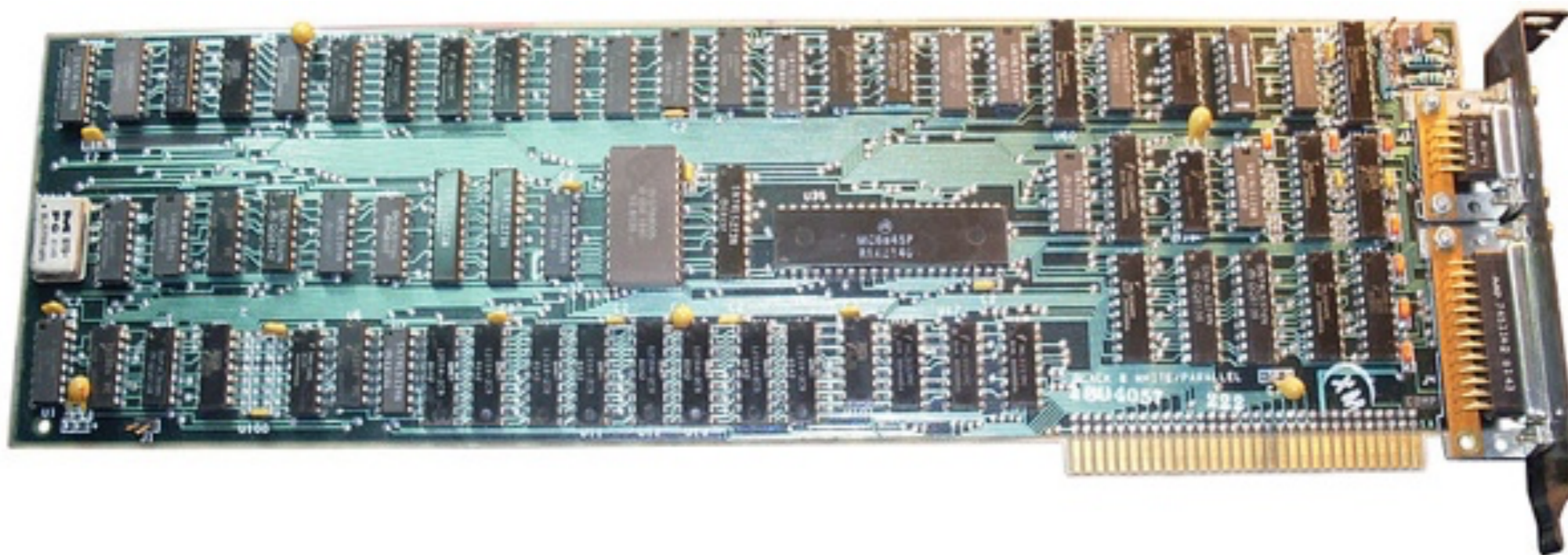


- Museum (this Wednesday)
- No class \*next\* Wednesday
- Data center visit on the 24th (more details to come)
- Midterm II May 31st
- 2 more class days :(

GPUs



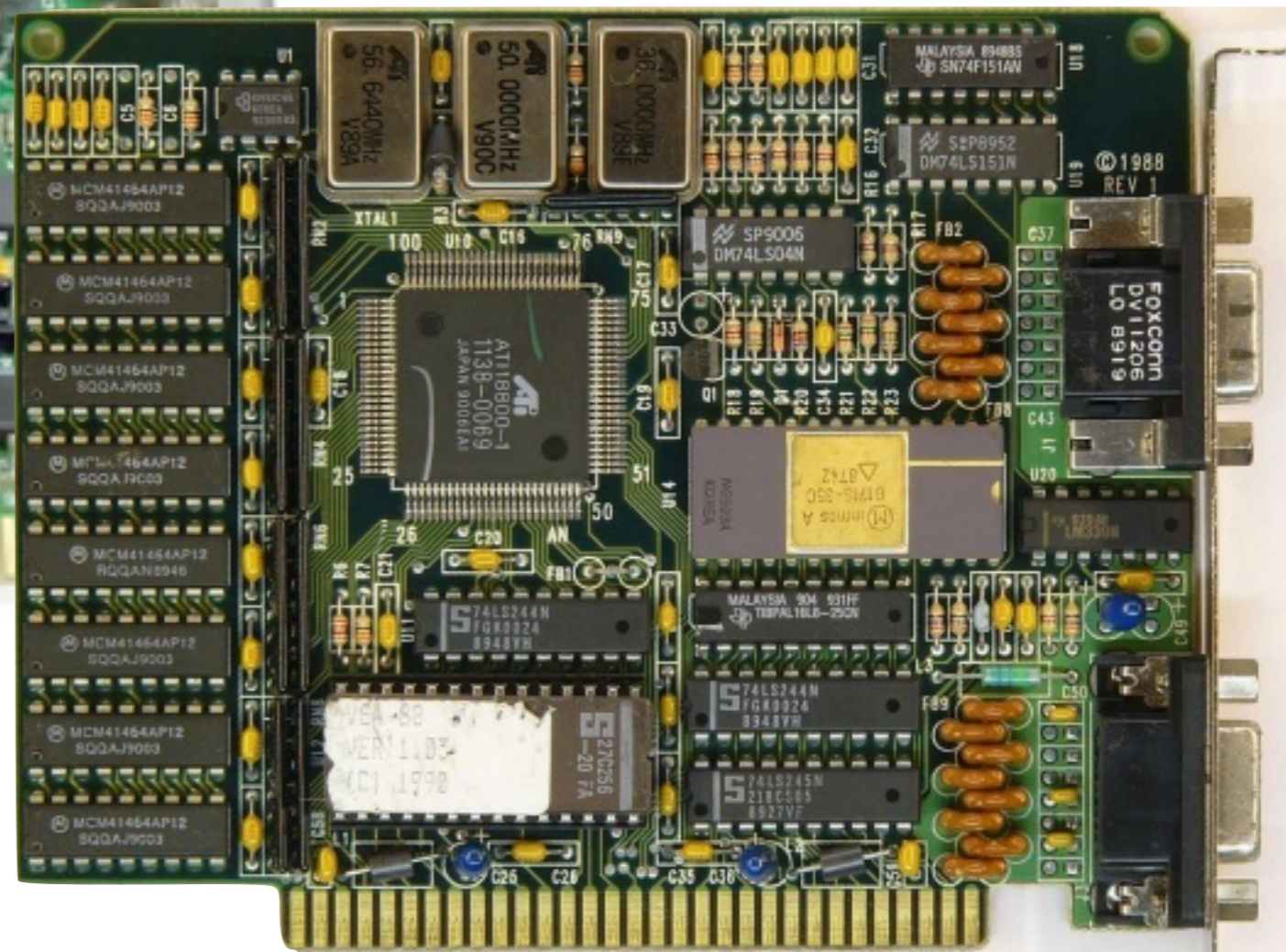
IBM monochrome adapter (1981)  
*w/parallel printer port :)*

Characters and raster mode





ATI (now AMD) EGA  
Wonder (1987)  
2D Graphics



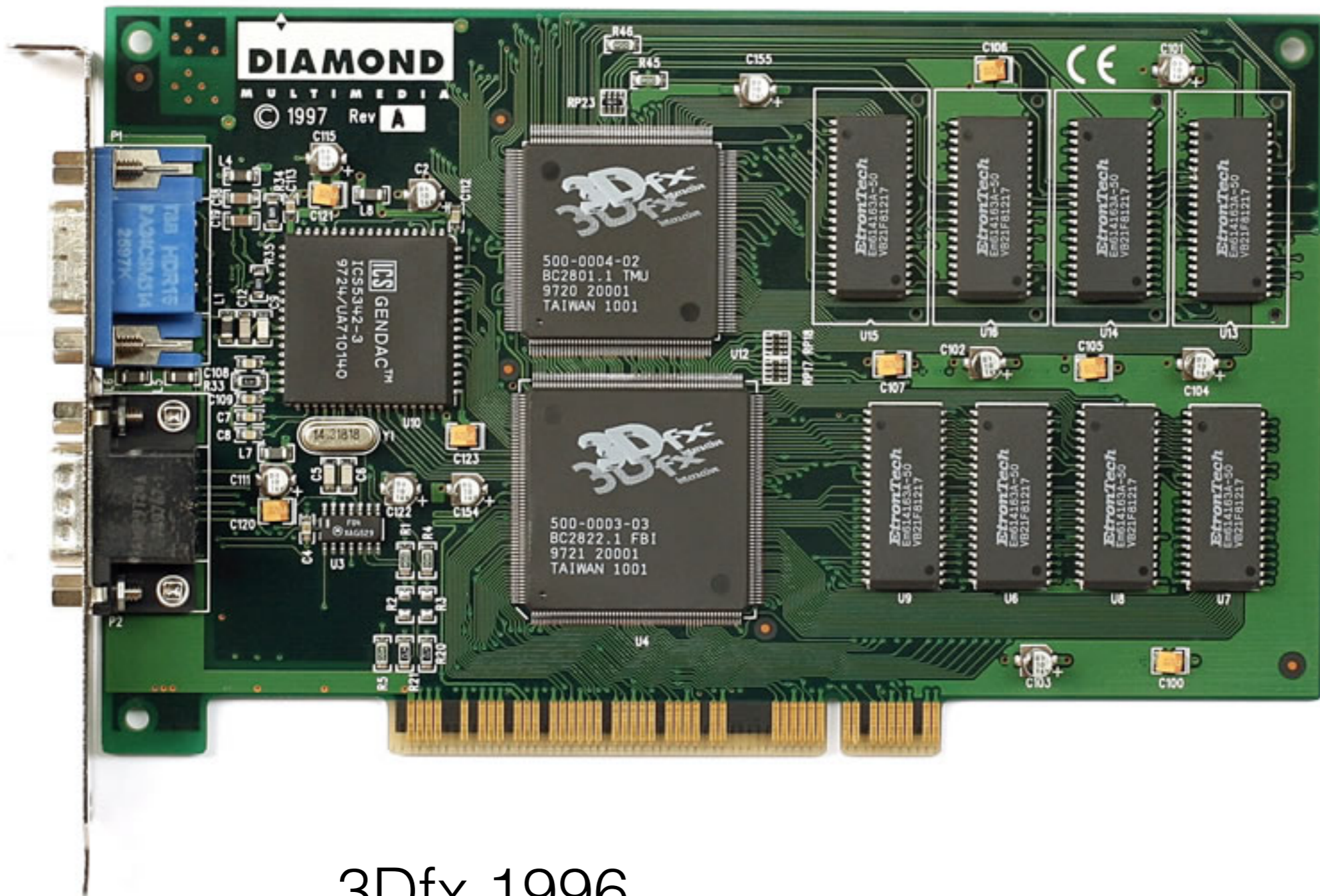
ATI VGA Wonder 1988





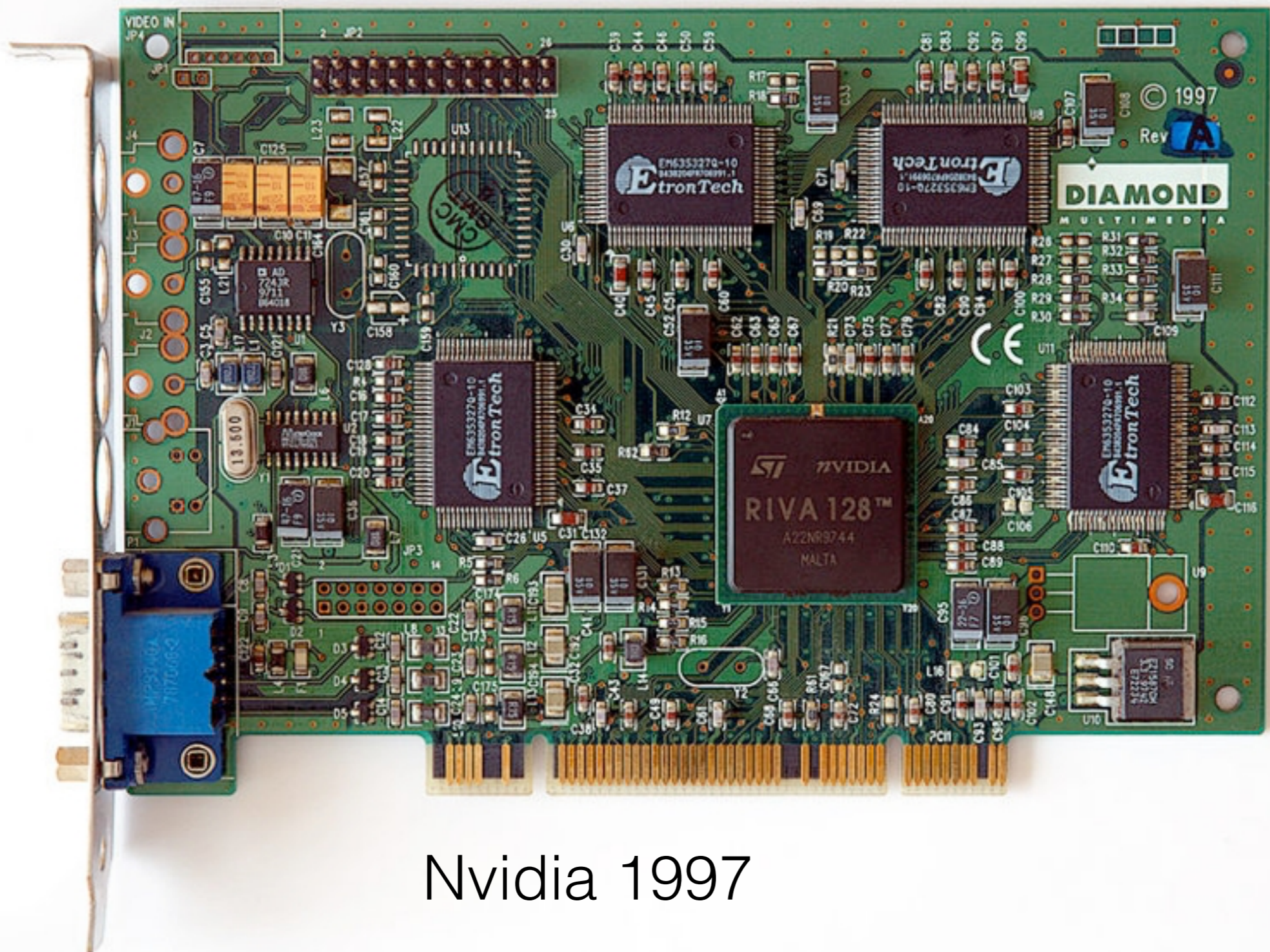
S3 Trio (1995) 3D Graphics





3Dfx 1996





Nvidia 1997





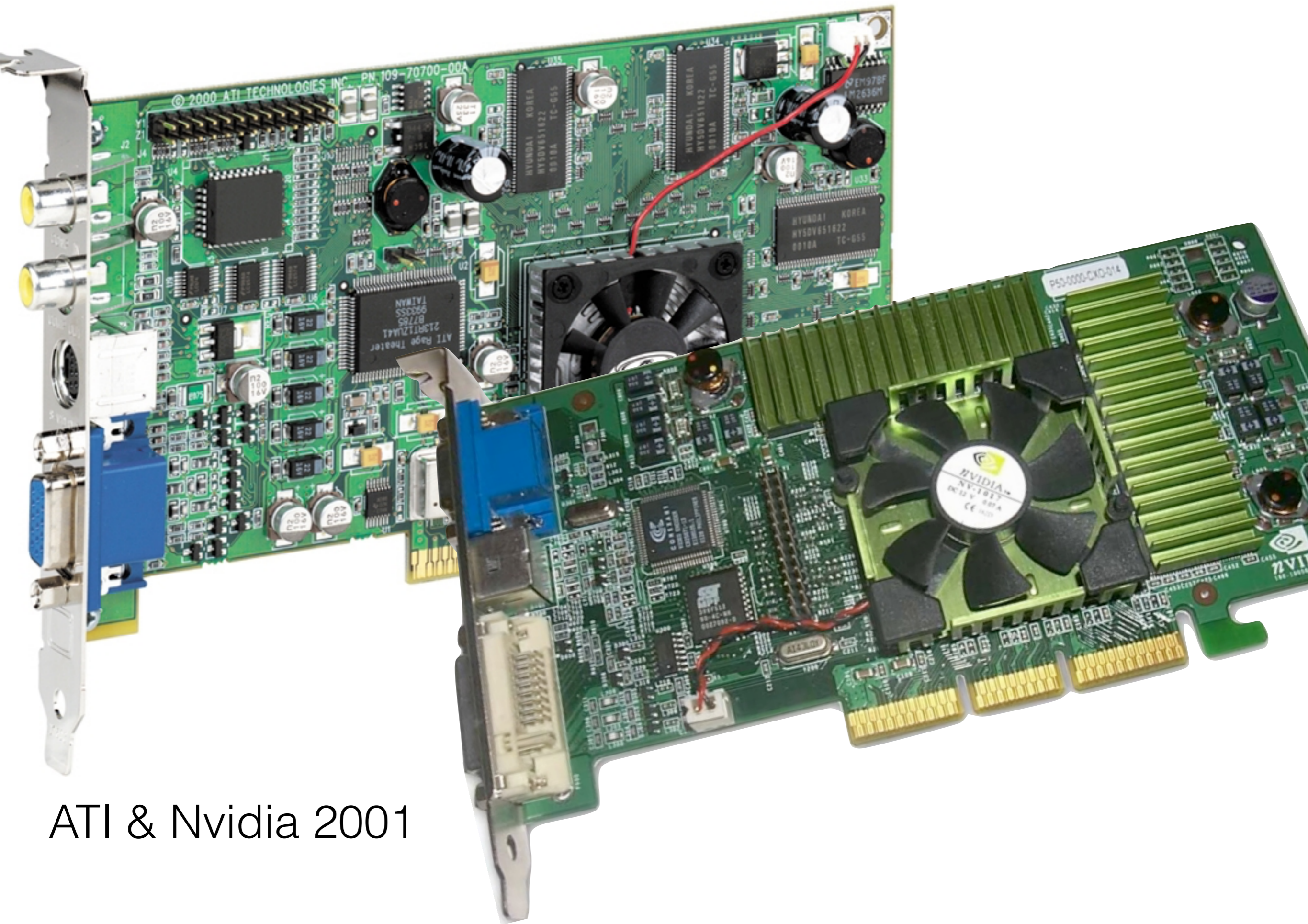
1998 *Note fans!*





3Dfx Voodoo 5 1999  
Note fans & auxiliary power connector



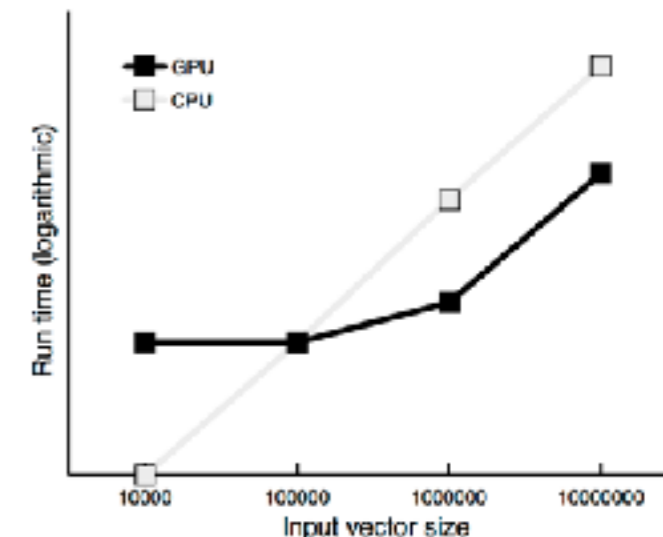


ATI & Nvidia 2001



# Using Modern Graphics Architectures for General-Purpose Computing: A Framework and Analysis

Chris J. Thompson    Sahngyun Hahn    Mark Oskin  
*Department of Computer Science and Engineering*  
*University of Washington*  
 {cthomp, syhahn, oskin}@cs.washington.edu



Winter 2002, published Fall 2002

Table 3. The vertex program instruction set.

Opcode	Description
ARL	Address register load
MOV	Move
MUL	Multiply
ADD	Add
SUB	Subtract
MAD	Multiply and add
ABS	Absolute value
RCP	Reciprocal
RCC	Reciprocal (clamped)
RSQ	Reciprocal square root
DP3	3-component dot product
DP4	4-component dot product
DPH	Homogenous dot product
DST	Cartesian distance
MIN	Minimum
MAX	Maximum
SLT	Set on less than
SGE	Set on greater/equal than
EXP	Exponential base 2
LOG	Logarithm base 2
LIT	Light coefficient formula

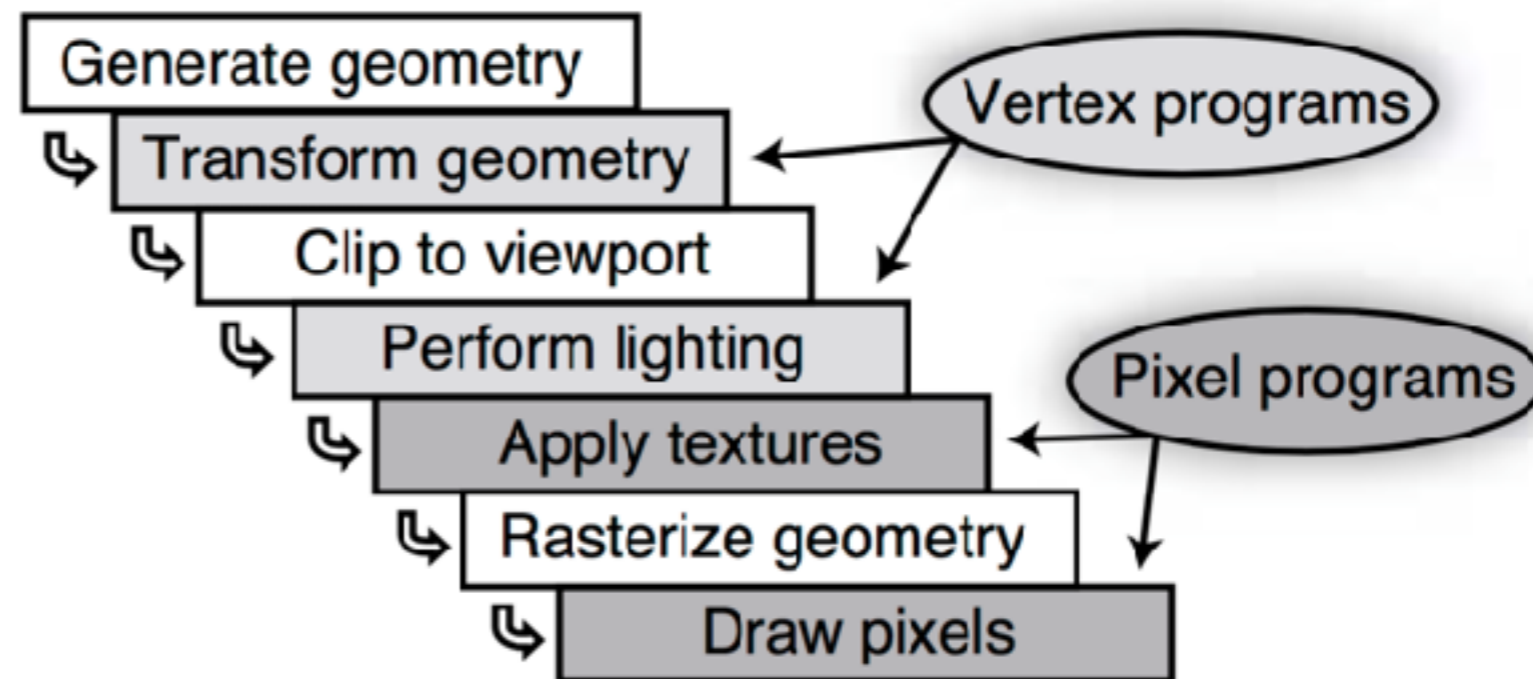


Figure 2. A programmable graphics pipeline.



2005



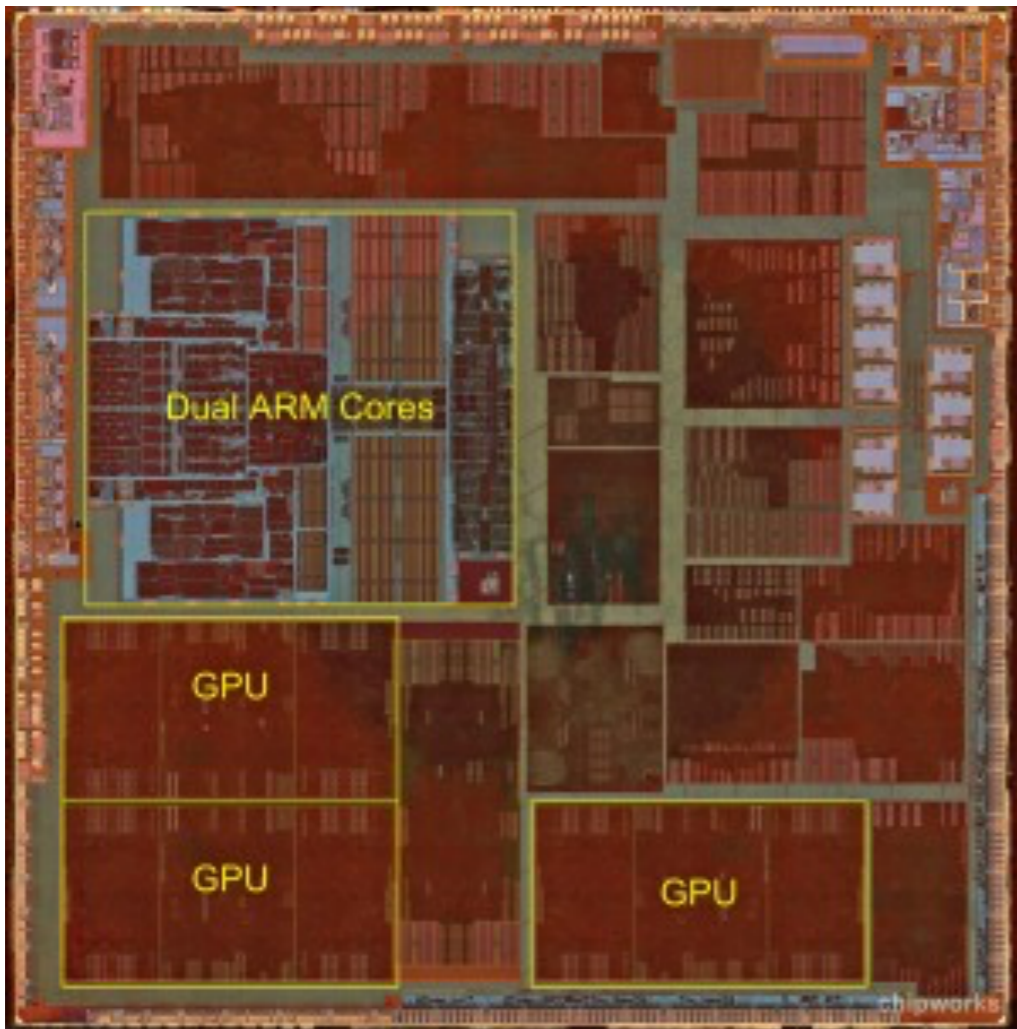


2007

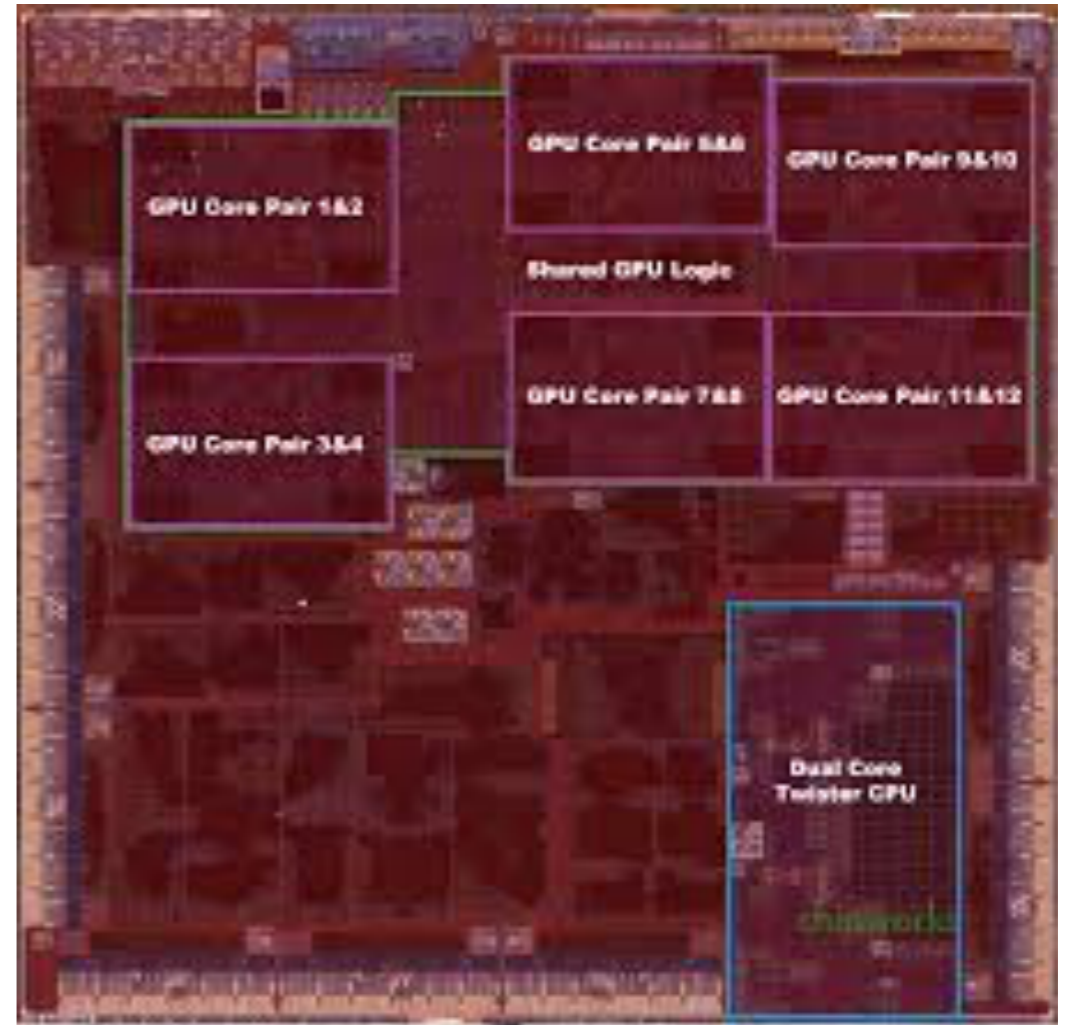


2009





A6

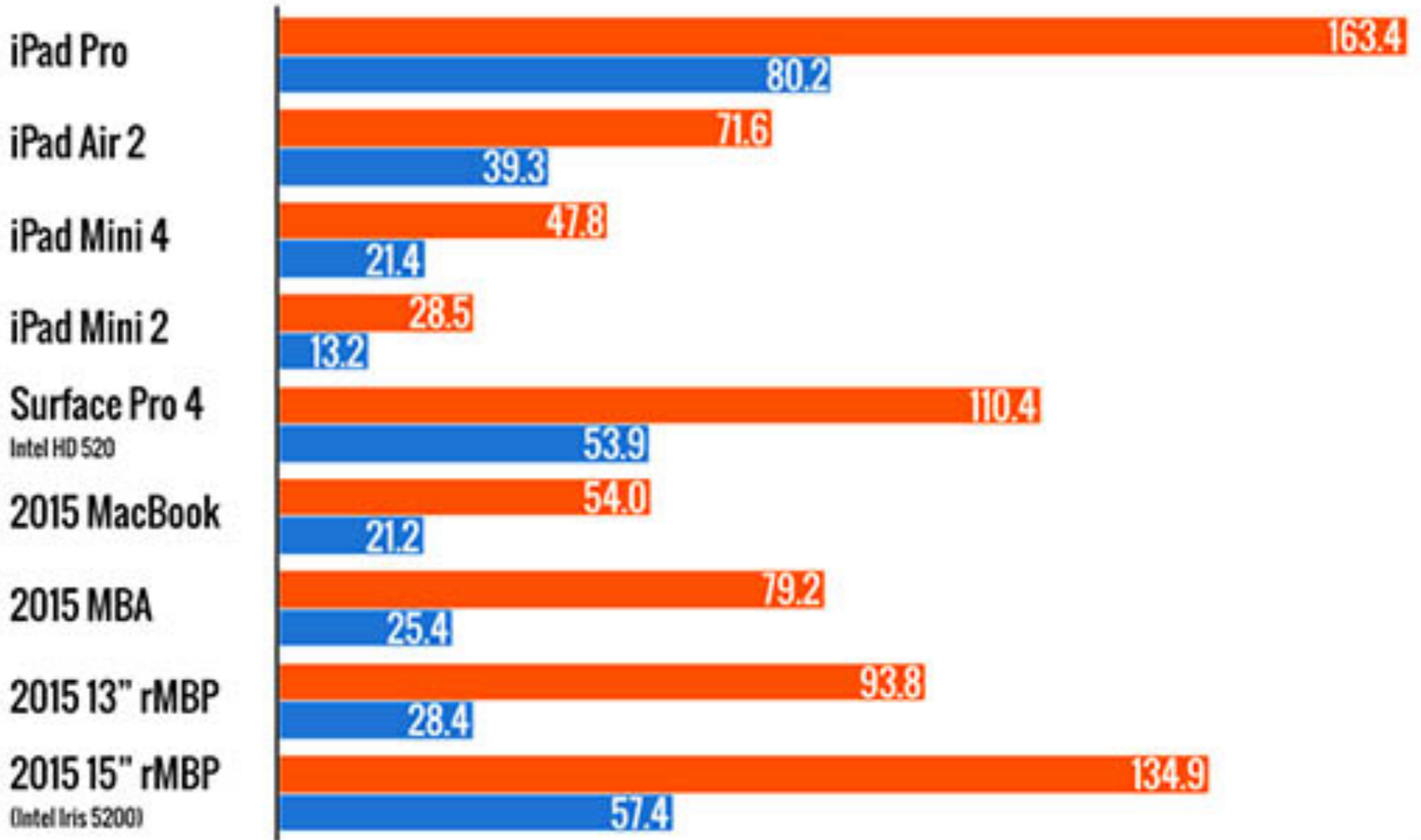


A9X

# GFXBENCH GL: OFFSCREEN

Frames per second (higher is better)

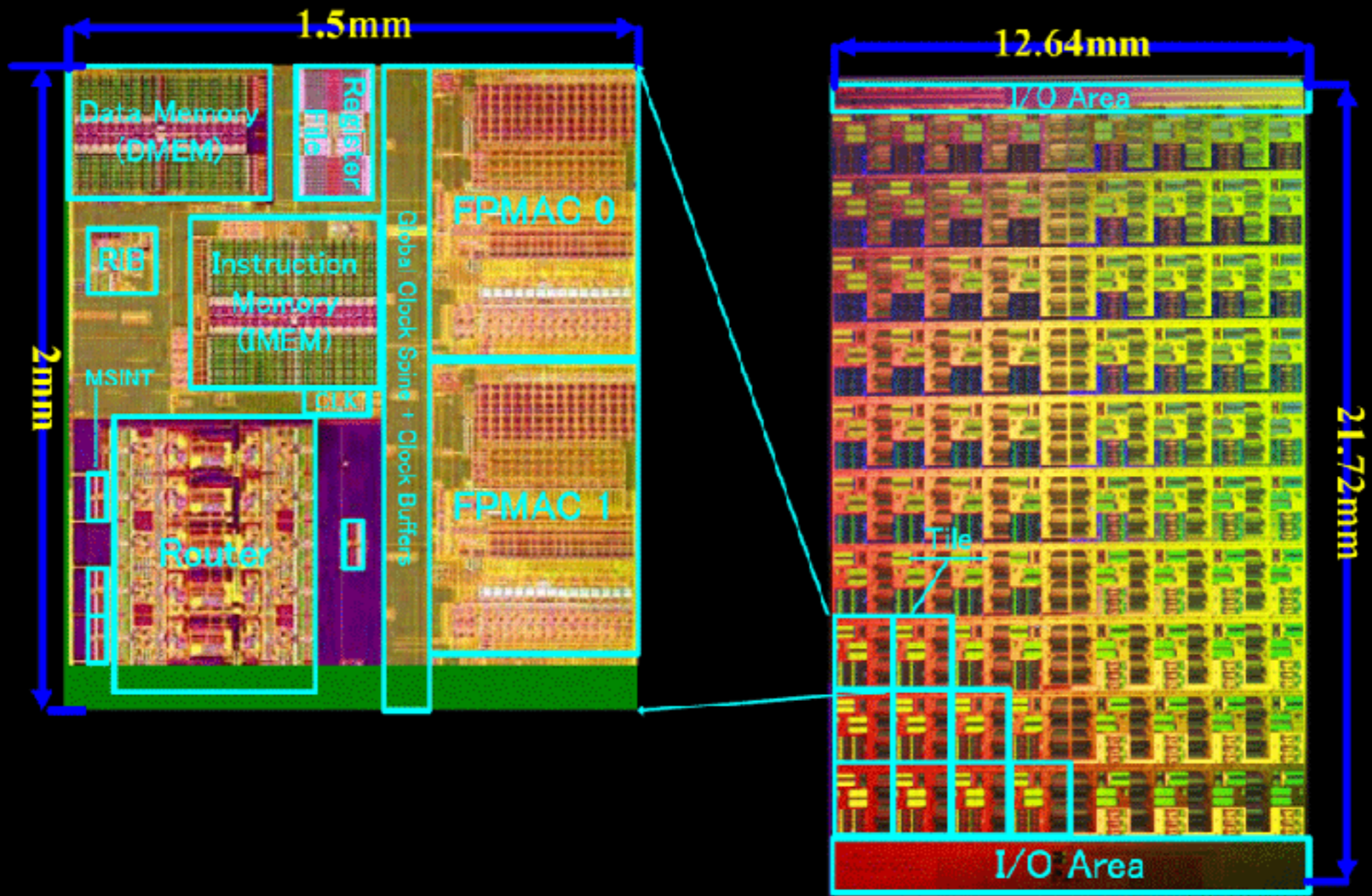
T-Rex HD  
Manhattan HD



Apple A9X



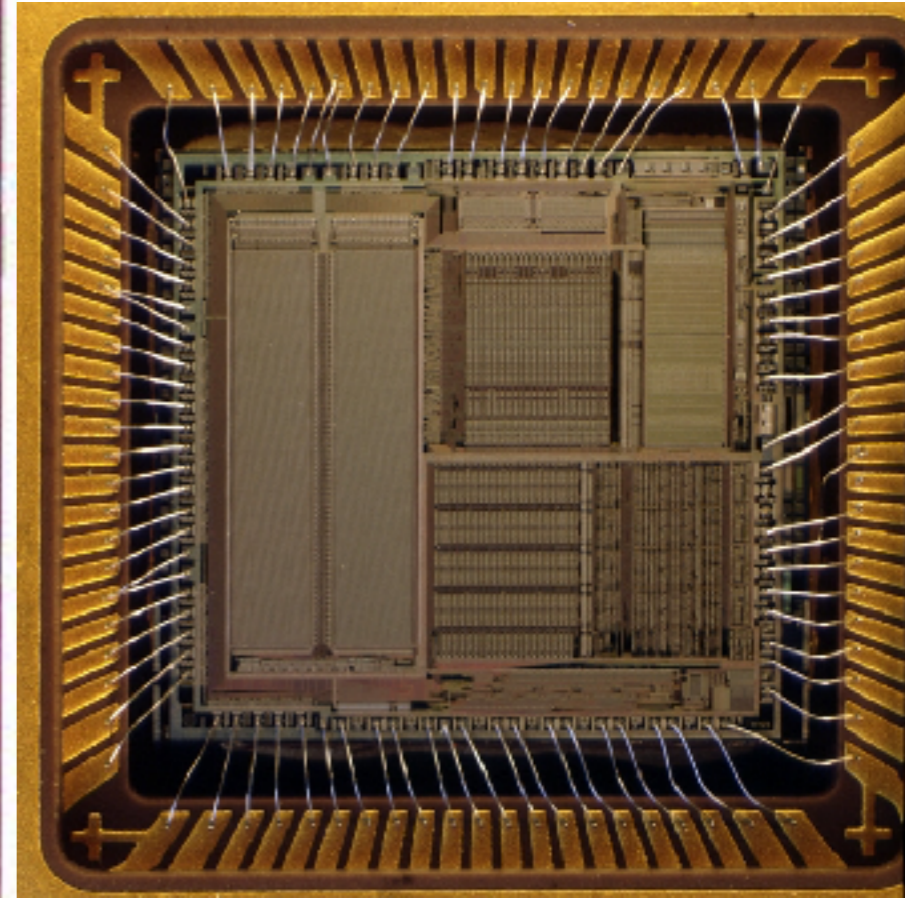
# NoC Die Overview



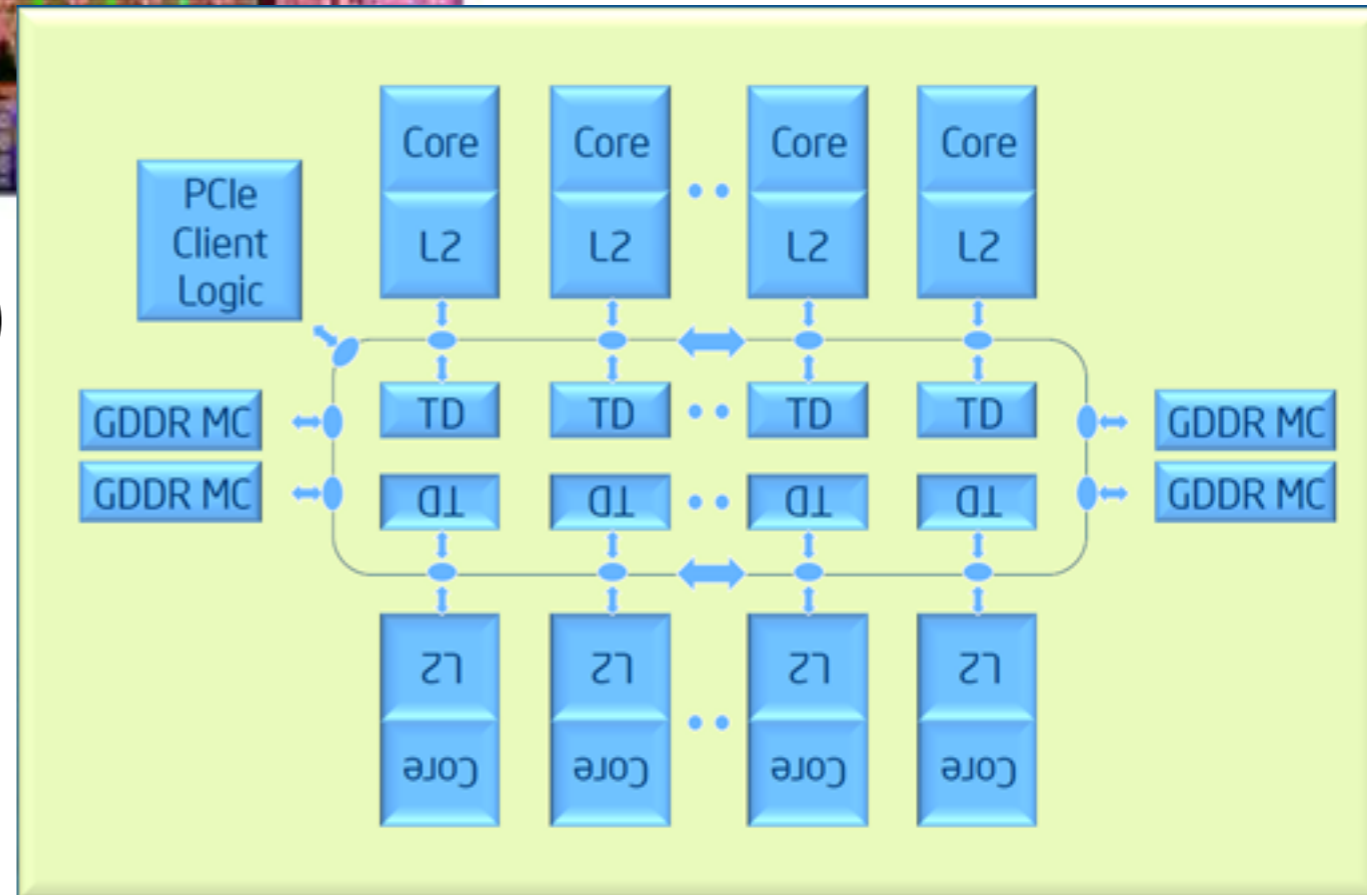
Tile	
Tile Transistors	1.2 M Transistors
Tile Area	3mm <sup>2</sup> (1.5mm x 2mm)
Router Transistors	210k Transistors
Router Area	0.34mm <sup>2</sup>

Full chip	
Process Technology	65nm CMOS
Interconnect	1 poly, 8 metal(Cu)
Transistors	100 M Transistors
Die Area	275mm <sup>2</sup> (12.64mm x 21.72mm)
C4 bumps	8390
Package	1248 pin LGA, 14 layers, 343 signal pins





(the poorly named) Intel MIC (2010)





# What are the key enabling technologies behind GPUs?

- Programmable pipeline
- Abstract ISA / API
- It's all about the memory
- High bandwidth PCIe is key for GPGPU
- Why can they use SIMD?
  - it is data parallel computation
  - the control flow is largely the same
- Need a lot of parallel tasks



Figure 3: GCN Compute Unit

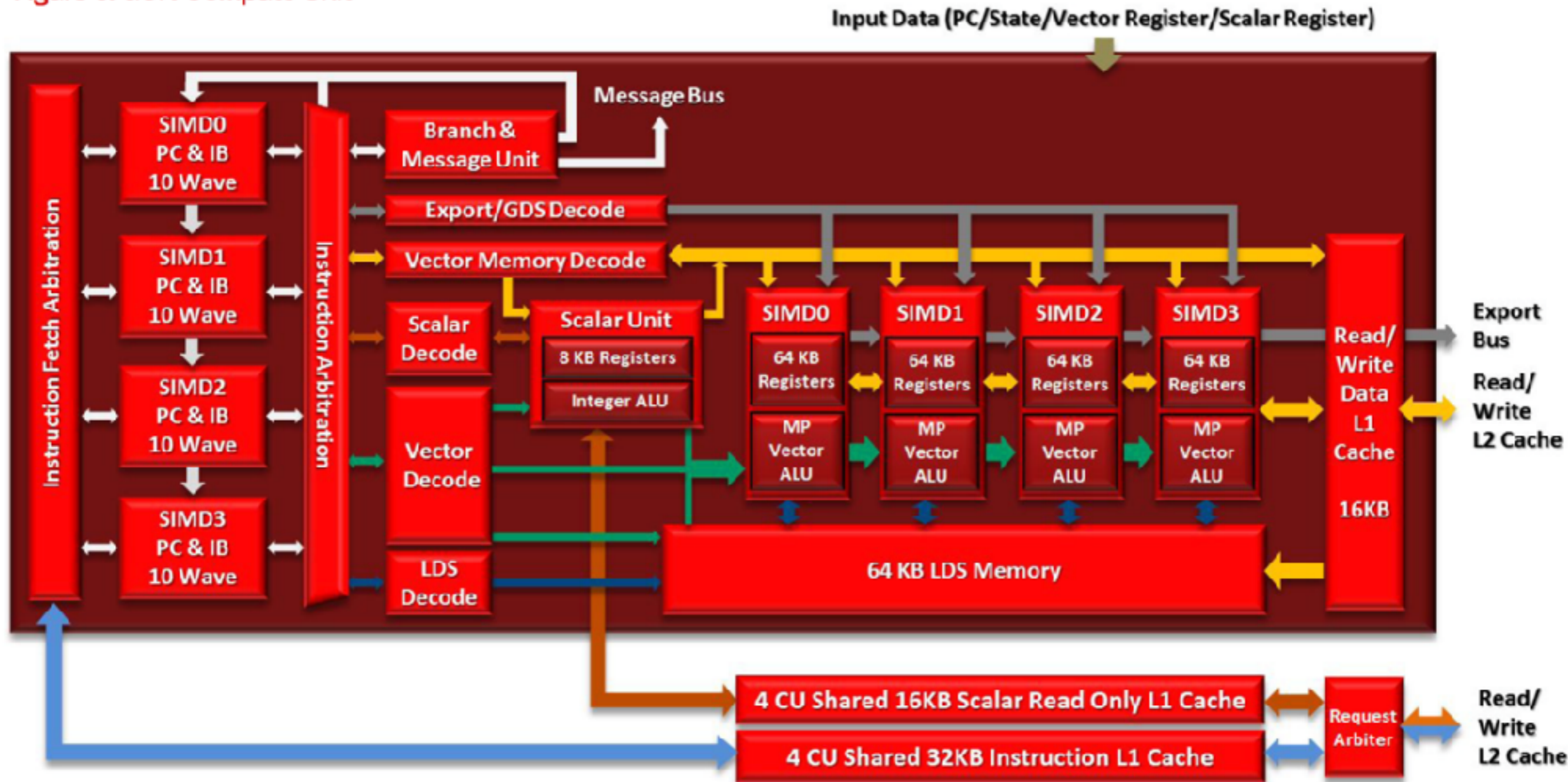




Figure 4: Local Data Share (LDS)

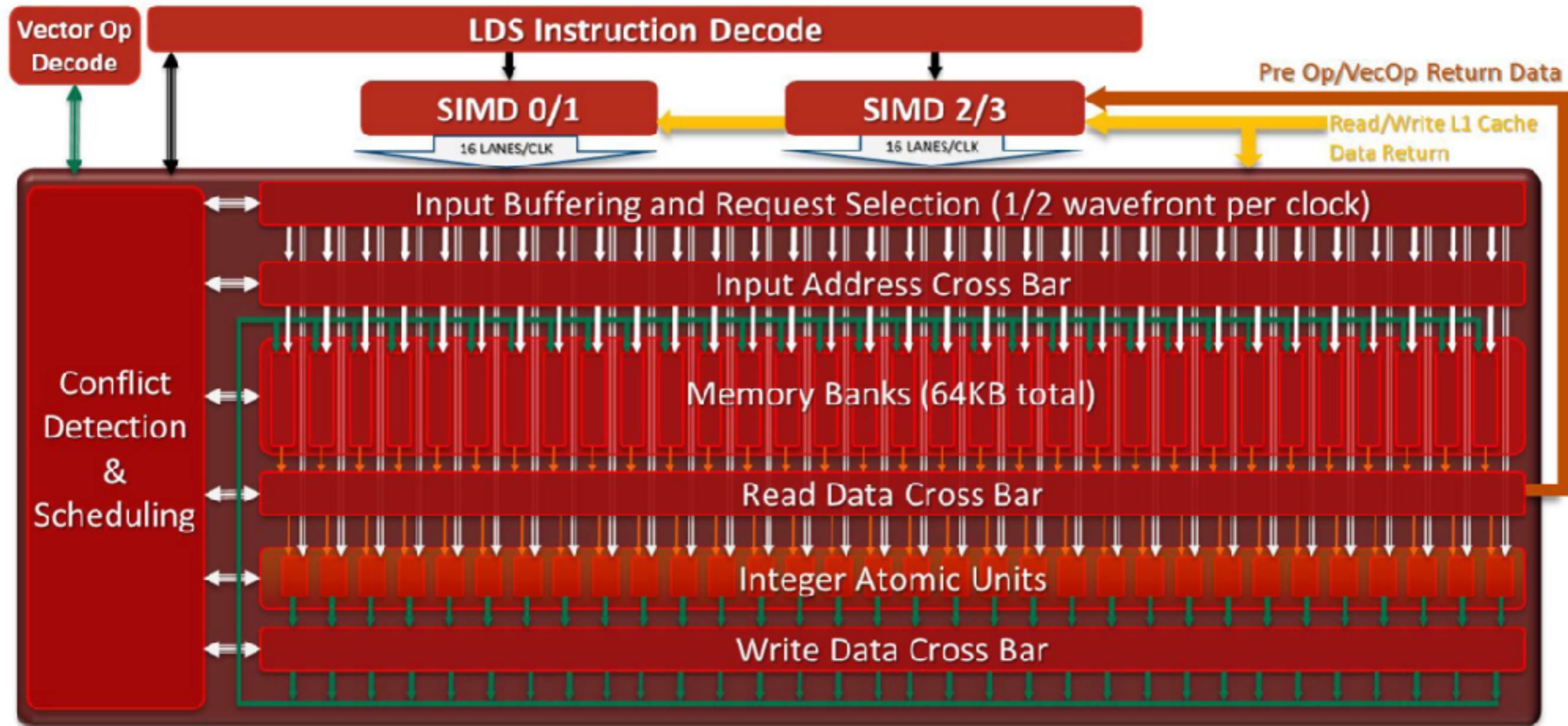




Figure 6: Cache Hierarchy

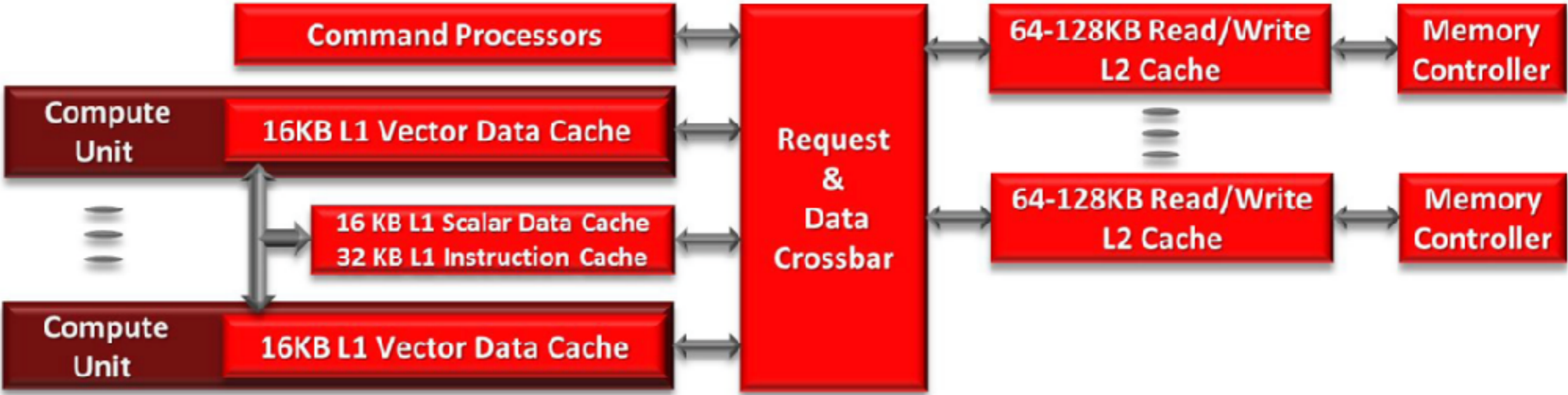
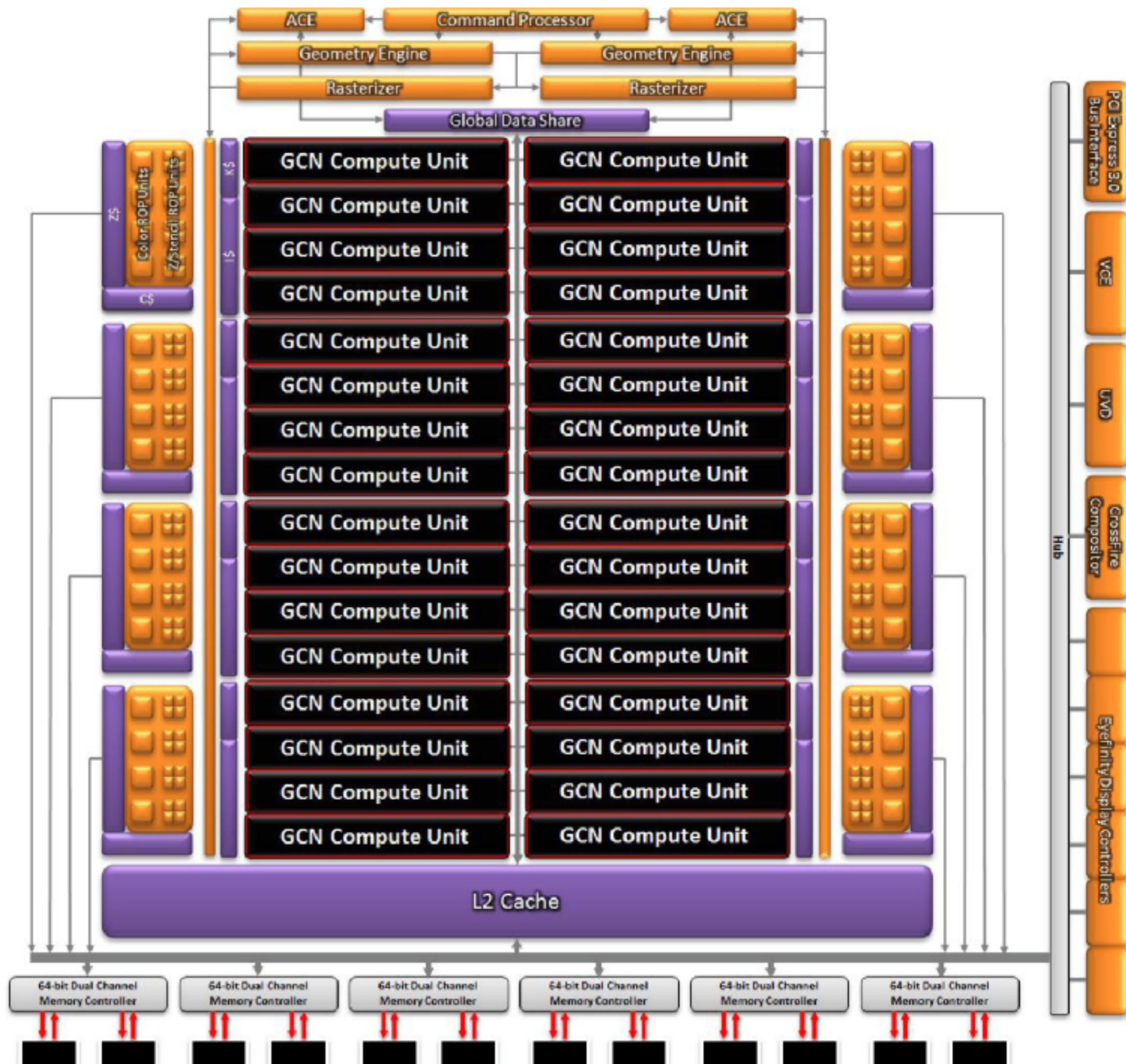
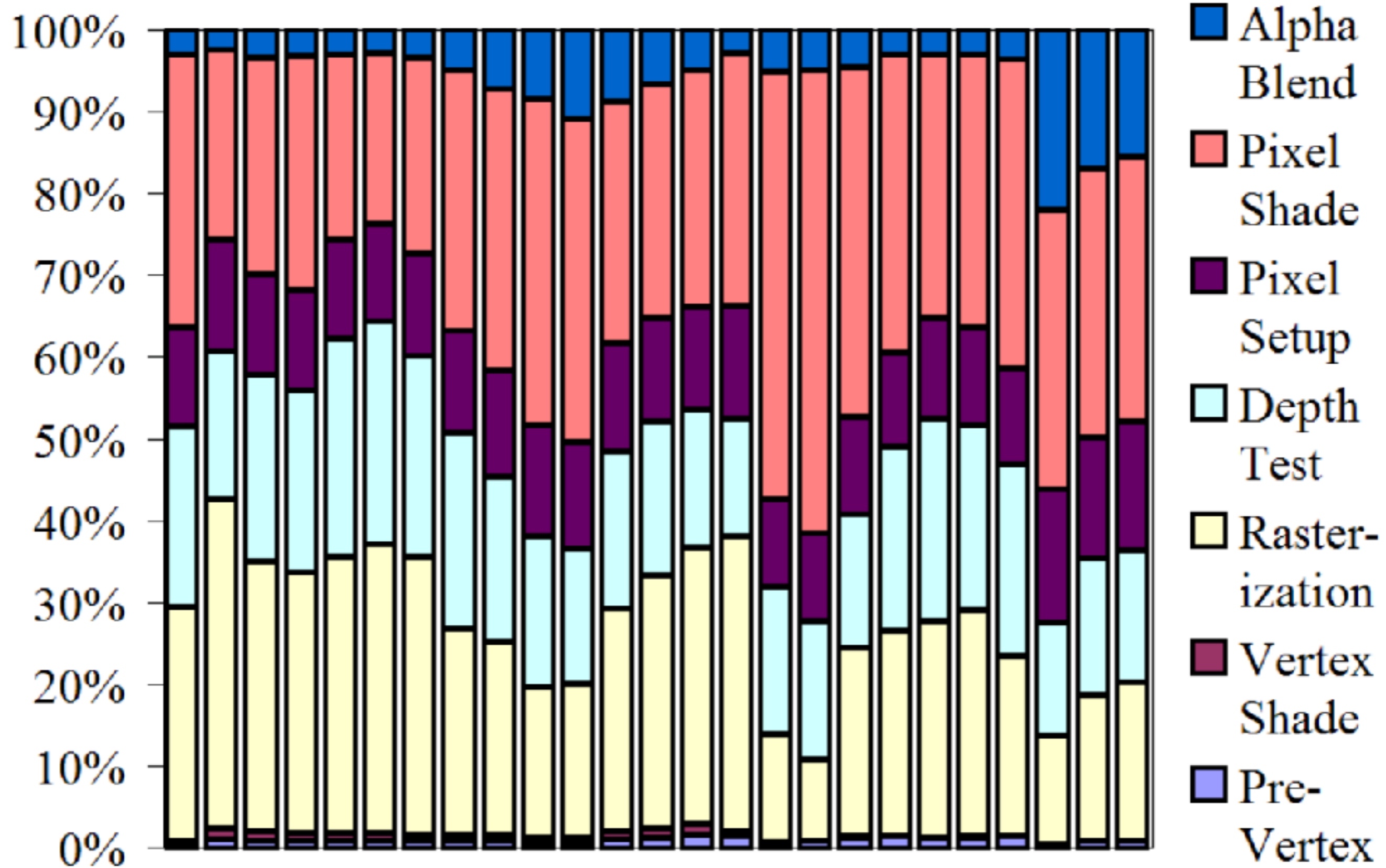




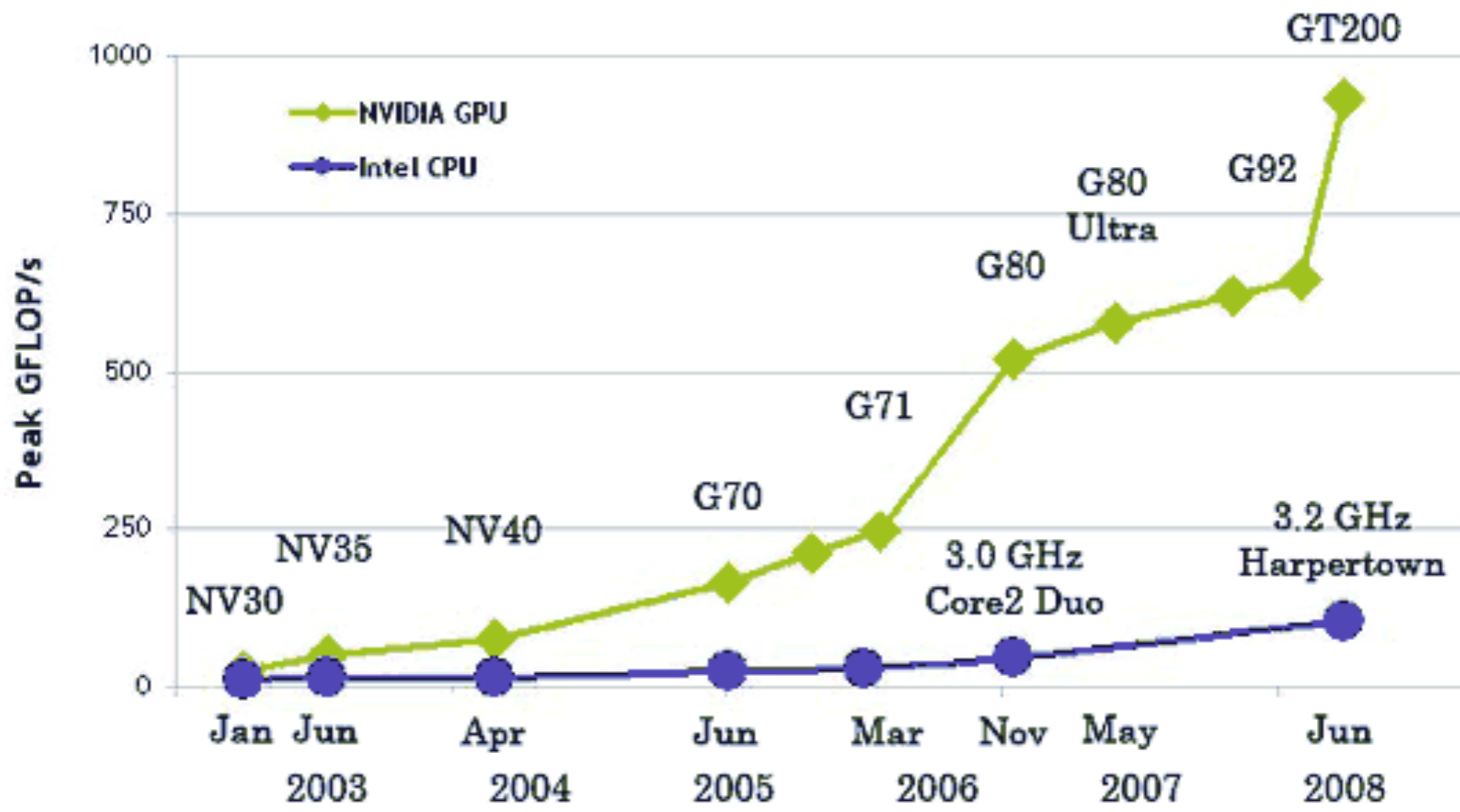
Figure 7: AMD Radeon™ HD 7970













What next?