

# Robust Combination of Distributed Gradients Under Adversarial Perturbations

## – Supplemental Document

Kwang In Kim  
UNIST

**Experiments with random gradient perturbations.** In the main paper, we simulated *random labeling* perturbations. Here, we report the results of experiments performed with *random gradient* perturbations which replace the local gradients with independently and identically distributed Gaussian random vectors. Table 1 summarizes the results. The results are consistent with the experiments with the random labeling perturbations. Federated averaging (*FedAvg*) performance rapidly degraded as the number of affected clients increased. The *Median* estimators were comparable to or slightly worse than *FedAvg*. The geometric median estimators (*GeoMed*) were significantly better than *Median*, yet it was noticeably inferior to *FedAvg* in a few cases. Overall, *MedTh* and ours constantly outperformed *FedAvg* and achieved the best performance.

**Combining *MedTh* and our algorithm.** Our final algorithm estimates the continuous gradient combination weights  $\mathbf{a}$ . This helps suppress perturbed gradients, but it does not completely rule out their contributions. Here, we demonstrate that we can achieve this goal by combining ours with the thresholding approach of *MedTh*: Each element of  $\mathbf{a}$  is set at zero when it is smaller than a threshold. For the experiments with *MedTh*, we determined this threshold such that the number of selected gradients is the same as the number of genuine gradients. As such information is not available in practice, we determine this number using a separate validation set at the server: We performed experiments with a validation set of the same size as the test set. Since the accuracy of the learner network is low at the early stages of training, we exercised this strategy only after 200 iterations of the training process. This (*MedTh + Ours* in Table 2) led to significant performance improvements over our original design. However, as *MedTh + Ours* requires validation sets at the server, its application domain can be limited.

Dataset	Method	Accuracy					
		10	20	30	40	50	60
<i>Caltech256</i>	<i>GTC</i>	76.54 (0.34)	75.96 (0.20)	75.47 (0.19)	74.84 (0.41)	73.83 (0.16)	72.26 (0.50)
	<i>FedAvg</i>	76.48 (0.32)	75.77 (0.51)	74.81 (0.19)	73.63 (0.24)	72.12 (0.19)	69.61 (0.42)
	<i>Median</i>	<b>72.05 (0.49)</b>	<b>70.13 (0.79)</b>	<b>68.19 (0.52)</b>	<b>64.41 (0.55)</b>	<b>58.92 (1.24)</b>	<b>50.24 (1.60)</b>
	<i>GeoMed</i>	<b>73.17 (0.37)</b>	<b>73.20 (0.31)</b>	<b>73.11 (0.26)</b>	<b>73.05 (0.31)</b>	<b>73.22 (0.27)</b>	<b>73.33 (0.45)</b>
	<i>MedTh</i>	<b>76.84 (0.36)</b>	<b>75.98 (0.31)</b>	<b>75.10 (0.25)</b>	<b>74.45 (0.38)</b>	73.01 (0.40)	71.13 (0.52)
	Ours	76.63 (0.33)	75.89 (0.44)	<b>75.23 (0.32)</b>	<b>74.50 (0.34)</b>	<b>73.58 (0.28)</b>	72.25 (0.41)
<i>CIFAR10</i>	<i>GTC</i>	89.53 (0.12)	89.51 (0.18)	89.45 (0.19)	89.53 (0.23)	89.41 (0.14)	89.20 (0.11)
	<i>FedAvg</i>	89.26 (0.15)	89.21 (0.13)	88.87 (0.07)	88.42 (0.27)	87.86 (0.16)	87.10 (0.11)
	<i>Median</i>	<b>88.02 (0.15)</b>	<b>87.88 (0.19)</b>	<b>87.34 (0.11)</b>	<b>86.98 (0.17)</b>	<b>86.30 (0.39)</b>	<b>85.47 (0.31)</b>
	<i>GeoMed</i>	<b>87.11 (0.08)</b>	<b>87.22 (0.16)</b>	<b>87.11 (0.08)</b>	87.78 (1.17)	<b>87.09 (0.27)</b>	87.97 (1.34)
	<i>MedTh</i>	89.27 (0.25)	89.12 (0.20)	<b>89.18 (0.15)</b>	<b>88.99 (0.20)</b>	<b>88.82 (0.16)</b>	<b>88.70 (0.12)</b>
	Ours	<b>89.33 (0.28)</b>	<b>89.47 (0.27)</b>	<b>89.40 (0.14)</b>	<b>89.37 (0.15)</b>	<b>89.27 (0.27)</b>	<b>88.95 (0.17)</b>
<i>CIFAR100</i>	<i>GTC</i>	66.86 (0.23)	66.50 (0.21)	66.10 (0.33)	65.81 (0.30)	65.14 (0.42)	64.41 (0.45)
	<i>FedAvg</i>	66.33 (0.26)	65.66 (0.22)	64.68 (0.24)	63.45 (0.47)	61.77 (0.47)	59.33 (0.50)
	<i>Median</i>	66.47 (0.31)	<b>66.25 (0.26)</b>	<b>65.59 (0.23)</b>	<b>65.11 (0.50)</b>	<b>64.62 (0.36)</b>	<b>63.58 (0.32)</b>
	<i>GeoMed</i>	<b>63.09 (0.44)</b>	<b>62.90 (0.33)</b>	<b>62.83 (0.24)</b>	<b>62.76 (0.39)</b>	63.03 (0.26)	62.98 (0.39)
	<i>MedTh</i>	<b>66.67 (0.31)</b>	<b>66.29 (0.19)</b>	<b>66.03 (0.43)</b>	<b>65.66 (0.39)</b>	64.96 (0.29)	<b>64.29 (0.37)</b>
	Ours	66.66 (0.22)	66.27 (0.19)	66.00 (0.27)	65.55 (0.51)	<b>65.02 (0.22)</b>	64.17 (0.37)
<i>FashionMNIST</i>	<i>GTC</i>	86.96 (0.18)	86.84 (0.23)	86.84 (0.11)	86.92 (0.19)	86.77 (0.19)	86.77 (0.20)
	<i>FedAvg</i>	86.46 (0.15)	85.83 (0.23)	85.21 (0.24)	84.26 (0.29)	83.32 (0.27)	81.78 (0.38)
	<i>Median</i>	<b>84.26 (0.38)</b>	<b>83.82 (0.32)</b>	<b>83.05 (0.27)</b>	<b>82.07 (0.47)</b>	<b>81.14 (0.36)</b>	<b>80.18 (0.35)</b>
	<i>GeoMed</i>	<b>85.32 (0.91)</b>	85.79 (1.52)	<b>84.74 (0.23)</b>	84.85 (0.19)	85.63 (0.98)	84.76 (0.33)
	<i>MedTh</i>	<b>86.58 (0.11)</b>	<b>86.45 (0.20)</b>	<b>86.33 (0.23)</b>	<b>86.39 (0.20)</b>	<b>86.30 (0.21)</b>	86.32 (0.29)
	Ours	<b>86.92 (0.18)</b>	<b>86.78 (0.20)</b>	<b>86.67 (0.23)</b>	<b>86.64 (0.13)</b>	<b>86.65 (0.11)</b>	<b>86.42 (0.21)</b>
<i>CINIC10</i>	<i>GTC</i>	79.69 (0.08)	79.73 (0.14)	79.60 (0.16)	79.67 (0.06)	79.69 (0.13)	79.59 (0.11)
	<i>FedAvg</i>	79.35 (0.13)	<b>79.20 (0.05)</b>	78.67 (0.10)	78.03 (0.07)	77.30 (0.17)	75.96 (0.43)
	<i>Median</i>	<b>76.99 (0.06)</b>	<b>76.41 (0.19)</b>	<b>76.04 (0.23)</b>	<b>74.86 (0.25)</b>	<b>74.36 (0.45)</b>	<b>73.18 (0.39)</b>
	<i>GeoMed</i>	<b>77.49 (0.03)</b>	<b>77.72 (0.16)</b>	<b>77.65 (0.32)</b>	<b>77.51 (0.13)</b>	77.54 (0.19)	<b>77.41 (0.08)</b>
	<i>MedTh</i>	<b>79.45 (0.15)</b>	79.18 (0.03)	<b>79.19 (0.25)</b>	<b>79.22 (0.06)</b>	<b>79.30 (0.08)</b>	79.20 (0.07)
	Ours	<b>80.02 (0.09)</b>	<b>79.95 (0.07)</b>	<b>79.81 (0.12)</b>	<b>79.88 (0.06)</b>	<b>79.84 (0.06)</b>	<b>79.65 (0.14)</b>

Table 1. Results of classification under random gradient perturbations: Mean classification accuracy and standard deviation in parenthesis (both in %; higher is better) on random gradient perturbations. The best and second-best results (except for GTC) are highlighted with **bold** and *italic*, respectively. For *Median*, *MedTh*, and *Ours*, the results of statistical significance test with respect to *FedAvg* (t-test with  $\alpha = 0.95$ ) are shown: Blue and orange respectively represents significantly better and worse results than *FedAvg*.

Dataset	Method	Accuracy					
		10	20	30	40	50	60
Caltech256	<i>GTC</i>	76.54 (0.34)	75.96 (0.20)	75.47 (0.19)	74.84 (0.41)	73.83 (0.16)	72.26 (0.50)
	<i>FedAvg</i>	76.29 (0.32)	75.16 (0.35)	74.14 (0.45)	72.29 (0.33)	70.50 (0.27)	66.75 (0.44)
	<i>Median</i>	<b>64.80 (0.47)</b>	<b>63.58 (0.97)</b>	<b>61.82 (0.74)</b>	<b>59.15 (0.98)</b>	<b>55.57 (0.85)</b>	<b>48.15 (0.99)</b>
	<i>GeoMed</i>	76.29 (0.36)	75.68 (0.41)	74.69 (0.44)	<b>71.70 (0.53)</b>	66.49 (0.36)	60.77 (0.82)
	<i>MedTh</i>	76.28 (0.27)	75.59 (0.38)	74.70 (0.85)	72.51 (1.17)	69.74 (1.04)	46.31 (5.19)
	Ours	76.39 (0.20)	75.95 (0.29)	75.06 (0.26)	73.72 (0.31)	71.90 (0.50)	69.25 (0.47)
	MedTh + Ours	<b>76.76 (0.34)</b>	<b>76.01 (0.08)</b>	<b>75.56 (0.33)</b>	<b>74.88 (0.17)</b>	<b>73.93 (0.20)</b>	<b>72.48 (0.22)</b>
CIFAR10	<i>GTC</i>	89.57 (0.11)	89.49 (0.15)	89.44 (0.15)	89.55 (0.20)	89.36 (0.16)	89.22 (0.10)
	<i>FedAvg</i>	89.23 (0.10)	87.34 (0.23)	86.58 (0.15)	85.98 (0.37)	85.24 (0.41)	84.83 (0.37)
	<i>Median</i>	<b>87.21 (0.25)</b>	<b>86.89 (0.26)</b>	86.74 (0.82)	86.30 (1.37)	84.65 (1.19)	<b>82.01 (0.62)</b>
	<i>GeoMed</i>	<b>87.38 (0.47)</b>	87.49 (1.06)	87.14 (1.07)	86.89 (1.83)	84.99 (0.85)	<b>82.46 (0.32)</b>
	<i>MedTh</i>	89.39 (0.18)	89.27 (0.10)	89.01 (0.19)	88.57 (0.25)	78.10 (19.90)	<b>31.40 (18.43)</b>
	Ours	<b>89.51 (0.06)</b>	<b>89.43 (0.17)</b>	89.22 (0.19)	89.14 (0.13)	<b>88.95 (0.16)</b>	88.85 (0.20)
	MedTh + Ours	89.47 (0.09)	<b>89.47 (0.09)</b>	<b>89.49 (0.17)</b>	<b>89.25 (0.18)</b>	<b>89.27 (0.16)</b>	<b>89.19 (0.16)</b>
CIFAR100	<i>GTC</i>	66.86 (0.23)	66.50 (0.21)	66.10 (0.33)	65.81 (0.30)	65.14 (0.42)	64.41 (0.45)
	<i>FedAvg</i>	66.50 (0.40)	65.95 (0.54)	64.67 (0.28)	63.25 (0.40)	61.99 (0.32)	60.53 (0.33)
	<i>Median</i>	<b>65.76 (0.23)</b>	<b>64.53 (0.29)</b>	<b>62.73 (0.31)</b>	<b>59.51 (0.75)</b>	<b>53.51 (0.85)</b>	<b>47.89 (0.65)</b>
	<i>GeoMed</i>	<b>67.11 (0.40)</b>	<b>66.80 (0.63)</b>	65.75 (0.69)	65.20 (0.23)	<b>58.63 (0.57)</b>	49.79 (1.05)
	<i>MedTh</i>	66.80 (0.14)	66.35 (0.29)	65.78 (0.24)	64.11 (0.33)	38.34 (1.43)	<b>36.93 (2.52)</b>
	Ours	66.80 (0.23)	66.41 (0.38)	65.87 (0.38)	65.44 (0.27)	64.80 (0.20)	63.56 (0.62)
	MedTh + Ours	66.99 (0.32)	66.46 (0.31)	<b>65.91 (0.31)</b>	<b>65.70 (0.36)</b>	<b>65.14 (0.28)</b>	<b>64.54 (0.38)</b>
FashionMNIST	<i>GTC</i>	86.96 (0.18)	86.84 (0.23)	86.84 (0.11)	86.92 (0.19)	86.77 (0.19)	86.77 (0.20)
	<i>FedAvg</i>	86.48 (0.19)	83.12 (0.24)	81.93 (0.35)	81.30 (0.31)	80.80 (0.50)	80.47 (0.26)
	<i>Median</i>	<b>83.23 (1.44)</b>	82.77 (1.30)	82.68 (1.71)	82.02 (1.57)	80.44 (1.05)	<b>78.40 (0.44)</b>
	<i>GeoMed</i>	<b>83.29 (1.65)</b>	<b>82.37 (0.40)</b>	82.42 (1.72)	81.20 (0.36)	80.41 (0.45)	<b>78.37 (0.30)</b>
	<i>MedTh</i>	86.75 (0.11)	86.55 (0.22)	<b>86.27 (0.22)</b>	83.57 (4.56)	76.16 (14.07)	<b>52.65 (23.84)</b>
	Ours	<b>87.01 (0.23)</b>	86.73 (0.16)	86.63 (0.27)	<b>86.40 (0.51)</b>	<b>86.24 (0.16)</b>	86.41 (0.27)
	MedTh + Ours	86.92 (0.28)	<b>86.82 (0.23)</b>	<b>86.73 (0.22)</b>	<b>86.64 (0.19)</b>	<b>86.47 (0.16)</b>	<b>86.43 (0.27)</b>
CINIC10	<i>GTC</i>	79.69 (0.07)	79.66 (0.11)	79.65 (0.12)	79.70 (0.08)	79.66 (0.12)	79.60 (0.15)
	<i>FedAvg</i>	78.20 (0.11)	76.88 (0.21)	76.04 (0.16)	75.41 (0.29)	74.81 (0.27)	74.38 (0.13)
	<i>Median</i>	77.16 (1.59)	76.91 (1.61)	77.03 (1.95)	<b>76.84 (1.33)</b>	73.83 (1.44)	<b>71.04 (0.72)</b>
	<i>GeoMed</i>	78.13 (1.95)	<b>79.01 (1.86)</b>	77.49 (2.61)	<b>78.92 (1.20)</b>	75.64 (1.28)	<b>71.51 (0.75)</b>
	<i>MedTh</i>	79.48 (0.11)	79.31 (0.15)	79.06 (0.10)	78.57 (0.15)	<b>57.10 (23.87)</b>	42.19 (20.96)
	Ours	<b>80.03 (0.13)</b>	<b>79.95 (0.17)</b>	<b>79.85 (0.16)</b>	<b>79.64 (0.32)</b>	<b>79.19 (0.50)</b>	78.69 (0.52)
	MedTh + Ours	79.68 (0.09)	79.69 (0.13)	79.25 (0.41)	79.00 (0.60)	79.06 (0.57)	<b>79.15 (0.42)</b>

Table 2. Results of classification under random labeling perturbations: Mean accuracy and standard deviation (in parenthesis) are shown in % (higher is better). The best and second-best results excluding *GTC* are highlighted with **bold** and *italic*, respectively. For *Median*, *GeoMed*, *MedTh*, and ours, the results of statistical significance test with respect to *FedAvg* (t-test with  $\alpha = 0.95$ ) are shown: **Blue** and **orange** respectively represents significantly better and worse results than *FedAvg*.