1    **Comparison of different linear-combination modelling algorithms for short-**

2    **TE proton spectra**

3    Helge J. Zöllner[1,2], Michal Považan[1,2], Steve C. N. Hui[1,2], Sofie Tapper[1,2], Richard A. E. Ed-

4    den[1,2], Georg Oeltzschner[1,2,*]

5    *[1] Russell H. Morgan Department of Radiology and Radiological Science, The Johns Hopkins*

6    *University School of Medicine, Baltimore, MD, United States*

7    *[2] F. M. Kirby Research Center for Functional Brain Imaging, Kennedy Krieger Institute, Balti-*

8    *more, MD, United States*

9

10   **\*Corresponding author:**

11   Georg Oeltzschner, Ph.D.

12   Division of Neuroradiology, Park 367G

13   The Johns Hopkins University School of Medicine

14   600 N Wolfe St

15   Baltimore, MD 21287

16   goeltzs1@jhmi.edu

17

## Graphical Abstract

## **Abstract**

Short-TE proton MRS is used to study metabolism in the human brain. Common analysis methods model the data as linear combination of metabolite basis spectra. This large-scale multi-site study compares the levels of the four major metabolite complexes in short-TE spectra estimated by three linear-combination modelling (LCM) algorithms.

277 medial parietal lobe short-TE PRESS spectra (TE = 35 ms) from a recent 3T multi-site study were pre-processed with the Osprey software. The resulting spectra were modelled with Osprey, Tarquin and LCModel, using the same three vendor-specific basis sets (GE, Philips, and Siemens) for each algorithm. Levels of total N-acetylaspartate (tNAA), total choline (tCho), myo-inositol (mI), and glutamate+glutamine (Glx) were quantified with respect to total creatine (tCr).

Group means and CVs of metabolite estimates agreed well for tNAA and tCho across vendors and algorithms, but substantially less so for Glx and mI, with mI systematically estimated lower by Tarquin. The cohort mean correlation coefficient for all pairs of LCM algorithms across all datasets and metabolites was $\overline{R^2}$ =0.39, indicating generally only moderate agreement of individual metabolite estimates between algorithms. There was a significant correlation between local baseline amplitude and metabolite estimates (cohort mean $\overline{R^2}$ =0.10).

While mean estimates of major metabolite complexes broadly agree between linear-combination modelling algorithms at group level, correlations between algorithms are only weak-to-moderate, despite standardized pre-processing, a large sample of young, healthy and cooperative subjects, and high spectral quality. These findings raise concerns about the comparability of MRS studies, which typically use one LCM software and much smaller sample sizes.

## Introduction

49

50    Proton MRS allows in-vivo research studies of metabolism[1,2]. Single-voxel MR spectra from the

51    human brain are frequently acquired using PRESS localization[3] , and can be modelled to esti-

52    mate metabolite levels. Accurate modelling is hampered by poor spectral resolution at clinical

53    field strengths, and for short-echo-time spectra, metabolite signals overlap with a broad back-

54    ground consisting of fast-decaying macromolecule and lipid signals. Linear-combination model-

55    ling (LCM) of the spectra maximizes the use of prior knowledge to constrain the model solution,

56    and is recommended by recent consensus[4]. LCM algorithms model spectra as a linear combina-

57    tion of (metabolite and macromolecular (MM)) basis functions, and typically also include terms

58    to account for smooth baseline fluctuations.

59

60    Several LCM algorithms are available to quantify MR spectra (**Table 1** describes some of the

61    most widely used: Osprey[5], INSPECTOR[6], Tarquin[7], AQSES[8], Vespa[9], QUEST[10], LCModel[11]).

62    The implementations (open-source vs. compiled 'black-box'), modelling approaches (modelling

63    domain and baseline model), and their licensure practices are diverse.

64    Surprisingly few studies have compared the performance of different LCM algorithms. Cross-

65    validation of quantitative results has almost exclusively been performed in the context of bench-

66    marking new algorithms against existing solutions. In-vivo comparisons are often limited to

67    small sample sizes, whether analyzing spectra from animal models[7,12,13] or human subjects[7,8,12].

68    To the best of our knowledge, two exceptions compared the LCM performance of different algo-

69    rithms in rat brain[14] and human body[15], respectively. Most studies report good agreement be-

70    tween results from different algorithms, inferring this from group-mean comparisons, or observ-

71    ing that differences between clinical groups are consistent regardless of the algorithm ap-

72    plied[14,16]. Correlations of estimates from different algorithms are rarely reported; however, a high

73    correlation between LCModel and Tarquin results was found in the rat brain at ultra-high field[14].

74    Despite the fact that LCM has been used to analyze thousands of studies (**Table 1**), a compre-

75    hensive assessment of the agreement between the algorithms is lacking, and the relationship be-

76    tween the choice of model parameters and quantitative outcomes is poorly understood. To begin

77    to address this gap, we conducted a large-scale comparison of short-TE in-vivo MRS data using

78    three LCM algorithms with standardized pre-processing. While recent expert consensus recom-

79    mends using measured MM background spectra, data for different sequences are not broadly

80    available or integrated in LCM software. This manuscript investigates current common practice,

81    and therefore all models included simulated MM basis functions as defined in LCModel. We

82    compared group-mean quantification results of four major metabolite complexes from each LCM

83    algorithm, performed between-algorithm correlation analyses, and investigated local baseline

84    power and creatine modelling as potential sources of differences between the algorithms.

85

86    ***Table 1****. Overview of linear-combination modelling algorithms. The domain (time TD or fre-*
87    *quency FD) of modelling and the baseline model approach are specified. *Citations re-*
88    *ported from Google Scholar on September 14, 2020.*

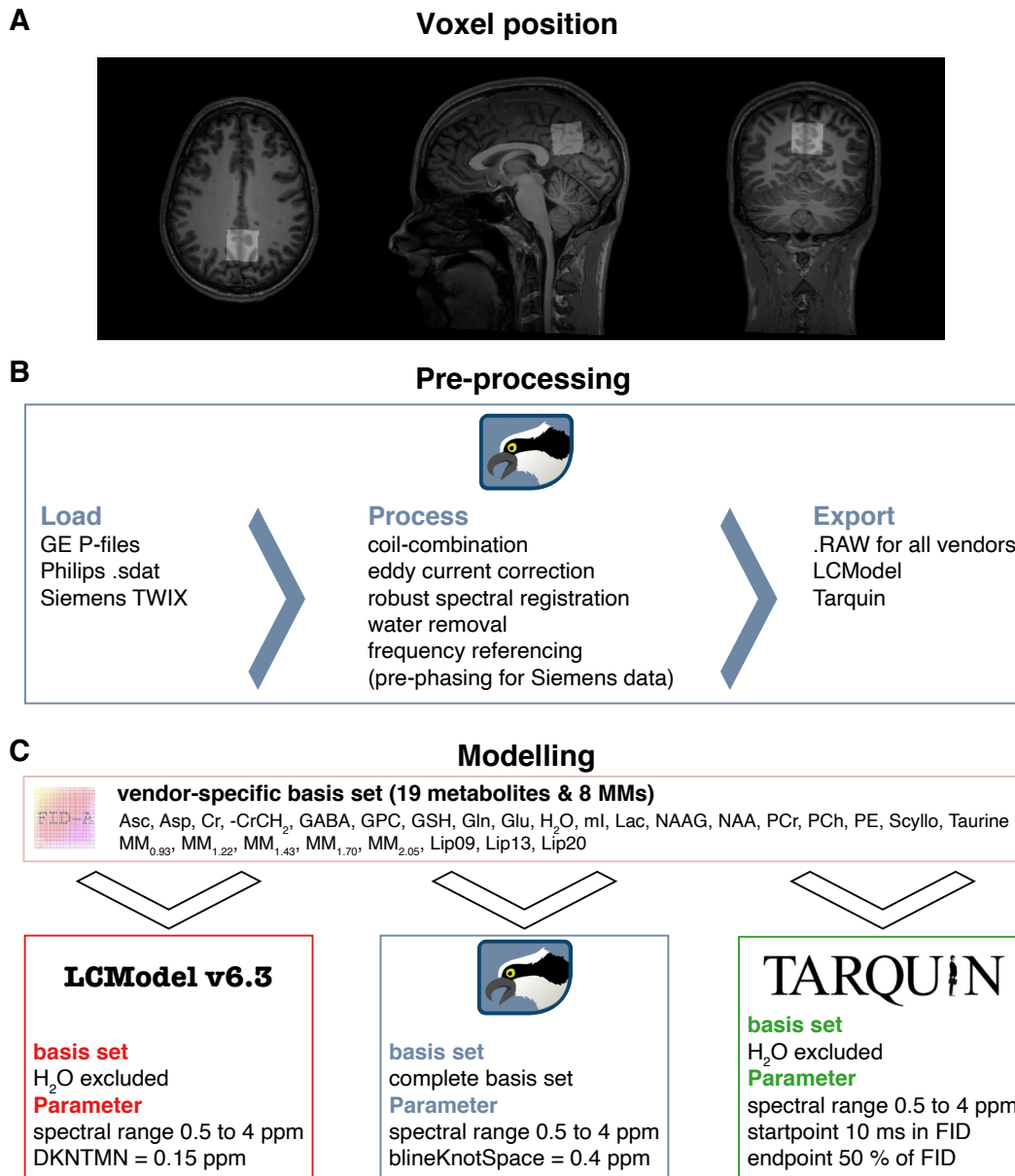| Name | Modelling Domain, Baseline approach | Cost | Code Availability | Published | Cita-tions* |
|---|---|---|---|---|---|
| Osprey | FD, spline baseline | free | open | 2020 | 1 |
| INSPECTOR | FD, $1^{st}$-order polynomial | free | open | 2018 | 0 |
| Tarquin | TD, smooth baseline | free | open | 2011 | 259 |
| AQSES (jMRUI) | TD, spline baseline | free | closed | 2007 | 141 |
| Vespa | FD, wavelet baseline | free | open | 2006 | 72 |
| QUEST (jMRUI) | TD, spline baseline | free | closed | 2004 | 34 |
| LCModel | FD, spline baseline | $13,300 | closed | 1992 | 3482 |

89

## Methods

### Participants & acquisition

277 single-voxel short-TE PRESS brain datasets from healthy volunteers acquired in a recent 3T multisite-study[17] were included in this analysis. Data were acquired at 25 sites (with up to 12 subjects per site) on scanners from three different vendors (GE: 8 sites with n = 91; Philips: 10 sites with n = 112; and Siemens: 7 sites with n = 74) with the following parameters: TR/TE = 2000/35 ms; 64 averages; 2, 4 or 5 kHz spectral bandwidth; 2048-4096 data points; acquisition time = 2.13 min; $3\times3\times3$ cm$^3$ voxel in the medial parietal lobe (**Figure 1A**). The water suppression pulse bandwidth was 140 Hz for Philips, 50 Hz for Siemens, and 150 Hz for GE. Reference spectra were acquired with similar parameters, but without water suppression and 8-16 averages. No more acquisition parameters were specified (for more details, please refer to [17]). Data were saved in vendor-native formats (GE P-files, Philips .sdat, and Siemens TWIX). In the initial study[18], written informed consent was obtained from each participant and the study was approved by local institutional review boards. Anonymized data were shared securely and analyzed at Johns Hopkins University with local IRB approval. Due to site-based data privacy guidelines, only a subset of these data (GE: 7 sites with n = 79; Philips: 9 sites with n = 100; and Siemens: 4 sites with n = 48) is publicly available[19].

### Data pre-processing

MRS data were pre-processed in Osprey[5], an open-source MATLAB toolbox, following recent peer-reviewed pre-processing recommendations[2], as summarized in **Figure 1B**. First, the vendor-native raw data were loaded, including the metabolite (water-suppressed) data and unsuppressed water reference data. Second the raw data were pre-processed into averaged spectra. Receiver-coil combination[20] and eddy-current correction[21] of the metabolite data were performed using the water reference data. Individual transients in Siemens and GE data were frequency-and-phase aligned using robust spectral registration[22], while Philips data had been averaged on the scanner. After averaging the individual transients, the residual water signal was removed with a Hankel singular value decomposition (HSVD) filter[23]. For Siemens spectra, an additional pre-phasing step was introduced by modelling the signals from creatine and choline-containing compounds at 3.02 and 3.20 ppm with a double Lorentzian model and applying the inverted model phase to the

120    data. This step corrected a zero-order phase shift in the data arising from the HSVD water re-

121    moval, likely because the Siemens water suppression introduced asymmetry to the residual water

122    signal. Finally, the pre-processed spectra were exported in .RAW format.

123



124

***Figure 1****. Voxel position and overview of the MRS analysis pipeline. (A) Representative voxel position in the medial parietal lobe extracted with 'OspreyCoreg' (B) Pre-processing pipeline implemented in Osprey including 'OspreyLoad' to load the vendor-native spectra, 'OspreyProcess' to process the raw data and to export the averaged spectra. (C) Modelling of the averaged spectra with details of the basis set and parameters of each LCM (LCModel, Osprey, and Tarquin).*

125 Data modelling

126 Fully localized 2D density-matrix simulations implemented in the MATLAB toolbox FID-A [24]

127 with vendor-specific refocusing pulse information, timings, and phase cycling were used to gen-

128 erate three vendor-specific basis sets (GE, Philips, and Siemens) including 19 spin systems:

129 ascorbate, aspartate, Cr, negative creatine methylene (-CrCH$_2$), γ-aminobutyric acid (GABA),

130 glycerophosphocholine (GPC), glutathione, glutamine (Gln), glutamate (Glu), water (H$_2$O), myo-

131 inositol (mI), lactate, NAA, N-acetylaspartylglutamate (NAAG), phosphocholine (PCh), PCr,

132 phosphoethanolamine, scyllo-inositol, and taurine. The -CrCH$_2$ term is a simulated negative cre-

133 atine methylene singlet at 3.95 ppm, included as a correction term to account for effects of water

134 suppression and relaxation. It is not included in the tCr model, which is used for quantitative ref-

135 erencing.

136 8 additional Gaussian basis functions were included in the basis set to simulate broad macromol-

137 ecules and lipid resonances[25] (simulated as defined in section 11.7 of the LCModel manual[26]):

138 MM$_{0.94}$, MM$_{1.22}$, MM$_{1.43}$, MM$_{1.70}$, MM$_{2.05}$, Lip09, Lip13, Lip20. The Gaussian amplitudes were

139 scaled relative to the 3.02 ppm creatine CH$_3$ singlet in each basis set (details in **Supplementary**

140 **Material 1**). Finally, to standardize the basis set for each algorithm, basis sets were stored as

141 .mat files for use in Osprey and as .BASIS-files for use in LCModel and Tarquin. In the follow-

142 ing paragraphs, each LCM algorithm investigated in this study is described briefly (for details,

143 please refer to the original publications[5,7,11]).

144

145 *LCModel v6.3*

146 The LCModel (6.3-0D) algorithm[11] models data in the frequency-domain. First, time-domain

147 data and basis functions are zero-filled by a factor of two. Second, frequency-domain spectra are

148 frequency-referenced by cross-correlating them with a set of delta functions representing the ma-

149 jor singlet landmarks of NAA (2.01 ppm), Cr (3.02 ppm), and Cho (3.20 ppm). Third, starting

150 values for phase and linebroadening parameters are estimated by modelling the data with a re-

151 duced basis set containing NAA, Cr, PCh, Glu, and mI, with a smooth baseline. Fourth, the final

152 modelling of the data is performed with the full basis set, regularized lineshape model and base-

153 line, with starting values for phase, linebroadening, and lineshape parameters derived from the

154 previous step. Model parameters are determined with a Levenberg-Marquardt[27,28] non-linear

155 least-squares optimization implementation that allows bounds to be imposed on the parameters.

156     Metabolite amplitude bounds are defined to be non-negative, and determined using a non-nega-

157     tive linear least-squares (NNLS) fit at each iteration of the non-linear optimization. Amplitude

158     ratio constraints on macromolecule and lipid amplitude, as well as selected pairs of metabolite

159     amplitudes (e.g. NAA+NAAG), are defined as in Osprey and Tarquin. LCModel constrains the

160     model with three additional regularization terms. Two of these terms penalize a lack of smooth-

161     ness in the baseline and lineshape models using the second derivative operator, preventing unrea-

162     sonable baseline flexibility and lineshape irregularity. The third term penalizes deviations of the

163     metabolite Lorentzian linebroadening and frequency shift parameters from their expected values.

164

165     *Osprey*

166     The Osprey (1.0.0) algorithm[5] adopts several key features of the LCModel and Tarquin algo-

167     rithms. Osprey follows the four-step workflow of LCModel including zero-filling, frequency ref-

168     erencing, preliminary optimization to determine starting values, and final optimization over the

169     real part of the frequency-domain spectrum. The model parameters are zero- and first-order

170     phase correction, global Gaussian linebroadening, individual Lorentzian linebroadening, and in-

171     dividual frequency shifts, which are applied to each basis function before Fourier transformation.

172     The frequency-domain basis functions are then convolved with an arbitrary, unregularized line-

173     shape model to account for deviations from a Voigt profile. The length of this lineshape model is

174     estimated during the initial referencing step and set to 2.5 times the FWHM estimate. The line-

175     shape model is normalized, so that the convolution does not impact the integral of basis func-

176     tions.

177     The spline baseline is constructed from cubic B-spline basis functions, including one additional

178     knot outside either end of the user-specified fit range, as in LCModel. In contrast to LCModel,

179     the baseline curvature is not regularized. Therefore, the baseline knot spacing is set to 0.15 ppm

180     for preliminary modelling step with a reduced basis set and increased to 0.4 ppm for the final full

181     model. Similar to LCModel, model parameters are determined with a Levenberg-Marquardt[27,28]

182     non-linear least-squares optimization algorithm and a NNLS fit to determine the non-negative

183     metabolite amplitudes at each step of the non-linear optimization.

184     *Tarquin*

185     Tarquin (4.3.10)[7] uses a four-step approach in the time domain to model spectra. First, residual

186     water is removed using singular value decomposition. Second, the global zero-order phase is

187    determined by minimizing the difference between the magnitude and the real spectra in the fre-

188    quency domain. Third, zero-filling to double the number of points and frequency referencing are

189    performed, as in the other algorithms. This step also estimates a starting value for the Gaussian

190    linebroadening used in the fourth step, the final modelling. The model includes common Gauss-

191    ian linebroadening, individual Lorentzian linebroadening, individual frequency-shifts, and zero-

192    and first-order phase correction factors applied in the frequency domain.

193    Optimization is performed in the time domain with a constrained non-linear least-squares Leven-

194    berg-Marquardt solver, allowing bounds and constraints on the parameters. In addition, the range

195    of time-domain datapoints is limited by removing the first 10 ms of the FID, so as to omit the

196    fast-decaying macromolecule and lipid signals. Finally, the baseline is estimated in the frequency

197    domain by convolving the model residual with a Gaussian filter with a width of 100 points.

198

199    *Model parameters*

200    The parameters chosen for each tool are summarized in **Figure 1C**. The fit range was limited to

201    0.5 to 4 ppm in all tools to reduce effects of differences in water suppression techniques. For the

202    baseline handling, the default and most commonly used parameters were chosen, i.e. bLine-

203    KnotSpace = 0.4 ppm for Osprey, DKNMNT = 0.15 ppm for LCModel, and an FID range from

204    10 ms to 50% of the FID for Tarquin.

205

206    Quantification, visualization, and secondary analyses

207    *Quantifaction*

208    The four major metabolite complexes tNAA (NAA + NAAG), tCho (GPC + PCh), mI, and Glx

209    (Glu + Gln) were quantified as basis-function amplitude ratios relative to total creatine (tCr = Cr

210    + PCr). Since the primary purpose was to compare performance of the core LCM algorithms, no

211    additional relaxation correction or partial volume correction was performed.

212    Model visualizations were generated with the *OspreyOverview* module, which allows LCModel

213    and Tarquin results files (.coord and .txt) to be imported. For each algorithm, the visualization

214    includes site-mean spectra, cohort-mean spectra (i.e. the mean of all spectra), and site- and co-

215    hort-mean modelling results (complete model, spline baseline, spline baseline + MM compo-

216    nents, and the separate models of the major metabolite complexes).

217

218    *Visualization*

219    As in the default visualizations for the LCModel and Tarquin software interfaces, inverse phase

220    estimates were applied to the spectra and final models. For the visualization, spectra were nor-

221    malized to the amplitude of the 3-ppm creatine singlet, and a DC offset was added to each site

222    mean spectrum to align the mean frequency-domain amplitude between 1.85 and 4.0 ppm, to aid

223    visual comparison between algorithms and sites.

224

225    *Secondary analyses*

226    To investigate potential vendor differences in linewidth and SNR based on the different export

227    formats of the data, NAA linewidth and SNR were investigated.

228    To investigate potential interactions between baseline power and metabolite estimates unbiased

229    by DC offsets, the MM + baseline models were first aligned vertically according to the fre-

230    quency-domain minimum of the acquired spectra between 2.66 and 2.7 ppm (i.e. between the as-

231    partyl signals, which is the region with the highest consistency between the baseline models).

232    Baseline models were normalized to the frequency-domain amplitude of each metabolite spec-

233    trum between 2.9 and 3.1 ppm to account for differences in the scaling of the model outputs of

234    LCModel and Tarquin. Baseline power beneath each major metabolite was then defined as the

235    range-normalized integral of the baseline model between 1.9 and 2.1 ppm for the tNAA baseline;

236    3.1 and 3.3 ppm for the tCho baseline; 3.33 and 3.75 ppm for mI; and 1.9 to 2.5 ppm and 3.6 to

237    3.8 ppm for the Glx baseline.

238    The contribution of variance in modelling of the creatine reference signal to metabolite ratios

239    was also investigated. To this end, each individual total creatine model (Cr + PCr) was normal-

240    ized to the frequency-domain amplitude of each metabolite spectrum between 1.9 and 2.1 ppm to

241    account for differences in the scaling of the total creatine model outputs of LCModel and Tar-

242    quin. Finally, the integral over the individual creatine model was calculated.

243

244    Data analysis

245    Quantitative metabolite estimates (tNAA/tCr, tCho/tCr, mI/tCr, Glx/tCr) were statistically ana-

246    lyzed and visualized using R[29] in RStudio (Version 1.2.5019, RStudio Inc.). The functions are

247    publicly available[30]. The supplemental materials with MATLAB- and R-files, example LCModel

248    control files (one for each vendor), and Tarquin batch-files for this study are publicly available[31].

249    The results from each LCM algorithm were imported into R with the *spant* package[32].

250

251    *Distribution analysis*

252    The results are presented  as raincloud plots[33] and Pearson's correlation analysis using the

253    *ggplot2* package[34]. The raincloud plots include individual data points, boxplots with median and

254    $25^{th}/75^{th}$ percentiles, a smoothed distribution, and mean $\pm$ SD error bars to identify systematic

255    differences between the LC algorithms. In addition, the coefficient of variation (CV = SD/mean)

256    and the mean $\overline{CV} = \frac{(CV_{tNAA}+CV_{tCho}+CV_{Ins}+CV_{Glx})}{4}$ across all four metabolites of each algorithm are

257    calculated.

258

259    *Correlation analysis*

260    The correlation analysis featured different levels, including pair-wise correlations between algo-

261    rithms, as well as correlations between baseline power and metabolite estimates of each algo-

262    rithm. The pair-wise correlation on the global level (black $R^2$), as well as within-vendor correla-

263    tions (color-coded $R^2$) with different color shades for different sites are reported. Furthermore,

264    mean $\overline{R^2}$ for each pair-wise correlation (e.g. Osprey vs LCModel) and metabolite, estimated by

265    row or column means e.g. $\overline{R^2} = \frac{(R^2_{tNAA}+R^2_{tCho}+R^2_{Ins}+R^2_{Glx})}{4}$, and a cohort mean $\overline{R^2}$ (across all pair-

266    wise correlations) are calculated. The correlations were Bonferroni corrected for the number of

267    correlation tests. The cohort mean $\overline{R^2}$ was used to identify global associations across all corre-

268    lation analysis, while the mean $\overline{R^2}$ allowed the identification of algorithm-specific (row means)

269    and  metabolite-specific (column means) interactions across all correlation analysis. Associations

270    between the outcome of specific algorithms were identified by the pair-wise correlation analysis

271    ($R^2$). Vendor-specific effects were identified by differentiating between global level and within-

272    vendor correlations.

273

274    *Statistical analysis*

275    In the statistical analysis, the presence of significant differences in the mean and the variance of

276    the metabolite estimates was assessed. Global metabolite estimates were compared between al-

277    gorithms with parametric tests, following recommendations for large sample sizes[35]. Differences

278    of variances were tested with Fligner-Killeen's test with a post-hoc pair-wise Fligner-Killeen's

279    test and Bonferroni correction for the number of pair-wise comparisons. Depending on whether

280    variances were different or not, an ANOVA or Welch's ANOVA was used to compare means

281    with a post-hoc paired t-test with equal or non-equal variances, respectively.

282

283

## **<u>Results</u>**

285 All 277 spectra were successfully processed, exported, and quantified with the three LCM algo-

286 rithms; no modelled spectra were excluded from further analysis.

287 <u>Summary and visual inspection of the modelling results</u>

288 A site-level averaged summary of the 277 spectra is shown in **Figure 2A**, **B** and **C**, for analyses

289 in LCModel, Osprey, and Tarquin, respectively**.** The averaged data, models and residuals for

290 each of the 25 sites are color-coded by vendor. The cohort-mean of all analyses for each vendor

291 is shown in **Figure 2D**, **E** and **F** (GE, Philips and Siemens, respectively). Data, models and re-
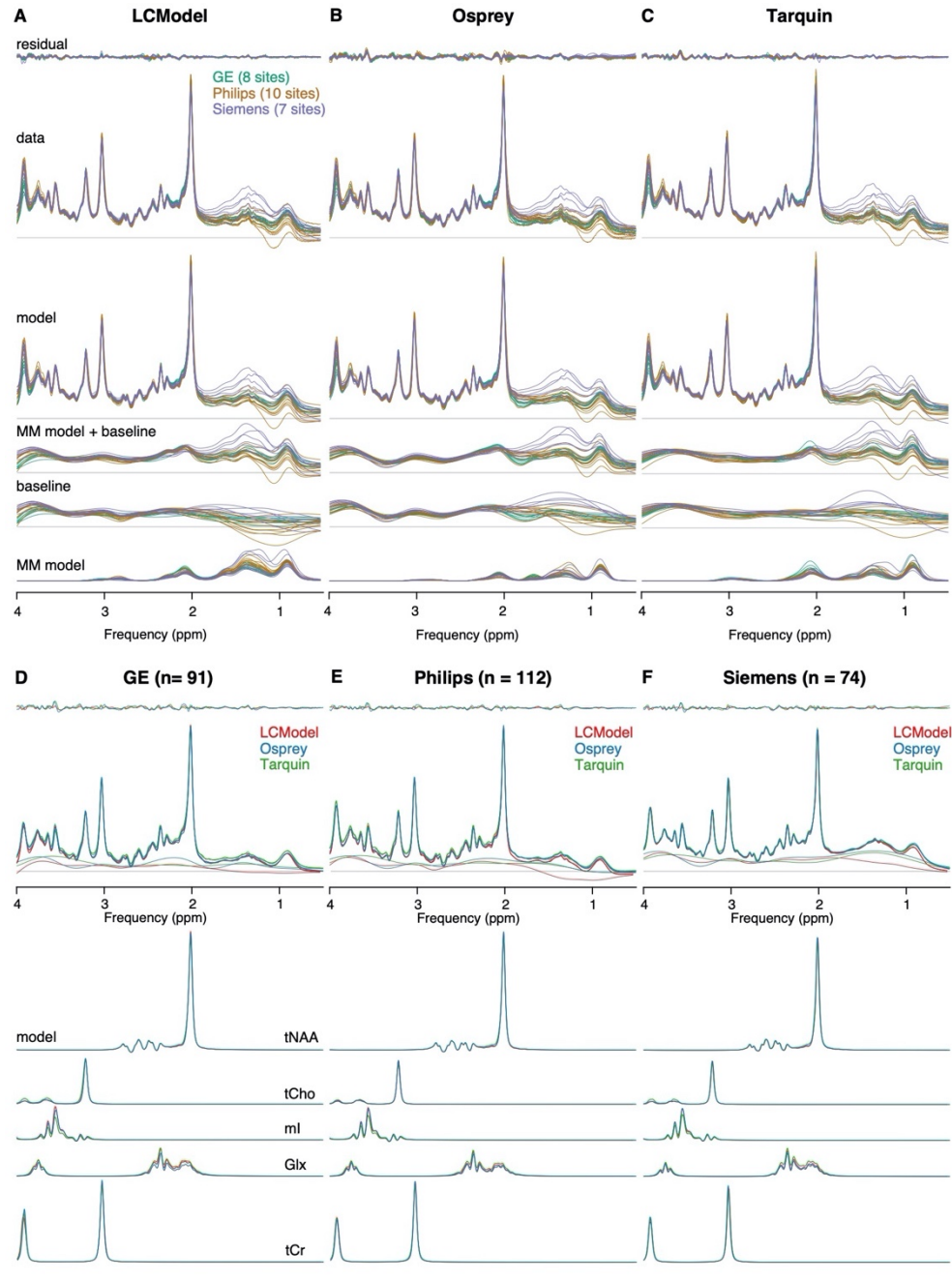
292 siduals are color-coded by algorithm.

293

294 In general, the phased spectra and models agreed well between vendors for all algorithms. Com-

295 paring the algorithms, notable differences in spectral features in the estimated baseline models

296 appeared between 0.5 and 1.95 ppm (degree of variability: Osprey > LCModel > Tarquin) and

297 between 3.6 and 4 ppm (degree of variability: LCModel > Osprey > Tarquin) (as shown in **Fig-**

298 **ure 2A-C**).

299 Cohort-mean spectra and models agreed well across all vendors and algorithms (**Figure 2D-F**).

300 The greatest differences in the spectral features of the baseline between algorithms occur be-

301 tween 0.5 and 1.95 ppm, with closer agreement between Osprey and Tarquin than with

302 LCModel. The amplitude of the residual over the whole spectral range is highest for Osprey, and

303 similar for Tarquin and LCModel. **Supplementary Material 2** shows individual data, models

304 and residuals for each algorithm color-coded by vendor.

305 NAA linewidth was significantly lower ($p < 0.001$) for Philips ($6.3 \pm 1.3$ Hz) compared to GE

306 ($7.3 \pm 1.5$ Hz), while no differences in the linewidth were found for the other comparisons (Sie-

307 mens $6.6 \pm 2.4$ Hz). SNR was significantly higher for Siemens ($285 \pm 72$ ) compared to both

308 other vendors ($p < 0.001$) and significantly higher ($p < 0.001$) for Philips ($226 \pm 58$) compared to
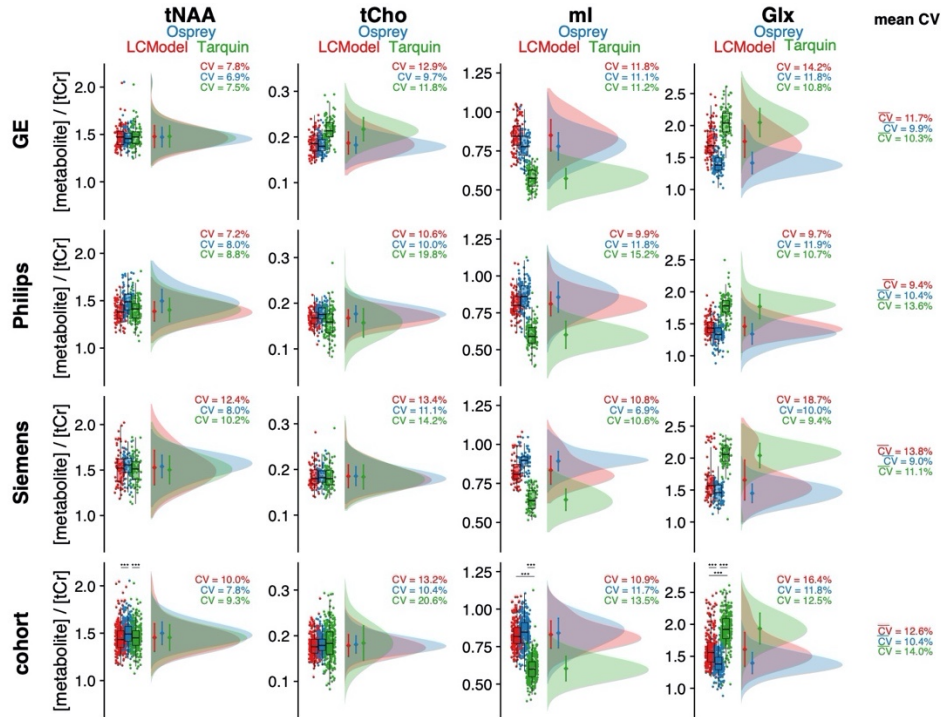
309 GE ($154 \pm 37$).

310

311

*Figure 2. Summary of the modelling results. (A–C) site-level averaged residual, data, model, MM model + baseline, baseline and MM model for each LCM algorithm, color-coded by vendor. (D–F) cohort-mean residual, data, model, MM model + baseline, and metabolite models for each vendor, color-coded by LCM algorithm.*

## Metabolite level distribution

The tCr ratio estimates and CVs of the four metabolites are summarized in **Table 2**. Distributions and group statistics are visualized in **Figure 3**, with the four rows corresponding the three vendors and a cohort summary across all datasets.



***Figure 3***. *Metabolite level distribution. Raincloud plots of the metabolite estimates of each LCM algorithm (color-coded). The four metabolites are reported in the columns, and the three vendors in rows, with a cohort summary in the last row. The coefficient of variation is reported for each distribution, as well as a mean $\overline{CV}$ reported in the last column, which is calculated across each row. Asterisks indicate significant differences (adjusted p < 0.001 = \*\*\*).*

Between-algorithm agreement was greatest for the group means and CVs of tNAA and tCho. The cohort-mean CV was lowest for Osprey (10.4%), followed by LCModel (12.6%) and Tarquin (14.0%). Group means and CVs for tNAA are relatively consistent. As a result, the cohort-mean tNAA/tCr was 1.45 ± 0.15 for LCModel, 1.50 ± 0.12 for Osprey, and 1.45 ± 0.14 for Tarquin, with significant differences between Osprey and both other LCM algorithms. Cohort means for tCho showed a high agreement between all algorithms. The global CV of tCho estimates was significantly higher for Tarquin compared to both other algorithms, and

325     significantly lower for Osprey compared to LCModel. Global tCho/tCr was $0.18 \pm 0.02$ for

326     LCModel, $0.18 \pm 0.02$ for Osprey, and $0.18 \pm 0.04$ for Tarquin.

327     ***Table 2*** *– Metabolite level distribution. Mean, standard deviation and coefficient of variation*
328     *(CV) of each metabolite-to-creatine ratio, listed by algorithm and vendor as well as global*
329     *summary values. Asterisks indicate significant differences (adjusted p < 0.01 = \*\* and adjusted*
330     *p < 0.001 = \*\*\* or ### or ''') in the mean (for the metabolite ratios) or the variance (for the CV)*
331     *compared to the algorithm in the next row (LCModel vs Osprey = \*\* or \*\*\*, Osprey vs Tarquin*
332     *= ###, and Tarquin vs LCModel = ''').*

| | [metabolite] / [tCr] (mean ± SD) | | | |
|---|---|---|---|---|
| | tNAA | tCho | mI | Glx |
| **GE** | | | | |
| *LCModel* | $1.48 \pm 0.12$ | $0.19 \pm 0.02$ | $0.85 \pm 0.10$ | $1.75 \pm 0.25$ |
| *Osprey* | $1.47 \pm 0.10$ | $0.18 \pm 0.02$ | $0.78 \pm 0.09$ | $1.42 \pm 0.17$ |
| *Tarquin* | $1.48 \pm 0.11$ | $0.22 \pm 0.03$ | $0.57 \pm 0.07$ | $2.05 \pm 0.22$ |
| **Philips** | | | | |
| *LCModel* | $1.38 \pm 0.10$ | $0.17 \pm 0.02$ | $0.81 \pm 0.08$ | $1.46 \pm 0.14$ |
| *Osprey* | $1.50 \pm 0.12$ | $0.18 \pm 0.02$ | $0.86 \pm 0.10$ | $1.34 \pm 0.16$ |
| *Tarquin* | $1.40 \pm 0.12$ | $0.16 \pm 0.03$ | $0.60 \pm 0.09$ | $1.78 \pm 0.19$ |
| **Siemens** | | | | |
| *LCModel* | $1.52 \pm 0.19$ | $0.19 \pm 0.02$ | $0.83 \pm 0.09$ | $1.65 \pm 0.31$ |
| *Osprey* | $1.54 \pm 0.12$ | $0.19 \pm 0.02$ | $0.89 \pm 0.06$ | $1.45 \pm 0.14$ |
| *Tarquin* | $1.50 \pm 0.15$ | $0.18 \pm 0.03$ | $0.65 \pm 0.07$ | $2.04 \pm 0.19$ |
| **global** | | | | |
| *LCModel* | $1.45 \pm 0.15$\*\*\* | $0.18 \pm 0.02$ | $0.83 \pm 0.09$ | $1.45 \pm 0.15$\*\*\* |
| *Osprey* | $1.50 \pm 0.12$### | $0.18 \pm 0.02$ | $0.84 \pm 0.09$### | $1.50 \pm 0.12$### |
| *Tarquin* | $1.46 \pm 0.14$ | $0.18 \pm 0.04$ | $0.60 \pm 0.08$''' | $1.93 \pm 0.24$''' |
| | CV (SD/mean) | | | |
| | tNAA | tCho | mI | Glx |
| **GE** | | | | |
| *LCModel* | 7.9% | 12.9% | 11.8% | 14.2% |
| *Osprey* | 6.9% | 9.7% | 11.1% | 11.8% |
| *Tarquin* | 7.5% | 11.7% | 11.2% | 10.8% |
| **Philips** | | | | |
| *LCModel* | 7.2% | 10.6% | 9.9% | 9.7% |
| *Osprey* | 8.0% | 10.0% | 11.8% | 11.9% |
| *Tarquin* | 8.8% | 19.8% | 15.2% | 10.7% |
| **Siemens** | | | | |
| *LCModel* | 12.4% | 13.4% | 10.8% | 18.7% |
| *Osprey* | 8.0% | 11.1% | 6.9% | 10.0% |
| *Tarquin* | 10.1% | 14.3% | 10.5% | 9.3% |
| **global** | | | | |
| *LCModel* | 10.0% | 13.2%\*\* | 10.9% | 16.4%\*\*\* |
| *Osprey* | 7.8% | 10.4%### | 11.7%### | 11.8%### |
| *Tarquin* | 9.3% | 20.5%''' | 13.6% | 12.3% |

333

334

335     For mI, group means and CVs were comparable for Osprey and LCModel, while Tarquin esti-

336     mates were lower by about 25%. Global CVs were significantly lower for Osprey compared to

337     Tarquin, while no significant differences in the CV were found for the other comparisons. Global

338     mI/tCr was $0.83 \pm 0.09$ for LCModel, $0.84 \pm 0.09$ for Osprey, and $0.60 \pm 0.08$ for Tarquin, with

339     significant mean differences between all Tarquin and both other algorithms.

340     Group means and CVs for Glx were comparable between Osprey and LCModel, while estimates

341     were about 30% higher in Tarquin. Global CV was significantly lower for Osprey compared to

342     both other algorithms. Global Glx/tCr was $1.45 \pm 0.15$ for LCModel, $1.50 \pm 0.12$ for Osprey, and

343     $1.93 \pm 0.24$ for Tarquin, with significant differences between all algorithms. Mean $\overline{CVs}$, esti-

344     mated by the row-mean, were between 9.0 and 13.8% for all algorithms and vendors.

345

346     <u>Correlation analysis: pairwise comparison between LCM algorithms</u>

347     The correlation analysis for each metabolite and algorithm pair is summarized in **Figure 4**. $\overline{R^2}$

348     for each algorithm pair and metabolite are reported in the corresponding row and column, re-
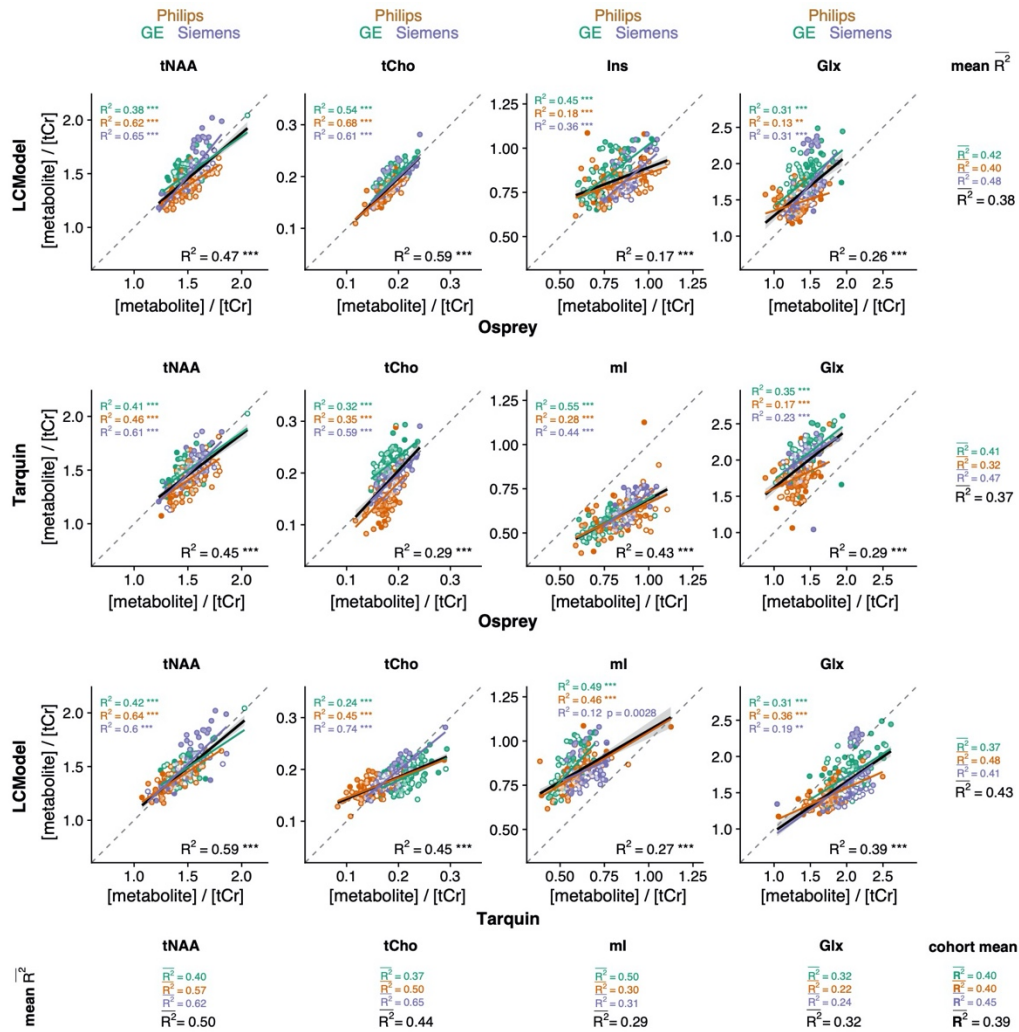
349     spectively.

350     The cohort-mean $\overline{\mathbf{R}^2} = 0.39$ suggests an overall moderate agreement between metabolite esti-

351     mates from different algorithms. The agreement between algorithms, estimated by the row-mean

352     $\overline{R^2}$, was highest for Tarquin-vs-LCModel ($\overline{R^2} = 0.43$), followed by Osprey-vs-LCModel ($\overline{R^2}$

353     $= 0.38$), and Osprey-vs-Tarquin ($\overline{R^2} = 0.37$).

354     The agreement between algorithm for each metabolite, estimated by the column-mean $\overline{R^2}$, was

355     highest for tNAA ($\overline{R^2} = 0.50$), followed by tCho ($\overline{R^2} = 0.44$), Glx ($\overline{R^2} = 0.32$), and mI ($\overline{R^2} =$

356     $0.29$). The cohort-mean $\overline{\mathbf{R}^2}$ for each vendor was higher for Siemens ($\overline{\mathbf{R}^2} = 0.45$) than for GE

357     ($\overline{\mathbf{R}^2} = 0.40$) and Philips ($\overline{\mathbf{R}^2} = 0.40$).

358

359     While the within-metabolite mean $\overline{R^2}$ (average down the columns in Figure 4) are comparable

360     between vendors, there is substantially higher variability of the $R^2$ values with increasing

361  granularity of the analysis. **Supplementary Material 3** includes an additional layer of correla-
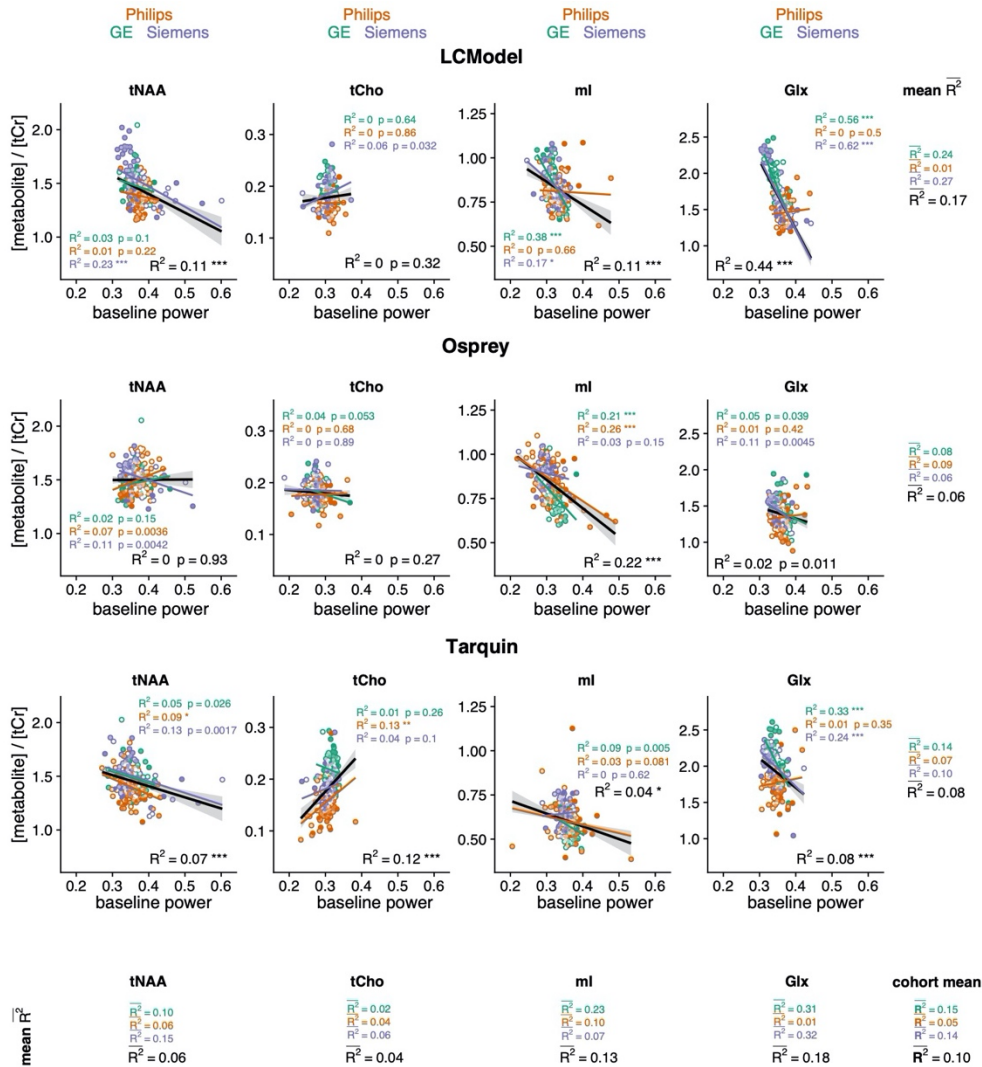
362  tions at the site level.



363

364

***Figure 4.*** *Pairwise correlational comparison of algorithms. LCModel and Osprey are compared in the first row, Tarquin and Osprey in the second row, and LCModel and Tarquin in the third row. Each column corresponds to a different metabolite. Within-vendor correlations are color-coded; global correlations are shown in black. The $\overline{R^2}$ values are calculated along each dimension of the grid with mean $R^2$ for each metabolite and each correlation. A cohort-mean $\overline{R^2}$ value is also calculated across all twelve pair-wise correlations. Asterisks indicate significant correlations (adjusted $p < 0.01$ = ** and adjusted $p < 0.001$ = ***).*

365 <u>Correlation analysis: baseline and metabolite estimates</u>



366

***Figure 5.*** *Correlation analysis between metabolite estimates and local baseline power for each algorithm, including global (black) and within-vendor (color-coded) correlations. The mean $\overline{R^2}$ values are calculated along each dimension of the grid for each metabolite and each algorithm. Similarly, a cohort-mean $\overline{R^2}$ value is calculated across all twelve pair-wise correlations. Asterisks indicate significant correlations (adjusted $p < 0.05$ = \*, adjusted $p < 0.01$ = \*\* , adjusted $p < 0.001$ = \*\*\*).*

367

368 The correlation analysis between local baseline power and metabolite estimates for each algo-

369 rithm is summarized in **Figure 5**. The cohort-mean $\overline{R^2} = 0.10$ suggests that overall, there is an

370 association between local baseline power and metabolite estimates, that is weak but statistically

371 significant. The influence of baseline on metabolite estimates differs between metabolites, as
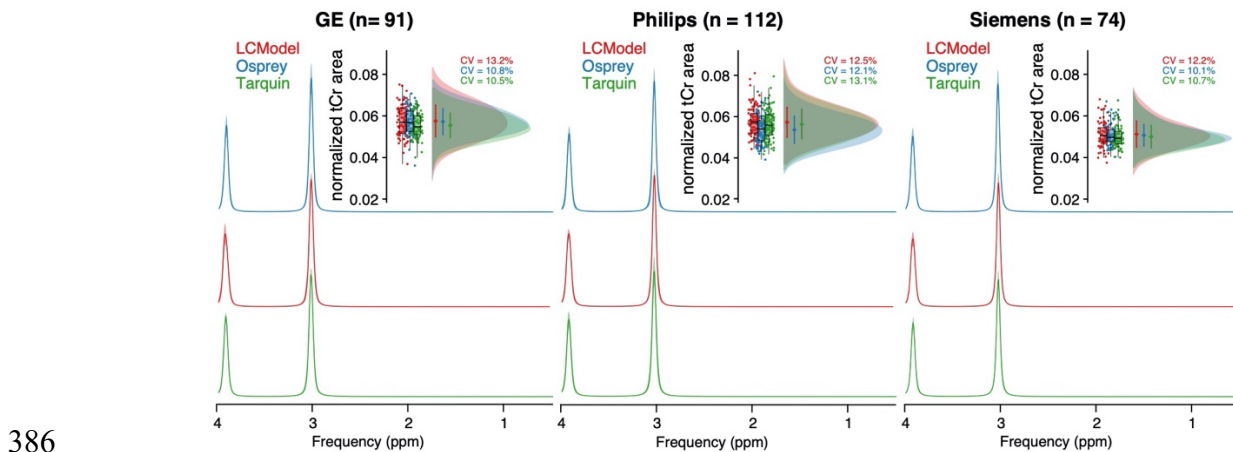
372    reflected by the column-mean $\overline{R^2}$ which was lowest for tCho ($\overline{R^2} = 0.04$) and tNAA ($\overline{R^2} =$

373    0.06), and higher for mI ($\overline{R^2} = 0.13$) and Glx ($\overline{R^2} = 0.18$). The global baseline correlations all

374    had negative slope, except for tCho estimates of Tarquin.

375    The mean $\overline{R^2}$ across metabolites for each algorithm, calculated as the row mean, were low for

376    all algorithms with LCModel ($\overline{R^2} = 0.17$) showing a greater effect than Tarquin ($\overline{R^2} = 0.08$)

377    and Osprey ($\overline{R^2} = 0.06$). Comparing between vendors, the cohort-mean $\overline{\mathbf{R}^2}$ was higher for GE

378    ($\overline{\mathbf{R}^2} = 0.15$) and Siemens ($\overline{\mathbf{R}^2} = 0.14$) than for Philips ($\overline{\mathbf{R}^2} = 0.05$) spectra.


379    <u>Variability of total creatine models</u>

380    Mean tCr model spectra (± one standard deviation) are summarized in **Figure 6** for each vendor

381    and LCM algorithm, along with distribution plots of the area under the model.

382    The agreement in mean and CV is greatest between Osprey and Tarquin for all vendors, while

383    tCr areas for LCModel appear slightly higher. Differences in water suppression are accounted for

384    with the -CrCH$_2$ correction term, which is not included in the tCr model used for quantitative ref-

385    erencing.

386


*Figure 6. Variability of tCr models. Mean models +/- standard deviation (shaded areas) are presented column-wise by vendor and color-coded by LCM algorithm. The distribution and CV of the areas under the models are inset.*

387

## **Discussion**

388

389   We have presented a three-way comparison of LCM algorithms applied to a large dataset of

390   short-TE in-vivo human brain spectra. The aims at the onset were to compare metabolite esti-

391   mates obtained with different LCM algorithms, as applied in the literature, and to identify poten-

392   tial sources of differences between the algorithms. The major findings are:

393   •   Group means and CVs for tNAA and tCho agreed well across vendors and algorithms.

394      For mI and Glx, group means and CVs were less consistent between algorithms, with a

395      higher degree of agreement between Osprey and LCModel than with Tarquin.

396   •   The strength of the correlations between individual metabolite estimates from different

397      algorithms was moderate. In general, tNAA and tCho estimates from different algorithms

398      agreed better than Glx and mI. With each sub-level of analysis, the variability of

399      correlation strength increased, i.e. correlations grew increasingly variable when

400      calculated separately for each vendor, or even each site.

401   •   Overall, the association between metabolite estimates and the local baseline power was

402      significant, with mI and Glx showing stronger associations than tNAA and tCho, and

403      LCModel showing greater effects than Tarquin and Osprey.

404     The strong agreement of group means and CVs for metabolites with prominent singlets

405   (tNAA/tCho) and inconsistency for lower-intensity coupled signals (mI/Glx) are in line with pre-

406   vious two-tool comparisons of simulated data [7,15] and in-vivo studies with smaller sample sizes

407   [7,14,16].

408   While previous work highlighted group means and standard deviations, the between-algorithm

409   agreement of individual metabolite estimates has not been extensively studied. Our results sug-

410   gest that substantial variability is introduced by the choice of the analysis software itself, indi-

411   cated by only moderate between-algorithm correlation strength (between-algorithm mean $\overline{R^2}$ <=

412   0.5 for all investigated metabolites), even for the well-established LCM algorithms LCModel and

413   Tarquin ($R^2$ between 0.27 and 0.59 for all metabolites). This finding raises concerns about the

414   generalizability and reproducibility of MRS study results. MRS studies typically suffer from low

415   sample sizes (~20 per comparison group is common). Considering the moderate between-tool

416   correlation of individual estimates, it is likely that marginally significant group effects and corre-

417   lations found with one analysis tool will not be found with another tool, even if the exact same

418    dataset is used. This is exacerbated by the substantial variability of correlation strengths at ven-

419    dor- or even site-level, and is even more likely to be the case for 'real-life' clinical data, given

420    the relatively high quality of the dataset in this study (standardized pre-processing; large sample

421    size; high SNR; low linewidth; young, healthy, cooperative subjects). While two previous studies

422    found that some differences between clinical groups remained significant independent of the

423    LCM algorithm [14,16], this is questionable as a default assumption. The lack of comparability aris-

424    ing from the additional variability originating in the choice of analysis tool is rarely recognized

425    or acknowledged. If choice of analysis tool is a significant contributor to measurement variance,

426    it could be argued that modelling of data with more than one algorithm will improve the

427    robustness and power of MRS studies. It should also be investigated whether the reduction of the

428    degrees of freedom by improving MM and baseline models (e.g. by using acquired MM data)

429    increases between-tool agreement and consistency between sites and vendors.

430    <u>Sources of variance</u>

431    In order to understand the substantial variability introduced by the choice of analysis tool, the in-

432    fluence of modelling strategies and parameters on quantitative results needs to be better under-

433    stood. Previous investigations have shown that, within a given LCM algorithm, metabolite esti-

434    mates can be affected by the choice of baseline knot spacing[36,37], the modelling of MM and lipids

435    [36,38], and SNR and linewidth[39–42]. In this study, we focused on the comparison of each LCM with

436    their default and commonly used parameters, and observed differences resulting both from the

437    default parameters and from differences in the core algorithm. Minor differences in spectral qual-

438    ity (SNR and LW) were found between vendors. The agreement between vendors was high for

439    the mean metabolite levels and the cohort-mean correlations. Further vendor-specific effects on

440    the LCM estimation of this dataset are described elsewhere[17].

441    LCM relies on the assumption that broad background and baseline signals can be separated from

442    narrower metabolite signals. This is true to a limited degree, and the choice of MM and baseline

443    modelling influences the quantification of metabolite resonances[4]. Our secondary analysis of the

444    relationship between baseline power and metabolite estimates showed a stronger interaction for

445    the broader coupled signals of Glx and mI than the singlets. tCho showed the weakest effect, and

446    the three LCMs showed the highest agreement between the MM+baseline models around 3.2

447    ppm. The higher variance of Glx and mI estimates may at least partly be explained by the ab-

448    sence of MM basis functions for frequencies >3 ppm in the model. MM signal must therefore

449     either be modelled by metabolite basis functions or the spline baseline. Including experimental

450     MM acquisitions into studies may reduce the degrees of freedom of modelling, but introduce

451     other sources of variance, such as age-dependency[43] or tissue composition[38,44]. While consensus

452     is emerging that such approaches are recommended many open questions must be resolved be-

453     fore the recommendations can be broadly implemented[25].

454     For all three LCM algorithms, optimization between the model and the data is solved by local

455     optimization. Algorithms could converge on a local minimum, if the search space of the non-

456     linear parameters is of high dimensionality, or if the starting values of the parameters are far

457     away from the global optimum[45]. The availability of open-source LCM such as Tarquin and Os-

458     prey will allow further investigation of the relationship between optimization starting values and

459     modelling outcomes.

460

461     Since this study focused on reporting tCr ratios, it is important to consider the variance of the

462     creatine model of each algorithm. With MRS only quantitative in a relative sense, separating the

463     variance contribution of the reference signal is a challenge. While mean tCr model areas were

464     slightly higher for LCModel than for Osprey and Tarquin, there was no generalizable observa-

465     tion of lower tCr ratios from LCModel. CVs of the tCr model areas were comparable across

466     LCM algorithms for each vendor. Vendor differences in water suppression of each vendor were

467     accounted for by limiting the analysis range to 0.5 to 4 ppm, and by including a -$CrCH_2$ correc-

468     tion term (omitted from calculations of the tCr ratios and the secondary analysis of the tCr mod-

469     els). The contribution of the reference signal to the variance of metabolite estimates is unclear

470     and hard to isolate. Nevertheless, tCr referencing was preferred in this study, since water refer-

471     encing is likely to add additional tool-specific variance resulting from water amplitude estima-

472     tion.

473

474     <u>Limitations</u>

475     As mentioned in greater detail above, there is currently no widely adopted consensus on the defi-

476     nition of MM basis functions, and measured MM background data are not widely available to

477     non-expert users. To reflect common practice in current MRS applications, the default MM basis

478     function definitions from LCModel were adapted for each algorithm in this study. These basis

479     functions only included MMs for frequencies < 3.0 ppm, which is likely insufficient for the

480    modelling of MM signals between 3 and 4 ppm[46], and will have repercussions for the estimation

481    of tCho, mI, and Glx. Second, standard modelling parameters were chosen for each LCM, which

482    ensure a broader comparability to the current literature, but may not be ideal. Third, there is ob-

483    viously no 'gold standard' of metabolite level estimation to validate MRS results against. The

484    performance of an algorithm is often judged based on the level of variance, but low variance

485    clearly does not reflect accuracy and may indicate insufficient responsiveness of a model to the

486    data. In comparing multiple algorithms, it is tempting to infer algorithms that show a higher de-

487    gree of correlation in results are more reliable, but it could equally be the case that shared algo-

488    rithm-based sources of variance increase such correlations. Efforts to use simulated spectra as a

489    gold-standard, including those applying machine learning [47,48], can only be successful to the ex-

490    tent that simulated data are truly representative of in-vivo data. Fourth, another criterion to judge

491    the performance of an algorithm is the residual. For example, a small residual indicates a higher

492    agreement between the complete model and the data for LCModel, it does not infer a better esti-

493    mation of individual metabolites, and may result from the higher degree of freedom in the base-

494    line of LCModel (higher number of splines) compared to Osprey and Tarquin. This is empha-

495    sized by the high agreement of the mean mI models, but lower agreement of the baseline models

496    around 3.58 ppm between LCModel and Osprey. Fifth, this study was limited to the two most

497    widely used algorithms LCModel and Tarquin, as well as the Osprey algorithm that is under on-

498    going development in our group. While including additional algorithms would increase the gen-

499    eral understanding of different algorithms, the complexity of the resulting analysis and interpre-

500    tation would be overwhelming and beyond the scope of a single publication.

501

## **<u>Conclusion</u>**

502

503     This study presents a comparison of three LCM algorithms applied to a large short-TE PRESS

504     dataset. While different LCM algorithms' estimates of major metabolite levels agree broadly at a

505     group level, correlations between results are only weak-to-moderate, despite standardized pre-

506     processing, a large sample of young, healthy and cooperative subjects, and high spectral quality.

507     The variability of metabolite estimates that is introduced by the choice of analysis software is

508     substantial, raising concerns about the robustness of  MRS research findings, which typically use

509     a single algorithm to draw inferences from much smaller sample sizes.

510

## **<u>Acknowledgement</u>**

# References

1. Öz G, Alger JR, Barker PB, et al. Clinical Proton MR Spectroscopy in Central Nervous System Disorders. *Radiology*. 2014;270(3):658-679. doi:10.1148/radiol.13130531

2. Wilson M, Andronesi O, Barker PB, et al. Methodological consensus on clinical proton MRS of the brain: Review and recommendations. *Magn Reson Med*. 2019;82(2):527–550. doi:10.1002/mrm.27742

3. Bottomley P. *Selective Volume Method for Performing Localized NMR Spectroscopy*. Vol 3.; 1985. doi:10.1016/0730-725X(85)90032-3

4. Near J, Harris AD, Juchem C, et al. Preprocessing, analysis and quantification in single-voxel magnetic resonance spectroscopy: experts' consensus recommendations. *NMR Biomed*. 2020;n/a(n/a):e4257. doi:10.1002/nbm.4257

5. Oeltzschner G, Zöllner HJ, Hui SCN, et al. Osprey: Open-source processing, reconstruction & estimation of magnetic resonance spectroscopy data. *J Neurosci Methods*. 2020;343:108827. doi:10.1016/j.jneumeth.2020.108827

6. Juchem C. INSPECTOR - A Tool for Teaching Magnetic Resonance Spectroscopy. In: *26th Annual Meeting of the International Society for Magnetic Resonance in Medicine (ISMRM)*. Paris, France; 2018.

7. Wilson M, Reynolds G, Kauppinen RA, Arvanitis TN, Peet AC. A constrained least-squares approach to the automated quantitation of in vivo 1 H magnetic resonance spectroscopy data. *Magn Reson Med*. 2011;65(1):1–12. doi:10.1002/mrm.22579

8. Poullet J-B, Sima DM, Simonetti AW, et al. An automated quantitation of short echo time MRS spectra in an open source software environment: AQSES. *NMR Biomed*. 2007;20(5):493–504. doi:10.1002/nbm.1112

9. Soher BJ, Semanchuk P, Todd D, Steinberg J, Young K. VeSPA: Integrated applications for RF pulse design, spectral simulation and MRS data analysis. In: *19th Annual Meeting of the International Society for Magnetic Resonance in Medicine (ISMRM)*. Montreal, Canada; 2011. https://cds.ismrm.org/protected/11MProceedings/files/1410.pdf. Accessed May 19, 2020.

10. Graveron-Demilly D. Quantification in magnetic resonance spectroscopy based on semi-parametric approaches. *Magn Reson Mater Phys Biol Med*. 2014;27(2):113-130. doi:10.1007/s10334-013-0393-4

11. Provencher SW. Estimation of metabolite concentrations from localized in vivo proton NMR spectra. *Magn Reson Med*. 1993;30(6):672–679. doi:10.1002/mrm.1910300604

12. Osorio-Garcia MI, Sima DM, Nielsen FU, Himmelreich U, Huffel SV. Quantification of magnetic resonance spectroscopy signals with lineshape estimation. *J Chemom*. 2011;25(4):183-192. doi:10.1002/cem.1353

552 13. Shen ZW, Chen YW, Wang HY, et al. Quantification of Metabolites in Swine Brain by ^1H
553 MR Spectroscopy Using LCModel and QUEST: A Comparison Study. In: *2008 Congress*
554 *on Image and Signal Processing*. Vol 5. ; 2008:299-302. doi:10.1109/CISP.2008.478

555 14. Kossowski B, Orzeł J, Bogorodzki P, Wilson M, Setkowicz Z, P. Gazdzinski S. Follow-up
556 analyses on the effects of long-term use of high fat diet on hippocampal metabolite concen-
557 trations in Wistar rats: Comparing Tarquin quantification of 7.0T rat metabolites to
558 LCModel. *Biol Eng Med*. 2017;2(4). doi:10.15761/BEM.1000129

559 15. Mosconi E, Sima DM, Garcia MIO, et al. Different quantification algorithms may lead to
560 different results: a comparison using proton MRS lipid signals. *NMR Biomed*.
561 2014;27(4):431-443. doi:10.1002/nbm.3079

562 16. Scott J, Underwood J, Garvey LJ, Mora-Peris B, Winston A. A comparison of two post-pro-
563 cessing analysis methods to quantify cerebral metabolites measured via proton magnetic
564 resonance spectroscopy in HIV disease. *Br J Radiol*. 2016;89(1060):20150979.
565 doi:10.1259/bjr.20150979

566 17. Považan M, Mikkelsen M, Berrington A, et al. Comparison of Multivendor Single-Voxel
567 MR Spectroscopy Data Acquired in Healthy Brain at 26 Sites. *Radiology*.
568 2020;295(1):191037. doi:10.1148/radiol.2020191037

569 18. Mikkelsen M, Barker PB, Bhattacharyya PK, et al. Big GABA: Edited MR spectroscopy at
570 24 research sites. *NeuroImage*. 2017;159:32–45. doi:10.1016/j.neuroimage.2017.07.021

571 19. Big GABA repository. Big GABA repository. https://www.nitrc.org/projects/biggaba/. Pub-
572 lished 2018. Accessed May 27, 2020.

573 20. Hall EL, Stephenson MC, Price D, Morris PG. Methodology for improved detection of low
574 concentration metabolites in MRS: Optimised combination of signals from multi-element
575 coil arrays. *NeuroImage*. 2014;86:35-42. doi:10.1016/j.neuroimage.2013.04.077

576 21. Klose U. In vivo proton spectroscopy in presence of eddy currents. *Magn Reson Med*.
577 1990;14(1):26–30. doi:10.1002/mrm.1910140104

578 22. Mikkelsen M, Tapper S, Near J, Mostofsky SH, Puts NAJ, Edden RAE. Correcting fre-
579 quency and phase offsets in MRS data using robust spectral registration. *NMR Biomed*. July
580 2020:e4368. doi:10.1002/nbm.4368

581 23. Barkhuijsen H, de Beer R, van Ormondt D. Improved algorithm for noniterative time-do-
582 main model fitting to exponentially damped magnetic resonance signals. *J Magn Reson*
583 *1969*. 1987;73(3):553–557. doi:10.1016/0022-2364(87)90023-0

584 24. Simpson R, Devenyi GA, Jezzard P, Hennessy TJ, Near J. Advanced processing and simu-
585 lation of MRS data using the FID appliance (FID-A)—An open source, MATLAB-based
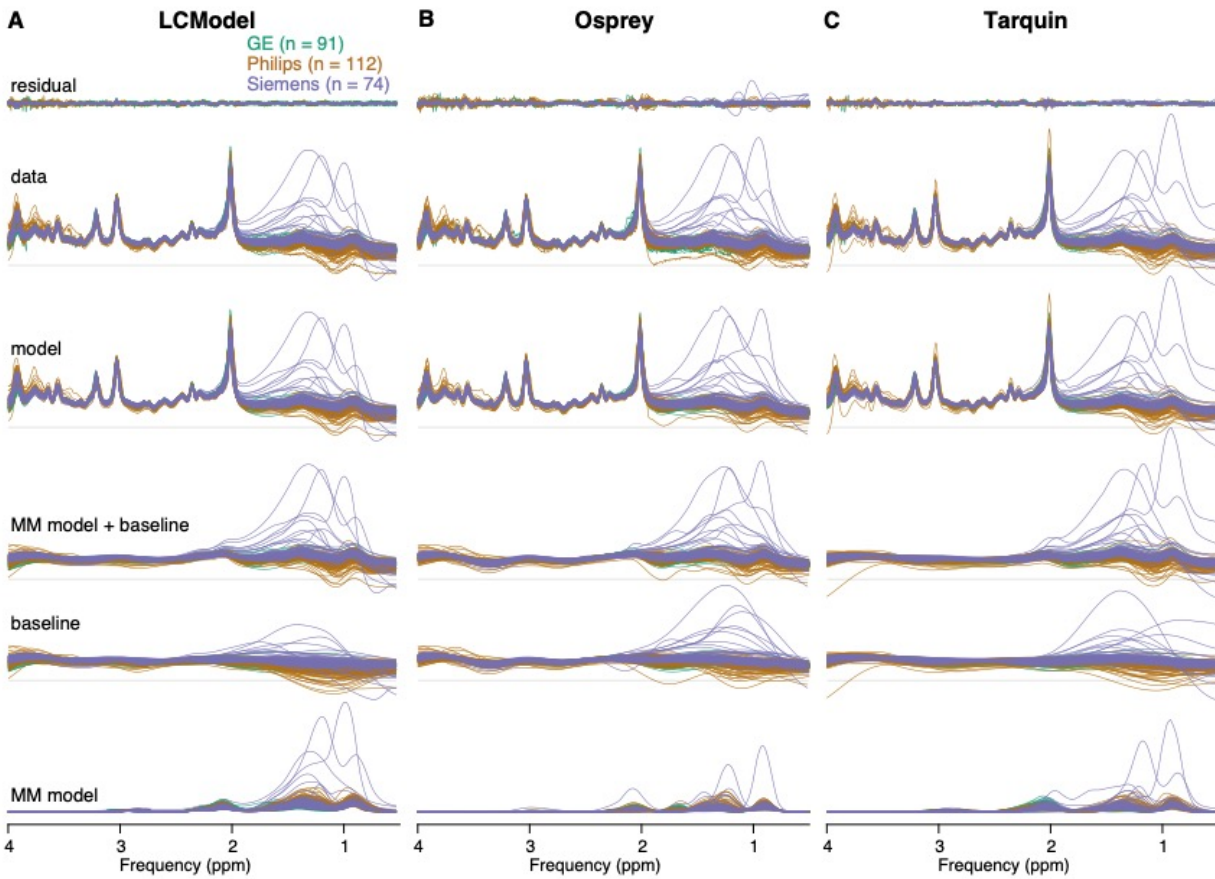586 toolkit. *Magn Reson Med*. 2017;77(1):23–33. doi:10.1002/mrm.26091

587  25.  Cudalbu C, Behar KL, Bhattacharyya PK, et al. Contribution of macromolecules to brain
588       1H MR spectra: Experts' consensus recommendations. *NMR Biomed Revis*. 2020.

589  26.  Provencher S. LCModel & LCMgui User's Manual. LCModel & LCMgui User's Manual.
590       http://s-provencher.com/pub/LCModel/manual/manual.pdf. Published 2020. Accessed July
591       15, 2020.

592  27.  Levenberg K. A method for the solution of certain non-linear problems in least squares. *Q*
593       *Appl Math*. 1944;2(2):164-168. doi:10.1090/qam/10666

594  28.  Marquardt DW. An Algorithm for Least-Squares Estimation of Nonlinear Parameters. *J Soc*
595       *Ind Appl Math*. 1963;11(2):431-441. doi:10.1137/0111030

596  29.  R Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria:
597       R Foundation for Statistical Computing; 2017. https://www.R-project.org/.

598  30.  SpecVis GitHub repository. SpecVis GitHub repository.
599       https://github.com/hezoe100/SpecVis. Published 2020. Accessed May 27, 2020.

600  31.  Zöllner HJ. Comparison of algorithms for linear-combination modelling of short-echo-time
601       magnetic resonance spectra. https://osf.io/3ekq4/. Published June 1, 2020. Accessed June 2,
602       2020.

603  32.  https://github.com/martin3141/spant. spant GitHub repository. https://github.com/mar-
604       tin3141/spant. Published 2017. Accessed May 27, 2020.

605  33.  Allen M, Poggiali D, Whitaker K, Marshall TR, Kievit RA. Raincloud plots: a multi-plat-
606       form tool for robust data visualization. *Wellcome Open Res*. 2019;4:63. doi:10.12688/well-
607       comeopenres.15191.1

608  34.  Wickham H. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York;
609       2009. http://ggplot2.org.

610  35.  Fagerland MW. T-tests, non-parametric tests, and large studiesa paradox of statistical prac-
611       tice? *BMC Med Res Methodol*. 2012;12(1):78. doi:10.1186/1471-2288-12-78

612  36.  Marjańska M, Terpstra M. Influence of fitting approaches in LCModel on MRS quantifica-
613       tion focusing on age-specific macromolecules and the spline baseline. *NMR Biomed*. No-
614       vember 2019. doi:10.1002/nbm.4197

615  37.  Wenger KJ, Hattingen E, Harter PN, et al. Fitting algorithms and baseline correction influ-
616       ence the results of non-invasive in vivo quantitation of 2-hydroxyglutarate with 1H-MRS.
617       *NMR Biomed*. 2019;32(1):e4027. doi:10.1002/nbm.4027

618  38.  Schaller B, Xin L, Gruetter R. Is the macromolecule signal tissue-specific in healthy human
619       brain? A \textlesssup\textgreater1\textless/sup\textgreater H MRS study at 7 tesla in the oc-
620       cipital lobe. *Magn Reson Med*. 2014;72(4):934–940. doi:10.1002/mrm.24995

621   39. Bartha R. The Effect of Signal to Noise Ratio and Linewidth On 4T Short Echo Time 1H
622          MRS Metabolite Quantification. *Proc 13th Sci Meet Int Soc Magn Reson Med*.
623          2005;216(1):2459–2459.

624   40. Near J. Investigating the effect of spectral linewidth on metabolite measurement bias in
625          short-TE MRS. In: *21th Annual Meeting of the International Society for Magnetic Reso-*
626          *nance in Medicine (ISMRM)*. Milan, Italy; 2014.

627   41. Wijtenburg SA, Knight-Scott J. The Impact of SNR on the Reliability of LCModel and
628          QUEST Quantitation in 1 H-MRS. In: *17th Annual Meeting of the International Society for*
629          *Magnetic Resonance in Medicine (ISMRM)*. ; 2009.

630   42. Zhang Y, Shen J. Effects of noise and linewidth on in vivo analysis of glutamate at 3 T. *J*
631          *Magn Reson*. 2020;314. doi:10.1016/j.jmr.2020.106732

632   43. Marjańska M, Deelchand DK, Hodges JS, et al. Altered macromolecular pattern and con-
633          tent in the aging human brain. *NMR Biomed*. 2018;31(2):e3865. doi:10.1002/nbm.3865

634   44. Považan M, Strasser B, Hangel G, et al. Simultaneous mapping of metabolites and individ-
635          ual macromolecular components via ultra-short acquisition delay 1H MRSI in the brain at
636          7T. *Magn Reson Med*. 2018;79(3):1231-1240. doi:10.1002/mrm.26778

637   45. Poullet J-B, Sima DM, Van Huffel S. MRS signal quantitation: A review of time- and fre-
638          quency-domain methods. *J Magn Reson*. 2008;195(2):134-144.
639          doi:10.1016/j.jmr.2008.09.005

640   46. Giapitzakis I-A, Avdievich N, Henning A. Characterization of macromolecular baseline of
641          human brain using metabolite cycled semi-LASER at 9.4T. *Magn Reson Med*.
642          2018;80(2):462-473. doi:10.1002/mrm.27070

643   47. Lee HH, Kim H. Deep learning-based target metabolite isolation and big data-driven meas-
644          urement uncertainty estimation in proton magnetic resonance spectroscopy of the brain.
645          *Magn Reson Med*. 2020;n/a(n/a). doi:10.1002/mrm.28234

646   48. Lee HH, Kim H. Intact metabolite spectrum mining by deep learning in proton magnetic
647          resonance spectroscopy of the brain. *Magn Reson Med*. 2019;82(1):33-48.
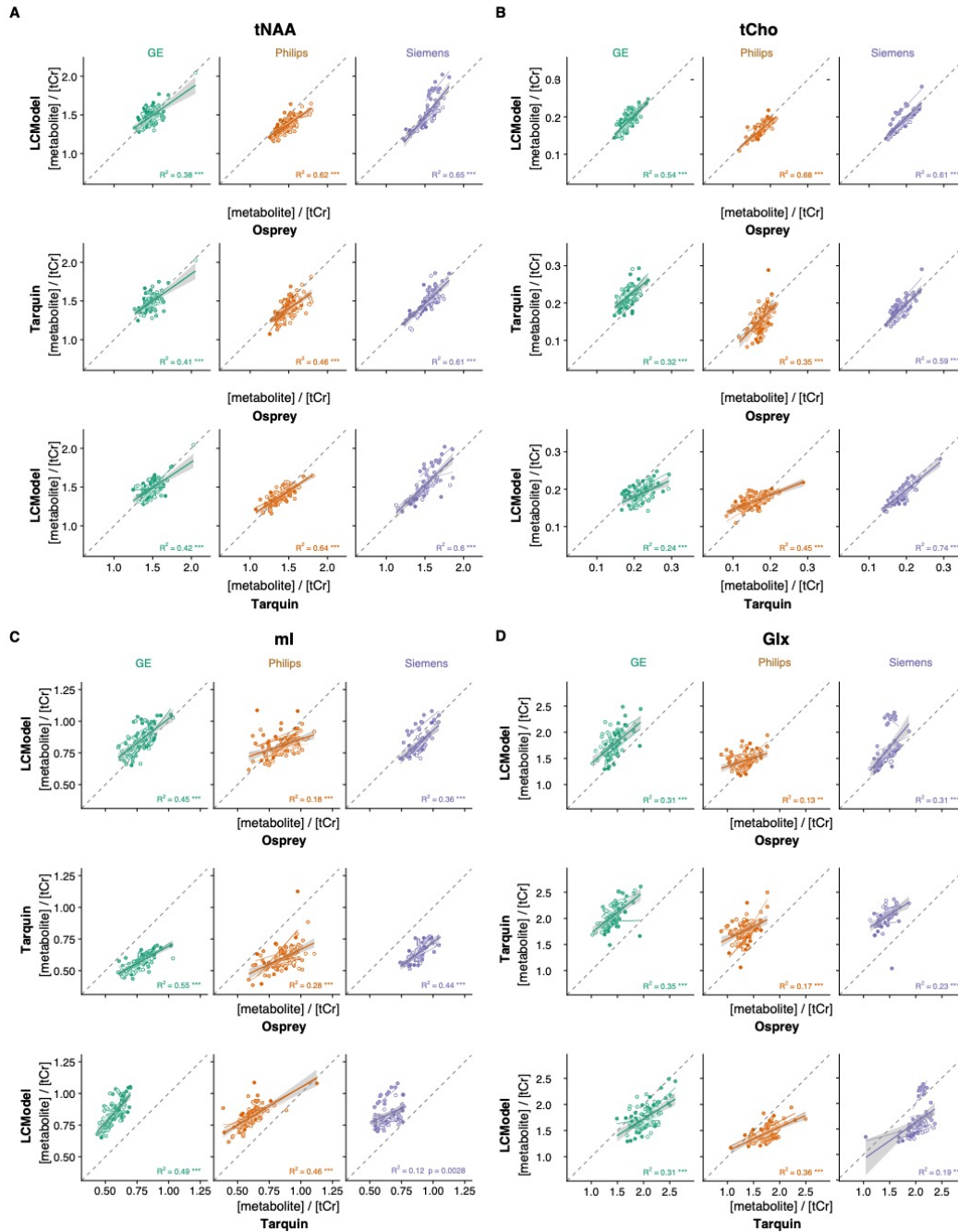648          doi:10.1002/mrm.27727

649

# Supplementary Material

| Name | Frequencies [ppm] | FWHM [ppm] | Amplitude |
|---|---|---|---|
| MM09 | 0.91 | 0.14 | 3.00 |
| MM12 | 1.21 | 0.15 | 2.00 |
| MM14 | 1.43 | 0.17 | 2.00 |
| MM17 | 1.67 | 0.15 | 0.20 |
| MM20 | 2.08 | 0.15 | 1.33 |
|  | 2.25 | 0.20 | 0.33 |
|  | 1.95 | 0.15 | 0.33 |
|  | 3.00 | 0.20 | 0.40 |
| Lip09 | 0.89 | 0.14 | 3.00 |
| Lip13a | 1.28 | 0.15 | 2.00 |
| Lip13b | 1.28 | 0.089 | 2.00 |
| Lip20 | 2.04 | 0.15 | 1.33 |
|  | 2.25 | 0.15 | 0.67 |
|  | 2.80 | 0.20 | 0.87 |

**Supplementary Material 1**. *Properties of the Gaussian functions of the broad macromolecule and lipid resonances included in the basis sets, taken from section 11.7 of the LCModel manual. The amplitude values are scaled relative to the CH$_3$ singlet of creatine with amplitude 3.*

***Supplementary Material 2.*** *Summary of the individual modelling results. (A–C) individual residuals, data, models, MM models + baseline, baseline and MM models for each LCM algorithm, color-coded by vendor.*

***Supplementary Material 3.*** *Facetted pair-wise correlational comparison of algorithms. LCModel and Osprey are compared in the first row, Tarquin and Osprey are compared in the second row, and LCModel and Tarquin are compared in the third row. Each sub-plot (A-D) corresponds to a different metabolite. Within-vendor (bold line with confidence interval) and within-site (thin line) correlations are color-coded. Asterisks indicate significant correlations (adjusted p < 0.01 = \*\* and adjusted p < 0.001 = \*\*\*).*