize concept forgetting.

## C.6 VISUALIZING TABLE 1

In this section, we present the results of Table 1 using an alternate visualization. For each of the 9 downstream datasets, we plot the LP accuracy on the dataset itself on the x-axis and the average LP accuracy on all other datasets on the y-axis. The best baselines should lie on the top right corner of this plot as that indicates that as the model is fine-tuned on the dataset, it still retains performance on other datasets. Results are in Figure 12. Clearly, we consistently see the LDIFS baselines to lie at the top right, again indicative of their performance in not just minimizing concept forgetting but also obtaining excellent downstream task performance.

## C.7 GENERALITY OF CONCEPT FORGETTING BEYOND CLIP ViT-B/32

To study the relevance of concept forgetting beyond a CLIP ViT model, in Table 8 and Table 9, we apply the end-to-end fine-tuning methods to fine-tune a CLIP RN-50 model and a FLAVA Singh et al. (2022) ViT-B/16 model respectively on EuroSAT. From both these tables, we observe both concept forgetting to exist and LP-init-LDIFS to minimize it, thereby validating our assumptions on models beyond a CLIP ViT model.

Figure 14: **Comparing ZS Accuracy $\mathcal{A}_{\mathrm{ZS}}$ of different fine-tuning methods with LDIFS for models fine-tuned on EuroSAT (first row), GTSRB (second row) and SVHN (third row) and evaluated on 8 datasets different from their fine-tuning dataset.** LP-init-LDIFS almost consistently beats all other baselines including LP-init-L2SP in preserving the model's transferability.



Figure 15: **Test set LP Accuracy $\mathcal{A}_{\mathrm{LP}}$ for different fine-tuning methods** on 9 image classification tasks: (a) Stanford Cars, (b) CIFAR-10, (c) CIFAR-100, (d) DTD, (e) EuroSAT, (f) GTSRB, (g) MNIST (h) RESISC45 and (i) SVHN. $\mathcal{A}_{\mathrm{LP}}$ generally rises over the course of fine-tuning. However, for the ZS-init-L2SP and LP-init-L2SP baselines, the gain in $\mathcal{A}_{\mathrm{LP}}$ is relatively lower. Furthermore, for LP-init baselines, the performance is consistently lower compared to other baselines



Figure 16: **LP Accuracy $\mathcal{A}_{\mathrm{LP}}$ for models fine-tuned on EuroSAT (first row), GTSRB (second row), and SVHN (third row) and evaluated on 8 different datasets different from their fine-tuning dataset.** Most fine-tuning methods show a drop in $\mathcal{A}_{\mathrm{LP}}$ performance over the course of fine-tuning indicating a reduction in the model's transferability.



Figure 17: **Comparing test set LP Accuracy $\mathcal{A}_{\mathrm{LP}}$ for different fine-tuning methods with LDIFS** on 9 image classification tasks.

Figure 18: **Comparing LP Accuracy $\mathcal{A}_{LP}$ of different fine-tuning methods with LDIFS for models fine-tuned on EuroSAT (first row), GTSRB (second row) and SVHN (third row) and evaluated on 8 datasets different from their fine-tuning dataset.** LP-init-LDIFS almost consistently beats all other baselines including LP-init-L2SP in preserving the model's transferability.



Figure 19: $\ell_2$ **distance in weight space** $||\theta - \theta_0||_2^2$ between pre-trained image encoder $f_{\theta_0}$ and fine-tuned image encoder $f_\theta$ over the course of fine-tuning. Apart from ZS-init-L2SP and LP-init-L2SP, all fine-tuning baselines diverge in weight space over the course of fine-tuning.



Figure 20: $\ell_2$ **distance in feature space** $d(\theta, \theta_0, \mathcal{D})$, between fine-tuned and pre-trained image encoders, $f_\theta$ and $f_{\theta_0}$ computed over different fine-tuning methods for models fine-tuned on EuroSAT.

| Dataset | LP | CoOp | CLIP-Adapter | Tip-Adapter | E2E Fine-tuning |
|---|---|---|---|---|---|
| Cars | 80.80 | 80.88 | 81.24 | 81.32 | **83.48** |
| CIFAR10 | 94.92 | 95.02 | 94.87 | 95.13 | **97.73** |
| CIFAR100 | 79.62 | 80.12 | 80.06 | 80.82 | **88.6** |
| DTD | 72.08 | 72.01 | 71.87 | 72.62 | **77.18** |
| EuroSAT | 95.56 | 95.14 | 95.38 | 96.23 | **98.76** |
| GTSRB | 86.70 | 86.81 | 87.45 | 88.04 | **98.52** |
| MNIST | 98.65 | 98.84 | 99.03 | 99.11 | **99.67** |
| RESISC45 | 91.86 | 91.79 | 91.82 | 91.80 | **95.76** |
| SVHN | 65.47 | 67.28 | 69.76 | 69.23 | **97.3** |

Table 5: Test set accuracy obtained from linear probing (LP), CoOp, CLIP-Adapter, Tip-Adapter and end-to-end (E2E) fine-tuning using ZS-init-CE loss.

Table 6: $\Delta_{ZS}$, $\Delta_{LP}$ and $\mathcal{A}_{LP}$ **for models fully fine-tuned on 9 different classification datasets.** LP-init-LDIFS outperforms other fine-tuning baselines and even LP-init-L2SP in minimizing concept forgetting, while also outperforming LP-init-L2SP on the fine-tuned task performance.

Table 7: $\Delta_{ZS}$, $\Delta_{LP}$ and $\mathcal{A}_{LP}$ **for models fine-tuned on SVHN → CIFAR10 and SVHN → CIFAR100.** LP-init-LDIFS consistently outperforms other baselines even in the sequence setup.

Table 8: $\Delta_{ZS}$, $\Delta_{LP}$ and $\mathcal{A}_{LP}$ for CLIP RN50 fully fine-tuned on EuroSAT.

Table 9: $\Delta_{LP}$ and $\mathcal{A}_{LP}$ for FLAVA ViT-B/16 fully fine-tuned on EuroSAT.