

1 First and foremost, we want to thank all the reviewers for their thorough reviews. We are especially happy about the
2 high scores we have received from *Reviewer #3* and *Reviewer #4*, but also appreciate the in-depth analysis by *Reviewer*
3 *#1* and the suggestions of *Reviewer #2*.

4 The problem we (and the NPU) are solving is to learn a subset of arithmetic operations (AOs), namely (\times, \div, x^w) . Prior
5 art either has problems with small and negative numbers (NALU) or implements only a subset of AOs (NMU). We have
6 demonstrated these limitations in experiments 2 & 3. *Reviewer #2* was concerned about sufficient novelty because we
7 (just) perform the algebraic operations in *complex* rather than *real* space. We believe that it is precisely this simplicity,
8 which makes our approach convincing, because it does not require additional, more complicated mechanisms to work
9 (like e.g. *iNALU*), and it fixes the shortcomings of prior work e.g. by a mathematically correct treatment of negative
10 inputs. *Reviewer #3* is correct that relevance gates that are not exactly zero or one lead to a complex treatment of real
11 inputs, which introduces an error in the result. We believe this error is actually an advantage because the optimization
12 has to push the gate to either zero or one to remove the error. However, we agree that gates should converge to exactly
13 zero or one, and we plan to address this in future work (e.g. by an additional penalization).
14 In the following three paragraphs, we address the criticism expressed concerning our experiments as they appear in the
15 paper.

16 **Experiment 1 - Equation Discovery** Since the NALU reported problems with convergence, this experiment was set to
17 demonstrate that the NPU can learn models of sequential data, which is notoriously difficult due to vanishing/exploding
18 gradients, and bifurcations (also faced e.g. by RNNs). Additionally, small changes in the parameters can lead to
19 substantial changes in the output after a few time steps. The experiment shows that the NPU can be used for sequential
20 tasks, which we believe is owed to the combination of complex weights and the relevance gate.

21 Our second intention with this experiment was to outline how to utilize the transparency of the NPU as one part within
22 a broader equation discovery framework. Since the NPU can represent more AOs than prior art, it is possible to fit the
23 same data with simpler models (in terms of the number of parameters). This improves *transparency by decomposability*
24 (see [Lipton](#)), which enables to extract a *practically useful* equation containing only a few terms. In the paper, we called
25 this *interpretability*, which we will correct as suggested by *Reviewer #1*. We did not claim to recover the correct fSIR
26 model, but rather something close to the original model. Finalization of the equation discovery would require more
27 post-processing such as a detailed search of models near the obtained result, e.g., by using methods from binary neural
28 networks or the regularization by penalizing entropy, as used in the [GNN explainer](#). Considering the feedback on this
29 experiment, we should have been more explicit about its purpose. We propose to change the title to "*A Step Towards*
30 *Equation Discovery*", and to clarify in the text that the experiment is a proof of concept.

31 In response to *Reviewer #3*: The large MSE that some NPU realizations exhibit in Fig. 4 is caused by diverging models
32 at the beginning of training (due to the difficulties named above), from which the models sometimes do not recover. We
33 used ADAM to obtain a good initial solution further refined by LBFGS, which can improve MSE by around an order of
34 magnitude.

35 **Experiment 2 - Simple Arithmetic** In this experiment, we wanted to show that the NPU can learn fractional power
36 functions, and additionally, that it finds a solution where the NALU fails. While the NALU can learn $+$, \times , and $\sqrt{\cdot}$ in
37 isolation (see [Trask et al.](#)), it does not converge for the more complicated task of learning those AOs in one model, as in
38 our setup of this task. [Schlör et al.](#) suggested that this might be due to the gating between addition and multiplication
39 paths in the NALU. We only tested one range in this experiment, because \div and $\sqrt{\cdot}$ quickly become approximately
40 linear away from zero, making the gradient signal to the correct solution (i.e., equation) weak. Training models on such
41 a weak gradient signal is very sensitive to step-size in SGD and requires increasing the batch sizes. We have therefore
42 considered training on a range with stronger gradient signals to be a fairer comparison between models.

43 **Experiment 3 - Large Scale Arithmetic** In this experiment, we aimed to stay close to the established benchmark
44 by [Madsen & Johansen](#). We have slightly modified it since the original setup of the task makes it difficult to learn \div ,
45 which is again due to weak gradient signals at the effectively sampled ranges. Specifically, their experiment applies a
46 single AO to the sum of two overlapping subsets of a vector. For example, summing 50 numbers from a vector with
47 samples from $\mathcal{U}(0, 1)$ - via the central limit theorem - results in samples for the AO of interest from a narrow Gaussian
48 centered at $\mu = 25$. For $+$ and \times , which were the interest of the NMU experiments, the signal is still strong, but for \div ,
49 this is results in a weak signal and thus a very hard task for optimization. Therefore, we chose an easier inverse x^{-1} ,
50 which also proves our point (that the NPU can represent and learn division by a variable).

51 Finally, we would like to thank the reviewers again for their effort and offer to update our description of experiment 3
52 as stated above, and to add a negative example to the broader impact section (the misuse of the setup in experiment 3
53 for real-world COVID-19 predictions) in case our paper should be accepted.