
A High Performance and Low Latency Deep Spiking Neural Networks Conversion Framework

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Spiking Neural Networks (SNN) are promised to be energy-efficient and achieve
2 Artificial Neural Networks (ANN) comparable performance through conversion
3 processes. However, a converted SNN relies on large timesteps to compensate
4 for conversion errors, which as a result compromises its efficiency in practice.
5 In this paper, we propose a novel framework to convert an ANN to its SNN
6 counterpart losslessly with minimal timesteps. By studying the errors introduced
7 by the whole conversion process, an overlooked inference error is revealed besides
8 the coding error occurred during converting. Inspired by the quantization aware
9 training, a QReLU activation is introduced during training to eliminate the coding
10 error theoretically. Furthermore, a buffered non-leaky-integrate-and-fire neuron
11 that utilizes the same basic operations as in conventional neurons is designed to
12 reduce the inference error. Experiments on classification and detection tasks show
13 that our proposed method attains ANNs level performance using only 16 timesteps.
14 To the best of our knowledge, it is the first time converted SNNs with low latency
15 demonstrate their capability to achieve high performance on nontrivial vision tasks.
16 Source code will be released later.

17 1 Introduction

18 Recently, Spiking Neural Networks (SNN) have attracted a great deal of researchers' attention
19 due to their essential advantages, such as spatio-temporal information processing capability and
20 energy-efficient with low power consumption. Unlike conventional Artificial Neural Networks (ANN)
21 communicating with every neuron between layers during every inference step, SNN only passes
22 sparse events when particular neurons are fired. This compute-on-demand property of SNN is
23 very friendly to low-power hardware, and suitable for neuromorphic platforms, such as event-based
24 camera [11] and brain-inspired chips [32]. Such dedicated bionic systems can significantly reduce
25 memory usage and energy consumption, which is very appealing to a wide range of applications like
26 autonomous mobile robots.

27 Despite the promising characteristics, it is still a challenging problem to train high-performance deep
28 SNNs for practical tasks. Directly applying the backpropagation approach is infeasible because of the
29 non-differentiable nature of the spike activity. Taking advantage of the ANN's success, conversion
30 methods have been proven to be a very promising direction in achieving high performance SNN. Cao
31 et al. [4] first convert ANNs to SNNs by demonstrating the relationship between spiking neuron and
32 ReLU function for rate-based methods. Diehl et al. [9] propose a weight normalization approach
33 to mapping weights from ANN to SNN. Later, Rueckauer et al. [39] give a theoretical analysis on
34 ANN-SNN conversion. A reset-by-subtraction IF neuron is proposed to better handle the accuracy
35 degradation during the conversion. They also propose the weight normalization algorithm that uses
36 percentile function instead of max value to achieve better performance. Kim et al. [23] propose a
37 channel-wise normalization to further reduce the performance gap between ANNs and SNNs .

38 Although great progress has been made, those converted SNNs suffer from *accuracy-latency trade-*
39 *off* [9, 31]. As a result, the converted SNNs need longer time to achieve comparable precision
40 with their counterpart ANN during inference and compromises the efficiency in practice. To deal
41 with this problem, hybrid methods [37, 36] are proposed to further finetune a converted SNN with
42 approximated gradients to decrease the inference time. On the other hand, for the theoretical training
43 gaps between ANNs and SNNs, Ding et al. [10] propose Rate Norm to better choose the scale factor
44 and introduce extra loss to decrease the latency. Deng and Gu [8] utilize training samples to transfer
45 the weights to the target SNN by combining threshold balance and soft-reset mechanisms. However,
46 these works only focus on reducing the errors introduced by mapping feature values to spike trains.

47 With an in-depth analysis of the entire conversion process, we recognize two types of conversion
48 errors, namely coding error and inference error, that degrade the performance of SNN. As other
49 conversion approaches mainly focus on the coding errors, our motivation is to reduce both the coding
50 and inference errors. Inspired by the analysis in Rueckauer et al. [39], we notice that the SNN
51 with step function resembles quantized ANNs. By establishing a strict equivalency between spiking
52 neuron and quantized activation function, the coding error can be eliminated completely under limited
53 timesteps. Further, we observe that the assumption in rate-based conversion is inconsistent with
54 the actual neuron behavior and the SNN’s dynamics during the inference. Most of the previous
55 conversion-related literature overlook this problem, which limits the performance of the algorithms
56 and leads to longer simulation time to achieve comparable results of the original ANN.

57 In this paper, we propose a novel High Performance Conversion (HPC) method to drastically close
58 the gap between ANNs and SNNs. Specifically, a QReLU activation is proposed to eliminate the
59 coding error. Since QReLU is applied to the ANN training phase, no extra effort is required during
60 the conversion. Besides, a novel buffered non-leaky IF neuron is designed to compensate for the
61 inference error and retain SNN’s efficiency during inference. Utilizing a new form of activation
62 function and novel efficient spiking neuron, SNN obtained by our proposed approach achieves high
63 performance and low latency simultaneously. With only 16 timesteps, HPC achieves state-of-the-art
64 performance across classification and detection tasks. Our major contributions are summarized as:

- 65 • An overlooked inference error is discovered. The conversion errors including both coding
66 errors and inference errors are theoretically analyzed to clarify the reason of performance
67 degradation during conversion.
- 68 • A novel high performance conversion method is proposed to deal with the two types of
69 conversion errors. The coding errors are eliminated completely in theory and the inference
70 error problem is relieved to a great extent.
- 71 • Comprehensive experiments on different tasks demonstrate the efficacy of the proposed
72 method. To the best of our knowledge, it is the first time converted SNN achieves ANN
73 level performance while maintaining high efficiency and low latency across various tasks.

74 2 Related Work

75 Up to now, there exist two main categories of supervised training strategies for SNN: approximated
76 error backpropagation methods, and conversion methods.

77 Our proposed work is mainly related to ANN-SNN conversion methods with rate-based coding [9,
78 35, 4, 39]. The goal of conversion approaches is to remain the high performance of ANN in SNN
79 framework. Recently, Xing et al. [44] extend the territory of conversions by adding supports to
80 more structures such as softmax activation and residual block. Instead of using longer timesteps
81 to approximate ANN accurately, conversion errors are analyzed and compensated by adjusting the
82 parameters of converted SNN [27]. Rate Norm is proposed [10] to better choose the scale factor.
83 Deng and Gu [8] utilize training samples to transfer the weights to the target SNN by combining
84 threshold balance and soft-reset mechanisms. As a general approach, Kim et al. [23] show that
85 conversion can be applied to detection tasks. However, only coding error are considered in those
86 work. We notice that the inference error or “unevenness error” is described in the parallel work [3].
87 However, the “unevenness error” is ignored during the analysis of conversion error in their work.

88 Other researchers [2, 29, 21] try to approximate the gradient to train the SNN by replacing the thresh-
89 old function with other functions. A comprehensive summary of surrogated gradient backpropagation
90 methods can be found in [30]. Huh [20] introduce differentiable synapse and neuron models to

91 SNNs. Since spike trains are normally sparse, gradients in SNNs are naturally sparse compared
 92 to traditional ANNs. An sparse spiking gradient descent approach is introduced [34] to speed up
 93 the backpropagation. To enable very deep spiking neural network, modified residual module with
 94 element-wise functions [14] is proposed to mimic directly gradient path of ResNet [17]. Rather
 95 than compute the approximated gradient, a neuron with extra inner state is introduced in [42] to
 96 estimate the gradient. Different from [42], the buffered neuron proposed in this literature utilizes extra
 97 inner state to compensate for inference error. Backpropagation trained SNNs show the potential of
 98 achieving ANN level performance. However, extra efforts are required to train the SNN properly,
 99 making it hard to train and limit its application on real world tasks.

100 Recently, combined methods [37, 36] emerges to alleviate the drawbacks that conversion methods
 101 require larger timesteps to achieve competitive performance. These methods often contain two steps:
 102 conversion and fine-tuning. The SNN converted in the first step is served as a weight initialization for
 103 a further fine-tuning procedure using backpropagation. With the help of roughly mapped weights,
 104 several epochs of fine-tuning yield nearly loss-less performance with fewer timesteps [37]. How-
 105 ever, compared to conversion methods, extra efforts are required to fine-tune the SNN. Also, the
 106 backpropagation imposes additional limitations to the hybrid approach, which further restrains its
 107 application.

108 3 Methods

109 In this section, we first analyse the theoretical equivalency and gaps between a target SNN and an
 110 original ANN. Then, a quantized activation is proposed to minimize the coding error during the
 111 conversion. A buffered non-leaky IF neuron is presented to reduce the inference error. Finally, we
 112 integrate the proposed conversion framework to further decrease the difference between SNN and
 113 ANN.

114 3.1 Theoretical Errors and Analysis

115 Following Rueckauer et al. [39], we here illustrate the equivalency between SNN and ANN and
 116 introduce two types of errors when using the conversion methods with rate-based coding. The target
 117 SNN shares a similar overall architecture with the original ANN except two significant discrepancies.
 118 First, neurons in the SNN communicate with each other by spike trains $S(t)$ where $S(t) \in \{0, 1\}$ for
 119 each timestep $t \in \{1, \dots, T\}$, while the activation functions in ANN utilize real-value. T is the total
 120 timesteps. Second, a non-leaky integrate-and-fire (IF) neuron is used instead of the ReLU activation
 121 in conventional ANN. The membrane function of the non-leaky IF neuron j can be described as:

$$V_j(t) = V_j(t-1) + V_{th} \sum_{i \in \mathcal{N}_j} w_{ij}^s S_i(t) - V_{th} S_j(t). \quad (1)$$

122 where $V_j(t)$ is the membrane potential of the j th neuron and t is the timestep. \mathcal{N}_j is the set of all
 123 input neurons in the last layer connecting to j th neuron. w_{ij}^s is the weight and connection strength
 124 between the input i th neuron and current j th neuron. V_{th} is the threshold for the neuron to fire a
 125 signal. $S_i(t)$ and $S_j(t)$ are input and the resulting output spike trains of current neuron, respectively.
 126 The firing signals are generated as follow:

$$S_j(t) = \begin{cases} 1, & \text{if } V_j'(t) \geq V_{th} \\ 0, & \text{otherwise.} \end{cases} \quad (2)$$

127 where $V_j'(t) = V_j(t-1) + V_{th} \sum_{i \in \mathcal{N}_j} w_{ij}^s S_i(t)$ denotes the membrane potential before signal firing
 128 and potential reset operation.

129 Considering ANN-SNN conversion and rate-based coding, we can compute spike rate as $r =$
 130 $\sum_{t=1}^T S(t)/T$ for a total timestep T . The goal of a conversion process is to derive a function transfers
 131 the parameters of an ANN to the SNN as $w^s = f(w^a)$ so that each neuron in the SNN produces a
 132 spike rate $r \propto a$ after T timesteps, where a is the corresponding activation value of original networks.

133 To this end, we can derive the spike count of the output spike train by summing Eqn. 1 over total
 134 timesteps T :

$$\sum_{t=1}^T S_j(t) = \sum_{i \in \mathcal{N}_j} w_{ij}^s \sum_{t=1}^T S_i(t) - \frac{V_j(T) - V_j(1)}{V_{th}}, \quad (3)$$

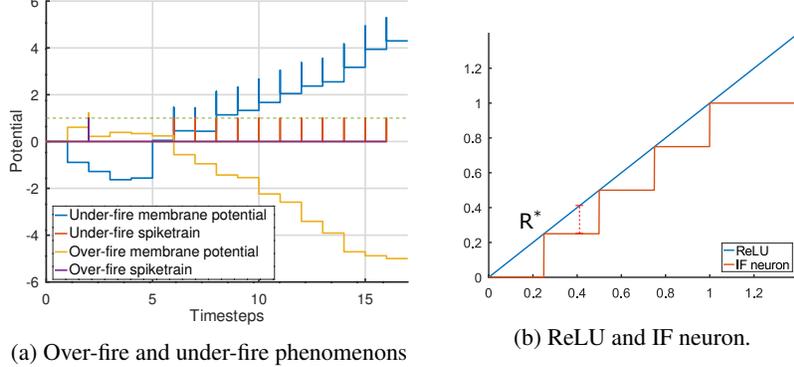


Figure 1: For each step of the membrane potential line in the figure, the left horizontal line indicates $V(t-1)$. The peak point indicates $V'(t)$, and the right horizontal line indicates $V(t)$ which is transferred to the next step. For over-fire spike train (purple), $\sum_{t=1}^T S_j(t) = 1$ while $\sum_{t=1}^T S_j^*(t) = 0$. For under-fire spike train (orange), $\sum_{t=1}^T S_j(t) = 11$ while $\sum_{t=1}^T S_j^*(t) = 15$.

135 where $V_j(1)$ is the initial membrane potential. For convenience, we set $V_j(1) = 0$ in this work
 136 without loss of generality. Eqn. 3 gives the relation between the total spike count of current neuron
 137 and its input neurons' spike count besides the internal membrane potential.

138 For an ideal conversion, the spiking neuron integrates membrane potential uniformly across time
 139 and introduces no error during inference. Under this ideal condition, the output spike count should
 140 fulfill the following equation,

$$\sum_{t=1}^T S_j^*(t) = \left\lfloor \max \left(\frac{V_a}{V_{th}}, 0 \right) \right\rfloor, \quad (4)$$

141 where $V_a = V_{th} \sum_{i \in \mathcal{N}_j} w_{ij}^s \sum_{t=1}^T S_i(t)$ denotes the accumulated membrane potential. Please refer
 142 to the appendix for detail proof. Assuming Eqn. 4 hold, we can derive the ideal output spike rate by
 143 averaging Eqn. 4 over timesteps T :

$$r_j^* = \max \left(\sum_{i \in \mathcal{N}_j} w_{ij}^s r_i, 0 \right) - R^*. \quad (5)$$

144 where R^* is defined as follow.

$$R^* = \begin{cases} \frac{V_a - \lfloor \frac{V_a}{V_{th}} \rfloor}{T}, & \text{if } V_a > 0 \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

145 For the ideal ANN-SNN conversion situation, when R^* is approaching zero, the formal expression of
 146 a simple IF neuron becomes similar to the ReLU function in processing information. Thus, the weight
 147 of a SNN can be transformed from its counterpart ANN directly. Since spike rate r^* is bounded, a
 148 linear mapping function is utilized to convert ANN weights $W^s = f(W^a) = \frac{\lambda_i}{\lambda_j} W^a$, in which λ_i is
 149 the scale factor for layer i .

150 3.1.1 Coding Error

151 As shown in Eqn. 5, R^* can be viewed as a *coding error* term for the ideal conversion process
 152 during the ANN-SNN weight mapping. As $\frac{V_a}{V_{th}} - \left\lfloor \frac{V_a}{V_{th}} \right\rfloor < 1$ for $V_a > 0$, the upper bound of the
 153 conversion error $R^* < 1/T$ continuously decreases when the timesteps is set increasingly. Due to
 154 errors accumulated between layers, the converted SNNs with deep layers require large timesteps T to
 155 maintain high fidelity and performance [39, 23]. On the other hand, the large timesteps impair the
 156 efficiency of the SNN. This phenomenon is also known as accuracy-latency trade-off [9, 31].

157 **3.1.2 Inference Error**

158 Besides the conversion error, the assumption in Eqn. 4 can be contravened easily due to essential
 159 nature of SNN neurons during the inference. Intuitively, the inference process of SNN is quite
 160 dynamic and spiking neurons are unable to accumulate membrane potential uniformly across
 161 timesteps as Eqn. 4, which leads to the divergence of overall output spike counts. We refer to this
 162 type of error as *inference error*.

163 We conclude two situations where the Eqn. 4 is violated during inference: over-fire and under-fire
 164 phenomenons. Over-fire indicates the situation where number of total generated spikes is greater
 165 than number of expected spikes $\sum_{t=1}^T S_j(t) > \sum_{t=1}^T S_j^*(t)$. Volatile membrane potential may
 166 cross the threshold $V_j^i(t) \geq V_{th}$ and cause spike $S_j(t) = 1$ that makes the number of generating
 167 spikes exceeding the number of expected spikes. Conversely, under-fire refers to the situation where
 168 insufficient spikes have been generated $\sum_{t=1}^T S_j(t) < \sum_{t=1}^T S_j^*(t)$. With inadequate spikes, under-
 169 fire neurons will have membrane potential large than V_{th} at the end of inference. Fig. 1a shows an
 170 example of these two situations.

171 **3.2 Quantized ReLU Activation**

172 To deal with the coding error, we propose a Quantized ReLU activation function in ANN training
 173 stage to accurately transform the weights to SNN with limited timesteps. According to Eqn.4, we can
 174 reformulate Eqn. 5 in a different way :

$$r_j^* = \frac{\lfloor T \max(\sum^{N_j} w_{ij}^s r_i, 0) \rfloor}{T}. \quad (7)$$

175 Let $x_s = \sum^{N_j} w^s r_i$ be the weighted sum of input spike rates in SNNs. Taking maximal spike rate
 176 into consideration, IF neuron(Fig. 1b) can be described with step function as follow,

$$IF(x_s) = \begin{cases} 0, & \text{if } x_s \leq 0 \\ r^{max}, & \text{if } x_s \geq r^{max} \\ \frac{\lfloor T x_s \rfloor}{T}, & \text{otherwise.} \end{cases} \quad (8)$$

177 The r^{max} is the maximum spike rate of a spike train, which is 1 for conventional spike trains. The
 178 proposed function is strictly equivalent to the IF neuron instead of increasing total timesteps T to
 179 approximate the ReLU function with a fine-grain step function. By mapping spike rate and activation
 180 value $a = \lambda r$, the proposed activation function is derived through Eqn. 8 :

$$QReLU(x_a, \lambda) = \lambda IF\left(\frac{x_a}{\lambda}\right) = \begin{cases} 0, & \text{if } x_a \leq 0 \\ \lambda, & \text{if } x_a \geq \lambda \\ \frac{\lfloor \frac{T}{\lambda} x_a \rfloor}{T}, & \text{otherwise.} \end{cases} \quad (9)$$

181 We name the proposed activation function as QReLU, since it can be seen as a quantized version of
 182 the ReLU function. The approximate gradients are given by the PACT algorithm [6] originated from
 183 STE [1]:

$$\frac{\partial QReLU(x_a, \lambda)}{\partial \lambda} = \begin{cases} 0, & \text{if } x_a < \lambda \\ 1, & \text{otherwise.} \end{cases} \quad (10)$$

$$\frac{\partial QReLU(x_a, \lambda)}{\partial x_a} = \begin{cases} 1, & \text{if } 0 < x_a < \lambda \\ 0, & \text{otherwise.} \end{cases} \quad (11)$$

184 The scale factor λ here is also referred to as quantization boundary and is learned during the ANN
 185 training process. Guaranteed by the success of quantization techniques [1, 6], a neural network
 186 consisting of QReLU is capable of achieving similar performance compared to the original ANN. By

187 mapping weights from an ANN equipped with QReLU activation, the coding error R^* is eliminated in
 188 theory.

189 3.3 Buffered Non-leaky IF Neuron

190 For the inference error, we observe that reducing over-fire problem alone can boost performance
 191 significantly based on our experiments. To further improve our conversion framework, we propose a
 192 novel buffered non-leaky IF neuron to relieve the over-fire situation to compensate for the inference
 193 errors.

194 The proposed buffered neuron utilizes a recurrent mechanism to regulate the spikes. In the proposed
 195 neuron, an additional buffered potential V^b is introduced as follows:

$$V_j^b(t) = V_j^b(t-1) + V_{th}S_j(t). \quad (12)$$

196 Similar to the membrane potential, the initial buffer potential $V_j^b(1)$ is set to zero for convenience.
 197 The spike generating process is described as

$$S_j(t) = \begin{cases} 1, & \text{if } V_j'(t) > V_{th} \\ -1, & \text{if } V_j'(t) < 0 \text{ and } V_j^b(t-1) \geq V_{th} \\ 0, & \text{otherwise.} \end{cases} \quad (13)$$

198 Intuitively, the proposed Buffered Non-leaky IF Neuron generates negative spikes to regulate the total
 199 spikes count for the inference errors. Thus, the over-fire problem can be alleviated. Please refer to
 200 appendix for detail proof. As same as ordinary IF neurons, only addition operation and threshold
 201 functions are involved in the proposed neuron. The inference efficiency is preserved since the extra
 202 power required by transmitting sign bit between neurons is negligible [45].

203 3.4 Integration with Multiple Bits Spike Train

204 To further improve the performance of the proposed conversion framework, we integrate shift-
 205 based multi-bits spike train into the proposed framework. Unlike the previous work [45] which
 206 compressed an existing SNN to speed up the inference, we extend the definition of the spike as
 207 $S(t) \in \{2^n | n \in \mathbb{N}, 2^n \leq S^{max}\}$ to facilitate the shift operation directly during the conversion. The
 208 multi-bits spike train version of Eqn.13 is as follow:

$$S_j(t) = \begin{cases} A\left(\left\lfloor \frac{V_j'(t)}{V_{th}} \right\rfloor\right), & \text{if } V_j'(t) > V_{th} \\ A\left(\min\left(\left\lfloor \frac{V_j'(t)}{-V_{th}} \right\rfloor, \left\lfloor \frac{V_j^b(t-1)}{V_{th}} \right\rfloor\right)\right), & \text{if } V_j'(t) < 0 \\ 0, & \text{and } V_j^b(t) \geq V_{th} \\ & \text{otherwise,} \end{cases} \quad (14)$$

209 where the adjust function $A(\cdot)$ is defined as

$$A(x) = 2^{\lfloor \log_2(\min(x, S^{max})) \rfloor}. \quad (15)$$

210 Multiple bits spike train carries more information than ubiquitous on-off spikes [5, 22, 33, 47].
 211 However, expensive multiplications are required during the integration with non-uniform spikes. To
 212 avoid multiplication, we employ the \log_2 adjust function and shift operation [45] as an alternative and
 213 achieve high energy efficiency. The adjust function can be implemented with a conditional function.

214 For QReLU, as the maximal ratio coding is equal to maximal spike counts, $r^{max} = \sum_{t=1}^T S^{max}/T =$
 215 S^{max} , the multi-bits version of Eqn. 9 can be described as follow:

Table 1: Classification performance compared to other methods in CIFAR10 and ImageNet-1k

Model	Arch	SNN Top-1	Δ	T	Model	Arch	SNN Top-1	Δ	T
CIFAR10					ImageNet				
Rueckauer et al. [39]	4 Conv, 2 Linear	90.85%	-1.06%	400	Rueckauer et al. [39]	VGG16	49.61%	-14.28%	400
Sengupta et al. [40]	ResNet20	87.46%	-1.64%	2500	Sengupta et al. [40]	VGG16	69.96%	-0.56%	2500
Lee et al. [26]	VGG9	90.45%	-	100	Han et al. [16]	VGG16	73.09%	-0.40%	4096
Kim and Panda [24]	VGG9	90.50%	-	25	Rathi et al. [37]	VGG16	65.19%	-4.16%	250
Wu et al. [43]	5 Conv, 2 Linear	90.53%	-	12	Rathi and Roy [36]	VGG16	66.52%	-3.56%	25
Han et al. [16]	ResNet20	91.36%	-0.11%	2048	Deng and Gu [8]	VGG16	55.80%	-16.6%	16
Rathi et al. [37]	ResNet20	92.22%	-0.93%	250	Bu et al. [3]	VGG16	50.97%	-23.32%	16
Rathi and Roy [36]	ResNet20	92.14%	-0.65%	25	Bu et al. [3]	VGG16	68.47%	-5.82%	32
Zheng et al. [48]	ResNet19	93.16%	-	6	Ours	VGG16	73.65%	-1.69%	8
Yu et al. [47]	6 Conv, 2 Linear	93.90%	-0.23%	300	Ours	VGG16	75.96%	-0.07%	16
Yan et al. [46]	VGG19	92.48%	-0.08%	600	Sengupta et al. [40]	ResNet34	65.47%	-5.22%	2000
Deng and Gu [8]	ResNet20	92.41%	0.09%	16	Han et al. [16]	ResNet34	65.47%	%	4096
Ding et al. [10]	PreActResNet18	91.96%	-1.10%	64	Fang et al. [14]	ResNet34	67.04%	-	4
Bu et al. [3]	ResNet18	94.82%	-1.22%	8	Zheng et al. [48]	ResNet34	63.72%	-	6
Bu et al. [3]	ResNet18	95.92%	-0.12%	16	Li et al. [27]	ResNet34	64.51%	-11.12%	32
Ours	ResNet18	95.13%	-0.17%	8	Li et al. [27]	ResNet34	74.61%	-1.05%	256
Ours	ResNet18	95.14%	-0.04%	16	Bu et al. [3]	ResNet34	59.35%	-14.97%	16
					Bu et al. [3]	ResNet34	69.37%	-4.95%	32
					Ours	ResNet34	73.20%	-0.52%	8
					Ours	ResNet34	74.46%	0.08%	16

Table 2: Detection performance

Model	Arch	SNN AP50	Δ	T
PASCAL VOC				
Kim et al. [23]	tiny-yolo	51.83%	-1.18%	5000
Ours	tiny-yolo	65.20%	-0.13%	16
Ours	tiny-yolo	65.73%	-0.06%	32
MS COCO				
Kim et al. [23]	tiny-yolo	25.66%	-0.58%	5000
Ours	tiny-yolo	38.6%	-1.1%	16
Ours	tiny-yolo	39.2%	-0.6%	32

Table 3: SNN computation efficiency.

Methods	Architecture (timesteps)	Operation ratio	SNN/ANN energy ratio
CIFAR10			
Rathi and Roy [36]	ResNet20(25)	2.55	0.51
Ours	ResNet18(8)	0.92	0.19
Ours	ResNet18(16)	1.69	0.34
ImageNet			
Rathi and Roy [36]	VGG16(25)	1.62	0.32
Ours	VGG16(8)	1.92	0.39
Ours	VGG16(16)	3.03	0.62

$$QReLU(x_a, \lambda^b) = \begin{cases} 0, & \text{if } x_a \leq 0 \\ \lambda^b, & \text{if } x_a \geq \lambda^b \\ \lfloor \frac{P}{\lambda^b} x_a \rfloor, & \text{otherwise.} \end{cases} \quad (16)$$

216 Here the $\lambda^b = \lambda S^{max}$ is the quantization boundary and $P = TS^{max}$ is the *representation power*.
 217 As implied by the definition of the representation power, even the minimal shift extension with
 218 $S^{max} = 2$ can maintain the same representation power with half timesteps. This phenomenon is
 219 named *strength-latency trade-off*.

220 4 Experiments

221 In this section, we conduct experiments on both classification and detection tasks to demonstrate the
 222 efficiency and effectiveness of the proposed method. Since the equivalency of features and spike rates
 223 are universal throughout the entire network, real-valued input is directly utilized in the converted
 224 SNN. We select a unity threshold $V_{th} = 1$ for all our experiments. Max spike value $S^{max} = 2$ is
 225 used for all experiments if not stated otherwise. Other implementation details can be found in the
 226 appendix.

227 4.1 Visual Object Recognition and Detection Performance

228 We test our proposed High Performance Conversion framework (HPC) on visual object recognition
 229 and detection tasks. A ResNet18 [17] like architecture is utilized as backbone networks to con-
 230 duct experiments on the CIFAR10 [25] dataset and a VGG16 [41] like network is trained on the
 231 ImageNet [7] dataset. Table 1 reports the classification accuracy of the proposed HPC compared

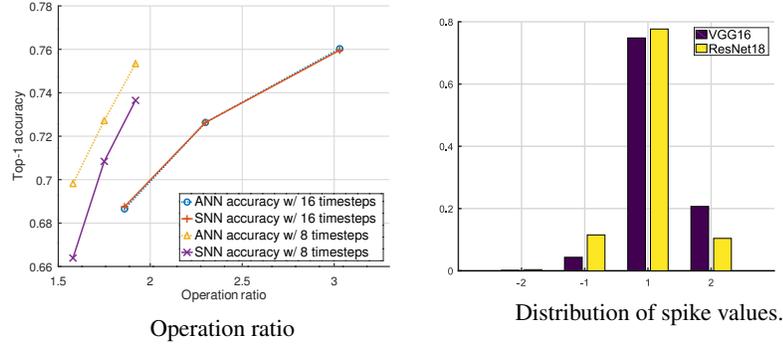


Figure 2: Left: Operation ratio and performance with ImageNet dataset under 16 timesteps. Data points with same operation ratio are converted SNN and its corresponding ANN. Right: Distribution of spike values.

232 with other methods. Top-1 accuracy is used to evaluate the performance. Since the performance of
 233 converted SNN is closely related to the ANN baseline, a performance drop Δ is also reported. As
 234 illustrated in Table 1, the proposed HPC can perform nearly lossless conversion while maintaining
 235 low inference latency.

236 As a general method, we also test our HPC with tiny-yolo [38] structure on both PASCAL VOC [13]
 237 and MS COCO [28] datasets. Table 2 summarizes the AP50 performance of the proposed method.
 238 We can see that our proposed HPC can achieve state-of-the-art conversion performance within 32
 239 timesteps, which is over 100 times faster than the previous SNN detection work [23].

240 4.2 Energy Efficiency

241 To evaluate the energy efficiency of the proposed method, *operation ratio* is adopted to measure the
 242 overall efficiency as described in Eqn. 17. In contrast to actual energy consumption which may vary
 243 between hardware devices [32], counting-based metric [36] is straightforward and representative
 244 for different methods. For a fair comparison, integration operations triggered by spikes for both
 245 membrane potential and buffer potential are counted in our experiments.

$$operation\ ratio = \frac{\#(addition\ ops\ in\ SNN)}{\#(flops\ in\ ANN)}. \quad (17)$$

246 Operation ratios, as well as estimated energy consumption ratio, are summarized in Table 3. Here
 247 we roughly estimate that multiplication consumes 4 times more energy than addition according
 248 to Horowitz [18]. The results on the CIFAR10 dataset show a strong correlation between timesteps
 249 and energy efficiency. To further study energy efficiency, we test the relationship between operation
 250 ratio and networks performance. As depicted in Fig. 2a, under the same timesteps setting, the
 251 spike ratio continuously decreases with the overall performance. Our experiments confirm a general
 252 relationship between inference energy and performance, which holds even under constant latency.

253 Fig. 2b depicts the distributions of spike values in classification experiments. Without loss of gener-
 254 ality, we conclude that only 20% of the spikes require shift operation. As measured by Gudovskiy
 255 and Rigazio [15], the shift operation only cost 10% of energy consumed by addition operation. Thus,
 256 multi-bits spike only use 2% extra energy, which is negligible.

257 4.3 Ablation Study

258 To reveal the effectiveness of the proposed framework, we conduct ablation studies on both clas-
 259 sification and detection tasks. Specifically, we remove the quantized training technique and the
 260 buffered IF neuron sequentially and compare the results. The baseline ANN is trained and fixed for
 261 experiments without QReLU component, meanwhile quantized ANN networks are trained separately
 262 for different timestep configurations. Without the QReLU, the scale factor is calculated using the
 263 percentile function with the parameter empirically set to 99.0. Channel-wise normalization is utilized
 264 to ensure competitive performances. Since conventional converted SNN requires large timesteps to

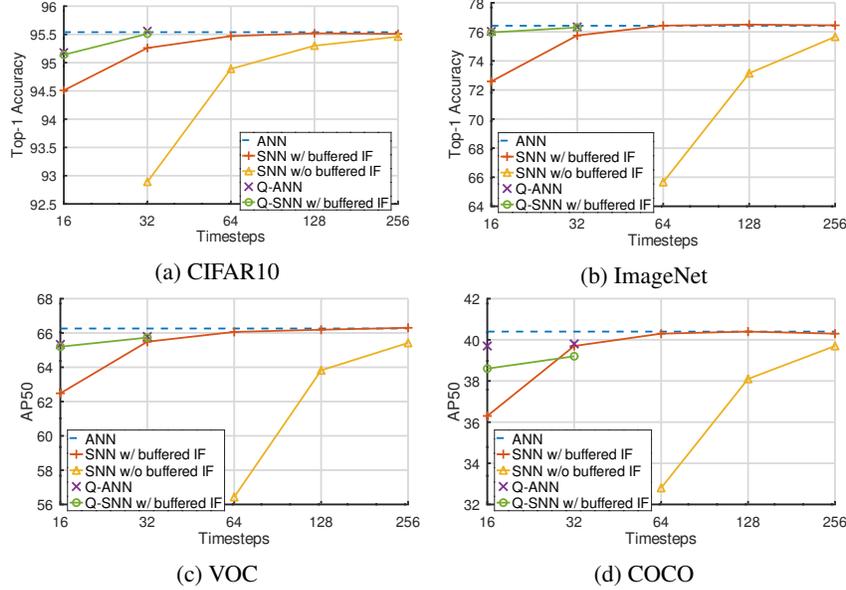


Figure 3: Ablation study of proposed techniques. Q-ANN denotes ANNs trained with QReLU, Q-SNN denotes SNNs converted from ANN trained with QReLU. SNN w/ buffered IF means that buffered non-leaky integrated and fire neuron are used. SNN w/o buffered IF are conventional SNNs with out the proposed neuron.

265 achieve plausible performance, we only report those performances with 32 or more timesteps for the
 266 CIFAR10 dataset and 64 for others.

267 Fig. 3 validates the effectiveness of the proposed techniques. Compared to conventional SNNs,
 268 the SNN integrates with buffered neurons only requires a quarter of timesteps to reach the same
 269 performance. As the buffered neuron mitigates only the over-fire problem, we conclude that the
 270 over-fire situation is the primary cause of performance degradation during inference. The introduction
 271 of quantized training further enhances the performance of converted SNN under limited timesteps.
 272 Similar trends can be observed across experiments.

273 5 Conclusion and Limitation

274 In this paper, we proposed a novel High Performance Conversion (HPC) method to simultaneously
 275 achieve high-performance SNN inference with low latency. Efforts had been made to reduce both
 276 coding errors and often overlooked inference errors. A novel QReLU activation was proposed to
 277 eliminates the conversion error. Meanwhile, a buffered non-leaky IF neuron was designed to mitigate
 278 the over-fire problem and boost the inference performance while maintaining its simplicity and
 279 efficiency. For the first time, efficient SNN converted without extra fine-tuning revealed its capability
 280 to achieve state-of-the-art performances in nontrivial vision tasks.

281 However, there are limitations. Since no actual hardware is involved, the efficiency is estimated using
 282 energy consumed by simple operations. As computing platform evolves rapidly, it is beyond our
 283 reach to accomplish real neuromorphic hardware which can be researched in further work.

284 For future work, we believe further study on novel structures such as attention mechanisms SE [19]
 285 and transformer [12] can extend the current framework and may lead to better understanding of the
 286 underlying perception mechanism. Moreover, applications of SNN on video-like input that achieves
 287 ANN level performance with great efficiency are another promising direction.

288 References

289 [1] Y. Bengio, N. Léonard, and A. Courville. Estimating or propagating gradients through stochastic
 290 neurons for conditional computation. *CoRR*, 2013.

- 291 [2] S. M. Bohte, J. N. Kok, and H. La Poutré. Error-backpropagation in temporally encoded
292 networks of spiking neurons. *Neurocomputing*, 48(1):17 – 37, 2002. ISSN 0925-2312. doi:
293 10.1016/S0925-2312(01)00658-0.
- 294 [3] T. Bu, W. Fang, J. Ding, P. DAI, Z. Yu, and T. Huang. Optimal ANN-SNN conversion for
295 high-accuracy and ultra-low-latency spiking neural networks. In *International Conference on*
296 *Learning Representations*, 2022.
- 297 [4] Y. Cao, Y. Chen, and D. Khosla. Spiking deep convolutional neural networks for energy-efficient
298 object recognition. *International Journal of Computer Vision*, 113(1):54–66, May 2015. ISSN
299 1573-1405. doi: 10.1007/s11263-014-0788-3.
- 300 [5] R. Chen, H. Ma, S. Xie, P. Guo, P. Li, and D. Wang. Fast and efficient deep sparse multi-strength
301 spiking neural networks with dynamic pruning. In *2018 International Joint Conference on*
302 *Neural Networks (IJCNN)*, pages 1–8, 2018. doi: 10.1109/IJCNN.2018.8489339.
- 303 [6] J. Choi, Z. Wang, S. Venkataramani, P. I.-J. Chuang, V. Srinivasan, and K. Gopalakrishnan.
304 Pact: Parameterized clipping activation for quantized neural networks. *CoRR*, 2018.
- 305 [7] J. Deng, W. Dong, R. Socher, L. Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical
306 image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages
307 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- 308 [8] S. Deng and S. Gu. Optimal conversion of conventional artificial neural networks to spiking
309 neural networks. In *International Conference on Learning Representations*, 2021.
- 310 [9] P. U. Diehl, D. Neil, J. Binas, M. Cook, S. Liu, and M. Pfeiffer. Fast-classifying, high-accuracy
311 spiking deep networks through weight and threshold balancing. In *2015 International Joint*
312 *Conference on Neural Networks (IJCNN)*, pages 1–8, 2015. doi: 10.1109/IJCNN.2015.7280696.
- 313 [10] J. Ding, Z. Yu, Y. Tian, and T. Huang. Optimal ann-snn conversion for fast and accurate
314 inference in deep spiking neural networks. In *IJCAI*, 2021. doi: 10.24963/ijcai.2021/321.
- 315 [11] S. Dong, T. Huang, and Y. Tian. Spike camera and its coding methods. In *Data Compression*
316 *Conference (DCC)*, 2017.
- 317 [12] A. Dosovitskiy and et al. An image is worth 16x16 words: Transformers for image recognition
318 at scale. In *ICLR*, 2021.
- 319 [13] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal
320 visual object classes (voc) challenge. *International Journal of Computer Vision*, 2010. doi:
321 10.1007/s11263-009-0275-4.
- 322 [14] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, and Y. Tian. Deep residual learning in
323 spiking neural networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors,
324 *Advances in Neural Information Processing Systems*, 2021.
- 325 [15] D. Gudovskiy and L. Rigazio. Shiftcnn: Generalized low-precision architecture for inference of
326 convolutional neural networks. *CoRR*, 2017.
- 327 [16] B. Han, G. Srinivasan, and K. Roy. Rmp-snn: Residual membrane potential neuron for enabling
328 deeper high-accuracy and low-latency spiking neural network. In *CVPR*, 2020.
- 329 [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In
330 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
331 June 2016.
- 332 [18] M. Horowitz. 1.1 computing’s energy problem (and what we can do about it). In *2014 IEEE*
333 *International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, pages 10–14,
334 2014. doi: 10.1109/ISSCC.2014.6757323.
- 335 [19] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *CVPR*, pages 7132–7141,
336 2018. doi: 10.1109/CVPR.2018.00745.
- 337 [20] D. Huh and T. J. Sejnowski. Gradient descent for spiking neural networks. In S. Bengio,
338 H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, *Advances in*
339 *Neural Information Processing Systems*, volume 31, pages 1433–1443. Curran Associates, Inc.,
340 2018.
- 341 [21] Y. Jin, W. Zhang, and P. Li. Hybrid macro/micro level backpropagation for training deep spiking
342 neural networks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and

- 343 R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 31, pages
344 7005–7015. Curran Associates, Inc., 2018.
- 345 [22] J. Kim, H. Kim, S. Huh, J. Lee, and K. Choi. Deep neural networks with weighted spikes.
346 *Review of Economic Dynamics*, 311:373–386, Oct. 2018. ISSN 1094-2025. doi: 10.1016/j.
347 neucom.2018.05.087. Funding Information: This work was supported by the KIST Institutional
348 Program (Project No. 2E27330-17-P026).
- 349 [23] S. Kim, S. Park, B. Na, and S. Yoon. Spiking-yolo: Spiking neural network for energy-efficient
350 object detection. *AAAI*, 2020.
- 351 [24] Y. Kim and P. Panda. Revisiting batch normalization for training low-latency deep spiking
352 neural networks from scratch, 2020.
- 353 [25] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. *Tech*
354 *Report*, 2009.
- 355 [26] C. Lee, S. S. Sarwar, P. Panda, G. Srinivasan, and K. Roy. Enabling spike-based backpropagation
356 for training deep neural network architectures. *Frontiers in Neuroscience*, 14, Feb 2020. ISSN
357 1662-453X. doi: 10.3389/fnins.2020.00119.
- 358 [27] Y. Li, S. Deng, X. Dong, R. Gong, and S. Gu. A free lunch from ann: Towards efficient, accurate
359 spiking neural networks calibration. In *ICML*, pages 6316–6325, 2021.
- 360 [28] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick.
361 Microsoft coco: Common objects in context. In *ECCV*, 2014. ISBN 978-3-319-10602-1.
- 362 [29] H. Mostafa. Supervised learning based on temporal coding in spiking neural networks. *IEEE*
363 *Transactions on Neural Networks and Learning Systems*, 29(7):3227–3235, 2018. doi: 10.1109/
364 TNNLS.2017.2726060.
- 365 [30] E. O. Neftci, H. Mostafa, and F. Zenke. Surrogate gradient learning in spiking neural networks:
366 Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal*
367 *Processing Magazine*, 36(6):51–63, 2019. doi: 10.1109/MSP.2019.2931595.
- 368 [31] D. Neil, M. Pfeiffer, and S.-C. Liu. Learning to be efficient: Algorithms for training low-latency,
369 low-compute deep spiking neural networks. In *Proceedings of ACM Symposium on Applied*
370 *Computing*, 2016.
- 371 [32] D. E. Nikonov and I. A. Young. Benchmarking physical performance of neural inference
372 circuits. *CoRR*, 2019.
- 373 [33] S. Park, S. Kim, H. Choe, and S. Yoon. Fast and efficient information transmission with burst
374 spikes in deep spiking neural networks. In *2019 56th ACM/IEEE Design Automation Conference*
375 *(DAC)*, pages 1–6, 2019.
- 376 [34] N. Perez-Nieves and D. Goodman. Sparse spiking gradient descent. In M. Ranzato, A. Beygelz-
377 imer, Y. Dauphin, P. Liang, and J. W. Vaughan, editors, *Advances in Neural Information*
378 *Processing Systems*, volume 34, pages 11795–11808. Curran Associates, Inc., 2021.
- 379 [35] J. A. Pérez-Carrasco, B. Zhao, C. Serrano, B. Acha, T. Serrano-Gotarredona, S. Chen, and
380 B. Linares-Barranco. Mapping from frame-driven to frame-free event-driven vision systems by
381 low-rate rate coding and coincidence processing—application to feedforward convnets. *IEEE*
382 *Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2706–2719, 2013. doi:
383 10.1109/TPAMI.2013.71.
- 384 [36] N. Rathi and K. Roy. Diet-snn: Direct input encoding with leakage and threshold optimization
385 in deep spiking neural networks. *CoRR*, 2020.
- 386 [37] N. Rathi, G. Srinivasan, P. Panda, and K. Roy. Enabling deep spiking neural networks with
387 hybrid conversion and spike timing dependent backpropagation. In *International Conference on*
388 *Learning Representations*, 2020.
- 389 [38] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time
390 object detection. In *CVPR*, June 2016.
- 391 [39] B. Rueckauer, I. Lungu, Y. Hu, M. Pfeiffer, and S. Liu. Conversion of continuous-valued deep
392 networks to efficient event-driven networks for image classification. *Frontiers in Neuroscience*,
393 2017. ISSN 1662-453X. doi: 10.3389/fnins.2017.00682.

- 394 [40] A. Sengupta, Y. Ye, R. Wang, C. Liu, and K. Roy. Going deeper in spiking neural networks:
395 Vgg and residual architectures. *Frontiers in Neuroscience*, 13:95, 2019. ISSN 1662-453X. doi:
396 10.3389/fnins.2019.00095.
- 397 [41] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image
398 recognition. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning*
399 *Representations*, 2015.
- 400 [42] J. C. Thiele, O. Bichler, and A. Dupret. Spikegrad: An ann-equivalent computation model
401 for implementing backpropagation with spikes. In *International Conference on Learning*
402 *Representations*, 2020.
- 403 [43] Y. Wu, L. Deng, G. Li, J. Zhu, Y. Xie, and L. Shi. Direct training for spiking neural networks:
404 Faster, larger, better. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01):
405 1311–1318, Jul. 2019. doi: 10.1609/aaai.v33i01.33011311.
- 406 [44] F. Xing, Y. Yuan, H. Huo, and T. Fang. Homeostasis-based cnn-to-snn conversion of inception
407 and residual architectures. In T. Gedeon, K. W. Wong, and M. Lee, editors, *Neural Information*
408 *Processing*, pages 173–184, Cham, 2019. Springer International Publishing. ISBN 978-3-030-
409 36718-3.
- 410 [45] C. Xu, W. Zhang, Y. Liu, and P. Li. Boosting throughput and efficiency of hardware spiking
411 neural accelerators using time compression supporting multiple spike codes. *Frontiers in*
412 *Neuroscience*, 14:104, 2020. ISSN 1662-453X. doi: 10.3389/fnins.2020.00104.
- 413 [46] Z. Yan, J. Zhou, and W.-F. Wong. Near lossless transfer learning for spiking neural networks.
414 *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12):10577–10584, May 2021.
- 415 [47] Q. Yu, C. Ma, S. Song, G. Zhang, J. Dang, and K. C. Tan. Constructing accurate and efficient
416 deep spiking neural networks with double-threshold and augmented schemes, 2020.
- 417 [48] H. Zheng, Y. Wu, L. Deng, Y. Hu, and G. Li. Going deeper with directly-trained larger
418 spiking neural networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(12):
419 11062–11070, May 2021.

420 Checklist

421 The checklist follows the references. Please read the checklist guidelines carefully for information on
422 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
423 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
424 the appropriate section of your paper or providing a brief inline description. For example:

- 425 • Did you include the license to the code and datasets? **[Yes]** See Section ??.
- 426 • Did you include the license to the code and datasets? **[No]** The code and the data are
427 proprietary.
- 428 • Did you include the license to the code and datasets? **[N/A]**

429 Please do not modify the questions and only use the provided macros for your answers. Note that the
430 Checklist section does not count towards the page limit. In your paper, please delete this instructions
431 block and only keep the Checklist section heading above along with the questions/answers below.

- 432 1. For all authors...
 - 433 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
434 contributions and scope? **[Yes]**
 - 435 (b) Did you describe the limitations of your work? **[Yes]**
 - 436 (c) Did you discuss any potential negative societal impacts of your work? **[No]**
 - 437 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
438 them? **[Yes]**
- 439 2. If you are including theoretical results...
 - 440 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
 - 441 (b) Did you include complete proofs of all theoretical results? **[Yes]**

- 442 3. If you ran experiments...
- 443 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 444 mental results (either in the supplemental material or as a URL)? [No] Source code
- 445 will be released later.
- 446 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 447 were chosen)? [Yes] Training details can be found in the appendix.
- 448 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 449 ments multiple times)? [Yes] Performance stability is reported in the appendix.
- 450 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 451 of GPUs, internal cluster, or cloud provider)? [Yes] Training details can be found in
- 452 the appendix.
- 453 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 454 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 455 (b) Did you mention the license of the assets? [No] Popular datasets are used.
- 456 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 457 (d) Did you discuss whether and how consent was obtained from people whose data you're
- 458 using/curating? [N/A]
- 459 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 460 information or offensive content? [N/A]
- 461 5. If you used crowdsourcing or conducted research with human subjects...
- 462 (a) Did you include the full text of instructions given to participants and screenshots, if
- 463 applicable? [N/A]
- 464 (b) Did you describe any potential participant risks, with links to Institutional Review
- 465 Board (IRB) approvals, if applicable? [N/A]
- 466 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 467 spent on participant compensation? [N/A]