

1
2
3
4
5
6
7

A multi-dataset evaluation of frame censoring for task-based fMRI

8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Michael S. Jones, Zhenchen Zhu*, Aahana Bajracharya, Austin Luor, Jonathan E. Peelle

Department of Otolaryngology, Washington University in St. Louis, St. Louis MO USA

* Current affiliation: MD Program, Chinese Academy of Medical Sciences and Peking Union
Medical College, Beijing China

27
28
29
30
31
32
33

Running title: Motion correction in task-based fMRI

34
35
36
37
38

Keywords: motion correction, head movement, frame censoring, scrubbing, FD, DVARS, task-based fMRI

39
40
41
42
43
44

Please address correspondence to:

Dr. Michael Jones
Department of Otolaryngology
Washington University in St. Louis
660 South Euclid, Box 8115
St. Louis, MO 63110
email: jones.mike@wustl.edu

Dr. Jonathan Peelle
Department of Otolaryngology
Washington University in St. Louis
660 South Euclid, Box 8115
St. Louis, MO 63110
email: jpeelle@wustl.edu

45

Abstract

46 Subject motion during fMRI can affect our ability to accurately measure signals of interest. In
47 recent years, frame censoring—that is, statistically excluding motion-contaminated data within
48 the general linear model using nuisance regressors—has appeared in several task-based fMRI
49 studies as a mitigation strategy. However, there have been few systematic investigations
50 quantifying its efficacy. In the present study, we compared the performance of frame censoring
51 to several other common motion correction approaches for task-based fMRI using open data
52 and reproducible workflows. We analyzed eight datasets available on OpenNeuro.org
53 representing eleven distinct tasks in child, adolescent, and adult participants. Performance was
54 quantified using maximum t-values in group analyses, and ROI-based mean activation and split-
55 half reliability in single subjects. We compared frame censoring to the use of 6 and 24 canonical
56 motion regressors, wavelet despiking, robust weighted least squares, and untrained ICA-based
57 denoising. Thresholds used to identify censored frames were based on both motion estimates
58 (FD) and image intensity changes (DVARS). Relative to standard motion regressors, we found
59 consistent improvements for modest amounts of frame censoring (e.g., 1–2% data loss),
60 although these gains were frequently comparable to what could be achieved using other
61 techniques. Importantly, no single approach consistently outperformed the others across all
62 datasets and tasks. These findings suggest that although frame censoring can improve results,
63 the choice of a motion mitigation strategy depends on the dataset and the outcome metric of
64 interest.
65

66

Introduction

67 Obtaining high-quality neuroimaging data depends on minimizing artifacts. Although
68 advancements in hardware and pulse sequence design have reduced many types of noise
69 inherent to functional MRI, other sources remain (Bianciardi et al., 2009). One prominent
70 challenge is artifacts caused by subject head motion. Among other effects, head motion
71 changes the part of the brain sampled by a particular voxel and can introduce changes in signal
72 intensity through interactions with the magnetic field, which add noise to the data and make it
73 harder to identify signals of interest.

74 The effects of head motion have received recent scrutiny in the context of resting state
75 functional connectivity. Because motion-related artifacts occur in many voxels simultaneously,
76 they can introduce correlations in fMRI time series that are unrelated to BOLD activity, leading
77 to inaccurate estimates of functional connectivity (Power et al., 2015; Satterthwaite et al., 2019).
78 However, spurious activation is also of concern in task-based functional neuroimaging. Rigid
79 body realignment—a mainstay of fMRI analysis for decades—goes some way towards
80 improving correspondence across images (Ashburner and Friston, 2004), but does not remove
81 extraneous signal components introduced by movement (Friston et al., 1996). A common
82 approach for mitigating motion-related artifacts is to include the 6 realignment parameters
83 (translation and rotation around the X, Y, and Z axes) as nuisance regressors in first-level
84 models.

85 Alternatively, several data-driven strategies have been developed to reduce the
86 influence of high-motion scans on estimated activations. Wavelet decomposition identifies
87 artifacts by exploiting their non-stationarity across different temporal scales (Patel et al., 2014).
88 The method has been applied in resting state studies but is also applicable to task-based data.
89 Independent component analysis (Pruim et al., 2015) identifies artifacts based on the spatial
90 distribution of shared variance. In robust weighted least squares (Diedrichsen and Shadmehr,
91 2005), a two-pass modeling procedure is used to produce a collection of nuisance regressors
92 which are then included in the final analysis to weight frames by the inverse of their variance
93 (that is, downweighting frames with high error).

94 An alternative motion correction strategy is “scrubbing” or “frame censoring” (Lemieux et
95 al., 2007; Siegel et al., 2014). In this approach, bad scans are identified and excluded from
96 statistical analysis. One approach is to do so by modeling them in the general linear model
97 using nuisance regressors (i.e. “scan-nulling regressors” or “one-hot encoding”). Although frame
98 censoring has received considerable interest in resting state fMRI over the past several years
99 (Power et al., 2012; Gratton et al., 2020a), it has not seen widespread use in the task-based
100 fMRI literature. Censoring approaches involve some effective data loss, in that censored frames
101 do not contribute to the task-related parameter estimates, and that columns introduced to the
102 design matrix to perform censoring reduce the available degrees of freedom. Choosing an
103 appropriate metric and associated threshold for identifying bad scans can also be challenging.
104 Thus, additional information over what threshold should be used for identifying bad frames—and
105 relatedly, how much data is lost vs. retained—is necessary to make informed decisions.

106 Although several published studies comparing differing correction strategies exist
107 (Ardekani et al., 2001; Oakes et al., 2005; Johnstone et al., 2006), a drawback of prior work is
108 that evaluation was often limited to a single dataset (see **Supplemental Table 1**). The degree to
109 which an optimal strategy for one dataset generalizes to other acquisition schemes, tasks, or
110 populations is not clear. With the increased public availability of neuroimaging datasets
111 (Poldrack et al., 2013; Markiewicz et al., 2021), the possibility of evaluating motion correction
112 approaches across a range of data has become more feasible.

113 In the present work, we sought to compare the performance of identical pipelines on a
114 diverse selection of tasks, using data from different sites, scanners, and subject pools.

115 Although our primary interest was frame censoring, we considered seven different motion-
116 correction approaches:

- 117 1. six canonical head motion estimates (RP6)
- 118 2. 24-term expansions of head motion estimates (RP24)
- 119 3. wavelet despiking (WDS)
- 120 4. robust weighted least squares (rWLS)
- 121 5. untrained independent component analysis (ICA)
- 122 6. frame censoring based on frame displacement (FD)
- 123 7. frame censoring based on variance differentiation (DVARs)

124 This list is not exhaustive but representative of approaches that are currently used and feasible
125 to include in an automated processing pipeline.

126 Because it is impossible to determine a “ground truth” result with which to compare the
127 effectiveness of these approaches, we instead considered three complementary outcome
128 metrics: 1) the maximum group t-statistic both across the whole-brain and in a region-of-interest
129 relevant to the task; 2) the average parameter estimates from within the same ROI (that is,
130 effect size); and 3) the degree of test-retest consistency exhibited by subject-level parametric
131 maps. These metrics are simple to define yet functionally meaningful, and can be applied to
132 data from almost any fMRI study.

133 **Methods**

134 **Datasets**

135 We analyzed eight studies obtained from OpenNeuro (Markiewicz et al., 2021), several of which
136 included multiple tasks or multiple participant groups. As such, the eight selected studies
137 provided a total of 15 datasets. The selection process was informal, but studies given priority
138 included 1) a clearly-defined task, 2) a sufficient number of subjects to allow second-level
139 modeling, 3) sufficient data to make test-retest evaluation possible, and 4) a publication
140 associated with the data describing a result to which we could compare our own analysis.

141 A summary of the eight datasets selected is shown in **Table 1** (acquisition details
142 provided in **Supplemental Table 2**). Additional information, including task details,
143 modeling/contrast descriptions compiled from publication(s) associated with a given study, and
144 any data irregularities encountered during analysis, is provided in the **Supplemental Materials**.

145 **Analysis**

146 Analysis was performed using Automatic Analysis version 5.4.0 (Cusack et al., 2015) (RRID:
147 SCR_003560), which scripted a combination of SPM12 (Wellcome Trust Centre for
148 Neuroimaging) version 7487 (RRID: SCR_007037) and FMRIB Software Library (FSL; FMRIB
149 Analysis Group; (Jenkinson et al., 2012) version 6.0.1 (RRID: SCR_002823). BrainWavelet
150 Toolbox v2.0 (Patel et al., 2014) was used for wavelet despiking, and rWLS version 4.0
151 (Diedrichsen and Shadmehr, 2005) for robust weighted least squares. Analysis scripts used in
152 the study are available at <https://osf.io/n5v3w/>.

153 To the extent possible, we used the same preprocessing pipeline for all datasets (**Figure**
154 **1a**). Briefly, structural and functional images were translated to the center of the scanned
155 volume and the first four frames of each session were removed in functional images to allow for
156 signal stabilization. This was followed by bias correction of the structural image, realignment,
157 coregistration of the functional and structural images, normalization into MNI space using a
158 unified segmentation approach (Ashburner and Friston, 2005) resampled at 2 mm, and
159 smoothing of the functional images using an 8 mm FWHM Gaussian kernel.

160

Table 1. Summary of datasets analyzed

Dataset	Reference	Task	Age group*	# subs	FD (median \pm SD)	frames per subject
ds000102	Kelly et al. (2008)	flanker	YA	22	0.11 \pm 0.12	284
ds000107	Duncan et al. (2009)	1-back	YA	43	0.08 \pm 0.14	323
ds000114	Gorgolewski et al. (2013a)	motor (lips)	YA	10	0.14 \pm 0.16	360
		covert verb	YA	10	0.11 \pm 0.11	338
		overt word	YA	10	0.13 \pm 0.12	144
		line bisection	YA	9	0.13 \pm 0.18	468
ds000228	Richardson et al. (2018)	movie viewing	C	122	0.21 \pm 0.93	164
			YA	33	0.18 \pm 0.27	164
ds001497	Lewis-Peacock and Postle (2008)	face perception	YA	10	0.11 \pm 0.12	1146
ds001534	Courtney et al. (2018)	food images	YA	42	0.10 \pm 0.16	552
ds001748	Fynes-Clinton et al. (2019)	memory retrieval	C	21	0.16 \pm 0.36	438
			T	20	0.12 \pm 0.17	438
			YA	21	0.08 \pm 0.17	438
ds002382	Rogers et al. (2020)	word recognition	YA	29	0.14 \pm 0.35	710
			OA	32	0.30 \pm 0.34	710

161 Note: * OA = older adults; YA = young adults; T = teens; C = children

162
163 Functional images were corrected for motion artifacts using each of the following
164 approaches: 1) inclusion of six canonical motion estimates in the first-level model as nuisance
165 regressors, 2) inclusion of 24 nuisance regressors based on a second-order expansion of the
166 motion estimates and first derivatives, 3) wavelet despiking, 4) robust weighted least squares, 5)
167 ICA denoising, 6) frame censoring based on framewise displacement (FD) or 7) differential
168 variance (DVARS) thresholding (FD/DVARS thresholding is described below).

169 Statistical modeling was performed in SPM in all motion correction approaches. First-
170 level modeling included a contrast of interest described in a publication associated with the
171 dataset for evaluation, followed by second-level analysis to produce group-level statistical maps.
172 All first- and second-level t-maps were thresholded at a voxelwise threshold of $p < 0.001$
173 (uncorrected).

174 Minor pipeline modifications were required for robust weighted least squares, wavelet
175 despiking, and ICA denoising. As recommended by developers of the rWLS toolbox,
176 unsmoothed data was used for variance estimation and contrast maps were smoothed after
177 modeling. For wavelet despiking, functional images were rescaled to a whole-brain median of
178 1000 across all frames prior to processing. The default toolbox settings (wavelet: d4, threshold:
179 10, boundary: reflection, chain search: moderate, scale number: liberal) were used. Finally, ICA-

180 based denoising was implemented using ICA-AROMA (Pruim et al., 2015) with additional
 181 processing steps performed within FSL. Briefly, the unsmoothed coregistered functional image
 182 was demeaned, detrended, smoothed, and then nonlinearly warped to the FSL 2 mm MNI152
 183 template using FNIRT. The normalized functional image was then passed to AROMA for
 184 denoising.
 185

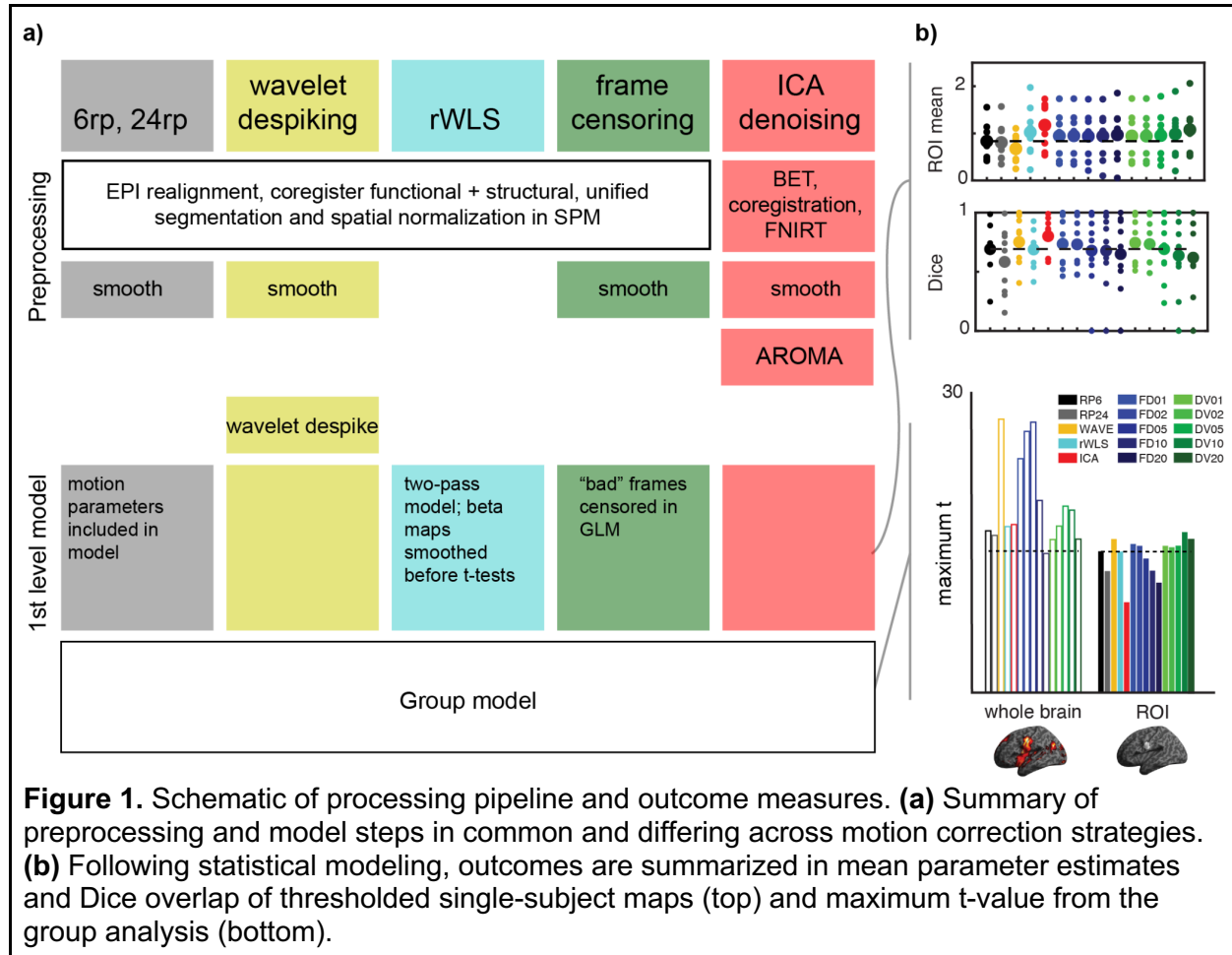


Figure 1. Schematic of processing pipeline and outcome measures. **(a)** Summary of preprocessing and model steps in common and differing across motion correction strategies. **(b)** Following statistical modeling, outcomes are summarized in mean parameter estimates and Dice overlap of thresholded single-subject maps (top) and maximum t-value from the group analysis (bottom).

186 Evaluation of Motion Correction Performance

187 Three measures were used to quantify the performance of each motion correction strategy,
 188 illustrated in **Figure 1b**: 1) maximum t-value, 2) effect size, and 3) subject replicability. In the
 189 first measure, the maximum t-value occurring in the group level parametric map was extracted
 190 both at the whole-brain level and also within a region-of-interest relevant to the task. The effect
 191 size was quantified as the mean of all voxels within the ROI for each subject using the first-level
 192 beta maps. To evaluate subject replicability, session data were treated as a test-retest paradigm
 193 (the first session versus the second session in studies having fewer than three sessions; even-
 194 numbered versus odd-numbered sessions otherwise). Replicability was quantified as the Dice
 195 coefficient of thresholded first-level t-maps (0.001, uncorrected) in each subject (restricted to the
 196 ROI).

197 FD and DVARS Thresholding

198

199 Motion correction approaches based on frame censoring required quantification of motion
200 artifacts which could then be subjected to thresholding. Both framewise displacement (FD) and
201 differential variance (DVARS) were used. Framewise displacement was calculated as the sum
202 of the six head motion estimates obtained from realignment, with a dimensional conversion of
203 the three rotations assuming the head is a 50 mm sphere (Power et al., 2012). DVARS was
204 calculated as the root-mean-squared of the time difference in the BOLD signal calculated across
205 the entire brain (Smyser et al., 2011). As shown in **Figure 2A**, both metrics closely tracked
206 artifacts apparent in voxel intensities and also each other. Although FD and DVARS in a given
207 session tended to be correlated (**Figure 2B**), they were not identical and could exhibit slightly
208 different time courses and relative peak amplitudes. As such, we explored the use of both
209 measures.

210 Thresholds were determined by calculating FD and DVARS across all sessions in all
211 subjects, which allowed values to be identified that resulted in 1%, 2%, 5%, 10%, and 20%
212 frame violations across the entire dataset (**Figure 2C**). We adopted this strategy rather than
213 using a fixed value of FD or DVARS for several reasons. First, FD and DVARS magnitudes
214 change with the TR of the data, because the TR is the sampling rate (for a given movement,
215 sampling more rapidly will give smaller FD values, even though the total motion is the same).
216 Secondly, different calculations of FD provide different values (Jenkinson et al., 2002; Power et
217 al., 2012; Van Dijk et al., 2012), and thus any absolute threshold would necessarily be metric-
218 specific. Finally, datasets differ in their tasks and populations, and we anticipated that a
219 constant threshold would not be suitable for all datasets. We, therefore, employed the frame-
220 percent thresholding strategy in order to obtain an informative range of results in all studies
221 examined. Because the threshold is chosen to limit data loss in the whole group, it allows high-
222 motion subjects to have more frames censored than low-motion subjects, which was one of our
223 primary goals.

224 The threshold values that resulted from percent data loss targeting in these datasets are
225 shown in **Supplemental Figure 1**.

226

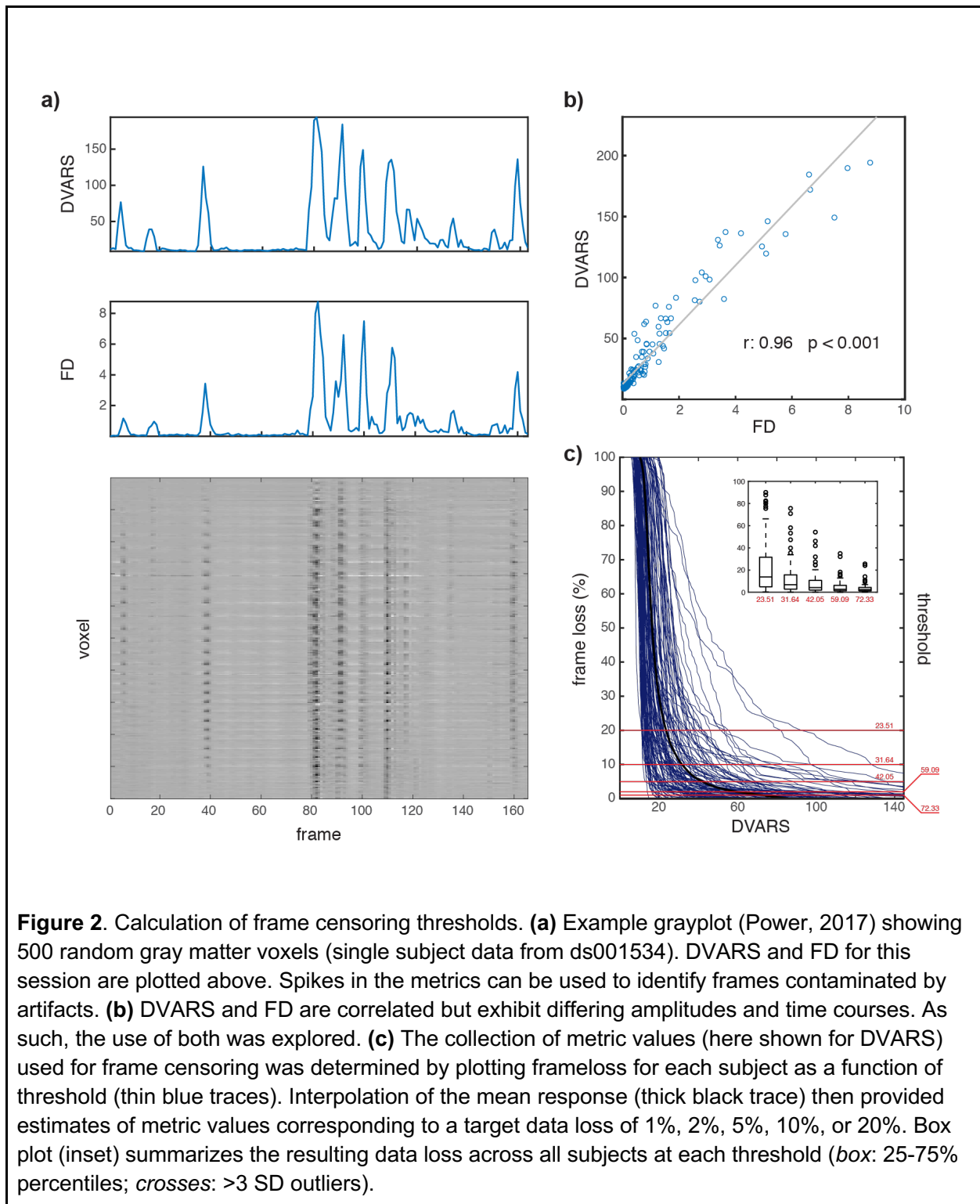


Figure 2. Calculation of frame censoring thresholds. **(a)** Example grayplot (Power, 2017) showing 500 random gray matter voxels (single subject data from ds001534). DVARS and FD for this session are plotted above. Spikes in the metrics can be used to identify frames contaminated by artifacts. **(b)** DVARS and FD are correlated but exhibit differing amplitudes and time courses. As such, the use of both was explored. **(c)** The collection of metric values (here shown for DVARS) used for frame censoring was determined by plotting frame loss for each subject as a function of threshold (thin blue traces). Interpolation of the mean response (thick black trace) then provided estimates of metric values corresponding to a target data loss of 1%, 2%, 5%, 10%, or 20%. Box plot (inset) summarizes the resulting data loss across all subjects at each threshold (*box*: 25-75% percentiles; *crosses*: >3 SD outliers).

227
228
229
230

To impose frame censoring, first-level modeling was repeated for each threshold with a delta function (i.e. a scan-nulling regressor) included in the design matrix at the location of each violation, which effectively removes the contribution of the targeted frame from the analysis.

231 **Region of Interest Definition**

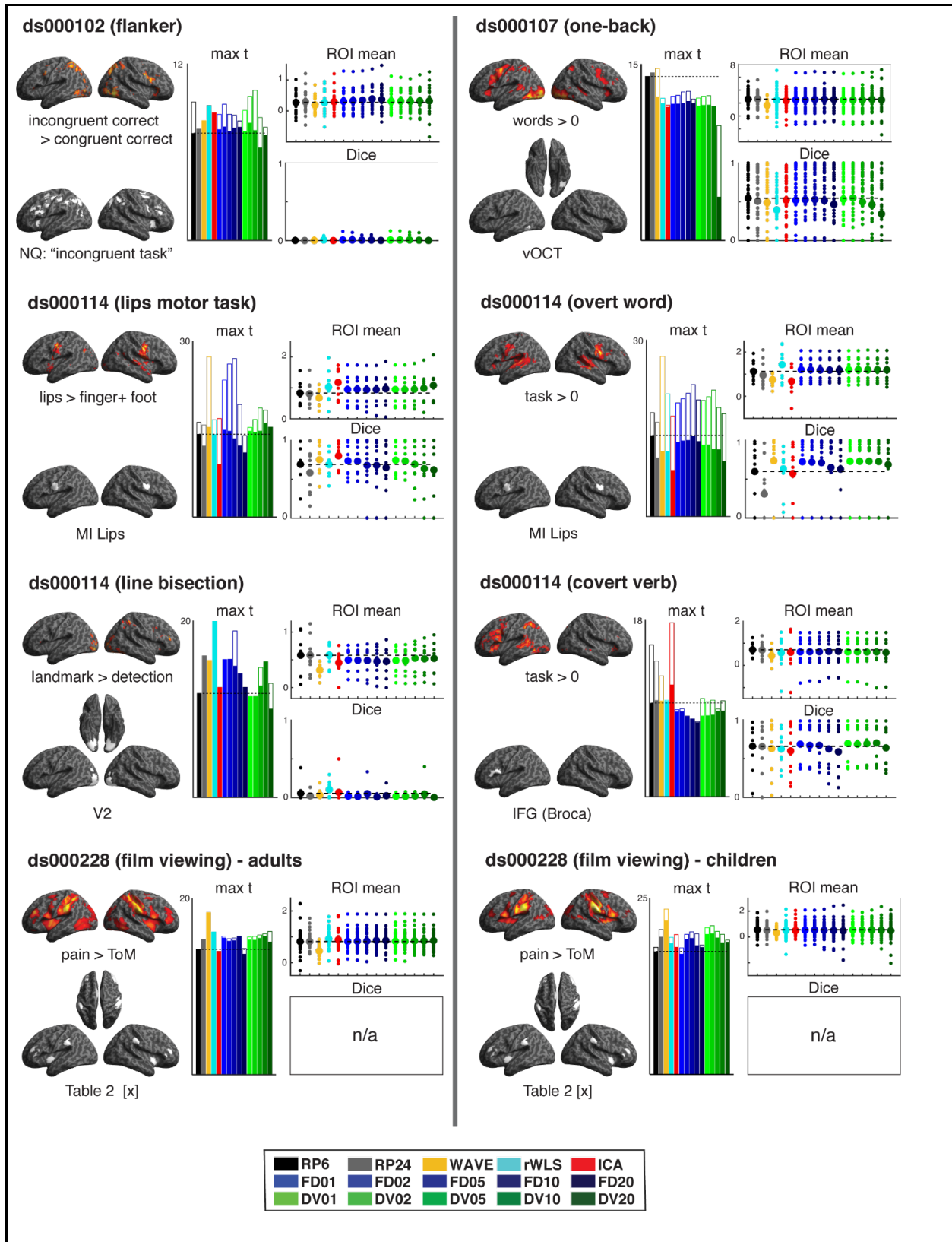
232 A task-relevant ROI for each study/task was defined in one of three ways: 1) a 5-mm sphere (or
233 spheres) centered at coordinates reported in a publication associated with the dataset, 2) a
234 whole-brain Z-mask generated by a task-relevant search term (e.g., "incongruent task") in
235 NeuroQuery (Dockès et al., 2020) and thresholded $z > 3$, or 3) a binarized tissue probability
236 map in the SPM Anatomy Toolbox (Eickhoff et al., 2005) for a task-relevant brain structure or
237 anatomical region (e.g., "V2").

238 **Results**

239 Performance of the motion correction strategies organized by dataset is shown in **Figure 3**.
240 Each panel includes a representative second-level thresholded t-map at the upper left ($p <$
241 0.001, uncorrected) using the "RP6" approach (six canonical motion parameters included as
242 nuisance regressors). A contrast descriptor is given below the map. The ROI used for
243 evaluation is shown at lower left with the source listed under the rendered image; "NQ"
244 indicates search term from NeuroQuery (Dockès et al., 2020); all other labels indicate either an
245 Anatomy Toolbox tissue probability map (Eickhoff et al., 2005) or a 5 mm sphere. Additional
246 details on ROI definition used in each analysis are provided in the **Supplemental Materials**.

247 These results show there is substantial variability in motion correction approaches, with
248 performance depending both on the data under consideration and the chosen performance
249 metric. However, some general trends are apparent. Wavelet despiking tended to offer the best
250 maximum t-value in both the whole-brain and ROI-constrained evaluation, with robust weighted
251 least squares also exhibiting good performance (note the ROI-restricted maximum t-value,
252 shown in filled bars, are superimposed on the whole-brain results, shown in open bars in Figure
253 3 due to space restrictions). Conversely, ICA gave consistently poorer results although it offered
254 the best maximum t-value in the ds000114 covert verb task. Performance of FD and DVARS
255 frame censoring were highly variable, with the application of increasingly stringent thresholds
256 improving performance in some datasets while decreasing it in others. A somewhat consistent
257 behavior is a loss of performance at the highest (20%) FD or DVARS threshold. As a rule, frame
258 censoring performed better than RP6 and RP24 motion correction, although RP6 is competitive
259 (if not optimal) in both ds000107 and ds001748.

260
261



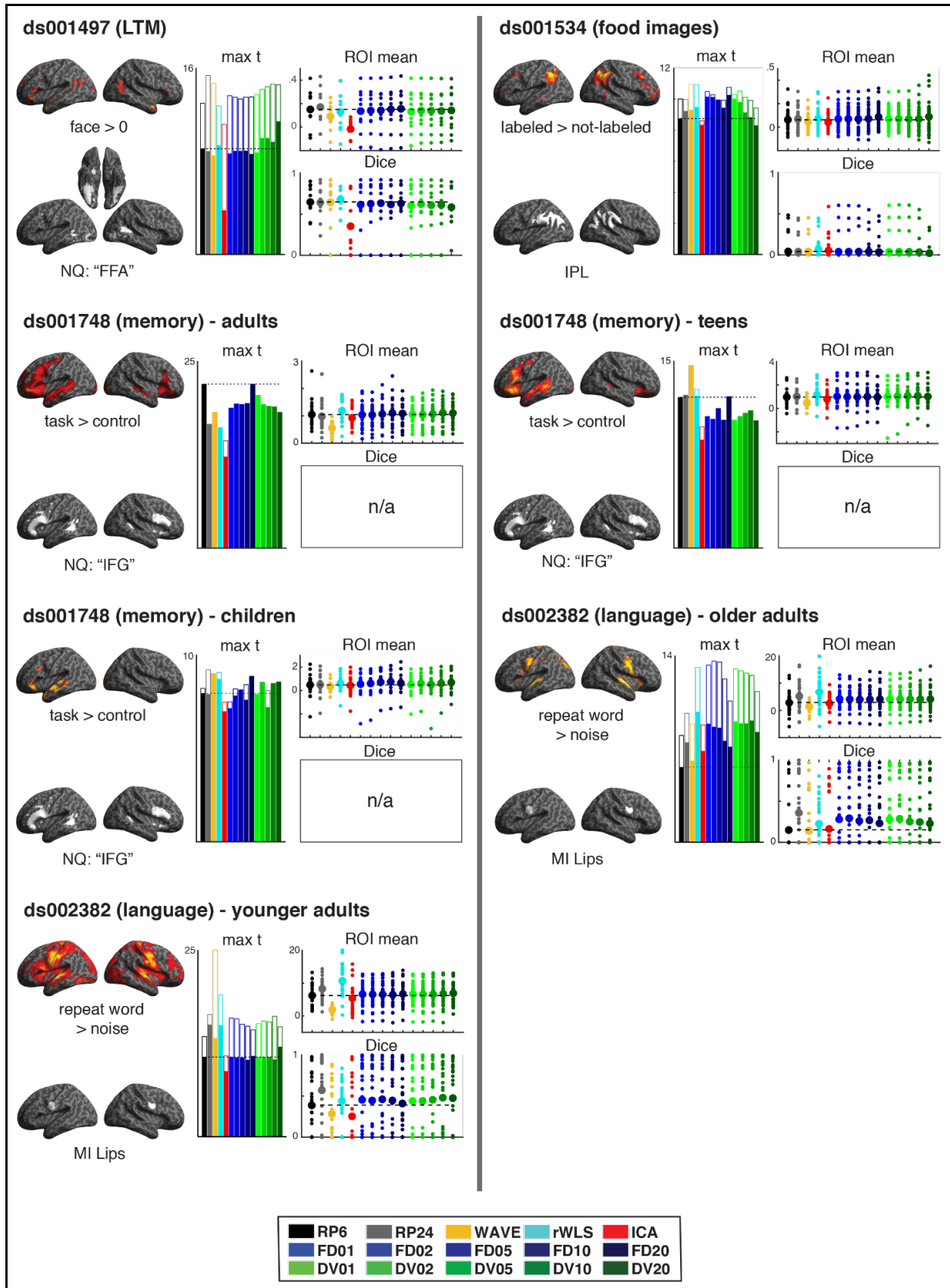


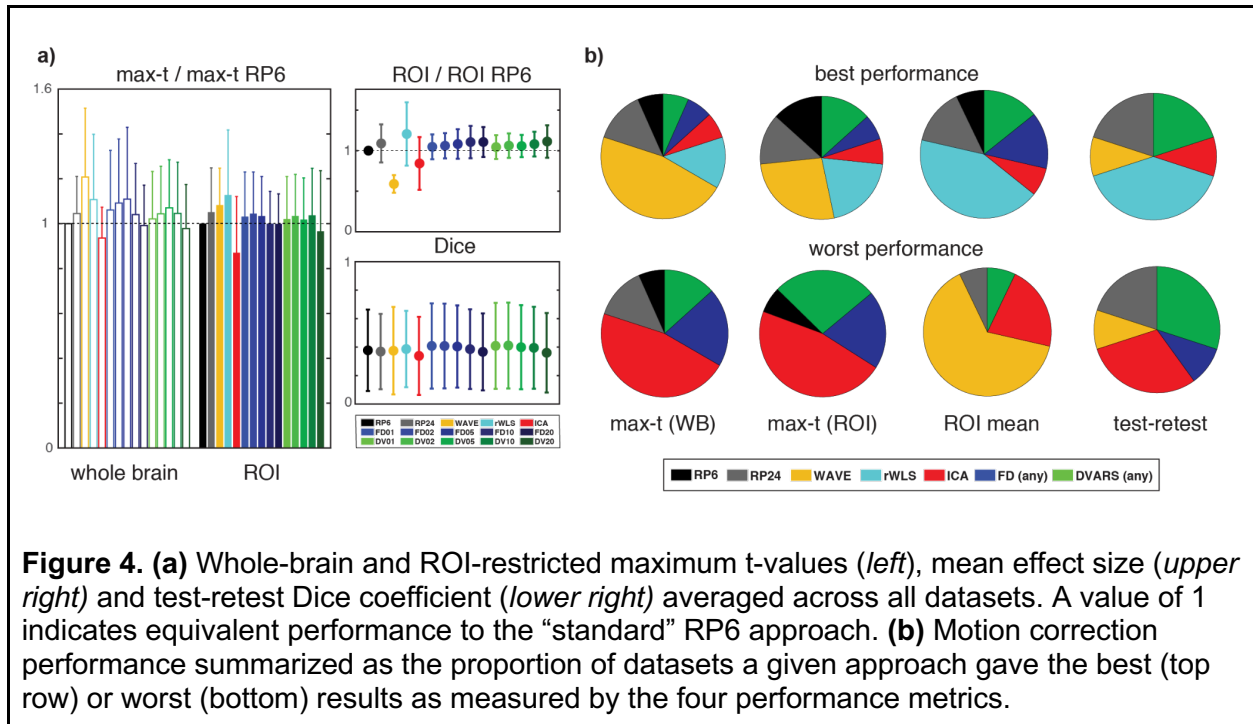
Figure 3. Summary of denoising algorithm performance for all datasets examined in the study. Each panel includes a representative thresholded group t-map at left ($p=0.001$, uncorrected) for the given contrast with the ROI used for evaluation plotted below. At center, ROI-restricted maximum t-values are superimposed on whole-brain results for each denoising approach. Plots at right show individual-subject mean ROI effect size (*top*) and Dice coefficient for a split-half test-retest evaluation (*bottom*). Datasets that did not permit test-retest evaluation are noted "n/a." The horizontal dotted line in a given plot indicates RP6 results for reference.

264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292

The mean effect size shown in these results is largely insensitive to the selected motion correction approach. The two exceptions are wavelet despiking and ICA, which produced consistently smaller values than the other approaches. This may reflect suboptimal parameter selection in these algorithms (see Discussion). Robust weighted least squares offered competitive results in all datasets and notably superior results in ds002382 and the ds000114 overt word task. FD and DVARS frame censoring neither improved nor degraded results regardless of threshold, producing a mean effect size indistinguishable from both the RP6 and RP24 approaches save for a few isolated individual subjects.

The test-retest results also demonstrate a great deal of variability. The Dice coefficients exhibit substantial inter-subject differences, resulting in a mean performance that is similar across all motion correction strategies. However, excluding ds000102, ds001534, and the ds000114 line bisection task which provide an uninformative test-retest quantification, some trends can be identified. There is a detectable decrease in both the FD and DVARS frame censoring results, especially at 20% thresholding. In general, all differences are minor, save for ICA which performs notably better in the ds000114 motor task and notably worse in ds001487.

A summary of these results is shown in **Figure 4a**, in which average values of the four performance metrics are plotted for all 15 datasets/tasks. Several of the trends noted above remain apparent. Wavelet despiking gives the largest whole-brain maximum t-value. Robust weighted least squares resulted in the best ROI-constrained performance. Light-to-moderate frame censoring results in improvement which then declines as more aggressive thresholding is applied. Robust weighted least squares produces the largest average effect size. Wavelet despiking and ICA produce poor results as measured by this metric. Finally, the averaged Dice coefficient is less than 0.5 in all datasets. A decline of FD and DVARS frame censoring performance with increasing threshold is apparent. However, all of the test-retest results exhibit substantial variability (error bars denote +1 SD in the maximum t-value plot; +/- 1 SD in ROI mean effect size and Dice).



293
294
295
296
297
298
299
300
301
302
303
304
305
306

An alternate summary of algorithm performance is presented in **Figure 4b**, in which the best and worst performer measured by each metric was identified in each of the 15 datasets and the resulting proportions plotted as pie charts. The trends and variability evident in the grand averages are also apparent in these data. Robust weighted least squares offered best performance on many datasets and worst performance on none. Wavelet despiking gave the best maximum t-value in about half (whole-brain) or one quarter (ROI-restricted) of the studies, but the worst ROI mean effect size in over half. ICA denoising was often the worst performer, yet gave best results across all four metrics in at least one dataset. Frame censoring performed roughly equally well (or equally poorly) using either FD or DVARS, with one notable result being FD was never the worst performer for effect size. Finally, the performance of the RP6 and RP24 approaches are middling, producing best or worse maximum t-value on only one or two datasets and, with one exception, never producing best nor worst ROI mean or test-retest results.

307

Discussion

308
309
310
311
312
313
314
315
316

We explored the performance of a variety of approaches to correcting motion-related artifacts in task-based fMRI. The studies examined represent a broad range of task domains, including sensory, motor, language, memory, and other cognitive functions, with participants varying in age, sex, and other characteristics. Although we set out expecting to find converging evidence for an optimal strategy, instead our results demonstrate that the performance of motion correction approaches depends on both the data and the outcome of interest. We review our selected metrics below—whole-brain and ROI-restricted maximum t-value, mean effect size, and test-retest repeatability—followed by some general comments on each motion correction approach.

317 **Comparing outcome metrics**

318 The use of whole-brain maximum t-value measured in group-level statistical maps has the
319 advantage that it requires few assumptions about the data or the expected pattern of activity.
320 However, we did not observe a consistent pattern regarding which motion correction approach
321 optimized the whole-brain maximum t-value. Disparity was even evident between different
322 participant groups within a given study. For example, wavelet despiking had the highest whole-
323 brain t statistic in ds0001748 in teens but RP6 offered better performance in adults.

324 In addition to whole-brain statistics, we examined maximum t-values within a selected
325 region of interest. Our rationale for doing so was that researchers interested in task-based
326 effects frequently have prior intuitions about where the most informative results are localized. A
327 potential downside of this approach is the need to specify an ROI to examine. We found that
328 motion correction approaches can exhibit substantially different whole-brain and ROI-restricted
329 performance. In the ds000114 overt word task, for example, RP6 offered best performance
330 within the motor cortex but poor performance in a whole-brain evaluation. Furthermore, frame
331 censoring performance improved in some datasets but degraded in others as more stringent
332 thresholding was applied. Obviously, a challenge inherent in such an evaluation is the actual
333 ROI selection. Although we believe our choices are sensible, selection of a different ROI set
334 may well result in a different overall view of performance.

335 To complement these group-level measures, we also considered two single-subject
336 metrics: mean effect size and test-retest repeatability measured by Dice overlap in thresholded t
337 maps. Effect size permits an examination of parameter estimates, and our use of averaging
338 offers a direct and simple quantification. However, with the exceptions of wavelet despiking and
339 aggressive frame censoring (revisited below), we observed that effect size was largely
340 insensitive to the choice of motion correction strategy, although less than the variability
341 observed in maximum t-value. This suggests the main effect of different motion correction
342 approaches is a differential reduction in model error variance. If parameter estimation is the
343 primary result of interest, then the choice of motion correction strategy may not be critical.

344 The test-retest evaluation was perhaps the least helpful result, with the performance of
345 all motion correction approaches essentially indistinguishable under this metric. Although the
346 outcome is disappointing, it should be noted that many of the studies included here were not
347 designed to include a split-half repeatability analysis. It may be that more data per subject may
348 be needed for this metric to be informative. In that sense, our analyses speak to the general
349 challenges of obtaining reliable single-subject data in fMRI (Smith et al., 2005; Bennett and
350 Miller, 2010; Gorgolewski et al., 2013b; Elliott et al., 2020), at least under conventional scanning
351 protocols (Gratton et al., 2020b). Investigators of resting state fMRI have confronted a similar
352 issue, and recommendations have appeared in the resting state fMRI literature outlining the
353 minimal scan time required for reproducible results (Birn et al., 2013; Laumann et al., 2015).
354 Perhaps an analogous standard might be possible for task-based fMRI, although any guideline
355 would necessarily require the cognitive complexity of the task under investigation to be
356 considered.

357 **Comparing Motion Correction Approaches**

358 No single denoising approach exhibited optimal performance on all datasets and all metrics.
359 Algorithm performance did not appear to be systematically related to the nature of the task,
360 acquisition parameters, nor any feature of the data that could be identified.

361 Interestingly, computationally-intensive approaches did not necessarily perform better
362 than basic corrective measures. For some datasets, including six motion estimates as
363 continuous nuisance regressors—a standard approach used in functional imaging for
364 decades—could perform as well or better than more sophisticated algorithms that have
365 emerged in recent years. Increasing the head motion estimate from a 6- to a 24-parameter

366 expansion led to an improvement in some data but degraded results in others. Although such
367 results are rather counterintuitive, we can provide a few observations, even if these data do not
368 currently permit conclusive recommendations.

369 Two motion correction approaches that showed generally strong performance were
370 wavelet despiking (WDS) and robust weighted least squares (rWLS). Together, these
371 approaches offered best performance in approximately half of the datasets across all
372 performance metrics (**Figure 4**). Additionally, in no evaluation did rWLS produce the worst
373 results. In a statistical sense, robust weighted least squares might be seen as an optimal
374 solution, in that it uses the error in the model to re-weight time points, reducing the influence of
375 motion on parameter estimates. However, we also found that other motion correction strategies
376 supplied similar, or superior, performance in several instances. One reason might be that rWLS
377 linearly weights time points inversely related to their variance. To the degree that motion
378 artifacts include a nonlinear component, linear weighting may not adequately (or at least, not
379 optimally) remove all of the artifact.

380 In contrast to good performance of wavelet despiking as measured by maximum t-value,
381 it gave notably low scores on mean effect size. However, this finding may simply reflect data
382 scaling specific to the WDS implementation. It should also be noted the WDS toolbox offers 20
383 wavelets, and additional options that control algorithm behavior such as thresholding and chain
384 search selection. The results obtained here are what can be expected using the default settings
385 recommended by the toolbox developers, which includes a median 1000 rescaling of the
386 functional data (and hence the lower effect size). Thus, numeric comparison to other
387 approaches (that do not include rescaling) are problematic. It also may be possible to improve
388 performance—including obtaining effect sizes concomitant with other motion correction
389 approaches if that is judged critical by tuning the algorithm, although it is unclear how that
390 process could be automated.

391 One unexpected result was the relatively poor performance of ICA denoising. Although
392 individual exceptions exist, the approach produced consistently low scores on all evaluation
393 metrics. However, ICA denoising was implemented here using FSL's ICA-AROMA. This
394 package was selected because it does not require classifier training. More sophisticated ICA
395 denoising tools such as MELODIC or ICA-FIX involve a visual review of training data to
396 generate a set of noise classifiers based on the temporal, spatial, and frequency characteristics
397 of identified artifacts (Salimi-Khorshidi et al., 2014; Griffanti et al., 2017). These options were not
398 considered here because we sought to evaluate tools for motion correction that could be
399 implemented in an automated pipeline. The potential of ICA for denoising task-based data
400 should not be dismissed; rather, our results only indicate that the use of untrained ICA is
401 probably suboptimal compared to other options, many of which are also less computationally
402 intensive.

403 Frame censoring has appeared in several recent task-based studies (O'Hearn et al.,
404 2016; Bakkour et al., 2017; Davis et al., 2017). In fact, it was an experience with frame
405 censoring in the analysis of in-scanner speech production (Rogers et al., 2020) that motivated
406 our interest in comparing motion correction approaches. We found that modest levels of frame
407 censoring (e.g., 2–5% data loss) revealed a regional activation in high-motion subjects that
408 appeared in low-motion subjects but was not apparent when standard (RP6) motion
409 compensation was used. This suggested that use of a discrete rather than a continuous
410 nuisance regressor may better preserve task variance in some applications. However, a more
411 nuanced picture emerges from the present results, which suggest frame censoring is neither
412 universally superior to nor worse than RP6. One possibility is that frame censoring performance
413 involves a complex interaction between data quantity and quality. As each censored frame
414 introduces an additional regressor to the design matrix, eventually the reduction in error
415 variance may be overwhelmed by a loss of model degrees of freedom. This is anecdotally
416 supported by a decline in many of the metric results observed here at the most stringent FD or

417 DVARS thresholds, an effect that was even more pronounced when 40% maximal censoring
418 was explored in pilot work (data not shown). A prerequisite to improving frame censoring
419 performance in future work would be to quantify this tradeoff.

420 One might argue that frame censoring should be based on a selected fixed threshold
421 rather than a targeted percent data loss. The present results offer only mixed support for such a
422 position. We investigated applying a fixed FD threshold of 0.9 to these data (**Supplemental**
423 **Figure 1**). This value was used by Siegel and colleagues (2014) in their exploration of frame
424 censoring and has since been used in published functional studies (e.g., Davis et al., 2017). In
425 most of the datasets considered here, a 0.9 FD threshold would have resulted in less than 1%
426 of frames being censored. This would be a reasonable amount of data loss, and might lead to
427 some improvements compared to a standard RP6 approach (although we did not test this
428 directly). However, ds000228/adults, ds001748/teens, and ds002382/YA would have incurred a
429 1-2% data loss, ds001748/child and ds002382/OA approximately 5% data loss, and
430 ds000228/child approximately 13% data loss. These outcomes do not correspond to the best
431 performance obtained across all approaches. Whole-brain or ROI-constrained maximum-t
432 metrics peak at these values in some, but not all, datasets. Mean effect size and Dice
433 coefficients add little to the evaluation as they appear largely insensitive to frame censoring
434 thresholds in this range. Taken together, these results suggest there is no single threshold value
435 that will optimize frame censoring for all applications.

436 Finally, it should be noted that we have focused on retrospective correction—that is,
437 strategies for dealing with motion in existing data. A complementary approach would be to
438 reduce head motion during acquisition. Protocols have been developed that offer promise to
439 reduce subject motion, including movie viewing (Greene et al., 2018), custom head molds
440 (Power et al., 2019), and providing feedback to participants (Dosenbach et al., 2017; Krause et
441 al., 2019). However, these have not yet been widely adopted, nor are all compatible with task-
442 based fMRI. With increasing awareness of the challenges caused by participant motion,
443 perhaps greater interest in motion reduction (as opposed to mitigation) will follow.

444 Clearly, the present results do not identify unequivocal guidelines to select a motion
445 correction strategy. Given the variability observed across datasets, with identical processing
446 pipelines, exploring multiple strategies in a given dataset may be the best way of reducing
447 motion artifacts, adding another set of parameters to an already large space of possible
448 analyses (Carp, 2012; Poldrack et al., 2017; Botvinik-Nezer et al., 2020). Our results suggest
449 that—frustratingly—no single motion correction strategy will give optimal results on every metric
450 in every study, and that choices require considering both the nature of the specific data of
451 interest and the most relevant outcome measure.

452
453

454

455

Acknowledgments

456 This work was supported by grants R01 DC014281, R01 DC016594, R21 DC016086, and T32
457 EB014855 (to A.B.) from the US National Institutes of Health. OpenNeuro is supported by NSF
458 Grant OCI-1131441.

459

460

461

462

References

- 463 Ardekani BA, Bachman AH, Helpert JA (2001) A quantitative comparison of motion detection
464 algorithms in fMRI. *Magn Reson Imaging* 19:959–963.
- 465 Ashburner J, Friston KJ (2004) Rigid Body Registration. In: *Human Brain Function, Second*.
466 (Frackowiak RSJ, Friston KJ, Frith CD, Dolan RJ, Price CJ, Zeki S, Ashburner J, Penny W,
467 eds), pp 635–653. New York: Elsevier.
- 468 Ashburner J, Friston KJ (2005) Unified segmentation. *Neuroimage* 26:839–851.
- 469 Bakkour A, Lewis-Peacock JA, Poldrack RA, Schonberg T (2017) Neural mechanisms of cue-
470 approach training. *Neuroimage* 151:92–104.
- 471 Bennett CM, Miller MB (2010) How reliable are the results from functional magnetic resonance
472 imaging? *Ann N Y Acad Sci* 1191:133–155.
- 473 Bianciardi M, Fukunaga M, van Gelderen P, Horowitz SG, de Zwart JA, Shmueli K, Duyn JH
474 (2009) Sources of functional magnetic resonance imaging signal fluctuations in the human
475 brain at rest: a 7 T study. *Magn Reson Imaging* 27:1019–1029.
- 476 Birn RM, Molloy EK, Patriat R, Parker T, Meier TB, Kirk GR, Nair VA, Meyerand ME,
477 Prabhakaran V (2013) The effect of scan length on the reliability of resting-state fMRI
478 connectivity estimates. *Neuroimage* 83:550–558.
- 479 Botvinik-Nezer R et al. (2020) Variability in the analysis of a single neuroimaging dataset by
480 many teams. *Nature* 582:84–88.
- 481 Carp J (2012) On the plurality of (methodological) worlds: estimating the analytic flexibility of
482 fMRI experiments. *Front Neurosci* 6:149.
- 483 Courtney AL, PeConga EK, Wagner DD, Rapuano KM (2018) Calorie information and dieting
484 status modulate reward and control activation during the evaluation of food images. *PLoS*
485 *One* 13:e0204744.
- 486 Cusack R, Vicente-Grabovetsky A, Mitchell DJ, Wild CJ, Auer T, Linke AC, Peelle JE (2015)
487 Automatic analysis (aa): Efficient neuroimaging workflows and parallel processing using
488 Matlab and XML. *Front Neuroinform* 8:90.
- 489 Davis T, Goldwater M, Giron J (2017) From Concrete Examples to Abstract Relations: The
490 Rostrolateral Prefrontal Cortex Integrates Novel Examples into Relational Categories.
491 *Cereb Cortex* 27:2652–2670.
- 492 Diedrichsen J, Shadmehr R (2005) Detecting and adjusting for artifacts in fMRI time series data.
493 *Neuroimage* 27:624–634.
- 494 Dockès J, Poldrack RA, Primet R, Gözükan H, Yarkoni T, Suchanek F, Thirion B, Varoquaux G
495 (2020) NeuroQuery, comprehensive meta-analysis of human brain mapping. *Elife* 9
496 Available at: <http://dx.doi.org/10.7554/eLife.53385>.
- 497 Dosenbach NUF, Koller JM, Earl EA, Miranda-Dominguez O, Klein RL, Van AN, Snyder AZ,
498 Nagel BJ, Nigg JT, Nguyen AL, Wesevich V, Greene DJ, Fair DA (2017) Real-time motion
499 analytics during brain MRI improve data quality and reduce costs. *Neuroimage* 161:80–93.

- 500 Duncan KJ, Pattamadilok C, Knierim I, Devlin JT (2009) Consistency and variability in functional
501 localisers. *Neuroimage* 46:1018–1026.
- 502 Eickhoff SB, Stephan KE, Mohlberg H, Grefkes C, Fink GR, Amunts K, Zilles K (2005) A new
503 SPM toolbox for combining probabilistic cytoarchitectonic maps and functional imaging
504 data. *Neuroimage* 25:1325–1335.
- 505 Elliott ML, Knodt AR, Ireland D, Morris ML, Poulton R, Ramrakha S, Sison ML, Moffitt TE, Caspi
506 A, Hariri AR (2020) What Is the Test-Retest Reliability of Common Task-Functional MRI
507 Measures? *New Empirical Evidence and a Meta-Analysis*. *Psychol Sci* 31:792–806.
- 508 Friston KJ, Williams S, Howard R, Frackowiak RS, Turner R (1996) Movement-related effects in
509 fMRI time-series. *Magn Reson Med* 35:346–355.
- 510 Fynes-Clinton S, Marstaller L, Burianová H (2019) Differentiation of functional networks during
511 long-term memory retrieval in children and adolescents. *Neuroimage* 191:93–103.
- 512 Gorgolewski KJ, Storkey A, Bastin ME, Whittle IR, Wardlaw JM, Pernet CR (2013a) A test-retest
513 functional MRI dataset for motor, language and spatial attention functions. Available at:
514 <http://gigadb.org/dataset/100051>.
- 515 Gorgolewski KJ, Storkey AJ, Bastin ME, Whittle I, Pernet CR (2013b) Single subject fMRI test-
516 retest reliability metrics and confounding factors. *Neuroimage* 69:231–243.
- 517 Gratton C, Dworetzky A, Coalson RS, Adeyemo B, Laumann TO, Wig GS, Kong TS, Gratton G,
518 Fabiani M, Barch DM, Tranel D, Dominguez OM-, Fair DA, Dosenbach NUF, Snyder AZ,
519 Perlmutter JS, Petersen SE, Campbell MC (2020a) Removal of high frequency
520 contamination from motion estimates in single-band fMRI saves data without biasing
521 functional connectivity. *Neuroimage*:116866.
- 522 Gratton C, Kraus BT, Greene DJ, Gordon EM, Laumann TO, Nelson SM, Dosenbach NUF,
523 Petersen SE (2020b) Defining Individual-Specific Functional Neuroanatomy for Precision
524 Psychiatry. *Biol Psychiatry* 88:28–39.
- 525 Greene DJ, Koller JM, Hampton JM, Wesevich V, Van AN, Nguyen AL, Hoyt CR, McIntyre L,
526 Earl EA, Klein RL, Shimony JS, Petersen SE, Schlaggar BL, Fair DA, Dosenbach NUF
527 (2018) Behavioral interventions for reducing head motion during MRI scans in children.
528 *Neuroimage* 171:234–245.
- 529 Griffanti L, Douaud G, Bijsterbosch J, Evangelisti S, Alfaro-Almagro F, Glasser MF, Duff EP,
530 Fitzgibbon S, Westphal R, Carone D, Beckmann CF, Smith SM (2017) Hand classification
531 of fMRI ICA noise components. *Neuroimage* 154:188–205.
- 532 Jenkinson M, Bannister PR, Brady JM, Smith SM (2002) Improved optimisation for the robust
533 and accurate linear registration and motion correction of brain images. *Neuroimage*
534 17:825–841.
- 535 Jenkinson M, Beckmann CF, Behrens TEJ, Woolrich MW, Smith SM (2012) FSL. *Neuroimage*
536 62:782–790.
- 537 Johnstone T, Ores Walsh KS, Greischar LL, Alexander AL, Fox AS, Davidson RJ, Oakes TR
538 (2006) Motion correction and the use of motion covariates in multiple-subject fMRI analysis.

- 539 Hum Brain Mapp 27:779–788.
- 540 Kelly AMC, Uddin LQ, Biswal BB, Castellanos FX, Milham MP (2008) Competition between
541 functional brain networks mediates behavioral variability. *Neuroimage* 39:527–537.
- 542 Krause F, Benjamins C, Eck J, Lührs M, van Hoof R, Goebel R (2019) Active head motion
543 reduction in magnetic resonance imaging using tactile feedback. *Hum Brain Mapp*
544 40:4026–4037.
- 545 Laumann TO, Gordon EM, Adeyemo B, Snyder AZ, Joo SJ, Chen M-Y, Gilmore AW,
546 McDermott KB, Nelson SM, Dosenbach NUF, Schlaggar BL, Mumford JA, Poldrack RA,
547 Petersen SE (2015) Functional System and Areal Organization of a Highly Sampled
548 Individual Human Brain. *Neuron* 87:657–670.
- 549 Lemieux L, Salek-Haddadi A, Lund TE, Laufs H, Carmichael D (2007) Modelling large motion
550 events in fMRI studies of patients with epilepsy. *Magn Reson Imaging* 25:894–901.
- 551 Lewis-Peacock JA, Postle BR (2008) Temporary activation of long-term memory supports
552 working memory. *J Neurosci* 28:8765–8771.
- 553 Markiewicz CJ, Gorgolewski KJ, Feingold F, Blair R, Halchenko YO, Miller E, Hardcastle N,
554 Wexler J, Esteban O, Goncalves M, Jwa A, Poldrack RA (2021) OpenNeuro: An open
555 resource for sharing of neuroimaging data. *bioRxiv*:2021.06.28.450168 Available at:
556 <https://www.biorxiv.org/content/10.1101/2021.06.28.450168v1.full.pdf+html> [Accessed July
557 5, 2021].
- 558 Oakes TR, Johnstone T, Ores Walsh KS, Greischar LL, Alexander AL, Fox AS, Davidson RJ
559 (2005) Comparison of fMRI motion correction software tools. *Neuroimage* 28:529–543.
- 560 O’Hearn K, Velanova K, Lynn A, Wright C, Hallquist M, Minshew N, Luna B (2016)
561 Abnormalities in brain systems supporting individuation and enumeration in autism. *Autism*
562 *Res* 9:82–96.
- 563 Patel AX, Kundu P, Rubinov M, Jones PS, Vértes PE, Ersche KD, Suckling J, Bullmore ET
564 (2014) A wavelet method for modeling and despiking motion artifacts from resting-state
565 fMRI time series. *Neuroimage* 95:287–304.
- 566 Poldrack RA, Baker CI, Durnez J, Gorgolewski KJ, Matthews PM, Munafò MR, Nichols TE,
567 Poline JB, Vul E, Yarkoni T (2017) Scanning the horizon: towards transparent and
568 reproducible neuroimaging research. *Nat Rev Neurosci* 18:115–126.
- 569 Poldrack RA, Barch DM, Mitchell JP, Wager TD, Wagner AD, Devlin JT, Cumba C, Koyejo O,
570 Milham MP (2013) Toward open sharing of task-based fMRI data: the OpenfMRI project.
571 *Front Neuroinform* 7:12.
- 572 Power JD (2017) A simple but useful way to assess fMRI scan qualities. *Neuroimage* 154:150–
573 158.
- 574 Power JD, Barnes KA, Snyder AZ, Schlaggar BL, Petersen SE (2012) Spurious but systematic
575 correlations in functional connectivity MRI networks arise from subject motion. *Neuroimage*
576 59:2142–2154.
- 577 Power JD, Schlaggar BL, Petersen SE (2015) Recent progress and outstanding issues in

- 578 motion correction in resting state fMRI. *Neuroimage* 105:536–551.
- 579 Power JD, Silver BM, Silverman MR, Ajodan EL, Bos DJ, Jones RM (2019) Customized head
580 molds reduce motion during resting state fMRI scans. *Neuroimage* 189:141–149.
- 581 Pruim RHR, Mennes M, van Rooij D, Llera A, Buitelaar JK, Beckmann CF (2015) ICA-AROMA:
582 A robust ICA-based strategy for removing motion artifacts from fMRI data. *Neuroimage*
583 112:267–277.
- 584 Richardson H, Lisandrelli G, Riobueno-Naylor A, Saxe R (2018) Development of the social brain
585 from age three to twelve years. *Nat Commun* 9:1027.
- 586 Rogers CS, Jones MS, McConkey S, Spehar B, Van Engen KJ, Sommers MS, Peelle JE (2020)
587 Age-related differences in auditory cortex activity during spoken word recognition.
588 *Neurobiology of Language* 1:452–473.
- 589 Salimi-Khorshidi G, Douaud G, Beckmann CF, Glasser MF, Griffanti L, Smith SM (2014)
590 Automatic denoising of functional MRI data: combining independent component analysis
591 and hierarchical fusion of classifiers. *Neuroimage* 90:449–468.
- 592 Satterthwaite TD, Ciric R, Roalf DR, Davatzikos C, Bassett DS, Wolf DH (2019) Motion artifact
593 in studies of functional connectivity: Characteristics and mitigation strategies. *Hum Brain*
594 *Mapp* 40:2033–2051.
- 595 Siegel JS, Power JD, Dubis JW, Vogel AC, Church JA, Schlaggar BL, Petersen SE (2014)
596 Statistical improvements in functional magnetic resonance imaging analyses produced by
597 censoring high-motion data points. *Hum Brain Mapp* 35:1981–1996.
- 598 Smith SM, Beckmann CF, Ramnani N, Woolrich MW, Bannister PR, Jenkinson M, Matthews
599 PM, McGonigle DJ (2005) Variability in fMRI: a re-examination of inter-session differences.
600 *Hum Brain Mapp* 24:248–257.
- 601 Smyser CD, Snyder AZ, Neil JJ (2011) Functional connectivity MRI in infants: exploration of the
602 functional organization of the developing brain. *Neuroimage* 56:1437–1452.
- 603 Van Dijk KRA, Sabuncu MR, Buckner RL (2012) The influence of head motion on intrinsic
604 functional connectivity MRI. *Neuroimage* 59:431–438.