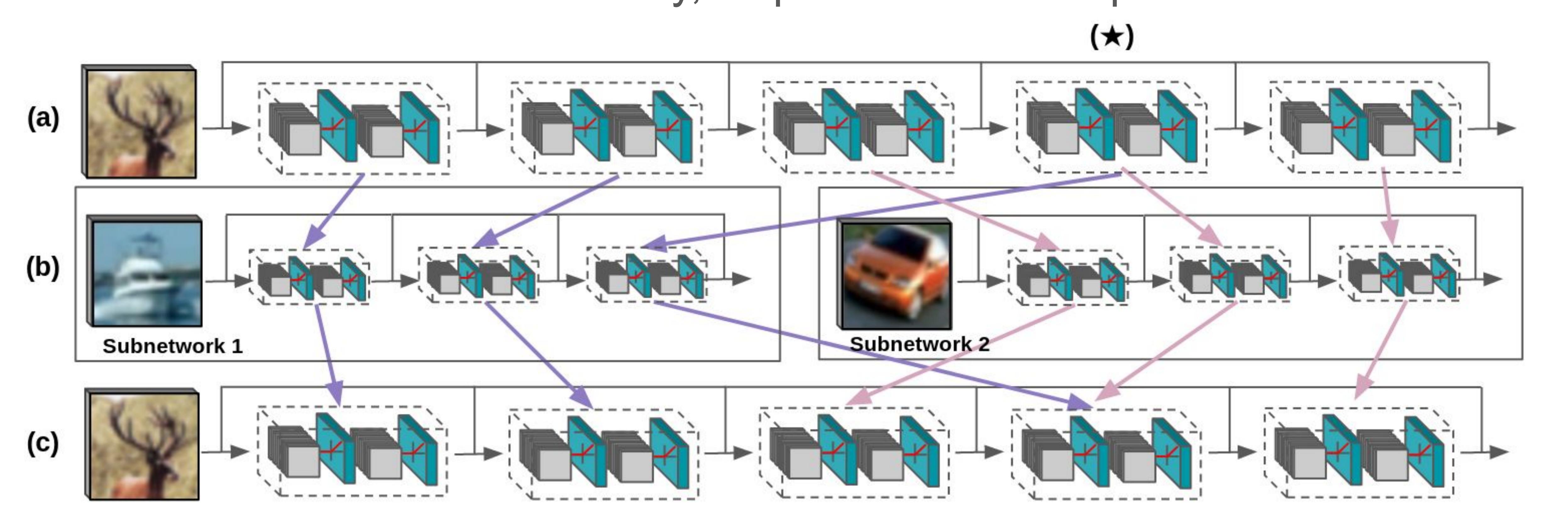
ResIST: Layer-wise Decomposition of ResNets for Distributed Training RICE (intel)

Chen Dun, Cameron Wolfe, Christopher M. Jermaine, Anastasios Kyrillidis Rice University, Department of Computer Science



```
Algorithm 1 ResIST Meta AlgorithmParameters: T synchronization iterations, S sub-ResNets, \ell local iterations, W ResNet weights.h(W) \leftarrow randomly initialized ResNet.for t = 0, ..., T - 1 do\{h_s(W_s)\}_{s=1}^S = \text{subResNets}(h(W), S).Distribute each h_s(W_s) to a different worker.for s = 1, ..., S doTrain h_s(W_s) for \ell iterations using local SGD.end forh(W) = \text{aggregate}\left(\{h_s(W_s)\}_{s=1}^S\right).
```

Motivation

Computer vision (CV) task is one of the most important real life scenarios for Federated learning, such as face recognition on mobile device and medical image classification on remote device.

While deep ResNet[1] is the most successful model for CV, when applied in Federated Learning scenario, its large model size results in high communication cost, high computation cost and high memory cost on each edge device.

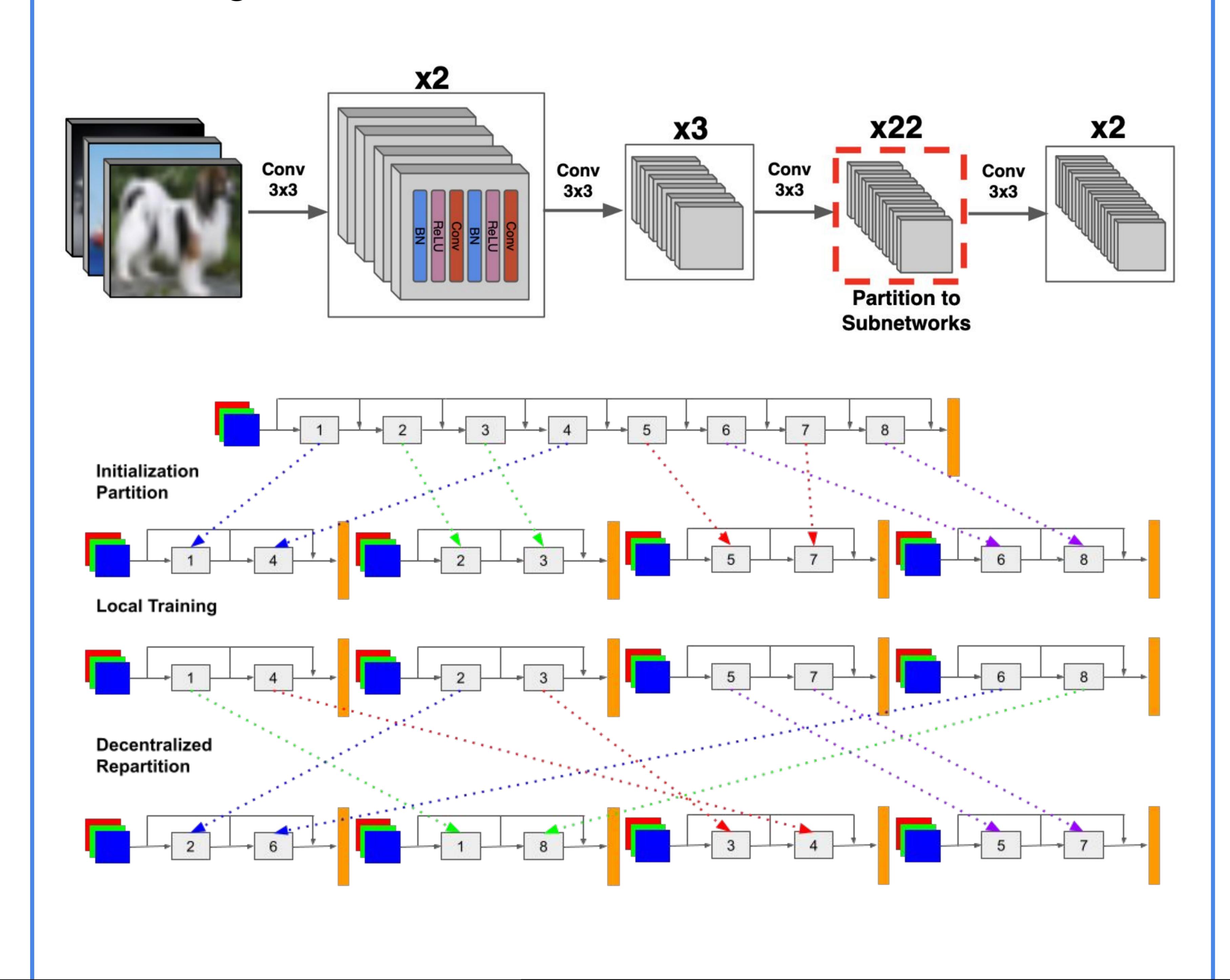
Contribution

We propose ResIST: A novel efficient Federated Learning algorithm for ResNet, which significantly reduces communication, computation and memory cost without losing accuracy.

	Data Parallel	Model Parallel	ResiST
Memory Cost	High	LOW	
Comm Frequency	Low	High	LOW
Comm Volume in each round	High	LOW	LOW
Computation Cost	High	High	LOW

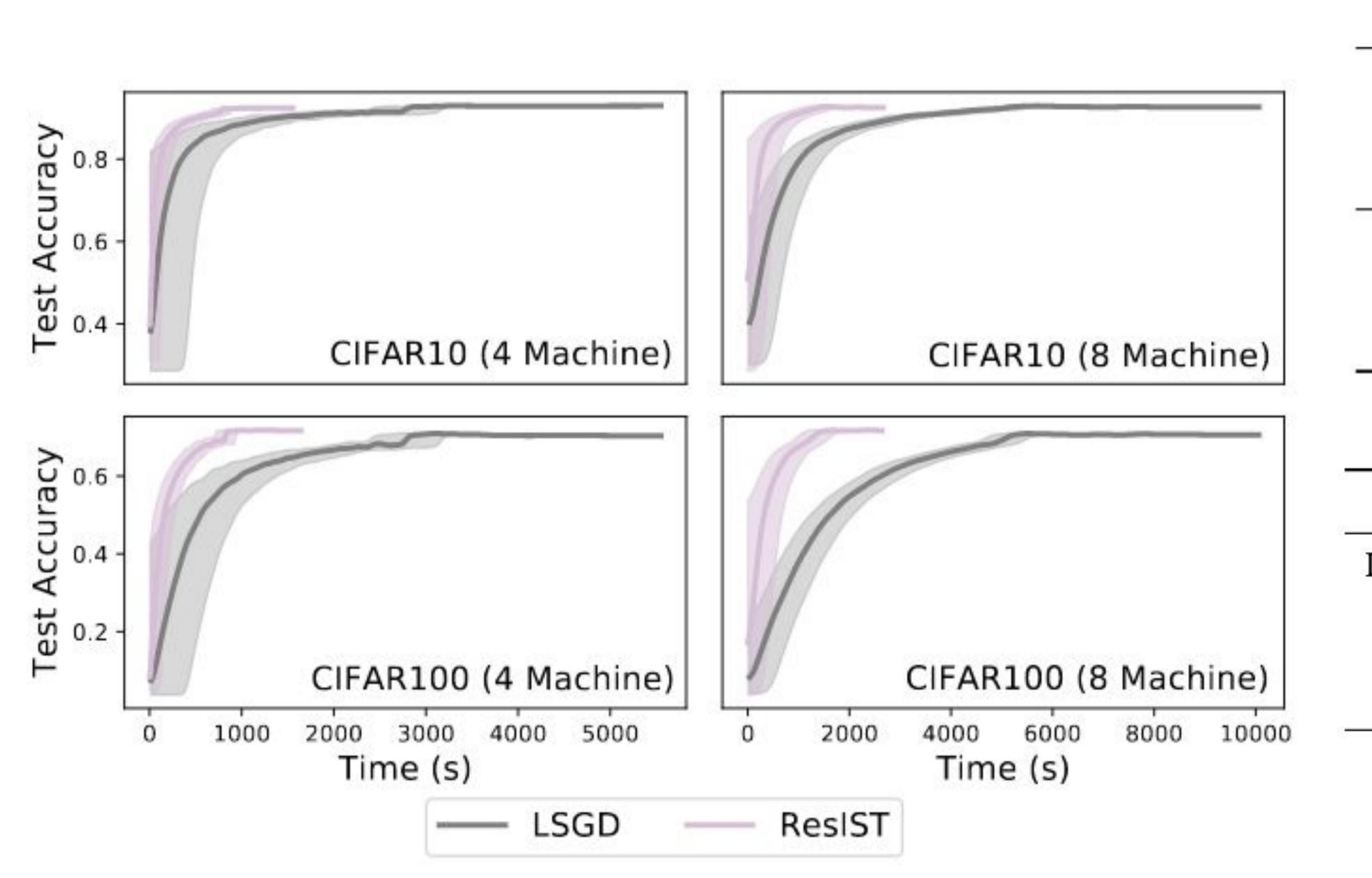
Method

ResIST randomly decomposes a global ResNet into several shallow sub-ResNets that are trained independently on each edge device for several local iterations, before having their updates synchronized and aggregated into the global model. In the next round, new sub-ResNets are randomly generated and the process repeats. By construction, per iteration, ResIST communicates only a small portion of network parameters to each machine and never uses the full model during training.



Result

ResIST is tested on Federated Learning scenario with (1) image classification datasets CIFAR10, CIFAR100 and Imagenet and (2) object detection dataset Pascal VOC. Local SGD (LSGD) is used as baseline.

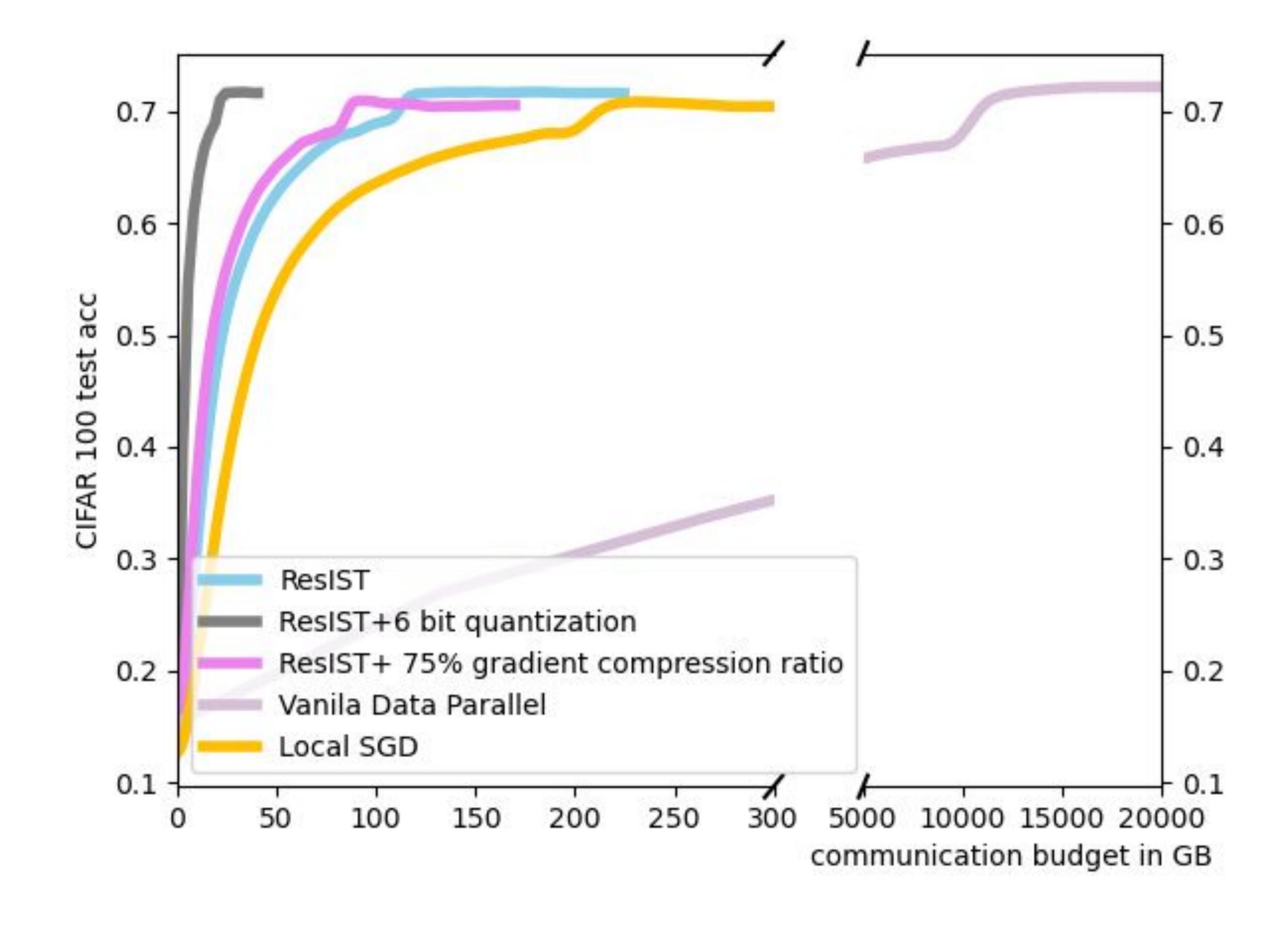


	4	92.90%	6 ± 0.06	71.5	$1\% \pm 0.04$
	8	92.00%	6 ± 0.07	69.6	$4\% \pm 0.05$
ResIST	2	91.95%	% ± 0.32	70.0	6% ± 0.51
	4	92.35%	6 ± 0.22	71.3	$0\% \pm 0.20$
	8	91.45%	6 ± 0.30	70.2	$6\% \pm 0.21$
	# Machines	Dataset	Total Ti	ime	Speedur
Local SGD	4	C10	5486 ± 7.05		- P
20000	-	C100	5528 ± 6		_
	8	C10	$10072 \pm$	5.12	_
		C100	$10058~\pm$	8.71	-
ResIST	4	C10	1532 ± (0.83	3.60×
		C100	1545 ± 1	1.27	$3.58 \times$
	8	C10	2671 ± 3	3.25	$3.77 \times$

C100 2639 ± 3.89 3.81×

	# Machines	Imagenet	Image MF	enet V2 Te T-0.7	est Set TI	Training Time	Speedup	Communication	Cost Ratio
Local SGD	2 4	73.32% 72.66%	60.72% 59.88%	69.47% 68.34%	75.48% 74.27%	48.61 hours 29.29 hours	_	7546.80 GB 7546.80 GB	=
ResIST	2	71.60% 70.74%	58.92% 57.56%	67.51% 66.46%	73.56% 72.65%	36.79 hours 22.37 hours	1.32× 1.31×	5831.2 GB 6007.6 GB	1.29× 1.26×

ResIST is also fully compatible with Quantization[2] and Sparsification[3] method in Federated Learning. When combined, ResIST can further reduces the communication cost.



^[1] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 770-778).

[2] Alistarh, D., Grubic, D., Li, J., Tomioka, R., & Vojnovic, M. (2017). QSGD: Communication-efficient SGD via gradient

quantization and encoding. *Advances in Neural Information Processing Systems*, 30, 1709-1720.

[3] Aji, A. F., & Heafield, K. (2017). Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021*.