

Guest Lecture 1: February 20

*Lecturer: Ahmed Ali-Eldin**Scribe: Jonathan Westin 2019, Daniel Thiyagu 2018*

This lecture will cover a high level of how data centers are designed and used. Datacenters, in general, are used for internet services, batch jobs and AI training (to mention a few)

1.1 Traditional vs modern

Traditional Data Center:

- Sys admins monitor and manage servers.
- They schedule a process to be run.
- Applications run on physical servers.
- NAS/SAN

Modern Data Center:

- Dynamic Larger Scale
- They transfer process to get better efficiency.
- Applications run on virtual machines.
- Increased automation allows larger scales.

The traditional datacenters use a static approach where system administrators would monitor and manage the servers. The server would often have the framework approach (such as xterminal) with connected Storage Array Networks (SAN) or Network Attached Storage (NAS).

The modern approach is a virtualization approach where it is easy to scale and offers other dynamic trade-offs. This cuts down on the workforce needed as seen in the traditional approach where data center needs more system administrators. In the modern approach, it is possible to focus on having a seasonal workforce that can handle applications (; clients take care of their own applications due to the hardware part is virtualized with flexible mappings). This allows bigger data centers, for instance since 2014, Amazon has between 50k-80k servers per data centers. And Facebook has a data center in northern Sweden where they only have 19 system administrators.

Resource Management:

- We need to keep it as highly utilized as possible
- Apps have a variable/unpredictable workload.
- Want High Performance and Low Cost
- Automated Resource management

1.2 Inside the data center

The data centers are filled with racks and storage arrays. Giving it a modular approach. This makes it more dynamic and possible for bigger bulks of hardware. Like GPU programmable arrays, tensors (Tensorflow) or other components.

Google bought an old watermill in Finland to make it possible to use the cold weather but also hydro-cooling for heat generated by the data center. Another factor in places like north Scandinavia is that they have an excess of hydroelectricity, making it cheap but also green.

High Performance Computing:

- Generally managed and used by the scientific community.
- They focus on high computational workloads.
- The types of machines are Highly Parallel Mainframes.
- Example: Used in weather forecasting.

Data Center:

- Generally managed and associated with enterprises.
- They focus on collecting data.
- They don't focus on high computational workloads.

Question: What is the difference between supercomputers vs cloud computers?

Answer: Supercomputer solves problems that take a long time and can be batched, they will also be optimized for supercomputers. The optimization is totally different. Cloud does have spikes while supercomputers don't. Well, not entirely true, supercomputer show spikes before conferences and such. Also, supercomputers doesn't need virtualization.

1.3 Modular Data Centers

The technology has made it possible for a more modular approach. One example of this is big companies can buy ship containers and have a "plug and play" to any data center they need. This makes data center easy to geographically moved or expandable. The main concern will be the electricity when installing whole ship containers of machines.

Question: There exist data center in the sea, how is this maintained?

Answer: This is still a research project where using water as cooling for the data center. But all communication and such run through the Atlantic broadband (Transatlantic communications cable) where the container is also stored, and the Atlantic broadband is also maintained.

Question: How does the shipping container "plug and play" really work?

Answer: Do not think of shipping containers "as is". It will need cooling etc. It is not as simple as putting down the container and it works, there will need some installment that is taken care of personnel, a data center isn't solely reliable of system administrators, it needs other types of engineers also.

1.4 Server Virtualization

Virtual Servers:

- Balance/consolidate load
- Faster Deployment
- Easy Maintenance

Today's data centers have virtualization at heart. They make it possible to have such big data centers. For instance, Akamai has between 15% and 30% of all web traffic, and this will take a big data center to handle. Therefore it is necessary to be able to use virtualization to be able to install and migrate VMs. There is much research in this area, Amazon rolled out the project called Firecracker where 125 VMs can be started in a second. Virtualization makes it easier to install a new machine and get it running, you don't need to have a USB stick to put into the machine to install the server. Virtualization gives isolation at low cost.

If you were to create a new game as Pokemon go, or you are the US president and you roll out affordable health care act where you can access it through a web platform for health insurance. How many servers should you buy? In both these cases, someone miscalculated the demand for the servers and it resulted in a non-working environment where customers couldn't access the platform. There are multiple of these examples. So if you have the possibilities of having virtual servers where you can deploy them fast and conciliate them.

Virtual servers make server consolidation¹ since we can have several virtual servers on the same host.

Another way people have used virtualization is having a big company (like UMass) to use remote desktops as your workstation and give the personnel a screen and a keyboard. This also makes the workstation accessible anywhere and it is easy to maintain while having reduced cost.

1.5 Data center challenges

One of the challenges for big data centers is that you're turning multiple servers into one big server in essence. How does scheduling work? How does load balancing work? Where should I put things? And now you have to face the issue of heat/cooling problem. So you will have the issue of having all your having a heavy computational load in the same area resulting in too much heat resulting in cascading failures. You also create the problem of having access to lots of parts that may result in catastrophic errors by humans. How do you manage these resources? How do you trust your system administrator with all the privileges that come with the role? This is the problem of consolidation.

Automated resource management and performance prediction/profiling help with this. For instance, Netflix will probably be used in the evening more general since people are working in the day (in general). But there are also unpredictable peaks, one example is when Michael Jackson died Google thought that someone was DDOS'ing their web services because it got so many requests that it more or less made "internet break" from Google's point of view.

Question: *The slides show the problems for the first year in a Google data center, how true are they?*

¹Server consolidation - efficient use of server resources in order to reduce the total number of servers.

Answer: There is a strong culture of secrecy and biased. It is hard to get a full picture from a company that actually earns money from having working data centers.

1.6 Reliability Challenges

More servers imply more challenges, all the errors possible becomes a problem. The prediction of having 80k machines is a 100% rate that you will have a failure. You also get new problems, like rats since they love cables. So when they eat your wire you do not only have to find it, change it, pest control, it can be troublesome. There are several failure possibilities and fixes takes time (see slide),

Reliability challenges:

- 0.5% overhear
- Power Distribution unit failures
- Network Rewiring - as you add/remove systems, etc.

1.7 Data Center Cost & Economy at scale

- 5-10 million dollars a month.

Effectiveness is often calculated as Power Usage Effectiveness; how much does each computational cycle cost.

$$PowerUsageEffectiveness = \frac{ITPower}{TotalPower}$$

Bulk buying equipment often reduces cost. Also able to get cheaper electricity rate. And automation allows a smaller number of system administrators. This gives a trend towards mega data centers (> 100k servers).

The data centers have three considerations for electricity:

- Cooling
- Power converters
- Backup generators

Location impacts:

- presence of customers nearby
- cooling
- Generation of power.
- presence of already existing grid power infrastructure.

Energy Efficiency:

- Servers consume a lot of energy

- Be Green
- Save Money

In general, it takes 50% power to run the computers and the remaining power to cool them.

1.8 The cloud

What is Cloud? It is remote, you pay as you go, we get high scalability, shared infrastructure.

IaaS Infrastructure as a Service: Eg: Google, Aws

PaaS Platform as a Service: Eg: Azure, Google App Engine

SaaS Software as a Service: Eg: Applications like Salesforce, Gmail.

Hybrid Cloud: It is a mix of private and public cloud usage.

Programming Models:

- Client Server(Interactive)
- Batch Processing(Not Interactive)
- Map reduce.(Not Interactive)

Future Challenges:

- Privacy/Security
- Extreme Scalability
- Programming Models