
Learning Division with Neural Arithmetic Logic Modules

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 To achieve systematic generalisation, it first makes sense to master simple tasks
2 such as arithmetic. Of the four fundamental arithmetic operations (+, -, \times , \div),
3 division is considered the most difficult for both humans and computers. In this
4 paper we show that robustly learning division in a systematic manner remains a
5 challenge even at the simplest level of dividing two numbers. We propose two
6 novel approaches for division which we call the Neural Reciprocal Unit (NRU) and
7 the Neural Multiplicative Reciprocal Unit (NMRU), and present improvements for
8 an existing division module, the Real Neural Power Unit (Real NPU). Experiments
9 in learning division with input redundancy on 225 different training sets, find that
10 our proposed modifications to the Real NPU obtains an average success of 85.3%
11 improving over the original by 15.1%. In light of the suggestion above, our NMRU
12 approach can further improve the success to 91.6%.

13 1 Introduction

14 Imagine you must learn to divide 2 numbers, but are only given 10 numbers and the target value. This
15 task requires finding the 2 relevant operands, the order to divide the operands, and learning to divide.
16 In machine learning, this is equivalent to a supervised regression task where the aim is to learn the
17 underlying function between the inputs and output such that the solution is generalisable to any input.
18 The ability to select relevant features is a desirable property of neural networks, useful for improved
19 interpretability, reduced pre-processing costs and greater generalisation [Chandrashekar and Sahin,
20 2014]. The ability to model division, one of the four fundamental arithmetic operations, is necessary
21 for expressing dynamical systems [Sahoo et al., 2018], and physics-based formulas [Udrescu and
22 Tegmark, 2020]. However, even recent models still struggle to learn division when there is input
23 redundancy [Schlör et al., 2020].

24 The main challenge of the above task comes from learning the selection and operation at the same
25 time, which can lead to conflicting priorities when learning network weights. Furthermore, the natural
26 properties of division of values around zero leads to undesirable gradients. Models which deal with
27 this naively (e.g. MLPs) are unable to deal with the fluctuant gradients caused by the asymptotic
28 nature and discontinuities in division [Trask et al., 2018].

29 Can we build models which can learn division in the presence of its undesirable, yet valid, properties?
30 We aim to address this question in this paper. Specifically, we contribute the following:¹

- 31 • **Improvements to the Real NPU** [Heim et al., 2020] including: clipping, discretisation and
32 constrained initialisation to improve performance in learning division on different training ranges.

¹Code (MIT license) available at: <https://anonymous.4open.science/r/nalu-stable-exp-neurips-review-2E4C/>.

- 33 • **Two novel division modules**, the NRU and the NMRU. The NRU explores extending the NMU
34 weight ranges from $[0,1]$ to $[-1,1]$ to include division, where we find a weakness in learning from
35 negative ranges. Learning from the weaknesses of the NRU, the NMRU extends the NMU to learn
36 division while keeping weights values between $[0,1]$. We further boost performance by using a
37 Real NPU inspired sign retrieval mechanism, enabling the NMRU to gain the best performance
38 when using a mean squared error (MSE) loss.
- 39 • **New understanding into the hindrances in learning division** including: training on mixed-sign
40 inputs, training on negative ranges, and division on extremely small values. We find these difficulties
41 can be sufficiently identified using synthetic division tasks.

42 The broader impact of our work relates to interpretable Artificial Intelligence where our modules
43 can be included in larger networks for applications such as image classification or analogy creation,
44 whilst retaining the ability to produce transparent generalisable solutions. However, there are possible
45 negative societal impacts. Such modules can be viewed as specialised feature selectors/aggregators
46 which do not require integrating domain knowledge. Therefore, if a non-domain-expert tries inter-
47 preting relations in the input data, they may incorrectly interpret causality, which can be especially
48 harmful if such a case occurs on medical or financial data. Mitigating against such downstream issues
49 requires to first focus efforts on producing robust modules to different distributions and understand
50 their affect on learning other networks architectures (e.g. CNN). Understanding this will enable
51 recognising situations where these modules can aid and where they should avoid being used.

52 2 Related Work

53 One approach to learn division would be symbolic regression networks [Sahoo et al., 2018]. However,
54 a symbolic approach pre-defines the operations, which is not a limitation of using Neural Arithmetic
55 Logic Modules (NALMs).

56 NALMs are neural networks which learn arithmetic operations and input selection [Mistry et al.,
57 2021]. The weights of these networks are interpretable such that a discrete value represents a specific
58 operation. For example, ‘-1’ to represent division and ‘0’ for no selection. From this research field,
59 we focus on the Real NPU and the NMU. Until now, the Real NPU only has learned division on
60 training ranges of either $\mathcal{U}[0.1,2]$ or Sobol(0,0.5) [Heim et al., 2020]. It remains unclear if this
61 module is robust to other training ranges even as a stand-alone unit. Robustness to training ranges is
62 important as these module’s applicational use comes from being part of larger end-to-end networks,
63 where the input range into the module cannot be controlled. The NMU is a multiplication module
64 which we extend to also do division. The authors of the NMU believe such an extension incurs too
65 many limitations for learning [Madsen and Johansen, 2020]. We use this paper as an opportunity to
66 explore this belief.

67 Trask et al. [2018] developed the Neural Arithmetic Logic Unit (NALU) which can model all four
68 arithmetic operations. However, studies show this module to be unstable in learning division [Schlör
69 et al., 2020, Heim et al., 2020]. In particular, their gating method responsible for selecting an operation
70 cannot learn consistently [Madsen and Johansen, 2020]. Schlör et al. [2020] developed iNALU
71 additionally applying weight and gradient clipping, sign retrieval, regularisation, reinitialisation and
72 separating shared parameters to the NALU. Even with these modifications, they still find consistently
73 learning division to a high precision to remain unattainable. Furthermore, Heim et al. [2020]’s results
74 imply iNALU is outperformed by the Real NPU for division.

75 3 Architectures

76 This section introduces the architectures for the (Real) NPU, NRU, and the NMRU. The (Real) NPU
77 is an existing module, which we improve in Section 5. The NRU and NMRU are novel contributions.
78 Appendix A summarises the important properties of these division modules.

79 3.1 Real Neural Power Unit

80 Heim et al. [2020] develop a module to learn to multiply and divide, using the intuition from Trask
81 et al. [2018] that *multiplicative operations are additive operations in log space*. Their work extends
82 this idea into complex space. The NPU can be used with its complex form (Equation 1) requiring both

83 a complex and real weight matrix ($\mathbf{W}^{(i)}$, $\mathbf{W}^{(r)}$), or only its real form the Real NPU (Equation 2).
 84 For improved gradients, a relevance gate \mathbf{r} (Equation 3) is used which converts inputs close to 0 (i.e.
 85 irrelevant features) to 1 to avoid the resulting output evaluating to 0. A gating vector \mathbf{g} , learns to
 86 select relevant input elements, where gate values are clipped between [0,1] during training.

$$\text{NPU} := \exp(\mathbf{W}^{(r)} \log(\mathbf{r}) - \mathbf{W}^{(i)} \mathbf{k}) \odot \cos(\mathbf{W}^{(i)} \log(\mathbf{r}) + \mathbf{W}^{(r)} \mathbf{k}), \quad (1)$$

87

$$\text{RealNPU} := \exp(\mathbf{W}^{(r)} \log(\mathbf{r})) \odot \cos(\mathbf{W}^{(r)} \mathbf{k}) \quad (2)$$

88

$$\text{where } \mathbf{r} = \mathbf{g} \odot (|\mathbf{x}| + \epsilon) + (\mathbf{1} - \mathbf{g}) \quad \text{and} \quad k_i = \begin{cases} 0 & x_i \geq 0 \\ \pi g_i & x_i < 0 \end{cases}. \quad (3)$$

89 A weighted L1 penalty is used when training. The weight value β grows between predefined values
 90 β_{start} to β_{end} and is increased every $\beta_{step} = 10,000$ iterations by a growth factor $\beta_{growth} = 10$.
 91 We focus on the Real NPU over the NPU as the solution of the tasks in this paper can be captured
 92 using only real values meaning that the complex form is not required.

93 3.2 Neural Reciprocal Unit

94 We propose the NRU, which can model multiplication and division. We extend the NMU, motivated
 95 by *division being multiplication of reciprocals*. The range which weight values can be is extended
 96 from [0,1] to [-1,1], where -1 represents applying the reciprocal on the corresponding input element.
 97 A NRU output element z_o is defined as

$$\text{NRU} : z_o = \prod_{i=1}^I (\text{sign}(x_i) \cdot |x_i|^{W_{i,o}} \cdot |W_{i,o}| + 1 - |W_{i,o}|), \quad (4)$$

98 where I is the number of inputs. Assuming weights are either 1 (multiply) or -1 (reciprocal), $|x_i|^{W_{i,o}}$
 99 will apply the operation on an input element. The absolute value is used so that the module only
 100 operates in the space of real numbers, as $x_i^{W_{i,o}}$ for a negative input (x_i) when $-1 < W_{i,o} < 1$ results
 101 in a complex number. The use of absolute means the sign of the input must be reapplied. For the
 102 no-selection case $W_{i,o} = 0$, we want the input element to convert to 1 (the identity value), resulting
 103 in applying $\cdot |W_{i,o}| + 1 - |W_{i,o}|$. The derivative of the absolute function at 0 is undefined meaning the
 104 gradients of Equation 4 can contain points of discontinuity. To alleviate this issue, we approximate
 105 the absolute function using a scaled tanh (inspired by Faber and Wattenhofer [2020]). More formally,

$$|W_{i,o}| = \begin{cases} \tanh(1000 \cdot W_{i,o})^2 & \text{if training} \\ |W_{i,o}| & \text{otherwise} \end{cases}.$$

106 The scale factor (1000) controls how close to the absolute function the approximation is, where larger
 107 values give a more accurate approximation. For clipping and regularisation, the same scheme as the
 108 Neural Addition Unit (NAU) (see Appendix B) is used.

109 3.3 Neural Multiplicative Reciprocal Unit

110 An alternate extension of the NMU, also motivated by *division being multiplication of reciprocals*
 111 is the NMRU (Equation 5). We concatenate the reciprocal of the input (plus a small ϵ) to the input
 112 resulting in a module which only needs to learn selection. Hence, weights can be in the range [0,1].

$$\text{NMRU} : z_o = \prod_{i=1}^{2I} (W_{i,o} \cdot |x_i| + 1 - W_{i,o}) \cdot \sum_{i=1}^{2I} (\cos(W_{i,o} \cdot k_i)), \quad \text{where } k_i = \begin{cases} 0 & x_i \geq 0 \\ \pi & x_i < 0 \end{cases}. \quad (5)$$

113 The iteration over $2I$ represents the going through all inputs and their reciprocals. We calculate the
 114 magnitude and sign separately, joining the result at the end. The magnitude is calculated passing
 115 absolute of the concatenated input through an NMU architecture and the sign by using a cosine
 116 mechanism similar to the Real NPU. However, unlike the Real NPU only the weight matrix is
 117 required. The norm of the weight's gradients are clipped to 1 prior to being updated by the optimiser.
 118 This is done to alleviate the issue of exploding gradients caused by including the reciprocal to the
 119 inputs. For clipping and regularisation, the same scheme as the NMU (see Appendix B) is used.

Table 1: Interpolation (train/validation) and extrapolation (test) ranges used. Data (as floats) is drawn from a Uniform distribution with the range values as the lower and upper bounds.

Interpolation	[-20, -10)	[-2, -1)	[-1.2, -1.1)	[-0.2, -0.1)	[-2, 2)
Extrapolation	[-40, -20)	[-6, -2)	[-6.1, -1.2)	[-2, -0.2)	[[[-6, -2), [2, 6)]
Interpolation	[0.1, 0.2)	[1, 2)	[1.1, 1.2)	[10, 20)	
Extrapolation	[0.2, 2)	[2, 6)	[1.2, 6)	[20, 40)	

120 4 Experiment Setup

121 We introduce the two main experiments used to evaluate modules, including: default parameters,
 122 train and test ranges, and evaluation metrics. The tasks evaluate the ability of a single module to
 123 divide two numbers from an input vector in two settings: **no redundancy** and **with redundancy**.

124 **Default parameters:** All experiments use a mean squared error (MSE) loss with an Adam optimiser
 125 [Kingma and Ba, 2015], with 10,000 samples for the validation and test sets. The best model for
 126 evaluation is taken using early stopping on the validation set. All runs are over 25 different seeds. All
 127 inputs are required in the *no redundancy* setting, i.e., input size of 2. Training takes 50,000 iterations
 128 where each iteration consists of a different batch of size 128. The Real NPU uses a learning rate of
 129 $5e-3$ with sparsity regularisation scaling during iterations 40,000 to 50,000. The NRU and NMRU
 130 use sparsity regularisation scaling during iterations 20,000 to 35,000 and a learning rate of 1 and $1e-2$
 131 respectively. In contrast, the *redundancy* setting uses an input size of 10, where 8 input values are not
 132 required for the final output. The total training iterations are extended to 100,000 with batch sizes
 133 of 128. The learning rates for the Real NPU, NRU and NMRU are $5e-3$, $1e-3$ and $1e-2$ respectively.
 134 Sparsity regularisation scaling occurs during iteration 50,000 to 75,000 for all modules. A summary
 135 of all relevant parameters is found in Appendix C.

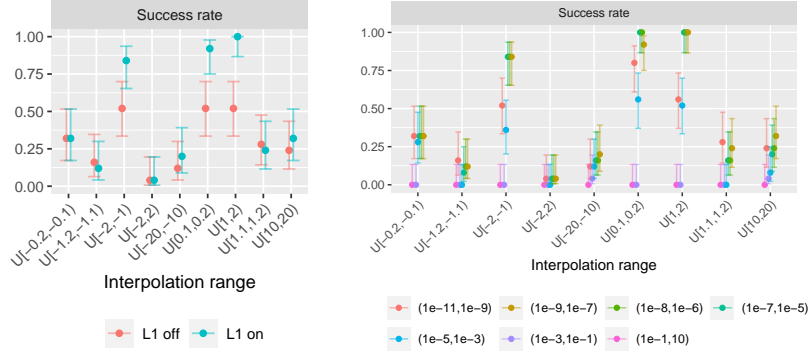
136 **Ranges:** The interpolation (train/validation) and extrapolation (test) ranges, are found in Table 1.
 137 The chosen ranges are influenced by Madsen and Johansen [2020].

138 **Evaluation metrics:** We use the Madsen and Johansen [2019]’s evaluation scheme, consisting of
 139 three evaluation metrics: the success on the extrapolation dataset against a near optimal solution
 140 (*success rate*), the first iteration which the task is considered solved (*speed of convergence*), and
 141 the extent of discretisation towards the weights’ inductive biases (*sparsity error*). Sparsity error
 142 calculated by $\max_{i,o}(\min(|W_{i,o}|, 1 - |W_{i,o}|))$, measures the weight element which is the furthest away
 143 from the acceptable discrete weights for the module. A success means the MSE of the trained model
 144 is lower than a threshold value (i.e. the MSE of a near optimal solution). We differ from Madsen
 145 and Johansen [2019] by using a fixed threshold value $1e-5$ rather than a simulated MSE, as there
 146 are no intermediate layers to accumulate numerical errors. We choose this precision as it can be
 147 guaranteed when working with 32-bit PyTorch Tensors. 95% confidence intervals (over the 25 seeds)
 148 are calculated from a specific family of distributions dependant on the metric. The success rate uses
 149 Binomial distribution because trials (i.e. run on a single seed) are either pass/ fail situations. The
 150 convergence metric uses a Gamma distribution and sparsity error uses a Beta distribution. Both Beta
 151 and Gamma can easily approximate the normal distribution and support its corresponding metric.

152 5 Improving the Real NPU’s Robustness

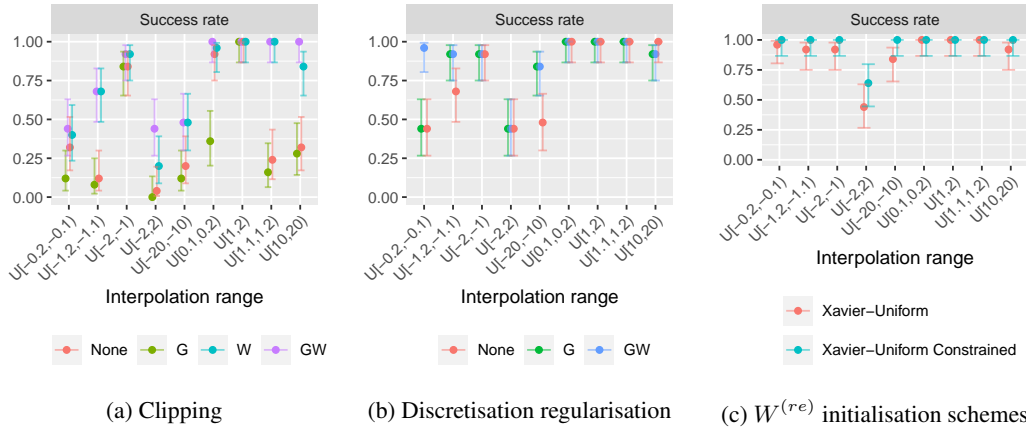
153 We first improve the robustness of the Real NPU on different training ranges. We use the Single
 154 Module Task with no redundancy (see Section 4) to investigate the following questions:

- 155 1. Is L1 regularisation required, and if so, do the regularisation parameters require tuning?
- 156 2. Does clipping the weight matrix aid learning?
- 157 3. Does enforcing discretisation on parameters improve convergence?
- 158 4. Can the weight matrix initialisation be improved?



(a) L1 regularisation (b) Sweep over L1 (start,end) beta parameters

Figure 1: Exploring the effect and sensitivity of L1 regularisation on the Real NPU



(a) Clipping (b) Discretisation regularisation (c) $W^{(re)}$ initialisation schemes

Figure 2: Effect of clipping, discretisation, and the NAU initialisation scheme on the Real NPU.

159 To address each question in order, we propose applying incremental modifications to the Real NPU.
 160 These modifications include: ablation study on the L1 regularisation (including a sweep over the
 161 scaling range hyperparameters), clipping, enforcing discretisation, and a more restrictive initialisation
 162 scheme. We assume that we are optimising the Real NPU to perform multiplication or division.
 163 Therefore, we trade-off the flexibility of having non-discretised weights, which enables the success of
 164 modelling the SIR data in Heim et al. [2020, Section 4.1], in favour of sparse models with discrete
 165 weight values. All the modifications suggested can also be generalised for the NPU architecture.

166 **Is L1 regularisation required? (Yes)** L1 encourages sparsity (i.e., zero weights) in solutions.
 167 Zero-valued weights means not to select an input and return the identity value 1. For the task, the
 168 optimal weight values require selecting all inputs and therefore non-zero values, suggesting the
 169 application of L1 could be damaging. Therefore, we compare against a model which does not use
 170 L1 regularisation, shown in Figure 1a. Removing L1 proves to be detrimental in five of the nine
 171 cases shown and only shows minor improvements in two of the nine ranges (i.e., $\mathcal{U}[-1.2,-1.1]$ and
 172 $\mathcal{U}[1.1,1.2]$). Hence, we keep L1 regularisation. The L1 regularisation scaling (see Section 3.1),
 173 requires setting the hyperparameters for the start (β_{start}) and end (β_{end}) scaling values. We run a
 174 sweep over six different start and end values, denoted ($\langle start \rangle$, $\langle end \rangle$), displaying results in Figure 1b.
 175 We find the configuration (1e-9, 1e-7) is the most successful when considering performance on all
 176 the ranges, and larger scaling values perform worse.

177 **Does clipping the learnable parameters help? (Yes)** Division and multiplication operations are
 178 represented by weight values of -1 and 1 respectively. The current architecture does not constrain the
 179 weights which can result in large weight values. The gate weights do get clipped and saved to another

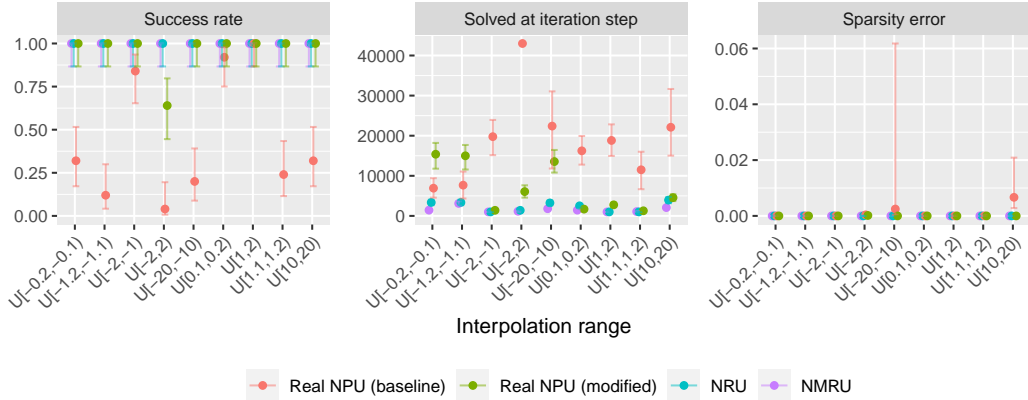


Figure 3: Division without redundancy (input size 2).

180 variable during the forward pass, meaning after an update step the gate values can also be out of the
 181 range $[-1,1]$. Hence, we investigate the effect of applying clipping directly to the weight and gate
 182 values after every optimisation step. Results, shown in Figure 2a, show clipping is beneficial, with
 183 clipping on both weight and gate (or just on the weights) to improve over the baseline on all ranges
 184 (excluding $\mathcal{U}[1,2]$ where the baseline has already achieved full success).

185 **Does enforcing discretisation help? (Yes)** Modelling division in a generalisable manner requires
 186 all learnable parameters to be discrete i.e., a value from $\{-1, 0, 1\}$. Using Madsen and Johansen
 187 [2020]’s regularisation scaling scheme, we penalise weights for not being discrete. We modify the
 188 scaling factor to be $\hat{\lambda} = 1$ and the regularisation to go from ‘off’ to ‘on’ between iterations 40,000 to
 189 50,000. Results, shown in Figure 2b, show discretising the gate improves over the baseline but also
 190 discretising the weights is additionally beneficial (especially for range $\mathcal{U}[-0.2,-0.1]$). $\mathcal{U}[10,20]$ is the
 191 only range where the baseline outperforms using discretisation, succeeding on two additional seeds.

192 **Does using a more constrained initialisation help? (Yes)** $W^{(r)}$ uses a Xavier-Uniform initial-
 193 isation [Glorot and Bengio, 2010]. This can result in weights initialised out of the range $[-1,1]$.
 194 Therefore, we use the initialisation for the Neural Addition Unit which is a constrained form of
 195 the Xavier-Uniform that does not allow the fan values of the uniform distribution to go beyond 0.5,
 196 meaning that no weight value will be out of the range $[-1,1]$ [Madsen and Johansen, 2020]. Figure 2c
 197 shows using the constrained initialisation provides improvements over multiple ranges.

198 6 Results: Single Module Task

199 We analyse the results for the: Real NPU without using the modifications of Section 5, Real NPU
 200 with modifications, NRU, and NMRU.

201 6.1 No Redundancy

202 Figure 3 shows the baseline Real NPU without modifications struggles with all ranges except $\mathcal{U}[1,2]$,
 203 struggling with sparsity on the larger ranges. Applying the modifications deals with the sparsity issue
 204 and improves the robustness such that only range $\mathcal{U}[-2,2]$ struggles (with a success rate of 0.64). The
 205 NRU and NMRU achieve full success over all ranges while solving the problem consistently fast and
 206 with low sparsity error. The success of the NRU is correlated with the learning rate (see Appendix E).

207 6.1.1 Mixed-signed Inputs

208 The remaining failure range of the Real NPU is $\mathcal{U}[-2,2]$ where inputs can consist of arbitrary signed
 209 values (e.g. all positives, all negatives, or a mixture of positive and negative values). *We question*
 210 *if the failure is due to the input samples in a batch having different signs from each other, or if the*
 211 *problem is due to the fact data samples can be close to 0 (leading to singularity issues).* To investigate

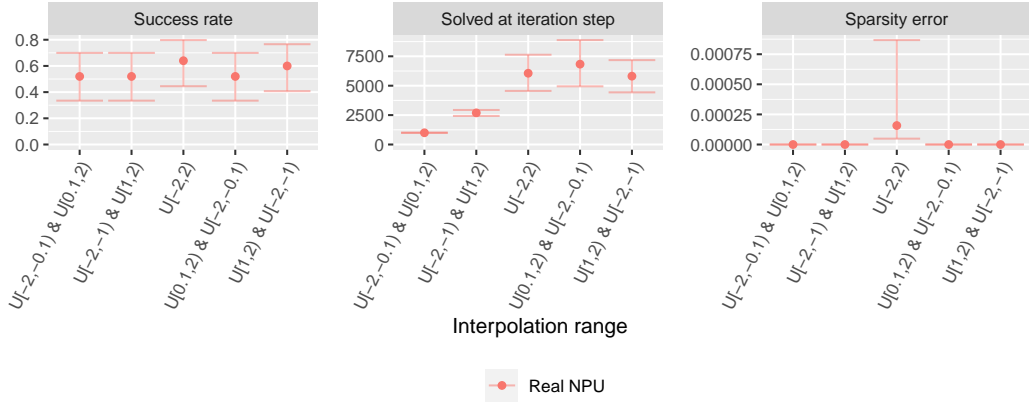


Figure 4: Extrapolation results on training the Real NPU using mixed-sign datasets that control the sign of the input elements. The ranges are in order of the datasets (i.e. dataset 1 to 5).

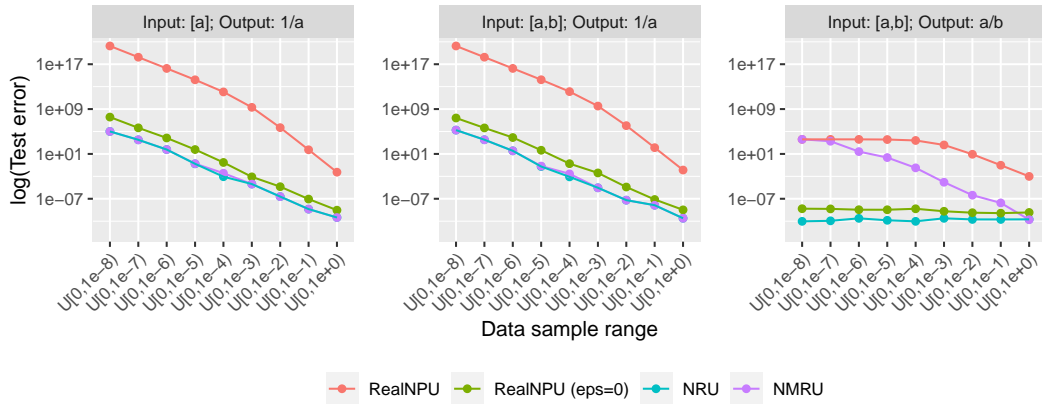


Figure 5: Effect of the singularity issue on the Real NPU, NRU and NMRU over increasing input ranges. Left: Reciprocal for an input size of 1 (no redundancy). Middle: Reciprocal for an input size of 2 (with redundancy). Right: Division for an input size of 2 (no redundancy).

212 this, we create additional mixed-sign datasets, controlling the range for each element in the input. The
 213 interpolation and extrapolation ranges for the different datasets can be found in Appendix C. Datasets
 214 1, 2, 4 and 5 sample a positive value for one input element and a negative value for the other element.
 215 Dataset 3 samples the signs randomly. Datasets 2 and 5 avoid sampling close to 0 values to mitigate
 216 the singularity issue. As shown by Figure 4, the Real NPU struggles on all these ranges, implying that
 217 the core issue is not from different input samples having different signs or due to the input samples
 218 being able to contain small values close to 0. The underlying issue is therefore most likely correlated
 219 to the each element in an input having different signs. When the denominator of the output is positive
 220 (dataset 1 or 2), the solution is found faster than when the denominator is a negative value (dataset 4
 221 or 5). When the signs for an input element are controlled, discretisation/sparsity is no problem, in
 222 contrast when the signs are arbitrary the sparsity error are slightly (though not significantly) higher.

223 6.2 Division by Small Numbers

224 Division by zero remains a challenge to model due to the inability to provide an computational value
 225 for the output and gradient. Furthermore, the discontinuous nature at zero causes its neighbouring
 226 values to have large gradients. To understand the extent of this issue when learning, we explore
 227 learning to divide by values close to zero using three tasks with increasing difficulty: 1) learning to
 228 take the reciprocal of a single input, 2) taking the reciprocal of the first input given two inputs, and 3)
 229 diving the first input by the second given two inputs. Figure 5 plots the test error for different modules

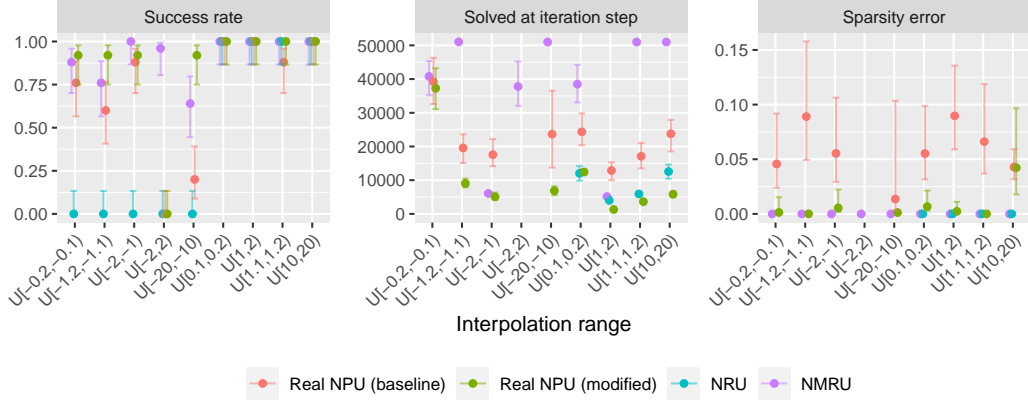


Figure 6: Division with redundancy (input size 10).

230 assuming the module weights are set to the ‘gold’ solution for the three tasks. As the range values
 231 become closer to zero, the test error thresholds become increasingly large. Therefore, even with the
 232 correct weights, relying on the test errors alone as an indicator become increasingly deceptive with
 233 values close to zero. The Real NPU has larger test errors for all tasks and ranges, caused by adding ϵ
 234 to the input (see Equation 3). Setting $\epsilon = 0$ reduces the test error at the cost of the ability to deal with
 235 zero-valued inputs. Appendix F provides the corresponding experimental results for these tasks.

236 6.3 With Redundancy

237 Introducing redundancy (Figure 6) causes failure modes to arise. Failures on range $\mathcal{U}[-2,2]$ become
 238 more prevalent. The baseline Real NPU produces high sparsity errors relative to the other modules
 239 suggesting struggle with discretisation. Using the modified Real NPU improves over all ranges of the
 240 baseline (which were not already at full success) in terms of success, speed and sparsity.² To ensure
 241 that complex weights do not fix the issue, we test the NPU module with all the modifications used on
 242 the real weight matrix (see Appendix G). Complex weights hinders success and convergence speeds
 243 of negative ranges. Assuming the global solution only uses the real weights, we enforce the complex
 244 weights to be clipped between $[-1,1]$ and to go to 0 during the regularisation stage using a L1 penalty.
 245 This did not result in any significant improvements against the Real NPU results. Input redundancy
 246 effects the NRU the most, resulting in full failures on all the negative ranges. The NMRU is the
 247 only module with success for the range $\mathcal{U}[-2,2]$, which is a result of using the sign mechanism (see
 248 Appendix H). It performs well over all ranges though can be outperformed by the modified Real NPU
 249 for negative ranges. Multiple ranges for the NMRU are solved around 50,000 iterations correlating to
 250 the sparsity regularisation being turned on.

251 6.3.1 Gradient Difficulties with the NRU

252 The partial derivative for the NRU weights, Equation 6, can give insight to the struggles of the NRU.

$$\frac{\partial \hat{\mathbf{y}}}{\partial w_i} = \tanh(1000w_i)(\text{sign}(x_i)|x_i|(\tanh(1000w_i) \log(|x|) + 2000 \text{sech}(1000w_i)^2) - 2000 \text{sech}(1000w_i)^2) \times \text{NRU}_{\tilde{\mathbf{x}} \in \mathbf{x} \setminus \{x_i\}}(\tilde{\mathbf{x}}). \quad (6)$$

253 $\text{NRU}_{\tilde{\mathbf{x}} \in \mathbf{x} \setminus \{x_i\}}(\tilde{\mathbf{x}})$ applies the NRU to all inputs excluding x_i influencing the gradient values between
 254 subsequent update steps. Factoring out this term, the following observations are made. If $x_i \approx 0$ and
 255 $w_i \approx 0$ then gradients become increasingly large. If $x_i \approx 0$ and $-1 \leq w_i < 0$ then as $w_i \rightarrow -1$ all
 256 gradients for x_i where $|x_i| \gg 1$ become increasingly small. The gradients for $x_i = -1$ and $x_i = 1$
 257 are 0 regardless the value of w_i . If $w_i = 0$ then the gradient is 0 for all x_i , a result of using the tanh
 258 approximation. Even if the sign and magnitude are calculated separately and then combined (see
 259 Appendix I) to try to control the gradient better, the problem remains. Therefore, we conclude that
 260 extending the NMU to divide using a weight of -1 is a poor choice when there are redundant inputs.

²Except for the sparsity error for range $\mathcal{U}[10,20]$.

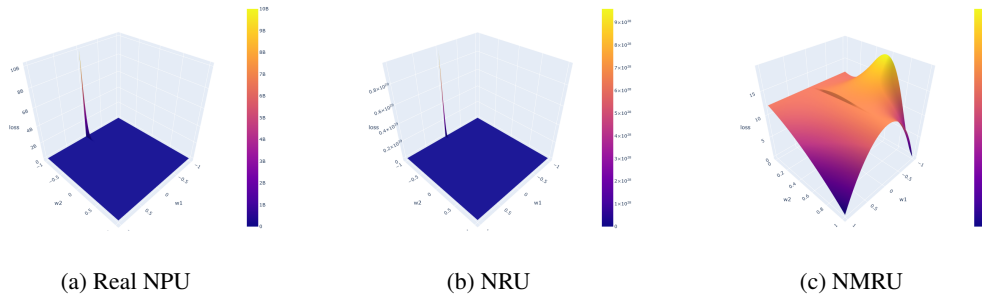


Figure 7: Root Mean Squared loss curvature for the NAU stacked with either a RealNPU, NRU, or NMRU. "The weight matrices are constrained to $\mathbf{W}_1 = \begin{bmatrix} w_1 & w_1 & 0 & 0 \\ w_1 & w_1 & w_1 & w_1 \end{bmatrix}$, $\mathbf{W}_2 = [w_2 \ w_2]$. The problem is $(x_1 + x_2) \cdot (x_1 + x_2 + x_3 + x_4)$ for $x = (1, 1.2, 1.8, 2)$ " [Madsen and Johansen, 2020]. The ideal solution is $w_1 = w_2 = 1$, though other valid solutions do exist e.g., $w_1 = -1, w_2 = 1$. (The NMRU's weight matrix would be $\mathbf{W}_2 = [w_2 \ w_2 \ 0 \ 0]$, and the Real NPU's $\mathbf{g} = [1 \ 1]$.)

261 6.3.2 The Real NPU's and NMRU's Exploitation of Multiplicative Rules

262 The NMRU solutions exploit the inverse rule of division in that $a_i \cdot \frac{1}{a_i} = 1$. Since the input also
 263 contains the reciprocals, numerous extrapolative solutions exist. However this comes at the cost of
 264 finding a 'simple' solution which contains ones only for relevant inputs. The Real NPU exploits the
 265 rules $a_i \cdot 0 = 0$ and $1^{a_i} = 1$ enabling non-zero weight values if the corresponding gate value is 0.
 266 However, we can avoid this by allowing 0 to also not be penalised during sparsity regularisation stage
 267 (see Appendix G). We find this alleviates the exploitation issue with no cost to performance.

268 7 Discussion

269 In this paper, we demonstrate the limitations of interpretable neural networks in learning to divide.
 270 Using the no redundancy setting (size 2), we find that the Real NPU is challenged when training data
 271 consists of mixed-signed inputs even with our applied improvements. Increasing the difficulty to
 272 have an input redundancy (with 8 redundant and 2 relevant input values) magnifies this issue, but
 273 also introduces failure modes for the NRU and NMRU for negative ranges. The NRU is unable to
 274 handle any negative ranges, in which we conclude it is not wise to use with MSE. Alternate losses
 275 can improve certain failure cases though sometimes at the cost of performance on other ranges. For
 276 further details see Appendix J which displays results on a correlation and scale-invariant based loss.

277 Our NMRU is the only module with reasonable success over all tested ranges, requiring only $2I \times O$
 278 learnable parameters. However, this comes at the cost of the simplicity of the solution due to its
 279 exploitation of the identity rule; an issue the Real NPU does not have.

280 Once robust modules are attainable in a single layer setting, the next step would be to question
 281 performance when learning stacked modules, e.g. learning a stacked additive and multiplicative
 282 module. Previously, Madsen and Johansen [2020, Figure 2] illustrates the troubles for multiplicative
 283 models with the capacity for division. They show how a stacked summative-multiplicative module can
 284 lead to an exploding loss when the output of the summative module is close to 0 and the multiplicative
 285 model tries to divide. In Figure 7, we recreate their setup to produce the loss surfaces for the NAU-
 286 Real NPU³, NAU-NRU and NAU-NMRU respectively.⁴ We find a similar issue with the Real-NPU
 287 and NRU, as both these units use a weight range of $[-1,1]$. In contrast, the NMRU, whose weight's
 288 range is limited to $[0,1]$ does not have exploding losses.

289 In conclusion, division remains a challenge to learn using interpretable neural networks, even for the
 290 simplest tasks. Nevertheless, by identifying the specific areas causing difficulty (e.g., training ranges),
 291 and useful architecture properties (e.g., using a sign retrieval mechanism), we hope the community
 292 has better intuition for dealing with division and develop more robust modules to learn division.

³The NAU is a summative module [Madsen and Johansen, 2020].

⁴Appendix K displays larger versions of these plots.

293 References

- 294 Girish Chandrashekar and Ferat Sahin. A survey on feature selection methods. *Computers &*
295 *Electrical Engineering*, 40(1):16–28, 2014. ISSN 0045-7906. doi: <https://doi.org/10.1016/j.comp>
296 [eleceng.2013.11.024](https://www.sciencedirect.com/science/article/pii/S045790613003066). URL [https://www.sciencedirect.com/science/article/pii/S0](https://www.sciencedirect.com/science/article/pii/S045790613003066)
297 [45790613003066](https://www.sciencedirect.com/science/article/pii/S045790613003066). 40th-year commemorative issue.
- 298 Lukas Faber and Roger Wattenhofer. Neural status registers. *CoRR*, abs/2004.07085, 2020. URL
299 <https://arxiv.org/abs/2004.07085>.
- 300 Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural
301 networks. In *Proceedings of the thirteenth international conference on artificial intelligence and*
302 *statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010. URL [http:](http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf)
303 [//proceedings.mlr.press/v9/glorot10a/glorot10a.pdf](http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf).
- 304 Niklas Heim, Tomáš Pevný, and Václav Šmídl. Neural power units. *Advances in Neural Information*
305 *Processing Systems*, 33, 2020. URL [https://papers.nips.cc/paper/2020/file/48e5900](https://papers.nips.cc/paper/2020/file/48e59000d7dfcf6c1d96ce4a603ed738-Paper.pdf)
306 [0d7dfcf6c1d96ce4a603ed738-Paper.pdf](https://papers.nips.cc/paper/2020/file/48e59000d7dfcf6c1d96ce4a603ed738-Paper.pdf).
- 307 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*,
308 abs/1412.6980, 2015. URL <https://arxiv.org/pdf/1412.6980.pdf>.
- 309 Andreas Madsen and Alexander Rosenberg Johansen. Measuring arithmetic extrapolation perfor-
310 mance. In *Science meets Engineering of Deep Learning at 33rd Conference on Neural Information*
311 *Processing Systems (NeurIPS 2019)*, volume abs/1910.01888, Vancouver, Canada, October 2019.
312 URL <https://arxiv.org/pdf/1910.01888.pdf>.
- 313 Andreas Madsen and Alexander Rosenberg Johansen. Neural arithmetic units. In *International*
314 *Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=H1gN0eHKPS)
315 [H1gN0eHKPS](https://openreview.net/forum?id=H1gN0eHKPS).
- 316 Bhumika Mistry, Katayoun Farrahi, and Jonathon Hare. A primer for neural arithmetic logic modules,
317 2021. URL <https://arxiv.org/pdf/2101.09530.pdf>.
- 318 Subham Sahoo, Christoph Lampert, and Georg Martius. Learning equations for extrapolation and
319 control. In *International Conference on Machine Learning*, pages 4442–4450. PMLR, 2018.
- 320 Daniel Schlör, Markus Ring, and Andreas Hotho. inalu: Improved neural arithmetic logic unit.
321 *Frontiers in Artificial Intelligence*, 3:71, 2020. ISSN 2624-8212. doi: 10.3389/frai.2020.00071.
322 URL <https://www.frontiersin.org/article/10.3389/frai.2020.00071>.
- 323 Andrew Trask, Felix Hill, Scott E Reed, Jack Rae, Chris Dyer, and Phil Blunsom. Neural arithmetic
324 logic units. In *Advances in Neural Information Processing Systems*, pages 8035–8044, 2018. URL
325 <https://openreview.net/pdf?id=H1gN0eHKPS>.
- 326 Silviu-Marian Udrescu and Max Tegmark. Ai feynman: A physics-inspired method for symbolic
327 regression. *Science Advances*, 6(16), 2020. doi: 10.1126/sciadv.aay2631. URL [https:](https://advances.sciencemag.org/content/6/16/eaay2631)
328 [//advances.sciencemag.org/content/6/16/eaay2631](https://advances.sciencemag.org/content/6/16/eaay2631).

329 Checklist

- 330 1. For all authors...
- 331 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
332 contributions and scope? [Yes] See contributions in the introduction. For Real NPU
333 improvements see Section 5. For two novel modules see Section 3.2 and 3.3 and results
334 in Section 6. For hindrances in learning division, see Section 6.
- 335 (b) Did you describe the limitations of your work? [Yes] See Section 7.
- 336 (c) Did you discuss any potential negative societal impacts of your work? [Yes] This
337 paper focuses on the ML techniques and the foundational research required to learn
338 division in a systematic manner. Once robust modules for the arithmetic operations (i.e.
339 NALMs) are achievable the community will possess trainable modules with significant
340 advantages regarding model transparency and generalisability. That being said, we
341 discuss how this leads to a negative societal impact in the end of Section 1.

- 342 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
343 them? [Yes] We have read the guidelines and our work does not use human-derived
344 data.
- 345 2. If you are including theoretical results...
- 346 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
347 (b) Did you include complete proofs of all theoretical results? [N/A]
- 348 3. If you ran experiments...
- 349 (a) Did you include the code, data, and instructions needed to reproduce the main ex-
350 perimental results (either in the supplemental material or as a URL)? [Yes] See link
351 provided to the code base. Data is generated in real time and the code to generate it is
352 available in the linked repository.
- 353 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
354 were chosen)? [Yes] See Section 4 and Appendix C.
- 355 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
356 ments multiple times)? [Yes] See plots and experiment details in Section 4 for further
357 details.
- 358 (d) Did you include the total amount of compute and the type of resources used (e.g., type
359 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix D.
- 360 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 361 (a) If your work uses existing assets, did you cite the creators? [Yes] See footnote link
362 provided to the code base and Section 4.
- 363 (b) Did you mention the license of the assets? [Yes] We state the MIT licence when giving
364 the link to the codebase.
- 365 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
366 All code for model, data and experiments are available through the link to the codebase.
- 367 (d) Did you discuss whether and how consent was obtained from people whose data you're
368 using/curating? [N/A]
- 369 (e) Did you discuss whether the data you are using/curating contains personally identifiable
370 information or offensive content? [No] All data is synthetic, containing no such
371 information.
- 372 5. If you used crowdsourcing or conducted research with human subjects...
- 373 (a) Did you include the full text of instructions given to participants and screenshots, if
374 applicable? [N/A]
- 375 (b) Did you describe any potential participant risks, with links to Institutional Review
376 Board (IRB) approvals, if applicable? [N/A]
- 377 (c) Did you include the estimated hourly wage paid to participants and the total amount
378 spent on participant compensation? [N/A]