



Crossing Guard: Mediating Host-Accelerator Coherence Interactions



Lena E. Olson*, Mark D. Hill, David A. Wood

University of Wisconsin-Madison

* Now at Google

ASPLOS 2017
April 10th, 2017





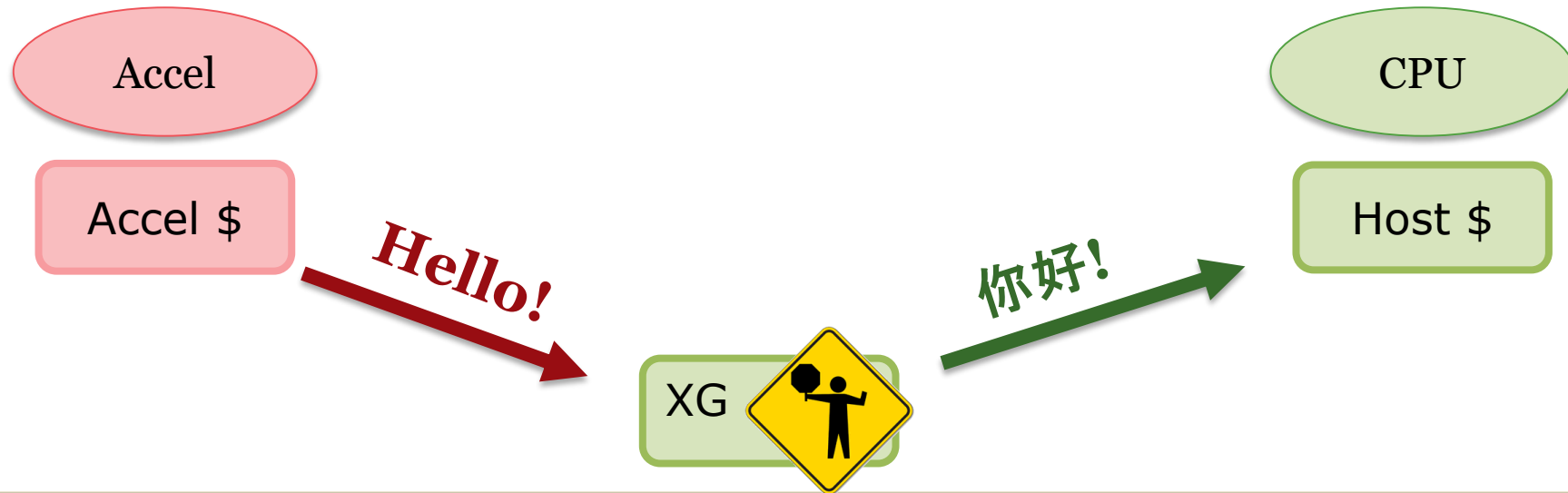
Accelerators are here!

- Complex, programmable accelerators increasingly prevalent
- **Many applications:** graphics, scientific computing, video encoding, machine learning, etc...
- Accelerators may benefit from **cache coherent shared memory**
- May be designed by **third parties**



However...

- Host coherence protocols may be **proprietary** and **complex**
- Bugs in accelerator implementations might **crash host system!**
- **Crossing Guard**: coherence interface to **safely** translate accelerator \leftrightarrow host protocol





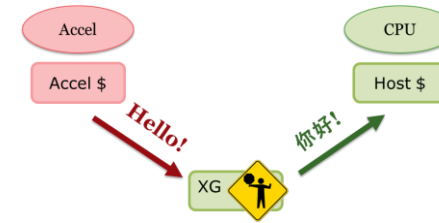
Outline



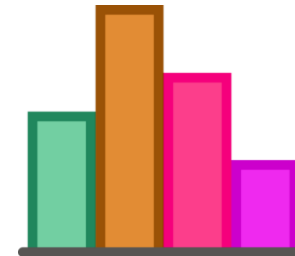
Goals



Guarantees



Design



Evaluation



Crossing Guard Goals

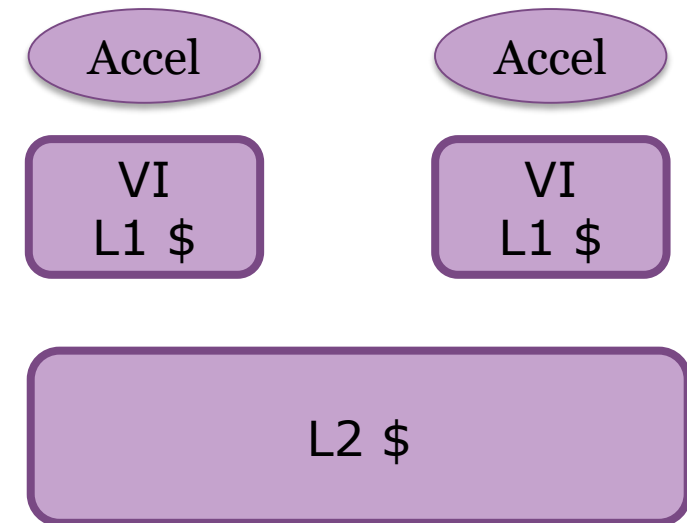
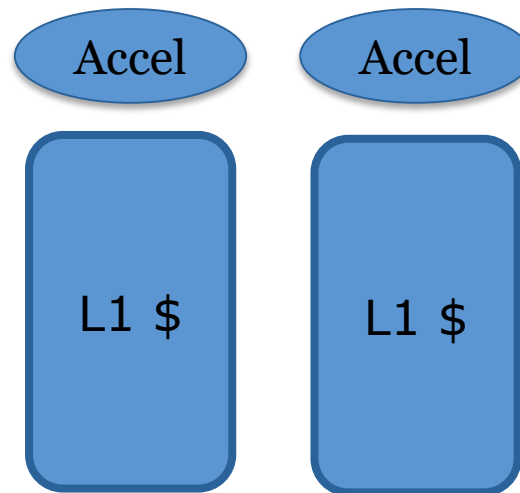
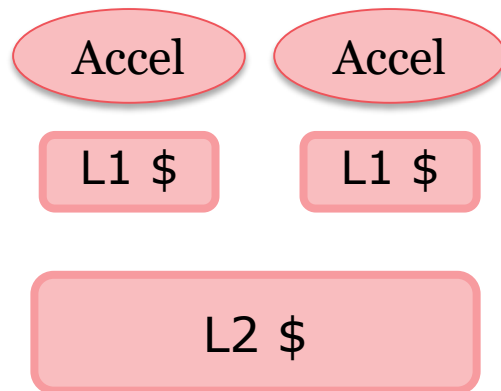
When adding accelerators to host coherence protocol:

1. Allow accelerators **customized caches**
2. **Simple, standardized** accelerator coherence interface
3. Guarantee **safety** for the host system



1. Why Customize Caches?

- CPU caches have to work with **most types** of **workloads**
- Accelerators may only run **some workloads**!
 - Optimize caches for likely data access patterns
 - Number of levels, writeback vs. writethrough, MSI vs VI, etc.

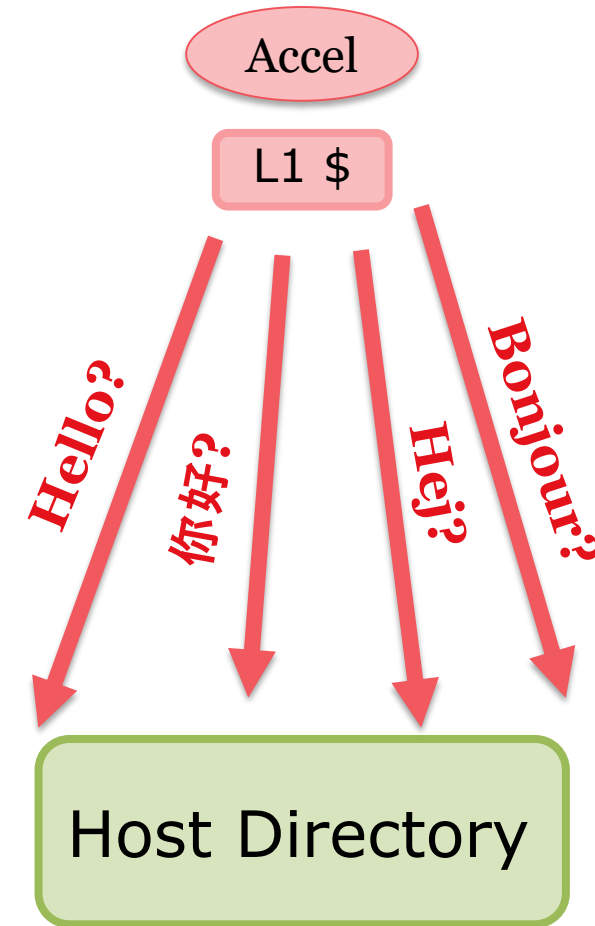




2. Why Simple, Standardized Interface?

Host systems speak different protocols...

- Expensive to redesign for each one!
 - Intel, AMD, ARM, IBM, Oracle...
 - CCIX shows industry cares!





2. Why Simple, Standardized Interface?

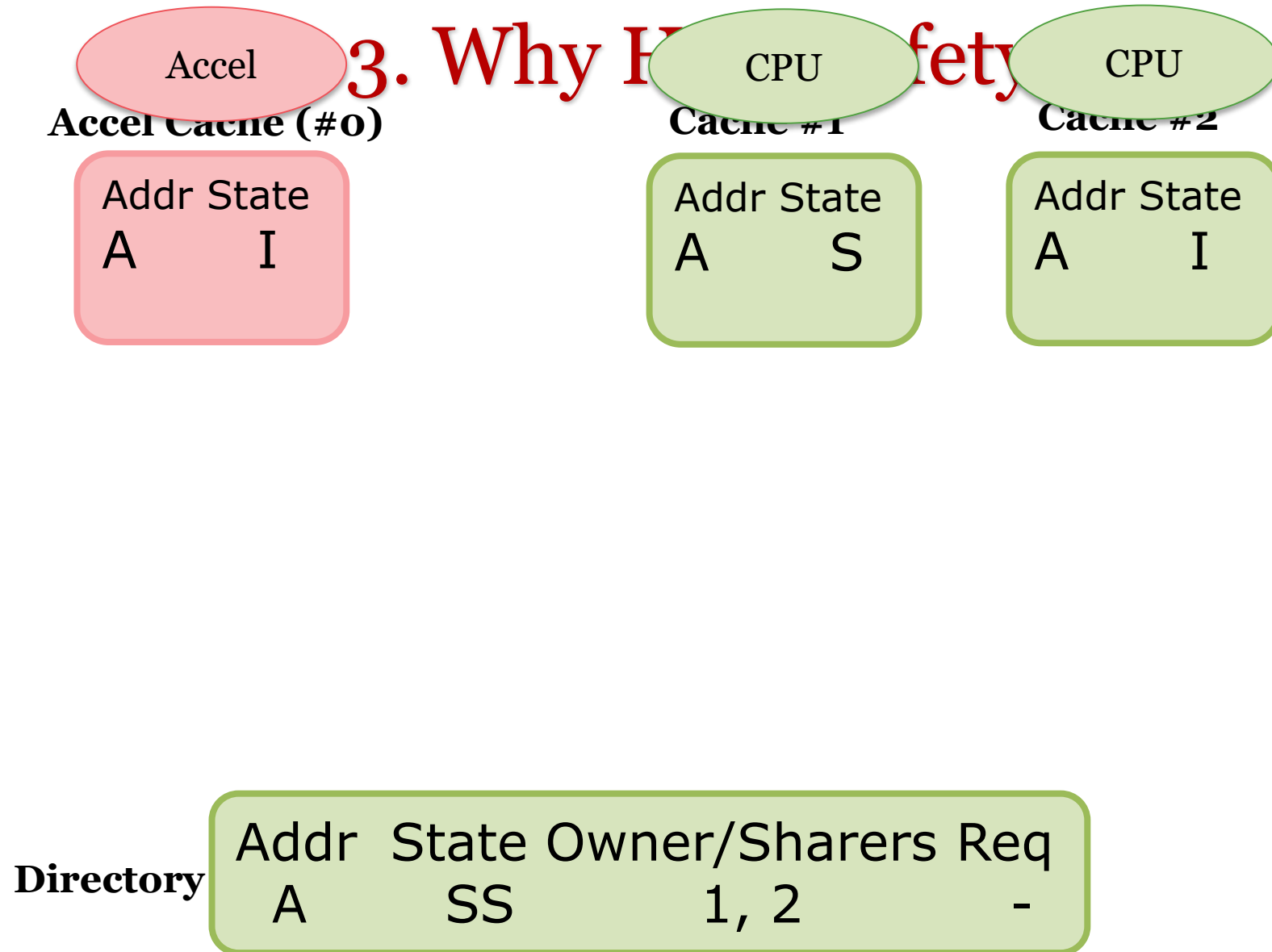
L1 controller from gem5's MOESI_hammer

Events

States

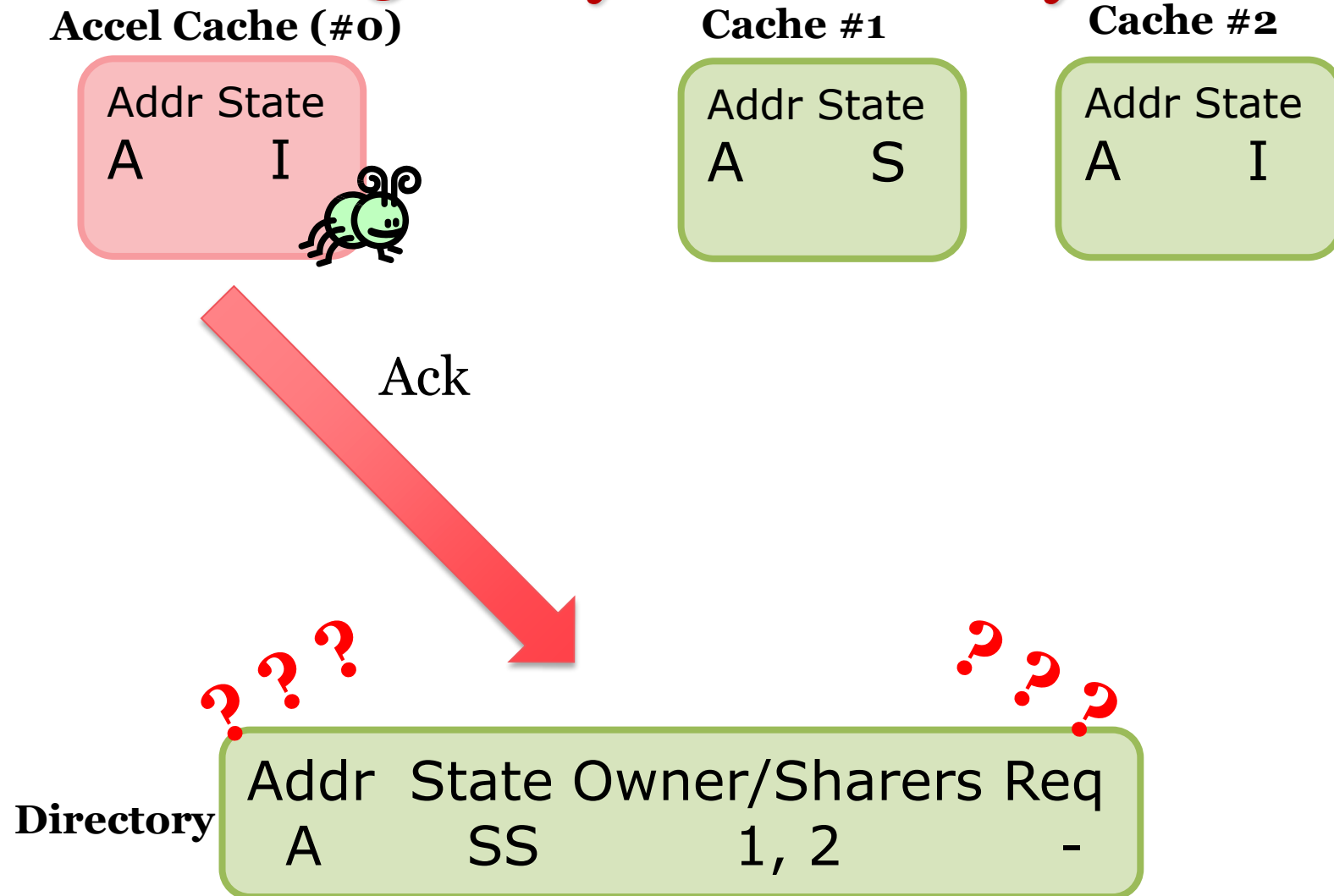
	Fetch	Load	Store	Invalidate	Other GETS only	Other GETS	Merged GETS	Other GETX	Ack	Shared Ack	Data	Shared Data	Exclusive Data	Writeback Ack	Writeback Nack	All acks	All acks no sharers	L2 Replacement	L1 to L2	Trigger L2 to L1	Trigger L2 to L1	Complete to L1
I	Invalid/IS	Invalid/IS	Invalid/IM	I	I	I		I														
S	Invalid/S	Invalid/S	Invalid/SM	I	I	I		I										Invalid/S	Invalid/S	Invalid/S	Invalid/S	
O	Invalid/O	Invalid/O	Invalid/OM	I	I	I	SM	I										Invalid/O	Invalid/O	Invalid/O	Invalid/O	
M	Invalid/M	Invalid/M	Invalid/MM	I	I/O	I/O	SM/O	I										Invalid/M	Invalid/M	Invalid/M	Invalid/M	
MM	Invalid/MM	Invalid/MM	Invalid/MM	I	I/O	I/O	SM/O	I										Invalid/MM	Invalid/MM	Invalid/MM	Invalid/MM	
IR	Invalid/IS	Invalid/IS	Invalid/IM	Z	Z	Z	Z	Z														
SR	Invalid/S	Invalid/S	Invalid/SM	Z	Z	Z	Z	Z														
OR	Invalid/O	Invalid/O	Invalid/OM	Z	Z	Z	Z	Z														
MR	Invalid/M	Invalid/M	Invalid/MM	Z	Z	Z	Z	Z														
MMR	Invalid/MM	Invalid/MM	Invalid/MM	Z	Z	Z	Z	Z														
IM	Z	Z	Z	I	I	I		I	mon		mon/ISM		mon/ISM					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
SM	Invalid/S	Invalid/S	Z	I	I	I		I	mon		mon/ISM		mon/ISM					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
OM	Invalid/O	Invalid/O	Z	I	I	I	SM	I	mon		mon/ISM		mon/ISM					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
ISM	Invalid/M	Invalid/M	Z						mon		mon/ISM		mon/ISM					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
MW	Invalid/M	Invalid/M	Invalid/MM						mon	mon	mon	mon	mon					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
MMW	Invalid/MM	Invalid/MM	Invalid/MM						mon	mon	mon	mon	mon					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
IS	Z	Z	Z	I	I	I		I	mon	mon	mon/ISM	mon/ISM	mon/ISM					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
SS	Invalid/S	Invalid/S	Z						mon	mon	mon	mon	mon					mon/ISM	mon/ISM	mon/ISM	mon/ISM	
OI	Z	Z	Z	I/O	I/O	I/O	SM	I/O										mon/ISM	mon/ISM	mon/ISM	mon/ISM	
MI	Z	Z	Z	I/O	I/O	I/O	SM/O	I/O										mon/ISM	mon/ISM	mon/ISM	mon/ISM	
II	Z	Z	Z	I	I	I		I										mon/ISM	mon/ISM	mon/ISM	mon/ISM	
ST	Z	Z	Z	Z	Z	Z	Z	Z													mon/ISM	
OT	Z	Z	Z	Z	Z	Z	Z	Z													mon/ISM	
MT	Z	Z	Z	Z	Z	Z	Z	Z													mon/ISM	
MMT	Z	Z	Z	Z	Z	Z	Z	Z													mon/ISM	

(Transition table in style of Sorin et al.)



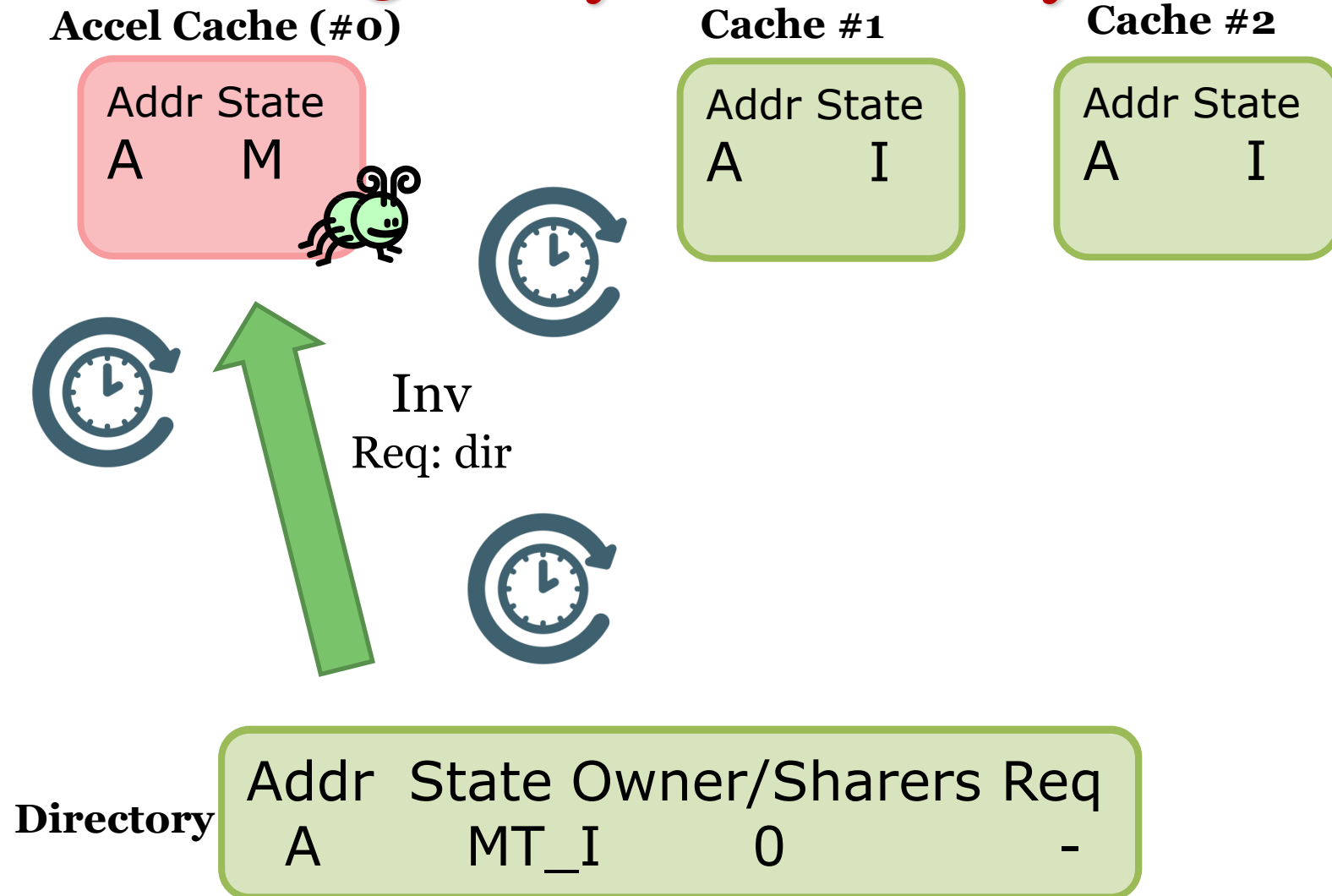


3. Why Host Safety?





3. Why Host Safety?





Outline



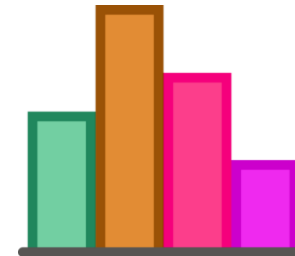
Goals



Guarantees



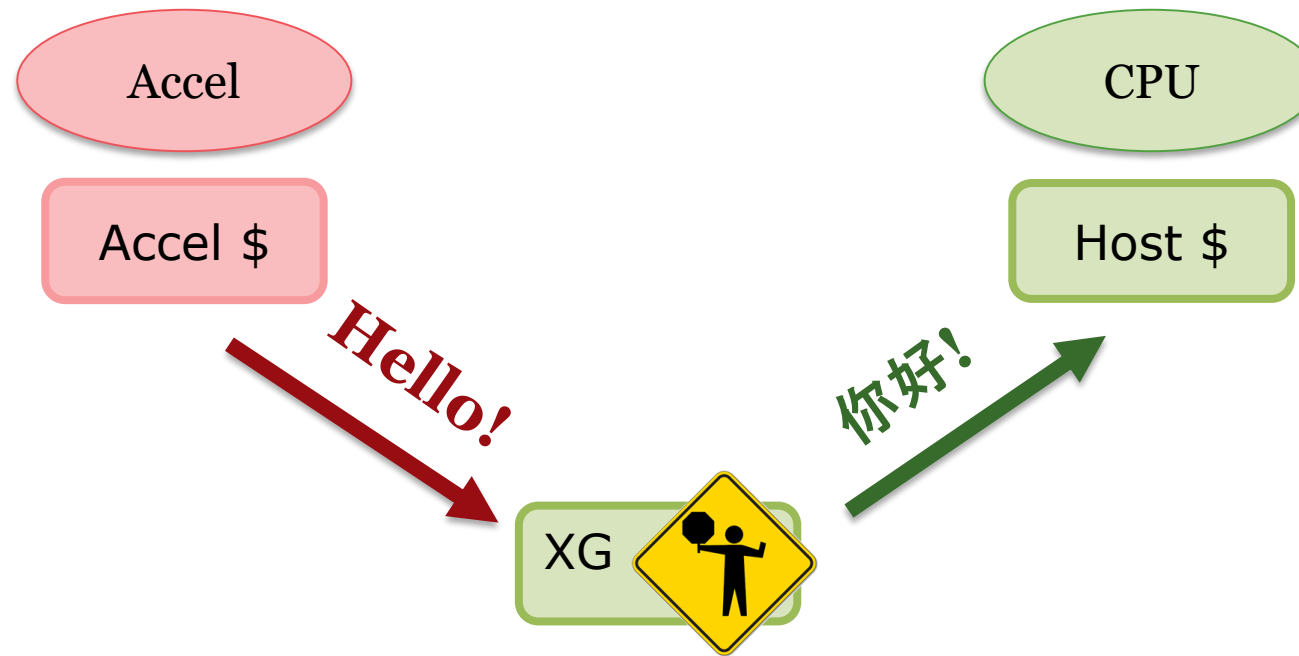
Design



Evaluation

Crossing Guard

- Hardware translating between host and accelerator protocols



- Set of accelerator \leftrightarrow host coherence messages (like an API)



Crossing Guard Interface

Accelerator → Host Requests

- GetS, GetM
- PutS, PutE, PutM

Host → Accelerator Requests

- Invalidate

Host → Accelerator Responses

- DataS, DataE, DataM
- Writeback Ack

Accelerator → Host Responses

- InvAck, Clean Writeback, Dirty Writeback



Crossing Guard

- Hides implementation details of host protocol
 - No counting acks, sending unblocks, handling races, etc.
- Moves protocol complexity into Crossing Guard hardware
 - Only implemented once per host system
 - By experts!



Experimental Implementation

- Coherence controllers / protocols implemented in slicc
- Simulations using gem5
- Code and transition tables available online

<http://research.cs.wisc.edu/multifacet/xguard/>



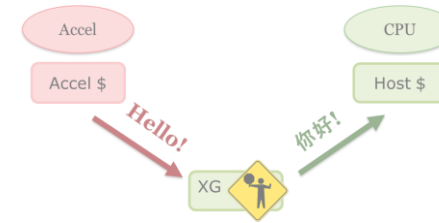
Outline



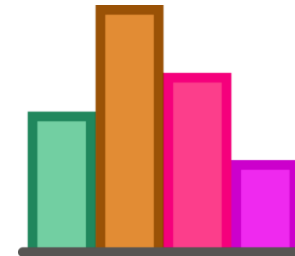
Goals



Guarantees



Design

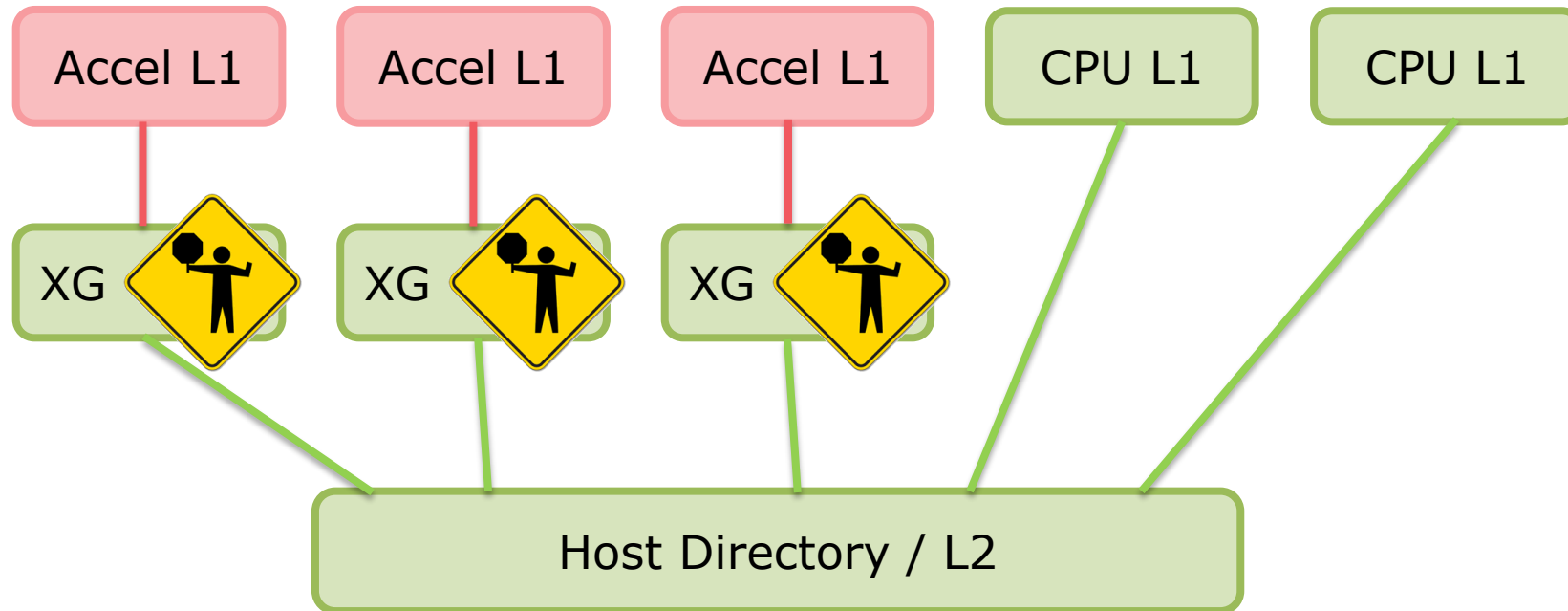


Evaluation

1. Customize Caches ✓

- Designed + implemented two sample systems

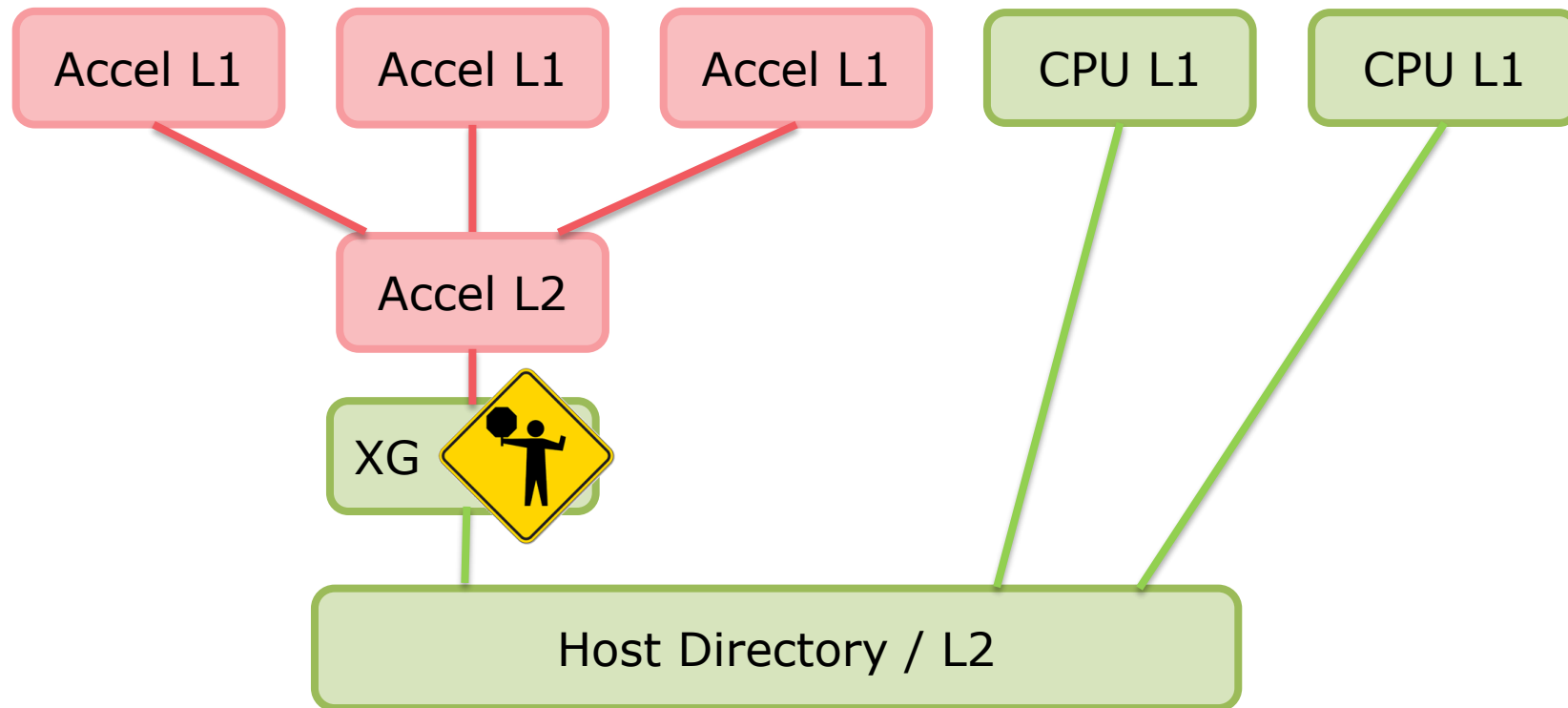
Private Per-Core L1 at Accelerator



1. Customize Caches ✓

- Designed + implemented two sample systems

Private L1s + Shared L2 at Accelerator





2. Simple, Standardized Interface ✓

	Fetch	Load	Store	Invalidate	Other GETS only	Other GETS	Merged GETS	Other GETN	Ack	Shared Ack	Data	Shared Data	Exclusive Data	Writeback Ack	Writeback Nack	All acks	All acks no sharers	L2 Replacement	L1 to L2	Trigger L2 to L1D	Trigger L2 to L1I	Complete L2 to L1
	<u>Load</u>	<u>Store</u>	<u>Replacement</u>	<u>Invalidate</u>																		
M	h q	hh q	j v c m l / B	d l m / I																		
E	h q	hh q / M	j v c e l / B	d c l m / I																		
S	h q	s j b / B	c s j v l / B	f l m / I																		
I	i j a / B	i j b / B		f m																		
B	z z	z z		f m							u k n / M	w u k n / E	w u k n / S	k n / I								
OM	h m h k	h m h k	z	c c c l / I M	e l	e l	c m l	c c c l / I M	m o	l l	k k m o n		k k m o n			e s t s m s i k d / M M	e s t s m s i k d / M M	z	z			
ISM	h m h k	h m h k	z						m o	l l	k k m o n		k k m o n			e s t s m s i k d / M M	e s t s m s i k d / M M	z	z			
MW	h m h k	h m h k	h m h k						m o	l l	k k m o n		k k m o n			e s t s m s i k d / M M	e s t s m s i k d / M M	z	z			
MMW	h m h k	h m h k	h m h k						m o	l l	k k m o n		k k m o n			e s t s m s i k d / M M	e s t s m s i k d / M M	z	z			
IS	z	z	z	l l					m o	l l	k k m o n		k k m o n			e s t s m s i k d / M M	e s t s m s i k d / M M	z	z			
SS	h m h k	h m h k	z						m o	l l	k k m o n		k k m o n			e s t s m s i k d / M M	e s t s m s i k d / M M	z	z			
OI	z	z	z	u l / I l	u l	u l	u m l	u l / I l						q s i k d / I	k k s i k d / I			z	z			
MI	z	z	z	u l / I l	u l / O l	u l / O l	u m l / O l	u l / I l						q s i k d / I	k k s i k d / I			z	z			
II	z	z	z	l l	l l	l l	l l	l l						q s i k d / I	k k s i k d / I			z	z			
ST	z	z	z	z	z	z	z	z										z	z			i k d / S R
OI	z	z	z	z	z	z	z	z										z	z			i k d / O R
MI	z	z	z	z	z	z	z	z										z	z			i k d / M R
MMI	z	z	z	z	z	z	z	z										z	z			i k d / M M R

Single-level Accelerator Cache using
Crossing Guard Interface

Controller	States	Transitions
AMD Hammer-like Private \$\$	24	148

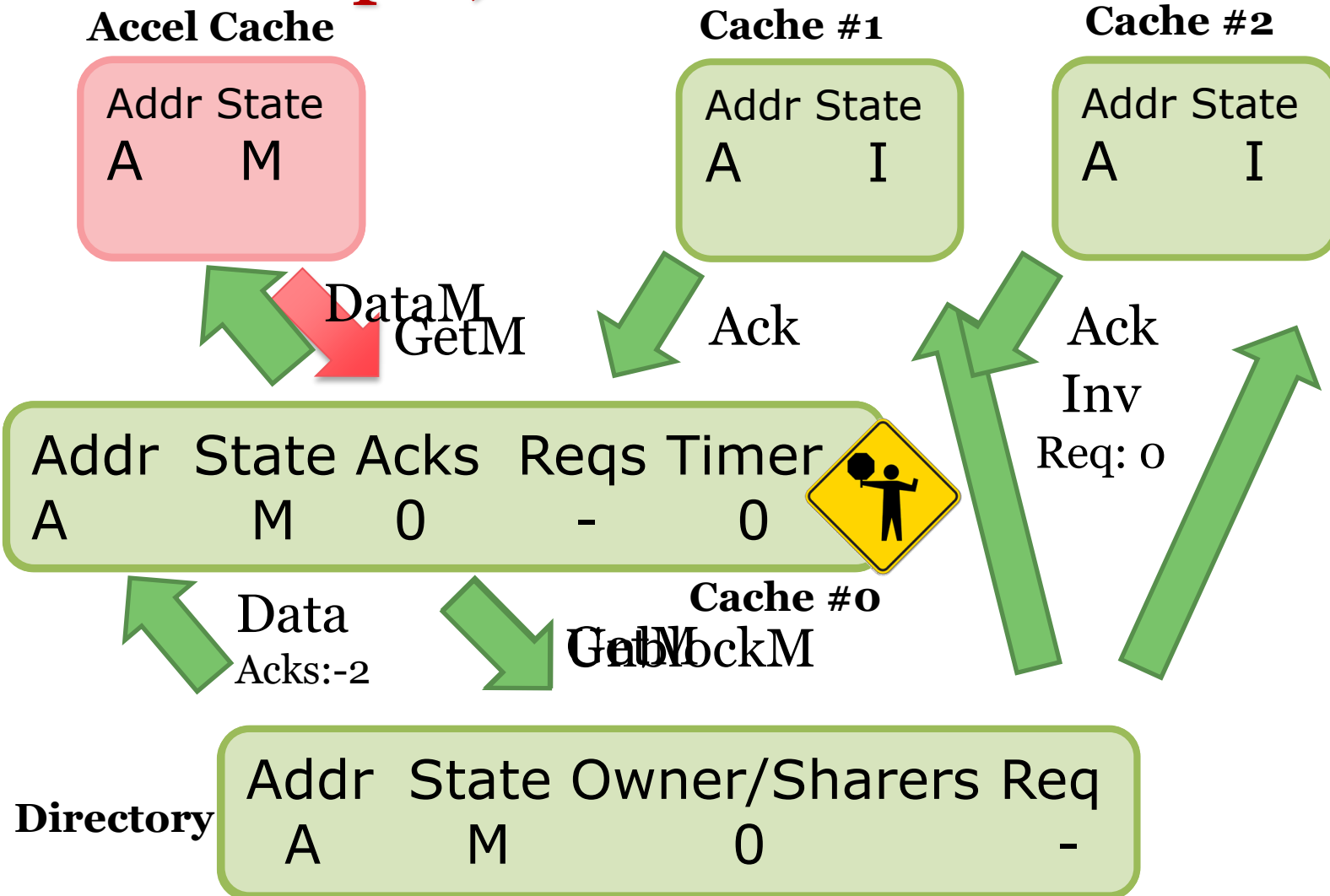


2. Simple, Standardized Interface ✓

- Implemented Crossing Guard controller for two host protocols
 - AMD Hammer-like Exclusive MOESI
 - Two-Level MESI Inclusive
- **Modularity:** Host and Accelerator protocol choice is completely independent

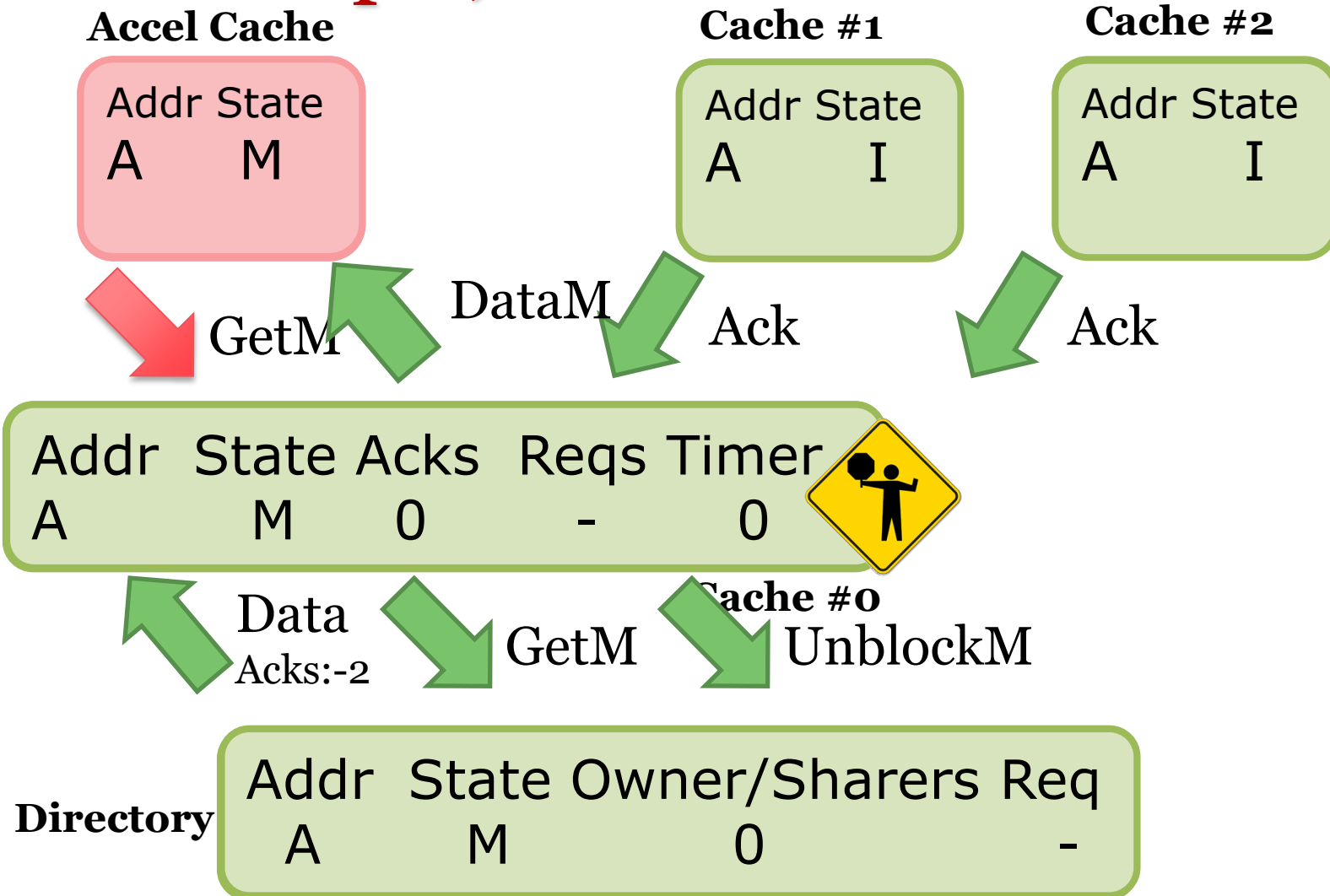


2. Simple, Standardized Interface ✓



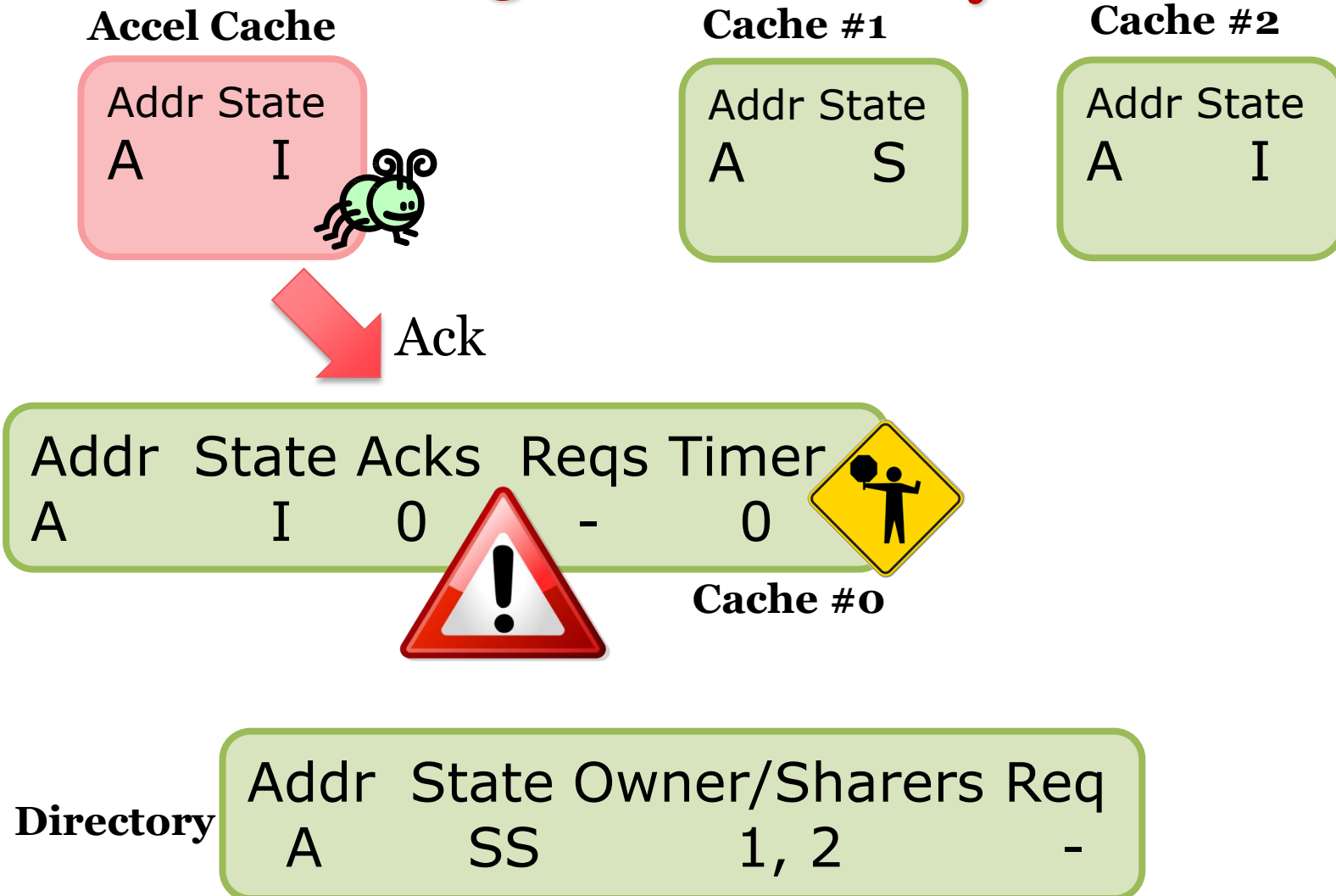


2. Simple, Standardized Interface ✓



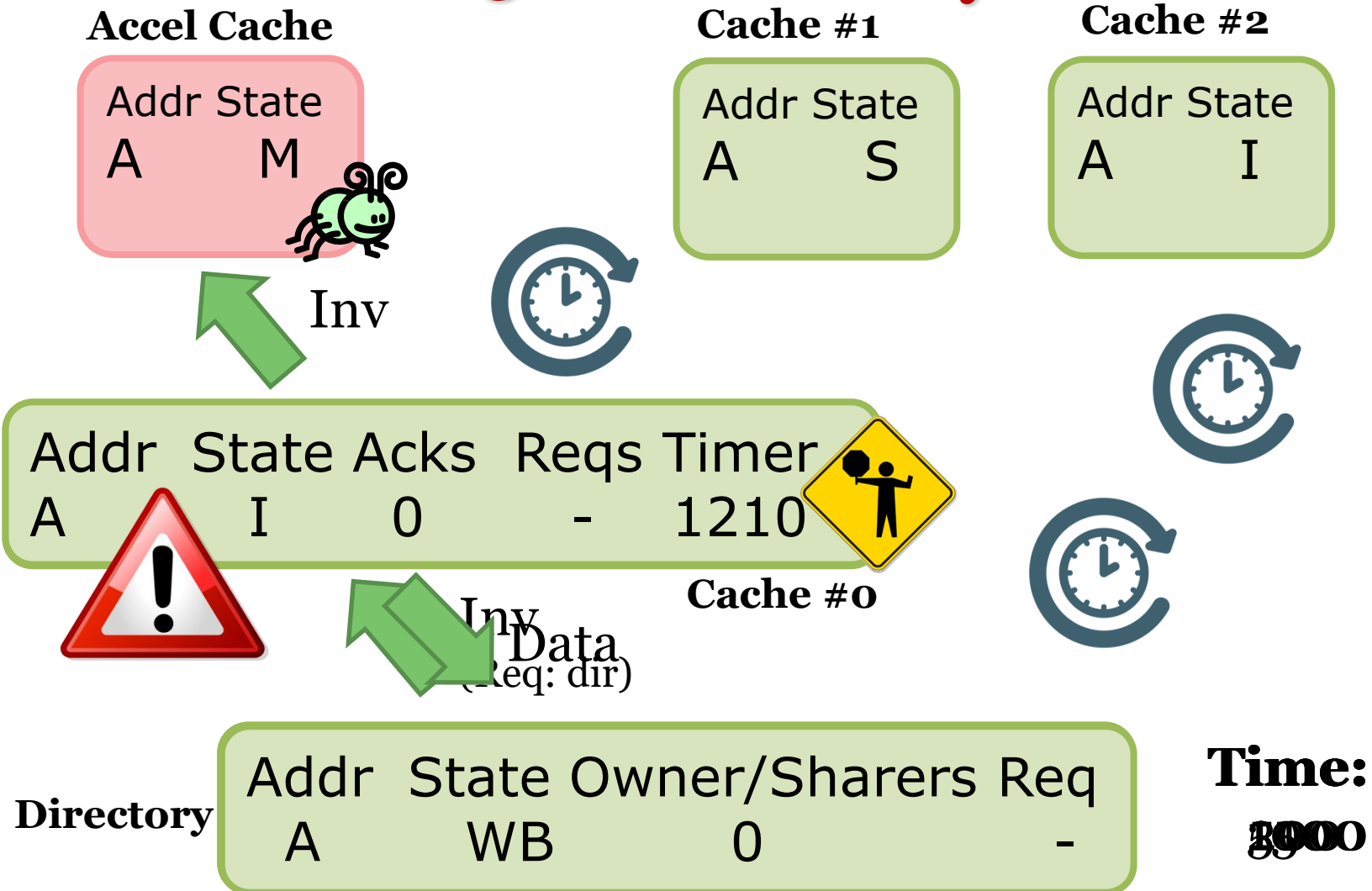


3. Host Safety





3. Host Safety





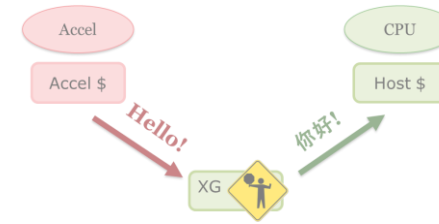
Outline



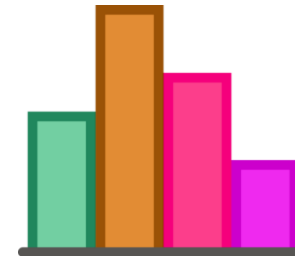
Goals



Guarantees



Design



Evaluation



Evaluation

- I. Does it provide coherence to correct accelerator?
- II. Does it provide safety to host?
- III. Does it allow high performance?



I. Correctness Testing

- Are **coherence invariants** are maintained when accelerator is acting **correctly**?
- How? **Random** tester
 - Store-Load pairs to random addresses
 - Check integrity of data
- Ran for **160 billion** load/store pairs
- Local coverage: **100% states, 100% events, > 99% transitions**



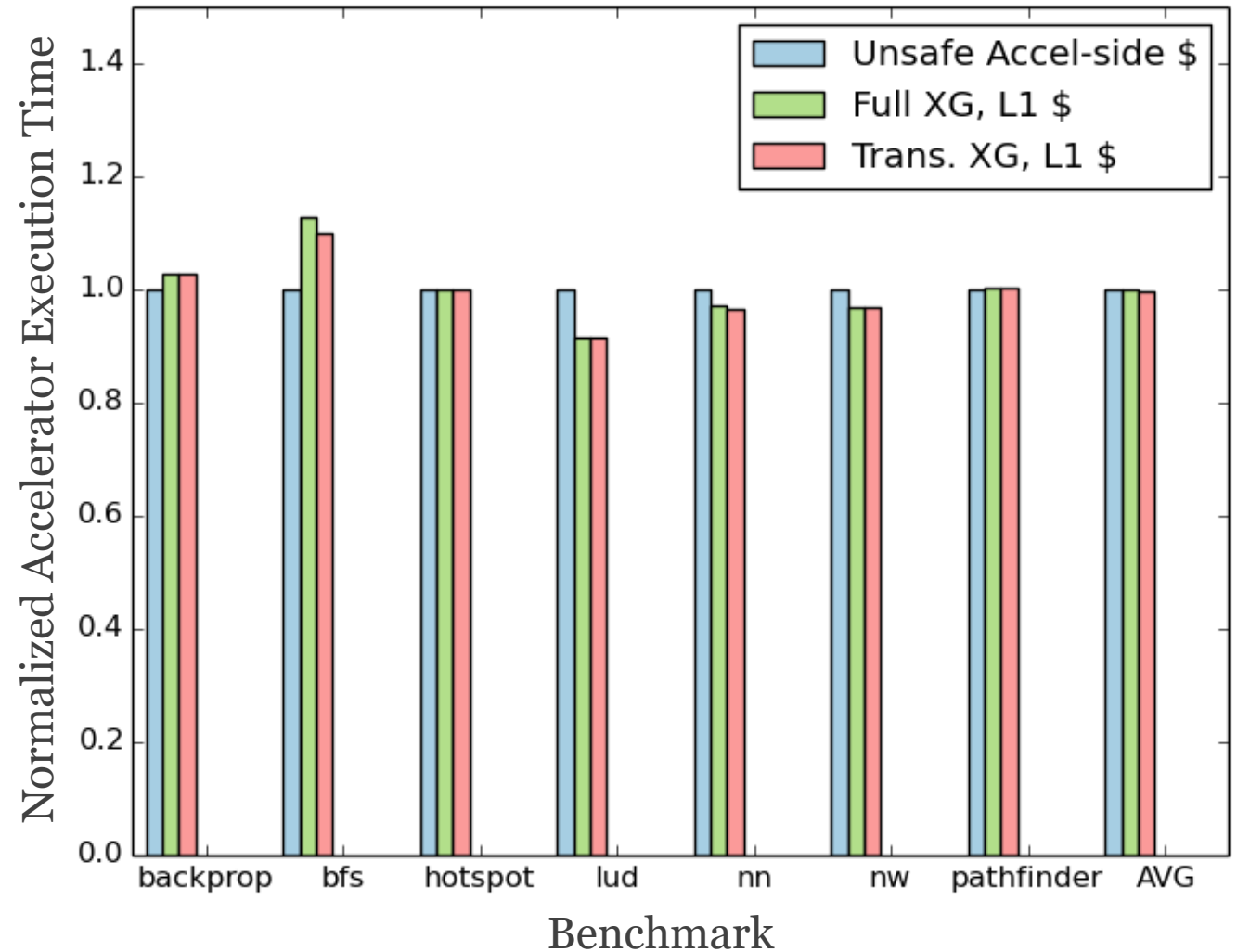
II. Fuzz Testing

- Is **host safety** maintained when accelerator **misbehaves**?
- How? Replace accelerator cache with **evil controller**
 - Generates random coherence messages to random addresses
 - Desired outcome: No deadlocks / crashes
- Ran for **7 billion** load/store pairs
- Local Coverage: **100% states, 100% events, > 99% transitions**



III. Performance Testing

- gem5-gpu
- Rodinia workloads
- MESI Inclusive host protocol





Crossing Guard Summary

- Provides **simple, standardized interface** to ease accelerator development
- **Correctness** when accelerator is correct
- **Host safety** when accelerator is incorrect
- **Low performance overhead**

Questions?



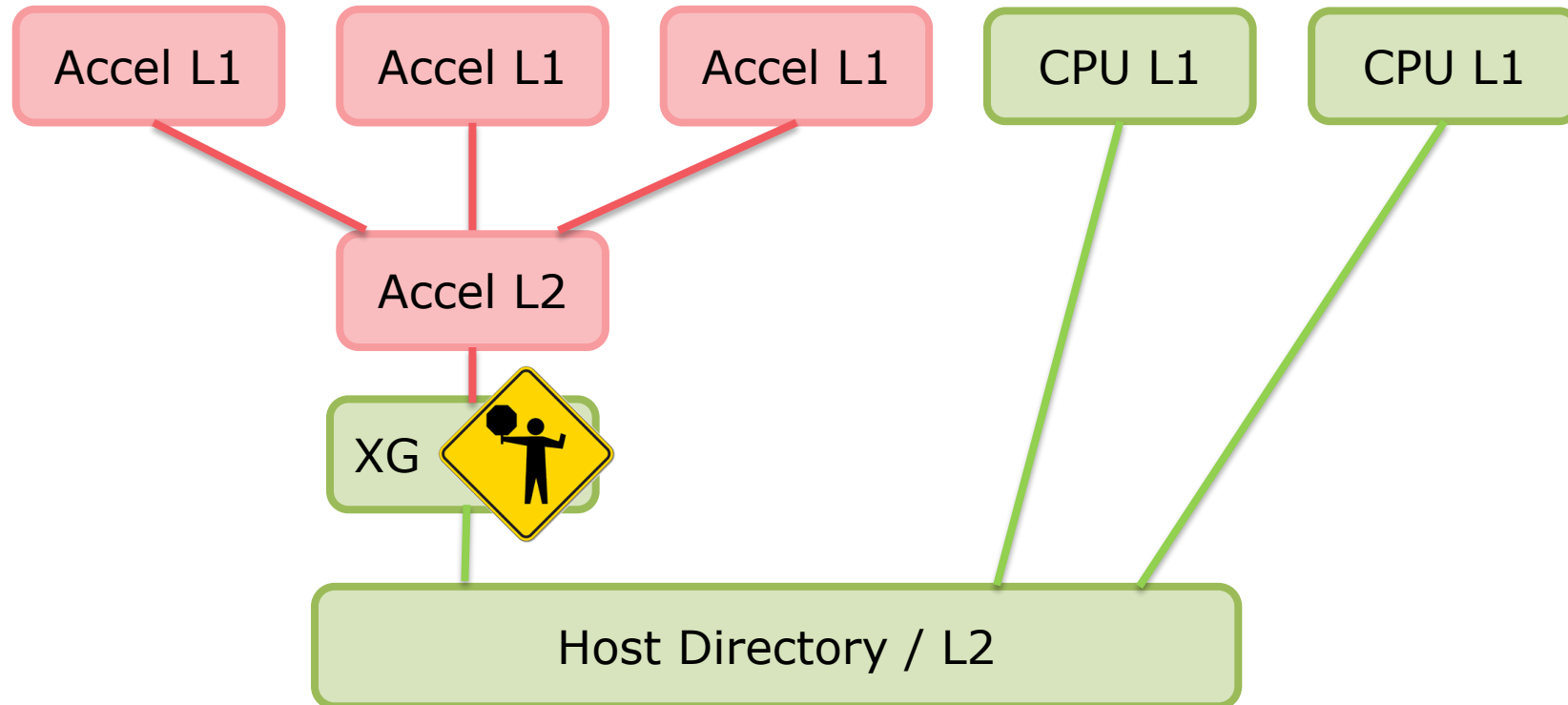
Backup Follows





Two-Level Accelerator Protocol (1)

Private L1s + Shared L2 at Accelerator





Two-Level Accelerator Protocol (2)

L1 Controller (M state contains dirty/clean bit)

	<u>Load</u>	<u>Store</u>	<u>Replacement</u>	<u>Invalidate</u>	<u>DataM</u>	<u>DataS</u>	<u>Writeback Ack</u>
<u>M</u>	<u>h q</u>	<u>h h q</u>	<u>j c v l / M I</u>	<u>d l m / I</u>			
<u>S</u>	<u>h q</u>	<u>j b q / I M</u>	<u>l / I</u>	<u>f l m / I</u>			
<u>I</u>	<u>i j a q / I S</u>	<u>i j b q / I M</u>		<u>f m</u>			
<u>I S</u>	<u>z</u>	<u>z</u>	<u>z</u>	<u>f m / I S I</u>		<u>w u k x x l h n / S</u>	
<u>I M</u>	<u>z</u>	<u>z</u>	<u>z</u>	<u>f m</u>	<u>u k x x s h n / M</u>		
<u>M I</u>	<u>z</u>	<u>z</u>	<u>z</u>	<u>f m</u>			<u>k n / I</u>
<u>I S I</u>	<u>z</u>	<u>z</u>	<u>z</u>	<u>f m</u>		<u>w u x x l h k l n / I</u>	



Two-Level Accelerator Protocol (3)

L2 Controller (Coordinates Sharing among Accelerator L1s)

	getM	getS	putM	InvAck	Writeback	Inv	DataM	DataS	WBAck	All Acks	L2 Replacement	L2 Replacement Clean
I	a r k qreq1 / IM	b r m qreq1 / IS				d qreqBC						
S	t k i p rrr qreq1 / SM	m fr rrr uu h qreq1				t w i s qreqBC / SI						t w i s / SR
MO	e r g rrr uu h qreq1 / M	m fr rrr uu h qreq1 / MOS				t w c s u qreqBC / I					t w c w m s / MRI	
M	h k qreq1 / MM	h m qreq1 / MMOS	n l p y o j qput1 / MO			t w h s qreqBC / MI					t w h s / MR	
MOS	t w k i p rrr qreq1 / MOSM	m fr uu h rrr qreq1				t w i s qreqBC / MOSI					t w i s / MOSR	t w i s / ER
IM	zz 1r	zz 1r				d qreqBC	n g rrr uu m qrspBC / M				z	z
IS	zz 1r	m qreq1				d qreqBC	n f uu m rrr qrspBC / MOS	n f uu m rrr qrspBC / S			z	z
SM	zz 1r	zz 1r		o p qrsp1		zz BCcr				y u a qt / IM	z	z
SI	z	z		o p qrsp1						d u qt / I		
SR	z	z		o p qrsp1		qreqBC / SI				c w s qt / SRI		
SRI	z	z				d qreqBC			u qrspBC / I			
MR	z	z	n l t p j qput1 / MRI	h r qrsp1	n t y o c w m qrsp1 / MRI	qreqBC / MI						
MRI	z	z				d qreqBC			u qrspBC / I			
MI	z	z	n l t p j qput1 / Mli	h r qrsp1	n t y o c u qrsp1 / I							
MOSI	z	z		o p qrsp1						c u qt / I		
MOSR	z	z		o p qrsp1		qreqBC / MOSI				c w m qt / MRI		
ER	z	z		o p qrsp1		qreqBC / MOSI				c w e qt / MRI		
MM	zz 1r	zz 1r	n l p j qput1 / MMi	h r qrsp1	n w g uu m qrsp1 / M	zz BCcr					z	z
MMOS	zz 1r	zz 1r	n l p j qput1 / MMOSi	h r qrsp1	n w f uu m y o qrsp1 / MOS	zz BCcr					z	z
MOSM	zz 1r	zz 1r		o p qrsp1		zz BCcr				g rrr y u uu m qt / M	z	z
Mli	z	z		y o c u qrsp1 / I								
MRI	z	z		y o c w m qrsp1 / MRI		qreqBC / Mli						
MMi	zz 1r	zz 1r		g uu m qrsp1 / M		zz BCcr					z	z
MMOSi	zz 1r	zz 1r		f uu m y o qrsp1 / MOS		zz BCcr					z	z



Crossing Guard Invariants

Crossing Guard Guarantees to Host:

1. Accelerator **requests** must be correct
 - a) Consistent with block **stable** state at accelerator
 - b) Consistent with block **transient** state at accelerator
2. Accelerator **responses** must be correct
 - a) Consistent with block **stable** state at accelerator
 - b) Consistent with block **transient** state at accelerator
 - c) Received within a reasonable **time**

(+ Border Control Protections!)



Crossing Guard Variants

- Full State Crossing Guard
 - Inclusive directory of accelerator state
 - + Places few restrictions on host protocol
 - + Can hide all errors
 - - Requires tag + metadata storage for all blocks
- Transactional Crossing Guard
 - Stores only data for in-flight transactions
 - + Small storage
 - + Provides most safety properties
 - - Requires some host tolerance



Single-Level Cache

States	Accelerator Events			XG Requests Invalidate	XG Responses			
	Load	Store	Replacement		DataM	DataE	DataS	WB Ack
M	hit	hit	issue PutM / B	send Dirty WB / I	-	-	-	-
E	hit	hit / M	issue PutE / B	send Clean WB / I	-	-	-	-
S	hit	issue GetM / B	issue PutS / B	send InvAck / I	-	-	-	-
I	issue GetS / B	issue GetM / B	-	send InvAck	-	-	-	-
B	<i>stall</i>	<i>stall</i>	<i>stall</i>	send InvAck	/ M	/ E	/ S	/ I



Simulation Parameters

CPU		
CPU Cores		1
CPU Frequency		3 GHz
Host Caches		
	Hammer-like	MESI Inclusive
L1I & L1D	32kB each	32kB each
L2	128kB private	512kB shared w/ GPGPU

Table 4: CPU simulation configuration details.

GPGPU			
Cores			4
GPU Frequency			700 MHz
GPGPU Caches (Hammer-like)			
	Accel-side	Host-side / 1-level	2-level
L1	16kB I and D	160kB	32kB
L2	128kB private	-	512kB shared
GPGPU Caches (MESI Inclusive)			
	Accel-side	Host-side / 1-level	2-Level
L1	32kB I and D	64kB	16kB
L2	-	-	192kB shared
Cache-to-Cache Latency			
Accelerator L1 to L2			10 cycles
Accelerator L2 to XG			200 cycles
XG to Directory/Shared L2			10 cycles
Accelerator to Host-side Cache			210 cycles

Table 5: GPGPU simulation configuration details.



Time Spent Simulating (Random)

Configuration	Time
XG Full + Hammer + 1 Level	5.28 years
XG Full + Hamer + 2 Level	2.51 years
XG Full + MESI Inc + 1 Level	133 days
XG Full + MESI Inc + 2 Level	223 days
XG Trans. + Hammer + 1 Level	3.17 years
XG Trans. + Hammer + 2 Level	1.38 years
XG Trans + Inc + 1 Level	90 days
XG Trans + Inc + 2 Level	103 days
TOTAL	13.9 years



Full Coverage %s (Random)

Full State XG	Single-level	Two-level
Hammer-like	99	99.8
MESI Inclusive	100	99.4
Transactional XG	Single-level	Two-level
Hammer-like	99.3	99.5
MESI Inclusive	100	99.7



Time Spent Simulating (Fuzz)

Configuration	Time
XG Full + Hammer-like	1.62 years
XG Full + MESI Inclusive	287 days
XG Transactional + Hammer-like	5.3 years
XG Transactional + MESI Inclusive	41 days
Total	7.82 years



Full Coverage %s (Fuzz)

Full State Crossing Guard	Fuzz Tester
Hammer-like	99.3
MESI Inclusive	99.7
Transactional Crossing Guard	Fuzz Tester
Hammer-like	99.7
MESI Inclusive	100



PutS Accelerator Messages

- Why?
 - Some host protocols use them
 - Simplify management of Full State Crossing Guard
 - Cannot implement Transactional Crossing Guard + host protocol with PutS without them
- Bandwidth Impact
 - Carry no data
 - Only between accelerator cache → Crossing Guard, not host system
 - ~1-4% of that bandwidth in experiments.
 - Could be reduced by setting a flag at Crossing Guard.

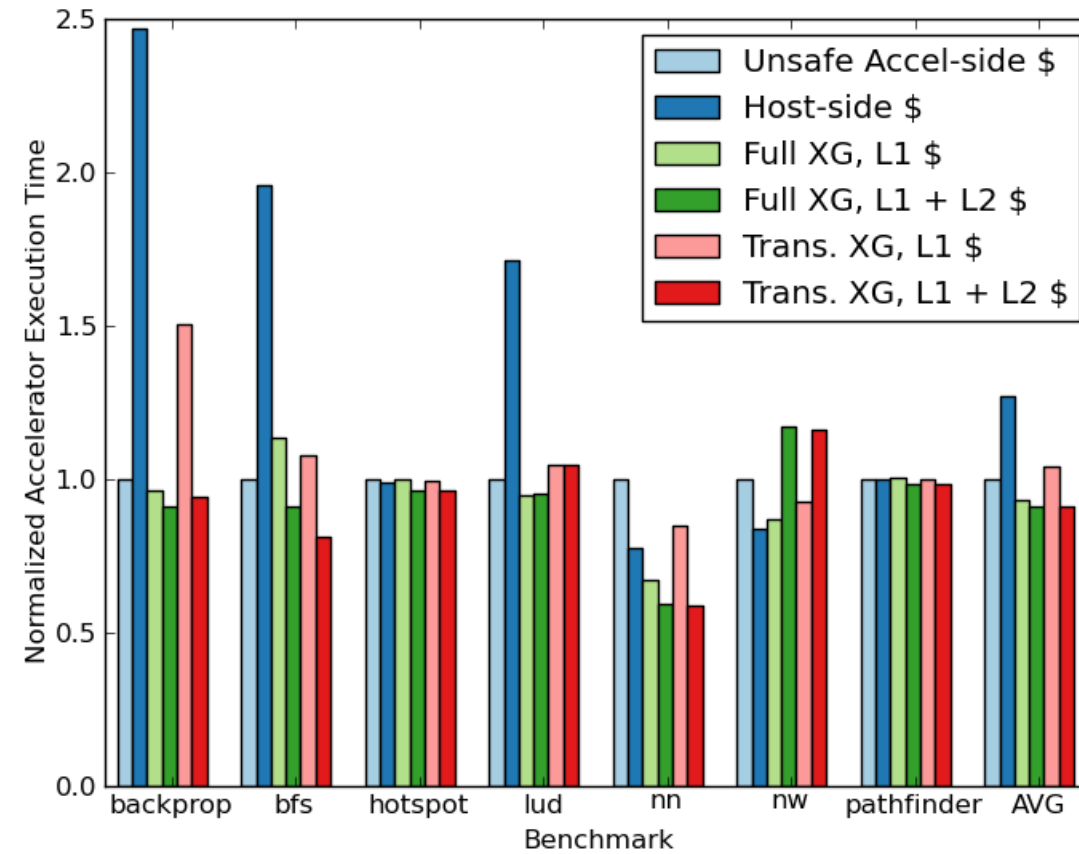


Why not Model Checking?

- Model checking is useful! Industrial implementation of Crossing Guard would use.
- Academic tools have limitations ☹
 - Benefit from symmetry, but Crossing Guard system asymmetric
 - May only work with one block in system
 - Substantial implementation overhead
- This work was a proof of concept
 - Random / Fuzz testing not perfect, but results suggestive.
 - Even models can have mistakes!

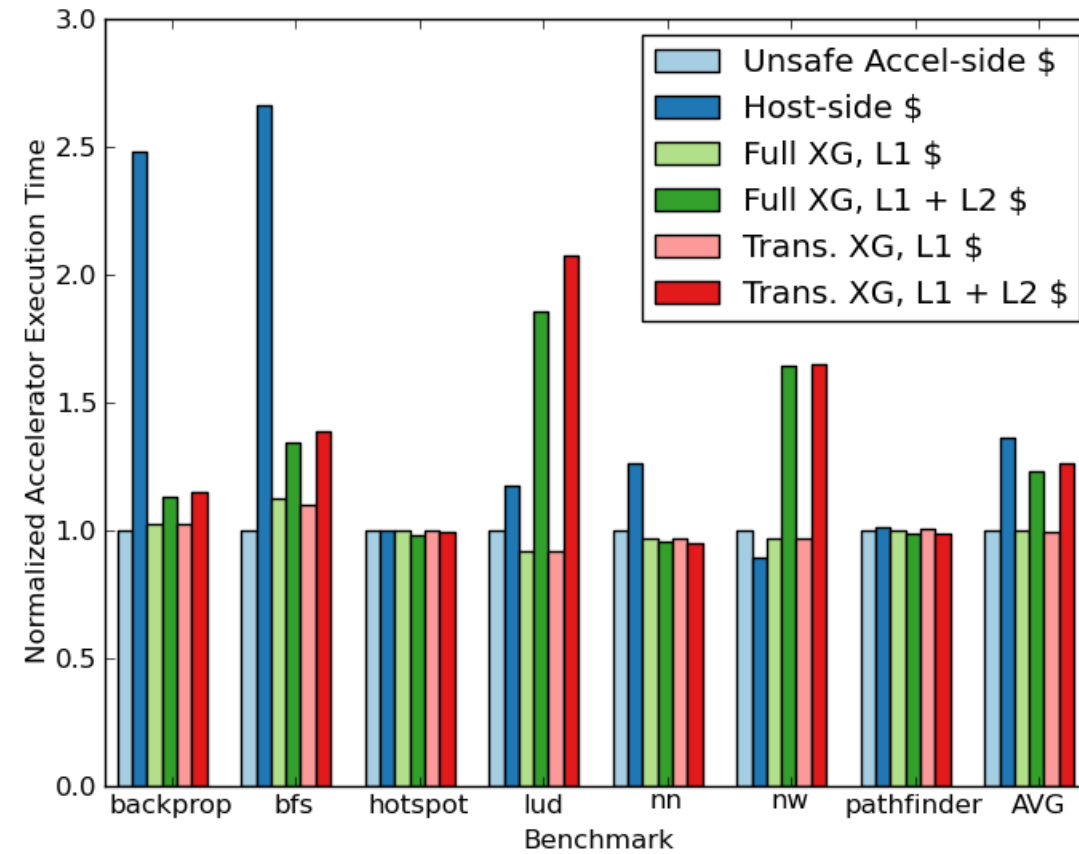


Performance: Hammer-like



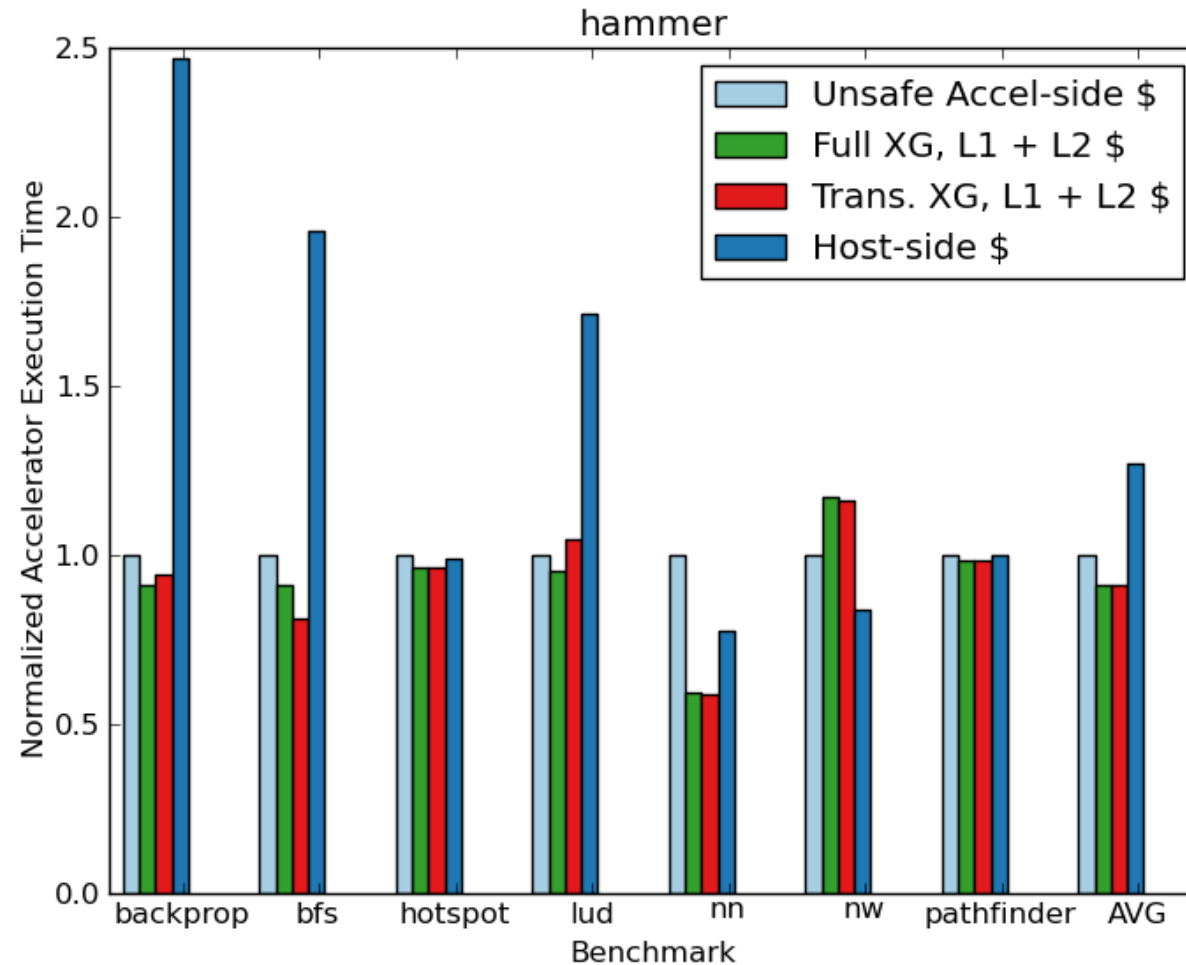


Performance: MESI Inclusive



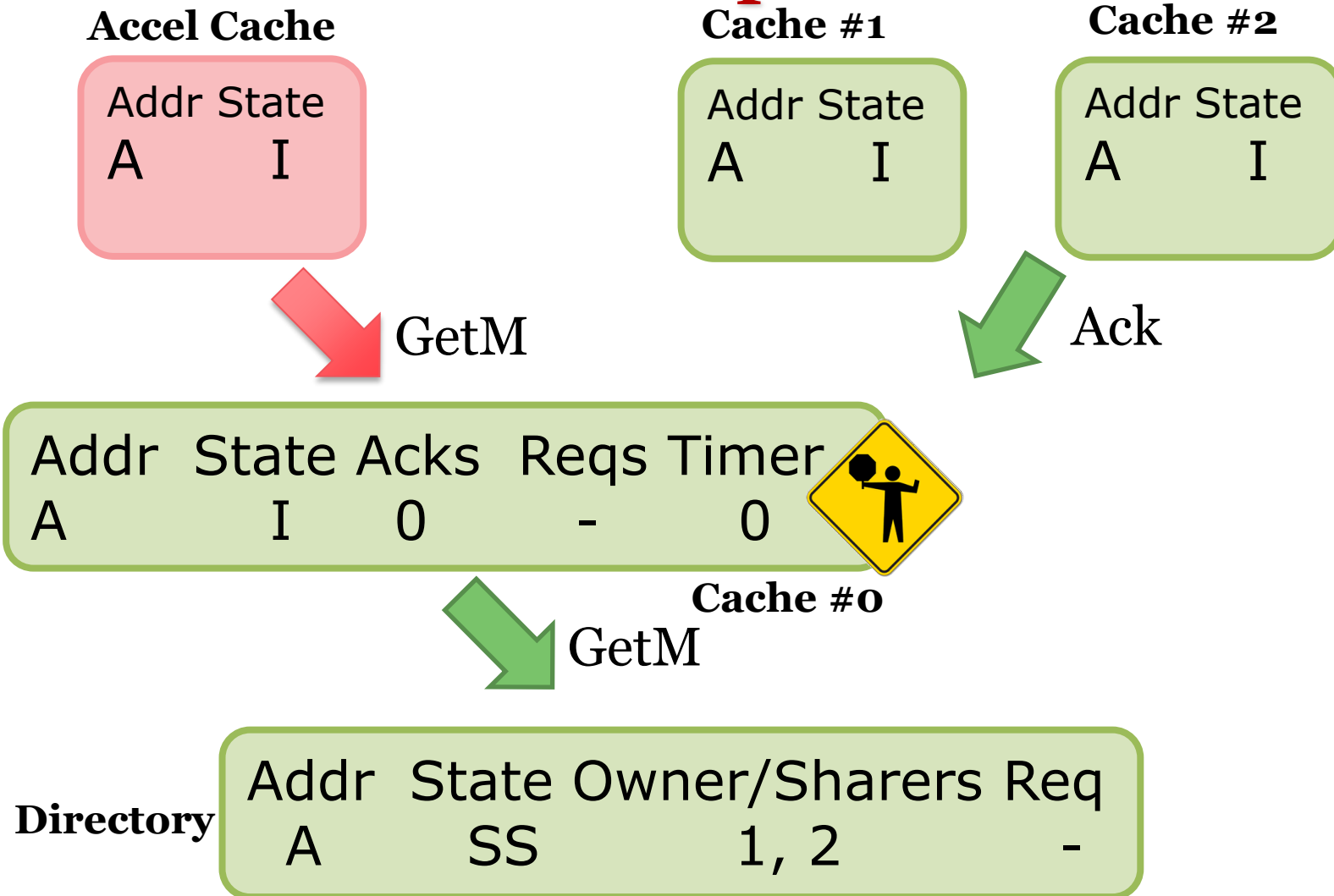


Performance (Hammer-like)





Template

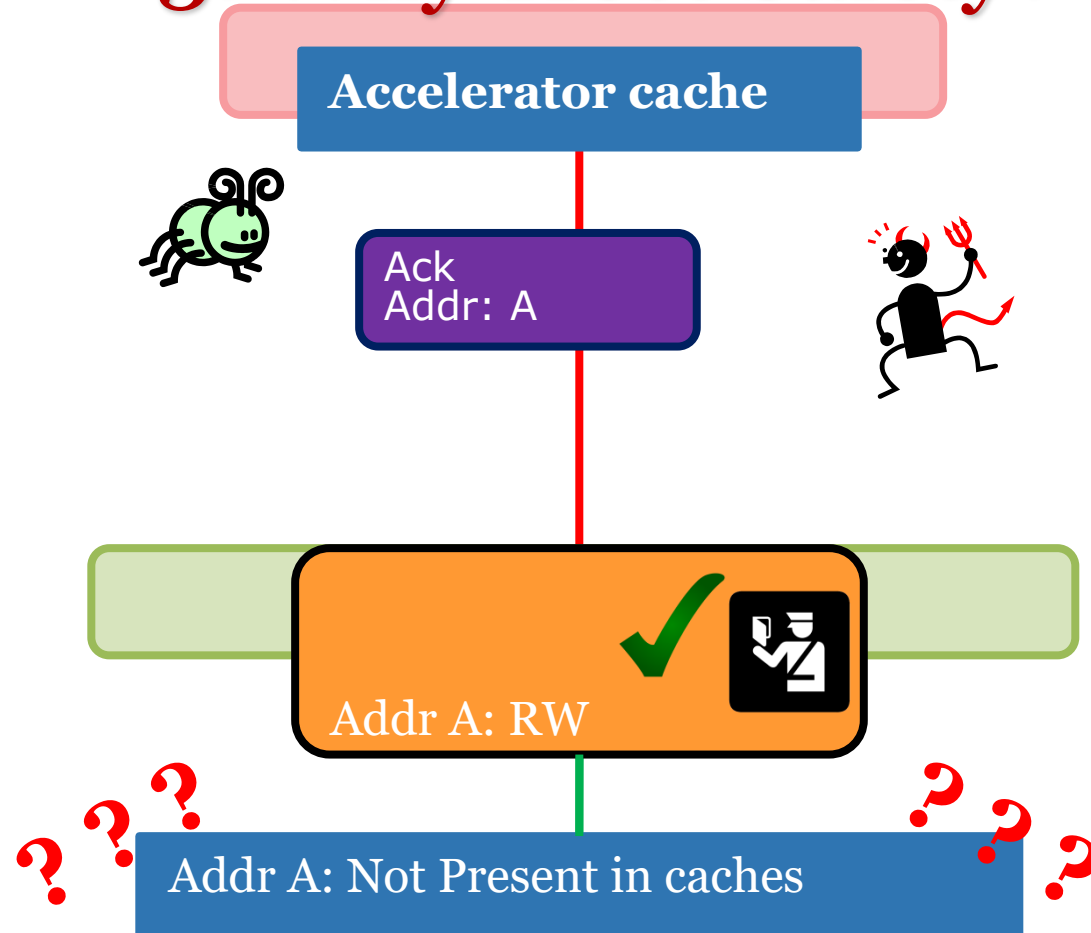


Old Slides



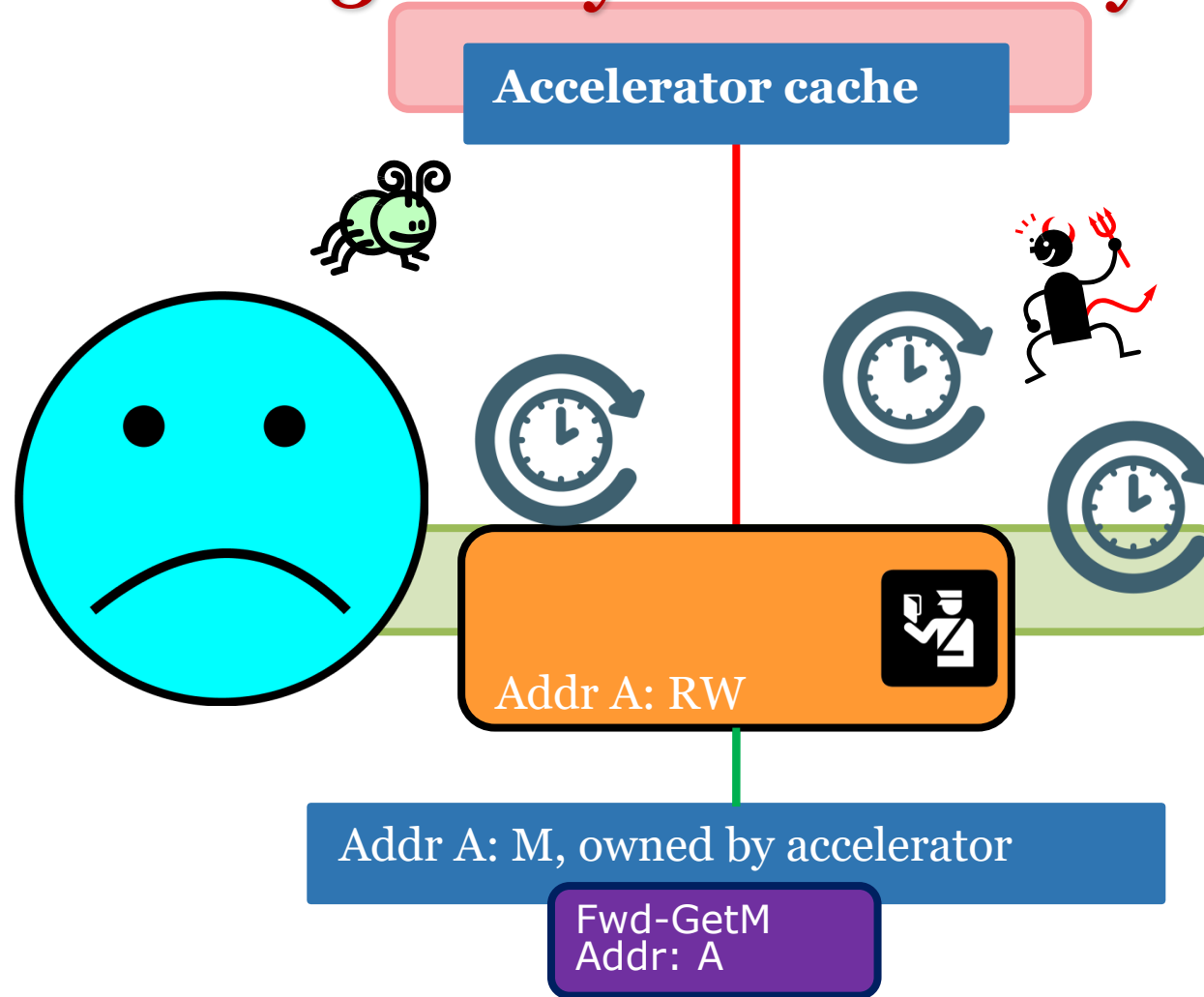


3. Why Host Safety?



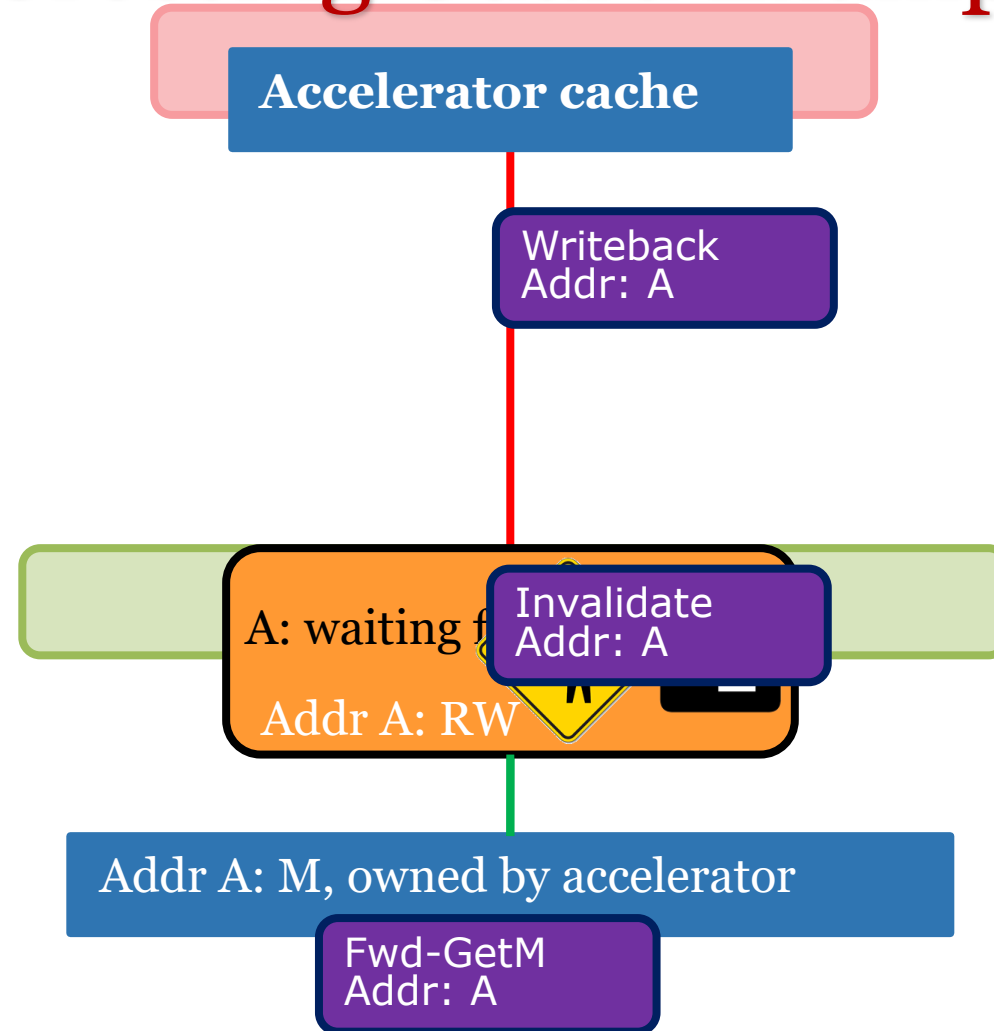


3. Why Host Safety?





Crossing Guard Example





Crossing Guard Example

