



WP377 (v1.1) February 17, 2012

Xilinx 7 Series FPGAs Embedded Memory Advantages

By: Nick Mehta

The architectures of the Xilinx® 7 series FPGAs feature flexible internal memory resources that can be configured in a variety of different sizes. This white paper details the available features, illustrating the wide array of memory sizes available and shows the trade-off of using different resources to perform memory functions of different sizes.

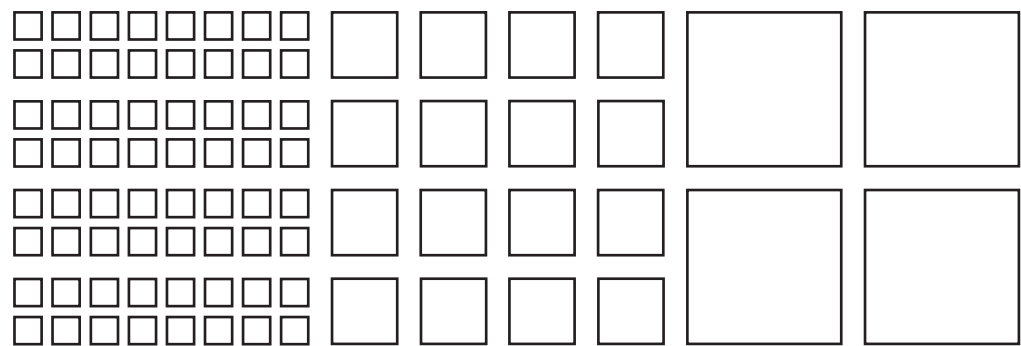
Usage of Internal Memory Resources in Designs

Almost every design built within an FPGA requires the use of internal memory resources of some size for storage of coefficients, buffering of data, and a variety of other uses. Typical systems need a combination of small, medium, and large memory arrays to meet all of their requirements, with the overall power consumption of the memories, and therefore the FPGA, being a primary concern.

When laying out an FPGA, it is important to create a device that meets the majority of customer needs. If an FPGA is built with small, medium, and large memory resources that perfectly suit one application, the solution will be optimal for some customers while others wanting to use the same part will likely need to make considerable trade-offs.

Users trying to get the best value out of their FPGAs can be concerned about wasted bits in larger RAMs. Building finer RAM granularity, though, requires additional FPGA interconnect lines, at a cost. This white paper explains the trade-offs: why the cost of finer RAM granularity is often higher than the cost of otherwise wasted bits.

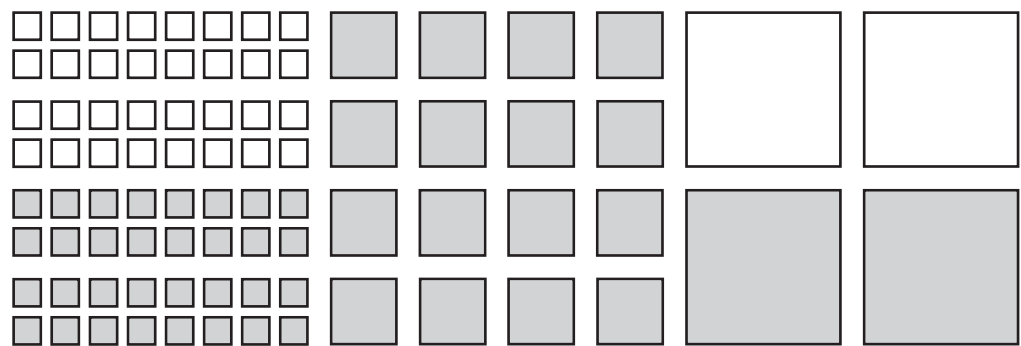
Figure 1 shows a theoretical distribution of small, medium, and large memory blocks within an FPGA (not drawn to any particular scale).



WP377_01_090310

Figure 1: Representation of Memory Arrays of Different Size

A design that requires exactly this combination of blocks is able to make perfect use of the available resources (see Figure 2).



WP377_02_090310

Figure 2: Representation of Memories Mapping Perfectly into Available Resources

However, imagine a scenario in which a user requires just four more of the medium sized resources. One possibility is to build the medium-sized memory array out of lots of small components, consuming lots of resources and incurring the complexity of connecting them together. Another alternative is to use one large block as a medium block, rendering all the remaining resources in the large block unusable while keeping them powered on, therefore consuming power (see [Figure 3](#)).

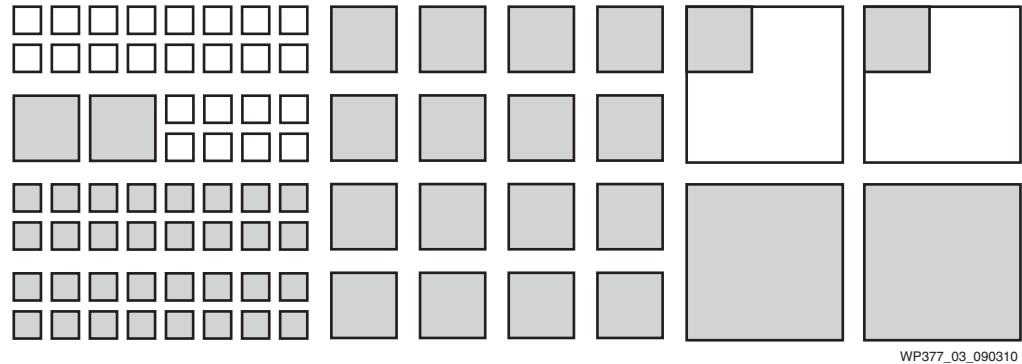


Figure 3: Representation of Memories Having to Use Non-Ideal Resources

The challenge for an FPGA manufacturer is to build devices with the most flexible combination of memory resources, giving all users the ability to fit their desired memory array size into the device while achieving the desired performance—without wasting vast quantities of resources and power.

FPGA Memory Resources

Xilinx FPGAs all use a variety of memory resources to give the best-in-class combination of flexibility and low cost—or cost per bit. All 7 series FPGAs, which include Artix™-7, Kintex™-7 and Virtex®-7 families, use the same memory blocks, enabling seamless migration from one 7 series FPGA family to another.

It is evident that building memory resources to fulfill the requirements of every user is a difficult challenge. The solution implemented in Xilinx 7 series FPGAs is to create blocks called block RAM (see [Figure 4](#)) that can either be combined together to make larger arrays or divided to make smaller arrays. The ability to combine block RAMs with 6-input look-up tables (LUTs) in the FPGA logic as small memory arrays provides the user the most flexible resources from which to create their various-sized memory arrays.

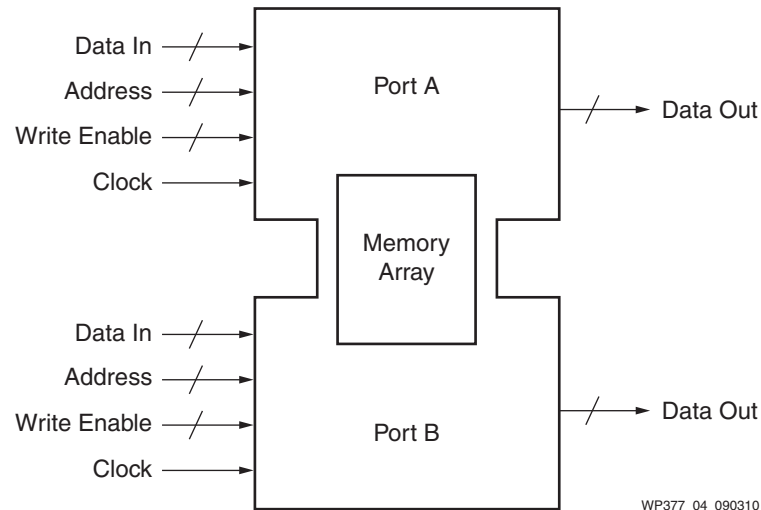


Figure 4: A Block RAM

Block RAM

Every 7 series FPGA has between 135 and 1,880 dual-port block RAMs, each capable of storing 36 Kb, 32 Kb of which is allocated to data storage and, in some memory configurations, an additional 4 Kb allocated to parity bits. Each block RAM has two completely independent ports that share nothing but the stored data.

Each port can be configured as:

- 32K x 1
- 16K x 2
- 8K x 4
- 4K x 9 (or 8)
- 2K x 18 (or 16)
- 1K x 36 (or 32)
- 512 x 72 (or 64)

Each block RAM can be divided into two completely independent 18 Kb block RAMs that can each be configured to any aspect ratio from 16K x 1 to 512 x 36. When a 36K block RAM is split into two independent block RAMs, each of the two independent block RAMs behaves exactly like a one 36 Kb block RAM, just half the size.

Conversely, if the user requires larger memory arrays, two adjacent 36 Kb block RAMs can be configured like one cascaded 64K x 1 dual-port RAM without any additional logic.

The block RAM components in Xilinx 7 series FPGAs can be configured in single port, simple dual-port, or true dual-port modes. Additionally, data can be read from the block RAM in one of three ways: READ_FIRST, WRITE_FIRST, or NO_CHANGE mode. For more information on the permitted configurations and write modes, refer to [UG473](#), *7 Series FPGAs Memory Resources User Guide*.

Splitting Block RAM

If the user only requires single port memories, rather than implementing full true dual-port functionality, it is possible to divide the block RAMs into smaller memory arrays. When a block RAM is in true dual-port mode (the default mode), Port A and Port B can implement separate, independent delay line memories, single port memories, or ROMs by connecting the most significant bit of the ADDRA address bus to V_{CC} and the most significant bit of the ADDR B address bus to ground, thus creating two single port block RAMs. For example, a RAMB36E1, the 36K block RAM primitive, can be divided into two 18K single port memories, and the RAMB18E1 can be divided into two 9K single port memories. Using this method, it is possible to create four delay line memories per block RAM array in 7 series FPGAs. To implement delay lines in the block RAM in this fashion, READ_FIRST or Read-Before-Write mode must be used. (In Read-Before-Write mode, data previously stored at the write address appears at the outputs while input data is being stored in the memory.) If implementing single port memories, there is no such limitation on the permitted modes; any supported mode (READ_FIRST, WRITE_FIRST, or NO_CHANGE) can be used, and the different memories within one block RAM can have different port widths. The resulting memory schemes perform one operation per clock cycle for each port, and therefore, four operations per block RAM per clock cycle.

Synchronous Operation

Each memory access, whether a read or a write, is controlled by the clock. All inputs, data, address, clock enables, and write enables are registered. Nothing happens without a clock. The input address is always clocked, retaining data until the next operation. An optional output data pipeline register allows higher clock rates at the cost of an extra cycle of latency. During a write operation, the data output can reflect either the previously stored data, the newly written data, or remain unchanged.

Byte-Wide Write Enable

The byte-wide write enable feature of block RAM gives the capability to write eight bit (one byte) portions of incoming data. There are up to four independent byte-wide write enable inputs to the true dual-port RAM. Each byte-wide write enable is associated with one byte of input data and one parity bit. This feature is useful when using block RAM to interface with a microprocessor. Byte-wide write enable is further described in [UG473](#), *7 Series FPGAs Memory Resources User Guide*.

Error Detection and Correction

Each 64-bit-wide block RAM in Xilinx 7 series FPGAs can generate, store, and utilize eight additional Hamming-code bits, and perform single-bit error correction and double-bit error detection (ECC) during the read process. The ECC logic can also be used when writing to or reading from external 64/72-bit-wide memories. This works in simple dual-port mode and does not support read-during-write.

FIFO Controller

The built-in FIFO controller in Xilinx 7 series FPGAs for single-clock (synchronous) or dual-clock (asynchronous or multirate) operation increments the internal addresses and provides four handshaking flags: full, empty, almost full, and almost empty. The almost full and almost empty flags are freely programmable. Similar to the block RAM, the FIFO width and depth are programmable, but the write and read ports always have identical width. First-word fall-through mode presents the first

written word on the data output even before the first read operation. After the first word has been read, there is no difference between this mode and the standard mode.

Distributed RAM

The logic of the Xilinx 7 series FPGAs consists of, among other elements, 6-input LUTs. The LUTs are arranged in groups of four and combine to make a slice. The 7 series FPGAs have two types of slices, SLICEM and SLICEL. LUTs in a SLICEM, which compose 25–50% of the total number of slices in 7 series FPGAs, can be implemented as a synchronous RAM resource called a distributed RAM element. Each 6-input LUT can be configured as a 64 x 1-bit RAM or two 32 x 1-bit RAMs. The 6-input LUTs within a SLICEM can be cascaded to form larger elements up to 64 x 3-bit in simple dual-port configuration or 256 x 1-bit in single port configuration. See [Figure 5](#).

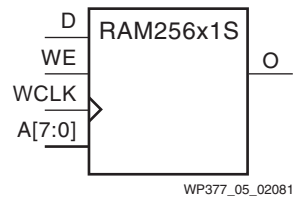


Figure 5: 256-Deep by 1-Wide Random Access Memory

Distributed RAM modules are synchronous (write) resources. A synchronous read can be implemented with a storage element or a flip-flop in the same slice. By placing this flip-flop, the distributed RAM performance is improved by decreasing the delay into the clock-to-out value of the flip-flop. However, additional clock latency is added. The distributed elements share the same clock input. For a write operation, the Write Enable (WE) input, driven by either the CE or WE pin of a SLICEM, must be set High.

A Typical User Design Utilization

[Table 1](#) shows the resource mapping in a typical high memory utilization design targeting a Kintex-7 XC7K410T FPGA. This data is based on a real-world example from a user concerned about wasted bits.

Table 1: Resource Mapping in a Typical High Memory Utilization Design

Requirement	Implementation
16 memories at (32K x 32)	512 block RAMs (36K configuration)
300 memories at (512 x 18)	150 block RAMs (18K configuration)
100 memories at (64 x 4)	100 SLICEMs configured as distributed RAM

The multiple different resources take advantage of the ability to configure the block RAM resources in different sizes.

The Cost of Infinite Granularity

Having established that the Xilinx FPGA architectures offer a great variety of different memory depth/width granularities, it is important to understand the trade-offs involved in adding finer granularity to the architecture. If, for example, the 36K block RAM in the 7 series FPGAs is divided into not only two 18K blocks but further into four 9K true dual-port blocks, there is a penalty to pay. Doubling the number of unique memories within a single block RAM means that the maximum number of signals to be routed into and out of each block RAM also needs to double, which, in turn, necessitates double the number of interconnect resources (or tiles) to accommodate routing of the ~400 signals.

For example, an 8K block RAM can be configured as 16 bits wide and 512 bits deep, requiring 16 data lines and 9 address lines for a total of 25 input signals. Then the 8K block RAM are divided into two 4K block RAMs. Each of these block RAMs can be configured as 16 bits wide and 256 bits deep. This configuration requires 16 data lines and 8 address lines per block, for a total of 24 input signals per 4K block RAM or 48 signals total. The 48 signals are roughly twice the 25 input signals needed for the 8K block RAM.

The effect of doubling the number of interconnect tiles associated with each block RAM increases the silicon area by 25% — a penalty that is applied to all blocks regardless of the configuration in which they are used. Therefore, the ability of every block RAM to be configured as four 9K blocks means that each block increases from four area units to five area units. On the positive side, it allows smaller memories to take advantage of the smaller block size (see [Table 2](#)).

Referring to the same design example, the impact amounts to the smaller 512 x 18 memories consuming fewer resources if they are able to be packed efficiently into adjacent 9K blocks. However, the larger blocks all become 25% larger, adding a significant area penalty.

Table 2: Area Impact on Typical High Memory Utilization Design

Resources and Silicon Area Used		
	If Supporting 36K/18K	If Supporting 36K/18K/9K
16X (32K x 32)	512 blocks of 4 area units	512 blocks of 5 area units
300X (512 x 8)	150 blocks of 4 area units	75 blocks of 5 area units
	Difference	+512 area units -(3 x 75) area units
	Result	+287 area units

Thus, if the block RAMs in the Kintex-7 FPGA can be divided into four independent memories and the user design can take full advantage of these smaller granularity blocks, this typical design still consumes a far greater area than the current 36K/18K configuration.

It is, however, unlikely that a design can take full advantage of having four independent 9K true dual-port memories by packing them into the same 36K block RAM due to the routing congestion of trying to connect so many signals to a single location. The area penalty in [Table 2](#) is, therefore, a best case scenario; in reality, not all 36K block RAMs would be populated by four 9K memories.

Impact on Device Resources

The impact of the block RAM resources increasing in size can have one of three results in context of overall device resources:

- Keep the device size the same and lose block RAM bits
- Keep the number of block RAM bits the same and reduce the quantity of other resources, i.e., CLBs
- Leave all resources as they are and allow the device to become larger

There are obvious penalties associated with each of these choices:

- Losing block RAM bits is not ideal because most users are utilizing the majority of the available memory. If the block RAM are all 25% larger, the number of blocks in a Kintex-7 XC7K410T FPGA decreases from 795 to 636, equating to a reduction from 28,620 Kb to 22,896 Kb.
- A 25% increase in block RAM size means the loss of one column of CLBs per block RAM column. The Kintex-7 XC7K410T FPGA has 12 block RAM columns, so in this scenario, it loses 12 columns of CLBs from the device, equating to a loss of 4,200 CLBs, which translates into approximately 54,000 logic cells. This would have the impact of reducing the Kintex-7 XC7K410T FPGA from 406,720 logic cells to 352,720 logic cells, resulting in a much less capable device at the same cost.
- Keeping all the resources the same and allowing the silicon area to increase has several drawbacks. The obvious first issue is that a larger device means a more expensive device, but the physically larger silicon also has a significant impact on power consumption.

Power Optimization of Block RAM Arrays

The Importance of Power Consumption

Power consumption is critical in the majority of modern applications, and reducing the power consumed within a system is a challenge for every designer. There are myriad approaches to reducing power, but many come with a significant performance penalty. With the 7 series FPGAs, Xilinx has taken an innovative approach to reducing power associated with the block RAM memory arrays. The two primary methods by which power is reduced are:

- Identifying and acting upon the areas in which power is consumed unnecessarily
- Providing the ability to trade-off lower power consumption for a slight reduction in maximum performance

Power of Unused Block RAM Sites

All block RAM components consume power when they are powered on, regardless of whether or not they are used in a design. Present in Xilinx 7 series FPGAs is a unique feature that enables the software to automatically recognize when block RAM sites are unused. When unused sites have been identified, they are automatically disabled and put into a zero-power state, thereby significantly reducing the overall power consumption of the FPGA.

Inferring RAMs with XST

XST infers block RAM if the read address is registered inside the module. Conversely, it infers distributed RAM if the RAM network uses an asynchronous read operation. The `ram_style` attribute can be used to specify whether block RAM or distributed RAM is used. When the RAM network is larger than is supported by one component, multiple block RAM can be used. The default strategy is to optimize for performance. In this strategy, small RAMs are implemented using distributed RAM. RAM networks can also be optimized for power and area; these are described in the [Building RAMs with CORE Generator Software](#) section of this white paper. Anything smaller than 128-bits deep or 512-bits depth x width is implemented in distributed RAM unless the user specifies otherwise with the `ram_style` attribute. To ensure the most efficient use of both block RAM and distributed RAM resources, XST first implements the largest RAMs in block RAM, then places the smaller RAMs in block RAM if there are still block RAM resources available. XST can also pack small single-port RAMs into a single block RAM.

Techniques to reduce block RAM power dissipation can be enabled in XST. They are part of a larger set of optimizations controlled with the Power Reduction (POWER) synthesis option and can be specifically enabled through the RAM Style (RAM_STYLE) constraint. RAM power optimization techniques in XST are primarily aimed at reducing the number of simultaneously active block RAMs on the device. They are applicable only to inferred memories that require a decomposition on several block RAM primitives, and take advantage of the enable capability of block RAM resources. Additional enable logic is created to ensure that only one block RAM primitive used to implement an inferred memory is simultaneously enabled. Activating power reduction techniques for an inferred memory that fits in single block RAM primitive has no effect. When enabled, power reduction is sought in conjunction with both Area and Speed optimization goals. Two optimization trade-offs are available through the RAM Style (RAM_STYLE) constraint:

- Mode `block_power1` achieves some level of power reduction while minimally impacting circuit performance. In this mode, the default, performance-oriented, block RAM decomposition algorithm is preserved. XST simply adds block RAM enable logic. Depending on the memory characteristics, power can be impacted in only a limited way.
- Mode `block_power2` provides more significant power reduction but can have a slight impact on performance. Additional slice logic can also be induced. This mode uses a different block RAM decomposition method. This method first aims to reduce the number of block RAM primitives required to implement an inferred memory. The number of active block RAMs is then minimized by inserting block RAM enable logic. Multiplexing logic is also created to read the data from active block RAMs. If the primary concern for the design is power reduction, and speed and area optimization are secondary, Xilinx recommends using the `block_power2` mode.

ROMs can be inferred by the use of large case statements. XST can also implement finite state machines (FSM) and logic in block RAM to maximize the available logic resources.

Building RAMs with CORE Generator Software

The CORE Generator™ tool has three algorithms by which it can optimize block RAM networks. The Minimum Area scheme uses the fewest possible resources (block RAMs) but also minimizes output multiplexing to maintain the highest performance at the smallest area (see Figure 6). The Low Power scheme can use more resources but ensures the fewest blocks are enabled during each read and write operation. This can result in a small quantity of additional decoding logic on the enable signal, but it is a small real estate penalty compared to the power saving made. When operating in this mode, the CORE Generator tool performs the same function as XST in the RAM_STYLE = block_power2 mode. Note that if the block RAM network is very large, splitting using the Low Power scheme can result in many additional signals needing to be routed, which can deplete the routing resources for other portions of the design.

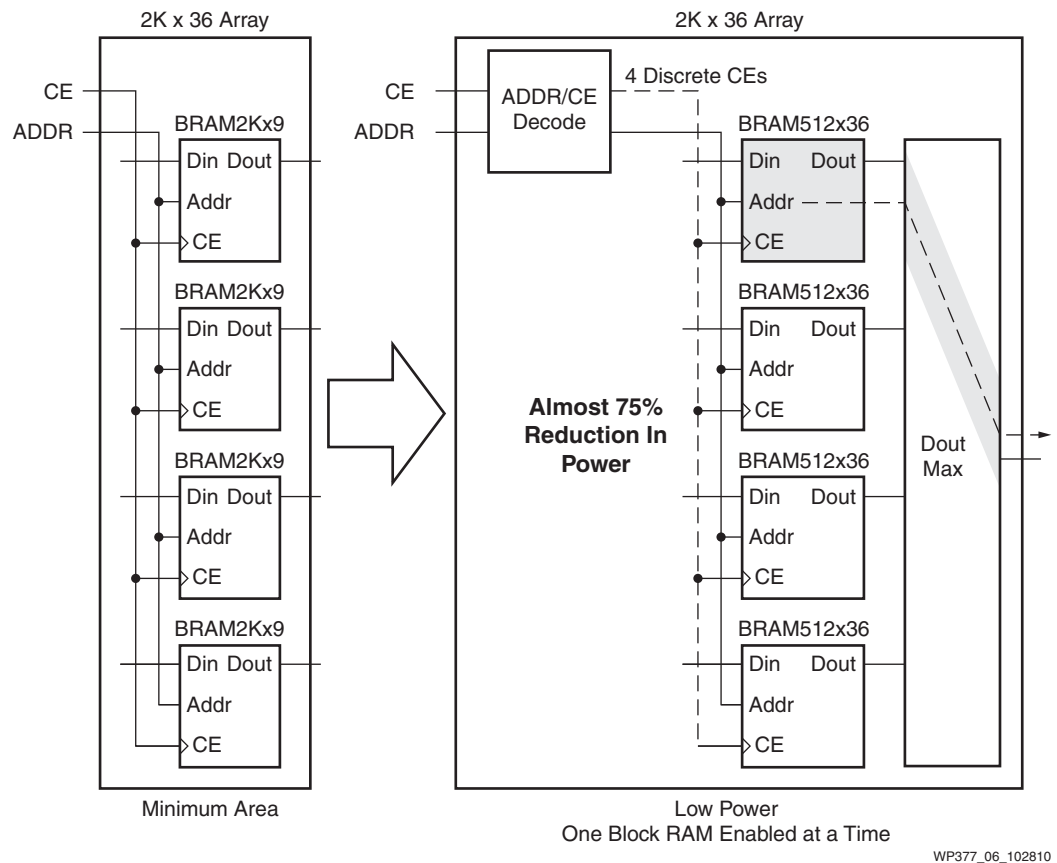


Figure 6: CORE Generator Software Method of Lowering Total Block RAM Power Consumption

The third optimization scheme is the Fixed Primitive, in which the user can choose a specific primitive, e.g., 4K x 4, from which to build their RAM network. It is up to the user to decide which primitive is best for their application.

The CORE Generator tool also provides an option to register the outputs of the RAM network to improve performance. If multiple block RAMs are used in the network, the user can choose whether to register at the output of each block RAM primitive or at the output of the core.

Conclusion

Combining block RAMs that can be configured in a variety of data width/depth combinations with distributed RAM to support smaller memory arrays offers the most flexible way of building memories of different sizes at the lowest cost per bit of memory. Adding additional functionality to a component, such as adding ports to a block RAM, can initially appear to be the best solution if that exact functionality is required. However, there are area penalties to both the component and its corresponding interconnect for adding new features. These area penalties have a significant impact on device resources. Many years of experience developing and building FPGA embedded memories have led to efficient solutions across a wide range of applications.

Additional Resources

The following resources provide additional information useful to this white paper:

- [UG474](#), *7 Series FPGAs Configurable Logic Block User Guide*
- [DS512](#), *Block Memory Generator v5.2*
- [DS322](#), *Distributed Memory Generator v6.1*
- [UG687](#), *XST User Guide for Virtex-6 and Spartan-6 Devices*

Revision History

The following table shows the revision history for this document:

Date	Version	Description of Revisions
03/01/11	1.0	Initial Xilinx release.
02/17/12	1.1	Updated Block RAM section.

Notice of Disclaimer

The information disclosed to you hereunder (the “Materials”) is provided solely for the selection and use of Xilinx products. To the maximum extent permitted by applicable law: (1) Materials are made available “AS IS” and with all faults, Xilinx hereby DISCLAIMS ALL WARRANTIES AND CONDITIONS, EXPRESS, IMPLIED, OR STATUTORY, INCLUDING BUT NOT LIMITED TO WARRANTIES OF MERCHANTABILITY, NON-INFRINGEMENT, OR FITNESS FOR ANY PARTICULAR PURPOSE; and (2) Xilinx shall not be liable (whether in contract or tort, including negligence, or under any other theory of liability) for any loss or damage of any kind or nature related to, arising under, or in connection with, the Materials (including your use of the Materials), including for any direct, indirect, special, incidental, or consequential loss or damage (including loss of data, profits, goodwill, or any type of loss or damage suffered as a result of any action brought by a third party) even if such damage or loss was reasonably foreseeable or Xilinx had been advised of the possibility of the same. Xilinx assumes no obligation to correct any errors contained in the Materials or to notify you of updates to the Materials or to product specifications. You may not reproduce, modify, distribute, or publicly display the Materials without prior written consent. Certain products are subject to the terms and conditions of the Limited Warranties which can be viewed at <http://www.xilinx.com/warranty.htm>; IP cores may be subject to warranty and support terms contained in a license issued to you by Xilinx. Xilinx products are not designed or intended to be fail-safe or for use in any application requiring fail-safe performance; you assume sole risk and liability for use of Xilinx products in Critical Applications: <http://www.xilinx.com/warranty.htm#critapps>.