# Data reduction and processing for photon science detectors

**David Pennicard** [1,*]**, Vahid Rahmani** [1] **and Heinz Graafsma** [1,2]

[1]*Deutsches Elektronen-Synchrotron DESY, Notkestr. 85, 22607 Hamburg, Germany*
[2]*Mid Sweden University, Sundsvall, Sweden*

Correspondence*:
Corresponding Author
david.pennicard@desy.de

## ABSTRACT

New detectors in photon science experiments produce rapidly-growing volumes of data. For detector developers, this poses two challenges; firstly, raw data streams from detectors must be converted to meaningful images at ever-higher rates, and secondly, there is an increasing need for data reduction relatively early in the data processing chain. An overview of data correction and reduction is presented, with an emphasis on how different data reduction methods apply to different experiments in photon science. These methods can be implemented in different hardware (e.g. CPU, GPU or FPGA) and in different stages of a detector's data acquisition (DAQ) chain; the strengths and weaknesses of these different approaches are discussed.

**Keywords: Photon science, detectors, X-rays, data processing, data reduction, hardware acceleration, DAQ**

## 1 INTRODUCTION

Developments in photon science sources and detectors have led to rapidly-growing data rates and volumes [1]. For example, experiments at the recently-upgraded ESRF EBS can potentially produce a total of a petabyte of data per day, and future detectors targeting frame rates over 100 kHz will have data rates (for raw data) exceeding 1 Tbit/s [2].

These improvements not only allow a much higher throughput of experiments, but also make new measurements feasible. For example, by focusing the beam and raster-scanning it across a sample at high speed, essentially any X-ray technique can be used as a form of microscopy, obtaining atomic-scale structure and chemical information about large samples. But naturally, these increasing data rates pose a variety of challenges for data storage and analysis. In particular, there is increasing demand for data reduction, to ensure that the volume of data that needs to transferred and stored is not unfeasibly large. From the perspective of detector developers, there are two key issues that need to be addressed.

Firstly, the raw data stream from a detector needs to be converted into meaningful images, as discussed in section 2, and this becomes increasingly challenging at high data rates. This conversion process is detector-specific, so implementing it requires detailed knowledge of the detector's characteristics. At the same time, since the correction process is relatively fixed, there's a lot of potential to use accelerators such as GPUs or FPGAs to do this. In addition, the complexity of this conversion process depends on the detector design, so this is something that should be considered during detector development.

29    Secondly, it can be beneficial to perform data reduction on-detector, or as part of the detector's DAQ
30    system. For a variety of reasons, the useful information in a dataset can be captured with a smaller number
31    of bits than the original raw data size; for example, by taking advantage of patterns or redundancy in the
32    raw data, or by rejecting non-useful images in the dataset. In the DAQ system, the data will typically
33    pass through a series of stages, as illustrated in Fig. 1 - firstly, from the ASIC or monolithic sensor to a
34    custom board in the detector, then to a specialized DAQ PC (or similar hardware), and finally to a more
35    conventional computing environment. By performing data reduction early on in this chain, it is possible to
36    reduce the bandwidth required by later stages. Not only can this reduce the cost and complexity of later
37    stages, it can also potentially enable the development of faster detectors by overcoming data bandwidth
38    bottlenecks. Performing this early data reduction often ties together with the process of converting the raw
39    data stream to real images, since these can be easier to compress. Conversely, though, performing data
40    reduction early in the chain can be more challenging, since the hardware in these early stages tends to offer
41    less flexibility, and there are constraints on space and power consumption within the detector.

42    This paper presents an overview of image correction and data reduction in photon science, with a
43    particular focus on how the characteristics of different detectors and experiments affect the choice of data
44    reduction algorithm. After a discussion of the features of common data processing hardware - CPUs, GPUs
45    and FPGAs - the paper returns to the topic of how different compression methods may be implemented in
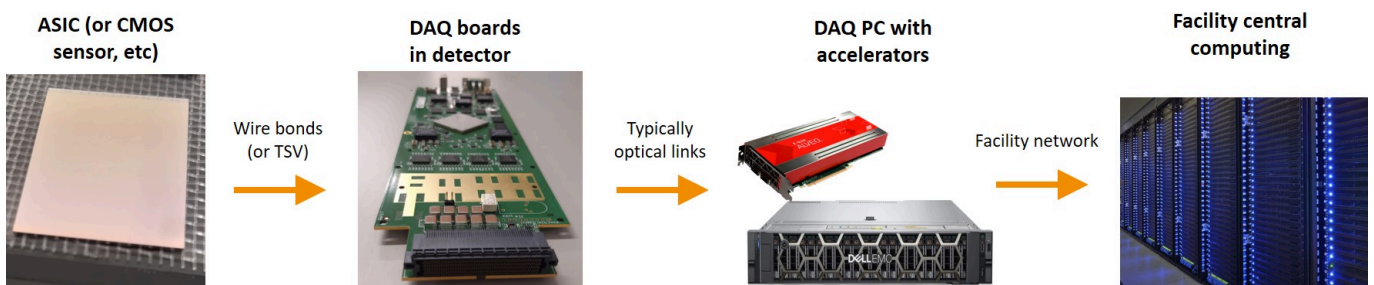46    different stages of a detector's data acquisition chin.



**Figure 1.** Illustration of typical elements in a detector's DAQ chain.

## 2   DETECTOR DATA CORRECTION

47    The raw data stream produced by a detector generally requires processing in order to produce a meaningful
48    image. The steps will of course depend on the design; here, two common cases of photon counting
49    and integrating pixel detectors are considered. Although it can be possible to compress data before all
50    corrections are applied, corrected data can be more compressible - for example, correcting pixel-to-pixel
51    variations can result in a more uniform image.

52    Firstly, we want the pixels in an image to follow a straightforward ordering; typically this is row-by-row
53    or column-by-column, though some image formats represent images as a series of blocks for performance
54    reasons. However, data streams from detectors often have a more complex ordering. One reason for this
55    is that data readout from an ASIC or monolithic detector is typically parallelized across multiple readout
56    signal lines which can result in interleaving of data. In addition, in detectors composed of multiple chips or
57    modules, there may be gaps in the image, or some parts of the detector may be rotated - this is illustrated
58    for the AGIPD 1M detector [3] in Fig 2. So, data reordering is a common first step.
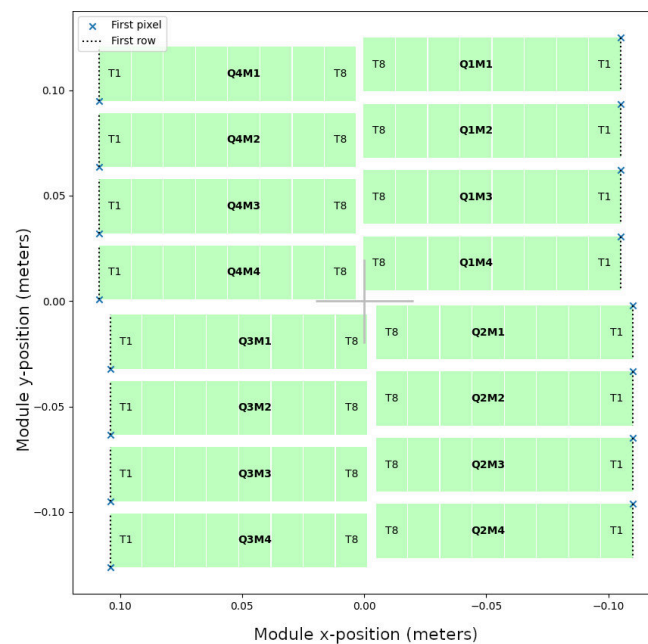
**Figure 2.** Module layout of the AGIPD 1M detector, which has a central hole for the beam, gaps between modules, and different module orientations (with the first pixel of each module indicated by a cross).

59   In the case of photon counting detectors, the value read out from each pixel correspond relatively directly
60   to the number of photons hitting the pixel. Nevertheless, at higher count rates losses occur due pulse
61   pileup, when photons hit a pixel in quick succession. So, pileup correction is needed. Although there are
62   two well-known models for pileup - paralyzable and non-paralyzable - in practice the behaviour of pixel
63   detectors can be somewhere in-between, and in addition to this some modern detectors have additional
64   pileup compensation, for example by detecting longer pulses that would indicate pileup [4] . So, the pileup
65   correction process can vary between different detectors.

66   In integrating detectors, each pixel's amplifier produces an analog value, which is then digitized. This
67   digitized value then needs to be converted to the energy deposited in the pixel. In many detectors, this is
68   a linear relationship which can be described in terms of a baseline and gain; these parameters can vary
69   from pixel-to-pixel. As discussed later, if the incoming X-ray beam is monochromatic then the measured
70   energy may then be converted to photons. In some integrating detectors designed for large dynamic range,
71   the response may not be a simple linear one. For example, in dynamic gain switching detectors [3] each
72   pixel can adjust its gain in response to the magnitude of the incoming signal, and the output of each pixel
73   consists of a digitized value plus information on which gain setting was used; for each gain setting, there is
74   a different baseline and gain that need to be applied to calculate the energy deposited in the pixel.

75   In practice, detectors may deviate from the ideal response, and additional corrections may be required.
76   For example, the response of a detector may not be fully linear, and more complex functions may be need
77   to describe their response. There may also be phenomena such as common-mode signals or crosstalk,
78   where corrections need to be applied that are dependent on the magnitude and pattern of the incoming
79   signal rather than independent for each pixel. It is also common to treat malfunctioning pixels, for example
80   by setting them to some special value.

81    An additional aspect of detector data processing is combining the image data with metadata. Some
82  metadata may be directly incorporated into the detector's data stream. For example, in Free Electron Laser
83  (FEL) experiments the detector needs to be synchronised with the X-ray bunches, and bunches can vary in
84  quality, so each image will be accompanied by a bunch ID, fed to the detector from the facility's control
85  system. Other metadata may be added later. For example, the NeXuS data format [5] has been adopted
86  by many labs; this is based on the HDF5 format [6], and specifies how metadata should be structured in
87  experiments in photon science and other fields.

## 3 DATA COMPRESSION AND PHOTON SCIENCE DATASETS

88  In general, data compression algorithms reduce the number of bits needed to represent data, by encoding it
89  in a way that takes advantage of patterns or redundancy in the data. These algorithms can be subdivided
90  into lossless compression [7], where the original data can be reconstructed with no error, and lossy
91  compression [8], where there is some difference between the original and reconstructed data. Currently,
92  lossless compression is often applied to photon science data before storage, whereas lossy compression
93  is uncommon, due to concerns about degrading or biasing the results of later analysis. However, lossy
94  compression can achieve higher compression ratios.

95    The compressibility of an image naturally depends on its content, and in photon science this can vary
96  depending on both the type of experiment and the detector characteristics. In particular, since random
97  noise will not have any particular redundancy or pattern, the presence of noise in an image will reduce the
98  amount of compression that can be achieved with lossless algorithms.

99    Compared to typical visible light images, X-ray images from pixel detectors can have distinctive features
100  that affect their compressibilty, as discussed further below. Firstly, individual X-ray photons have much
101  greater energy and thus can more easily be discriminated with a suitable detector. Also, X-ray diffraction
102  patterns have characteristics that can make lossless compression relatively efficient. However, imaging
103  experiments using scintillators and visible light cameras produce images more akin to conventional visible
104  light imaging, which do not losslessly compress well.

### 3.1 Noise in X-ray images, and its effects on compression

106    Random noise in an image can potentially come from different sources; firstly, electronic noise introduced
107  by the detector, and secondly, inherent statistical variation in the experiment itself.

108    Any readout electronics will inevitably have electronic noise. However, the signal seen by a detector
109  consists of discrete X-ray photons, and the inherent discreteness of our signal makes it possible to reject
110  the electronic noise, provided that it is small compared to the signal produced by a single photon. In a
111  silicon sensor without gain, we will generate on average one electron-hole pair per 3.6 eV of energy, e.g. a
112  12 keV photon will generate approximately 3300 electrons. In turn, if for example the electronic noise is
113  Gaussian with a standard deviation corresponding to 0.1 times the photon energy (330 electrons here) then
114  the probability of a pixel having a noise fluctuation corresponding to 0.5 photons or more is approximately
115  1 in 1.7 million [9].

116    In integrating pixel detectors, each pixel measures the total energy deposited in the pixel during the
117  integration time of the image. If the photon source is monochromatic, then it is possible to convert the
118  integrated signal to an equivalent number of photons in postprocessing. Given sufficiently low noise,
119  this value can be quantized to the nearest whole number of photons to eliminate electronic noise. In
120  photon counting detectors, a hit will be recorded in a pixel if the pulse produced by the photon exceeds

121  a user-defined threshold; once again, if the threshold level is far enough from the noise, we will have
122  effectively noise-free counting.

123     In both cases, the electronic noise in a pixel will be dependent on the integration time for an image, or
124  the shaping time in the case of a photon-counting detector. There are two main competing effects here.
125  On the one hand, for longer timescales, the shot noise due to fluctuations in leakage current will be larger.
126  Conversely, to achieve shorter integration times or shaping times, a higher amplifier bandwidth is required,
127  and this will increase the amount of thermal noise detected [10]. Thermal noise is the noise associated with
128  random thermal motion of electrons, which is effectively white noise.

129     In addition to electronic noise, however, the physics of photon emission and interaction are inherently
130  probabilistic, and so even if an experiment were repeated under identical conditions there would be
131  statistical fluctuations in the number of photons impinging on each pixel. These fluctuations follow Poisson
132  statistics, and if the expected number of photons arriving in a pixel during a measurement is $N$, then the
133  standard deviation of the corresponding Poisson distribution will be $\sqrt{(N)}$. On the one hand, this can
134  make it easier to develop "low noise" detectors; provided the detector noise for a given photon flux is
135  significantly smaller than $\sqrt{(N)}$, then it will have little effect on the final result. However, this makes
136  the data less compressible with lossless algorithms, since random noise introduces variation in the image
137  that is not patterned or redundant. (As discussed later, quantizing pixel values with a variable step size,
138  increasing with $\sqrt{(N)}$, can be a way of achieving lossy compression.)

139  ## 3.2  Applying lossless compression to diffraction data

140     As noted above, data compression relies on patterns or redundancy in data, and this will vary depending
141  on the experiment. Diffraction patterns differ from conventional images in a variety of respects, with the
142  pattern being mathematically related to the Fourier transform of the object. An example of a diffraction
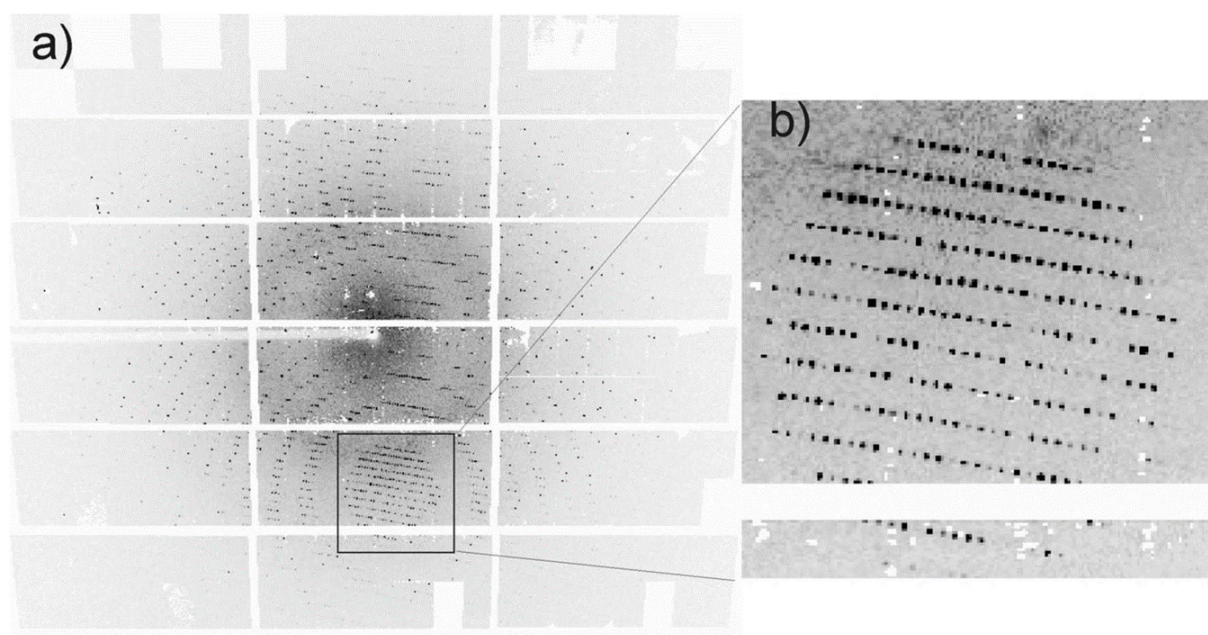143  pattern from macromolecular crystallography is shown in Fig. 3.



**Figure 3.** A diffraction pattern from a thaumatin crystal, recorded with a Pilatus 1M detector. Reproduced
from [11] with permission of the International Union of Crystallography.

144  Diffraction patterns tend to have various features that are advantageous for lossless compression:

145  • Hybrid pixel detectors with sensitivity to single photons are the technology of choice for these
146    experiments, making the measurement effectively free from electronic noise as described above.
147    (Detectors for X-ray diffraction require high sensitivity, but pixel sizes in the range of 50-200 $\mu$m are
148    generally acceptable.)

149  • The intensity values measured in the detector typically have a nonuniform distribution, with most
150    pixels measuring relatively low or even zero photons, and a small proportion of pixels having higher
151    values. This is partly because the diffracted intensity drops rapidly at higher scattering angles, and
152    partly because of interference phenomena that tend to produce high intensity in certain places (e.g.
153    Bragg spots or speckles) and low intensity elsewhere.

154  • As with many images, nearby pixels will tend to have similar intensities.

155  • Combining these points, in images with fewer photons, there can be patches in the image with many
156    neighbouring zero-value pixels.

157  Empirically, a range of lossless compression algorithms can achieve good results with diffraction data,
158  especially in high-frame-rate measurements where the number of photons per image will tend to be
159  lower. Experiments applying the DEFLATE [12] algorithm used in GZIP to datasets acquired with photon
160  counting detectors running at high speed showed compression ratios of 19 for high-energy X-ray diffraction,
161  70 for ptychography and 350 for XPCS experiments with dilute samples (where most pixel values are zero)
162  [13]. Experiments with the Jungfrau integrating detector, applying different compression algorithms to the
163  same data, found that multiple compression methods such as GZIP, LZ4 with a bitshuffle filter and Zstd
164  gave similar compression ratios to one another, but varied greatly in speed, with GZIP being a factor of 10
165  slower [14].

166  3.2.1  Example - lossless compression with DEFLATE (GZIP) and Bitshuffle/LZ4

167  As illustrative examples, we consider the DEFLATE algorithm used in GZIP [15], and the Bitshuffle/LZ4
168  algorithm [16], firstly to discuss how they take advantage of redundancy to compress diffraction data, and
169  secondly how algorithms with different computational cost can achieve similar performance.

170  DEFLATE [12] consists of two stages. Firstly, the LZ77 algorithm [17] is applied to the data, which
171  searches for recurring sequences of characters in a file (such as repeated words or phrases in text, or
172  long runs of the same character) and encodes them efficiently. In the output of this algorithm, the first
173  instance of a sequence of characters is written in full, but then later instances are replaced with special
174  codes referring back to the previous instance. In diffraction datasets, long recurring sequences of characters
175  are generally rare, but there is one big exception; long runs of zeroes. So, the pattern-matching in LZ77
176  will efficiently compress long runs of zeroes, but the computational work the algorithm does to find more
177  complex recurring sequences is largely wasted.

178  Secondly, DEFLATE takes the output of the LZ77 stage, and applies Huffman coding [18] to it. Normally,
179  different characters in a dataset (e.g. integers in image data) are all represented with the same number
180  of bits. Huffman coding looks at the frequencies of different characters in the dataset, and produces a
181  new coding scheme that represents common characters with shorter sequences of bits and rare ones with
182  longer sequences. This is analogous to Morse code, where the most common letter in English, "E", is
183  represented by a single dot, whereas rare letters have longer sequences. For diffraction data, this stage will
184  achieve compression due to the nonuniform statistics of pixel values, where low pixel values are much
185  more common than high ones.

186    As a contrasting example, the Bitshuffle LZ4 algorithm [16] implicitly takes advantage of the knowledge
187    that the higher bits of pixel values are mostly zero and strongly correlated between neighbouring pixels.
188    In the first step, Bitshuffle, the bits in the data stream are rearranged so that the first bit from every pixel
189    are all grouped together, then the second bit, etc. After this regrouping, the result will often contain long
190    runs of bytes with value zero, as illustrated in Fig. 4. The LZ4 algorithm, which is similar to LZ77, will
191    then efficiently encode these long runs of zero bytes. As mentioned above, this algorithm gives similar
192    performance to GZIP for X-ray diffraction data while being less computationally expensive.

| Original pixel values | 1 | 0 | 2 | 5 | 7 | 3 | 9 | 5 |
|---|---|---|---|---|---|---|---|---|
| | 00000001 | 00000000 | 00000010 | 00000101 | 00000111 | 00000011 | 00001001 | 00000101 |

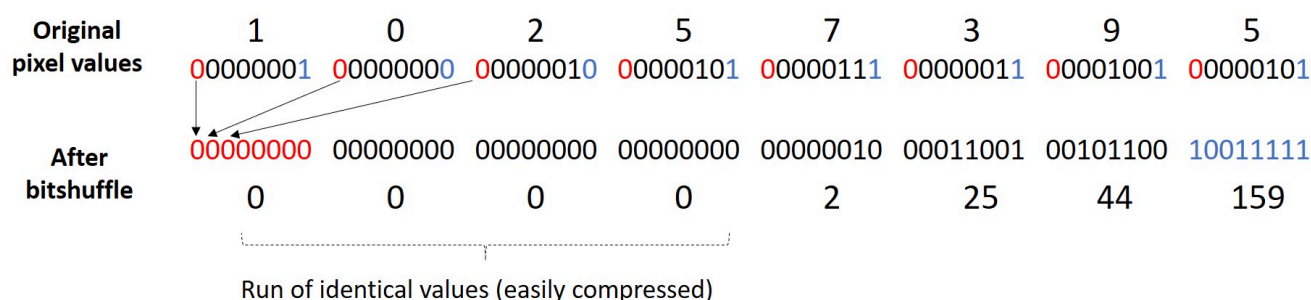| After bitshuffle | 00000000 | 00000000 | 00000000 | 00000000 | 00000010 | 00011001 | 00101100 | 10011111 |
|---|---|---|---|---|---|---|---|---|
| | 0 | 0 | 0 | 0 | 2 | 25 | 44 | 159 |

Run of identical values (easily compressed)

**Figure 4.** Illustration of the Bitshuffle process, used in Bitshuffle/LZ4. This reorders the bits in the data stream so that the LZ4 compression stage can take advantage of the fact that the upper bits in diffraction data are mostly zero, and strongly correlated between neighbouring pixels.

## 3.3   Applying lossy compression to imaging data

194    While lossless compression can achieve good compression ratios for diffraction data, there is increasing
195    demand for lossy data compression. Firstly, data from some experiments such as imaging do not losslessly
196    compress well, due to noise, and secondly, the increasing data volumes produced by new detectors and
197    experiments mean that higher compression ratios are desirable. The key challenge of lossy compression
198    is ensuring that the compression does not significantly change the final result of analysis. This requires
199    evaluating the results of the compression with a variety of datasets, either by directly comparing the
200    compressed and uncompressed images with a metric of similarity, or by performing the data analysis and
201    applying some appropriate metric of quality to the final result.

202    In X-ray imaging and tomography experiments, the detector measures X-ray transmission through the
203    sample, and perhaps also additional effects such as enhancement of edges through phase contrast. This
204    means that most pixels will receive a reasonably high X-ray flux, in contrast to X-ray diffraction where
205    most pixels see few or even zero photons. So, noise due to Poisson statistics will be relatively high in
206    most pixels. Likewise, a key requirement for detectors in these experiments is a small effective pixel size,
207    whereas single photon sensitivity is less critical. To achieve this, the typical approach is to couple a thin
208    scintillator to a visible light camera with small pixels such as a CMOS sensor, using magnifying optics
209    [19]. This means that the detector noise is also non-negligible.

210    Since an X-ray transmission image is a real-space image of an object, and broadly resembles a
211    conventional photograph (unlike a diffraction pattern), widely-used lossy compression algorithms for
212    images can potentially be used for compression. In particular, JPEG2000 [20] is already widely used
213    in medical imaging, and in tomography at synchrotrons it has been demonstrated to achieve a factor of
214    3-4 compression without significantly affecting the results of the reconstruction [21]. To test this, the
215    reconstruction of the object was performed both before and after compression, and the two results compared

216 using the Mean Structural Similarity Index Measure (MSSIM) metric. A compression factor of 6-8 was
217 possible in data with a high signal-to-noise ratio.

### 3.3.1 Example - lossy compression with JPEG2000

219 One common approach to both lossy and lossless compression is apply a transform to the data that
220 results in a sparse representation, i.e. most of the resulting values are low or zero. (The choice of transform
221 naturally depends on the characteristics of the data.) The sparse representation can then be compressed
222 efficiently.

223 In the case of JPEG2000 [20], the Discrete Wavelet Transform [22] is used; in effect, this transformation
224 represents the image of a sum of wave packets with different positions and spatial frequencies. This tends
225 to work well for real space images, which tend to consist of a combination of localized objects and smooth
226 gradients, which can be found on different length scales. After applying the transform, the resulting values
227 are typically rounded off to some level of accuracy, allowing for varying degrees of lossy compression.
228 (By not applying this rounding, lossless compression is also possible.) Then, these values are encoded by
229 a method called arithmetic coding, which is comparable to Huffman coding; it achieves compression by
230 taking advantage of the nonuniform statistics of the transformed values, which in this case are mostly small
231 or zero.

232 JPEG2000 can be compared with the older JPEG standard, where the transformation consists of splitting
233 the image into $8 \times 8$ blocks, and applying the Discrete Cosine Transform to each block, thus taking advantage
234 of the fact that images tend to be locally smooth. This is computationally cheaper than JPEG2000, but one
235 key drawback is that this can lead to discontinuities between these $8 \times 8$ blocks after compression.

## 3.4 Novel methods for lossy compression

237 As mentioned previously, increasing data volumes mean there is demand for achieving increasing
238 compression, even for experiments such as X-ray diffraction where lossless compression works reasonably
239 well. Naturally, this can be approached by testing a variety of well-established lossy compression algorithms
240 on data, and experimenting with with methods such as rounding the data to some level of accuracy. However,
241 there are also new lossy compression methods being developed specifically for scientific data. One particular
242 point of contrast is that most image compression algorithms focus on minimizing the perceptible difference
243 to a human viewer, whereas for scientific data other criteria can be more important, such as imposing limits
244 on the maximum error allowed between original values and compressed values.

245 One simple example is quantizing X-ray data with a step size smaller than the Poisson noise, which
246 is proportional to $\sqrt{(N)}$. For ptychography, for example, it has been demonstrated that quantizing pixel
247 values with a step size of $0.5 \times \sqrt{(N)}$ [23] does not degrade the quality of the reconstruction.

248 A more sophisticated example of error-bounded lossy compression is the SZ algorithm [24], which is a
249 method for compressing a series of floating-point values. For each new element in the series, the algorithm
250 checks if its value can be successfully "predicted" (within a specified margin of error) using any one of
251 three methods: directly taking the previous value; linear extrapolation from the previous two values; or
252 quadratically extrapolating from the previous three values. If so, then the pixel value is represented by a
253 2-bit code indicating the appropriate prediction method. If not, then this "unpredictable" value needs to be
254 stored explicitly. This is illustrated in Fig. 5. After the algorithm runs, further compression is applied to the
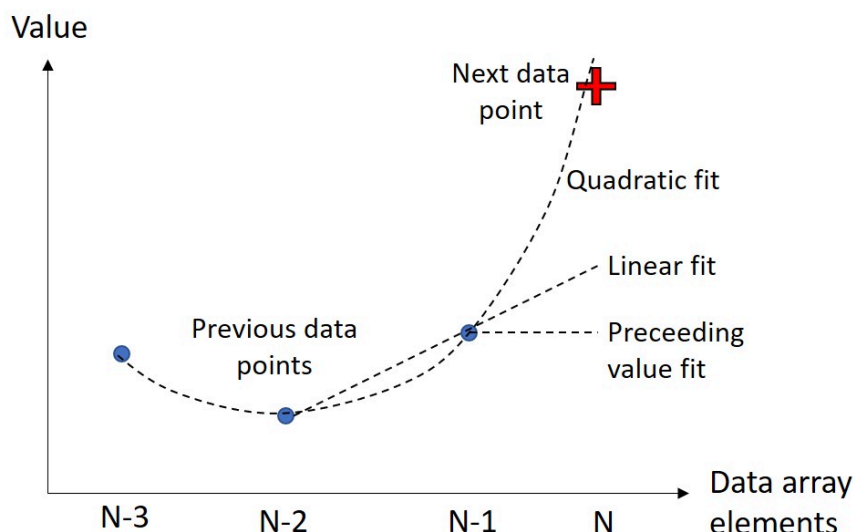255 list of unpredictable values.

**Figure 5.** Illustration of the "predictive" approach used in SZ compression of floating-point data. If the value of the next data point can be extrapolated (within a user-defined margin of error) from previous data points using one of 3 methods, the point is represented by a 2-bit code indicating the method, e.g. quadratic in this case.

256    Underwood et al. in [25] reports on applying SZ to serial crystallography data from LCLS, where the
257   images consist of floating point values obtained from an integrating detector. In this method, lossless
258   compression is applied to regions of interest consisting of Bragg peaks detected in the image, while binning
259   followed by lossy compression is applied to the rest of the image, with the goal of ensuring that any weak
260   peaks that went undetected will still have their intensities preserved sufficiently well. It is reported that
261   when using this strategy, using SZ for the lossy compression can achieve a compression ratio of 190 while
262   still achieving acceptable data quality, whereas other lossy compression methods tested gave a compression
263   ratio of 35 at best.

## 4   OTHER FORMS OF DATA REDUCTION

264   The data compression methods discussed thus far work by representing a given file using fewer bits;
265   the original file can be reconstructed from the compressed file (albeit with some inaccuracy in the case
266   of lossy compression). More broadly, though, there are other methods of data reduction which rely on
267   eliminating non-useful data entirely, or transforming the data in a non-recoverable way. These methods
268   have the potential to greatly reduce the amount of data needing to be stored, though of course it is crucial
269   to establish that these methods are reliable before putting them into practice.

### 4.1   Data rejection / vetoing

270

271    In some experiments, a large fraction of the data collected does not contain useful information. For
272   example, in experiments at FELs such as serial crystallography and single particle imaging, the sample can
273   consist of a liquid jet containing protein nanocrystals or objects such as viruses passing through the path
274   of the beam. However, only a small fraction of X-ray pulses (in some cases, of order 1%) actually hit a
275   sample to produce an image containing a useful diffraction pattern. So, in these kind of experiments, data
276   volumes can be greatly reduced by rejecting images where the beam did not the target.

277     A variety of methods have been developed for doing this. In serial crystallography, images where the
278 beam hit a crystal will contain Bragg peaks, whereas miss images will only have scattering from the liquid
279 jet. So, a common approach is to search for Bragg peaks in each image, and only accept images where the
280 total number of peaks exceeds some threshold value; this approach is used in software such as Cheetah
281 [26] and DIALS [27]. Since X-ray diffraction intensity varies with scattering angle, these methods often
282 rely on estimating the background signal in the surrounding area of the image when judging whether a
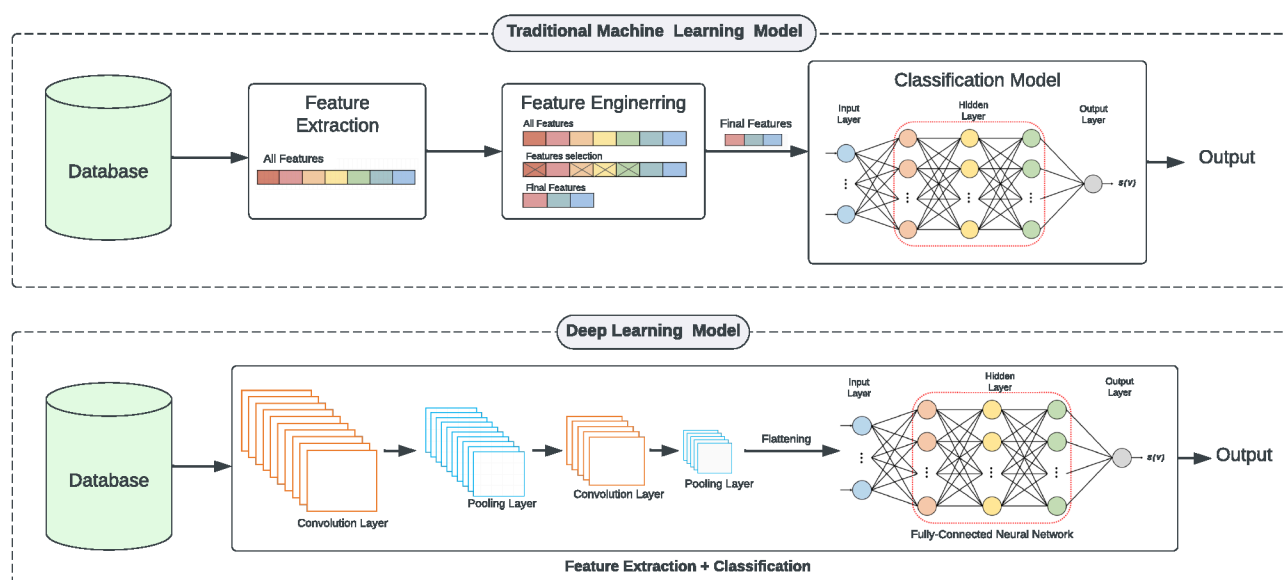283 peak is present or not.



**Figure 6.** Classical Machine Learning and Deep Learning

284     Machine learning techniques have also been applied to this task, where supervised learning is used to
285 distinguish between good and bad images based on training data from simulation or previous experiments
286 [28, 29, 30, 31].As illustrated in Fig. 6, there are a variety of methods for doing this; for example, in
287 "traditional" machine learning, features are first extracted from images using algorithms from computer
288 vision, then machine learning is applied to classify images based on these features, whereas in deep
289 learning, the features themselves are also learned using convolutional layers. These methods can potentially
290 be computationally cheaper than peak-counting methods, or require less expert fine-tuning of parameters.
291 However, in some cases a machine learning algorithm trained on one dataset may fail to generalize to
292 others, so it is important for example to validate machine learning methods using different datasets.

### 4.2   Data reduction by processing

294     Often, the final result of data analysis is much more compact than the original dataset; in determining a
295 protein structure, for example, the data may consist of tens of thousands of diffraction images, while the
296 resulting structure can be described as a list of atomic positions in the molecule.

297     In some cases the full analysis of a dataset by the user may take months or years, and a lot of this analysis
298 is very experiment-specific, so this is not suitable for achieving fast data reduction close to the detector.
299 Nevertheless, there can be initial processing steps that are common to multiple experiments and which
300 could be applied quickly as part of the DAQ system.

301     One example of this is azimuthal integration. In some X-ray diffraction experiments, such as powder
302 diffraction, the signal on the detector has rotational symmetry, and so the data can be reduced to a 1-D
303 profile of X-ray intensity as a function of diffraction angle. For a multi-megapixel-sized detector, this
304 corresponds to a factor of 1000 reduction in data size. Hardware-accelerated implementations of azimuthal
305 integration have been developed for GPUs [32], and more recently for FPGAs [33].

## 4.3   Connections between data reduction and on-the-fly analysis

306

307     If data reduction involves performing analysis steps on data (e.g. azimuthal integration) or makes the data
308 easier to process (e.g. making the dataset smaller by rejecting bad data) then this synergises with on-the-fly
309 data processing, where the data is processed immediately during the course of the experiment rather than
310 being saved to disk and processed at a later point.

311     Typically, these sort of data reduction methods have the drawback that if some or all of the original data
312 is discarded, then any errors in the processing cannot be corrected. For example, performing azimuthal
313 integration correctly requires accurate calibration of the centre of the diffraction pattern and any tilt of the
314 detector relative to the beam.

315     However, since on-the-fly data analysis is useful to users, the trustworthiness of these methods can
316 be established in stages. For example, as a first stage all the raw data may be saved to disk for later
317 analysis, with on-the-fly processing providing quick feedback during the experiment. Once the reliability
318 of the automatic processing is better-established, raw data might only be stored for a shorter time before
319 deletion, to give the opportunity for the user to cross-check for correctness. Finally, once the on-the-fly data
320 processing is fully trusted, it may be possible to only save a small fraction of the raw data for validation
321 purposes.

# 5   HARDWARE-ACCELERATED IMAGE CORRECTION AND DATA REDUCTION

322 ”Off the shelf” data processing hardware of course naturally plays a key role in building data acquisition
323 systems for detectors; even custom hardware such as circuit boards for detector control will incorporate
324 commercial parts such as microcontrollers or FPGAs.

325     Algorithms can be designed or adapted to take advantage of parallelism on both CPUs and GPUs or
326 pipelining of FPGAs. For instance, algorithms like Fast Fourier Transform (FFT) [34], image processing
327 filters [35], and some statistical calculations [36] can be parallelized effectively. Furthermore, in the realm
328 of data reduction, machine learning models rely significantly on hardware acceleration to address their
329 computational requirements [37, 38, 39].

330     Here, we give an overview of the distinctive features of CPUs, GPUs and FPGAs for data processing.
331 In particular, a major aspect of high-performance computing is parallelization of work, where a data
332 processing task is divided between many processing units, but how this is achieved varies between devices,
333 which can affect the best choice of hardware for the task.

## 5.1   CPUs - Central Processing Units

334

335     The architecture of a typical CPU is shown in Fig. 7. A key feature of modern CPUs is that they're
336 designed to be able to run many different programs in a way that appears simultaneous to the user. Most
337 modern CPUs are composed of a number of cores; for example, CPUs in the high-end AMD EPYC
338 series have from 32 to 128, though desktop . Each core has a control unit, which is able to independently
339 interpret a series of instructions in software and execute them, along with an arithmetic unit for performing

340 operations on data, and cache memory for temporarily storing data. Furthermore, CPUs cores can very
341 rapidly switch between executing different tasks, so even a single CPU core can perform many tasks in a
342 way that seems parallel to the user. In most cases, the CPU will be running an operating system, which
343 handles the sharing of resources such as memory and peripherals between the tasks in a convenient way for
344 programmers. High performance computing with CPUs generally involves parallelization of work across
345 multiple cores; in doing this, the different cores can operate relatively independently.



**Figure 7.** Architecture of a Quad-core hyperthreading CPU.

346 Compared to GPUs and FPGAs, CPUs are generally the easiest to use and program. Indeed, GPUs are
347 virtually always used as accelerator cards as part of a system with a CPU, and it's common for systems
348 based on FPGAs to make use of a CPU (e.g. as part of a System On Chip). In addition, if a task is inherently
349 sequential and cannot be parallelized, then a CPU core may give better performance than a GPU or FPGA.
350 However, CPUs generally have much less parallel processing power than GPUs or FPGAs.

## 351 **5.2 GPUs - Graphics Processing Units**

352 GPUs are designed to allow massively parallel processing for tasks such as graphics processing or
353 general-purpose computing. The basic units of GPUs are called CUDA cores (in NVIDIA GPUs) or stream
354 processors (in AMD GPUs). These are much more numerous than CPU cores, numbering in the thousands,
355 allowing much greater parallelism and processing power. However, rather than each core being able to
356 independently interpret a stream of instructions, these cores are organised into blocks, with all the cores in
357 a block simultaneously executing the same instruction on different data elements - for example, applying
358 the same mathematical operation to different elements in an array. The structure of a typical GPU is shown
359 in Fig. 8.

360 GPU programming relies on specialized frameworks such as CUDA or OpenCL. GPUs are generally
361 used as accelerator cards in a system with a CPU, and these frameworks firstly allow the CPU to control the
362 GPU's operation (by passing data to and from it, and starting tasks) and secondly to program functions that
363 will run on the GPU. Typically, GPU code contains additional instructions controlling how different GPU
364 cores in a block access data; e.g. when performing an operation on an array, the first core in a block might
365 perform operations on the first element, and so forth. This can be relatively easy to apply to a lot of image
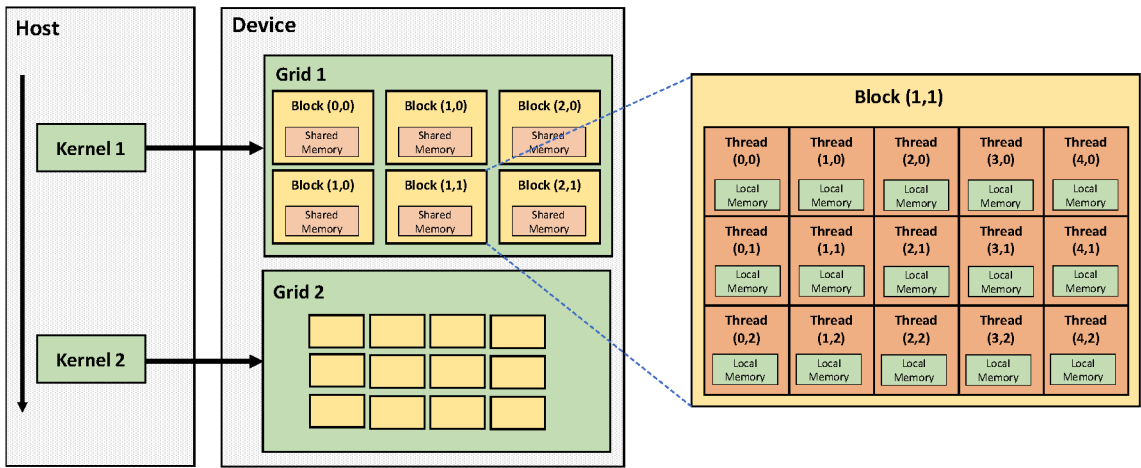
**Figure 8.** NVIDIA GPU CUDA Block and thread architecture.

366 processing tasks, where each processing step can be applied to pixels relatively independently, but GPU
367 implementation can be more challenging for some compression algorithms that act sequentially on data.

## 5.3   FPGAs - Field Programmable Gate Arrays

369   FPGAs are a type of highly configurable hardware. As shown in Fig. 9, within an FPGA there are many
370 blocks providing functionality such as combinatorial logic, digital signal processing and memory, with
371 programmable interconnects between them [40]. While software for a GPU or CPU consists of a series of
372 instructions for the processor to follow, an FPGA's firmware configures the programmable interconnects to
373 produce a circuit providing the required functionality. The large number of blocks in an FPGA can then
374 provide highly parallel processing and high throughput of data.



**Figure 9.** The architecture of the Xilinx UltraScale FPGA accelerator card.

375 Traditionally, firmware is developed by using a hardware description language to describe the required
376 functionality; a compiler then finds a suitable configuration of blocks and interconnects to achieve this.
377 While hardware description languages differ substantially from conventional programming languages,
378 high-level-synthesis tools for FPGAs have been more recently developed to make it possible to describe an
379 algorithm with a more conventional programming language (such as C++), with special directives in the
380 code being used to indicate how parallelism can be achieved with the FPGA.

381 A distinctive form of parallelism in FPGAs is pipelining. A series of processing steps will often be
382 implemented as a series of blocks forming a pipeline, with data being passed along from one block to another.
383 This is analogous to a production line in a factory, where each station carries out a particular fabrication
384 step, and goods pass from one station to another. As with the production line, at any given moment there
385 will be multiple data elements in different stages of processing being worked on simultaneously. In turn,
386 to make algorithms run efficiently in an FPGA, they need to be designed to make effective use of this
387 pipeline parallelism. This means that FPGAs tend to be better-suited to algorithms that operate on relatively
388 continuous streams of data, rather than algorithms that rely on holding large amounts of data in memory
389 and accessing them in an arbitrary way. As well as achieving a high throughput, algorithms on FPGA can
390 give a reliably low latency, which can be important in cases where fast feedback is required. However, it is
391 generally more challenging to program FPGAs than CPUs or GPUs.

## 6 DATA TRANSFER AND PROCESSING CHAINS IN DETECTORS

392 A typical detector and DAQ system consist of series of stages, as illustrated previously in Fig. 1. At each
393 stage, data transfer is needed, and we face a series of bottlenecks which can limit the possible data rate. By
394 doing data reduction earlier on in this chain, the data transfer demands on later stages are reduced, and
395 we can take advantage of this to achieve higher detector speeds. Conversely, the hardware used in earlier
396 stages of this chain is typically more "custom" or specialized, making the implementation of data reduction
397 more challenging, and typically there is also less flexibility in these stages.

398 The stages in the readout and processing chain are as follows, and the potential for data reduction at each
399 stage will be discussed in more detail in the following section.

400 • The readout chip (hybrid pixel) or monolithic sensor.
401 • Readout electronics within the detector itself, which typically incorporate an FPGA, microcontroller
402 or similar.
403 • Detector-specific data processing hardware. Typically, this consists of one or more PCs, which can
404 include peripherals such as FPGAs or GPUs, but it is also possible to build more specialized systems,
405 such as crates of FPGA boards.
406 • Generic high-performance-computing hardware, typically in the facility's computing cluster.

### 6.1 On-chip data reduction

408 The first data bottleneck encountered in the detector is transferring data out of the hybrid pixel readout
409 chip or monolithic sensor. In a pixel detector, we have a 2D array of pixels generating data, and within
410 the chip it's possible to have a high density of signal routing. However, data output takes place across a
411 limited number of transceivers, which are typically located at the periphery of the chip and connect to a
412 circuit board via wire bonds (though there is increasing effort in using Through Silicon Vias - TSVs - for
413 interconnect [41]). In turn, in connecting these to the rest of the readout system, further bottlenecks are

414  faced. To build large detector areas with minimal gaps, it's typical to adopt a modular detector design, with
415  each module's electronics placed behind it; this naturally creates limits on the space available for PCB
416  traces and components such as FPGAs and optical transceivers. So, performing data reduction on-chip
417  could potentially allow the chips to overcome these bottleneck and run at higher speeds.

418  Naturally, any circuitry for on-chip compression must be designed to not occupy excessive amounts of
419  space in the pixels or periphery. Additionally, readout chips are developed in technology scales that are
420  relatively large compared to commercial data processing hardware like GPUs and FPGAs, due to the high
421  cost of smaller nodes. So, rather than implementing general-purpose processing logic into detector chips,
422  specific algorithms are implemented.

423  Hammer et al. [13] present a design for on-chip data compression for photon counting detectors, using
424  techniques similar to those discussed previously in Section 3. Firstly, within the pixels, count values are
425  encoded with a varying step size, with the step size getting larger for larger pixel values such that the step
426  does not exceed the $\sqrt{(N)}$ Poisson noise. This is lossy compression, but the additional noise introduced by
427  this encoding should be less than the Poisson noise. Then, in the chip's periphery during readout, a lossless
428  compression scheme is applied that applies bit shuffling to the data (much as described for bitshuffle LZ4)
429  then encodes runs of zeroes efficiently. Applied to example datasets, this computationally-cheap approach
430  achieved a compression ratio of around 6 for XRD data, compared to 19 obtained with GZIP.

431  One important limitation of on-chip data compression is that in a large tiled detector composed of multiple
432  chips, the compressibility of the data can vary a great deal between different chips. In X-ray diffraction
433  experiments in particular, the X-ray intensity close to the beam is much higher than at large scattering
434  angles, leading to less compression. So, chips close to the beam may encounter data bottlenecks even if a
435  high overall level of compression is achieved for the detector as a whole.

436  Another example of on-chip reduction is data vetoing by rejecting bad images, as discussed previously.
437  The Sparkpix-ED chip is one of a family of chips with built-in data reduction being developed by SLAC
438  [42]. It is an integrating pixel detector with two key features. Firstly, each pixel has built-in memory which
439  is used to store recent images. Secondly, there is summing circuitry that can add together the signals in
440  groups of $3 \times 3$ pixels, to produce a low-resolution image with $1/9$ of the size. During operation, each
441  time a new image is acquired the detector will store the full-resolution image in memory and send out the
442  low-resolution image to the readout system. The readout system will then analyse the low-resolution images
443  to see if an interesting event occurred (e.g. diffraction from a protein crystal in a serial crystallography
444  experiment). If so, the detector can be triggered to read out the corresponding full-resolution image from
445  memory; if not, the image will be discarded. A small prototype has demonstrated reading out low-resolution
446  images at 1 MHz frame rate, and full-resolution images at 100 kHz.

## 6.2  On-detector processing and compression

448  Data processing and compression can potentially be done within the detector, before data is transferred to
449  the control system (typically over optical links). Once again, this can reduce the bandwidth required for
450  this data transfer.

451  Typically, on-detector electronics are needed both to control the detector's operation, and to perform
452  serialization and encoding of data into some standard format so that it can be sent efficiently back to the
453  control system and received using off-the-shelf hardware. (In some cases, it's also necessary to interface to
454  additional components such as on-board ADCs.) One particular benefit of serialization is that individual
455  transceivers on a readout chip typically run at a lower data rate than can be achieved by modern optical

456  links, so serialization can allow these links to be used more efficiently; for example, even the fastest on-chip
457  transceivers typically do not have data rates above 5 Gbit/s, whereas data sent using 100 Gigabit Ethernet
458  with QSFP28 transceivers consists of 4 channels with a data rate of 25 Gbit/s each. These tasks of control
459  and serialization are typically implemented in an FPGA, or a System-On-Chip (SoC) with an FPGA fabric.
460  So, the most common approach to on-detector data processing is to use the FPGA's processing resources.

461  As discussed in more detail earlier, FPGAs can provide highly-parallelized data processing, but firmware
462  development can be challenging and time consuming. So, detector-specific processing routines that will
463  always need to be performed on the data are better-suited to FPGA implementation than ones that vary a
464  lot between experiments.

465  For example, FPGAs in the EIGER photon counting detector perform count-rate correction and image
466  summing [43]. Since the counters in the pixel have a depth of 12 bits, this makes it possible for the detector
467  to acquire images with a depth of up to 32 bits by acquiring a series of images and internally summing
468  them, rather than needing to transfer many 12-bit images to the DAQ system.

469  One drawback of on-detector processing is that if the FPGA is used for control, serialization and data
470  processing, then it can become more difficult to change the data processing routines. When FPGA firmware
471  is compiled, the compiler will route together blocks in the FPGA to produced the desired functionality. So,
472  changing the data processing routines can change the routing of other functionality in the FPGA, potentially
473  affecting reliability. So, careful re-testing is required after changing the data processing. This is another
474  reason why on-detector processing with FPGAs is mostly used for fixed, detector-specific routines.

### 6.3   Data acquisition hardware such as PCs with accelerator cards

476  Data sent out from a detector module will normally be received by either one or more DAQ PCs, or more
477  specialized hardware, located either at the beamline or in the facility's computing centre. These parts of the
478  DAQ system can have a range of functions:

479  • Detector configuration and acquisition control, which requires both monitoring the detector's state and
480    data output, and receiving commands from the control system.

481  • Ensuring reliable data reception. It is common for a detector's output to be a continuous flow of data,
482    transferred by a simple protocol like UDP without the capability to re-send lost packets [44], So the
483    DAQ system is required to reliably receive and buffer this data; this typically requires, for example,
484    having dedicated, high performance network or receiver cards.

485  • Data correction and reduction.

486  • Transferring data to where it's needed, for example the facility's storage system, online processing
487    and/or a user interface for feedback. This can include tasks like adding metadata and file formatting.

488  Compared to detectors themselves, these systems tend to be built with relatively off-the-shelf hardware
489  components - for example, standard network cards or accelerator cards. The appeal of using off-the-shelf
490  hardware, in addition to reducing development costs, is that it is easier to upgrade the hardware to take
491  advantage of rapid improvements in technology. These DAQ systems may use conventional CPUs, hardware
492  accelerators such as GPUs and FPGAs, or a mixture of these.

493  One recent example of this approach is the Jungfraujoch processing system, developed by PSI for the
494  Jungfrau 4-megapixel detector [45]. (A similar approach is taken by the CITIUS detector [46].) The
495  Jungfraujoch system is based on an IBM IC922 server PC, equipped with accelerator cards, and can handle
496  17 GB/s data when the detector is running at 2 kHz frame rate.

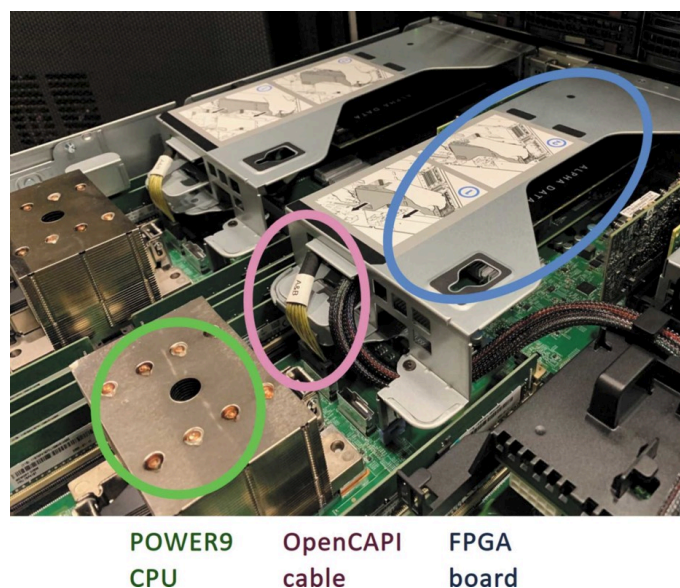POWER9   OpenCAPI   FPGA
CPU       cable      board

**Figure 10.** A photograph of the IC922 server (IBM) used in Jungfraujoch, showing the FPGA board used for data reception and processing, and the OpenCAPI link allowing coherent access to the CPU's memory. Reproduced from [45] with permission of the International Union of Crystallography.

The Jungfraujoch server is shown in Fig. 10 (reproduced from [45] with permission of the International Union of Crystallography). Firstly, it is equipped with two "smart network cards" from Alpha Data, each with a Xilinx Virtex Ultrascale+ FPGA and a 100 Gigabit Ethernet network link. These receive data directly from the detector over UDP. The FPGA then converts the raw data into images, as described in [47]. Jungfrau is an integrating detector with gain switching, so the raw data consists of ADC values, and the process of converting this to either photons or energy values includes subtracting the dark current and scaling by an appropriate gain factor. The FPGA also performs the Bitshuffle algorithm, which is the first step in Bitshuffle LZ4 compression described previously; since this algorithm involves reordering bits in a data stream, this is well-suited to implementation on FPGA. The FPGA is primarily programmed by using High Level Synthesis (HLS) based on C++.

The data from each FPGA is transferred to the host server PC using OpenCAPI interconnects, which both allows high speed data transfer at 25 GB/s, and makes it simpler for the FPGA to access memory on the server. The CPU in the server performs the LZ4 part of the Bitshuffle LZ4 lossless compression algorithm, then forwards this data to the file system using ZeroMQ so it can be written to the file system as an HDF5 file. There is also a GPU in the system that can perform tasks such as spot finding in macromolecular crystallography, and monitoring aspects of the detector's behaviour such as the dark current. (These tasks are well suited to GPUs, since they involve performing operations on full images in parallel.)

One promising technological development in this field is that FPGA and GPU vendors are increasingly developing accelerator cards with built-in network links for the datacenter market. For example, Xilinx have the Alveo line of FPGA cards with two 100 Gigabit Ethernet links and optional high-bandwidth memory, and likewise NVIDIA is currently developing GPU cards with network links. So, this will increase the availability of powerful, standardized hardware for building DAQ systems.

## 6.4  Computing clusters

Modern light sources are generally supported by large computer clusters for data processing and large storage systems [48]. Data from DAQ PCs or similar hardware at the beamlines can be transferred to them over the facilty's network. Data processing in computing clusters can be divided into online analysis, where the data is transferred directly to the cluster to provide fast results to the experiment, and offine analysis where data is read back from storage for processing at a later point. (Naturally, there may not be a sharp dividing line between these two approaches.) In addition to server PCs with CPUs, computing clusters can also incorporate GPUs or less commonly FPGAs.

In a computing cluster, there is generally a scheduling system that controls how tasks are assigned to the computers. This has the advantage of allowing sharing of resources between different experiments and users according to needs, whereas dedicated hardware installed at a beamline may be idle much of the time. Computing clusters are also much better-suited to allowing users to remotely access to their data and providing the software tools required to analyse it. However, this flexibility in task scheduling and usage can make it more difficult to ensure that we can reliably receive and process images from the detector at the required rate; this is one reason for having dedicated DAQ computers at the beamline for receiving detector data.

## 7  CONCLUSIONS

The increasing data rates of detectors for photon science mean there is a strong need for high-speed detector data correction, and data reduction.

There are strong ties between these tasks. On the one hand, data reduction can often yield better results on properly-corrected data, since the correction process can reduce spurious variation in images (e.g. pixel-to-pixel variation in response) and make it easier to exploit redundancy in the data. Conversely, after lossy data reduction it is impossible to perfectly recover the original data, so it is crucial to ensure that the quality of the data correction is as good as possible.

Both of these tasks can benefit from making better use of hardware accelerators such as GPUs and FPGAs for highly parallel processing. The increasing use of accelerators in other areas, such as datacenters, means that we can take advantage of improvements in both their hardware and in tools for programming them. In particular, as FPGAs and GPUs with built-in network links become available "off the shelf", this increases the potential for different labs to build their DAQ systems with compatible hardware, and share the algorithms they develop. Although this paper emphasizes data processing hardware, it is also important to note that well-designed and coded algorithms can deliver much better performance.

Data reduction is a growing field in photon science. To date, lossless compression has mainly been used, since this ensures that there is no loss in data quality, and in some experiments lossless methods can achieve impressive compression ratios. However, lossless compression is relatively ineffective for methods such as imaging, and as data volumes increase there is demand for even greater data reduction in diffraction experiments. So, there will be an increasing need for lossy compression and other methods of reduction. In doing so, it is crucial to use appropriate metrics to test that the data reduction does not significantly reduce the quality of the final analysis. By using well-established methods of lossy compression, for example image compression with JPEG2000, it is easier to incorporate data reduction into existing data analysis pipelines. However, novel methods of data reduction tailored to photon science experiments have the potential for better performance.

559 Data reduction can be incorporated into different stages of a detector's DAQ chain; generally, performing
560 data reduction earlier in the chain is more challenging and less flexible, but has the advantage of reducing
561 the bandwidth requirements of later stages, and can enable greater detector performance by overcoming
562 bottlenecks in bandwidth. The development of on-chip data reduction is a particularly exciting development
563 for enabling higher-speed detectors, though this would typically require the development of different chips
564 for different classes of experiment. As the demand for data reduction increases, we can expect detectors
565 and experiments to incorporate a series of data reduction steps, beginning with simpler or more generic
566 reduction early in the processing chain, and then more experiment-specific data reduction taking place in
567 computer clusters.

## CONFLICT OF INTEREST STATEMENT

568 The authors declare that the research was conducted in the absence of any commercial or financial
569 relationships that could be construed as a potential conflict of interest.

## AUTHOR CONTRIBUTIONS

570 DP: Writing – original draft, review & editing, Conceptualization, Investigation; VR: Writing – review &
571 editing, Visualization; HG: Funding acquisition, Supervision.

## FUNDING

## ACKNOWLEDGMENTS

## REFERENCES

581 [1] Rao R. Synchrotrons face a data deluge. *Physics Today* (2020). doi:10.1063/PT.6.2.20200925a.
582 [2] Marras A, Klujev A, Lange S, Laurus T, Pennicard D, Trunk U, et al. Development of CoRDIA: An
583     imaging detector for next-generation photon science X-ray sources. *Nuclear Instruments and Methods*
584     *in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*
585     **1047** (2023) 167814. doi:https://doi.org/10.1016/j.nima.2022.167814.
586 [3] Allahgholi A, Becker J, Bianco L, Delfs A, Dinapoli R, Goettlicher P, et al. AGIPD, a high dynamic
587     range fast detector for the European XFEL. *Journal of Instrumentation* **10** (2015) C01023. doi:10.
588     1088/1748-0221/10/01/C01023.
589 [4] Hsieh SS, Iniewski K. Improving paralysis compensation in photon counting detectors. *IEEE*
590     *Transactions on Medical Imaging* **40** (2021) 3–11. doi:10.1109/TMI.2020.3019461.

[5] Könnecke M, Akeroyd FA, Bernstein HJ, Brewster AS, Campbell SI, Clausen B, et al. The NeXus data format. *Journal of Applied Crystallography* **48** (2015) 301–305. doi:10.1107/S1600576714027575.

[6] [Dataset] The HDF Group. Hierarchical data format version 5 (2000-2010).

[7] Sayood K. *Lossless compression handbook* (Elsevier) (2002).

[8] Al-Shaykh OK, Mersereau RM. Lossy compression of noisy images. *IEEE Transactions on Image Processing* **7** (1998) 1641–1652.

[9] Becker J, Greiffenberg D, Trunk U, Shi X, Dinapoli R, Mozzanica A, et al. The single photon sensitivity of the Adaptive Gain Integrating Pixel Detector. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* **694** (2012) 82–90. doi:https://doi.org/10.1016/j.nima.2012.08.008.

[10] Ballabriga R, Alozy J, Campbell M, Frojdh E, Heijne E, Koenig T, et al. Review of hybrid pixel detector readout ASICs for spectroscopic X-ray imaging. *Journal of Instrumentation* **11** (2016) P01007. doi:10.1088/1748-0221/11/01/P01007.

[11] Broennimann C, Eikenberry EF, Henrich B, Horisberger R, Huelsen G, Pohl E, et al. The PILATUS 1M detector. *Journal of Synchrotron Radiation* **13** (2006) 120–130. doi:10.1107/S0909049505038665.

[12] Deutsch LP. *DEFLATE Compressed Data Format Specification version 1.3*. No. 1951 in Request for Comments (RFC Editor) (1996). doi:10.17487/RFC1951.

[13] Hammer M, Yoshii K, Miceli A. Strategies for on-chip digital data compression for X-ray pixel detectors. *Journal of Instrumentation* **16** (2021) P01025. doi:10.1088/1748-0221/16/01/P01025.

[14] Leonarski F, Mozzanica A, Brückner M, Lopez-Cuenca C, Redford S, Sala L, et al. JUNGFRAU detector for brighter x-ray sources: Solutions for IT and data science challenges in macromolecular crystallography. *Structural Dynamics* **7** (2020) 014305. doi:10.1063/1.5143480.

[15] Gailly J, Adler M. Gzip documentation and sources. *available as gzip-\*. tar in ftp://prep. ai. mit. edu/pub/gnu* (1993).

[16] Masui K, Amiri M, Connor L, Deng M, Fandino M, Höfer C, et al. A compression scheme for radio data in high performance computing. *Astronomy and Computing* **12** (2015) 181–190. doi:https://doi.org/10.1016/j.ascom.2015.07.002.

[17] Ziv J, Lempel A. A universal algorithm for sequential data compression. *IEEE Transactions on information theory* **23** (1977) 337–343.

[18] Huffman DA. A method for the construction of minimum-redundancy codes. *Proceedings of the IRE* **40** (1952) 1098–1101.

[19] Mittone A, Manakov I, Broche L, Jarnias C, Coan P, Bravin A. Characterization of a sCMOS-based high-resolution imaging system. *Journal of Synchrotron Radiation* **24** (2017) 1226–1236. doi:10.1107/S160057751701222X.

[20] Skodras A, Christopoulos C, Ebrahimi T. The JPEG 2000 still image compression standard. *IEEE Signal Processing Magazine* **18** (2001) 36–58. doi:10.1109/79.952804.

[21] Marone F, Vogel J, Stampanoni M. Impact of lossy compression of X-ray projections onto reconstructed tomographic slices. *Journal of Synchrotron Radiation* **27** (2020) 1326–1338. doi:10.1107/S1600577520007353.

[22] Shensa MJ, et al. The discrete wavelet transform: wedding the a trous and mallat algorithms. *IEEE Transactions on signal processing* **40** (1992) 2464–2482.

[23] Huang P, Du M, Hammer M, Miceli A, Jacobsen C. Fast digital lossy compression for X-ray ptychographic data. *Journal of Synchrotron Radiation* **28** (2021) 292–300. doi:10.1107/S1600577520013326.

[24] Di S, Cappello F. Fast error-bounded lossy HPC data compression with SZ. *2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS)* (2016), 730–739. doi:10.1109/IPDPS.2016. 11.

[25] Underwood R, Yoon C, Gok A, Di S, Cappello F. ROIBIN-SZ: Fast and science-preserving compression for serial crystallography (2022). doi:10.48550/arXiv.2206.11297.

[26] Barty A, Kirian RA, Maia FR, Hantke M, Yoon CH, White TA, et al. Cheetah: software for high-throughput reduction and analysis of serial femtosecond x-ray diffraction data. *Journal of applied crystallography* **47** (2014) 1118–1131.

[27] Winter G, Waterman DG, Parkhurst JM, Brewster AS, Gildea RJ, Gerstel M, et al. *DIALS*: implementation and evaluation of a new integration package. *Acta Crystallographica Section D* **74** (2018) 85–97. doi:10.1107/S2059798317017235.

[28] Rahmani V, Nawaz S, Pennicard D, Setty SPR, Graafsma H. Data reduction for X-ray serial crystallography using machine learning. *Journal of Applied Crystallography* **56** (2023) 200–213. doi:10.1107/S1600576722011748.

[29] Ke TW, Brewster AS, Yu SX, Ushizima D, Yang C, Sauter NK. A convolutional neural network-based screening tool for X-ray serial crystallography. *Journal of synchrotron radiation* **25** (2018) 655–670.

[30] Blaj G, Chang CE, Kenney CJ. Ultrafast processing of pixel detector data with machine learning frameworks. *AIP Conference Proceedings* (AIP Publishing) (2019), vol. 2054.

[31] Chen L, Xu K, Zheng X, Zhu Y, Jing Y. Image distillation based screening for x-ray crystallography diffraction images. *2021 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom)* (IEEE) (2021), 517–521.

[32] Kieffer J, Wright J. PyFAI: a python library for high performance azimuthal integration on GPU. *Powder Diffraction* **28** (2013) S339–S350. doi:10.1017/S0885715613000924.

[33] Matěj Z, Skovhede K, Johnsen C, Barczyk A, Weninger C, Salnikov A, et al. Azimuthal integration and crystallographic algorithms on field-programmable gate arrays. *Acta Crystallographica Section A* **77** (2021) C1185. doi:10.1107/S0108767321085263.

[34] Lin JM. Python non-uniform fast fourier transform (pynufft): An accelerated non-cartesian mri package on a heterogeneous platform (cpu/gpu). *Journal of Imaging* **4** (2018) 51.

[35] Qasaimeh M, Denolf K, Lo J, Vissers K, Zambreno J, Jones PH. Comparing energy efficiency of cpu, gpu and fpga implementations for vision kernels. *2019 IEEE international conference on embedded software and systems (ICESS)* (IEEE) (2019), 1–8.

[36] Becker D, Streit A. Real-time signal identification in big data streams bragg-spot localization in photon science. *2015 International Conference on High Performance Computing & Simulation (HPCS)* (IEEE) (2015), 611–616.

[37] Becker D, Streit A. A neural network based pre-selection of big data in photon science. *2014 IEEE Fourth International Conference on Big Data and Cloud Computing* (IEEE) (2014), 71–76.

[38] Souza A, Oliveira LB, Hollatz S, Feldman M, Olukotun K, Holton JM, et al. Deepfreak: learning crystallography diffraction patterns with automated machine learning. *arXiv preprint arXiv:1904.11834* (2019).

[39] Branco S, Ferreira AG, Cabral J. Machine learning in resource-scarce embedded systems, fpgas, and end-devices: A survey. *Electronics* **8** (2019) 1289.

[40] Abuowaimer Z, Maarouf D, Martin T, Foxcroft J, Gréwal G, Areibi S, et al. Gplace3. 0: Routability-driven analytic placer for ultrascale fpga architectures. *ACM Transactions on Design Automation of Electronic Systems (TODAES)* **23** (2018) 1–33.

680  [41] Hügging F, Owtscharenko N, Pohl DL, Wermes N, Ehrmann O, Fritzsch T, et al. Advanced through
681       silicon vias for hybrid pixel detector modules. *Nuclear Instruments and Methods in Physics Research*
682       *Section A: Accelerators, Spectrometers, Detectors and Associated Equipment* **936** (2019) 642–643.
683       doi:https://doi.org/10.1016/j.nima.2018.08.067. Frontier Detectors for Frontier Physics: 14th Pisa
684       Meeting on Advanced Detectors.
685  [42] Rota L, Perez AP, Habib A, Dragone A, Miceli A, Markovic B, et al. X-ray detectors for LCLS-II
686       with real-time information extraction: the SparkPix family (2023). 24th International Workshop On
687       Radiation Imaging Detectors (IWORID 2023), Oslo, Norway.
688  [43] Bruckner M, Bergamaschi A, Cartier S, Dinapoli R, Frojdh E, Greiffenberg D, et al. A multiple 10
689       Gbit Ethernet data transfer system for EIGER (2016). 20th IEEE Real Time Conference, Padova,
690       Italy.
691  [44] Gottlicher P, Sheviakov I, Zimmer M. 10G-Ethernet prototyping for 2-D X-Ray detectors at the XFEL.
692       *2009 16th IEEE-NPSS Real Time Conference* (2009), 434–437. doi:10.1109/RTC.2009.5321620.
693  [45] Leonarski F, Brückner M, Lopez-Cuenca C, Mozzanica A, Stadler HC, Matěj Z, et al. Jungfraujoch:
694       hardware-accelerated data-acquisition system for kilohertz pixel-array X-ray detectors. *Journal of*
695       *Synchrotron Radiation* **30** (2023) 227–234. doi:10.1107/S1600577522010268.
696  [46] Grimes M, Pauwels K, Schülli TU, Martin T, Fajardo P, Douissard PA, et al. Bragg coherent diffraction
697       imaging with the CITIUS charge-integrating detector. *Journal of Applied Crystallography* **56** (2023)
698       1032–1037. doi:10.1107/S1600576723004314.
699  [47] Redford S, Andrä M, Barten R, Bergamaschi A, Brückner M, Dinapoli R, et al. First full dynamic
700       range calibration of the JUNGFRAU photon detector. *Journal of Instrumentation* **13** (2018) C01027.
701       doi:10.1088/1748-0221/13/01/C01027.
702  [48] Thayer J, Damiani D, Ford C, Dubrovin M, Gaponenko I, O'Grady CP, et al. Data systems for the
703       Linac Coherent Light Source. *Advanced Structural and Chemical Imaging* **3** (2017) 3. doi:10.1186/
704       s40679-016-0037-7.

## FIGURE CAPTIONS