

Comparison of seven modelling algorithms for GABA-edited ^1H -MRS

Evaluation of data from 20 sites, with: FSL-MRS, Gannet, jMRUI (AMARES, QUEST), LCModel, Osprey and Tarquin

Alexander R. Craven^{a,b,c,*}, Pallab K. Bhattacharyya^d, William T. Clarke^{e,f}, Ulrike Dydak^g, Richard A. E. Edden^{h,i}, Lars Ersland^{b,a}, Pravat K. Mandal^{l,k}, Mark Mikkelsen^{h,i,l}, James B. Murdoch, Jamie Near^{m,n,o}, Reuben Rideaux^p, Deepika Shukla^{q,r,j}, Min Wang^s, Martin Wilson^t, Helge Zöllner^{h,i}, Kenneth Hugdahl^{a,u,v}, Georg Oeltzschner^{h,i}

^a Department of Biological and Medical Psychology, University of Bergen, Bergen, Norway

^b Department of Clinical Engineering, Haukeland University Hospital, Bergen, Norway

^c NORMENT Center of Excellence, Haukeland University Hospital, Bergen, Norway

^d Cleveland Clinic Foundation, Imaging Institute, Cleveland, Ohio, USA

^e Wellcome Centre for Integrative Neuroimaging, FMRIB, Nuffield Department of Clinical Neurosciences, University of Oxford, Oxford, United Kingdom

^f MRC Brain Network Dynamics Unit, University of Oxford, Oxford, United Kingdom

^g School of Health Sciences, Purdue University, West Lafayette, Indiana, USA

^h Russell H. Morgan Department of Radiology and Radiological Science, The Johns Hopkins University School of Medicine, Baltimore, Maryland, USA

ⁱ F. M. Kirby Research Center for Functional Brain Imaging, Kennedy Krieger Institute, Baltimore, Maryland, USA

^j NeuroImaging and NeuroSpectroscopy (NINS) Laboratory, National Brain Research Centre, Gurgaon, India

^k Florey Institute of Neuroscience and Mental Health, Parkville, Melbourne, Victoria, Australia

^l Department of Radiology, Weill Cornell Medicine, New York, New York, USA

^m Centre d'Imagerie Cérébrale, Douglas Mental Health University Institute, Montreal, Canada

ⁿ Department of Biomedical Engineering, McGill University, Montreal, Canada

^o Department of Psychiatry, McGill University, Montreal, Canada

^p Queensland Brain Institute, The University of Queensland, Brisbane, Australia

^q Perinatal Trials Unit Foundation, Bengaluru, India

^r Centre for Perinatal Neuroscience, Imperial College London, London, United Kingdom

^s College of Biomedical Engineering and Instrument Science, Zhejiang University, Hangzhou, China

^t Centre for Human Brain Health and School of Psychology, University of Birmingham, Birmingham, United Kingdom

^u Division of Psychiatry, Haukeland University Hospital, Bergen, Norway

^v Department of Radiology, Haukeland University Hospital, Bergen, Norway

* Corresponding author:

Alexander R. Craven (M.Sc.)

Department of Biological and Medical Psychology, University of Bergen

Jonas Lies vei 91, 5009 Bergen, Norway

alex.craven@uib.no

Word count: 6637 (main text); 242 (abstract); 76 (graphical abstract)

1 **Abstract:**

2
3 Edited MRS sequences are widely used for studying GABA in the human brain.
4 Several algorithms are available for modelling these data, deriving metabolite concentration
5 estimates through peak fitting or a linear combination of basis spectra. The present study
6 compares seven such algorithms, using data obtained in a large multi-site study.

7 GABA-edited (GABA+, TE = 68 ms MEGA-PRESS) data from 222 subjects at 20 sites
8 were processed via a standardised pipeline, before modelling with FSL-MRS, Gannet,
9 AMARES, QUEST, LCModel, Osprey and Tarquin, using standardised vendor-specific basis
10 sets (for GE, Philips and Siemens) where appropriate.

11 After referencing metabolite estimates (to water or creatine), systematic differences
12 in scale were observed between datasets acquired on different vendors' hardware,
13 presenting across algorithms. Scale differences across algorithms were also observed.

14 Using the correlation between metabolite estimates and voxel tissue fraction as a
15 benchmark, most algorithms were found to be similarly effective in detecting differences in
16 GABA+. An inter-class correlation across all algorithms showed single-rater consistency for
17 GABA+ estimates of around 0.38, indicating moderate agreement. Upon inclusion of a basis
18 set component explicitly modelling the macromolecule signal underlying the observed 3.0
19 ppm GABA peaks, single-rater consistency improved to 0.44. Correlation between discrete
20 pairs of algorithms varied, and was concerningly weak in some cases.

21 Our findings highlight the need for consensus on appropriate modelling parameters
22 across different algorithms, and for detailed reporting of the parameters adopted in
23 individual studies to ensure reproducibility and meaningful comparison of outcomes
24 between different studies.

25

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

Keywords:

MRS; GABA; spectral editing; MEGA-PRESS; quantification; macromolecule

Highlights:

- GABA-edited MRS data from 222 healthy adults across 20 research sites were analyzed
- Data were modelled using seven different algorithms, yielding GABA+ and Glx estimates
- Moderate agreement was seen across all the tested algorithms
- Adding a component to represent co-edited macromolecule signals improved concordance
- Baseline modelling emerged as major factor differentiating outcomes

List of Abbreviations

1		
2		
3		
4	Cho	<i>choline</i>
5	CI	<i>confidence interval</i>
6	Cr	<i>creatine</i>
7	CRLB	<i>Cramér-Rao lower bound for uncertainty</i>
8	diff	<i>difference (edited) spectrum</i>
9	ECC	<i>eddy-current correction</i>
10	FD	<i>frequency domain</i>
11	FID	<i>free induction decay (observed time-domain signal)</i>
12	FWHM	<i>linewidth: full width at half maximum</i>
13	GABA	<i>γ-aminobutyric acid</i>
14	GABA+	<i>total edited signal at 3 ppm; GABA with underlying coedited signals</i>
15	Gln	<i>glutamine</i>
16	Glu	<i>glutamate</i>
17	Glx	<i>combined signal of glutamate + glutamine</i>
18	GSH	<i>glutathione</i>
19	H ₂ O (noTC)	<i>water (noTC: referenced without tissue class correction)</i>
20	HSVD	<i>Hankel singular value decomposition</i>
21	i.u.	<i>institutional units</i>
22	ICC	<i>intra-class correlation coefficient</i>
23	LCM	<i>linear combination modelling</i>
24	MAD	<i>median absolute deviation</i>
25	MEGA-PRESS	<i>Mescher–Garwood point-resolved spectroscopy</i>
26	MMx(y)	<i>macromolecule signal around x(.y) ppm</i>
27	NAA	<i>N-acetylaspartate</i>
28	NAAG	<i>N-acetylaspartylglutamate</i>
29	p _{holm}	<i>Holm-Bonferroni adjusted p-value</i>
30	ppm	<i>parts per million</i>
31	Q-Q	<i>quantile-quantile</i>
32	R1-4	<i>adopted rejection criteria; see methods section</i>
33	SD	<i>standard deviation</i>
34	SNR	<i>signal-to-noise ratio</i>
35	tCr	<i>total creatine (creatine+phosphocreatine)</i>
36	TD	<i>time domain</i>
37	VPC	<i>variance partition coefficients</i>
38		
39		

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

1 Introduction

Several software packages and modelling algorithms are available for processing and quantifying MR spectroscopy (MRS) data. While they are all designed to extract quantitative estimates of metabolite levels from spectra, the packages differ significantly in their approach to processing and modelling the underlying data, and isolating the components of interest from any artefactual signals therein. This may give rise to systematic differences in metabolite estimates between different software packages. While an effect of “choice of software” has been documented for short-echo-time data ¹⁻⁴, similar studies for GABA-edited MRS quantification are lacking.

Spectral editing experiments ^{5,6}, such as the widely used MEGA-PRESS for the selective detection of GABA, present a special case for quantification. In a typical MEGA-PRESS editing sequence, two interleaved sub-spectra are acquired: the edit-ON sub-spectrum in which coupling to GABA spins at 3 ppm is refocused, and the edit-OFF sub-spectrum in which it is not. Subtracting the edit-ON and edit-OFF sub-spectra yields a relatively sparse difference spectrum, featuring prominent broad signal for GABA (with underlying macromolecule contributions) at 3 ppm and co-edited signals including glutamate (Glu) and glutamine (Gln) peaks (usually reported collectively as Glx) around 3.75 ppm, and strong negative peaks close to the editing frequency (primarily N-acetylaspartate (NAA) and N-acetylaspartylglutamate (NAAG)).

Most notable among the challenges for modelling edited spectra are co-edited macromolecular signals coupled to spins near the editing frequency ⁶, some of which appear in the same frequency range as the GABA and Glx signals and therefore interfere with their unambiguous modelling. As they are broad and poorly characterized, no consensus currently exists on how they should be accounted for in the modelling stage. Constrained by the inability to reliably separate GABA and macromolecules, their composite (GABA+) is commonly reported.

A rigorous assessment of the comparability of GABA+ estimates obtained across a range of different analysis software packages is currently lacking. Several prior studies,

1 including⁷⁻¹⁰, have investigated test-retest reproducibility of GABA+ estimates using a small
2 selection of available software packages, but without detailed examination of the
3 differences in estimates arising between software packages. Each considered data from a
4 single site only. Another study¹¹ has investigated GABA+ estimates from Gannet and Tarquin
5 compared to a simulated “ground truth”, specifically with respect to the influence of signal-
6 to-noise ratio (SNR) and linewidth on estimates, showing that the two algorithms agreed
7 under favourable conditions of linewidth and signal-to-noise ratio but diverged under poorer
8 conditions – however, only two algorithms were included in this analysis. A recent
9 conference paper¹² has reported early findings from GABA-edited MEGA-PRESS data
10 showing moderate associations between five different algorithms, with data from four sites
11 (representing two scanner vendors), albeit with divergent processing. A more thorough
12 examination, covering a broader range of sites and an extended selection of contemporary
13 algorithms is required to better characterise the noted discrepancies.

14 Therefore, to establish the degree to which different software packages agree in
15 estimating GABA+ from MEGA-PRESS data, this study compares GABA+ estimates from
16 seven modelling algorithms: FSL-MRS¹³, Gannet¹⁴, LCMoDel¹⁵, Osprey¹⁶, Tarquin^{17,18},
17 AMARES¹⁹, and QUEST^{20,21}, with the last two implemented in the jMRUI software package
18^{22,23}. Estimates of Glx from the difference spectra are also considered. Detailed
19 characterisation of differences observable across algorithms is essential for meaningful
20 comparison of findings reported from different tools, and particularly in reconciling any
21 discrepancies therein.

22 2 Methods

23 2.1 Data

24 Data from twenty 3T MRI scanners from the three major manufacturers (GE, Philips,
25 Siemens), each at a different site, were obtained from the Big GABA^{24,25} repository on
26 NITRC, <https://www.nitrc.org/projects/biggaba>. GABA-edited spectra (TR/TE = 2000/68 ms,
27 320 averages, editing at 1.9/7.46 ppm for edit-ON/-OFF, respectively) and corresponding
28 water-unsuppressed reference data (8 or 16 averages) were obtained from a 3 x 3 x 3 cm³
29 voxel in the posterior cingulate region, from 222 consenting adult volunteers (18-36 years of
30 age, approximately even female/male split, having no known neurological or psychiatric

1 illness), in accordance with ethical standards of their respective local institutional review
2 boards (IRB). Subjects consented to the sharing of anonymized data, with allowance for
3 further study. Datasets also included T_1 -weighted structural MR images, which were used for
4 tissue segmentation.

5 This extensive collection of datasets was acquired in an international collaborative
6 study; several aspects have been previously reported ²⁴⁻²⁶, with a focus on comparability
7 across sites and vendors. Full details on the acquisition protocol, software and hardware
8 configurations and sample composition may be found in these papers, and are summarised
9 in Supplementary Table 1. Detailed vendor-specific parameters have been reported
10 previously ^{25,27}.

11 2.2 Processing

12 To the maximum extent practical, data were prepared for each algorithm using a
13 common pipeline, to avoid variations in processing that might otherwise confound
14 observations regarding the model fit. Original data in vendor-specific format were imported
15 using the GannetLoad function from Gannet (v3.1). The GannetLoad function was chosen
16 due to it having the broadest support for the diverse file formats and sequence
17 implementations present in these datasets. This function applies coil combination where
18 necessary, and initial categorisation of individual FIDs into edit-ON/OFF sub-spectra and
19 water reference spectra. Although GannetLoad also incorporates a full processing pipeline,
20 we did not make use of this, instead electing to implement a generalised pipeline in
21 accordance with current consensus recommendations ²⁸, using the processing tools from
22 FID-A ²⁹. The rationale for this was to provide a common, neutral starting point for
23 quantification across all the algorithms to be assessed, rather than one which may have
24 been tuned for a particular quantification algorithm. Additionally, the standard Gannet
25 pipeline performs line-broadening and zero-filling, which invalidates assumptions for error
26 calculations in linear-combination modelling algorithms such as LCModel.

27 Initially, motion-corrupted transients were removed by comparing the root mean
28 square of the difference between each transient and the median of transients in the time
29 domain, rejecting those which differed from the mean by more than four standard
30 deviations. This unlikeness metric was calculated independently within the edit-ON and edit-

1 OFF sub-spectra but applied pairwise. To correct for frequency and phase drift, the
2 remaining FIDs were then aligned in the frequency domain using the spectral registration
3 method ³⁰, iteratively on a variable, restricted frequency range (~1.6-5.5 ppm in the first
4 iteration, reducing to ~1.8-4.0 ppm in subsequent iterations), before averaging within the
5 edit-ON and edit-OFF sub-spectra. Eddy-current correction (ECC) was applied ³¹, before zero-
6 order phase adjustment of each sub-spectrum according to a dual-Lorentzian model for
7 creatine and choline defined in Gannet ¹⁴ and implemented in Osprey. Thereafter, edit-ON
8 and edit-OFF sub-spectra were aligned by spectral registration ³⁰, and sum and difference
9 spectra were calculated arithmetically on the time domain data, dividing by the number of
10 sub-spectra. All resultant spectra and sub-spectra were frequency-shifted such that the main
11 creatine peak (from the dual-Lorentzian model for creatine and choline) in the sum spectrum
12 appeared at 3.027 ppm. After calculation of the difference spectrum, residual water was
13 filtered from edit-ON and edit-OFF data separately with the HSVD method ³². A final
14 automated check was performed to ensure correct ON/OFF ordering and orientation of all
15 resultant spectra, flipping where necessary. Processed data were exported with the same
16 resolution (number of samples and sweep width/dwell time) as the incoming data, without
17 line-broadening or zero-filling. Processing and modelling workflow are summarised in Figure
18 1.

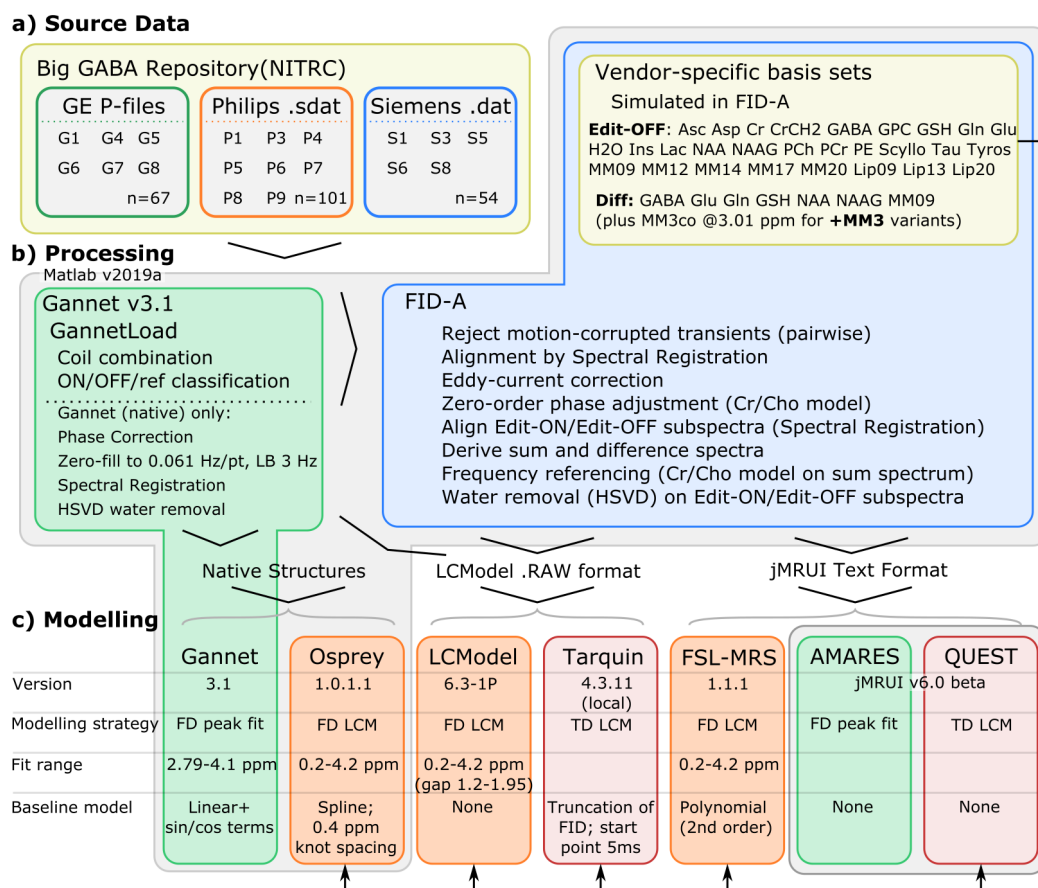


Figure 1: Processing (b) and modelling (c) workflow, summarising key differences between the algorithms assessed.

2.2.1 Quality Control: Processing

Processed spectra were tested against two rejection criteria, designated **R1** and **R2** in subsequent usage:

- R1 captures spectra having strongly aberrant features in the fit range: processing was deemed to have failed if the 0-lag cross-correlation of the normalized, reconstructed frequency domain difference spectrum in the metabolite range (2.6-4.2 ppm) with the normalized mean of all other difference spectra was below 0.5 or differed from the group mean by more than three standard deviations.
- R2 establishes thresholds on basic signal quality metrics: SNR (< 80, defined by maximum peak height around NAA_{diff} in the [1.8, 2.2] ppm interval, over standard deviation over the [-2, 0] ppm range) and linewidth (FWHM > 10 Hz, ³³) measured from NAA_{diff}.

1 Data deemed to have failed at the processing stage were still passed to the fit
2 algorithms but flagged as having failed and excluded from evaluation of groupwise statistics
3 (such as median estimates) in further analysis.

4 2.3 Initial Fit and Quantification

5 Identically processed data were fed into each algorithm. To the maximum extent
6 practical, data were modelled using the developer-supplied default or recommended
7 configuration parameters for GABA-edited MEGA-PRESS data, to yield outcomes
8 representative of those which researchers could expect without extensive local optimisation.

9 Batch processing for all algorithms was automated in Matlab (v2019a), with the
10 exception of the jMRUI-based algorithms for which processed data were exported then
11 processed as batches (grouped by manufacturer and spectral resolution) in a standardised
12 but manual procedure through the jMRUI user interface. As the commonly used default
13 processing pipeline for Gannet incorporates zero-fill and line-broadening factors not present
14 in the standardised pipeline adopted here, we report outcomes both from the standardised
15 processing pipeline (hereafter denoted “Gannet”), and from data processed with Gannet’s
16 own default pipeline, denoted “Gannet (native)”. Tarquin fitting is often performed with an
17 internally simulated basis set; we also assess outcomes from this mode of operation,
18 hereafter denoted “Tarquin (internal)”.

19 Full details on the operation of each method are supplied in the Supplementary
20 Material, section C for quantification of the edited difference spectra. To facilitate
21 concentration scaling to an internal creatine reference, corresponding edit-OFF sub-spectra
22 are also modelled; this is described in the Supplementary Material, section E. Concentration
23 estimates are reported both relative to total creatine (tCr), and with respect to an internal
24 water reference; the complexities and relative merits of each approach are described in ²⁸.

25 As specifics of each algorithm’s water referencing procedure varied considerably,
26 scaling as documented for the respective algorithms was first reversed to yield a raw ratio of
27 signal intensities, before applying tissue-class correction ³⁴ using previously derived tissue
28 fractions ²⁴. Full details on the adjustment for each algorithm are provided in the
29 Supplementary Material, section D. Water-scaled, tissue-class corrected molar concentration

1 estimates are hereafter denoted “/H₂O”; concentration estimates scaled to water with no
2 adjustment for tissue class (assuming pure water concentration per eq(3) of ³⁵) are also
3 calculated, denoted “/H₂O_{noTC}”.

4 2.3.1 Basis Set Preparation and Prior Knowledge

5 All the algorithms examined require some degree of prior knowledge to describe
6 expected spectral features, either in the form of parameter constraints or simulated basis
7 sets. In the present study, prior knowledge was standardised as far as possible: all algorithms
8 requiring a basis set were supplied with the same simulated basis set appropriate to the
9 dataset, whilst both algorithms parameterising individual peaks (Gannet and AMARES) were
10 supplied with similar model parameters.

11 For comparison of the basis set algorithms (FSL-MRS, LCMModel, Osprey, QUEST and
12 Tarquin), a standard simulated basis set specific to each hardware vendor was adopted. As a
13 starting point, vendor-specific basis sets for GABA-edited MEGA-PRESS (TE = 68 ms) that are
14 distributed with Osprey were used; these are derived from fast spatially resolved 2D density-
15 matrix simulations ³⁶ implemented in FID-A using ideal excitation pulses and vendor-specific
16 refocusing pulses and timings, and using chemical shifts and J-coupling coefficients from
17 Kaiser et al. ³⁷. These incorporated metabolite basis functions for GABA, Glu, Gln, glutathione
18 (GSH), NAA, NAAG and a Gaussian component (FWHM = 10.9 Hz) representing co-edited
19 macromolecules around 0.91 ppm (MM09ex).

20 A variation of this basis set was created, incorporating an additional Gaussian
21 component at 3.0 ppm (simulated with FWHM = 14 Hz and scaled intensity equivalent to
22 two protons) to represent co-edited macromolecule signal underlying the GABA peak around
23 3.0 ppm. This component, denoted MM3co, allowed the influence of macromolecule
24 modelling on the various algorithms to be examined; subsequent use of this basis set is
25 annotated with “+MM3”. The interaction of this component with baseline stiffness and soft
26 constraint models similar to those of ^{38–40} is explored for Osprey and LCMModel in a
27 supplementary analysis, section B. All basis set algorithms were run both with and without
28 the MM3 component; in all cases, the reported “GABA+” values include contributions from
29 the underlying macromolecule signal, either explicitly in cases where the MM3 component
30 was modelled (i.e., GABA + MM3co), or implicitly in cases where it was not.

1 2.3.2 Quality Control: Modelling

2 The available quality metrics vary between algorithms; all except Osprey report some
3 form of modelling uncertainty (%SD, %CRLB of metabolite estimates, or % Fit Error for the
4 model), and most report SNR and linewidth of water and/or some metabolite components.
5 As specifics of each algorithm's SNR and linewidth calculation vary, independently derived
6 values are assessed at the processing stage (R2, section 2.2.1). Adopting rather liberal
7 criteria, individual fits were flagged as having failed if either of the following additional
8 criteria were met; in cases where a given metric was not available, the condition is ignored.
9 Criteria below are designated **R3** and **R4** for subsequent usage.

- 10 • R3: %SD, CRLB or FitError for GABA+, Glx_{diff} or tCr_{edit_off} estimate exceeded
11 50% (per ^{33,41}, acknowledging that this strategy must be used with caution ⁴²)
- 12 • R4: Final, scaled estimate for any target metabolite (GABA+/H₂O, Glx_{diff}/H₂O,
13 GABA+/tCr_{edit_off}, Glx_{diff}/tCr_{edit_off}) differing from the median value by more
14 than 5 times the median absolute deviation (MAD) ⁴³ for that algorithm; this
15 was intended to capture any poor fits not flagged by any other criteria.

16 Visual inspection of data, fit outcomes and residuals was also performed, to confirm
17 that no grossly aberrant outcomes eluded the defined rejection criteria. All subsequent
18 analyses are performed after exclusion of individual algorithms' fits (not entire subject
19 datasets) per these criteria.

20 2.4 Statistical Analysis of Modelling Outcomes

21 After batch modelling, statistical analysis was performed using locally implemented
22 scripts written in Python (v3.7.3), using the pandas ⁴⁴ (v0.23.3) data analysis framework, with
23 numeric methods from NumPy ⁴⁵, and statistical methods from the SciPy ⁴⁶ (v1.1.0), pingouin
24 ⁴⁷ and statsmodels ⁴⁸ (v0.12.1) libraries.

25 Scaled estimates for target metabolite, grouped by algorithm, were tested for
26 normality using the Shapiro-Wilk method ⁴⁹, and for comparable distribution of variance
27 between algorithms by the Fligner-Killeen's test ⁵⁰, both implemented in SciPy. Limits of
28 agreement between pairs of algorithms were derived, along with their 95% confidence

1 intervals, in accordance with the Bland-Altman method ⁵¹. Estimates grouped by algorithm
2 and manufacturer were compared using Welch's t-test ⁵², with Holm-Bonferroni correction
3 ^{53,54} for multiple comparisons. An adjusted p-value less than 0.05 was considered significant.

4 An unconditional linear mixed-effects model was fit to water-referenced GABA+
5 estimates, using R version 3.5.3 ⁵⁵ with the *lme4* package ⁵⁶ and an implementation derived
6 from ²⁵. Vendor, site, algorithm and subject factors were incorporated into the fit, to
7 calculate variance partition coefficients (VPCs) estimating the proportion of total variance
8 attributable to each factor. Significance testing was performed using chi-square likelihood
9 ratio tests, against a null hypothesis simulated by parametric bootstrapping (2000
10 simulations) ⁵⁷. The Big GABA dataset described herein has previously been assessed with
11 respect to demographics and signal quality ²⁴.

12 For each metabolite of interest, a global median was calculated across all subjects
13 and all algorithms. Subsequently, estimates grouped by site and algorithm were linearly
14 scaled to match the global median, thereby removing broad scaling differences observed
15 between certain sites, vendors and algorithms which would otherwise bias inter-algorithm
16 correlations.

17 The degree of correlation between MRS-detected GABA estimates and voxel tissue
18 fraction has been shown to be an effective index of GABA estimation accuracy ³⁵, given the
19 differing GABA concentrations between grey and white matter ⁵⁸⁻⁶⁰. Building on this
20 approach, robust Spearman correlation coefficients between voxel grey matter fraction and
21 GABA+/H₂O_{noTC} estimates were calculated using the 'skipped' method ^{61,62}, implemented in
22 pingouin, to exclude bivariate outliers. To determine whether correlation coefficients
23 obtained for any one of the algorithms differed significantly from the correlation obtained
24 across all algorithms, z-scored coefficients were compared (two-tailed), with a significance
25 level of $p_{\text{holm}} < 0.05$; similarly, z-scored coefficients were compared between variants without
26 and with the MM3 macromolecule component in the basis set.

27 Finally, intraclass correlation coefficients (ICCs) were calculated between all
28 algorithms, separately without and with the MM3 component for basis set algorithms, and
29 between pairs of algorithms, using a two-way mixed-effects model for single-rater
30 consistency (ICC(3,1) implemented in pingouin).

3 Results

3.1 Fit and Residuals

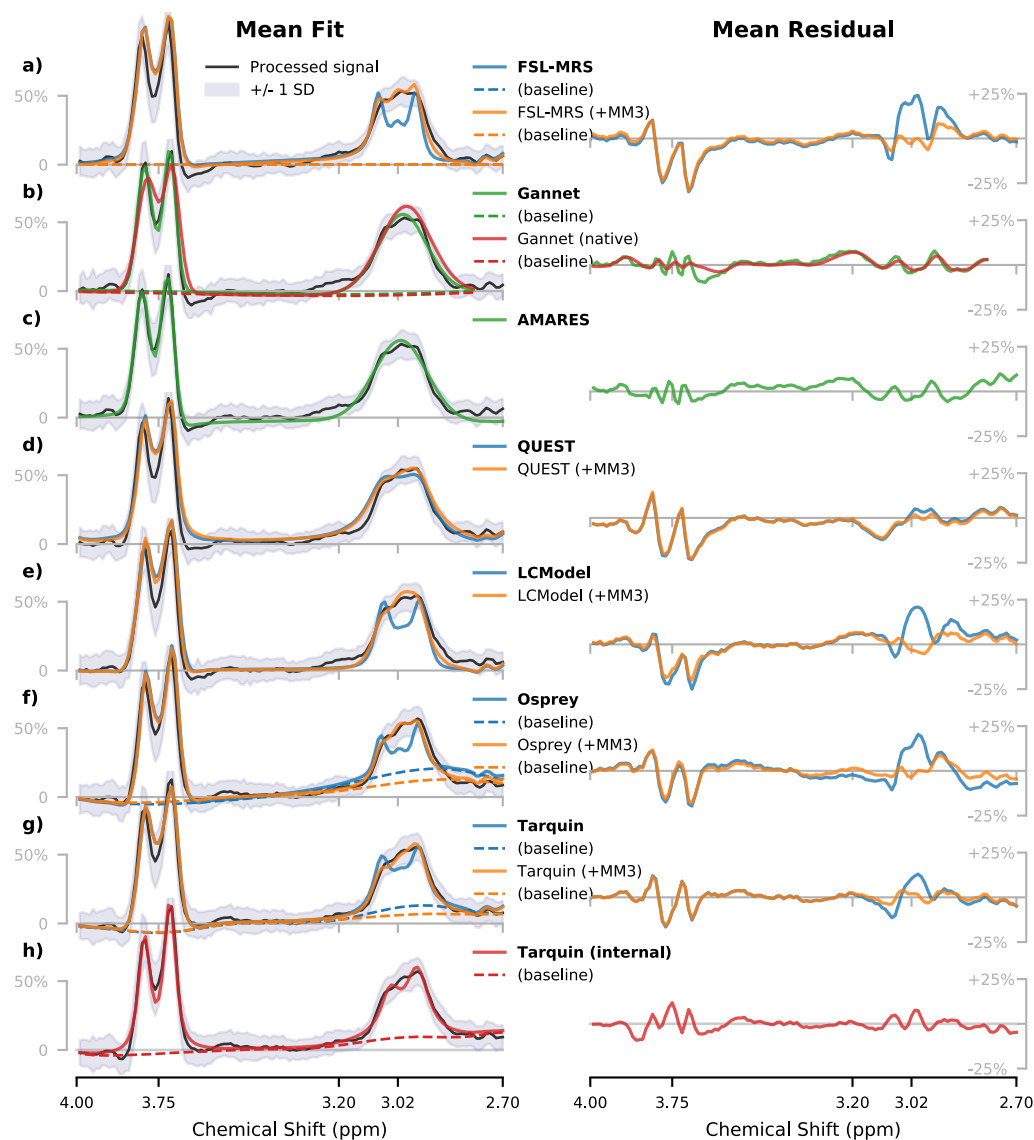
Mean fit outcomes for each algorithm are presented in Figure 2. Note that all basis set algorithms without MM3, except QUEST, show strong residuals in the 3 ppm range; in the case of Osprey (Figure 2f), this appears to have a strong influence on baseline in the vicinity. Variants with MM3 show generally reduced residuals in that range, indicating that the inclusion of a dedicated MM3 basis function leads to more appropriate modelling of the data.

A characteristic hump around 3.2 ppm is handled differently by the various algorithms: peak fitting algorithms Gannet and AMARES, (Figure 2b,c) are largely unperturbed, QUEST (Figure 2d) envelopes the entire signal with broader 3.0 ppm peak, while other algorithms fall somewhere in between.

A notable difference between algorithms arises from the differences in baseline estimation practices. While AMARES, QUEST, and LCModel do not include a baseline term in their default settings for MEGA-PRESS and Gannet and FSL-MRS adopt relatively stiff, low-order models, both Tarquin and Osprey attribute a considerable fraction of the edited 3-ppm signal to the baseline. This tendency is mitigated upon the inclusion of the MM3 model.

Finally, there is a distinct pattern to the residuals around the Glx peaks from all basis set algorithms, not present in the fits applying simple peak fitting on a restricted frequency range (Gannet, AMARES).

A summary of basic quality metrics from the fit spectra is presented in Supplementary Figure 6, along with the number of spectra rejected according to the defined criteria (per sections 2.2.1 and 2.3.2).



1
 2 *Figure 2 Average metabolite and baseline (where applicable) models with corresponding residuals for the GABA+ edited*
 3 *spectra, for each algorithm. Vertical scaling is normalised; outcomes over the full fit range are presented in Supplementary*
 4 *Figure 8; outcomes split by vendor are presented in Supplementary Figure 9.*

5
 6 **3.2 Statistical Analysis**

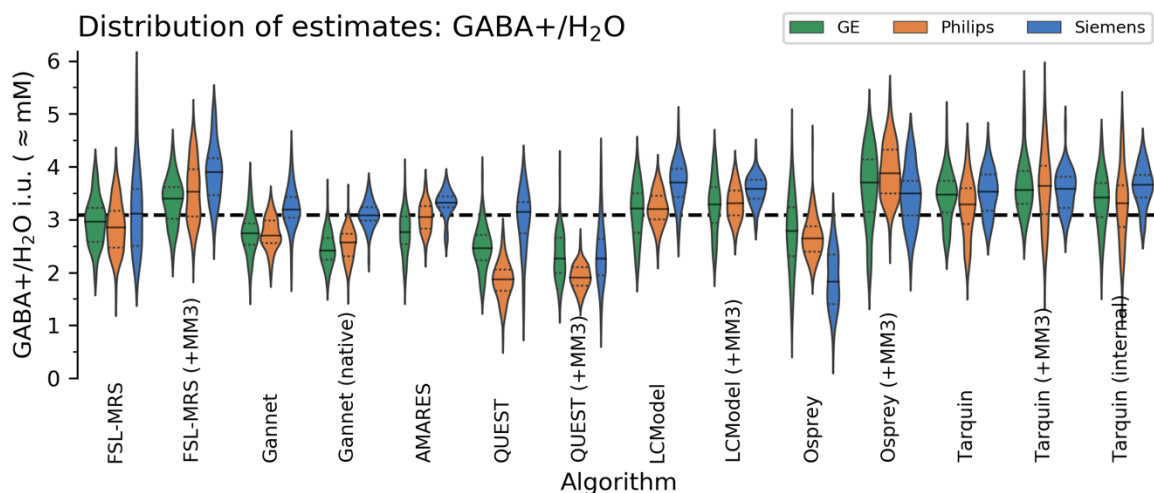
7 Shapiro-Wilk testing and subsequent inspection of quantile-quantile (Q-Q) plots
 8 indicated that while concentration estimates from most algorithms satisfied the assumption
 9 of a normal distribution, several (predominantly $Gl_{x_{diff}}/tCr$ estimates) deviated slightly from
 10 this. Fligner-Killeen tests revealed mismatched variances between several sets of estimates
 11 (predominantly relating to QUEST and LCModel +MM3), which motivated the subsequent
 12 adoption of Welch's t-test for groupwise comparisons.

1 3.2.1 Water-referenced concentration estimates

2 Comparisons between algorithms for GABA+/H₂O are summarised in Figure 3, with
3 full details in Supplementary Table 6, and Bland-Altman plots describing limits of agreement
4 in Supplementary Figure 10. The global median estimate for GABA+/H₂O across all
5 algorithms and subjects was found to be 3.2 ± 0.4 i.u..

6 For several quantification algorithms, water-referenced estimates for GABA+/H₂O
7 were found to be significantly higher from Siemens datasets than from other manufacturers,
8 by a factor of 9-17% ($p_{\text{holm}} < 0.001$, depending on the algorithm) for FSL-MRS (+MM3),
9 Gannet, Gannet (native), AMARES and LCModel, increased to 41% ($p_{\text{holm}} < 0.001$) for QUEST.
10 Osprey gave significantly lower estimates for Siemens datasets (-28.0%, $p_{\text{holm}} < 0.001$). QUEST
11 (and the +MM3 variant) gave significantly lower estimates for Philips datasets (-16.1% and -
12 7.7% respectively, $p_{\text{holm}} < 0.001$), and AMARES and Gannet (native) gave lower estimates for
13 GE datasets (-9.5%, $p_{\text{holm}} < 0.01$ and -8.5%, $p_{\text{holm}} < 0.05$ respectively). Median GABA+ estimate
14 across all algorithms was 5.8% higher for Siemens sites ($p_{\text{holm}} < 0.01$). All differences are
15 expressed relative to the mean across all subjects for the respective algorithm. No other
16 variants showed significant effects.

17 Water-referenced Glx_{diff} estimates from all algorithms were significantly higher for
18 Siemens sites: median Glx_{diff}/H₂O across algorithms +15.7%, $p_{\text{holm}} < 0.001$ relative to group
19 mean. Estimates from Philips sites were somewhat lower (-10.1%, $p_{\text{holm}} < 0.01$).



20

21 Figure 3: Distribution of GABA+/H₂O estimates from each algorithm, grouped by manufacturer. Global median is shown in
22 dashed black.

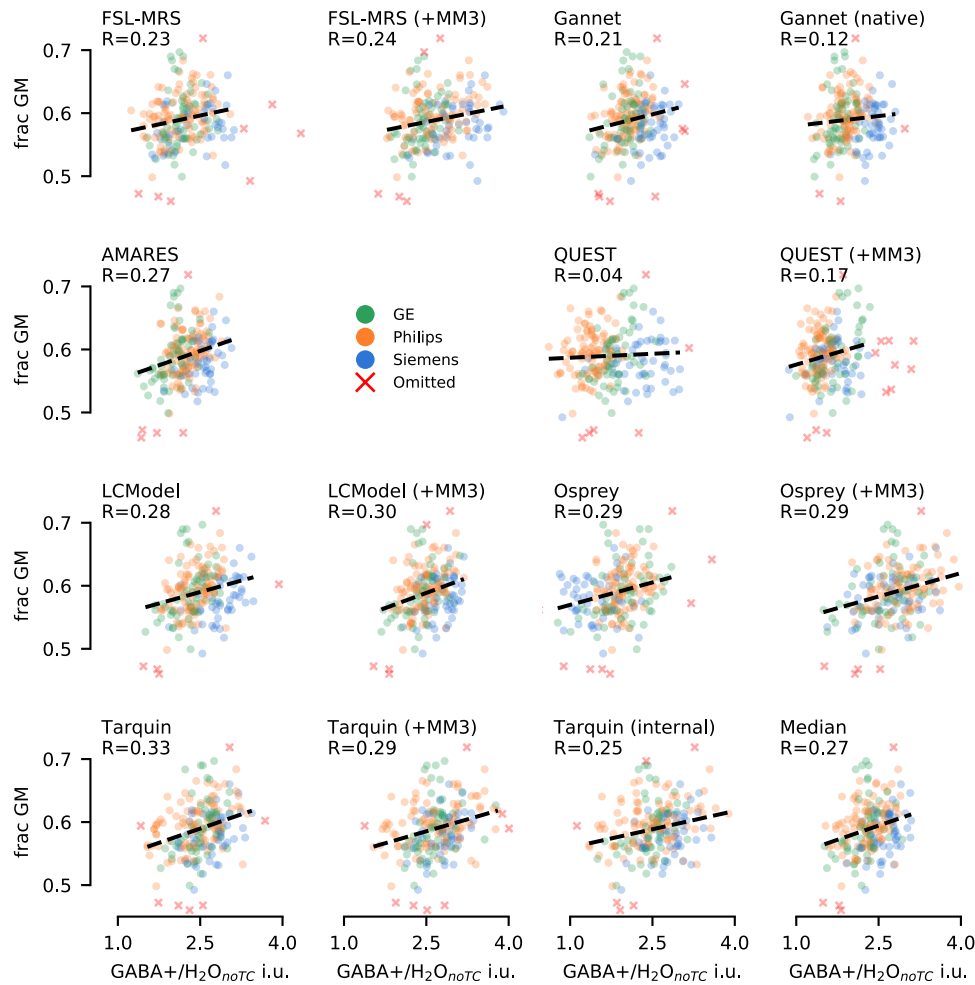
1 For data fit without the explicit MM3 component, the unconditional linear mixed-
2 effects model yielded VPCs of [33.8, 16.4, 6.4, 4.0%] for algorithm, site, subject and vendor
3 factors respectively. In this context, the “subject” factor reflects systematic *within*-subject
4 variation in estimates, while the residual 39.4% accounts for inherent, systematic *between*-
5 subject variation, as well as any other variance which could not be accounted for in the
6 model. Parametric bootstrap testing showed all factors to be significant ($p_{\text{holm}} < 0.001$).

7 3.2.2 Metabolite-referenced concentration estimates

8 Estimates for GABA+/tCr_{edit_off} were consistently higher for GE datasets (+17.3%
9 across algorithms, $p_{\text{holm}} < 0.001$) and lower for Siemens datasets (-14.3%, $p_{\text{holm}} < 0.001$).
10 Glx_{diff}/tCr_{edit_off} ratios were higher in GE datasets (21.9%, $p_{\text{holm}} < 0.001$) and slightly lower in
11 Philips (-5.0%, $p_{\text{holm}} < 0.05$) and Siemens (-6.7%, $p_{\text{holm}} < 0.05$) datasets. As in section 3.2.1,
12 differences are quoted relative to the mean estimate across all subjects for the respective
13 algorithm. All these trends presented similarly across all modelling algorithms, albeit with
14 varying magnitudes and significance levels.

15 3.2.3 Grey matter volume fraction correlation

16 The relationship between estimated GABA+ and grey matter volume fraction is
17 reported in Figure 4, as an index of estimation accuracy. The accuracy of QUEST (without
18 MM3) was found to be significantly below that of other algorithms ($p_{\text{holm}} < 0.01$), while the
19 QUEST +MM3 variant performed comparably with other algorithms. Otherwise, slight
20 differences observable between algorithms were not statistically significant, and no
21 particular trend is evident between algorithm variants without and with MM3 components.



1

2 *Figure 4 Relationship between GABA+ and grey matter, with different modelling strategies for GABA+. Robust (skipped)*
 3 *correlation coefficients are reported, with line-of-best-fit in dashed black.*

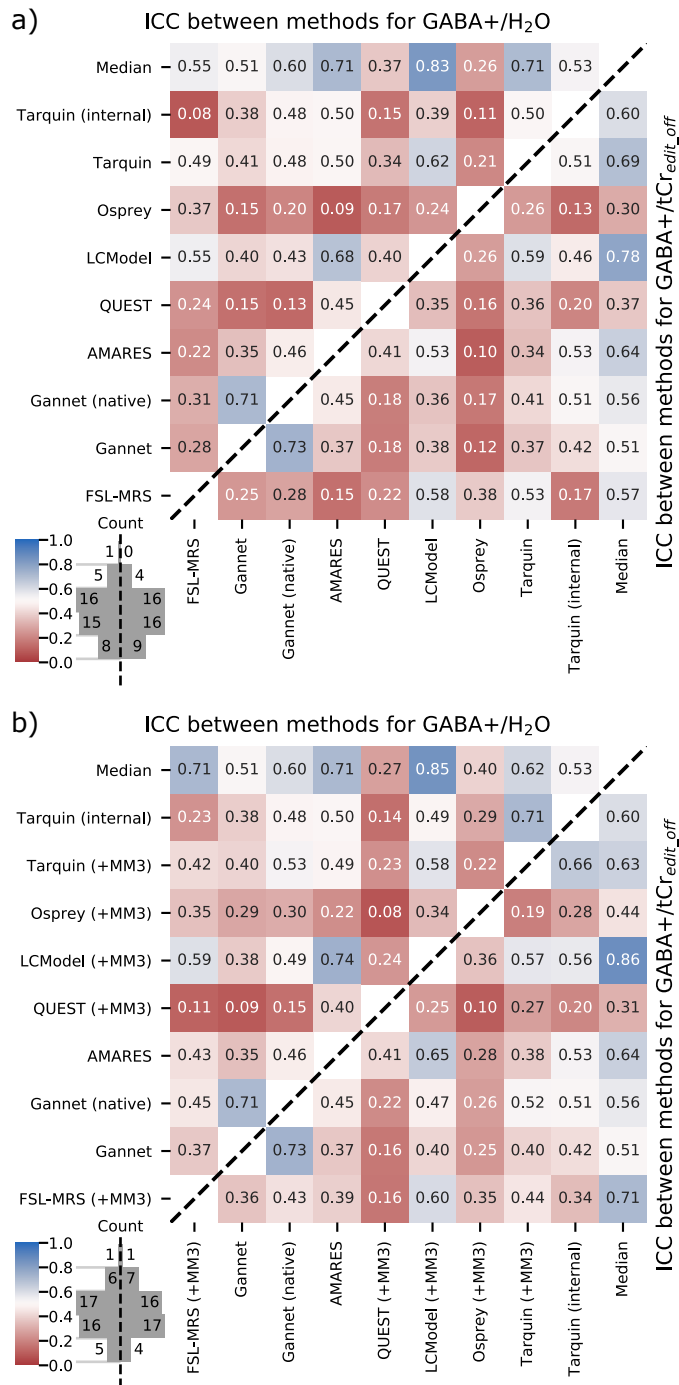
4 3.2.4 Correlational analysis

5 ICC (single-rater consistency) for GABA+ across all algorithms was 0.38 (95% CI 0.32-
 6 0.44) without the MM3 component for basis set algorithms, and increased to 0.44 (95% CI
 7 0.39-0.5) with MM3 included, supporting that the inclusion of this dedicated component is
 8 warranted. ICCs between all pairs of algorithms are presented in Figure 5. For fits performed
 9 without the MM3 component, GABA+/H₂O estimates showed moderate correlation between
 10 most algorithms (typically on the range $r=0.4-0.6$; slightly lower when referenced to
 11 tCr_{edit_off}). Correlations for AMARES, LCMoDel and Tarquin were significantly stronger
 12 ($p_{holm}<0.01$) than the group mean, those for QUEST and Osprey somewhat lower. Inclusion
 13 of an MM3 basis set component generally improved concordance with other algorithms for
 14 FSL-MRS ($p_{holm}<0.001$) and Osprey, the latter at trend level. However, both time domain

1 basis set algorithms (QUEST and Tarquin) showed reduced concordance (at trend level) upon
2 inclusion of the MM3 component.

3 ICCs for additional metabolites and ratios are presented in Supplementary Figure 12;
4 $\text{Gl}_{\text{diff}}/\text{H}_2\text{O}$ estimates from the edited spectrum correlated more strongly between
5 algorithms (typically on the range $r=0.6-0.8$, slightly lower when referenced to $t\text{Cr}_{\text{edit_off}}$).

6



1

2 *Figure 5* Intraclass correlation coefficients between algorithms, scaled to water (upper left triangle) and tCr_{edit_off} (lower right

3 triangle), with basis set algorithms excluding (a) and including (b) a component representing co-edited macromolecule

4 contribution. "Median" data denotes correlation with the median estimate across all algorithms.

5

4 Discussion and Conclusions

4.1 Quality Control

Basic signal quality metrics (such as SNR and linewidth) and reliability-of-fit estimates (%CRLB, %Fit Error) are often used as the basis for rejecting poor fits. However, as seen in Supplementary Figure 6, these are often not sufficient. Whilst four datasets were deemed to have failed at the processing stage (R1), yielding output barely recognisable as GABA-edited difference spectra, all algorithms “successfully” fit some of these (Supplementary Figure 7), with quality metrics that satisfied all other criteria. We therefore repeat the observation that simply filtering results based on these basic signal metrics is inadequate as a means of quality control; the metrics themselves may have limited reliability, particularly in cases where the model does not accurately reflect the experimental data⁶³. Consideration must also be given to the shape of the data, fit and residuals themselves – either by algorithmic assessment or, if feasible, visual inspection.

4.2 Scale Differences between Manufacturers

Previous studies, including⁶⁴, have explored systematic differences in reported GABA+ estimates between different manufacturers, and their relation to GABA editing efficiency and the contribution of co-edited macromolecules to the measured signal. Furthermore, a previous examination of water-scaled GABA+ estimates on a superset of data also incorporating the present study’s subjects²⁴ identified systematically higher GABA+ estimates from datasets acquired on the Siemens platform, by approximately 29%, which could not be explained in terms of editing efficiency or macromolecule contribution. Observations in the present analysis corroborate this, with most algorithms yielding higher GABA+/H₂O estimates from Siemens datasets, and all algorithms yielding higher estimates for both Glx_{diff}/H₂O and tC_{redit_off}/H₂O for Siemens datasets. The fact that this trend is seen across different metabolites and within the edit-OFF sub-spectra gives support to the notion that water reference data from the Siemens implementation may not be optimal for scaling purposes, although it cannot be ruled out that both GE and Philips share a similar mis-scaling.

4.3 Macromolecule Fitting

The observed signal around 3 ppm includes substantial contamination (around 50%⁶⁵) from co-edited signals, including homocarnosine⁶⁶ with coupling between 1.89 and 3.00 ppm spins, and poorly-characterised components tentatively attributed to lysine-containing macromolecules^{6,67,68}, with coupling between 1.71 and 3.01 ppm spins. This may give rise to large residuals or biased baseline estimates if not considered in the model. Although their impact has been studied by some^{38–40,69,70}, there is currently no consensus on how these should be handled.

The default LCMoel configuration (flat baseline) performs surprisingly well by simply ignoring such signals, giving rise to characteristic peaks in the residuals. Meanwhile, the case of Osprey (without MM3co component, Figure 2f) exemplifies the potential baseline distortions and correspondingly reduced GABA+ estimates resulting from these residuals. The result of our ICC analysis across all algorithms suggests that more consistent GABA+ estimates may be obtained by explicitly parametrizing the MM3 contribution in the model. However, using correlation between GABA+ estimates and grey matter fraction as a benchmark, incorporation of the MM3 component did not significantly impact the effectiveness of individual algorithms in measuring differences in GABA+ levels. Moreover, while the supplementary analysis (supplementary section B) suggested improved effectiveness for LCMoel after incorporation of soft constraints on MM3 amplitude (whether to GABA or MM0.9), similar performance for this algorithm was obtained by simply modelling a stiff but non-zero baseline, allowing this to absorb some of the MM3 signal. This configuration was also found to be effective for modelling GABA alone, consistent with previously reported findings⁶⁹ wherein modelling a more flexible baseline in LCMoel to selectively remove a portion of the MM3 contribution allowed for closer measurement of GABA rather than GABA+. Estimates from Osprey in similar configurations were comparable, and it has been shown that higher degrees of baseline flexibility cause greater fractions of the 3-ppm signal to be absorbed into the baseline⁴⁰.

Peak-fitting algorithms (such as Gannet and AMARES) may circumvent this issue somewhat by considering the entire 3.0 ppm GABA+MM signal with a single broad Gaussian model as examined herein; this approach performed comparably with more elaborate

1 models. Furthermore, Tarquin with its internally simulated basis set has two separate
2 Gaussian components representing the GABA+ signal, which are seen to shift and broaden to
3 conform to the shape of the observed GABA+MM signal (Figure 2h). QUEST, similarly,
4 broadened the GABA basis function substantially to more closely envelop the entire
5 GABA+MM signal and (unfortunately) adjacent artefacts (Figure 2d), perhaps accounting for
6 its somewhat lower correlation with grey matter fraction and agreement with other
7 algorithms.

8 4.4 Artefact Rejection around 3.2 ppm

9 MEGA-edited GABA spectra often exhibit a slight artefactual feature around 3.2 ppm,
10 which can be problematic for fitting algorithms. The origin is uncertain, but potentially
11 related to incomplete subtraction of choline⁷¹ or contribution from undetermined other co-
12 edited signals (such as, perhaps, valine-containing macromolecules⁶⁷ or arginine⁷²). In the
13 present study, the baseline for the Osprey fit (without MM3co) tends to respond to this
14 artefact, inducing a bend in the baseline which appears to cut out a significant part of the
15 real GABA peak (Figure 2f), leading to a likely underestimation of GABA+ area. QUEST
16 appears to broaden the GABA and/or MM3co basis components, incorporating the artefact
17 into the GABA+ estimate and most likely over-estimating the GABA+ signal area (Figure 2d);
18 this effect was most pronounced for Siemens datasets (see Supplementary Figure 9d), where
19 the feature manifests more prominently. FSL-MRS and LCModel both handle the artefact
20 well in the general case, largely rejecting it from both the baseline and metabolite models
21 (Figure 2a,e); this is likely owed to the fact that their default MEGA-PRESS settings prescribe
22 a low-order polynomial baseline (FSL-MRS) or no baseline at all (LCModel). Whilst other
23 basis set algorithms end up somewhere in between (with a degree of contamination from
24 the artefact), peak-fitting algorithms AMARES and Gannet both perform well in this area.
25 Indeed, Gannet (Figure 2b) is the only algorithm to explicitly deal with this artefact, down-
26 weighting some residuals in this region. We suggest that comparably rigid baseline
27 estimation as well as incorporating a Gaussian basis component around 3.2 ppm, with tight
28 constraints on linewidth, shift and amplitude to avoid inadvertently fitting part of the GABA
29 peak, may yield some benefits in this area for other algorithms. Ultimately, further
30 investigation into the underlying signal, and more complete profiling of the co-edited
31 metabolite and macromolecule signals in the region would be preferable.

1 4.5 Glx

2 Although quantification of Glx from the difference spectrum has been demonstrated
3 to be reliable given suitable quality constraints ⁷³, several researchers have highlighted the
4 relatively low concordance between estimates from short-TE PRESS spectra and GABA-
5 edited difference spectra, with estimates from the edit-OFF sub-spectra often found to
6 agree better with the short-TE PRESS ^{74–76}; this is unsurprising given that the Glx signal in the
7 short-TE and edit-OFF sub-spectra are subject to similar underlying uncertainties, including
8 MM background and overlapping GSH and aspartyl signals. In comparing Glx quantification
9 between MEGA-PRESS difference and edit-OFF sub-spectra, recent studies ^{76,77} report a
10 correlation around $r=0.8$; results in the present study show a more moderate correlation,
11 between $r=0.34$ and 0.69 depending on algorithm: see Supplementary Figure 13; a linear
12 scaling factor is also observed, consistent with recent findings ⁷⁸. It is notable that agreement
13 between algorithms is higher for co-edited Glx than for GABA+, reflecting the better-defined
14 signal seen in the difference spectrum (Figure 2).

15 With reference to Supplementary Figure 8, all basis set algorithms showed a distinct
16 structure in the residuals around 3.7 ppm, with the model peaks appearing a little to the
17 right of the peaks observed in the data. This is most likely due to the complicated signal
18 patterns around 2.3 ppm in the edited spectrum (resulting from overlapping signals of GABA
19 and co-edited Glu, Gln and GSH), which interact critically with the 3.75 ppm modelling. It is
20 likely that there is a poorly understood baseline fluctuation arising from co-edited
21 macromolecular signals appearing between 1.5 and 2.5 ppm ³⁹, which will bias the correct
22 phase estimation of the 2.25 ppm signals, at the expense of getting the phase of the related
23 3.75 ppm signals right. The peak-fitting algorithms tested, where modelling around 3.7 ppm
24 is not bound to features in other parts of the spectrum (such as around 2.3 ppm), show
25 much lower residuals in the region. It is possible that basis set fitting on a constrained range
26 would mitigate this effect, at the expense of throwing away useful spectral information and
27 hence detracting from the utility of the basis set approach in general. A model which shares
28 lineshape information between the 2.3 and 3.7 ppm Glx peaks but allows a tightly-
29 constrained frequency shift between them may present a reasonable alternative.

1 4.6 Limitations

2 The basis set adopted in the present study was simulated with ideal excitation pulses,
3 and therefore may not fully model subtle variations in spectral structure between
4 manufacturers. However, the impact of excitation is likely negligible compared to the impact
5 of refocusing, which is appropriately accounted for in the 2D simulations. Vendor-specific
6 excitation pulses may also contribute to subtly different shape and asymmetry of the 3.0
7 ppm peak and varying manifestation of the 3.2 ppm feature, which may be observed in
8 Supplementary Figure 9.

9 Whilst the present analysis examines a variety of commonly used implementations
10 representing a range of modelling strategies for GABA-edited spectroscopy data, we note
11 that several other algorithms and implementations are also available to the MRS
12 community, including AQSES ⁷⁹, INSPECTOR ⁸⁰, KALPANA ⁸¹, OXSA ⁸², spant ⁸³ and Vespa ⁸⁴.

13 Furthermore, many of the packages examined offer extended functionality which
14 may well lead to improved performance in certain circumstances, but this did not align with
15 our approach of adopting recommended/default configurations for all algorithms. Most
16 significantly, many software packages offer the fine-tuning of several aspects of the
17 modelling process, for instance the baseline parametrization. As further examples, jMRUI
18 QUEST offers flexible baseline modelling strategies; FSL-MRS offers independent shift groups
19 which were not assessed; Osprey can additionally simultaneously optimise difference and
20 sum spectra, potentially benefiting from additional spectral information and improved SNR.
21 Both peak-fitting algorithms examined are flexible in their choice of model, with (for
22 example) dual-Gaussian models for the GABA+ signal readily available. An inevitable
23 consequence of adopting default settings in this analysis is that these might not be optimal
24 for data processed with the standard pipeline adopted herein. All the tools tested are highly
25 configurable and offer expert users many possibilities to tune performance optimally for
26 particular datasets, offering the potential for further invention and protection against
27 establishing a possibly incorrect orthodoxy. Nonetheless, this flexibility comes with the
28 caveat that it may also lead to misuse, selective reporting and inappropriate modelling.
29 Moreover, such variability runs counter to efforts to standardize analysis methodology ¹.
30 While more research into optimized modelling strategies is needed to improve the

1 comparability and robustness of MEGA-PRESS fitting, this work highlights that the complete
2 and accurate reporting of all decisions made during analysis and modelling is immensely
3 important, going even beyond the recently published minimum reporting standards for MRS
4 ⁸⁵.

5 The present study examines metabolite estimates for each algorithm, relative to
6 water and creatine references obtained using that same algorithm – this reflects typical
7 usage, but means that variations discussed herein are not necessarily driven purely by the
8 metabolite modelling.

9 Finally, since the findings documented herein are substantiated entirely by in-vivo
10 data, there is no “ground truth” available by which to directly assess the algorithms’
11 accuracy. Whilst this limitation has been partially addressed by considering the strength of
12 correlation between GABA+ estimates and grey matter fraction as an index of relative
13 accuracy ³⁵, a more direct assessment of accuracy could be achieved in further studies
14 involving carefully prepared phantom or synthetic data ⁸⁶, each approach having its own
15 inherent limitations. In either case, meticulous attention to macromolecule baseline, SNR
16 and line shape would be required to ensure transferability of findings to in-vivo applications.

17 4.7 Key Recommendations

18 Based on these findings, we recommend the following for future studies:

- 19 • When applying basis set modelling approaches, special consideration must be
20 given to the co-edited macromolecular signal underlying the 3.0 ppm GABA
21 peak. While appropriate modelling outcomes may be obtained in some cases
22 by entrusting a carefully tuned baseline to capture the entirety of the signal ⁴⁰,
23 a more generalisable approach is simply to routinely incorporate a simulated
24 basis component to represent this signal.
- 25 • Care must be taken to ensure consistent behaviour in the presence of
26 commonly observed artefacts, such as the signal around 3.2 ppm. This
27 artefact could be explicitly incorporated into the model, or mitigated with a
28 rigid baseline model, which is less likely to follow the local signal curvature.

- 1 • More generally, when inspecting fit outcomes the behaviour of the baseline
2 (where modelled) demands close attention, to ensure that it does not unduly
3 bias modelling of the GABA peak.

- 4 • When assessing data acquired at different sites, systematic differences in
5 scale are to be expected and must be considered, regardless of the algorithm
6 applied.

7 Additionally, we propose four key areas for further systematic investigation in future
8 studies:

- 9 1. Robust methods for the generation of large synthetic datasets for validation
10 are necessary to facilitate direct assessment of modelling accuracy. These
11 synthetic data need to be truly representative of in-vivo data, hence the
12 design of their underlying physical signal models will require great attention
13 to detail. Subsequent interpretation must bear in mind that the outcome is
14 likely to be determined by the degree of similarity between the physical
15 model used to generate the data and the model used to decompose it during
16 linear-combination fitting.

- 17 2. Detailed exploration of the co-edited macromolecule profile which underlies
18 typical GABA-edited data is required. Metabolite-nulled edited spectra
19 obtained on different hardware platforms could provide the basis for a more
20 informed parametrization of these signals during modelling, and also
21 contribute to accurate and in-vivo-like representation of synthetic data.

- 22 3. The origin of the spectral feature around 3.2 ppm requires further
23 investigation: it needs to be determined whether this is a subtraction artefact
24 or an actual real signal (for example from valine-containing macromolecules
25 or other signals hitherto not routinely included in modelling, such as arginine
26 which has a compatible spin system). Clarification of this feature would allow
27 for more appropriate modelling and parametrization of synthetic data.

1 4. When considering Glx estimates from the difference spectrum, further
2 investigation into the complex interactions of co-edited Glu, Gln, GSH and
3 possibly other signals around 2.3 ppm, and their impact on the 3.75 ppm Glx
4 modelling, would be beneficial. This focus area will benefit from increased
5 insight into the macromolecular background of the GABA-edited spectrum.

6 4.8 Conclusions

7 Although the observed consistency across algorithms was generally moderate, with
8 pairwise correlation in some cases concerningly weak, we emphasize that more *consistent*
9 estimates are not necessarily more *accurate* estimates: all the algorithms tested (except for
10 QUEST without MM3) were shown to be comparable in their effectiveness in detecting
11 differences in GABA+ concentration. This does however raise some concerns regarding the
12 comparability of findings between different studies, each of which will typically employ a
13 single modelling algorithm, often with divergent processing and prior knowledge and with
14 significantly smaller sample sizes than tested here. This means that the choice of analysis
15 already adds considerable uncertainty and variability.

16 Improved standardisation of sequence implementation ^{27,87}, and adoption of
17 standardised processing pipelines and prior knowledge (e.g. in the form of publicly available
18 basis set libraries) may reduce sources of variation between studies. However, the
19 interaction of baseline, co-edited macromolecule and metabolite signals, and other
20 artefactual signals remains a critical source of variation between algorithms, and within
21 different configurations of the same algorithm. Better characterisation of these signals
22 would allow for more complete modelling (at the risk of over-fitting). Consensus on optimal
23 (or at least, appropriate) control parameters for the respective algorithms would also be
24 beneficial. It may further be of benefit to conceive a “consensus algorithm” to be
25 implemented across software environments, and used as a shared starting point to refine
26 the algorithmic decision making in future iterations of the algorithm. Meanwhile, careful
27 attention to the behaviour of the model with regards to such signals, and rigorous reporting
28 of the configuration employed, are necessary to facilitate meaningful comparisons between
29 studies.

1 5 Acknowledgements

2 The authors wish to thank Dr Stephen Provencher for his assistance and constructive
3 feedback on the application of the LCModel algorithm.

4 Analysis was performed within a project funded by ERC grant #249516, which
5 additionally supports the contributions of ARC, LE, KH. Data used in this analysis were
6 previously collected through an international collaborative study funded under NIH grant
7 R01 EB016089.

8 WTC is supported by funding from Wellcome Trust and the Royal Society
9 (102584/Z/13/Z). PM thanks India –Australia Strategic Biotechnology Funding (BT/Indo-
10 Aus/10/31/ 2016/PKM). RR was supported by the Australian Research Council
11 (DE210100790). GO and RE acknowledge funding support from K99/R00 AG062230, S10
12 OD021648, P41 EB031771, P41 EB015909, R01 EB016089, R01 EB023963, R01 EB028259 and
13 R21 HD100869.

14 Graphical abstract illustrated by Laura Garrison (University of Bergen)

15 6 Declaration of Interest

16 The authors declare no conflicting interests.

17 7 Data Availability Statement

18 Scripts used for automation and reporting contained in the present manuscript are
19 publicly available here; further dependencies are described within:

20 <https://git.app.uib.no/bergen-fmri/analyzing-big-gaba>

21 Spectra analysed in this manuscript were obtained from the publicly available Big
22 GABA repository on NITRC, https://www.nitrc.org/projects/big_gaba

23 Basis sets used in the primary analysis were obtained from the publicly available
24 Osprey package, <https://schorschinho.github.io/osprey>

25

8 References

1. Bhogal AA, Schür RR, Houtepen LC, et al. 1 H-MRS processing parameters affect metabolite quantification: The urgent need for uniform and transparent standardization. *NMR Biomed.* 2017;30(11):e3804. doi:10.1002/nbm.3804
2. Kanowski M, Kaufmann J, Braun J, Bernarding J, Tempelmann C. Quantitation of simulated short echo time 1H human brain spectra by LCMoDel and AMARES. *Magn Reson Med.* 2004;51(5):904-912. doi:10.1002/mrm.20063
3. Mullins PG, Rowland L, Bustillo J, Bedrick EJ, Lauriello J, Brooks WM. Reproducibility of 1H-MRS measurements in schizophrenic patients. *Magn Reson Med.* 2003;50(4):704-707. doi:10.1002/mrm.10598
4. Zöllner HJ, Považan M, Hui SCN, Tapper S, Edden RAE, Oeltzschner G. Comparison of different linear-combination modeling algorithms for short-TE proton spectra. *NMR Biomed.* 2021;34(4). doi:10.1002/nbm.4482
5. Mescher M, Merkle H, Kirsch J, Garwood M, Gruetter R. Simultaneous in vivo spectral editing and water suppression. *NMR Biomed.* 1998;11(6):266-272. doi:10.1002/(sici)1099-1492(199810)11:6<266::aid-nbm530>3.0.co;2-j
6. Rothman DL, Petroff OA, Behar KL, Mattson RH. Localized 1H NMR measurements of gamma-aminobutyric acid in human brain in vivo. *Proc Natl Acad Sci.* 1993;90(12):5662-5666. doi:10.1073/pnas.90.12.5662
7. Brix MK, Erslund L, Hugdahl K, et al. Within-and between-session reproducibility of GABA measurements with MR spectroscopy. *J Magn Reson Imaging.* 2017;46(2):421-430.
8. Duda JM, Moser AD, Zuo CS, et al. Repeatability and reliability of GABA measurements with magnetic resonance spectroscopy in healthy young adults. *Magn Reson Med.* Published online November 20, 2020:mrm.28587. doi:10.1002/mrm.28587
9. O’Gorman RL, Michels L, Edden RA, Murdoch JB, Martin E. In vivo detection of GABA and glutamate with MEGA-PRESS: Reproducibility and gender effects. *J Magn Reson Imaging.* 2011;33(5):1262-1267. doi:10.1002/jmri.22520
10. Baeshen A, Wyss PO, Henning A, et al. Test–Retest Reliability of the Brain Metabolites GABA and Glx With JPRESS, PRESS, and MEGA-PRESS MRS Sequences in vivo at 3T. *J Magn Reson Imaging.* 2020;51(4):1181-1191. doi:10.1002/jmri.26921
11. Zöllner HJ, Oeltzschner G, Schnitzler A, Wittsack HJ. In silico GABA+ MEGA-PRESS: Effects of signal-to-noise ratio and linewidth on modeling the 3 ppm GABA+ resonance. *NMR Biomed.* Published online September 28, 2020:e4410. doi:10.1002/nbm.4410
12. Mikkelsen M, Bhattacharyya PK, Mandal PK, et al. Analyzing big GABA: Comparison of five software packages for GABA-edited MRS. In: ; 2019.
13. Clarke WT, Stagg CJ, Jbabdi S. FSL-MRS: An end-to-end spectroscopy analysis package. *Magn Reson Med.* 2021;85(6):2950-2964. doi:10.1002/mrm.28630
14. Edden RAE, Puts NAJ, Harris AD, Barker PB, Evans CJ. Gannet: A batch-processing tool for the quantitative analysis of gamma-aminobutyric acid-edited MR spectroscopy spectra:

- 1 Gannet: GABA Analysis Toolkit. *J Magn Reson Imaging*. 2014;40(6):1445-1452.
2 doi:10.1002/jmri.24478
- 3 15. Provencher SW. Estimation of metabolite concentrations from localized in vivo proton
4 NMR spectra. *Magn Reson Med*. 1993;30(6):672-679. doi:10.1002/mrm.1910300604
- 5 16. Oeltzschner G, Zöllner HJ, Hui SCN, et al. Osprey: Open-source processing,
6 reconstruction & estimation of magnetic resonance spectroscopy data. *J Neurosci*
7 *Methods*. 2020;343:108827. doi:10.1016/j.jneumeth.2020.108827
- 8 17. Reynolds G, Wilson M, Peet A, Arvanitis TN. An algorithm for the automated
9 quantitation of metabolites in in vitro NMR signals. *Magn Reson Med*. 2006;56(6):1211-
10 1219. doi:10.1002/mrm.21081
- 11 18. Wilson M, Reynolds G, Kauppinen RA, Arvanitis TN, Peet AC. A constrained least-squares
12 approach to the automated quantitation of in vivo 1 H magnetic resonance
13 spectroscopy data: Automated Quantitation of In Vivo 1 H MRS Data. *Magn Reson Med*.
14 2011;65(1):1-12. doi:10.1002/mrm.22579
- 15 19. Vanhamme L, van den Boogaart A, Van Huffel S. Improved Method for Accurate and
16 Efficient Quantification of MRS Data with Use of Prior Knowledge. *J Magn Reson*.
17 1997;129(1):35-43. doi:10.1006/jmre.1997.1244
- 18 20. Graveron-Demilly D. Quantification in magnetic resonance spectroscopy based on semi-
19 parametric approaches. *Magn Reson Mater Phys Biol Med*. 2014;27(2):113-130.
20 doi:10.1007/s10334-013-0393-4
- 21 21. Ratiney H, Coenradie Y, Cavassila S, van Ormondt D, Graveron-Demilly D. Time-domain
22 quantitation of 1 H short echo-time signals: background accommodation. *MAGMA Magn*
23 *Reson Mater Phys Biol Med*. 2004;16(6):284-296. doi:10.1007/s10334-004-0037-9
- 24 22. Naressi A, Couturier C, Devos JM, et al. Java-based graphical user interface for the MRUI
25 quantitation package. *Magma Magn Reson Mater Phys Biol Med*. 2001;12(2-3):141-152.
26 doi:10.1007/BF02668096
- 27 23. Stefan D, Cesare FD, Andrasescu A, et al. Quantitation of magnetic resonance
28 spectroscopy signals: the jMRUI software package. *Meas Sci Technol*.
29 2009;20(10):104035. doi:10.1088/0957-0233/20/10/104035
- 30 24. Mikkelsen M, Rimbault DL, Barker PB, et al. Big GABA II: Water-referenced edited MR
31 spectroscopy at 25 research sites. *NeuroImage*. 2019;191:537-548.
32 doi:10.1016/j.neuroimage.2019.02.059
- 33 25. Mikkelsen M, Barker PB, Bhattacharyya PK, et al. Big GABA: Edited MR spectroscopy at
34 24 research sites. *NeuroImage*. 2017;159:32-45. doi:10.1016/j.neuroimage.2017.07.021
- 35 26. Považan M, Mikkelsen M, Berrington A, et al. Comparison of multivendor single-voxel
36 MR spectroscopy data acquired in healthy brain at 26 sites. *Radiology*. 2020;295(1):171-
37 180.
- 38 27. Saleh MG, Rimbault D, Mikkelsen M, et al. Multi-vendor standardized sequence for
39 edited magnetic resonance spectroscopy. *NeuroImage*. 2019;189:425-431.
40 doi:10.1016/j.neuroimage.2019.01.056

- 1 28. Near J, Harris AD, Juchem C, et al. Preprocessing, analysis and quantification in single-
2 voxel magnetic resonance spectroscopy: experts' consensus recommendations. *NMR*
3 *Biomed.* 2021;34(5). doi:10.1002/nbm.4257
- 4 29. Simpson R, Devenyi GA, Jezzard P, Hennessy TJ, Near J. Advanced processing and
5 simulation of MRS data using the FID appliance (FID-A)—An open source, MATLAB -
6 based toolkit. *Magn Reson Med.* 2017;77(1):23-33. doi:10.1002/mrm.26091
- 7 30. Near J, Edden R, Evans CJ, Paquin R, Harris A, Jezzard P. Frequency and phase drift
8 correction of magnetic resonance spectroscopy data by spectral registration in the time
9 domain: MRS Drift Correction Using Spectral Registration. *Magn Reson Med.*
10 2015;73(1):44-50. doi:10.1002/mrm.25094
- 11 31. Klose U. In vivo proton spectroscopy in presence of eddy currents. *Magn Reson Med.*
12 1990;14(1):26-30. doi:10.1002/mrm.1910140104
- 13 32. Barkhuijsen H, de Beer R, van Ormondt D. Improved algorithm for noniterative time-
14 domain model fitting to exponentially damped magnetic resonance signals. *J Magn*
15 *Reson* 1969. 1987;73(3):553-557. doi:10.1016/0022-2364(87)90023-0
- 16 33. Kreis R. Issues of spectral quality in clinical1H-magnetic resonance spectroscopy and a
17 gallery of artifacts. *NMR Biomed.* 2004;17(6):361-381. doi:10.1002/nbm.891
- 18 34. Gasparovic C, Song T, Devier D, et al. Use of tissue water as a concentration reference
19 for proton spectroscopic imaging. *Magn Reson Med.* 2006;55(6):1219-1226.
20 doi:10.1002/mrm.20901
- 21 35. Rideaux R, Mikkelsen M, Edden RAE. Comparison of methods for spectral alignment and
22 signal modelling of GABA-edited MR spectroscopy data. *NeuroImage.* 2021;232:117900.
23 doi:10.1016/j.neuroimage.2021.117900
- 24 36. Zhang Y, An L, Shen J. Fast computation of full density matrix of multispin systems for
25 spatially localized in vivo magnetic resonance spectroscopy. *Med Phys.* 2017;44(8):4169-
26 4178. doi:10.1002/mp.12375
- 27 37. Kaiser LG, Young K, Meyerhoff DJ, Mueller SG, Matson GB. A detailed analysis of
28 localized J-difference GABA editing: theoretical and experimental study at 4 T. *NMR*
29 *Biomed.* 2008;21(1):22-32. doi:10.1002/nbm.1150
- 30 38. Murdoch JB, Dydak U. Modeling MEGA-PRESS macromolecules for a better grasp of
31 GABA. In: ; 2011:1394. <https://cds.ismrm.org/protected/11MProceedings/files/1394.pdf>
- 32 39. Bhagwagar Z, Wylezinska M, Jezzard P, et al. Reduction in Occipital Cortex γ -
33 Aminobutyric Acid Concentrations in Medication-Free Recovered Unipolar Depressed
34 and Bipolar Subjects. *Biol Psychiatry.* 2007;61(6):806-812.
35 doi:10.1016/j.biopsych.2006.08.048
- 36 40. Zöllner HJ, Tapper S, Hui SCN, Barker PB, Edden RAE, Oeltzschner G. Comparison of
37 Linear Combination Modeling Strategies for GABA-Edited MRS at 3T. *Biochemistry;*
38 2021. doi:10.1101/2021.05.26.445817
- 39 41. Pedrosa de Barros N, Slotboom J. Quality management in in vivo proton MRS. *Anal*
40 *Biochem.* 2017;529:98-116. doi:10.1016/j.ab.2017.01.017

- 1 42. Kreis R. The trouble with quality filtering based on relative Cramér-Rao lower bounds:
2 The Trouble with Quality Filtering Based on Relative CRLB. *Magn Reson Med*.
3 2016;75(1):15-18. doi:10.1002/mrm.25568
- 4 43. Rousseeuw PJ, Hubert M. Robust statistics for outlier detection. *WIREs Data Min Knowl*
5 *Discov*. 2011;1(1):73-79. doi:10.1002/widm.2
- 6 44. McKinney W. Data Structures for Statistical Computing in Python. In: ; 2010:56-61.
7 doi:10.25080/Majora-92bf1922-00a
- 8 45. Harris CR, Millman KJ, van der Walt SJ, et al. Array programming with NumPy. *Nature*.
9 2020;585(7825):357-362. doi:10.1038/s41586-020-2649-2
- 10 46. SciPy 1.0 Contributors, Virtanen P, Gommers R, et al. SciPy 1.0: fundamental algorithms
11 for scientific computing in Python. *Nat Methods*. 2020;17(3):261-272.
12 doi:10.1038/s41592-019-0686-2
- 13 47. Vallat R. Pingouin: statistics in Python. *J Open Source Softw*. 2018;3(31):1026.
14 doi:10.21105/joss.01026
- 15 48. Seabold S, Perktold J. statsmodels: Econometric and statistical modeling with python. In:
16 ; 2010. <https://www.statsmodels.org/>
- 17 49. Shapiro SS, Wilk MB. An Analysis of Variance Test for Normality (Complete Samples).
18 *Biometrika*. 1965;52(3/4):591. doi:10.2307/2333709
- 19 50. Fligner MA, Killeen TJ. Distribution-Free Two-Sample Tests for Scale. *J Am Stat Assoc*.
20 1976;71(353):210-213. doi:10.1080/01621459.1976.10481517
- 21 51. Bland JM, Altman DG. Statistical methods for assessing agreement between two
22 methods of clinical measurement. *Lancet Lond Engl*. 1986;1(8476):307-310.
- 23 52. Welch BL. The generalisation of student's problems when several different population
24 variances are involved. *Biometrika*. 1947;34(1-2):28-35. doi:10.1093/biomet/34.1-2.28
- 25 53. Holm S. A Simple Sequentially Rejective Multiple Test Procedure. *Scand J Stat*.
26 1979;6(2):65-70.
- 27 54. Bonferroni CE. Il calcolo delle assicurazioni su gruppi di teste. *Studi Onore Profr*
28 *Salvatore Ortu Carboni*. Published online 1935:13-60.
- 29 55. R Core Team. R: A Language and Environment for Statistical Computing. R Foundation
30 for Statistical Computing; 2018. <https://www.R-project.org/>
- 31 56. Bates D, Mächler M, Bolker B, Walker S. Fitting Linear Mixed-Effects Models Using lme4.
32 *J Stat Softw*. 2015;67(1). doi:10.18637/jss.v067.i01
- 33 57. Halekoh U, Højsgaard S. A Kenward-Roger Approximation and Parametric Bootstrap
34 Methods for Tests in Linear Mixed Models - The R Package pbkrtest. *J Stat Softw*.
35 2014;59(9). doi:10.18637/jss.v059.i09
- 36 58. Harris AD, Puts NAJ, Edden RAE. Tissue correction for GABA-edited MRS: Considerations
37 of voxel composition, tissue segmentation, and tissue relaxations: Tissue Correction for
38 GABA-Edited MRS. *J Magn Reson Imaging*. 2015;42(5):1431-1440.
39 doi:10.1002/jmri.24903

- 1 59. Jensen JE, deB. Frederick B, Renshaw PF. Grey and white matter GABA level differences
2 in the human brain using two-dimensional, J-resolved spectroscopic imaging. *NMR*
3 *Biomed.* 2005;18(8):570-576. doi:10.1002/nbm.994
- 4 60. Mikkelsen M, Singh KD, Brealy JA, Linden DEJ, Evans CJ. Quantification of γ -aminobutyric
5 acid (GABA) in 1 H MRS volumes composed heterogeneously of grey and white matter:
6 GABA Quantification in Heterogeneous Volumes. *NMR Biomed.* 2016;29(11):1644-1655.
7 doi:10.1002/nbm.3622
- 8 61. Pernet CR, Wilcox R, Rousselet GA. Robust Correlation Analyses: False Positive and
9 Power Validation Using a New Open Source Matlab Toolbox. *Front Psychol.* 2013;3.
10 doi:10.3389/fpsyg.2012.00606
- 11 62. Rousselet GA, Pernet CR. Improving standards in brain-behavior correlation analyses.
12 *Front Hum Neurosci.* 2012;6. doi:10.3389/fnhum.2012.00119
- 13 63. Landheer K, Juchem C. Are Cramér-Rao lower bounds an accurate estimate for standard
14 deviations in in vivo magnetic resonance spectroscopy? *NMR Biomed.* Published online
15 April 19, 2021. doi:10.1002/nbm.4521
- 16 64. Harris AD, Puts NA, Wijtenburg SA, et al. Normalizing data from GABA-edited MEGA-
17 PRESS implementations at 3 Tesla. *Magn Reson Imaging.* 2017;42:8-15.
18 doi:10.1016/j.mri.2017.04.013
- 19 65. Mullins PG, McGonigle DJ, O’Gorman RL, et al. Current practice in the use of MEGA-
20 PRESS spectroscopy for the detection of GABA. *NeuroImage.* 2014;86:43-52.
21 doi:10.1016/j.neuroimage.2012.12.004
- 22 66. Rothman DL, Behar KL, Prichard JW, Petroff OAC. Homocarnosine and the measurement
23 of neuronal pH in patients with epilepsy. *Magn Reson Med.* 1997;38(6):924-929.
24 doi:10.1002/mrm.1910380611
- 25 67. Behar KL, Ogino T. Characterization of macromolecule resonances in the 1H NMR
26 spectrum of rat brain. *Magn Reson Med.* 1993;30(1):38-44.
27 doi:10.1002/mrm.1910300107
- 28 68. Cudalbu C, Behar KL, Bhattacharyya PK, et al. Contribution of macromolecules to brain 1
29 H MR spectra: Experts’ consensus recommendations. *NMR Biomed.* Published online
30 November 25, 2020. doi:10.1002/nbm.4393
- 31 69. Dydak U, Jiang YM, Long LL, et al. In vivo measurement of brain GABA concentrations by
32 magnetic resonance spectroscopy in smelters occupationally exposed to manganese.
33 *Environ Health Perspect.* 2011;119(2):219-224. doi:10.1289/ehp.1002192
- 34 70. Deelchand DK, Marjańska M, Henry P, Terpstra M. MEGA-PRESS of GABA+: Influences of
35 acquisition parameters. *NMR Biomed.* 2021;34(5). doi:10.1002/nbm.4199
- 36 71. Evans CJ, Puts NAJ, Robson SE, et al. Subtraction artifacts and frequency (Mis-)alignment
37 in J -difference GABA editing: J-Difference GABA Editing. *J Magn Reson Imaging.*
38 2013;38(4):970-975. doi:10.1002/jmri.23923
- 39 72. Jofre F, Anderson ME, Markley JL. L_arginine. Published online 2006.
40 doi:10.13018/BMSE000029

- 1 73. Sanaei Nezhad F, Anton A, Michou E, Jung J, Parkes LM, Williams SR. Quantification of
2 GABA, glutamate and glutamine in a single measurement at 3 T using GABA-edited
3 MEGA-PRESS. *NMR Biomed.* 2018;31(1). doi:10.1002/nbm.3847
- 4 74. Bell T, Boudes ES, Loo RS, et al. In vivo Glx and Glu measurements from GABA-edited
5 MRS at 3 T. *NMR Biomed.* Published online January 28, 2020:e4245.
6 doi:10.1002/nbm.4245
- 7 75. Maddock RJ, Caton MD, Ragland JD. Estimating glutamate and Glx from GABA-optimized
8 MEGA-PRESS: Off-resonance but not difference spectra values correspond to PRESS
9 values. *Psychiatry Res Neuroimaging.* 2018;279:22-30.
10 doi:10.1016/j.pscychresns.2018.07.003
- 11 76. van Veenendaal TM, Backes WH, van Bussel FCG, et al. Glutamate quantification by
12 PRESS or MEGA-PRESS: Validation, repeatability, and concordance. *Magn Reson*
13 *Imaging.* 2018;48:107-114. doi:10.1016/j.mri.2017.12.029
- 14 77. Dhamala E, Abdelkefi I, Nguyen M, Hennessy TJ, Nadeau H, Near J. Validation of in vivo
15 MRS measures of metabolite concentrations in the human brain. *NMR Biomed.*
16 2019;32(3):e4058. doi:10.1002/nbm.4058
- 17 78. Cheng H, Wang A, Newman S, Dydak U. An investigation of glutamate quantification
18 with PRESS and MEGA-PRESS. *NMR Biomed.* 2021;34(2). doi:10.1002/nbm.4453
- 19 79. Pouillet JB, Sima DM, Simonetti AW, et al. An automated quantitation of short echo time
20 MRS spectra in an open source software environment: AQSES. *NMR Biomed.*
21 2007;20(5):493-504. doi:10.1002/nbm.1112
- 22 80. Gajdošík M, Landheer K, Swanberg KM, Juchem C. INSPECTOR: free software for
23 magnetic resonance spectroscopy data inspection, processing, simulation and analysis.
24 *Sci Rep.* 2021;11(1):2094. doi:10.1038/s41598-021-81193-9
- 25 81. Mandal PK, Shukla D. KALPANA: Advanced Spectroscopic Signal Processing Platform for
26 Improved Accuracy to Aid in Early Diagnosis of Brain Disorders in Clinical Setting. *J*
27 *Alzheimers Dis.* 2020;75(2):397-402. doi:10.3233/JAD-191351
- 28 82. Purvis LAB, Clarke WT, Biasioli L, Valkovič L, Robson MD, Rodgers CT. OXSA: An open-
29 source magnetic resonance spectroscopy analysis toolbox in MATLAB. Motta A, ed. *PLOS*
30 *ONE.* 2017;12(9):e0185356. doi:10.1371/journal.pone.0185356
- 31 83. Wilson M. Adaptive baseline fitting for MR spectroscopy analysis. *Magn Reson Med.*
32 2021;85(1):13-29. doi:10.1002/mrm.28385
- 33 84. Soher B, Semanchuk P, Todd D, Steinberg J, Young K. VeSPA: integrated applications for
34 RF pulse design, spectral simulation and MRS data analysis. In: *Proc Int Soc Magn Reson*
35 *Med.* Vol 19. ; 2011:1410.
- 36 85. Lin A, Andronesi O, Bogner W, et al. Minimum Reporting Standards for in vivo Magnetic
37 Resonance Spectroscopy (MRSinMRS): Experts' consensus recommendations. *NMR*
38 *Biomed.* Published online February 9, 2021. doi:10.1002/nbm.4484
- 39 86. Henning A. Advanced Spectral Quantification: Parameter Handling, Nonparametric
40 Pattern Modeling, and Multidimensional Fitting. In: Harris RK, Wasylishen RL, eds.
41 *EMagRes.* John Wiley & Sons, Ltd; 2016:981-994.
42 doi:10.1002/9780470034590.emrstm1472

- 1 87. Deelchand DK, Berrington A, Noeske R, et al. Across-vendor standardization of semi-
2 LASER for single-voxel MRS at 3T. NMR Biomed. 2021;34(5). doi:10.1002/nbm.4218

3

4