

Neural Arithmetic Units

By **Andreas Madsen** and **Alexander Rosenberg Johansen**
@andreas_madsen @AlexRoseJo

ICLR 2020, spotlight awarded paper

paper: <https://openreview.net/forum?id=H1gNOeHKPS>

code: <https://github.com/AndreasMadsen/stable-nalu>

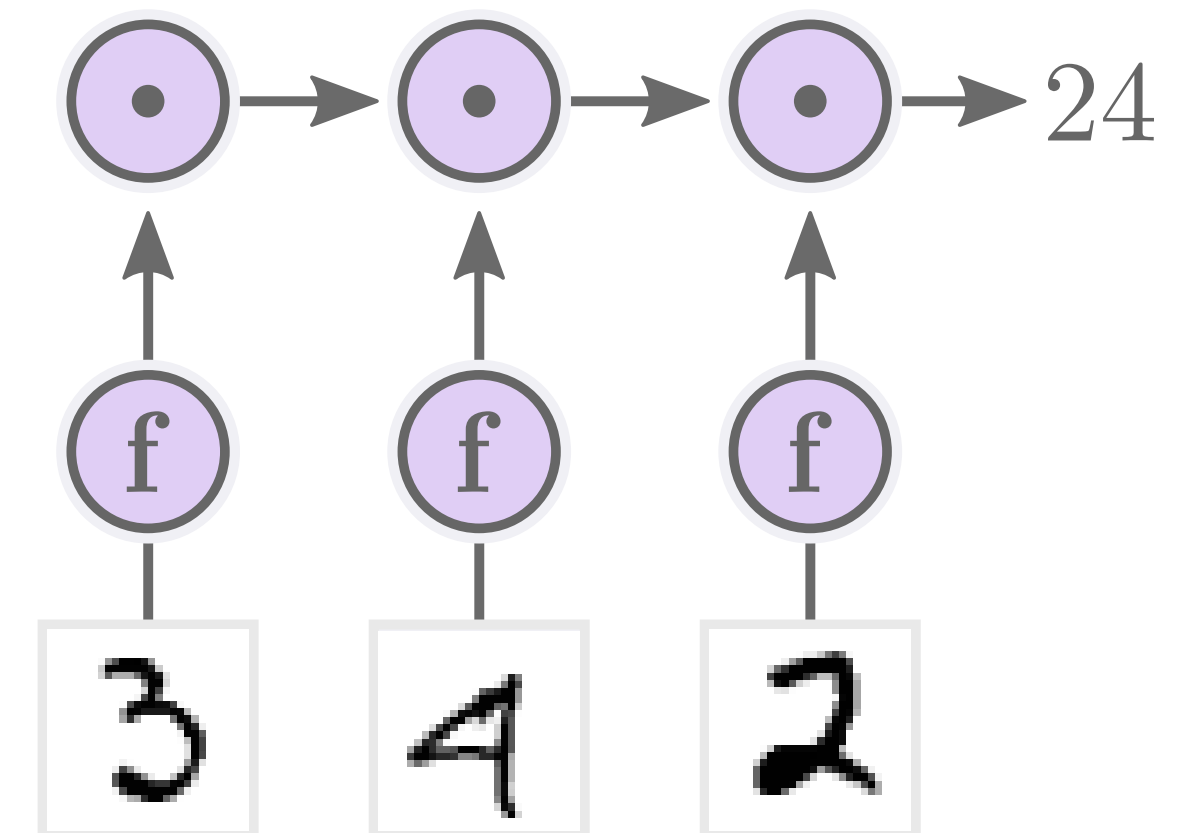
Arithmetic Extrapolation

- Neural networks are great at interpolation but can rarely extrapolate.
- Arithmetic extrapolation assumes there is an underlying function partially composed of simple arithmetics.
- Simple arithmetics might occur in:
 - Physical modelling
 - Financial modelling
 - NLP Q&A tasks

Direct arithmetic example:

$$t = (x_1 + x_2) \cdot (x_1 + x_2 + x_3 + x_4) \text{ for } x \in \mathbb{R}^4$$

MNIST example:



NALU by Andrew Trask, et. al

- Central idea: Learn the underlying function, by being able to represent it exactly and having constrained (and biased) weights.
- NALU has 3 components¹:
 - Addition/Subtraction component
 - Multiplication/Division component
 - Gating component

σ controls scale (0, 1)
 \tanh controls sign (-1, 1)

$$W_{h_\ell, h_{\ell-1}} = \tanh(\hat{W}_{h_\ell, h_{\ell-1}}) \sigma(\hat{M}_{h_\ell, h_{\ell-1}})$$

to avoid issues with $\log(\cdot)$

$$\text{NAC}_+ : z_{h_\ell} = \sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h_{\ell-1}} z_{h_{\ell-1}}$$

$$\text{NAC}_\bullet : z_{h_\ell} = \exp \left(\sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_\ell, h_{\ell-1}} \log(|z_{h_{\ell-1}}| + \epsilon) \right)$$

$$g_{h_\ell} = \sum_{h_{\ell-1}=1}^{H_{\ell-1}} G_{h_\ell, h_{\ell-1}} z_{h_{\ell-1}}$$

$$\text{NALU} : z_{h_\ell} = g_{h_\ell} \odot \text{NAC}_+(\mathbf{z}_{\ell-1})_{h_\ell} + (1 - g_{h_\ell}) \odot \text{NAC}_\bullet(\mathbf{z}_{\ell-1})_{h_\ell}$$

gating choose which operation to use

¹) Neural Arithmetic Logic Units (NALU), by Andrew Trask et. al, NeurIPS 2018.

Measuring Arithmetic Extrapolation

To extrapolate in an arithmetic task, the exact solution must be found.

In our NeurIPS workshop paper, we therefore argue consistency is a primary concern ².

Direct arithmetic example:

$$t = (x_1 + x_2) \cdot (x_1 + x_2 + x_3 + x_4) \text{ for } x \in \mathbb{R}^4$$

Op	Model	Success	Solved at iteration step		Sparsity error
		Rate	Median	Mean	Mean
×	NAC•	13% ^{+8%} _{-5%}	$5.5 \cdot 10^4$	$5.9 \cdot 10^4$ ^{+7.8·10³} _{-6.6·10³}	$7.5 \cdot 10^{-6}$ ^{+2.0·10⁻⁶} _{-2.0·10⁻⁶}
	NALU	26% ^{+9%} _{-8%}	$7.0 \cdot 10^4$	$7.8 \cdot 10^4$ ^{+6.2·10³} _{-8.6·10³}	$9.2 \cdot 10^{-6}$ ^{+1.7·10⁻⁶} _{-1.7·10⁻⁶}
	NMU	94% ^{+3%} _{-6%}	$1.4 \cdot 10^4$	$1.4 \cdot 10^4$ ^{+2.2·10²} _{-2.1·10²}	$2.6 \cdot 10^{-8}$ ^{+6.4·10⁻⁹} _{-6.4·10⁻⁹}

*2) Measuring Arithmetic Extrapolation Performance by
Andreas Madsen and Alexander R. Johansen. SEDL at
NeurIPS 2019.*

Analysis of issues

- weight issues:
 - gradient w.r.t. **M** is expected to be zero
 - assumed bias towards (-1, 0, 1) does not exist
- Multiplication issues:
 - singularities for $w < 0$ in **NAC**.
 - **NAC** can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

Weight issues

- weight issues:
 - gradient w.r.t. **M** is expected to be zero
 - assumed bias towards (-1, 0, 1) does not exist
- Multiplication issues:
 - singularities for $w < 0$ in **NAC**.
 - **NAC** can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

Expectation of the gradient w.r.t. **M**:

$$E \left[\frac{\partial \mathcal{L}}{\partial \hat{M}_{h_{\ell-1}, h_{\ell}}} \right] = E \left[\frac{\partial \mathcal{L}}{\partial W_{h_{\ell-1}, h_{\ell}}} \right] E \left[\tanh(\hat{W}_{h_{\ell-1}, h_{\ell}}) \right] E \left[\sigma'(\hat{M}_{h_{\ell-1}, h_{\ell}}) \right] = 0$$

Weight issues

- weight issues:
 - gradient w.r.t. **M** is expected to be zero
 - assumed bias towards (-1, 0, 1) does not exist
- Multiplication issues:
 - singularities for $w < 0$ in **NAC**.
 - NAC** can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

Expectation of the gradient w.r.t. **M**:

$$E \left[\frac{\partial \mathcal{L}}{\partial \hat{M}_{h_{\ell-1}, h_{\ell}}} \right] = E \left[\frac{\partial \mathcal{L}}{\partial W_{h_{\ell-1}, h_{\ell}}} \right] E \left[\tanh(\hat{W}_{h_{\ell-1}, h_{\ell}}) \right] E \left[\sigma'(\hat{M}_{h_{\ell-1}, h_{\ell}}) \right] = 0$$

Evaluate sparsity error on valid solutions:

$$E_{\text{sparsity}} = \max_{h_{\ell-1}, h_{\ell}} \min(|W_{h_{\ell-1}, h_{\ell}}|, |1 - |W_{h_{\ell-1}, h_{\ell}}||)$$

Op	Model	Success	Solved at iteration step		Sparsity error
		Rate	Median	Mean	Mean
+	NAC ₊	100% ^{+0%} _{-4%}	$2.5 \cdot 10^5$	$4.9 \cdot 10^5$ ^{+5.2·10⁴} _{-4.5·10⁴}	$2.3 \cdot 10^{-1}$ ^{+6.5·10⁻³} _{-6.5·10⁻³}
	Linear	100% ^{+0%} _{-4%}	$6.1 \cdot 10^4$	6.3 · 10⁴ ^{+2.5·10³} _{-3.3·10³}	$2.5 \cdot 10^{-1}$ ^{+3.6·10⁻⁴} _{-3.6·10⁻⁴}
	NALU	14% ^{+8%} _{-5%}	$1.5 \cdot 10^6$	$1.6 \cdot 10^6$ ^{+3.8·10⁵} _{-3.3·10⁵}	$1.7 \cdot 10^{-1}$ ^{+2.7·10⁻²} _{-2.5·10⁻²}
	NAU	100% ^{+0%} _{-4%}	1.8 · 10⁴	$3.9 \cdot 10^5$ ^{+4.5·10⁴} _{-3.7·10⁴}	3.2 · 10⁻⁵ ^{+1.3·10⁻⁵} _{-1.3·10⁻⁵}
-	NAC ₊	100% ^{+0%} _{-4%}	$9.0 \cdot 10^3$	$3.7 \cdot 10^5$ ^{+3.8·10⁴} _{-3.8·10⁴}	$2.3 \cdot 10^{-1}$ ^{+5.4·10⁻³} _{-5.4·10⁻³}
	Linear	7% ^{+7%} _{-4%}	$3.3 \cdot 10^6$	$1.4 \cdot 10^6$ ^{+7.0·10⁵} _{-6.1·10⁵}	$1.8 \cdot 10^{-1}$ ^{+7.2·10⁻²} _{-5.8·10⁻²}
	NALU	14% ^{+8%} _{-5%}	$1.9 \cdot 10^6$	$1.9 \cdot 10^6$ ^{+4.4·10⁵} _{-4.5·10⁵}	$2.1 \cdot 10^{-1}$ ^{+2.2·10⁻²} _{-2.2·10⁻²}
	NAU	100% ^{+0%} _{-4%}	5.0 · 10³	1.6 · 10⁵ ^{+1.7·10⁴} _{-1.6·10⁴}	6.6 · 10⁻² ^{+2.5·10⁻²} _{-1.9·10⁻²}

Neural Addition Unit

- weight issues:
 - gradient w.r.t. **M** is expected to be zero
 - assumed bias towards (-1, 0, 1) does not exist
- Multiplication issues:
 - singularities for $w < 0$ in **NAC**.
 - **NAC** can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

Our solution, **N**eural **A**ddition **U**nit:

$$\begin{aligned}
 & \text{clipped linear weights} \\
 W_{h_{\ell-1}, h_{\ell}} &= \min(\max(W_{h_{\ell-1}, h_{\ell}}, -1), 1), \\
 \mathcal{R}_{\ell, \text{sparse}} &= \frac{1}{H_{\ell} \cdot H_{\ell-1}} \sum_{h_{\ell}=1}^{H_{\ell}} \sum_{h_{\ell-1}=1}^{H_{\ell-1}} \min(|W_{h_{\ell-1}, h_{\ell}}|, 1 - |W_{h_{\ell-1}, h_{\ell}}|) \\
 & \text{sparsity regularizer} \\
 \text{NAU : } z_{h_{\ell}} &= \sum_{h_{\ell-1}=1}^{H_{\ell-1}} W_{h_{\ell}, h_{\ell-1}} z_{h_{\ell-1}}
 \end{aligned}$$

Multiplication issues

- weight issues:
 - gradient w.r.t. \mathbf{M} is expected to be zero
 - assumed bias towards $(-1, 0, 1)$ does not exist
- Multiplication issues:
 - singularities for $w < 0$ in NAC_\bullet .
 - NAC_\bullet can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

Division causes a singularity in the loss curvature

singularity

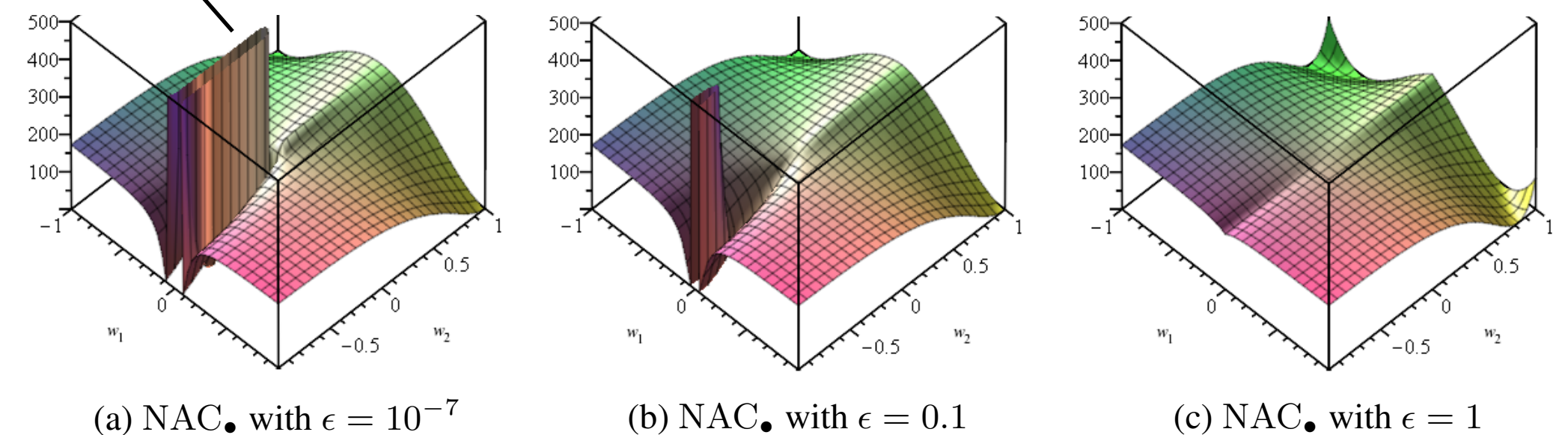


Figure 2: RMS loss curvature for a NAC_+ unit followed by a NAC_\bullet . The weight matrices are constrained to $\mathbf{W}_1 = \begin{bmatrix} w_1 & w_1 & 0 & 0 \\ w_1 & w_1 & w_1 & w_1 \end{bmatrix}$, $\mathbf{W}_2 = \begin{bmatrix} w_2 & w_2 \end{bmatrix}$. The problem is $(x_1 + x_2) \cdot (x_1 + x_2 + x_3 + x_4)$ for $x = (1, 1.2, 1.8, 2)$. The solution is $w_1 = w_2 = 1$ in (a), with many unstable alternatives.

Multiplication issues

- weight issues:
 - gradient w.r.t. \mathbf{M} is expected to be zero
 - assumed bias towards $(-1, 0, 1)$ does not exist
- Multiplication issues:
 - singularities for $w < 0$ in NAC_\bullet .
 - NAC_\bullet can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

Division causes a singularity in the loss curvature

singularity

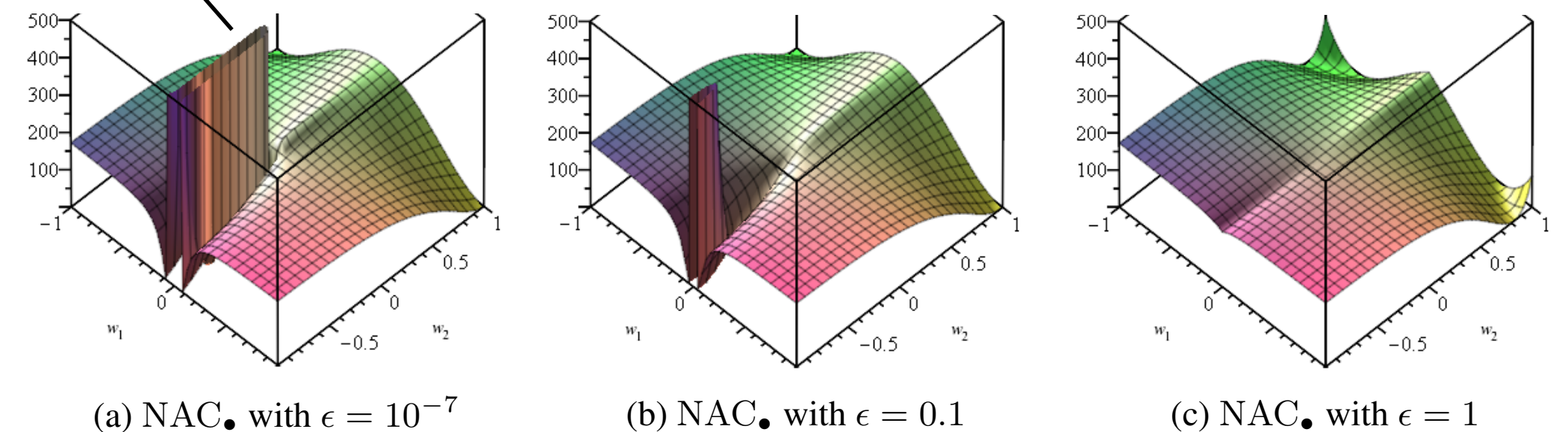


Figure 2: RMS loss curvature for a NAC_+ unit followed by a NAC_\bullet . The weight matrices are constrained to $\mathbf{W}_1 = \begin{bmatrix} w_1 & w_1 & 0 & 0 \\ w_1 & w_1 & w_1 & w_1 \end{bmatrix}$, $\mathbf{W}_2 = \begin{bmatrix} w_2 & w_2 \end{bmatrix}$. The problem is $(x_1 + x_2) \cdot (x_1 + x_2 + x_3 + x_4)$ for $x = (1, 1.2, 1.8, 2)$. The solution is $w_1 = w_2 = 1$ in (a), with many unstable alternatives.

Second order Taylor approximation

$$E[z_{h_\ell}] \approx \left(1 + \frac{1}{2} \text{Var}[W_{h_\ell, h_{\ell-1}}] \log(|E[z_{h_{\ell-1}}]| + \epsilon)^2 \right)^{H_{\ell-1}} \Rightarrow E[z_{h_\ell}] > 1.$$

The desirable is $E[z] = 0$

Neural Multiplication Unit

- weight issues:
 - gradient w.r.t. \mathbf{M} is expected to be zero
 - assumed bias towards $(-1, 0, 1)$ does not exist
- Multiplication issues:
 - singularities for $w < 0$ in NAC.
 - NAC₊ can't be initialized optimally
 - (no multiplication of negative numbers)
- Gating issues:
 - gating does not converge consistently

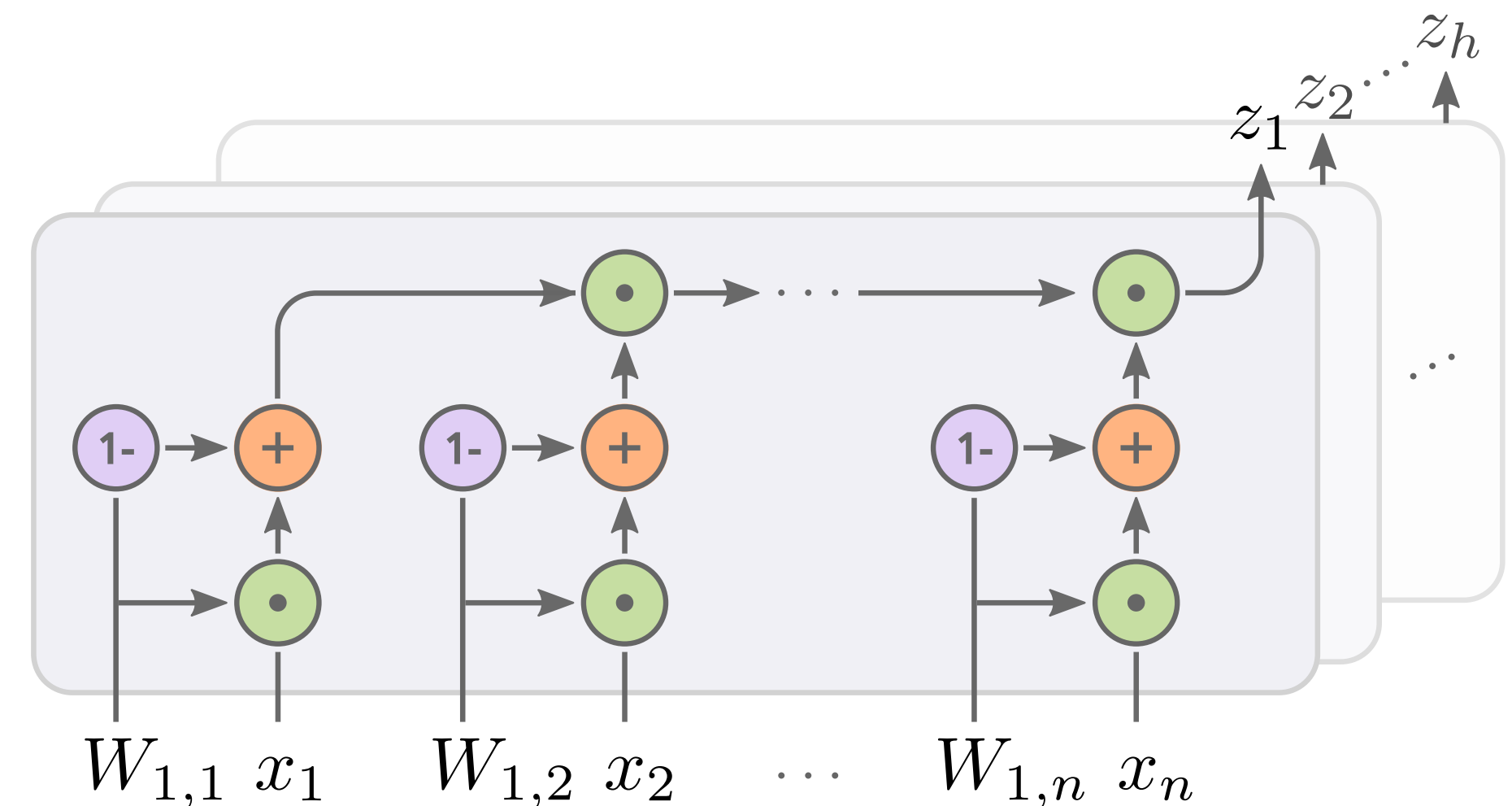
Our solution, **Neural Multiplication Unit**:

$$W_{h_{\ell-1}, h_{\ell}} = \min(\max(W_{h_{\ell-1}, h_{\ell}}, 0), 1),$$

$$\mathcal{R}_{\ell, \text{sparse}} = \frac{1}{H_{\ell} \cdot H_{\ell-1}} \sum_{h_{\ell}=1}^{H_{\ell}} \sum_{h_{\ell-1}=1}^{H_{\ell-1}} \min(W_{h_{\ell-1}, h_{\ell}}, 1 - W_{h_{\ell-1}, h_{\ell}})$$

$$\text{NMU} : z_{h_{\ell}} = \prod_{h_{\ell-1}=1}^{H_{\ell-1}} (W_{h_{\ell-1}, h_{\ell}} z_{h_{\ell-1}} + 1 - W_{h_{\ell-1}, h_{\ell}})$$

essentially a linear gate between 1 and z



Results

higher success-rate

Table 2: Comparison of: success-rate, first iteration reaching success, and sparsity error, all with 95% confidence interval on the “arithmetic datasets” task. Each value is a summary of 100 different seeds.

Op	Model	Success	Solved at iteration step		Sparsity error
		Rate	Median	Mean	Mean
×	NAC _•	31% $\begin{smallmatrix} +10\% \\ -8\% \end{smallmatrix}$	$2.8 \cdot 10^6$	$3.0 \cdot 10^6$ $\begin{smallmatrix} +2.9 \cdot 10^5 \\ -2.4 \cdot 10^5 \end{smallmatrix}$	$5.8 \cdot 10^{-4}$ $\begin{smallmatrix} +4.8 \cdot 10^{-4} \\ -2.6 \cdot 10^{-4} \end{smallmatrix}$
	NALU	0% $\begin{smallmatrix} +4\% \\ -0\% \end{smallmatrix}$	—	—	—
	NMU	98% $\begin{smallmatrix} +1\% \\ -5\% \end{smallmatrix}$	$1.4 \cdot 10^6$	$1.5 \cdot 10^6$ $\begin{smallmatrix} +5.0 \cdot 10^4 \\ -6.6 \cdot 10^4 \end{smallmatrix}$	$4.2 \cdot 10^{-7}$ $\begin{smallmatrix} +2.9 \cdot 10^{-8} \\ -2.9 \cdot 10^{-8} \end{smallmatrix}$
+	NAC ₊	100% $\begin{smallmatrix} +0\% \\ -4\% \end{smallmatrix}$	$2.5 \cdot 10^5$	$4.9 \cdot 10^5$ $\begin{smallmatrix} +5.2 \cdot 10^4 \\ -4.5 \cdot 10^4 \end{smallmatrix}$	$2.3 \cdot 10^{-1}$ $\begin{smallmatrix} +6.5 \cdot 10^{-3} \\ -6.5 \cdot 10^{-3} \end{smallmatrix}$
	Linear	100% $\begin{smallmatrix} +0\% \\ -4\% \end{smallmatrix}$	$6.1 \cdot 10^4$	$6.3 \cdot 10^4$ $\begin{smallmatrix} +2.5 \cdot 10^3 \\ -3.3 \cdot 10^3 \end{smallmatrix}$	$2.5 \cdot 10^{-1}$ $\begin{smallmatrix} +3.6 \cdot 10^{-4} \\ -3.6 \cdot 10^{-4} \end{smallmatrix}$
	NALU	14% $\begin{smallmatrix} +8\% \\ -5\% \end{smallmatrix}$	$1.5 \cdot 10^6$	$1.6 \cdot 10^6$ $\begin{smallmatrix} +3.8 \cdot 10^5 \\ -3.3 \cdot 10^5 \end{smallmatrix}$	$1.7 \cdot 10^{-1}$ $\begin{smallmatrix} +2.7 \cdot 10^{-2} \\ -2.5 \cdot 10^{-2} \end{smallmatrix}$
	NAU	100% $\begin{smallmatrix} +0\% \\ -4\% \end{smallmatrix}$	$1.8 \cdot 10^4$	$3.9 \cdot 10^5$ $\begin{smallmatrix} +4.5 \cdot 10^4 \\ -3.7 \cdot 10^4 \end{smallmatrix}$	$3.2 \cdot 10^{-5}$ $\begin{smallmatrix} +1.3 \cdot 10^{-5} \\ -1.3 \cdot 10^{-5} \end{smallmatrix}$
−	NAC ₊	100% $\begin{smallmatrix} +0\% \\ -4\% \end{smallmatrix}$	$9.0 \cdot 10^3$	$3.7 \cdot 10^5$ $\begin{smallmatrix} +3.8 \cdot 10^4 \\ -3.8 \cdot 10^4 \end{smallmatrix}$	$2.3 \cdot 10^{-1}$ $\begin{smallmatrix} +5.4 \cdot 10^{-3} \\ -5.4 \cdot 10^{-3} \end{smallmatrix}$
	Linear	7% $\begin{smallmatrix} +7\% \\ -4\% \end{smallmatrix}$	$3.3 \cdot 10^6$	$1.4 \cdot 10^6$ $\begin{smallmatrix} +7.0 \cdot 10^5 \\ -6.1 \cdot 10^5 \end{smallmatrix}$	$1.8 \cdot 10^{-1}$ $\begin{smallmatrix} +7.2 \cdot 10^{-2} \\ -5.8 \cdot 10^{-2} \end{smallmatrix}$
	NALU	14% $\begin{smallmatrix} +8\% \\ -5\% \end{smallmatrix}$	$1.9 \cdot 10^6$	$1.9 \cdot 10^6$ $\begin{smallmatrix} +4.4 \cdot 10^5 \\ -4.5 \cdot 10^5 \end{smallmatrix}$	$2.1 \cdot 10^{-1}$ $\begin{smallmatrix} +2.2 \cdot 10^{-2} \\ -2.2 \cdot 10^{-2} \end{smallmatrix}$
	NAU	100% $\begin{smallmatrix} +0\% \\ -4\% \end{smallmatrix}$	$5.0 \cdot 10^3$	$1.6 \cdot 10^5$ $\begin{smallmatrix} +1.7 \cdot 10^4 \\ -1.6 \cdot 10^4 \end{smallmatrix}$	$6.6 \cdot 10^{-2}$ $\begin{smallmatrix} +2.5 \cdot 10^{-2} \\ -1.9 \cdot 10^{-2} \end{smallmatrix}$

faster convergence *sparser solution*

Results

higher success-rate

Table 2: Comparison of: success-rate, first iteration reaching success, and sparsity error, all with 95% confidence interval on the “arithmetic datasets” task. Each value is a summary of 100 different seeds.

Op	Model	Success	Solved at iteration step		Sparsity error
		Rate	Median	Mean	Mean
×	NAC _•	31% $+10\%$ -8%	$2.8 \cdot 10^6$	$3.0 \cdot 10^6$ $+2.9 \cdot 10^5$ $-2.4 \cdot 10^5$	$5.8 \cdot 10^{-4}$ $+4.8 \cdot 10^{-4}$ $-2.6 \cdot 10^{-4}$
	NALU	0% $+4\%$ -0%	—	—	—
	NMU	98% $+1\%$ -5%	$1.4 \cdot 10^6$	$1.5 \cdot 10^6$ $+5.0 \cdot 10^4$ $-6.6 \cdot 10^4$	$4.2 \cdot 10^{-7}$ $+2.9 \cdot 10^{-8}$ $-2.9 \cdot 10^{-8}$
+	NAC ₊	100% $+0\%$ -4%	$2.5 \cdot 10^5$	$4.9 \cdot 10^5$ $+5.2 \cdot 10^4$ $-4.5 \cdot 10^4$	$2.3 \cdot 10^{-1}$ $+6.5 \cdot 10^{-3}$ $-6.5 \cdot 10^{-3}$
	Linear	100% $+0\%$ -4%	$6.1 \cdot 10^4$	$6.3 \cdot 10^4$ $+2.5 \cdot 10^3$ $-3.3 \cdot 10^3$	$2.5 \cdot 10^{-1}$ $+3.6 \cdot 10^{-4}$ $-3.6 \cdot 10^{-4}$
	NALU	14% $+8\%$ -5%	$1.5 \cdot 10^6$	$1.6 \cdot 10^6$ $+3.8 \cdot 10^5$ $-3.3 \cdot 10^5$	$1.7 \cdot 10^{-1}$ $+2.7 \cdot 10^{-2}$ $-2.5 \cdot 10^{-2}$
	NAU	100% $+0\%$ -4%	$1.8 \cdot 10^4$	$3.9 \cdot 10^5$ $+4.5 \cdot 10^4$ $-3.7 \cdot 10^4$	$3.2 \cdot 10^{-5}$ $+1.3 \cdot 10^{-5}$ $-1.3 \cdot 10^{-5}$
−	NAC ₊	100% $+0\%$ -4%	$9.0 \cdot 10^3$	$3.7 \cdot 10^5$ $+3.8 \cdot 10^4$ $-3.8 \cdot 10^4$	$2.3 \cdot 10^{-1}$ $+5.4 \cdot 10^{-3}$ $-5.4 \cdot 10^{-3}$
	Linear	7% $+7\%$ -4%	$3.3 \cdot 10^6$	$1.4 \cdot 10^6$ $+7.0 \cdot 10^5$ $-6.1 \cdot 10^5$	$1.8 \cdot 10^{-1}$ $+7.2 \cdot 10^{-2}$ $-5.8 \cdot 10^{-2}$
	NALU	14% $+8\%$ -5%	$1.9 \cdot 10^6$	$1.9 \cdot 10^6$ $+4.4 \cdot 10^5$ $-4.5 \cdot 10^5$	$2.1 \cdot 10^{-1}$ $+2.2 \cdot 10^{-2}$ $-2.2 \cdot 10^{-2}$
	NAU	100% $+0\%$ -4%	$5.0 \cdot 10^3$	$1.6 \cdot 10^5$ $+1.7 \cdot 10^4$ $-1.6 \cdot 10^4$	$6.6 \cdot 10^{-2}$ $+2.5 \cdot 10^{-2}$ $-1.9 \cdot 10^{-2}$

sparser solution

faster convergence

supports negative values

better at larger hidden-sizes

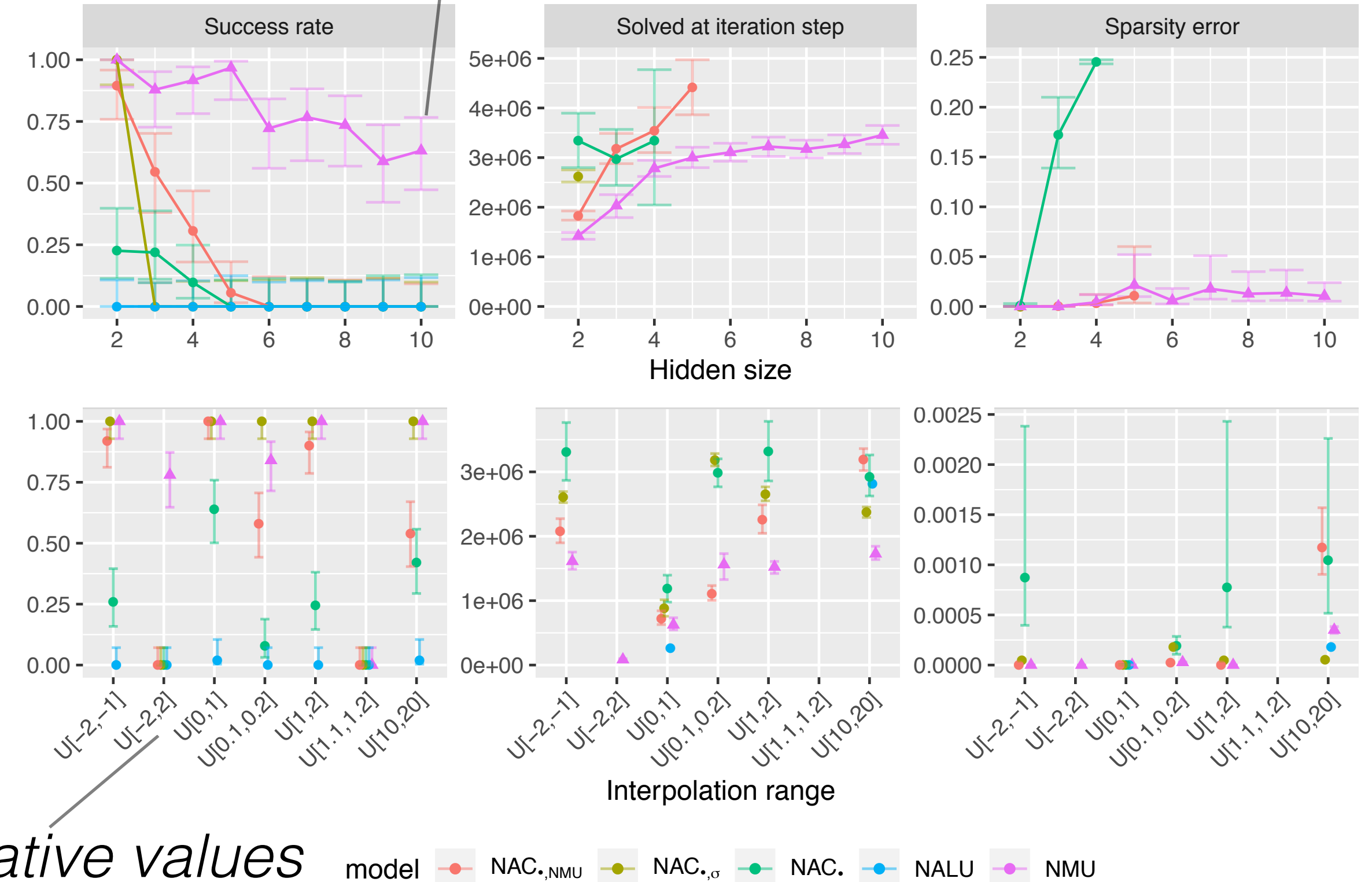


Figure 3: Multiplication task results when varying the hidden input-size and when varying the input-range. Extrapolation ranges are defined in Appendix C.4.

Not covered in this presentation

- Theory:
 - Initialization of NMU
 - Gradients of NMU
 - Regularization scaling
- Discussions:
 - Gating issues (Appendix C.5)
 - Effect of shared weights in NALU
 - Measuring performance
- Experiments:
 - Detailed ablation study
 - Effect of dataset parameters
 - Results for MNIST experiment
 - Hyperparameter optimization
 - Complete comparison of all models on all arithmetic problems

Neural Arithmetic Units

By **Andreas Madsen** and **Alexander Rosenberg Johansen**
@andreas_madsen @AlexRoseJo

ICLR 2020, spotlight awarded paper

paper: <https://openreview.net/forum?id=H1gNOeHKPS>

code: <https://github.com/AndreasMadsen/stable-nalu>