



On The Scalability of Storage Sub-System Back-end Network

Yan Li, Roland Ibbett, Nigel Topham and Tim Courtney[‡]

School of Informatics
University of Edinburgh
[‡]Xyratex, UK



Project Motivation

- Theoretically the more disks used in a disk array, the higher the degree of parallelism, so leading to larger the performance potential performance benefits.
- However, in a real system there is a limitation on the scale of RAID systems due to the limitation of interconnection network.
- The more disks are added to the system, the higher the contention for the shared media.
- When the number of disks and cache size in a RAID system reaches a certain threshold, there will be no further gain in performance by adding more disk or cache due to the saturation of the back-end network.



Project Goal

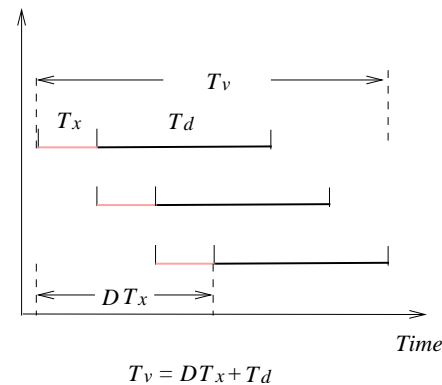
- Investigate the capacity of interconnection networks in terms of the numbers of disks that can be included in one chain. In particular, Fibre Channel (FC) SBOD is chosen as the interconnection network.
 - Give a certain number of disks and cache size, how much bandwidth is necessary to support them to get the maximum performance?
 - Likewise, how many disks (and cache) can be connected to a 2GFC (4GFC) port?



- Assumptions:
 - The request size of random access workload is equal to the size of stripe unit.
 - The access address of each request is aligned to the stripe unit boundary.
 - The capacity of the RAID system keeps fixed for all the study.
 - The queue length of disk is 1, ie., no disk command waits in the disk for service.
 - The FC SBOD is chosen as the research subject.
- Disk command transmission time $T_x = \frac{S}{B} + \text{overhead}$.
- Disk command execution time $T_d = \frac{S}{B_{port}} + \frac{S}{B_{media}} + T_{seek}$.
- S size of stripe unit; B network bandwidth;



Large Sequential Access

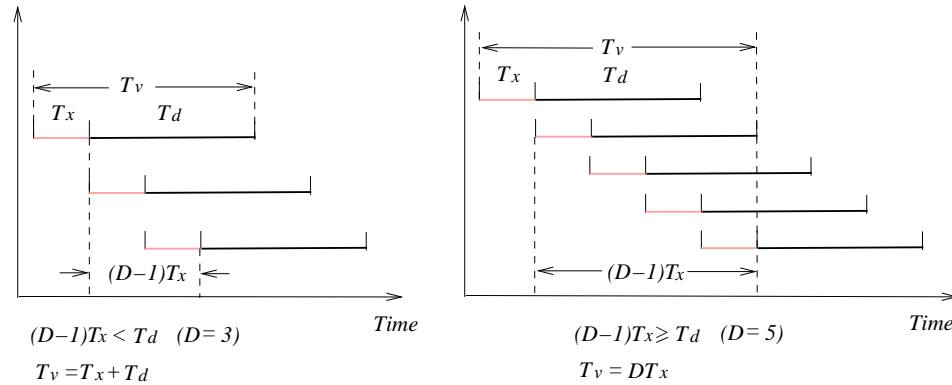


- For sequential access workload, a large volume command is divided into D disk commands, so that the response time for that command $T_v = DT_x + T_d$.
- Throughput is the major performance metric:

$$\text{Throughput} = \frac{K}{T_v} = \frac{K}{DT_x + T_d}.$$



Small Random Access



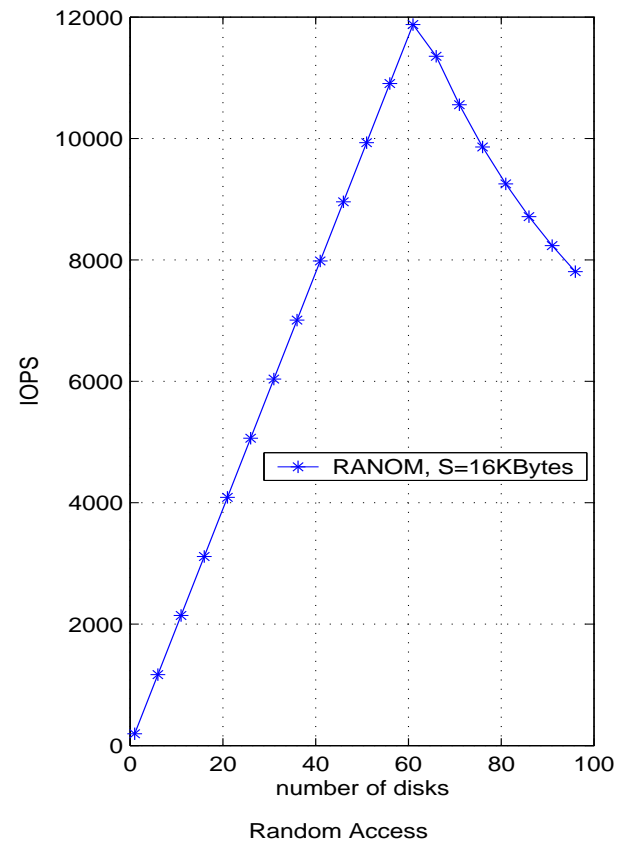
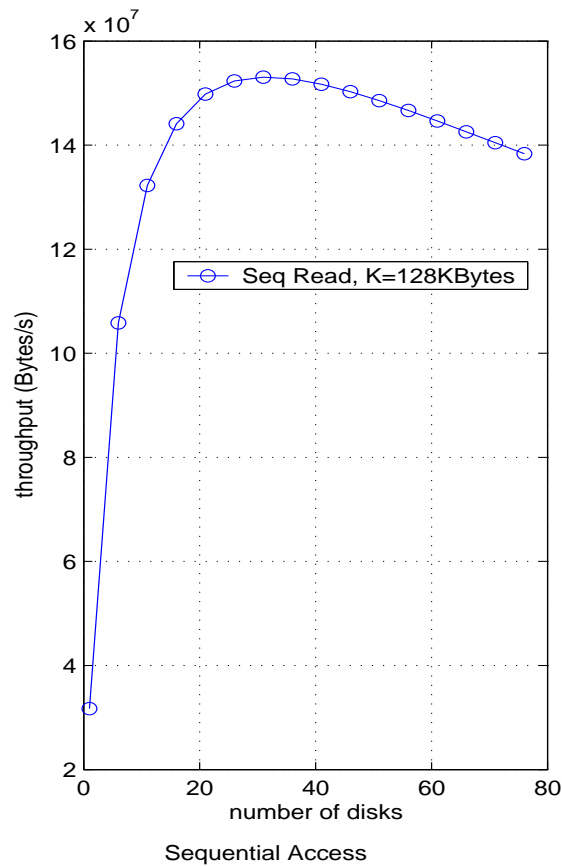
$$T_v = \begin{cases} T_x + T_d & (D-1) * T_x < T_d \\ D * T_x & (D-1) * T_x \geq T_d \end{cases}$$

- For random access workload, IOPS is the major performance metric,

$$\text{IOPS} = \begin{cases} \frac{D}{T_x + T_d} & (D-1) * T_x < T_d \\ \frac{T_d}{T_x} * \frac{1}{D * T_x} & (D-1) * T_x \geq T_d \end{cases}$$



Analytical Results (no cache)





General Model

- $B = F(D, C, S, L, P) = Num_d * S + overhead$

B , the bandwidth required to achieve the maximum performance with D disks and C cache in system.

D : number of disks

C : cache size

S : size of stripe unit

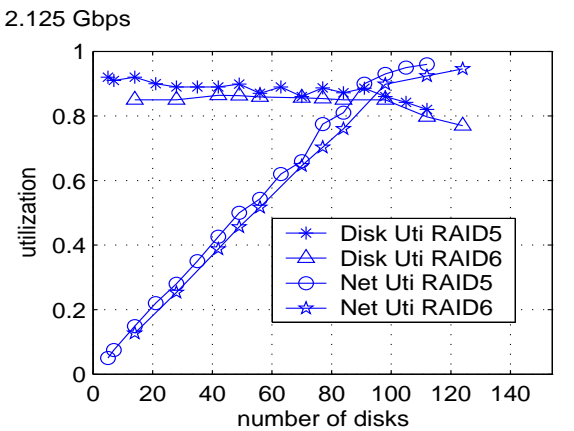
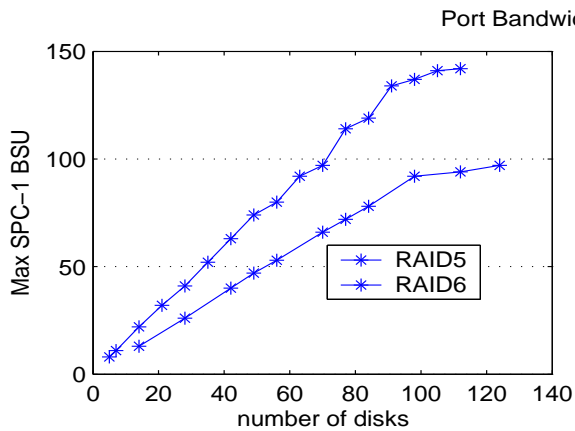
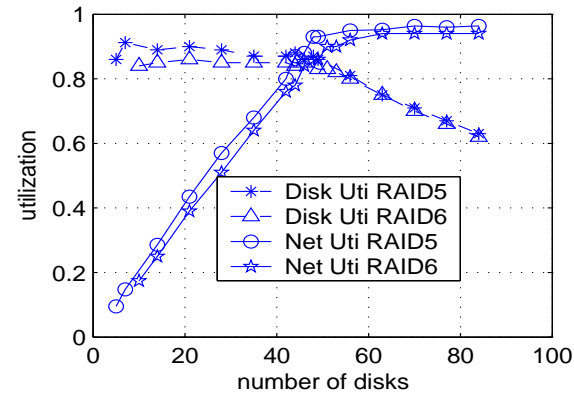
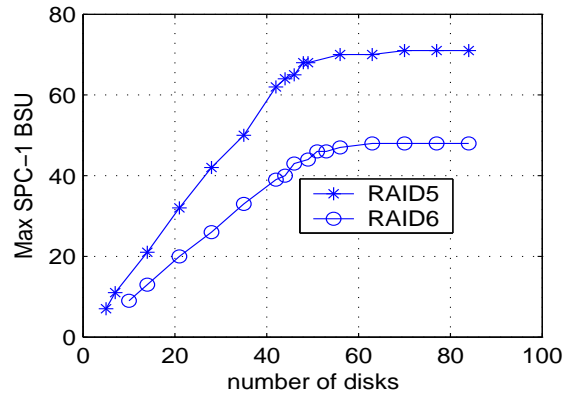
L : workload characteristic

P : cache destage threshold, $P=0$ for our study

Num_d : number of disk commands send to disk per second

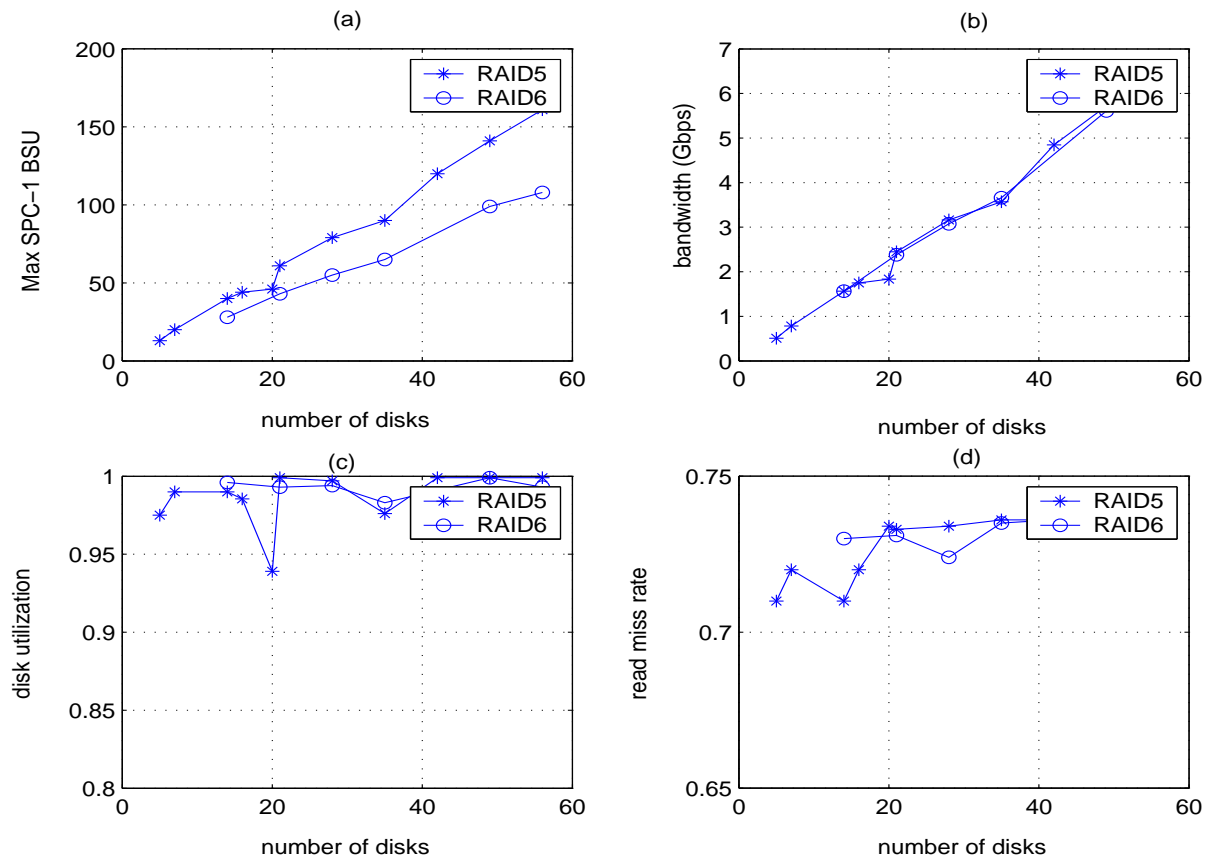


Systems without Cache ($S = 16KB$)





Systems with Cache ($S = 32KB$)





Summary

- When this is no cache, a 2G FC port is able to support up to 46 disks for RAID5 and 53 disks for RAID6 (size of stripe unit = 16kBytes).
- With enough cache, a 2G FC port is able to support up to 18 disks under OLTP like workload. (size of stripe unit = 32KBytes).
- When there is enough cache, the bandwidth required to support a certain number of disks is fixed. It is irrelevant with protection level and cache size.



Future Work

- Study the scalability of back-end network when the size of stripe unit is 16k and the system performance.
- Study the network bandwidth requirement when there is cache coherency between two controllers.