

Stochastic Rounding: Algorithms and Hardware Accelerator

Mantas Mikaitis, University of Manchester, UK

Contact: mantas.mikaitis@manchester.ac.uk

Virtual talk. July 2021.

IJCNN 2021 : International Joint Conference on Neural Networks

Background and motivation

- **SpiNNaker** is a 1M-core **digital neuromorphic computer**
- Optimized to simulate **Spiking Neural Networks (SNNs)**
- Situated in Manchester and has been in use for more than 10 years
- **SpiNNaker2** will be a 10M-core next-gen system @ TU Dresden
- Being built in collaboration between Manchester and Dresden



SpiNNaker



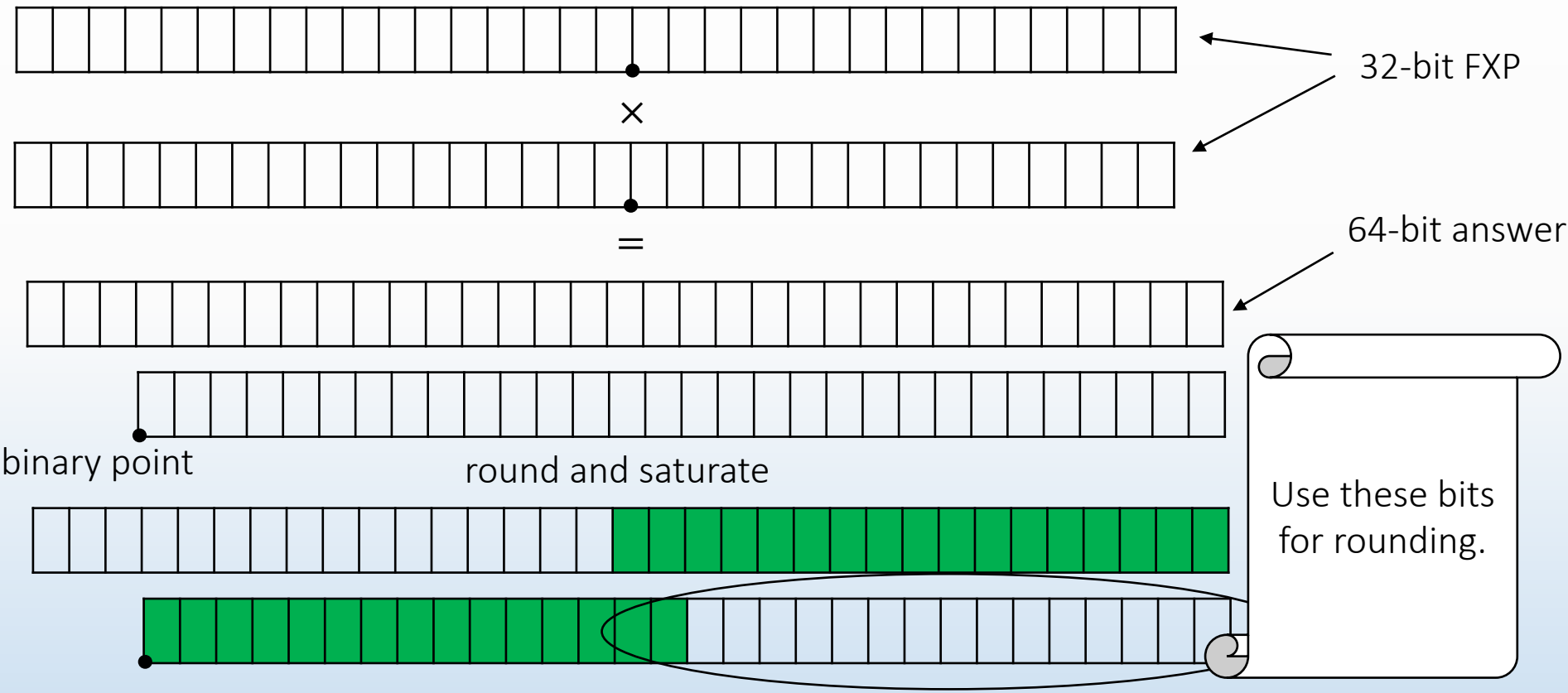
SpiNNaker2 prototype board

Background and motivation

- SpiNNaker chip uses ARM integer processor
- Real number computation implemented using fixed-point (FXP) representation
- SpiNNaker2 will have binary32 (*float*) floating-point (FLP)
- However, could still favour lower precision in parts of code
- Such as 16-bit synaptic weights
- Custom-precision FXP arith is not supported in HW
- Arith operators implemented in software
- The goal of this work: design custom-prec arith accelerator for SpiNNaker2
- FXP custom prec + binary32 -> bfloat16 rounding

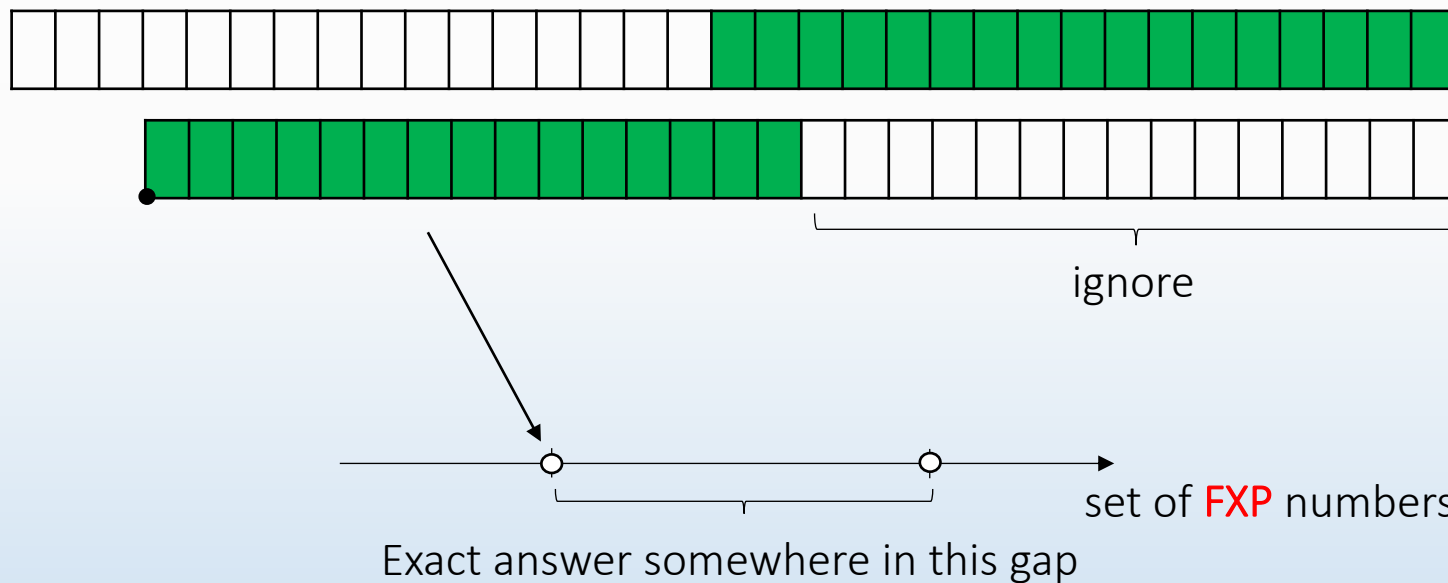
Programmable precision round-and-saturate accelerator for SpiNNaker2

Example: fixed-point 15-bit fraction multiplier



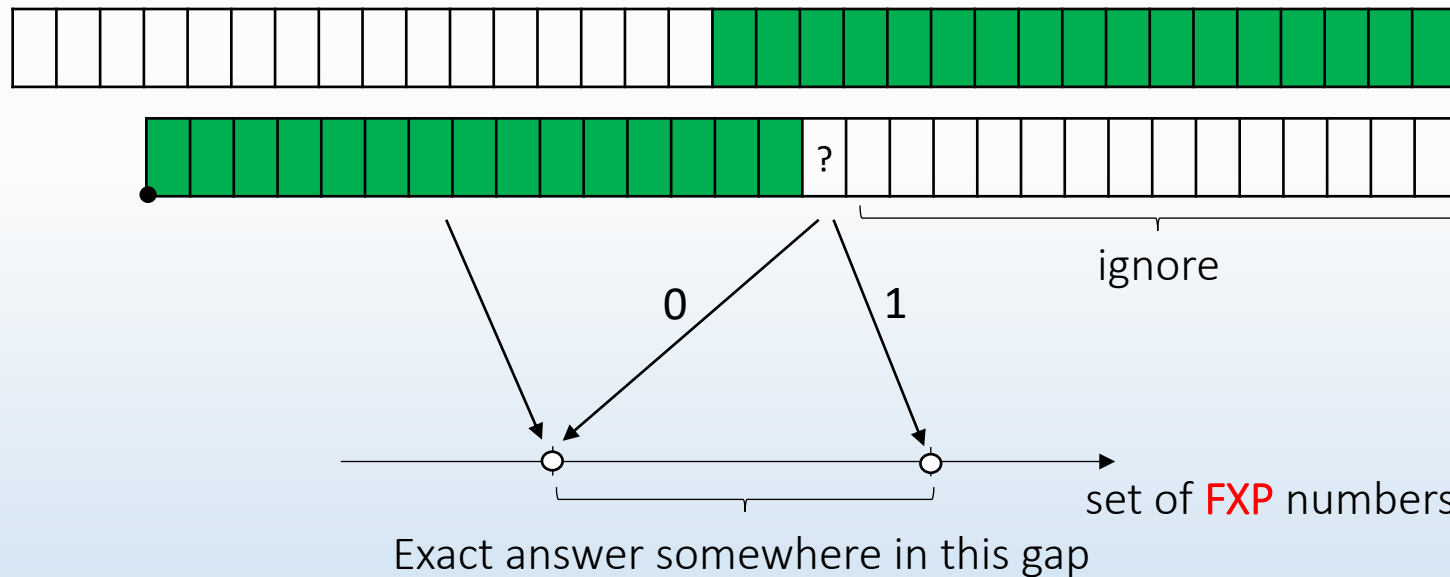
Round-down (RD)

Given the output from the fixed-point multiplier:



Round-to-nearest (RN)

Given the output from the fixed-point multiplier:



Stochastic rounding (SR)

Given a real number x , a random number $P \in [0,1)$ from a uniform distribution and a fixed-point destination format $\langle s, i, p \rangle$ with $\epsilon = 2^{-p}$,

$$\text{SR}(x, \langle s, i, p \rangle) = \begin{cases} \lfloor x \rfloor & \text{if } P \geq \frac{x - \lfloor x \rfloor}{\epsilon}, \\ \lfloor x \rfloor + \epsilon & \text{if } P < \frac{x - \lfloor x \rfloor}{\epsilon}. \end{cases}$$

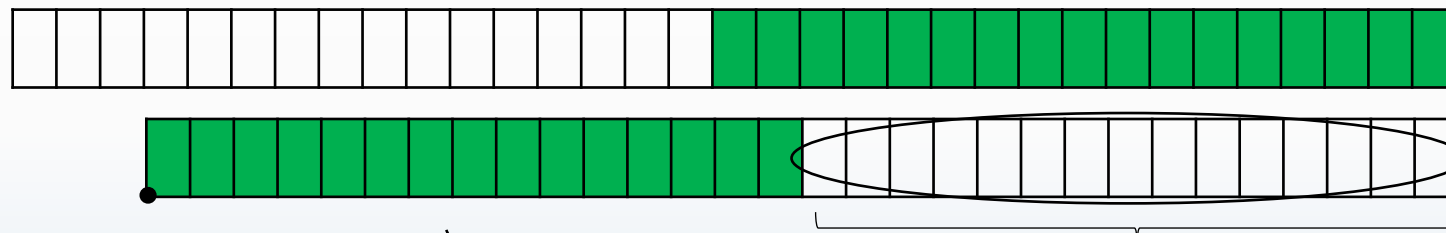
Over many roundings:

$$\mathbb{E}(\text{SR}(x, \langle s, i, p \rangle)) = x.$$

Specific precision SR present in some HW, for example Intel and Graphcore

SR at bit level

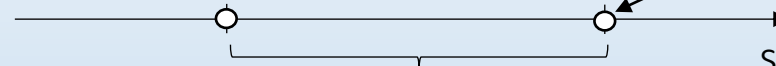
Given the output from the fixed-point multiplier:



Use these bits as probability of rounding up, $[0,1)$.

Round-off bits

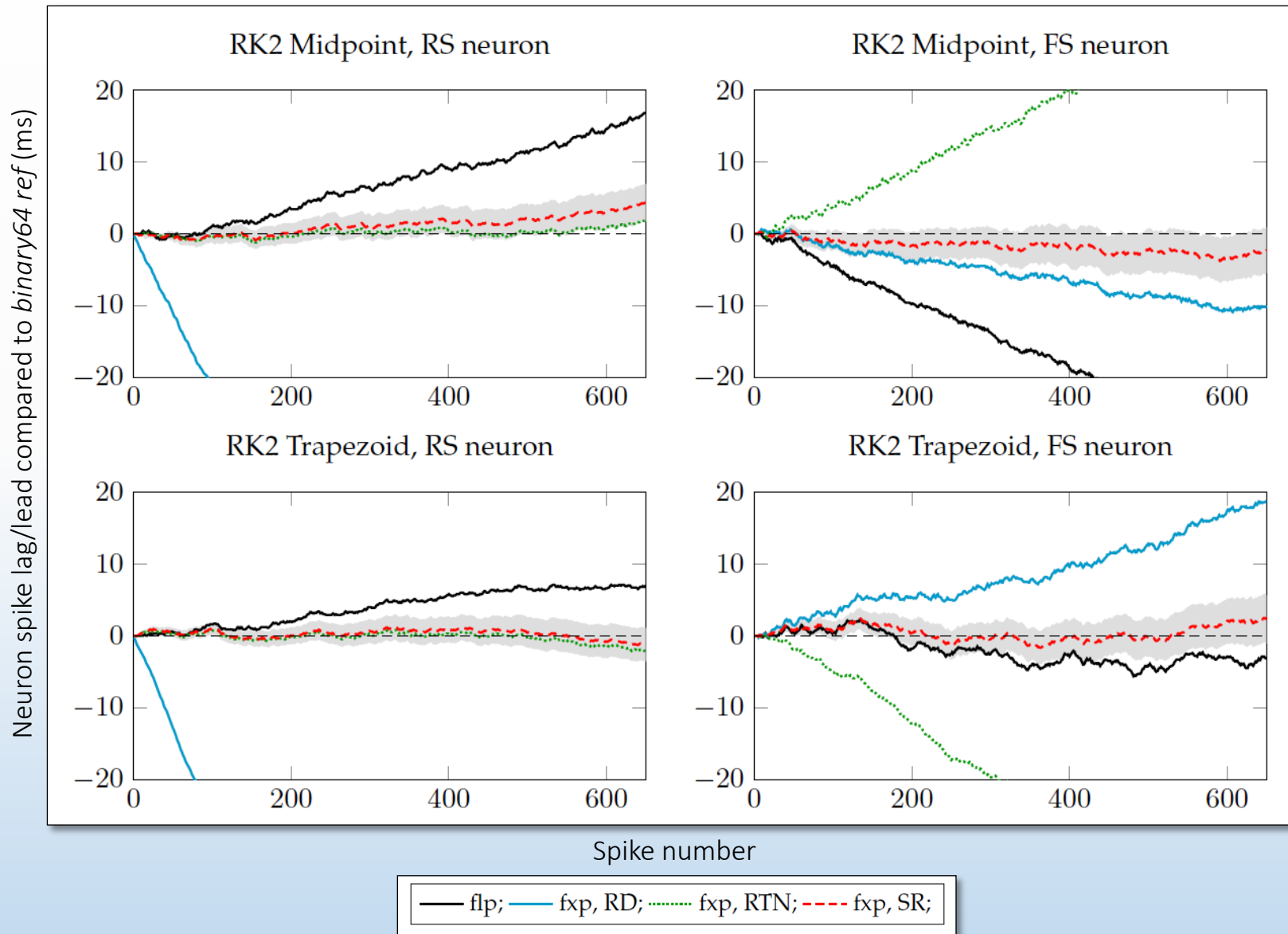
If $\text{die} < \text{round-off value}$
round up, else round down.



set of **FXP** numbers

Exact answer somewhere in this gap

SR in ODE solvers on SpiNNaker



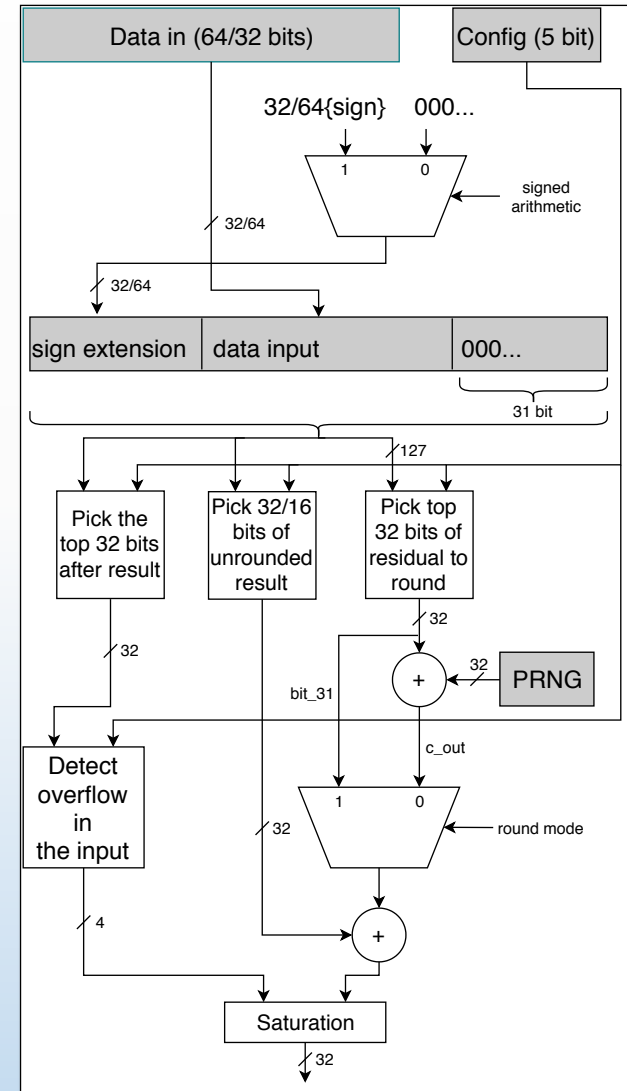
SR accelerator on SpiNNaker2

- Round and saturate 64-, 32-, or 16-bit to 32- or 16-bit **FXP** numbers
- **SR** and **RN** (ties up) modes.
- **Rounding bit position programmable**
- **Signed and unsigned formats**
- Rounding and saturation of **binary32** values to **bfloat16** values
- Uses **SpiNNaker2** hardware pseudo-random 32-bit streams (**PRNG**)
- Up to four threads with different **PRNG** seeds are supported
- Accelerator is general purpose: not specialized to **SpiNNaker2** hardware or **SNN** applications
- 3 or 4 clock cycle latency

Use: perform arith ops on ARM, round/saturate to custom prec with the accelerator

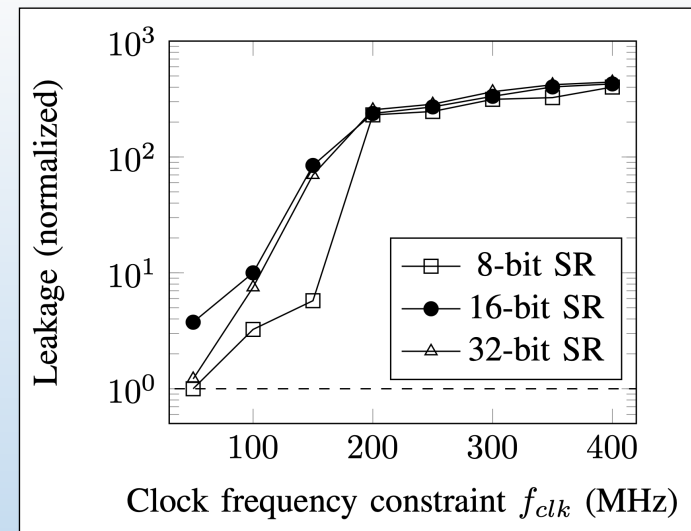
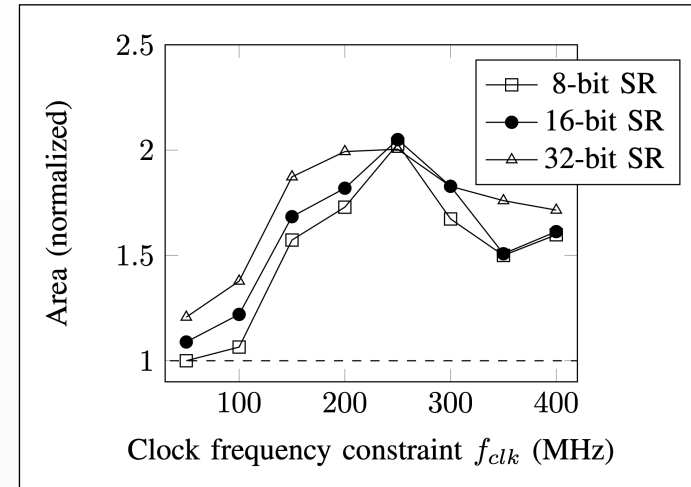
SR accelerator on SpiNNaker2

- Data and config written to registers
- Config contains **round bit**: 0 to 31
- Depending on address, determine
 1. Source and dest **word widths**
 2. **Signed/unsigned** arith
 3. **Rounding mode**
- Data is sign extended and padded with 0's on the right
- Rounding is performed by utilizing **PRNG (SR)** or MSB of chopped bits (**RN**)
- Detect **overflow** using top bits, **saturate** to target width if needed



Evaluation of accelerators

- **Three accelerators** were evaluated, with 8-, 16-, and 32-bit adders
- Influences **SR** precision and **PRNG** bits needed for each rounding
- Synthesis in **SpiNNaker2** 22nm library
- **Worst case speed conditions** used
- Target frequency 50 to 400Mhz
- 8-bit **SR** can provide lower leakage than 32-bit in some cases
- Some circuit area savings with 8-bit
- Results only include **SR** accelerators
- Further savings in **PRNG** would be possible if 8-bit is good enough



Summary

- This paper presents a design of **SR** rounding unit—first with programmable precision
- Arithmetic can be done in **ARM** and rounded with this unit
- Accelerator replaces majority of steps in SW arithmetic operations of **SpiNNaker**
- It will speed up both stochastic and standard arithmetic libraries on **SpiNNaker2**
- **Binary32** to **bfloat16** rounding potentially useful for memory savings: operate on **bfloat16** as **binary32** using the **FPU**

SpiNNaker2 is scheduled for tape out later this year

References

- S. B. Furber, F. Galluppi, S. Temple, and L. A. Plana. [The SpiNNaker Project](#). Proc. IEEE, vol. 102. May 2014
- S. Höppner *et al.* [The SpiNNaker 2 Processing Element Architecture for Hybrid Digital Neuromorphic Computing](#). arXiv:2103.08392 [cs.AR]. Mar. 2021
- S. Höppner and C. Mayr. [SpiNNaker2 - Towards Extremely Efficient Digital Neuromorphics and Multi-scale Brain Emulation](#). Proc. Neuro Inspired Computational Elements Conference. 2018
- M. Hopkins, M. Mikaitis, D. R. Lester, and S. B. Furber. [Stochastic rounding and reduced-precision fixed-point arithmetic for solving neural ordinary differential equations](#). Phil. Trans. R. Soc. A, vol. 378. Jan. 2020