

# Minimal Data Needed for Valid & Accurate Image-Based fMRI Meta-Analysis

Camille Maumet, Thomas Nichols

WMG, University of Warwick, Coventry, UK

Statistics Department, University of Warwick, Coventry, UK.

---

## Abstract

Meta-analysis provides a quantitative approach to summarise the rich functional Magnetic Resonance Imaging literature (fMRI). When image data is available for each study, a number of approaches have been proposed to perform such meta-analysis including combination of standardised statistics, just effect estimates or both effects estimates and their sampling variance. While the latter is the preferred approach in the statistical community, its properties are only guaranteed under large sample sizes. Additionally, often only standardised estimates are shared, reducing the possible meta-analytic approaches. Finally, because the BOLD signal is non-quantitative care has to be taken in order to insure that effect estimates are expressed in the same units, especially when combining data from different software packages. Given the growing interest in data sharing in the neuroimaging community there is a need to identify what is the minimal data to be shared in order to allow for future image-based meta-analysis. In this paper, we compare the validity and the accuracy of nine meta-analytic approaches on simulated and real data.

*Keywords:* Meta-analysis, Neuroimaging, Mixed-effects

---

## 1. Introduction

A growing literature is focusing on the lack of statistical power in neuroimaging studies (see, e.g. [2]), feeding the debate on the validity and reproducibility of published neuroimaging results. Meta-analysis, by providing inference based on the results of previously conducted studies, provides an essential method to increase power and hence confidence in neuroimaging.

A number of methods have been proposed for neuroimaging meta-analysis (see [13] for a review). As the results of neuroimaging studies are usually conveyed by providing a table of peak coordinate and statistics, most of these meta-analyses are restricted to combining coordinate-based information. Nevertheless the best practice method is an Intensity-Based Meta-Analysis (IBMA) that combines the effect estimates and their standard errors from each study [1].

In order for IBMA to be possible in neuroimaging, tools for sharing 3D volumes obtained as a result of a statistical analysis are needed. Various efforts are currently underway to facilitate sharing of neuroimaging data but emphasis is usually on statistical maps (see, e.g. [2]). There are three evident approaches to sharing summary data from each study  $i$ :

1. the contrast estimates  $\hat{\beta}_i$  and contrast variance estimates  $s_i^2$ .
2. the contrast estimates  $\hat{\beta}_i$ .
3. the standardized statistical maps  $Z_i$ .

Depending on how much data is shared, different strategies can be used to combine the available results into a meta-analysis. While the first option is the best practice, leading to statistically optimal estimates [5], it requires the contrasts to be expressed with in the same units and inference relies on asymptotic results (i.e under large sample sizes). In fMRI, units will depends on the field strength [4] as well as data, model and contrast vector scaling [11] and the number of samples included in a meta-analysis is usually small.

Given the growing interest in data sharing in the neuroimaging community, and the relative easiness of sharing just (unitless) statistic maps, there is a need to identify what is the minimal data to be shared in order to allow for future IBMA.

Here we compare the use of IMBA using 9 meta-analytic approaches: 2 approaches use  $\hat{\beta}_i$ 's and  $s_i^2$ 's, 2  $\hat{\beta}_i$ 's only and 5  $Z_i$ 's. We compare the validity and the accuracy of the nine meta-analytic approaches on simulated and real data including 21 studies of pain in control subjects.

Section 2 describes the meta-analytic estimates along with the experiments undertaken on simulated and real data to assert their validity. The results are described in section 3. Finally, we conclude in section 4.

Figure 1: False positive rates of the meta-analytic estimators under the null hypothesis for  $p < 0.05$ .

	$\hat{\gamma}$	$\text{Var}(\hat{\gamma})$	Assumptions
MFX GLM	$(\sum \kappa_i \hat{\beta}_i) / (\sum \kappa_i)$ with $\kappa_i = 1/(\hat{\tau}^2 + s_i^2)$	$1/\sum \kappa_i$	IGE.
RFX GLM	$\sum \hat{\beta}_i/k$	$\sigma_C^2/k$	IGE; $\tau^2 + \sigma_i^2 = \sigma_C^2 \forall i$
FFX GLM	$(\sum \hat{\beta}_i \times n_i/\sigma_i^2) / (\sum n_i/\sigma_i^2)$	$1/(\sum n_i/\sigma_i^2)$	IGE; $\tau^2 = 0$ .
Contrast Perm.	$\sum \hat{\beta}_i/k$	Empirical	ISE.
Z MFX	$\sum Z_i/k$	$\sigma_C^2/k$	IGE; $1 + \tau^2/\sigma_i^2$ cst.
Z Perm.	$(\sum_{i=1}^k Z_i) / \sqrt{k}$	Empirical	ISE.

Table 1: One-sample meta-analytic estimates, sampling variance and associated assumptions. Note:  $P_i = \Phi(-Z_i)$

## 2. Methods

### 2.1. Theory

For study  $i = 1, \dots, k$  we have contrast estimate  $\hat{\beta}_i$ , its contrast variance estimate  $s_i^2$  (i.e. squared standard error), its statistic map  $Z_i$  and its sample size  $n_i$ .

*Combining contrast estimates and their standard error.* The gold standard approach is to fit contrast estimates and their standard error with a hierarchical general linear model (GLM) [5], creating a third-level (level 1: subject; level 2: study; level 3: meta-analysis). The general formulation for the study-level data is:

$$\hat{\beta} = X\gamma + \epsilon \quad (1)$$

where  $\gamma$  is the meta-analytic parameter to estimate,  $\hat{\beta} = [\hat{\beta}_1 \dots \hat{\beta}_k]^T$  is the vector of contrast estimates,  $X$  is the  $k \times p$  study-level matrix (typically just a column of ones for a one-sample test) and  $\epsilon \sim \mathcal{N}(0, W)$  is the residual error term.

In the most general case of a random-effects (RFX) meta-analysis, we have  $W = \text{diag}(\sigma_1^2/n_1 + \tau^2, \dots, \sigma_k^2/n_k + \tau^2)$  where  $\tau^2$  denotes the between-study variance and  $\sigma_i^2/n_i$  denotes the contrast variance for study  $i$ . Eq. (1) can be solved by weighted least squares giving:

$$\hat{\gamma} = (X^T W^{-1} X)^{-1} X^T W^{-1} \hat{\beta} \quad (2)$$

$$\text{Var}(\hat{\gamma}) = (X^T W^{-1} X)^{-1} \quad (3)$$

But in practice, the weight matrix  $W$  is unknown and has to be estimated from the data. Given  $\hat{W}$  a consistent estimate of  $W$ , the feasible generalized least squares (FGLS) estimator is computed as:

$$\hat{\gamma} = (X^T \hat{W}^{-1} X)^{-1} X^T \hat{W}^{-1} \hat{\beta} \quad (4)$$

$$\text{Var}(\hat{\gamma}) = (X^T \hat{W}^{-1} X)^{-1} \quad (5)$$

Approximating  $\sigma_i^2/n_i$  by  $s_i^2$  and given  $\hat{\tau}^2$  an estimate of  $\tau^2$  we obtain the estimate detailed in Table 1 for a one-sample test. Asymptotic theory shows that inference can be carried out by comparing the statistic to a Student distribution with  $k - 1$  degrees of freedom + ADDREF as depicted in Table 2. This reference approach will be referred to as **Mixed-effects (MFX) GLM**.

In a **fixed-effects (FFX) GLM**, i.e. assuming no or negligible between-study variance,  $W = \text{diag}(\sigma_1^2/n_1, \dots, \sigma_k^2/n_k)$ . Approximating  $\sigma_i^2/n_i$  by  $s_i^2$  we obtain the feasible generalized least squares estimate detailed in Table 1 for a one-sample test. Asymptotic theory shows that inference can be carried out by comparing the statistic to a Student distribution with  $(\sum_i n_i * k) - 1$  degrees of freedom + ADDREF as depicted in Table 2.

*Combining contrast estimates.* If the  $s_i^2$  are unavailable, the contrast estimates  $\hat{\beta}_i$  can be combined by assuming that the within-study contrast variance  $\sigma_i^2/n_i$  is constant ( $\sigma_i^2/n_i = \sigma^2 \forall i$ ) or negligible in comparison to the between-study variance ( $\sigma_i^2/n_i \ll \tau^2$ ). Then  $W = \text{diag}(\sigma_C^2, \dots, \sigma_C^2)$  where  $\sigma_C^2$  is the combined within and between-subject variance, i.e.  $\sigma_C^2 \simeq \tau^2$  or  $\sigma_C^2 \simeq \tau^2 + \sigma^2$  (note, however, in this setting we do not separately estimate  $\tau^2$  or  $\sigma^2$ ). Under these assumptions, Eq. (1) can be solved by ordinary least squares giving:

	Meta-analytic statistic	Nominal $H_0$ distrib.	Inputs	Properties
MFX GLM	$\left(\sum \kappa_i \hat{\beta}_i\right) / \sqrt{\sum_{i=1}^k \kappa_i}$ with $\kappa_i = 1/(\hat{\tau}^2 + s_i^2)$	$\mathcal{T}_{k-1}$	$\hat{\beta}_i, s_i^2$	Asymptotic.
RFX GLM	$\left(\sum_{i=1}^k \frac{\hat{\beta}_i}{\sqrt{k}}\right) / \widehat{\sigma}_C^2$	$\mathcal{T}_{k-1}$	$\hat{\beta}_i$	Finite sample.
Ctrst Perm.	$\left(\sum_{i=1}^k \frac{\hat{\beta}_i}{\sqrt{k}}\right) / \widehat{\sigma}_C^2$	Empirical	$\hat{\beta}_i$	??.
FFX GLM	$\left(\sum_{i=1}^k \frac{\hat{\beta}_i}{s_i^2}\right) / \sqrt{\sum_{i=1}^k 1/s_i^2}$	$\mathcal{T}_{(\sum_{i=1}^k n_i - 1) - 1}$	$\hat{\beta}_i, s_i^2$	Asymptotic.
Fisher's	$-2 \sum_i \ln P_i$	$\chi^2_{(2k)}$	$Z_i$	??.
Stouffer's	$\sqrt{k} \times \frac{1}{k} \sum_i Z_i$	$\mathcal{N}(0, 1)$	$Z_i$	??.
Wght Stouff.	$\frac{1}{\sqrt{\sum_i n_i}} \sum_i \sqrt{n_i} Z_i$	$\mathcal{N}(0, 1)$	$Z_i, n_i$	??.
Z MFX	$\left(\sum_{i=1}^k Z_i\right) / \sqrt{k} \hat{\sigma}$	$\mathcal{T}_{k-1}$	$Z_i$	??.
Z Perm.	$\left(\sum_{i=1}^k Z_i\right) / \sqrt{k}$	Empirical	$Z_i$	??.

Table 2: Statistics for one-sample meta-analysis tests and their sampling distributions under the null hypothesis  $H_0$ . Empirical null distributions are determined using permutations with sign flipping. IGE=Independent Gaussian Errors, ISE=Independent Symmetric Errors. Note:  $P_i = \Phi(-Z_i)$ ,  $\widehat{\sigma}_C^2$  is the unbiased sample variance.

$$\hat{\gamma} = (X^T X)^{-1} X^T \hat{\beta} \quad (6)$$

$$\text{Var}(\hat{\gamma}) = (X^T X)^{-1} \sigma_C^2 \quad (7)$$

Given  $\widehat{\sigma}_C^2$  the unbiased sample variance, we obtain the statistics presented in Table 1 for one sample tests. This approach will be referred to **RFX GLM** in the following. Inference can be carried out by comparing the RFX GLM statistic to a Student distribution with  $k-1$  degrees of freedom, this result holds asymptotically as well as in small samples +*ADDERF*.

As an alternative to parametric approaches, non-parametric inference [7, 10] can be performed by comparing the RFX GLM T-statistic to the distribution obtained with “sign flipping”, i.e. randomly multiplying each study’s data by 1 or -1, justified by an assumption of independent studies and symmetrically distributed random error. This approach will be referred to as **Contrast permutation**.

*Combining standardised statistics.* When only test statistic images are available there are several alternative approaches available. **Fisher's** meta-analysis provide a statistic to combine the associated p-values [6]. **Stouffer's** approach combines directly the standardised statistic [16]. In [17] following [8], the author proposed a weighted method that weights each study’s  $Z_i$  by the square root of its sample size [3,7]. This approach will be referred to as **Weighted Stouffer's**. All these meta-analytic statistics assumes no or negligible between-study variance and are suited only for one-sample tests. The corresponding statistics are presented in Table 2. As suggested in [1], to get a kind of MFX with Stouffer’s approach, the standardised statistical estimates  $Z_i$  can be combined in an OLS analysis. The corresponding estimate, referred as **Z MFX** is also provided in 2

Non-parametric inference [7, 10] can also be obtained by sign flipping on the  $Z_i$ 's. This approach will be referred to as **Z permutation**.

*Approximations.* In practice, all of the methods based on contrast data have approximate parametric null distributions. The nominal distributions listed in Table 2 are under the (unrealistic) assumption of homogeneous standard errors over studies; even if all studies are ‘clean’ and conducted at the same center, variation in sample size will induce differences in  $s_i^2$ 's. Further, even under homoscedasticity, MFX GLM's null is approximate in small samples. Furthermore, all contrast methods require the contrasts to be expressed with in the same units. In fMRI, units will depends on the field strength [4] as well as data, model and contrast vector scaling [11].

## 2.2. Experiments

### 2.2.1. Simulations

We used Monte Carlo simulations to empirically investigate the validity of each estimator. We simulated a set of studies contrast estimates  $\hat{\beta}_i$  and contrast variance estimates  $s_i^2$  according to:

$$\hat{\beta}_i \sim \mathcal{N}(0, \frac{\sigma_i^2}{n} + \tau^2) \quad (8)$$

$$s_i^2 \sim \frac{\sigma_i^2}{n} \frac{\chi_{(n-1)}^2}{n-1} \quad (9)$$

where  $\sigma_i^2 = \sigma^2 \alpha_i$  with  $\sigma^2 \in n \times [0.25, 0.5, 1, 2, 4]$  and  $\alpha_i$  is either equal to 1 for all studies or is taken from  $\alpha_i \in [1, 2, 4, 8, 16]$  to simulate varying within-study variances,  $\tau^2 \in [0, 1]$  is the between-study variance. Four different number of studies per meta-analysis were used:  $k \in [5, 10, 25, 50]$ . We set the number of subjects per studies  $n = 20$  which is a common sample size in existing neuroimaging studies [12]. A total of 32 parameter sets ( $4 \sigma_i^2 \times 2 \tau^2 \times 4 k$ ) was therefore tested and a total of 1 026 000 realisations were performed for each parameter set.

Three types of analyses were computed: a one-sample meta-analysis (testing significance on the mean effect in 1 group of  $k$ ), a two-sample meta-analysis (testing significance in mean differences between two groups of  $k$ ) and an unbalanced two-sample meta-analysis (testing significance in mean differences between two groups of  $2^*k/5$  and  $2^*k*4/5$  respectively).

We conducted simulations to evaluate the validity of each estimator in small samples and under violations of their assumptions, namely inhomogeneity of contrast variances  $s_i^2$ , presence of non-negligible between-study variance.

Furthermore, we studied the robustness of contrast-based methods to the presence of mismatched units across studies. To simulate units mismatch, each contrast estimates  $\hat{\beta}_i$  and contrast variance estimates  $s_i^2$  was replaced by a rescaled version:  $\hat{\beta}_i i^* = \hat{\beta}_i a_i$  and  $s_i^{2*} = s_i^2 a_i^2$ . 2 types of unit mismatched were investigated:

- Mismatch in scaling of the contrast vectors:  $a_i$  linearly sampled between 0.4 and 1.6 (mean is 1).
- Mismatch in data scaling from different software: for each group  $a_i=1$  for  $i \leq n_{soft1}$   $a_i=1$  for  $i > n_{soft1}$  with  $n_{soft1} \in \{1/5, 1/2\} * n_g$  where  $n_g$  is the number of studies in the group.

The full set of simulations is summarised in table ???. Code is available at: [https://github.com/cmaumet/zmeta\\_buster](https://github.com/cmaumet/zmeta_buster).

### 2.2.2. Real data

We then compared the nine meta-analytic estimators to the reference approach, MFX GLM, on a dataset of 21 studies of pain. Comparability of contrast estimates depends on equivalent scaling of the data, models, and contrast vectors. Data scaling was consistently performed by FSL, setting median brain intensity to 10,000; model were all created by FSL's Feat tool; and contrasts were constructed to preserve units, with sum of positive elements equal to 1, sum of negative elements equal to -1.

To investigate the presence of between-study variation, we computed the ratio of the between-study variance (estimated using FSL's FLAME [15]) to the total variance (sum of between- and within-study variances), as suggested in [3]. Here we use the average (across study) within-study variance as an estimate of within-study variance in the denominator:  $\hat{\tau}^2 / (\hat{\tau}^2 + \sum_{i=1}^k s_i^2)$ . Using this metric, voxels with values close to 0 present negligible between-study variance and values close to 1 outline appreciable study heterogeneity and the importance of RFX models.

Then for each estimator we compared the standardised meta-analytic statistic to the z-statistic obtained with the reference approach. Overestimation of z-statistic leads to overly optimistic detections while underestimation outline a reduced sensitivity of the approach.

## 3. Results

### 3.1. Robustness to small samples

Fig. 2 presents the one-sample simulation results in small samples, i.e. under small number of studies or small number of subjects. We focus here on methods that are guaranteed only in large samples: FFX and MFX. When the number of subjects is small, FFX is invalid for all within-study variances investigated, regardless of the number of studies included in the meta-analysis. On the other hand MFX GLM is conservative for small number of studies and constant within-study variance. In the presence of large variations in the within-study variances and under small number of studies, MFX is invalid, regardless of the number of subjects included in each study.

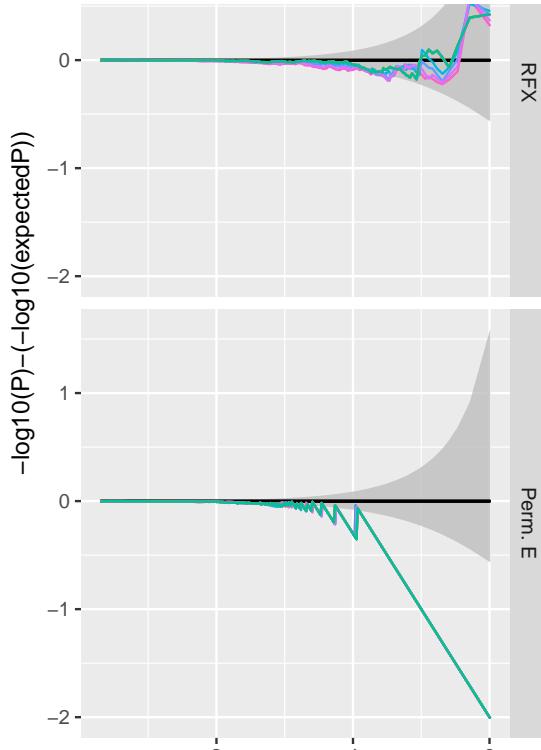


Figure 2: Deviation from theoretical P-values in one-sample tests under ideal circumstances with respect to heterogeneity for each statistical approach ( $\tau^2 = 0$  for FFX and  $\tau^2 = 1$  for MFX) and  $k = 5, 25, 50$  with matched (“nominal”) units.

For small number of studies, permutation methods (including Contrast Permutation and Z permutation) are conservative as expected due to the discrete nature of their distribution (cf. supplementary figure TODO).

Other approaches (TODO) have a nominal behaviour under small sample sizes as expected according to theory. Only Stouffers MFX presents some invalidity, which can be explained by the fact that this is an ad hoc not recommended statistic (cf. supplementary figure TODO).

### 3.2. Robustness to heteroscedasticity

RFX and Perm E. are robust to heteroscedasticity for all settings studied.

### 3.3. Robustness to heterogeneity

#### 3.4. Robustness to units mismatch

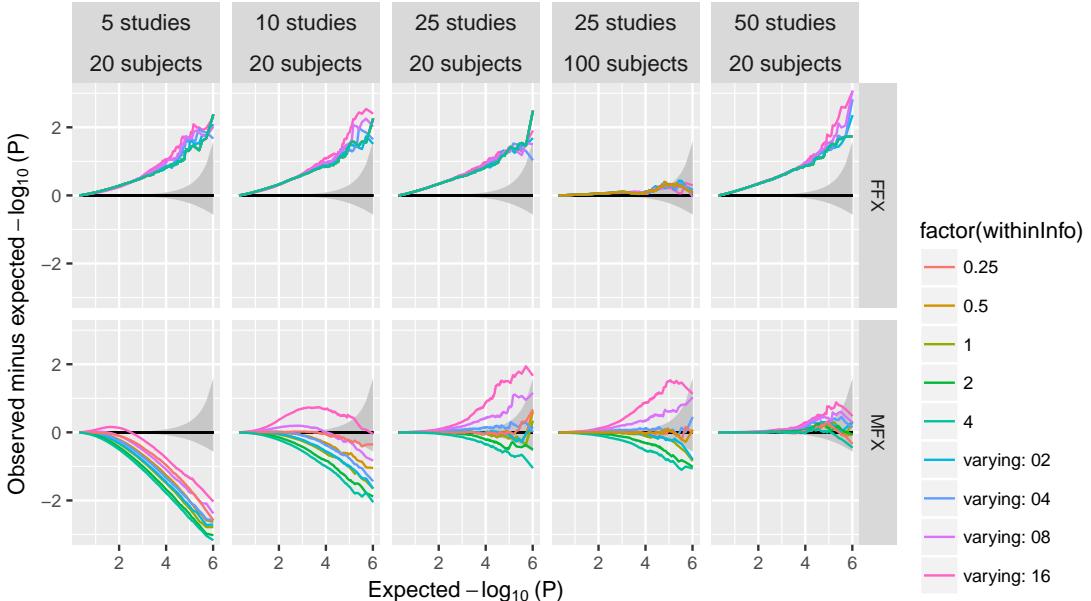
##### 3.4.1. Group meta-analysis

Fig. 4 presents the simulation results under unit mismatches for one-sample tests. In the nominal case, i.e. when the units are matched across studies and contrasts, RFX GLM and Contrast Permutation are valid, as expected. Surprisingly, while MFX GLM is valid for constant within-study variances it is invalid when the within-study variances are varying (TODO: Need to understand this point better!!). For large values of Z, Contrast Permutation is conservative as expected due to the discrete nature of its distribution. In the presence of a high within-study variance, MFX GLM also appears to be conservative. RFX GLM displays the best behaviour with a pattern that is within the 95% confidence interval of the theoretical Z for all within-study variance studied and only slightly conservative when the within-study variances are varying across studies.

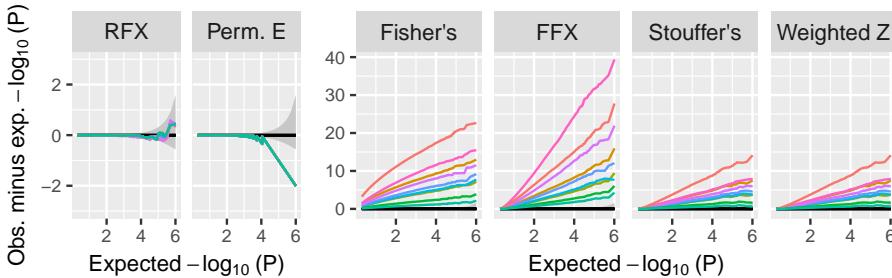
When different scaling algorithm are used, i.e. with different neuroimaging software packages (provided that the scaling target has been accounted for), Contrast Permutation still has a behaviour that is very similar to nominal. For small or high within-study variances ( $\sigma_i^2 = 5, \sigma_i^2 = 10, \sigma_i^2 = 80$ ), MFX GLM and RFX GLM display invalidity for small Z’s and conservativeness for large Z’s. When within-study and between-study variances are similar, MFX GLM is slightly invalid for all Z values. When the contrast are scaled differently, we observe a very similar pattern than for different scaling algorithm.

## Robustness of the meta-analytic estimators

### A Small sample sizes



### B Heteroscedasticity



### C Heterogeneity

Figure 3: Deviation from theoretical P-values for a range of varying within-study variances with  $k = 25$  and  $n = 20$  with matched (“nominal”) units.

#### 3.4.2. Balanced between-group meta-analysis

Fig. 5 presents the simulation results under unit mismatch for two-sample tests. For the nominal case, GLM RFX, GLM MFX and contrast permutation provide valid estimates. Contrast Permutation is conservative for large Z values. When within-study and between-study variances are similar, MFX GLM is slightly invalid for all Z values. RFX GLM display the best behaviour with a pattern that is within the 95% confidence interval of the theoretical Z.

When different scaling algorithm are used Contrast Permutation still has a behaviour that is very similar to nominal. For small or high within-study variances ( $\sigma_i^2 = 5, \sigma_i^2 = 10, \sigma_i^2 = 80$ ), MFX GLM and RFX GLM display invalidity for small Z’s and conservativness for large Z’s. When within-study and between-study variances are similar, MFX GLM is invalid for all Z values. When the contrast are scaled differently, we observe a very similar pattern than for different scaling algorithm.

#### 3.4.3. Unbalanced between-group meta-analysis

Fig. 6 presents the simulation results under unit mismatch for unbalanced two-sample tests. For the nominal case, GLM RFX and contrast permutation provide valid estimates. Contrast Permutation is conservative for large Z values. For the nominal case, GLM MFX is slightly invalid for balanced within and between-study variances. RFX GLM display the best behaviour with a pattern that is mostly within the 95% confidence interval of the theoretical Z.

For small or high within-study variances ( $\sigma_i^2 = 5, \sigma_i^2 = 10, \sigma_i^2 = 80$ ), MFX GLM and RFX GLM display invalidity for

## Unit mismatch: one-sample test

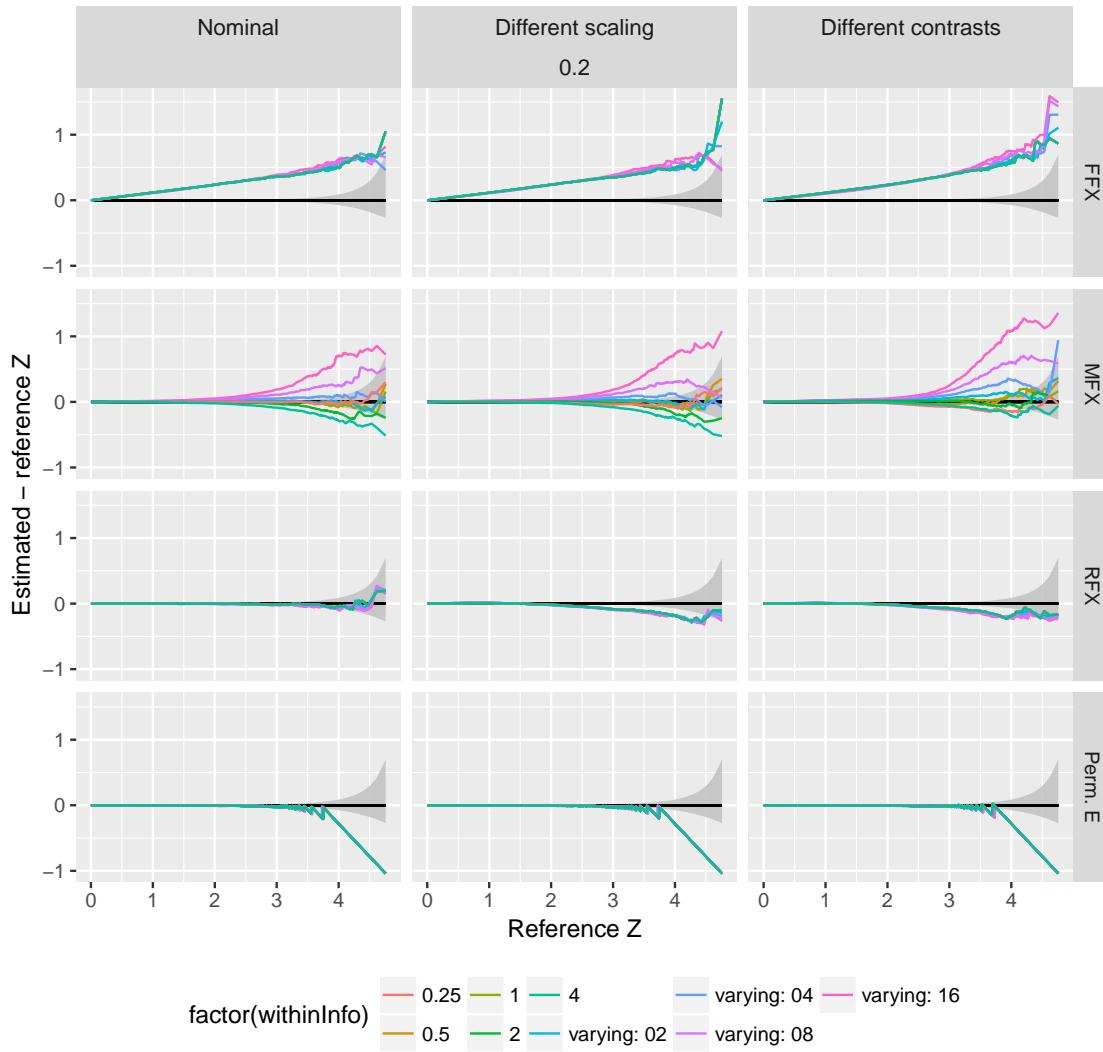


Figure 4: Deviation from theoretical Z in one-sample tests with unit mismatch, under ideal circumstances for each statistical approach ( $\tau^2 = 1$  and  $k = 5, 25, 50$  with matched (“nominal”) or mismatched (“different scaling target”, “different scaling algorithm”, “different contrast vector scaling”) units.

small Z’s and conservativeness for large Z’s. A pattern that is also seen for Contrast permutation (across all within-study variances). When within-study variance is large, MFX GLM is invalid for all Z values.

When the contrast are scaled differently, MFX GLM display the best behaviour but is conservative. RFX GLM and contrastpermutation display a very similar pattern than with different scaling.

### 3.5. Small sample properties

#### 3.5.1. Group meta-analysis

Fig. 11 presents the simulation results for a one-sample test with  $\tau^2 = 1$  and a sample size  $k = 5, 25, 50$ . For the nominal case, i.e. when the units are matched across studies and contrasts, MFX GLM, RFX GLM and Contrast Permutation are all valid, as expected. For small sample sizes ( $k = 5$ ), MFX GLM and contrast permutation are both very conservative. For large values of Z, Contrast Estimation is conservative as expected due to the discrete nature of its distribution. More surprising, in the presence of a high within-study variance, MFX GLM also appears to be conservative. RFX GLM displays the best behaviour with a pattern that is always within the 95% confidence interval of the theoretical Z.

## Unit mismatch: two-sample test

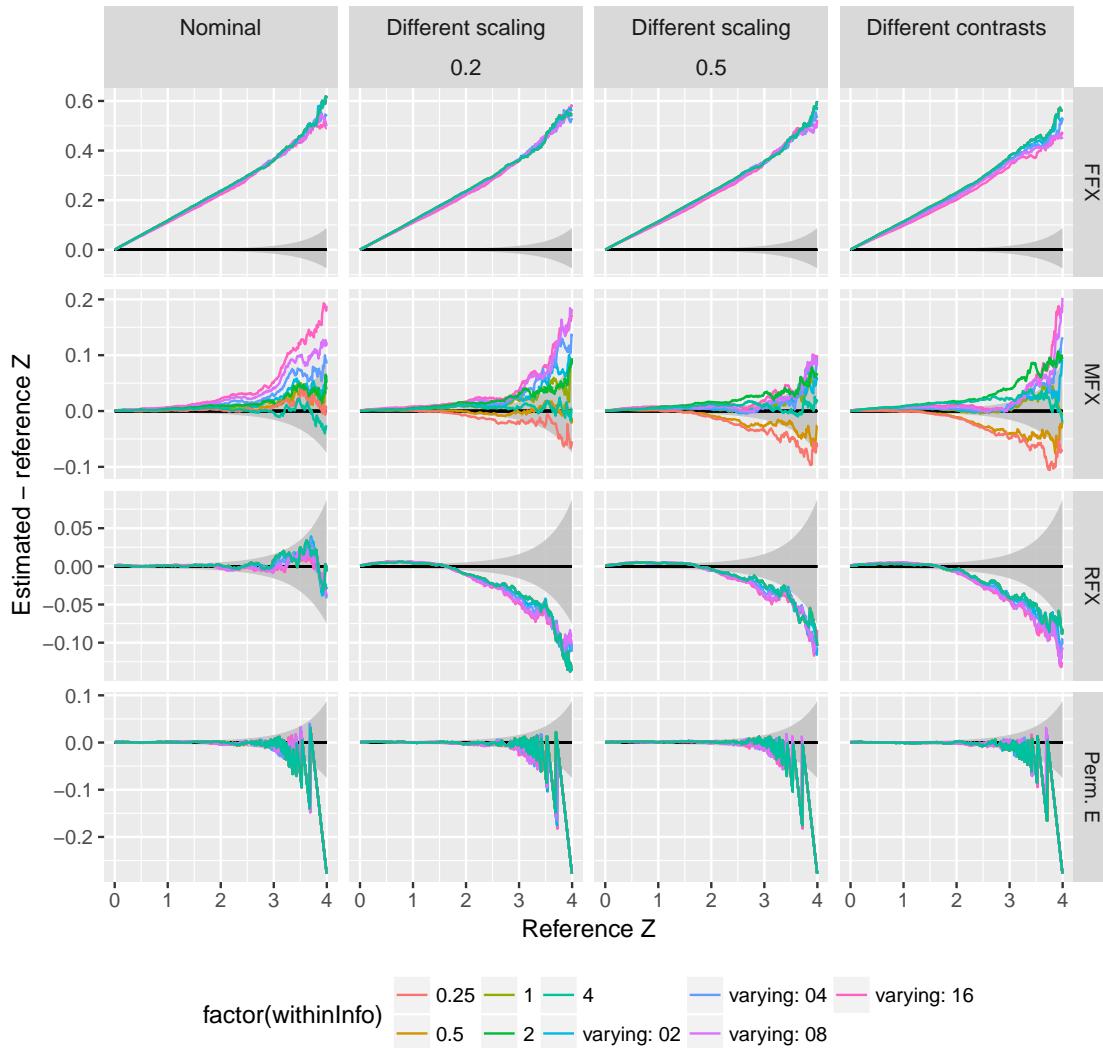


Figure 5: Deviation from theoretical Z in two-sample tests with unit mismatch, under ideal circumstances for each statistical approach ( $\tau^2 = 1$  and  $k = 5, 25, 50$  with matched (“nominal”) or mismatched (“different scaling target”, “different scaling algorithm”, “different contrast vector scaling”) units.

### 3.5.2. Balanced between-group meta-analysis

Fig. TODO presents the simulation results for a two-sample meta-analyses with  $\tau^2 = 1$  and a sample size  $k = 25$ . For the nominal case, GLM RFX, GLM RFX and contrast estimation provide valid estimates. Contrast Permutation is conservative for large Z values. Both RFX GLM and MFX GLM display the best behaviour with a pattern that is within the 95% confidence interval of the theoretical Z.

In the extreme case of different scaling target, contrast permutation is always valid with a pattern very similar than its nominal behaviour. GLM RFX is valid for Z values greater than 1.5, which is the area of interest in detections, but display a strong conservativness, more pronounced than the Contarst Permutation. GLM MFX is slightly invalid for all within-study variances except the largest one when 20% of the studies come from the second software.

### 3.5.3. Unbalanced between-group meta-analysis

Fig. TODO presents the simulation results for unbalanced two-sample meta-analyses with  $\tau^2 = 1$  and a sample size  $k = 25$ . For the nominal case, MFX GLM, GLM RFX and contrast permutation provide valide estimate. As expected due to the discrete nature of its ampling distribution, contrast permutation is conservative for large Z value. GLM RFX is conservative. RFX GLM is closest to the theoretical behaviour with Z-values that are always within the 95% confidence interval.

### Unit mismatch: unbalanced two-sample test

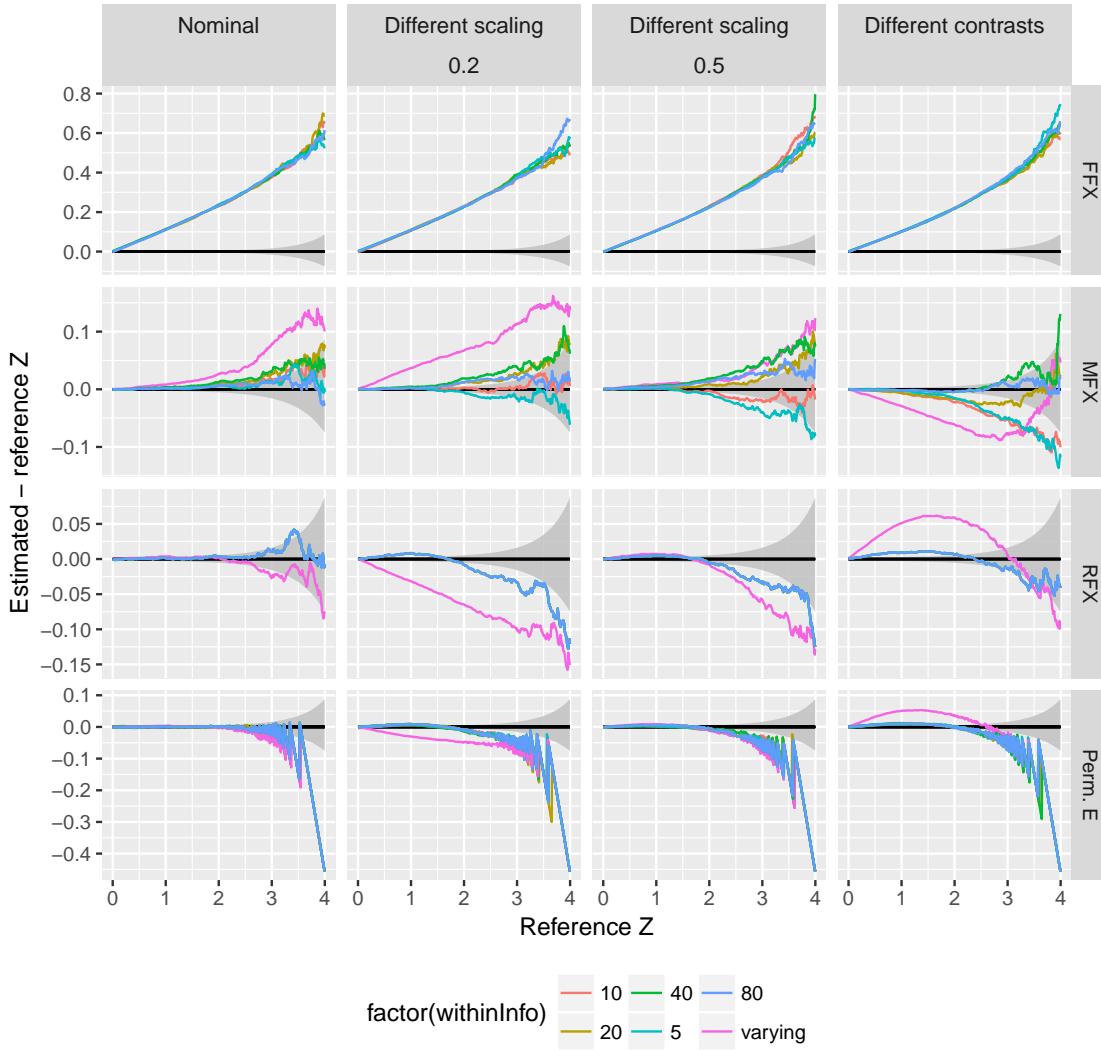


Figure 6: Deviation from theoretical Z in unbalanced two-sample tests with unit mismatch, under ideal circumstances for each statistical approach ( $\tau^2 = 1$  and  $k = 5, 25, 50$  with matched (“nominal”) or mismatched (“different scaling target”, “different scaling algorithm”, “different contrast vector scaling”) units.

In the extreme case of different scaling target, MFX GLM is always valid but slightly conservative. RFX GLM is valid for Z values greater than 1.5 (area of interest in detections) but conservative. Similarly contrast permutation is invalid for Z smaller than 1.5 and conservative otherwise. This can be explained by the violation of the exchangeability condition.

When different scaling algorithm are used, (same paragraph as for one-sample test)

When the contrast are scaled differently, we observe a very similar pattern than for different scaling algorithm with higher variance of the estimates.

#### 3.6. Simulations

Fig. 1 displays the false positive rate at  $p < 0.05$  obtained for the nine estimators over all set of parameters in the absence and presence of between-study variation. As expected, the fixed-effects meta-analytic summary statistics, i.e. Fisher’s, Stouffer’s and weighted Stouffer’s estimates, are liberal in the presence of study heterogeneity. The original Fisher’s approach is the most invalid. More surprising, FFX GLM is also invalid with homogeneous studies. The explanation is over-estimation of degrees-of-freedom (DF); while DF is computed as  $(\sum n_i - 1) - 1$ , under heteroscedasticity (from  $\sigma_i$  or  $n_i$ ) it will be much lower [14]. Z MFX and GLM RFX provide valid estimates, and the permutation estimates are valid but tend to be conservative with greater variation in false positive rates.

## Small sample sizes: one-sample test

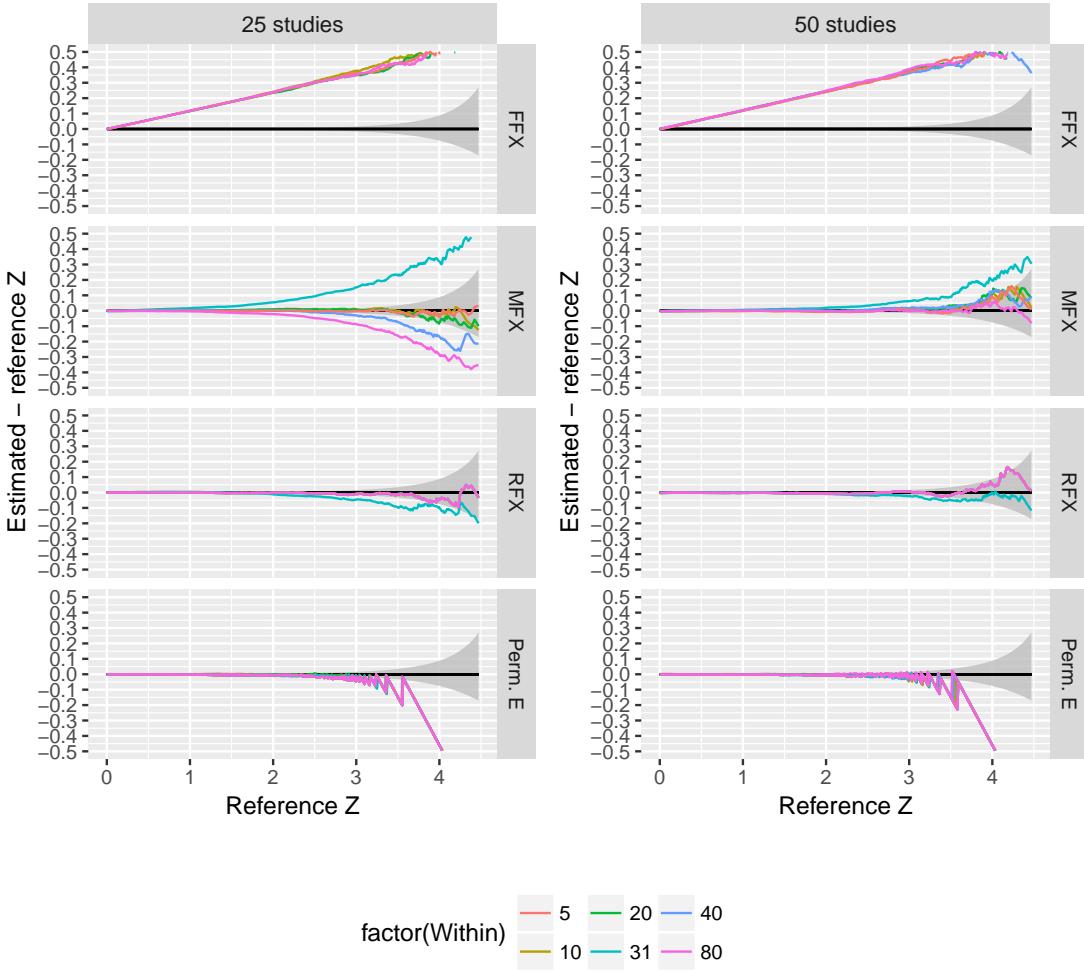


Figure 7:

The impact of the number of studies involved in the meta-analysis and of the size of the within-study variance are investigated in Fig. TODO. Permutation inference is valid but conservative when 5 studies are used; this is because there are only  $2^5 = 32$  possible permutations and thus  $1/32 = 0.03125$  is largest attainable valid P-value. All approaches perform equally as soon as 10 or more studies are included in the meta-analysis.

### 3.7. Real data

Fig. 14 plots the difference between the z-score estimated by each meta-analytic approach against the reference z-score computed with MFX GLM. All FFX statistics provide overly optimistic z-estimate suggesting, again, that study heterogeneity is present in the studied dataset. Among the RFX meta-analytic approaches, GLM RFX and contrast permutations provide z-scores estimate that are equal or smaller than the reference. Z permutation provides slightly larger z-scores between 1 and 3 (reference p-values between 0.16 and 0.0013) but is mostly in agreement with the reference z-scores. On the other hand, Z MFX is more liberal than the reference for z-score ranging from 3 to 5 (reference p-values between 0.0013 and 2.9e-07) and more stringent for z-scores smaller than 5.

## 4. Conclusion

We have compared nine meta-analytic approaches in the context of one-sample test. Through simulations, we found the expected invalidity of standard FFX approaches in the presence of study heterogeneity, but also of FFX GLM even with no

## Small sample sizes: two-sample test

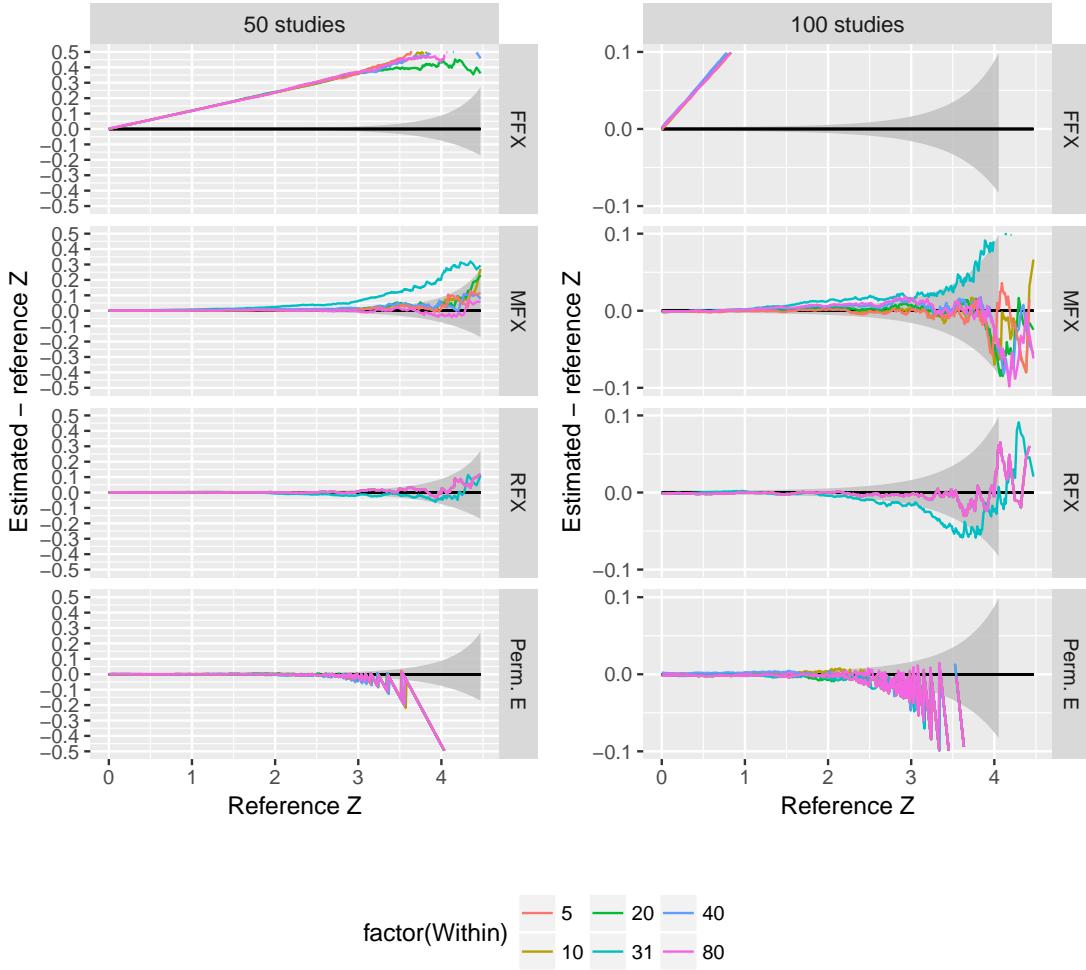


Figure 8:

between-study variation. In a real dataset of 21 studies of pain, there was evidence for substantial between-study variation that supports the use of RFX meta-analytic statistics. When only contrast estimates are available, RFX GLM was valid. This is in line with previous results on within-group one-sample t-tests studies [9]. When only standardised estimates are available, permutation is the preferred option as the one providing the most faithful results. Further investigations are needed in order to assess the behaviour of these estimators in other configurations, including meta-analyses focusing on between-study differences.

### 5. Acknowledgements

We gratefully acknowledge the use of this data from the Tracey pain group, FMRI, Oxford.

- [1] Meta-analysis of neuroimaging data: a comparison of image-based and coordinate-based pooling of studies. *NeuroImage*, 45(3):810–23, 2009.
- [2] Katherine S Button, John P a Ioannidis, Claire Mokrysz, Brian a Nosek, Jonathan Flint, Emma S J Robinson, and Marcus R Munafò. Power failure: why small sample size undermines the reliability of neuroscience. *Nature reviews Neuroscience*, 14(5):365–76, 2013.

## Small sample sizes: unbalanced two-sample test

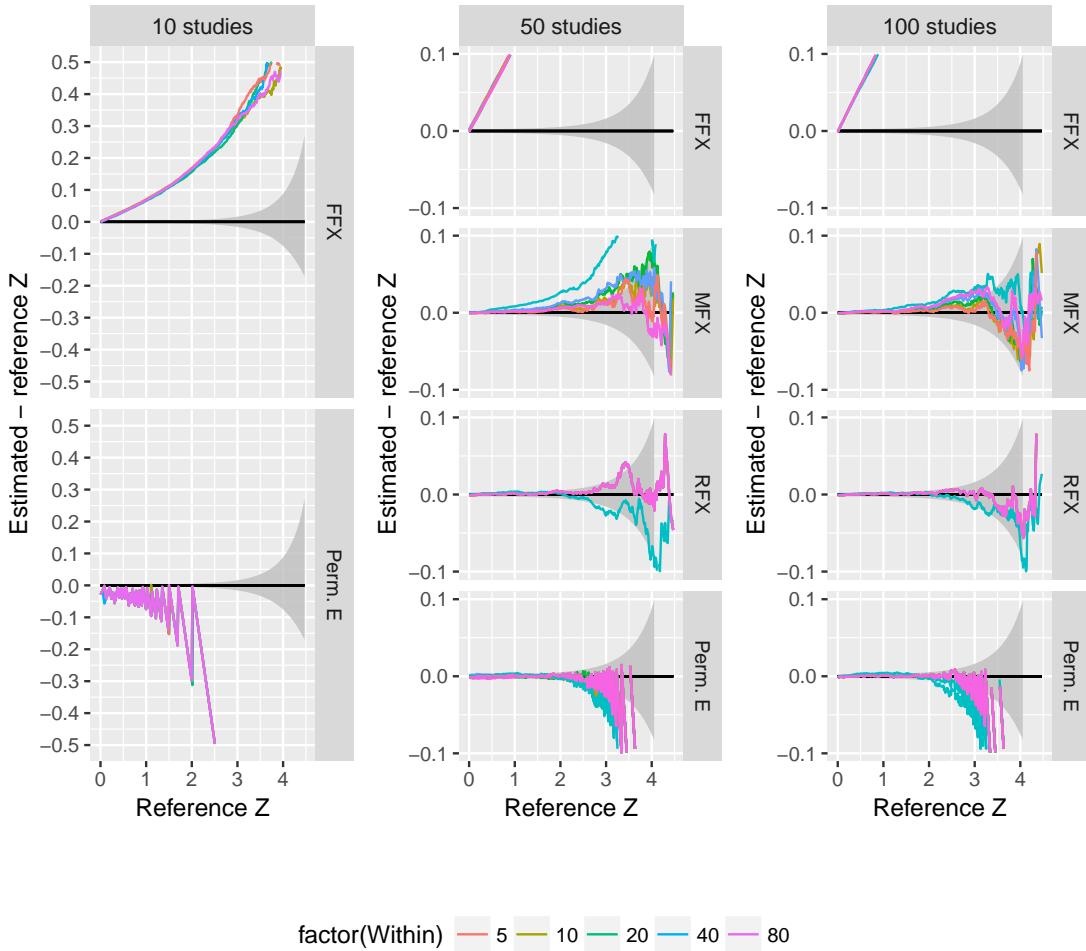


Figure 9:

- [3] G. Chen, Z. S. Saad, A. R. Nath, M. S. Beauchamp, and R. W. Cox. FMRI group analysis combining effect estimates and their variances. *NeuroImage*, 60(1):747–65, March 2012.
- [4] Gang Chen, Paul A Taylor, and Robert W Cox. Is the Statistic Value All We Should Care about in Neuroimaging? *bioRxiv*, (September):064212, 2016.
- [5] P. Cummings. Meta-analysis based on standardized effects is unreliable. *Archives of pediatrics & adolescent medicine*, 158(6):595–7, 2004.
- [6] R.A. Fisher. *Statistical Methods for Research Workers*. Oliver and Boyd, Edinburgh, 1932.
- [7] A. P. Holmes, R. C. Blair, G. J.D. Watson, and I. Ford. Nonparametric Analysis of Statistic Images from Functional Mapping Experiments. *Journal of Cerebral Blood Flow and Metabolism*, 16:7–22, 1996.
- [8] T. Liptak. On the combination of independent tests. *Magyar Tud. Akad. Mat. Kutato Int. Kozl.*, 3:171–197, 1958.
- [9] J. A. Mumford and T. E. Nichols. Simple group fMRI modeling and inference. *NeuroImage*, 47(4):1469–75, 2009.
- [10] T. E. Nichols and A. P Holmes. Nonparametric permutation tests for functional neuroimaging: a primer with examples. *Human brain mapping*, 15(1):1–25, 2002.

### RFX: under varying within-study variance

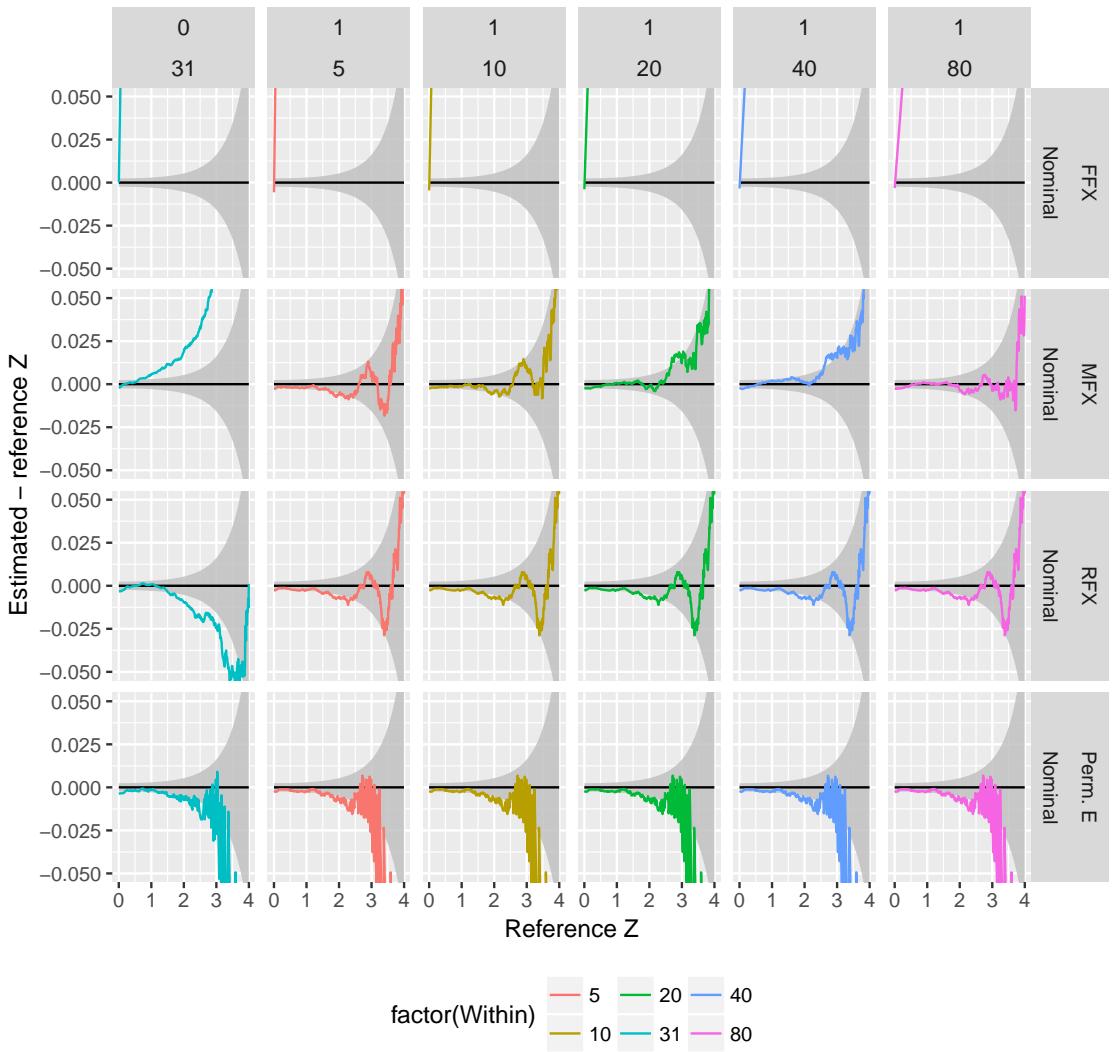


Figure 10:

- [11] Thomas E. Nichols. Spm plot units, 2012.
- [12] Russell A Poldrack, Chris I Baker, Joke Durnez, Krzysztof J Gorgolewski, Paul M Matthews, Marcus R Munafò, Thomas E Nichols, Jean-baptiste Poline, Edward Vul, and Tal Yarkoni. Scanning the horizon : towards. *Nature Publishing Group*, 2017.
- [13] J. Radua and D. Mataix-Cols. Meta-analytic methods for neuroimaging data explained. *Biology of mood & anxiety disorders*, 2(1):6, 2012.
- [14] F. E. Satterthwaite. An approximate distribution of estimates of variance components. *Biometrics Bulletin*, 2:110–114, 1946.
- [15] S. Smith, P. R. Bannister, C. Beckmann, M. Brady, S. Glare, D. Flitney, P. Hansen, M. Jenkinson, D. Leibovici, B. Ripley, M. Woolrich, and Y. Zhang. FSL : New Tools for Functional and Structural Brain Image Analysis. (6):2001, 2001.
- [16] S. Stouffer, L. DeVinney, and E. Suchmen. *The American Soldier: Adjustment During Army Life*, volume 1. Princeton University Press, Princeton, NJ, 1949.
- [17] D V Zaykin. Optimally weighted Z-test is a powerful method for combining probabilities in meta-analysis. *Journal of evolutionary biology*, 24(8):1836–41, 2011.

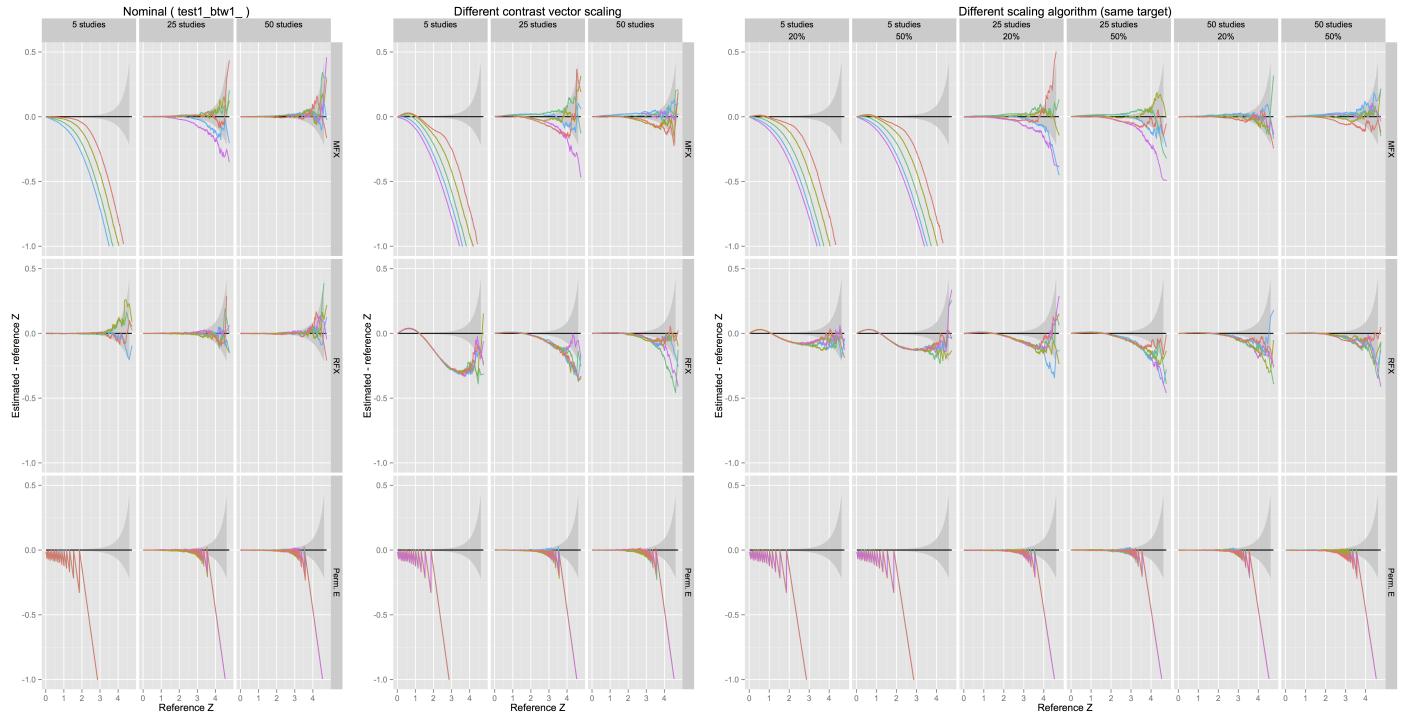


Figure 11: Deviation from theoretical Z in one-sample tests with  $\tau^2 = 1$  and  $k = 5, 25, 50$  with matched (“nominal”) or mismatched (“different scaling target”, “different scaling algorithm”, “different contrast vector scaling”) units.

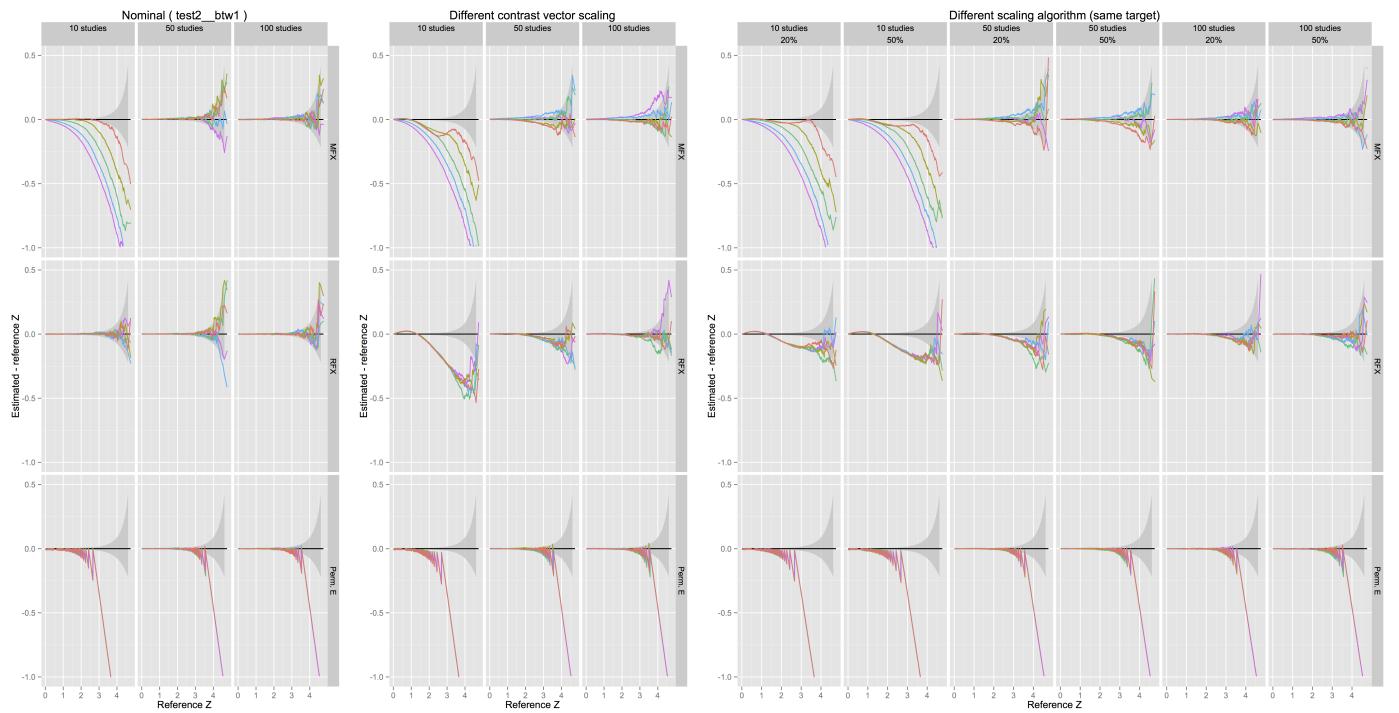


Figure 12: Deviation from theoretical Z in balanced two-sample tests with  $\tau^2 = 1$  and  $k = 25$  with matched (“nominal”) or mismatched (“different scaling target”, “different scaling algorithm”, “different contrast vector scaling”) units.

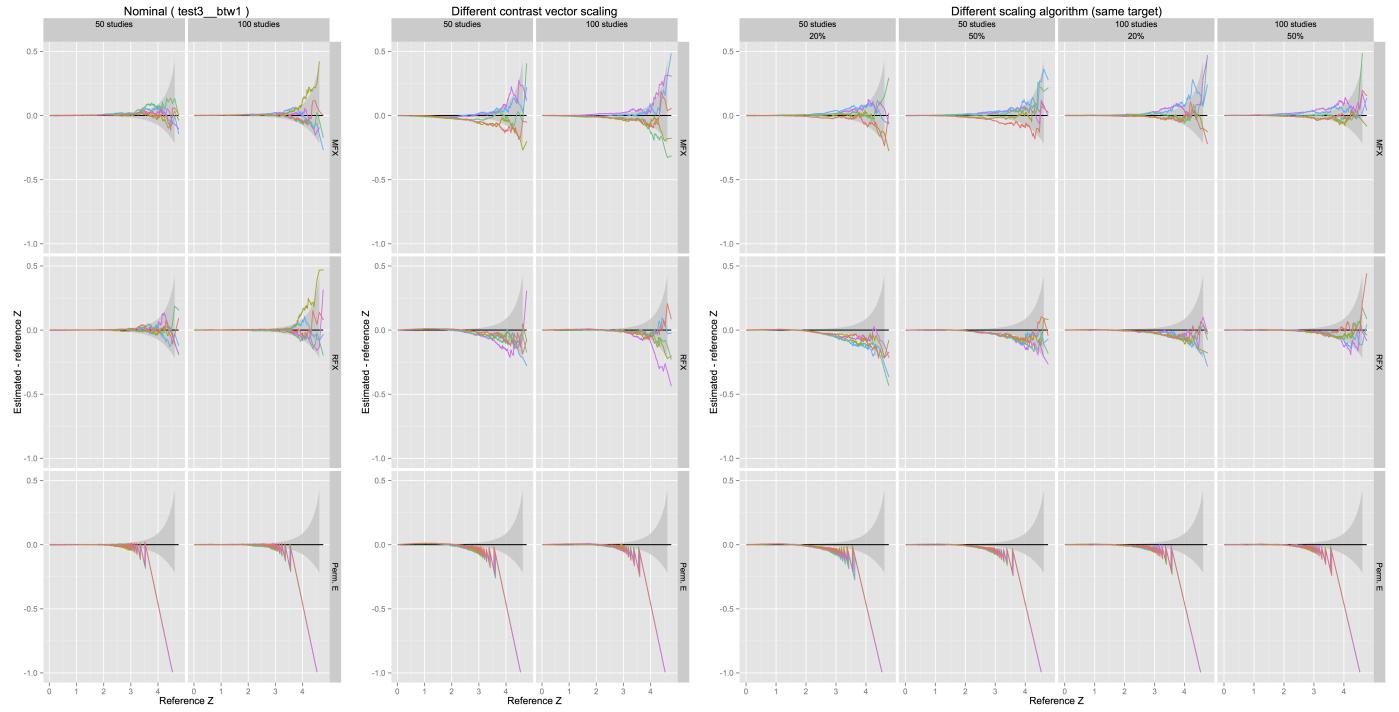


Figure 13: Deviation from theoretical Z in unbalanced two-sample tests with  $\tau^2 = 1$  and  $k = 25$  with matched (“nominal”) or mismatched (“different scaling target”, “different scaling algorithm”, “different contrast vector scaling”) units.

Figure 14: Difference between the z-score estimated from each meta-analytic approach and the reference z-score from MFX GLM as a function of reference z-score.