



(51) International Patent Classification:

G06F 9/30 (2006.01) G06F 12/08 (2006.01)  
G06F 9/38 (2006.01) G06T 1/20 (2006.01)

(21) International Application Number:

PCT/US2016/036632

(22) International Filing Date:

9 June 2016 (09.06.2016)

(25) Filing Language:

English

(26) Publication Language:

English

(30) Priority Data:

14/810,684 28 July 2015 (28.07.2015) US

(71) Applicant: INTEL CORPORATION [US/US]; 2200 Mission College Boulevard, Santa Clara, California 95054 (US).

(72) Inventors: TSENG, Janet; 16511 NW Canton St., Portland, Oregon 97229 (US). DEGROOD, Felix J.; 625 SW 90th Ave., Portland, Oregon 97225 (US). MIN, Alexander W.; 5566 NW Primino Ave., Portland, Oregon 97229 (US). TSAI, Jr-Shian; 6657 NW 165th Ave., Portland,

Oregon 97229 (US). TAI, Tsung-Yuan C.; 12709 NW Majestic Sequoia Way, Portland, Oregon 97229 (US).

(74) Agents: TROP, Timothy N. et al.; Trop, Pruner & Hu, P.C., 1616 S. Voss Rd., Ste. 750, Houston, Texas 77057-2631 (US).

(81) Designated States (unless otherwise indicated, for every kind of national protection available): AE, AG, AL, AM, AO, AT, AU, AZ, BA, BB, BG, BH, BN, BR, BW, BY, BZ, CA, CH, CL, CN, CO, CR, CU, CZ, DE, DK, DM, DO, DZ, EC, EE, EG, ES, FI, GB, GD, GE, GH, GM, GT, HN, HR, HU, ID, IL, IN, IR, IS, JP, KE, KG, KN, KP, KR, KZ, LA, LC, LK, LR, LS, LU, LY, MA, MD, ME, MG, MK, MN, MW, MX, MY, MZ, NA, NG, NI, NO, NZ, OM, PA, PE, PG, PH, PL, PT, QA, RO, RS, RU, RW, SA, SC, SD, SE, SG, SK, SL, SM, ST, SV, SY, TH, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, ZA, ZM, ZW.

(84) Designated States (unless otherwise indicated, for every kind of regional protection available): ARIPO (BW, GH, GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU,

[Continued on next page]

(54) Title: PACKET PROCESSING ON GRAPHICS PROCESSING UNITS USING CONTINUOUS THREADS

(57) Abstract: In accordance with some embodiments, a continuous thread is operated on the graphics processing unit. A continuous thread is launched one time from the central processing unit and then it runs continuously until an application on the central processing unit decides to terminate the thread. For example, the application may decide to terminate the thread in one of a variety of situations which may be programmed in advance. For example, upon error detection, a desire to change the way that the thread on the graphics processing unit operates, or in power off, the thread may terminate. But unless actively terminated by the central processing unit, the continuous thread generally runs uninterrupted.

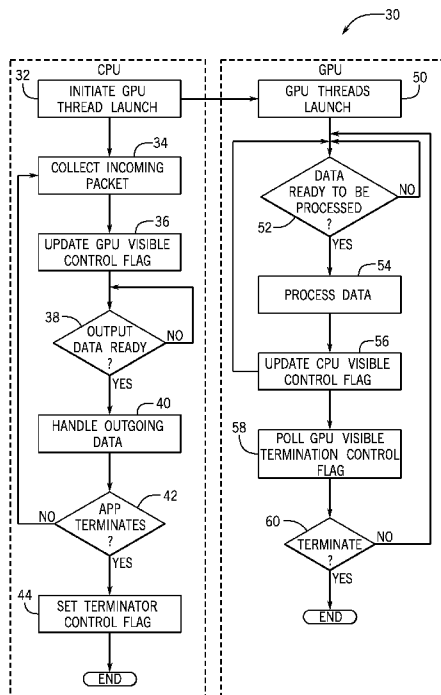


FIG. 2

WO 2017/019183 A1



LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK,  
SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ,  
GW, KM, ML, MR, NE, SN, TD, TG).

— *as to the applicant's entitlement to claim the priority of  
the earlier application (Rule 4.17(iii))*

**Declarations under Rule 4.17:**

- *as to the identity of the inventor (Rule 4.17(i))*
- *as to applicant's entitlement to apply for and be granted  
a patent (Rule 4.17(ii))*

**Published:**

- *with international search report (Art. 21(3))*

PACKET PROCESSING ON GRAPHICS  
PROCESSING UNITS USING CONTINUOUS THREADS

Background

[0001] This relates to graphics processing units and particularly to the use of graphics processing units to handle latency sensitive applications.

[0002] Because of advances in graphics processing unit architecture, graphics processing units are being relied upon to handle ever more complex operations. However, in connection with latency sensitive applications, graphics processing units have some drawbacks. Typically in order for a task to be handled by the graphics processing unit, it must be assigned from a central processing unit. This assignment operation involves the job passing through schedulers and command buffers that are part of a vertical stack that generally increases the time that many simple operations require. Because there may be a large number of these simple operations such as launching threads on graphics processing units, assignment of such tasks to graphics processing units in latency sensitive applications may not be effective.

[0003] Generally, despite the high programmability and parallel computation available with graphics processing units, accelerated packet processing on graphics processing is difficult because of the high central processing to graphics processing unit communication overhead and high costs of launching threads on a graphics processing unit.

[0004] The nature of network packet processing application is latency sensitivity. Typically graphics processing unit working threads rely on the host or central processing unit to notify the graphics processing unit when the producer data is ready to be processed. The latency introduced from thousands of thread launches per second plus the communication overhead between the central processing unit and the graphics processing unit is inconsistent with the latency requirements of network applications such as packet forwarding.

[0005] While some techniques such as batching do let kernels process a large number of tasks in order to amortize this overhead, with streaming-type applications

like packet forwarding, which is sensitive to latency, batch processing is not practical. One reason is that one can only batch a certain number of packets before being processed.

### Brief Description Of The Drawings

[0006] Some embodiments are described with respect to the following figures:

Figure 1 is a hardware depiction of one embodiment;

Figure 2 is a flow chart for one embodiment;

Figure 3 is a block diagram of a processing system according to one embodiment;

Figure 4 is a block diagram of a processor according to one embodiment;

Figure 5 is a block diagram of a graphics processor according to one embodiment;

Figure 6 is a block diagram of a graphics processing engine according to one embodiment;

Figure 7 is a block diagram of another embodiment of a graphics processor;

Figure 8 is a depiction thread execution logic according to one embodiment;

Figure 9 is a block diagram of a graphics processor instruction format according to some embodiments;

Figure 10 is a block diagram of another embodiment of a graphics processor;

Figure 11A is a block diagram of a graphics processor command format according to some embodiments;

Figure 11B is a block diagram illustrating a graphics processor command sequence according to some embodiments;

Figure 12 is a depiction of an exemplary graphics software architecture according to some embodiments;

Figure 13 is a block diagram illustrating an IP core development system according to some embodiments; and

Figure 14 is a block diagram showing an exemplary system on chip integrated circuit according to some embodiments.

### Detailed Description

[0007] In accordance with some embodiments, a continuous thread is operated on the graphics processing unit. A continuous thread is launched one time from the central processing unit and then it runs continuously until an application on the central processing unit decides to terminate the thread. For example, the application may decide to terminate the thread in one of a variety of situations which may be programmed in advance. For example, upon error detection, a desire to change the way that the thread on the graphics processing unit operates, or in power off, the thread may terminate. But unless actively terminated by the central processing unit, the continuous thread generally runs uninterrupted.

[0008] In some embodiments, the uninterrupted continuous thread may be implemented by fine grained termination procedures available on some operating systems such as Windows 8 available from Microsoft Corporation. In some cases, the thread may be comparable to a hardware thread, at least in terms of life span, in that the thread continues to live until such time the application decides to terminate it. Thus it can continuously operate to handle latency sensitive operations without a lot of communication between the central processing unit and the graphics processing unit, which communications add to latency. Since there is only one launch of a continuous thread latency may be reduced in some cases.

[0009] Referring to **Figure 1**, in some embodiments, a continuous thread is made possible by using a shared virtual memory 26 between the graphics processing unit 12 and the central processing unit 14 in a processor-based system 10. In Figure 1, the graphics processing unit 12 communicates with the central processing unit 14 through a level 3 or L3 cache 16 and a shared low level cache (LLC) 18. System memory 20 may include the shared virtual memory (SVM) 26. Thus a shared virtual memory may be accessed by both the general processing unit 12 and the central processing unit 14 are both read and write. Communications between the shared virtual memory 26 and the central processing unit take place by way of the shared LLC cache 18, the level 2 cache 22 and the level one cache 24 in some embodiments.

[0010] In some embodiments, the operating system supports fine grained graphics processing unit preemption to allow a long running kernel to coexist with other

rendering display kernels. In the hardware view of the control model, the graphics processing unit and the central processing unit communicate through a shared address space model instead of the traditional graphics processing unit command driven programming model. Thus in some embodiments, the graphics processing unit can handle continuous streaming workloads. It is not limited to handling one task at a time. In some embodiments, the model allows a graphics processing unit and the central processing unit to interact continuously without a lot of overhead.

[0011] In one model, the central processing unit is the producer and the graphics processing unit is the consumer. Two-way communication is possible.

[0012] When the graphics processing unit finishes processing a task assigned by the central processing unit, the graphics processing unit makes its output available to the central processing unit by updating a flag and making the information available to the central processing unit in the shared virtual memory in one embodiment. At the same time, when the central processing unit has a workload for the graphics processing unit to handle, it can update a graphics processing unit visible control flag that causes the graphics processing unit to obtain the workload from the shared virtual memory.

[0013] Thus in some embodiments, only one copy of the workload and the results is ever made. In addition, a control mechanism for communication between the graphics processing unit and the central processing unit is possible without going through a latency hungry driver stack. In some embodiments, as opposed to batch processing, the continuous thread continues to exist beyond any one workload. Persistent threads generally only last for one workload, commonly called the batch. One advantage of the continuous thread is that it reduces the overhead involved in repeatedly launching threads on the graphics processing unit and through the operation of the shared virtual memory, can reduce the amount of communication overhead between the graphics processing unit and the central processing unit.

[0014] The central processing unit application initiates determination of the continuous thread on the graphics processing unit through application scheduling of

workloads. Thus, the continuous thread continues to process work without being switched on.

[0015] In some embodiments, the sequence 30 shown in **Figure 2**, may be implemented in software, hardware and/or firmware. In software and firmware embodiments, it may be implemented in computer executed instructions stored in one or more non-transitory computer readable media such as magnetic, optical or semiconductor storage.

[0016] The central processing unit portion of the sequence 30 begins by initiating a graphics processing unit thread launch as indicated in block 32. Then incoming packets of a packet-based workload are collected as indicated in block 34. These packets are then made available to the graphics processing unit in a shared virtual memory, such as the SVM 26 in Figure 1. Then the graphics processing unit visible control flag is updated as indicated at block 36. This signals the graphics processing unit that a workload is available in the shared virtual memory for its use. Once the output data is ready as determined in diamond 38, the outgoing data is handled as indicated in block 40.

[0017] Then a check at diamond 42 determines whether there is any application-based termination in the central processing unit. A termination is a decision that a condition exists that would cause the continuous thread to be terminated. If not, the flow iterates and otherwise a termination control flag that is graphics processing unit visible is updated as indicated in block 44.

[0018] The graphics processing unit flow begins when it receives the initiation of the thread launch from the central processing unit. The graphics processing unit thread launches in response as indicated in block 50. Then in diamond 52, the graphics processing unit determines whether data is ready to be processed. This may be determined from the GPU visible control flag in one embodiment. If so, the data is processed as indicated in block 54 and the results are returned to the shared virtual memory. Then a central processing unit visible control flag is updated as indicated in block 56 to indicate that the workload is ready and to provide its location within the shared virtual memory.

[0019] Next the graphics processing unit visible determination control flag is pulled as indicated in block 58 to determine whether central processing unit based application desires for some reason to terminate the continuous thread. If there is no termination control flag, then the flow iterates as determined in diamond 60 and otherwise the flow ends.

[0020] The graphics processing unit threads spin-wait for the central processing unit flags to be updated. The central processing unit threads spin-waits for the graphics processing unit flags to be updated. Thus, coherency for a central processing unit and graphics processing unit is supplied by the shared virtual memory.

[0021] In some embodiments, processing throughput for graphics processing unit accelerated applications can be improved. In some embodiments, host to device communication may be reduced.

[0022] While the preceding discussion is in terms of graphics and central processing units, the concepts described herein can be applied to any group of two more different processing units including any processor, controller or accelerator.

[0023] **Figure 3** is a block diagram of a processing system 100, according to an embodiment. In various embodiments the system 100 includes one or more processors 102 and one or more graphics processors 108, and may be a single processor desktop system, a multiprocessor workstation system, or a server system having a large number of processors 102 or processor cores 107. In one embodiment, the system 100 is a processing platform incorporated within a system-on-a-chip (SoC) integrated circuit for use in mobile, handheld, or embedded devices.

[0024] An embodiment of system 100 can include, or be incorporated within a server-based gaming platform, a game console, including a game and media console, a mobile gaming console, a handheld game console, or an online game console. In some embodiments system 100 is a mobile phone, smart phone, tablet computing device or mobile Internet device. Data processing system 100 can also include, couple with, or be integrated within a wearable device, such as a smart watch wearable device, smart eyewear device, augmented reality device, or virtual



reality device. In some embodiments, data processing system 100 is a television or set top box device having one or more processors 102 and a graphical interface generated by one or more graphics processors 108.

[0025] In some embodiments, the one or more processors 102 each include one or more processor cores 107 to process instructions which, when executed, perform operations for system and user software. In some embodiments, each of the one or more processor cores 107 is configured to process a specific instruction set 109. In some embodiments, instruction set 109 may facilitate Complex Instruction Set Computing (CISC), Reduced Instruction Set Computing (RISC), or computing via a Very Long Instruction Word (VLIW). Multiple processor cores 107 may each process a different instruction set 109, which may include instructions to facilitate the emulation of other instruction sets. Processor core 107 may also include other processing devices, such a Digital Signal Processor (DSP).

[0026] In some embodiments, the processor 102 includes cache memory 104. Depending on the architecture, the processor 102 can have a single internal cache or multiple levels of internal cache. In some embodiments, the cache memory is shared among various components of the processor 102. In some embodiments, the processor 102 also uses an external cache (e.g., a Level-3 (L3) cache or Last Level Cache (LLC)) (not shown), which may be shared among processor cores 107 using known cache coherency techniques. A register file 106 is additionally included in processor 102 which may include different types of registers for storing different types of data (e.g., integer registers, floating point registers, status registers, and an instruction pointer register). Some registers may be general-purpose registers, while other registers may be specific to the design of the processor 102.

[0027] In some embodiments, processor 102 is coupled to a processor bus 110 to transmit communication signals such as address, data, or control signals between processor 102 and other components in system 100. In one embodiment the system 100 uses an exemplary 'hub' system architecture, including a memory controller hub 116 and an Input Output (I/O) controller hub 130. A memory controller hub 116 facilitates communication between a memory device and other components of system 100, while an I/O Controller Hub (ICH) 130 provides connections to I/O

devices via a local I/O bus. In one embodiment, the logic of the memory controller hub 116 is integrated within the processor.

[0028] Memory device 120 can be a dynamic random access memory (DRAM) device, a static random access memory (SRAM) device, flash memory device, phase-change memory device, or some other memory device having suitable performance to serve as process memory. In one embodiment the memory device 120 can operate as system memory for the system 100, to store data 122 and instructions 121 for use when the one or more processors 102 executes an application or process. Memory controller hub 116 also couples with an optional external graphics processor 112, which may communicate with the one or more graphics processors 108 in processors 102 to perform graphics and media operations.

[0029] In some embodiments, ICH 130 enables peripherals to connect to memory device 120 and processor 102 via a high-speed I/O bus. The I/O peripherals include, but are not limited to, an audio controller 146, a firmware interface 128, a wireless transceiver 126 (e.g., Wi-Fi, Bluetooth), a data storage device 124 (e.g., hard disk drive, flash memory, etc.), and a legacy I/O controller 140 for coupling legacy (e.g., Personal System 2 (PS/2)) devices to the system. One or more Universal Serial Bus (USB) controllers 142 connect input devices, such as keyboard and mouse 144 combinations. A network controller 134 may also couple to ICH 130. In some embodiments, a high-performance network controller (not shown) couples to processor bus 110. It will be appreciated that the system 100 shown is exemplary and not limiting, as other types of data processing systems that are differently configured may also be used. For example, the I/O controller hub 130 may be integrated within the one or more processor 102, or the memory controller hub 116 and I/O controller hub 130 may be integrated into a discreet external graphics processor, such as the external graphics processor 112.

[0030] **Figure 4** is a block diagram of an embodiment of a processor 200 having one or more processor cores 202A-202N, an integrated memory controller 214, and an integrated graphics processor 208. Those elements of **Figure 4** having the same reference numbers (or names) as the elements of any other figure herein can

operate or function in any manner similar to that described elsewhere herein, but are not limited to such. Processor 200 can include additional cores up to and including additional core 202N represented by the dashed lined boxes. Each of processor cores 202A-202N includes one or more internal cache units 204A-204N. In some embodiments each processor core also has access to one or more shared cache units 206.

[0031] The internal cache units 204A-204N and shared cache units 206 represent a cache memory hierarchy within the processor 200. The cache memory hierarchy may include at least one level of instruction and data cache within each processor core and one or more levels of shared mid-level cache, such as a Level 2 (L2), Level 3 (L3), Level 4 (L4), or other levels of cache, where the highest level of cache before external memory is classified as the LLC. In some embodiments, cache coherency logic maintains coherency between the various cache units 206 and 204A-204N.

[0032] In some embodiments, processor 200 may also include a set of one or more bus controller units 216 and a system agent core 210. The one or more bus controller units 216 manage a set of peripheral buses, such as one or more Peripheral Component Interconnect buses (e.g., PCI, PCI Express). System agent core 210 provides management functionality for the various processor components. In some embodiments, system agent core 210 includes one or more integrated memory controllers 214 to manage access to various external memory devices (not shown).

[0033] In some embodiments, one or more of the processor cores 202A-202N include support for simultaneous multi-threading. In such embodiment, the system agent core 210 includes components for coordinating and operating cores 202A-202N during multi-threaded processing. System agent core 210 may additionally include a power control unit (PCU), which includes logic and components to regulate the power state of processor cores 202A-202N and graphics processor 208.

[0034] In some embodiments, processor 200 additionally includes graphics processor 208 to execute graphics processing operations. In some embodiments, the graphics processor 208 couples with the set of shared cache units 206, and the

system agent core 210, including the one or more integrated memory controllers 214. In some embodiments, a display controller 211 is coupled with the graphics processor 208 to drive graphics processor output to one or more coupled displays. In some embodiments, display controller 211 may be a separate module coupled with the graphics processor via at least one interconnect, or may be integrated within the graphics processor 208 or system agent core 210.

[0035] In some embodiments, a ring based interconnect unit 212 is used to couple the internal components of the processor 200. However, an alternative interconnect unit may be used, such as a point-to-point interconnect, a switched interconnect, or other techniques, including techniques well known in the art. In some embodiments, graphics processor 208 couples with the ring interconnect 212 via an I/O link 213.

[0036] The exemplary I/O link 213 represents at least one of multiple varieties of I/O interconnects, including an on package I/O interconnect which facilitates communication between various processor components and a high-performance embedded memory module 218, such as an eDRAM module. In some embodiments, each of the processor cores 202-202N and graphics processor 208 use embedded memory modules 218 as a shared Last Level Cache.

[0037] In some embodiments, processor cores 202A-202N are homogenous cores executing the same instruction set architecture. In another embodiment, processor cores 202A-202N are heterogeneous in terms of instruction set architecture (ISA), where one or more of processor cores 202A-N execute a first instruction set, while at least one of the other cores executes a subset of the first instruction set or a different instruction set. In one embodiment processor cores 202A-202N are heterogeneous in terms of microarchitecture, where one or more cores having a relatively higher power consumption couple with one or more power cores having a lower power consumption. Additionally, processor 200 can be implemented on one or more chips or as an SoC integrated circuit having the illustrated components, in addition to other components.

[0038] **Figure 5** is a block diagram of a graphics processor 300, which may be a discrete graphics processing unit, or may be a graphics processor integrated with a

plurality of processing cores. In some embodiments, the graphics processor communicates via a memory mapped I/O interface to registers on the graphics processor and with commands placed into the processor memory. In some embodiments, graphics processor 300 includes a memory interface 314 to access memory. Memory interface 314 can be an interface to local memory, one or more internal caches, one or more shared external caches, and/or to system memory.

[0039] In some embodiments, graphics processor 300 also includes a display controller 302 to drive display output data to a display device 320. Display controller 302 includes hardware for one or more overlay planes for the display and composition of multiple layers of video or user interface elements. In some embodiments, graphics processor 300 includes a video codec engine 306 to encode, decode, or transcode media to, from, or between one or more media encoding formats, including, but not limited to Moving Picture Experts Group (MPEG) formats such as MPEG-2, Advanced Video Coding (AVC) formats such as H.264/MPEG-4 AVC, as well as the Society of Motion Picture & Television Engineers (SMPTE) 421M/VC-1, and Joint Photographic Experts Group (JPEG) formats such as JPEG, and Motion JPEG (MJPEG) formats.

[0040] In some embodiments, graphics processor 300 includes a block image transfer (BLIT) engine 304 to perform two-dimensional (2D) rasterizer operations including, for example, bit-boundary block transfers. However, in one embodiment, 2D graphics operations are performed using one or more components of graphics processing engine (GPE) 310. In some embodiments, graphics processing engine 310 is a compute engine for performing graphics operations, including three-dimensional (3D) graphics operations and media operations.

[0041] In some embodiments, GPE 310 includes a 3D pipeline 312 for performing 3D operations, such as rendering three-dimensional images and scenes using processing functions that act upon 3D primitive shapes (e.g., rectangle, triangle, etc.). The 3D pipeline 312 includes programmable and fixed function elements that perform various tasks within the element and/or spawn execution threads to a 3D/Media sub-system 315. While 3D pipeline 312 can be used to perform media operations, an embodiment of GPE 310 also includes a media pipeline 316 that is

specifically used to perform media operations, such as video post-processing and image enhancement.

[0042] In some embodiments, media pipeline 316 includes fixed function or programmable logic units to perform one or more specialized media operations, such as video decode acceleration, video de-interlacing, and video encode acceleration in place of, or on behalf of video codec engine 306. In some embodiments, media pipeline 316 additionally includes a thread spawning unit to spawn threads for execution on 3D/Media sub-system 315. The spawned threads perform computations for the media operations on one or more graphics execution units included in 3D/Media sub-system 315.

[0043] In some embodiments, 3D/Media subsystem 315 includes logic for executing threads spawned by 3D pipeline 312 and media pipeline 316. In one embodiment, the pipelines send thread execution requests to 3D/Media subsystem 315, which includes thread dispatch logic for arbitrating and dispatching the various requests to available thread execution resources. The execution resources include an array of graphics execution units to process the 3D and media threads. In some embodiments, 3D/Media subsystem 315 includes one or more internal caches for thread instructions and data. In some embodiments, the subsystem also includes shared memory, including registers and addressable memory, to share data between threads and to store output data.

[0044] **Figure 6** is a block diagram of a graphics processing engine 410 of a graphics processor in accordance with some embodiments. In one embodiment, the GPE 410 is a version of the GPE 310 shown in **Figure 5**. Elements of **Figure 6** having the same reference numbers (or names) as the elements of any other figure herein can operate or function in any manner similar to that described elsewhere herein, but are not limited to such.

[0045] In some embodiments, GPE 410 couples with a command streamer 403, which provides a command stream to the GPE 3D and media pipelines 412, 416. In some embodiments, command streamer 403 is coupled to memory, which can be system memory, or one or more of internal cache memory and shared cache

memory. In some embodiments, command streamer 403 receives commands from the memory and sends the commands to 3D pipeline 412 and/or media pipeline 416. The commands are directives fetched from a ring buffer, which stores commands for the 3D and media pipelines 412, 416. In one embodiment, the ring buffer can additionally include batch command buffers storing batches of multiple commands. The 3D and media pipelines 412, 416 process the commands by performing operations via logic within the respective pipelines or by dispatching one or more execution threads to an execution unit array 414. In some embodiments, execution unit array 414 is scalable, such that the array includes a variable number of execution units based on the target power and performance level of GPE 410.

[0046] In some embodiments, a sampling engine 430 couples with memory (e.g., cache memory or system memory) and execution unit array 414. In some embodiments, sampling engine 430 provides a memory access mechanism for execution unit array 414 that allows execution array 414 to read graphics and media data from memory. In some embodiments, sampling engine 430 includes logic to perform specialized image sampling operations for media.

[0047] In some embodiments, the specialized media sampling logic in sampling engine 430 includes a de-noise/de-interlace module 432, a motion estimation module 434, and an image scaling and filtering module 436. In some embodiments, de-noise/de-interlace module 432 includes logic to perform one or more of a de-noise or a de-interlace algorithm on decoded video data. The de-interlace logic combines alternating fields of interlaced video content into a single frame of video. The de-noise logic reduces or removes data noise from video and image data. In some embodiments, the de-noise logic and de-interlace logic are motion adaptive and use spatial or temporal filtering based on the amount of motion detected in the video data. In some embodiments, the de-noise/de-interlace module 432 includes dedicated motion detection logic (e.g., within the motion estimation engine 434).

[0048] In some embodiments, motion estimation engine 434 provides hardware acceleration for video operations by performing video acceleration functions such as motion vector estimation and prediction on video data. The motion estimation engine determines motion vectors that describe the transformation of image data

between successive video frames. In some embodiments, a graphics processor media codec uses video motion estimation engine 434 to perform operations on video at the macro-block level that may otherwise be too computationally intensive to perform with a general-purpose processor. In some embodiments, motion estimation engine 434 is generally available to graphics processor components to assist with video decode and processing functions that are sensitive or adaptive to the direction or magnitude of the motion within video data.

[0049] In some embodiments, image scaling and filtering module 436 performs image-processing operations to enhance the visual quality of generated images and video. In some embodiments, scaling and filtering module 436 processes image and video data during the sampling operation before providing the data to execution unit array 414.

[0050] In some embodiments, the GPE 410 includes a data port 444, which provides an additional mechanism for graphics subsystems to access memory. In some embodiments, data port 444 facilitates memory access for operations including render target writes, constant buffer reads, scratch memory space reads/writes, and media surface accesses. In some embodiments, data port 444 includes cache memory space to cache accesses to memory. The cache memory can be a single data cache or separated into multiple caches for the multiple subsystems that access memory via the data port (e.g., a render buffer cache, a constant buffer cache, etc.). In some embodiments, threads executing on an execution unit in execution unit array 414 communicate with the data port by exchanging messages via a data distribution interconnect that couples each of the sub-systems of GPE 410.

[0051] **Figure 7** is a block diagram of another embodiment of a graphics processor 500. Elements of **Figure 7** having the same reference numbers (or names) as the elements of any other figure herein can operate or function in any manner similar to that described elsewhere herein, but are not limited to such.

[0052] In some embodiments, graphics processor 500 includes a ring interconnect 502, a pipeline front-end 504, a media engine 537, and graphics cores 580A-580N. In some embodiments, ring interconnect 502 couples the graphics processor to other



processing units, including other graphics processors or one or more general-purpose processor cores. In some embodiments, the graphics processor is one of many processors integrated within a multi-core processing system.

[0053] In some embodiments, graphics processor 500 receives batches of commands via ring interconnect 502. The incoming commands are interpreted by a command streamer 503 in the pipeline front-end 504. In some embodiments, graphics processor 500 includes scalable execution logic to perform 3D geometry processing and media processing via the graphics core(s) 580A-580N. For 3D geometry processing commands, command streamer 503 supplies commands to geometry pipeline 536. For at least some media processing commands, command streamer 503 supplies the commands to a video front end 534, which couples with a media engine 537. In some embodiments, media engine 537 includes a Video Quality Engine (VQE) 530 for video and image post-processing and a multi-format encode/decode (MFX) 533 engine to provide hardware-accelerated media data encode and decode. In some embodiments, geometry pipeline 536 and media engine 537 each generate execution threads for the thread execution resources provided by at least one graphics core 580A.

[0054] In some embodiments, graphics processor 500 includes scalable thread execution resources featuring modular cores 580A-580N (sometimes referred to as core slices), each having multiple sub-cores 550A-550N, 560A-560N (sometimes referred to as core sub-slices). In some embodiments, graphics processor 500 can have any number of graphics cores 580A through 580N. In some embodiments, graphics processor 500 includes a graphics core 580A having at least a first sub-core 550A and a second core sub-core 560A. In other embodiments, the graphics processor is a low power processor with a single sub-core (e.g., 550A). In some embodiments, graphics processor 500 includes multiple graphics cores 580A-580N, each including a set of first sub-cores 550A-550N and a set of second sub-cores 560A-560N. Each sub-core in the set of first sub-cores 550A-550N includes at least a first set of execution units 552A-552N and media/texture samplers 554A-554N. Each sub-core in the set of second sub-cores 560A-560N includes at least a second set of execution units 562A-562N and samplers 564A-564N. In some embodiments,

each sub-core 550A-550N, 560A-560N shares a set of shared resources 570A-570N. In some embodiments, the shared resources include shared cache memory and pixel operation logic. Other shared resources may also be included in the various embodiments of the graphics processor.

[0055] **Figure 8** illustrates thread execution logic 600 including an array of processing elements employed in some embodiments of a GPE. Elements of **Figure 8** having the same reference numbers (or names) as the elements of any other figure herein can operate or function in any manner similar to that described elsewhere herein, but are not limited to such.

[0056] In some embodiments, thread execution logic 600 includes a pixel shader 602, a thread dispatcher 604, instruction cache 606, a scalable execution unit array including a plurality of execution units 608A-608N, a sampler 610, a data cache 612, and a data port 614. In one embodiment the included components are interconnected via an interconnect fabric that links to each of the components. In some embodiments, thread execution logic 600 includes one or more connections to memory, such as system memory or cache memory, through one or more of instruction cache 606, data port 614, sampler 610, and execution unit array 608A-608N. In some embodiments, each execution unit (e.g. 608A) is an individual vector processor capable of executing multiple simultaneous threads and processing multiple data elements in parallel for each thread. In some embodiments, execution unit array 608A-608N includes any number individual execution units.

[0057] In some embodiments, execution unit array 608A-608N is primarily used to execute "shader" programs. In some embodiments, the execution units in array 608A-608N execute an instruction set that includes native support for many standard 3D graphics shader instructions, such that shader programs from graphics libraries (e.g., Direct 3D and OpenGL) are executed with a minimal translation. The execution units support vertex and geometry processing (e.g., vertex programs, geometry programs, vertex shaders), pixel processing (e.g., pixel shaders, fragment shaders) and general-purpose processing (e.g., compute and media shaders).

[0058] Each execution unit in execution unit array 608A-608N operates on arrays of data elements. The number of data elements is the “execution size,” or the number of channels for the instruction. An execution channel is a logical unit of execution for data element access, masking, and flow control within instructions. The number of channels may be independent of the number of physical Arithmetic Logic Units (ALUs) or Floating Point Units (FPUs) for a particular graphics processor. In some embodiments, execution units 608A-608N support integer and floating-point data types.

[0059] The execution unit instruction set includes single instruction multiple data (SIMD) instructions. The various data elements can be stored as a packed data type in a register and the execution unit will process the various elements based on the data size of the elements. For example, when operating on a 256-bit wide vector, the 256 bits of the vector are stored in a register and the execution unit operates on the vector as four separate 64-bit packed data elements (Quad-Word (QW) size data elements), eight separate 32-bit packed data elements (Double Word (DW) size data elements), sixteen separate 16-bit packed data elements (Word (W) size data elements), or thirty-two separate 8-bit data elements (byte (B) size data elements). However, different vector widths and register sizes are possible.

[0060] One or more internal instruction caches (e.g., 606) are included in the thread execution logic 600 to cache thread instructions for the execution units. In some embodiments, one or more data caches (e.g., 612) are included to cache thread data during thread execution. In some embodiments, sampler 610 is included to provide texture sampling for 3D operations and media sampling for media operations. In some embodiments, sampler 610 includes specialized texture or media sampling functionality to process texture or media data during the sampling process before providing the sampled data to an execution unit.

[0061] During execution, the graphics and media pipelines send thread initiation requests to thread execution logic 600 via thread spawning and dispatch logic. In some embodiments, thread execution logic 600 includes a local thread dispatcher 604 that arbitrates thread initiation requests from the graphics and media pipelines and instantiates the requested threads on one or more execution units 608A-608N.

For example, the geometry pipeline (e.g., 536 of **Figure 7**) dispatches vertex processing, tessellation, or geometry processing threads to thread execution logic 600 (**Figure 8**). In some embodiments, thread dispatcher 604 can also process runtime thread spawning requests from the executing shader programs.

[0062] Once a group of geometric objects has been processed and rasterized into pixel data, pixel shader 602 is invoked to further compute output information and cause results to be written to output surfaces (e.g., color buffers, depth buffers, stencil buffers, etc.). In some embodiments, pixel shader 602 calculates the values of the various vertex attributes that are to be interpolated across the rasterized object. In some embodiments, pixel shader 602 then executes an application programming interface (API)-supplied pixel shader program. To execute the pixel shader program, pixel shader 602 dispatches threads to an execution unit (e.g., 608A) via thread dispatcher 604. In some embodiments, pixel shader 602 uses texture sampling logic in sampler 610 to access texture data in texture maps stored in memory. Arithmetic operations on the texture data and the input geometry data compute pixel color data for each geometric fragment, or discards one or more pixels from further processing.

[0063] In some embodiments, the data port 614 provides a memory access mechanism for the thread execution logic 600 output processed data to memory for processing on a graphics processor output pipeline. In some embodiments, the data port 614 includes or couples to one or more cache memories (e.g., data cache 612) to cache data for memory access via the data port.

[0064] **Figure 9** is a block diagram illustrating a graphics processor instruction formats 700 according to some embodiments. In one or more embodiment, the graphics processor execution units support an instruction set having instructions in multiple formats. The solid lined boxes illustrate the components that are generally included in an execution unit instruction, while the dashed lines include components that are optional or that are only included in a sub-set of the instructions. In some embodiments, instruction format 700 described and illustrated are macro-instructions, in that they are instructions supplied to the execution unit, as opposed

to micro-operations resulting from instruction decode once the instruction is processed.

[0065] In some embodiments, the graphics processor execution units natively support instructions in a 128-bit format 710. A 64-bit compacted instruction format 730 is available for some instructions based on the selected instruction, instruction options, and number of operands. The native 128-bit format 710 provides access to all instruction options, while some options and operations are restricted in the 64-bit format 730. The native instructions available in the 64-bit format 730 vary by embodiment. In some embodiments, the instruction is compacted in part using a set of index values in an index field 713. The execution unit hardware references a set of compaction tables based on the index values and uses the compaction table outputs to reconstruct a native instruction in the 128-bit format 710.

[0066] For each format, instruction opcode 712 defines the operation that the execution unit is to perform. The execution units execute each instruction in parallel across the multiple data elements of each operand. For example, in response to an add instruction the execution unit performs a simultaneous add operation across each color channel representing a texture element or picture element. By default, the execution unit performs each instruction across all data channels of the operands. In some embodiments, instruction control field 714 enables control over certain execution options, such as channels selection (e.g., predication) and data channel order (e.g., swizzle). For 128-bit instructions 710 an exec-size field 716 limits the number of data channels that will be executed in parallel. In some embodiments, exec-size field 716 is not available for use in the 64-bit compact instruction format 730.

[0067] Some execution unit instructions have up to three operands including two source operands, src0 722, src1 722, and one destination 718. In some embodiments, the execution units support dual destination instructions, where one of the destinations is implied. Data manipulation instructions can have a third source operand (e.g., SRC2 724), where the instruction opcode 712 determines the number of source operands. An instruction's last source operand can be an immediate (e.g., hard-coded) value passed with the instruction.

[0068] In some embodiments, the 128-bit instruction format 710 includes an access/address mode information 726 specifying, for example, whether direct register addressing mode or indirect register addressing mode is used. When direct register addressing mode is used, the register address of one or more operands is directly provided by bits in the instruction 710.

[0069] In some embodiments, the 128-bit instruction format 710 includes an access/address mode field 726, which specifies an address mode and/or an access mode for the instruction. In one embodiment the access mode to define a data access alignment for the instruction. Some embodiments support access modes including a 16-byte aligned access mode and a 1-byte aligned access mode, where the byte alignment of the access mode determines the access alignment of the instruction operands. For example, when in a first mode, the instruction 710 may use byte-aligned addressing for source and destination operands and when in a second mode, the instruction 710 may use 16-byte-aligned addressing for all source and destination operands.

[0070] In one embodiment, the address mode portion of the access/address mode field 726 determines whether the instruction is to use direct or indirect addressing. When direct register addressing mode is used bits in the instruction 710 directly provide the register address of one or more operands. When indirect register addressing mode is used, the register address of one or more operands may be computed based on an address register value and an address immediate field in the instruction.

[0071] In some embodiments instructions are grouped based on opcode 712 bit-fields to simplify Opcode decode 740. For an 8-bit opcode, bits 4, 5, and 6 allow the execution unit to determine the type of opcode. The precise opcode grouping shown is merely an example. In some embodiments, a move and logic opcode group 742 includes data movement and logic instructions (e.g., move (mov), compare (cmp)). In some embodiments, move and logic group 742 shares the five most significant bits (MSB), where move (mov) instructions are in the form of 0000xxxxb and logic instructions are in the form of 0001xxxxb. A flow control instruction group 744 (e.g., call, jump (jmp)) includes instructions in the form of 0010xxxxb (e.g., 0x20). A

miscellaneous instruction group 746 includes a mix of instructions, including synchronization instructions (e.g., wait, send) in the form of 0011xxxxb (e.g., 0x30). A parallel math instruction group 748 includes component-wise arithmetic instructions (e.g., add, multiply (mul)) in the form of 0100xxxxb (e.g., 0x40). The parallel math group 748 performs the arithmetic operations in parallel across data channels. The vector math group 750 includes arithmetic instructions (e.g., dp4) in the form of 0101xxxxb (e.g., 0x50). The vector math group performs arithmetic such as dot product calculations on vector operands.

[0072] **Figure 10** is a block diagram of another embodiment of a graphics processor 800. Elements of **Figure 10** having the same reference numbers (or names) as the elements of any other figure herein can operate or function in any manner similar to that described elsewhere herein, but are not limited to such.

[0073] In some embodiments, graphics processor 800 includes a graphics pipeline 820, a media pipeline 830, a display engine 840, thread execution logic 850, and a render output pipeline 870. In some embodiments, graphics processor 800 is a graphics processor within a multi-core processing system that includes one or more general purpose processing cores. The graphics processor is controlled by register writes to one or more control registers (not shown) or via commands issued to graphics processor 800 via a ring interconnect 802. In some embodiments, ring interconnect 802 couples graphics processor 800 to other processing components, such as other graphics processors or general-purpose processors. Commands from ring interconnect 802 are interpreted by a command streamer 803, which supplies instructions to individual components of graphics pipeline 820 or media pipeline 830.

[0074] In some embodiments, command streamer 803 directs the operation of a vertex fetcher 805 that reads vertex data from memory and executes vertex-processing commands provided by command streamer 803. In some embodiments, vertex fetcher 805 provides vertex data to a vertex shader 807, which performs coordinate space transformation and lighting operations to each vertex. In some embodiments, vertex fetcher 805 and vertex shader 807 execute vertex-processing instructions by dispatching execution threads to execution units 852A, 852B via a thread dispatcher 831.

[0075] In some embodiments, execution units 852A, 852B are an array of vector processors having an instruction set for performing graphics and media operations. In some embodiments, execution units 852A, 852B have an attached L1 cache 851 that is specific for each array or shared between the arrays. The cache can be configured as a data cache, an instruction cache, or a single cache that is partitioned to contain data and instructions in different partitions.

[0076] In some embodiments, graphics pipeline 820 includes tessellation components to perform hardware-accelerated tessellation of 3D objects. In some embodiments, a programmable hull shader 811 configures the tessellation operations. A programmable domain shader 817 provides back-end evaluation of tessellation output. A tessellator 813 operates at the direction of hull shader 811 and contains special purpose logic to generate a set of detailed geometric objects based on a coarse geometric model that is provided as input to graphics pipeline 820. In some embodiments, if tessellation is not used, tessellation components 811, 813, 817 can be bypassed.

[0077] In some embodiments, complete geometric objects can be processed by a geometry shader 819 via one or more threads dispatched to execution units 852A, 852B, or can proceed directly to the clipper 829. In some embodiments, the geometry shader operates on entire geometric objects, rather than vertices or patches of vertices as in previous stages of the graphics pipeline. If the tessellation is disabled the geometry shader 819 receives input from the vertex shader 807. In some embodiments, geometry shader 819 is programmable by a geometry shader program to perform geometry tessellation if the tessellation units are disabled.

[0078] Before rasterization, a clipper 829 processes vertex data. The clipper 829 may be a fixed function clipper or a programmable clipper having clipping and geometry shader functions. In some embodiments, a rasterizer/depth 873 in the render output pipeline 870 dispatches pixel shaders to convert the geometric objects into their per pixel representations. In some embodiments, pixel shader logic is included in thread execution logic 850. In some embodiments, an application can bypass the rasterizer 873 and access un-rasterized vertex data via a stream out unit 823.



[0079] The graphics processor 800 has an interconnect bus, interconnect fabric, or some other interconnect mechanism that allows data and message passing amongst the major components of the processor. In some embodiments, execution units 852A, 852B and associated cache(s) 851, texture and media sampler 854, and texture/sampler cache 858 interconnect via a data port 856 to perform memory access and communicate with render output pipeline components of the processor. In some embodiments, sampler 854, caches 851, 858 and execution units 852A, 852B each have separate memory access paths.

[0080] In some embodiments, render output pipeline 870 contains a rasterizer and depth test component 873 that converts vertex-based objects into an associated pixel-based representation. In some embodiments, the rasterizer logic includes a windower/masker unit to perform fixed function triangle and line rasterization. An associated render cache 878 and depth cache 879 are also available in some embodiments. A pixel operations component 877 performs pixel-based operations on the data, though in some instances, pixel operations associated with 2D operations (e.g. bit block image transfers with blending) are performed by the 2D engine 841, or substituted at display time by the display controller 843 using overlay display planes. In some embodiments, a shared L3 cache 875 is available to all graphics components, allowing the sharing of data without the use of main system memory.

[0081] In some embodiments, graphics processor media pipeline 830 includes a media engine 837 and a video front end 834. In some embodiments, video front end 834 receives pipeline commands from the command streamer 803. In some embodiments, media pipeline 830 includes a separate command streamer. In some embodiments, video front-end 834 processes media commands before sending the command to the media engine 837. In some embodiments, media engine 337 includes thread spawning functionality to spawn threads for dispatch to thread execution logic 850 via thread dispatcher 831.

[0082] In some embodiments, graphics processor 800 includes a display engine 840. In some embodiments, display engine 840 is external to processor 800 and couples with the graphics processor via the ring interconnect 802, or some other

interconnect bus or fabric. In some embodiments, display engine 840 includes a 2D engine 841 and a display controller 843. In some embodiments, display engine 840 contains special purpose logic capable of operating independently of the 3D pipeline. In some embodiments, display controller 843 couples with a display device (not shown), which may be a system integrated display device, as in a laptop computer, or an external display device attached via a display device connector.

[0083] In some embodiments, graphics pipeline 820 and media pipeline 830 are configurable to perform operations based on multiple graphics and media programming interfaces and are not specific to any one application programming interface (API). In some embodiments, driver software for the graphics processor translates API calls that are specific to a particular graphics or media library into commands that can be processed by the graphics processor. In some embodiments, support is provided for the Open Graphics Library (OpenGL) and Open Computing Language (OpenCL) from the Khronos Group, the Direct3D library from the Microsoft Corporation, or support may be provided to both OpenGL and D3D. Support may also be provided for the Open Source Computer Vision Library (OpenCV). A future API with a compatible 3D pipeline would also be supported if a mapping can be made from the pipeline of the future API to the pipeline of the graphics processor.

[0084] **Figure 11A** is a block diagram illustrating a graphics processor command format 900 according to some embodiments. **Figure 11B** is a block diagram illustrating a graphics processor command sequence 910 according to an embodiment. The solid lined boxes in **Figure 11A** illustrate the components that are generally included in a graphics command while the dashed lines include components that are optional or that are only included in a sub-set of the graphics commands. The exemplary graphics processor command format 900 of **Figure 11A** includes data fields to identify a target client 902 of the command, a command operation code (opcode) 904, and the relevant data 906 for the command. A sub-opcode 905 and a command size 908 are also included in some commands.

[0085] In some embodiments, client 902 specifies the client unit of the graphics device that processes the command data. In some embodiments, a graphics

processor command parser examines the client field of each command to condition the further processing of the command and route the command data to the appropriate client unit. In some embodiments, the graphics processor client units include a memory interface unit, a render unit, a 2D unit, a 3D unit, and a media unit. Each client unit has a corresponding processing pipeline that processes the commands. Once the command is received by the client unit, the client unit reads the opcode 904 and, if present, sub-opcode 905 to determine the operation to perform. The client unit performs the command using information in data field 906. For some commands an explicit command size 908 is expected to specify the size of the command. In some embodiments, the command parser automatically determines the size of at least some of the commands based on the command opcode. In some embodiments commands are aligned via multiples of a double word.

[0086] The flow diagram in **Figure 11B** shows an exemplary graphics processor command sequence 910. In some embodiments, software or firmware of a data processing system that features an embodiment of a graphics processor uses a version of the command sequence shown to set up, execute, and terminate a set of graphics operations. A sample command sequence is shown and described for purposes of example only as embodiments are not limited to these specific commands or to this command sequence. Moreover, the commands may be issued as batch of commands in a command sequence, such that the graphics processor will process the sequence of commands in at least partially concurrence.

[0087] In some embodiments, the graphics processor command sequence 910 may begin with a pipeline flush command 912 to cause any active graphics pipeline to complete the currently pending commands for the pipeline. In some embodiments, the 3D pipeline 922 and the media pipeline 924 do not operate concurrently. The pipeline flush is performed to cause the active graphics pipeline to complete any pending commands. In response to a pipeline flush, the command parser for the graphics processor will pause command processing until the active drawing engines complete pending operations and the relevant read caches are invalidated. Optionally, any data in the render cache that is marked 'dirty' can be

flushed to memory. In some embodiments, pipeline flush command 912 can be used for pipeline synchronization or before placing the graphics processor into a low power state.

[0088] In some embodiments, a pipeline select command 913 is used when a command sequence requires the graphics processor to explicitly switch between pipelines. In some embodiments, a pipeline select command 913 is required only once within an execution context before issuing pipeline commands unless the context is to issue commands for both pipelines. In some embodiments, a pipeline flush command is 912 is required immediately before a pipeline switch via the pipeline select command 913.

[0089] In some embodiments, a pipeline control command 914 configures a graphics pipeline for operation and is used to program the 3D pipeline 922 and the media pipeline 924. In some embodiments, pipeline control command 914 configures the pipeline state for the active pipeline. In one embodiment, the pipeline control command 914 is used for pipeline synchronization and to clear data from one or more cache memories within the active pipeline before processing a batch of commands.

[0090] In some embodiments, return buffer state commands 916 are used to configure a set of return buffers for the respective pipelines to write data. Some pipeline operations require the allocation, selection, or configuration of one or more return buffers into which the operations write intermediate data during processing. In some embodiments, the graphics processor also uses one or more return buffers to store output data and to perform cross thread communication. In some embodiments, the return buffer state 916 includes selecting the size and number of return buffers to use for a set of pipeline operations.

[0091] The remaining commands in the command sequence differ based on the active pipeline for operations. Based on a pipeline determination 920, the command sequence is tailored to the 3D pipeline 922 beginning with the 3D pipeline state 930, or the media pipeline 924 beginning at the media pipeline state 940.

[0092] The commands for the 3D pipeline state 930 include 3D state setting commands for vertex buffer state, vertex element state, constant color state, depth buffer state, and other state variables that are to be configured before 3D primitive commands are processed. The values of these commands are determined at least in part based the particular 3D API in use. In some embodiments, 3D pipeline state 930 commands are also able to selectively disable or bypass certain pipeline elements if those elements will not be used.

[0093] In some embodiments, 3D primitive 932 command is used to submit 3D primitives to be processed by the 3D pipeline. Commands and associated parameters that are passed to the graphics processor via the 3D primitive 932 command are forwarded to the vertex fetch function in the graphics pipeline. The vertex fetch function uses the 3D primitive 932 command data to generate vertex data structures. The vertex data structures are stored in one or more return buffers. In some embodiments, 3D primitive 932 command is used to perform vertex operations on 3D primitives via vertex shaders. To process vertex shaders, 3D pipeline 922 dispatches shader execution threads to graphics processor execution units.

[0094] In some embodiments, 3D pipeline 922 is triggered via an execute 934 command or event. In some embodiments, a register write triggers command execution. In some embodiments execution is triggered via a 'go' or 'kick' command in the command sequence. In one embodiment command execution is triggered using a pipeline synchronization command to flush the command sequence through the graphics pipeline. The 3D pipeline will perform geometry processing for the 3D primitives. Once operations are complete, the resulting geometric objects are rasterized and the pixel engine colors the resulting pixels. Additional commands to control pixel shading and pixel back end operations may also be included for those operations.

[0095] In some embodiments, the graphics processor command sequence 910 follows the media pipeline 924 path when performing media operations. In general, the specific use and manner of programming for the media pipeline 924 depends on the media or compute operations to be performed. Specific media decode

operations may be offloaded to the media pipeline during media decode. In some embodiments, the media pipeline can also be bypassed and media decode can be performed in whole or in part using resources provided by one or more general purpose processing cores. In one embodiment, the media pipeline also includes elements for general-purpose graphics processor unit (GPGPU) operations, where the graphics processor is used to perform SIMD vector operations using computational shader programs that are not explicitly related to the rendering of graphics primitives.

[0096] In some embodiments, media pipeline 924 is configured in a similar manner as the 3D pipeline 922. A set of media pipeline state commands 940 are dispatched or placed into in a command queue before the media object commands 942. In some embodiments, media pipeline state commands 940 include data to configure the media pipeline elements that will be used to process the media objects. This includes data to configure the video decode and video encode logic within the media pipeline, such as encode or decode format. In some embodiments, media pipeline state commands 940 also support the use one or more pointers to “indirect” state elements that contain a batch of state settings.

[0097] In some embodiments, media object commands 942 supply pointers to media objects for processing by the media pipeline. The media objects include memory buffers containing video data to be processed. In some embodiments, all media pipeline states must be valid before issuing a media object command 942. Once the pipeline state is configured and media object commands 942 are queued, the media pipeline 924 is triggered via an execute command 944 or an equivalent execute event (e.g., register write). Output from media pipeline 924 may then be post processed by operations provided by the 3D pipeline 922 or the media pipeline 924. In some embodiments, GPGPU operations are configured and executed in a similar manner as media operations.

[0098] **Figure 12** illustrates exemplary graphics software architecture for a data processing system 1000 according to some embodiments. In some embodiments, software architecture includes a 3D graphics application 1010, an operating system 1020, and at least one processor 1030. In some embodiments, processor 1030

includes a graphics processor 1032 and one or more general-purpose processor core(s) 1034. The graphics application 1010 and operating system 1020 each execute in the system memory 1050 of the data processing system.

[0099] In some embodiments, 3D graphics application 1010 contains one or more shader programs including shader instructions 1012. The shader language instructions may be in a high-level shader language, such as the High Level Shader Language (HLSL) or the OpenGL Shader Language (GLSL). The application also includes executable instructions 1014 in a machine language suitable for execution by the general-purpose processor core 1034. The application also includes graphics objects 1016 defined by vertex data.

[0100] In some embodiments, operating system 1020 is a Microsoft® Windows® operating system from the Microsoft Corporation, a proprietary UNIX-like operating system, or an open source UNIX-like operating system using a variant of the Linux kernel. When the Direct3D API is in use, the operating system 1020 uses a front-end shader compiler 1024 to compile any shader instructions 1012 in HLSL into a lower-level shader language. The compilation may be a just-in-time (JIT) compilation or the application can perform shader pre-compilation. In some embodiments, high-level shaders are compiled into low-level shaders during the compilation of the 3D graphics application 1010.

[0101] In some embodiments, user mode graphics driver 1026 contains a back-end shader compiler 1027 to convert the shader instructions 1012 into a hardware specific representation. When the OpenGL API is in use, shader instructions 1012 in the GLSL high-level language are passed to a user mode graphics driver 1026 for compilation. In some embodiments, user mode graphics driver 1026 uses operating system kernel mode functions 1028 to communicate with a kernel mode graphics driver 1029. In some embodiments, kernel mode graphics driver 1029 communicates with graphics processor 1032 to dispatch commands and instructions.

[0102] One or more aspects of at least one embodiment may be implemented by representative code stored on a machine-readable medium which represents and/or defines logic within an integrated circuit such as a processor. For example, the

machine-readable medium may include instructions which represent various logic within the processor. When read by a machine, the instructions may cause the machine to fabricate the logic to perform the techniques described herein. Such representations, known as "IP cores," are reusable units of logic for an integrated circuit that may be stored on a tangible, machine-readable medium as a hardware model that describes the structure of the integrated circuit. The hardware model may be supplied to various customers or manufacturing facilities, which load the hardware model on fabrication machines that manufacture the integrated circuit. The integrated circuit may be fabricated such that the circuit performs operations described in association with any of the embodiments described herein.

[0103] **Figure 13** is a block diagram illustrating an IP core development system 1100 that may be used to manufacture an integrated circuit to perform operations according to an embodiment. The IP core development system 1100 may be used to generate modular, re-usable designs that can be incorporated into a larger design or used to construct an entire integrated circuit (e.g., an SOC integrated circuit). A design facility 1130 can generate a software simulation 1110 of an IP core design in a high level programming language (e.g., C/C++). The software simulation 1110 can be used to design, test, and verify the behavior of the IP core using a simulation model 1112. The simulation model 1112 may include functional, behavioral, and/or timing simulations. A register transfer level (RTL) design can then be created or synthesized from the simulation model 1112. The RTL design 1115 is an abstraction of the behavior of the integrated circuit that models the flow of digital signals between hardware registers, including the associated logic performed using the modeled digital signals. In addition to an RTL design 1115, lower-level designs at the logic level or transistor level may also be created, designed, or synthesized. Thus, the particular details of the initial design and simulation may vary.

[0104] The RTL design 1115 or equivalent may be further synthesized by the design facility into a hardware model 1120, which may be in a hardware description language (HDL), or some other representation of physical design data. The HDL may be further simulated or tested to verify the IP core design. The IP core design can be stored for delivery to a 3<sup>rd</sup> party fabrication facility 1165 using non-volatile



memory 1140 (e.g., hard disk, flash memory, or any non-volatile storage medium). Alternatively, the IP core design may be transmitted (e.g., via the Internet) over a wired connection 1150 or wireless connection 1160. The fabrication facility 1165 may then fabricate an integrated circuit that is based at least in part on the IP core design. The fabricated integrated circuit can be configured to perform operations in accordance with at least one embodiment described herein.

[0105] **Figure 14** is a block diagram illustrating an exemplary system on a chip integrated circuit 1200 that may be fabricated using one or more IP cores, according to an embodiment. The exemplary integrated circuit includes one or more application processors 1205 (e.g., CPUs), at least one graphics processor 1210, and may additionally include an image processor 1215 and/or a video processor 1220, any of which may be a modular IP core from the same or multiple different design facilities. The integrated circuit includes peripheral or bus logic including a USB controller 1225, UART controller 1230, an SPI/SDIO controller 1235, and an I<sup>2</sup>S/I<sup>2</sup>C controller 1240. Additionally, the integrated circuit can include a display device 1245 coupled to one or more of a high-definition multimedia interface (HDMI) controller 1250 and a mobile industry processor interface (MIPI) display interface 1255. Storage may be provided by a flash memory subsystem 1260 including flash memory and a flash memory controller. Memory interface may be provided via a memory controller 1265 for access to SDRAM or SRAM memory devices. Some integrated circuits additionally include an embedded security engine 1270.

[0106] Additionally, other logic and circuits may be included in the processor of integrated circuit 1200, including additional graphics processors/cores, peripheral interface controllers, or general purpose processor cores.

[0107] The following clauses and/or examples pertain to further embodiments:

One example embodiment may be a first processing unit implemented method comprising initiating a continuous thread on a second processing unit, collecting packets to be processed on the second processing unit, storing the packets in a shared virtual memory, and notifying the second processing unit of the availability of the packets. The method may also include wherein said first

processing unit is a central processing unit and said second processing unit is a graphics processing unit. The method may also include wherein notifying includes updating a graphics processing unit visible flag. The method may also include determining whether to terminate the continuous thread. The method may also include terminating the thread by setting a graphics processing unit visible flag. The method may also include terminating the thread from an application running on the central processing unit.

[0108] Another example embodiment may be a first processing unit implemented method comprising receiving a signal from a second processing unit and in response initiating a continuous thread, in response to a signal from the first processing unit, accessing packets to be processed in a shared virtual memory, and notifying the second processing unit when the packets have been processed. The method may also include wherein said first processing unit is a graphics processing unit and said second processing unit is a central processing unit. The method may also include wherein notifying includes setting a flag visible to the central processing unit. The method may also include storing packet processing results in said shared virtual memory.

[0109] In another example embodiment one or more non-transitory computer readable media storing instructions executed by a first processing unit to perform a sequence comprising initiating a continuous thread on a second processing unit, collecting packets to be processed on the second processing unit, storing the packets in a shared virtual memory, and notifying the graphics processing unit of the availability of the packets. The media may include wherein said first processing unit is a graphics processing unit and said second processing unit is a central processing unit. The media may include wherein notifying includes updating a graphics processing unit visible flag. The media may include said sequence including determining whether to terminate the continuous thread. The media may include said sequence including terminating the thread by setting a graphics processing unit visible flag. The media may include said sequence including terminating the thread from an application running on the central processing unit.

[0110] Another example embodiment may be an apparatus comprising a first processing unit to receive a signal from a second processing unit and in response initiate a continuous thread, in response to a signal from the second processing unit, access packets to be processed in a shared virtual memory, and notify the second processing unit when the packets have been processed, and a storage coupled to said graphics processing unit. The apparatus may include wherein said first processing unit is a graphics processing unit and said second processing unit is a central processing unit. The apparatus may include wherein notifying includes setting a flag visible to the central processing unit. The apparatus may include said processor to store packet processing results in said shared virtual memory. The apparatus may include a central processing unit to initiate a continuous thread on the graphics processing unit, collect packets to be processed on the graphics processing unit, store the packets in the shared virtual memory, and notify the graphics processing unit of the availability of the packets. The apparatus may include wherein notifying includes updating a graphics processing unit visible flag. The apparatus may include said central processing unit to determine whether to terminate the continuous thread. The apparatus may include said central processing unit to terminate the thread by setting a graphics processing unit visible flag. The apparatus may include said central processing unit to terminate the thread from an application running on the central processing unit.

[0111] The graphics processing techniques described herein may be implemented in various hardware architectures. For example, graphics functionality may be integrated within a chipset. Alternatively, a discrete graphics processor may be used. As still another embodiment, the graphics functions may be implemented by a general purpose processor, including a multicore processor.

[0112] References throughout this specification to “one embodiment” or “an embodiment” mean that a particular feature, structure, or characteristic described in connection with the embodiment is included in at least one implementation encompassed within the present disclosure. Thus, appearances of the phrase “one embodiment” or “in an embodiment” are not necessarily referring to the same embodiment. Furthermore, the particular features, structures, or characteristics may

be instituted in other suitable forms other than the particular embodiment illustrated and all such forms may be encompassed within the claims of the present application.

[0113] While a limited number of embodiments have been described, those skilled in the art will appreciate numerous modifications and variations therefrom. It is intended that the appended claims cover all such modifications and variations as fall within the true spirit and scope of this disclosure.

What is claimed is:

- 1 1. A first processing unit implemented method comprising:  
2 initiating a continuous thread on a second processing unit;  
3 collecting packets to be processed on the second processing unit;  
4 storing the packets in a shared virtual memory; and  
5 notifying the second processing unit of the availability of the packets.
- 1 2. The method of claim 1 wherein said first processing unit is a central  
2 processing unit and said second processing unit is a graphics processing unit.
- 1 3. The method of claim 2 wherein notifying includes updating a graphics  
2 processing unit visible flag.
- 1 4. The method of claim 2 including determining whether to terminate the  
2 continuous thread.
- 1 5. The method of claim 4 including terminating the thread by setting a graphics  
2 processing unit visible flag.
- 1 6. The method of claim 2 including terminating the thread from an application  
2 running on the central processing unit.
- 1 7. A first processing unit implemented method comprising:  
2 receiving a signal from a second processing unit and in response initiating a  
3 continuous thread;  
4 in response to a signal from the first processing unit, accessing packets to be  
5 processed in a shared virtual memory; and  
6 notifying the second processing unit when the packets have been processed.
- 1 8. The method of claim 7 wherein said first processing unit is a graphics  
2 processing unit and said second processing unit is a central processing unit.

- 1 9. The method of claim 8 wherein notifying includes setting a flag visible to the  
2 central processing unit.
- 1 10. The method of claim 8 including storing packet processing results in said  
2 shared virtual memory.
- 1 11. One or more non-transitory computer readable media storing instructions  
2 executed by a first processing unit to perform a sequence comprising:  
3 initiating a continuous thread on a second processing unit;  
4 collecting packets to be processed on the second processing unit;  
5 storing the packets in a shared virtual memory; and  
6 notifying the graphics processing unit of the availability of the packets.
- 1 12. The media of claim 11 wherein said first processing unit is a graphics  
2 processing unit and said second processing unit is a central processing unit.
- 1 13. The media of claim 12 wherein notifying includes updating a graphics  
2 processing unit visible flag.
- 1 14. The media of claim 12, said sequence including determining whether to  
2 terminate the continuous thread.
- 1 15. The media of claim 14, said sequence including terminating the thread by  
2 setting a graphics processing unit visible flag.
- 1 16. The media of claim 12, said sequence including terminating the thread from  
2 an application running on the central processing unit.
- 1 17. An apparatus comprising:  
2 a first processing unit to receive a signal from a second processing unit and in  
3 response initiate a continuous thread, in response to a signal from the second

4 processing unit, access packets to be processed in a shared virtual memory, and  
5 notify the second processing unit when the packets have been processed; and  
6 a storage coupled to said graphics processing unit.

1 18. The apparatus of claim 17 wherein said first processing unit is a graphics  
2 processing unit and said second processing unit is a central processing unit.

1 19. The apparatus of claim 18 wherein notifying includes setting a flag visible to  
2 the central processing unit.

1 20. The apparatus of claim 18, said processor to store packet processing results  
2 in said shared virtual memory.

1 21. The apparatus of claim 18 including a central processing unit to initiate a  
2 continuous thread on the graphics processing unit, collect packets to be processed  
3 on the graphics processing unit, store the packets in the shared virtual memory, and  
4 notify the graphics processing unit of the availability of the packets.

1 22. The apparatus of claim 21 wherein notifying includes updating a graphics  
2 processing unit visible flag.

1 23. The apparatus of claim 21, said central processing unit to determine whether  
2 to terminate the continuous thread.

1 24. The apparatus of claim 23, said central processing unit to terminate the thread  
2 by setting a graphics processing unit visible flag.

1 25. The apparatus of claim 21, said central processing unit to terminate the thread  
2 from an application running on the central processing unit.

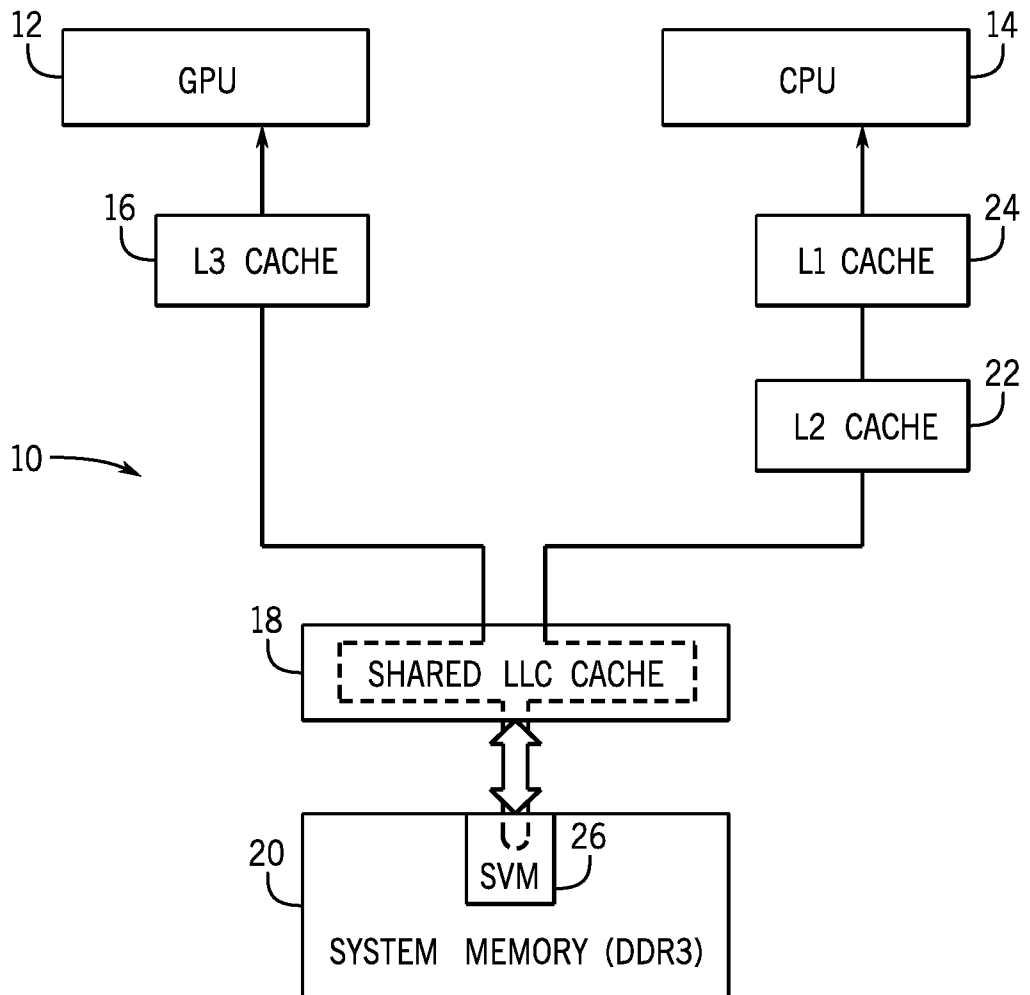


FIG. 1



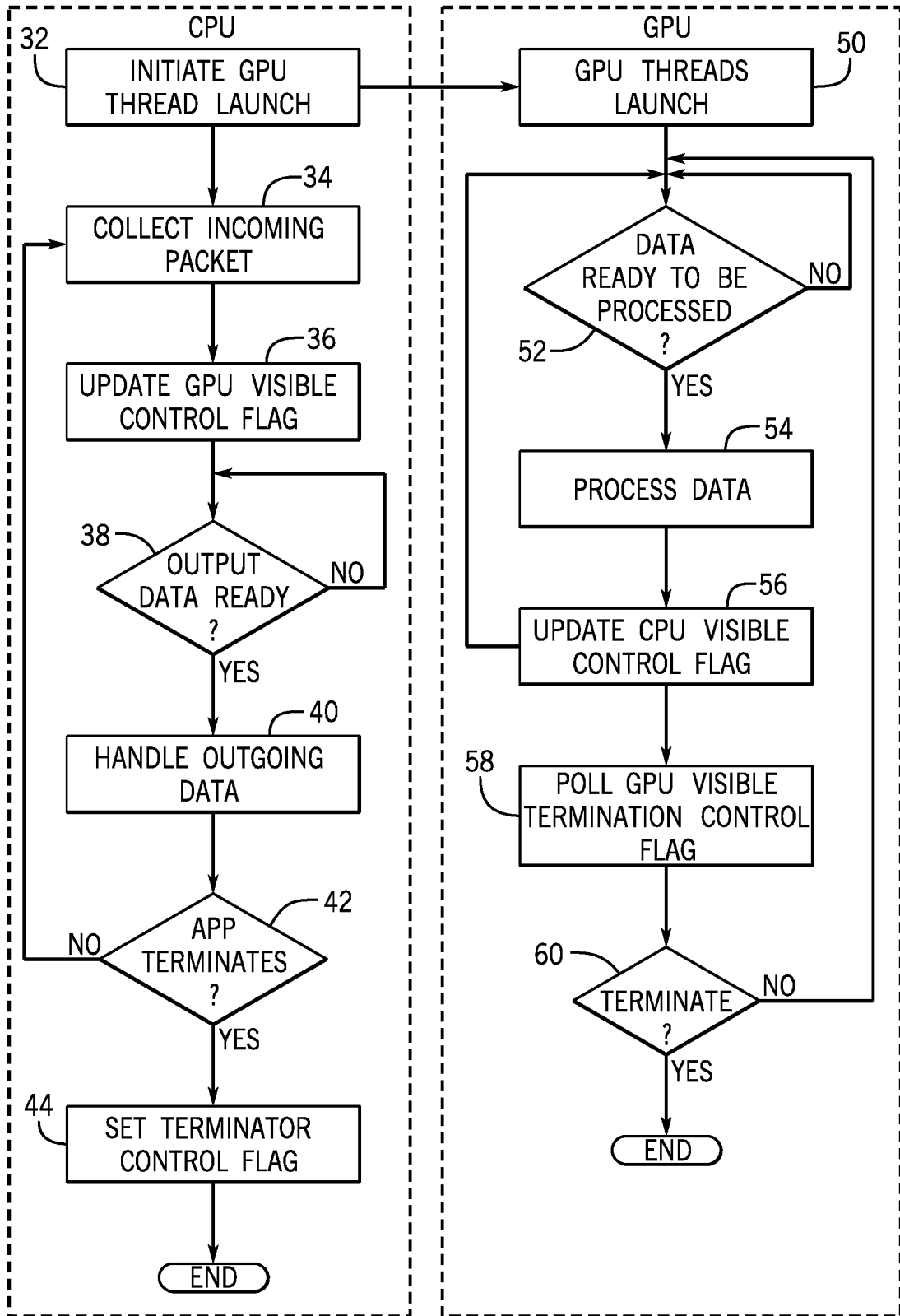


FIG. 2

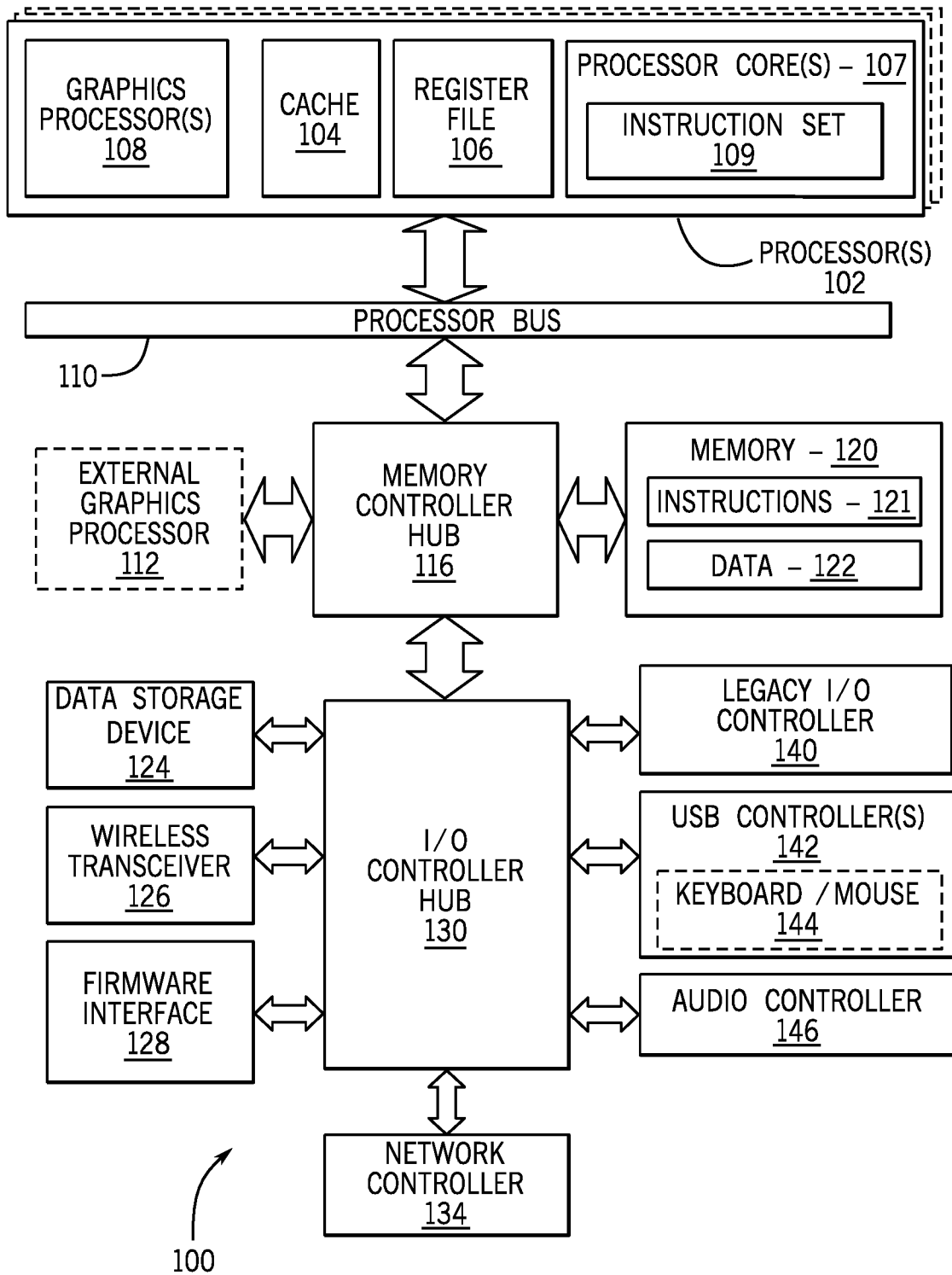


FIG. 3

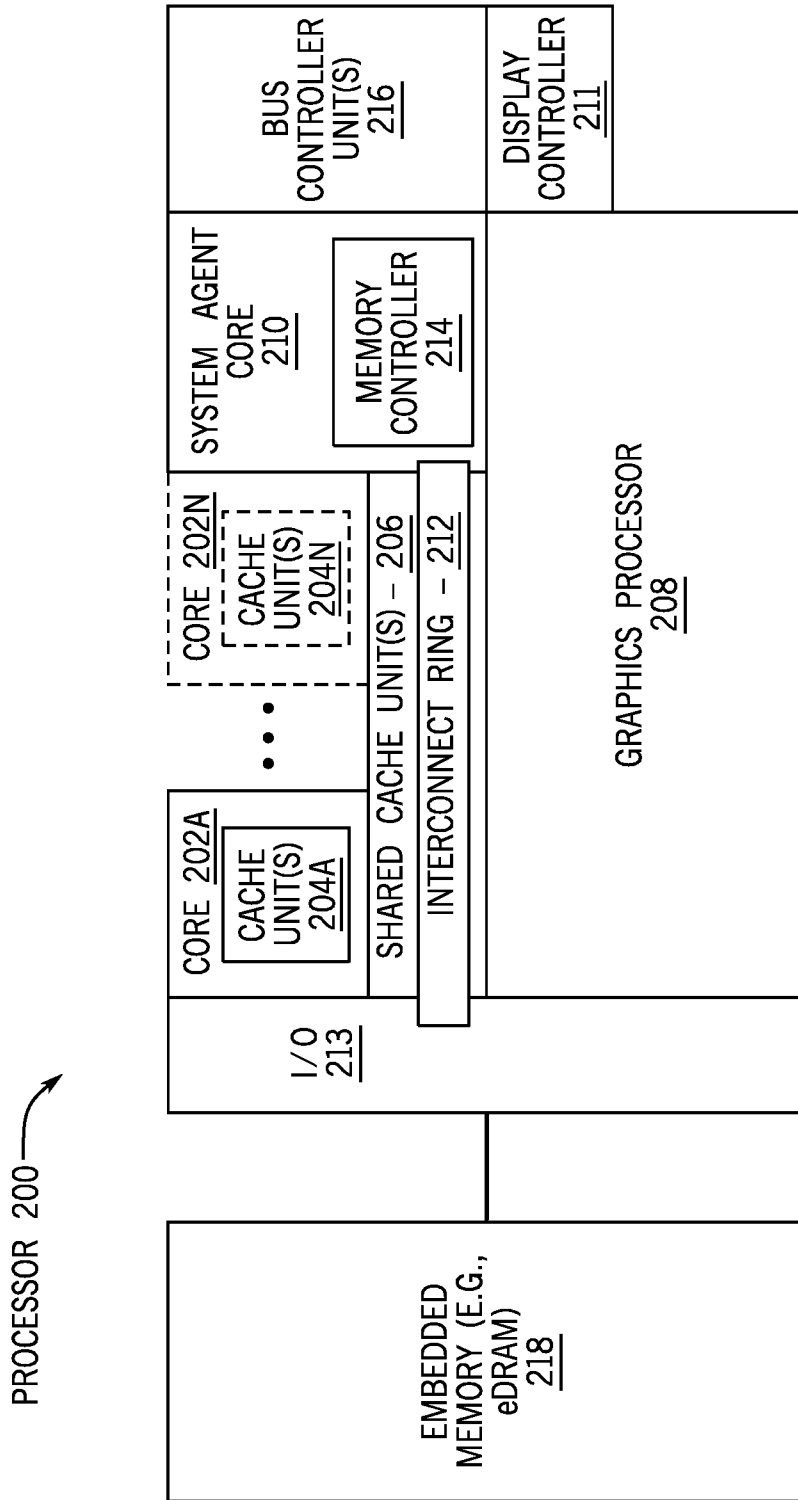


FIG. 4

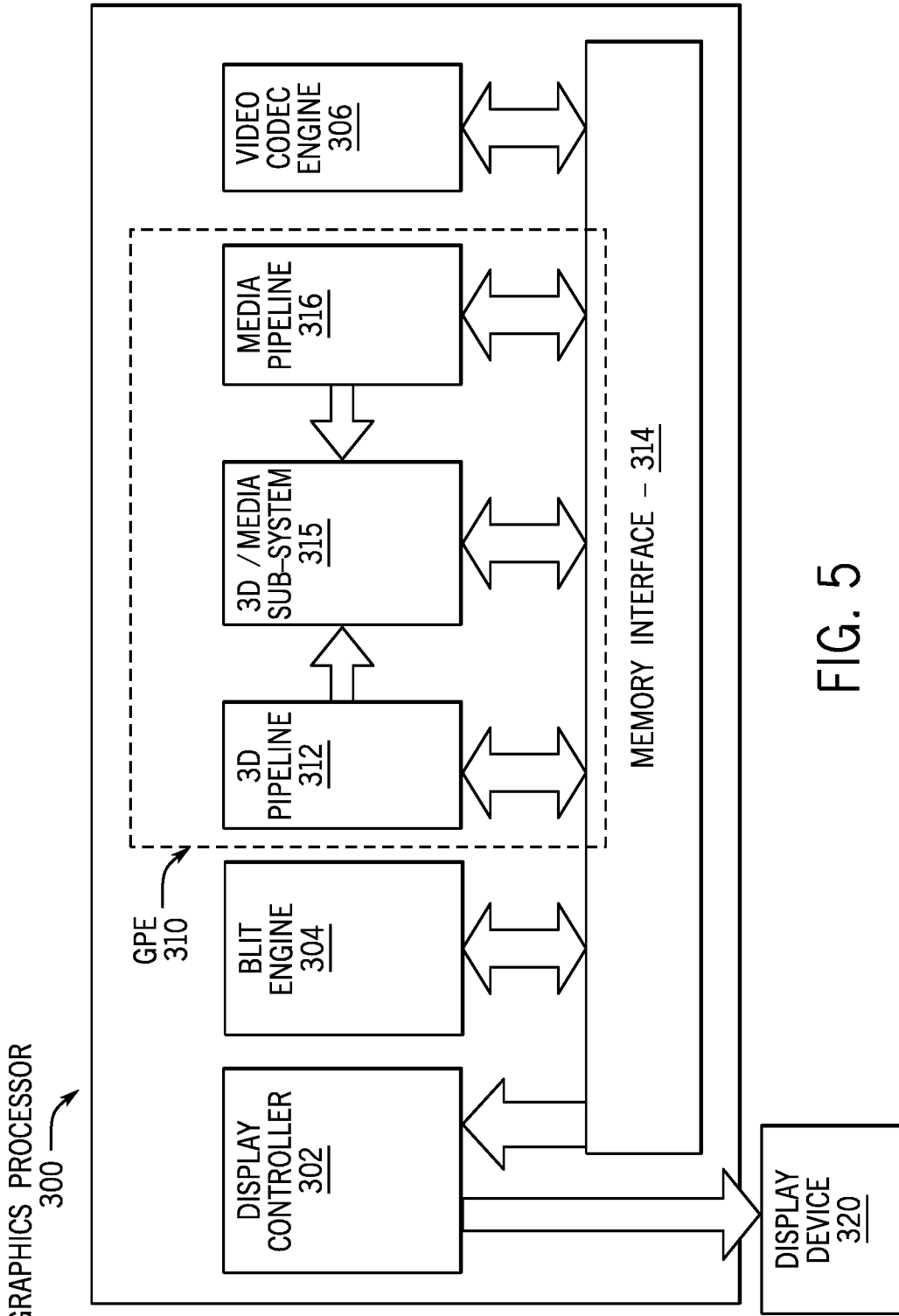


FIG. 5

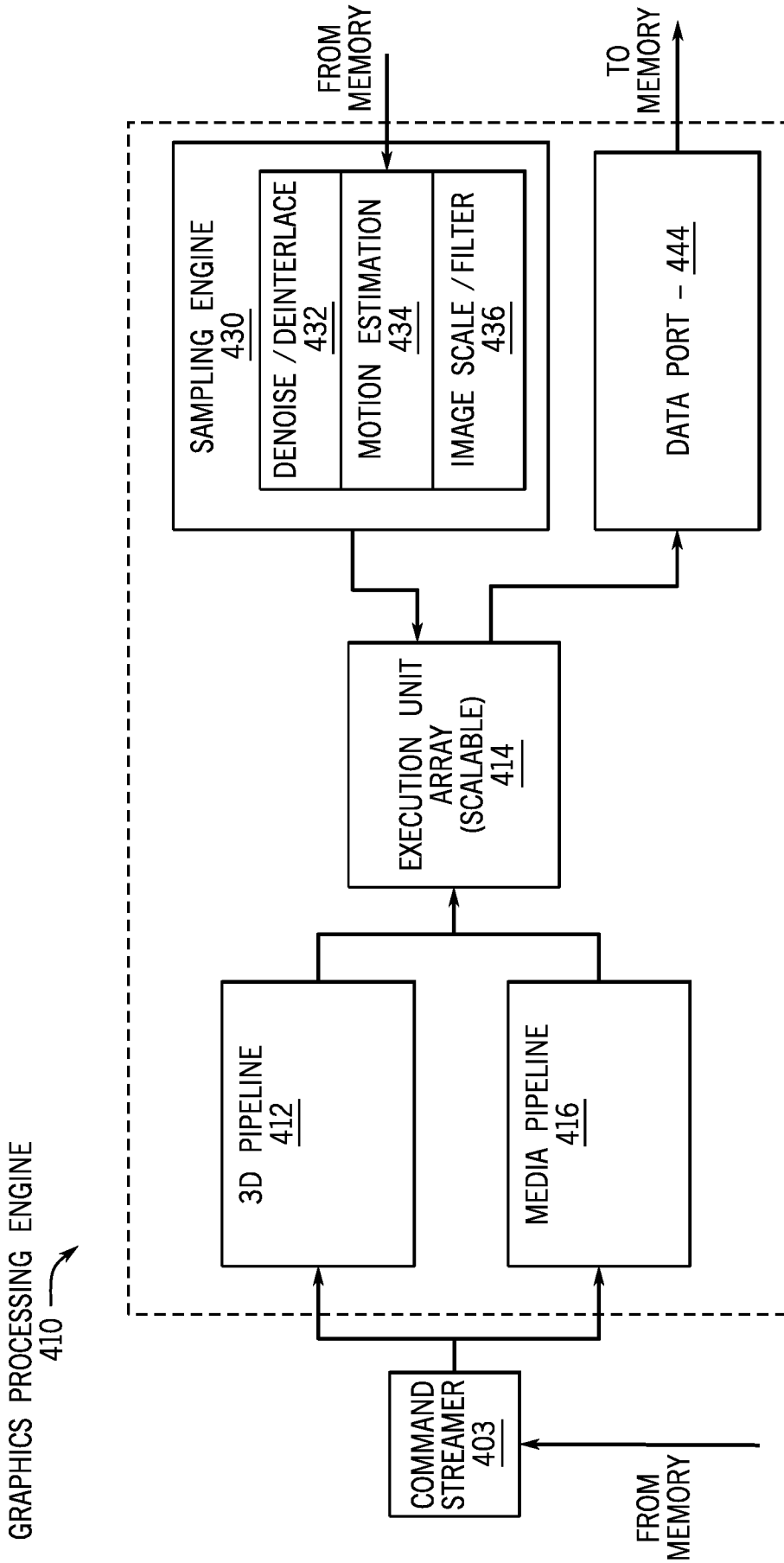


FIG. 6

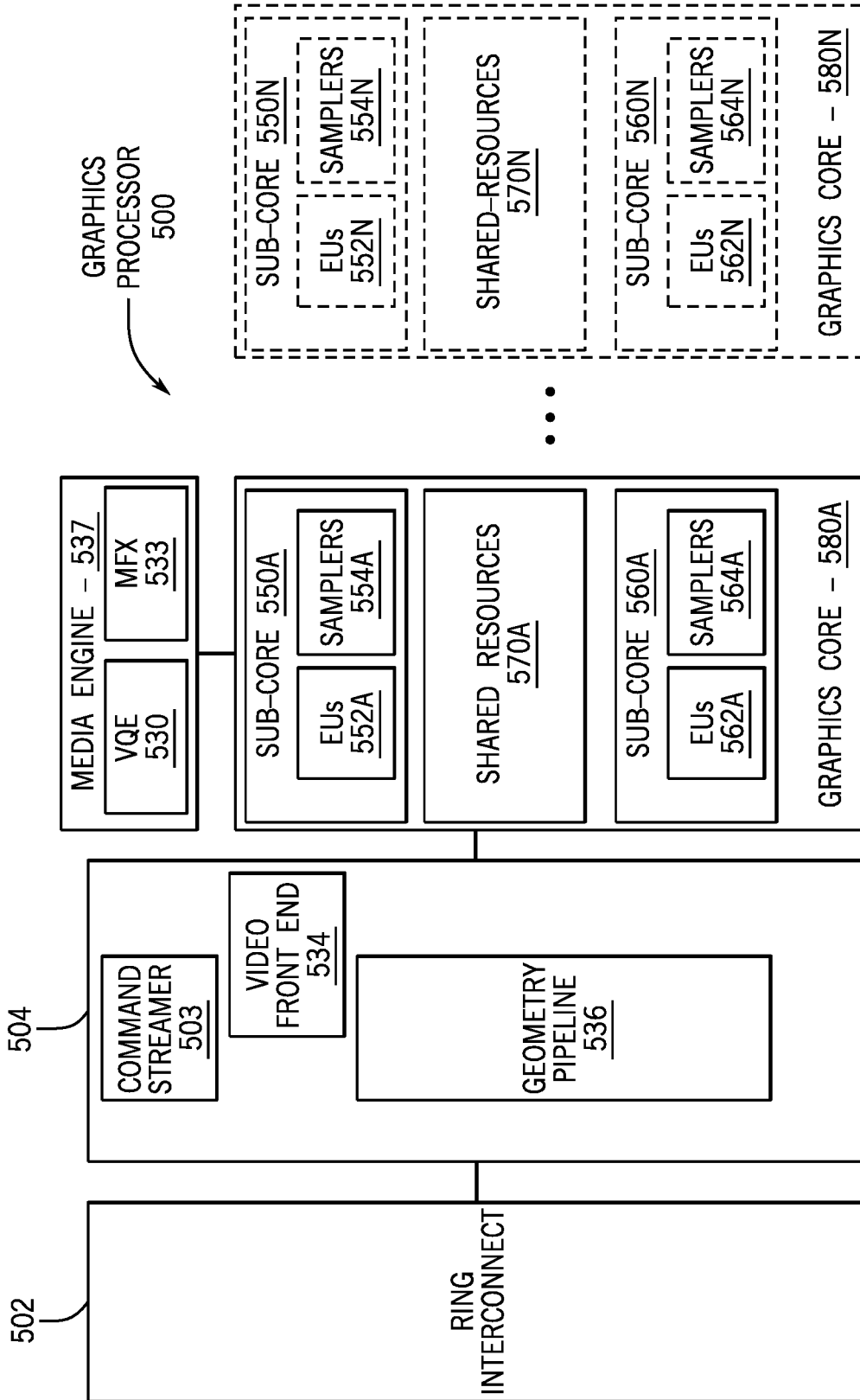


FIG. 7

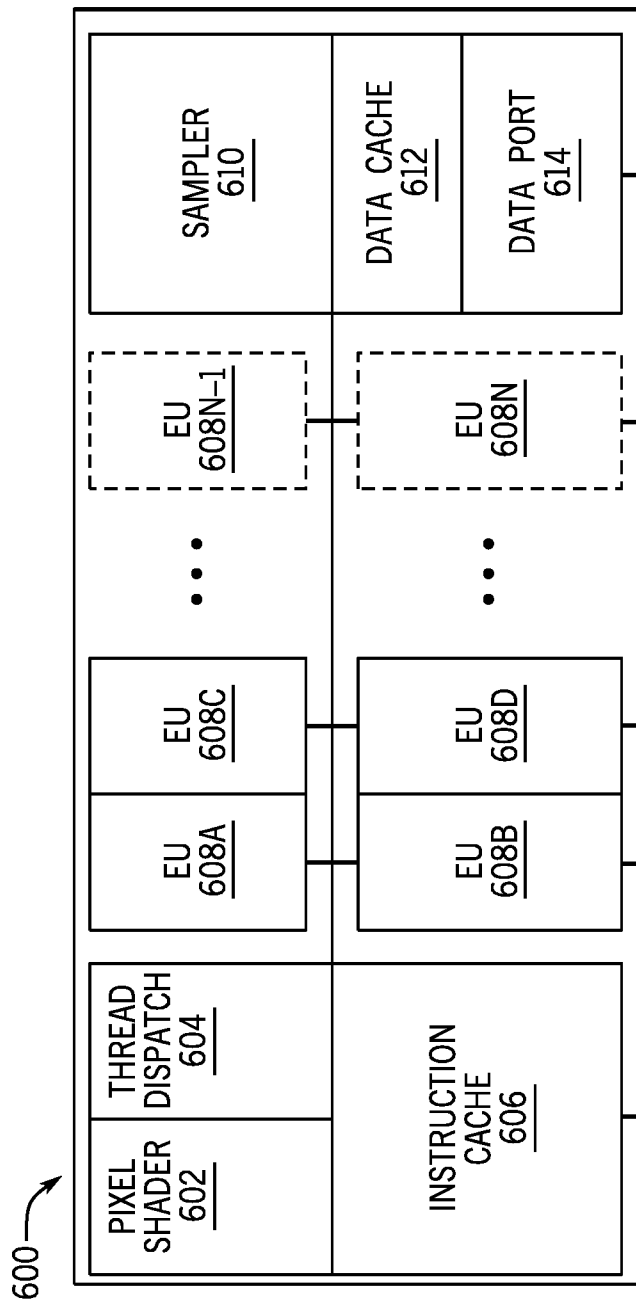


FIG. 8

GRAPHICS CORE INSTRUCTION FORMATS

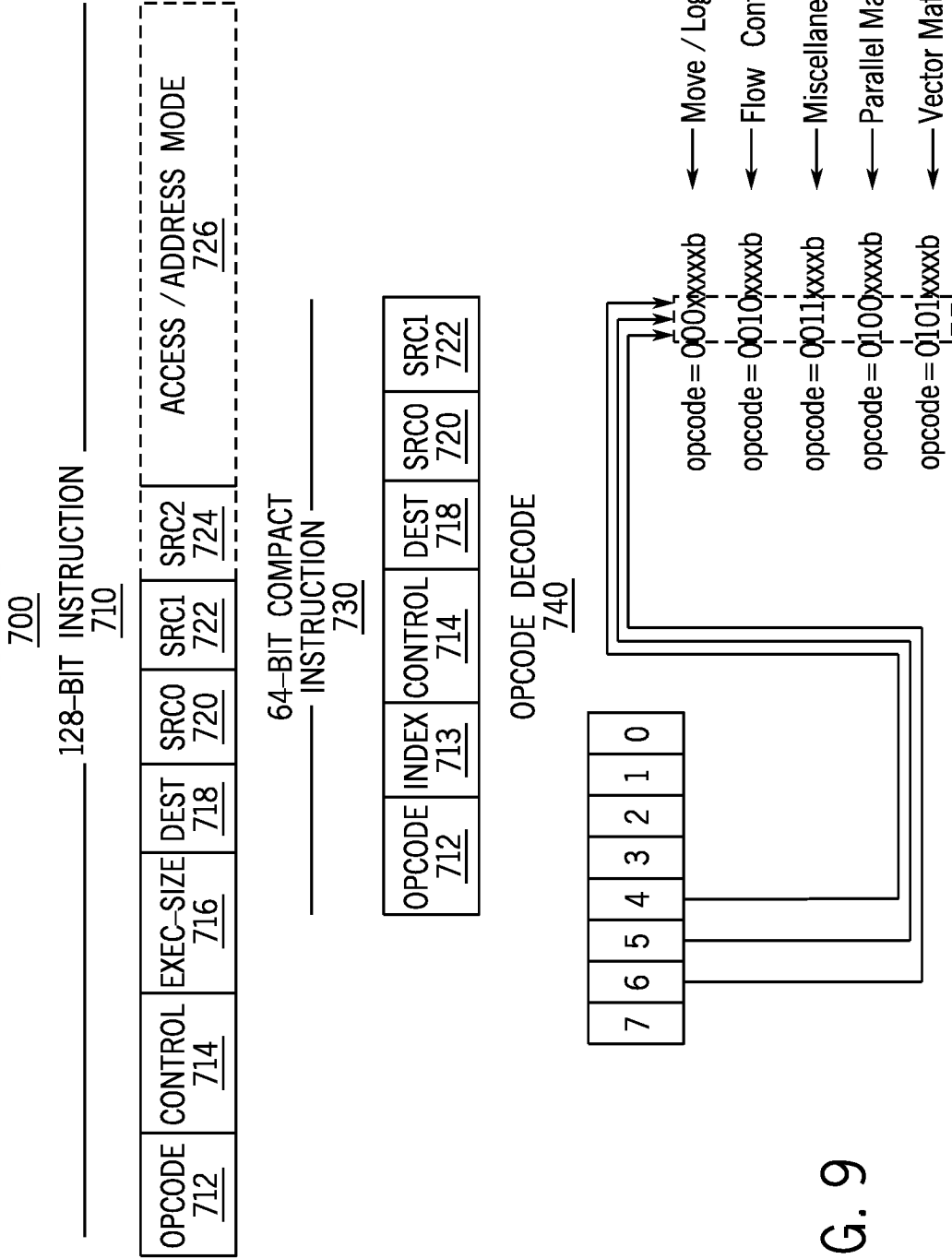


FIG. 9



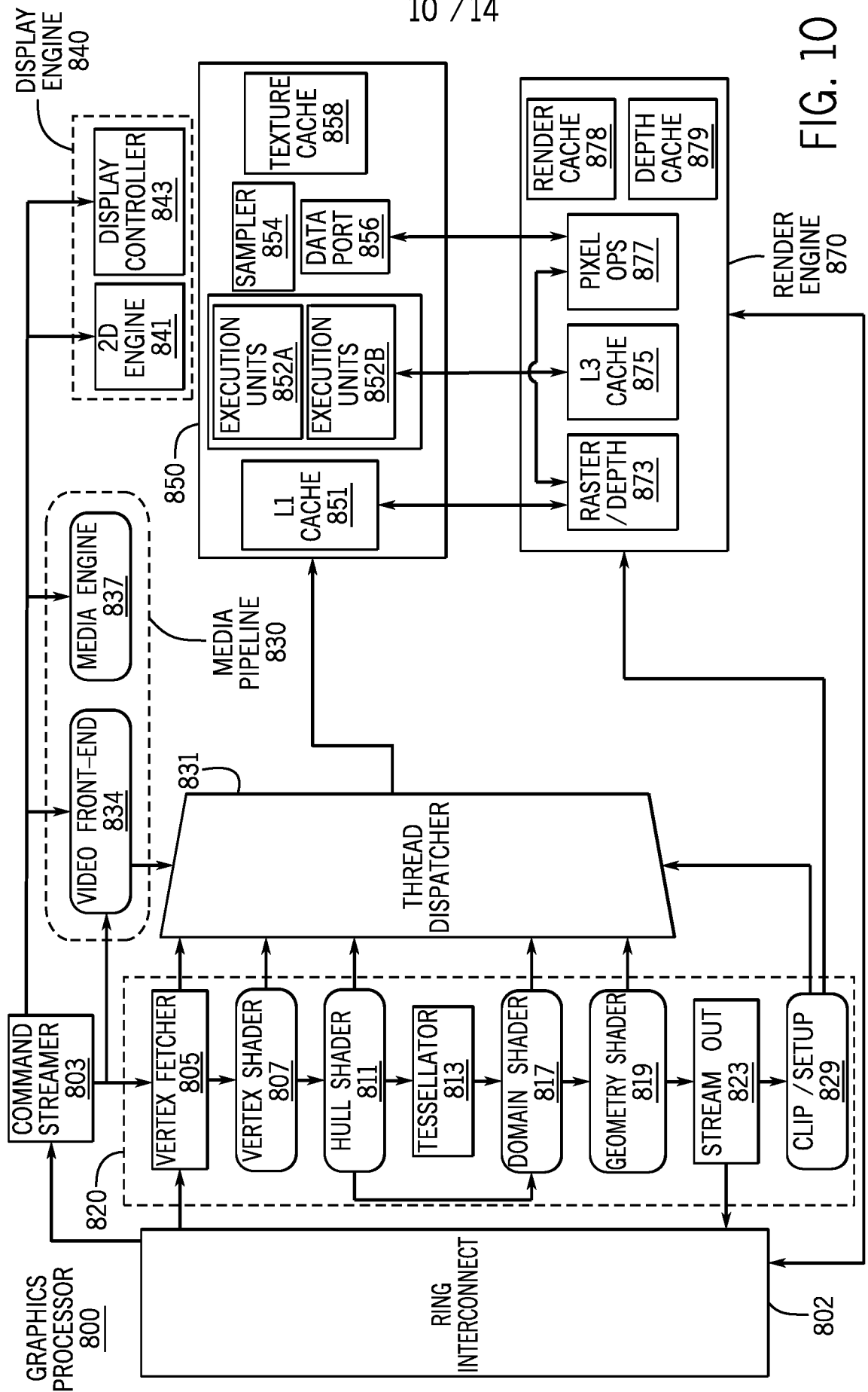


FIG. 10

GRAPHICS PROCESSOR COMMAND FORMAT  
900

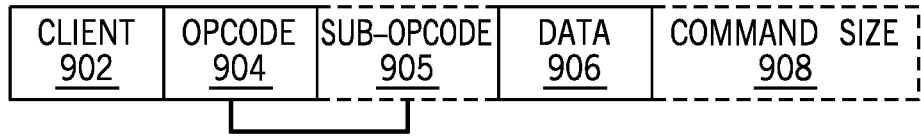


FIG. 11A

GRAPHICS PROCESSOR COMMAND SEQUENCE  
910

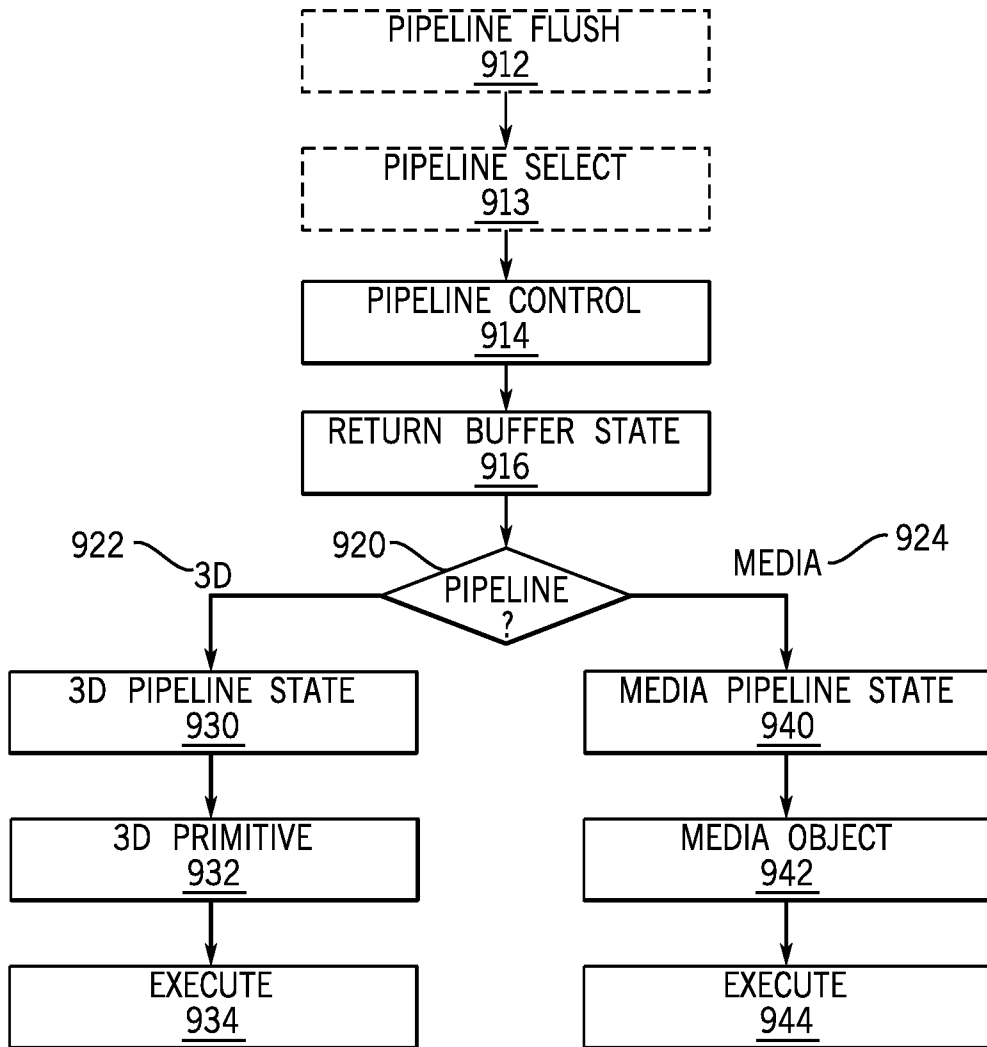


FIG. 11B

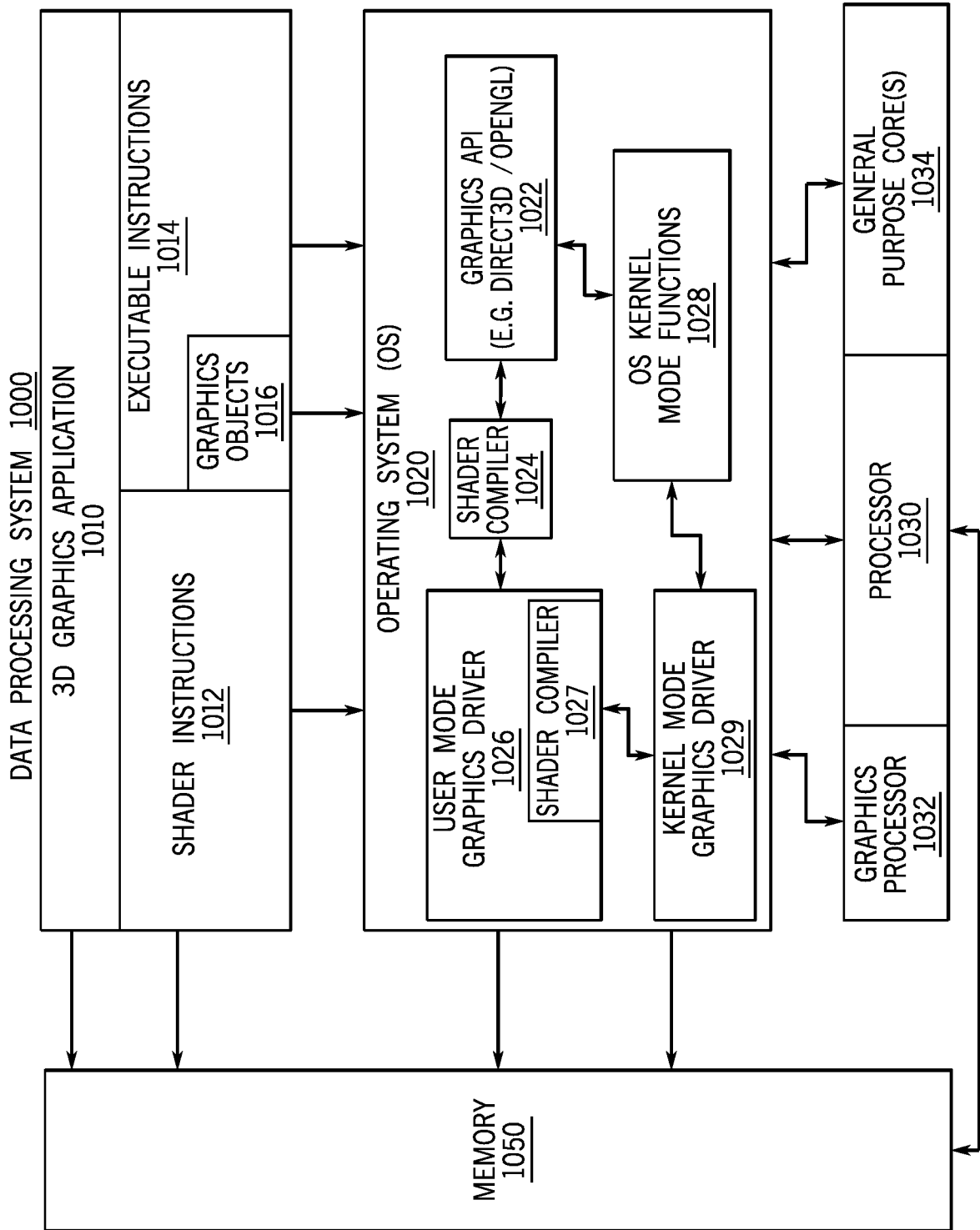


FIG. 12

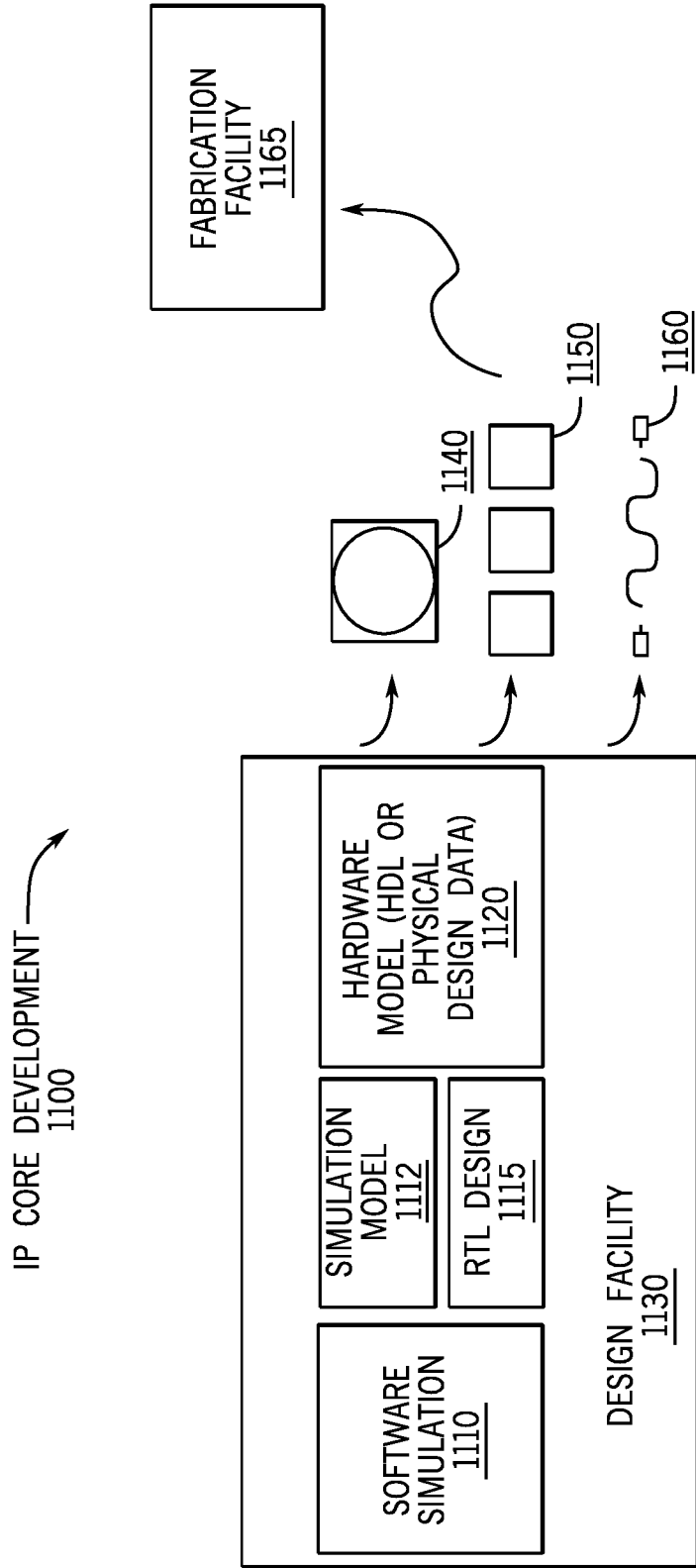


FIG. 13

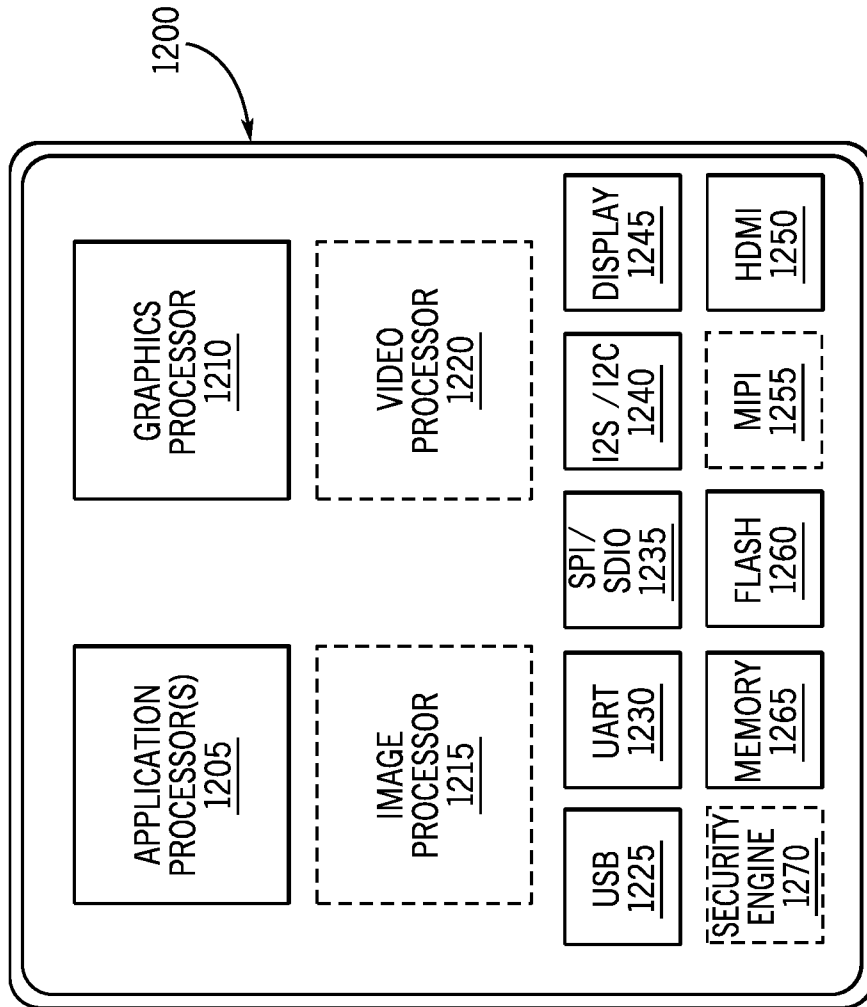


FIG. 14

**A. CLASSIFICATION OF SUBJECT MATTER****G06F 9/30(2006.01)I, G06F 9/38(2006.01)I, G06F 12/08(2006.01)I, G06T 1/20(2006.01)I**

According to International Patent Classification (IPC) or to both national classification and IPC

**B. FIELDS SEARCHED**

Minimum documentation searched (classification system followed by classification symbols)

G06F 9/30; G06F 15/16; G06F 12/02; G09G 5/36; G06F 12/00; G06F 15/167; G06T 1/00; G06F 9/38; G06F 12/08; G06T 1/20

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Korean utility models and applications for utility models

Japanese utility models and applications for utility models

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

eKOMPASS(KIPO internal) &amp; keywords: continuous thread, CPU, GPU, visible flag, launch, shared virtual memory, packet, terminate

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	US 2006-0098022 A1 (JEFFERY A. ANDREWS et al.) 11 May 2006 See paragraphs [0063], [0068]-[0075], [0091]-[0095], [0102]-[0105], [0118]-[0121]; claims 9, 18; and figures 5, 10-12, 15.	1-3, 7-9, 11-13, 17-19, 21-22
Y		4-6, 10, 14-16, 20, 23-25
Y	US 2009-0313440 A1 (YOUNG LAK KIM et al.) 17 December 2009 See paragraphs [0064]-[0071], [0093]-[0095]; and figures 13-16, 22.	4-6, 14-16, 23-25
Y	US 2010-0045682 A1 (SIMON ANDREW FORD et al.) 25 February 2010 See paragraphs [0003]-[0008], [0048]-[0054]; claim 1; and figures 1-3B.	10, 20
A	US 2013-0027410 A1 (BORIS GINZBURG et al.) 31 January 2013 See paragraphs [0011]-[0027]; claim 1; and figures 1-3.	1-25
A	US 2008-0276064 A1 (AFTAB MUNSHI et al.) 06 November 2008 See paragraphs [0033]-[0055]; and figures 2-7.	1-25

 Further documents are listed in the continuation of Box C. See patent family annex.

\* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier application or patent but published on or after the international filing date

"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&amp;" document member of the same patent family

Date of the actual completion of the international search

22 September 2016 (22.09.2016)

Date of mailing of the international search report

**23 September 2016 (23.09.2016)**

Name and mailing address of the ISA/KR

International Application Division

Korean Intellectual Property Office

189 Cheongsa-ro, Seo-gu, Daejeon, 35208, Republic of Korea

Facsimile No. +82-42-481-8578

Authorized officer

CHIN, Sang Bum

Telephone No. +82-42-481-8398



**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No.

**PCT/US2016/036632**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
US 2006-0098022 A1	11/05/2006	EP 1498824 A2	19/01/2005
		EP 1498824 A3	20/05/2009
		JP 2005-025749 A	27/01/2005
		JP 3886988 B2	28/02/2007
		US 2004-0263519 A1	30/12/2004
		US 2005-0122339 A1	09/06/2005
		US 2006-0095672 A1	04/05/2006
		US 6862027 B2	01/03/2005
		US 7333114 B2	19/02/2008
		US 7355601 B2	08/04/2008
		US 2009-0313440 A1	17/12/2009
KR 10-2009-0128605 A	16/12/2009		
US 8230180 B2	24/07/2012		
US 2010-0045682 A1	25/02/2010	CN 101667284 A	10/03/2010
		CN 101667284 B	05/06/2013
		GB 0815442 D0	01/10/2008
		GB 0818203 D0	12/11/2008
		GB 2462860 A	24/02/2010
		GB 2462860 B	16/05/2012
		JP 2010-050970 A	04/03/2010
		US 8675006 B2	18/03/2014
US 2013-0027410 A1	31/01/2013	CN 103718156 A	09/04/2014
		EP 2737396 A2	04/06/2014
		EP 2737396 A4	10/06/2015
		JP 2014-522038 A	28/08/2014
		JP 5933000 B2	08/06/2016
		WO 2013-019350 A2	07/02/2013
		WO 2013-019350 A3	10/05/2013
US 2008-0276064 A1	06/11/2008	AU 2008-239696 A1	23/10/2008
		AU 2008-239696 B2	08/09/2011
		AU 2008-239697 A1	23/10/2008
		AU 2008-239697 B2	13/10/2011
		CN 101657795 A	24/02/2010
		CN 101657795 B	23/10/2013
		CN 101802789 A	11/08/2010
		CN 101802789 B	07/05/2014
		CN 103927150 A	16/07/2014
		EP 2135163 A2	23/12/2009
		EP 2140352 A2	06/01/2010
		EP 2146283 A2	20/01/2010
		EP 2146283 A3	25/05/2016
		US 2008-0276220 A1	06/11/2008
		US 2008-0276261 A1	06/11/2008
		US 2008-0276262 A1	06/11/2008
		US 8108633 B2	31/01/2012

**INTERNATIONAL SEARCH REPORT**

Information on patent family members

International application No.

**PCT/US2016/036632**

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
		US 8276164 B2	25/09/2012
		US 8286196 B2	09/10/2012
		US 8341611 B2	25/12/2012
		WO 2008-127604 A2	23/10/2008
		WO 2008-127604 A3	05/02/2009
		WO 2008-127610 A2	23/10/2008
		WO 2008-127622 A2	23/10/2008
		WO 2008-127623 A2	23/10/2008