



Боремся с сетевым оверхедом в распределённых системах: альтернативный подход к работе с данными

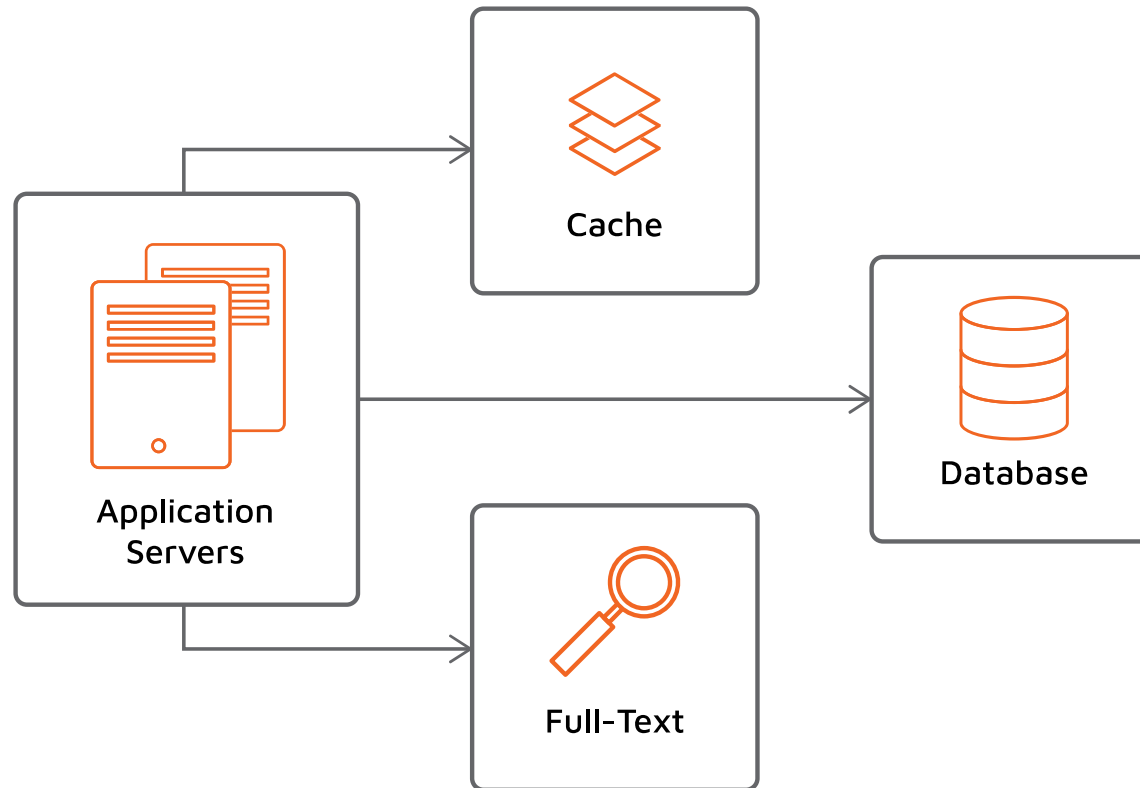
Pavel Tupitsyn

Lead .NET engineer at GridGain

.NET client maintainer at Apache Ignite

ptupitsyn@apache.org, [ptupitsyn.github.io](https://github.com/ptupitsyn)

Traditional Approach: Bring Data To Code



A lot of network calls for one user search query:

- Get item ids from ElasticSearch
- Get item descriptions from Redis
- If not in Redis, get from DB, update Redis

But Network Is Fast! Or Is It?



Latency Numbers Every Programmer Should Know



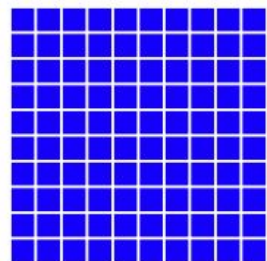
Main memory reference:
100ns



1,000ns \approx 1 μ s



Compress 1KB wth Zippy:
2,000ns \approx 2 μ s



10,000ns \approx 10 μ s = ■

~~Send 2,000 bytes over
commodity network: 44ns~~



SSD random read: 16,000ns \approx
16 μ s



Read 1,000,000 bytes
sequentially from memory:
3,000ns \approx 3 μ s



Round trip in same
datacenter: 500,000ns \approx
500 μ s

Connections Are Not Free



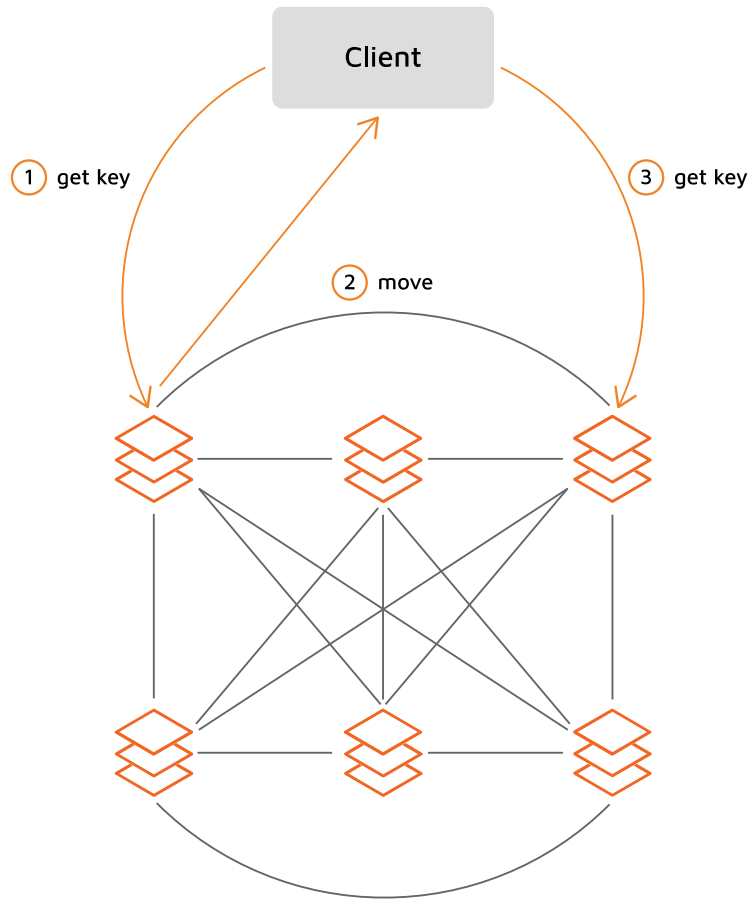
Problem: Connections kill

- 80 – 100 processes on each web server
- Each datacenter has thousands of servers
- **400K** connections to any memcached server
- Solution: Client and Server use **UDP**

UDP means never having to say ACK

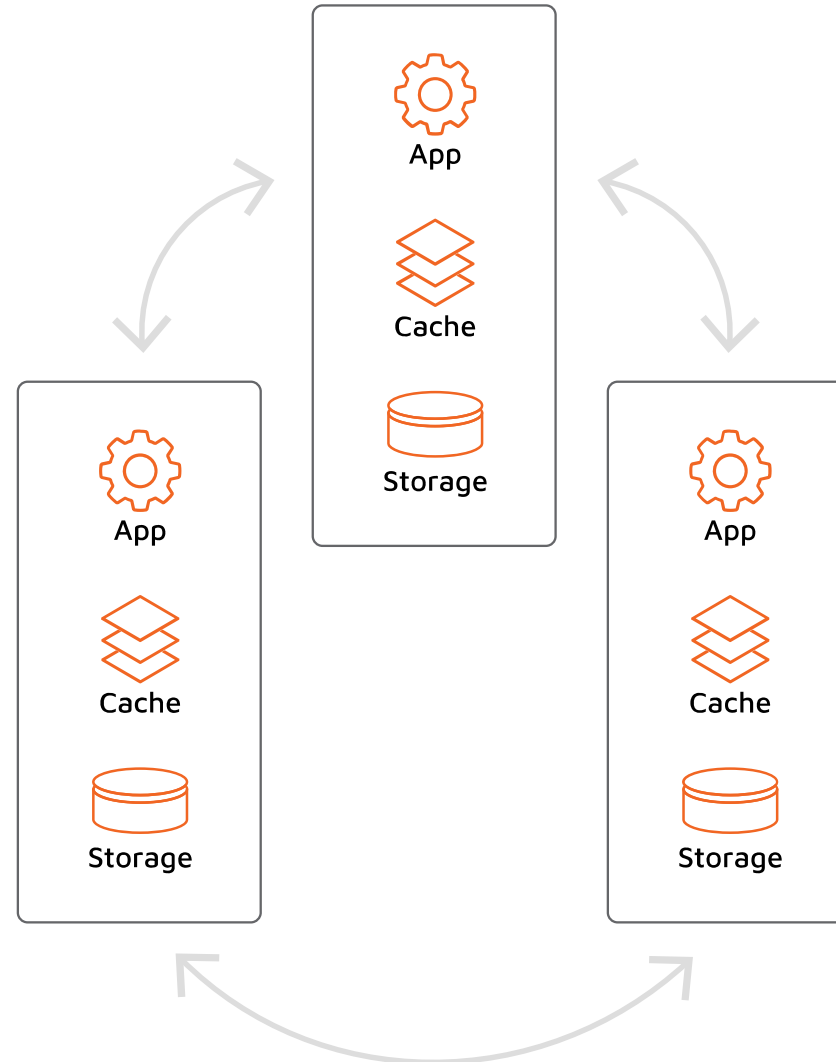
Facebook and memcached - Tech Talk
<https://www.youtube.com/watch?v=UH7wkvcf0ys>

What if we store the data together with code?



gossip protocol

VS



Local Redis vs In-Process Ignite

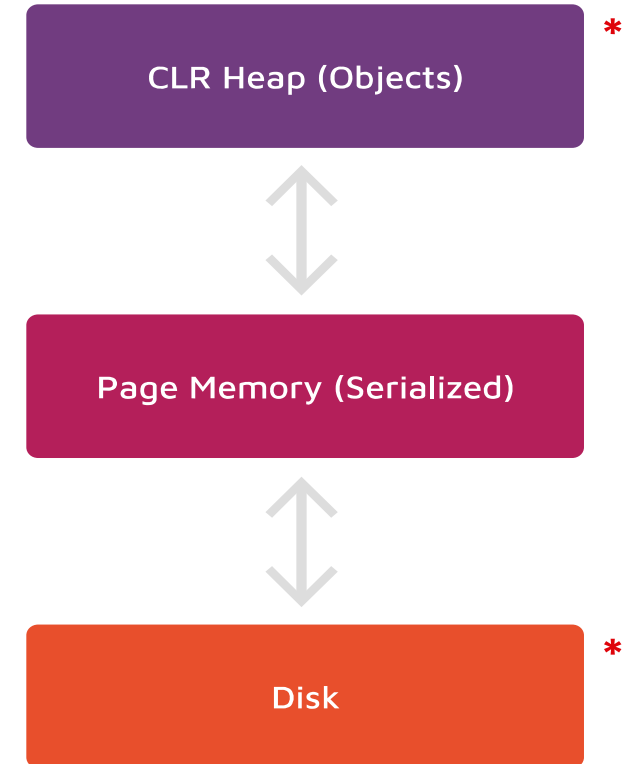


Method	Mean	Error	StdDev	Ratio	RatioSD
GetIgniteClrHeap	46.33 ns	0.305 ns	0.286 ns	1.00	0.00
GetIgnite	1,929.56 ns	36.238 ns	33.897 ns	41.65	0.84
GetRedis	60,080.39 ns	536.235 ns	475.358 ns	1,296.69	11.13

What is Apache Ignite?



- Distributed Database
- Transactional
- Automatic data partitioning (sharding)
- In-Memory + On Disk
- NoSQL + SQL / LINQ + FullText
- Streaming, messaging (pub/sub)
- Compute (map/reduce)
- Embedded Mode for .NET, Java, C++



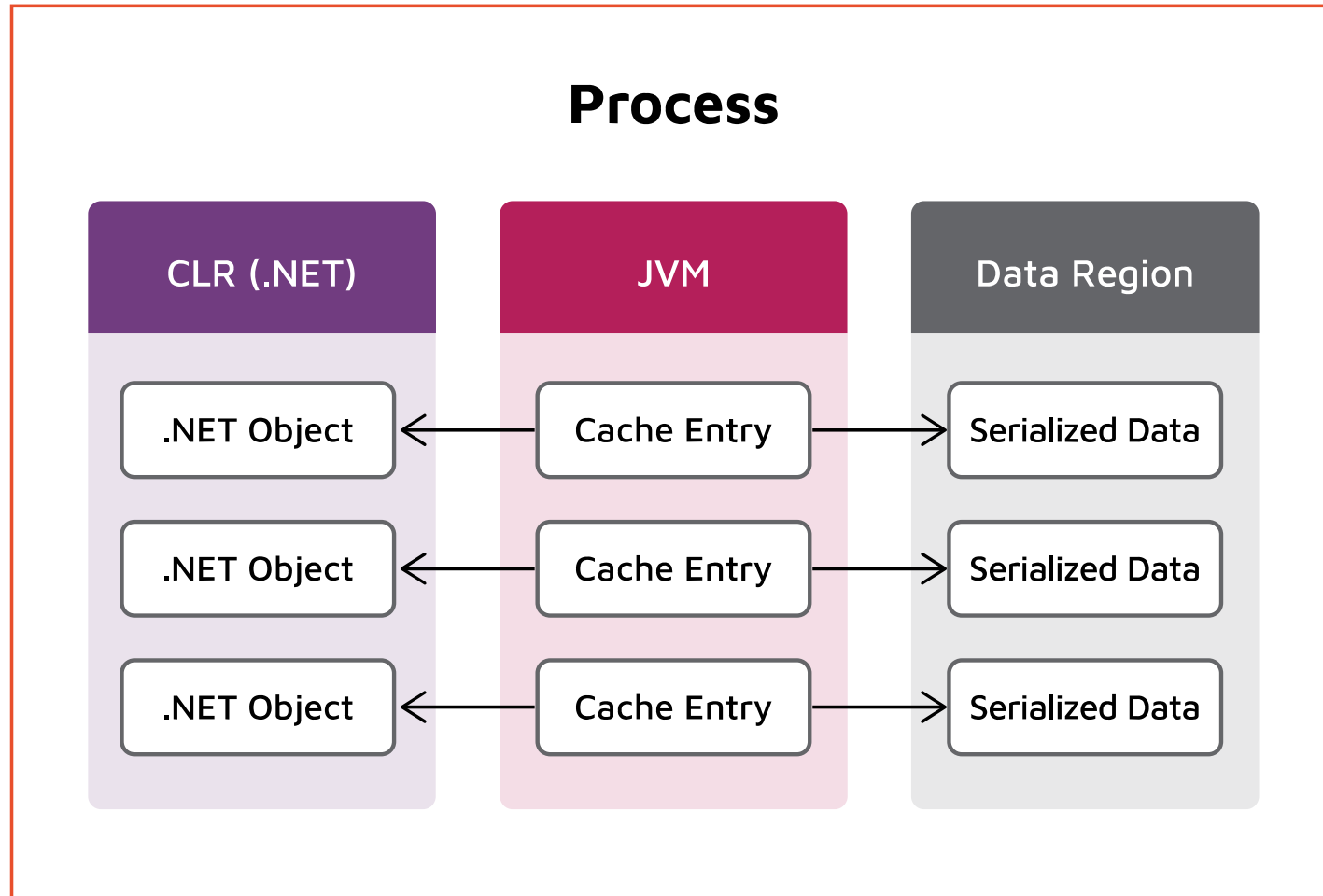
* Optional

Apache Ignite: Embedded Mode



- App + DB in one process (like SQLite)
- Data is in the process memory (“Page Memory” - unmanaged)
- Optionally on local disk
- Ignition.Start(config) to run
- Your app is the node: run on multiple machines to form a cluster
- Nodes find each other (address list / multicast / k8s discovery)
- Embedded mode is not the only choice:
 - Thin clients (.NET, C++, Java, JS, Python, PHP, Rust..)
 - ODBC
 - REST

Demo: Caching

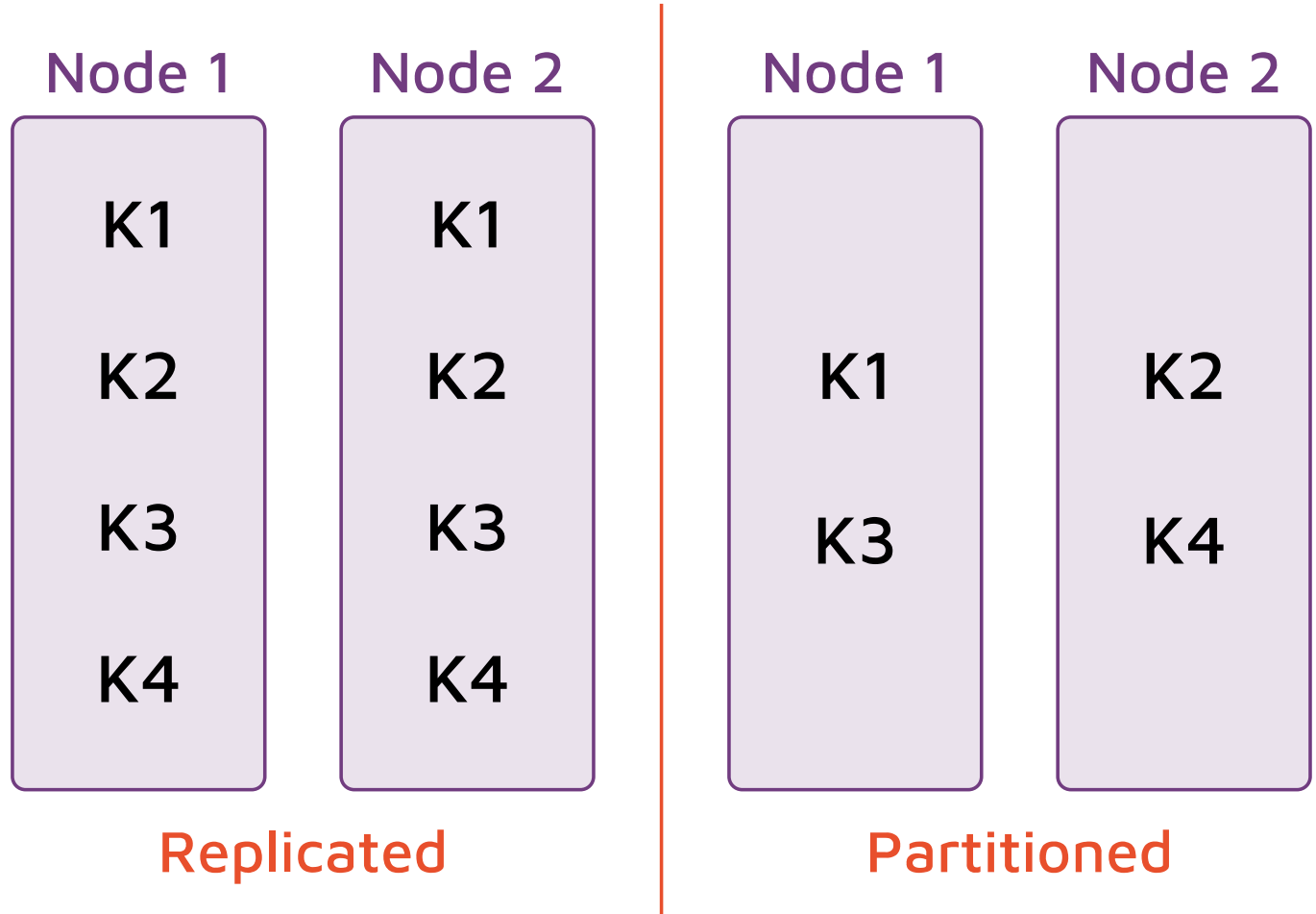


Caching Summary



- In-process cache avoids I/O costs
- Platform (CLR) cache reduces serialization costs
 - When writes < reads
- Replicated cache survives node loss

Replicated and Partitioned Storage



<https://ignite.apache.org/docs/latest/data-modeling/data-partitioning>
https://en.wikipedia.org/wiki/Rendezvous_hashing

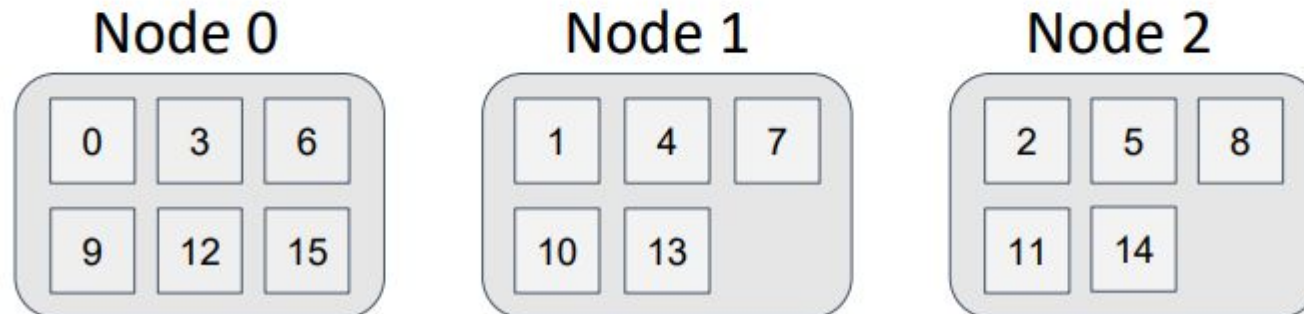
Partitioned Mode Explained



$$h(K) \bmod N$$

where

- K - object key
- h - hash function
- N - amount of nodes

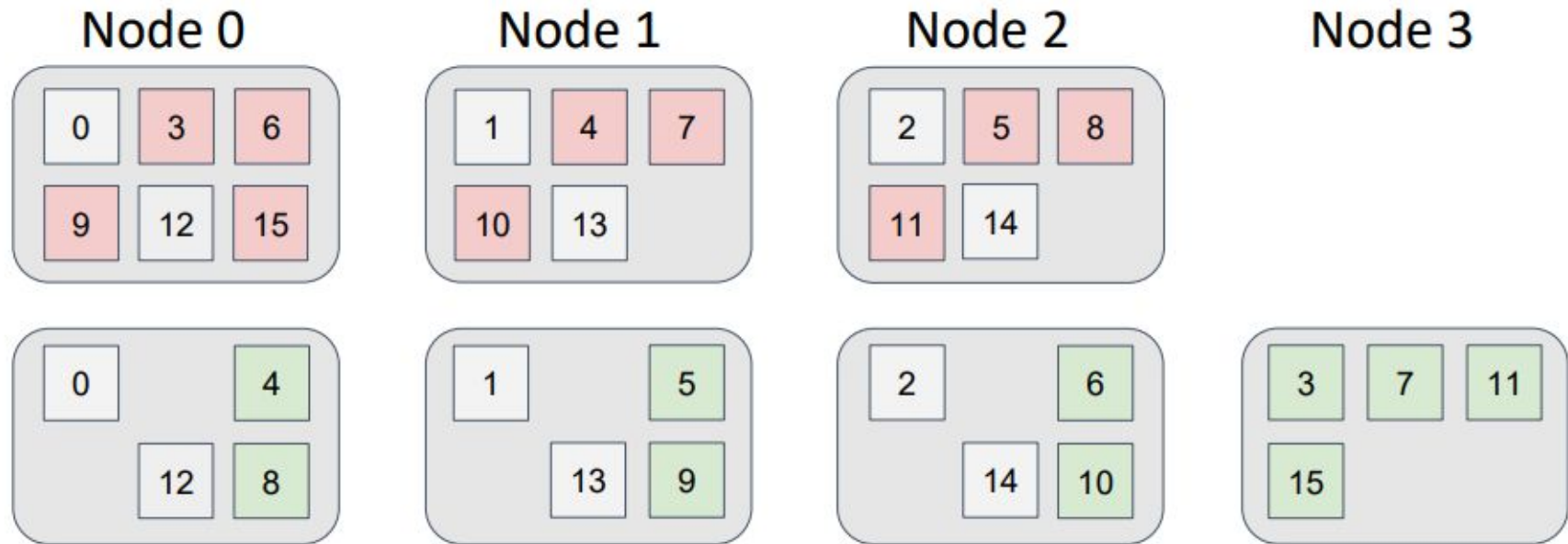


Partitioned Mode Explained



$h(K) \bmod N$

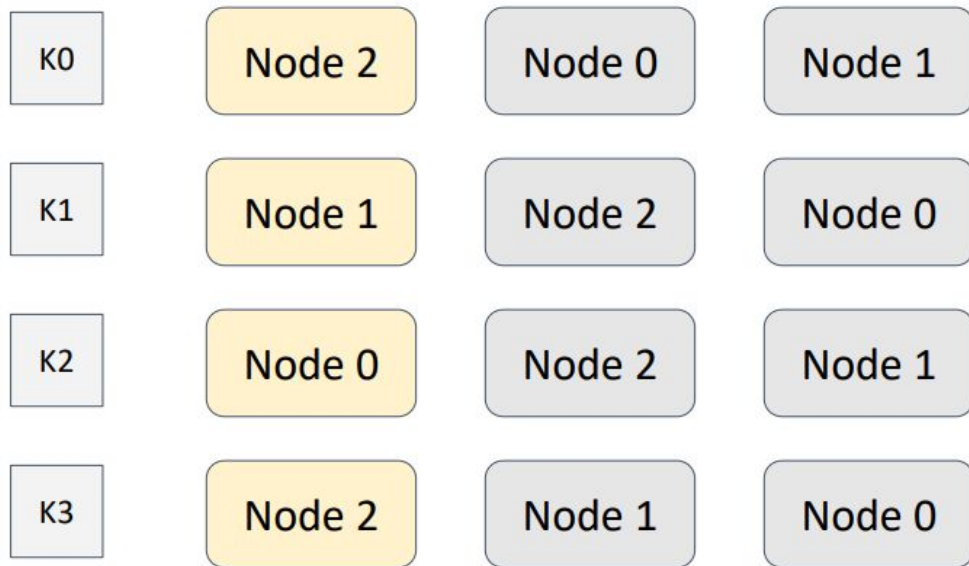
Add new node and ...



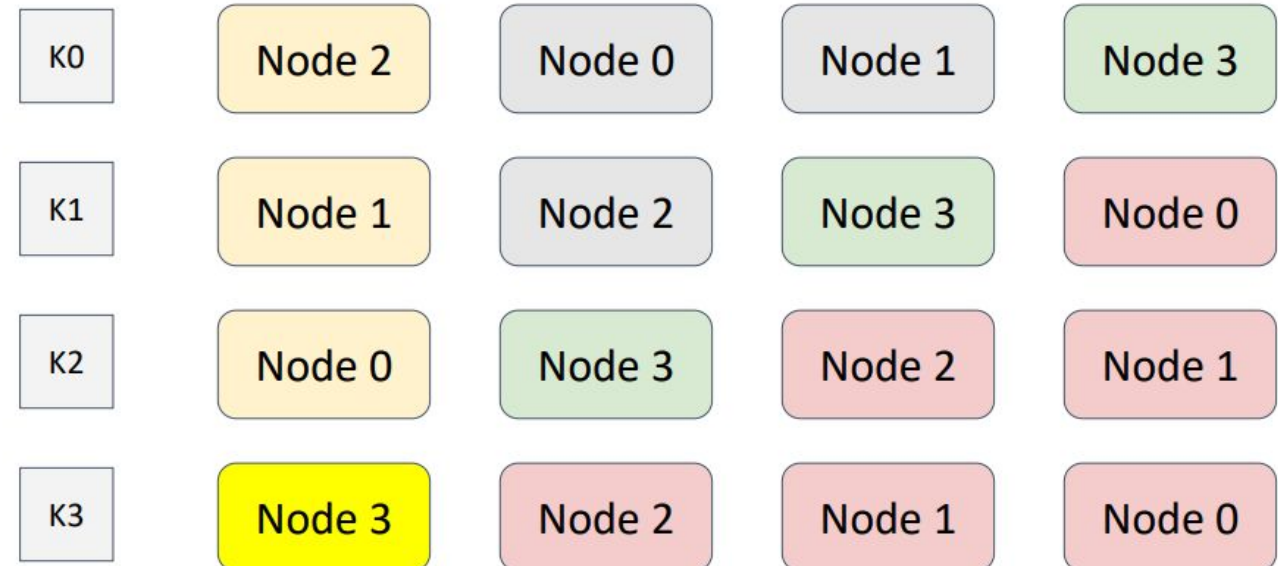
Partitioned Mode Explained



Distribution example: 4 keys, 3 nodes initially

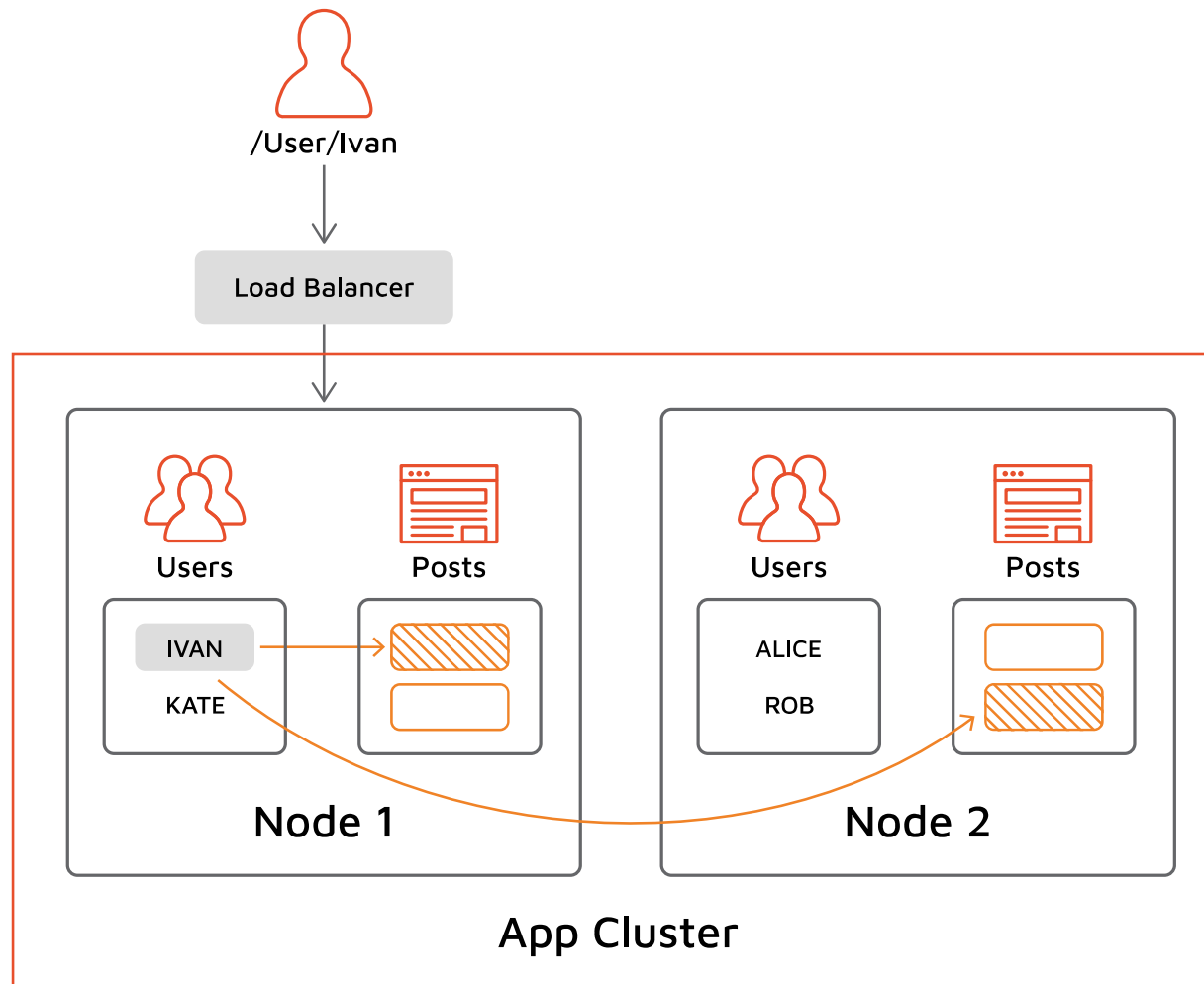


Distribution example: 4 keys, new node added



```
nodes.OrderBy(n => HashCode.Combine(n, key))
```

Demo: Data Colocation

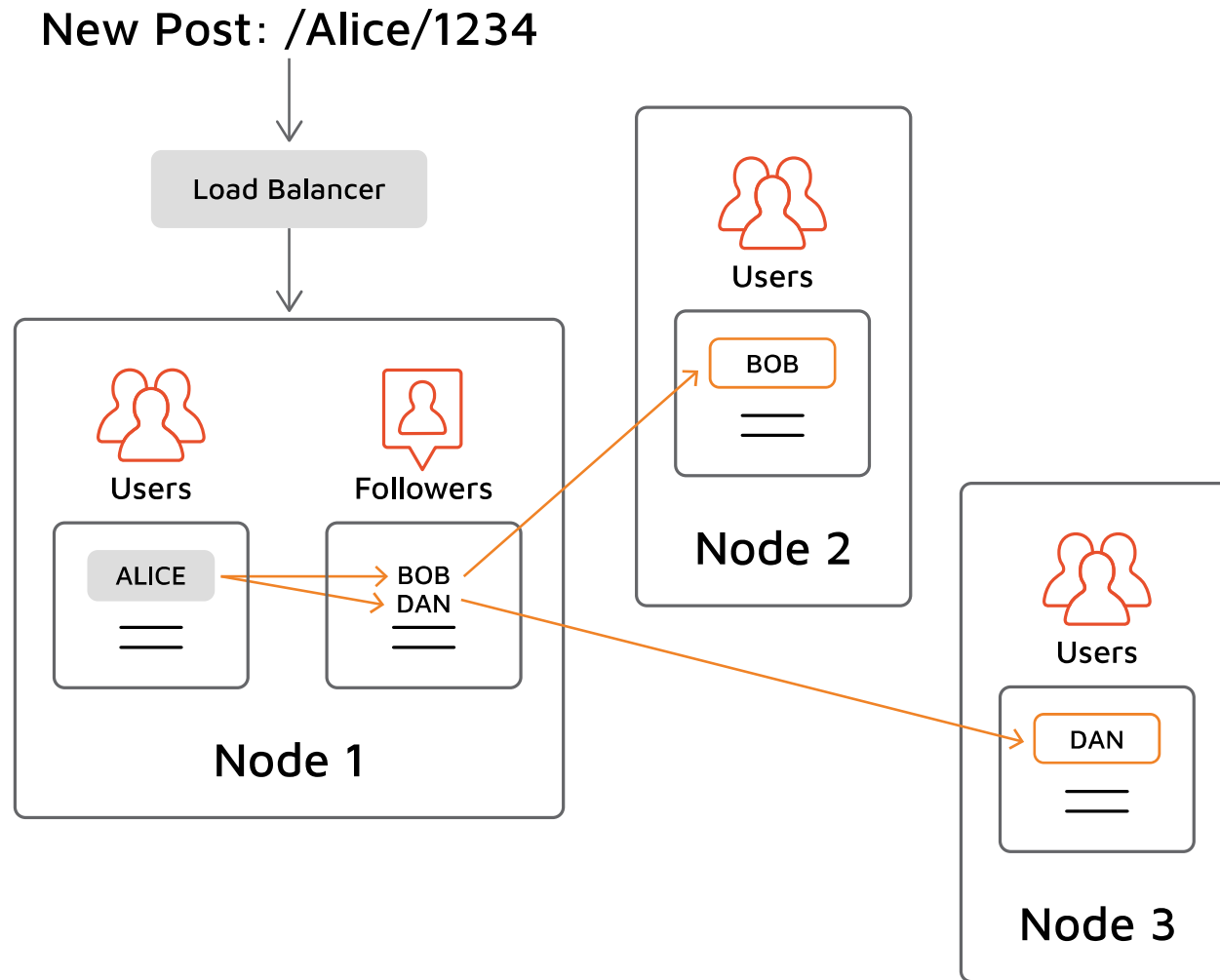


Data Colocation Summary



- Partitioned (sharded) mode distributes data across nodes based on Key by default
- Use AffinityKey or [AffinityKeyMapped] to distribute data based on a different field
- Keep related data together
- You can't colocate everything - optimize for important scenarios

Demo: Send Code To Data - Invoke

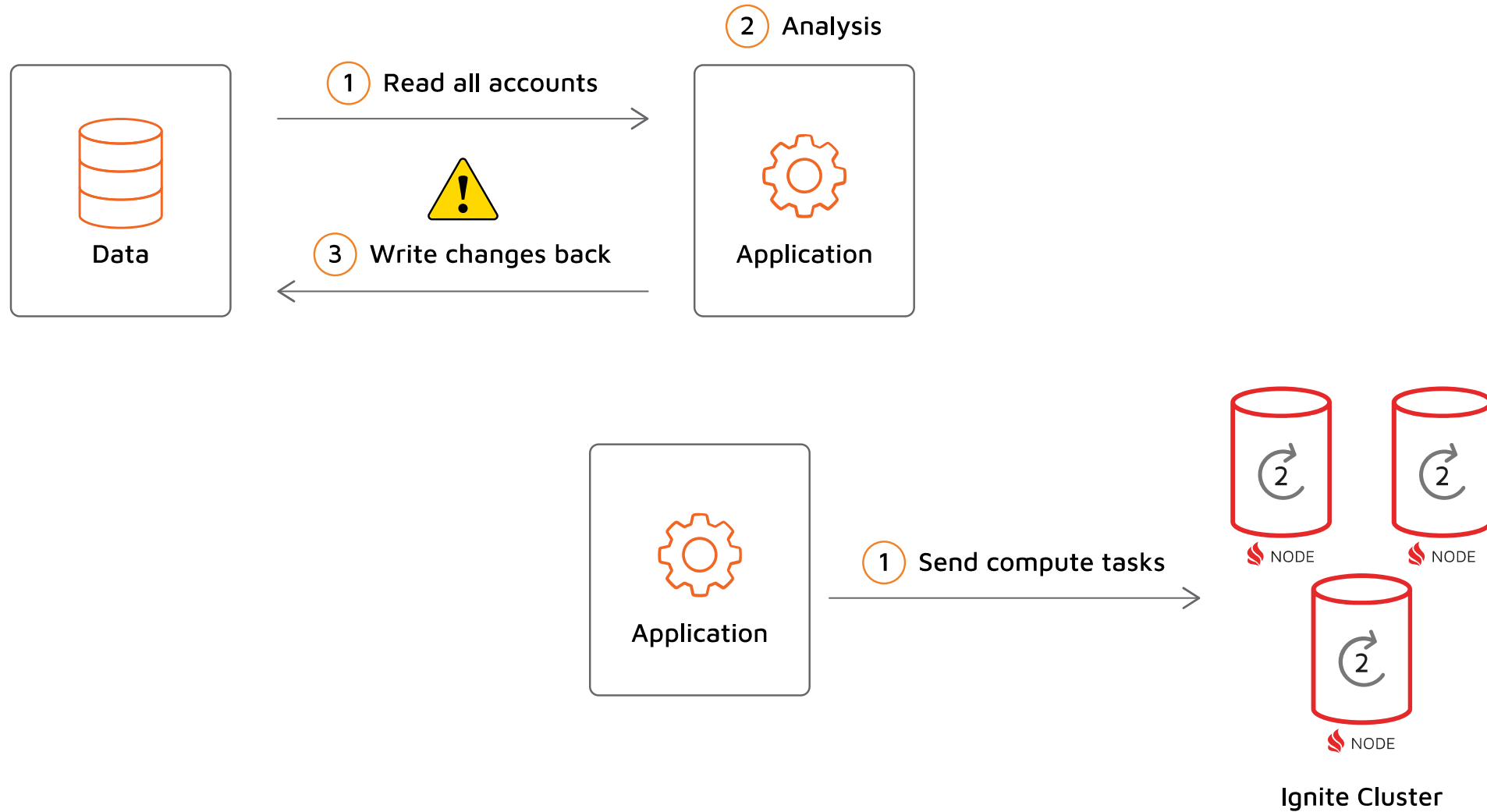


Send Code To Data Summary

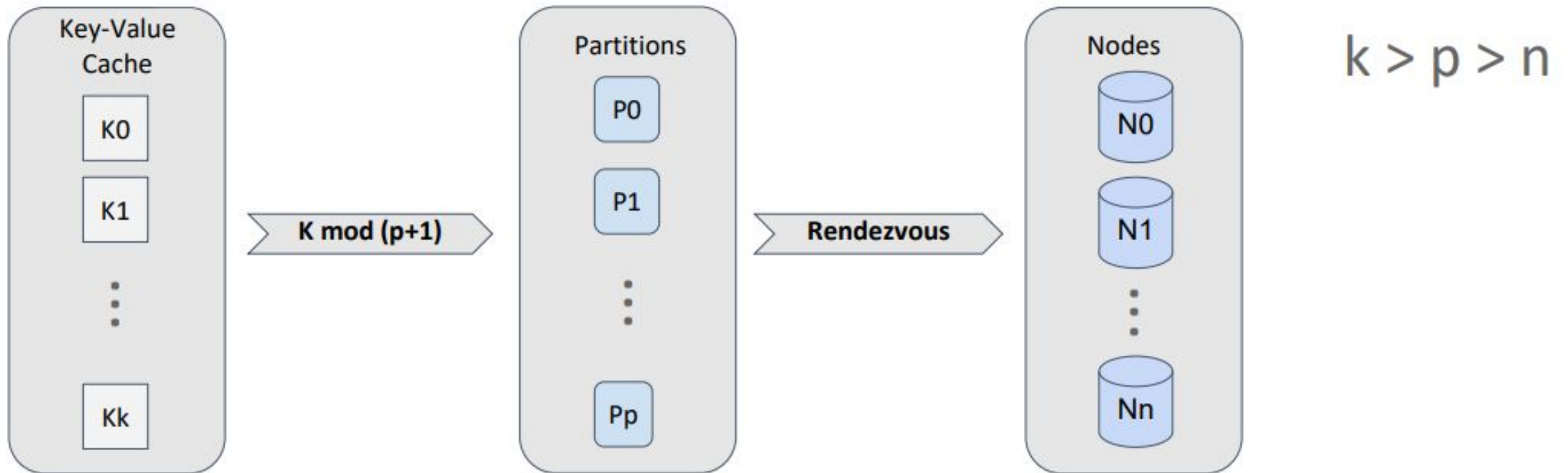


- Similar to SQL, e.g. “SELECT AVG ...”
- Not all computations can be expressed in SQL

Send Code To Data: Compute



Partitioned Cache: What is a Partition?

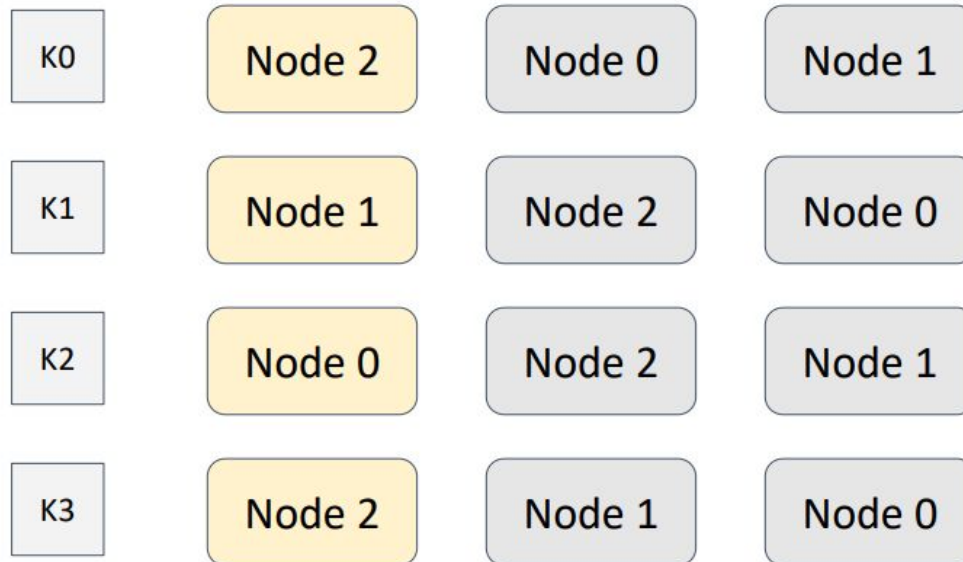


- Similar to hash table bucket

Partitioned Mode Explained - Reminder

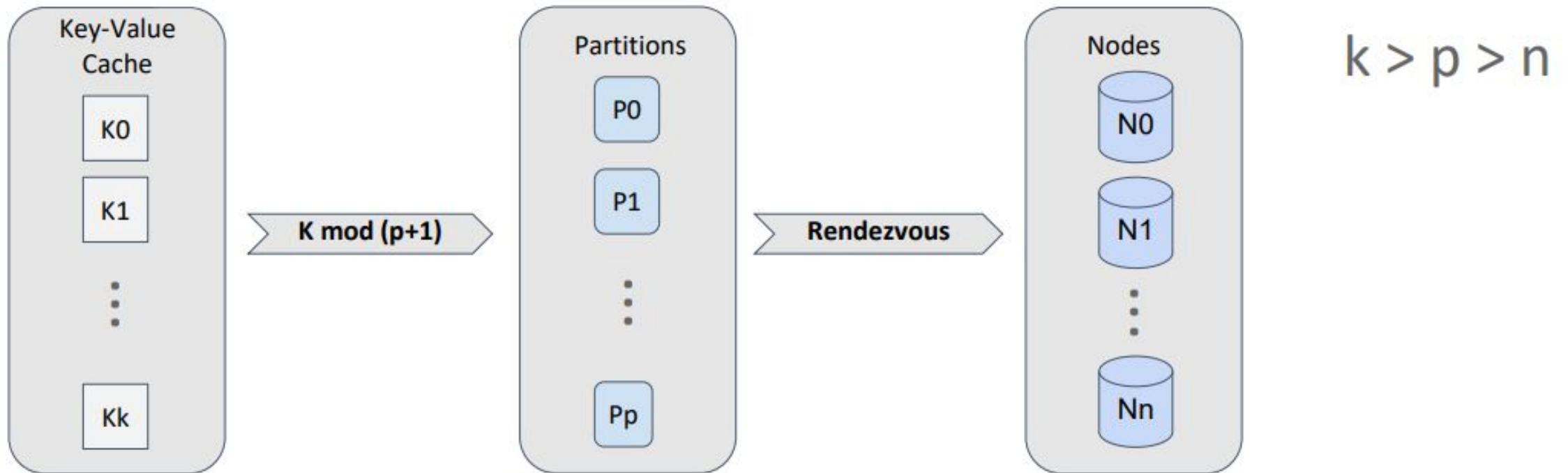


Partitions



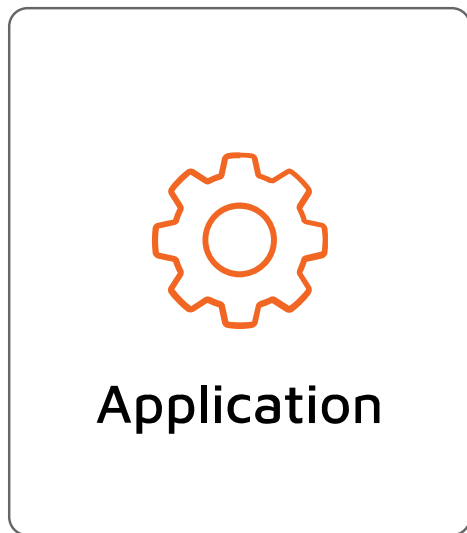
```
nodes.OrderBy(n => GetHashCode.Combine(n, partition))
```

Partitioned Cache: What is a Partition?

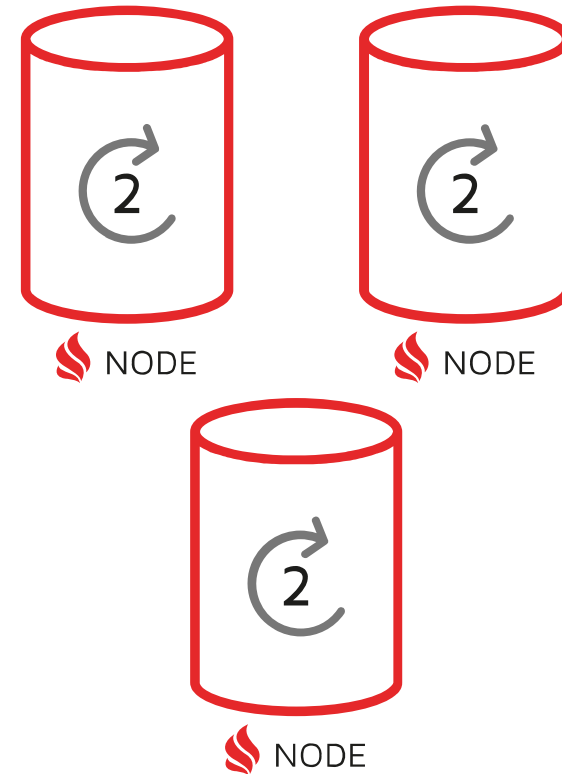


- Simplifies per-key operations: **mod**
- Fixed number of partitions => distribution can be cached once per topology change
- Partition is moved as a batch
- Partition can be locked in place during the computation to guarantee consistency

Send Code To Data: Demo



1 Send compute tasks

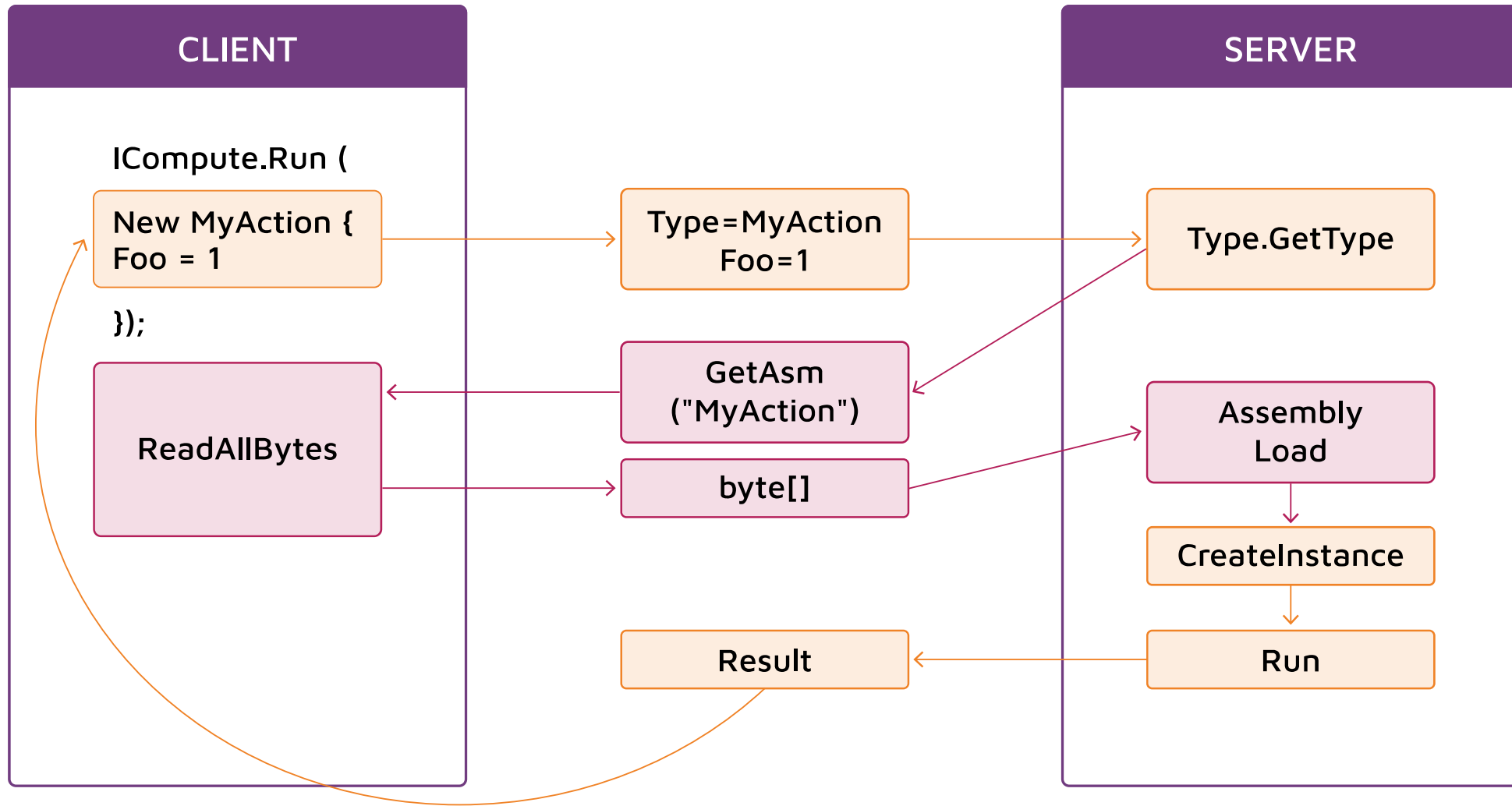


Send Code To Data: Summary



- Send code to when doing batch updates, changing large records
- Do not use for small changes - always take data size into account
- Topology may change and will change: use Affinity* APIs - they provide guarantees

Send Code To Data: How It Works



Summary



- Network calls are not free - far from it
- Usual approach with App and Data in separate clusters = overhead
- Combined (embedded) approach reduces this overhead
- Ignite simplifies combined approach
 - Automatic data distribution
 - Fault tolerance
 - Efficient APIs
- Choose data distribution strategy per use case
 - Replicated - more reads than writes, data fits in memory
 - Partitioned - data does not fit, more writes than reads
 - Backup copies
- Colocation: keep related data together
- Send code to data

Q&A



- Распределение данных в Apache Ignite
 - <https://habr.com/ru/company/gridgain/blog/489962/>
- Rendezvous Hashing
 - https://en.wikipedia.org/wiki/Rendezvous_hashing
- Demos and Benchmarks
 - <https://github.com/ptupitsyn/talks/>
- Apache Ignite
 - <https://ignite.apache.org>

Pavel Tupitsyn, GridGain

ptupitsyn@apache.org

ptupitsyn.github.io