

## Introduction

The RELEASE project aims to improve the scalability of Erlang on emergent commodity architectures with 10<sup>5</sup> cores.

We anticipate that such architectures require scalable and available persistent storage on up to 100 hosts.

This research investigates the provision of persistent data structures by studying the ability of Erlang distributed DBMS to scale on our target 10<sup>5</sup> core architectures.

## Theoretical Analysis

**Mnesia** and **CouchDB** have scalability limitations:

- Explicit placement of replicas and fragments
- A single point of failure due to the lack of a P2P model

**Dynamo-style** NoSQL DBMS like Riak and Cassandra:

- Do have a potential to provide scalable persistent storage for large distributed architectures

## Scalability of Erlang DBMSs

### The requirements for scalable and available persistent storage

#### Replication

- Decentralized model
- Location transparency
- Asynchronous replication

#### Data Fragmentation

- Decentralized model
- Systematic load balancing
- Location transparency

#### Query Processing

- Location Transparency
- Local Execution
- Parallelism



#### Availability

- Eventual consistency
- Reconciling conflicts

	Mnesia	CouchDB	Riak	Cassandra
Fragmentation	<ul style="list-style-type: none"><li>• <b>Explicit placement</b></li><li>• Client-server</li><li>• Automatic by using a hash function</li></ul>	<ul style="list-style-type: none"><li>• <b>Explicit placement</b></li><li>• Multi-server</li><li>• <b>Lounge</b> is not part of each CouchDB node</li></ul>	<ul style="list-style-type: none"><li>• Implicit placement</li><li>• Peer to peer</li><li>• Automatic by using consistent hash technique</li></ul>	<ul style="list-style-type: none"><li>• Implicit placement</li><li>• Peer to peer</li><li>• Automatic by using consistent hash technique</li></ul>
Replication	<ul style="list-style-type: none"><li>• <b>Explicit placement</b></li><li>• Client-server</li><li>• Asynchronous ( Dirty operation)</li></ul>	<ul style="list-style-type: none"><li>• <b>Explicit placement</b></li><li>• Multi-server</li><li>• Asynchronous</li></ul>	<ul style="list-style-type: none"><li>• Implicit placement</li><li>• Peer to peer</li><li>• Asynchronous</li></ul>	<ul style="list-style-type: none"><li>• Implicit placement</li><li>• Peer to peer</li><li>• Asynchronous</li></ul>
Partition Tolerant	<ul style="list-style-type: none"><li>• <b>Strong consistency</b></li></ul>	<ul style="list-style-type: none"><li>• Eventual consistency</li><li>• Multi-Version Concurrency Control for reconciliation</li></ul>	<ul style="list-style-type: none"><li>• Eventual consistency</li><li>• Vector clocks for reconciliation</li></ul>	<ul style="list-style-type: none"><li>• Eventual consistency</li><li>• Use timestamp to reconcile</li></ul>
Backend Storage & Query Processing	<ul style="list-style-type: none"><li>• The largest possible Mnesia table is <b>4Gb</b></li></ul>	<ul style="list-style-type: none"><li>• No limitation</li><li>• Support Map/Reduce queries</li></ul>	<ul style="list-style-type: none"><li>• <b>Bitcask has memory limitation</b></li><li>• LevelDB has no limitation</li><li>• Support Map/Reduce queries</li></ul>	<ul style="list-style-type: none"><li>• No limitation</li><li>• Support Map/Reduce queries</li></ul>

## Riak Scalability

We investigate the scalability of Riak version 1.1.1 using the Basho Bench benchmarking tool on the Kalkyl cluster at Uppsala University.

We measure throughput rises vs. the number of Riak nodes

- Every experiment is repeated 3 times
- The scalability diagram depicts the mean values
- The green line represents variation from the mean

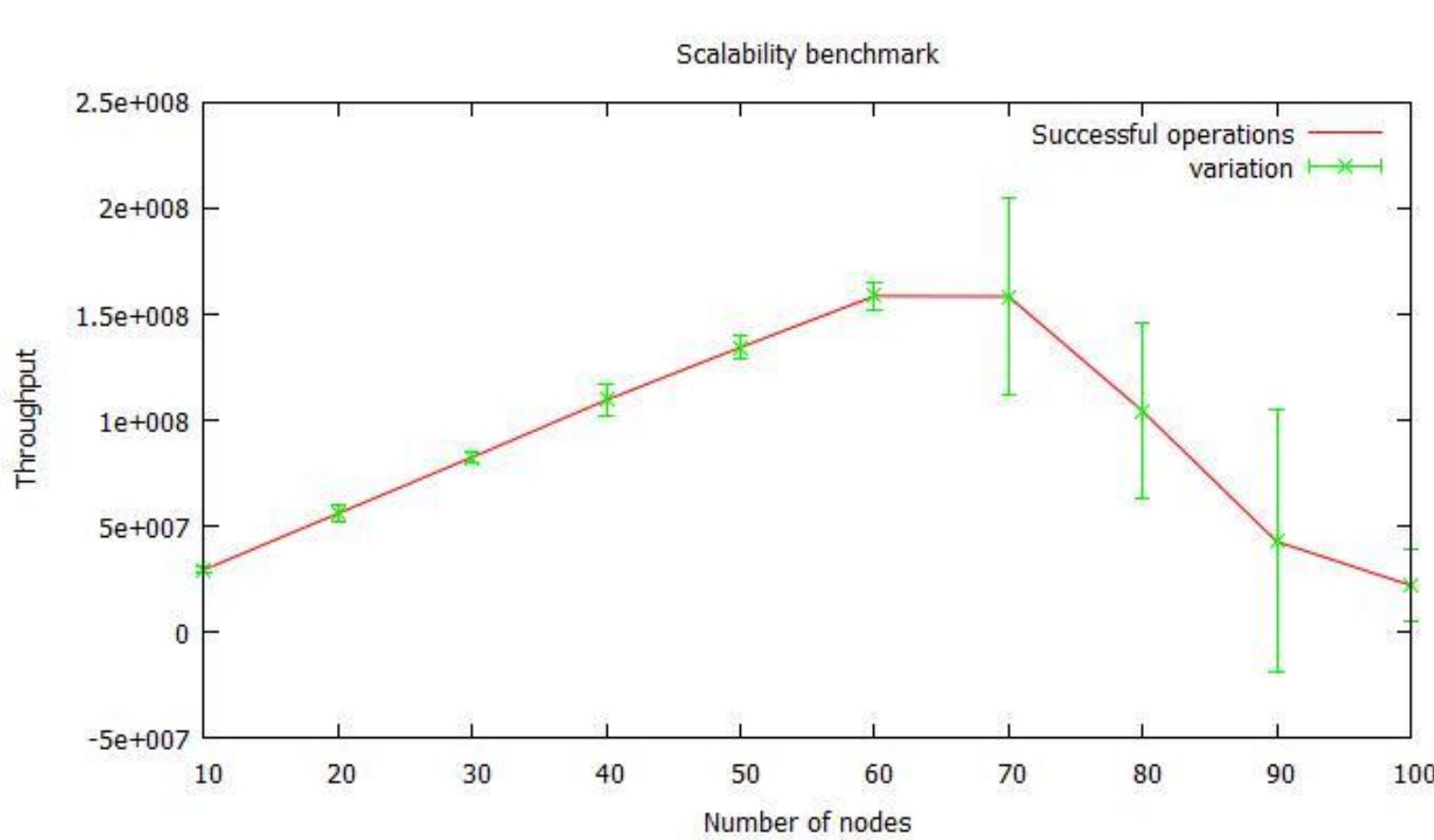
Riak scales linearly up to 60 nodes, but it does not scale beyond 60 nodes.

**Resource usage**

- Maximum RAM usage -- 3%
- Maximum disc usage -- 10%
- Maximum core usage -- 5.5 of 8 cores

**Network profiling**

The number of *retransmitted* packets (200 packets) is negligible in comparison with the total number of successfully transmitted packets 5\*10<sup>8</sup> packets.



### The scalability of Riak software.

- Riak makes no *global.erl* calls
- Of the 15 most time-consuming *gen\_server.erl* operations only one, *rpc:call* grows with cluster size
- Of the 5 Riak RPC calls only *start\_put\_fsm* function from module *riak\_kv\_put\_fsm\_sup* grows with cluster size
- *riak\_kv\_get/put\_fsm\_sup* and statistics reporting are supervisor processes bottleneck
- Riak version 1.3 and 1.4 employ the library *sidejob* to tackle the problem
- *sidejob* library is available at: <https://github.com/basho/sidejob>

## Conclusion & Future Work

We have analysed Erlang DBMSs against the requirements for scalable persistent storage

- *Dynamo-style* NoSQL DBMS have a potential to provide scalable persistent storage for large distributed architectures
- We have shown that Riak 1.1.1 already provides scalable persistent storage on up to 60 nodes
- We further show that resources like disc and network do not limit scalability, and identify two bottlenecks for improvement
- The Riak single-process bottleneck issues are addressed in versions 1.3 and 1.4.

We are investigating the scalability limitations of Distributed Erlang, and developing techniques to further improve the scalability of persistent storage engines implemented in distributed Erlang

## References

- Ghaffari, N. Chechina, P. Trinder, and J. Meredith. Scalable Persistent Storage for Erlang: Theory and Practice. In Proceedings of the Twelfth ACM SIGPLAN Workshop on Erlang, pages 73-74, September 2013.
- N. Chechina, P. Trinder, A. Ghaffari, R. Green, K. Lundin, and R. Virding. The Design of Scalable Distributed Erlang. In Draft Proceedings of the Symposium on Implementation and Application of Functional Languages 2012, July 2012.