Running Head: FMRI RELIABILITY

How reliable are the results from functional magnetic resonance imaging?

Craig M. Bennett¹ and Michael B. Miller¹

¹ Department of Psychology, University of California at Santa Barbara, Santa Barbara, California, 93106

Word Count (excluding references, abstract): 9,615

Keywords: fMRI statistics reliability

Correspondence should be addressed to:

Craig M. Bennett Department of Psychology University of California, Santa Barbara Santa Barbara, CA 93106

415.994.7747 bennett@psych.ucsb.edu http://prefrontal.org

Abstract

Functional magnetic resonance imaging is one of the most important methods for in vivo investigation of cognitive processes in the human brain. Within the last two decades an explosion of research has emerged using fMRI, revealing the underpinnings of everything from motor and sensory processes to the foundations of social cognition. While these results have revealed the potential of neuroimaging, important questions regarding the reliability of these results remain unanswered. In this chapter we take a close look at what is currently known about the reliability of fMRI findings. First, we examine the many factors that influence the quality of acquired fMRI data. We also conduct a review of the existing literature to determine if some measure of agreement has emerged regarding the reliability of fMRI. Finally, we provide commentary on ways to improve fMRI reliability and what questions remain unanswered. Reliability is the foundation on which scientific investigation is based. How reliable are the results from fMRI?

Reliability is the cornerstone of any scientific enterprise. Issues of research validity and significance are relatively meaningless if the results of our experiments are not trustworthy. It is the case that reliability can vary greatly depending on the tools being used and what is being measured. Therefore, it is imperative that any scientific endeavor be aware of the reliability of its measurements.

Surprisingly, most fMRI researchers have only a vague idea of how reliable their results are. Reliability is not a typical topic of conversation between most investigators and only a small fraction of papers investigating fMRI reliability have been published. This became an important issue in 2009 as a paper by Vul, Harris, Winkielman, and Pashler set the stage for debate (2009). Their paper, originally entitled "Voodoo Correlations in Social Neuroscience", was focused on a statistical problem known as the 'non-independence error'. Critical to their argument was the reliability of functional imaging results. Vul et al. argued that test-retest variability of fMRI results placed an 'upper bound' on the strength of possible correlations between fMRI data and behavioral measures:

 $r_{(ObservedA,ObservedB)} = r_{(A,B)} * sqrt(reliability_A * reliability_B)$

This calculation reflects that the strength of a correlation between two measures is a product of the measured relationship and the reliability of the measurements (Nunnally, 1970; Vul et al., 2009). Vul et al. specified that behavioral measures of personality and emotion have a reliability of around 0.8 and that fMRI results have a reliability of around 0.7. Not everyone agreed. Across several written exchanges multiple research groups debated what the "actual reliability" of fMRI was. Jabbi et al. stated that the reliability of fMRI could be as high as 0.98 (2009). Lieberman et al. split the difference and argued that fMRI reliability was likely around 0.90 (2009). While much ink was spilled debating the reliability of fMRI results, very little consensus was reached regarding an appropriate approximation of its value.

The difficulty of detecting signal (what we are trying to measure) from amongst a sea of noise (everything else we don't care about) is a constant struggle for all scientists. It influences what effects can be examined and is directly tied to the reliability of research results. What follows in this chapter is a multifaceted examination of fMRI reliability. We examine why reliability is a critical metric of fMRI data, discuss what factors influence the quality of the blood oxygen level dependent (BOLD) signal, and investigate the existing reliability literature to determine if some measure of agreement has emerged across studies. Fundamentally, there is

one critical question that this chapter seeks to address: if you repeat your fMRI experiment, what is the likelihood you will get the same result?

Pragmatics of Reliability

Why worry about reliability at all? As long as investigators are following accepted statistical practices and being conservative in the generation of their results, why should the field be bothered with how reproducible the results might be? There are, at least, four primary reasons why test-retest reliability should be a concern for all fMRI researchers.

<u>Scientific truth.</u> While it is a simple statement that can be taken straight out of an undergraduate research methods course, an important point must be made about reliability in research studies: it is the foundation on which scientific knowledge is based. Without reliable, reproducible results no study can effectively contribute to scientific knowledge. After all, if a researcher obtains a different set of results today than they did yesterday, what has really been discovered? To ensure the long-term success of functional neuroimaging it is critical to investigate the many sources of variability that impact reliability. It is a strong statement, but if results do not generalize from one set of subjects to another or from one scanner to another then the findings are of little value scientifically.

<u>Clinical and Diagnostic Applications.</u> The longitudinal assessment of changes in regional brain activity is becoming increasingly important for the diagnosis and treatment of clinical disorders. One potential use of fMRI is for the localization of specific cognitive functions before surgery. A good example is the localization of language function prior to tissue resection for epilepsy treatment (Fernandez et al., 2003). This is truly a case where an investigator does not want a slightly different result each time they conduct the scan. If fMRI is to be used for surgical planning or clinical diagnostics then any issues of reliability must be quantified and addressed.

<u>Evidentiary Applications</u>. The results from functional imaging are increasingly being submitted as evidence into the United States legal system. For example, results from a commercial company called No Lie MRI (San Diego, CA; http://www.noliemri.com/) were introduced into a juvenile sex abuse case in San Diego during the spring of 2009. The defense was attempting to introduce the fMRI results as scientific justification of their client's claim of innocence. A concerted effort from imaging scientists, including in-person testimony from Marc

Raichle, eventually forced the defense to withdraw the request. While the fMRI results never made it into this case, it is clear that fMRI evidence will be increasingly common in the courtroom. What are the larger implications if the reliability of this evidence is not as trustworthy as we assume?

Scientific Collaboration. A final pragmatic dimension of fMRI reliability is the ability to share data between researchers. This is already a difficult challenge, as each scanner has its own unique sources of error that become part of the data (Jovicich et al., 2006). Early evidence has indicated that the results from a standard cognitive task can be quite similar across scanners (Casey et al., 1998; Friedman et al., 2008). Still, concordance of results remains an issue that must be addressed for large-scale, collaborative inter-center investigations. The ultimate level of reliability is the reproducibility of results from any equivalent scanner around the world and the ability to integrate this data into larger investigations.

What Factors Influence fMRI Reliability?

The ability of fMRI to detect meaningful signals is limited by a number of factors that add error to each measurement. Some of these factors include thermal noise, system noise in the scanner, physiological noise from the subject, non-task related cognitive processes, and changes in cognitive strategy over time (Huettel et al., 2008; Kruger and Glover, 2001). The concept of reliability is, at its core, a representation of the ability to routinely detect relevant signals from this background of meaningless noise. If a voxel timeseries contains a large amount of signal then the primary sources of variability are actual changes in blood flow related to neural activity within the brain. Conversely, in a voxel containing a large amount of noise the measurements are dominated by error and would not contain meaningful information. By increasing the amount of signal, or decreasing the amount of noise, a researcher can effectively increase the quality and reliability of acquired data.

The quality of data in magnetic resonance imaging is typically measured using the signalto-noise ratio (SNR) of the acquired images. The goal is to maximize this ratio. Two kinds of SNR are important for functional MRI. The first is the image SNR. It is related to the quality of data acquired in a single fMRI volume. Image SNR is typically computed as the mean signal value of all voxels divided by the standard deviation of all voxels in a single image:

$SNR_{image} = \mu_{image} / \sigma_{image}$

Increasing the image SNR will improve the quality of data at a single point in time. However, most important for functional neuroimaging is the amount of signal present in the data across time. This makes the temporal SNR (tSNR) perhaps the most important metric of data for functional MRI. It represents the signal-to-noise ratio of the timeseries at each voxel:

$SNR_{temporal} = \mu_{timeseries} \; / \; \sigma_{timeseries}$

The tSNR is not the same across all voxels in the brain. Some regions will have higher or lower tSNR depending on location and constitution. For example, there are documented differences in tSNR between gray matter and white matter (Bodurka et al., 2005). The typical tSNR of fMRI can also vary depending on the same factors that influence image SNR.

Another metric of data quality is the contrast-to-noise ratio (CNR). This refers to the ability to maximize differences between signal intensity in different areas in an image (image CNR) or to maximize differences between different points in time (temporal CNR). With regard to functional neuroimaging, the temporal CNR represents the maximum relative difference in signal intensity that is represented within a single voxel. In a voxel with low CNR there would be very little difference between two conditions of interest. Conversely, in a voxel with high CNR there would be relatively large differences between two conditions of interest. The image CNR is not critical to fMRI, but having a high temporal CNR is very important for detecting task effects.

It is generally accepted that fMRI is a rather noisy measurement with a characteristically low tSNR, requiring extensive signal averaging to achieve effective signal detection (Murphy et al., 2007). The following sections provide greater detail on the influence of specific factors on the SNR/tSNR of functional MRI data. We break these factors down by the influence of differences in image acquisition, the image analysis pipeline, and the contribution of the subjects themselves.

SNR influences of MRI acquisition

The typical high-field MRI scanner is a precision superconducting device constructed to very exact manufacturing tolerances. Still, the images it produces can be somewhat variable depending on a number of hardware and software variables. With regard to hardware, one well-known influence on the signal to noise ratio of MRI is the strength of the primary B0 magnetic

field (Bandettini et al., 1994; Ogawa et al., 1993). Doubling this field, such as moving from 1.5 Tesla to a 3.0 Tesla field strength, can theoretically double the SNR of the data. The B0 field strength is especially important for fMRI, which relies on magnetic susceptibility effects to create the blood oxygen level dependent (BOLD) signal (Turner et al., 1993). Hoenig et al. showed that, relative to a 1.5 Tesla magnet, a 3.0 Tesla fMRI acquisition had 60-80% more significant voxels (2005). They also demonstrated that the CNR of the results was 1.3 times higher than those obtained at 1.5 Tesla. The strength and slew rate of the gradient magnets can have a similar impact on SNR. Advances in head coil design are also notable, as parallel acquisition head coils have increased radiofrequency reception sensitivity.

It is important to note that there are negative aspects of higher field strength as well. Artifacts due to physiological effects and susceptibility are all increasingly pronounced at higher fields. The increased contribution of physiological noise reduces the expected gains in SNR at high field (Kruger and Glover, 2001). The increasing contribution of susceptibility artifacts can virtually wipe out areas of orbital prefrontal cortex and inferior temporal cortex (Jezzard and Clare, 1999). Also, in terms of tSNR there are diminishing returns with each step up in B0 field strength. At typical fMRI spatial resolution values tSNR approaches an asymptotic limit between 3 Tesla and 7 Tesla (Kruger and Glover, 2001; Triantafyllou et al., 2005).

Looking beyond the scanner hardware, the parameters of the fMRI acquisition can also have a significant impact on the SNR/CNR of the final images. For example, small changes in the voxel size of a sequence can dramatically alter the final SNR. Moving from 1.5 mm³ voxels to 3.0 mm³ voxels can potentially increase the acquisition SNR by a factor of eight, but at a cost of spatial resolution. Some other acquisition variables that will influence the acquired SNR/CNR are : repetition time (TR), echo time (TE), bandwidth, slice gap, and k-space trajectory. For example, Moser et al. found that optimizing the flip angle of their acquisition could approximately double the SNR of their data in a visual stimulation task (1996). Further, the effect of each parameter varies according to the field strength of the magnet (Triantafyllou et al., 2005). The optimal parameter set for a 3 Tesla system may not be optimal with a 7 Tesla system.

The ugly truth is that any number of factors in the control room or magnet suite can increase noise in the images. A famous example from one imaging center was when the broken filament from a light bulb in a distant corner of the magnet suite started causing visible sinusoidal striations in the acquired EPI images. This is an extreme example, but it makes the point that the scanner is a precision device that is designed to operate in a narrow set of well-defined circumstances. Any deviation from those circumstances will increase noise, thereby reducing SNR and reliability.

SNR considerations of analysis methods

The methods used to analyze fMRI data will affect the reliability of the final results. In particular, those steps taken to reduce known sources of error are critical to increasing the final SNR/CNR of preprocessed images. For example, spatial realignment of the EPI data can have a dramatic effect on lowering movement-related variance and has become a standard part of fMRI preprocessing (Oakes et al., 2005; Zhilkin and Alexander, 2004). Recent algorithms can also help remove remaining signal variability due to magnetic susceptibility induced by movement (Andersson et al., 2001). Temporal filtering of the EPI timeseries can reduce undesired sources of noise by frequency. The use of a high-pass filter is a common method to remove low-frequency noise, such as signal drift due to the scanner (Kiebel and Holmes, 2007). Spatial smoothing of the data can also improve the SNR/CNR of an image. There is some measure of random noise added to the true signal of each voxel during acquisition. Smoothing across voxels can help to average out error across the area of the smoothing filter (Mikl et al., 2008). It can also help account for local differences in anatomy across subjects. Smoothing is most often done using a Gaussian kernel of approximately 6-12 mm³ FWHM.

There has been some degree of standardization regarding preprocessing and statistical approaches in fMRI. For instance, Mumford and Nichols found that approximately 92% of group fMRI results were computed using an ordinary least squares (OLS) estimation of the general linear model (2009). Comparison studies with carefully standardized processing procedures have shown that the output of standard software packages can be very similar (Gold et al., 1998; Morgan et al., 2007). However, in actual practice the diversity of tools and approaches in fMRI increases the variability between sets of results. The functional imaging analysis contest (FIAC) in 2005 demonstrated that prominent differences existed between fMRI results generated by different groups using the same original dataset. On reviewing the results the organizers concluded that brain regions exhibiting robust signal changes could be quite similar across analysis techniques, but the detection of areas with lower signal was highly

variable (Poline et al., 2006). It remains the case that decisions made by the researcher regarding how to analyze the data will impact what results are found.

Strother et al. have done a great deal of research into the influence of image processing pipelines using a predictive modeling framework (2004; 2002; Zhang et al., 2009). They found that small changes in the processing pipeline of fMRI images have a dramatic impact on the final statistics derived from that data. Some steps, such as slice timing correction, were found to have little influence on the results from experiments with a block design. This is logical, given the relative insensitivity of block designs to small temporal shifts. However, the steps of motion correction, high-pass filtering, and spatial smoothing were found to significantly improve the analysis. They reported that the optimization of preprocessing pipelines improved both intra-subject and between-subject reproducibility of results (Zhang et al., 2009). Identifying an optimal set of processing steps and parameters can dramatically improve the sensitivity of an analysis.

SNR influences of participants

The MRI system and fMRI analysis methods have received a great deal of attention with regard to SNR. However, one area that may have the greatest contribution to fMRI reliability is how stable/unstable the patterns of activity within a single subject can be. After all, a test-retest methodology involving human beings is akin to hitting a moving target. Any discussion of test-retest reliability in fMRI has to take into consideration the fact that the cognitive state of a subject is variable over time.

There are two important ways that a subject can influence reliability within a test-retest experimental design. The first involves within-subject changes that take place over the course of a single session. For instance, differences in attention and arousal can significantly modulate subsequent responses to sensory stimulation (Munneke et al., 2008; Peyron et al., 1999; Sterr et al., 2007). Variability can also be caused by evolving changes in cognitive strategy used during tasks like episodic retrieval (Miller et al., 2001; Miller et al., 2002). If a subject spontaneously shifts to a new decision criterion midway during a session then the resulting data may reflect the results of two different cognitive processes. Finally, learning will take place with continued task experience, shifting the pattern of activity as brain regions are engaged and disengaged during task-relevant processing (Grafton et al., 1995; Poldrack et al., 1999; Rostami et al., 2009). For

studies investigating learning this is a desired effect, but for others this is an undesired source of noise.

The second influence on reliability is related to physiological and cognitive changes that may take place within a subject between the test and retest sessions. Within 24 hours an infinite variety of reliability-reducing events can take place. All of the above factors may show changes over the days, weeks, months, or years between scans. These changes may be even more dramatic depending on the amount of time between scanning sessions.

Estimates of fMRI Reliability

A diverse array of methods have been created for measuring the reliability of fMRI. What differs between them is the specific facet of reliability they are intended to quantify. Some methods are only concerned with significant voxels. Other methods address similarity in the magnitude of estimated activity across all voxels. The choice of how to calculate reliability often comes down to which aspect of the results are desired to remain stable over time.

<u>Measuring stability of super-threshold extent.</u> Do you want the voxels that are significant during the test scan to still be significant during the retest scan? This would indicate that super-threshold voxels are to remain above the threshold during subsequent sessions. The most prevalent method to quantify this reliability is the cluster overlap method. The cluster overlap method is a measure revealing what set of voxels are considered to be super-threshold during both test and retest sessions.

Two approaches have been used to calculate cluster overlap. The first, and by far most prevalent, is a measure of similarity known as the Dice coefficient. It was first used to calculate fMRI cluster overlap by Rombouts et al. and has become a standard measure of result similarity (1997). It is typically calculated by the following equation:

$$R_{overlap} = 2 \times (V_{overlap}) / (V_1 + V_2)$$

Results from the Dice equation can be interpreted as the number of voxels that will overlap divided by the average number of significant voxels across sessions. Another approach to calculating similarity is the Jaccard index. The Jaccard index has the advantage of being readily interpretable as the percent of voxels that are shared, but is infrequently used in the investigation of reliability. It is typically calculated by the following equation:

$$\mathbf{R}_{\text{overlap}} = \mathbf{V}_{\text{overlap}} / (\mathbf{V}_1 + \mathbf{V}_2 - \mathbf{V}_{\text{overlap}})$$

Results from the Jaccard equation can be interpreted as the number of overlapping voxels divided by the total number of unique voxels in all sessions. For both the Dice and Jaccard methods a value of 1.0 would indicate that all super-threshold voxels identified during the test scan were also active in the retest scan, and vice-versa. A value of 0.0 would indicate that no voxels in either scan were shared between the test and retest sessions. See Figure 1 for a graphical representation of overlapping results from two runs in an example dataset.

- FIGURE 1 ABOUT HERE -

The main limitation of all cluster overlap methods is that they are highly dependent on the statistical threshold used to define what is 'active'. Duncan et al. demonstrated that the reported reliability of the cluster overlap method decreases as the significance threshold is increased (2009). Similar results were reported by Rombouts et al., who found nonlinear changes in cluster overlap reliability across multiple levels of significance (1998).

These overlap statistics seek to represent the proportion of voxels that remain significant across repetitions relative to the proportion that are significant in only a subset of the results. Another, similar approach would be to conduct a formal conjunction analysis between the repetitions. The goal of this approach would be to uniquely identify those voxels that are significant in all sessions. One example of this approach would be the 'Minimum Statistic compared to the Conjunction Null' (MS/CN) of Nichols et al (2005). Using this approach a researcher could threshold the results, allowing for the investigation of reliability with a statistical criterion.

A method similar to cluster overlap, called voxel counting, was reported in early papers. The use of voxel counting simply evaluated the total number of activated voxels in the test and retest images. This has proven to be a suboptimal approach for the examination of reliability, as it is done without regard to the spatial location of significant voxels (Cohen and DuBois, 1999). An entirely different set of results could be observed in each image yet they could contain the same number of significant voxels. As a result this method is no longer used.

<u>Measuring stability of activity in significant clusters</u>. Do you want the estimated magnitude of activity in each cluster to be stable between the test scan and the retest scan? This

is a more stringent criteria than simple extent reliability, as it is necessary to replicate the exact degree of activation and not simply what survives thresholding. The most standard method to quantify this reliability is through an intra-class correlation (ICC) of the time1-time2 cluster values. The intra-class correlation is different from the traditional Pearson product-moment correlation as it is specialized for data of one type, or class. While there are many versions of the ICC, it is typically taken to be a ratio of the variance of interest divided by the total variance (Bartko, 1966; Shrout and Fleiss, 1979). The ICC can be computed as follows:

$ICC = \sigma^{2}_{between} / (\sigma^{2}_{between} + \sigma^{2}_{within})$

One of the best reviews of the ICC was completed by Shrout and Fleiss, who detailed six types of ICC calculation and when each is appropriate to use (1979). One advantage of the ICC is that it can be interpreted similarly to the Pearson correlation. A value of 1.0 would indicate near-perfect agreement between the values of the test and retest sessions, as there would be no influence of within-subject variability. A value of 0.0 would indicate that there was no agreement between the values of the test and retest sessions, since within-subject variability would dominate the equation.

Studies examining reliability using intra-class correlations are often computed based on summary values from regions of interest (ROIs). Caceras et al. compared four methods commonly used to compute ROI reliability using intraclass correlations (2009). The median(ICC) is the median of the ICC values from within a ROI. ICC_{med} is the median ICC of the contrast values. ICC_{max} is the calculation of ICC values at the peak activated voxel within an activated cluster. ICC_v is defined the intra-voxel reliability, a measure of the total variability that can be explained by the intra-voxel variance.

There are several notable weaknesses to the use of ICC in calculating reliability. First, the generalization of ICC results is limited because calculation is specific to the dataset under investigation. An experiment with high inter-subject variability could have different ICC values relative to an experiment with low inter-subject variability, even if the stability of values over time is the same. As discussed later in this chapter, this can be particularly problematic when comparing the reliability of clinical disorders to that of normal controls. Second, because of the variety of ICC subtypes there can often be confusion regarding which one to use. Using an incorrect subtype can result in quite different reliability estimates (Muller and Buttner, 1994).

<u>Measuring voxelwise reliability of the whole brain</u>. Do you want to know the reliability of results on a whole-brain, voxelwise basis? Completing a voxelwise calculation would indicate that the level of activity in all voxels should remain consistent between the test and retest scans. This is the strictest criterion for reliability. It yields a global measure of concordance that indicates how effectively activity across the whole brain is represented in each test-retest pairing. Very few studies have examined reliability using this approach, but it may be one of the most valuable metrics of fMRI reliability. This is one of the few methods that gives weight to the idea that the estimated activity should remain consistent between test and retest, even if the level of activity is close to zero.

- FIGURE 2 ABOUT HERE -

Figure 2 is an example histogram plot from our own data that shows the frequency of ICC values for all voxels across the whole brain during a two-back working memory task (Bennett et al., 2009). The mean and mode of the distribution is plotted. It is quickly apparent that there is a wide range of ICC reliability values across the whole brain, with some voxels having almost no reliability and others approaching near perfect reliability.

<u>Other reliability methods</u>. Numerous other methods have also been used to measure the reliability of estimated activity. Some of these include maximum likelihood (ML), coefficient of variation (CV), and variance decomposition. While these methods are in the minority by frequency of use, this does not diminish their utility in examining reliability. This is especially true with regard to identifying the sources of test-retest variability that can influence the stability of results.

One particularly promising approach for the quantification of reliability is predictive modeling. Predictive modeling measures the ability of a training set of data to predict the structure of a testing set of data. One of the best established modeling techniques within functional neuroimaging is the nonparametric prediction, activation, influence, and reproducibility sampling (NPAIRS) approach by Strother et al. (2004; 2002). Within the NPAIRS modeling framework separate metrics of prediction and reproducibility are generated (Zhang et al., 2008). The first, prediction accuracy, evaluates classification in the temporal domain, predicting which condition of the experiment each scan belongs to. The second metric,

reproducibility, evaluates the model in the spatial domain, comparing patterns of regional brain activity over time. While this approach is far more complicated than the relatively simple cluster overlap or ICC metrics, predictive modeling does not suffer from many of the drawbacks that these methods have. NPAIRS, and other predictive modeling approaches, enable a much more thorough examination of fMRI reliability.

Some studies have investigated fMRI reliability using the Pearson product-moment (r) correlation. Intuitively this is a logical method to use, as it measures the relationship between two variables. However, it is generally held that the Pearson product-moment correlation is not an ideal measure of test-retest reliability. Safrit identified three reasons why the product-moment correlation should not be used to calculate reliability (1976). First, the Pearson product-moment correlation is setup to determine the relationship between two variables, not the stability of a single variable. Second, it is difficult to measure reliability with the Pearson product-moment correlation beyond a single test-retest pair. It becomes increasingly awkward to quantify reliability with two or more retest sessions. One can try to average over multiple pairwise Pearson product-moment correlations between the multiple sessions, but it is far easier to take the ANOVA approach of the ICC and examine it from the standpoint of between- and withinsubject variability. Third, the Pearson product-moment correlation cannot detect systematic error. This would be the case when the retest values deviate by a similar degree, such as adding a constant value to all of the original test values. The Pearson product-moment correlation would remain the same, while an appropriate ICC would indicate that the test-retest agreement is not exact. While the use of ICC measures has its own set of issues, it is generally a more appropriate tool for the investigation of test-retest reliability.

Review of Existing Reliability Estimates

Since the advent of fMRI some results have been common and quite easily replicated. For example, activity in primary visual cortex during visual stimulation has been thoroughly studied. Other fMRI results have been somewhat difficult to replicate. What does the existing literature have to say regarding the reliability of fMRI results?

There have been a number of individual studies investigating the test-retest reliability of fMRI results, but few articles have reviewed the entire body of literature to find trends across studies. To obtain a more effective estimate of fMRI reliability we conducted a survey of the

existing literature on fMRI reliability. To find papers for this investigation we searched for "testretest fMRI" using the NCBI PubMed database (www.pubmed.gov). This search yielded a total of 183 papers, 37 of which used fMRI as a method of investigation, used a general linear model to compute their results, and provided test-retest measures of reliability. To broaden the scope of the search we then went through the reference section of the 37 papers found using PubMed to look for additional works not identified in the initial search. There were 26 additional papers added to the investigation through this secondary search method. The total number of papers retrieved was 63. Each paper was examined with regard to the type of cognitive task, kind of fMRI design, number of subjects, and basis of reliability calculation.

We have separated out the results into three groups: those that used the voxel overlap method, those that used intraclass correlation, and papers that used other calculation methods. The results of this investigation can be seen in Tables 1, 2, and 3. In the examination of cluster overlap values in the literature we attempted to only include values that were observed at a similar significance threshold across all of the papers. The value we chose as the standard was p(uncorrected) < 0.001. Other deviations from this standard approach are noted in the tables.

- TABLES 1, 2, AND 3 ABOUT HERE -

Conclusions From the Reliability Review

What follows are some general points that can be taken away from the reliability survey. Some of the conclusions that follow are quantitative results from the review and some are qualitative descriptions of trends that were observed as we conducted the review.

<u>A diverse collection of methods have been used to assess fMRI reliability</u>. The first finding mirrors the above discussion on reliability calculation. A very diverse collection of methods has been used to investigate fMRI reliability. This list includes: intra-class correlation (ICC), cluster overlap, voxel counts, receiver operating characteristic (ROC) curves, maximum likelihood (ML), conjunction analysis, Cohen's kappa index, coefficient of variation (CV), Kendall's W, laterality index (LI), variance component decomposition, Pearson correlation, predictive modeling, and still others. While this diversity of methods has created converging evidence of fMRI reliability, it has also limited the ability to compare and contrast the results of existing reliability studies. Intra-class correlation and cluster overlap methods dominate the calculation of test-retest reliability. While there have been a number of methods used to investigate reliability, the two that stand out by frequency of use are cluster overlap and intra-class correlation. One advantage of these methods is that they are easy to calculate. The equations are simple to understand, easy to implement, and fast to process. A second advantage of these methods is their easy interpretation by other scientists. Even members of the general public can understand the concept behind the overlapping of clusters and most everyone is familiar with correlation values. While these techniques certainly have limitations and caveats, they seem to be the emerging standard for the analysis of fMRI reliability.

Most previous studies of reliability and reproducibility have been done with relatively few subjects. What sample size is necessary to conduct effective reliability research? Most of the studies that were reviewed used less than 10 subjects to calculate their reliability measures, with 11 subjects being the overall average across the investigation. Should reliability studies have more subjects? Since a large amount of the error variance is coming from subject-specific factors it may be wise to use larger sample sizes when assessing study reliability, as a single anomalous subject could sway study reliability in either direction. Another notable factor is that a large percentage of studies using fMRI are completed with a restricted range of subjects. Most samples will typically be recruited from a pool of university undergraduates. These samples may have a different reliability than a sample pulled at random from the larger population. Because of sample restriction the results of most test-retest investigations may not reflect the true reliability of other populations, such as children, the elderly, and individuals with clinical disorders.

<u>Reliability varies by test-retest interval</u>. Generally, increased amounts of time between the initial test scan and the subsequent retest scan will lower reliability. Still, even back-to-back scans are not perfectly reliable. The average Jaccard overlap of studies where the test and retest scans took place within the same hour was 33%. Many studies with intervals lasting three months or more had a lower overlap percentage. This is a somewhat loose guideline though. Notably, the results reported by Aron et al. had one of the longest test-retest intervals but also possessed the highest average ICC score (2006).

<u>Reliability varies by cognitive task and experimental design</u>. Motor and sensory tasks seem to have greater reliability than tasks involving higher cognition. Caceras et al. found that

the reliability of an N-back task was generally higher than that of an auditory target detection task (2009). Differences in the design of an fMRI experiment also seem to affect the reliability of results. Specifically, block designs appear to have a slight advantage over event-related designs in terms of reliability. This may be a function of the greater statistical power inherent in a block design and its increased SNR.

Significance is related to reliability, but it is not a strong correlation. Several studies have illustrated that super-threshold voxels are not necessarily more reliable than sub-threshold voxels. Caceras et al. examined the joint probability distribution of significance and reliability (2009). They found that there were some highly activated ROIs with low reliability and some sub-threshold regions that had high reliability. These ICC results fit in well with the data from cluster overlap studies. The average cluster overlap was 29%. This means that, across studies, the average number of significant voxels that will replicate is roughly one-third. This evidence speaks against the assumption that significant voxels will be far more reliable in an investigation of test-retest reliability.

An optimal threshold of reliability has not been established. There is no consensus value regarding what constitutes an acceptable level of reliability in fMRI. Is an ICC value of 0.50 enough? Should studies be required to achieve an ICC of 0.70? All of the studies in the review simply reported what the reliability values were. Few studies proposed any kind of criteria to be considered a 'reliable' result. Cicchetti and Sparrow did propose some qualitative descriptions of data based on the ICC-derived reliability of results (1981). They proposed that results with an ICC above 0.75 be considered 'excellent', results between 0.59 and 0.75 be considered 'good', results between .40 and .58 be considered 'fair', and results lower than 0.40 be considered 'poor'. More specifically to neuroimaging, Eaton et al. (2008) used a threshold of ICC > 0.4 as the mask value for their study while Aron et al. (2006) used an ICC cutoff of ICC > 0.5 as the mask value.

Inter-individual variability is consistently greater than intra-individual variability. Many studies reported both within-subject and between-subject reliability values in their results. In every case the within-subject reliability far exceeded the between-subjects reliability. Miller et al. explicitly examined variability across subjects and concluded that there are large-scale, stable differences between individuals on almost any cognitive task (2001; 2002). More recently, Miller et al. directly contrasted within- and between-subject variability (2009). They concluded

that between-subject variability was far higher than any within-subject variability. They further demonstrated that the results from one subject completing two different cognitive tasks are typically more similar than the data from two subjects doing the same task. These results are mirrored by those of Costafreda et al. who found that well over half (57%) of the variability in their fMRI data was due to between-subject variation (2007). It seems to be the case that within-subject measurements over time may vary, but they vary far less than differences in the overall pattern of activity between individuals.

There is little agreement regarding the true reliability of fMRI results. While we mention this as a final conclusion from the literature review, it is perhaps the most important point. Some studies have estimated the reliability of fMRI data to be quite high, or even close to perfect for some tasks and brain regions (Aron et al., 2006; Maldjian et al., 2002; Raemaekers et al., 2007). Other studies have been less enthusiastic, showing fMRI reliability to be relatively low (Duncan et al., 2009; Rau et al., 2007). Across the survey of fMRI test-retest reliability we found that the average ICC value was 0.50 and the average cluster overlap value was 29% of voxels (Dice overlap = 0.45, Jaccard overlap = 0.29). This represents an average across many different cognitive tasks, fMRI experimental designs, test-retest time periods, and other variables. While these numbers may not be representative of any one experiment, they do provide an effective overview of fMRI reliability.

Other Issues and Comparisons

Test-Retest Reliability in Clinical Disorders

There have been few examinations of test-retest reliability in clinical disorders relative to the number of studies with normal controls. A contributing factor to this problem may be that the scientific understanding of brain disorders is still in its infancy. It may be premature to examine clinical reliability if there is only a vague understanding of anatomical and functional abnormalities in the brain. Still, some investigators have taken significant steps forward in the clinical realm. These few investigations suggest that reliability in clinical disorders is typically lower than the reliability of data from normal controls. Some highlights of these results are listed below, categorized by disorder.

Epilepsy. Functional imaging has enormous potential to aid in the clinical diagnosis of epileptiform disorders. Focusing on fMRI, research by Di Bonaventura et al. found that the spatial extent of activity associated with fixation off sensitivity (FOS) was stable over time in epileptic patients (2005). Of greater research interest for epilepsy has been the reliability of combined EEG/fMRI imaging. Symms et al. reported that they could reliably localize interictal epileptiform discharges using EEG-triggered fMRI (1999). Waites et al. also reported the reliable detection of discharges with combined EEG/fMRI at levels significantly above chance (2005). Functional imaging also has the potential to assist in the localization of cognitive function prior to resection for epilepsy treatment. One possibility would be to use noninvasive fMRI measures to replace cerebral sodium amobarbital anesthetization (Wada Test). Fernandez et al. reported good reliability of lateralization indices (whole-brain test-retest r = 0.82) and cluster overlap measures (Dice overlap = .43, Jaccard overlap = 0.27) (2003).

Stroke. Many aspects of stroke recovery can impact the results of functional imaging data. The lesion location, size, and time elapsed since the stroke event each have the potential to alter function within the brain. These factors can also lead to increased between-subject variability relative to groups of normal controls. This is especially true when areas proximal to the lesion location contribute to specific aspects of information processing, such as speech production. Kimberley et al. found that stroke patients had generally higher ICC values relative to normal controls (2008). This mirrors the findings of Eaton et al., who showed that the average reliability of aphasia patients was approximately equal to that of normal controls as measured by ICC (2008). These results may be indicative of equivalent fMRI reliability in stroke victims, or it may be an artifact of the ICC calculation. Kimberly et al. state that increased between-subject variability of stroke patients likely falls within the moderate range of values (0.4 < ICC < 0.6).

Schizophrenia. Schizophrenia is a multidimensional mental disorder characterized by a wide array of cognitive and perceptual dysfunctions (Freedman, 2003; Morrison and Murray, 2005). While there have been a number of studies on the reliability of anatomical measures in schizophrenia there have been few that have focused on function. Manoach et al. demonstrated that the fMRI results from schizophrenic patients on a working memory task were less reliable overall than that of normal controls (2001). The reliability of significant ROIs in the schizophrenic group ranged from ICC values of -0.20 to 0.57. However, the opposite effect was

found by Whalley et al. in a group of subjects at high genetic risk for schizophrenia (no psychotic symptoms) (2009). The ICC values for these subjects were equally reliable relative to normal controls on a sentence completion task. More research is certainly needed to find consensus on reliability in schizophrenia.

<u>Aging</u>. The anatomical and functional changes that take place during aging can increase the variability of fMRI results at all levels (MacDonald et al., 2006). Clement et al. reported that cluster overlap percentages and the cluster-wise ICC values were not significantly different between normal elderly controls and patients with mild cognitive impairment (MCI) (2009). On an episodic retrieval task healthy controls had ICC values averaging 0.69 while patients diagnosed with MCI had values averaging 0.70. However, they also reported that all values for the older samples were lower than those reported for younger adults on similar tasks. Marshall et al. found that while the qualitative reproducibility of results was high, the reliability of activation magnitude during aging was quite low (2004).

It is clear that the use of intra-class correlations in clinical research must be approached carefully. As mentioned by Bosnell et al. and Kimberly et al., extreme levels of between-subject variability will artificially inflate the resulting ICC reliability estimate (Bosnell et al., 2008; Kimberley et al., 2008). Increased between-subject variability is a characteristic found in many clinical populations. Therefore, it may be the case that comparing two populations with different levels of between-subject variability may be impossible when using an ICC measure.

Reliability Across Scanners / Multicenter Studies

One area of increasing research interest is the ability to combine the data from multiple scanners into larger, integrative data sets (Van Horn and Toga, 2009). There are two areas of reliability that are important for such studies. The first is subject-level reliability, or how stable the activity of one person will be scan-to-scan. The second is group-level reliability, or how stable the group fMRI results will be from one set of subjects to another or from one scanner to another. Given the importance of multi-center collaboration it is critical to evaluate how results will differ when the data comes from a heterogeneous group of MRI scanners as opposed to a single machine. Generally, the concordance of fMRI results from center to center is quite good, but not perfect.

Casey et al. was one of the first groups to examine the reliability of results across scanners (1998). Between three imaging centers they found a 'strong similarity' in the location and distribution of significant voxel clusters. More recently, Friedman et al. found that inter-center reliability was somewhat worse than test-retest reliability across several centers with an identical hardware configuration (2008). The median ICC of their inter-center results was ICC = 0.22. Costafreda et al. also examined the reproducibility of results from identical fMRI setups (2007). Using a variance components analysis they determined that the MR system accounted for roughly 8% of the variation in the BOLD signal. This compares favorably relative to the level of between-subject variability (57%).

The reliability of results from one scanner to another seems to be approximately equal to or slightly less than the values of test-retest reliability with the same MRI hardware. Special calibration and quality control steps can be taken to ensure maximum concordance across scanners. For instance, before conducting anatomical MRI scans in the Alzheimer's Disease Neuroimaging Initiative (ADNI, http://www.loni.ucla.edu/ADNI/) a special MR phantom is typically scanned. This allows for correction of magnet-specific field inhomogeneity and maximizes the ability to compare data from separate scanners. Similar calibration measures are being discussed for functional MRI (Chiarelli et al., 2007; Friedman and Glover, 2006; Thomason et al., 2007). It may be the case that as calibration becomes standardized it will lead to increased inter-center reliability.

Other Statistical Issues in fMRI

It is important to note that a number of important fMRI statistical issues have gone unmentioned in this chapter. First, there is the problem of conducting thousands of statistical comparisons without an appropriate threshold adjustment. Correction for multiple comparisons is a necessary step in fMRI analysis that is often skipped or ignored (Bennett et al., in press). Another statistical issue in fMRI is temporal autocorrelation in the acquired timeseries. This refers to the fact that any single timepoint of data is not necessarily independent of the acquisitions that came before and after (Smith et al., 2007; Woolrich et al., 2001). Autocorrelation correction is widely available, but is not implemented by most investigators. Finally, throughout the last year the 'non-independence error' has been discussed at length. Briefly, this refers to selecting a set of voxels to create a region of interest (ROI) and then using the same measure to evaluate some statistical aspect of that region. Ideally, an independent data set should be used after the ROI has been initially defined. It is important to address these issues because they are still debated within the field and often ignored in fMRI analysis. Their correction can have a dramatic impact on how reproducible the results will be from study to study.

Conclusions

How can a researcher improve fMRI reliability?

The generation of highly reliable results requires that sources of error be minimized across a wide array of factors. An issue within any single factor can significantly reduce reliability. Problems with the scanner, a poorly designed task, or an improper analysis method could each be extremely detrimental. Conversely, elimination of all such issues is necessary for high reliability. A well maintained scanner, well designed tasks, and effective analysis techniques are all prerequisites for reliable results.

There are a number of practical ways that fMRI researchers can improve the reliability of their results. For example, Friedman and Glover reported that simply increasing the number of fMRI runs improved the reliability of their results from ICC = 0.26 to ICC = 0.58 (2006). That is quite a large jump for an additional ten or fifteen minutes of scanning. Below are some general areas where reliability can be improved.

Increase the SNR and CNR of the acquisition. One area of attention is to improve the signal-to-noise and contrast-to-noise ratios of the data collection. An easy way to do this would be to simply acquire more data. It is a zero-sum game, as increasing the number of TRs that are acquired will help improve the SNR but will also increase the task length. Subject fatigue, scanner time limitations, and the diminishing returns with each duration increase will all play a role in limiting the amount of time that can be dedicated to any one task. Still, a researcher considering a single six-minute EPI scan for their task might add additional data collection to improve the SNR of the results. With regard to the magnet, every imaging center should verify acquisition quality before scanning. Many sites conduct quality assurance scans (QA) at the beginning of each day to ensure stable operation. This has proven to be an effective method of detecting issues with the MR system before they cause trouble for investigators. It is a hassle to

cancel a scanning session when there are subtle artifacts present, but this is a better option than acquiring noisy data that does not make a meaningful contribution to the investigation. As a final thought, research groups can always start fundraising to purchase a new magnet with improved specifications. If data acquisition is being done on a 1.5 Tesla magnet with a quadrature head coil enormous gains in SNR can be made by moving to 3.0 Tesla or higher and using a parallel-acquisition head coil (Simmons et al., 2009; Zou et al., 2005).

<u>Minimize individual differences in cognitive state, both across subjects and over time.</u> Because magnet time is expensive and precious the critical component of effective task instruction can often be overlooked. Researchers would rather be acquiring data as opposed to spending additional time giving detailed instructions to a subject. However, this is a very easy way to improve the quality of the final data set. If it takes ten trials for the participant to really 'get' the task then those trials have been wasted, adding unnecessary noise to the final results. Task training in a separate laboratory session in conjunction with time in a mock MRI scanner can go a long way toward homogenizing the scanner experience for subjects. It may not always be possible to fully implement these steps, but they should not be avoided simply to reduce the time spent per subject.

For multi-session studies steps can be taken to help stabilize intra-subject changes over time. Scanning test and retest session at the same time of day can help due to circadian changes in hormone level and cognitive performance (Carrier and Monk, 2000; Huang et al., 2006; Salthouse et al., 2006). A further step to consider is minimizing the time between sessions to help stabilize the results. Much more can change over the course of a month than over the course of a week.

<u>Maximize the experiment's statistical power</u>. Power represents the ability of an experiment to reject the null hypothesis when the null hypothesis is indeed false (Cohen, 1977). For fMRI this ability is often discussed in terms of the number of subjects that will be scanned and the design of the task that will be administered, including how many volumes of data will be acquired from each subject. More subjects and more volumes almost always contribute to increasing power, but there are occasions when one may improve power more than the other. For example, Mumford and Nichols demonstrated that, when scanner time was limited, different combinations of subjects and trials could be used to achieve high levels of power (2008). For their hypothetical task it would take only five 15 second blocks to achieve 80% power if there

were 23 subjects, but it would take 25 blocks if there were only 18 subjects. These kinds of power estimations are quite useful in determining the best use of available scanner time. Tools like fmripower (http://fmripower.org) can utilize data from existing experiments to yield new information on how many subjects and scans a new experiment will require to reach a desired power level (Mumford and Nichols, 2008; Mumford et al., 2007 2007; Van Horn et al., 1998).

The structure of the stimulus presentation has a strong influence on an experiment's statistical power. The dynamic interplay between stimulus presentation and inter-stimulus jitter are important, as is knowing what contrasts will be completed once the data has been acquired. Each of these parameters can influence the power and efficiency of the experiment, later impacting the reliability of the results. Block designs tend to have greater power relative to event-related designs. One can also increase power by increasing block length, but care should be exercised not to make blocks so long that they approach the low frequencies associated with scanner drift. There are several good software tools available that will help researchers create an optimal design for fMRI experiments. OptSeq is a program that helps to maximize the efficiency of an event-related fMRI design (1999). OptimizeDesign is a set of Matlab scripts that utilize a genetic search algorithm to maximize specific aspects of the design (Wager and Nichols, 2003). Researchers can separately weight statistical power, HRF estimation efficiency, stimulus counterbalancing, and maintenance of stimulus frequency. These two programs, and others like them, are valuable tools for ensuring that the ability to detect meaningful signals is effectively maximized.

It is important to state that the reliability of a study in no way implies that an experiment has accurately assessed a specific cognitive process. The validity of a study can be quite orthogonal to its reliability – it is possible to have very reliable results from a task that mean little with regard to the cognitive process under investigation. No increase in SNR or optimization of event timing can hope to improve an experiment that is testing for the wrong thing. This makes task selection of paramount importance in the planning of an experiment. It also places a burden on the researcher in terms of effective interpretation of fMRI results once the analysis is done.

Where does neuroimaging go next?

In many ways cognitive neuroscience is still at the beginning of fMRI as a research tool. Looking back on the last two decades it is clear that functional MRI has made enormous gains in both statistical methodology and popularity. However, there is still much work to do. With specific regard to reliability, there are some specific next steps that must be taken for the continued improvement of this method.

Better Characterization of the Factors that Influence Reliability. Additional research is necessary to effectively understand what factors influence the reliability of fMRI results. The field has a good grasp of the acquisition and analysis factors that influence SNR. Still, there is relatively little knowledge regarding how stable individuals are over time and what influences that stability. Large-scale studies specifically investigating reliability and reproducibility should therefore be conducted across several cognitive domains. The end goal of this research would be to better characterize the reliability of fMRI across multiple dimensions of influence within a homogeneous set of data. Such a study would also create greater awareness of fMRI reliability in the field as a whole. The direct comparison of reliability analysis methods, including predictive modeling, should also be completed.

Meta/Mega Analysis. The increased pooling of data from across multiple studies can give a more generalized view of important cognitive processes. One method, meta-analysis, refers to pooling the statistical results of numerous studies to identify those results that are concordant and discordant with others. For example, one could obtain the MNI coordinates of significant clusters from several studies having to do with response inhibition and plot them in the same stereotaxic space to determine their concordance. One popular method of performing such an analysis is the creation of an Activation Likelihood Estimate, or ALE (Eickhoff et al., 2009; Turkeltaub et al., 2002). This method allows for the statistical thresholding of meta-analysis results, making it a powerful tool to examine the findings of many studies at once. Another method, mega-analysis, refers to reprocessing the raw data from numerous studies in a new statistical analysis with much greater power. Using this approach any systematic error introduced by any one study will contribute far less to the final statistical result (Costafreda, in press). Mega-analyses are far more difficult to implement since the raw imaging data from multiple studies must be obtained and reprocessed. Still, the increase in detection power and the greater generalizability of the results are strong reasons to engage in such an approach.

One roadblock to collaborative multi-center studies is the lack of data provenance in functional neuroimaging. Provenance refers to complete detail regarding the origin of a dataset and the history of operations that have been preformed on the data. Having a complete history of

the data enables analysis by other researchers and provides information that is critical for replication studies (Mackenzie-Graham et al., 2008). Moving forward there will be an additional focus on provenance to enable increased understanding of individual studies and facilitate integration into larger analyses.

<u>New Emphasis on Replication</u>. The non-independence debate of 2009 was less about effect sizes and more about reproducibility. The implicit argument made about studies that were 'non-independent' was that if researchers ran a non-independent study over again the resulting correlation would be far lower with a new, independent dataset. There should be a greater emphasis on the replicability of studies in the future. This can be frustrating because it is expensive and time consuming to acquire and process a replication study. However, moving forward this may become increasingly important to validate important results and conclusions.

General Conclusions

One thing is abundantly clear: fMRI is an effective research tool that has opened broad new horizons of investigation to scientists around the world. However, the results from fMRI research may be somewhat less reliable than many researchers implicitly believe. While it may be frustrating to know that fMRI results are not perfectly replicable, it is beneficial to take a longer-term view regarding the scientific impact of these studies. In neuroimaging, as in other scientific fields, errors will be made and some results will not replicate. Still, over time some measure of truth will accrue. This chapter is not intended to be an accusation against fMRI as a method. Quite the contrary, it is meant to increase the understanding of how much each fMRI result can contribute to scientific knowledge. If only 30% of the significant voxels in a cluster will replicate then that value represents an important piece of contextual information to be aware of. Likewise, if the magnitude of a voxel is only reliable at a level of ICC = 0.50 then that value represents important information when examining scatter plots comparing estimates of activity against a behavioral measure.

There are a variety of methods that can be used to evaluate reliability, and each can provide information on unique aspects of the results. Our findings speak strongly to the question of why there is no agreed-upon average value for fMRI reliability. There are so many factors spread out across so many levels of influence that it is almost impossible to summarize the reliability of fMRI with a single value. While our average ICC value of 0.50 and our average overlap value of

30% are effective summaries of fMRI as a whole, these values may be higher or lower on a study-to-study basis. The best characterization of fMRI reliability would be to give a window within which fMRI results are typically reliable. Breaking up the range of 0.0 to 1.0 into thirds, it is appropriate to say that most fMRI results are reliable in the ICC = 0.33 to 0.66 range.

To conclude, functional neuroimaging with fMRI is no longer in its infancy. Instead it has reached a point of adolescence, where knowledge and methods have made enormous progress but there is still much development left to be done. Our growing pains from this point forward are going to be a more complete understanding of its strengths, weaknesses, and limitations. A working knowledge of fMRI reliability is key to this understanding. The reliability of fMRI may not be the high relative to other scientific measures, but it is presently the best tool available for the in vivo investigation of brain function.

References

Andersson, J.L., Hutton, C., Ashburner, J., Turner, R., Friston, K., 2001. Modeling geometric deformations in EPI time series. Neuroimage 13, 903-919.

Aron, A.R., Gluck, M.A., Poldrack, R.A., 2006. Long-term test-retest reliability of functional MRI in a classification learning task. Neuroimage 29, 1000-1006.

Bandettini, P.A., Wong, E.C., Jesmanowicz, A., Hinks, R.S., Hyde, J.S., 1994. Spin-echo and gradient-echo EPI of human brain activation using BOLD contrast: a comparative study at 1.5 T. NMR Biomed 7, 12-20.

Bartko, J., 1966. The intraclass correlation coefficient as a measure of reliability. Psychological Reports 19, 3-11.

Bennett, C.M., Guerin, S.A., Miller, M.B., 2009. The impact of experimental design on the detection of individual variability in fMRI. Cognitive Neuroscience Society, San Francisco, CA.

Bennett, C.M., Wolford, G.L., Miller, M.B., in press. The principled control of false positives in neuroimaging. Social Cognitive and Affective Neuroscience.

Bodurka, J., Ye, F., Petridou, N., Bandettini, P.A., 2005. Determination of the brain tissuespecific temporal signal to noise limit of 3 T BOLD-weighted time course data., Proc. Intl. Soc. Mag. reson. Med., Miami.

Bosnell, R., Wegner, C., Kincses, Z.T., Korteweg, T., Agosta, F., Ciccarelli, O., De Stefano, N., Gass, A., Hirsch, J., Johansen-Berg, H., Kappos, L., Barkhof, F., Mancini, L., Manfredonia, F., Marino, S., Miller, D.H., Montalban, X., Palace, J., Rocca, M., Enzinger, C., Ropele, S., Rovira, A., Smith, S., Thompson, A., Thornton, J., Yousry, T., Whitcher, B., Filippi, M., Matthews, P.M., 2008. Reproducibility of fMRI in the clinical setting: implications for trial designs. Neuroimage 42, 603-610.

Caceres, A., Hall, D.L., Zelaya, F.O., Williams, S.C., Mehta, M.A., 2009. Measuring fMRI reliability with the intra-class correlation coefficient. Neuroimage 45, 758-768.

Carrier, J., Monk, T.H., 2000. Circadian rhythms of performance: new trends. Chronobiol Int 17, 719-732.

Casey, B.J., Cohen, J.D., O'Craven, K., Davidson, R.J., Irwin, W., Nelson, C.A., Noll, D.C., Hu, X., Lowe, M.J., Rosen, B.R., Truwitt, C.L., Turski, P.A., 1998. Reproducibility of fMRI results across four institutions using a spatial working memory task. Neuroimage 8, 249-261.

Chen, E.E., Small, S.L., 2007. Test-retest reliability in fMRI of language: group and task effects. Brain Lang 102, 176-185.

Chiarelli, P.A., Bulte, D.P., Wise, R., Gallichan, D., Jezzard, P., 2007. A calibration method for quantitative BOLD fMRI based on hyperoxia. Neuroimage 37, 808-820.

Cicchetti, D., Sparrow, S., 1981. Developing criteria for establishing interrater reliability of specific items: Applications to assessment of adaptive behavior. Am J Ment Defic 86, 127-137.

Clement, F., Belleville, S., 2009. Test-retest reliability of fMRI verbal episodic memory paradigms in healthy older adults and in persons with mild cognitive impairment. Hum Brain Mapp.

Cohen, J., 1977. Statistical power analysis for the behavioral sciences., (revised edition) ed. Academic Press, New York, NY.

Cohen, M.S., DuBois, R.M., 1999. Stability, repeatability, and the expression of signal magnitude in functional magnetic resonance imaging. J Magn Reson Imaging 10, 33-40.

Costafreda, S.G., in press. Pooling fMRI data: meta-analysis, mega-analysis and multi-center studies. . Frontiers in Neuroinformatics.

Costafreda, S.G., Brammer, M.J., Vencio, R.Z., Mourao, M.L., Portela, L.A., de Castro, C.C., Giampietro, V.P., Amaro, E., Jr., 2007. Multisite fMRI reproducibility of a motor task using identical MR systems. J Magn Reson Imaging 26, 1122-1126.

Dale, A., 1999. Optimal Experimental Design for Event-Related fMRI. Human Brain Mapping 8, 109-114.

Di Bonaventura, C., Vaudano, A.E., Carni, M., Pantano, P., Nucciarelli, V., Garreffa, G., Maraviglia, B., Prencipe, M., Bozzao, L., Manfredi, M., Giallonardo, A.T., 2005. Long-term reproducibility of fMRI activation in epilepsy patients with Fixation Off Sensitivity. Epilepsia 46, 1149-1151.

Duncan, K.J., Pattamadilok, C., Knierim, I., Devlin, J.T., 2009. Consistency and variability in functional localisers. Neuroimage 46, 1018-1026.

Eaton, K.P., Szaflarski, J.P., Altaye, M., Ball, A.L., Kissela, B.M., Banks, C., Holland, S.K., 2008. Reliability of fMRI for studies of language in post-stroke aphasia subjects. Neuroimage 41, 311-322.

Eickhoff, S.B., Laird, A.R., Grefkes, C., Wang, L.E., Zilles, K., Fox, P.T., 2009. Coordinatebased activation likelihood estimation meta-analysis of neuroimaging data: a random-effects approach based on empirical estimates of spatial uncertainty. Hum Brain Mapp 30, 2907-2926.

Feredoes, E., Postle, B.R., 2007. Localization of load sensitivity of working memory storage: quantitatively and qualitatively discrepant results yielded by single-subject and group-averaged approaches to fMRI group analysis. Neuroimage 35, 881-903.

Fernandez, G., Specht, K., Weis, S., Tendolkar, I., Reuber, M., Fell, J., Klaver, P., Ruhlmann, J., Reul, J., Elger, C.E., 2003. Intrasubject reproducibility of presurgical language lateralization and mapping using fMRI. Neurology 60, 969-975.

Freedman, R., 2003. Schizophrenia. N Engl J Med 349, 1738-1749.

Freyer, T., Valerius, G., Kuelz, A.K., Speck, O., Glauche, V., Hull, M., Voderholzer, U., 2009. Test-retest reliability of event-related functional MRI in a probabilistic reversal learning task. Psychiatry Res.

Friedman, L., Glover, G.H., 2006. Reducing interscanner variability of activation in a multicenter fMRI study: controlling for signal-to-fluctuation-noise-ratio (SFNR) differences. Neuroimage 33, 471-481.

Friedman, L., Stern, H., Brown, G.G., Mathalon, D.H., Turner, J., Glover, G.H., Gollub, R.L., Lauriello, J., Lim, K.O., Cannon, T., Greve, D.N., Bockholt, H.J., Belger, A., Mueller, B., Doty, M.J., He, J., Wells, W., Smyth, P., Pieper, S., Kim, S., Kubicki, M., Vangel, M., Potkin, S.G., 2008. Test-retest and between-site reliability in a multicenter fMRI study. Hum Brain Mapp 29, 958-972.

Gold, S., Christian, B., Arndt, S., Zeien, G., Cizadlo, T., Johnson, D.L., Flaum, M., Andreasen, N.C., 1998. Functional MRI statistical software packages: a comparative analysis. Hum Brain Mapp 6, 73-84.

Gountouna, V.E., Job, D.E., McIntosh, A.M., Moorhead, T.W., Lymer, G.K., Whalley, H.C., Hall, J., Waiter, G.D., Brennan, D., McGonigle, D.J., Ahearn, T.S., Cavanagh, J., Condon, B., Hadley, D.M., Marshall, I., Murray, A.D., Steele, J.D., Wardlaw, J.M., Lawrie, S.M., 2009. Functional Magnetic Resonance Imaging (fMRI) reproducibility and variance components across visits and scanning sites with a finger tapping task. Neuroimage.

Grafton, S., Hazeltine, E., Ivry, R., 1995. Functional mapping of sequence learning in normal humans. Journal of Cognitive Neuroscience 7, 497-510.

Harrington, G.S., Buonocore, M.H., Farias, S.T., 2006a. Intrasubject reproducibility of functional MR imaging activation in language tasks. AJNR Am J Neuroradiol 27, 938-944.

Harrington, G.S., Tomaszewski Farias, S., Buonocore, M.H., Yonelinas, A.P., 2006b. The intersubject and intrasubject reproducibility of FMRI activation during three encoding tasks: implications for clinical applications. Neuroradiology 48, 495-505.

Havel, P., Braun, B., Rau, S., Tonn, J.C., Fesl, G., Bruckmann, H., Ilmberger, J., 2006. Reproducibility of activation in four motor paradigms. An fMRI study. J Neurol 253, 471-476.

Hoenig, K., Kuhl, C.K., Scheef, L., 2005. Functional 3.0-T MR assessment of higher cognitive function: are there advantages over 1.5-T imaging? Radiology 234, 860-868.

Huang, J., Katsuura, T., Shimomura, Y., Iwanaga, K., 2006. Diurnal changes of ERP response to sound stimuli of varying frequency in morning-type and evening-type subjects. J Physiol Anthropol 25, 49-54.

Huettel, S.A., Song, A.W., McCarthy, G., 2008. Functional Magnetic Resonance Imaging, 2nd ed. Sinauer Associates, Sunderland, MA.

Jabbi, M., Keysers, C., Singer, T., Stephan, K.E., 2009. Response to "Voodoo Correlations in Social Neuroscience" by Vul et al.

Jansen, A., Menke, R., Sommer, J., Forster, A.F., Bruchmann, S., Hempleman, J., Weber, B., Knecht, S., 2006. The assessment of hemispheric lateralization in functional MRI--robustness and reproducibility. Neuroimage 33, 204-217.

Jezzard, P., Clare, S., 1999. Sources of distortion in functional MRI data. Hum Brain Mapp 8, 80-85.

Johnstone, T., Somerville, L.H., Alexander, A.L., Oakes, T.R., Davidson, R.J., Kalin, N.H., Whalen, P.J., 2005. Stability of amygdala BOLD response to fearful faces over multiple scan sessions. Neuroimage 25, 1112-1123.

Jovicich, J., Czanner, S., Greve, D., Haley, E., van der Kouwe, A., Gollub, R., Kennedy, D., Schmitt, F., Brown, G., Macfall, J., Fischl, B., Dale, A., 2006. Reliability in multi-site structural MRI studies: effects of gradient non-linearity correction on phantom and human data. Neuroimage 30, 436-443.

Kiebel, S., Holmes, A., 2007. The general linear model. In: Friston, K., Ashburner, J., Kiebel, S., Nichols, T., Penny, W. (Eds.), Statistical Parametric Mapping: The Analysis of Functional Brain Images. Academic Press, London.

Kiehl, K.A., Liddle, P.F., 2003. Reproducibility of the hemodynamic response to auditory oddball stimuli: a six-week test-retest study. Hum Brain Mapp 18, 42-52.

Kimberley, T.J., Khandekar, G., Borich, M., 2008. fMRI reliability in subjects with stroke. Exp Brain Res 186, 183-190.

Kong, J., Gollub, R.L., Webb, J.M., Kong, J.T., Vangel, M.G., Kwong, K., 2007. Test-retest study of fMRI signal change evoked by electroacupuncture stimulation. Neuroimage 34, 1171-1181.

Kruger, G., Glover, G.H., 2001. Physiological noise in oxygenation-sensitive magnetic resonance imaging. Magn Reson Med 46, 631-637.

Leontiev, O., Buxton, R.B., 2007. Reproducibility of BOLD, perfusion, and CMRO2 measurements with calibrated-BOLD fMRI. Neuroimage 35, 175-184.

Lieberman, M.D., Berkman, E.T., Wager, T.D., 2009. Correlations in social neuroscience aren't voodoo: Commentary on Vul et al. (2009). Perspectives on Psychological Science 4.

Liou, M., Su, H.R., Savostyanov, A.N., Lee, J.D., Aston, J.A., Chuang, C.H., Cheng, P.E., 2009. Beyond p-values: averaged and reproducible evidence in fMRI experiments. Psychophysiology 46, 367-378.

Liu, J.Z., Zhang, L., Brown, R.W., Yue, G.H., 2004. Reproducibility of fMRI at 1.5 T in a strictly controlled motor task. Magn Reson Med 52, 751-760.

Loubinoux, I., Carel, C., Alary, F., Boulanouar, K., Viallard, G., Manelfe, C., Rascol, O., Celsis, P., Chollet, F., 2001. Within-session and between-session reproducibility of cerebral sensorimotor activation: a test--retest effect evidenced with functional magnetic resonance imaging. J Cereb Blood Flow Metab 21, 592-607.

MacDonald, S.W., Nyberg, L., Backman, L., 2006. Intra-individual variability in behavior: links to brain structure, neurotransmission and neuronal activity. Trends Neurosci 29, 474-480.

Machielsen, W.C., Rombouts, S.A., Barkhof, F., Scheltens, P., Witter, M.P., 2000. FMRI of visual encoding: reproducibility of activation. Hum Brain Mapp 9, 156-164.

Mackenzie-Graham, A.J., Van Horn, J.D., Woods, R.P., Crawford, K.L., Toga, A.W., 2008. Provenance in neuroimaging. Neuroimage 42, 178-195.

Magon, S., Basso, G., Farace, P., Ricciardi, G.K., Beltramello, A., Sbarbati, A., 2009. Reproducibility of BOLD signal change induced by breath holding. Neuroimage 45, 702-712.

Maitra, R., 2009. Assessing certainty of activation or inactivation in test-retest fMRI studies. Neuroimage 47, 88-97.

Maitra, R., Roys, S.R., Gullapalli, R.P., 2002. Test-retest reliability estimation of functional MRI data. Magn Reson Med 48, 62-70.

Maldjian, J.A., Laurienti, P.J., Driskill, L., Burdette, J.H., 2002. Multiple reproducibility indices for evaluation of cognitive functional MR imaging paradigms. AJNR Am J Neuroradiol 23, 1030-1037.

Manoach, D.S., Halpern, E.F., Kramer, T.S., Chang, Y., Goff, D.C., Rauch, S.L., Kennedy, D.N., Gollub, R.L., 2001. Test-retest reliability of a functional MRI working memory paradigm in normal and schizophrenic subjects. Am J Psychiatry 158, 955-958.

Marshall, I., Simonotto, E., Deary, I.J., Maclullich, A., Ebmeier, K.P., Rose, E.J., Wardlaw, J.M., Goddard, N., Chappell, F.M., 2004. Repeatability of motor and working-memory tasks in healthy older volunteers: assessment at functional MR imaging. Radiology 233, 868-877.

Mayer, A.R., Xu, J., Pare-Blagoev, J., Posse, S., 2006. Reproducibility of activation in Broca's area during covert generation of single words at high field: a single trial FMRI study at 4 T. Neuroimage 32, 129-137.

McGonigle, D.J., Howseman, A.M., Athwal, B.S., Friston, K.J., Frackowiak, R.S., Holmes, A.P., 2000. Variability in fMRI: an examination of intersession differences. Neuroimage 11, 708-734.

Meindl, T., Teipel, S., Elmouden, R., Mueller, S., Koch, W., Dietrich, O., Coates, U., Reiser, M., Glaser, C., 2009. Test-retest reproducibility of the default-mode network in healthy individuals. Hum Brain Mapp.

Miki, A., Liu, G.T., Englander, S.A., Raz, J., van Erp, T.G., Modestino, E.J., Liu, C.J., Haselgrove, J.C., 2001. Reproducibility of visual activation during checkerboard stimulation in functional magnetic resonance imaging at 4 Tesla. Jpn J Ophthalmol 45, 151-155.

Miki, A., Raz, J., van Erp, T.G., Liu, C.S., Haselgrove, J.C., Liu, G.T., 2000. Reproducibility of visual activation in functional MR imaging and effects of postprocessing. AJNR Am J Neuroradiol 21, 910-915.

Mikl, M., Marecek, R., Hlustik, P., Pavlicova, M., Drastich, A., Chlebus, P., Brazdil, M., Krupa, P., 2008. Effects of spatial smoothing on fMRI group inferences. Magn Reson Imaging 26, 490-503.

Miller, M.B., Donovan, C.L., Van Horn, J.D., German, E., Sokol-Hessner, P., Wolford, G.L., 2009. Unique and persistent individual patterns of brain activity across different memory retrieval tasks. Neuroimage 48, 625-635.

Miller, M.B., Handy, T.C., Cutler, J., Inati, S., Wolford, G.L., 2001. Brain activations associated with shifts in response criterion on a recognition test. Can J Exp Psychol 55, 162-173.

Miller, M.B., Van Horn, J.D., Wolford, G.L., Handy, T.C., Valsangkar-Smyth, M., Inati, S., Grafton, S., Gazzaniga, M.S., 2002. Extensive individual differences in brain activations associated with episodic retrieval are reliable over time. J Cogn Neurosci 14, 1200-1214.

Morgan, V.L., Dawant, B.M., Li, Y., Pickens, D.R., 2007. Comparison of fMRI statistical software packages and strategies for analysis of images containing random and stimulus-correlated motion. Comput Med Imaging Graph 31, 436-446.

Morrison, P.D., Murray, R.M., 2005. Schizophrenia. Curr Biol 15, R980-984.

Moser, E., Teichtmeister, C., Diemling, M., 1996. Reproducibility and postprocessing of gradient-echo functional MRI to improve localization of brain activity in the human visual cortex. Magn Reson Imaging 14, 567-579.

Muller, R., Buttner, P., 1994. A critical discussion of intraclass correlation coefficients. Stat Med 13, 2465-2476.

Mumford, J.A., Nichols, T., 2009. Simple group fMRI modeling and inference. Neuroimage 47, 1469-1475.

Mumford, J.A., Nichols, T.E., 2008. Power calculation for group fMRI studies accounting for arbitrary design and temporal autocorrelation. Neuroimage 39, 261-268.

Mumford, J.A., Poldrack, R.A., Nichols, T., 2007. FMRIpower: A Power Calculation Tool for 2-Stage fMRI models. Human Brain Mapping, Chicago, IL.

Munneke, J., Heslenfeld, D.J., Theeuwes, J., 2008. Directing attention to a location in space results in retinotopic activation in primary visual cortex. Brain Res 1222, 184-191.

Murphy, K., Bodurka, J., Bandettini, P.A., 2007. How long to scan? The relationship between fMRI temporal signal to noise ratio and necessary scan duration. Neuroimage 34, 565-574.

Neumann, J., Lohmann, G., Zysset, S., von Cramon, D.Y., 2003. Within-subject variability of BOLD response dynamics. Neuroimage 19, 784-796.

Nichols, T., Brett, M., Andersson, J., Wager, T., Poline, J.B., 2005. Valid conjunction inference with the minimum statistic. Neuroimage 25, 653-660.

Nunnally, J., 1970. Introduction to psychological measurement. McGraw Hill, New York.

Oakes, T.R., Johnstone, T., Ores Walsh, K.S., Greischar, L.L., Alexander, A.L., Fox, A.S., Davidson, R.J., 2005. Comparison of fMRI motion correction software tools. Neuroimage 28, 529-543.

Ogawa, S., Menon, R.S., Tank, D.W., Kim, S.G., Merkle, H., Ellermann, J.M., Ugurbil, K., 1993. Functional brain mapping by blood oxygenation level-dependent contrast magnetic resonance imaging. A comparison of signal characteristics with a biophysical model. Biophys J 64, 803-812.

Peelen, M.V., Downing, P.E., 2005. Within-subject reproducibility of category-specific visual activation with functional MRI. Hum Brain Mapp 25, 402-408.

Peyron, R., Garcia-Larrea, L., Gregoire, M.C., Costes, N., Convers, P., Lavenne, F., Mauguiere, F., Michel, D., Laurent, B., 1999. Haemodynamic brain responses to acute pain in humans: sensory and attentional networks. Brain 122 (Pt 9), 1765-1780.

Phan, K.L., Liberzon, I., Welsh, R.C., Britton, J.C., Taylor, S.F., 2003. Habituation of rostral anterior cingulate cortex to repeated emotionally salient pictures. Neuropsychopharmacology 28, 1344-1350.

Poldrack, R.A., Prabhakaran, V., Seger, C.A., Gabrieli, J.D., 1999. Striatal activation during acquisition of a cognitive skill. Neuropsychology 13, 564-574.

Poline, J.B., Strother, S.C., Dehaene-Lambertz, G., Egan, G.F., Lancaster, J.L., 2006. Motivation and synthesis of the FIAC experiment: Reproducibility of fMRI results across expert analyses. Hum Brain Mapp 27, 351-359.

Raemaekers, M., Vink, M., Zandbelt, B., van Wezel, R.J., Kahn, R.S., Ramsey, N.F., 2007. Testretest reliability of fMRI activation during prosaccades and antisaccades. Neuroimage 36, 532-542.

Ramsey, N., Tallent, K., van Gelderen, P., Frank, J., Moonen, C., Weinberger, D., 1996. Reproducibility of Human 3D fMRI Brain Maps Acquired During a Motor Task. Human Brain Mapping 4, 113-121.

Rau, S., Fesl, G., Bruhns, P., Havel, P., Braun, B., Tonn, J.C., Ilmberger, J., 2007. Reproducibility of activations in Broca area with two language tasks: a functional MR imaging study. AJNR Am J Neuroradiol 28, 1346-1353.

Rombouts, S.A., Barkhof, F., Hoogenraad, F.G., Sprenger, M., Scheltens, P., 1998. Withinsubject reproducibility of visual activation patterns with functional magnetic resonance imaging using multislice echo planar imaging. Magn Reson Imaging 16, 105-113. Rombouts, S.A., Barkhof, F., Hoogenraad, F.G., Sprenger, M., Valk, J., Scheltens, P., 1997. Test-retest analysis with functional MR of the activated area in the human visual cortex. AJNR Am J Neuroradiol 18, 1317-1322.

Rostami, M., Hosseini, S.M., Takahashi, M., Sugiura, M., Kawashima, R., 2009. Neural bases of goal-directed implicit learning. Neuroimage 48, 303-310.

Rutten, G.J., Ramsey, N.F., van Rijen, P.C., van Veelen, C.W., 2002. Reproducibility of fMRIdetermined language lateralization in individual subjects. Brain Lang 80, 421-437.

Safrit, M., 1976. Reliability theory. American Alliance for Health, Physical Education, and Recreation, Washington, DC.

Salli, E., Korvenoja, A., Visa, A., Katila, T., Aronen, H.J., 2001. Reproducibility of fMRI: effect of the use of contextual information. Neuroimage 13, 459-471.

Salthouse, T.A., Nesselroade, J.R., Berish, D.E., 2006. Short-term variability in cognitive performance and the calibration of longitudinal change. J Gerontol B Psychol Sci Soc Sci 61, P144-151.

Schunck, T., Erb, G., Mathis, A., Jacob, N., Gilles, C., Namer, I.J., Meier, D., Luthringer, R., 2008. Test-retest reliability of a functional MRI anticipatory anxiety paradigm in healthy volunteers. J Magn Reson Imaging 27, 459-468.

Shehzad, Z., Kelly, A.M., Reiss, P.T., Gee, D.G., Gotimer, K., Uddin, L.Q., Lee, S.H., Margulies, D.S., Roy, A.K., Biswal, B.B., Petkova, E., Castellanos, F.X., Milham, M.P., 2009. The resting brain: unconstrained yet reliable. Cereb Cortex 19, 2209-2229.

Shrout, P., Fleiss, J., 1979. Intraclass Correlations: Uses in Assessing Rater Reliability. Psychological Bulletin 86, 420-428.

Simmons, W.K., Reddish, M., Bellgowan, P.S., Martin, A., 2009. The Selectivity and Functional Connectivity of the Anterior Temporal Lobes. Cereb Cortex.

Smith, A.T., Singh, K.D., Balsters, J.H., 2007. A comment on the severity of the effects of non-white noise in fMRI time-series. Neuroimage 36, 282-288.

Smith, S.M., Beckmann, C.F., Ramnani, N., Woolrich, M.W., Bannister, P.R., Jenkinson, M., Matthews, P.M., McGonigle, D.J., 2005. Variability in fMRI: a re-examination of inter-session differences. Hum Brain Mapp 24, 248-257.

Specht, K., Willmes, K., Shah, N.J., Jancke, L., 2003. Assessment of reliability in functional imaging studies. J Magn Reson Imaging 17, 463-471.

Stark, R., Schienle, A., Walter, B., Kirsch, P., Blecker, C., Ott, U., Schafer, A., Sammer, G., Zimmermann, M., Vaitl, D., 2004. Hemodynamic effects of negative emotional pictures - a test-retest analysis. Neuropsychobiology 50, 108-118.

Sterr, A., Shen, S., Zaman, A., Roberts, N., Szameitat, A., 2007. Activation of SI is modulated by attention: a random effects fMRI study using mechanical stimuli. Neuroreport 18, 607-611.

Strother, S., La Conte, S., Kai Hansen, L., Anderson, J., Zhang, J., Pulapura, S., Rottenberg, D., 2004. Optimizing the fMRI data-processing pipeline using prediction and reproducibility performance metrics: I. A preliminary group analysis. Neuroimage 23 Suppl 1, S196-207.

Strother, S.C., Anderson, J., Hansen, L.K., Kjems, U., Kustra, R., Sidtis, J., Frutiger, S., Muley, S., LaConte, S., Rottenberg, D., 2002. The quantitative evaluation of functional neuroimaging experiments: the NPAIRS data analysis framework. Neuroimage 15, 747-771.

Swallow, K.M., Braver, T.S., Snyder, A.Z., Speer, N.K., Zacks, J.M., 2003. Reliability of functional localization using fMRI. Neuroimage 20, 1561-1577.

Symms, M.R., Allen, P.J., Woermann, F.G., Polizzi, G., Krakow, K., Barker, G.J., Fish, D.R., Duncan, J.S., 1999. Reproducible localization of interictal epileptiform discharges using EEG-triggered fMRI. Phys Med Biol 44, N161-168.

Tegeler, C., Strother, S.C., Anderson, J.R., Kim, S.G., 1999. Reproducibility of BOLD-based functional MRI obtained at 4 T. Hum Brain Mapp 7, 267-283.

Thomason, M.E., Foland, L.C., Glover, G.H., 2007. Calibration of BOLD fMRI using breath holding reduces group variance during a cognitive task. Hum Brain Mapp 28, 59-68.

Triantafyllou, C., Hoge, R.D., Krueger, G., Wiggins, C.J., Potthast, A., Wiggins, G.C., Wald, L.L., 2005. Comparison of physiological noise at 1.5 T, 3 T and 7 T and optimization of fMRI acquisition parameters. Neuroimage 26, 243-250.

Turkeltaub, P.E., Guinevere, F.E., Jones, K.M., Zeffiro, T.A., 2002. Meta-Analysis of the Functional Neuroanatomy of Single-Word Reading: Method and Validation. Neuroimage 16, 765-780.

Turner, R., Jezzard, P., Wen, H., Kwong, K.K., Le Bihan, D., Zeffiro, T., Balaban, R.S., 1993. Functional mapping of the human visual cortex at 4 and 1.5 tesla using deoxygenation contrast EPI. Magn Reson Med 29, 277-279.

Van Horn, J.D., Ellmore, T.M., Esposito, G., Berman, K.F., 1998. Mapping voxel-based statistical power on parametric images. Neuroimage 7, 97-107.

Van Horn, J.D., Toga, A.W., 2009. Multisite neuroimaging trials. Curr Opin Neurol 22, 370-378.

Vul, E., Harris, C., Winkielman, P., Pashler, H., 2009. Puzzlingly high correlations in fMRI studies of emotion, personality, and social cognition. Perspectives on Psychological Science 4.

Wager, T.D., Nichols, T., 2003. Optimization of experimental design in fMRI: a general framework using a genetic algorithm. Neuroimage 18, 293-309.

Wagner, K., Frings, L., Quiske, A., Unterrainer, J., Schwarzwald, R., Spreer, J., Halsband, U., Schulze-Bonhage, A., 2005. The reliability of fMRI activations in the medial temporal lobes in a verbal episodic memory task. Neuroimage 28, 122-131.

Waites, A.B., Shaw, M.E., Briellmann, R.S., Labate, A., Abbott, D.F., Jackson, G.D., 2005. How reliable are fMRI-EEG studies of epilepsy? A nonparametric approach to analysis validation and optimization. Neuroimage 24, 192-199.

Waldvogel, D., van Gelderen, P., Immisch, I., Pfeiffer, C., Hallett, M., 2000. The variability of serial fMRI data: correlation between a visual and a motor task. Neuroreport 11, 3843-3847.

Wei, X., Yoo, S.S., Dickey, C.C., Zou, K.H., Guttmann, C.R., Panych, L.P., 2004. Functional MRI of auditory verbal working memory: long-term reproducibility analysis. Neuroimage 21, 1000-1008.

Whalley, H.C., Gountouna, V.E., Hall, J., McIntosh, A.M., Simonotto, E., Job, D.E., Owens, D.G., Johnstone, E.C., Lawrie, S.M., 2009. fMRI changes over time and reproducibility in unmedicated subjects at high genetic risk of schizophrenia. Psychol Med 39, 1189-1199.

White, T., O'Leary, D., Magnotta, V., Arndt, S., Flaum, M., Andreasen, N.C., 2001. Anatomic and functional variability: the effects of filter size in group fMRI data analysis. Neuroimage 13, 577-588.

Woolrich, M.W., Ripley, B.D., Brady, M., Smith, S.M., 2001. Temporal autocorrelation in univariate linear modeling of FMRI data. Neuroimage 14, 1370-1386.

Yetkin, F.Z., McAuliffe, T.L., Cox, R., Haughton, V.M., 1996. Test-retest precision of functional MR in sensory and motor task activation. AJNR Am J Neuroradiol 17, 95-98.

Yoo, S.S., O'Leary, H.M., Lee, J.H., Chen, N.K., Panych, L.P., Jolesz, F.A., 2007. Reproducibility of trial-based functional MRI on motor imagery. Int J Neurosci 117, 215-227.

Yoo, S.S., Wei, X., Dickey, C.C., Guttmann, C.R., Panych, L.P., 2005. Long-term reproducibility analysis of fMRI using hand motor task. Int J Neurosci 115, 55-77.

Zandbelt, B.B., Gladwin, T.E., Raemaekers, M., van Buuren, M., Neggers, S.F., Kahn, R.S., Ramsey, N.F., Vink, M., 2008. Within-subject variation in BOLD-fMRI signal changes across repeated measurements: quantification and implications for sample size. Neuroimage 42, 196-206.

Zhang, J., Anderson, J.R., Liang, L., Pulapura, S.K., Gatewood, L., Rottenberg, D.A., Strother, S.C., 2009. Evaluation and optimization of fMRI single-subject processing pipelines with NPAIRS and second-level CVA. Magn Reson Imaging 27, 264-278.

Zhang, J., Liang, L., Anderson, J.R., Gatewood, L., Rottenberg, D.A., Strother, S.C., 2008. A Java-based fMRI processing pipeline evaluation system for assessment of univariate general linear model and multivariate canonical variate analysis-based pipelines. Neuroinformatics 6, 123-134.

Zhilkin, P., Alexander, M.E., 2004. Affine registration: a comparison of several programs. Magn Reson Imaging 22, 55-66.

Zou, K.H., Greve, D.N., Wang, M., Pieper, S.D., Warfield, S.K., White, N.S., Manandhar, S., Brown, G.G., Vangel, M.G., Kikinis, R., Wells, W.M., 3rd, 2005. Reproducibility of functional MR imaging: preliminary results of prospective multi-institutional study performed by Biomedical Informatics Research Network. Radiology 237, 781-789.

Results of exam	ined pap	ers us	ing intra-class correlation as a re	liability me	etric.							
									Approximate	Min	Mean	Max
First Author	Year		Task	Design	Type	Basis	Contrast	# Subs	T-R Interval	ICC	ICC	ICC
ţ					м 	Contrast	Task vs	-	, ,			
Caceres ^{+/}	2009	ಕ	Auditory Target Detection	Block	Sig. Voxels	Values	Rest	10	3 Months	·	0.35	
Cacarac ⁴⁷	0000	8	M. Book Working Mamony	Direk	Cia Wovele	Values	Control	10	3 Monthe		070	
Caceles	6007		Probabilistic Reversal	DIUCK	DIE. VUXEIS	Contrast	Task vs	10	STITUTOTAL C	I	0.47	I
Freyer ⁹⁴	2009	ą	Learning	Event	All Voxels	Values	Control	10	16 Weeks			
			•			Contrast	Task vs					
Gountouna ⁹⁵	2009		Finger tapping	Block	ROI	Values	Rest	14	Unknown	0.23	0.53	0.72
			1			Percent Signal	Task vs					
Bosnell ⁷¹	2008		Hand Tapping	Block	ROI	Change	Rest	22	< 1 Day		0.82	
1						Percent Signal	Task vs					
Friedman ⁸	2008	a,c	Finger Tapping	Block	ROI	Change	Rest	5	1 Day	0.47	0.74	0.85
						Percent Signal	Task vs					
Schunck ⁹⁶	2008		Anticipatory Anxiety	Block	ROI	Change	Rest	14	10 Days	-0.06	0.34	0.66
!						Percent Signal	Task vs					
$Kong^{97}$	2007		Finger tapping	Block	ROI	Change	Rest	8	1 Week	0.00	0.37	0.76
1						Percent Signal	Task vs					
$Kong^{97}$	2007		Acupuncture	Block	ROI	Change	Rest	œ	1 Week	0.00	0.16	0.54
						t-Statistic	Task vs					
Raemaekers ⁵⁸	2007		Prosaccade/Antisaccade	Event	All Voxels	Values	Rest	12	1 Week	-0.08		0.79
:			Probabilistic classification			Contrast	Task vs					
Aron ⁵²	2006		learning	Event	ROI	Values	Rest	~	59 Weeks	0.76	0.88	0.99
:			Amygdala-Facial affect		Amygdala	Contrast	Task vs					
Johnstone ⁹⁸	2005		localizer	Block	ROI	Values	Rest	15	8 weeks	0.02	0.38	0.63
:						Activation	Task vs					
Wei ⁹⁹	2004		Auditory two-back	Block	ROI	Index	Rest	~	9 Weeks	0.14	0.43	0.71
						Percent Signal	Task vs					
Specht ¹⁰⁰	2003	q	Visual Attention	Event	Sig. Voxels	Change	Rest	S	8 Weeks			
:						Percent Signal	Task vs					
Manoach ⁶⁶	2001		Sternberg item recognition	Event	ROI	Change	Control	7	14 Weeks	0.23	0.52	0.81
Mean Value										0.17	0.50	0.75
^a - Median valu	e given											

 $^{\rm b}$ - Data presented as graphs or figures, unable to quantify values $^{\circ}$ - Data acquired from multiple scanners

FMRI RELIABILITY

Table 1

Table 2 Results of exam	nined nar	I SADO	usino cluster overlan as a re	liahility 1	netric								
the second s			C C									Dice Overlap	
First Author	Year		Task	Design	Calculation	Basis	Contrast	Threshold	# Subs	Approximate T-R Interval	Min Overlap	Avg Overlap	Max Overlap
=			One-back object/word				Task vs	p(uncorr) <					
Duncan ⁴¹	2009	75	localizer	Block	Dice	ROI Sig	Rest Task ve	0.001	45	<1 hour	0.380	0.435	0.490
Gountouna ⁹⁵	2009		Finger tapping	Block	Dice	Voxels	Rest	p(corr) < 0.05	14	Not Given	0.410	0.455	0.500
Meindl ¹⁰¹	2009		Default Mode	Нгее	Dice	τυν	Component	p(uncorr) < 0.05	18	1 Week	0.080	0300	0.760
	6007	þ,		221		Sig	Task vs	0.0	10		0000	0600	00/.0
Feredoes ¹⁰²	2007	e t	Button Press	Event	Custom	Voxels	Rest	p(corr) < 0.05	9	1 Week		0.245	
Feredoes ¹⁰²	2007	ວັຍ	Delayed Recognition	Event	Custom	Sig Voxels	Task vs Control	p(corr) < 0.05	9	1 Week	0.000	0.210	0.413
Raemaekers ⁵⁸	2007		Prosaccade/Antisaccade	Event	Dice	Sig Voxels	Task vs Rest	p(corr) < 0.005	12	1 Week	0.760	0.785	0.810
Raii ⁵⁹	2007		Naming / Noun Generation	Block	Dice	Sig Voxels	Task vs Rest	n(corr) < 0.05	<u>"</u>	9 Davs	0.000	0350	0.820
Harrington ¹⁰³	2006a	0	Multinle Language	Event	Dice	ROI	Task vs Rest	p(corr) < 0.05	10	4 Weeks			ı
104	12000	0	-ooooooooooooo-		, C		Task vs		-	1.0 10.01			
Harrington	70000		Multiple Encoding	Block	Dice	Sig	Tack ve	p(corr) < 0.05 (mcorr) <	h	IU Weeks			
Havel ¹⁰⁵	2006		Motor movement	Block	Dice	Voxels	Rest	0.001	15	6 Days	0.000	0.230	0.710
Wagner ¹⁰⁶	2005		Verbal Encoding	Block	Dice	Sig Voxels	Task vs Control	Individualized	20	33 Weeks		0.362	ı
Wagner ¹⁰⁶	2005		Verhal Recomition	Block	Dice	Sig	Task vs Control	Individualized	00	33 Weeks	I	0.420	I
1211311	2004			NACIO	2017		Task vs	p(uncorr) <	2			071-0	
Yoo^{107}	2005	9	Finger Tapping	Block	Dice	ROI	Rest	0.005	8	8 weeks	·	,	ı
Specht ¹⁰⁰	2003		Visual Attention	Event	Dice	Voxelwise	l ask vs Rest	p(uncorr) < 0.01	S	2 Weeks	0.420	0.583	0.692
901			Visual FEF and MT				Task vs	z(uncorr) >					
Swallow ¹⁰⁶	2003	•	localizers	Block	Jaccard	ROI	Rest	4.5	=	Not Given	0.416	0.463	0.507
Maldjian ⁵⁷	2002	q	Word Generation	Block	Jaccard	Sig Voxels	lask vs Rest	p(uncorr) < 0.005	8	1 Week	0.748	0.856	0.993
,			Forward-Backward			Sig	Task vs	p(uncorr) <					
Maldjian ⁵⁷	2002	٩	Listening	Block	Jaccard	Voxels	Rest	0.005	8	1 Week	0.410	0.662	0.817
Rutten ¹⁰⁹	2002	þ,d	Combined language	Block	Cuetom	Sig Vovele	Task vs Pact	z(uncorr) >	0	5 months	,	0.420	
011						Sig	Task vs	z(uncorr) >					
Miki	2001		Visual checkerboard	Block	Dice	Voxels c:~	Rest Tech we	4.5	4	<1 hour	0.560	0.610	0.660
Machielsen ¹¹¹	2000		Visual Encoding	Block	Dice	Voxels	Lask vs Control	p(corr) < 0.05	10	14 Days		0.507	

						Task vs						
Machielsen ¹¹¹	2000	Visual Encoding	Block	Dice	ROI	Control	p(corr) < 0.05	10	14 Days	0.211	0.374	0.514
					Sig	Task vs	z(uncorr) >					
Miki ¹¹²	2000	Visual light stimulation	Block	Dice	Voxels	Rest	4.5	7	5 Days	0.020	0.480	0.770
					Sig	Task vs	Top 2% of					
Tegeler ¹¹³	1999	Finger tapping	Block	Dice	Voxels	Rest	voxels	9	<1 Hour		0.410	
					Sig	Task vs						
Rombouts ⁴²	1998	Visual light stimulation	Block	Dice	Voxels	Rest	p(corr) < 0.05	10	2 Weeks	0.460	0.640	0.760
					Sig	Task vs	r(uncorr) >					
Rombouts ⁴⁰	1997	Visual light stimulation	Block	Dice	Voxels	Rest	0.50	14	2 weeks	0.150	0.310	0.500
					Sig	Task vs						
Ramsey ¹¹⁴	1996	^b Finger tapping	Block	Jaccard	Voxels	Rest	p(corr) < 0.05	7	11 weeks		0.333	
					Sig	Task vs	r(uncorr) >					
Yetkin ¹¹⁵	1996	^b Finger tapping	Block	Jaccard	Voxels	Rest	09.0	4	< 1 hour		0.742	
					Sig	Task vs	r(uncorr) >					
Yetkin ¹¹⁵	1996	^b Somatosensory Touch	Block	Jaccard	Voxels	Rest	0.60	4	< 1 hour		0.621	
Mean Value										0.314	0.476	0.670
^a - Overlap valı	tes estimé	ated from figure										

^b - Results recalculated to represent Dice statistic

 $^{\circ}$ - Data presented as graphs or figures, unable to quantify values

^d - Overlap only calculated for a single region

 $^{\rm e}$ - Calculated using total voxels in first session only, not average

								Approximate
First Author	Year		Task	Design	Method	Туре	# Subs	T-R Interval
Liou ¹¹⁶	2009		Multiple Tasks	Event	Cohen's Kappa Index	Voxelwise	12	< 1 Hour
Magon ¹¹⁷	2009		Breath Holding	Block	Variation Maximum	ROI	11	3 Weeks
Maitra ¹¹⁸	2009		Finger Tapping Multiple Memory	Block	Likelihood	Voxelwise	1	5 Days
Miller ⁵⁵	2009		Tasks	Block/Event	Pearson Correlation	Voxelwise	14	12 Weeks 1 Hour / 5
Shehzad ¹¹⁹	2009		Resting State	Free	Connectivity ICC	ROI	26	Months 1 Hour / 5
Shehzad ¹¹⁹	2009		Resting State	Free	Kendall's W	ROI	26	Months
Zhang ³¹	2009		Static Force	Block	NPAIRS	Components	16	< 1 Hour
Zandbelt ¹²⁰	2008		Go/NoGo	Block	Signal Change SD Reliability Mans /	ROI	10	1 Week
Chen ¹²¹	2007		Language Imaging Retinotopic	Block	ROC Coefficient of	Voxelwise	12	Variable
Leontiev ¹²²	2007		Mapping	Block	Variation	ROI	10	1 Day
Y00 ¹²³	2007		Motor Imagery	Block	Signal Change SD	ROI	10	< 1 Hour
Jansen ¹²⁴	2006		Multiple Tasks Covert Word	Block	Laterality Index	ROI	10	2 Hours
Mayer ¹²⁵	2006		Generation	Event	Active Voxel Count	ROI	8	<1 Hour
Peelen ¹²⁶	2005		Categorization Visual, Motor.	Block	Sign Test Variance	ROI	6	3 Weeks
Smith ¹²⁷	2005		Cognitive	Block	Components	ROI	1	1 Day
Wagner ¹⁰⁶	2005		Verbal Memory	Block	Pearson Correlation General Linear	All Voxels	20	33 Weeks
Liu ¹²⁸	2004		Handgrip Task Emotional	Block	Model	Voxelwise	8	1 Month
Stark ¹²⁹	2004		Pictures	Block	Kappa Index	Sig Voxels	24	1 Week
Strother ³⁰	2004		Static Force	Block	NPAIRS General Linear	Components	16	< 1 Hour
Phan ¹³⁰	2003		Aversive Images	Block	Model Conjunction	ROI	8	<1 Hour
Kiehl ¹³¹	2003		Auditory Oddball	Event	Analysis	Voxelwise	10	6 Weeks
Neumann ¹³²	2003		Stroop Task	Event	BOLD Dynamics Maximum	ROI	4	1 Week
Maitra ¹³³	2002		Finger Tapping	Block	Likelihood	All Voxels	1	~ 1 Week
Miller ³⁶	2002		Episodic Retrieval	Block	Pearson Correlation Coefficient of	All Voxels	6	6 Months
Loubinoux ¹³⁴	2001		Sensorimotor	Block	Variation	Sig Voxels	21	Variable
Salli ¹³⁵	2001		Wrist Flexing	Block	Reliability Maps	Voxelwise	1	< 1 Hour
White ¹³⁶	2001		Finger Tapping Visual Motor	Block	ROI Discrimination General Linear	ROI	6	3 Weeks
McGonigle ¹³⁷	2000		Cognitive	Block	Model Signal Change	Voxelwise	1	1 Day
Waldvogel ¹³⁸	2000		Checkerboard	Block	Stability	All Voxels	6	1 Week
Cohen ⁴⁴	1999	a	Visual and Motor	Block	Voxel Counting	Sig Voxels	6	< 1 Hour
Tegeler ¹¹³	1999		Finger tapping	Block	Pearson Correlation	All Voxels	6	<1 Hour
Moser ¹⁹	1996		Visual stimulation	Block	Signal Change SD	Single Slice	18	< 1 Hour

Table 3Results of examined papers using other forms of reliability calculation.

^a - Cohen et al. conducted the experiment to argue against voxel counting

Figure Captions

Figure 1. Visualization of cluster overlap using two runs of data from a two-back working memory task. The regions in red represent significant clusters from the first run and regions in blue represent significant clusters from the second run. The crosshatched region represents the overlapping voxels that were significant in both runs. Important to note is that not all significant voxels remained significant across the two runs. One cluster in the cerebellum did not replicate at all. Data is from Bennett, Guerin, and Miller (2009).

Figure 2. Histogram showing the frequency of voxelwise ICC values during a two-back working memory task. The histogram was computed from a dataset of sixteen subjects using 100 bins between ICC values of 1.0 and -1.0. The distribution of values is negatively skewed, with a mean ICC value of ICC = 0.44 and the most frequently occurring value of ICC = 0.57. Data is from Bennett, Guerin, and Miller (2009).



