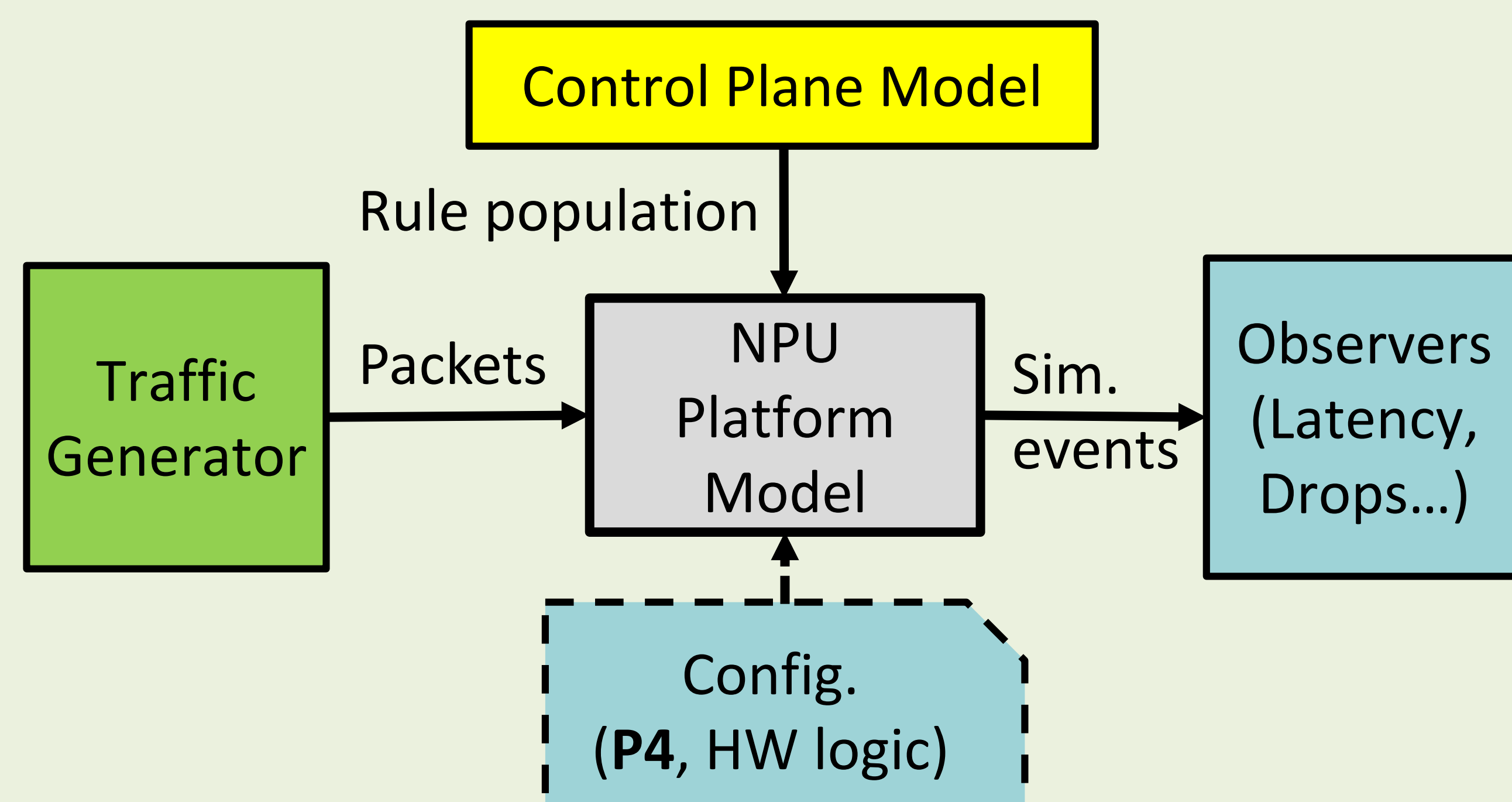
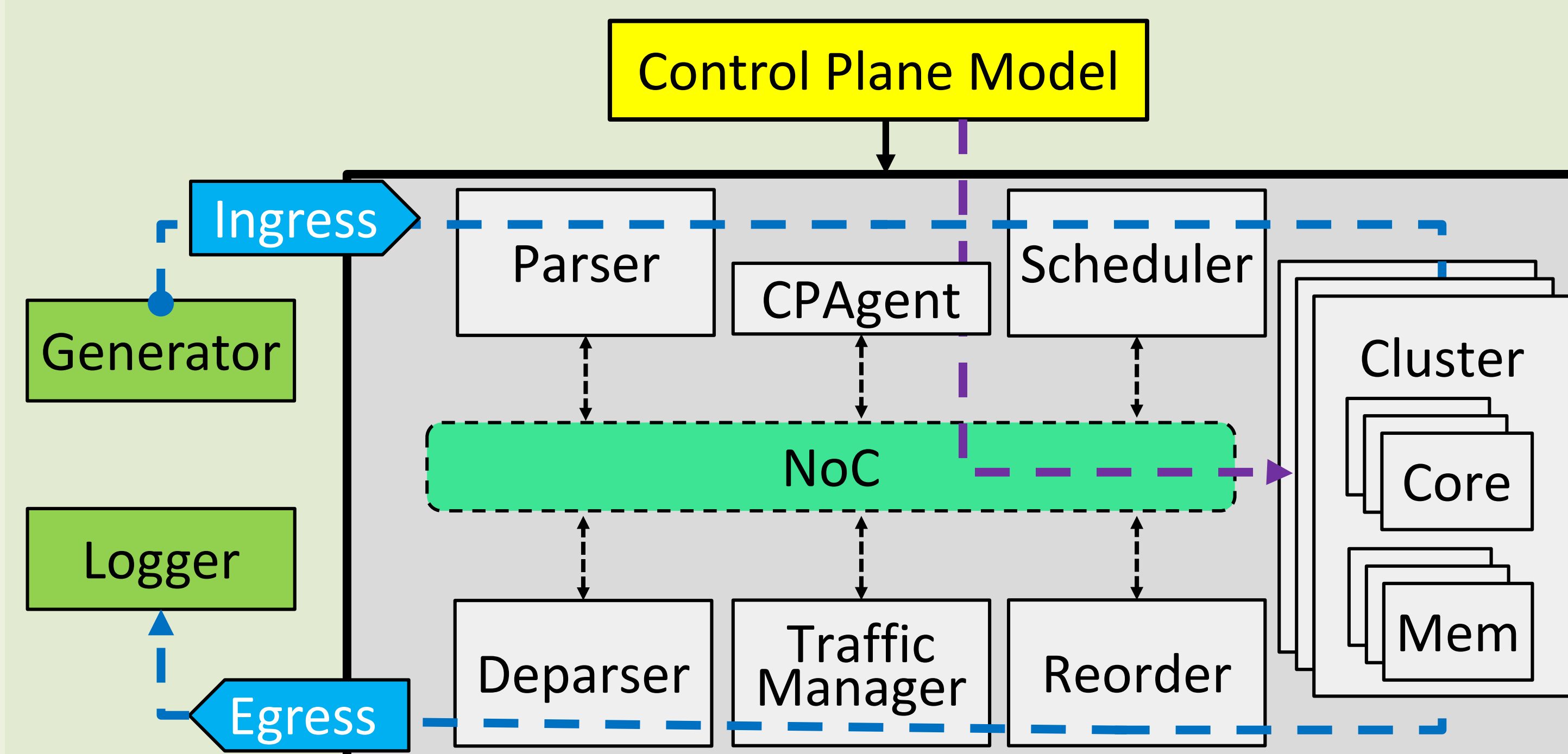


Modeling Goals



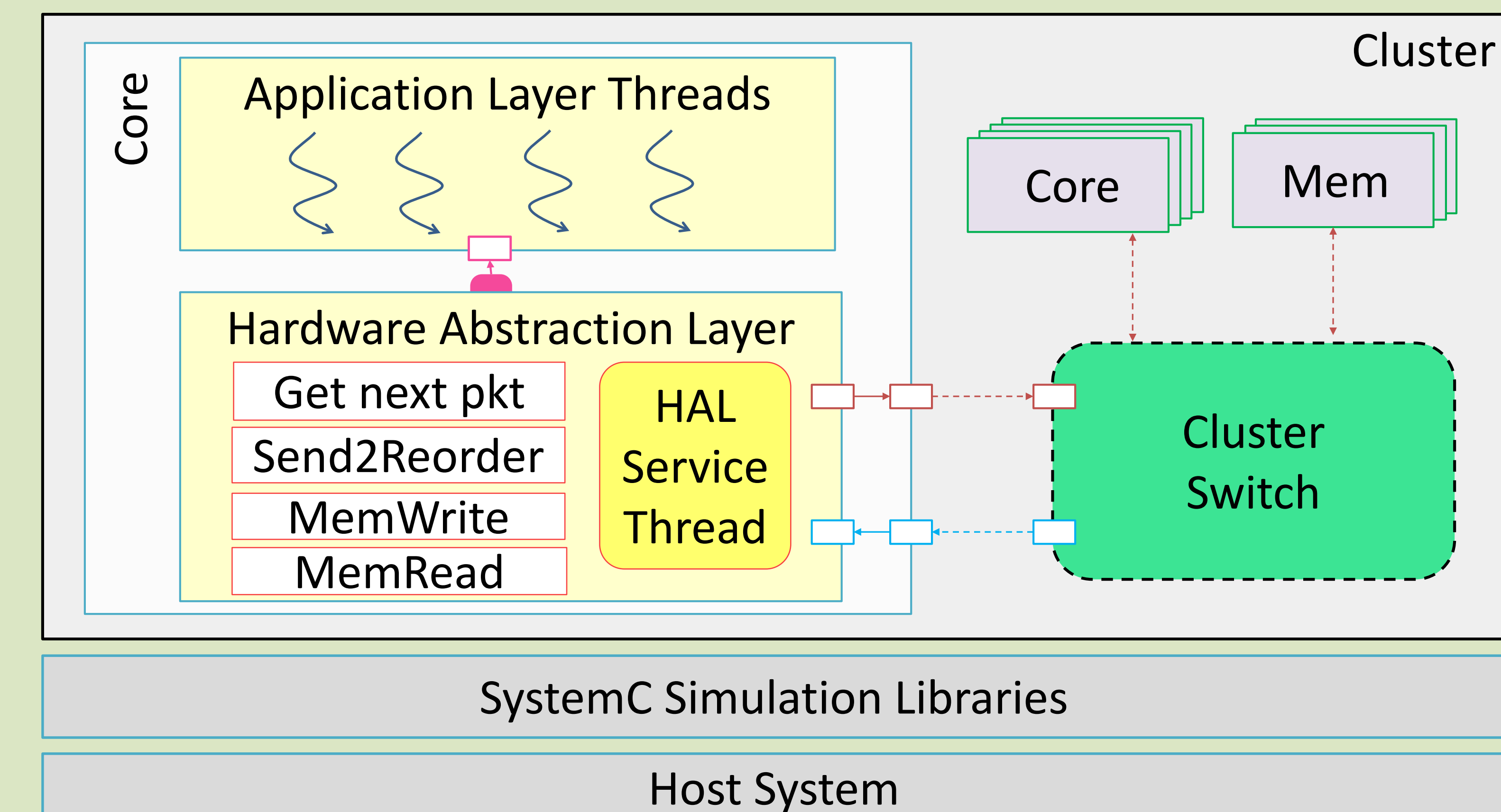
- System-level simulation of NPU architecture
- High simulation speed for functional validation
- (Reasonably) accurate performance estimation
- Easy modification of architectural parameters
- Reliable evaluation of architectural decisions
- Easy debug of a P4/C application on the NPU model

SystemC Model of NPU



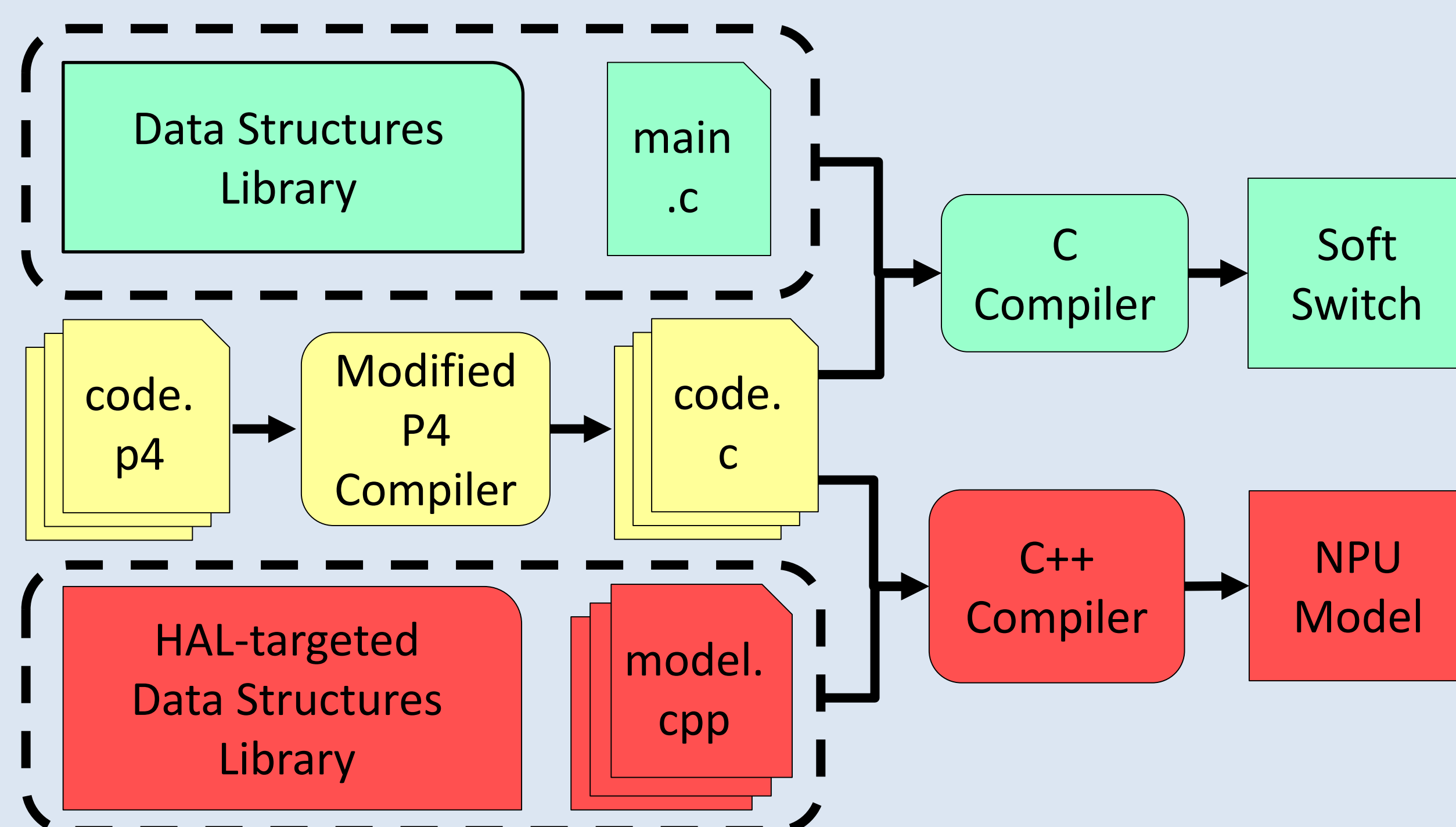
- Generic NPU architecture inspired by Ericsson SNP4000
- Hierarchy of SC_MODULES for HW structure
 - SC_MODULES contain SC_THREADS to model logic on PE
 - Run-to-complete application threads on core modules
- Memories → passive transaction-level SC_MODULE
 - Provide service through implemented interface
- On-chip routers and switches → special Router SC_MODULE
 - Routing config. entered separately during elaboration

Hardware Abstraction Layer Model



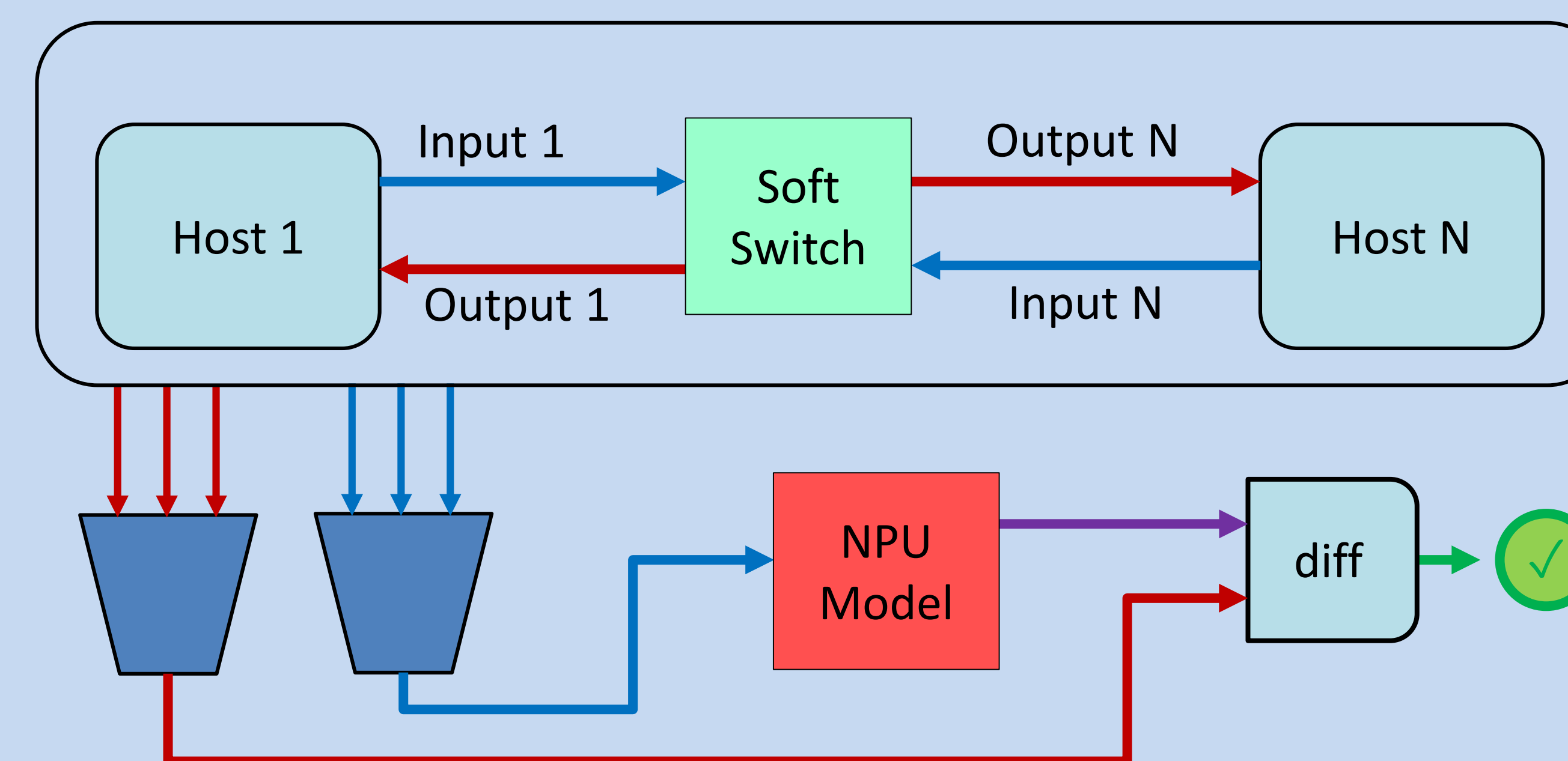
- Application threads are host-compiled with SystemC platform model
- HAL hides memory management from application
 - TLMVAR: base class for objects residing in target memory
 - TLMVAR operators are overloaded to access target memory model
 - All objects residing in target memories must derive from TLMVAR
- Trie library based on TLMVAR provided for *match* operations

P4 Application on NPU Model



- P4 compiler back end modified to expose simple API
 - Packet parsing function called from Parser module
 - Table application function called from app. threads
 - Deparsing function called from Deparser module
- P4 generated code linked to HAL-targeted Trie and HashTable
 - Used to simulate memory accesses during MAT processing

Model Validation

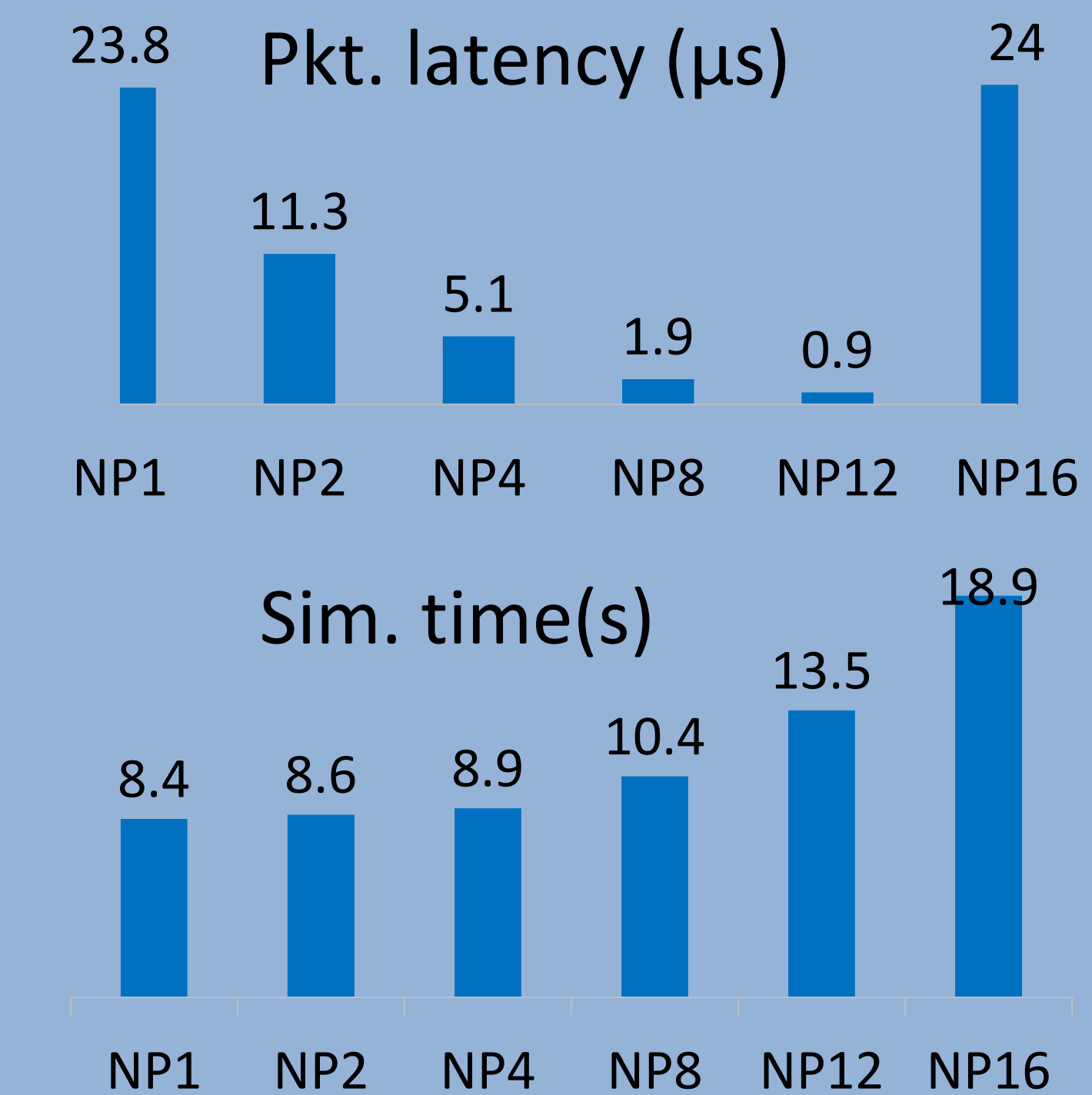


- Create a P4 soft switch inside Mininet
- Harpoon traffic generator to create test traffic
- Capture each interface's ingress and egress separately
- Merge ingress and egress streams chronologically
- Use ingress stream as input to NPU model
- Use egress stream to validate NPU model output

Experimental Results

Design	NP1	NP2	NP4	NP8	NP12	NP16
#cores	4	8	16	32	48	64
eDRAM	64K	32K	16K	8K	6K	4K

- 2.7 GHz, 8GB RAM simulation host
- Sample P4 application with 5 MATs
- Table size = 2048, 5K packets simulated
- NP16 tries spill over to off-chip memory
- Simulation time 17X – 40X of soft-switch



Next Steps

- Modeling of hardware accelerators
- Automatic timing/energy annotation in application threads
- Accuracy comparison vs. cycle-accurate simulator
- Improvements to simulation speed
- A P4 debugger for target model