



Heterogeneous, Distributed and Scalable Cache-Coherent Interconnect

Scale system performance faster than Moore's Law will currently allow



MPSoC Conference 2016
Nara, Japan, July 13, 2016

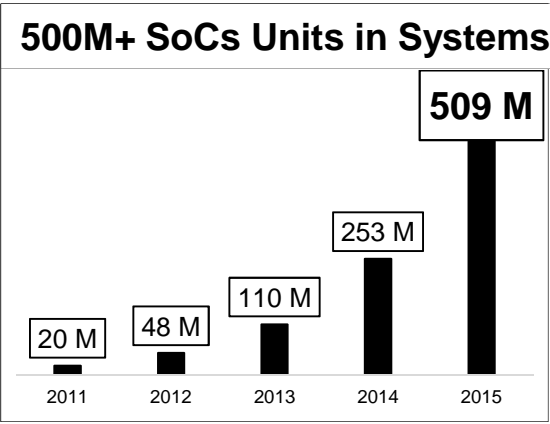
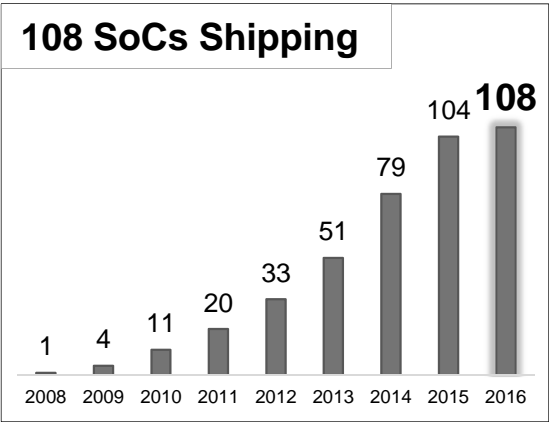
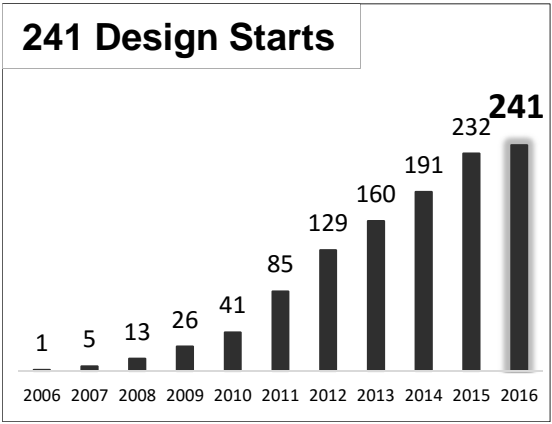


K. Charles Janac
President and CEO, Arteris

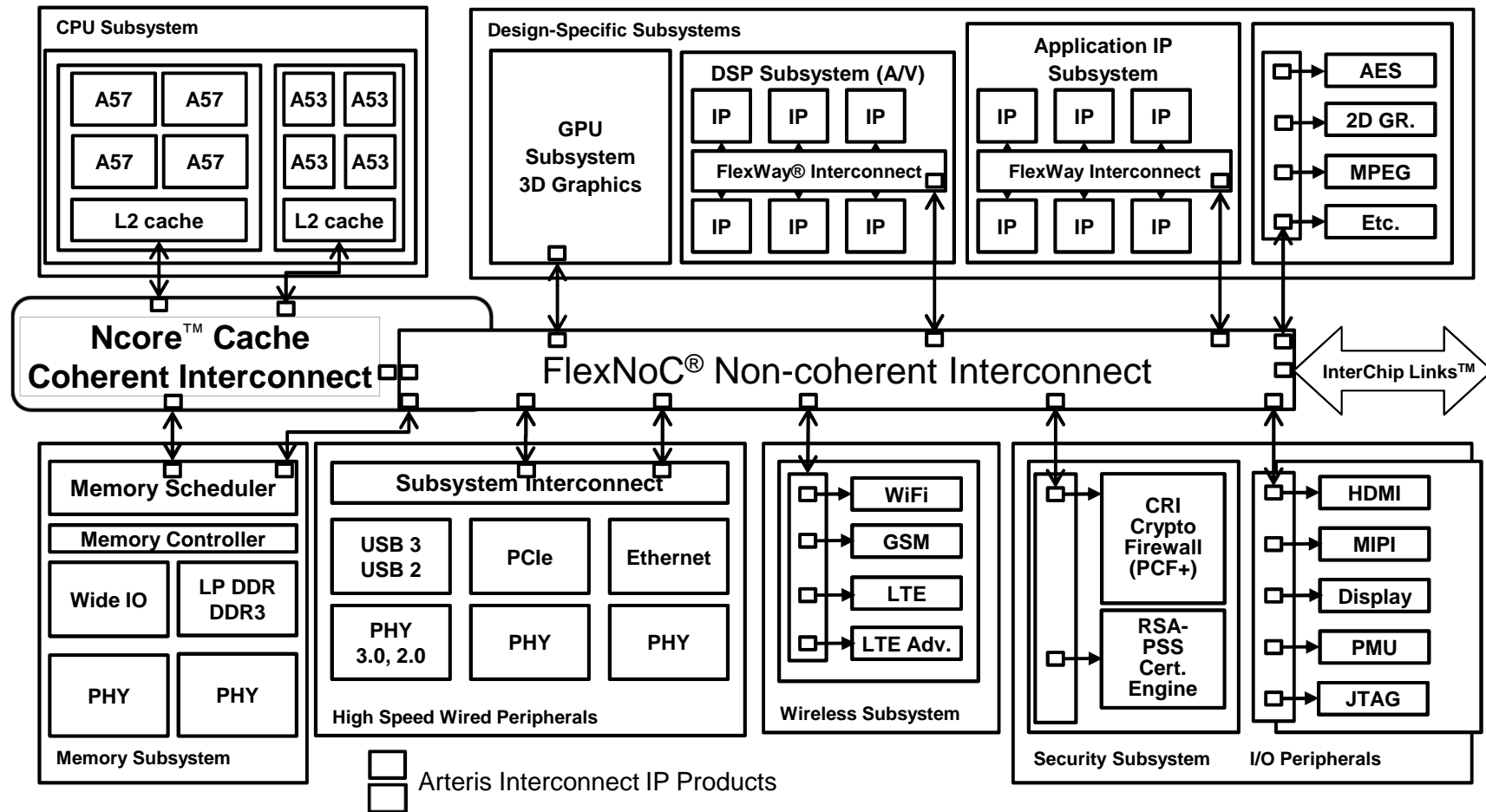
Arteris customers and their SoCs

Current as of 8 June 2016

					Very Large SoC Maker		
	Major Automotive OEM	Major Auto & CE SoC Maker					
		Major System OEM		Japan System OEM	Automotive SoC Maker		
Japan Tier 1 SoC Maker					Large Drone Maker		
Major SSD Vendor						Major IP Provider	
Defense Contractor	Defense Contractor	Major SSD Vendor	Defense Contractor	Silicon Foundry			

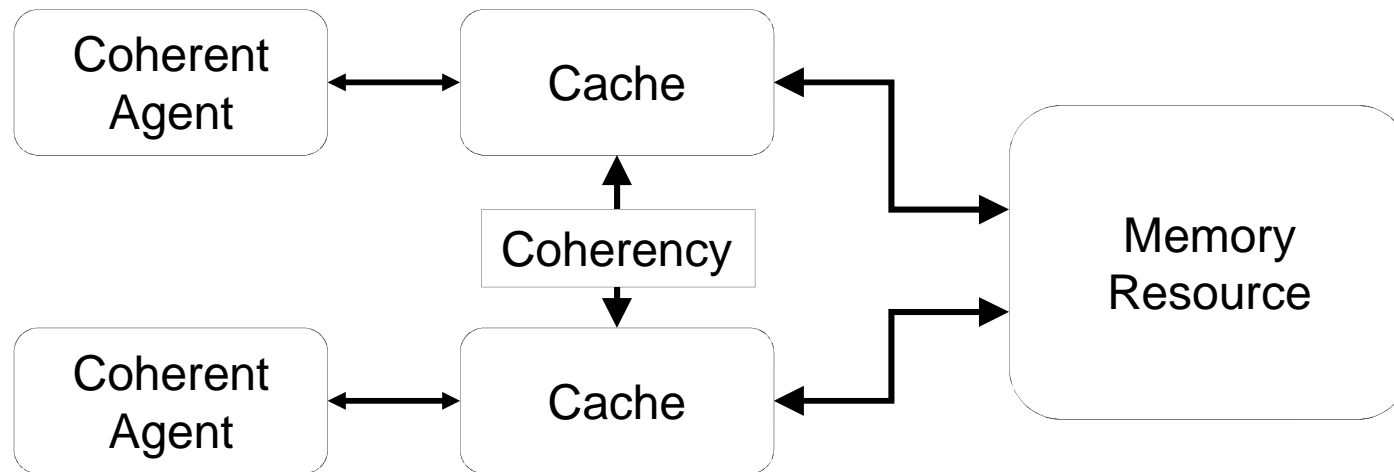


Arteris Interconnect IP: The easiest and most cost effective way to build differentiated SoCs



Cache Coherency Primer

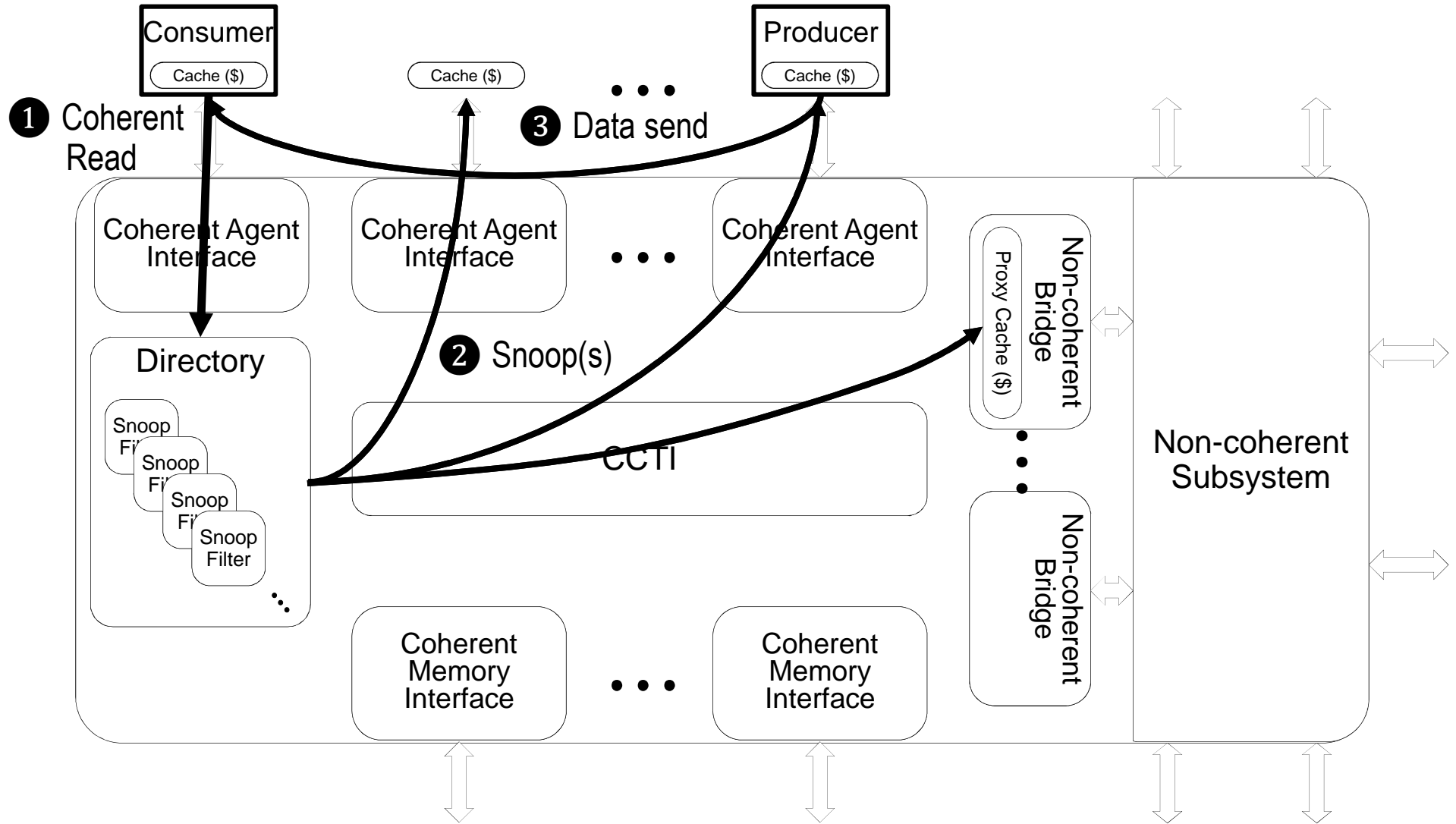
Cache coherency concept



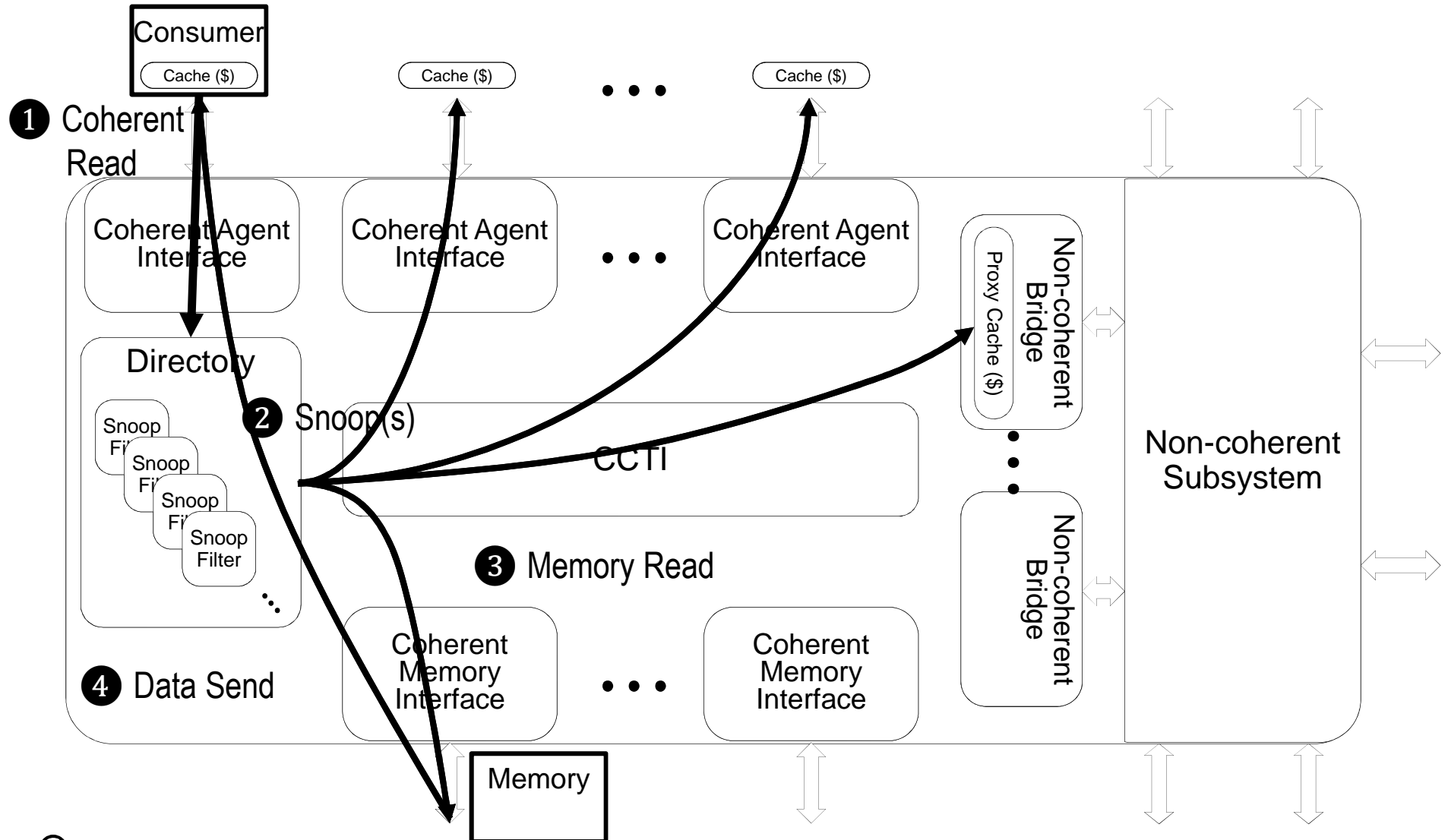
○ Source: Wikipedia



Coherent read example – cache hit

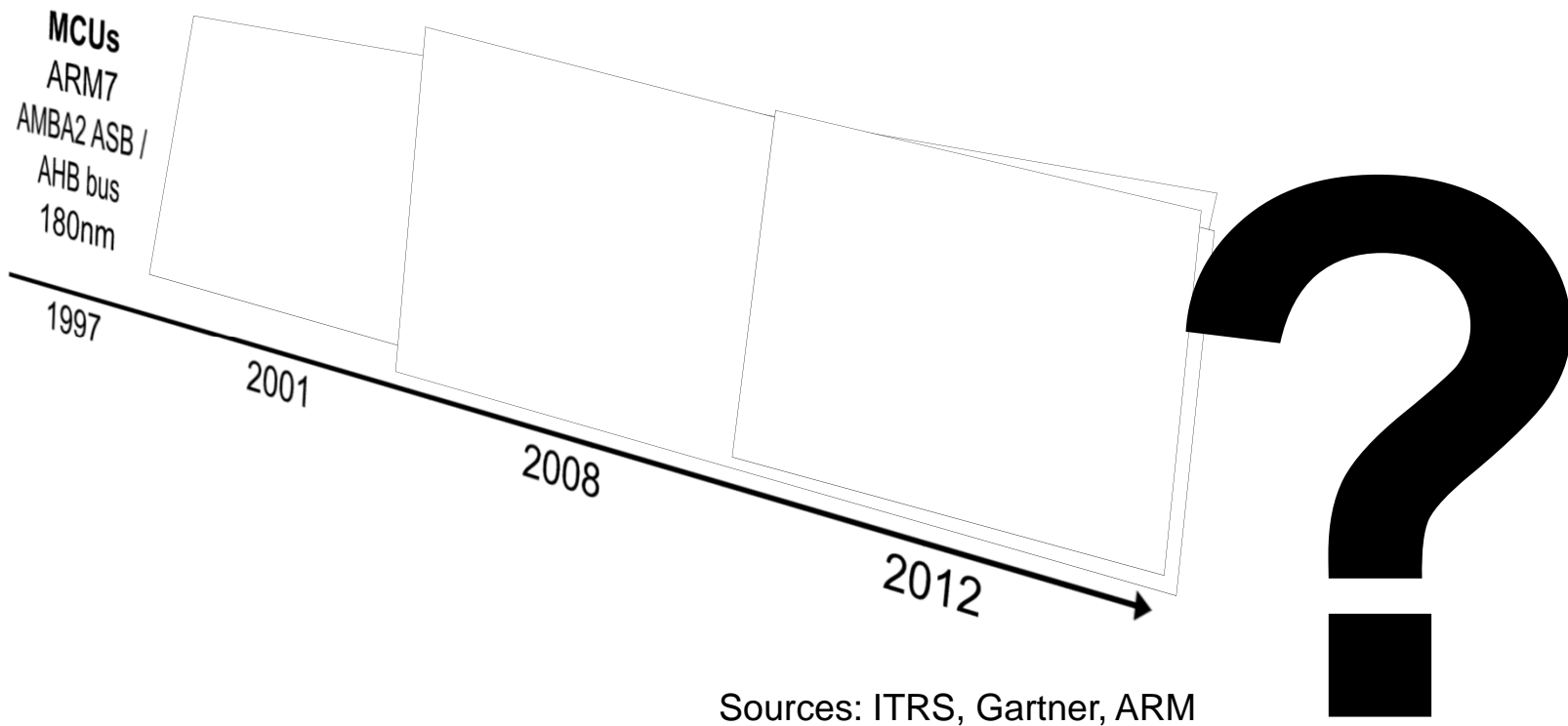


Coherent read example – cache miss



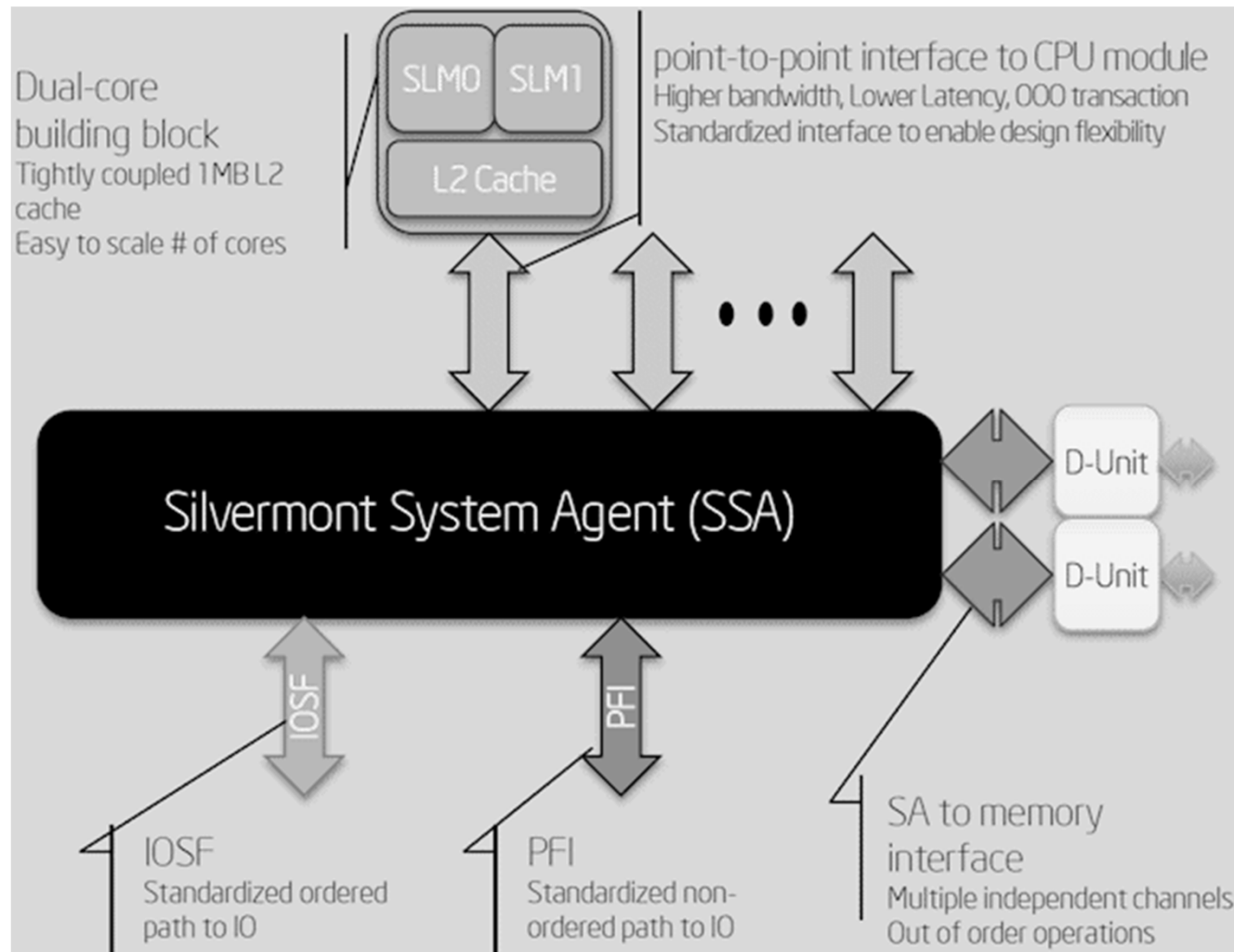
Who needs heterogeneous
cache coherency?

How did we get here?



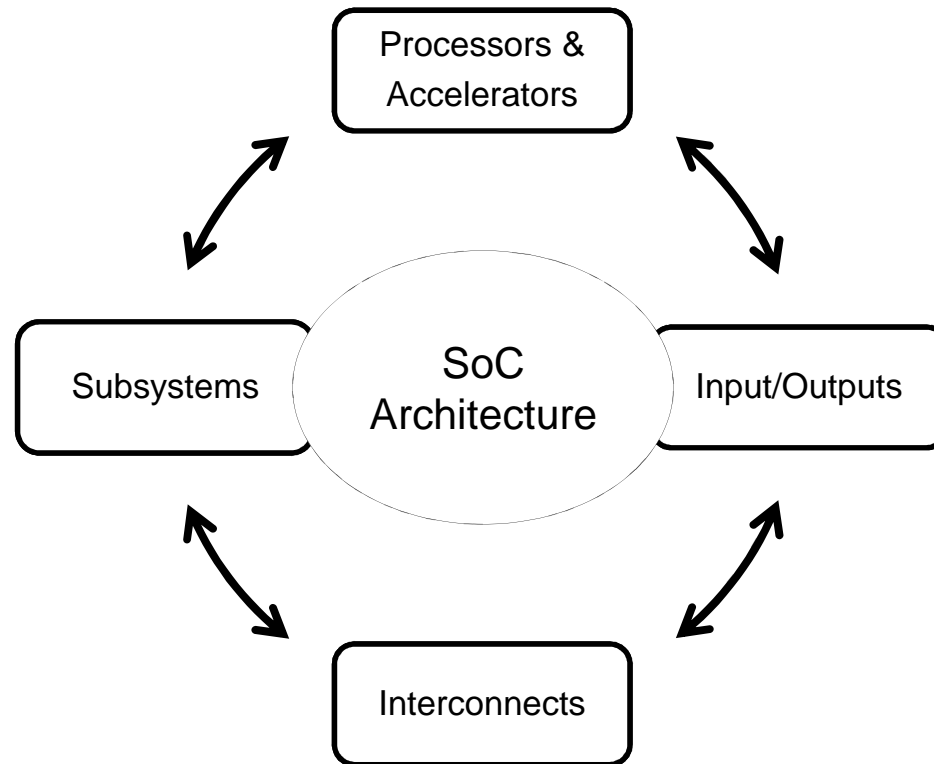
We can't improve performance and economics by transistor scaling. No more free lunch!

Cache coherency has been available



Keep this in mind...

Software Simplicity

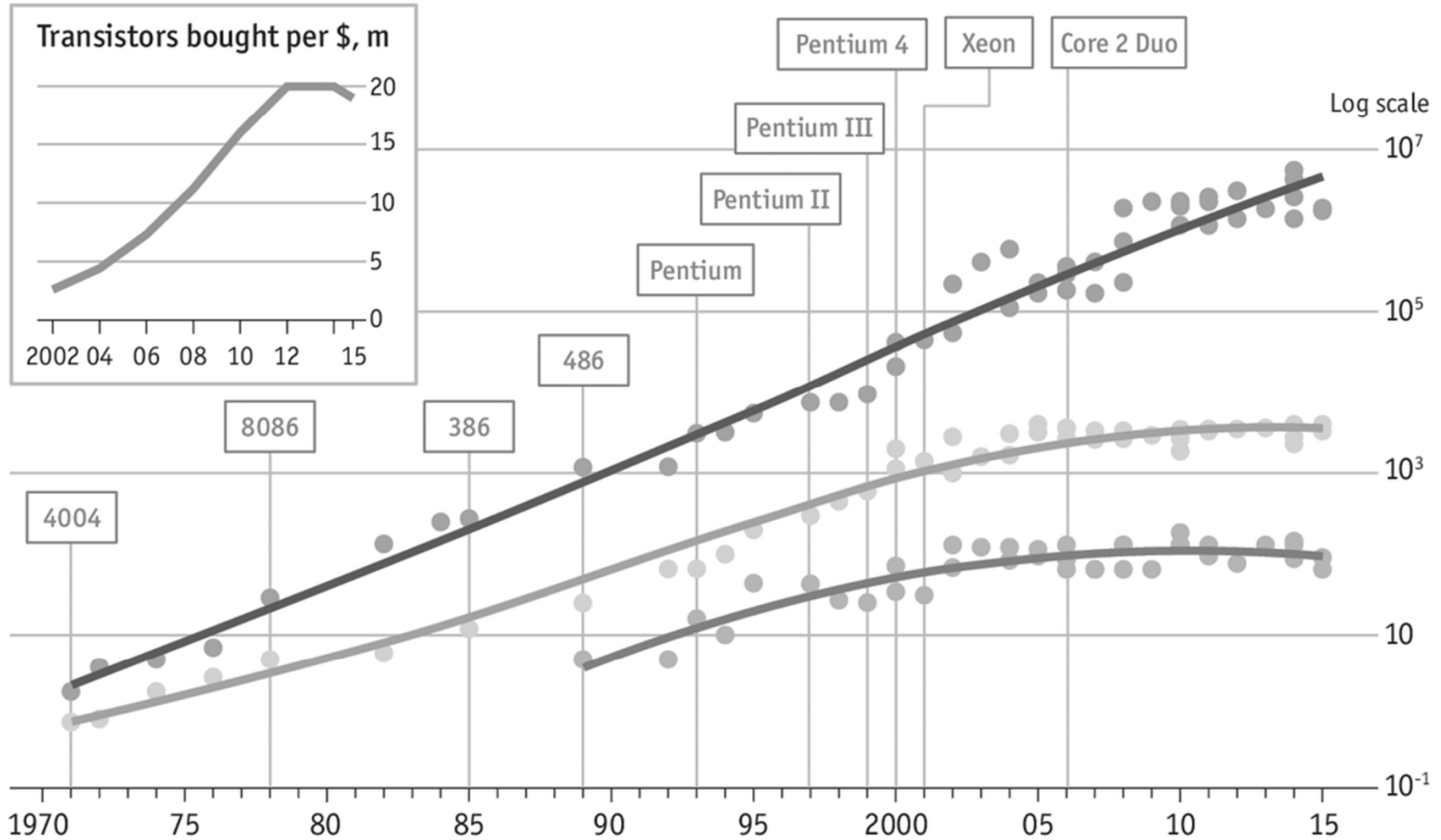


Moore's Law

Stuttering

● Transistors per chip, '000 ● Clock speed (max), MHz ● Thermal design power*, w

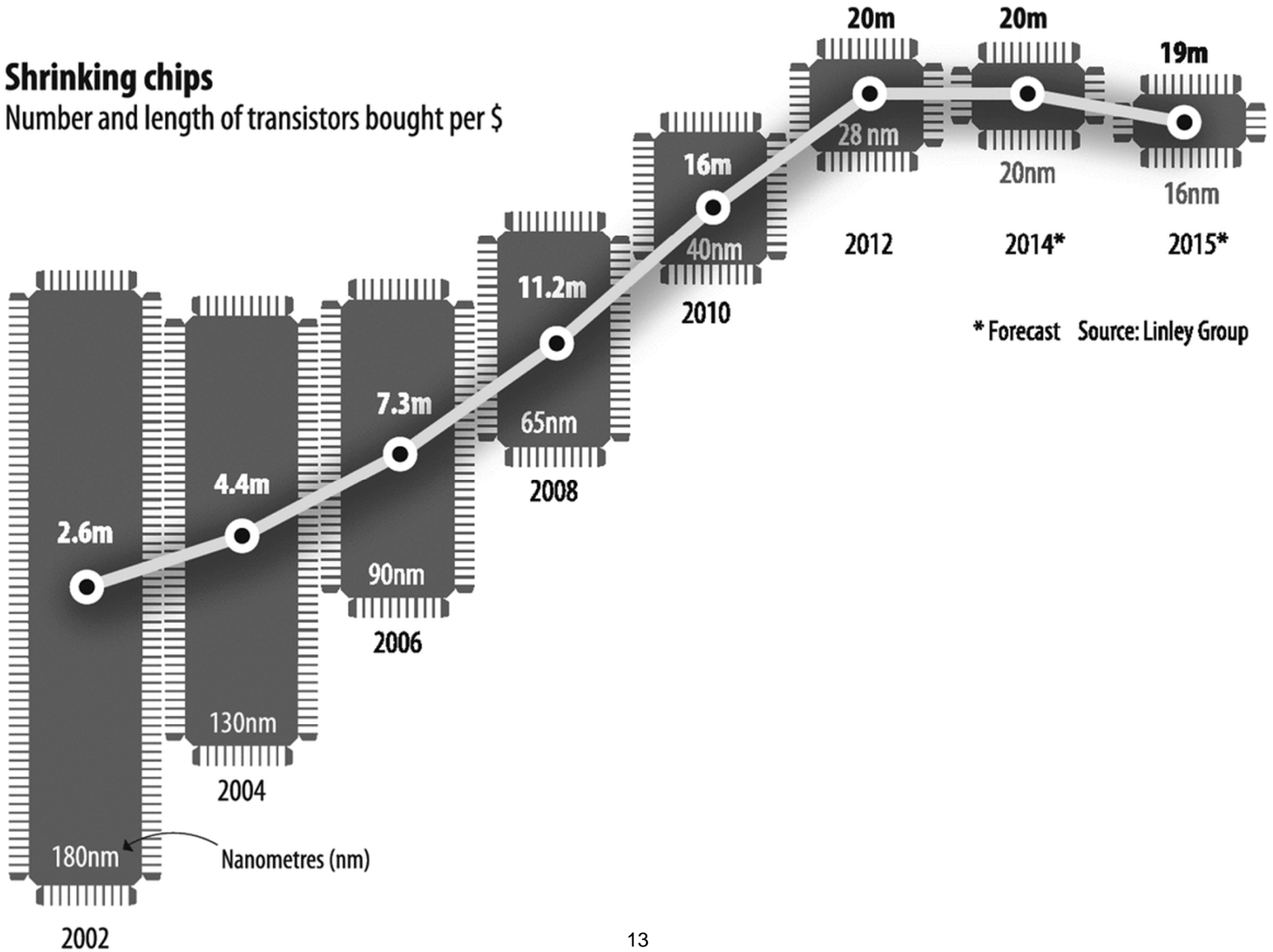
□ Chip introduction dates, selected



Sources: Intel; Bob Colwell; Linley Group; International Business Strategies; *The Economist* *Maximum safe power consumption

Shrinking chips

Number and length of transistors bought per \$



We need parallel computing to squeeze the most out of a process node

<h2>System</h2>	<ul style="list-style-type: none">• HW + SW efficiency• Acceleration / CPU offload• Bandwidth & latency efficiency
<h2>Software</h2>	<ul style="list-style-type: none">• Simplicity: Single view of memory• Reuse existing SW• “Software yield”: Useful work per LOC
<h2>Hardware</h2>	<ul style="list-style-type: none">• Optimize for dissimilar HW – data sharing• Efficiency: Useful work per area (or mAh)

But we need to address **systemic** complexity in addition to scale complexity

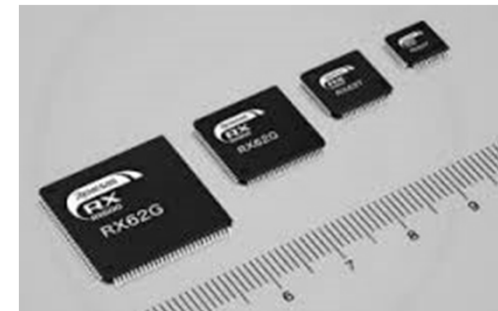
Heterogeneous cache coherency enables simpler parallel computing

Multiple Coherence Models	Dissimilar Caching Agents	Non-coherent IP as coherent peers
<ul style="list-style-type: none">• IP from different vendors / teams• Logical – coherence protocols / models• Generic system coherence model using lightweight messaging layer	<ul style="list-style-type: none">• Physical – cache organization, transaction table sizes• Multiple configurable snoop filters save area	<ul style="list-style-type: none">• Data sharing between non-coherent agents and coherent agents• Data sharing between non-coherent agents• Multiple configurable proxy caches minimize communication through DRAM



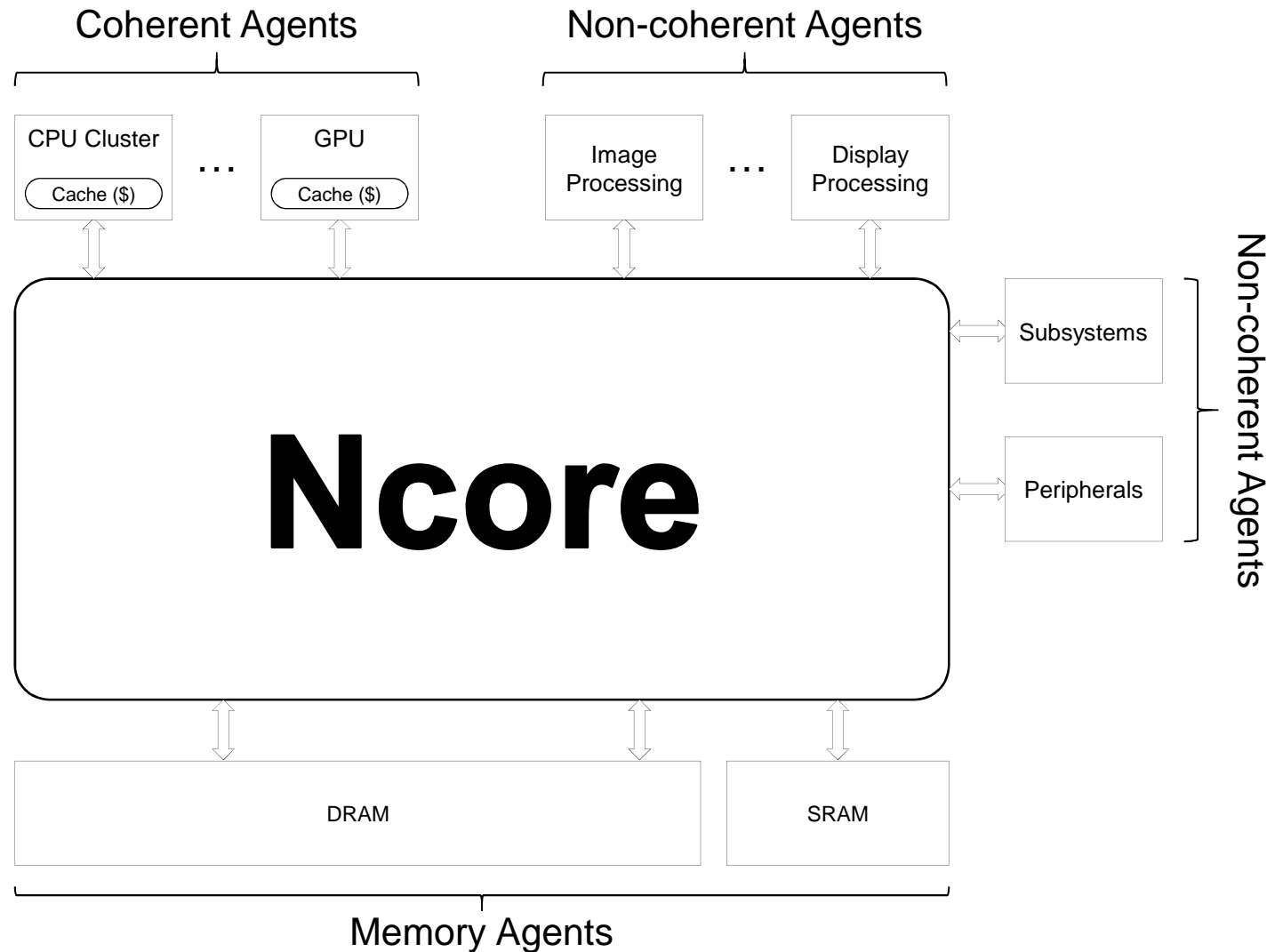
Where is heterogeneous cache coherency relevant?

- Next generation automotive (ADAS)
- Mobility applications processors
- Virtual reality
- IoT
- Servers and others

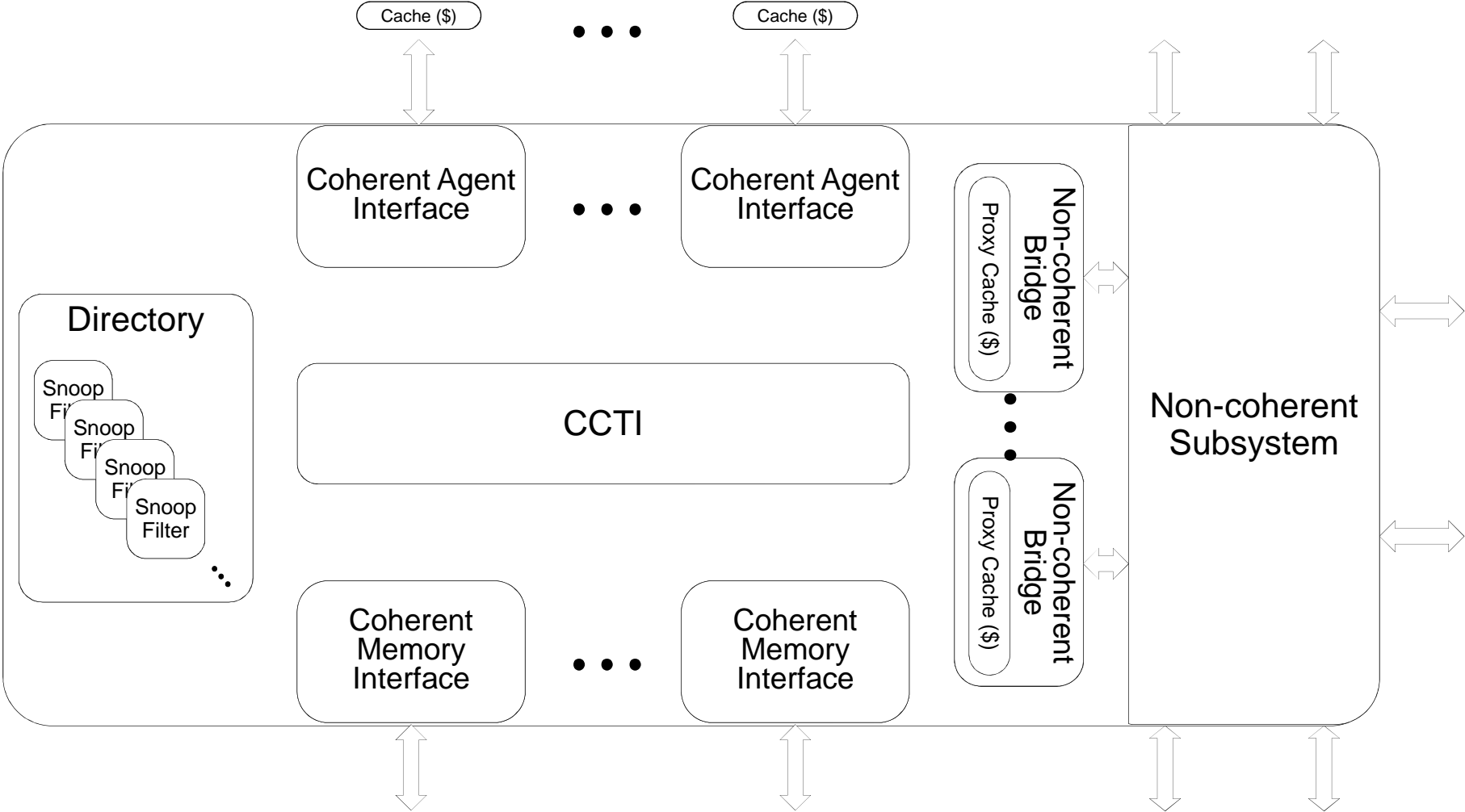


Heterogeneous Cache Coherency Implementation

Ncore Cache Coherent Interconnect IP

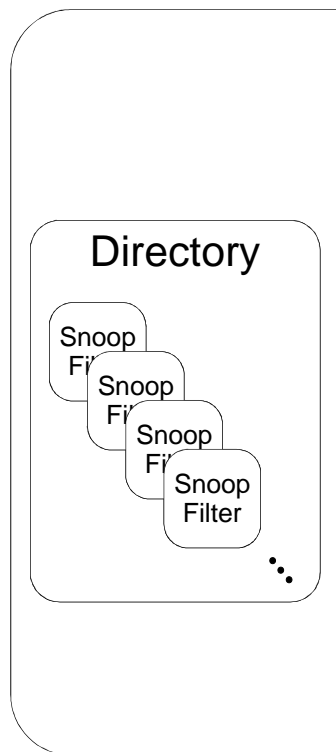


Ncore interconnect architecture



True heterogeneous coherency

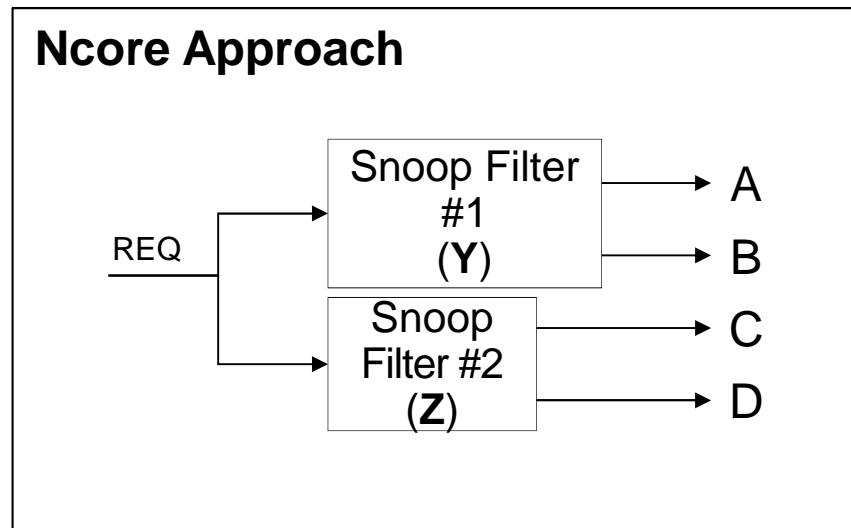
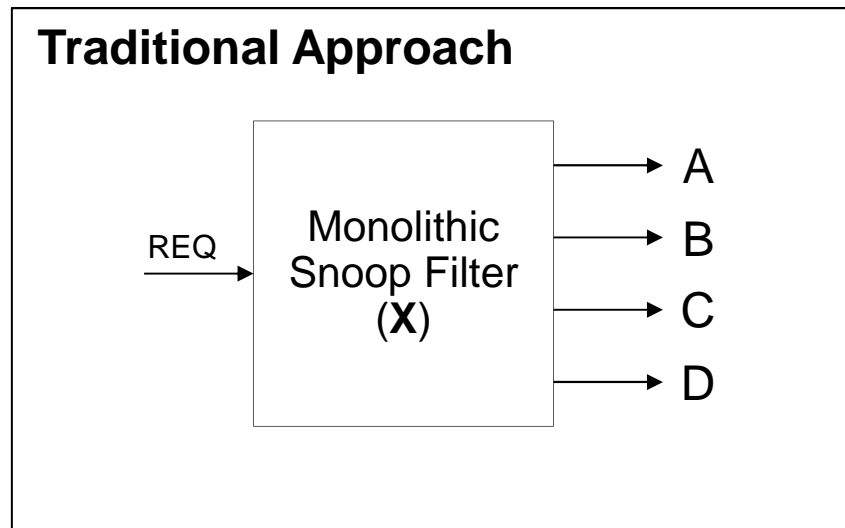
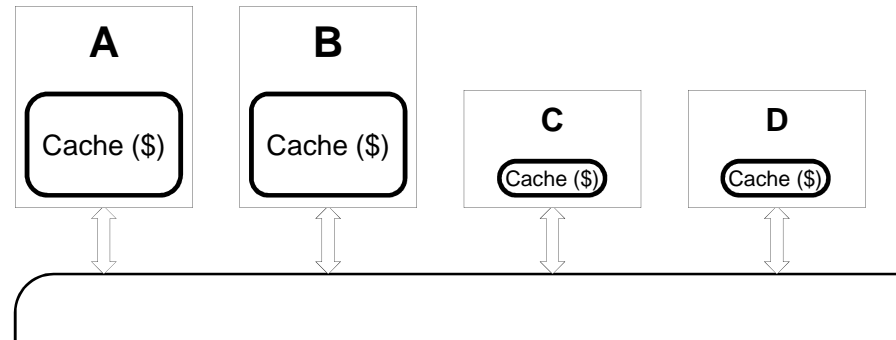
With multiple configurable snoop filters



- Cache coherent agents can have very different behaviors
 - Cache organization
 - Coherency models
 - Workloads

- Associating caching agents that share common properties with individual snoop filters can consume **less die area** than a monolithic snoop filter

Multiple snoop filters are more area-efficient than one

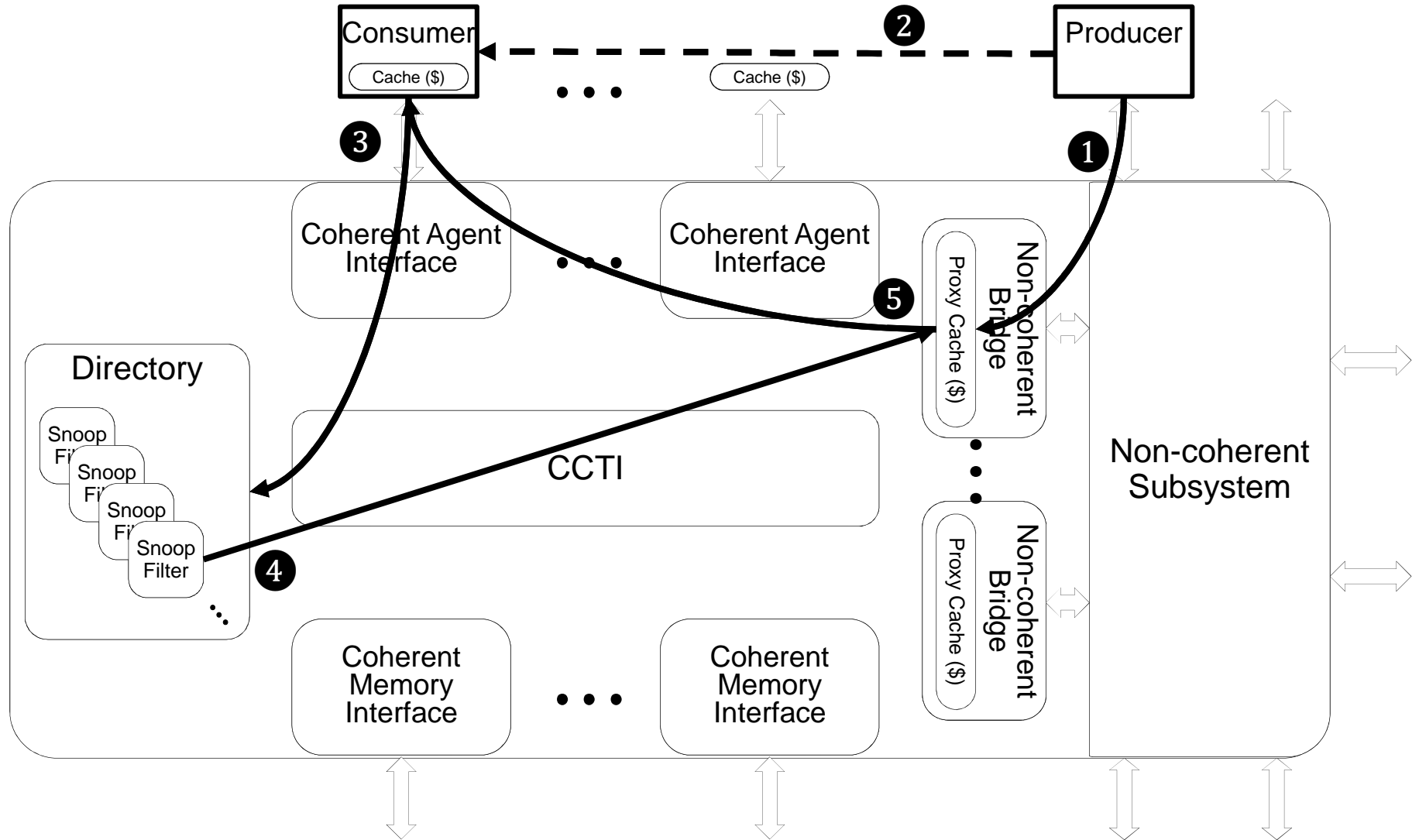


Multiple snoop filters are smaller: $area(Y+Z) < area(X)$

Higher performance with non-coherent IP

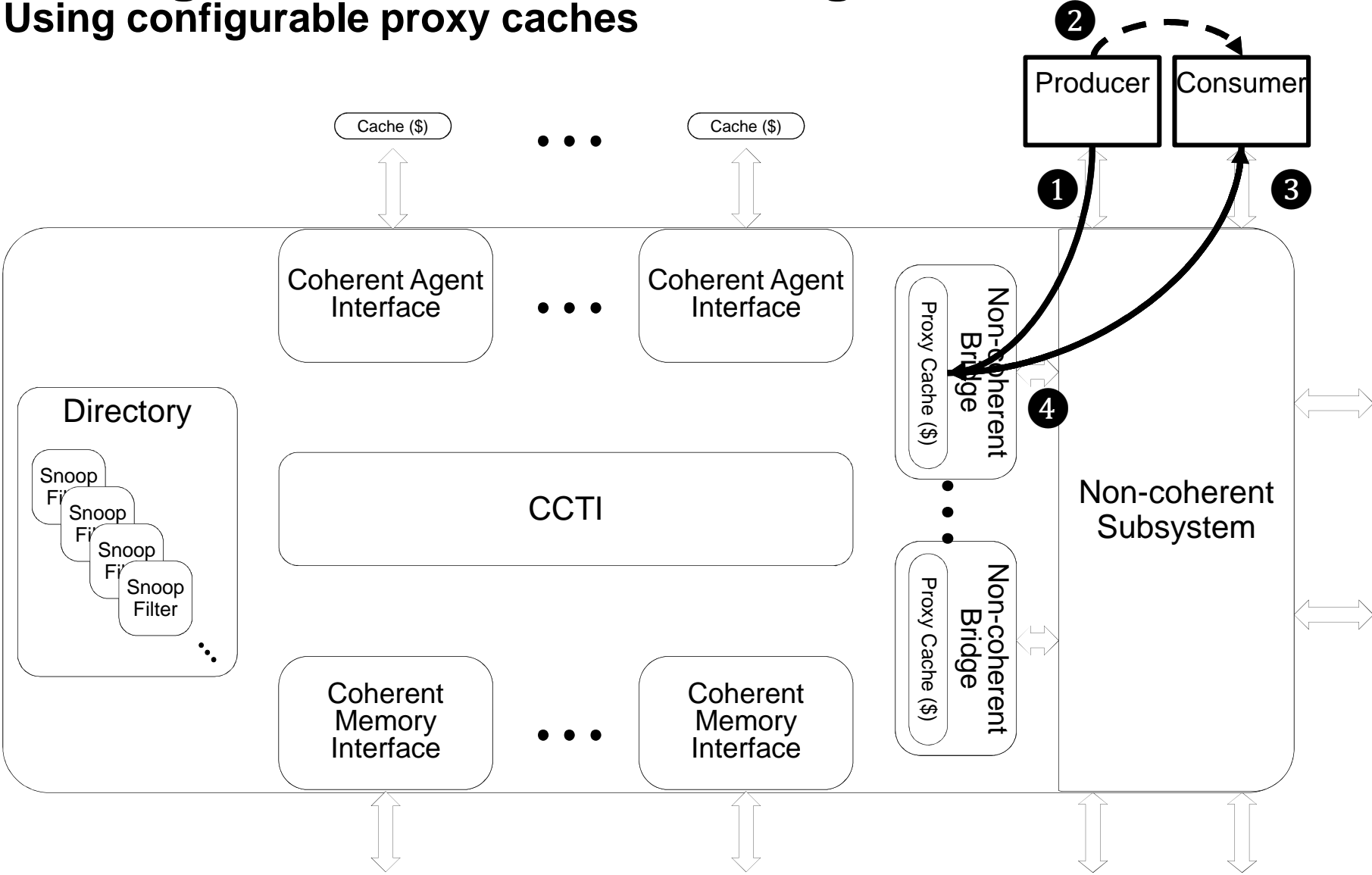
Sharing between non-coherent & coherent agents

Using configurable proxy caches

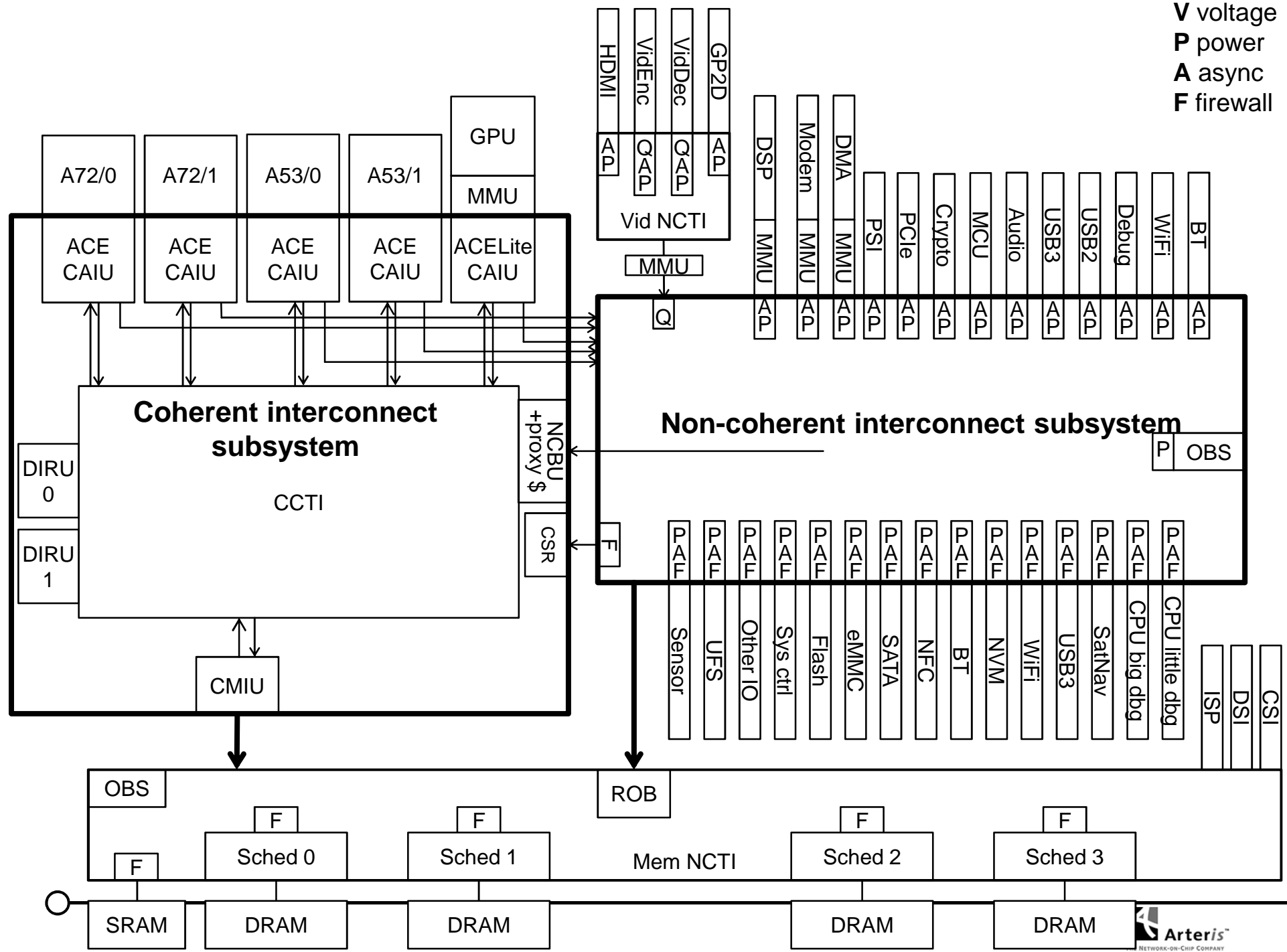


Sharing between non-coherent agents

Using configurable proxy caches



V voltage
P power
A async
F firewall



Heterogeneous cache coherency required to maintain system power/performance trends

Ncore™ Cache Coherent Interconnect IP is targeted at heterogeneous SoCs.

Benefits

- Scalability
- Configurability
- Area efficiency
- High performance
- Optimal power consumption

Major Technologies

- Multiple configurable **snoop filters**
- Multiple configurable **proxy caches**
- Modular **distributed architecture**

RESULT:

Custom-configured interconnect IP that meets exact system requirements

What's next?

Upcoming challenges and trends

- Management of physical constraints
- Consolidation/de-consolidation of coherency protocols
- Hardware layer security
- Hardware layer resilience
- Lower power, lower power, lower power
- High performance
- Off-chip coherency
- 2.5/3D
- Constantly lowering the cost of coherency





Arteris™

THE NETWORK-ON-CHIP COMPANY

Thank YOU

