




A data reduction and compression description for high throughput time-resolved electron microscopy

Abhik Datta^{1,2}, Kian Fong Ng^{1,2}, Deepan Balakrishnan ^{1,2}, Melissa Ding³, See Wee Chee ^{1,2,4}, Yvonne Ban^{1,2}, Jian Shi^{1,2} & N. Duane Loh ^{1,2,4}✉

Fast, direct electron detectors have significantly improved the spatio-temporal resolution of electron microscopy movies. Preserving both spatial and temporal resolution in extended observations, however, requires storing prohibitively large amounts of data. Here, we describe an efficient and flexible data reduction and compression scheme (ReCoDe) that retains both spatial and temporal resolution by preserving individual electron events. Running ReCoDe on a workstation we demonstrate on-the-fly reduction and compression of raw data streaming off a detector at 3 GB/s, for hours of uninterrupted data collection. The output was 100-fold smaller than the raw data and saved directly onto network-attached storage drives over a 10 GbE connection. We discuss calibration techniques that support electron detection and counting (e.g., estimate electron backscattering rates, false positive rates, and data compressibility), and novel data analysis methods enabled by ReCoDe (e.g., recalibration of data post acquisition, and accurate estimation of coincidence loss).

¹Centre for Bioluminescence Sciences, National University of Singapore, Singapore, Singapore. ²Department of Biological Sciences, National University of Singapore, Singapore, Singapore. ³Department of Computer Science and Engineering, Ohio State University, Columbus, OH, USA. ⁴Department of Physics, National University of Singapore, Singapore, Singapore. ✉email: duaneloh@nus.edu.sg

Fast, back-thinned direct electron detectors are rapidly transforming electron microscopy. These detectors ushered in a “resolution revolution” for electron cryo-microscopy (cryo-EM), and the prospect of seeing sub-millisecond dynamics for in-situ electron microscopy. These transformations are driven by three key factors: (1) improved detection efficiency, (2) shorter detector readout times to better resolve individual electron events, and (3) algorithms that translate these advances into improved spatial and temporal resolution. Whereas the first two factors have received considerable attention, it remains impractical for many existing algorithms to process the very large raw output produced by these movie-mode detectors. Fortunately, the useful information on these raw data are typically sparse, hence a suitable data reduction and compression scheme should allow us to fully reap the advantages offered by these detectors.

Nearly all the useful information in a single raw detector image is contained within “secondary electron puddles”, each of which is digitized from the cloud of secondary charged particles formed in the wake of individual high energy electrons passing through the detector’s sensor. While the size and shape of secondary electron puddles contain some information¹, localizing the entry point of the incident electron from its electron cloud already noticeably improves the spatial resolution of the image. To accurately localize these electron puddles they must be spatio-temporally well separated (by increasing the frame rate or reducing the incident electron flux), thereby reducing the so-called coincidence loss². This separation creates a very high raw data load when acquiring images that add up to a desired accumulated electron dose. For example, the memory needed to store the incident electron entry points, in a low coincidence-loss image (~6%) acquired at 0.01 e/pixel/frame is approximately a hundredth that of the raw detector readout, with the remainder holding only thermal readout noise.

Currently, there are three popular options to manage the large raw data loads that a high-throughput electron detector generates. First, which is typical in cryo-EM, is to employ a higher internal frame rate on the detector for counting electrons at low coincidence loss, but add many of these frames together before they are stored to disk. The downside here is the loss of temporal resolution in the added, stored images. The second option is to reduce the total data acquisition time. Here, an experimenter may fill terabytes of local hard disk with raw data for 10 min, then wait at least twice as long to offload this data to a larger networked drive before more data acquisition can proceed. The third option is to collect data at the maximum detector frame rate but only store the frames that contain significant information. However, this strategy only works at high dose rates where individual pre-selected frames still show sufficient contrast for the experimenter to judge whether to keep or discard them. At such high dose rates, the experimenter has to either sacrifice spatial resolution or be limited to atomic resolution only for radiation-hard samples.

None of these three options are ideal, especially since the vast majority of these high data loads are storing only the detector’s thermal and readout noise. Furthermore, these options also limit us from using faster detectors³ to study dynamics at even shorter timescales. Naturally, reducing and compressing the raw data would obviate the need to choose between these three compromising options. If we stored only electron arrival events, we can enjoy high temporal and spatial resolution, while continuously acquiring movies of dose-sensitive samples at very low dose rates for practically hours, uninterrupted.

While more experiments that require both high temporal and spatial resolution are emerging^{4–6}, acquiring such movies for long time scales remains expensive, and in many cases, infeasible. For perspective, a 4 TB hard drive only accommodates about 21 min of data collection at a DE-16 detector’s maximum data rate of

3.08 GB/s. A prominent example that exploits fast detectors is motion-correction in TEM (transmission electron microscopy). Here, the imaging resolution is demonstrably improved when fast detectors fractionate the total electron dose on a sample onto a time series of micrographs that are individually corrected for relative dose-induced motion⁷. In fact, recent work suggests that using more efficient data representation to further increase dose fractionation, hence finer time resolution, can improve spatial resolution⁸.

As electron microscopy becomes increasingly reliant on larger datasets and more complex processing and analysis workflows, there is an ever-greater push for publications to include raw data necessary for others to validate and reproduce the analyses⁹. Improved compression will make public archives like EMPIAR more accessible and increase their adoption, and encourage deposition of raw micrographs facilitating validation of the structures produced using them¹⁰. Without an effective data reduction and compression scheme, storing raw detector data will be costly: at ~US\$20 per terabyte (TB) of archival storage on commodity HDDs and ~US\$400 per TB on SSDs (based on the prices of lower end external hard disk drives and solid-state drives, as of April 2019^{11,12}, just 15 min of continuous data acquisition per day on the DE-16 (Direct Electron, LP) detector at its maximum frame rate (3 GB/s throughput) will cost between US\$ 20,000 to US\$ 400,000 per year, respectively.

Here, we propose a data reduction and compression scheme capable of file size reductions that are as high as 100× for realistic electron-counting scenarios. The output of this scheme is a file format known as ReCoDe (Reduced Compressed Description). For simplicity, we refer to the reduction compression scheme as the ReCoDe scheme (or simply ReCoDe when the context is clear). In this scheme, the original raw data is first reduced to keep only the information regarding identified electron puddles, which are then further compressed. The ReCoDe scheme permits four reduction levels and several different types of lossless compression algorithms, whose combinations are discussed in this work. We show how data stored in the least lossy ReCoDe reduction level can be re-processed post-acquisition to remove detector artifacts, especially those owing to slow drifts in the thermal background. Moreover, storing data at this reduction level retains the puddle shape and intensity information. Through several use cases, we show the benefits of retaining this information. One of these is coincidence loss estimation, where we show that puddle shape information is essential for accurate estimation. We also develop methods for estimating the prevalence of back scattered electrons and estimating the false positive rates of electron events using this information. For the DE-16 detector we estimated the ratio of primary to backscattered electrons to be ~8.6. The ReCoDe scheme is sufficiently parallelizable such that data streams from even the fastest current detectors can be reduced and compressed on-the-fly onto networked storage disks using only modest computing resources, provided the raw data can be accessed before it is written to disk. For instance, the raw data stream of a low dose experiment (0.8 e/pixel/s) collected on a DE-16 detector (~3.08 GB/s throughput) can be reduced, compressed by 10 Intel Xeon CPU cores, then written to network-attached storage devices via a modest 10 gigabit ethernet connection. Furthermore, the ReCoDe data format has been designed for fast sequential and random access, so that frames can be decompressed into the reduced representation on demand. Using ReCoDe in-situ electron microscopy movies can retain the high dose fractionation and sub-millisecond time resolution while extending acquisition time from minutes to hours. Giving users the flexibility to fractionate their doses from tens to thousands of frames per second for more precise temporal resolution and drift-correction where possible.

Finally, using several publicly available EMPIAR⁹ datasets with moderate to high dose rates, ranging from 0.5 to 5.0 electrons/pixel/frame, we show that ReCoDe achieves 2–8× compression, outperforming existing compression approaches.

Results

Data reduction levels. A secondary electron puddle is characterized by its spatial location, its two-dimensional (2D) shape, and the pixel intensities within this shape. To accommodate various downstream processing needs we define four logical data reduction levels: L1, L2, L3, and L4, with progressively higher levels of file size reduction and hence information loss from L1 to L4 (Fig. 1).

All four data reduction schemes begin with a thresholding step, which produces a binary map identifying pixels as containing useful signal or not. The ADU (analog-digital unit) threshold used to label signal pixels are independent for each pixel and decided based on the signal-noise calibration procedure (discussed below). In ReCoDe level L1, the sparsified signal pixel intensities are then bit-packed into a dense format. Bit packing removes unused bits and converts the list of ADU values into a continuous string of bits. The binary map and the bit packed intensity values are independently compressed and the two compressed data are stacked to create an L1 reduced compressed frame. As L1 reduction retains all the information about electron puddles, electron counting can be performed long afterward, should the user wish to explore different counting algorithms and/or parameters. Both thresholding and packing are sufficiently fast to make L1 suitable for on-the-fly processing (discussed in “Demonstration of on-the-fly Reduction and Compression” section). Even for relatively high electron flux data (0.05 e/pixel/frame) L1 reduction alone achieves a 10× reduction in file size. This reduced data can be further compressed to achieve an overall 25× file size reduction (see below).

In L3 reduction, the pixel intensities are discarded during thresholding and only the binary map is retained and compressed. L3 is therefore optimized for speed, at the expense of puddle specific ADU (analog-digital unit, or pixel intensity) information.

To compute puddle specific features, in L2 and L4 reductions, the clusters of connected pixels that constitute individual puddles are identified from the binary map using a connected components labeling algorithm, discussed in the Methods section. In L4 reduction, each puddle in the binary map is further reduced to the single pixel, where the primary electron was likely incident. L4 reduction, therefore, results in a highly sparse binary map that is optimized for maximum compression. At the same electron flux (0.05 e/pixel/frame) L4 reduction and compression results in 45× file size reduction. This increased compression comes at the cost of throughput since counting has to be performed as part of the reduction step.

In L2 reduction, a summary statistic, such as mean, maximum or sum of ADU (analog-digital unit), is extracted for each electron puddle. Preliminary studies suggest that such information may correlate with whether a measured electron was elastically or inelastically scattered¹. The sparse puddle features are then packed into a dense format and the binary map and the dense puddle features are independently compressed. Several applications that record diffraction patterns benefit from a high dynamic range but do not necessarily need to retain the entire signal as done in L1. L2 is designed for such applications.

In L1 and L2 reductions, the binary maps and the packed intensity summary statistics are independently compressed and then stacked. As the binary maps and intensity values have very different characteristics, compressing them independently results in optimal compression (Fig. 1b–d).

The reduced compressed data formats are detailed in Supplementary Method 1.

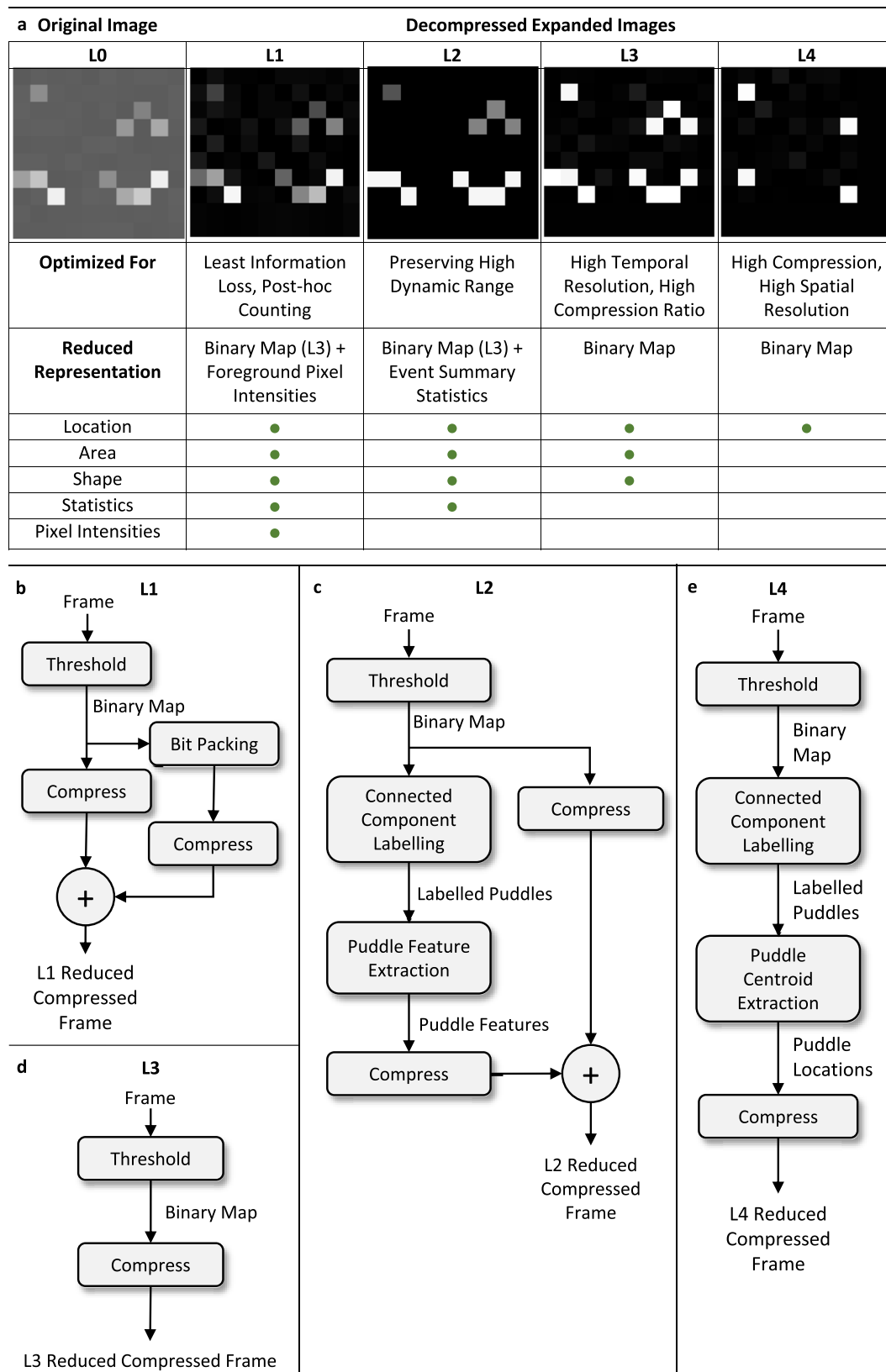
All four data reduction schemes in ReCoDe first reduce the data by removing primarily readout noise (thresholding) and then compressing the signal. Accurate signal-noise separation is therefore critical. To remove pixel-level differences in dark noise and gain that can bias the identification of isolated electron puddles, individual thresholds are calculated per pixel based on calibration data (see “Methods” section). For the DE-16 detector, this calibration can be done with a single dataset with flat-field illumination at a low dose rate and extended exposure times. Since different detectors may require custom calibration, ReCoDe only requires the per pixel thresholds for separating signal-noise as input and is agnostic of the calibration method used. These thresholds are specified in a single image frame, which is reloaded by ReCoDe at user-specified intervals. External programs can update the thresholds intermittently for on-the-fly recalibration to accommodate changing detector response.

Calibrating parameters for data reduction. An appropriate threshold separating signal from noise is critical for electron counting to be effective. Typically, this threshold is established through calibration, based on dark and gain references obtained during data acquisition. In most imaging software these calibrations depend on several hyper-parameters that are predetermined (for instance, the number of frames used in the dark reference). Once the calibrated frames are reduced to electron counted images, the calibration cannot be revised, and the effects of the hyper-parameters are permanent. The L1 reduction presents an alternative, where the data can be recalibrated post-acquisition without having to store the entire dataset, as long as a sufficiently permissive threshold is used. In low dose rate experiments, during data acquisition, the quality of images cannot be verified through visual inspection. The effectiveness of the calibration can, therefore, be difficult to judge. The ability to recalibrate datasets in such cases can significantly improve image quality, as shown in Fig. 2. Here, the data was recalibrated by using a higher threshold for separating dark noise and signal and pixel gains were recalculated after removing single pixel puddles (see Fine calibration in “Methods” section for details). We observed that such recalibration can significantly reduce the number of false positive electron events.

Even small deviations in calibration can significantly bias counting and therefore recalibration (or at least a quality assessment) should be a necessary step in ensuring accurate counting. L1 reduced data facilitates such post-hoc analysis. This includes using the electron puddle size/shape distributions to estimate realistic coincidence losses specific to the detector and imaging conditions (Table 1).

Table 1 shows coincidence losses estimated using five different techniques. In the first three (columns from left to right) puddles are assumed to be of fixed shape and size, whereas, in the last two, the actual puddle shape and size information are included in the calculation (see “Methods” section). Clearly, the knowledge of puddle shape and size is essential for accurate coincidence loss estimation. Therefore, accurately estimating coincidence loss requires retaining data at reduction levels L1–L3.

A recent study⁸ has proposed storing L4 reduced data in a sparse format to benefit from higher dose fractionation without overwhelming acquisition systems with storage requirements. To achieve super-resolution electron counting, which is critical for improving reconstruction resolution in cryo-EM, they propose subdividing each pixel before counting and storing the higher-resolution spatial locations of electron events using a higher bit-depth. ReCoDe’s L1 reduction scheme enables super-resolution electron counting without the need to subdivide pixels at the time



of acquisition, thus eliminating the need to predetermine to what extent pixels should be partitioned.

Reducibility and compressibility with increasing electron fluxes. With increasing electron flux, the data naturally becomes

less reducible and less compressible. To quantify this change, we simulated images at eight electron fluxes between 0.0025 to 0.07 e/pixel/frame (Fig. 3). This range was chosen for tolerable coincidence loss during electron counting (Table 1). For data without any reduction (unreduced compression line in Fig. 3), the compression ratio remains similar across all fluxes (~4×), because

Fig. 1 Data reduction levels and scheme. **a** The leftmost image (L0) depicts a 10×10 pixel image (the raw detector output) with four secondary electron puddles. The remaining four images from left to right correspond to the four data reduction levels, L1 to L4, respectively. Each image represents a reconstruction of the original image (L0) using only the information retained at that level (see table at the bottom). The L1 image retains all the useful information about the secondary puddles by first removing detector readout/thermal noise from L0. In L2, the spatial location of the four puddles, the number of pixels (area) in each puddle, the shape of the four puddles and an intensity summary statistic (sum, maximum or mean) for each puddle are retained. Each reduction level offers different advantages in terms of speed, compression, information loss, spatial or temporal resolution, etc (see row labeled “Optimized For”). The row labeled “Reduced Representation” describes how the information retained at each level is packed in the reduced format. These packings are tuned to provide a good balance between reduction speed and compressibility. In L3, the puddle area, shape and location information are all encoded in a single binary image, which is easily computed and highly compressible. These three aspects in L1 and L2 are packed as the binary image used in L3. Only the most likely locations of incident electrons are saved as binary maps in L4. Panels **b**, **c**, **d**, and **e** are the reduction compression pipelines for reduction levels L1, L2, L3, and L4, respectively. Here, the thresholding step produces a binary map identifying pixels as signal or noise. Bit packing removes unused bits and converts the list of ADU values into a continuous string of bits. The connected components labeling algorithm identifies clusters of connected pixels that constitute individual electron puddles from this binary map. Puddle centroid extraction further reduces each puddle to a single representative pixel; and puddle feature extraction computes puddle specific features such as mean or maximum ADU.

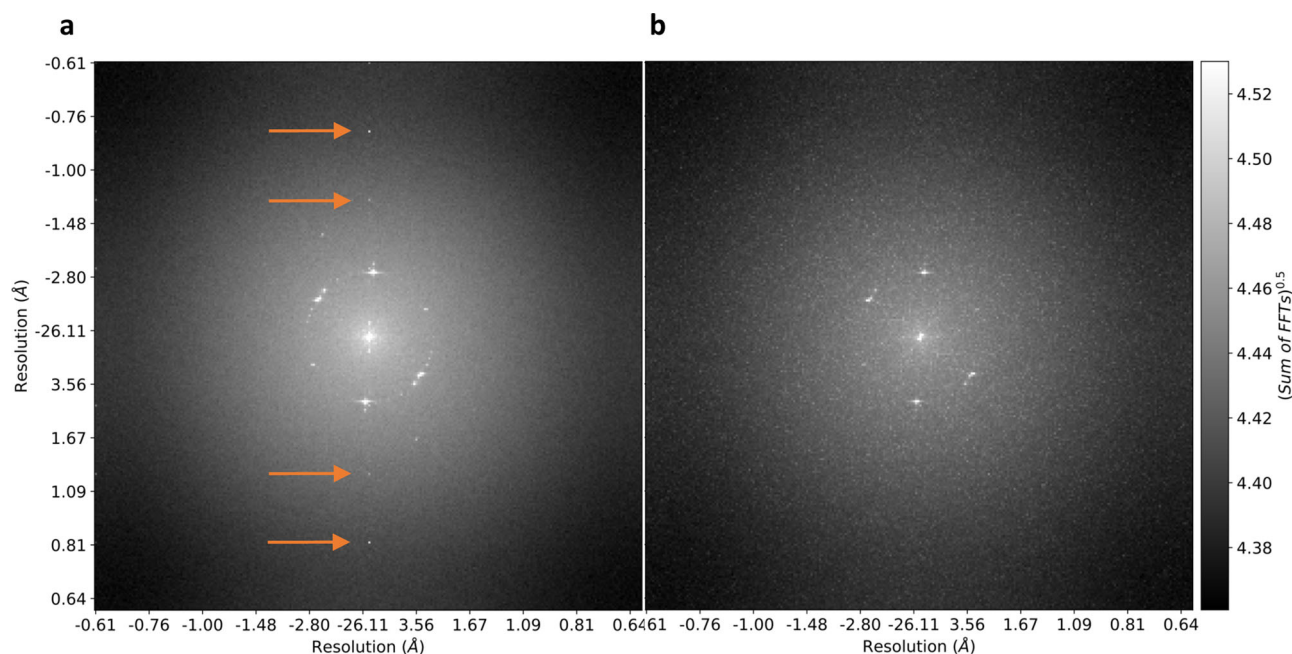


Fig. 2 Recalibration of L1 reduced data to remove artifacts. Panels **a** and **b** are Fourier transforms (FT) of summed L1 reduced frames of HRTEM movies of a molybdenum disulfide 2-D crystal, acquired using a JEOL 2200 microscope operating at 200 keV and a DE-16 detector running at 300 fps, with a pixel resolution of 0.2 \AA . **(a)** is L1 reduced using fast on-the-fly calibration using a 3 \AA threshold (see “Methods” section) **(b)** is the result of recalibrating **(a)** with a more stringent fine calibration that uses an area threshold and a 4 \AA threshold (see “Methods” section). The Fourier peaks indicated with orange arrows in **a** are due to detector artifacts, which are not readily visible in the image but can severely impact drift correction. **a** and **b** are the sum of FFTs of 9000 frames.

of readout dark noise. L3 and L4 reduced data are essentially binary images with 1-bit per pixel (Fig. 1). Therefore, if the input data uses n bits to represent each pixel’s intensity, a factor of n reduction is achieved using L3 or L4 reduction alone. In Fig. 3, a $16\times$ reduction is seen for the 16-bit simulated data. In L1 and L2, pixel intensity information and event summary statistics are retained in addition to the L3 binary map. As electron flux increases, more pixel intensities/event statistics need to be stored. However, due to coincidence loss the number of counted electron events and L1 and L2 file sizes increase only sub-linearly.

With increasing electron flux the binary images used to store location and shape information in the reduced format, also become less compressible. This is evident from the L3 and L4 “reduction + compression” lines in Fig. 3. At the same time, for L1 reduction, the proportion of reduced data containing pixel intensities increases rapidly with increasing electron flux. As a result, the compressibility of L1 reduced data falls very quickly with increasing electron flux.

At moderate ($0.01 \text{ e/pixel/frame}$) and low ($0.001 \text{ e/pixel/frame}$) electron flux L1 reduction compression results in $60\times$ and $170\times$ data reduction, respectively.

Reduction L4, where only puddle locations are retained, is optimized for maximum compression and can achieve reduction compression ratios as high as $45\times$, $100\times$, and $250\times$ at high ($0.05 \text{ e/pixel/frame}$), moderate ($0.01 \text{ e/pixel/frame}$) and low ($0.001 \text{ e/pixel/frame}$) electron flux, respectively.

Compression algorithms exploit the same basic idea: by representing the more frequently occurring symbols with fewer bits the total number of bits needed to encode a dataset is effectively reduced. Consequently, data is more compressible when symbols are sparsely distributed. Such sparse distributions are readily present in the back-thinned DE-16 electron detector, where nearly 80% of the digitized secondary electron puddles span fewer than three pixels (Supplementary Fig. 8). Even for puddles that span four pixels (of which there are 110 possibilities) nearly half (48.3%) are the 2×2 -pixel square motif.

Table 1 Coincidence loss estimation methods.

Dose rate (e^- /pixel/frame)	Coincidence loss (Fraction of e^- events lost)				
	Simulated with fixed size and shape			Analytically computed using size distribution	Simulated using size and shape distributions
	Assuming 3 × 3 pixel PSF	Assuming 2 × 2 pixel PSF	Assuming 1 pixel PSF		
0.0025	0.060	0.031	0.011	0.0140	0.0156
0.005	0.117	0.061	0.022	0.0354	0.0317
0.01	0.725	0.119	0.044	0.0784	0.0619
0.02	0.407	0.227	0.087	0.1524	0.1207
0.03	0.556	0.324	0.128	0.2115	0.1767
0.04	0.676	0.413	0.167	0.2590	0.2296
0.05	0.772	0.492	0.205	0.2975	0.2801
0.06	0.846	0.563	0.242	0.3288	0.3272
0.07	0.902	0.626	0.278	0.3544	0.3726
0.08	0.942	0.684	0.312	0.3753	0.4155
0.09	0.969	0.734	0.345	0.3923	0.4558
0.1	0.983	0.779	0.376	0.4061	0.4940

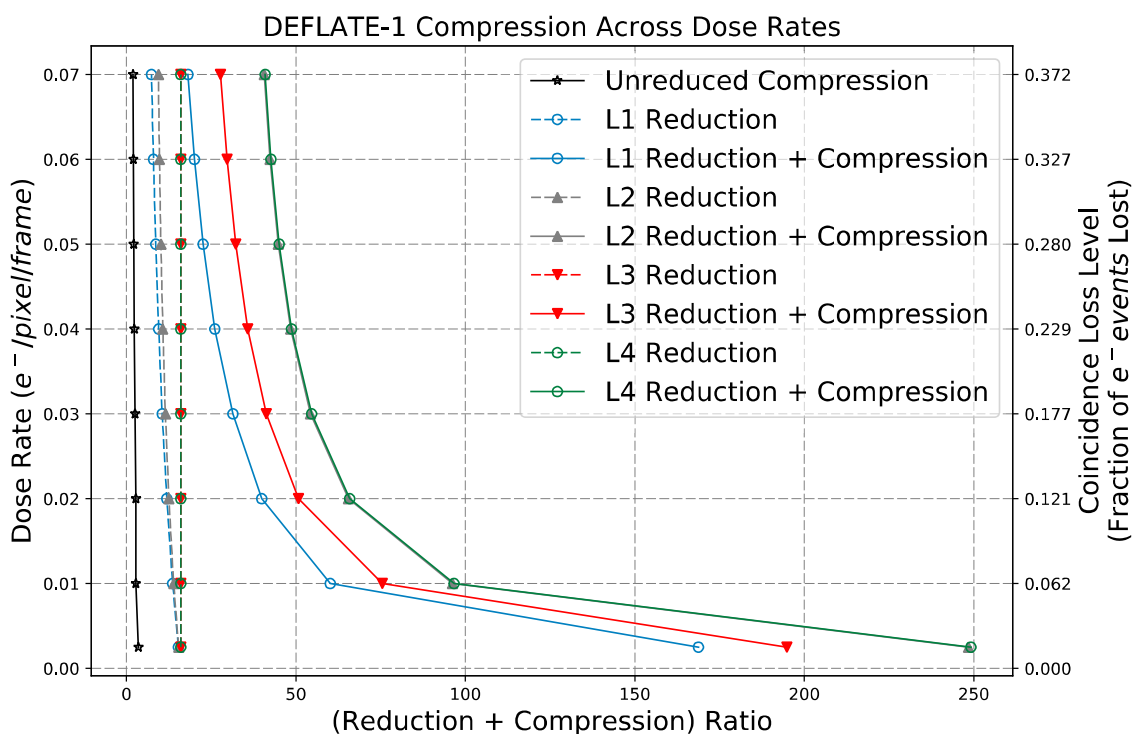


Fig. 3 Reducibility and compressibility of data with increasing electron flux. The solid black line (“unreduced compression”) shows the compression ratios achieved on unreduced raw data (including dark noise) using Deflate-1. The dashed lines show the compression ratios achieved with just the four levels of data reduction and without any compression. The solid lines show the compression ratios after compressing the reduced data using Deflate-1. The coincidence loss levels corresponding to the electron fluxes label the second y-axis on the right.

The randomly distributed centroids of secondary electron puddles account for the largest fraction of memory needed to store the reduced frames. We considered three representations for storing these centroids and ultimately adopted a binary image representation (Methods section).

Compression algorithms. Any lossless compression algorithm can operate on the reduced data levels in Fig. 1. Compression algorithms are either optimized for compression power or for compression speed, and the desired balance depends on the application. For on-the-fly compression, a faster algorithm is preferable even if it sub-optimally compresses the data, whereas an archival application may prefer higher compression power at the expense of compression speed.

We evaluated the compression powers and speeds of six popular compression algorithms that are included by default in the ReCoDe package: Deflate¹³, Zstandard (Zstd), bzip2 (Bzip)^{14,15}, LZMA¹⁶, LZ4¹⁷, and Snappy¹⁸ (Fig. 4). Each algorithm offers different advantages; bzip, for instance, is optimized for compression power whereas Snappy is optimized for compression and decompression speed. All five algorithms can be further parameterized to favor compression speed or power. We evaluated the two extreme internal optimization levels of these algorithms: fastest but sub-optimal compression, and slowest but optimal compression.

Data reduction schemes similar to L1 have been previously used to compress astronomical radio data in Masui et al.¹⁹. They proposed the bitshuffle algorithm for compressing radio data after removing thermal noise from it. We experimented with

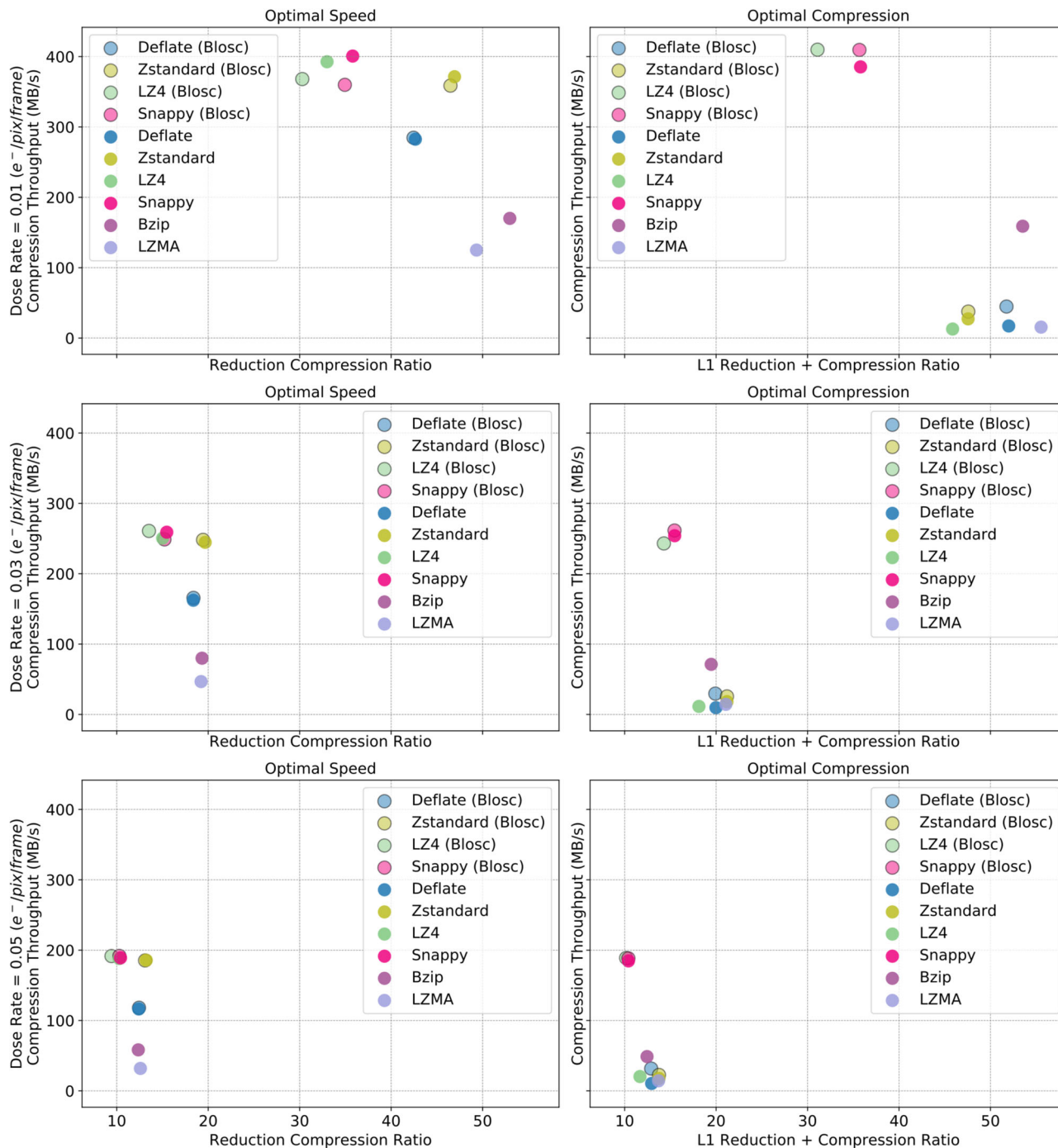


Fig. 4 Comparison of compression algorithms with L1 reduction at three dose rates. Each scatter plot shows the reduction compression ratios and the compression throughputs of six compression algorithms (Deflate, Zstandard (Zstd), bzip2 (Bzip), LZ4, LZMA, and SNAPPY), plus the BloSC variants of Deflate, Zstandard (Zstd), LZ4, and SNAPPY. Reduction compression ratio (horizontal axes in all panels) is the ratio between the raw (uncompressed) data and the reduced compressed data sizes. The three rows of scatter plots correspond to three different electron fluxes: 0.01, 0.03, and 0.05 e/pixel/frame, from top to bottom. The left and right columns of scatter plots correspond to the two most extreme internal optimization levels of the compression algorithms: fastest but suboptimal compression labeled “Optimal Speed” (left column), and optimal but slow compression labeled “Optimal Compression” (right column). The data throughputs (vertical axes in all panels) are based on single threaded operation of ReCoDe and include the time taken for both reduction and compression. The decompression throughputs of the six algorithms are presented in Supplementary Fig. 2.

BloSC, a meta-compressor for binary data, that implements bitshuffle in addition to breaking the data into chunks that fit into the system’s L1 cache, to improve compression throughputs (Fig. 4).

LZ4 and SNAPPY have the highest throughputs across all reduction levels and electron fluxes, with reduction compression

ratios slightly worse than the remaining five algorithms. At the lowest dose rate (0.01 e/pixel/frame) Bzip results in the best reduction compression ratios, regardless of the internal optimization level. At higher dose rates (0.03 and 0.05 e/pixel/frame) Zstd has the highest compression ratio. Considering all dose rates and internal optimization levels, Zstd on average offers the best

balance between compression ratio and throughput. The choice of internal optimization level only marginally affects the reduction compression level but significantly improves throughput.

Deflate optimized for speed for instance is almost $\sim 25\times$ faster than Deflate optimized for compression, across the three dose rates. We use Deflate optimized for speed (referred to as Deflate-1) as the reference compression algorithm for the rest of the paper, as it represents a good average case performance among all the compression algorithms. In subsequent sections, we will show that Deflate-1 is fast enough for on-the-fly compression. All algorithms have higher decompression throughput than compression throughput (Supplementary Fig. 2). Deflate-1 has ten times higher decompression throughput than compression throughput, which means the same computing hardware for reduction and compression can support on-the-fly retrieval and decompression of frames for downstream data processing.

For most compression algorithms Blosc marginally improves compression throughput (Fig. 4), except in the case of optimal compression with LZ4, where Blosc improves throughput by as much as 400 MB/s.

Demonstration of on-the-fly reduction and compression.

Electron microscopy imaging often has to be performed at low electron flux to reduce beam induced damage, if the sample is dose sensitive, as well as to minimize beam induced reactions. Observing rare events or slow reactions in such cases require extended acquisition, that is not feasible with current detector software without compromising temporal resolution. Loss of temporal resolution, in turn, degrades drift correction and therefore limits spatial resolution. ReCoDe's on-the-fly reduction compression fills this critical gap, enabling hours long continuous acquisition without overwhelming storage requirements, compromising on temporal resolution, or losing puddle information.

ReCoDe is easily parallelized, with multiple threads independently reducing and compressing different frames in a frame stack. In this multithreaded scheme, each thread reduces and compresses the data to an intermediate file, which are merged when data collection is complete. The merging algorithm reads from the intermediate files and writes to the merged file sequentially and is therefore extremely fast. The intermediate and merged (ReCoDe) file structures and the merging process are described in Supplementary Method 1.

With this multithreaded scheme, ReCoDe can achieve throughputs matching that of the detectors enabling on-the-fly reduction and compression. In addition, intermediate files can be accessed sequentially in both forward and reverse directions, with frames indexed by frame number, time stamp, and optionally scan position. Owing to the small size of the reduced compressed frames, they can be read from intermediate files by external programs for live processing and feedback during acquisition even without merging them back into a single file. Users also have the option of retaining raw (unreduced and uncompressed) frames at specified intervals for validation or for on-the-fly recalibration. In electron microscopy facilities data is often archived in high capacity network-attached storage (NAS) servers. A schematic of this on-the-fly reduction compression pipeline is shown in Fig. 5a. We evaluated the feasibility of directly collecting the reduced-compressed data onto NAS servers, to avoid the overhead of transferring data after collecting it on the microscope's local computer.

With the DE-16 detector running at 400 fps, at a dose rate of 0.001 e/pixel/frame and ReCoDe using 10 CPU cores of the acquisition computer that shipped with the DE-16 detector, we continuously captured data directly onto NAS servers connected

by a 10 gigabits/s Ethernet (10 GbE) connection, for 90 min (see "Methods" subsection: On-the-fly Compression Pipeline). To further evaluate this multithreaded scheme, we simulated a series of on-the-fly data reduction and compression at different electron fluxes. The implementation used for these simulations emulates the worst-case write performance of ReCoDe, where a single thread sequentially accesses the disk (see Supplementary Discussion 3 for details). At relatively low electron flux (0.01 e/pixel/s) we are able to achieve throughputs as high as 8.3 gigabytes per second (GB/s, Fig. 5b) using 50 threads on a 28 core system. At the same dose rate, to keep up with the DE-16 detector (which has a throughput of ~ 3.08 GB/s) only 10 CPU cores are sufficient. For perspective, another popular direct electron detector, the K2-IS (Gatan Inc.), nominally outputs bit-packed binary files at approximately 2.2 GB/s. However, since we did not have to incur extra computation time to unpack bits on the raw data from DE-16, the DE-16 benchmarks on Fig. 5 will not directly apply to K2-IS data.

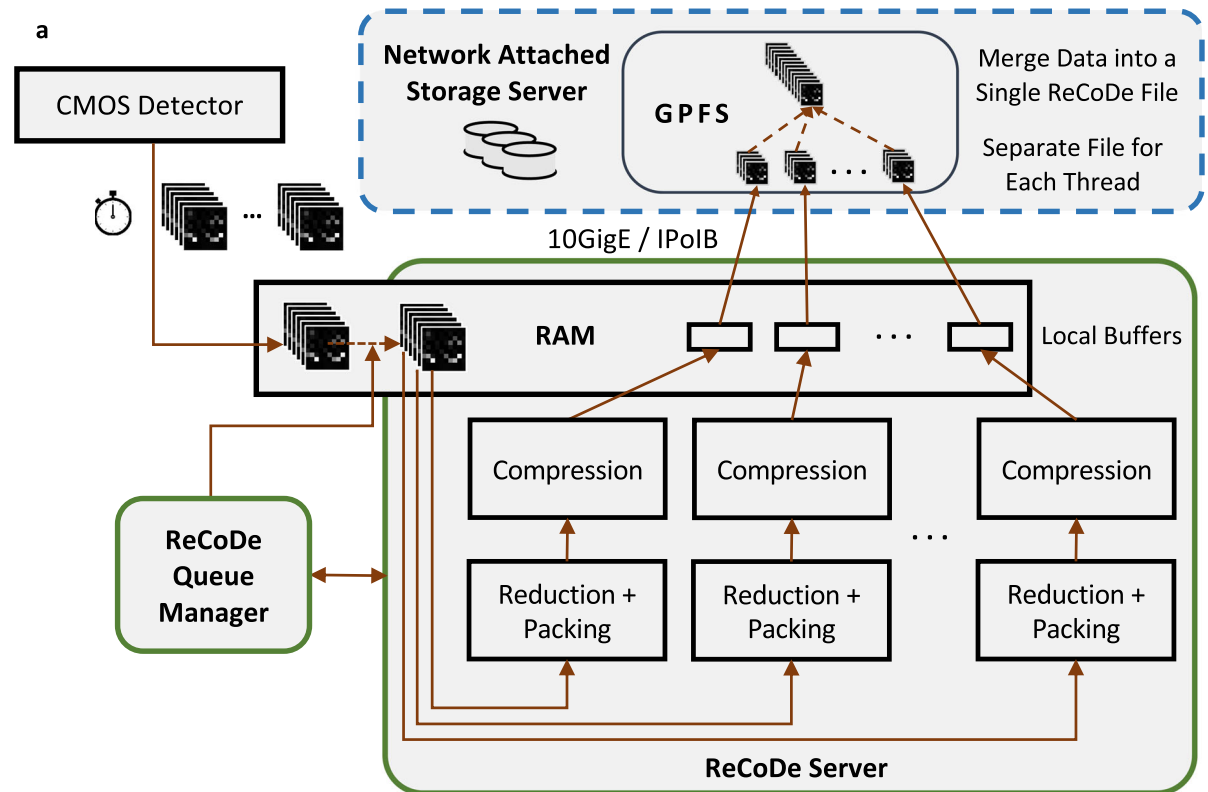
At moderate electron flux, writing directly to GPFS NAS servers using both 10 GbE and IPoIB (Internet Protocol over InfiniBand) has comparable throughputs to that of collecting data locally on the microscope's computer (Fig. 5c,d). However, at very low electron flux writing directly to the NAS server with IPoIB has slightly higher throughput. This is likely due to the reduced communication overhead per call in IPoIB and the distributed data access (IBM GPFS) supported by NAS servers, both of which are optimized to handle multiple simultaneous small write requests. In the absence of such a parallel data access ReCoDe still executes at close to 89% parallel (Supplementary Discussion 3).

Both the reduction and compression steps are essential for high throughput on-the-fly processing. Without compression, the reduced data is still too large to write over 10 GbE, particularly at moderate electron flux (Fig. 5e). Without reduction, the data is not compressible enough; the throughput of Deflate-1 compression without any data reduction (Fig. 5f) is abysmally low even when using 50 threads.

Effects of reduction on counted image quality. In many applications, the L2 and L3 reduced data has to be ultimately reduced to L4 (electron-counted image). Here, we consider how the information lost in L2 and L3 reductions affect the resolution in L4 images. In L4 puddles are reduced to a single pixel, which ideally contains the entry point of the incident electron. However, there is no clear consensus on the best approximation strategy for determining an electron's entry point, given a secondary electron puddle¹. The three common strategies are, to reduce the puddle to (1) the pixel that has the maximum intensity, (2) the pixel intensity weighted centroid (center of mass) or (3) the unweighted centroid of the puddle. Unlike L1 reduction, where all the information needed for counting with any of these strategies are retained, with L2 and L3 reductions, the pixel intensity information is either partially or completely lost. The puddles can then only be reduced to the unweighted centroid of the puddle using the third strategy. With L4 reduction, the approximation strategy has to be chosen prior to data acquisition. To evaluate how this information loss affects image quality we performed a knife-edge test using a beam blanker (see "Methods" section for implementation details). The results show (Supplementary Fig. 4) that the choice of approximation strategy, and therefore the choice of reduction level, has little consequence on image resolution.

Discussion

Studying millisecond in-situ dynamics with TEM, such as surface-mediated nanoparticle diffusion in water²⁰, requires us to



b L1 Reduction + Compression + Write to Local Drive

No. of Threads	6.9	20.0	39.9	48.8
50	13.42	8.31	4.63	2.67
45	13.40	8.14	4.49	2.65
40	13.36	8.19	4.54	2.65
35	13.26	8.22	4.45	2.65
30	13.00	8.18	4.43	2.60
25	11.84	7.27	3.97	2.32
20	9.83	5.89	3.15	1.88
15	7.50	4.54	2.54	1.43
10	4.93	2.95	1.72	0.97
5	2.53	1.50	0.85	0.51

c L1 Reduction + Compression + Write to Network Drive using 10GigE

No. of Threads	6.9	20.0	39.9	48.8
50	12.58	7.81	4.11	2.54
45	12.45	7.84	4.19	2.58
40	12.75	7.90	4.27	2.58
35	12.45	7.82	4.22	2.53
30	12.12	7.58	4.02	2.41
25	11.15	6.76	3.68	2.18
20	9.46	5.67	3.08	1.79
15	7.15	4.37	2.41	1.38
10	4.89	2.97	1.64	0.96
5	2.44	1.48	0.86	0.50

d L1 Reduction + Compression + Write to Network Drive using IPoIB

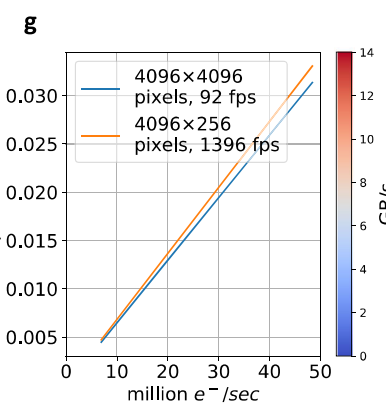
No. of Threads	6.9	20.0	39.9	48.8
50	13.45	8.30	4.50	2.70
45	13.19	8.15	4.51	2.63
40	13.03	7.99	4.45	2.63
35	12.86	8.07	4.44	2.65
30	13.25	8.07	4.44	2.63
25	11.92	7.33	4.12	2.37
20	9.92	5.93	3.20	1.88
15	7.59	4.50	2.52	1.47
10	5.09	3.04	1.70	1.00
5	2.54	1.54	0.89	0.52

e L1 Reduction + Write to Network Drive using IPoIB

No. of Threads	6.9	20.0	39.9	48.8
50	5.67	4.54	3.02	1.67
45	5.64	4.58	2.92	1.72
40	5.78	4.55	3.04	1.73
35	5.80	4.54	2.96	1.76
30	5.88	4.59	3.05	1.75
25	5.84	4.62	3.02	1.76
20	5.93	4.70	3.00	1.77
15	5.83	4.58	2.87	1.63
10	4.73	3.78	2.30	1.33
5	2.98	2.23	1.31	0.72

f Compression of Unreduced Data + Write to Network Drive using IPoIB

No. of Threads	6.9	20.0	39.9	48.8
50	0.76	0.76	0.72	0.76
45	0.78	0.76	0.78	0.72
40	0.80	0.78	0.79	0.76
35	0.75	0.75	0.76	0.71
30	0.73	0.72	0.70	0.73
25	0.65	0.65	0.64	0.63
20	0.58	0.57	0.58	0.54
15	0.49	0.46	0.50	0.45
10	0.37	0.38	0.39	0.34
5	0.23	0.21	0.21	0.21



operate at the maximum frame rates of these detectors. In addition, longer total acquisition times would be beneficial for studying reactions such as spontaneous nucleation²¹ where the experimenter systematically searches a large surface for samples. Several pixelated TEM electron detectors are now able to achieve sub-millisecond temporal resolutions, with the downside that the

local buffer storage accessible to these detectors fills up very quickly. Figure 6 shows that current TEM detectors running at maximum frame rates produce 1 TB of data in several minutes. When the temporal resolution is critical for an imaging modality, reducing the frame rate is not an option. An example is fast operando electron tomography²². To capture how the 3D

Fig. 5 Pipeline and data throughput of on-the-fly reduction and compression. **a** ReCoDe’s multithreaded reduction compression pipeline used for live data acquisition. The CMOS detector writes data into the RAM-disk in timed chunks, which the ReCoDe server processes onto local buffers and then moves to NAS servers. The ReCoDe Queue Manager synchronizes interactions between the ReCoDe server and the detector. **b** L1 reduction and compression throughput (GB/s) of Deflate-1, with multiple cores at four electron fluxes. The throughput of ReCoDe depends only on the number of electron events every second, hence the four dose rates (horizontal axis) are labeled in million electrons/second. The simulations were performed on a 28-core system, as a result, throughput scales non-linearly when using more than 28 cores (Supplementary Fig. 3). **c, d** Show throughputs when using 10 GbE and IPoB connections to write directly to NAS, respectively. In **e**, throughput of L1 reduction without any compression; **(f)** throughput of Deflate-1 when compressing the unreduced raw data. **g** Shows the conversion between million e/s and e/pixel/frame for two different frame size-frame rate configurations of the DE-16 detector.

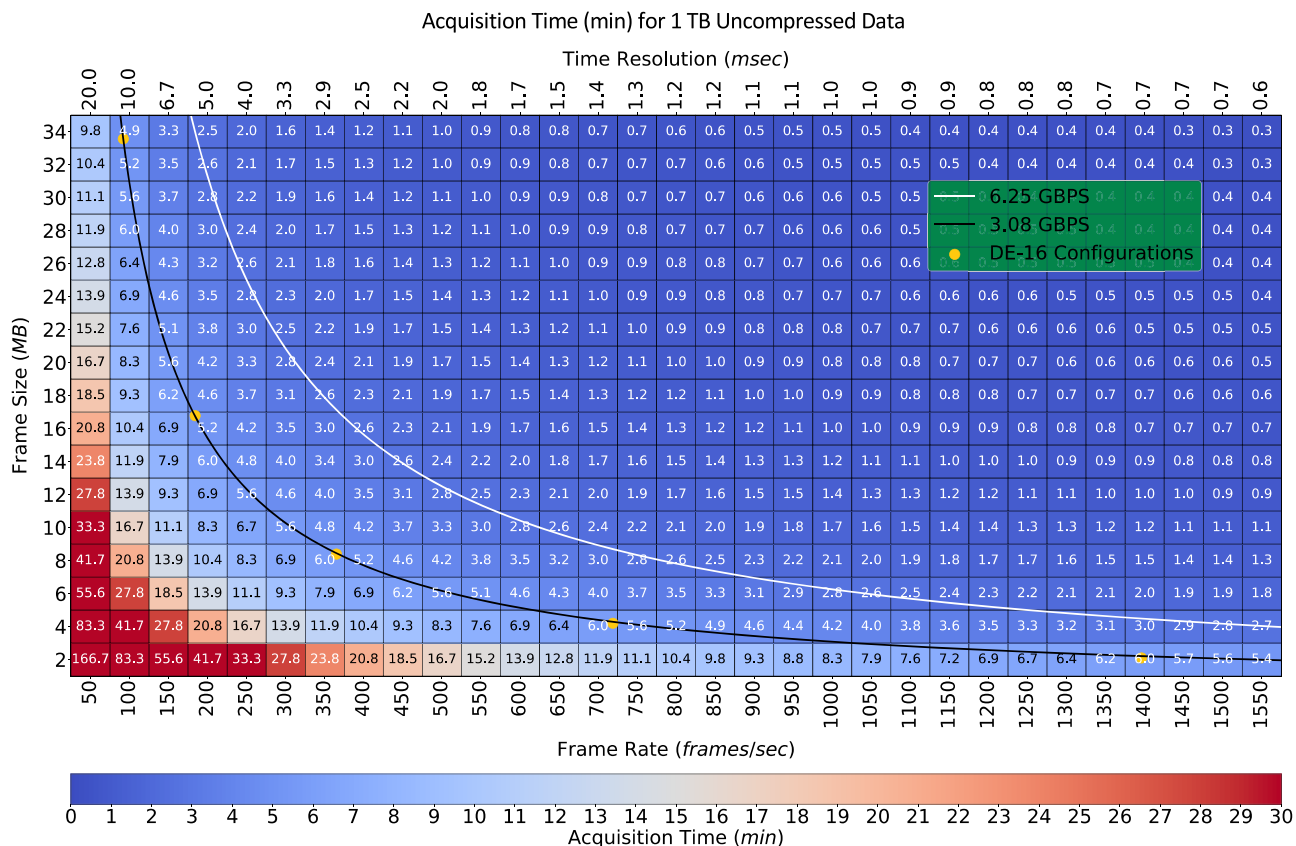


Fig. 6 Maximum data acquisition time of 1 TB of movie-mode TEM without data reduction and compression. Each cell’s horizontal and vertical grid position marks the temporal resolution (or, equivalently, frame rate) and frame size of a hypothetical movie-mode data acquisition scenario, respectively. A cell’s text and color indicates the time taken to acquire one terabyte (TB) of data at that frame size and temporal resolution without reduction and compression. For larger frames and high temporal resolution (top right corner), acquisitions lasting merely tens of seconds already produce 1 TB of data. With a 95× reduction in data size the same experiment can span 20 times longer, enabling the observation of millisecond dynamics in reactions that span several minutes. The yellow dots show a few of the frame size-frame rate combinations available for the DE-16 detector.

morphology of an object evolves over several seconds, a full-tilt series of the object has to be rapidly acquired at the detector’s peak frame rate. Here again, the duration of these observations can be significantly extended by substantially reducing the output data load with ReCoDe.

4D scanning transmission electron microscopy (4D STEM) techniques including Ptychography were used to image weak phase objects and beam sensitive samples such as Metal Oxide Frameworks (MOFs)²³. Here, a converged electron probe raster scans a sample collecting 2D diffraction patterns at each scan point. Although these experiments can produce hundreds of gigabytes of data in minutes²⁴, the diffraction patterns tend to be sparse outside of the central diffraction spot. As noise-robust STEM-Ptychography becomes a reality²⁵, their convergent beam electron diffraction patterns will be even sparser. ReCoDe level L1 reduction and compression, which preserves the patterns’

dynamic range while removing only dark noise, are likely to be useful for such data. Once the large datasets in 4D STEM are reduced they will readily fit into the RAM of desktop workstations, which also facilitates sparse and efficient implementations of processing algorithms.

Electron beam-induced damage is a major limitation for all cryo-EM modalities. In single-particle analysis (SPA) the energy deposited by inelastically scattered electrons manifests as sample damage and ice drift, where global and site-specific sample damage is detectable even at exposures as low as 0.1 e/Å²²⁶. Here higher electron dose fractionation improves resolution in two ways: (1) by reducing coincidence loss and thereby improving detection efficiency²⁷ and (2) by enabling more accurate estimation of sample drift at a higher temporal resolution⁸. Increasing detector frame rates can reduce the average displacement of each particle captured in each dose-fractionated frame,

but doing so further inflates the already large amounts of movie-mode data collected (see Supplementary Fig. 6). On-the-fly reduction and compression can significantly reduce the storage costs of movie-mode data, to accommodate image correction algorithms that operate at a degree of dose fractionation that is higher than current practice.

The recently proposed compressed MRCZ format²⁸ and ReCoDe offer complementary strategies to reduce file sizes generated by electron detectors. MRCZ is ideal for compressing information-dense images of electron counts integrated over longer acquisition times. ReCoDe, however, excels in reducing and compressing the much sparser raw detector data that are used to produce the integrated images typically meant for MRCZ. By doing so ReCoDe can preserve the arrival times of incident electrons that are lost when they are integrated into a single frame. Applying an MRCZ-like scheme on the raw un-reduced signal is inefficient, as shown with the “Unreduced Compression” line in Fig. 3. Figure 7a compares compression ratios obtained by MRCZ and ReCoDe on publicly available EMPIAR datasets from multiple published results^{29–32}, spanning a range of dose rates and detectors as listed in Table 2. Figure 7b shows the compression ratios achieved by the two approaches on simulated images. Across the range of dose rates ReCoDe produces better compression ratios on the EMPIAR datasets. As low dose rate datasets (below 0.58 electrons/pixel/frame) are likely to be sparse, ReCoDe as expected, achieves higher compression ratios than MRCZ (Fig. 7b). However, surprisingly, ReCoDe outperforms MRCZ even for some datasets with much higher average dose rates (EMPIAR-10346 in Fig. 7a). These datasets have particularly high contrast resulting in higher average dose rates but are still very sparse (Supplementary Fig. 7c), making them well suited for compression with ReCoDe.

We have described three novel analysis methods that demonstrate the necessity of reduction levels L1–L3. These methods cannot be applied on counted (L4 reduced) data, as they rely on puddle shape and intensity information. The first is the recalibration of L1 reduced data post acquisition, to improve counting accuracy. The second analysis uses puddle shape information to accurately estimate coincidence loss. When counting electrons, coincidence loss adversely affects spatial resolution (Supplementary Fig. 5). However, as we have shown, estimates of coincidence loss from the counted data can be inaccurate (Table 1). As reduction levels L1–L3 retain puddle shape information these can be used when accurate coincidence loss estimates are desired. In the third analysis we use a series of L1 reduced data sets with diminishing dose rates and extremely sparse electron events to estimate false positive rates of detecting electron events (Supplementary Note 12).

We also describe a novel method for estimating the proportion of backscattered electrons, using counted (L4 reduced) data (Supplementary Note 13). Using this analysis we estimated the proportion of primary to backscattered electrons for the DE-16 detector is ~8.6. In the future, it may be possible to even classify and eliminate backscattered electrons based on their sizes, shapes and proximity to primary electrons. Development of such techniques requires retaining more information than is currently done using counted data. The L1–L3 reduction levels in ReCoDe are designed to facilitate such future developments.

In summary, we present the ReCoDe data reduction and compression framework for high-throughput electron-counting detectors, which comprises interchangeable components that can be easily configured to meet application-specific requirements. ReCoDe supports four data reduction levels to balance application-specific needs for information preservation and processing speed. We further tested three electron localization strategies, and show that they produce similar spatial resolutions even

when the electron puddle intensity information is absent. By comparing five candidate compression algorithms on reduced electron data, we found that although LZ4 is the fastest, Deflate-1 offers the best compromise between speed and compressibility.

Remarkably, we demonstrated on-the-fly data reduction and compression with ReCoDe on the DE-16 detector for 90 min. Using only a desktop workstation, we continuously converted a 3 GB/s raw input data stream into a ~200 MB/s output that was, in turn, streamed onto networked drives via 10 Gbit ethernet. Crucially, this demonstration showed that on-the-fly data reduction and compression at low dose rates on our fastest S/TEM detectors is not compute-limited if the detector’s raw data stream is accessible (via a RAM-disk) before it is stored to SSDs. Even higher throughputs will be achievable with direct in-memory access to this raw data stream without the need for a RAM-disk. In the absence of fast simultaneous read-write, there is a critical lack of feedback in low dose rate, long time experiments. The experimenter is left blind in such situations, as individual frames do not have sufficient contrast and the frames available on disk cannot be read to produce a summed image with sufficient contrast. On-the-fly data reduction and compression with ReCoDe enables continuous feedback without interrupting data acquisition for hours.

The ReCoDe scheme can dramatically increase the throughput of electron microscopy experiments. Furthermore, the quality of observations for electron microscopy experiments can also improve. In cryo-EM, ReCoDe can support movies of higher frame rates, which can lead to better drift correction and lower coincidence loss. For in-situ experiments, higher frame rates can also improve the temporal and spatial resolution of the imaged samples.

Currently, a clear barrier for commercial vendors to produce higher throughput detectors is that users cannot afford to store the increased raw data that these faster detectors will bring. Going forward, the readout rates of CMOS detectors may increase to their internal megahertz clock rates³³, or even into the gigahertz regime³⁴. This uptrend is troubling if one considers, by default, that a detector’s raw data output rate increases linearly with its readout rate. However, because the ReCoDe format has very little storage overhead per frame, in principle, its processing and storage rate only scales with the total electron dose when the detector readout rate is fast enough to resolve individual electron puddles. Consequently, the ReCoDe output rate will not increase substantially with megahertz frame rates when the total electron dose is held constant. By efficiently reducing raw data into compact representations, ReCoDe prepares us for an exciting future of megahertz electron detectors in three crucial ways: it limits the storage costs of electron microscopy experiments, facilitates much longer data acquisition experimental runs, and very efficient processing algorithms that only compute on the essential features. More broadly, making ReCoDe open source encourages its own development by the community and incentivizes commercial vendors to specialize in much-needed hardware innovation. The full impact of electron counting detectors, quite possibly, is still ahead of us.

Methods

Data acquisition. All experimental data were collected on a DE-16 detector (Direct Electron Inc., USA) installed on the JEM-2200FS microscope (JEOL Inc., Tokyo, Japan) equipped with a field emission gun and operating at 200 keV accelerating voltage. StreamPix (Norpix Inc., Montreal, Canada) acquisition software was used to save the data in sequence file format without any internal compression. Data for puddle shape and size analysis (Supplementary Fig. 8) and MTF characterization with the knife-edge method (Fig. 5) were collected at 690 frames per second and 400 frames per second respectively with an electron flux of ~0.8 e/pixel/s.

All simulations of on-the-fly data collection were performed on a 28-core (14 core × 2 chips) system with 2.6 GHz E5-2690v4 Intel Broadwell Xeon processors and 512 GB DDR4-2400 RAM.

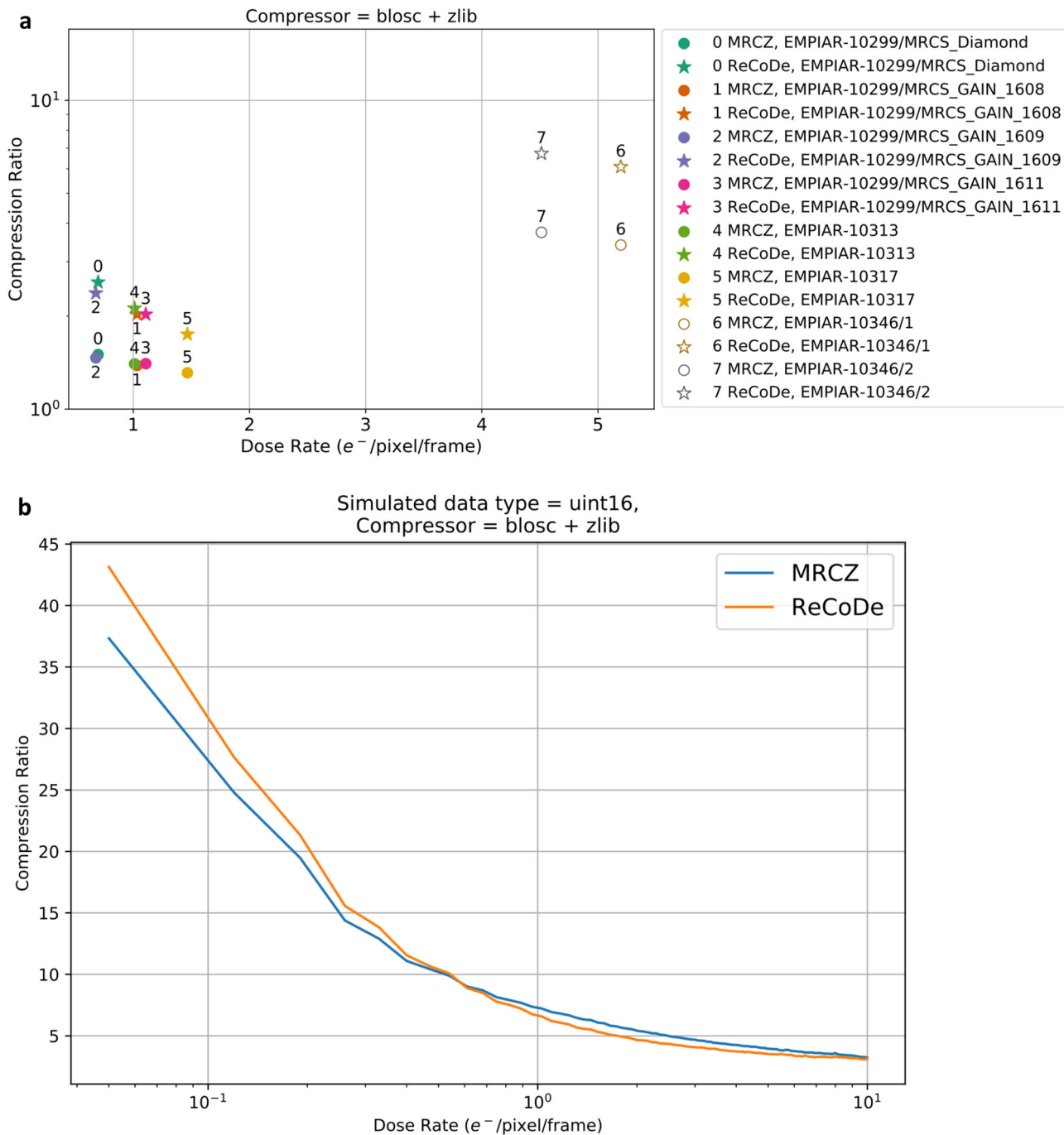


Fig. 7 Comparison of ReCoDe and MRCZ for archival datasets in EMPIAR. **a** Shows that the compression ratios obtained by ReCoDe (filled stars) on relatively low dose rate EMPIAR datasets are higher than those due to MRCZ (filled circles). **b** Compression ratios obtained using MRCZ and ReCoDe on simulated 16-bit unsigned integer data. The crossover point for performance occurs at 0.58 electron/pix/frame. At dose rates below this ReCoDe achieves higher compression ratios than MRCZ, whereas at dose rates above this MRCZ achieves slightly higher compression ratios. The number of electron events per pixel follows a Poisson distribution in these simulated datasets. The underlying compression algorithms used in **a** and **b** is Blosc + Deflate (zlib) for both MRCZ and ReCoDe. Table 2 lists a short description of the seven EMPIAR dataset used to generate **(a)**. Overall in the simulated data, for both compression algorithms, compression ratios reduce as dose rate increases, as expected. However, for the EMPIAR datasets, there are two groups, one for the floating-point data (datasets 0–5) and another for integer data (datasets 6 and 7). Although the floating-point data have lower dose rates than the integer type data, the former is less compressible because they are naturally less sparse than the latter. Nevertheless, within each group, the expected trend (reduction in compression ratio with increasing dose rate) holds true and ReCoDe outperforms MRCZ. A comparison where all the datasets are standardized to the same integer data type, presented in Supplementary Fig. 6, shows that the results from EMPIAR datasets and simulated data are quite similar.

Table 2 EMPIAR datasets used to compare ReCoDe and MRCZ in Fig. 7.

No.	EMPIAR ID/dataset	Dose rate (e/pixel/s)	Detector	Image size	Data type	Ref.
0	10299/MRCS_Diamond	0.69	K2 Summit	7676 × 7420	Float 32-bit	29
1	10299/MRCS_GAIN_1608	1.03	Falcon II	7676 × 7420	Float 32-bit	
2	10299/MRCS_GAIN_1609	0.67	Falcon II	7676 × 7420	Float 32-bit	
3	10299/MRCS_GAIN_1611	1.107	Falcon II	7676 × 7420	Float 32-bit	
4	10313	1.01	K2 Summit	3838 × 3710	Float 32-bit	30
5	10317	1.46	K2 Summit	3838 × 3710	Float 32-bit	31
6	10346	5.19	K3	11520 × 8184	Unsigned byte	32
7	10346	4.51	K3	11520 × 8184	Unsigned byte	

Connected components labeling. To compute the features specific to each electron puddle (e.g., centroids in L4 and the user-chosen summary statistics (ADU sum or maximum) in L2), the set of connected pixels (components) that constitute individual puddles have to be identified from the thresholded image. This connected components labeling can be computationally expensive for large puddles. Fortunately, puddle sizes tend to be small for most back-thinned direct electron detectors. For the DE-16 detector, 90% of the puddles are fewer than five pixels in size (Supplementary Fig. 8). Therefore, we use a greedy approach similar to the watershed segmentation algorithm³⁵ to perform connected components labeling. The algorithm assigns unique labels to each connected component or puddle, and the pixels associated with a given puddle are identified by the label of that puddle. In L2 reduction, these labels are used to extract the chosen summary statistics from the puddle and in L4 reduction, these labels are used to approximate the secondary electron puddle to a single pixel, by computing the centroid or center of mass, etc.).

Representation of puddle centroids. The randomly distributed centroids of secondary electron puddles account for the largest fraction of memory needed to store the reduced frames. We considered three representations for storing these centroids. In the first representation, a centroid is encoded as a single 2n-bit linear index. In the second representation, these linear indices are sorted and run-length encoded (RLE) since the ordering of centroids in a single frame is inconsequential. In the third representation, the centroids are encoded as a binary image (similar to L4). The RLE and binary image representations were found to be much more compressible than linear indices (Supplementary Fig. 9). Ultimately, the binary image representation was adopted in ReCoDe because the sorting needed for RLE is computationally expensive, for only a marginally higher compression.

Signal-noise calibration. In the current implementation, ReCoDe requires as input a single set of pre-computed calibration parameters, comprising each pixel's dark and gain corrected threshold for separating signal and noise at that pixel. Any calibration method can be used to compute this calibration frame. The "On-the-fly Calibration" methods subsection below describes a fast routine, for estimating this calibration frame, which we applied to the DE-16 detector. The "Fine Calibration" methods subsection thereafter details a more deliberate data collection approach, where additional diagnostics on the detector are also measured. Both calibration approaches yield practically similar results at dose rates above 10^{-4} e/pix/frame (Supplementary Fig. 14).

On-the-fly calibration. First, a flat-field illuminated dataset, comprising many raw detector frames preferably at the same low dose rate targeted for actual imaging afterward, is collected. An estimate of the incident dose rate is then computed using this dataset. Whereas this could be obtained with an independent measurement (e.g., Faraday cup), the dose rate computed by the procedure described here factors in the detector's detective quantum efficiency. Ideally, a pixel's intensity across the calibration frames would follow a mixture of two well-separated normal distributions, corresponding to either dark noise or signal. However, in practice, because the detector PSF is larger than a pixel, charge sharing from fast electrons incident on neighboring pixels will contribute to a single pixel's intensity, which causes the noise and signal distributions to overlap severely (Supplementary Fig. 10).

The calibration (summarized in Supplementary Note 11) begins by first estimating a single global threshold that separates signal from noise for all pixels. Assuming the histogram of dark values are normally distributed, this global threshold is estimated based on a user-specified upper limit on the tolerable false positive rate of a surmised normal distribution (r). However, because individual pixels behave differently from each other, using the same threshold for all pixels can severely bias electron counting. To remove this bias the global threshold has to be adapted for each pixel individually based on the pixel's gain and dark noise level.

Now we are ready to estimate the effective detectable electron count on the detector from the dataset directly. Given the low dose rate in this calibration dataset, only in a small fraction of frames does an individual pixel see electron events. Therefore, a pixel's median across all calibration frames is effectively its dark noise level, at this dose rate. Given the compact PSF and high SNR of DE-16

detectors, to calculate each pixel's gain, we assume that direct electron hits result in larger intensities than those due to charge sharing, even when the pixels have different gains. If the calibration dataset has a total dose of N e/pixel, where N is sufficiently small such that the probability of two electrons hitting the same pixel is negligible, then a pixel's gain is the median of the N largest intensities it has across all calibration frames. Therefore, we first estimate the total dose per pixel in the calibration dataset using a few randomly selected small two-dimensional (2D) patches. Separate thresholds are identified for individual pixels in these patches in a similar manner to the global threshold (i.e., assuming normality in the dark distribution and using a false positive rate parameter r). These thresholds are used to identify the connected components in each selected 2D patch across all frames in the calibration dataset. The number of connected components emanating from the central pixel of a 2D patch across all calibration frames gives an estimate of the number of electron events (n_c) at the central pixel of that patch. The average of these values across all randomly selected patches (\bar{n}_c) is used as the estimated total dose per pixel in the calibration dataset. Here, a puddle is assumed to emanate from the pixel that has the maximum value in the puddle. Finally, using the per-pixel dark noise levels and gains the global threshold is adapted to compute each pixel's independent threshold. To compute a pixel's threshold the global threshold is first shifted such that the pixel's dark noise level matches the global mean dark noise level and then scaled such that the pixel's gain matches the global gain.

For sufficiently sparse calibration data, even mean pixel intensity, which is much more efficiently computed than median, can be used to estimate the dark noise level for the pixel, although at the expense of a slightly higher false positive rate.

Fine calibration. To further assess the fast, on-the-fly calibration, a slower and more intricate calibration, referred to as fine calibration, was also implemented. The fine calibration adds two steps to the on-the-fly calibration procedure described above, a common-mode correction and a puddle area based filtering. The common-mode correction eliminates dynamically fluctuating biases in electron counting due to correlated thermal fluctuations between pixels that are connected to the same local voltage (hence thermal) source; the area threshold filters electron puddles to reduce false positive puddle detection. Analysis of the temporal response of DE-16's pixels revealed local detector regions of size 4×256 pixels that have correlated responses (Supplementary Note 12 and Supplementary Fig. 12). The fine calibration steps are summarized in Supplementary Fig. 13.

A series of datasets with increasingly sparse data was used to compare the two calibrations. The controlled reduction in dose rate was achieved by increasing magnification in successive datasets by 200× while keeping the electron flux constant. A comparison of the rate of decay of the estimated dose rates with counting following the two calibration strategies revealed that on-the-fly calibration includes a substantial number of false positive puddles (Supplementary Fig. 14A-B). However, this can be easily remedied with area-based filtering following L1 reduced data acquisition. While the common-mode correction has only a marginal effect on counting with the DE-16 detector (Supplementary Fig. 14C), it might be essential for other detectors.

Backthinning of direct electron detectors was instrumental in reducing noise from backscattered electrons while also shrinking electron puddle sizes³⁶. The smaller puddle sizes and improved signal-to-noise ratio, in turn, made electron counting feasible. By comparing the distribution of neighboring puddle distances in the ultra-low dose rate datasets with those in simulated images, we were able to estimate the ratio of primary to backscattered electrons to be ~8.6 (Supplementary Note 13). A feature of L1 reduction is that all the puddle shape information is retained. In the future this shape information may be useful in algorithmically distinguishing backscattered electrons from primary electrons, leading to a further reduction in noise due to backscattered electrons.

On-the-fly compression pipeline. Continuous on-the-fly data reduction and compression for 90 min were performed using 10 cores of the computer shipped with the DE-16 detector. This computer has two E5-2687v4 Intel Xeon processors (24-cores, 12 cores per chip, each core running at an average of 3.0 GHz base clock

rate), 128 GB DDR4 RAM, and is connected to a 1 Petabyte IBM GPFS NAS via a 10 GbE connection.

If the raw data stream coming from the detector is accessible in-memory (RAM), reduction compression can be performed directly on the incoming data stream. However, many detectors (including the DE-16 detector) make the raw data available only after it is written to disk. Reduction compression then requires simultaneously reading the data from disk back into RAM while more data from the detector is being written to disk. While sufficiently fast SSDs in RAID 0 can support multiplexed reads and writes to different parts of the RAID partition, a more scalable solution is to use a virtual file pointer to a location in fast DDR RAM (using RAM-disk). DDR RAMs, in fact, have sufficient read-write throughputs such that multiple ReCoDe threads can read different sections of the available data stream parallelly, while new data coming from the detector is written.

For continuous on-the-fly data collection, the StreamPix software was configured with a script to acquire data in five-second chunks and save each chunk in a separate file. The DE-16 acquisition software does not allow direct in-memory access to data coming from the detector to RAM, restricting access to data only after it has been written to disk. While SSDs have fast enough write speeds to keep up with the throughput of the DE-16 detector, on-the-fly reduction compression requires simultaneous write and read, each at 3 GB/s, which is not possible even with SSDs. To overcome this problem, a virtual file pointer to a location in fast DDR RAM (using RAM-disk) was used. When StreamPix finishes writing a five-second file to RAM-disk, the ReCoDe queue manager adds the file to the queue and informs the ReCoDe server. The ReCoDe server then picks off the next five-second file in the queue, where each processing thread in the server independently reads a different subset of frames within this file. Subsequently, each thread independently appends its reduced and compressed output to its own intermediate file on the NAS server via the 10 GbE connection. When the ReCoDe server is finished processing a five-second file in the queue it informs ReCoDe queue manager, which then deletes this file from RAM-disk. When the acquisition is complete all intermediate files are automatically merged into a single ReCoDe file where the reduced and compressed frames are time-ordered.

While the RAM-disk based approach bypasses read-writes to SSDs, it requires copying the same data in RAM twice. First, the data stream from the detector is written to an inaccessible partition on the RAM, then copied to the readable RAM-disk partition. If we had direct access to first copy in the currently inaccessible partition on the RAM, the subsequent copy to the RAM-disk can be eliminated, hence freeing up important read-write bandwidth on the RAM. At the DE-16 detector's throughput of 3.08 GB/s, this copying (read and write) uses a significant portion of the RAM's bandwidth (6.16 GB/s out of DDR4 RAM's 21–27 GB/s, or 20–25 GiB/s, transfer rate). Direct access to the detector's data stream without such copying will, therefore, enable reduction compression at even higher throughputs.

Data availability

Raw detector data that were used to generate the figures in this manuscript are available upon email request to the corresponding author.

Code availability

A fully parallelized Pythonic implementation of ReCoDe with features for on-the-fly reduction compression is available at the Github code repository³⁷: <https://github.com/NDLOHGRP/pyReCoDe>. Python notebooks used to generate the figures in this manuscript are also included in this code repository.

Received: 19 November 2019; Accepted: 9 December 2020;

Published online: 28 January 2021

References

- Datta, A., Chee, S. W., Bammes, B., Jin, L. & Loh, D. What can we learn from the shapes of secondary electron puddles on direct electron detectors? *Microsc. Microanal.* **23**, 190–191 (2017).
- Li, X., Zheng, S. Q., Egami, K., Agard, D. A. & Cheng, Y. Influence of electron dose rate on electron counting images recorded with the K2 camera. *J. Struct. Biol.* **184**, 251–260 (2013).
- Johnson, I. J. et al. Development of a fast framing detector for electron microscopy. In *2016 IEEE Nuclear Science Symposium, Medical Imaging Conference and Room-Temperature Semiconductor Detector Workshop (NSS/MIC/RTSD)* 1–2 (IEEE 2016).
- Chee, S. W., Anand, U., Bisht, G., Tan, S. F. & Mirsaidov, U. Direct observations of the rotation and translation of anisotropic nanoparticles adsorbed at a liquid-solid interface. *Nano Lett.* **19**, 2871–2878 (2019).
- Levin, B. D. A., Lawrence, E. L. & Crozier, P. A. Tracking the picoscale spatial motion of atomic columns during dynamic structural change. *Ultramicroscopy* **213**, 112978 (2020).
- Liao, H.-G. et al. Nanoparticle growth. Facet development during platinum nanocube growth. *Science* **345**, 916–919 (2014).
- Zheng, S. Q. et al. MotionCor2: anisotropic correction of beam-induced motion for improved cryo-electron microscopy. *Nat. Methods* **14**, 331–332 (2017).
- Guo, H. et al. Electron-event representation data enable efficient cryoEM file storage with full preservation of spatial and temporal resolution. *IUCr* **7**, 860–869 (2020).
- Iudin, A., Korir, P. K., Salavert-Torres, J., Kleywegt, G. J. & Patwardhan, A. EMPIAR: a public archive for raw electron microscopy image data. *Nat. Methods* **13**, 387–388 (2016).
- Henderson, R. et al. Outcome of the first electron microscopy validation task force meeting. *Structure* **20**, 205–214 (2012).
- McCallum, J. C. Disk Drive Prices (1955–2019). <https://jcmr.net/diskprice.htm>. (2019).
- Klein, A. The Cost of Hard Drives Over Time. *Backblaze Blog|Cloud Storage & Cloud Backup* <https://www.backblaze.com/blog/hard-drive-cost-per-gigabyte/> (2017).
- Deutsch, L. P. DEFLATE Compressed Data Format Specification Version 1.3. <http://zlib.net/> (1996).
- Burrows, M. & Wheeler, D. J. Block-sorting Lossless Data Compression Algorithm, SRC Research Report 124, Digital Systems Research Center, Palo Alto., (1994).
- bzip2. bzip2: Home. <https://www.sourceware.org/bzip2/> (1996).
- Pavlov, I. LZMA SDK (Software Development Kit). <https://www.7-zip.org/sdk.html> (2013).
- Collet, Y. lz4. <https://github.com/lz4> (2011).
- Dean, J., Ghemawat, S. & Gunderson, S. H. snappy. <https://github.com/google/snappy> (2011).
- Masui, K. et al. A compression scheme for radio data in high performance computing. *Astron. Comput.* **12**, 181–190 (2015).
- Chee, S. W., Baraissov, Z., Loh, N. D., Matsudaira, P. T. & Mirsaidov, U. Desorption-mediated motion of nanoparticles at the liquid-solid interface. *J. Phys. Chem. C* **120**, 20462–20470 (2016).
- Duane Loh, N. et al. Multistep nucleation of nanocrystals in aqueous solution. *Nat. Chem.* **9**, 77–82 (2016).
- Koneti, S. et al. Fast electron tomography: Applications to beam sensitive samples and in situ TEM or operando environmental TEM studies. *Mater. Charact.* **151**, 480–495 (2019).
- Jiang, Y. et al. Electron ptychography of 2D materials to deep sub-ångström resolution. *Nature* **559**, 343–349 (2018).
- Ophus, C. Four-Dimensional Scanning Transmission Electron Microscopy (4D-STEM): from scanning nanodiffraction to ptychography and beyond. *Microsc. Microanal.* **25**, 563–582 (2019).
- Pelz, P. M., Qiu, W. X., Bücker, R., Kassier, G. & Miller, R. J. D. Low-dose cryo electron ptychography via non-convex Bayesian optimization. *Sci. Rep.* **7**, 9883 (2017).
- Hattne, J. et al. Analysis of global and site-specific radiation damage in Cryo-EM. *Structure* **26**, 759–766.e4 (2018).
- Clough, R. & Kirkland, A. I. Chapter One-direct digital electron detectors. In *Advances in Imaging and Electron Physics* (ed. Hawkes, P. W.) Vol. 198, 1–42 (Elsevier, 2016).
- McLeod, R. A., Diogo Righetto, R., Stewart, A. & Stahlberg, H. MRCZ-A file format for cryo-TEM data with fast compression. *J. Struct. Biol.* **201**, 252–257 (2018).
- Casañal, A. et al. Architecture of eukaryotic mRNA 3'-end processing machinery. *Science* **358**, 1056–1059 (2017).
- Falcon, B. et al. Novel tau filament fold in chronic traumatic encephalopathy encloses hydrophobic molecules. *Nature* **568**, 420–423 (2019).
- Hofmann, S. et al. Conformation space of a heterodimeric ABC exporter under turnover conditions. *Nature* **571**, 580–583 (2019).
- Zhao, P. et al. Activation of the GLP-1 receptor by a non-peptidic agonist. *Nature* **577**, 432–436 (2020).
- Allahgholi, A. et al. AGIPD, a high dynamic range fast detector for the European XFEL. *J. Instrum.* **10**, C01023 (2015).
- El-Desouki, M. et al. CMOS image sensors for high speed applications. *Sensors* **9**, 430–444 (2009).
- Vincent, L. & Soille, P. Watersheds in digital spaces: an efficient algorithm based on immersion simulations. *IEEE Trans. Pattern Anal. Mach. Intell.* **13**, 583–598 (1991).
- McMullan, G., Faruqi, A. R. & Henderson, R. Direct electron detectors. *Methods Enzymol.* **579**, 1–17 (2016).
- Datta, A. et al. NDLOHGRP/pyReCoDe v0.1.0 (Version v0.1.0). Zenodo. <https://doi.org/10.5281/zenodo.4267160> (2020).

Acknowledgements

The authors would like to thank Xiaoxu Zhao for his kind contribution of the molybdenum disulfide 2d crystals. We also thank Liang Jin for helpful discussions, and Benjamin Bammes regarding the DE-16 detector and data processing pipeline, plus his careful reading of the manuscript. We are grateful to Ming Pan, Ana Pakzad, and Cory Czarnik for details about the K2-IS detector and associated binary data formats. The authors would also like to

acknowledge Chong Ping Lee from the Centre for Bio-Imaging Sciences (CBIS) at the National University of Singapore (NUS) for training and microscope facility management support, and Bai Chang from CBIS for IT infrastructure support. A.D. and Y.B. were funded by the Singapore National Research Foundation (grant NRF-CRP16-2015-05), with additional support from the NUS Startup grant (R-154-000-A09-133), NUS Early Career Research Award (R-154-000-B35-133), and the Singapore Ministry of Education Academic Research Fund Tier 1 Grant (R-154-000-C01-114).

Author contributions

D.L. and A.D. conceived the project. A.D., D.B., K.F.N., S.W.C., and J.S. collected data, which A.D., K.F.N., M.D., and Y.B. analyzed. ReCoDe framework was programmed by A.D., under D.L.'s advice. A.D. and D.L. wrote the manuscript with inputs from all the authors.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41467-020-20694-z>.

Correspondence and requests for materials should be addressed to N.D.L.

Peer review information *Nature Communications* thanks Robert McLeod and the other anonymous reviewers for their contribution to the peer review of this work.

Reprints and permission information The online version contains supplementary material available at <http://www.nature.com/reprints>

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021