# Ceph Code Optimization

-----

## **Compile-time Optimizations**

Most of ceph is built with "-O2". Selected parts either contain -m<extension> or contain assembler code with instruction extensions. All of this is enabled at compile time based strictly on what options the build compiler and assembler support. The bulk of this is limited to the erasure code logic. There is limited support for crc32 acceleration and sha256 acceleration in several other parts.

\*\*Special notes: zstd (14 files) is built with -O0. Some parts of the intel is that should support avx512 are built as empty files, because we don't supply the right option to enable avx512 support.

#### 1. General

- a. radosgw Like most of ceph, built almost entirely w/ -O2.
- b. ceph-osd Rocksdb, used by bluestore, has its own build system, and looks for crc32.
  The erasure code is built as separate plugins, and the "ec-isa" plugin can take advantage of Intel storage acceleration.

#### 2. Modules of Special Interest

- a. *openssl* We don't build this, so no compile flags here.
- b. *jerasure* Most of the erasure code logic can take advantage of up to sse4.1 logic. The erasure code plug-in "libec\_isa.so" can use up to AVX2 instructions.
- c. *rocksdb* This code uses crc32, and includes support to look for and use SSE4.2 for this.
- d. *compressor code (zstd)* Used by s3 interface to compress objects in buckets. Bad optimization, but contains possible use of crc32 acceleration.

#### **Object Code As Packaged (Ceph and platform)**

#### 3. Objdump Inspection of Selected Objects

- a. radosgw The bulk of this is built with just -O2. It links against libceph\_common, which looks for and can use crc32 instructions.
- b. ceph-osd The bluestore code can use BMI instructions, part of haswell.
- c. openssl This is supplied by Ubuntu, and loaded at runtime by radosgw. On Ubuntu, it is also used by libcurl. It should be able to take advantage of most CPU instruction optimizations, including use of sha1 and avx2.
- d. jerasure Most parts can use up to sse4.1. Selected parts (esp. libec\_isa.so) can use additional intel features, up to avx2.

#### 4. Implementation of Key Algorithms

- a. CRC (Messenger, bluestore? checksums) Used by bluestore and messenger. The bluestore call can definitely use extensions. The messenger code should be able to as well.
- b. MD5 (radosgw ETAG) No special optimizations.
- c. SHA-256 (radosgw AWS4\_HMAC\_SHA256) There are various optimizations versions available up to AVX2. Ceph code should be able to take advantage of this.

d. AES (civetweb/beast/libcurl SSL/TLS) - For incoming connections, civetweb & beast will both use openssl, which contains aes acceleration. For outgoing connections, libcurl on ubuntu will also use openssl, so ditto. There is some internal support for aes acceleration in ceph itself, such as for cephx.

### Test environments:

- a. QE uses "magna" machines which are ivy bridge.
- b. a<sup>^</sup>2 has mostly haswell machines for development.
- c. Scalelab machine has skylake machines, for very temporary testing at scale.