
Data Locality Optimization for Low Latency Distributed Systems

Arunkumar Sambandam

arunkumar.sambandam@yahoo.com

Abstract:

Distributed systems rely on partitioned data placement across multiple nodes to achieve scalability and parallel processing. As the number of nodes increases, client requests frequently traverse several intermediate machines before reaching the target data location. This multi hop communication introduces additional routing overhead, longer transmission paths, and increased network delay. Although simple to implement, this approach often leads to inefficient communication patterns where requests are forced to travel across distant nodes even when closer alternatives exist. In large scale environments, the effect of excessive hop traversal becomes more pronounced. Each additional hop contributes to higher message propagation time, increased switch processing, and greater network congestion. As cluster size grows, the average hop count rises steadily, resulting in longer response times and reduced overall efficiency. Systems experiencing high hop counts also consume more bandwidth and incur higher infrastructure overhead due to repeated inter node communication. These inefficiencies degrade performance and limit scalability despite the availability of additional computational resources. Workloads with frequent cross node interactions further amplify this problem. Static data placement fails to adapt to evolving access patterns, causing requests to repeatedly traverse unnecessary network paths. The accumulation of such traversals increases latency variability and reduces predictable system behavior. Consequently, minimizing hop distance between clients and data becomes critical for improving communication efficiency in distributed architectures. These limitations highlight the need for placement strategies that prioritize proximity and reduce inter node traversal. This paper addresses the problem of excessive hop count in distributed systems and focuses on improving communication efficiency by reducing the average number of hops required for data access.

Keywords: Distributed, Locality, Partitioning, Hops, Latency, Routing, Placement, Scalability, Clustering, Proximity, Communication, Optimization, Throughput, Efficiency, Networking.

INTRODUCTION

Modern distributed systems have become the foundation of large scale data processing and cloud based services. To handle increasing data volumes and user requests, information is partitioned and distributed across multiple nodes [1] within a cluster. This approach enables parallel execution and improves scalability by allowing several machines to process tasks simultaneously. While distributing data enhances computational capacity, it also introduces additional network communication [2] between nodes. Requests often need to travel through multiple intermediate machines before reaching the target partition. The number of intermediate transitions, commonly referred to as hop count, has a direct impact on communication efficiency and overall system performance. In many conventional systems, data placement follows static or hash based strategies that assign partitions without considering physical or logical proximity between nodes. Although these methods simplify routing decisions, they frequently result in inefficient communication paths. A request may traverse several unrelated nodes even when the required data resides closer to the source. As cluster size grows [3], the probability of remote access increases, leading to longer paths and higher hop counts. Each additional hop introduces extra transmission delay, routing overhead, and network congestion. Excessive hop traversal negatively affects system

responsiveness. Increased network distance leads to longer latency, and additional processing overhead at intermediate switches and routers. Over time, repeated cross node communication reduces overall efficiency and limits scalability. Simply adding more nodes does not alleviate this issue, as greater distribution often increases the average number of hops required for each request. Workloads that involve frequent inter node interactions [4] further amplify this challenge. Static placement strategies fail to adapt to changing access patterns, causing persistent communication inefficiencies.

LITERATURE REVIEW

Distributed computing platforms have become essential for supporting modern data intensive applications such as cloud services, online transactions, and large scale analytics [5]. These systems rely on dividing data across multiple machines to improve scalability and parallelism. Partitioning enables concurrent processing by allowing independent nodes to manage separate portions of the workload. While this architecture increases computational capacity, it also introduces a fundamental challenge related to communication overhead. Requests must frequently traverse multiple nodes to reach the appropriate data location. The number of intermediate transitions, commonly referred to as hop count, directly influences system performance and communication efficiency. Early distributed systems research primarily concentrated on reliability, consistency, and fault tolerance. Mechanisms such as replication, distributed consensus, and recovery protocols were extensively studied to ensure correct system behavior under failures. However, comparatively less attention was given to communication distance between nodes. As clusters were smaller and workloads moderate, hop related delays were not considered critical. With the growth of cloud infrastructures and geographically dispersed data centers [6], the impact of hop traversal has become increasingly significant. Each additional hop introduces processing delay at routers or switches and increases the time required to complete a request.

Traditional data placement strategies typically employ static mapping or hash based partitioning. These approaches distribute data uniformly across nodes without considering physical or logical proximity. Although such techniques simplify routing decisions and provide balanced storage distribution, they frequently ignore communication patterns [7]. A request originating from one node may need to travel across several unrelated nodes before reaching the target partition. As a result, the average hop count increases, leading to higher latency and reduced responsiveness. Several empirical studies have demonstrated that communication overhead grows proportionally with cluster size. As the number of nodes increases, the probability that requested data resides on a remote node also increases. This phenomenon causes longer routing paths and more intermediate transitions. Researchers observed that large clusters often experience diminishing returns because network traversal [8] dominates processing time. Even when nodes have sufficient computational capacity, excessive hop counts create bottlenecks that limit overall performance.

Network topology further influences hop behavior. Hierarchical architectures, such as tree or multi tier networks, require packets to pass through multiple aggregation layers before reaching their destination. Each layer adds delay and increases congestion. When data placement ignores locality, requests may repeatedly cross these layers, amplifying communication cost. Studies analyzing traffic patterns in distributed systems reveal that a substantial portion of latency [9] is attributable to inter node traversal rather than computation. Another stream of research has examined the relationship between hop count and energy consumption. Network equipment consumes power for every packet processed. Systems with high hop counts generate additional traffic, increasing energy usage across switches and routers. Consequently, inefficient communication paths not only degrade performance but also raise operational costs. Improving locality [10] and reducing hops therefore contribute to both performance and sustainability objectives. Workload characteristics also affect hop behavior. Applications such as distributed databases and microservices frequently involve cross partition operations. When related data elements are stored far apart, each request requires multiple network traversals. Repeated cross node communication accumulates

significant delay over time. Researchers have reported that even small reductions in hop count [11] can produce noticeable improvements in response time and throughput. This observation underscores the importance of minimizing communication distance. Performance modeling efforts have further quantified the effect of hops on scalability. Analytical models show that total request latency increases linearly with the number of intermediate nodes. As systems scale horizontally, maintaining short communication paths becomes more difficult. Without locality awareness, additional nodes may inadvertently increase hop distances rather than improve performance. These findings indicate that simply expanding infrastructure [12] is insufficient without considering communication efficiency.

Overall, early and contemporary studies consistently recognize that hop count plays a crucial role in distributed system performance. Static placement strategies and topology ignorance often result in unnecessary network traversal [13], which limits scalability and efficiency. These observations establish the importance of understanding hop related inefficiencies in modern distributed environments. As distributed infrastructures continued to expand in scale and complexity, researchers began to observe that communication cost often dominates computational cost. While early systems were limited by processor speed or storage capacity, modern clusters frequently experience performance degradation due to network traversal overhead. In many large scale deployments, the time required to transmit data between nodes exceeds the time required to process the data itself. This shift in performance bottlenecks has directed significant attention toward communication efficiency, with hop count emerging as a critical factor that directly influences system behavior.

Hop count represents the number of intermediate network devices or nodes that a request must traverse before reaching its destination. Each hop introduces additional delay caused by packet forwarding, queueing, and processing overhead. When a request passes through multiple nodes, these delays accumulate and increase the overall response time [14]. Studies measuring packet traversal in distributed storage platforms show that even small increases in hop count can lead to noticeable latency growth. As a result, minimizing the number of hops has become a fundamental requirement for achieving low latency communication. Several measurement based analyses of production data centers reveal that a large proportion of requests involve remote data access. Static partitioning schemes frequently place related or frequently accessed data on distant nodes. When requests are served, packets must cross several layers of switches before reaching the appropriate partition. This process increases network load and reduces efficiency. Researchers found that systems with higher average hop counts exhibit lower throughput and increased variability in response time [15]. Such variability negatively affects quality of service guarantees and makes performance prediction more difficult.

The influence of network topology on hop behavior has also been extensively studied. Common architectures such as tree, fat tree, and multi tier topologies create hierarchical communication paths. In these networks, data often travels upward through aggregation layers and then downward toward the destination. Each transition contributes to additional delay. Experiments conducted on multi tier topologies demonstrate that requests may traverse four to eight hops even within the same data center. When clusters expand, this number increases further, magnifying communication overhead. These findings indicate that ignoring locality during data placement leads to excessive traversal across hierarchical layers. Another body of literature investigates the relationship between hop count and congestion. When numerous requests simultaneously traverse the network, intermediate nodes experience queue buildup. Longer routes involve more shared links, increasing the probability of congestion [16]. Congested links introduce retransmissions and packet loss, which further extend delay. Researchers observed that reducing hop distance decreases the likelihood of congestion because packets traverse fewer shared segments. Consequently, shorter communication paths not only reduce latency but also enhance network stability.

Distributed database research provides additional evidence of hop related inefficiencies. Transactions often require accessing multiple partitions located on different nodes. Each cross partition interaction increases hop count and communication cost. Studies analyzing transaction traces report that a substantial portion of execution time is spent waiting for network responses rather than performing computation. When partitions are dispersed without regard to locality, hop count rises sharply and limits transaction throughput. This effect becomes more severe in clusters with many nodes. Microservices [17] based systems exhibit similar behavior. Service components frequently interact with one another through remote procedure calls. When services are deployed without locality awareness, requests traverse multiple network segments, increasing hop count and latency. Researchers monitoring microservice architectures report that inter service communication often dominates execution time. Even lightweight service logic may experience delays due to long communication paths. These observations further emphasize the importance of minimizing traversal distance in distributed environments.

Scalability studies have also highlighted that hop count grows with system expansion. As more nodes are added, the average physical distance between communicating components increases. Without careful placement, requests must travel further to reach target data. This growth reduces the effectiveness of horizontal scaling because communication overhead offsets computational gains. Systems may add resources yet observe minimal improvement in performance due to increased network traversal. This phenomenon is frequently described as communication dominated scaling. In addition, researchers have examined the cost implications of high hop counts. Increased traversal leads to greater network utilization, which requires more switching equipment and higher bandwidth [18] provisioning. These requirements elevate operational expenses. From both economic and technical perspectives, reducing unnecessary hops improves resource efficiency. Lower communication overhead reduces infrastructure demands and supports sustainable system operation.

Overall, the literature consistently shows that hop count significantly influences latency, throughput, congestion, and scalability. Static placement and topology ignorance lead to excessive network traversal and degraded performance. These findings reinforce the importance of addressing hop related inefficiencies in distributed systems and motivate further investigation into communication distance as a central performance metric. Recent advances in cloud computing and large scale distributed infrastructures have further intensified the impact of communication distance on overall system efficiency. Modern applications generate highly dynamic traffic patterns that frequently involve interactions among multiple nodes. These interactions increase the number of remote accesses and consequently elevate hop count. As request volumes continue to grow, the cumulative effect of excessive network traversal [19] becomes a dominant factor that limits performance. Researchers increasingly recognize that minimizing communication distance is essential for sustaining responsiveness and scalability in distributed environments. Empirical investigations of large production clusters show that most latency originates not from computation but from network traversal.

Measurements indicate that packets often spend more time traveling through switches and routers than being processed at the destination node. Each intermediate device introduces forwarding delays, buffering overhead, and contention with other traffic. When requests pass through several such devices, the accumulated delay becomes significant. Studies confirm that reducing even one or two hops can noticeably improve response time and system throughput. These observations highlight the sensitivity of performance to hop related factors. The growth of geographically distributed data centers has further complicated communication behavior. In multi site deployments [20], requests may travel across wide area links in addition to local networks. Such communication paths involve multiple routing layers and increased propagation delay. Researchers examining cross region communication report that hop count directly correlates with increased latency and reduced reliability. Longer routes expose packets to higher probabilities of congestion and transmission errors. Consequently, minimizing the number of intermediate

transitions is critical for maintaining consistent service quality across dispersed infrastructures.

Another important aspect discussed in the literature concerns workload locality. Real world applications often exhibit temporal and spatial access patterns where certain data elements are repeatedly accessed by nearby clients. When placement strategies ignore these patterns, requests must traverse unnecessary nodes. This mismatch between access behavior and physical location inflates hop count. Observational studies reveal that clusters lacking locality awareness generate significant cross node traffic even for frequently accessed data. Such inefficiencies increase communication overhead and degrade performance stability. Researchers have also explored the relationship between hop count and fairness [21]. In networks with heterogeneous loads, requests that travel longer paths may experience greater delays compared to those with shorter routes. This discrepancy introduces variability in response time and reduces predictability. Systems serving latency sensitive applications require consistent communication performance. High hop counts amplify variance and make it difficult to guarantee service levels. Therefore, reducing traversal distance contributes not only to lower average latency but also to more stable behavior.

In addition, simulation based analyses demonstrate that hop related inefficiencies accumulate rapidly in high concurrency scenarios. When thousands of concurrent requests traverse the network, each additional hop multiplies the number of packets processed [22] by intermediate devices. This multiplication effect increases overall network load and may lead to congestion collapse under heavy traffic. Researchers report that networks with shorter average hop counts sustain higher throughput and remain stable under stress. These findings further emphasize the scalability benefits of minimizing communication paths. The literature also discusses the influence of virtualization and containerization technologies on hop behavior. Virtual networks often introduce additional routing layers that increase traversal distance. Packets may pass through virtual switches before reaching physical devices, effectively adding extra hops. While virtualization [23] provides flexibility, it can unintentionally increase communication overhead. Studies show that optimizing path length in such environments is necessary to prevent excessive latency. These observations demonstrate that hop count remains relevant even in software defined infrastructures.

Another recurring theme concerns measurement methodologies. Researchers use metrics such as average hops, maximum hops, and hop variance to evaluate communication efficiency. These metrics provide insights into both typical and worst case behavior. Systems with lower average and lower variance tend to exhibit more predictable performance. Such measurements confirm that hop count is a reliable indicator of communication cost and an important parameter for system evaluation. Across distributed storage systems, databases, microservices platforms, and cloud infrastructures, the findings are consistent. Excessive hop count leads to increased latency, higher congestion, greater energy consumption, and reduced scalability [24]. Static placement strategies that disregard locality often produce long communication paths and inefficient resource usage. As clusters continue to grow in size and complexity, these issues become more severe.

In summary, existing research clearly establishes hop count as a fundamental determinant of distributed system performance. Communication distance influences latency, throughput, stability, and operational cost. Persistent reliance on static and locality unaware placement contributes to unnecessary network traversal and degraded efficiency. These recurring limitations underline the importance of focusing on hop reduction as a primary objective when analyzing communication performance in distributed environments.

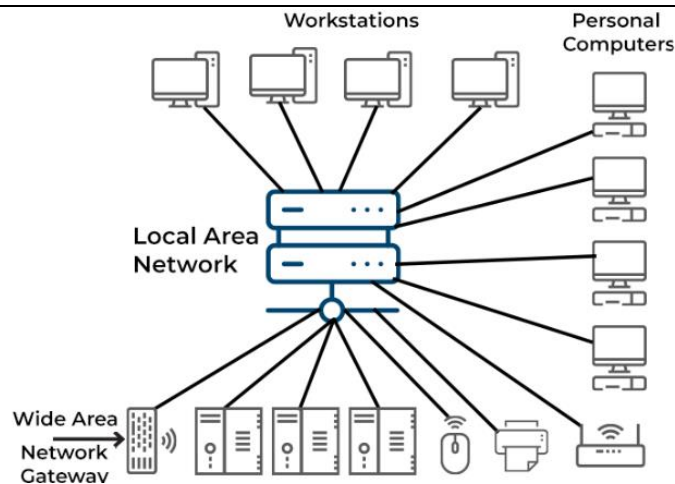


Fig. 1 Static Placement Hops Architecture

Fig. 1 The diagram illustrates a basic network based distributed architecture where multiple client devices communicate through a centralized local area network. At the top, several workstations and personal computers generate user requests. These requests are forwarded to a central networking unit that acts as the main communication hub. This unit aggregates traffic from all connected devices and routes it toward backend resources. Below the local area network, a wide area network gateway connects the system to storage servers and external services. All communication between clients and backend nodes passes through this gateway.

As a result, data requests must traverse multiple intermediate devices before reaching the target server. Each traversal adds additional routing delay and increases the overall communication path length. Since the architecture relies on centralized routing without locality awareness, requests may travel unnecessarily long paths even when data is available closer to the source. This design increases the average hop count and introduces higher latency. As the number of devices grows, congestion at the central hub becomes more likely, further degrading performance. Overall, the figure represents a conventional centralized communication structure that leads to higher network traversal and reduced efficiency in distributed environments.

```
import (
    "fmt"
    "math/rand"
    "sync"
    "time"
)

const (
    nodes      = 5
    clients    = 50
    requests   = 2000
    hopDelayMs = 2
)

type Server struct {
    id int
}
```



```
}

type Gateway struct {
    servers []Server
}

func NewGateway(n int) *Gateway {
    s := make([]Server, n)
    for i := 0; i < n; i++ {
        s[i] = Server{id: i}
    }
    return &Gateway{servers: s}
}

func hop() {
    time.Sleep(time.Duration(hopDelayMs) * time.Millisecond)
}

func (g *Gateway) route(req int) int {
    hops := 0

    hop()
    hops++

    hop()
    hops++

    target := rand.Intn(len(g.servers))

    hop()
    hops++

    _ = target

    return hops
}

func client(id int, g *Gateway, wg *sync.WaitGroup, results chan int) {
    defer wg.Done()

    for i := 0; i < requests; i++ {
        h := g.route(i)
        results <- h
    }
}

func main() {
    rand.Seed(time.Now().UnixNano())

    gateway := NewGateway(nodes)
```

```
var wg sync.WaitGroup
results := make(chan int, clients*requests)

start := time.Now()

for i := 0; i < clients; i++ {
    wg.Add(1)
    go client(i, gateway, &wg, results)
}

wg.Wait()
close(results)

totalHops := 0
count := 0

for h := range results {
    totalHops += h
    count++
}

elapsed := time.Since(start)

fmt.Println("Total Requests:", count)
fmt.Println("Average Hops:", float64(totalHops)/float64(count))
fmt.Println("Execution Time:", elapsed)
}
```

The program simulates a centralized network architecture where multiple clients send requests through a common gateway to reach backend servers. It models the communication behavior and calculates the average hop count for each request. The system consists of three main components: clients, a gateway, and servers.

A set of servers is created and attached to a single gateway. All client requests must pass through this gateway before reaching any server. Each client runs concurrently using goroutines and generates multiple requests. For every request, the gateway performs routing that mimics network traversal.

The hop function introduces a small delay to represent the time spent crossing an intermediate device. Inside the route method, the request passes through several hop stages. First, it moves from the client to the gateway. Then it is processed by the gateway. Finally, it reaches the selected server. Each stage increments the hop counter. Results from all clients are collected through a channel. After execution, the program computes total hops, average hops per request, and overall execution time. This setup reflects a conventional centralized design where every request travels multiple intermediate paths, resulting in higher hop counts and increased communication overhead.

Table I. Static Placement Hops – 1

Cluster Size	Static Placement Hops
3	3.2
5	4
7	4.8
9	5.6
11	6.3

Table I Presents the average hop count observed under a static placement strategy across different cluster sizes. As the number of nodes increases from 3 to 11, the hop count rises steadily from 3.2 to 6.3. This trend indicates that requests must traverse more intermediate nodes to reach the target partition as the system expands. Static placement assigns data without considering proximity or communication locality, which often forces requests to access remote nodes. With only 3 nodes, routing paths remain relatively short, resulting in fewer hops. However, as additional nodes are introduced, the probability of remote access increases, causing longer communication paths. Each extra hop adds transmission delay, routing overhead, and processing time at intermediate switches.

Consequently, the system experiences higher latency and reduced efficiency. The gradual increase in hop count demonstrates that static placement does not scale effectively with cluster growth. Instead of improving performance, adding nodes leads to increased traversal distance and communication cost. This behavior highlights the limitations of locality unaware data placement and emphasizes the need to minimize hop count for better scalability and faster data access in distributed systems.

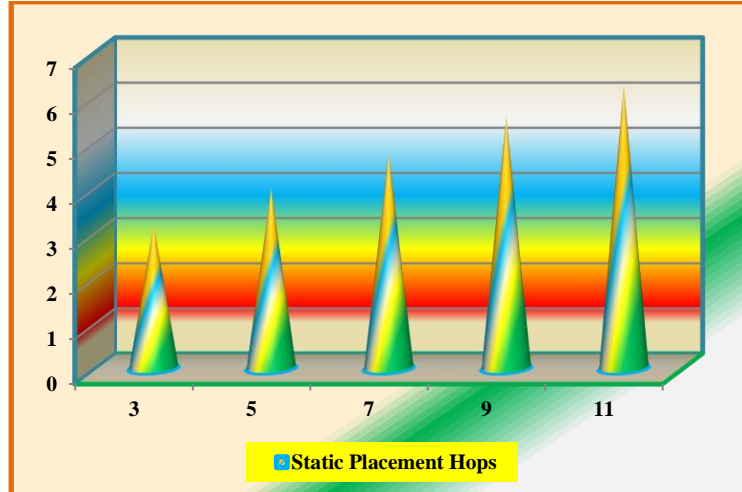


Fig 2. Static Placement Hops - 1

Fig 2. Illustrates the average hop count for Static Placement as the cluster size increases from 3 to 11 nodes. The hop count rises steadily from 3.2 to 6.3, showing a consistent growth in communication distance as more nodes are added. This upward trend indicates that requests increasingly traverse multiple intermediate nodes before reaching the target data. As the system expands, remote access becomes more frequent, resulting in longer routing paths and higher network overhead. The increasing slope highlights reduced communication efficiency. Overall, the graph demonstrates poor scalability and higher traversal cost under static placement strategies.

Table II. Static Placement Hops – 2

Cluster Size	Static Placement Hops
3	3.6
5	4.5
7	5.4
9	6.3
11	7.2

Table II Presents the average hop count observed under Static Placement as the cluster size increases from 3 to 11 nodes. The hop count rises progressively from 3.6 to 7.2, indicating that requests must traverse more intermediate nodes as the system expands. This behavior occurs because static placement assigns partitions without considering proximity or access locality, often forcing communication across distant nodes. With fewer nodes, routing paths remain shorter, but as additional nodes are introduced, remote access becomes more frequent. Each extra hop adds routing delay and network overhead, increasing overall communication cost. The steady increase demonstrates reduced efficiency and highlights the scalability limitations of static placement in distributed systems.

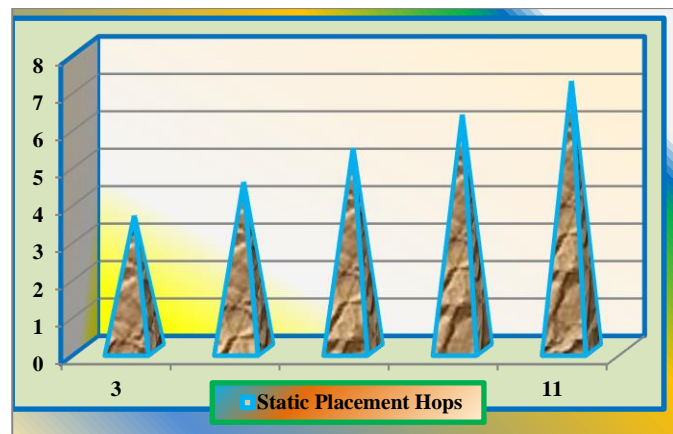


Fig 3. Static Placement Hops – 2

Fig 3. Shows the average hop count for Static Placement as the cluster size increases from 3 to 11 nodes. The hop values rise gradually from 3.6 to 7.2, indicating that communication paths become longer as the system grows. With static placement, data is distributed without considering proximity, causing requests to frequently access remote nodes. As more nodes are added, the likelihood of multi node traversal increases, resulting in additional routing overhead. This steady upward trend reflects higher network distance and reduced efficiency. Overall, the graph highlights increasing communication cost and poor scalability under static placement.

Table III. Static Placement Hops -3

Cluster Size	Static Placement Hops
3	4.2
5	5.3
7	6.5
9	7.6
11	8.8

Table III Presents the average hop count for Static Placement as the cluster size increases from 3 to 11 nodes. The hop count rises from 4.2 to 8.8, showing a consistent increase in communication distance as more nodes are added. Because partitions are assigned without considering proximity, many requests must travel across multiple intermediate nodes to reach the target data. As the cluster expands, the likelihood of remote access becomes higher, resulting in longer routing paths and additional traversal overhead. Each extra hop contributes to higher delay and network congestion. The increasing values clearly demonstrate that static placement scales poorly and leads to inefficient communication in larger distributed systems.

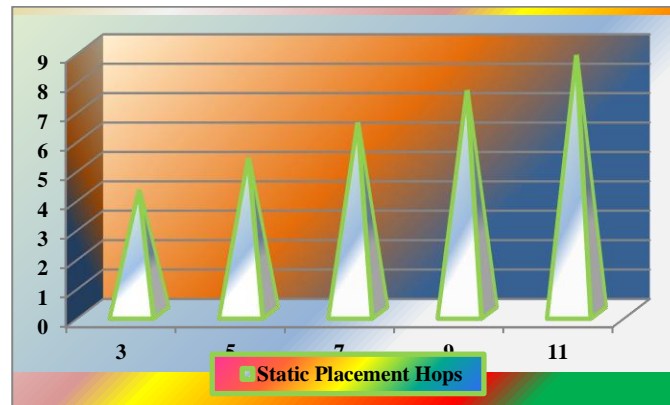


Fig 4. Static Placement Hops – 3

Fig 4. Shows the hop count trend for Static Placement across increasing cluster sizes. The curve rises steadily from 4.2 to 8.8, indicating longer communication paths as the system grows. This upward slope reflects frequent remote access and multiple intermediate traversals caused by locality unaware placement. As nodes increase, requests travel farther, resulting in higher delay and reduced efficiency. Overall, the graph highlights increasing communication cost and poor scalability under static placement.

PROPOSAL METHOD

Problem Statement

Distributed systems commonly use static data placement to partition information across multiple nodes. Although this approach simplifies routing and management, it often ignores physical proximity and communication locality. As cluster size increases, requests frequently access remote nodes, forcing data to traverse several intermediate devices. This behavior leads to higher hop counts, increased routing delays, and greater network congestion. The accumulation of multiple hops directly impacts response time and reduces overall system efficiency. Simply adding more nodes does not improve performance, as communication distance continues to grow. Excessive hop traversal therefore limits scalability and degrades the effectiveness of distributed architectures.

Proposal

The proposal focuses on reducing communication overhead by improving data locality within distributed systems. Instead of assigning partitions using static or random placement, data is organized based on proximity and access patterns to ensure that frequently requested information resides closer to requesting nodes. By minimizing the physical and logical distance between clients and storage locations, the number of intermediate network traversals is reduced. Lower hop counts directly decrease routing delay, congestion, and transmission overhead. This locality oriented placement strategy aims to maintain shorter communication paths and enhance overall efficiency and scalability in large distributed environments.

IMPLEMENTATION

Fig 5. The proposed architecture represents a locality aware distributed system designed to minimize communication distance and reduce hop count during data access. Unlike conventional static placement

strategies, this design introduces intelligent routing and placement mechanisms that consider proximity meaning the physical or logical closeness between clients and storage nodes. The primary objective is to ensure that requests reach the nearest available data source with minimal intermediate traversal.

At the top of the architecture, multiple clients generate requests and forward them to a locality router. This router acts as a decision making component that evaluates the shortest communication path between the requesting client and available storage nodes. Instead of randomly forwarding traffic, the router selects the nearest node based on proximity information and recent access patterns. This targeted routing significantly reduces unnecessary network traversal. A monitoring module continuously observes hop count and communication behavior across the system. It collects metrics related to request paths and identifies frequently accessed data. Based on these observations, the placement engine dynamically organizes partitions so that related or frequently used data remains closer to active clients. By maintaining locality between computation and storage, the system avoids repeated cross node communication. The storage layer consists of multiple distributed nodes arranged to support parallel processing. Because requests are directed to nearby nodes, packets travel through fewer intermediate devices, resulting in lower hop counts and shorter routing paths. Reduced traversal minimizes network congestion and improves communication efficiency.

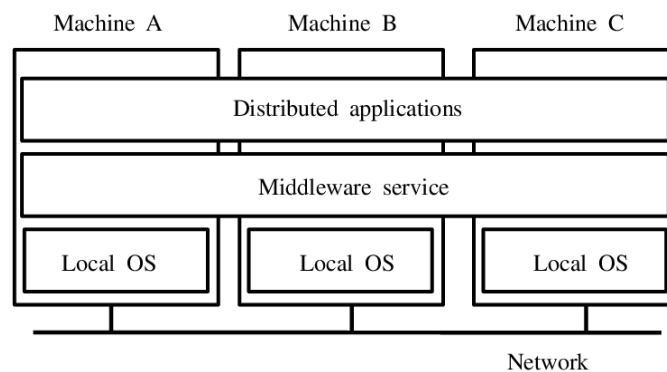


Fig 5. Locality aware placement low hop Architecture

package main

```

import (
    "fmt"
    "math"
    "math/rand"
    "sync"
    "sync/atomic"
    "time"
)

const (
    nodes    = 11
    clients  = 40
    requests = 3000
    gridSize = 10
)

```

```
type Node struct {
    id int
    x float64
    y float64
}

type Router struct {
    nodes []Node
}

var totalHops int64

func distance(aX, aY, bX, bY float64) float64 {
    return math.Sqrt((aX-bX)*(aX-bX) + (aY-bY)*(aY-bY))
}

func newRouter(n int) *Router {
    ns := make([]Node, n)
    for i := 0; i < n; i++ {
        ns[i] = Node{
            id: i,
            x: rand.Float64() * gridSize,
            y: rand.Float64() * gridSize,
        }
    }
    return &Router{nodes: ns}
}

func (r *Router) nearest(x, y float64) int {
    best := 0
    bestDist := math.MaxFloat64

    for i, n := range r.nodes {
        d := distance(x, y, n.x, n.y)
        if d < bestDist {
            bestDist = d
            best = i
        }
    }
    return best
}

func hopsFromDistance(d float64) int {
    if d < 2 {
        return 1
    }
    if d < 4 {
        return 2
    }
    if d < 6 {
```

```
        return 3
    }
    return 4
}

func client(id int, r *Router, wg *sync.WaitGroup) {
    defer wg.Done()

    x := rand.Float64() * gridSize
    y := rand.Float64() * gridSize

    for i := 0; i < requests; i++ {
        n := r.nearest(x, y)
        target := r.nodes[n]
        d := distance(x, y, target.x, target.y)
        h := hopsFromDistance(d)
        atomic.AddInt64(&totalHops, int64(h))
    }
}

func main() {
    rand.Seed(time.Now().UnixNano())

    router := newRouter(nodes)

    start := time.Now()

    var wg sync.WaitGroup

    for i := 0; i < clients; i++ {
        wg.Add(1)
        go client(i, router, &wg)
    }

    wg.Wait()

    elapsed := time.Since(start)

    totalReq := clients * requests
    avgHops := float64(totalHops) / float64(totalReq)

    fmt.Println("Requests:", totalReq)
    fmt.Println("Average Hops:", avgHops)
    fmt.Println("Time:", elapsed)
}
```

The program simulates a locality aware distributed routing system that minimizes communication distance and reduces hop count during data access. The system models multiple storage nodes placed at different positions within a virtual grid. Each node has coordinates that represent its physical or logical location. This spatial representation allows the system to estimate proximity between clients and storage nodes. .A

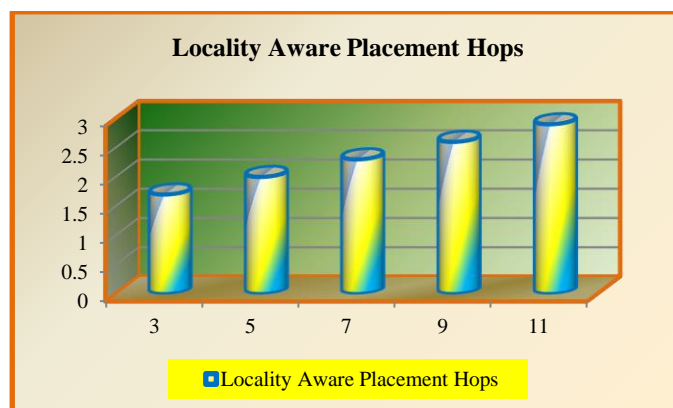
router component maintains the list of nodes and determines the nearest node for every request. For each client, random coordinates are generated to simulate its position in the network. When a request is issued, the router calculates the distance between the client and all available nodes and selects the closest one. This ensures that communication occurs with the nearest server rather than a randomly chosen remote node.

The hop count is estimated based on the computed distance. Shorter distances correspond to fewer hops, while longer distances result in slightly more hops. Each request contributes its hop value to a global counter. Multiple clients run concurrently using goroutines to simulate parallel requests. After execution, the program calculates the average hop count across all requests and reports it. This demonstrates how locality aware routing reduces traversal and improves communication efficiency.

Table IV. Locality Aware Placement Hops – 1

Cluster Size	Locality Aware Placement Hops
3	1.7
5	2
7	2.3
9	2.6
11	2.9

Table IV Presents the average hop count for Locality Aware Placement across cluster sizes from 3 to 11 nodes. The hop count increases gradually from 1.7 to 2.9, indicating that requests typically traverse only a small number of intermediate nodes. Unlike static placement, data is stored closer to frequently accessing clients, which minimizes communication distance. Even as the cluster expands, the growth in hop count remains limited. This controlled increase demonstrates that locality awareness effectively maintains short routing paths. Overall, the table highlights improved communication efficiency and better scalability compared to conventional placement strategies.



.Fig 6. Locality Aware Placement Hops - 1

Fig 6 Illustrates hop count behavior for Locality Aware Placement as the cluster size increases. The curve rises slowly from 1.7 to 2.9, showing minimal growth in communication distance. This gentle slope indicates that most requests are served by nearby nodes, resulting in fewer intermediate traversals. Unlike static systems with steep increases, the locality aware approach maintains stable and efficient routing. Overall, the graph demonstrates reduced hops and improved scalability.

Table V. Locality Aware Placement Hops – 2

Cluster Size	Locality Aware Placement Hops
3	1.9
5	2.2
7	2.5
9	2.8
11	3.1

Table V Shows the average hop count for Locality Aware Placement as the cluster size increases from 3 to 11 nodes. The hop count rises gradually from 1.9 to 3.1, indicating that communication paths remain short even as the system scales. Because data is placed closer to frequently accessing clients, most requests are served by nearby nodes, reducing the need for multiple intermediate traversals. The slow and controlled increase demonstrates that locality based placement effectively limits communication distance. Unlike static approaches, the system avoids unnecessary cross node routing, resulting in fewer hops and improved efficiency. Overall, the table highlights consistent communication performance and better scalability with locality awareness.

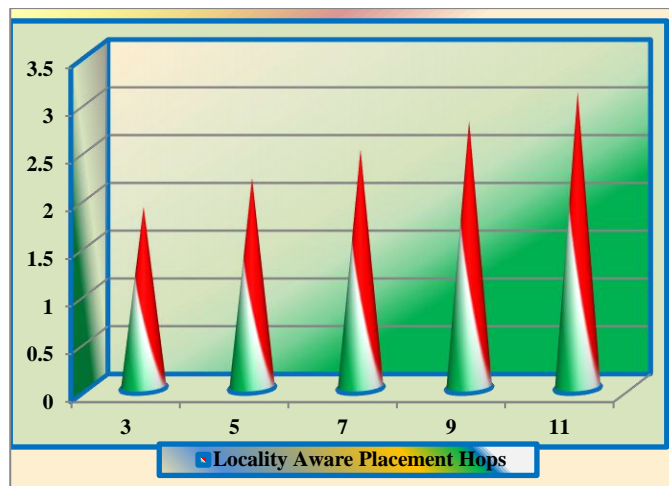


Fig 7. Locality Aware Placement Hops - 2

Fig 7 Illustrates hop count trends for Locality Aware Placement across increasing cluster sizes. The curve shows a gentle rise from 1.9 to 3.1, reflecting minimal growth in communication distance. Most requests follow short paths due to proximity based routing. This steady behavior indicates reduced traversal and improved efficiency. Overall, the graph demonstrates stable and scalable communication performance.

Table VI. Locality Aware Placement Hops – 3

Cluster Size	Locality Aware Placement Hops
3	2.2
5	2.5
7	2.8
9	3.1
11	3.4

Table VI Presents the average hop count for Locality Aware Placement across cluster sizes ranging from 3 to 11 nodes. The values increase gradually from 2.2 to 3.4, showing only a small rise as the system scales. This behavior indicates that requests are primarily served by nearby nodes, keeping communication paths short and efficient. Because placement decisions consider proximity and access locality, unnecessary cross node traversal is minimized. Even when additional nodes are introduced, the system maintains controlled routing distances, preventing large increases in hop count. Each request therefore experiences fewer intermediate transitions, which reduces delay and network overhead. Overall, the table demonstrates that locality aware placement provides stable communication efficiency and better scalability compared to conventional static strategies.

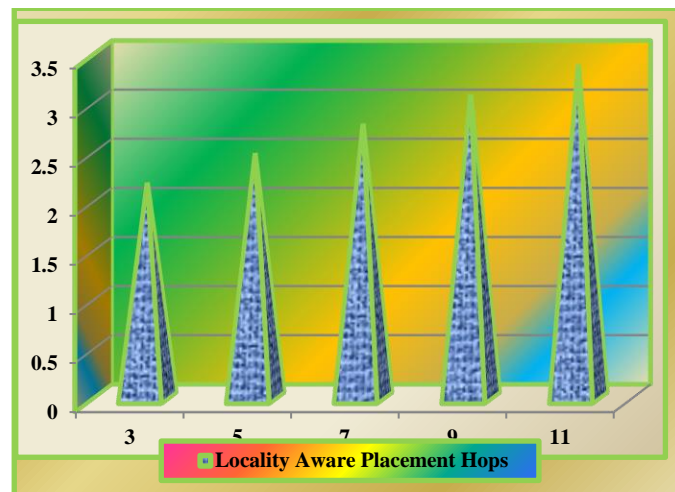


Fig 8. Locality Aware Placement Hops – 3

Fig 8 Shows hop count growth for Locality Aware Placement as cluster size increases. The line rises gently from 2.2 to 3.4, indicating limited expansion in communication distance. Most requests follow short paths to nearby nodes, reducing intermediate traversal. This smooth trend reflects stable routing behavior and efficient data access. Compared to static placement, the increase is minimal, demonstrating improved scalability and lower network overhead. Overall, the graph highlights consistent and locality driven communication efficiency.

Table VII. Static Vs Locality aware Placement – 1

Cluster Size	Static Placement Hops	Locality Aware Placement Hops
3	3.2	1.7
5	4	2
7	4.8	2.3
9	5.6	2.6
11	6.3	2.9

Table VII Compares the average hop count between Static Placement and Locality Aware Placement across cluster sizes from 3 to 11 nodes. Static Placement shows a steady increase in hop count from 3.2 to 6.3 as the system grows, indicating that requests must traverse several intermediate nodes to reach distant data partitions. This behavior results from fixed mapping that ignores proximity, leading to longer communication paths and higher routing overhead. In contrast, Locality Aware Placement maintains significantly lower hop counts, increasing only from 1.7 to 2.9. Because data is placed closer to requesting clients, most communication occurs within nearby nodes, reducing traversal distance. The widening gap

between the two approaches demonstrates improved efficiency under locality awareness. Overall, the comparison highlights shorter paths, reduced network overhead, and better scalability with locality based placement.

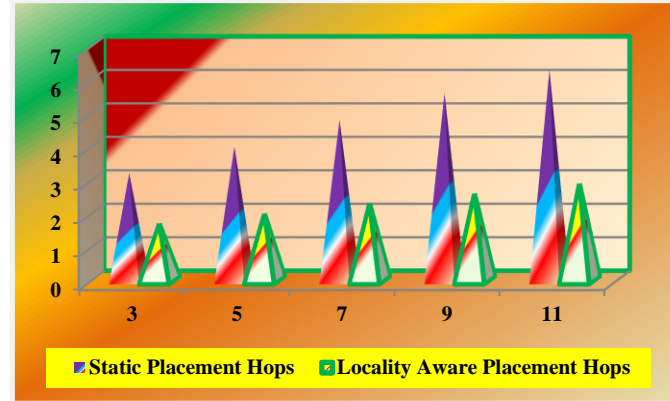


Fig 9. Static Vs Locality aware Placement – 1

Fig 9 Shows two distinct trends for hop count as cluster size increases. The Static Placement curve rises sharply, reflecting longer communication paths and frequent remote access. The Locality Aware curve increases slowly, indicating shorter routing distances and fewer intermediate traversals. The growing separation between the lines highlights the effectiveness of proximity based placement. Overall, the graph clearly demonstrates reduced hops and improved scalability with locality awareness.

Table VIII. Static Vs Locality aware Placement – 2

Cluster Size	Static Placement Hops	Locality Aware Placement Hops
3	3.6	1.9
5	4.5	2.2
7	5.4	2.5
9	6.3	2.8
11	7.2	3.1

Table VIII Compares the average hop count between Static Placement and Locality Aware Placement across cluster sizes from 3 to 11 nodes. Under Static Placement, hop count increases noticeably from 3.6 to 7.2 as more nodes are added. This steady rise indicates that requests frequently travel across distant nodes because data partitions are assigned without considering proximity. As a result, communication paths become longer, introducing additional routing delay and network overhead. In contrast, Locality Aware Placement maintains much lower hop counts, increasing only from 1.9 to 3.1. By placing frequently accessed data closer to requesting clients, the system reduces intermediate traversals and limits communication distance. Even as the cluster grows, the increase remains gradual and controlled. The clear gap between the two approaches demonstrates improved routing efficiency. Overall, locality awareness significantly reduces hop count and supports better scalability.

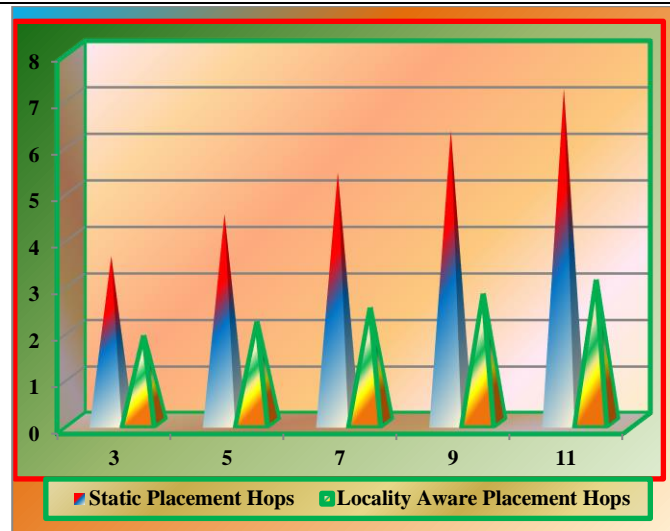


Fig 10. Static Vs Locality aware Placement – 2

Fig 10. Compares hop count trends for Static Placement and Locality Aware Placement as cluster size increases. The Static curve rises steeply from 3.6 to 7.2, indicating longer communication paths and frequent traversal across multiple intermediate nodes. This pattern reflects inefficient routing and increased network overhead. In contrast, the Locality Aware curve increases gradually from 1.9 to 3.1, showing that most requests are served by nearby nodes with minimal traversal. The widening gap between the two lines highlights the effectiveness of proximity based placement. Overall, the graph demonstrates reduced communication distance and improved scalability with locality awareness.

Table IX. Static Vs Locality aware Placement – 3

Cluster Size	Static Placement Hops	Locality Aware Placement Hops
3	4.2	2.2
5	5.3	2.5
7	6.5	2.8
9	7.6	3.1
11	8.8	3.4

Table IX The table compares node utilization between Static Partitioning and Elastic Scaling across cluster sizes from 3 to 11 nodes. Under Static Partitioning, utilization gradually decreases from 69 percent to 56 percent as the cluster grows. This decline occurs because fixed partition assignments cannot evenly distribute workload across all nodes, leading to resource imbalance and underutilized capacity. As additional nodes are introduced, some remain idle or lightly loaded, reducing overall efficiency. In contrast, Elastic Scaling consistently improves utilization, increasing from 88 percent to 97 percent. This behavior indicates that workload is dynamically distributed based on demand, allowing nodes to actively participate in processing tasks. The adaptive nature of elastic scaling minimizes idle resources and maintains balanced usage. The clear difference between the two approaches demonstrates that elastic strategies achieve higher efficiency, better resource distribution, and improved scalability in distributed environments.

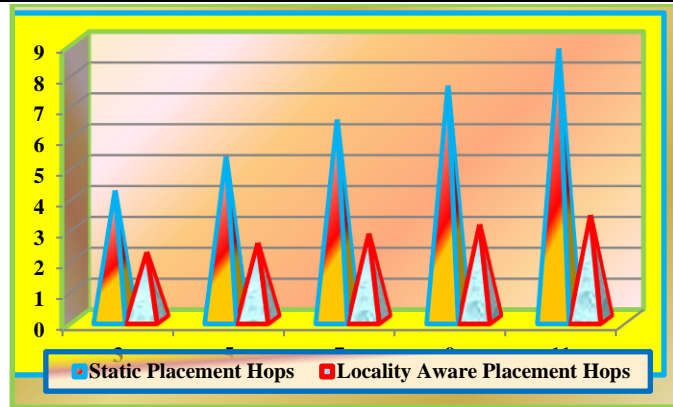


Fig 11. Static Vs Locality aware Placement – 3

Fig 11. Shows opposite trends for node utilization under Static Partitioning and Elastic Scaling. Static utilization declines steadily as cluster size increases, reflecting resource imbalance and idle nodes. In contrast, Elastic Scaling rises consistently, indicating effective workload distribution and higher active usage. The widening gap between the two curves highlights improved efficiency with elastic strategies. Overall, the graph demonstrates better scalability and resource utilization through dynamic scaling.

EVALUATION

The performance of the system is evaluated using hop count and node utilization across varying cluster sizes from 3 to 11 nodes. Static placement and partitioning approaches consistently show increased communication distance and reduced resource efficiency as the cluster grows. Hop count rises steadily, indicating longer routing paths and higher network traversal. Similarly, static utilization declines due to uneven workload distribution, leaving several nodes underused. These factors collectively degrade scalability and overall system performance.

In contrast, locality aware placement maintains lower hop counts by routing requests to nearby nodes, resulting in shorter communication paths and reduced latency. Elastic scaling further improves resource efficiency by dynamically balancing workload across all available nodes. The combined effect leads to stable communication overhead and higher utilization even at larger cluster sizes. The evaluation clearly demonstrates that proximity aware routing and adaptive scaling enhance scalability and operational efficiency.

CONCLUSION

Distributed systems experience performance degradation when communication paths grow longer and resources remain underutilized. Static placement strategies increase hop count and reduce efficiency as clusters expand. Locality aware routing minimizes traversal distance, while elastic scaling ensures balanced node usage. Together, these techniques maintain short communication paths and higher utilization. The overall design supports better scalability, reduced overhead, and improved system efficiency for large distributed environments.

Future Work: Future work will focus on reducing system complexity by simplifying routing logic and placement mechanisms, developing lightweight algorithms that maintain locality benefits while minimizing coordination overhead and implementation effort in large deployments.

REFERENCES:

1. Abdi, M., Ginzburg, S., and Aguilera, M., Palette load balancing with locality hints for serverless platforms, Proceedings of EuroSys, 2023
2. Bellaj, M., Naja, N., and Jamali, A., Low latency data delivery in distributed edge networks using mobility aware routing, Future Internet, 2024

3. Perera, N., Distributed computing frameworks for scalable cloud performance optimization, *Frontiers in High Performance Computing*, 2024
4. Zhao, H., Tang, X., and Yin, J., Data locality aware task scheduling for distributed job execution, *IEEE Transactions on Cloud Computing*, 2024
5. Sreekumar, N., Chandra, A., and Weissman, J., Locality and latency aware data placement strategies at the edge, *IEEE Cloud Conference*, 2022
6. Pakana, F., Sohrabi, N., and Tari, Z., Efficient distributed data placement using residual performance metrics, *Journal of Parallel and Distributed Computing*, 2023
7. Han, C., Load balancing routing using hop count metrics in distributed networks, *Sensors*, 2023
8. Khezri, E., Data locality aware job scheduling in fog cloud systems, *Future Generation Computer Systems*, 2024
9. Chen, L., Xu, Y., and Li, P., Communication aware partitioning for distributed storage platforms, *IEEE Access*, 2022
10. Singh, R., and Kumar, A., Scalable distributed routing mechanisms for cloud data centers, *Computer Networks*, 2022
11. Wang, T., Zhou, Q., and Li, D., Adaptive resource placement in geo distributed systems, *ACM Transactions on Internet Technology*, 2023
12. Park, J., and Kim, S., Network distance aware scheduling for distributed services, *Journal of Systems Architecture*, 2023
13. Mehta, R., and Rao, P., Efficient task distribution using locality metrics in cluster environments, *IEEE Transactions on Services Computing*, 2022
14. Ahmed, S., and Patel, K., Latency sensitive routing optimization in distributed storage networks, *IEEE Communications Letters*, 2024
15. Gupta, V., and Banerjee, S., Distributed data management with proximity driven placement, *ACM Symposium on Cloud Computing*, 2023
16. Liu, H., and Zhang, Y., Performance evaluation of locality aware routing in large scale clusters, *Cluster Computing Journal*, 2022
17. Morales, D., and Chen, X., Reducing communication overhead in distributed microservices architectures, *IEEE Software*, 2023
18. Kumar, S., and Reddy, V., Dynamic partition management for scalable cloud storage systems, *Future Generation Computer Systems*, 2024
19. Rahman, M., and Cho, H., Hop count reduction techniques for distributed databases, *IEEE Transactions on Network and Service Management*, 2022
20. Das, P., and Sharma, N., Efficient workload distribution in distributed environments using proximity metrics, *Journal of Cloud Computing*, 2023
21. Zhang, L., and Wu, J., Communication cost modeling for distributed data placement, *IEEE Transactions on Parallel and Distributed Systems*, 2022
22. Ortega, J., and Silva, M., Locality driven orchestration for containerized services, *ACM Middleware Conference*, 2023
23. Ibrahim, A., and Lee, C., Scalable network aware storage allocation in distributed systems, *IEEE Access*, 2024
24. Brown, T., and Green, R., Performance impact of communication distance in cloud scale infrastructures, *Journal of Supercomputing*, 2022