

1 A Properties of a Division Module

2 When building a division module, the following properties should be included:

3 **Ability to multiply:** Without multiplication the module is limited to expressing reciprocals.

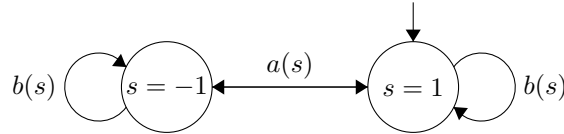
4 **Interpretable weights:** Having a discrete set of weight values to represent specific operations, e.g.,
 5 -1 to divide, 1 to multiply, 0 to not select. Doing this also has the additional benefit of producing
 6 generalisable solutions to out-of-bounds data.

7 **Calculating the output:** This can be decomposed into three tasks: input selection, magnitude
 8 calculation, and sign calculation.

9 **Magnitude:** This is achieved using discrete weights. The Real NPU and NRU use -1 for reciprocals,
 10 1 for multiplication. The NMRU uses 1 for selecting an input element which represents either a
 11 multiplication or reciprocal depending on its position.

12 **Sign of the output:** Calculating the sign value (1/-1) of the output can occur at an element level in
 13 which the sign is calculated for each intermediary value as each input element is being processed, or
 14 at the higher input level in which the sign is calculated separately for the magnitude and then applied
 15 once the final output magnitude is calculated. The NRU uses the prior method while the Real NPU
 16 and NMRU use the latter method. If an input is 0 or considered irrelevant then the output sign will be
 17 1. (Ablation studies on the NMRU, Figure 7, suggest the latter option which separately calculates the
 18 sign to be more beneficial).

19 The Real NPU and NMRU use the cosine function to calculate the final sign of the module's output
 20 neuron. Below shows the state diagram of how the sign value (i.e. the state) of the output would
 21 change depending on the inputs and relevant parameters being processed. We only consider the
 22 discrete parameters for simplicity. Both the Real NPU and NMRU use the same state diagram but
 23 have different conditions for a state transition to occur.



24

25 The conditions for the Real NPU transition functions $a(s) = -s$ and $b(s) = s$, where s is the state
 26 value -1, or 1, are defined as follows:

$$a(s) : x_i < 0 \wedge w_{i,o} \in \{-1, 1\} \wedge g_i = 1 ,$$

$$b(s) : x_i \geq 0 \vee w_{i,o} = 0 \vee g_i = 0 .$$

27 Transitioning from one sign to another only occurs if the input element (x_i) is negative and is
 28 considered relevant i.e. the gate (g_i) and weight value ($w_{i,o}$) is non-0. In contrast, to remain at a state
 29 requires either the input element to be ≥ 0 or not be considered relevant.

30 The conditions for the NMRU transition functions $a(s) = -s$ and $b(s) = s$, where s is the state value
 31 -1, or 1, are defined as follows:

$$a(s) : x_i < 0 \wedge w_{i,o} = 1 ,$$

$$b(s) : x_i \geq 0 \vee w_{i,o} = 0 .$$

32 Transitioning from one sign to another only occurs if the input element (x_i) is negative and is
 33 considered relevant i.e. the weight value ($w_{i,o}$) is 1. To remain at a state requires either the input
 34 element to be ≥ 0 or the weight value to not select the input.

35 **Selection:** Not all inputs are relevant for the output value. To process any irrelevant input elements
 36 can be interpreted as converting to the identity value of multiplication/division (=1). The identity
 37 property means that any value multiplied/divided by the identity value remains at the original number.
 38 Hence, irrelevant inputs are converted into 1 (rather than being masked out to 0). For the multiplication
 39 case, this stops the output becoming 0, and for division it avoids the divide by 0 case. For all the
 40 explored modules, a weight value of 0 will deal with the irrelevant input case. However, the Real
 41 NPU goes a step further by also having an additional gate vector with the purpose of learning to
 42 select relevant inputs. Such gating has been proven to be helpful for an NPU based module [Heim
 43 et al., 2020], but may not be necessary when dealing with weights between [0,1] like in the NMRU
 44 (see Appendix H).

45 **B Neural Addition and Neural Multiplication Units' (NAU & NMU)**

46 Madsen and Johansen [2020] develop two modules: one for dealing with addition and subtraction
 47 (the NAU) and the other for multiplication (the NMU). NAU output element a_o is defined as

$$\text{NAU} : a_o = \sum_{i=1}^I (W_{i,o} \cdot x_i) \quad (1)$$

48 where I is the number of inputs. The NMU output element m_o is defined as

$$\text{NMU} : m_o = \prod_{i=1}^I (W_{i,o} \cdot x_i + 1 - W_{i,o}). \quad (2)$$

49 Before passing a input through a module, the weight matrix is clamped to [-1,1] for the NAU or
 50 [0,1] for the NMU. Weights are ideally discrete values, where the NAU is 0, 1, or -1, representing no
 51 selection, addition and subtraction, and the NMU is 0 or 1, representing no selection and multiplication.
 52 To enforce discretisation of weights both units have a regularisation penalty for a given period of
 53 training. The penalty is

$$\lambda \cdot \frac{1}{I \cdot O} \sum_{o=1}^O \sum_{i=1}^I \min(|W_{i,o}|, 1 - |W_{i,o}|), \quad (3)$$

54 where O is the number of outputs and λ is defined as

$$\lambda = \hat{\lambda} \cdot \max\left(\min\left(\frac{\text{iteration}_i - \lambda_{start}}{\lambda_{end} - \lambda_{start}}, 1\right), 0\right). \quad (4)$$

55 Regularisation strength is scaled by a predefined $\hat{\lambda}$. The regularisation will grow from 0 to $\hat{\lambda}$ between
 56 iterations λ_{start} and λ_{end} , after which it plateaus and remains at $\hat{\lambda}$.

57 **C Experiment Parameters**

58 Tables 1 and 2 for the breakdown of parameters used in the Single Module Tasks. Table 3 gives the
 59 interpolation and extrapolation ranges used in the mixed-sign datasets tasks.

Table 1: Parameters which are applied to all modules. Parameters have been split based on the experiment. *Validation and test datasets generate one batch of samples at the start which gets used for evaluation for all iterations. † the Real NPU modules use a value of 1.

Parameter	Without redundancy	With redundancy
Layers	1	1
Input size	2	10
Total iterations	50,000	100,000
Train samples	128 per batch	128 per batch
Validation samples*	10000	10000
Test samples*	10000	10000
Seeds	25	25
Optimiser	Adam (with default parameters)	Adam (with default parameters)
$\hat{\lambda}^\dagger$	10	10

Table 2: Parameters specific to the Real NPU modules for the Single Module Tasks.

Parameter	Value
$(\beta_{start}, \beta_{end})$	(1e-9, 1e-7)
β_{growth}	10
β_{step}	10000
$\hat{\lambda}$	1

Table 3: Mixed-Sign Datasets: The interpolation and extrapolation ranges to sample the two input elements for a single data sample. The target expression to learn is: input 1 \div input 2.

DATASET	INTERPOLATION		EXTRAPOLATION	
	INPUT 1	INPUT 2	INPUT 1	INPUT 2
1	U[-2, -0.1)	U[0.1, 2)	U[-6, -2)	U[2, 6)
2	U[-2, -1)	U[1, 2)	U[-6, -2)	U[2, 6)
3	U[-2, 2)	U[-2, 2)	U[-6, -2)	U[2, 6)
4	U[0.1, 2)	U[-2, -0.1)	U[2, 6)	U[-6, -2)
5	U[1, -2)	U[-2, -1)	U[2, 6)	U[-6, -2)

60 **C.1 Parameter Initialisation**

61 We give the initialisations used on the different module parameters:

62 **Real NPU:** The real weight matrix uses the Pytorch’s Xavier Uniform initialisation. The gate vector
 63 initialises all values to 0.5. (This is the same initialisation used in Heim et al. [2020].)

64 **NPU:** The imaginary weight matrix is initialised to 0. The rest of the parameters are initialised same
 65 as the Real NPU. (This is the same initialisation used in Heim et al. [2020].)

66 **NRU:** The weight matrix uses a Xavier Uniform initialisation which can have a maximum range
 67 between -0.5 to 0.5 (depending on the network sizes). (This is the same initialisation the Neural
 68 Addition Unit uses [Madsen and Johansen, 2020].)

69 **NMRU:** The weight matrix uses a Uniform initialisation which can have a maximum range between
 70 0.25 to 0.75 (depending on the network sizes). (This is the same initialisation the Neural Multiplication
 71 unit uses [Madsen and Johansen, 2020].)

72 **D Hardware and Time to Run Experiments**

73 All experiments were trained on the CPU, as training on GPUs takes considerably longer. All Real
74 NPU experiments were run on XXX (the University of XXX 's supercomputer), where a compute
75 node has 40 CPUs with 192 GB of DDR4 memory which uses dual 2.0 GHz Intel Skylake processors.
76 All NRU and NMRU experiments were run on a 16 core CPU server with 125 GB memory 1.2 GHz
77 processors.

78 Table 4 displays time taken for each experiment to run a single seed for a single range. Timings are
79 based on a single run rather than the runtime of a script execution because the queuing time from
80 jobs when executing scripts is not relevant to the experiment timings. For a single model, a single
81 experiment would have 225 runs (for 9 training ranges and 25 seeds).

Table 4: Timings of experiments.

Experiment	Model	Approximate time for completing 1 seed (mm:ss)
No redundancy (size 2)	Real NPU	03:20
	NRU	02:00
	NMRU	03:00
With redundancy (size 10)	Real NPU	05:30
	NRU	05:00
	NMRU	05:15

82 **E NRU on the Single Module Task (no redundancy): Effect of Learning**
 83 **Rate**

84 Figure 1 displays the effect of different learning rates for the NRU. An learning rate of 1 gets full
 85 success on all ranges with performance deteriorating as the learning rate reduces.

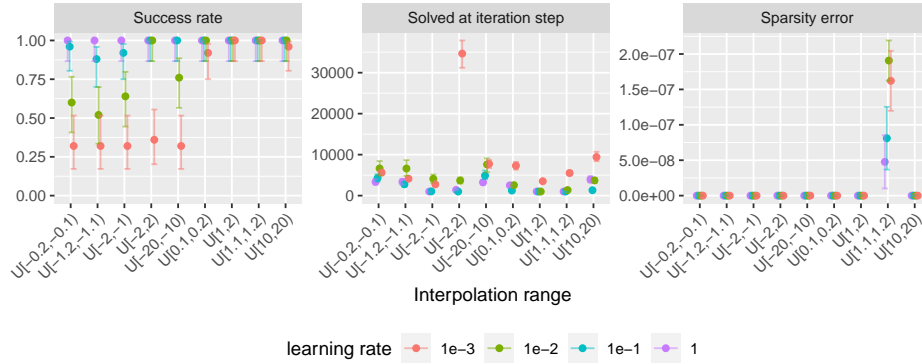


Figure 1: Different learning rates on the NRU for the Single Module Task (no redundancy)

86 **F Division by small values: Experimental Results**

87 This section shows the results on trying to learn the reciprocal/division of values close to zero using
 88 the Real NPU, NRU and NMRU. We train and test on the ranges where the lowest bound is 0 and
 89 the upper bounds are: 1e-4, 1e-3, 1e-2, 1e-1 and 1. Unless stated otherwise, the hyperparameters
 90 of a model are set to what is used for the Single Layer Task without redundancy. The first task runs
 91 for 5,000 iterations with no regularisation for any module. The second and third tasks both run for
 92 50,000 iterations.

93 Due to precision errors, a solution with the ideal parameters will not evaluate to a MSE of 0. Therefore,
 94 we calculate thresholds which the test MSE should be within. A threshold value for a task is calculated
 95 from evaluating the MSE of each range’s test dataset for each module, using the ‘golden’ weight
 96 values and adding an epsilon term¹ to the resulting error which takes into account precision errors.
 97 All experiments are run using 32-bit precision.

98 In general, successful runs take longer to solve as the input ranges become smaller. The simplest
 99 task, of taking the reciprocal when the input size is 1 (Figure 2) is achieved with ease for all modules,
 though for $\mathcal{U}[0, 1e-4]$, we find the NRU begins to start struggling.

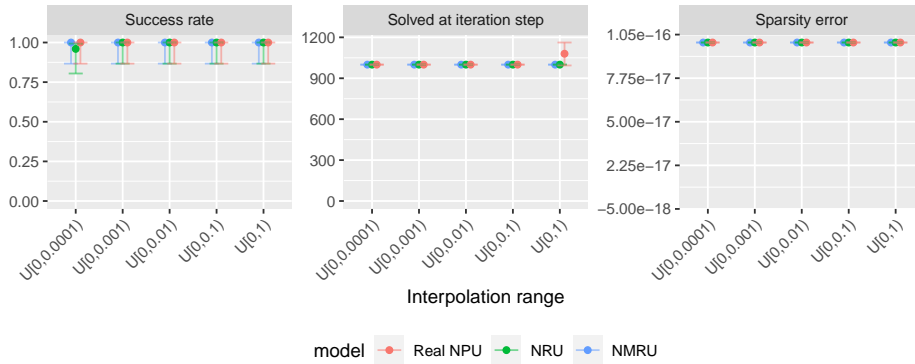


Figure 2: Input: $[a]$, output $\frac{1}{a}$. Learns reciprocal when there is no input redundancy.

100

101 Introducing a redundant input (Figure 3) greatly impacts performance with only the NMRU able to
 102 achieve reasonable success for the larger ranges. The successes shown for the Real NPU at range
 103 $\mathcal{U}[0, 1e-4]$ are false positives caused by the ϵ in the architecture used for stability. Test false positives
 can also be indicated by the high sparsity error of the weights.

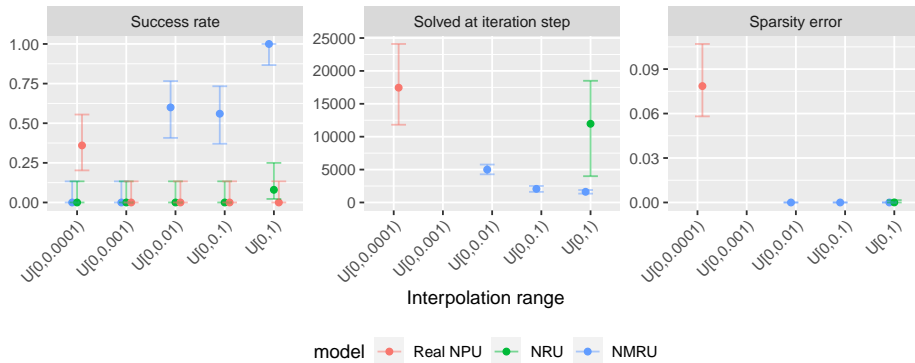


Figure 3: Input: $[a, b]$, output $\frac{1}{a}$. Learns reciprocal of the first input when there is redundancy.

104

105 Modifying the task to division (Figure 4), meaning the redundant input is now relevant, shows
 106 improvement for the NMRU and NRU for the larger ranges.

¹The term is the pytorch default eps value, torch.finfo().eps

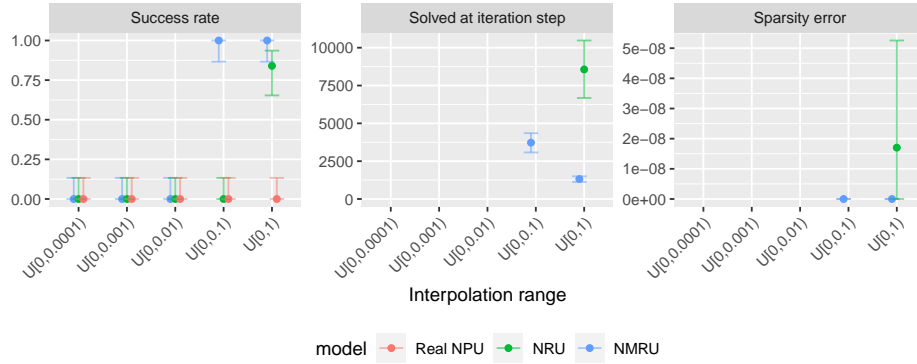


Figure 4: Input: $[a,b]$, output $\frac{a}{b}$. Learns division of the first and second value when there is no redundancy.

107 **G Real NPU; Single Module Task (with Redundancy): Additional**
 108 **Experiments**

Figure 5 shows results of using the NPU for the task with redundancy.

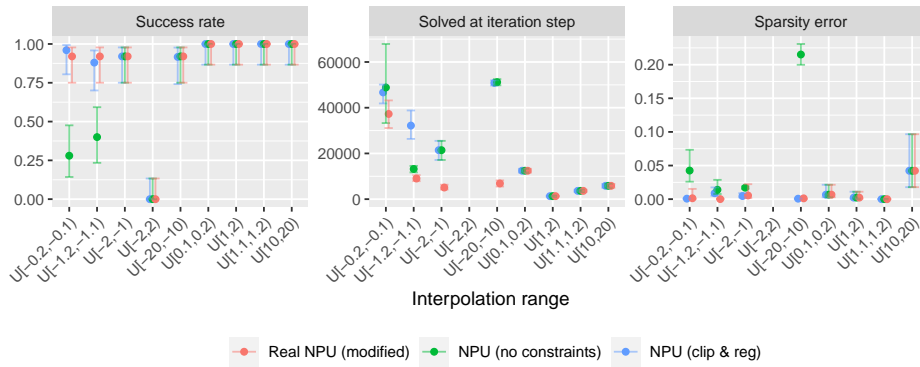


Figure 5: Adapting the Real NPU to use complex weights (NPU) on the Single Module Task with redundancy. Compares the NPU architecture with the Real NPU modifications (i.e. NPU (no constraints)) and the same model but with the imaginary weights clipped to $[-1,1]$ and L1 sparsity regularisation on the complex weights (i.e. NPU (clip & reg)).

109

110 Figure 6 shows how modifying the weight discretisation to not penalise weights at 0 does not effect
 111 success.

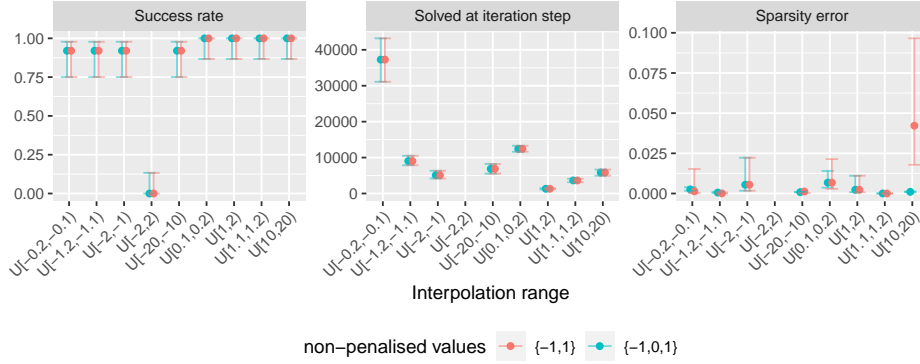


Figure 6: Comparing weight discretisation on the NPU weights which penalises not having weight of $\{-1, 1\}$ vs $\{-1, 0, 1\}$.

112 H NMRU; Single Module Task with Redundancy (Additional Experiments)

113 This section further explores the NMRU architecture.

114 Figure 7 shows an ablation study on different components of the NMRU architecture. Removing both
 115 the sign retrieval and grad norm clipping performs poorly over a majority of ranges (including positive ranges). Gradient norm clipping alone is unable to solve the issue in learning negative ranges, however

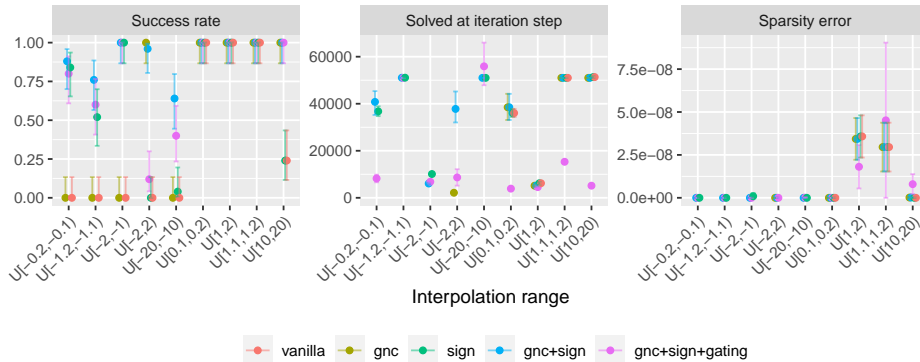


Figure 7: Ablation study for the NMRU.

116 fully succeeds on the $\mathcal{U}[-2,2]$ range. Using the sign retrieval without the gradient clipping gains
 117 successes for the negative ranges, though performance on $\mathcal{U}[2,-2]$ is effected. However, including
 118 both gradient clipping and sign retrieval results in separating the calculation of the magnitude of
 119 the output and its sign while having reasonable gradients, gaining the most improvement over the
 120 vanilla NMRU. Further including a learnable gate vector (like the Real NPU), which is applied to
 121 the input vector, hinders performance. The largest solved at iteration step seems to be bounded at
 122 approximately 50,000 iterations which correlates to the point at which the sparsity regularisation
 123 begins, which highlights the importance of discretisation. Even with the different ablations, the
 124 sparsity errors of the successful seeds remain extremely low (which is not always the case for the
 125 Real NPU (see Figure 6)).
 126

127 Figure 8 shows the effect of using different learning rates on the NMRU (with grad norm clipping
 128 and sign retrieval) using an Adam optimiser. Too low a learning rate struggles on the mixed-sign
 129 range $\mathcal{U}[-2,2]$. Too high a learning leads to no success on multiple ranges.

130 Figure 9 compares training the NMRU with either an Adam and SGD optimiser. As expected,
 131 Adam outperforms SGD in all ranges (except two, where both perform equally). This difference in
 132 performance can be accounted for by Adam’s ability to scale the step size of each weight, which can
 133 compliment the clipped gradient norm of the NMRU, in contrast to the SGD’s global step size.

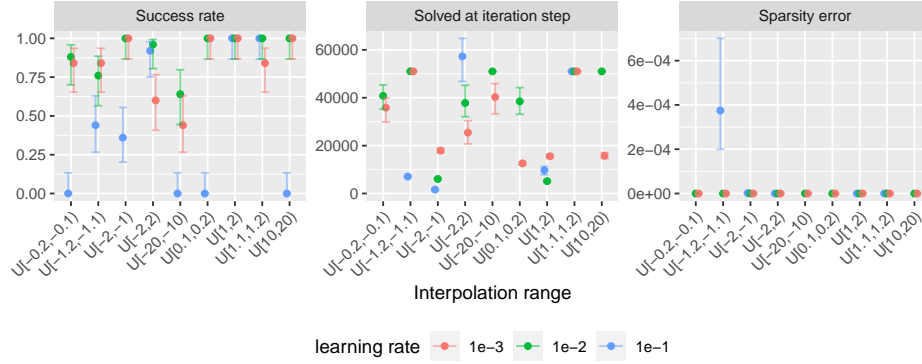


Figure 8: Effect of different learning rates on the NMRU

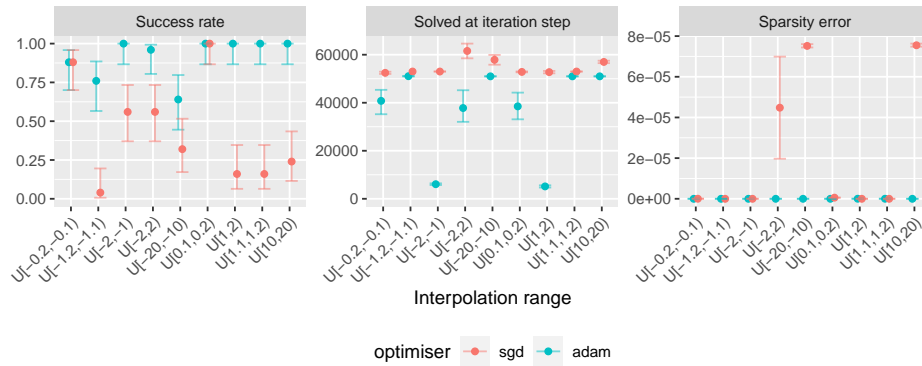


Figure 9: Effect of optimiser on the NMRU. SGD = Stochastic Gradient Descent.

134 **I NRU; Single Module Task (with Redundancy): Calculating the Sign**
 135 **Separately**

136 The ‘separate NRU’ module calculates the magnitude and sign separately and then combines them
 137 using multiplication together once all input elements are accounted for. The following definition is
 138 used to calculate a NRU with separate magnitude and sign calculation,

$$z_o = \prod_{i=1}^I (|x_i|^{W_{i,o}} \cdot |W_{i,o}| + 1 - |W_{i,o}|) \cdot \prod_{i=1}^I \text{sign}(x_i)^{\text{round}(W_{i,o})}. \quad (5)$$

139 Figure 10 shows results, where the separate sign method shows no difference in success to the original
 140 NRU architecture.

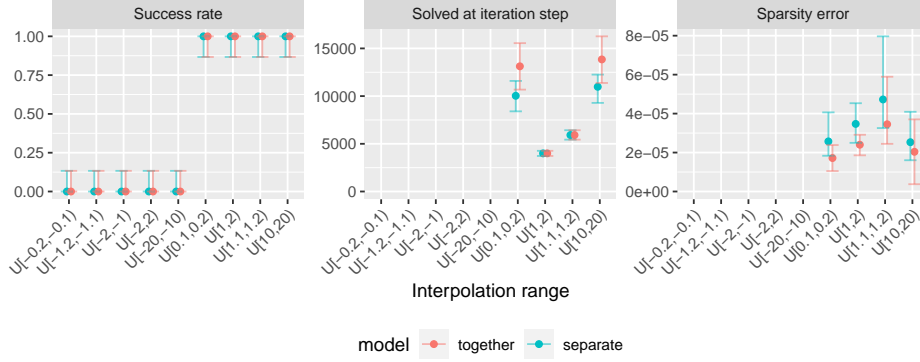


Figure 10: NRU on the redundancy experiment comparing a module which calculates the magnitude and sign together vs calculating the magnitude and sign separately and then combining them.

141 J Effect of Different Losses on the Single Module Task (with Redundancy)

Table 5: The properties of different loss functions.

	MSE	PCC	MAPE
Batch mean	✓	✓	✓
Standardisation		✓	✓
Difference of prediction from target	✓		✓
Projection		✓	
Mean centering		✓	

142 Different losses induce different loss landscapes impacting the areas of success for a module. We
 143 explore the effects of three different losses including the MSE, Pearson’s Correlation Coefficient
 144 (Equation 7), and the Mean Absolute Precision Error (Equation 8). We use the division task with 10
 145 inputs. The properties of each loss is summarised in Table 5. All experiment parameters match the
 146 original MSE runs in the main experiments. The only difference is the loss used.

$$\begin{aligned}
 v_{x,i} &= (\hat{y}_i - \bar{\hat{y}}), & s_x &= \sqrt{\text{clamp}\left(\frac{1}{N} \sum_i v_{x,i}^2, \epsilon\right)} \\
 v_{y,i} &= (y_i - \bar{y}), & s_y &= \sqrt{\text{clamp}\left(\frac{1}{N} \sum_i v_{y,i}^2, \epsilon\right)}
 \end{aligned} \tag{6}$$

$$r = \frac{1}{N} \sum_i \left(\frac{v_{x,i}}{s_x + \epsilon} \cdot \frac{v_{y,i}}{s_y + \epsilon} \right)$$

147

$$\text{pcc loss} := 1 - r \tag{7}$$

148 where N is the batch size, and the means ($\bar{\hat{y}}$ and \bar{y}) are taken over the batch. ϵ is used to provide better
 149 numerical stability.

$$\text{mape loss} := \frac{1}{N} \sum_i \left(\frac{|y_i - \hat{y}_i|}{y_i} \right) \tag{8}$$

150 **Real NPU (Figure 11)** Both the Real NPU and MAPE are able to get success on the $\mathcal{U}[-2,2]$
 151 range, which the MSE completely fails on, implying that having a loss with standardisation is useful.
 152 However, in order to gain successes in the mixed-sign range, the other negative ranges have reduced

153 in success for both PCC and MAPE. Both speed and sparsity retain similar performance to MSE in a
 154 majority of cases, with PCC solving especially fast for all tested ranges.

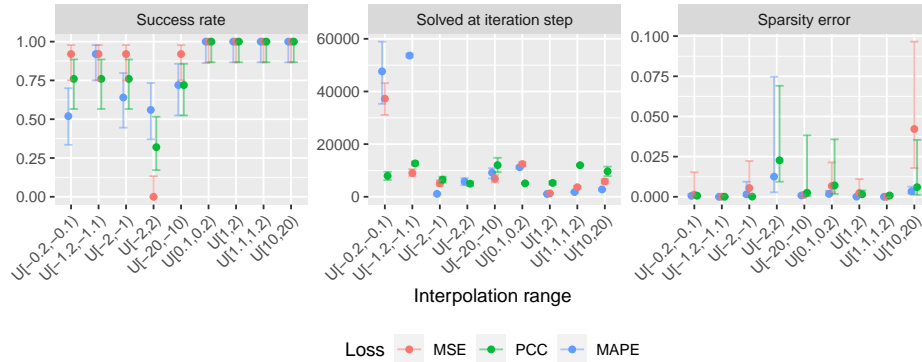


Figure 11: Single Module Task with redundancy on the Real NPU, comparing different loss functions.

154

155 **NRU (Figure 12)** Different losses have little effect on the NRU. All three losses perform well on
 156 the positive ranges. Compared to the Real NPU, the PCC loss on the NRU takes longer to converge
 157 to a success for negative ranges.

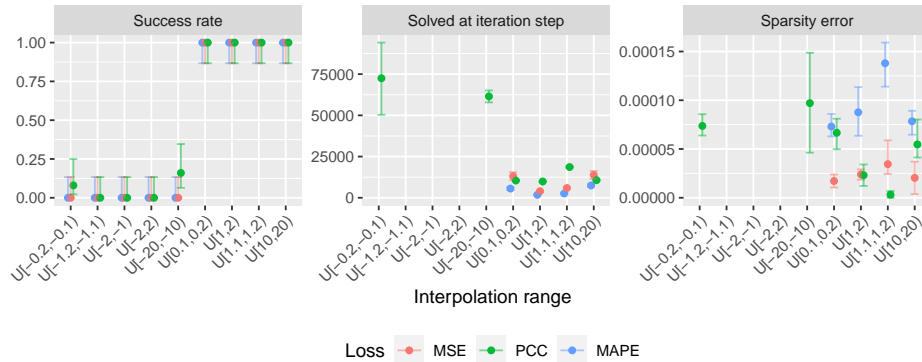


Figure 12: Single Module Task with redundancy on the NRU, comparing different loss functions.

157

158 **NMRU (Figure 13)** All three losses perform reasonably well, with the PCC struggling the most.
 159 Unlike the other units, $\mathcal{U}[-20,-10]$ causes the most trouble, whereas $\mathcal{U}[-2,2]$ gains near to full success
 160 on two of the three losses.

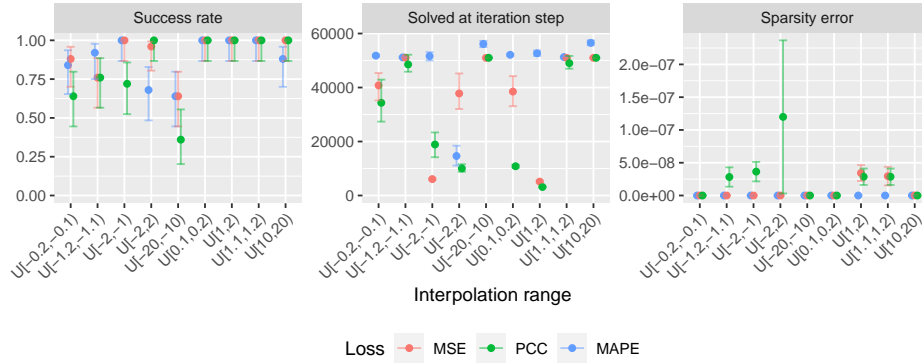


Figure 13: Single Module Task with redundancy on the NMRU, comparing different loss functions.

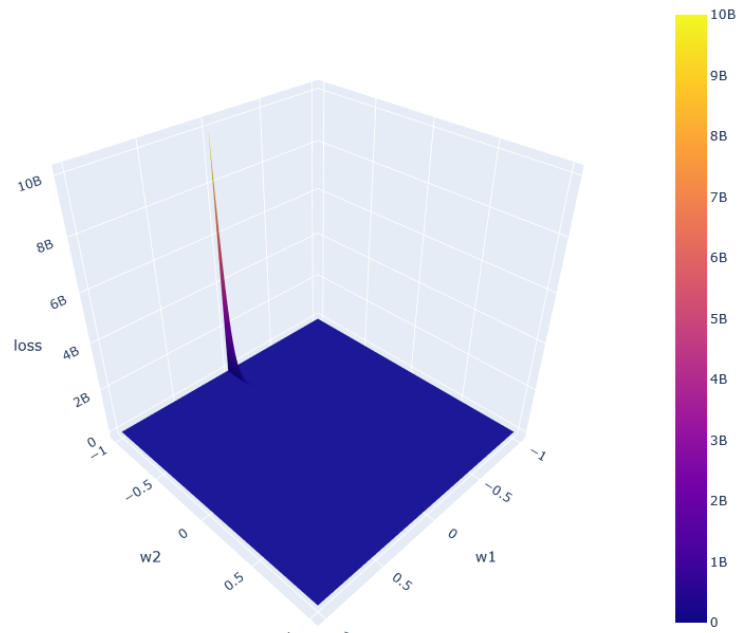
161 K RMSE Loss Landscapes

162 For clarity, we show bigger versions of each subplot from Figure 7.

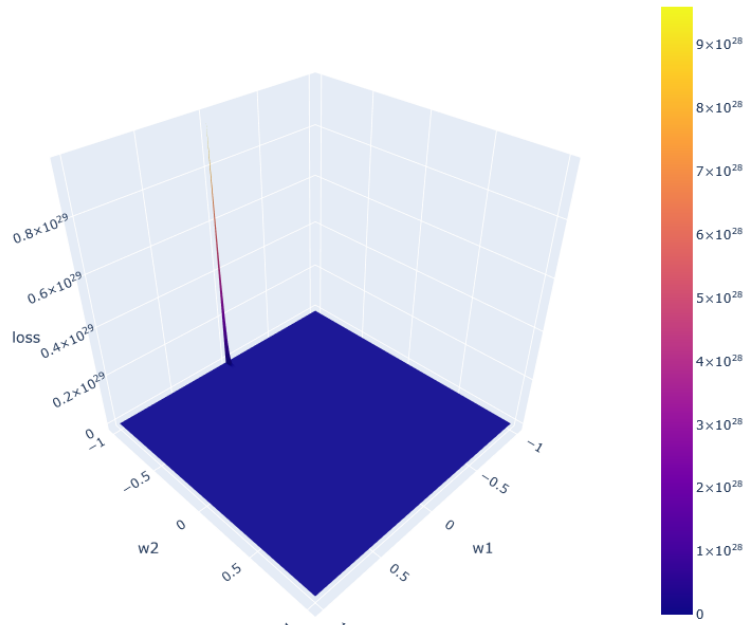
163 References

164 Niklas Heim, Tomáš Pevný, and Václav Šmídl. Neural power units. *Advances in Neural Infor-*
 165 *mation Processing Systems*, 33, 2020. URL [https://papers.nips.cc/paper/2020/file/](https://papers.nips.cc/paper/2020/file/48e59000d7dfcf6c1d96ce4a603ed738-Paper.pdf)
 166 [48e59000d7dfcf6c1d96ce4a603ed738-Paper.pdf](https://papers.nips.cc/paper/2020/file/48e59000d7dfcf6c1d96ce4a603ed738-Paper.pdf).

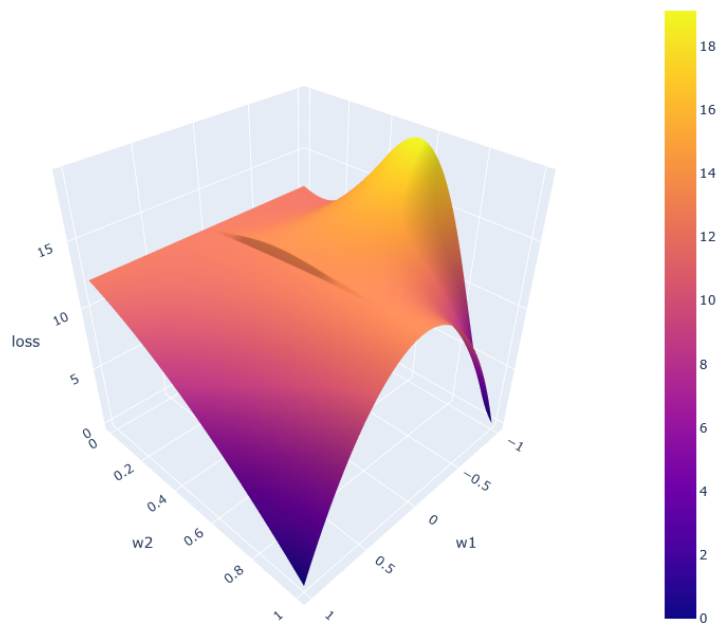
167 Andreas Madsen and Alexander Rosenberg Johansen. Neural arithmetic units. In *International*
 168 *Conference on Learning Representations*, 2020. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=H1gNOeHKPS)
 169 [H1gNOeHKPS](https://openreview.net/forum?id=H1gNOeHKPS).



(a) NAU-Real NPU (where $\epsilon = 1e - 5$)



(b) NAU-NRU



(c) NAU-NMRU

Figure 14: Enlarged loss landscapes of different stacked summative-multiplicative units.