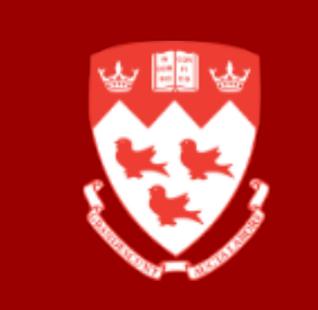


COMET: A Novel Memory-Efficient Deep Learning Training Framework by Using Error-Bounded Lossy Compression



Evaluation







Sian Jin¹, Chengming Zhang¹, Xintong Jiang², Yunhe Feng³, Hui Guan⁴, Guanpeng Li⁵, Shuaiwen Leon Song⁶, Dingwen Tao¹

¹Indiana University, ²McGill University, ³University of Washington, ⁴University of Massachusetts, ⁵University of Iowa, ⁶University of Sydney

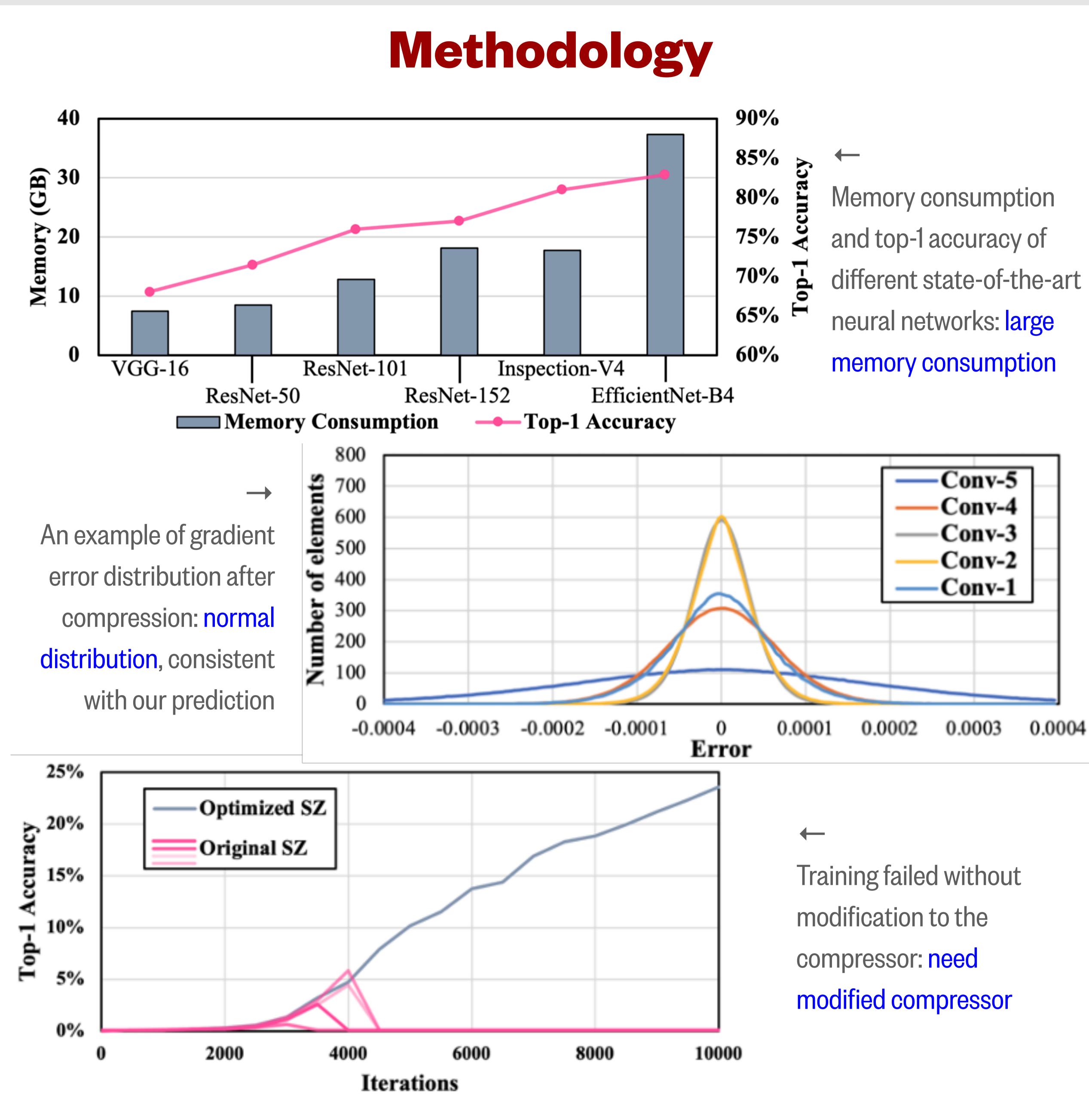
Abstract

Training wide and deep neural networks require large amounts of storage resources such as memory because the intermediate activation data must be saved in the memory during forward propagation and then restored for backward propagation. However, state-of-the-art accelerators such as GPUs are only equipped with very limited memory capacities due to hardware design constraints, which significantly limits the maximum batch size and hence performance speedup when training large-scale DNNs. In this paper, we propose a novel memory-efficient CNN training framework (called COMET) that leverages error-bounded lossy compression to significantly reduce the memory requirement for training in order to allow training larger models or to accelerate training. Our framework purposely adopts error-bounded lossy compression with a strict error-controlling mechanism.

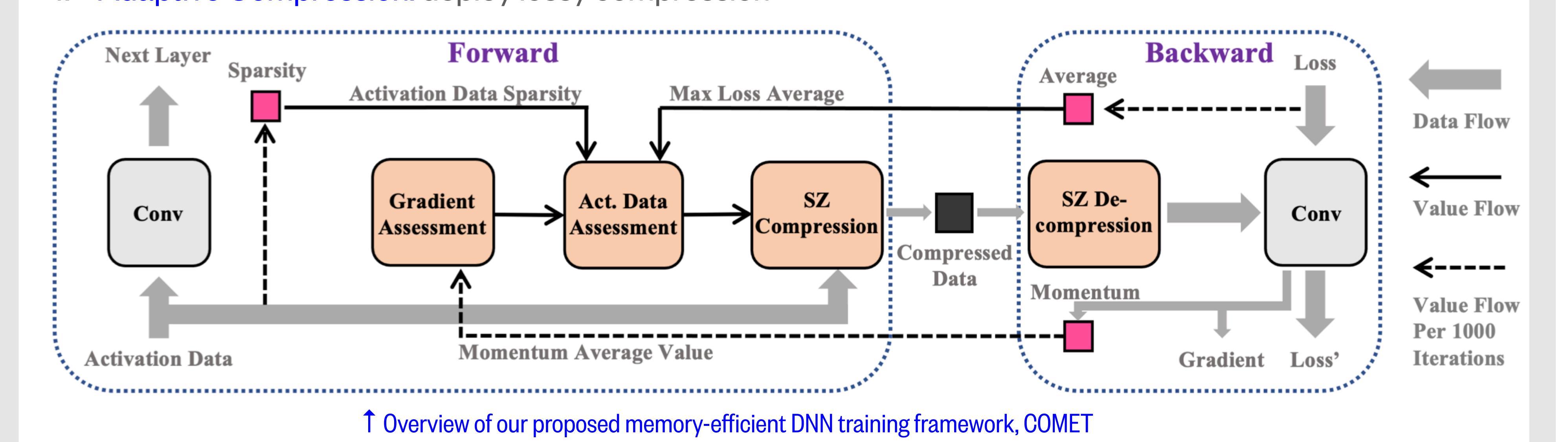
Contribution

This work has mainly five contributions:

- 1. A novel memory-efficient CNN training framework via dynamically compressing the intermediate activation data through error-bounded lossy compression
- 2. A thorough analysis of the impact of compression error propagation during DNN training from both theoretical and empirical perspectives
- 3. An adaptive scheme to adaptively configure the errorbounded lossy compression based on a series of current training status data
- 4. Improved SZ error-bounded lossy compression to handle compressing continuous zeros
- 5. Reduce the memory consumption by up to 13.5× and 1.8× compared to the original training framework and the state-of-the-art method, respectively. Improve the end-to-end training performance by up to 2×



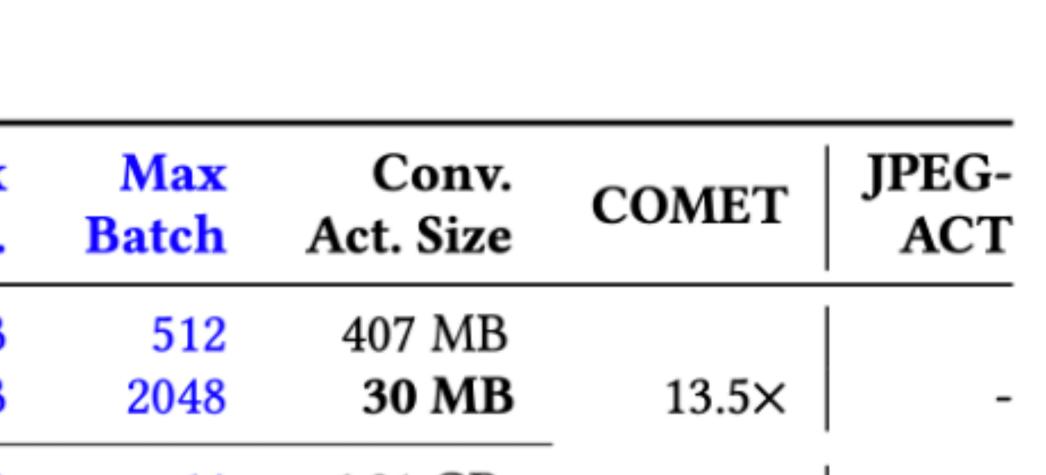
- 1. Parameter Collection: collect parameters for analysis and updating compression configurations
- 2. Gradient Assessment: estimate acceptable variance in the gradient
- 3. Activation Assessment: estimate acceptable error introduced for compressing activation data
- 4. Adaptive Compression: deploy lossy compression



Accurate theoretical prediction to the gradient error distribution

Real Distribution

2.5E-04

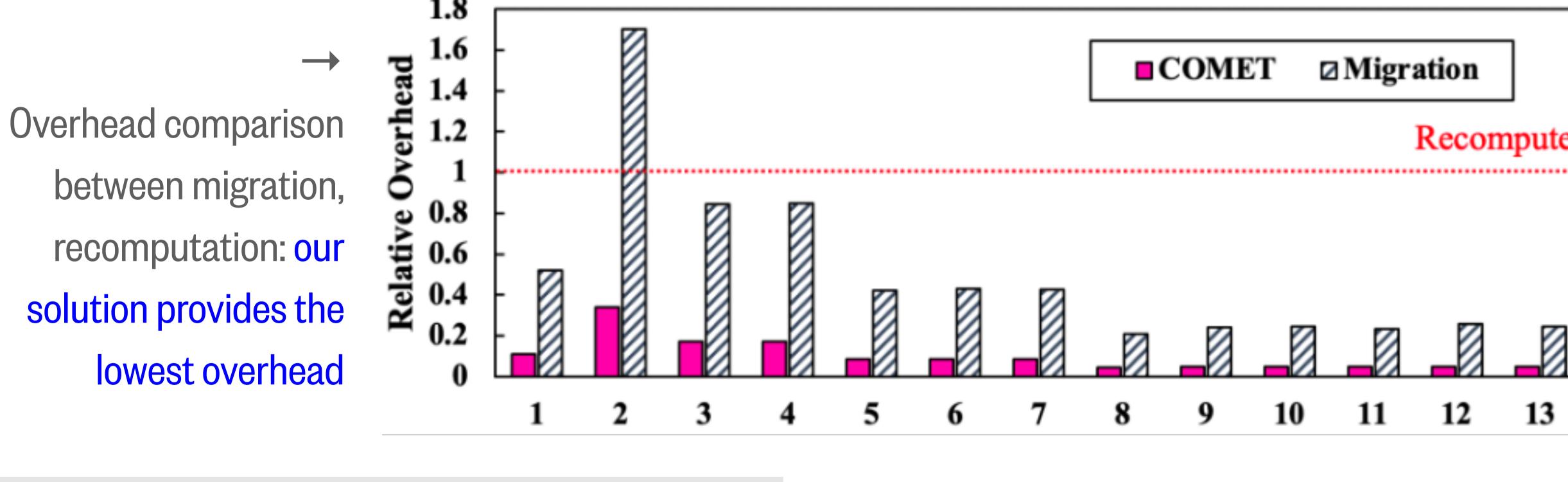


 b.
 57.41%
 2.17 GB
 512
 407 MB
 407 MB

b.= baseline, c.= compressed

Neural Nets

24 Z curve comparison 20 E between the baseline 16 and our proposed framework: high compression ratio -ResNet-50 Ori with almost no AlexNet COMET ——AlexNet Ori 10% accuracy AlexNet Ratio ResNet-50 Ratio degradation **Epochs**



0.0E+00

0.0E + 00