# Big Data Systems

- Big Data Technology
  - Previously:

Application → Database

*MVC

- Issues
  - DB can't handle > $x$ writes/sec

App
Queue ← Buffer
Workers
Failures
Recovery Mgr →
Database $$$$ ← Scale-Up (Vertical)

* Sharding (DB-Level)

DB $$$$ = 
A-F DB    G-M DB
DB        DB
N-R       S-Z

⤷ Recovery Mgr for DB

- System Rethinking
  - Fault-Tolerance
  - Low-Latency (IOPS)
  - High-Throughput (IOPS)
  - Scalability (Infinite)
    - Compute ⎤
    - Storage ⎦ HW / Cloud
  - Flexibility
    - Varietys (Unstructured)
    - Interactions (Queries)
    - Change Mgmt (Schema Drift)
  - Extensibility
  - Ease of Use


- Big Data Attributes
  - Volume (Throughput)
    - Size Scale: GB, TB, PB ....
    - "Data at Rest"
    - Storage Tech Limitations
      - Size of disk
      - Capacity per Server
      - Storage I/O

  - Velocity (Latency)
    - "Data in Motion"
    - Stream Scale: 1d, 1h, 1m, 1s, <1s

- Stream Scale: 1d, 1h, 1m, 1s, < 1s
- Compute Tech Limitations
  · Speed of GPU/CPU/TPU...
  · Interconnect speed
  · Memory speed/capacity

· Variety
  - Unstructured Data
  - "Schemaless"
  - Significant storage/compute overhead

  ＊ Structured Data: Schema → Fields/Types
  ＊ Semi-Structured : Common Case

- File Systems (Storage)
  · Tied to running OS (POSIX)     ＊NFS
  · Coupled to specific computer/hardware
  · Filesystems over Network → Internet
    - HTTP : GET/PUT
              read   write

- Object Store
  · Independent of individual machines/architectures
  · Independent of client software/platform

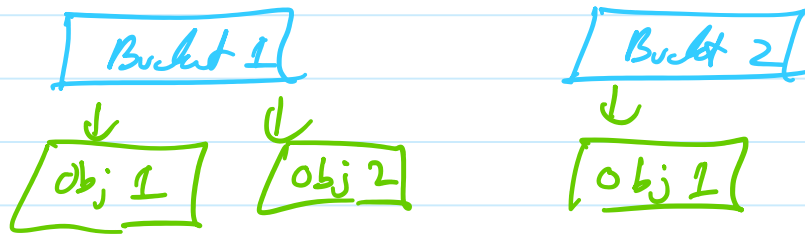  ＊ Accessible via Web tech → REST

  ＊ AWS: S3

- Simple Storage Service
- Scalable
- Integrates w/ other Aws services

* Ceph
* Minio

- Concept
  • Buckets
  • Objects
  • keys
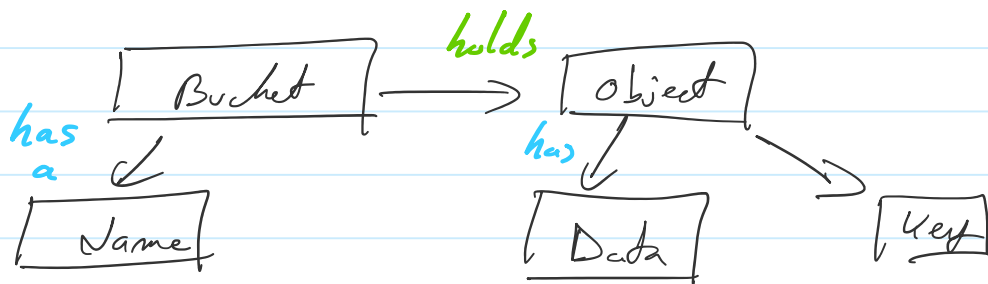
Global Uniqueness

Bucket Name + Object Name (key)

Bucket 1          Bucket 2

Obj 1    Obj 2     Obj 1

* File ⟷ Object (key)
  ↳ WORM

has a ⟶ Bucket ⟶ holds ⟶ Object
         │                   has ↓        ↘
        Name                Data          Key

mybucket / hello.txt
mybucket / temp / hello.txt