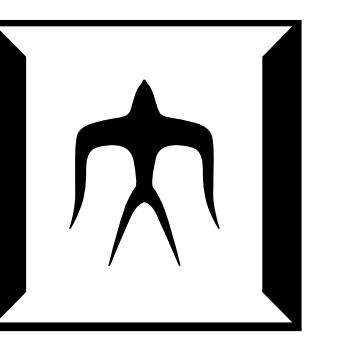


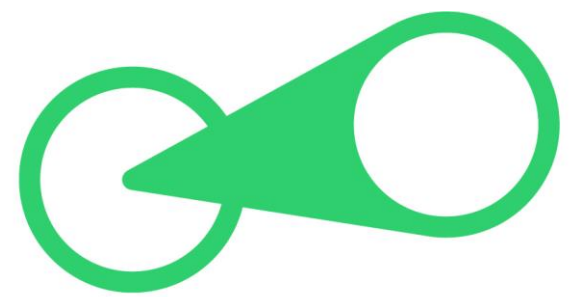
# Learning Non-Uniform Step-Sizes for Neural Network Quantization



東京工業大学  
Tokyo Institute of Technology

J.R. Liang<sup>1</sup>, S. Gongyo<sup>2</sup>, M. Ambai<sup>2</sup>, R. Kawakami<sup>1</sup>, I. Sato<sup>1,2</sup>

<sup>1</sup>Tokyo Institute of Technology, <sup>2</sup>Denso IT Laboratory, Inc.



Recognition and Learning  
Algorithm Laboratory

## Background & Contribution

### Background

- Trends in increasing DNN model size.
- Industrial applications often demand:
  - Real-time DNN inference
  - Use of low-end device

### Contribution

- We propose a novel non-uniform LSQ quantizer (nuLSQ) for DNN compression.
- nuLSQ outperforms LSQ on CIFAR-10 and -100.

## Existing Method

**Learned Step-Size Quantization (LSQ)** [Esser, Steven K.+., ICLR 2020]

### Forward Pass

- Uniform quantization process in activations:

$$Q_{LSQ}(x, s) = \sum_{n=1}^N s \sigma\left(x - ns + \frac{s}{2}\right)$$

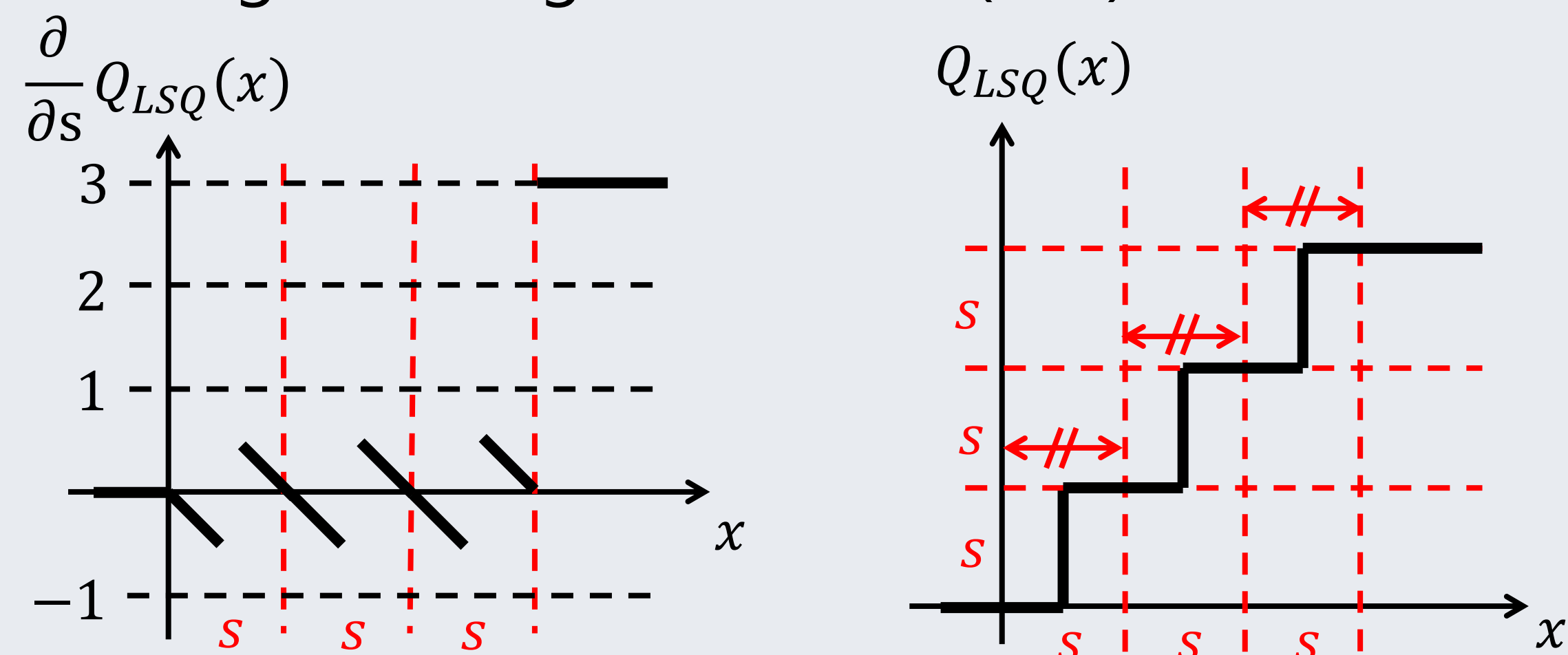
The number of step-size  $N$

Uniform step size  $s$

Unit step function  $\sigma(\cdot)$

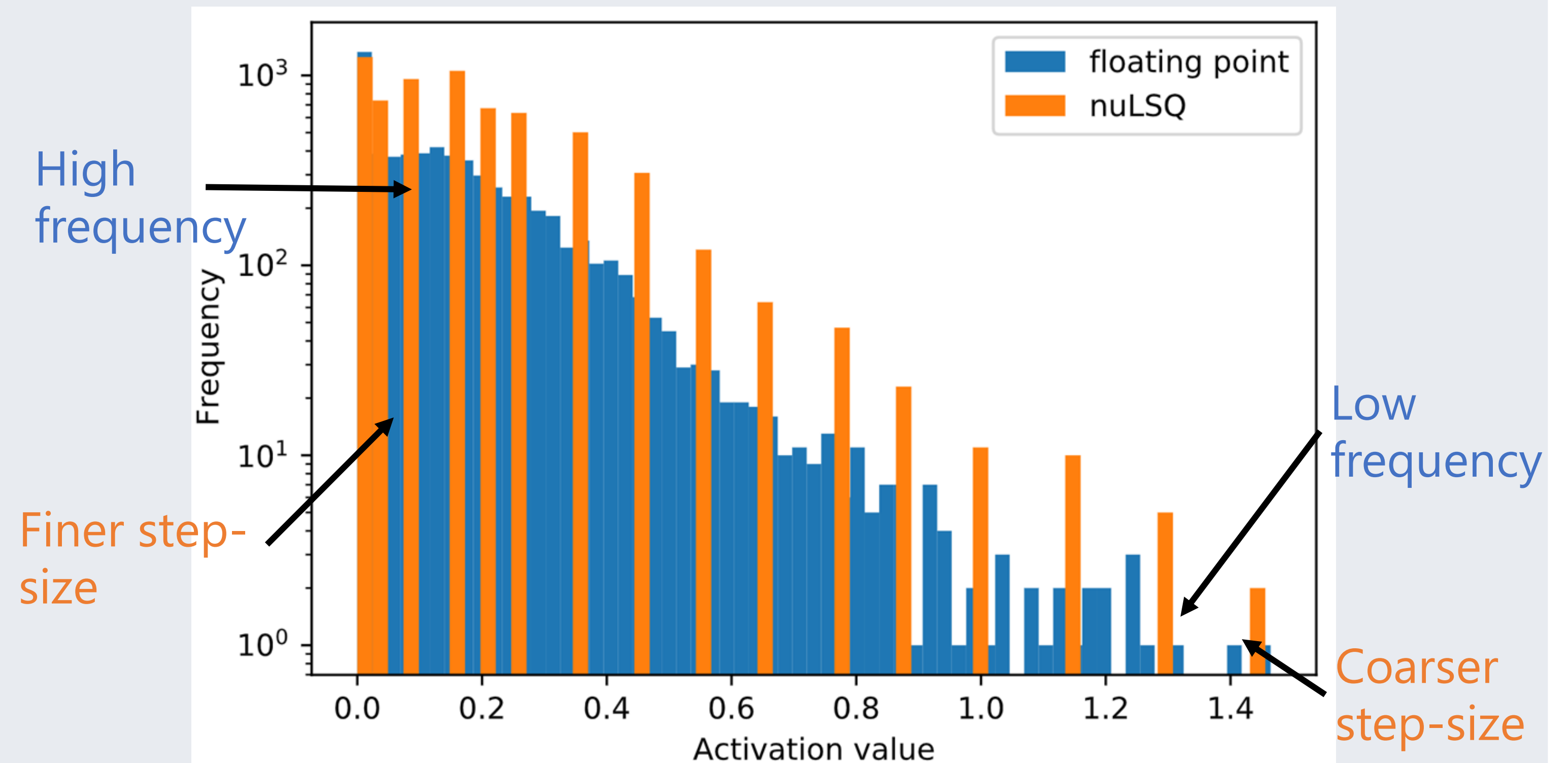
### Backward Pass

- Uniform step-size gradient approximated with straight through estimator (STE).



## Proposed Method

**non-uniform Learned Step-Size Quantization (nuLSQ)**



☺ nuLSQ can better approximate real activations by introducing non-uniform quantization.

### Forward Pass

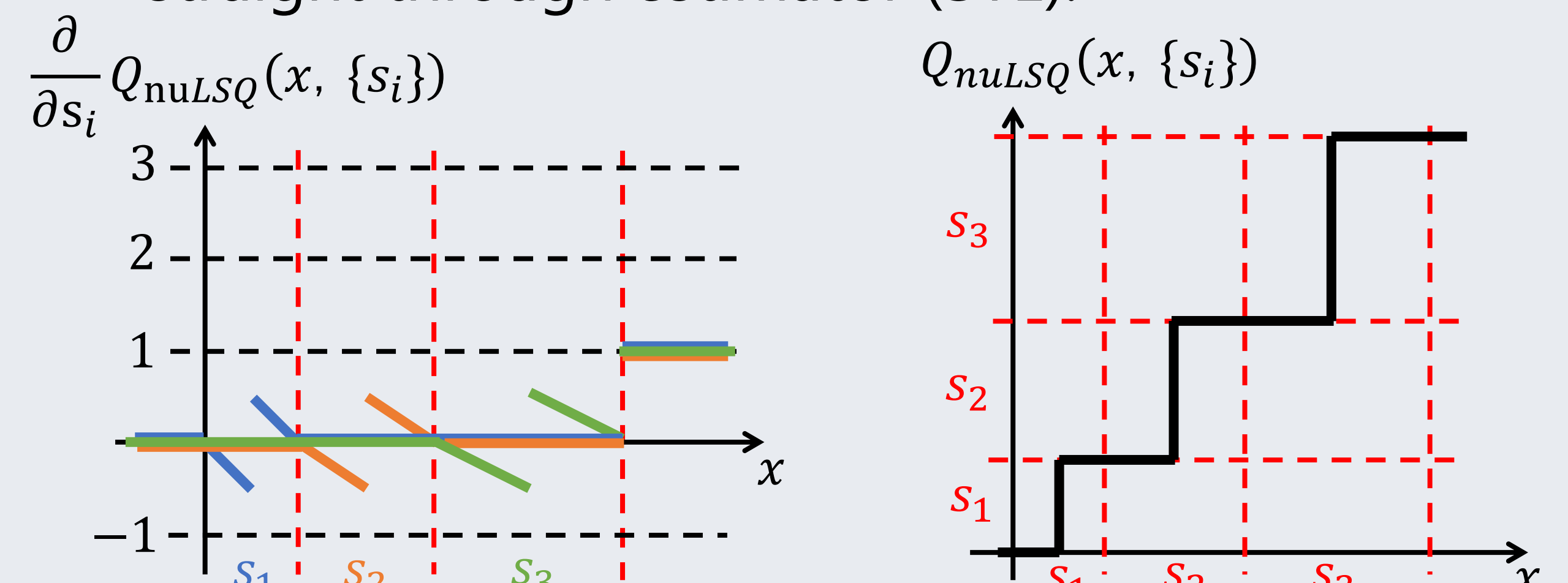
- Non-uniform quantization process in activations:

$$Q_{nuLSQ}(x, \{s_i\}) = \sum_{n=1}^N s_n \sigma\left(x - \left(\sum_{m=1}^{n-1} s_m + \frac{s_n}{2}\right)\right)$$

Non-uniform step sizes  $\{s_i\}$

### Backward Pass

- Non-uniform step-sizes gradient approximated with straight through estimator (STE).



## Exp-1: Performance

### Settings

- Quantized weight and activation.
- Measured mean test accuracy over the last 10 training epochs.

### Results

- ☺ nuLSQ outperforms LSQ under 2-, 3- and 4-bit quantization.

#### Test accuracy of ResNet-20 on CIFAR10

	2-bit	3-bit	4-bit	Float
LSQ	84.5%	88.0%	88.7%	89.0%
nuLSQ (ours)	85.2%	88.2%	88.9%	

#### Test accuracy of ResNet-56 on CIFAR100

	2-bit	3-bit	4-bit	Float
LSQ	63.4%	65.6%	65.7%	66.4%
nuLSQ (ours)	64.1%	65.7%	66.8%	

## Evaluation

### Exp-2: Information Entropy

### Settings

- Measured Shannon entropy of the quantized activation outputs (the first and last layers are omitted).

### Results

- ☺ nuLSQ demonstrates an overall **18%** information gain over the uniform LSQ.
  - nuLSQ has a more uniformly distributed patterns.

