

Numerical Computing on Graphics Hardware

임인성

서강대학교 공과대학 컴퓨터학과

최근 일반 범용 PC 에 장착되고 있는 ATI 나 NVIDIA 등의 그래픽스 가속기의 성능은 수년전과 비교할 때 비교가 안 될 정도의 빠른 속도를 자랑하고 있다. 이러한 속도 향상과 함께 급격하게 일어나고 있는 변화 중의 하나는 바로 기존의 고정된 기능의 그래픽스 파이프라인(fixed-function graphics pipeline)과는 달리 프로그래머가 가속기의 기능을 자유자재로 프로그래밍할 수 있도록 해주는 프로그래밍이 가능한 파이프라인(programmable graphics pipeline)의 출현이라 할 수 있다. 이러한 가속기에 장착되고 있는 GPU (Graphics Processing Unit)는 간단한 형태의 SIMD 프로세서라 할 수 있는데, 특히 GPU 의 한 부분인 픽셀 셰이더는 그 처리 속도가 매우 높기 때문에 이를 통하여 기존의 수치 알고리즘을 병렬화 하려는 시도가 활발히 일어나고 있다. 본 강연에서는 다양한 수치 계산을 그래픽스 가속기를 사용하여 해결하려는 시도에 대하여 간단히 살펴본다.

강연자 전자메일
ihm@sogang.ac.kr

Numerical Computing on Graphics Hardware

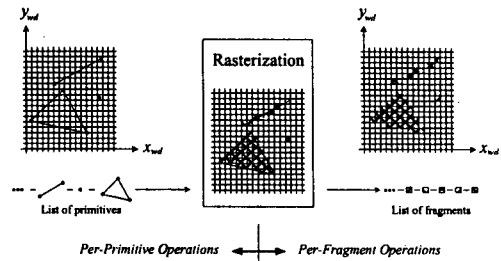
2004. 4. 23

서강대학교 공과대학 컴퓨터학과
임인성

이메일: jinlim@scnu.ac.kr

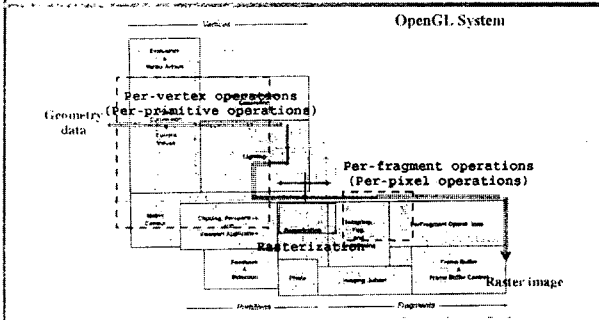
서강대학교 컴퓨터 그래픽스 연구실

Real-Time Rendering Pipeline on GPU



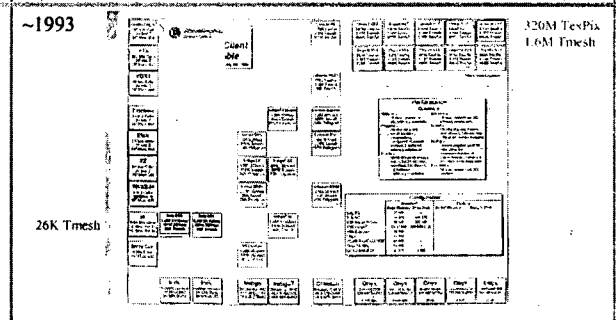
서강대학교 컴퓨터 그래픽스 연구실

Fixed-Function Graphics Pipeline



서강대학교 컴퓨터 그래픽스 연구실

Advances in Graphics Processors



서강대학교 컴퓨터 그래픽스 연구실

2000
~ 2001

	Xbox	Sony PlayStation 2	Nintendo Game Cube
Graphics Processor	250 MHz Custom-designed chip by MS and NVIDIA	147.456 MHz	202.5 MHz Custom chip "Flipper"
Total memory	64 MB	32 MB	43 MB
Memory Bandwidth	6.4 GB/s	3.2 GB/s	3.2 GB/s
Polygon Performance	125 M/s	66 M/s	6-12M/s
Simultaneous Textures	4	1	N/A
Pixel Fill Rate - No Texture	4.0 G/s	2.4 G/s	N/A
Pixel Fill Rate - 1 Texture	4.0 G/s	1.2 G/s	N/A
Pixel Fill Rate - 2 Texture	4.0 G/s	0.6 G/s	N/A

서강대학교 컴퓨터 그래픽스 연구실

2002
~ 2003

	ATI Radeon 9700 PRO	NVIDIA GeForce4 Ti6600	NVIDIA GeForceFX 5800	NVIDIA GeForceFX 5800 Ultra
Chip Technology	150nm	150nm	150nm	150nm
Transistors	~107 Million	~128 Million	~125 Million	~125 Million
Memory Bus	256-bit DDR	128-bit DDR	128-bit DDR2	128-bit DDR2
Memory Bandwidth	19.8 GB/s	10.4 GB/s	12.8 GB/s	16 GB/s
Pixel Fill Rate	2.6 Gpixels/s	1.24 Gpixels/s	~3.2 Gpixels/s	~4 Gpixels/s
Anti-Aliased Fill Rate	15.6 Billion AA Samples/s	4.8 Billion AA Samples/s	~12.8 Billion AA Samples/s	~16 Billion AA Samples/s
Triangle Transform Rate	325M Triangles/s	69M Triangles/s	280M Triangles/s	350M Triangles/s
Graphics Memory	128/256MB	128MB	128/256MB	128/256MB
GPU Clock	325 MHz	300 MHz	~400 MHz	~500 MHz
Memory Clock	310 MHz (620 DDR)	325 MHz (650 DDR)	400 MHz (800 DDR2)	500 MHz (1000 DDR2)
Textures per Texture Unit	8	4	16	16

서강대학교 컴퓨터 그래픽스 연구실

Programmable GPU(Graphics Processing Unit)

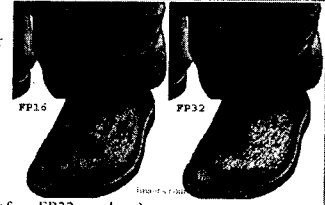
- ❖ GPUs have become programmable!
 - NVIDIA GeForce FX and ATI Radeon 9700
- ❖ The recent GPUs support 4-wide full floating-point SIMD operations.
 - 4-component vector in FP32 → 128bit processing
- ❖ In particular, fragment(pixel) shader hardware supports high pixel rates through its flexible texturing and blending units.



서강대학교 컴퓨터 그래픽스 연구실

Some Words on Latest Programmable GPUs

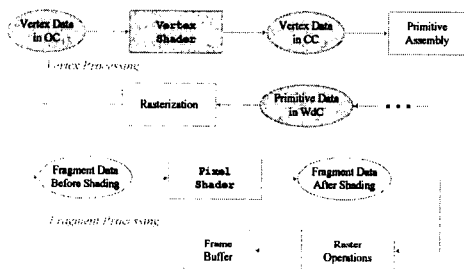
- ❖ Speed: faster than ever
- ❖ Graphics memory: larger than ever
- ❖ Bandwidth: wider than ever



- ❖ 4-wide SIMD processing
- ❖ 128-bit vertex processing pipeline (four FP32 numbers)
- ❖ 128-bit pixel processing pipeline (four FP32 numbers)
- ❖ 4 parallel geometry pipelines
- ❖ 8 parallel pixel pipelines
- ❖ Programmability

서강대학교 컴퓨터 그래픽스 연구실

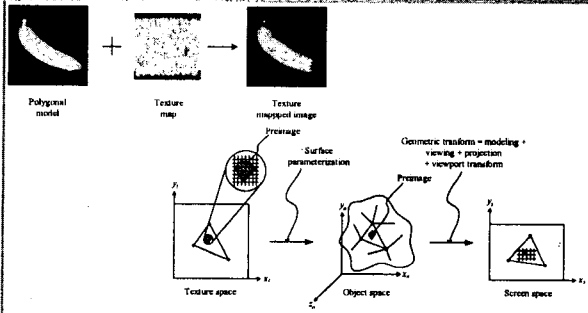
(Current) Programmable Graphics Pipeline



서강대학교 컴퓨터 그래픽스 연구실

9

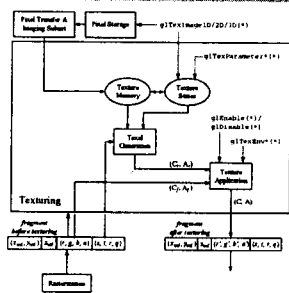
Pixel Shader and Texture Mapping



서강대학교 컴퓨터 그래픽스 연구실

10

Conventional Texture Mapping in OpenGL



서강대학교 컴퓨터 그래픽스 연구실

11

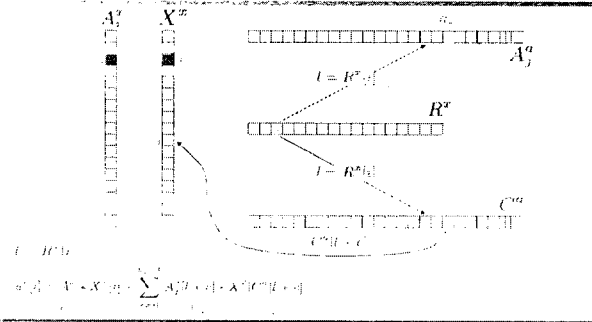
Numerical Computing on GPU

- ❖ Parallelization of numerical algorithms on multi-processor architectures has been an active research area in applied numerical analysis.
- ❖ Many physically-based computer graphics (and CFD) applications require high-intensity numerical solvers.
- ❖ Various numerical algorithms are being actively mapped onto GPUs by graphics people.

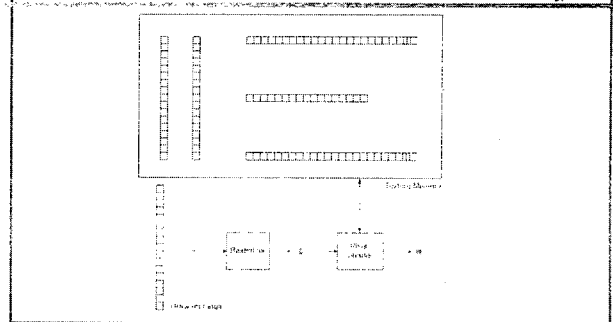
서강대학교 컴퓨터 그래픽스 연구실

12

Texture Layout and Indirection



Mapping onto Pixel Shader

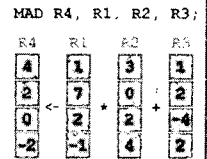


SIMD Optimization of Linear Expressions [Bajaj, Ihm, Oh, Min??]

- Minimization of 4-wide SIMD instructions for computing $y = Ax + b$ through matrix transformation
- Abstract SIMD shader model
 - The shader supports four-wide SIMD parallelism.
 - Its instruction set includes the instructions shown in the table.
 - Any component of the source registers may swizzle and/or replicate into any other component. Furthermore, destination registers may be masked. These register modifiers do not harm shader performance.
 - Every instruction executes in a single clock cycle. Hence, the number of instructions in a vertex shader program is the major factor affecting shader performance.

Instruction Set of the SIMD Model

Operation	Usage	Description
ADD	ADD D, S0, S1	$D \leftarrow S0 + S1$
MUL	MUL D, S0, S1	$D \leftarrow S0 * S1$
MAD	MAD D, S0, S1, S2	$D \leftarrow S0 * S1 + S2$
DP4	DP4 D, S0, S1	$D \leftarrow S0 \cdot S1$
MOV	MOV D, S	$D \leftarrow S$



Partition of a Linear Expression

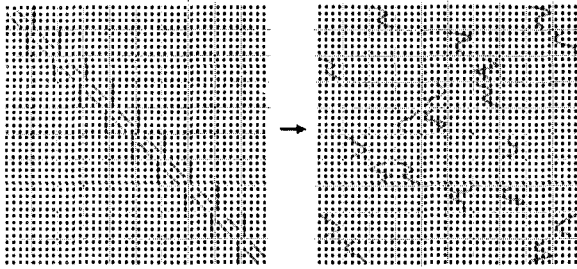
$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_{m-2} \\ a_{m-1} \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} A_0^0 & A_0^1 & A_0^2 & A_0^3 & A_0^4 & \dots & A_0^{m-2} & A_0^{m-1} & 0 & 0 \\ A_1^0 & A_1^1 & A_1^2 & A_1^3 & A_1^4 & \dots & A_1^{m-2} & A_1^{m-1} & 0 & 0 \\ A_2^0 & A_2^1 & A_2^2 & A_2^3 & A_2^4 & \dots & A_2^{m-2} & A_2^{m-1} & 0 & 0 \\ A_3^0 & A_3^1 & A_3^2 & A_3^3 & A_3^4 & \dots & A_3^{m-2} & A_3^{m-1} & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ A_{m-2}^0 & A_{m-2}^1 & A_{m-2}^2 & A_{m-2}^3 & A_{m-2}^4 & \dots & A_{m-2}^{m-2} & A_{m-2}^{m-1} & 0 & 0 \\ A_{m-1}^0 & A_{m-1}^1 & A_{m-1}^2 & A_{m-1}^3 & A_{m-1}^4 & \dots & A_{m-1}^{m-2} & A_{m-1}^{m-1} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} V^0 \\ V^1 \\ V^2 \\ V^3 \\ \vdots \\ V^{m-2} \\ V^{m-1} \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} A_0^m \\ A_1^m \\ A_2^m \\ A_3^m \\ \vdots \\ A_{m-2}^m \\ A_{m-1}^m \\ 0 \\ 0 \end{bmatrix}$$

Decomposition

$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \\ \vdots \\ a_{m-2} \\ a_{m-1} \\ 0 \\ 0 \end{bmatrix} = M_{0,0} \cdot \begin{bmatrix} V^0 \\ V^1 \\ V^2 \\ V^3 \\ \vdots \\ V^{m-2} \\ V^{m-1} \\ 0 \\ 0 \end{bmatrix} + M_{0,1} \cdot \begin{bmatrix} V^4 \\ \vdots \\ V^7 \end{bmatrix} + \dots + M_{0,m-1} \cdot \begin{bmatrix} V^{m-2} \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} A_0^m \\ A_1^m \\ A_2^m \\ A_3^m \\ \vdots \\ A_{m-2}^m \\ A_{m-1}^m \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} a_{m-2} \\ a_{m-1} \\ 0 \\ 0 \end{bmatrix} = M_{m-1,0} \cdot \begin{bmatrix} V^0 \\ V^1 \\ V^2 \\ V^3 \\ \vdots \\ V^{m-2} \\ V^{m-1} \\ 0 \\ 0 \end{bmatrix} + M_{m-1,1} \cdot \begin{bmatrix} V^4 \\ \vdots \\ V^7 \end{bmatrix} + \dots + M_{m-1,m-1} \cdot \begin{bmatrix} V^{m-2} \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} A_{m-2}^m \\ A_{m-1}^m \\ 0 \\ 0 \end{bmatrix}$$

Solving 260,000 X 260,000 Poisson Equation on GPU



Timing Performance



		tolerance		
		10^{-3}	10^{-4}	10^{-5}
GS	time	101.6	364.8	1178.6
on CPU	iter	16	61	200
EBJ	time	62.8	142.4	419.1
on GPU	iter	20	84	291
EBJ	time	49.1 (9.8)	111.6 (68.3)	315.9 (272.6)
on GPU (optimized)	iter	19	77	253