# Online Compression with Intel® QAT in DAOS

Weigang Li, Intel

Johann Lombardi, Intel

November, 2020

**intel.**

# DAOS compression



**Compute Nodes**

Application

DAOS library

Compressor

■ Compress on Server

- Optimize storage bandwidth/usage and not network bandwidth

- Able to use hardware acceleration

**DAOS Nodes**

DAOS service

Compressor

Compressed Data

Storage

Compressed Data

■ Compress on Client

- Reduce amount of data transferred over the network

- Consume compute CPU cycles

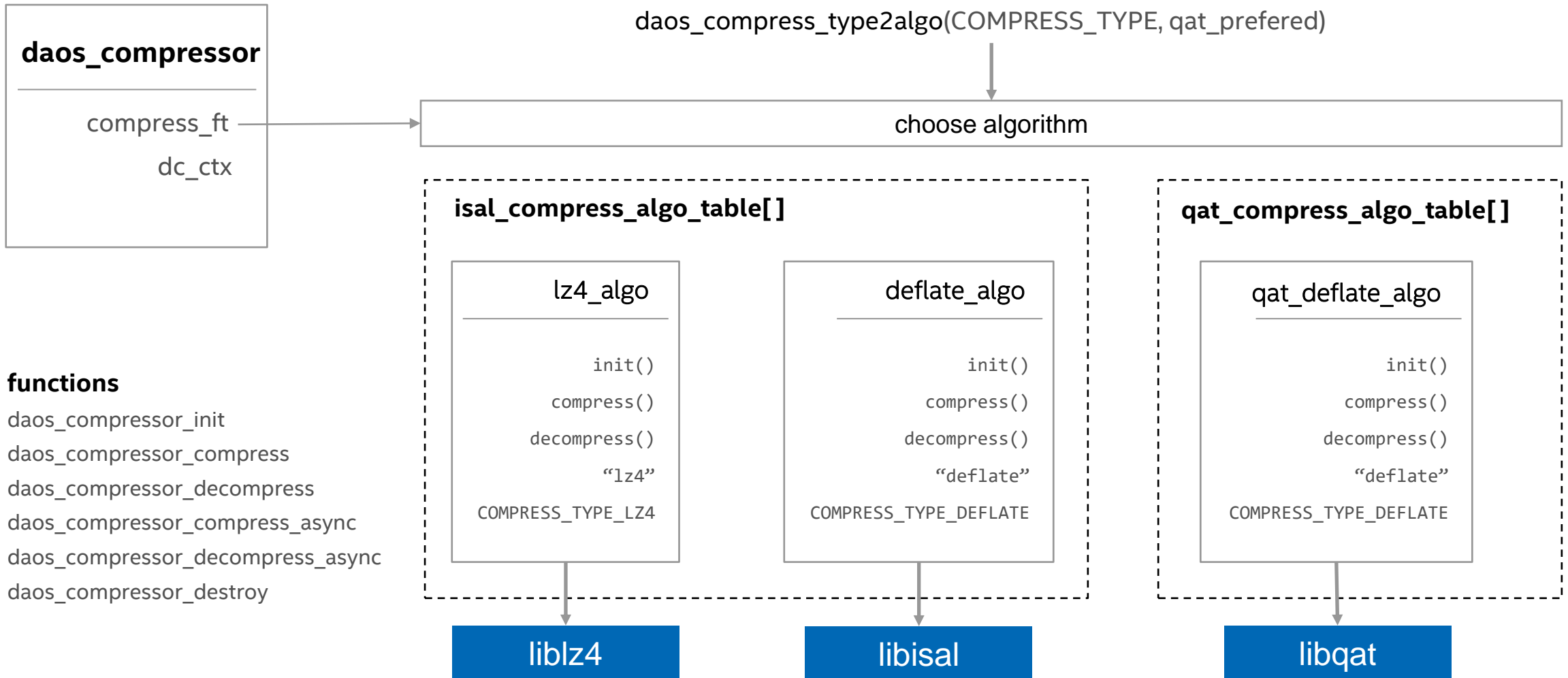*Note: Only data will be compressed for now and not metadata*

# Current Status

- All the infrastructure to enable compression and encryption has been landed.

- It is set at container creation time via properties DAOS_PROP_CO_COMPRESS and DAOS_PROP_CO_ENCRYPT.

*# daos cont get-prop --path /mnt/lustre/mycontainer --svc $SVC*
*Container properties for d28d9df6-5ee7-4030-b20f-*
*bf06c13e2226 :*
*label:              container label not set*
*layout type:         POSIX (1)*
*layout version:      1*
*checksum type:       off*
*checksum chunk-size:  32768*
*cksum verif. on server: off*
*deduplication:       off*
*dedup threshold:     4096*
*redundancy factor:   rf1*
*redundancy level:    rack*
*max snapshots:       0*
<span style="color:red">*compression type:     off*</span>
<span style="color:red">*encryption type:      off*</span>
*owner:              root@*
*owner-group:         root@*

- DAOS-5605 common: add basic infrastructure for QAT #3402 (Merged)

- DAOS-5605 build: add lz4 dependency #3403 (Merged)

- DAOS-5719 compression: add framework and lz4/deflate support #3561 (Merged)

- DAOS-5719 compression: add qat support for deflate #3621 (Merged)
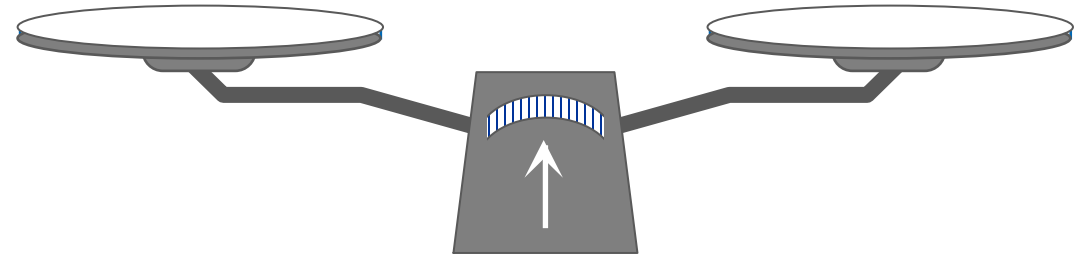
- QAT async support and compress_timing test (WIP)

# DAOS Compressor: Abstraction layer in DAOS

**daos_compressor**
- compress_ft
- dc_ctx

daos_compress_type2algo(COMPRESS_TYPE, qat_prefered)

choose algorithm

**isal_compress_algo_table[ ]**

**lz4_algo**
- init()
- compress()
- decompress()
- "lz4"
- COMPRESS_TYPE_LZ4

**deflate_algo**
- init()
- compress()
- decompress()
- "deflate"
- COMPRESS_TYPE_DEFLATE

**qat_compress_algo_table[ ]**

**qat_deflate_algo**
- init()
- compress()
- decompress()
- "deflate"
- COMPRESS_TYPE_DEFLATE

**functions**

daos_compressor_init
daos_compressor_compress
daos_compressor_decompress
daos_compressor_compress_async
daos_compressor_decompress_async
daos_compressor_destroy

liblz4

libisal

libqat

# Compression algorithms

- LZ4
  - Fast, low compression ratio
- Deflate
  - Slow, high compression ratio
  - Implementation: ZLIB, GZIP
- zStandard
  - New algorithm
  - Fast, high compression ratio
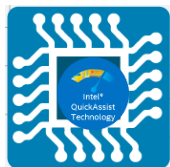  - May be added to DAOS in future

- Performance
  - › Throughput
  - › Compression Ratio

- Cost

Hardware acceleration
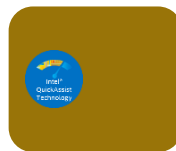
# Intel® QuickAssist Technology (QAT)

- Intel® QuickAssist Technology (QAT) integrates hardware acceleration for compute intensive workloads:
  - ✓ Bulk Cryptography
  - ✓ Public Key Exchange
  - ✓ Compression

- Formfactors: Chipset, PCIE card, SoC
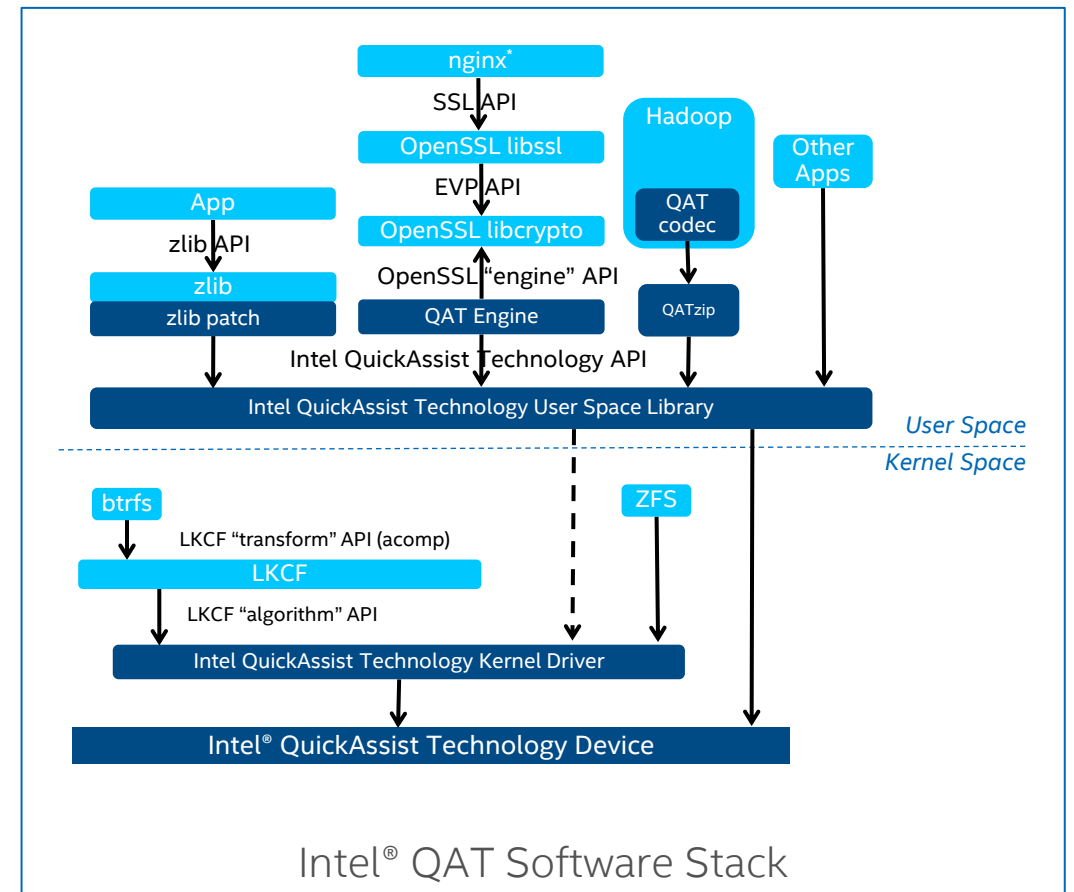
Intel® C620 Series Chipsets
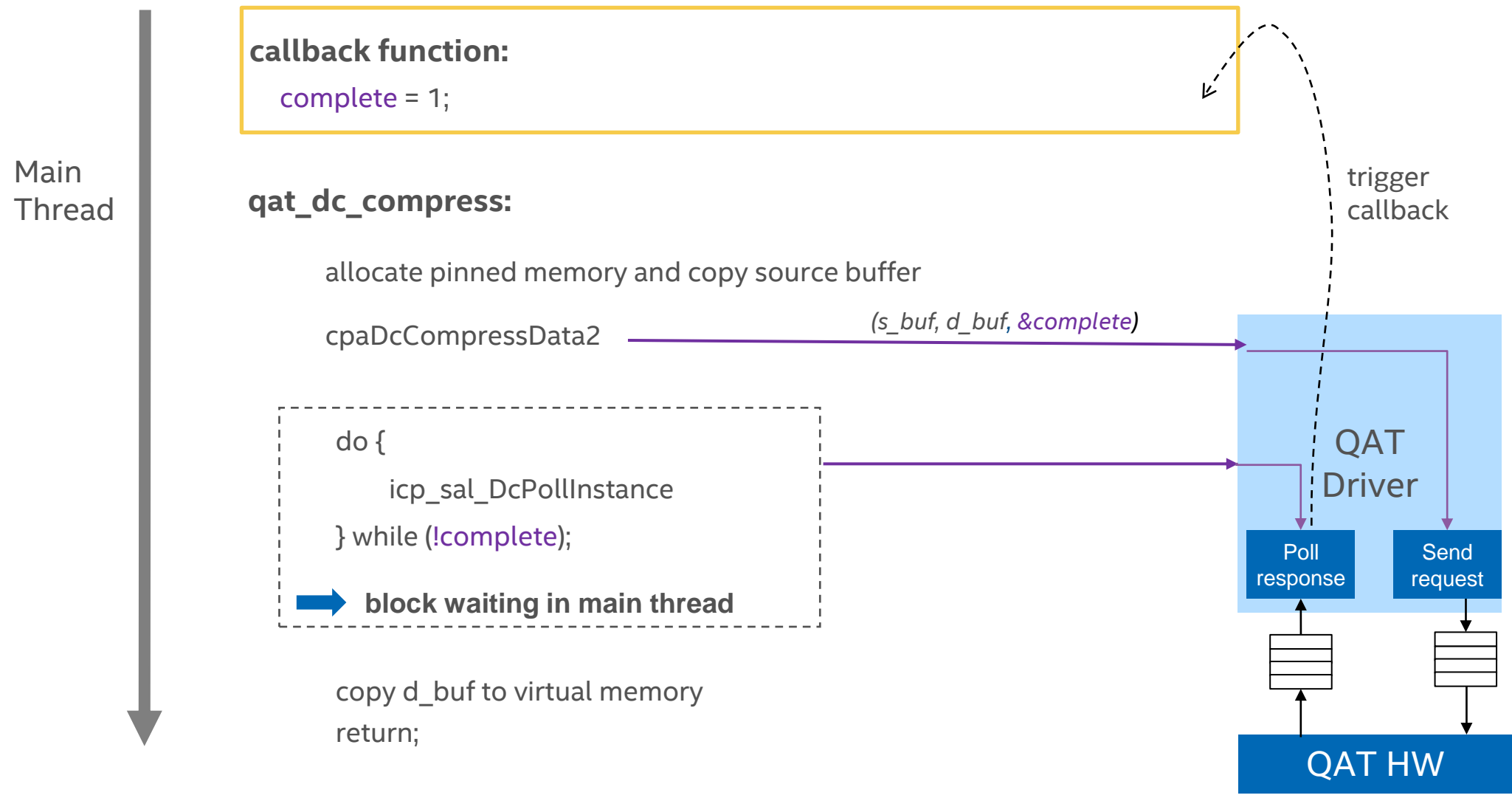
Intel® QuickAssist Adapter 89xx

Rangeley, Denverton

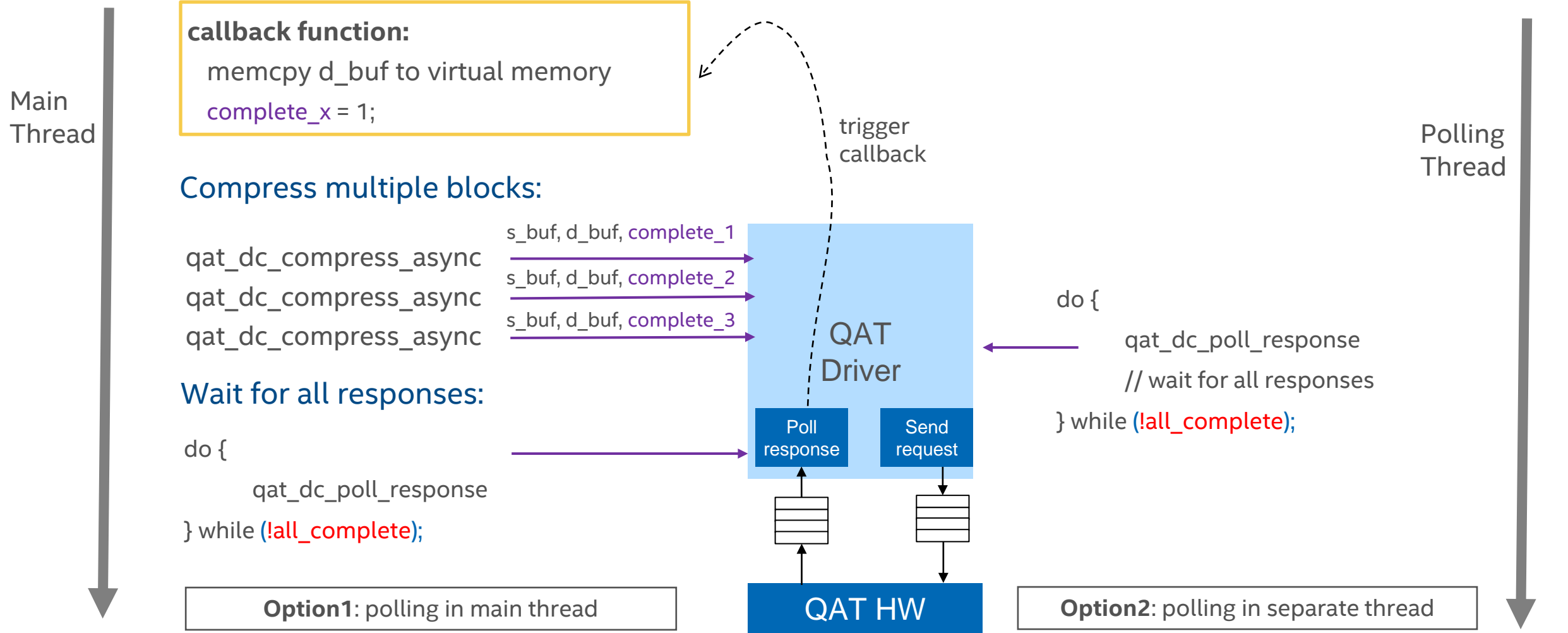- Integrated with many mainstream software framework



Intel® QAT Software Stack

https://01.org/intel-quickassist-technology

# QAT Compression Sync Mode

Main
Thread

**callback function:**

    complete = 1;

trigger
callback

**qat_dc_compress:**

    allocate pinned memory and copy source buffer

    cpaDcCompressData2 ————————— *(s_buf, d_buf, &complete)* ————→

    do {

        icp_sal_DcPollInstance

    } while (!complete);

    ➡ **block waiting in main thread**

    copy d_buf to virtual memory
    return;

QAT
Driver

Poll
response

Send
request

QAT HW

# QAT Compression Async Mode

Main Thread

**callback function:**
memcpy d_buf to virtual memory
complete_x = 1;

Polling Thread

trigger callback

## Compress multiple blocks:

qat_dc_compress_async → s_buf, d_buf, complete_1
qat_dc_compress_async → s_buf, d_buf, complete_2
qat_dc_compress_async → s_buf, d_buf, complete_3

QAT Driver

do {

qat_dc_poll_response

// wait for all responses

} while (!all_complete);

## Wait for all responses:

do {

    qat_dc_poll_response

} while (!all_complete);

Poll response

Send request

QAT HW

**Option1**: polling in main thread

**Option2**: polling in separate thread

# Compress_timing

For BS = 4KB, 8KB, 16KB, ... 512KB:

Divide 3.2MB Calgary Corpus buffer to block size:

| BS | BS | BS | BS | ...... | BS |

daos_compressor_compress(_async)

Iterations = 1000

Calculate performance: throughput, ratio
Decompress and verify result

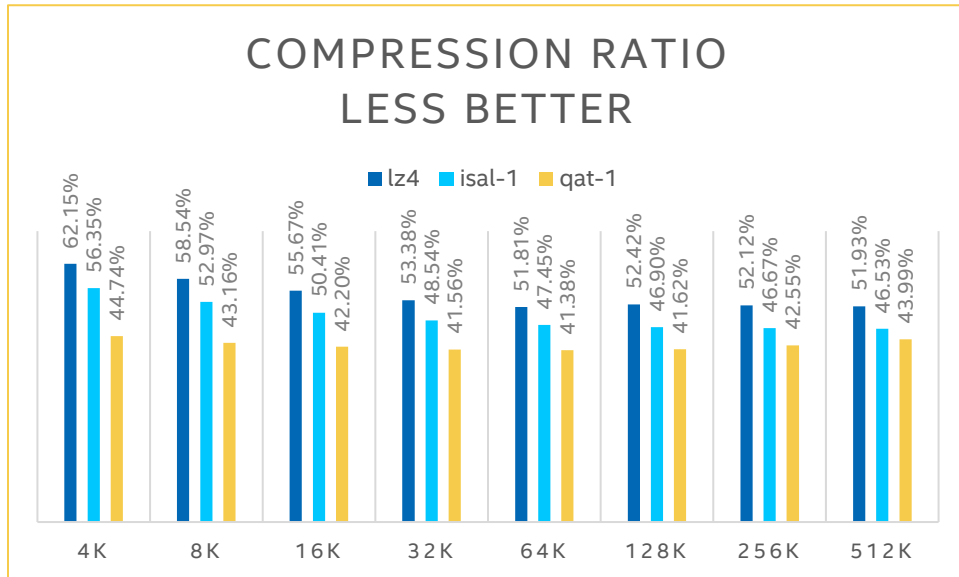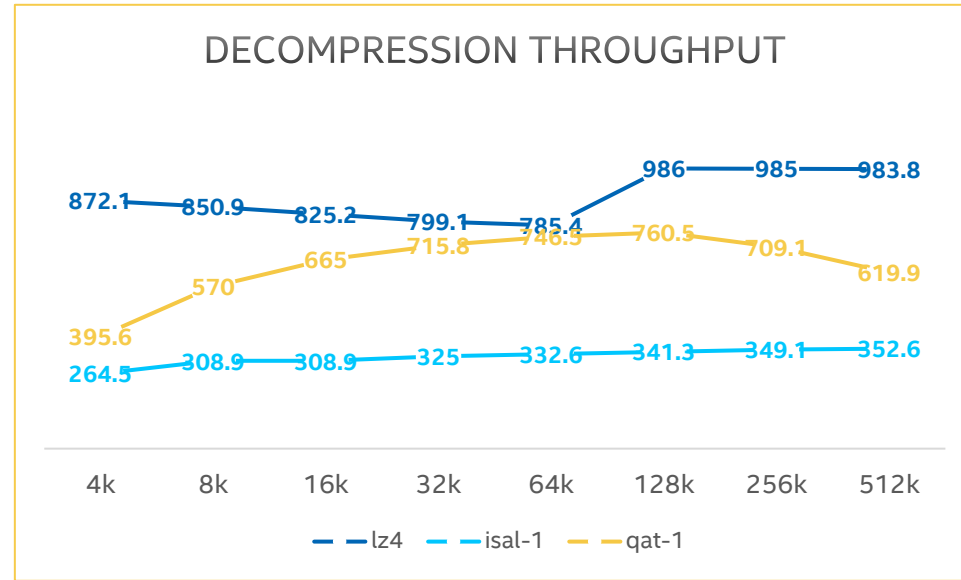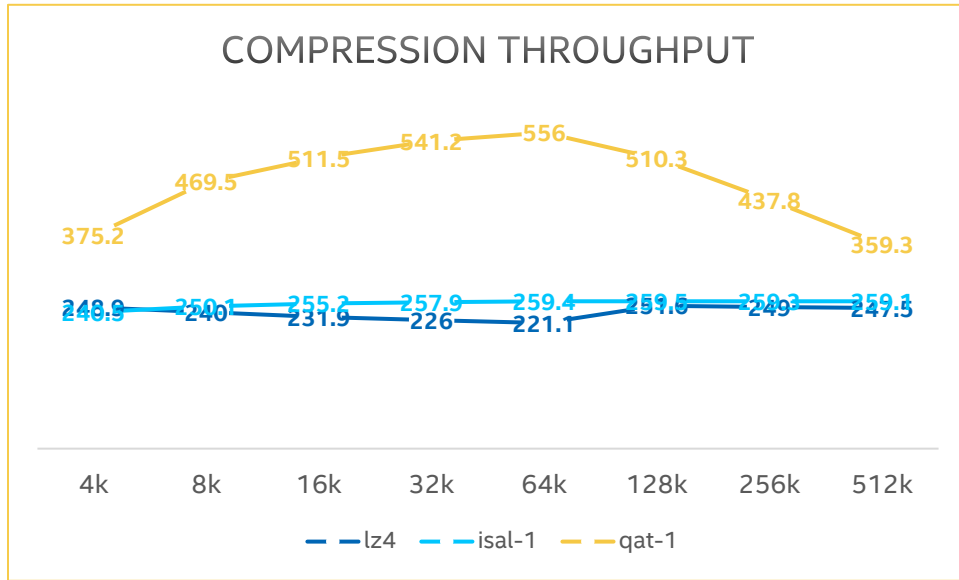Repeat every block size

Single Thread

```
[root@qat-weigang-wildcat1 daos]# build/dev/gcc/src/common/tests/compress_timing
File size:      3175KB
Block Size:     4 KB
        lz4:            comp    15.8 usec       247.5 MB/s      62.15%
        lz4:            decomp  4.5 usec        870.0 MB/s      Pass
        deflate:        comp    10.5 usec       371.8 MB/s      44.74%
        deflate:        decomp  10.0 usec       392.5 MB/s      Pass
        deflate1:       comp    10.5 usec       373.5 MB/s      44.74%
        deflate1:       decomp  9.9 usec        393.3 MB/s      Pass
        deflate2:       comp    10.2 usec       382.3 MB/s      43.13%
        deflate2:       decomp  9.8 usec        397.5 MB/s      Pass
        deflate3:       comp    10.2 usec       383.8 MB/s      42.85%
        deflate3:       decomp  9.8 usec        396.7 MB/s      Pass
        deflate4:       comp    10.1 usec       386.3 MB/s      42.74%
        deflate4:       decomp  9.8 usec        397.6 MB/s      Pass
Block Size:     8 KB
        lz4:            comp    32.6 usec       239.4 MB/s      58.54%
        lz4:            decomp  9.2 usec        848.5 MB/s      Pass
        deflate:        comp    16.7 usec       468.9 MB/s      43.16%
        deflate:        decomp  13.7 usec       568.7 MB/s      Pass
        deflate1:       comp    16.7 usec       469.0 MB/s      43.16%
        deflate1:       decomp  13.7 usec       569.1 MB/s      Pass
        deflate2:       comp    17.1 usec       458.1 MB/s      41.09%
        deflate2:       decomp  14.1 usec       553.5 MB/s      Pass
        deflate3:       comp    17.2 usec       453.2 MB/s      40.65%
        deflate3:       decomp  14.2 usec       551.4 MB/s      Pass
        deflate4:       comp    17.0 usec       460.6 MB/s      40.45%
        deflate4:       decomp  13.5 usec       576.7 MB/s      Pass
```

# Compress_timing results @64KB block size

### Compression throughput (MB/s) @64KB

| Category | Value |
|----------|-------|
| lz4 | 221.1 |
| isal-1 | 259.4 |
| isal-2 | 227.4 |
| qat-1 | 555.8 |
| qat-2 | 513 |

### Decompression throughput (MB/s) @64KB

| Category | Value |
|----------|-------|
| lz4 | 785.4 |
| isal-1 | 332.2 |
| isal-2 | 374.9 |
| qat-1 | 746.8 |
| qat-2 | 746.3 |

### Compression ratio @64KB (less is better)

| Category | Value |
|----------|-------|
| lz4 | 51.81% |
| isal-1 | 47.45% |
| isal-2 | 38.24% |
| qat-1 | 41.38% |
| qat-2 | 38.48% |

- Tested with QAT Gen-2 DH8950 plugin card
- Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz
- Ratio = compressed_size / origin_size
- Single Thread

# Compress_timing results for all block sizes

### COMPRESSION THROUGHPUT



| | 4k | 8k | 16k | 32k | 64k | 128k | 256k | 512k |
|---|---|---|---|---|---|---|---|---|
| qat-1 | 375.2 | 469.5 | 511.5 | 541.2 | 556 | 510.3 | 437.8 | 359.3 |
| lz4/isal-1 | 248.9 / 240.9 | 250.1 / 240 | 255.2 / 231.9 | 257.9 / 226 | 259.4 / 221.1 | 259.5 / 251.8 | 250.3 / 249 | 259.1 / 247.5 |

lz4 — isal-1 — qat-1

### DECOMPRESSION THROUGHPUT



| | 4k | 8k | 16k | 32k | 64k | 128k | 256k | 512k |
|---|---|---|---|---|---|---|---|---|
| lz4 | 872.1 | 850.9 | 825.2 | 799.1 | 785.4 | 986 | 985 | 983.8 |
| qat-1 | 395.6 | 570 | 665 | 715.8 | 746.5 | 760.5 | 709.1 | 619.9 |
| isal-1 | 264.5 | 308.9 | 308.9 | 325 | 332.6 | 341.3 | 349.1 | 352.6 |

lz4 — isal-1 — qat-1

### COMPRESSION RATIO LESS BETTER



lz4 ■ isal-1 ■ qat-1

| | 4K | 8K | 16K | 32K | 64K | 128K | 256K | 512K |
|---|---|---|---|---|---|---|---|---|
| lz4 | 62.15% | 58.54% | 55.67% | 53.38% | 51.81% | 52.42% | 52.12% | 51.93% |
| isal-1 | 56.35% | 52.97% | 50.41% | 48.54% | 47.45% | 46.90% | 46.67% | 46.53% |
| qat-1 | 44.74% | 43.16% | 42.20% | 41.56% | 41.38% | 41.62% | 42.55% | 43.99% |

- QAT
    - ✓ Best compression throughput
    - ✓ Best compression ratio
- LZ4
    - ✓ Best decompression throughput

# Future Work

- Implement actual extents/records compression on both client and server

- Change aggregation to support compression

- Support partial read in the middle of a compressed extent

  - Similar to end-to-end checksum support where we have to read the whole extent to validate the checksum. We will just reuse this code.

- Automatically discover QAT at build time and runtime, intelligently switch between software and hardware acceleration API

- Chaining of compression and hash in one function call

# Legal Disclaimer