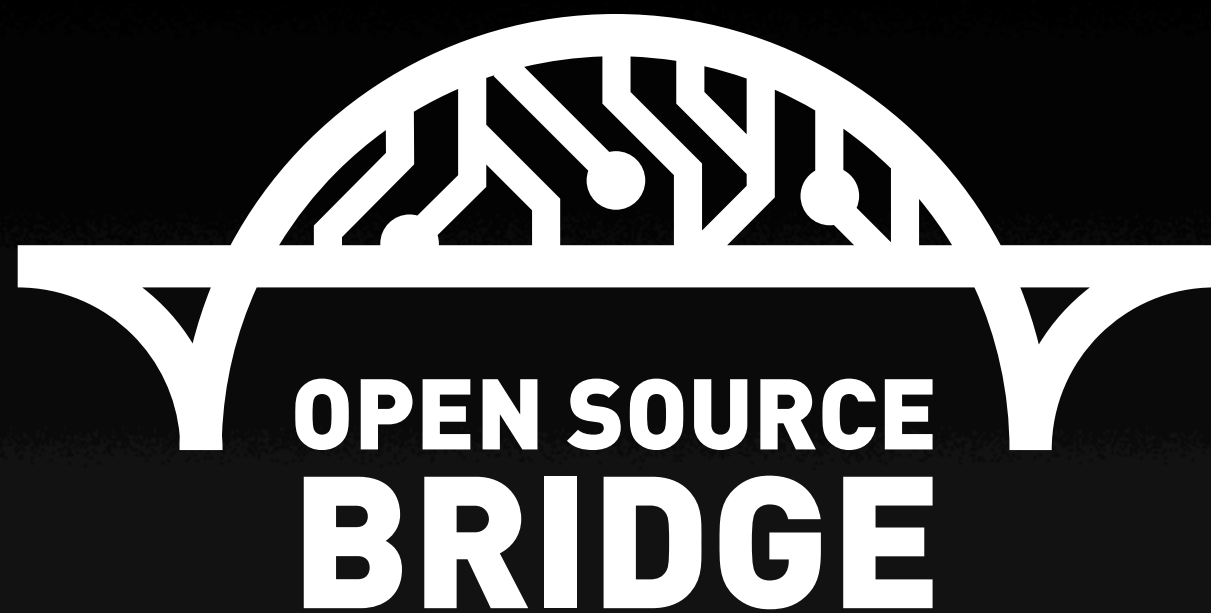


Linux Filesystems Performance for Databases

Portland PostgreSQL Performance Pad

Selena Deckelmann
selena@postgresql.org
PostgreSQL Global Development Group
twitter: @selenamarie



Do filesystems do
what we expect?

We are volunteers.

We think you should run
these tests.

We are:
DBAs
Sysadmins
Performance tuners

How will this hardware perform?

How will this filesystem perform?

Why should you care about
filesystem-specific
performance?

Expectations

Where to start?

The Defaults.

HOW TO START A FIGHT



PERFORMANCE

HOW TO START A FIGHT

Not addressing reliability

~~PERFORMANCE~~

HOW TO START A FIGHT

Very Narrow Use Case:
A Relational Database

~~PERFORMANCE~~

HOW TO START A FIGHT

Need for periodic testing.
(And we've got some
hardware!)

~~PERFORMANCE~~

HOW TO START A FIGHT

- ★Kernel differences
- ★FS patch-level differences
- ★Mount options
- ★mkfs options

PERFORMANCE

HOW TO START A FIGHT

Focused on
THROUGHPUT

(Because that's what people who
buy large systems look for)

~~PERFORMANCE~~

HOW TO START A FIGHT

Later:
Response Time
Operations per second

~~PERFORMANCE~~



No, we will not
be testing ZFS.



BtrFS
(nope, not yet)

What do we expect?

Some conventional wisdom:

“RAID5 is the
worst choice
for a database.”

“LVM incurs
too much overhead
to use.”

“Striping doubles
performance.”

“Turning off 'atime'
is a big
performance gain.”

Replacing Atime With Relatime in the Kernel

Posted by [ScuttleMonkey](#) on Wed Aug 08, 2007 05:08 PM
from the [results-apparently-too-much-to-ask-for](#) dept.

[eldavojohn](#) writes

"Our friend Jeremy at the Kernel Mailing List's [criticism of atime](#) from Linus to improve relatime he noted: 'I don't think we can deal it is in practice. Atime updates would give us more performance deficiency that pagecache speedups of the last 10 years. This is a design idea of all times. Unix has always been a 'For every file that is read from disk, there is a file that is already cached and doesn't need to be read from disk!' Well, I guess I can ex

“Getting rid of atime updates would give us more everyday Linux performance than all the pagecache speedups of the last 10 years, _combined_.”

“Journaling filesystems (ext3) will have worse performance than non-journaling filesystems (ext2).”

“Your read-ahead
buffer
is big enough.”

Now... on to the good stuff.



**INTERNET.
SERIOUS BUSINESS.**

PostgreSQL's Portland Performance Pad



Hosted by CommandPrompt, Inc.

Our machine:

HP ProLiant DL380G5
Smart Array p800

72GB 15,000 RPM SAS (up to 25 disks)
32GB RAM

Linux:
2.6.25-gentoo-r6
*New tests being run with 2.6.28

Our machine:
Chosen because
of it's low, low price.

Thank you, HP.

Our tests:

fio

64 GB working set

8 threads

no fadvise

no direct i/o

8KB blocksize

I/O elevator: deadline

Our tests:

fio

read (sequential, random)

write (sequential, random)

read-write (50/50 mix)

Our stats:

sar

mpstat

iostat

vmstat

readprofile

Our tests:
Chosen because of their
relevance to PostgreSQL

Filesystems Tested:

ext2

ext3

jfs

xfs

reiserfs

ext4 (but had trouble)

Disk configs tested:

Single disk

RAID-0

RAID-1

RAID-5

RAID-10

RAID-6

The Data:
<http://moourl.com/fsperf>

Confessions:

- May be high standard deviation with results (don't know yet!)
- No filesystem tuning, all default create and mount options
- No software raid comparison or lvm (volume management test) for 2.6.28 tests

Confessions:

- Some xfs runs had to be repeated and some ext4 runs did not complete successfully
- Only presenting throughput
- Interested in system performance for a specific application, not code performance

Confessions:

- I/O profiles don't exhibit a time or partition alignment issues
- Disk controller firmware not at the latest version in 2.6.25 tests
- Software RAID is on top of 1 disk RAID 0 devices (HP SmartArray doesn't have JBOD option)

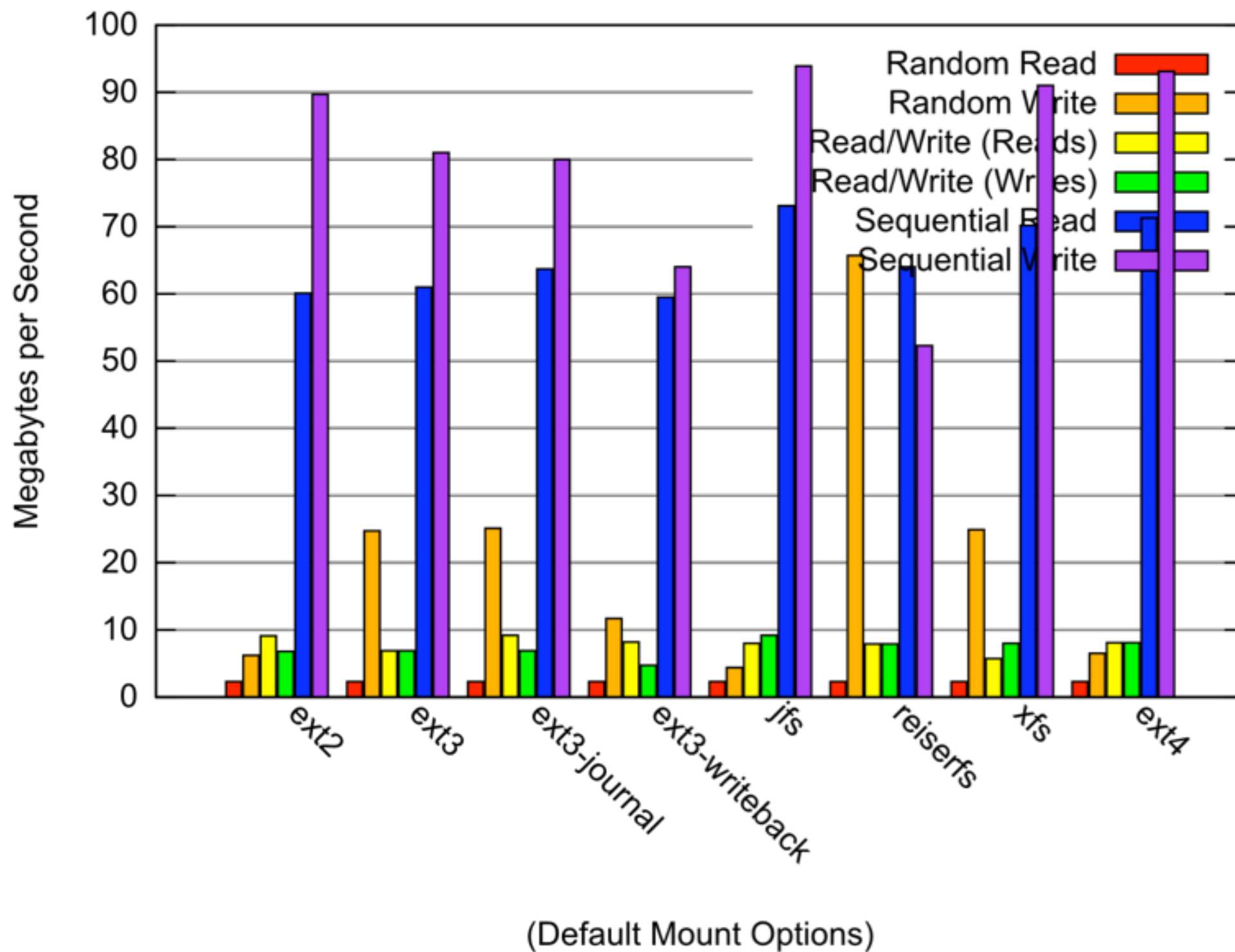
HOW TO START A FIGHT

AUDIENCE PARTICIPATION

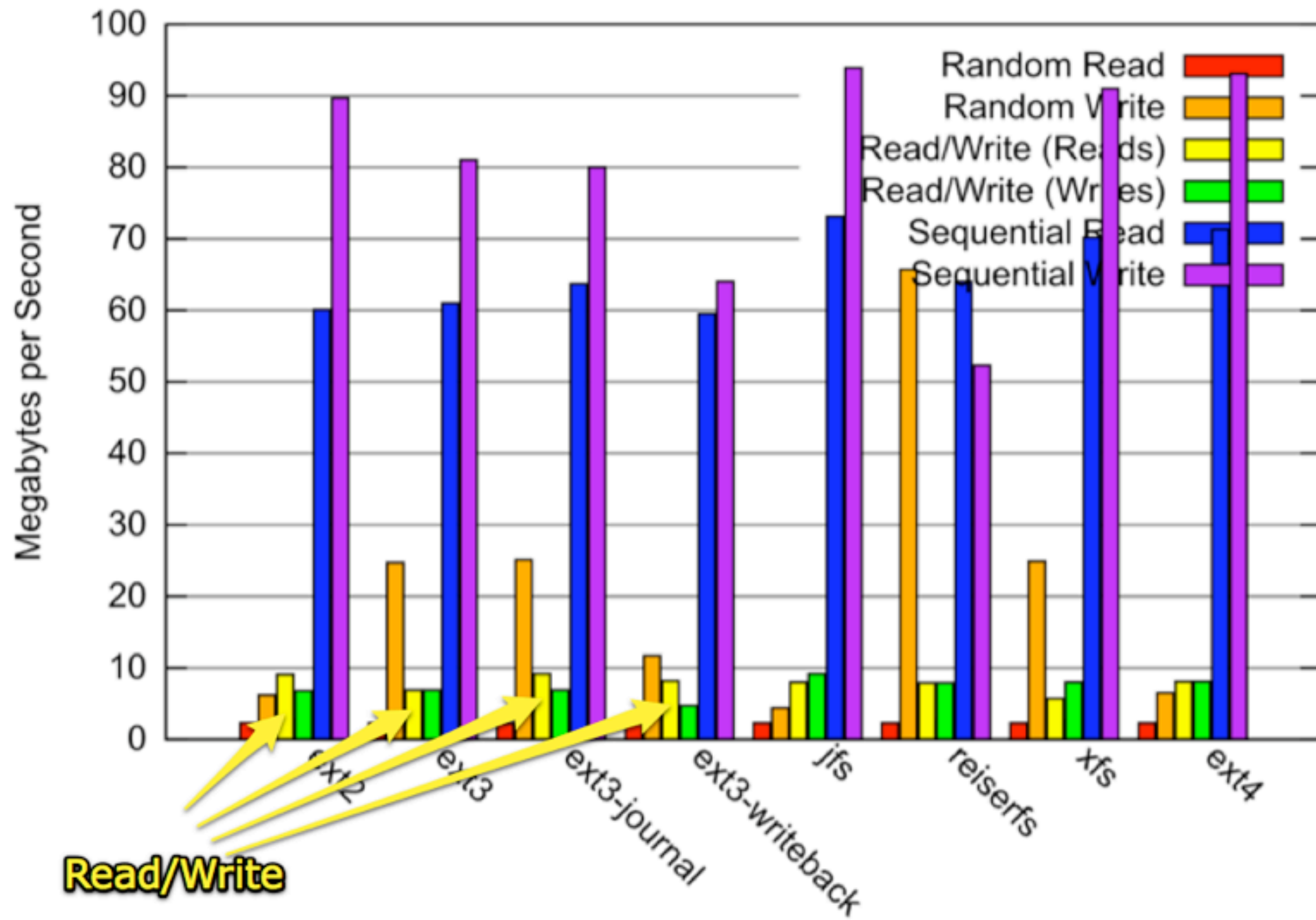
Higher throughput:
ext2 or ext3?

~~PERFORMANCE~~

1 Disk RAID 0

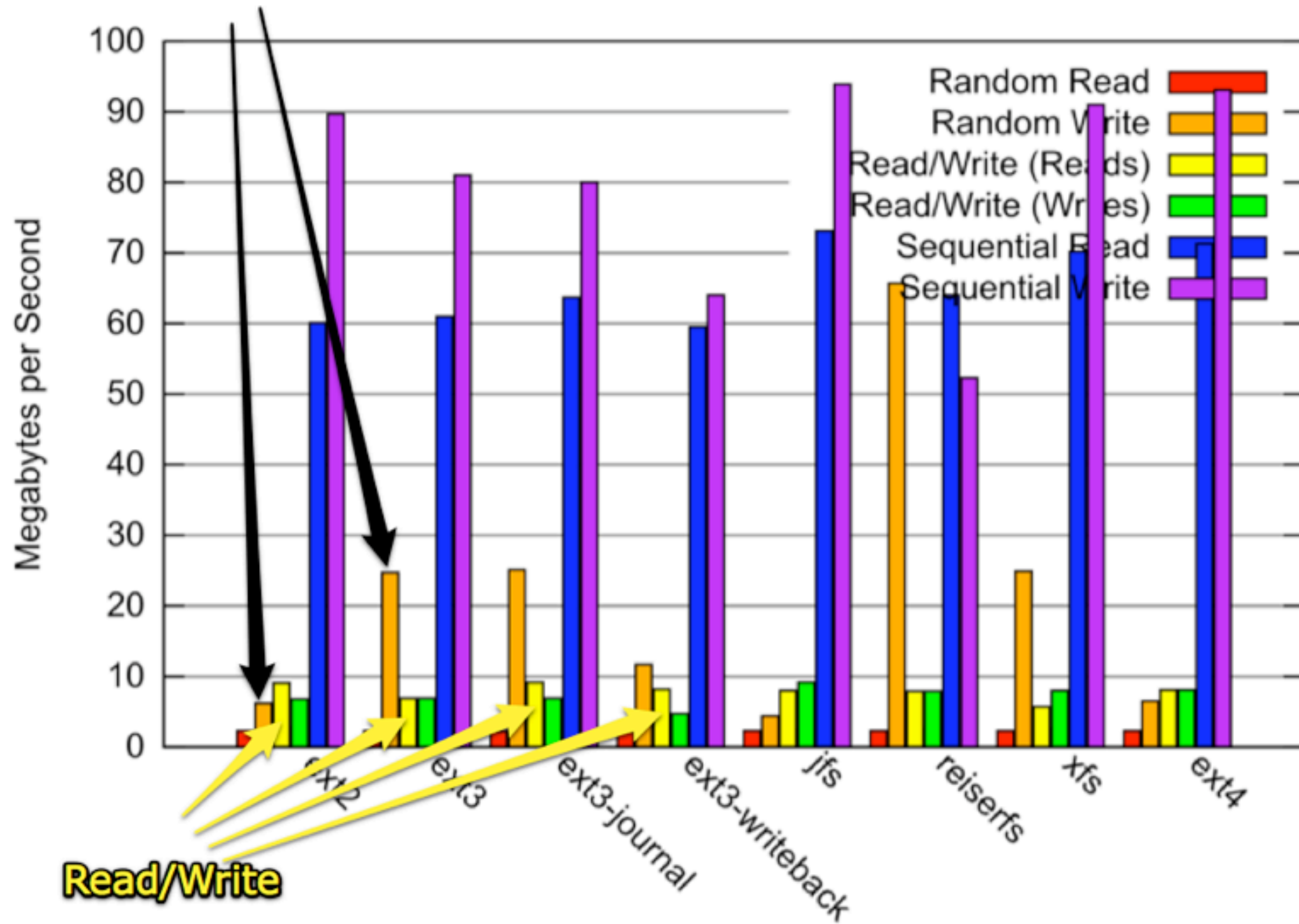


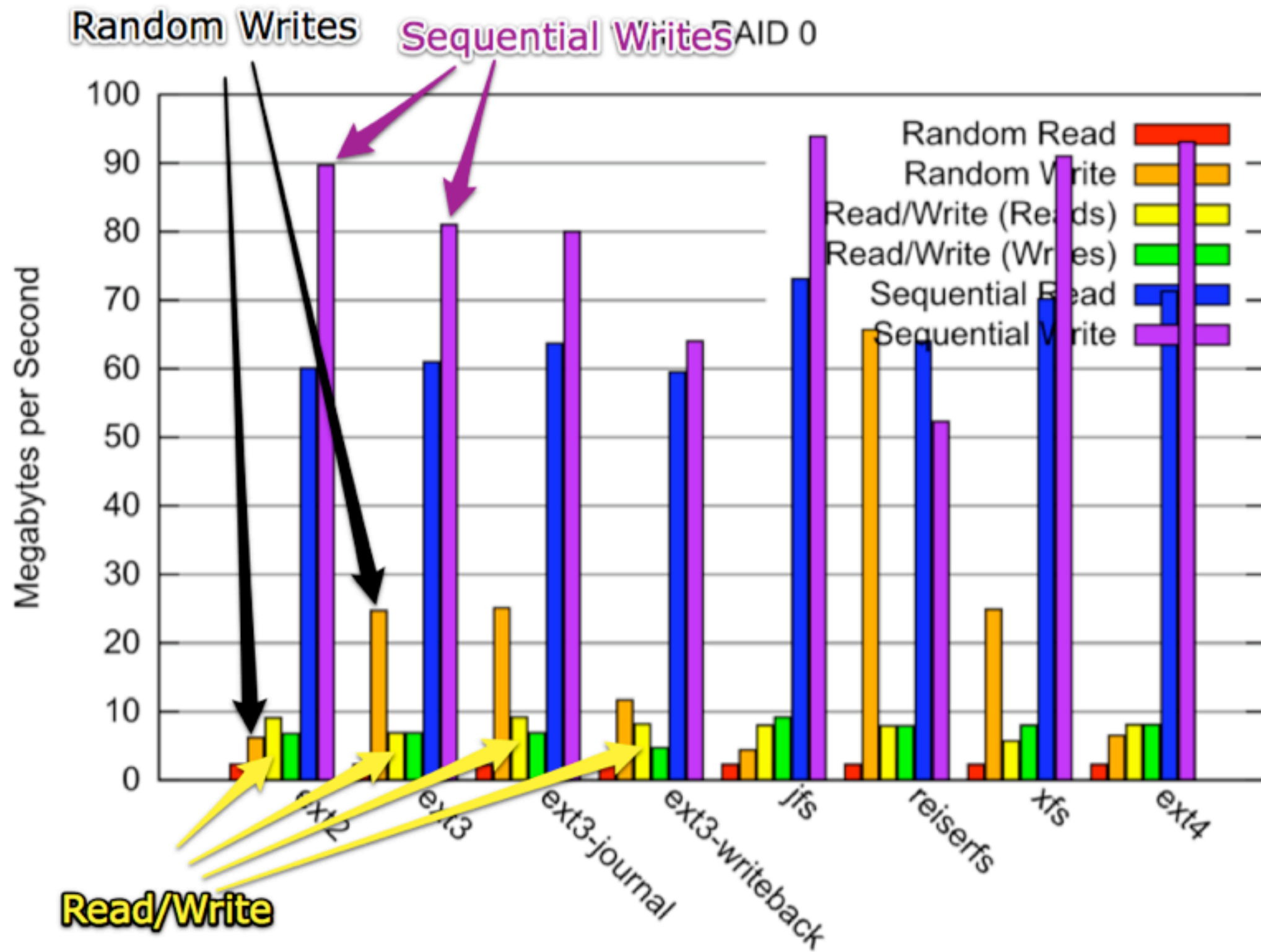
1 Disk RAID 0



Random Writes

1 Disk RAID 0

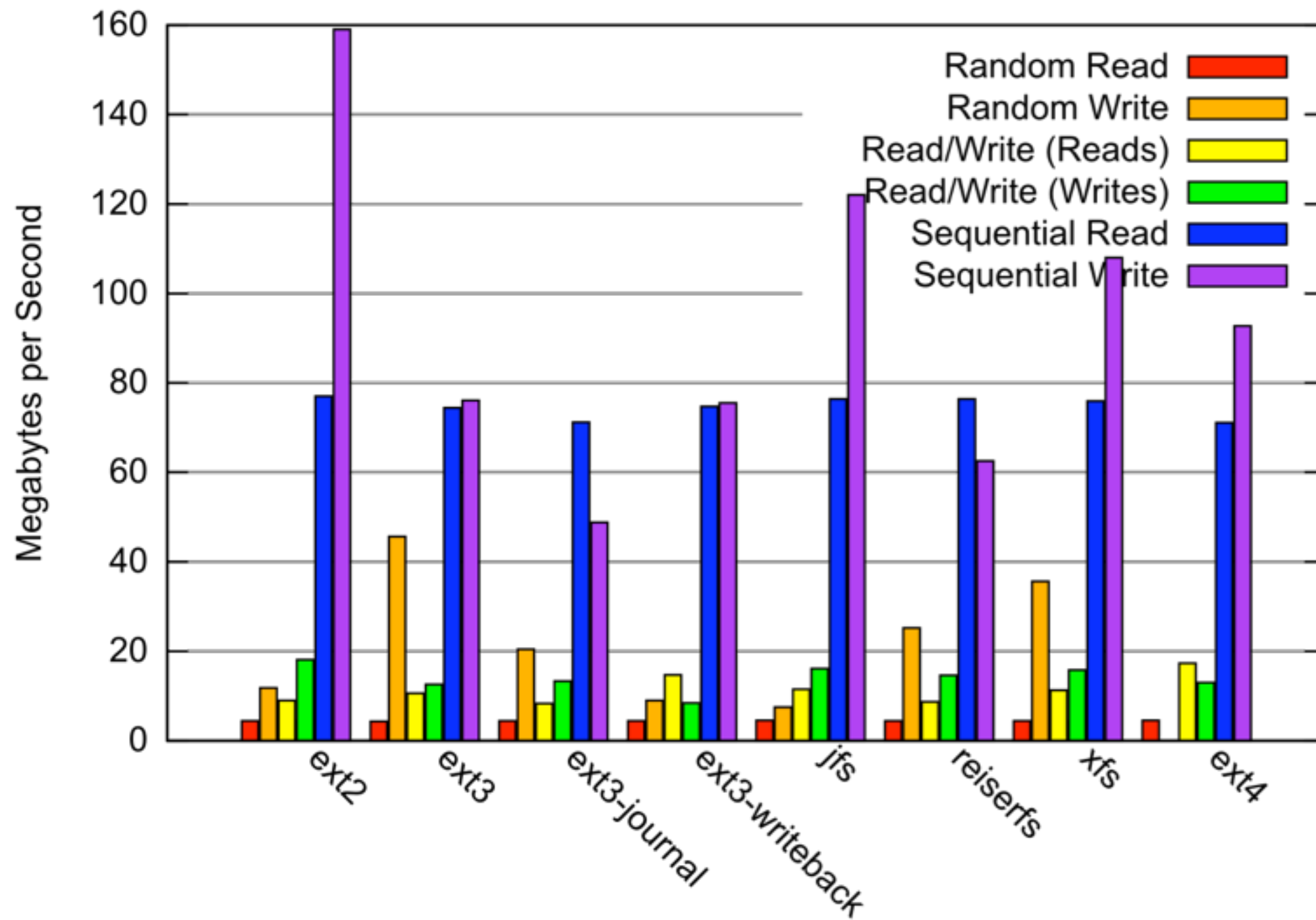




Seek bundling/batching
in ext3 is better?

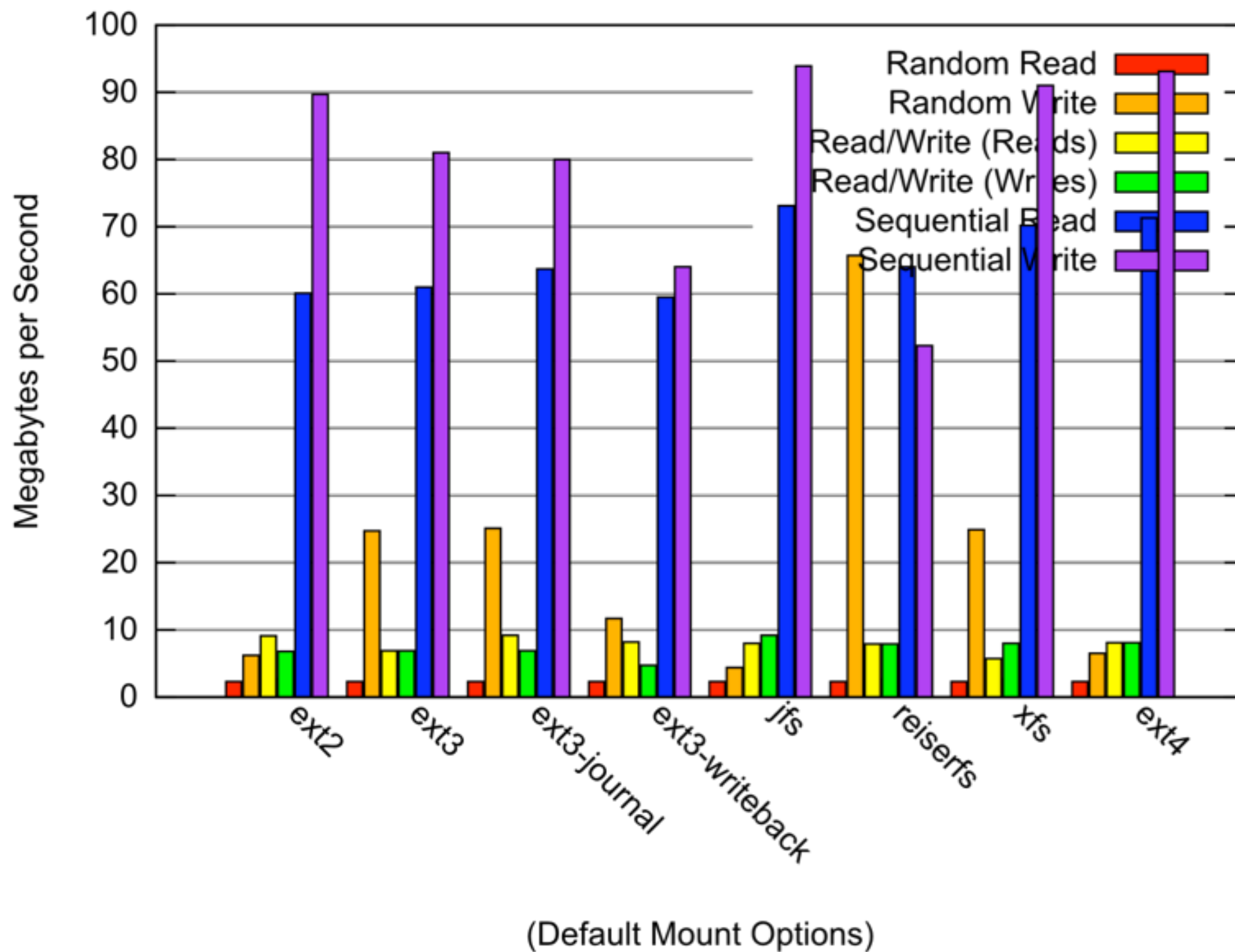
What if we add a disk?

2 Disk RAID 0

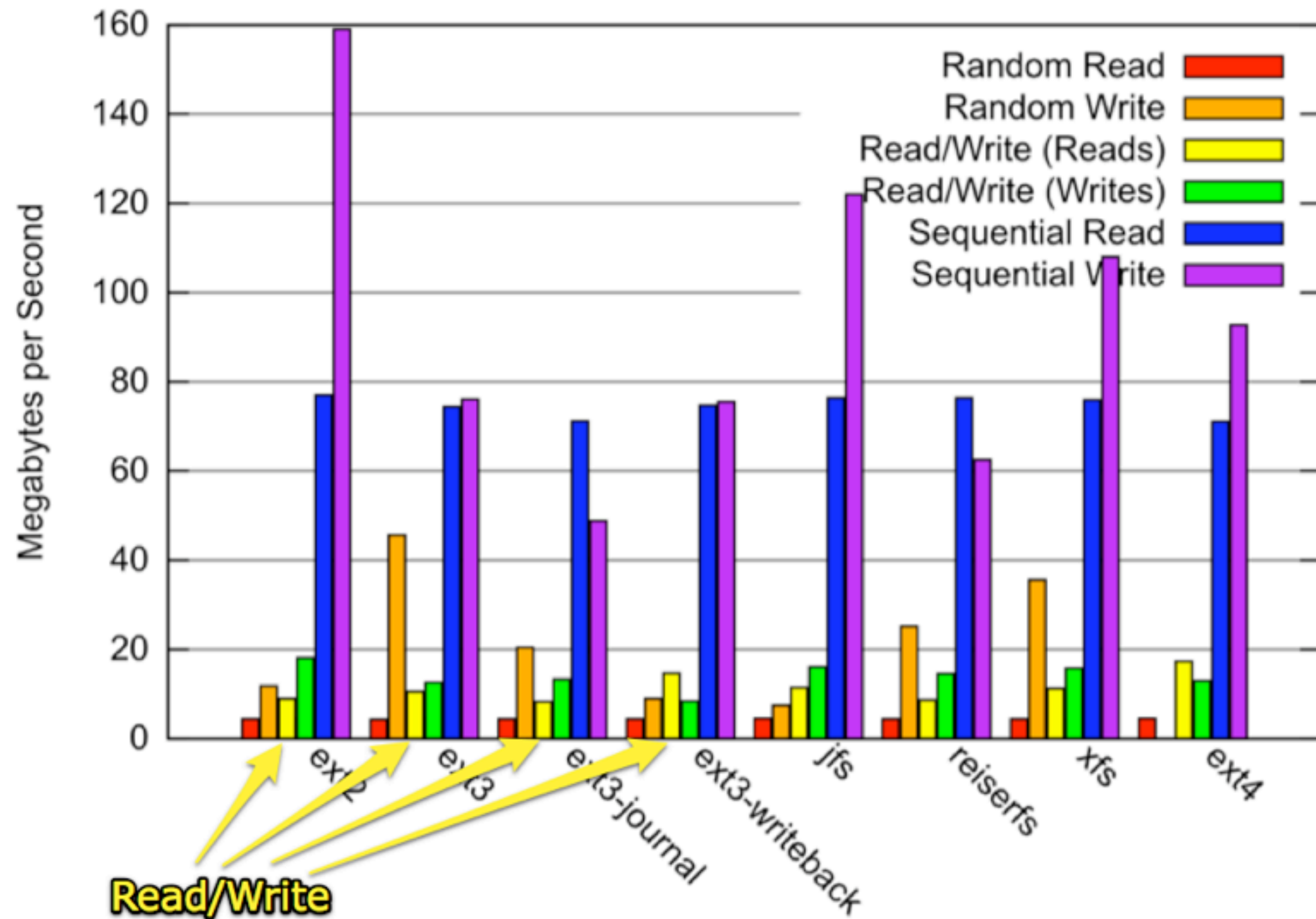


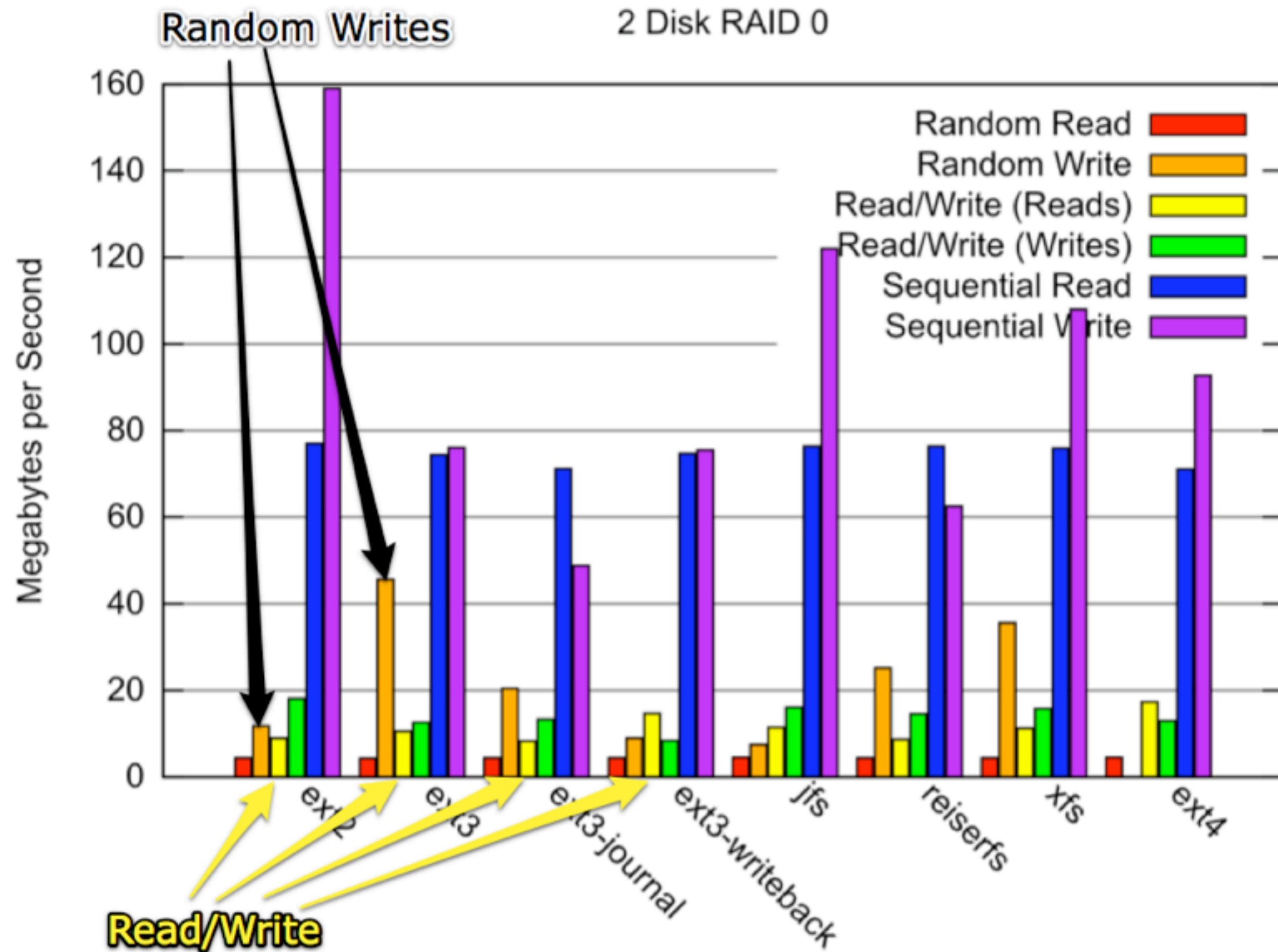
(Default Mount Options)

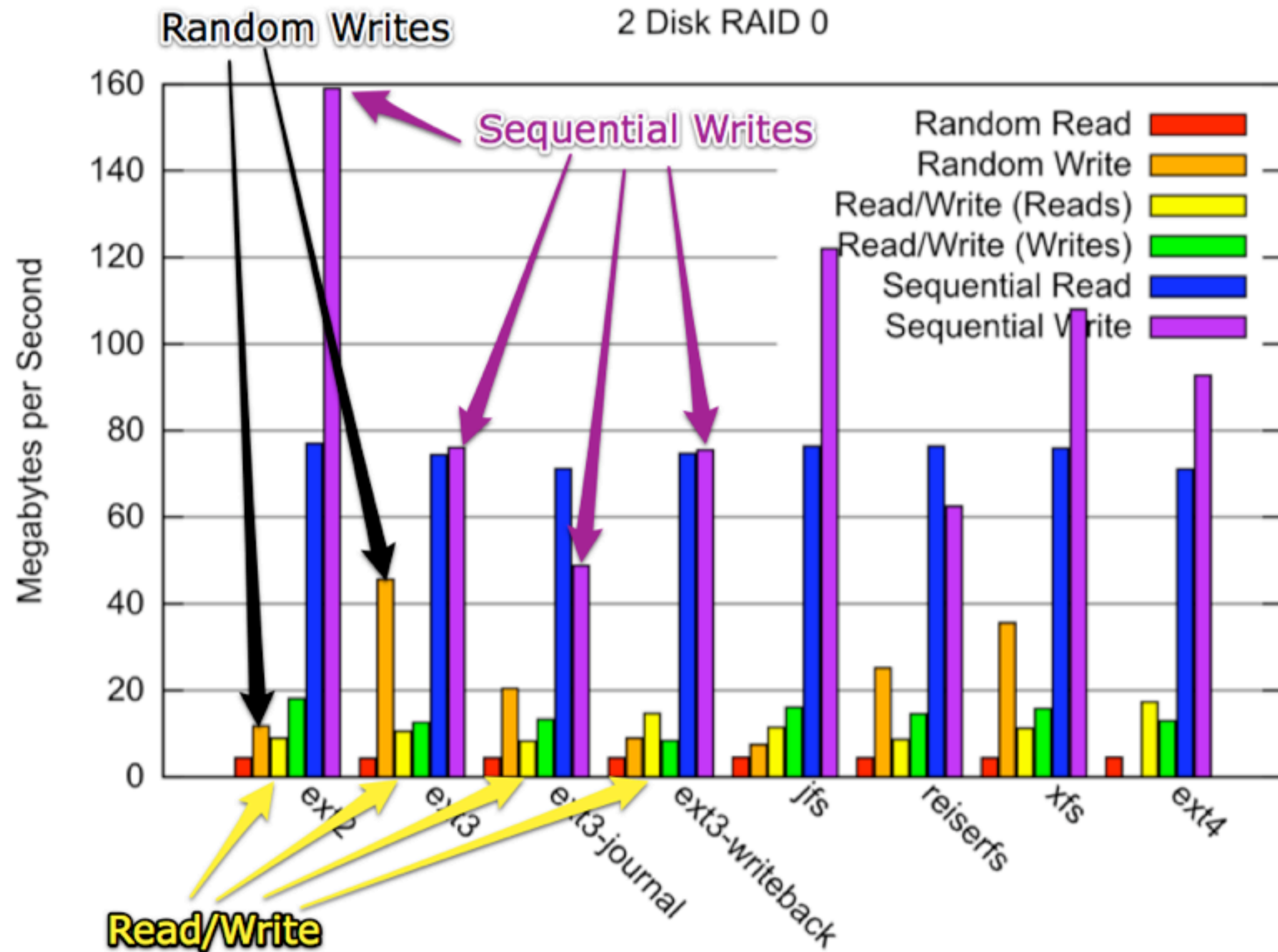
1 Disk RAID 0



2 Disk RAID 0







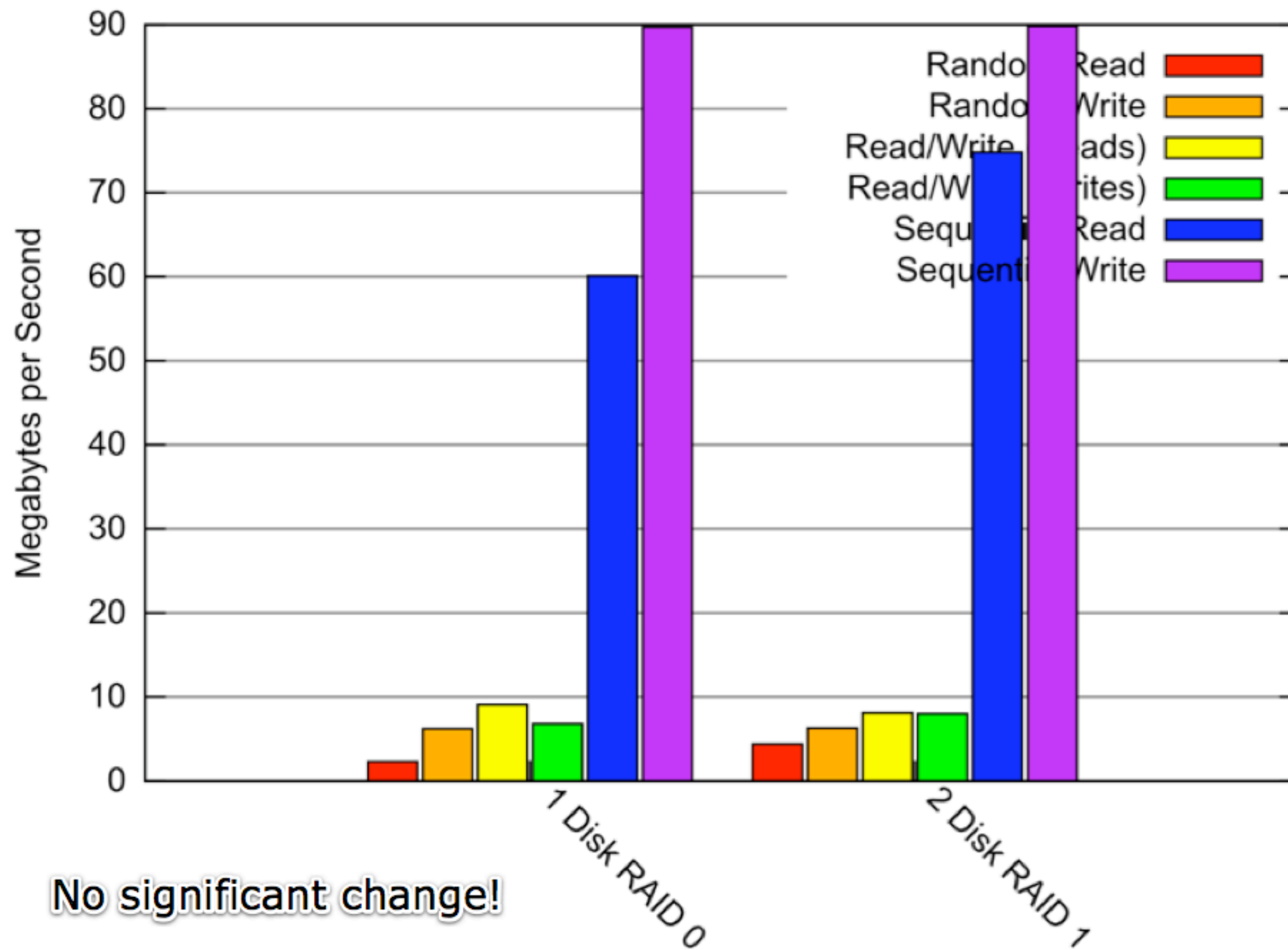
HOW TO START A FIGHT

AUDIENCE PARTICIPATION

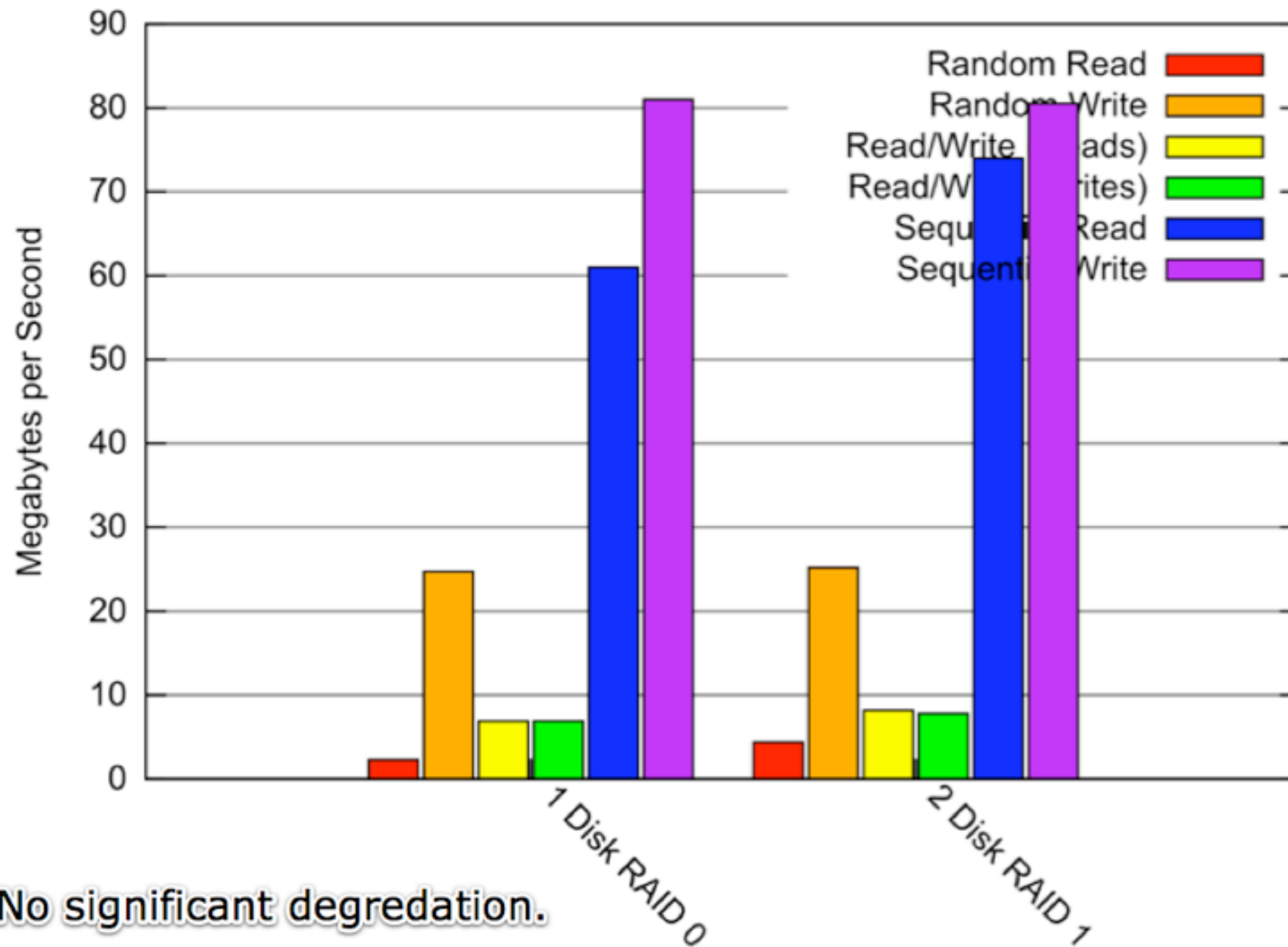
RAID 0 (stripe) versus
RAID 1 (mirroring)
performance?

~~PERFORMANCE~~

ext2 on 1 Disk RAID 0 vs 2 Disk RAID 1

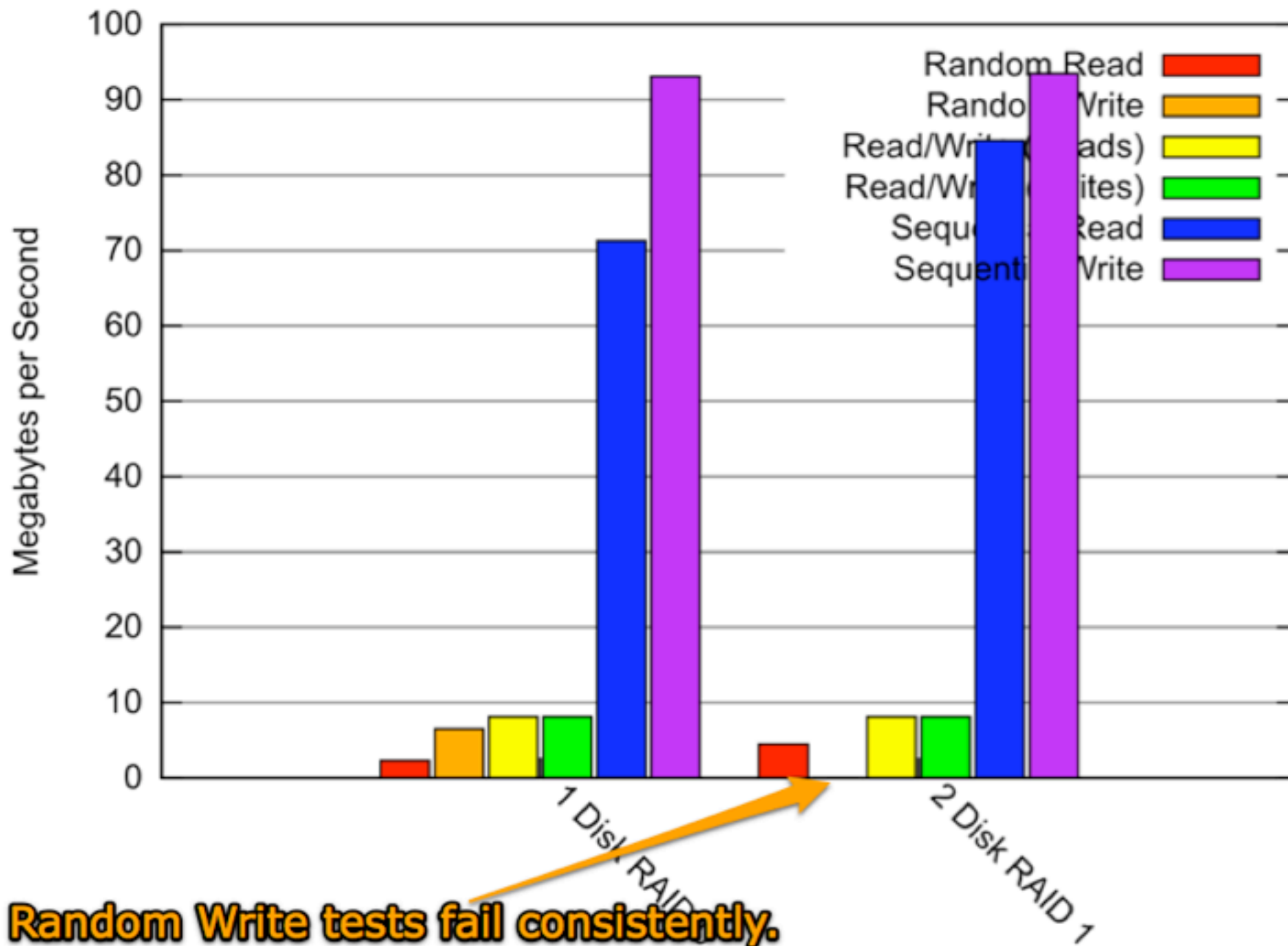


ext3 on 1 Disk RAID 0 vs 2 Disk RAID 1



No significant degradation.

ext4 on 1 Disk RAID 0 vs 2 Disk RAID 1



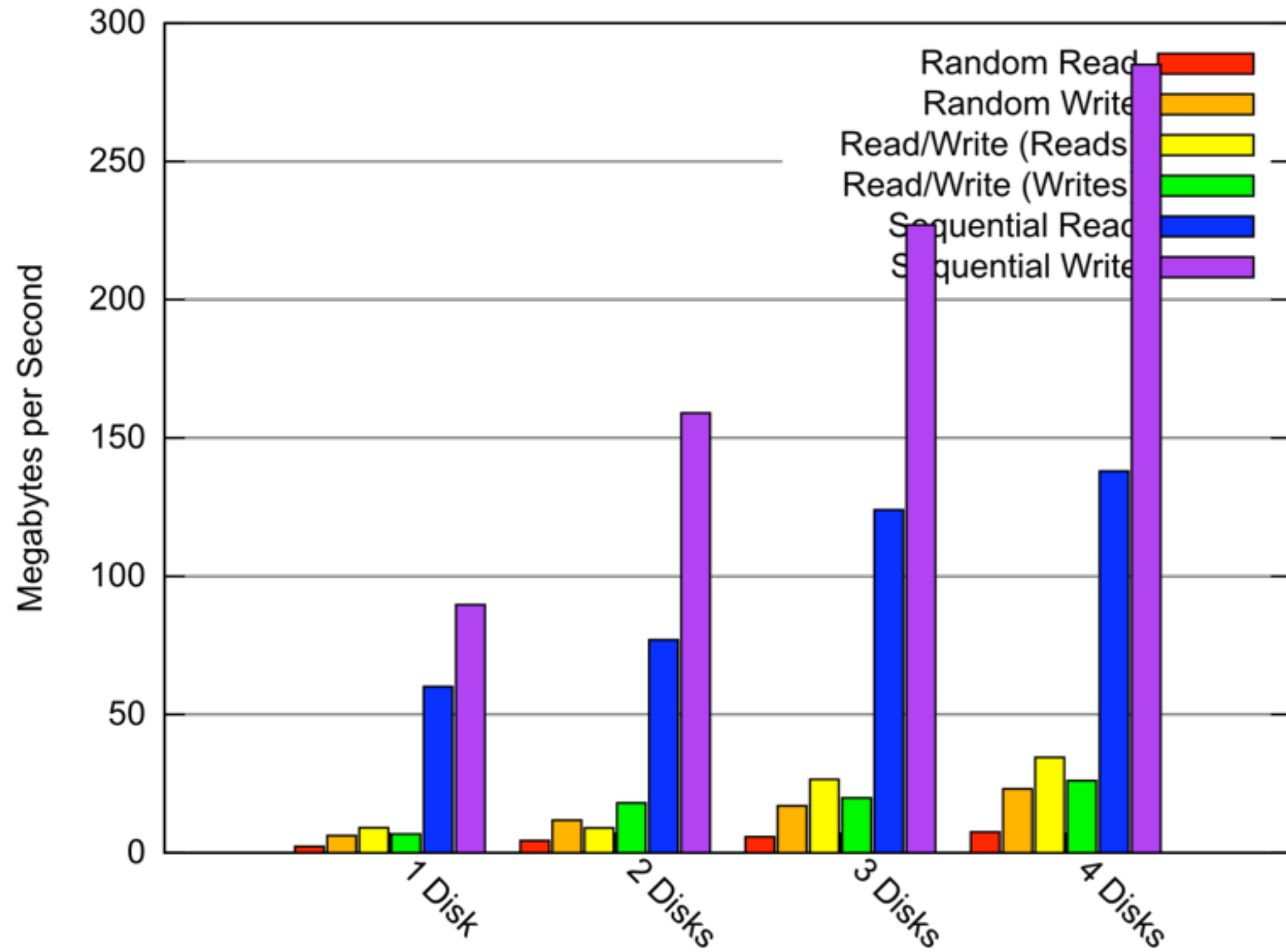
Random Write tests fail consistently.

HOW TO START A FIGHT

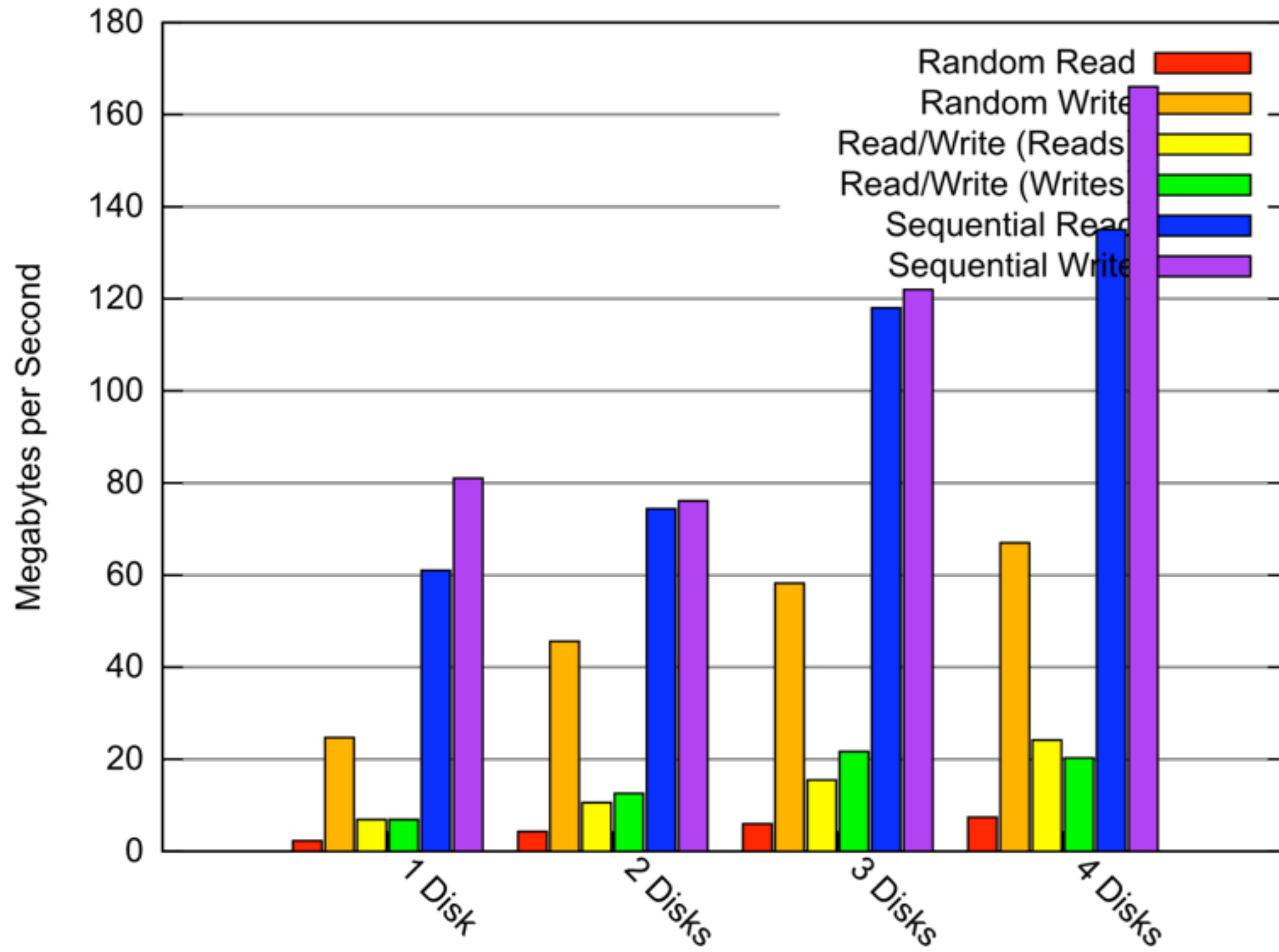
What happens when we:
add disks to a
RAID 0 (stripe) LUN?

PERFORMANCE

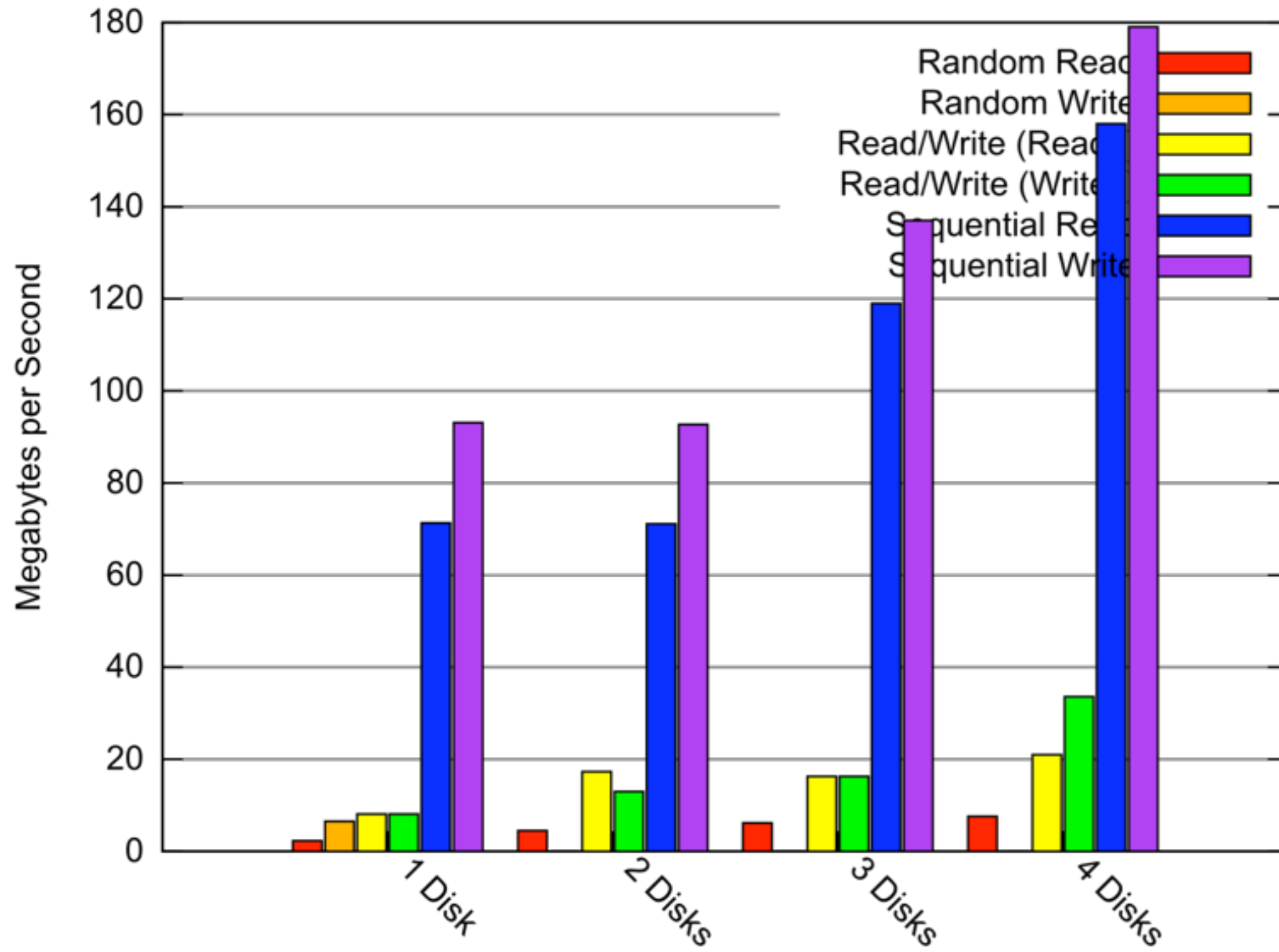
ext2 on RAID 0



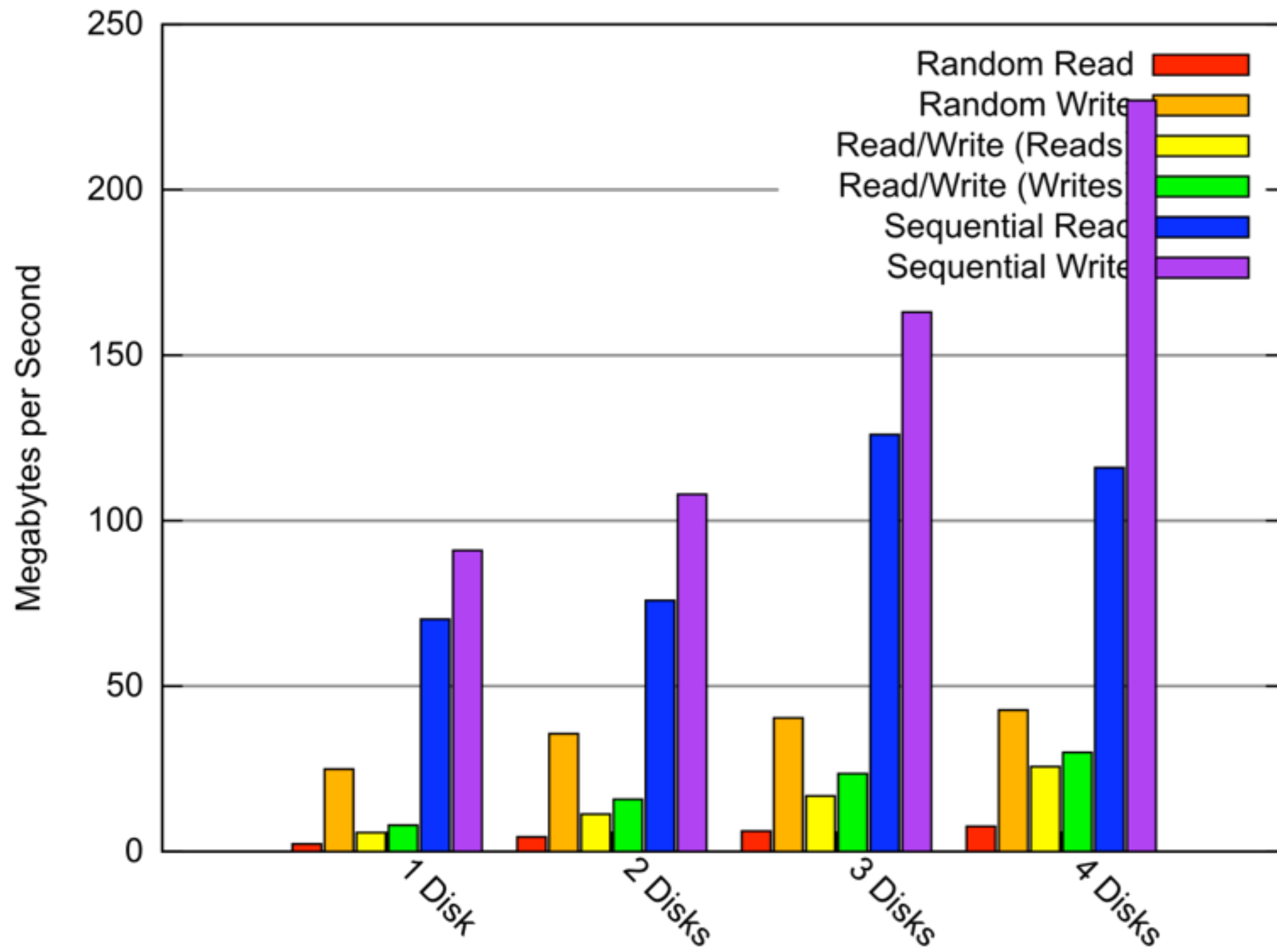
ext3 on RAID 0



ext4 on RAID 0



xfst on RAID 0

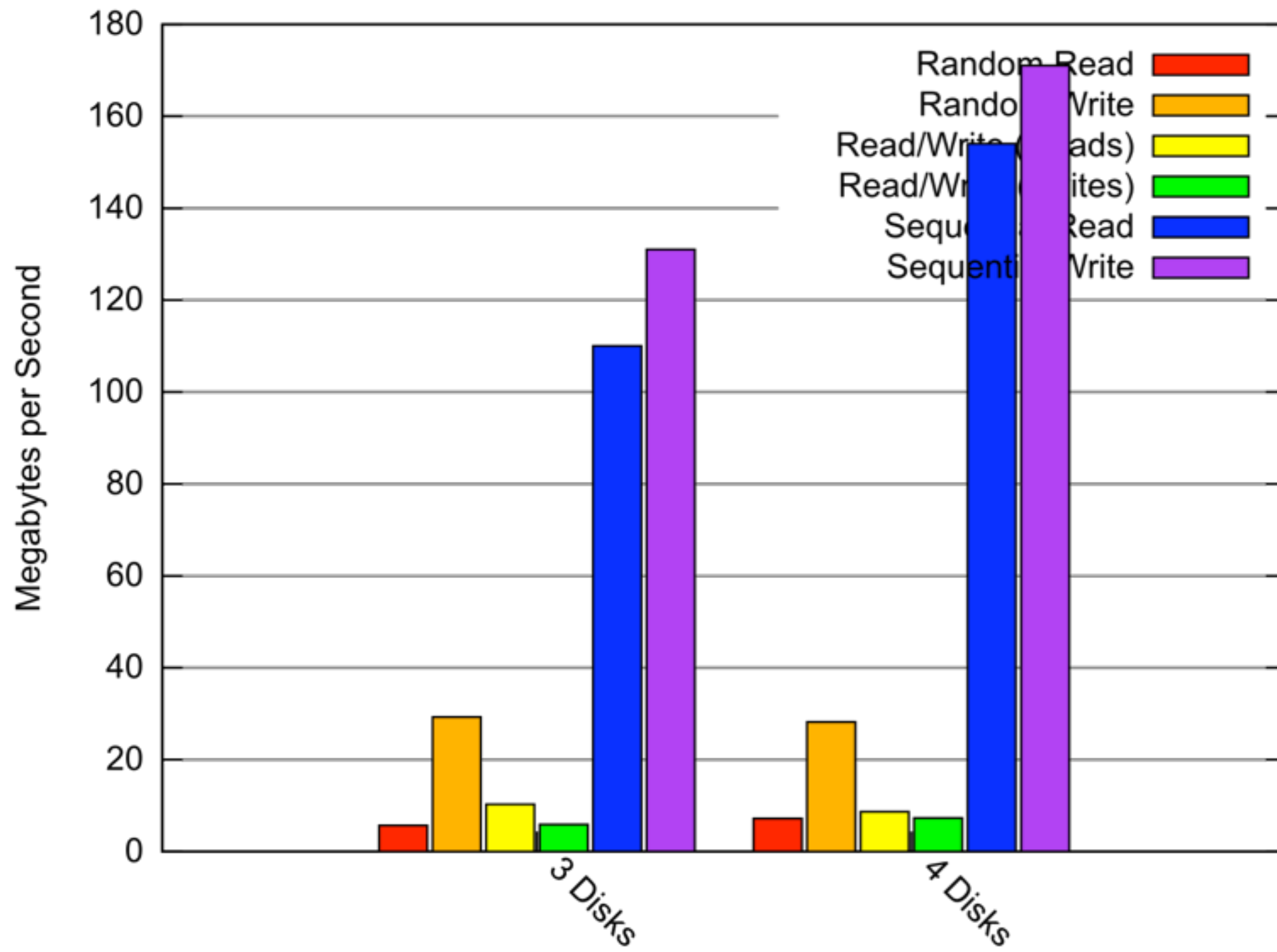


HOW TO START A FIGHT

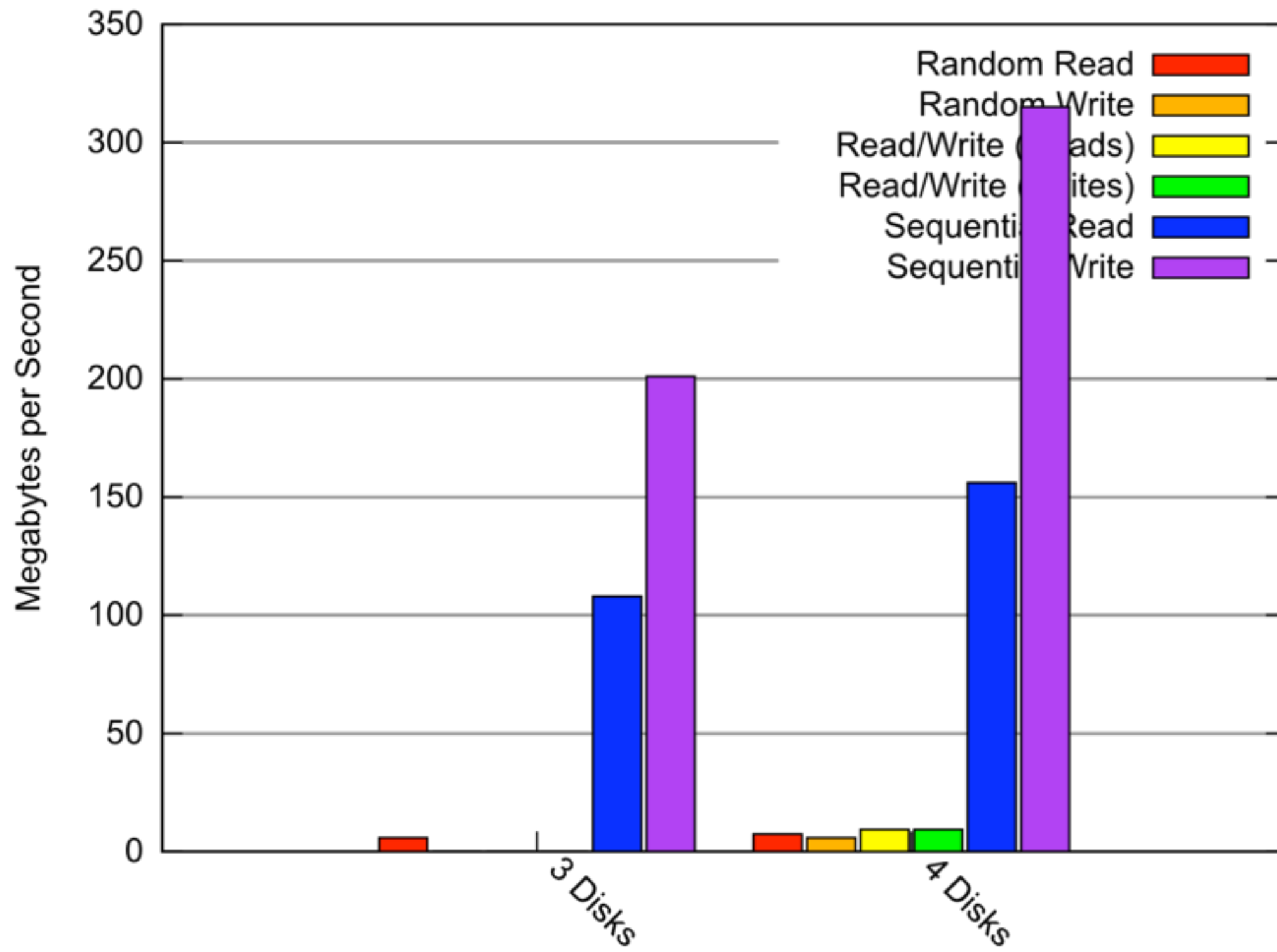
Adding disks to a
RAID 5 LUN

~~PERFORMANCE~~

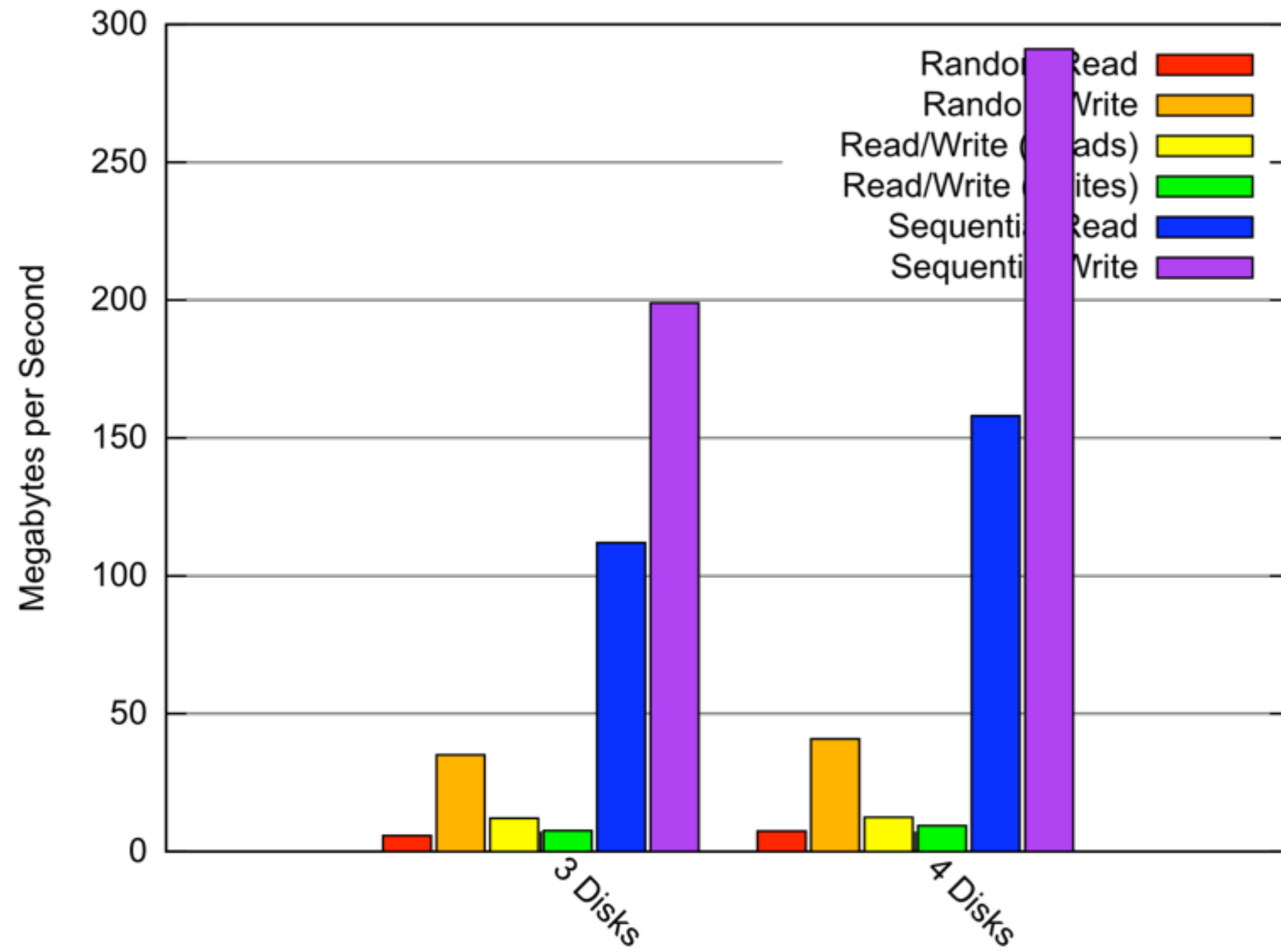
ext3 on RAID 5



ext4 on RAID 5



xfst on RAID 5

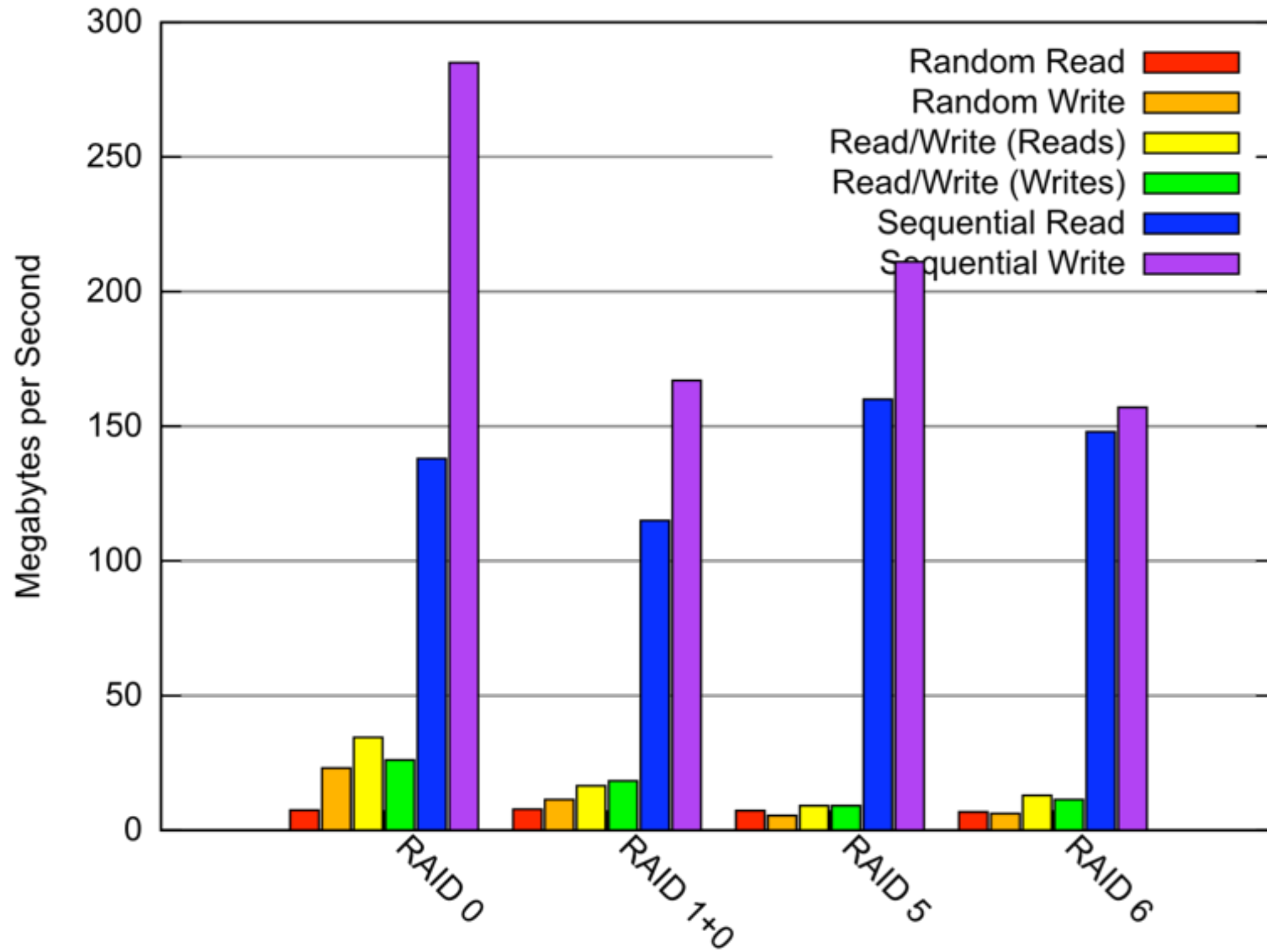


HOW TO START A FIGHT

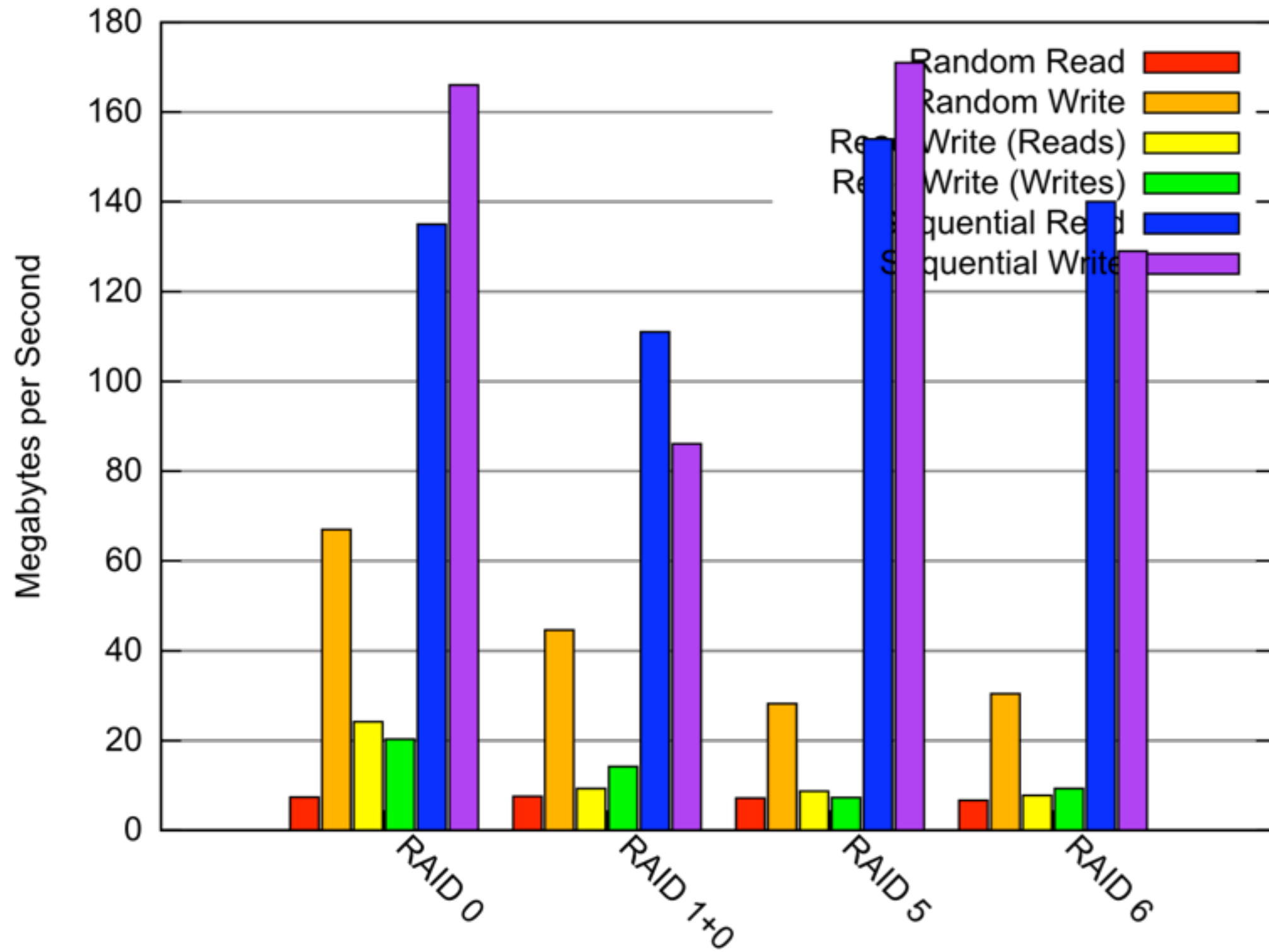
Only have 4 disks?
What should you do?

~~PERFORMANCE~~

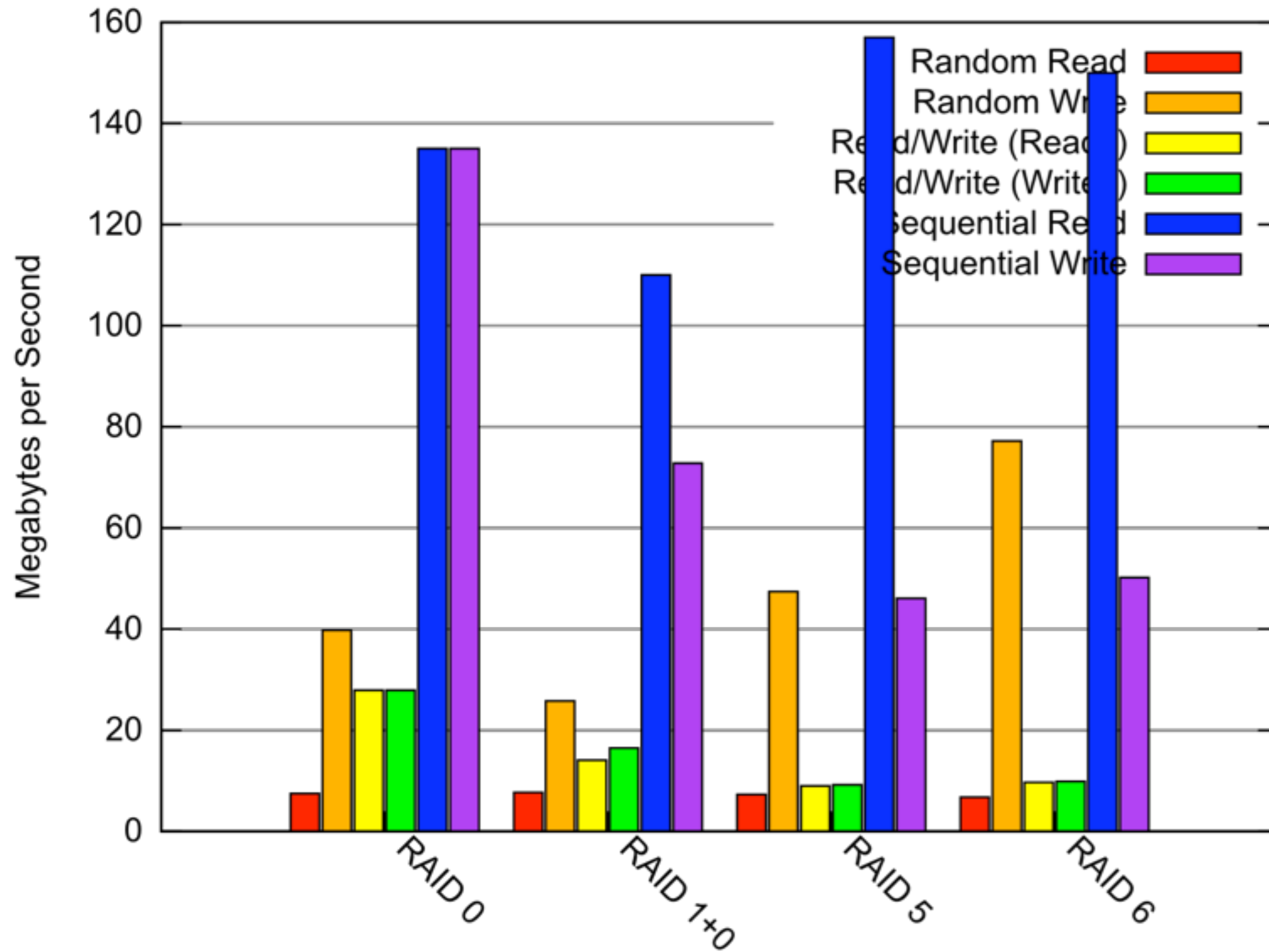
ext2 on 4 Disks



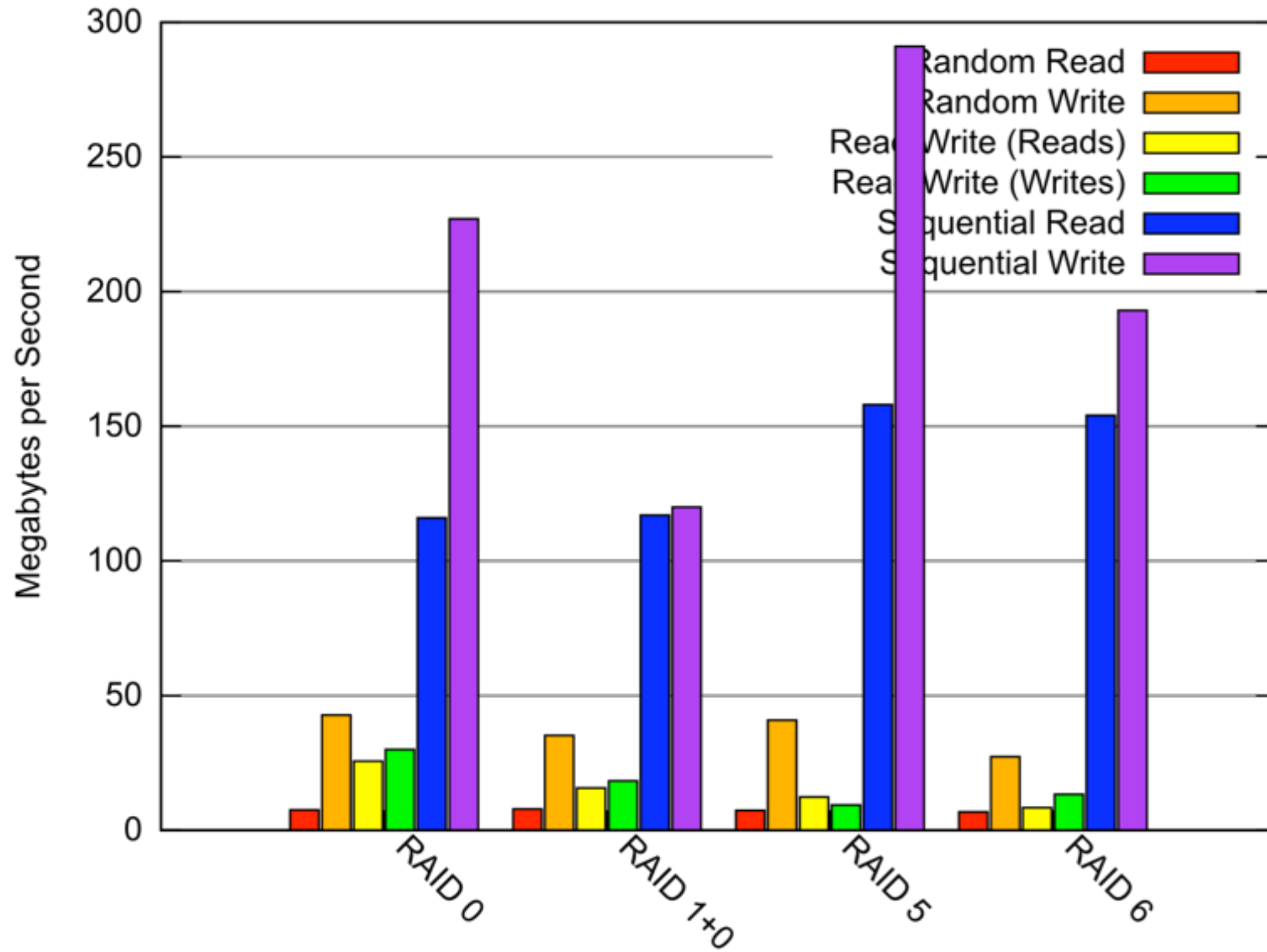
ext3 on 4 Disks



reiserfs on 4 Disks



xfst on 4 Disks



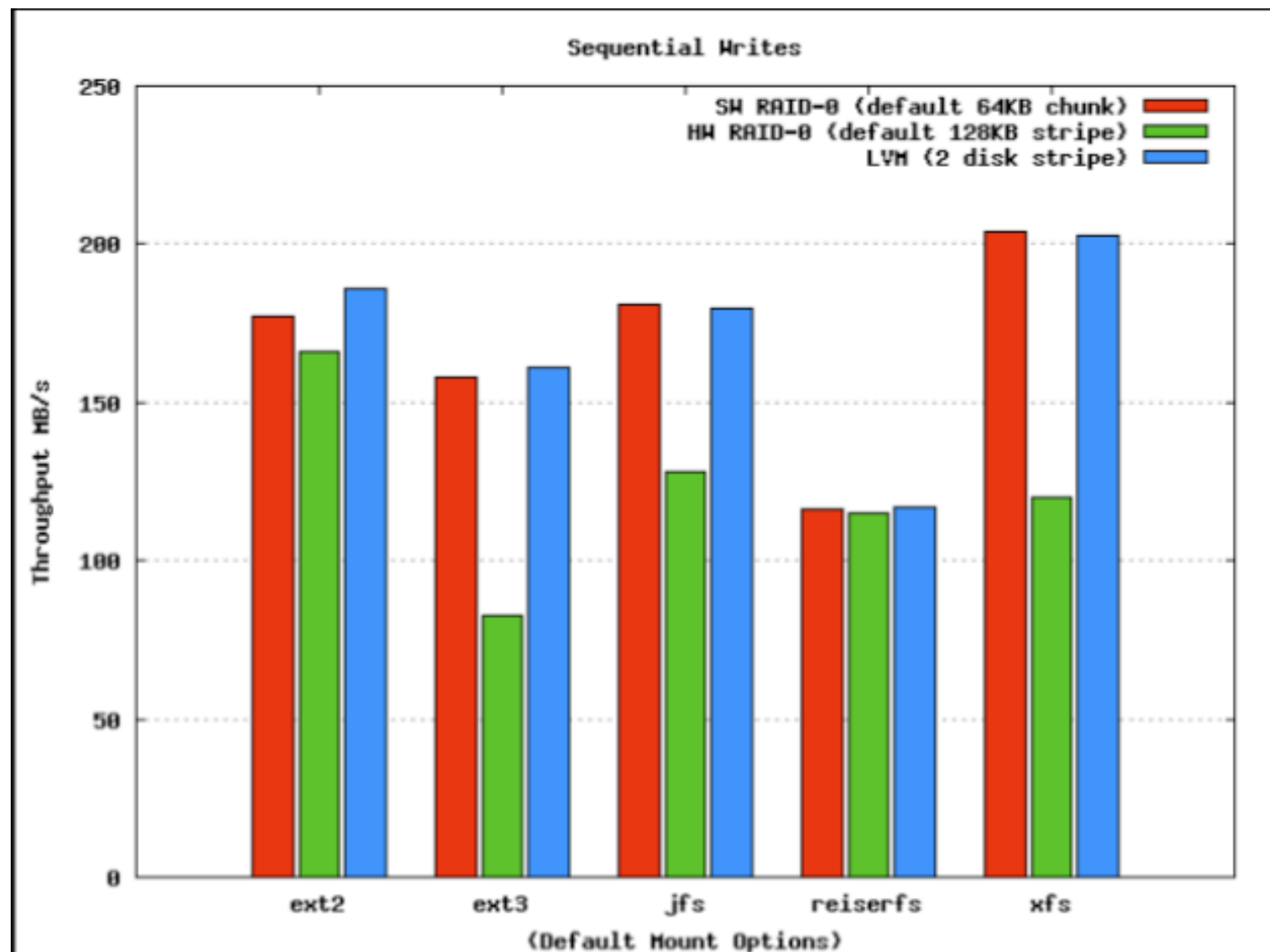
In most cases, RAID 5 out-performs
on sequential writes (xlog).

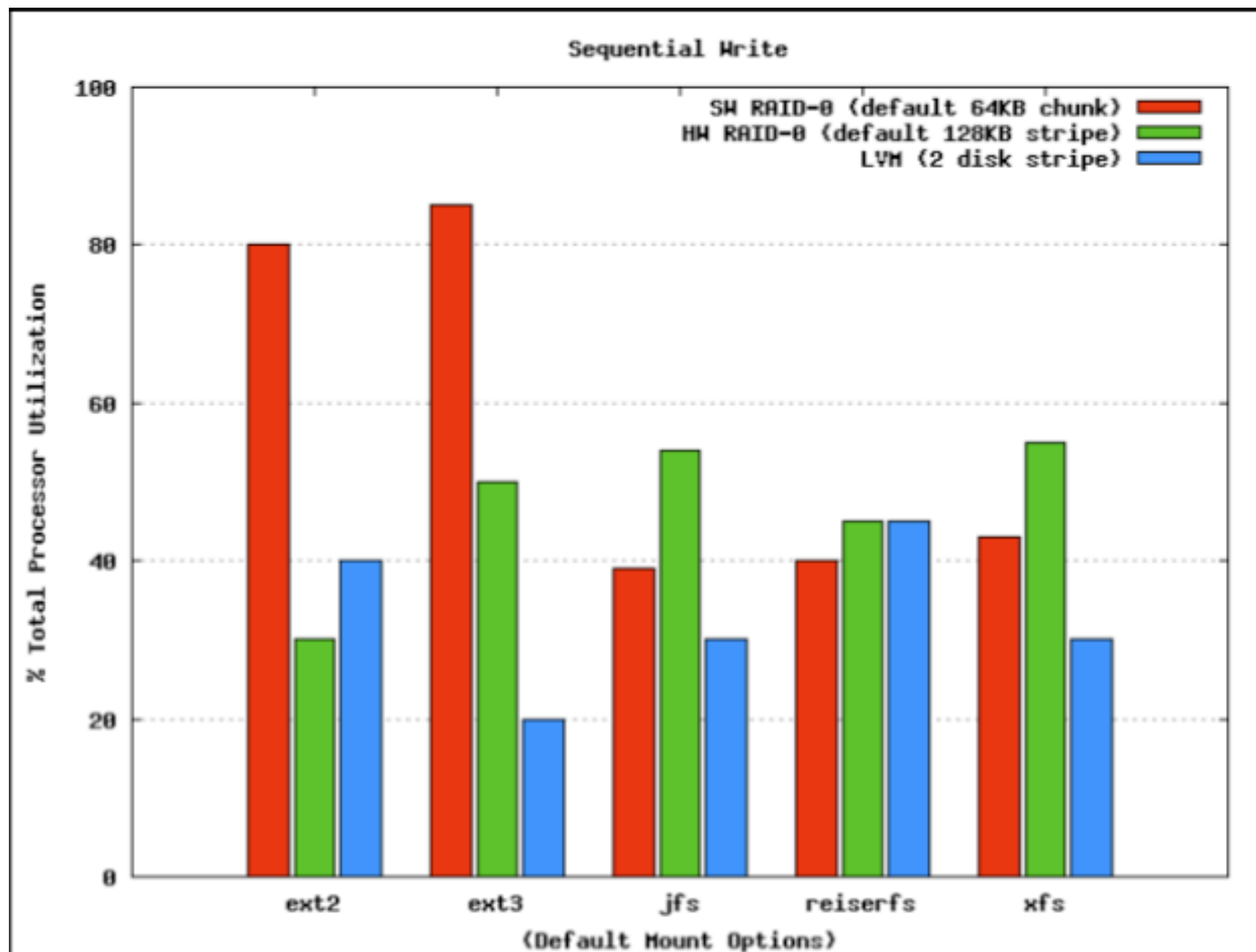
Random writes is only an improvement
on xfs and reiserfs.

HOW TO START A FIGHT

Are software RAID
and LVM slow?

~~PERFORMANCE~~





The Read-ahead buffer

HOW TO START A FIGHT

AUDIENCE PARTICIPATION

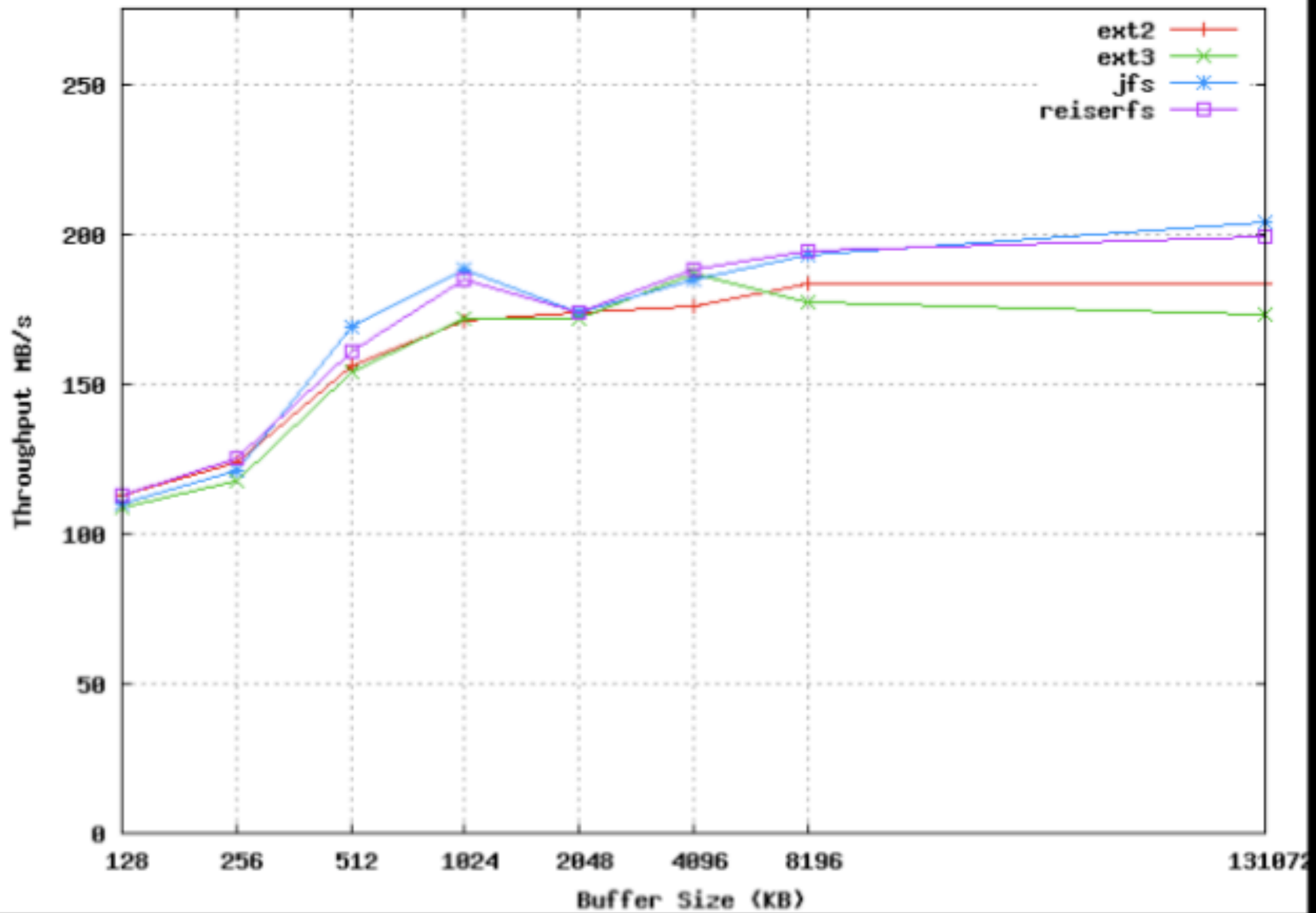
Readahead buffer:

Default is 128 K

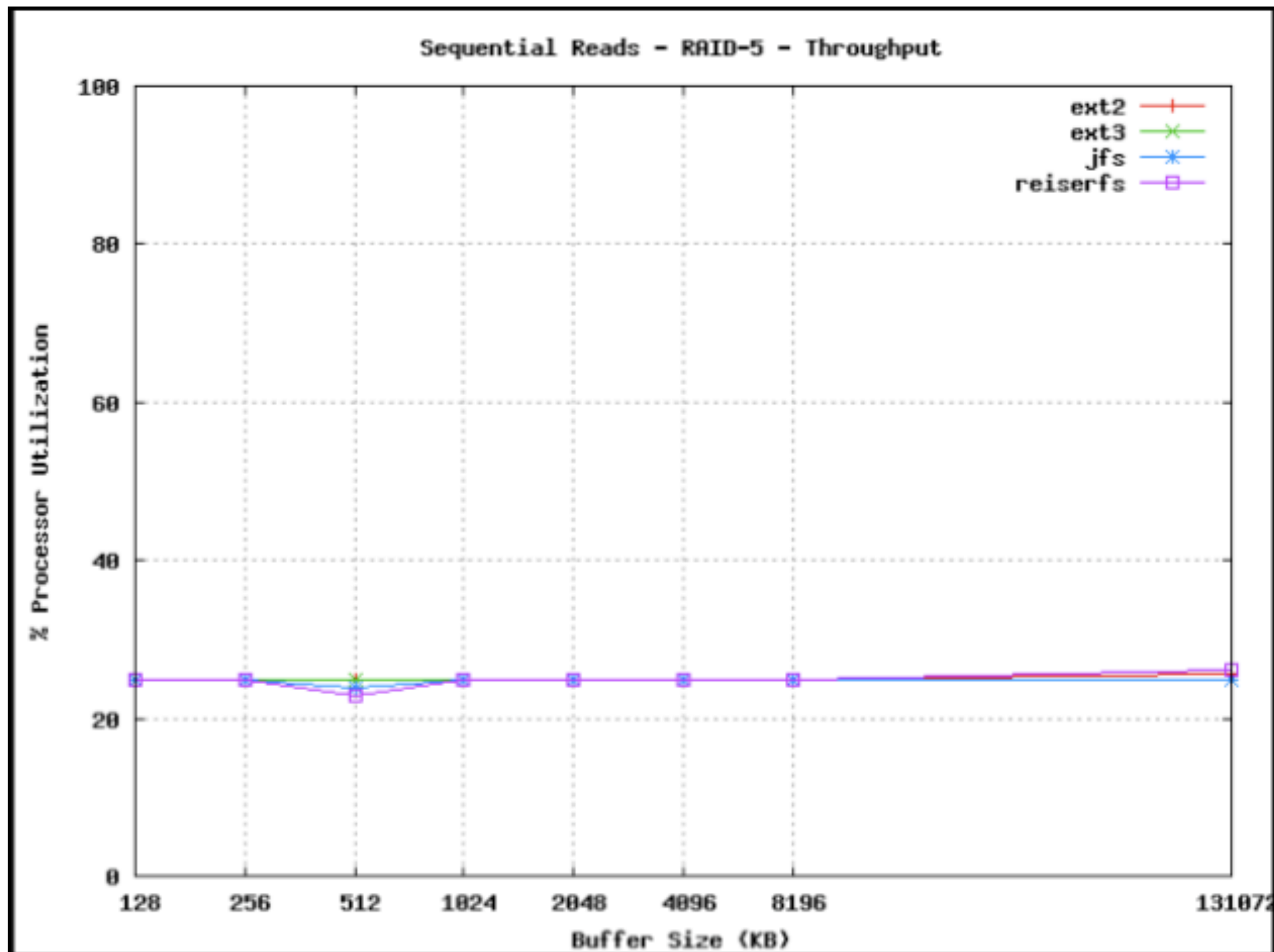
What do you think it should be?

~~PERFORMANCE~~

Sequential Reads - RAID-5

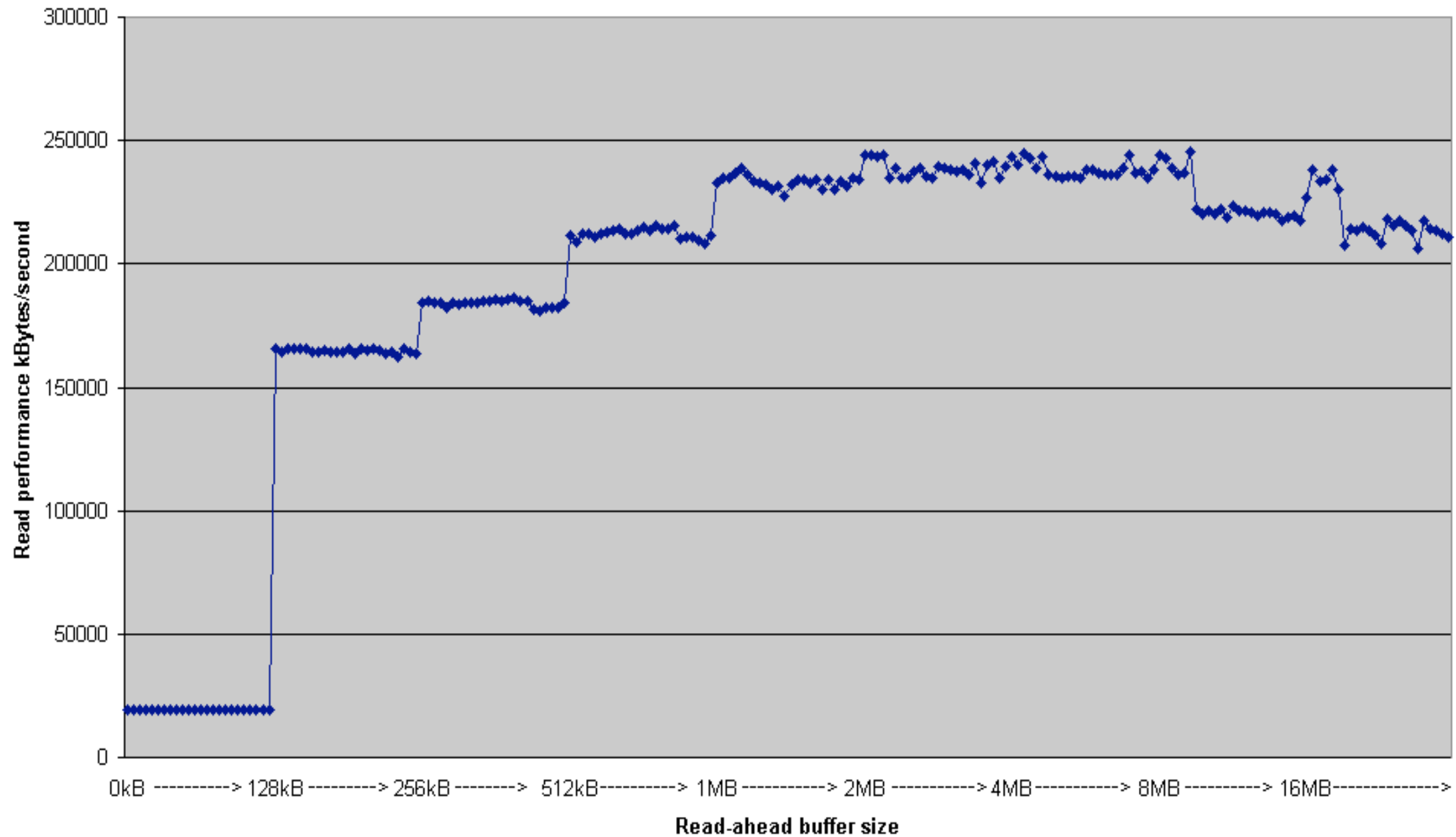


And is there a cost to
increasing the buffer
that much?



<http://moourl.com/readaheadconfirm>

Read performance versus Read-ahead buffer size (for f10b with 128kB Raid stripe size)



OLTP workload

- DBT-2 toolkit
(Fair-use derivative of TPC-C)
- Used 35 drives ultimately
- pgtune:
<http://pgfoundry.org/projects/pgtune/>

Out of the Box - From a 25 disk RAID 0 device¹

Transaction	%	Response Time (s)		Total	Rollbacks	%
		Average :	90th %			
Delivery	3.99	11.433 :	12.647	45757	0	0.00
New Order	45.24	10.257 :	11.236	518945	5224	1.02
Order Status	4.00	9.998 :	11.023	45926	0	0.00
Payment	42.81	9.983 :	11.022	491102	0	0.00
Stock Level	3.95	9.855 :	10.837	45344	0	0.00

8574.99 new-order transactions per minute (NOTPM)

59.3 minute duration

0 total unknown errors

1041 second(s) ramping up

This result is from before we ran pg tune to show if it'll help.

¹<http://207.173.203.223/~markwkm/community6/dbt2/baseline.1000.2/>

pgtune - From a 25 disk RAID 0 device²

Transaction	%	Response Time (s)		Total	Rollbacks	%
		Average :	90th %			
Delivery	3.99	8.715 :	10.553	48961	0	0.00
New Order	45.22	8.237 :	9.949	554565	5425	0.99
Order Status	3.95	8.037 :	9.828	48493	0	0.00
Payment	42.84	8.026 :	9.795	525387	0	0.00
Stock Level	3.99	7.829 :	9.563	48879	0	0.00

9171.46 new-order transactions per minute (NOTPM)

59.3 minute duration

0 total unknown errors

1041 second(s) ramping up

This result is from after running pgtune 0.3.

²<http://207.173.203.223/~markwkm/community6/dbt2/pgtune.1000.100.3/>

7% improvement! :)

For more info...

- See Mark Wong's blog:
<http://pugs.postgresql.org/blog/92>
- Takeaway: for DBT-2, increasing checkpoint_segments had the largest impact (fewer checkpoints :)

Future Work

- OLTP system characterization, sizing (ongoing)
- Daily OLTP regression testing
- More presentations
- P5 - PostgreSQL Portland Performance Pad PRACTICE (done!)

MOAR Hardware?

Thanks again, HP!

MSA70, DL380 in late 2009 ??

Let's recap...

“RAID5 is the worst choice for a database.” Fast for sequential writes in our tests.

“LVM incurs too much overhead to use. Software RAID is slower.” For reads - throughput is about the same, but saw higher CPU.

“Turning off 'atime' is a big performance gain.” Not in our tests. But, 2-3% for “free”.

“Journaling filesystems will have worse performance than non-journaling filesystems.” Turn the data journaling off on ext3, and you do see better performance, but there are edge cases and performance differences we could not explain.

“Striping doubles performance.”
Performance is better, but no where near double. Why?

“Your read-ahead buffer is big enough.”
Your read-ahead buffer IS NOT big
enough. Make it 8MB. And can we make
that the default?

Thank you!

Results:

[http://wiki.postgresql.org/wiki/
HP_ProLiant_DL380_G5_Tuning_Guide](http://wiki.postgresql.org/wiki/HP_ProLiant_DL380_G5_Tuning_Guide)

<http://moourl.com/fsperf>

Selena Deckelmann
selena@postgresql.org
twitter: @selenamarie