

Productive Design of Extensible Cache Coherence Protocols

Henry Cook, Jonathan Bachrach,
Krste Asanovic

Par Lab Summer Retreat, Santa Cruz

June 1, 2011

- Cache coherence is important
 - Not just functionality, but performance and energy
 - Major implications for programming models
- Cache coherence is difficult
 - To implement and verify a single protocol
 - To explore design space of multiple protocols
- Hardware design in general is not productive
 - Often lacking modularity, extensibility, composability

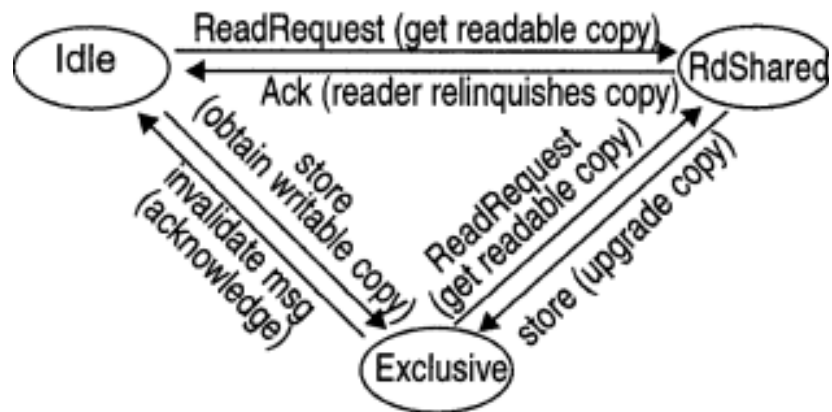
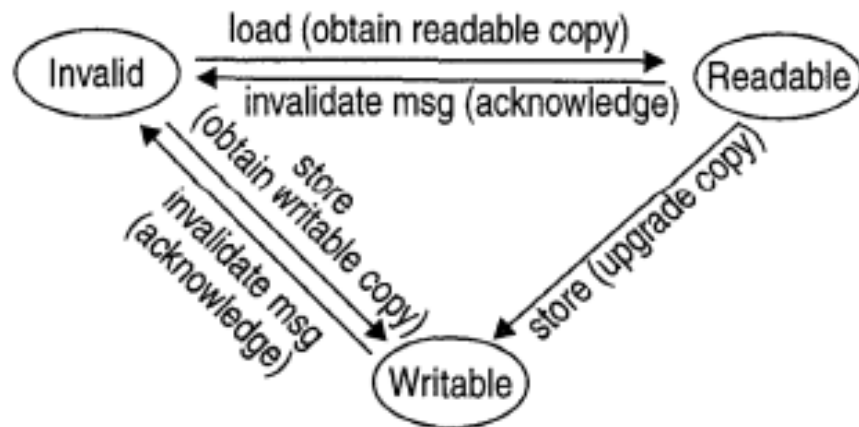
- Hypothesis: We can write protocols using succinct, declarative descriptions, and generate effective hardware implementations
 - Produce verified implementations from verified specifications
 - Experiment with more designs
- Hypothesis: Customization of protocol behavior is important for energy efficiency
 - On a per-motif or per-specializer basis
 - Heterogeneity in memory hierarchy

- **Constructing Hardware In a Scala Embedded Language**
- Embed a hardware-description language in Scala, using Scala's extension facilities
 - A hardware module is just a data structure in Scala
 - Different backends can generate different types of output (C, Verilog) from same Chisel representation
- Full power of Scala for writing hardware generators
 - Object-Oriented: Factory objects, traits, overloading etc
 - Functional: Higher-order funcs, anonymous funcs, currying
 - Compiles to JVM: Good performance, Java interoperability

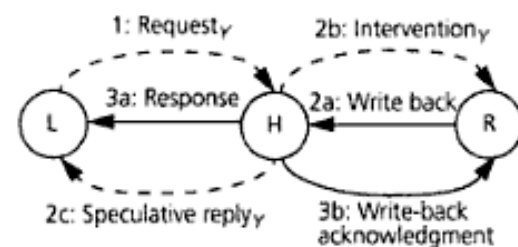
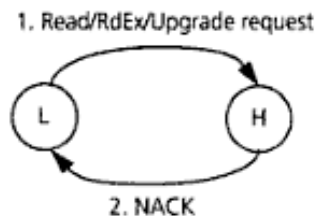
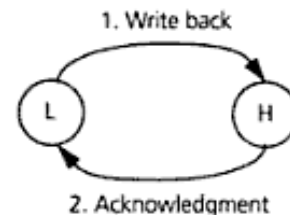
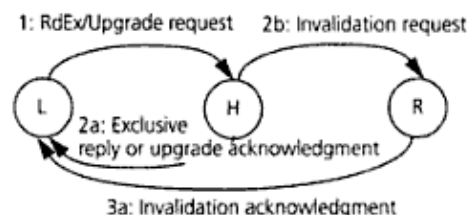
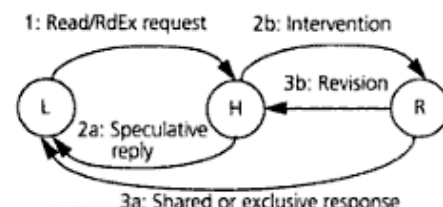
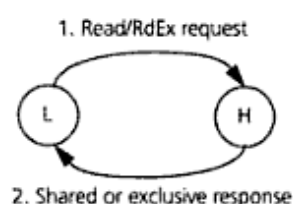
- Apply the best of SW design practices to HW design
 - Write reusable modules
 - Capture design patterns as generators
 - Declarative design and search
- Write it the way you do on the whiteboard

- Three views
- First view:
coherence protocol
as abstract state
machine

- Node types
- States
- Invariants



- Second view:
coherence protocol
as set of sequences
of request/reply
messages



- Set of all sequences
- Order of messages
in each sequence
- Messages, payloads

- Third view:
coherence protocol
as rule tables
- Given state and
input, emit messages
and update state
- Set of rules for each
node type

P2C request ⁴	deferQ State	CState	Action	Next CState
Load(<i>a</i>)	$a \in \text{deferQ}$	-	$\text{req} \rightarrow \text{deferQ}^1$	-
	$a \notin \text{deferQ}$	$\text{Cell}(a, v, \text{Sh})$	retire^2	$\text{Cell}(a, v, \text{Sh})$
	$a \notin \text{deferQ}$	$\text{Cell}(a, v, \text{Ex})$	retire^2	$\text{Cell}(a, v, \text{Ex})$
	$a \notin \text{deferQ}$	$\text{Cell}(a, -, \text{Pen})$	$\text{req} \rightarrow \text{deferQ}^1$	$\text{Cell}(a, -, \text{Pen})$
	$a \notin \text{deferQ}$	$a \notin \text{cache}$	$\text{if } \text{cmissQ.isNotFull then}$ $\langle \text{ShReq}, a, L \rangle \rightarrow \text{Mem},$ $\text{req} \rightarrow \text{cmissQ}$ else $\text{req} \rightarrow \text{deferQ}^1$	$\text{Cell}(a, -, \text{Pen})$
Store(<i>a, v</i>)	$a \in \text{deferQ}$	-	$\text{req} \rightarrow \text{deferQ}^1$	-
	$a \notin \text{deferQ}$	$\text{Cell}(a, -, \text{Sh})$	$\langle \text{Inv}, a, H \rangle \rightarrow \text{Mem},$ Keep req	$a \notin \text{cache}$
	$a \notin \text{deferQ}$	$\text{Cell}(a, -, \text{Ex})$	retire^2	$\text{Cell}(a, v, \text{Ex})$
	$a \notin \text{deferQ}$	$\text{Cell}(a, -, \text{Pen})$	$\text{req} \rightarrow \text{deferQ}^1$	$\text{Cell}(a, -, \text{Pen})$
	$a \notin \text{deferQ}$	$a \notin \text{cache}$	$\text{if } \text{cmissQ.isNotFull then}$ $\langle \text{ExReq}, a, L \rangle \rightarrow \text{Mem},$ $\text{req} \rightarrow \text{cmissQ}$ else $\text{req} \rightarrow \text{deferQ}^1$	$\text{Cell}(a, -, \text{Pen})$
voluntary rule	-	$\text{Cell}(a, -, \text{Sh})$	$\langle \text{Inv}, a, H \rangle \rightarrow \text{Mem}^3$	$a \notin \text{cache}$
	-	$\text{Cell}(a, v, \text{Ex})$	$\langle \text{WBI}, a, v, H \rangle \rightarrow \text{Mem}^3$	$a \notin \text{cache}$
	-	$\text{Cell}(a, v, \text{Ex})$	$\langle \text{WB}, a, v, H \rangle \rightarrow \text{Mem}^3$	$\text{Cell}(a, v, \text{Sh})$

¹ deferQ must not be full for this operation, otherwise, req will remain in the p2cQ

² retire means a response is sent to the requesting processor and the input request is deleted

³ c2mQ_H must not be full for this operation

⁴ The rules for handling deferQ requests are almost identical and not shown

Figure 3: Rules for Handling P2C Requests at Cache-site

C2M Message	Priority	MState	MDIR	Action	Next MState	Next MDIR
ShReq(<i>c, a</i>)	Low	-	\emptyset^1	$\langle \text{ShResp}, a, \text{Mem}[a] \rangle \rightarrow c,$ deq c2m Message	S	$\{c\}$
		S	$c \notin \text{MDIR}$	$\langle \text{ShResp}, a, \text{Mem}[a] \rangle \rightarrow c,$ deq c2m Message	S	$\text{MDIR} + \{c\}$
		E	$\{c'\}, c' \neq c$	$\langle \text{WBReq}, a \rangle \rightarrow c$	T	$\{c'\}$
ExReq(<i>c, a</i>)	Low	-	\emptyset^1	$\langle \text{ExResp}, a, \text{Mem}[a] \rangle \rightarrow c,$ deq c2m Message	E	$\{c\}$
		S	$c \notin \text{MDIR}$	$\forall c' \in \text{MDIR}. \langle \text{InvReq}, a \rangle \rightarrow c',$	T	MDIR
		E	$\{c'\}, c' \neq c$	$\langle \text{WBReq}, a \rangle \rightarrow c'$	T	$\{c'\}$
Inv(<i>c, a</i>)	High	<i>mstate</i>	$c \in \text{MDIR}$	deq c2m Message	<i>mstate</i>	$\text{MDIR} - \{c\}$
WBI(<i>c, a, v</i>)	High	T E	$\{c\}$	$\text{Mem}[a] := v,$ deq c2m Message	S	\emptyset
WB(<i>c, a, v</i>)	High	T E	$\{c\}$	$\text{Mem}[a] := v,$ deq c2m Message	S	$\{c\}$

¹ any state with MDIR = \emptyset is treated as S with \emptyset

Figure 4: Rules for Handling Cache Messages at Memory-site

- We can verify this protocol's rules, but are there additional sources of complexity?
- Turning message sequences into transactions
- Making multi-step, intra-node behavior atomic

- Complexity at the inter-node message sequence level
 - Interleave messages re: particular block
 - Add transient/busy states to protocol
 - Handle races
 - Provide write serializability and atomicity
 - Avoid deadlock, livelock, starvation
 - Address externalities: type of network used, amount of message buffering available

- Complexity at the intra-node atomicity level
 - Arbitrating for finite number of SRAM ports
 - Dequeueing and buffering requests
 - Enqueueing requests and responses
 - Filling and draining MSHRs
 - Multi-cycle ops with potentially conflicting updates lead to additional transient states
- Any modification to deal with the above (or area/timing constraints) could render original verification work useless

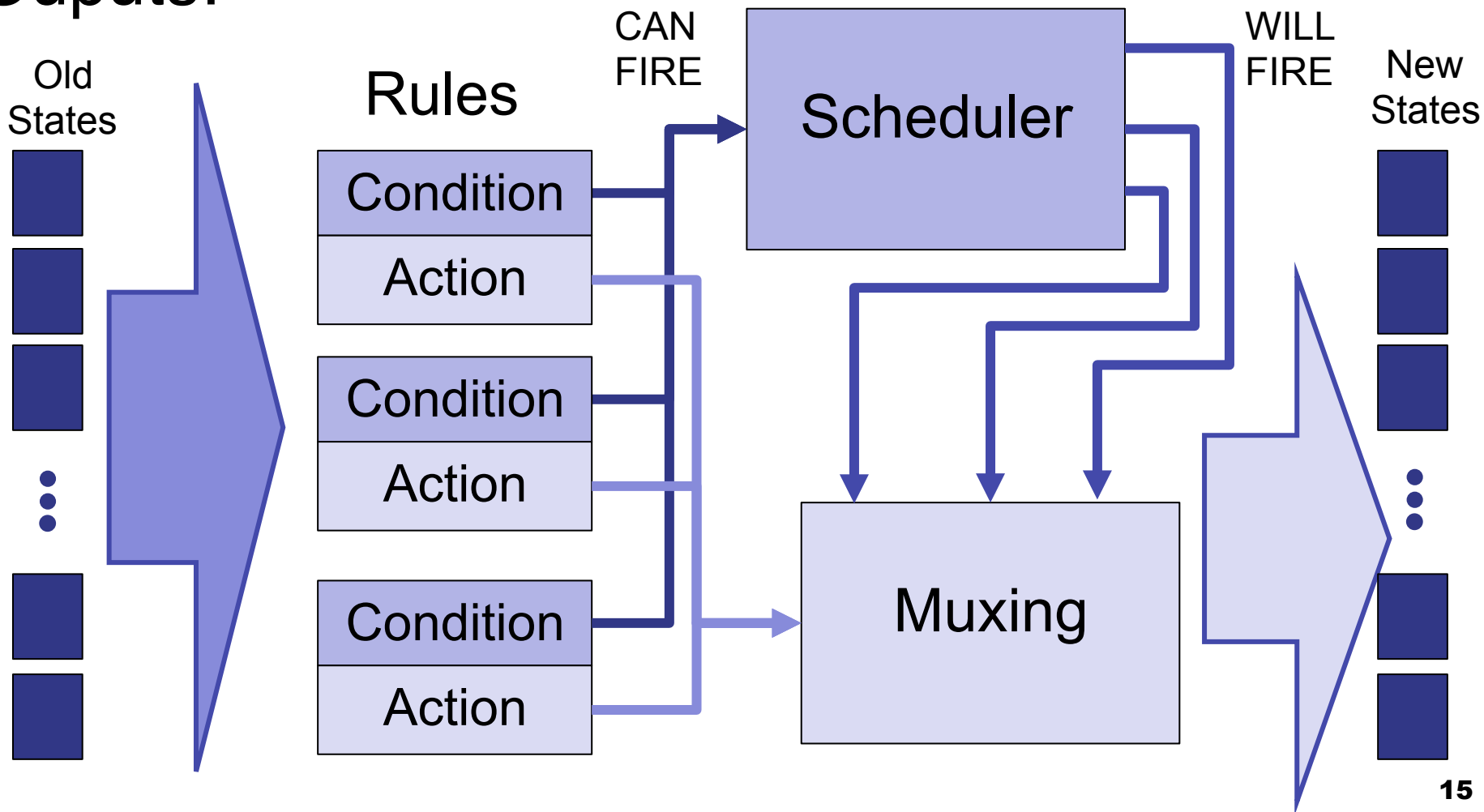
- Address intra-node complexity using **BlueChisel**, a declarative, embedded DSL built on top of Chisel
- Goal: Generate not only control logic for protocol-defined activity but also:
 - Arbitration logic for access to SRAM ports
 - Skid buffers and queuing logic
 - Logic implementing intermediate/transient protocol states

- High-level, functional HDL compiled to a term rewriting system and translated into HW
 - Natural way to describe many HW devices
 - Understandable, well-defined semantics
 - Conditional atomic execution of state updates, based on rules
 - Guarded atomic actions
 - Scheduler dynamically tries to fire as many as possible
- Limitations:
 - In general, guarded atomic actions are a productive abstraction in some cases, but not in others
 - Rules can only express actions that take single cycle

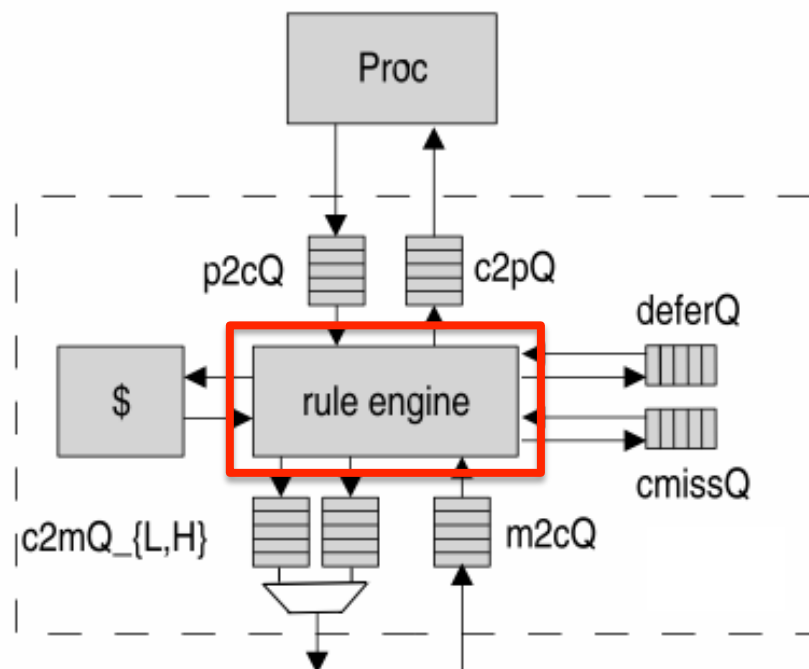
- Core functionality: conditional evaluation of rules leads to atomic state updates
- For cache coherence, automatically generate:
 - Extra transient/implicit protocol states
 - Additional rules to govern multi-cycle operations
 - Fairness and rule priority with urgency annotations
- Extension built on top of Chisel, which we choose to apply where appropriate

Inputs: `rule (cond) { updates ... }`

Outputs:

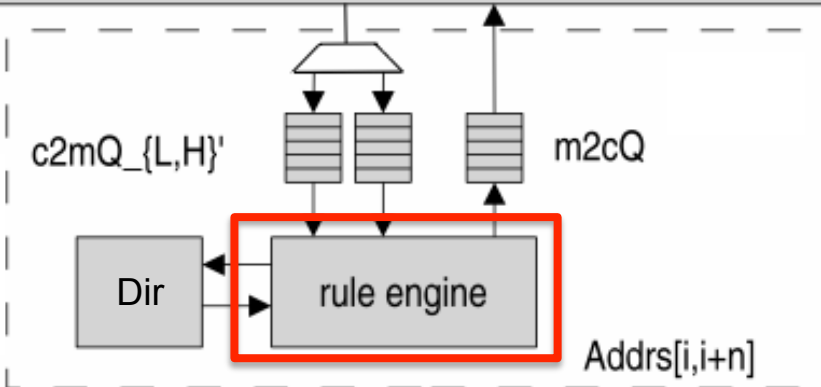


Cache controller
protocol rule engine



...

Directory controller
protocol rule engine



...



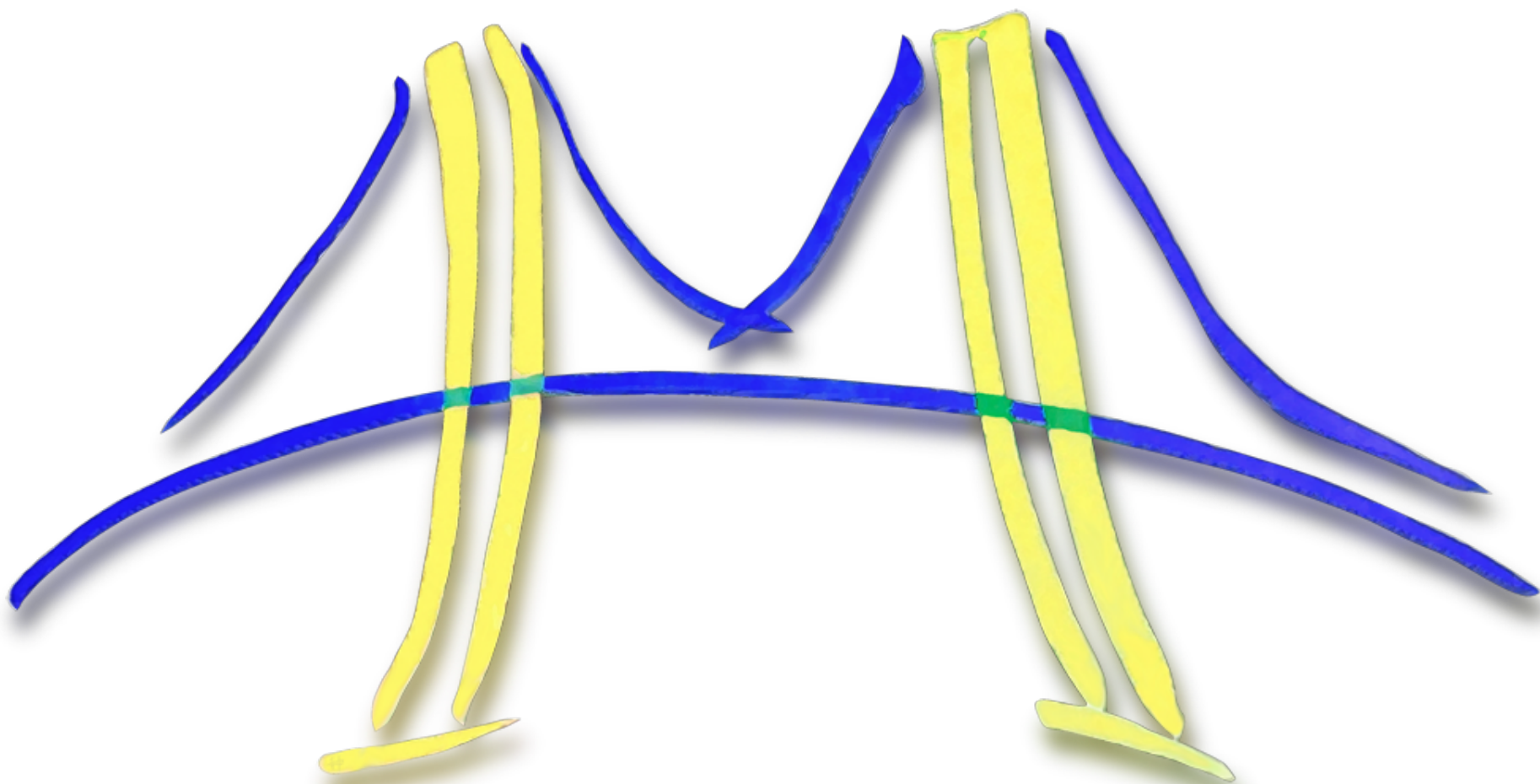
- Try to address inter-node, message sequence interleaving complexity
 - Generate sufficient transient/busy states
 - Compatible physical networks and buffering
- High-level composition of new transactions with existing protocols
 - “MESI” = “MI” + “E” + “S” + ?
- Patterns in message sequences



- Protocol extensions that
 - Exploit SW knowledge via explicit SW->HW directives
 - Allow for heterogeneous memory hierarchy behavior
- Assignment of data to particular state
 - Clean this cache block
- Assignment of data to particular sub-protocol
 - Keep this block coherent using an update protocol
- Exemption of data from any protocol other than SW-directed actions
 - Only move this block in response to a DMA command

- Programs have been found to employ a set of common types of sharing behavior
 - Write-once, Private, Write-many, Result, Synchronization, Migratory, Producer/Consumer, Read-mostly, Streaming
- Sharing behavior is often known by expert or even application programmer
 - Utilize high-level info instead of reconstructing in HW
- SEJITS: emit optimized code from high level abstractions
 - Code that can control cache behavior according to pattern
 - Even define user-level protocols
- Bloom: Use consistency-analysis to inform decisions about which sub-protocols to employ

- We can write protocols using succinct, declarative descriptions, and generate effective hardware implementations
 - Declarative extension to Chisel based on ideas from Bluespec
 - Produce verified implementations from verified specifications
 - Experiment with more designs
- Explore space of protocols that use SW input to create heterogeneous behavior



- **Murphi**
 - DSL for finite-state analysis/model checking
- **Teapot**
 - DSL for user-space software coherence protocols
- **SLICC**
 - DSL for emitting SW modules for GEMS simulator
- **Bluespec SystemVerilog**
 - HDL based on guarded atomic actions
- **Bloom**
 - DSL for high-level consistency analysis of distributed parallel algorithms

- For now, left up to enforcement at processor by compiler-issued ISA constructs (e.g. fences)
- In the future, would like to consider protocols that exploit very relaxed consistency models (e.g. location consistency)
- What patterns? What whiteboard design?

- Inputs:

```
rule ( cond ) { update ... }
```

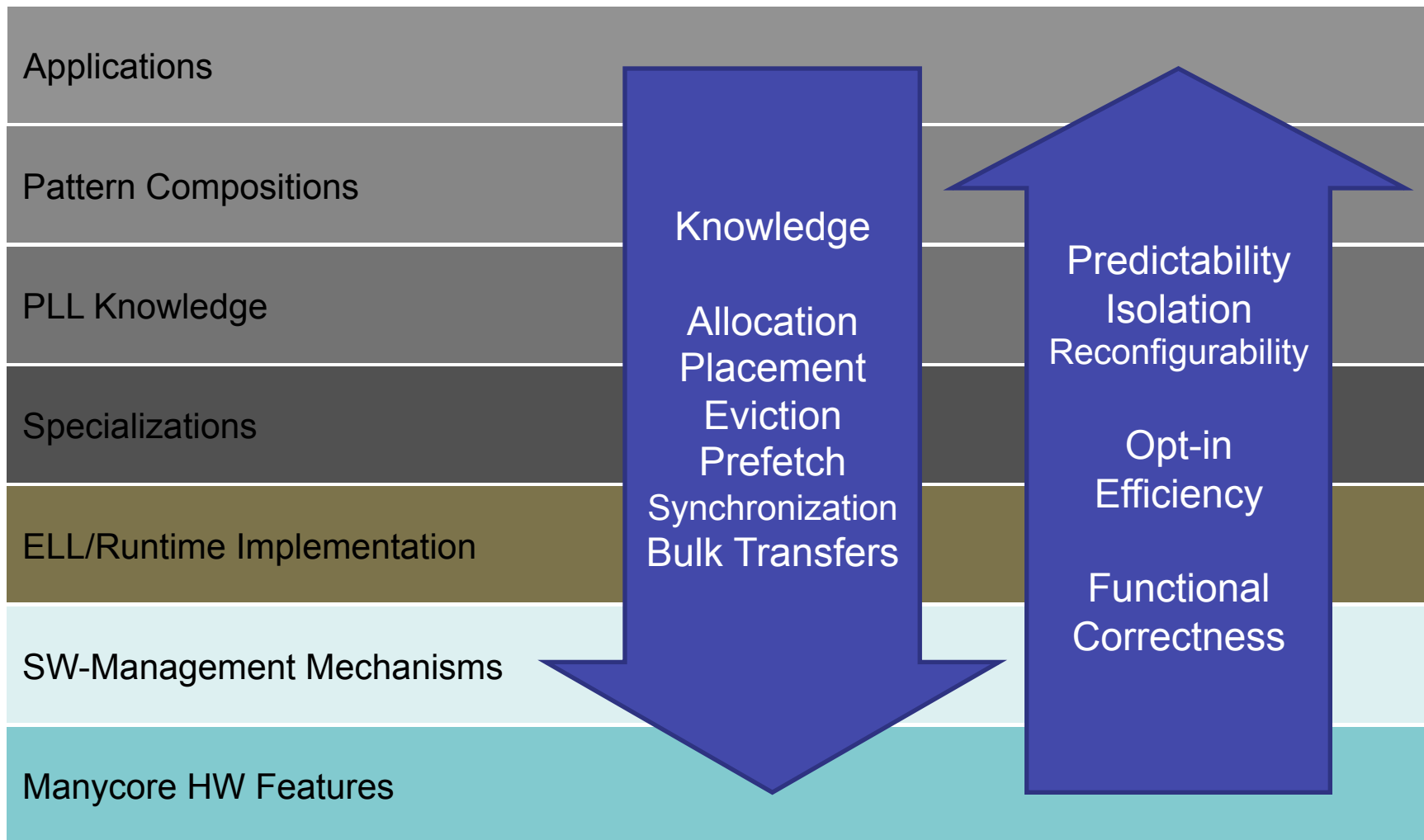
- Outputs:

- Rule engines, consisting of scheduler and muxing logic, that conditional modify state

- CHISEL internals:

- During creation, log rule with component
- During elaboration, analyze domain and range of rules to create scheduler and mux code

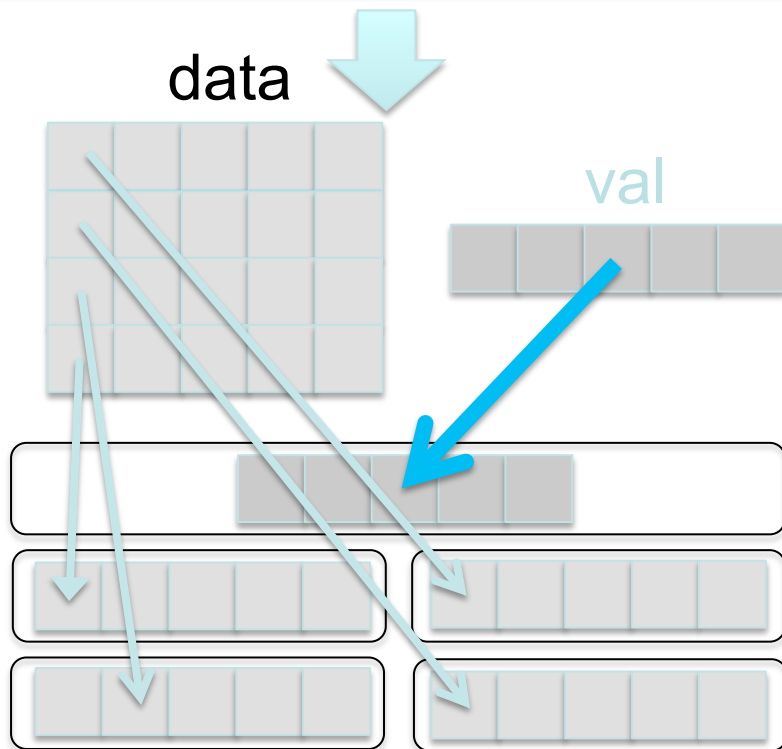
- What are the abstractions?
 - Data (states, messages, payloads)
 - Actions (send messages, update states)
 - Rules
- What are the reusable modules?
 - Collections of rules encapsulated in rule engines
- What are the design patterns?
 - Control logic for queues, arbiters, MSHRs



- **Write-once**
 - Initialized but then only read
 - Best supported by replication (selected portions of large objects)
- **Private**
 - Need not be managed
 - On violation, demote/activate management?
- **Write-many**
 - Frequently modified by multiple threads between synch points
 - Use delayed-update protocol
- **Result**
 - Restricted subset of write-many, no reads until all writes complete
 - Lack of conflicts allows maximum utilization of delayed-update protocol
- **Synchronization**
 - Distributed locks, atomic operands
- **Migratory**
 - Read and written by single thread at a time, as object in critical sections of code
 - Associate w/ lock movement, look for signature pattern
- **Producer/Consumer**
 - Produced by one thread and consumed by fixed set of other threads
 - Eager object movement in update protocol
- **Read-mostly**
 - Replicate and update infrequently via broadcast
- **Streaming**
 - Read once (or few times), too large to keep in particular level for reuse
- **General read/write**
 - Default, rare

Case Study 1: SVM Training

```
def train (val, data)
  def val_compare(x):
    return compare(x, val)
  z = map(val_compare, data)
```

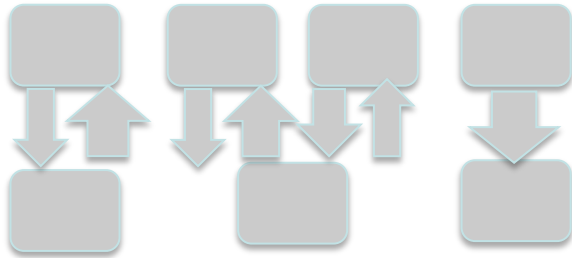


- Producer/consumer
- Update protocol with proactive transmission of ghost cells to static neighbors

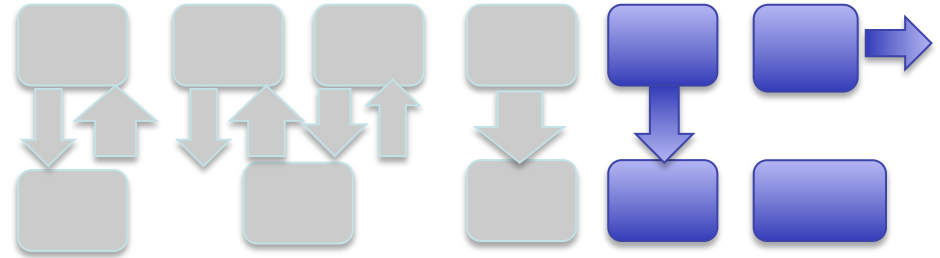
Case Study 3: nbody

- Migratory
- Migrate on read miss rather than replicate

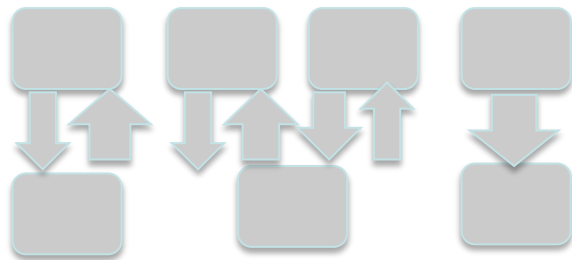
- Known at compile time (prod/con)
- Fixed per run (prod/con)
- Dynamic (migratory)



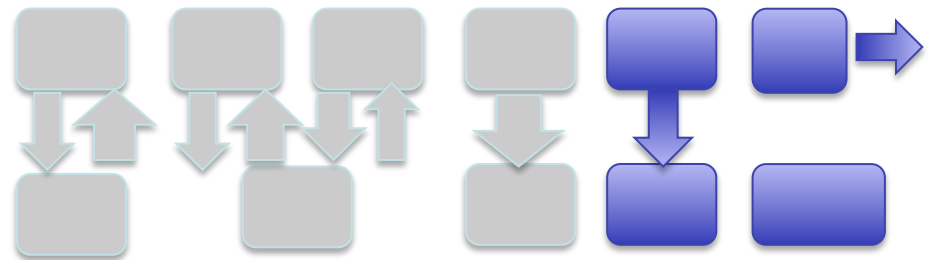
MI



MEI



MSI



MESI

