

Neural Networks are "Hot Hot Hot!!!"

- Neural networks are layered structures consisting of neurons useful for:
 - Classification/regression
 - Soft computing
- Broad use in high tech, e.g., Google [1]
- Enable alternative computing methodologies
 - Automatic parallelization, e.g., ASC [2]
 - Approximate computing, e.g., NPU [3]

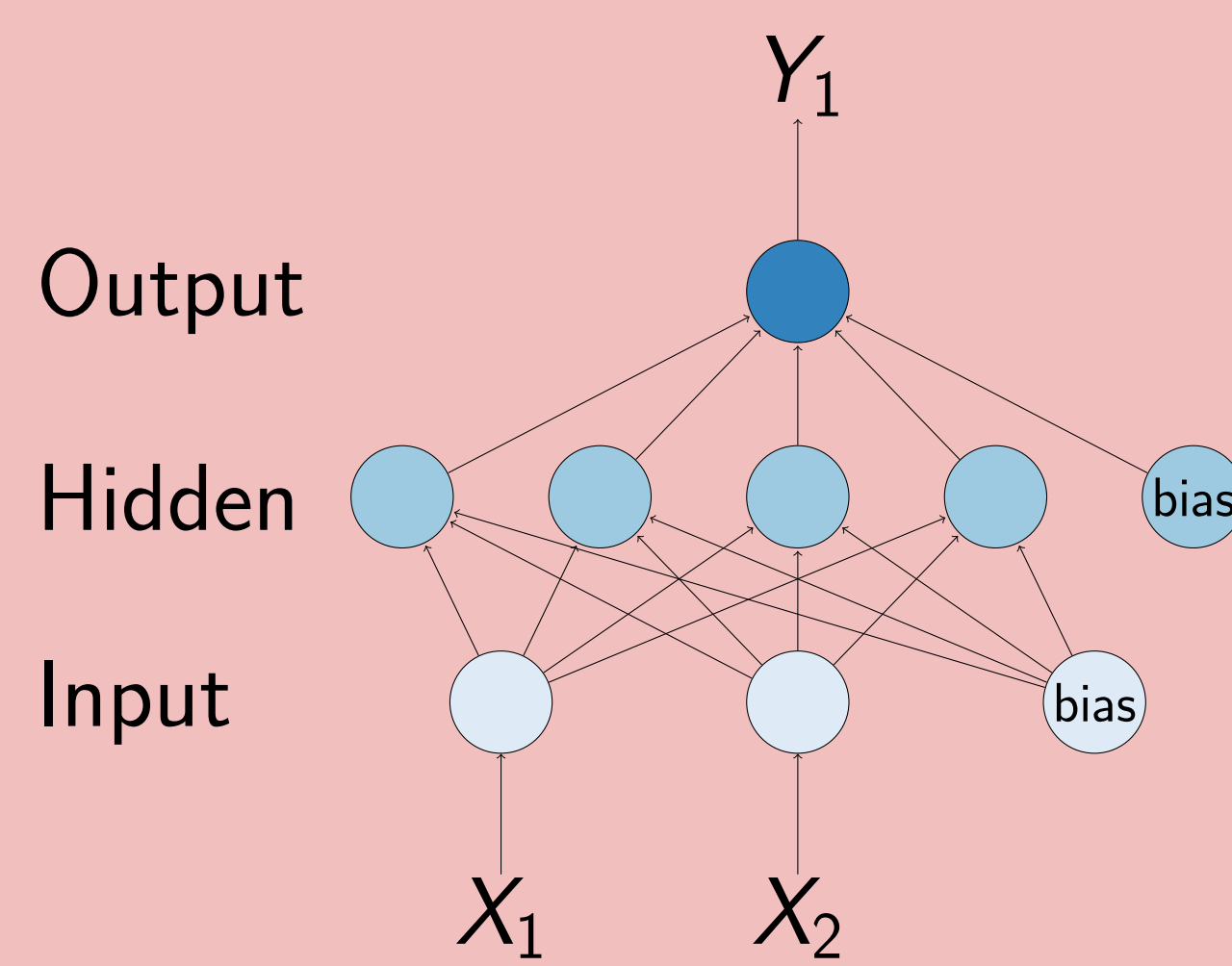


Figure 1: An example neural network

- [1] Google TensorFlow, <https://github.com/tensorflow/tensorflow>
 [2] A. Waterland, E. Angelino et al., "Asc: Automatically scalable computation," in *Proc. ASPLOS*, 2014.
 [3] H. Esmailzadeh, A. Sampson et al., "Neural acceleration for general-purpose approximate programs," in *Proc. MICRO*, 2012.

Our Vision of General-Purpose Neural Network Computing

- Treat neural networks as *functional primitives* backed by accelerators [1]
- Follow a transaction model for neural network computation
- Exploit anticipated sharing of neural networks across applications
 - We break this into two distinct contributions:
 - X-FILES – Software/Hardware extensions for transaction management
 - DANA – One possible backend accelerator that interfaces with X-FILES

- [1] S. Eldridge, A. Waterland et al., "Towards general-purpose neural network computing," in *Proc. PACT*, 2015.

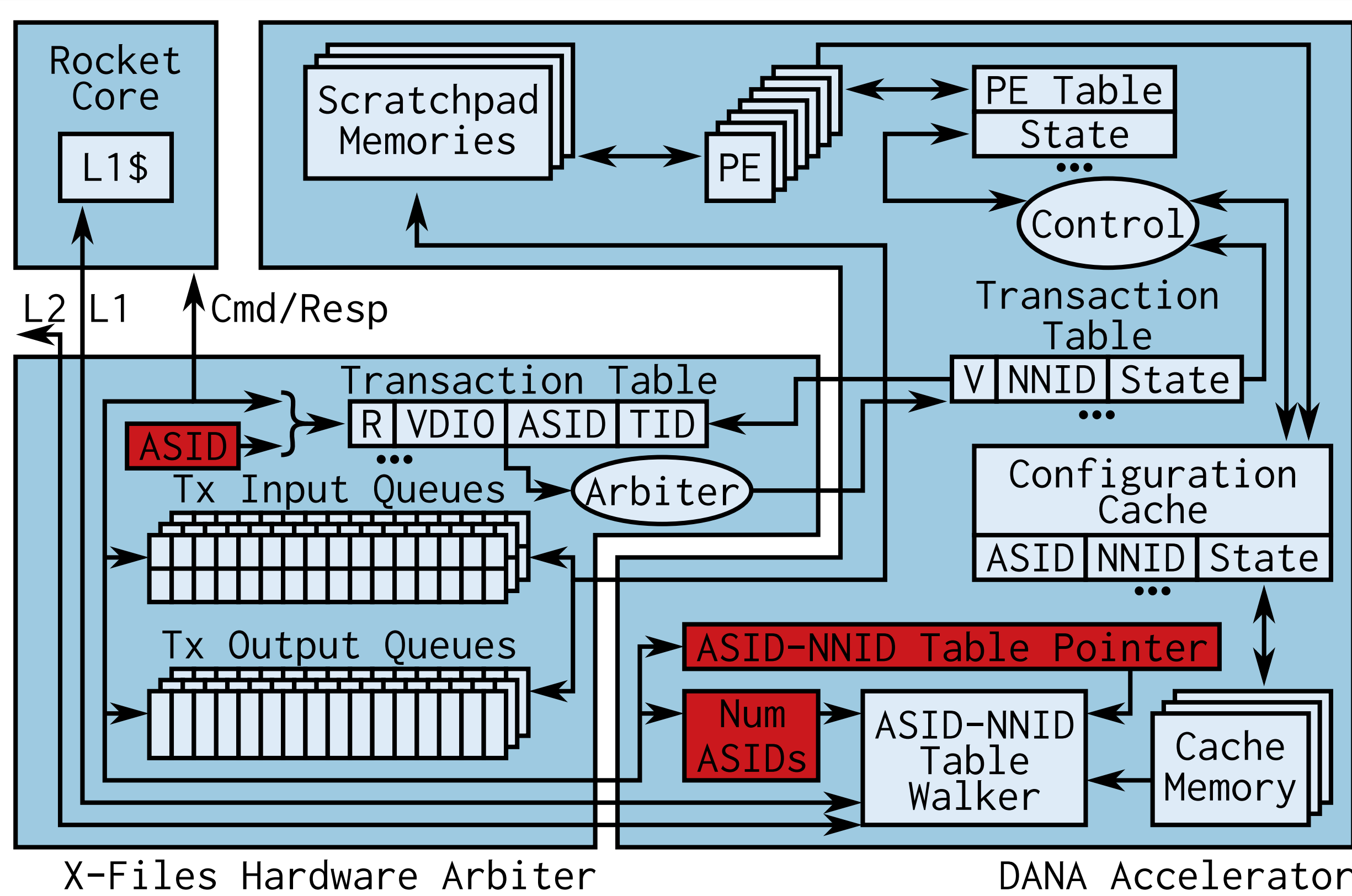


Figure 2: The X-FILES/DANA hardware architecture. Neural network transactions are managed by the X-FILES hardware arbiter and executed on the backend accelerator, DANA.

Table 1: X-FILES user, supervisor, and Proxy Kernel syscall API

User/Supervisor Function	Code
supervisor	old asid = set asid(new asid) old antp = set antp(*asid nnid table entry, size) csr = xf read csr(csr_index)
user	tid = new write request(nnid, learning_type, num output) error_code = write data(tid, *inputs, num input) error_code = read data spinlock(tid, *output, num output) id = xfiles dana_id(flag print) error_code = kill transaction(tid)
user (PK only)	old asid = pk syscall set asid(new asid) old antp = pk syscall set antp(new antp) asid nnid table create(**table, num asids, num configs) asid nnid table destroy(**table) attach nn configuration(**table, asid, *nn_config) attach garbage(**table, asid)

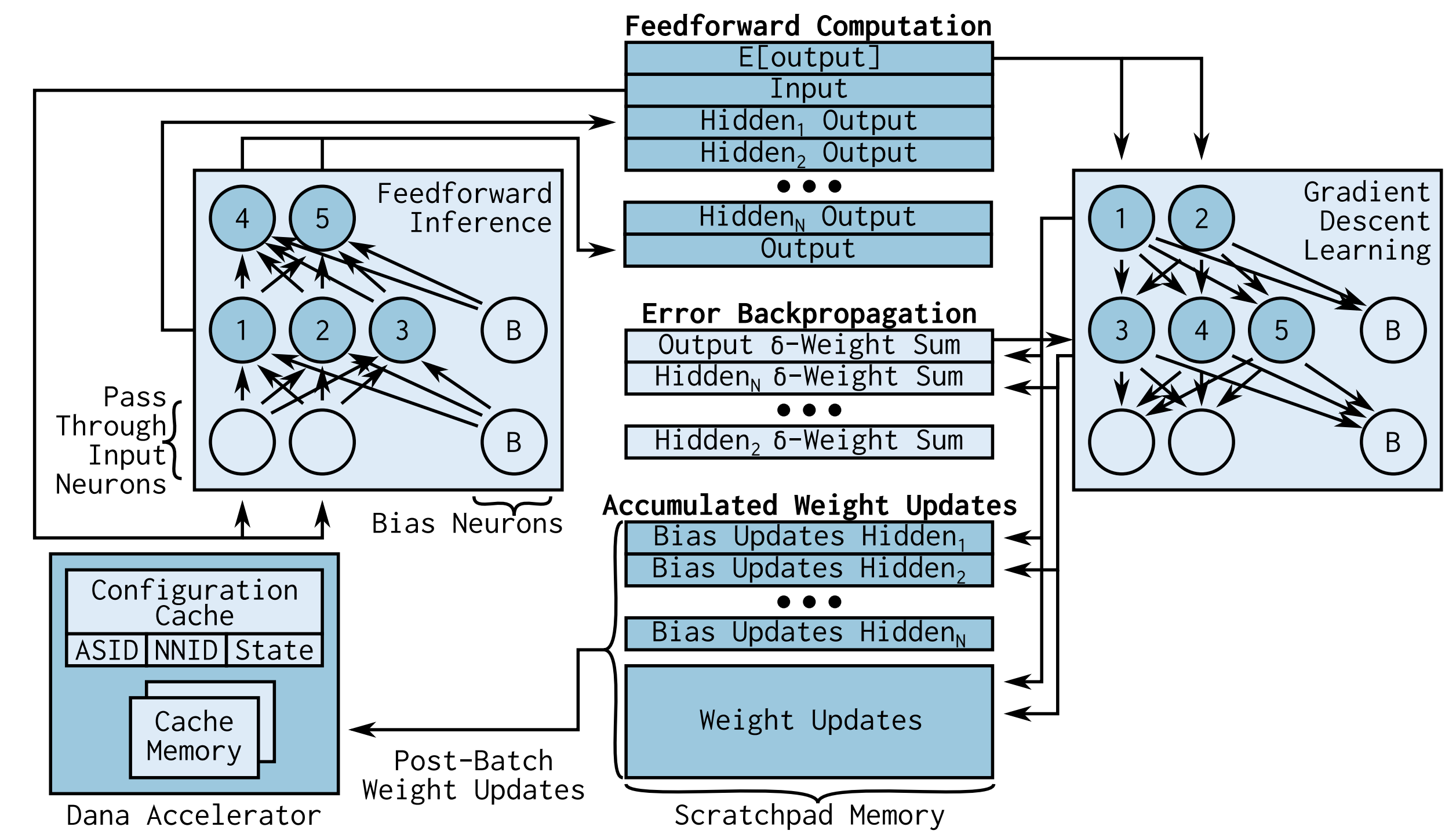


Figure 3: Processing Element (PE) ordering and the interactions of PEs with one of DANA's per-transaction scratchpad memories for feedforward and learning transactions

Safe Multi-Transaction Management using ASIDs

- An ASID defines an address space of neural networks, identified with NNIDs
- An OS-managed ASID-NNID Table enables NNID dereferencing (Figure 4)
- All transparently handled by the OS

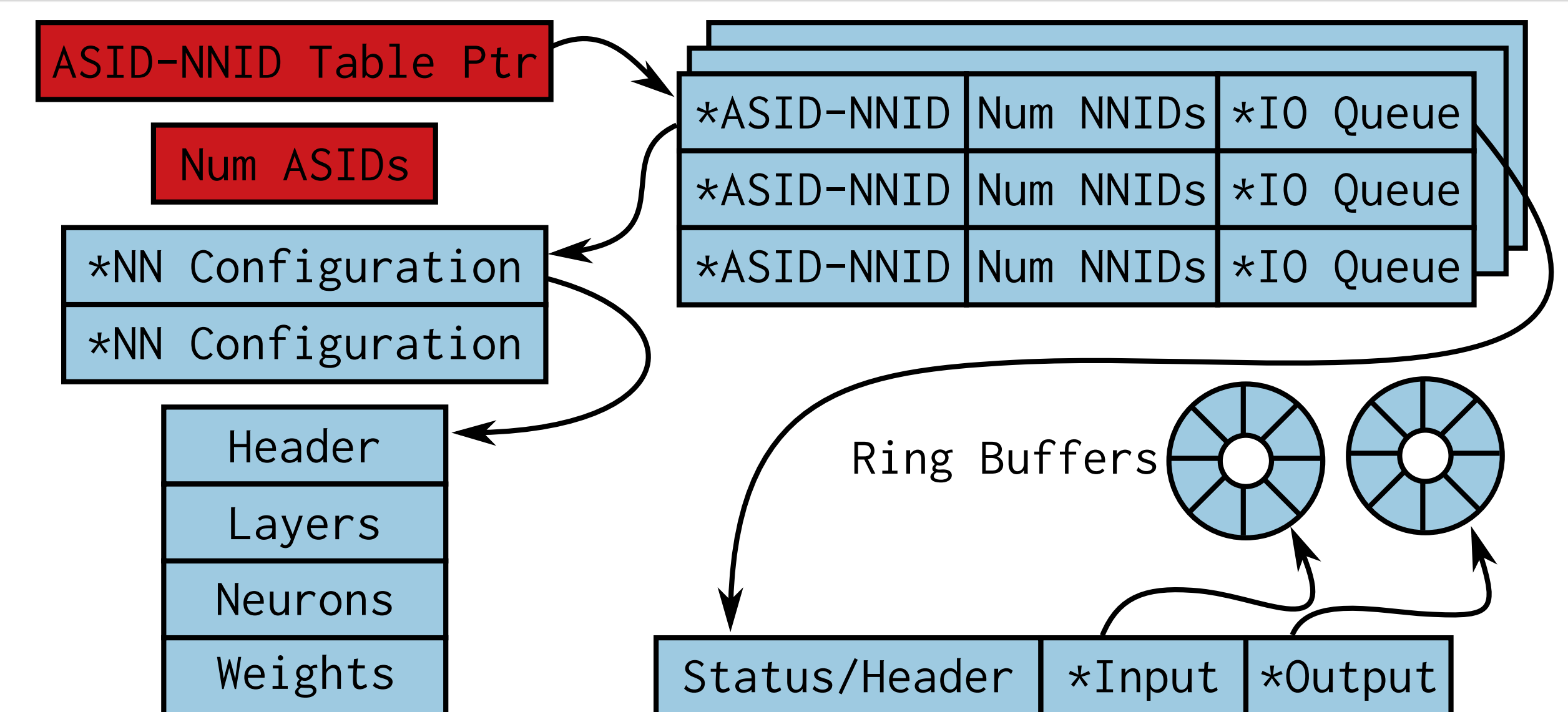


Figure 4: An ASID-NNID Table for dereferencing neural network configurations from an ASID and an NNID.

Usage with a RISC-V Microprocessor

- Grab a Rocket Chip RISC-V Microprocessor [1]
- Build a RISC-V toolchain
- Grab a copy of our X-FILES/DANA accelerator [2]
- Build Rocket + X-FILES/DANA for FPGA
- User processes can throw *transactions* at X-FILES

- [1] Rocket Chip git repository, UC Berkeley, Online: github.com/ucb-bar/rocket-chip
 [2] X-FILES/DANA git repository, Boston University, Online: github.com/bu-icsg/xfiles-dana

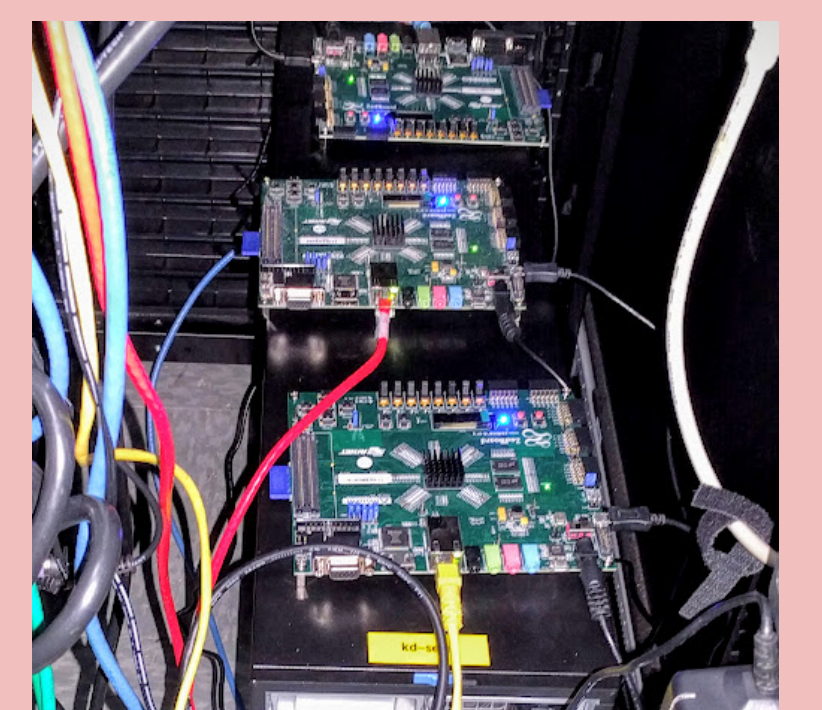


Figure 5: FPGA Setup

Initial Evaluation

- Comparison against a RISC-V microprocessor running FANN [1]
 - Figure 6 shows a 30x improvement over software for networks in Table 2

- [1] S. Nissen, "Implementation of a fast artificial neural network library (fann)," DIKU, Tech. Rep., 2003, <http://fann.sf.net>.

Open Source Availability and Acknowledgments

- On GitHub:** github.com/bu-icsg/xfiles-dana
- This work was supported by the following:
 - A NASA Space Technology Research Fellowship, NSF CAREER Awards

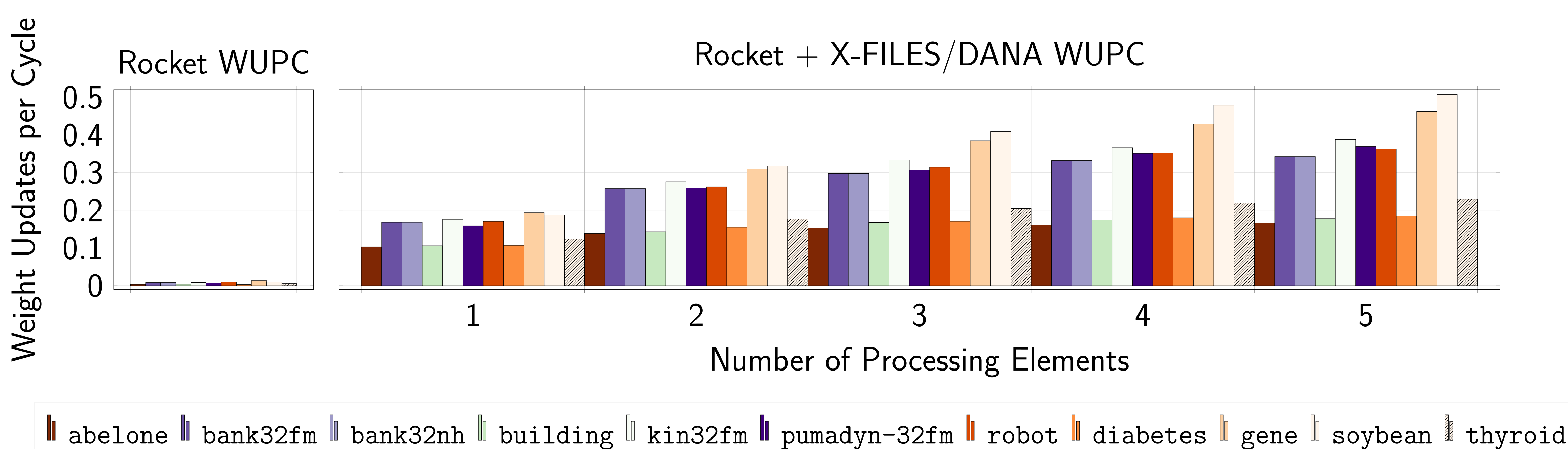


Figure 6: The weight updates per cycle (WUPC) when executing batch learning requests for different NN configurations

Table 2: Evaluated FANN-provided datasets

Type	Name	Topology
Regression	abelone	10 × 8 × 1
	bank32fm	32 × 16 × 1
	bank32nh	32 × 16 × 1
	building	14 × 8 × 3
	kin32fm	32 × 20 × 1
	pumadyn_32fm	32 × 16 × 8 × 4 × 1
Classification	robot	48 × 16 × 3
	diabetes	8 × 10 × 2
	gene	120 × 20 × 3
	soybean	82 × 32 × 19
	thyroid	21 × 10 × 3