

1 **Evaluating the effect of data merging and post-acquisition normalization on statistical analysis of**  
2 **untargeted high-resolution mass spectrometry based urinary metabolomics data**

3 Fynn Brix,<sup>†,\*</sup> Tobias Demetrowitsch,<sup>†</sup> Julia Jensen-Kroll,<sup>†</sup> Helena Zacharias,<sup>∇,\*</sup> Silke Szymczak,<sup>§</sup>  
4 Matthias Laudes,<sup>#,||</sup> Stefan Schreiber,<sup>#,||</sup> Karin Schwarz<sup>†</sup>

5 <sup>†</sup>Institute of Human Nutrition and Food Science, Kiel University, Kiel, Heinrich-Hecht-Platz 10, 24118  
6 Kiel, Germany

7 <sup>∇</sup>Peter L. Reichertz Institute for Medical Informatics of TU Braunschweig and Hannover Medical School,  
8 30625 Hannover, Germany

9 <sup>§</sup>Institute of Medical Biometry and Statistics, University of Luebeck and Medical Centre Schleswig-  
10 Holstein, Campus Luebeck, 23562 Luebeck, Germany

11 <sup>#</sup>Department of Internal Medicine I, University Medical Centre Schleswig-Holstein, Campus Kiel, Kiel,  
12 Germany

13 <sup>‡</sup>Institute of Clinical Molecular Biology, Kiel University and University Medical Center Schleswig-  
14 Holstein, Campus Kiel, 24105 Kiel, Germany

15 <sup>||</sup>Institute of Diabetes and Clinical Metabolic Research, Kiel University, Düsternbrooker Weg 17, 24105,  
16 Kiel, Germany

17 **ABSTRACT:** Urine is one of the most widely used biofluids in metabolomic studies, because it can be  
18 collected non-invasively and is available in large quantities. However, it shows large heterogeneity in  
19 sample concentration and consequently requires normalization to reduce unwanted variation and  
20 extract meaningful biological information. Biological samples like urine are commonly measured with  
21 electrospray ionization (ESI) coupled to a mass spectrometer, producing datasets for positive and  
22 negative mode. Combining these gives a more complete picture of the total metabolites present in a  
23 sample. However, the effect of this data merging on subsequent data analysis, especially in  
24 combination with normalization, has not yet been analysed. To address this issue, we conducted a  
25 neutral comparison study to evaluate the performance of eight post-acquisition normalization  
26 methods under different data merging procedures using 1029 urine samples from the Food Chain plus  
27 (FoCUS) cohort. Samples were measured by a Fourier transform ion cyclotron resonance mass  
28 spectrometer (FT-ICR-MS). Normalization methods were evaluated by five criteria capturing the ability  
29 to remove sample concentration variation and preserve relevant biological information. Merging data  
30 after normalization was generally favourable for quality control (QC) sample similarity, sample  
31 classification and feature selection for most of the tested normalization methods. Merging data after  
32 normalization and the usage of probabilistic quotient normalization (PQN) in a similar setting are  
33 generally recommended. Relying on a single analyte to capture sample concentration differences, like  
34 with post-acquisition creatinine normalization, seems to be a less preferable approach, especially  
35 when data merging is applied.

36

37 Urine is one of the most widely used biofluids in metabolomic studies, because it can be collected non-  
38 invasively and is available in large quantities.<sup>1</sup> However, it has a large heterogeneity in sample  
39 concentration<sup>2</sup> and volume may change up to 15-fold under normal conditions.<sup>3</sup>

40 Consequently, numerous normalization methods have been developed to reduce variation originating  
41 from unwanted factors.<sup>1</sup> These normalization methods may be categorized as being pre- or post-  
42 acquisition.<sup>2</sup> Pre-acquisition methods adjust the sample volumes based on measured reference  
43 quantities. These parameters can also be used in post-acquisition methods, but requires an additional  
44 workflow step and information on sample volume is prerequisite.<sup>4</sup>

45 Post-acquisition normalization methods are applied after data collection.<sup>5</sup> Several evaluation studies  
46 have been conducted for human urine,<sup>2,6-8</sup> serum,<sup>9</sup> and plasma,<sup>10</sup> as well as animal urine,<sup>11,3,12</sup>  
47 measured by nuclear magnetic resonance spectroscopy or liquid chromatography-mass spectrometry.  
48 Frequently used methods include PQN,<sup>11,13,2,6,8,5</sup> as well as variance stabilization normalization  
49 (VSN),<sup>6,8,14</sup> and quantile normalization.<sup>6,7,14</sup>

50 Most post-acquisition normalization methods adjust for sample-to-sample variation, whereas VSN  
51 additionally adjusts for variation on the metabolite level.<sup>7,8</sup> All of the above normalization methods  
52 make specific assumptions about the data. Whether or not these assumptions are met for a specific  
53 data set may influence the performance of the normalization methods.<sup>15</sup> Thus, the selection of a  
54 particular normalization method should be made based on the data characteristics, research question  
55 and the subsequent data analysis methods.<sup>13</sup>

56 Studies comparing post-acquisition normalization methods with urine samples and mass spectrometry  
57 (MS) with ESI are either using positive ionization of molecules only,<sup>16-18</sup> or include datasets of both  
58 polarities.<sup>19,8,12</sup> However, none used merged datasets from both polarities. Furthermore, most studies  
59 are based on small sample sizes with < 100 samples,<sup>19,16,17,2,6,7,20,5</sup> and only few are based on data sets  
60 with > 1000 samples.<sup>8,21</sup> The combined effect of data merging and normalization has not yet been  
61 evaluated in the workflow for pre-processing metabolomics data. Thus, in this neutral comparison  
62 study we aim to objectively evaluate the performance of eight post-acquisition normalization methods  
63 under different data merging procedures based on a large-scale urinary metabolomics dataset. Based  
64 on our results we will provide recommendations for different data merging procedures and  
65 normalization methods in different scenarios.

## 66 MATERIALS AND METHODS

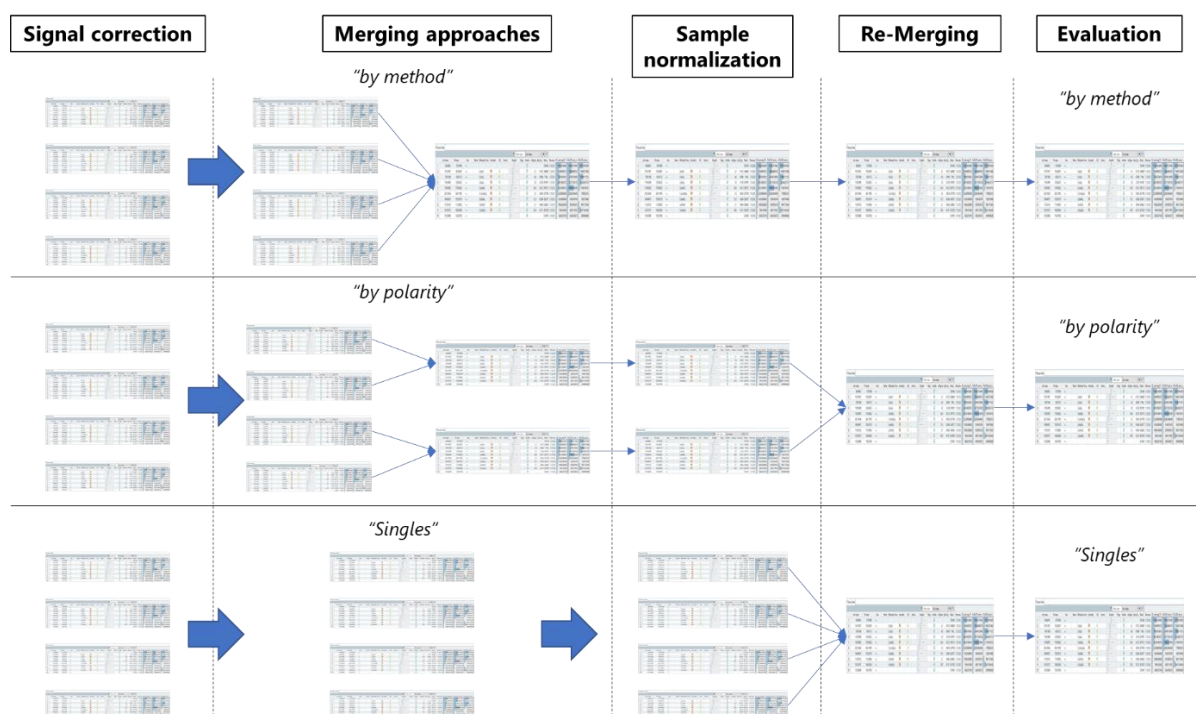
67 **Human Urine Samples and Sample preparation.** Urine metabolomics profiling was performed in the  
68 Food Chain Plus (FoCUS) cohort, which has been published recently.<sup>22</sup> The cohort was established in  
69 2011 for population-based research with a focus on metabolic inflammation. The study was approved  
70 by the local ethics committee of the Kiel University (A156-03/Date 2011/07/28) and was registered  
71 under the clinical trial number DRKS00005285 at the German Clinical Trials Register in Cologne.<sup>22</sup> The  
72 average age of the study participants who gave urine is 52 years with biological sex of 40% males and  
73 60% females.

74 Spot urine samples of 1031 participants were available; two samples were excluded, because  
75 creatinine was lacking or data for one analytical method were missing. Samples were diluted 1:500  
76 with methanol and water (50:50, v/v) prior to analysis. The preparation of the quality control (QC)  
77 samples followed a procedure modified from a previous publication.<sup>23</sup>

78 **Data acquisition.** Data were acquired using a 1260 Infinity HPLC (Agilent, Waldbronn, Germany) for  
79 direct injection of samples. The HPLC was linked to an ultrahigh-resolution Fourier transform ion  
80 cyclotron resonance mass spectrometer (FT-ICR-MS) (7T, SolarixR, Bruker, Bremen, Germany). Mass  
81 spectra were acquired with electrospray ionization (ESI) source in both modes (positive and negative  
82 ionization) and two methods with a mass range between 65 and 1500 Da (USM, SM), leading to four  
83 data sets per sample. The intensity threshold was 10<sup>6</sup> counts. Data were calibrated using an in-house  
84 database in the quadratic mode with a tolerated mass error < 0.5 mDa. Detailed instrumental  
85 parameters can be found in an earlier publication.<sup>24</sup> Blank samples were injected prior to each batch.  
86 Pooled QC samples were injected at the start and end of each batch.

87 **Data processing and merging.** Bruker raw data (.d) were processed in the MetaboScape 2021b  
88 software (Bruker, Bremen, Germany). Ion deconvolution and other settings are given in Table S1. Most  
89 probable chemical formulas were assigned based on accurate measured masses, isotopic patterns, as  
90 well as the seven golden rules.<sup>25</sup> Potential compound names were assigned based on the Human  
91 Metabolome Database entries (version 5.0) for all human biospecimens.<sup>26</sup> The raw Bucket tables were

92 exported and further processed in R (version 4.1.2). Figure 1 shows a schematic depiction of the  
 93 different analysis workflows. Signal correction, peak filtering and imputation of missing values were  
 94 conducted using the R package *statTarget* version 1.24.0.<sup>27</sup> Signal correction was applied using QC-  
 95 based random forest signal correction (QC-RFSC). Peak filtering was conducted by only including  
 96 compounds that were detected in at least 80% of the samples. Imputation of missing values was done  
 97 using the k-Nearest Neighbour (kNN) method.<sup>27</sup>  
 98 After the signal correction the datasets were merged. Data merging was done in three different ways:  
 99 merging all four analytical datasets into one (“by-method”), merging positive and negative data into  
 100 two datasets (“by-polarity”) and no merging at all prior to normalisation (“singles”). Merging of  
 101 compounds detected in multiple datasets was done by selecting the compound with either the lowest  
 102 number of missing values, the highest median intensity, or the highest mean intensity calculated using  
 103 all samples. In case of ties the criteria are applied in the described order. Compounds only detected in  
 104 one of the datasets are simply added to the respective merged dataset. For the next step, each of the  
 105 merged datasets was corrected with each of the normalisation methods. Subsequently, all datasets  
 106 were merged into one dataset for comparison.



107

108

109 **Figure 1:** Schematic depiction of the analysis workflows. First, signal corrected data are merged  
 110 either to one set, two sets or kept as four individual sets. Different normalization methods are  
 111 applied and data are merged subsequently and evaluated by different criteria.

112

113 **Evaluated normalization methods.** Post-acquisition normalization methods evaluated were:  
 114 creatinine, sum, MS total useful signal (MSTUS),<sup>12</sup> PQN<sup>28</sup> with median of all QC samples as reference  
 115 spectrum and all samples as reference, quantile,<sup>29</sup> median, VSN,<sup>30</sup> and cubic spline normalization  
 116 (CSN)<sup>31</sup>. A description for each method and the R packages which have been used can be found in  
 117 Table S2. Baseline data refers to data corrected using a QC-based random forest signal correction  
 118 method.<sup>27</sup>

119 **Evaluation criteria.** The performance evaluation of the normalization methods emphasized on  
 120 removal of sample concentration variation and preservation of biological information. Sample  
 121 concentration variation removal was assessed by QC clustering in principal component analysis (PCA)

122 and the median of the relative standard deviations (RSD) of all metabolites in QC samples. QC  
 123 clustering was determined in a quantitative way by calculating the average auto scaled distances of  
 124 each QC sample to the centroid of all QC samples in the score plot.

125 Preservation of biological information was evaluated by training random forest (RF) models using the  
 126 R package ranger<sup>32</sup> with default parameter settings to predict subjects' sex. To identify compounds  
 127 differentiating for sex a Web of Science literature search was conducted. For literature search  
 128 keywords and the resulting 41 compounds see Tables S3 and S4.

## 129 RESULTS AND DISCUSSION

130 Common post-acquisition normalization methods were evaluated under different merging procedures  
 131 by the ability to remove sample concentration variation and to preserve biological information. The  
 132 first criterion was tested via reduction of QC sample variation and the latter via comparing prediction  
 133 accuracy, number of significant metabolites, and matches with previously reported differing  
 134 metabolites for biological sex.

135 **Comparison of merging procedures.** All normalization methods except for creatinine achieved the  
 136 lowest average distance among QC samples after the “singles” merging procedure (Table 1).  
 137 Furthermore, almost all normalization methods except for creatinine exhibited a monotonic decline of  
 138 the QC distance. This trend is also reflected in the total distance score for all normalisation methods  
 139 combined, which ranged from 85 for “by-method” to 56 for “by-polarity” and down to 26 for “singles”  
 140 merging procedure. Results for the RSD of QCs were similar to those of the QC distance, since all  
 141 normalization methods except for one (here CSN) achieved the lowest average RSD value after  
 142 “singles” merging. In addition, a monotonic decline in the RSD was also observed for five of the  
 143 normalization methods, namely the MSTUS, creatinine, PQN(QC), sum, and median normalization. A  
 144 systematic decline was also observed, calculating the total of the RSD values. Merging procedures “by-  
 145 method” and “by-polarity” did not differ much, with values of 120 and 118 respectively. However, this  
 146 value decreased to 107 for the “singles” merging procedure. The smallest change in the QC distance  
 147 and average RSD in QCs was found for VSN normalized data. This is because the variance stabilization  
 148 reduces variability and the data is on a different scale after the normalization.

149  
 150 **Table 1: Average distance of QCs in the PCA score plots and average RSD of features in QCs, for each**  
 151 **normalization method and merging procedure<sup>a</sup>**

Normalization methods	PCA QC aDist			QC RSD		
	By-method	By-polarity	Singles	By-method	By-polarity	Singles
Baseline	2.57	2.6	2.61	5.55	5.55	5.55
Creatinine	<b>0.6</b>	4.65	1.91	11.98	10.58	8.33
Sum	2.35	1.4	0.75	6.94	6.00	4.27
MSTUS	2.06	1.41	1.21	6.51	5.16	5.13
Median	2.02	1.32	0.56	5.73	5.22	4.39
PQN	1.93	1.61	0.60	6.84	7.26	4.42
PQN(QC)	2.03	<b>1.28</b>	<b>0.53</b>	5.72	5.15	4.31
Quantile	32.22	16.32	7.5	53.14	54.56	52.28
CSN	36.93	22.50	7.9	17.68	18.14	17.73
VSN	3.06	3.04	3.00	<b>0.83</b>	<b>0.85</b>	<b>0.82</b>
<i>Total:</i>	<i>85.77</i>	<i>56.13</i>	<i>26.57</i>	<i>120.92</i>	<i>118.47</i>	<i>107.23</i>

152 <sup>a</sup>QC aDist, Average distance of QCs; QC RSD, average RSD of features in QCs. Bold letters indicate best  
 153 values for the respective merging procedure (column).

154 The balanced accuracy of the random forest models predicting the sex of subjects using the differently  
 155 merged and normalized data is shown in Table 2. The highest balanced accuracy was achieved by  
 156 models built using data from the “singles” merging procedure. MSTSU and VSN were the exceptions,  
 157 which had the best values after merging according to polarity. Also, the majority of normalization  
 158 methods attained the most significant metabolites for predicting sex (Table 5) after the “singles”  
 159 merging procedure. Moreover, there is also a positive increasing trend for the average number of  
 160 significant metabolites calculated including all normalization methods for each respective merging  
 161 procedure. Normalization methods benefitting from “singles” merging include creatinine, PQN(QC),  
 162 sum, median and CSN. Merging “by-polarity” was the best procedure for MSTUS normalization.  
 163 Merging “by-method” resulted in the highest number of significant features for PQN and VSN  
 164 normalized data.

165

166 **Table 2: Test set based average balanced accuracy for predicting the subjects’ sex<sup>a</sup>**

Normalization methods	Average balanced accuracy (SD)		
	By method	By polarity	Singles
Baseline	74.68 (3.50)	75.58 (2.99)	74.67 (3.05)
Creatinine	69.81 (2.82)	71.91 (1.31)	73.05 (3.90)
Sum	76.33 (2.60)	76.02 (3.59)	76.49 (3.08)
MSTUS	74.77 (4.46)	75.77 (3.58)	74.26 (3.34)
Median	76.26 (3.70)	76.05 (3.56)	<b>77.98</b> <b>(3.37)</b>
PQN	76.35 (2.71)	75.91 (3.61)	77.42 (2.93)
PQN(QC)	75.86 (3.19)	76.14 (2.83)	76.64 (3.24)
Quantile	<b>76.65</b> <b>(3.02)</b>	<b>76.96</b> <b>(2.91)</b>	77.62 (3.18)
CSN	76.35 (2.19)	76.34 (3.43)	77.46 (2.99)
VSN	74.75 (3.29)	75.97 (3.63)	75.38 (3.93)
<i>Average:</i>	<i>75.18</i> <i>(2.04)</i>	<i>75.66</i> <i>(1.37)</i>	<i>76.10</i> <i>(1.67)</i>

167 <sup>a</sup>Bold letters indicate best values for each merging procedure. Merging procedures are compared  
 168 row-wise and normalization methods column-wise.

169

170 To summarize the above results, a comparison was made for each normalization method as to which  
 171 merging procedure lead to the best values for each of the evaluation criteria (Table S6). “Singles”  
 172 merging led to the highest values for four out of five evaluation criteria for all normalization methods  
 173 except for MSTUS. Median normalization achieved best values with the same merging procedure  
 174 (“singles”) for all evaluation criteria. VSN and MSTUS normalization had the highest prediction values  
 175 with polarity merging and for RSD and PCA, “singles” merging led to better results.

176

177 **Table 3: Number of shared significant metabolites for predicting the subjects' sex<sup>a</sup>**

Normalization methods	# of shared significant metabolites		
	By method	By polarity	Singles
Baseline	98	96	95
Creatinine	107	110	116
Sum	103	106	108
MSTUS	103	105	97
Median	102	111	<b>122</b>
PQN	<b>123</b>	<b>116</b>	112
PQN(QC)	107	104	116
Quantile	106	100	105
CSN	105	113	117
VSN	98	97	89
<i>Total:</i>	<i>1137</i>	<i>1145</i>	<i>1164</i>

178 <sup>a</sup>Significant metabolites were identified as described in the Material and Methods section. Bold  
 179 letters indicate best values for the respective merging procedure.

180  
 181 The above results showed that the “singles” merging procedure was optimal for almost all  
 182 normalisation methods. One possible explanation is, that with separate normalization of the datasets,  
 183 the assumptions of the normalization methods are better met. This is supported by the fact that  
 184 methods with similar assumptions benefit from the same merging procedure. The methods sum,  
 185 MSTUS, median, and both PQN variants all share the premise that only a small proportion of the  
 186 compounds (i.e. metabolites) is up and down regulated in equal frequencies.<sup>6,10</sup> These methods all  
 187 perform the best with “singles” merging procedures, except for one of the evaluation criteria (Table  
 188 S5). The MSTUS method seems to favour merging “by-polarity”, however, the performance for this  
 189 particular method may also depend on the number of compounds detected in all samples, which  
 190 differs between merging procedures. Better performance for single merging for CSN and quantile  
 191 normalization may be due to compound intensities being more comparable from sample to sample, if  
 192 datasets are kept separate for normalization.<sup>7</sup>

193 Another influencing factor may be, that compound signal levels could differ between analytical  
 194 methods or ESI modes, due to different ionization efficiencies.<sup>33</sup> Therefore, a single compound, like in  
 195 creatinine normalization, may not be representative enough for the whole dataset, which includes  
 196 peaks generated with differing ionizations or analytical methods. This would be a possible explanation  
 197 for the bad performance of creatinine normalization. Nevertheless, creatinine is still frequently used.<sup>1,5</sup>

198 **Reported gender related metabolites.** In order to evaluate the potential loss of biological  
 199 information by normalization and merging, the number of matches of the significant metabolites  
 200 with the reported sex-specific metabolites was calculated and compared to the number of the  
 201 baseline data (Table 5). A total of eight studies<sup>34-41</sup> were used for the compilation of the sex related  
 202 metabolites and initially 65 compounds were included. A chemical formula was identified for 56 of  
 203 the compounds and 41 compounds were matched by chemical formula in the study datasets. It has  
 204 to be noted that it is not expected for any normalization method to match with all the differential  
 205 metabolites as they are not established biomarkers themselves. The “singles” merging procedure led  
 206 to the highest number of matches for eight out of nine normalization methods. For PQN, VSN, and  
 207 sum normalization, the numbers were equal between singles and “by-polarity” merging procedures.  
 208 For PQN(QC), “singles” and “by-method” merging were the best merging procedures. For the MSTUS  
 209 method, the merging “by-method” was the optimal procedure.

210 **Table 5: Number of matches among the significant metabolites with reported sex metabolites<sup>a</sup>**

Normalization methods	# matches with reported sex metabolites		
	By method	By polarity	Singles
Baseline	7	9	8
Creatinine	5	4	8
Sum	8	9	9
MSTUS	11	10	8
Median	10	8	<b>12</b>
PQN	11	<b>12</b>	<b>12</b>
PQN(QC)	<b>12</b>	9	<b>12</b>
Quantile	10	9	<b>12</b>
CSN	9	10	11
VSN	7	9	9
<i>Total:</i>	<i>96</i>	<i>94</i>	<i>106</i>

211 <sup>a</sup>Bold letters indicate best values for the respective merging procedure.

212  
 213 Comparing the overlap of the significant features with the reported sex related metabolites showed  
 214 that, generally, normalization methods with a high number of significant metabolites also showed a  
 215 high number of matches with sex metabolites, compared to methods with a lower number of  
 216 significant metabolites (Table 3). Creatinine is the exception, showing a high number of significant  
 217 metabolites, which is not reflected by the number of matches for the sex metabolites. This may be  
 218 indicative of possibly more low-quality features being selected among the significant metabolites for  
 219 creatinine normalized data in comparison to the other. PQN and PQN(QC) achieved the highest  
 220 number of matches and showed stable performance across merging procedures. Creatinine and VSN  
 221 normalized data led to a low number of matches across merging procedures. Creatinine normalization  
 222 did not improve baseline data for all merging procedures.

223  
 224 Gender/sex have been used to evaluate urinary sample normalization methods in the past, e.g., for  
 225 sample clustering in PCA,<sup>42</sup> relation of model variance,<sup>1</sup> and as group variable to determine differential  
 226 metabolites.<sup>15</sup> PQN normalization gave the highest matches of significant metabolites with reported  
 227 sex metabolites, however, none of the normalization methods was superior. Similarly, Li and  
 228 colleagues<sup>8</sup> also found comparable performance of the normalization methods PQN, CSN, MSTUS,  
 229 VSN, and quantile normalization in terms of the overlap between experimentally validated biomarkers  
 230 and spiked-in biomarkers with determined statistically significant metabolites. Quantile normalization  
 231 performed slightly lower based on a spike-in dataset, which was not observed in our study. It is  
 232 possible, that intensity differences due to the spiking-in of compounds with different concentrations  
 233 are diminished in quantile normalization, because of the assignment of average values during the  
 234 normalization.

235 In this study, all normalization methods exceeded or at least equalled the baseline data in terms of  
 236 matches with the reported gender metabolites except for creatinine and median normalization. Kohl  
 237 and colleagues<sup>7</sup> evaluated the retention of genuine biological information by relating variation of  
 238 expected constant features with that of varied spiked-in features. In line with this work, PQN  
 239 performed the best and was comparable to the non-normalized data. Quantile, CSN, and VSN  
 240 performed fairly comparable and did not match the non-normalized ones, however, the spiked-in  
 241 signals still clearly stood out. Mervant and colleagues<sup>1</sup> also used gender and assessed biological  
 242 information when comparing normalization methods. They evaluated the explained variance

243 associated with sex after normalization using a modified partial least squares (PLS) method. Contrary  
 244 to the results here, they found post-acquisition creatinine to slightly increase variance related to sex,  
 245 while PQN and MSTUS did not increase variance, compared to the reference data. Possible reasons for  
 246 differing results may include differences in the PQN application, statistical workflow, and that  
 247 explained variance may not be directly correlate with matches for reported sex metabolites.

248 **Comparison of normalization methods.** To determine which normalization methods performed  
 249 optimal considering all evaluation criteria and merging procedures, individual ranks for each evaluation  
 250 criterion were assigned (Table 4). These ranks were equally weighted by the same factor and the  
 251 resulting weighted ranks summed up. The method with the lowest weighted sum is placed on rank 1  
 252 overall.

253  
 254 **Table 4: Ranking of normalization methods based on evaluation criteria for each merging procedure<sup>a</sup>**

Merged by method							
Method	Rank	Weighted Sum	PCA	RSD	RF (BA <sup>1</sup> )	RF (sig mets <sup>2</sup> )	RF (lit matches <sup>3</sup> )
<i>PQN</i>	1	2.80	<b>2</b>	6	<b>2.5</b>	<b>1</b>	<b>2.5</b>
<i>PQN(QC)</i>	2	3.30	4	<b>3</b>	6	<b>2.5</b>	<b>1</b>
<i>Median</i>	3	4.90	<b>3</b>	4	5	8	4.5
<i>MSTUS</i>	4	5.20	5	5	7	6.5	<b>2.5</b>
<i>Quantile</i>	5	6.10	9	10	<b>1</b>	4	4.5
<i>Sum</i>	6	6.30	6	7	4	6.5	7
<i>Creatinine</i>	7	6.70	<b>1</b>	8	10	<b>2.5</b>	10
<i>CSN</i>	8	6.90	10	9	<b>2.5</b>	5	6
<i>VSN</i>	9	7.20	8	<b>1</b>	8	9.5	8.5
<i>Baseline</i>	10	7.40	7	<b>2</b>	9	9.5	8.5
Merged by polarity							
<i>PQN(QC)</i>	1	3.84	<b>1</b>	<b>2</b>	<b>3</b>	7	6.2
<i>PQN</i>	2	4.20	5	7	7	<b>1</b>	<b>1</b>
<i>Median</i>	3	4.40	<b>2</b>	4	4	<b>3</b>	9
<i>MSTUS</i>	4	4.70	4	<b>3</b>	8	6	<b>2.5</b>
<i>Sum</i>	5	5.04	<b>3</b>	6	5	5	6.2
<i>CSN</i>	6	5.50	10	9	<b>2</b>	<b>2</b>	<b>2.5</b>
<i>VSN</i>	7	5.84	7	<b>1</b>	6	9	6.2
<i>Baseline</i>	8	7.24	6	5	9	10	6.2
<i>Quantile</i>	9	7.24	9	10	<b>1</b>	8	6.2
<i>Creatinine</i>	10	8.20	8	8	10	4	10
Singles							
<i>Median</i>	1	1.79	<b>2</b>	4	<b>1</b>	<b>1</b>	<b>2.5</b>
<i>PQN(QC)</i>	2	3.71	<b>1</b>	<b>3</b>	5	<b>3.5</b>	<b>2.5</b>
<i>PQN</i>	3	3.79	<b>3</b>	5	4	5	<b>2.5</b>
<i>Sum</i>	4	5.07	4	<b>2</b>	6	6	6.5
<i>CSN</i>	5	5.43	10	9	<b>3</b>	<b>2</b>	5
<i>Quantile</i>	6	5.21	9	10	<b>2</b>	7	<b>2.5</b>
<i>VSN</i>	7	6.79	8	<b>1</b>	7	10	6.5
<i>MSTUS</i>	8	8.00	5	6	9	8	9
<i>Creatinine</i>	9	8.50	6	8	10	<b>3.5</b>	9
<i>Baseline</i>	10	8.14	7	7	8	9	9

255 <sup>a</sup>BA, balanced accuracy; sig mets, significant metabolites; lit matches, matches with literature  
 256 metabolites for sex. Bold letters indicate best values for each evaluation criterion.



257 Methods with top ranks for QC clustering in PCA include PQN, PQN(QC), median, MSTUS and sum  
258 normalization (Table 4). The QC distance values for these methods were similar across merging  
259 procedures (Table 1). Normalization methods with low ranking performance include VSN, quantile,  
260 CSN. The results for the RSD were very similar to those of the PCA analyses. PQN(QC), PQN, median,  
261 sum and MSTUS showed similar performance and are at the top of the ranking across merging  
262 procedures. Creatinine, CSN and quantile normalized data consistently achieved lower ranks and none  
263 of them achieved a lower RSD than the baseline data. Remarkably, for merging “by-method”, none of  
264 the normalization methods (except VSN) were able to achieve a lower RSD than the value of 5.55 from  
265 the QC-RFSC corrected data. However, VSN normalized data are on a generalised logarithm scale with  
266 base 2 and thus not immediately comparable to the average RSD of the other normalization methods.  
267

268 The balanced accuracy values of the sex prediction were similar across normalization methods, expect  
269 for creatinine normalization, which yielded the lowest values (Table 2). Creatinine normalized data did  
270 not achieve a higher balanced accuracy than the baseline data, which was also true for MSTUS if  
271 merged by “singles” procedure. Quantile and CSN had the highest accuracies across the three merging  
272 procedures. The number of significant metabolites across all normalization methods were also fairly  
273 similar. Only PQN with merging “by-method” exhibited more significant metabolites (123), followed  
274 by creatinine with 107 (Table 3). CSN, creatinine, and PQN normalization led to a high number of  
275 significant metabolites across all merging procedures compared to the other methods. VSN showed  
276 the lowest number of significant metabolites across all merging procedures.  
277

278 PQN and PQN(QC) normalizations were among the best methods, resulting in the lowest weighted sum  
279 of individual ranks for each evaluated criterion across all merging procedures (Table 4). This is in line  
280 with other studies,<sup>11,9,13,6,8,43</sup> which also found PQN to perform best and recommended it as optimum  
281 normalization method. Median normalization despite its simplicity was also among the top methods  
282 across merging procedures. Other studies have also found Median normalization comparable to PQN  
283 in terms of sample clustering in PCA and pooled RSD using test samples,<sup>15</sup> as well as RSD in QC  
284 samples,<sup>43</sup> which is in accordance with findings of this work. The similarity in the performance between  
285 PQN and median normalization may be due to the fact that both methods operate in a similar manner  
286 by relating each sample to a median spectrum.<sup>43</sup> The performance of a normalization method may  
287 depend on how well the assumptions of the respective method are met by the data.<sup>15</sup> Therefore, one  
288 possible explanation for the performance of PQN and median may be, that the current data meets the  
289 assumptions of these two methods more, compared to those of the other methods. Quantile  
290 normalization varied in its overall ranking across merging procedures, but showed particular good  
291 performance for the sample classification, which is consistent with other research.<sup>7</sup> Post-acquisition  
292 creatinine normalization showed low performance in agreement with earlier studies.<sup>16,17,7,1,20,36,5,18,12,42</sup>  
293 In the PCA analysis, creatinine normalized QC samples showed greater dispersion in comparison to  
294 other normalization methods, in line with previous reports.<sup>17,20,12</sup> Sum normalization showed a  
295 mediocre performance in this work, ranging in the middle of the tested normalisation methods. Sum  
296 normalization may be susceptible to compounds with large abundance,<sup>10</sup> which would explain the  
297 spread of QC samples in PCA. Therefore, some authors questioned the usage of sum normalization for  
298 metabolomics data.<sup>6,10</sup>

## 299 **CONCLUSION**

300 The present study shows for the first time that data merging has an effect on normalization  
301 performance and subsequent analysis steps and must be considered when planning the data analysis.  
302 Merging data after normalization was generally favourable for QC similarity, sample classification and  
303 feature selection for most of the tested normalization methods. PQN and Median normalization  
304 showed the best performance overall, considering all tested criteria. Based on this, several

305 recommendations can be provided. Merging data after normalization (“singles”) and the usage of PQN  
306 in a similar setting are generally recommended. PQN is preferred here over median normalisation  
307 because of the more suitable assumptions made about the data. Relying on a single analyte to capture  
308 sample concentration differences, like with post-acquisition creatinine normalization, seems to be a  
309 less preferable approach, especially when data merging is applied. The results of this study may have  
310 broader implications, since other biological matrices like saliva, sweat or faeces also show  
311 heterogeneity in sample concentration or metabolite signals.  
312

### 313 ■ ASSOCIATED CONTENT

314 \*S Supporting Information

315 Additional information as noted in the text. This material is available free of charge via the Internet at  
316 <http://pubs.acs.org>.

### 317 ■ AUTHOR INFORMATION

318 Corresponding Author

319 \*F.B. E-mail: [fbrix@foodtech.uni-kiel.de](mailto:fbrix@foodtech.uni-kiel.de). Tel.: +49 (0)431 880 5365.

### 320 Author Contributions

321 The paper was written through contributions of all authors. All authors have given approval to the final  
322 version of the paper.

### 323 Notes

324 The authors declare no competing financial interest.

### 325 ■ ACKNOWLEDGEMENTS

326 This work was supported by the German Ministry of Education and Research (BMBF) grant 0315540C  
327 and by the Cluster of Excellence “Precision Medicine in Chronic Inflammation” of the German Research  
328 Foundation grant EXC2167. HUZ is supported by the BMBF within the framework of the e:Med research  
329 and funding concept grant 01ZX1912A.

330

### 331 References

- 332 (1) Mervant, L.; Tremblay-Franco, M.; Jamin, E. L.; Kesse-Guyot, E.; Galan, P.; Martin, J.-F.; Guéraud,  
333 F.; Debrauwer, L. *Metabolomics : Official journal of the Metabolomic Society* **2021**, DOI:  
334 10.1007/s11306-020-01758-z.
- 335 (2) Gagnebin, Y.; Tonoli, D.; Lescuyer, P.; Ponte, B.; Seigneux, S. de; Martin, P.-Y.; Schappler, J.;  
336 Boccard, J.; Rudaz, S. *Analytica chimica acta* **2017**, DOI: 10.1016/j.aca.2016.12.029.
- 337 (3) Veselkov, K. A.; Vingara, L. K.; Masson, P.; Robinette, S. L.; Want, E.; Li, J. V.; Barton, R. H.;  
338 Boursier-Neyret, C.; Walther, B.; Ebbels, T. M.; Pelczer, I.; Holmes, E.; Lindon, J. C.; Nicholson, J. K.  
339 *Analytical chemistry* **2011**, DOI: 10.1021/ac201065j.
- 340 (4) Wu, Y.; Li, L. *Journal of chromatography. A* **2016**, DOI: 10.1016/j.chroma.2015.12.007.
- 341 (5) Rosen Vollmar, A. K.; Rattray, N. J. W.; Cai, Y.; Santos-Neto, Á. J.; Deziel, N. C.; Jukic, A. M. Z.;  
342 Johnson, C. H. *Metabolites* **2019**, DOI: 10.3390/metabo9100198.
- 343 (6) Hochrein, J.; Zacharias, H. U.; Taruttis, F.; Samol, C.; Engelmann, J. C.; Spang, R.; Oefner, P. J.;  
344 Gronwald, W. *Journal of proteome research* **2015**, DOI: 10.1021/acs.jproteome.5b00192.
- 345 (7) Kohl, S. M.; Klein, M. S.; Hochrein, J.; Oefner, P. J.; Spang, R.; Gronwald, W. *Metabolomics : Official*  
346 *journal of the Metabolomic Society* **2012**, DOI: 10.1007/s11306-011-0350-z.
- 347 (8) Li, B.; Tang, J.; Yang, Q.; Cui, X.; Li, S.; Chen, S.; Cao, Q.; Xue, W.; Chen, N.; Zhu, F. *Scientific reports*  
348 **2016**, DOI: 10.1038/srep38881.
- 349 (9) Di Guida, R.; Engel, J.; Allwood, J. W.; Weber, R. J. M.; Jones, M. R.; Sommer, U.; Viant, M. R.;  
350 Dunn, W. B. *Metabolomics* **2016**, DOI: 10.1007/s11306-016-1030-9.
- 351 (10) Wulff, J. E.; Mitchell, M. W. *ABB* **2018**, DOI: 10.4236/abb.2018.98022.

352 (11) Chen, J.; Zhang, P.; Lv, M.; Guo, H.; Huang, Y.; Zhang, Z.; Xu, F. *Analytical chemistry* **2017**, DOI:  
353 10.1021/acs.analchem.6b05152.

354 (12) Warrack, B. M.; Hnatyshyn, S.; Ott, K.-H.; Reilly, M. D.; Sanders, M.; Zhang, H.; Drexler, D. M.  
355 *Journal of chromatography. B, Analytical technologies in the biomedical and life sciences* **2009**, DOI:  
356 10.1016/j.jchromb.2009.01.007.

357 (13) Filzmoser, P.; Walczak, B. *Journal of chromatography. A* **2014**, DOI:  
358 10.1016/j.chroma.2014.08.050.

359 (14) Walach, J.; Filzmoser, P.; Hron, K., Chapter Seven - Data Normalization and Scaling. In  
360 *Comprehensive Analytical Chemistry : Data Analysis for Omic Sciences: Methods and Applications*;  
361 Jaumot, Joaquim; Bedia, Carmen; Tauler, Romà, Eds.; Elsevier, 2018.

362 (15) Yu, H.; Huan, T. *Bioinformatics* **2022**, DOI: 10.1093/bioinformatics/btac355.

363 (16) Chetwynd, A. J.; Abdul-Sada, A.; Holt, S. G.; Hill, E. M. *Journal of chromatography. A* **2016**, DOI:  
364 10.1016/j.chroma.2015.12.056.

365 (17) Cook, T.; Ma, Y.; Gamagedara, S. *Journal of pharmaceutical and biomedical analysis* **2020**, DOI:  
366 10.1016/j.jpba.2019.112854.

367 (18) Vogl, F. C.; Mehrl, S.; Heizinger, L.; Schlecht, I.; Zacharias, H. U.; Ellmann, L.; Nürnberger, N.;  
368 Gronwald, W.; Leitzmann, M. F.; Rossert, J.; Eckardt, K.-U.; Dettmer, K.; Oefner, P. J. *Analytical and*  
369 *bioanalytical chemistry* **2016**, DOI: 10.1007/s00216-016-9974-1.

370 (19) Chen, G.; Liao, H.; Tseng, Y. J.; Tsai, I.; Kuo, C. *Analytica chimica acta* **2015**, DOI:  
371 10.1016/j.aca.2015.01.022.

372 (20) Nam, S. L.; La Mata, A. P. d.; Dias, R. P.; Harynuk, J. J. *Metabolites* **2020**, DOI:  
373 10.3390/metabo10090376.

374 (21) Yang, Q.; Bo, L.; Zhu, F., Eds. *Comparison of Standardization Approaches Applied to*  
375 *Metabolomics Data*; ACM: New York, 2017.

376 (22) Geisler, C.; Schlicht, K.; Knappe, C.; Rohmann, N.; Hartmann, K.; Türk, K.; Settgast, U.; Schulte, D.  
377 M.; Demetrowitsch, T.; Jensen-Kroll, J.; Pisarevskaja, A.; Brix, F.; Gruber, B.; Rimbach, G.; Döring, F.;  
378 Rosenstiel, P.; Franke, A.; Schreiber, S.; Henning, C. H. C. A.; Lieb, W.; Nöthlings, U.; Schwarz, K.;  
379 Laudes, M. *European journal of epidemiology* **2022**, DOI: 10.1007/s10654-022-00924-y.

380 (23) Demetrowitsch, T. J.; Petersen, B.; Keppler, J. K.; Koch, A.; Schreiber, S.; Laudes, M.; Schwarz, K.  
381 *Bioanalysis* **2015**, DOI: 10.4155/bio.14.270.

382 (24) Seoudy, A. K.; Schlicht, K.; Kulle, A.; Demetrowitsch, T.; Beckmann, A.; Geisler, C.; Türk, K.;  
383 Rohmann, N.; Hartmann, K.; Brandes, J.; Schulte, D. M.; Schreiber, S.; Schwarz, K.; Holterhus, P.-M.;  
384 Laudes, M. *Neuroendocrinology* **2023**, DOI: 10.1159/000529146.

385 (25) Kind, T.; Fiehn, O. *BMC bioinformatics* **2007**, DOI: 10.1186/1471-2105-8-105.

386 (26) Wishart, D. S.; Guo, A.; Oler, E.; Wang, F.; Anjum, A.; Peters, H.; Dizon, R.; Sayeeda, Z.; Tian, S.;  
387 Lee, B. L.; Berjanskii, M.; Mah, R.; Yamamoto, M.; Jovel, J.; Torres-Calzada, C.; Hiebert-Giesbrecht, M.;  
388 Lui, V. W.; Varshavi, D.; Varshavi, D.; Allen, D.; Arndt, D.; Khetarpal, N.; Sivakumaran, A.; Harford, K.;  
389 Sanford, S.; Yee, K.; Cao, X.; Budinski, Z.; Liigand, J.; Zhang, L.; Zheng, J.; Mandal, R.; Karu, N.;  
390 Dambrova, M.; Schiöth, H. B.; Greiner, R.; Gautam, V. *Nucleic acids research* **2022**, DOI:  
391 10.1093/nar/gkab1062.

392 (27) Luan, H.; Ji, F.; Chen, Y.; Cai, Z. *Analytica chimica acta* **2018**, DOI: 10.1016/j.aca.2018.08.002.

393 (28) Dieterle, F.; Ross, A.; Schlotterbeck, G.; Senn, H. *Analytical chemistry* **2006**, DOI:  
394 10.1021/ac051632c.

395 (29) Bolstad, B. M.; Irizarry, R. A.; Astrand, M.; Speed, T. P. *Bioinformatics* **2003**, DOI:  
396 10.1093/bioinformatics/19.2.185.

397 (30) Lin, S. M.; Du, P.; Huber, W.; Kibbe, W. A. *Nucleic acids research* **2008**, DOI:  
398 10.1093/nar/gkm1075.

399 (31) Workman, C.; Jensen, L. J.; Jarmer, H.; Berka, R.; Gautier, L.; Nielser, H. B.; Saxild, H.-H.; Nielsen,  
400 C.; Brunak, S.; Knudsen, S. *Genome biology* **2002**, DOI: 10.1186/gb-2002-3-9-research0048.

401 (32) Wright, M. N.; Ziegler, A. *J. Stat. Soft.* **2017**, DOI: 10.18637/jss.v077.i01.  
402 (33) Liigand, P.; Kaupmees, K.; Haav, K.; Liigand, J.; Leito, I.; Girod, M.; Antoine, R.; Kruve, A.  
403 *Analytical chemistry* **2017**, DOI: 10.1021/acs.analchem.7b00096.  
404 (34) Kochhar, S.; Jacobs, D. M.; Ramadan, Z.; Berruex, F.; Fuerholz, A.; Fay, L. B. *Analytical*  
405 *biochemistry* **2006**, DOI: 10.1016/j.ab.2006.02.033.  
406 (35) Psihogios, N. G.; Gazi, I. F.; Elisaf, M. S.; Seferiadis, K. I.; Bairaktari, E. T. *NMR in biomedicine*  
407 **2008**, DOI: 10.1002/nbm.1176.  
408 (36) Rasmussen, L. G.; Savorani, F.; Larsen, T. M.; Dragsted, L. O.; Astrup, A.; Engelsen, S. B.  
409 *Metabolomics* **2011**, DOI: 10.1007/s11306-010-0234-7.  
410 (37) Rist, M. J.; Roth, A.; Frommherz, L.; Weinert, C. H.; Krüger, R.; Merz, B.; Bunzel, D.; Mack, C.;  
411 Egert, B.; Bub, A.; Görling, B.; Tzvetkova, P.; Luy, B.; Hoffmann, I.; Kulling, S. E.; Watzl, B. *PloS one*  
412 **2017**, DOI: 10.1371/journal.pone.0183228.  
413 (38) Thévenot, E. A.; Roux, A.; Xu, Y.; Ezan, E.; Junot, C. *Journal of proteome research* **2015**, DOI:  
414 10.1021/acs.jproteome.5b00354.  
415 (39) Trimigno, A.; Khakimov, B.; Savorani, F.; Tenori, L.; Hendrixson, V.; Čivilis, A.; Glibetic, M.;  
416 Gurinovic, M.; Pentikäinen, S.; Sallinen, J.; Garduno Diaz, S.; Pasqui, F.; Khokhar, S.; Luchinat, C.;  
417 Bordoni, A.; Capozzi, F.; Balling Engelsen, S. *Molecular nutrition & food research* **2019**, DOI:  
418 10.1002/mnfr.201800216.  
419 (40) Wang, T.; Tang, L.; Lin, R.; He, D.; Wu, Y.; Zhang, Y.; Yang, P.; He, J. *Molecular genetics & genomic*  
420 *medicine* **2021**, DOI: 10.1002/mgg3.1738.  
421 (41) Zhang, S.; Liu, L.; Steffen, D.; Ye, T.; Raftery, D. *Metabolomics* **2012**, DOI: 10.1007/s11306-011-  
422 0315-2.  
423 (42) Yan, Z.; Yan, R. *Journal of chromatography. A* **2016**, DOI: 10.1016/j.chroma.2016.03.023.  
424 (43) Noonan, M. J.; Tinnesand, H. V.; Buesching, C. D. *BioEssays : news and reviews in molecular,*  
425 *cellular and developmental biology* **2018**, DOI: 10.1002/bies.201700210.  
426