

# Checksum

In a distributed system, while moving data between components, it is possible that the data fetched from a node may arrive corrupted.

This corruption can occur because of faults in a storage device, network, software etc.

How can a distributed system ensure data integrity so that the client receives an error instead of corrupt data?



Calculate a checksum and store it with data.

To calculate a checksum, a cryptographic hash function like MD5, SHA-1, SHA-256 or SHA-512 is used.

The hash function takes the input data and produces a string (containing letters and numbers) of fixed length; this string is called the checksum.

When a system is storing some data, it computes a checksum of the data, and stores the checksum with the data.

When a client retrieves the data, it verifies that the data it received from the server matches the checksum stored.

If the checksum does not match, then the client can opt to retrieve that data from another replica.

Of course, it may also be the case that the checksum itself gets corrupted.

However, given the avalanche effect property of the cryptographic hash functions, it is highly unlikely that both data and checksum will be corrupted in such a way that hashing the data to a checksum will match the corrupted checksum.

## Examples

**HDFS** and **Chubby** store the checksum of each file with the data.