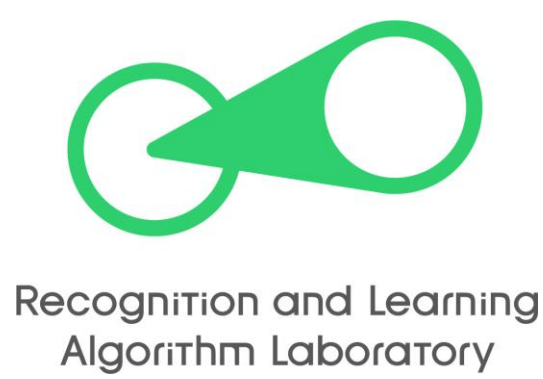


# Learning Non-Uniform Step-Sizes for Neural Network Quantization

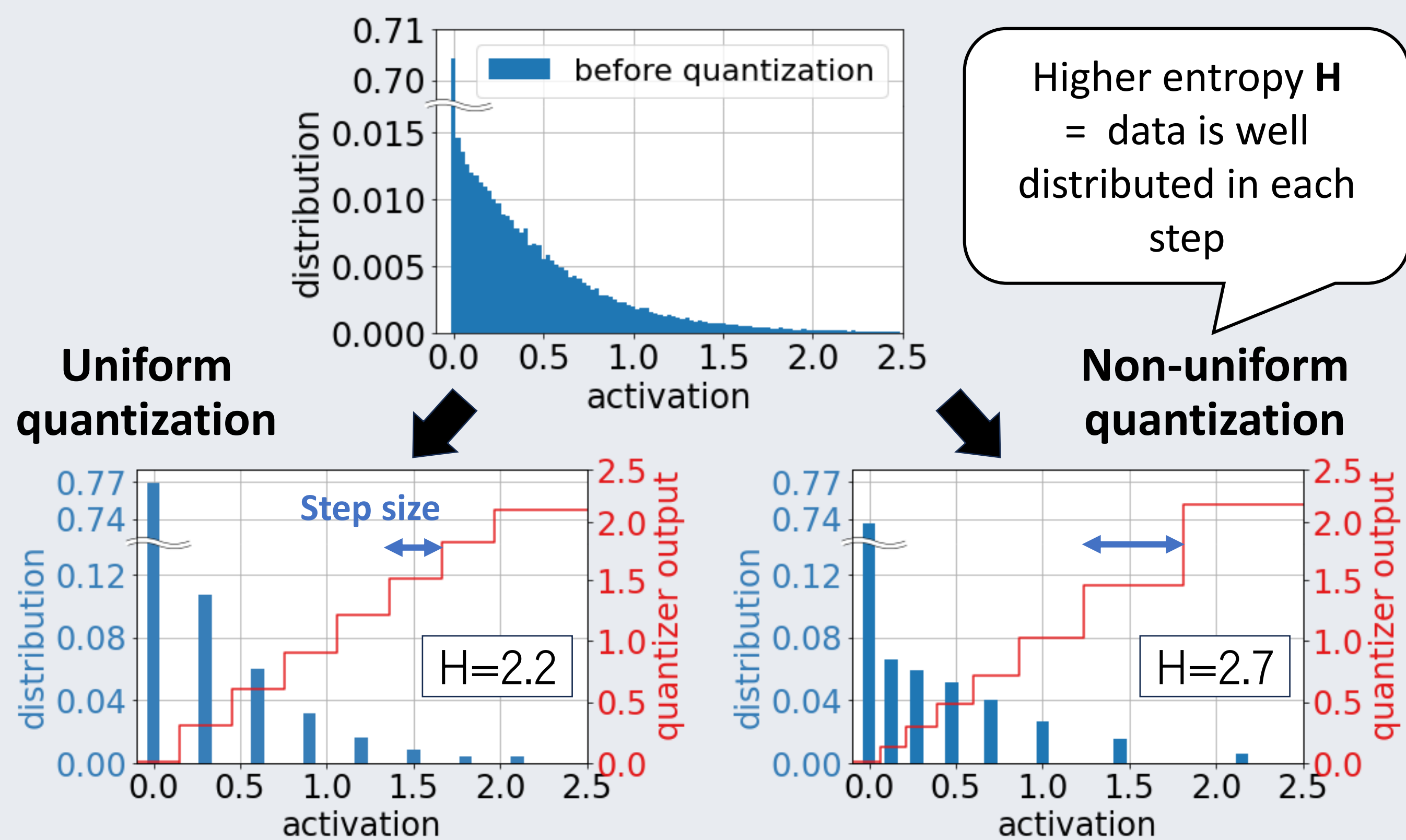
S. Gongyo<sup>1</sup>, J.R. Liang<sup>2</sup>, M. Ambai<sup>1</sup>, R. Kawakami<sup>2</sup>, I. Sato<sup>1,2</sup>

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

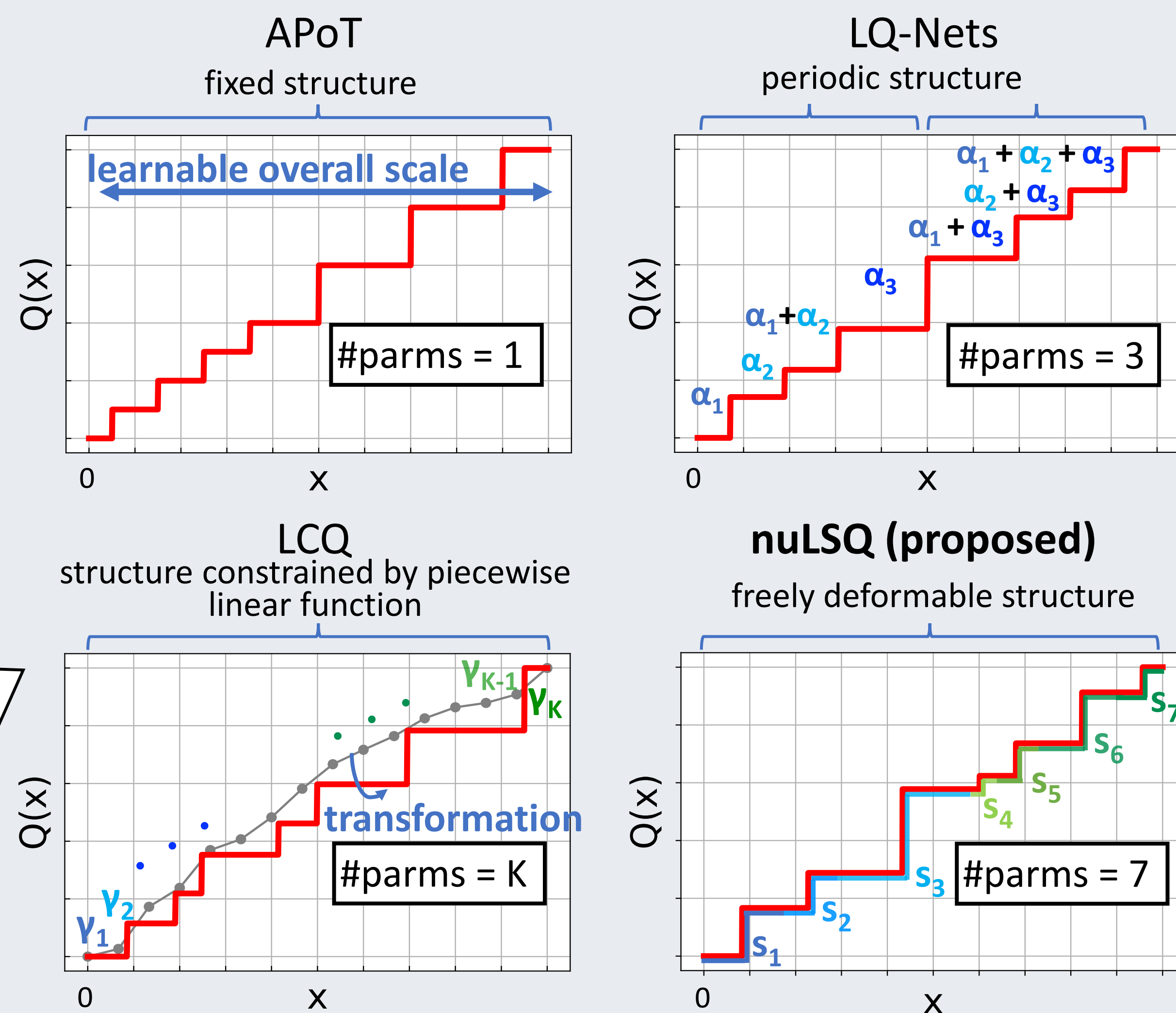


## Introduction

- Quantization typically assumes a uniform step size
- Use of **non-uniform** step size can better capture the distribution



nuLSQ has high flexibility among non-uniform methods



## Learned Step-Size Quantization (LSQ)

[Esser, Steven K.+. ICLR 2020]

### Forward Pass

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

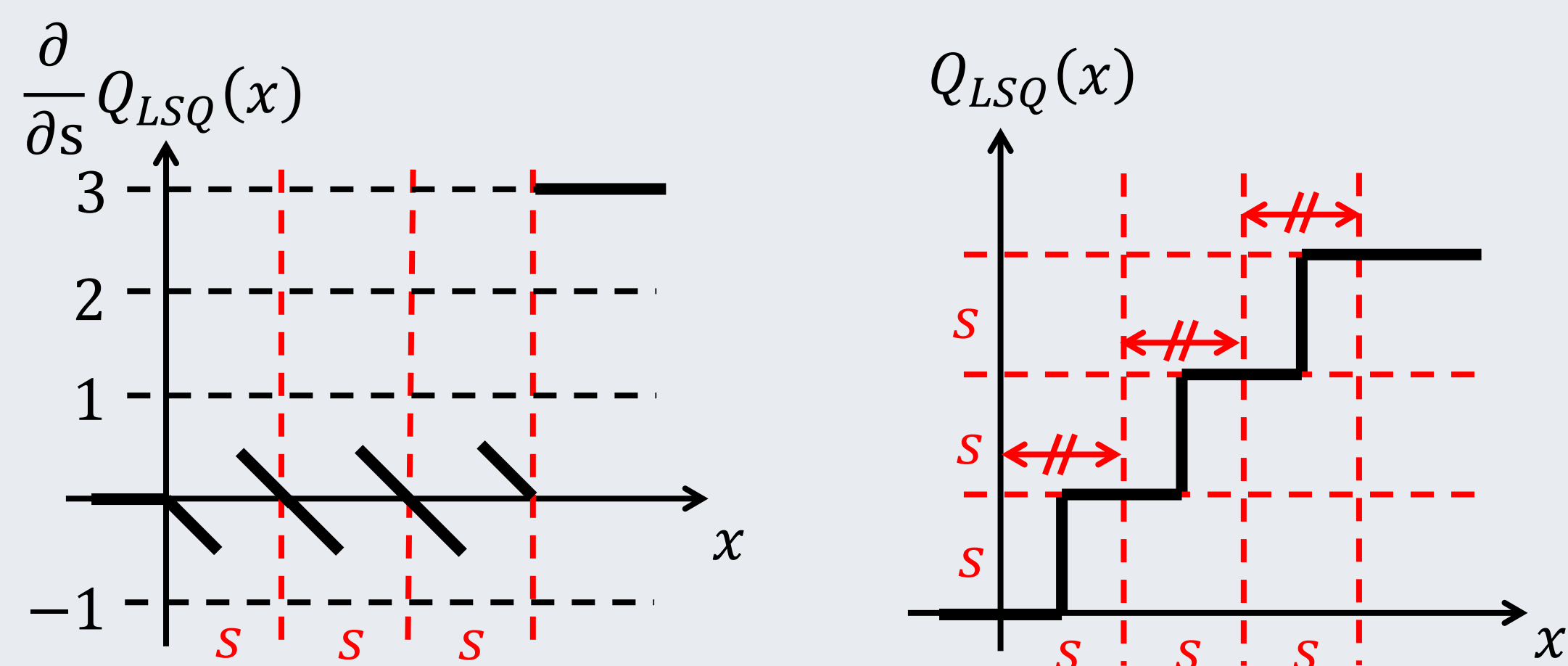
Uniform step size  $s$

The number of step-size  $N$

Unit step function  $\sigma(\cdot)$

### Backward Pass

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



## Proposed Method

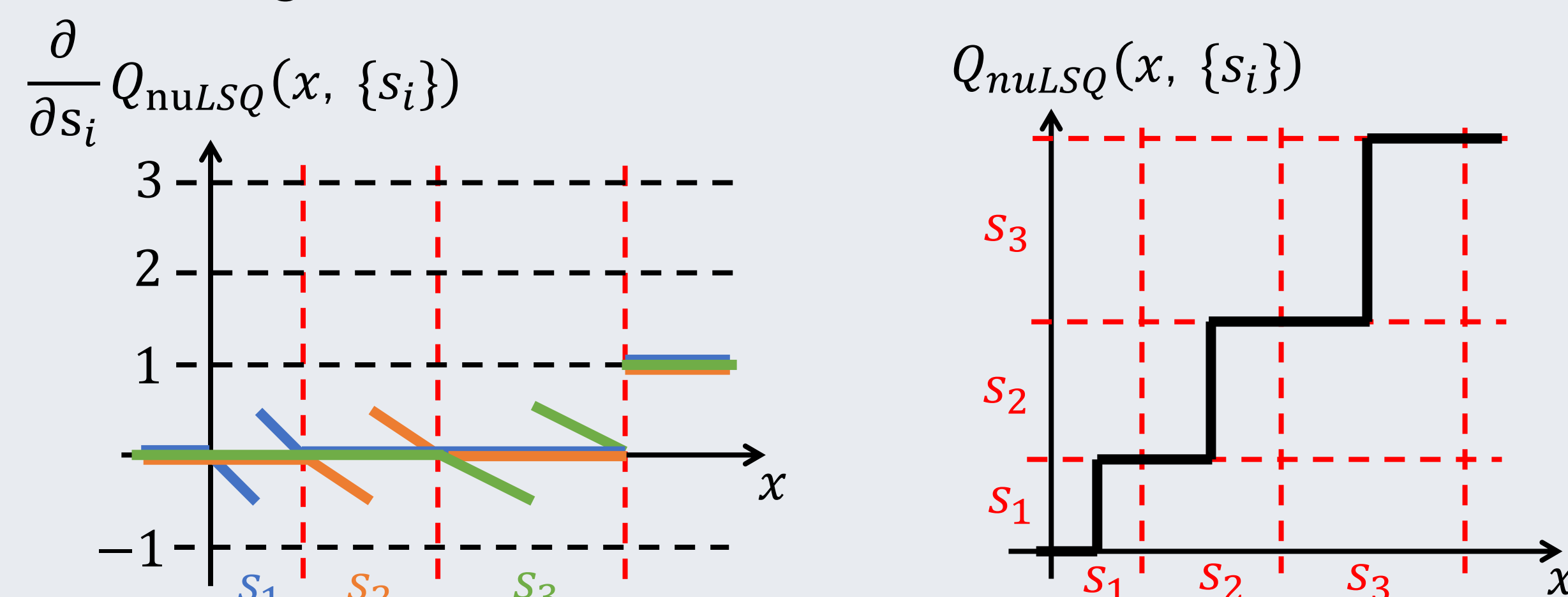
### Forward Pass

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

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

### Backward Pass

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



## Experiments

### ImageNet top-1 accuracy (%) of MobileNetV2 (FP:72.9)

Bits (WA)	LSQ	LSQ +BR	LSQ +Dp/Fz	UniQ	LCQ	LLSQ	EWGS	nuLSQA
2/2	46.7	50.6	-	50.5	-	-	-	<b>58.4</b>
3/3	65.3	67.4	67.8	65.0	-	-	-	<b>67.9</b>
4/4	69.5	70.4	70.6	68.2	70.8	67.4	70.3	<b>71.1</b>

### ImageNet top-1 accuracy (%) comparison

Methods	Swin-T (FP:81.2)			ConvNext (FP:81.87)
	2/2	3/3	4/4	3/3
LSQ*	74.58	77.48	78.33	72.9
nuLSQ-WA	<b>74.91</b>	<b>77.71</b>	<b>78.37</b>	<b>73.39</b>

### ImageNet top-1 accuracy (%) of ResNet-18 (FP:69.76)

PACT*	DoReFa*	LSQ*	APoT†	LQ-Nets*	LCQ*	nuLSQA
62.48	63.28	64.51	64.41	63.71	64.67	<b>64.89</b>

### Comparison to QAT methods with distillation

Network	Methods	Top-1 Acc (%)	Top-5 Acc (%)
MobileNet V2	QKD	67.4	87.0
	PROFIT	71.56	90.40
	nuLSQ-A	<b>71.89</b>	<b>90.44</b>

### CIFAR100 top-1 accuracy (%) comparison to LSQ

Network	Methods	Bits (W/A)		
		2/2	3/3	4/4
ResNet-20 (FP:69.8)	LSQ*	65.82	68.60	<b>69.56</b>
	nuLSQ-A	<b>66.02</b>	68.58	69.42
	nuLSQ-W	66.00	<b>68.70</b>	69.40
ResNet-56 (FP:74.9)	LSQ*	70.50	72.62	<b>73.48</b>
	nuLSQ-A	70.66	<b>72.98</b>	<b>73.48</b>
	nuLSQ-W	<b>70.82</b>	72.82	73.42

### Comparison of Shannon Entropy

