



GPU Computing with fragment shaders

”Classic GPGPU”

Use graphics shaders for general-purpose computing.

Adapt your data and computing to fit the graphics pipeline.

Hot until CUDA arrived, now overshadowed by CUDA and OpenCL.



Why is classic GPGPU interesting?

- **Highly suited to all problems dealing with images, computer vision, image coding etc**
- **Parallelization "comes natural", you can't avoid it and good speedups are likely. Fewer pitfalls.**
- **Highly optimized (for graphics performance).**
 - **Compatibility is vastly superior!**
 - **Very much easier to install!**



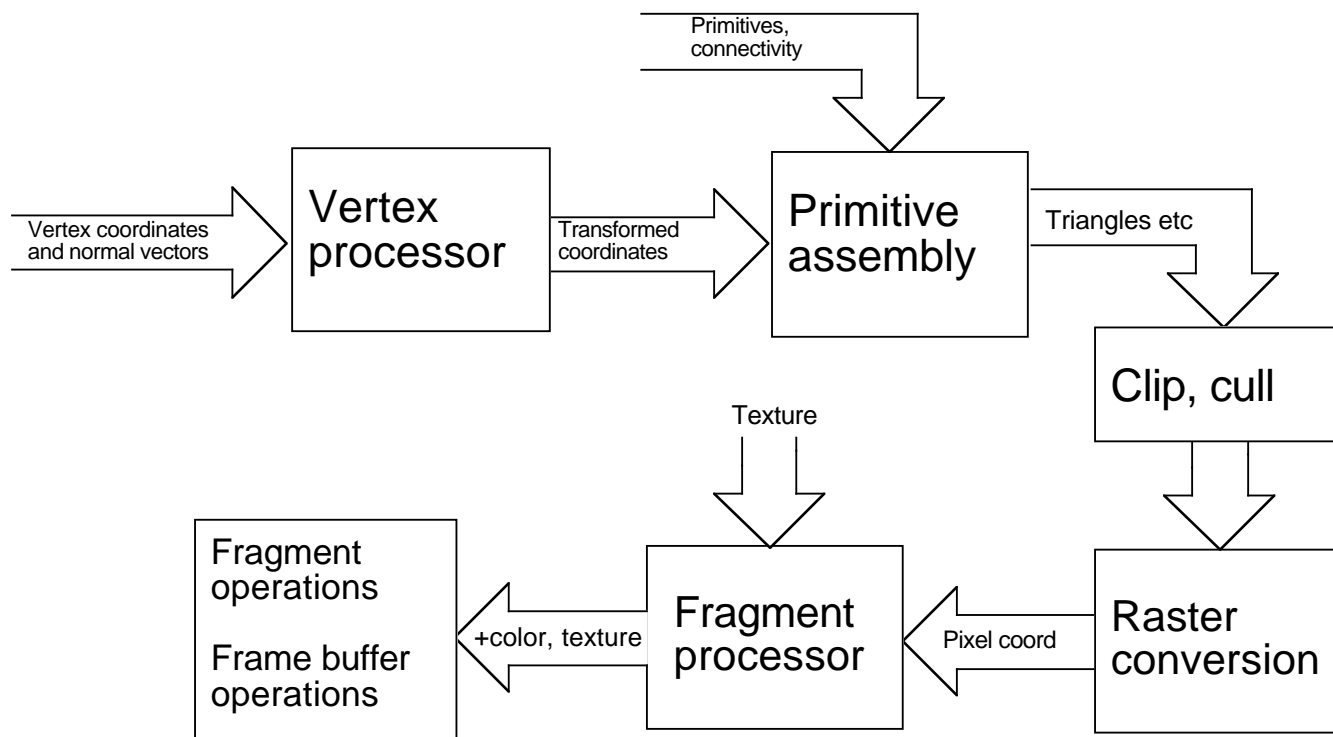
So what is not so good?

- Must map data to image data
- Computing controlled by pixels in output image
- No shared memory access (except for compute shaders)

However: OpenGL 4 adds much flexibility, moves closer to CUDA and (especially) OpenCL. Writable textures, atomics, synchronization...

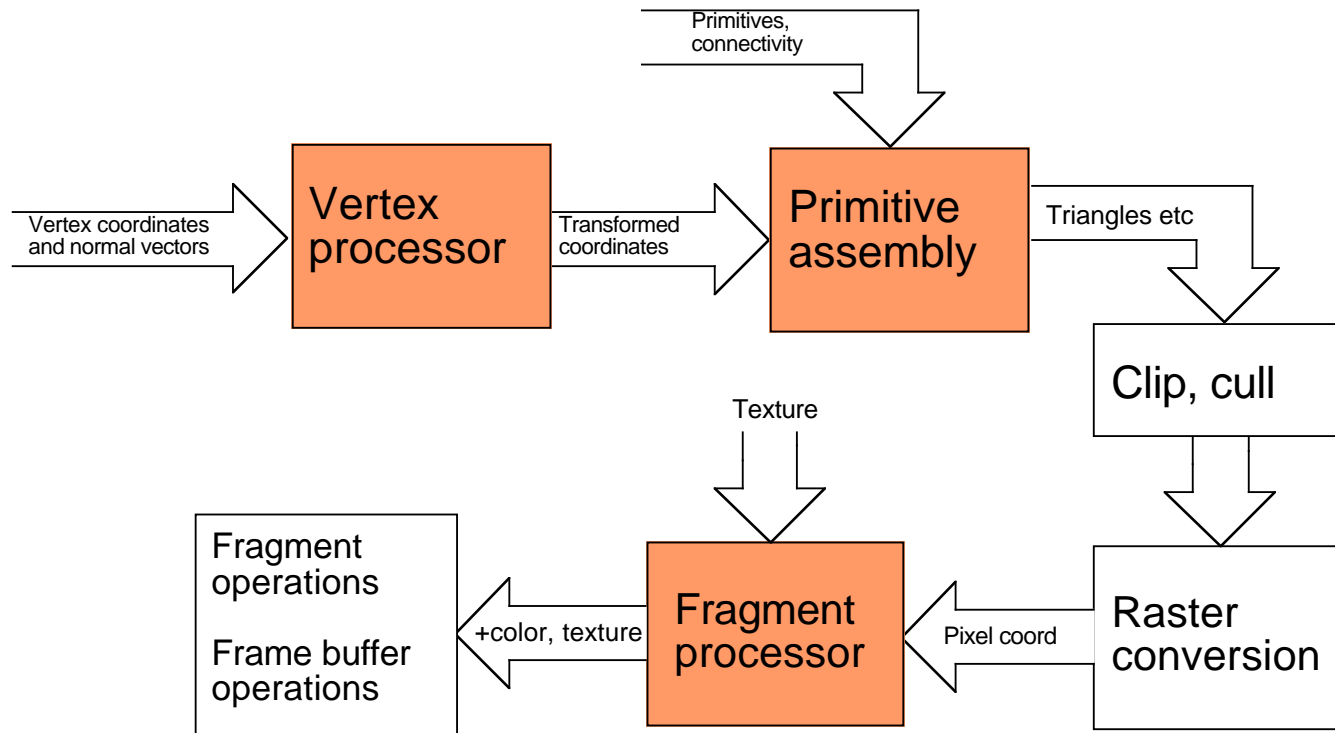


The OpenGL pipeline



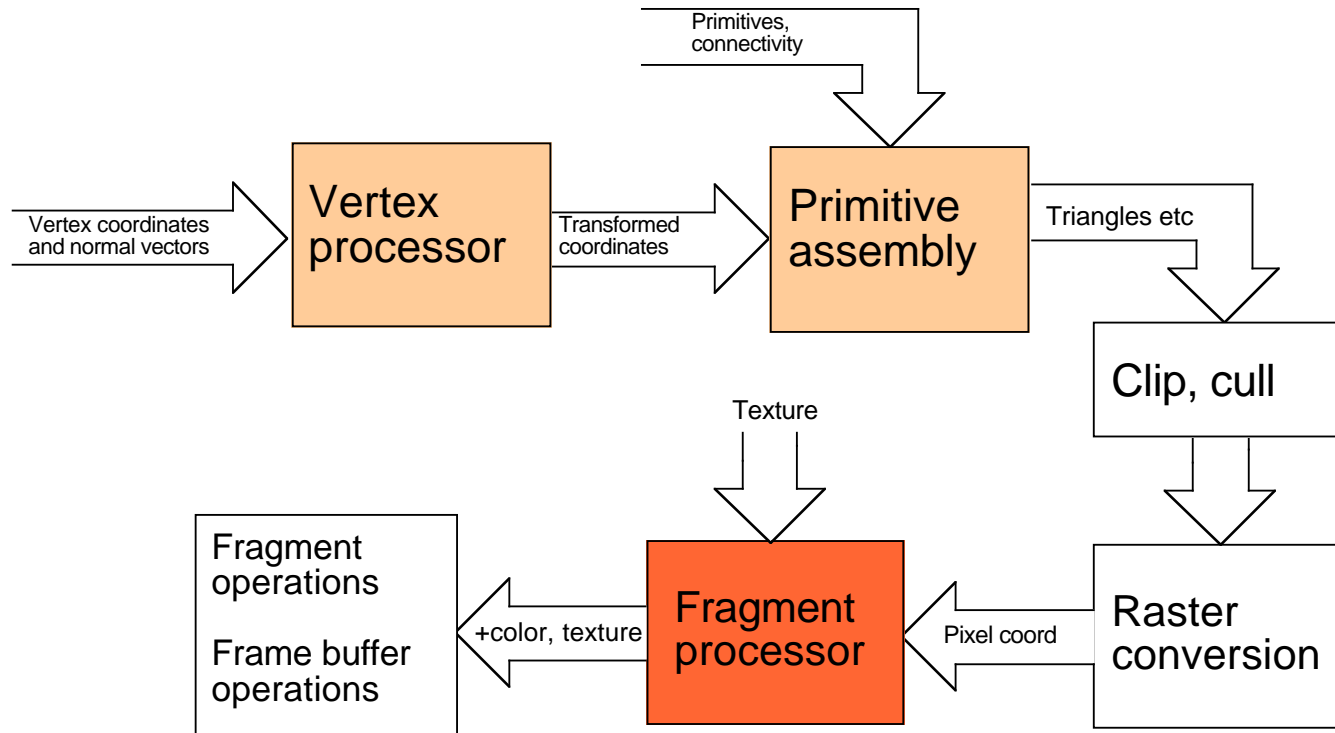


Out of these, three are programmable!





But only one creates easily accessible output data!





GPGPU

Problem:

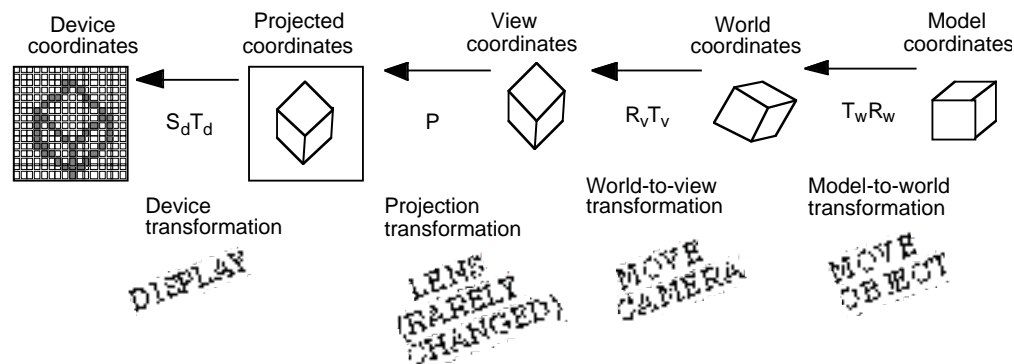
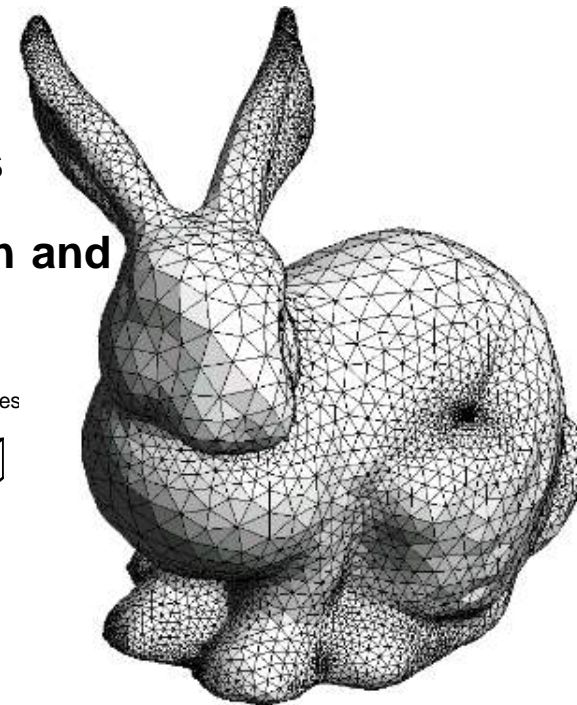
- Algorithms must be parallelized - more than with CUDA. No intermediate results from neighbors can be used. (Some possibilities with GL4.)
- No access to shared memory. (But access to constant memory and easy access to texture memory.)

Does it pay to use shaders for GPU Computing?



Typical OpenGL situation

- Complex geometry
- Many transformations
- Perspective projection
- Lighting and material calculations for the surfaces
- Many texture accesses for interpolation and supersampling



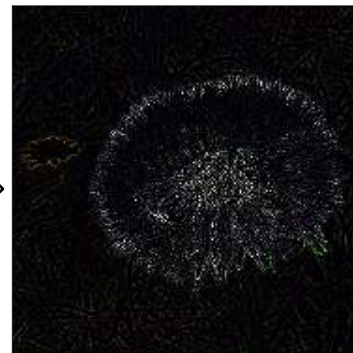
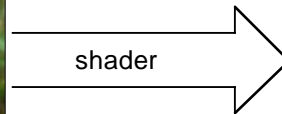


Typical GPGPU processing (also used in filtering in graphics):

- Render to a single rectangle covering the entire image buffer.
- Use FBOs for effective feedback
- Floating-point buffers
- Ping-ponging, many pass with different shaders



Render image 1:1



Output



The GPGPU/shaders model

- Array of input data = texture
- Array of output data = resulting frame buffer
- Computation kernel = shader
- Computation = rendering
- Feedback = switch between FBO's or copy frame buffer to texture



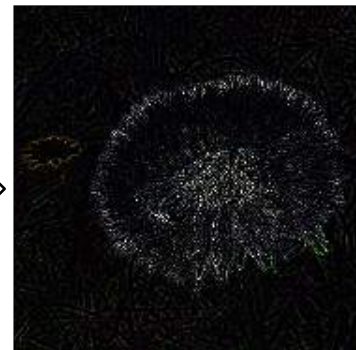
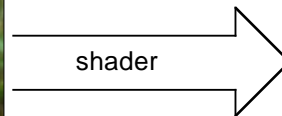
Computation = rendering

Typical situation:

- **Texture and frame buffer same size**
- **Render the polygon over the entire frame buffer**



Texture



Frame buffer



Kernel = shader

Shaders are read and compiled to one or more program objects. A GPGPU application can use several shaders in conjunction!

Activate desired shader as needed using `glUseProgram()`;

The fragment shader performs the computation:

```
uniform sampler2D texUnit;  
in vec2 texCoord;  
out vec4 fragColor;  
  
void main(void)  
{  
    vec4 texVal = texture(texUnit, texCoord);  
    fragColor = sqrt(texVal);  
}
```

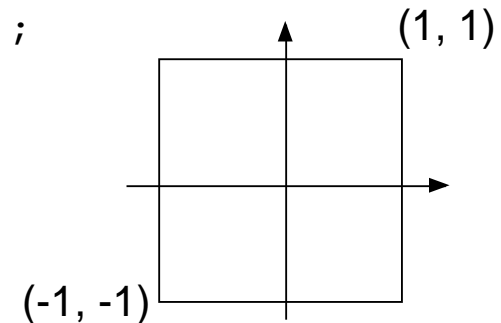


Render a single polygon

- Texture and frame buffer same size
- Render polygon over entire frame buffer

```
GLfloat quadVertices[] = { -1.0f, -1.0f, 0.0f,  
                           -1.0f, 1.0f, 0.0f,  
                           1.0f, 1.0f, 0.0f,  
                           1.0f, -1.0f, 0.0f};
```

```
GLuint quadIndices[] = {0, 1, 2, 0, 2, 3};
```





Program structure:

- Set up OpenGL
- Upload data to texture
- Load shaders from file and compile
- Draw quad on screen using OpenGL
- Data is computed by the fragment shader, per pixel
- Output can be downloaded as image data

Examples...



Feedback

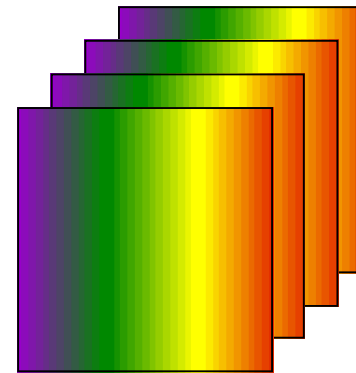
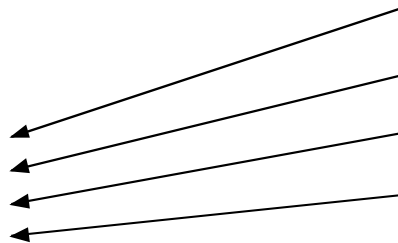
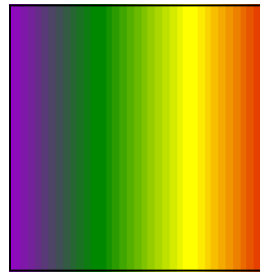
We must be able to pass output from one operation as input of the next!

Solution: Render to texture, "framebuffer objects", create a texture used as input for a later stage



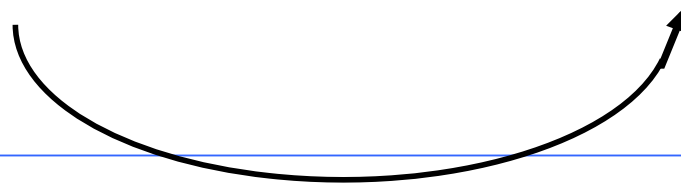
“Ping-pong”-ing

The kernel reads from one or more texture, writes into the frame buffer



Using “framebuffer objects” the output image can be a texture

Input data is a number of textures. Limited by the number of texturing units available.





Filtering, convolution

Common problem, highly suited for shaders.

All kinds of linear filters:

- **Low-pass filtering (smoothing)**
 - **Gradient, embossing**

Must be done by gather operations, not scatter!



3x1 filter

```
uniform sampler2D texUnit;  
uniform float texSize;  
in vec2 texCoord;  
out vec4 fragColor;
```

```
void main(void)  
{  
    float offset = 1.0 / texSize;  
    vec2 tc = texCoord;  
    vec4 c = texture(texUnit, tc);  
    tc.x = tc.x + offset;  
    vec4 l = texture(texUnit, tc);  
    tc.x = tc.x - 2.0*offset;  
    vec4 r = texture(texUnit, tc);  
    tc.x = tc.x - offset;  
    fragColor = (c + c + l + r) * 0.25;  
}
```

1	2	1
---	---	---

More graphics heritage:
Index data by steps of
1/size, not 1!



Separable filters

1	4	6	4	1
4	16	24	16	4
6	24	36	24	6
4	16	24	16	4
1	4	6	4	1

$$=$$

1	2	1
---	---	---

$$\otimes$$

1	2	1
---	---	---

$$\otimes$$

1
2
1

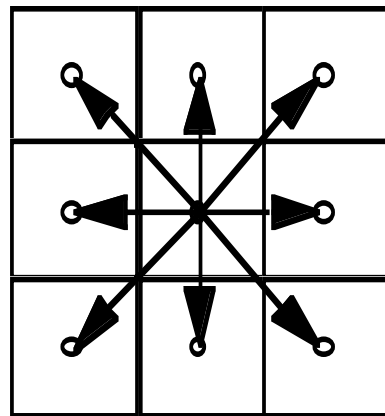
$$\otimes$$

1
2
1

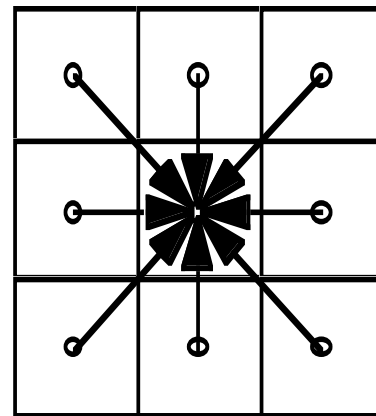
**Implemented as ping-ponging passes.
Optimization possibilities!**



Scatter vs gather



Scatter



Gather

Shaders give output for *one* pixel -> gather only!



How about CUDA/OpenCL?

Scatter vs gather: You usually prefer gather. Less synchronization! (Remember, synchronization comes for a cost!)

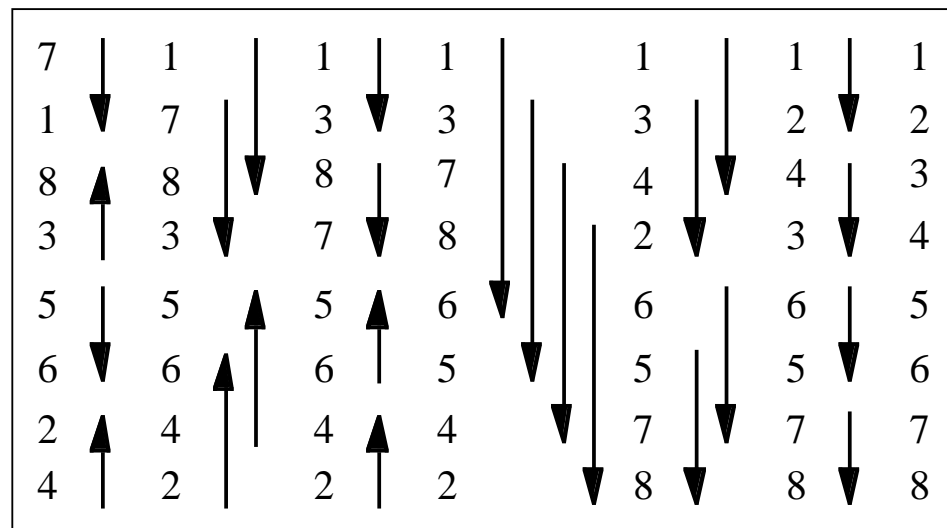
Separable filters: Optimization just as valid for all techniques! (But particularly common in shaders, for images.)



Sorting

QuickSort hard to implement in shaders

Bitonic Sort fits shaders well (see earlier lectures)



One
rendering
pass per
stage!



Reduction

Reduction algorithms are implemented by a ping-ponging pyramid

Maximum, minimum, global average...

Output smaller than input

47	2	3	57	5	12	7	8
10	20	6	13	14	15	16	17
19	11	21	22	23	68	25	26
38	29	64	31	32	33	35	34
37	28	39	49	53	42	41	52
46	1	48	40	61	51	44	43
55	71	4	58	69	62	50	60
30	65	66	67	24	59	70	56

→

47	57	15	17
38	64	68	35
46	49	61	52
71	67	69	70

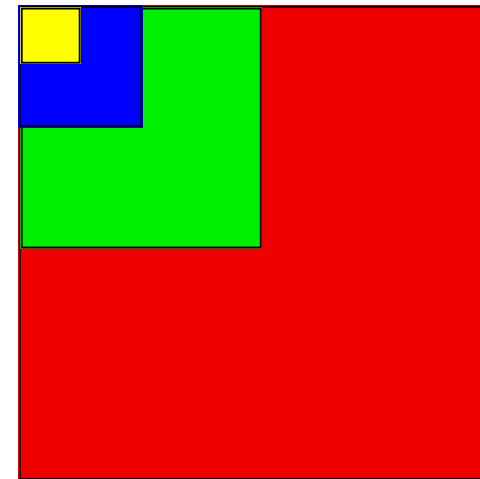
(Images by Dominik Göttsche)



Reduction

- 1) Texture pyramid, typically 2x2
- 2) Constant texture size, use smaller and smaller parts of the texture!

Same performance! The geometry coverage is what counts!



Method just as relevant for CUDA/OpenCL! But there you can make the passes with less overhead!



Conclusions:

- **Shader-based GPGPU is not dead, it is just not hyped**

Superior compatibility and ease of installation makes it highly interesting for the foreseeable future. Especially suitable for all image-related problems.

- **How to do GPGPU with shaders**

FBOs, Ping-ponging, algorithms, special considerations.

But stay tuned for Compute Shaders to change things...