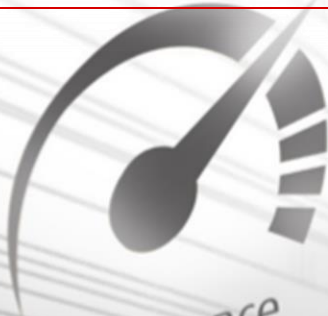


# White Paper

## FUJITSU Server PRIMERGY & PRIMEQUEST RAID Controller Performance 2016

This technical documentation is aimed at the persons responsible for the disk I/O performance of Fujitsu PRIMERGY and PRIMEQUEST servers. The document is intended to help you become acquainted - from a performance viewpoint - with the options and application areas of various RAID controllers for internal disk subsystems. Depending on the requirements for data security and performance as well as planned or existing server configuration, specific recommendations arise for the selection and parameterization of controllers. Controllers of the current generation that are available for PRIMERGY and PRIMEQUEST systems in 2016 are to be considered here.



Version
1.0d
2016-08-29



## Contents

Document history .....	2
Introduction .....	3
RAID controllers for PRIMERGY and PRIMEQUEST: Basics .....	4
Presentation of the RAID controllers .....	4
Controller interfaces and their throughput limits .....	6
Safeguarding the controller cache against power failure .....	9
FastPath .....	10
Performance-relevant cache settings .....	10
Other settings .....	13
Characteristics of the onboard controllers .....	15
Measurement context .....	16
Measurement method .....	16
Measurement environment .....	17
Controller comparison .....	19
RAID 1 (two SATA hard disks) .....	20
RAID 0 and 10 (four SATA hard disks) .....	22
RAID 0, 10 and 5 (eight SAS hard disks) .....	28
RAID 0, 10 and 5 (more than eight SAS-SSDs) .....	41
Influence of the size of the controller cache .....	47
Lower load levels .....	48
Conclusion .....	49
Literature .....	50
Contact .....	50

## Document history

### **Version 1.0 (2016-03-07)**

- Initial version

### **Version 1.0a (2016-03-23)**

- Minor changes

### **Version 1.0b (2016-03-23)**

- Minor changes

### **Version 1.0c (2016-05-20)**

- Minor changes

### **Version 1.0d (2016-08-29)**

- Minor changes

## Introduction

Hard disks are a security factor as well as critical performance components in the server environment. It is thus important to bundle the performance of such components via intelligent organization so that they do not cause a system bottleneck. They should simultaneously compensate for any failure of an individual component. Methods exist for arranging several hard disks in logical drive so that any hard disk failure can be compensated. This is known as a “Redundant Array of Independent Disks” or in short RAID. Special RAID controllers are normally used.

The PRIMERGY and PRIMEQUEST servers are available in a wide range of internal configuration versions with different RAID controller and hard disk configurations. The “Modular RAID” concept that is offered as a standard for all servers in the PRIMERGY and PRIMEQUEST family consists of a modular controller family and standardized management via the Fujitsu RAID Manager software known as “ServerView RAID Manager”. The comprehensive offer of RAID solutions enables the user to select the appropriate controller for a particular application scenario. The performance of the disk subsystem is defined by the controller, the selected hard disks and the features of the RAID level.

Several documents have been created in the PRIMERGY & PRIMEQUEST white paper series which illustrate all aspects of “Modular RAID” regarding performance:

- We recommend - as a comprehensive introduction to disk I/O performance - the White Paper “[Basics of Disk I/O Performance](#)”.
- This document “[RAID Controller Performance 2016](#)” covers all RAID controllers of the current generation, including their performance, that are on offer for PRIMERGY and PRIMEQUEST.
- This predecessor document “[RAID Controller Performance 2013](#)” covers the RAID controllers of the generation of that time and their performance.

When sizing internal disk subsystems for PRIMERGY and PRIMEQUEST servers you can proceed in such a way that a suitable hard disk type is selected and the necessary number of hard disks for the required RAID level is estimated using rules of thumb. Due to the number and technology of the hard disks that are to be connected as well as the required RAID level the RAID controller is self-evident. This may be adequate for years in order to accurately size a disk subsystem.

However, the technology of storage media (for example Solid State Drives, or in short SSDs) or in the internal interfaces of the server has progressed over the years and the new disk subsystem no longer meets the increased requirements. Or, in a productive server configuration the application scenario changes and the achieved disk I/O performance is - despite an adequate number of hard disks - not as desired. In both these cases it can be worthwhile to look at the influence of the RAID controller on performance more closely. Sometimes the right controller, or even simply the correctly configured controller, is prerequisite for the best possible performance.

That outlines the objective of this document. First, there will be an overview of the current internal RAID controllers that are available for the PRIMERGY and PRIMEQUEST systems. The throughput limits of the involved controller interfaces will then be presented under the aspects of performance. After a brief introduction into the measurement context, the different RAID controllers will be compared at various RAID levels and in different application scenarios, which will be substantiated by the measurement results.

In the past the terms “Hard Disk” and also “Hard Disk Drive” (HDD) were used for a hard magnetic-coated, rotating, digital, non-volatile storage medium that could be directly addressed. Technical development has now seen new “hard disk” versions introduced as storage media; they use the same interface to the server and are accordingly handled as hard disks by the server. An SSD, which as an electronic storage medium does not contain any moving parts, can be stated as a typical example, but which nevertheless is also colloquially referred to as a hard disk. Throughout this document the term “hard disk” is used as a generic term, with the names “SSD” and “HDD” being used as a means of differentiation.

This document specifies hard disk capacities on a basis of 10 (1 TB =  $10^{12}$  bytes) while all other capacities, file sizes, block sizes and throughputs are specified on a basis of 2 (1 MB/s =  $2^{20}$  bytes/s).

## RAID controllers for PRIMERGY and PRIMEQUEST: Basics

Firstly, the RAID controllers that are available for PRIMERGY and PRIMEQUEST servers are to be presented with their essential functions in this section. Then the throughput limits, which result for the individual controllers due to their interfaces in the server, will be looked at more closely. We will subsequently deal with the possible settings of the controllers, and ultimately discuss the characteristics of onboard controllers.

### Presentation of the RAID controllers

The following table summarizes the most important data with regard to the functionality of the available RAID controllers.

To simplify the naming of these controllers this white paper will for the most part only use the short name from the column "Alias", thus for example C610.

Controller name	Alias	FF	Cache	Fre- quency	Supported interfaces		Max. # disks	RAID levels	FBU
LSI SW RAID on Intel C220 (Onboard SATA)	C220	I	-	6G	SATA 6G	DMI 2.0 x4	4	JBOD, 0, 1, 10	-
LSI SW RAID on Intel C236 (Onboard SATA)	C236	I	-	6G	SATA 6G	DMI 3.0 x4	4	JBOD, 0, 1, 10	-
LSI SW RAID on Intel C610 (Onboard SATA) <sup>3)</sup>	C610	I	-	6G	SATA 6G	DMI 2.0 x4	4 (6)	JBOD, 0, 1, 10	-
PRAID CM400i <sup>1)</sup>	PRAID CM400i	M	-	12G	SATA 6G SAS 12G	PCIe 3.0 x8	8	0, 1	-
PRAID CP400i	PRAID CP400i	P	-	12G	SATA 6G SAS 12G	PCIe 3.0 x8	8	0, 1, 1E, 5, 10, 50	-
PSAS CP400i	PSAS CP400i	P	-	12G	SATA 6G SAS 12G	PCIe 3.0 x8	8	JBOD, 0, 1	-
PRAID EM400i <sup>1)2)</sup>	PRAID EP400i	M	1 GB	12G	SATA 6G SAS 12G	PCIe 3.0 x8	8	0, 1, 1E, 5, 6, 10, 50, 60	✓
PRAID EP400i		P							
PRAID EP420i	PRAID EP420i	P	2 GB	12G	SATA 6G SAS 12G	PCIe 3.0 x8	8	0, 1, 1E, 5, 6, 10, 50, 60	✓

<sup>1)</sup> These controllers are only available for PRIMERGY BX2560 M1.

<sup>2)</sup> All RAID levels except for 0 and 1 can only be used for the PRAID EM400i controller in connection with more disk slots in PRIMERGY SX980 Storage Blades.

<sup>3)</sup> Two onboard RAID controllers are integrated in the Intel C610 chipset. Each of these controllers can be used via PRIMERGY RAID Management to form logical drives consisting of up to four hard disks (in the event of the CX25x0 M1 up to six hard disks are even possible with one of these controllers). The alias of the onboard controller is used in this document to refer to one controller instance in the chipset.

The column "FF" expresses the form factor; "I" means "integrated", "P" means "PCIe slot" and "M" means "mezzanine card". The column "Max. # disks" specifies the maximum number of hard disks that can be directly run on the controller within the context of the RAID Management concept of PRIMERGY and PRIMEQUEST servers. This information can be of help in detecting whether the controller could be a theoretical bottleneck. In some PRIMERGY models so-called "expanders" (special components defined in the SAS standard) are used – in connection with specific controller models – in order to further increase the maximum number of hard disks. In so doing, the expander cannot increase the bandwidth of the existing ports, but makes it available in total to all connected hard disks.

In the evaluation of the performance of disk subsystems, processor performance and memory configuration do not for the most part play a significant role in today's systems - a possible bottleneck usually affects the hard disks and the RAID controller, and not CPU or memory of the server system. Thus the various RAID controllers can be compared independently of the PRIMERGY or PRIMEQUEST models in which they are used - even if all the configurations are not possible in all PRIMERGYs or PRIMEQUESTs due to their expandability with hard disks.

The following table is a compilation of which RAID controllers of the current generation are released in the individual PRIMERGY and PRIMEQUEST systems for the connection of hard disks at the time this white paper was written and how many hard disks the single RAID controllers support in these models at most. Please see the configurators of the systems for the possible combinations of PRIMERGY and PRIMEQUEST configuration versions and controllers.

System	Expander	Onboard controller			Controller with PCIe interface					
		C220	C236	C610	PRAID CM400i	PRAID CP400i	PSAS CP400i	PRAID EM400i	PRAID EP400i	PRAID EP420i
PRIMERGY BX2560 M1				2	2			2		
PRIMERGY BX2580 M1				2						
PRIMERGY CX2550 M1				6		6	6		6	6
PRIMERGY CX2570 M1				6		6			6	6
PRIMEQUEST 2800B2 (DU)										4
PRIMEQUEST 2x00E2 (SB)										4
PRIMEQUEST 2x00E2 (DU)										4
PRIMERGY RX1330 M1	-/✓	4				8	10		10	10
PRIMERGY RX1330 M2	-/✓		4			8			10	10
PRIMERGY RX2530 M1	-/✓			4 (8)		8	10		10	10
PRIMERGY RX2540 M1	-/✓			4 (8)		8	24		24	24
PRIMERGY RX2560 M1	-/✓					8	32		32	32
PRIMERGY RX4770 M2						8			8	8
PRIMERGY SX960 S1	-/✓						10		10	10
PRIMERGY TX1310 M1		4								
PRIMERGY TX1320 M1		4								
PRIMERGY TX1320 M2			4			6			6	6
PRIMERGY TX1330 M1		4								
PRIMERGY TX1330 M2	-/✓		4			8			24	24
PRIMERGY TX2560 M1	-/✓					8	32		32	32

RAID controllers of previous generations (SAS-6G) can also be ordered for some systems. Since these controllers have already been dealt with in the previous document [RAID Controller Performance 2013](#), they will not be analyzed again here.

The abbreviation "DU" stands for "Disk Unit", and "SB" stands for "System Board" for PRIMEQUEST systems. The figures in the corresponding table lines specify in each case the maximum number of hard disks in such a sub-unit.

This white paper only examines the previously mentioned mezzanine cards in connection with internal hard disks in the same server blade.

In connection with hard disks, the PSAS CP400i is essentially planned for Microsoft Windows Server 2012 Storage Spaces. For this purpose, this controller passes on the physical drives to the operating system in an unchanged state. The hardware RAID support that is also available in the controller offers RAID 0 and RAID 1 and is intended for a boot drive.

## Controller interfaces and their throughput limits

A RAID controller needs an interface on the one hand to the hard disks and on the other hand to the CPU. The first one is typically SAS or SATA, the second one is typically PCIe or, in the event of integrated onboard controllers, DMI. The upper limits for the throughputs of SAS, SATA, PCIe and DMI have been put together below.

### SAS and SATA

“Serial Attached SCSI” (SAS) and “Serial Advanced Technology Attachment” (SATA) are serial interfaces, whose data throughput depends on the frequency. These interfaces are used to connect non-volatile storage media, such as hard disks, optical drives and tape drives.

Type	Frequency	Theoretical throughput	Practical throughput (90%)
SAS 3G / SATA 3G	3000 MHz	286 MB/s	257 MB/s
SAS 6G / SATA 6G	6000 MHz	572 MB/s	515 MB/s
SAS 12G	12000 MHz	1144 MB/s	1030 MB/s

Alternatively, a version number is also used with SAS - 1.0 for 3G, 2.0 for 6G and 3.0 for 12G. Alternatively, version number 2.0 is used for 3G and 3.0 for 6G with SATA.

The theoretically achievable throughput is calculated as follows: 1 bit per 1 Hz, minus 20% redundancy of the serial transfer due to the so-called 8b/10b coding. The throughput that can be achieved in practice can be estimated by multiplying this with 0.90. This 90% is a mean empirical value taken from the values that have been observed over the years for various components.

All the components of a connection between end devices must use the same version of the SAS or SATA protocol. In addition to the hard disks, these also include the controllers and any expanders that are possibly used. If different components come together here, the most high-performance standard that is jointly supported by all components is automatically used, i.e. a lower frequency is possible. In this respect, the higher protocols are downwards compatible.

Whereas each port with SATA is often individually connected to a hard disk, four SAS connections and cables are frequently put together and referred to as an “x4 SAS” or “x4 wide port”. This makes it possible to directly connect a maximum of four SAS hard disks via a backplane. The throughput of x4 SAS is four times that of the corresponding individual SAS connection; this also applies similarly for SATA.

Interface	Connection	Frequency	Theoretical throughput	Practical throughput (90%)
SAS 3G / SATA 3G	1 x x4	3000 MHz	1144 MB/s	1030 MB/s
SAS 3G / SATA 3G	2 x x4	3000 MHz	2289 MB/s	2060 MB/s
SAS 6G / SATA 6G	1 x x4	6000 MHz	2289 MB/s	2060 MB/s
SAS 6G / SATA 6G	2 x x4	6000 MHz	4578 MB/s	4120 MB/s
SAS 12G	1 x x4	12000 MHz	4578 MB/s	4120 MB/s
SAS 12G	2 x x4	12000 MHz	9155 MB/s	8240 MB/s

Some PRIMERGY models can be expanded with a larger number of hard disks than the controller has hard disk connections. In this case, the number of connectable hard disks is increased by means of an expander. As already mentioned, an expander can only distribute the data flow, not increase the throughput.

The SAS protocol is defined in such a way that it can also transport the SATA protocols of the same or a lower frequency (tunneling). This enables the controllers of both SAS versions to communicate with SATA hard disks. Conversely, it is not possible to connect SAS hard disks via a SATA interface.



## PCIe and DMI

PCIe is also a serial interface between the controller and the motherboard. The connectors are designed with a different width and number of lanes. x4 (four lanes) and x8 (eight lanes) are normal, whereby the actual number of electrically used lanes is the important thing (here referred to below as the “functional PCIe width”). The throughput of a lane is also determined by the frequency.

Interface	Connection	Frequency	Theoretical throughput	Practical throughput (90%)
PCIe.1.0, PCIe Gen1	x4	2500 MHz	954 MB/s	858 MB/s
PCIe.1.0, PCIe Gen1	x8	2500 MHz	1907 MB/s	1716 MB/s
PCIe.2.0, PCIe Gen2	x4	5000 MHz	1907 MB/s	1716 MB/s
PCIe.2.0, PCIe Gen2	x8	5000 MHz	3815 MB/s	3433 MB/s
PCIe 3.0, PCIe Gen3	x4	8000 MHz	3756 MB/s	3380 MB/s
PCIe 3.0, PCIe Gen3	x8	8000 MHz	7512 MB/s	6761 MB/s

PCIe 1.0 is also often referred to as “PCIe Gen1”, PCIe 2.0 as “PCIe Gen2” and PCIe 3.0 as “PCIe Gen3”.

The theoretically achievable throughput is calculated as follows: 1 bit per 1 Hz multiplied by the number of connections (x4 or x8), minus 20% redundancy of the serial transfer due to the so-called 8b/10b coding for PCIe 1.0 and 2.0 or minus 1.54% redundancy due to the 128b/130b coding for PCIe 3.0 respectively. The throughput that can be achieved in practice can be estimated by multiplying this with 0.90. This 90% value is a mean empirical value taken from the values for various components that have been observed over the years.

All PRIMERGY servers, beginning with the generation introduced in 2010 (e. g. PRIMERGY RX300 S5), support PCIe 2.0 and from the generation introduced in 2012 (e. g. PRIMERGY RX300 S7) PCIe 3.0. If different components come together here, the highest frequency jointly supported by all components is used.

The Direct Media Interface, or in its abbreviated form DMI, is closely related to PCIe. This is an Intel-specific standard for connecting a CPU to the chipset. The corresponding statements apply for DMI with regard to the throughputs, as do those for PCIe in the above table. Thus, for example DMI 2.0, x4, permits a maximum practical throughput of 1716 MB/s. On the input side (CPU side) this throughput value is relevant for the onboard controllers, as these are accommodated in the chipsets.

### Application to the RAID controllers

The next table presents the performance-determining key data for all RAID controllers. The throughput limits listed here are obtained with the help of the two previous subsections “[SAS and SATA](#)” and “[PCIe and DMI](#)”. The significant throughput limit in each case is highlighted in the table in bold print.

Controller alias	# CPU cores	Cache memory type	# Disk side data channels	Limit for throughput of disk interface	# CPU-side data channels	Limit for throughput of CPU-side interface
C220			4 x SATA 6G	2060 MB/s	4 x DMI 2.0	<b>1716 MB/s</b>
C236			4 x SATA 6G	<b>2060 MB/s</b>	4 x DMI 3.0	3433 MB/s
1 x C610			1 x 4 x SATA 6G	2060 MB/s	4 x DMI 2.0	<b>1716 MB/s</b>
2 x C610			2 x 4 x SATA 6G	4120 MB/s	4 x DMI 2.0	<sup>1)</sup> <b>1716 MB/s</b>
PRAID CM400i	1 x 1.2 GHz		2 x SAS 12G	<b>2060 MB/s</b>	8 x PCIe 3.0	6761 MB/s
				<sup>2)</sup> <b>1030 MB/s</b>		6761 MB/s
PRAID CP400i	1 x 1.2 GHz		8 x SAS 12G	8240 MB/s	8 x PCIe 3.0	<b>6761 MB/s</b>
				<sup>2)</sup> <b>4120 MB/s</b>		6761 MB/s
PSAS CP400i	1 x 1.2 GHz		8 x SAS 12G	8240 MB/s	8 x PCIe 3.0	<b>6761 MB/s</b>
				<sup>2)</sup> <b>4120 MB/s</b>		6761 MB/s
PRAID EM400i	2 x 1.2 GHz	DDR3 / 1866 MHz	2 x SAS 12G	<b>2060 MB/s</b>	8 x PCIe 3.0	6761 MB/s
				<sup>2)</sup> <b>1030 MB/s</b>		6761 MB/s
PRAID EP400i	2 x 1.2 GHz	DDR3 / 1866 MHz	8 x SAS 12G	8240 MB/s	8 x PCIe 3.0	<b>6761 MB/s</b>
				<sup>2)</sup> <b>4120 MB/s</b>		6761 MB/s
PRAID EP420i	2 x 1.2 GHz	DDR3 / 1866 MHz	8 x SAS 12G	8240 MB/s	8 x PCIe 3.0	<b>6761 MB/s</b>
				<sup>2)</sup> <b>4120 MB/s</b>		6761 MB/s

1) The second controller instance does not increase the throughput limit of the CPU-side interface.

2) This halved throughput limit applies for the case, in which only hard disks with a 6G interface are connected to the controller.

In the majority of cases the throughput limits do not represent a bottleneck. In practice, the application scenarios with random access to conventional hard disks prevail in particular, in which no high throughputs are achieved.

The throughput values in the column “Limit for throughput of disk interface” apply for the connections between the controller and the hard disks in their entirety. The throughputs via this SAS/SATA interface are only in the case of RAID 0 identical with the throughputs from the viewpoint of the application. With other RAID levels the throughput via the SAS/SATA interface is from the viewpoint of the application multiplied by a specific factor compared with the throughput. This factor is always  $\geq 1$  and depends on the RAID level and several characteristics of the access pattern. The real throughput limits are therefore always lower than the values in the column “Limit for throughput of disk interface” by the mentioned specific factor.



## **Safeguarding the controller cache against power failure**

In order to back up data in the cache in the event of a power failure there are as a matter of principle two options for PRIMERGY and PRIMEQUEST servers:

### ***Battery Backup Unit (BBU)***

The conventional method works with a battery backup unit (BBU). In this case, the power supply of the volatile cache memory is backed up during a power failure by means of a rechargeable battery (accumulator). Since the battery has a limited capacity, which also decreases with time due to physical and chemical processes, it can only ensure supply to the cache memory for a limited amount of time. This time is not constant, but depends on several influencing factors, such as the age of the BBU, charging status, temperature, etc. Thus, the BBU is subject to a limited warranty.

As soon as the power is available again and the server has been powered up, the RAID controller can continue to work with the content of the cache memory and the data remains consistent.

### ***Flash Backup Unit (FBU)***

With a flash backup unit (FBU), which is based on more recent technology, the data is not kept in the cache memory in case of a power failure; the content of the cache memory is in contrast copied to a non-volatile flash memory. The data can remain in this flash memory for almost as long as you like, which means that the retention time known from the BBU is no longer a problem.

The energy required for the copying process from the cache memory to the flash memory in case of a power failure comes from a super capacitor.

As soon as the power is available again and the server has been powered up, the cache content is written from the flash memory back to the cache memory. The RAID controller can now continue to work again and the data remains consistent.

The FBU version is offered for all the RAID controllers with controller cache that are dealt with in this white paper.

## FastPath

FastPath is a high-performance IO accelerator for logical drives that are made up of SSDs. This optimized version of LSI MegaRAID technology permits a clear-cut increase in the performance of applications with a high IO load for random access if SSDs are used.

FastPath used to be part of the RAID option “RAID Advanced Software Options int.” which could be ordered in addition to a RAID controller.

From firmware package version 24.7.0-0061, FastPath has automatically been active in the 12G-enabled RAID controllers with cache (PRAID EM400i, PRAID EP400i and PRAID EP420i) and effective for newly created logical drives or ones that were already created with older firmware versions. You should merely ensure that there are generally optimal prerequisites for SSDs as far as the cache settings are concerned. This means that when creating a logical drive with the ServerView RAID Manager the cache settings must be set en bloc to “Fast Path optimum”, and with the existing logical drives you should ensure that the settings are as follows:

- Read Mode “No read-ahead”
- Write Mode “Write-through”
- Cache Mode “Direct”
- Disk Cache “Enabled”

In the remainder of this document it is assumed that on account of the firmware status FastPath is active.

## Performance-relevant cache settings

Accurate parameter setting for the respective logical drive is essential to the optimal performance of the RAID controller. Depending on the controller there is a varying number of parameters that can be set. For the purpose of easy and reliable handling of the settings for RAID controllers and hard disks it is advisable to use the software “ServerView RAID Manager” (version  $\geq 6.3.3$ ), which is supplied for PRIMERGY and PRIMEQUEST servers. All the cache settings for controllers and hard disks can usually be made en bloc when creating a logical drive – specifically for the application – by using the pre-defined modi “Performance”, “Data Protection” or “Fast Path optimum”.

In “Data Protection” mode protection against data loss is ensured in the case of a power failure. This means that the write caches of the RAID controller and hard disks are normally disabled. If the RAID controller has a write cache, which is safeguarded against power failure by an operational FBU, this write cache is enabled. As a result, there are usually major performance advantages.

Any existing controller and hard disk caches are enabled by the “Performance” mode, which is why the cache of the RAID controller should be protected in this mode against data loss in the event of a power failure by means of a flash back unit (FBU). Furthermore, the hard disk caches should also be safeguarded by the use of an uninterruptible power supply (UPS). In this case, the server in use should also have redundant power supply units. The “Performance” mode ensures the best possible performance settings for the majority of the application scenarios with HDDs.

The “Fast Path optimum” mode is only displayed if the option FastPath is in fact active in the RAID controller. It should be selected if maximum transaction rates with SSDs are to be achieved for random accesses with small blocks ( $\leq 8$  kB, e. g. OLTP operation of databases). In such cases, the write and read cache of the RAID controller – based on the already very short access times of SSDs – act predominantly as a brake and are thus disabled by this mode.

In special cases, a parameter setting that deviates from the standard setting of the “Performance” mode can make sense. Reference is made to whether it makes sense in the appropriate part of the section “[Controller comparison](#)”.

The cache settings of the “ServerView RAID Manager” software contain – depending on the controller – all or part of the following setting options of the RAID controller and the hard disks. The first three setting options control the RAID controller, and the last one controls the hard disks of the logical drive. All the parameters can be specifically set for each logical drive.

**Read mode**

The “Read mode” parameter can be used to control whether reading is done in advance. Two options “No read-ahead” and “Read-ahead” are available. Reading in advance is not done with “No read-ahead”. Blocks that sequentially follow directly requested blocks are read and transferred to the controller cache in the case of “Read-ahead”. This is done with the expectation that the blocks are also required in one of the next requests.

In the case of the “Read-ahead” option the onboard controllers (e.g. C220) generally read blocks in advance. The PCIe controllers with a cache work in a more differentiated way for this option: The requested blocks are continuously analyzed to see whether there is sequential read access. If the controller detects such an access, it starts to also read the sequentially following blocks – in addition to the requested block – in the cache in order to have them available for the expected, next requests. The current option “Read-ahead” is in other words adaptive. This is a merger of the two previous options “Read-ahead” and “Adaptive”.

**Write mode**

The setting options of the controller cache that control the handling of write requests are summarized under the term “Write mode”. There are three options for setting the write cache: “Write-through”, “Write-back” and “Always Write-back (independent of BBU state)”. The “Write-through” option ensures that each write request from the controller is only reported back as completed when it has been acknowledged by the hard disk. With the “Write-back” and “Always Write-back” options the requests are cached in the controller cache, immediately acknowledged to the application as completed and only transferred to the hard disk later. This procedure enables optimal utilization of controller resources, faster succession of the write requests and therefore higher throughput. Any power failures can be bridged by means of an optional FBU, thus guaranteeing the integrity of the data in the controller cache. The “Always Write-back” option enables the write cache on a permanent basis; it is also used if the FBU is not operational. On the other hand, the “Write-back” option automatically switches to “Write-through” as long as the controller cache is not safeguarded by the FBU.

**Cache mode**

The “Cache Mode” parameter is also sometimes referred to as the “I/O Cache”. The “Direct” option specifies that the data to be read are transferred directly from the hard disk to the RAM of the server. The alternative “Cached” causes all the data to be read and written on its way between the server memory and the hard disks to pass the controller cache. “Direct” is the recommended setting. The Read-Ahead functionality is not influenced by the Cache Mode setting.

**Disk cache mode**

The possible values here are “enabled” and “disabled”. In most cases, the enabling of the hard disk cache entails an increase in throughput for write access. If the system is safeguarded by a UPS, the enabling of the hard disk cache is recommended for performance reasons.

The next table shows which of these setting options exist for the individual controllers.

Controller alias	Read mode	Write mode	Cache mode
C220, C236, C610	✓		
PRAID CM400i, PRAID CP400i, PSAS CP400i			
PRAID EM400i PRAID EP400i, PRAID EP420i	✓	✓	✓

To complete matters the following table also provides a compilation of the settings that are currently implemented in the modi “Data Protection”, “Performance” and “Fast Path optimum” in ServerView RAID Manager. It should be noted that the settings for the controllers with a controller cache also depend on the existence of a FBU, but are independent of the selected RAID level.

Controller alias		C220, C236, C610	PRAID CM400i, PRAID CP400i, PSAS CP400i	PRAID EM400i, PRAID EP400i, PRAID EP420i	
FBU?					✓
Data Protection	Read mode	Read-ahead		Read-ahead	Read-ahead
	Write mode			Write-through	Write-back
	Cache mode			Direct	Direct
	Disk cache	off	off	off	off
Performance	Read mode	Read-ahead		Read-ahead	Read-ahead
	Write mode			Always Write-back	Write-back
	Cache mode			Direct	Direct
	Disk cache	on	on	on	on
Fast Path optimum	Read mode			No read-ahead	No read-ahead
	Write mode			Write-through	Write-through
	Cache mode			Direct	Direct
	Disk cache		on	on	on

## Other settings

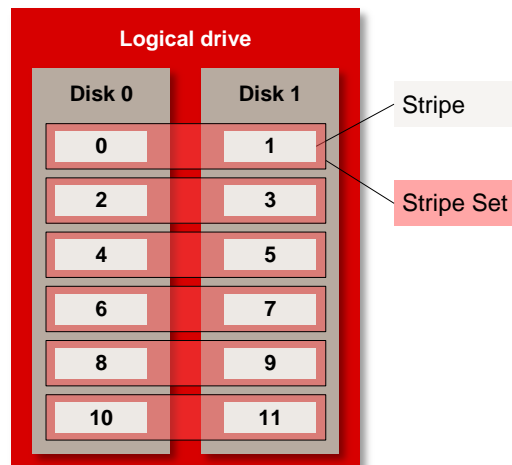
In addition to the setting options for the caches of RAID controllers and hard disks, there are further setting options in the “ServerView RAID Manager” (version  $\geq$  6.3.3) for logical drives, and knowledge of these options is worthwhile from a performance viewpoint.

### Stripe size

The first interesting parameter is the Stripe size. It can only be set when you create a logical drive. Various values are possible for the RAID controllers with cache (e. g. PRAID EP400i); the default value for all other controllers is 64 kB.

The significance of the stripe size is to be explained in detail below using the example of the simplest case RAID 0.

The stripe size controls the design of logical drives that are made up of physical hard disks. The controller implements access to a data block of a logical drive by converting the addresses in the logical drive by means of a specific rule to addresses in the involved physical hard disks. This conversion is based on a division of each of the involved hard disks – beginning in each case from the start of the hard disk - in equally sized blocks of N bytes each. The first N bytes of the logical drive are now assigned to block 0 on hard disk 0, the next N bytes are then assigned to block 0 on hard disk 1. This continues successively until assignment to block 0 has taken place on all the involved hard disks. It then continues with block 1 on hard disk 0, block 1 on hard disk 1, etc. The conversion rule is illustrated by the following diagram:



Each of these blocks on one of these hard disks is called a stripe, and its size in bytes is called the stripe size. All the stripes that lie horizontally adjacent to each other in the above diagram are known as a stripe set.

The stripe size influences performance. On the one hand the stripe size must be small enough to distribute - with a high degree of probability - accesses to the logical drive evenly over the hard disks. On the other hand it must also be large enough to prevent the requested blocks of the logical drive from mostly being divided at the hard disk limits. This would result in an unwanted multiplication of hard disk accesses and thus an unnecessarily early overload of the hard disks.

Normally, the default of the stripe size is optimal. Merely in the case of random accesses the previously described block division should for the most part be avoided. In other words, the stripe size should

- either be large compared with the blocks requested by the application (example: requested blocks of 8 kB for a 64 kB stripe size)
- or be exactly the same size as the blocks requested by the application if the latter aligns them at the stripe limits

The possible values of the stripe size for the RAID controllers with cache that are dealt with here are 64 kB, 128 kB, 256 kB, 512 kB and 1 MB, the default value is 256 kB.

### **Emulation type**

The second interesting parameter is the emulation type. The handling of 512e hard disks should be associated with emulation. The internal structure of such hard disks has a sector size of 4096 B. Externally, however, they emulate a sector size of 512 B. In other words, the physical sector size for such hard disks is 4096 B, and the logical sector size is 512 B. Detailed information on the topic of 512e HDDs is available in the white paper [512e HDDs: Technology, Performance, Configurations](#).

The emulation type can not only be set for the creation of a logical drive; a subsequent change is also possible and this has an effect after the next reboot. There are three possible values:

- |            |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Default    | If only 512n hard disks are contained in a logical drive, it is given the property “logical sector size = 512 B” for the operating system. As soon as at least one 512e hard disk is included, a logical drive has the property “physical sector size = 4096 B”. This default should normally be retained. It provides meaningful parameter information to the accessing software layers that are located above: If the logical drive contains a hard disk with a physical sector size of 4096 B, the software layers located above receive the information and can align their accesses to the logical drive to the physical sectors of 4096 B with an optimal performance level. |
| None       | The logical drive always has the property “physical sector size = 512 B”, even if one of the affected hard disks has the physical sector size 4096 B. This mode does not make sense in productive use.                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| Force 512e | The logical drive always has the property “physical sector size = 4096 B”, even if the physical sector size is only 512 B. In the case of an existing logical drive consisting of 512n hard disks this setting can make sense if you want to ensure that replacing a failed hard disk with a 512e hard disk does not result in losses in performance, either.                                                                                                                                                                                                                                                                                                                      |



## Characteristics of the onboard controllers

Some PRIMERGY models offer with an onboard controller a simple, low-priced entry-level solution for operation with up to four hard disks. With the support of RAID 0, RAID 1 and RAID 10 such an onboard controller offers a range of common RAID levels and does not occupy a PCIe slot.

All the controllers that are connected via the PCIe interface are not treated as onboard controllers here, even if they are integrated on the motherboard (see table in the section "[Presentation of the RAID controllers](#)").

The onboard controller is implemented as a firmware/driver-based software RAID solution. It is integrated in the "Platform Controller Hub" chip, which belongs to the motherboard chip set. During the boot phase accesses to the logical drive are implemented by the firmware. As soon as the operating system is active, suitable drivers take on this task.

The onboard controller does not have a processor of its own, but uses the CPU of the server system for RAID functionality. The pro rata consumption of the server's processor performance is increasingly less important in newer servers.

### **C220, C236 and C610**

These onboard controllers are pure SATA controllers. The C220 onboard controller is included in the chipset of the generation of 1-socket servers that was introduced in 2013, the C610 onboard controller is in the chipset of the generation of 2-socket servers introduced in 2014 and the C236 onboard controller is in the chipset of the generation of 1-socket servers introduced in 2015. These controllers can be set in various modes via the BIOS. Although only the "RAID" mode is suitable for the effective use of this type of controller, all the modes of these SATA controllers are presented for the sake of completeness. Support is not provided for the advanced SATA features Native Command Queuing (NCQ) and "hot swapping" in all cases. There are three modes:

RAID	Recommended mode on account of its flexibility. A trouble-free migration of a SATA-HDD from a non-RAID to a RAID configuration is only possible here. All the functionalities of SATA are supported, i.e. also NCQ and "hot swapping". A firmware named "LSI Logic Embedded MegaRAID" is integrated in the controller BIOS in PRIMERGY servers for the supported RAID levels. Only in this mode are logical drives already possible during the boot phase, and only in this mode can controllers and hard disks be seen in the "ServerView RAID Manager", where they can be administered. Special drivers are needed.
AHCI	AHCI stands for "Advanced Host Controller Interface" and is a cross-manufacturer interface standard for SATA controllers. Support is provided for NCQ and "hot swapping". Special drivers in the operating system are also necessary for AHCI.
IDE	In this operating mode the SATA ports as such are made visible to the operating system. NCQ is not supported. Appropriate SATA drivers are required, which are supplied for various operating systems on the "ServerStart DVD".

## Measurement context

Now that the various controllers have been presented and their technical features explained, it is our intention in the following section "[Controller comparison](#)" to discuss the controllers in various application scenarios and to back this up on the basis of measurement results. Hence, a brief introduction to begin with of the measurement method and the measurement environment.

All the details of the measurement method and the basics of disk I/O performance are described in the white paper "[Basics of Disk I/O Performance](#)".

## Measurement method

As standard, performance measurements of disk subsystems in PRIMERGY and PRIMEQUEST servers are carried out with a defined measurement method, which models the hard disk accesses of real application scenarios on the basis of specifications.

The essential specifications are:

- Share of random accesses / sequential accesses
- Share of read / write access types
- Block size (kB)
- Number of parallel accesses (# of outstanding I/Os)

A given value combination of these specifications is known as "load profile". The following five standard load profiles can be allocated to typical application scenarios:

Standard load profile	Access	Type of access		Block size [kB]	Application
		read	write		
File copy	random	50%	50%	64	Copying of files
File server	random	67%	33%	64	File server
Database	random	67%	33%	8	Database (data transfer) Mail server
Streaming	sequential	100%	0%	64	Database (log file), Data backup; Video streaming (partial)
Restore	sequential	0%	100%	64	Restoring of files

In order to model applications that access in parallel with a different load intensity, the "# of Outstanding I/Os" is increased from 1 to 512 (in steps to the power of two).

The measurements of this document are based on these standard load profiles.

The main results of a measurement are:

- Throughput [MB/s]      Throughput in megabytes per second
- Transactions [IO/s]      Transaction rate in I/O operations per second
- Latency [ms]              Average response time in ms

The data throughput has established itself as the normal measurement variable for sequential load profiles, whereas the measurement variable "transaction rate" is mostly used for random load profiles with their small block sizes. Data throughput and transaction rate are directly proportional to each other and can be transferred to each other according to the formula

<i>Data throughput [MB/s]</i>	$= \text{Transaction rate [IO/s]} \times \text{Block size [MB]}$
<i>Transaction rate [IO/s]</i>	$= \text{Data throughput [MB/s]} / \text{Block size [MB]}$

## Measurement environment

All the measurement results discussed in this document were determined using the hardware and software components listed below:

System Under Test (SUT)		
<b>Hardware</b>		
Model	PRIMERGY RX1330 M1 PRIMERGY RX2540 M1 PRIMERGY RX2560 M1 PRIMERGY TX1330 M2	
Controller	<p>C220: Intel C226 PCH, Code name Lynx Point (in PRIMERGY RX1330 M1) Driver name: megasr1.sys, Driver version: 16.01.2014.0318 BIOS version: A.14.02121826R</p> <p>C236: Intel C236 PCH, Code name Sunrise Point (in PRIMERGY TX1330 M2) Driver name: megasr1.sys, Driver version: 17.01.2015.0716 BIOS version: A.15.08211538R</p> <p>C610: Intel C610 PCH, Code name Wellsburg (in PRIMERGY RX2560 M1) Driver name: megasr1.sys, Driver version: 16.02.2014.0811 BIOS version: A.14.02121826R</p> <p>PRAID CM400i, PRAID CP400i, PRAID EM400i, PRAID EP400i, PRAID EP420i: Driver name: megasas2.sys, Driver version: 6.706.06 Firmware package: 24.7.0-0061</p> <p>PSAS CP400i: Driver name: lsi_sas3.sys, Driver version: 2.50.85.00 Firmware: 05.00.00.00</p>	
Storage media	SSDs	HDDs
	<p>SAS-12G: Toshiba PX02SMF040</p> <p>SATA-6G: Intel SSDSC2BA400G3C</p>	<p>SAS-12G: HGST HUC156045CSS204</p> <p>SATA-6G: Seagate ST91000640NS</p>
<b>Software</b>		
Operating system	Microsoft Windows Server 2012 Standard R2	
RAID Manager software	ServerView RAID Manager 6.3.4	
Benchmark version	3.0	
RAID type	Logical drive of type RAID 0, 1, 5 or 10	
Stripe size	Controller default (i.e. 256 kB for 12G controllers with cache, 64 kB otherwise)	
Measuring tool	Iometer 1.1.0	
Measurement area	The first 10% of the usable LBA area is used for sequential accesses; the next 25% for random accesses.	
File system	raw	
Total number of Iometer workers	1	
Alignment of Iometer accesses	Aligned to whole multiples of 4096 bytes	

The hard disk models used for the controller comparison are summarized again below in detail together with their fundamental performance data, because this is important for your understanding of the performance values achieved with the controllers. A high-performance SATA-6G and SAS-12G hard disk were chosen in each case for the classic hard disks (HDDs), and a SAS-12G-SSD and a SATA-6G-SSD represent the SSD class.

The table depicts the maximum values measured with a single hard disk for the five standard load profiles that were shown in the previous subsection "[Measurement method](#)". The hard disk cache is enabled in all cases, because this almost always ensures optimal performance.

Hard disk type	Short name (alias)	Sequential maximum throughput [MB/s] 64 kB block size		Maximum transaction rates for random accesses [IO/s]		
				8 kB block size	64 kB block size	
		Read	Write	Read share:		
				67%	67%	50%
HDD SATA, 6 Gb/s, 2.5" 1000 GB, 7200 rpm hot-pluggable category: Business-Critical (BC)	SATA-6G-HDD	108 MB/s	108 MB/s	302 IO/s	258 IO/s	243 IO/s
HDD SAS, 12 Gb/s, 2.5" 450 GB, 15000 rpm, hot-pluggable category: Enterprise (EP)	SAS-12G-HDD	237 MB/s	237 MB/s	744 IO/s	608 IO/s	631 IO/s
SSD SATA, 6 Gb/s, 2.5" 400 GB, hot-pluggable category: Enterprise (EP)	SATA-6G-SSD	468 MB/s	436 MB/s	41005 IO/s	5268 IO/s	5206 IO/s
SSD SAS, 12 Gb/s, 2.5" 400 GB, hot-pluggable category: Enterprise (EP)	SAS-12G-SSD	950 MB/s	420 MB/s	55865 IO/s	7599 IO/s	6715 IO/s

## Controller comparison

All the important preliminary information about controllers has been provided in the previous sections. This information will in many cases already narrow down the choice of controller for a given application. If further customer information about the planned use of the controller is added, a great deal more can be said about the performance to be expected with the individual controllers. Thus this section is to compare the controllers differentiated for various RAID levels, application scenarios, load intensities, numbers of hard disks as well as hard disk technologies. The statements are illustrated with the help of measurement results. The comparisons are divided into the following subsections, which can be read independently of each other:

- [RAID 1 \(two SATA hard disks\)](#)
- [RAID 0 and 10 \(four SATA hard disks\)](#)
- [RAID 0, 10 and 5 \(eight SAS hard disks\)](#)
- [RAID 0, 10 and 5 \(more than eight SAS-SSDs\)](#)

General preliminary remarks about the comparisons:

- The five load profiles described in the section "[Measurement method](#)", "File copy", "Database", "File server", "Streaming" and "Restore", are mostly used in the comparisons, thus enabling the random and sequential application scenarios to be reasonably covered. If the customer load profile significantly differs from this, the statements made here no longer apply without restrictions.
- As the benchmark for the performance of a disk subsystem the transaction rate is specified - as is common practice - in IO/s for random load profiles, and throughput in MB/s for sequential load profiles.
- All the controllers that support the RAID level and hard disk type that have just been considered are to be discussed.
- To make things more easily understandable, the diagrams in this section are for the most part restricted to the maximum values achievable. These are usually only achieved with a high load intensity of the disk subsystem.
- The four hard disks that are dealt with more closely in the section "[Measurement environment](#)" (SATA-6G-HDD, SAS-12G-HDD, SATA-6G-SSD and SAS-12G-SSD) are used as example for hard disk technologies. Their key performance data is also presented there. In some places of the following comparisons the achieved performance values are explained on the basis of the performance data of these hard disk types.
- To achieve maximum performance for the measurements the cache and hard disk settings were made as follows using the "ServerView RAID Manager" modes:
  - For SATA-6G-HDDs: "Performance" mode (usually the highest performance mode for HDDs)
  - For SAS-12G-HDDs: "Performance" mode with the only difference: Disk Cache Disabled
  - For SSDs: "Fast Path optimum" mode (if available for the controller), otherwise "Performance".

If there were any exceptions for a certain measurement, this should be mentioned at the point concerned.

- Conventional hard disks (in contrast to SSDs) are now only referred to in short as "HDDs" in the following controller comparisons.

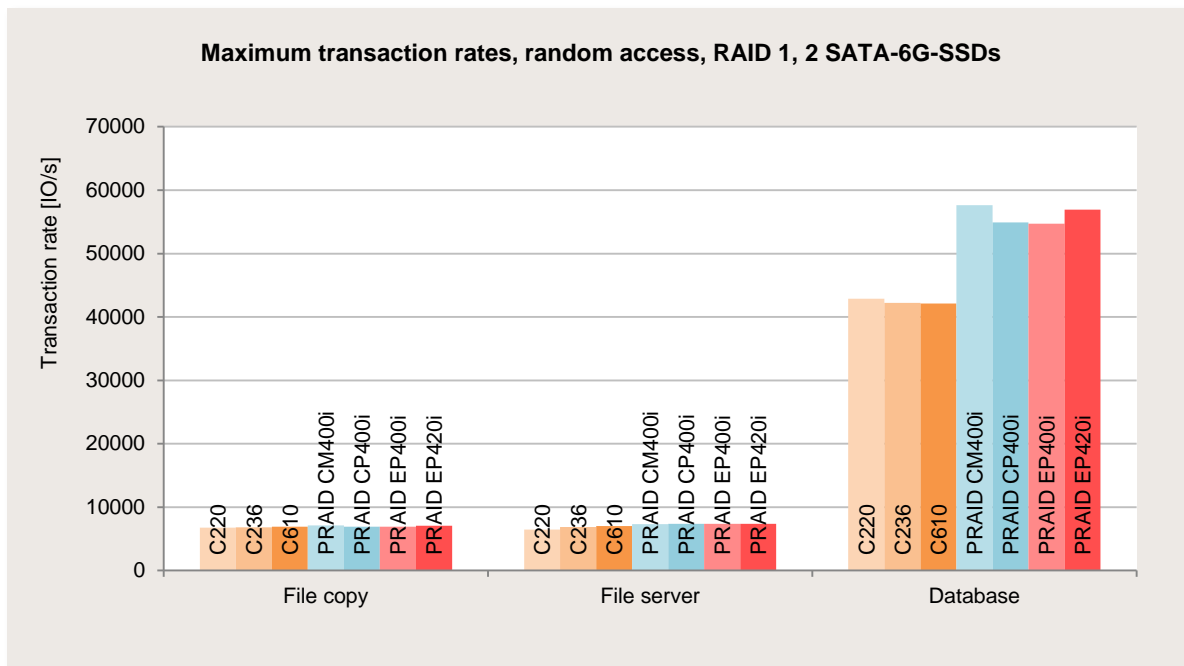
## RAID 1 (two SATA hard disks)

It is worth looking at this case separately, because all onboard controllers can only be compared in this way with most of the PCIe controllers. This is done using the SATA-6G-SSD, which was described more closely in the section [“Measurement environment”](#).

### Random accesses

RAID 1 with two SATA-6G-SSDs

The diagram shows a controller comparison for two SATA-6G-SSDs configured as RAID 1. The three groups of columns in the diagram represent the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



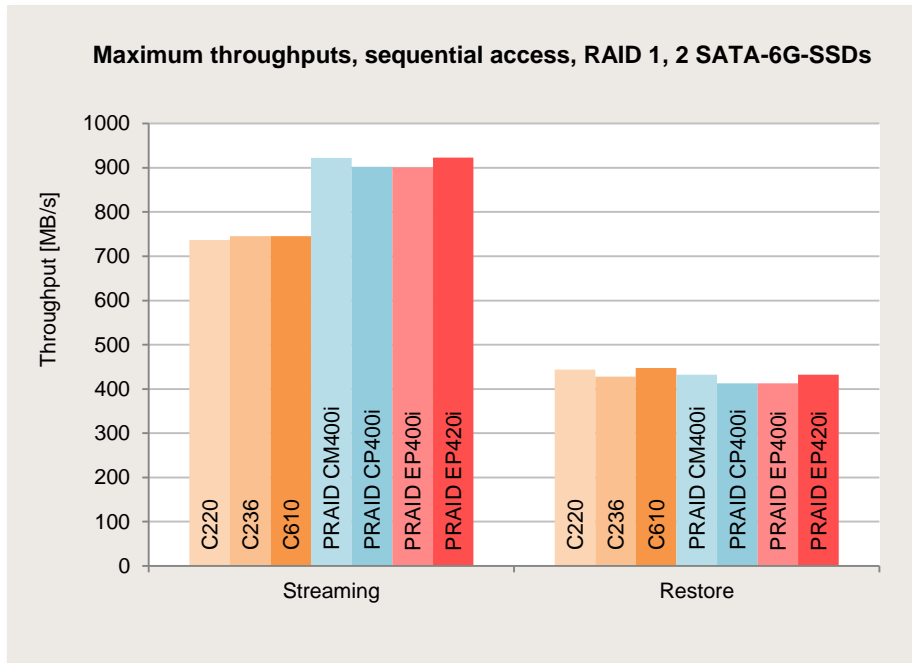
The PCIe controllers provide the highest overall transaction rates here.



## Sequential accesses

### RAID 1 with two SATA-6G-SSDs

The next diagram shows a controller comparison for two SATA-6G-SSDs configured as RAID 1. The two groups of columns in the diagram represent the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).



When reading with higher load intensities the PCIe controllers use both hard disks to a greater extent than the onboard controllers and consequently show a higher maximum throughput.

## RAID 0 and 10 (four SATA hard disks)

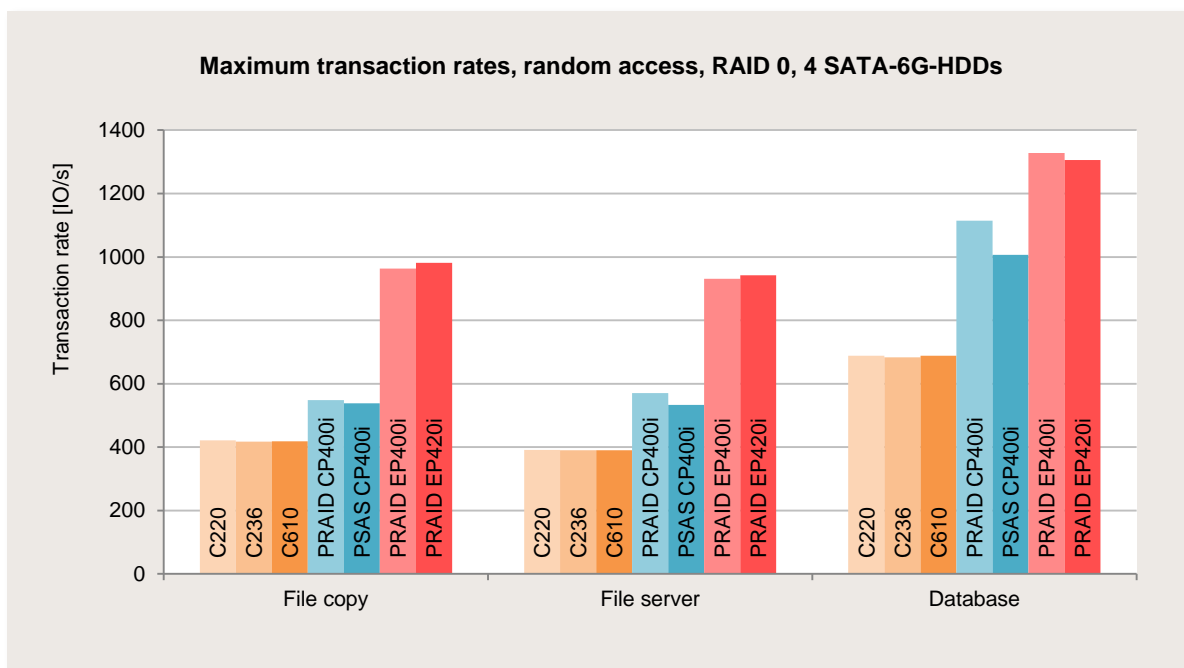
This subsection continues the controller comparisons made in [RAID 1 \(two SATA hard disks\)](#) for four hard disks. It makes sense here to differentiate between HDDs and SSDs, because in the case of the HDDs one controller model more (PSAS CP400i) is released and with the SSDs it is possible to compare the higher performance range of the controllers.

### Random accesses

#### HDDs

##### RAID 0 with four SATA-6G-HDDs

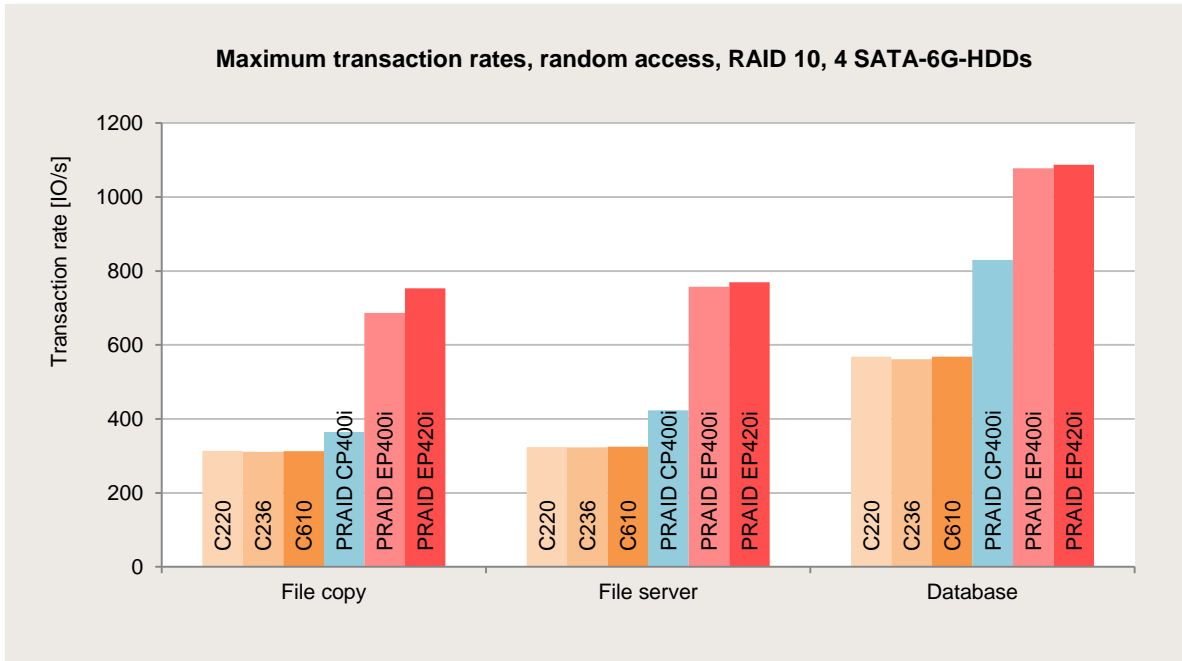
The next diagram shows the transaction rates of the logical drive of type RAID 0 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



It can be clearly seen that the transaction rates are higher if the quality of the controller is higher.

### RAID 10 with four SATA-6G-HDDs

The next diagram shows the transaction rates of the logical drive of type RAID 10 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).

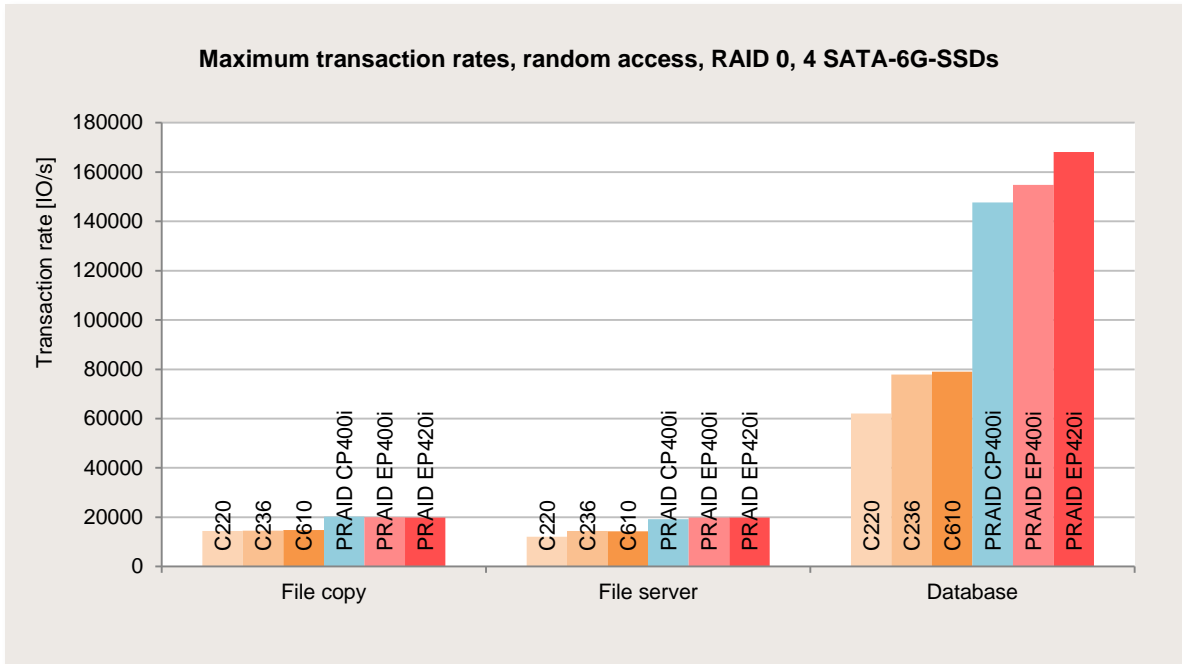


The fact that the transaction rates are higher if the quality of the controller is higher can also be clearly seen here.

## SSDs

### RAID 0 with four SATA-6G-SSDs

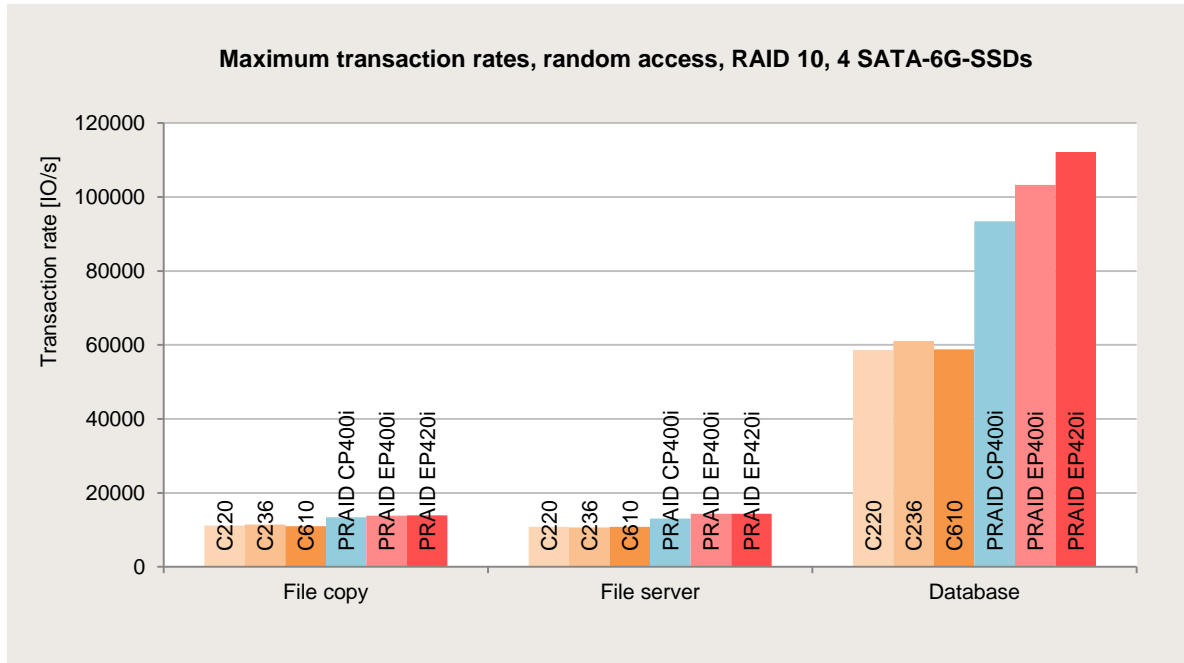
The next diagram shows the transaction rates of the logical drive of type RAID 0 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



It can be clearly seen that the transaction rates are higher if the quality of the controller is higher.

**RAID 10 with four SATA-6G-SSDs**

The next diagram shows the transaction rates of the logical drive of type RAID 10 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



The fact that the transaction rates are higher if the quality of the controller is higher can also be clearly seen here.

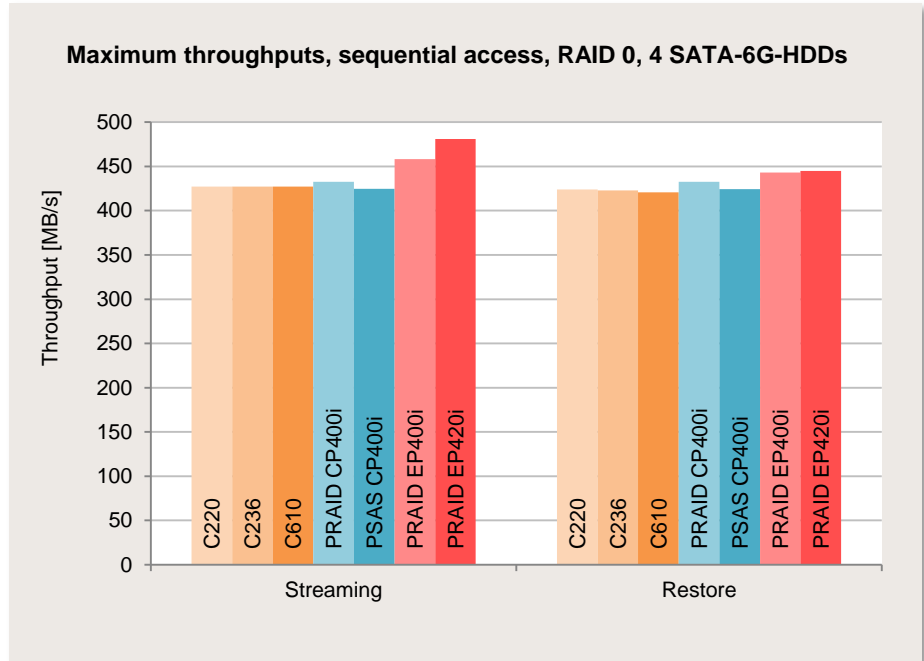
## Sequential accesses

### HDDs

#### RAID 0 with four SATA-6G-HDDs

The next diagram shows the maximum throughputs of the logical drive of type RAID 0 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

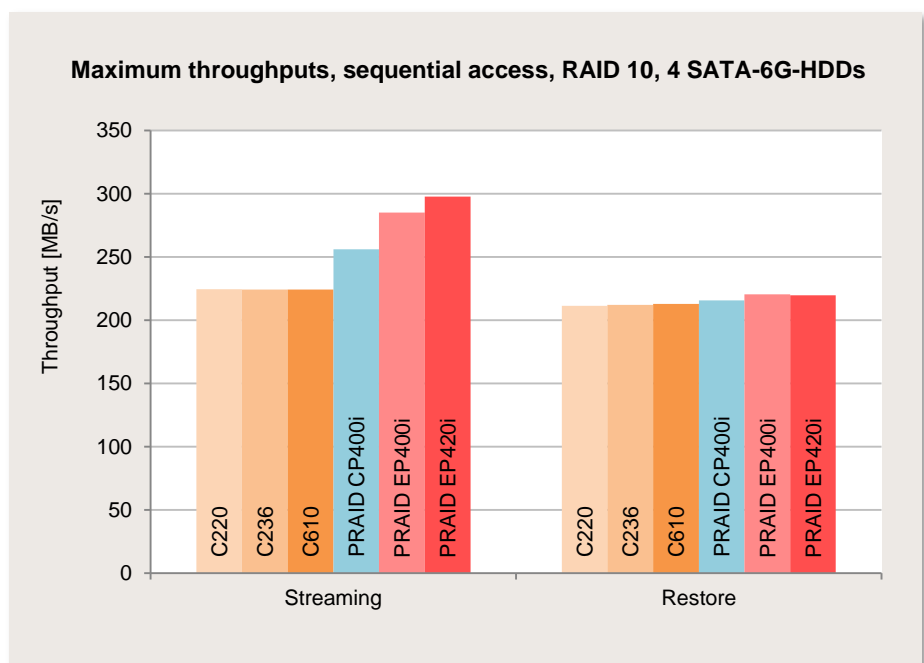
All the controllers deliver approximately the same performance in these cases.



#### RAID 10 with four SATA-6G-HDDs

The next diagram shows the maximum throughputs of the logical drive of type RAID 10 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

In the case of the standard load profile “Streaming” you can see that the PCIe controllers do not to some extent read simultaneously from just two but from four HDDs and can as a result increase the throughput by more than the double value for a single HDD.

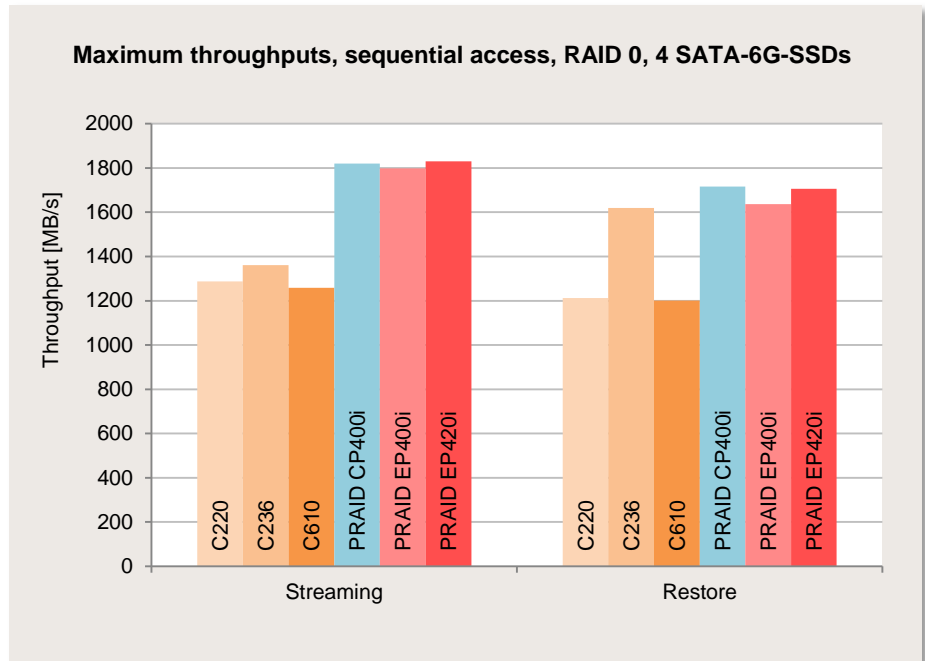




**SSDs**

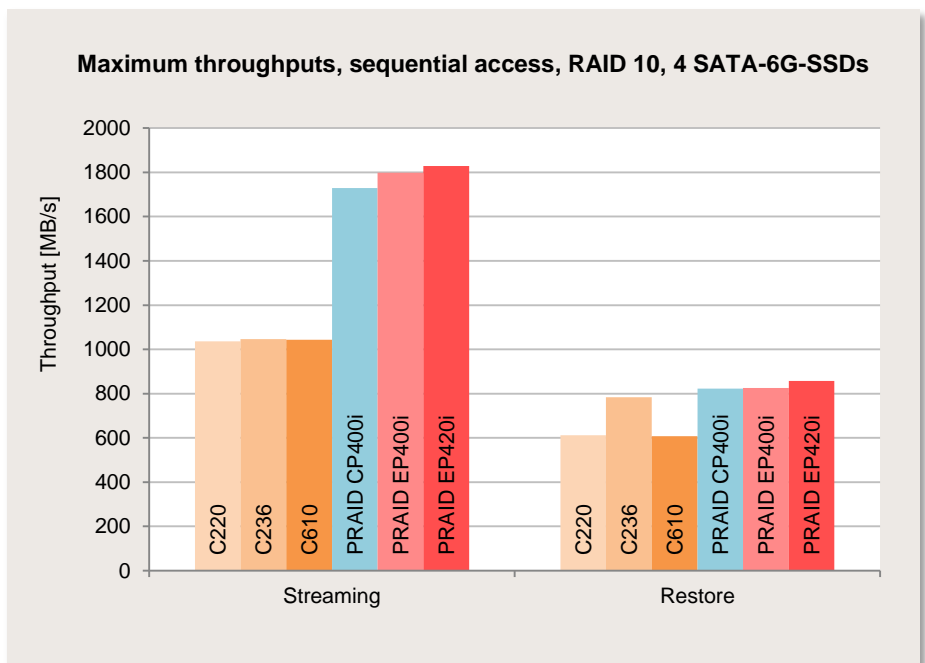
*RAID 0 with four SATA-6G-SSDs*

The next diagram shows the maximum throughputs of the logical drive of type RAID 0 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).



*RAID 10 with four SATA-6G-SSDs*

The next diagram shows the maximum throughputs of the logical drive of type RAID 10 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).



## **RAID 0, 10 and 5 (eight SAS hard disks)**

The onboard controllers are no longer sufficient for the operation of more than six hard disks in the current PRIMERGY servers, therefore only the PCIe controllers are compared below. As one of the PCIe controllers (PRAID CP400i) is released for a maximum of eight hard disks, it makes sense to compare all the controllers for the eight connected hard disks. At the same time the performance values presented in this subsection cover the range of medium numbers of hard disks on a representative basis. Since the maximum values for higher performance requirements are of particular interest here, measurements with high-performance SAS-12G-HDDs or SAS-12G-SSDs are used as a means of illustration. These hard disks are described in more detail in the section "[Measurement environment](#)".

### ***Random accesses***

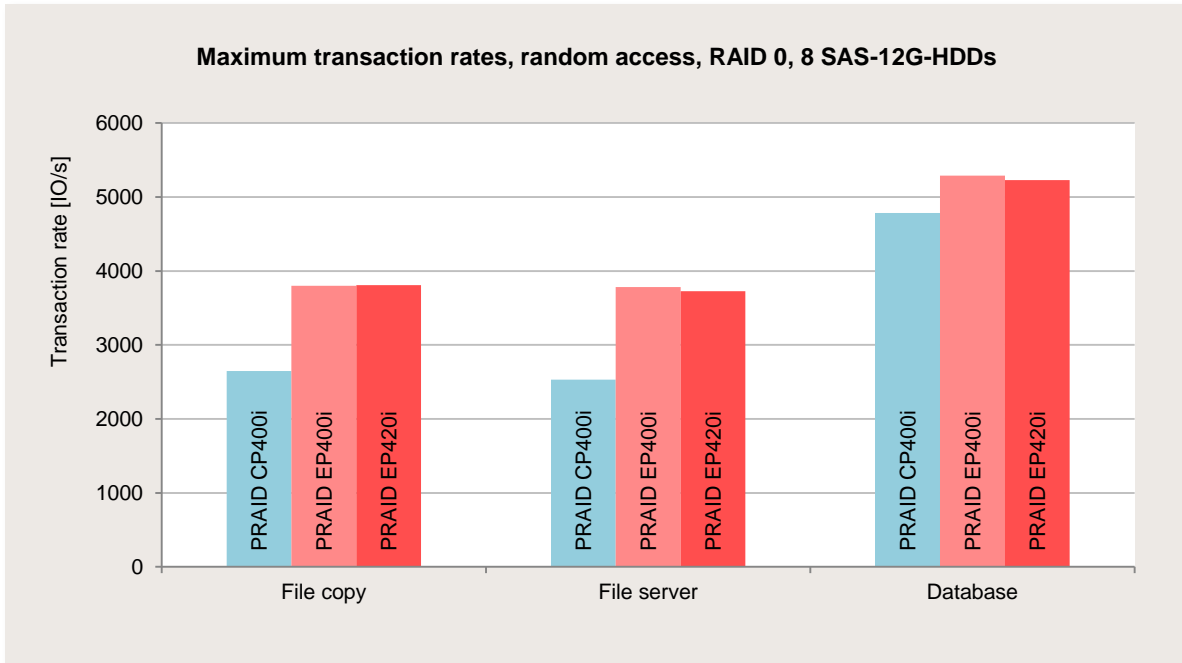
When considering random accesses for larger numbers of hard disks it makes sense to distinguish between HDDs and SSDs, because the maximum values for SSDs are of a quite different magnitude.

#### ***HDDs***

The controllers are compared below with random accesses to HDDs. The maximum transaction rates of the storage medium for the load profile used are the most important limiting factor here. Nevertheless, performance in such cases is not fully independent of the controller. Although the following results were acquired with eight SAS-12G-HDDs, they can also be used to estimate the maximum transaction rates to be expected for other types and numbers ( $\leq 8$ ) of hard disks.

**RAID 0 with eight SAS-12G-HDDs**

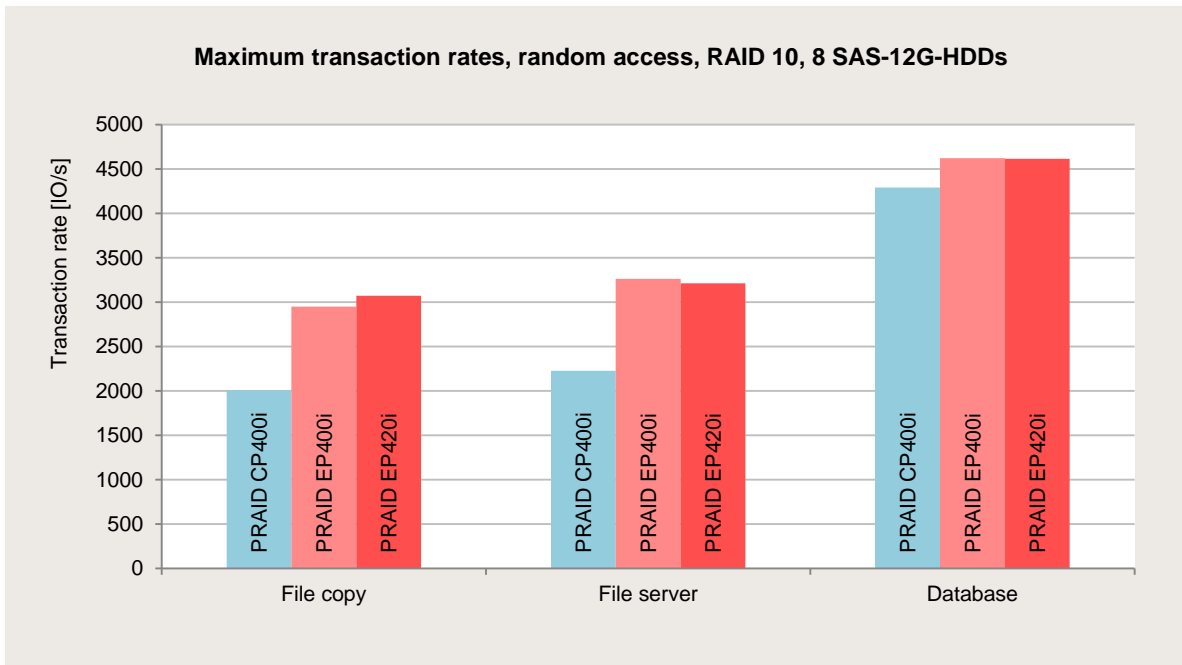
The next diagram shows the transaction rates of the logical drive of type RAID 0 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



The two right-hand columns in each of the three groups of columns in this diagram represent the two controllers with a cache (PRAID EP400i and PRAID EP420i). The superiority of these two controllers is made possible on the one hand by the controller cache, and on the other hand by the higher default value of the stripe size compared with the PRAID CP400i.

**RAID 10 with eight SAS-12G-HDDs**

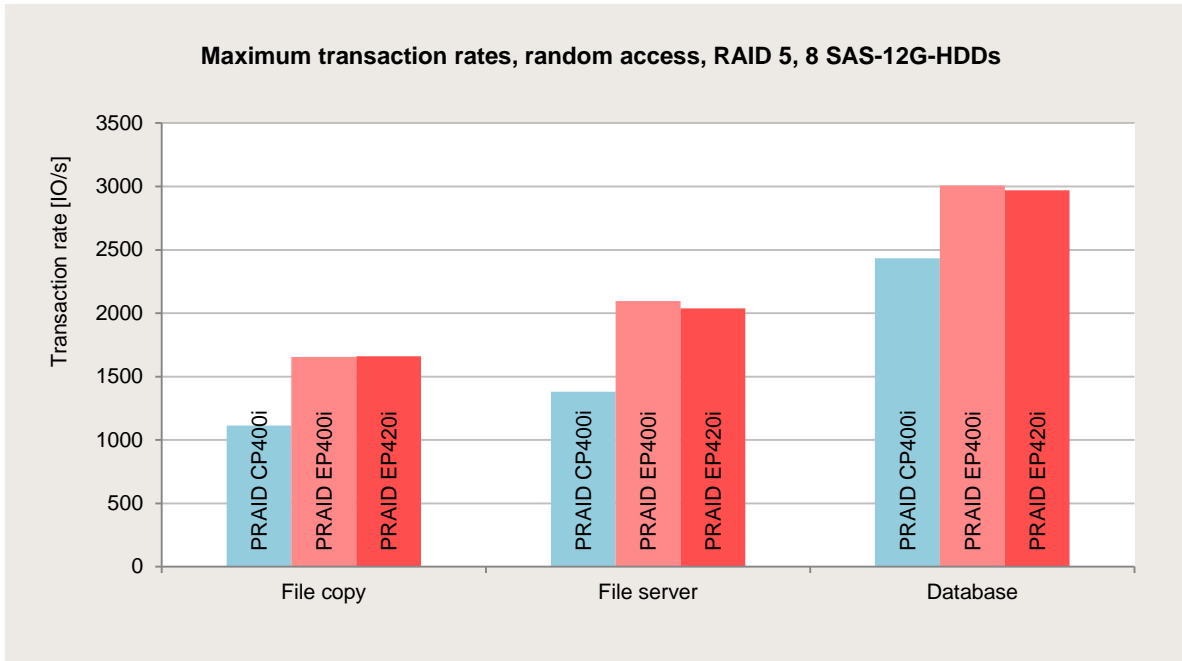
The next diagram shows the transaction rates of the logical drive of type RAID 10 for random load profiles that can be achieved with various controllers. The three groups of columns in the diagram show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



The diagram shows the same principal behavior as with RAID 0.

### RAID 5 with eight SAS-12G-HDDs

The next diagram shows the transaction rates of the logical drive of type RAID 5 for random load profiles that can be achieved with various controllers. The three groups of columns in the diagram show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



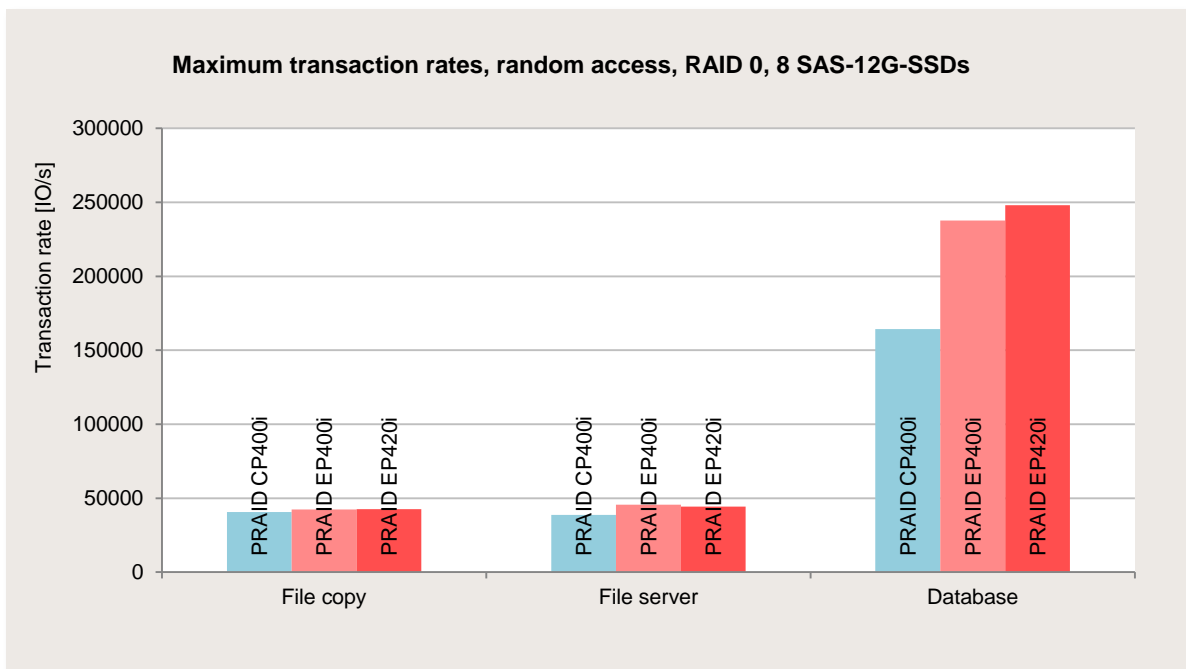
The diagram shows the same principal behavior as with RAID 0.

## SSDs

For the number of SSDs under consideration here the possible transaction rates of a logical drive are so high that the FastPath option, which is enabled as standard in the latest controller firmware, has a distinct influence. This can be seen below by the superiority of the controllers PRAID EP400i and PRAID EP420i compared with the PRAID CP400i. The latter does not support the FastPath option.

### RAID 0 with eight SAS-12G-SSDs

The next diagram shows the transaction rates of the logical drive of type RAID 0 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



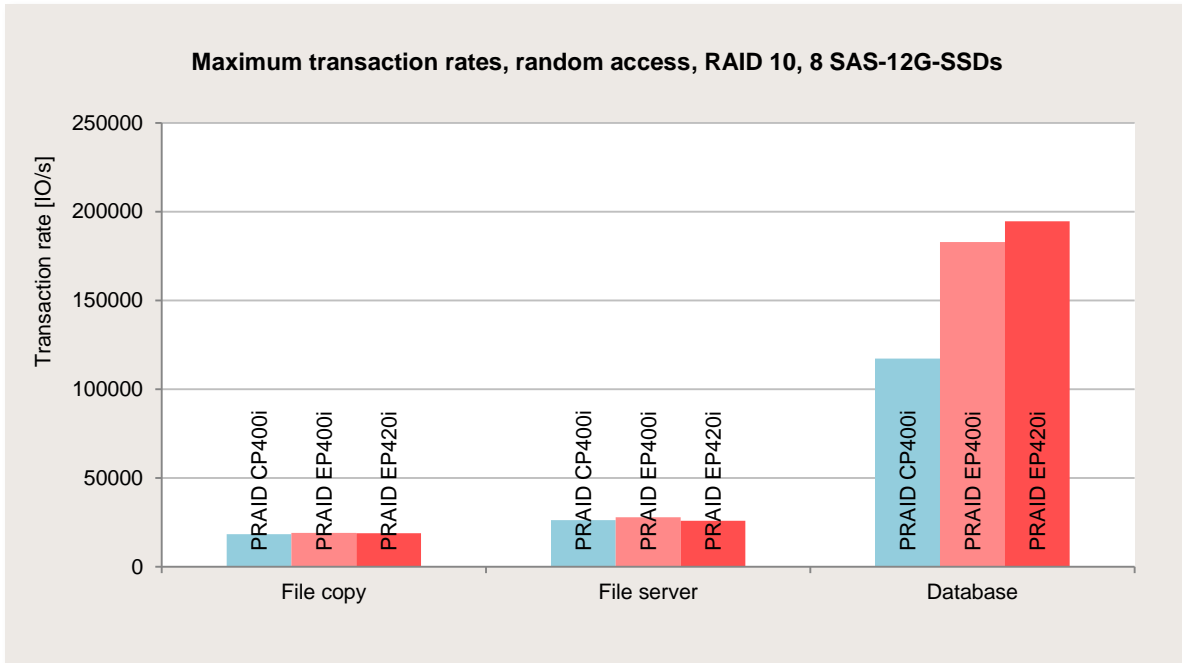
The PRAID EP420i controller is the one with top performance here.

The controllers with cache have a clear advantage for the load profile “Database” (8 kB block size). They also achieve their maximum transaction rate here.

It is also interesting to understand the throughput values that are associated with these transaction rates. Despite the lower transaction rates, the two load profiles with a 64 kB block size have the higher throughputs. For example, the PRAID EP400i controller handles a throughput of about 2848 MB/s with the load profile “File server”. The controllers do not yet have a limiting effect here for the two load profiles “File copy” and “File server” (64 kB block size).

**RAID 10 with eight SAS-12G-SSDs**

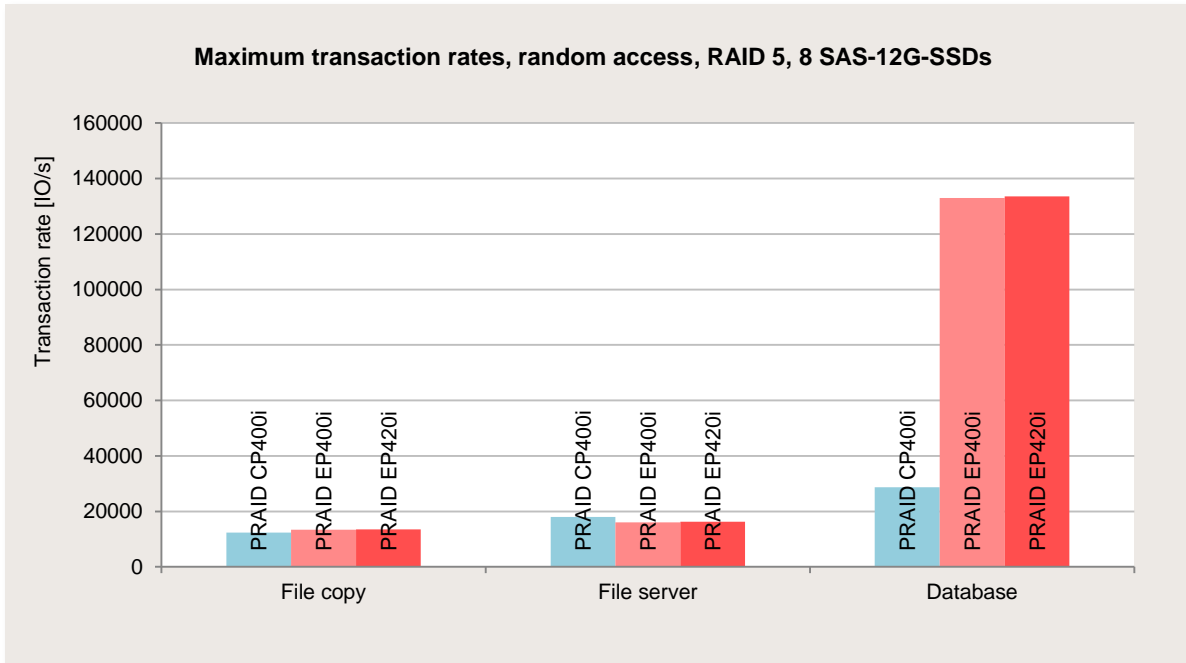
The next diagram shows the transaction rates of the logical drive of type RAID 10 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



In the case of the load profile with the small blocks (“Database”) the controllers with cache have an advantage here, too.

**RAID 5 with eight SAS-12G-SSDs**

The next diagram shows the transaction rates of the logical drive of type RAID 5 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



In the case of the load profile with the small blocks (“Database”) the controllers with cache have an advantage here, too.



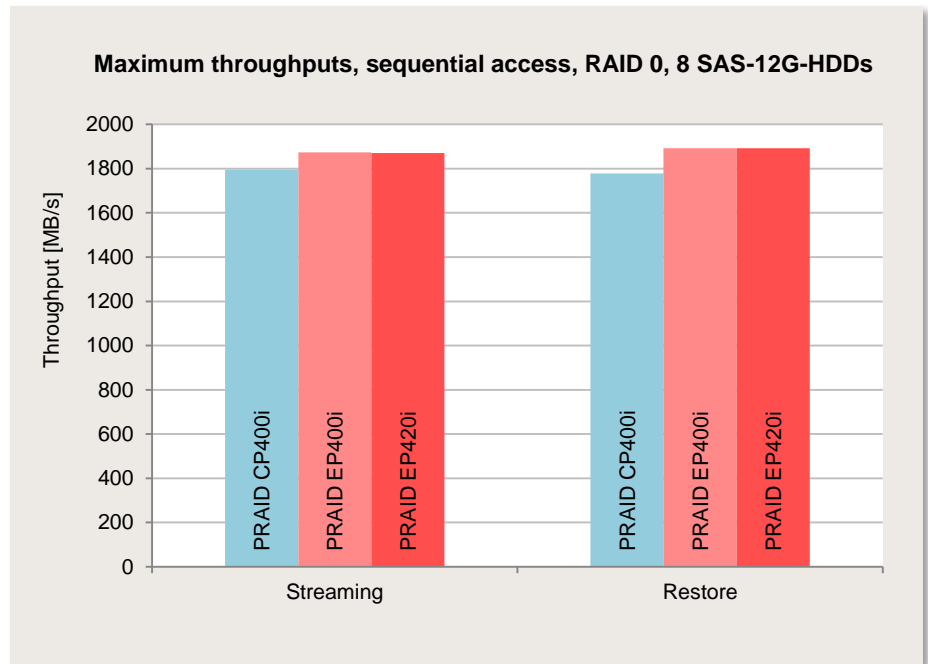
## Sequential accesses

### HDDs

#### RAID 0 with eight SAS-12G-HDDs

The next diagram shows the throughputs of the logical drive of type RAID 0 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

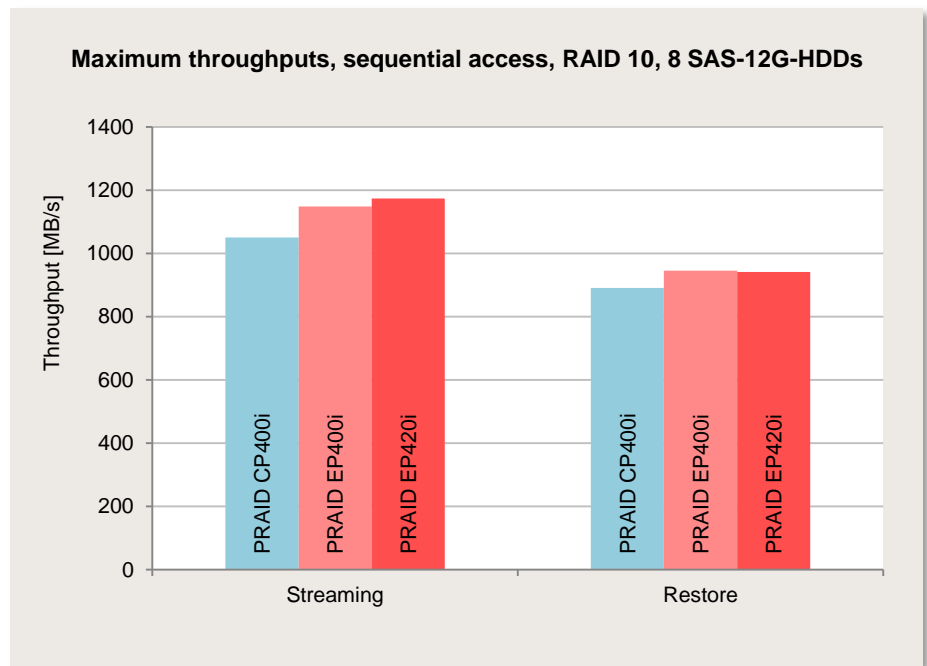
The throughputs here are clearly limited by HDD type and number.



### RAID 10 with eight SAS-12G-HDDs

The next diagram shows the throughputs of the logical drive of type RAID 10 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

For sequential read and write all three of the controllers under consideration for this logical drive achieve or exceed a throughput of approximately four times the maximum throughput of a single HDD (i.e. about 940 MB/s in this case).



**RAID 5 with eight SAS-12G-HDDs**

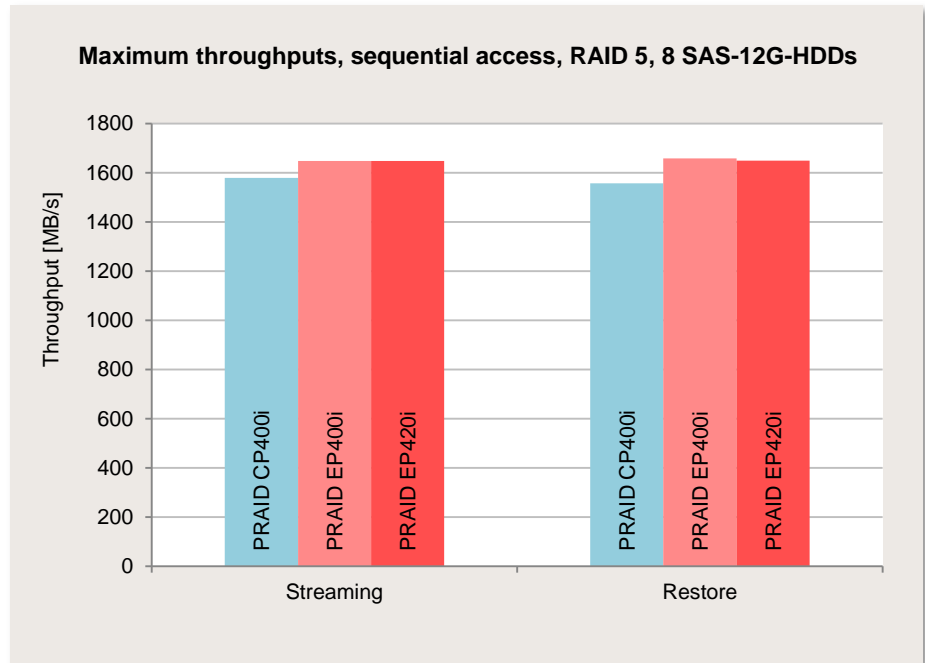
The next diagram represents the throughputs of the logical drive of type RAID 5 that can be achieved with the various controllers for sequential load profiles. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

For logical drives of type RAID 5 consisting of N HDDs the rule of estimation holds that the usable throughput is at most (N-1) times the maximum sequential throughput of the HDD.

In the case viewed here no controller threshold values are yet effective. Thus, the maximum throughput based on the HDD number and type could be estimated as:

$$7 \times 237 \text{ MB/s} = 1659 \text{ MB/s}$$

The diagram confirms this.



**Explanation for the rule of estimation:**

Consider any one of the HDDs from which the logical drive is built. For this HDD precisely one of a number (N) of successive stripes is a parity stripe. These parity stripes are data areas without any user data and in the case of HDDs reduce the usable data throughput for both read and write.

When reading from the HDD, the parity block is simply ignored. Since the write/read head requires a certain amount of time to move over the parity block (due to spindle rotation), the usable data throughput is reduced precisely according to this proportion of time. When writing, it is necessary in the case of N stripes for precisely one parity stripe to be written. Also in this case therefore the amount, by which the usable data throughput is smaller, is the share of the parity stripes. The sequential throughput of actual data that can be supplied by this HDD during both read and write can thus be at most only a share (N-1)/N of the maximum sequential throughput of the HDD. Therefore, all the N HDDs of the logical drive can at most only provide a throughput of actual data that is (N-1) times the maximum sequential throughput of the HDD.

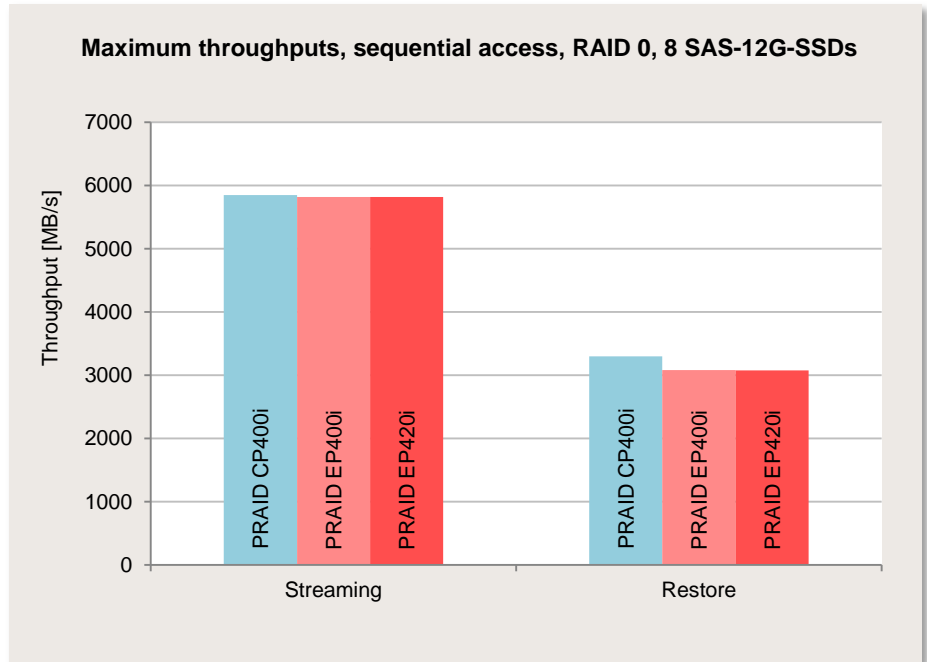
When reading SSDs, the stripes which consist of user data can be directly addressed. No loss in time occurs from the fact that a write/read head would have to move over and ignore a stripe of parity data. However, when writing on SSDs the usable data throughput is also accordingly smaller, because in addition to the user data a time interval is needed for the writing of the parity stripe.

**SSDs**

*RAID 0 with eight SAS-12G-SSDs*

The next diagram shows the throughputs of the logical drive of type RAID 0 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

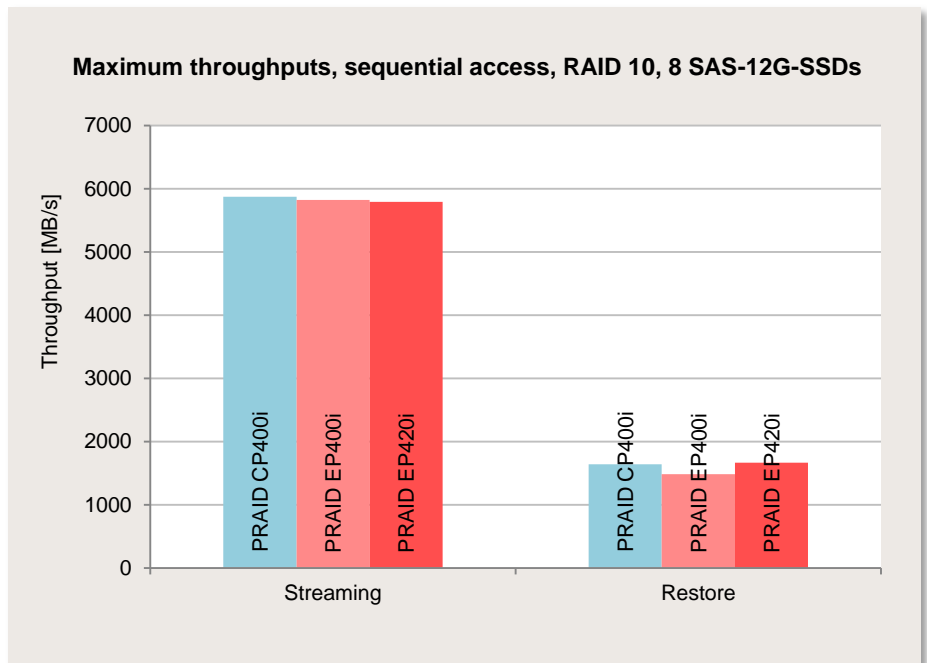
For “Streaming” the throughput limit of the controllers in reading direction (approx. 5900 MB/s) is achieved with eight SAS-12G-SSDs as a RAID 0. The maximum performance of the SSDs and not that of the controllers is decisive here for “Restore”.



*RAID 10 with eight SAS-12G-SSDs*

The next diagram shows the throughputs of the logical drive of type RAID 10 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

For “Streaming” the throughput limit of the controllers in reading direction (approx. 5900 MB/s) is also in this case the limiting factor with RAID 10.

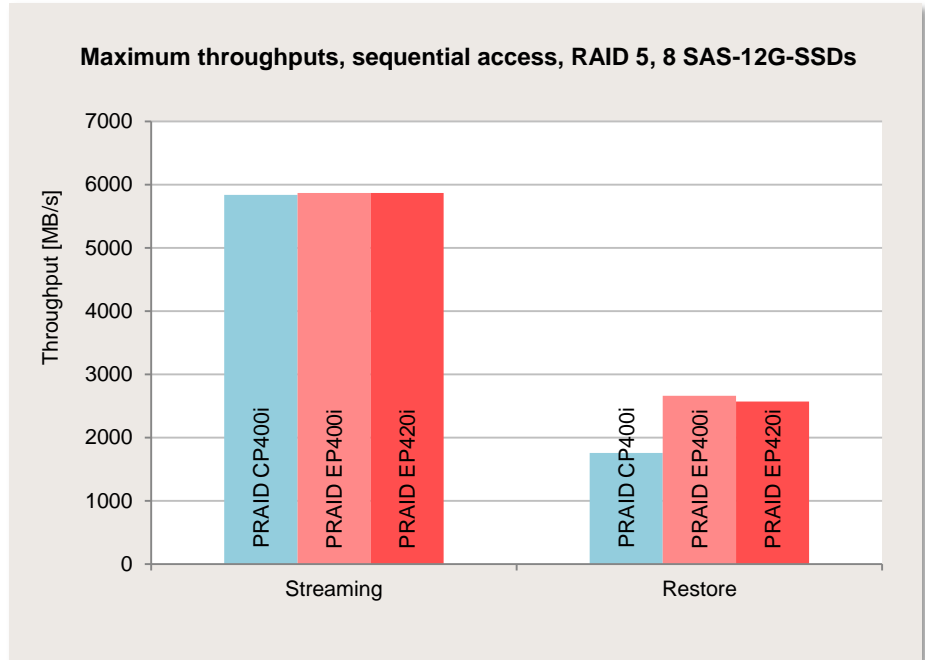


**RAID 5 with eight SAS-12G-SSDs**

The next diagram represents the throughputs of the logical drive of type RAID 5 that can be achieved with the various controllers for sequential load profiles. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

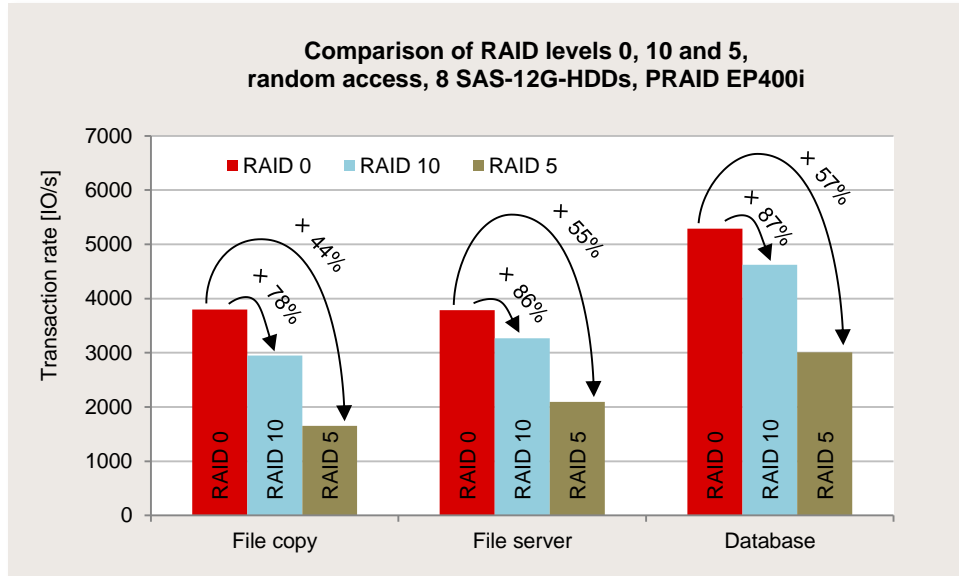
For “Streaming” the throughput limit of the controllers in reading direction (approx. 5900 MB/s) is also in this case the limiting factor with RAID 5.

In order to achieve the maximum data throughputs here with RAID 5 – despite the SSDs – “Performance” mode is required in the ServerView RAID Manager. The enabled write cache of the controllers is vital here in order to achieve the maximum data throughput for “Restore”.



### Random accesses to HDDs: Interrelations between RAID levels

In the case of random accesses to logical drives consisting of HDDs it is possible for a given load profile to estimate the maximum transaction rate for another RAID level by multiplying the maximum transaction rate for RAID 0 with a suitable factor. In this case, the HDD number and type, stripe size and controller must be identical. First of all, these interrelations should be illustrated in the following diagram on the basis of measurement values.



These percentages can also be theoretically estimated if you use a multiplication factor for random write accesses. This is a matter of a so-called “write penalty”, which is defined as:

$$\frac{\text{\# of accesses caused from the viewpoint of all the physical hard disks}}{\text{\# of causing write accesses from the viewpoint of the application}}$$

This “write penalty”<sup>1</sup> has a value of 1 for RAID 0, 2 for RAID 10 and a value of 4 for RAID 5. Taken together with the read share (which does not multiply) contained in the respective load profile, the result is a specific multiplication factor between the accesses from the viewpoint of the application and the accesses from the viewpoint of all the hard disks. For example with RAID 5 compared to RAID 0, this factor causes the hard disks to already come under maximum load from an application viewpoint with much lower transaction rates. These theoretical percentage differences between the various RAID levels are listed in following table for the three random standard load profiles (and thus ultimately write shares).

RAID levels compared	Theoretical ratio of maximum transaction rates for load profile		
	File copy (50% write)	File server (33% write)	Database (33% write)
RAID 10 / RAID 0	67%	75%	75%
RAID 5 / RAID 0	40%	50%	50%

In a comparison of these theoretical values with the percentages from the above diagram, which result from the measurement values, you will find that the percentages in the diagram are somewhat higher. This is due to optimization measures of the controllers through cache usage.

<sup>1</sup> In the case of RAID 10 the value 2 expresses the double writing of each data block on account of disk mirroring. And in the case of RAID 5 one write access must take place as follows on account of the random load profile: 1) Read the old data stripe; 2) Read the old parity stripe; 3) Calculate the new parity stripe from the read stripes; 4) Write the new data stripe; 5) Write the new parity stripe. Thus, in total the random writing of a data stripe means read twice and write twice. So this is why value 4 is for the “write penalty”.

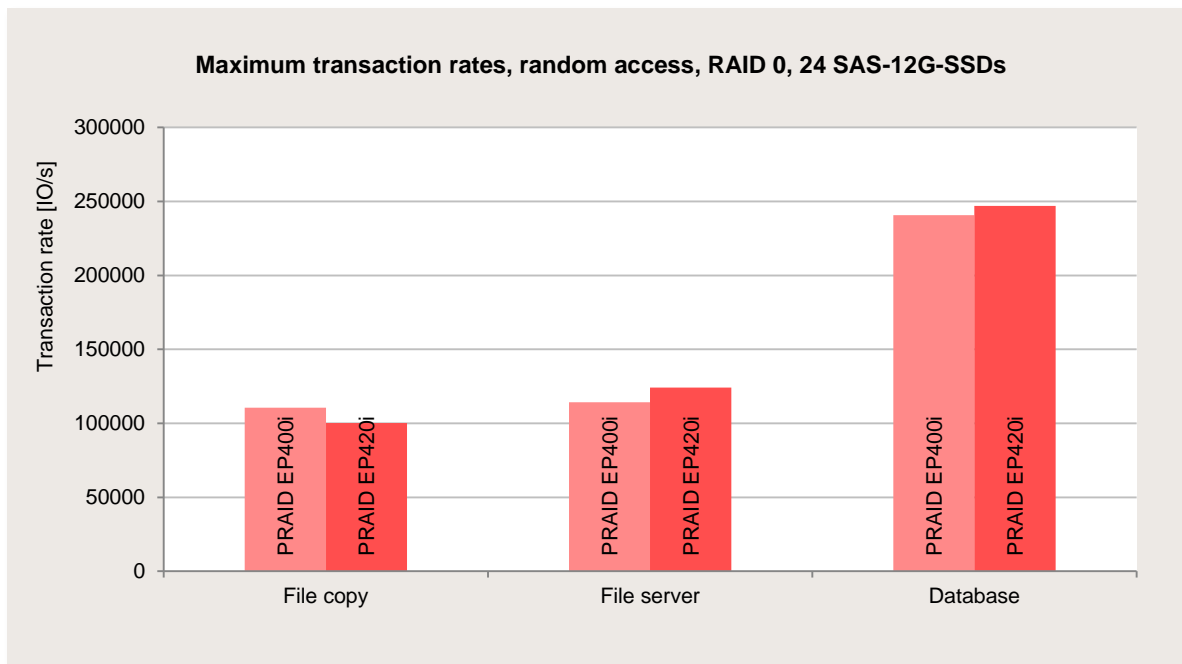
## RAID 0, 10 and 5 (more than eight SAS-SSDs)

The two controllers with the top performance are available for the operation of more than eight hard disks in the current PRIMERGY servers: the PRAID EP400i and PRAID EP420i. They currently enable server configurations with up to 32 internal hard disks (e. g. in the PRIMERGY RX2560 M1). The list below for the five standard load profiles is to include the maximum throughputs and transaction rates of the controllers for various RAID levels. To this end, measurements with high numbers of SSDs are to be used: In the case of RAID 0 and RAID 5 the measurements are carried out with 24 SSDs, and for RAID 10 with 16 SSDs (for RAID 10 the latter is the current upper limit for a single logical drive). As a result, the range of large numbers of hard disks is covered on a representative basis. As was also the case in the previous subsection, measurements with high-performance SAS-12G-SSDs are used as a means of illustration. These hard disks are described in more detail in the section [“Measurement environment”](#).

### Random accesses

#### RAID 0 with 24 SAS-12G-SSDs

The next diagram shows the transaction rates of the logical drive of type RAID 0 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



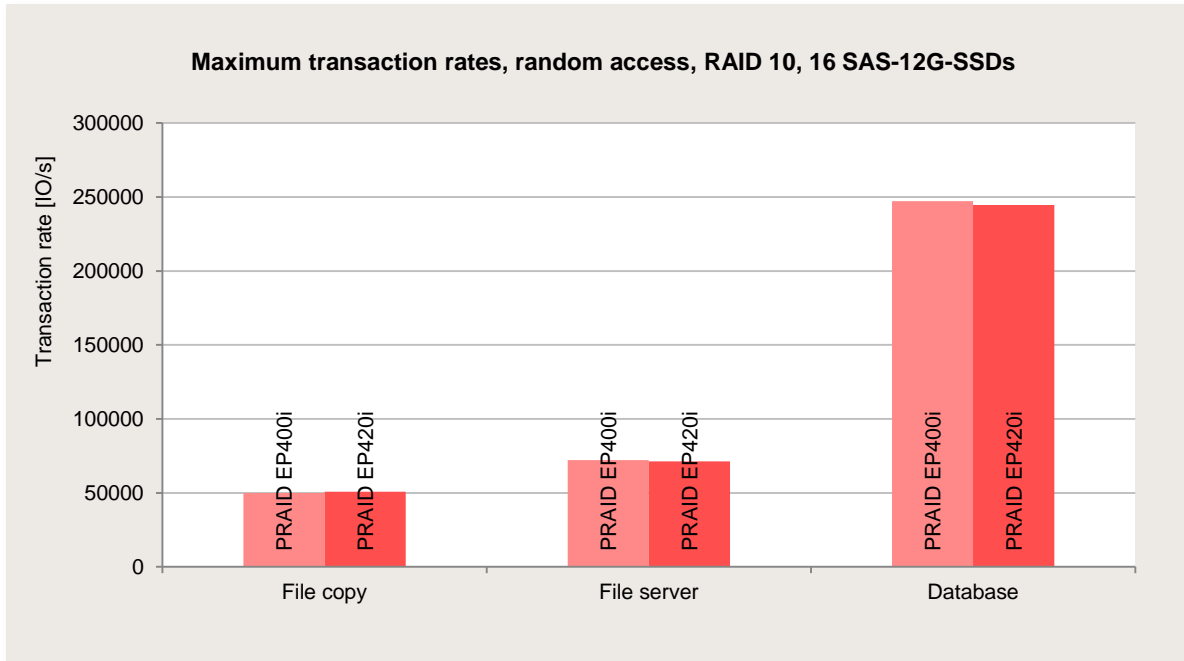
The most important information received here is the very high transaction rate of about 240000 IO/s that can be achieved for the load profile with a small block size (“Database”). The impact of the FastPath option, which is enabled as standard in the latest controller firmware, is especially apparent for this load profile.

Expressed in the form of SAS-12G-SSD numbers: In order to make full use of the possibilities of the PRAID EP400i for RAID 0 it is necessary to have between five (8 kB block size) and 17 (64 kB block size) fully loaded SAS-12G-SSDs - depending on the random load profile. With a smaller SSD load or other SSD types these numbers must be suitably modified.

It is also interesting for us to recognize the throughput values that result through conversion from these transaction rates. Despite the lower transaction rates, the two load profiles with a 64 kB block size have the higher throughputs. For example, the PRAID EP420i handles a throughput of about 7763 MB/s with the load profile “File server”. This value is remarkable, because it is higher than the two sequential maximum throughputs of the controller for 100% read and 100% write with this RAID level. This value would not have been reached without real bidirectional use of the SAS connections.

### RAID 10 with 16 SAS-12G-SSDs

The next diagram shows the transaction rates of the logical drive of type RAID 10 for random load profiles that can be achieved with various controllers. The three groups of columns show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



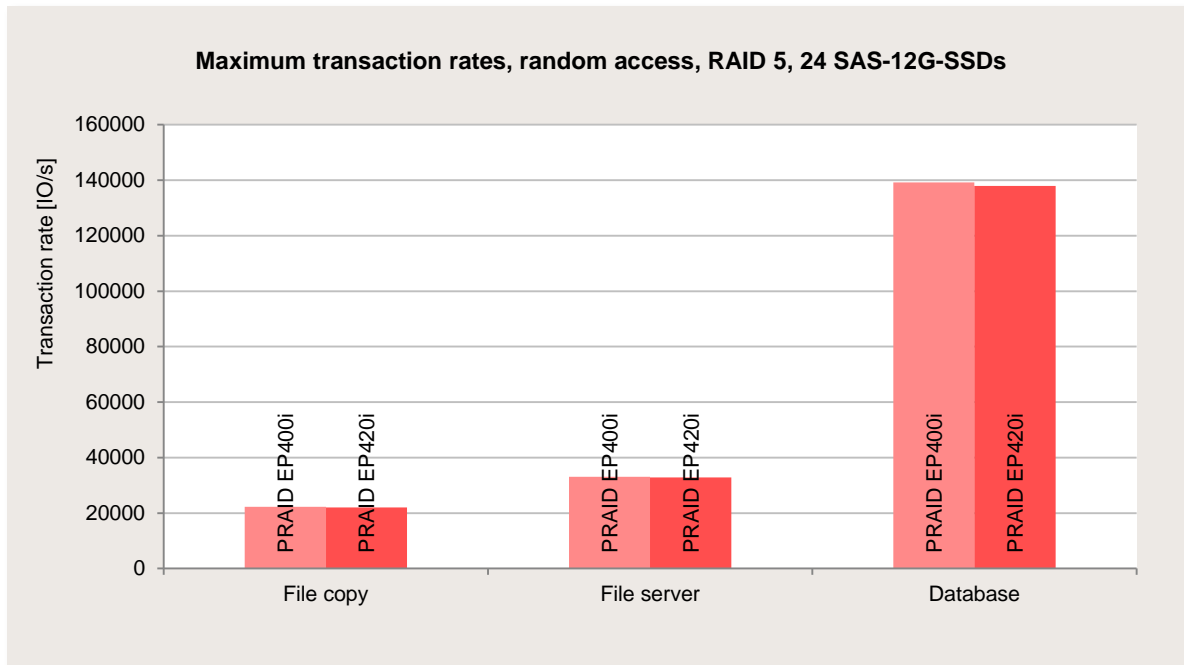
**Remark:**

The configurations consisting of numerous SAS-12G-SSDs as described here can be used to achieve transaction rates of several hundred thousand IO/s with small block sizes ( $\leq 8$  kB). Handling so many I/Os can utilize the processing CPU core at almost 100% of its capacity. As a result, the real frequency of the server CPU can become the limiting factor. A CPU of medium nominal frequency (Xeon E5-2660 v3 @ 2.60 GHz) and BIOS settings for optimal performance were used to obtain measurement results that are also valid for an average CPU configuration and a wide selection of server models. In Xeon E5-2600 v4 based servers the transaction rates presented here are also achieved for example with the processor type Xeon E5-2623 v4 @ 2.60 GHz. If a CPU with optimal frequency is used, it is possible to clearly surpass the transaction rates presented here. For example, it is possible with a CPU Xeon E5 2637 v3 @ 3.50 GHz or Xeon E5 2637 v4 @ 3.50 GHz to achieve more than 300000 IO/s instead of about 250000 IO/s for a logical drive of type [RAID 0 consisting of 24 SAS-SSDs](#) for the load profile “Database” (random access, 67% read, 8 kB block size).



**RAID 5 with 24 SAS-12G-SSDs**

The following diagram lists the maximum transaction rates for the various controllers using the example of a logical drive consisting of 24 SAS-12G-SSDs. The three groups of columns in the diagram show the transaction rates for the standard load profiles “File copy” (random access, 50% read, 64 kB block size), “File server” (random access, 67% read, 64 kB block size) and “Database” (random access, 67% read, 8 kB block size).



To express this in numbers of SAS-12G-SSDs, it is necessary to have between seven (8 kB block size) and 17 (64 kB block size) fully loaded SAS-12G-SSDs in order to make full use of the possibilities of the PRAID EP400i for RAID 5 - depending on the random load profile. With a smaller SSD load or other SSD types these numbers must be suitably modified.

## Sequential accesses

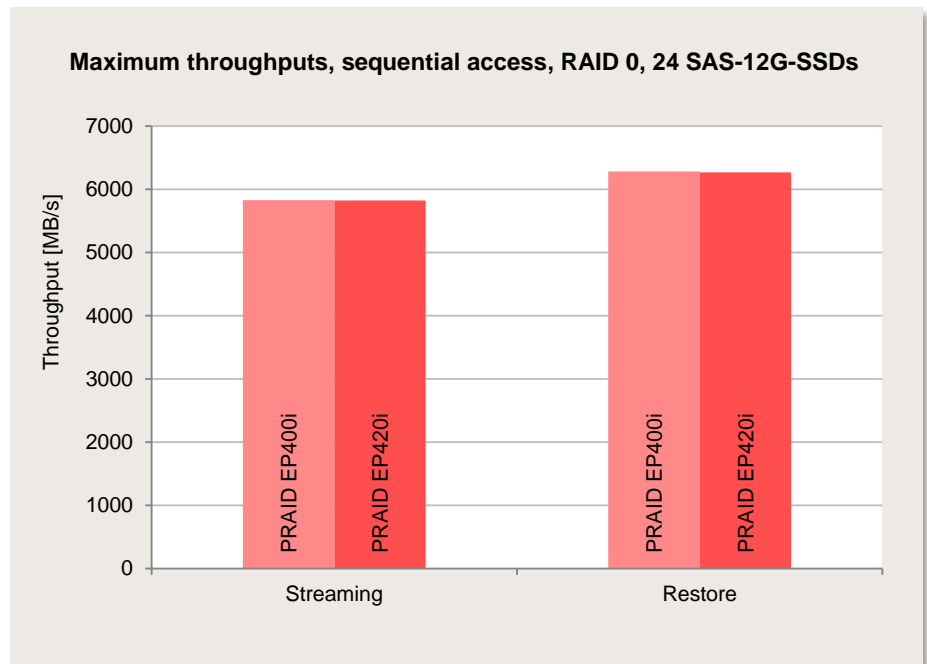
Generally applicable statements about the controllers are listed below based on measurements with 24 (or 16) SAS-12G-SSDs. It is possible to calculate the anticipated maximum throughputs for other hard disk types and -numbers by appropriately multiplying the basic performance values of the hard disk. If the throughput calculated in this way exceeds the threshold value of the controller, the controller threshold value becomes effective.

### RAID 0 with 24 SAS-12G-SSDs

The next diagram shows the throughputs of the logical drive of type RAID 0 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

The data throughputs presented here for “Restore” (approx. 6280 MB/s) are the limit of the controllers used here for RAID 0.

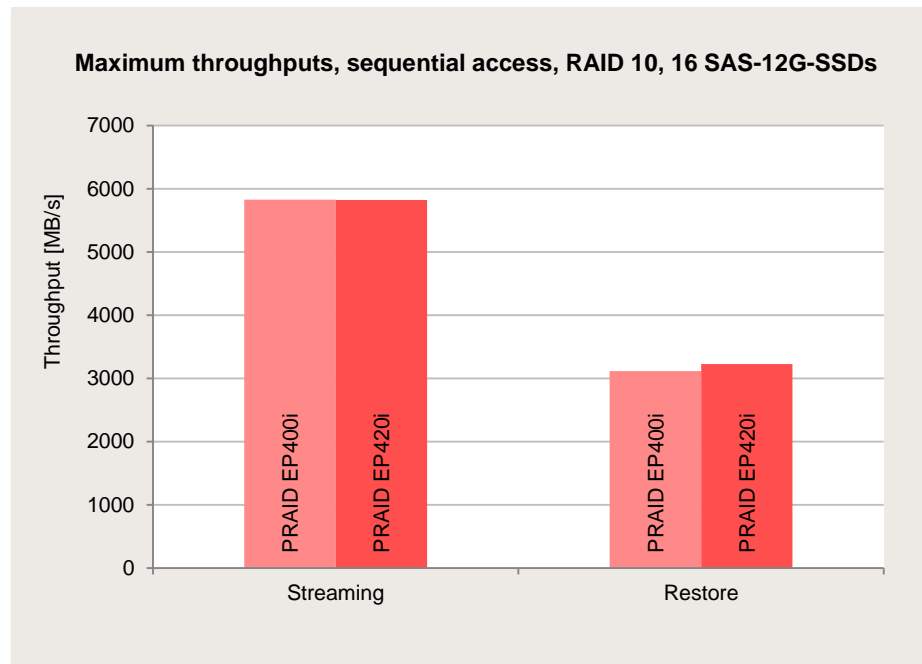
The limit for “Streaming” (approx. 5900 MB/s) was already achieved with eight SAS-12G-SSDs as a RAID 0.



### RAID 10 with 16 SAS-12G-SSDs

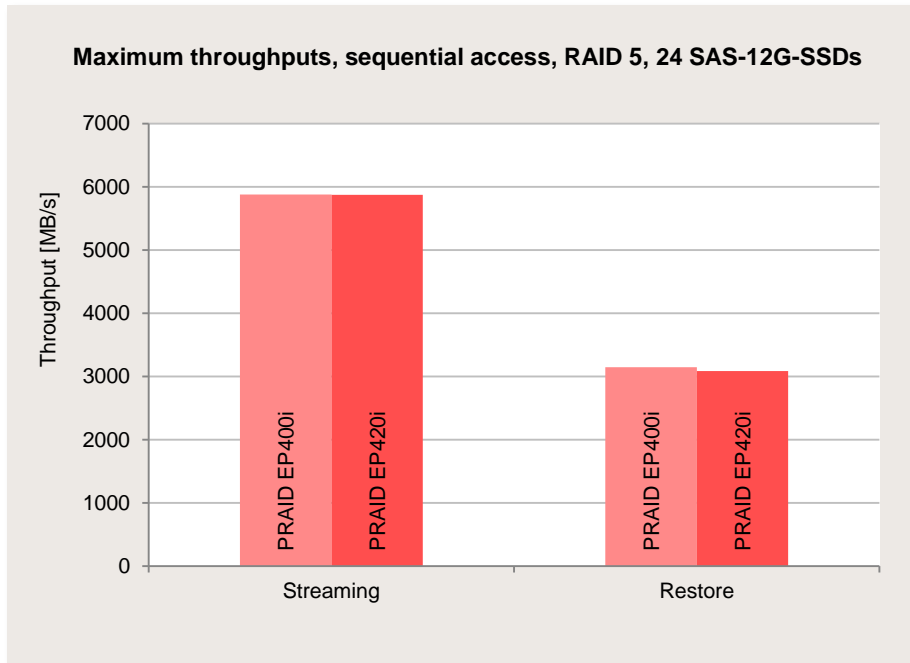
The next diagram shows the throughputs of the logical drive of type RAID 10 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).

The limit for “Streaming” (approx. 5900 MB/s) was already achieved – as with RAID 0 – with eight SAS-12G-SSDs as a RAID 10.



**RAID 5 with 24 SAS-12G-SSDs**

The next diagram shows the throughputs of the logical drive of type RAID 5 for sequential load profiles that can be achieved with various controllers. The two groups of columns in the diagram show the throughputs for the standard load profiles “Streaming” (sequential access, 100% read, 64 kB block size) and “Restore” (sequential access, 100% write, 64 kB block size).



The data throughputs presented here for “Restore” (approx. 3100 MB/s) and “Streaming” (approx. 5900 MB/s) are the limits of the controllers used here for RAID 5.

In the case of RAID 5 this maximum value for sequential write is a significant indicator for the performance of a RAID controller, as the speed of the controller is reflected in a relatively undistorted way here in the calculation of the parity blocks.

In order to achieve the maximum data throughputs here with RAID 5 – despite the SSDs – “Performance” mode is required in the ServerView RAID Manager. The enabled write cache of the controllers is vital here in order to achieve the maximum data throughput for “Restore”.

## Influence of the size of the controller cache

In principle the larger controller cache of the PRAID EP420i (2 GB) in comparison with the PRAID EP400i (1 GB) has performance advantages with logical drives consisting of HDDs that are used for random load profiles with a high write share. In other use cases with HDDs and generally with SSDs the size of the controller cache doesn't have a significant influence.

The advantage of the larger cache is most distinct if a logical drive of the type RAID 5 consisting of HDDs is accessed with small blocks ( $\leq 8$  kB), randomly, solely with writes (0% random read). Here the larger cache typically provides an advantage of 20% to 30%. Such load profiles can be expected in practice less frequently; for example data restores with a high number of small files belong to these. Most of the application scenarios for servers contain a read share of at least 50% (database servers, mail servers, web servers). For these cases the performance advantage is less.

Overall, the performance advantage of the larger cache in case of random access to a logical drive consisting of HDDs depends on the following factors:

- The higher the write share, the higher the advantage
- Stronger with RAID 5 than with RAID 10
- Stronger with low HDD number than with high HDD number
- Stronger with small blocks than with large blocks
- The higher the load, the stronger the advantage

To illustrate how the percentage of the performance advantage depends on the essential factors the following comparison table looks at a suitable selection of load profiles and load intensities. The percentages contained herein are typical values which have been calculated by averaging the measuring values for a wide selection of HDD types.

Load profile	Load intensity	Performance advantage (percentage of IO/s) of PRAID EP420i over PRAID EP400i	
		RAID 5 (up to 8 HDDs)	RAID 10 (up to 8 HDDs)
67% random read, 4 kB	Low load up to acceptable high load	6-11%	<6%
	Overload	6-11%	6-11%
67% random read, 64 kB	Low load up to acceptable high load	<6%	<6%
	Overload	6-11%	6-11%
0% random read, 4 kB	Low load up to acceptable high load	20-30%	11-20%
	Overload	20-30%	11-20%
0% random read, 64 kB	Low load up to acceptable high load	<6%	<6%
	Overload	6-11%	6-11%

### Remarks:

- In the cases under consideration here the disk cache was "Disabled". The percentages for the setting "Enabled" are in each case very similar to the corresponding percentages for "Disabled"
- The ranges of load intensities mentioned here have been modeled within the measurements as follows: "Low load up to acceptable high load" corresponds to 1-32 outstanding I/Os, "Overload" corresponds to 64-512 outstanding I/Os
- RAID 10 also stands as an example for RAID 0 and RAID 1 (RAID levels without parity calculation); RAID 5 also stands as an example for RAID 6, RAID 50 and RAID 60 (RAID levels with parity calculation)

## Lower load levels

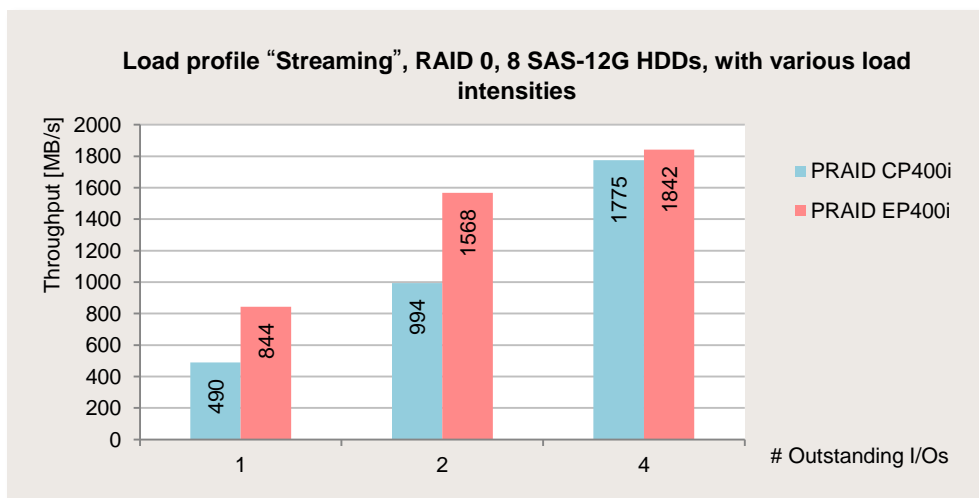
The differences between the controllers can not only be seen in the achievable maximum values under high load, but also under a lower load. Thus, the following table compiles for the individual controllers the maximum values for the sequential throughputs of a single application, which works without parallelism for its accesses (referred to below with “1 outstanding IO”). These maximum values apply for all the RAID levels that are supported by the controllers and cannot be exceeded. These throughputs can only be achieved with RAID 0; and with the other RAID levels the throughputs usually remain below the maximum values.

Controller	Maximum values of sequential throughput, 1 outstanding IO, 64 kB block size	
	100% read (load profile “Streaming”)	100% write (load profile “Restore”)
PRAID CP400i	490 MB/s	315 MB/s
PRAID EP400i	1300 MB/s	1375 MB/s
PRAID EP420i	1300 MB/s	1375 MB/s

The “Read-ahead” setting is prerequisite to achieving these values for sequential read, as is the setting “Write-back” for sequential write. These maximum throughputs also depend very much on the block size, whereby the interrelations of the table values as regards size are similar for other block sizes.

The differences between the controllers in the table become significant at the latest when the logical drive used is as a matter of principle in a position to enable more than 500 MB/s of sequential throughput for “1 outstanding IO”. In such cases, an inappropriately selected controller can have the effect of a restriction.

The following example illustrates this on the basis of throughput measurements with a logical drive of type RAID 0 consisting of eight SAS-12G-HDDs for the load profile “Streaming” (sequential access, 100% read, 64 kB block size). The comparison is made between the PRAID CP400i and the PRAID EP400i with differing numbers of parallel accesses (“# Outstanding IOs”).



You can clearly see that in this case the PRAID CP400i does not achieve the same throughput for 1 and 2 outstanding IOs as the PRAID EP400i. The latter already achieves a throughput of 844 MB/s, whereas the PRAID CP400i controller only achieves a little more than half.

From the viewpoint of response times, this means that it is possible for low load intensities to approximately halve the response times with the PRAID EP400i compared with the PRAID CP400i.

## Conclusion

The PRIMERGY and PRIMEQUEST servers use the “Modular RAID” concept to offer a plethora of opportunities to meet the requirements of various application scenarios.

An onboard controller is a low-priced entry-level alternative for the RAID levels 0, 1 and 10, which saves one PCIe slot but is restricted to six hard disks. The pro rata consumption of the server's processor performance is increasingly less important in newer servers.

On the SATA side the current onboard controllers support the standards up to frequency 6G.

In the case of PCIe controllers the current generation supports the standard SAS-12G. As a result, the maximum real data throughputs were increased from 3800 MB/s to 6280 MB/s compared with the predecessor generation.

The PRAID CP400i is the PCIe controller suited for average requirements. This controller does not have a cache, permits up to eight hard disks and supports the RAID solutions RAID 0, RAID 1, RAID 1E, RAID 10 and – contrary to the predecessor controller – also RAID 5.

The PRAID EP400i and PRAID EP420i controllers offer all the current standard RAID solutions RAID 0, RAID 1, RAID 1E, RAID 5, RAID 6, RAID 10, RAID 50 and RAID 60 in the High-End sector. These controllers have a controller cache and can as an option be backed up using an FBU. Manifold options to set the use of the cache make it possible to flexibly adapt the controller performance to suit the RAID level used. A further optimization option here is the adjustable stripe size. In many application scenarios, for example if random accesses take place on conventional hard disks with a high load intensity, these controllers enable a 75% higher transaction rate than the PRAID CP400i (example: RAID 0 with four SATA-6G-HDDs, random access, 50% read, 64 kB block size).

The RAID controllers PRAID EP400i and PRAID EP420i only differ as far as cache size is concerned. The first controller has a 1 GB cache, and the second one has a 2 GB cache. The larger cache is recommended for HDDs that are used for random load profiles with a high write share.

The majority of the application scenarios that put a load on the disk subsystem come along with a random read / write access. If SSDs are used to manage very high IO rates, the controller has considerable influence on the maximum transaction rate. In the case of a logical drive of type RAID 0 and typical accesses for a database (67% read, random, block size 8 kB) the PRAID CP400i for example permits up to 164000 IO/s, and the PRAID EP420i on the other hand permits up to 248000 IO/s, in other words 1.5 times that amount. The differences are particularly large for a logical drive of type RAID 5: the PRAID CP400i achieves up to 28700 IO/s for typical accesses for a database, while the PRAID EP420i has up to 133000 IO/s, in other words about 4.6 times that amount. Thus, in the case of RAID 5 in conjunction with SSDs it is imperative to choose the PRAID EP400i or the PRAID EP420i.

Regardless of the hard disk type, the various controllers each have maximum sequential throughputs that are specific to the RAID level and the load profile. These maximum values have in part increased substantially in comparison to the predecessor generation. In the case of sequential write on a logical drive of type RAID 5 the PRAID EP420i for example achieves approx. 3100 MB/s, whereas the predecessor controller only achieved approx. 2200 MB/s.

If a higher transaction rate or higher throughput is required for the planned application scenario than a single controller can provide, two controllers can be used. A number of PRIMERGY servers provide this option (e. g. PRIMERGY RX2540 M1).

A further aspect of faster controllers with sequential access profiles is the increased throughput that is already achieved with low access parallelism. If the logical drive is efficient enough, it means that more than 1300 MB/s is possible for read and write with the PRAID EP400i in this special application. Compared with controllers of the previous generation this also means a significant increase in the maximum throughput for these special cases.

The RAID-Manager software “ServerView RAID Manager” that is supplied for PRIMERGY servers is recommended for the configuration of controllers and hard disks. This utility program makes it possible to conveniently adapt controller and hard disk settings to meet customer requirements regarding performance and data security in a controller-independent way for the majority of the application scenarios. If FBUs and UPSs are used as buffers in the case of power failures, maximum performance can be reconciled with data security.

## Literature


### PRIMERGY & PRIMEQUEST Servers


<http://www.fujitsu.com/fts/products/computing/servers/>


### Performance of Server Components

<http://www.fujitsu.com/fts/products/computing/servers/mission-critical/benchmarks/x86-components.html>

This White Paper:

 <http://docs.ts.fujitsu.com/dl.aspx?id=9845be50-7d4f-4ef7-ac61-bbde399c1014>

 <http://docs.ts.fujitsu.com/dl.aspx?id=7826d783-bc71-4cd7-8486-d74f4dc2509c>

 <http://docs.ts.fujitsu.com/dl.aspx?id=3075886a-3c79-4b5b-8d9f-e9269e083bef>

### BIOS optimizations for Xeon E5-2600 v4 based systems

<http://docs.ts.fujitsu.com/dl.aspx?id=eb90c352-8d98-4f5a-9eed-b5aade5ccae1>

### BIOS optimizations for Xeon E5-2600 v3 based systems

<http://docs.ts.fujitsu.com/dl.aspx?id=f154aca6-d799-487c-8411-e5b4e558c88b>

### RAID-Controller-Performance 2013 (previous white paper)

<http://docs.ts.fujitsu.com/dl.aspx?id=e2489893-cab7-44f6-bff2-7aeea97c5aef>

### 512e HDDs: Technology, Performance, Configurations

<http://docs.ts.fujitsu.com/dl.aspx?id=f5550c48-d4db-47f6-ab9d-ce135eaac81>

### Basics of Disk I/O Performance

<http://docs.ts.fujitsu.com/dl.aspx?id=65781a00-556f-4a98-90a7-7022feacc602>

### Information about Iometer

<http://www.iometer.org>

## Contact

### FUJITSU

Website: <http://www.fujitsu.com/>

### PRIMERGY & PRIMEQUEST Product Marketing

<mailto:Primergy-PM@ts.fujitsu.com>

### PRIMERGY Performance and Benchmarks

<mailto:primergy.benchmark@ts.fujitsu.com>