

1 Thanks for the constructive comments. Let us categorize them into groups and answer due to the page limit.
 2 **Reasoning behind ANTLR:** We would like to emphasize that our main contribution is to discover the complementary
 3 nature of two main approaches for SNN training and provide a justification to combine them, not just randomly adding
 4 two different gradients. A spike-train of the neuron i can be represented in two different ways: $\Phi_i = \{S_i[t] \text{ for } t =$
 5 $0, 1, \dots, T - 1\}$ (*activation-based representation*) and $\Phi_i = \{\hat{t}_{i,k} \text{ for } k = 0, 1, \dots, N - 1\}$ (*timing-based representation*).
 6 One of the main arguments of our paper (in Section 3.1) is that the change in the state of neuron i can be described as
 7 $\Delta\Phi_i = \Delta\Phi_{i,\text{gen}} + \Delta\Phi_{i,\text{remov}} + \Delta\Phi_{i,\text{shift}}$, where $\Delta\Phi_i = \{\Delta S_i[t] \text{ for } t = 0, 1, \dots, T - 1\}$, $\Delta\Phi_{i,\text{gen}} = \{\Delta S_i[t] = +1 \text{ for } t \in \mathcal{T}_{\text{gen}}\}$,
 8 $\Delta\Phi_{i,\text{remov}} = \{\Delta S_i[t] = -1 \text{ for } t \in \mathcal{T}_{\text{remov}}\}$, and $\Delta\Phi_{i,\text{shift}} = \{\Delta S_i[t_{\text{before}}] = -1, \Delta S_i[t_{\text{after}}] = +1 \text{ for } \{t_{\text{before}}, t_{\text{after}}\} \in \mathcal{T}_{\text{shift}} =$
 9 $\{\hat{t}_{i,k} \text{ for } k = 0, 1, \dots, N - 1\}\}$. To update any parameter x for desired network output, $\frac{\partial\Phi_i}{\partial x}$ must be precisely
 10 computed so that it satisfies $\Delta x \cdot \frac{\partial\Phi_i}{\partial x} \approx \Delta\Phi_i$. Due to the limitation in considering the reset path for gradient

11 computation, $\frac{\Delta\Phi_{i,\text{shift}}}{\Delta x}$ cannot be precisely estimated/predicted by the activation-based gradient $\frac{\partial\Phi_i}{\partial x}|_{\text{act}}$. In contrast,
 12 the timing-based gradient $\frac{\partial\Phi_i}{\partial x}|_{\text{tim}}$ cannot estimate $\Delta\Phi_{i,\text{gen}}$ and $\Delta\Phi_{i,\text{remov}}$. Those changes cannot even be described
 13 in the *timing-based representation*. The solution we proposed with ANTLR is to compensate the inaccuracy in
 14 gradients by adding two imperfect and complementary values as: $\frac{\partial\Phi_i}{\partial x}|_{\text{act}} \approx \frac{\Delta\Phi_{i,\text{gen}}}{\Delta x} + \frac{\Delta\Phi_{i,\text{remov}}}{\Delta x}$, $\frac{\partial\Phi_i}{\partial x}|_{\text{tim}} \approx \frac{\Delta\Phi_{i,\text{shift}}}{\Delta x}$,
 15 $\frac{\partial\Phi_i}{\partial x}|_{\text{ant}} = \lambda_{\text{act}} \frac{\partial\Phi_i}{\partial x}|_{\text{act}} + \lambda_{\text{tim}} \frac{\partial\Phi_i}{\partial x}|_{\text{tim}} \approx \lambda_{\text{act}} \frac{\Delta\Phi_{i,\text{gen}}}{\Delta x} + \lambda_{\text{act}} \frac{\Delta\Phi_{i,\text{remov}}}{\Delta x} + \lambda_{\text{tim}} \frac{\Delta\Phi_{i,\text{shift}}}{\Delta x}$

16 **Layer-wise vs network-wise summation:** As Reviewer #2 pointed out, one can run the activation-based method and
 17 the timing-based method separately and use the summation of those two gradients for parameter update. However, there
 18 are major differences between that approach and ANTLR. First of all, depending on the type of the loss function, one
 19 of learning methods may provide zero gradients for every parameter (Table 1). Then combining the two methods by
 20 network-wise summation can be meaningless. Moreover, if the activation- and timing-based gradients are not combined
 21 in each layer, errors in the gradients are accumulated through the back-propagation along multiple layers.

22 **Coefficients λ :** We introduced the coefficients $\lambda_{\text{act}}, \lambda_{\text{tim}}$ to balance the gradients from two methods because the scale of
 23 the activation-based gradient can be arbitrarily changed by the hyper-parameters of the surrogate derivative. Even
 24 though we simply used $\lambda_{\text{act}}, \lambda_{\text{tim}} = 1$ in this work for convenience, the optimal coefficients should further be studied.

25 **Adding/removing spikes:** Reviewer #2 mentioned that timing-based methods can also add/remove spikes when a single
 26 spike constraint is relaxed. Timing-based methods may unintentionally generate/remove spikes as a result of parameter
 27 update while they try to shift the spike timing to reduce the loss. However, these *unintended* generation/removal of
 28 spikes do not contribute to training, as the timing-based method cannot estimate them in gradient computation.

29 **No-spike penalty:** The no-spike penalty (Line 191-192) has been used in the timing-based methods because they cannot
 30 train parameters when a neuron does not emit any spike (dead neuron problem). It is implemented by encouraging
 31 *every neuron* to emit at least one spike. One of the main advantages of ANTLR is that it can solve the dead neuron
 32 problem more efficiently, using variant of the count loss $\{\min(\sum_{\tau} S_d[\tau], 1) - 1\}^2$ (d is the index of the desired class
 33 label) that penalizes only *output neurons* without no spike. With the help of the activation-based part, ANTLR can
 34 add/remove spikes in both hidden and output neurons while allowing some neurons not to emit any spikes. Reviewer #2
 35 commented that addition of the count loss would decrease the firing activity but it actually *increases* the spiking activity
 36 because it still penalizes some neurons without spikes.

37 **Comparison with related works:** In the submission, we did not report the comparison with other works which did
 38 not report the sparsity of spikes since it is not fair to compare each method solely by accuracy results. The accuracy
 39 of SNNs highly depends on the number of spikes used. To support the argument, we compared previous results of
 40 fully-connected SNNs on N-MNIST classification tasks with our experimental results (Table R1). Even though related
 41 works did not report exact results of spike numbers, our experiments using similar settings imply that previous works
 42 with high accuracy were benefited from the large amount of spike usage.

43 **Experimental Settings:** As the reviewers men-
 44 tioned, our experimental results used different set-
 45 tings for each method. We reported the best (in
 46 terms of accuracy and efficiency) results from each
 47 method after we tested every option available to
 48 each method. ANTLR can provide similar accu-
 49 racy compared to the activation-based method
 50 when the same setting is used (Table R1). However,
 51 this situation is not desirable for efficient use of
 52 SNNs because it requires larger number of spikes.

53 **Experiments using larger datasets:** We agree
 54 that our experimental results are only from rela-
 55 tively small datasets. We wanted to focus on fun-
 56 damental limitations of existing learning methods
 57 in gradient computation in this work. Experiments and analysis of ANTLR on larger datasets remain as a future study.

Table R1: Comparison of fully-connected SNNs on N-MNIST

Method	Type*	Accuracy [%]	Loss**	# Target spikes	# Spikes/sample
Lee et al. [27]	S	98.66	C	not fixed	N/A
Jin et al. [26]	S	98.84±0.02	C	35 / 5	N/A
SLAYER [10]	A	98.89±0.06	C	60 / 10	N/A
STBP [11]	A	98.78	C	300 / 0	N/A
SRM-based***	A	97.73±0.14	C	10 / 0	436±17
SRM-based***	A	98.30±0.06	C	60 / 10	6536±120
ANTLR	A&T	97.73±0.09	C	10 / 0	415±14
ANTLR	A&T	98.05±0.10	C	60 / 10	6638±130
Timing***	T	94.10±0.51	L	-	2166±294
ANTLR	A&T	96.58±0.25	L	-	111±8

* A (activation-based), T (timing-based), and S (scalar-mediated, refer to Section 5), ** C (count loss) and L (latency loss), *** Our implementation of existing approaches