# Accumulation Bit-Width Scaling for Ultra-Low Precision Training of Deep Neural Networks
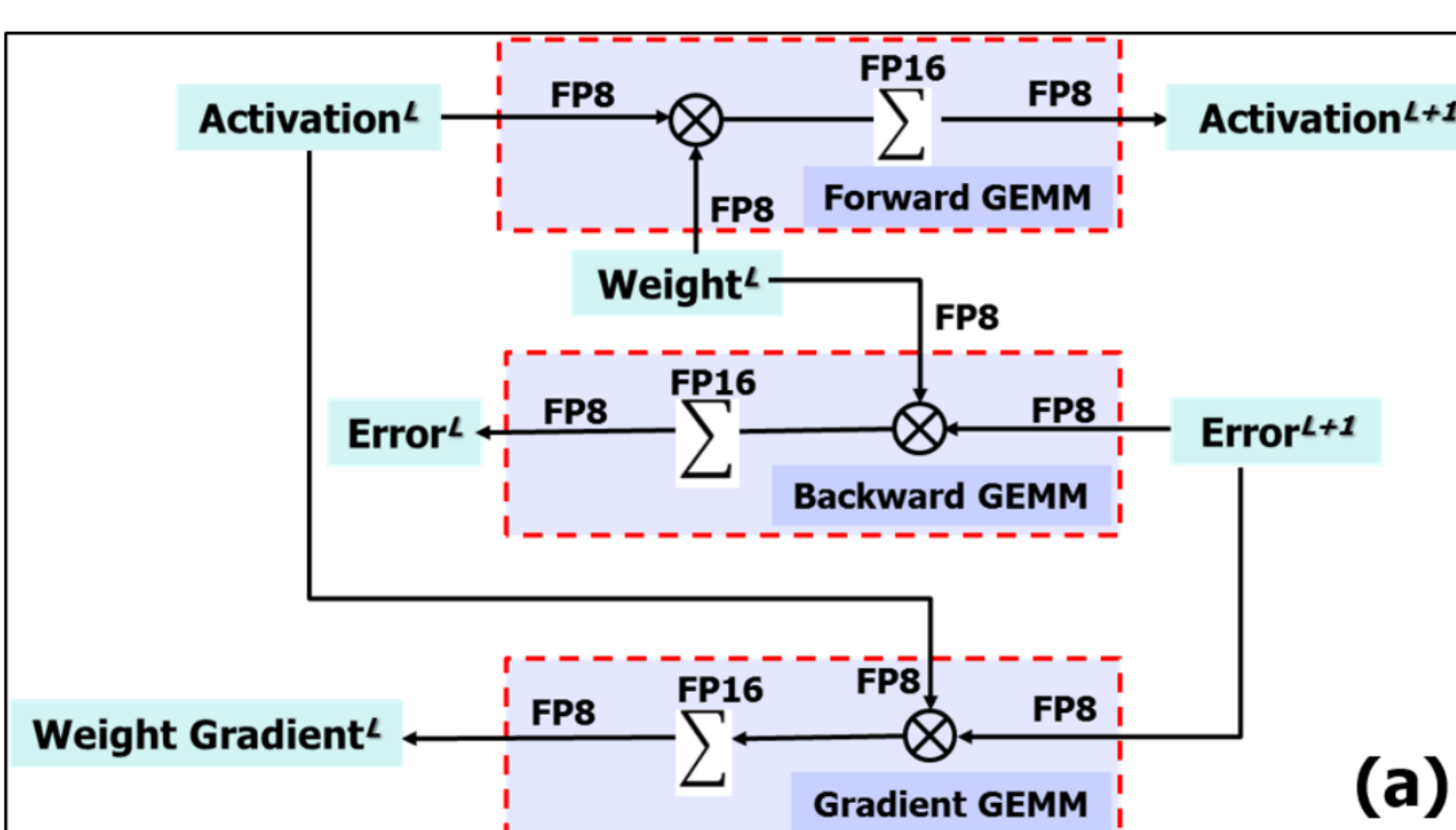
Charbel Sakr[1], Naigang Wang[2], Chia-Yu Chen[2], Jungwook Choi[2], Ankur Agrawal[2], Naresh Shanbhag[1], Kailash Gopalakrishnan[2]

[1]University of Illinois at Urbana-Champaign, [2]IBM T.J. Watson Research Center

## Reduced Precision FP Training

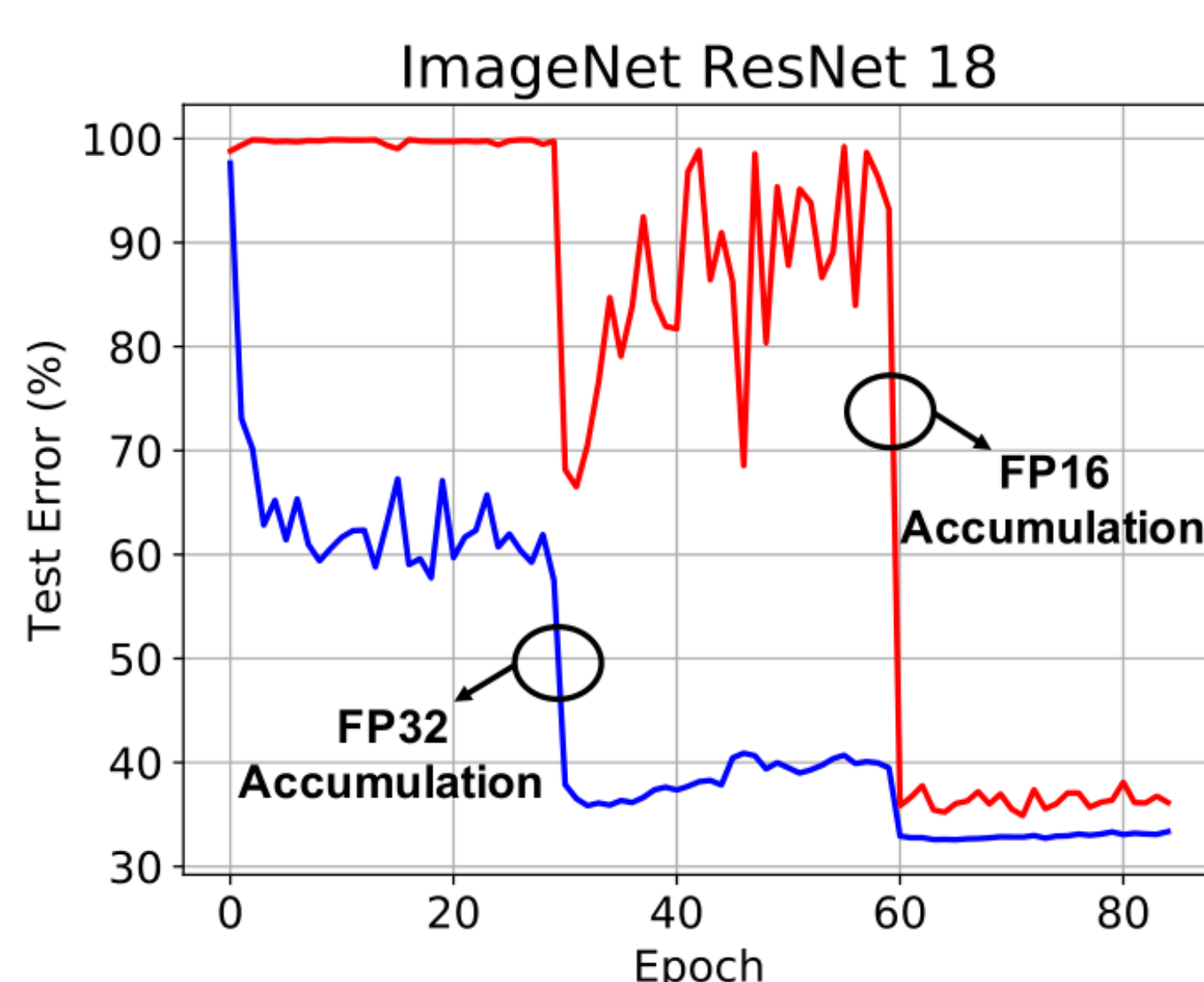### Representation & Accumulation Quantization



[Wang et al., NeurIPS 2018]

### Chunk-Based Accumulation

```
Input: {x_n}_{n=1:N}, {y_n}_{n=1:N} (FP_mult),
Parameter: chunk size CL
Output: sum (FP_acc)
sum = 0.0; idx = 0; num_ch = N/CL
for n=1:num_ch {
    sum_ch = 0.0
    for i=1:CL {
        idx++
        tmp = x_idx · y_idx  (in FP_mult)
        sum_ch += tmp  (in FP_acc))
    sum += sum_ch  (in FP_acc)}     (a)
```
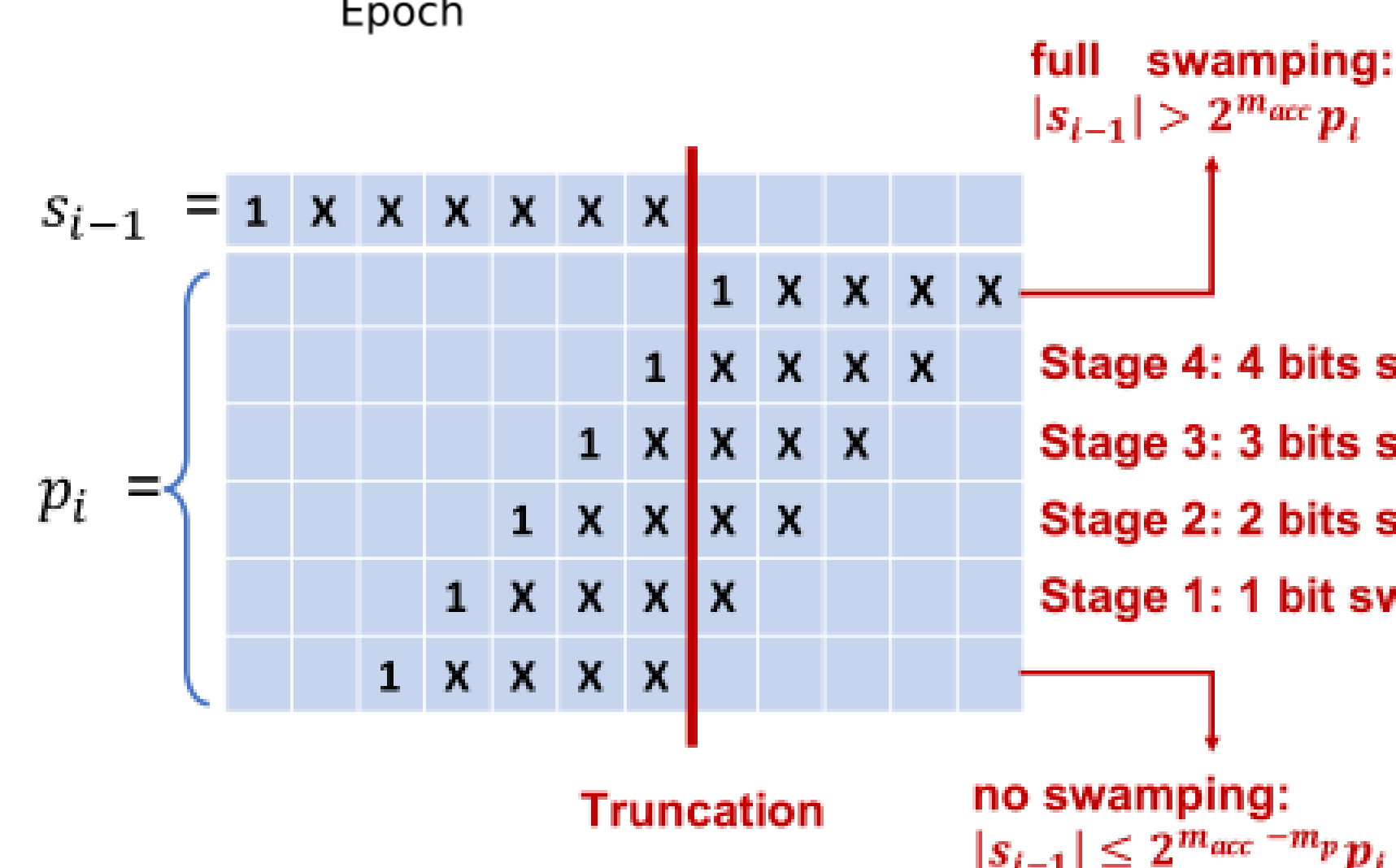
[Wang et al., NeurIPS 2018]

- tensors are quantized to FP8 = (1,5,2)
- accumulations are in FP16 = (1,6,9) but require chunking
- what is the accumulation precision required?

## Problem Setup

### ImageNet ResNet 18



- reducing representation precision in FP format is well studied [Wang et al., NeurIPS'2018]
- problem of reducing accumulation precision largely overlooked



full swamping: $|s_{i-1}| > 2^{m_{acc}} p_i$

Stage 4: 4 bits swamped
Stage 3: 3 bits swamped
Stage 2: 2 bits swamped
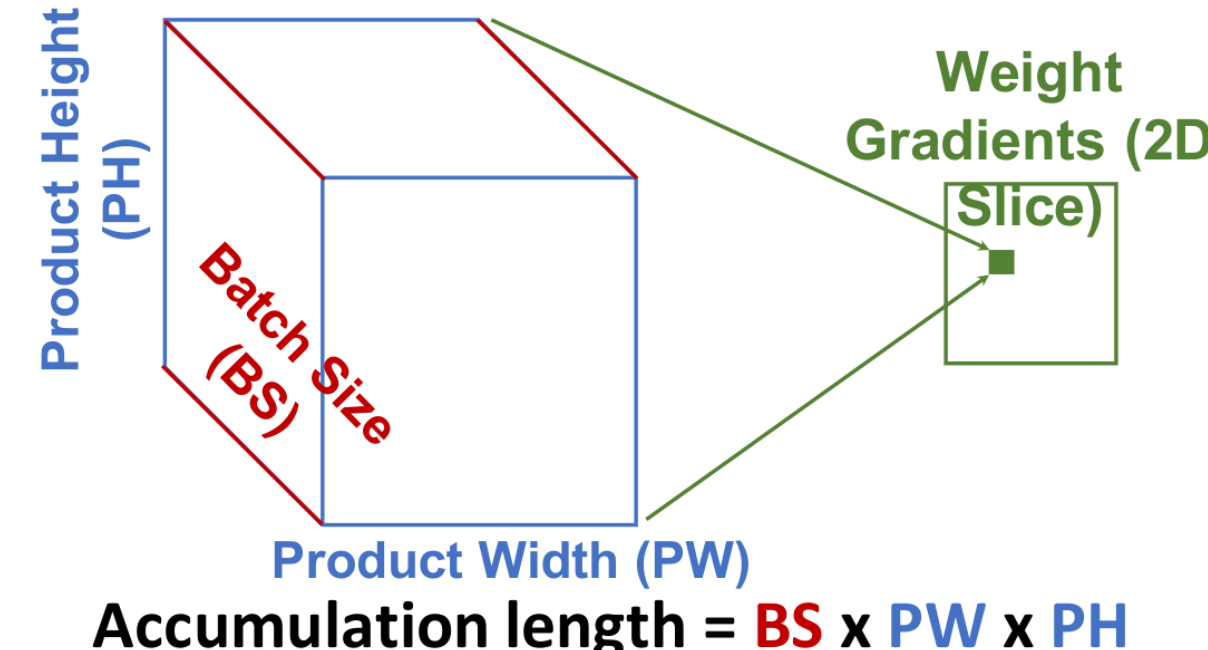Stage 1: 1 bit swamped

no swamping: $|s_{i-1}| \leq 2^{m_{acc}-m_p} p_i$

- accumulators typically designed conservatively because swamping effects are very destructive and intractable (hard to analyze)
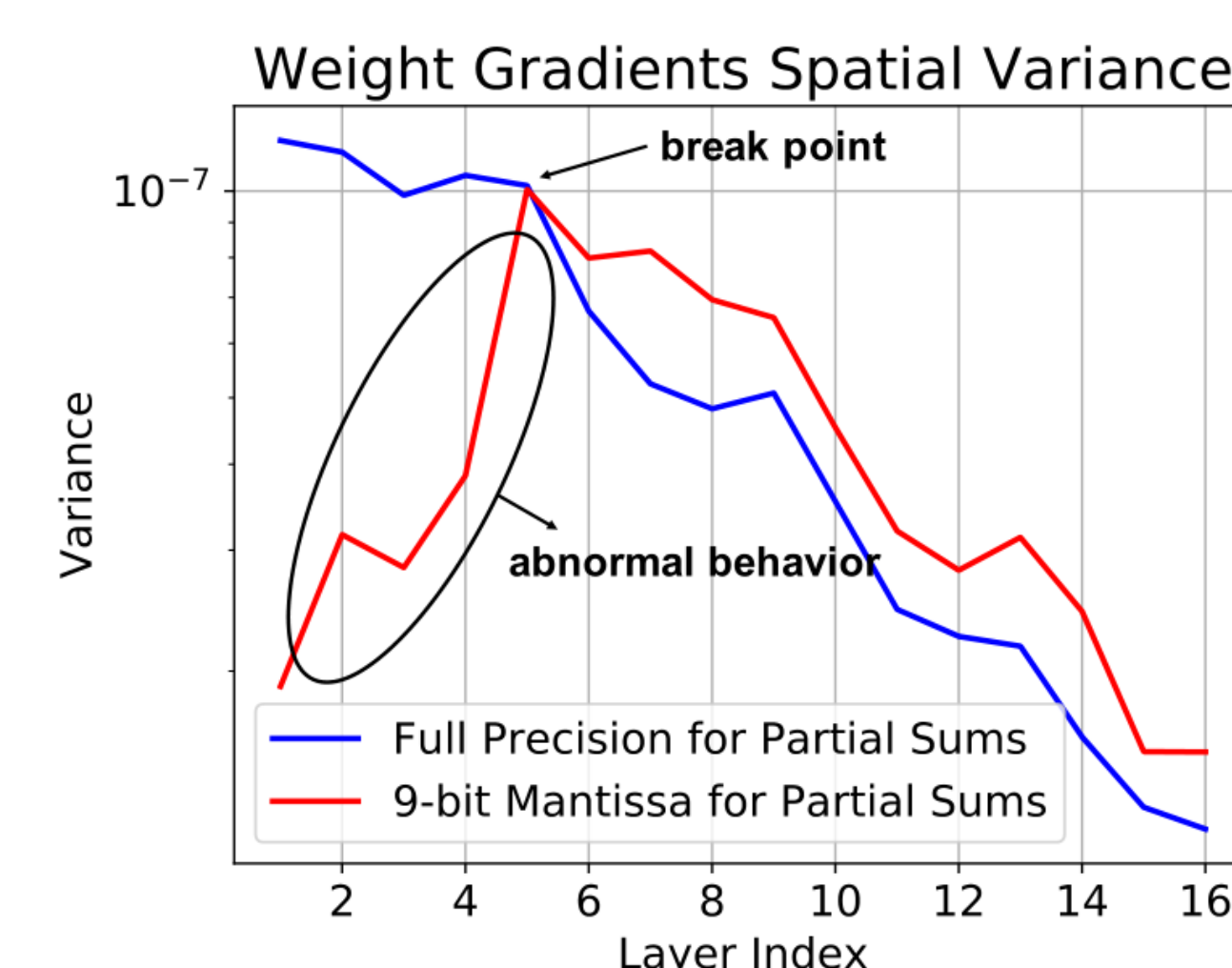
## Accumulation Variance

### Gradient Accumulation



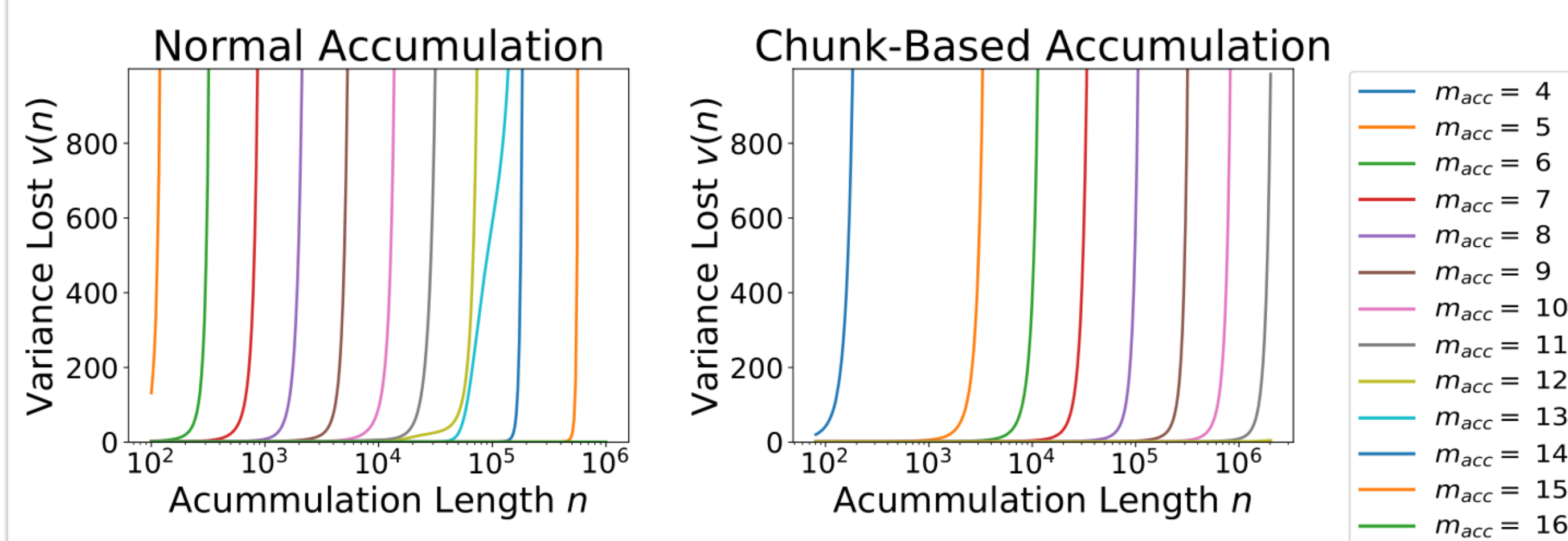Elementwise Product of Input Activations and Output Derivatives

Weight Gradients (2D Slice)

Product Height (PH)
Batch Size (BS)
Product Width (PW)

Accumulation length = BS x PW x PH

### Variance Lost

Weight Gradients Spatial Variance



break point
abnormal behavior
Full Precision for Partial Sums
9-bit Mantissa for Partial Sums

- analysis reveals correlation between accumulated variance and convergence behavior
- break point corresponds to increase in accumulation length
- hypothesis: there is a relationship between precision, accumulation variance, and accumulation length

## Variance Retention Ratio



Normal Accumulation

Chunk-Based Accumulation

$m_{acc} = 4$
$m_{acc} = 5$
$m_{acc} = 6$
$m_{acc} = 7$
$m_{acc} = 8$
$m_{acc} = 9$
$m_{acc} = 10$
$m_{acc} = 11$
$m_{acc} = 12$
$m_{acc} = 13$
$m_{acc} = 14$
$m_{acc} = 15$
$m_{acc} = 16$

$$VRR = \frac{\sum_{i=2}^{n-1}(i-\alpha)_+ q_i \mathbf{1}_{\{i>\alpha\}} + \sum_{j_r=2}^{m_p}(n-\alpha_{j_r})_+ q'_i \mathbf{1}_{\{n>\alpha_{j_r}\}} + nk_3}{kn}$$

where $(x)_+ = \begin{cases} x \text{ if } x > 0 \\ 0 \text{ otherwise} \end{cases}$, $\mathbf{1}_A = \begin{cases} 1 \text{ if } A \text{ is true} \\ 0 \text{ otherwise} \end{cases}$,

$\alpha = \frac{2^{m_{acc}-3m_p}}{3}\sum_{j=1}^{m_p} 2^j(2^j-1)(2^{j+1}-1)$, $q_i = 2Q\left(\frac{2^{m_{acc}}}{\sqrt{i}}\right)\left(1-2Q\left(\frac{2^{m_{acc}}}{\sqrt{i-1}}\right)\right)$,
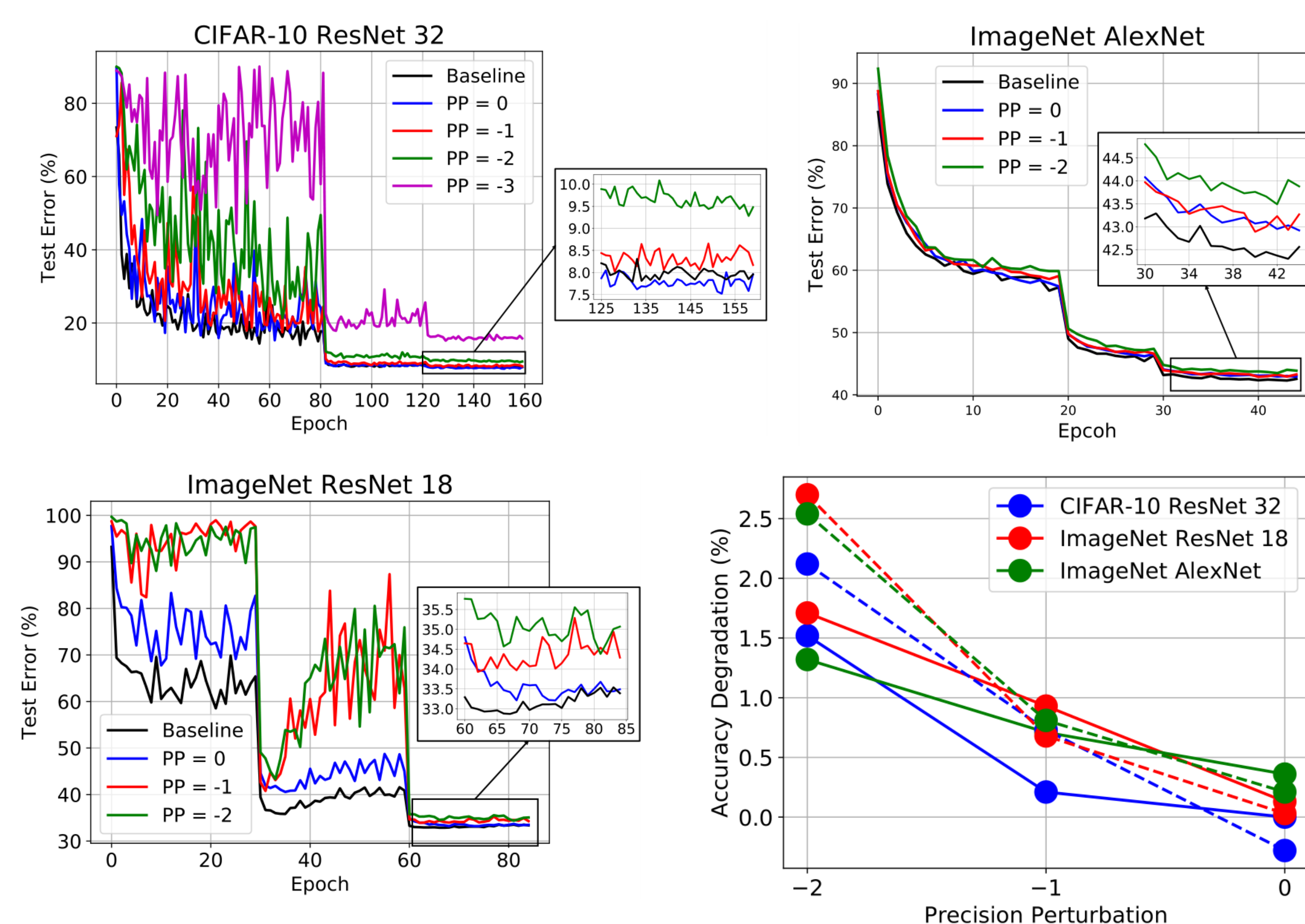
$\alpha_{j_r} = \frac{2^{m_{acc}-3m_p}}{3}\sum_{j=1}^{j_r-1} 2^j(2^j-1)(2^{j+1}-1)$,

$q'_{j_r} = N_{j_r-1}2Q\left(\frac{2^{m_{acc}-m_p+j_r-1}}{\sqrt{n}}\right)\left(1-2Q\left(\frac{2^{m_{acc}-m_p+j_r}}{\sqrt{n}}\right)\right), k = k_1 + k_2 + k_3,$

$k_1 = \sum_{i=2}^{n-1} q_i \mathbf{1}_{\{i>\alpha\}}$, $k_2 = \sum_{j_r=2}^{m_p} q'_i \mathbf{1}_{\{n>\alpha_{j_r}\}}$, and $k_3 = 1 - 2Q\left(\frac{2^{m_{acc}-m_p+1}}{\sqrt{n}}\right)$.
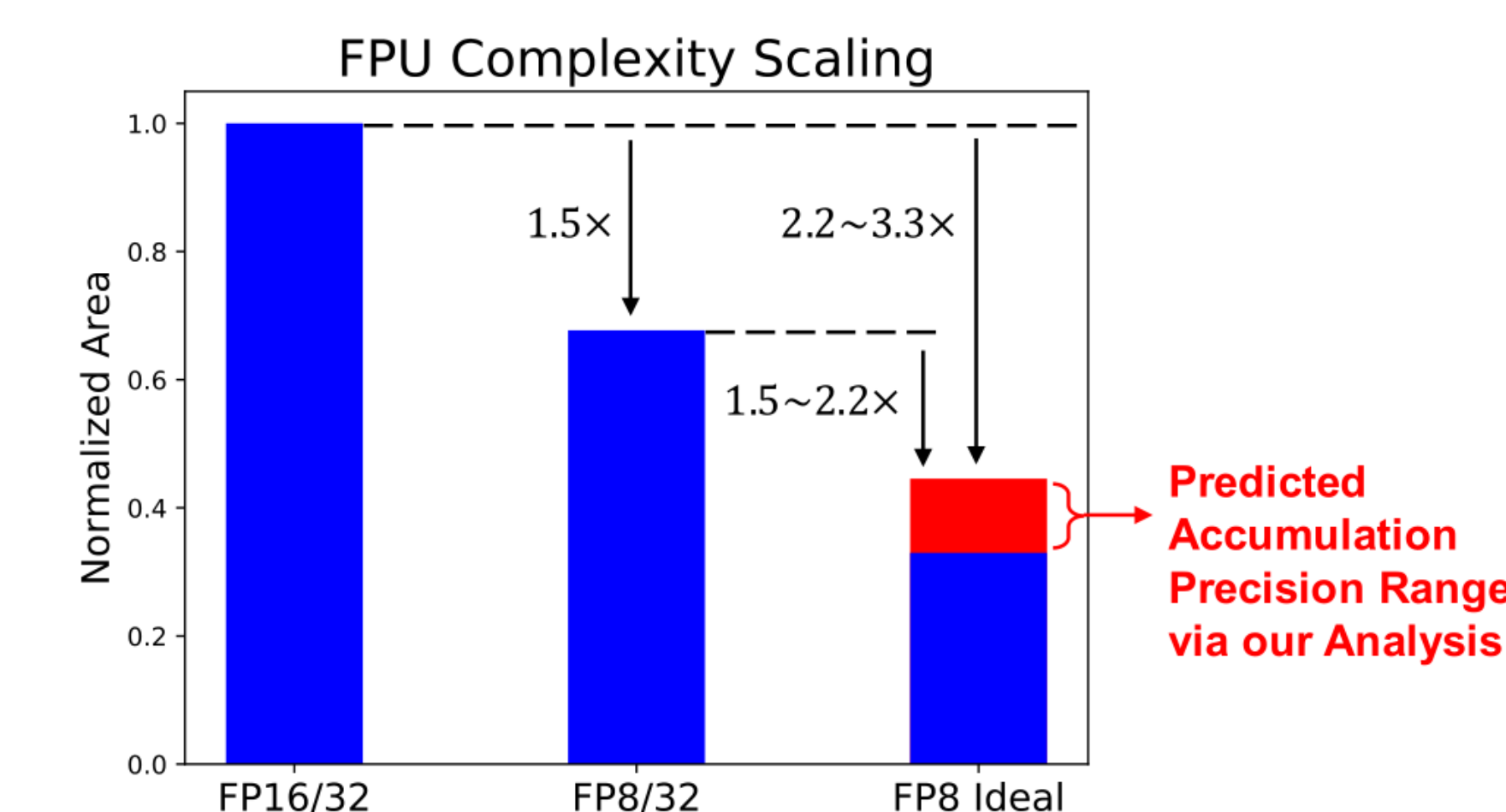
- swamping causes loss of accumulated variance in partial sums
- variance retention ratio (VRR) is derived analytically as a function of precision and accumulation length
- VRR 'knee' corresponds to the maximum accumulation length allowed for a given precision

## Convergence with Low-Precision Accumulation



CIFAR-10 ResNet 32
Baseline
PP = 0
PP = -1
PP = -2
PP = -3

ImageNet AlexNet
Baseline
PP = 0
PP = -1
PP = -2

ImageNet ResNet 18
Baseline
PP = 0
PP = -1
PP = -2

CIFAR-10 ResNet 32
ImageNet ResNet 18
ImageNet AlexNet

- VRR-based analysis enables convergence with low-precision accumulation and is tight

## Hardware Benefits



FPU Complexity Scaling

1.5×
2.2~3.3×
1.5~2.2×

Predicted Accumulation Precision Range via our Analysis

FP16/32    FP8/32    FP8 Ideal

- low-precision accumulation reduces hardware cost over by ~2× compared to representation quantization

## Acknowledgement