

# Qualcomm High Performance Processor Core and Platform for Mobile Applications

Lou Mallia, Senior Staff Engineer, Qualcomm Inc.

# The objective

How to get from here .....> To here?



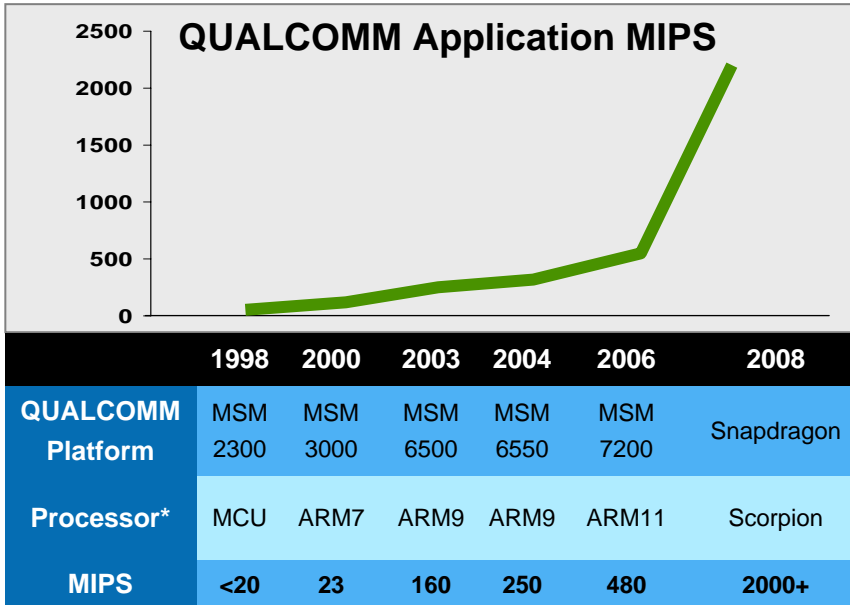
## Features

- QVGA screens
- Fancy ring tones
- Limited graphics
- Snapshot-quality camera
- Single-mode modem



## Features

- WVGA video
- Surround-sound audio
- MP3 player
- PC gaming-quality graphics
- Professional 12MP photos
- Multi-mode modem
- Bluetooth
- GPS w/ navigation
- Web browsing with DRM
- Secure financial transactions
- Live TV
- WiFi
- ...all running at the same time!!



snapdragon™

In short:

*How do we get 2000+ DMIPS and still keep the power below 500mW ???*

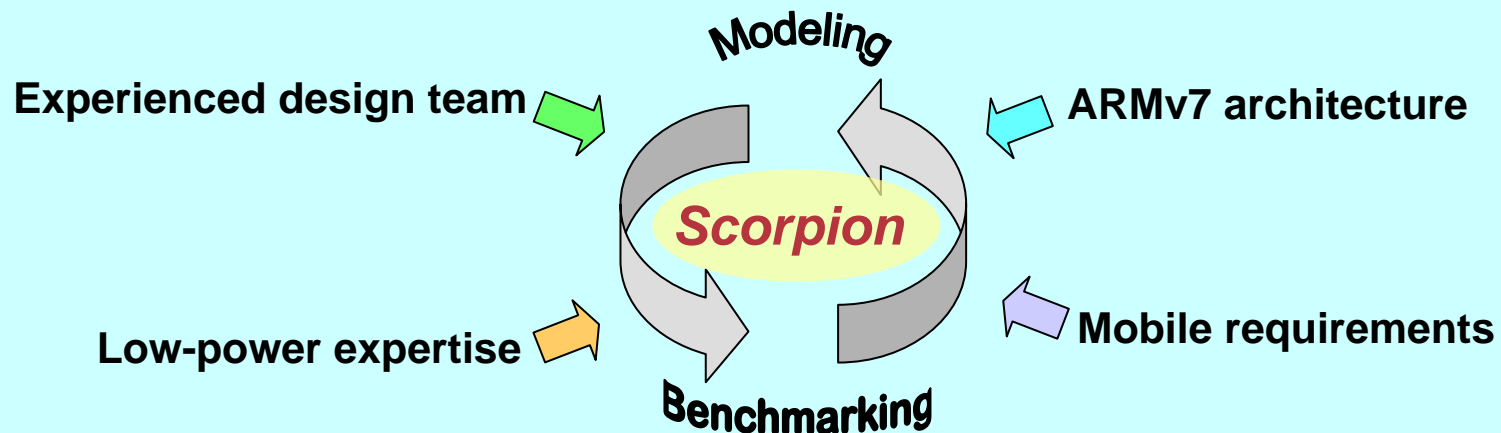
\*QUALCOMM implementations

# The challenges

- Need more MIPS, MFLOPS, and data bandwidth...
  - Performance approaching the level of PC's
- ...but also want lower power and smaller form factor
  - Smaller batteries
  - Always on
- Other design team challenges
  - ARMv7 architecture new for design team
  - Architectural complexity
  - Energy-efficient data movement
  - New OS standards

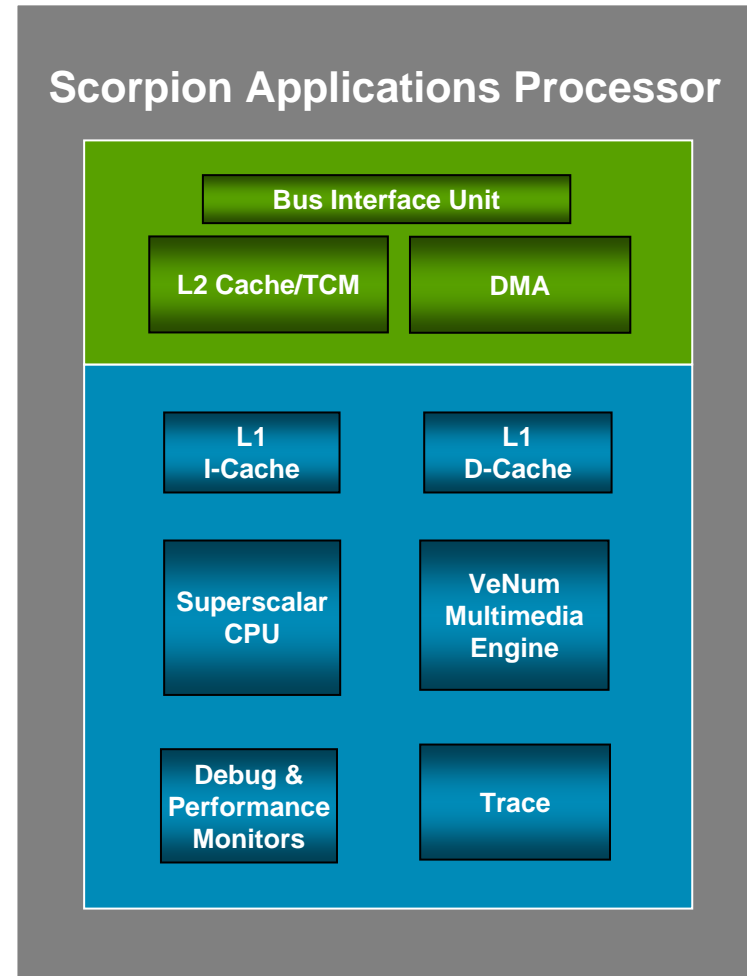
# The plan of attack

- Experienced design team
  - 5 generations of low-power RISC processors
- ARMv7 architecture license
  - New design, not a standard product
  - Partner with ARM on evolution
- Target cycle time of 20-25 FO4
  - Using low-power technology
- Aggressive energy management
- Design, model, and iterate



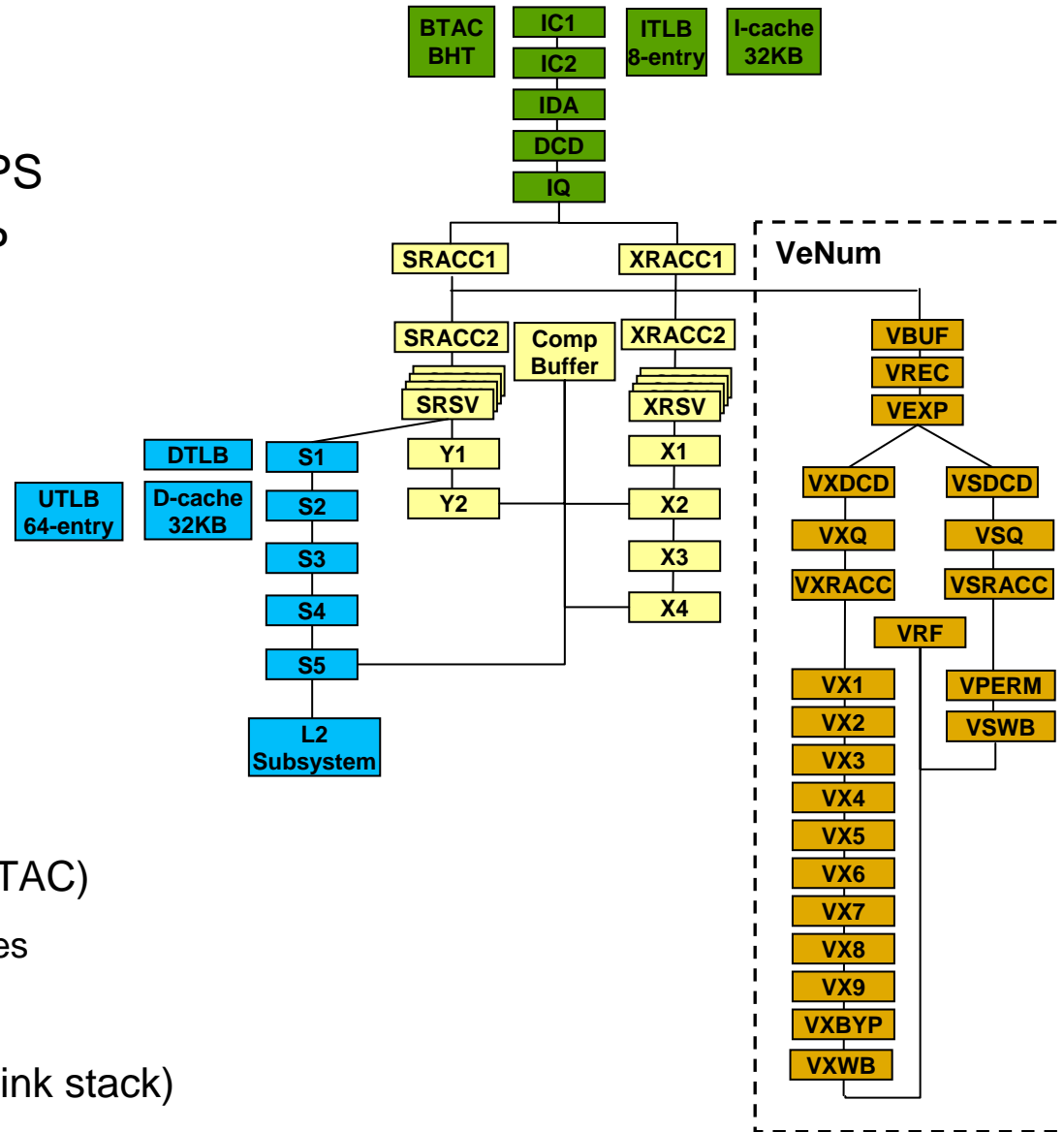
# Scorpion overview

- Highly integrated design
  - Superscalar CPU
  - Tightly-coupled multimedia engine
  - L1 and L2 caches
  - Built-in DMA channels
  - Debug, trace, and performance monitors
- All the latest ARMv7 architectural features, including:
  - Multimedia enhancements
    - Neon 128-bit Advanced SIMD extensions (ASE)
    - VFPv3 floating-point (32 double-precision registers)
    - FP-16 half-precision floating-point format
  - TrustZone™ security extensions



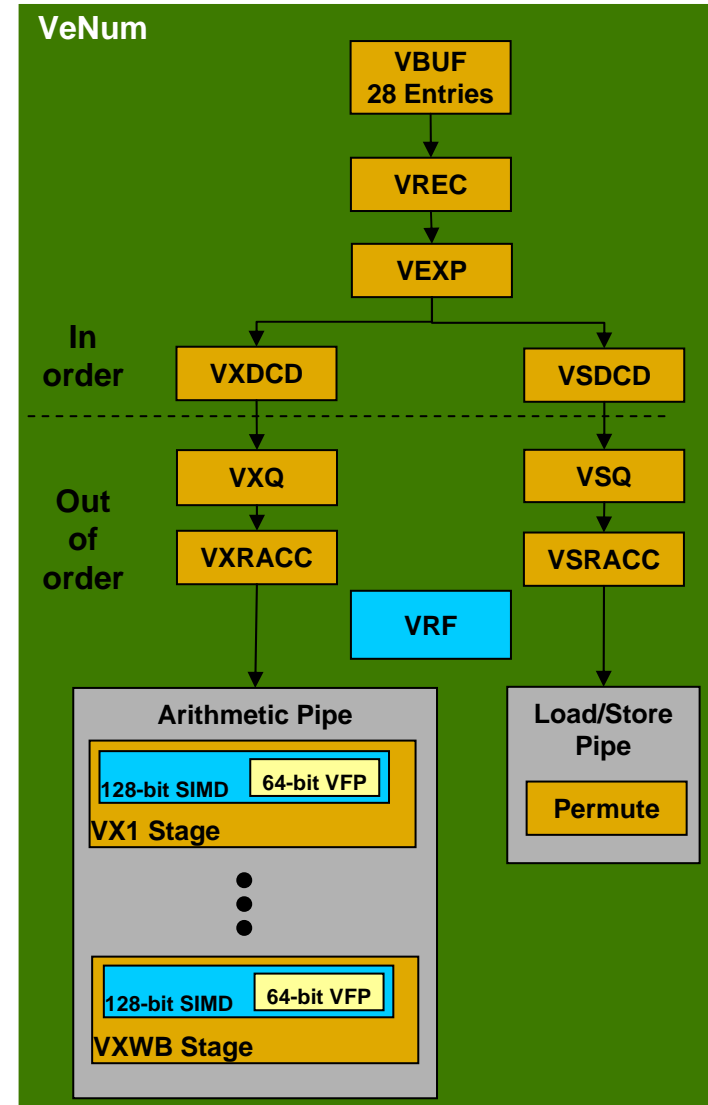
# Scorpion details

- 1.0 GHz (65nm LP technology)
- 2,100 DMIPS and 8,000 MFLOPS
- Power as low as 0.14 mW/DMIP
- Dual-Issue
- Speculative out-of-order issue
- Deeply pipelined (24 FO4)
  - 13-stage load/store pipe
  - 10-/12-stage integer pipes
  - 23-stage floating-point pipe
- Dynamic branch prediction
  - Branch history table (BHT)
  - Branch target address cache (BTAC)
    - 1-cycle penalty on taken branches
  - Global history register (GHR)
  - Subroutine return acceleration (link stack)



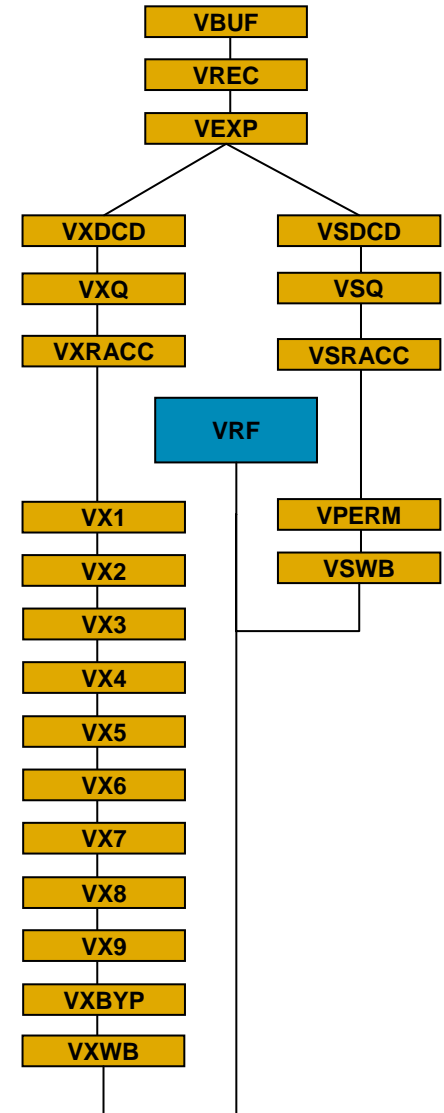
# Scorpion VeNum multimedia engine (1)

- VeNum = “**V**ector **N**umerics”
  - VFPv3 floating-point
  - Neon advanced SIMD extensions (ASE)
  - 11-stage, 128-bit arithmetic and load/store pipelines
    - VFP operations merged with low-order 64-bits of SIMD
    - Unified multiplier (integer and floating-point)
    - No “trundling” of data (early-out bypass and writeback)
  - Dual-issue, out-of-order execution
    - 128-bit load/store plus arithmetic operation
  - No speculative execution (reduces wasted energy)
  - Separate clock domain from CPU (synchronous)



# Scorpion VeNum multimedia engine (2)

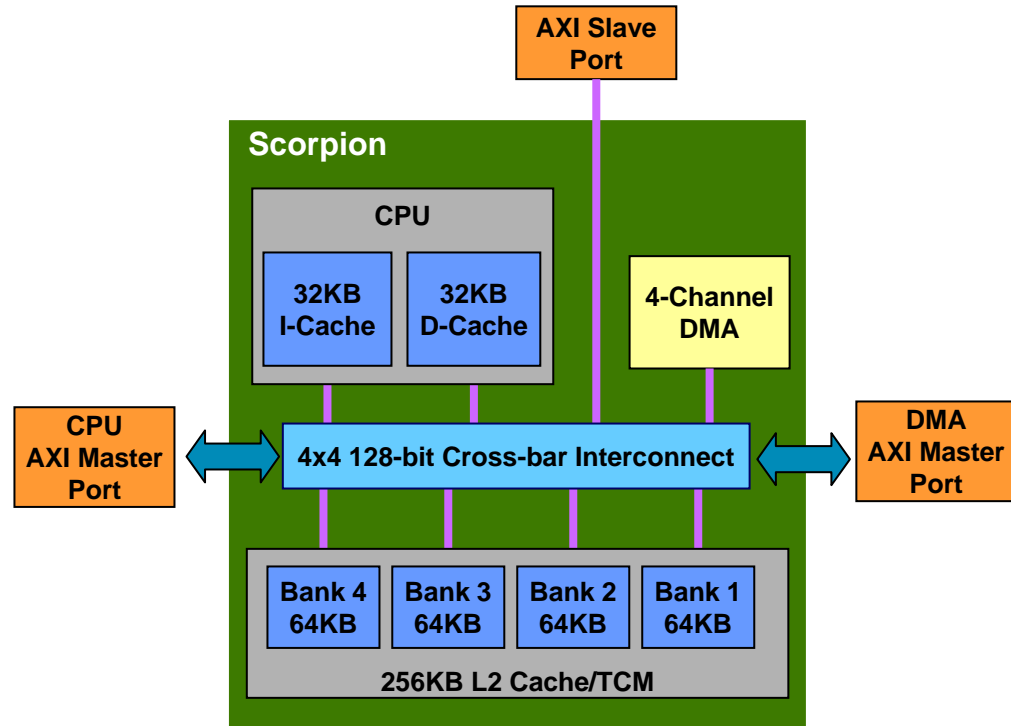
- VFPv3 floating-point
  - Register file expanded to 32x64-bit
  - Pipelined for both double- and single-precision
    - Including subnormals, NaNs, and multiply-add
  - Divide and square root
  - Precise trapping w/ syndromes on IEEE exceptions
- Neon advanced SIMD extension (ASE)
  - Fully-pipelined 128-bit datapath
  - Shares VFP register file (accessed as 16x128-bit)
  - Integer SIMD (16x8-bit, 8x16-bit, 4x32-bit)
  - Floating-point SIMD (4x32-bit single-precision)
    - 8000 MFLOPS
- FP-16 half-precision support
  - Doubles load/store bandwidth and saves energy
  - Conversion operations between half- and single-precision
    - Supports both OpenGL ES 2.0 formats





# Scorpion memory subsystem

- 32KB/32KB L1 instruction/data caches
- 256KB unified L2 array
  - Cache or TCM
    - Configurable by 64KB bank
  - Multi-port TCM access
    - CPU
    - DMA
    - AXI slave port
  - TCM coherent with L1 data cache
- Internal four-channel DMA controller
- Enhanced AXI-based bus architecture
  - Out-of-order transactions
  - Barrier operations
  - Semaphore operation protocol
  - Three ports
    - 64-bit CPU master port
    - 64-bit DMA master port
    - 64-bit slave port (TCM access)
  - Up to 4.8 GB/s throughput



# QUALCOMM energy management (1)

- Technology
  - Low-leakage, multi-Vt CMOS process (65nm, 45nm)
  - Customized by QUALCOMM with multiple fab partners
- Logic design process
  - Low-power front-end design
    - Optimal per-cycle local clock gate for **every** register in design
    - Unused dataflow stabilized on per-cycle basis
  - VeNum 64-bit SIMD multiply mode
    - Limits peak power of 128-bit operations
- Multiple clock domains (Clock-do-Mania™)
  - Dynamic regional clock gating
    - CPU, SIMD/FPU, L2 cache, trace logic
  - Dynamic domain clock gating
    - Processor core, system interfaces, debug

# QUALCOMM energy management (2)

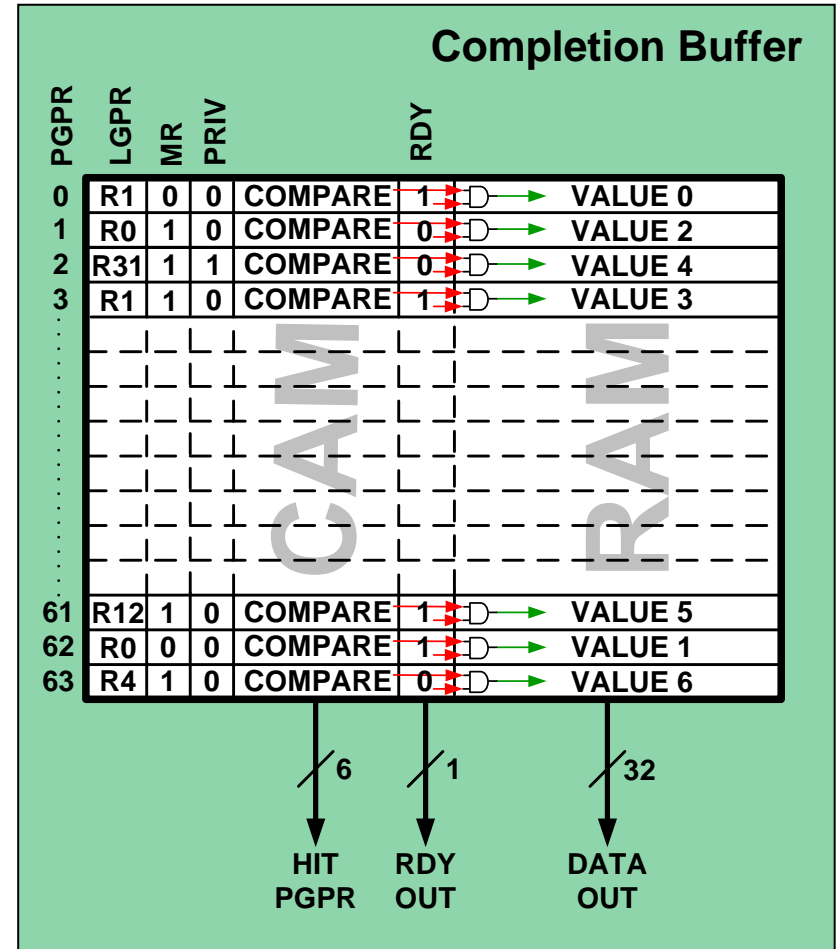
- Dynamic power reduction
  - Selective use of HVT, NVT, LVT devices
  - Low operational voltages
- Leakage power reduction
  - Head/foot switches
  - Ultra-low retention voltage
- Multiple voltage and frequency realms
  - Supported by level shifters and clamps
  - Configured by software
  - Adjusted by hardware
- System-in-package (SIP) stacked memory approach
  - Reduces SDRAM access power

# Completion buffer overview

- ARM architecture defines 32 logical general-purpose registers (LGPRs)
  - 15 user-mode GPRs
  - 17 privileged GPRs
- Completion buffer (CB) supports 64 physical registers (PGPRs)
  - LGPRs “renamed” to PGPRs
- Higher performance
  - Allows pipelining of register hazards
  - Allows out-of-order writeback of results
- Lower power
  - Completion buffer IS the register file
    - No need to move from completion buffer to register file as with reorder buffer
    - Early pipeline results written directly to CB
      - » No trundling through later pipeline stages

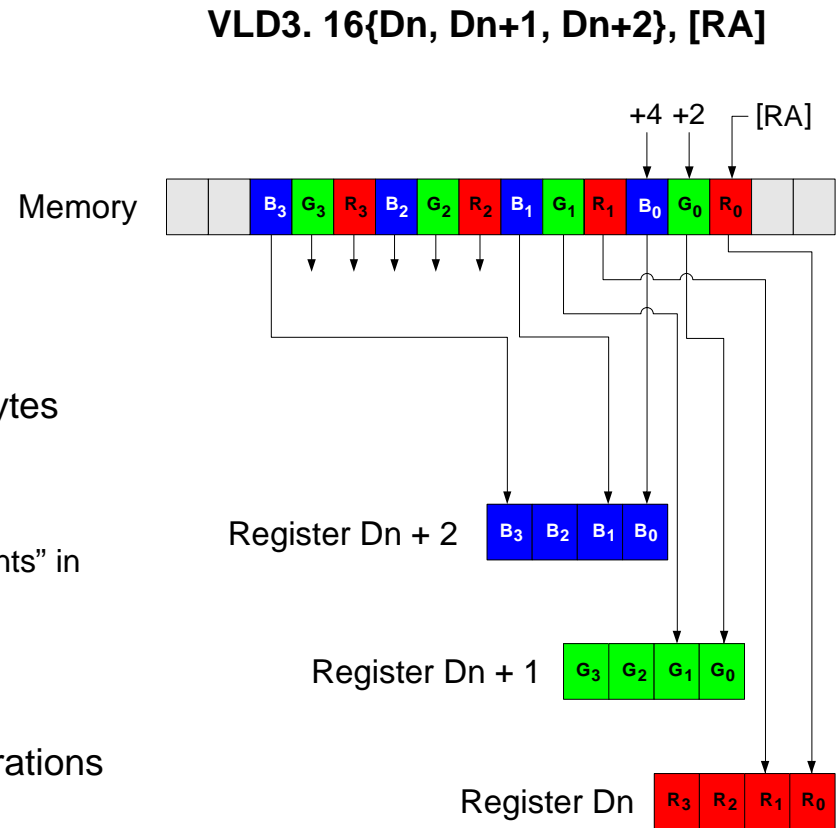
# Completion buffer power-saving features

- CAM search
  - Only search PRIV=0 entries in user mode
  - Only search Most recent (MR=1) entries
  - Gate off comparator for PRIV=1 and/or MR=0
- RAM read
  - Only read matched entry if RDY=1
  - RDY=0 entry won't fire RAM read wordlines
    - Bitlines and outputs unswitched
- Common case for deep pipelines
  - Most source values forwarded from pipeline
    - Not read from completion buffer



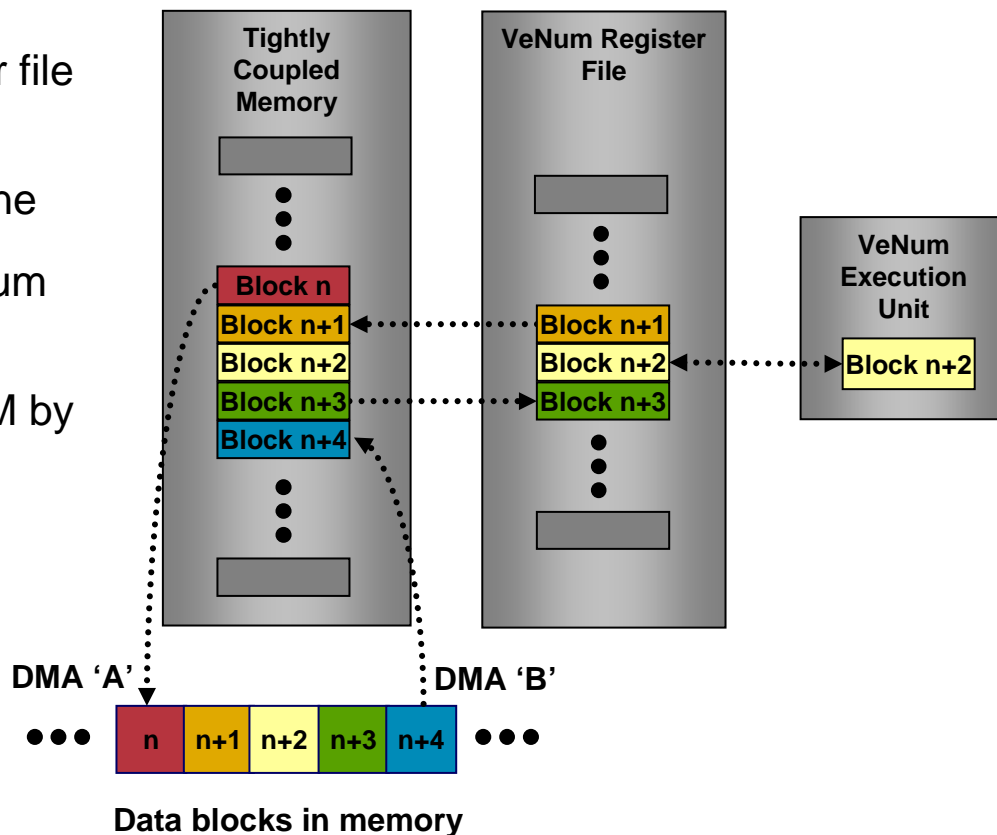
# Multimedia data processing (structure load/store)

- Multimedia data is a series of structures
- Each structure has 1-4 elements
  - 1-element structure (e.g., sampled values)
  - 2-element structure (e.g., coordinates)
  - 3-element structure (e.g., color space)
  - 4-element structure (e.g., 3D graphics)
- Dilemma
  - Registers normally filled in-order with sequential bytes
    - First register gets filled first, then next register
  - But, processing algorithms require different order
    - Put all “first elements” in register 1, all “second elements” in register 2, etc.
- Solution – Auto-permuting load/store operations
  - Elements auto-permuted into registers “on-the-fly”
  - Saves energy by avoiding read-permute-write operations
- Example
  - Four, 3-element structures (16-bits per element)
  - Loaded into three doubleword registers (Dn, Dn+1, Dn+2)



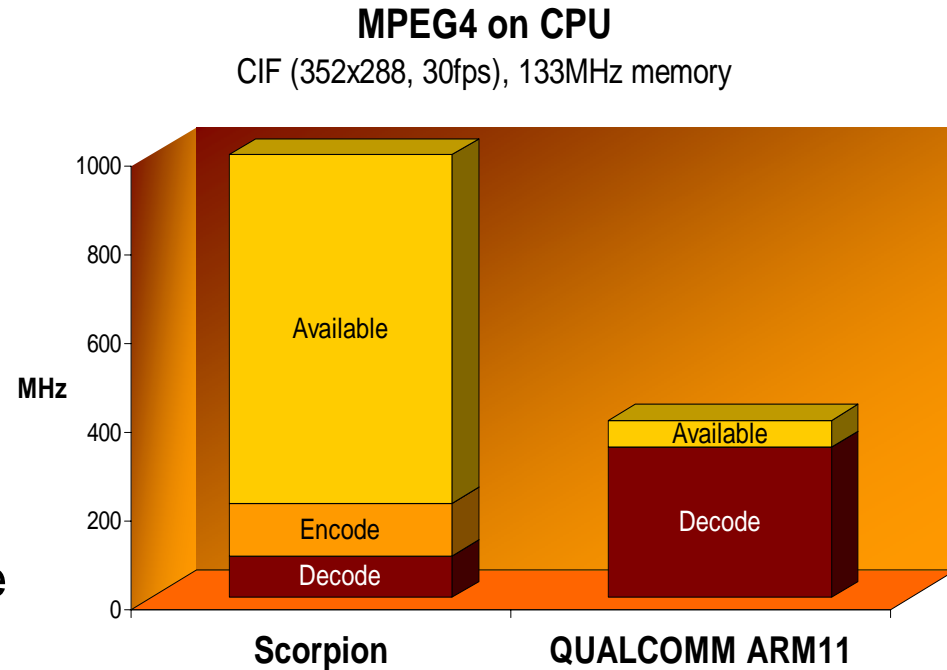
# Multimedia data processing (pipelined DMA + TCM)

- Leverage Scorpion DMA and TCM
  - Block n written from TCM back to memory by DMA channel “A”
  - Block n+1 stored from VeNum register file back to TCM
  - Block n+2 processed in VeNum pipeline
  - Block n+3 loaded from TCM into VeNum register file
  - Block n+4 read from memory into TCM by DMA channel “B”
- TCM and L1 D-cache kept coherent



# MPEG-4 encode/decode performance

- Leverages VeNum SIMD engine
  - Cosine transforms
  - Deblocking filters
  - Motion estimation
  - Color space conversion
- Full-duplex video encode/decode
  - Almost 800MHz of headroom
    - Available for OS and other tasks

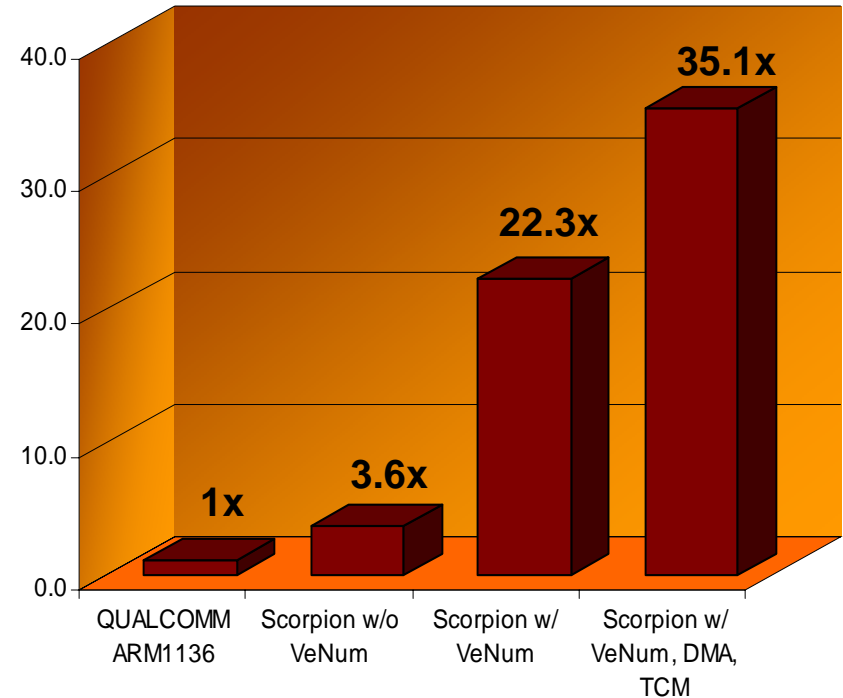




# Graphics vertex processing performance

- Vertex processing gaming application
- Improvement from QUALCOMM ARM1136
  - 3.6x due to clock speed & micro-architecture
  - 6.2x more from VeNum SIMD engine
  - 1.6x more from DMA plus TCM
- Total Scorpion improvement: **35x !!**
- Complements graphics processing unit (GPU)
  - Software chooses where to split the graphics processing chain

Graphics Vertex Processing  
(vertices/sec)  
Normalized to QUALCOMM 400MHz ARM1136



# Scorpion Summary

- Scorpion processor core for mobile applications
  - A unique micro-architectural realization of the ARMv7 architecture
    - Provides maximum energy efficiency at high performance levels
    - Performance up to 2100 DMIPS and 8000 MFLOPS at 1GHz
      - » Using 65nm LP technology
    - Power as low as 0.14mW/DMIP
  - The cornerstone of QUALCOMM's Snapdragon technology platform
    - QUALCOMM creating a variety of Scorpion-based Snapdragon products for different applications

# High Performance Mobile Platform Challenges

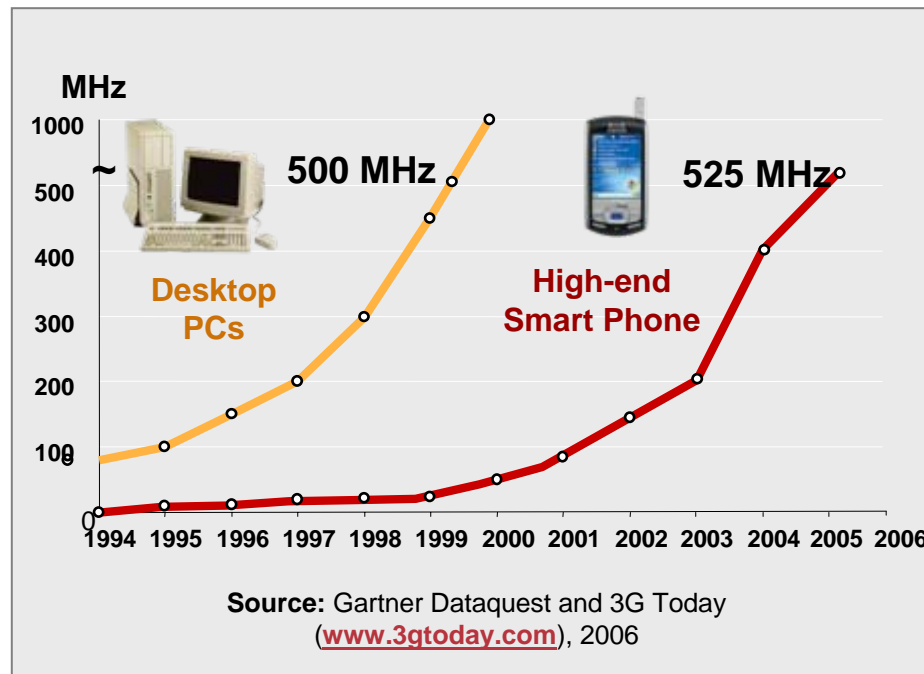


- Mobile platforms
  - Complex systems with a wide range of applications
  - Trade-offs between generic and customized processing solutions
  - Complex system modeling and architecture trade-offs
- Key challenges
  - Optimizing system bandwidth for multiple concurrent applications
  - Extremely power efficient designs for maximum battery life
  - Minimum footprint to enable small, lightweight mobile form factors

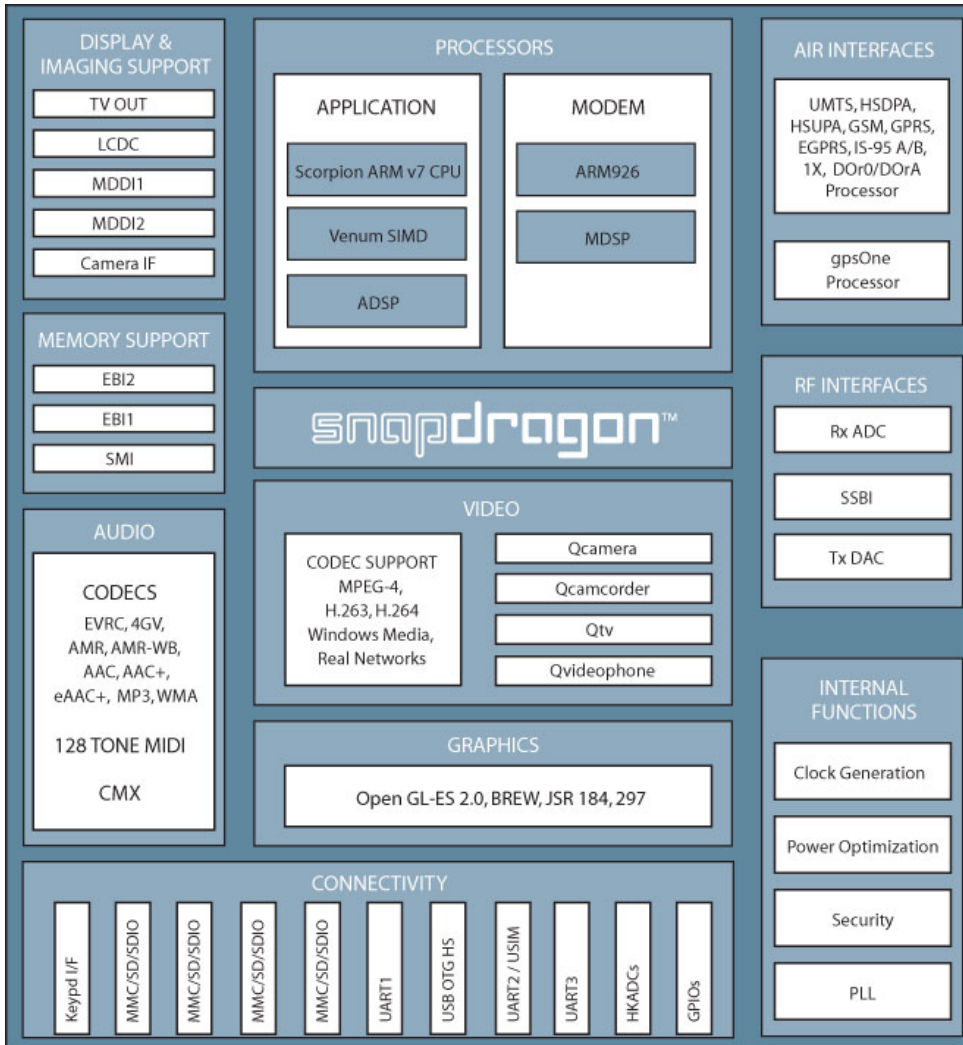
# Key Mobile Platform Trends

- Multimedia and connected applications driving need for more MIPS
  - Performance approaching the level of desktop PC's
- More features in smaller form factors drives need for lower power and higher integration
  - Smaller batteries
  - Always on

## Computing Performance Trend



# Snapdragon Highly Integrated Platform



- **Always On**

- Low power consumption through custom CPU and DSP cores
- All the performance of a laptop in your pocket and much more battery life

- **Industry leading Performance**

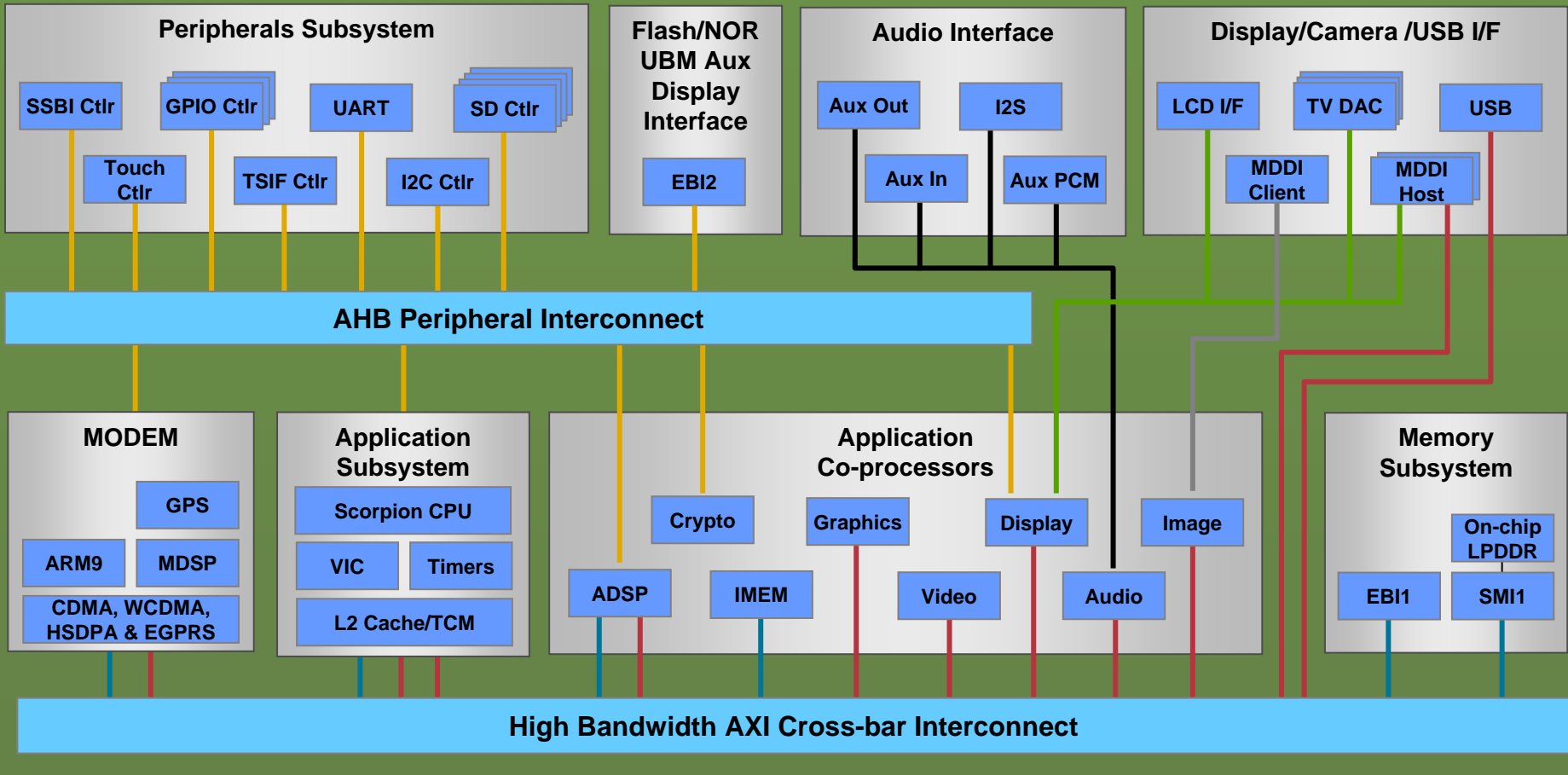
- Superscalar CPU: Scorpion surpasses 2100 DMIPS at 1 GHz speeds
- Next Generation DSP running at 600MHz
- High resolution up to XGA support for uncompromised Video and Computing

- **Ubiquitous Connectivity**

- CDMA, WCDMA, HSPA, GPS, Bluetooth, WiFi, Broadcast (MediaFLO, DVB-H, etc.)

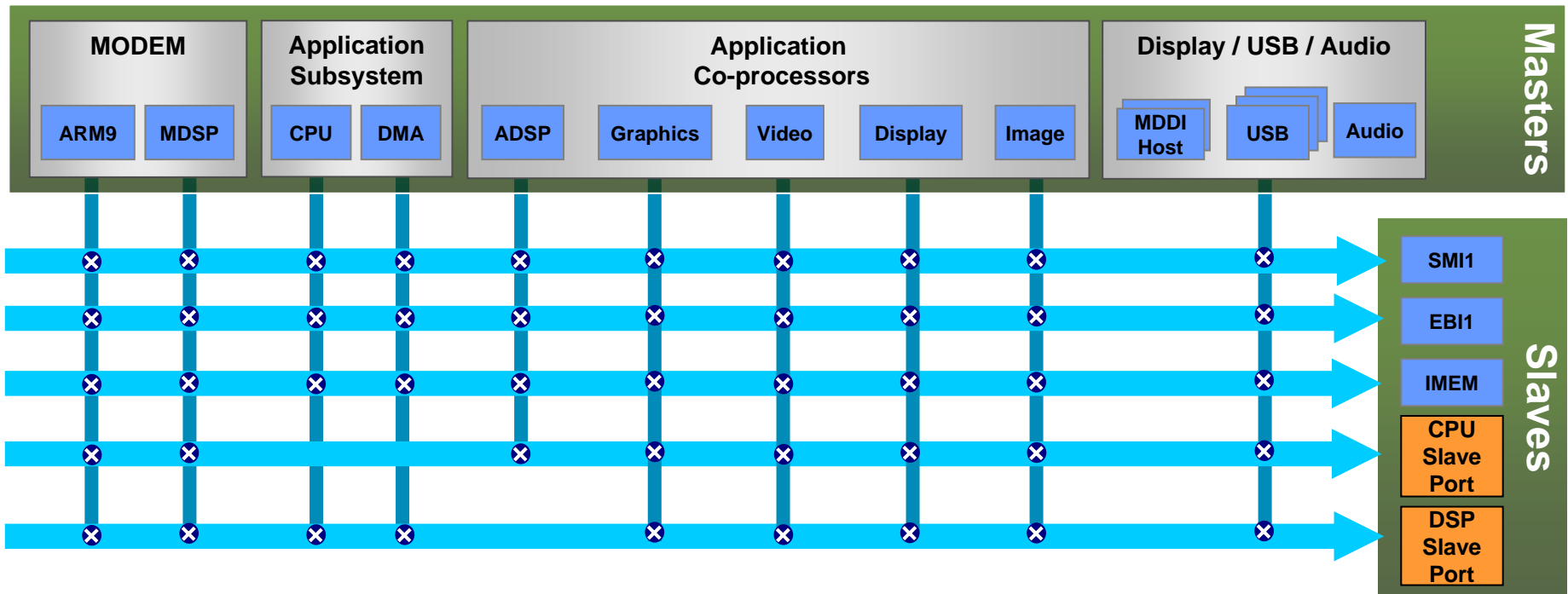
# Snapdragon Platform Connectivity

- Cross-bar interconnect to enable simultaneous traffic
- Balanced interconnect enabling any master to access any slave
- Tiered bus structure to off-load low bandwidth/latency tolerant traffic



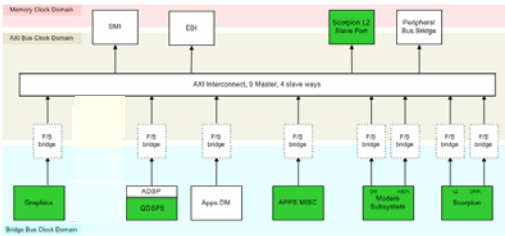
# Snapdragon High Bandwidth AXI Interconnect

- Enhanced AXI architecture
  - Memory barriers, memory type and attributes, ordering
- Configurable full cross-bar implementation
  - Arbiter for each slave interface, parallel R/W data paths
  - Configurable number of masters & slaves, queue depths, pipeline depths
  - Simultaneous access to all slaves
- Integrated performance monitor and bus trace



# Snapdragon Bandwidth Modeling Methodology

## 1. System C Model



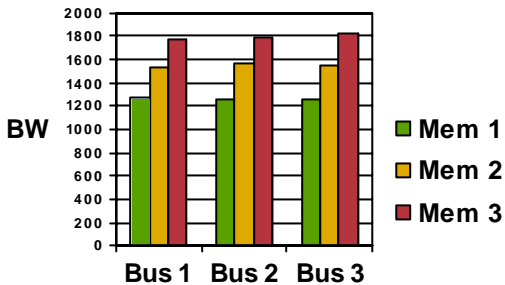
- **Goal**
  - Optimal interconnect for mobile applications
- **Methodology**
  - Model bus & memory application traffic
  - Characterize various bus and memory alternatives
  - Key criteria includes bandwidth, latency, and utilization

## 2. Application Traffic BW

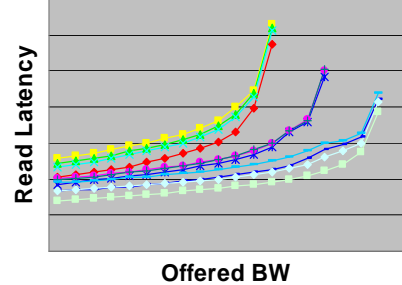
System	Mem	AXI	GPU	ADSP	Apps	MPPE	Graphics	Camera
System 1	100	100	100	100	100	100	100	100
System 2	100	100	100	100	100	100	100	100
System 3	100	100	100	100	100	100	100	100
System 4	100	100	100	100	100	100	100	100
System 5	100	100	100	100	100	100	100	100
System 6	100	100	100	100	100	100	100	100
System 7	100	100	100	100	100	100	100	100
System 8	100	100	100	100	100	100	100	100
System 9	100	100	100	100	100	100	100	100
System 10	100	100	100	100	100	100	100	100

- **Conclusions**
  - System bandwidth was limited by memory bandwidth
  - AXI Bus utilization was less than 50% of theoretical maximum
  - Memory controller enhancements delivered up to 40% system bandwidth improvement without increased AXI bus frequency

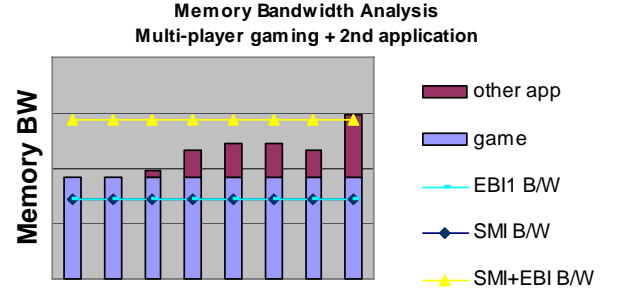
## 3. Memory / Bus Freq. Analysis



## 4. Latency/BW Analysis



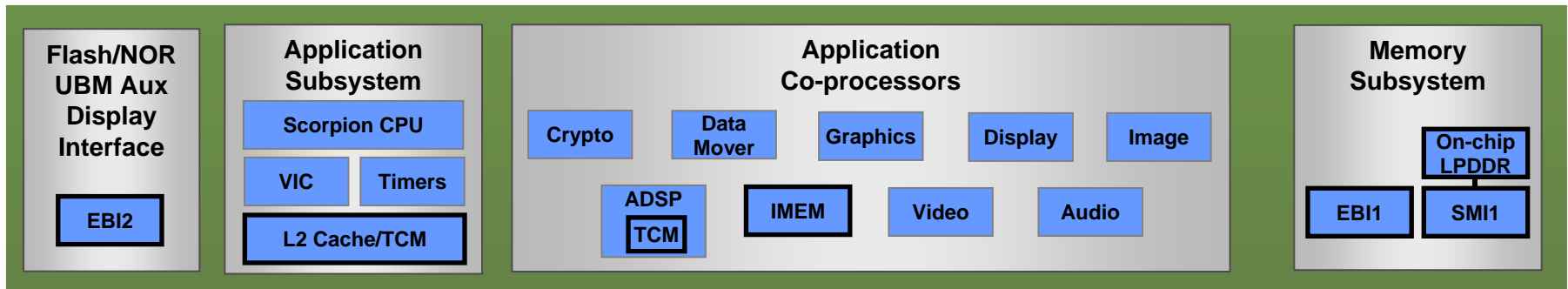
## 5. Concurrent Application BW Analysis





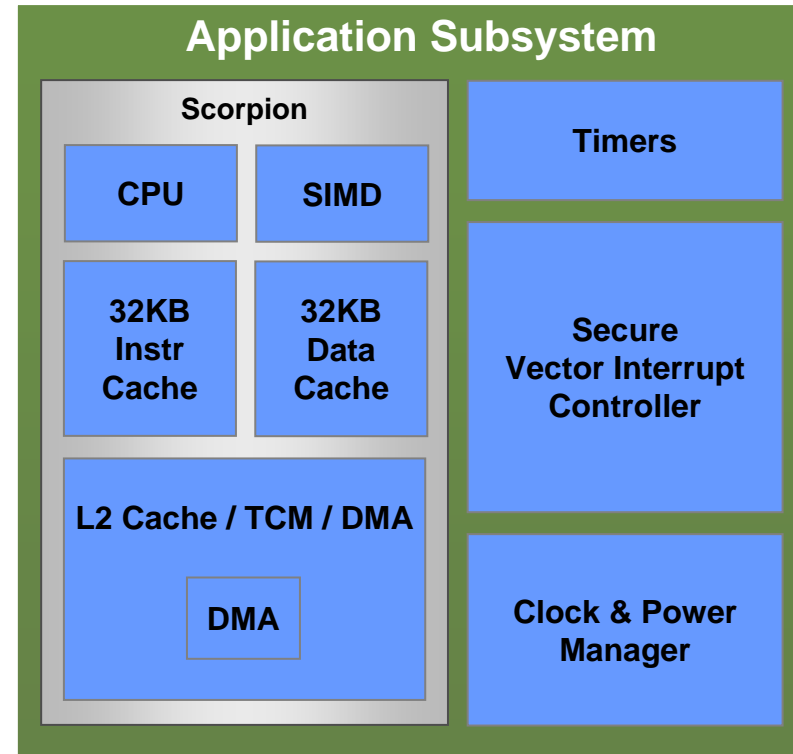
# Snapdragon Memory Overview

- **Multiple memory interfaces**
  - Stacked LP DDR SDRAM reduces power dissipation and board area
  - External LP DDR SDRAM provides additional bandwidth
- **Highly power optimized external memory controller**
  - Power-down, deep power down, clock stop,
  - Self refresh, auto refresh, directed auto-refresh, temperature adjusted refresh rates
  - IO calibration to adjust IO impedance
- **External bus interface controller (EBI2) supports multiple memory options**
  - NAND, OneNAND/M-systems, burst NOR support
- **Integrated on-chip memory**
  - IMEM reduces off-chip memory accesses
  - CPU and DSP L2 caches can be configured as tightly coupled on-chip memories



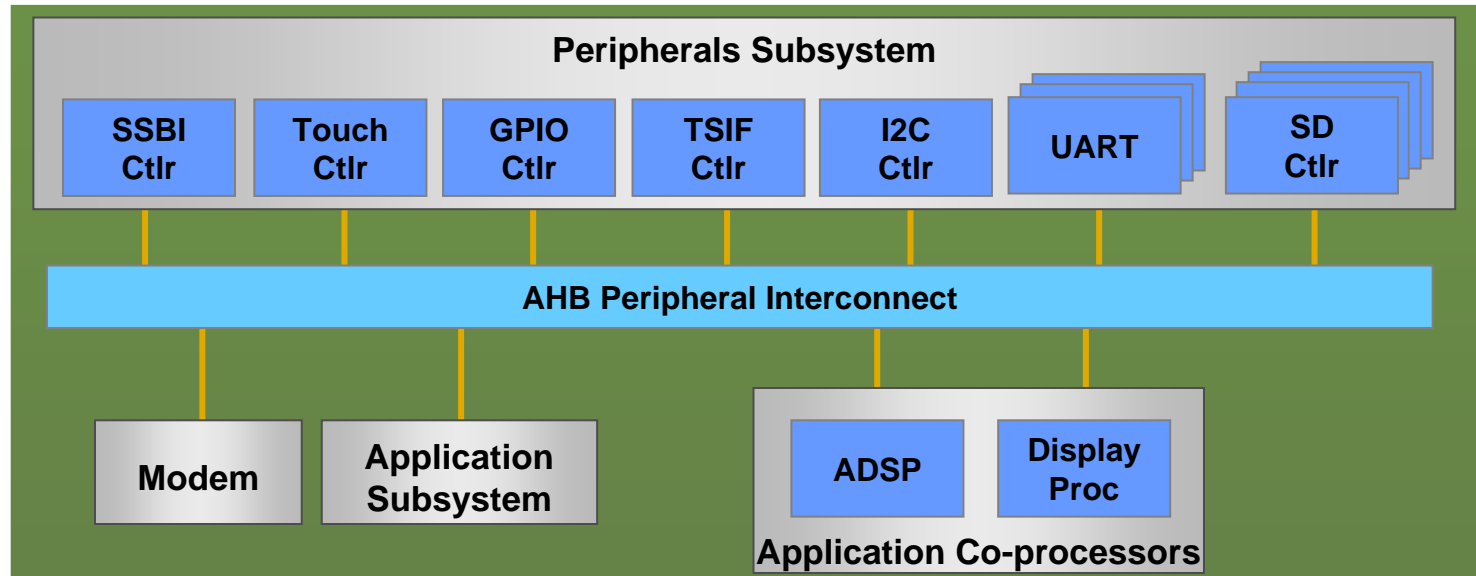
# Snapdragon Application CPU Subsystem

- Secure vector interrupt controller
  - Configurable up to 64 primary interrupts
  - 8-level Prioritized Interrupts to FIQ/IRQ
  - 32-bit IRQ/FIQ vector address
  - Trustzone™ compliant security mechanism
- Clock and power manager
  - High frequency, low jitter PLL
  - Clock source selection, gating and routing
  - Power collapse and voltage scaling
- RTOS, general purpose & secure timers



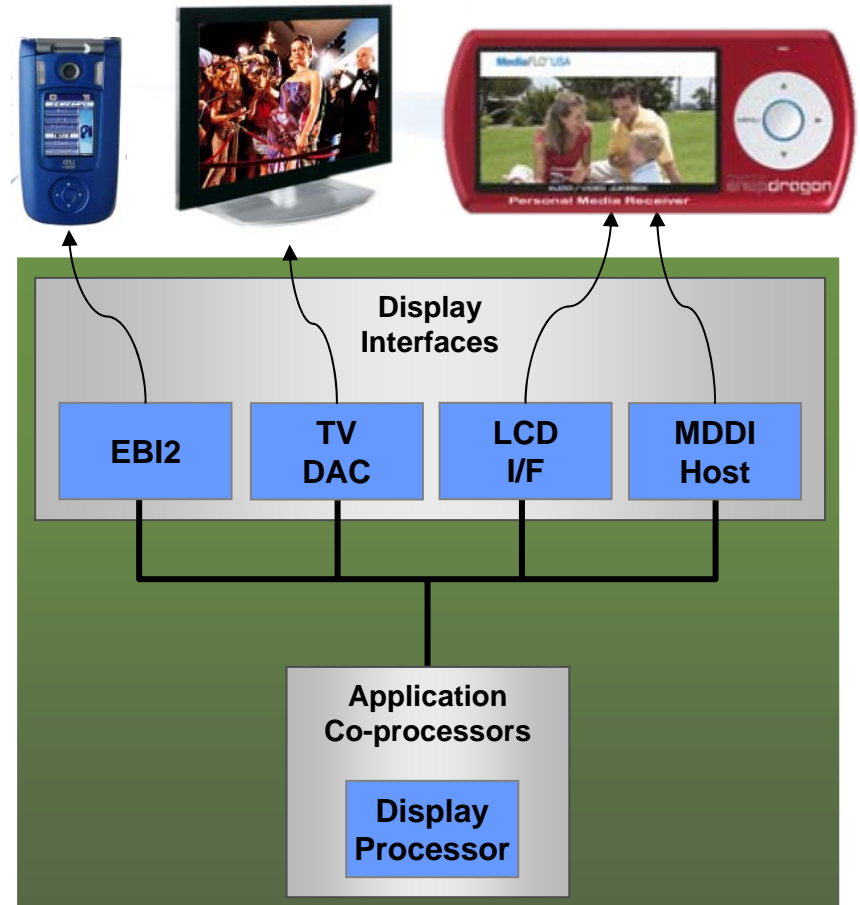
# Snapdragon Peripheral Subsystem

- Provides connectivity to peripheral devices
  - Direct access from processors
  - Off-loads peripheral traffic from memory bus
  - Round robin arbitration with 5 levels of priority
  - Memory protection for secure peripherals



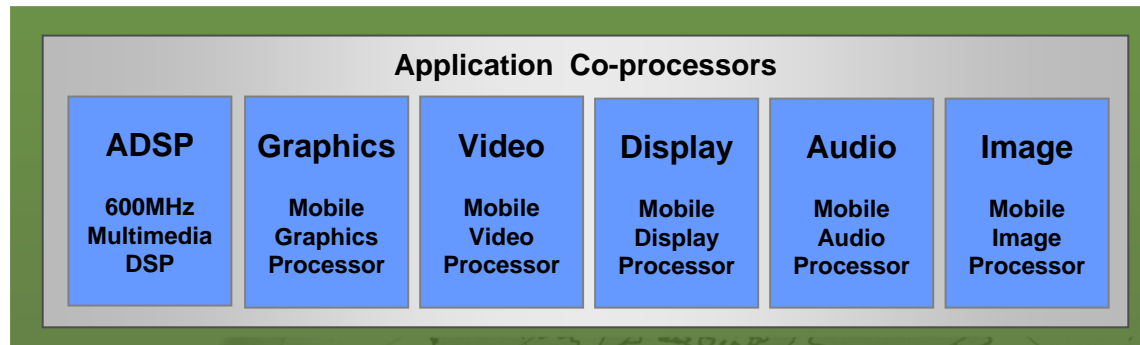
# Snapdragon Display Support

- RGB LCD Controller
  - Supports direct-attach LCD panels up to XGA at 60Hz
  - 24 bit RGB outputs, programmable refresh rates and display sizes
- Mobile Display Digital Interface (MDDI)
  - High-speed serial communication for displays and sensors
  - Type II (1 Gbps) MDDI interface
  - Displays up to XGA
- TV Out
  - Composite and S-Video output supported
  - Integrated 10-bit DAC, NTSC or PAL
- Auxiliary LCD interface for sub displays



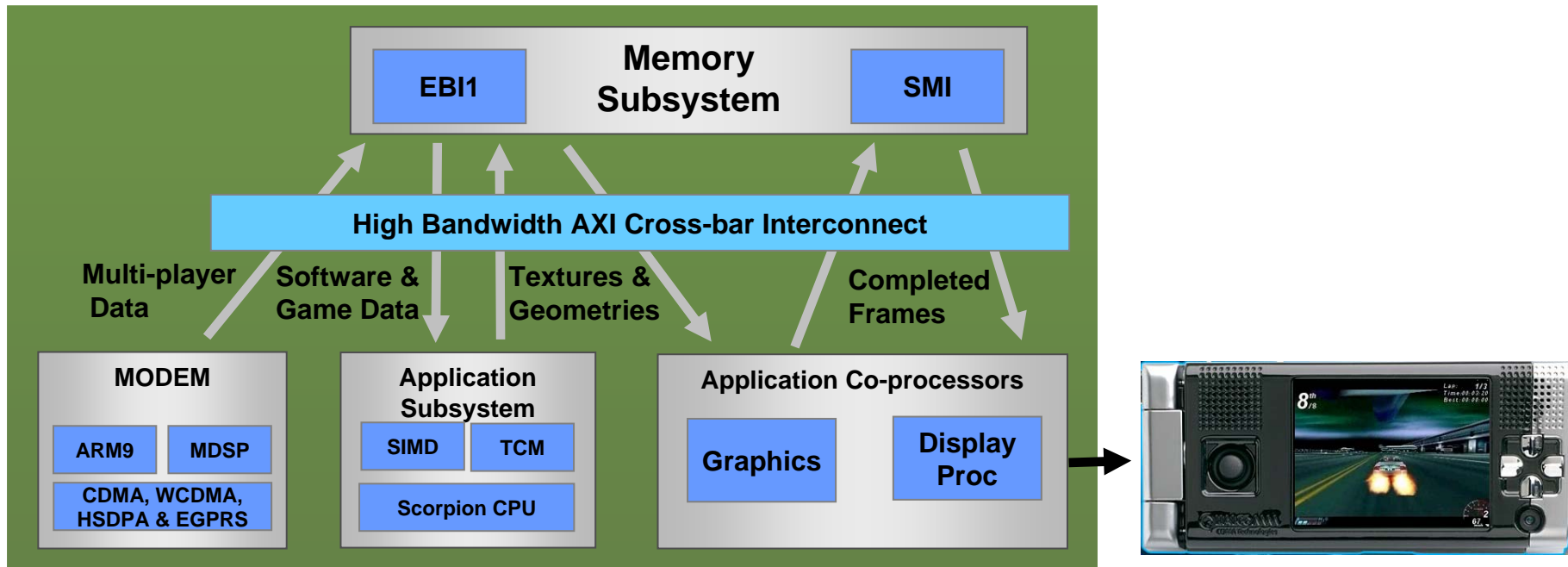
# Snapdragon Multimedia Co-Processors

- **Integrated high-performance, power-optimized multimedia processing engines**
  - Programmable to enable adaptation of emerging standards
- **Multimedia DSP**
  - Custom QUALCOMM designed 600MHz DSP
- **Mobile Graphics Processor**
  - Support for OpenGL ES 2.0
  - 133M pixel/sec or 21M triangles/sec
- **Mobile Video Processor**
  - Video encoding and decoding
  - Supporting H.263 H.264, and MPEG-4
- **Mobile Display Processor**
  - Integrated LCD controller
  - Image Processing (e.g. rotate, scale)
- **Mobile Audio Processor**
  - wideband stereo CODEC,
  - I2S(Inter-IC Sound),
  - PCM, and Dual Microphone Support
- **Mobile Image Processor**
  - Camera sensor image processor
  - Viewfinder
  - Video and image capture
  - Snapshot processing
  - Encoding
  - Image display, and image processing



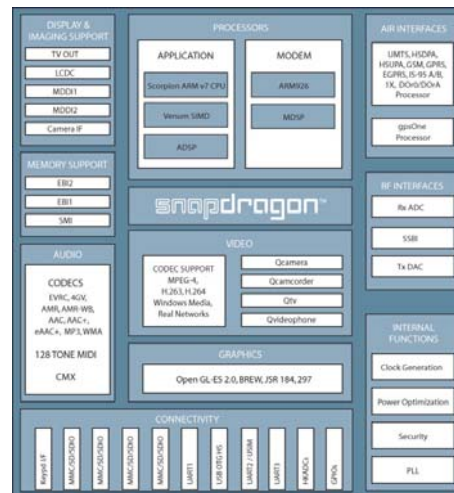
# Network Gaming Example

- Application requires power efficient data movement
- Crossbar interconnect provides parallel data paths
- Tightly coupled memory is used to store intermediate results
- Separate read and write data paths allow concurrent data transfers



# Snapdragon Summary

- Snapdragon platform is a highly integrated, high performance, power optimized mobile solution
  - High performance, power efficient applications processor, multimedia DSP, multimedia applications co-processors
  - Configurable, power efficient interconnect optimized to enable advanced, concurrent mobile applications
- Snapdragon platform enables more advanced applications in smaller, longer lasting, always connected mobile devices



**Thank you!**