

Tutorial 1: Introduction to Data Science

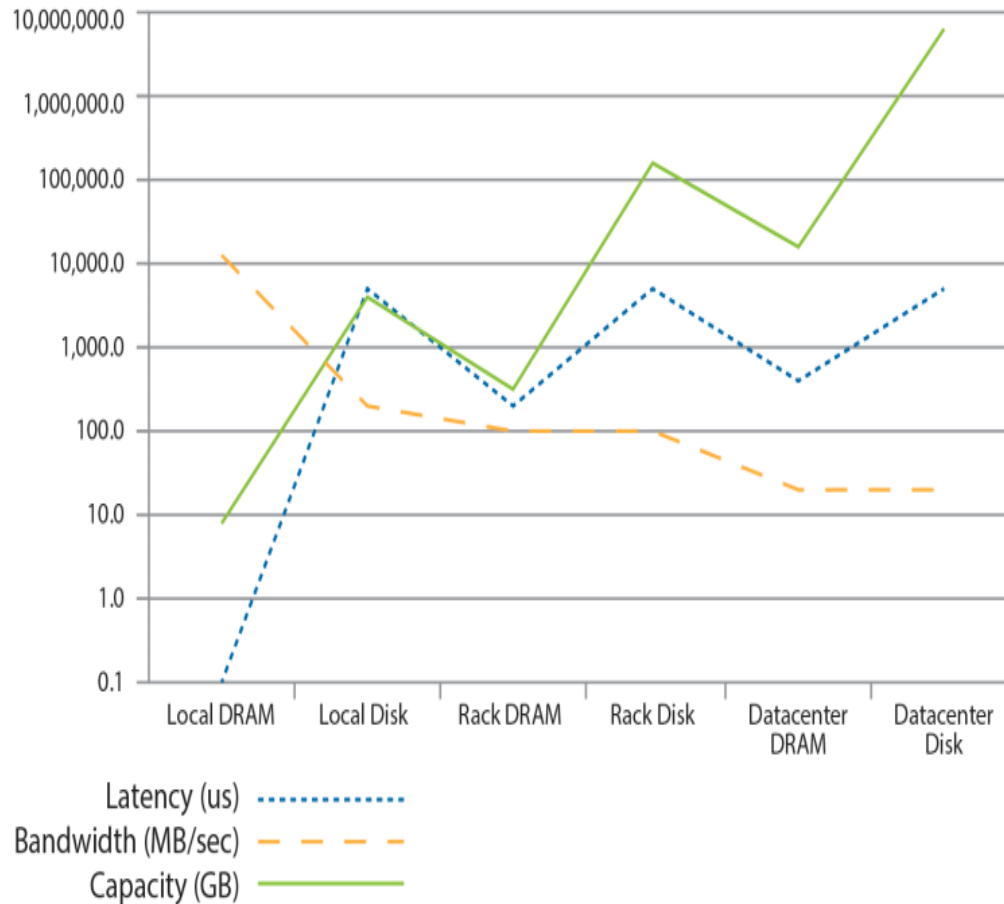
Li Yuan

li.yuan@u.nus.edu

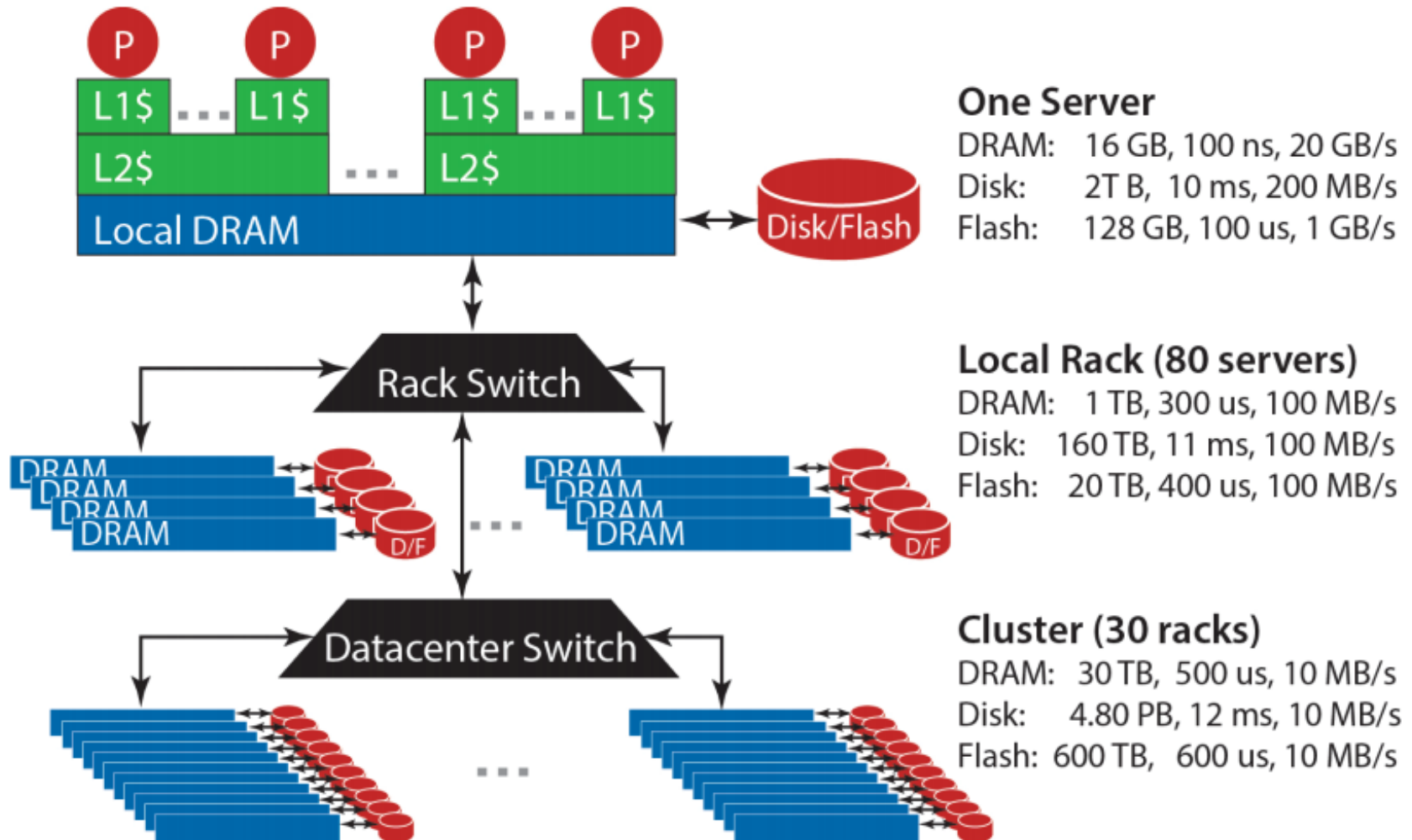


Question 1a

Why the latency of Rack DRAM is much higher than that of Local DRAM?

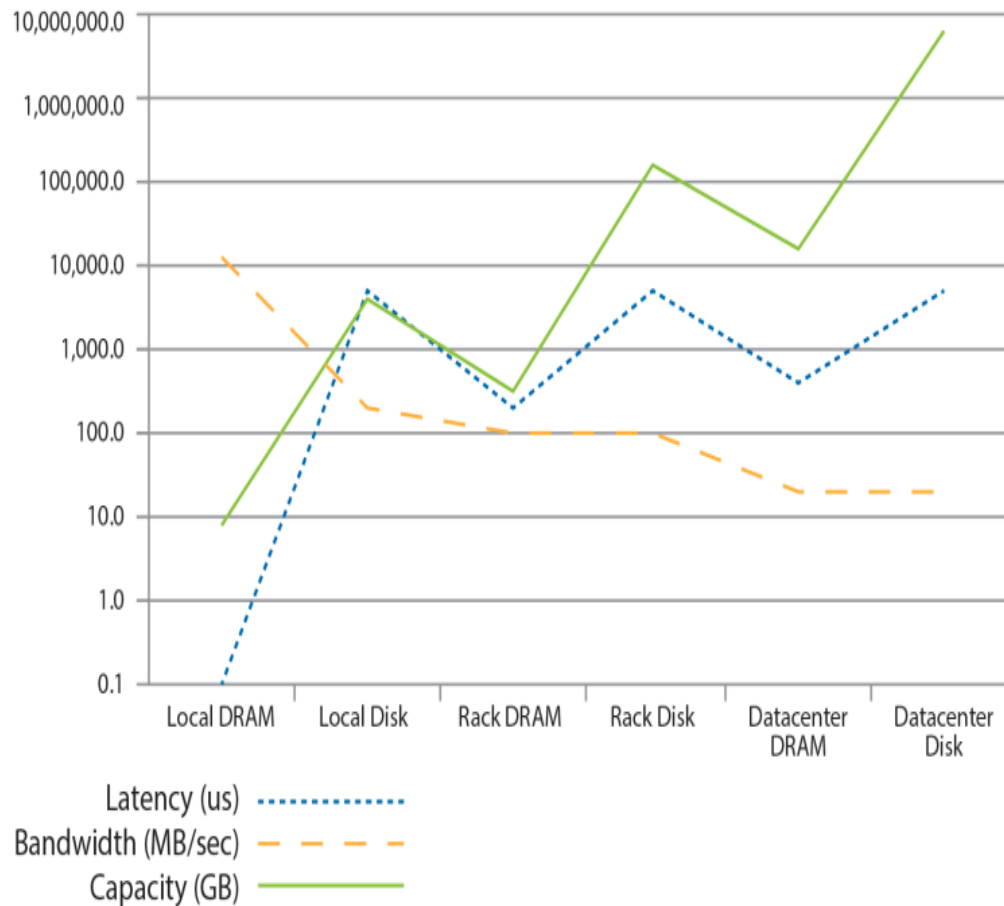


Storage Hierarchy



Question 1a

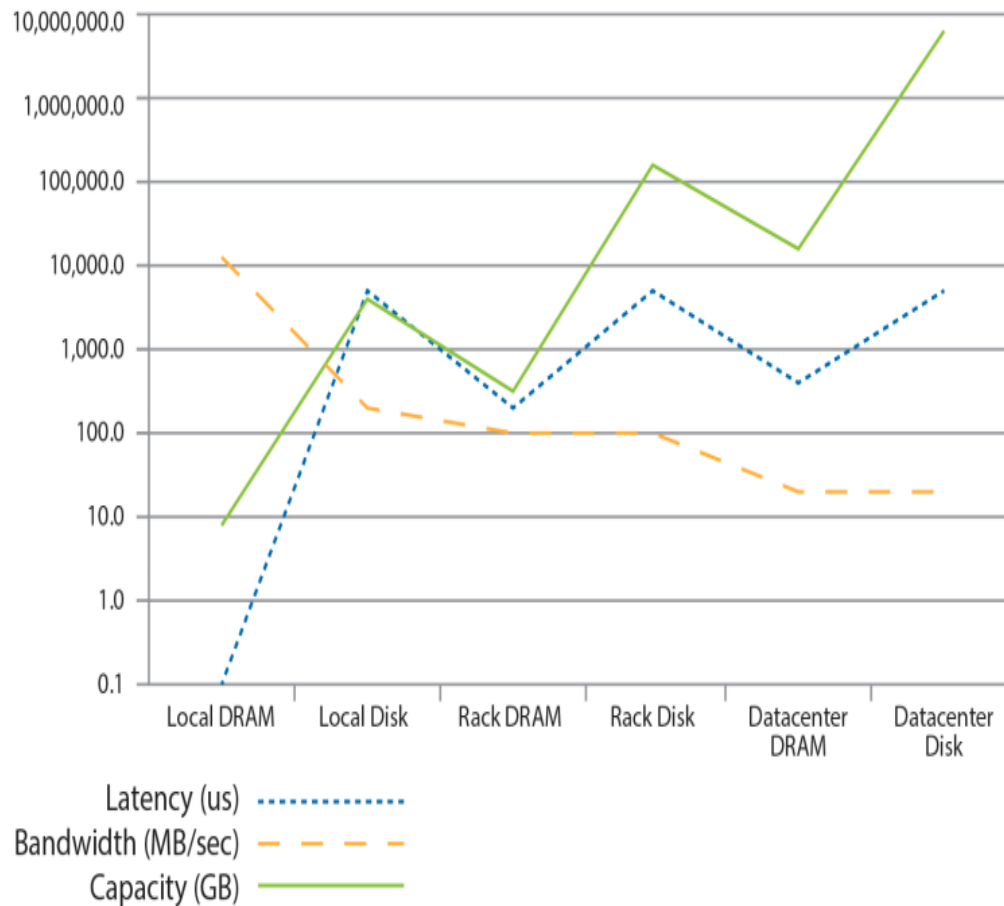
Why the latency of Rack DRAM is much higher than that of Local DRAM?



Answer: Networking.

Question 1b

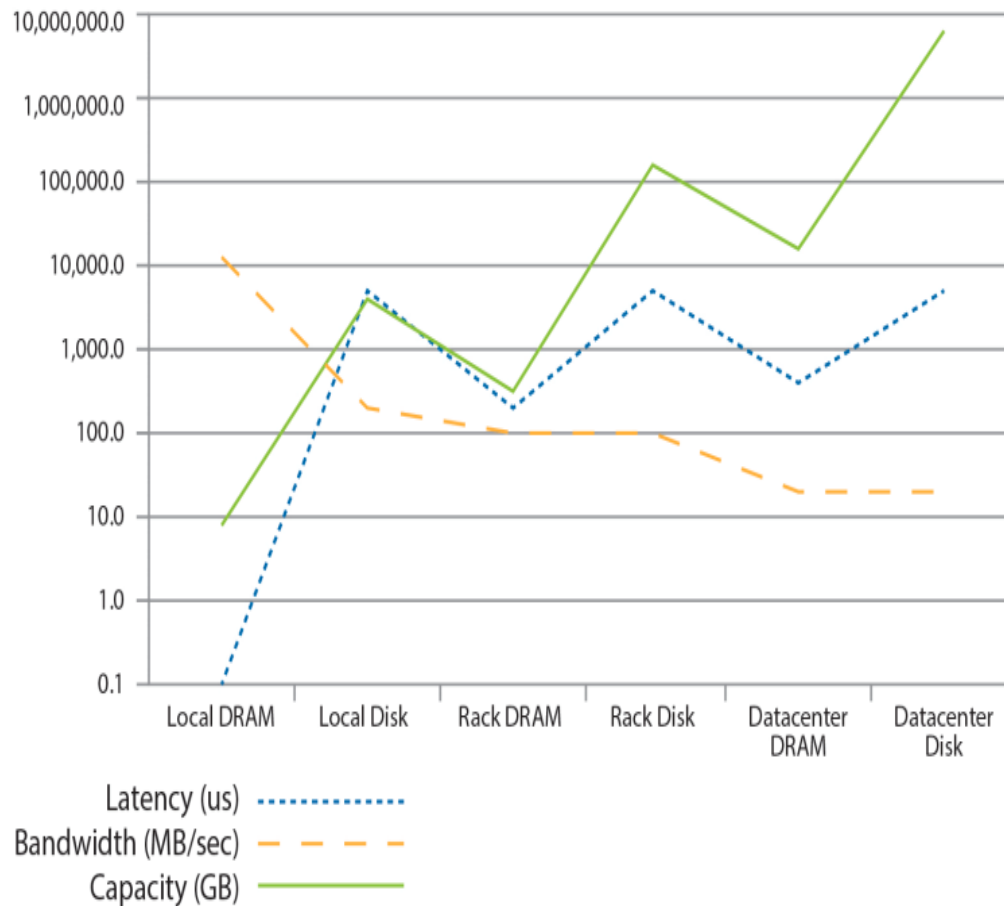
Why the latency of Rack Disk is similar to that of local Disk?



Answer: The network in the Rack is relatively fast, and the disk performance becomes the bottleneck.

Question 1c

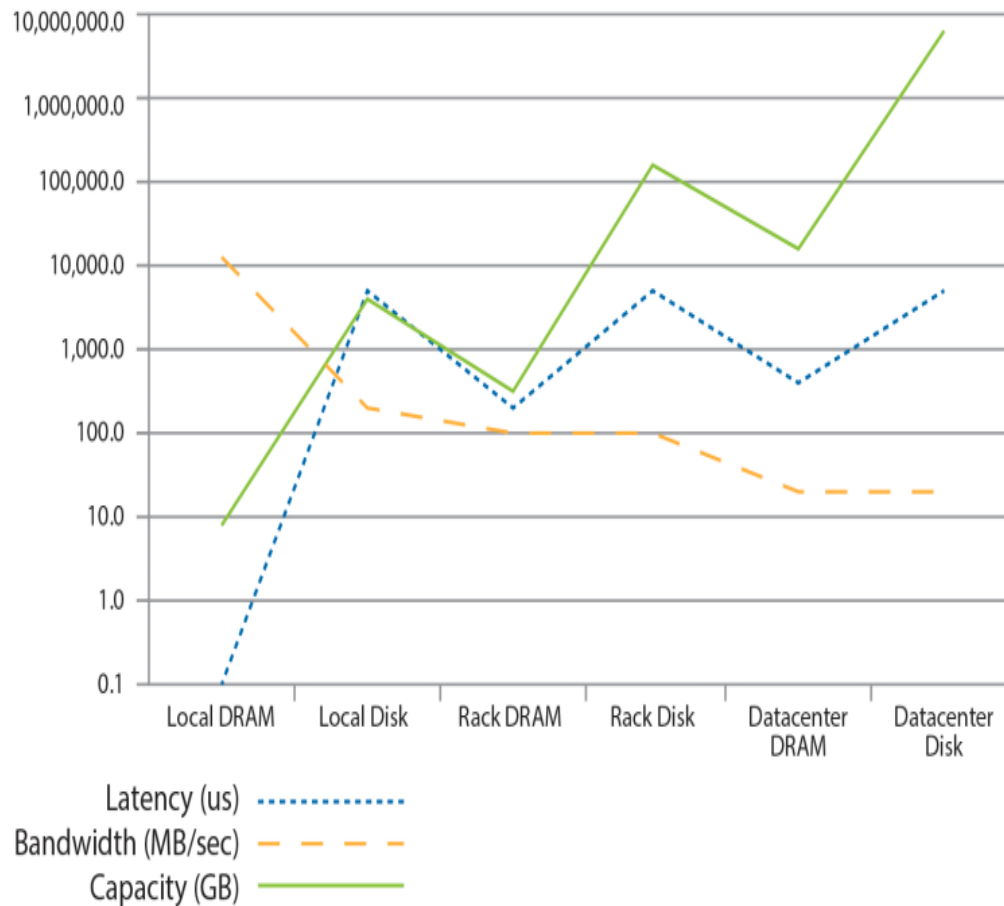
If we have an application with the working set smaller than Local DRAM, should we always put the working set into a single machine for efficiency



Answer: No. Depending on the gain of performing parallel computation on multiple machines.

Question 1d

Should we always put the data into DRAM?



Answer: No. DRAM is volatile (the data is off when power off), and disk can be used for persistency, and reliability execution.

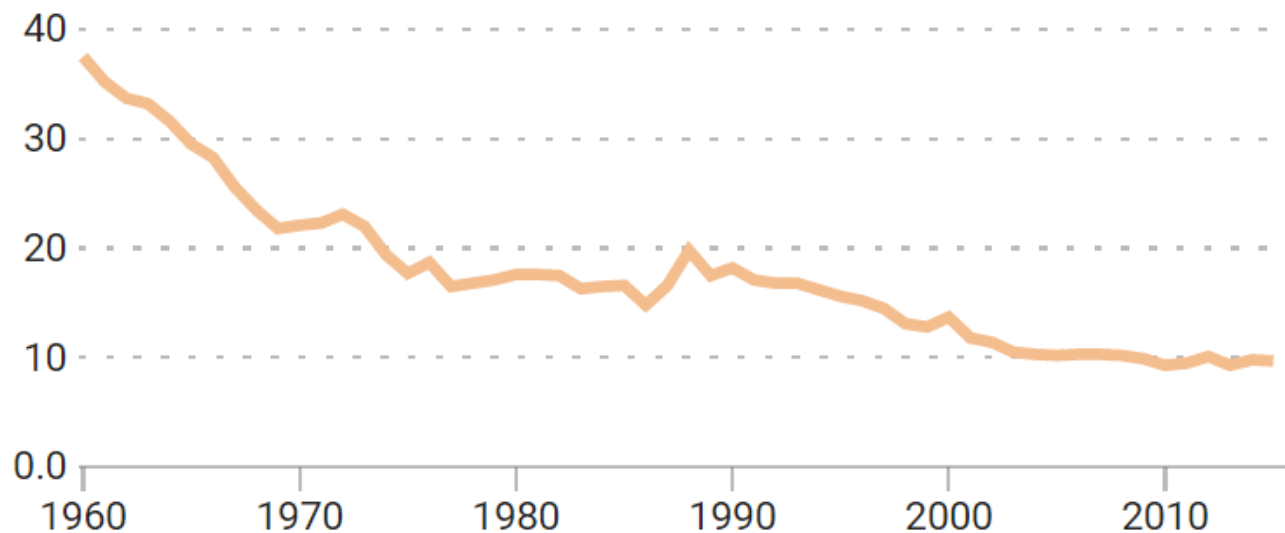
Question 2

2. "A picture is worth a thousand words". An important aspect of big data is the presentation of the answers. In this question, we examine a number of plots from <https://data.gov.sg/>, which is an effort of public data analytics in Singapore. For each plot, write down your findings, share/discuss your findings with peer classmates and try to dig out more insights from discussion.

- a) Crude birth rate**
- b) GDP at current market prices**
- c) Graduates from university first degree courses**
- d) Rainfall - monthly total**

9.7

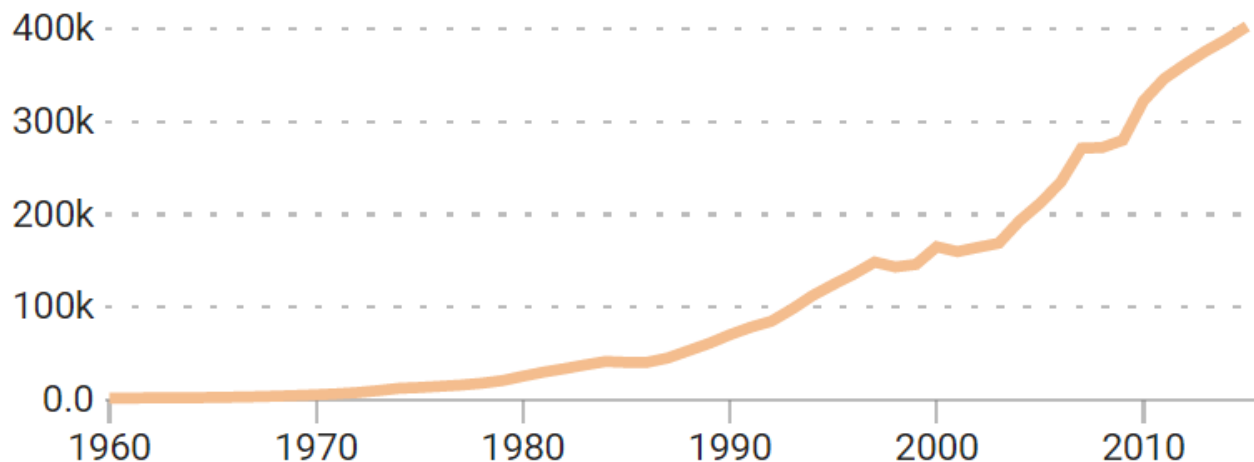
Per Thousand Population in 2015



Answer: The birth rate has been going down since 1960. Now it is becoming very stably low.

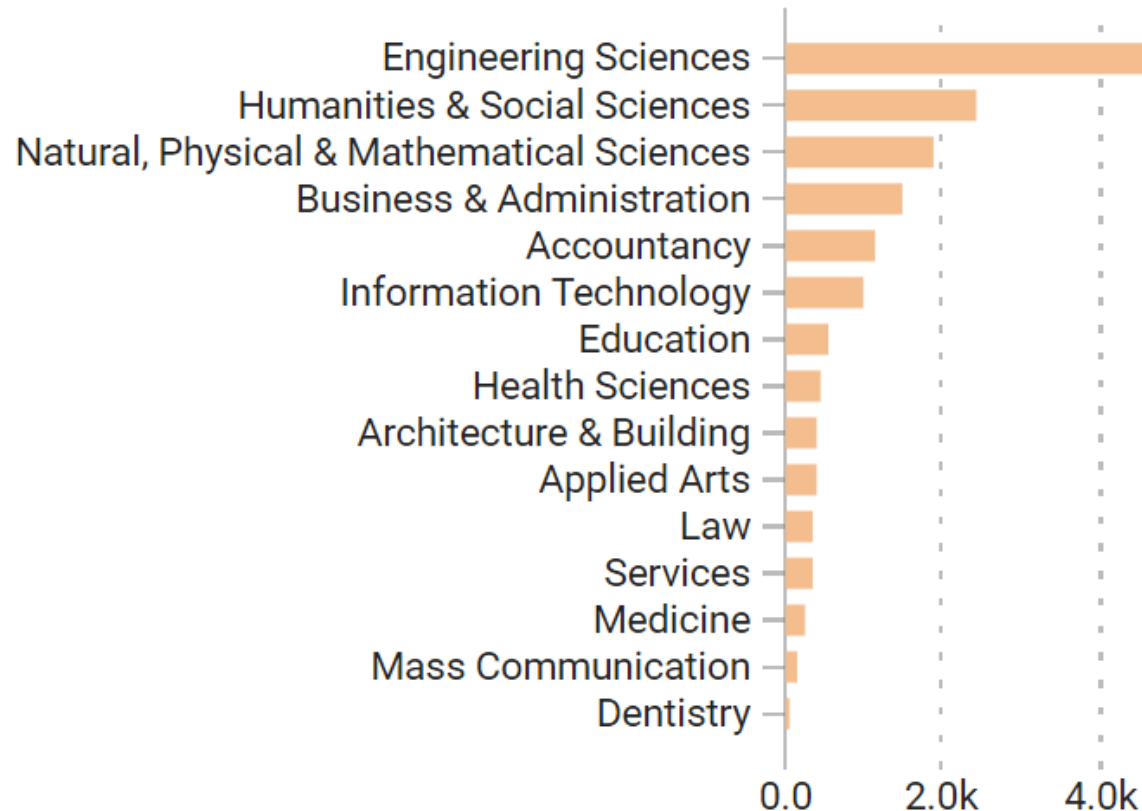
s\$402,457.9

Million in 2015



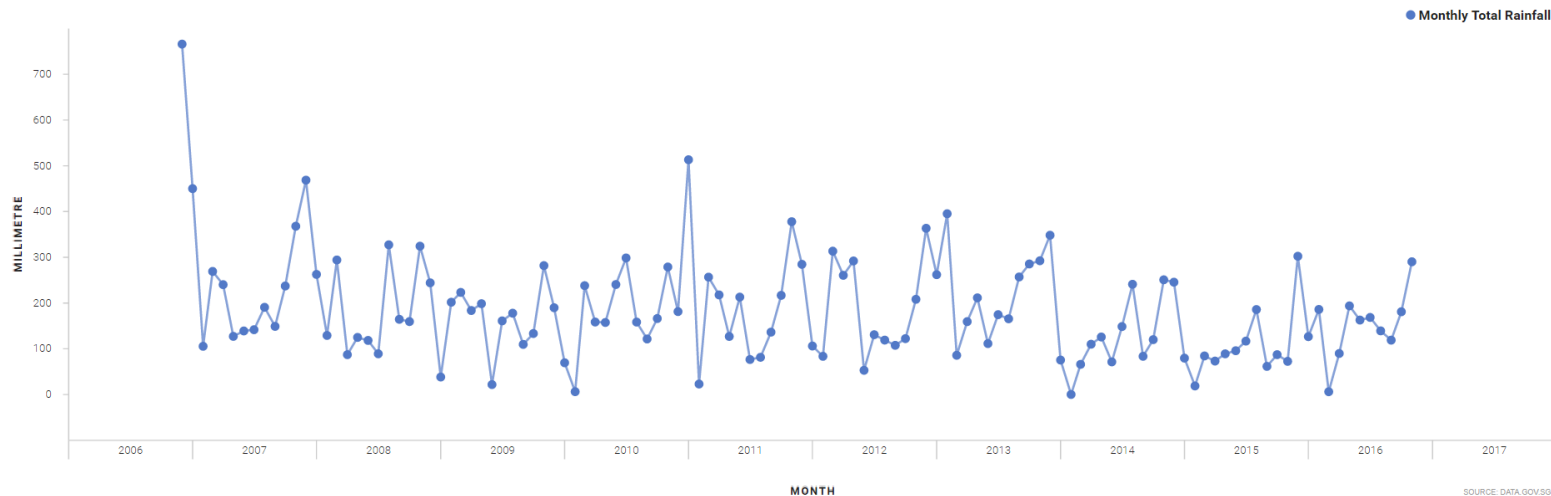
Answer: The GDP has been growing really fast. The latter years grows even at a higher pace.

Graduates From University First Degree Courses (2014) - Data



Answer: Long tails. The head is on engineering.





Answer: Within a year, the monthly rainfall fluctuates dramatically. Across years, the monthly rainfall also fluctuates quite significantly on the same month.

More Exercise Questions

Index	Question
1	<p>Which challenges are important for gene sequencing in precision medicine?</p> <ul style="list-style-type: none">a) Volumeb) Velocityc) Veracityd) Variety
2	<p>When the number of tasks in a job is small, we should move the data to the task, rather than moving task to the data.</p> <ul style="list-style-type: none">a) Trueb) False
3	<p>The major advantage of MapReduce is on its scalability and programmability. But, we may come out with non-optimal solutions for an application by using MapReduce programming.</p> <ul style="list-style-type: none">a) Trueb) False

Index	Question
1	<p>Which challenges are important for gene sequencing in precision medicine?</p> <ul style="list-style-type: none">a) Volumeb) Velocityc) Veracityd) Variety
2	<p>When the number of tasks in a job is small, we should move the data to the task, rather than moving task to the data.</p> <ul style="list-style-type: none">a) Trueb) False
3	<p>The major advantage of MapReduce is on its scalability and programmability. But, we may come out with non-optimal solutions for an application by using MapReduce programming.</p> <ul style="list-style-type: none">a) Trueb) False

Latency Numbers Every Programmer Should Know



https://people.eecs.berkeley.edu/~rcs/research/interactive_latency.html

Acknowledgement



Thanks to Li Qinbin for making these slides.

liqinbin@u.nus.edu