702 703	А	Appendix						
704	A.1	Reproducibility Statement						
705	4 11							
706	All C	All our code and model weights will be open-sourced post the conference anonymity period and will be available to the larger community to use under open source licensing.						
708	be a	anable to the target community to use under open source neensing.						
709	A.2	Equal Opportunity in Fairness						
710	F							
711	Equa	al Opportunity ensures that the True Positive Rates (TPR) are equal across demographic ns. Mathematically:						
712	grou	$TPR_{Group 1} = TPR_{Group 2} = \dots$						
714	when	e:						
715		$TPR = \frac{True Positives (TP)}{TPR}$						
716		True Positives (TP) + False Negatives (FN)						
717	Εv	μαι ε. Ι σανι Δααροναι						
718	LAA	MILL. LOAN MITROVAL						
720	Cons	sider a loan model evaluating two demographic groups with the following outcomes:						
721								
722		Group True Positives (TP) False Negatives (FN) TPR $\frac{1}{2}$						
723		$\begin{array}{c ccccccccccccccccccccccccccccccccccc$						
724								
726	The	model violates Equal Opportunity because $\text{TPR}_A \neq \text{TPR}_B$.						
727	Env.							
728	FIX:	THRESHOLD ADJUSTMENT						
729 730	Adju	sting decision thresholds can equalize the TPR:						
731		• Group A: Keep the threshold at 0.5						
732		• Group B: Lower the threshold to 0.4						
733		Group D. Lower the threshold to 0.1.						
734	Afte	r adjustment, the outcomes are:						
736								
737		GroupIrue Positives (IP)Faise Negatives (FN)IPRA80200.80						
738		$\begin{array}{c ccccccccccccccccccccccccccccccccccc$						
739								
740	Now	$TPR_A = TPR_B = 0.80$, satisfying Equal Opportunity.						
742	FAIR	NESS CONSTRAINT						
743	1 / 111							
744	Duri	During training, Equal Opportunity can be enforced as:						
745		$ \mathrm{TPR}_{\mathrm{Group} \mathrm{A}} - \mathrm{TPR}_{\mathrm{Group} \mathrm{B}} \leq \epsilon$						
747	or by	or by adding a penalty to the loss function:						
748	-	$c = c + \lambda$ TDD TDD						
749		$\mathcal{L}_{\text{fair}} = \mathcal{L}_{\text{original}} + \lambda \cdot \mathbf{IPK}_{\text{Group A}} - \mathbf{IPK}_{\text{Group B}} $						
750 751	when	where λ controls the trade-off between fairness and accuracy.						
752	A.3	A.3 FRÉCHET INCEPTION DISTANCE (FID) SCORES - DISTRIBUTION SHIFT						
753	T							
754	The	Frechet Inception Distance (FID) is a widely-used metric to evaluate the quality of generated						
(00	distr	ibutions of two datasets (real and generated images) using the activations of the InceptionV3						

model. Specifically, the FID score is computed using the mean and covariance of these activations, assuming a multivariate Gaussian distribution.

In our setup, we employed the Clean FID library Parmar et al. (2022) to calculate the FID between various datasets. For large datasets like LAION-5B, we sampled random subsets and computed the FID on these smaller samples to manage computational constraints. For smaller datasets such as Waterbirds, GeoDe, and Aircrafts, we used the entire dataset for the calculation.

Fine-Tuning Data	ImageNet	LAION
Aircrafts	181.98	229.72
Waterbirds	117.61	142.46
GeoDE	54.38	56.91

767 768

769 770 771

772 773

774

Table 4: FID comparison between the fine-tuning and pre-training data.

A.4 TRAINABLE PARAMETER RATIOS WITH LORA/FAIRLORA

The number of trainable parameter remain the same for LoRA as well as FairLoRA for the same rank. All models have similar ratios for % of Trainable parameters.

Model	Rank	Trainable Params	Total Params	(% of Trainable)
	8	325,672	86,155,088	0.38
	16	666,724	86,155,088	0.77
DiNO	32	1,256,548	86,155,088	1.44
	64	2,436,196	86,155,088	2.76
	128	4,795,492	86,155,088	5.29

Table 5: Summary of model parameters for DiNO. The table includes the rank, number of trainable parameters, total parameters, and the percentage of trainable parameters.

788

792 793 794

796

797

798

783

A.5 GRADIENT UPDATES IN LORA

The gradient updates in LoRA apply only to the low-rank matrices A and B, while the pre-trained parameters θ_0 remain fixed. Given an objective function \mathcal{L} , the gradients of the loss with respect to A and B are computed as:

$\partial \mathcal{L}$	$\partial \mathcal{L}_{P^{\top}}$	$\partial \mathcal{L}$ _	$_{A^{\top}}\partial \mathcal{L}$
∂A	$= \overline{\partial \theta}^{D}$,	$\overline{\partial B} =$	$A \overline{\partial \theta}$

Here, $\frac{\partial \mathcal{L}}{\partial \theta}$ is the gradient of the loss with respect to the full parameter matrix θ . These updates allow the model to adapt to the downstream task with far fewer trainable parameters, preserving most of the pre-trained knowledge while fine-tuning for the new task.

LoRA's parameterization is particularly effective in large models where the parameter matrices are high-dimensional, as it avoids the computational cost of updating the entire matrix. By focusing on the low-rank updates, LoRA achieves a balance between fine-tuning flexibility and resource efficiency.

For further details, the original LoRA formulation and its theoretical justification can be found in Hu et al. (2021).

806

808

807 A.6 GRADIENTS IN FAIRLORA

Let $\theta = \theta_0 + AB$, where $\Delta \theta = AB$ is the LoRA low-rank update. We need to compute the gradients of $\mathcal{J}(\theta)$ with respect to both A and B.

GRADIENT WITH RESPECT TO *A*:

Gradient with respect to B:

$$\frac{\partial \mathcal{J}}{\partial B} = A^{\top} \frac{\partial \mathcal{L}}{\partial \theta} + \lambda \sum_{g \in \mathcal{G}} 2 \left(\mathcal{L}_g(\theta) - \frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \mathcal{L}_{g'}(\theta) \right) A^{\top} \frac{\partial \mathcal{L}_g(\theta)}{\partial \theta}$$

 $\frac{\partial \mathcal{J}}{\partial A} = \frac{\partial \mathcal{L}}{\partial \theta} B^{\top} + \lambda \sum_{g \in \mathcal{G}} 2 \left(\mathcal{L}_g(\theta) - \frac{1}{|\mathcal{G}|} \sum_{g' \in \mathcal{G}} \mathcal{L}_{g'}(\theta) \right) \frac{\partial \mathcal{L}_g(\theta)}{\partial \theta} B^{\top}$

822 A.7 EMPIRICAL EVALUATIONS CONTINUED



(a) Model: DiNO. We notice a dominant pattern for
FairLoRA across metrics apart from EOD, where it is
comparable to LoRA

(b) Model: DiNO. We notice a dominant pattern for FairLoRA across metrics.

Figure 6: All metrics are normalized to the same scale and adjusted such that higher is better. Comparison of FairLoRA performance on DiNO model across datasets.





918						
919	M_1_1		A (A)			
920	Nodel	Method	Accuracy (†)	F1 Min (†)	Recall Min (†)	Δ F1 (\downarrow)
021		LoRA	81.98 ± 1.38	17.60 ± 5.05	14.97 ± 5.10	80.92 ± 5.00
921	CLiP	FairLoRA	86.67 ± 2.56	27.44 ± 15.30	24.81 ± 14.07	71.06 ± 13.78
922		FFT	81.92 ± 1.21	21.28 ± 5.31	18.78 ± 3.14	77.21 ± 5.29
923		FairFFT	87.03 ± 2.81	22.82 ± 12.47	17.94 ± 10.30	76.17 ± 10.72
924		LoRA	69.39 ± 0.87	21.43 ± 2.81	19.99 ± 4.57	77.59 ± 1.97
925	D:NO	FairLoRA	68.28 ± 0.32	20.64 ± 4.67	18.75 ± 5.92	77.88 ± 4.65
926	DINO	FFT	71.26 ± 0.06	20.37 ± 7.65	16.67 ± 6.79	79.15 ± 6.81
927		FairFFT	70.95 ± 0.95	23.46 ± 0.11	20.59 ± 2.94	75.10 ± 1.32
928		LoRA	70.06 ± 0.92	25.99 ± 3.11	24.24 ± 3.03	72.03 ± 2.28
929	ViT	FairLoRA	70.45 ± 0.48	27.11 ± 2.11	25.49 ± 3.40	72.41 ± 1.27
930		FFT	74.17 ± 0.31	32.18 ± 5.83	27.75 ± 6.28	66.34 ± 5.80
931		FairFFT	74.18 ± 0.64	35.36 ± 5.04	29.71 ± 2.99	63.15 ± 4.99

Table 6: The table compares FFT vs LoRA and FairFFT vs FairLoRA for Aircrafts. Metrics include: Accuracy, the mean classification accuracy; F1 Min, the minimum F1 score across classes; Recall Min, the minimum Recall across classes; Δ F1, the difference between the maximum and minimum F1 score across classes.



Figure 9: Comparison on the impact of rank on performance as well as fairness across models for Aircrafts. Higher value is better in both the graphs.



Figure 10: Comparison on the impact of rank on performance as well as fairness across models for
Waterbirds. Higher value is better in both the graphs.



Figure 11: On the left, we plot the performance of the model with and without fairness regularizers for both full finetuning (FFT) as well as LoRA. We notice that using the regularizer improves overall performance as well in this particular example. On the right, we visualize the effect on the variance of loss across classes and notice that this variance is lower for the ones with regularizer and the combination of LoRA and the regularizer yeilds the best results.



Figure 12: On the left, we plot the performance of the model with and without fairness regularizers for both full finetuning (FFT) as well as LoRA. On the right, we visualize the effect on the variance of loss across classes.