

IndiLog

Distributed Indexing for Dynamic Scaling
in Stateful Serverless Computing with Shared Logs

Maximilian Wiesholler

Advisor: Dr. Florin Dinu

Supervisor: Prof. Dr.-Ing. Pramod Bhatotia

Chair of Decentralized Systems Engineering

<https://dse.in.tum.de/>



15.01.2022 – 15.07.2022

- Serverless functions
 - Ephemeral nature
 - Functions rely on service for state management
- Cloud storage services perform bad for state-sharing
 - May not offer consistency guarantees
 - Performance and costs trade-offs

- Recent development of new systems to improve state management
 - Cloudburst¹, Pocket², FaaS³, ...
- Boki: distributed shared logs as a promising serverless storage substrate⁴
 - Failure resilience
 - Consistency guarantees

¹V. Sreekanti et al. “Cloudburst: Stateful Functions-as-a-Service.”

²A. Klimovic et al. “Pocket: Elastic Ephemeral Storage for Serverless Analytics.”

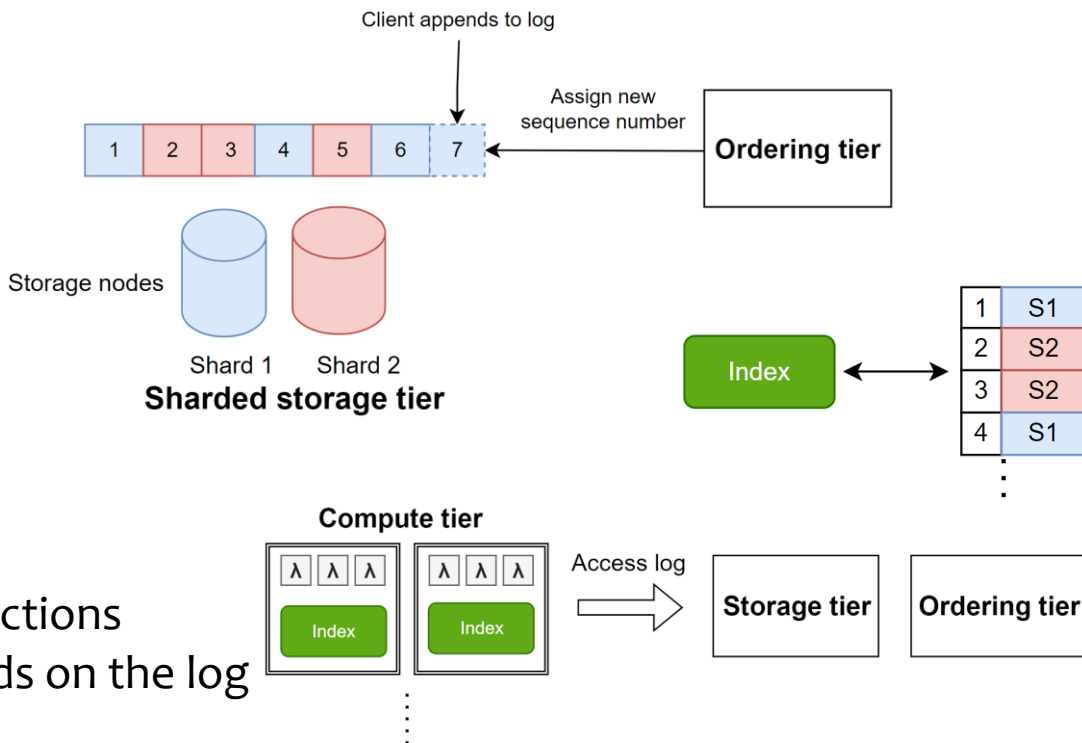
³F. Romero et al. “FaaS^T: A Transparent Auto-Scaling Cache for Serverless Applications.”

⁴Z. Jia et al. “Boki: Stateful Serverless Computing with Shared Logs.”

Background: distributed shared logs

Log: sequence of immutable log records; append-only

2 main tiers:
storage and ordering



For serverless

- Compute tier to run functions
- Indexes to locate records on the log

Distributed shared logs for serverless state management shift **indexing** into the focus

The current approach to indexing has the following limitations:

- Functions share resources with indexes
- Scalability of the compute tier is impeded by the design of indexing

How to design an efficient indexing architecture

- Do not impede the scaling of the compute tier by the design of indexing
- Limit resources for indexes co-located with functions on compute nodes

IndiLog:

a distributed indexing architecture for shared logs

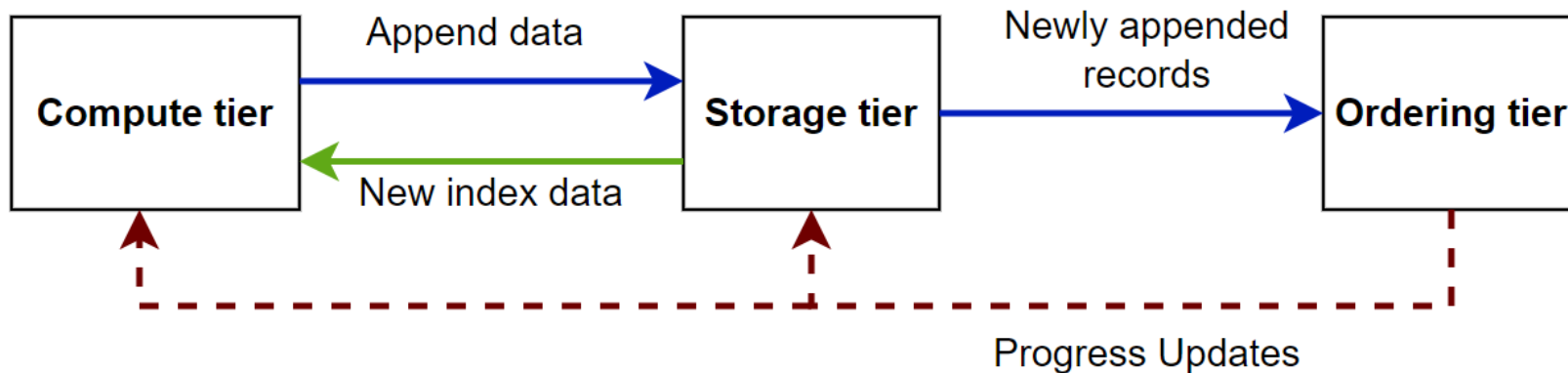
Compute nodes maintain optional, local and incomplete indexes
A sharded index tier balances the index data across index nodes and is complete

System design goals:

Performance:	many index lookups are captured locally
Resource efficiency:	local indexes are size-bounded
Scalability:	compute tier scalability is not impeded
Functions run anywhere:	no constraints where functions run

- ~~Motivation 1: high level~~
- Background
 - Boki
- Motivation 2: Boki's limitations
- Design: IndiLog
- Implementation
- Evaluation

State-of-the-art distributed shared log for state management of serverless functions¹

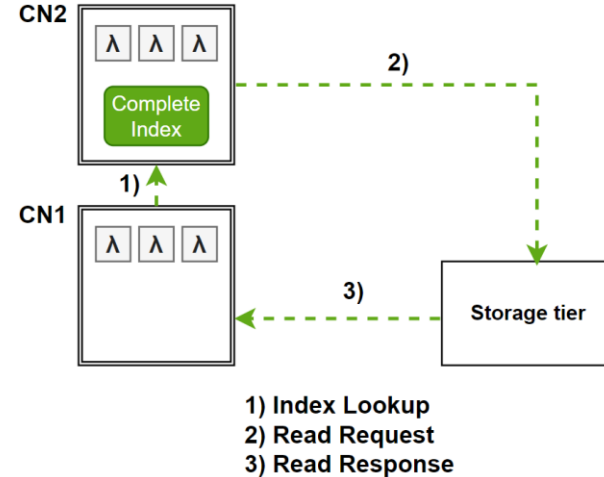
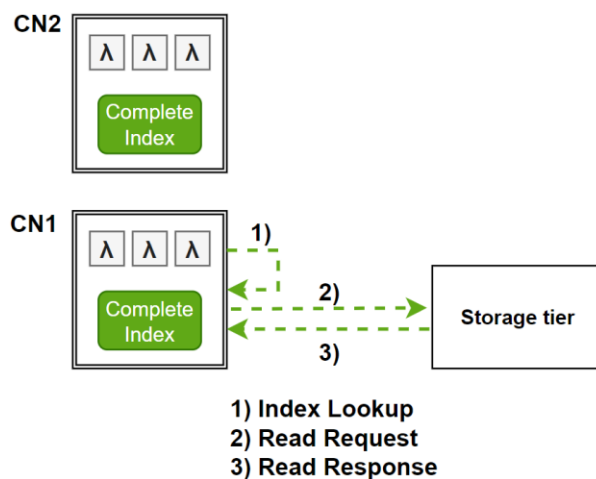


¹Z. Jia et al. "Boki: Stateful Serverless Computing with Shared Logs."

Boki: indexing

Indexes are complete and co-located with serverless functions on compute nodes

The complete index lets functions locate records on the storage tier



Boki uses **tags** to create logical sub-streams over the log¹

- Function needs only reading records that belong to the sub-streams

tag 1	1	4	6		
tag 2	2	3	8	9	10
tag 3	5	7			

Boki uses **bounded reads** which may not target a specific sequence number

Read tag 1 with $6 \leq X$

tag 1	1	4	6
-------	---	---	---

6 is **exact** match

Read tag 1 with $5 \geq X$

tag 1	1	4	6
-------	---	---	---

4 is **closest** match

Exact match

Bound is in sub-stream

Closest match

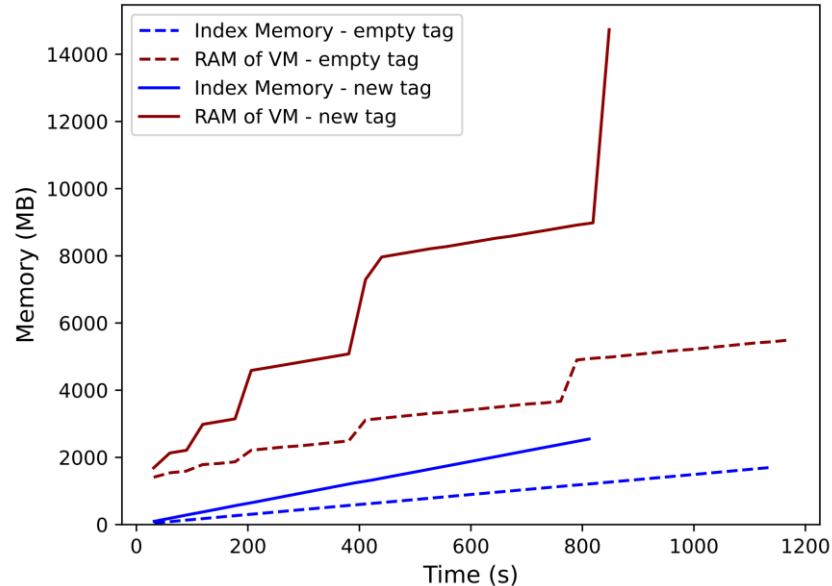
Bound is not in sub-stream

¹M. Balakrishnan et al. "Tango: Distributed Data Structures over a Shared Log."

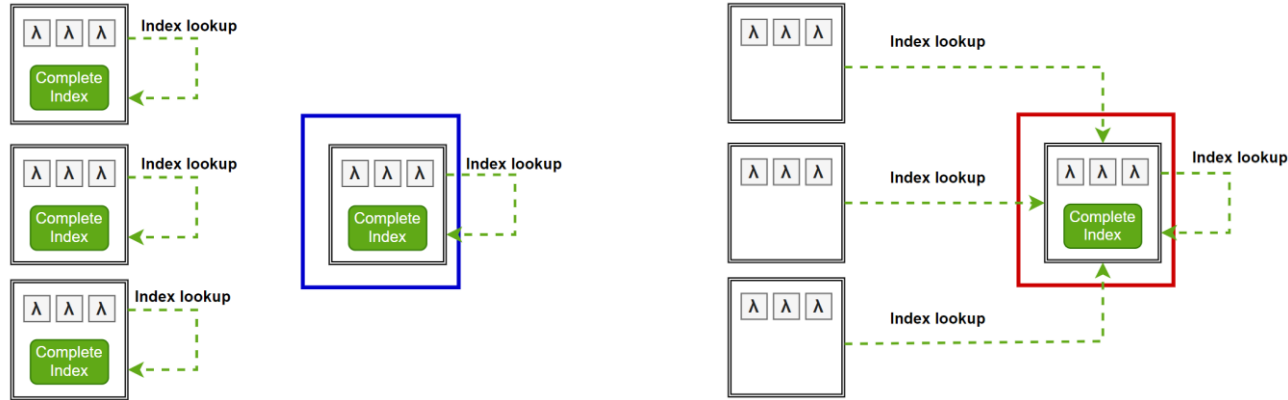
- ~~Motivation 1: high-level~~
- ~~Background~~
- Motivation 2: Boki's limitations
- Design
- Implementation
- Evaluation

Boki: out-of-memory crashes in compute nodes

Memory usage of indexes increase over time and the compute nodes eventually crash



Boki: index lookups add high contention after scaling



Increase workload →

Throughput [kOp/s]	40.73	60.44	71.48	77.11	93.09	100.82	102.9
Throughput [kOp/s]	33.84	46.28	52.93	56.66	63.65	63.28	63.78
Ratio [%]	83.1	76.6	74.0	72.1	68.4	62.8	62.0

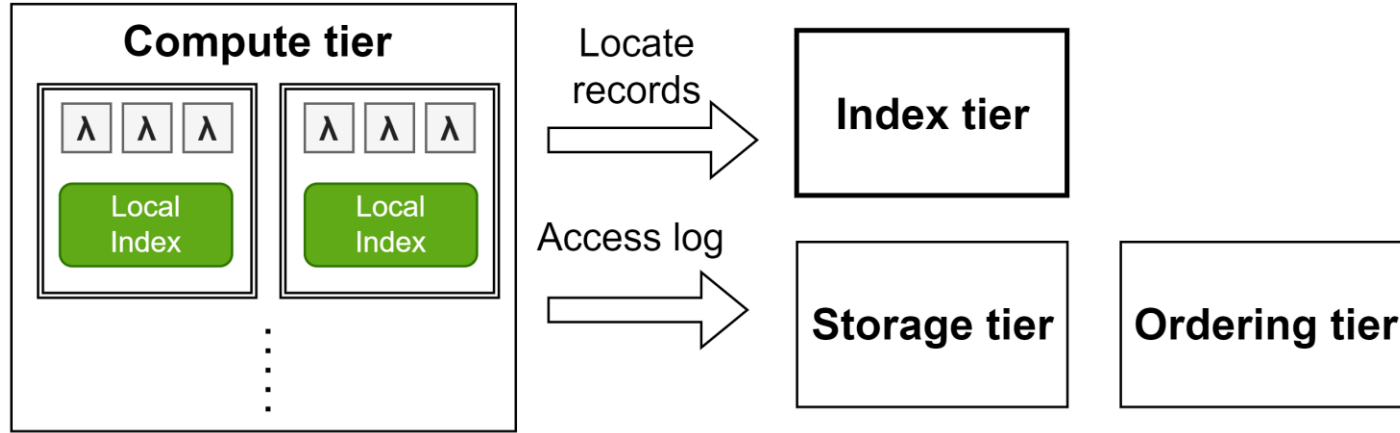
- ~~Motivation 1: high-level~~
- ~~Background~~
- ~~Motivation 2: Boki's limitations~~
- Design: IndiLog
- Implementation
- Evaluation

IndiLog: system overview

4 tiers: compute, ordering, storage and **index**

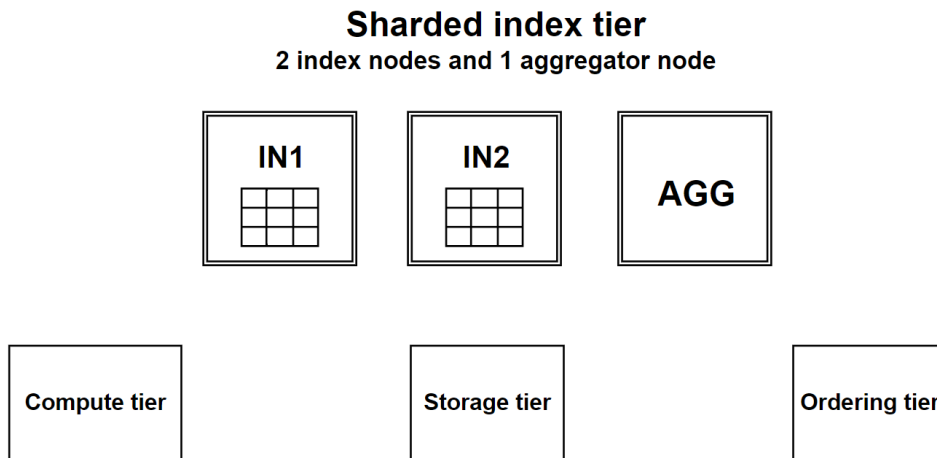
Compute nodes have optional, **local indexes** to capture most of the index lookups

- size-limited by eviction policies



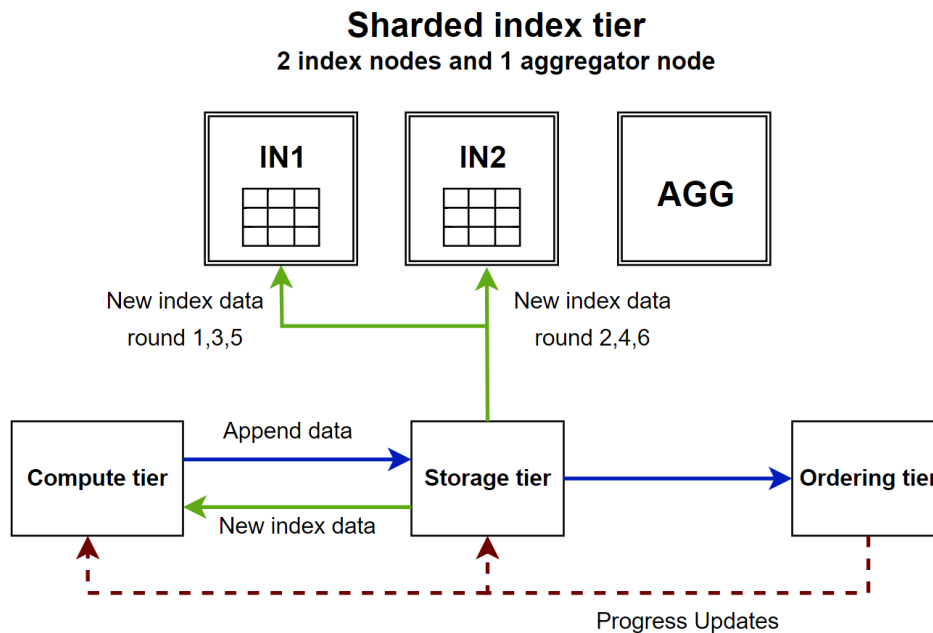
Design of the index tier

- Sharded:** indexes on any two index shards do not intersect
- Completeness:** the index tier can serve all index lookups
- Aggregating:** aggregator determines the closest match from index lookups



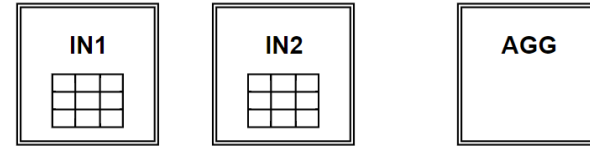
Append path

New index data is sent to the index tier in round-robin fashion to balance the distribution

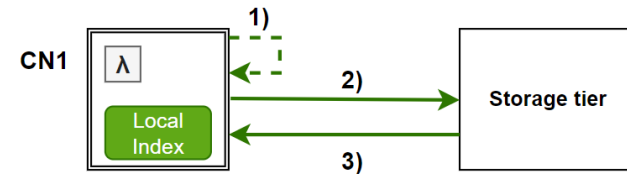


In IndiLog we have three types of reads

Type 1: the local index of a compute node has a match for the index lookup



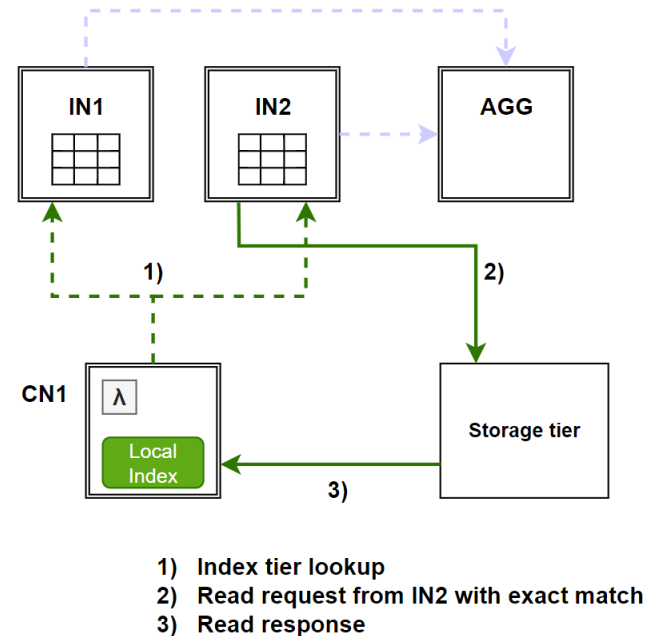
No lookup in index tier:
2 * one-way network latency



- 1) Index lookup with local hit
- 2) Read request
- 3) Read response

Type 2: the index lookup goes to the index tier – one index shard has an **exact match**

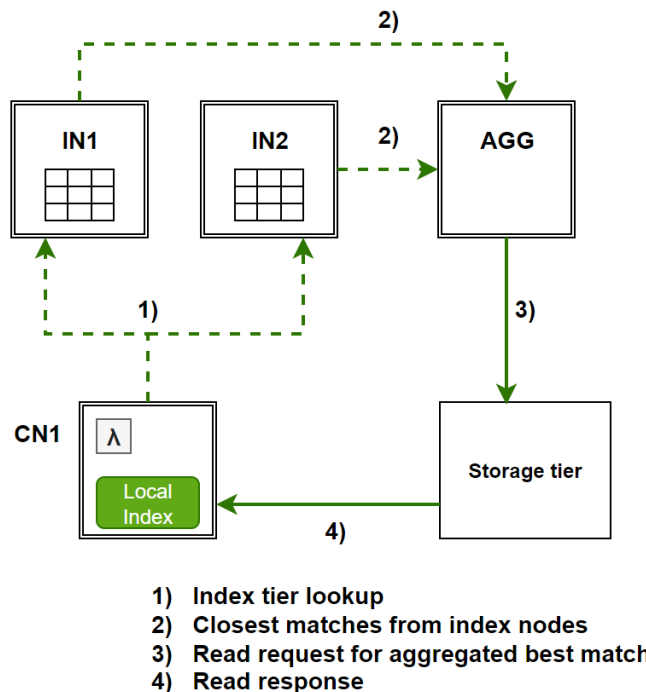
Index node with exact match can
safely forward read request:
3 * one-way network latency



Type 3: the index lookup goes to the index tier – all index shards have **closest match**

Aggregator **aggregates the best match**

Aggregator forwards read request:
 $4 * \text{one-way network latency}$



- ~~Motivation 1: high-level~~
- ~~Background~~
- ~~Motivation 2: Boki's limitations~~
- ~~Design: IndiLog~~
- Implementation
- Evaluation

IndiLog is built in C++ by adding an indexing architecture on top of Boki

- Index data structures implemented with the C++ Standard library and Abseil libraries¹
- Index data structures use a coarse-grained locking
- Aggregator uses fine-grained, per-request locks (data structures from oneTBB²)
- ZooKeeper³ for configuration and dynamical scaling of the compute tier

IndiLog: 7,892 LoC | **Benchmarks:** 19,131 LoC

¹Abseil libraries: <https://abseil.io/>

²Threading Building Blocks (oneTBB): <https://oneapi-src.github.io/oneTBB/>

³Apache ZooKeeper: <https://zookeeper.apache.org/>

- ~~Motivation 1: high-level~~
- ~~Background~~
- ~~Motivation 2: Boki's limitations~~
- ~~Design: IndiLog~~
- ~~Implementation~~
- Evaluation

- What effects do we observe when we scale the compute tier of IndiLog?
 - Different workloads, various scaling sizes
- How does the index tier of IndiLog perform?
 - Type 2 and type 3 reads
- How does IndiLog behaves for real applications?
 - Object storage used on top of IndiLog

- Experimental setup:
 - Cloud VMs: RAM 16GB | vCPUs 4 | Ubuntu 20.04 | Kernel version 5.10.0
 - Latency: 180 μ s +/- 40 μ s Bandwidth: 2,127 Mbps
- Workload: append new records (1 KB) and mix reads for persisted records

We compare
IndiLog against Boki

	IndiLog	Boki
Compute Tier	4 VMs	4 VMs
Storage Tier	4 VMs	4 VMs
Ordering Tier	3 VMs	3 VMs
Index Tier	3 VMs (2 IN, 1 AGG)	-
Local Index	20 MB	complete i.e., 16 GB

Scaling the compute tier from 1 to 4

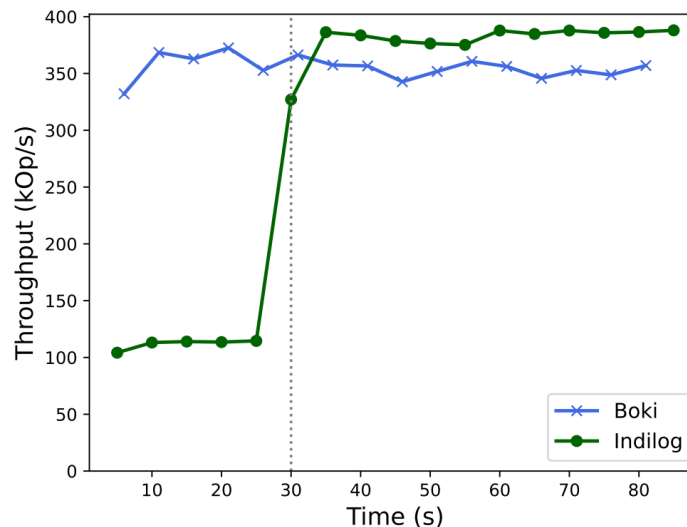
Workload: 50/50 Append/Read + 87% Local Index Hit Ratio in IndiLog

Indilog

- Starts with 1 compute node
- Scales to 4 compute nodes after 30 sec

Boki

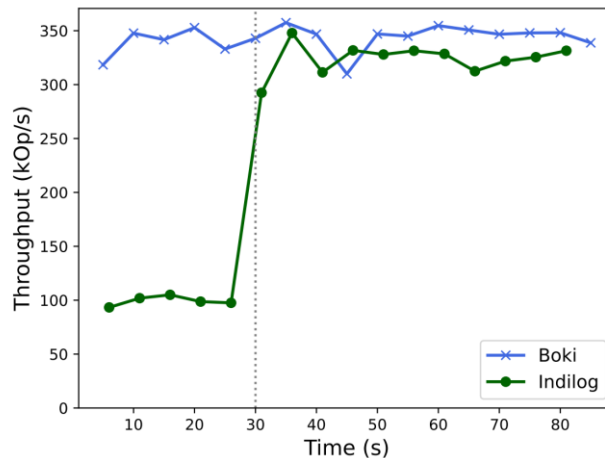
- Cannot scale dynamically
 - Configured as scaling happened
- only 1 compute node has a complete index
- 3 compute nodes send remote index lookups



IndiLog beats Boki when IndiLog captures many index lookups locally

Scaling the compute tier from 1 to 4

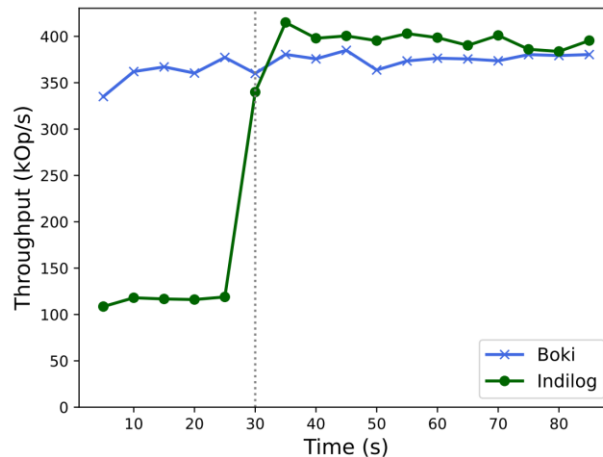
Workload: 50/50 Append/Read + 20% Local Index Hit Ratio in IndiLog



A low index hit ratio in IndiLog lowers the overall throughput

Scaling the compute tier from 1 to 4

Workload: 5/95 Append/Read + 20% Local Index Hit Ratio in IndiLog

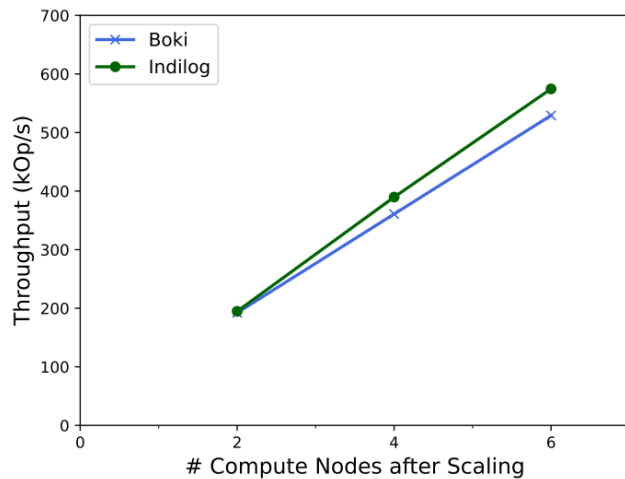


A read-heavy workload favors IndiLog:
Fewer write updates on index data structures

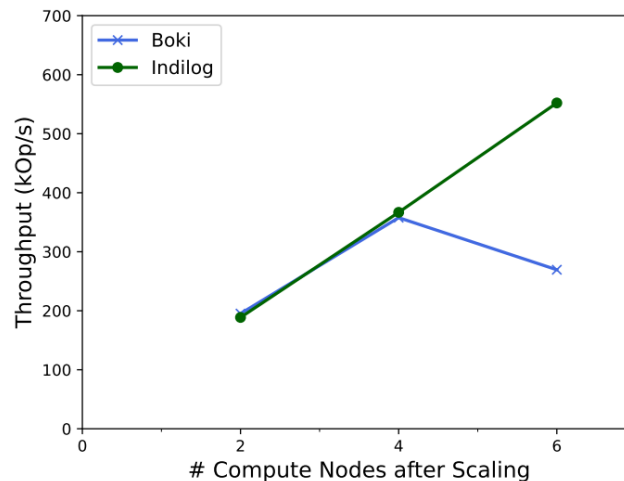
Scaling the compute tier from 1 to 2/4/6

87% Local Index Hit Ratio in IndiLog

50/50 Append/Read



5/95 Append/Read



Boki's node with the complete index gets under heavy contention

Read latencies of the index tier

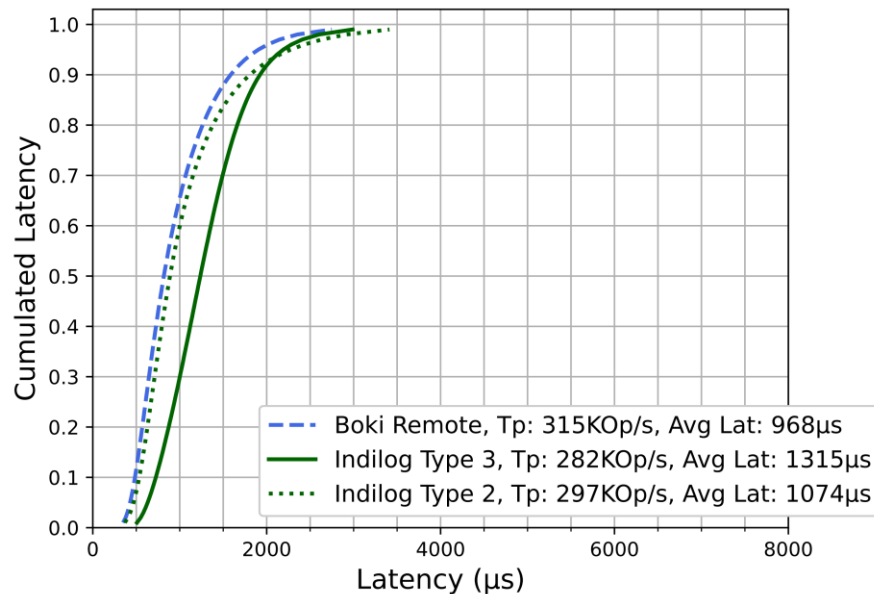
IndiLog

Local index disabled

- All index lookups go the index tier

Boki

- Few compute nodes maintain complete indexes but do not run functions
- Compute nodes with functions do remote index lookups only



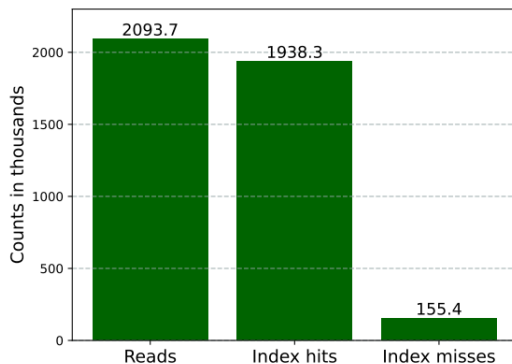
IndiLog's sharded index tier comparable to remote complete indexes

Real application: object storage

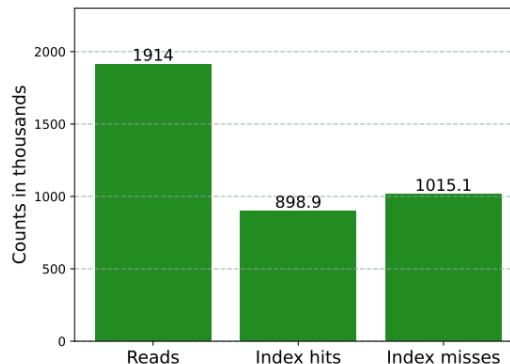
IndiLog as infrastructure layer of an object storage library with transaction support

- Workload: functions doing CRUD operations on key-value objects

Default index size (~20 MB)



Small index size (~0.2 MB)



System	Throughput [Op/s]
Boki (complete index)	8950.9
IndiLog (default)	8700.1
IndiLog (small)	8381.4

IndiLog is able to capture 93% of index lookups locally

Even with a very small index IndiLog captures almost 50% of lookups locally

Current state-of-the-art shared logs neglect efficient indexing

- Boki's complete index:
 - Leads to high RAM consumption and eventually OOM crash
 - Impedes scalability of the compute tier

IndiLog

- Local indexes + index tier for efficient indexing of log records
- Dynamic scaling of the compute tier

IndiLog source code:

<https://github.com/MaxWies/IndiLog>

Benchmark source code:

<https://github.com/MaxWies/indilog-benchmarks>

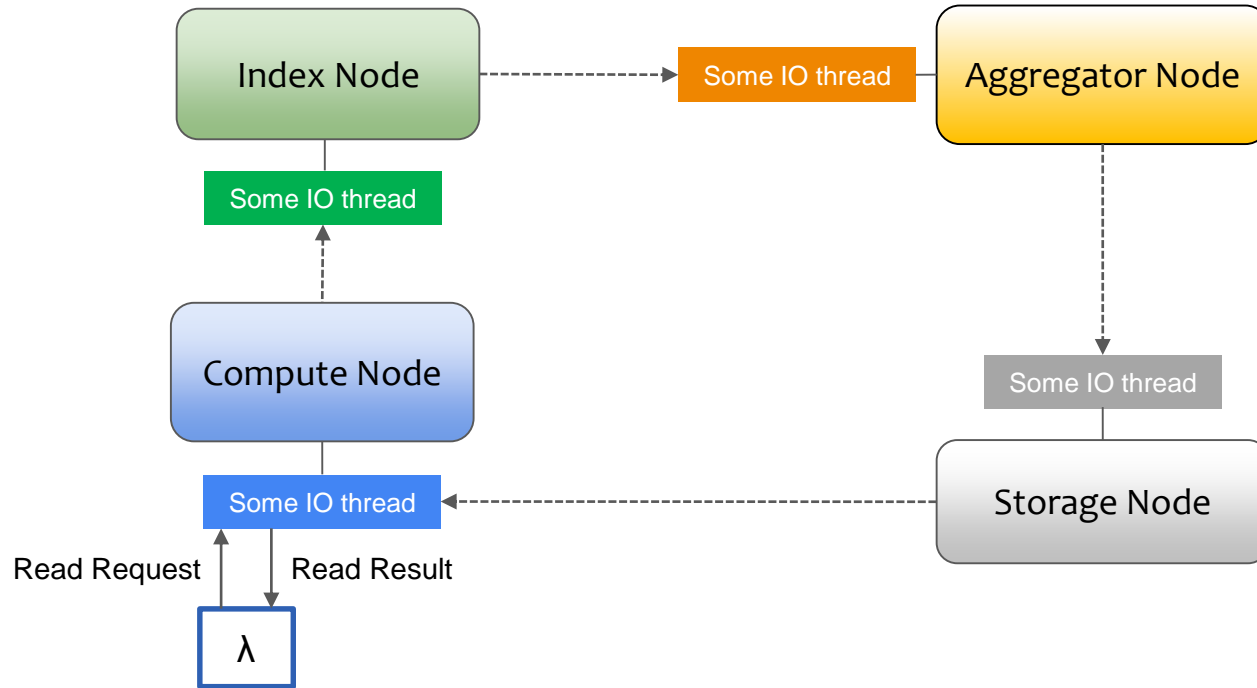
Discussion & Questions

Backup

Fine-grained functions run developers code on a Function-as-a-Service (FaaS) platform

Configuration and management of the platform is done by the platform provider

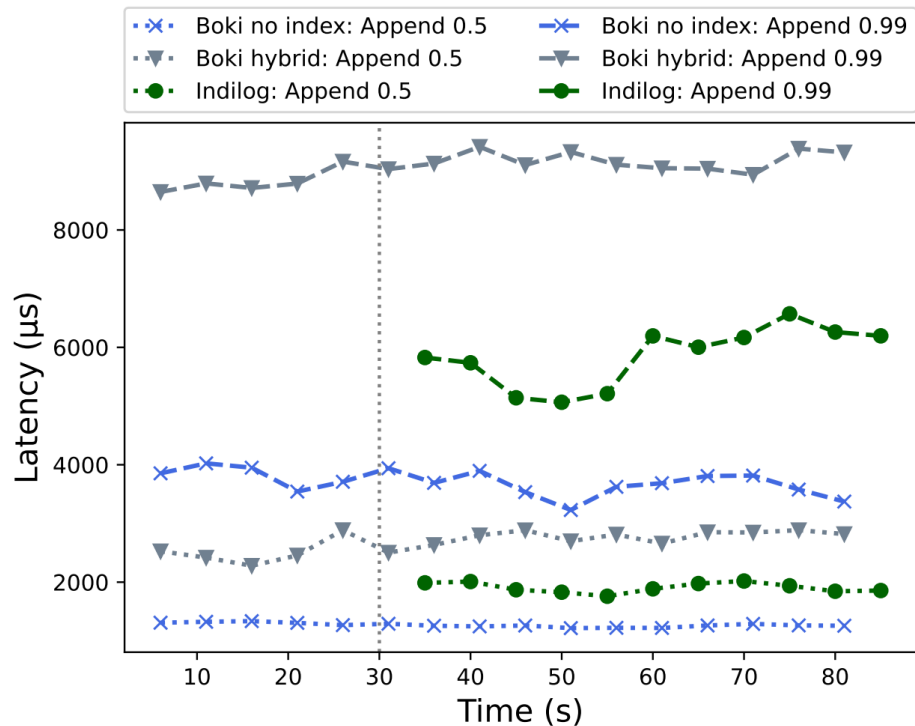
Threading model for reads (here: type 3)



Function waits until it receives result - Nodes process requests asynchronously

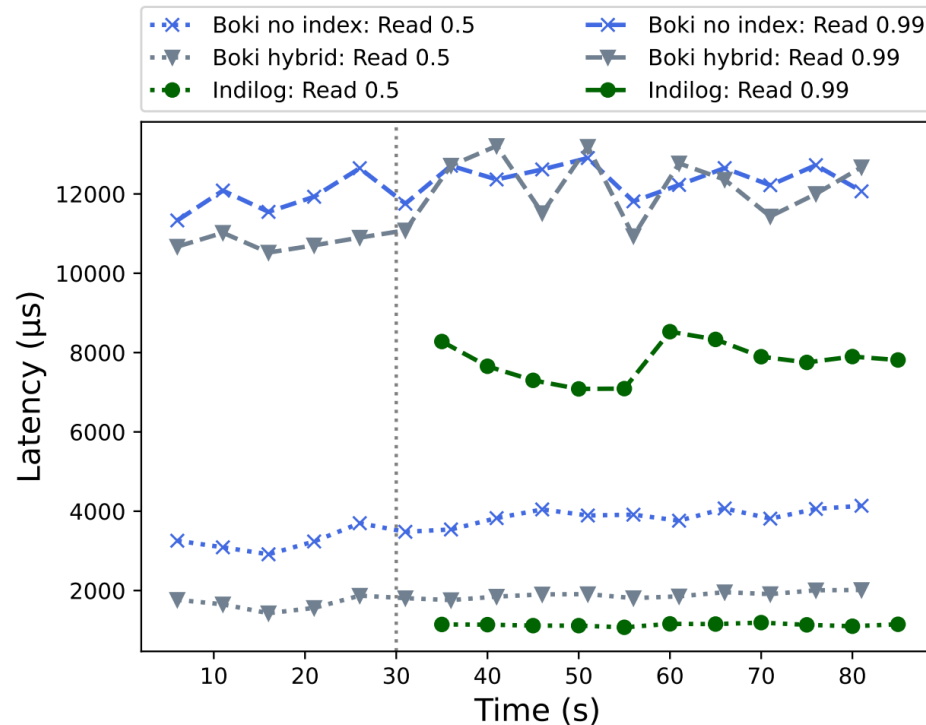
Scaling the compute tier from 1 to 4

50/50 Append/Read + 87% Local Index Hit Ratio in IndiLog

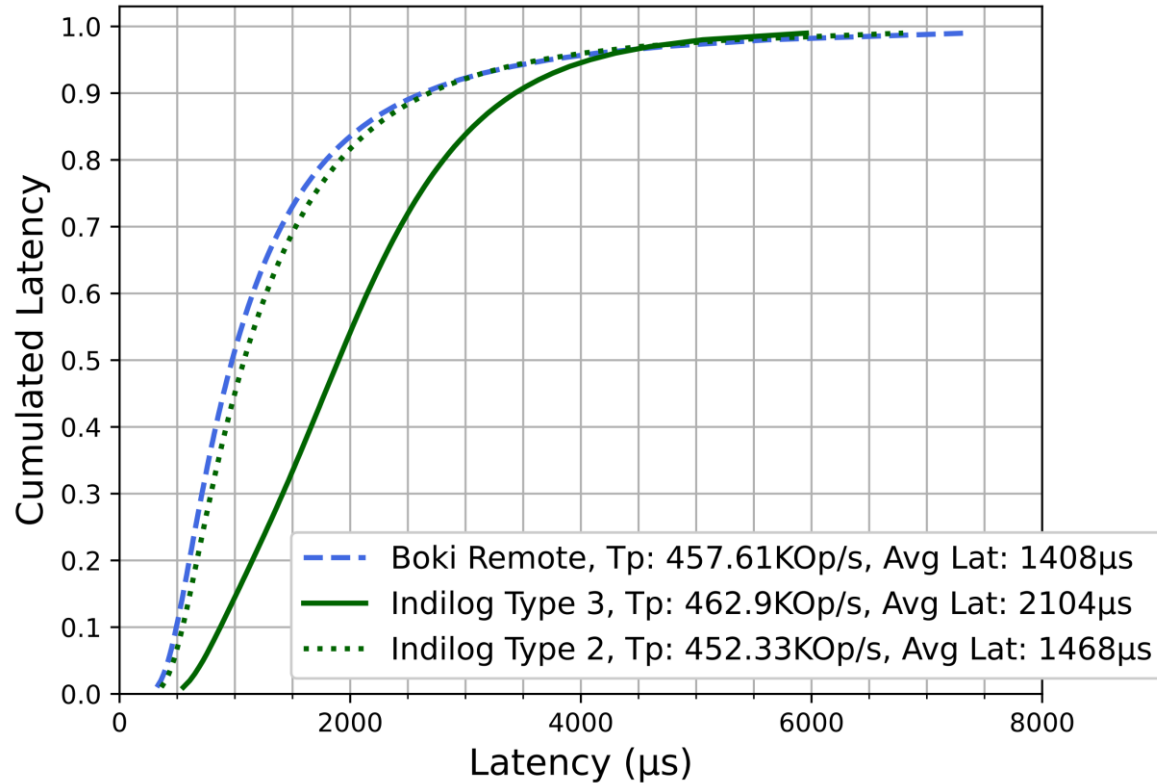


Scaling the compute tier from 1 to 4

50/50 Append/Read + 87% Local Index Hit Ratio in IndiLog



Read latencies of the index tier (high concurrency)



Scaling the index tier

Type 3 reads - all index lookups go to the index tier

S: Index shards

R: Replication factor

A: Aggregators

Index tier configuration	Throughput [kOp/s]
S:2 R:1 A:1	465.2
S:2 R:1 A:2	461.4
S:6 R:1 A:1	385.4
S:6 R:1 A:2	389.5

Throughput Comparison

System	Configuration	Throughput [Op/s]
Boki	4 compute nodes with complete indexes	8950.9
IndiLog	4 compute nodes, 20 MB local index	8700.1
IndiLog	4 compute nodes, 0.2 MB local index	8430.5
Boki	4 compute nodes with no indexes 2 compute nodes with complete indexes but no functions	8381.4