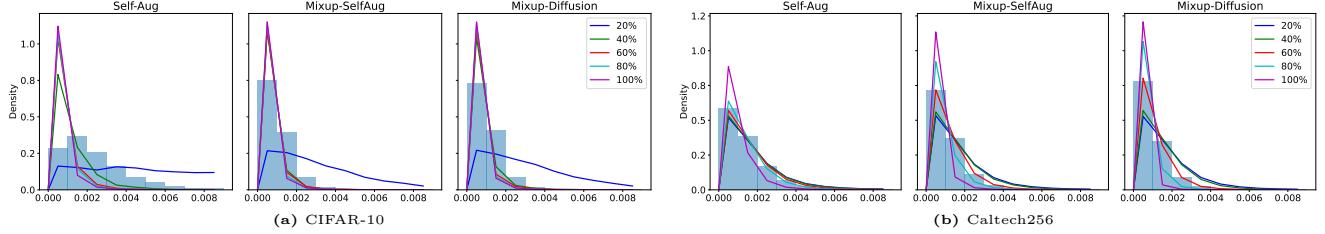


**Figure 1:** Norms of gradients before clipping and adding noise on CIFAR-10(a) and Caltech256(b) with fining tuning Vit-B-16 model with  $\varepsilon = 1$  and  $\delta = 10^{-5}$ . The different curves represent different training stages. The histogram shows average norm of gradients. Gradients with our proposed MIXUP-SELAUG and Mixup-Diffusion are more concentrated suggesting more stable training and faster convergence.



**Table 1:** Fine-tune Vit-B-16 model and ConvNext models on Caltech256, SUN397 and Oxford-IIIT Pet datasets with different  $\varepsilon$ . We set  $\delta = 10^{-5}$ . We can observe that our proposed methods improves performance in all cases.

Model	Dataset	Method	$\varepsilon=1$	$\varepsilon=2$	$\varepsilon=4$	$\varepsilon=8$
Caltech256	Vit-B-16	Mixed Ghost clipping	79.74% $(\pm 0.15\%)$	88.21% $(\pm 0.21\%)$	91.42% $(\pm 0.24\%)$	92.27% $(\pm 0.16\%)$
		Self-Aug	80.36% $(\pm 0.11\%)$	89.67% $(\pm 0.16\%)$	92.01% $(\pm 0.08\%)$	93.17% $(\pm 0.15\%)$
		Mixup-SelfAug	<b>81.21%</b> $(\pm 0.15\%)$	<b>90.12%</b> $(\pm 0.17\%)$	<b>92.17%</b> $(\pm 0.21\%)$	<b>93.39%</b> $(\pm 0.08\%)$
SUN397	Vit-B-16	Mixup-Diffusion	85.82% $(\pm 0.12\%)$	91.32% $(\pm 0.11\%)$	92.86% $(\pm 0.14\%)$	93.87% $(\pm 0.10\%)$
		Mixed Ghost clipping	70.67% $(\pm 0.17\%)$	71.21% $(\pm 0.13\%)$	72.24% $(\pm 0.15\%)$	72.51% $(\pm 0.19\%)$
		Self-Aug	72.65% $(\pm 0.09\%)$	76.02% $(\pm 0.14\%)$	78.05% $(\pm 0.11\%)$	79.54% $(\pm 0.15\%)$
Oxford-IIIT Pet	Vit-B-16	Mixup-SelfAug	<b>73.19%</b> $(\pm 0.13\%)$	<b>76.45%</b> $(\pm 0.17\%)$	<b>78.67%</b> $(\pm 0.16\%)$	<b>79.57%</b> $(\pm 0.14\%)$
		Mixup-Diffusion	75.12% $(\pm 0.17\%)$	77.78% $(\pm 0.12\%)$	79.47% $(\pm 0.18\%)$	80.57% $(\pm 0.09\%)$
		Mixed Ghost clipping	71.21% $(\pm 0.18\%)$	79.12% $(\pm 0.15\%)$	80.37% $(\pm 0.22\%)$	81.02% $(\pm 0.19\%)$
Caltech256	ConvNext	Self-Aug	72.21% $(\pm 0.21\%)$	82.11% $(\pm 0.19\%)$	85.84% $(\pm 0.25\%)$	88.23% $(\pm 0.13\%)$
		Mixup-SelfAug	<b>72.45%</b> $(\pm 0.24\%)$	<b>82.51%</b> $(\pm 0.21\%)$	<b>86.75%</b> $(\pm 0.17\%)$	<b>88.70%</b> $(\pm 0.15\%)$
		Mixup-Diffusion	73.67% $(\pm 0.25\%)$	84.57% $(\pm 0.18\%)$	88.14% $(\pm 0.21\%)$	89.25% $(\pm 0.17\%)$
SUN397	ConvNext	Mixed Ghost clipping	79.21% $(\pm 0.18\%)$	87.37% $(\pm 0.19\%)$	87.96% $(\pm 0.14\%)$	88.17% $(\pm 0.17\%)$
		Self-Aug	79.98% $(\pm 0.15\%)$	88.15% $(\pm 0.11\%)$	91.02% $(\pm 0.14\%)$	92.17% $(\pm 0.10\%)$
		Mixup-SelfAug	<b>81.07%</b> $(\pm 0.11\%)$	<b>88.69%</b> $(\pm 0.12\%)$	<b>91.34%</b> $(\pm 0.08\%)$	<b>92.38%</b> $(\pm 0.09\%)$
Oxford-IIIT Pet	ConvNext	Mixup-Diffusion	85.08% $(\pm 0.14\%)$	90.74% $(\pm 0.16\%)$	92.27% $(\pm 0.11\%)$	93.27% $(\pm 0.06\%)$
		Mixed Ghost clipping	64.27% $(\pm 0.22\%)$	65.05% $(\pm 0.15\%)$	65.25% $(\pm 0.12\%)$	65.33% $(\pm 0.14\%)$
		Self-Aug	72.21% $(\pm 0.14\%)$	76.45% $(\pm 0.10\%)$	77.95% $(\pm 0.16\%)$	78.93% $(\pm 0.11\%)$
SUN397	ConvNext	Mixup-SelfAug	<b>72.47%</b> $(\pm 0.07\%)$	<b>76.82%</b> $(\pm 0.11\%)$	<b>78.52%</b> $(\pm 0.04\%)$	<b>79.45%</b> $(\pm 0.08\%)$
		Mixup-Diffusion	74.95% $(\pm 0.16\%)$	77.51% $(\pm 0.13\%)$	79.33% $(\pm 0.10\%)$	79.98% $(\pm 0.12\%)$
		Mixed Ghost clipping	65.21% $(\pm 0.23\%)$	78.19% $(\pm 0.18\%)$	79.12% $(\pm 0.24\%)$	79.85% $(\pm 0.21\%)$
Oxford-IIIT Pet	ConvNext	Self-Aug	68.11% $(\pm 0.18\%)$	81.29% $(\pm 0.21\%)$	85.47% $(\pm 0.14\%)$	86.96% $(\pm 0.11\%)$
		Mixup-SelfAug	<b>68.75%</b> $(\pm 0.24\%)$	<b>81.65%</b> $(\pm 0.22\%)$	<b>86.25%</b> $(\pm 0.17\%)$	<b>87.67%</b> $(\pm 0.16\%)$
		Mixup-Diffusion	71.91% $(\pm 0.19\%)$	83.55% $(\pm 0.22\%)$	87.94% $(\pm 0.21\%)$	88.56% $(\pm 0.18\%)$

**Table 2:** FID values between the train and test sets, and train set and text-to-images diffusion model generated images.

FID model	Measurement	CIFAR-10	CIFAR-100	EuroSAT	Caltech256
InceptionV3 / Vit-B-16	Between training and testing	3.15 / 0.42	3.57 / 0.49	7.41 / 7.26	6.78 / 1.64
	Between training and text-diffusion	30.53 / 30.10	19.79 / 21.85	<b>164.00</b> / <b>82.88</b>	25.14 / 37.39

**Table 3:** FID values for different methods' generated images compared to the training set and test set.

Compared with	Method	CIFAR-10	CIFAR-100	EuroSAT	Caltech256
Train / Test	Self-Aug	2.64 / 3.08	3.09 / 3.53	2.18 / 4.12	1.02 / 2.53
	Self-Aug +	5.95 / 6.37	6.05 / 6.47	6.08 / 8.27	2.88 / 4.23
	Mixup-SelfAug	3.07 / 3.49	3.34 / 3.79	2.99 / 6.14	1.31 / 2.77
	Mixup-Diffusion	3.26 / 3.67	3.46 / 3.92	6.57 / 10.89	1.54 / 2.85

**Table 4:** Vit-B-16 model performance with Self-Augs+ on CIFAR-10,CIFAR-100 and EuroSAT with different  $\varepsilon$

Dataset	Method	$\varepsilon=1$	$\varepsilon=2$	$\varepsilon=4$	$\varepsilon=8$
CIFAR-10	Mixed Ghost clipping	95.08% $(\pm 0.14\%)$	95.00% $(\pm 0.23\%)$	95.28% $(\pm 0.42\%)$	95.33% $(\pm 0.21\%)$
	Self-Aug	96.49% $(\pm 0.12\%)$	96.98% $(\pm 0.02\%)$	97.06% $(\pm 0.03\%)$	97.23% $(\pm 0.03\%)$
	MIXUP-SELF AUG	<b>97.21%</b> $(\pm 0.27\%)$	<b>97.35%</b> $(\pm 0.18\%)$	<b>97.42%</b> $(\pm 0.22\%)$	<b>97.55%</b> $(\pm 0.25\%)$
	Self-Aug+	96.30% $(\pm 0.09\%)$	96.49% $(\pm 0.14\%)$	96.80% $(\pm 0.10\%)$	96.88% $(\pm 0.08\%)$
CIFAR-100	Mixed Ghost clipping	78.16% $(\pm 0.35\%)$	78.53% $(\pm 0.05\%)$	78.36% $(\pm 0.32\%)$	78.43% $(\pm 0.10\%)$
	Self-Aug	79.28% $(\pm 0.18\%)$	83.18% $(\pm 0.33\%)$	83.47% $(\pm 0.30\%)$	84.18% $(\pm 0.08\%)$
	MIXUP-SELF AUG	<b>81.75%</b> $(\pm 0.15\%)$	<b>83.54%</b> $(\pm 0.12\%)$	<b>84.52%</b> $(\pm 0.05\%)$	<b>84.58%</b> $(\pm 0.20\%)$
	Self-Aug+	78.15% $(\pm 0.21\%)$	81.10% $(\pm 0.24\%)$	81.52% $(\pm 0.17\%)$	83.20% $(\pm 0.15\%)$
EuroSAT	Mixed Ghost clipping	84.02% $(\pm 0.09\%)$	84.82% $(\pm 0.15\%)$	84.85% $(\pm 0.07\%)$	85.04% $(\pm 0.15\%)$
	Self-Aug	93.28% $(\pm 0.15\%)$	94.13% $(\pm 0.23\%)$	95.34% $(\pm 0.21\%)$	95.49% $(\pm 0.15\%)$
	MIXUP-SELF AUG	<b>94.32%</b> $(\pm 0.11\%)$	<b>94.91%</b> $(\pm 0.15\%)$	<b>95.56%</b> $(\pm 0.21\%)$	<b>95.58%</b> $(\pm 0.13\%)$
	Self-Aug+	92.20% $(\pm 0.14\%)$	93.87% $(\pm 0.11\%)$	94.27% $(\pm 0.13\%)$	94.33% $(\pm 0.15\%)$

**Table 5:** Change the number of diffusion samples( $k''$ ) for Mixup-Diffusion by using Vit-B-16 model.We set  $\delta = 10^{-5}$  and  $\varepsilon = 1$ .

# of diffusion samples	CIFAR-100	EuroSAT	Caltech 256	Oxford-IIIT PET
2	82.02% $(\pm 0.11\%)$	92.58% $(\pm 0.14\%)$	85.82% $(\pm 0.12\%)$	73.67% $(\pm 0.25\%)$
4	81.91% $(\pm 0.15\%)$	92.31% $(\pm 0.18\%)$	86.68% $(\pm 0.15\%)$	77.07% $(\pm 0.22\%)$
8	81.86% $(\pm 0.12\%)$	90.93% $(\pm 0.08\%)$	88.59% $(\pm 0.07\%)$	78.59% $(\pm 0.18\%)$
16	81.26% $(\pm 0.18\%)$	89.38% $(\pm 0.12\%)$	89.28% $(\pm 0.11\%)$	83.32% $(\pm 0.23\%)$

**Table 6:** Increase K from 16 to 36 for Self-Aug. We conduct experiments on CIFAR-10 with  $\varepsilon = 8$  and  $\delta = 10^{-5}$  in two settings: train a WRN-16-4 model from scratch and fine-tune a pretrained Vit-B-16. We can observe that for both cases, there are no substantial performance improvements when increasing K.

Training method	K	WRN-16-4	Vit-B-16 (pretrained)
Self-Aug	16	78.74% $(\pm 0.45\%)$	97.23% $(\pm 0.03\%)$
	24	78.49% $(\pm 0.42\%)$	96.95% $(\pm 0.11\%)$
	32	78.40% $(\pm 0.34\%)$	97.04% $(\pm 0.05\%)$
	36	78.55% $(\pm 0.38\%)$	96.96% $(\pm 0.08\%)$
Mixup-SelfAug	32	<b>79.83%</b> $(\pm 0.32\%)$	<b>97.55%</b> $(\pm 0.25\%)$