

# Formal Verification and Synthesis for Quality-of-Service in On-Chip Networks

**August 27, 2013**

**Daniel E. Holcomb**

**UC Berkeley**

**Prof. Sanjit A. Seshia**



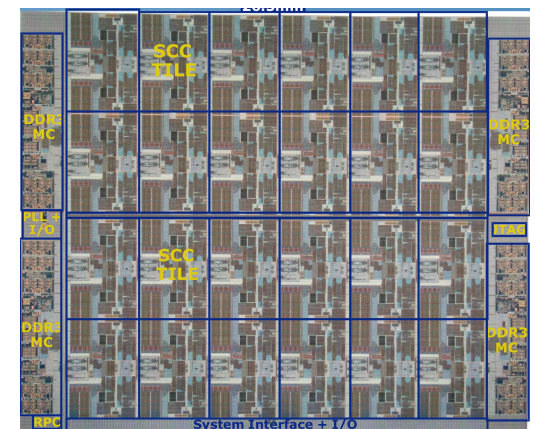
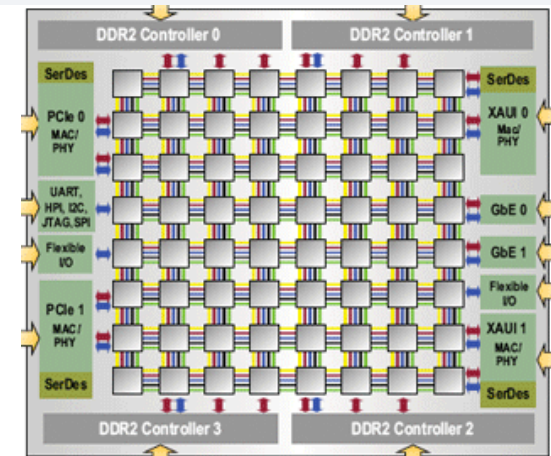
# Acknowledgments

- **Bryan Brady**
- **Michael Kishinevsky**
- **Alexander Gotmanov**
- **Yuriy Viktorov**
- **Satrajit Chatterjee**
- **Bob Brayton's group**
- **Jiang and Alan**
- **etc**

**Gigascale Systems  
Research Center, one of  
six research centers  
funded under the Focus  
Center Research  
Program (FCRP), a  
Semiconductor Research  
Corporation entity.**

# Motivation

- On-chip networks due to scaling
  - Intel Teraflops Research chip
    - 80 simple cores
  - Tiler TILE-Gx100
    - 100 general purpose cores
  - Intel Single-Chip cloud computer
    - 48 IA cores
- Significant Cost/Performance considerations
- Complex to reason about
  - Arbitration -- including unfair policies
  - Reservations
- State of the art for Quality of Service (QoS)
  - Reasoning through extensive simulation
  - Analytical techniques in some simple cases
  - Can we do better with formal methods?



# Thesis Statement

**Leverage model checking for solving NoC QoS latency problems. Address capacity limitations by extending well-known formal techniques including abstraction and compositional reasoning into the NoC domain**

- This talk concerns 3 specific QoS contributions
    - Workload abstraction of traffic models
    - Latency proofs by property strengthening
    - Optimal buffer sizing for QoS
- } Compositional Latency Verification



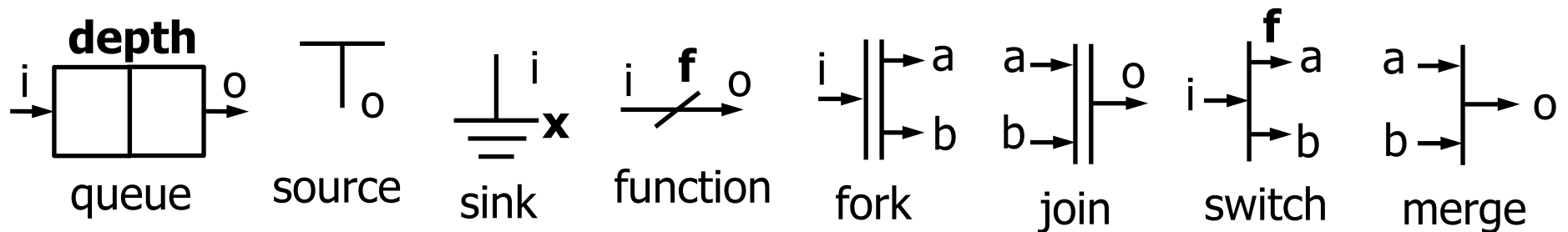
# Outline

- **Background**
  - **xMAS: Formal model of NoC**
  - **Verification technology**
- Compositional latency verification
  - Abstraction and traffic modeling
  - Inductive proof by property strengthening
- Buffer sizing using counterexamples

# xMAS Modeling Language

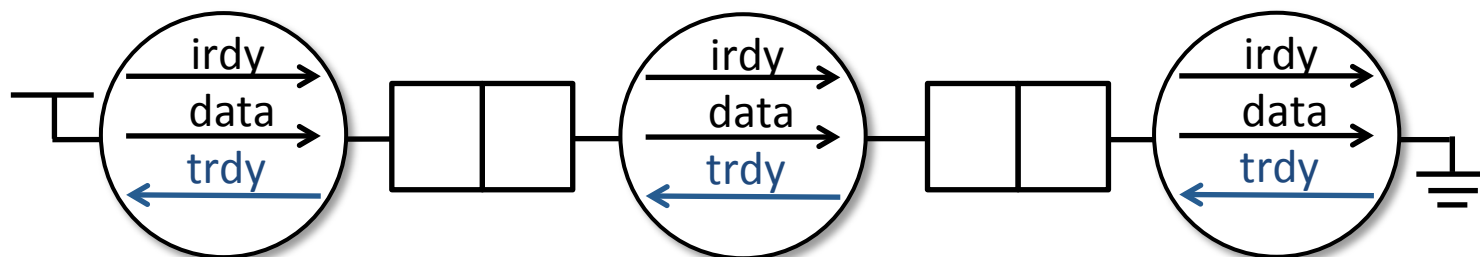
## Executable Microarchitectural Specifications (xMAS)

- Represent communication fabrics as compositions of kernel primitives [Chatterjee et al., HLDVT'10]



- Primitives connected by channels

transfer = irdy & trdy



# Why use xMAS Formalism?

- Effective abstraction for real designs
  - Expressive enough to model interesting behaviors
  - Hides messy details of arbitrary RTL
- Finite set of primitives allows reuse of reasoning
- Convenient entry point for verification tools
  - RTL or UCLID modeling language
  - Can augment xMAS with arbitrary sequential logic as needed

# Formal Verification / Model Checking

- Model checking
  - Given a finite state model with initial state(s), inputs and transition relation
  - Formally check whether any sequence of inputs can drive the model to a bad state
- Used when bugs are unacceptable
  - Safety-critical applications
  - VLSI designs
- Adoption limited by lack of appropriate models, and inability to scale

# Model Checking an xMAS network

- SAT-based verification using And-Inverter-Graphs (AIGs)
- State-of-the-art model checking engines in ABC [Brayton and Mishchenko, CAV'10]
  - Bounded Model Checking (BMC)
  - Induction / K-induction
  - Property directed reachability (PDR) [Een and Mishchenko, IWLS'11]
    - ABC's Implementation of IC3 [Bradley VMCAI'11]
- VeriABC converts verilog to AIG [Long et al., IWLS'11]

# Model Checking an xMAS network

- Satisfiability Modulo Theories (SMT) solving
  - An approach for coping with state explosion
  - Decide validity of term-level formulas
- UCLID model checker [Bryant et al., CAV'02]
  - Verifier for systems expressed in a combination of logical theories
  - Symbolic simulation
  - SMT solving as the underlying engine

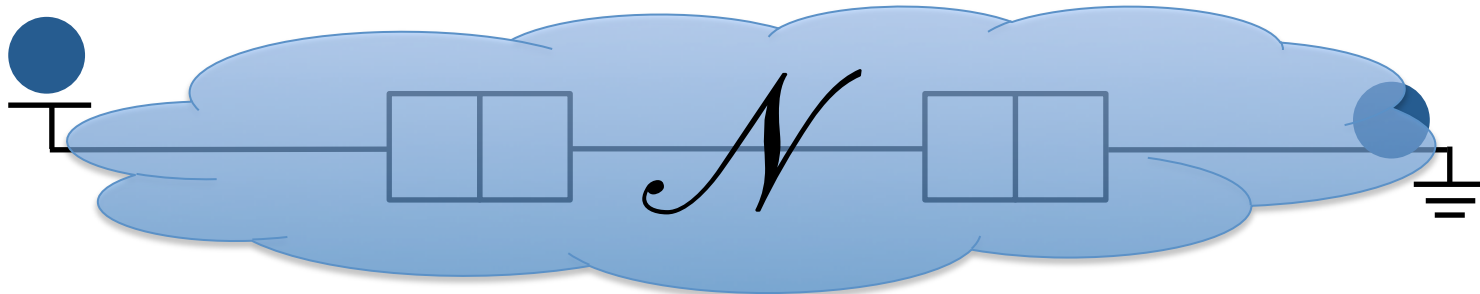
# Verifying Latency Bounds

- xMAS network  $\mathcal{N}$
- Denote latency property  $\Phi$
- Property for global latency bound  $T$

$$\Phi_T^G \approx \mathbf{G}(src \implies \mathbf{F}^{<T} sink)$$

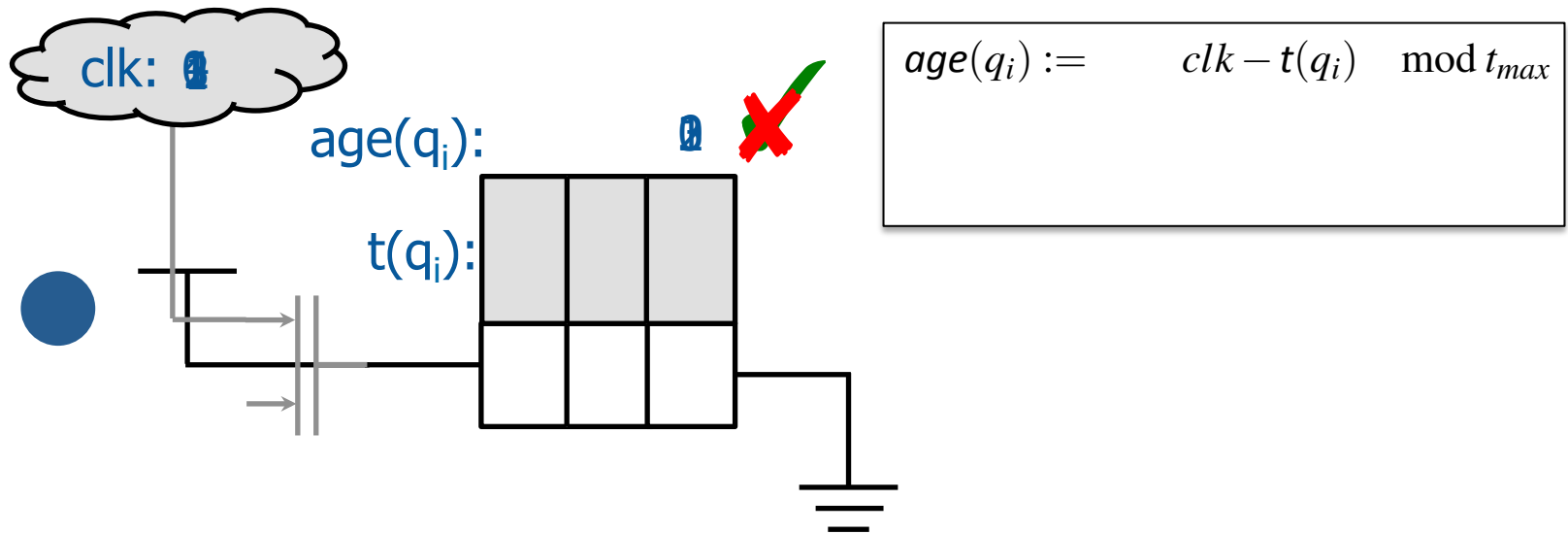
– A “prompt-LTL” property [Kupferman et al. ‘09]

- Goal is to verify  $\mathcal{N} \models \Phi_T^G$



# Checking Latency as Simple Safety Property

- Represent latency since injection as **age** of packet
  - Global clock (*n*-bit counter)
  - Injection timestamps on packets (*widen queue slots by n*-bits)





# Outline

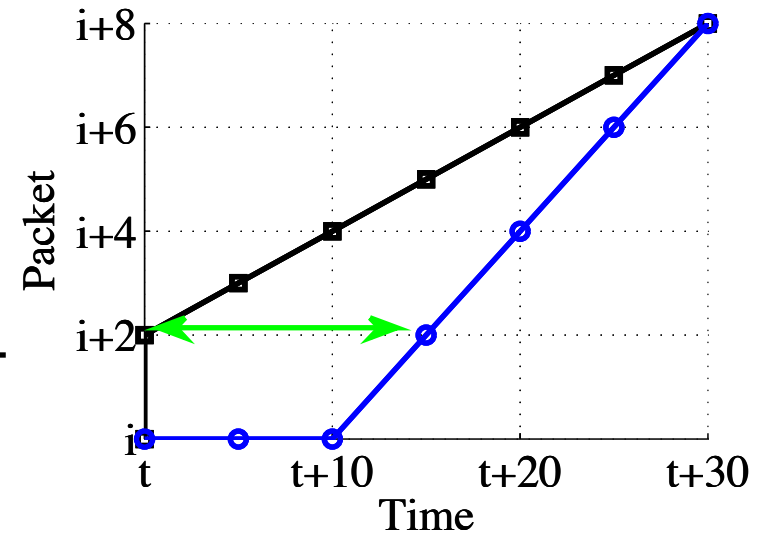
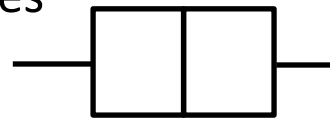
- Background
  - xMAS: Formal model of NoC
  - Verification technology
- **Compositional latency verification**
  - Abstraction and traffic modeling
  - Inductive proof by property strengthening
- Buffer sizing using counterexamples

# Compositional Latency Verification

- Size of interesting latency problems
  - 10s or 100s of cycles
  - Similar number of queues
- Without compositional reasoning, verification difficulty grows with both
  - Beyond scope of current tools
- We give two approaches to compositional latency verification of NoC

# Related Work

- **Network Calculus** [Cruz '91] for analyzing networks with traffic rates bounded by linear envelopes
  - Arrival and **service** curves
  - **Latency** bounds
  - Algebraic composition rules



- Model checking to verify of Router Latency [Krishna et al. Haifa '11]
- **Static Timing Analysis** of combinational circuits
  - Use graph as abstraction of gate-level timing

# Outline

- Background
  - xMAS: Formal model of NoC
  - Verification technology
- Compositional latency verification
  - **Abstraction and traffic modeling**
  - Inductive proof by property strengthening
- Buffer sizing using counterexamples

# Traffic Abstraction

- Decompose both model and property
- Reduce scope of verification problem by focusing on one router in isolation
- Abstract rest of network into precise environment model

Daniel Holcomb, Bryan Brady, and Sanjit A. Seshia. **Abstraction-Based Performance Analysis of NoCs**. In *Proceedings of the Design Automation Conference (DAC)*, pp. 492–497, June 2011.

<http://www.eecs.berkeley.edu/~holcomb/dac11-noc.pdf>

## Abstraction-Based Performance Analysis of NoCs

Daniel E. Holcomb  
UC Berkeley  
holcomb@eecs.berkeley.edu

Bryan A. Brady  
UC Berkeley  
brady@eecs.berkeley.edu

Sanjit A. Seshia  
UC Berkeley  
seshia@eecs.berkeley.edu

**ABSTRACT**

We present a problem decomposition and analysis of network QoS performance in NoC designs. To handle realistic performance metrics, we use a detailed model of the network, including an abstract model of the rest of the network. This model is used to generate a precise environment model for the router under analysis. This model is used to generate a precise environment model for the router under analysis. This model is used to generate a precise environment model for the router under analysis.

### 1. INTRODUCTION

Network-on-chip (NoC) is a paradigm for communication within large-scale systems. It consists of a network of routers, each with a local processor, and a set of links connecting them. The network is used to transport data between different parts of the system. The network is used to transport data between different parts of the system. The network is used to transport data between different parts of the system.

Performance metrics such as latency and throughput are critical for NoC designs. This paper presents a method for analyzing these metrics in a precise environment model. The method is used to generate a precise environment model for the router under analysis.

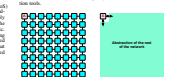


Figure 1. Challenge of performance verification. We wish to verify a bound on the latency from node A to node B.

The abstract performance approach to apply formal methods to this problem is to abstract the network. In this case, we can abstract the network into a precise environment model. This model is used to generate a precise environment model for the router under analysis. This model is used to generate a precise environment model for the router under analysis.

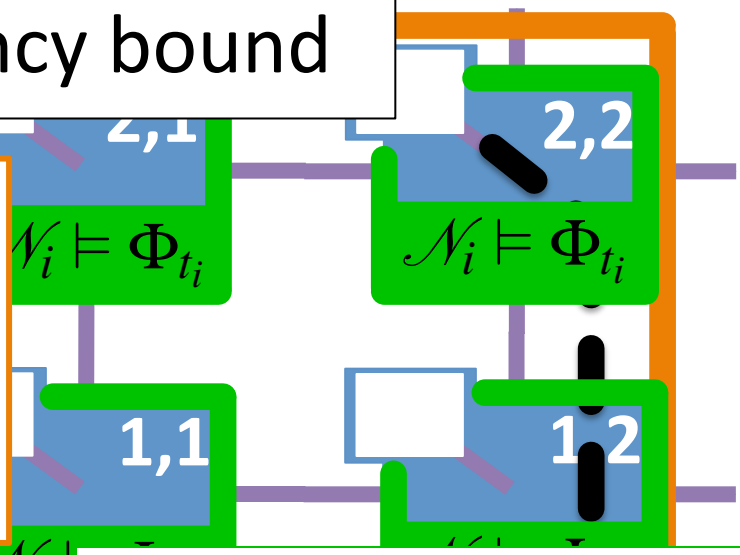
Performance metrics such as latency and throughput are critical for NoC designs. This paper presents a method for analyzing these metrics in a precise environment model. The method is used to generate a precise environment model for the router under analysis.

# Abstraction Enables Formal Approach

- Goal: prove a global latency bound

- ~~Model in xMAS~~
- ~~Find bounds by sweeping latency  $T$  and checking each~~
- Model too large!

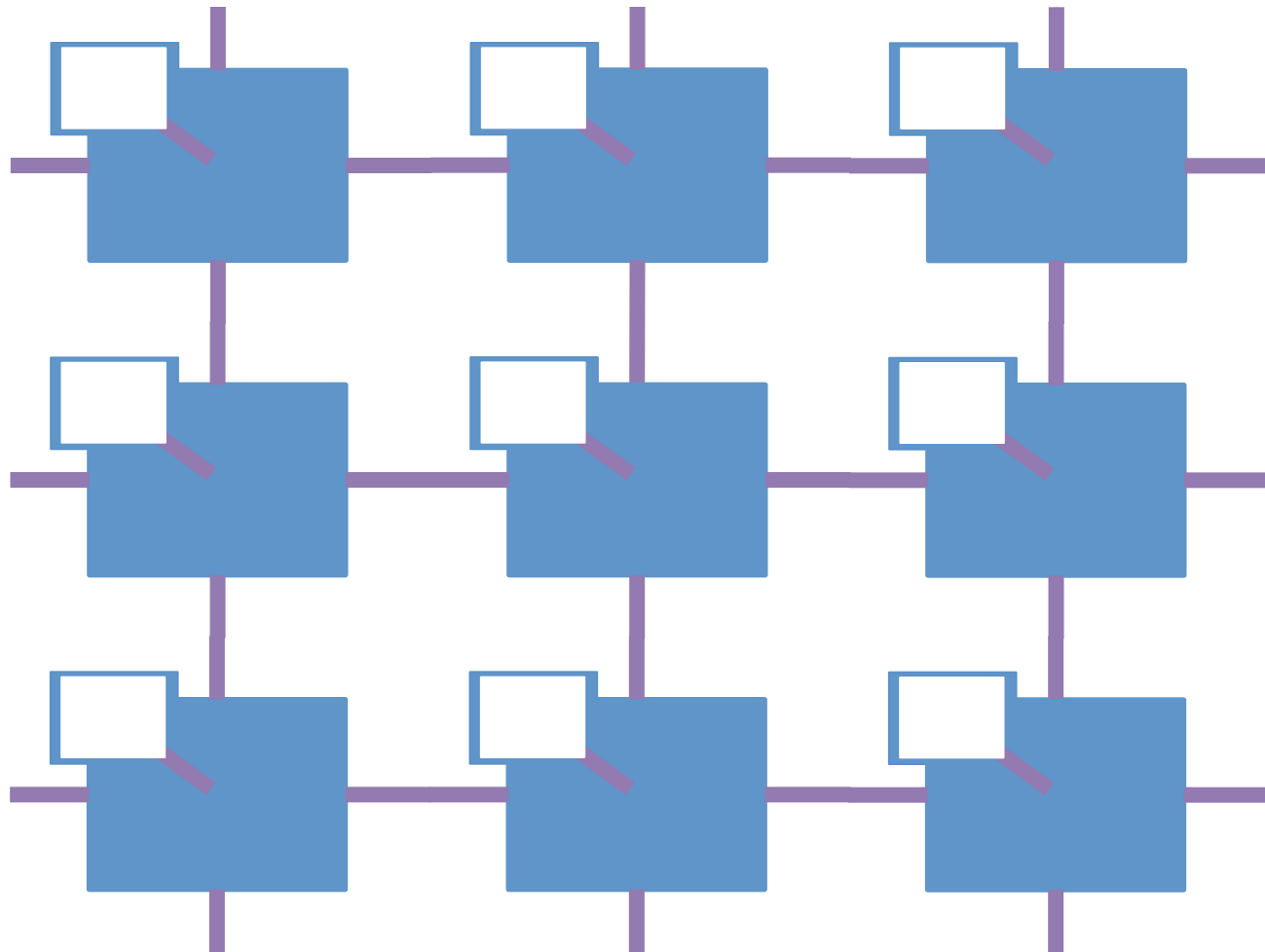
- Prove conservative latency bounds for each router
- Compose using path routing constraints



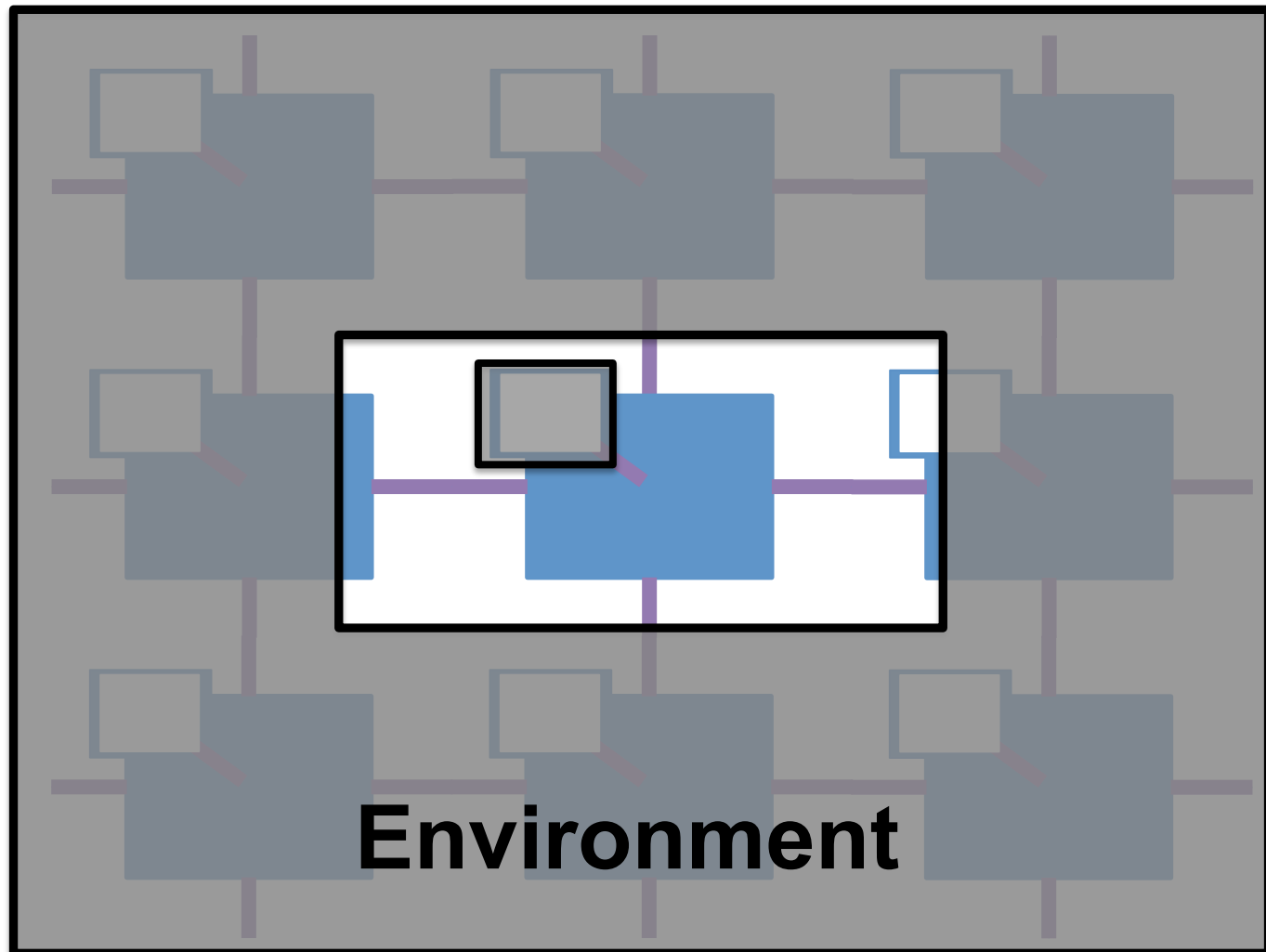
$$T := \sum_{i \in path} t_i$$

$$\left( \bigwedge_i \mathcal{N}_i \models \Phi_{t_i} \right) \implies \mathcal{N} \models \Phi_T^G$$

# Abstraction Enables Formal Approach

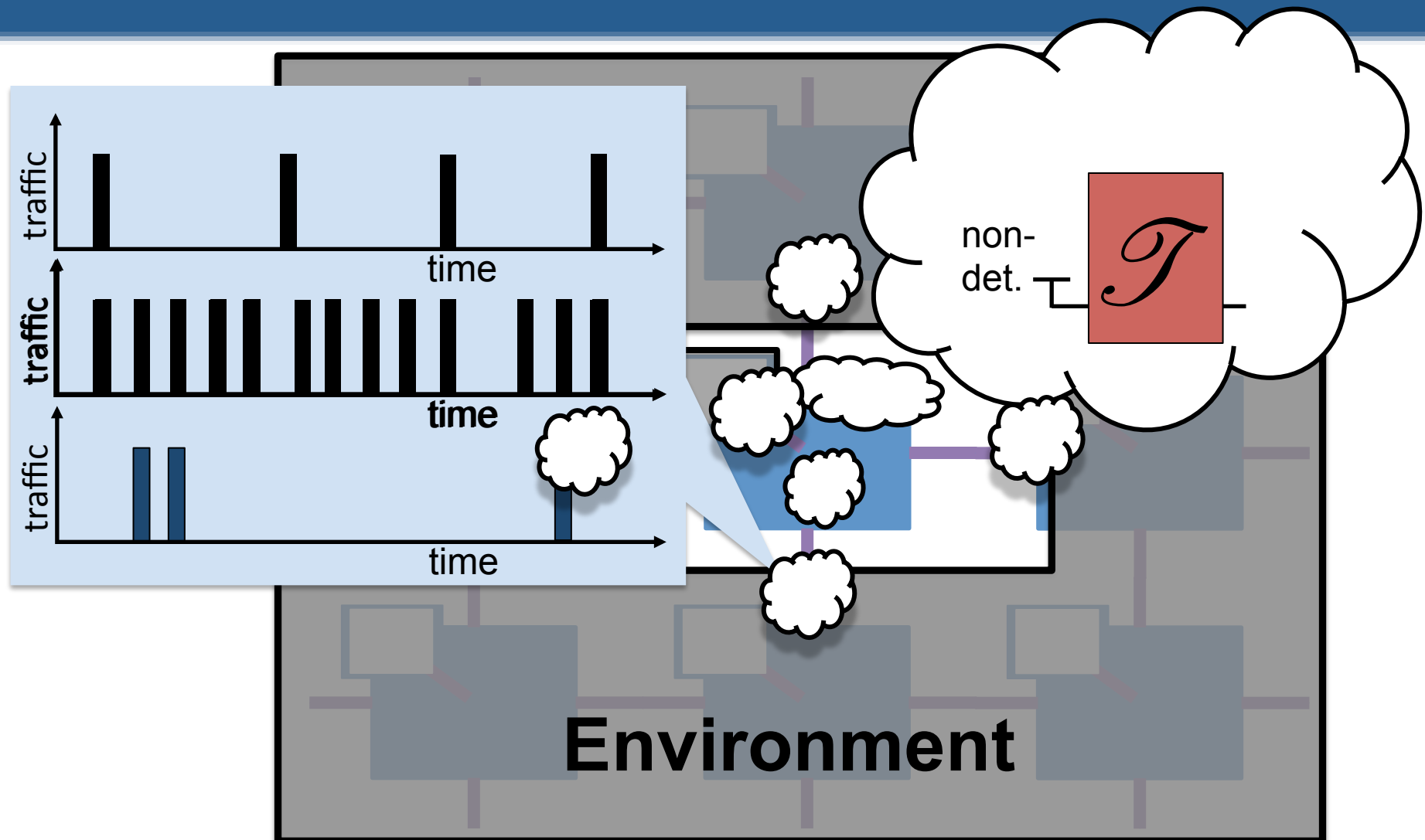


# Abstraction Enables Formal Approach



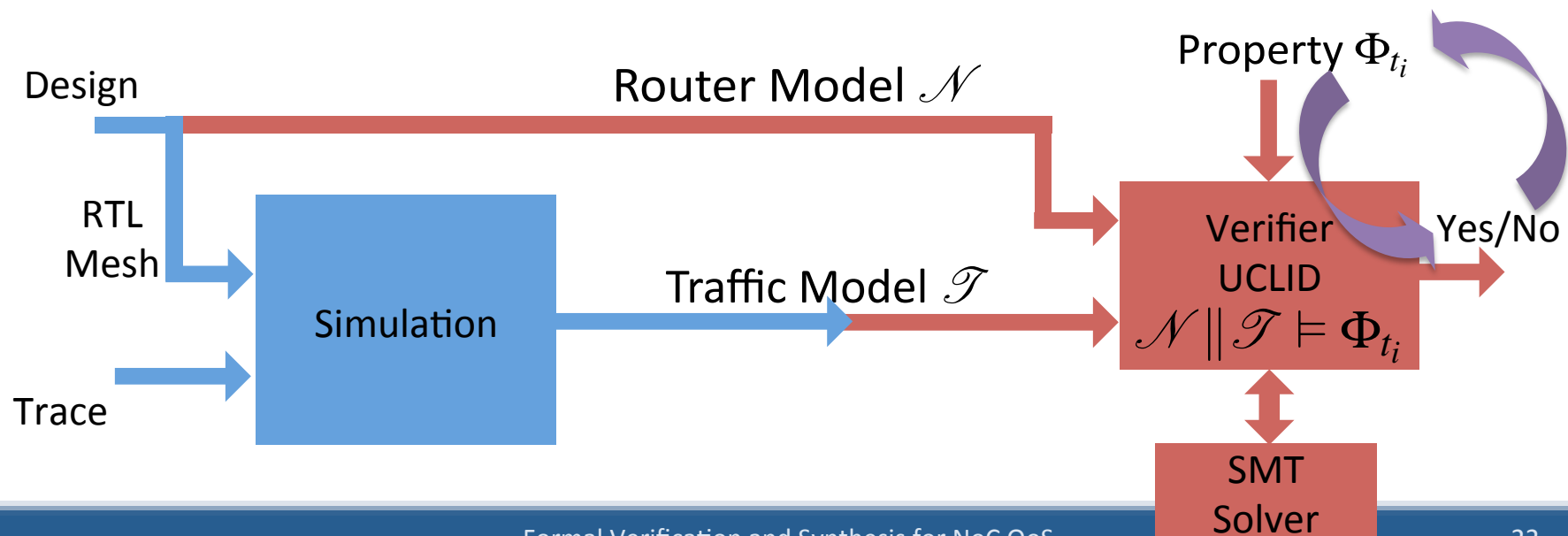


# Abstraction Enables Formal Approach



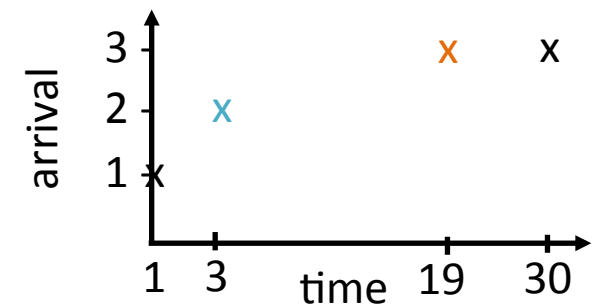
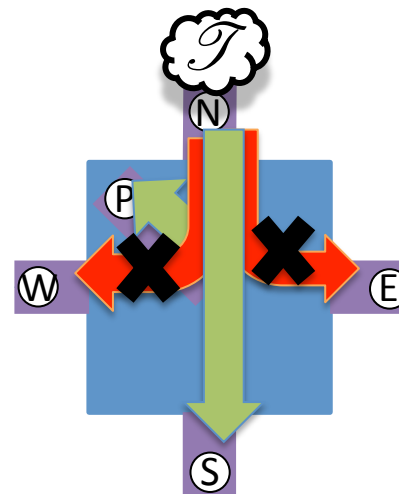
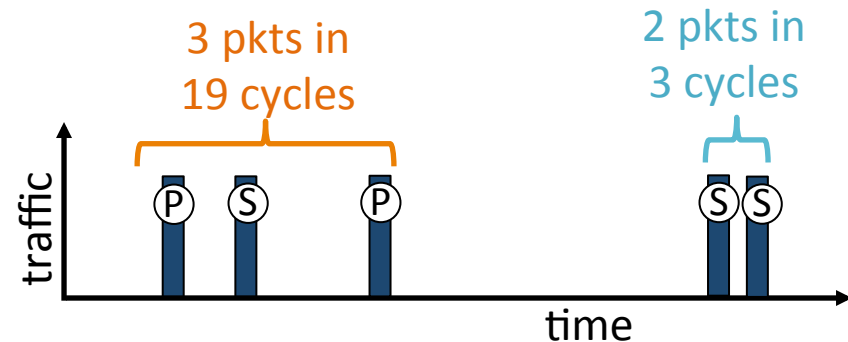
# Contributions

- Precise abstraction for formal NoC analysis
  - Extension of network calculus [Cruz, Tran. Info Theory '91]
- Inferring traffic model from simulation data
- Verifying performance using inferred traffic models



# Inferring Formal Traffic Model

- Inferred from RTL simulation
  - CMP router  
[Peh PhD Thesis 2002]
  - Random samples from PARSEC traces  
[Soteriou et al., MASCOTS 2006]
- Predicates on channel behaviors
  - Boolean constraints
  - Conjunctions of rate constraints



# Inferring Formal Traffic Model

- Inferred from RTL

simulation

- CMP

[Peh Ph

- Rand

PARS

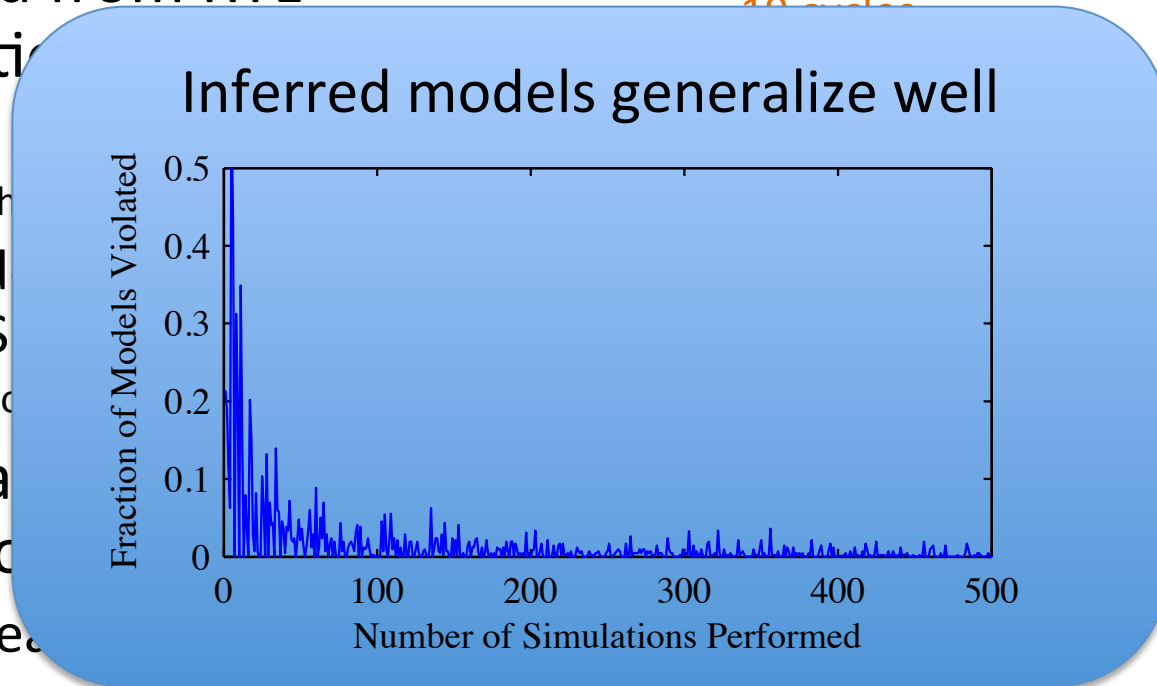
[Soteric

- Predica

behavior

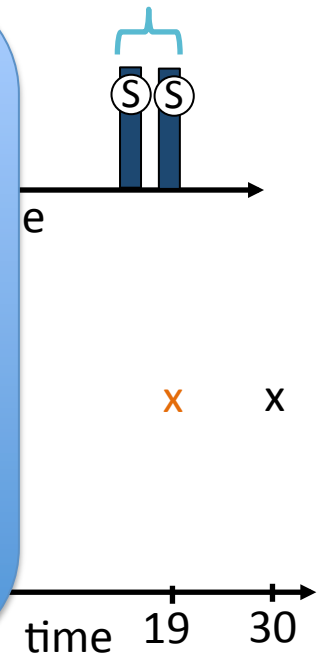
- Boolean

- Conjunctions of rate constraints



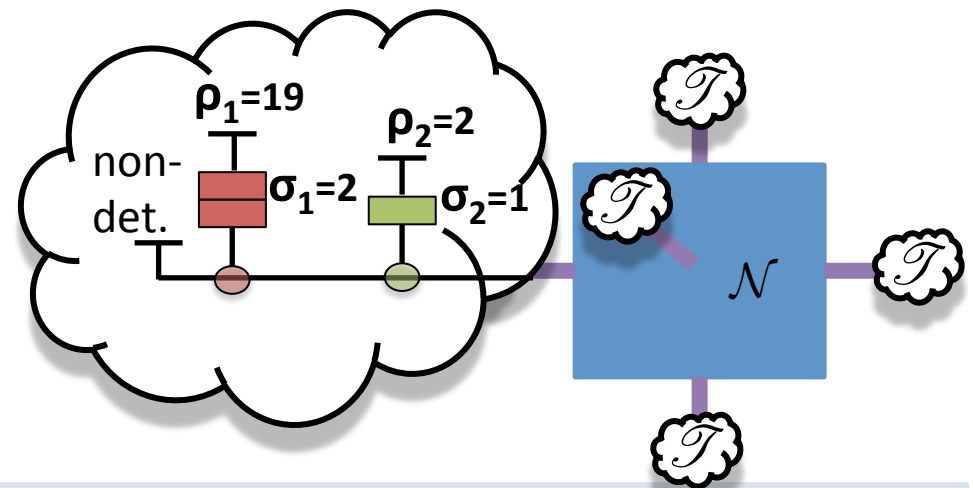
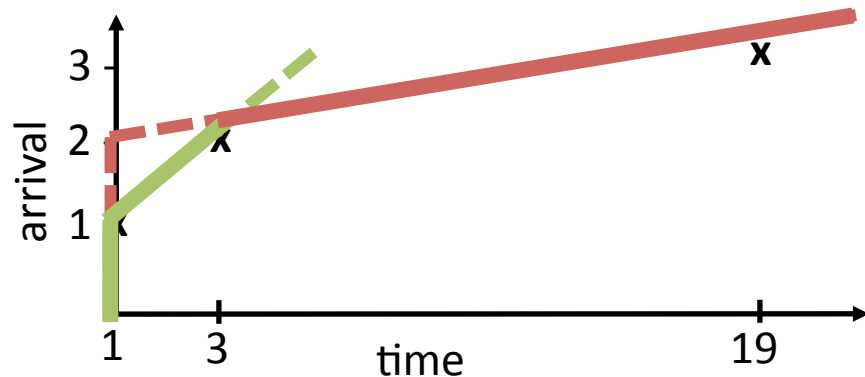
3 pkts in  
10 cycles

2 pkts in  
3 cycles



# Applying Formal Traffic Model

- Non-deterministic choice of allowed destinations
- Rate constraints in xMAS enforced by **token-bucket** regulation of non-deterministic sources
  - Size of token bucket ( $\sigma$ ) constrains traffic bursts
  - Token injection frequency ( $\rho$ ) constrains long-time average rate ( $1/\rho$ )

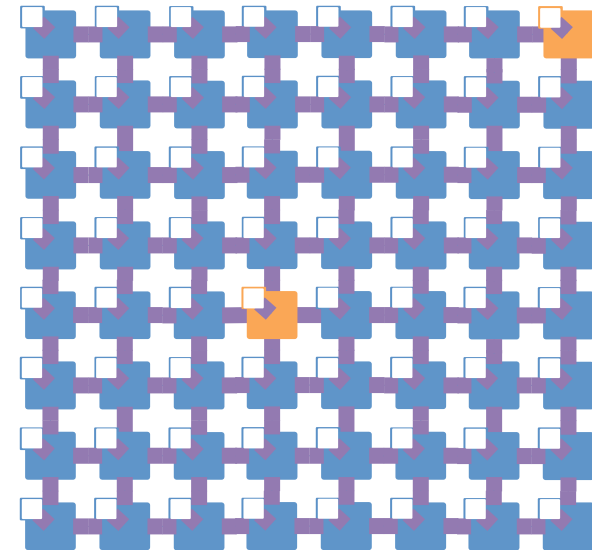


# Traffic modeling -- Conclusions

- UCLID Symbolic Simulation – SMT theory of Bitvectors
  - 30 cycle BMC with traffic model for sources and eager sinks

Node $i$	Simple Model		Fortified Model	
	Runtime [seconds]	Latency [cycles]	Runtime [seconds]	Latency [cycles]
0,0	2160	18	3083	15
1,1	7626	23	17454	20
2,2	5413	24	4444	24
<b>3,3</b>	<b>12060</b>	<b>24</b>	<b>16014</b>	<b>23</b>
4,4	2851	25	4880	24
5,5	5848	24	18555	20
6,6	6486	24	9927	23
<b>7,7</b>	<b>3372</b>	<b>17</b>	<b>4468</b>	<b>14</b>
<b>Non-Det</b>	<b>1621</b>	<b>25</b>	<b>1621</b>	<b>25</b>

Runtimes using Boolector 1.4.1 as solver  
[Brummayer & Biere '09]



Traffic model leads to tighter verified latency bounds

# Summary of Approach

- One strategy for compositional reasoning
  - Decompose both model and latency property along router boundaries
  - Traffic models as interface specs
- Limitations
  - Finding router latency bound is brute force
  - Bounded model checking for local proofs
  - Traffic abstraction has no guarantee of soundness
- Limitations addressed in next work
  - Automated, sound compositions
  - Inductive proofs of latency

# Outline

- Background
  - xMAS: Formal model of NoC
  - Verification technology
- Compositional latency verification
  - Abstraction and traffic modeling
  - **Inductive proof by property strengthening**
- Buffer sizing using counterexamples



# Compositional Verification

- Strengthen latency property with subgoals
- Reduce required unrolling
- Make latency bounds inductive
- Orders of magnitude speedup

Daniel Holcomb, Alexander Gotmanov, Michael Kishinevsky, and Sanjit A. Seshia. **Compositional Performance Verification of NoC Designs**. In *Proceedings of the 10th ACM/IEEE International Conference on Formal Methods and Models for Codesign (MEMOCODE)*, July 2012.

<http://www.eecs.berkeley.edu/~holcomb/memocode12-noc.pdf>

## Compositional Performance Verification of NoC Designs

Daniel E. Holcomb      Alexander Gotmanov      Michael Kishinevsky      Sanjit A. Seshia  
UC Berkeley                  UC Berkeley                  UC Berkeley                  UC Berkeley

**Abstract**—We present a compositional approach to formally verify quality-of-service (QoS) properties of network-on-chip (NoC) designs. A single latency bound is sufficient to describe the worst-case latency for hardware or firmware cycles, which are beyond the scope of state-of-the-art model checkers. We address this challenge by a compositional form of induction. The overall latency bound problem is reduced to a number of subproblems, each of which is tractable. Each latency bound term that a packet must traverse is verified at a particular “stage” of progress. We present a partially-abstracted method of operating these stages based on the topology of the network and a subset of relevant state, and verify the latency bound inductively. The effectiveness of the compositional technique is demonstrated on illustrative examples as well as an industrial chip interconnection network.

### 1. INTRODUCTION

Network on chip (NoC) is a paradigm for communication within large multi-core systems on chip (SoC) designs. In NoC architectures is a network of interconnected nodes, where each node has networking logic and is associated with a processor core, memory controller, specialized IP block, etc. Industrial examples include the TeraXilinx™ processor [1], ST2 Core BE [2], and Intel Larcher [3]. NoCs typically offer quality of service (QoS) guarantees on worst case latency, jitter, and throughput for some classes of network traffic. QoS violations are often similar to starvation and deadlock scenarios and can be easily noticed in performance simulations.

High-level modeling of NoCs [4, 5] and semantic abstraction [6] help to hide unnecessary details, easing the way for formal analysis. Even so, verifying QoS properties can be challenging for industrial NoC designs, both due to the scale of the designs and the property to be verified. Consider for instance the problem of proving an upper bound on the latency of sending a packet from one node in the network to another. In principle, this property can be expressed in linear temporal logic (LTL), and the problem can be solved using model checking. The LTL property expresses a bounded *diverge* property, written in English as “every packet from source A gets to its destination B within  $N$  cycles”. Bounded liveness is essentially a safety assertion where one adds some extra logic to track the progress of time. One can use state-of-the-art model checking strategies such as *k*-abstraction [7], interpolative [8], and IC3-property directed reachability (PDR) [9, 10] to verify this property. However, regardless of the strategy, it is necessary to unroll at least  $k$  consecutive cycles in order to prove or disprove the latency bound, assuming that  $N$  is tight. Typical latency bounds for NoCs can be in hundreds or thousands of clock cycles. Unrolling of model transition relations to such depths is beyond the capacity of most of the state model checking engines. Property directed reachability [9, 10], while avoiding explicit unrolling of the transition relations, still does not scale past one run of clock cycle.

In order to address this challenge, in this paper we use a trial-and-error approach to formal verification: compositional reasoning. In compositional reasoning, one breaks up the overall proof obligation (proving a latency bound of  $N$  cycles) into a number of “smaller” proof sub-goals, which are much easier to verify, such that if all of the sub-goals are proved, then the overall property is proved. The key is to derive a decomposition that is well-suited to the verification task at hand. A natural approach for latency bounds is to first prove smaller bounds on a packet’s progress through the network, e.g., how long does it take to progress past the network, how much time does it spend along a particular abstract, etc. We term these small sub-goals *latency invariants*. Methods to discover and apply them are the core contribution of this paper.

Specifically, we show that for some common network topologies, one can enumerate latently many stages that a packet can progress through. Each location in the network belongs to at least one stage at every time moment. Stages are arranged into a directed, acyclic stage graph, to represent the order in which they can be visited by a packet. A latency bound bounds the number of cycles between when a packet is injected into the network and when the same packet exits some particular stage. By proving the latency invariants for all stages, one proves the bound, corresponding to all paths through the network.

In summary, we make the following novel contributions in this paper:

- A compositional approach to proving latency bound properties in NoC by decomposing into latency invariants.
  - Methods of enumerating latency invariants using a stage graph based on the topology of the network and current state.
  - Experimental results on two illustrative examples show that our approach can reduce the number of induction verification of latency bounds by 50x, and reduce the runtime of the state-of-the-art IC3PDR technique by 2x. However, using a technique to verify latency with ksums added gives a 1x-1.5x overall reduction in runtime relative to verifying the same property using PDR.
  - Experimental results on an industrial-style ring interconnection network show that latency invariants give a significant speedup, and in all cases offer latency bounds to be proved inductively. Several variants of the ring can only be verified within the reduced runtime when latency invariants are used.
- The rest of the paper is organized as follows. Sec. II introduces basic terminology and defines the latency invariants using a simple example. Sec. III describes our compositional approach in detail. Sec. IV describes our strategy for creating a stage graph. Results for illustrative examples are presented in Sec. V, and for the ring

# Scalability of Induction

- Inductive proofs are efficient because of reduced unrolling

- No model decomposition

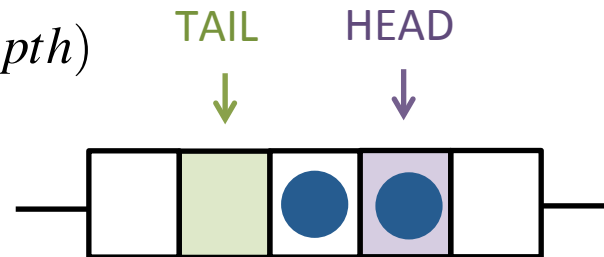
$$p(s_i) \wedge R(s_i, s_{i+1}) \implies p(s_{i+1})$$
$$\left( \bigwedge_{i \in [j, j+k-1]} p(s_i) \wedge \bigwedge_{i \in [j, j+k-1]} R(s_i, s_{i+1}) \right) \implies p(s_{j+k})$$

- Two challenges:
  - General initial state
  - Proof may require large k

# Auxiliary Invariants

- Auxiliary invariant  $\Psi$  blocks off unreachable states from verifier
  - Helps restrict general initial state to good states
  - Otherwise can include deadlock
- Simplicity of xMAS enables helps create  $\Psi$ 
  - Numeric invariants [Chatterjee et al., CAV'10]
  - Channel persistency [Gotmanov et al., VMCAI'11]
  - Structural queue invariants

$$\begin{aligned}tail &= head + num \pmod{depth} \\tail &< depth \\head &< depth\end{aligned}$$



# Verifying Global Latency Bounds

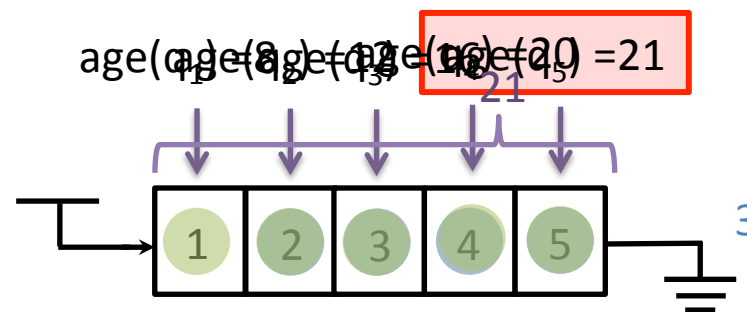
- Prove end-to-end latency of 21 cycles

$$\Phi_{21}^G := \bigwedge_{used_i} age(q_i) < 21$$

- Strengthened with auxiliary invariant

$$\mathcal{N} \models \Phi_{21}^G \wedge \Psi$$

*Induction depth  $k$  is too large to be of much use*



# Strengthening Latency Property

- Strengthen property with Age Lemmas
  - Impose “stage graph” onto xMAS network
  - Precise composable subgoals
- Each age lemma  $\phi_j$  defined by:
  - $t_j$  Precise age bound for stage  $j$
  - $p_{i,j}$  =True iff packet in slot  $i$  maps to stage  $j$

$$\mathcal{N} \models \Phi^G \wedge \Psi \wedge \Phi^L$$

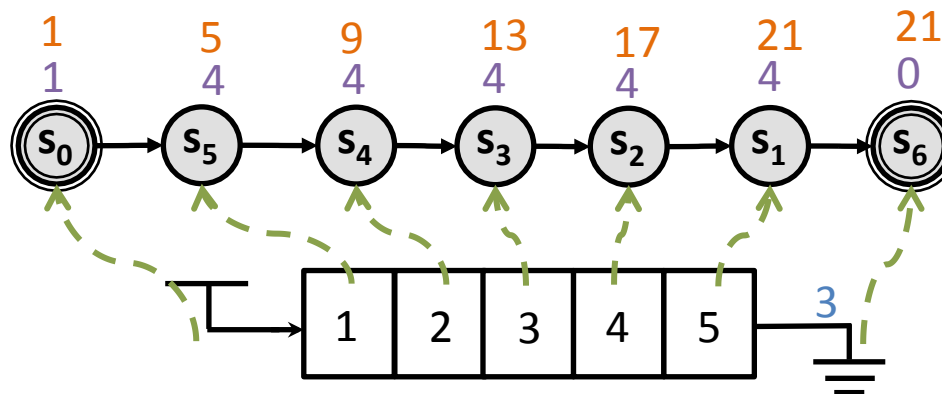
$$\phi_j := \bigwedge_i (p_{i,j} \implies \text{age}(q_i) < t_j)$$

$$\Phi^L := \bigwedge_j \phi_j$$

$$T_L := \max_{j \in \text{stages}} (t_j)$$

$$\Phi^L \implies \Phi_{T_L}^G$$

$d_j = \max$  time  
at stage  $j$

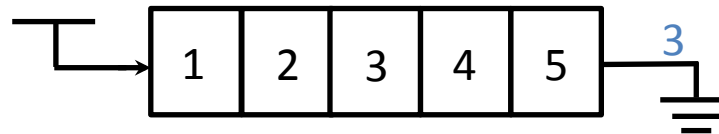
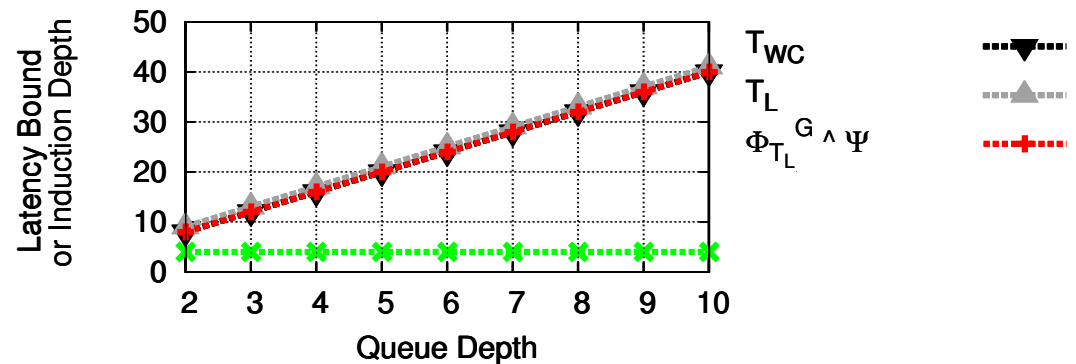
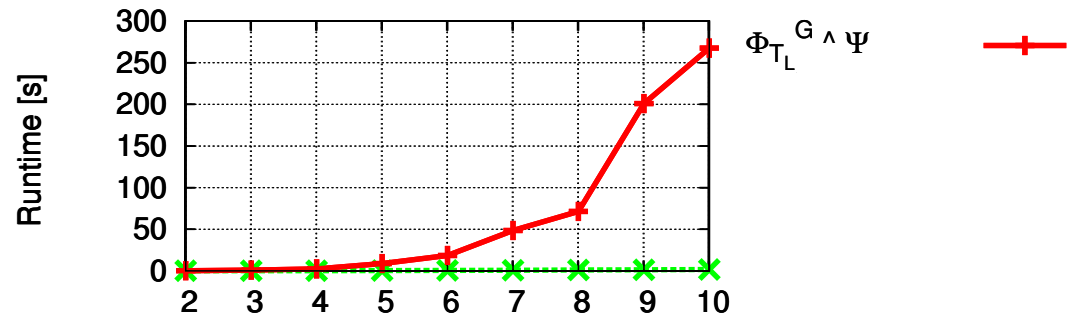


0.37s,  $k=4$   
vs  
8.70s,  $k=20$

# Why Compositional?

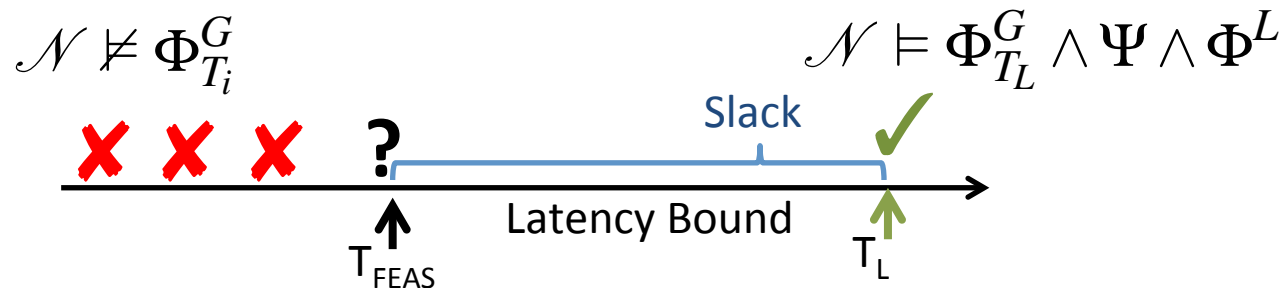
Global bound: Induction depth proportional to worst-case latency

Compositional: Induction depth proportional to size of subgoal



# Tightness of Bound from Stage Graph?

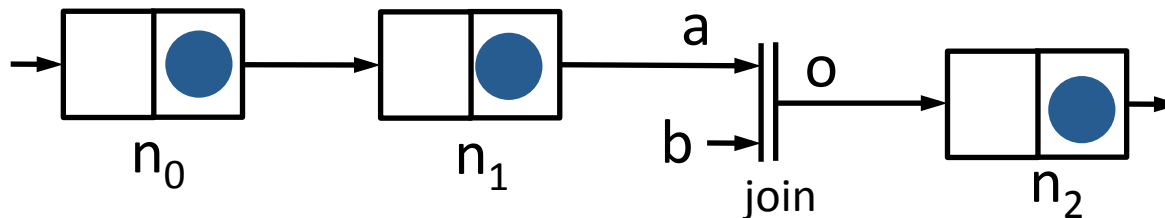
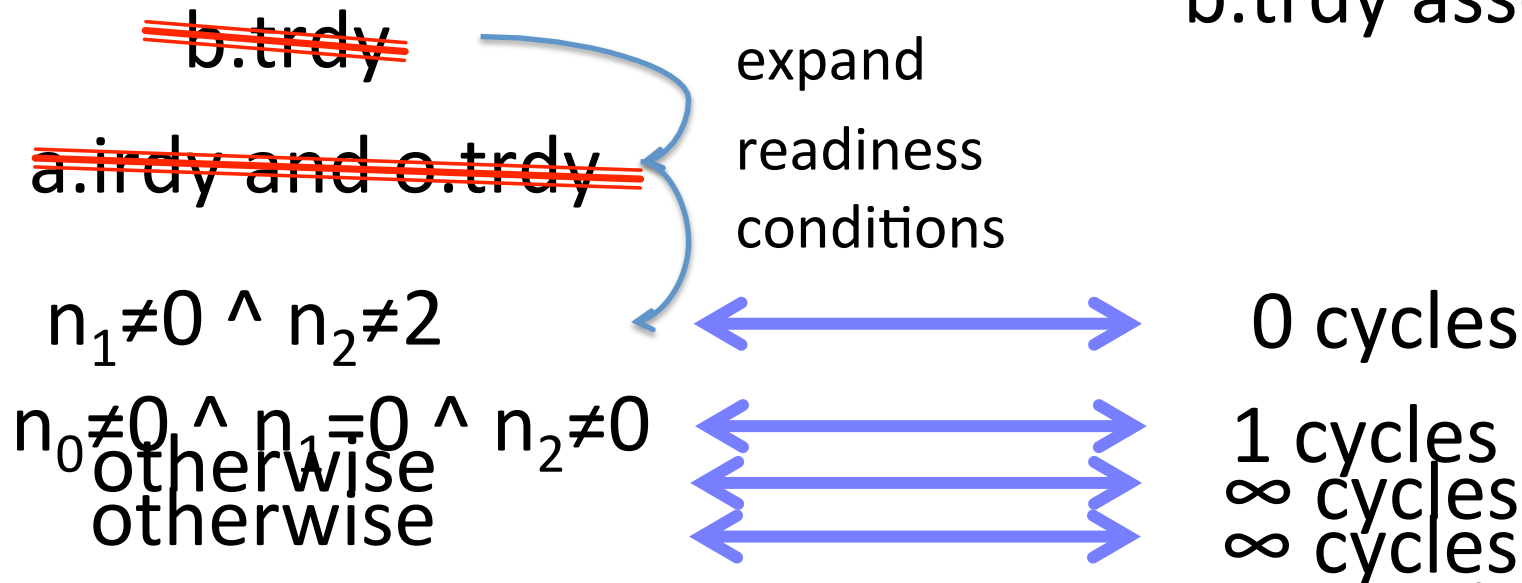
- Age lemmas are conservative with respect to time at each stage of progress
  - But how loose is implied bound  $T_L$ ?
- Use BMC to disprove smaller latency bounds
  - Sweep to find largest disprovable
  - Conservative estimate of slack



# Timing Implications of xMAS Transfer Signals

Boolean condition for channel b target ready?

**How long** until b.trdy asserted?



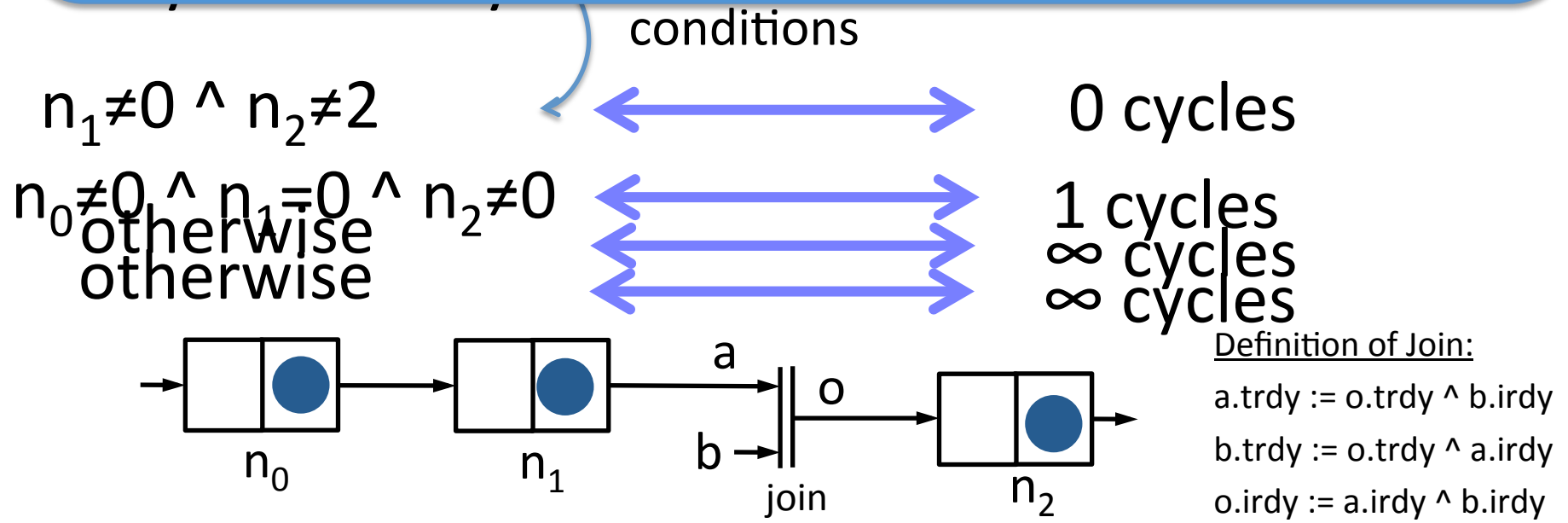
Definition of Join:

$a.trdy := o.trdy \wedge b.irdy$   
 $b.trdy := o.trdy \wedge a.irdy$   
 $o.irdy := a.irdy \wedge b.irdy$



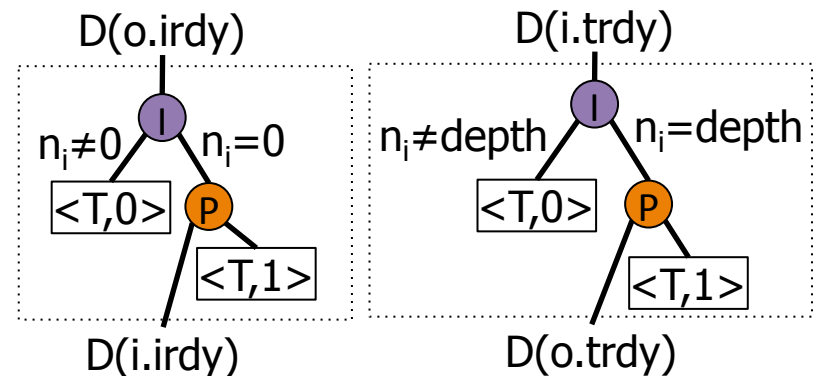
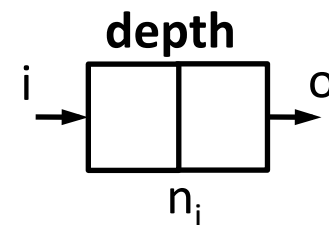
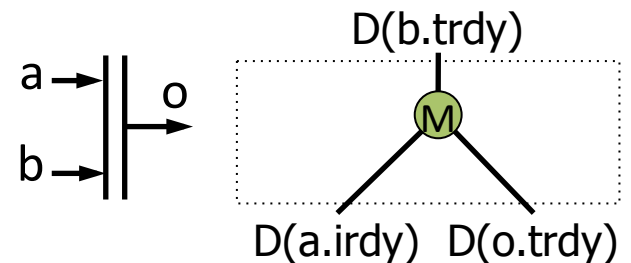
# Timing Implications of xMAS Transfer Signals

Idea: reason about future readiness of signals in way that is analogous to current readiness of Boolean signals

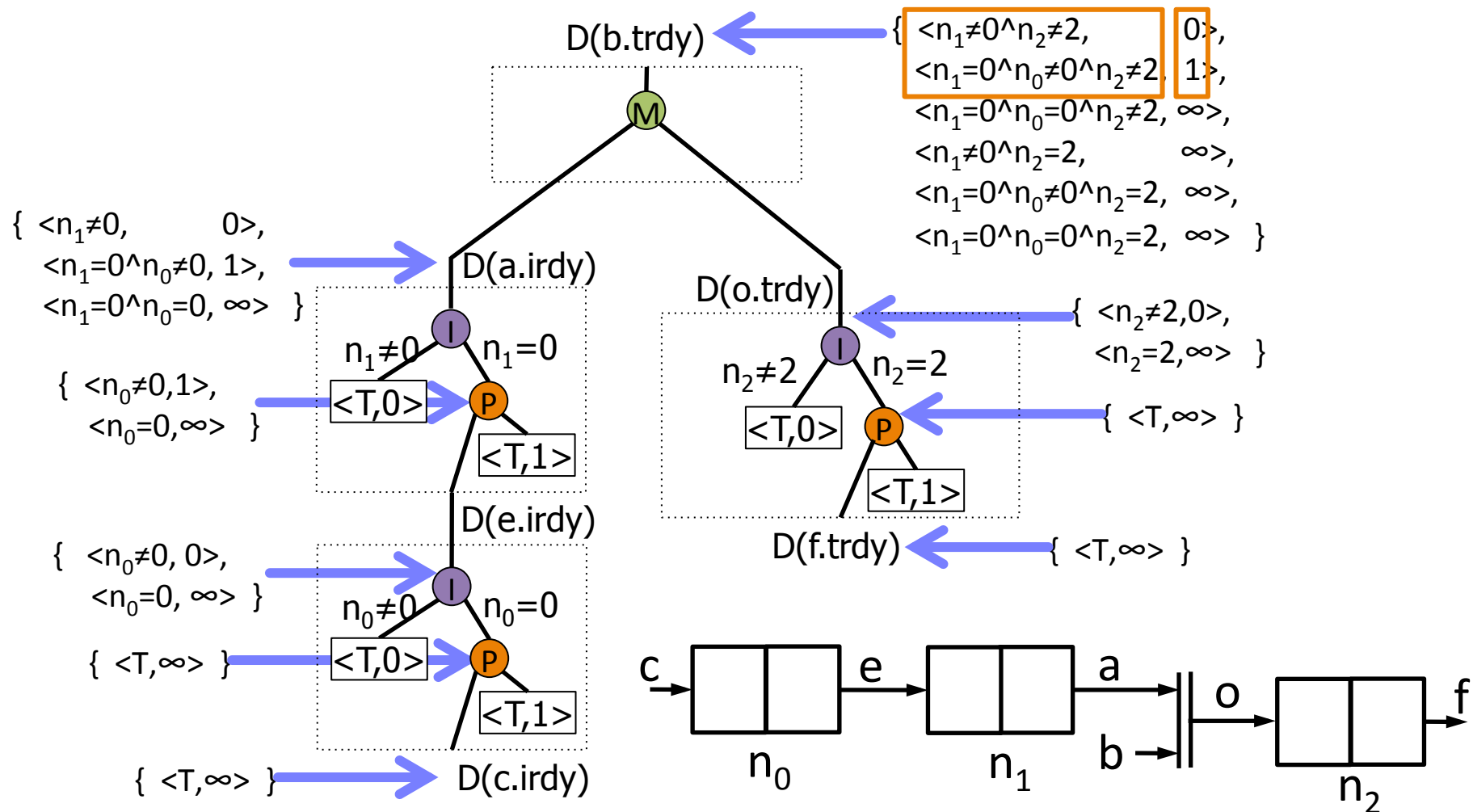


# Future Readiness Bounds

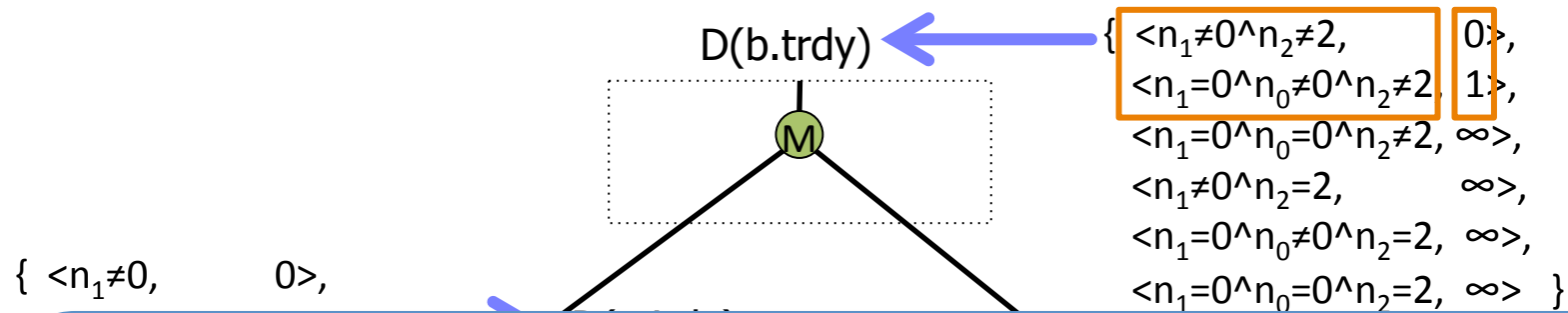
- Want to reason about future readiness time of irdy/trdy
  - Depends on network state
- Represent conditionally:
  - $D(\text{trdy}) = \{ \langle g_0, \delta_0 \rangle, \langle g_1, \delta_1 \rangle, \dots, \langle g_N, \delta_N \rangle \}$
  - $g_i$  is a conjunction of queue predicates
  - $\delta_i$  is a numeric bound on readiness from state satisfying  $g_i$
$$g_i \implies \mathbf{F}^{<\delta_i} \text{trdy}$$
- Rule-based propagation
  - Analogous to Boolean transfer formulas for current readiness in xMAS
  - Manipulated using operations Max, Plus, ITE



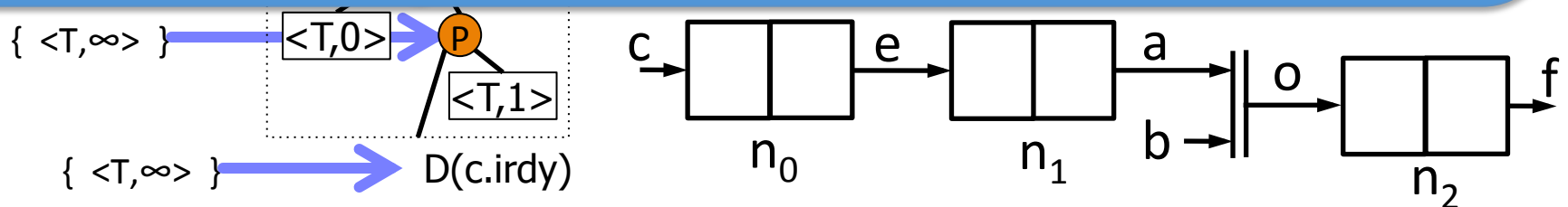
# Deriving Progress Lemmas



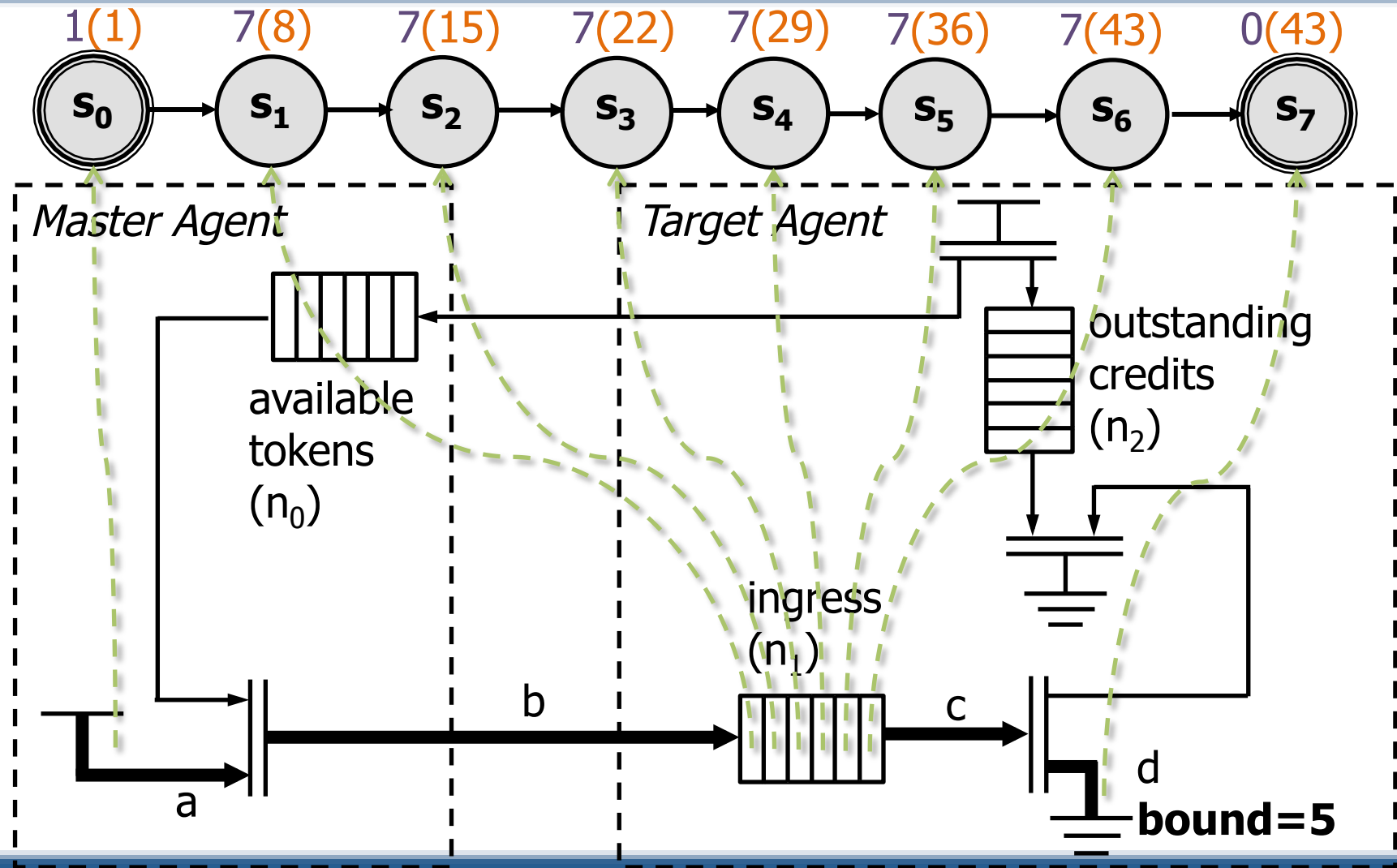
# Deriving Progress Lemmas



- Property  $\Theta$  asserts that network state satisfies a condition with a finite bound
- If  $\Theta$  is valid, then no more than 1 cycles of blocking occur

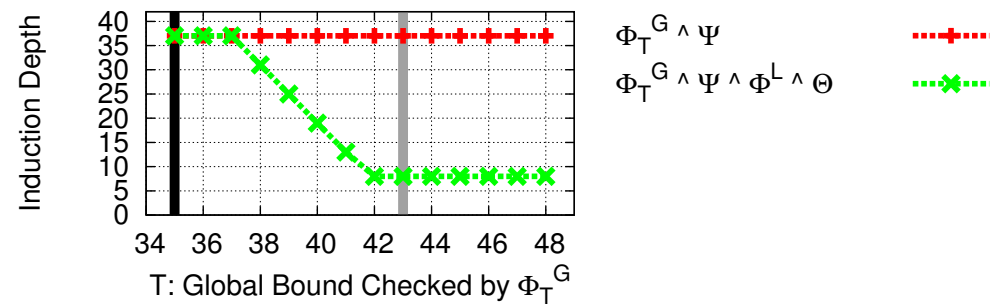
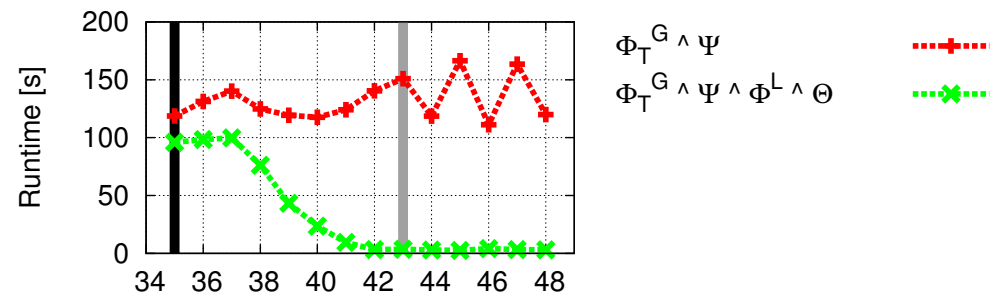


# Credit Loop Network



# Trading Looseness against Runtime

find bound	Runtime (s)	Frames	Cex	Engine	Property
$T_{FEAS} \equiv 35$	52.79	42	Y	bmc	$\Phi_{34}^G$
$\equiv 35$	1045.62	200	-	bmc	$\Phi_{35}^G$

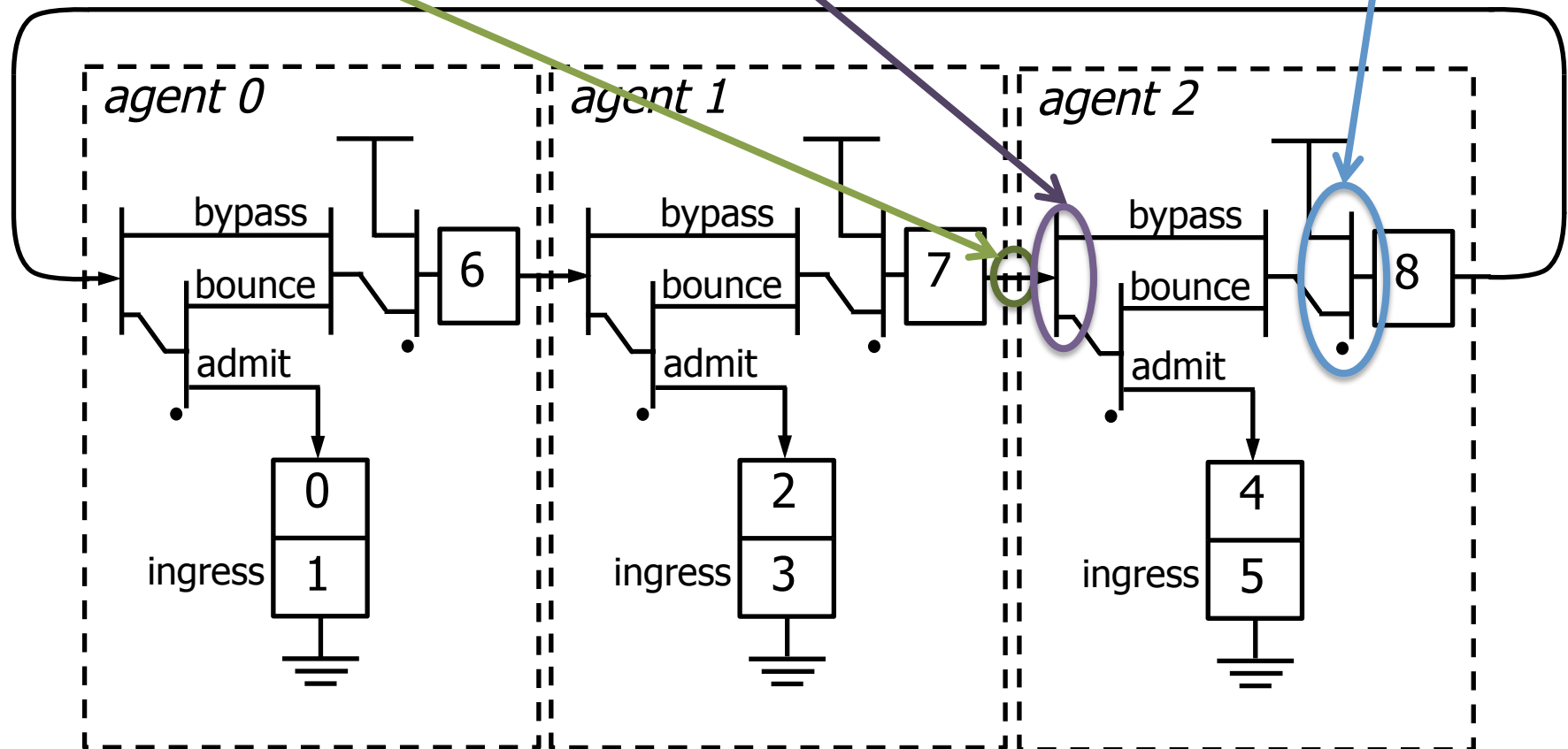


# Non-Stalling Ring Interconnect

Non-blocking channel

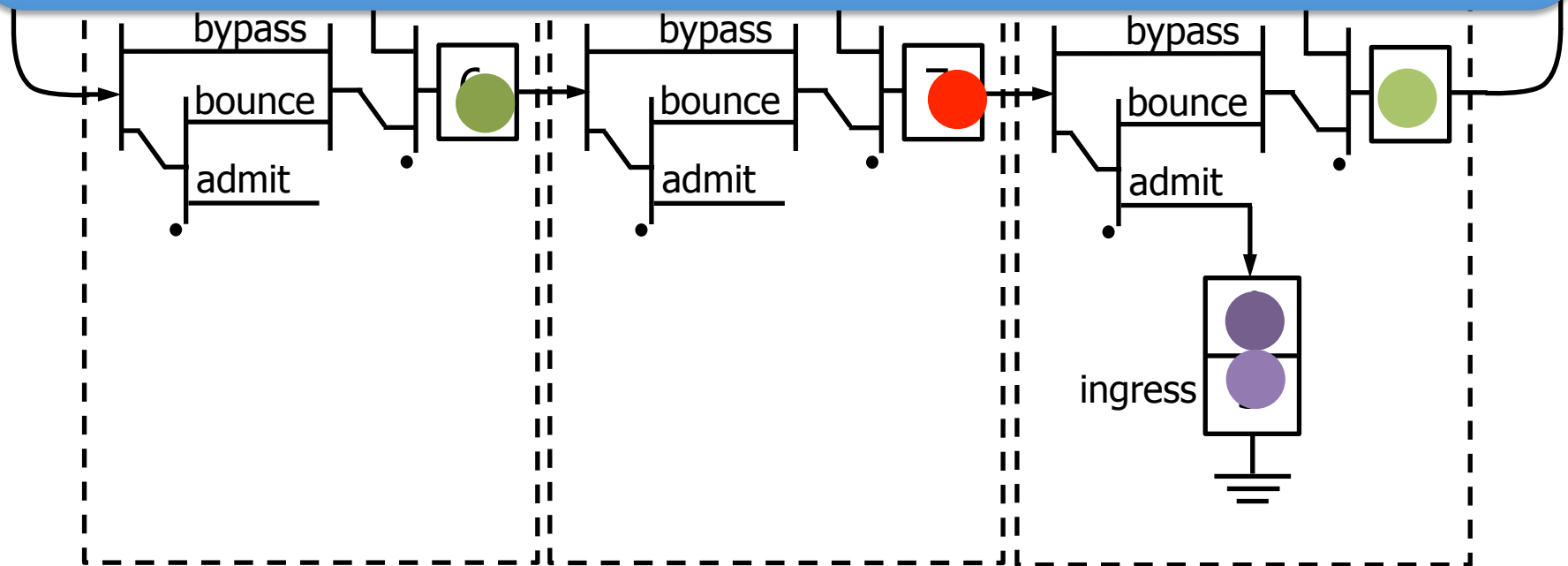
Route to ingress queue if dest=2

Packets on ring get priority



# Non-Stalling Ring Interconnect

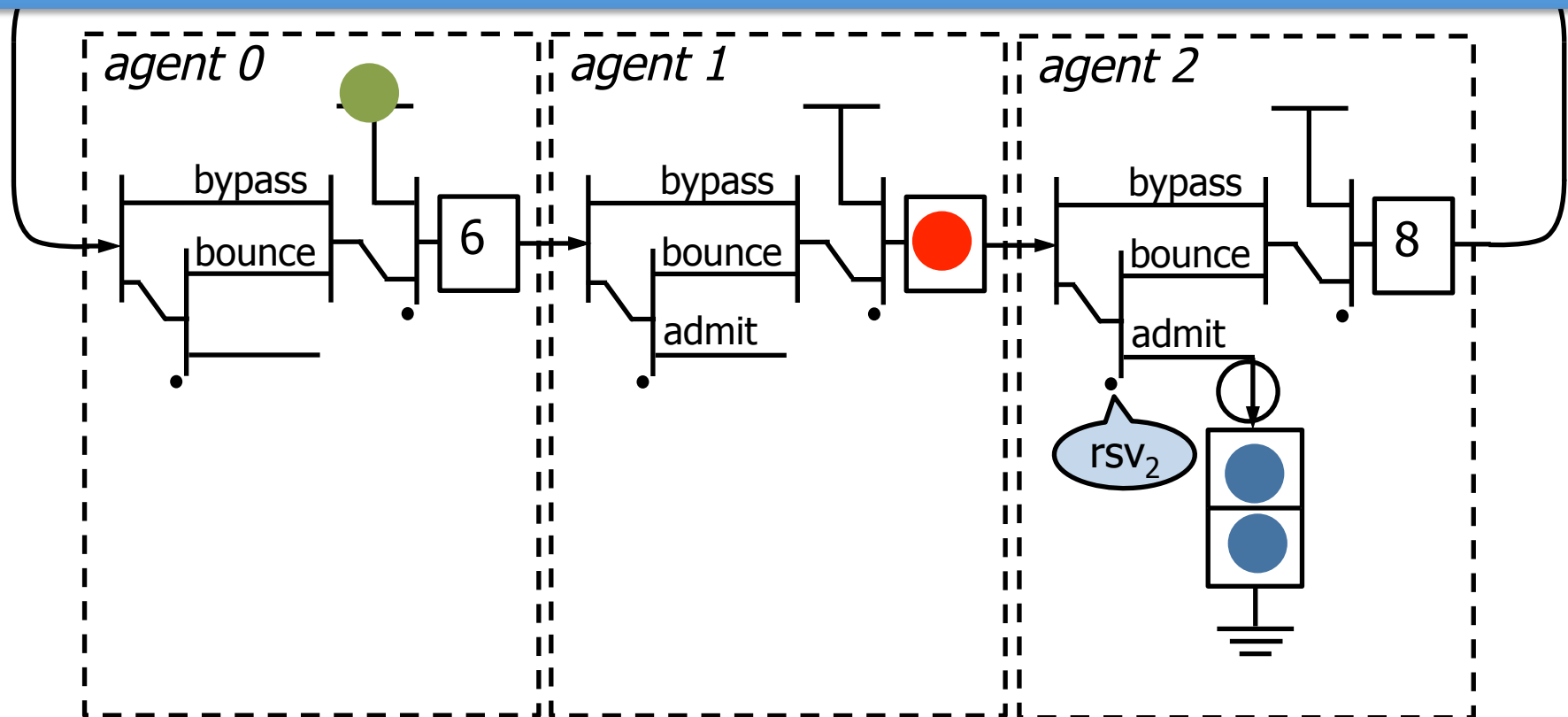
- No progress without reservations



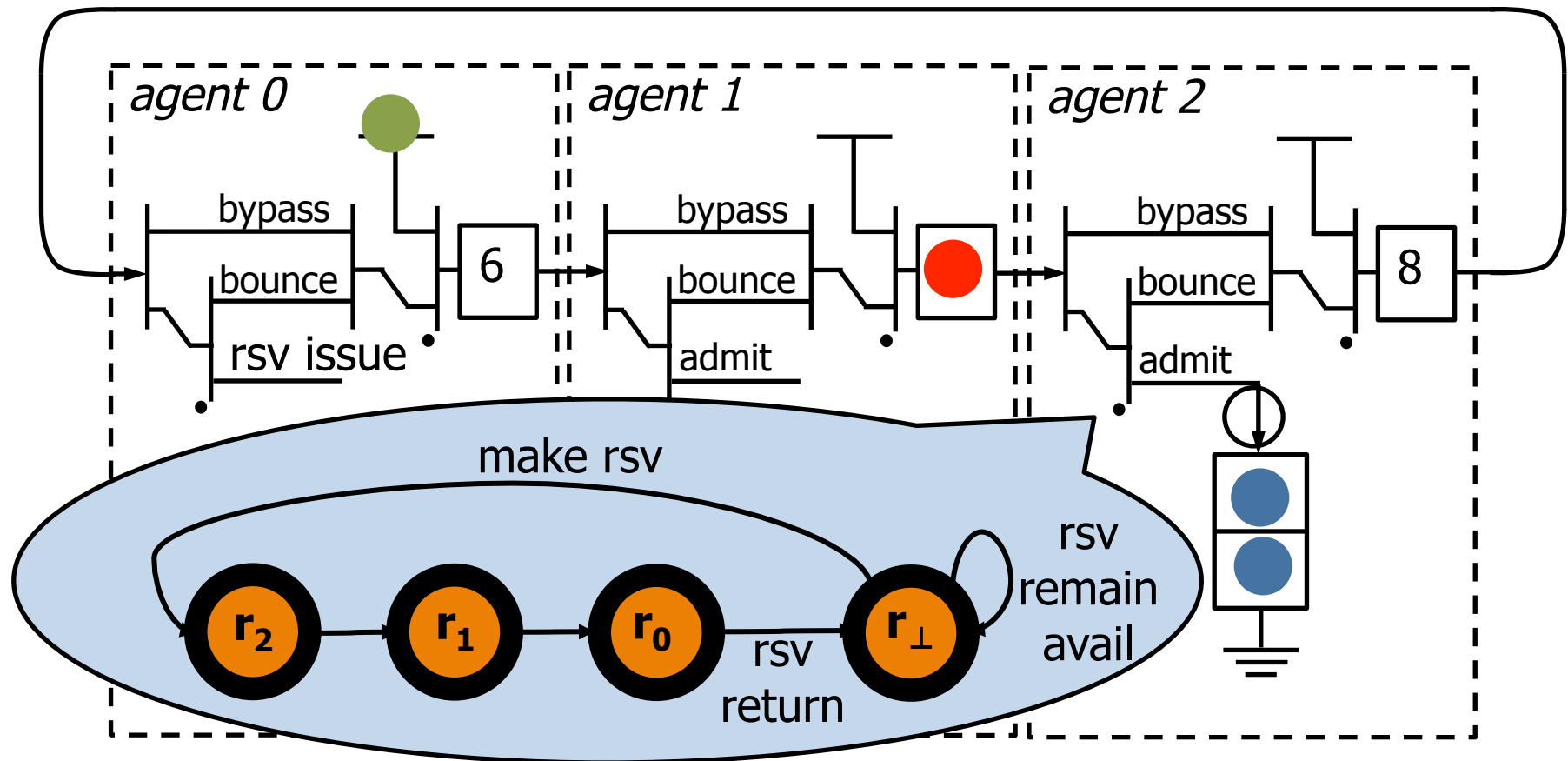


# Receive Reservation in Ring

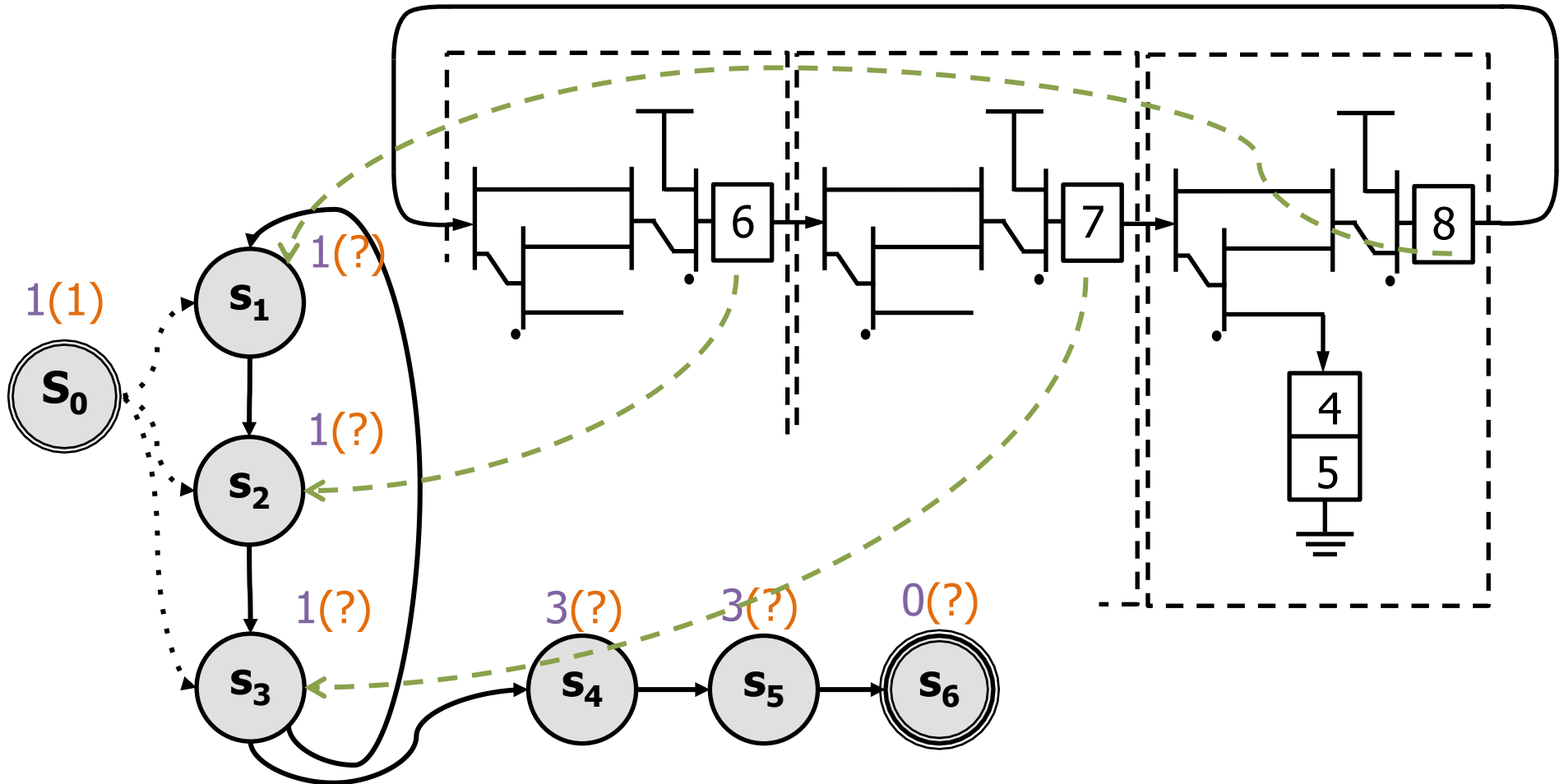
- Livelock prevention to ensure finite latency [Mattina et al]
- Fair ingress admission with respect to packets on the rings



# Receive Reservation in Ring

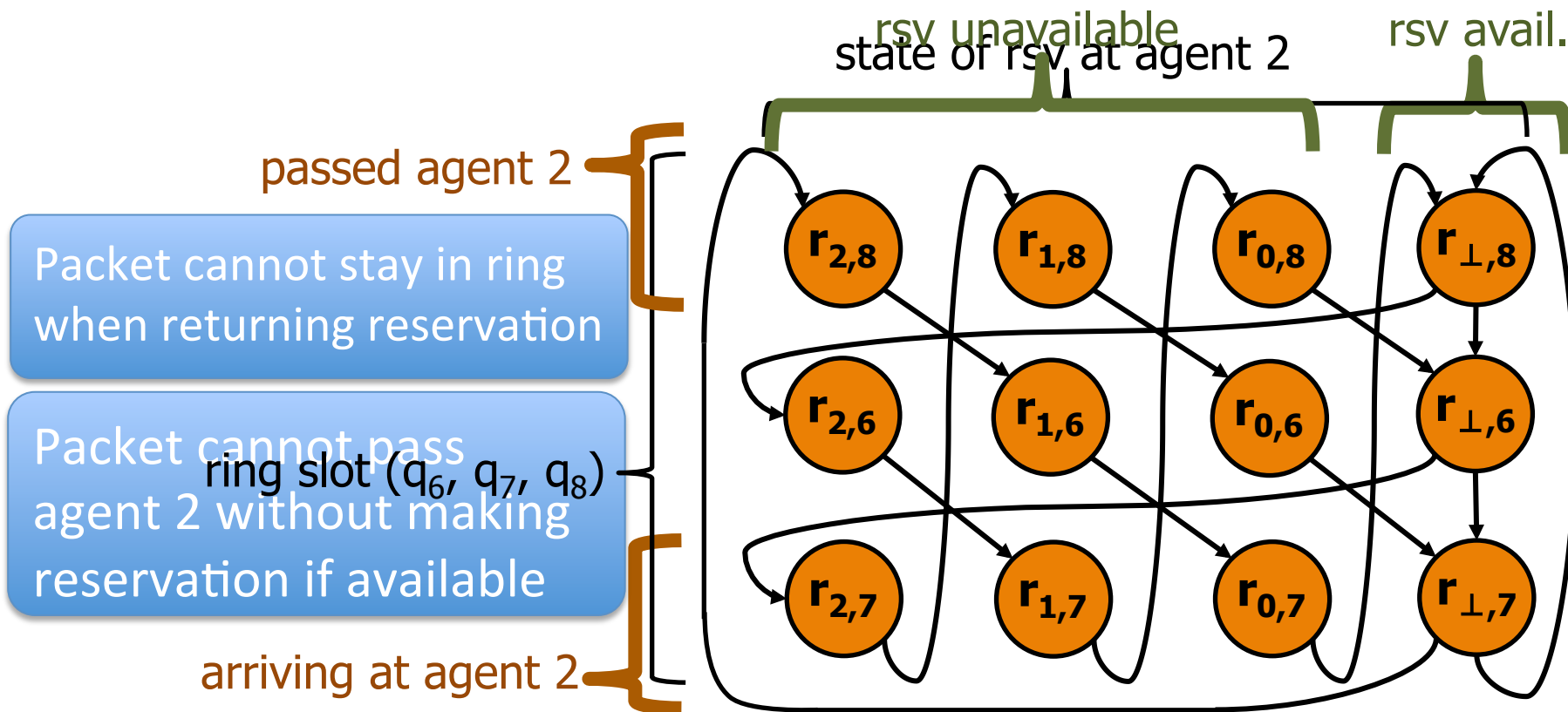


# Naïve Stage Graph of Ring

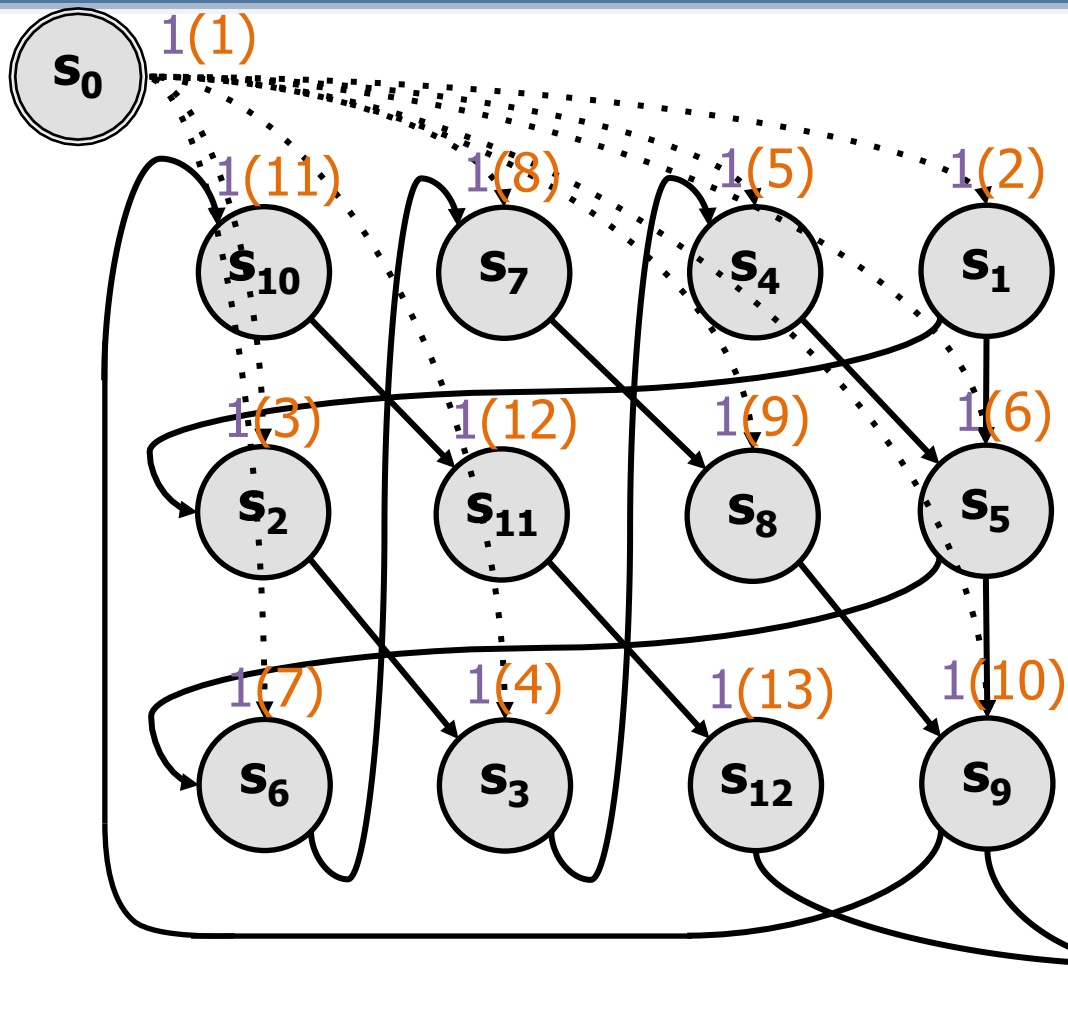


# Toward Acyclic Stage Graph for Ring

- Product automaton of ring and reservation logic



# Stage Graph for 3-Agent Ring



stage $s_j$	$d_j$	$i$	age lemma $\phi_j$ $P_{i,j}$	$t_j$
$s_0$	1	-	-	1
$s_1$	1	8	$dst(q_8) = 2 \wedge rsv_2 = \perp$	2
$s_2$	1	6	$dst(q_6) = 2 \wedge rsv_2 = 2$	3
$s_3$	1	7	$dst(q_7) = 2 \wedge rsv_2 = 1$	4
$s_4$	1	8	$dst(q_8) = 2 \wedge rsv_2 = 0$	5
$s_5$	1	6	$dst(q_6) = 2 \wedge rsv_2 = \perp$	6
$s_6$	1	7	$dst(q_7) = 2 \wedge rsv_2 = 2$	7
$s_7$	1	8	$dst(q_8) = 2 \wedge rsv_2 = 1$	8
$s_8$	1	6	$dst(q_6) = 2 \wedge rsv_2 = 0$	9
$s_9$	1	7	$dst(q_7) = 2 \wedge rsv_2 = \perp$	10
$s_{10}$	1	8	$dst(q_8) = 2 \wedge rsv_2 = 2$	11
$s_{11}$	1	6	$dst(q_6) = 2 \wedge rsv_2 = 1$	12
$s_{12}$	1	7	$dst(q_7) = 2 \wedge rsv_2 = 0$	13
$s_{13}$	3	4	<b>true</b>	16
$s_{14}$	3	5	<b>true</b>	19
$s_{15}$	0	-	-	19

# Ring Verification Results

- Prove 79 cycle bound in 10 frames of induction
- Proved bound is 1 cycle loose
- Speedup of 65x to >130x

find bound	Runtime (s)	Frames	Cex	Engine	Property
$T_{FEAS}$	24.12	20	Y	bmc	$\Phi_{17}^G$
$\equiv 18$	868.28	200	-	bmc	$\Phi_{18}^G$

verify bound	Runtime (s)	Frames	Proved	Engine	Property
$T_L$	62.34	18	Y	kind	$\Phi_{19}^G \wedge \Psi$
$\equiv 19$	<b>1.31</b>	<b>4</b>	<b>Y</b>	<b>kind</b>	$\Phi_{19}^G \wedge \Psi \wedge \Phi^L \wedge \Theta$
	88.39	12	Y	pdr	$\Phi_{19}^G \wedge \Psi$
	6.57	14	Y	pdr	$\Phi_{19}^G \wedge \Psi \wedge \Phi^L \wedge \Theta$

3-agent ring





find bound	Runtime (s)	Frames	Cex	Engine	Property
$T_{FEAS}$	3901.95	80	Y	bmc	$\Phi_{77}^G$
$\equiv 78$	10,000.00	111	-	bmc	$\Phi_{78}^G$

verify bound	Runtime (s)	Frames	Proved	Engine	Property
$T_L$	10,000.00	-	-	kind	$\Phi_{79}^G \wedge \Psi$
$\equiv 79$	<b>75.28</b>	<b>10</b>	<b>Y</b>	<b>kind</b>	$\Phi_{79}^G \wedge \Psi \wedge \Phi^L \wedge \Theta$
	10,000.00	-	-	pdr	$\Phi_{79}^G \wedge \Psi$
	662.15	73	Y	pdr	$\Phi_{79}^G \wedge \Psi \wedge \Phi^L \wedge \Theta$

8-agent ring

# Inductive Proof via Strengthening -- Conclusions

- Inductive verification of latency property on a bit-level RTL network from xMAS model
  -  Orders of magnitude verification speedup
  -  Models complex arbitration behaviors
  -  Sound composition
  -  Not automated for cyclic networks

Knowledge of higher-level network structure produces useful invariants for bit-level verification

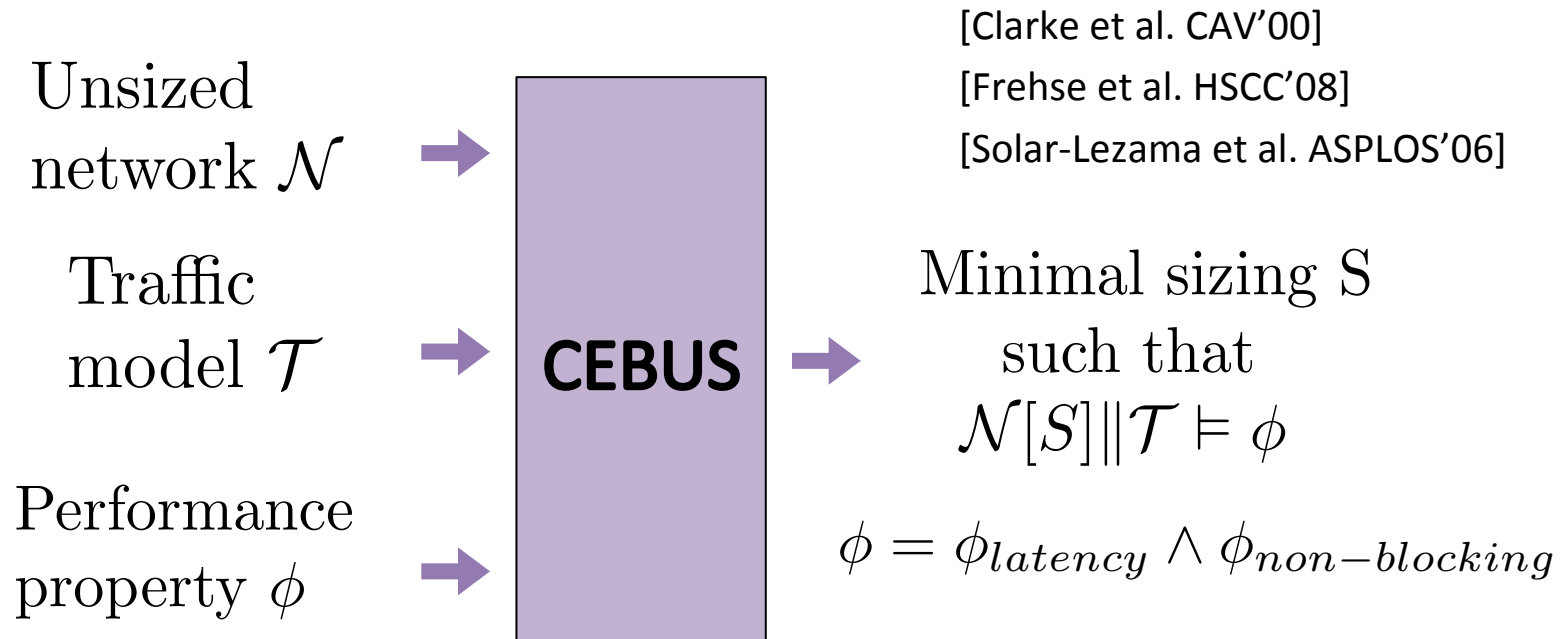
# Outline

- Background
  - xMAS: Formal model of NoC
  - Verification technology
- Compositional latency verification
  - Abstraction and traffic modeling
  - Inductive proof by property strengthening
- **Buffer sizing using counterexamples**





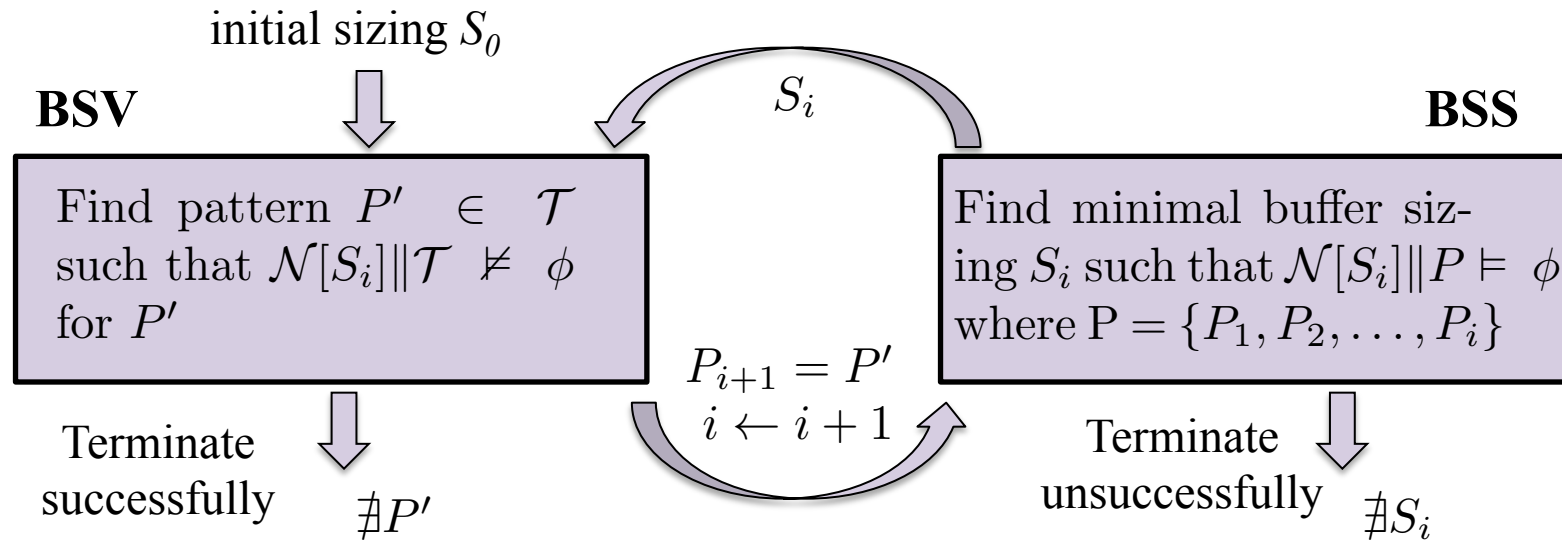
# Buffer Sizing



- Need to address quantifier alternation

$$\exists S_i \forall P' \in \mathcal{T}. (\mathcal{N}[S_i] \parallel P' \models \phi)$$

# Quantifier Instantiation



Check validity of formula

$$\mathcal{N}[S_i] \parallel \mathcal{T} \models \phi$$

- **Valid:** sizing  $S_i$  ensures  $\phi$
- **Invalid:** counterexample is traffic pattern  $P'$  that causes  $\neg\phi$

Check validity of formula

$$|S_i| < size \implies \neg(\mathcal{N}[S_i] \parallel P \models \phi)$$

- **Valid:** no  $S_i$  ensures  $\phi$  for all  $P$
- **Invalid:** Counterexample is sizing  $S_i$  that ensures  $\phi$

# Buffer Sizing -- Conclusions

Completeness has a cost -- even for network of 3 queues

## Verification

- Checking same formula at each iteration  $i$
- More difficult SAT problem at later iterations
  - Few satisfying assignments to be found

Iteration $i$	<i>Buffer-Size Verification (BSV)</i>					
	CNF Size		Runtime (sec)			
	Vars	Clauses	SAT	Enc	SSim	Total
0	158K	473K	1.4	17.5	1.3	20.2
2	158K	473K	0.9	17.7	1.2	19.9
4	158K	473K	8.9	17.7	1.3	28.0
6	158K	473K	16.8	17.7	1.3	35.8
8	158K	473K	91.7	17.9	1.3	111.0
10	158K	473K	242.8	17.7	1.3	261.8
11	158K	473K	106.0	17.9	1.2	125.3
12*	158K	473K	373.5	16.2	2.7	392.5

## Synthesis

- Problem size linear in  $i$
- Binary search to minimize
- Difficulty of SAT problem not necessarily proportional to size
  - Symbolic sim. and decision proc. encoding dominate runtime

Iteration $i$	<i>Buffer-Size Synthesis (BSS)</i>					
	CNF Size		Runtime (sec)			
	Vars	Clauses	SAT	Enc	SSim	Total
0	79K	236K	5.6	69.0	11.0	85.7
2	100K	300K	2.8	60.8	16.0	79.7
4	217K	650K	5.7	79.7	25.1	110.6
6	259K	778K	15.0	195.3	60.2	270.6
8	345K	1037K	27.2	234.5	91.4	353.2
10	429K	1287K	39.7	342.2	126.7	508.8
11	486K	1458K	61.0	392.9	155.6	609.6

# Summary

- Formal methods are promising approach for synthesizing and verifying NoC QoS
- Achieved several orders of magnitude speedup over monolithic model checking of latency properties without strengthening
- Not yet a push-button solution for general network models, but xMAS methodology helps

# Thank you

**Leverage model checking for solving NoC QoS latency problems. Address capacity limitations by extending well-known formal techniques including abstraction and compositional reasoning into the NoC domain**

- 3 specific QoS contributions
    - Workload abstraction of traffic models
    - Latency proofs by property strengthening
    - Optimal buffer sizing for QoS
- } Compositional Latency Verification