

# Memory Efficient Principal Component Analysis for the Dimensionality Reduction of Large Mass Spectrometry Imaging Datasets

Race, Alan; Steven, Rory; Palmer, Andrew; Styles, Iain; Bunch, Josephine

DOI:  
[10.1021/ac302528v](https://doi.org/10.1021/ac302528v)

*Document Version*  
Peer reviewed version

*Citation for published version (Harvard):*  
Race, A, Steven, R, Palmer, A, Styles, I & Bunch, J 2013, 'Memory Efficient Principal Component Analysis for the Dimensionality Reduction of Large Mass Spectrometry Imaging Datasets', *Analytical Chemistry*, vol. 85, no. 6, pp. 3071-3078. <https://doi.org/10.1021/ac302528v>

[Link to publication on Research at Birmingham portal](#)

## General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

## Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

# Supporting information for: Memory Efficient Principal Component Analysis for the Dimensionality Reduction of Large Mass Spectrometry Imaging Datasets

Alan M. Race,<sup>†,‡,¶</sup> Rory T. Steven,<sup>†,¶</sup> Andrew D. Palmer,<sup>†,‡,¶</sup> Iain B. Styles,<sup>\*,‡</sup> and  
Josephine Bunch<sup>\*,†,¶</sup>

*Centre for Physical Sciences of Imaging in the Biomedical Sciences, School of Chemistry,  
University of Birmingham, Edgbaston, Birmingham, B15 2TT, United Kingdom., School of  
Computer Science, University of Birmingham, Edgbaston, Birmingham, B15 2TT, United  
Kingdom., and School of Chemistry, University of Birmingham, Edgbaston, Birmingham, B15  
2TT, United Kingdom.*

E-mail: i.b.styles@cs.bham.ac.uk; j.bunch@bham.ac.uk

## Abstract

Here we provide all mass spectrometry imaging experimental details, a detailed description of the peak detection algorithm, a figure demonstrating a comparison between the proportion of an entire organ that can be processed with *princomp* and the proposed method and a figure to show the clustering results for  $k=2\dots 10$ . We also provide a MATLAB implementation of the algorithm.

---

\*To whom correspondence should be addressed

<sup>†</sup>PSIBS, University of Birmingham

<sup>‡</sup>School of Computer Science, University of Birmingham

<sup>¶</sup>School of Chemistry, University of Birmingham

## Materials

Methanol (HPLC grade) used in preparation of matrix and analyte solutions was purchased from Fisher Scientific (Leicestershire, UK). Trifluoroacetic acid (TFA, 99% purity) and MALDI matrices  $\alpha$ -cyano-4-hydroxycinnamic acid (CHCA) and 4-nitroaniline (PNA) were purchased from Sigma Aldrich (Dorset, UK).

## Single rat brain image

MALDI MS imaging of formalin fixed rat brain was recently described.<sup>1</sup> In summary, a formalin fixed rat brain was frozen and bisected along the midline. 12  $\mu\text{m}$  sagittal sections were thaw mounted onto stainless steel MALDI target plates and coated in 5  $\text{mgmL}^{-1}$  CHCA prepared in 80% methanol (0.1% trifluoroacetic acid) using an automated matrix deposition system (TM sprayer from HTX Technologies, NC, U.S.A.). Sections were coated using 8 deposition cycles at 150 °C, 10 psi, a flow rate of 0.25 mL/min with a stage velocity of 500 mm/min. Positive ion MALDI MS imaging was carried out on a QqTOF (Qstar Elite) mass spectrometer (Applied Biosystems, Foster City, U.S.A.), using an Nd:YAG (355 nm) laser operated at 20% available power (2.1  $\mu\text{J}$ ) and a repetition rate of 500 Hz. Sequential spectra were acquired at a resolution of 100  $\mu\text{m}^2$  using the “dynamic pixel” setting (oMALDI, 5.1). External calibration was performed using Analyst QS 1.1.

## 3-D mouse brain image

Mouse brain samples were stored at 193 K until sectioned. Beginning within 2 mm parallel to the sagittal midline and working toward the midline, 14  $\mu\text{m}$  serial sections were collected (Leica CM1810) and thaw mounted onto MALDI imaging plates. The serial sections were thaw mounted onto a single standard stainless steel imaging plate (Applied Biosystems). Samples were coated in matrix (PNA) solution (20  $\text{mgmL}^{-1}$  in 80%  $\text{CH}_3\text{OH}$ , 0.1 % TFA) using an artist airbrush (Draper

Air Tools Airbrush Kit (Hampshire, UK)) propelled by dry N<sub>2</sub>. MALDI TOF MSI analysis was carried out on a QSTAR XL QqTOF instrument using Analyst QS 1.1 with oMALDI server 5.1 (Applied Biosystems). An Nd:YVO<sub>4</sub> (Elforlight: SPOT-10-100-355) DPSS laser (355 nm), was triggered by a Thurlby Thandar Instruments (Huntingdon, Cambridgeshire) TGP110 10MHz Pulse Generator. All imaging data were acquired at a pulse energy of 12 mJ and rep rate of 5kHz. All MALDI MS images collected in positive ion mode using a acquisition at a speed of 1 mms<sup>-1</sup>, with pixel dimensions of 100 μm. External calibration was performed using Analyst QS 1.1. Two sections were excluded from further analysis due to a tissue tear which affected the registration significantly.

## Second Derivative Gradient Peak Detection

Where the first derivative of a spectrum crosses the y-axis ( $y = 0$ ), the spectrum trace has altered from increasing to decreasing, or vice-versa indicating the presence of a maxima (peak) or minima (trough). To distinguish between a peak and a trough, the second derivative of the spectrum is calculated and if the value of the second derivative at the peak/trough  $m/z$  bin is negative then it is a peak, however if it is positive then it is a trough. Only peaks were retained; troughs introduced through the Savitzky-Golay smoothing were discarded.

## References

- (1) Carter, C.; McLeod, C.; Bunch, J. *Journal of the American Society for Mass Spectrometry* **2011**, *22*, 1991–1998.
- (2) Fonville, J. M.; Carter, C.; Cloarec, O.; Nicholson, J. K.; Lindon, J. C.; Bunch, J.; Holmes, E. *Analytical Chemistry* **2012**, *84*, 1310–1319.

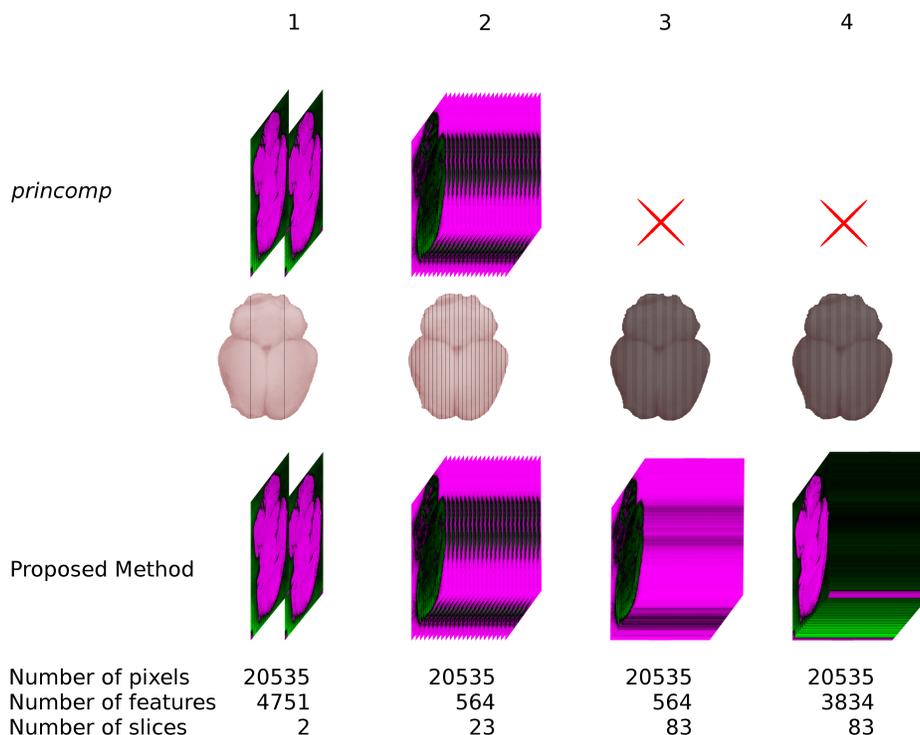


Figure 1: Simulated 3D MALDI MSI data of a 2 cm x 1 cm rat brain through repetition of a single 12  $\mu\text{m}$  section image and the corresponding first principal component score image when using *princomp* (top) and the proposed method (bottom). Column 1 shows the maximum number of sections (2 sections) that could be retained if the data was binned at 0.2  $m/z$  (4751 bins), the standard bin width used in BIOMAP (Novartis), and then analysed with *princomp* with 8 GB RAM. Distance between sections would be 325  $\mu\text{m}$ . Column 2 shows the maximum number of sections (23 sections) that could be retained if informative peaks were extracted<sup>2</sup> (564 extracted  $m/z$  bins) prior to PCA using *princomp*. Distance between sections would be 30  $\mu\text{m}$ . Column 3 shows the PCA results if the entire brain was sectioned and analysed (83 sections) and informative peaks were extracted.<sup>2</sup> The red cross indicates *princomp* failed due to memory limitations. Column 4 shows the PCA results if the entire brain was sectioned and analysed (83 sections) and all detected peaks (3834 detected peaks) were retained when determined from the entire data set.

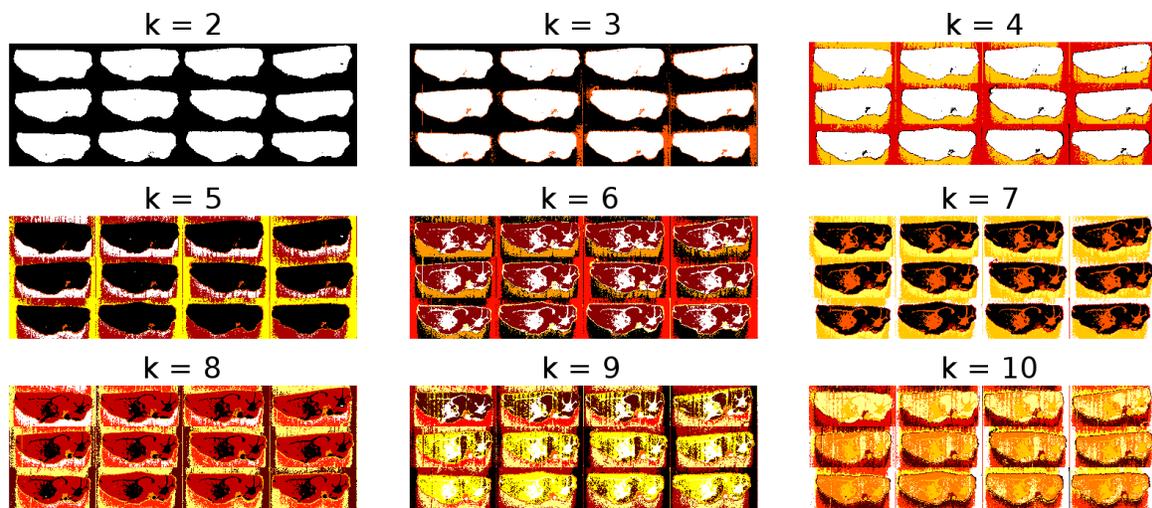


Figure 2:  $k$ -means applied to PC 1-40 scores of serial mouse brain sections with varying values for  $k = 2 \dots 10$ .