



**HAL**  
open science

## The genetic architecture of language functional connectivity

Yasmina Mekki, Vincent Guillemot, Hervé Lemaître, Amaia Carrión-Castillo, Stephanie Forkel, Vincent Frouin, Cathy Philippe

► **To cite this version:**

Yasmina Mekki, Vincent Guillemot, Hervé Lemaître, Amaia Carrión-Castillo, Stephanie Forkel, et al.. The genetic architecture of language functional connectivity. *NeuroImage*, 2022, 249, pp.118795. 10.1016/j.neuroimage.2021.118795 . hal-03566120

**HAL Id: hal-03566120**

<https://hal.sorbonne-universite.fr/hal-03566120v1>

Submitted on 11 Feb 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License



# The genetic architecture of language functional connectivity

Yasmina Mekki<sup>a,\*</sup>, Vincent Guillemot<sup>b</sup>, Hervé Lemaître<sup>c</sup>, Amaia Carrión-Castillo<sup>d</sup>,  
Stephanie Forkel<sup>c,e,f</sup>, Vincent Frouin<sup>a</sup>, Cathy Philippe<sup>a,\*</sup>

<sup>a</sup> NeuroSpin, Institut Joliot, CEA - Université Paris-Saclay, Gif-Sur-Yvette 91191, France

<sup>b</sup> Hub de Bioinformatique et Biostatistique, Département Biologie Computationnelle, Institut Pasteur, Paris, France

<sup>c</sup> Groupe d'Imagerie Neurofonctionnelle, Institut des Maladies Neurodégénératives, CNRS UMR 5293, Université de Bordeaux, Centre Broca Nouvelle-Aquitaine, Bordeaux, France

<sup>d</sup> Basque Center on Cognition, Brain and Language, San Sebastian, Spain

<sup>e</sup> Brain Connectivity and Behavior Laboratory, Sorbonne Universities, Paris, France

<sup>f</sup> Department of Neuroimaging, Institute of Psychiatry, Psychology and Neurosciences, King's College London, UK

## ARTICLE INFO

### Keywords:

Imaging-genetics

Resting-state functional MRI

Language

GWAS

UK Biobank

Multivariate analysis

## ABSTRACT

Language is a unique trait of the human species, of which the genetic architecture remains largely unknown. Through language disorders studies, many candidate genes were identified. However, such complex and multifactorial trait is unlikely to be driven by only few genes and case-control studies, suffering from a lack of power, struggle to uncover significant variants. In parallel, neuroimaging has significantly contributed to the understanding of structural and functional aspects of language in the human brain and the recent availability of large scale cohorts like UK Biobank have made possible to study language via image-derived endophenotypes in the general population. Because of its strong relationship with task-based fMRI (tbMRI) activations and its easiness of acquisition, resting-state functional MRI (rsfMRI) have been more popularised, making it a good surrogate of functional neuronal processes. Taking advantage of such a synergistic system by aggregating effects across spatially distributed traits, we performed a multivariate genome-wide association study (mvGWAS) between genetic variations and resting-state functional connectivity (FC) of classical brain language areas in the inferior frontal (pars opercularis, triangularis and orbitalis), temporal and inferior parietal lobes (angular and supramarginal gyri), in 32,186 participants from UK Biobank. Twenty genomic loci were found associated with language FCs, out of which three were replicated in an independent replication sample. A locus in 3p11.1, regulating *EPHA3* gene expression, is found associated with FCs of the semantic component of the language network, while a locus in 15q14, regulating *THBS1* gene expression is found associated with FCs of the perceptual-motor language processing, bringing novel insights into the neurobiology of language.

## 1. Introduction

Language is a unique trait of the human species. Although its genetic origins are broadly accepted, they remain largely unknown. Since the seminal study that revealed the major role of *FOXP2* in language processing (Fisher et al., 1998), several candidate genes related to language disorders were identified (Landi and Perdue, 2019). Human language is a complex system – both structurally and functionally. As such a complex and multifactorial trait, it is unlikely to be associated with only a few genes but rather with many genes that are also interact-

ing with each other. These genes likely contribute to the development of neural pathways involved in language development, together with experience-dependent contributions from the environment (Fisher and Vernes, 2015). In parallel, neuroimaging techniques provided innovative, quantitative and non-invasive ways to study language and thus, widened the doors to the investigation of language neurobiological organization in the brain in general population. Anatomically, the language system comprises perisylvian cortical regions predominantly - but not exclusively - in the left hemisphere. Among these regions the prominent regions are in the pars orbitalis and triangularis in the in-

**Abbreviations:** eQTL, expression Quantitative Trait Locus; FA, Fractional Anisotropy; FC, Functional Connectivity; FDR, False Discovery Rate; GWAS, Genome-Wide Association Study; ICA, Independent Component Analysis; ICVF, IntraCellular Volume Fraction; ISOVF, Isotropic Volume Fraction; LD, Linkage Disequilibrium; MD, Mean Diffusivity; MO, Tensor Mode; MOSTest, Multivariate Omnibus Statistical Test; mvGWAS, multivariate GWAS; OD, Orientation Dispersion index; ROI, Region Of Interest; rsfMRI, resting-state functional MRI; SNP, Single Nucleotide Polymorphism; STAP, Superior Temporal Asymmetrical Pit; tbfMRI, task-based functional MRI; dMRI, Diffusion Magnetic Resonance Imaging; DW-MRI, Diffusion-weighted Magnetic Resonance Imaging.

\* Corresponding authors.

E-mail addresses: [yas.mekki@gmail.com](mailto:yas.mekki@gmail.com) (Y. Mekki), [cathy.philippe@cea.fr](mailto:cathy.philippe@cea.fr) (C. Philippe).

<https://doi.org/10.1016/j.neuroimage.2021.118795>.

Received 30 July 2021; Received in revised form 11 November 2021; Accepted 8 December 2021

Available online 18 December 2021.

1053-8119/© 2021 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license

(<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

ferior frontal gyrus (also referred to as ‘Broca’s’ region), the angular and supramarginal gyri in the inferior parietal lobe (also referred to as ‘Geschwind’s’ region), and the posterior temporal regions (‘Wernicke’s’ region). These cortical regions are interconnected by a network of brain connections, most prominently the arcuate fasciculus (Catani et al., 2005; Forkel and Catani, 2018). These regions also connect to the sensory-motor system (auditory, visual, and motor cortex). Functionally, phonology, semantics, and syntax are three main language components and form a tripartite interconnected architecture, with regions that are specialized in one or two components (e.g. Jackendoff and Jackendoff (2002); Bates et al. (2003); Vigneau et al. (2006); Price (2012)). Several studies have used neuroimaging-derived features to extend the understanding of how the brain supports language (Price, 2012; Friederici, 2017; Ardila et al., 2016; Leroy et al., 2015; Labache et al., 2020). These MRI features give access to biologically relevant measurements of individual variability (Forkel et al., 2014a, 2020b, 2020a; Uddén et al., 2019; Dubois and Adolphs, 2016; Seghier and Price, 2018; Fedorenko, 2021) and are consequently suitable for the search of genetic associations such as the planum temporale volume asymmetry (Carrion-Castillo et al., 2020), the superior temporal asymmetrical pit (STAP) (Le Guen et al., 2020), or the left–right hemispheric asymmetry of the brain (Sha et al., 2021), and are then called endophenotypes. Despite low frequency blood-oxygen-level-dependent (BOLD) signals acquired in resting state condition, rsfMRI does reflect large-scale circuit organization (Biswal et al., 1995; Fox and Raichle, 2007), they are less used to study specific brain functions. However, a growing body of evidence suggests that there is a close correspondence between resting state networks and known cognitive task activation maps (Smith et al., 2009; Cole et al., 2014). (Tavor et al., 2016) showed that a prediction model relating FC estimated from rsfMRI to task-based activations can accurately predict individual differences in brain activity for participants not used in the training. More specifically regarding language, several studies sought the brain dynamics and showed that resting-state FC map could map language networks previously identified in task-based imaging studies (Hampson et al., 2006; Kelly et al., 2010; Koyama et al., 2010; Xiang et al., 2010). Especially, Xiang et al. (2010) used resting-state FC to infer the functional organization of Broca’s area and the perisylvian language networks by investigating their functional correlations and reported a clear topographical FC pattern in the left middle frontal, parietal, and temporal areas. On the other side, several studies has examined the relationship between resting-state FC and behavioral scores. Relevant to our work, these studies have shown that language-related performances like basic reading abilities are correlated to the strength of resting-state FC (Koyama et al., 2011; Stevens et al., 2017; Cross et al., 2021; Cheema et al., 2021). Koyama et al. (2011) reported a positive correlation between the reading scores and the resting-state FC between the fusiform gyrus and the inferior frontal gyrus (IFG) -pars opercularis-, belonging to the reading network (RdN). Cheema et al. (2021) reported a positive correlation between reading fluency and the strength of FC between the left IFG and left supplementary motor area in a group without reading impairment. They also found that the stronger the FC between the supramarginal and the angular gyri, the higher is the word recognition accuracy score. Moreover, Cross et al. (2021) showed that the FC between the pars triangularis and the right fusiform gyrus was positively correlated with rapid automatized naming performance. Finally, Stevens et al. (2017) reported a positive correlation between accuracy on word classification and resting-state FC strength of the left occipitotemporal sulcus and Wernicke’s area. Next, recent works suggest that mapping perisylvian language regions can be accomplished using either tbfMRI or rsfMRI (Tie et al., 2014; Branco et al., 2016; Lu et al., 2017; Jones et al., 2017; Lemée et al., 2019; Park et al., 2020). Indeed, rsfMRI has been applied to mapping eloquent areas in surgical patients (Jones et al., 2017) and stroke patients (Klingbeil et al., 2019), where performing tbfMRI is not an option. Overall, these studies investigated the relevance of the rsfMRI approach in studying specific brain functions such as language, including its relationship with behavioral scores,

task-based network at both individual and group levels. They also highlighted its advantages over task-based design, when certain population has difficulties to engage in specific tasks. rsfMRI is paradigm free and as such it is easier to implement in very large cohorts like the UK Biobank. It can recover specific brain functional activations, making it a good surrogate of functional neuronal processes. Here, we propose to use task-free FC in UK Biobank (Bycroft et al., 2018), using perisylvian cortical areas as regions of interest (ROI) serving as a proxy for language.

As rsfMRI FC are low amplitudes and correlated to each other, we anticipate it would be really difficult to disentangle the genetic associations with each FC signal in a massively univariate manner. Language related brain regions share information across components and scales, and genetic variants are supposed to have distributed effects across regions. Thus, we take advantage of this synergistic system and perform a multivariate approach with MOSTest (van der Meer et al., 2020). This method considers the distributed nature of genetic signals shared across brain regions and aggregates effects across spatially distributed traits of interest. This approach tests each SNP independently for its simultaneous association with the brain endophenotypes, making it half multivariate and half univariate. For convenience, we will use the term “multivariate GWAS” (mvGWAS) while being aware that correlations between SNPs are not accounted for in this approach.

In this study, we use rsfMRI in a discovery sample of 32,186 healthy participants from the UK Biobank (Sudlow et al., 2015) and the compiled information of a large-scale meta-analysis on language components (Vigneau et al., 2006, 2011). First, we used a region of interest (ROI) informed calculation of the functional connectivity (FC) and more precisely the ROI-to-ROI FC estimation approach. We retained that differences in such FCs are related to language differences. In our work, we studied the associations between the FC differences and genotypes differences across the participants, to indirectly attempt to reveal genes associated to language. For this purpose, we performed a mvGWAS of language specific FC, filtered based on heritability significance. The results from this analysis were subjected to a replication study in an independent sample ( $N = 4,754$ ). Additionally, as the connections between different language regions are ensured by the white matter bundles (Catani et al., 2005; Catani and Forkel, 2019), we tested the potential associations of the hit SNPs with the neuroanatomical tracts underlying the hit endophenotypes using diffusion-based white matter analysis. Finally, the extensive functional annotations of each genomic risk locus allowed us to suggest two new genes with a role in different aspects of the language system.

## 2. Materials and methods

### 2.1. Demographics and neuroimaging data from the UK Biobank

#### 2.1.1. UK Biobank cohort

The UK Biobank is an open-access longitudinal populationwide cohort study that includes 500k participants from all over the United Kingdom (Sudlow et al., 2015). Data collection comprises detailed genotyping and a wide variety of endophenotypes ranging from health/activity questionnaires, extended demographics to neuroimaging and clinical health records. All participants provided informed consent and the study was approved by the North West Multi-Center Research Ethics Committee (MREC). This study used the February 2020 release (UK Biobank application number #64984). This release consisted of 36,940 participants, age range between 40 and 70 years (mean age =  $54 \pm 7.45$  years), with genotyping and resting-state functional MRI. To avoid any possible confounding effects related to ancestry, we restricted our analysis to individuals with British ancestry using the sample quality control information provided by UK Biobank (Bycroft et al., 2018). A final cohort of 32,186 participants (16,952 females, mean age =  $55 \pm 7.51$  years) were included in the study. We made use of the first ten principal components (Data field 22009) of the genotyping data’s principal component analysis capturing population genetic diversity to account for population

stratification. An independent replication dataset of 4,754 non-British individuals was also drawn from the UK Biobank. The age range of these participants was 40 to 70 (mean age =  $53 \pm 7.55$  years), 2,601 were female.

### 2.1.2. Resting-state functional MRI data

The MRI data available from the UK Biobank are described in the UK Biobank Brain Imaging Documentation (v.1.7, January 2020) as well as in Miller et al (2016) and Alfaro-Almagro et al. (2018). Briefly, rsfMRI data were acquired using the following parameters: 3T Siemens Skyra scanner, TR = 0.735 s, TE = 39 ms, duration = 6 min (490 time points), resolution:  $2.4 \times 2.4 \times 2.4$  mm, Field-of-view =  $88 \times 88 \times 64$  matrix. During the resting-state scan, participants were instructed to keep their eyes fixated on a crosshair, to relax, and to think of nothing particular (Miller et al., 2016). The preprocessing of the UK Biobank data includes motion correction, grand-mean intensity normalization, high-pass temporal filtering including EPI (echo-planar imaging) unwarping with alignment to the T1 template and gradient non-linearity distortion correction (GDC) unwarping, brain masking, and registration to MNI (Montreal Neurological Institute-Hospital) space. The rsfMRI volumes were further cleaned using FMRIB's ICA-based Xnoiseifier (Functional Magnetic Resonance Imaging of the Brain's Independent Component Analysis Xnoiseifier) for automatically identifying and removing artefacts.

### 2.1.3. Diffusion-weighted MRI data

The Diffusion-weighted MRI (DW-MRI) images were acquired using the following parameters; isotropic voxel size (resolution):  $2 \times 2 \times 2$  mm, five non diffusion-weighted image  $b = 0$  s/mm<sup>2</sup>, diffusion-weighting of  $b = 1000$ , and 2000s/mm<sup>2</sup> with 50 directions each, acquisition time: 7 min. Tensor fits utilize the  $b = 1000$ s/mm<sup>2</sup> data and the NODDI (Nurite Orientation Dispersion and Density Imaging) (Zhang et al., 2012) model is fit using AMICO (Accelerated Microstructure Imaging via Convex Optimization) (Daducci et al., 2015) tool, creating outputs including nine diffusion indices maps. These ones were subject to a TBSS-style (Tract-based spatial statistics) analysis using FSL (FMRIB Software Library) tool resulting in a white matter skeleton.

## 2.2. Genetic quality control

Genotyping was performed using the UK BiLEVE Axiom array by Affymetrix (Wain et al., 2015) on a subset of 49,950 participants (807,411 markers) and the UK Biobank Axiom array on 438,427 participants (825,927 markers). Both arrays are extremely similar and share 95% of common SNP probes. The imputed genotypes were obtained from the UK Biobank repository Bycroft et al. (2018). These genetic data underwent a stringent quality control protocol, excluding participants with unusual heterozygosity, high missingness (Data field 22027). Variants with minor allele frequency lower than 0.01 were filtered out from the imputed genotyping data using PLINK 1.9 (Chang et al., 2015) to retain the common variants only. Overall, 9,812,367 autosomal SNPs were considered.

## 2.3. Regions of interest for rsfMRI functional connectivity

We leverage a large-scale meta-analysis of 946 activation peaks (728 peaks in the left hemisphere, 218 peaks in the right hemisphere) obtained from a meta-analysis of 129 task-based fMRI language studies (Vigneau et al., 2006, 2011). The identified fronto-parietal-temporal activation foci revealed via a hierarchical clustering analysis, 50 distinct, albeit partially overlapping, clusters of activation foci for phonology, semantics, and sentence processing: 30 clusters in the left hemisphere and 20 in the right hemisphere.

Because this overlap could unduly increase the co-activation between regions and to avoid a deconvolution bias in the estimation of the functional connectivity, we proceeded as follow: First, because the

clustering process was performed for each component independently, we checked whether pairs of clusters belonging to different language-component networks were spatially distinct considering the significance of their mean Euclidean distance with paired t-tests. We identified areas that are common to multiple language components; in the temporal lobe, the anterior part of the Superior temporal gyrus (T1a) area appears to be common to all three language components, the anterior part of the superior temporal sulcus (Pole) and Lateral/middle part of the middle temporal gyrus (T2ml) are common to semantic and sentence's clusters and the posterior part of the left inferior temporal gyrus (T3p) to semantic and phonology clusters. In the frontal lobe, the L—R dorsal part of the pars opercularis (F3opd) and the ventral part of the pars triangularis (F3tv) are common to semantic and syntactic clusters. In these cases, we retained the larger cluster and assigned multiple labels. Second, ROIs were obtained for each cluster by building a 3D convex-hull of the peaks in the MNI space and were then subjected to a morphological opening operation. Third, overlapping areas between the convex-hull ROIs were processed as follows: the common region between two ROIs was attributed to the most representative in terms of the number of peaks. Finally, we excluded regions with less than 100 voxels. This preprocessing resulted in 25 multilabelled ROIs: 19 in the left hemisphere and 6 in the right hemisphere which are summarised in Fig. 1A and Table S11.

## 2.4. Neuroimaging endophenotypes

### 2.4.1. FC endophenotypes

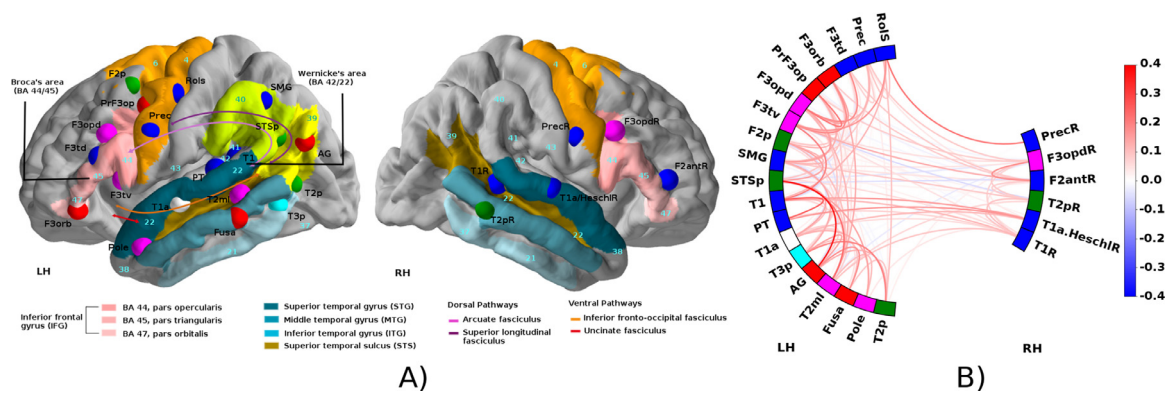
The preprocessed resting-state BOLD signal was masked using the 25 ROIs and averaged at each time volume. A connectome matrix was computed using Nilearn (Abraham et al., 2014) for each participant using a shrunk (Ledoit and Wolf, 2004) estimate of partial correlation (Marrelec et al., 2006). This resulted in 300 (=  $25 \times 24/2$ ) edges connecting language ROIs for each individual. Each edge -also denoted functional connectivity (FC)- is further considered as a candidate endophenotype. See the Fig. 1B.

### 2.4.2. Diffusion MRI endophenotypes

We hypothesised that the hit SNPs associated with the hit FCs could be associated with neuroanatomical white matter tracts that supports the information transmission between the regions that compose these hit-FCs. Therefore, we tested the potential associations between the hit SNPs with the following white matter bundles: the corpus callosum, the left frontal aslant tract, the left arcuate anterior/long/posterior segment, the left inferior fronto-occipital fasciculus, the left uncinate tract (See Section 3.3 for more details). The resulting skeletonised images are averaged across the set of 7 brain white matter structures defined by the probabilistic atlas (Rojkova et al., 2016) thresholded at 90% of probabilities. These structural white matter tracts are assessed by 9 indices: fractional anisotropy (FA) maps, tensor mode (MO), mean diffusivity (MD), intracellular volume fraction (ICVF), isotropic volume fraction (ISOVF), mean eigenvectors (L1, L2, L3), and orientation dispersion index (OD) yielding  $63 = 7 \times 9$  dMRI endophenotypes.

## 2.5. SNP-based heritability and genetic correlation analysis

The proportion of additive genetic variance in the FC phenotypic variance, also called narrow-sense heritability, was estimated using the genotyped SNPs information using genome-based restricted maximum likelihood (GREML) (Yang et al., 2010) for each FCs, controlling for the above-mentioned covariates (refer to Section 2.4). To define significantly heritable FCs, a 0.05 threshold on False Discovery Rate (FDR) adjusted p-values was applied to account for multiple testing on the 300 FCs. Similarly, the proportion of additive genetic variance in the covariance of pairs of FCs was estimated using the bivariate GREML (Lee et al., 2012). Both heritability and part of covariance explained by



**Fig. 1.** (A) Overview of the regions obtained from the meta-analysis. Each language seed is color-coded according to its language category: phonology (blue), semantic (red), and syntax (green). ROIs of different components that were not spatially distinct are color-coded as pink (semantic/syntax), cyan (phonology/semantic) and white for the three language component. For the sake of ROIs figure visibility, the coordinates were modified. The exact coordinates for each ROI are available in Table SI9. Different gyri and sulcus, known to be relevant for language: the inferior frontal gyrus (IFG), middle temporal gyrus (MTG), superior temporal gyrus (STG), and superior temporal sulcus (STS), are color-coded. Numbers in the left hemisphere (LH) represents language-relevant Brodmann areas (BA) which were defined on the basis of cytoarchitectonic characteristics. Numbers in the right hemisphere (RH) represents the language-relevant BA counterpart. The pars opercularis (BA 44), the pars triangularis (BA 45) represents Broca's area. The pars orbitalis (BA 47) is located anterior to Broca's area. BA 42 and BA 22 represents Wernicke's area [Friederici \(2011\)](#). Both supramarginal gyrus (BA40) and angular gyrus (BA39), also known as Geschwind's territory, are represented by green/yellow colors, respectively. The primary motor cortex (BA4), the premotor cortex and the supplementary motor area (B6) are colored in orange. Within the left hemisphere, dorsal and ventral long-range fiber bundles connect language areas and are indicated by color-coded arrows. (B) Mean functional connectivity of the 142 heritable endophenotypes, calculated using a shrinked estimate of partial correlation [Marrelec et al. \(2006\)](#) (estimated with a Ledoit-Wolf estimator ([Ledoit and Wolf, 2004](#))) over 32,186 UKB rs fMRI participants.

genetics were obtained using GCTA (Genome-wide Complex Trait Analysis) ([Yang et al., 2011](#)).

## 2.6. Multivariate genome-wide association studies (mvGWAS)

We performed a mvGWAS between the filtered imputed genotypes and the 142 significantly heritable FC endophenotypes, using the Multivariate Omnibus Statistical Test (MOSTest) ([van der Meer et al., 2020](#)). All endophenotypes were pre-residualised controlling for covariates including sex, genotype array type, age, recruitment site, and ten genetic principal components provided by UK Biobank. In addition, MOSTest performs a rank-based inverse-normal transformation of the residualised endophenotypes to ensure that the inputs are normally distributed. The distributions across the participants of all endophenotypes were visually inspected before and after covariate adjustment. MOSTest generated summary statistics that capture the significance of the association across all heritable 142 language FC endophenotypes. To account for multiple testing over the whole genome, statistically significant SNPs were considered as those reaching the genome-wide threshold  $p = 5e-8$ .

## 2.7. mvGWAS replication

The mvGWAS results were replicated in an independent non-British sample considering the nominal significance threshold  $p < 0.05$ . Following the same pre-processing steps as for the primary sample, the non-British replication sample consists in 4754 individuals with a mean age of 53 years ( $\pm 7.55$ ) and 2,601 female. Despite the different ancestry in the replication sample, we still consider the replication design as valid, as there is no evidence for ancestry being important in language and it is a common practice in the field ([Sha et al., 2021](#); [van der Meer et al., 2020](#)). Moreover, replicated hits can then be considered relevant for humans regardless of ancestry.

## 2.8. Fine-mapping: identification of genomic risk loci and functional annotation

We performed functional annotation analysis using the FUMA (functional mapping and annotation) online platform v1.3.6a ([Watanabe et al., 2017](#)) with default parameters. The genomic

positions are reported according to the GRCh37 reference. SNPs were annotated for functional consequences on gene functions using ANNOVAR ([Wang et al., 2010](#)), Combined Annotation Dependent Depletion (CADD) scores ([Kircher et al., 2014](#)), and 15-core chromatin state prediction by ChromHMM ([Ernst and Kellis, 2012](#)). In addition, they were annotated for their effects on gene expression using expression Quantitative Trait Loci (eQTLs) of various tissue types. The eQTL module queried data from different tissue-datasets using GTEx v8 ([Consortium et al., 2017](#)), Blood eQTL browser ([Westra et al., 2013](#)), BIOS QTL browser ([Zhernakova et al., 2017](#)), BRAINEAC ([Ramasamy et al., 2014](#)), eQTLGen ([Vösa et al., 2018](#)), PsychENCODE ([Wang et al., 2018](#)), DICE ([Schmiedel et al., 2018](#)). RegulomeDB v2.0 ([Boyle et al., 2012](#)) was queried externally. Coding hit SNPs are also annotated with polymorphism phenotyping v2 (Polyphen-2) ([Ramensky et al., 2002](#)).

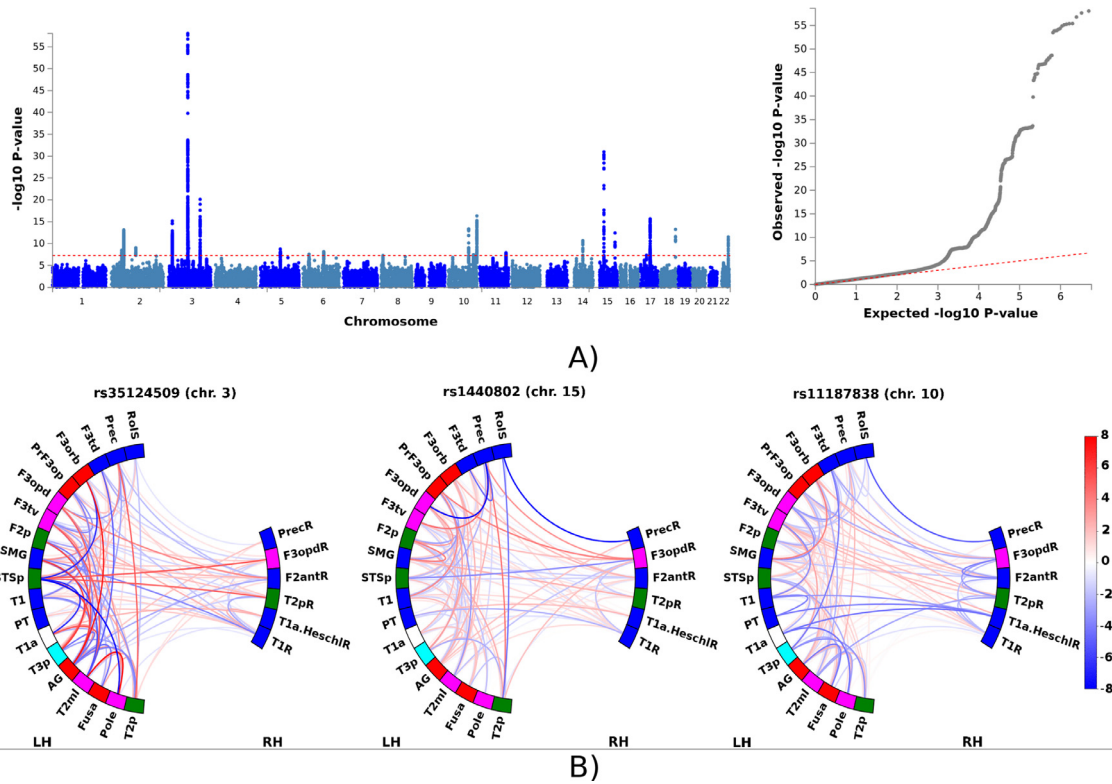
## 3. Results

### 3.1. SNP-based heritability of functional connectivity measures

The SNP-based heritability ( $h^2$ ) was estimated for each of the 300 FCs endophenotypes.  $P$ -values correction for multiple testing revealed 142 FCs with significant SNP-based heritabilities (Table SI2), ranging from 14% for the SMG $\leftrightarrow$ F3opd to 3% for the SMG $\leftrightarrow$ T1 FCs.

### 3.2. Multivariate genome-wide association analysis

We performed a mvGWAS using the MOSTest ([van der Meer et al., 2020](#)) method, with the 142 FCs with significant SNP-based heritabilities. This analysis tested each SNP separately for its simultaneous association with the 142 FCs and yielded 4566 significant SNPs at a genomic threshold (see Table SI3), distributed on chromosomes 2, 3, 5, 6, 10, 11, 14, 15, 17, 18 and 22. FUMA ([Watanabe et al., 2017](#)) software was used to analyze mvGWAS results and identify lead SNPs at each associated locus. Considering the genome-wide significance threshold  $p = 5e-8$ , there were 20 distinct genomic loci distributed on the 11 chromosomes, associated with different aspects of language FC (Figs. 2A, SI1, SI2, 2B, SI3 and Table 1) and represented by 20 lead SNPs.



**Fig. 2.** (A) Multivariate GWAS analysis of the resting state functional connectivity in 32,186 participants. Manhattan plot for multivariate GWAS across 142 FCs. The red dashed line indicates the genome-wide significance threshold  $p = 5e-8$ . The Quantile-quantile plot is also shown. (B) Circos plots illustrating the 3 lead SNPs identified from the mvGWAS. Z-values from the univariate GWAS for each FC are mapped. The absolute Z-values scaling is clipped at 8 ( $p = 1.2e-15$ ). Positive effects of carrying the minor allele are shown in red, and negative effects in blue.

**Table 1**

Lead SNP: ID of the lead SNPs within each locus. Position: position of the SNP in the hg19 human reference genome. mvgwasP discovery -British-: MOSTest association p-value obtained using the discovery sample. mvgwasP replication -non British-: MOSTest association p-value obtained using the independent replication sample. Functional category: Functional consequence of the SNP on the gene obtained from ANNOVAR. 'Central' phenotypes: the phenotypes that contributed most to the multivariate association considering the genome-wide association threshold ( $5e-8$ ).

Loci	Lead SNP	Chr	Position	Functional Category	Non Effect Allele	Effect Allele	MAF	mvgwasP (discovery -British-)	mvgwasP (replication -non British-)	Nearest Gene	"central" Phenotypes
1	rs62141276	2	48,214,217	NcRNA Intronic	A	G	0.367	$p = 3.26e-9$	$p = 0.27$	AC079807.4	-
2	rs2717046	2	58,041,936	intergenic	T	C	0.380	$p = 7.50e-14$	$p = 0.95$	CTD-2026C7.1	-
3	rs62158166	2	114,077,218	intergenic	C	G	0.223	$p = 8.69e-10$	$p = 0.38$	PAX8	-
4	rs67851870	3	17,554,860	intronic	G	A	0.322	$p = 6.57e-16$	$p = 0.35$	TBC1D5	-
5	rs35124509	3	89,521,693	exonic	C	T	0.401	$p = 8.95e-59$	$p = 3.25e-3$	EPHA3	AG↔F3orb; pSTS↔Pole; Pole↔T2ml; T1a↔STSp; STSp↔F3orb; SMG↔T3p; AG↔STSp; F2p↔AG; T2ml↔SMG
6	rs62266110	3	93,537,923	intergenic	A	G	0.319	$p = 1.17e-09$	$p = 0.93$	RNU6-488P	-
7	rs2279829	3	147,106,319	UTR3	T	C	0.212	$p = 7.57e-21$	$p = 0.68$	ZIC4	-
8	rs145120402	5	93,174,765	intronic	C	A	0.0433	$p = 1.83e-9$	$p = 0.10$	FAM172A	-
9	5:94,068,140_AC_A	5	94,068,140	intronic	A	AC	0.209	$p = 6.79e-9$	$p = 0.30$	ANKRD32:MCTP1	-
10	rs4262195	6	96,929,475	NcRNA Intronic	C	T	0.181	$p = 7.19e-9$	$p = 0.70$	UFL1-AS1	-
11	rs11187838	10	96,038,686	intronic	A	G	0.435	$p = 4.29e-14$	$p = 2.92e-2$	PLCE1	-
12	rs11146399	10	134,308,479	intergenic	T	C	0.457	$p = 5.50e-16$	$p = 0.28$	RP11-432J24.5	-
13	rs11218557	11	122,099,839	NcRNA Intronic	C	T	0.4579	$p = 1.24e-8$	$p = 0.77$	RP11-820L6.1	-
14	rs186347	14	59,072,226	intergenic	T	G	0.458	$p = 2.08e-11$	$p = 0.92$	DACT1	-
15	rs1440802	15	39,635,124	ncRNA Intronic	C	T	0.090	$p = 1e-31$	$p = 9.58e-3$	RP11-624L4.1	PrecR↔F3opd; PrecR↔Rols
16	rs4702	15	91,426,560	UTR3	A	G	0.442	$p = 3.77e-13$	$p = 0.42$	FURIN	-
17	rs34039488	17	27,320,232	intronic	A	G	0.162	$p = 4.74e-8$	$p = 0.46$	PIPOX:SEZ6	-
18	17:44,270,659_G_A	17	44,270,659	intronic	A	G	0.399	$p = 5.36e-16$	$p = 0.45$	KANSL1	-
19	rs7234875	18	73,114,340	intergenic	C	T	0.399	$p = 5.71e-14$	$p = 0.82$	RP11-321M21.3	-
20	rs2542028	22	47,196,524	intronic	G	A	0.268	$p = 3.06e-12$	$p = 0.60$	TBC1D22A	-

**Table 1:** Genomic loci associated with language FCs using the multivariate genome-wide association study.

### 3.2.1. Validation of lead SNPs associated with rsfMRI FCs

The three lead SNPs were replicated at the nominal significance level ( $p < 5e-2$ ) on multivariate test in the independent non-British replication dataset: rs1440802 ( $p = 9.58e-3$ ), rs35124509 ( $p = 3.25e-3$ ), rs11187838 ( $p = 2.92e-2$ ). Table S14 summarises these results. Moreover, these lead SNPs showed association at  $p < 0.05$  on univariate testing of all but three specific central traits identified in the discovery mvGWAS. Here, we present three of these loci that were replicated in an independent data set (refer to Section 2.3). MOSTest results highlighted the three following genomic risk regions: (i) 15q14 locus (chr15, start=39,598,529, length=260 kb) with its strongest association related to the imputed SNP rs1440802 ( $p = 1e-31$ ); (ii) 3p11.1 locus (chr3, start=89,121,389, length=1381 kb) with its strongest association related to the imputed SNP rs35124509 ( $p = 8.95e-59$ ); (iii) 10q23.33 locus (chr10, start=95,988,042, length=139 kb) with its strongest association related to the imputed SNP rs11187838 ( $p = 4.29e-14$ ). See Fig. 2A and Table 1.

### 3.2.2. Identification of central endophenotypes associated with genomic risk regions

For each lead SNP, we defined the 'central' endophenotypes that contributed the most in the multivariate association by using the individual univariate summary statistics performed by MOSTest and by considering the genome-wide significance threshold ( $p < 5e-8$ ) (Table S15).

On 15q14, the lead SNP rs1440802 had two central FCs: the minor allele was associated with the partial correlation between (i) the precentral gyrus and the dorsal pars opercularis (Prec $\leftrightarrow$ F3opd). Both connected regions are in the left frontal lobe, and are labelled with a phonological linguistic component (Prec) and multi-labelled with semantic and sentence language processing (F3opd). (ii) The (PrecR $\leftrightarrow$ RolS) corresponds to the partial correlation between the precentral gyrus and the Rolandic sulcus. Both regions are identified in the right and left frontal lobes, respectively, and are labelled as phonological linguistic component (Fig. 3A and Table S15). These edges have previously been described in FC studies dedicated to language and more specifically in the perceptual motor interactions (Schwartz et al., 2008; Fridriksson et al., 2009; Turner et al., 2009; Nishitani and Hari, 2000; Schwartz et al., 2012). At the univariate level, these loci associated to central endophenotypes display an important overlap; See Fig. 3D.

On 3p11.1, the lead SNP rs35124509 had nine central FCs: the minor allele was associated with the partial correlation between the left posterior part of the superior temporal sulcus and the left temporal pole (Pole $\leftrightarrow$ STSp), the left temporal pole and the lateral/middle part of the middle temporal gyrus (Pole $\leftrightarrow$ T2ml), the angular gyrus and the pars orbitalis of the left inferior frontal gyrus (AG $\leftrightarrow$ F3orb), the anterior part of the Superior temporal gyrus and the left posterior part of the superior temporal sulcus (T1a $\leftrightarrow$ STSp), the left posterior part of the superior temporal sulcus and the pars orbitalis of the left inferior frontal gyrus (STSp $\leftrightarrow$ F3orb), the supramarginal gyrus and the posterior part of the left inferior temporal gyrus (SMG $\leftrightarrow$ T3p), the angular gyrus and the left posterior part of the superior temporal sulcus (AG $\leftrightarrow$ STSp), the Posterior part of the middle frontal gyrus and the angular gyrus (F2p $\leftrightarrow$ AG), the lateral/middle part of the middle temporal gyrus and the supramarginal gyrus (T2ml $\leftrightarrow$ SMG) (Fig. 4A and Table S15). These connected regions are located across the left parieto-frontal-temporal lobe, and are mainly labelled as semantic language processing. These edges have previously been described in FC studies dedicated to language and especially to the semantic component. This component typically includes the inferior frontal gyrus, the left temporal cortex (i.e. temporal pole, middle temporal gyrus, fusiform gyrus) and the left angular gyrus (Binder et al., 2009; Jackson et al., 2016; Vigneau et al., 2006). At the univariate level, these

loci associated to central endophenotypes display an important overlap; See Fig. 4D.

A locus in 10q23.33 was highlighted by the mvGWAS. At the univariate level, no endophenotype reached the genome-wide significance threshold for the lead SNP in this locus (rs11187838).

As a conclusion of the mvGWAS, we retained: (i) a multifold link between two FCs and a locus in 15q14 region; and (ii) a multifold link between nine FCs and a locus in 3p11.1 region. Such a multivariate approach has the advantage of leveraging the distributed nature of genetic effects and the presence of pleiotropy across endophenotypes. Loci, respectively, identified by MOSTest as associated with several FCs made clear that these SNPs have distributed effects, often with mixed directions, across regions and FCs. Fig. 2B shows the FCs associations with both 15q14, and 3p11.1 lead SNPs. The regional effects of all other lead SNPs can be appreciated in the Fig. S13.

## 3.3. Downstream analyses

### 3.3.1. SNP-based genetic correlation of functional connectivity measures

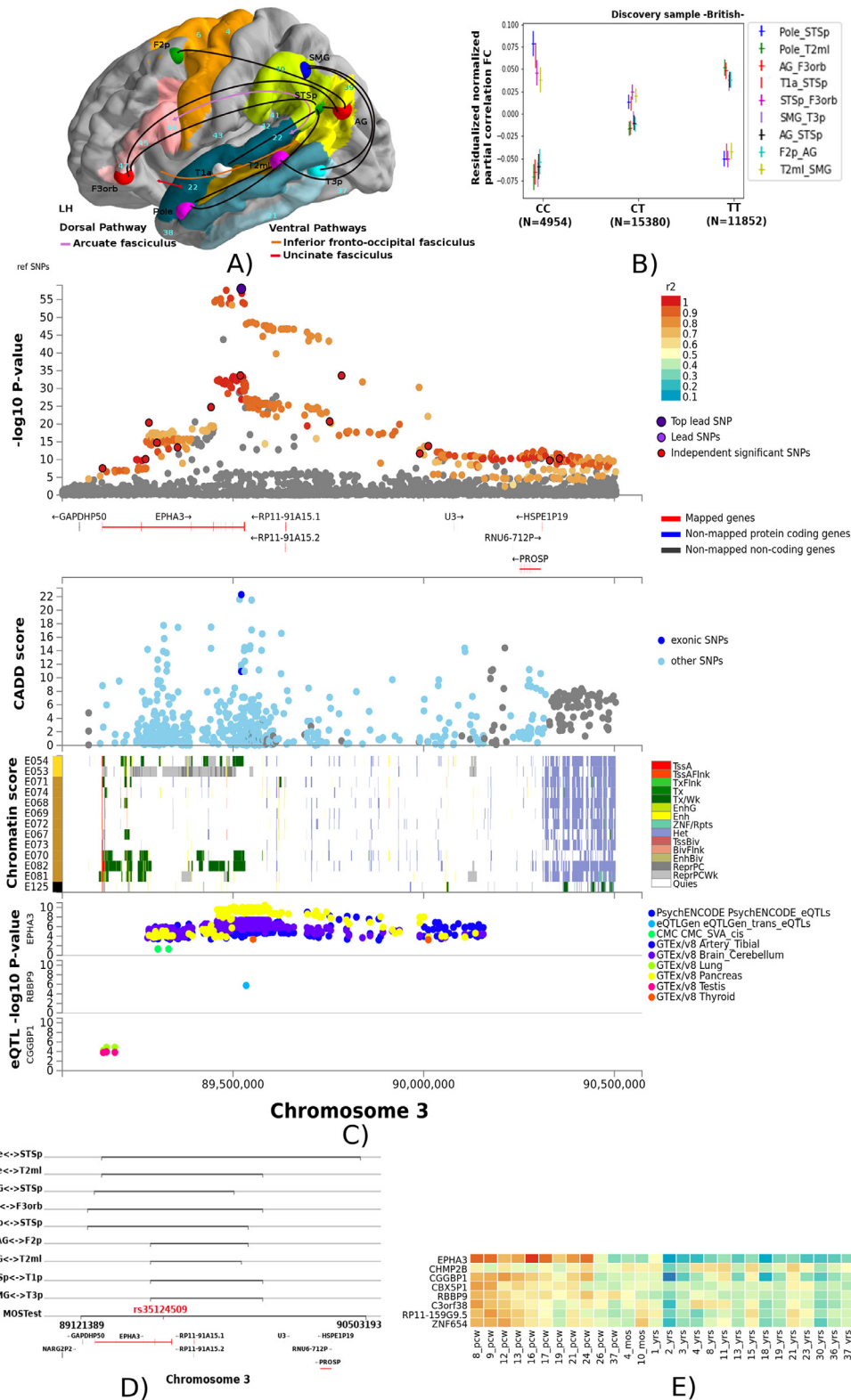
The SNP-based genetic correlation analysis was estimated for each pair of central FCs associated to 15q14 or 3p11.1 genetic loci, indicating overlapping genetic contributions among several FCs (Table S16). For central endophenotypes associated with 3p11.1 locus, a negative genetic correlation between some FCs has been observed which indicates that variants can have antagonistic effects on the co-activations of these regions.

### 3.3.2. Validation of lead SNPs using diffusion imaging derived endophenotypes

We hypothesised that the genetic variants significantly associated with the language FCs could be associated with neuroanatomical tracts that support the information transmission between language areas. Therefore, we tested the potential associations between the hit SNPs with the average values of dMRI relevant white matter tracts:

Three white matter tracts to be tested with locus on 15q14: the white matter tracts linking the regions of the (Prec $\leftrightarrow$ F3opd) consists of the i) arcuate anterior segment fasciculus (AF) dorsal pathway (Catani et al., 2005), ii) the frontal aslant tract (FAT) which is reported as connecting Broca's region (BA44/45) with dorsal medial frontal areas including supplementary and pre-supplementary motor area (BA6) (Rojkova et al., 2016; Catani and Forkel, 2019) while the anatomical connectivity underlying the (PrecR $\leftrightarrow$ RolS) FC endophenotype consists of the corpus callosum which interconnects both hemispheres.

Five white matter tracts to be tested with locus on 3p11.1: the nine central endophenotypes associated with the 3p11.1 locus, the anatomical connectivity underlying these connections consists of the (i) inferior fronto-occipital fasciculus (IFOF) which connects the inferior frontal regions with the temporal and occipital cortex (Forkel et al., 2014b), (ii) uncinate fasciculus (UF) which is reported to connect the anterior temporal lobe to the orbital region and part of the inferior frontal (Vigneau et al., 2006; Catani and De Schotten, 2008; Friederici, 2017; Catani and Forkel, 2019), and the (iii) arcuate long/anterior/posterior segment fasciculus (AF) (Catani et al., 2005). As the anterior segment of AF is tested with both loci, this yields a Bonferroni-corrected threshold of  $p = 6.94e-3(0.05/(3 * 9 + 5 * 9))$  (Table S17). The MO measured in the FAT and the OD measured in the anterior segment of AF are associated with the rs1440802 SNP with  $p = 3.33e-6$  and  $p = 2.47e-65$ , respectively. The corpus callosum exhibits no significant association. The MO measured in the IFOF and UF is associated with the rs35124509 SNPs with  $p = 2.49e-7$  and  $p = 2.12e-7$ , respectively. Both long and posterior segment of AF are associated with rs35124509 SNPs with  $p = 5.88e-6$  (OD) and  $p = 3.90e-7$  (L3), while the anterior segment of AF exhibits no significant association (See Table S17).



**Fig. 3.** Main results for the 15q14 locus. (A) The two pairs of ROIs that forms the endpoints of the associated FCs reported as black bold lines. (B) Effect sizes of the SNP rs1440802 for the two connections: (Prec $\leftrightarrow$ F3opd) FC in green and (PrecR $\leftrightarrow$ RoIS) FC in yellow. (C) Locus Zoom of the genomic region identified by the mvGWAS. Chromatin state of the genomic region. Brain tissue name abbreviations are the following: E054:Ganglion Eminence derived primary cultured neurospheres, E053: Cortex derived primary cultured neurospheres, E071: Brain Hippocampus Middle, E074: Brain Substantia Nigra, E068: Brain Anterior Caudate, E069: Brain Cingulate Gyrus, E072: Brain Inferior Temporal Lobe, E067:Brain Angular Gyrus, E073: Brain Dorsolateral Prefrontal Cortex, E070: Brain Germinal Matrix, E082: Fetal Brain Female, E081: Fetal Brain Male, E125: NH-A Astrocytes Primary Cells. The state abbreviations are the following: TssA: active transcription start site (TSS), TssFlnk: Flanking Active TSS, TxFlnk: Transcription at gene 50 and 30, Tx: Strong transcription, TxWk: Weak transcription, EnhG: Genic enhancers, Enh: Enhancers, ZNF/Rpts: ZNF genes & repeats, Het: Heterochromatin, TssBiv: Bivalent/Poised TSS, BivFlnk: Flanking Bivalent TSS/Enh, EnhBiv: Bivalent Enhancer, ReprPC: Repressed PolyComb, ReprPCWk: Weak Repressed PolyComb, Quies: Quiescent/Low. Expression quantitative trait loci (eQTL) associations (data source: eQTLGen (Vosa et al., 2018), PsychENCODE (Wang et al., 2018), DICE (Schmiedel et al., 2018), BIOS QTL browser (Zhernakova et al., 2017), GTEx/v8 (Consortium et al., 2017), eQTLcatalogue). (D) Overlap of the genomic region risk region identified from FUMA for MOSTest results, (Prec $\leftrightarrow$ F3opd) and (PrecR $\leftrightarrow$ RoIS). (E) Gene expression from BrainSpan for the interesting genes prioritised by FUMA.





**Fig. 4.** Main results for the 3p11.1 locus. (A) The pairs of ROIs that forms the endpoints of the associated FCs reported as black bold lines. (B) Effect sizes of the SNP rs35124509 for the nine connections: (AG↔F3orb), (Pole↔STSp), (Pole↔T2ml), (T1a↔STSp), (STSp↔F3orb), (SMG↔T3p), (AG↔STSp), (F2p↔AG) and (T2ml↔SMG) FCs. (C) Locus Zoom of the genomic region identified by the mvGWAS. Chromatin state of the genomic region. Brain tissue name abbreviations are the following; E054:Ganglion Eminence derived primary cultured neurospheres, E053: Cortex derived primary cultured neurospheres, E071: Brain Hippocampus Middle, E074: Brain Substantia Nigra, E068: Brain Anterior Caudate, E069: Brain Cingulate Gyrus, E072: Brain Inferior Temporal Lobe, E067:Brain Angular Gyrus, E073: Brain Dorsolateral Prefrontal Cortex, E070: Brain Germinal Matrix, E082: Fetal Brain Female, E081: Fetal Brain Male, E125: NH-A Astrocytes Primary Cells. The state abbreviations are the following; TssA: active transcription start site (TSS), TssFlnk: Flanking Active TSS, TxFlnk: Transcription at gene 50 and 30, Tx: Strong transcription, TxWk: Weak transcription, EnhG: Genic enhancers, Enh: Enhancers, ZNF/Rpts: ZNF genes & repeats, Het: Heterochromatin, TssBiv: Bivalent/Poised TSS, BivFlnk: Flanking Bivalent TSS/Enh, EnhBiv: Bivalent Enhancer, ReprPC: Repressed PolyComb, ReprPCWk: Weak Repressed PolyComb, Quies: Quiescent/Low.

### 3.4. Functional annotations of genomic loci associated with language

#### 3.4.1. Locus in 15q14 associated to (Prec↔F3opd) and (PrecR↔RolS) endophenotypes

Four independent SNPs were identified in locus 15q14 (rs1440802, rs11629938, rs773225188, rs34680120) (Fig. 3C). Regarding eQTL annotations, we explored tissue-specific gene expression resources, including both brain tissues and blood, considered as a good proxy when brain tissues are not available (Qi et al., 2018). Significant results were obtained.

The four independent SNPs are cis-eQTL of *THBS1* gene in eQTL-Gen, BIOSQTL and GTEx/v8. Additionally, rs34680120 is eQTL of *RP11-37C7.1* gene ( $p_{\text{adj}} < 1.02e-3$ ) in PsychENCODE and eQTL of *CTD-2033D15.1* gene ( $p_{\text{adj}} < 6.0e-6$ ) in BIOSQTL; see Fig. 3C. Overall, the variants of this genomic risk region are found 72 times as eQTL of genes from different data sources. All eQTL associations are presented in more detail in Table S17. Based on the human gene expression data from the Brainspan database, we found that *THBS1* gene has relatively high mRNA (messenger Ribonucleic Acid) expression during early mid-prenatal to late prenatal stages, from 16 to 37 post-conceptual weeks; see Fig. 3E. Indirect predictions might be added from the following annotation. *RASGRP1*, identified by chromatin interaction mapping and which also appears to be under control of temporal expression during neurodevelopment, is reported as over-expressed in the perisylvian language areas (Johnson et al., 2009) and as up-regulated in the dorsal striatum (Cimaru et al., 2020). Fig. 3 summarises these results, found by mvGWAS, associated to (Prec↔F3opd) and (PrecR↔RolS) FC endophenotypes. These pinpoint *THBS1* as the possible gene underlying this association signal.

#### 3.4.2. Locus in 3p11.1 associated to semantic-language related endophenotypes

Fourteen independent SNPs were identified in locus 3p11.1 (Fig. 4C). The rs35124509 SNP is a non-synonymous variant within exon 16 of *EPHA3* protein-coding gene. The subregion around rs35124509 and rs113141104 has its chromatin state annotated as (weak) actively transcribed states (Tx, TxWk) in the brain tissues, specifically in the Brain Germinal Matrix, the Ganglion Eminence derived primary cultured neurospheres, and in the Fetal Brain Female. Concerning the subregion around rs6551410, it has its chromatin state annotated as Weak transcription (TxWk) in the Fetal Brain Female, enhancer (enh) in the Brain Germinal Matrix and Repressed PolyComb (ReprPC) in both the Ganglion Eminence and Cortex derived primary cultured neurospheres. Additionally, the subregion around rs6551407 has its chromatin state annotated as Weak transcription (TxWk) in the Brain Germinal Matrix, Fetal Brain Male and Fetal Brain Female. Overall, this reveals a genomic region involved in fine regulation mechanisms of brain development.

Considering the rs35124509 SNP and variants in linkage disequilibrium (LD) with it in the genomic risk region, we scrutinised CADD and RDB scores, precise genomic positions and risk prediction, and we noticed some remarkable SNPs. We observed two exonic variants: (i) The SNP rs1054750 ( $LD_{\text{rs35124509}} r^2 > 0.99$ ,  $p_{\text{mvGWAS}} = 6.65e-34$ ) is a synonymous variant within exon 16 of *EPHA3*. (ii) the already mentioned non-synonymous lead SNP rs35124509 ( $p_{\text{mvGWAS}} = 8.95e-59$ ), the minor allele results in a substitution in the protein from tryptophan (W) residue (large size and aromatic) into an arginine (R) (large size and basic) at position 924 (W924R, p.Trp924Arg) in the Sterile Alpha Motif (SAM) domain. This SNP is not predicted to alter protein function (Polyphen-2="benign") but is predicted to be potentially a regulatory element by several tools (RDB score = 3a, CADD=22.3 - when  $CADD_{\text{thresh}} = 12.37$  for deleterious effect as suggested by Kircher et al. (2014)) Moreover, we

observed eight SNPs (rs28623022, rs7650184, rs7650466, rs73139147, rs3762717, rs73139144, rs73139148, rs566480002) ( $LD_{\text{rs35124509}} r^2 > 0.73$ ,  $p_{\text{mvGWAS}} < 4.46e-20$ ) located in 3'-UTR of *EPHA3* which could affect its expression by modulating miRNA (micro Ribonucleic Acid) binding (Popp et al., 2016). The hit-SNP rs35124509 and the rest of highlighted SNPs act as eQTL for *EPHA3* in different tissues including brain cerebellum ( $p_{\text{FDR}} < 5e-2$  in GTEx/v8 data source). The exhaustive eQTL associations are presented in Table S18. Fig. 4 summarises the functional annotations in 3p11.1 associated to multiple FC endophenotypes in semantic component of language. These functional characterization supports *EPHA3* as a possible gene with a key role in language development in humans.

#### 3.4.3. Locus in 10q23.33

Four independent SNPs were identified in locus 10q23.33 (rs11187838, rs17109875, rs11187844, rs20772180). The subregion around all four SNPs has its chromatin state annotated as (weak) actively-transcribed states (Tx, TxWk) in the brain tissues, specifically in the ganglion eminence and cortex derived primary cultured neurospheres, hippocampus (middle), substantia nigra, anterior caudate, angular gyrus, Dorsolateral/Prefrontal cortex, brain germinal matrix, fetal brain female/male and NH-A (normal human astrocytes) primary cells. Two exonic variants are noteworthy: The rs2274224 ( $LD_{\text{rs11187838}} r^2 > 0.99$ ,  $p_{\text{mvGWAS}} = 5.04e-14$ ) and rs11187895 ( $LD_{\text{rs17109875}} r^2 > 0.6$ ,  $p_{\text{mvGWAS}} = 3.08e-7$ ) SNPs are nonsynonymous variants within exon 19 of *PLCE1* and exon 11 of *NOC3L* and are both not predicted to alter protein function (Polyphen-2="benign") but are predicted to have a deleterious effect (CADD = 17.35, CADD = 19.24). Moreover, we observed three SNP (rs11187870, rs11187877, rs145707916) ( $LD_{\text{rs17109875}} r^2 > 0.66$ ,  $p_{\text{mvGWAS}} < 7.546e-7$ ) located in 3'-UTR of *PLCE1:NOC3L*. Regarding eQTL annotations, the variants in the 10q23.33 locus act as eQTL for *HELLS*, *NOC3L* and *PLCE1* genes in different brain tissues including brain cerebellum, brain cerebellar hemisphere, Brain nucleus accumbens basal ganglia, hippocampus ( $p_{\text{FDR}} < 5e-2$  in GTEx/v8 data source). The exhaustive eQTL associations are presented in Table S18. These functional characterizations highlight these three genes (*HELLS*, *NOC3L* and *PLCE1*) that may influence the FCs related to language processing in humans.

## 4. Discussion

In this study, we extracted individual language FC endophenotypes from the rsfMRI data of 32,186 participants from the UK Biobank cohort and conducted a mvGWAS. We found 4566 significantly associated SNPs distributed over 11 chromosomes. Three multivariate associations with lead SNPs were replicated in the non-British cohort, highlighting the robustness of these signals across different ancestries. Two functional connections, contributing in the *perceptual motor* interaction, associated with 15q14 locus located in the *RP11-624L4.1* antisense gene with modulatory effects on the expression of the *THBS1* gene. Multiple FCs in the *fronto-temporal semantic* language network were found to be associated with SNPs regulating *EPHA3* gene expression in 3p11.1 locus. Each lead SNP was found to be associated with the neuroanatomical white matter tracts that support each of these FCs.

### 4.1. Locus regulating *THBS1* associated with the *perceptual motor* interactions process

A locus in 15q14 was associated with the precentral-opercular FC (Prec↔F3opd) and the precentral-Rolandic FC endophenotypes (PrecR↔RolS). The L—R Prec regions in the ventral precentral gyrus

are both associated with phonology language component and considered relevant for pharynx and tongue fine-movement coordination in the human and nonhuman primates (Vigneau et al., 2006; Kumar et al., 2016; Belyk and Brown, 2017). RoLS in the dorsal Rolandic sulcus is attributed to the phonology component and matches the mouth primary motor area but also the perception of syllables (Vigneau et al., 2006; Wilson et al., 2004; Fadiga et al., 2002). F3opd in the dorsal pars opercularis (BA44/45) is associated with semantic/sentence processing. The motor theory of speech perception has been quite an old debate (Liberman and Mattingly, 1985; Galantucci et al., 2006; Flinker et al., 2015; Schwartz et al., 2008; Whalen, 2019). In this study, we report a locus in 15q14 (lead SNP rs1440802) associated with both this FC between the motor and Broca's areas and the frontal aslant tract connecting directly (pre)supplementary motor area with the opercular part of inferior frontal gyrus (Vergani et al., 2014; Catani et al., 2012), in line with this perception–motor link.

SNPs in high linkage disequilibrium (LD) with rs1440802 in the genomic region have been linked to several other structural features (surface area and cortical thickness) including primary motor cortex, primary somatosensory cortex (Elliott et al., 2018; van der Meer et al., 2020), supramarginal, and pars opercularis (van der Meer et al., 2020), supporting a common genetic influence of the sensory-motor interaction.

The lead SNP rs1440802 and SNPs in LD uncovered to be associated with both (Prec $\leftrightarrow$ F3opd) and (PrecR $\leftrightarrow$ RoLS) are found to be eQTL of *THBS1* gene in the blood with high confidence. The thrombospondin-1 protein encoded by *THBS1* gene is a member of the thrombospondin family, a glycoprotein expressed in the extracellular matrix. It has been implicated in synaptogenesis (Christopherson et al., 2005) and regulates the differentiation and proliferation of neural progenitor cells (Lu and Kipnis, 2010), and has been involved in human neocortical evolution (C'aceres et al., 2003, 2007). Other members of the thrombospondin family, *THBS2* and *THBS4*, have been shown to be over-expressed in the adult human cerebral cortex compared to chimpanzees and macaques (C'aceres et al., 2007). Their increased expression suggests that human brain might display distinctive features involving enhanced synaptic plasticity in adulthood which may contribute to cognitive and linguistic abilities (Sherwood et al., 2008). From a developmental point of view, *THBS1* appears to be under control of temporal expression during development, as revealed by BrainSpan data (See Figs. 3E and SI4). *THBS1* expression was studied from the longitudinal transcriptomic profile resource of the developing human brain (18, 19, 21, 23 weeks of gestation) (Johnson et al., 2009). Its expression is reported as over-expressed in the neocortex, including the perisylvian language areas, compared to phylogenetically older parts of the brain such as the striatum, thalamus and cerebellum (Johnson et al., 2009). Thrombospondin-1 have been linked to Autism spectrum disorder (Lu et al., 2014), Alzheimer's disease (Ko et al., 2015), and Schizophrenia (Park et al., 2012).

Taken together, these results indicate that *THBS1*, modulated by a lead SNP in the 15q14 locus, could be prioritised in the study of key genes playing a role in the functional connectivity part of the perceptual motor interaction required for language, and with the anatomical connectivity, support of their interactions.

#### 4.2. Locus in *EPHA3* associated with the fronto-temporal semantic network

A locus in 3p11.1 is found associated with nine fronto-parietal-temporal endophenotypes. The angular gyrus (AG) has been shown to activate during functional imaging tasks probing semantics and involved in conceptual knowledge (Vigneau et al., 2006). F3orb in the pars orbitalis in the inferior frontal gyrus is labelled semantic for its involvement in semantic retrieval in spoken and sign language (Rönnerberg et al., 2004). It has also been associated with categorization, association, and word generation tasks (Noppeney and Price, 2004; Booth et al., 2002; Gurd et al., 2002). The temporal pole region, located in the anterior temporal lobe, is associated with semantic and sentence processing

(Vigneau et al., 2006) and the posterior superior temporal sulcus (pSTS) is reported to be implicated in syntactic complexity (Constable et al., 2004) but also process the semantic integration of complex linguistic material (Vigneau et al., 2006). Both pSTS and the angular gyrus overlap with the Geschwind's territory (See Fig. 1a). The lateral/middle part of the middle temporal gyrus is devoted to verbal knowledge (Vigneau et al., 2006). These regions and their corresponding endophenotypes fit rather well with the *fronto-temporal semantic system* described in (Vigneau et al., 2006) facilitating the association of integrated input messages with internal knowledge. The anterior part of the superior temporal gyrus and the posterior part of the inferior temporal gyrus are phonological–semantic interface areas processing. (Vigneau et al., 2006) propose that these ones are transitional zones between the perception and semantic integration of language stimuli and are crucial during the development of language.

SNPs of this genomic region in high LD with the lead SNP rs35124509 have already been found associated with: rsfMRI ICA FC (edge 387, 383, 399, and ICA-features 3); see (Elliott et al., 2018). The ICA maps used for these FC estimations partially-overlap semantic language areas including the angular gyrus, the most anterior part of the STS, the anterior fusiform gyrus, the lateral-middle part of T2, the ventral part of the pars triangularis and the pars orbitalis of the left inferior frontal gyrus. Regarding cognitive traits, this locus was associated to intelligence (Savage et al., 2018). Finally, other SNPs, in strong LD with the lead SNP rs35124509, consistently act as an eQTL of *EPHA3* in brain tissues.

The ephrin type-A receptor 3 protein encoded by *EPHA3* gene belongs to the ephrin receptor family that can bind the ephrins subfamily of the tyrosine kinase protein family. EPH (erythropoietin-producing hepatocellular carcinoma) receptors and their ligands were found to play important roles in multiple developmental processes, including tissue morphogenesis, embryogenesis, neurogenesis, vascular network formation, neural crest cell migration, axon fasciculation, axon guidance, and topographic neural map formation (Pasquale, 2008; Gibson and Ma, 2011; Gerstmann and Zimmer, 2018). *EPHA3* binds predominantly EFNA5 and plays a role in the segregation of motor and sensory axons during neuromuscular circuit development (Lawrenson et al., 2002). In Johnson et al. (2009), *EPHA3* is reported as over-expressed in the fetal rhesus macaque monkey neocortex (NCTX) and especially in the occipital lobe compared to the other NCTX areas. Noticeably, its ligand *EFNA5* is overexpressed in perisylvian areas and is located in a human accelerated conserved non-coding sequence (haCNS704) (Johnson et al., 2009). EPH receptors have been linked to neurodevelopmental disorders, including schizophrenia (Zhang et al., 2010) and autism spectrum disorder (Casey et al., 2012). Moreover, in Rudov et al. (2013), *EPHA3* is found *in silico*, as putative gene implicated in dyspraxia, dyslexia and specific language impairment (SLI). Finally, we observed that *EPHA3* is expressed in the human brain, in a consistent manner across developmental stages from early prenatal to late-mid prenatal (8–24 pcw, BrainSpan; see Figs. 4E and SI4).

Taken together, these results indicate that *EPHA3* in the 3p11.1 locus, could be prioritised in the study of key genes playing a role in the fronto-temporal semantic network, and with the anatomical connectivity support of this network.

#### 4.3. Locus related to *PLCE1*, *NOC3L* and *HELLS*

A locus in 10q23.33 was highlighted by the mvGWAS. At the univariate level, no endophenotype reached the genome-wide significance threshold. But looking at the suggestive threshold  $p = 1e-5$ , we pinpoint putative 'central' endophenotypes to aid interpretation of the processes underlying this association signal. Two bilateral fronto-temporal endophenotypes were the most associated to rs11187838: the precentral-Rolandic FC endophenotypes (PrecR $\leftrightarrow$ RoLS,  $p = 1.85e-07$ ) and the right anterior part of the superior temporal gyrus (T1aR) overlapping Heschl's gyrus (T1a/HeschlR) and its homotopic areas of left hemisphere pri-

primary auditory regions (T1a↔T1a/HeschlR,  $p = 9.61e - 06$ ). All these regions participate in an elementary audio-motor loop involved in both comprehension and production of syllables forming a bilateral fronto-temporal network activated by the auditory representation of speech sounds (Vigneau et al., 2006, 2011). SNPs of this genomic region in high LD with the lead SNP rs11187838 have already been found associated with: rsfMRI ICA FC network within the perisylvian area (Elliott et al., 2018). These act as an eQTL of *HELLS*, *NOC3L*, *PLCE1* genes in multiple brain tissues (Table S18). The *HELLS* gene encodes the lymphoid-specific helicase (Lsh), a member of the SNF2 helicase family of chromatin remodeling proteins. Patients with a genetic mutation of *HELLS* present psychomotor retardation including slow cognitive, motor development and psychomotor impairment (Thijssen et al., 2015). The Lsh protein might play a role as epigenetic regulator in neural cells (Han et al., 2017). Finally, we observed that the three genes (*NOC3L*, *PLCE1*, *HELLS*) are expressed in the human brain, across developmental stages from early prenatal to early mid prenatal (8–17 pcw, BrainSpan).

Taken together, these results indicate that the three highlighted genes (*PLCE1*, *NOC3L* and *HELLS*) in the 10q23.33 locus, as potential candidates in the study of key genes playing a role in the *bilateral fronto-temporal auditory-motor* network.

#### 4.4. Limitations

Although functional MR imaging is recognised to produce valuable endophenotypes (Elliott et al., 2018), resting-state FC represent remote measures of language as the participants are not engaged in any language task-experiments and simply stay at rest. As stated in the introduction, previous works reported that resting-state derived endophenotypes are correlated with behavioral language processing (Koyama et al., 2011; Stevens et al., 2017; Cross et al., 2021; Cheema et al., 2021). However, with the release at hand in this study, no behavioral language scores were provided in the UK Biobank cohort. We were thus unable to check the correlations between our resting-state FC endophenotypes and language behavioral scores. Moreover, we observed a hierarchy of language heritability estimates from the Human Connectome Project (twin study) with 51–67% for language-behavioral scores, 22–55% for language-task activations (Le Guen et al., 2018) and 10–50% for resting-state FC (*Data not shown*), making resting-state FC a less suitable endophenotype than language task activations for an association study about language, and therefore less suitable than language scores. The lack of a large, age-matched replication sample represents one major limitation of the present study in the sense that we could not reproduce all our results. Although multivariate methods have shown to substantially increase statistical power and gene discovery compared to univariate approaches, the results are less straightforward to interpret. We have addressed this issue by assessing each of the prioritized loci at the univariate level, to pinpoint at central endophenotypes that are contributing the most to the multivariate signal. Moreover, as such a complex trait as language may be driven by a lot of interacting genes, a multivariate approach on the SNPs side is highly desired to uncover relevant gene pathways in language development and processing. Compared to structural endophenotypes, the FCs have low amplitude which hinders the study in terms of statistical power. This observation constitutes a third limitation that is somehow surpassed when working on large scale cohorts and using multivariate approaches. Another potential limitation is the UK Biobank dataset in which this study is based. It should be noted that the UKB constitutes a relatively old sample. Future studies in other developmental stages (i.e. children, adolescent, young-adult) will inform us whether the observed associations are stable across development, or whether they reflect some age-related specificity.

#### 5. Conclusions

In this study, we extracted language endophenotypes from rsfMRI acquisitions in the largest imaging-genetic cohort on general population

to date. To make these endophenotypes as closer as possible to language function in the brain, we adopted a ROI-based approach for FC estimation with language ROIs derived from a comprehensive meta-analysis of tbfMRI studies on language. This approach makes the endophenotypes comparable across individuals and thus suitable for an association study. After filtering on heritability significance, we performed a multivariate GWAS technique in order to take advantage of the correlation structure among rsfMRI FC and uncover potential pleiotropic loci. Thereby, we highlighted potential key genes related to language processing including *EPHA3* gene in the 3p11.1 locus with a role in the *fronto-temporal semantic* network, *THBS1* gene, modulated by the 15q14 locus, associated with the functional connectivity part of the *perceptual motor* interaction required for language, and *PLCE1* gene in the 10q23.33 locus with a potential role in the *bilateral fronto-temporal auditory-motor* network. These genes could be prioritized to study language in suitable genetic model. Furthermore, two results that have not shown significant association in the replication sample, have been reported in previous works; a 3q24 locus in *ZIC4* which is involved in visual and auditory pathway development (Hornig et al., 2009) has been associated with brain asymmetry (Sha et al., 2021); a 14q23.1 locus near *DACT1* has been reported to be associated with the STAP (Le Guen et al., 2018) and superior temporal sulcus surface area (Sha et al., 2021). Altogether, these results provide an novel insight into the genetic architecture of the language in humans. A growing number of works claim that language studies should consider task-free fMRI data in general population in order to really focus on the neurobiological organization of language and how it supports natural language, i.e. as it is used in everyday life (Hasson et al., 2018). On a clinical side, such research settings are highly desirable as there is a need to map language areas in patients unable to perform language tasks (Ramage et al., 2020; Branco et al., 2016; Klingbeil et al., 2019; Park et al., 2020). By validating our findings in rsfMRI FC endophenotypes with dMRI-derived endophenotypes, we tried to provide both a functional and structural connectome point of view of the language organization in the brain. Finally, recent imaging genetic approaches such as multivariate GWAS, allowed us to increase in statistical power and circumvent the small effect sizes of these original endophenotypes. We believe that such approaches are really promising and will broadly disseminate in the imaging-genetic field and beyond. With the presented approach, we tried to contribute to these innovative trends and to pave the way to other alternative task-free approaches to study natural language and its genetic underpinnings.

#### Declaration of Competing Interest

The authors declare that they have no conflict of interest.

#### Credit authorship contribution statement

**Yasmina Mekki:** Investigation, Formal analysis, Validation, Visualization, Writing – original draft. **Vincent Guillemot:** Writing – original draft, Writing – review & editing, Supervision. **Hervé Lemaître:** Writing – review & editing. **Amaia Carrión-Castillo:** Writing – review & editing. **Stephanie Forkel:** Writing – review & editing. **Vincent Frouin:** Writing – original draft, Writing – review & editing, Supervision. **Cathy Philippe:** Conceptualization, Writing – original draft, Supervision, Writing – review & editing.

#### Acknowledgments

This research was conducted using the UK Biobank resource under application #64984. This project was supported by the Marie Skłodowska-Curie program awarded to Stephanie J. Forkel (Grant agreement No. 101028551). Amaia Carrion-Castillo was supported by a Juan de la Cierva fellowship from the Spanish Ministry of Science and Innovation, and a Gipuzkoa Fellows fellowship from the Basque Government.

## Data availability

The FC endophenotypes are available here : <https://biobank.ndph.ox.ac.uk/ukb/docs.cgi?id=1> with project #64984. The summary statistics can be accessed via GWAS Catalog : <https://www.ebi.ac.uk/gwas/studies/GCP000274>

## Ethics statement

UK Biobank dataset: informed consent is obtained from all UK Biobank participants; ethical procedures are controlled by a dedicated Ethics and Guidance Council (<http://www.ukbiobank.ac.uk/ethics>) that has developed with UK Biobank an Ethics and Governance Framework (given in full at <http://www.ukbiobank.ac.uk/wp-content/uploads/2011/05/EGF20082.pdf>), with IRB approval also obtained from the North West Multi-center Research Ethics Committee.

## Code availability

This study used openly available software and codes, specifically GCTA (<https://cnsgenomics.com/software/gcta/#GREML>), PLINK (<http://zzz.bwh.harvard.edu/plink/>), MOSTest (<https://github.com/precimed/mostest>), and FUMA (<https://fuma.ctglab.nl/>). The anatomical connectivity atlas used is available at ([http://www.bcblab.com/BCB/Atlas\\_of\\_Human\\_Brain\\_Connections.html](http://www.bcblab.com/BCB/Atlas_of_Human_Brain_Connections.html)).

## Supplementary materials

Supplementary material associated with this article can be found, in the online version, at [doi:10.1016/j.neuroimage.2021.118795](https://doi.org/10.1016/j.neuroimage.2021.118795).

## References

- Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., Varoquaux, G., 2014. Machine learning for neuroimaging with scikit-learn. *Front. Neuroinform.* 8, 14. doi:10.3389/fninf.2014.00014.
- Alfaro-Almagro, F., Jenkinson, M., Bangerter, N.K., Andersson, J.L., Griffanti, L., Douaud, G., Sotiropoulos, S.N., Jbabdi, S., Hernandez-Fernandez, M., Vallee, E., et al., 2018. Image processing and quality control for the first 10,000 brain imaging datasets from UK Biobank. *Neuroimage* 166, 400–424.
- Ardila, A., Bernal, B., Rosselli, M., 2016. How localized are language brain areas? a review of brodmann areas involvement in oral language. *Arch. Clin. Neuropsychol.* 31, 112–122.
- Bates, E., Wilson, S.M., Saygin, A.P., Dick, F., Sereno, M.I., Knight, R.T., Dronkers, N.F., 2003. Voxel-based lesion-symptom mapping. *Nat. Neurosci.* 6, 448–450.
- Belyk, M., Brown, S., 2017. The origins of the vocal brain in humans. *Neurosci. Biobehav. Rev.* 77, 177–193.
- Binder, J.R., Desai, R.H., Graves, W.W., Conant, L.L., 2009. Where is the semantic system? a critical review and meta-analysis of 120 functional neuroimaging studies. *Cereb. Cortex* 19, 2767–2796.
- Biswal, B., Zerrin Yetkin, F., Haughton, V.M., Hyde, J.S., 1995. Functional connectivity in the motor cortex of resting human brain using echo-planar MRI. *Magn. Reson. Med.* 34, 537–541.
- Booth, J.R., Burman, D.D., Meyer, J.R., Gitelman, D.R., Parrish, T.B., Mesulam, M.M., 2002. Modality independence of word comprehension. *Hum. Brain Mapp.* 16, 251–261.
- Boyle, A.P., Hong, E.L., Hariharan, M., Cheng, Y., Schaub, M.A., Kasowski, M., Karczewski, K.J., Park, J., Hitz, B.C., Weng, S., et al., 2012. Annotation of functional variation in personal genomes using regulomedb. *Genome Res.* 22, 1790–1797.
- Branco, P., Seixas, D., Deprez, S., Kovacs, S., Peeters, R., Castro, S.L., Sunaert, S., 2016. Resting-state functional magnetic resonance imaging for language preoperative planning. *Front. Hum. Neurosci.* 10, 11.
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O'Connell, J., et al., 2018. The UK Biobank resource with deep phenotyping and genomic data. *Nature* 562, 203–209.
- Cáceres, M., Lachuer, J., Zapala, M.A., Redmond, J.C., Kudo, L., Geschwind, D.H., Lockhart, D.J., Preuss, T.M., Barlow, C., 2003. Elevated gene expression levels distinguish human from non-human primate brains. *Proc. Natl. Acad. Sci.* 100, 13030–13035.
- Cáceres, M., Suwyn, C., Maddox, M., Thomas, J.W., Preuss, T.M., 2007. Increased cortical expression of two synaptogenic thrombospondins in human brain evolution. *Cereb. Cortex* 17, 2312–2321.
- Carrion-Castillo, A., Pepe, A., Kong, X.Z., Fisher, S.E., Mazoyer, B., Tzourio-Mazoyer, N., Crivello, F., Francks, C., 2020. Genetic effects on planum temporale asymmetry and their limited relevance to neurodevelopmental disorders, intelligence or educational attainment. *Cortex* 124, 137–153.
- Casey, J.P., Magalhaes, T., Conroy, J.M., Regan, R., Shah, N., Anney, R., Shields, D.C., Abrahams, B.S., Almeida, J., Bacchelli, E., et al., 2012. A novel approach of homozygous haplotype sharing identifies candidate genes in autism spectrum disorder. *Hum. Genet.* 131, 565–579.
- Catani, M., De Schotten, M.T., 2008. A diffusion tensor imaging tractography atlas for virtual *in vivo* dissections. *Cortex* 44, 1105–1132.
- Catani, M., Dell'Acqua, F., Vergani, F., Malik, F., Hodge, H., Roy, P., Valabregue, R., De Schotten, M.T., 2012. Short frontal lobe connections of the human brain. *Cortex* 48, 273–291.
- Catani, M., Forkel, S.J., 2019. Diffusion Imaging Methods in Language Sciences. *The Oxford Handbook of Neurolinguistics*, p. 212.
- Catani, M., Jones, D.K., Ffytche, D.H., 2005. Perisylvian language networks of the human brain. *Ann. Neurol.* 57, 8–16 Official Journal of the American Neurological Association and the Child Neurology Society.
- Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., Lee, J.J., 2015. Second-generation plink: rising to the challenge of larger and richer datasets. *Gigascience* 4, s13742–s14015.
- Cheema, K., Ostevik, A.V., Westover, L., Hodgetts, W.E., Cummine, J., 2021. Resting-state networks and reading in adults with and without reading impairments. *J. Neurolinguist.* 60, 101016.
- Christopherson, K.S., Ullian, E.M., Stokes, C.C., Mallowney, C.E., Hell, J.W., Agah, A., Lawler, J., Mosher, D.F., Bornstein, P., Barres, B.A., 2005. Thrombospondins are astrocyte-secreted proteins that promote CNS synaptogenesis. *Cell* 120, 421–433.
- Ciraru, M.D., Song, S., Tshilenge, K.T., Corwin, C., Mleccko, J., Aguirre, C.G., Benhabib, H., Bendl, J., Fullard, J., Apontes, P., et al., 2020. Transcriptional and epigenetic characterization of early striosomes identifies foxf2 and olig2 as factors required for development of striatal compartmentation and neuronal phenotypic differentiation. *bioRxiv*.
- Cole, M.W., Bassett, D.S., Power, J.D., Braver, T.S., Petersen, S.E., 2014. Intrinsic and task-evoked network architectures of the human brain. *Neuron* 83, 238–251.
- Consortium, G., et al., 2017. Genetic effects on gene expression across human tissues. *Nature* 550, 204–213.
- Constable, R.T., Pugh, K.R., Berroya, E., Mencl, W.E., Westerveld, M., Ni, W., Shankweiler, D., 2004. Sentence complexity and input modality effects in sentence comprehension: an fMRI study. *Neuroimage* 22, 11–21.
- Cross, A.M., Ramdajal, R., Peters, L., Vandermeer, M.R., Hayden, E.P., Frijters, J.C., Steinbach, K.A., Lovett, M.W., Archibald, L.M., Joannisse, M.F., 2021. Resting-state functional connectivity and reading subskills in children. *Neuroimage* 243, 118529.
- Daducci, A., Canales-Rodríguez, E.J., Zhang, H., Dyrby, T.B., Alexander, D.C., Thiran, J.P., 2015. Accelerated microstructure imaging via convex optimization (amico) from diffusion MRI data. *Neuroimage* 105, 32–44.
- Dubois, J., Adolphs, R., 2016. Building a science of individual differences from fMRI. *Trends Cogn. Sci.* 20, 425–443.
- Elliott, L.T., Sharp, K., Alfaro-Almagro, F., Shi, S., Miller, K.L., Douaud, G., Marchini, J., Smith, S.M., 2018. Genome-wide association studies of brain imaging phenotypes in UK Biobank. *Nature* 562, 210–216.
- Ernst, J., Kellis, M., 2012. ChromHMM: automating chromatin-state discovery and characterization. *Nat. Methods* 9, 215–216.
- Fadiga, L., Craighero, L., Buccino, G., Rizzolatti, G., 2002. Speech listening specifically modulates the excitability of tongue muscles: a TMS study. *Eur. J. Neurosci.* 15, 399–402.
- Fedorenko, E., 2021. The early origins and the growing popularity of the individual-subject analytic approach in human neuroscience. *Curr. Opin. Behav. Sci.* 40, 105–112.
- Fisher, S.E., Vargha-Khadem, F., Watkins, K.E., Monaco, A.P., Pembrey, M.E., 1998. Localisation of a gene implicated in a severe speech and language disorder. *Nat. Genet.* 18, 168–170.
- Fisher, S.E., Vernes, S.C., 2015. Genetics and the language sciences. *Annu. Rev. Linguist.* 1, 289–310.
- Flinker, A., Korzeniewska, A., Sheshyuk, A.Y., Franszczuk, P.J., Dronkers, N.F., Knight, R.T., Crone, N.E., 2015. Redefining the role of Broca's area in speech. *Proc. Natl. Acad. Sci.* 112, 2871–2875.
- Forkel, S., Catani, M., 2018. Structural Neuroimaging. *Research Methods in Psycholinguistics*. Wiley & Sons, pp. 288–308.
- Forkel, S.J., Friedrich, P., Thiebaut de Schotten, M., Howells, H., 2020a. White matter variability, cognition, and disorders: a systematic review.
- Forkel, S.J., Rogalski, E., Sancho, N.D., D'Anna, L., Laguna, P.L., Sridhar, J., Dell'Acqua, F., Weintroub, S., Thompson, C., Mesulam, M.M., 2020b. Anatomical evidence of an indirect pathway for word repetition. *Neurology* 94, e594–e606.
- Forkel, S.J., Thiebaut de Schotten, M., Dell'Acqua, F., Kalra, L., Murphy, D.G., Williams, S.C., Catani, M., 2014a. Anatomical predictors of aphasia recovery: a tractography study of bilateral Perisylvian language networks. *Brain* 137, 2027–2039.
- Forkel, S.J., de Schotten, M.T., Kawadler, J.M., Dell'Acqua, F., Danek, A., Catani, M., 2014b. The anatomy of fronto-occipital connections from early blunt dissections to contemporary tractography. *Cortex* 56, 73–84.
- Fox, M.D., Raichle, M.E., 2007. Spontaneous fluctuations in brain activity observed with functional magnetic resonance imaging. *Nat. Rev. Neurosci.* 8, 700–711.
- Fridriksson, J., Moser, D., Ryalls, J., Bonilha, L., Rorden, C., Baylis, G., 2009. Modulation of frontal lobe speech areas associated with the production and perception of speech movements.
- Friederici, A.D., 2011. The brain basis of language processing: from structure to function. *Physiol. Rev.* 91, 1357–1392.
- Friederici, A.D., 2017. *Language in Our Brain: The origins of a Uniquely Human Capacity*. MIT Press.
- Galantucci, B., Fowler, C.A., Turvey, M.T., 2006. The motor theory of speech perception reviewed. *Psychon. Bull. Rev.* 13, 361–377.

- Gerstmann, K., Zimmer, G., 2018. The role of the eph/ephrin family during cortical development and cerebral malformations. *Med. Res. Arch.* 6 (3). doi:10.18103/mra.v6i3.1694.
- Gibson, D.A., Ma, L., 2011. Developmental regulation of axon branching in the vertebrate nervous system. *Development* 138, 183–195.
- Gurd, J.M., Amunts, K., Weiss, P.H., Zafiris, O., Zilles, K., Marshall, J.C., Fink, G.R., 2002. Posterior parietal cortex is implicated in continuous switching between verbal fluency tasks: an fMRI study with clinical implications. *Brain* 125, 1024–1038.
- Hampson, M., Tokoglu, F., Sun, Z., Schafer, R.J., Skudlarski, P., Gore, J.C., Constable, R.T., 2006. Connectivity–behavior analysis reveals that functional connectivity between left ba39 and broca’s area varies with reading ability. *Neuroimage* 31, 513–519.
- Han, Y., Ren, J., Lee, E., Xu, X., Yu, W., Muegge, K., 2017. Lsh/hells regulates self-renewal/proliferation of neural stem/progenitor cells. *Sci. Rep.* 7, 1–14.
- Hasson, U., Egidi, G., Marelli, M., Willems, R.M.A., 2018. Grounding the neurobiology of language in first principles: the necessity of non-language-centric explanations for language comprehension. *Cognition* 180, 135–157.
- Hornig, S., Kreiman, G., Ellsworth, C., Page, D., Blank, M., Millen, K., Sur, M., 2009. Differential gene expression in the developing lateral geniculate nucleus and medial geniculate nucleus reveals novel roles for *zic4* and *foxp2* in visual and auditory pathway development. *J. Neurosci.* 29, 13672–13683.
- Jackendoff, R., Jackendoff, R.S., 2002. *Foundations of Language: Brain, meaning, grammar, Evolution.* Oxford University Press, USA.
- Jackson, R.L., Hoffman, P., Pobric, G., Ralph, M.A.L., 2016. The semantic network at work and rest: differential connectivity of anterior temporal lobe subregions. *J. Neurosci.* 36, 1490–1501.
- Johnson, M.B., Kawasawa, Y.I., Mason, C.E., Krsnik, Z., Coppola, G., Bogdanović, D., Geschwind, D.H., Mane, S.M., State, M.W., Sestan, N., 2009. Functional and evolutionary insights into human brain development through global transcriptome analysis. *Neuron* 62, 494–509.
- Jones, O.P., Voets, N., Adcock, J., Stacey, R., Jbabdi, S., 2017. Resting connectivity predicts task activation in pre-surgical populations. *NeuroImage Clin.* 13, 378–385.
- Kelly, C., Uddin, L.Q., Shehzad, Z., Margulies, D.S., Castellanos, F.X., Milham, M.P., Petrides, M., 2010. Broca’s region: linking human brain functional connectivity data and non-human primate tracing anatomy studies. *Eur. J. Neurosci.* 32, 383–398.
- Kircher, M., Witten, D.M., Jain, P., O’Roak, B.J., Cooper, G.M., Shendure, J., 2014. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat. Genet.* 46, 310–315.
- Klingbeil, J., Wawrzyniak, M., Stockert, A., Saur, D., 2019. Resting-state functional connectivity: an emerging method for the study of language networks in post-stroke aphasia. *Brain Cogn.* 131, 22–33.
- Ko, C.Y., Chu, Y.Y., Narumiya, S., Chi, J.Y., Furuyashiki, T., Aoki, T., Wang, S.M., Chang, W.C., Wang, J.M., 2015. The *ccaat/enhancer-binding protein delta/mir135a/thrombospondin 1 axis* mediates *pge2*-induced angiogenesis in alzheimer’s disease. *Neurobiol. Aging* 36, 1356–1368.
- Koyama, M.S., Di Martino, A., Zuo, X.N., Kelly, C., Mennes, M., Jutagir, D.R., Castellanos, F.X., Milham, M.P., 2011. Resting-state functional connectivity indexes reading competence in children and adults. *J. Neurosci.* 31, 8617–8624.
- Koyama, M.S., Kelly, C., Shehzad, Z., Penesetti, D., Castellanos, F.X., Milham, M.P., 2010. Reading networks at rest. *Cereb. Cortex* 20, 2549–2559.
- Kumar, V., Croxson, P.L., Simonyan, K., 2016. Structural organization of the laryngeal motor cortical network and its implication for evolution of speech production. *J. Neurosci.* 36, 4170–4181.
- Labache, L., Mazoyer, B., Joliot, M., Crivello, F., Hesling, I., Tzourio-Mazoyer, N., 2020. Typical and atypical language brain organization based on intrinsic connectivity and multitask functional asymmetries. *eLife* 9, 1–31.
- Landi, N., Perdue, M.V., 2019. Neuroimaging genetics studies of specific reading disability and developmental language disorder: a review. *Lang. Linguist. Compass* 13, e12349.
- Lawrenson, I.D., Wimmer-Kleikamp, S.H., Lock, P., Schoenwaelder, S.M., Down, M., Boyd, A.W., Alewood, P.F., Lackmann, M., 2002. Ephrin-a5 induces rounding, blebbing and *de875* adhesion of *epha3*-expressing 293t and melanoma cells by *crkii* and rho-mediated signalling. *J. Cell Sci.* 115, 1059–1072.
- Le Guen, Y., Amalric, M., Pinel, P., Pallier, C., Frouin, V., 2018. Shared genetic aetiology between cognitive performance and brain activations in language and math tasks. *Sci. Rep.* 8, 1–11.
- Le Guen, Y., Leroy, F., Philippe, C., Consortium, I., Mangin, J.F., Dehaene-Lambertz, G., Frouin, V., 2020. Enhancer locus in *ch14q23.1* modulates brain asymmetric temporal regions involved in language processing. *Cereb. Cortex* 30 (10), 5322–5332. doi:10.1093/cercor/bhaa112.
- Ledoit, O., Wolf, M., 2004. A well-conditioned estimator for large-dimensional covariance matrices. *J. Multivar. Anal.* 88, 365–411.
- Lee, S.H., Yang, J., Goddard, M.E., Visscher, P.M., Wray, N.R., 2012. Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinformatics* 28, 2540–2542.
- Lemée, J.M., Berro, D.H., Bernard, F., Chinier, E., Leiber, L.M., Menei, P., Ter Minasian, A., 2019. Resting-state functional magnetic resonance imaging versus task-based activity for language mapping and correlation with perioperative cortical mapping. *Brain Behav.* 9, e01362.
- Leroy, F., Cai, Q., Bogart, S.L., Dubois, J., Coulon, O., Monzalvo, K., Fischer, C., Glasel, H., Van der Haegen, L., B’én’ezit, A., Lin, C.P., Kennedy, D.N., Ihara, A.S., Hertz-Pannier, L., Moutard, M.L., Poupon, C., Brysbaert, M., Roberts, N., Hopkins, W.D., Mangin, J.F., Dehaene-Lambertz, G., 2015. New human-specific brain landmark: the depth asymmetry of superior temporal sulcus. In: *Proceedings of the National Academy of Sciences*, 112 (4). National Academy of Sciences, pp. 1208–1213.
- Liberman, A.M., Mattingly, I.G., 1985. The motor theory of speech perception revised. *Cognition* 21, 1–36.
- Lu, J., Zhang, H., Hameed, N.F., Zhang, J., Yuan, S., Qiu, T., Shen, D., Wu, J., 2017. An automated method for identifying an independent component analysis-based language-related resting-state network in brain tumor subjects for surgical planning. *Sci. Rep.* 7, 1–16.
- Lu, L., Guo, H., Peng, Y., Xun, G., Liu, Y., Xiong, Z., Tian, D., Liu, Y., Li, W., Xu, X., et al., 2014. Common and rare variants of the *thbs1* gene associated with the risk for autism. *Psychiatr. Genet.* 24, 235–240.
- Lu, Z., Kipnis, J., 2010. Thrombospondin 1—a key astrocyte-derived neurogenic factor. *FASEB J.* 24, 1925–1934.
- Marrelec, G., Krainik, A., Duffau, H., P’el’egrini-Issac, M., Leh’ericy, S., Doyon, J., Bernali, H., 2006. Partial correlation for functional brain interactivity investigation in functional MRI. *Neuroimage* 32, 228–237.
- van der Meer, D., Frei, O., Kaufmann, T., Shadrin, A.A., Devor, A., Smeland, O.B., Thomp912 son, W.K., Fan, C.C., Holland, D., Westlye, L.T., et al., 2020. Understanding the genetic determinants of the brain with mostest. *Nat. Commun.* 11, 1–9.
- Miller, K.L., Alfaro-Almagro, F., Bangerter, N.K., Thomas, D.L., Yacoub, E., Xu, J., Bartsch, A.J., Jbabdi, S., Sotiropoulos, S.N., Andersson, J.L., et al., 2016. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nat. Neurosci.* 19, 1523–1536.
- Nishitani, N., Hari, R., 2000. Temporal dynamics of cortical representation for action. *Proc. Natl. Acad. Sci.* 97, 913–918.
- Noppeney, U., Price, C.J., 2004. Retrieval of abstract semantics. *Neuroimage* 22, 164–170.
- Park, H.J., Kim, S.K., Kim, J.W., Kang, W.S., Chung, J.H., 2012. Association of thrombospondin 1 gene with schizophrenia in Korean population. *Mol. Biol. Rep.* 39, 6875–6880.
- Park, K.Y., Lee, J.J., Dierker, D., Marple, L.M., Hacker, C.D., Roland, J.L., Marcus, D.S., Milchenko, M., Miller-Thomas, M.M., Benzinger, T.L., et al., 2020. Mapping language function with task-based vs. resting-state functional MRI. *PLoS One* 15, e0236423.
- Pasquale, E.B., 2008. Eph-ephrin bidirectional signaling in physiology and disease. *Cell* 133, 38–52.
- Popp, N.A., Yu, D., Green, B., Chew, E.Y., Ning, B., Chan, C.C., Tuo, J., 2016. Functional single nucleotide polymorphism in *il-17* a 3’ untranslated region is targeted by *mi-r-4480* *in vitro* and may be associated with age-related macular degeneration. *Environ. Mol. Mutagen.* 57, 58–64.
- Price, C.J., 2012. A review and synthesis of the first 20 years of pet and fMRI studies of heard speech, spoken language and reading. *Neuroimage* 62, 816–847.
- Qi, T., Wu, Y., Zeng, J., Zhang, F., Xue, A., Jiang, L., Zhu, Z., Kemper, K., Yengo, L., Zheng, Z., et al., 2018. Identifying gene targets for brain-related traits using transcriptomic and methylomic data from blood. *Nat. Commun.* 9, 1–12.
- Ramage, A.E., Aytur, S., Ballard, K.J., 2020. Resting-state functional magnetic resonance imaging connectivity between semantic and phonological regions of interest may inform language targets in aphasia. *J. Speech Lang. Hear. Res.* 63, 3051–3067.
- Ramasamy, A., Trabzuni, D., Guelfi, S., Varghese, V., Smith, C., Walker, R., De, T., Coin, L., De Silva, R., Cookson, M.R., et al., 2014. Genetic variability in the regulation of gene expression in ten regions of the human brain. *Nat. Neurosci.* 17, 1418–1428.
- Ramensky, V., Bork, P., Sunyaev, S., 2002. Human non-synonymous SNPs: server and survey. *Nucleic Acids Res.* 30, 3894–3900.
- Rojkova, K., Volle, E., Urbanski, M., Humbert, F., Dell’Acqua, F., De Schotten, M.T., 2016. Atlasing the frontal lobe connections and their variability due to age and education: a spherical deconvolution tractography study. *Brain Struct. Funct.* 221, 1751–1766.
- Rönnberg, J., Rudner, M., Ingvar, M., 2004. Neural correlates of working memory for sign language. *Cogn. Brain Res.* 20, 165–182.
- Rudov, A., Rocchi, M.B.L., Accorsi, A., Spada, G., Procopio, A.D., Olivieri, F., Rippon, M.R., Albertini, M.C., 2013. Putative mirnas for the diagnosis of dyslexia, dyspraxia, and specific language impairment. *Epigenetics* 8, 1023–1029.
- Savage, J.E., Jansen, P.R., Stringer, S., Watanabe, K., Bryois, J., De Leeuw, C.A., Nagel, M., Awasthi, S., Barr, P.B., Coleman, J.R., et al., 2018. Genome-wide association meta-analysis in 269,867 individuals identifies new genetic and functional links to intelligence. *Nat. Genet.* 50, 912–919.
- Schmiedel, B.J., Singh, D., Madrigal, A., Valdovino-Gonzalez, A.G., White, B.M., Zapardiel-Gonzalo, J., Ha, B., Altay, G., Greenbaum, J.A., McVicker, G., et al., 2018. Impact of genetic polymorphisms on human immune cell gene expression. *Cell* 175, 1701–1715.
- Schwartz, J.L., Basirat, A., M’énard, L., Sato, M., 2012. The perception-for-action-control theory (pact): a perceptuo-motor theory of speech perception. *J. Neurolinguist.* 25, 336–354.
- Schwartz, J.L., Sato, M., Fadiga, L., 2008. The common language of speech perception and action: a neurocognitive perspective. *Revue française de linguistique appliquée* 13, 9–22.
- Seghier, M.L., Price, C.J., 2018. Interpreting and utilising intersubject variability in brain function. *Trends Cogn. Sci.* 22, 517–530.
- Sha, Z., Schijven, D., Carrion-Castillo, A., Joliot, M., Mazoyer, B., Fisher, S.E., Crivello, F., Francks, C., 2021. The genetic architecture of structural left-right asymmetry of the human brain. *Nat. Hum. Behav.* 1–14.
- Sherwood, C.C., Subiaul, F., Zawadzki, T.W., 2008. A natural history of the human mind: tracing evolutionary changes in brain and cognition. *J. Anat.* 212, 426–454.
- Smith, S.M., Fox, P.T., Miller, K.L., Glahn, D.C., Fox, P.M., Mackay, C.E., Filippini, N., Watkins, K.E., Toro, R., Laird, A.R., et al., 2009. Correspondence of the brain’s functional architecture during activation and rest. *Proc. Natl. Acad. Sci.* 106, 13040–13045.
- Stevens, W.D., Kravitz, D.J., Peng, C.S., Tessler, M.H., Martin, A., 2017. Privileged functional connectivity between the visual word form area and the language system. *J. Neurosci.* 37, 5288–5297.
- Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., et al., 2015. UK Biobank: an open access resource for identify-

- ing the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* 12, e1001779.
- Tavor, I., Jones, O.P., Mars, R., Smith, S., Behrens, T., Jbabdi, S., 2016. Task-free MRI predicts individual differences in brain activity during task performance. *Science* 352, 216–220.
- Thijssen, P.E., Ito, Y., Grillo, G., Wang, J., Velasco, G., Nitta, H., Unoki, M., Yoshihara, M., Suyama, M., Sun, Y., et al., 2015. Mutations in *cdca7* and *hells* cause immunodeficiency-centromeric instability-facial anomalies syndrome. *Nat. Commun.* 6, 1–8.
- Tie, Y., Rigolo, L., Norton, I.H., Huang, R.Y., Wu, W., Orringer, D., Mukundan, S., Golby, A.J., 2014. Defining language networks from resting-state fMRI for surgical planning—a feasibility study. *Hum. Brain Mapp.* 35, 1018–1030.
- Turner, T.H., Fridriksson, J., Baker, J., Eoute, D., Bonilha, L., Rorden, C., 2009. Obligatory broca's area modulation associated with passive speech perception. *Neuroreport* 20, 492.
- Uddén, J., Hult'én, A., Bendtz, K., Mineroff, Z., Kucera, K.S., VINO, A., Fedorenko, E., Hagoort, P., Fisher, S.E., 2019. Toward robust functional neuroimaging genetics of cognition. *J. Neurosci.* 39, 8778–8787.
- Vergani, F., Lacerda, L., Martino, J., Attems, J., Morris, C., Mitchell, P., de Schotten, M.T., Dell'Acqua, F., 2014. White matter connections of the supplementary motor area in humans. *J. Neurol. Neurosurg. Psychiatry* 85, 1377–1385.
- Vigneau, M., Beaucousin, V., Herve, P.Y., Duffau, H., Crivello, F., Houde, O., Mazoyer, B., Tzourio-Mazoyer, N., 2006. Meta-analyzing left hemisphere language areas: phonology, semantics, and sentence processing. *Neuroimage* 30, 1414–1432.
- Vigneau, M., Beaucousin, V., Herv'e, P.Y., Jobard, G., Petit, L., Crivello, F., Mellet, E., Zago, L., Mazoyer, B., Tzourio-Mazoyer, N., 2011. What is right-hemisphere contribution to phonological, lexico-semantic, and sentence processing?: insights from a meta-analysis. *Neuroimage* 54, 577–593.
- Võsa, U., Claringbould, A., Westra, H.J., Bonder, M.J., Deelen, P., Zeng, B., Kirsten, H., Saha, A., Kreuzhuber, R., Kasela, S., et al., 2018. Unraveling the polygenic architecture of complex traits using blood eqtl metaanalysis. *BioRxiv*, 447367.
- Wain, L.V., Shrine, N., Miller, S., Jackson, V.E., Ntalla, I., Artigas, M.S., Billington, C.K., Kheirallah, A.K., Allen, R., Cook, J.P., et al., 2015. Novel insights into the genetics of smoking behaviour, lung function, and chronic obstructive pulmonary disease (UK bi leve): a genetic association study in UK Biobank. *Lancet Respir. Med.* 3, 769–781.
- Wang, D., Liu, S., Warrell, J., Won, H., Shi, X., Navarro, F.C., Clarke, D., Gu, M., Emani, P., Yang, Y.T., et al., 2018. Comprehensive functional genomic resource and integrative model for the human brain. *Science* 362 (6420). doi:10.1126/science.aat8464.
- Wang, K., Li, M., Hakonarson, H., 2010. Annovar: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* 38, e164–e164.
- Watanabe, K., Taskesen, E., Van Bochoven, A., Posthuma, D., 2017. Functional mapping and annotation of genetic associations with fuma. *Nat. Commun.* 8, 1–11.
- Westra, H.J., Peters, M.J., Esko, T., Yaghootkar, H., Schurmann, C., Kettunen, J., Chris1033 tiansen, M.W., Fairfax, B.P., Schramm, K., Powell, J.E., et al., 2013. Systematic identification of trans eQTLs as putative drivers of known disease associations. *Nat. Genet.* 45, 1238–1243.
- Whalen, D., 2019. The motor theory of speech perception. *Oxford Research Encyclopedia of Linguistics*. Oxford University Press.
- Wilson, S.M., Saygin, A.P., Sereno, M.I., Iacoboni, M., 2004. Listening to speech activates motor areas involved in speech production. *Nat. Neurosci.* 7, 701–702.
- Xiang, H.D., Fonteijn, H.M., Norris, D.G., Hagoort, P., 2010. Topographical functional connectivity pattern in the perisylvian language networks. *Cereb. Cortex* 20, 549–560.
- Yang, J., Benyamin, B., McEvoy, B.P., Gordon, S., Henders, A.K., Nyholt, D.R., Madden, P.A., Heath, A.C., Martin, N.G., Montgomery, G.W., Goddard, M.E., Visscher, P.M., 2010. Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* 42, 565–569.
- Yang, J., Lee, S.H., Goddard, M.E., Visscher, P.M., 2011. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* 88, 76–82.
- Zhang, H., Schneider, T., Wheeler-Kingshott, C.A., Alexander, D.C., 2012. Noddi: practical *in vivo* neurite orientation dispersion and density imaging of the human brain. *Neuroimage* 61, 1000–1016.
- Zhang, R., Zhong, N.N., Liu, X.G., Yan, H., Qiu, C., Han, Y., Wang, W., Hou, W.K., Liu, Y., Gao, C.G., et al., 2010. Is the *efnb2* locus associated with schizophrenia? single nucleotide polymorphisms and haplotypes analysis. *Psychiatry Res.* 180, 5–9.
- Zhernakova, D.V., Deelen, P., Vermaat, M., Van Iterson, M., Van Galen, M., Arindarto, W., Van't Hof, P., Mei, H., Van Dijk, F., Westra, H.J., et al., 2017. Identification of context-dependent expression quantitative trait loci in whole blood. *Nat. Genet.* 49, 139–145.