

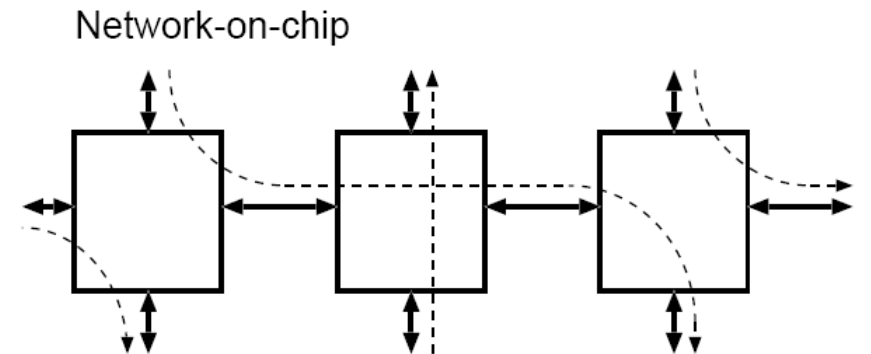
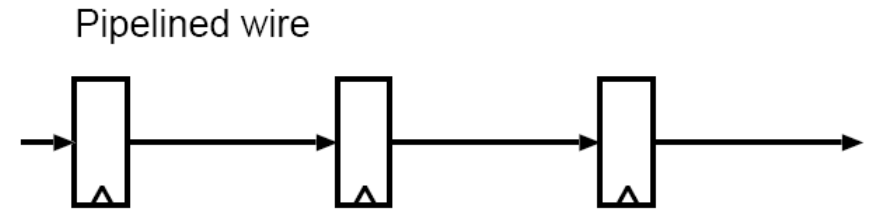
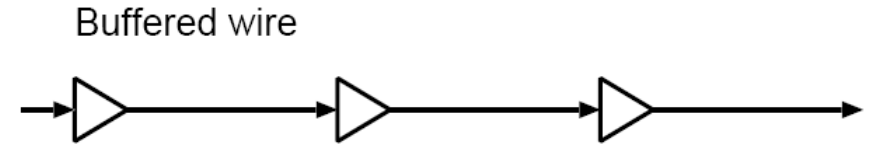
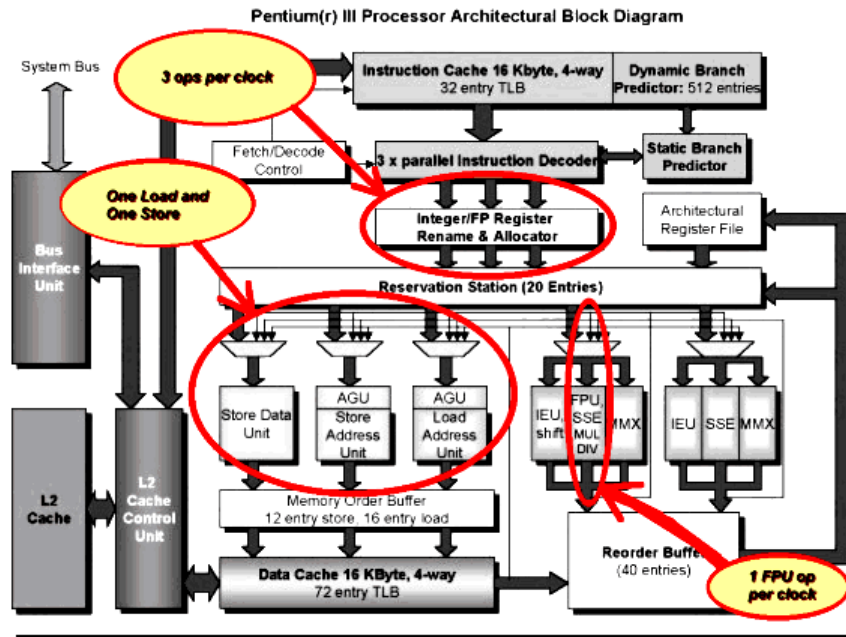
Foundations of On-chip Communication: Latency Issues in Multicore Platforms

Radu Marculescu

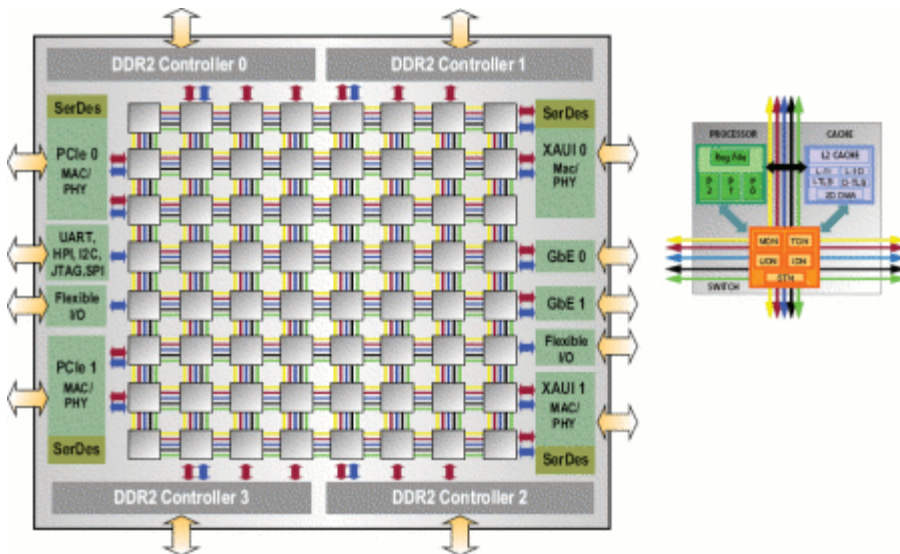
**Carnegie Mellon University
Pittsburgh, PA 15213, USA**

**NSF Wireless Workshop
March 26, 2015**

Low-power Design. Computation vs. communication



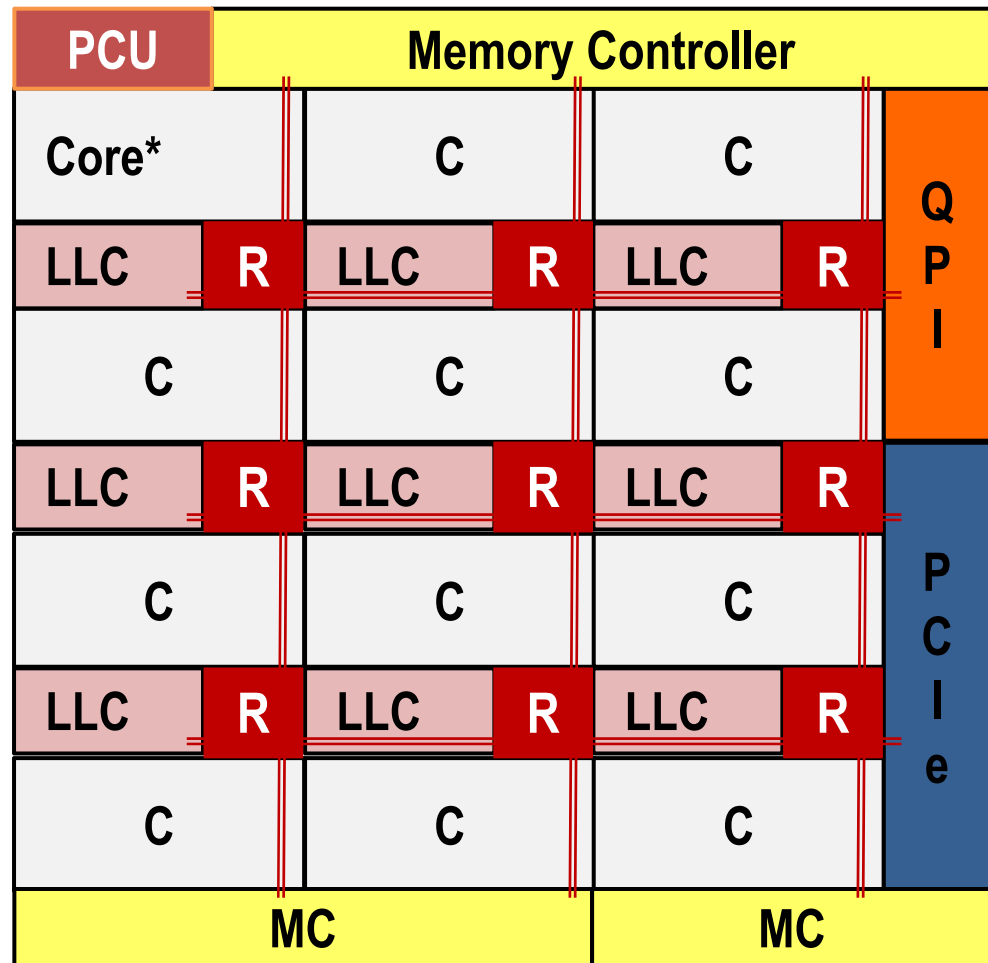
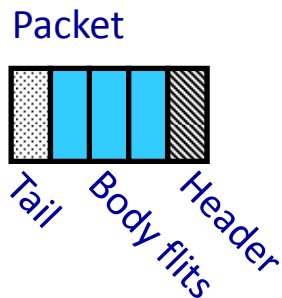
The Intel Pentium !!!



[D. Greenfield, S. Moore, Computer J., 2008]

Multicore platforms are large scale distributed systems at nanoscale; they are dominated by communication costs

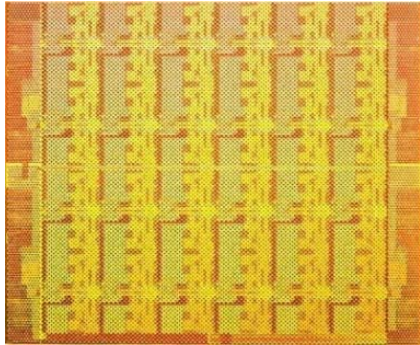
Last level cache (LLC)
 Memory controller (MC) & channels
 I/O controller(s)
 QPI controller,
 Power control unit (PCU),
 etc.



[U. Ogras, Tutorial ASPLOS 2012]

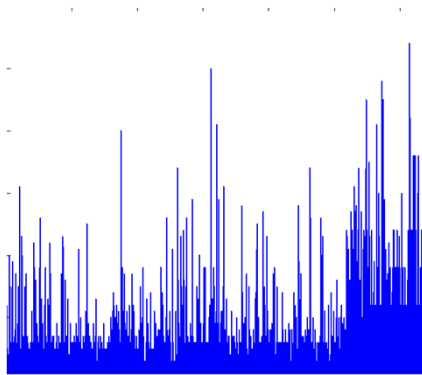
- 3 Need to understand the behavior of thousand core systems. Network (routers+links) is the missing link in understanding.

This presentation focuses on the new CPS paradigm and its impact on future integrated systems



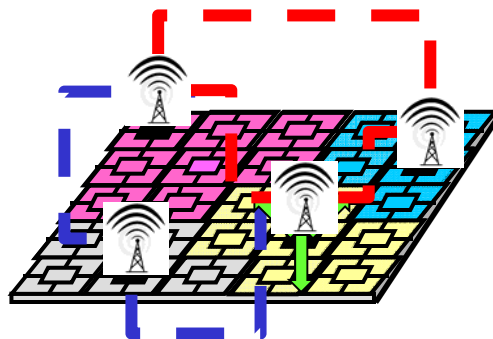
Structure

Architecture and small world effects



Dynamics

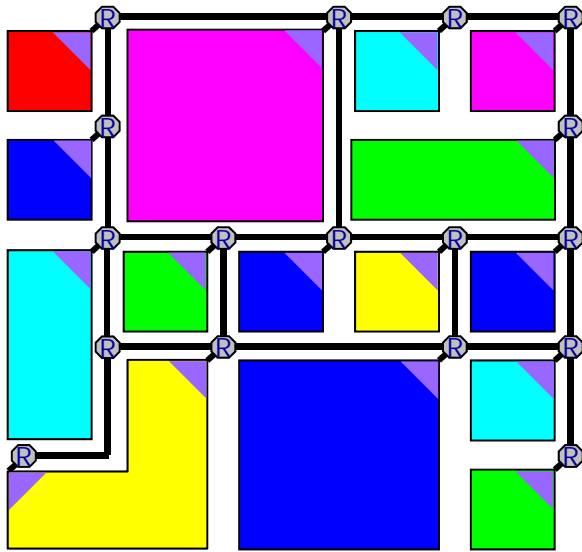
Workloads and multiscale behavior



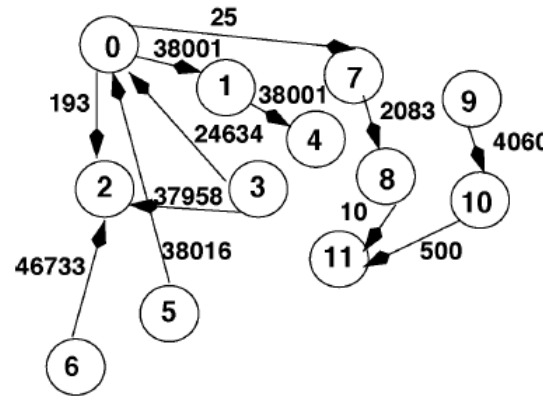
Control

Power and resource management

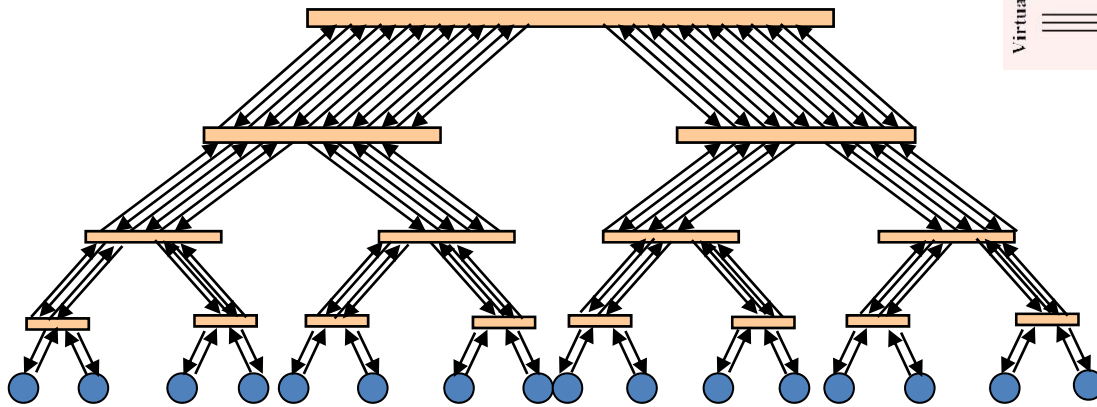
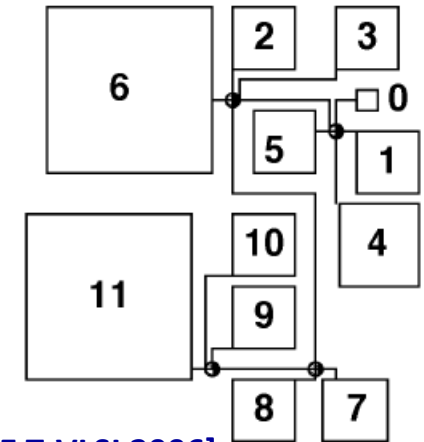
Our first insight into communication-based design came through architecture (topology, buffer, etc.) optimization



[Bolotin et al. DATE 2007]

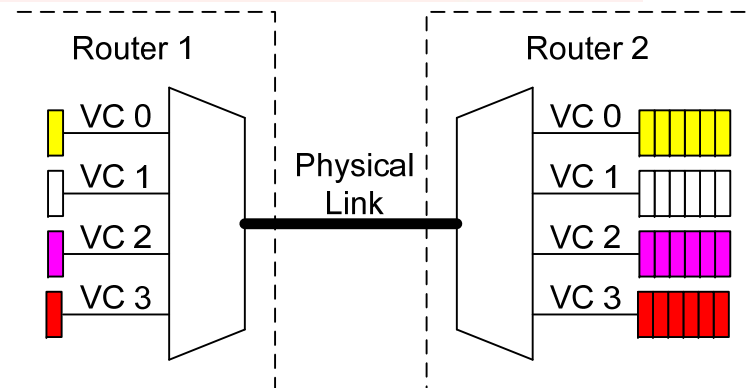
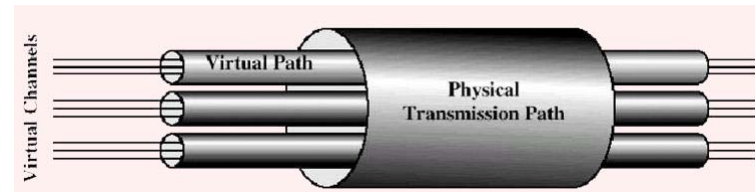


[K. Srinivasan et al. IEEE T-VLSI 2006]

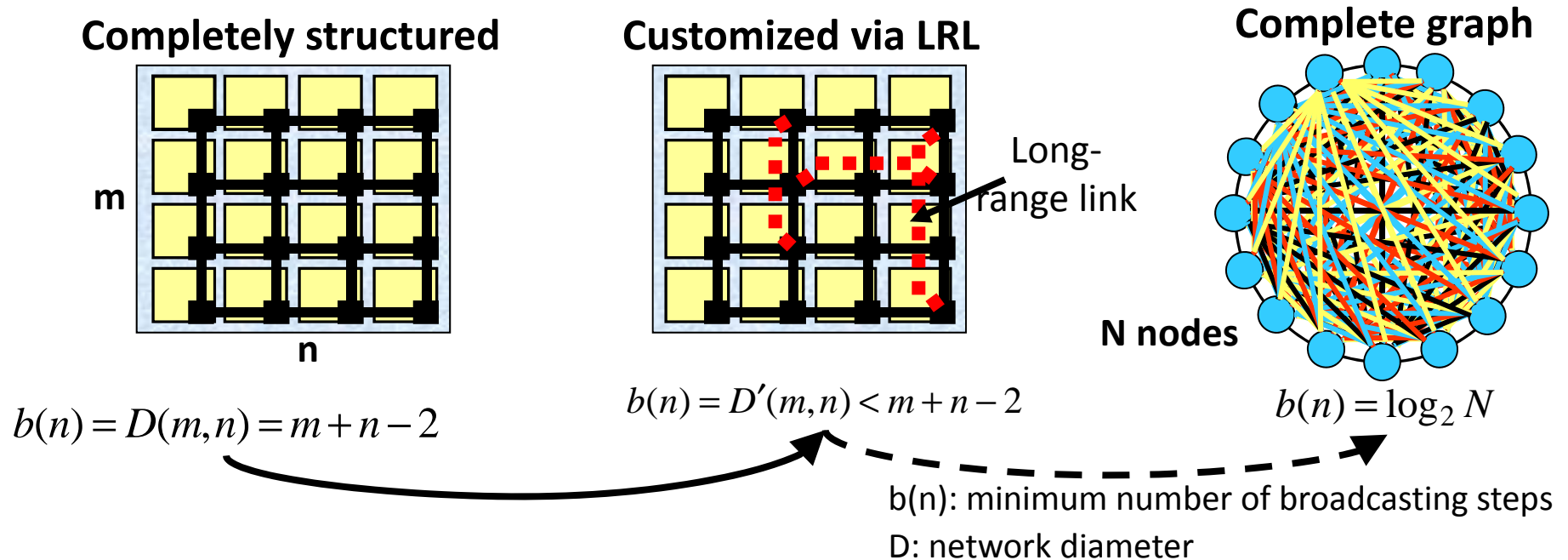


5

[J. Kim et al. ISCA 2007]

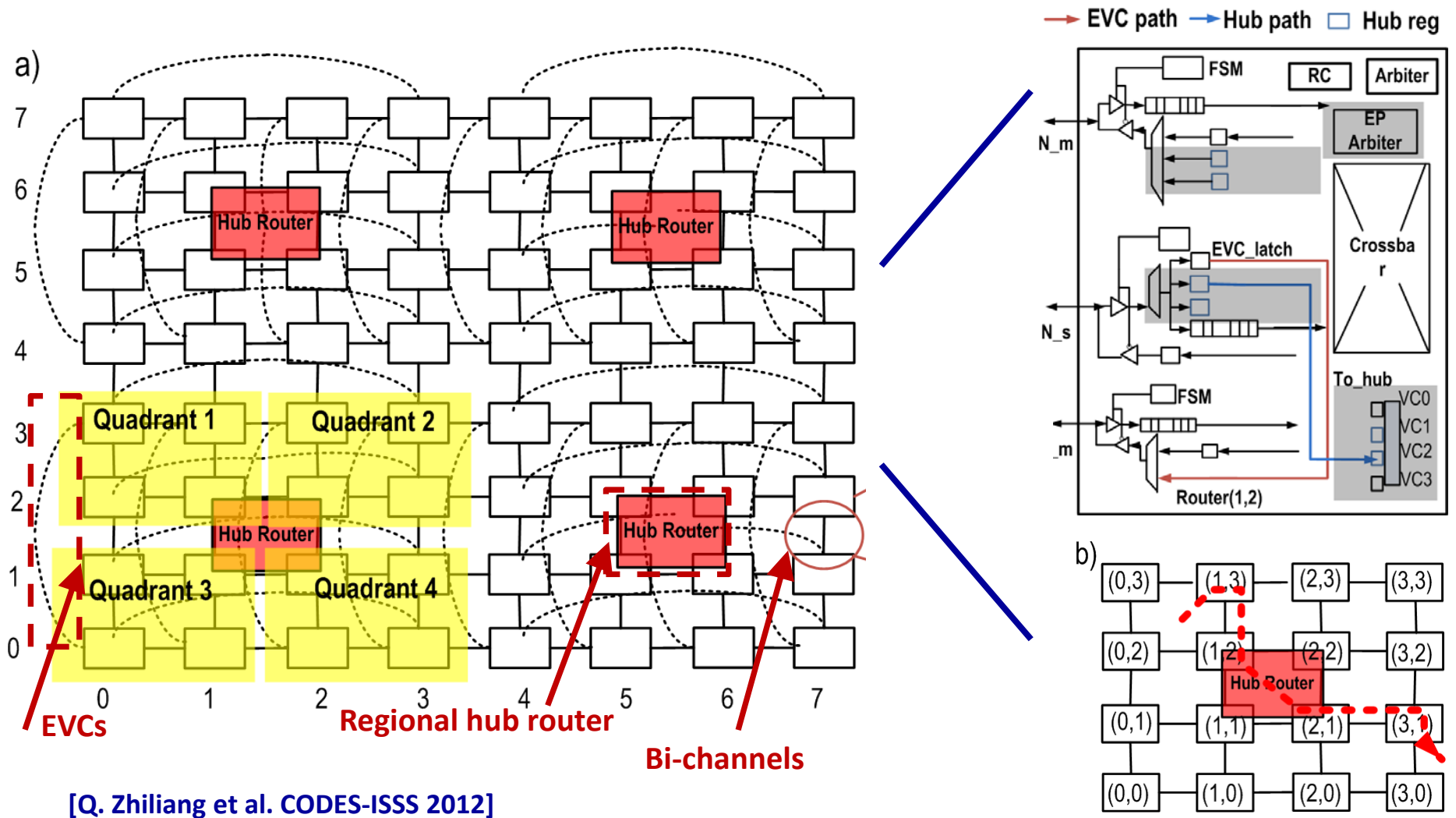


Can induce small-world effects in regular NoCs. This brings huge performance improvements



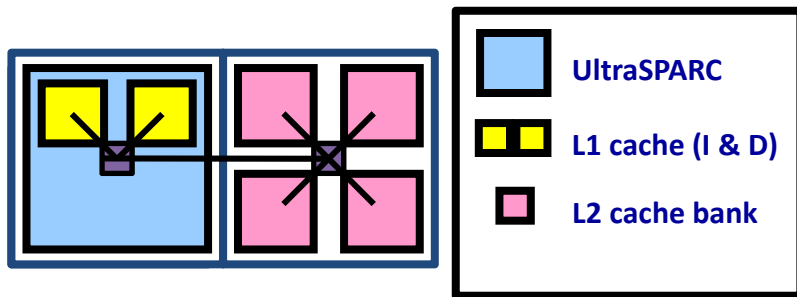
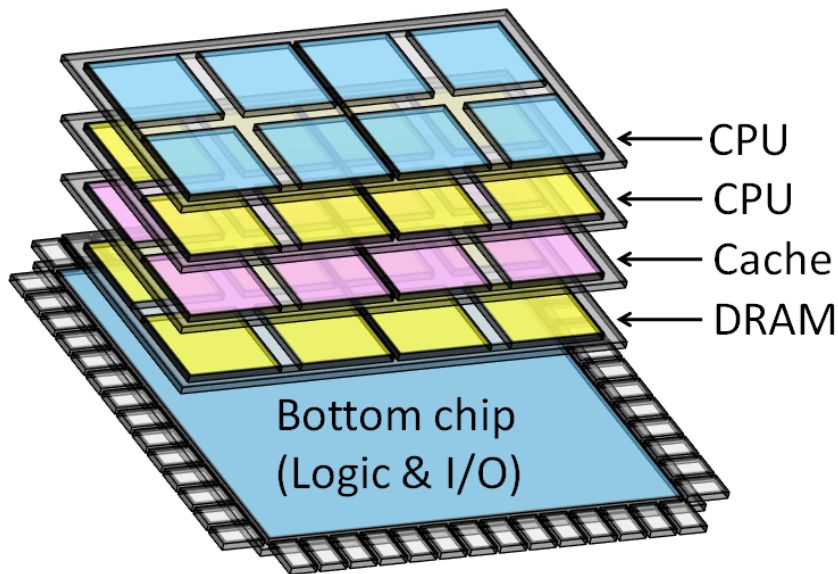
6 This way, the fundamental idea of Small World networks (aka “six degrees of separation”) enters the multicore world

Flow-control mechanisms, reconfigurability, adaptive routing have all been used to improve performance

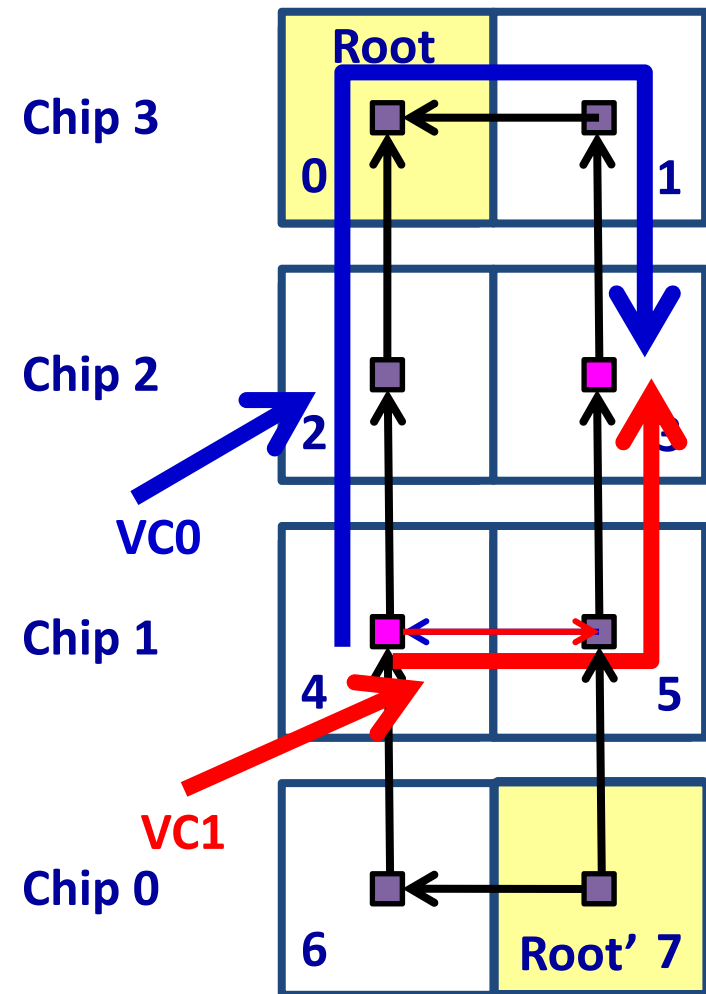


7 One way or another, they all exploit the small world effects...

Small world effects can also be exploited to reduce hop count in 3D wireless NoCs and improve performance

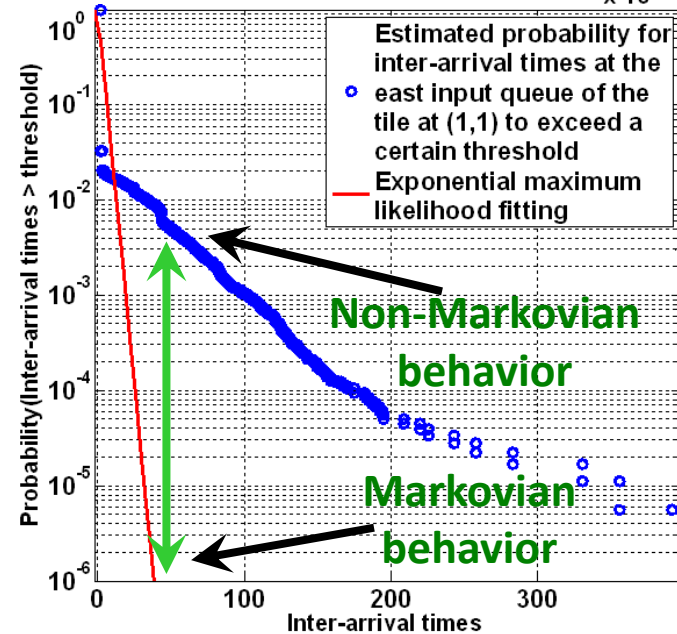
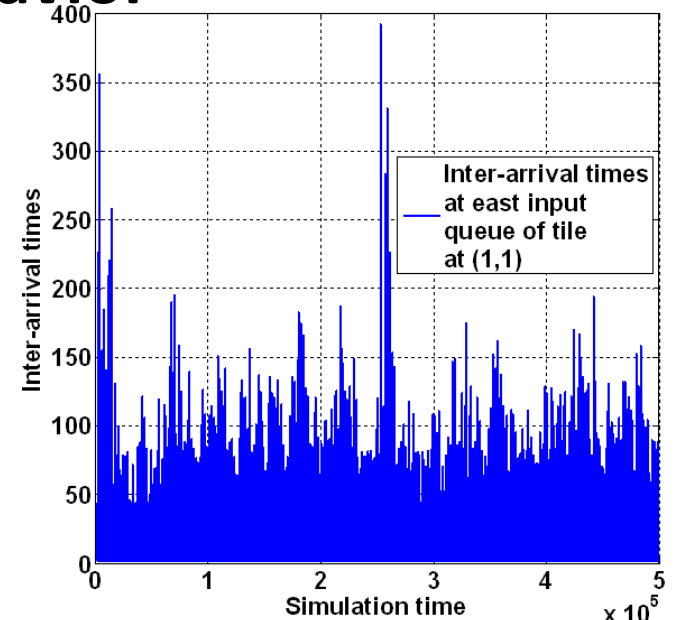
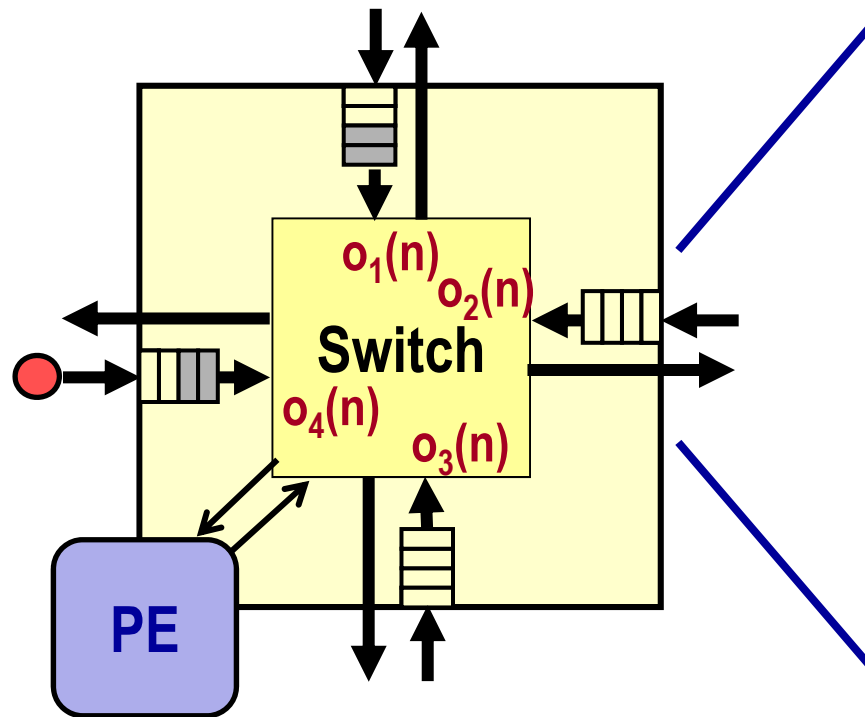


[H. Matsutani et al. ASPDAC 2013]

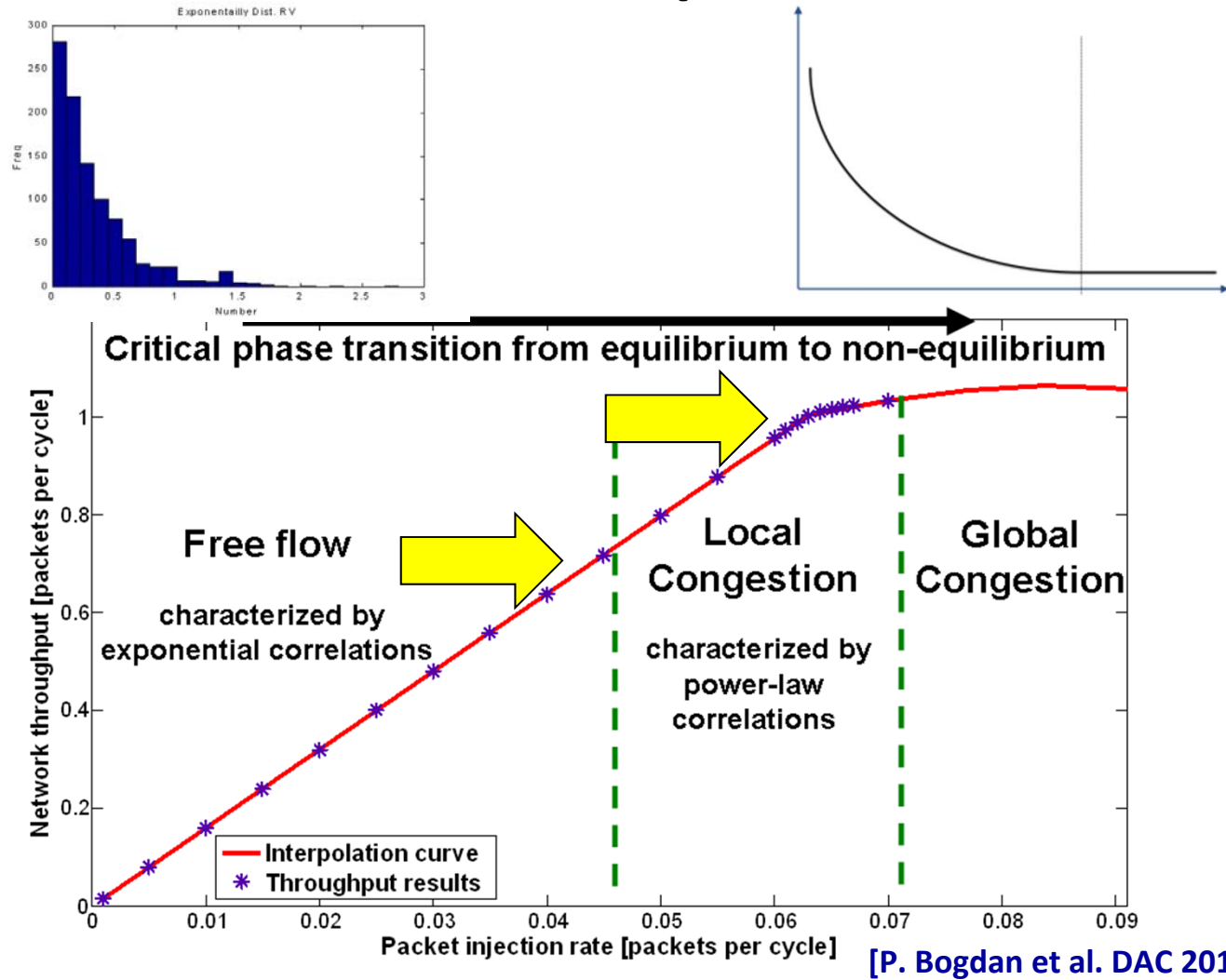


Wired and wireless NoCs can be used intra-chip, while inter-chip communication is based on wireless inductive-coupling.

Packet inter-arrival times at interface queues play a fundamental part in network behavior

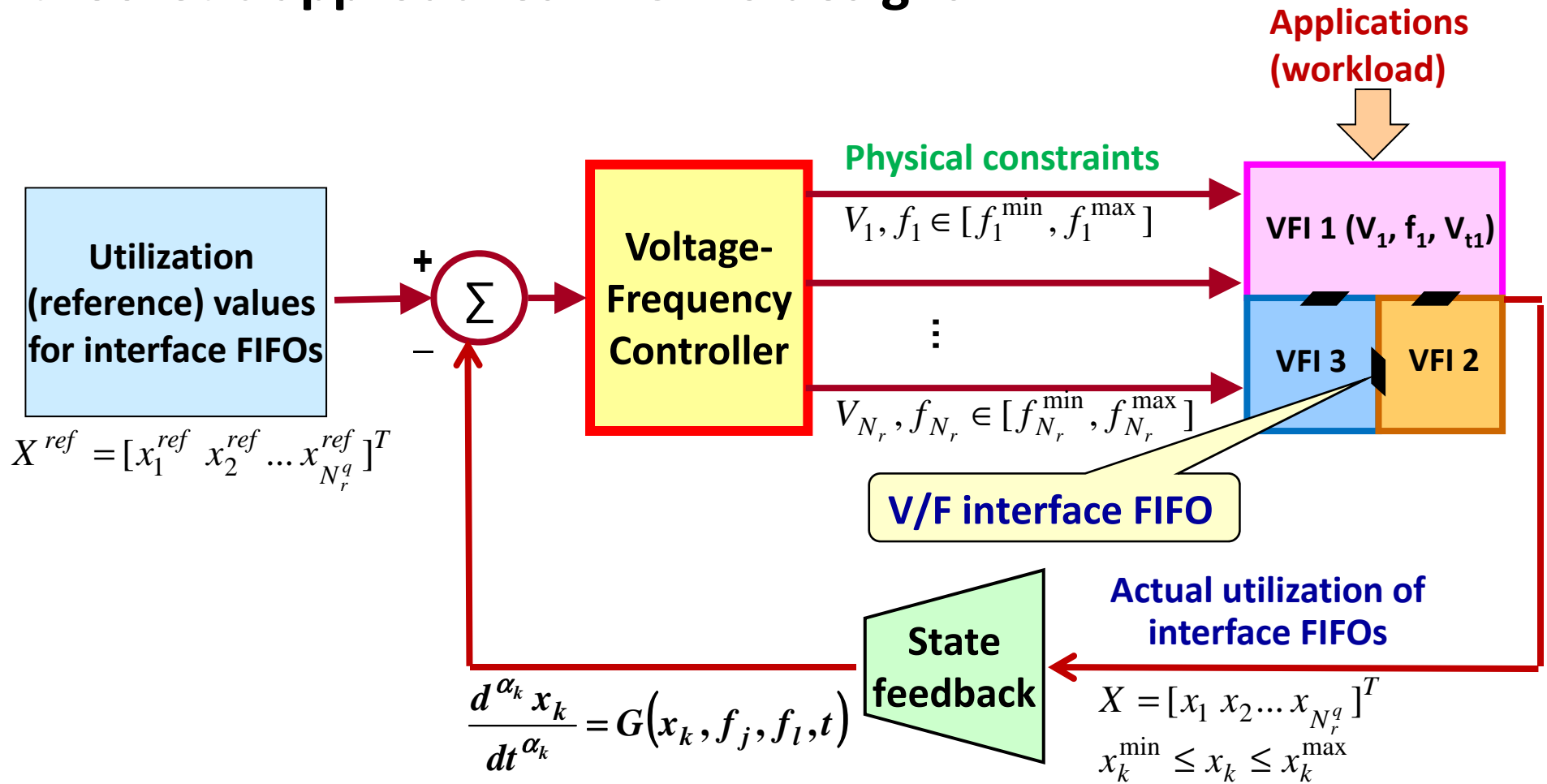


High injection rates cause inter-arrival times deviate from exponential distrib. and exhibit power law correlations.



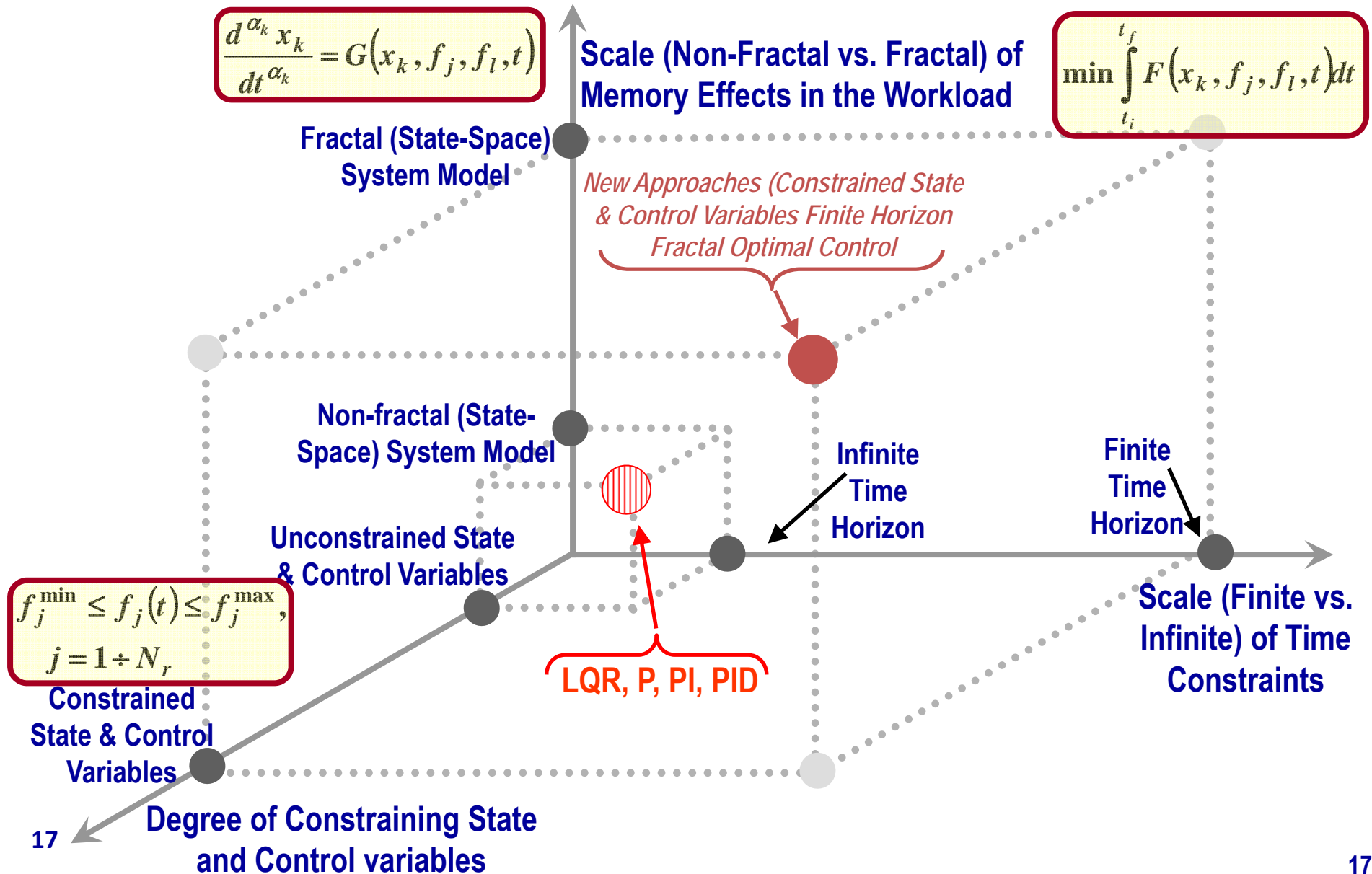
Workload analysis should not be an afterthought. In real platforms network traffic is neither Poisson, nor stationary.

Power management can be implemented via control-theoretic approaches in GALS designs

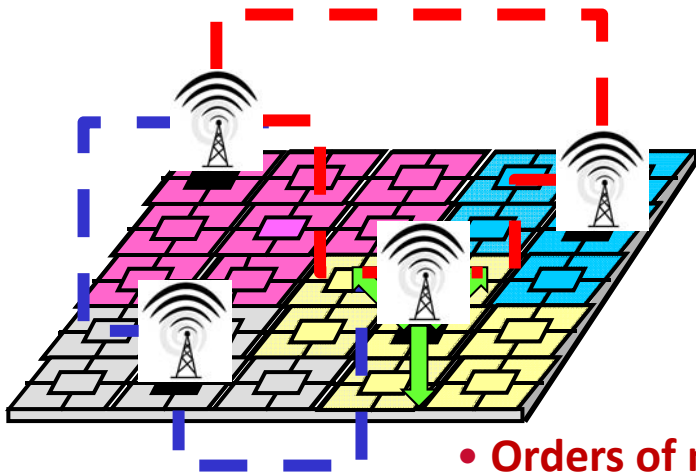


Fine-grain power management becomes possible by exploiting workload variations.

Need a new paradigm shift, i.e., incorporate correlation structure of traces into dynamical state equations



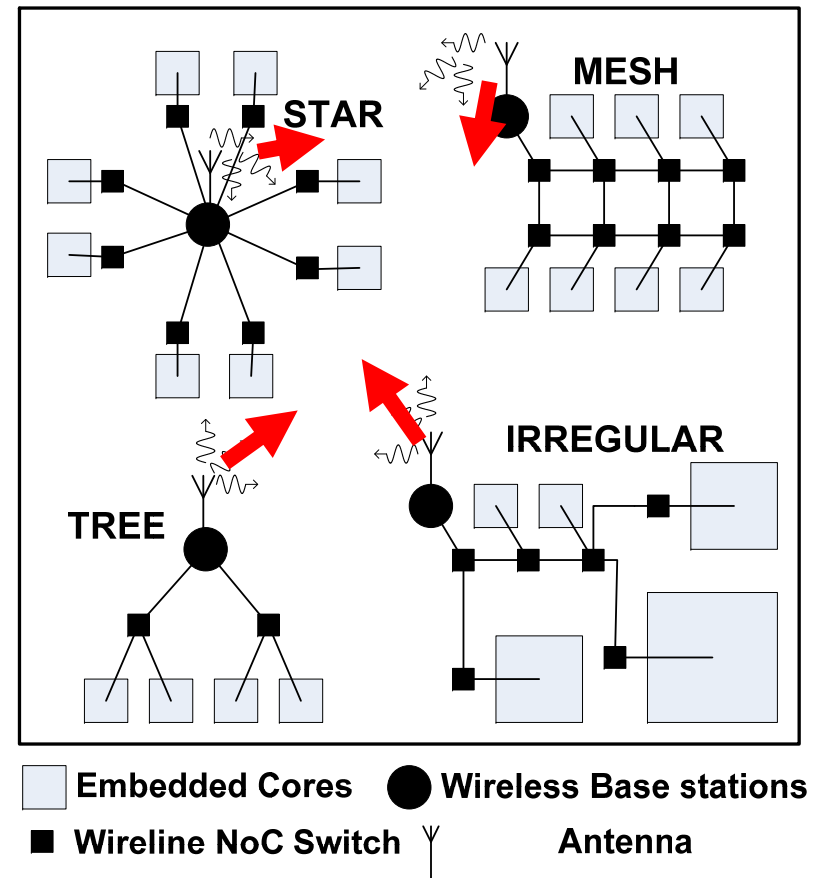
An hierarchy of globally distributed locally centralized control may help the system self-organize



- Orders of magnitude!
- Gets better with size

System Size	Flat Mesh (n)	WiNoC (n)	Factor
128	1319	22.57	58x
256	2936	24.02	122x
512	4992	37.48	133x

[Ganguly et al. IEEE Trans. Comp., 2010]



Local control w/ full state information, global control w/ partial information. Small world effects help convergence

Finally...

Contributors (in no particular order...)

Paul Bogdan (Univ. of Southern Calif.), Umit Y. Ogras (Intel/ASU), Partha Pande (Washington State Univ.), Diana Marculescu (Carnegie Mellon Univ), Qian Zhiliang (Hong Kong Univ Sci&Tech), Chi-Ying Tsui (Hong Kong Univ Sci&Tech), Hiroki Matsutani (Keio Univ).

Relevant papers - www.ece.cmu.edu/~sld

Sponsors



Intel Corporation