

Qimonda GDDR5 – White Paper

August 2007

1. Introduction

GDDR5 introduces features and functions that go beyond previous GDDR standards and enables GDDR5 to operate at data rates of 5Gbps. The GDDR5 combines highest performance with stable system operation and low implementation costs. The article discusses the new GDDR5 features and highlights their motivation and benefit for the application.

GDDR5 – key elements for reliable high speed data transmission

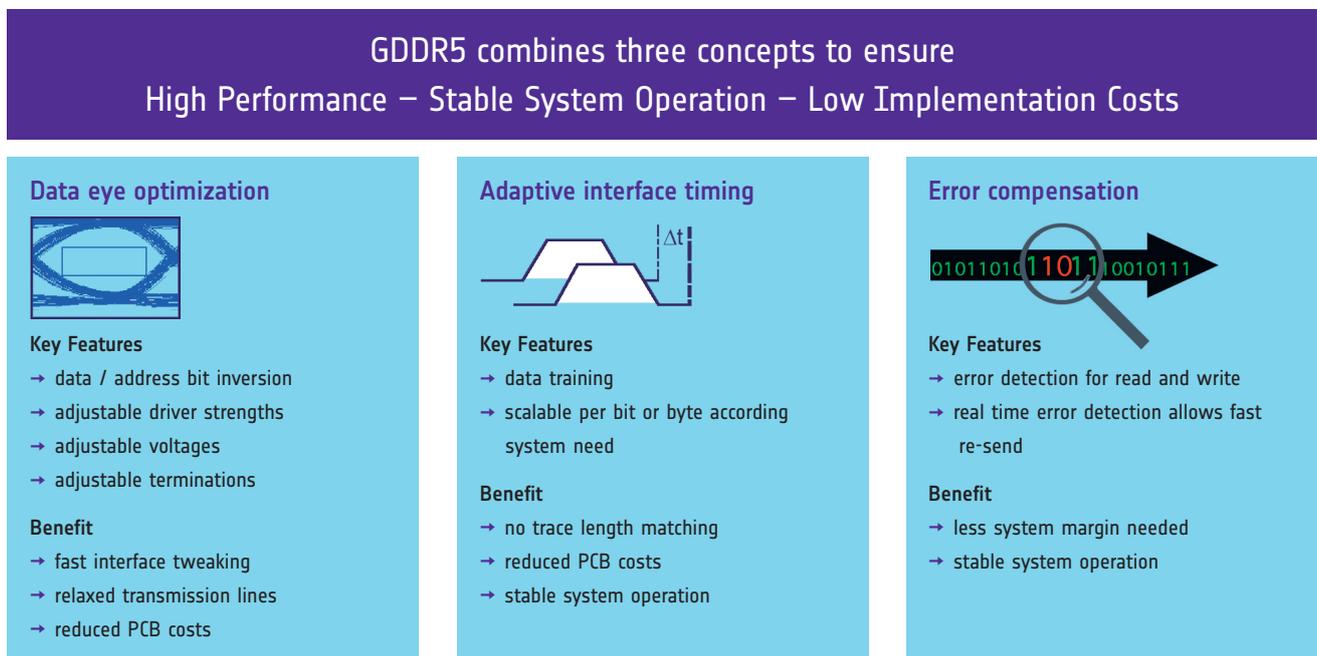


Figure 1

2. Signaling

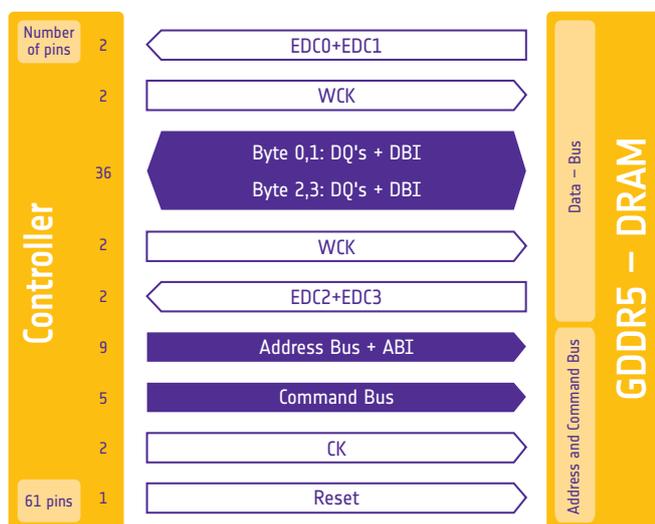


Figure 2

Qimonda GDDR5 SGRAMs extend the proven and reliable single-ended signalling concept of previous graphics standards to very high frequencies. This allows system integrators to leverage all their experience in system design and bring-up when migrating to GDDR5. GDDR5 continues the proven termination concept of high level termination, whereas voltage levels of V_{dd} and V_{ddq} are as low as 1.5 V, with an option for further voltage reductions. The GDDR5 interface is optimized for systems with 32 or 64 bit channels. Clocks, commands and addresses may be shared between two devices, while DQs are routed point-to-point to ensure the high data rates.

2.1 Clocking

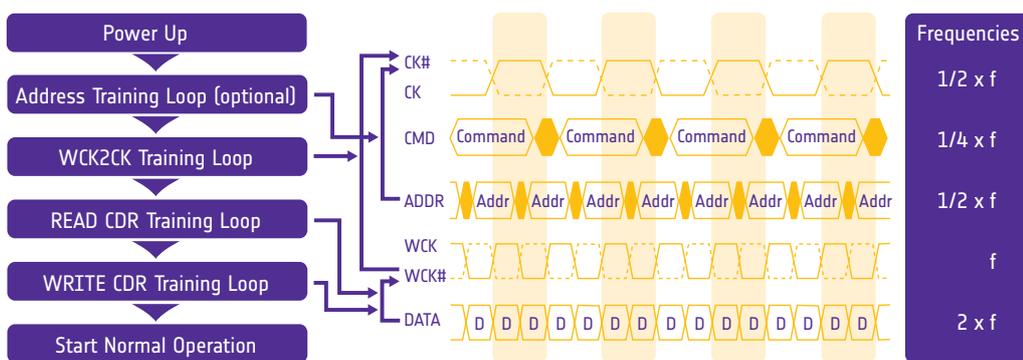
GDDR5 operates with two different clock types. A differential command clock (CK) to where address and command inputs are referenced, and a forwarded differential write clock (WCK) where read and write data are referenced to. Being more precise, the GDDR5 SGRAM uses two write clocks, each of them assigned to two bytes. The WCK runs at twice the CK frequency. Taking a GDDR5 with 5 Gbps data rate per pin as an example, the CK clock runs with 1.25 GHz and WCK with 2.5 GHz. The CK and WCK clocks will be aligned during the initialization and training sequence. This alignment allows read and write access with minimum latency.

The concept of using separated command (CK) and write clocks (WCK) was chosen for an optimal data transmission with minimal noise and jitter at very high data rates like 5 Gbps and higher.

2.2 Command and addresses

Addresses and commands are referenced to CK. Commands are latched single data rate with the rising CK edge. Addresses are latched double data rate with the rising CK and CK# edges. This allows the address information to be transmitted within the same clock cycle as the command. For improved timing margins and higher board routing flexibility users may optionally train the address inputs during power-up. Multiplexed addressing significantly reduces the pin count without impacting the performance of the GDDR5 device.

The GDDR5 SGRAM has five dedicated command pins and eight address pins for 512Mb and 1Gb densities.



Training sequence and signal alignment

Figure 3

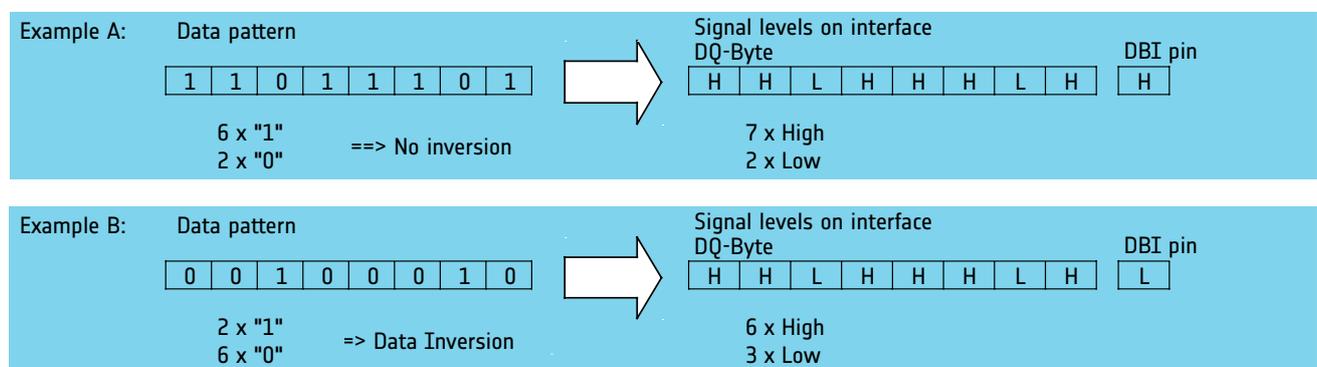
2.3 Data transmission

GDDR5 SGRAMs are organized as x32 interface having 32 single ended data lines. The data is transmitted double data rate relative to the differential write clock (WCK). Since there are two write clocks, always 2 bytes of data are aligned to one of the WCKs. The alignment is done during the initial training period and can be repeated during operation e.g. to track temperature and voltage changes. The data training feature of GDDR5 allows hidden re-training e.g. during refresh operation without impacting the memory bandwidth.

2.4 Data- and address bus inversion

Data inversion is a technique to reduce the number of zeros that are transmitted. The data is inverted if more than 50% of the data bits within a byte are zeros. The inversion is indicated via an additional DBI# pin. There is one DBI pin for each data byte. Since the GDDR5 transmission lines have high level termination, reducing the amount of signal lines driving a low level (= zero) results in reduced power dissipation in the termination resistors and output drivers. Additionally, data inversion improves the signal quality by reducing the supply noise induced jitter on the data lines.

GDDR5 provides address bus inversion (ABI) in the same fashion as DBI which reduces power consumption and noise on the address lines as well.



Data bit inversion - operating mechanism

Figure 4

2.5 Address, clock and data training

Signal alignment is one of the big challenges of high speed applications. The GDDR5 signal training concept offers an easy to implement and reliable solution to system builders. The GDDR5 concept allows system designers to drastically increase the performance of their systems while maintaining the current level of constraints for board design and layout accuracy or even managing to reduce them. Therefore GDDR5 offers the possibility to build cost efficient applications running at very high data rates, and reduces product development times for systems operating robust and reliable.

The GDDR5 signal training concept offers the possibility to train several signals and clocks relative to each other. Training means phase adjustment of the various signals.

GDDR5 clock and data training is performed in 3 steps.

Step 1 – Address training: The address training aligns the address bus to the CK clock. This training step is performed once during memory initialization. The training is performed in a way that the memory returns the latched address to the controller. By sweeping through various phase relations between clock and address, the controller can find the optimal phase setting. Address training is optional.

Training of the command signals is not necessary due to the lower frequency of the command bus.

Step 2 – Write clock (WCK) to clock (CK) alignment. The WCK-to-CK training adjusts the WCK clock to the CK clock at the DRAM's internal phase detector. During the training procedure, the GDDR5 DRAM indicates to the controller the relative phase alignment between CK and WCK as "early" or "late". This allows the controller to easily find the optimum settings for each

DRAM. The WCK-to-CK training is also performed once during initialization of the memory and after frequency changes.

Step 3 – Data training: The data training finally aligns the data with the respective WCK clock. Depending on system cost and performance requirements, the controller can perform this training “per bit” for each single data line separately, “per byte” for a group of 8 data lines together or per “double byte” for 16 data lines together.

Data training is thought to be repeated during system operation, especially in high performance systems. GDDR5 offers the possibility to do a “hidden” re-training of the data lines without impacting the memory bandwidth.

As a result of the training procedure, all signals and clocks are aligned in the memory as shown in figure (3)

3. I/O cells

GDDR5 SGRAMs offer several features that let the controller perfectly adapt the device's input and output characteristics to the actual system impedance and thus improve the data eye for a reliable data transmission.

- Auto calibration for process, voltage and temperature drift compensation
- Software controlled adjustable drive strengths
- Software controlled adjustable data, address and command termination impedances
- Software controlled adjustable data input reference voltage

The output drive strength and on-die termination are auto-calibrated to cancel out effects of process variations and voltage or temperature drift during device operation. The auto-calibrated values may further be offset by the controller to adjust the device to the individual channel and controller characteristics.

All values can get adjusted on a software level typically via the BIOS of the graphics card. This makes GDDR5 system bring-up extremely powerful and flexible and finally results in quick learning cycles and an improved time to market.

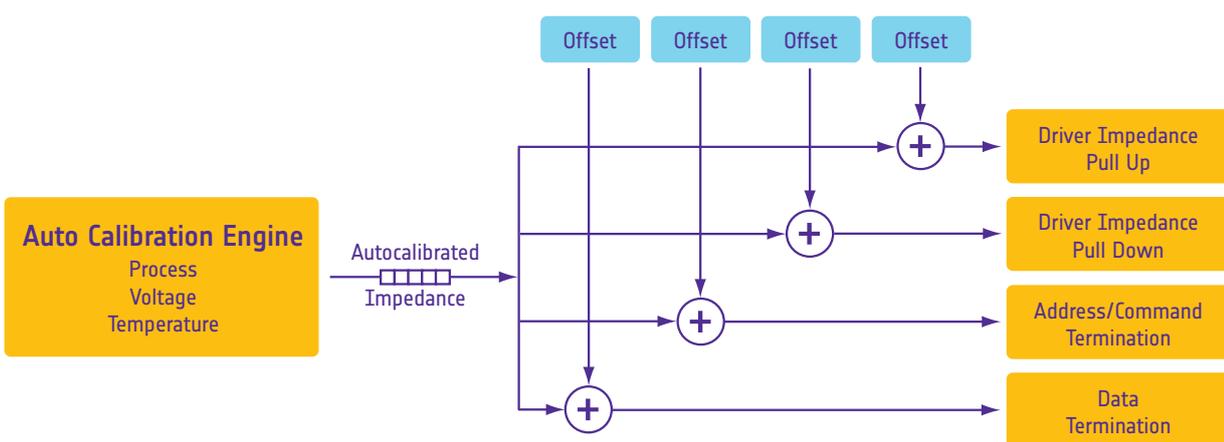


Figure 5

4. Read / Write transaction

The GDDR5 continues the proven concept of for read and write transactions enriching it towards higher data rates, efficient and reliable data transmission. The GDDR5 burst length of only 8 bit allows an access granularity of 256 bit. For larger data chunks continuous gapless read or write bursts are supported.

4.1 Data mask

A GDDR5 data burst transmits always 256 data bits. Sometimes part of this data burst should not overwrite the content of the DRAM. With data mask, GDDR5 offers a very efficient feature to mask parts of the data burst. The masked bytes will not be written to the DRAM, only the unmasked bytes are written to the DRAM and are overwriting the previous content. Without data mask the controller would be forced to execute a read-modify-write cycle which takes time and reduces bus efficiency. GDDR5 uses a new concept of masking data which reduces the number of high speed signal pins. GDDR5 uses the address bus to transmit the masking information. For most efficient usage of the data mask, the GDDR5 offers the possibility to mask the write data per each byte or per double bytes.

4.2 Error detection mechanism

A new feature of GDDR5 is the capability for detection of transmission errors occurring on the high speed signal lines. As graphics systems store increasingly more code in the DRAM, error detection becomes essential, as random bit fails associated with any high speed data transmission would lead to unacceptable system failures.

In GDDR5 the transmitted data is secured via CRC (cycle redundancy check) using an algorithm that is well established within high quality communication environments like ATM networks. The algorithm detects all single and double errors with 100% probability. The CRC scheme is implemented on a per byte basis, securing all DQ and DBI# lines. A eight bit checksum is calculated by the DRAM on each data burst (8 DQs + 1 DBI# x burst of 8 = 72 bit) and returned to the controller via dedicated EDC pins. When the DRAM controller detects an error, the command that caused the error can be repeated. Error detection can be used to trigger re-training of the data transmission line which allows the system to dynamically adapt to changing conditions like e.g. temperature and voltage drift.

5. GDDR5 memory core and addressing scheme

5.1 Core addressing scheme

The GDDR5 core architecture is optimized to sustain the high interface bandwidth of GDDR5 and in parallel ensure low latency random access throughout the core. This is achieved by having up to 16 banks and a prefetch of 8 which makes GDDR5 a superior high performance standard. The basic core parameters are:

- Number of banks: 8 for 512 Mbit GDDR5 and 16 for 1 or 2 Gbit
- prefetch of 8
- fixed page size for all densities of 2 kB

	512M 16M x 32 GDDR5	1G 32M x 32 GDDR5	2G 64M x 32 GDDR5
Row Address	A0-A11	A0-A11	A0-A12
Column Address	A0-A5	A0-A5	A0-A5
Page Size	2kByte	2kByte	2kByte
Bank Address	BA0-BA2	BA0-BA3	BA0-BA3
Bank Groups	4	4	4

Figure 6

5.2 Bank grouping

Additional to the performance advantage of having 16 real physical banks, the GDDR5 adds the bank group concept on top. Bank groups allows the GDDR5 to stay with a prefetch of 8 while going to interface data rates higher than 4 Gbps, the approximate speed limitation of the individual memory bank. Low prefetch is essential for high performance graphics applications as they need fast random access to small chunks of data. Low prefetch increases the efficiency of data bus usage and minimizes the amount of useless data transmitted. Above a data rate of 4 Gbps the GDDR5 can be operated using the bank group concept. Therefore the GDDR5 is divided into 4 different bank groups where interleaving access to banks from different bank groups is possible, allowing gapless read and write access beyond 4Gbps data rate.

6. Clamshell mode (x16 mode)

Graphics system designers expect GDDR5 standard to offer high flexibility in terms of frame buffer and bandwidth variation. GDDR5 supports this need for flexibility in an outstanding way with its clamshell mode. The clamshell mode allows 32 controller-I/Os to be shared between two GDDR5 components. In clamshell mode each GDDR5 DRAM's interface is reduced to 16 I/Os. 32 controller I/Os can, therefore, be populated with two GDDR5 DRAMs, while DQ's are single loaded and the address and command bus is shared between the two components. Operation in clam shell mode has no impact on system band width.

Example: System configurations with 512M GDDR5 device using a controller with 256 bit interface:

- A) 8pcs of 512M GDDR5 in standard mode → Frame buffer: 512 MB
- B) 16pcs of 512M GDDR5 in clamshell mode → Frame buffer: 1 GB

Every GDDR5 component supports the clam shell mode. In this way, multiple frame-buffer variants can be built up using only one component type which drastically reduces the number of different inventory positions and increases flexibility in a very dynamic market environment.

Standard – x32 mode

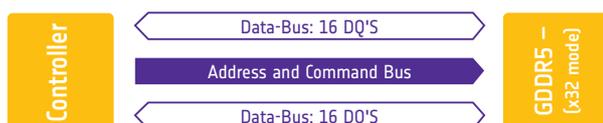


Figure 7

Clamshell – x16 mode

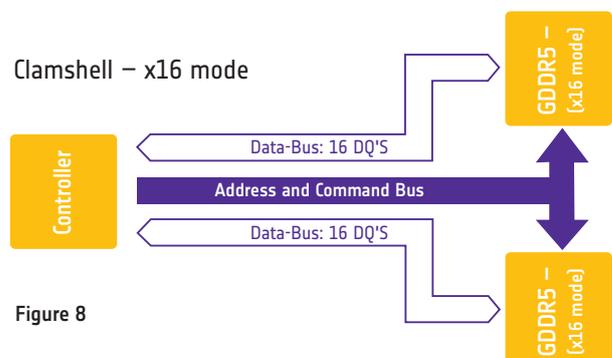


Figure 8

7. GDDR5 Power management

The GDDR5 is designed in a way to only consume power when really needed. Several features and methods are implemented in a way to allow a demand driven power management.

- Extreme wide clock frequency range and data rates
- Multi level, demand driven termination enabling
- Low power modes for DRAM core
- Low supply voltage of 1.5V
- Data and address bit inversion
- Power-down and self refresh modes

Scalable clock frequency and data rate

GDDR5 allows the system to dynamically scale the memory I/O data rate according to the workload. The I/O data rate of GDDR5 can be gaplessly varied from 5 Gbps down to 200 Mbps (50 MHz clock frequency). Towards lower frequencies the PLL is turned off for additional power saving.

Low power strobe mode

For lower frequencies the GDDR5 interface can be set to a low power strobe mode. In this case, the GPU can turn off the clock data recovery reducing the power consumption. For data alignment in the strobe mode, the DRAM sends out a strobe signal via the EDC pins together with the data.

Low power mode for DRAM core

Additionally to the I/O frequency scaling, the GDDR5 DRAM core offers a low power mode for operation at lower frequencies. The low power mode is initiated by a low power bit set from the controller.

Multi level termination

Signal termination is necessary to match impedances of transmission lines and improve signal quality. But termination consumes also power. At lower data rates, the signalling gains more and more margin which allows to operate the system with partially matched termination impedance. GDDR5 allows to double the termination impedance for slower data rates or even turn it off. Termination can be adjusted separately for data bus, command & address bus and WCK clock to take maximum advantage of the power saving potential.

Data and address bit inversion

The data and address bit inversion feature does reduce the power consumed by the termination resistors and output drivers. The GDDR5 high level termination scheme only consumes power if the transmission line is driven to Low. As data and address bit inversion effectively reduces the amount of LOW signals on the data and address bus, it reduces the average power consumption as well.

8. GDDR5 target applications

Memory bandwidth is a key factor in the rapidly increasing 3D rendering performance of PC graphics systems and game consoles. Over the last years, the memory bandwidth of graphics DRAMs grew by nearly 30 % per year. Bandwidth of graphics memory systems exceeds those of PC main memory by far.

3D Graphics – always hungry for band width

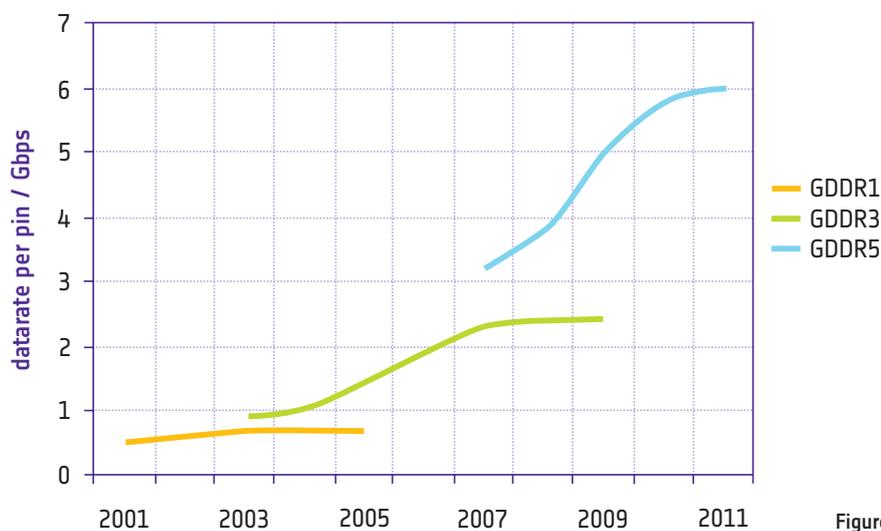


Figure 9

GDDR5 Target applications:

- PC graphics solutions
- Game Consoles
- High End Consumer applications (e.g. DTV, Set top boxes)
- Networking

PC graphics solutions

PC graphics is one of the most dynamic areas. Hungry for highest band width and very short innovation cycles requires always leading edge memory solutions. GDDR5 is ideally suited to keep the pace of PC graphics solutions for several years.

Game Consoles

Future Game consoles will be more than just a gaming device. They are most likely a complete entertainment center fully connected to other devices of the electronic home. But like today they will provide leading edge 3D graphics performance and will be one of the most powerful processing devices in the household. GDDR5 is designed to play a prominent role in future game console architectures.

High end consumer applications

Digital broadcasting opens a new space of content provision and interactivity.

Advanced video compression, picture scaling and rendering, content transmission, interactive user interfaces are only a few techniques entering TV and set top box applications.

These real time processing requirements leads to increasing processing power, memory amount and memory performance in high end consumer applications. GDDR5 provides the ideal feature set to satisfy performance requirements and ensure system robustness.

Networking applications

High bandwidth, fast access times and reliable data transmission are key requirements for memories in networking applications. GDDR5 therfor provides an attractive feature with outstanding I/O performance and as first graphics memory coming with an error detection mechanism.

GDDR memory drives leading edge graphics

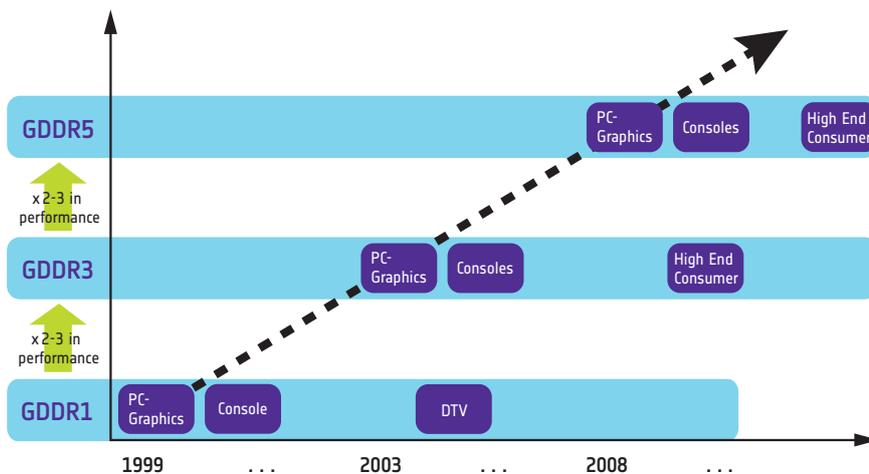


Figure 10

Glossary

ABI:	Address Bit Inversion	Gbps:	Giga bit per second
ATM:	Asynchronous Transfer Mode	I/O:	Input / Output
CK / CK#:	Clock / inverted clock	NOP:	No Operation
CRC:	Cyclic Redundancy Check	PLL:	Phase Locked Loop
DBI:	Data Bit Inversion	Vdd:	DRAM core voltage
DM:	Data Mask	Vddq:	I/O voltage
DQ:	Data I/O	WCK / WCK#:	Write clock / inverted write clock
DTV:	Digital Television		