

Revisiting the assessment of inter-individual differences in fMRI activations-behavior relationships.

Maël Lebreton^{1,2,*} & Stefano Palminteri^{3,4}

¹Amsterdam Brain and Cognition (ABC), and

²Amsterdam School of Economics (ASE), Universiteit van Amsterdam, 1018 WB Amsterdam, the Netherlands.

³Institute of Cognitive Sciences (ICN), University College London, WC1N 3AR, London, United Kingdom.

⁴Laboratoire de Neurosciences Cognitives (LNC), INSERM U960, École Normale Supérieure, 75005 Paris, France.

* To whom correspondence should be addressed (m.p.lebreton@uva.nl)

Abstract:

Characterizing inter-individual differences is critical to realize neuroimaging full potential, but can hardly be achieved without accurately assessing the statistical dependencies between inter-individual differences in behavior and inter-individual differences in neural activity. In this manuscript, we consider two hypotheses: 1) BOLD signal scales linearly with behavioral variables across individuals and 2) BOLD signal encodes behavioral variables on a similar scale across individuals. We formally show that these two hypotheses induce opposite brain-behavior correlational results in group-level analyses, illustrating the importance of explicitly testing inter-individual brain-behavior scaling before engaging in the study of functional inter-individual differences. To further evidence the relevance of this framework, we illustrate its practical consequences for model-based fMRI using computational simulations, and demonstrate its empirical robustness in four fMRI studies investigating values coding in the prefrontal cortex. This may constitute an important step forward in our conceptualization, analysis and interpretation of inter-individual differences in cognitive neurosciences.

INTRODUCTION:

“There is very little difference between one man and another; but what little there is, is very important.”ⁱ They are two complementary goals in cognitive neuroscience: understanding the average – *typical* - brain by linking its structure and functions with cognition and behavior, and understanding how individuals differ from each-others from the normal to the pathological ranges. With respect to these two quests, inter-individual differences in brain-behavior relationships are fundamentally important either because they constitute a statistical challenge to understand the typical brain, or because they represent the very object of interest (Braver et al., 2010; Gabrieli et al., 2015).

Task-related functional neuroimaging (fMRI) constitutes a tool of choice to investigate the neurobiological underpinnings of cognition. Initially confined to the mapping of the typical brain, fMRI is increasingly used to investigate the neural bases of differential cognition. Individual brain activations are linked to external heterogeneity factors, such as the diagnostic criterion for some pathology, psychosocio-economic measures -i.e. “traits”-, or to behavioral measures recorded during the experiment – i.e. “task performance” (**Figure.1.A** and **1.B**). Significant associations are then

interpreted in term of individual neural resource mobilized to complete the task – sometimes with opposite post-hoc rationalization. It is paradigmatic in the example of executive control literature, linking activations in the frontal regions and individual performance: positive associations are typically interpreted as an effective increase in cognitive control mobilization, whereas negative associations are interpreted as an increase in neural efficiency. This represents a critical failure of the current inter-individual difference framework (Poldrack, 2015; Yarkoni and Braver, 2010).

In the present paper, we suggest that these inconsistencies are possible because the hypotheses concerning the inter-individual relative scaling of brain signal (typically Blood Oxygen Level Dependent signal, BOLD) with behavior are not explicitly stated and tested. We therefore present the scaling issue and evidence its impact on inter-individual fMRI results by intuitively and mathematically demonstrating that two simple - and equally plausible - scaling hypotheses can produce opposite results and conclusions. Using computational simulations, we show how the scaling issue can percolate in model-based fMRI, and we provide a priori principle to process model-based variables. Finally, in order to prove the interest of our framework and validate its robustness, we assess -and replicate in four imaging datasets- the statistical law describing the relationship between inter-individual differences in prefrontal BOLD signal

ⁱWilliam James (1897). The Importance of Individuals, In *The Will to Believe and Other Essays in Popular Philosophy*

and inter-individual differences in the critical variable of economic decision-making: values (Lebreton et al., 2009, 2012, 2015; Palminteri et al., 2015).

RESULTS

Theoretical considerations on scaling-laws

fMRI analysis background. The classical fMRI analysis scheme relies on the general linear model (GLM) framework (Friston et al., 1994), and follow a *multi-level summary statistics* approach (Beckmann et al., 2003; Friston et al., 2005; Holmes and Friston, 1998; Woolrich et al., 2004; Worsley et al., 2002). This provides a practically good approximation of mixed-effects designs (Beckmann et al., 2003; Friston et al., 2005; Mumford and Nichols, 2009), and results have been shown to be consistent across software packages (Bennett and Miller, 2010; Gold et al., 1998; Morgan et al., 2007). This approach can be briefly summarized as follows: in a first step, the linear relation between the time series of BOLD signal within a specific brain region (or voxel) and the time series of the different explanatory variables are assessed at the individual level. For each individual k , this entails designing a first-level GLM: $Y_k = \beta_k \cdot x_k + u_k$, where Y_k is the BOLD time-series, x_k contains the explanatory variables time-series, u_k is a Gaussian noise, and β_k is the vector of the unstandardized linear coefficient of regressions to be estimated. First-level summary statistics, i.e. estimated individual betas $\widehat{\beta}_k$ or contrasts of betas, are then used in a population level analysis. For studies aiming at mapping a cognitive function in the typical brain, this second-level analysis is usually a random effect – a “one-sample t-test”. However, second-level analyses can also aim at investigating differences in activations between different categories subject (e.g. pathological or non-pathological sub-populations), or across a continuum of subjects (e.g. following individual “traits”). These differential analyses can respectively be implemented with “two-sample t-test” or ANOVAs for the categorical case, and with second-level “multiple regression” for the continuous case (**Figure.1.B.**).

Scaling laws and inter-individual differences.

Critically, second level analyses relies on *unstandardized* first-level betas ($\widehat{\beta}_k$). As recalled in the **Experimental procedures**, the magnitude of $\widehat{\beta}_k$ depends on the ratio between the standard deviations of the BOLD ($\sigma(Y_k)$) and of the experimental factor time series ($\sigma(x_k)$). In the following, we refer to the statistical relationships between $\sigma(Y_k)$ and $\sigma(x_k)$ as *scaling laws*, and investigate how they impact assessments of inter-individual differences in brain activity as indexed by $\widehat{\beta}_k$. We consider two opposite, though neuro-

biologically plausible, scaling laws between a behavioral measure and the BOLD signal that we called the *proportional* and the *normalization* hypotheses. Without loss of generality, we consider those scaling laws in an idealized situation, where the BOLD signal in a brain region (Y_k) encodes a behavioral parametric measure of interest (x_k) whose distribution ($\sigma(x_k)$) varies between individuals k due to a heterogeneity source (see **Figure.1.A** and **B**, and **Experimental procedures** for formal definitions). The *proportional* hypothesis represents the typical theoretical case underlying random-effect analyses in neuroimaging: the brain-behavior linear relation captured by the $\widehat{\beta}_k$ is kept constant across individual, notwithstanding random variations. As a consequence, no difference in “activations” -as measured by $\widehat{\beta}_k$ - accounts for the differences in behavior x_k (**Figure 1.C.a**). In other words, under the proportional hypothesis, *behavioral differences are reflected by underlying differences in brain activity, but not by differences in brain activations, as indexed by the standard fMRI analysis output*. Under the *normalization* hypothesis, the BOLD signal also encodes the variable of interest, but on an identical scale across individuals, independent of the behavioral output. Then, inter-individual or between-group differences in “activations” -as measured by $\widehat{\beta}_k$ -, logically derive from differences in behavior x_k (**Figure 1.C.b**). In other words, under the normalization hypothesis, *behavioral differences are not reflected by underlying differences in brain activity, but by differences in brain activations, as indexed by the standard fMRI analysis output*. These derivations mean that the interpretations of inter-individual differences in $\widehat{\beta}_k$ should be made with caution. For instance, let's assume that the activity in a region of interest (ROI) Y_k is causally responsible for a behavioral measure x_k , and that a heterogeneity factor causes changes in Y_k , inducing proportional changes in x_k . In the proportional context, comparisons between $\widehat{\beta}_k$ in the ROI might be inconclusive, potentially misleading to false-negative conclusions about the role of the ROI in the observed inter-individual differences due to the heterogeneity factor (**Figure.1.D**).

Neutralizing inter-individual differences in behavior. To circumvent the scaling issue, one might be tempted to neutralize inter-individual differences in behavioral by individually normalizing –Z-transforming- the measure of interest. However, as shown in the **Experimental procedures**, under the *proportional* hypothesis, Z-scoring induces a monotonic relationship between the “activation” ($\widehat{\beta}_k^Z$) and the standard deviation of the behavior, which was absent before Z-scoring (**Figure 1.C.c**).

Conversely, under the *normalization* hypothesis, Z-scoring cancels the monotonic relationship between the “activation” ($\widehat{\beta}_k^Z$) and the standard deviation of the behavior, leaving no inter-individual or between-group differences (**Figure 1.C.d**). These derivations critically demonstrate that the significance and interpretation of typical inter-individual correlational results largely depend on an interaction between the behavioral variable pre-processing and the scaling hypothesis. The interaction between the scaling hypothesis and the Z-scoring of the behavioral variable also has a significant impact on the results of second-level random-effect analyses (i.e. one sample t-test on the $\widehat{\beta}_k$): under the proportional scaling hypothesis, Z-scoring \mathbf{x}_k introduces a systematic variance in the $\widehat{\beta}_k^Z$ due to the dependence of $\widehat{\beta}_k^Z$ to $\sigma(\mathbf{x}_k)$, which can decrease the significance of random-effect model (see **Figure 1.C.a** vs **1.C.c**). On the opposite, under the normalization scaling hypothesis, Z-scoring \mathbf{x}_k erases the systematic variance initially due to the difference in $\sigma(\mathbf{x}_k)$, hence can increase the significance of random-effect model (see **Figure 1.C.b** vs **1.C.d**). Systematic investigation of brain-behavior scaling laws might therefore provide precious information about how to best pre-process variables of interest to perform classical fMRI analyses.

Task designs and inter-individual differences.

Consider a typical fMRI investigation attempting to link inter-individual differences in a performance “trait” (IQ, psychometric measure or socioeconomic status) to some neural activations, derived from a trial-by-trial measure of a task-performance explanatory variable \mathbf{x}_k (e.g. confidence, decision-value, or reaction time). If the task is too easy, high-performing people exhibit ceiling task-performances, (generating smaller $\sigma(\mathbf{x}_k)$), inducing a negative correlation between individual performance “trait” and $\sigma(\mathbf{x}_k)$. Symmetrically, if the task is too hard, low-performing people might exhibit flooring task performances, (generating smaller $\sigma(\mathbf{x}_k)$), inducing a positive correlation between individual performance “trait” and $\sigma(\mathbf{x}_k)$. Given the dependencies between $\widehat{\beta}_k$ and $\sigma(\mathbf{x}_k)$, this will lead to opposite correlations between “activations” ($\widehat{\beta}_k$), and individual performances “trait” (proxied by $\sigma(\mathbf{x}_k)$) (**Figure 1.E**). This example illustrates the potential of the scaling framework to explain inconsistent results in the cognitive neuroscience literature of inter-individual differences.

Activation measures derivable from $\widehat{\beta}_k$ exist, which do not depend on $\sigma(\mathbf{x}_k)$ but also carry a different meaning (e.g. “t-values” or “Z-values”, see **Experimental procedures**). It seems critical, for future studies of inter-individual differences, to

clarifies the question they are attempting to answer -is the brain region linearly coding the variable on a *different scale* in different subjects? Or: is the region linearly coding the variable with a different *reliability* in different subjects?- and make appropriate choices, descriptions and interpretations of their activation measure.

Scaling-laws and model-based fMRI

Computational models and latent variable distributions.

In this second part, we will develop some consequences of the scaling issue for model-based fMRI. Model-based fMRI typically uses as dependent variables in first-level GLMs (\mathbf{x}_k) *latent variables* derived from individual choice patterns: a computational model is selected, its free-parameters are adjusted to account at best for behavioral data, and the model parameters are used to generate the *latent variables* of interest \mathbf{x}_k (O’Doherty et al., 2007). Importantly, the model free-parameters can be either considered as fixed (FFX) -i.e. shared across individuals - or random-effects (RFX) -i.e. each individual’s parameters are drawn from a common population distribution (Daw, 2011). In the case of random-effects, model free-parameters almost inevitably impact the latent variable distribution properties, including $\sigma(\mathbf{x}_k)$, with consequences for the scaling law. Consider the example of using an expected-utility model in decision-making under risk: for choice situations involving simple prospects combining potential gain(s) g with a probability p of winning, expected utility theory stipulates that agents choose the option which maximize the expected utility $\mathbf{eu} = p \times g^r$. In this model r is the utility curvature free-parameter and captures the individual attitude toward risk (see e.g. (Bernoulli, 1954)). By simulating a task, where individual are confronted with several options -i.e. combinations of g and p -, and by computing \mathbf{eu} for different and plausible values of r , we can unambiguously show that the model free-parameter r monotonically determines $\sigma(\mathbf{eu})$ (**Figure.2.A**).

Task designs and computational models.

In many cases however, the link between the model parameters and the standard deviation of the latent variable $\sigma(\mathbf{x}_k)$ is not trivial and largely depends on the task setting and the stimuli space. To illustrate this second point, we take the example of the hyperbolic delay-discounting model. This model states that in choice situations involving prospects combining future gain(s) g deferred by delays D , decision-makers are choosing the option with the highest discounted value $\mathbf{dv} = \frac{g}{1+k \times D}$. In this model k is the discounting free-parameter, and captures individual patience or impulsivity (see e.g. (Ainslie

and Haslam, 1992)). We simulated a task, where individual are confronted with 2 different sets of options combining g and D , and computed the discounted value of those options for plausible values of k . We then show that, depending of the option set, the model free-parameter k and $\sigma(\mathbf{d}\mathbf{v})$ can be either positively or negatively monotonically related (**Figure.2.B**).

From latent variables to choice functions. Finally, the associations between the model free-parameters and the standard deviation of the latent variable $\sigma(\mathbf{x}_k)$ are not limited to parameters involved in the computation of the latent variable \mathbf{x}_k . In the case of value-based decision-making, this means that the standard deviation of the value variable $\sigma(\mathbf{v}_k)$ may not only be linked to parameters of the value function, but can be linked to other parameters, e.g. controlling the decision policy. To illustrate this third point, we ran a last set of simulations, taking for example a simple reinforcement-learning situation, where decision-makers have to learn, by trial and errors, to select the stimulus which is associated to a higher reward rate. The Rescola-Wagner model proposes that agents learn the value of stimuli (Q-values), thanks to a trial-by-trial iterative process: the value of the chosen option is updated at each trial with a prediction-error, i.e.; the difference between the predicted outcome (the preceding Q-value), and the actual outcome (R): $\mathbf{q}\mathbf{v}_{t+1} = \mathbf{q}\mathbf{v}_t + \alpha \times (R - \mathbf{q}\mathbf{v}_t)$. The update is pondered by a parameter: the learning rate α , which quantifies how much individuals “learn” from their errors. The decision rule between two options A and B is often implemented as a soft-max (e.g. logistic) function:
$$p(A) = \frac{1}{1 + \exp(-\vartheta \times (\mathbf{q}\mathbf{v}(A) - \mathbf{q}\mathbf{v}(B)))}$$
. ϑ , the second model free-parameter, is called the temperature, and indexes the trade-off exploration/exploitation. We first simulate learning sequences with different values of α and R , and show that these two parameters are strongly associated with $\sigma(\mathbf{q}\mathbf{v})$ (**Figure.2.C.a**). However, R is rarely set as a free-parameter, and is usually set to the outcome value, despite potential individual differences in the sensitivity to the outcome magnitude. We therefore simulated choices occurring from learning sequences with a fixed ϑ but different values of α and R , and estimated a classical model, where only α and ϑ are free-parameters. In that case, we show that those two model parameters are strongly associated with $\sigma(\mathbf{q}\mathbf{v})$, despite the fact that ϑ is not a parameters governing the computation of $\mathbf{q}\mathbf{v}$, but rather governs the choice process (**Figure.2.C.b**).

Model-based fMRI and random-effects. Although

treating model free-parameters as random-effects often seem to provide the best account of individuals’ behavior as assessed by rigorous model-comparisons, a common practice in the literature is to treat them as a fixed-effect –i.e. use a population parameter- to generate the latent variables for fMRI analysis (Daw et al., 2006; Gershman et al., 2009; Gläscher et al., 2009, 2010; O’Doherty et al., 2004; Palminteri et al., 2009; Pessiglione et al., 2008). This is justified by the fact that individual free-parameter estimates are “noisy” and using the data from the full population is an efficient way to regularize them. However, when individual parameters still provide a better account of the population behavioral data according to rigorous model-comparison procedures, one might argue that the variance modeled in the individual free-parameters actually captures a true inter-individual variability in the cognitive process at stake, hence might contribute to give a better account of individual neurophysiological data. In the light of the scaling issue raised in this paper, we suggest that the use of population free-parameters actually constrains $\sigma(\mathbf{x}_k)$ to a unique population value, provided that individuals are given the same input. Under the *normalization* scaling hypothesis, this can substantially increase the statistical power of subsequent second-level random effects analyses. In this case, a better way to model brain activation (i.e. accounting for individual differences) would be to use individual model free-parameters, and Z-score the latent variables generated by these individual models. This discussion, again, raises the interest of better documenting scaling-laws in fMRI, so as to provide *a priori* principled rational to process independent variables of interest, in order to increase the sensitivity and replicability of model-based fMRI.

Model-based fMRI and inter-individual correlations. Another consequence of the associations between the model free-parameters and $\sigma(\mathbf{x}_k)$, is that inter-individual correlations between model free-parameters and activations ($\widehat{\beta}_k$) in an ROI encoding \mathbf{x}_k should be interpreted with much caution. Indeed, they may rely on simple mathematical dependencies between the free parameters, $\sigma(\mathbf{x}_k)$, and $\widehat{\beta}_k$, and they largely depend on interactions between the underlying brain-behavior scaling hypothesis and the processing (Z-scoring) of \mathbf{x}_k . In other words, observing e.g. a significant correlation between individual learning rates and individual Q-values *activations* in a given area may simply reflect the fact that differences in learning rate induced differences in the individual Q-value standard deviation. Besides, given that the link between model-parameters and $\sigma(\mathbf{x}_k)$ can reverse

depending on the task design (**Figure 2.B**), one can anticipate reports of opposite inter-individual correlations (positive or negative) between model-parameters (e.g. a discount factor or a learning rate) and the value “activations”, as measured by $\widehat{\beta}_k$.

Empirical testing of scaling laws

In this section, we assess the practical impact of BOLD-behavior scaling laws on fMRI analysis – random effects, and inter-individual correlations.

Scaling laws in value-rating fMRI studies. The first experimental data consists in three published fMRI datasets investigating “values” –the presumed determinant of decision-making (Camerer, 2008; Rangel et al., 2008). Functional neuroimaging measures were recorded while subjects were performing similar tasks (**Figure 3.A.a, B.a and C.a**): judging the pleasantness of pictures of paintings, houses and faces (Study 1), the desirability of objects depicted in short videos (Study 2) or the desirability of events described in sentences (Study 3), and reporting those evaluations on a rating scale (Lebreton et al., 2009, 2012, 2015).

As previously and extensively reported (Bartra et al., 2013; Clithero and Rangel, 2014; Peters and Büchel, 2010; Sescousse et al., 2013), we found that random-effect analyses on the parametric native (\mathbf{v}) and individually Z-scored values (\mathbf{v}^Z) independent variable “value rating” are very significant in a large ventral prefrontal region, including ventromedial prefrontal cortex (VMPFC) and medial orbitofrontal cortex (MOFC) ($P_{FWE} < 0.05$ **Figure 3.A.b, B.b and C.b**).

The mathematical derivations predict that the random effect should be more (resp. less) significant with \mathbf{v}^Z under the normalization (resp. proportional) hypothesis. We used two measures to assess the significance of the random effects: 1) VMPFC-k: the size of the VMPFC clusters ($p_{FWE-clu} < 0.05$, with a voxel-wise cluster-generating threshold $p_{UNC} < 0.001$), and 2) ROI-log(P): the negative log of the p-value of the random effects in an anatomical independent VMPFC ROI (one-sample t-test on the individual $\widehat{\beta}_k$ and $\widehat{\beta}_k^Z$ averaged over the voxels of the ROI). These two measures in the 3 datasets consistently indicated that using \mathbf{v}^Z produces more significant random-effects (**Table 1.a**). Next, in order to formally assess the two scaling laws in this context, we tested the relationships derived in the **Experimental procedure**, and linking “activation”, as measured by $\widehat{\beta}$ and the standard deviation of the value $\sigma(\mathbf{v})$. Again, a very consistent pattern emerged across all 3 studies: we found no significant positive correlations between $\widehat{\beta}_k^Z$ and $\sigma(\mathbf{v}_k)$, whereas

correlation between $\widehat{\beta}_k$ and $1/\sigma(\mathbf{v}_k)$ were systematically significantly positive (**Table 1.a and Figure 3.A.c, B.c and C.c**). Importantly, the same correlations assessed with the t-values (\widehat{t}_k^Z or \widehat{t}_k – which are practically identical: $R > .99$ in the three studies) did not exhibit the same pattern (**Table 1.a**). Hence, inter-individual differences in “activations”, as measured by $\widehat{\beta}_k$ are likely due to scaling issues rather than differences in the linear dependencies between the BOLD signal and the behavioral measure \mathbf{v} .

The results suggest that the inter-individual representation of values in the VMPFC, in such rating tasks, follows a normalization scaling rule: despite individual differences in the range (variance) of the behavioral value ratings, individuals exhibit similar range of BOLD signal in the core of the brain valuation system. This can be given concurrent interpretations: 1) the “true” underlying value signal range is similar across individuals - accurately captured by the fMRI analysis- despite individual differences in the behavior, e.g. due to calibration differences on the experimental rating scale; or 2) the underlying “true” value signal range is actually different across individuals –following the differences in the range of ratings reported on the experimental scale-, but there are experimental limitations which prevent the correct assessment of this inter-individual variability at the neural level. This raise new questions –e.g.: can we infer whether an option is more valuable to an individual than to another from fMRI data? -, whose answers will determine our ability to fulfill some of the promises of fMRI applications.

Scaling laws in a model-based study of reinforcement learning.

In this section, we illustrate the implications of scaling law for model-based fMRI, using a fourth experimental fMRI dataset investigating value-based learning (Palminteri et al., 2015). Participants were faced with repeated choices between abstract stimuli, which were probabilistically paired with different outcomes (neural, reward or punishment). The goal was to learn to select the stimuli, which maximize reward occurrences and minimize punishment occurrences (**Figure 4.A**). This task can be efficiently modelled with a variant of the Rescola-Wagner reinforcement-learning rule: participants learn, by trial and error, the value (Q-values) of the stimuli and make their choices by soft-maximizing expected value (see (Palminteri et al., 2015) and **Figure 4.B**). Two core free-parameters of the model capture the individuals’ learning dynamics and choice variance: the temperature θ , and the learning rate α . These free-parameters are typically set to maximize the likelihood of observed choices under the considered model.

Assessing inter-individual variability in brain-behavior relationship

As outlined in the first section, the implications of scaling laws for model-based fMRI root in the differences arising from using population (FFX) versus individual (RFX) sets of free-parameters. We therefore explored the consequences of those two modelling options in the behavioral data of (Palminteri et al., 2015). A model-comparison approach first unambiguously indicated that individual choices are more likely accounted for by individual free-parameters than by a single set of population free-parameters, even after accounting for the extra degrees of freedom (lack of parsimony) engendered by this procedure ($AIC_{RFX} = 8742$ vs $AIC_{FFX} = 9586$ and $BIC_{RFX} = 8960$ vs $BIC_{FFX} = 9615$). Importantly for the scaling-law issue, both individual free-parameters (θ and α) are very strongly associated with the individual standard deviation of the Q-value of the chosen option $\sigma(\mathbf{q}_{c_k})$ (inter-individual correlations, respectively $R = .65$; $p < .001$, and $R = .67$; $p < .001$, **Figure 4.C**), while being uncorrelated with one-another ($R = .31$, $p = .11$). Intuitively, α is positively associated to $\sigma(\mathbf{q}_{c_k})$ because it directly affects the amplitude of the learned values, whereas the θ affects the stochasticity of choices, such that subjects with higher θ_k frequently alternate between the best and the worst option, thus indirectly increasing the variance of \mathbf{q}_{c_k} .

Turning to neuroimaging data, we ran 3 GLM, differing only in the way the parametric regressor “chosen Q-value” (Q_C) was generated: we used population (\mathbf{q}_c^P) or individual model free-parameters. In this latter case, the Q_C could be entered in the GLM in their native scale (\mathbf{q}_c), or Z-scored per individual and session (\mathbf{q}_c^Z). We first ran a whole-brain random-effect analysis on the parametric regressor Q_C in the 3 GLMs (i.e. using \mathbf{q}_c , \mathbf{q}_c^P , or \mathbf{q}_c^Z as the independent variable). Replicating numerous findings, we found that Q_C are represented in the VMPFC (**Figure 4.D**). In order to assess the quality of individual fit, we extracted t-values \hat{t}_k in an anatomical VMPFC ROI. A random effect showed that these statistics are bigger when using Q_C generated with individual model free-parameters (\mathbf{q}_c or \mathbf{q}_c^Z) than group model free-parameters (\mathbf{q}_c^P) (one sided one-sample t-test, $t_{27} = 1.67$ $p = .05$). This means that, regardless of any scaling issue, the BOLD signal is better fitted with individual model free-parameters. This parallels the model comparison approach with the behavioral data, and might indicate that individual-fit lead to latent variables estimates, which are closer to the variables actually represented in subjects’ brains. We then compared the statistical significance of the random-effects in our 3 GLMs. First we noted that our whole-brain analysis resulted in a large and very significant cluster when using \mathbf{q}_c^P or \mathbf{q}_c^Z as the

independent variable, using \mathbf{q}_c only generated weak, sub-threshold activations (**Figure 4.D**). Then, paralleling the previous section, we used our two measures to assess the significance of the fMRI random effects (VMPFC-k and ROI-log(P)). These two measures gave similar conclusions: while using \mathbf{q}_c^P as an independent variable seems to improve random-effect models compared to using \mathbf{q}_c , using \mathbf{q}_c^Z provide the most significant random-effects, supporting the normalization hypothesis (**Table 1.b**, see also **Figure 4.E**).

In order to assess more specifically the scaling law in this new dataset, we next tested the correlation between $\widehat{\beta}_k^Z$ (i.e. computed with \mathbf{q}_c^Z) and $\sigma(\mathbf{q}_{c_k})$ and correlation between $\widehat{\beta}_k$ (i.e. computed with \mathbf{q}_c) and $1/\sigma(\mathbf{q}_{c_k})$ in the anatomical VMPFC ROI. Again, whereas the first correlation is not significantly positive the second is (**Table 1.b** and **Figure 4.F**). Importantly, the same correlations assessed with t-values (\hat{t}_k^Z and \hat{t}_k) did not exhibit the same pattern. Hence, inter-individual differences in Q_C representations, as measured by $\widehat{\beta}_k$ are likely due to scaling issues (namely normalization) rather than differences in the linear dependencies between the BOLD signal and the behavioral measure \mathbf{q}_c . Note that one subject has α and θ close to 0. Due to the fact that $\widehat{\beta}_k$ correlate with $1/\sigma_{X,k}$, and that close-to-zero model free-parameters generate latent variables X_k with close-to-zero $\sigma_{X,k}$, the corresponding $\widehat{\beta}_k$ of this subject take an outlying value. Importantly, however, our test of scaling law holds when excluding the potential outlier ($P < .05$, **Figure 4.F**) again strongly favoring the normalization hypothesis.

Altogether, our results suggest that scaling issues might explain the apparent contradictory observations that fMRI random-effects are more significant using \mathbf{q}_c^P than \mathbf{q}_c despite the superiority of individually-fitted models to account for individual behavioral choices. Overall the results of these analyses advocate for the use of individual parameters in value-related model-based fMRI, together with a Z-scoring of the model-estimated latent variable *-value-* to account at best for the inter-individual normalization effect occurring in the VMPFC.

Finally, we assessed the correlation between model free-parameters and VMPFC activations –as measure with $\widehat{\beta}_k$. The rational would be to take the model free-parameters as a traits-of-interest, to support a statement like: “individual who are better learners –i.e. higher learning-rates- have a stronger value-related activations in the VMPFC”. In our case, model free-parameters positively correlate with $\sigma(\mathbf{q}_{c_k})$, and the normalization scaling law implies that $\widehat{\beta}_k$ scale with $1/\sigma(\mathbf{q}_{c_k})$. We therefore expect

individual learning-rates and soft-max-temperatures to be negatively correlated to $\widehat{\beta}_k$ in the VMPFC. Experimental data support this prediction (respectively $R=-.44$, $P<.05$ and $R=-.38$, $P<.05$). Excluding the outlier precluded the significance of the correlation with the learning rate, but not with the soft-max temperature (respectively $R=-.24$, $P=.22$ and $R=-.39$, $P<.05$). This analysis justifies our cautious note about the interpretations of correlations between model-parameters and activations –as measured by $\widehat{\beta}_k$ –, as they may be dependent on the statistical relationship between model-parameters and $\sigma(\mathbf{q}_{c_k})$.

DISCUSSION

Researchers are increasingly interested in inter-individual variability in cognitive neurosciences, in the normal and pathological ranges. The ability to assess and predict individual differences from neural measures –neuromarkers- has emerged as one the most promising application of fMRI in society (Gabrieli et al., 2015; Wang and Krystal, 2014) In this manuscript, we explored a specific type of neuromarker: task-dependent fMRI “activations”, indexed by unstandardized coefficients of regression $\widehat{\beta}_k$ between individual behavioral variables (\mathbf{x}_k) and BOLD signal (Y_k). We recalled that $\widehat{\beta}_k$ depend on the ratio of the standard deviation of \mathbf{x}_k and Y_k , making task-dependent fMRI neuromarkers partly reflecting *scaling laws* between the BOLD signal and the behavioral variable of interest. Documenting those scaling laws is therefore paramount to correctly interpret assessments of inter-individual differences in cognitive neuroscience. With this goal in mind, we proposed a new taxonomy –proportional/normalization- to qualify such inter-individual brain-behavior scaling relationship. Importantly, this taxonomy is based on a formalized description of the statistical dependency between the BOLD signal and the behavioral variable, rather than on a functional (over-)interpretation of such statistical quantities –like in the current efficiency vs. activation taxonomy- (Poldrack, 2015). By doing so, it aims at helping the building of a cumulative cognitive science, based on the falsification of precise predictions. Although we acknowledge that the present paper does not cover the full range of potential link between brain activation and behavior, we think that this new perspective might contribute to reconcile previous contradictory findings, and foster a fruitful discussion on the way to interpret and assess investigations of individual-difference in neuroimaging.

Notably, the level of interpretation of the scaling relationships may depend on the extent to which state-of-the art fMRI technics can capture

inter-individual variations in the range of BOLD activations $\sigma(Y_k)$. Indeed, the *proportional* scaling hypothesis can only be supported if it is possible to link inter-individual variations in the extent of the explanatory variable $\sigma(\mathbf{x}_k)$ to inter-individual variations in the range of BOLD activations $\sigma(Y_k)$; failing to do so will *de facto* provide evidence for the *normalization* hypothesis, regardless of the underlying neuro-cognitive scaling hypothesis. However, this should not be viewed as a failure of the proposed framework, which aims at better describing inter-individual brain-behavior relationships and whose validity is therefore independent of the level of description considered; Still, the ability to reliably assess inter-individual differences in the range of BOLD activations $\sigma(Y_k)$ appears critical, notably to provide evidence in favor of the *proportional* scaling laws. Although fMRI reliability has been the focus of extensive research (Bennett and Miller, 2010), this specific question has received little attention so far. Overall, besides the sensitivity of MRI measures to inter-individual differences, numerous factors are suspected to play a role in our ability to correctly estimate inter-individual differences in the range of BOLD activations $\sigma(Y_k)$, such as e.g. individual differences in vascularization (Logothetis, 2008), s) or preprocessing and analytic strategies.

We next attempted to decipher the consequences of scaling laws on random-effect models, inter-individual differences analyses, and computational modelling. We particularly stressed the fact that, in event-related parametric designs, differences in $\widehat{\beta}_k$ –which are used to quantify brain “activations”- can trivially derive from differences in the range of the behavioral measure $\sigma(\mathbf{x}_k)$. This is the case under the normalization hypothesis when neuroimaging data are analyzed with native behavioral variables \mathbf{x}_i , and under the proportional hypothesis, when neuroimaging data are analyzed with Z-scored behavioral data. This poses serious challenges to current interpretations of inter-individual differences in $\widehat{\beta}_k$ in the absence of knowledge about the underlying scaling law. This warning is not limited to inter-individual differences claims, but generalizes to within-individual inter-session claims: between-sessions scaling laws may have important consequences for design involving between-session manipulations –such as brain stimulation or pharmacological modulations- which often also impact the behavior, hence $\sigma(\mathbf{x}_k)$. Therefore, similarly to the inter-individual; case, we recommend that researchers start documenting the impact of their between-session manipulation on $\sigma(\mathbf{x}_k)$, provide a more detailed account of the processing of independent variables \mathbf{x}_k (Z-scoring or not) used for neuroimaging analysis, and investigate scaling-laws. In the value-based

Assessing inter-individual variability in brain-behavior relationship

community, a recent paper reported that normalization of BOLD signal under different value ranges actually occurs between different sessions of the same individual. (Cox and Kable, 2014), giving credit to our warning.

Although we introduced the scaling issue in the context of fMRI event-related parametric designs, most conclusions and warnings can be extended to fMRI categorical designs, notably in the critical situation where categorical events are constructed from individual reaction times. Indeed, modelling categorical events with individual time-varying boxcars introduces an inter-individual difference in the modelling of BOLD-signal amplitude, therefore inevitably generates scaling issue. (Grinband et al., 2008; Poldrack, 2015).

We propose that a good practice before engaging in the study of fMRI inter-individual variability is to start documenting the statistical relationship between traits of interest (individual clinical scores, psycho-social measures, model free-parameters) and the standard deviation $\sigma(x_k)$ of fMRI regressors. Ideally, researchers might explicitly test brain-behavior scaling laws for the cognitive function of interest, in the brain region of interest, using their specific task – indeed, one can expect that different cognitive processes, elicited with different tasks could follow different scaling law, in different brain regions. In order to improve the reproducibility of fMRI findings, it is paramount to formulate clear a priori hypothesis about inter-individual-differences and to use an appropriate operationalization.

Finally, we initiated this practice by documenting inter-individual normalization of values representation in the VMPFC using four datasets. This parallels recent findings reporting within-individual range adaptation of value coding in the same area (Cox and Kable, 2014; Padoa-Schioppa, 2009). This finding might contribute to improve our understanding of the valuation process, and provide principled rational to preprocess variables of interest and carry out model-based fMRI in the value-based decision-making community.

Acknowledgments:

We are thankful to M. Pessiglione (MP) and G. Coricelli (GC) for granting us unlimited and unrestricted access to the fMRI datasets. These datasets were collected thanks to the European Research Council (ERC Starting Grant BioMotiv to MP and ERC Consolidator Grant 617629 to GC), a Research Grant from the Schlumberger Foundation to MP, and an Agence National de la Recherche (ANR-11-EMCO-010) grant to GC. ML is supported by an EU Marie Skłodowska-Curie Individual

Fellowship (IF-2015 Grant 657904), a Universiteit van Amsterdam – Amsterdam Brain and Cognition Talent Grant, and acknowledge the support of the Bettencourt-Schueller Foundation. SP is also supported by an EU Marie Skłodowska-Curie Individual Fellowship (PIEF-GA-2012 Grant 328822).

References

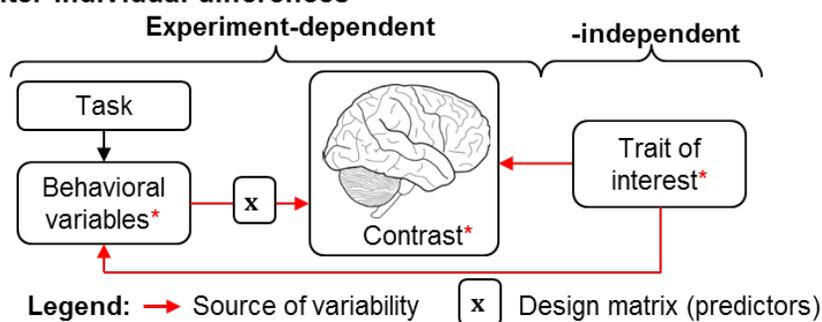
- Ainslie, G., and Haslam, N. (1992). Hyperbolic discounting. In *Choice over Time*, G. Loewenstein, and J. Elster, eds. (New York, NY, US: Russell Sage Foundation), pp. 57–92.
- Bartra, O., McGuire, J.T., and Kable, J.W. (2013). The valuation system: A coordinate-based meta-analysis of BOLD fMRI experiments examining neural correlates of subjective value. *NeuroImage* 76, 412–427.
- Beckmann, C.F., Jenkinson, M., and Smith, S.M. (2003). General multilevel linear modeling for group analysis in fMRI. *Neuroimage* 20, 1052–1063.
- Bennett, C.M., and Miller, M.B. (2010). How reliable are the results from functional magnetic resonance imaging? *Ann. N. Y. Acad. Sci.* 1191, 133–155.
- Bernoulli, D. (1954). Exposition of a New Theory on the Measurement of Risk. *Econometrica* 22, 23–36.
- Braver, T.S., Cole, M.W., and Yarkoni, T. (2010). Vive les differences! Individual variation in neural mechanisms of executive control. *Curr. Opin. Neurobiol.* 20, 242–250.
- Camerer, C.F. (2008). *Neuroeconomics: Opening the Gray Box*. *Neuron* 60, 416–419.
- Clithero, J.A., and Rangel, A. (2014). Informatic parcellation of the network involved in the computation of subjective value. *Soc. Cogn. Affect. Neurosci.* 9, 1289–1302.
- Cohen, J., Cohen, P., West, S.G., and Aiken, L.S. (2013). *Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences* (Routledge).
- Cox, K.M., and Kable, J.W. (2014). BOLD Subjective Value Signals Exhibit Robust Range Adaptation. *J. Neurosci.* 34, 16533–16543.
- Daw, N.D. (2011). Trial-by-trial data analysis using computational models. *Decis. Mak. Affect Learn. Atten. Perform.* XXIII 23, 3–38.
- Daw, N.D., O'Doherty, J.P., Dayan, P., Seymour, B., and Dolan, R.J. (2006). Cortical substrates for exploratory decisions in humans. *Nature* 441, 876–879.
- Deichmann, R., Gottfried, J.A., Hutton, C., and Turner, R. (2003). Optimized EPI for fMRI studies of the orbitofrontal cortex. *NeuroImage* 19, 430–441.
- Friston, K.J., Holmes, A.P., Worsley, K.J., Poline, J.-P., Frith, C.D., and Frackowiak, R.S.J. (1994). Statistical parametric maps in functional imaging: A general linear approach. *Hum. Brain Mapp.* 2, 189–210.
- Friston, K.J., Stephan, K.E., Lund, T.E., Morcom, A., and Kiebel, S. (2005). Mixed-effects and fMRI studies. *Neuroimage* 24, 244–252.
- Gabrieli, J.D.E., Ghosh, S.S., and Whitfield-Gabrieli, S. (2015). Prediction as a Humanitarian and Pragmatic Contribution from Human Cognitive Neuroscience. *Neuron* 85, 11–26.
- Gershman, S.J., Pesaran, B., and Daw, N.D. (2009). Human Reinforcement Learning Subdivides Structured Action Spaces by Learning Effector-Specific Values. *J. Neurosci.* 29, 13524–13531.
- Gläscher, J., Hampton, A.N., and O'Doherty, J.P. (2009). Determining a Role for Ventromedial Prefrontal Cortex in Encoding Action-Based Value Signals During Reward-Related Decision Making. *Cereb. Cortex* 19, 483–495.
- Gläscher, J., Daw, N., Dayan, P., and O'Doherty, J.P. (2010). States versus Rewards: Dissociable Neural Prediction Error Signals Underlying Model-Based and Model-Free Reinforcement Learning. *Neuron* 66, 585–595.
- Gold, S., Christian, B., Arndt, S., Zeien, G., Cizadlo, T., Johnson, D.L., Flaum, M., and Andreasen, N.C. (1998). Functional

Assessing inter-individual variability in brain-behavior relationship

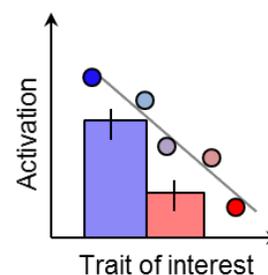
- MRI statistical software packages: a comparative analysis. *Hum. Brain Mapp.* 6, 73–84.
- Grinband, J., Wager, T.D., Lindquist, M., Ferrera, V.P., and Hirsch, J. (2008). Detection of time-varying signals in event-related fMRI designs. *NeuroImage* 43, 509–520.
- Holmes, A.P., and Friston, K.J. (1998). Generalisability, Random Effects & Population Inference. *Neuroimage* 7, S754.
- Lebreton, M., Jorge, S., Michel, V., Thirion, B., and Pessiglione, M. (2009). An Automatic Valuation System in the Human Brain: Evidence from Functional Neuroimaging. *Neuron* 64, 431–439.
- Lebreton, M., Kawa, S., d'Arc, B.F., Daunizeau, J., and Pessiglione, M. (2012). Your Goal Is Mine: Unraveling Mimetic Desires in the Human Brain. *J. Neurosci.* 32, 7146–7157.
- Lebreton, M., Abitbol, R., Daunizeau, J., and Pessiglione, M. (2015). Automatic integration of confidence in the brain valuation signal. *Nat. Neurosci.* 18, 1159–1167.
- Logothetis, N.K. (2008). What we can do and what we cannot do with fMRI. *Nature* 453, 869–878.
- Morgan, V.L., Dawant, B.M., Li, Y., and 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.
- Mumford, J.A., and Nichols, T. (2009). Simple group fMRI modeling and inference. *Neuroimage* 47, 1469–1475.
- O'Doherty, J., Dayan, P., Schultz, J., Deichmann, R., Friston, K., and Dolan, R.J. (2004). Dissociable Roles of Ventral and Dorsal Striatum in Instrumental Conditioning. *Science* 304, 452–454.
- O'Doherty, J.P., Hampton, A., and Kim, H. (2007). Model-Based fMRI and Its Application to Reward Learning and Decision Making. *Ann. N. Y. Acad. Sci.* 1104, 35–53.
- Padoa-Schioppa, C. (2009). Range-Adapting Representation of Economic Value in the Orbitofrontal Cortex. *J. Neurosci.* 29, 14004–14014.
- Palminteri, S., Boraud, T., Lafargue, G., Dubois, B., and Pessiglione, M. (2009). Brain Hemispheres Selectively Track the Expected Value of Contralateral Options. *J. Neurosci.* 29, 13465–13472.
- Palminteri, S., Khamassi, M., Joffily, M., and Coricelli, G. (2015). Contextual modulation of value signals in reward and punishment learning. *Nat. Commun.* 6.
- Pessiglione, M., Petrovic, P., Daunizeau, J., Palminteri, S., Dolan, R.J., and Frith, C.D. (2008). Subliminal Instrumental Conditioning Demonstrated in the Human Brain. *Neuron* 59, 561–567.
- Peters, J., and Büchel, C. (2010). Neural representations of subjective reward value. *Behav. Brain Res.* 213, 135–141.
- Poldrack, R.A. (2015). Is “efficiency” a useful concept in cognitive neuroscience? *Dev. Cogn. Neurosci.* 11, 12–17.
- Rangel, A., Camerer, C., and Montague, P.R. (2008). A framework for studying the neurobiology of value-based decision making. *Nat. Rev. Neurosci.* 9, 545–556.
- Sescousse, G., Caldú, X., Segura, B., and Dreher, J.-C. (2013). Processing of primary and secondary rewards: a quantitative meta-analysis and review of human functional neuroimaging studies. *Neurosci. Biobehav. Rev.* 37, 681–696.
- Wang, X.-J., and Krystal, J.H. (2014). Computational Psychiatry. *Neuron* 84, 638–654.
- Woolrich, M.W., Behrens, T.E., Beckmann, C.F., Jenkinson, M., and Smith, S.M. (2004). Multilevel linear modelling for FMRI group analysis using Bayesian inference. *Neuroimage* 21, 1732–1747.
- Worsley, K.J., Liao, C.H., Aston, J., Petre, V., Duncan, G.H., Morales, F., and Evans, A.C. (2002). A general statistical analysis for fMRI data. *Neuroimage* 15, 1–15.
- Yarkoni, T., and Braver, T.S. (2010). Cognitive Neuroscience Approaches to Individual Differences in Working Memory and Executive Control: Conceptual and Methodological Issues. In *Handbook of Individual Differences in Cognition*, A. Gruszka, G. Matthews, and B. Szymura, eds. (Springer New York), pp. 87–107.

Figure 1

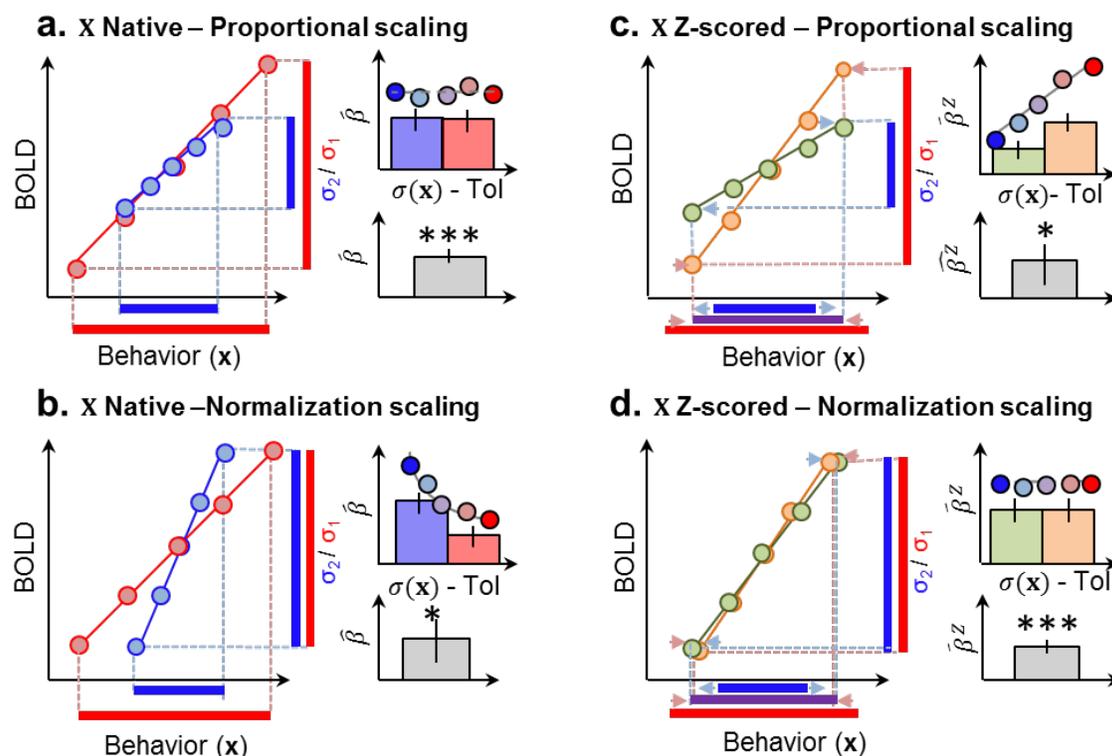
A. Inter-individual differences



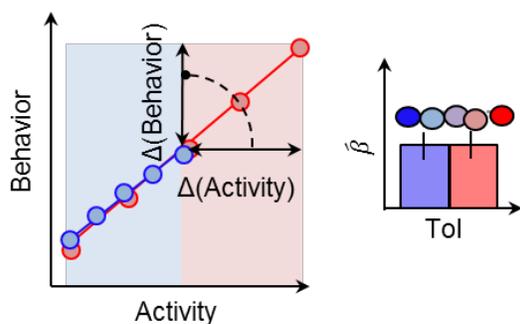
B. Typical result



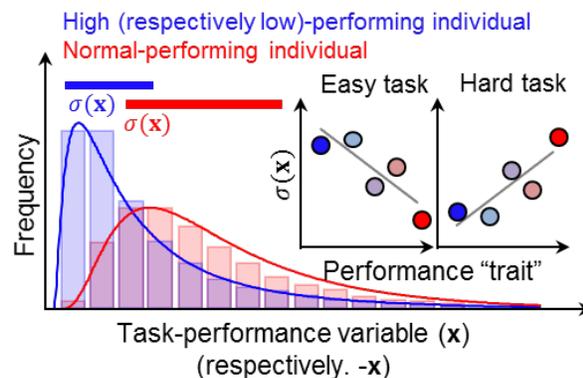
C. The relationship between scaling laws and second-level statistics



D. Interpreting (in-)differences in $\hat{\beta}$



E. Task design, behavior, and traits of interest



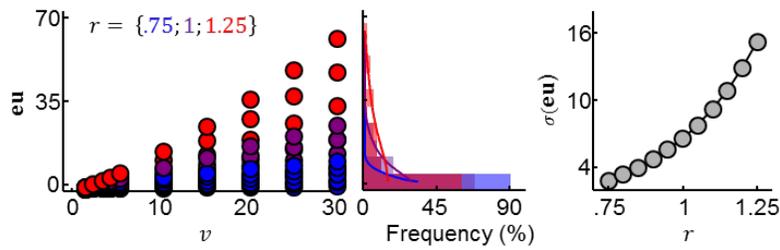
Assessing inter-individual variability in brain-behavior relationship

Figure 1: Inter-individual in brain-behavior relationships and scaling laws

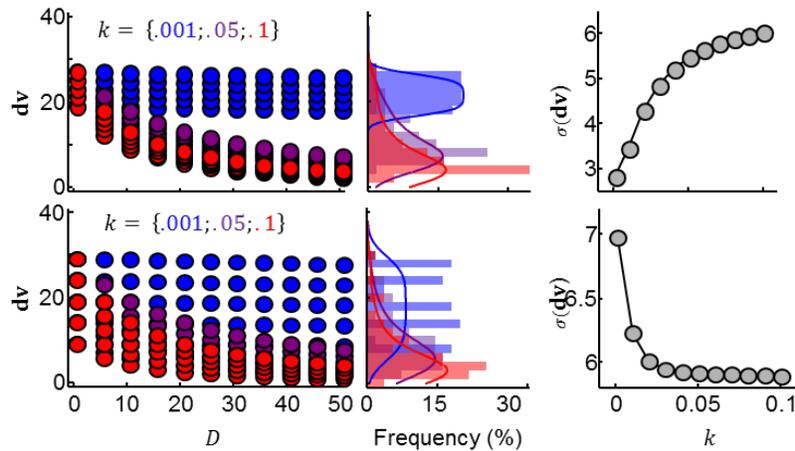
A. Inter-individual differences in brain-behavior relationships are paramount in virtually all fMRI designs, regardless of whether they address these differences explicitly or not. They are typically assessed by linking individual fMRI contrast values with a trait of interest. However they also often leak in fMRI data analysis when behavioral variables are used as independent variables in the design matrix. Problems can arise when the trait of interest also generates differences in the variance of the behavioral variables. **B.** Typical results illustrating inter-individual differences in brain-behavior relationships, as they are reported in the literature, in the form of categorical contrast between group of subjects (bars \pm s.e.m) or group level linear correlation with a continuous trait (dots and regression line). **C.** The relationship between scaling laws and second-level statistics. The panels describe the impact of inter-individual differences in the standard deviation $\sigma(\mathbf{x})$ of the behavioral variable \mathbf{x} , under different preprocessing -native (a,b) or Z-scored variable (c,d)- and under different scaling-laws hypotheses -proportional (a,c) and normalization (a,d). Each sub panel contains three graphs. On the left we illustrate how \mathbf{x} is related to the BOLD signal in two individuals with different initial $\sigma(\mathbf{x})$ (blue vs. red or green vs. orange). The individual unstandardized coefficients of regression $\hat{\beta}$, corresponds the slope of the corresponding lines. On the upper right corner, we illustrate the statistical relations between individual brain activations and $\sigma(\mathbf{x})$ (presumably linked to the trait of interest T_{ol}) as a between-group analysis (histograms) or continuous inter-individual correlation (dots). On the bottom-right corner we illustrate the consequences for the significance of second-level random-effect analysis (*: lower significance, vs. ***: higher significance). $\hat{\beta}$ and $\hat{\beta}^2$ respectively refer to fMRI unstandardized coefficients of regressions computed with a native scaling or an individual Z-scoring of the parametric regressor \mathbf{x} . **D.** Interpreting differences in unstandardized $\hat{\beta}$. Consider a brain region, which causally and proportionally causes a behavior (i.e. the more activation, the higher the behavioral variable, within and across subjects). In case a trait of interest (e.g. pathology) directly impacts the range of activation of this region (e.g. due to degeneration), this cannot be assessed/detected by differences in $\hat{\beta}$ (right, inset). Misunderstanding the signification of $\hat{\beta}$ and ignoring scaling laws can lead to erroneous negative conclusions. **E.** Linking task design, behavior, and traits of interest. The red/blue histograms depict the distribution of a behavioral variable (reaction time, decision value, confidence) in two individual with variable performance. If the task is easy (respectively difficult), the high (respectively low)-performing individual can exhibit a ceiling (respectively floor) effect on performance. This creates statistical dependencies between the performance trait performance) and the standard deviation ($\sigma(\mathbf{x})$) of the behavioral task-performance variable (see graphical insets). Given the dependencies between fMRI $\hat{\beta}$ and $\sigma(\mathbf{x})$, this can lead to opposite correlations between “activations” –as measured by β -, and “individual performances”, proxied by $\sigma(\mathbf{x})$.

Figure 2

A. Utility - $eu = p \cdot v^r$

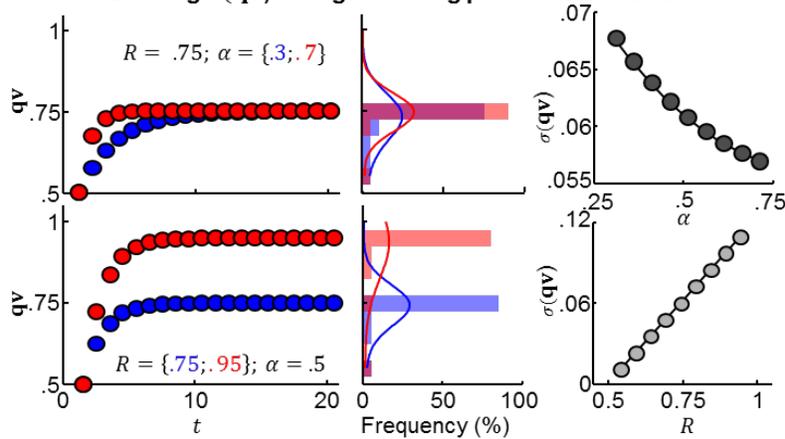


B. Delay Discounting - $dv = v / (1 + k \times D)$



C. Reinforcement learning - $qv_{t+1} = qv_t + \alpha \times (R - qv_t)$

a. Linking $\sigma(qv)$ and generating parameters α and R



b. Linking $\sigma(qv)$ and estimated parameters $\hat{\alpha}$ and $\hat{\theta}$

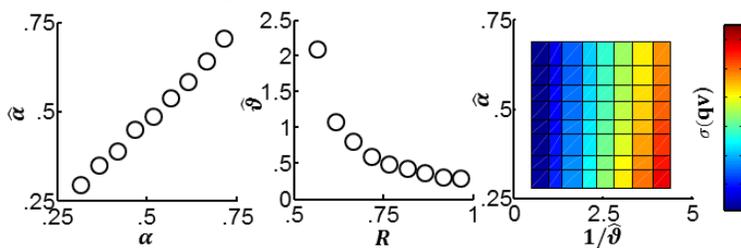
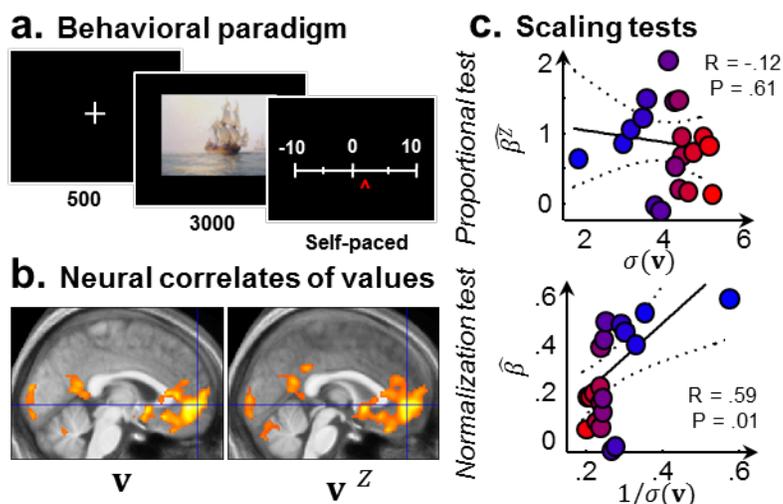


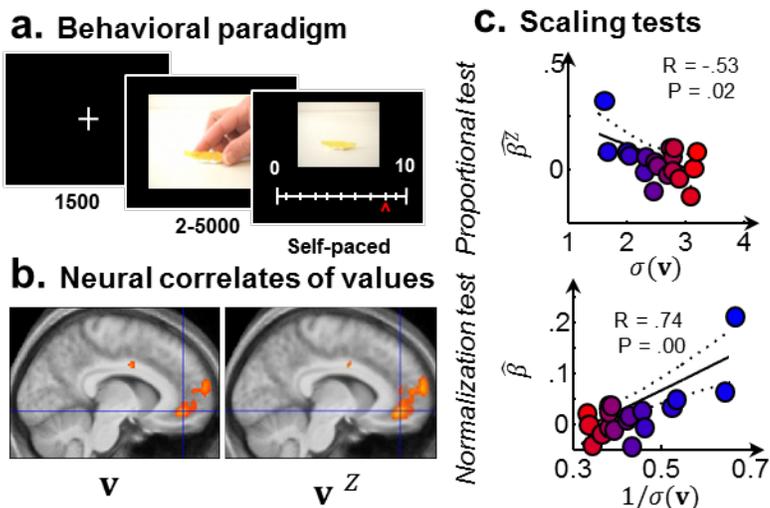
Figure 2: Model-based fMRI and the impact of model-parameters on inter-individual differences

A. Linking individual model-parameters and $\sigma(\mathbf{v})$. We computed the expected utility $\mathbf{eu} = p \times g^r$ of different prospects mixing gains g and probabilities p - $g \in \{1:5:10:5:30\}\text{€}$, $p \in \{10:20:90\}\%$ -, for different values of the utility curvature parameter r (blue: $r = .75$, purple: $r = 1.0$ and red: $r = 1.25$). Left: Expected utility of the task-stimuli, Middle: histogram of computed expected-utility and corresponding scaled density functions. Right: estimated $\sigma(\mathbf{eu})$ as a function of r . **B.** Linking individual model-parameters, task design and $\sigma(\mathbf{v})$. We computed the discounted value $\mathbf{dv} = \frac{g}{1+k \times D}$ of different prospects mixing gains g and delays D , for different values of the discount parameter k (blue: $k = .001$, purple: $k = .04$ and red: $k = .1$), and for 2 task designs (i.e. 2 sets of prospects). Top: $g \in \{20:2:28\}\text{€}$, $D \in \{0:5:50\}$ Bottom: $g \in \{10:5:30\}\text{€}$, $D \in \{0:5:50\}$. Left: Discounted value of the task-stimuli. Middle: histogram of computed expected-utility and corresponding scaled density functions. Right: estimated $\sigma(\mathbf{dv})$ as a function of k . **C.** Illustrating the link between non-value related model-parameters and $\sigma(\mathbf{v})$. **a.** We computed the Q-value $\mathbf{qv}_{t+1} = \mathbf{qv}_t + \alpha \times (R - \mathbf{qv}_t)$ of a 20-trials learning sequence, for different learning rates α (Top, $R = .75$, and blue: $\alpha = .3$ or red: $\alpha = .7$) and different outcome magnitude R (Bottom; $\alpha = .5$, and blue: $R = .75$, and red: $R = 1$) Left: Q-values as a function of the trial number. Middle: histogram of computed Q-values and corresponding scaled density functions. Right: estimated $\sigma(\mathbf{qv})$ as a function of α (top) or R (bottom). **b.** We simulated a task, implementing binary choices between a fixed-option of known value (0.5), and an option whose value had to be learned through trial-and error. We generated sequences of Qv of the unknown option, with different values of α - $\alpha \in \{.3:.05:.7\}$ - and R $R \in \{.55:.05:.95\}$ -, and corresponding stochastic choices. For each set of parameters, we generated 50 learning-sequences of 20 trials. We then estimated the parameter of the model, but with a fixed $R (=1)$ and a softmax (logistic) choice function $p(A) = \frac{1}{1 + \exp(-\vartheta \times (\mathbf{qv}(A) - \mathbf{qv}(B)))}$, setting α and ϑ as the model free parameters. Left: estimated $\hat{\alpha}$ as a function of the true α . Middle: estimated $\hat{\vartheta}$ as a function of the true R . Right: estimated $\sigma(\mathbf{qv})$ as a function of $\hat{\alpha}$ and $\hat{\vartheta}$.

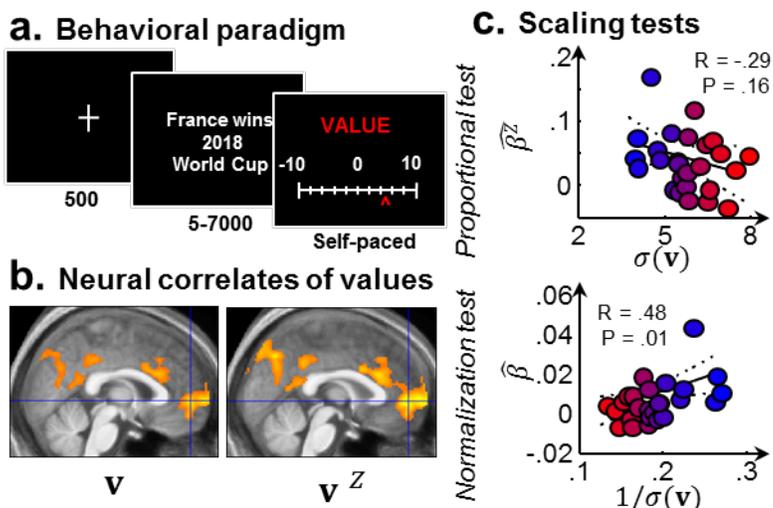
A. Study 1: Lebreton, et al. 2009



B. Study 2: Lebreton, et al. 2012



C. Study 3: Lebreton, et al. 2015

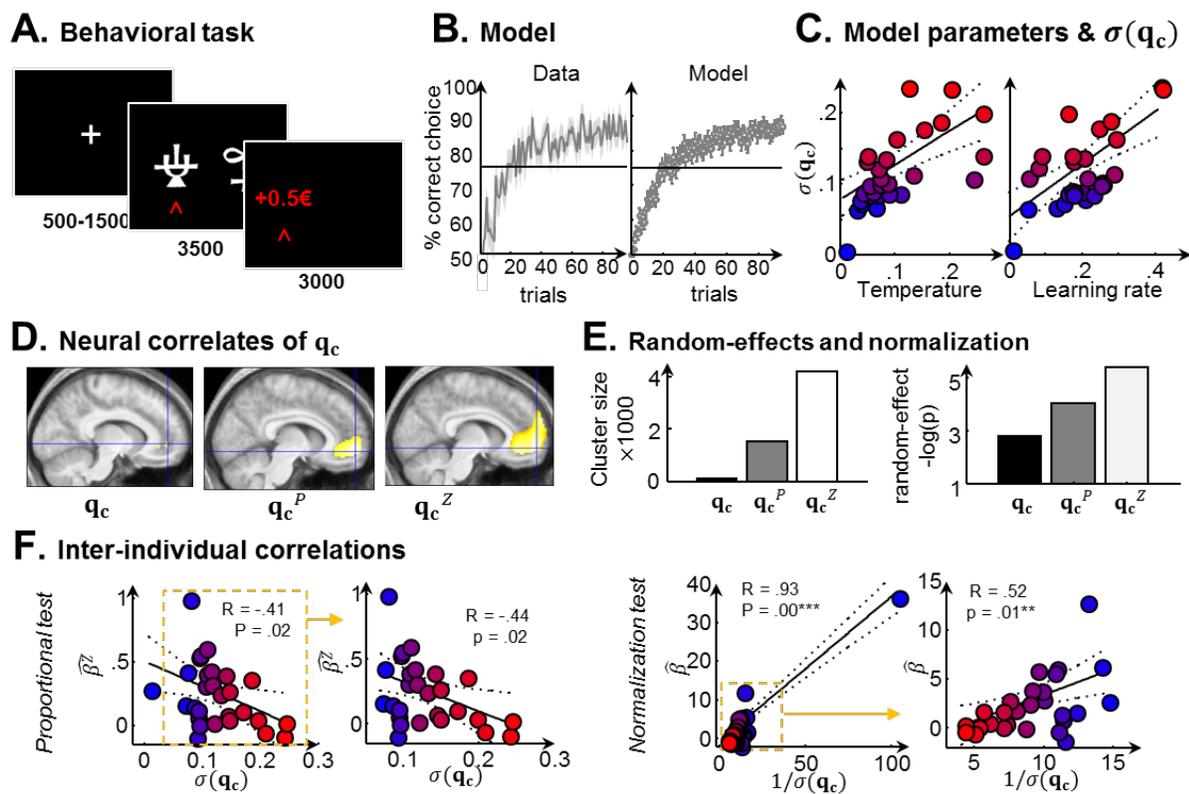


Assessing inter-individual variability in brain-behavior relationship

Figure 3: Assessing BOLD-Behavior scaling relationships in value rating tasks

We used 3 published datasets: **A.** Study 1: Lebreton, et al. 2009; **B.** Study 2: Lebreton, et al. 2012; **C.** Study 3: Lebreton, et al. 2015 (right); **a.** Rating task. Successive screens displayed in one trial are shown from left to right, with duration in milliseconds. **b.** Group-level neural correlates of values, using native values v (left) or individually Z-scored values v^z (right). The color code on glass brains (left maps) and sagittal slices (right) indicate the statistical significance of clusters that survived the whole-brain family-wise error (FWE) correction for multiple comparisons, computed at the cluster level ($P_{FWE} < .05$, with a voxel-wise cluster-generating threshold $P_{UNCORR} < .001$). **c.** Inter-individual correlations between value rating standard deviation $\sigma(v)$ and unstandardized coefficients of regressions estimated using Z-scored values $\hat{\beta}^z$ (top), and between the inverse of the value rating standard deviation ($1/\sigma(v)$) and unstandardized coefficients of regressions estimated using native-scored values $\hat{\beta}$ (bottom). Solid lines indicate the best linear fit, and dotted lines the 95% confidence interval. The dots color codes the relative $\sigma(v)$, from low (blue) to high (red).

Figure 4



Assessing inter-individual variability in brain-behavior relationship

Figure 4: Assessing BOLD-Behavior scaling relationships in a learning tasks

A. Behavioral task. Successive screens displayed in one trial are shown from left to right, with duration in milliseconds (Palminteri et al., 2015). **B.** Average behavior (left) and model fit (right), displayed as the correct choice rate. Shaded area (left) and error bars (right) correspond to $\text{mean} \pm \text{sem}$. **C.** Statistical relations between model free-parameters - temperature θ (left) and learning rate α (right)- and the standard deviation of the chosen Q-values $\sigma(\mathbf{q}_c)$. Solid line indicate the best linear fit, and dotted line the 95% confidence interval. The dots color codes the relative $\sigma(\mathbf{q}_c)$ value, from low (blue) to high (red). **D.** Group-level neural correlates of chosen Q-values. The color code on glass brains (left maps) and sagittal slices (right) indicate the statistical significance of clusters that survived the whole-brain family-wise error (FWE) correction for multiple comparisons, computed at the cluster level ($P_{\text{FWE}} < .05$, with a voxel-wise cluster-generating threshold $P_{\text{UNCORR}} < .001$), except for \mathbf{q}_c , where activations are presented at a more lenient threshold (voxel-wise threshold $P_{\text{UNCORR}} < .001$) **E.** Random-effects comparisons. VMPFC k refers to the size of the VMPFC cluster (in voxels, cluster-generating voxel threshold $P_{\text{UNCORR}} < .001$); ROI $-\log(P)$ refers to the negative logarithm of the P-value of a random effect analysis performed on the individual averaged coefficient of regression extracted from an anatomical VMPFC ROI. **E.** Scaling law assessment: inter-individual correlations between value rating standard deviation $\sigma(\mathbf{v})$ and unstandardized coefficients of regressions estimated using Z-scored values $\hat{\beta}^Z$ (left), and between the inverse of the value rating standard deviation ($1/\sigma(\mathbf{v})$) and unstandardized coefficients of regressions estimated using native-scored values $\hat{\beta}$ (right). The second plot presents the same correlation(s), excluding the potential outlier. Solid lines indicate the best linear fit, and dotted lines the 95% confidence interval. The dots color codes the relative $\sigma(\mathbf{v})$, from low (blue) to high (red). \mathbf{q}_c^P , \mathbf{q}_c and \mathbf{q}_c^Z indicate that the fMRI GLMs are designed with the variable chosen Q-values generated with population (\mathbf{q}_c^P) or individual (\mathbf{q}_c and \mathbf{q}_c^Z) model free parameters, and using a native scaling (\mathbf{q}_c^P and \mathbf{q}_c) or an individual Z-scoring (\mathbf{q}_c^Z) of the variable. $\hat{\beta}$ and $\hat{\beta}^Z$ respectively refer to fMRI unstandardized coefficients of regressions computed with a native scaling (\mathbf{q}_c) or an individual Z-scoring (\mathbf{q}_c^Z) of the parametric regressor \mathbf{Q}_c . $\sigma(\mathbf{q}_c)$ refers to the standard deviation of the native \mathbf{Q}_c . *: $P < .05$; ***: $P < .001$;

Assessing inter-individual variability in brain-behavior relationship

		Random effects		Correlations: R (P)			
		VMPFC k	ROI -log(P)	$\widehat{\beta}^Z \propto \sigma(\mathbf{v})$	$\widehat{\beta} \propto 1/\sigma(\mathbf{v})$	$\widehat{t}^Z \propto \sigma(\mathbf{v})$	$\widehat{t} \propto 1/\sigma(\mathbf{v})$
a. Rating studies							
Lebreton, et al. 2009	\mathbf{v}	3301	5.33	-	.59 (.01**)	-	.08 (.74)
	\mathbf{v}^Z	4442	5.82	-.12 (.61)	-	-.13 (.48)	-
Lebreton, et al. 2012	\mathbf{v}	586	1.80	-	.74 (.00***)	-	.50 (.03*)
	\mathbf{v}^Z	784	2.14	-.53 (.02)	-	-.44 (.07)	-
Lebreton, et al. 2015	\mathbf{v}	2116	3.87	-	.48 (.01*)	-	.14 (.49)
	\mathbf{v}^Z	2606	4.53	-.29 (.16)	-	-.12 (.57)	-
b. Learning study							
Palminteri, et al. 2015	q_c	123	2.18	-	.93 (.00***)	-	.25 (.21)
	q_c^P	1547	3.73	-	-	-	-
	q_c^Z	4189	5.37	-.41 (.02)	-	-.41 (.02)	-

Table 1.

VMPFC k refers to the size of the VMPFC cluster (in voxels, cluster-generating voxel threshold $P_{\text{UNCORR}} < .001$); ROI $-\log(P)$ refers to the negative logarithm of the P-value of a random effect analysis performed on the individual averaged coefficient of regression extracted from an anatomical VMPFC ROI. $\widehat{\beta}$ and $\widehat{\beta}^Z$ respectively refer to fMRI unstandardized coefficients of regressions computed with \mathbf{v} and \mathbf{v}^Z . \widehat{t} refers to fMRI t-statistics derived from $\widehat{\beta}$. $\sigma(\mathbf{v})$ refers to the standard deviation of the native value-rating measure. *: $P < .05$; **: $P < .01$; ***: $P < .001$ one-sample t-test.

EXPERIMENTAL PROCEDURES

Regression coefficients in fMRI

At the first level, in each individual, , we can write, for each independent variable x_i , $\hat{\beta}_i = \frac{\sigma(Y)}{\sigma(x_i)} \rho(x_i, Y) \sqrt{VIF_i}$ (1) (Cohen et al., 2013). Here, $\rho(x_i, Y)$ is the semi-partial correlation between x_i and Y , and $VIF_i = \frac{1}{1-R^2_{x_i x_j}}$, (with $R^2_{x_i x_j}$ indexing the variance explained by a regression with x_i as a dependent variable, and all $x_j, j \neq i$ as independent variables) is the Variance Inflation Factor, which quantifies $\hat{\beta}_i$ over-estimation due to multicollinearity issues (i.e. due to the correlations between x_i and the other independent variables $x_j, j \neq i$). Hence, a fundamental property of $\hat{\beta}_i$ is that their value is proportional to linear dependency between the dependent and the independent variable $-\rho(x_i, Y)$, but also to the ratio of their standard deviation $\frac{\sigma(Y)}{\sigma(x_i)}$ – i.e. to the scaling of those variables. $\hat{\beta}_i$ then quantifies the change in Y (in Y unit) for an increase of 1 x_i (in x_i unit). Outside the neuroimaging community, it is common to also report and use *standardized* betas \hat{b}_i . Those are computed with normalized –or Z-scored- dependent Y^Z and independent x_i^Z variables: $x_i^Z = \frac{x_i - \mu(x_i)}{\sigma(x_i)}$ and $Y^Z = \frac{Y - \mu(Y)}{\sigma(Y)}$. This implies $\sigma(x_k^Z) = \sigma(Y^Z) = 1$, and therefore $\hat{b}_i = \rho(x_i, Y) \sqrt{VIF_i}$. \hat{b}_i quantifies the change in Y (in Y standard deviation) for an increase of 1 x_i (in x_i standard deviation).

Finally, one can compute \hat{t}_i , the *Student t-statistic* of $\hat{\beta}_i$, relative to the null hypothesis $\hat{\beta}_i = 0$. Assuming a Gaussian noise u in (1), we have $\hat{t}_i = \frac{\hat{\beta}_i}{s(\hat{\beta}_i)}$. Here, $s(\hat{\beta}_i) = \sqrt{\hat{\sigma}^2 VIF_i} \frac{\sigma(Y)}{\sigma(x_i)}$ is the standard error of the estimate $\hat{\beta}_i$, and $\hat{\sigma}^2 = \frac{SSE}{n-p}$, (where n is the sample size, p is the number of coefficients in the model including intercept, and SSE is the sum of squared errors) estimates σ^2 , the variance of the errors of the regression model. Hence, one can easily show that $\hat{t}_i = \frac{\rho(x_i, Y)}{\hat{\sigma}}$, or in other words, that \hat{t}_i only depends on the linear dependency between the dependent and the independent variable, and on the overall quality of the regression model. \hat{t}_i follow a Student's t-distribution with $(n - p)$ degrees of freedom, from which the P-value corresponding to the null hypothesis $\hat{\beta}_i = 0$ can be computed.

Z-values of $\hat{\beta}_i$, which are sometimes preferred to *t-values* because they are independent of the sample size (number of degrees-of-freedom), are typically re-computed from these P-values.

Scaling laws

We consider an idealized situation, where, in each individual k the BOLD signal in a brain region/voxel (Y) encodes one behavioral parametric measure of interest (x_k) whose distribution ($\sigma(x_k)$) varies between individuals due to a heterogeneity source. We also assumes that $\rho(x_k, Y)$ and $\sqrt{VIF_k}$ are independent of $\sigma(x_k)$ and $\sigma(Y_k)$, i.e. that the quality of this encoding is similar across subjects and does not depend on the individual brain activation or behavioral variable or range. We can formalize the *proportional* hypothesis by setting $\hat{\beta}_k \propto \alpha$ (2). By neglecting inter-individual differences in SR_k and VIF_k (1) and (2) imply, $\sigma(x_k) \propto \sigma(Y_k)$ (3), i.e. the BOLD activation scale proportionally with the behavior. Likewise, we can formalize the *normalization* hypothesis by setting $\sigma(Y_k) \propto \alpha$ (4). In this case, (1) and (4) imply $\hat{\beta}_k \propto \frac{\alpha}{\sigma(x_k)}$ (5), i.e. the activation summary statistics is inversely correlated with the standard deviation of the behavior.

Scaling laws and Z-scoring

Z-transforming individually a measure of interest entails subtracting its original mean $\mu_{X,k}$ and dividing the resulting centered variable by its original standard deviation $\sigma(x_k)$, i.e. $x_k^Z = \frac{x_k - \mu(x_k)}{\sigma(x_k)}$. Let us assume that the two scaling hypothesis are still related to the original variables, i.e. (2) and (4) still hold, but that first level analysis are conducted with normalized variables, i.e. $\sigma(x_k^Z) = 1$, hence $\hat{\beta}_k^Z = \sigma(Y_k)$ (6). Hence, under the *proportional* hypothesis, (3) and (6) imply $\hat{\beta}_k^Z \propto \sigma(x_k)$, while under the *normalization* hypothesis, (4) and (6) imply $\hat{\beta}_k^Z \propto \alpha$ (7).

Subjects

Studies were approved by the local Ethics Committee: the Ethics Committee for Biomedical Research of the Pitié-Salpêtrière Hospital for the 3 rating studies, and the local Ethical Committee of the University of Trento for the learning study. All subjects gave informed consent prior to partaking in the, study.

Subjects of the *rating* studies were paid 100€ for the fMRI experiments. A total of 65 subjects were included in the 3 different rating studies (Study 1: n=20, 10 males, age=22.0±2.7; Study 2: n=19, 11 males, age=23.9±4.0;

Assessing inter-individual variability in brain-behavior relationship

Study 3: $n=26$, 12 males, age 25.3 ± 5.5).

Subjects of the *learning* study were remunerated according to the exact amount of money won in the experiment plus a fixed amount for their travel to the MRI center. A total of 28 subjects (16 females; age 25.6 ± 5.4 years) were included in this study.

Task and behavioral analyses

Rating tasks. The behavioral tasks involved rating procedures on a Likert scale that were implemented as follows: subjects could move the cursor by pressing a button with the right index to go left or with the right middle finger to go right. Ratings were all self-paced, and subjects had to press a button with the left index finger to validate their response and go to the next trial. The initial position of the cursor on the scale was randomized to avoid confounding the ratings with the movements they involved.

Details specific to the different tasks are described below. 2009 (see Lebreton, et al. 2009, Lebreton, et al. 2012, and Lebreton, et al. 2015 for detailed methods)

Study 1: This fMRI study is a re-analysis of data obtained in Lebreton, et al. (2009). Stimuli were 120 faces, 120 houses and 120 paintings, for a total of 360 pictures which were randomly distributed over 6 sessions of 60 trials each (20 faces, 20 houses and 20 pictures). In every trial, the picture was first displayed on the screen for 3 seconds, following a fixation cross. Then a -10-10 rating scale appeared, and participants had to indicate on this scale how pleasant or how old the presented stimulus was.

Study 2: This fMRI study is a re-analysis of data obtained in Lebreton, et al. (2012). Stimuli were 240 short (2-5 sec) videos featuring different objects (food, toys, clothes, and tools), randomly distributed over four 60-trial sessions. In every trial, the video was first played on the screen, following a fixation cross. Then a 0-10 rating scale appeared, and participants had to indicate “how much they would like to acquire the object”.

Study 3: This fMRI study is a re-analysis of data obtained in Lebreton, et al. (2015). Stimuli were 270 potential events from various domains (politics, sport, society, culture, media, economics, diplomacy, science, technology, etc...). They were randomly distributed over 5 sessions of 54 trials each. Subjects were instructed to read the text depicting the event and think of how pleased they would feel should this event happen in the next 5 years (desirability rating). On every trial one prospect was displayed alone on the screen (5-7 seconds), following a 1s fixation cross. The desirability (-10-10) or probability (0-100%) scale only appeared after prospect display.

Learning task. Subjects were repeatedly presented with fixed pairs of abstract symbols. Over 4 sessions, eight novel options, defining four novel fixed pairs, were presented 24 times for a total of 96 trials. The four option pairs corresponded to four contexts (reward/partial, reward/complete, punishment/partial and punishment/complete), associated with different pairs of outcomes (reward contexts: winning 0.5€ versus nothing; punishment contexts: losing 0.5€ versus nothing) and different quantities of information being given at feedback (partial and complete). In the partial feedback contexts, only the outcome about the chosen option was provided, while in the complete feedback contexts both the outcome of the chosen and the unchosen option were provided. Within each pair, the two options were associated to the two possible outcomes with reciprocal probabilities (0.75/0.25 and 0.25/0.75).

Computational model. For this manuscript, we used the model named RELATIVE in Relative from (Palminteri et al., 2015). The model is an adaptation of a Q-learning model. It stipulates that subjects learn by trial and error to compute a value $Q(s)$ for each option. These values are learned via a Rescorla-Wagner rule (also called delta-rule): they are updated at each trial, by integrating an error term, which compare this expected value $Q(s)$ to the actual outcome – a so-called prediction-error δ . As reported in Palminteri, et al (2015), this model provide a very good account of the subjects' behavior (see the original paper for an extensive description and justification of the RELATIVE model parameters and its relation to other models, such as the actor-critic model). In this paper, we focus on two of RELATIVE model individual free-parameters: the “factual” learning-rate $\alpha_{1,k}$, and the inverse temperature θ_k . The choice of $\alpha_{1,k}$, is justified by the fact that this learning rate is involved in both conditions (complete and incomplete information, as opposed to $\alpha_{2,k}$ which is only involved in counterfactual learning i.e. in the complete information condition) and directly impact the Q-values (as opposed to $\alpha_{3,k}$, which is only used in the centering process).

Parameter optimization. We optimized the model free-parameters, the temperature (*temp*), the factual (α_1), the counterfactual (α_2) and the contextual (α_3) learning rates, by minimizing the negative log likelihood (Lmax) of the participant choices under the model using Matlab's *fmincon* function, initialized at multiple starting points of the parameter space.

For the population parameter condition, a single set of parameters was estimated to account for the behavior of all 28 subjects. This set of parameter could then be used to generate individual time-series of Q-values (\mathbf{q}_c^P).

Assessing inter-individual variability in brain-behavior relationship

For the individual parameter conditions, a set of parameters was estimated per subject. This set of parameter could then be used to generate individual time-series of Q-values (q_c), which could also be subsequently Z-scored (q_c^Z).

Model Comparison. Negative log-likelihoods (Llmax) were used to compute classical model selection criteria. We computed the Akaike's information criterion (AIC) at the individual level (RFX), and at the group level (FFX):

$$AIC_{FFX} = 2 \times (\text{LlmaxGroup} + \text{DF});$$
$$AIC_{RFX} = \sum_{\text{subjects}} 2 \times (\text{LlmaxGroup} + \text{DF});$$

We also computed the Bayesian information criterion (BIC) at the individual level (RFX), and at the group level (FFX):

$$BIC_{FFX} = 2 \times (\text{LlmaxGroup}) + \text{DF} \times \log(n_{\text{trial}} \times n_{\text{subjects}});$$
$$BIC_{RFX} = \sum_{\text{subjects}} 2 \times (\text{LlmaxSub}) + \text{DF} \times \log(n_{\text{trial}});$$

Where "DF" is the number of free parameters.

Neuroimaging

Data acquisition. For all imaging studies, T2*-weighted echo planar images (EPI) were acquired with blood oxygen-level dependent (BOLD) contrast. All studies employed a tilted plane acquisition sequence designed to optimize functional sensitivity in the orbitofrontal cortex and medial temporal lobes (Deichmann et al., 2003).

The *rating* studies were imaged with a 3.0 Tesla magnetic resonance scanner. To cover the whole brain with good spatial resolution, we used the following parameters: Study 1: TR=2.29s, 35 slices, 2 mm slice thickness, 1 mm inter-slice gap ; Study 2: TR=2.0s, 35 slices, 2 mm slice thickness, 1.5 mm inter-slice gap; Study 3: TR=2.03s, 35 slices, 2 mm slice thickness, 1.6 inter-slice gap.

The *learning* study was imaged with a 4.0 Tesla magnetic resonance scanner (4T Bruker MedSpec Biospin MR scanner – CiMEC, Trento, Italy). To cover the whole brain with good spatial resolution, we used the following parameters: TR=2.20s, 47 slices, 2 mm slice thickness, 1 mm inter-slice gap.

For all studies, T1-weighted structural images were also acquired, co-registered with the mean EPI, normalized to a standard T1 template, and averaged across subjects to allow group level anatomical localization. EPI data were analyzed in an event-related manner, within a general linear model, using the statistical parametric mapping software SPM8 (Wellcome Trust center for NeuroImaging, London, UK) implemented in MATLAB®. The first 5 volumes of each session were discarded to allow for T1 equilibration effects. Preprocessing consisted of spatial realignment, normalization using the same transformation as structural images, and spatial smoothing using a Gaussian kernel with a full-width at half-maximum (FWHM) of 8 mm. Preprocessed images were subsequently analyzed in an event related manner within the general linear model (GLM) framework.

GLMs – rating studies. For each *rating* study; we used two similar GLM to explain subject level time-series: the only difference was that in the first GLM, the parametric regressor "value" was entered in the native form $-v$ -, whereas in the second GLM, it was normalized (i.e. Z-scored) per subject and session (and category for Study 1) $-v^Z$ -.

- Study 1: Events were image onsets, corresponding to the 3 categories of stimuli (face, house, painting), modeled as a stick functions. These 3 categorical regressors were modulated by the parametric regressor accounting for the pleasantness rating. We also modeled the rating period in another regressor with a stick function modulated by response time.

- Study 2: Events were video display, modeled as boxcar function. This categorical regressor was modulated by the parameters accounting for the desirability ratings. We also modeled the rating period in another regressor with a stick function modulated by response time.

- Study 3: Desirability rating trials were modeled as boxcar functions covering stimulus presentation. This event was modulated the parameter accounting for the desirability ratings. We also modeled the rating period in another regressor with a stick function modulated by response time.

GLMs – learning study. For the *learning* study; we used three similar GLM to explain subject level time-series: the only difference was that the parametric regressor "chosen Q-value" (Q_c) used in the GLM could be generated using the population parameter (q_c^P) or the individual model free-parameters. In this latter case, the Q_c could be entered in the GLM in their native scale (q_c), or Z-scored per individual and session (q_c^Z) In all

Assessing inter-individual variability in brain-behavior relationship

GLMs, each trial was modelled as having two time points, corresponding to choice and outcome display onsets, modelled by two separate regressors. Choice onset was then modulated with a parametric regressors accounting for the chosen option Q-value, and the outcome onset was modulated with a parametric regressors accounting for the actual outcome (+0.5; 0; or - 0.5).

Whole-brain analysis. All regressors of interest were convolved with a canonical hemodynamic response function (HRF). To correct for motion artifacts, subject-specific realignment parameters were modeled as covariates of no interest. Linear contrasts of regression coefficients (betas) were computed at the session level, averaged at the subject level, and taken to a group-level random effect analysis, using one-sample t-tests. Unless otherwise specified, all activations maps were threshold using family-wise correction for multiple comparison (FWE) at the cluster level ($P_{FWE} < 0.05$). This cluster-wise correction was estimated by SPM8 using cluster-generating voxel-level thresholds of $P_{UNCORR} < 0.001$.

Region of interest (ROI). The VMPFC anatomical ROI was generated using WFU PickAtlas, and include a bilateral mask of the Frontal Medial Orbital cortex.