

High-Performance Software Rasterization on GPUs

Samuli Laine

Tero Karras

NVIDIA Research



NVIDIA



HIGH-PERFORMANCE GRAPHICS
VANCOUVER, CANADA
AUGUST 5-7, 2011

Graphics and Programmability

- Graphics pipeline (OpenGL/D3D)
 - Driven by dedicated hardware
 - Executes user code in shaders
- What's next in programmability?



Programmable Graphics Pipeline

- But what does it mean?
 - Another API?
 - More programmable stages?
 - Coupling between CUDA/OpenCL and OpenGL/D3D?
 - Or "it's just a program"?



Our Approach

- Try and implement a full pixel pipeline using CUDA
 - From triangle setup to ROP
- Obey fundamental requirements of gfx pipe
 - **Maintain input order**
 - Hole-free rasterizer with correct rasterization rules
- Prefer speed over features



Project Goals

- Establish a firm data point through careful experiment
- Provide fully programmable platform for exploring algorithms that extend the hardware gfx pipe
 - Programmable ROP
 - Stochastic rasterization
 - Non-linear rasterization
 - Non-quad derivatives
 - Quad merging
 - Decoupled sampling
 - Compact after discard
 - etc.
- Ideas for future hardware
 - Ultimate goal = flexibility of software, performance of fixed-function hardware

Previous Work: FreePipe [Liu et al. 2010]

- Very simple rasterization pipeline
 - Each triangle processed by one thread
 - Blit pixels directly to DRAM using atomics
- Limitations
 - Cannot retain inputs order
 - Limited support for ROP operations (dictated by atomics)
 - Large triangles → game over
- We are 15x faster on average
 - Larger difference for high resolutions

Design Considerations

- Run everything in parallel
 - We need a lot of threads to fill the machine
- Minimize amount of synchronization
 - Avoid excessive use of atomics
- Focus on load balancing
 - Graphics workloads are wild



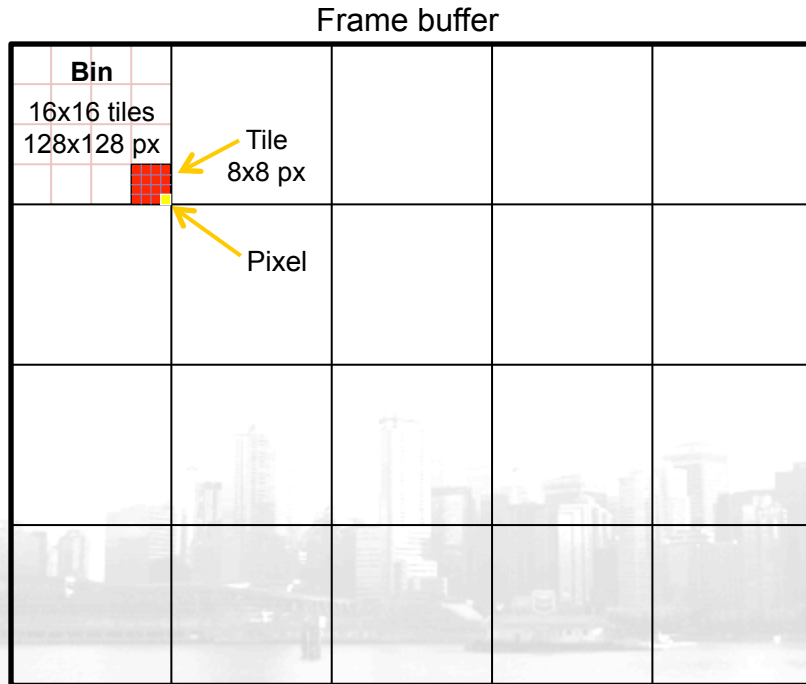
Pipeline Structure

- Chunker-style pipeline with four stages

Triangle setup → Bin raster → Coarse raster → Fine raster

- Run data in large batches
 - Separate kernel launch for each stage
- Keep data in input order all the time

Chunking to Bins and Tiles

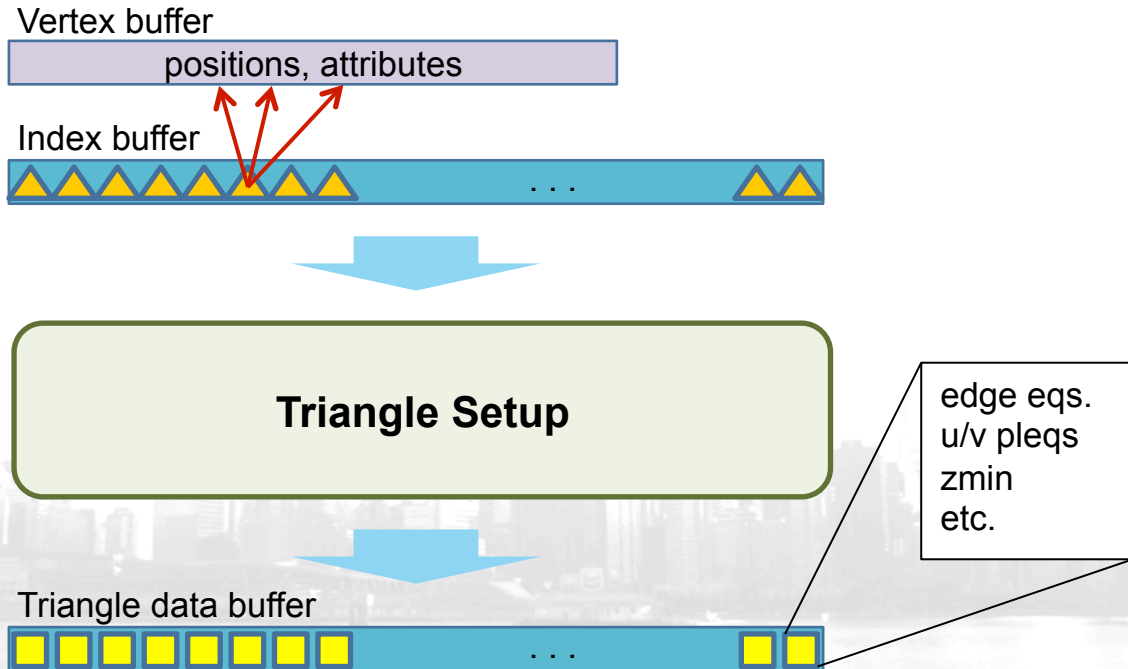


Pipeline Stages

- Quick look at each stage
- More details in the paper



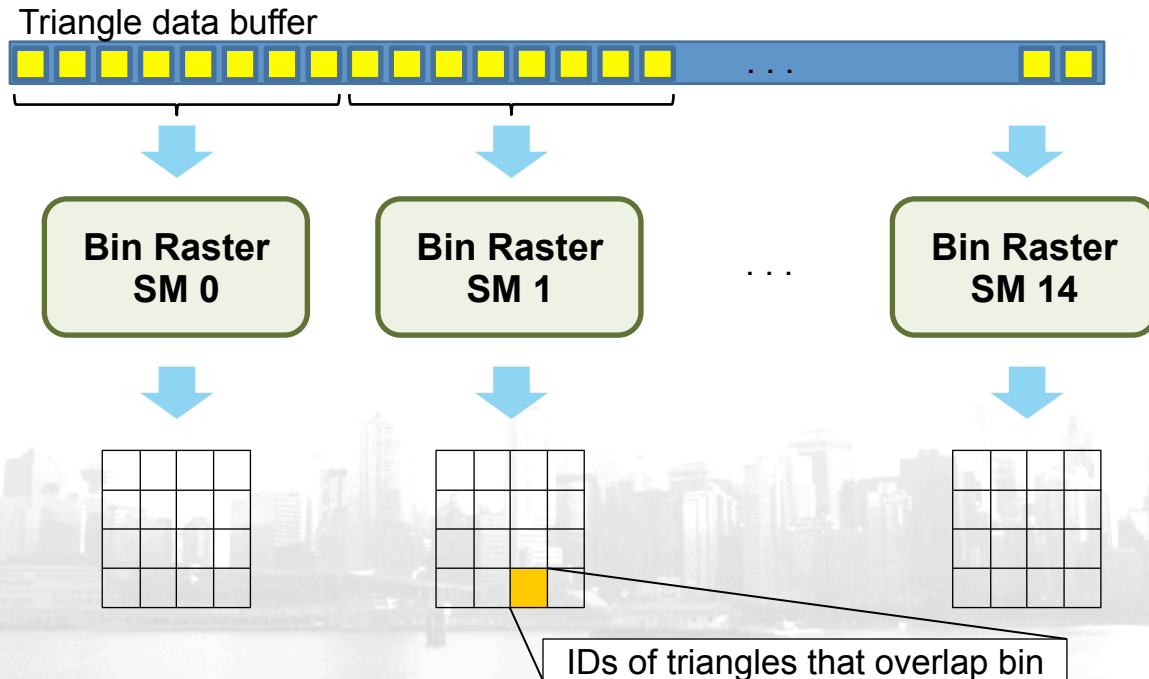
Triangle Setup



Triangle Setup

- Fetch vertex indices and positions
 - Clip if necessary (has guardband)
 - Frustum, backface and between-samples cull
 - Setup screen-space pleqs for u/w , v/w , z/w , $1/w$
 - Compute z_{min} for early depth cull in fine raster
- One-to-one mapping between input and output
 - Trivial to employ full chip while preserving ordering

Bin Raster



Bin Raster, First Phase

- Pick triangle batch (atomic, 16k tris)
- Read 512 set-up triangles
- Compact/expand according to culling/clipping results
 - Efficient within thread block
- Repeat until has enough triangles to utilize all threads

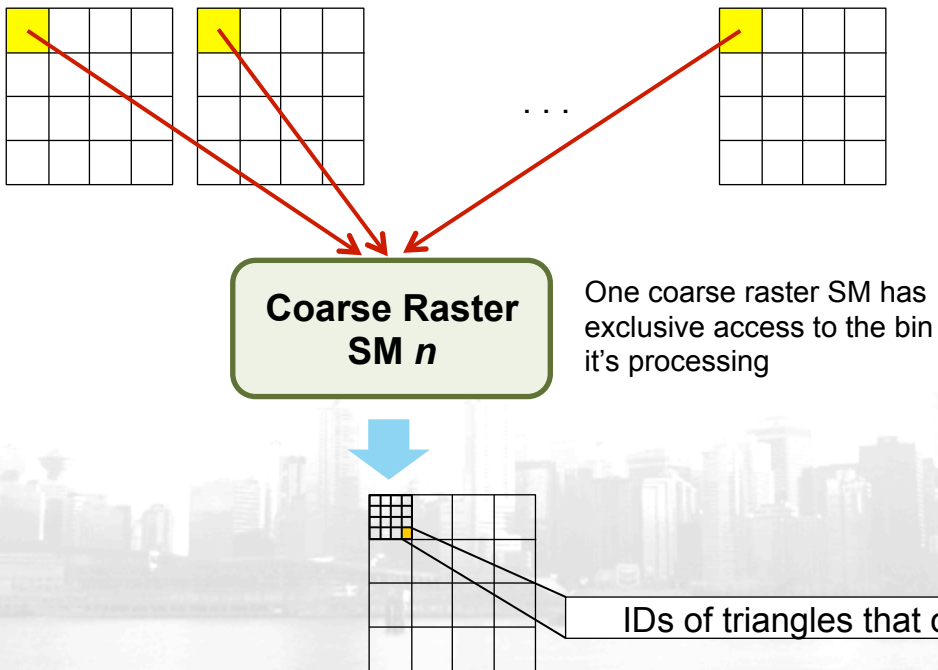


Bin Raster, Second Phase

- Rasterize
 - Determine bin coverage for each triangle (1 thread per tri)
 - Fast paths for 2x2 and smaller footprints
- Output
 - Output to per-SM queue → no sync between SMs



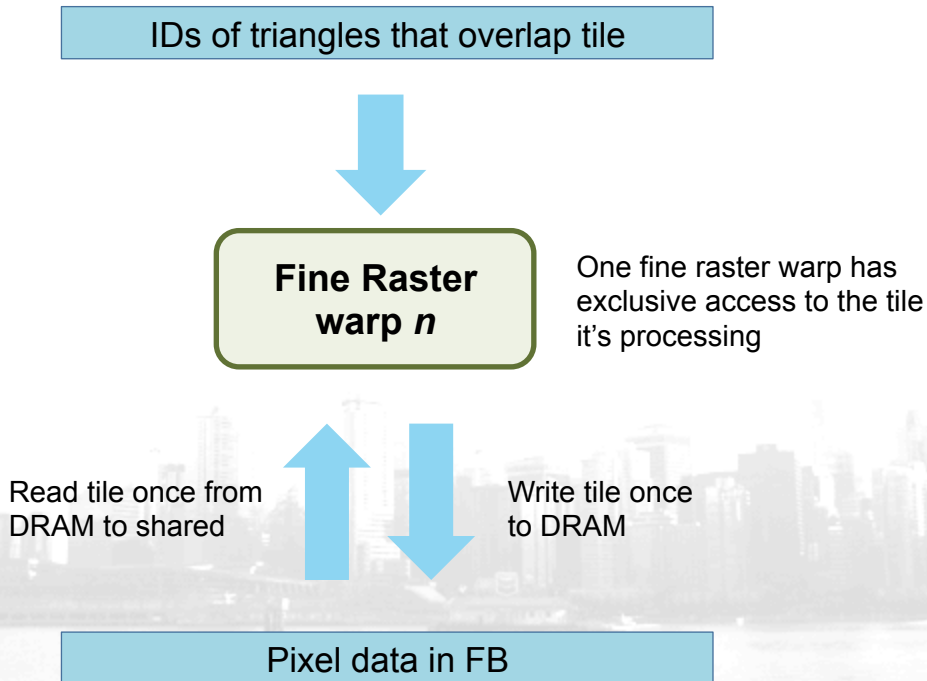
Coarse Raster



Coarse Raster

- Input Phase
 - Merge from 15 input queues (one per bin SM)
 - Continue until enough triangles to utilize all threads
- Rasterize
 - Determine tile coverage for each triangle (1 thread per tri)
 - Fast paths for small / largely overlapping triangles
- Output
 - Highly varying number of output items → divide evenly to threads
 - Only one SM outputs to tiles of any given bin → no sync needed

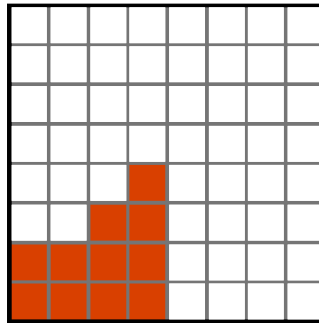
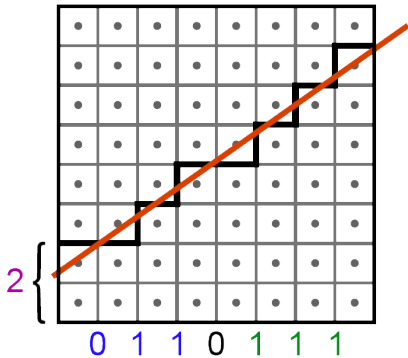
Fine Raster



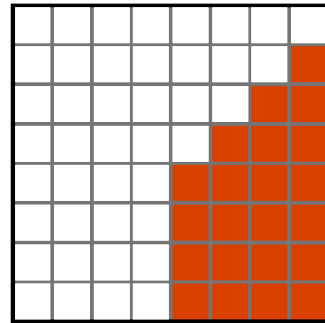
Fine Raster

- Pick tile, read FB tile, process, write FB tile
- Input
 - Read 32 triangles in shared memory
 - Early z kill based on triangle zmin and tile zmax
 - Calculate pixel coverage using LUTs (153 instr. for 8x8 stamp)
 - Repeat until has at least 32 fragments
- Raster
 - Process one fragment per thread, full utilization
 - Shade and ROP

Tidbit 1: Coverage Calculation



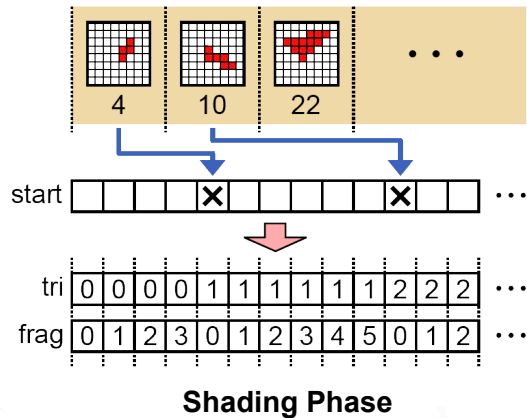
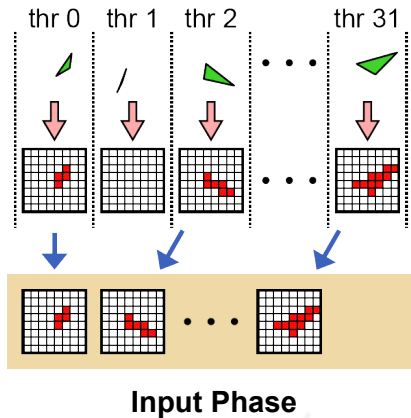
LUT A [2][011]



LUT B [2+2][111]

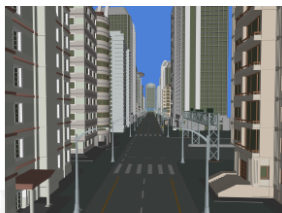
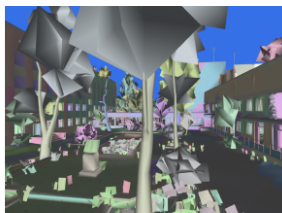
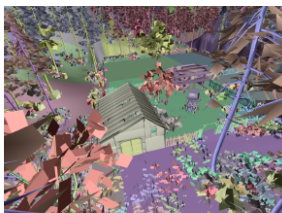
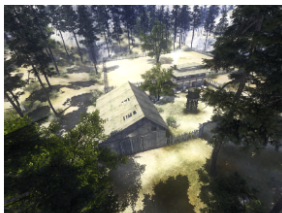
- Step along edge (Bresenham-like)
- Use look-up tables to generate coverage masks
- ~50 instructions for 8x8 stamp, one edge

Tidbit 2: Fragment Distribution



- In input phase, calculate coverage and store in list
- In shading phase, detect triangle changes and calculate triangle index and fragment in triangle

Test Scenes



SAN MIGUEL, 189MB
5.44M tris, 25% visible
2.4 pixels / triangle

JUAREZ, 24MB
546K tris, 37% visible
14.6 pixels / triangle

STALKER, 11MB
349K tris, 41% visible
14.1 pixels / triangle

CITY, 51MB
879K tris, 21% visible
16.3 pixels / triangle

BUDDHA, 29MB
1.09M tris, 32% visible
1.4 pixels / triangle

Call of Juarez scene courtesy of Techland

S.T.A.L.K.E.R.: Call of Pripyat scene courtesy of GSC Game World

Results: Performance

Scene	Resolution	HW	Our (SW)	FreePipe (FP)	SW:HW ratio	FP:SW ratio
SAN MIGUEL	512×384	5.37	7.82	130.14	1.46	16.65
	1024×768	5.43	9.48	510.20	1.74	53.84
	2048×1536	5.86	15.44	1652.52	2.64	107.06
JUAREZ	512×384	0.59	2.71	5.34	4.56	1.97
	1024×768	0.67	3.28	18.63	4.87	5.69
	2048×1536	1.03	7.06	72.45	6.84	10.26
STALKER	512×384	0.31	1.81	23.47	5.91	12.96
	1024×768	0.39	2.31	92.73	5.96	40.14
	2048×1536	0.67	5.41	386.07	8.10	71.36
CITY	512×384	0.93	2.16	64.56	2.32	29.88
	1024×768	1.04	3.13	251.86	3.01	80.54
	2048×1536	1.42	6.79	1032.83	4.77	152.13
BUDDHA	512×384	1.06	2.09	2.14	1.98	1.02
	1024×768	1.07	2.66	3.08	2.50	1.16
	2048×1536	1.11	4.01	6.96	3.62	1.73

Frame rendering time in ms (depth test + color, no MSAA, no blending)

Results: Memory Usage

	San Miguel	Juarez	Stalker	City	Buddha
Scene data	189	24	11	51	29
Triangle setup data	420.0	42.2	26.9	67.9	84.0
Bin queues	4.0	1.5	1.2	0.9	2.0
Tile queues	4.4	2.9	2.2	2.2	1.5

Memory usage in MB

Comparison to Hardware (1/3)

— Resolution

- Cannot match hardware in raster, z kill + compact
- Currently support max 2K x 2K frame buffer, 4 subpixel bits

— Attributes

- Fetched when used → bad latency hiding
- Expensive interpolation

— Antialiasing

- Hardware nearly oblivious to MSAA, we much less so

Comparison to Hardware (2/3)

- Memory usage, buffering through DRAM
 - Performance implications of reduced buffering unknown
 - Streaming through on-chip memory would be much better
- + Shader complexity
 - Shader performance theoretically the same as in graphics pipe
- + Frame buffer bandwidth
 - Each pixel touched only once in DRAM

Comparison to Hardware (3/3)

+ Extensibility

- Need one stage to do something extra?
- Need a new stage altogether?
- *You can actually implement it*

+ Specialization to individual applications

- Rip out what you don't need, hard-code what you can



Exploration Potential

- Shader performance boosters
 - Compact after discard, quad merging, decoupled sampling, ...
- Things to do with programmable ROP
 - A-buffering, order-independent transparency, ...
- Stochastic rasterization
- Non-linear rasterization
- (Your idea here)

The Code is Out There

- The entire codebase is open-sourced and released



The screenshot shows the SourceForge project page for **cudaraster**, an open-source implementation of "High-Performance Software Rasterization on GPUs". The page includes navigation links for Project Home, Downloads, Source, Checkout, Browse, and Changes. A large yellow box highlights the URL <http://code.google.com/p/cudaraster/>. Below the URL, a snippet of C++ code is displayed, showing various GPU-related functions and macros. On the right side of the code snippet, there are links for Project members, Older revisions, All revisions of this file, File info, and View raw file.

```
10 * Unless required by applicable law or agreed to in writing, software
11 * distributed under the License is distributed on an "AS IS" BASIS,
12 * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
13 * See the License for the specific language governing permissions and
14 * limitations under the License.
15 */
16
17 #include "util.hpp"
18 namespace FM
19 {
20 //-----
21 __device__ __inline__ U32 idiv_fast(U32 a, U32 b)
22 {
23     return f32_to_u32_sat_rmi(((F32)a * 0.5f) / (F32)b);
24 }
25 //-----
26 __device__ __inline__ U32 toABGR(Float4 color)
27 {
28     // 11 instructions: 4*FMA, 4*F2I, 3*PRMT
29     U32 x = f32_to_u32_sat_rmi(fma_rmi(color.x, (1 << 24) * 255.0f, (1 << 24) * 0.5f));
30     U32 y = f32_to_u32_sat_rmi(fma_rmi(color.y, (1 << 24) * 255.0f, (1 << 24) * 0.5f));
31     U32 z = f32_to_u32_sat_rmi(fma_rmi(color.z, (1 << 24) * 255.0f, (1 << 24) * 0.5f));
32     U32 w = f32_to_u32_sat_rmi(fma_rmi(color.w, (1 << 24) * 255.0f, (1 << 24) * 0.5f));
33     return prmt(prmt(x, y, 0x007f), prmt(z, w, 0x007f), 0x5430);
34 }
35 //-----
36 __device__ __inline__ U32 blendABGR(U32 src, U32 dst, U32 srcColorFactor, U32 dstColorFactor, U32 srcAlphaFactor, U32 dstAlphaFactor)
37 {
38     U32 a = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
39     U32 b = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
40     U32 c = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
41     U32 d = vmad_b3_b3(src, srcAlphaFactor, vmad_b3_b3(dst, dstAlphaFactor, 0)) * 0x0010101 + 0x800000;
42     return prmt(prmt(a, b, 0x007f), prmt(c, d, 0x007f), 0x5430);
43 }
44 //-----
45 __device__ __inline__ U32 blendABGR16(U32 src, U32 dst, U32 srcColorFactor, U32 dstColorFactor, U32 srcAlphaFactor, U32 dstAlphaFactor)
46 {
47     U32 a = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
48     U32 b = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
49     U32 c = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
50     U32 d = vmad_b3_b3(src, srcAlphaFactor, vmad_b3_b3(dst, dstAlphaFactor, 0)) * 0x0010101 + 0x800000;
51     return prmt(prmt(a, b, 0x007f), prmt(c, d, 0x007f), 0x5430);
52 }
53 //-----
54 __device__ __inline__ U32 blendABGR16(U32 src, U32 dst, U32 srcColorFactor, U32 dstColorFactor, U32 srcAlphaFactor, U32 dstAlphaFactor)
55 {
56     U32 a = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
57     U32 b = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
58     U32 c = vmad_b3_b3(src, srcColorFactor, vmad_b3_b3(dst, dstColorFactor, 0)) * 0x0010101 + 0x800000;
59     U32 d = vmad_b3_b3(src, srcAlphaFactor, vmad_b3_b3(dst, dstAlphaFactor, 0)) * 0x0010101 + 0x800000;
60     return prmt(prmt(a, b, 0x007f), prmt(c, d, 0x007f), 0x5430);
61 }
```

Thank You

- Questions

