# Adding file-level CRC32C support to Hadoop Distributed FileSystem

Authors: dhuo@google.com,

# Objective

Define a new "FileChecksum" type in Hadoop Distributed FileSystem (HDFS) which is the raw CRC32C of the entire file contents, to enable checksum comparison between HDFS instances with very different underlying block configurations, between replicated and striped HDFS files, and even comparison with non-HDFS implementations of Hadoop's FileSystem interface.

# Background

HDFS uses CRC32Cs to maintain data integrity in several different contexts:

- At rest, DataNodes continuously verify data against stored CRCs to detect and repair bit-rot
- In transit, the DataNodes send known CRCs alongside the corresponding bulk data, and HDFS client libraries cooperatively compute per-"chunk" CRCs to compare against the CRCs received from the DataNodes
- For HDFS administrative purposes, "block"-level checksums are used for low-level manual integrity checks of individual "block" files on DataNodes
- For arbitrary application-layer use cases, the FileSystem interface defines getFileChecksum, and the HDFS implementation uses its stored fine-grained CRCs to define such a "file-level" checksum

For most day-to-day uses, the CRCs are used transparently with respect to the application layer, and only use per-"chunk" CRC32Cs which are already precomputed and stored in "metadata" files alongside block data. The "chunk" size is defined by `dfs.bytes-per-checksum` and has a default value of 512 bytes. All API-exposed checksums currently take the form of an MD5 of a concatenation of chunk CRC32Cs, either at the "block" level through the low-level `DataTransferProtocol`, or at the "file" level through the top-level `FileSystem` interface. The latter is defined as the MD5 of the concatenation of all the block checksums, each of which is an MD5 of a concatenation of chunk CRCs, and is therefore referred to as an "`MD5MD5CRC32FileChecksum`". This is effectively an on-demand three-layer [Merckle tree](#).

This definition of the "file-level" checksum is sensitive to implementation and data-layout details of HDFS, namely the "chunk" size (default 512 bytes) and the "block" size (default 128MB). As such, it is not usable in any of the following situations:

- Two different copies of the same files in HDFS but with different per-file block sizes configured
- Two different instances of HDFS with different block or chunk sizes configured

- Copying across non-HDFS "Hadoop compatible filesystem implementations" such as Google Cloud Storage, AWS S3, Azure Blob Storage, etc.

# Overview

Since CRC32C [can be efficiently composed](#), it is possible to define new "composite block CRCs" and "composite file CRCs" as the mathematically composed CRC across the stored chunk CRCs rather than using MD5 of the component CRCs to calculate a single CRC that is representative of the entire block or file and independent of the lower-level granularity of "chunk" CRCs.

Given the sensitivity of data integrity and the vast volume of data stored in existing HDFS deployments, it is desirable to minimize changes to existing behaviors, even if hidden from the application-layer interface. This means the added functionality should avoid changing the way `BlockScanner` maintains data integrity at rest or the way `BlockReaderRemote` verifies chunk CRCs in transit.

# Design Details

## Modifying BlockChecksum in DataTransferProtocol

The `DataTransferProtocol` defines the low-level protocol-buffer-based interface over TCP for HDFS clients to access DataNode data or metadata. Checksum information is available in certain mutation requests and as a pre-computed "MD5 of CRC" block checksum, in addition to providing the complete chunk-granularity stream of chunk CRCs in streaming reads. While the client could reconstruct comparable composite CRCs from the read stream, it is necessary to provide a means of computing composite CRCs without incurring the cost of ingesting the complete actual data contents from disk.

As such, there doesn't currently exist an efficient accessor for CRC metadata in the DataTransferProtocol, so any complete implementation requires modification to the DataNode service to modify the DataTransferProtocol. To reuse the framework for dealing with any partial chunks or range requests, this feature will modify both the existing `BlockChecksum` and `BlockGroupChecksum` methods. In theory, the remote caller of this method needs only the single composite CRC and CRC type in the response. Notably, in contrast to MD5-based block checksums, the response does *not* need to expose internal details about bytes-per-CRC or crcs-per-block. However, since the FileChecksum doubles up to be used for file-attribute propagation in certain cases, the bytes-per-CRC is still needed in the response.

The behavior of BlockChecksum will be determined by an additional option in the OpBlockChecksumProto to indicate whether MD5CRC or COMPOSITE_CRC is desired at the

block level. The option itself will not distinguish between COMPOSITE_CRC32 vs COMPOSITE_CRC32C, since the option is a runtime property, while CRC32 vs CRC32C is a sticky property of the underlying data.

In contrast to adding a completely separate DataTransferProtocol Op composite-crc paths:

Pros
- Avoids the proliferation of protocol Op codes
- Easier reuse of ReplicatedBlockChecksumComputer and BlockGroupNonStripedChecksumComputer

Cons
- Failure modes for mismatched client/server versions less clean
- More changes to existing codepaths leading to increased risk of bugs impacting the old behavior
- Requires changing existing method signatures of nominally public interfaces

Client-side support will additionally require shared logic for hierarchically merging multiple block CRCs into a file-level CRC.

For better support of custom client-side definitions of range CRCs, the idea of adding raw accessors for possible [client-side aggregation algorithms for new striped erasure-coded files](#) has been discussed before. As it turns out, we can simply generalize the logic for being able to specify a "stripe length" for block checksums to use; when this length aligns with cell size for striped encodings, block-group checksums can be reassembled efficiently by the parent datanode of block-group checksum requests, and a full-block CRC can be considered equivalent to an unlimited "stripe length". This accessor then leaves open the possibility of new client-side protocols to fetch smaller CRC stripes, even if not being used for a striped erasure-code format. This is discussed in more depth in the "Striped Erasure-Coded formats" section below.

## Legacy "gzip" CRC32 support

Prior to [HADOOP-7443](#), HDFS used the same CRC polynomial used in [java.util.zip.CRC32](#), with little-endian bit-representation 0xEDB88320. In file-checksum contexts, this is internally/colloquially referred to as the "Gzip" variant of the MD5-composed checksum, i.e. the [MD5MD5CRC32GzipFileChecksum](#)), not to be confused with checksums over gzipped contents, but simply named as such due to the legacy polynomial being the same one used by gzip.

It is desirable for this feature to support both the Castagnoli variant as well as the "Gzip" variant, for two reasons:

- A significant advantage of this proposed design is in-place backwards compatibility, so supporting the legacy format which may still be in use by some older HDFS deployments is in-line with this goal
- It is a good driving use case to ensure the implementation is done in a general manner to better accommodate new CRC polynomials in the future, especially when likely moving to 64-bit CRCs if/when CPU-native support for a 64-bit standard is introduced

Nonetheless, it is expected that CRC32C will remain standard/preferred for the foreseeable future, and the primary variant to be used in compatibility across heterogeneous storage systems for the same reasons it was introduced in the first place (superior error-detection semantics and CPU intrinsic support since [SSE 4.2](#)).

In practice, this means the file-level composition strategy will be kept distinct and orthogonal to choice of underlying component CRCs; rather than introducing COMPOSITE-CRC32C as a "peer" of "MD5MD5CRC32", we should think of "COMPOSITE-CRC" and "MD5MD5CRC" as different composition strategies applicable to different underlying CRC options. Additionally, protocol definitions will continue to return variable-length composite checksum definitions (i.e. "bytes" in the protocol buffer definition) instead of uint32.

# CRC32C strength and possible future CRC64 support

## Error detection vs tamper resistance

In assessing the role of file-level checksums in various aspects of data integrity, it is important to note that the same properties which make CRCs well-suited for distributed storage systems (composability, reversibility) mean it is fundamentally not tamper-resistant. At the same time, MD5 is also considered insecure in the context of adversarial tamper resistance. As such, we can recognize that "secure" data integrity is already a problem which must be solved out-of-band from existing internal error-detection/correction mechanisms.

This observation helps focus the driving requirements for the category of file-level checksums discussed here. Importantly, the theoretical existence of collisions and/or the triviality of being able to construct intentional collisions is *not* a driving concern, and instead the protection strength can be assessed in the context of actual random-error sources in the use-case at hand.

## Hierarchical error detection for data transfers

In general, 32-bit CRCs are expected to be valuable for generalized detection of transfer-time errors in files being migrated between separate storage instances, since each hierarchical protocol layer provides certain guarantees on the nature of errors that may go undetected. In particular, since HDFS continues to verify per-chunk CRCs at the transfer layer, the minimum number of bit-errors required to generate a failed error detection in an arbitrary-length payload is

lower-bounded by the [CRC's minimum Hamming distance](#) for an undetected error in a chunk-sized payload. For a Hamming distance of N,

1. If all N bit-flips occur within a single chunk, this is trivially true by being detected by the chunk verifier
2. If the N bit-flips occur spread across several chunks to result in a file-level CRC collision, then some number of chunks would have held K-bit errors ranging from 1 <= k <= N - 1, and by definition of CRC Hamming distance, all bit-errors <= N bits will be detected by the chunk verifier

The second case is where the hierarchical inclusion of per-chunk verification distinguishes the possible error types from a pure-file-checksum based integrity check. In terms of concrete numbers, CRC32C has a minimum Hamming [distance of 6 for the default chunk size of 512 bytes](#), whereas, for example, a standalone CRC32C of a 2GB payload (greater than the period of the polynomial) would be vulnerable to 2-bit errors.

For certain dense bit-errors that go undetected within a single chunk, the approach of using a composite CRC is no worse than the existing approach of using MD5-of-CRCs, since in both cases the strength of the aggregated checksum is no better than that of individual chunks; importantly, MD5-of-CRCs is not equivalent to MD5 of the raw byte contents.

Ultimately, the types of errors most likely to benefit from the use of file-level composite CRCs are those caused by software bugs that may introduce errors independently of transfer-layer chunk checksums, such as rebroadcast/duplication/out-of-order buffering bugs or other software-layer memory corruption introduced after chunk-level checks are performed.

## Use cases vulnerable to collisions

Certain legitimate use cases may call for longer CRCs, such as data de-duplication across a large number of files, since in such a case the birthday paradox applies and we'd expect only something on the order of $2^{16}$ files to reach a ~50% chance of a collision in a 32-bit space.

In anticipation of future extension to longer CRCs, protocol definitions will be length-agnostic.

# Supporting file prefix-range checksums

The same approach for prefix-range checksums currently used in MD5MD5CRC combine mode will apply to new the COMPOSITE-CRC combine mode, where the final partial-chunk will need to be explicitly fetched from disk to obtain a new CRC32C of the partial chunk on-demand.

As long as the length of the partial chunk is accounted for when composing CRCs, the composed CRCs will uniquely support further composition with a file suffix-range while preserving comparison with whole-file checksums. In contrast, MD5MD5CRC mode fundamentally does not support extending a prefix-range checksum with a suffix while retaining equality with the whole-file checksum, because there is no way to "back out" the partial-chunk checksum digested into the block MD5.

## Striped Erasure-Coded formats

HDFS-7285 introduces a new striped, erasure-coded file format to HDFS in Hadoop 3+, adding a new hierarchical layer of data granularity called a "cell", defaulting to 64kB, arranged into a "block groups" which are logically analogous to normal "blocks" in a non-striped file, are now striped across multiple datanodes.

For purposes of file-level checksum support, HDFS-8430 implements a three-layer MD5-of-CRCs approach, whereby checksums of individual contiguous sections (striped blocks) are calculated through the same non-striped "blockChecksum" method but are aggregated at the "block group" layer before then being combined at the file-level. Each "blockChecksum" is thus an MD5-of-CRC, each "blockGroupChecksum" is MD5-of-MD5-of-CRC, and the FileChecksum is an "MD5-of-MD5-of-MD5-of-CRC", even though it is exposed as a comparable MD5MD5CRC checksum type, which implies compatibility with regular non-striped file checksums.

This incompatibility between FileChecksums of striped files vs replicated files was indeed identified as one of the key shortcomings of the straightforward MD5 approach, and was discussed in-depth in HDFS-8430. Maintaining compatibility with the existing replicated-file MD5BD5CRC was deemed infeasible due to the approach requiring fetching all *chunk*-level CRCs from all sibling DataNodes sharing cells of a single block group, and having a single block-group mediator combine chunk CRCs in order to make blockGroupChecksum analogous to replicated-file blockChecksums.

This document's new COMPOSITE-CRC32 approach allows both backwards-compatibility and comparability between replicated and striped files independently of cell layout. This will require adding API support for CRC composition to the blockGroupChecksum method of the DataXceiver, along with more extensive changes to the way underlying blockChecksums are aggregated at the BlockGroup level.

Specifically, even though the same blockChecksum method will be called as children of the BlockGroup for convenience, the blockChecksum protocol for contributing to a striped file will be somewhat different from the protocol for computing a complete block-level CRC. The BlockGroup parent must indicate a desired CRC "stripeLength", and in such cases the blockChecksum will return a list of N == (requestedDataLength - 1) /

`stripeLength + 1` CRCs. The BlockGroupChecksum parent must then reconstitute the striped CRCs in the correct logical order.

This approach does increase network traffic, but since cell sizes are much larger than chunk sizes, and a datanode parent is responsible for aggregating at the BlockGroup level, the amount of network traffic is negligible. For example, at a default 128MB blocksize and 1MB stripe size, the blockChecksum response would contain 128 individual CRCs, amounting to a response size of only 512 bytes.

Much of the logic already implemented for hierarchical MD5 composition can be reused to handle cases of missing data blocks requiring cell reconstruction against parity blocks, but some refactoring is still needed to abstract out the CRC composition logic from the parity-block repair logic in helper classes like [BlockChecksumHelper](#) to eliminate hard-coded assumptions about accumulating an MD5 of underlying checksums.

# Other Implementation Details

## DataChecksum vs FileChecksum options

The [DataChecksum](#) class encapsulates "internal" checksum options pertaining to transfer-level checksums, while the [FileChecksum](#) constitutes the public interface. At the moment, the nature of the FileChecksum is fully a function of the underlying DataChecksum options, being implicitly defined by an effective combination of `dfs.checksum.type`, `dfs.bytes-per-checksum`, and `dfs.blocksize`, all defined in [HdfsClientConfigKeys](#).

For this new feature, it is desirable to allow a client-side configuration definition to choose the CRC combine strategy at runtime, so we introduce a new key `dfs.checksum.combine.mode`, configured orthogonally to the transfer-level configuration options.

## ChecksumOpt preservation of low-level FileAttributes

Since the value of `dfs.bytes-per-checksum` is part of the definition of the FileChecksum algorithm name when using MD5MD5CRC combine mode, the embedded [FileChecksum.getChecksumOpt](#) doubles as a mechanism to obtain per-file low-level checksum configs for use with cases like [FileAttribute preservation in DistCp](#).

Since the new combine mode makes FileChecksum agnostic to underlying chunk or block representation, `dfs.bytes-per-checksum` has no reason to be propagated into

FileChecksum, and in particular will not be part of the "algorithm name" to ensure comparable checksums between different HDFS instances with different underlying chunk configurations.

However, since attribute-preservation is in theory unrelated to the comparability of FileChecksum instances, the interface seems to dictate propagating the chunk configuration into the ChecksumOpt despite being unnecessary to the checksum computation itself. For copies from HDFS to HDFS, this will thus behave as expected even when "COMPOSITE-CRC" is used as the combine mode.

If copying from a storage system which is unable to expose underlying chunk configuration, the ChecksumOpt may either be set to an uninitialized value of `dfs.bytes-per-checksum`, or could inherit runtime HDFS settings.

## Abstracting out FileChecksum implementation from DFSClient

Though the [FileChecksum](#) interface declared as the return value for FileSystem and AbstractFileSystem's getFileChecksum method is sufficiently generic to accommodate this new COMPOSITE-CRC format, the lower-level [DFSClient.getFileChecksum](#) method explicitly returns an "MD5MD5CRC32FileChecksum". In general, the DFSClient is only supposed to be used as an internal implementation detail of the DistributedFileSystem and in other HDFS-internal contexts, but it is declared as a public class and DFSClient.getFileChecksum is a public method (mitigated by the class-level [annotation "@InterfaceAudience.Private"](#)).

In order to return the COMPOSITE-CRC as a different subclass of FileChecksum, either DFSClient must change its public method signature (and break any users assuming MD5MD5CRC32FileChecksum to be the concrete class returned), or the configuration option must be applied one level higher, in the [DistributedFileSystem](#). In the interest of backwards-compatability, the preferred approach will be to apply the configuration in the DistributedFileSystem, at the cost of leading to slightly more code duplication in the DFSClient.

# Performance

## Amortization across fixed-size chunks

While the critical-path computation of CRCs from raw data benefits from SSE intrinsics, there is no such native "compose-crc" support. Though in many cases the order of magnitude of CRC compose operations performed is small enough to make efficiency considerations negligible (for example, a single compose operation on "append", or concatenating on the order of 10s to 100s of block CRCs into a file CRCs), in this case the reuse of chunk CRCs to compute the aggregate CRC introduces non-negligible efficiency requirements. At a default `dfs.bytes-per-checksum` of 512 and a default `dfs.blocksize` of 128MB, this translates

to 250,000 CRC-composition operations per block, and larger block sizes are commonly used in large deployments, where 512MB blocks mean 1,000,000 operations.

As indicated in [this whitepaper](#), the composition of two CRCs $CRC(M_1, 0)$ and $CRC(M_2, 0)$ can be modeled as a special case of computing a "change of initialization polynomial" given a source with initialization polynomial of **0** and target polynomial of $CRC(M_1)$. Applying the formula, we see:

$CRC(M_2, CRC(M_1)) = CRC(M_2, 0) + ((CRC(M_1, 0) - 0) \, x^{|M2|} \, mod \, P$

The bulk of the computation thus lies in calculating the monomial $x^{|M2|} \, mod \, P$ and multiplying it by the CRC of $M_1$. As usual, by expressing the length L in terms of its binary representation,

$$L \; = \; |M_2| \; = \; \Sigma b_i 2^i$$
$$x^L = x^{\Sigma b_i 2^i}$$
$$= \Pi x^{b_i 2^i}$$

and each power-of-two monomial $x^{2^i}$ can be efficiently computed independently by repeated squaring. We thus see that the basic runtime of a single composition operation is logarithmic in the length L of $M_2$. While powers-of-two monomials can be precomputed, we must always support arbitrary lengths of $M_2$ when it is an incomplete chunk or incomplete block. Furthermore, there is no fundamental constraint that chunk sizes are exact powers of 2.

Naively, if we have $N$ chunks each of size $L$, computing the total composite CRC is **O(log(L) * N)**. Using sample numbers of 1,000,000 operations, letting log(L) ~= 64, each multiplication taking 32 operations, and 1ns per operation, we see this approaching 1M * 2048ns ~= 2 CPU-seconds for a single composite CRC.

In contrast, since the chunk layout is largely homogenous, we can have the block-checksum loop precompute the monomial associated with the given block's chunk size (even if this differs between different files, within a block the chunk size is always constant) and thus reduce every chunk composition to a single 32-bit polynomial multiplication (and one XOR) except for the last partial chunk, to achieve overall **O(log(L) + N)** time. Using the same example numbers, this reduces the CPU-cost to ~30ms from 2s.

In practice, chunk sizes aligned with powers of two should have the same time complexity as precomputing the monomial. However, looping over unset bits in recomputing the monomial still introduces overhead, and more importantly, since chunk sizes are configurable per-file, it is desirable to enforce consistent CPU performance of file checksums, rather than allowing an ill-conceived or intentionally-malicious chunk size to suddenly cause DataNodes to spend i.e. 10-20x the typical CPU cycles on checksumming.

Benchmarks indicate performance characteristics in-line with the theoretical values, taking ~0.5s to compose 1,000,000 CRCs with data length 511 using a shared monomial vs ~4.5s to do the same recomputing the monomial from a powers-of-two lookup table for each composition (and more than 10s to do the same without a powers-of-two lookup table, instead performing repeated squaring on-demand).

### Skipping chunk CRCs during parity-block reconstruction

Since existing MD5-of-CRC checksums are sensitive to chunk size, the [StripedBlockChecksumReconstructor](#) must still compute individual chunk CRCs independently before combining them into an MD5; in this case, chunk CRCs aren't used other than in the aggregate checksum, so it would be more efficient for a COMPOSITE-CRC to simply compute the contiguous running CRC across the reconstructed block to avoid a second CRC-composition phase. However, the code is not well structured to support this divergent branch of logic, so it would either require refactoring or a custom block-reconstruction implementation. Since the reconstruction of data blocks involves reading orders of magnitude more data from disk than plain CRC metadata, the additional inefficiency of a CRC composition phase is negligible anyways, and it is likely not worth the maintenance overhead to perform this optimization.

# Augmenting FileSystem interfaces with data integrity checks

An immediate benefit of implementing COMPOSITE-CRC is that several existing FileSystem interfaces can be augmented to apply low-level data integrity checks transparently and efficiently where this wasn't previously possible. Notably:

- **concat** - The namenode could share the CRC-combine logic and compare against the datanode-aggregated values to ensure no out-of-order issues occured in assigning new block index mappings
- **append** - The client can pre-fetch the existing file-level CRC and incrementally extend the CRC with newly appended bytes without requiring knowledge of underlying chunk or block layout; on completion, a full-file checksum can be requested and compared against the checksum computed from the continued stream
    - In contrast, this is not possible with the MD5MD5CRC mode even with client-side knowledge of chunk/block layout, because a partial chunk and/or partial block would be factored into the original file which is not present at all in the full-file checksum post-append

# Auxiliary tooling

## BlockReader accumulator without DataTransferProtocol support

While the existing DataTransferProtocol doesn't expose CRCs directly, it does embed them into block read streams, so at least for live data migrations it is possible to opportunistically preserve pure client-computed aggregate CRCs. Modifications would need to be made to the BlockReaderRemote and BlockReaderLocal to accumulate per-chunk CRCs and expose an accessor to fetch the composite value upon completion of each block. The DFSInputStream would then accumulate per-block CRCs directly from the BlockReaders and would itself compose a file-level checksum.

This type of modification could be used in cases where it is infeasible to upgrade DataNodes on a large HDFS cluster, and a pure client change is needed. For single-stream use cases, a custom HDFS client bundled locally should work, but running a distributed job with custom HDFS clients may run into classpath collisions with existing HDFS client classes under i.e. /usr/lib/hadoop-hdfs/hadoop-hdfs-client*.jar. In such a case, it would be necessary to build a clean end-to-end dependency stack (including i.e. DistCp itself) using Maven shade plugin to relocate the entire HDFS client package.

## Block metadata file validator

Given mappings to existing block metadata files and not wanting to update DataNode daemons in-place, it could still be possible to overlay metadata file readers as separate ad-hoc daemons running only for the duration of a data-verification effort. This tooling would involve implementing a fully standalone and lightweight client/server pair, with servers responsible for reading block metadata and returning composed block CRCs, and the client using the real namenode to fetch block locations before connecting to the dedicated block-CRC daemons instead of to the datanode ports indicated by the namenode.

# Version History

- v3 - Updated DataTransferProtocol and Striped Erasure-Coded formats sections to correct some inaccuracies and expand on the way striped composite CRCs are reconstituted within BlockGroupChecksums

- v2 - Updated section about modifying DataTransferProtocol to modify existing BlockChecksum and BlockGroupChecksum methods with new BlockChecksumType parameters instead of adding separate composite-crc specific methods
- v1 - Initial draft