



# Buffer-aware Worst-Case Timing Analysis of Wormhole NoCs Using Network Calculus\*

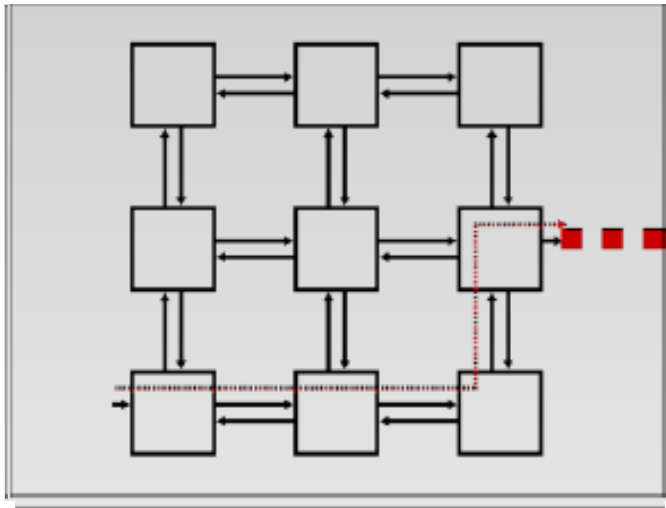
*Frédéric Giroudot and Ahlem Mifdaoui*

WONECA 2018

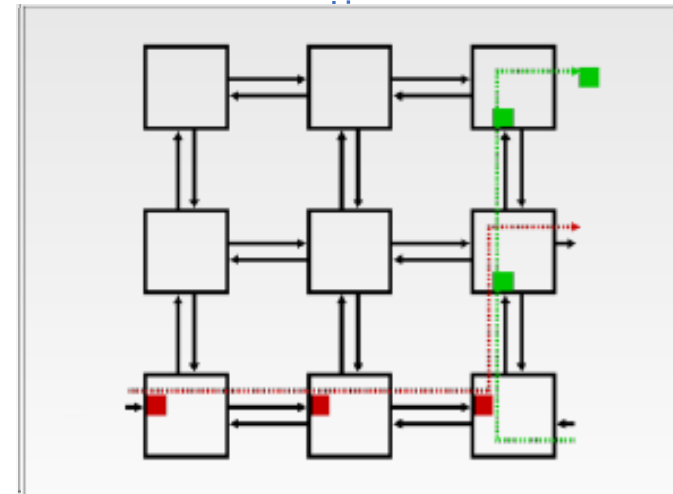
\*Accepted to appear in RTAS'2018

# What is challenging in Wormhole NoCs?

No contention



Under contention



The wormhole routing:

- + Reduce drastically the storage buffers in routers
  - + The contention-free packet latency becomes insensitive to the path length
  - Complicate the congestion pattern
  - Introduce indirect blocking delays due to buffer backpressure
- ⇒ Need appropriate timing analysis to compute safe delay bounds

# Related work

Approach	Contribution	wormhole	multiple VCs	priority sharing	VCs sharing	flows serialization	buffer size $B$ (vs packet length $L$ )		
							1 flit	$L \leq B$	$B \leq L$
Scheduling Theory	[2]	X	X				X	X	X
	[4]	X	X	X					
	[5]	X	X	X			X		
CPA	[6]	X	-	X				X	
	[7]	X	X	X	X			X	
Network Calculus	[8]	X		X			X	X	X
	[9]		X	X		X		X	
	<b>our approach</b>	X	X	X	X	X	X	X	X

- [2] Q. Xiong, F. Wu, Z. Lu, and C. Xie, "Extending real-time analysis for wormhole nocs," *IEEE Transactions on Computers*, vol. PP, no. 99, pp. 1–1, 2017.
- [3] Z. Shi and A. Burns, "Real-time communication analysis for on-chip networks with wormhole switching," in *Networks-on-Chip, Second ACM/IEEE International Symposium on*, April 2008.
- [4] Z. Shi and A. Burns, "Real-time communication analysis with a priority share policy in on-chip networks," in *21st Euromicro Conference on Real-Time Systems*, pp. 3–12, July 2009.
- [5] M. Liu, M. Becker, M. Behnam, and T. Nolte, "Tighter time analysis for real-time traffic in on-chip networks with shared priorities," in *10th IEEE/ACM International Symposium on Networks-on-Chip*, 2016.
- [6] S. Tobuschat and R. Ernst, "Real-time communication analysis for networks-on-chip with backpressure," in *Design, Automation Test in Europe Conference Exhibition*, 2017.
- [7] E. A. Rambo and R. Ernst, "Worst-case communication time analysis of networks-on-chip with shared virtual channels," in *Proceedings of Design, Automation Test in Europe Conference Exhibition*, 2015.
- [8] Y. Qian, Z. Lu, and W. Dou, "Analysis of worst-case delay bounds for best-effort communication in wormhole networks on chip," in *Networks-on-Chip, 3rd ACM/IEEE International Symposium on*, May 2009.
- [9] F. Jafari, Z. Lu, and A. Jantsch, "Least upper delay bound for vbr flows in networks-on-chip with virtual channels," *ACM Trans. Des. Autom. Electron. Syst.*, vol. 20, pp. 35:1–35:33, June 2015.

⇒ **None** of the existing approaches cover **all** the implemented **mechanisms** and/or **phenomena**

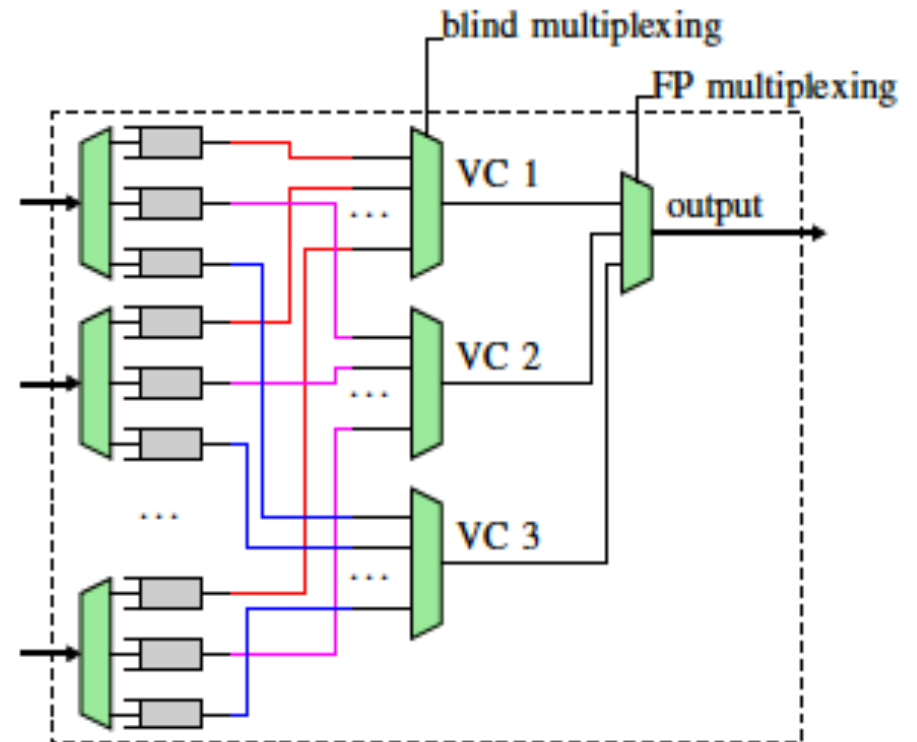
⇒ **Our proposal**: a new **buffer-aware** timing analysis considering the flows **serialization** phenomena based on NC

# Outline

- Context & Problematic
- **Main Contributions**
  - ✓ **System Model and Assumptions**
  - ✓ **Main Steps of the Buffer-aware Timing Analysis**
  - ✓ **Performance Evaluation**
- Conclusions & Perspectives

# System Model and Assumptions

- **Input-buffered** routers
- **VCs sharing**, i.e., a VC supports many traffic classes
- **Priority sharing**, i.e., many flows mapped on the same priority-level
- **Arbitrary<sup>(1)</sup> multiplexing** of flows within the same VC
- **Priority-based** arbitration of VCs
- **Flit-level preemption**
- **Rate-latency** service curve for each output port
- **Leaky-bucket** arrival curve for each flow



(1) To cover different service policies, such as FIFO and RR

# Main Steps of the Buffer-aware Timing Analysis

**Main idea:** to compute upper bound on end-to-end delay for a *flow*  $f$ , we need the granted **end-to-end service curve** to  $f$ :

$$\beta_f(t) = R_f (t - T_f )^+$$

Where:

- $R_f$ : the bottleneck rate along the flow path
- $T_f$ : the service latency

$$T_f = T_{hp} + T_{sp} + T_{lp} + T_{IB} + T_{P_f}$$

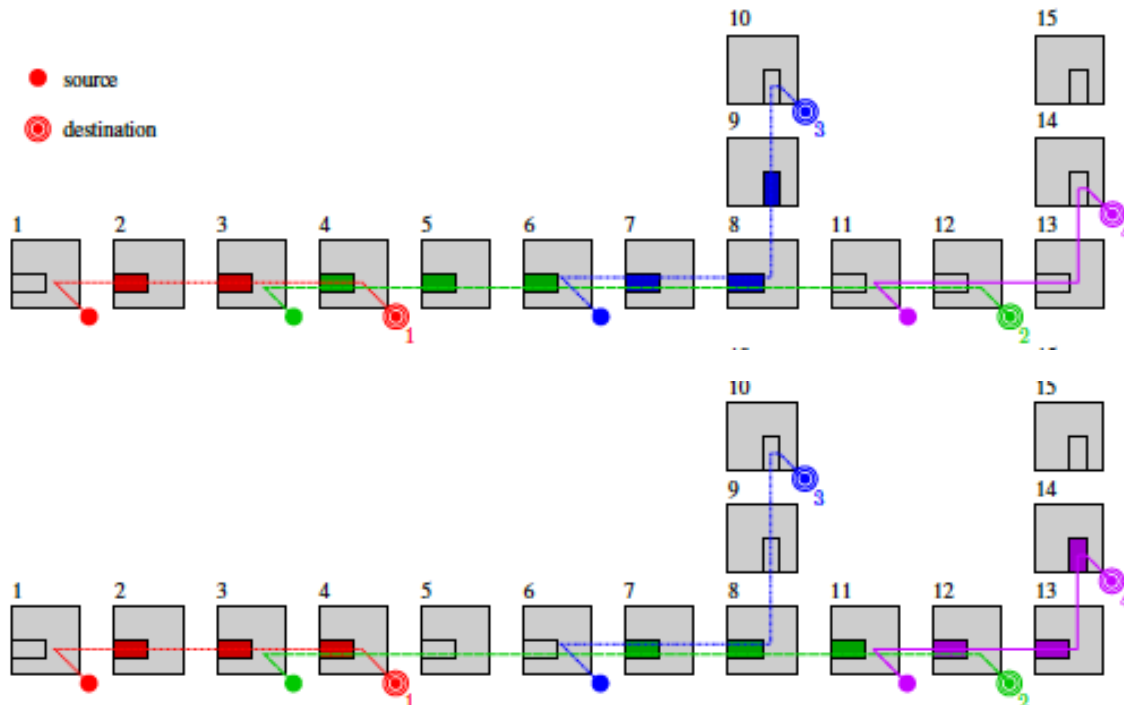
Direct blocking  
due to interfering flows

Indirect blocking  
due to backpressure

Technological latencies  
within routers

# #Step1 of the Buffer-aware Timing Analysis

## Indirect Blocking flows set



- One-flit buffers
- 3-flit long packets
- One VC
- *foi* flow 1
- $IB_1 = \{3,4\}$  without taking into account the buffer size
- $IB_1' = \{3\}$  under **buffer-aware** analysis
- **The buffer size may limit the indirect blocking set (delay)**

# #Step1 of the Buffer-aware Timing Analysis

- Find flows blocking the foi even though they do not share resources with it ( $IB_{foi}$ )
- Determine which **section** of the IB flow's path is **involved**
- Quantify the packet **spread index** of each IB flow  $f$ :

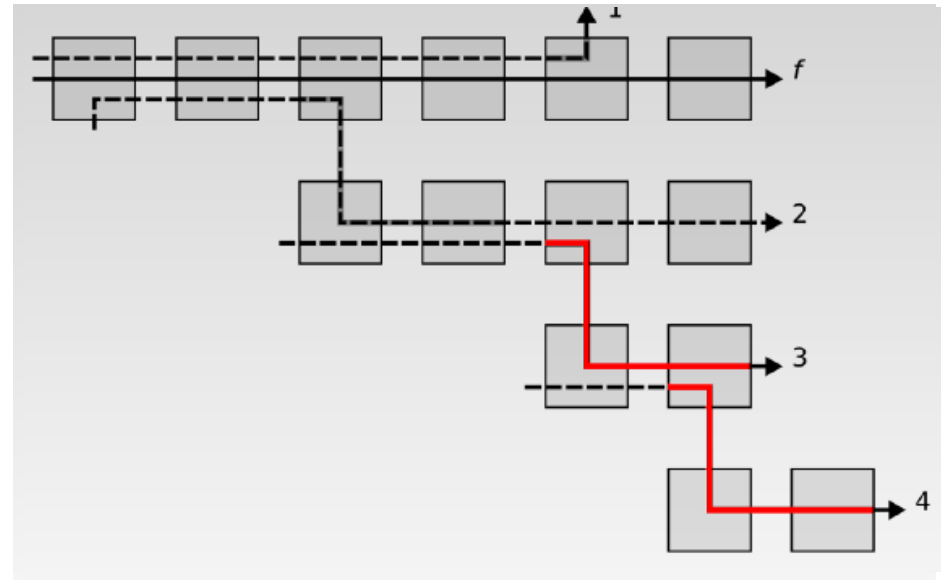
$$N_f = \left\lceil \frac{L_f}{B} \right\rceil$$

- **Propagate** spread sections from the **divergence** point to compute  $IB_{foi}$

⇒ *The complexity of the computation algorithm of  $IB_{foi}$  is linear:*

$$\mathcal{C}(|\mathcal{F}|) = \mathcal{O} \left( \underbrace{|\mathcal{F}|}_{\text{Number of flows}} \cdot \left( K \cdot \max_{f \in \mathcal{F}} \underbrace{|\mathbb{P}_f|}_{\text{f-path length}} \right) \right)$$

Number of flows      f-path length





## #Step2 of the Buffer-aware Timing Analysis

### Theorem [Maximum Direct Blocking Latency]

The maximum direct blocking latency for a *foi*  $f$  along its path  $\mathbb{P}_f$ , in a NoC under **flit-level preemptive FP** multiplexing with strict service curve nodes of the **rate-latency type** and **leaky-bucket** constrained arrival curves is:

$$\begin{aligned}
 T_{hp} &= \sum_{i \in DB_f \cap hp(f)} \frac{T_{\mathbb{P}_f} + T_{hp} + T_{sp} + T_{lp} + \sigma_i^{cv(i,f)} + \rho_i \cdot \sum_{r \in \mathbb{P}_f \cap \mathbb{P}_i} \left( T^r + \frac{L_{slp(f)}^r}{R^r} \right)}{R_\varepsilon} \\
 T_{lp} &= \sum_{r \in \mathbb{P}_f} \frac{L_{slp(f)}^r}{R^r} \quad T_{sp} = \sum_{i \in DB_f \cap sp(f)} \frac{\sigma_i^{cv(i,f)} + \rho_i \cdot \sum_{r \in \mathbb{P}_f \cap \mathbb{P}_i} \left( T^r + \frac{L_{slp(f)}^r}{R^r} \right)}{R_f} \\
 L_{slp(f)}^r &= \max \left( \max_{j \in sp(f)} \left( L_j \cdot \mathbf{1}_{\{sp(f) \supset r\}} \right), S_{flit} \cdot \mathbf{1}_{\{lp(f) \supset r\}} \right) \quad R_f = \min_{r \in \mathbb{P}_f} \left\{ R^r - \sum_{j \ni r, j \in shp(f)} \rho_j \right\}
 \end{aligned}$$

## #Step3 of the Buffer-aware Timing Analysis

### Theorem [Indirect Blocking Latency]

The maximum indirect blocking latency for a *foi*  $f$  along its path  $P_f$ , in a NoC under **flit-level preemptive FP** multiplexing with strict service curve nodes of the **rate-latency** type and **leaky-bucket** constrained arrival curves is:

$$T_{IB} = \sum_{(k, subP) \in IB_f} \frac{\sigma_k^{subP[0]}}{\tilde{R}_k^{subP}} + \tilde{T}_k^{subP}$$

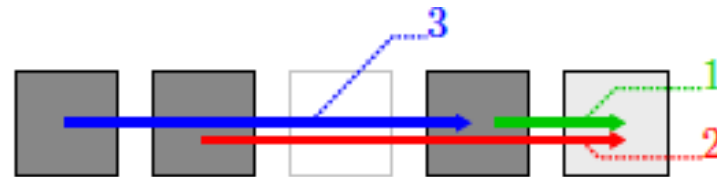
$$\tilde{R}_k^{subP} = \min_{r \in subP} \left\{ R^r - \sum_{j \ni r, j \in hp(f)} \rho_j \right\}$$

$$\tilde{T}_k^{subP} = \sum_{r \in subP} \left( T^r + \frac{S_{flit} \mathbf{1}_{\{lp(k) \supset r\}}}{R^r} \right)$$

$$+ \sum_{i \in DB_k^{subP} \cap hp(k)} \frac{\sigma_i^{cv(i,k)} + \rho_i \sum_{r \in subP \cap \mathbb{P}_i} \left( T^r + \frac{S_{flit} \mathbf{1}_{\{lp(k) \supset r\}}}{R^r} \right)}{\tilde{R}_k^{subP}}$$

# Performance Evaluation (1)

## Comparative analysis vs Scheduling Theory approaches

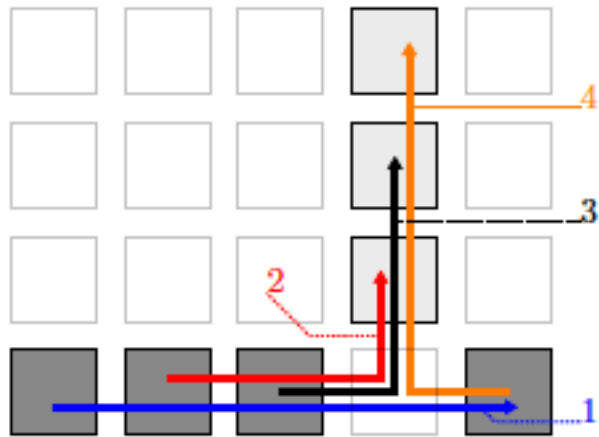


Flow index	1	2	3
Priority	1	2	3
Period	100	100	100
Deadline	100	100	40
Release jitter	0	0	0
Base latency (no contention)	21	24	14
Packet size	19	20	10
Cycle accurate scenario in [2]	21	43	43
Upper bound by [3]	21	45	38
Upper bound (our approach)	23	57	44

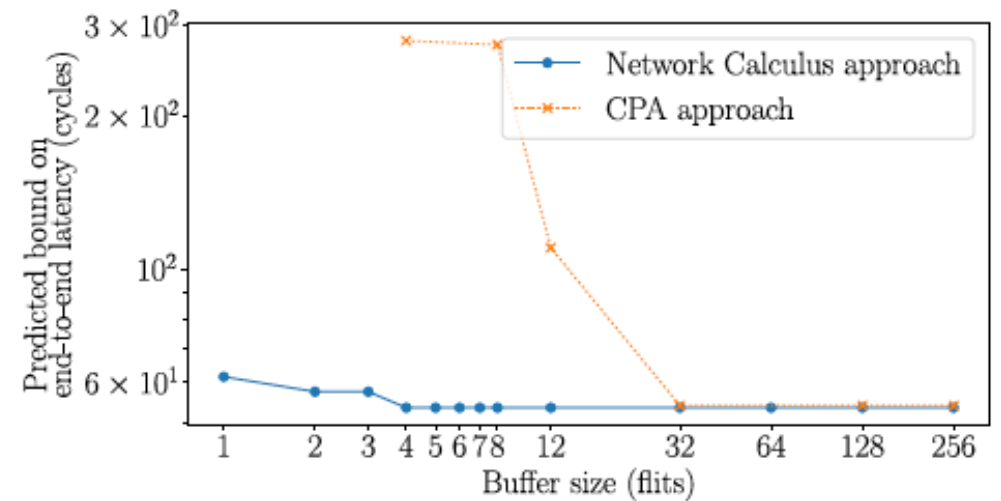
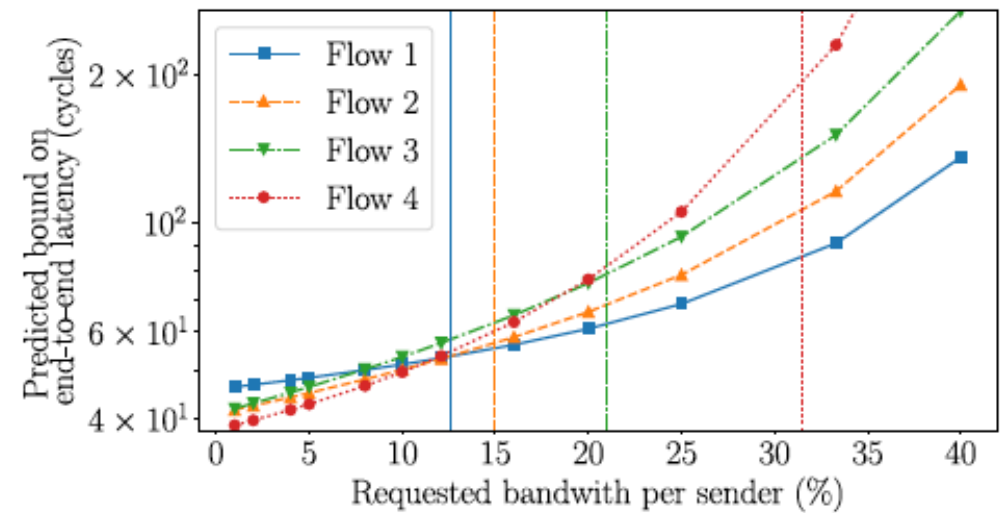
⇒ Safe delay bounds with our approach, in comparison to existing ST ones

# Performance Evaluation (2)

## Comparative analysis vs CPA

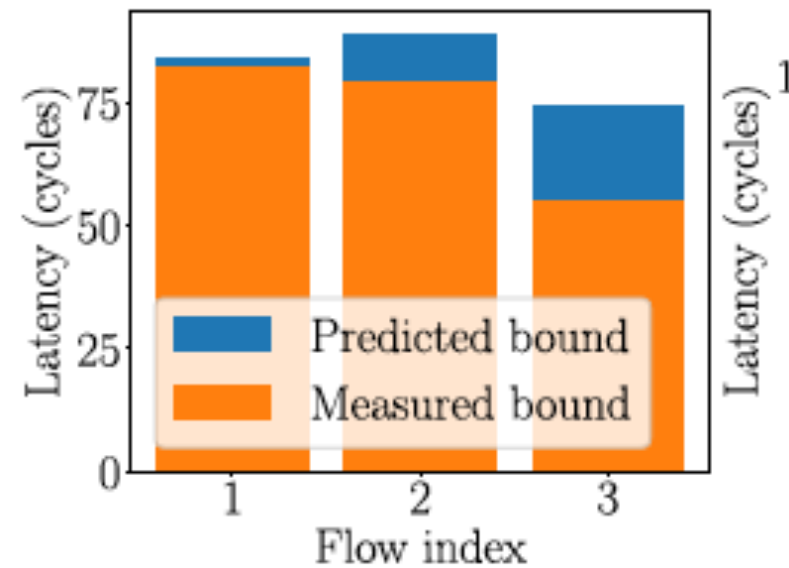
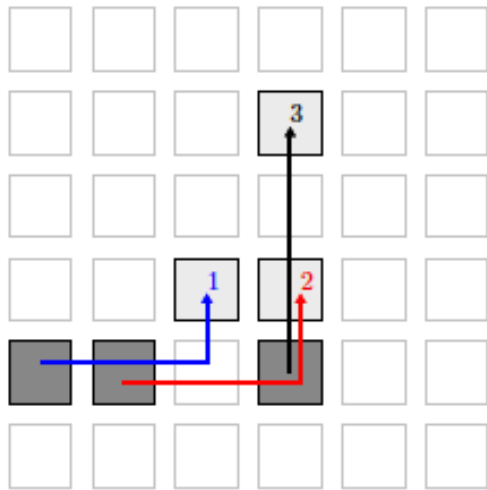


⇒ When **increasing** the network **congestion** or the **buffer size**, the delay bounds with our approach are **tighter**, in comparison to CPA



# Performance Evaluation (3)

## Experiments on a Physical Platform



⇒ The model tightness is high, with reference to experimental results

# Outline

- Context & Problematic
- Main Contributions
  - ✓ System Model and Assumptions
  - ✓ Main Steps of the Buffer-aware Timing Analysis
  - ✓ Performance Evaluation
- **Conclusions & Perspectives**

# Conclusions

## Proposed timing analysis of wormhole NoCs

- ✓ Covering a **large panel of NoCs** architectures
- ✓ Taking into-account the **buffer size** (backpressure)
- ✓ Taking into-account the **flows serialization** phenomena

## Results show:

- The **safety** of the obtained bounds, in comparison to **Scheduling Theory** approaches
- The **tightness** of the obtained bounds, in comparison to **CPA** and **experimental** results

# Perspectives

- To conduct a deeper **sensitivity analysis** of our model (network size, utilization rate, the buffer size, the flow burst and rate...)
- **Refining the model** when specifying a service policy between classes of the same VC and flows of the same class
- **Further experiments** with more complex congestion patterns to measure the tightness of our model



THANK

YOU

