## Technical Note

# A Space Efficient Direct Access Data Compression Approach for Mass Spectrometry Imaging

Patrik Källback, Anna Nilsson, Mohammadreza Shariatgorji, and Per E. Andren

### Just Accepted

1
2
3
4    # A Space Efficient Direct Access Data Compression

5
6
7
8
9    # Approach for Mass Spectrometry Imaging

10
11
12
13    3
14
15
16    4    (Technical Note)
17
18    5
19
20
21    6    Patrik Källback, Anna Nilsson, Mohammadreza Shariatgorji, and Per E. Andrén*
22
23
24    7    Biomolecular Mass Spectrometry Imaging, National Resource for Mass Spectrometry Imaging, Science
25
26    8    for Life Laboratory, Department of Pharmaceutical Biosciences, Uppsala University, Box 591 BMC,
27
28    9    75124, Uppsala, Sweden.
29
30
31    10
32
33    11   * Corresponding author: Tel: +46-70-167 9334; E-mail address: per.andren@farmbio.uu.se
34
35
36    12
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56                                   1
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

## Abstract

Advances in mass spectrometry imaging that improve both spatial and mass resolution are resulting in increasingly larger data files that are difficult to handle with current software. We have developed a novel near-lossless compression method with data entropy reduction that reduces the file size significantly. The reduction in data size can be set at four different levels (coarse, medium, fine, and superfine) prior to running the data compression. This can be applied to spectra, spectrum-by-spectrum, or on the transpose arrays, array-by-array, to efficiently read the data without decompressing the whole dataset. The results show that a compression ratio of up to 5.9:1 was achieved for data from commercial mass spectrometry software programs and 55:1 for data from our in-house developed msIQuant program. Comparing the average signals from regions of interest, the maximum deviation was 0.2% between compressed and uncompressed datasets with coarse accuracy for the data entropy reduction. In addition, when accessing the compressed data by selecting a random $m/z$ value using msIQuant, the time to update an image on the computer screen was only slightly increased from 92 ($\pm$32) ms (uncompressed) to 114 ($\pm$13) ms (compressed). Furthermore, the compressed data can be stored on readily accessible servers for data evaluation without further data reprocessing. We have developed a space efficient, direct access data compression algorithm for mass spectrometry imaging, which can be used for various data-demanding mass spectrometry imaging applications.

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

2

## Introduction

Mass spectrometry imaging (MSI) is a powerful technique for studying the spatial distribution of biomolecules [1] such as lipids [2,3], proteins [1,4], peptides [5], endogenous metabolites and neurotransmitters [6,7], as well as pharmaceutical substances [8-10] directly in tissue sections. Powerful MS techniques, such as Fourier transform ion cyclotron resonance (FTICR), present new opportunities for the imaging of molecules because of the high mass resolution capabilities that enable the extraction of thousands of ion images during a single MSI experiment. Furthermore, the spatial resolution of matrix-assisted laser desorption ionization (MALDI) MSI is currently <5 μm [11] and, for high resolution nanoscale secondary ion mass spectrometry (nanoSIMS), it is <100 nm [12]. This means that MSI instruments that are capable of high mass resolution and high spatial resolution may produce several hundred gigabytes (GB) of data in an experiment on a large tissue section (>100,000 pixels) [13].

It is a significant challenge to handle such large amounts of data and to extract data for visualization efficiently [13,14]. Several strategies to reduce the amount of data have been developed by removing zero intensity data in the mass spectra [15], by converting profile mode spectra into centroid mode data [16] or by reducing the dimensionality of the data by compressed PCA process [17], or by wavelet transform [18]. It can also be compressed using lossless compression to minimize storage size [14]. It is also reported that intensity values from mass spectra is typecast to integer values before compression to reduce the datasize [19]. Normally, the size of MSI data is also reduced using different strategies to fit the data in the internal memory of a computer. Strategies include re-binning profile data [20], creating centroid data [16] or by reducing the number of displayed pixels [20]. Other strategies involve transposing the MSI data, which allows efficient access without performing any data reduction [13] or to divide the MSI data into a number of chunks with limited mass range [14].

Our in-house developed software msIQuant [13] stores data, non-reduced and non-compressed, spectrum-by-spectrum; their transposes are utilized for rapid visualization in real-time. The drawback of this data

3

56 structure is that a large storage space is needed. Here we report a method that efficiently compresses MSI

57 data and their transposes, enabling fast visualization of images in real-time and space efficient data

58 archiving. Each spectrum in the data is compressed using our compression method. To store data with a

59 high degree of compression, low data entropy is needed [21]. Our algorithm reduces the data entropy before

60 compression by reducing the number of significant figures in the intensity values [22], while leaving the

61 number of *m/z* channels and the precision of the *m/z* values in the spectra unaffected.  This strategy is

62 generic and can be implemented on many types of spectral data. The advantage of our method is that both

63 baseline noise and high intensity values are maintained.

64

65

## Material and Methods

67 *MSI datasets*

68 Four datasets were used for this study (Table 1), with datasets 1 and 2 acquired from mouse brain tissue

69 sections (unpublished data). Datasets 3 and 4 are from our previously published work [6,23]. All animal

70 studies were carried out in accordance with the European Communities Council Directives of 1986

71 (86/609/EEC) and 2010 (2010/63/EU) and were approved by the local ethics committees on animal

72 experiments. The MSI experiments used MALDI time-of-flight (TOF) (Ultraflextreme, Bruker Daltonics,

73 Bremen, Germany; dataset 1), MALDI quadrupole-TOF (Q-TOF) (Synapt G2si, Waters Corp.,

74 Manchester, UK; dataset 2), desorption electrospray ionization (DESI)-Q-TOF (Synapt G2si, Waters

75 Corp.; dataset 3), and MALDI FTICR instruments (Solarix XR, 12T, Bruker Daltonics; dataset 4).

76 *Data Processing*

77 All datasets were converted to imzML format [15] with the software used in the MSI experiment (see Table

78 1). The imzML converter used for dataset 1 and 4 was the software flexImaging (Bruker Daltonics) which

79 generated double precision floating-point for both *m/z* and intensity values. Dataset 2 and 3 were

4

Analytical Chemistry

80 generated by the imzML converter available in Masslynx (Waters Corp) which produced imzML data

81 with single precision for both *m/z* and intensity values. Each imzML file was converted into the msIQuant

82 format (uncompressed) where the *m/z* values are presented as double precision and the intensities as

83 single precision.

84 *Algorithms for data entropy reduction and compression*

85 Next, the data entropy was reduced and compressed datasets were created with the msIQuant data

86 compression algorithm. The algorithms for data entropy reduction and compression are described in detail

87 in the Supporting Information and are divided into five different algorithms. The data entropy reduction

88 can be carried out at four different levels: coarse accuracy (10 bits), medium accuracy (13 bits), fine

89 accuracy (16 bits), and superfine accuracy (20 bits). The uncompressed accuracy is defined using a single

90 precision floating point value.

91 *Compression and Decompression Method*

92 The data compression and decompression was derived from the lossless compression method Lempel-

93 Ziv-Markov chain algorithm (LZMA) Software Development Kit [24], which is based on LZ77 compression

94 [25].

95 *Test Software*

96 Two versions of msIQuant were used to investigate the changes in performance following data

97 compression: msIQuant (version 2.0.1.14) [13] and a modified version of msIQuant, adapted for the

98 compressed data structure.

99 *Test Procedure*

100 The four datasets were used to evaluate the performance and quality of the data entropy reduction and the

101 compression algorithm. The following evaluations were carried out: a) the compression ratio was

102 compared for both raw data from commercial MS software and msIQuant uncompressed data, and b) the

103 relative deviation between an uncompressed and a compressed spectrum, which was evaluated by

5

104 calculating the root mean square (RMS). The spectrum with the highest total ion current (TIC) was

105 selected from each dataset because they contain the highest degree of information. RMS was used instead

106 of average values since the relative deviation was close to zero and produced both positive and negative

107 values. The relative deviation between the uncompressed and compressed spectrum was calculated as

108 $\frac{(I_c[i]-I_u[i])}{I_u[i]}$, where $I_c[i]$ is the intensity from the compressed dataset and $I_u[i]$ is the intensity from the

109 uncompressed dataset. Thus, the resulting equation for calculating the RMS for the relative deviation was

110 $\text{RMS} = \sqrt{\frac{1}{n}\sum_{i=0}^{n-1}\left(\frac{(I_c[i]-I_u[i])}{I_u[i]}\right)^2}$. A further evaluation was carried out: the image of an ion of interest in

111 datasets 3 and 4 was selected and compared between the compressed and uncompressed datasets. The

112 similarity ratio between the compressed and uncompressed dataset was calculated as $\left(1 - \frac{RMS}{256}\right)$, and the

113 RMS was defined as the pixel value difference in a grayscale image, which has 256 discrete levels

114 (Supporting Information, Figure S1). The ions of interest, dipalmitoylphosphatidylcholine (DPPC)

115 [M+K]$^+$ (*m/z* 772.5) and derivatized dopamine (DA-DPP) (*m/z* 368.165) (Figure 1) were selected, since

116 both have important biological functions in the brain. More specifically, the average signals from

117 annotated ROIs of the tissue were compared between the compressed and uncompressed datasets. Finally,

118 the compressed dataset 4, stored on a server at Uppsala University, was accessed using a virtual private

119 network (VPN) client from a remote location and the time to read data was measured using a timer

120 function from the msIQuant software, after selecting a random m/z value. The broadband speed at the

121 remote location was 14 Mbit/s across the VPN connection.

122

## Results

124   To estimate the robustness of the data entropy reduction and compression method, it was tested on

125   various MSI datasets acquired using different ionization methods, mass analyzers, and software (Table 1).

126   Compression ratios using the four data entropy reduction levels used were estimated from both

127   commercial and the msIQuant software. Results from MSI experiments were investigated qualitatively

128   and quantitatively. Qualitative comparison was achieved by estimating a similarity ratio between a

129   compressed and an uncompressed image, and quantitative comparison was achieved by annotating

130   different ROIs in a tissue section and comparing the ROI average ion intensity between the compressed

131   and uncompressed dataset. Finally, individual spectra from compressed and uncompressed datasets were

132   compared to investigate the ion intensity deviation.

### *Data compression ratios and compression time*

134   Our results showed that after carrying out the data entropy reduction and compression, the compression

135   ratio (using coarse accuracy) ranged from 4.2:1 to 55:1 (Table 1) compared to the uncompressed

136   msIQuant dataset, and 0.9:1 to 5.9:1 compared with raw data not converted to msIQuant. The

137   compression time varies based on the spectra size, data entropy and the compression method (Table 1).

138   The coarse, medium and fine accuracy compression have similar compression times for each data set.

139   Compression times using superfine accuracy and 'LZMA only' are twice as long. The 'zlib only' method

140   has the fastest compression times but also the lowest compression ratios compared to the other methods.

141   When selecting an *m/z* value for image visualization, the time to update the image on the computer screen

142   increased by 22 ms (24%) for the coarse accuracy compression, from an average of 92 (±32) to 114 (±13)

143   ms (measured using the timer in msIQuant).

### *Intensity deviation using data compression*

145   When comparing the relative deviation of the intensity in a spectrum with the high TIC intensity from

146   each dataset, the RMS values ranged from 2,894 to 4,652 ppm for the coarse accuracy, 278 – 580 ppm for

7

147    the medium accuracy, 27 – 75 ppm for the fine accuracy and 1.4 – 4.6 ppm for the superfine accuracy

148    (Table 1). The relative unit ppm is used instead of percentage since the relative deviation is low (i.e. 1,000

149    ppm is equal to 0.1%).

### Qualitative and quantitative effects of data compression

151    The images generated from the uncompressed and compressed datasets are compared by calculating a

152    similarity ratio which is described in Supporting Information, Figure S1. The similarity ratio between an

153    uncompressed and a compressed ion image of DPPC (Figures 1A, B) ranged from 99.61% (coarse

154    accuracy) to 99.99 % (superfine accuracy) (Supporting Information, Figure S1).

155    Three brain structures, cortex (CTX), caudate-putamen (CPu) and anterior commissure (aca), were

156    selected in an image acquired using DESI-Q-TOF MSI and the average intensities were evaluated for

157    uncompressed and compressed (coarse accuracy) data (Figure 1C). The uncompressed dataset provided

158    the following average intensities (% of maximum intensity): CTX 64.959, CPu 51.181, and aca 25.584.

159    The coarse accuracy compression produced these average intensities (% of maximum intensity): CTX

160    64.962, CPu 51.182, and aca 25.561. This resulted in a deviation of the coarse from the uncompressed

161    dataset by CTX 0.005%, CPu 0.002%, and aca -0.089%. The deviation between the average intensities for

162    the three brain structures was less than 0.1% (absolute value).

163    Comparison of an uncompressed and compressed (course accuracy) MALDI FTICR ion image of

164    dopamine (Figures 1D, E) showed that the similarity ratio was 99.77%, and the image quality increased

165    with lower compression rates (medium accuracy, 99.92%, fine accuracy 99.98%, and superfine accuracy

166    100.00%, see Supporting Information, Figure S-1).

167    The ion intensities for dopamine were compared for three brain structures in the MALDI FTICR

168    experiment (Figure 1F). The uncompressed dataset provided average intensities (%) of: CPu 66.332, aca

169    13.290 and nucleus accumbens shell (AcbSh) 59.616. The coarse accuracy compression level provided

170    average intensities (%) of: CPu 66.330, aca 13.261 and AcbSh 59.622. This resulted in a deviation of the

8

171 coarse from the uncompressed dataset by CPu -0.004%, aca -0.223%, and AcbSh 0.010%. The deviation

172 in this dataset was less than 0.2% (absolute value).

173 To investigate how the ion intensity is affected by our compression method, spectra from three different

174 pixels in dataset 3 were compared with the uncompressed and the decompressed course accuracy

175 compression (Figure 2). The result showed that the uncompressed and decompressed spectra were

176 consistent through the whole m/z range and from low to high intensity values. For example, the average

177 relative deviation was 4,642 ppm. Additional result (Supporting Information, Figure S2) showed that the

178 uncompressed and decompressed coarse accuracy compression spectra had 99.56% similarity. When

179 compiling a quantile-quantile (Q-Q) plot of the two spectra, the slope was 1.000850, intercept 0.000984

180 and $R^2$ 0.999961, which verified that the two datasets have a common m/z distribution.

181 ***Access of compressed dataset housed on a remote server***

182 A compressed data array (coarse accuracy level) on the server at Uppsala University was accessed by a

183 computer running msIQuant at a remote location, for visualization on the computer screen. After selecting

184 a random *m/z* value from dataset 4, the time to update the image on the computer screen was between 94

185 – 1563 ms, with an average time of 279 ms (measured using the timer function in msIQuant). These

186 values were, as expected, slightly higher compared to the time it takes to visualize the same m/z ratio

187 directly from the computer (92 ms, uncompressed data and 114 ms, coarse accuracy compression).

188

## Discussion

190 Data files from MSI experiments are rapidly increasing in size, hence our aim was to develop an effective

191 way of compressing the data using a near-lossless [19] compression method. Modern methods for lossless

192 compression are based on the Lempel–Ziv 1977 algorithm (LZ77) compression method [25]. The

193 compression method implemented in our msIQuant software uses the Lempel–Ziv–Markov chain

194 algorithm (LZMA) [24]. The main attribute for a lossless compression algorithm is that it cannot carry out

195 any compression of a dataset with true random numbers because of the Shannon limit [21]. That is, in a

196 mass spectrum, there are sequences that repeat themselves, and hence a lossless compression can be used.

197 However, when a spectrum contains noise, the obtained compression rate is lower and the data entropy

198 for such a spectrum is relatively high [21]. This was demonstrated by processing using a 5-point Savitsky-

199 Golay smoothing filter [26] which was enough to increase the compression ratios, which had the highest

200 effect when coarse accuracy precision was implemented (Supporting Information, Table S1). The

201 exception was the spectrum from dataset 4 where the compression ratio was decreased. This was due to

202 the smoothing filter that introduces more non-zero data in the spectrum for datasets with zero baseline,

203 which increases the data entropy.

204 The software used to generate MSI raw data in this study (Fleximaging and Masslynx) both use lossless

205 compression methods, but are undocumented since the software structures are proprietary information.

206 The standardized imzML data format [15] can be compressed with a method called zlib [27], but the software

207 OpenMSI [14] can also reduce data with gzip compression [28]. Both compression methods are based on a

208 technique called DEFLATE [29], which in turn is based on the LZ77 algorithm [25].

209 Our way of reducing data entropy is to reduce the number of significant figures, e.g., changing from

210 representing a floating point value as double precision to a single precision or event to half precision [22]

211 (defined in standard IEEE 754 [30]). This method reduces the number of significant figures from

212 approximately 16 decimal digits to 7 or 3 decimal digits. A mass array with *m/z* values should always be

10

213    expressed with double precision as the mass accuracy of modern mass spectrometers is increasing. An

214    intensity array, however, may be expressed with less significant figures since an intensity signal may vary

215    from zero to $10^7 - 10^9$ counts. As an example, if an ion intensity has the value of $9.813854 \times 10^7$ and this

216    value is represented by a binary value, the exponent can be expressed with 3 bits, while the mantissa

217    needs to be expressed with 24 bits, a total of 27 bits. If the acceptable accuracy of the same ion intensity

218    can be expressed as $9.81 \times 10^7$, or even as $9.8 \times 10^7$, then the mantissa can be expressed with 7 bits and

219    the exponent with 3 bits, so the total size is now 10 bits. This process of reducing the number of bits of a

220    binary value from 27 to 10 bits is the same as reducing the data entropy. Details about the algorithms for

221    data entropy reduction and compression and implementation are described in Supporting Information

222    (Algorithms including pseudo-code). The algorithms were written in C++ to show the pseudo-code as

223    programming code.

224        There are other compression methods for MSI data available, e.g., 'randomized approximation

225    compressed PCA process' and 'numerical compression schemes'. The compressed PCA process reported

226    by Palmer et al. [17] reduces the dimensionality of the MSI data and has a compression ratio > 170:1, which

227    is much higher than our compression ratio (~5:1). The difference between the two methods is that the

228    compressed PCA process is lossy. The method can recreate spectrum peaks but cannot recreate the noise

229    and the fine structure in the spectra.

230    The numerical compression schemes for MS data reported by Teleman et al. [19], converts the intensities to

231    integer values before compressing the data with gzip or zlib compression methods. The two main

232    algorithms described in the paper are MS Numpress positive integer compression (numPic), and MS

233    Numpress short logged float compression (numSlof). The numPic method truncates the intensity to

234    nearest integer, while numSlof method takes the natural logarithm of the intensity multiplied by a scaling

235    factor before truncating the resulting value to nearest integer. The numSlof method is similar to our

236    algorithm but with the difference that our method use four different ranges for the integer that optimizes

11

237   compression. The compression method that is used in our algorithm is LZMA which give higher

238   compression ratio than gzip or zlib.

239   Our results showed that compressing data at the coarse accuracy level produced a compression ratio up to

240   5.9:1 when compared to data obtained using commercial software (Table 1). However, for dataset 1, the

241   compression ratio was only 0.9:1 for the msIQuant software, meaning that the compressed data was

242   larger. The reason for this is that the compressed msIQuant data contain both spectra and transposed data

243   [13] but the commercial software raw data contain only spectral data. When comparing just the spectrum

244   data size, the file size from commercial software was 3.53 GB and the msIQuant coarse accuracy

245   compressed spectra data file size was 2.21 GB (not listed in Table 1), giving a compression ratio of 1.6:1.

246   The maximum deviation was 0.2% when comparing the average ion intensity from annotated ROIs

247   between uncompressed and coarse accuracy compression (Figure 1).

248   To investigate the effectiveness of the data entropy reduction and compression algorithm, datasets

249   compressed with LZMA and zlib methods without any reduction of the precision were compared with our

250   compression method. LZMA compression only produced $1 - 42\%$ larger file sizes compared with

251   superfine accuracy, while zlib compression produced $39 - 52\%$ larger files. The benefit with zlib

252   compared with LZMA is that the method is $100 - 350$ % faster, but with the cost of less effective

253   compression ratio.

254   Images compared using image differential analysis to show differences (Figures 1A to 1B and Figures 1D

255   to 1E and Supporting Information, Figure S-1) resulted in an image similarity ratio of 99.61% and

256   99.77% between uncompressed and coarse accuracy compression. To compare the technical variation

257   with our compression algorithm, three consecutive coronal rat brain tissue sections were investigated

258   (Supporting Information, Figure S3A-F). The average intensities of the lipid PC(32:0) $[M+K]^+$ (*m/z* 772.5

259   $\pm$ 0.1) from two different brain areas (striatum and cortex) were measured. The standard deviations from

260   intensity measurement of the technical replicates were 0.519 and 0.821, while the average absolute

12

261   deviations between the uncompressed and the decompressed coarse accuracy compression were 0.000466

262   and 0.000387 respectively. This demonstrates that the loss in accuracy using our compression method is

263   1,000-fold lower than the technical variation in such experiment. Further, analyzing the ion intensity of

264   individual spectra from each dataset showed that the relative deviation decreased by one order of

265   magnitude for each degree of accuracy i.e. the medium level has a higher degree of accuracy than the

266   coarse level (Table 1, Relative deviation from uncompressed spectrum). Apart from the reduced data size,

267   another important feature of the compression algorithm is that the data can be directly accessed without

268   the need for uncompressing or reprocessing the dataset. This means that the data can be stored on a server

269   while still being accessible without moving them to a local computer for evaluation.

270   The data entropy reduction and compression algorithm described in this study is not only valid for

271   compressing MSI spectra, but could also be implemented in all various scientific methodologies where

272   spectra are obtained, such as Raman spectroscopy [31], multi-modal imaging analysis [32], and hyperspectral

273   imaging surveillance by satellites [33]. A similar approach to reducing the number of significant figures has

274   been described for computerized calculations of simulating material stress, which generates large amount

275   of data [34].

276   In summary, we have developed a space efficient, direct access data compression algorithm for mass

277   spectrometry imaging, which could also be used for various other applications where spectra are

278   generated.  We have tested both the efficiency and the accuracy of the compression, which can be set to

279   four different levels: coarse, medium, fine and superfine. We have also proven that the compressed

280   datasets can be stored on a server and still be accessed for data analysis without moving them to a local

281   computer for data processing.

282

283   **Associated Content**

284   ***Supporting Information***

13

285    The Supporting Information is available free of charge on the publication's website at DOI:.

286

# Author Information

288    *Corresponding Author*

289    Email:  per.andren@farmbio.uu.se

290    **ORCID**

291    Per E. Andrén: 0000-0002-4062-7743

292    *Author Contribution*

293    The study and manuscript was carried out and written with contributions from all the authors. All authors

294    approved the final version of the manuscript.

295    *Notes*

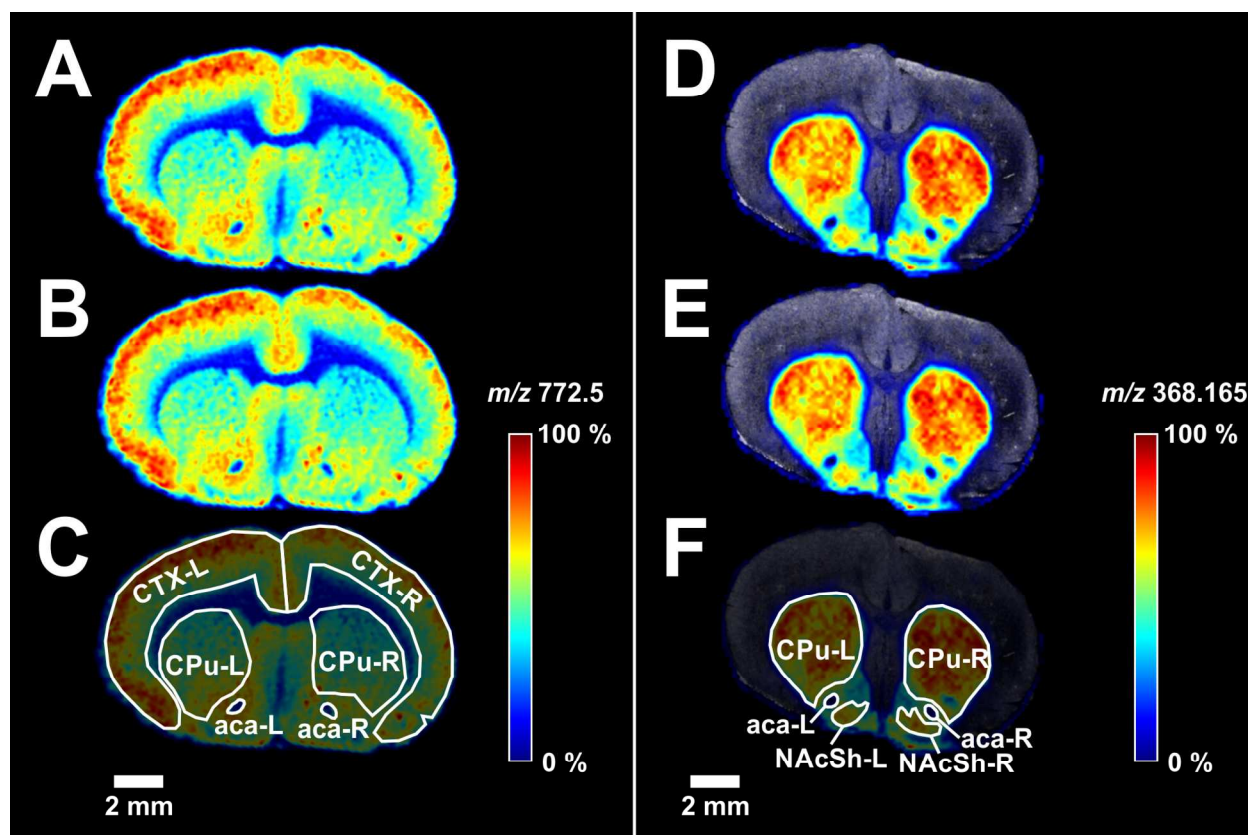296    The authors declare no competing financial interest.

297

# Acknowledgments

14

304    **Figures**

305    **Figure 1. Quality comparisons between images extracted from the uncompressed dataset and**

306    **images from the coarse accuracy (10 bit) dataset of a coronal rat brain tissue section.**
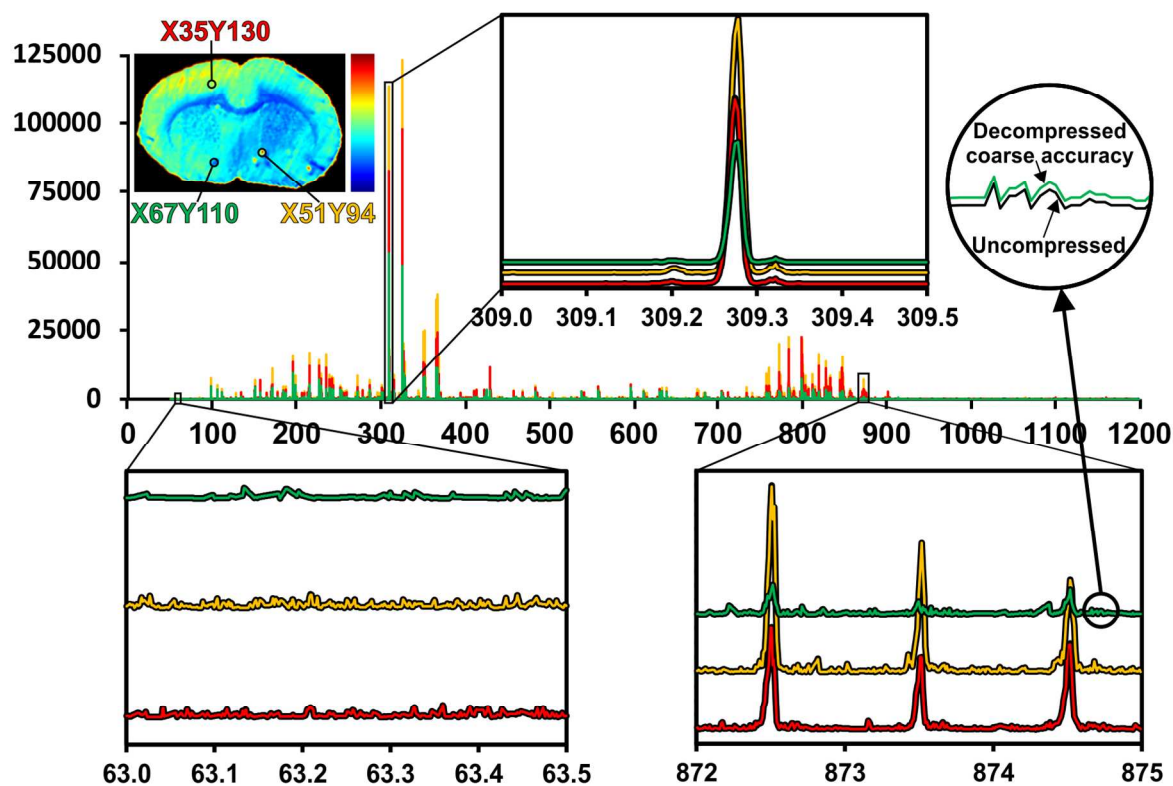
307    Tissue sections showing dipalmitoylphosphatidylcholine (DPPC) $[M+K]^+$, *m/z* 772.5 from (A)

308    uncompressed and (B) coarse accuracy compressed data acquired at a spatial distribution of 150 μm. The

309    images were generated using DESI coupled to a Q-TOF instrument. (C) Three brain regions (cerebral

310    cortex (CTX), caudate-putamen (CPu), and the anterior part of the anterior commissure (aca) were

311    selected to measure the average intensity of DPPC $[M+K]^+$. The average intensity deviations were

312    compared between the uncompressed and the coarse accuracy level compression (0.005%, 0.002% and -

313    0.089% for CTX, CPu and aca, respectively). Dopamine (DA) derivatized with DPP-TFB, m/z 368.165 is

314    shown from the (D) uncompressed dataset and (E) coarse accuracy compressed data acquired at a spatial

315    distribution of 150 μm using a MALDI FTICR mass spectrometer. (F) Three brain regions (CPu, aca and

316    nucleus accumbens shell (AcbSh) were selected to measure the average intensity of derivatized DA, and

317    the average intensity deviations were compared between the uncompressed and the coarse accuracy level

318    compression (-0.004%, -0.223% and -0.010% for CPu, aca and AcbSh respectively). The blue-red

319    rainbow color scale represents the intensity from 0 – 100 %.

15

320

321  **Figure 2. Comparison of mass spectra from three pixels in a brain tissue section before and after**

322  **compression using the coarse accuracy level compression level.**

323  Overlay of uncompressed (solid black) and decompressed coarse accuracy (red, yellow and green) spectra

324  from three pixels shown in the inset tissue image, displaying TIC-normalization-factor distribution [35].

325  Three rectangular enlargements, from different parts of the *m/z* axis show that decompressed coarse

326  accuracy spectra maintain the same fine structure as the uncompressed spectra. The spectra from the three

327  pixels are displayed with intensity offset to distinguish them from each other. Each uncompressed

328  spectrum is displayed with a thicker line than the decompressed coarse accuracy spectrum. The circular

329  enlargement displays a small part of the uncompressed and the decompressed spectrum with a small

330  intensity offset.



331

332

333  **Table**

334  **Table 1. Properties of four MSI datasets acquired from different MSI platforms.**

335  The original sizes of the data files are shown together with the compressed files at four different levels of

336  accuracy. The compression ratio in relation to the instrument-specific software used to acquire the data is

337  shown in brackets, while the compression ratio in relation to the msIQuant files is shown in bold in

338  brackets. The compression times of msIQuant files are shown as hours and minutes (hr:min) and the

339  measurements were performed on an HP Z440 workstation. The relative deviation of ion intensity

340  between compressed and uncompressed data for a single spectrum is shown as RMS (in ppm) at the four

341  different levels of accuracy.

| | Dataset 1 | Dataset 2 | Dataset 3 | Dataset 4 |
|---|---|---|---|---|
| Ion Source | MALDI | MALDI | DESI | MALDI |
| Mass Analyzer | TOF | Q-TOF | Q-TOF | FTICR |
| Spatial resolution [µm] | 50 | 100 | 150 | 150 |
| Number of pixels | 54,808 | 19,370 | 28,718 | 12,411 |
| Number of m/z channels | 84,821 | 110,376 | 310,734 | 372,516 |
| imzML converter | Fleximaging | Masslynx | Masslynx | Fleximaging |
| Raw data size [GB]* | 4.17 | 11.28 | 37.35 | 3.7 |
| msIQuant data size [GB] | *34.67* | *15.94* | *66.51* | *34.46* |
| Coarse data size [GB] | 4.57 (0.9:1) *(7.6:1)* | 3.82 (3.0:1) *(4.2:1)* | 14.49 (2.6:1) *(4.6:1)* | 0.63 (5.9:1) *(54.7:1)* |
| Medium data size [GB] | 4.71 (0.9:1) *(7.4:1)* | 5.36 (2.1:1) *(3.0:1)* | 22.11 (1.7:1) *(3.0:1)* | 0.79 (4.7:1) *(43.6:1)* |
| Fine data size [GB] | 4.79 (0.9:1) *(7.2:1)* | 6.87 (1.6:1) *(2.3:1)* | 29.05 (1.3:1) *(2.3:1)* | 0.93 (4.0:1) *(37.1:1)* |
| Superfine data size [GB] | 5.05 (0.8:1) *(6.9:1)* | 8.80 (1.3:1) *(1.8:1)* | 38.51 (1.0:1) *(1.7:1)* | 1.24 (3.0:1) *(27.8:1)* |
| LZMA only [GB] | 5.08 (0.8:1) *(6.8:1)* | 11.79 (1.0:1) *(1.4:1)* | 54.54 (0.7:1) *(1.2:1)* | 1.63 (2.3:1) *(21.1:1)* |
| Zlib only [GB] | 7.01 (0.6:1) *(4.9:1)* | 12.90 (0.9:1) *(1.2:1)* | 58.25 (0.6:1) *(1.1:1)* | 1.72 (2.2:1) *(20.0:1)* |
| Coarse comp. time [hh:mm] | 01:29 | 00:29 | 01:44 | 00:19 |
| Medium comp. time [hh:mm] | 01:19 | 00:30 | 01:52 | 00:20 |
| Fine comp. time [hh:mm] | 01:15 | 00:28 | 01:43 | 00:20 |
| Superfine comp. time [hh:mm] | 02:32 | 00:57 | 03:51 | 00:30 |
| LZMA only comp. time | 02:46 | 00:54 | 03:27 | 00:28 |
| Zlib only comp. time | 00:37 | 00:22 | 01:41 | 00:06 |
| Coarse, RMS [ppm] | 3,222 | 4,637 | 4,652 | 2,894 |
| Medium, RMS [ppm] | 529 | 580 | 580 | 278 |

18

| | | | | |
|---|---|---|---|---|
| Fine, RMS [ppm] | 75 | 72 | 72 | 27 |
| Superfine, RMS [ppm] | 1.4 | 4.5 | 4.6 | 1.4 |

*\* The raw data size refers to vendors' proprietary data.*
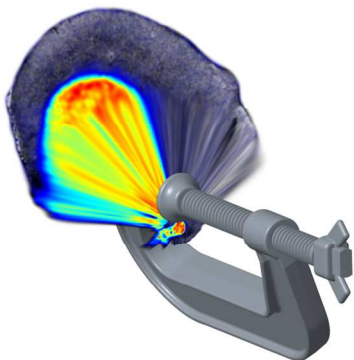
342

## References

343 **References**

344 (1) Caprioli, R. M.; Farmer, T. B.; Gile, J. *Anal Chem* **1997**, *69*, 4751-4760.

345 (2) Touboul, D.; Piednoel, H.; Voisin, V.; De La Porte, S.; Brunelle, A.; Halgand, F.; Laprevote, O.

346 *European journal of mass spectrometry* **2004**, *10*, 657-664.

347 (3) Woods, A. S.; Jackson, S. N. *The AAPS journal* **2006**, *8*, E391-395.

348 (4) Stoeckli, M.; Chaurand, P.; Hallahan, D. E.; Caprioli, R. M. *Nature medicine* **2001**, *7*, 493-496.

349 (5) Skold, K.; Svensson, M.; Nilsson, A.; Zhang, X.; Nydahl, K.; Caprioli, R. M.; Svenningsson, P.;

350 Andren, P. E. *J Proteome Res* **2006**, *5*, 262-269.

351 (6) Shariatgorji, M.; Nilsson, A.; Goodwin, R. J.; Kallback, P.; Schintu, N.; Zhang, X.; Crossman, A. R.;

352 Bezard, E.; Svenningsson, P.; Andren, P. E. *Neuron* **2014**, *84*, 697-707.

353 (7) Sugiura, Y.; Setou, M. *J Neuroimmune Pharm* **2010**, *5*, 31-43.

354 (8) Nilsson, A.; Fehniger, T. E.; Gustavsson, L.; Andersson, M.; Kenne, K.; Marko-Varga, G.; Andren, P.

355 E. *PLoS One* **2010**, *5*, e11411.

356 (9) Nilsson, A.; Goodwin, R. J.; Shariatgorji, M.; Vallianatou, T.; Webborn, P. J.; Andren, P. E. *Anal*

357 *Chem* **2015**, *87*, 1437-1455.

358 (10) Reyzer, M. L.; Hsieh, Y.; Ng, K.; Korfmacher, W. A.; Caprioli, R. M. *J Mass Spectrom* **2003**, *38*,

359 1081-1092.

360 (11) Guenther, S.; Rompp, A.; Kummer, W.; Spengler, B. *Int J Mass Spectrom* **2011**, *305*, 228-237.

361 (12) Herrmann, A. M.; Ritz, K.; Nunan, N.; Clode, P. L.; Pett-Ridge, J.; Kilburn, M. R.; Murphy, D. V.;

362 O'Donnell, A. G.; Stockdale, E. A. *Soil Biol Biochem* **2007**, *39*, 1835-1850.

363    (13) Kallback, P.; Nilsson, A.; Shariatgorji, M.; Andren, P. E. *Anal Chem* **2016**, *88*, 4346-4353.

364    (14) Rubel, O.; Greiner, A.; Cholia, S.; Louie, K.; Bethel, E. W.; Northen, T. R.; Bowen, B. P. *Anal Chem*

365    **2013**, *85*, 10354-10361.

366    (15) Rompp, A.; Schramm, T.; Hester, A.; Klinkert, I.; Both, J. P.; Heeren, R. M.; Stockli, M.; Spengler,

367    B. *Methods Mol Biol* **2011**, *696*, 205-224.

368    (16) Bielow, C.; Gropl, C.; Kohlbacher, O.; Reinert, K. *Methods Mol Biol* **2011**, *719*, 331-349.

369    (17) Palmer, A. D.; Bunch, J.; Styles, I. B. *Anal Chem* **2013**, *85*, 5078-5086.

370    (18) Barclay, V. J.; Bonner, R. F.; Hamilton, I. P. *Anal Chem* **1997**, *69*, 78-90.

371    (19) Teleman, J.; Dowsey, A. W.; Gonzalez-Galarza, F. F.; Perkins, S.; Pratt, B.; Rost, H. L.; Malmstrom,

372    L.; Malmstrom, J.; Jones, A. R.; Deutsch, E. W.; Levander, F. *Mol Cell Proteomics* **2014**, *13*, 1537-1542.

373    (20) Klinkert, I.; Chughtai, K.; Ellis, S. R.; Heeren, R. M. A. *Int J Mass Spectrom* **2014**, *362*, 40-47.

374    (21) Shannon, C. E. *Bell Syst Tech J* **1948**, *27*, 379-423, DOI: 10.1002/j.1538-7305.1948.tb01338.x.

375    (22) Maass, C.; Baer, M.; Kachelriess, M. *Med Phys* **2011**, *38 Suppl 1*, S95.

376    (23) Shariatgorji, M.; Strittmatter, N.; Nilsson, A.; Kallback, P.; Alvarsson, A.; Zhang, X.; Vallianatou,

377    T.; Svenningsson, P.; Goodwin, R. J.; Andren, P. E. *Neuroimage* **2016**, *136*, 129-138.

378    (24) Pavlov, I. *LZMA SDK (Software Development Kit)* **1998**, http://www.7-zip.org/.

379    (25) Ziv, J.; Lempel, A. *IEEE Trans Inf Theory* **1977**, *23*, 337-343.

380    (26) Savitzky, A.; Golay, M. J. E. *Anal Chem* **1964**, *36*, 1627-1639.

381    (27) Gailly, J. l.; Adler, M. *zlib* **1995**, https://zlib.net/.

21

382    (28) Gailly, J. l.; Adler, M. *gzip* **1992**, https://www.gnu.org/software/gzip/.

383    (29) Katz, P. W., U.S.A., *Patent: US5051745 A* **1991**.

384    (30) IEEE Computer Society, IEEE Std 754-2008, *IEEE* **2008**, p 70.

385    (31) Freudiger, C. W.; Min, W.; Saar, B. G.; Lu, S.; Holtom, G. R.; He, C.; Tsai, J. C.; Kang, J. X.; Xie,

386    X. S. *Science* **2008**, *322*, 1857-1861.

387    (32) Van de Plas, R.; Yang, J.; Spraggins, J.; Caprioli, R. M. *Nat Methods* **2015**, *12*, 366-372.

388    (33) Cudahy, T. J.; Hewson, R.; Huntington, J. F.; Quigley, M. A.; Barry, P. S. *Int Geosci Remote Se*

389    **2001**, 314-316.

390    (34) Thole, C. A. *5th European LSDYNA Conference* **2005**, 1-6.

391    (35) Deininger, S. O.; Cornett, D. S.; Paape, R.; Becker, M.; Pineau, C.; Rauser, S.; Walch, A.; Wolski, E.

392    *Analytical and Bioanalytical Chemistry* **2011**, *401*, 167-181.

393

394

22

395 **For TOC only**



396