

# **Elasticsearch Service**

## **ES Kernel Enhancement**

### **Product Documentation**



## Copyright Notice

©2013-2019 Tencent Cloud. All rights reserved.

Copyright in this document is exclusively owned by Tencent Cloud. You must not reproduce, modify, copy or distribute in any way, in whole or in part, the contents of this document without Tencent Cloud's the prior written consent.

## Trademark Notice



All trademarks associated with Tencent Cloud and its services are owned by Tencent Cloud Computing (Beijing) Company Limited and its affiliated companies. Trademarks of third parties referred to in this document are owned by their respective proprietors.

## Service Statement

This document is intended to provide users with general information about Tencent Cloud's products and services only and does not form part of Tencent Cloud's terms and conditions. Tencent Cloud's products or services are subject to change. Specific products and services and the standards applicable to them are exclusively provided for in Tencent Cloud's applicable terms and conditions.

# ES Kernel Enhancement

Last updated : 2021-03-12 14:58:58

While maintaining full compatibility with the open-source Elasticsearch kernel, Tencent Cloud ES team has been continuously and thoroughly exploring and optimizing its kernel based on its rich experience in large-scale applications in various use cases, so as to enhance cluster performance, improve stability, and reduce costs. In addition, it has been maintaining close communications with the open-source community. This document describes kernel optimizations.

The following table summarizes the key kernel optimizations that the ES team has made by **July 2020** since the start of its kernel research:

Optimization Dimension	Optimization Category	Optimization Policy	Supported Versions
Performance	Write performance	The translog lock mechanism is optimized, increasing the overall write performance by 20%. Write deduplication and segment file cropping are optimized, increasing the performance of writes with primary keys by over 50%.	7.5.1
	Query performance	<ul style="list-style-type: none"><li>The aggregation performance is optimized, making query pruning more efficient and improving the composite aggregation performance by 3-7 times in sorting scenarios.</li><li>The query cache is optimized by canceling data caches with high overheads and low hit rates, reducing query glitches from 750 ms to 50 ms in actual use cases.</li><li>The merge policies are optimized by developing proprietary merge policies based on time series and size similarity and auto warm shard merge policy, improving the query performance by over 40% in search scenarios.</li><li>Sequence capture in the query fetch phase is optimized, increasing the cache hit rate and improving the performance by over 10% in scenarios where the result set is large.</li></ul>	6.4.3, 6.8.2, 7.5.1
Stability	Availability	<ul style="list-style-type: none"><li>Traffic can be limited through a smooth line curve at the access layer.</li><li>The coordinator node performs memory bloat estimation after receiving results returned by the</li></ul>	6.4.3, 6.8.2, 7.5.1

		<p>data node to check whether the estimated memory will exceed the limit.</p> <ul style="list-style-type: none"> <li>Result sets of large aggregated queries are checked in a streaming manner, and requests will be canceled if the used memory reaches the threshold.</li> <li>The proprietary single request circuit breaker can prevent a large query from occupying excessive resources and thus affecting other queries.</li> <li>Node crashes and cluster avalanches caused by high-concurrency writes and large queries are significantly reduced, and the overall availability is increased to 99.99%.</li> </ul>	
	Balancing policy	<ul style="list-style-type: none"> <li>Balancing policies based on index and node distribution are introduced, alleviating the serious uneven allocation of shards caused by new nodes added to the cluster.</li> <li>The uneven allocation of shards among multiple disks (multiple data directories) is alleviated.</li> <li>The balance of shards of newly created indices in cluster scale-out scenarios and multiple-disk scenarios is improved, reducing OPS costs.</li> </ul>	5.6.4, 6.4.3, 6.8.2, 7.5.1
	Rolling restart speed	<ul style="list-style-type: none"> <li>The logic of reusing local data for shards in case of node restart is optimized.</li> <li>The restoration of shard copies within a scheduled delay time period can be precisely controlled. The time to restart one single node in a large cluster is reduced from over 10 minutes to 1 minute.</li> </ul>	6.4.3, 6.8.2, 7.5.1
	Online master switch	<p>The proprietary online master switch feature allows you to switch the master online in seconds by specifying the preferred master through APIs. Typical use cases include:</p> <ul style="list-style-type: none"> <li>You can switch online from the current heavily loaded master to a node with a higher specification and a lower load during manual OPS.</li> <li>During rolling restart, you can restart the master node last and quickly switch the master role to another node before the restart, which</li> </ul>	6.4.3, 6.8.2, 7.5.1

		helps reduce the service interruption from minutes to seconds.	
Costs	Memory	<ul style="list-style-type: none"><li>• The proprietary off-heap cache helps achieve FST off-heap optimization.</li><li>• The off-heap cache ensures that the FST reclaim policy is controllable.</li><li>• The precise eviction policy improves the cache hit rate.</li><li>• Zero-copy and multi-level caches guarantee high access performance.</li><li>• The heap memory overheads are significantly reduced, the GC time is decreased by over 10%, and the disk capacity of a single node can reach 50 TB, with read/write performance generally unaffected.</li></ul>	6.8.2, 7.5.1
	Storage	<ul style="list-style-type: none"><li>• The proprietary ID field-based row storage cropping algorithm reduces storage overheads by over 20% in time series scenarios.</li><li>• A new compression algorithm is introduced, increasing the compression ratio by 30%-50% and the compression performance by 30%.</li></ul>	5.6.4, 6.4.3, 6.8.2, 7.5.1