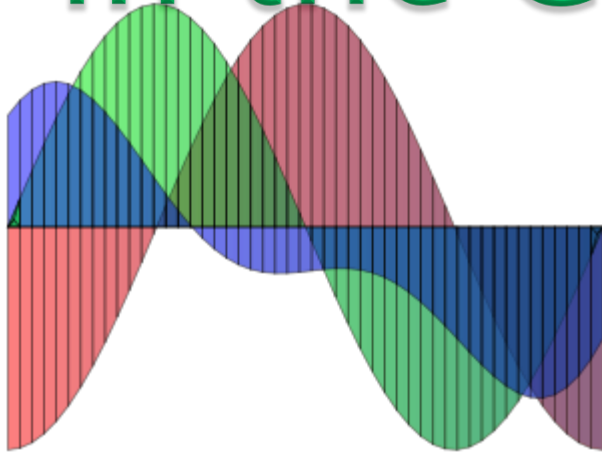


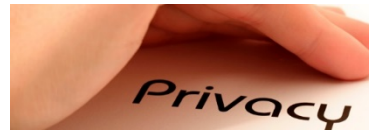
# Managing Large Scale Transactional Data in The Cloud

**Divykant Agrawal and Amr El Abbadi**  
**University of California, Santa Barbara**

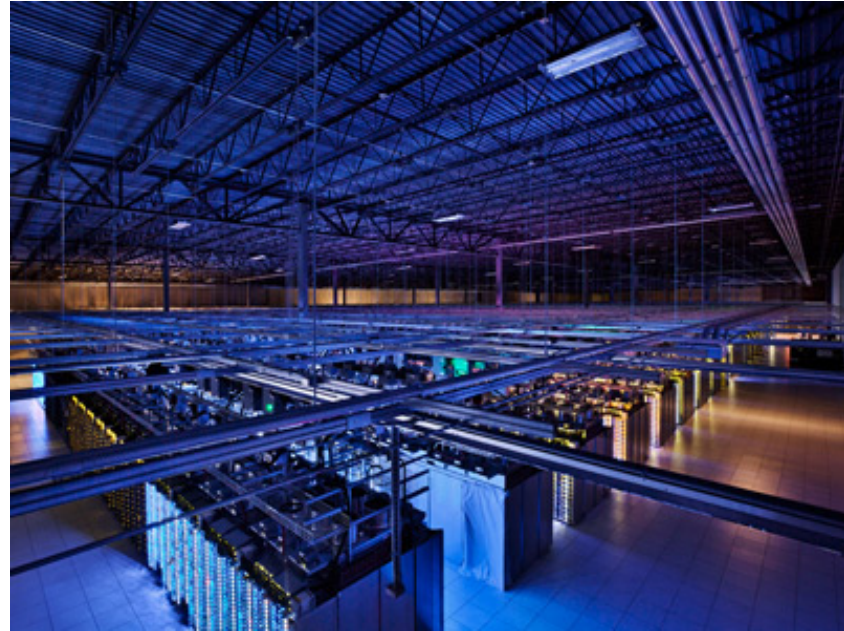
# The Big Data Eco-System in the Cloud



Infrastructure



# Inside a Data Center



# MaaT: Distributed Transaction Processing



Pessimistic locking

Less **aborts** compared to OCC

**Blocking** limits throughput



Optimistic Concurrency Control  
(OCC)

High **Throughput**

**More aborts** with contention

# Maat Design Principles

High Throughput

Conflict resolution at **fine granularity**

**Avoid blocking** transactions

Resolve conflict with **less aborts**

Scalability

**Distributed** verification

Only involve the **nodes accessed**



# Data Center Data Management: MaaT



# Catastrophic Failures: Geo-Replication



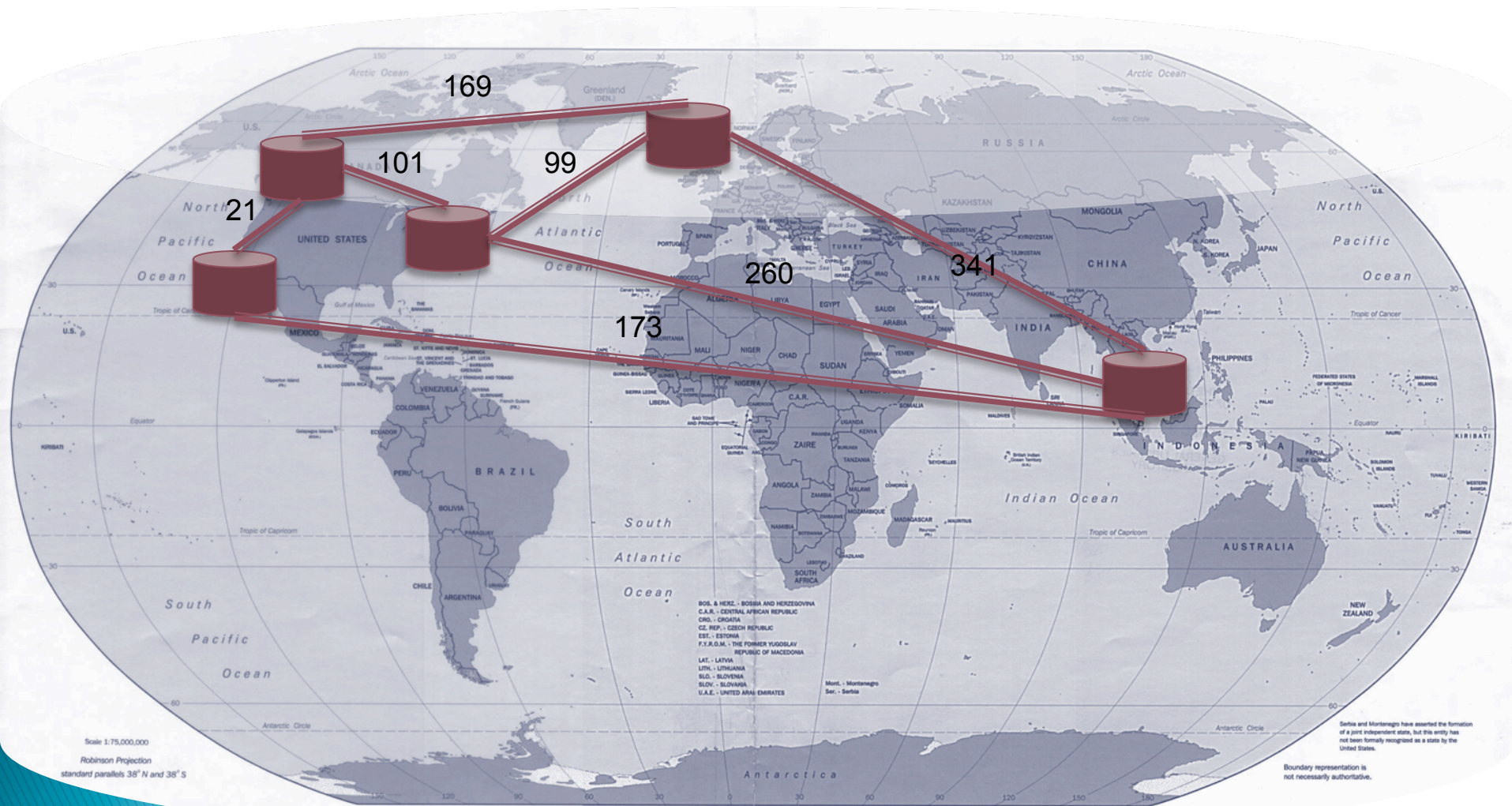
# Google Spanner

- ▶ **Global-scale** data infrastructure
- ▶ Data is **partitioned** within data center
- ▶ **Replication** across data centers using **Paxos**
- ▶ **Transactions** execute on data using **2-phase commit**



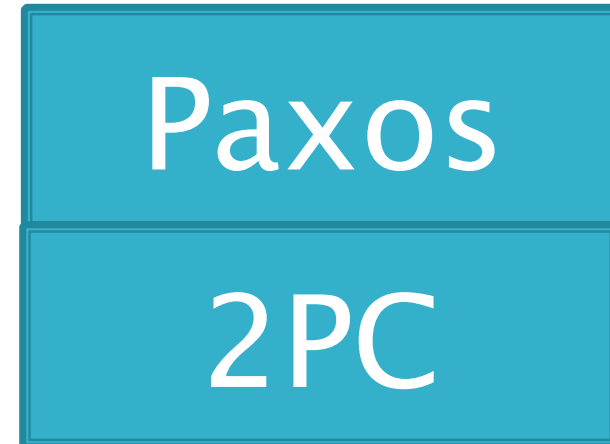


# Communication Overhead

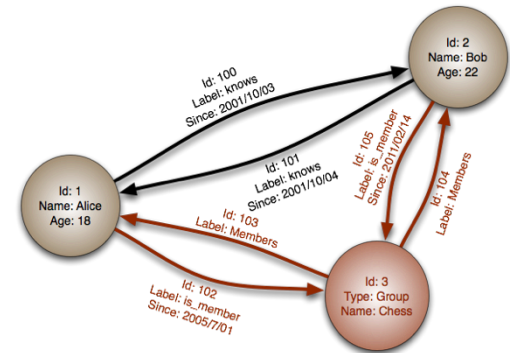
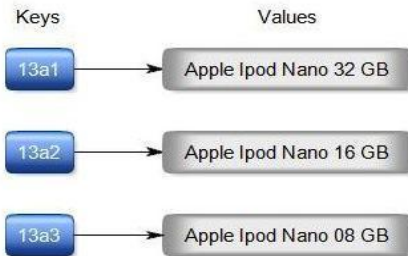
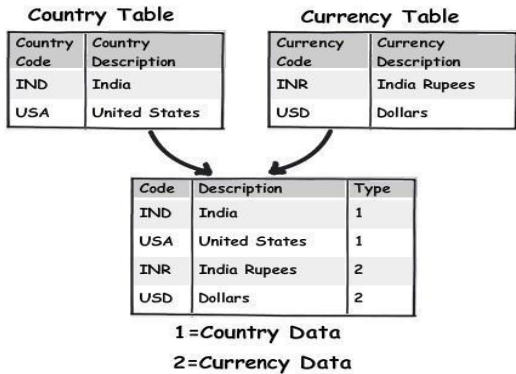


# Replicated commit (VLDB 2013)

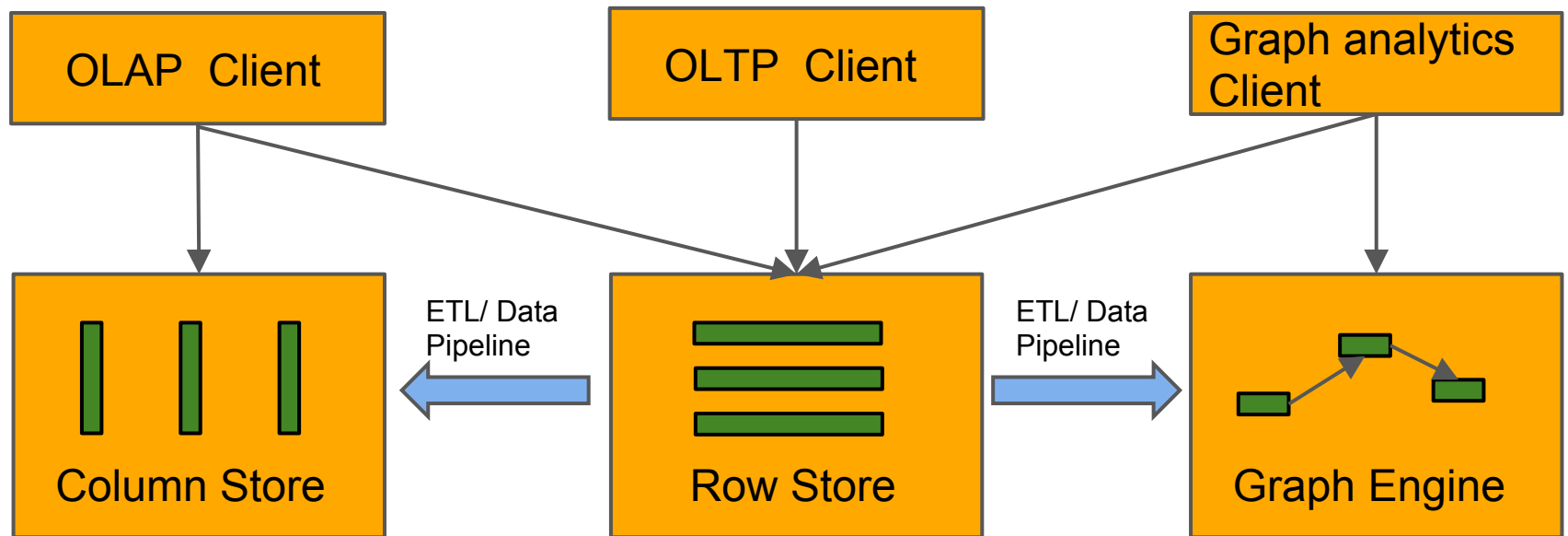
- ▶ Execute communication expensive transactions within a data center.
- ▶ Fault-tolerance across data centers using Paxos
- ▶ Consistency within data center using 2PC
- ▶ Significant reduction in communication costs.



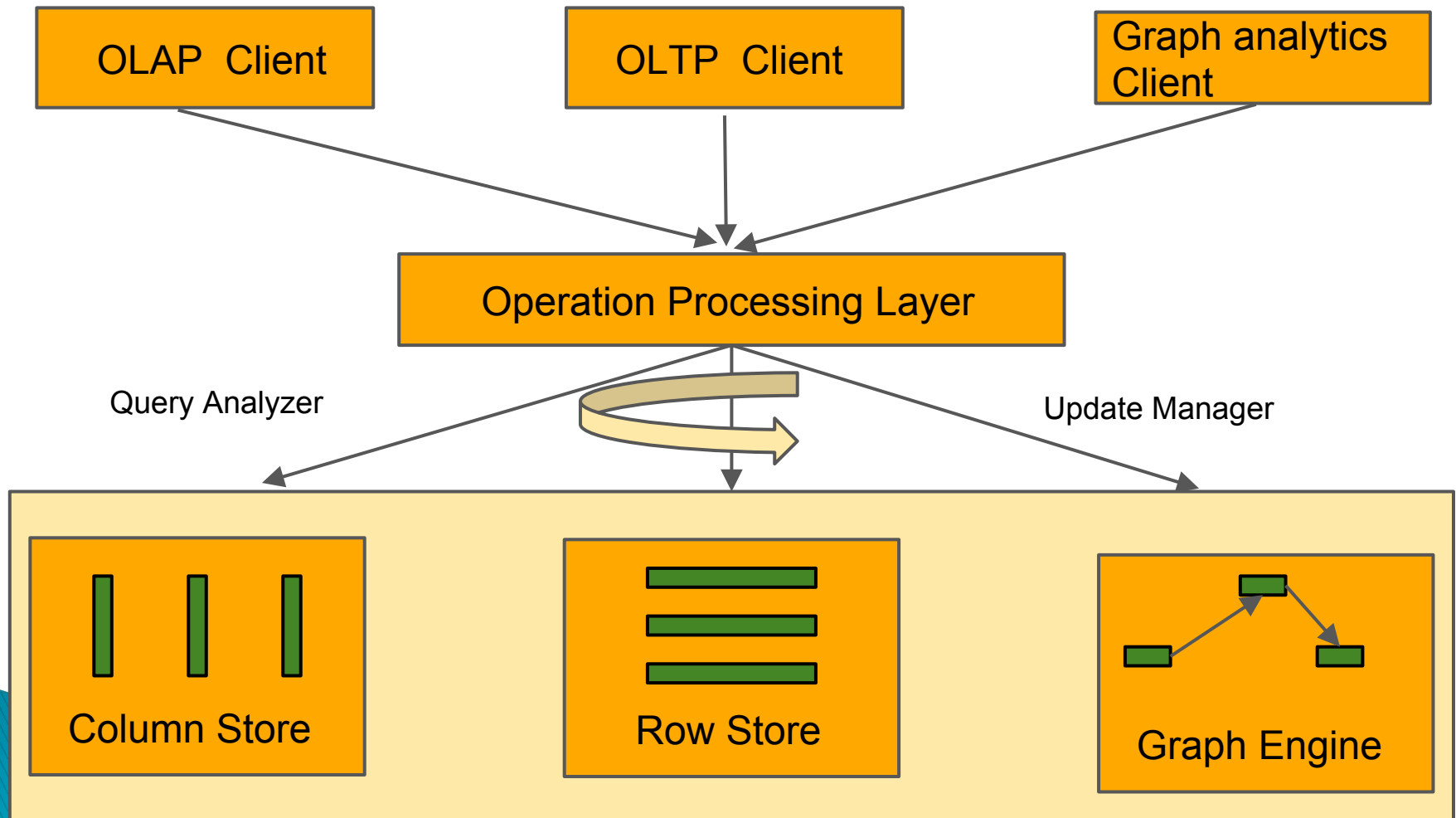
# Data Variety



# System Architecture



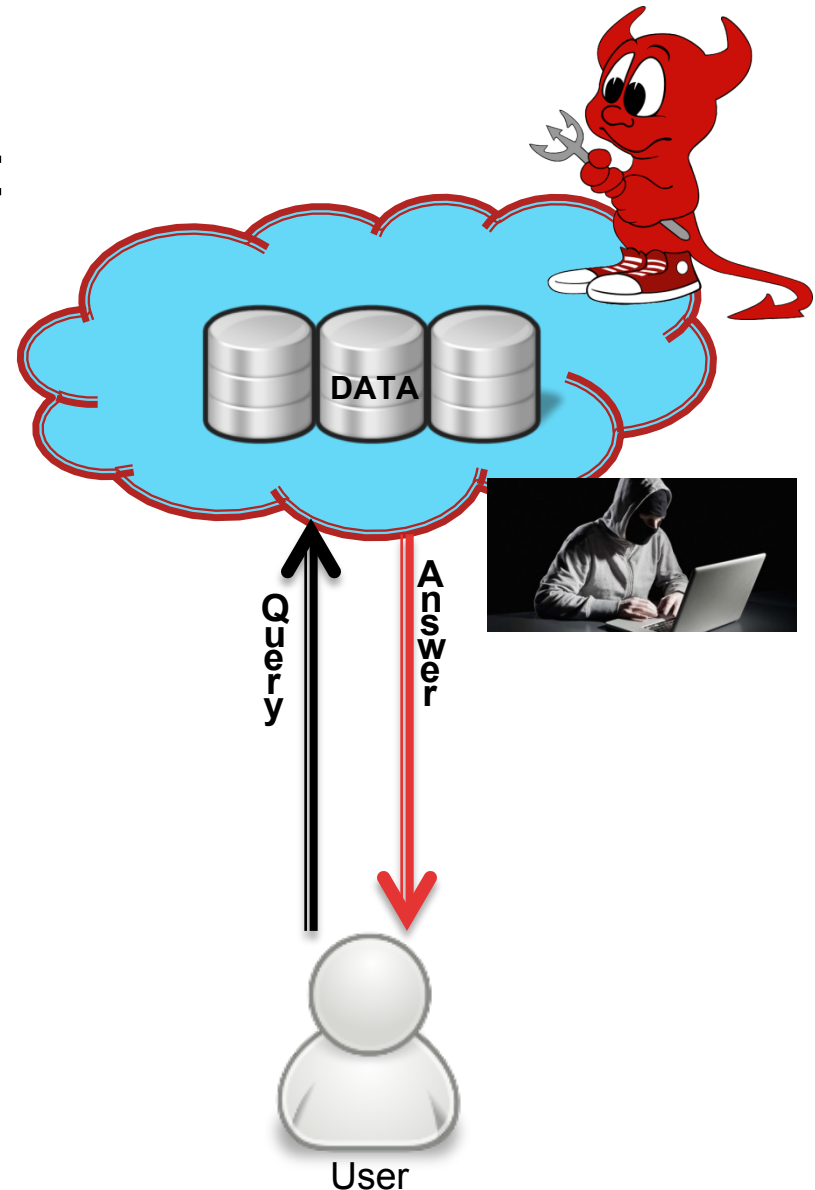
# Replication Driven Solution





# Privacy-Preserving Data Services In the Cloud

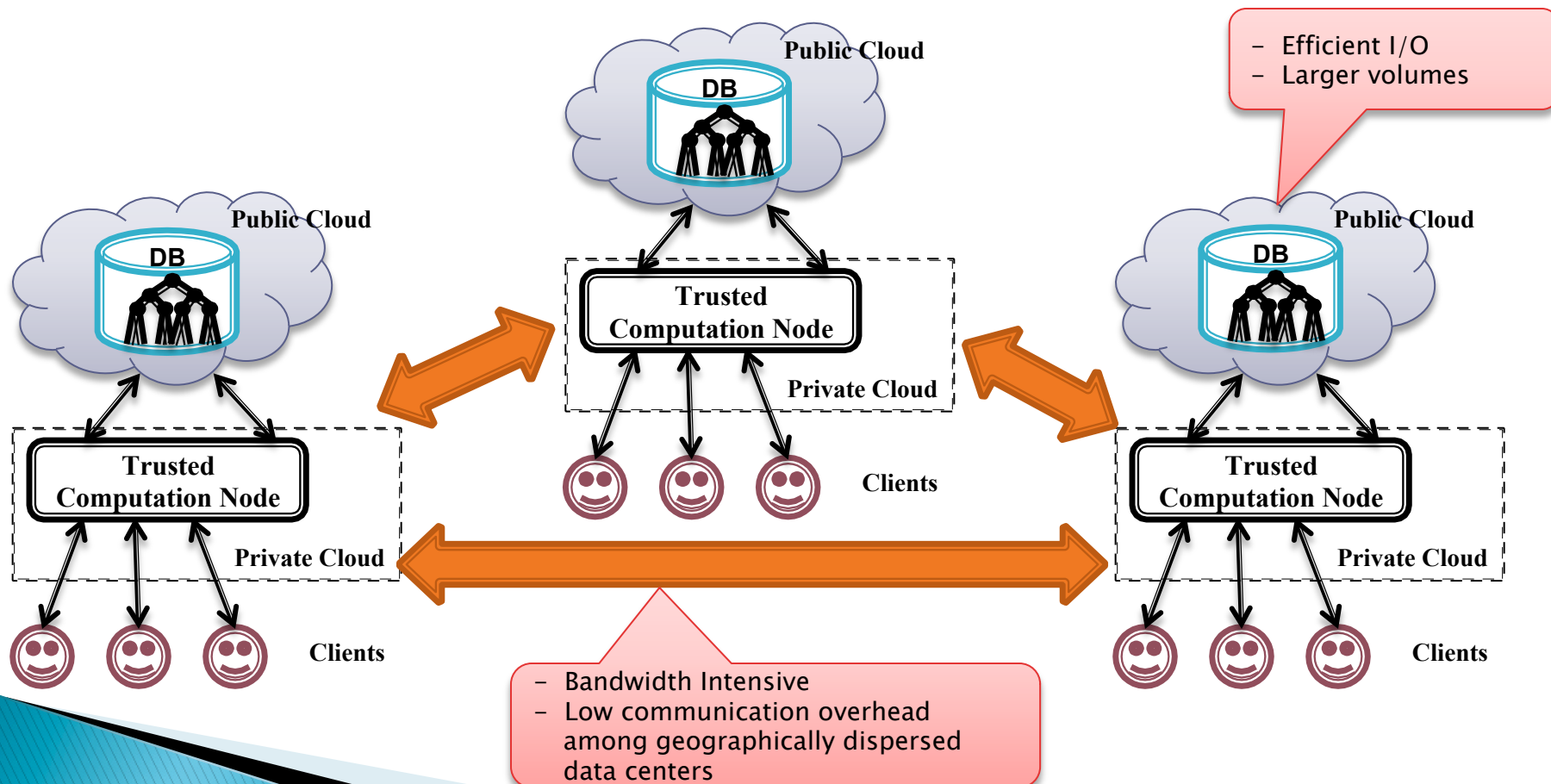
- Data security and privacy in the cloud vulnerable to:
  - Curious/Snooping system administrators
  - Hackers with illegal access
- **GOAL:** Functionality and performance of database systems while preserving data privacy and security.



# Cloud for Privacy-Preserving Data Storage

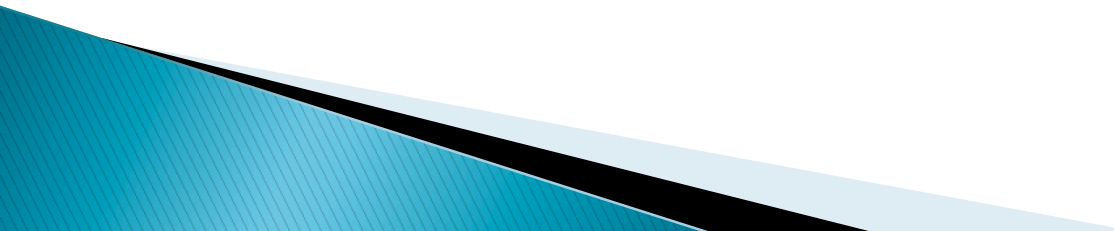
Secure data storage requires **encryption** before outsourcing the data

- ❖ More **space** is required to store
- ❖ Data is **transferred encrypted**





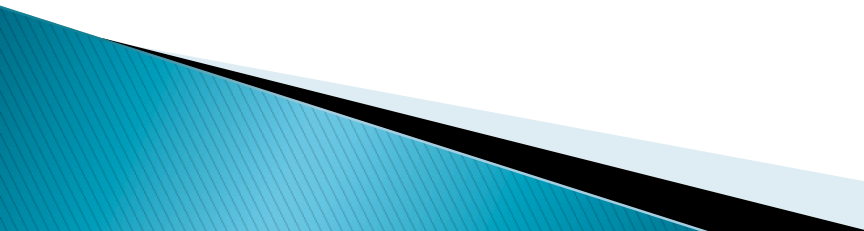
# Experimental Setup

- **Single DataCenter**
    - Servers, Racks and Clusters
    - Different Compute, Memory and Network configs
  - **Multi DataCenter**
    - Datacenters in physically different locations
  - Virtual Machine Access as given by Amazon EC2, Microsoft Azure
- 

# Evaluation Workloads

- **TPC-C** : Evaluating Single partitions and distributed Transactional Processing
- **Transactional YCSB** : Geo-Replication
- **TPC-H** and **Graph Workloads** - Variety

# Evaluation Scenarios

- Contention
  - Throughput Evaluation
  - Scale-Out
  - Handling Failures - Node and Data Centers
- 

# Wish List

- ▶ VM Placement Control
  - ▶ Beyond the Virtual Machine Statistics:
    - NW Utilization
    - Physical Machine Utilization
    - Disc Utilization
  - ▶ Infrastructure for Benchmark Generation
  - ▶ Overall, EC2 + more control and statistics.
- 