

# An end-to-end Spark based data stack in the hybrid cloud

Farhan Abrol  
Product Lead, Pure Storage

[fabrol92@gmail.com](mailto:fabrol92@gmail.com)  
@F\_Abrol  
[www.linkedin.com/in/fabrol](http://www.linkedin.com/in/fabrol)

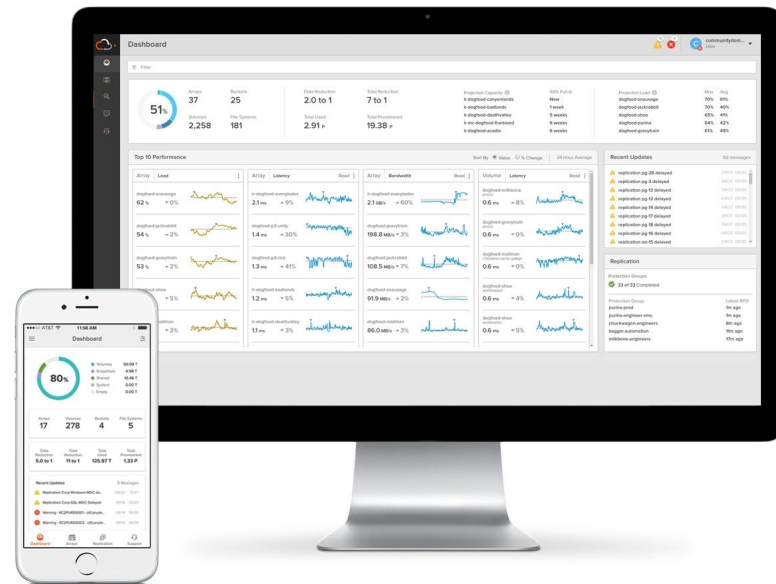
**#HWCSAIS12**

# Outline

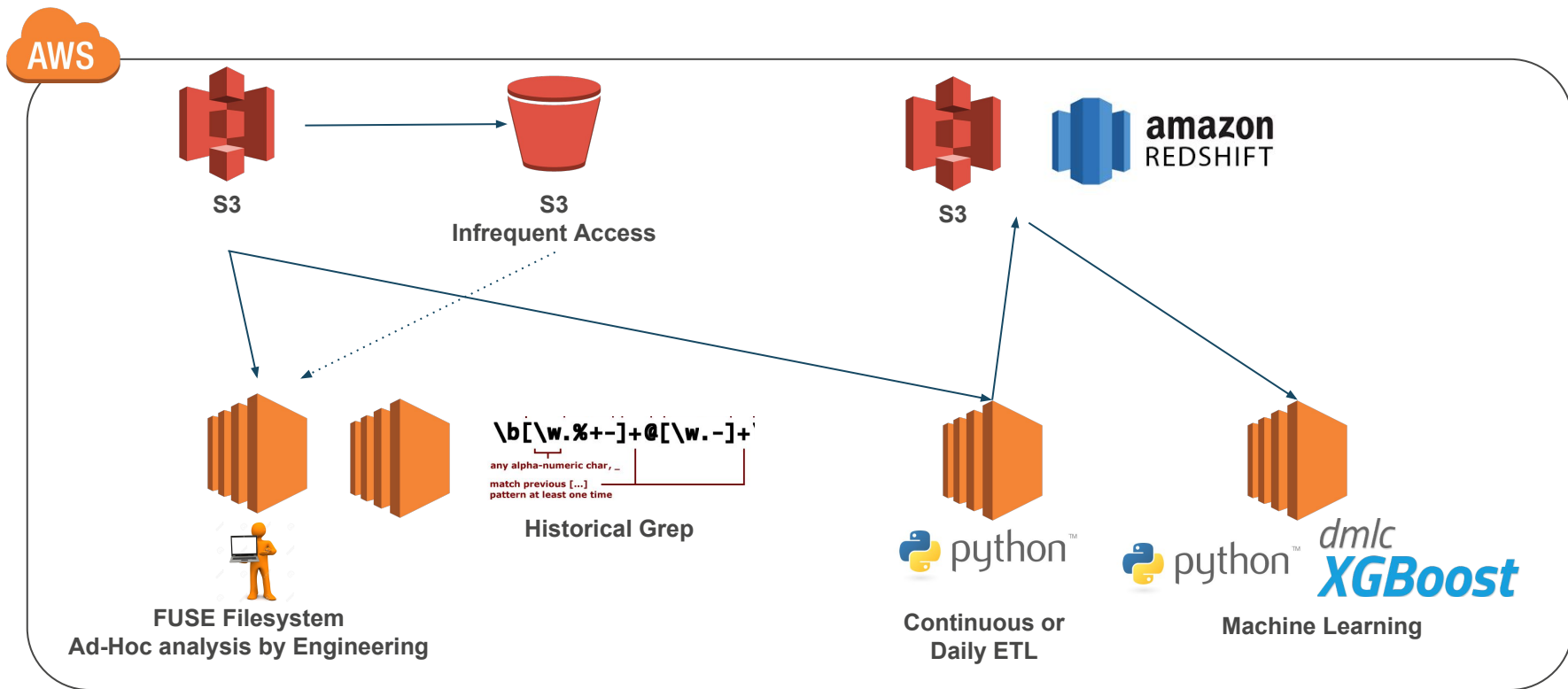
- Environment overview & problems
- Solutions - Hint : Spark
- More Spark More Problems
- Hybrid Cloud
  - Options & Performance comparison
  - Should you do it ?
  - Basics of datacenter

# Pure1

- Fleet dashboard for IoT devices
  - Storage arrays
  - VM's
- Real-time log/metric streaming
- 16 TB logs/metrics ingested daily
- Intelligence
  - Proactive scanning for issues
  - Predictive alerting
  - Machine learned forecasting



# Logs are king



# Problems

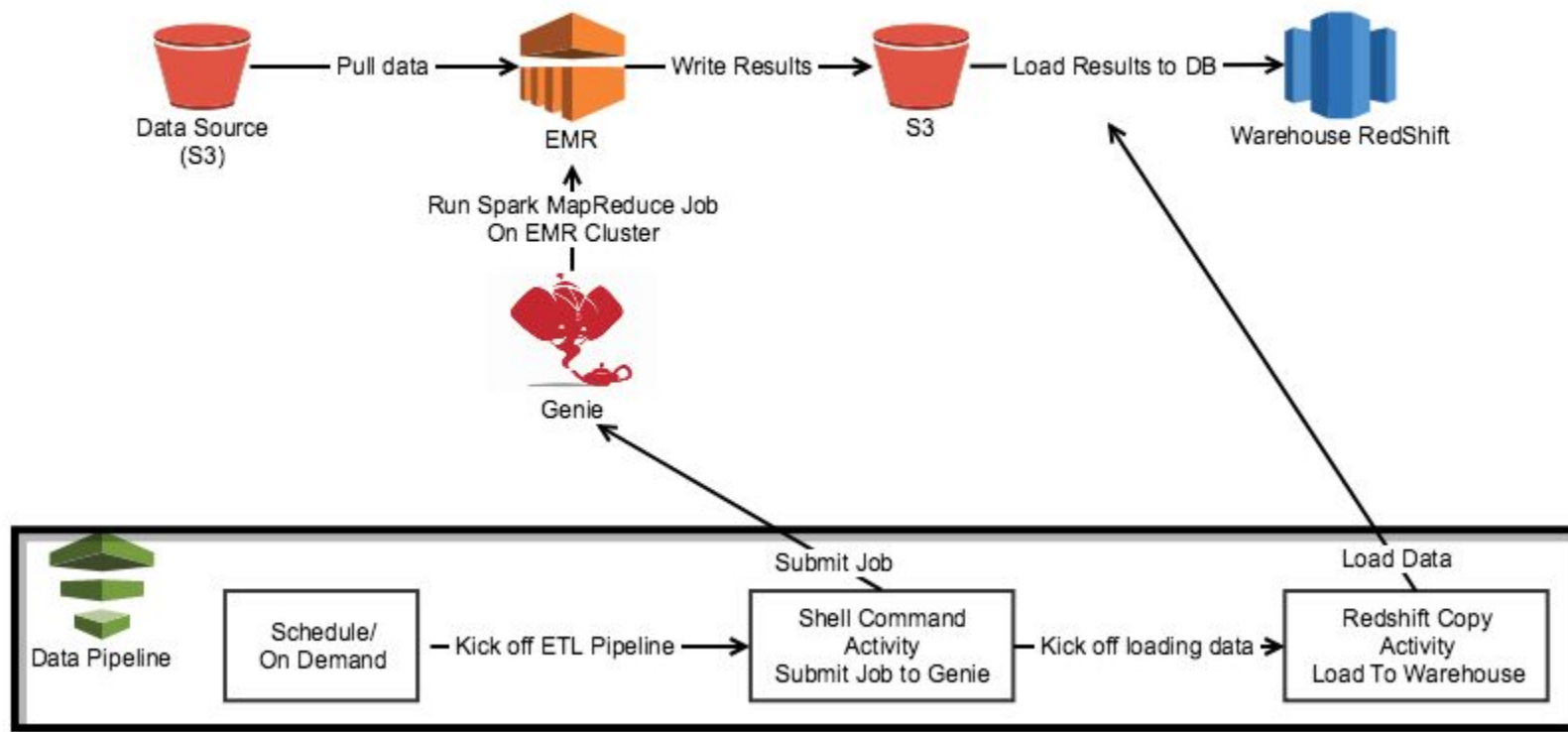
- Speed of running historical greps
  - Bottlenecked on single machine throughput
- Resource wastage for ETL machines
- Code/maintenance for new ETL jobs
  - Becoming a monolith
- ML training time
  - As data grows, taking 8-12 hours

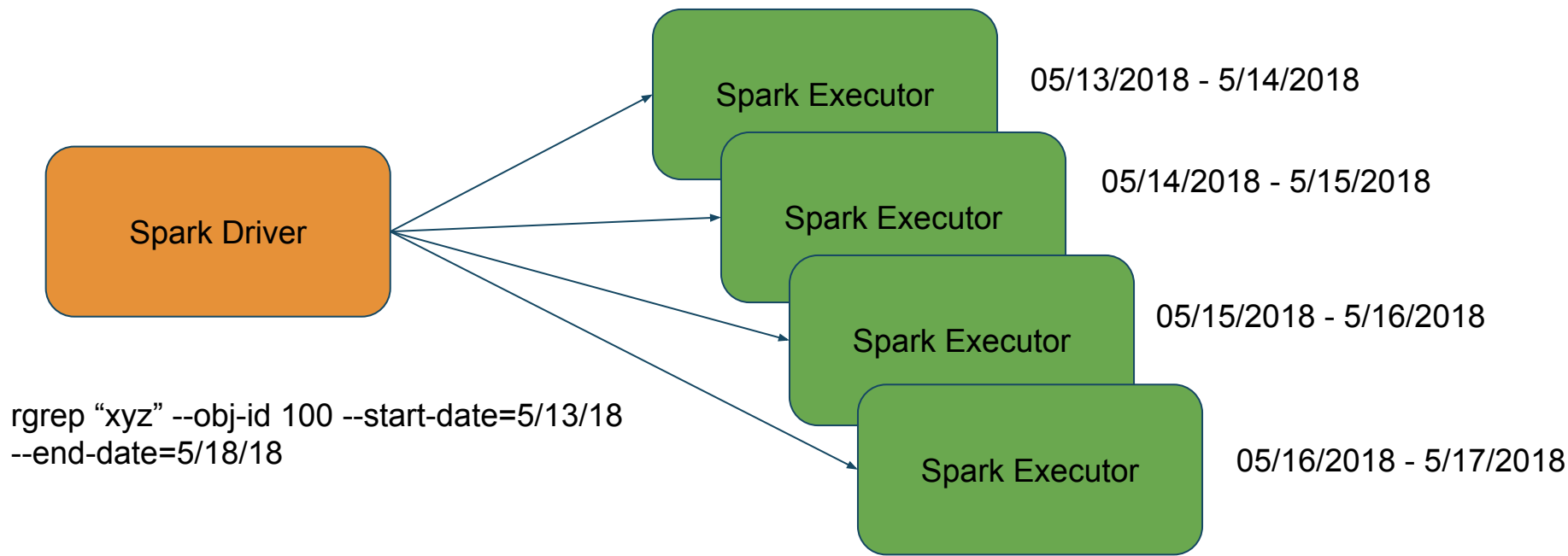


all the things !

- Faster\*
- Better resource utilization
- Uniform language and tooling
- Streaming / batch jobs
- One infra to maintain

## Automated ETL Workflow with AWS DATA Pipeline



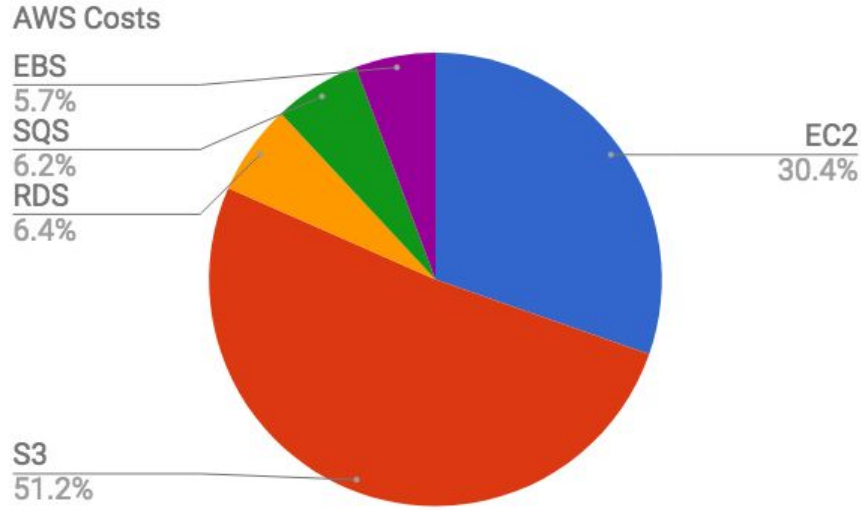


**Grep -> Distributed grep on Spark**

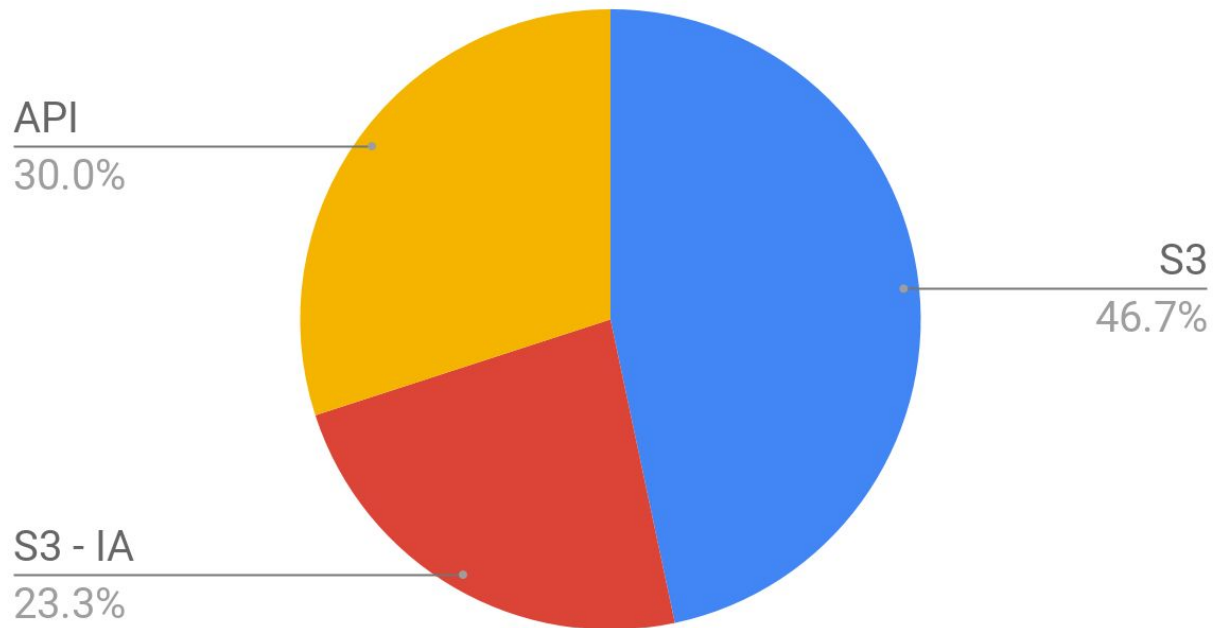


Done !

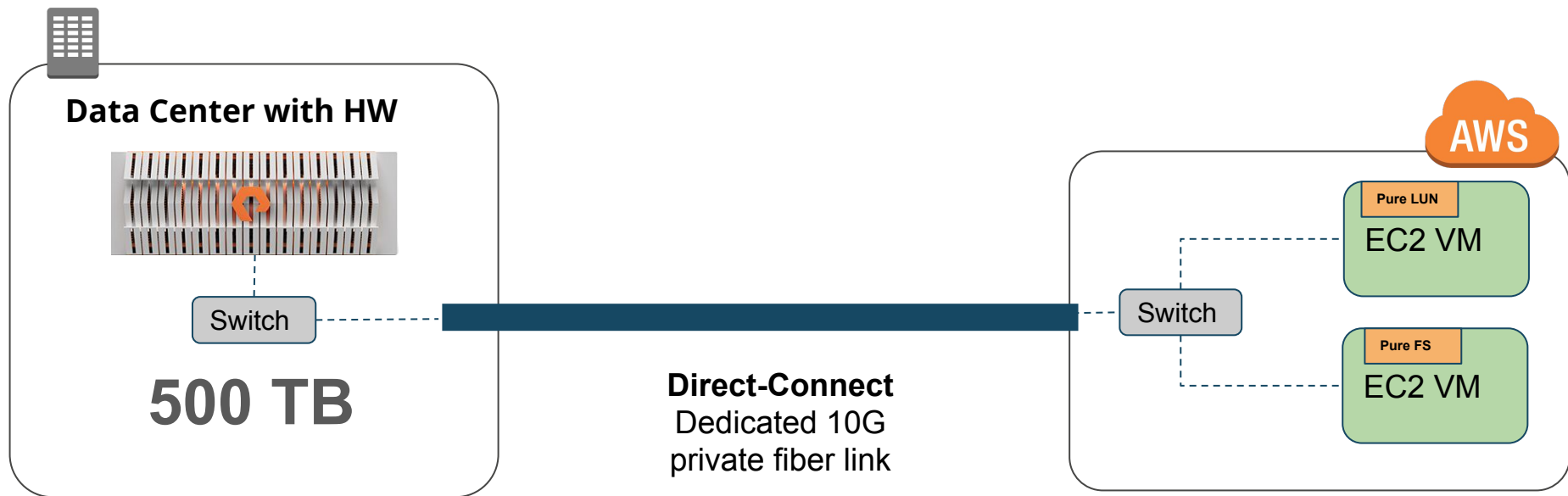
# Problem - AWS Cost trend



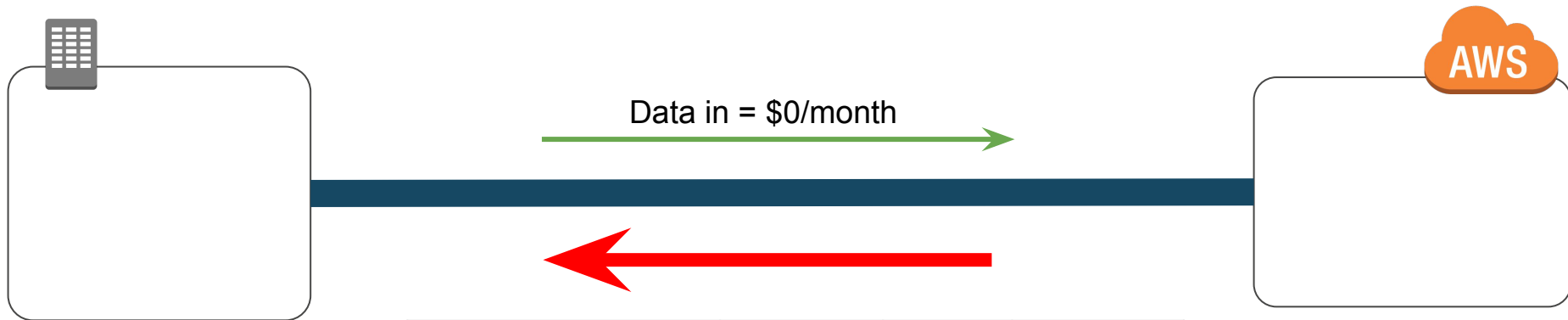
## AWS Storage Costs



# Hybrid Cloud

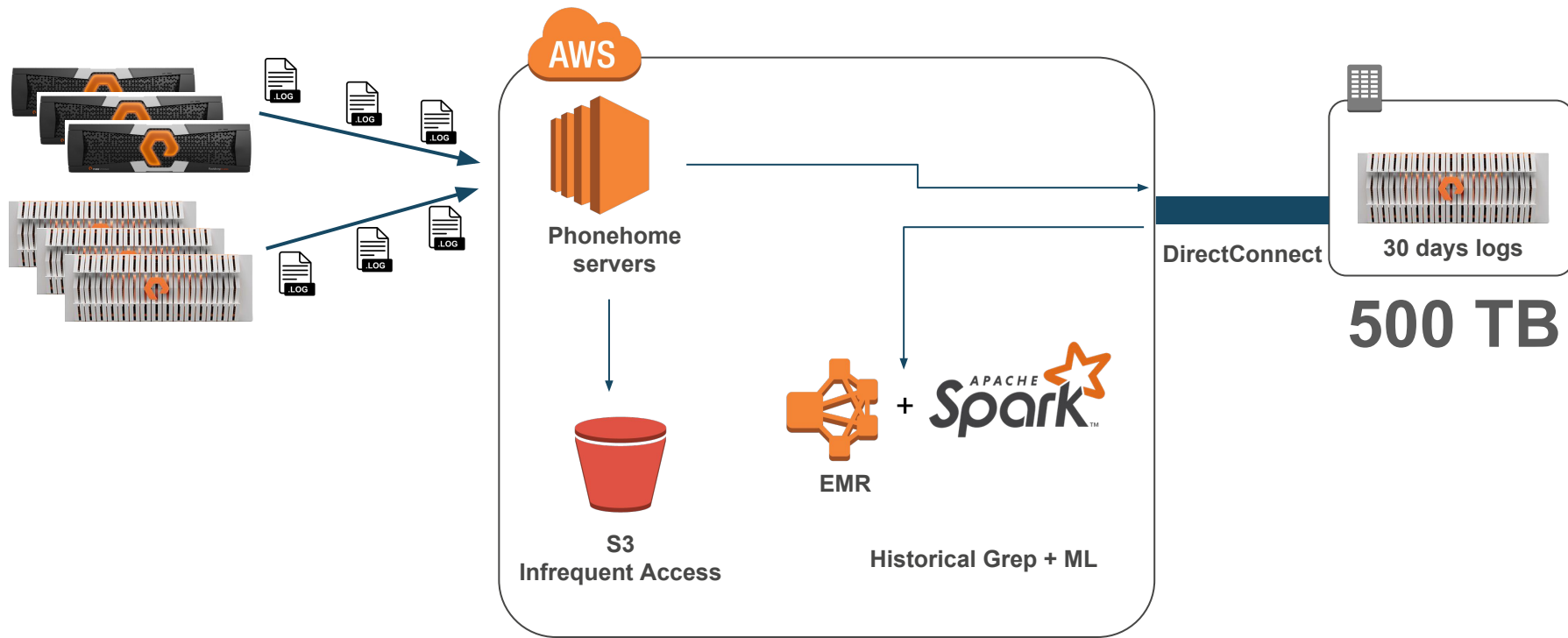


# Hybrid Cloud - Pricing



Utility	Price	Usage	Total per month
10G port	\$2.25/hr	720 hr	\$1620
Data transfer out of AWS	\$0.020/GB	500 TB	\$10000
AWS Cost			\$11620

# Log analysis pipeline - Smoke test



# Aside Storage Protocols

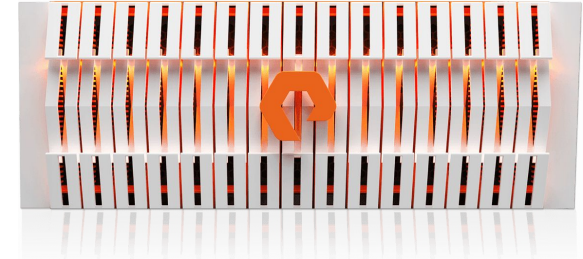
## Storage system



### Generic

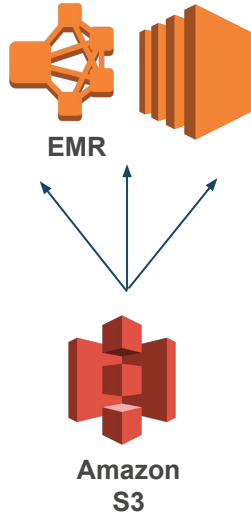


### Optimized

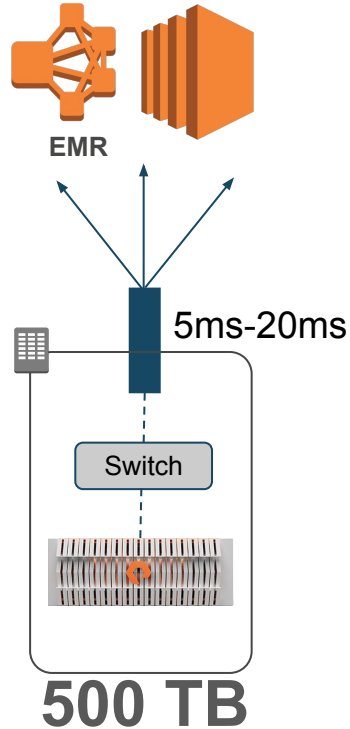


Flashblade

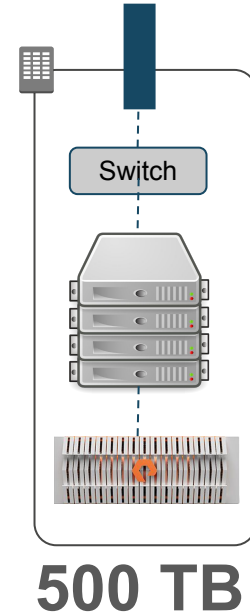
## AWS Only



## Hybrid with EC2



## Hybrid with Local Compute



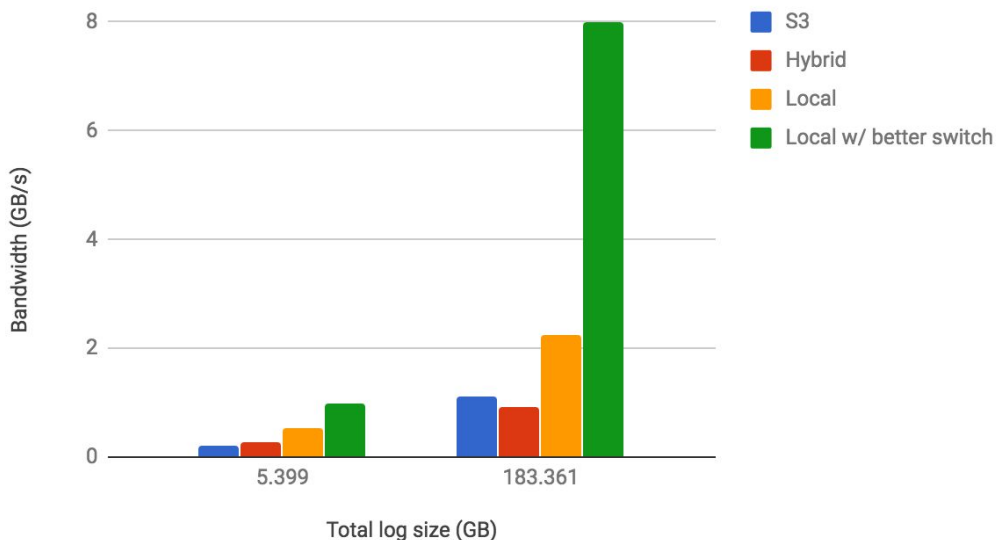


# 144 node spark cluster

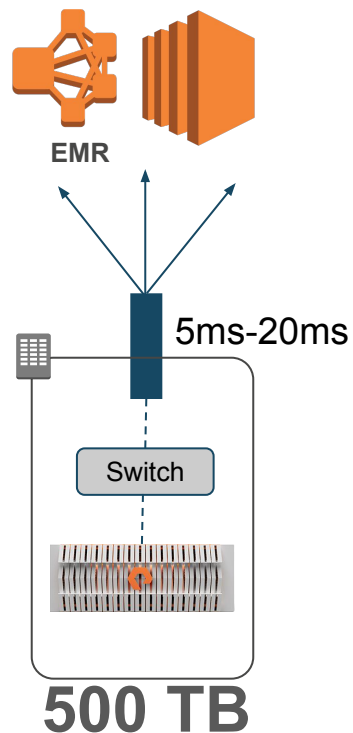
## Workload - Distributed grep

~3x-10x better throughput

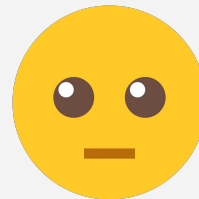
Read throughput comparison



## Hybrid with EC2



Performance



- Link latency
- Cloud networking stack

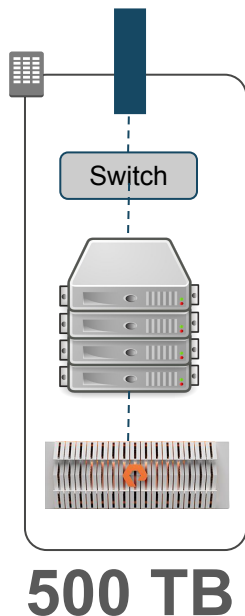
Costs



Good for

- Read heavy workloads
- Latency insensitive workloads
- Low Bandwidth workloads

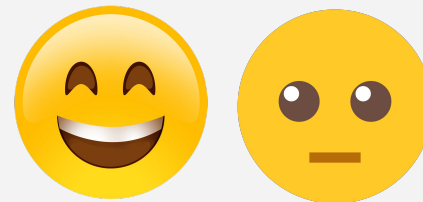
## Hybrid with Local Compute



Performance



Costs



Good for

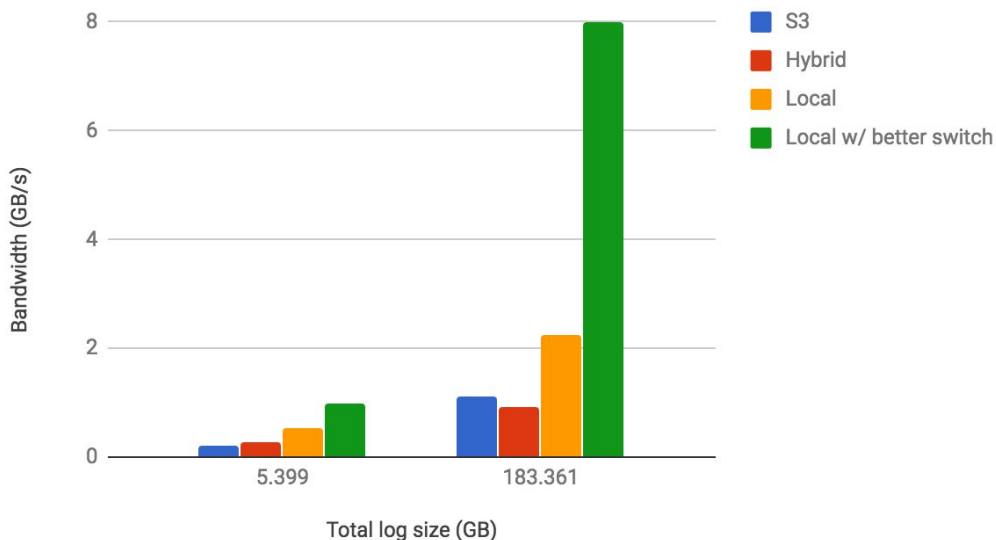
- Read heavy workloads
- Latency sensitive workloads
- High bandwidth workloads

# 144 node spark cluster

## Workload - Distributed grep

~3x-10x better throughput

Read throughput comparison



# Datacenter setup



# Conclusion

- Best use cases: Workloads with **higher read, lower write** requirements
- When write portion of read/write ratio increases, be cognizant of cumulative **AWS transfer costs**
- **High performance cloud services** can be expensive, on-prem can alleviate this cost
- Unique capabilities of on-prem storage & compute:
  - **Instant snapshots**
  - **All kind of workloads on one platform**
  - **Resilience**

