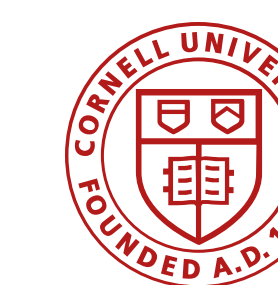# SWALP: Stochastic Weight Averaging for Low-Precision Training

Guandao Yang, Tianyi Zhang, Polina Kirichenko, Junwen Bai, Andrew Gordon Wilson, Christopher De Sa
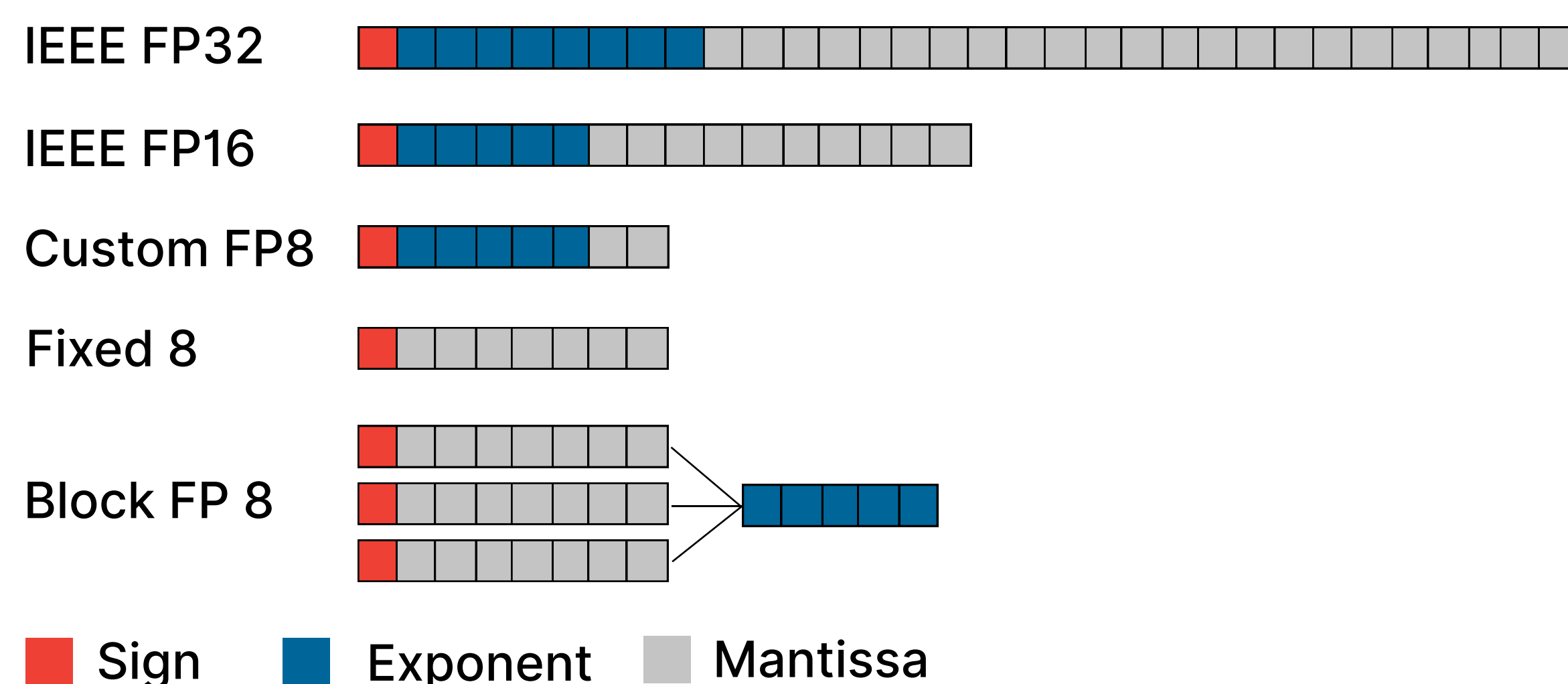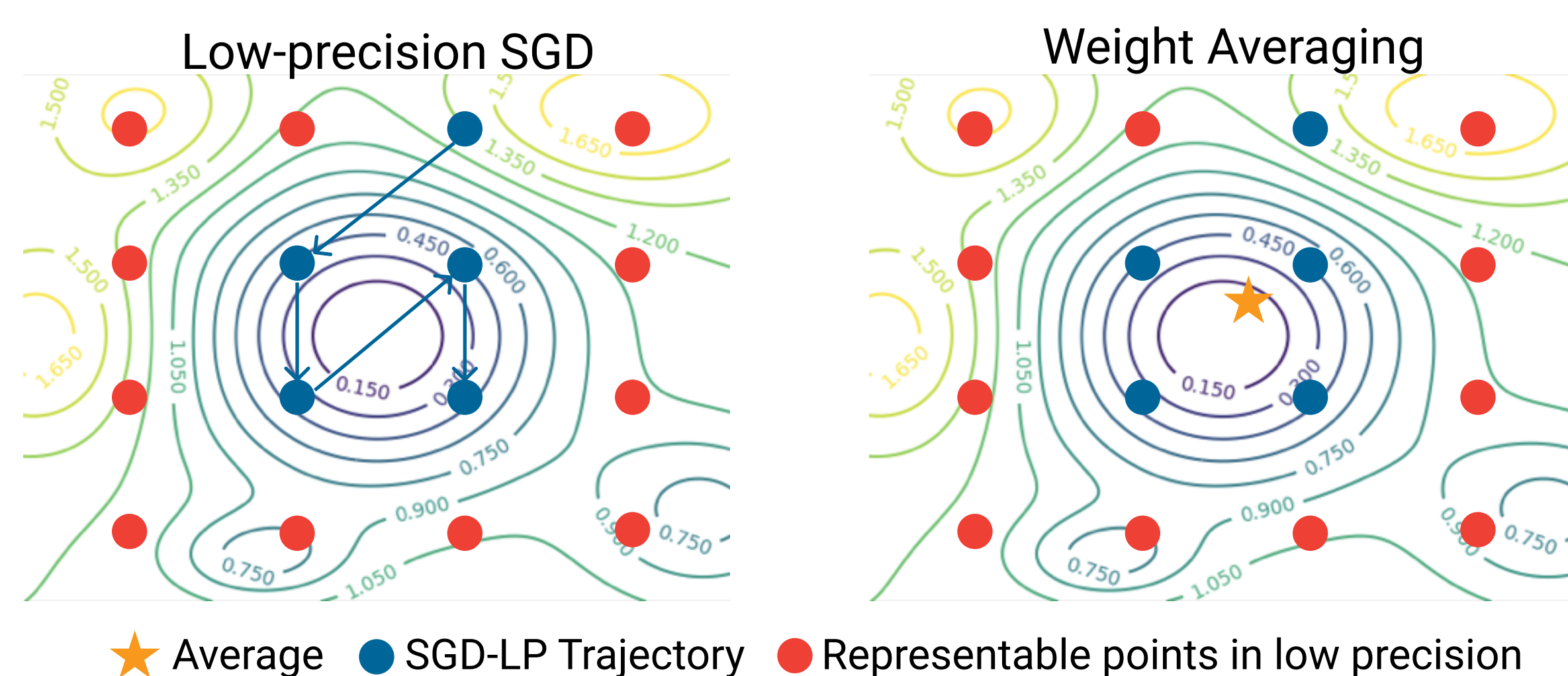
Cornell University

## This work

- Studies how to leverage low-precision training to obtain a high-accuracy model, which may be higher-precision.
- Proposes a principled approach to using stochastic weight averaging in low-precision training (SWALP).
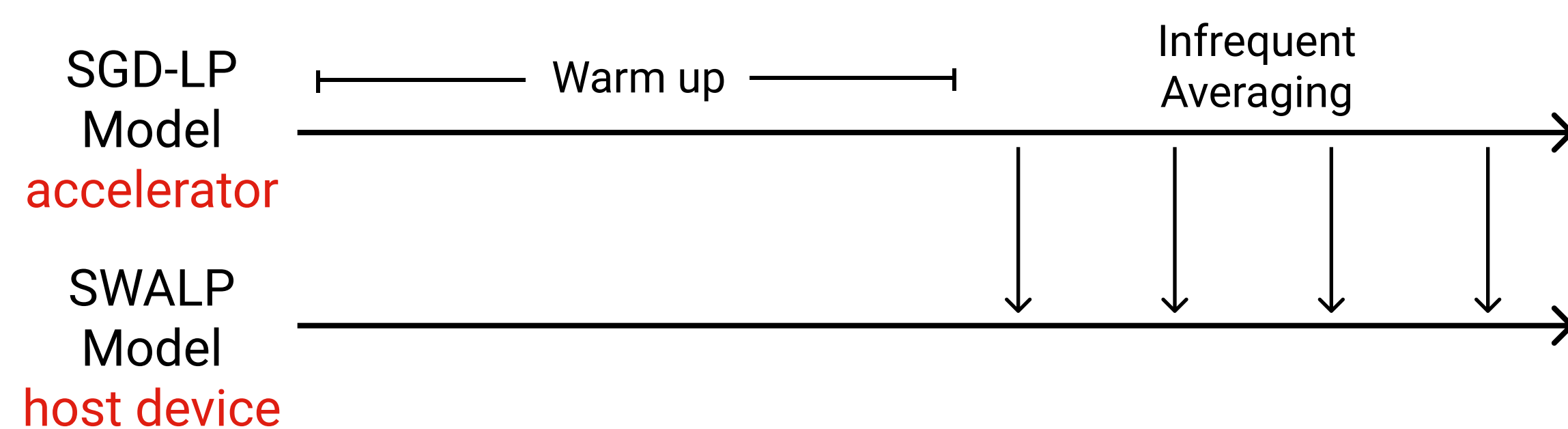- Shows SWA sigificantly reduce the performance gap between low-precision and full-precision training.

## Low-Precision Computation



IEEE FP32
IEEE FP16
Custom FP8
Fixed 8
Block FP 8

Sign   Exponent   Mantissa

## SWALP



Low-precision SGD          Weight Averaging

★ Average   ● SGD-LP Trajectory   ● Representable points in low precision

- Low-precision representation inherently limits the accuracy.
- By averaging, we hope to recover a better solution.



SGD-LP Model — accelerator — Warm up
SWALP Model — host device — Infrequent Averaging

## Convergence Analysis

Let T be the number of iterations and δ be the quantization gap (the difference between two successive representable numbers). With standard assumptions and fixed point quantization, we can prove the following statements.

### Theorem 1 (Quadratic)
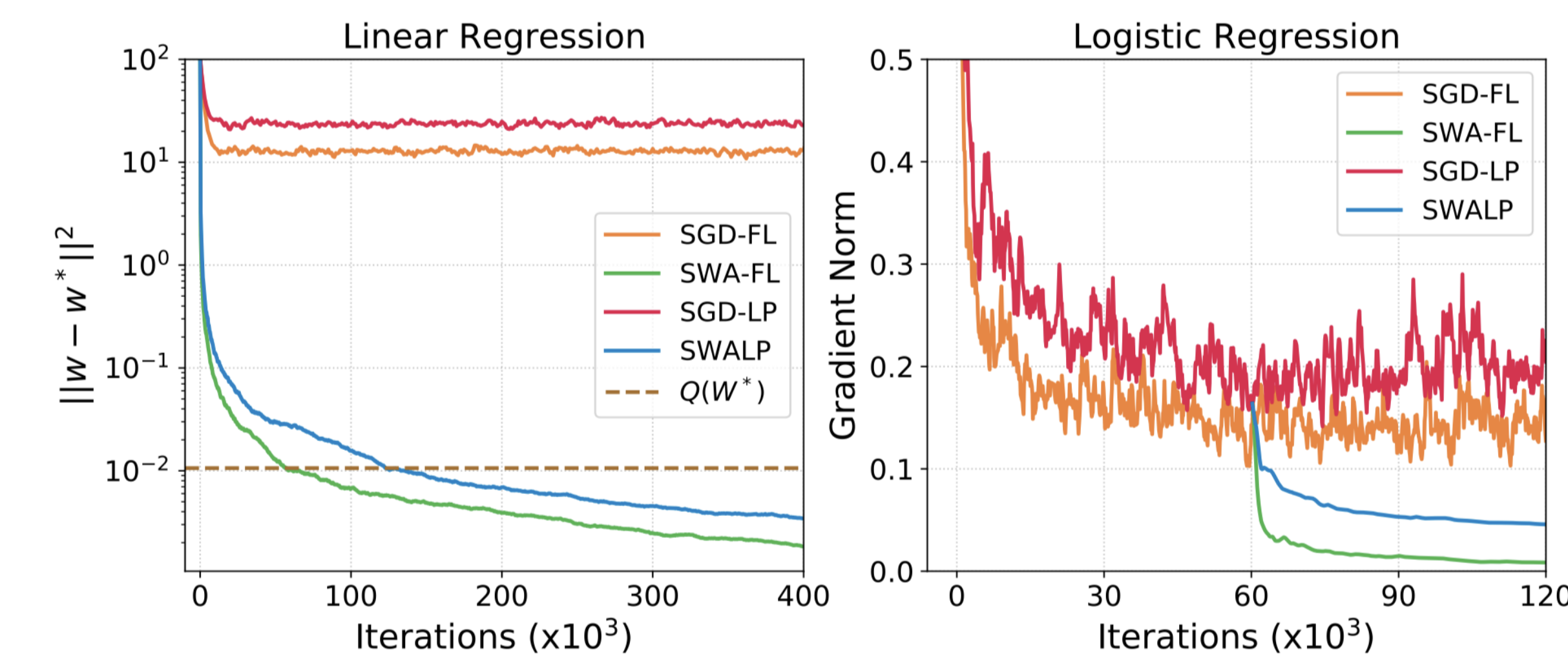The expected squared distance between the SWALP solution and the optimal one converges to 0 at a O(1/T) rate.

- SWALP has the same O(1/T) convergence rate with full-precision SGD.
- SWALP converges to the optimal solution regardless of the numerical precision.

### Theorem 2 (Strongly Convex)
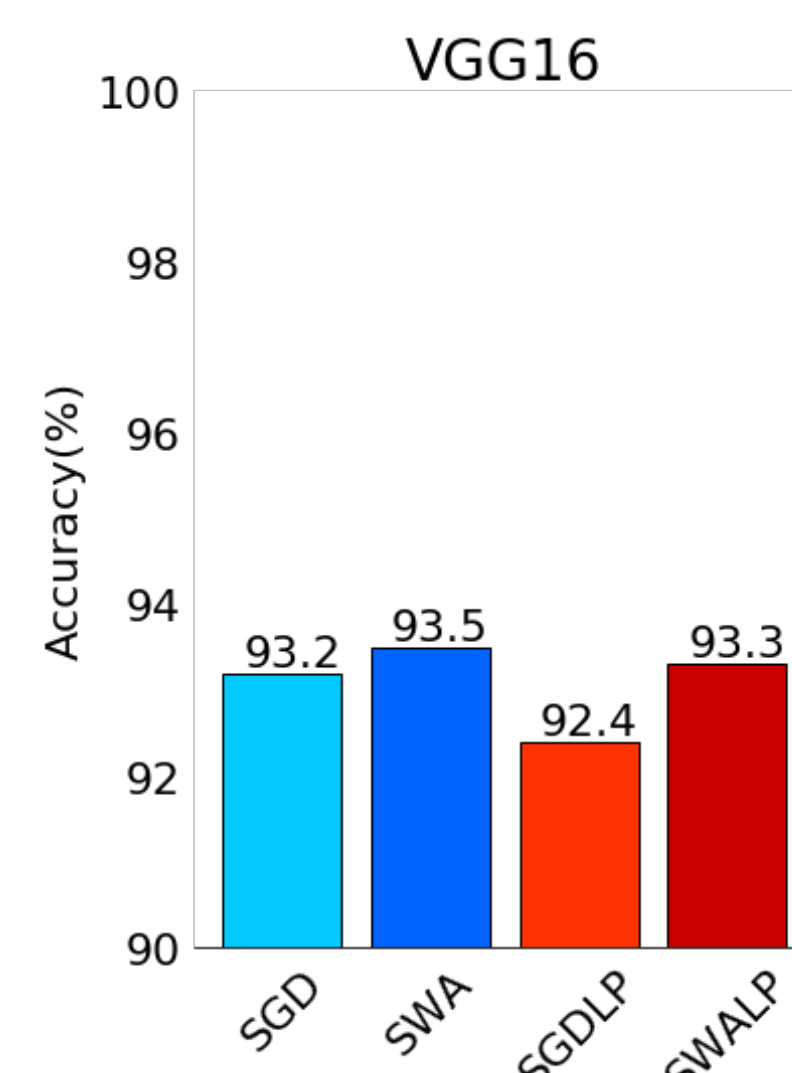The expected squared distance between the SWALP solution and the optimal one has a O(δ²).

- The best bound for low-precision SGD is O(δ) (Li et al, 2017).
- SWALP requires half of the number of bits to reduce the noise ball by the same factor.
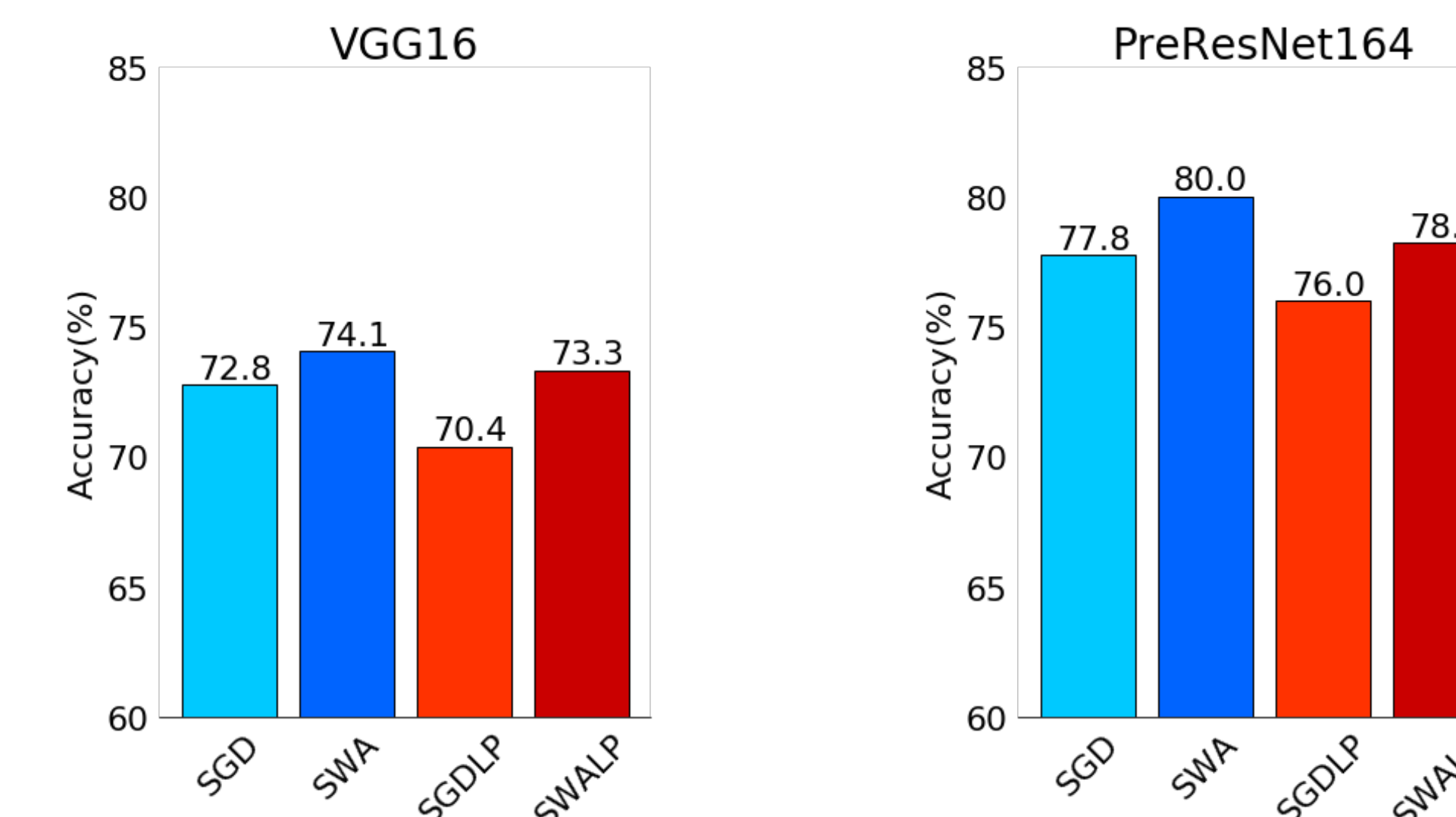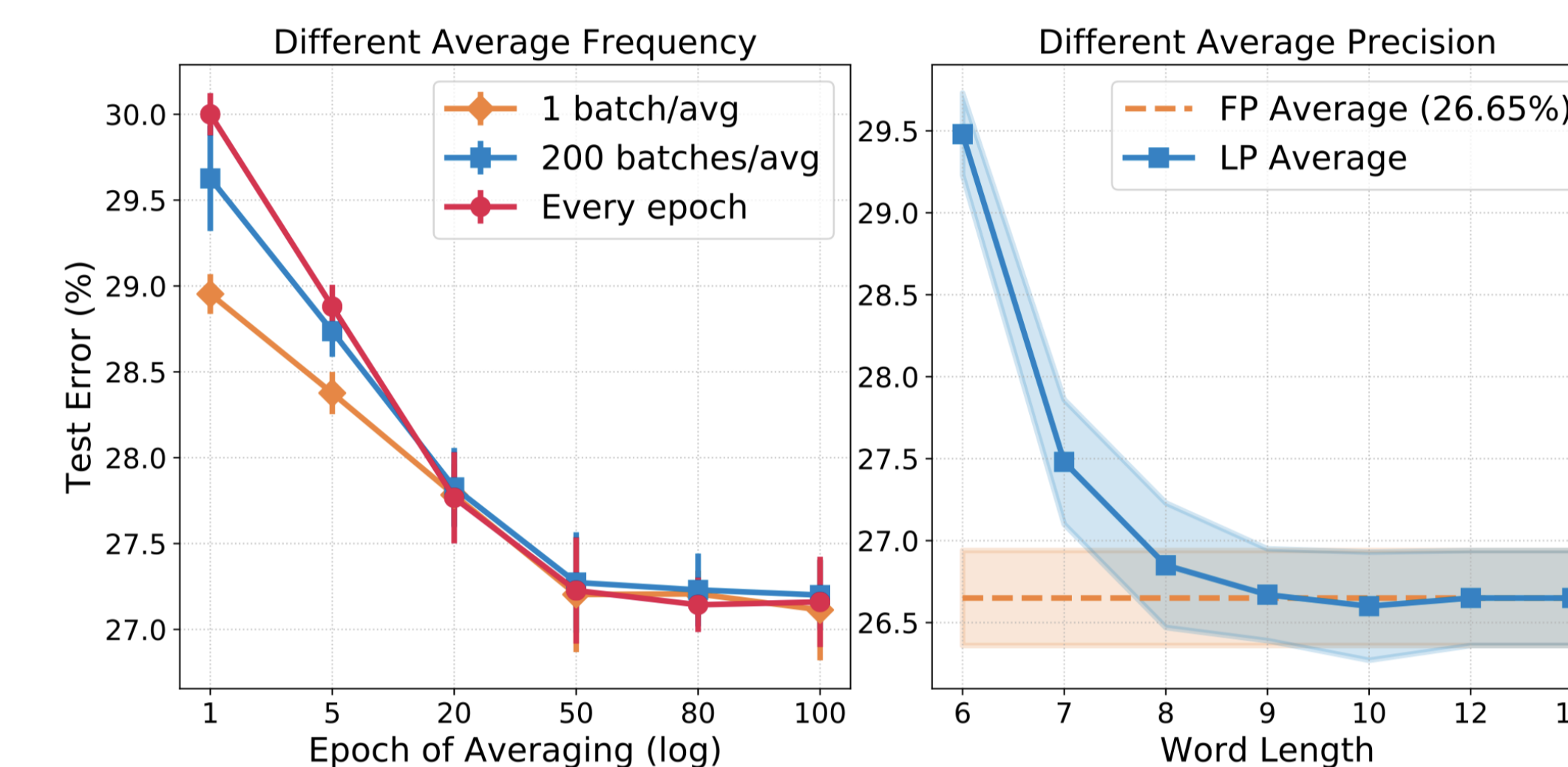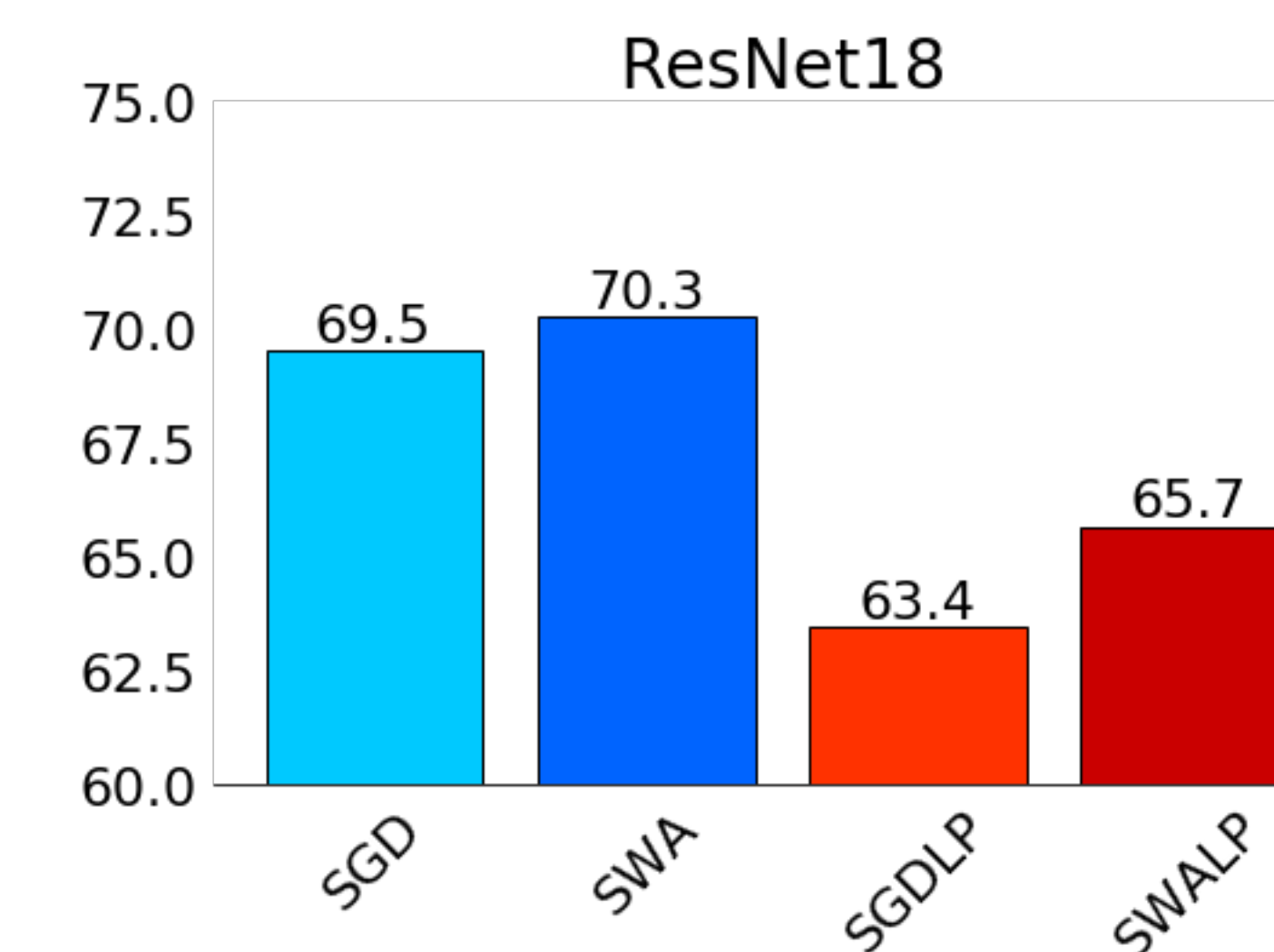
### Experimental Validation



Linear Regression — $\|w - w^*\|^2$ vs Iterations (x10³); legend: SGD-FL, SWA-FL, SGD-LP, SWALP, $Q(W^*)$

Logistic Regression — Gradient Norm vs Iterations (x10³); legend: SGD-FL, SWA-FL, SGD-LP, SWALP

## Experiments

### Results: CIFAR10



VGG16: SGD 93.2, SWA 93.5, SGDLP 92.4, SWALP 93.3

PreResNet164: SGD 95.4, SWA 96.0, SGDLP 94.2, SWALP 95.0

## Results: CIFAR100



VGG16: SGD 72.8, SWA 74.1, SGDLP 70.4, SWALP 73.3

PreResNet164: SGD 77.8, SWA 80.0, SGDLP 76.0, SWALP 78.2

### Averaging in Different Precision and Frequency



Different Average Frequency — Test Error (%) vs Epoch of Averaging (log); legend: 1 batch/avg, 200 batches/avg, Every epoch

Different Average Precision — vs Word Length; legend: FP Average (26.65%), LP Average

## Results: ImageNet



ResNet18: SGD 69.5, SWA 70.3, SGDLP 63.4, SWALP 65.7

## QPyTorch

We release QPyTorch, a low-precision arithmetic simulation package in PyTorch. A diverse range of quantization methods is supported with GPU acceleration.